# HEC MONTRÉAL

## Master in Financial Engineering - Supervised Projet

---

# Developing Category-specific Economic Policy Uncertainty Indices through Genetic Algorithm

---

*Author:*
Jun Wang 11238309

*Professeur:*
David Ardia

April 6, 2022

# Contents

# 1  Introduction

There is a growing interest in building uncertainty indices from texts to measure various economic outcomes. The use of uncertainty indices was introduced by Baker et al. (2016) who developed a new index of economic policy uncertainty based on newspapers coverage frequency for the United States. Their work shows that frequency counting methods for keywords on newspaper articles can be used to construct valuable proxies for policy-related economic uncertainty movements. Manela et al. (2017) create a VIX-optimized uncertainty index using SVM regression and n-gram counts from newspapers. They let the index be constructed from how well it captures the dynamics of VIX.

This project presents a new approach to building category-specific Economic Policy Uncertainty (EPU) indices by extracting keywords from texts via the genetic algorithm. Knowing that the EPU indices vary with different sets of EPU terms, we let the EPU terms vary to fit a specific target variable to obtain the final optimized EPU index. We proceed in three steps: first, with initial EPU terms, we use a pre-trained word embedding space to extend the set of candidate keywords.[1] This allows us to quickly access similar keywords in the nearest dimension in each topic. Second, using the genetic algorithm, we search for the best subset of candidate keywords that matches the dynamics of a pre-determined target variable. Last, following the same steps proposed by Baker et al. (2016), we build the new category-specific EPU index.

Compare with the approach by Baker et al. (2016) , our method does not require the careful human audit process to select the keywords. Since we derive keywords from newspaper sources, this method is less costly and can be easily adapted to other uncertainty measures and economic outcomes. We find that the extracted keywords match our intuition and our final VIX-optimized index achieves high out-of-sample correlation. Our approach is also a proof-of-concept for the Baker et al. (2016) approach, showing the feasibility of using frequency counts of newspaper articles to construct reasonable proxies for economic variables. For different target variables, we are able to construct their own newspaper coverage-based uncertainty indices, which capture their unique movements.

# 2  Literature Review

This section covers recent investigations of building newspaper-based uncertainty indices.

## 2.1  Building EPU keywords

To construct the newspaper coverage-based uncertainty indices, one of the critical steps is to find the set of EPU keywords to form the search expressions.

In their pioneering work, Baker et al. (2016) select the EPU terms that minimize the gross error rate relative to the human benchmark. Their construction of a human standard requires extensive work of human reading and auditing of newspapers, which is essential to reduce human bias.

Ghirelli et al. (2019) use a richer set of keywords to build the EPU index for Spain. Their results suggest that keywords richness is one of the keys to improving the quality of the EPU index, in addition to widening press and time coverage. In the same spirit, Algaba et al. (2020) use a pre-trained word embeddings space to expand the candidate keywords. Azqueta-Gavaldon et al. (2020) use word embedding techniques to retrieve articles related to economic uncertainty and implement topic modeling to assign articles into specific categories to build the EPU index without using candidate keywords.

Tobback et al. (2018) apply two different text mining techniques to improve the original policy uncertainty index. The modality annotation counts the number of occurrences of words that express uncertainty across the entire data and assigns modality scores to every article to narrow the scope of articles addressing economic policy uncertainty. The support vector machine technique automatically looks for patterns in the text documents and selects the words with the most significant discriminative power. This technique replaces the predefined keywords of the original methodology, thus reducing potential audit bias. Their

---

[1]A word embedding space maps words into high dimensional vectors. Words with similar semantic meanings are close to each other in the space. We use the pre-trained embedding space from Pennington et al. with GloVe Technique

results show that the EPU index using text mining techniques outperforms the original EPU index proposed by Baker et al. (2016) in terms of the predictive power of macroeconomic and financial variables.

## 2.2 Uncertainty indices construction methodology

There are several ways in the existing literature to construct uncertainty indices. The first is based on the frequency count of pre-determined EPU keywords (see, e.g., Baker et al. (2016), Ghirelli et al. (2019)). The second uses text mining techniques to model topics of articles and calculates the topic frequency (Azqueta-Gavaldon et al. (2020)). The third uses machine learning techniques (e.g., Support Vector Machine classifier) to classify the EPU mentions (Tobback et al. (2018)). The fourth uses a token-distance-based triple approach to identify EPU mentions in the textual documents (Ardia et al. (2021)).

