

Agricultural Systems

Explaining farming systems spatial patterns: a farm-level choice model based on socioeconomic and biophysical drivers --Manuscript Draft--

Manuscript Number:	AGSY-D-20-00630R2
Article Type:	Research Paper
Keywords:	Farming systems; Territorial context; Choice modeling; spatial analysis; random forest
Corresponding Author:	Paulo Flores Ribeiro, PhD University of Lisbon Higher Institute of Agronomy: Universidade de Lisboa Instituto Superior de Agronomia Lisboa, PORTUGAL
First Author:	Paulo Flores Ribeiro, PhD
Order of Authors:	Paulo Flores Ribeiro, PhD José Lima Santos, PhD Maria João Canadas, PhD Ana Novais, PhD Francisco Moreira, PhD Angela Lomba, PhD
Abstract:	<p>CONTEXT: Efforts to bring together landscape analysis and farming systems have failed to explain the drivers behind their spatial distribution. Since agricultural landscapes are an outcome of farmers' decisions, understanding the role of socioeconomic and biophysical drivers of such decisions is essential for policy-making targeting landscape-level provision of public goods and ecosystem services from agriculture.</p> <p>OBJECTIVE: Aiming to better understand the role of these drivers, we focused on a region dominated by agricultural use, with extensive variability in biophysical and socioeconomic conditions. A typology of farming systems was derived from spatially explicit farm-level data provided by the Portuguese agency responsible for Common Agricultural Policy payments, for 2017. Farms were thoroughly characterized through relevant biophysical and socioeconomic variables considered as potential drivers of farming systems.</p> <p>METHODS: A random forest approach was used to develop a farming system choice-model, dependent on those biophysical and socioeconomic variables. Variable importance measures and partial dependence plots were used to explore the role of these variables in explaining the spatial distribution of farming systems and to predict spatial patterns at the landscape scale.</p> <p>RESULTS AND CONCLUSIONS: Results showed that both biophysical and socioeconomic drivers play a significant role in the spatial distribution of most agricultural systems. Its importance, however, varies significantly across farming systems, being crucial for some and almost irrelevant for others. Farm size and climate have proved to be the most relevant drivers for most farming systems. Overall, our approach proved to be quite accurate in predicting patterns of farming systems at the landscape scale.</p> <p>SIGNIFICANCE: The proposed framework has shown great potential as a tool to support information-based policy design to improve agricultural landscape planning, by linking farm-level management decisions with the provision of socially valued public goods from agriculture, perceived at the landscape-level.</p>
Suggested Reviewers:	Clunie Keenleyside ckeenleyside@ieep.eu Guy Beaufoy guy@efnecp.org Juan Oñate

	juan.onate@uam.es
	Eric Ruto eruto@lincoln.ac.uk
Response to Reviewers:	Dear Dr Jagadish Timsina, Editor of Agricultural Systems, As requested, in this new revision I added: 1) a new (clean) version of the manuscript; 2) a new version of Supplementary Information, and 3) a new cover letter explaining the changes that have been made. Best regards, Paulo Flores Ribeiro

Dr. Jagadish Timsina
Editor
Agricultural Systems

March 21, 2021

Dear Dr. Jagadish Timsina,

As requested in your email of 15 February 2021, I inform you that I have now revised the article by making the changes requested by the reviewer 1, namely:

- added a "north" arrow to figure 1;
- I separated the mean and the standard deviation into two columns in table 1;
- I separated the mean, standard deviation, minimum and maximum into different columns in figure 2;
- I increased the size of the partial dependence plots in Supplementary Information.

Hoping that the article is now up to the high standards of Agricultural Systems, I would like to take the opportunity to thank the editor and the reviewers for all the revision work that has contributed so much to improving the article.

Best regards,



Paulo Flores Ribeiro
(pfribeiro@isa.ulisboa.pt)

HIGHLIGHTS

Context/Background:

Farmland landscape patterns can show substantial spatial variations, resulting from farmers' adaptive responses to the environmental context

Objective:

Understanding how distinct biophysical and socioeconomic factors influence farming system choice enables to anticipate farm management decisions

Results:

Farms' biophysical and socioeconomic features are key drivers of farmers' decisions, resulting in the choice of the farming system

Conclusions:

Considering farm-level drivers of farming systems choice to predict landscape patterns is crucial to an informed agricultural planning

Significance:

An innovative tool to assess scenarios of policy or climate change is presented, based on a farm-level farming systems approach

FARMING SYSTEMS SPATIAL PATTERNS IN ALENTEJO, PORTUGAL



Observed

**Farming System
choice modelling**

based on
**Biophysical and
socioeconomic drivers**

Random forest



Predicted

Abstract

CONTEXT: Efforts to bring together landscape analysis and farming systems have failed to explain the drivers behind their spatial distribution. Since agricultural landscapes are an outcome of farmers' decisions, understanding the role of socioeconomic and biophysical drivers of such decisions is essential for policy-making targeting landscape-level provision of public goods and ecosystem services from agriculture.

OBJECTIVE: Aiming to better understand the role of these drivers, we focused on a region dominated by agricultural use, with extensive variability in biophysical and socioeconomic conditions. A typology of farming systems was derived from spatially explicit farm-level data provided by the Portuguese agency responsible for Common Agricultural Policy payments, for 2017. Farms were thoroughly characterized through relevant biophysical and socioeconomic variables considered as potential drivers of farming systems.

METHODS: A random forest approach was used to develop a farming system choice-model, dependent on those biophysical and socioeconomic variables. Variable importance measures and partial dependence plots were used to explore the role of these variables in explaining the spatial distribution of farming systems and to predict spatial patterns at the landscape scale.

RESULTS AND CONCLUSIONS: Results showed that both biophysical and socioeconomic drivers play a significant role in the spatial distribution of most agricultural systems. Its importance, however, varies significantly across farming systems, being crucial for some and almost irrelevant for others. Farm size and climate have proved to be the most relevant drivers for most farming systems. Overall, our approach proved to be quite accurate in predicting patterns of farming systems at the landscape scale.

SIGNIFICANCE: The proposed framework has shown great potential as a tool to support information-based policy design to improve agricultural landscape planning, by linking farm-level management decisions with the provision of socially valued public goods from agriculture, perceived at the landscape-level.

Explaining farming systems spatial patterns: a farm-level choice model based on socioeconomic and biophysical drivers

Response to Reviewers

Reviewer #1:

Manuscript title: Explaining farming systems spatial patterns: a farm-level choice model based on socioeconomic and biophysical drivers.

This paper explores the farming systems typologies and potential drivers with their roles, using spatial farm level data collected from the Alentejo region in Portugal. The author(s) have used bio-physical and socio-economics features to characterize the farming systems. They have used random forest - a machine learning algorithm to identify the key drivers through the measure of variable of relative importance and potential role of these drivers using partial dependence plots that generates the response of these drives over the different sub-systems.

Overall, at my first read, I enjoy reading this paper as there are no such papers that explores the farming systems typologies using on spatially explicit models. Therefore, this paper is some how novel (quite impressive) and it can be considered in Agricultural Systems journal. However, while reading the paper second time, I found several issues associated with the analytical approach and potential drivers of farming systems included in the model. Here are some of my comments and I encourage author(s) to address these comments as this paper is important and it will add value to the literature.

R: We appreciate the kind words of the Reviewer and her/his work in revising the MS, as well as the valuable comments and suggestions that helped to improve the paper.

1. The variables used in the model are highly time variant, but the analytical approach has included data from a single year (2017). However, the author(s) mentioned that the data are collected on yearly basis (line 78-79 of the MS). If the data are accessible, I encourage authors to include multiple years data and check consistency in results across different sub-systems.

R: We understand the Reviewer's comment and agree that it would be interesting to include data from multiple years. However, although a huge amount of spatially explicit, yearly updated, farm-level data has been gathered by agricultural agencies across EU countries for administrative and farm payments purposes, access to these data is notoriously difficult, allegedly due to data confidentiality issues. This has been an important hurdle in our research lab, and very recently (December 2020) we published an article (Santos et a., 2020 – reference now included in the MS) where we make an explicitly plea for greater openness in making this data available, particularly for research purposes. Nevertheless, in this study we had the opportunity to access a set of these data which, although for a single year (2017), covered a region large enough to explore the spatial variability of the territory and farming systems, and

thus advancing the existing understanding of why farmers' productive decisions, made at the same point in time and under the same market and policy context, are so different in space.

An explicit reference to the shortcomings of using data from just one year has now been added to the MS in a new section named "4.5. Shortcomings of the approach and recommendations for future research" (L544-551):

"Because our farm characterization variables report to a single year, the effect of economic or policy variables such as prices or subsidies can only be assumed as underpinning the farmers' choices reflected on the observed 2017 IACS/LPIS data. However, the use of this type of variables in the model, provided that time-series of farm-level data can be made available, would significantly extend the scope of this approach, allowing its use to evaluate policy and price change scenarios. Even without additional temporal data, the framework can take advantage of the wide extension of the study area to perform, e.g., climate-change scenarios assessment, by adopting a space-for-time substitution approach."

In any case, we agree that this is indeed an important issue, so we reviewed the MS by adding the following appeal in section "4.6. Concluding remarks" (L609-615):

"The use of IACS / LPIS data proved to be an invaluable asset for the research, enabling a high-detailed farm-level analysis, not achievable using official statistics and usually only possible through expensive and time-consuming farm surveys, often unfeasible for research works developed at regional scales like the one used in this study. Therefore, it is worth renewing an appeal previously made (Santos et al., 2020; Tóth and Kučas, 2016), addressed at the EU bodies responsible for maintaining the IACS databases, to make them more accessible to the scientific community, while safeguarding confidentiality duties."

2. Don't the farming practices across these sub-systems (different farming systems) affect FS? I even didn't see any sort of farming systems management features in the model. (Ref: Table 2).

R: We understand the Reviewer's concern, which stems from the fact that the MS lacked a definition of the underlying concept of farming system (a concern also raised by Reviewer #3). Here, we followed the farming system (FS) definition proposed by Santos et al (2020), as a set of farms roughly practicing the same crops and agricultural activities, using similar technological processes and input endowments. In this context, and considering the information available in the IACS agricultural database (which basically describes livestock and land use/cover patterns), we assume that the same crops and agricultural activities, when practiced by farms in the same FS, resort to approximately the same production means and techniques. To clarify this issue in the revised version of the MS, we introduced a paragraph in the Introduction section with the definition of farming system (FS) (L58-64):

"The FS concept used in this study follows that proposed by Santos et al. (2020), according to which a FS can be defined as a set of farms roughly practicing the same crops and agricultural activities, using similar technological processes and input endowments. A key aspect in this concept is that only variables resulting from farm management decisions are considered, when defining a FS; all variables that may influence these decisions but do not result from them, at least in the short run (e.g. farm size or fragmentation level, climate, slopes, market or policy), should be considered as exogenous to the FS and, therefore, as potential drivers of the FS choice (Silva et al., 2020)."

3. I wonder why results show the less significance of irrigation? Is there any special reason in Alentejo region, as water is often a limiting factor in agriculture and of course that drives the farming system in the global south?

R: We tested 2 variables describing farm access to irrigation water: WPUBLIC and WPRIVATE. The first described the proportion of the farm area (UAA) inside public irrigation systems and it showed up very significant (high “variable importance” in current Fig. 3). The second was just a dummy variable indicating if the farm had access to surface water (mostly from small water streams or pounds) and, in fact, it showed little significance in all models. We attributed this finding to the fact that, in the Alentejo region, surface water in small/medium water streams is mostly of torrential regime, drying out during summer, and thus not suitable to support most irrigated agricultural systems (see also our response to minor comment #2, below).

4. Model did not include any of the features related with the marketing and it is not discussed anywhere in the MS. Market is the vital component of the FS.

R: From our point of view, the market is an exogenous variable to the farming system, being a potential key driver of farmers' productive choices, but it is not a defining characteristic of the FS, as it does not result from the decisions of individual farmers, as explained in the definition of FS now added to the MS (see reply to comment 2). However, at the same point in time, market conditions (especially product and input prices) are (virtually) the same for all farmers in a given region. Therefore, in a non-temporal model there is no variation in these variables and for that reason, their effects cannot be evaluated (the same goes for the effects of policies). This important limitation of non-temporal models with an economic feature was addressed in the Discussion in the original version of the MS (L533-540), and has now been moved into the new section "4.5. Shortcomings of the approach and recommendations for future research" added in response to the Reviewer's comment #5 (L544-551):

“Because our farm characterization variables report to a single year, the effect of economic or policy variables such as prices or subsidies can only be assumed as underpinning the farmers' choices reflected on the observed 2017 IACS/LPIS data. However, the use of this type of variables in the model, provided that time-series of farm-level data can be made available, would significantly extend the scope of this approach, allowing its use to evaluate policy and price change scenarios. Even without additional temporal data, the framework can take advantage of the wide extension of the study area to perform, e.g., climate-change scenarios assessment, by adopting a space-for-time substitution approach.”

5. I would suggest to the author(s) to add the "limitation of the study" section before the conclusion. You can write all the methodological limitations (mentioned above, if not addressed properly) and you have also removed several samples (Line 92). Also, you can write why the classification error rate for some of the FS sub-components are much higher, here in this section.

R: We appreciate the recommendation and proceed as suggested. A new section «4.5 Shortcomings of the approach and recommendations for future research» has been added to the MS.

Minor comments:

1. Please include the legends of the graphs in Annex-I for all the features.

R: Done

2. Referring variable WPRIVATE: Isn't there a spatially distributed surface water spatial map? I am not sure about using yes/no type of raster as a feature in spatial prediction for water. Yes, we can use for example for land types (e.g., uplands or lowlands). But for irrigation I am not sure.

R: We understand the concern raised by the Reviewer. However, the only information available that we were able to access was a surface water map. This map is presented in Annex I (supplementary information), in which we highlighted the UAA with direct contact with small private surface water sources, like small streams or ponds (we discarded large rivers and reservoirs, since large rivers are hardly a direct irrigation water source for most farms, and reservoirs in this region are mostly associated with dams of public irrigation systems, whose irrigation areas were considered in the WPUBLIC variable). Thereby, this map only shows UAA where the possibility (but not the verification) of growing irrigated crops from private surface water sources exists ("yes"), regardless of whether it was actually used for irrigation. We have now added an explanation on this in the new section «4.5 Shortcomings of the approach and recommendations for future research» (L554-560).

"(...) The problems observed with variable WPRIVATE may be one such case, as this variable only reported access to small private surface water sources, which are mostly torrential regime in this region, with insufficient water guarantees to encourage investing in irrigation systems, and not taking into account that a significant portion of private irrigation in this region is probably resorting to groundwater sources. This premise, which we could not test due to lack of data, would be worth further investigation, should spatially explicit data on groundwater uptakes becomes available."

3. Less focus on the productivity with the farm system (FS).

R: Productivity (whether land productivity or labour productivity) were not used to define the farming systems, but they were later used to characterize the FS (Table 4), where we present an "intensity" indicator, resulting from an estimate of land productivity (in 10^3 euros per hectare). Therefore, (land) productivity was used in order to characterize the result of the system and not the FS itself.

4. Economic status and economical parameters are not used extensively as the components of FS as agriculture mechanization fully depends on economic parameters.

R: Our data did not provide any information on the mechanization levels of the farms, nor their economic profitability. We tried to circumvent these limitations by estimating the total gross product per land unit (in €/ha UAA) for each farm, following the EU "standard output"

approach (Commission Regulation (EC) No 1242/2008 of 8 December 2008), as explained in L125-127.

“(...) The intensity variable was calculated following the EU “standard output” approach (Commission Regulation (EC) No 1242/2008 of 8 December 2008) by estimating the total gross product per land unit (in €/ha UAA) for each farm.”

Reviewer #2:

The manuscript is well written while using rich data set.

I have following suggestions on the manuscript

R: We thank the Reviewer for revising this paper and we welcome her/his valuable comments and suggestions.

- Present a map of study location

R: We added a map of the study-area in the revised version of our manuscript.

- Variables category (Table 1) is not clear. For example what are the crops under cereal (as rice is also cereal crop). Same with horticultural crops and livestock. For clarification, present the description of the category of the variable in a separate table as supplementary material or present in the text.

R: We understand the Reviewer’s comment. To clarify this issue, we changed Table 1 by adding a description of the crops included in each category, whenever appropriate.

- In Partial dependent plots (Figure 3), present the variable unit in the X - axis for clarification.

R: We significantly changed this section of the MS and merged the two previous figures 2 and 3 into a single figure (current Fig. 3), also following Reviewer #3 suggestions (comment #7). The information on the marginal effects of the drivers on the FS, provided by the partial dependence plots, was also included in Fig. 3, as described in the caption. The partial dependence plots are now presented only in supplementary information (Annex IV).

Reviewer #3:

The authors aim to link the farm-level drivers of the choice of farming system (FS) and predict their spatial distribution at the landscape level in a large geographical area in Portugal dominated by agricultural use. The authors delineate a typology of farming systems (22 no.) and characterize them by several biophysical and socioeconomic variables (assumed to explain the choice of FS). These factors/drivers were used in a random forest-based FS choice model to

identify the critical drivers (and their roles) for all identified farm types. Finally, based on these farm-level drivers, the authors predict the spatial patterns of FS at the landscape scale. The article is well-written and the analytical approaches employed are novel in a sense that they link the farm-level drivers to the landscape level predictions, which is essential for practical purposes (e.g. provision of public goods). However, I have a few comments and suggestions for the authors to improve the readability for the readers having a different disciplinary background. My comments may also help other researchers to replicate the methods of this study in other contexts.

R: We are very grateful to the Reviewer for revising this MS so thoroughly and for conceding the contribution and novelty of the work. Thank you for your comments and suggestions, which have greatly improved the MS.

1. I am not very clear about the statement "There has been, however, a recent surge in the development of proposals to bring together landscape analysis and farming systems (FS) to understand agricultural landscapes, which can establish the FS geography but struggle to explain the drivers behind their spatial distribution" - Is it that no existing report on landscape analysis and farming systems explain such spatial distribution? If yes, I am interested to know why do they struggle to do that - is it linked to the constrained access to data or methodological limitation?

R: We agree that the statement was misleading and therefore re-wrote it as "(...) but do not go into explaining the (...)".

2. The concept of 'farming system' needs to be defined somewhere in the methods section since it has an established technical definition in the related literature. I guess farm types (identified by PCA and Cluster Analysis) are analogous to FS. This needs to be mentioned to avoid speculation and confusion.

R: We agree that a definition of the FS concept used in this study was missing, so we added a paragraph in the Introduction to clarify the concept early in the manuscript. We also revised the Methods section to clarify methodological issues related to the concept whenever needed. (L58-64):

"The FS concept used in this study follows that proposed by Santos et al. (2020), according to which a FS can be defined as a set of farms roughly practicing the same crops and agricultural activities, using similar technological processes and input endowments. A key aspect in this concept is that only variables resulting from farm management decisions are considered, when defining a FS; all variables that may influence these decisions but do not result from them, at least in the short run (e.g. farm size or fragmentation level, climate, slopes, market or policy), should be considered as exogenous to the FS and, therefore, as potential drivers of the FS choice (Silva et al., 2020)."

According to this definition, farm type and farming system are not necessarily analogous: for example, in building a farm typology we would probably use the size of the farm, to distinguish different farm types (e.g. small, medium and large farms); in a typology of farming systems we do not use this variable, as we consider that the size of the farm is, in fact, a driver that expands or restricts the farmer's productive options, that is, the choice of the FS. For this

reason, we do not use the term "farm type" in this paper (except in L133, when referring to the context of Commission Regulation (EC) No 1242/2008 of 8 December 2008).

3. I am curious to know the process of identifying the socioeconomic and biophysical drivers. Is this ad-hoc or selected through an internal review process or by expert consultation? Also, how were they screened? Please mention.

R: Yes, the screening of the potential drivers was mostly based on literature review and the experience of the authors from previous studies. We added a clarification on this at the beginning of section 2.3 – *Socioeconomic and biophysical drivers* (L167-170):

“Potential socioeconomic and biophysical drivers of farming system choice were screened from literature (e.g. Grigg, 2005; Hazell and Wood, 2008; Kristensen et al., 2016; Plieninger et al., 2016; Reboul, 1989; van Vliet et al., 2015) and the authors’ experience from previous studies where similar approaches were applied (Ribeiro et al., 2018, 2014; Silva et al., 2020).”

4. "...by testing different stratified sampling approaches to deal with anticipated unbalanced data" - This is an exciting strategy to optimize the prediction accuracy of the model across FS. But, I think the authors need to explain this in details in the Supplementary Information.

R: Agreed (please, see response to comment #12).

5. Although PDP gives a hint of marginal effects of the drivers on individual FS, I am interested to know how the interactions of drivers (especially powerful drivers) affected the choice of FS and how that could affect the landscape level prediction of the model. Because it is often common that certain combinations of the drivers (e.g. landholding and access to irrigation) emerge as powerful influencing factors for certain FS. Related discussions are placed here and there in the MS, but I would love to see the authors say explicitly how this complexity (i.e. possible interactions) was handled. Or at least provide the readers with an indication of such interactions operative at the farm level. You may also suggest that accommodating these interactions in the landscape level prediction models could be a future scope of research.

R: This is, indeed, an interesting issue. Interaction effects in random forest models has been the subject of discussion among experts and a final conclusion still seems to have not been reached. Some argue that the random forest model, being based on decision trees where variables are analysed sequentially, is therefore able to deal with interactions without having to specify them. Others draw attention to the fact that interactions can be masked by marginal effects, making it impossible to differentiate between interactions and marginal effects. Others still suggest working on interactions in advance, creating new variables in the database, but this can lead to a significant unfolding in the number of predictor variables. In models that are characterized by the high dimensionality of the data, this can increase the complexity of the analysis to undesirable levels.

In this study, we were mostly focused on exploring the direct (marginal) effects of biophysical and socioeconomic drivers on the spatial distribution of FS. However, although exploring the effects of interaction between drivers is beyond the scope of this study, it must be recognized that they are both possible and probable, and thus it should be mentioned in the MS. Therefore,

we left a reference to this in section «4.5 Shortcomings of the approach and recommendations for future research» (L565-568):

“Also, one aspect that has not been explored in the present study and should merit further investigation is the occurrence of interaction effects between drivers. Although the way random forests deal with these effects is still subject to discussion (Wright et al., 2016), its likely existence recommends additional analysis.”

6. Authors have made critical reflections regarding the high error rate in the model for certain FS. Can a list of reasons be prepared for FS with very high error rates so that readers do not run across the text to explore them?

R: The high error rate in some FS is attributed to the existence of factors that are not being controlled by the considered variables. These may include the effect of farmers' individual desires, attitudes or motivations, or her/his socioeconomic profile, which cannot be controlled based on IACS data. It is therefore difficult, if not impossible, to present a list of the potential reasons behind the high error rate observed in some FS, based on the available data. Only in some cases we could speculate on the possible reasons for this mismatch, such as the "Pastures without livestock" system, to which we referred in the MS. Following a suggestion from Reviewer #1, we have gathered references to these issues in the new section «4.5 Shortcomings of the approach and recommendations for future research» and changed the text hoping to make it clearer (L569-585).

“Finally, the fact that the prediction error rate has shown significant disparities across the FS suggests that the choice of some of these FS may be due to effects not measured by the variables examined, including factors related to farmers' desires, attitudes and motivations, or with their socioeconomic profile which, as mentioned above, cannot be assessed on the basis of IACS data. One such case would be the Pastures without livestock system, whose choice is probably mostly determined by the presence of livestock farms in the nearby, with whom the farm can negotiate grazing land renting, rather than by the biophysical characteristics of the farm or its socioeconomic context. On the other hand, FS with lower error rates in the model were those who most depend on the chosen socioeconomic or biophysical factors, such as the Rice, Irrigated cereals and horticulture or Rainfed cereals and oilseed systems (where cereals are an autumn-winter rainfed crop and oilseeds are grown in spring-summer season, often irrigated) that highly depend on irrigation water provided by public irrigation systems in this region. The same applies to the Vineyards system, whose location is highly dependent on the availability of regional labour supply, to meet peaks of labour needs at certain times of the year, related to certain crop operations (e.g. harvesting or pruning). In the present market, policy and technological context, these FS revealed greater dependence on farm structure and “territorial embeddedness” (sensu Cerceau et al., 2018).”

7. I strongly recommend that the authors prepare a summary table or (ideally) a value-added visualization (e.g. heatmap) to show all the critical drivers for 22 FS together. That will be great for the readers to get an overview of variable importance covering all 22 FS instead of scanning through the figures and their discussions. A large number of related figures in the main text may be sent to the Supplementary Information.

R: We agreed with the Reviewer's suggestion and merged the two previous figures 2 and 3 into a single figure (current Fig. 3) showing a heatmap of variable importance of drivers in each FS, and including an information on the marginal effects of the drivers on the FS, provided by the partial dependence plots, as described in the caption. The partial dependence plots are now presented only in supplementary information (Annex IV).

8. I was curious to know (to teach myself) why hexagonal units were taken as the unit of analysis? I see an explanation in the methods. But, my interest is in the usefulness of the shape of the unit. Can any reference be cited here?

R: Hexagonal grids have been widely used in ecology studies exploring landscape metrics (Birch et al., 2007; Grif et al., 2000; Schindler et al., 2008). Compared to square grids, hexagonal grids can be seen as more "natural" shapes, as they potentially reduce bias due to edge effects by decreasing the perimeter-area ratio (which is minimal in the circle, although it does not allow continuous grids to be formed) (Birch et al., 2007; Elkie et al., 1999). We have introduced a brief explanation of this option in the MS (L246-248):

"A hexagonal grid was preferred over a square grid because it is less subject to bias from the edge effects when computing landscape metrics (Birch et al., 2007)."

9. Is it the case that a hexagon where an FS (say FS-X) is predominant will be more accurately predicted if the error rate associated with FS-X in the RF Choice-Model is low? If that is the case, the prediction accuracy is dependent on FS composition in the hexagon. Then, I assume that the hexagons should be marked as 'more predictable' and 'less predictable' for informing the planners. Does marking hexagons as 'less error-prone' add value to the landscape-level prediction and their use in practice? Please think.

R: In our view, several factors may influence the accuracy rate in predicting the landscape pattern in each hexagon. On the one hand, the hit rate on each hexagon will depend on the hit rate of the FS that make it up, as mentioned by the Reviewer. On the other hand, it will also depend on the number of FS in that hexagon (a higher number of different FS may contribute to lower the hit rate). Not to mention the influence that the size of the hexagons can also have on the predictive ability of the model, as noted by the Reviewer in his comment #10.

Additionally, it should be noted that the differences were not calculated at the farm level, but at the level of the hexagon as a whole, as described in L249-253. Therefore, the FS hit rate may not evolve in parallel with the landscape hit rate. For example, in a hypothetical landscape with only two competing FS, each occupying 50% of the total area, if the error rate in the prediction of the FS is 100%, the success rate in predicting the landscape pattern will be also 100%, as the same 50-50% ratio is maintained in the hexagon's FS composition. This resembles to the *fuzzy kappa* statistic approach, which is often used for pattern recognition in raster map comparisons (which is not, however, our case) (Hagen-Zanker, 2009; Hagen, 2002; Visser and De Nijs, 2006). Fuzzy map comparisons stand out as they resemble the way human observers compare maps (e.g. two chess boards on top of each other result in a complete mismatch of colours if one of them is rotated 90 degrees; however, for a human observer the patch pattern will be the same in both boards).

