

**European Political Boundaries as  
the Outcome of a Self-Organizing Process**

**Eric Weese**

**September 2016**

**Discussion Paper No.1629**

**GRADUATE SCHOOL OF ECONOMICS**

**KOBE UNIVERSITY**

**ROKKO, KOBE, JAPAN**

# European Political Boundaries as the Outcome of a Self-Organizing Process

Eric Weese\*

## Abstract

Political economy theories predict certain configurations of national boundaries, but these have not been calculated because of computational difficulties. Taking advantage of advances in mixed integer programming algorithms, we compute predicted political boundaries for Europe using a simple theoretical model taken from the literature: the size and arrangement of countries is determined by a tradeoff between efficiencies of scale and geographic heterogeneity. The model shows that the “natural borders” that lead to states emerging in certain configurations do not need to be particularly extreme, and a small number of these geographic features can influence the configuration of boundaries over a larger area. Our results show how real-world political boundaries can be described by a simple one parameter theoretical model that ignores many of the proximate causes of boundary changes.

## Significance

Why do international boundaries appear in certain places rather than others? Inspection of a map suggests they are not random, and theoretical models are plentiful. Applying any of these models to actual geographic data, however, presents a substantial combinatorial issue: if we have  $N$  small chunks of land that we wish to partition into countries, the number of possible ways of doing this grows superexponentially in  $N$ . We use modern mixed integer programming techniques and standard political economy theory to avoid this combinatorial issue, and produce simulated political boundaries for Europe. The simulation results qualitatively match historical boundaries, and include contemporary controversies: Northern Ireland is generally predicted to be a part of England, and Northern and Southern Italy are separate.

---

\*Department of Economics, Yale University. Email: eric.weese@yale.edu

“The history of Europe is that of its borders.” [1] From the historical perspective, current European international boundaries are the result of a combination of inheritance [2], war [3], nationalist sentiment [4], and other causes [5]. Given this enormous historical variation, it is clear that no single model of political jurisdiction formation will be able to completely explain observed political boundaries. We present a very simple model that describes international boundary formation in a democratic system, and show how it predicts boundaries that match some important features of the data. At the end of the paper, we discuss some reasons why this model might also give reasonable results for boundaries determined in an autocratic system.

A model of democratic political boundaries consists of two parts. First, there needs to be a description of the relative desirability of various configurations of political jurisdictions, and how this varies across different individuals. Second, there needs to be a rule for how political boundaries are determined given these individual preferences. We use existing models [6, 7, 8, 9] in the political economy literature for both of these parts. This approach differs substantially from an older literature based on Voronoi diagrams [10]. This previous literature considered where boundaries might form *given* a certain set of capital cities or other generating points for the countries in question. In contrast, the model in this paper does not assume a set of generating points, instead considering what countries might form given characteristics of the overall population in question.

We begin by presenting our model of the desirability of different configurations of political jurisdictions. The starting assumption is that each individual in Europe is located at a fixed physical location, will belong to exactly one country, and cares only about the characteristics of the country to which they belong. This is obviously an extreme simplification: it ignores, for example, migration, wars between countries, and the possibility that an area might be governed by multiple countries or none. Our theoretical model is based on individuals; however, for computation, we will aggregate similar individuals into small geographic units, and conduct simulations based on these units.

The first ingredient of our model of desirability is based on the idea that a jurisdiction with a small population is expensive to operate in per capita terms, and thus there are advantages to larger jurisdictions. In the economics literature, this is commonly referred to as *efficiencies of scale*. We use the standard form[11]:

$$C(P) = F + VP. \tag{1}$$

Here  $C(P)$  is the basic cost of operating a jurisdiction with a population of  $P$  people. There is a fixed cost,  $F$ , that does not depend on the size of the jurisdiction, and a variable cost,  $V$ , that is proportional to population. It will be easier to work with the per capita version of  $C$ :

$$C(P)/P = F/P + V. \tag{2}$$

The units of  $F$  can either be money or time, where the interpretation in the latter case is that in order to provide services the government requires a certain number of person-hours of labour each year. The functional form chosen for  $C$  is a very simple one, but there is some evidence[12] that it is close to the true functional form. We assume that cost is borne equally by all individuals in the jurisdiction, and thus each of them is responsible for the same share  $C(P)/P$  of the total cost.

A model considering only the basic per capita cost  $C(P)/P$  of providing government services would predict trivial boundaries for Europe: per capita cost is always decreasing in population, and thus all of Europe should be a single country. We observe a more complicated set of boundaries in Europe, however, and qualitative evidence also suggests that large countries face certain difficulties, including in collecting information and in providing services. In the literature, the additional difficulty faced by large countries is usually described as a *heterogeneity* cost. A functional form often used to model this cost is

$$L_i(S) = P_S^{-1} \sum_{i' \in S} \ell_{i,i'}, \tag{3}$$

where  $S$  is the set of individuals making up a jurisdiction,  $P_S$  is the total number of individuals in the jurisdiction, and  $\ell_{i,i'}$  is the distance between individuals  $i$  and  $i'$ .  $L_i(S)$  is thus the average distance between individual  $i$  and a randomly selected individual in the jurisdiction described by  $S$ . This  $L$  function originated in the linguistics literature[13, 14], where  $\ell_{i,i'} = 1$  if individuals  $i$  and  $i'$  speak the same language, and 0 if they do not. In this paper, rather than a  $\{0, 1\}$  coding, we let  $\ell_{i,i'}$  be the geographic distance between  $i$  and  $i'$ . We calculate these distances using a shortest path algorithm, with a few exceptions restricting travel to land rather than water. Details are provided in the Supplemental Methods section. A variety of similar distance concepts have been used in the applied political economy literature[15, 16]: we choose this one because it substantially simplifies computation.

Combining the two terms just presented, we see that the total benefit to individual  $i$  in the case where they are a part of jurisdiction  $S$  is

$$U_i(S) = -F/P_S - V - \gamma P_S^{-1} \sum_{i' \in S} \ell_{i,i'}, \quad (4)$$

where  $\gamma$  is a parameter describing the relative importance of heterogeneity, and everything has been multiplied by  $-1$  because it is standard to discuss individual decisions in terms of their benefits but both terms being considered are costs.

Given this model of the benefit to individual  $i$  of belonging to country  $S$ , we consider what countries will emerge given the preferences described by  $U$ . We use an algorithm based on ideas from matching theory[17]. Begin with an arbitrary starting partition  $\pi_0$  (we will use as  $\pi_0$  the partition where there is a separate country for every individual, but this choice is not important). Consider then whether there is any alternative country  $S' \notin \pi_0$  such that, for every individual  $i \in S'$ ,  $U_i(S')$  is higher than the benefit  $i$  was receiving from the country they were a part of in  $\pi_0$ . If so, then  $S'$  will form: create a new partition  $\pi_1$  by adding  $S'$  to  $\pi_0$  and removing or modifying other countries as necessary. Then repeat this process using  $\pi_1$  instead of  $\pi_0$ , and continue in this fashion until no  $S'$  can be found. There is no

theoretical guarantee that this algorithm will terminate[18, 19, 20], but in previous work[21] we find that it does. Additional details are provided in the Supplemental Methods section.

The intuition for this algorithm is that  $S'$  is a (potential) country where all of its constituents would want to leave whatever country that they are currently a member of, and form  $S'$  instead. Although not always followed in practice, ideals of “self determination” suggest that we should see partitions where this sort of  $S'$  does not exist. Partitions that have this property are known as *core partitions* in economics, and received considerable attention. At each intermediate step in the algorithm described above, there are often many potential choices for  $S'$ : different core partitions can be found by using a different rule for selecting  $S'$  in cases where there are multiple options available.

The model just described depends crucially on the distribution of population across Europe; however, this population distribution has experienced recent dramatic shifts due to urbanization. We would like detailed historical data regarding population distribution for all of Europe, but this sort of data is simply not available. We will thus use the agricultural suitability of land[22], shown in Supplemental Figure 3 as a proxy for historical population. This proxy is reasonable because historically the vast majority of the population was rural and agricultural, and was living in Malthusian conditions where the carrying capacity of a location determined the population living in it.<sup>1</sup>

In theory, the model just presented could be applied directly to the grid square data on agricultural suitability. In practice, however, serious numerical issues arise when considering grid squares with low suitability. We thus aggregate grid squares up to the polygons shown in Supplemental Figure 4, as described in the Methods section. This aggregation somewhat complicates the statistical interpretation of the results, as discussed in the Statistical Analysis appendix, but it is required for simulation to be feasible at all.