Apart from the construction of general-purpose uncertainty indices, the use of external variables as drivers to build the category-specific uncertainty indices is also represented in the literature. For example, Manela et al. (2017) estimate the relationship between VIX and n-gram counts using support vector regression and construct a text-based measure of disaster concerns, called News implied volatility. They show that the News implied volatility captures well the disasters concerns of the average investors, and disasters concerns drive expected returns.

# 3  Data

We consider articles from the Wall Street Journal from January 2001 to August 2021. The articles are converted to a word frequency matrix where the rows represent the article IDs, the columns represent the words, and the values represent the frequency counts of words. Stop words are removed, and all words are converted to lowercase. We use pre-trained Global Vectors from stanfordnlp/GloVe to find the candidate keywords. This word vectors space has 6 billion tokens, with 300-dimension-per-word and 400k vocabulary word vectors. We use the VIX[2], unemployment rate, and industrial production as target economic variables for constructing the category-specific EPU indices. The time series of VIX is from CBOE exchange. The other economic variables are from U.S. Bureau of Labor Statistics.

# 4  Methodology

## 4.1  Extension of candidate keywords

Following a human audit process, Baker et al. (2016) propose the following EPU triple: "economic" or "economy" for the topic Economy; "Congress", "deficit", "Federal Reserve", "legislation", "regulation" or "White House" for the topic Policy; "uncertainty" or "uncertain" for the topic Uncertainty. These three topics document the uncertainty about who will make economic policy decisions, what economic policy actions will be undertaken and when, and the economic effects of policy actions.

A large and appropriate set of candidates must be created for every topic to search for the most discriminative keywords towards the target variable. To achieve this, we first split the initial n-gram keywords into uni-gram keywords, then take the average of the pre-trained keyword vectors in every category as its vector representation of the topic. We lastly search in the space of GloVe Vocabulary for the top $n$ keywords with the highest cosine similarity to the vector representation of each topic.

$$Sc(A, B) = \frac{A \times B}{||A|| \times ||B||} = \frac{\sum_{i=1}^{n} A_i B_i}{\sqrt{\sum_{i=1}^{n} A_i^2} \sqrt{\sum_{i=1}^{n} B_i^2}}$$

where $A$ and $B$ denote the vector representation of keyword A and keyword B, respectively.

In this way, we manage to construct the $3 \times n_c$ candidate keywords ($n_c$ for every topic). We also add the original keywords to the candidate keywords set. In our case, the total number of candidate keywords is 63 (see Table 1).

---

[2]We use a monthly average of daily VIX, which is the same as in Baker et al.'s work

| EPU subcategory | $n_c$ | candidate keywords |
|---|---|---|
| Economy | 20 | **economic**, **economy**, growth, economies, recession, slowdown, recovery, downturn, financial, sector, crisis, global, slowing, inflation,reforms, market, boost, unemployment, markets, stability |
| Policy | 23 | **federal reserve**, **federal**, **congress**, **legislation**, **deficit**, **regulation**, **white house**, **house**, congressional, senate, bill, administration, lawmakers, spending, budget, government, republicans, proposal, measure, legislature, democrats, policy, reform |
| Uncertainty | 20 | **uncertainty**, **uncertain**, uncertainties, unclear, turmoil, doubt, unsettled, worries, prospects, outcome, unsure, instability, confusion, certainty, volatility, outlook, remains, concern, remain, worried |

Table 1: The universe of keywords. The keywords in bold are the original keywords proposed by Baker et al. The rest of the keyword sequence is ordered from highest to lowest by cosine similarity score.

## 4.2 Genetic Algorithm

We want to select a subset of keywords with the most discriminative power towards the target economic variable among these candidate keywords. This problem is essentially a subset selection problem which can be naturally handled through the genetic algorithm using a binary vector, with 1 indicating the presence of candidate keywords and 0 indicating their absence. We use the GA package in R; see https://10.18637/jss.v053.i04. Therefore, our initial population would be many binary vectors, where each vector represents each subset of the keywords.

### 4.2.1 EPU indices construction

With the population of binary vectors, we build the EPU indices following the same procedures as in Baker et al. (2016). The steps are:

1. Transform the binary vector into a list of selected keywords (the new EPU triple).

2. Count the number of articles containing at least one keyword in all three categories (Economy, Policy, and Uncertainty), namely the raw counts.

3. Divide the raw counts by the total number of articles in the same newspaper and month.

4. Standardize the new monthly series to unit standard deviation before a specific date and average across different newspapers.