It is therefore not easy to state how assigning a “predictability-tag” to individual hexagons could be useful for planning purposes, given the uncertainty about its actual meaning. The approach, as proposed, is particularly suited for its ability to establish a direct link between the farm and the landscape scales, enabling to easily assess the impact of local (farm) policies on landscape patterns. It can thus be used to find the right balance between farmers’ private goals and societal demand for agricultural public goods by, e.g., simulating policies that encourage farmers to adopt particular high nature value farming systems.

10. "For agricultural landscape planning focused on agroecosystem services provision, this may be the right scale of analysis" - I appreciate this insight. However, how can we say this without experimenting with the size of the hexagon? I do not expect that everything will be done in a single study. However, identifying the appropriate scale for landscape-level prediction through the RF choice-model may be mentioned as the scope of future research.

R: We agree that this is a relevant research issue and it should be further investigated. We have therefore inserted a reference to this question in the new section "4.5. Shortcomings of the approach and recommendations for future research" (L561-564):

“Another issue deserving further investigation concerns the dimension of the grid of landscape analysis units. It is possible that the size of these units (i.e. the hexagons, in the current case) influences the accuracy of the model, so future investigation focused on determining its optimal size could prove to be of high value.”

11. I surmise the assumed equivalence of farm-level decision for all 22 Farm Types. Can the choice of pasture and olive be judged with the same set of indicators? I understand the random allocation of variables (at individual nodes of the tree) during the random forest exercise, but it is tricky to equate the existence of a rice field with an eight-year-old orchard. I request the author to reflect on this.

R: On a first note, it is important to bear in mind that the model is choosing between farming systems (FS) and not between crops or land uses/covers. Each FS is represented by the centroid of the group (cluster) that includes all the farms that were classified in this class (FS). The land use/cover of this synthetic “representative farm” assumes the average land use/cover of all farms in the same class. Therefore, the question of choosing between pasture or olive groves, or between rice and orchards, does not arise as such.

The estimation of a random forest model involves training hundreds of individual decision trees, each one using a slightly different (random) subset of the observations, and selecting for each splitting node the best predictor variable from a (random) subset of the total predictors in the dataset. The final prediction will be the mode of the predictions from all individual trees.

Using a random forest model in practice, to classify a new observation (i.e. to make a prediction), involves making this observation go through all the nodes of each tree, until reaching a terminal node (or "leaf") that will assign it with a single class of the dependent variable (i.e. a FS in our case). After going through all the trees in the forest, this observation will be assigned to the class verified in most of the trees (the most “voted” FS).

Therefore, each observation will be classified in a completely independent way, based only on the values of its own predictor variables (i.e., based on its intrinsic characteristics), and without any relation to the values of the other observations. Observations with similar characteristics will thus tend to be classified in the same category and, as these characteristics result from spatial components, they will also tend to be spatially arranged.

So, the FS assigned by the model to each farm depends only on the characteristics (biophysical and socioeconomic) of that farm, without any relation to its previous use. The fact that the model is able to predict correctly most of the times just means that, in fact, these biophysical and socioeconomic factors have been effectively taken into account by farmers in their production decisions.

12. Details of the RF model

I prefer the authors to detail the RF exercise to enhance the reproducibility of the Methods. The authors do not detail their data pretreatments nor the implementation of the machine learning methods. Were the comparisons made 'on equal footing'? Can the author describe how they pretreated the data - this is important. One of the tricky bits about implementing machine learning is the data pretreatments and the data engineering that is required before fitting the models - I suggest the authors to clearly state what they did (for instance, any filtering? transformations? centering? scaling? etc.). Then, when they fit the models, how did they optimize? For RF what are the hyperparameters (the number of trees, the minimum number of samples in a leaf node, the minimum number of samples required to split an internal node, etc.) they used and how they chose those hyperparameters for RF? It is insufficient to state that the "randomForest" package was used. I request the authors to describe all preprocessing and hyperparameter tuning and their effects in the modelling, etc. I do not mean that they should include textbook descriptions of the algorithms, but instead, they should detail their specific implementation, optimizations etc., maybe in the Supplementary Information).

R: We agree with the Reviewer and added a new annex in supplementary information (Annex II) describing the procedures followed in the parametrization of the random forest model. Accordingly, some of the more technical parts of the text in the original MS, describing the random forest model, have now been moved to this annex.

Having said all these things, I must submit that every aspect of a complex phenomenon cannot be addressed (and thus controlled) in a single study. No one can do that/ expect that. Many of my suggestions are to address the raised issues in the MS either by rationalizing them or by mentioning them as the limitation/future scope of research. I congratulate the authors on coming up with a novel analytical approach that links farm-level drivers to the landscape-level prediction, something more acceptable to the policymakers. Overall, I have a positive impression of the paper and believe that it can make a significant contribution to the existing literature.

Cited references:

Birch, C.P.D., Oom, S.P., Beecham, J.A., 2007. Rectangular and hexagonal grids used for

observation, experiment and simulation in ecology. *Ecol. Modell.* 206, 347–359.
doi:10.1016/j.ecolmodel.2007.03.041

Elkie, P., Rempel, R., Carr, A., 1999. Patch analyst user's manual: a tool for quantifying landscape structure, Northwest. ed, Ont. Min. Natur. Resour. Northwest Sci. & Technol. Ontario Ministry of Natural Resources Northwest, Ontario.

Grif, J.A., Martinko, E.A., Price, K.P., 2000. Landscape structure analysis of Kansas at three scales 52.

Hagen-Zanker, A., 2009. An improved Fuzzy Kappa statistic that accounts for spatial autocorrelation. *Int. J. Geogr. Inf. Sci.* 23, 61–73. doi:10.1080/13658810802570317

Hagen, A., 2002. Multi-method assessment of map similarity. 5th Agil. Conf. *Geogr. Inf. Sci.* 1–8.

Schindler, S., Poirazidis, K., Wrbka, T., 2008. Towards a core set of landscape metrics for biodiversity assessments: A case study from Dadia National Park, Greece. *Ecol. Indic.* 8, 502–514. doi:10.1016/j.ecolind.2007.06.001

Visser, H., De Nijs, T., 2006. The map comparison kit. *Environ. Model. Softw.* 21, 346–358. doi:10.1016/j.envsoft.2004.11.013

Wright, M.N., Ziegler, A., König, I.R., 2016. Do little interactions get lost in dark random forests? *BMC Bioinformatics* 17, 145. doi:10.1186/s12859-016-0995-8

1 Explaining farming systems spatial patterns: a farm-level choice 2 model based on socioeconomic and biophysical drivers

3

4 Abstract

5

6 **CONTEXT:** Efforts to bring together landscape analysis and farming systems have failed to explain the
7 drivers behind their spatial distribution. Since agricultural landscapes are an outcome of farmers'
8 decisions, understanding the role of socioeconomic and biophysical drivers of such decisions is
9 essential for policy-making targeting landscape-level provision of public goods and ecosystem services
10 from agriculture.

11 **OBJECTIVE:** Aiming to better understand the role of these drivers, we focused on a region dominated
12 by agricultural use, with extensive variability in biophysical and socioeconomic conditions. A typology
13 of farming systems was derived from spatially explicit farm-level data provided by the Portuguese
14 agency responsible for Common Agricultural Policy payments, for 2017. Farms were thoroughly
15 characterized through relevant biophysical and socioeconomic variables considered as potential
16 drivers of farming systems.

17 **METHODS:** A random forest approach was used to develop a farming system choice-model, dependent
18 on those biophysical and socioeconomic variables. Variable importance measures and partial
19 dependence plots were used to explore the role of these variables in explaining the spatial distribution
20 of farming systems and to predict spatial patterns at the landscape scale.

21 **RESULTS AND CONCLUSIONS:** Results showed that both biophysical and socioeconomic drivers play a
22 significant role in the spatial distribution of most agricultural systems. Its importance, however, varies
23 significantly across farming systems, being crucial for some and almost irrelevant for others. Farm size
24 and climate have proved to be the most relevant drivers for most farming systems. Overall, our
25 approach proved to be quite accurate in predicting patterns of farming systems at the landscape scale.

26 **SIGNIFICANCE:** The proposed framework has shown great potential as a tool to support information-
27 based policy design to improve agricultural landscape planning, by linking farm-level management
28 decisions with the provision of socially valued public goods from agriculture, perceived at the
29 landscape-level.

30

31 1. Introduction

32 Agriculture is a dominant land use in many parts of the world, resulting from human interaction with
33 nature over time. This interaction is mostly regulated by two main types of drivers: biophysical
34 (climate, soil, topography...) and socioeconomic (farm structure, characteristics of farmers, markets,
35 policies...). The way each of these drivers affects agricultural landscapes has attracted the interest of
36 researchers (Grigg, 2005; Hazell and Wood, 2008; Kristensen et al., 2016; Plieninger et al., 2016; van
37 Vliet et al., 2015), but many unanswered questions still persist (Plieninger et al., 2016; Wilson, 2009).
38 Advancing knowledge about the role played by each of these factors in shaping agricultural landscapes
39 can thus improve our understanding of human/environment interactions, allowing to anticipate farm
40 management decisions and supporting evidence-based public intervention (Levers et al., 2016; van de
41 Steeg et al., 2010).

42 Such issues have recently been raised in the context of the provision of public goods and
43 agroecosystem services in general, including biodiversity conservation (Landis, 2017; Schaller et al.,
44 2018). Much literature resort to aggregated data concerning land use or to agriculture intensification
45 or specialization indicators, privileging landscape-dynamics analysis over landscape regional
46 differentiation, and seldom take the farm as the unity of inquiry (Debolini et al., 2018; Ruiz-Martinez
47 et al., 2015). There has been, however, a pressing need to the development of approaches linking
48 landscape analysis and farming systems (FS) to understand agricultural landscapes, which are able to
49 establish the FS geography but do not go into explaining the drivers behind their spatial distribution
50 (Andersen, 2017; Benoît et al., 2012; Martel et al., 2019; Rizzo et al., 2013; van de Steeg et al., 2010).
51 Indeed, considering the mismatch between the farm-scale, where management decisions take place,
52 and the landscape-scale, where ecosystem services are perceived, landscape analysis can greatly
53 benefit from a deeper understanding of the factors that influence farm management decisions. Thus,
54 understanding the multiple production decisions of adjacent farmers and combining these decisions
55 at the landscape-scale is key to explain the landscape mosaic and the ecological disturbance regimes
56 (fire, grazing, ploughing...) that shape the habitats of wild species and the provision of diverse
57 ecosystem services.

58 The FS concept used in this study follows that proposed by Santos et al. (2020), according to which a
59 FS can be defined as a set of farms roughly practicing the same crops and agricultural activities, using
60 similar technological processes and input endowments. A key aspect in this concept is that only
61 variables resulting from farm management decisions are considered, when defining a FS; all variables
62 that may influence these decisions but do not result from them, at least in the short run (e.g. farm size

63 or fragmentation level, climate, slopes, market or policy), should be considered as exogenous to the
64 FS and, therefore, as potential drivers of the FS choice (Silva et al., 2020).

65 To explain the spatial distribution of FS, distinct groups of drivers can be considered according to
66 distinct disciplinary perspectives or theoretical approaches. The analysis of farm biophysical
67 endowments to explain spatial patterns of FS has largely been explored by geography and
68 geo-agronomy (Deffontaines, 2004; Deffontaines et al., 1995; Grigg, 2005; Lacoste et al., 2018).
69 Climate, soil, and slope are often considered to establish a range of restrictions to the choice of the
70 farming system. But FS are also dependent on farmland structure and social context. Farmland
71 structure covers an ensemble of constraints such as farm size, fragmentation and spatial composition
72 which potentially restrict farmers decisions (Grigg, 2005; Latruffe and Piet, 2014; Reboul, 1976; Ribeiro
73 et al., 2018). The influence of territorial socioeconomic context on FS location may be grounded in the
74 notion of local embeddedness, supported by local sociocultural, demographic and economic structures
75 (Canadas and Novais, 2014; Debolini et al., 2018).

76 Using farm-level data collected in 2017 in a large-scale study area, we developed an innovative
77 methodological approach to: 1) derive a spatially-explicit FS typology; 2) assess the role of
78 socioeconomic and biophysical factors in explaining the spatial distribution of those FS; 3) assess the
79 extent to which we can predict FS patterns based on biophysical and socioeconomic variables. Results
80 were used to discuss the role of these drivers on the choice of the FS and their potential to predict
81 landscape patterns, seeking to draw conclusions to better inform policy design for landscape-level
82 provision of public goods from agriculture and prediction of landscape patterns in face of biophysical
83 or socioeconomic changes.

84

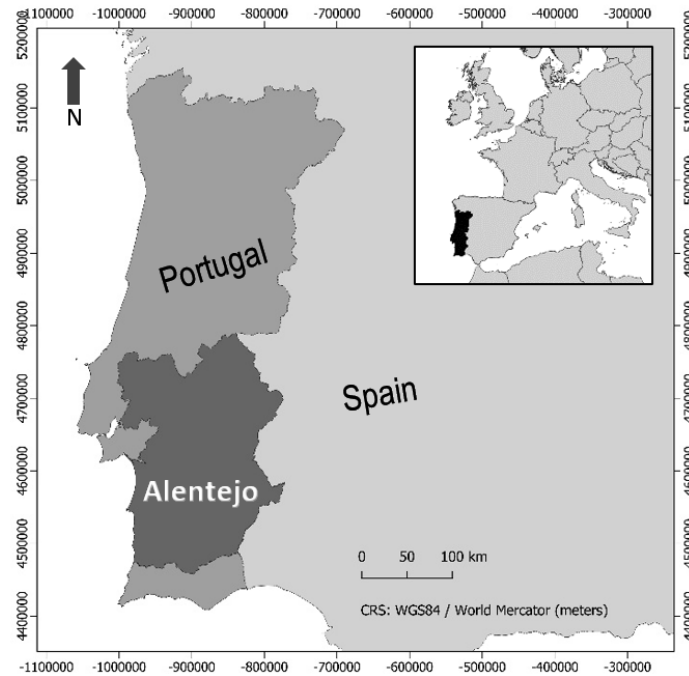
85 2. Methods

86 2.1. Study area

87 The study focused on the Alentejo region, in southern Portugal (Fig. 1), corresponding to the EU
88 statistical region PT18, at the NUTS2 level (Nomenclature of Territorial Units for Statistics). Covering
89 about 31,551 km² (ca. 1/3 of Portugal), the region has a Mediterranean climate, with hot dry summers
90 and mild rainy winters. The annual average temperature is about 16.3°C, ranging from 9.9°C to 23.4°C
91 in January and August, respectively, and the total annual rainfall is about 619 mm, largely concentrated
92 in the rainy season (approx. October to March). The relief is predominantly smooth (47% of the land

93 with slope < 5%; but 14% with slope > 15%), with few mountain areas (average altitude is 176 m a.s.l.,
94 ranging from 0 to 1020 m).

95



96

97 Fig. 1 - Location of the study area in the Alentejo region (NUTS2), Portugal

98

99 According to the latest agricultural census in Portugal (2009), the utilized agricultural area (UAA) in
100 Alentejo (NUT 2) was then ca. 2.2 million hectares, covering almost 70% of the region and making it
101 the dominant land use. Official statistics report that in 2016 the utilized agricultural area (UAA) was
102 dominated by permanent pastures (64%), followed by annual crops (24%) and permanent crops (11%).
103 Cereals, forages and olive groves were the main crops, with roughly equal shares of 8% in total UAA
104 (making ca. 70% of the UAA excluding permanent pastures). Nearly 40% of the UAA is under the canopy
105 of scattered trees, mainly cork and holm oaks (*Quercus suber* and *Q. rotundifolia* respectively),
106 originating an agroforestry system locally named "montado", which is largely acknowledge for its high
107 nature value (Ferraz-de-Oliveira et al., 2016). Cropland in these undercover areas are mainly
108 permanent pastures (70%) and annual crops (30%). Most of the UAA is rainfed (ca. 90%) and irrigated
109 areas are mostly located within state-promoted irrigation systems, often depending on large dams.
110 The region is dominated by large holdings, with almost 90% of the UAA in farms with more than 50 ha.

111

112 2.2. Farming systems identification

113 To build a farming systems typology for the study area we used data from the EU Integrated
114 Administration and Control System (IACS) for 2017, associated with spatially explicit farm parcel data
115 from the Land Parcel Identification System (LPIS), provided by the Portuguese agency responsible for
116 Common Agricultural Policy (CAP) payments. These data are collected on a yearly basis from farmers
117 declarations when applying for CAP payments and its usefulness for FS research has been
118 demonstrated by previous studies (Lomba et al., 2017; Ribeiro et al., 2018, 2016, 2014).

119 The raw dataset identified 26,648 CAP beneficiaries in the study area, covering a total of 2,221,816 ha
120 distributed over 208,338 parcels which, in turn, included 560,213 subparcels for which land use/crop
121 cover was described. Livestock declared by each beneficiary was also provided, describing species
122 composition, gender, age groups and an indication of whether they were kept in stables or grazing.

123 First, all parcels declared by the same CAP beneficiary were taken as a single farm. However, we found
124 that some beneficiaries reported very scattered parcels, sometimes separated by hundreds of
125 kilometres, where the farm concept (as an agricultural management unit) would not apply. In these
126 cases, we decided to regroup these parcels into new (sub)farms by forcing the distance between them
127 not to exceed 25 km, which increased the total number of farms to 28,739. This decision also helped
128 to narrow down the range of biophysical variability within each farm, and thus to better link farm units
129 to their biophysical context, described in the next section. We also discarded farms with total area
130 equal or below 2 ha (4409 farms, representing less than 1% of total UAA) because the land use in
131 smaller farms is likely to be highly sensitive to crop rotations, which cannot be properly captured with
132 one-year data.

133 The raw data included 129 land use/cover categories, which were simplified by aggregation into
134 broader categories, while maintaining the distinction between irrigated and rainfed crops, when
135 applicable (e.g., irrigated and rainfed cereals). We also included two variables describing the
136 proportion of the UAA under the cover of cork and holm oaks, respectively, because their presence is
137 prone to influence farm management, as the first is a major source of income for farmers (cork
138 production) and the later provides shade and food (acorns) to livestock grazing, in addition to valuable
139 firewood. These two variables were computed on a geographical information system (GIS)
140 environment by intersecting the farms map (derived from the LPIS spatial data) with digital information
141 on cork and holm oak distribution and computing, for each farm, the share of the UAA covered with
142 both land cover classes.

143 Livestock numbers were converted into livestock units (LU) using EU standard conversion factors, and
 144 these were used to describe the percentage composition of livestock by species, as well as livestock
 145 density in each farm. Thus, a set of 28 variables was defined to characterize the land use/cover and
 146 livestock patterns for each farm (Table 1).

147 A principal component analysis (PCA) was performed on a correlation matrix of these 28 variables to
 148 reduce variable redundancy and the principal axes with eigenvalues above 1 entered a hierarchical
 149 cluster analysis (Ward method) to derive the FS typology. The number of clusters to retain was decided
 150 based on a visual analysis of the dendrogram and on expert knowledge of the study area.

151 To help interpreting the resulting FS, we calculated three variables indicating the level of agricultural
 152 intensity, specialization and dependence on labour. The intensity variable was calculated following the
 153 EU “standard output” approach (Commission Regulation (EC) No 1242/2008 of 8 December 2008) by
 154 estimating the total gross product per land unit (in €/ha UAA) for each farm. The specialization variable
 155 was computed as the highest proportion of standard output from a single farm activity. The labour
 156 indicator aims to differentiate the FS based on their specific labour needs, in annual work units per
 157 land unit (AWU/ha UAA). Due to data limitations, we had to resort to official statistics on the “EU farm
 158 typology by economic size and type of farming” (in the sense of the above-mentioned legal text), at
 159 NUT2 level (Alentejo) for the year of 2013, from which we extracted the number of annual work units
 160 per hectare for each farm type, to be directly associated to each of the resultant FS on a similarity base.
 161 Thereby, this indicator was not computed at farm level, but directly at FS level.

162

163 Table 1 – Summary statistics for the land use/cover and livestock farm characterization variables (n =
 164 24313 farms)

Variable	Mean	SD
<i>Land use/cover variables (proportion of total UAA)</i>		
Rice (both Indica and Japonica)	0.012	0.1
Cereals Irrigated (corn, wheat, oats, barley, triticale)	0.018	0.104
Cereals rainfed (wheat, corn, oats, barley, rye and triticale)	0.056	0.165
Orchards (orange, apple, plum, fig, loquat, cherries, blackberry, raspberry)	0.013	0.078
Forages Irrigated (ryegrass, lucerne, silage maize, sorghum, vetch)	0.006	0.051
Forages Rainfed (ryegrass, oats, corn, sorghum, lupine)	0.049	0.153
Horticultural (potatoes, carrots, onions, cabbages, beans, chickpeas)	0.017	0.089
Industrial horticulture (tomato and pepper)	0.011	0.092
Oilseeds (sunflower and rapeseed)	0.01	0.067

Pastures (temporary grass and permanent grasslands)	0.511	0.41
Fallows	0.043	0.146
Olive groves Irrigated	0.034	0.156
Olive groves Rainfed	0.171	0.291
Vineyards	0.034	0.145
Walnuts and almond trees	0.003	0.048
Stone pine	0.009	0.079
Other dry fruits (hazelnut, chestnut, pistachios, carob)	0.001	0.019
Cork oak cover	0.149	0.265
Holm oak cover	0.111	0.229
<i>Livestock variables (proportion in total LU)</i>		
Cattle grazing	0.168	0.34
Cattle stabled	0.003	0.04
Fattening cattle grazing	0.018	0.054
Fattening cattle stabled	0.002	0.037
Sheep grazing	0.205	0.386
Goat grazing	0.024	0.131
Dairy cows	0.004	0.047
Pigs grazing	0.008	0.076
Livestock density (LU/ha UAA) (includes all farm animals, added-up in LU)	0.526	3.506

165

166 2.3. Socioeconomic and biophysical drivers

167 Potential socioeconomic and biophysical drivers of farming system choice were screened from
168 literature (e.g. [Grigg, 2005](#); [Hazell and Wood, 2008](#); [Kristensen et al., 2016](#); [Martel et al., 2019](#);
169 [Plieninger et al., 2016](#); [Reboul, 1989](#); [van Vliet et al., 2015](#)) and the authors' experience from previous
170 studies where similar approaches were applied ([Ribeiro et al., 2018, 2014](#); [Silva et al., 2020](#)).
171 Subsequently, each farm was characterized according to a set of socioeconomic and biophysical
172 variables thus identified, considered as potential drivers of FS spatial patterns ([Table 2](#)). These
173 variables vary spatially but are mostly constant over time (at least for the time scale of most farm
174 management decisions).

175 Socioeconomic variables included seven farm structure variables (farm and block size, farm
176 fragmentation and dispersion, access to public and private water sources for irrigation, nature
177 conservation constraints on farm use), and six local context variables computed from official statistics
178 at the administrative parish level (one demographic variable, population density, and five agricultural
179 variables, e.g. AWU availability or the share of rented UAA); all farms in the same parish where

180 assigned the same value in these variables; when farms had areas in more than one parish, these
 181 variables were computed through average-weighting by farm-area shares in each parish. Biophysical
 182 variables included three climatic variables (describing temperature and precipitation), eight soil quality
 183 variables (describing soil depth, texture and pH) and three topographic variables (slope categories).
 184 (Table 2).

185 Values for explanatory variables were derived for each farm using a GIS (maps for explanatory variables
 186 are provided in supplementary information, Annex I). Farms with missing values resulting from map
 187 mismatches were discarded, dropping the number of valid observations to 23,416 farms.

188

189 Table 2 – Summary statistics for the socioeconomic and biophysical drivers (n = 23416 farms)

Variable	Description	Mean	SD	Min	Max
<i>Socioeconomic variables – farm structure variables</i>					
F SIZE	Farm size – Total UAA (ha) (1)	84.09	184.46	2.01	7191.16
BLKSIZE	Average farm-block size (ha) (1)	23.15	45.37	0.20	1109.93
JANUS	Januszewski index (adimensional) (1) (2)	0.65	0.23	0.13	1.00
BLKDIST	Average area-weighted block distances to farm centroids (m) (1)	1571.00	2128.00	0.00	56951.00
WPRIVATE	Access to water from private ponds or small streams (yes=1; no=0) (5)	0.16	0.37	0.00	1.00
WPUBLIC	Proportion of UAA in public irrigation systems (6)	0.15	0.31	0.00	1.00
NATURE	Proportion of UAA included in areas classified for nature conservation (7)	0.22	0.39	0.00	1.00
<i>Socioeconomic variables – local socioeconomic variables</i>					
INCAGRI	Proportion of farms where agriculture is the main household income source (3)	0.23	0.14	0.00	0.84
INCOTH	Proportion of farms where household income is mostly from outside the farm, but not pensions (3)	0.26	0.07	0.00	0.67
PDENS	Population density (inhabitants/km ²) (4)	32.5	77.8	0.89	1084.24
AWU	Number of annual work units (AWU) per km ² of total parish area (3)	1.96	1.70	0.21	17.95
AWU hired	Proportion of hired work in total labour (3)	0.26	0.15	0.00	0.93
RENT	Proportion of rented land in total UAA (3)	0.18	0.12	0.00	1.00
<i>Biophysical variables</i>					
TMIN	Average minimum temperature in the coldest month 1970-2000 (°C) (8)	4.71	0.59	3.01	8.40

TMAX	Average maximum temperature in the warmest month 1970-2000 (°C) (8)	31.56	1.95	20.24	35.68
PREC	Average annual rainfall 1970-2000 (mm) (8)	592.89	107.28	376.83	1195.51
SDEPTH	Soil depth (cm) (5)	52.74	29.80	0.00	150.00
SMOOTH	Proportion of UAA with smooth slopes (<5%) (5)	0.51	0.32	0.00	1.00
MODERATE	Proportion of UAA with moderate slopes (5-16%) (5)	0.38	0.24	0.00	1.00
STEEP	Proportion of UAA with steep slopes (>16%) (5)	0.11	0.19	0.00	1.00
HEAVY_S	Proportion of UAA with heavy texture soils (5)	0.33	0.37	0.00	1.00
MEDIUM_S	Proportion of UAA with medium texture soils (5)	0.42	0.38	0.00	1.00
LIGHT_S	Proportion of UAA with light texture soils (5)	0.24	0.36	0.00	1.00
VERYACID	Proportion of UAA with very acid soils (pH<5) (5)	0.27	0.33	0.00	1.00
ACID	Proportion of UAA with acid soils (5<pH<6) (5)	0.41	0.38	0.00	1.00
NEUTRAL	Proportion of UAA with pH neutral soils (6<pH<7) (5)	0.21	0.30	0.00	1.00
ALKALINE	Proportion of UAA with alkaline soils (pH>7) (5)	0.11	0.24	0.00	1.00

190 Sources: (1) Computed from LPIS data; (2) Farm spatial fragmentation index, varying from 0 to 1 with higher values
191 indicating a higher degree of farmland consolidation (Januszewski, 1968); (3) Agricultural census 2009 - parish level; (4)
192 Population census 2011 - parish level; (5) EPIC WebGIS Portugal (<http://epic-webgis-portugal.isa.ulisboa.pt/>); (6) DGADR -
193 Direção-Geral de Agricultura e Desenvolvimento Rural (<http://sir.dgadr.gov.pt/expl-alentejo>); (7) ICNF – Instituto de
194 Conservação da Natureza e das Florestas (<http://www2.icnf.pt/portal/pn/ap>); (8) IPMA - Instituto Português do Mar e da
195 Atmosfera (<https://www.ipma.pt/pt/oclima/normais.clima/>)

196

197 2.4. Model design

198 We developed a random forest FS choice model to explore the farm-level relationships between the
199 typologies of FS derived from cluster analysis and the socioeconomic and biophysical variables.
200 Random forest is a popular machine learning method that can be used both for regression and
201 classification, and is well-suited for high dimensional data (Strobl et al., 2009). Random forest use
202 bootstrap and aggregation (bagging), building multiple decision trees based on random subsets of the
203 data and using a random subset of predictor variables candidates for each node, in each decision tree
204 (Liaw and Wiener, 2002). On a classification problem, each observation is assigned to a class according
205 to the majority of votes from all trees. Both the number of trees and the number of predictor variables
206 sampled for each node are user-defined and can be used to tune the model. The mean out-of-bag
207 (OOB) error rate computed across all trees provides a measure of model prediction accuracy (Breiman,
208 2001). Random forests have been widely used in many scientific fields and have proved to be one of
209 the best machine learning techniques currently available, including for predictive modelling of spatial
210 and spatio-temporal data (Hengl et al., 2018).

