When did author B take over for author A? Confidence Distributions for Change-Points

Céline Cunen, Gudmund Hermansen, Nils Lid Hjort Department of Mathematics, University of Oslo cmlcunen@math.uio.no

Introduction

The 15th century Catalan chivalry romance *Tirant lo Blanch* is an important work in medieval literature, and a perfect subject for a change-point analysis. The first author, Joanot Martorell, died before the completion of the novel and his friend and fellow knight, Marti Joan de Galba, took over the task. Many analyses, both statistical and literary, have attempted to identify exactly where the change-of-author takes place. Here, we propose our solution, using a new method which both estimates the change-point and provides the associated confidence curve.

Average word length

For the average word length in each chapter we assume the following model,

 $y_i \sim N(\mu_L, \sigma_L^2/m_i)$ for $i \leq \tau$ $y_i \sim N(\mu_R, \sigma_R^2/m_i)$ for $i \geq \tau + 1$.

According to this model, our best guess is that the change of author takes place in chapter 345. Author B uses longer words on average, and is also more variable across chapters. (The mean number of words per chapter is $\bar{m} = 927$.)

	$\hat{\mu}_L = 3.96$	$\hat{\mu}_R = 4.13$ o	
4. 8. –	$\hat{\sigma}_{L} = 3.46$	$\hat{\sigma}_{R} = 4.40$	

Data

pos ala mía penía de aquest tre ball tustamet excusar me por ues Empero contiant en lo louiran be bonabos be tota los bens qui siuda ale bons defige fuppline lo defalliment dels definante. E pozta los bona propolitas degu des fins. E voltra lenyozia qui per la virtut comportara los de allimets ari en fui com en ozoe: en lo prefet tractat per mi pofate per inabuertencia:e pus verbabe ramet iunozancia me atteutre er ponore:no folamet de lengua an alefa en portoguefa. ADas enca ta de postoguela en vulgar vale ciana:perco que la nacio Don vo fo natural le pura alegrar e molt aiuvar per los tats e ta infignes actes coz bi fon. Supplicant vo ftra pirtuofifima fenyozia accep teu com de feruídoz affectat la fet obza:car fi defallimete alguni of fonscertament fenyoznes en part coula la oita lengua anglela bela qual en algunes partibes es impolfible poper be girar los vo cables attenet ala afectio e belig oue continuamet unch de feruis ofra reouptable fenyozia. 10 bauer fguaro ala ruoitat bela oz Dinacio e Diferencia De fentecies efi que per voltra virtut la comu

Effenete lo libre bel j valeros e firenu caus ller Etrant lo blancio Etfenetic lo libre bel j valeros e firenu cana ller Etrant lo blanch Brincep ze Lefar bel Briperi grech be Ko testinoble. Lo qual fon traduit de Angles en lengua porto aucía. E a pres en vulgar lengua valèciona p lo magnifich : e virtuos cauallet/ molte jobanot martozell. Lo qual per most fua non poque acabar be trabute fino les tres parts. La quat ta part que ce la fi del libre ce ftada travatoa apzegaries bela noble fen yoza boa plabel de lozic plo mag mfich canaller molten Darri joba o galbaze it defait bi fera trobarvol fie atribuit ala fua ignozancia. Ell qual noître lenyoz Befu crift per la fua mmenía bondat vulla donar en premi de fos treballs la gioría ó pa radis. E proteíta que fi en lo de li / bre boura polabes algunes coles i no fien cariooliques que no les vo bauer bace.ans les remet a cozres cio ocla fancta catholica fglefia.

Our data is the text of an ancient Catalan book with 487 chapters and we search for the chapter where author B takes over from author A.

- From the text we can measure, for all chapters $i = 1, \ldots, 487$:
 - $y_i = \text{average word length}$ in chapter i, $m_i = \text{number of words}$ in chapter i
 - \mathbf{z}_i = the vector of proportions of words of length 1 to 10 in chapter *i*.