5. Normalize the series to a mean of 100 before a certain date.[3]

### 4.2.2 Objective Function

We aim to identify the most relevant keywords for explaining variations of the target variable through genetic algorithm. To evaluate the fitness at every iteration stage, we need to develop an objective function. The objective function we defined is the absolute value of correlation between the percentage change of the EPU index and the target variable, plus some penalty on the total number of selected keywords. We also test other fitness functions like mean square root error or tracking error; see more in Appendix.

$$|cor(rets(EPU_t), rets(Y))| - \lambda \times n_s$$

---

[3]Since we only have a newspapers source in our case, the standardization and normalization may not be needed. However, we still keep the same steps for better comparison.

where $Y$ denotes the uncertainty measure of interest, $rets$ denotes the percentage change of the series, $n_s$ denotes the number of selected keywords, and $\lambda$ denotes the penalization parameter.

In this way, we extract only the valuable keywords that shift the EPU index significantly towards the selected $Y$ measure.

### 4.2.3 Calibration of penalization parameter

The objective function has a penalization parameter $\lambda$, which requires careful calibration. On the one hand, the penalization parameter cannot be too small; otherwise, it will result in too many unimportant keywords being selected. On the other hand, the penalization cannot be too big; otherwise, it will cause essential keywords not to be chosen, and the optimized EPU will poorly fit the target variable. Therefore, we set the $\lambda_{max}$ as the smallest value such that only one keyword is selected per dimension, and $\lambda_{min}$ as a fraction of $\lambda_{max}$ with a pre-specified ratio. In our case, we set the ratio of $\lambda$ to 0.1. So, $\lambda_{min} = \lambda_{max} \times 0.1$.

To find the $\lambda_{max}$, we first assign the $\lambda_{max}$ to a large number like ten and run the genetic algorithm to obtain the final set of selected keywords. If the final set has only one keyword per topic, we think the penalization is large enough, and we scale down the $\lambda_{max}$ with a factor of 10 to test its downside limit. So, we repeat the process until the final set of keywords has more than one keyword in any dimension and set the $\lambda_{max}$ as the current $\lambda_{max} \times 10$.

After the selection of $\lambda_{min}$ and $\lambda_{max}$, we then log scale from $\lambda_{min}$ to $\lambda_{max}$ to obtain the parameters sequence of length $n$ ($n = 10$ in our case).

