

# An LSTM approach to Forecast Migration using Google Trends

Nicolas Golenvaux\*  
Pablo Gonzalez Alvarez\*  
UCLouvain, Belgium

Harold Silvère Kiossou  
UCLouvain, Belgium

Pierre Schaus  
UCLouvain, Belgium  
pierre.schaus@uclouvain.be

## ABSTRACT

Being able to model and forecast international migration as precisely as possible is crucial for policymaking. Recently Google Trends data in addition to other economic and demographic data have been shown to improve the forecasting quality of a gravity linear model for the one-year ahead forecasting. In this work, we replace the linear model with a long short-term memory (LSTM) approach and compare it with two existing approaches: the linear gravity model and an artificial neural network (ANN) model. Our LSTM approach combined with Google Trends data outperforms both these models on various metrics in the task of forecasting the one-year ahead incoming international migration to 35 Organization for Economic Co-operation and Development (OECD) countries: for example the root mean square error (RMSE) and the mean average error (MAE) have been divided by 5 and 4 on the test set. This positive result demonstrates that machine learning techniques constitute a serious alternative over traditional approaches for studying migration mechanisms.

## CCS CONCEPTS

• Applied computing → Forecasting; • Computing methodologies → Neural networks.

## KEYWORDS

forecasting, Google Trends, long short-term memory, migration, recurrent neural network

## ACM Reference Format:

Nicolas Golenvaux, Pablo Gonzalez Alvarez, Harold Silvère Kiossou, and Pierre Schaus. 2020. An LSTM approach to Forecast Migration using Google Trends. In *San Diego '20: KDD 2020 Conference Workshop ACM SIGKDD Conference on Knowledge Discovery and Data Mining August 24, 2020, San Diego, California USA*. ACM, New York, NY, USA, 8 pages. <https://doi.org/10.1145/1122445.1122456>

## 1 INTRODUCTION

Mobility has always been part of human history. In 2017, there were about 258 million international migrants worldwide, of which 150.3 million are migrant workers [36]. Modeling and forecasting human mobility is therefore important not only to help formulate

\*Both authors contributed equally to this research.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

*San Diego '20, August 24, 2020, San Diego, California USA*

© 2020 Association for Computing Machinery.  
ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00  
<https://doi.org/10.1145/1122445.1122456>

effective governance strategies but also to deliver insight at scale to humanitarian responders and policymakers. However, developing reliable methods able to forecast  $T_{i,j}$ , the number of people moving at the next time step from a given region  $i$  to another region  $j$  among  $m$  origin regions and  $n$  destination regions is extremely challenging due to the absence of, low frequency and long lags in recent migration data, especially for developing countries [4, 6].

One way to mitigate this lack of timely data is the use of real-time geo-referenced data on the internet like the Global Database of Events, Language, and Tone (GDELT Project) or Google Trends. Both have been successfully used to make forecasting in various fields [1, 7, 15]. Recently, Böhme et al. [6] demonstrated that adding geo-referenced online search data to forecast migration flows yields better performance compared to only using common economic and demographic indices like gross domestic product (GDP), and population size. The authors propose to forecast bilateral migration flows of the following year with a linear model relying on the Google Trends data captured the previous year.

In this work, we use the same data, but we replace the linear model by a recurrent neural network (LSTM [22]) that can consider the whole history to make forecasts. We demonstrate that the forecasting quality can be drastically improved by capturing better complex migration dynamics [26] and complex interactions between the many features.

The outline of our work is the following. We first introduce the related work in section 2. In order to make this article more self-contained, we explain in section 3 how the Google Trends features are extracted in Böhme et al. [6] and also briefly introduce recurrent neural networks. We then describe our recurrent neural network approach in section 4. Finally, our approach is evaluated and compared with the previous approaches in section 5.

## 2 RELATED WORK

In traditional models, the problem of forecasting  $T_{i,j}$  (i.e., to find  $\hat{T}_{i,j}$ ) is usually divided into two sub-problems: (a) forecast  $G_i$  the number of people leaving a region  $i$  (also known as the production function); and (b) forecast  $P_{i,j}$  the probability of a movement from  $i$  to  $j$ . Thus we have  $\hat{T}_{i,j} = G_i P_{i,j}$ .

There are two conventional models: (a) the gravity model; and (b) the radiation model. Gravity models, inspired by Newton's law, evaluate the probability of a movement between two regions to be proportional to the population size of the two regions  $i$  and  $j$ , and inversely proportional to the distance between them [2, 24, 29]. In radiation models, inspired by diffusion dynamics, a movement is emitted from a region  $i$  and has a certain probability of being absorbed by a neighboring region  $j$ . The subtlety here is that this probability is dependent on the population of origin, the population of the destination, and the population inside a circle centered in  $i$  with a radius equal to the distance from  $i$  to  $j$  [33]. Gravity is usually

**Table 1: The input features used for the different models. Each feature spans from 2004 to 2014 for a pair of origin-destination country. Refer to subsection 3.1 for a detailed explanation of the Google Trends Index (GTI).**

Input features $_{i,j,t}$	Description
$GDP_{i,t}$	Gross Domestic Product for origin country $i$ during the year $t$
$GDP_{j,t}$	Gross Domestic Product for destination country $j$ during the year $t$
$pop_{i,t}$	Population size for origin country $i$ during year $t$
$pop_{j,t}$	Population size for destination country $j$ during year $t$
$fixed_i$	Origin country $i$ fixed effects, encoded as a one-hot vector
$fixed_j$	Destination country $j$ fixed effects, encoded as a one-hot vector
$fixed_t$	Year $t$ fixed effects, encoded as a one-hot vector
$GTI_{bi,j,t}$	Bilateral GTI for a pair origin country $i$ and destination country $j$ during a year $t$
$GTI_{uni,i,t} \times GTI_{dest,i,j,t}$	Unilateral and destination GTI for an origin country $i$ , a destination country $j$ during a year $t$
$T_{i,j,t}$	Current year migration flow from country $i$ to country $j$

better to capture short distance mobility behaviors, while radiation is usually better to capture long-distance mobility behaviors [26].

