

Classifying New-York Times comments using Hierarchical Optimal Topic Transport

Machine Learning for Natural Language Processing 2020

Clément Guillo

ENSAE, ENS Paris-Saclay
clement.guillo@ensae.fr

Hugo Thimonier

ENSAE, ENS Paris-Saclay
hugo.thimonier@ensae.fr

Abstract

The main goal of our work¹ is to predict the category to which comments posted on New-York Times articles belong. We used the HOTT methodology which allows us to evaluate the distances between identified topics in the corpus and to predict NYT categories using topic proximity. Our results suggest that overall comments posted in different categories address similar topics, nevertheless some categories can still be identified thanks to the topics discussed in their comments.

1 Problem Framing

The objective of our work² is to identify the main topics that are covered in the New-York Times comments, to infer differences between article categories in the topics addressed in the comments and to use these differences to predict the category of the article. In order to do so we used HOTT (Hierarchical Optimal Topic Transport) (Yurochkin et al., 2019) (1) which relies on identifying topic similarities between documents to predict the categories to which documents belongs.

This approach benefits from several assets such as smaller computational complexity in comparison with other Optimal Transport (OT) methods to evaluate documents similarity, but also interpretability of the results. This methodology will also allow us to identify the main topics treated in the comments since HOTT requires to use Latent Dirichlet Allocation, a hierarchical Bayesian model that permits to identify the topics discussed in the comments.

The problem at stake is thus a multi-class classification, as we first considered 31 categories to

which articles can belong and then reduced it to 7 categories.

2 Experiments Protocol

2.1 Dataset

The dataset we studied is a balanced subsample of all comments posted in April 2018 on articles of the NYT website. The dataset is composed of 5971 comments, the average number of words per document is 401. After removing stopwords, words not contained in the Glove Embedding dataset, and words appearing less than twice (see Figure 1 and 2) the latter number is reduced to 65. The vocabulary is comprised of 6250 words in total.

The dataset was constructed so that each category we wish to predict has more or less the same number of comments. For that matter, we removed the categories which did not have enough comments inside them, reducing the number of categories from 39 to 31.

2.2 Model : HOTT (Hierarchical Optimal Topic Transport)

We refer the reader to the Appendix A for more detail on the HOTT methods, however we give in this section a brief summary of the procedure.

The Word mover's distance (Kusner et al., 2015) (2) is defined thanks to the 1-Wassertein distance and given an embedding of a vocabulary. A corpus is represented using a distribution over the vocabulary via a normalized bag of words. Then the word mover's distance between two documents d^1, d^2 is defined as the 1-Wasserstein distance between them :

$$\text{WMD}(d^1, d^2) \equiv W_1(d^1, d^2)$$

The HOTT method consists in first estimating using LDA (Latent Dirichlet Allocation) topics (de-

¹https://github.com/hugothimonier/HOTT_NLP_ENSAE

²Google Colab

defined as distribution over words) among the documents. Documents are now defined as distribution over topics, in other words a document can be summarized by the percentage of each topic present within the document. Hierarchical Optimal Topic Transport is thus defined as

$$HOTT(d^1, d^2) = W_1\left(\sum_{k=1}^T \bar{d}_k^1 \delta_{t_k}, \sum_{k=1}^T \bar{d}_k^2 \delta_{t_k}\right)$$

where, $\bar{d}^i \in \Delta^T$ is the document distribution, the Diracs δ_{t_k} are probability distribution on the corresponding topic t_k and the metric $d(., .)$ used is the WMD between topics.

3 Implementation

3.1 K-Nearest Neighbours Classification

The metric described above is used to perform knn classification. We proceed as follows : (i) we compute³ for each document i in the test sample its HOTT distance with every document in the training sample, (ii) we only keep the k nearest documents for each document i (where k is a hyperparameter to optimize) and finally, (iii) we check the most recurrent label in the k nearest documents of document i and assign it to be its prediction.

3.2 Multiclass prediction : 31 categories

When considering the 31 classes the model performs very poorly even after parameter tuning (less than 20% accuracy as one can see in Figure 3). This suggests that when considering a high number of categories, there is no sufficient differences in terms of topics addressed between categories for our predictions to be accurate. This intuition is confirmed by the t-sne representation⁴ (Van der Maaten et al. 2008) (4) : one can note on this graph that there is no cluster by category. In other words, it is unclear whether comments posted in a particular article category are closer to comments posted in another article category than any other category as all categories seem to be mixed. Figure 7 also confirms such intuition as there is no clear pattern of confusion : when a document is not assigned to the right category, it is not assigned to another category in particular which could be considered closer to the true category.