## 5 Empirical results

We divide our data into three parts: training set, cross-validation set, and test set. The training period is from January 2001 to December 2011, the cross-validation period is from January 2012 to December 2016, and the test period is from January 2017 to August 2021. First, we run the genetic algorithm in the training set to obtain the best keywords for the target variables with the sequence of possible $\lambda$. After having a list of keywords set, we then evaluate the performance based on the absolute correlation with the percentage change of the target variable on the cross-validation set. The set of keywords with the highest fitness in the cross-validation set is then selected as the final keywords, by which a new target-specific EPU index can be built.
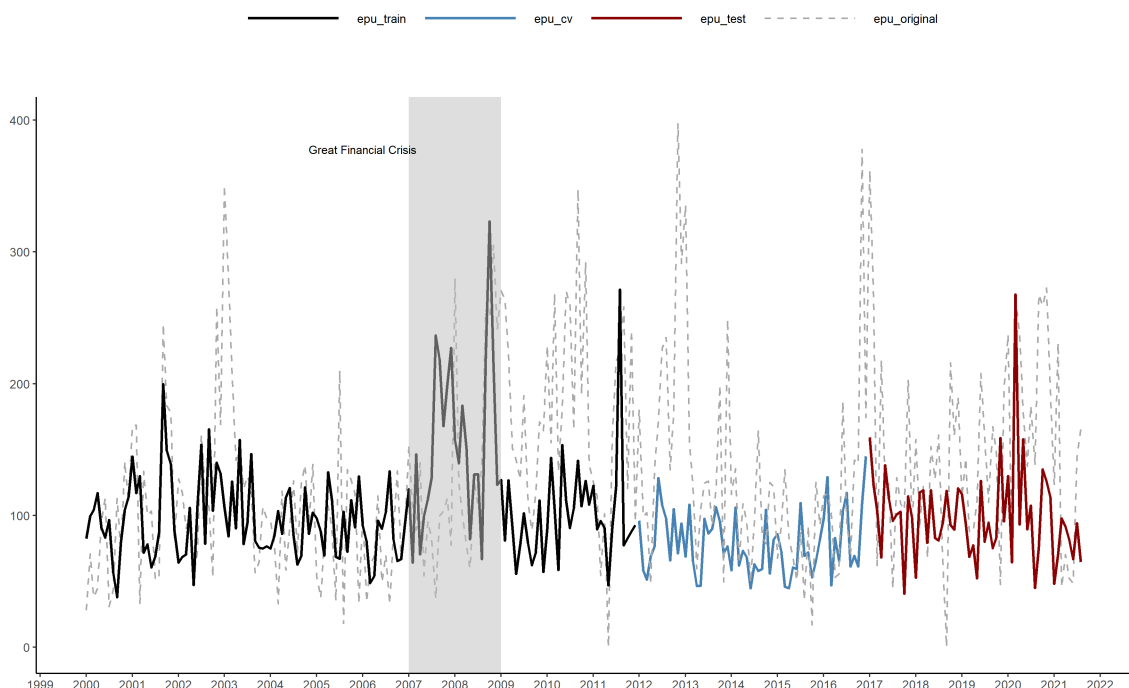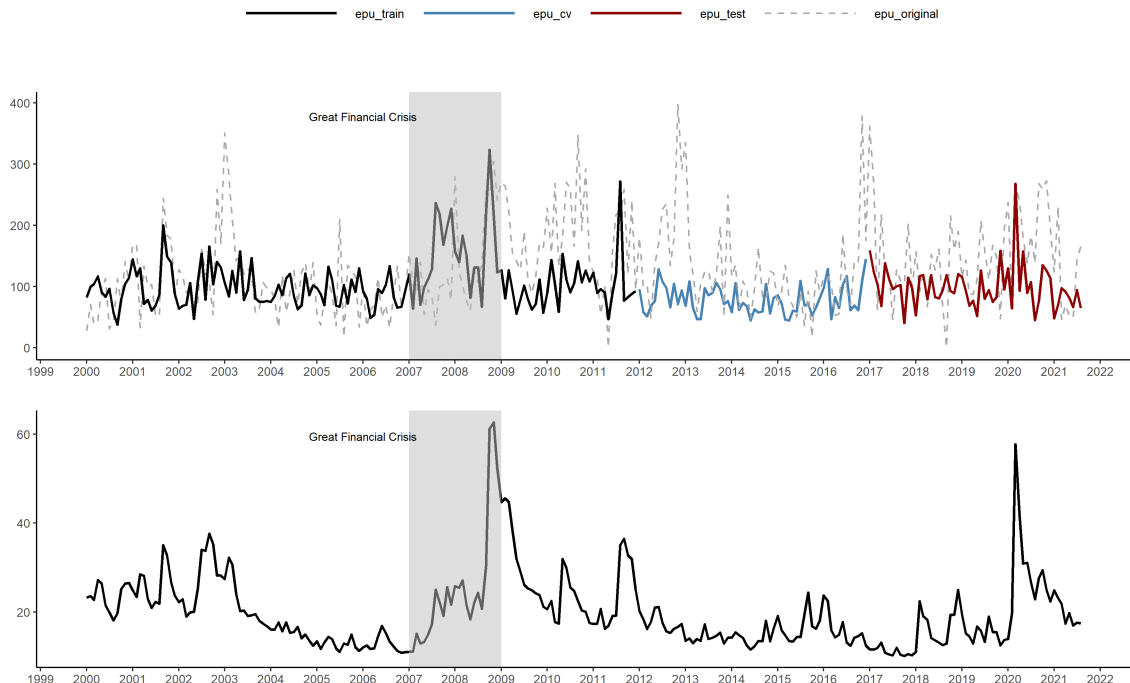


Figure 1: VIX-specific EPU Index

Figure 2: Comparison between the VIX index and VIX-specific EPU Index

Figure 1 shows the evolution of the VIX-specific EPU index from January 2000 to August 2020, together with the EPU index constructed from keywords proposed by Baker et al. (2016). We see an apparent deviation of our final optimized EPU index from the original EPU index (the grey dashed line). Our index responds more to the Great financial crisis, and we see a spike in covid 2020. As the stock market recovers, a decline in uncertainty corresponds to a drop in our VIX-specific EPU Index, which the original general-purpose EPU does not capture. Since we only use the information before 2016 to find our final keywords, the $EPU_{test}$ in the red line does not involve any forward-looking but is still able to capture the local movements in VIX. A correlation score of 0.83 between our EPU index and VIX in the test period further validates our index (See in Figure 2).

Table 2 shows the final chosen keywords for VIX. Out of the 63 keywords, 8 are selected to construct the final VIX-optimized uncertainty index. We see how the algorithm carefully selects the essential keywords for the target. The keywords "market", "stability", and "uncertainty" are crucial when talking about equity market uncertainty.

| EPU subcategory | $n_c$ | final optimized keywords |
| --- | --- | --- |
| Economy | 20 | **market**, **stability** |
| Policy | 23 | legislation, house, measure |
| Uncertainty | 20 | **uncertainty**, turmoil, worried |

Table 2: The final set of keywords for VIX-specified EPU index

## 5.1 Robustness

We look at three methods to further validate our approach to constructing the category-specific EPU indices. First, with 1,000 randomly generated lists of keywords set, we build the EPU indices and calculate the absolute correlation coefficients between the percentage change of EPU indices and the target variable in the testing period. We find that our VIX-specific EPU index has an absolute correlation score of 0.55 in the test set, which outperforms 95% of EPU indices in the same period. This proves the consistency of the selected keywords within time.

