

Statistical detection of spatial plant patterns under the effect of forest use

I. López^{a*}, T. Standovár^b, J. Garay^b, Z. Varga^c and M. Gámez^a

^a *Department of Statistics and Applied Mathematics, University of Almería. La Cañada de San Urbano 04120 Almería, Spain (milopez@ual.es , mgamez@ual.es)*

^b *Department of Plant Taxonomy and Ecology, L. Eotvos University, H-1117 Budapest, Pazmany Peter setany 1/c. , Hungary (garayj@ludens.elte.hu , standy@ludens.elte.hu)*

^c *Institute of Mathematics and Informatics, Szent István University, Páter K. u. 1., H-2103 Godollo, Hungary. (Varga.Zoltan@gek.szie.hu)*

Abstract: The analysis of the consequences of land use (in particular forest use) may be considered as a partial step towards an integrated modelling of a land system. In the paper a forest territory is considered, where a gap-cut is made, and after a given time period the eventual change in the spatial distribution of undergrowth plants and tree seedlings is to be detected. Floristic data are collected along a line transect. A method for the detection of the change in the plant distributions along the transect is proposed to see whether this occurs at the geometric frontier of the human intervention.

* *Corresponding author, phone: +34 950 015775, fax: +34 950 015167*

Since in the considered case the distribution of the change-point estimate is not known, as a substitute of its confidence interval, the so-called change-interval is calculated, using an adaptation of the bootstrap method. As an illustration, for a concrete plant species, the maximum likelihood estimation of the change-point and the calculation of the above mentioned change-interval is presented. Finally, the validation of the proposed method against some typical ecological situations is also presented, which provides a justification of the used algorithms.

Keywords: forest use, forest gap, plant patches, edge detection, change-point, change-interval, bootstrap.

1. INTRODUCTION

The analysis of the consequences of land use (in particular forest use) may be considered as a partial step towards an integrated modelling of a land system. Let us consider a forest territory, where a gap-cut is made, and after a given time period the eventual change in the spatial distribution of undergrowth plants and tree seedlings is to be detected (see [15] and [9]). If floristic data are collected along a line transect, we can try to detect the change in the plant distributions along the transect, the so-called *change-point*, and see whether this occurs at the geometric frontier of human intervention.

The problem, at a theoretical level, can be addressed using the methodology of *change-point analysis* which is a technically involved branch of mathematical statistics (see e.g.

[2], [4]), widely used to explore the possible temporal or spatial structure of local homogeneity from collected data. The main application fields of change-point analysis include meteorology, hydrology, or environmental studies, economy, quality control in industry, biology and medicine. In this paper we propose a practical, operative approach, using only technique of classical statistics.

One approach of treating vegetation patchiness is classification, i.e. distinguishing discrete entities based on resemblance. Early attempts classified vegetation based on similarities in physiognomy (the gross appearance) caused by the relative importance of different growth forms ([10]). Then classical phytosociology put the emphasis on resemblance of species composition. Papers in [14] illustrate how different schools of phytosociology developed their criteria by which the units of classification (associations) were recognised. From the mid 20th century methods of numerical classification have been used more and more widely (e.g. [16]).

Another approach is to study spatial patterns of individual plant species. Spatial pattern means the arrangement of plants or of patches of plants in space with certain amount of predictability ([5]).

In ecology the change-point problem is also known as problem of “boundary or edge detection”, see [3], [8] and [11]. In these papers further references to a large variety of applications of edge detection can be also found. The methodology of the change point has already been applied with success to the analysis of plant patterns, see [12].

In our case, given a plant species, along a line transect quadrats have been located and in each quadrat the individuals have been counted. We consider these data as samples of

two distributions of the same type but with different parameters, separated by a change-point K . Based on the maximum likelihood approach, an algorithm is given to estimate K .

Since the distribution of the change-point estimate is not known, as a substitute of its confidence interval, the so-called *change-interval* will be calculated, using an adaptation of the *bootstrap* method. For this widely applied simulation method see e.g. [6], a justification of the use of bootstrap in this case can be found in [7]. The implementation of the above algorithms was realized with the application of the statistical software “R”. As an illustration, for a concrete plant species, the maximum likelihood estimation of the change-point and the calculation of the above mentioned change-interval will be presented. Finally, for a justification of the proposed method, our algorithms are also tested against some typical ecological situations.