211

212 2.4.1. Explaining spatial distribution of farming systems

213 Since we were firstly interested in exploring causal theories on the main drivers of FS spatial
214 distribution, rather than using the model to make predictions on new data (e.g. to assess scenarios of
215 policy or climate change), we tuned the model to optimize its average prediction accuracy across FS,
216 rather than maximizing the overall prediction power, by testing different stratified sampling
217 approaches to deal with anticipated unbalanced data (high variance in group sizes) (see details of
218 model parametrization in supplementary information – [Annex II](#)). At this stage, model overfitting
219 should not be an issue, since the focus was on explaining our training data, rather than the
220 generalization of the model ([Shmueli, 2010](#)).

221 With this modelling outset, all FS are assumed to be competing simultaneously for each farm and the
222 choice is made dependent only on variables that vary in space, while keeping constant the effect of
223 temporal variables (such as prices or policies). The effect of these temporal variables on the choices
224 observed in the study year cannot be estimated, as we only have one observation on FS choice for
225 each farm, that is: the choice observed in the study year 2017.

226 We used variable importance measures to assess the relevance of each predictor variable in the model
227 and their marginal effect on each FS was examined using partial dependence plots ([Friedman, 2001](#))
228 (supplementary information – [Annex II](#)). We investigated the shape of the partial dependence plots
229 fitted functions for each class of the dependent variable (that is, for each FS) to infer their role as
230 drivers or constraints for each FS. In addition, we computed the correlation coefficient between the
231 level of farming intensity characterizing each FS with the corresponding prediction accuracy rate
232 obtained by the model, to test the hypothesis of a positive relationship between the levels of this
233 indicator and the degree of FS dependence on socioeconomic and biophysical drivers.

234 All statistical analyses were carried out in R 3.4.1 ([R Development Core Team, 2017](#)).

235

236 2.4.2. Predicting spatial patterns of farming systems

237 On a following step, we focused on exploring the predictive capacity of the model in the choice of the
238 FS, based on the socioeconomic and biophysical variables described above. Since we were mostly
239 interested in predicting FS choice at the landscape-scale rather than at farm-scale, taking into account
240 the importance of landscape patterns for biodiversity and public goods delivery, we focused the
241 analysis on the model's ability to predict FS spatial patterns at a scale comparable to that of the

242 landscape (Andersen, 2017). For this purpose, the study area was divided into a random network of
243 hexagons of about 54,125 ha each, corresponding to a hexagon apothem of 12.5 km which was chosen
244 with reference to the 25 km threshold used to define the farms. These hexagons were then used as
245 analysis units to compare, for each hexagon, the percentage distribution of the UAA by FS in the
246 observed situation with that predicted by the model. A hexagonal grid was preferred over a square
247 grid because it is less subject to bias from the edge effects when computing landscape metrics (Birch
248 et al., 2007). We rejected all hexagons with more than 66% of the area outside the LPIS data, due to
249 low significance for this purpose. In each hexagon, we calculated the difference between the observed
250 and predicted UAA shares for each FS and computed the half-sum of their absolute values. The average
251 of these results across all hexagons was interpreted as an estimate of the percentage of accuracy
252 obtained in model predictions, that is, the capacity of the model to predict spatial patterns of FS
253 composition at the landscape-scale. In addition, we also computed the determination coefficient (r^2)
254 between the observed and predicted values in each hexagon, taking its mean as a measure of the
255 quality of fit of the model. Model predictions were obtained by running the model on a random test-
256 set of the data with ca. 1/3 of the observations (farms), after estimating it in a train-set with the
257 remaining 2/3.

258

259 3. Results

260 3.1. Farming systems typology

261 A solution of 30 groups, representing farming systems, was selected from the cluster analysis. As some
262 groups included only a very small number of observations (farms), we anticipated potential problems
263 in the estimation of the predictive model and so we decided to eliminate groups with less than 0.7%
264 of the total number of observations, an arbitrary threshold mostly based on expert judgement. This
265 led to the removal of 8 non-representative FS, comprising 613 farms accounting for 3.1% of total UAA,
266 which were discarded for further analysis. Consequently, the final number of FS was set at 22 (Table
267 3).

268 By chance, these FS resulted equally divided into livestock-oriented systems and crop-oriented
269 systems. Both groups include similar shares in number of farms (51.5% and 48.5%, respectively),
270 although farms in livestock-oriented systems cover a much larger share of total UAA (78.2%) denoting
271 they are larger farms, on average.

272 Within the livestock systems, six are oriented to sheep, three to cattle, one to goats and one is mixed
273 with cattle and sheep. Among the six sheep-oriented systems, two are agroforestry grazing systems,
274 one associated with cork oak and the other with holm oak, a third one is related with open land
275 pastures, a fourth sheep system is mainly dependent on forage crops, the fifth depends both on
276 permanent pastures and forage crops, and the last is mostly a mixed-system combining rainfed olive
277 groves with sheep grazing. The three cattle systems also include two agroforestry grazing systems with
278 permanent pastures under the canopy of cork and holm oaks, respectively, and a third one depending
279 mainly on forage crops. The mixed cattle-sheep system is highly dependent on irrigated forages and
280 the last livestock-oriented system is the goat system, which is also a pasture-dependent grazing system
281 (Table 3).

282 Among the crop-oriented systems, five are dedicated to permanent crops, four to annual crops and
283 the last two refer to special situations, one including farms without livestock but with almost all UAA
284 under pasture, probably yearly rented to neighbours with cattle, and the other encompassing farms
285 with almost all UAA set to fallow. The permanent crops systems included two systems dedicated to
286 olive groves, one of which was irrigated and the other rainfed, one to vineyards, another to fruit trees
287 and the last one to stone pines (for pine nut production). The annual crops systems included two
288 rainfed systems, one dedicated to cereals and the other to cereals and oilseeds, one dedicated to
289 irrigated cereals and horticultural crops, and the last one to rice production (Table 3).

290 The average farming intensity across the 22 FS is about 1650 €/ha, with the Fruit trees system as the
291 most intensive, reaching ca. 12600 €/ha, and 15 systems below 1000 €/ha. Agricultural specialization
292 is relatively high, with more than half of the FS earning more than 80% of their standard output from
293 a single activity. Average farm specialization is higher in crop systems than in livestock systems (85%
294 and 74%, respectively), where most systems earn more than 90% from a single activity. Average labour
295 needs are also higher in crop systems than in livestock systems (0.039 and 0.004 AWU/ha, respectively,
296 i.e. nearly 10 times more), with a maximum of 0.259 AWU/ha found in the Irrigated cereals and
297 horticultural crops system and a minimum of 0.003 AWU/ha found in systems Pastures without
298 livestock and Fallows (Table 4).

299 Average farm size varies significantly across FS, with values going from ca. 11 ha in both rainfed olive
300 grove systems (with and without sheep) until over 200 ha, in cattle grazing – HO and – CO systems
301 (288 and 249 ha, respectively) (Table 4).

302 Almost 1/3 of all farms are included in only three FS, all with more than 2500 farms (systems Rainfed
303 olive groves, Pastures without livestock and Cattle grazing – CO). However, nearly 1/3 of total UAA is

304 concentrated in one single FS, the Cattle grazing – CO. The three cattle-oriented FS comprise more
305 than half of the total UAA (53.3%) ([Table 4](#)).

306

307 Table 3.a – Farming system description – Land cover composition (average values in proportion to the total UAA; values under 0.01 are omitted; values above 0.5
 308 are in bold)

Farming system	Rice	Cereals Irrigated	Cereals rainfed	Forages Irrigated	Forages Rainfed	Horticultural	Industrial horticulture	Oilseeds	Fallows	Pastures	Fruit trees	Olive groves Irrigated	Olive groves Rainfed	Vineyards	Walnuts and almond	Stone pine	Other dry fruits	UAA under Cork oak	UAA under Holm oak
Cattle grazing – CO*					0.039					0.884		0.026						0.329	0.082
Cattle grazing – HO*			0.020		0.035					0.906		0.020						0.048	0.590
Cattle grazing – forages		0.011	0.143		0.239	0.015		0.020	0.027	0.388		0.045	0.093					0.049	0.086
Grazing goats					0.027					0.923			0.023					0.314	0.207
Mixed Cattle and sheep - Irrigated forages		0.034	0.037	0.444	0.112	0.019		0.016	0.018	0.219	0.018	0.015	0.049					0.066	0.028
Sheep grazing – CO*					0.012					0.936		0.022			0.012			0.686	0.022
Sheep grazing – HO*			0.029		0.025				0.013	0.891		0.032						0.051	0.641
Sheep grazing - pastures			0.014		0.021					0.852		0.085						0.141	0.048
Sheep grazing - pastures and forages			0.169		0.139	0.015			0.027	0.437		0.156	0.030					0.056	0.036
Sheep grazing - forages			0.041		0.650				0.033	0.127		0.107						0.068	0.053
Rainfed olive groves with sheep					0.016					0.291		0.660	0.010					0.039	
Rainfed olive groves									0.016	0.099	0.013	0.823	0.018					0.011	
Irrigated olive groves			0.012						0.022	0.052		0.770	0.083		0.021			0.016	0.024
Vineyards		0.010			0.013				0.054	0.066	0.014	0.029	0.093	0.697				0.025	0.010
Fruit trees					0.014				0.025	0.243	0.548	0.012	0.083	0.020			0.025	0.142	0.040
Stone pine	0.038								0.019	0.189		0.023			0.713		0.000	0.249	0.010
Rice	0.850	0.012							0.039	0.068								0.024	
Irrigated cereals and horticultural crops		0.300	0.038			0.241	0.242		0.049	0.077								0.018	0.012
Rainfed cereals and oilseeds		0.087	0.300			0.038		0.430	0.063	0.030		0.011	0.026						0.018
Rainfed cereals			0.463		0.016	0.015		0.010	0.171	0.150			0.147					0.031	0.039
Pastures without livestock										0.785		0.015	0.152	0.012				0.082	0.088
Fallows			0.079		0.011	0.016			0.752	0.037	0.011		0.077					0.048	0.141

309 * CO – Under cover of cork oak; HO – Under cover of holm oak

310

311

312 Table 3.b – Farming system description – Livestock composition in livestock-oriented farming
 313 systems (average values in proportion to total LU; values under 0.01 are omitted; proportions
 314 above 0.5 are in bold) and livestock density

Farming system	Cattle grazing	Cattle stabled	Fattening steers grazing	Fattening steers stabled	Sheep grazing	Goat grazing	Dairy cows	Pigs grazing	Livestock density (LU/ha)
	Cattle grazing – CO*	0.872		0.084		0.028			
Cattle grazing – HO*	0.865		0.080		0.044				0.677
Cattle grazing – forages	0.859		0.083		0.041				0.967
Grazing goats					0.082	0.910			1.039
Mixed Cattle and sheep - Irrigated forages	0.480		0.073		0.300	0.069	0.077		0.618
Sheep grazing – CO*	0.144		0.010		0.799	0.043			0.245
Sheep grazing – HO*					0.963	0.028			0.387
Sheep grazing - pastures					0.973	0.017			1.004
Sheep grazing - pastures and forages					0.792	0.206			0.711
Sheep grazing - forages	0.049				0.709	0.236			0.378
Rainfed olive groves with sheep					0.974	0.024			1.212

315 * CO – Under cover of cork oak; HO – Under cover of holm oak

316

317 Table 4 – Characterization of farming systems according to the levels of farming intensity,
 318 specialization and labour needs, average farm size (in hectares of UAA and number of LU) and
 319 representativeness (in number of farms, UAA and LU)

	Characterization of farming systems			Average farm size		Representativeness					
	Intensity (10 ³ €/ha)	Speciali- zation (%)	Labour needs (AWU/ha)	UAA (ha)	LU (n.)	Number of farms		UAA		Livestock Units	
						Total	%	Total (10 ³ ha)	(%)	Total (10 ³ LU)	(%)
Cattle grazing – CO*	0.46	84.3	0.005	248.5	128.7	2515	10.6	625	31.1	323.8	36.5
Cattle grazing – HO*	0.31	84.0	0.005	288.2	157.0	1245	5.3	359	17.9	195.5	22.0

Cattle grazing – forages	0.75	68.3	0.005	186.1	88.1	463	2.0	86	4.3	40.8	4.6
Grazing goats	0.94	88.6	0.004	52.6	19.7	251	1.1	13	0.7	4.9	0.6
Mixed Cattle and sheep - Irrigated forages	1.14	75.0	0.005	58.8	24.3	171	0.7	10	0.5	4.1	0.5
Sheep grazing – CO*	0.24	54.6	0.004	89.0	19.0	2346	9.9	209	10.4	44.6	5.0
Sheep grazing – HO*	0.28	64.9	0.004	84.9	23.5	1391	5.9	118	5.9	32.6	3.7
Sheep grazing - pastures	0.71	84.1	0.004	54.5	22.4	1461	6.2	80	4.0	32.7	3.7
Sheep grazing - pastures and forages	0.79	70.6	0.004	52.8	18.7	745	3.1	39	2.0	13.9	1.6
Sheep grazing - forages	0.46	69.9	0.004	25.3	4.3	848	3.6	21	1.1	3.7	0.4
Rainfed olive groves with sheep	0.91	74.1	0.006	12.0	9.6	774	3.3	9	0.5	7.4	0.8
Rainfed olive groves	0.30	92.3	0.010	10.6	0.1	2626	11.1	28	1.4	0.3	0.0
Irrigated olive groves	1.45	93.0	0.023	82.4	0.8	864	3.6	71	3.5	0.7	0.1
Vineyards	1.84	90.3	0.050	24.3	0.5	928	3.9	23	1.1	0.4	0.0
Fruit trees	12.59	89.9	0.036	25.7	1.5	325	1.4	8	0.4	0.5	0.1
Stone pine	4.63	97.6	0.009	68.0	1.4	221	0.9	15	0.7	0.3	0.0
Rice	1.70	93.8	0.018	52.8	4.0	314	1.3	17	0.8	1.3	0.1
Irrigated cereals and horticultural crops	4.66	90.9	0.259	53.7	1.8	1070	4.5	57	2.9	1.9	0.2
Rainfed cereals and oilseeds	0.98	88.2	0.006	65.7	0.9	421	1.8	28	1.4	0.4	0.0
Rainfed cereals	0.42	82.0	0.010	33.8	0.4	1537	6.5	52	2.6	0.6	0.1
Pastures without livestock	0.20	61.5	0.003	48.5	8.5	2602	11.0	126	6.3	22.2	2.5
Fallows	0.59	54.1	0.003	20.8	0.0	582	2.5	12	0.6	0.0	0.0
Total	-	-	-	-	-	23700	100	2007	100	733	-

320 * CO – Under cover of cork oak; HO – Under cover of holm oak

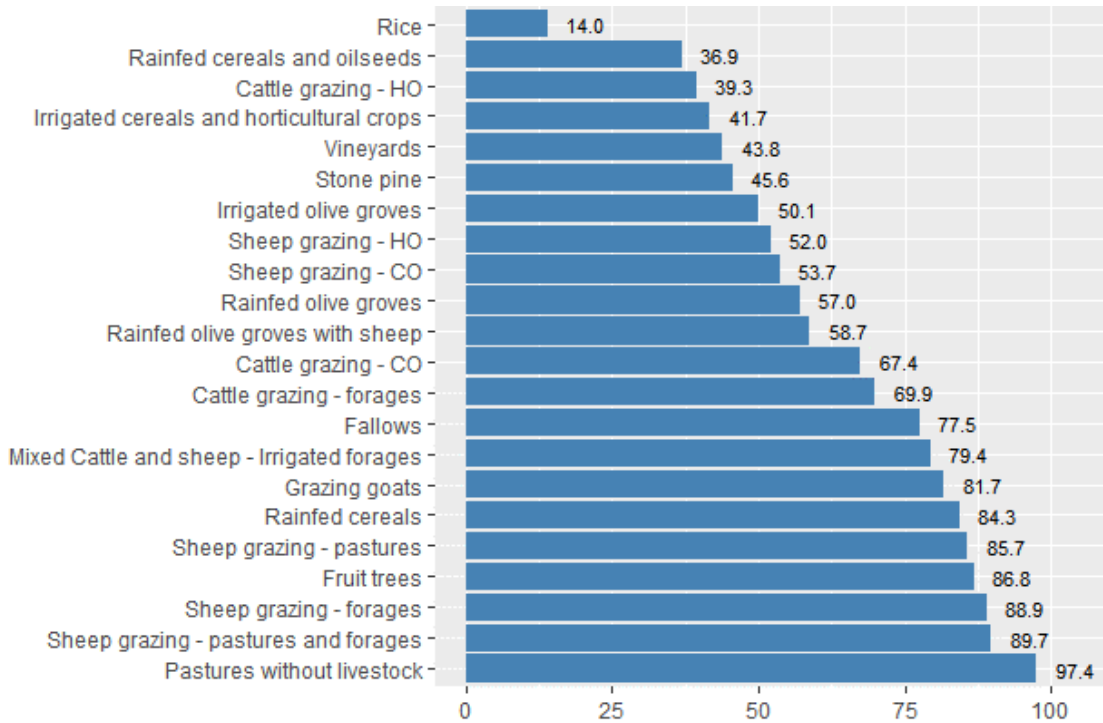
321

322 3.2. Spatial determinants of farming system choice

323 The tuning of the random forest model led to a 500 trees model, with 5 variables randomly
324 sampled as candidates at each split and using the “sampsiz” option to correct size differences
325 across the FS categories (see details in [Annex II](#) – supplementary information). The classification
326 error rates for each of the 22 FS ranged from 14.0% in the Rice system to 97.4% in the Pastures
327 without livestock system (Fig. 2), with an average of 63.7% across all FS, a value that should be

328 evaluated positively considering the high number of classes in the dependent variable (22 FS,
329 for which the random error rate would be about 95.4% with balanced data).

330



331

332 Fig. 2 - Classification error rates for the 22 farming systems (values in %)

333

334 The relative importance of socioeconomic and biophysical variables was very similar, and among
335 the top ten variables, in terms of mean decrease accuracy ([Annex II](#) – supplementary
336 information), six are socioeconomic and four are biophysical. The farm physical dimension
337 variables (FSIZE and BLKSIZE) and a local context of high dependence on family income in
338 agriculture (INCAGRI) proved to be the most relevant socioeconomic factors influencing the
339 choice of FS, while in the biophysical variables the most important were the climatic variables
340 (Fig. 3). The variable indicating access to surface water sources (WPRIVATE) was found to be the
341 least important, either in the global model or in most of the class-specific models.

342

	MODEL OVERALL	Cattle grazing - CO	Cattle grazing - HO	Cattle grazing - forages	Grazing goats	Mixed Cattle and sheep - Irrigated forages	Sheep grazing - CO	Sheep grazing - HO	Sheep grazing - pastures	Sheep grazing - pastures and forages	Sheep grazing - forages	Rainfed olive groves with sheep	Rainfed olive groves	Irrigated olive groves	Vineyards	Fruit trees	Stone pine	Rice	Irrigated cereals and horticultural crops	Rainfed cereals and oilseeds	Rainfed cereals	Pastures without livestock	Fallow
FSIZE	65,6	↑ 59,5	↑ 55,2	↑ 41,5	↓ 5,8	↓ 4,7	↑ -0,7	↑ 6,0	↓ 12,6	↓ 12,4	↓ 15,6	↓ 29,9	↓ 54,8	↓ 19,6	↓ 23,5	↓ 13,0	↑ 9,0	↑ 15,7	↓ 12,2	↑ 13,9	↓ 14,6	↑ 9,7	↓ 29,1
BLKSIZE	50,9	↑ 37,9	↑ 32,9	↑ 19,3	↓ 6,4	↑ 3,1	↑ 28,5	↑ 11,3	↓ 15,5	↓ 14,2	↓ 15,2	↓ 42,5	↓ 50,2	↑ 25,9	↓ 27,1	↓ 7,1	↑ 19,4	↑ 19,4	↓ 14,3	↓ 19,8	↓ 29,0	↑ 6,9	↓ 18,2
JANUS	32,5	↑ 3,3	↑ 6,0	↓ 13,9	↑ 4,3	↑ 1,3	↑ 7,0	↑ 5,8	↓ 7,4	↓ 10,6	↑ 4,8	↓ 5,3	↑ 7,8	↑ 6,3	↓ 12,5	↑ 5,3	↑ 3,5	↓ 8,5	↑ 12,5	↓ 21,4	↑ 3,5	↑ 2,4	↑ 12,8
BLKDIST	30,9	↑ 6,7	↑ 7,0	↓ 5,6	↓ 0,1	↑ 1,1	↑ 3,6	↑ 5,5	↓ 6,3	↓ 2,1	↓ 4,8	↓ 7,1	↓ 11,6	↑ 9,5	↓ 14,0	↓ 6,5	↑ 1,5	↑ 2,3	↑ 12,4	↑ 14,8	↑ 3,1	↓ 4,5	↓ 9,3
INCAGRI	52,7	↑ 18,4	↑ 5,4	↓ 13,6	↓ 1,4	↑ 4,4	↑ 14,0	↑ 26,8	↓ 14,3	↓ 11,6	↓ 12,6	↓ 23,9	↓ 30,2	↓ 17,9	↓ 27,0	↑ 10,3	↑ 29,6	↑ 27,3	↑ 16,6	↑ 24,7	↑ 18,3	↑ 5,8	↓ 26,2
RENT	47,3	↑ 12,6	↑ 6,1	↑ 7,9	↓ 2,1	↑ 3,8	↑ 8,7	↑ 17,8	↓ 10,3	↓ 5,3	↓ 11,7	↑ 11,6	↑ 18,5	↑ 17,7	↓ 23,1	↓ 8,9	↓ 11,8	↑ 15,5	↓ 16,7	↑ 14,3	↑ 11,6	↓ 7,0	↓ 17,1
INCOTH	47,3	↓ 14,8	↑ 8,7	↑ 5,4	↑ 4,8	↓ 4,6	↓ 14,6	↓ 15,8	↑ 12,3	↓ 6,2	↑ 11,3	↑ 10,6	↑ 18,8	↑ 13,8	↑ 20,6	↑ 7,6	↓ 19,1	↓ 16,1	↑ 14,9	↑ 16,0	↑ 10,3	↑ 5,6	↑ 19,6
PDENS	46,0	↓ 16,3	↓ 6,4	↓ 5,4	↓ 1,1	↑ 5,0	↓ 14,8	↓ 23,5	↑ 11,7	↓ 5,8	↑ 12,1	↑ 20,4	↑ 19,6	↓ 16,4	↑ 27,3	↓ 10,7	↓ 20,4	↑ 13,3	↑ 20,4	↑ 21,6	↓ 15,2	↓ 5,3	↓ 25,4
AWU	45,7	↓ 15,7	↓ 18,4	↓ 10,3	↓ 4,9	↑ 5,9	↓ 20,6	↓ 21,7	↓ 19,2	↓ 10,5	↑ 10,2	↑ 14,8	↑ 27,2	↑ 27,0	↑ 42,0	↑ 9,6	↓ 23,3	↑ 17,0	↑ 27,6	↑ 18,1	↑ 20,7	↓ 11,2	↑ 23,5
WPUBLIC	43,7	↓ 28,2	↓ 20,0	↓ 5,3	↓ 14,6	↑ 6,4	↓ 24,1	↓ 27,0	↓ 20,6	↓ 3,0	↑ 8,2	↑ 20,4	↑ 17,2	↑ 43,4	↑ 29,2	↑ 9,6	↓ 18,6	↑ 32,2	↑ 18,8	↑ 42,5	↑ 8,2	↓ 10,3	↑ 4,7
AWU_hired	43,6	↑ 11,9	↑ 11,2	↓ 9,0	↓ 3,1	↑ 3,7	↑ 14,1	↓ 17,3	↓ 7,8	↓ 7,8	↑ 8,8	↑ 9,3	↓ 18,2	↓ 14,4	↑ 19,7	↓ 5,4	↑ 12,4	↑ 12,0	↓ 18,2	↑ 19,0	↓ 14,7	↓ 6,8	↓ 15,0
NATURE	27,0	↑ 5,5	↑ 8,8	↑ 11,8	↓ 4,0	↓ 6,2	↓ 6,4	↓ 6,0	↓ 6,0	↓ 8,4	↓ 6,8	↓ 6,2	↑ 9,8	↓ 12,2	↓ 12,0	↓ 7,0	↑ 6,1	↑ 9,6	↓ 10,5	↓ 12,2	↓ 9,6	↓ 2,4	↓ 11,0
WPRIVATE	8,7	↑ 0,8	↑ 3,0	↓ -0,7	↓ 4,0	↑ 0,9	↓ -1,5	↑ 1,3	↓ 1,7	↓ 1,2	↓ 1,8	↓ 4,6	↓ 8,5	↓ 4,2	↓ 6,2	↑ 2,0	↑ 3,3	↑ 3,4	↑ 0,9	↑ 3,3	↓ 3,5	↓ 1,9	↓ 2,8
TMAX	56,4	↓ 32,0	↑ 21,0	↑ 16,5	↓ 12,2	↑ 6,8	↓ 31,7	↑ 25,6	↑ 11,7	↑ 9,7	↑ 13,6	↑ 22,4	↑ 36,8	↑ 43,7	↑ 28,1	↑ 15,9	↓ 26,2	↑ 31,6	↑ 25,9	↑ 40,2	↑ 28,4	↑ 9,0	↑ 25,6
TMIN	56,2	↑ 18,7	↑ 9,8	↑ 15,6	↓ 3,5	↓ 5,4	↓ 20,0	↓ 26,0	↑ 10,4	↓ 8,0	↓ 10,1	↑ 16,9	↓ 22,2	↓ 19,4	↓ 26,4	↓ 10,3	↓ 14,4	↓ 15,4	↓ 19,7	↓ 19,2	↓ 15,3	↓ 7,3	↓ 18,4
PREC	48,1	↑ 25,7	↓ 23,7	↓ 20,1	↑ 11,5	↑ 5,4	↑ 32,2	↓ 33,3	↑ 14,7	↑ 12,9	↑ 13,2	↑ 22,4	↑ 22,0	↓ 25,7	↑ 24,4	↑ 17,7	↑ 20,5	↑ 22,4	↑ 23,6	↓ 22,3	↓ 19,4	↑ 6,9	↓ 23,3
SDEPTH	47,3	↓ 32,2	↓ 13,8	↓ 10,5	↓ 15,9	↑ 6,2	↓ 23,4	↓ 33,2	↓ 9,1	↓ 9,0	↓ 5,3	↑ 16,7	↑ 20,0	↓ 14,0	↑ 22,7	↓ 5,5	↑ 27,9	↑ 25,7	↑ 35,9	↑ 34,5	↓ 14,9	↓ 8,5	↓ 15,7
LIGHT_S	41,5	↑ 15,8	↑ 19,0	↓ 14,2	↑ 7,0	↑ 5,8	↑ 26,4	↓ 18,5	↓ 8,2	↓ 9,2	↓ 2,4	↑ 13,2	↓ 17,1	↓ 16,9	↑ 14,3	↑ 9,7	↑ 26,8	↑ 22,8	↑ 15,8	↓ 20,6	↓ 16,7	↓ 4,0	↓ 16,6
ACID	41,2	↑ 20,8	↑ 12,5	↓ 7,7	↑ 5,9	↓ -0,4	↑ 13,5	↑ 12,0	↓ 6,0	↓ 4,0	↓ 4,3	↓ 9,4	↓ 24,6	↓ 17,1	↓ 18,8	↑ 3,4	↑ 8,7	↑ 10,6	↓ 12,0	↓ 20,0	↓ 10,7	↑ 2,9	↓ 11,2
NEUTRAL	40,9	↓ 14,2	↓ 10,8	↓ 10,9	↓ 9,7	↑ 2,0	↓ 15,9	↓ 24,3	↓ 8,9	↓ 5,0	↓ 4,8	↑ 11,1	↑ 28,6	↑ 21,9	↑ 19,7	↓ 9,7	↓ 15,6	↓ 19,0	↑ 15,6	↑ 34,3	↑ 12,0	↓ 0,8	↓ 10,9
MEDIUM_S	40,1	↓ 13,1	↑ 14,4	↑ 7,3	↓ 4,1	↓ 3,5	↑ 12,5	↑ 16,6	↑ 5,7	↑ 6,8	↑ 1,3	↓ 10,4	↓ 11,8	↑ 8,2	↓ 17,6	↓ 4,9	↑ 16,0	↑ 13,0	↓ 10,2	↑ 17,3	↑ 9,2	↓ 5,9	↑ 12,3
VERYACID	38,1	↑ 13,5	↑ 7,9	↓ 5,9	↑ 2,9	↑ 4,7	↑ 17,7	↑ 16,6	↑ 6,8	↑ 6,8	↓ 5,3	↓ 9,7	↓ 11,7	↓ 9,8	↑ 10,6	↑ 2,5	↑ 11,2	↓ 10,3	↑ 11,7	↑ 12,3	↑ 4,7	↓ 3,5	↑ 11,8
SMOOTH	37,7	↓ 20,3	↓ 13,4	↓ 10,5	↓ 14,6	↑ 3,6	↓ 22,7	↓ 26,7	↓ 10,6	↓ 6,4	↑ 5,9	↓ 9,9	↓ 12,9	↓ 14,2	↑ 7,9	↑ 18,9	↓ 7,4	↓ 13,2	↑ 20,1	↑ 24,8	↑ 19,7	↓ 11,9	↓ 3,0
HEAVY_S	34,2	↓ 14,3	↓ 10,5	↓ 5,7	↓ 4,9	↑ 5,9	↓ 17,2	↓ 15,3	↓ 10,4	↓ 6,1	↓ 5,8	↑ 11,6	↑ 15,6	↑ 13,0	↑ 16,3	↓ 10,2	↓ 20,6	↑ 20,4	↑ 16,3	↑ 17,3	↑ 15,8	↓ 5,1	↑ 14,6
STEEP	31,8	↑ 19,2	↑ 10,7	↓ 9,4	↑ 19,3	↓ 0,8	↑ 27,4	↑ 25,5	↑ 2,2	↓ 6,8	↓ 4,0	↑ 3,1	↑ 11,8	↓ 14,2	↓ 14,5	↑ 5,8	↑ 9,2	↓ 9,9	↓ 9,4	↓ 15,2	↓ 12,9	↑ 5,9	↓ 13,1
ALCALINE	31,7	↓ 12,1	↓ 11,5	↓ 8,6	↓ 3,5	↑ 4,2	↓ 18,1	↓ 16,6	↓ 9,4	↓ 1,1	↓ 8,1	↑ 9,7	↑ 11,0	↑ 7,9	↑ 12,0	↓ 0,6	↓ 12,6	↑ 15,1	↑ 19,8	↑ 12,3	↓ 12,4	↓ 5,7	↓ 9,2
MODERATE	29,9	↑ 18,0	↑ 12,0	↑ 6,1	↑ 5,2	↓ 2,9	↑ 18,0	↑ 18,4	↑ 11,2	↑ 3,0	↑ 4,6	↑ 10,0	↑ 7,1	↑ 5,2	↑ 14,2	↑ 4,4	↑ 7,7	↓ 17,8	↓ 22,1	↓ 13,4	↑ 4,6	↑ 0,3	↑ 11,2