The second approach is to add arbitrary unrelated keywords (Table 3 in Appendix) to the solution keywords and rerun the genetic algorithm to see if the irrelevant keywords are selected. In the context of the VIX-specific EPU index, none of the unrelated keywords are selected, proving that our approach is not picking any keywords but keywords towards the target.

Last, we set the percentage change of the final optimized index as the independent variable and the percentage change of VIX as the dependent variable. We run a regression on the test period and find a p-value less than 0.0001 proving the relevance between the two indices.

The three methods presented above show that our model can select the most influential keywords among the 63 candidate keywords, and our index for VIX may be considered a reasonable proxy for the movements in VIX.

## 5.2    Extension to other target variables

We further extend our work to construct the unemployment rate-specific EPU index and industrial production-specific EPU index with the exact initial seed keywords. The final keywords are selected by maximizing the absolute correlation plus penalization for the two returns series. The EPU level series are then constructed with the final chosen keywords (Figure 2 and Figure 3 in the Appendix). The two EPU indices show a similar pattern as the original EPU index but are less volatile. Unlike the high correlation for percentage change of VIX-specific EPU index, the test correlations for the percentage change of unemployment rate-specific EPU index and industrial production-specific EPU index are -0.08 and 0.07, respectively, which are less than the test 95% quantile. The industrial production-specific EPU index can distinguish the final keywords from the noisy keywords, but this is not the case for the unemployment rate-specific EPU index. The high test p values for the two series fail to reject the null hypothesis, saying that our target-specific EPU indices may be uncorrelated to the target variables.

## 5.3    Limitations

So, what has led to the poor performance of the unemployment rate and industrial production? In this section, we will discuss three possible reasons behind them.

The first reason could be the selection of initial keywords to form the search space of the genetic algorithm. For example, when we set industrial production as the target variable, there are not many keywords related to industrial production in the candidate keyword space. The words "industrial" and "production" are not involved per se. If we do not have relevant keywords space, how can we expect our algorithm to find the right keywords to build the target-optimized EPU index? Therefore, one possible improvement of this index construction method is to extend the initial seed keywords to a larger space. For example, we could include all existing category-specific EPU triples in the initial seed keywords. Or, at least, we add the name of the target variable to the initial seed keywords. In addition, when we expand our candidate keywords through word embedding techniques, we could set the $n_c$ to a more significant number than 20 to form a more decadent keywords candidate space. In this way, we can explore more keywords relevant to the target variable. However, we can not ignore the increase in dimensionality in the richer keyword space. It would be a trade-off between efficiency and accuracy.

The second possible reason is the choice of the target variables. Not all economic indicators are suitable for expressing the uncertainty in words, and we may need to be careful about the choice of target variables. For example, in this project, we try to construct the unemployment rate-specific EPU index but find that our created index cannot capture well the movement of the unemployment rate. This may be because the spikes in the unemployment rate show the certainty of an economic downturn but not uncertainty, and the well-known lagging indicator may not be a good target variable.

The third reason lies in the lack of newspaper sources. When we compare our original EPU index with the EPU index with the index from Baker et al. (2016), we see a significant deviation in the period of 2020 to 2021.[4] Since the only difference is that we take one newspaper source while Baker et al. (2016) take ten different sources to construct the US index, we may need more newspaper sources. As suggested by Ghirelli et al. (2019), the wideness of newspapers coverage is essential to construct a useful EPU index.

---

[4]The US EPU index proposed by Baker et al. (2016) can be found here.

## 5.4 Fast EPU indices Computation

Our method requires the construction of numerous EPU indices, so it is important to optimize the construction procedure to save computation time. Here is how we proceed to construct the EPU index:

1. Sort the document features matrix $mat_{docfeat}$ by month and alphabetical order.

2. Create the dummy variables $mat_{doc2month}$ of dimension $nm$ and $nd$ where columns $(nd)$ denote the document ids, and rows $(nm)$ denote the months in which the document is published.

3. Transform binary vector into a matrix $mat_{keywords}$ of dimension $nf$ and 3 where columns (3) denote the topic to which it belongs, rows $(nf)$ denote the keywords.

4. Multiply the $mat_{docfeat}$ by the $mat_{keywords}$ to get the sum of keyword counts for each topic (Economy, Policy, and Uncertainty) for each document. The new matrix of dimension $nd$ and 3 is named $mat_{counts}$.

