# Authorship Identification of the Speeches of US Presidents

Thong Bui, Jason Vantomme
*thongnbui@ischool.berkeley.edu, jvantomme@ischool.berkeley.edu*

Nirkhi and Dharaska (2013) define authorship identification (also known as authorship attribution) as a process that "determines the likelihood of a piece of writing to be produced by a particular author by examining other writings by that author." In fact, attribution of authorship has been a topic of inquiry since the late 1800's (Wikipedia, 2017) and, relevant to our context, has been applied to the decisions of the US Supreme Court5 as well as the pre-presidential radio addresses of Ronald Reagan (Airoldi et al., 2007). A common focus area of authorship identification in a present-day research is its use in criminal forensics. Furthermore, it is easy to understand in today's era of "fake news" how the verification of statements, policies or quotes emerges as a critical activity. In this project, we attempt to show how authorship identification of the speeches of US Presidents can be implemented using word and character-level models, with specific attention paid to the method described by Bagnall (2015).

## Data Sources

Primary sources for speech text from Presidents Hoover through Trump were obtained from The American Presidency Project @ UC Santa Barbara (https://goo.gl/rCJ8Gv) as well as from the Avalon Project at Yale (https://goo.gl/exPBFm). The majority of the data extraction comes from the JSON files downloaded from this website then cleaned (see: cleanup). Data preprocessing is then used to load the clean speech files and build the list of sentences for each president (see: preprocessing). Here are the statistics of the data of the 15 presidents' speeches:

```
How many speeches per president?          Approximately many words of text per president?
0  : Barack Obama            148           0  : Barack Obama            860977
1  : Donald J. Trump          22           1  : Donald J. Trump        101328
2  : Dwight D. Eisenhower    194           2  : Dwight D. Eisenhower   567639
3  : Franklin D. Roosevelt   224           3  : Franklin D. Roosevelt  323789
4  : George Bush              98           4  : George Bush            346032
5  : George W. Bush           56           5  : George W. Bush         323772
6  : Gerald R. Ford           40           6  : Gerald R. Ford         124695
7  : Harry S. Truman         302           7  : Harry S. Truman        368810
8  : Herbert Hoover          268           8  : Herbert Hoover         138249
9  : Jimmy Carter             60           9  : Jimmy Carter           224639
10 : John F. Kennedy          64           10 : John F. Kennedy        239684
11 : Lyndon B. Johnson       135           11 : Lyndon B. Johnson      417045
12 : Richard Nixon            41           12 : Richard Nixon          177545
13 : Ronald Reagan            48           13 : Ronald Reagan          184895
14 : William J. Clinton       64           14 : William J. Clinton     327027


How many sentences of text per president?  How many characters of text per president?
0  : Barack Obama          41919           0  : Barack Obama           4862045
1  : Donald J. Trump        8284           1  : Donald J. Trump         554978
2  : Dwight D. Eisenhower  24607           2  : Dwight D. Eisenhower   3084215
3  : Franklin D. Roosevelt 19223           3  : Franklin D. Roosevelt 1743457
4  : George Bush           21279           4  : George Bush           1896536
5  : George W. Bush        20326           5  : George W. Bush        1791014
6  : Gerald R. Ford         6053           6  : Gerald R. Ford         687272
7  : Harry S. Truman       29432           7  : Harry S. Truman       2001118
8  : Herbert Hoover         5899           8  : Herbert Hoover         786184
9  : Jimmy Carter          10752           9  : Jimmy Carter          1272152
10 : John F. Kennedy       10605           10 : John F. Kennedy       1344012
11 : Lyndon B. Johnson     23118           11 : Lyndon B. Johnson     2288606
12 : Richard Nixon          7357           12 : Richard Nixon          972275
13 : Ronald Reagan          9120           13 : Ronald Reagan         1010291
14 : William J. Clinton    16222           14 : William J. Clinton    1784166
```

**Design and Implementation**

In order to explore the ability to classify a given text to one of the US Presidents between Hoover and Trump, we implemented four models in Keras with a Tensorflow backend:

- Word-level RNN model
- Word-level LSTM model
- Character-level Simple RNN model
- Character-level, multi-headed softmax model

