

# Stylometric and Neural Features Combined Deep Bayesian Classifier for Authorship Verification

Yitao Sun, Svetlana Afanaseva and Kailash Patil

[ysun@pindrop.com](mailto:ysun@pindrop.com), [safanaseva@pindrop.com](mailto:safanaseva@pindrop.com)





## Outline

1. Idea
2. Model overview
3. Results
4. On-going work
5. Q&A

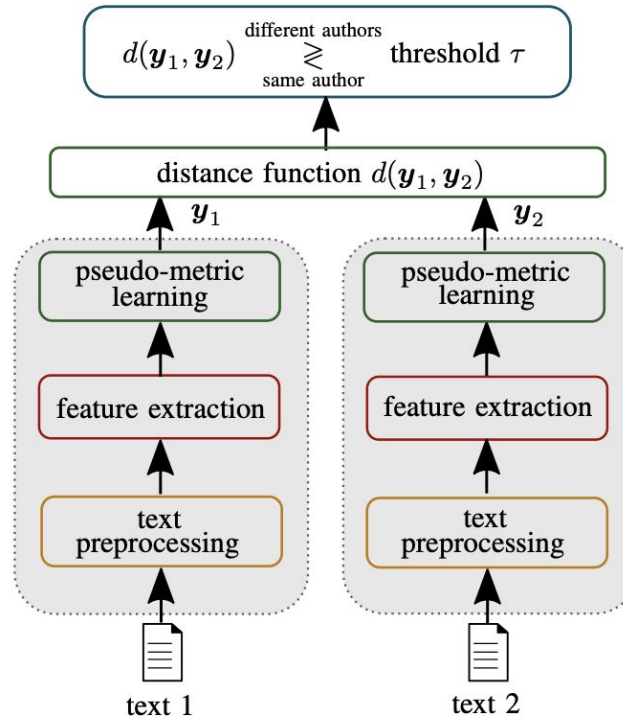
# Ideas

- 1) Siamese Network Adhominem system (B. Boenninghoff, 2019) lacks stylometric feature
- 2) To improve the Adhominem system:
  - a) Add Stylometric features (J. Weerasinghe, 2020), e.g frequency of function words, vocab richness, to enrich input data
  - b) Use probabilistic linear discriminant analysis (PLDA) as metric layer



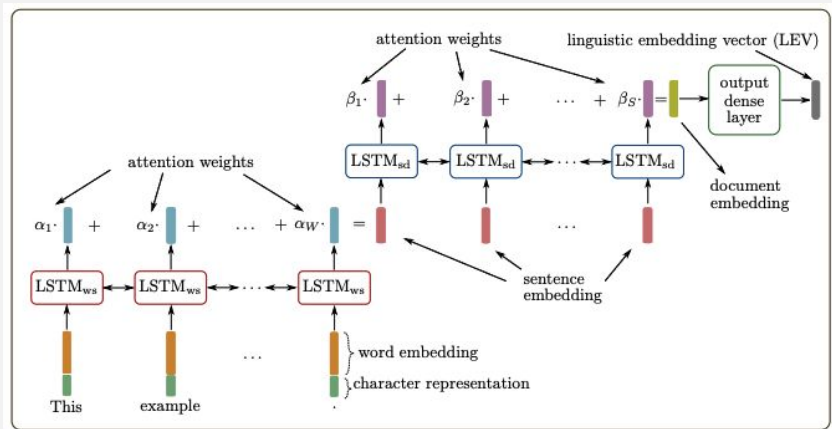
# Model Overview

Siamese Network

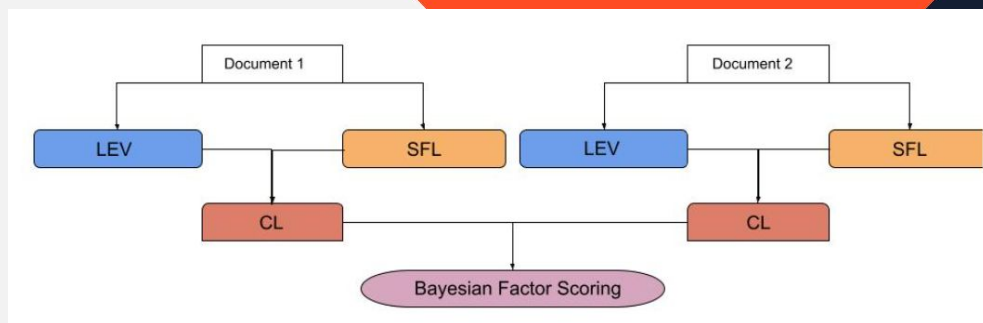


# Model Overview

## LEV (linguistic embedding vectors) in detail



## SFL (Stylometric Feature Layers)



# Bayesian Factor Scoring(B. Boenninghoff 2020)

$$\begin{array}{ccccc} \underline{y} & = & \underline{x} & + & \underline{\epsilon} \\ \text{combined layers} & & \text{author's writing style} & & \text{noise term} \end{array}$$

- $\mathcal{H}_s$ : The two documents were written by the same person,
- $\mathcal{H}_d$ : The two documents were written by two different persons.