5. Count EPU mentions only if each column is larger than 0 for each row in $mat_{counts}$.

6. Perform monthly aggregation by multiplying $mat_{doc2month}$ and $mat_{counts}$.

7. Divide the $mat_{rawcounts}$ by the total number of articles in the same newspaper and month.

8. Standardize the new monthly series to unit standard deviation before a specific date and average across different newspapers.

9. Normalize the series to a mean of 100 before a specific date.

The $mat_{docfeat}$ and $mat_{doc2month}$ remain fixed for constructing different EPU indices, so the first and second steps are done before the genetic algorithm. The computation bottleneck would now be the two matrix multiplications. One is count aggregation, and the other is monthly aggregation. Since the matrices involved are sparse, we write an Rcpp function to perform the sparse matrix multiplication. To further speed up the computation, we run this process in parallel. However, the complied Rcpp function is not exportable to the parallel workers. So, we create a package for the Rcpp function and export it to the workers. In this way, we can evaluate 10,000,000 EPU indices in less than 10 hours with three cores running.

## 6 Conclusion

In this report, we present a new method to construct category-specific uncertainty indices through genetic algorithm. We show that our final VIX-specific EPU index captures well the VIX movements. Our approach is easily adapted to other target variables and can construct a frequency-based uncertainty index in a less costly and automatic manner. Only we need is a triple of initial keywords to form the candidate keywords set. However, we argue that not all economic variables are suitable for constructing text-based indices. The future work would be to explore the importance of candidate keywords and design a nowcasting/forecasting framework that may help decision-makers understand the category-specific investors' concerns.

# References

Algaba, A., Borms, S., Boudt, K., Pelt, J.V., 2020. The economic policy uncertainty index for Flanders, Wallonia and Belgium. BFW digitaal / RBF num´erique 6. doi:10.2139/ssrn.3580000

Ardia, D., Bluteau, K., Boudt, K., 2019. Questioning the news about economic growth: Sparse forecasting using thousands of news-based sentiment values. International Journal of Forecasting 35, 1370–1386. doi:10.2139/ssrn.2976084

Algaba, A., Ardia, D., Bluteau, K., Borms, S., Boudt, K., 2020. Econometrics Meets Sentiment: An Overview of Methodology and Applications. Journal of Economic Surveys, Vol. 34, Issue 3, pp. 512-547, 2020. doi:10.2139/ssrn.2652876

Ardia, D., Bluteau, Kassem, A., 2021. A century of Economic Policy Uncertainty through the French–Canadian lens. Economics Letters, Volume 205, 109938, ISSN 0165-1765. doi:/10.1016/j.econlet.2021.109938

Ardia, D., Bluteau, K., Borms, S., Boudt, K., 2021. The R Package sentometrics to Compute, Aggregate and Predict with Textual Sentiment. Journal of Statistical Software, Vol. 99, Issue 2, pp. 1-40, 2021. doi:10.2139/ssrn.3067734

Azqueta-Gavaldon, A., 2017. Developing news-based economic policy uncertainty index with unsupervised machine learning. Economics Letters 158, 47–50. doi:10.1016/j.econlet.2017.06.032.

Azqueta-Gavaldon, A, Hirschbühl, D., Onorante, L., Saiz, L., 2020. Economic Policy Uncertainty in the Euro Area: An Unsupervised Machine Learning Approach. doi:10.2139/ssrn.3516756.

Baker, S.R., Bloom, N., Davis, S.J., 2016. Measuring economic policy uncertainty. The Quarterly Journal of Economics 131, 1593–1636. doi:10.1093/qje/qjw024.

Friedman, J. H., Hastie, T., Tibshirani, R., 2010. Regularization Paths for Generalized Linear Models via Coordinate Descent. Journal of Statistical Software, 33(1), 1–22. doi:10.18637/jss.v033.i01

Ghirelli, C., Perez, J.J., Urtasun, A., 2019. A new economic policy uncertainty index for Spain. EconomicsLetters 182, 64–67. doi:10.1016/j.econlet.2019.05.021.

Manela, A., Moreira, A., 2017. News implied volatility and disaster concerns. Journal of Financial Economics 123, 137–162. doi:10.1016/j.jfineco.2016.01.032.

Jeffrey P., Richard S., Christopher D.M., 2014. GloVe: Global Vectors for Word Representation. doi:10.3115/v1/D14-1162

Scrucca, L., 2013. GA: A Package for Genetic Algorithms in R. Journal of Statistical Software, 53(4), 1–37. doi:10.18637/jss.v053.i04.