³Using the pre-trained Glove (300d) embedding (3)

⁴This T-SNE representation is constructed using the HOTT metric.

3.3 Multiclass prediction : 7 categories

Considering only 7 categories of articles⁵, we manage to highly improve the accuracy of our prediction and the interpretability of the results. Figure 6 shows notable results as one can now observe how close categories' comments are : comments posted on the 'Games' section (bottom left of the graph) and those in the 'Science' section (top right of the graph) are apparently not discussing the same topics at all, on the other hand 'Climate' and 'Foreign' seem to discuss topics close to each other. One can also see that the 'Art' category is mixed with all other categories in the middle of the graph suggesting no major difference in terms of topics addressed with respect to the other categories. Figure 9 and 10 in the appendix also confirm those statements. Figure 8 also confirms this intuition : when 'Climate' is not correctly predicted, the most recurrent error is to predict 'Foreign', the converse is also true. Similarly, the most correctly predicted category is Science, which was the most isolated category on the graph. Finally, the 'Art' category which was the least isolated category on the graph is also the category which is the least correctly predicted.

4 Conclusion

When considering all the categories of the New-York Times articles, the HOTT - KNN classification performs poorly : this suggests that comments between so many categories are not different enough in terms of topics addressed to be able to differentiate each category from the others. On the other hand, when only considering 7 categories, it appears that the topics addressed in the comments do depend on the category of the article : 'Foreign' comments and 'Climate' comments seem to address the same type of topics, while 'Science' apparently addresses its own topics which differ from the other categories, while 'Arts and Leisure' comments seem to discuss topics that are also discussed in the other categories. Thus, the HOTT metric allowed us to have both good prediction result for a small number of categories but also an informative interpretation when the model performed badly for all the categories.

⁵Dining, Games, Foreign, Arts&Leisure, Science, Sports and Climate

Appendices

A HOTT Metric

Let us define the Word mover's distance (Kusner et al., 2015) defined thanks to the Wasserstein distance. Let $x = \{x_1, \dots, x_n\}$ and $y = \{y_1, \dots, y_m\}$ be two sets of points in a metric space. Let $\Delta^n \subset R^{n+1}$ denote the probability simplex on n elements, and let $p \in \Delta^n$ and $q \in \Delta^m$ be distributions over x and y . Then, the 1-Wasserstein distance between p and q is defined by the following OT problem

$$W_1(p, q) = \begin{cases} \min_{\Gamma \in R_+^{n \times m}} \sum_{i,j} C_{i,j} \Gamma_{i,j} \\ \text{subject to } \sum_j \Gamma_{i,j} = p_i \\ \text{and } \sum_i \Gamma_{i,j} = q_j \end{cases}$$

where the cost matrix C , has each entry $C_{i,j}$ defined by $d(x_i, y_j)$, where $d(\cdot, \cdot)$ is a distance metric. In that framework, Γ is interpreted as a transport plan between p and q .

The Word mover's distance is then defined given an embedding of a vocabulary, $V \subset R^n$. A corpus $D = \{d^1, d^2, \dots, d^{|D|}\}$ is represented using a distribution over V via a normalized bag of words. For $d^i \in \Delta^{l_i}$, where l_i denotes the number of unique words in d^i , then each entry $d_j^i = \frac{c_j^i}{|d^i|}$ where c_j^i is the count of word v_j in document d^i and $|d^i|$ is the number of words in d^i . Then the word mover's distance is defined using the euclidean metric between word embeddings as a distance between elements such that

$$WMD(d^1, d^2) \equiv W_1(d^1, d^2)$$

The HOTT method consists in first estimating using LDA (Latent Dirichlet Allocation) topics among the documents. Topics are hereby defined as distribution over words, thus $T = \{t^1, \dots, t^{|T|}\} \subset V$ meaning that each topic contains certain words of the whole vocabulary of the corpus. Similarly, documents are now defined as distribution over topics such that $\bar{d}^i \in \Delta^T$, in other words a document can be summarized by the percentage of each topic present within the document. Hence, WMD can be applied to topics such that HOTT is defined by WMD between the documents as

$$HOTT(d^1, d^2) = W_1\left(\sum_{k=1}^T \bar{d}_k^1 \delta_{t_k}, \sum_{k=1}^T \bar{d}_k^2 \delta_{t_k}\right)$$

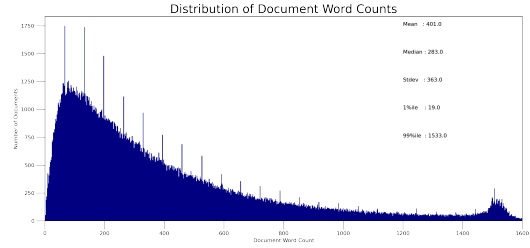
where the Diracs δ_{t_k} are probability distribution on the corresponding topic t_k and the metric $d(\cdot, \cdot)$ used is the WMD between topics.