With machine learning (ML) models, the approach is quite different as the goal is to directly forecast  $\hat{T}_{i,j} = f(\text{features})$  from a set of features<sup>1</sup>. To the best of our knowledge, [31] is the first attempt to use ML to forecast human migration. The authors use two ML techniques: (a) "extreme" gradient boosting regression (XGBoost) model; and (b) deep learning based artificial neural network (ANN) model. Similarly to us, this approach also attempts to directly forecast  $\hat{T}_{i,j}$  from the set of features without requiring any production function. But this approach exhibits two important differences with our approach: a) it uses traditional features for their forecasting model, which is composed of geographical and econometric properties such as the inter-country distance, and the median household income; and b) it does not capture the dynamic aspect since the forecast only relies on the previous time-step set of features.

More recently, Böhme et al. [6], use the Google Trends Index (GTI) of a set of keywords related to migration (e.g., visa, migrant, work) as a new feature set to make migration forecasts. Böhme et al. [6] rely on a bilateral gravity model to forecast the total number of migrant leaving a country of origin towards any of the OECD's destination countries during a specific year. The gravity models are estimated by a linear regression. Our approach uses the same input data and thus also relies on the GTI data. But instead of a linear least-squares estimation model, we propose a recurrent neural network (LSTM) that uses the complete set of historical features rather than only the ones coming from the previous time step.

### 3 BACKGROUND

In this section, first, we describe the data used for learning and forecasting the migration, giving more details about the Google Trends new set of features from [6]. Then, we present the performance metric used to compare the forecasting models also used in [31]. Finally, we briefly introduce recurrent neural networks.

### 3.1 Data and features sets

Table 1 gives an overview of the features used from the data provided by Böhme et al. [6]. More specifically, we use the following features: both GDP and population size for the origin and destination countries [37], the bilateral GTI, migration numbers from the previous year [28], as well as 3 one-hot vectors for encoding the origin, destination and the year.

*Google Trends Index features.* The Google Trends Index (GTI) is based on the Google Trends data freely accessible at [17]. The Google Trends tool allows collecting a daily measure of the relative quantities of web search of a precise keyword in a particular region of the world for a specified span of time<sup>2</sup> [12, 17]. To best represent the migration intentions of internet users via online searches, a set of terms related to the theme of migration is selected. It is composed of the 67 most semantically related terms to "immigration" and the 67 most semantically related terms to the word "economics" according to the website "Semantic Link"<sup>3</sup>. Every term is transcribed in 3 languages with Latin roots: English, Spanish and French to not complicate the extraction too much while covering a maximum of people, that is, about 841 million native speakers [10]. Tables 4 and 5 in the Appendix contain the set of main keywords [6].

The Google Trends Indexes of a precise keyword for a given country are then calculated by capturing the measures provided by Google Trends for the keyword in the geographical area corresponding to the country for the period spanning from 2004 to 2014<sup>4</sup>. Since the values provided by Google are provided as intervals of one month<sup>5</sup> and are normalized in a range between 0 and 100, the GTI are computed by taking the average of the values for each year to match the migration data. The indexes, therefore, reflect the variation of the number of searches for the keyword over the years.

The bilateral GTI data is made up of the two different forms of vectors:  $GTI_{bi,j,t}$  and  $GTI_{uni,i,t} \times GTI_{dest,i,j,t}$  for the unilateral and bilateral aspects. Three different forms of GTI values are then defined:

<sup>2</sup>The data can be downloaded from their website, or through an unofficial API [12].

<sup>3</sup><https://semantic-link.com/>

<sup>4</sup>Google Trends data only starts from 2004 and the migration data stops after 2015.

<sup>5</sup>This is specific to requests spanning from 2004 to the present.

<sup>1</sup>Notice that you could approach the problem the same way as with traditional models but it is not a common practice.

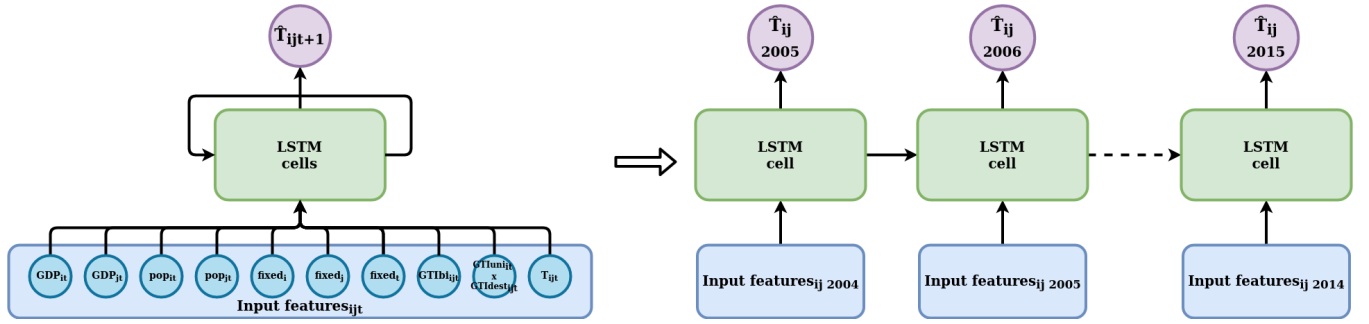


Figure 1: An unfolded gated-RNN with LSTM cells. The left-side corresponds to the folded RNN, while the right-side to the unfolded RNN.