*Word-level Models*

The word-level RNN employed 1 embedding layer, 1 RNN hidden layer, 1 dropout layer and an output layer with softmax activation; optimization was done with the Adagrad algorithm with a loss function of categorical cross entropy. The word-level LSTM model employed an identical graph structure except that the RNN hidden layer was replaced with an LSTM layer. In both cases, for each input sequence there will be 15 prediction probabilities (one per president) produced by the softmax layer.
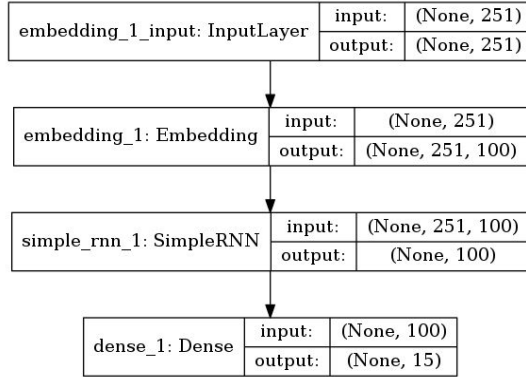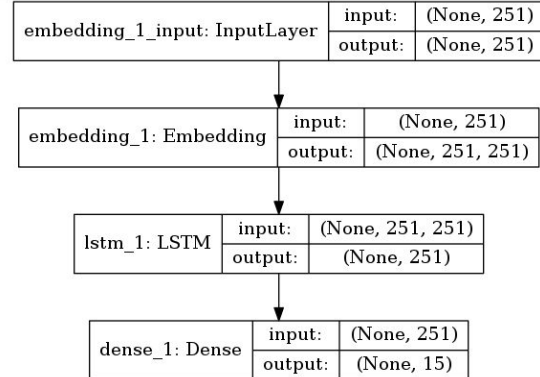
| embedding_1_input: InputLayer | input: | (None, 251) |
|---|---|---|
| | output: | (None, 251) |

| embedding_1: Embedding | input: | (None, 251) |
|---|---|---|
| | output: | (None, 251, 100) |

| simple_rnn_1: SimpleRNN | input: | (None, 251, 100) |
|---|---|---|
| | output: | (None, 100) |

| dense_1: Dense | input: | (None, 100) |
|---|---|---|
| | output: | (None, 15) |

*Fig. 1: Word-level RNN*

| embedding_1_input: InputLayer | input: | (None, 251) |
|---|---|---|
| | output: | (None, 251) |

| embedding_1: Embedding | input: | (None, 251) |
|---|---|---|
| | output: | (None, 251, 251) |

| lstm_1: LSTM | input: | (None, 251, 251) |
|---|---|---|
| | output: | (None, 251) |

| dense_1: Dense | input: | (None, 251) |
|---|---|---|
| | output: | (None, 15) |

*Fig. 2: Word-level LSTM*

For each of the president, the 80% of the sentences will be put in train dataset to be trained by the models, 20% will go into test dataset to test and predict which president a test sentence is from. A max length of sentence is determined and will be used as the input length for the model

During training of the word-level models, we have observed that large differences in both training and test loss and accuracy measures which indicated overfitting. After increasing the dropout from 0.2 to 0.5 the deltas between the training and tests sets decreased while increasing test accuracy; the character-level models also used a 0.5 drop-out rate.

For the RNN model, rates when broken down by president are generally proportionate with the size of the data for corresponding presidents even though class size was accounted for with class weighting; see Figure 3. However, there is a notable exception: Even with very low volume of sentences, accuracy for President Trump is much higher than presidents with more data such as Reagan, Carter, Kennedy, and Clinton. One might conclude that this is due to Trump's distinctive speech style even in speeches that are scripted.