343

344 Fig. 3 - Variable importance for the overall model and for each farming system. Socioeconomic farm structure variables in blue; local-socioeconomic variables
345 in orange; biophysical variables in green. Variables ordered by decreasing variable importance in the overall model and within each sub-group. Symbols ↑, ↓
346 and ⇕ indicate whether the marginal effect of the variable in each farming system is mostly positive, negative or non-monotonic, respectively, based on the
347 shape of the fitted function on the partial dependence plots (partial dependence plots are provided in supplementary information, Annex IV). Variable
348 description in Table 2

349 Farm size (FSIZE) and average farm-block size (BLKSIZE) were the most relevant variables for the
350 choice of Cattle grazing FS, positively influencing its choice (Fig. 3). The same variables also have
351 a relevant effect on most sheep systems but, in this case, predominantly on the opposite
352 direction (Fig. 3). The choice of the Cattle grazing – CO system is positively influenced by the
353 increase of the average annual rainfall (PREC) and negatively by high summer temperatures
354 (TMAX), which has a positive effect on the choice of the Cattle grazing – HO system. The Cattle
355 grazing – forages system is distinguished by a preference for warmer winters (Fig. 3).

356 The Grazing goats system is positively related to sloping terrain; its choice is favoured by
357 increasing the slope (STEEP), while avoiding flat land (SMOOTH). This system is also
358 characterized by avoiding public irrigation areas (WPUBLIC) and deep soils (SDEPTH). The choice
359 of the Mixed Cattle and sheep - Irrigated forages system is favoured by deeper soils, public
360 irrigation systems (WPUBLIC) and high local labour availability (AWU), which is probably related
361 to the irrigated forages component of this FS or with the labour needs associated with grazing
362 herds. The average annual rainfall (PREC) has opposite effects in Sheep grazing – HO and – CO
363 systems, with the first system being favoured by lower rainfall values, and the other way around
364 in the later system. Sheep grazing – CO is also favoured by areas with steeper slopes (STEEP) and
365 light soils (LIGHT_S), while the choice of Sheep grazing – HO decreases with deeper soils,
366 smoother terrain and public irrigation structures. Lower values of local labour availability (AWU)
367 seem to promote the choice of the Sheep grazing – pastures system, while the choice of Sheep
368 grazing – forages system is negatively influenced as the local values of agricultural income
369 dependence (INCAGRI) raises (Fig. 3).

370 Both Rainfed olive groves systems (with and without sheep) are strongly related to smaller farm
371 sizes, as these are also the two systems with lower average UAA (Table 4). Both are positively
372 related to high summer temperatures (TMAX) and negatively to higher regional values of
373 agricultural income dependence. The Rainfed olive groves with sheep system is favoured when
374 average annual rainfall increases, and the Rainfed olive groves system is positively related to
375 neutral pH soils (NEUTRAL) (Fig. 3).

376 The Irrigated olive groves system is positively related to high summer temperatures, public
377 irrigation systems, high local labour availability and high average farm-block size. It is negatively
378 related to high average annual rainfall. The choice of the Vineyards system tends to increase
379 with higher values of regional labour availability, public irrigation systems and population
380 density (PDENS). The Fruit trees system is positively associated with average annual rainfall and
381 negatively with high population density and warmer winters (TMIN). The choice of the Stone

382 pine system is favoured by light soils (LIGHT_S) and discouraged by high summer temperatures
383 and population density (Fig. 3).

384 In the annual crops, the Rice system is mostly favoured by the presence of public irrigation
385 systems, also by higher regional values of agricultural income dependence (INCAGRI) and soil
386 depth, while negatively influenced by high summer temperatures. The Irrigated cereals and
387 horticultural crops system is positively related to soil depth, regional labour availability and
388 smooth slope terrain. The choice of the Rainfed cereals and oilseeds system is encouraged with
389 public irrigation systems and higher values of soil depth, neutral pH and high summer
390 temperatures. The Rainfed cereals system is negatively related to bigger farm-block sizes and
391 average annual rainfall. The Pastures without livestock system seems to be promoted when
392 labour availability is lower and outside public irrigation systems, although this FS presented the
393 highest error rate (Fig. 2). The Fallows system also displays complex relations with the
394 predictors, though it seems to be more positively associated with small farms and areas of low
395 population density (Fig. 3).

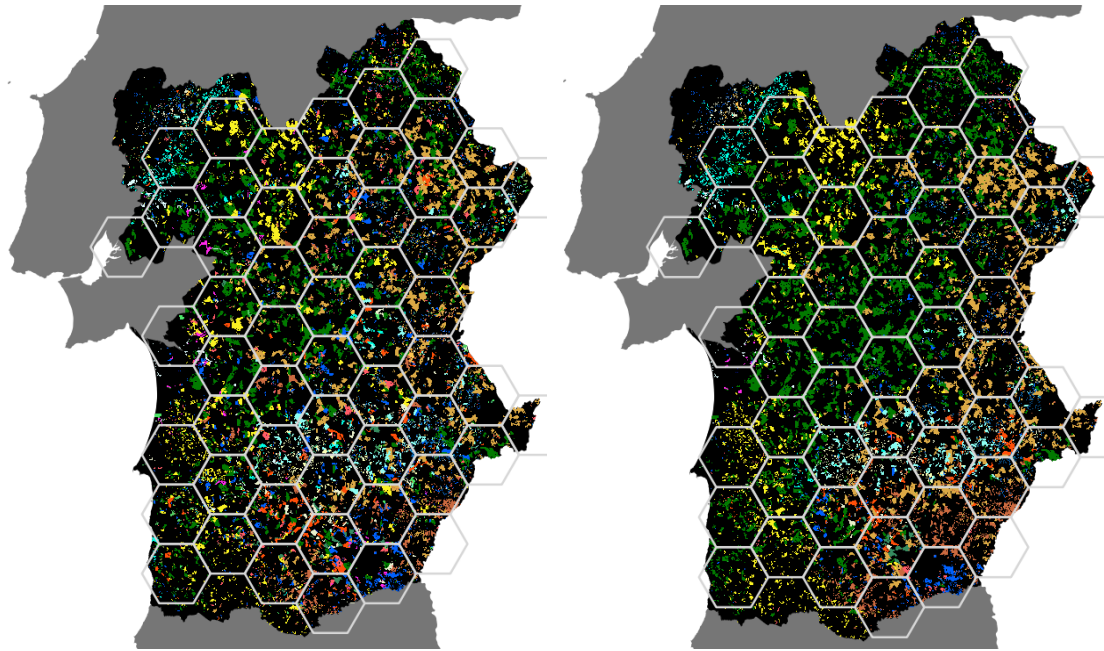
396 Finally, the prediction accuracy for the different farming systems (Fig. 2) showed a modest but
397 positive correlation with the corresponding levels of agricultural specialization and labour needs
398 (Table 4) (correlation coefficients of 0.44 and 0.26, respectively), and a virtually non-existent
399 relationship with the level of agricultural intensity (correlation coefficient -0.03).

400

401 3.3. Spatial patterns of landscape-scale farming systems composition

402 The hexagonal lattice resulted with 56 usable analysis units, i.e., hexagons with >33% of the area
403 overlapped with LPIS data (Fig. 4). The average error rate in the FS spatial pattern predictions
404 across all hexagons was 28.7% (max. 47.3%; min. 9.2%), which is substantially lower than the
405 error rate obtained with model predictions at the farm-level (67.3%). The average coefficient of
406 determination was 0.89 (max. 1.00; min. 0.28), revealing a good model fit.

407



408

409 Fig. 4 - Observed (left) and predicted (right) FS maps for the 1/3 observations used in the
 410 model validation dataset and the hexagons network used to assess model accuracy in FS
 411 spatial patterns prediction (different colours identify distinct FS; detailed maps showing the
 412 spatial distribution of each farming system are provided in supplementary information,
 413 [Annex III](#)).

414

415 4. Discussion

416 The use of farm-level data (IACS) provided by the national CAP paying agency proved to be a
 417 suitable approach to derive the FS typology for the study area, in line with previous studies
 418 ([Ribeiro et al., 2018, 2016, 2014](#)). The spatial-explicit nature of these data (LPIS) allowed a very
 419 fine characterization of farms, including in their biophysical, structural and socioeconomic
 420 features. As expected, the extent and heterogeneity of the study area, in both socioeconomic
 421 and biophysical features, led to a broad typology of 22 farming systems, which are a direct
 422 outcome of distinct farm-management adaptive-responses to a variety of farm features and
 423 contexts.

424 Although the FS typology was balanced in terms of crop- and livestock-oriented systems, the
 425 results showed that most of the study area is currently devoted to livestock systems, particularly
 426 cattle grazing. Although the present study does not allow this to be confirmed, farmers'
 427 preference for these systems may be due to an (at date) ongoing direct payment for suckler

428 cows (and partially to sheep and goats), a national agricultural policy option taken under the
429 2003 CAP reform that significantly impacted FS dynamics in the region (Ribeiro et al., 2014).

430

431 4.1. Farm structure drivers

432 Many of the effects of structural socioeconomic variables observed here are consistent with
433 those of previous studies. For example, the farm-size was found to positively influence the
434 choice of extensive livestock systems over crop systems, which was also observed in Ribeiro et
435 al. (2018), and also in the choice between cattle grazing over some sheep grazing specialized
436 systems, which was also observed in studies by Ribeiro et al. (2014).

437 Access to private sources of surface irrigation water showed very little importance in the FS
438 choice-models, which is apparently odd for a region where water is often a limiting factor. This
439 was probably due not only to the type of variable used (dummy variable, with 1 = “yes, the farm
440 has access to surface water sources” and 0 otherwise) but also to the fact of not including access
441 to groundwater from water wells, due to lack of data, which are a common source in parts of
442 the region. In contrast, water availability from public irrigation systems is essential in explaining
443 the spatial location of several irrigated FS (either cereals, oil seeds or intensive olive groves and
444 vineyards) showing the importance of public water management policy over other biophysical
445 constraints (Kahil et al., 2015). Not surprisingly, these farming systems most associated with
446 large public irrigation systems are among the most intensive ones.

447 Public intervention in nature conservation areas seems to be of little relevance for FS choice
448 since although a considerable share of agriculture area is classified for nature conservation, the
449 corresponding variable (NATURE) was one of the least relevant within a list of dimensions that
450 has farm and block size at the top.

451 An interesting side-result of our approach was the insight of an overall negative, though
452 moderate, relationship between farm size and the level of agricultural intensity, indicating that
453 larger farms tend to adopt less intensive FS, a finding that goes back to earlier works (Cornia,
454 1985; Grigg, 2005; Reboul, 1989, 1976). Exceptions, however, can be found when contrasting,
455 e.g., the Rainfed olive groves and the Irrigated olive groves systems, where large investments in
456 fixed capital (including irrigation systems), together with labour availability, seem to provide
457 increasing returns to scale, which was also reported in more recent studies (Deininger et al.,
458 2018; Rada and Fuglie, 2019).

459

460 4.2. Socioeconomic context drivers

461 Regarding the socioeconomic context of the farms, the level of agricultural professionalization
462 (inferred from the INCAGRI variable) and farm labour availability proved to be significant drivers
463 of FS. On one side, higher levels of professionalization, which in Portugal are considerably low
464 in average when comparing to non-South European countries (Arnalte-Alegre and Ortiz-
465 Miranda, 2013), are positively associated with Rice, Stone pine or Rainfed cereals and oilseeds
466 systems. On the other side, Vineyards and Irrigated cereals and horticultural crops, which show
467 the highest levels of labour intensity per hectare and the highest average of labour units per
468 farm, are positively associated with local availability of farm labour. Considering that
469 horticultural crops typically have the highest wage labour ratios compared to other crops
470 (Baptista and Rolo, 2017), it was surprising that it did not show up associated with high local
471 proportion of hired labour. A possible explanation may be the high geographic mobility of hired
472 workers (Baptista and Rolo, 2017), although it may also emerge from the heterogeneity in labour
473 intensity within this FS, since it encompasses irrigated cereals and industrial horticulture, with
474 considerable levels of mechanization, as well as horticultural crops with very high levels of
475 labour needs.

476 The fact that local labour availability has a more widespread importance as a FS driver than rural
477 population density, which only stands out in the single case of Vineyards, contradicts the idea of
478 permanent crops and horticulture as able of promoting rural population retention (Egea and
479 Pérez y Pérez, 2016), i.e., it points to the dissociation between farm labour dynamics and local
480 demographics (Baptista and Rolo, 2017). While vineyards remain located in higher populated
481 parishes, following deep-rooted institutional constraints by protected designations of origin,
482 olive groves (either irrigated or rainfed) show no relation with local demographics.

483 Land renting (RENT) did not appear in the top 5 drivers in any FS, suggesting that the size of the
484 land renting market does not appear to have much effect on the choice of FS in the study area.
485 However, the positive relationship observed between land renting and livestock grazing FS,
486 especially cattle, suggests that these systems, which have experienced marked growth in the
487 region in recent years (Ribeiro et al., 2018), expanded in part at the expense of this tenancy
488 regime.

489

490 4.3. Biophysical drivers

491 As anticipated, biophysical factors related to climate, soils and relief, proved to be strong
492 determinants of FS spatial distribution (Grigg, 2005). Summer heat and annual precipitation
493 came up as the main biophysical drivers of FS spatial distribution in the study area. High summer
494 temperatures seem to favour the choice of olive groves, vineyards, rainfed cereals and cattle
495 grazing systems associated to Holms oak, and to discourage livestock systems associated to Cork
496 oak, Stone pine or Rice systems. Winter cold increases the likelihood of fruit tree systems and
497 the opposite with forage systems.

498 Deep soils and smooth relief are positive drivers of the Rice and Irrigated cereals and
499 horticultural crops systems. The opposite effect is found towards the Grazing goats system,
500 which is strongly related to steeper slopes. Soil pH did not emerge as a major driver for the
501 distribution of any FS, except for rainfed cereals and olive groves systems which showed a
502 preference for neutral pH soils.

503 Following Cork and Holm oak distinct preferences for soil and climate (Surová and Pinto-Correia,
504 2008), livestock systems associated with these two species of oaks were found distributed
505 accordingly: Cork oak-associated systems prevail more to the coast and north of the study area,
506 where summer temperatures are milder, annual rainfall is higher and soils are sandy and light-
507 textured; Holm oak-associated systems are further inland and south, where summers are
508 warmer, annual rainfall is lower and soils are frequently poor and fairly thin.

509

510 4.4. Farming system prediction at the farm and landscape levels

511 Although the model's ability to predict individual FS was quite varied, depending on the FS, when
512 applied to predicting FS patterns at the landscape-level the model revealed a much higher hit
513 rate. The random forest approach applied in the model estimation proved to be a valuable
514 choice, particularly in dealing with such high dimensional data (Strobl et al., 2009). At the
515 landscape level, the model was very effective in predicting farming systems patterns, i.e., the
516 shares of FS composition within hexagon-shaped landscape units. For agricultural landscape
517 planning focused on agroecosystem services provision, this may be the right scale of analysis,
518 since a minimum share of farmland managed under the FS delivering those services should be
519 sufficient to ensure the socially desired level of service, rather than requiring the service to be
520 provided by a specific set of farms over a period of time (Andersen, 2017), as is typically the case
521 with many agri-environment schemes requiring multi-annual contracts with individual farmers.

522

523 4.5. Shortcomings of the approach and recommendations for future 524 research

525 Despite the valuable advantages evidenced by the proposed approach, there is still room for
526 future improvement. Improvements mostly relate to characteristics of the IACS and LPIS
527 datasets and methodological options that are dependent on the geographic context of our study
528 area.

529 While recognized as having high potential for supporting data driven research, the IACS / LPIS
530 datasets present limitations, such as the lack of information to characterize farmers'
531 socioeconomic profile, or information on complementarity relationships between farms, such
532 as the rental or sale of pastures, which can mislead the computation of farms' stock density.
533 Such information would be valuable to include in the FS choice models.

534 The fact that the empirical work was carried out in a region where the landscape is largely
535 dominated by agriculture, makes it possible to closely link FS choice with landscape modelling.
536 Where this is not the case, such as many mountain and less favoured regions across the EU, this
537 approach may not deliver the same results, given the smaller share of agriculture in the
538 landscape. Additionally, in such regions a significant part of agriculture is probably outside any
539 CAP support system, so that an approach based on IACS / LPIS data can only partially capture an
540 agricultural reality that is itself marginal at the landscape scale. Paradoxically, these regions
541 often include significant shares of high nature value farmlands at the EU level ([Lomba et al.,](#)
542 [2014](#)). Nevertheless, it should be worth trying to reproduce the approach in such regions in the
543 future, to test the generalization of the framework.

544 Because our farm characterization variables report to a single year, the effect of economic or
545 policy variables such as prices or subsidies can only be assumed as underpinning the farmers'
546 choices reflected on the observed 2017 IACS / LPIS data. However, the use of this type of
547 variables in the model, provided that time-series of farm-level data can be made available,
548 would significantly extend the scope of this approach, allowing its use to evaluate policy and
549 price change scenarios. Even without additional temporal data, the framework can take
550 advantage of the wide extension of the study area to perform, e.g., climate-change scenarios
551 assessment, by adopting a space-for-time substitution approach.

552 The selection of candidate variables to be tested as drivers of FS choice is also a key step in the
553 modelling approach. The misspecification or the absence of key variables can substantially

554 undermine models' performance. The problems observed with variable WPRIVATE may be one
555 such case, as this variable only reported access to small private surface water sources, which are
556 mostly torrential regime in this region, with insufficient water guarantees to encourage investing
557 in irrigation systems, and not taking into account that a significant portion of private irrigation
558 in this region is probably resorting to groundwater sources. This premise, which we could not
559 test due to lack of data, would be worth further investigation, should spatially explicit data on
560 groundwater uptakes becomes available.

561 Another issue deserving further investigation concerns the dimension of the grid of landscape
562 analysis units. It is possible that the size of these units (i.e. the hexagons, in the current case)
563 influences the accuracy of the model, so future investigation focused on determining its optimal
564 size could prove to be of high value.

565 Also, one aspect that has not been explored in the present study and should merit further
566 investigation is the occurrence of interaction effects between drivers. Although the way random
567 forests deal with these effects is still subject to discussion ([Wright et al., 2016](#)), its likely existence
568 recommends additional analysis.

569 Finally, the fact that the prediction error rate has shown significant disparities across the FS
570 suggests that the choice of some of these FS may be due to effects not measured by the variables
571 examined, including factors related to farmers' desires, attitudes and motivations, or with their
572 socioeconomic profile which, as mentioned above, cannot be assessed on the basis of IACS data.
573 One such case would be the Pastures without livestock system, whose choice is probably mostly
574 determined by the presence of livestock farms in the nearby, with whom the farm can negotiate
575 grazing land renting, rather than by the biophysical characteristics of the farm or its
576 socioeconomic context. On the other hand, FS with lower error rates in the model were those
577 who most depend on the chosen socioeconomic or biophysical factors, such as the Rice, Irrigated
578 cereals and horticulture or Rainfed cereals and oilseed systems (where cereals are an autumn-
579 winter rainfed crop and oilseeds are grown in spring-summer season, often irrigated) that highly
580 depend on irrigation water provided by public irrigation systems in this region. The same applies
581 to the Vineyards system, whose location is highly dependent on the availability of regional
582 labour supply, to meet peaks of labour needs at certain times of the year, related to certain crop
583 operations (e.g. harvesting or pruning). In the present market, policy and technological context,
584 these FS revealed greater dependence on farm structure and "territorial embeddedness" (*sensu*
585 [Cerceanu et al., 2018](#)).

586

587 4.6. Concluding remarks

588 Our framework proved to be a suitable approach to investigate the role of human and physical
589 factors in farmers' decisions regarding the choice of the FS, providing effective contributions to
590 improve our understanding of the spatial distribution of FS when observed at a regional scale.

591 This research led to a better understanding of how each of the considered socioeconomic and
592 biophysical factors influences the spatial location of a wide range of FS, a subject seldom
593 explored in such detail in the literature. Results showed that both socioeconomic and
594 biophysical factors exert a high influence on the spatial distribution of FS, clearly revealing the
595 shortcomings of planning proposals exclusively confined to the agroecological aptitude
596 perspective (Nguyen et al., 2015; Pirovani et al., 2018). That influence, however, is not
597 comparable across FS, being decisive for the location of some FS and marginal for others.

598 Contrasting relationships were found between the agricultural intensity level and the degree of
599 dependence on biophysical drivers among the FS, with the simultaneous existence of intensive
600 FS with strong connection to biophysical factors (e.g. Rice system), and others similarly intensive
601 FS but where this relation is much weaker (e.g. Fruit trees system). This finding shows the
602 shortcomings of the assimilation between agricultural intensity and degree of artificialization of
603 the farm's conditions, largely dominant in the literature on the relationship between agriculture
604 and biodiversity/natural resources (Keenleyside et al., 2014). This assimilation ignores the
605 distinction between land and labour productivity and the fact that intensity differences may be
606 due to labour intensity levels rather than higher levels of external outputs. Our results point thus
607 to the need of not reducing farming systems diversity to an intensity gradient, when comparing
608 across distinct productions (Ribeiro et al., 2016).