- the vector of unilateral GTI ( $GTI_{uni,i,t}$ ) contains the GTI values of the set of keywords for the country of origin  $i$  during the year  $t$ ;
- the vector of bilateral GTI ( $GTI_{bi,i,j,t}$ ) contains GTI values also specific to the country of destination  $j$ . The values are still captured in the country of origin  $i$  during the year  $t$  but the related keywords correspond to the combination of the terms with the name of the destination country (e.g., visa Spain, migrant Spain, work Spain);
- the destination GTI ( $GTI_{dest,i,j,t}$ ) contains only the GTI value of the keyword corresponding to the destination country—name  $j$  (e.g., Spain) for the country  $i$  and the year  $t$ .

*Migration Data.* The OECD International Migration Database [28] provides a yearly incoming migratory flow from 101 countries of origin to the 35 countries members of the OECD from the early 1980's until 2015. Demographic and economic data about each destination and origin countries have been gathered from the World Development Indicators [37].

### 3.2 Evaluating forecasting models

The performance of forecasting models can be evaluated with several metrics. We present below the metrics used in [31]. Their formulas are summarized in Table 2.

**Common Part of Commuters (CPC):** Its value is 0 when the ground matrix  $T$  and the forecasting matrix  $\hat{T}$  have no entries in common, and 1 when they are identical.

**Mean Absolute Error (MAE):** Its value is 0 when the values of both matrices are identical, and arbitrarily positive the worse the forecasting gets.

**Root Mean Square Error (RMSE):** Its value is 0 when the values of both matrices are identical, and arbitrarily positive the worse the forecasting gets. The main difference with the MAE is that the RMSE penalizes more strongly the large errors.

**Coefficient of determination ( $r^2$ ):** Its value is 1 when the forecasts perfectly fit the ground truth values, 0 when the forecasts are identical to the expectation of the ground truth values, and arbitrarily negative the worse the fit gets.

**Mean Absolute Error In ( $MAE_{in}$ ):** The MAE on total incoming migrant by destination countries  $v_j = \sum_{i=1}^m T_{i,j}$ .

To make fair comparisons, for our experiments, we use these metrics. Remember that the focus of the forecasting is on incoming international migration to OECD countries.

### 3.3 Recurrent Neural Networks and Long Short-Term Memory (LSTM)

Recurrent neural networks (RNN) [16, 18]) are types of artificial neural networks (ANN) architectures particularly well suited to predict time-series or sequential data. It allows sharing features learned across different parts of the sequential data to persist through the network. It is not required to have a fixed set of input vectors. Long short-term memory (LSTM) [5, 9, 18, 21] are special architectures of RNN improving their ability to learn properly long-term dependencies by limiting the risk of vanishing and exploding gradient problems.

LSTM has recently gained momentum for several applications including in forecasting. It performs well in forecasting compared to other ML techniques [13, 14, 19, 23, 25, 32, 35, 38].

## 4 OUR LSTM APPROACH

Figure 1 shows the architecture of our RNN. We use one RNN in charge of forecasting the bilateral flows  $\hat{T}_{i,j,t+1}$  with the origin and destination countries, one hot encoded for all pairs  $(i, j)$ . The RNN has a unique LSTM layer. Another approach would have been to use different networks to estimate the flow for each pair of countries. The amount of data to train each would have been very limited, though. We use an additional densely-connected neural network layer on top of the gated-RNN layer, which computes the output scalar value  $\hat{T}_{i,j,t+1}$  given hidden vectors provided by the LSTM<sup>6</sup>.

### 4.1 Learning models and hyper-parameter optimization

To train the ML models we proceed using three sets [16]: (a) a train set, gathering the input features from 2004 to 2012 (input  $\_features_{i,j,04..12} \forall i, j$ ) and also all the observed migration flows

<sup>6</sup>It is common to use an activation function such as softmax with LSTM. However, in this work the output is scalar and we do not use an activation function for the densely connected neural network layer.

**Table 2: The different metrics –  $T$  is the ground truth value,  $\hat{T}$  is the forecasting matrix,  $m$  is the number of origin countries,  $n$  is the number of destination countries,  $v_j = \sum_{i=1}^m T_{i,j}$  is the number of incoming migrants for a zone  $j$ ,  $\hat{v}_j$  its forecast.**