The paper is organized as follows. Section 2 recalls some general aspects of statistical analysis of human effects on a forest. Section 3 describes the experimental framework of our study. Section 4 is devoted to the set-up of the basic model, and also contains the description of the used algorithms and the obtained results. In Section 5 the proposed methodology is validated against some typical ecological situations. Finally, in Section 6 some conclusions are drawn.

2. STATISTICAL ANALYSIS OF HUMAN EFFECTS ON A FOREST

In forests, human management is aimed at only a few – though dominant – components (trees, wild game species) of the whole ecosystem. However, nowadays increasing attention is paid to the loss of biodiversity. As a consequence there is a need to assess the effects of different land use activities – in our case forest management – on original biodiversity. Experience shows that not all plant species show a clear, easily detectable reaction to management activities, even if they create steep gradients (like the opening of a small canopy gap in an old forest). Often one can only find difference in the distribution of a species among patches of different quality. Consequently, new methods – capable of detecting such minor changes – could be used to detect minor, not readily detectable causes of human management.

There is a very wide range of applications of statistical methods of change-point analysis in ecology (see [1], [18], [17] and their references). Our method not only provides an estimation for the location of a change-point and the different distributions laying in different patches, but also provides a so-called *change-interval* (C.I. for brevity) which localizes the distribution change with a high probability level. This interval can be considered as an estimate of a transient zone between patches. The latter has particular importance in plant ecology since the change between patches usually is not point-like. In a transient zone there may be a mix of two patches, or a special plant composition. Our aim is to study case when the transient zone is small and contains only a mix of the patches.

3. THE FIELD EXPERIMENT

To test the applicability of our method, we used data that were collected in the framework of a study aimed at investigating the effects of canopy gap size on the resulting spatial distributions of key abiotic environmental variables (light and soil moisture) in gaps, and at studying how light and soil moisture affect the abundance and distribution of herb layer species. The study site is located in the Börzsöny Mountains, northern Hungary (47.9° N, 18.9° E). Mean annual temperature is 8 °C, mean monthly temperature is -3.5 °C and 18 °C in January and in July, respectively. Annual precipitation is 700–800 mm. Bedrock is andesite, on which medium deep brown forest soil has developed. The study area is located at 540–610 m elevation, on a relatively steep east-northeast facing slope, that is covered by an almost pure stand of European beech (*Fagus sylvatica* L.). Average tree height is 25 m, mean diameter at breast height is 30 cm. Detailed site description is given in [9].

The selected stand was a good representation of even-aged, mature (86 year old), and dense forests – typical products of the common contemporary silvicultural system (uniform shelter-wood), [13]. Because of the dense tree canopy, understorey vegetation was extremely sparse before opening the experimental gaps. Three large gaps (the proportion of tree height of surrounding stand (H) to gap diameter (D) was 1:1.5) and five small gaps (H:D was 1:0.5) were created in February 2001 (Figures 1 and 2). We used a systematic sampling design with 5-meter grid resolution and 1x1 m quadrats. Each large gap contained 123 quadrats, whereas small gaps contained 64 quadrats each. Vegetation data were recorded on seven occasions (in September 2000 – before gap creation – May/ September 2001, and May/September 2002, August 2004 and 2006). On each occasion we determined the cover of each herbaceous species using visual estimation in each quadrat. Among other environmental variables, light conditions were

studied in each quadrat, so we could reliably decide if a quadrat was in a gap or non-gap environment.

For the present study we used the data of one species, bramble (*Rubus fruticosus* L.) collected in one of the large gaps in 2006, containing 25 quadrats, i.e. in the seventh growing season after the artificial gaps had been opened.