609 The use of IACS / LPIS data proved to be an invaluable asset for the research, enabling a high-
610 detailed farm-level analysis, not achievable using official statistics and usually only possible
611 through expensive and time-consuming farm surveys, often unfeasible for research works
612 developed at regional scales like the one used in this study. Therefore, it is worth renewing an
613 appeal previously made (Santos et al., 2020; Tóth and Kučas, 2016), addressed at the EU bodies
614 responsible for maintaining the IACS databases, to make them more accessible to the scientific
615 community, while safeguarding confidentiality duties.

616 Overall, the model's ability to perform scenario simulations and to predict patterns of farming
617 systems assigns this approach with a high potential to support information-based policy design

618 to improve agricultural landscape planning and ensure the provision of socially valued
619 agroecosystem services.

620

621 Acknowledgments

622 This work was funded by project “FARSYD– FARming SYstems as tool to support policies for
623 effective conservation and management of high nature value farmlanDs” – POCI-01-0145-
624 FEDER-016664 (PTDC/AAG-REC/5007/2014), supported by Norte Portugal Regional Operational
625 Programme (NORTE 2020), under the PORTUGAL 2020 Partnership Agreement, through the
626 European Regional Development Fund (ERDF). This research was also supported by the Forest
627 Research Centre, a research unit funded by Fundação para a Ciência e a Tecnologia I.P. (FCT),
628 Portugal (UID/AGR/00239/2019). FM was supported by FCT (contract IF/01053/2015). AL was
629 supported by national funds through FCT – Fundação para a Ciência e a Tecnologia, I.P., in the
630 context of the Transitory Norm - DL57/2016/CP1440/CT0001

631

5. References

- Andersen, E., 2017. The farming system component of European agricultural landscapes. *Eur. J. Agron.* 82, 282–291. doi:10.1016/j.eja.2016.09.011
- Arnalte-Alegre, E., Ortiz-Miranda, D., 2013. The “southern model” of european agriculture revisited: Continuities and dynamics, *Research in Rural Sociology and Development*. Emerald Group Publishing Limited. doi:10.1108/S1057-1922(2013)0000019005
- Baptista, F.O., Rolo, J.C., 2017. Trabalho agrícola: percursos e modelos. *Cultiv.* n.º 10 27–35.
- Benoît, M., Rizzo, D., Marraccini, E., Moonen, A.C., Galli, M., Lardon, S., Rapey, H., Thenail, C., Bonari, E., 2012. Landscape agronomy: A new field for addressing agricultural landscape dynamics. *Landsc. Ecol.* 27, 1385–1394. doi:10.1007/s10980-012-9802-8
- Birch, C.P.D., Oom, S.P., Beecham, J.A., 2007. Rectangular and hexagonal grids used for observation, experiment and simulation in ecology. *Ecol. Modell.* 206, 347–359. doi:10.1016/j.ecolmodel.2007.03.041
- Breiman, L., 2001. Random Forests. *Eur. J. Math.* 45, 5–32. doi:10.1023/A:1010933404324
- Canadas, M.J., Novais, A., 2014. Bringing local socioeconomic context to the analysis of forest owners’ management. *Land use policy* 41, 397–407. doi:10.1016/J.LANDUSEPOL.2014.06.017
- Cerceau, J., Mat, N., Junqua, G., 2018. Territorial embeddedness of natural resource management: A perspective through the implementation of Industrial Ecology. *Geoforum* 89, 29–42. doi:10.1016/j.geoforum.2018.01.001
- Cornia, G.A., 1985. Farm size, land yields and the agricultural production function: An analysis for fifteen developing countries. *World Dev.* 13, 513–534. doi:10.1016/0305-750X(85)90054-3
- Debolini, M., Marraccini, E., Dubeuf, J.P., Geijzendorffer, I.R., Guerra, C., Simon, M., Targetti, S., Napoléone, C., 2018. Land and farming system dynamics and their drivers in the Mediterranean Basin. *Land use policy* 75, 702–710. doi:10.1016/j.landusepol.2017.07.010
- Deffontaines, J.-P., 2004. L’objet dans l’espace agricole. Le regard d’un géoagronome. *Natures Sci. Sociétés* 12, 299–304. doi:10.1051/nss:2004041

- Deffontaines, J.P.P., Thenail, C., Baudry, J., 1995. Agricultural systems and landscape patterns: how can we build a relationship? *Landsc. Urban Plan.* 31, 3–10. doi:10.1016/0169-2046(94)01031-3
- Deininger, K., Jin, S., Liu, Y., Singh, S.K., 2018. Can Labor-Market Imperfections Explain Changes in the Inverse Farm Size–Productivity Relationship? Longitudinal Evidence from Rural India. *Land Econ.* 94, 239–258. doi:10.3368/le.94.2.239
- Egea, P., Pérez y Pérez, L., 2016. Sustainability and multifunctionality of protected designations of origin of olive oil in Spain. *Land use policy* 58, 264–275. doi:10.1016/j.landusepol.2016.07.017
- Ferraz-de-Oliveira, M.I., Azeda, C., Pinto-Correia, T., 2016. Management of Montados and Dehesas for High Nature Value: an interdisciplinary pathway. *Agrofor. Syst.* 90, 1–6. doi:10.1007/s10457-016-9900-8
- Friedman, J.H., 2001. Greedy Function Approximation: A Gradient Boosting Machine. *Ann. Stat.* 29, 1189–1232. doi:10.1214/009053606000000795
- Grigg, D., 2005. *An Introduction to Agricultural Geography, Second Edition, Second ed.* ed. Routledge, London and New York.
- Hazell, P., Wood, S., 2008. Drivers of change in global agriculture. *Philos. Trans. R. Soc. B Biol. Sci.* 363, 495–515. doi:10.1098/rstb.2007.2166
- Hengl, T., Nussbaum, M., Wright, M.N., Heuvelink, G.B.M., Gräler, B., 2018. Random forest as a generic framework for predictive modeling of spatial and spatio-temporal variables. *PeerJ* 6, e5518. doi:10.7717/peerj.5518
- Januszewski, J., 1968. Index of land consolidation as a criterion of the degree of concentration. *Geogr. Plonica* 14, 291–296.
- Kahil, M.T., Connor, J.D., Albiac, J., 2015. Efficient water management policies for irrigation adaptation to climate change in Southern Europe. *Ecol. Econ.* 120, 226–233. doi:10.1016/j.ecolecon.2015.11.004
- Keenleyside, C., Beaufoy, G., Tucker, G., Jones, G., 2014. High Nature Value farming throughout EU-27 and its financial support under the CAP. Report Prepared for DG Environment, Contract No ENV B.1/ETU/2012/0035, Institute for European Environmental Policy. London. doi:10.2779/91086

- Kristensen, S.B.P., Busck, A.G., van der Sluis, T., Gaube, V., 2016. Patterns and drivers of farm-level land use change in selected European rural landscapes. *Land use policy* 57, 786–799. doi:10.1016/j.landusepol.2015.07.014
- Lacoste, M., Lawes, R., Ducourtieux, O., Flower, K., 2018. Assessing regional farming system diversity using a mixed methods typology: the value of comparative agriculture tested in broadacre Australia. *Geoforum* 90, 183–205. doi:10.1016/j.geoforum.2018.01.017
- Landis, D.A., 2017. Designing agricultural landscapes for biodiversity-based ecosystem services. *Basic Appl. Ecol.* 18, 1–12. doi:10.1016/j.baae.2016.07.005
- Latruffe, L., Piet, L., 2014. Does land fragmentation affect farm performance? A case study from Brittany, France. *Agric. Syst.* 129, 68–80. doi:10.1016/j.agsy.2014.05.005
- Levers, C., Butsic, V., Verburg, P.H., Müller, D., Kuemmerle, T., 2016. Drivers of changes in agricultural intensity in Europe. *Land use policy* 58, 380–393. doi:10.1016/j.landusepol.2016.08.013
- Liaw, A., Wiener, M., 2002. Classification and Regression by randomForest. *R News* 2, 18–22.
- Lomba, A., Guerra, C., Alonso, J., Honrado, J.P., Jongman, R., McCracken, D., 2014. Mapping and monitoring High Nature Value farmlands: challenges in European landscapes. *J. Environ. Manage.* 143, 140–50. doi:10.1016/j.jenvman.2014.04.029
- Lomba, A., Strohbach, M., Jerrentrup, J.S., Dauber, J., Klimek, S., McCracken, D.I., 2017. Making the best of both worlds: Can high-resolution agricultural administrative data support the assessment of High Nature Value farmlands across Europe? *Ecol. Indic.* 72, 118–130. doi:10.1016/j.ecolind.2016.08.008
- Martel, G., Aviron, S., Joannon, A., Lalechère, E., Roche, B., Boussard, H., 2019. Impact of farming systems on agricultural landscapes and biodiversity: From plot to farm and landscape scales. *Eur. J. Agron.* 107, 53–62. doi:10.1016/j.eja.2017.07.014
- Nguyen, T.T., Verdoodt, A., Van Y, T., Delbecque, N., Tran, T.C., Van Ranst, E., 2015. Design of a GIS and multi-criteria based land evaluation procedure for sustainable land-use planning at the regional level. *Agric. Ecosyst. Environ.* 200, 1–11. doi:10.1016/j.agee.2014.10.015
- Pirovani, D.B., Pezzopane, J.E.M., Xavier, A.C., Pezzopane, J.R.M., de Jesus Júnior, W.C., Machuca, M.A.H., dos Santos, G.M.A.D.A., da Silva, S.F., de Almeida, S.L.H., de Oliveira Peluzio, T.M., Eugenio, F.C., Moreira, T.R., Alexandre, R.S., dos Santos, A.R., 2018. Climate change impacts on the aptitude area of forest species. *Ecol. Indic.* 95, 405–416.

doi:10.1016/j.ecolind.2018.08.002

Plieninger, T., Draux, H., Fagerholm, N., Bieling, C., Bürgi, M., Kizos, T., Kuemmerle, T., Primdahl, J., Verburg, P.H., 2016. The driving forces of landscape change in Europe: A systematic review of the evidence. *Land use policy* 57, 204–214.

doi:10.1016/j.landusepol.2016.04.040

R Development Core Team, 2017. R: A language and environment for statistical computing. [WWW Document]. R Found. Stat. Comput. URL <http://www.r-project.org> (accessed 11.12.18).

Rada, N.E., Fuglie, K.O., 2019. New perspectives on farm size and productivity. *Food Policy* 84, 147–152. doi:10.1016/j.foodpol.2018.03.015

Reboul, C., 1989. Monsieur le capital et madame la terre - Fertilité agronomique et fertilité économique 1989.

Reboul, C., 1976. Mode de production et systèmes de culture et d'élevage. *Économie Rural*. 112, 55–65. doi:10.3406/ecoru.1976.2413

Ribeiro, P.F., Nunes, L.C., Beja, P., Reino, L., Santana, J., Moreira, F., Santos, J.L., 2018. A Spatially Explicit Choice Model to Assess the Impact of Conservation Policy on High Nature Value Farming Systems. *Ecol. Econ.* 145, 331–338.

doi:10.1016/j.ecolecon.2017.11.011

Ribeiro, P.F., Santos, J.L., Bugalho, M.N., Santana, J., Reino, L., Beja, P., Moreira, F., 2014. Modelling farming system dynamics in High Nature Value Farmland under policy change. *Agric. Ecosyst. Environ.* 183, 138–144. doi:10.1016/j.agee.2013.11.002

Ribeiro, P.F., Santos, J.L., Santana, J., Reino, L., Leitão, P.J., Beja, P., Moreira, F., 2016. Landscape makers and landscape takers: links between farming systems and landscape patterns along an intensification gradient. *Landsc. Ecol.* 31, 791–803.

doi:10.1007/s10980-015-0287-0

Rizzo, D., Marraccini, E., Lardon, S., Rapey, H., Debolini, M., Benoît, M., Thenail, C., 2013. Farming systems designing landscapes: land management units at the interface between agronomy and geography. *Geogr. Tidsskr. J. Geogr.* 113, 71–86.

doi:10.1080/00167223.2013.849391

Ruiz-Martinez, I., Marraccini, E., Debolini, M., Bonari, E., 2015. Indicators of agricultural intensity and intensification: a review of the literature. *Ital. J. Agron.* 10, 74.

doi:10.4081/ija.2015.656

- Santos, J.L., Moreira, F., Ribeiro, P.F., Canadas, M.J., Novais, A., Lomba, A., 2020. A farming systems approach to linking agricultural policies with biodiversity and ecosystem services. *Front. Ecol. Environ.* in press, fee.2292. doi:10.1002/fee.2292
- Schaller, L., Targetti, S., Villanueva, A.J., Zasada, I., Kantelhardt, J., Arriaza, M., Bal, T., Fedrigotti, V.B., Giray, F.H., Häfner, K., Majewski, E., Malak-Rawlikowska, A., Nikolov, D., Paoli, J.-C., Piorr, A., Rodríguez-Entrena, M., Ungaro, F., Verburg, P.H., van Zanten, B., Viaggi, D., 2018. Agricultural landscapes, ecosystem services and regional competitiveness—Assessing drivers and mechanisms in nine European case study areas. *Land use policy* 76, 735–745. doi:10.1016/j.landusepol.2018.03.001
- Shmueli, G., 2010. To Explain or to Predict? *Stat. Sci.* 25, 289–310. doi:10.1214/10-STS330
- Silva, J.F., Santos, J.L., Ribeiro, P.F., Canadas, M.J., Novais, A., Lomba, A., Magalhães, M.R., Moreira, F., 2020. Identifying and explaining the farming system composition of agricultural landscapes: the role of socioeconomic drivers under strong biophysical gradients. *Landsc. Urban Plan.* 202, 103879. doi:10.1016/j.landurbplan.2020.103879
- Strobl, C., Malley, J., Tutz, G., 2009. An introduction to recursive partitioning: Rationale, application, and characteristics of classification and regression trees, bagging, and random forests. *Psychol. Methods* 14, 323–348. doi:10.1037/a0016973
- Surová, D., Pinto-Correia, T., 2008. Landscape preferences in the cork oak Montado region of Alentejo, southern Portugal: Searching for valuable landscape characteristics for different user groups. *Landsc. Res.* 33, 311–330. doi:10.1080/01426390802045962
- Tóth, K., Kučas, A., 2016. Spatial information in European agricultural data management. Requirements and interoperability supported by a domain model. *Land use policy* 57, 64–79. doi:10.1016/j.landusepol.2016.05.023
- van de Steeg, J.A., Verburg, P.H., Baltenweck, I., Staal, S.J., 2010. Characterization of the spatial distribution of farming systems in the Kenyan Highlands. *Appl. Geogr.* 30, 239–253. doi:10.1016/j.apgeog.2009.05.005
- van Vliet, J., de Groot, H.L.F., Rietveld, P., Verburg, P.H., 2015. Manifestations and underlying drivers of agricultural land use change in Europe. *Landsc. Urban Plan.* 133, 24–36. doi:10.1016/j.landurbplan.2014.09.001
- Wilson, G.A., 2009. The spatiality of multifunctional agriculture: A human geography

perspective. *Geoforum* 40, 269–280. doi:10.1016/j.geoforum.2008.12.007

Wright, M.N., Ziegler, A., König, I.R., 2016. Do little interactions get lost in dark random forests? *BMC Bioinformatics* 17, 145. doi:10.1186/s12859-016-0995-8

1 Explaining farming systems spatial patterns: a farm-level choice
2 model based on socioeconomic and biophysical drivers

3

4 Abstract

5

6 CONTEXT: Efforts to bring together landscape analysis and farming systems have failed to explain the
7 drivers behind their spatial distribution. Since agricultural landscapes are an outcome of farmers'
8 decisions, understanding the role of socioeconomic and biophysical drivers of such decisions is
9 essential for policy-making targeting landscape-level provision of public goods and ecosystem services
10 from agriculture.

11 OBJECTIVE: Aiming to better understand the role of these drivers, we focused on a region dominated
12 by agricultural use, with extensive variability in biophysical and socioeconomic conditions. A typology
13 of farming systems was derived from spatially explicit farm-level data provided by the Portuguese
14 agency responsible for Common Agricultural Policy payments, for 2017. Farms were thoroughly
15 characterized through relevant biophysical and socioeconomic variables considered as potential
16 drivers of farming systems.

17 METHODS: A random forest approach was used to develop a farming system choice-model, dependent
18 on those biophysical and socioeconomic variables. Variable importance measures and partial
19 dependence plots were used to explore the role of these variables in explaining the spatial distribution
20 of farming systems and to predict spatial patterns at the landscape scale.

21 RESULTS AND CONCLUSIONS: Results showed that both biophysical and socioeconomic drivers play a
22 significant role in the spatial distribution of most agricultural systems. Its importance, however, varies
23 significantly across farming systems, being crucial for some and almost irrelevant for others. Farm size
24 and climate have proved to be the most relevant drivers for most farming systems. Overall, our
25 approach proved to be quite accurate in predicting patterns of farming systems at the landscape scale.

26 SIGNIFICANCE: The proposed framework has shown great potential as a tool to support information-
27 based policy design to improve agricultural landscape planning, by linking farm-level management
28 decisions with the provision of socially valued public goods from agriculture, perceived at the
29 landscape-level.

30

31

32 1. Introduction

33 Agriculture is a dominant land use in many parts of the world, resulting from human interaction with
34 nature over time. This interaction is mostly regulated by two main types of drivers: biophysical
35 (climate, soil, topography...) and socioeconomic (farm structure, characteristics of farmers, markets,
36 policies...). The way each of these drivers affects agricultural landscapes has attracted the interest of
37 researchers (Grigg, 2005; Hazell and Wood, 2008; Kristensen et al., 2016; Plieninger et al., 2016; van
38 Vliet et al., 2015), but many unanswered questions still persist (Plieninger et al., 2016; Wilson, 2009).
39 Advancing knowledge about the role played by each of these factors in shaping agricultural landscapes
40 can thus improve our understanding of human/environment interactions, allowing to anticipate farm
41 management decisions and supporting evidence-based public intervention (Levers et al., 2016; van de
42 Steeg et al., 2010).

43 Such issues have recently been raised in the context of the provision of public goods and
44 agroecosystem services in general, including biodiversity conservation (Landis, 2017; Schaller et al.,
45 2018). Much literature resort to aggregated data concerning land use or to agriculture intensification
46 or specialization indicators, privileging landscape-dynamics analysis over landscape regional
47 differentiation, and seldom take the farm as the unity of inquiry (Debolini et al., 2018; Ruiz-Martinez
48 et al., 2015). There has been, however, a recent surge pressing need to ~~in~~ the development of proposals
49 approaches to bring together linking landscape analysis and farming systems (FS) to understand
50 agricultural landscapes, which are able to establish the FS geography but struggle do not go into
51 explaining the drivers behind their spatial distribution (Andersen, 2017; Benoit et al., 2012; Martel et
52 al., 2019; Rizzo et al., 2013; van de Steeg et al., 2010). Indeed, considering the mismatch between the
53 farm-scale, where management decisions take place, and the landscape-scale, where ecosystem
54 services are perceived, landscape analysis can greatly benefit from a deeper understanding of the
55 factors that influence farm management decisions. Thus, understanding the multiple production
56 decisions of adjacent farmers and combining these decisions at the landscape-scale is key to explain
57 the landscape mosaic and the ecological disturbance regimes (fire, grazing, ploughing...) that shape
58 the habitats of wild species and the provision of diverse ecosystem services.

59 This study builds on the Santos et al. (2020) conceptual framework considering for. Conversely, such are
60 should be considered as exogenous to the FS and, therefore, as potential drivers of the FS choice (Silva
61 et al., 2020).

Formatted: Font: Italic

Formatted: Font color: Blue

62 The FS concept used in this study follows that proposed by Santos et al. (2020), according to which a
63 FS can be defined as a set of farms roughly practicing the same crops and agricultural activities, using
64 similar technological processes and input endowments. A key aspect in this concept is that only
65 variables resulting from farm management decisions are considered, when defining a FS; all variables
66 that may influence these decisions but do not result from them, at least in the short run (e.g. farm size
67 or fragmentation level, climate, slopes, market or policy), should be considered as exogenous to the
68 FS and, therefore, as potential drivers of the FS choice (Silva et al., 2020).

69 To explain the spatial distribution of FS, distinct groups of drivers can be considered according to
70 distinct disciplinary perspectives or theoretical approaches. The analysis of farm biophysical
71 endowments to explain spatial patterns of FS has largely been explored by geography and
72 geo-agronomy (Deffontaines, 2004; Deffontaines et al., 1995; Grigg, 2005; Lacoste et al., 2018).
73 Climate, soil, and slope are often considered to establish a range of restrictions to the choice of the
74 farming system. But FS are also dependent on farmland structure and social context. Farmland
75 structure covers an ensemble of constraints such as farm size, fragmentation and spatial composition
76 which potentially restrict farmers decisions (Grigg, 2005; Latruffe and Piet, 2014; Reboul, 1976; Ribeiro
77 et al., 2018). The influence of territorial socioeconomic context on FS location may be grounded in the
78 notion of local embeddedness, supported by local sociocultural, demographic and economic structures
79 (Canadas and Novais, 2014; Debolini et al., 2018).

80

81

82 Using farm-level data collected in 2017 in a large-scale study area, we developed an innovative
83 methodological approach to: 1) derive a spatially-explicit FS typology; 2) assess the role of
84 socioeconomic and biophysical factors in explaining the spatial distribution of those FS; 3) assess the
85 extent to which we can predict FS patterns based on biophysical and socioeconomic variables. Results
86 were used to discuss the role of these drivers on the choice of the FS and their potential to predict
87 landscape patterns, seeking to draw conclusions to better inform policy design for landscape-level
88 provision of public goods from agriculture and prediction of landscape patterns in face of biophysical
89 or socioeconomic changes.

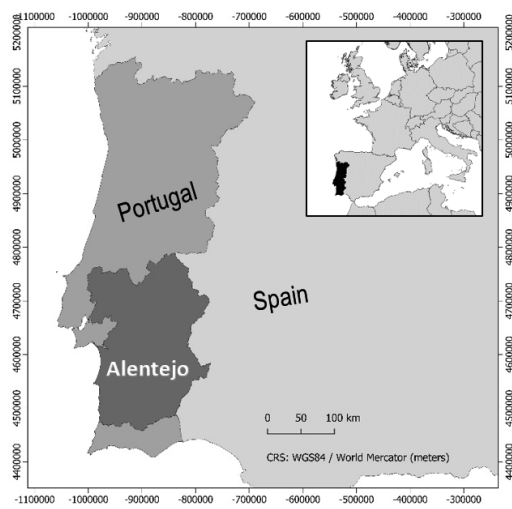
90

91 **2. Methods**

92 **2.1. Study area**

93 The study focused on the Alentejo region, in southern Portugal (Fig. 1), corresponding to the EU
94 statistical region PT18, at the NUTS2 level (Nomenclature of Territorial Units for Statistics). Covering
95 about 31,551 km² (ca. 1/3 of Portugal), the region has a Mediterranean climate, with hot dry summers
96 and mild rainy winters. The annual average temperature is about 16.3°C, ranging from 9.9°C to 23.4°C
97 in January and August, respectively, and the total annual rainfall is about 619 mm, largely concentrated
98 in the rainy season (approx. October to March). The relief is predominantly smooth (47% of the land
99 with slope < 5%; but 14% with slope > 15%), with few mountain areas (average altitude is 176 m a.s.l.,
100 ranging from 0 to 1020 m).

101



102

103 **Fig. 1 - Location of the study area in the Alentejo region (NUTS2), Portugal**

104

105

106 According to the latest agricultural census in Portugal (2009), the utilized agricultural area (UAA) in
107 Alentejo (NUT 2) was then ca. 2.2 million hectares, covering almost 70% of the region and making it
108 the dominant land use. Official statistics report that in 2016 the utilized agricultural area (UAA) was

Formatted: Caption, Left, Line spacing: single

109 dominated by permanent pastures (64%), followed by annual crops (24%) and permanent crops (11%).
110 Cereals, forages and olive groves were the main crops, with roughly equal shares of 8% in total UAA
111 (making ca. 70% of the UAA excluding permanent pastures). Nearly 40% of the UAA is under the canopy
112 of scattered trees, mainly cork and holm oaks (*Quercus suber* and *Q. rotundifolia* respectively),
113 originating an agroforestry system locally named "montado", which is largely acknowledge for its high
114 nature value (Ferraz-de-Oliveira et al., 2016). Cropland in these undercover areas are mainly
115 permanent pastures (70%) and annual crops (30%). Most of the UAA is rainfed (ca. 90%) and irrigated
116 areas are mostly located within state-promoted irrigation systems, often depending on large dams.
117 The region is dominated by large holdings, with almost 90% of the UAA in farms with more than 50 ha.

118

119 2.2. Farming systems identification

120 To build a farming systems typology for the study area we used data from the EU Integrated
121 Administration and Control System (IACS) for 2017, associated with spatially explicit farm parcel data
122 from the Land Parcel Identification System (LPIS), provided by the Portuguese agency responsible for
123 Common Agricultural Policy (CAP) payments. These data are collected on a yearly basis from farmers
124 declarations when applying for CAP payments and its usefulness for FS research has been
125 demonstrated by previous studies (Lomba et al., 2017; Ribeiro et al., 2018, 2016, 2014).

126 The raw dataset identified 26,648 CAP beneficiaries in the study area, covering a total of 2,221,816 ha
127 distributed over 208,338 parcels which, in turn, included 560,213 subparcels for which land use/crop
128 cover was described. Livestock declared by each beneficiary was also provided, describing species
129 composition, gender, age groups and an indication of whether they were kept in stables or grazing.

130 First, all parcels declared by the same CAP beneficiary were taken as a single farm. However, we found
131 that some beneficiaries reported very scattered parcels, sometimes separated by hundreds of
132 kilometres, where the farm concept (as an agricultural management unit) would not apply. In these
133 cases, we decided to regroup these parcels into new (sub)farms by forcing the distance between them
134 not to exceed 25 km, which increased the total number of farms to 28,739. This decision also helped
135 to narrow down the range of biophysical variability within each farm, and thus to better link farm units
136 to their biophysical context, described in the next section. We also discarded farms with total area
137 equal or below 2 ha (4409 farms, representing less than 1% of total UAA) because the land use in
138 smaller farms is likely to be highly sensitive to crop rotations, which cannot be properly captured with
139 one-year data.

Formatted: Font color: Blue

140 The raw data included 129 land use/cover categories, which were simplified by aggregation into
141 broader categories, while maintaining the distinction between irrigated and rainfed crops, when
142 applicable (e.g., ~~irrigated and all rainfed cereal crops were merged into a single category named~~
143 ~~"rainfed cereals"~~). We also included two variables describing the proportion of the UAA under the
144 cover of cork and holm oaks, respectively, because their presence is prone to influence farm
145 management, as the first is a major source of income for farmers (cork production) and the later
146 provides shade and food (acorns) to livestock grazing, in addition to valuable firewood. These two
147 variables were computed on a geographical information system (GIS) environment by intersecting the
148 farms map (derived from the LPIS spatial data) with digital information on cork and holm oak
149 distribution and computing, for each farm, the share of the UAA covered with both land cover classes.