For the word-level LSTM, the same patterns of accuracy are similar to the RNN model though accuracy rates generally exhibit an increase for most presidents, even the ones with smallest set of data such as Ford, Nixon, and Reagan; see Figure 4.
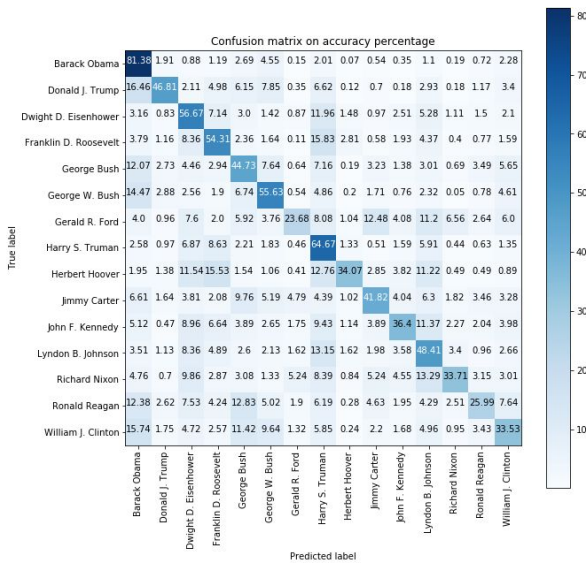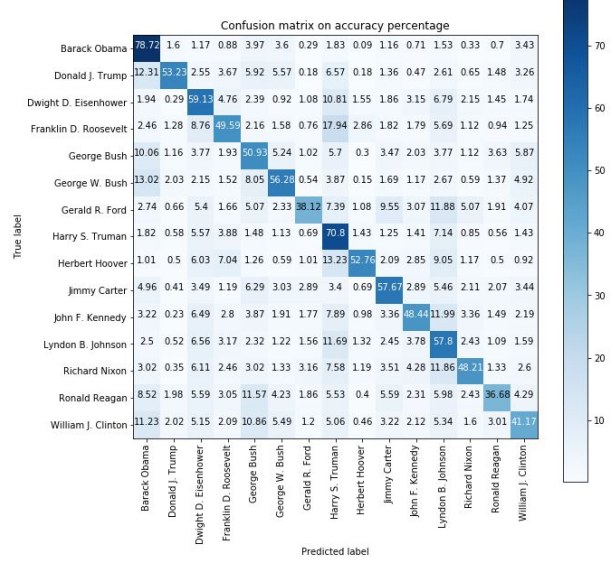
Fig. 3: Confusion Matrix for Word-level RNN



Fig. 4: Confusion Matrix for Word-level LSTM

Specific examples of the confusion in the word-level models are as follows:

| Predicted | Correct | Sentence Text |
|---|---|---|
| Bush | Nixon | i think it is a question of timeliness |
| Truman | Eisenhower | there was no offer made to the peruvians |
| Johnson | Roosevelt | they cannot get credit and private capital won't give them credit i am not saying improperly |
| Johnson | Carter | we are increasing every year the allocation of funds for that purpose |
| Carter | Ford | in order to get new york city to restore their credibility in the money markets they have taken these steps which have eliminated 3 95 billion cash deficit |
| Johnson | Truman | second the smaller amount of new obligational authority which i am recommending indicates the substantial portion of the financial requirements for our military buildup that has been met in the appropriations already made by the congress |
| Bush | Carter | and i would hope that from now on after this news conference that we could leave out references to allegations that anybody thinks that i'm a racist or that any of the other candidates in the race for president are racists |
| Obama | Trump | just dont worry about it |

*Character-level Models*

Similar to the word-level RNN, the character-level RNN employed 1 input layer, 1 RNN hidden layer, an output layer with softmax activation; optimization was also done with the Adagrad algorithm with a loss function of categorical cross entropy. The RNN hidden layer was wrapped in a function (Bidirectional) that ensure the RNN processes a sequence both forward and backward; this was shown to improve accuracy; see Figure 5.

Bagnall (2015) details an approach to the authorship attribution task at the CLEF 15 conference that, while simple in its notion, outperformed all other approaches in 3 of 4 languages tested (Stamatatos et al., 2015). In his approach, Bagnall employs a character-model RNN with a softmax output layer *per author*; the result: "Each softmax group is trained predominantly on one author's corpus, causing the recurrent layer to model a combination of all the author's texts, approximating the language as a whole. The output groups learn to weight the recurrent values in a way that best reflects their author's tendencies." (Bagnall 2015). This model proved very difficult to recreate in Keras due to the lack of dynamic graph execution which was critical to making the Bagnall (2015) method work.