Metrics	Equations
Common Part of Commuters	$CPC(T, \hat{T}) = \frac{2 \sum_{i,j=1}^{m,n} \min(T_{i,j}, \hat{T}_{i,j})}{\sum_{i,j=1}^{m,n} T_{i,j} + \sum_{i,j=1}^{m,n} \hat{T}_{i,j}} \quad (1)$
Mean Absolute Error	$MAE(T, \hat{T}) = \frac{1}{m \cdot n} \sum_{i,j=1}^{m,n}  T_{i,j} - \hat{T}_{i,j}  \quad (2)$
Root Mean Square Error	$RMSE(T, \hat{T}) = \sqrt{\frac{1}{m \cdot n} \sum_{i,j=1}^{m,n} (T_{i,j} - \hat{T}_{i,j})^2} \quad (3)$
Coefficient of determination	$r^2(T, \hat{T}) = 1 - \frac{\sum_{i,j=1}^{m,n} (T_{i,j} - \hat{T}_{i,j})^2}{\sum_{i,j=1}^{m,n} (T_{i,j} - \bar{T})^2} \quad (4)$
Mean Absolute Error In	$MAE_{in}(v, \hat{v}) = \frac{1}{n} \sum_j  v_j - \hat{v}_j  \quad (5)$

spanning from 2005 to 2013 as output ( $T_{i,j,05..13} \forall i, j$ ) since we forecast next year migration; (b) a validation set, containing input features on the year 2013 ( $input\_features_{i,j,13} \forall i, j$ ) and migration flows of 2014 ( $T_{i,j,14} \forall i, j$ ); and (c) a test set, on the year 2014 ( $input\_features_{i,j,14} \forall i, j$ ) and 2015 ( $T_{i,j,15} \forall i, j$ ). The hyperparameters are optimized accordingly for each model.

The data set is composed of 101 countries of origin, 34 countries of destination, 1997 time-series of length in the range from 2 to 11 for a total of 19 326 observations. Since the validation and the test sets each gather data for one of the 11 years available, each of these sets represents slightly less than 10% of the whole data.

A simplified version of our LSTM training is presented in Algorithm 1, while our LSTM evaluation is presented in Algorithm 2. Notice that the span of years presented in the algorithms corresponds to the one used once the validation is completed, that is, we fit our model on both the training and validation set.

---

#### Algorithm 1: Our training algorithm

---

**Data:** *model*: LSTM untrained model  
**Result:** Model is trained  
**for each epoch do**  
  **for each pair  $i, j$  of origin-destination countries do**  
    /\* gradient descent for each batch: \*/  
    *model.fit*(*input\_features* <sub>$i, j, 04..13$</sub> ,  $T_{i, j, 05..14}$ )  
  *evaluation(model)* /\* see algorithm 2 \*/

---

Due to the specificity of LSTM, we fit our LSTM time series by time series<sup>7</sup>. Therefore we use a batch of the size corresponding to the number of years present in the series. This implies that the gradient descent is applied and the LSTM's parameters are updated after each propagation of a time series through the LSTM cells (as

<sup>7</sup>By time series we mean the sequence of annual migration flows between a pair of origin-destination.

---

#### Algorithm 2: Our evaluation algorithm

---

**Data:** *model*: LSTM trained model  
**Result:** Model is evaluated  
**for each pair  $i, j$  of origin-destination countries do**  
   $\hat{T}_{i, j, 15} \leftarrow model.forecast(input\_features_{i, j, 04..14})$   
  error  $\leftarrow compute\_metrics(T_{15}, \hat{T}_{15})$

---

presented in Figure 1). Furthermore, the features have been normalized by time series of origin-destination using a min-max scaler [16]. Our LSTM model uses a bias of 1 for the LSTM forget gate since it has been shown to improve performances drastically [13, 23].

We make use of dropout regularization to reduce the overfitting for both models and ensure a better generalization [3, 11, 34]

For the experiments, we have used three different loss functions: (a) Mean Absolute Error (*MAE*); (b) Mean Square Error (*MSE*); and (c) Common Part of Commuters (*CPC*, as described in [31], that is,  $loss_{cpc} = 1 - cpc$  with *cpc* given by Equation (1)) and adapt them to handle time series.

We optimize the following hyperparameters and present them along with their optimal value: loss function - *MAE* using *Adam* optimizer, number and size of hidden layers - 1 layer of width 50, number of epochs - 50, and dropout - 0.15.

## 5 RESULTS AND DISCUSSION

We carry out experiments comparing the performance of our LSTM approach with two other models:

- (a) the bilateral gravity model estimated through an OLS model as presented in [6] whose gravity equation is represented

**Table 3: Comparison of the 3 models for the specified metrics. The values are shown by pair (train - test). Bold values indicate the best values per column. There are on average 742 migrants and 46 119 incoming migrants.**

Models	CPC		MAE		RMSE		$r^2$		MAE <sub>in</sub>	
	train	test	train	test	train	test	train	test	train	test
Gravity	0.871	0.866	819	877	6 100	5 239	0.800	0.773	24 128	28 737
ANN	<b>0.931</b>	0.834	119	306	818	1 553	0.975	0.921	3 257	9 664
LSTM	<b>0.945</b>	<b>0.892</b>	<b>96</b>	<b>225</b>	<b>639</b>	<b>1 028</b>	<b>0.985</b>	<b>0.967</b>	<b>2 261</b>	<b>4 827</b>

below:

$$\log T_{i,j,t+1} = \beta_1 GTIbi_{i,j,t} + \beta_2 GTIuni_{i,t} \times GTIdest_{i,j,t} + \beta_3 \log GDP_{i,t} + \beta_4 \log pop_{i,t} + \beta_5 \log GDP_{j,t} + \beta_6 \log pop_{j,t} + fixed_i + fixed_j + fixed_t + \epsilon_{i,j,t} \tag{6}$$

With  $\epsilon_{i,j,t}$  representing the robust error term;

- (b) a deep learning based artificial neural network model (ANN model) as proposed in [31]. The ANN is composed of densely connected with rectified linear units (ReLU) activation layers. We use the same model for all the forecasts with a time-step of 1 year. This means that the ANN receives as input the set of features  $input\ features_{i,j,t}$  described in Table 1 and outputs the forecasted next-year migration flow  $T_{i,j,t+1}$ . We optimize the following hyperparameters and present them along with their optimal value: loss function - MAE using Adam optimizer, number and width of hidden layers - 2 layers of width 200, training batch size - 32, number of epochs - 170 and dropout - 0.1.