150 Livestock numbers were converted into livestock units (LU) using EU standard conversion factors, and
151 these were used to describe the percentage composition of livestock by species, as well as livestock
152 density in each farm. Thus, a set of 28 variables was defined to characterize the land use/cover and
153 livestock patterns for each farm (Table 1).

154 A principal component analysis (PCA) was performed on a correlation matrix of these 28 variables to
155 reduce variable redundancy and the principal axes with eigenvalues above 1 entered a hierarchical
156 cluster analysis (Ward method) to derive the FS typology. The number of clusters to retain was decided
157 based on a visual analysis of the dendrogram and on expert knowledge of the study area.

158 To help interpreting the resulting FS, we calculated three variables indicating the level of agricultural
159 intensity, specialization and dependence on labour. The intensity variable was calculated following the
160 EU "standard output" approach (Commission Regulation (EC) No 1242/2008 of 8 December 2008) by
161 estimating the total gross product per land unit (in €/ha UAA) for each farm. The specialization variable
162 was computed as the highest proportion of standard output from a single farm activity. The labour
163 indicator aims to differentiate the FS based on their specific labour needs, in annual work units per
164 land unit (AWU/ha UAA). Due to data limitations, we had to resort to official statistics on the "EU farm
165 typology by economic size and type of farming" (in the sense of the above-mentioned legal text), at
166 NUT2 level (Alentejo) for the year of 2013, from which we extracted the number of annual work units
167 per hectare for each farm type, to be directly associated to each of the resultant FS on a similarity base.
168 Thereby, this indicator was not computed at farm level, but directly at FS level.

169

170 Table 1 – Summary statistics for the land use/cover and livestock farm characterization variables (n =
171 24313 farms)

Variable	Mean ± SD
<i>Land use/cover variables (proportion of total UAA)</i>	
Rice (<u>both Indica and Japonica</u>)	0.012 ± 0.1
Cereals Irrigated (<u>corn, wheat, oats, barley, triticale</u>)	0.018 ± 0.104
Cereals rainfed (<u>wheat, corn, oats, barley, rye and triticale</u>)	0.056 ± 0.165
Orchards (<u>orange, apple, plum, fig, loquat, cherries, blackberry, raspberry</u>)	0.013 ± 0.078
Forages Irrigated (<u>ryegrass, lucerne, silage maize, sorghum, vetch</u>)	0.006 ± 0.051
Forages Rainfed (<u>ryegrass, oats, corn, sorghum, lupine</u>)	0.049 ± 0.153
Horticultural (<u>potatoes, carrots, onions, cabbages, beans, chickpeas</u>)	0.017 ± 0.089
Industrial horticulture (<u>tomato and pepper</u>)	0.011 ± 0.092
Oilseeds (<u>sunflower and rapeseed</u>)	0.01 ± 0.067
Pastures (<u>temporary grass and permanent grasslands</u>)	0.511 ± 0.41
Fallows	0.043 ± 0.146
Olive groves Irrigated	0.034 ± 0.156
Olive groves Rainfed	0.171 ± 0.291
Vineyards	0.034 ± 0.145
Walnuts and almond trees	0.003 ± 0.048
Stone pine	0.009 ± 0.079
Other dry fruits (<u>hazelnut, chestnut, pistachios, carob</u>)	0.001 ± 0.019
Cork oak cover	0.149 ± 0.265
Holm oak cover	0.111 ± 0.229
<i>Livestock variables (proportion in total LU)</i>	
Cattle grazing	0.168 ± 0.34
Cattle stabled	0.003 ± 0.04
Fattening cattle grazing	0.018 ± 0.054
Fattening cattle stabled	0.002 ± 0.037
Sheep grazing	0.205 ± 0.386
Goat grazing	0.024 ± 0.131
Dairy cows	0.004 ± 0.047
Pigs grazing	0.008 ± 0.076
Livestock density (LU/ha UAA) (<u>includes all farm animals, added-up in LU</u>)	0.526 ± 3.506

Formatted Table

172

2.3. Socioeconomic and biophysical drivers

Potential socioeconomic and biophysical drivers of farming system choice were screened from literature (e.g. Grigg, 2005; Hazell and Wood, 2008; Kristensen et al., 2016; Martel et al., 2019; Plieninger et al., 2016; Reboul, 1989; van Vliet et al., 2015) and the authors' experience from previous

Field Code Changed

Formatted: Font color: Blue

173

174

175

176

177 studies where similar approaches were applied (Ribeiro et al., 2018, 2014; Silva et al., 2020).
 178 Subsequently, each farm was~~Each farm was~~ characterized according to a set of socioeconomic and
 179 biophysical variables thus identified, considered as potential drivers of FS spatial patterns (Table 2).
 180 These variables vary spatially but are mostly constant over time (at least for the time scale of most
 181 farm management decisions).

182 Socioeconomic variables included seven farm structure variables (farm and block size, farm
 183 fragmentation and dispersion, access to public and private water sources for irrigation, nature
 184 conservation constraints on farm use), and six local context variables computed from official statistics
 185 at the administrative parish level (one demographic variable, population density, and five agricultural
 186 variables, e.g. AWU availability or the share of rented UAA); all farms in the same parish were
 187 assigned the same value in these variables; when farms had areas in more than one parish, these
 188 variables were computed through average-weighting by farm-area shares in each parish. Biophysical
 189 variables included three climatic variables (describing temperature and precipitation), eight soil quality
 190 variables (describing soil depth, texture and pH) and three topographic variables (slope
 191 classes/categories). (Table 2).

192 Values for explanatory variables were derived for each farm using a GIS (maps for explanatory variables
 193 are provided in supplementary information, Fig-S1, Annex I). Farms with missing values resulting from
 194 map mismatches were discarded, dropping the number of valid observations to 23,416 farms.

195

196 Table 2 – Summary statistics for the socioeconomic and biophysical drivers (n = 23416 farms)

Variable	Description	Mean ± SD (Min-Max)
<i>Socioeconomic variables – farm structure variables</i>		
FSIZE	Farm size – Total UAA (ha) (1)	84.09 ± 184.46 (2.01-7191.16)
BLKSIZE	Average farm-block size (ha) (1)	23.15 ± 45.37 (0.20-1109.93)
JANUS	Januszewski index (adimensional) (1) (2)	0.65 ± 0.23 (0.13-1.00)
BLKDIST	Average area-weighted block distances to farm centroids (m) (1)	1571 ± 2128 (0-56951)
WPRIVATE	Access to water from private ponds or small streams (yes=1; no=0) (5)	0.16 ± 0.37 (0.00-1.00)
WPUBLIC	Proportion of UAA in public irrigation systems (6)	0.15 ± 0.31 (0.00-1.00)
NATURE	Proportion of UAA included in areas classified for nature conservation (7)	0.22 ± 0.39 (0.00-1.00)
<i>Socioeconomic variables – local socioeconomic variables</i>		

Formatted: Font color: Blue

INCAGRI	Proportion of farms where agriculture is the main household income source (3)	0.23 ± 0.14 (0.00-0.84)
INCOTH	Proportion of farms where household income is mostly from outside the farm, but not pensions (3)	0.26 ± 0.07 (0.00-0.67)
PDENS	Population density (inhabitants/km ²) (4)	32.5 ± 77.8 (0.89-1084.24)
AWU	Number of annual work units (AWU) per km ² of total parish area (3)	1.96 ± 1.70 (0.21-17.95)
AWU hired	Proportion of hired work in total labour (3)	0.26 ± 0.15 (0.00-0.93)
RENT	Proportion of rented land in total UAA (3)	0.18 ± 0.12 (0.00-1.00)
<i>Biophysical variables</i>		
TMIN	Average minimum temperature in the coldest month 1970-2000 (°C) (8)	4.71 ± 0.59 (3.01-8.40)
TMAX	Average maximum temperature in the warmest month 1970-2000 (°C) (8)	31.56 ± 1.95 (20.24-35.68)
PREC	Average annual rainfall 1970-2000 (mm) (8)	592.89 ± 107.28 (376.83-1195.51)
SDEPTH	Soil depth (cm) (5)	52.74 ± 29.80 (0.00-150.00)
SMOOTH	Proportion of UAA with smooth slopes (<5%) (5)	0.51 ± 0.32 (0.00-1.00)
MODERATE	Proportion of UAA with moderate slopes (5-16%) (5)	0.38 ± 0.24 (0.00-1.00)
STEEP	Proportion of UAA with steep slopes (>16%) (5)	0.11 ± 0.19 (0.00-1.00)
HEAVY_S	Proportion of UAA with heavy texture soils (5)	0.33 ± 0.37 (0.00-1.00)
MEDIUM_S	Proportion of UAA with medium texture soils (5)	0.42 ± 0.38 (0.00-1.00)
LIGHT_S	Proportion of UAA with light texture soils (5)	0.24 ± 0.36 (0.00-1.00)
VERYACID	Proportion of UAA with very acid soils (pH<5) (5)	0.27 ± 0.33 (0.00-1.00)
ACID	Proportion of UAA with acid soils (5<pH<6) (5)	0.41 ± 0.38 (0.00-1.00)
NEUTRAL	Proportion of UAA with pH neutral soils (6<pH<7) (5)	0.21 ± 0.30 (0.00-1.00)
ALKALINE	Proportion of UAA with alkaline soils (pH>7) (5)	0.11 ± 0.24 (0.00-1.00)

197 Sources: (1) Computed from LPI data; (2) Farm spatial fragmentation index, varying from 0 to 1 with higher values
198 indicating a higher degree of farmland consolidation (Januszewski, 1968); (3) Agricultural census 2009 - parish level; (4)
199 Population census 2011 - parish level; (5) EPIC WebGIS Portugal (<http://epic-webgis-portugal.isa.ulisboa.pt/>); (6) DGADR -
200 Direção-Geral de Agricultura e Desenvolvimento Rural (<http://sir.dgadr.gov.pt/expl-alentejo>); (7) ICNF – Instituto de
201 Conservação da Natureza e das Florestas (<http://www2.icnf.pt/portal/pn/ap>); (8) IPMA - Instituto Português do Mar e da
202 Atmosfera (<https://www.ipma.pt/pt/oclima/normais.clima/>)
203

204 2.4. Model design

205 We developed a random forest FS choice model to explore the **farm-level** relationships between the
206 typologies of FS derived from cluster analysis and the socioeconomic and biophysical variables.
207 Random forest is a popular machine learning method that can be used both for regression and

Formatted: Portuguese (Portugal)

208 classification, and is well-suited for high dimensional data (Strobl et al., 2009). Random forest use
209 bootstrap and aggregation (bagging), building multiple decision trees based on random subsets of the
210 data and using a random subset of predictor variables candidates for each node, in each decision tree
211 (Liaw and Wiener, 2002). On a classification problem, each observation is assigned to a class according
212 to the majority of votes from all trees. Both the number of trees and the number of predictor variables
213 sampled for each node are user-defined and can be used to tune the model. The mean out-of-bag
214 (OOB) error rate computed across all trees provides a measure of model prediction accuracy (Breiman,
215 2001). Random forests have been widely used in many scientific fields and have proved to be one of
216 the best machine learning techniques currently available, including for predictive modelling of spatial
217 and spatio-temporal data (Hengl et al., 2018).

218

219 2.4.1. Explaining spatial distribution of farming systems

220 Since we were firstly interested in exploring causal theories on the main drivers of FS spatial
221 distribution, rather than using the model to make predictions on new data (e.g. to assess scenarios of
222 policy or climate change), we tuned the model to optimize its average prediction accuracy across FS,
223 rather than maximizing the overall prediction power, by testing different stratified sampling
224 approaches to deal with anticipated unbalanced data (high variance in group sizes), (see details of
225 model parametrization in supplementary information – Annex II). At this stage, model overfitting
226 should not be an issue, since the focus was on explaining our training data, rather than the
227 generalization of the model (Shmueli, 2010).

228 With this modelling outset, all FS are assumed to be competing simultaneously for each farm and the
229 choice is made dependent only on variables that vary in space, while keeping constant the effect of
230 temporal variables (such as prices or policies). The effect of these temporal variables on the choices
231 observed in the study year cannot be estimated, as we only have one observation on FS choice for
232 each farm, that is: the choice observed in the study year 2017.

233 We used variable importance measures to assess the relevance of each predictor variable in the model
234 ~~and their ,which is computed by recording the error rate on the OOB portion of the data for each tree,~~
235 ~~and then repeat the calculation after permuting (shuffling) each predictor variable (Breiman, 2001).~~
236 ~~The difference between the two are then averaged over all trees and normalized by the standard~~
237 ~~deviation of the differences to obtain an adimensional measure of the mean decrease accuracy (MDA)~~
238 ~~for each predictor variable, computed both for the entire model and for each class (Liaw and Wiener,~~
239 ~~2002; Strobl et al., 2009).~~

Formatted: Not Highlight

Formatted: Font color: Blue, Not Highlight

Formatted: Not Highlight

240 ~~To examine the~~ marginal effect ~~of each variable~~ on each FS ~~was examined using~~~~we used~~ partial
241 dependence plots ~~(PDP)~~ (Friedman, 2001). ~~PDP indicate the model outcome in relation to each~~
242 ~~predictor variable, while considering the average effect of all other predictors in the model. These plots~~
243 ~~can show if the relationship between the target and a predictor variable is linear, monotonic or more~~
244 ~~complex~~ ~~–(supplementary information – Annex II).~~ We investigated the shape of the partial
245 dependence plots ~~-fitted functions of the top five variables-~~ for each class of the dependent variable
246 (that is, for each FS) to infer their role as drivers or constraints for each FS. In addition, we computed
247 the correlation coefficient between the level of farming intensity characterizing each FS with the
248 corresponding prediction accuracy rate obtained by the model, to test the hypothesis of a positive
249 relationship between the levels of this indicator and the degree of FS dependence on socioeconomic
250 and biophysical drivers.

251 All statistical analyses were carried out in R 3.4.1 (R Development Core Team, 2017). ~~Model estimation~~
252 ~~was performed with package “randomForest” (Liaw and Wiener, 2002). Partial dependence plots were~~
253 ~~conducted with package “pdp” (Greenwell, 2018).~~

Field Code Changed

Field Code Changed

254

255 2.4.2. Predicting spatial patterns of farming systems

256 On a following step, we focused on exploring the predictive capacity of the model in the choice of the
257 FS, based on the socioeconomic and biophysical variables described above. Since we were mostly
258 interested in predicting FS choice at the landscape-scale rather than at farm-scale, taking into account
259 the importance of landscape patterns for biodiversity and public goods delivery, we focused the
260 analysis on the model's ability to predict FS spatial patterns at a scale comparable to that of the
261 landscape (Andersen, 2017). For this purpose, the study area was divided into a random network of
262 hexagons of about 54,125 ha each, corresponding to a hexagon apothem of 12.5 km which was chosen
263 with reference to the 25 km threshold used to define the farms. ~~–These hexagons were then used as~~
264 ~~analysis units to compare, for each hexagon, the percentage distribution of the UAA by FS in the~~
265 ~~observed situation with that predicted by the model–. A hexagonal grid was preferred over a square~~
266 grid because it is less subject to bias from the edge effects when computing landscape metrics (Birch
267 et al., 2007). We rejected all hexagons with more than 66% of the area outside the LPIS data, due to
268 low significance for this purpose. In each hexagon, we calculated the difference between the observed
269 and predicted UAA shares for each FS and computed the half-sum of their absolute values. The average
270 of these results across all hexagons was interpreted as an estimate of the percentage of accuracy
271 obtained in model predictions, that is, the capacity of the model to predict spatial patterns of FS
272 composition at the landscape-scale. In addition, we also computed the determination coefficient (r^2)

273 between the observed and predicted values in each hexagon, taking its mean as a measure of the
274 quality of fit of the model. Model predictions were obtained by running the model on a random test-
275 set of the data with ca. 1/3 of the observations (farms), after estimating it in a train-set with the
276 remaining 2/3.

277

278 3. Results

279 3.1. Farming systems typology

280 A solution of 30 groups, representing farming systems, was selected from the cluster analysis. As some
281 groups included only a very small number of observations (farms), we anticipated potential problems
282 in the estimation of the predictive model and so we decided to eliminate groups with less than 0.7%
283 of the total number of observations, an arbitrary threshold mostly based on expert judgement. This
284 led to the removal of 8 non-representative FS, comprising 613 farms accounting for 3.1% of total UAA,
285 which were discarded for further analysis. Consequently, the final number of FS was set at 22 ([Table](#)
286 [3](#)).

287 By chance, these FS resulted equally divided into livestock-oriented systems and crop-oriented
288 systems. Both groups include similar shares in number of farms (51.5% and 48.5%, respectively),
289 although farms in livestock-oriented systems cover a much larger share of total UAA (78.2%) denoting
290 they are larger farms, on average.

291 Within the livestock systems, six are oriented to sheep, three to cattle, one to goats and one is mixed
292 with cattle and sheep. Among the six sheep-oriented systems, two are agroforestry grazing systems,
293 one associated with cork oak and the other with holm oak, a third one is related with open land
294 pastures, a fourth sheep system is mainly dependent on forage crops, the fifth depends both on
295 permanent pastures and forage crops, and the last is mostly a mixed-system combining rainfed olive
296 groves with sheep grazing. The three cattle systems also include two agroforestry grazing systems with
297 permanent pastures under the canopy of cork and holm oaks, respectively, and a third one depending
298 mainly on forage crops. The mixed cattle-sheep system is highly dependent on irrigated forages and
299 the last livestock-oriented system is the goat system, which is also a pasture-dependent grazing system
300 ([Table 3](#)).

301 Among the crop-oriented systems, five are dedicated to permanent crops, four to annual crops and
302 the last two refer to special situations, one including farms without livestock but with almost all UAA
303 under pasture, probably yearly rented to neighbours with cattle, and the other encompassing farms

304 with almost all UAA set to fallow. The permanent crops systems included two systems dedicated to
305 olive groves, one of which was irrigated and the other rainfed, one to vineyards, another to fruit trees
306 and the last one to stone pines (for pine nut production). The annual crops systems included two
307 rainfed systems, one dedicated to cereals and the other to cereals and oilseeds, one dedicated to
308 irrigated cereals and horticultural crops, and the last one to rice production (Table 3).

309 The average farming intensity across the 22 FS is about 1650 €/ha, with the Fruit trees system as the
310 most intensive, reaching ca. 12600 €/ha, and 15 systems below 1000 €/ha. Agricultural specialization
311 is relatively high, with more than half of the FS earning more than 80% of their standard output from
312 a single activity. Average farm specialization is higher in crop systems than in livestock systems (85%
313 and 74%, respectively), where most systems earn more than 90% from a single activity. Average labour
314 needs are also higher in crop systems than in livestock systems (0.039 and 0.004 AWU/ha, respectively,
315 i.e. nearly 10 times more), with a maximum of 0.259 AWU/ha found in the Irrigated cereals and
316 horticultural crops system and a minimum of 0.003 AWU/ha found in systems Pastures without
317 livestock and Fallows (Table 4).

318 Average farm size varies significantly across FS, with values going from ca. 11 ha in both rainfed olive
319 grove systems (with and without sheep) until over 200 ha, in cattle grazing – HO and – CO systems
320 (288 and 249 ha, respectively) (Table 4).

321 Almost 1/3 of all farms are included in only three FS, all with more than 2500 farms (systems Rainfed
322 olive groves, Pastures without livestock and Cattle grazing – CO). However, nearly 1/3 of total UAA is
323 concentrated in one single FS, the Cattle grazing – CO. The three cattle-oriented FS comprise more
324 than half of the total UAA (53.3%) (Table 4).

325

326 Table 3.a – Farming system description – Land cover composition (average values in proportion to the total UAA; values under 0.01 are omitted; values above 0.5
 327 are in bold)

Farming system	Rice	Cereals Irrigated	Cereals rainfed	Forages Irrigated	Forages Rainfed	Horticultural	Industrial horticulture	Oilseeds	Fallows	Pastures	Fruit trees	Olive groves Irrigated	Olive groves Rainfed	Vineyards	Walnuts and almond	Stone pine	Other dry fruits	UAA under Cork oak	UAA under Holm oak
Cattle grazing – CO*					0.039					0.884		0.026						0.329	0.082
Cattle grazing – HO*			0.020		0.035					0.906		0.020						0.048	0.590
Cattle grazing – forages		0.011	0.143		0.239	0.015		0.020	0.027	0.388		0.045	0.093					0.049	0.086
Grazing goats					0.027					0.923			0.023					0.314	0.207
Mixed Cattle and sheep - Irrigated forages		0.034	0.037	0.444	0.112	0.019		0.016	0.018	0.219	0.018	0.015	0.049					0.066	0.028
Sheep grazing – CO*					0.012					0.936			0.022			0.012		0.686	0.022
Sheep grazing – HO*			0.029		0.025				0.013	0.891			0.032					0.051	0.641
Sheep grazing - pastures			0.014		0.021					0.852			0.085					0.141	0.048
Sheep grazing - pastures and forages			0.169		0.139	0.015			0.027	0.437			0.156	0.030				0.056	0.036
Sheep grazing - forages			0.041		0.650				0.033	0.127			0.107					0.068	0.053
Rainfed olive groves with sheep					0.016					0.291			0.660	0.010				0.039	
Rainfed olive groves									0.016	0.099	0.013		0.823	0.018				0.011	
Irrigated olive groves			0.012						0.022	0.052		0.770	0.083		0.021			0.016	0.024
Vineyards		0.010			0.013				0.054	0.066	0.014	0.029	0.093	0.697				0.025	0.010
Fruit trees					0.014				0.025	0.243	0.548	0.012	0.083	0.020			0.025	0.142	0.040
Stone pine	0.038								0.019	0.189			0.023			0.713	0.000	0.249	0.010
Rice	0.850	0.012							0.039	0.068								0.024	
Irrigated cereals and horticultural crops		0.300	0.038			0.241	0.242		0.049	0.077								0.018	0.012
Rainfed cereals and oilseeds		0.087	0.300			0.038		0.430	0.063	0.030		0.011	0.026						0.018
Rainfed cereals			0.463		0.016	0.015		0.010	0.171	0.150			0.147					0.031	0.039
Pastures without livestock										0.785		0.015	0.152	0.012				0.082	0.088
Fallows			0.079		0.011	0.016			0.752	0.037	0.011		0.077					0.048	0.141

328 * CO – Under cover of cork oak; HO – Under cover of holm oak

330

331 Table 3.b – Farming system description – Livestock composition in livestock-oriented farming
 332 systems (average values in proportion to total LU; values under 0.01 are omitted; proportions
 333 above 0.5 are in bold) and livestock density

Farming system	Cattle grazing	Cattle stabled	Fattening steers grazing	Fattening steers stabled	Sheep grazing	Goat grazing	Dairy cows	Pigs grazing	Livestock density (LU/ha)
Cattle grazing – CO*	0.872		0.084		0.028				0.799
Cattle grazing – HO*	0.865		0.080		0.044				0.677
Cattle grazing – forages	0.859		0.083		0.041				0.967
Grazing goats					0.082	0.910			1.039
Mixed Cattle and sheep - Irrigated forages	0.480		0.073		0.300	0.069	0.077		0.618
Sheep grazing – CO*	0.144		0.010		0.799	0.043			0.245
Sheep grazing – HO*					0.963	0.028			0.387
Sheep grazing - pastures					0.973	0.017			1.004
Sheep grazing - pastures and forages					0.792	0.206			0.711
Sheep grazing - forages	0.049				0.709	0.236			0.378
Rainfed olive groves with sheep					0.974	0.024			1.212

334 * CO – Under cover of cork oak; HO – Under cover of holm oak

335

336

337 Table 4 – Characterization of farming systems according to the levels of farming intensity,
 338 specialization and labour needs, average farm size (in hectares of UAA and number of LU) and
 339 representativeness (in number of farms, UAA and LU)

	Characterization of farming systems			Average farm size		Representativeness					
	Intensity (10 ³ €/ha)	Speciali- zation (%)	Labour needs (AWU/ha)	UAA (ha)	LU (n.)	Number of farms		UAA		Livestock Units	
						Total	%	Total (10 ³ ha)	(%)	Total (10 ³ LU)	(%)
Cattle grazing – CO*	0.46	84.3	0.005	248.5	128.7	2515	10.6	625	31.1	323.8	36.5

Cattle grazing – HO*	0.31	84.0	0.005	288.2	157.0	1245	5.3	359	17.9	195.5	22.0
Cattle grazing – forages	0.75	68.3	0.005	186.1	88.1	463	2.0	86	4.3	40.8	4.6
Grazing goats	0.94	88.6	0.004	52.6	19.7	251	1.1	13	0.7	4.9	0.6
Mixed Cattle and sheep - Irrigated forages	1.14	75.0	0.005	58.8	24.3	171	0.7	10	0.5	4.1	0.5
Sheep grazing – CO*	0.24	54.6	0.004	89.0	19.0	2346	9.9	209	10.4	44.6	5.0
Sheep grazing – HO*	0.28	64.9	0.004	84.9	23.5	1391	5.9	118	5.9	32.6	3.7
Sheep grazing - pastures	0.71	84.1	0.004	54.5	22.4	1461	6.2	80	4.0	32.7	3.7
Sheep grazing - pastures and forages	0.79	70.6	0.004	52.8	18.7	745	3.1	39	2.0	13.9	1.6
Sheep grazing - forages	0.46	69.9	0.004	25.3	4.3	848	3.6	21	1.1	3.7	0.4
Rainfed olive groves with sheep	0.91	74.1	0.006	12.0	9.6	774	3.3	9	0.5	7.4	0.8
Rainfed olive groves	0.30	92.3	0.010	10.6	0.1	2626	11.1	28	1.4	0.3	0.0
Irrigated olive groves	1.45	93.0	0.023	82.4	0.8	864	3.6	71	3.5	0.7	0.1
Vineyards	1.84	90.3	0.050	24.3	0.5	928	3.9	23	1.1	0.4	0.0
Fruit trees	12.59	89.9	0.036	25.7	1.5	325	1.4	8	0.4	0.5	0.1
Stone pine	4.63	97.6	0.009	68.0	1.4	221	0.9	15	0.7	0.3	0.0
Rice	1.70	93.8	0.018	52.8	4.0	314	1.3	17	0.8	1.3	0.1
Irrigated cereals and horticultural crops	4.66	90.9	0.259	53.7	1.8	1070	4.5	57	2.9	1.9	0.2
Rainfed cereals and oilseeds	0.98	88.2	0.006	65.7	0.9	421	1.8	28	1.4	0.4	0.0
Rainfed cereals	0.42	82.0	0.010	33.8	0.4	1537	6.5	52	2.6	0.6	0.1
Pastures without livestock	0.20	61.5	0.003	48.5	8.5	2602	11.0	126	6.3	22.2	2.5
Fallows	0.59	54.1	0.003	20.8	0.0	582	2.5	12	0.6	0.0	0.0
Total	-	-	-	-	-	23700	100	2007	100	733	-

340 * CO – Under cover of cork oak; HO – Under cover of holm oak

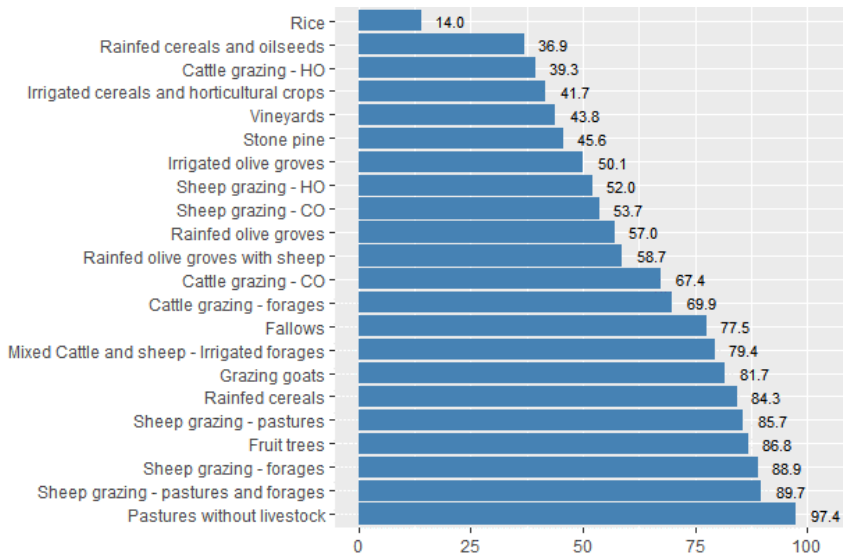
341

342 3.2. Spatial determinants of farming system choice

343 The tuning of the random forest model lead to a 500 trees model, with 5 variables randomly
344 sampled as candidates at each split and using the “sampsize” option to correct size differences
345 across the FS categories (see details in Annex II – supplementary information). The classification
346 error rates for each of the 22 FS ranged from 14.0% in the Rice system to 97.4% in the Pastures
347 without livestock system (Fig. 2 Fig. 1), with an average of 63.7% across all FS, a value that should

348 be evaluated positively considering the high number of classes in the dependent variable (22 FS,
 349 for which the random error rate would be about 95.4% with balanced data).