In both (1) and (2) pairwise distance between each document must be computed, however in (2) the computational time of $WMD(d^1, d^2)$ can be very long as it highly depends on $l = \max(l^1, l^2)$ where l^i stands for the number of unique words in document i . On the other hand, the complexity of $HOTT(d^1, d^2)$ only depends on the number of topics. Thus this allows to drastically reduce the computational time.

B Figures

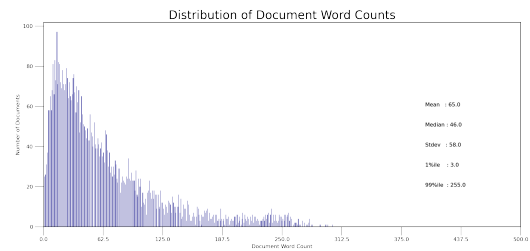
B.1 Word distributions

Figure 1: Distribution of Documents word counts before treatments



Note : Mean = 401; Median = 283; Standard Deviation = 363; 1st percentile = 19; 99th percentile = 1533.

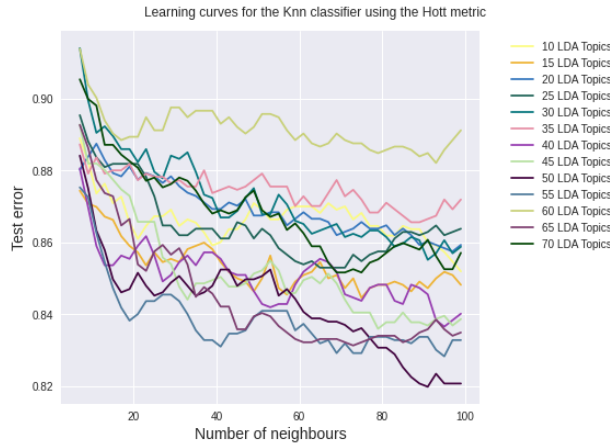
Figure 2: Distribution of Documents word counts after treatments



Note : Mean = 65; Median = 46; Standard Deviation = 58; 1st percentile = 3; 99th percentile = 255.

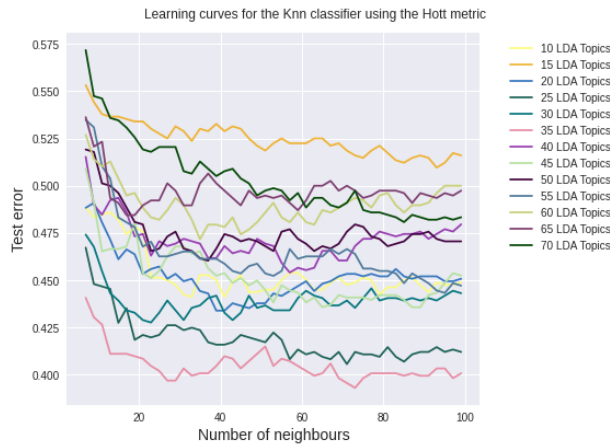
B.2 Model Performances

Figure 3: Model performance wrt to the number of LDA Topics and the number of neighbours considered (31 categories)



Note : The parameter tuning performed on the number of neighbours and the number of LDA Topics for the knn classification's performance when we consider 31 categories does not significantly improve the results of the prediction. The test error is the lowest at 82% for more than 80 neighbours and 50 LDA topics.

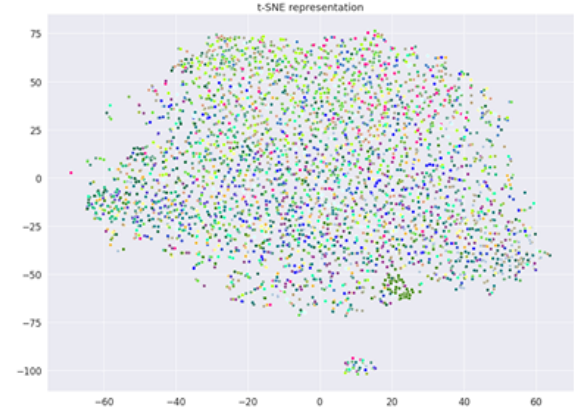
Figure 4: Model performance wrt to the number of LDA Topics and the number of neighbours considered (7 categories)



Note : The parameter tuning when we consider 7 categories improves the results of the prediction by 20%. The test error is the lowest at 40% for more than 75 neighbours and 35 LDA topics. However for the sake of interpretability we decided to choose the 20 LDA topics model which displays very satisfactory results as well.