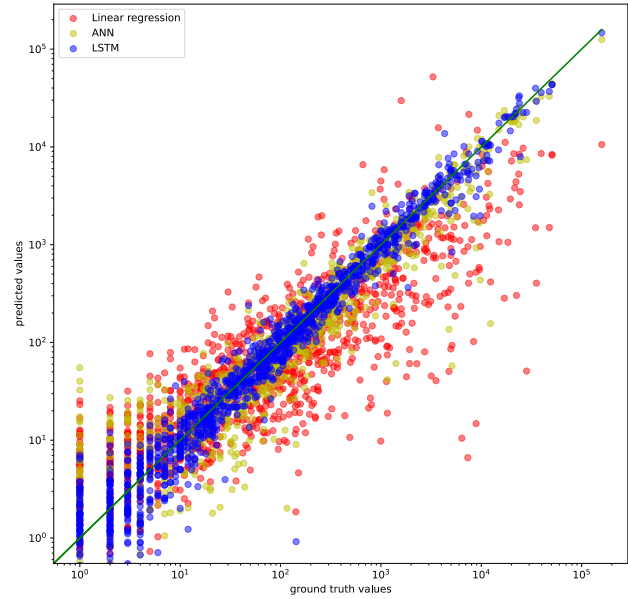
Our source code is available on the following git repository: <https://github.com/aia-uclouvain/gti-mig-paper>. It contains the script to extract the Google Trends Index, the Google Colab notebook to build the different models, as well as the data we used. The code is written in python and uses the Keras library, which runs on top of TensorFlow.

To assess the forecasting power of each model, we use a test set represented by every migration flow taking place in 2015 which represents a bit less than 10% of the whole data.

Table 3 shows the results of each model on both the training and test sets for the five metrics described in section 3.2.

Clearly, the ML models perform much better than the Böhme et al. [6]’s gravity model. Indeed, with the same data, the ANN is better than the first model in almost every metric while the LSTM model completely outperforms it in all the measures. The ANN model fits very well with the training data but it does not seem to generalize as well as the LSTM model as shown by their performance on the test set. On this data set, the LSTM is the best forecasting model among these three.

Note that, from Table 3 the RMSE values are always higher than the MAE (between 5 and 7 times larger). We can conclude that the models tend to make a few really large errors. This can be explained by analyzing the data. In the data set, the mean value of migration flows between 2 countries during a year is 742 but the median value is only 17 while the maximum is about 190 000. This indicates that our data set is very sparse: there is a lot of near-zero observations (40% are below 10) for a very few extremely important ones (less

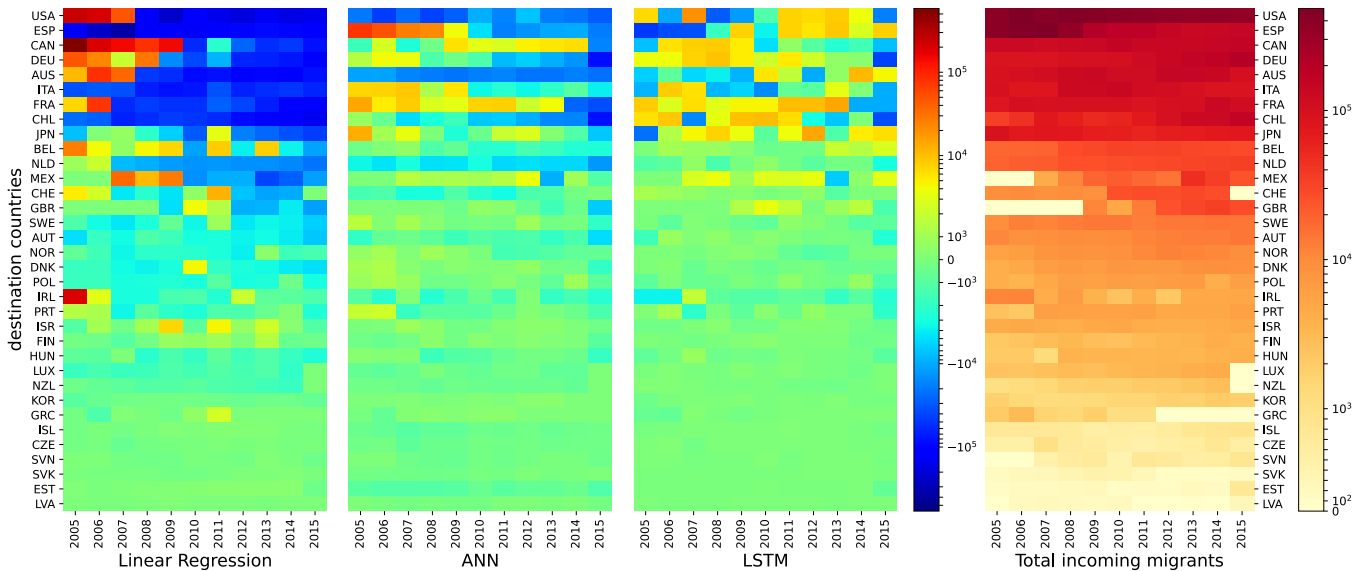


**Figure 2: Scatter plot for the 3 models on the test set (year 2015) – The coefficient of determination for the linear regression is 0.773, for the ANN 0.921, and for the LSTM 0.967 (see the Table 3 for more details).**

than 2% reach 10 000). Notice that the mean absolute errors of the different models are very important compared to the mean annual migration flows (742 and 46 119, see Table 3 caption) but these values are heavily biased by the sparsity of the data and by the large errors made on the really large migration flows (e.g., the USA and Spain).