Our attempt described in Figure 6 is clearly imperfect but allowed us to explore the possibility of other potential configurations as well as their inherent strengths, weaknesses and flaws. (We were able to explore other alternatives, but with no greater test accuracy.)
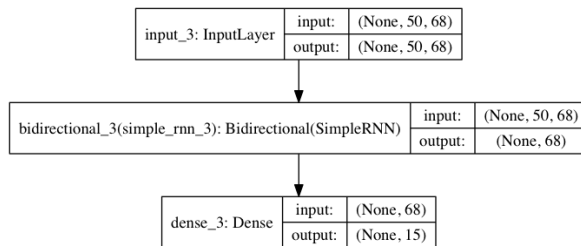


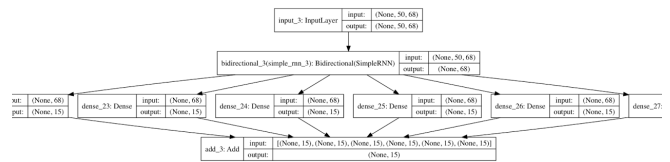*Fig. 5: Character-level Simple RNN*                    *Fig. 6: Character-level Multi-softmax RNN*

As can be seen in the confusion matrices for both the simple and multi-softmax models (Figs 7 & 8), the performance of these models is nowhere near the accuracy of the word-based models. It is notable that the most prominent "bands" of inaccurate labels in the multi-softmax case (e.g., Hoover in Fig. 8) does not obviously correlate to a president with the greatest number of characters. Eisenhower does show some over-attribution in the simple RNN model (Fig. 7) and is the president with the second-highest number of characters though class weighting should have accounted for volume differences.
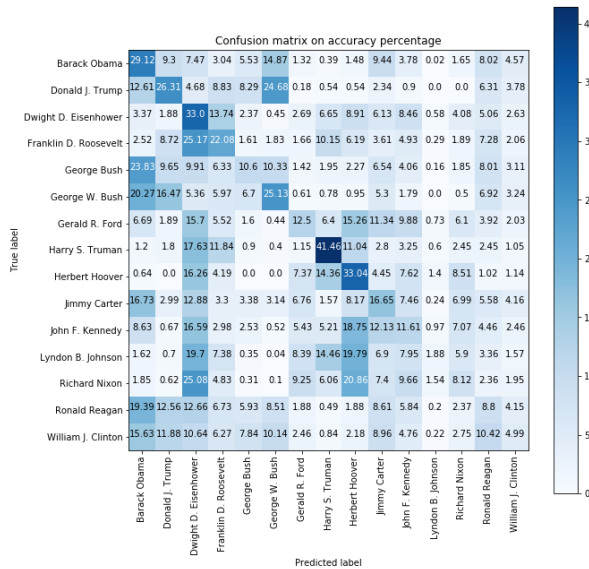


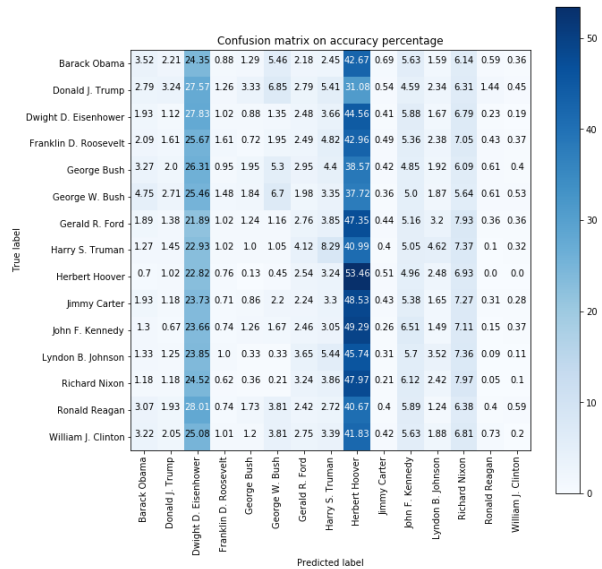Fig. 7: *Confusion Matrix for Character-level Simple RNN*    Fig. 8: *Confusion Matrix for Character-level Multi-softmax RNN*

Given the lacking performance on even a simple RNN, one might expect that there is something more fundamentally challenged in our basic assumptions regarding character-level NN set-up; this is an obvious area of exploration to consider.