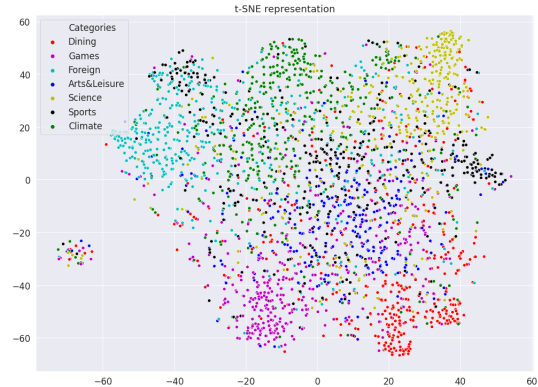
B.3 T-SNE representations

Figure 5: T-SNE representation using HOTT metric (20 LDA topics, 31 classes)



Note : This T-SNE representation illustrates how there is no clear separation between article categories in terms of topics addressed in the comments. For instance, we can see no clear separation between topics on this graph as all labels are mixed.

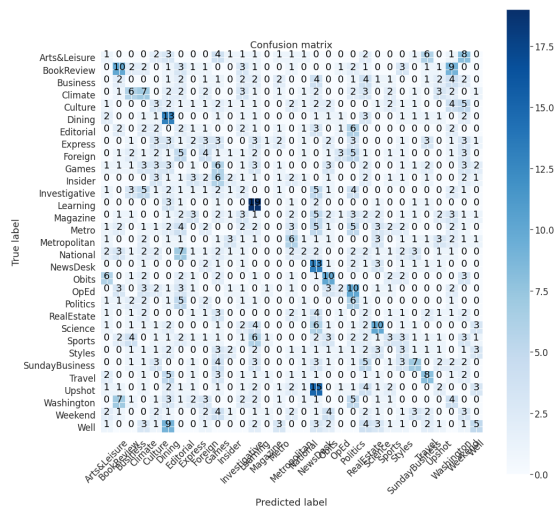
Figure 6: T-SNE representation using HOTT metric (20 LDA topics, 7 classes)



Note : Contrary to figure 5 this graph displays clear differences between article categories' comments in the topics space. For instance, one can see 'Dining' is concentrated at the bottom right of the graph, 'Science' at the top right or 'Games' at the bottom left. Moreover, one can see that 'Foreign' and 'Climate' seem to be concentrated in the same area of the graph which suggests the topics addressed in the comments of this category are quite similar. One can also see that the 'Art' category is mixed in all categories in the middle of the graph suggesting no major difference in terms of topics addressed with respect to the other categories.

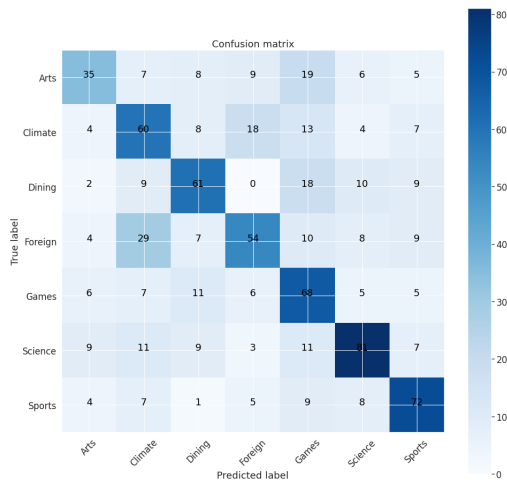
B.4 Confusion Matrices

Figure 7: Knn Classification Confusion matrix (20 LDA topics, 31 classes)



Note : As suggested by the t-sne representation in figure 5 as there are no clear differences in topics addressed between article categories, the model performs poorly and only very few categories are well predicted. The best (although very low) predictions include Learning and Dining categories.

Figure 8: Knn Classification Confusion matrix (20 LDA topics, 7 classes)



Note : This confusion matrix displaying the results of the knn classification for only 7 categories confirms the information that was suggested in figure 6 : Art is not very well predicted as the topics addressed in its comments are not sufficiently different from the other categories. Science is the most correctly predicted category : it was the most isolated category on figure 6. Finally, Climate and Foreign which seemed to share topics on the t-sne representation are often predicted to be in the other category.