## Same-author pair probability

$$p(y_1, y_2 | \mathcal{H}_s) = \frac{p(y_1, y_2 | x_0, \mathcal{H}_s) p(x_0 | \mathcal{H}_s)}{p(x_0 | y_1, y_2, \mathcal{H}_s)} = \frac{p(y_1 | x_0) p(y_2 | x_0) p(x_0)}{p(x_0 | y_1, y_2)}$$

## Different-author pair probability

$$p(y_1, y_2 | \mathcal{H}_d) = p(y_1 | \mathcal{H}_d) p(y_2 | \mathcal{H}_d) = \frac{p(y_1 | x_1) p(x_1)}{p(x_1 | y_1)}, \frac{p(y_2 | x_2) p(x_2)}{p(x_2 | y_2)}$$



# Bayesian Factor Scoring

## Final bayesian score generation

$$\begin{aligned}\text{score}(\mathbf{y}_1, \mathbf{y}_2) &= \log p(\mathbf{y}_1, \mathbf{y}_2 \mid \mathcal{H}_s) - \log p(\mathbf{y}_1, \mathbf{y}_2 \mid \mathcal{H}_d) \\ &= \log p(\mathbf{x}_0) - \log p(\mathbf{x}_1) - \log p(\mathbf{x}_2) \\ &\quad + \log p(\mathbf{y}_1 \mid \mathbf{x}_0) + \log p(\mathbf{y}_2 \mid \mathbf{x}_0) - \log p(\mathbf{y}_1 \mid \mathbf{x}_1) - \log p(\mathbf{y}_2 \mid \mathbf{x}_2) \\ &\quad - \log p(\mathbf{x}_0 \mid \mathbf{y}_1, \mathbf{y}_2) + \log p(\mathbf{x}_1 \mid \mathbf{y}_1) + \log p(\mathbf{x}_2 \mid \mathbf{y}_2)\end{aligned}$$



# Bayesian Factor Scoring

To learn the probability layer from the data

$$p(\mathcal{H}_s | \mathbf{y}_1, \mathbf{y}_2) = \frac{p(\mathbf{y}_1, \mathbf{y}_2 | \mathcal{H}_s)}{p(\mathbf{y}_1, \mathbf{y}_2 | \mathcal{H}_s) + p(\mathbf{y}_1, \mathbf{y}_2 | \mathcal{H}_d)} = \text{Sigmoid}(\text{score}(\mathbf{y}_1, \mathbf{y}_2))$$

The loss function

$$\mathcal{L}_\phi = l \cdot \log \{p(\mathcal{H}_s | \mathbf{y}_1, \mathbf{y}_2)\} + (1 - l) \cdot \log \{1 - p(\mathcal{H}_s | \mathbf{y}_1, \mathbf{y}_2)\}$$





# Results

Test Results of PAN 2023 Training Dataset

<i>Model</i>	<i>AUC</i>	<i>C@1</i>	<i>f<sub>0</sub>5<sub>u</sub></i>	<i>F1</i>	<i>brier</i>	<i>overall</i>
Naive, Distance-based	0.493	0.497	0.553	0.664	0.741	0.589
Method-based text compression	0.504	0.033	0.048	0.621	0.75	0.391
DML without SFL	0.503	0.523	0.492	0.357	0.603	0.495
UAL without SFL	0.499	0.52	0.477	0.336	0.593	0.485
BFS without SFL	0.47	0.502	0.474	0.37	0.597	0.483
DML with SFL	0.523	0.499	0.605	0.522	0.73	0.576
UAL with SFL	0.568	0.492	0.584	0.467	0.747	0.571
BFS with SFL	<b>0.658</b>	<b>0.662</b>	<b>0.739</b>	<b>0.735</b>	<b>0.762</b>	<b>0.711</b>



## On-going work and thoughts

1. Didn't finetune the threshold, can do better next year
2. Might consider EER(Equal Error Rate) as one of the metric
3. Does online text(word) distribute the same way as transcribed text
4. Use other popular pretrained transformer with BFS(Bayesian Factor Scoring) to boost performance



# Q&A