To have a better visualization of the forecasting power of the models, we represent in Figure 2 the scatter plot of the 3 models for the test set only. The graph reflects well the sparse nature of the data as shown by the density of points along the x-axis. As expected following the first results, we can observe that the gravity model does not provide very accurate forecasts. The ANN model, on the other hand, shows a stronger tendency to underestimate the ground truth values. Ultimately, the LSTM’s estimations are the ones sticking the most to the actual migration flows which confirm our first assumption.

Finally, Figure 3 shows the error of the total number of incoming migrants per destination country per year for each model. We can observe that whatever the model, for the majority of countries and



**Figure 3: Heat maps of the error on total incoming migrants for 34 OECD countries on the test year (2015) showing how well each model fits the data. From left to right: Gravity Model, ANN Model, and LSTM Model. The rightmost figure is the ground values for the total number of incoming migrants by destination countries. Countries are in descending order of total incoming migrants.**

years, the estimation error is close to null and that the big errors often appear in the same countries of destination. Knowing that we can see that the heat maps of the ANN model and the linear regression in Figure 3 highlight their tendency to underestimate the migration flows especially for the last year (the test year).

To compare these errors with the actual migration flows, we represent in the rightmost heat map in Figure 3 the ground truth values of the total number of incoming migrants per destination country and per year in descending order. With this figure, we can clearly see that the errors we make are mostly for the countries with important incoming migration flow.

In the case of Spain notice that there has been an important drop in incoming migration flows in 2008 due to the 2007–2008 financial crisis [8]. If we look at the LSTM model in Figure 3, we largely underestimate the forecasts for Spain from 2005 to 2007. From 2008 and onwards, the forecast errors are comparatively smaller to those before that pivot year. This might indicate a lack of complexity of our model as it does not take into account major past events like a financial crisis.

Table 3, Figures 2 and 3 confirm our empirical results that the LSTM approach is able to forecast better than both a gravity model and an ANN model on the data set using different metrics.

## 6 CONCLUSION

Böhme et al. [6] have recently demonstrated that including Google Trends data in the set of standard features could improve the migration forecasting models. In this work, relying exactly on the same data, we improved the quality of the forecast significantly by replacing the gravity model used in by a Long short-term memory (LSTM) artificial recurrent neural network (RNN) architecture. Our

experiments also demonstrated that the LSTM was outperforming a standard ANN on this task.

A limitation of our testing procedure is we do not test the model on unknown pairs of origin-destination countries. Instead of splitting the training and the test on the years, we could split them according to the countries such that there is no country overlap. This could help us analyze the universality of our model and whether it generalizes properly on unknown countries or not.

Another possible improvement could be to compute  $\hat{T}_{i,j,t+1}$  not from the hidden vector of the last cell, but rather from an interpolation of the  $n$  previous cells. It might better reduce the impact of the abnormalities due to a specific year.

Moreover, the results have shown that the models lack some information to acknowledge important variations of the number of incoming migrants in some destination countries. Adding some factors like the presence of catastrophic events (financial crisis, war, natural disasters, epidemics), unemployment, and the share of internet users could significantly improve our approaches.

In this work, we used categorical labels in form of one-hot vector ( $fixed_i$ ,  $fixed_j$ , and  $fixed_t$ ). One could use a single one-hot vector  $fixed_{i,j}$  containing a pair origin-destination countries, absorbing more complex time-invariant factors like the distance between the 2 countries, and the presence of common language. ML literature has presented different techniques to handle such inputs [20, 30].

Finally, a drawback of our machine learning approach is that we lose the interpretability of the model and forecasts despite the high interpretability potential of Google search keywords.

As future work, we would like to apply the latest interpretability techniques (see [27]) to better identify the most important features for making high-quality migration forecasts. This would equip

economists, demographers, and experts in migration with new tools to shed light on migration mechanisms.

## ACKNOWLEDGMENTS

The authors acknowledge financial support from the UCLouvain ARC convention on "New approaches to understanding and modelling global migration trends" (convention 18/23-091).