C Topic repartition

Figure 9: Category repartition per topic

category	Arts&Leisure	Climate	Dining	Foreign	Games	Science	Sports
topics							
1	4.87	3.98	7.30	9.07	4.42	60.84	9.51
2	5.52	21.38	4.14	42.76	5.52	6.90	13.79
3	2.08	9.69	1.73	74.39	1.38	3.81	6.92
4	21.85	4.20	26.05	9.24	15.97	11.76	10.92
5	7.51	1.02	10.24	1.37	72.01	2.05	5.80
6	1.03	78.08	5.14	3.08	3.08	7.53	2.05
7	4.72	0.47	3.30	0.94	2.36	4.25	83.96
8	20.19	6.73	14.42	15.38	10.58	17.31	15.38
9	41.41	7.07	12.12	9.09	18.18	3.03	9.09
10	17.33	18.67	17.33	14.67	14.67	8.00	9.33
11	10.43	2.61	57.39	1.74	17.39	6.09	4.35
12	2.86	28.57	7.14	14.29	20.00	11.43	15.71
13	19.75	14.81	7.41	14.81	18.52	12.35	12.35
14	30.47	14.06	8.59	11.72	12.50	10.94	11.72
15	12.99	20.78	10.39	19.48	14.29	10.39	11.69
16	4.59	19.27	11.93	20.18	15.60	11.93	16.51
17	45.26	7.37	9.47	9.47	13.68	7.37	7.37
18	8.84	2.21	63.54	3.87	11.05	6.08	4.42
19	5.13	9.40	16.24	13.68	21.37	9.40	24.79
20	6.58	22.37	14.47	10.53	10.53	9.21	26.32

Note : One reads the table as follows : 60.84% of topic 1 is addressed in the 'Science' category's comments. This table confirms the fact that the 'Science' category is very different from other categories as most of its content is concentrated in topic 1 which is not so much discussed in other categories' comments. One can also note that 'Foreign' and 'Climat' display a very similar pattern.

Figure 10: Topic repartition per Category

category	Arts&Leisure	Climate	Dining	Foreign	Games	Science	Sports
topics							
1	6.75	3.82	7.67	8.27	4.21	58.51	9.33
2	2.45	6.58	1.40	12.50	1.68	2.13	4.34
3	1.84	5.94	1.16	43.35	0.84	2.34	4.34
4	7.98	1.06	7.21	2.22	4.00	2.98	2.82
5	6.75	0.64	6.98	0.81	44.42	1.28	3.69
6	0.92	48.41	3.49	1.81	1.89	4.68	1.30
7	3.07	0.21	1.63	0.40	1.05	1.91	38.61
8	6.44	1.49	3.49	3.23	2.32	3.83	3.47
9	12.58	1.49	2.79	1.81	3.79	0.64	1.95
10	3.99	2.97	3.02	2.22	2.32	1.28	1.52
11	3.68	0.64	15.35	0.40	4.21	1.49	1.08
12	0.61	4.25	1.16	2.02	2.95	1.70	2.39
13	4.91	2.55	1.40	2.42	3.16	2.13	2.17
14	11.96	3.82	2.56	3.02	3.37	2.98	3.25
15	3.07	3.40	1.86	3.02	2.32	1.70	1.95
16	1.53	4.46	3.02	4.44	3.58	2.77	3.90
17	13.19	1.49	2.09	1.81	2.74	1.49	1.52
18	4.91	0.85	26.74	1.41	4.21	2.34	1.74
19	1.84	2.34	4.42	3.23	5.26	2.34	6.29
20	1.53	3.61	2.56	1.61	1.68	1.49	4.34

Note : One reads the table as follows : 6.75% of the topics addressed in the Arts and Leisure category is topic 1.

References

- [1] Mikhail YUROCHKIN, Sebastian CLAICI, Edward CHIEN, Farzaneh MIRZAZADEH and Justin SOLOMON. *Hierarchical Optimal Transport for Document Representation*. NeurIPS 2019.
- [2] KUSNER, M., SUN, Y., KOLKIN, N., and WEINBERGER, K. *From word embeddings to document distances*. In International Conference on Machine Learning, pp. 957–966, 2015.
- [3] Jeffrey PENNINGTON, Richard SOCHER, and Christopher D. MANNING. *GloVe: Global Vectors for Word Representation*, 2014.
- [4] L.J.P. VAN DER MAATEN and HINTON, G.E., *Visualizing High-Dimensional Data Using t-SNE*, Journal of Machine Learning Research, vol. 9, November 2008, p. 2579–2605.