4. MODEL DESCRIPTION, ALGORITHMS AND RESULTS

4.1. Model Description

We fix $0 < K < N$ and suppose that the number of plants in quadrats 1, 2, 3, ..., K are independent random variables with the same discrete probability distribution

$$\xi: \begin{cases} 0 & 1 & 2 & \dots & r \\ p_0 & p_1 & p_2 & \dots & p_r \end{cases}$$

whereas the number of plants in quadrats $K+1$, $K+2$, $K+3$, ..., N are independent random variables with the same discrete probability distribution

$$\eta: \begin{cases} 0 & 1 & 2 & \dots & l \\ q_0 & q_1 & q_2 & \dots & q_l \end{cases}$$

1	2	...	$K-1$	K	$K+1$	$K+2$...	N
ξ	ξ	...	ξ	ξ	η	η	...	η

First, from a given sample vector $X:=(x_1, x_2, \dots, x_N)$, for each possible K , we estimate distributions of ξ and η , and the probability of “realization” of the given sample. Then, from the possible values of K we obtain the required estimate for K , applying the maximum likelihood approach.

4.2. Estimation of distributions ξ and η

For given $1 \leq K \leq N-1$, a possibility to estimate ξ in terms of relative frequencies may be the following: Let

$$r = \max_{j=1, \dots, K} x_j,$$

and for each $i = 0, 1, \dots, r$, we define the probability that the variable ξ takes each of its possible values:

$$\widehat{p}_{K,i} = P(\xi = i) = \frac{\text{number of indices } j = 1, 2, \dots, K \text{ with } x_j = i}{K},$$

providing an estimate for the distribution of ξ .

In analogous way we estimate the probability distribution of η : let

$$l = \max_{j=K+1, K+2, \dots, N-1} x_j,$$

and for each $i = 0, 1, \dots, l$, we define

$$\widehat{q}_{K,i} = P(\eta = i) = \frac{\text{number of indices } j = K+1, \dots, N \text{ with } x_j = i}{N-K}.$$

Let P_K be the probability of “realization” of the sample $X:=(x_1, x_2, \dots, x_N)$, calculated with the above estimated probabilities:

$$P_K := \left(\prod_{i=1}^K \widehat{p}_{K,x_i} \right) \left(\prod_{s=K+1}^N \widehat{q}_{K,x_s} \right),$$

considered the “goodness” of K . Based on the given sample X , our purpose is to find a K which maximizes P_K , providing the “best” (i.e. the “most probable”) value of K . We shall deal with this in the next section.

4.3. Algorithms

Algorithm 1 (Estimation of the change-point K):

1. Introduce sample X . $N := \text{Size}(X)$.
2. FOR $K=1$ until $N-1$:
 - a) Calculate: $\hat{p}_{K,i}$ and $\hat{q}_{K,i}$, for each i .
 - b) Calculate: $\text{Log } P_K = \sum_{i=1}^K \text{Log } \hat{p}_{K,x_i} + \sum_{s=K+1}^N \text{Log } \hat{q}_{K,x_s}$.

(Logarithm is introduced to avoid too small probability values.)
3. $\text{LogProbSample} := (\text{Log}P_1, \dots, \text{Log}P_{N-1})$.
4. $\text{Estimate}K :=$ Position K with maximum value among the coordinates of LogProbSample .
5. Return $\text{Estimate}K$.

To find a change-interval for K we elaborate a resampling method based on the known Bootstrap, but with certain modifications in the choice of the elements of each sample of the simulations, in order to fit the method to our problem. The original sample is divided in two homogenous parts, such that the order of the elements of the new samples is important, since, by the linear arrangement, we must not mix *all elements* in a random way. The generated samples must keep the particularity of having two homogenous parts. The process to follow is explained below:

Algorithm 2 (Calculation of a 90% level Change-Interval):

1. Introduce the sample $X:=(x_1,x_2,\dots,x_N)$. $N:= \text{Size}(X)$.

2. FOR $K=1$ until $N-1$

a) Calculate a weight for each K :

$$W_K = \left(\prod_{i=1}^K 10 \widehat{p}_{K,x_i} \right) \left(\prod_{s=K+1}^N 10 \widehat{q}_{K,x_s} \right)$$

b) Normalize the weights (and denote them by WN_K).

c) FOR $L=1$ until m (we generate m samples for each K):

c1) Generate K random numbers $\{u_1, \dots, u_K\}$ of a discrete uniform distribution $U[1, K]$, and $N-K$ random numbers $\{u_{K+1}, \dots, u_N\}$ of distribution $U[K+1, N]$.

c2) We generate each sample with two homogenous zones, selecting the elements of the original sample according to the random positions obtained in c1) for both zones:

$$X_L := (x_{u_1}, \dots, x_{u_K}, x_{u_{K+1}}, \dots, x_{u_N}).$$

c3) We apply Algorithm 1 to the sample X_L , to obtain an estimate K_L for the change-point.

d) From the obtained values K_1, \dots, K_m , (m large enough) we calculate a distribution d_K of the change-point, for each fixed K .