## REFERENCES

- [1] Mohammed N Ahmed, Gianni Barlacchi, Stefano Braghin, Michele Ferretti, Vincent Lonij, Rahul Nair, Rana Novack, Jurij Paraszczak, and Andeep S Toor. 2016. A Multi-Scale Approach to Data-Driven Mass Migration Analysis. *SoGood@ ECML-PKDD* (2016), 17.
- [2] James E. Anderson. 2011. The Gravity Model. *Annual Review of Economics* 3, 1 (Sept. 2011), 133–160. <https://doi.org/10.1146/annurev-economics-111809-125114>
- [3] Devansh Arpit, Stanislaw Jastrzebski, Nicolas Ballas, David Krueger, Emmanuel Bengio, Maxinder S. Kanwal, Tegan Maharaj, Asja Fischer, Aaron Courville, Yoshua Bengio, and Simon Lacoste-Julien. 2017. A Closer Look at Memorization in Deep Networks. *arXiv:1706.05394 [cs, stat]* (July 2017). [arXiv:cs, stat/1706.05394](https://arxiv.org/abs/1706.05394)
- [4] Nikolaos Askitas and Klaus F. Zimmermann. 2015. The Internet as a Data Source for Advancement in Social Sciences. *International Journal of Manpower* 36, 1 (2015), 2–12.
- [5] Yoshua Bengio, Patrice Simard, and Paolo Frasconi. 1994. Learning Long-Term Dependencies with Gradient Descent Is Difficult. *IEEE transactions on neural networks* 5, 2 (1994), 157–166.
- [6] Marcus H. Böhme, André Gröger, and Tobias Stöhr. 2020. Searching for a Better Life: Predicting International Migration with Online Search Keywords. *Journal of Development Economics* 142 (Jan. 2020), 102347. <https://doi.org/10.1016/j.jdeveco.2019.04.002>
- [7] Hyunyoung Choi and Hal Varian. 2012. Predicting the Present with Google Trends. *Economic Record* 88, s1 (2012), 2–9. <https://doi.org/10.1111/j.1475-4932.2012.00809.x>
- [8] Andreu Domingo. 2017. El Sistema Migratorio Hispano-Americano del Siglo XXI México y España. *Revista de Ciencias y Humanidades - Fundación Ramón Areces* (Dec. 2017).
- [9] Kenji Doya. 1993. Bifurcations of Recurrent Neural Networks in Gradient Descent Learning. *IEEE Transactions on neural networks* 1, 75 (1993), 218.
- [10] David M. Eberhard, F. Simons Gary, and D. Fennig Charles. 2020. Ethnologue: Languages of the World. <https://www.ethnologue.com/>.
- [11] Yarin Gal and Zoubin Ghahramani. 2016. A Theoretically Grounded Application of Dropout in Recurrent Neural Networks. *arXiv:1512.05287 [stat]* (Oct. 2016). [arXiv:stat/1512.05287](https://arxiv.org/abs/1512.05287)
- [12] General Mills. 2019. GeneralMills/Pytrend. General Mills.
- [13] Felix A. Gers, Jürgen Schmidhuber, and Fred Cummins. 2000. Learning to Forget: Continual Prediction with LSTM. *Neural Computation* 12, 10 (Oct. 2000), 2451–2471. <https://doi.org/10.1162/089976600300015015>
- [14] C Lee Giles. 2001. Noisy Time Series Prediction Using Recurrent Neural Networks and Grammatical Inference. *Machine learning* 44, 1-2 (2001), 161–183.
- [15] Jeremy Ginsberg, Matthew H Mohebbi, Rajan S Patel, Lynnette Brammer, Mark S Smolinski, and Larry Brilliant. 2009. Detecting influenza epidemics using search engine query data. *Nature* 457, 7232 (2009), 1012–1014.
- [16] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. 2016. *Deep Learning*. MIT press.
- [17] Google. 2020. Google Trends. <https://www.google.com/trends>.
- [18] Alex Graves. 2012. *Supervised Sequence Labelling with Recurrent Neural Networks*. Studies in Computational Intelligence, Vol. 385. Springer Berlin Heidelberg, Berlin, Heidelberg. <https://doi.org/10.1007/978-3-642-24797-2>
- [19] Klaus Greff, Rupesh Kumar Srivastava, Jan Koutnik, Bas R. Steunebrink, and Jürgen Schmidhuber. 2017. LSTM: A Search Space Odyssey. *IEEE Transactions on Neural Networks and Learning Systems* 28, 10 (Oct. 2017), 2222–2232. <https://doi.org/10.1109/TNNLS.2016.2582924> [arXiv:1503.04069](https://arxiv.org/abs/1503.04069)
- [20] Cheng Guo and Felix Berkhahn. 2016. Entity Embeddings of Categorical Variables. *arXiv:1604.06737 [cs]* (April 2016). [arXiv:cs/1604.06737](https://arxiv.org/abs/1604.06737)
- [21] Sepp Hochreiter. 1991. Untersuchungen Zu Dynamischen Neuronalen Netzen. *Diploma, Technische Universität München* 91, 1 (1991).
- [22] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long Short-Term Memory. *Neural computation* 9, 8 (1997), 1735–1780.
- [23] Rafal Jozefowicz, Wojciech Zaremba, and Ilya Sutskever. 2015. An Empirical Exploration of Recurrent Network Architectures. In *International Conference on Machine Learning*. 2342–2350.
- [24] Emmanuel Letouzé, Mark Purser, Francisco Rodriguez, and Matthew Cummins. 2009. Revisiting the Migration-Development Nexus: A Gravity Model Approach. *Human Development Research Paper* 44 (2009).
- [25] Hao Liang, Meng Zhang, and Hailan Wang. 2019. A Neural Network Model for Wildfire Scale Prediction Using Meteorological Factors. *IEEE Access* 7 (2019), 176746–176755. <https://doi.org/10.1109/ACCESS.2019.2957837>
- [26] A. Paolo Masucci, Joan Serras, Anders Johansson, and Michael Batty. 2013. Gravity versus Radiation Models: On the Importance of Scale and Heterogeneity in Commuting Flows. *Physical Review E* 88, 2 (2013), 022812.
- [27] Christoph Molnar. 2019. *Interpretable machine learning*. Lulu. com.
- [28] OECD. 2020. International Migration Database. <https://www.oecd-ilibrary.org/content/data/data-00342-en>.
- [29] Jacques Poot, Omoniyi Alimi, Michael P Cameron, and David C Maré. 2016. The Gravity Model of Migration: The Successful Comeback of an Ageing Superstar in Regional Science. (2016), 27.
- [30] Kedar Potdar, Taher S., and Chinmay D. 2017. A Comparative Study of Categorical Variable Encoding Techniques for Neural Network Classifiers. *International Journal of Computer Applications* 175, 4 (Oct. 2017), 7–9. <https://doi.org/10.5120/ijca2017915495>
- [31] Caleb Robinson and Bistra Dilkina. 2018. A Machine Learning Approach to Modeling Human Migration. In *Proceedings of the 1st ACM SIGCAS Conference on Computing and Sustainable Societies*. 1–8.
- [32] Jürgen Schmidhuber, Daan Wierstra, and Faustino J. Gomez. 2005. Evolino: Hybrid Neuroevolution/Optimal Linear Search for Sequence Prediction. In *Proceedings of the 19th International Joint Conference on Artificial Intelligence (IJCAI)*.
- [33] Filippo Simini, Marta C. González, Amos Maritan, and Albert-László Barabási. 2012. A Universal Model for Mobility and Migration Patterns. *Nature* 484, 7392 (2012), 96–100.
- [34] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: A Simple Way to Prevent Neural Networks from Overfitting. *The Journal of Machine Learning Research* 15, 1 (Jan. 2014), 1929–1958.
- [35] Niek Tax, Ilya Verenich, Marcello La Rosa, and Marlon Dumas. 2017. Predictive Business Process Monitoring with LSTM Neural Networks. *arXiv:1612.02130 [cs, stat]* 10253 (2017), 477–492. [https://doi.org/10.1007/978-3-319-59536-8\\_30](https://doi.org/10.1007/978-3-319-59536-8_30) [arXiv:cs, stat/1612.02130](https://arxiv.org/abs/1612.02130)
- [36] United Nations Department of Economic and Social Affairs Population Division UN DESA. 2019. International Migrant Stock 2019. <https://www.un.org/en/development/desa/population/index.asp>.
- [37] World Bank. 2020. World Development Indicators. <https://datacatalog.worldbank.org/dataset/world-development-indicators>.
- [38] Kaisheng Yao, Trevor Cohn, Katerina Vylomova, Kevin Duh, and Chris Dyer. 2015. Depth-Gated LSTM. *arXiv:1508.03790 [cs]* (Aug. 2015). [arXiv:cs/1508.03790](https://arxiv.org/abs/1508.03790)