Scrucca, L., 2017. On some extensions to GA package: hybrid optimisation, parallelisation and islands evolution. The R Journal, 9(1), 187–206. doi:10.32614/RJ-2017-008

Tobback, E., Naudts, H., Daelemans w., Fortuny, E.J., Martens, D., 2018. Belgian economic policy uncertainty index: Improvement through text mining. International Journal of Forecasting,Volume 34, Issue 2, 2018, Pages 355-365, ISSN 0169-2070. doi:10.1016/j.ijforecast.2016.08.006.

# 7 Appendix

## 7.1 Alternative objective function

First, we do not introduce penalization, and it turns out that the noisy terms would be selected for some target variables. Therefore, we add a penalty on the number of keywords to let the algorithm select only the most discriminative keywords. In addition to the penalization, we still have choices for fitness between the EPU index and the objective variable. In this section, we discuss the fitness function.

In a simple linear regression, we have

$$y_i = \hat{y}_i + \epsilon_i = \hat{\alpha} + \hat{\beta} \times x_i + \epsilon_i$$

To minimize the sum of square residual errors,

$$\hat{\beta} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{S_{x,y}}{S_x^2} = r_{x,y} \frac{S_y}{S_x}$$

$$\hat{\alpha} = \bar{y} - \hat{\beta} \times \bar{x}$$

where $\bar{x}$ and $\bar{y}$ are the average of the $x_i$ and $y_i$ respectively, $r_{x,y}$ is the sample correlation coefficient, $S_x$ and $S_y$ are the sample standard deviation, $S_x^2$ is the sample variance, and $S_{x,y}$ is the sample covariance.

We get

$$\hat{y}_i = r_{x,y} \frac{S_y}{S_x}(x_i - \bar{x}) + \bar{y}$$

So,

$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n ((y_i - \bar{y}) - r_{x,y} \frac{S_y}{S_x}(x_i - \bar{x}))^2 = (n-1)S_y^2(1 - r^2)$$

Therefore,

$$MSE = \frac{1}{n}SSE = \frac{n-1}{n}S_y^2(1 - r^2)$$

where $n$ denotes the number of observations

From the above equation, we see that the MSE and correlation are related in a simple linear regression, but MSE also depends on the variance of the target variable Y. Since the variance of Y is unchanged when we change the subset of keywords, we can conclude that the MSE measure and the absolute correlation measure have similar measurement effects.

As for tracking errors, we know

$$TE = \sqrt{var(r_y - r_x)} = \sqrt{\frac{n}{n-1}(E[(r_y - r_x)^2] - E^2[r_y - r_x])}$$

where $\theta = E^2[r_y - r_x]$ captures the bias of the estimator.

It is hard to find the mathematical connection between the tracking error and the other two measures. The idea is that tracking error indicates how much a portfolio closely follows the index it is benchmarked against. In our context, tracking error measures the difference between the returns of two series. It does not require fitting $r_y$ with $r_x$, which is less flexible and may lead to underperformance.

## 7.2 Other results

Table 4 and Table 5 show the selected keywords for the unemployment rate and industrial production. Figure 3 and figure 4 show the comparison between constructed EPU index and the two target variables.

| EPU subcategory | selected keywords |
| --- | --- |
| Economy | apple, banana, fraction, normal, pen |
| Policy | book, green, yellow, girl, men |
| Uncertainty | cow, small, user, team, spider |

Table 3: Noisy keywords – arbitrary selected

| EPU subcategory | selected keywords |
| --- | --- |
| Economy | recession, recovery, inflation, boost |
| Policy | legislation, senate, government, measure |
| Uncertainty | concern, remain |

Table 4: Selected keywords for Unemployment rate - specific EPU index

| EPU subcategory | selected keywords |
| --- | --- |
| Economy | market |
| Policy | federal reserve |
| Uncertainty | remain |