3. We combine all these new distributions to obtain a unique distribution $\sum_K WN_K d_K$ for K , for which we calculate the 90%-level change-interval, with percentile

5 and percentile 95 as extremes.

If we have a small amount of data, we can increase the number of data in the following way: we uniformly divide each quadrat of the linear transect into 100 small quadrats along a straight line. In each small quadrat the species will be present (value 1) or not (value 0). Let w be the number of small quadrats where the species is present. Then, using the statistical software “R” (version 2.7.2) we generate randomly w values of a discrete uniform distribution between 1 and 100. These w values will indicate the positions of the small quadrats with one plant, and the remaining $100-w$ will be the small quadrats with no plant. Therefore, in these 100 small quadrats we represent a w % presence of the species. We denote this new data vector by S . Now, however, we may have too many data for a reasonable run time for the calculations. To reduce them, we sum the values of each 10 consecutive quadrats, and denote this new data vector by Z . We carry out this in the following:

Algorithm 3:

1. Introduce the original sample $X = (x_1, x_2, \dots, x_N)$. $N := \text{Size}(X)$.
2. We increase its size to $100N$, then X changes to a vector $S = (s_1, s_2, \dots, s_{100 \cdot N})$ of 1s and 0s, and the frequency of 1s, uniformly placed in a random way between the positions $(i-1) \cdot 100 + 1, i \cdot 100$ in S , will be x_i .
3. We sum every 10 values of S obtaining a sample vector $Z = (z_1, z_2, \dots, z_{10N})$, to which we apply Algorithm 1 and Algorithm 2 (with $m = 100$), obtaining the estimate for K and the 90%-level change-interval.

4.4. The linea

Consider the species *Rubus fruticosus* with data of 2006, taken from the following area, and described by Cartesian coordinates “ X - Y ”, as shown Table 1.

In Table 2, in the first three columns we present the data for the species at each location with coordinates X and Y . We know that the distribution at the centre of this area is different from those observed at the extremes. The change of distribution is observed around $X=20$ on the left, and around $X=50$ on the right. Now, by symmetrically “folding” a diameter of the gap, we practically get a radius of the gap. Let us consider the data originally obtained for the 25 quadrats along a diameter of the gap, and redistribute them along the corresponding radius, as shown in the 5th and 6th columns of Table 2. The new data set is given in the last column. In this way, on the one hand, instead of finding *two change-points*, we will estimate a *single change-point* (for which our statistical method was proposed), on the other hand, by the folding, the number of quadrats along the new linea is virtually doubled. This approach can be justified by the geometric symmetry of the sampling arrangement, and by the homogeneity of the surrounding forest, as described in Section 3. Having estimated the single change-point (and the change-interval) for the “folded” arrangement, just by “unfolding” we will be able to estimate both change-points and their respective change-intervals, too.

Now, from the data of the last column, we want to detect the change of distribution. i.e. the X - coordinate of the change-point. To this end we apply Algorithm 3 to the last data column as “original sample”, obtaining a vector Z of the following 250 data:

```
[1] 000000000000000000000000000000000000000000000000000000000
[39] 000000000000000000000000000000000000000000000000000000000
[77] 000000000000000000000000000000000000000000000000000000000
[115] 0000000000000000000000000000000000000000000000000000000010
[153] 103000003123314224332322431214330011111
[191] 0000001112132362320300000000000063335543
[229] 3534333525255375444454
```

4.5. Results

With these 250 data, algorithms 3, 1 and 2 (with $m=100$) provide an estimated K equal to 150, and 90 % level change interval [149,160]. These results would correspond in the large data (S) to 1500 for K and [1490,1600] for the C.I., which in terms of the original quadrats (in the sense of the reordering given in Table 2) would be 15 for K , and [15,16] for the C.I., as shown in Figure 3. As it can be read from Table 2, to the value $K=15$, in the unfolded data system there correspond $X=55$, and symmetrically, $X=20$; and to the C.I. [15,16], there correspond a left C.I. [20,25], and a right C.I. [50,55].