## A USED KEYWORDS

Tables 4 and 5 contain the set of main keywords: "For GTI data retrieval, both singular and plural as well as male and female forms of these keywords are used where applicable. In the English language, both British and American English spelling is used. All French and Spanish keywords were included with and without accents" [6, Table 1].



**Table 4: List of main keywords – First part [6, Table 1].**

English	French	Spanish
applicant	candidat	solicitante
arrival	arrivee	llegada
asylum	asile	asilo
benefit	allocation sociale	beneficio
border control	controle frontiere	control frontera
business	entreprise	negocio
citizenship	citoyennete	ciudadania
compensation	compensation	compensacion
consulate	consulat	consulado
contract	contrat	contrato
customs	douane	aduana
deportation	expulsion	deportacion
diaspora	diaspora	diaspora
discriminate	discriminer	discriminar
earning	revenu	ganancia
economy	economie	economia
embassy	ambassade	embajada
emigrant	emigre	emigrante
emigrate	emigrer	emigrar
emigration	emigration	emigracion
employer	employer	empleador
employment	emploi	empleo
foreigner	etranger	extranjero

**Table 5: List of main keywords – Second part [6, Table 1].**

English	French	Spanish
GDP	PIB	PIB
hiring	embauche	contratacion
illegal	illegal	ilegal
immigrant	immigre	inmigrante
immigrate	immigrer	inmigrar
immigration	immigration	inmigracion
income	revenu	ingreso
inflation	inflation	inflacion
internship	stage	pasantia
job	emploi	trabajo
labor	travail	mano de obra
layoff	licenciement	despido
legalization	regularisation	legalizacion
migrant	migrant	migrante
migrate	migrer	migrar
migration	migration	migracion
minimum	minimum	minimo
nationality	nationalite	nacionalidad
naturalization	naturalisation	naturalizacion
passport	pasport	pasaporte
payroll	paie	nomina
pension	retraite	pension
quota	quota	cuota
recession	recession	recesion
recruitment	recrutement	reclutamiento
refugee	refugie	refugiado
remuneration	remuneration	remuneracion
required documents	documents requis	documentos requisito
salary	salaire	sueldo
Schengen	Schengen	Schengen
smuggler	trafiquant	traficante
smuggling	trafic	contrabando
tax	tax	impuesto
tourist	touriste	turista
unauthorized	non autorisee	no autorizado
undocumented	sans papiers	indocumentado
unemployment	chomage	desempleo
union	syndicat	sindicato
unskilled	non qualifies	no capacitado
vacancy	poste vacante	vacante
visa	visa	visa
waiver	exemption	exencion
wage	salaire	salario
welfare	aide sociale	asistencia social