Fon acababa o empremptat la pre internationa en la Lintat de Elalencia en la Lintat de Elalencia en la Lintat de Elalencia ela nativitat de noître ienyor beu nu Feiu crist mil.occe.linte.

For example $z_1 = (0.08, 0.23, 0.17, 0.08, 0.13, 0.08, 0.06, 0.07, 0.04, 0.07)$, so 8% of the words in chapter 1 are one-letters words (length 1), and 7% of the words have more than 10 letters.

General method

We denote the unknown change-point by τ . To estimate τ and the associated uncertainty we use a method described in Cunen, Hermansen and Hjort (2016). In this case, we assume independent data y_i , following a parametric model with

$Y_i \sim f(y, \theta_L) \quad \text{for } i \leq \tau$



Proportion of words of different length

For the vector of proportions of words of length 1 to 9 (disregarding element no. 10, since the proportions sum to one) we assume the following multinormal distribution, $N_{i}(t - \sum_{i=1}^{n} t_{i}) \leq t_{i}$

 $oldsymbol{z}_i \sim \mathrm{N}_9(oldsymbol{\xi}_L, \Sigma_L/m_i) ext{ for } i \leq au$ $oldsymbol{z}_i \sim \mathrm{N}_9(oldsymbol{\xi}_R, \Sigma_R/m_i) ext{ for } i \geq au + 1.$

 $Y_i \sim f(y, \theta_R)$ for $i \ge \tau + 1$

We can then form the profile log-likelihood function

$$\ell_{\text{prof}}(\tau) = \max\{\ell(\tau, \theta_L, \theta_R) : \text{all } \theta_L, \theta_R\} \\ = \sum_{i \le \tau} \log f(y_i, \widehat{\theta}_L(\tau)) + \sum_{i \ge \tau+1} \log f(y_i, \widehat{\theta}_R(\tau)).$$

The maximiser $\hat{\tau}$ of this function is a good estimate of the change-point. From this we form the deviance function

 $D(\tau, Y) = 2\{\ell_{\text{prof}}(\widehat{\tau}) - \ell_{\text{prof}}(\tau)\}.$

In order to find the full confidence curve (essentially confidence intervals at different levels) we need the estimated distribution of $D(\tau, Y)$ for each position τ ,

$$cc(\tau) = \Pr_{\tau,\widehat{\theta}_L,\widehat{\theta}_R} \{ D(\tau, Y) < D(\tau, y_{obs}) \}.$$

This can be computed via stochastic simulation: we generate a large number B of simulated datasets Y^* from $f(y, \hat{\theta}_L)$ to the left of τ and $f(y, \hat{\theta}_R)$ to the right, at each candidate value τ , and calculate

$$cc(\tau) = \frac{1}{B} \sum_{j=1}^{B} I\{D(\tau, Y_j^*) < D(\tau, y_{obs})\}.$$

The estimated mean vectors indicate that author B differs from author A in using more one-letter words and also using more long words (with more than 7 letters). The change-of-author point-estimate according to this model is chapter 371, but the model also places some confidence in the change taking place earlier, in chapter 345.





Conclusion and references

Using word lengths as a measure of literary style, we have discovered two clear candidates for the chapter where Marti Joan de Galba took over from Joanot Martorell. Both our analyses indicate that the change in *Tirant lo Blanch* takes place in the last quarter of the book. This is consistent with previous studies, both statistical and literary.

Chen, H. & Zhang, N. (2015). Graph-based change-point detection. Annals of Statistics 43, 139–176.
Cunen, C., Hermansen, G. and Hjort, N.L. (2016). Confidence distributions for change-points and regime shifts. To appear in Journal of Statistical Planning and Inference.

Riba, A. & Ginebra, J. (2005). Change-point estimation in a multinomial sequence and homogeneity of literary style. *Journal of Applied Statistics* 32, 61–74.