5. VALIDATION OF THE PROPOSED METHOD AGAINST SOME BASIC ECOLOGICAL SITUATIONS

In this section the proposed method is validated against some basic ecological situations, providing at the same time a verification of the applied algorithms.

We shall consider some simple distribution changes which typically occur in plant ecology, dealing with the change from one discrete distribution to another, called for brevity *left* and *right distributions* and denoted by ξ and η , respectively.

In all illustrative situations, both the left and the right distributions will have five possible values (0, 1, 2, 3, 4) and we set $N=80$ and $K=30$. With these parameters we will generate $h=100$ random samples, testing our method on these samples, as explained below.

First, having fixed the above parameters, we generate $h=100$ random samples of size N , such that for each sample, the first K elements are taken from the given left distribution ξ and the rest of them from the given right distribution η . Then, we apply Algorithms 1 and 2 to every single randomly generated sample (taking $m=100$ in Algorithm 2). In this way, for each sample, Algorithm 1 will return an estimate for K , and Algorithm 2 will provide a change-interval for K . Finally, we shall have 100 estimated K values and 100 change-intervals with level 90%. In fact, we can check whether in 90 cases out of 100, the final change-interval includes the real, previously fixed value of K .

Finally, the mean of the 100 estimated K values will be accepted as change-point. Similarly, the final change-interval will be obtained from the means of the corresponding 100 estimated endpoints. We will also calculate the corresponding standard deviations.

It is intuitively clear that, the larger the Euclidean distance $|p - q|$ between the left and the right distributions $p=(p_0, p_1, p_2, p_3, p_4)$ and $q=(q_0, q_1, q_2, q_3, q_4)$ respectively, is, the smaller the obtained change-interval should be.

Now, in order to test our approach we will perform the above calculations with illustrative data, and with an a priori fixed change-point, dealing with distribution changes which typically occur in plant ecology.

5.1. The left distribution is symmetric and the right one is not

By considering the left and right distribution $p = (0.075, 0.125, 0.6, 0.125, 0.075)$ and $q = (0.1, 0.8, 0.05, 0.03, 0.02)$ respectively, (see Figure 4), our method gives the following results: for K we get 30.02, the change-interval is (26.17, 33.85) and the standard deviations are, respectively: 2.093665; 2.835757, 2.793842.

5.2. The abundance of a plant species changes in space

By considering the left and right distribution $p = (0.02, 0.03, 0.1, 0.6, 0.25)$ and $q = (0.1, 0.8, 0.05, 0.03, 0.02)$, respectively, (see Figure 5), our method provides the following results. The estimate for K is 30.08, the change-interval is (28.7, 31.61) and the respective standard deviations are: 0.9393744; 1.184922, 1.340096.

5.3 A non-uniform symmetric distribution if changed to a uniform one

By considering the left and right distribution $p = (0.05, 0.1, 0.7, 0.1, 0.05)$ and $q = (0.2, 0.2, 0.2, 0.2, 0.2)$, respectively, (see Figure 6), which biologically means that an

aggregation disappears, the proposed method gives an estimate for K equal to 31.3, a change-interval (22.33, 46.78) and the corresponding standard deviations are: 6.857128; 4.653781, 10.30208.

Summing up, we emphasize that in all considered cases (which occur very often in ecology) our method gives appropriate results in the sense the estimation of K is very close to its theoretical value 30. Moreover, the calculated change-interval always contains this theoretical value in its interior. These results not only validate our model, but at the same time verify the appropriateness of our algorithms, too.

6. CONCLUSIONS

Based on the data of a plant species, bramble (*Rubus fruticosus* L.) collected in an experimental forest gap, we have shown how a bootstrap method can be applied for the estimation of the changes in plant densities implied by human intervention. At this initial stage of our study we investigated a relatively small data set concerning a single species, in a real situation there may be about 100 plant species, and different species usually to respond differently to environmental changes.

We emphasize that our method is not only another approach for the estimation of a change-point; the estimate of the change-interval we offer can be applied not only based on the maximum likelihood principle we used in this paper, but any point estimation method known for edge detection can be developed in this way to get an estimate of the change interval.