Table 5: Selected keywords for Industrial Production - specific EPU index
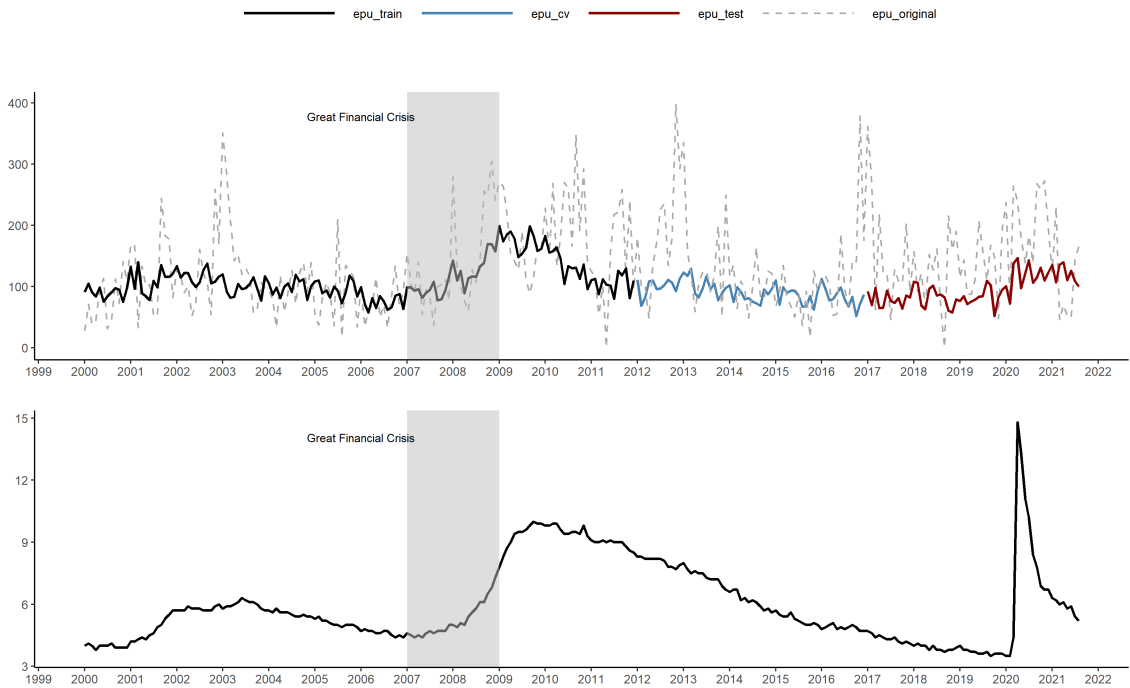


Figure 3: Comparison between the Unemployment Rate and Unemployment Rate-specific EPU Index
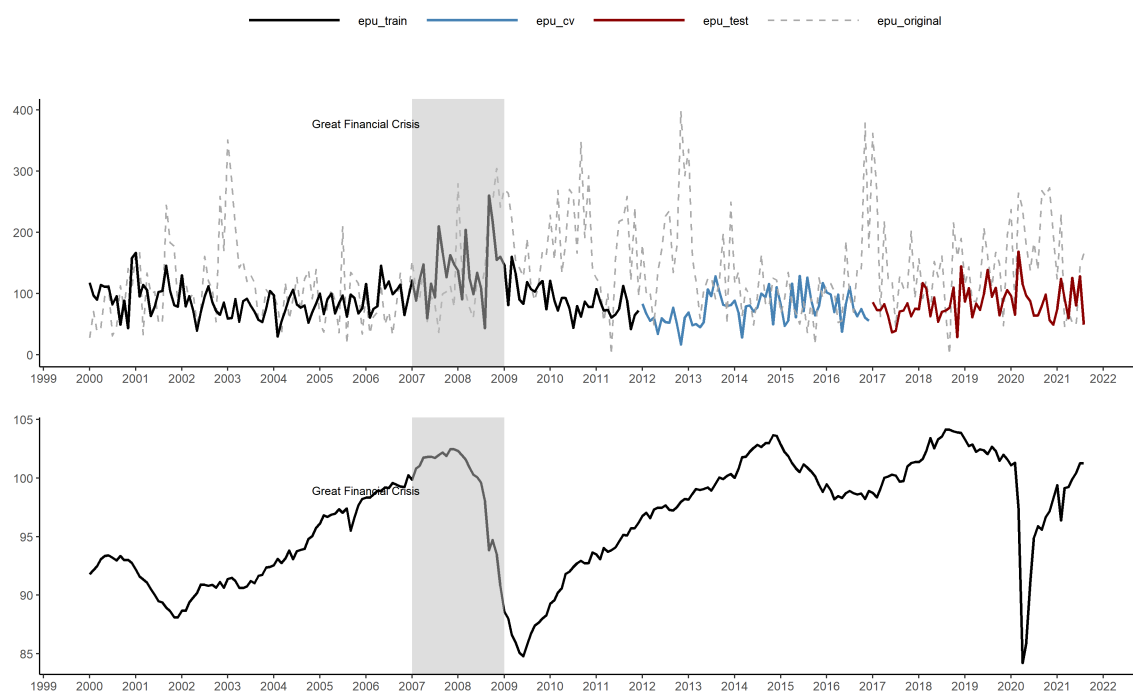
Figure 4: Comparison between the Industrial Production and the Industrial Production-specific EPU Index