---

<sup>1</sup>An additional advantage of using agricultural suitability data is that it is less vulnerable to reverse causality: actual historical population distributions might be concentrated near the center of actual countries because capitals and the resulting administrative infrastructure are generally geographically central. In this case, actual population data would very successfully predict actual country boundaries, but for a reason completely unrelated to that of the model proposed in this paper.

Simulation also requires a choice of the parameter  $\gamma$ , describing the relative importance of heterogeneity. Estimation of this parameter would be a challenging undertaking, and there is no generally accepted method. Rather than choose a complicated estimation method, which might arouse suspicion of having specially selected a “good” value for  $\gamma$ , we perform a simple order of magnitude calibration. With distance in km, and the units for agricultural suitability being “fraction of a grid square” taken directly from [22], we find that  $\gamma = 0.1$  gives an average of 23 simulated countries, while  $\gamma = 0.01$  gives 3, and  $\gamma = 1$  gives 113. We use  $\gamma = 0.1$ .

The final piece of data required is actual European boundaries, in order to evaluate the model. We use boundaries between 1000 and 2000 CE, at 100 year intervals, as provided by EurAtlas. These boundaries are shown in Figure 1. We use average boundaries over 1000 years, rather than current boundaries, because the model attempts to describe “natural borders”, that should repeatedly appear in the data, rather than “accidents of history”, such as internal Soviet and Yugoslavian boundaries that were never intended to serve as international borders but now do, due to *uti possidetis*[23].

Figure 2 shows the boundaries predicted by the model, averaging over 100 simulation runs. Supplemental Figures 5 and 6 show examples of the partitions generated in these runs. While simulated boundaries do not match the actual boundaries perfectly, key features are reflected.<sup>2</sup> Most notable are the mountain ranges (the Pyrenees, Alps, and Carpathians) shown in Supplemental Figures 7 and 8 that appear as boundaries in Figure 2, despite the fact that elevation does not directly enter into the model. In the model, the mountain ranges form natural boundaries not because they are difficult to cross, but simply because there is no agricultural land there: a country that spanned both sides of the mountains would have high geographic heterogeneity, and thus would tend not to be part of a core partition.<sup>3</sup>

---

<sup>2</sup>Quantitative tests based on additional simulations are presented in the Statistical Analysis appendix.

<sup>3</sup>In independent work, “Natural Borders” by S. Kitamura and N. Lagerlöf examines the relationship between mountains, rivers, and other geographic features, and political and ethnic boundaries. Despite the immediate intuitive reaction, it is not obvious that the Alps are a natural boundary because they are difficult to cross. First, note that the Ligurian Sea and its surroundings are navigable, and thus the Alps can be bypassed by coastal travel. Second, Northern Ireland is part of the same country as England both in the

Supplemental Figure 9 shows an overlay of the simulated boundaries on the agricultural suitability data.<sup>4</sup>

Some additional boundaries appear in both Figures 1 and 2, but do not correspond to obvious features in Supplemental Figure 3. These include the (historical) boundary between Germany and Poland, and the boundaries further east than the Carpathians. These boundaries appear to emerge from the simulations because of the natural boundaries just discussed: if one country begins at the Carpathians, and the model leads to countries being a certain size given agricultural suitability, then Carpathians also lead to another country tending to emerge at another point further to the east. Similarly, some boundaries for Poland emerge not because there are any natural features exactly at those boundaries, but rather because mountains further to the south lead to a certain configuration tending to form. In general, areas with low agricultural suitability lead not only to a boundary in the immediate vicinity, but lead to certain patterns of boundaries further away.<sup>5</sup>

Notable failures of the model are France and the Ukraine (both split into three pieces), and Portugal (far too large). On the other hand, Denmark exists reliably, although the boundary varies simulation by simulation. Yugoslavia appears fairly clearly in the simulation results, despite the fact that it has been regarded as an artificial country. The simulation results suggest that the country's creation may have been reasonable in theory, even if it ultimately failed in practice. The results for France and the Ukraine are interesting given that the former is frequently given as an example of a strong homogenizing state[26, 27], and the later has a troubled territorial history[28].

---

data and in many of the simulation results: thus, countries that are not connected by land can and do exist. A more sophisticated model could be proposed, that would include the relatively small size of Ireland, and perhaps its military weakness, and contrasting that with the relative strength of France, and the difficulty of a maritime invasion. The point of the simulations in this paper is not to prove that the height of the Alps is not important: rather, they merely suggest that another explanation is also available.

<sup>4</sup>Earlier empirical work includes [24]. The relationship between geography and the number of ethnic groups has previously been examined using a different technique[25].

<sup>5</sup>Some other simulated boundaries are interesting, but are dependent on the details of model implementation. These include Northern Ireland being a part of England, and there being a split between Northern and Southern Italy. These are both dependent on the particular treatment of water barriers, as discussed in the Supplemental Methods section, and these aspects of the simulation results might be different if a different rule were used.

The simulation results just presented used a model based on the idea of self-determination. For much of European history, however, boundaries were decided by wars between autocratic leaders, interested mainly in collecting taxes from the land they controlled. Appendix A.2 shows that the model presented above also describes the tradeoff that would be faced by despots attempting to collect a tax by travelling to and from a capital located at the centroid of their territory.

To formally model despotism, a change would also be required to the algorithm used to simulate partitions, in order to specify when a new despot would attempt to start his own country. Intuitively, however, the “self-determination” algorithm presented tends to generate coalitions  $S'$  which, in a despotic setting, would be associated with a new despot who is better positioned to rule the population in  $S'$  than the existing despots. It is not obvious how to formalize this intuition, but it suggests why a model based on democratic principles matches data that contains a large number of autocratic jurisdictions. Further research may yield a formal model of autocratic governments for which simulation of boundaries is computationally feasible.

## References

- [1] Pomian, K. *L'Europe et ses nations*. D'ebat, historie, politique, societ e (Gallimard, Paris, 1990).
- [2] Ganshof, F. L. On the genesis and significance of the Treaty of Verdun (843). In *The Carolingians and the Frankish monarchy*, Studies in Carolingian history (Longman, London, 1971).
- [3] Macmillan, M. O. *Paris 1919: six months that changed the world* (Random House, New York, 2002), 1st u.s. ed edn.
- [4] Beissinger, M. R. *Nationalist mobilization and the collapse of the Soviet State*. Cambridge studies in comparative politics (Cambridge University Press, Cambridge, UK ; New York, 2002).
- [5] Toschi, U. The Vatican City State: From the Standpoint of Political Geography. *Geographical Review* **21**, 529–538 (1931). URL <http://www.jstor.org/stable/209364>.

- [6] Cremer, H., De Kerchove, A.-M. & Thisse, J.-F. An economic theory of public facilities in space. *Mathematical Social Sciences* **9**, 249–262 (1985). URL <http://www.sciencedirect.com/science/article/pii/0165489685900599>.
- [7] Alesina, A. & Spolaore, E. On the Number and Size of Nations. *The Quarterly Journal of Economics* **112**, 1027–1056 (1997). URL <http://www.jstor.org/stable/2951265>.
- [8] Drèze, J., Le Breton, M., Savvateev, A. & Weber, S. “Almost” subsidy-free spatial pricing in a multi-dimensional setting. *Journal of Economic Theory* **143**, 275–291 (2008). URL <http://www.sciencedirect.com/science/article/pii/S0022053108000343>.
- [9] Desmet, K., Breton, M. L., Ortuño-Ortín, I. & Weber, S. The stability and breakup of nations: a quantitative analysis. *Journal of Economic Growth* **16**, 183–213 (2011). URL <http://link.springer.com/article/10.1007/s10887-011-9068-z>.
- [10] Okabe, A., Boots, B., Sugihara, K. & Chiu, S. N. *Spatial Tessellations: Concepts and Applications of Voronoi Diagrams* (Wiley, Chichester ; New York, 2000), 2 edition edn.
- [11] Stephan, G. E. Territorial Division: The Least-Time Constraint Behind the Formation of Subnational Boundaries. *Science* **196**, 523–524 (1977). URL <http://www.sciencemag.org/content/196/4289/523>.
- [12] Weese, E. Political mergers as coalition formation: An analysis of the Heisei municipal amalgamations. *Quantitative Economics* **6**, 257–307 (2015). URL <http://dx.doi.org/10.3982/QE442>.
- [13] Greenberg, J. H. The Measurement of Linguistic Diversity. *Language* **32**, 109–115 (1956). URL <http://www.jstor.org/stable/410659>.
- [14] Lieberman, S. An Extension of Greenberg’s Linguistic Diversity Measures. *Language* **40**, 526–531 (1964). URL <http://www.jstor.org/stable/411935>.
- [15] Brasington, D. M. Joint provision of public goods: the consolidation of school districts. *Journal of Public Economics* **73**, 373–393 (1999). URL <http://www.sciencedirect.com/science/article/B6V76-40V4THC-4/1/80cc8e6f7df6783b202b1>
- [16] Gordon, N. & Knight, B. A spatial merger estimator with an application to school district consolidation. *Journal of Public Economics* **93**, 752–765 (2009). URL <http://www.sciencedirect.com/science/article/B6V76-4VR9FFC-1/2/20468f81e0dac602acd09>
- [17] Roth, A. E. & Vate, J. H. V. Random Paths to Stability in Two-Sided Matching. *Econometrica* **58**, 1475–1480 (1990). URL <http://www.jstor.org/stable/2938326>.
- [18] Gale, D. & Shapley, L. S. College Admissions and the Stability of Marriage. *The American Mathematical Monthly* **69**, 9–15 (1962). URL <http://www.jstor.org/stable/2312726>.

- [19] Bogomolnaia, A. & Jackson, M. O. The Stability of Hedonic Coalition Structures. *Games and Economic Behavior* **38**, 201–230 (2002). URL <http://www.sciencedirect.com/science/article/B6FW-458NBXX-1/2/12b3e799527594e128bb7>
- [20] Arkin, E. M. *et al.* Geometric stable roommates. *Information Processing Letters* **109**, 219–224 (2009). URL <http://www.sciencedirect.com/science/article/pii/S0020019008003098>.
- [21] Weese, E., Hayashi, M. & Nishikawa, M. Inefficiency and Self-Determination: Simulation-Based Evidence from Meiji Japan. SSRN Scholarly Paper ID 2651855, Social Science Research Network, Rochester, NY (2015). URL <http://papers.ssrn.com/abstract=2651855>.
- [22] Ramankutty, N., Foley, J. A., Norman, J. & McSweeney, K. The global distribution of cultivable lands: current patterns and sensitivity to possible climate change. *Global Ecology and Biogeography* **11**, 377–392 (2002). URL <http://onlinelibrary.wiley.com/doi/10.1046/j.1466-822x.2002.00294.x/abstract>.
- [23] Lalonde, S. *Determining boundaries in a conflicted world the role of uti possidetis* (McGill-Queen’s University Press, Montréal, Que., 2002).
- [24] Alesina, A., Baqir, R. & Hoxby, C. Political Jurisdictions in Heterogeneous Communities. *The Journal of Political Economy* **112**, 348–396 (2004). URL <http://www.jstor.org/stable/3555176>.
- [25] Michalopoulos, S. The Origins of Ethnolinguistic Diversity. *American Economic Review* **102**, 1508–39 (2012). URL <https://www.aeaweb.org/articles.php?doi=10.1257/aer.102.4.1508>.
- [26] Weber, E. J. *Peasants into Frenchmen: The Modernization of Rural France, 1870-1914* (Stanford University Press, Stanford, Calif, 1976).
- [27] Lodge, R. A. *French: From Dialect to Standard* (Routledge, London, 1993).
- [28] Kuzio, T. Borders, symbolism and nation-state building: Ukraine and Russia. *Geopolitics and International Boundaries* **2**, 36–56 (1997). URL <http://dx.doi.org/10.1080/13629379708407589>.
- [29] Scheffe, H. *The analysis of variance*. A Wiley publication in mathematical statistics (Wiley, New York, 1959).
- [30] Bixby, R. E. A Brief History of Linear and Mixed-Integer Programming Computation. *Documenta Mathematica Extra Volume ISMP*, 107–121 (2012). URL <http://www.math.uiuc.edu/documenta/>.
- [31] Jetz, W. & Rahbek, C. Geometric constraints explain much of the species richness pattern in African birds. *Proceedings of the National Academy of Sciences* **98**, 5661–5666 (2001). URL <http://www.pnas.org/content/98/10/5661>.

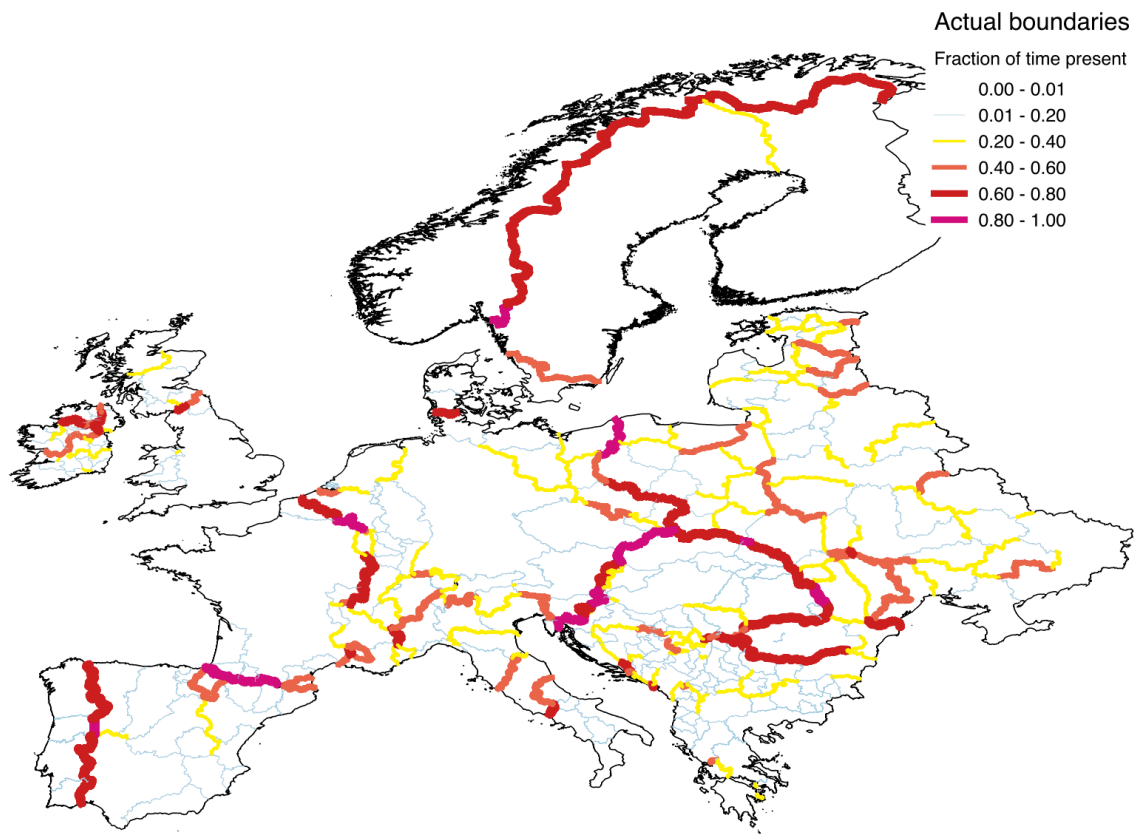


Figure 1: Actual Boundaries (see Methods section for note regarding Portugal)

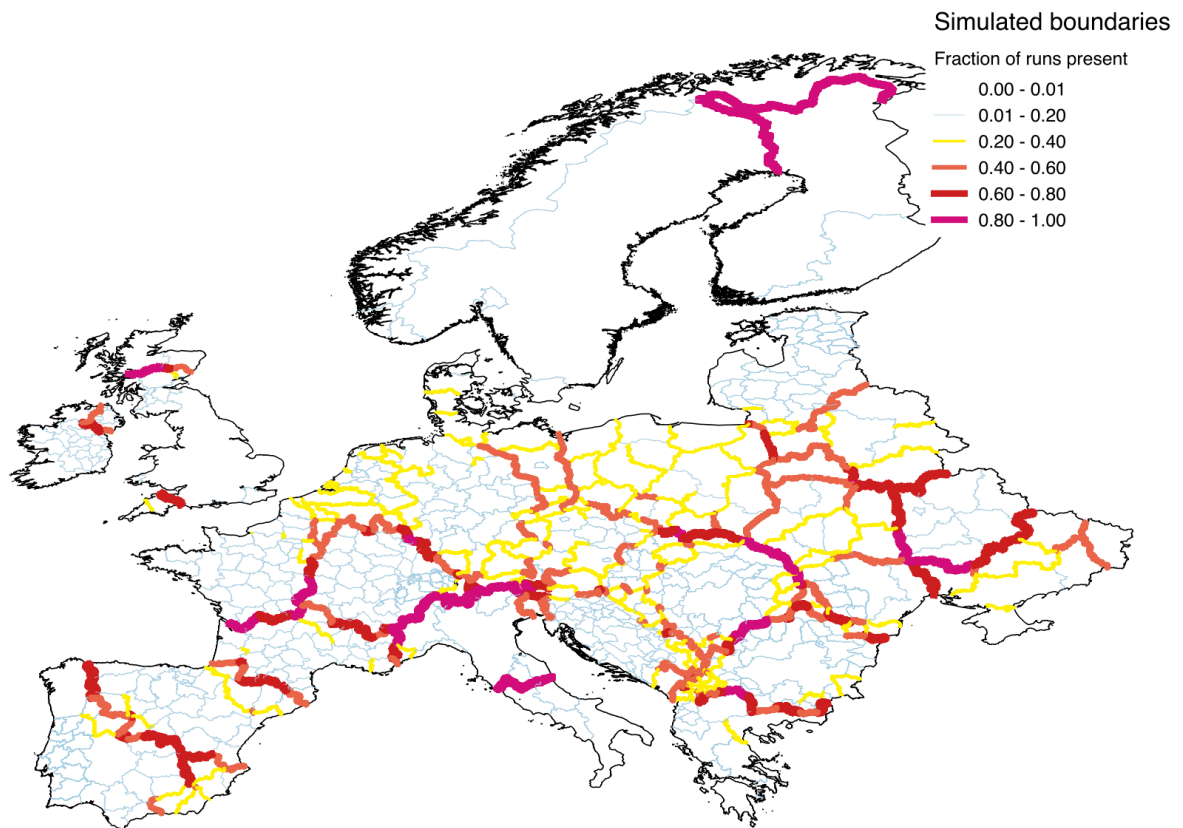


Figure 2: Simulated Boundaries

## A Methods

The method used to find the stable partitions of countries is taken from [21], which includes a description of the mixed integer programming technique used to find possible new countries  $S'$ . A discussion of the speed gains in mixed integer programming can be found in [30].

According to the theory presented, each person in Europe could be considered as a separate  $i$ , with a specific individual location on the continent. This direct approach is computationally infeasible. Instead, we aggregate the population of Europe into small geographic units, and then use  $i$  to represent these units. All simulations are then run based on these units.

The algorithm used to generate stable partitions is computationally stable when considering units of self-determination with populations that are non-trivial relative to the populations of the countries in stable partitions. If the grid square data on agricultural suitability shown in Supplemental Figure 3 were used directly, however, a substantial number of grid squares would have extremely small populations. If, during a simulation, these grid squares were to end up by themselves (i.e. in a country consisting only of that grid square), then the per capita cost faced by that grid square would potentially be many orders of magnitude higher than for any of the other grid squares. In mixed integer programming, large differences in magnitude of this sort generally lead to difficulties in obtaining accurate solutions, and thus need to be avoided. To avoid this situation, we aggregate grid squares into larger units, and use those units as the basic units for the simulation performed.

We begin by using the administrative divisions provided by the Global Administrative Boundaries dataset. We omit the European microstates (Andorra, Monaco, San Marino, Vatican City, and Liechtenstein) and Iceland. We use national boundaries (level 0 subdivisions) for four particularly small countries: Montenegro, Moldova, Macedonia, and Luxembourg. We use level 2 subdivisions for five countries: the United Kingdom, Germany, France, Spain, and Bosnia-Herzegovina. For all other European countries, we use level 1 administrative subdivisions. We eliminate Mediterranean islands, and all other holdings south of continen-

tal Europe (e.g. Ceuta and Melilla). The resulting set of administrative polygons forms the basis for our aggregations of agricultural suitability.

We continue to use the grid square as a unit of measure for agricultural suitability. Thus, a polygon that includes all of exactly one grid square where the grid square is rated as 10% suitable, will be counted as having a total suitability of 0.1. Similarly, a polygon that contained 50% of a grid square that was 100% suitable for agriculture, and all of a grid square that was 10% suitable will be counted as having a total suitability of  $0.5 \times 1.0 + 1.0 \times 0.1 = 0.6$ .

We then calculate distance between each of the polygons using the Dijkstra shortest-path algorithm. An alternative would be to simply calculate the straight-line distance between polygons. This straight-line distance, however, would not necessarily remain on land. On the one hand, historically travel was often faster by sea. However, most of this seafaring occurred within a short distance of the coast. A shortest-path distance does not allow straight-line voyages across, for example, the Ligurian Sea between Italy and France. Travel distance along the coast by sea, however, will be approximated very closely by the travel distance along the coast by land, and this will be calculated by the shortest-path algorithm.<sup>6</sup>

To calculate the shortest-path distance, begin with the set of polygons. Let these units form the vertices in a graph. Edges of the graph indicate geographic adjacency. The locations of the vertices are determined by the weighted mean location of the grid squares contained in the polygons in question.<sup>7</sup> The distance that will be used for any two units in the simulation is the shortest path distance between the corresponding vertices in the graph.

We then use an iterative process to further aggregate polygons that still have particularly low totals for agricultural suitability. We use as a threshold 0.1 units of agricultural suitability, using the grid square units of measure just described. We amalgamate the polygon with the lowest total amount of agriculturally suitable land with its closest neighbour. We then recalculate the distances using the shortest path algorithm just described, as the

---

<sup>6</sup>The major assumption made is assigning the same cost to these routes as “inland” routes. A more sophisticated model might allow for longer voyages at the same cost, so long as they were along the coast.

<sup>7</sup>This graph is not shown, but it is quite similar to the final graph shown in Supplemental Figure 10.

number of vertices in the graph has decreased by one. Repeat this process until the player with the lowest total amount of agriculturally suitable land has at least 0.1 units of agricultural suitability. The resulting set of polygons is shown in Supplemental Figure 4. This algorithm does not pay any attention to actual national boundaries, and thus the final polygons used for the simulation units may cross current national boundaries. A notable case where this happens is northern Portugal, which is grouped with Galicia because of a lack of agriculturally suitable land in the northwest corner of the Iberian peninsula.

A final issue arises because, with the data actually considered, the presence of islands would lead to the graph consisting of several connected components, rather than only one. While this would not affect the simulations, some actual countries, most notably the United Kingdom, include territory separated by water. We thus join some vertices across water, at the narrowest point of the water in question. Four of edges of this sort are added: they are Great Britain - (Northern) Ireland, Great Britain - France, Denmark - Sweden, and Finland - Estonia. The resulting graph is shown in Supplemental Figure 10.

Although the edges chosen are the shortest links for the general areas in question, the choice to make these four links is based on a heuristic examination of the graph in question, rather than some specific formal rule. We thus cannot draw any firm conclusions on the nature of water as a natural border from the simulations performed. The four “water edges” are not shown in Figures 1 and 2, but are shown in Supplemental Figures 11 and 12, which are an alternative visualization of the same data. In the actual data all four of the edges are often boundaries between countries. This is also true in the simulation results for three of the four edges, with the exception being the Great Britain - (Northern) Ireland edge. In the simulations, Northern Ireland is usually in the same country as England.

## A.1 Boundary Statistics

To determine whether a country boundary lies between two units, we use the locations shown in Supplemental Figure 10, which are weighted mean locations. If the vertices on either

side of an edge are in the same country, then we code this edge as having no boundary. Conversely, if these vertices are in different countries, then we code the edge in question as having a boundary. The exact location of the boundaries displayed in Figures 1 and 2 correspond to boundaries of the polygons in Supplemental Figure 4, coloured appropriately. The boundaries in Figure 1 thus do not match precisely with actual international borders: this is most obvious in the case of Spain and Portugal, where the lack of agriculturally suitable land in Galicia has resulted in a particularly large polygon, the weighted mean of which lies in Portugal. The entire western coast of the Iberian Peninsula is thus reported in Figure 1 as belonging to Portugal. This method is used because it allows easy comparisons between the actual boundaries reported in 1, and the simulated boundaries reported in Figure 2.

## A.2 Distance squared

Suppose that instead of using geographic distance for  $\ell$  in Equation 3, we used geographic distance squared. The last term of Equation 4 then becomes  $\gamma P_S^{-1} \sum_{i' \in S} d_{i,i'}^2$ , where  $d_{i,i'}$  is geographic distance squared. The average of this over all individuals in  $S$  is  $\gamma P_S^{-2} \sum_{i \in S} \sum_{i' \in S} d_{i,i'}^2$ . A standard analysis of variance result[29] is that this is equal to  $\gamma P_S^{-1} \sum_{i \in S} d_{i,s}^2$ , where  $s$  is the geographic location of the population-weighted centroid of  $S$ .

The sum  $\sum_{i \in S} U_i$ , where  $U_i$  is as in Equation 4, except using  $\ell = d^2$ , would thus describe the payoff to a despot located at  $s$  who collects a tax of  $-V$  from everyone in the country,<sup>8</sup> but who must pay a travel cost  $\gamma d_{i,s}^2$  to collect the tax from individual  $i$ , and in addition must pay a fixed cost  $F$  to run the country.

## B Statistical Analysis

A qualitative comparison of Figures 1 and 2 in the main text suggests that the model proposed in this paper successfully captures some aspects of the actual process of boundary

---

<sup>8</sup>Here  $V$  should be negative, whereas in the democratic case the explanation for the model suggested that  $V$  is a cost that would be positive.

formation. Two issues are best addressed via quantitative analysis, however: the possibility that these results are due to a mechanical effect due to differing surface areas of units used in the simulations, and the possibility that the results observed are no more surprising than what would be observed if random convex “countries” were selected so as to partition Europe.

The first of these issues arises because the simulation technique used requires that the units of self-determination considered contain approximately the same amount of agriculturally suitable land. This requires that some units be much larger than others in terms of surface area. If both actual and simulated boundaries were distributed randomly across space, this would result in a spurious correlation between the actual and simulated boundaries, because there would be more boundaries between units that were farther apart, both in the actual and simulated data. It is not possible to eliminate this bias, given the computational constraint that requires that there be no units in the simulation that have particularly small amounts of agriculturally suitable land.

To ensure that this bias is not responsible for the results displayed in Figures 1 and 2, we conduct a statistical analysis that includes as control variables functions of the distance between the units in question. Specifically, consider the regression

$$\text{BOUNDARY.ACTUAL}_{i,i'} = \beta_0 + \beta_1 \text{BOUNDARY.SIMULATED} + \beta_2 \text{DIST}_{i,i'} + \epsilon_{i,i'}$$

where  $\text{DIST}_{i,i'}$  is the geographic distance between  $i$  and  $i'$ , calculated as described in the Methods section, and  $\text{BOUNDARY.ACTUAL}_{i,i'}$  is the data from Figure 1: the fraction of the time there is a boundary between  $i$  and  $i'$ . Similarly,  $\text{BOUNDARY.SIMULATED}_{i,i'}$  is the data from Figure 2.

$\text{BOUNDARY.ACTUAL}$  is a limited dependent variable, as it is always between 0 and 1. We thus also consider a logistic regression. In addition, we consider similar regressions where  $\beta_3 \text{DIST}_{i,i'}^2$ , and potentially also a cubic term, are included.

Table 1 shows the results of this analysis. The estimate of the coefficient on `BOUNDARY.SIMULATED` is reduced somewhat by including the distance term. Additional higher order terms do not have any statistically significant effect on the coefficient estimate. In all columns it is clear that there is still a statistically significant relationship between the actual boundaries and the simulated boundaries, with a  $t$ -value of at least 8.

The  $t$ -values associated with the results in Table 1 might appear low given a qualitative comparison of Figures 1 and 2. A closer look at these figures, however, reveals that in several important cases the simulated results give a boundary that is very close, but not exactly the same, as the actual boundary. For example, most of the simulated boundary in the Pyrenees is slightly to the south of the actual boundary, while the simulated boundary between Germany and Poland is to the west of the (average historical) actual boundary. The results shown in Table 1 only compare simulated boundaries and actual boundaries for exactly the same  $i$  and  $i'$  pair. Thus, simulated boundaries that are qualitatively correct, but not in precisely the right place, are not accounted for in the coefficient estimates presented in the table.

One might wonder whether the statistics reported in Table 1 are themselves uninformative, because due to the particular geographic shape of Europe, certain countries will tend to emerge regardless of how agricultural land is distributed. In this case the proposed model is inappropriate, as it contributes nothing beyond an even simpler model such as one following [31]. To check this we maintain the same geographic base units for our simulation (as shown in Figure 4), but randomly shuffle total agricultural suitability among units. That is, a given unit  $i$  has a  $1/699$  probability of receiving the suitability of each of the other 699 units. We generate 100 such randomly shuffled versions of Europe, and compute a simulated partition for each of these. We then rerun the analysis shown in Table 1 for these randomly shuffled Europes. The results are reported in Table 2. After adding polynomial controls for the distance between units, the simulation results have no statistically significant relationship with the actual boundaries of Europe. This suggests that the statistically significant results

reported in Table 1 are not merely due to the “shape” of the units involved in the coalition formation game considered, as for example shown in Figure 10.

A remaining possibility is that the statistically significant results reported in Table 1 are merely the result of the fact that real countries generally consist of geographically contiguous territory, and in many cases have a roughly convex shape. The model presented will tend to produce simulated countries that are geographically contiguous and roughly convex, but if this alone is responsible for the results in Table 1, then a much simpler model could be proposed that has the same explanatory power.

To test this possibility, we construct random countries based on Voronoi cells. We choose 23 units as “generating points”, and make 23 countries by assigning each of the 699 units shown in Figure 10 to the closest of these generating points. The resulting Voronoi cells will be convex, and because distances were calculated via the shortest path algorithm, the resulting countries will be geographically contiguous. We generate 1000 of these “Voronoi partitions” of Europe, starting each time with a new random set of 23 generating points. We then redo the analysis in Table 1, using these Voronoi partitions. The results of this analysis are shown in Table 3. Again, we find that after controlling for a polynomial of geographic distance, there is no statistically significant relationship between the simulated boundaries and the actual boundaries of Europe. This suggests that the results of Table 1 are not due to basic geometric properties of actual countries.

Table 1: Regression Results

<i>Dependent variable: BOUDARY.ACTUAL (fraction of time there is a boundary in the Euratlas data)</i>								
	<i>OLS</i>				<i>Logistic</i>			
	I	II	III	IV	V	VI	VII	VIII
BOUNDARY . SIMULATED	0.294*** (0.022)	0.212*** (0.024)	0.199*** (0.024)	0.199*** (0.024)	1.810*** (0.074)	1.318*** (0.081)	1.163*** (0.084)	1.150*** (0.084)
DIST		-0.001*** (0.0001)	-0.001*** (0.0002)	-0.001*** (0.0004)		-0.005*** (0.0004)	-0.010*** (0.001)	-0.015*** (0.002)
DIST^2			0.00000** (0.00000)	0.00000 (0.00000)			0.00001*** (0.00000)	0.00004*** (0.00001)
DIST^3				0.000 (0.000)				-0.00000*** (0.000)
constant	0.595*** (0.019)	0.738*** (0.025)	0.774*** (0.029)	0.773*** (0.032)	0.202*** (0.060)	1.088*** (0.086)	1.531*** (0.106)	1.759*** (0.130)
Observations	1,837	1,837	1,837	1,837	1,837	1,837	1,837	1,837
R <sup>2</sup>	0.090	0.123	0.126	0.126				
Adjusted R <sup>2</sup>	0.089	0.122	0.124	0.124				

Note:

\*p&lt;0.1; \*\*p&lt;0.05; \*\*\*p&lt;0.01

Table 2: Regression Results (shuffled agricultural suitability)

<i>Dependent variable: BOUDARY.ACTUAL (fraction of time there is a boundary in the Euratlas data)</i>								
	<i>OLS</i>				<i>Logistic</i>			
	I	II	III	IV	V	VI	VII	VIII
BOUNDARY . SIMULATED	0.260*** (0.033)	0.048 (0.039)	0.008 (0.040)	0.008 (0.040)	1.751*** (0.117)	0.464*** (0.141)	0.072 (0.144)	0.032 (0.144)
DIST		-0.001*** (0.0001)	-0.002*** (0.0002)	-0.002*** (0.0004)		-0.007*** (0.0004)	-0.014*** (0.001)	-0.020*** (0.002)
DIST^2			0.00000*** (0.00000)	0.00000 (0.00000)			0.00002*** (0.00000)	0.0001*** (0.00001)
DIST^3				0.000 (0.000)				-0.00000*** (0.000)
constant	0.624*** (0.028)	0.901*** (0.038)	0.978*** (0.043)	0.976*** (0.046)	0.226** (0.095)	1.918*** (0.139)	2.711*** (0.158)	3.001*** (0.177)
Observations	1,837	1,837	1,837	1,837	1,837	1,837	1,837	1,837
R <sup>2</sup>	0.032	0.085	0.093	0.093				
Adjusted R <sup>2</sup>	0.031	0.084	0.092	0.091				

Note:

\*p<0.1; \*\*p<0.05; \*\*\*p<0.01

Table 3: Regression Results (randomly generated Voronoi cells)

<i>Dependent variable: BOUDARY.ACTUAL (fraction of time there is a boundary in the Euratlas data)</i>								
	<i>OLS</i>				<i>Logistic</i>			
	I	II	III	IV	V	VI	VII	VIII
BOUNDARY . SIMULATED	0.500*** (0.053)	0.137** (0.066)	0.048 (0.070)	0.052 (0.072)	3.485*** (0.195)	1.310*** (0.240)	0.512** (0.250)	0.314 (0.256)
DIST		-0.001*** (0.0001)	-0.002*** (0.0002)	-0.002*** (0.0005)		-0.006*** (0.0004)	-0.013*** (0.001)	-0.019*** (0.002)
DIST^2			0.00000*** (0.00000)	0.00000 (0.00000)			0.00002*** (0.00000)	0.00005*** (0.00001)
DIST^3				0.000 (0.000)				-0.00000*** (0.000)
constant	0.431*** (0.044)	0.822*** (0.061)	0.940*** (0.068)	0.932*** (0.075)	-1.159*** (0.156)	1.183*** (0.218)	2.298*** (0.244)	2.717*** (0.274)
Observations	1,837	1,837	1,837	1,837	1,837	1,837	1,837	1,837
R <sup>2</sup>	0.046	0.086	0.093	0.093				
Adjusted R <sup>2</sup>	0.045	0.085	0.092	0.091				

Note:

\*p<0.1; \*\*p<0.05; \*\*\*p<0.01

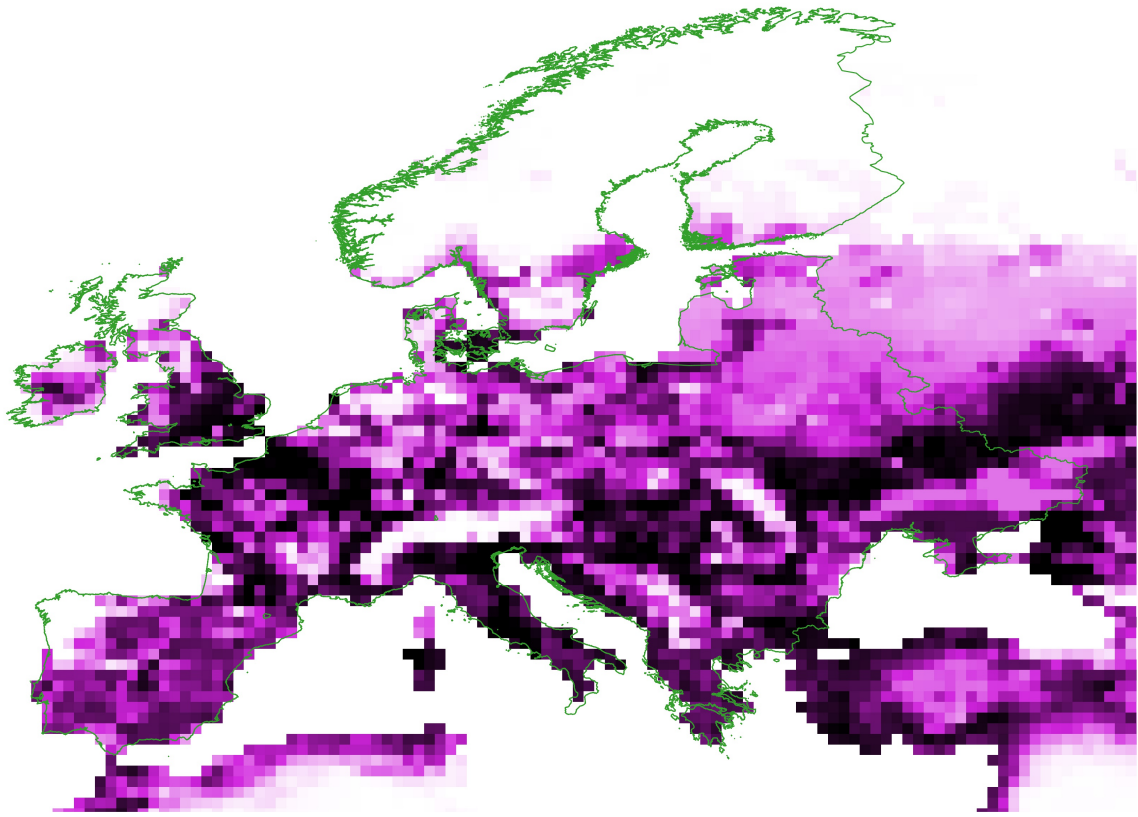


Figure 3: Agricultural Suitability

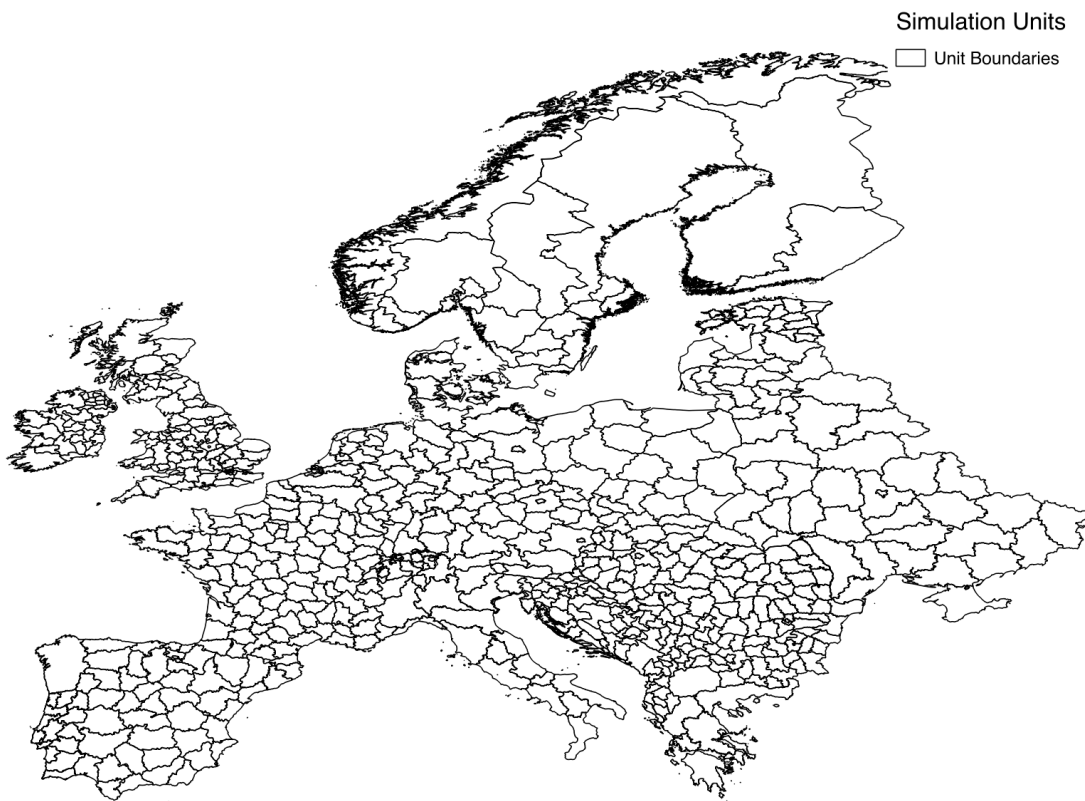


Figure 4: Units of Self-Determination for Simulation

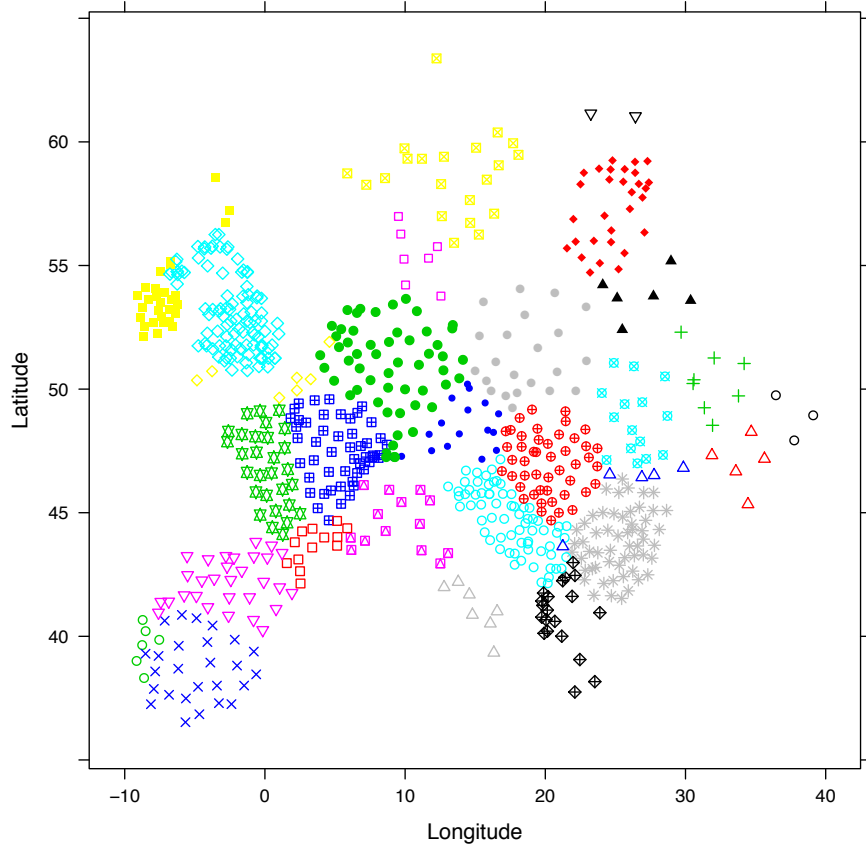


Figure 5: A Core Partition (coloured symbol indicates country)

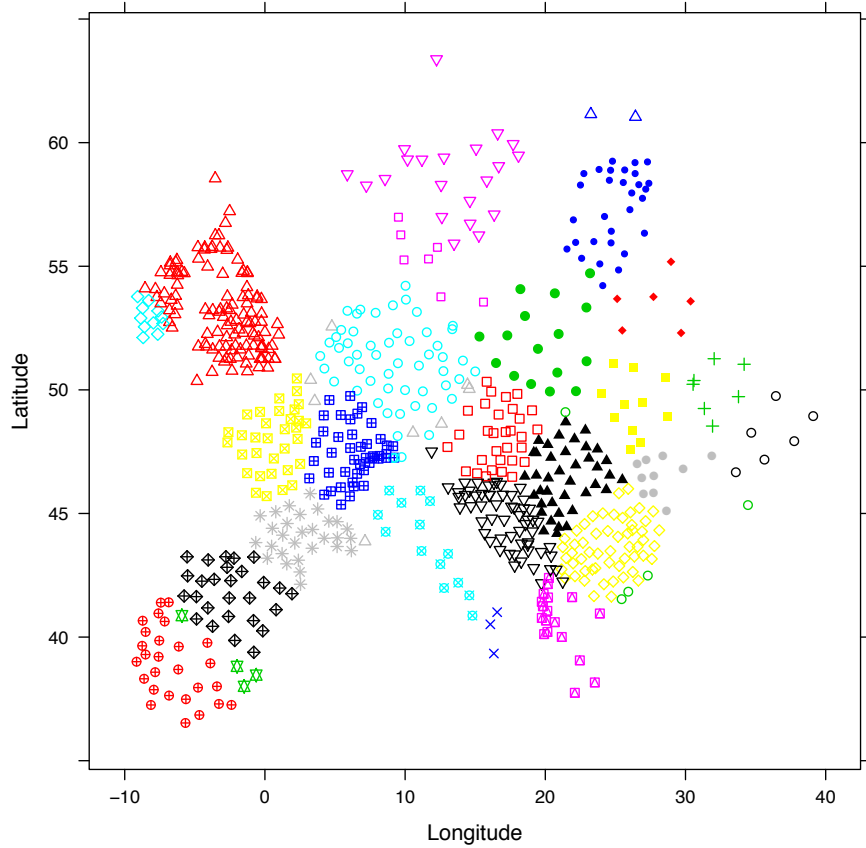


Figure 6: Another Core Partition (symbols have no relationship to those in previous figure)

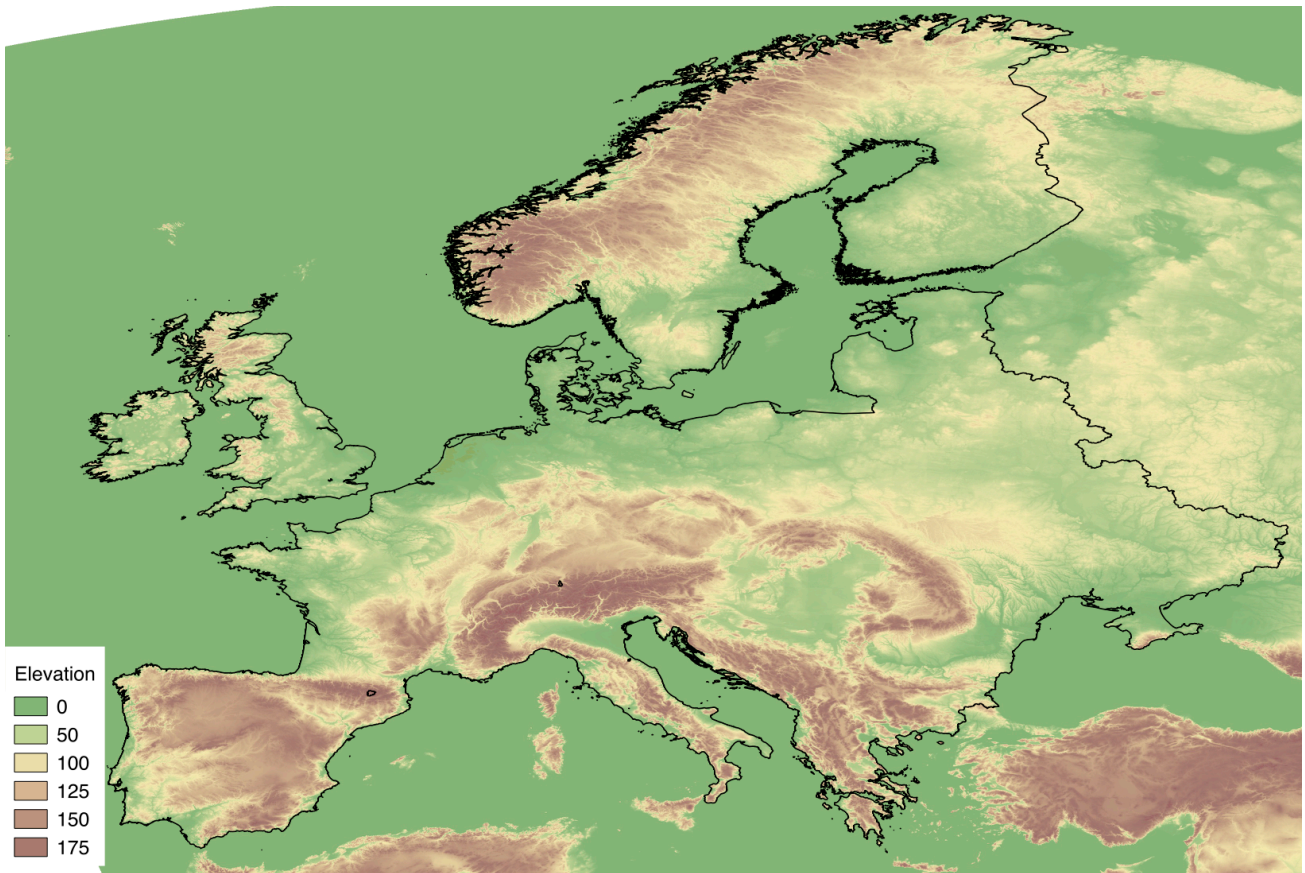


Figure 7: Elevation (GTOPO30 data, processed by European Environmental Agency)

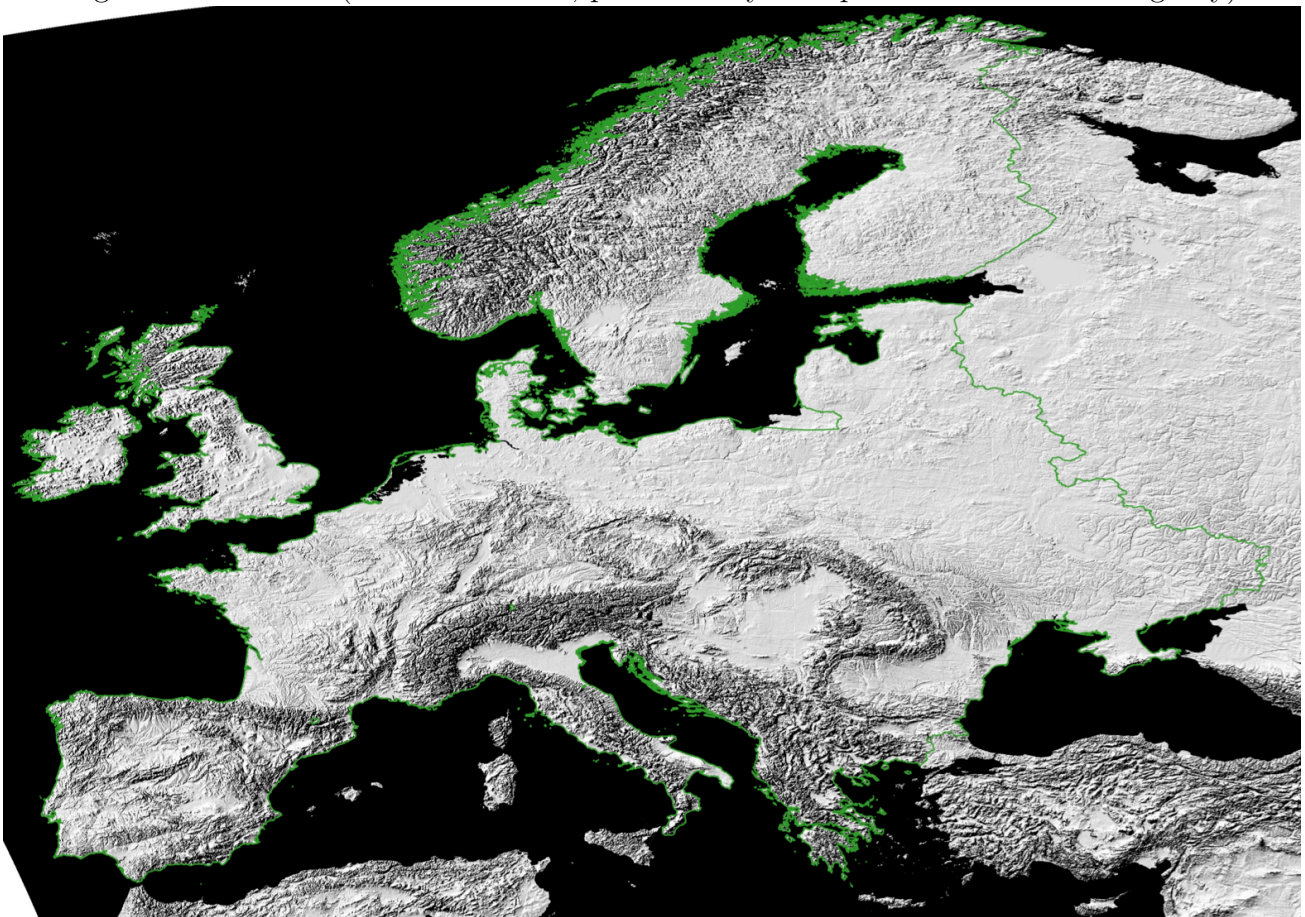


Figure 8: Hillshade (GTOPO30 data, processed by European Environmental Agency)

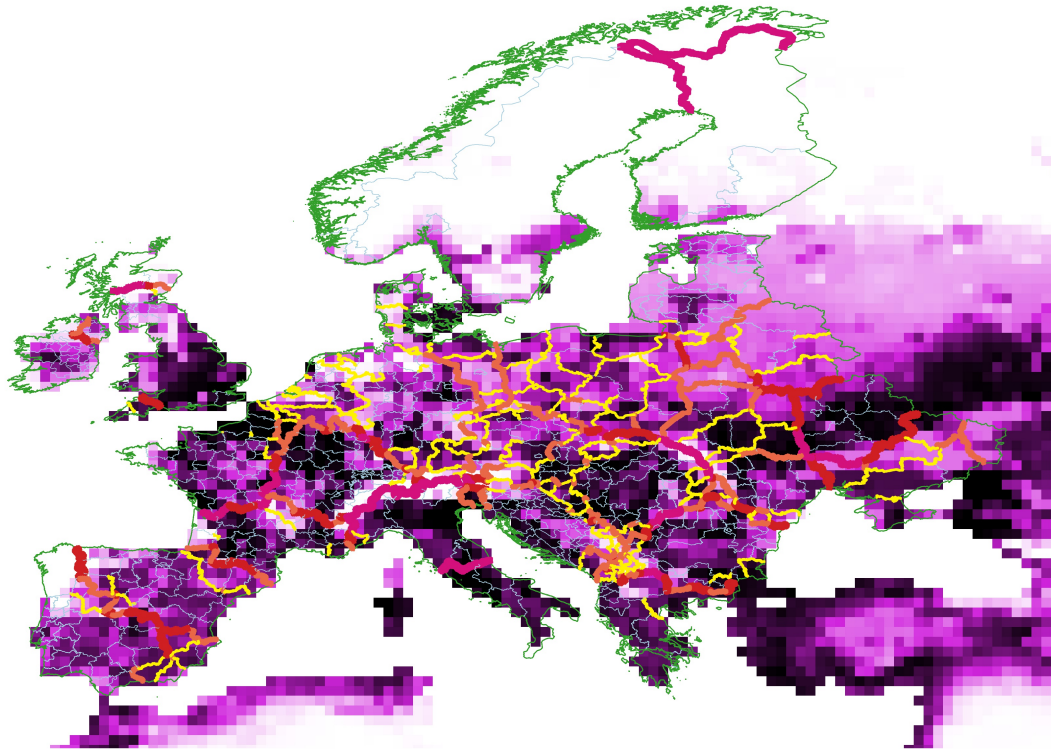


Figure 9: Agricultural Suitability and Simulated Boundaries



Figure 10: Units (vertices) and Geographic Adjacency (edges)

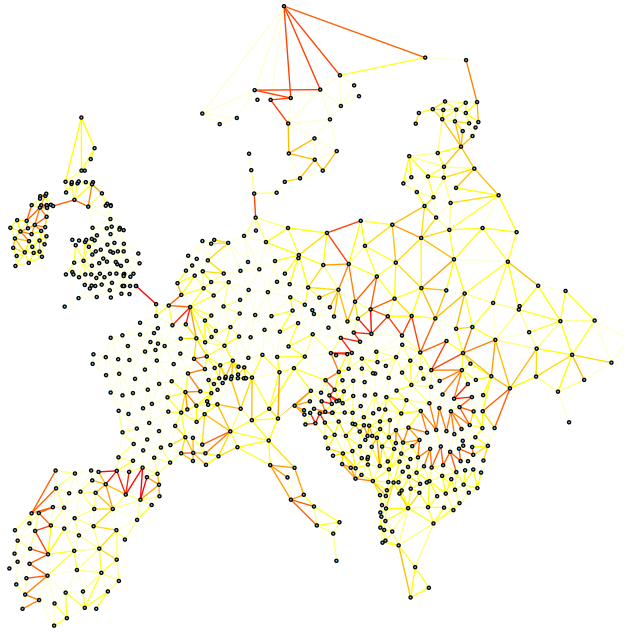


Figure 11: Actual Boundaries (different representation of same data as Figure 1, plus “water” edges)

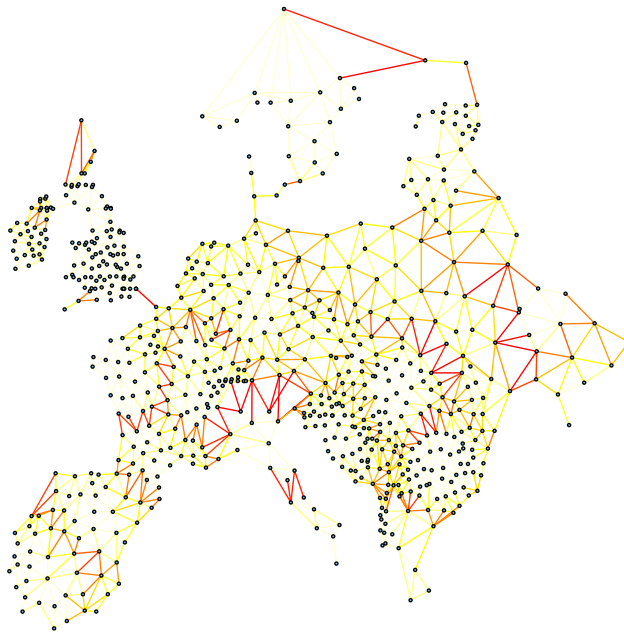


Figure 12: Simulated Boundaries (different representation of same data as Figure 2, plus “water” edges)

Legend: a dark red edge is one where the incident vertices are always in different countries; a white edge is one where they are always in the same country. An uncoloured version of the graph is shown in Supplemental Figure 10.