350



351

352 Fig. 2 - Classification error rates for the 22 farming systems (values in %)

Formatted: Caption, Left, Line spacing: single, Don't keep with next

353

354 The relative importance of socioeconomic and biophysical variables was very ~~elosesimilar~~, and
 355 among the top ten variables, in terms of mean decrease accuracy (Annex II – supplementary
 356 information), six are socioeconomic and four are biophysical. The farm physical dimension
 357 variables (FSIZE and BLKSIZE) and a local context of high dependence on family income in
 358 agriculture (INCAGRI) proved to be the most relevant socioeconomic factors influencing the
 359 choice of FS, while in the biophysical variables the most important were the climatic variables
 360 (Fig. 3 first plot in). The variable indicating access to surface water sources (WPRIVATE) was
 361 found to be the least important, either in the global model or in most of the class-specific
 362 models.

Formatted: Not Highlight

363

	MODEL OVERALL	Cattle grazing - CO	Cattle grazing - HO	Cattle grazing - forages	Grazing goats	Mixed Cattle and sheep - Irrigated forages	Sheep grazing - CO	Sheep grazing - HO	Sheep grazing - pastures	Sheep grazing - pastures and forages	Sheep grazing - forages	Rainfed olive groves with sheep	Rainfed olive groves	Irrigated olive groves	Vineyards	Fruit trees	Stone pine	Rice	Irrigated cereals and horticultural crops	Rainfed cereals and oilseeds	Rainfed cereals	Pastures without livestock	Fallows
FSIZE	65,6	↑	↑	↑	↑	↓	↑	↑	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
BLKSIZE	50,9	↑	↑	↑	↑	↓	↑	↑	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
JANUS	32,5	↑	↑	↑	↑	↓	↑	↑	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
BLKDIST	30,9	↓	↑	↑	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
INCAGRI	52,7	↓	↑	↑	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
RENT	47,3	↑	↑	↑	↑	↓	↑	↑	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
INCOTH	47,3	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
PDENS	46,0	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
AWU	45,7	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
WPUBLIC	43,7	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
AWU_hired	43,6	↑	↑	↑	↑	↓	↑	↑	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
NATURE	27,0	↑	↑	↑	↑	↓	↑	↑	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
WPRIVATE	8,7	↑	↑	↑	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
TMAX	56,4	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
TMIN	56,2	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
PREC	48,1	↑	↑	↑	↑	↓	↑	↑	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
SDEPTH	47,3	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
LIGHT_S	41,5	↑	↑	↑	↑	↓	↑	↑	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
ACID	41,2	↑	↑	↑	↑	↓	↑	↑	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
NEUTRAL	40,9	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
MEDIUM_S	40,1	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
VERYACID	38,1	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
SMOOTH	37,7	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
HEAVY_S	34,2	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
STEEP	31,8	↑	↑	↑	↑	↓	↑	↑	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
ALCALINE	31,7	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
MODERATE	29,9	↑	↑	↑	↑	↓	↑	↑	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓

Formatted: Keep with next
Formatted: Left: 0.98", Right: 0.98", Top: 0.79", Bottom: 0.59", Width: 11.69", Height: 8.27"

364
365
366
367
368
369

Fig. 3 - Variable importance for the overall model and for each farming system. Socioeconomic farm structure variables in blue; local-socioeconomic variables in orange; biophysical variables in green. Variables ordered by decreasing variable importance in the overall model and within each sub-group. Symbols ↑, ↓ and ↑ indicate whether the marginal effect of the variable in each farming system is mostly positive, negative or non-monotonic, respectively, based on the shape of the fitted function on the partial dependence plots (partial dependence plots are provided in supplementary information, Annex IV). Variable description in Table 2

Formatted: Caption, Left, Line spacing: single

Formatted: Font color: Blue

370

371 Farm size (FSIZE) and average farm-block size (BLKSIZE) were the most relevant variables for the
372 choice of Cattle grazing FS (Fig. 2), positively influencing its choice as revealed by the partial
373 dependence plots (Fig. 3 Fig. 3). The same variables also have a relevant effect on most sheep
374 systems but, in this case, predominantly on the opposite direction (Fig. 3 Fig. 3). The choice of
375 the Cattle grazing – CO system is positively influenced by the increase of the average annual
376 rainfall (PREC) and negatively by high summer temperatures (TMAX), which has a positive effect
377 on the choice of the Cattle grazing – HO system. The Cattle grazing – forages system is
378 distinguished by a preference for warmer winters, as revealed by the increasing trend of the
379 TMIN curve (Fig. 3 Fig. 3).

380 The Grazing goats system is positively related to sloping terrain; its choice is favoured by
381 increasing the slope (STEEP), while avoiding flat land (SMOOTH). This system is also
382 characterized by avoiding public irrigation areas (WPUBLIC) and deep soils (SDEPTH). The choice
383 of the Mixed Cattle and sheep - Irrigated forages system is favoured by deeper soils, public
384 irrigation systems (WPUBLIC) and high local labour availability (AWU), which is probably related
385 to the irrigated forages component of this FS or with the labour needs associated with grazing
386 herds. The average annual rainfall (PREC) has opposite effects in Sheep grazing – HO and – CO
387 systems, with the first system being favoured by lower rainfall values, and the other way around
388 in the later system. Sheep grazing – CO is also favoured by areas with steeper slopes (STEEP) and
389 light soils (LIGHT_S), while the choice of Sheep grazing – HO decreases with deeper soils,
390 smoother terrain and public irrigation structures. Lower values of local labour availability (AWU)
391 seem to promote the choice of the Sheep grazing – pastures system, while the choice of Sheep
392 grazing – forages system is negatively influenced as the local values of agricultural income
393 dependence (INCAGRI) raises (Fig. 3).

394 Both Rainfed olive groves systems (with and without sheep) are strongly related to smaller farm
395 sizes, as these are also the two systems with lower average UAA (Table 4). Both are positively
396 related to high summer temperatures (TMAX) and negatively to higher regional values of
397 agricultural income dependence. The Rainfed olive groves with sheep system is favoured when
398 average annual rainfall increases, and the Rainfed olive groves system is positively related to
399 neutral pH soils (NEUTRAL) (Fig. 3).

400 The Irrigated olive groves system is positively related to high summer temperatures, public
401 irrigation systems, high local labour availability and high average farm-block size. It is negatively
402 related to high average annual rainfall. The choice of the Vineyards system tends to increase

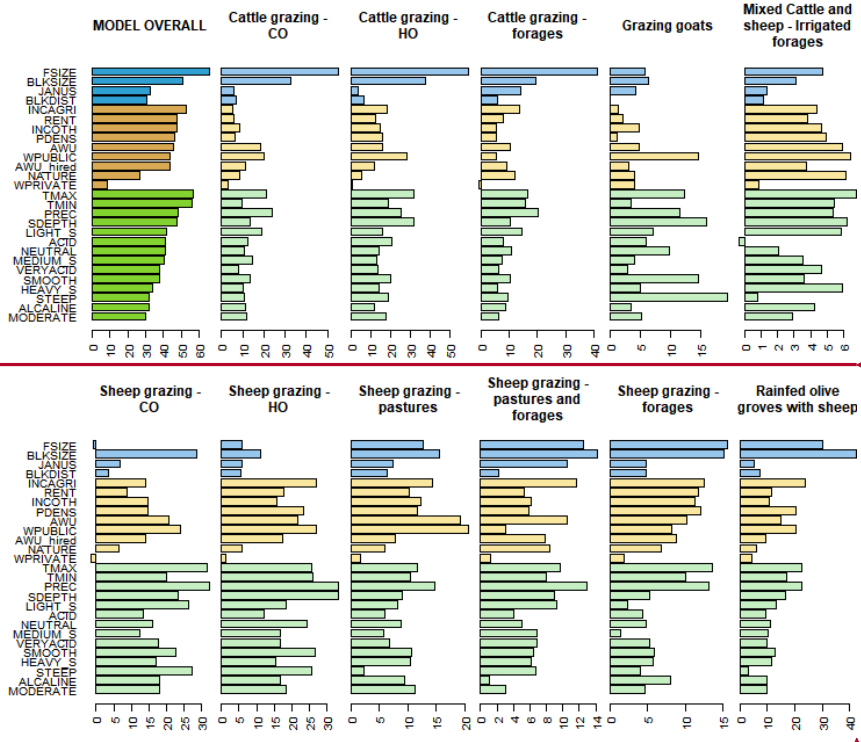
Formatted: English (United States)

403 with higher values of regional labour availability, public irrigation systems and population
404 density (PDENS). The Fruit trees system is positively associated with average annual rainfall and
405 negatively with high population density and warmer winters (TMIN). The choice of the Stone
406 pine system is favoured by light soils (LIGHT_S) and discouraged by high summer temperatures
407 and population density (Fig. 3).

408 In the annual crops, the Rice system is mostly favoured by the presence of public irrigation
409 systems, also by higher regional values of agricultural income dependence (INCAGRI) and soil
410 depth, while negatively influenced by high summer temperatures. The Irrigated cereals and
411 horticultural crops system is positively related to soil depth, regional labour availability and
412 smooth slope terrain. The choice of the Rainfed cereals and oilseeds system is encouraged with
413 public irrigation systems and higher values of soil depth, neutral pH and high summer
414 temperatures. The Rainfed cereals system is negatively related to bigger farm-block sizes and
415 average annual rainfall. The Pastures without livestock system seems to be promoted when
416 labour availability is lower and outside public irrigation systems, although this FS presented the
417 highest error rate (Fig. 2 Fig. 1), so the shape of the fitted function on the partial dependence
418 plots does not provide so clear relationships. The Fallows system also displays complex relations
419 with the predictors, though it seems to be more positively associated with small farms and areas
420 of low population density (Fig. 3).

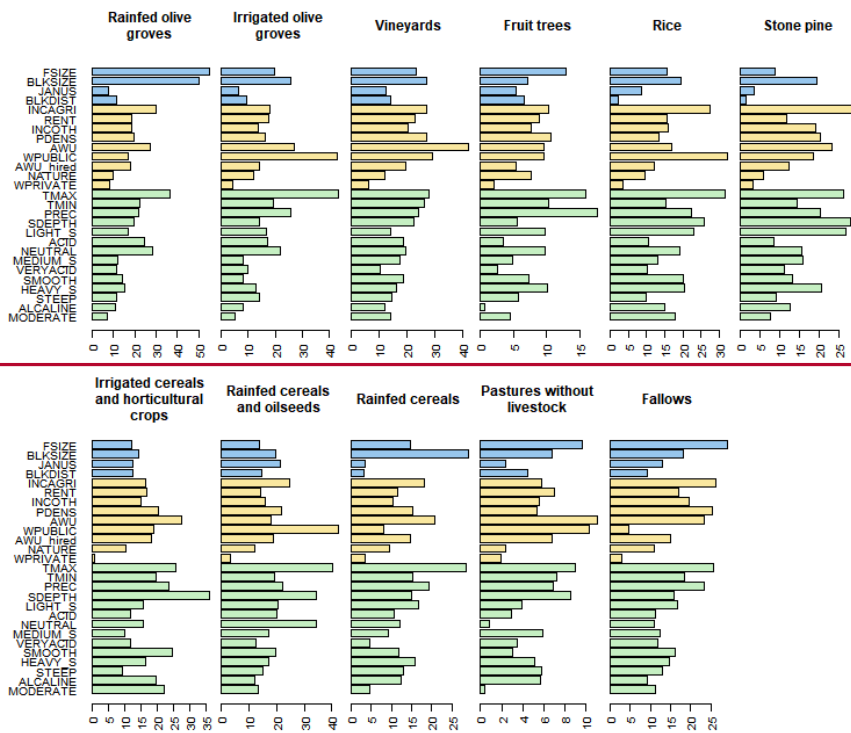
421 Finally, the prediction accuracy for the different farming systems (Fig. 2 Fig. 1) showed a modest
422 but positive correlation with the corresponding levels of agricultural specialization and labour
423 needs (Table 4) (correlation coefficients of 0.44 and 0.26, respectively), and a virtually non-
424 existent relationship with the level of agricultural intensity (correlation coefficient -0.03).

425



Formatted: Caption, Left, Line spacing: single, Don't keep with next

Formatted: English (United States)

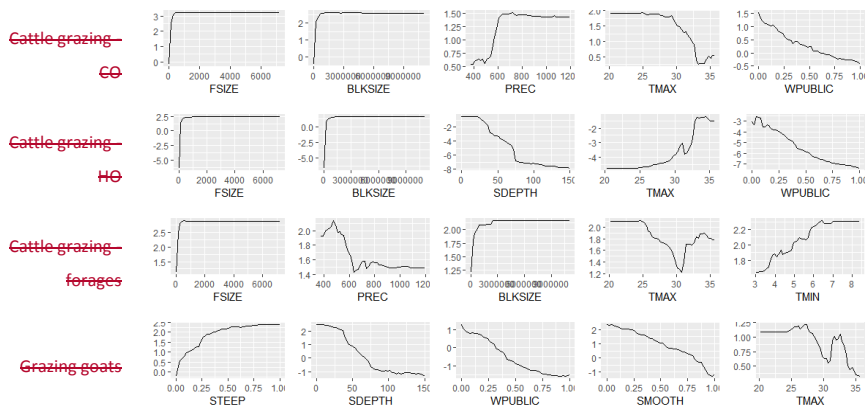


427

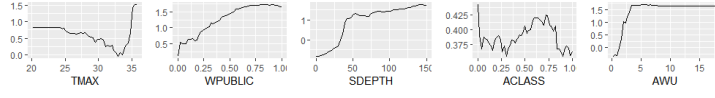
428 Fig. 2 (Continued)

429

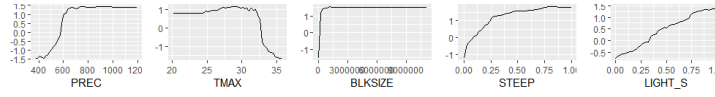
430



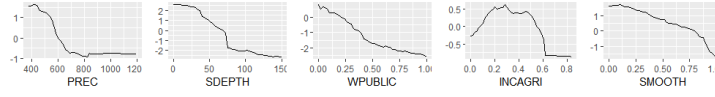
Mixed-Cattle and
sheep-Irrigated
forages



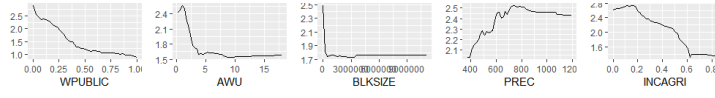
Sheep grazing-
EQ



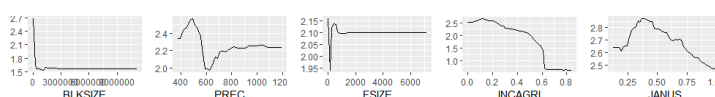
Sheep grazing-
HQ



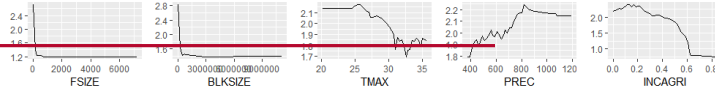
Sheep grazing-
pastures



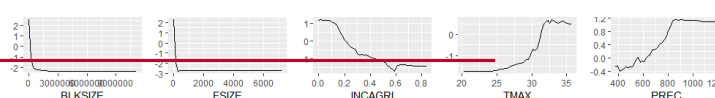
Sheep grazing-
pastures and
forages



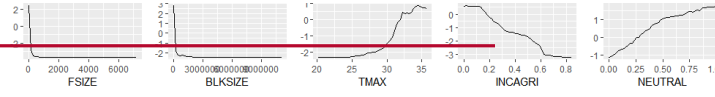
Sheep grazing-
forages



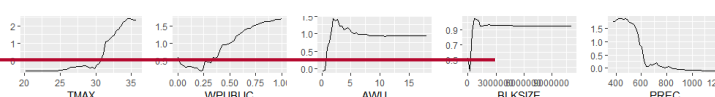
Rainfed olive
groves with
sheep



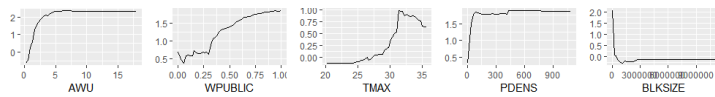
Rainfed olive
groves



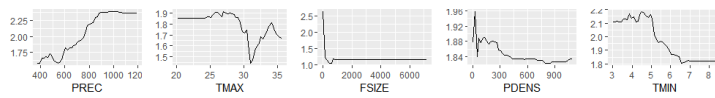
Irrigated olive
groves



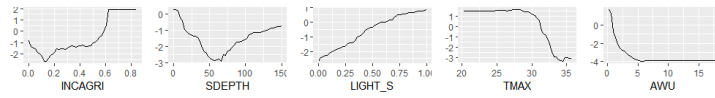
Vineyards



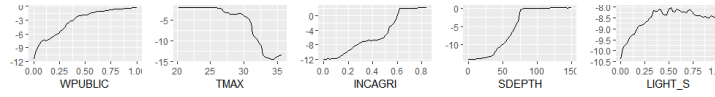
Fruit trees

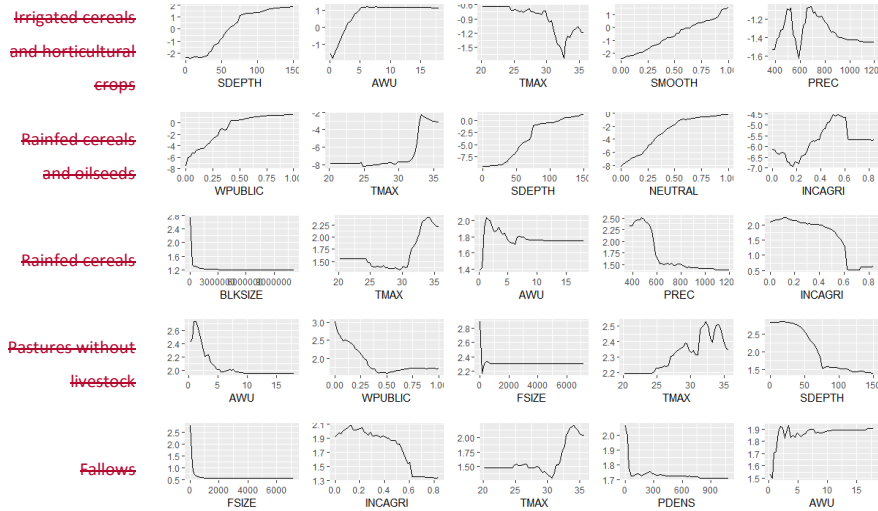


Stone pine



Rice



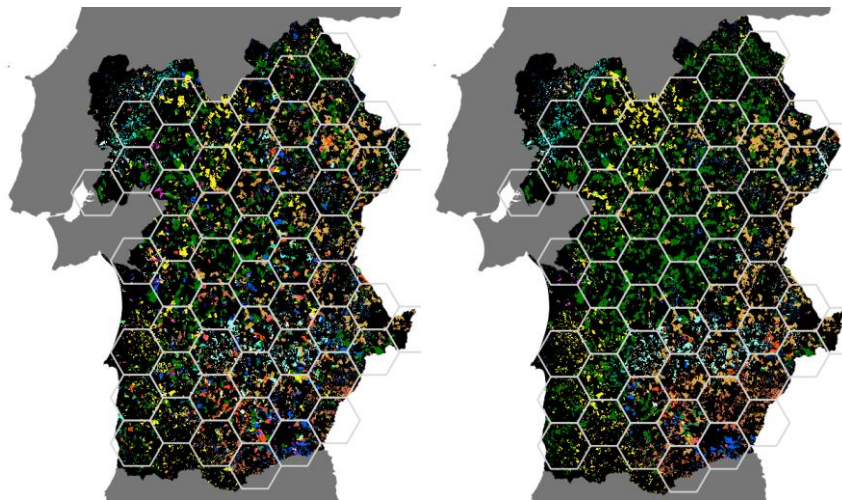


431

432 3.3. Spatial patterns of landscape-scale farming systems composition

433 The hexagonal lattice resulted with 56 usable analysis units, i.e., hexagons with >33% of the area
 434 overlapped with LPIS data (Fig. 4Fig-4). The average error rate in the FS spatial pattern
 435 predictions across all hexagons was 28.7% (max. 47.3%; min. 9.2%), which is substantially lower
 436 than the error rate obtained with model predictions at the farm-level (67.3%). The average
 437 coefficient of determination was 0.89 (max. 1.00; min. 0.28), revealing a good model fit.

438



Formatted: Centered, Keep with next

439

440 Fig. 4 - Observed (left) and predicted (right) FS maps for the 1/3 observations used in the
 441 model validation dataset and the hexagons network used to assess model accuracy in FS
 442 spatial patterns prediction (different colours identify distinct FS; detailed maps showing the
 443 spatial distribution of each farming system are provided in supplementary information,
 444 Annex III).

Formatted: Caption, Left, Line spacing: single

445

446 4. Discussion

447 The use of farm-level data (IACS) provided by the national CAP paying agency proved to be a
 448 suitable approach to derive the FS typology for the study area, in line with previous studies
 449 (Ribeiro et al., 2018, 2016, 2014). The spatial-explicit nature of these data (LPIS) allowed a very
 450 fine characterization of farms, including in their biophysical, structural and socioeconomic
 451 features. As expected, the extent and heterogeneity of the study area, in both socioeconomic
 452 and biophysical features, led to a broad typology of 22 farming systems, which are a direct
 453 outcome of distinct farm-management adaptive-responses to a variety of farm features and
 454 contexts.

455 Although the FS typology was balanced in terms of crop- and livestock-oriented systems, the
 456 results showed that most of the study area is currently devoted to livestock systems, particularly
 457 cattle grazing. Although the present study does not allow this to be confirmed, farmers'
 458 preference for these systems may be due to an (at date) ongoing direct payment for suckler
 459 cows (and partially to sheep and goats), a national agricultural policy option taken under the

460 2003 CAP reform that ~~led to significantly impacted~~ FS dynamics in the region (Ribeiro et al.,
461 2014).

462

463 4.1. Farm structure drivers

464 Many of the effects of structural socioeconomic variables observed here are consistent with
465 those of previous studies. For example, the farm-size was found to positively influence the
466 choice of extensive livestock systems over crop systems, which was also observed in Ribeiro et
467 al. (2018), and also in the choice between cattle grazing over some sheep grazing specialized
468 systems, which was also observed in studies by Ribeiro et al. (2014).

469 Access to private sources of surface irrigation water showed very little importance in the FS
470 choice-models, which is apparently odd for a region where water is often a limiting factor. This
471 was probably due not only to the type of variable used (dummy variable, with 1 = “yes, the farm
472 has access to surface water sources” and 0 otherwise) but also to the fact of not including access
473 to groundwater from water wells, due to lack of data, which are a common source in parts of
474 the region. In contrast, water availability from public irrigation systems is essential in explaining
475 the spatial location of several irrigated FS (either cereals, oil seeds or intensive olive groves and
476 vineyards) showing the importance of public water management policy over other biophysical
477 constraints (Kahil et al., 2015). Not surprisingly, these farming systems most associated with
478 large public irrigation systems are among the most intensive ones.

479 Public intervention in nature conservation areas seems to be of little relevance for FS choice
480 since although a considerable share of agriculture area is classified for nature conservation, the
481 corresponding variable (NATURE) was one of the least relevant within a list of dimensions that
482 has farm and block size at the top.

483 An interesting side-result of our approach was the insight of an overall negative, though
484 moderate, relationship between farm size and the level of agricultural intensity, indicating that
485 larger farms tend to adopt less intensive FS, a finding that goes back to earlier works (Cornia,
486 1985; Grigg, 2005; Reboul, 1989, 1976). Exceptions, however, can be found when contrasting,
487 e.g., the Rainfed olive groves and the Irrigated olive groves systems, where large investments in
488 fixed capital (including irrigation systems), together with labour availability, seem to provide
489 increasing returns to scale, which was also reported in more recent studies (Deininger et al.,
490 2018; Rada and Fuglie, 2019).

491

492 4.2. Socioeconomic context drivers

493 Regarding the socioeconomic context of the farms, the level of agricultural professionalization
494 (inferred from the INCAGRI variable) and farm labour availability proved to be significant drivers
495 of FS. On one side, higher levels of professionalization, which in Portugal are considerably low
496 in average when comparing to non-South European countries (Arnalte-Alegre and Ortiz-
497 Miranda, 2013), are positively associated with Rice, Stone pine or Rainfed cereals and oilseeds
498 systems. On the other side, Vineyards and Irrigated cereals and horticultural crops, which show
499 the highest levels of labour intensity per hectare and the highest average of labour units per
500 farm, are positively associated with local availability of farm labour. Considering that
501 horticultural crops typically have the highest wage labour ratios compared to other crops
502 (Baptista and Rolo, 2017), it was surprising that it did not show up associated with high local
503 proportion of hired labour. A possible explanation may be the high geographic mobility of hired
504 workers (Baptista and Rolo, 2017), although it may also emerge from the heterogeneity in labour
505 intensity within this FS, since it encompasses irrigated cereals and industrial horticulture, with
506 considerable levels of mechanization, as well as horticultural crops with very high levels of
507 labour needs.