**Comparison of Models**

As can be seen in the results in Figure 9, the two word-based models significantly outperformed character-based models in predicting classifications for our current data set; in particular, LSTM model is the best in identifying the authors of these president speeches. The relatively strong performance of the word-based models with remarkably little tuning was very encouraging.

| Model | Train loss | Test loss | Train accuracy | Test accuracy |
|---|---|---|---|---|
| Word-based RNN | 0.6221 | 1.8764 | 0.8165 | 0.5204 |
| Word-based LTSM | 0.7485 | 1.5055 | 0.7591 | 0.5831 |
| Char-based Simple RNN | 2.2708 | 2.3008 | 0.2190 | 0.2093 |
| Char-based Multi RNN | 2.6327 | 2.6901 | 0.1385 | 0.0804 |

*Figure 9: Comparison of models in the classification task*

**Conclusion**

We were challenged in implementing our target character-based approach described by Bagnall (2015), including using the ReSQRT activation function described in the paper which we were never able to utilize with stable results. Unfortunately, Bangall (2015) was light on details code available online was a bespoke solution written in C with little-to-no documentation.

Over the course of the project, we learned that simple character-based architectures may not be effective and, in fact, state-of-the-art character models are vastly more complicated and employ CNNs, highway networks, RNN+LSTM layers and other innovations and often many of these methods used together. These more complex models provide a path to continue exploring this problem.

Despite the challenges we encountered implementing the target character model, we were able to learn a significant amount about how to apply each of these methods for problems we will encounter in the future.

**Bibliography**

Airoldi, E., S. Fienberg, and K. Skinner. "Whose Ideas? Whose Words? Authorship of Ronald Reagan's Radio Addresses." *PS: Political Science & Politics* 40, no. 3 (n.d.): 501–6. doi:10.1017/S1049096507070874.

Bagnall, Douglas. "Author Identification Using Multi-Headed Recurrent Neural Network." In *CLEF 2015 Labs and Workshops, Notebook Papers; CEUR Workshop Proceedings*. September 8-11, 2015, 2015. http://ceur-ws.org/Vol-1391/150-CR.pdf.

———. "Authorship Clustering Using Multi-Headed Recurrent Neural Networks." In *CLEF 2016 Labs and Workshops, Notebook Papers; CEUR Workshop Proceedings*, 791–804. Évora, Portugal, 2016. http://ceur-ws.org/Vol-1609/16090791.pdf.

Graves, Alex. "Generating Sequences With Recurrent Neural Networks." *arXiv:1308.0850 [cs]*, August 4, 2013. http://arxiv.org/abs/1308.0850.

Nirkhi, Smita, and R. V. Dharaskar. "Comparative Study of Authorship Identification Techniques for Cyber Forensics Analysis." *arXiv:1401.6118 [cs]*, December 24, 2013. http://arxiv.org/abs/1401.6118.

Pascual, Santiago, and Antonio Bonafonte. "Multi-Output RNN-LSTM for Multiple Speaker Speech Synthesis and Adaptation." In *Signal Processing Conference (EUSIPCO), 2016 24th European*, 2325–29. IEEE, 2016. http://ieeexplore.ieee.org/abstract/document/7760664/.

Rosenthal, Jeffrey, and Albert H Yoon. "Detecting Multiple Authorship Of United States Supreme Court Legal Decisions Using Function Words." *The Annals of Applied Statistics* 5, no. 1 (2011): 283–308. doi:10.1214/10-AOAS378.

Stamatatos, Efstathios, Walter Daelemans, Ben Verhoeven, Patrick Juola, Aurelio López-López, Martin Potthast, and Benno Stein. "Overview of the Author Identification Task at PAN 2015." In *CLEF 2015 Labs and Workshops, Notebook Papers; CEUR Workshop Proceedings*. Toulouse, France, 2015. http://www.uni-weimar.de/medien/webis/events/pan-15/pan15-papers-final/pan15-authorship-verification/stamatatos15-overview.pdf.

"Stylometry." *Wikipedia*, April 27, 2017. https://en.wikipedia.org/w/index.php?title=Stylometry&oldid=777553870.