Once we have estimated where the densities of different plant species change, we will be able to investigate whether these plant species change in the same zone, or as a response to a changed environment, a special “plant community” has been formed. However, this may be the topic of further studies.

ACKNOWLEDGEMENTS

The authors wish to thank the Ministry of Education and Science of Spain for the financial support of the project TIN2007-67418-C03-02. The Hungarian Scientific Research Fund OTKA (Grant No. K62000) also supported the research, the final version of the paper has been realized in the framework of the Hungarian–Spanish intergovernmental scientific and technological collaboration, with the support of the Scientific and Technological Innovation Fund (of Hungary) and the Ministry of Education and Sciences (of Spain HH2008-0023). During the research TS and JG were grantees of the János Bolyai Scholarship.

REFERENCES

[1] Becker, B., Lawrence J., Belisle P., Wolfson, D.B., Platt W.J., 2007. Bayesian change point analyses in ecology. *New Phytologist*. 174, 456-467.

- [2] Brodsky, B.E. and Darkhovsky, B.S., 1993. *Nonparametric Methods in Change-point Problems*. Kluwer Academic Publishers, The Netherlands.
- [3] Camarero, J. J, Gutiérrez E., Fortin, M. J., 2000. Boundary detection in altitudinal treeline ecotones in the Spanish Central Pyrenees. *Arctic, Antarctic, and Alpine Research* 32, 117-126.
- [4] Csörgö, M. and Horváth, L., 1997. *Limit Theorems in Change-point Analysis*. Wiley, Chichester.
- [5] Dale, M.R.T., 1999. *Spatial pattern analysis in plant ecology*. Cambridge University Press, Cambridge.
- [6] Efron, B. and Tibshirani, R., 1993. *An Introduction to the Bootstrap*. Chapman & Hall, New York, 184-188.
- [7] Ferger, D., 1993. Asymptotic distribution theory of change-point estimators and confidence intervals based on bootstrap approximation. *Mathematical Methods of Statistics, Vol. 3, No. 4*, 362-378.
- [8] Fortin, M. J., Olson, R. J., Ferson, S., Iverson, I., Hunsaker, C., Edwards, G., Levine, D., Butera, K. & Klemas, V., 2000. Issues related to the detection of boundaries. *Landscape Ecology* 15, 453-466.
- [9] Gálhidy, L., Mihók, B., Hagyó, A., Rajkai, K. and Standovár, T., 2006. Effects of gap size and associated changes in light and soil moisture on the understorey vegetation of a temperate deciduous forest. *Plant Ecology* 183, 133-145.
- [10] Humboldt, A. von., 1805. *Essai sur la géographie des plantes*. Paris: Levrault, Schoell et Cie.
- [11] Laurance, W. F, Didham R. K & Power M. F., 2001. Ecological boundaries: A search for synthesis. *TREE* 16, 70-71.

- [12] López, I., Gámez, M., Garay, J., Standovár, T. and Varga, Z., 2010. Application of Change-Point Problem to the detection of plant patches. *Acta Biotheoretica*, 58, 51-63.
- [13] Matthews, J.D., 1991. *Silvicultural Systems*. Calderon Press, Oxford.
- [14] McIntosh, R.P. (ed.), 1978. *Phytosociology*. Dowden, Hutchinson, and Ross, Stroudsburg, Pennsylvania.
- [15] Mihók B., Gálhidy L., Kelemen K., and Standovár T., 2005. Study of Gap-phase Regeneration in a Managed Beech Forest: Relations between Tree Regeneration and Light, Substrate Features and Cover of Ground Vegetation. *Acta Silv. Lign. Hung.*, 1, 25-38.
- [16] Podani, J., 2000. *Introduction to the Exploration of Multivariate Biological Data*. Backhuys Publishers, Leiden.
- [17] Reed, W.J., 2000. Reconstruction their history of forest fire frequency: Identifying hazard rate change points the Bayesian information criterion. *Can. Jour. Stat.* 28, 353-365.
- [18] Schleip, C. Menzel, A., Estrella, N., Dose, V., 2006. The use of Bayesian analysis to detect recent change in phenological events throughout of years. *Agri. Forest Meteor.* 141, 179-191.