508 The fact that local labour availability has a more widespread importance as a FS driver than rural
509 population density, which only stands out in the single case of Vineyards, contradicts the idea of
510 permanent crops and horticulture as able of promoting rural population retention (Egea and
511 Pérez y Pérez, 2016), i.e., it points to the dissociation between farm labour dynamics and local
512 demographics (Baptista and Rolo, 2017). While vineyards remain located in higher populated
513 parishes, following deep-rooted institutional constraints by protected designations of origin,
514 olive groves (either irrigated or rainfed) show no relation with local demographics.

515 Land renting (RENT) did not appear in the top 5 drivers in any FS, suggesting that the size of the
516 land renting market does not appear to have much effect on the choice of FS in the study area.
517 However, the positive relationship observed between land renting and livestock grazing FS,
518 especially cattle, ~~(supplementary information, Annex III)~~ suggests that these systems, which
519 have experienced marked growth in the region in recent years (Ribeiro et al., 2018), expanded
520 in part at the expense of this tenancy regime.

521

4.3. Biophysical drivers

As anticipated, biophysical factors related to climate, soils and relief, proved to be strong determinants of FS spatial distribution (Grigg, 2005). Summer heat and annual precipitation came up as the main biophysical drivers of FS spatial distribution in the study area. High summer temperatures seem to favour the choice of olive groves, vineyards, rainfed cereals and cattle grazing systems associated to Holms oak, and to discourage livestock systems associated to Cork oak, Stone pine or Rice systems. Winter cold increases the likelihood of fruit tree systems and the opposite with forage systems.

Deep soils and smooth relief are positive drivers of the Rice and Irrigated cereals and horticultural crops systems. The opposite effect is found towards the Grazing goats system, which is strongly related to steeper slopes. Soil pH did not emerge as a major driver for the distribution of any FS, except for rainfed cereals and olive groves systems which showed a preference for neutral pH soils.

Following Cork and Holm oak distinct preferences for soil and climate (Surová and Pinto-Correia, 2008), livestock systems associated with these two species of oaks were found distributed accordingly: Cork oak-associated systems prevail more to the coast and north of the study area, where summer temperatures are milder, annual rainfall is higher and soils are sandy and light-textured; Holm oak-associated systems are further inland and south, where summers are warmer, annual rainfall is lower and soils are frequently poor and fairly thin.

4.4. Farming system prediction at the farm and landscape levels

Although the model's ability to predict individual FS was quite varied, depending on the FS, when applied to predicting FS patterns at the landscape-level the model revealed a much higher hit rate. The random forest approach applied in the model estimation proved to be a valuable choice, particularly in dealing with such high dimensional data (Strobl et al., 2009).

~~The fact that some FS showed little relation to the selected socioeconomic and biophysical drivers, presenting very high error rates in model FS predictions, may be due to effects not controlled by the variables examined. One such case would be the Pastures without livestock system, whose choice is probably mostly determined by the presence of livestock farms in the nearby, with whom the farm can negotiate grazing land renting, rather than by the~~

~~characteristics of the farm itself. As a consequence, it presents an apparently random spatial distribution, little associated with the socioeconomic or biophysical features.~~

~~Farming systems with lower error rates in the model were those who most depend on the chosen socioeconomic or biophysical factors, such as the Rice, Irrigated cereals and horticulture or Rainfed cereals and oilseed systems (where cereals are an autumn-winter rainfed crop and oilseeds are grown in spring-summer season, often irrigated) that highly depend on irrigation water provided by public irrigation systems in this region. Or the Vineyards system, whose location is highly dependent on the availability of regional labour supply, to meet peaks of labour needs at certain times of the year, related to certain crop operations (e.g. harvesting or pruning). In the present market, policy and technological context, these FS revealed greater dependence on farm structure and “territorial embeddedness” (sensu Cerceau et al., 2018).~~

At the landscape level, the model was very effective in predicting farming systems patterns, i.e., the shares of FS composition within hexagon-shaped landscape units. For agricultural landscape planning focused on agroecosystem services provision, this may be the right scale of analysis, since a minimum share of farmland managed under the FS delivering those services should be sufficient to ensure the socially desired level of service, rather than requiring the service to be provided by a specific set of farms over a period of time (Andersen, 2017), as is typically the case with many agri-environment schemes requiring multi-annual contracts with individual farmers.

4.5. Shortcomings of the approach and recommendations for future research

~~Despite the valuable advantages evidenced by the proposed approach, some limitations werethere is still room for future improvement. identified, mostly related to shortcomings in the baseline data and peculiarities of the study area, some of which may call for future work. Improvements mostly relate to characteristics of the IACS and LPIS datasets and methodological options that are dependent on the geographic context of our study area.~~

~~While recognized as having high potential for supporting data-driven research, the In defiance of their high potential to develop research work like in the present study, IACS / LPIS datasets present have some limitations, such as the lack of information to carefully characterize farmers’ socioeconomic profile, or information on complementarity relationships between farms, such as the rental or sale of pastures, which can mislead the computation of farms’ stock density. Such information, which would be valuable to include in the FS choice models.~~

584 The fact that the empirical work was carried out in a region where the landscape is largely
585 dominated by agriculture, makes it possible to closely link FS choice with landscape modelling.
586 Where this is not the case, such as many mountain and less favoured regions across the EU, this
587 approach may not deliver the same results, given the smaller share of agriculture in the
588 landscape. Additionally, in such regions a significant part of agriculture is probably outside any
589 CAP support system, so that an approach based on IACS / LPIS data can only partially capture an
590 agricultural reality that is itself marginal at the landscape scale. Paradoxically, these regions
591 often include significant shares of high nature value farmlands at the EU level (Lomba et al.,
592 2014). Nevertheless, it should be worth trying to reproduce the approach in such regions in the
593 future, to test the generalization of the framework.

Formatted: English (United Kingdom)

Formatted: Font color: Blue

594 Because our farm characterization variables report to a single year, the effect of economic or
595 policy variables such as prices or subsidies can only be assumed as underpinning the farmers'
596 choices reflected on the observed 2017 IACS / LPIS data. However, the use of this type of
597 variables in the model, provided that time-series of farm-level data can be made available,
598 would significantly extend the scope of this approach, allowing its use to evaluate policy and
599 price change scenarios. Even without additional temporal data, the framework can take
600 advantage of the wide extension of the study area to perform, e.g., climate-change scenarios
601 assessment, by adopting a space-for-time substitution approach.

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

602 The selection of candidate variables to be tested as drivers of FS choice is also a key step in the
603 modelling approach. The misspecification or the absence of key variables can substantially
604 undermine models' performance. The problems observed with variable WPRIVATE may be one
605 such case, as this variable only reported access to small private surface water sources, which are
606 mostly torrential regime in this region, with insufficient water guarantees to encourage investing
607 in irrigation systems, and not taking into account that a significant portion of private irrigation
608 in this region is probably resorting to groundwater sources. This premise, which we could not
609 test due to lack of data, would be worth further investigation, should spatially explicit data on
610 groundwater uptakes becomes available.

611 Another issue deserving further investigation concerns the dimension of the grid of landscape
612 analysis units. It is possible that the size of these units (i.e. the hexagons, in the current case)
613 influences the accuracy of the model, so future investigation focused on determining its optimal
614 size could prove to be of high value.

615 Also, one aspect that has not been explored in the present study and should merit further
616 investigation is the occurrence of interaction effects between drivers. Although the way random

Formatted: Font color: Blue

617 forests deal with these effects is still subject to discussion (Wright et al., 2016), its likely existence
618 recommends additional analysis.

619 Finally, the fact that the prediction error rate has shown significant disparities across the FS
620 should not be seen as a flaw in the conceptual approach, but rather as an indication suggests
621 that the choice of some of these FS may be due to effects not controlled measured by the
622 variables examined, including factors related to farmers' desires, attitudes and motivations, or
623 with their socioeconomic profile which, as mentioned above, cannot be assessed on the basis of
624 IACS data. One such case would be the Pastures without livestock system, whose choice is
625 probably mostly determined by the presence of livestock farms in the nearby, with whom the
626 farm can negotiate grazing land renting, rather than by the biophysical characteristics of the
627 farm or its socioeconomic context. On the other hand, Farming systems FS with lower error rates
628 in the model were those who most depend on the chosen socioeconomic or biophysical factors,
629 such as the Rice, Irrigated cereals and horticulture or Rainfed cereals and oilseed systems (where
630 cereals are an autumn-winter rainfed crop and oilseeds are grown in spring-summer season,
631 often irrigated) that highly depend on irrigation water provided by public irrigation systems in
632 this region. Or The same applies to the Vineyards system, whose location is highly dependent on
633 the availability of regional labour supply, to meet peaks of labour needs at certain times of the
634 year, related to certain crop operations (e.g. harvesting or pruning). In the present market, policy
635 and technological context, these FS revealed greater dependence on farm structure and
636 “territorial embeddedness” (sensu Cerceau et al., 2018).

637

638

639 4.5.4.6. Concluding remarks

640 Our framework proved to be a suitable approach to investigate the role of human and physical
641 factors in farmers' decisions regarding the choice of the FS, providing effective contributions to
642 improve our understanding of the spatial distribution of FS when observed at a regional scale.

643 This research led to a better understanding of how each of the considered socioeconomic and
644 biophysical factors influences the spatial location of a wide range of FS, a subject seldom
645 explored in such detail in the literature. Results showed that both socioeconomic and
646 biophysical factors exert a high influence on the spatial distribution of FS, clearly revealing the
647 shortcomings of planning proposals exclusively confined to the agroecological aptitude

648 perspective (Nguyen et al., 2015; Pirovani et al., 2018). That influence, however, is not
649 comparable across FS, being decisive for the location of some FS and marginal for others.

650 Contrasting relationships were found between the agricultural intensity level and the degree of
651 dependence on biophysical drivers among the FS, with the simultaneous existence of intensive
652 FS with strong connection to biophysical factors (e.g. Rice system), and others similarly intensive
653 FS but where this relation is much weaker (e.g. Fruit trees system). This finding shows the
654 shortcomings of the assimilation between agricultural intensity and degree of artificialization of
655 the farm's conditions, largely dominant in the literature on the relationship between agriculture
656 and biodiversity/natural resources (Keenleyside et al., 2014). This assimilation ignores the
657 distinction between land and labour productivity and the fact that intensity differences may be
658 due to labour intensity levels rather than higher levels of external outputs. Our results point thus
659 to the need of not reducing farming systems diversity to an intensity gradient, when- comparing
660 across distinct productions (Ribeiro et al., 2016).

661 The use of IACS / LPIS data proved to be an invaluable asset for the research, enabling a high-
662 detailed farm-level analysis, not achievable using official statistics and usually only possible
663 through expensive and time-consuming farm surveys, often unfeasible for research works
664 developed at regional scales like the one used in this study. Therefore, it is worth renewing an
665 appeal previously made (Santos et al., 2020; Tóth and Kučas, 2016), addressed at the EU bodies
666 responsible for maintaining the IACS databases, to make them more accessible to the scientific
667 community, while safeguarding confidentiality duties.

Formatted: Font color: Blue

668 ~~Because our farm characterization variables report to a single year, the effect of economic or~~
669 ~~policy variables such as prices or subsidies can only be assumed as underpinning the farmers'~~
670 ~~choices reflected on the observed 2017 IACS/LPIS data. However, the use of this type of variables~~
671 ~~in the model, provided that time series of farm level data can be made available, would~~
672 ~~significantly extend the scope of this approach, allowing its use to evaluate policy and price~~
673 ~~change scenarios. Even without additional temporal data, the framework can take advantage of~~
674 ~~the wide extension of the study area to perform, e.g., climate change scenarios assessment, by~~
675 ~~adopting a space for time substitution approach. All these possibilities, coupled with~~Overall,
676 the model's ability to perform scenario simulations and to predict patterns of farming systems,
677 assigns this approach with a high potential to support information-based policy design to
678 improve agricultural landscape planning and ensure the provision of socially valued
679 agroecosystem services.

680

681 Acknowledgments

682 This work was funded by project “FARSYD– FARming SYstems as tool to support policies for
683 effective conservation and management of high nature value farmlanDs” – POCI-01-0145-
684 FEDER-016664 (PTDC/AAG-REC/5007/2014), supported by Norte Portugal Regional Operational
685 Programme (NORTE 2020), under the PORTUGAL 2020 Partnership Agreement, through the
686 European Regional Development Fund (ERDF). This research was also supported by the Forest
687 Research Centre, a research unit funded by Fundação para a Ciência e a Tecnologia I.P. (FCT),
688 Portugal (UID/AGR/00239/2019). FM was supported by FCT (contract IF/01053/2015). AL was
689 supported by national funds through FCT – Fundação para a Ciência e a Tecnologia, I.P., in the
690 context of the Transitory Norm - DL57/2016/CP1440/CT0001

691

5. References

Andersen, E., 2017. The farming system component of European agricultural landscapes. *Eur. J. Agron.* 82, 282–291. doi:10.1016/j.eja.2016.09.011

Formatted: English (United States)

Arnalte-Alegre, E., Ortiz-Miranda, D., 2013. The “southern model” of european agriculture revisited: Continuities and dynamics, *Research in Rural Sociology and Development*. Emerald Group Publishing Limited. doi:10.1108/S1057-1922(2013)0000019005

Formatted: English (United States)

Baptista, F.O., Rolo, J.C., 2017. Trabalho agrícola: percursos e modelos. *Cultiv. n.º* 10 27–35.

Formatted: English (United States)

Benoît, M., Rizzo, D., Marraccini, E., Moonen, A.C., Galli, M., Lardon, S., Rapey, H., Thenail, C., Bonari, E., 2012. Landscape agronomy: A new field for addressing agricultural landscape dynamics. *Landsc. Ecol.* 27, 1385–1394. doi:10.1007/s10980-012-9802-8

Birch, C.P.D., Oom, S.P., Beecham, J.A., 2007. Rectangular and hexagonal grids used for observation, experiment and simulation in ecology. *Ecol. Modell.* 206, 347–359. doi:10.1016/j.ecolmodel.2007.03.041

Breiman, L., 2001. Random Forests. *Eur. J. Math.* 45, 5–32. doi:10.1023/A:1010933404324

Canadas, M.J., Novais, A., 2014. Bringing local socioeconomic context to the analysis of forest owners’ management. *Land use policy* 41, 397–407. doi:10.1016/J.LANDUSEPOL.2014.06.017

Formatted: English (United States)

Cerceau, J., Mat, N., Junqua, G., 2018. Territorial embeddedness of natural resource management: A perspective through the implementation of Industrial Ecology. *Geoforum* 89, 29–42. doi:10.1016/j.geoforum.2018.01.001

Cornia, G.A., 1985. Farm size, land yields and the agricultural production function: An analysis for fifteen developing countries. *World Dev.* 13, 513–534. doi:10.1016/0305-750X(85)90054-3

Debolini, M., Marraccini, E., Dubeuf, J.P., Geijzendorffer, I.R., Guerra, C., Simon, M., Targetti, S., Napoléone, C., 2018. Land and farming system dynamics and their drivers in the Mediterranean Basin. *Land use policy* 75, 702–710. doi:10.1016/j.landusepol.2017.07.010

Deffontaines, J.-P., 2004. L’objet dans l’espace agricole. Le regard d’un géoagronome. *Natures Sci. Sociétés* 12, 299–304. doi:10.1051/nss:2004041

Deffontaines, J.P.P., Thenail, C., Baudry, J., 1995. Agricultural systems and landscape patterns: how can we build a relationship? *Landsc. Urban Plan.* 31, 3–10. doi:10.1016/0169-2046(94)01031-3

Deininger, K., Jin, S., Liu, Y., Singh, S.K., 2018. Can Labor-Market Imperfections Explain Changes in the Inverse Farm Size–Productivity Relationship? Longitudinal Evidence from Rural India. *Land Econ.* 94, 239–258. doi:10.3368/le.94.2.239

Egea, P., Pérez y Pérez, L., 2016. Sustainability and multifunctionality of protected designations of origin of olive oil in Spain. *Land use policy* 58, 264–275. doi:10.1016/j.landusepol.2016.07.017

Ferraz-de-Oliveira, M.I., Azeda, C., Pinto-Correia, T., 2016. Management of Montados and Dehesas for High Nature Value: an interdisciplinary pathway. *Agrofor. Syst.* 90, 1–6. doi:10.1007/s10457-016-9900-8

Formatted: English (United States)

Friedman, J.H., 2001. Greedy Function Approximation: A Gradient Boosting Machine. *Ann. Stat.* 29, 1189–1232. doi:10.1214/009053606000000795

Grigg, D., 2005. *An Introduction to Agricultural Geography, Second Edition*, Second edi. ed. Routledge, London and New York.

Hazell, P., Wood, S., 2008. Drivers of change in global agriculture. *Philos. Trans. R. Soc. B Biol. Sci.* 363, 495–515. doi:10.1098/rstb.2007.2166

Hengl, T., Nussbaum, M., Wright, M.N., Heuvelink, G.B.M., Gräler, B., 2018. Random forest as a generic framework for predictive modeling of spatial and spatio-temporal variables. *PeerJ* 6, e5518. doi:10.7717/peerj.5518

Januszewski, J., 1968. Index of land consolidation as a criterion of the degree of concentration. *Geogr. Plonica* 14, 291–296.

Kahil, M.T., Connor, J.D., Albiac, J., 2015. Efficient water management policies for irrigation adaptation to climate change in Southern Europe. *Ecol. Econ.* 120, 226–233. doi:10.1016/j.ecolecon.2015.11.004

Keenleyside, C., Beaufoy, G., Tucker, G., Jones, G., 2014. High Nature Value farming throughout EU-27 and its financial support under the CAP. Report Prepared for DG Environment, Contract No ENV B.1/ETU/2012/0035, Institute for European Environmental Policy. London. doi:10.2779/91086

- Kristensen, S.B.P., Busck, A.G., van der Sluis, T., Gaube, V., 2016. Patterns and drivers of farm-level land use change in selected European rural landscapes. *Land use policy* 57, 786–799. doi:10.1016/j.landusepol.2015.07.014
- Lacoste, M., Lawes, R., Ducourtieux, O., Flower, K., 2018. Assessing regional farming system diversity using a mixed methods typology: the value of comparative agriculture tested in broadacre Australia. *Geoforum* 90, 183–205. doi:10.1016/j.geoforum.2018.01.017
- Landis, D.A., 2017. Designing agricultural landscapes for biodiversity-based ecosystem services. *Basic Appl. Ecol.* 18, 1–12. doi:10.1016/j.baae.2016.07.005
- Latruffe, L., Piet, L., 2014. Does land fragmentation affect farm performance? A case study from Brittany, France. *Agric. Syst.* 129, 68–80. doi:10.1016/j.agry.2014.05.005
- Levers, C., Butsic, V., Verburg, P.H., Müller, D., Kuemmerle, T., 2016. Drivers of changes in agricultural intensity in Europe. *Land use policy* 58, 380–393. doi:10.1016/j.landusepol.2016.08.013
- Liaw, A., Wiener, M., 2002. Classification and Regression by randomForest. *R News* 2, 18–22.
- Lomba, A., Guerra, C., Alonso, J., Honrado, J.P., Jongman, R., McCracken, D., 2014. Mapping and monitoring High Nature Value farmlands: challenges in European landscapes. *J. Environ. Manage.* 143, 140–50. doi:10.1016/j.jenvman.2014.04.029
- Lomba, A., Strohbach, M., Jerrentrup, J.S., Dauber, J., Klimek, S., McCracken, D.I., 2017. Making the best of both worlds: Can high-resolution agricultural administrative data support the assessment of High Nature Value farmlands across Europe? *Ecol. Indic.* 72, 118–130. doi:10.1016/j.ecolind.2016.08.008
- Martel, G., Aviron, S., Joannon, A., Lalechère, E., Roche, B., Boussard, H., 2019. Impact of farming systems on agricultural landscapes and biodiversity: From plot to farm and landscape scales. *Eur. J. Agron.* 107, 53–62. doi:10.1016/j.eja.2017.07.014
- Nguyen, T.T., Verdoodt, A., Van Y, T., Delbecque, N., Tran, T.C., Van Ranst, E., 2015. Design of a GIS and multi-criteria based land evaluation procedure for sustainable land-use planning at the regional level. *Agric. Ecosyst. Environ.* 200, 1–11. doi:10.1016/j.agee.2014.10.015
- Pirovani, D.B., Pezzopane, J.E.M., Xavier, A.C., Pezzopane, J.R.M., de Jesus Júnior, W.C., Machuca, M.A.H., dos Santos, G.M.A.D.A., da Silva, S.F., de Almeida, S.L.H., de Oliveira Peluzio, T.M., Eugenio, F.C., Moreira, T.R., Alexandre, R.S., dos Santos, A.R., 2018. Climate change impacts on the aptitude area of forest species. *Ecol. Indic.* 95, 405–416.

Formatted: English (United States)

doi:10.1016/j.ecolind.2018.08.002

Plieninger, T., Draux, H., Fagerholm, N., Bieling, C., Bürgi, M., Kizos, T., Kuemmerle, T., Primdahl, J., Verburg, P.H., 2016. The driving forces of landscape change in Europe: A systematic review of the evidence. *Land use policy* 57, 204–214.

doi:10.1016/j.landusepol.2016.04.040

R Development Core Team, 2017. R: A language and environment for statistical computing. [WWW Document]. R Found. Stat. Comput. URL <http://www.r-project.org> (accessed 11.12.18).

Rada, N.E., Fuglie, K.O., 2019. New perspectives on farm size and productivity. *Food Policy* 84, 147–152. doi:10.1016/j.foodpol.2018.03.015

Formatted: English (United States)

Reboul, C., 1989. Monsieur le capital et madame la terre - Fertilité agronomique et fertilité économique 1989.

Reboul, C., 1976. Mode de production et systèmes de culture et d'élevage. *Économie Rural*. 112, 55–65. doi:10.3406/ecoru.1976.2413

Ribeiro, P.F., Nunes, L.C., Beja, P., Reino, L., Santana, J., Moreira, F., Santos, J.L., 2018. A Spatially Explicit Choice Model to Assess the Impact of Conservation Policy on High Nature Value Farming Systems. *Ecol. Econ.* 145, 331–338.

doi:10.1016/j.ecolecon.2017.11.011

Formatted: English (United States)

Ribeiro, P.F., Santos, J.L., Bugalho, M.N., Santana, J., Reino, L., Beja, P., Moreira, F., 2014. Modelling farming system dynamics in High Nature Value Farmland under policy change. *Agric. Ecosyst. Environ.* 183, 138–144. doi:10.1016/j.agee.2013.11.002

Formatted: English (United States)

Ribeiro, P.F., Santos, J.L., Santana, J., Reino, L., Leitão, P.J., Beja, P., Moreira, F., 2016. Landscape makers and landscape takers: links between farming systems and landscape patterns along an intensification gradient. *Landsc. Ecol.* 31, 791–803.

doi:10.1007/s10980-015-0287-0

Formatted: English (United States)

Rizzo, D., Marraccini, E., Lardon, S., Rapey, H., Debolini, M., Benoît, M., Thenail, C., 2013. Farming systems designing landscapes: land management units at the interface between agronomy and geography. *Geogr. Tidsskr. J. Geogr.* 113, 71–86.

doi:10.1080/00167223.2013.849391

Ruiz-Martinez, I., Marraccini, E., Debolini, M., Bonari, E., 2015. Indicators of agricultural intensity and intensification: a review of the literature. *Ital. J. Agron.* 10, 74.

Formatted: English (United States)

doi:10.4081/ija.2015.656

Santos, J.L., Moreira, F., Ribeiro, P.F., Canadas, M.J., Novais, A., Lomba, A., 2020. A farming systems approach to linking agricultural policies with biodiversity and ecosystem services. *Front. Ecol. Environ.* in press, fee.2292. doi:10.1002/fee.2292

Formatted: English (United States)

Schaller, L., Targetti, S., Villanueva, A.J., Zasada, I., Kantelhardt, J., Arriaza, M., Bal, T., Fedrigotti, V.B., Giray, F.H., Häfner, K., Majewski, E., Malak-Rawlikowska, A., Nikolov, D., Paoli, J.-C., Piorr, A., Rodríguez-Entrena, M., Ungaro, F., Verburg, P.H., van Zanten, B., Viaggi, D., 2018. Agricultural landscapes, ecosystem services and regional competitiveness—Assessing drivers and mechanisms in nine European case study areas. *Land use policy* 76, 735–745. doi:10.1016/j.landusepol.2018.03.001

Shmueli, G., 2010. To Explain or to Predict? *Stat. Sci.* 25, 289–310. doi:10.1214/10-STS330

Silva, J.F., Santos, J.L., Ribeiro, P.F., Canadas, M.J., Novais, A., Lomba, A., Magalhães, M.R., Moreira, F., 2020. Identifying and explaining the farming system composition of agricultural landscapes: the role of socioeconomic drivers under strong biophysical gradients. *Landsc. Urban Plan.* 202, 103879. doi:10.1016/j.landurbplan.2020.103879

Formatted: English (United States)

Strobl, C., Malley, J., Tutz, G., 2009. An introduction to recursive partitioning: Rationale, application, and characteristics of classification and regression trees, bagging, and random forests. *Psychol. Methods* 14, 323–348. doi:10.1037/a0016973

Surová, D., Pinto-Correia, T., 2008. Landscape preferences in the cork oak Montado region of Alentejo, southern Portugal: Searching for valuable landscape characteristics for different user groups. *Landsc. Res.* 33, 311–330. doi:10.1080/01426390802045962

Tóth, K., Kučas, A., 2016. Spatial information in European agricultural data management. Requirements and interoperability supported by a domain model. *Land use policy* 57, 64–79. doi:10.1016/j.landusepol.2016.05.023

van de Steeg, J.A., Verburg, P.H., Baltenweck, I., Staal, S.J., 2010. Characterization of the spatial distribution of farming systems in the Kenyan Highlands. *Appl. Geogr.* 30, 239–253. doi:10.1016/j.apgeog.2009.05.005

van Vliet, J., de Groot, H.L.F., Rietveld, P., Verburg, P.H., 2015. Manifestations and underlying drivers of agricultural land use change in Europe. *Landsc. Urban Plan.* 133, 24–36. doi:10.1016/j.landurbplan.2014.09.001

Wilson, G.A., 2009. The spatiality of multifunctional agriculture: A human geography

perspective. *Geoforum* 40, 269–280. doi:10.1016/j.geoforum.2008.12.007

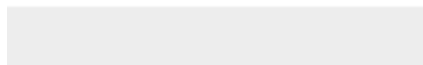
Wright, M.N., Ziegler, A., König, I.R., 2016. Do little interactions get lost in dark random forests? *BMC Bioinformatics* 17, 145. doi:10.1186/s12859-016-0995-8



[Click here to access/download](#)

Supplementary Material

[TitlePage_AS_PFR_22092020.docx](#)





[Click here to access/download](#)

Supplementary Material

[Revised_SupInformation_AS_PFR_21032021.docx](#)



Declaration of interests

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests: