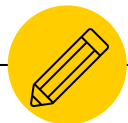


# Who wrote it?

Speech Technologies Project, Spring 2016



**Lucas Rodés**

**Sergi Liesegang**

# Outline

---

- Introduction
- Corpus
- Classifiers
- Results
- Conclusions. Future Work.

---

1

# Introduction

---

...propiedad, descomulgando  
quien mal lo quiere, y yo que  
satisfecho y ufano de haber  
primero que gozé el fruto de  
enteramente, como deseaba, y  
sido otro mi deseo que por  
abhorrecimiento de los hon  
fingidas y disparatada  
los libros de caballería  
mi verdaderamente don Qu

Spanish Literature





A detailed oil painting of William Shakespeare's face and upper torso. He has a full, greyish-brown beard and mustache, and is wearing a white ruffled collar. The background is dark and textured.

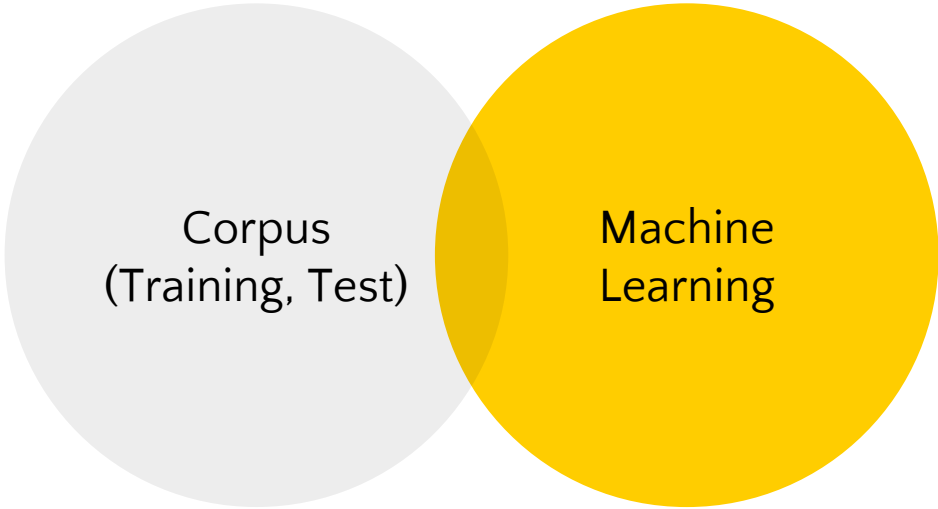
## Author Classification



1

## Scope

---



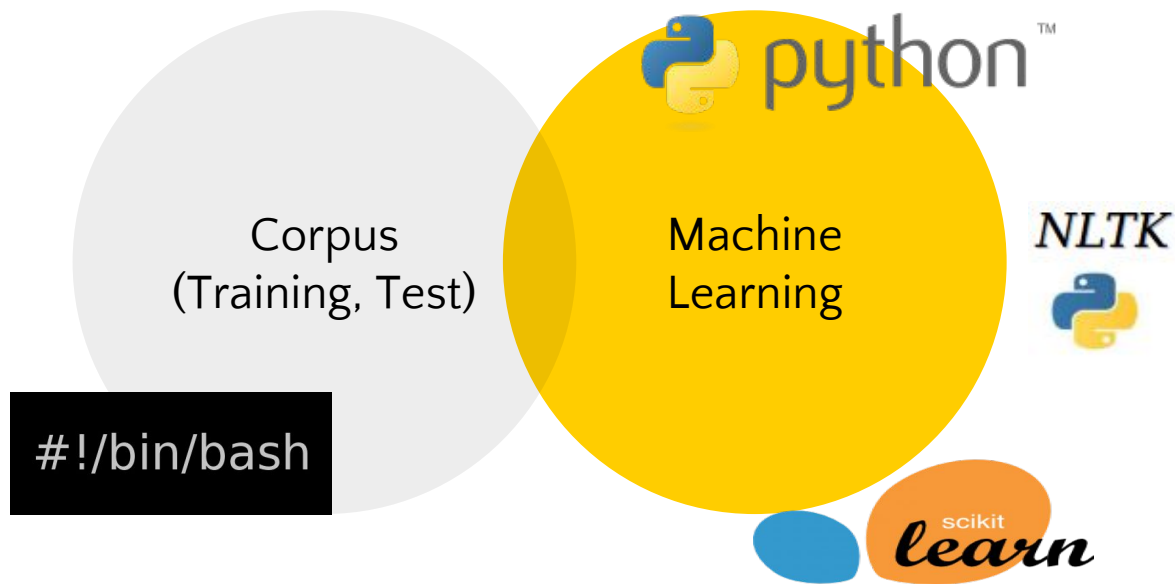
Corpus  
(Training, Test)

A Venn diagram consisting of two overlapping circles. The left circle is light gray and contains the text 'Corpus (Training, Test)'. The right circle is yellow and contains the text 'Machine Learning'. The overlapping area in the center is a darker shade of yellow.

Machine  
Learning

1

## Tools



---

2

# Corpus


---



<b>Author</b>	<b>Epoch</b>	<b>Genre</b>
Miguel de <b>Cervantes</b>	Renaissance (C16-17th)	Novel - Poetry - Drama
P. <b>Calderón de la Barca</b>	Golden Age - Baroque (C17th)	Drama - Poetry
Felix <b>Lope de Vega</b>	Golden Age - Baroque (C17th)	Novel - Poetry - Drama
Francisco de <b>Quevedo</b>	Baroque (C17th)	Poetry
G. Adolfo <b>Becquer</b>	Romanticism (C19th)	Poetry - Novel
Emilia <b>Pardo Bazan</b>	Realism - Naturalism (C19th)	Novel
Benito Pérez <b>Galdós</b>	Realism (C19th)	Novel - Drama
Federico <b>Garcia Lorca</b>	Generation of 27 (C20th)	Poetry - Drama

2

## Database

Correlated 

Author	Epoch	Genre
Miguel de <b>Cervantes</b>	Renaissance (C16-17th)	Novel - Poetry - Drama
<b>P. Calderón de la Barca</b>	Golden Age - Baroque (C17th)	Drama - Poetry
<b>Felix Lope de Vega</b>	Golden Age - Baroque (C17th)	Novel - Poetry - Drama
Francisco de <b>Quevedo</b>	Baroque (C17th)	Poetry
G. Adolfo <b>Becquer</b>	Romanticism (C19th)	Poetry - Novel
Emilia <b>Pardo</b> Bazan	Realism - Naturalism (C19th)	Novel
Benito Pérez <b>Galdós</b>	Realism (C19th)	Novel - Drama
Federico <b>Garcia Lorca</b>	Generation of 27 (C20th)	Poetry - Drama

---

3

# Classifiers

---

3

## Naive Bayes

**Mainstay** in text Processing

**Simple** Categorization

Based on **Bag of Words** model

Compares elements in texts

3

## Naive Bayes

**Mainstay** in text Processing

**Simple** Categorization

Based on **Bag of Words** model

Compares elements in texts

**Sparse Matrix** as input feature

Word occurrences



3

## Naive Bayes

Mainstay in text Processing    Simple Categorization

Based on Bag of Words model

Compares elements in texts

Sparse Matrix as input feature

Word occurrences

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

3

## Naive Bayes

**Mainstay** in text Processing    **Simple** Categorization

Based on **Bag of Words** model  
Compares elements in texts

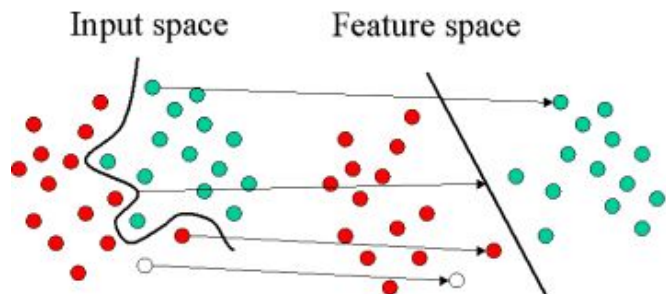
**Sparse Matrix** as input feature  
Word occurrences

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Implemented in Python  
Scikit-learn Library

3

## Support Vector Machine



**Well-known classifier**

Finds boundary with maximum class separability.

**Kernel**

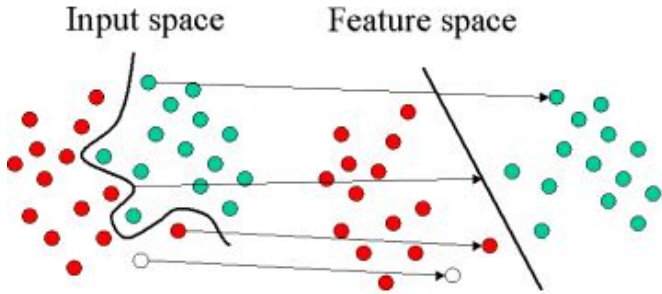
We used **Linear**

**Sparse Matrix** as input feature

Word occurrences

3

## Support Vector Machine



### Well-known classifier

Finds boundary with maximum class separability.

### Kernel

We used **Linear**

### Sparse Matrix as input feature

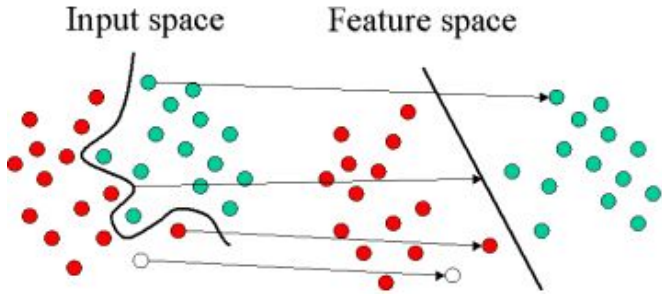
Word occurrences

### Grid Search

Optimization method to find best classifier parameters

3

## Support Vector Machine



### Well-known classifier

Finds boundary with maximum class separability.

### Kernel

We used **Linear**

### Sparse Matrix as input feature

Word occurrences

### Grid Search

Optimization method to find best classifier parameters

Implemented in Python

Scikit-learn Library



3

## Approaches

- Download books

Author 1



...

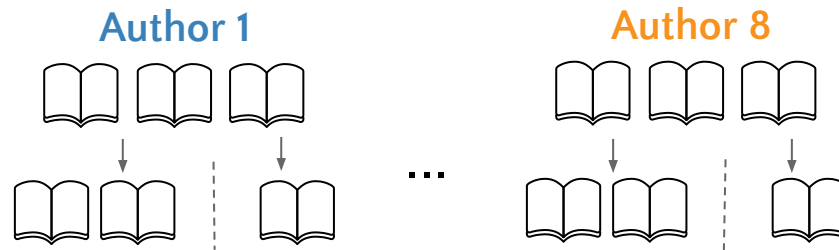
Author 8



3

## Approaches

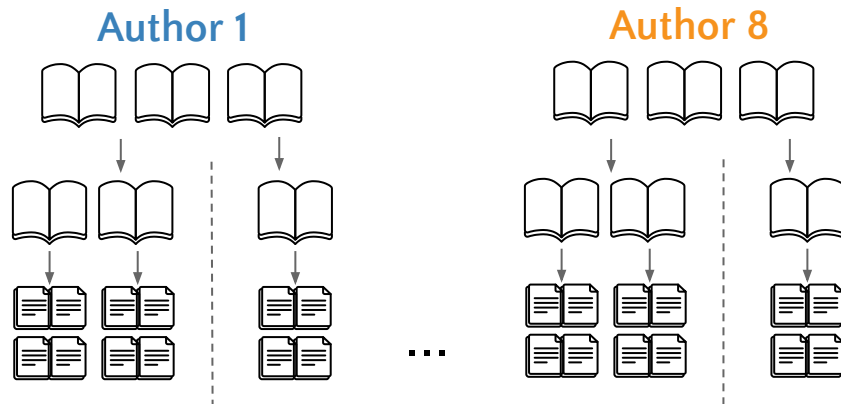
- **Download books**
- 66% Training, 33% Test



3

## Approaches

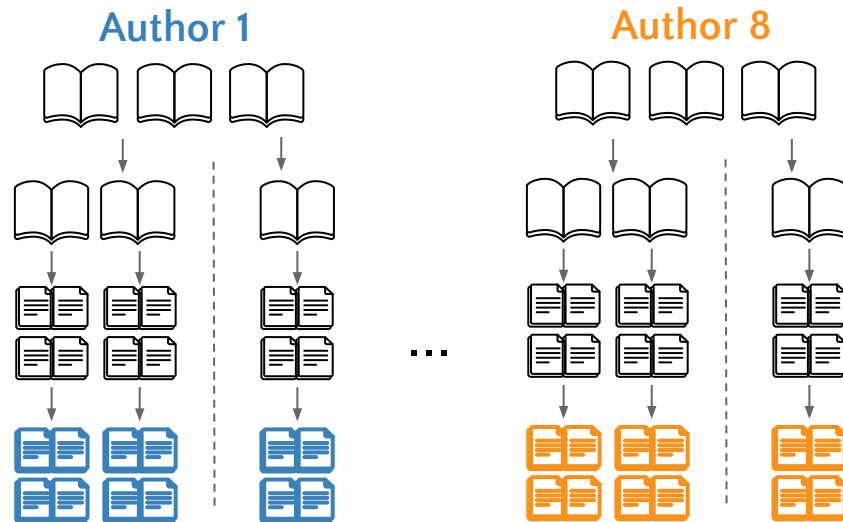
- Download books
- 66% Training, 33% Test
- Split in chunks of  $L$  lines



3

## Approaches

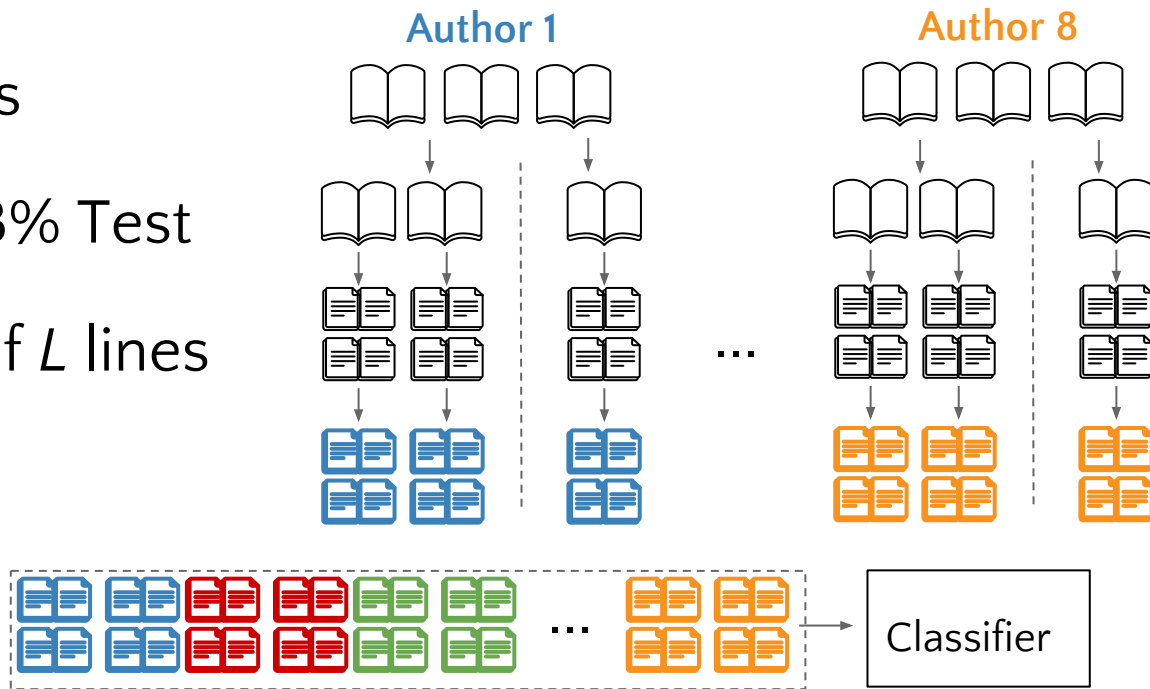
- Download books
- 66% Training, 33% Test
- Split** in chunks of  $L$  lines
- Labelling



3

## Approaches

- Download books
- 66% Training, 33% Test
- Split in chunks of  $L$  lines
- Labelling
- Train Classifier

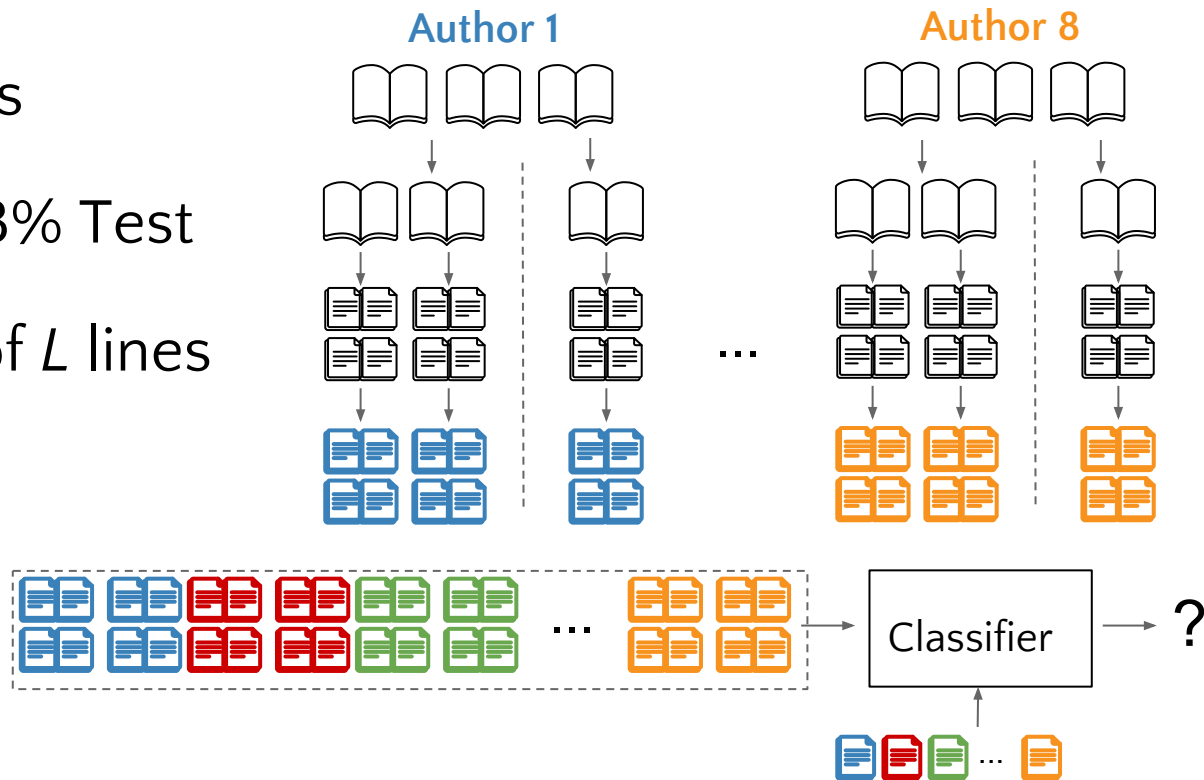




3

## Approaches

- Download books
- 66% Training, 33% Test
- Split in chunks of  $L$  lines
- Labelling
- Train Classifier
- Test Classifier



3

## Natural Language Modeling

Popular in speech processing

Based on Statistical Learning

Word statistics

3

## Natural Language Modeling

Popular in speech processing

Based on Statistical Learning

Word statistics

### N-gram Model

Unigram, Bigram,  
Trigram.

3

## Natural Language Modeling

Popular in speech processing

Based on Statistical Learning

Word statistics

**N-gram Model**

Unigram, Bigram,  
Trigram.

**Smoothing Methods**

We have used Add one Smoothing

3

## Natural Language Modeling

Popular in **speech processing**

Based on **Statistical Learning**

Word statistics

**N-gram Model**

Unigram, Bigram,  
Trigram.

**Smoothing Methods**

We have used Add one Smoothing

Implemented in Python

NLTK Library



3

## Approaches

- Download books

Author 1



...

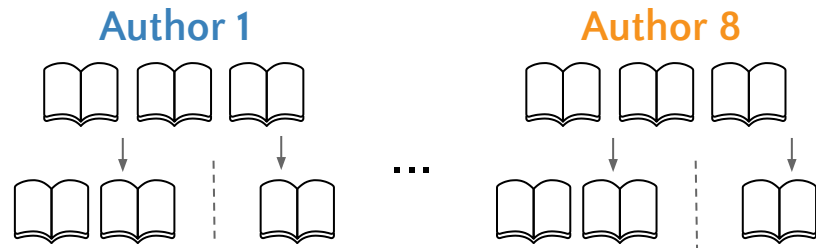
Author 8



3

## Approaches

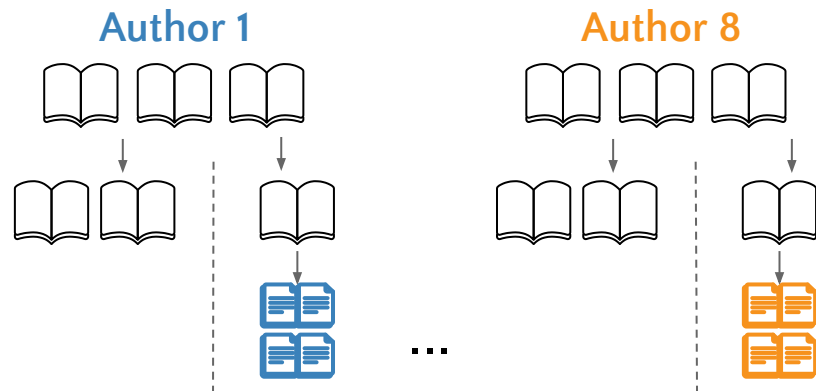
- Download books
- 66% Training, 33% Test



3

## Approaches

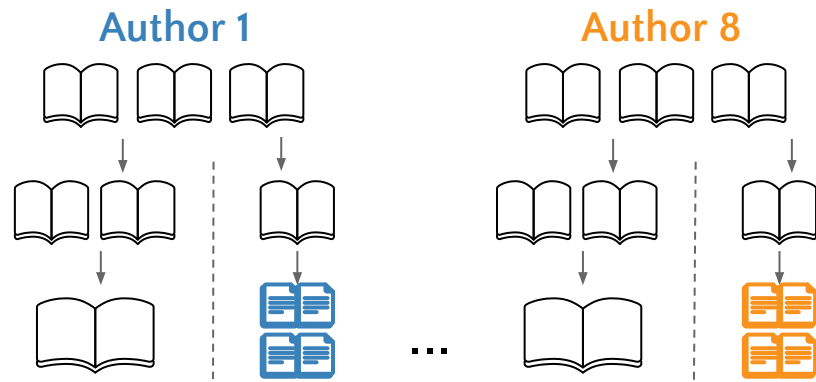
- Download books
- 66% Training, 33% Test
- Split Test in chunks & label



3

## Approaches

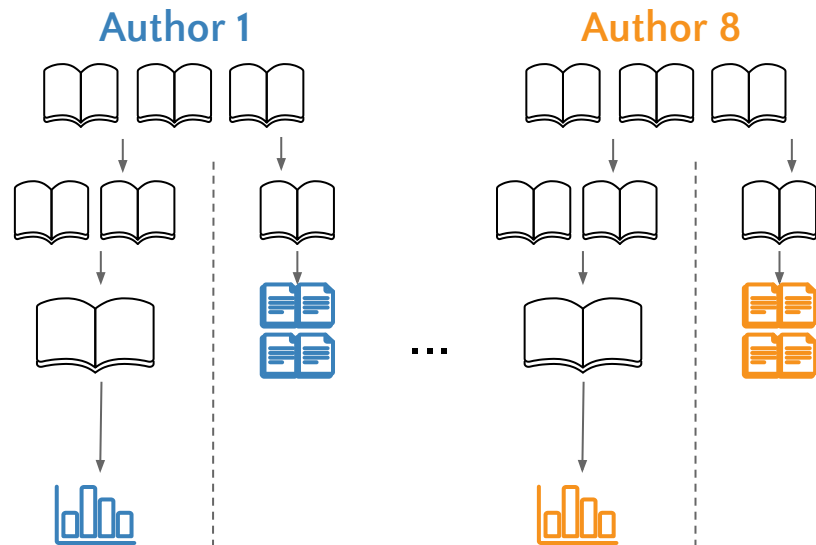
- Download books
- 66% Training, 33% Test
- Split Test in chunks & label
- Merge Training books



3

## Approaches

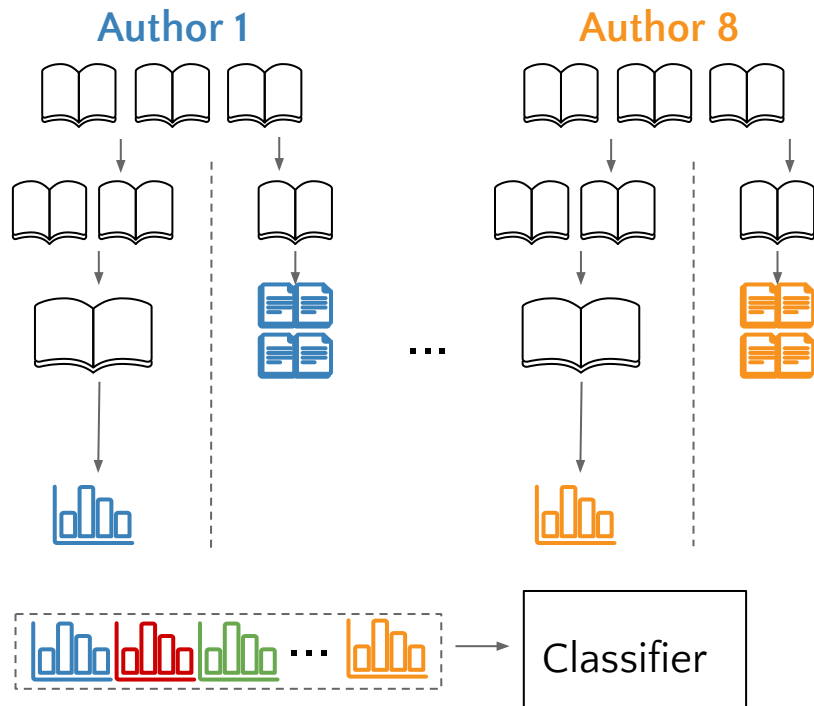
- Download books
- 66% Training, 33% Test
- Split Test in chunks & label
- Merge Training books
- Obtain **Language Model**



3

## Approaches

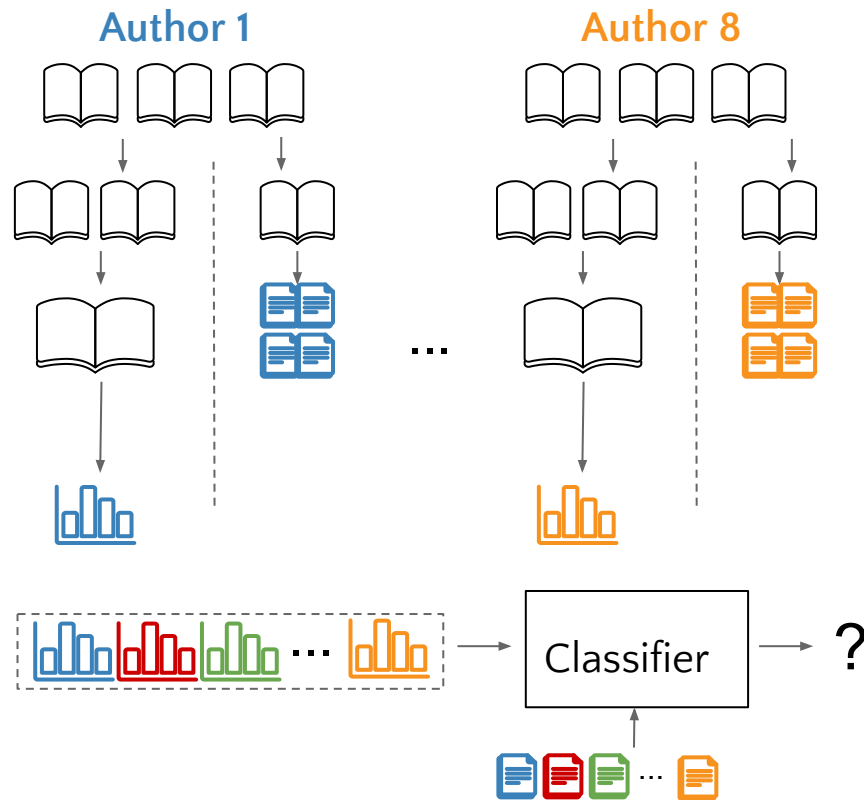
- Download books
- 66% Training, 33% Test
- Split Test in chunks & label
- Merge Training books
- Obtain **Language Model**
- Perplexity** based classifier



3

## Approaches

- Download books
- 66% Training, 33% Test
- Split Test in chunks & label
- Merge Training books
- Obtain **Language Model**
- Perplexity** based classifier
- Test Classifier



---

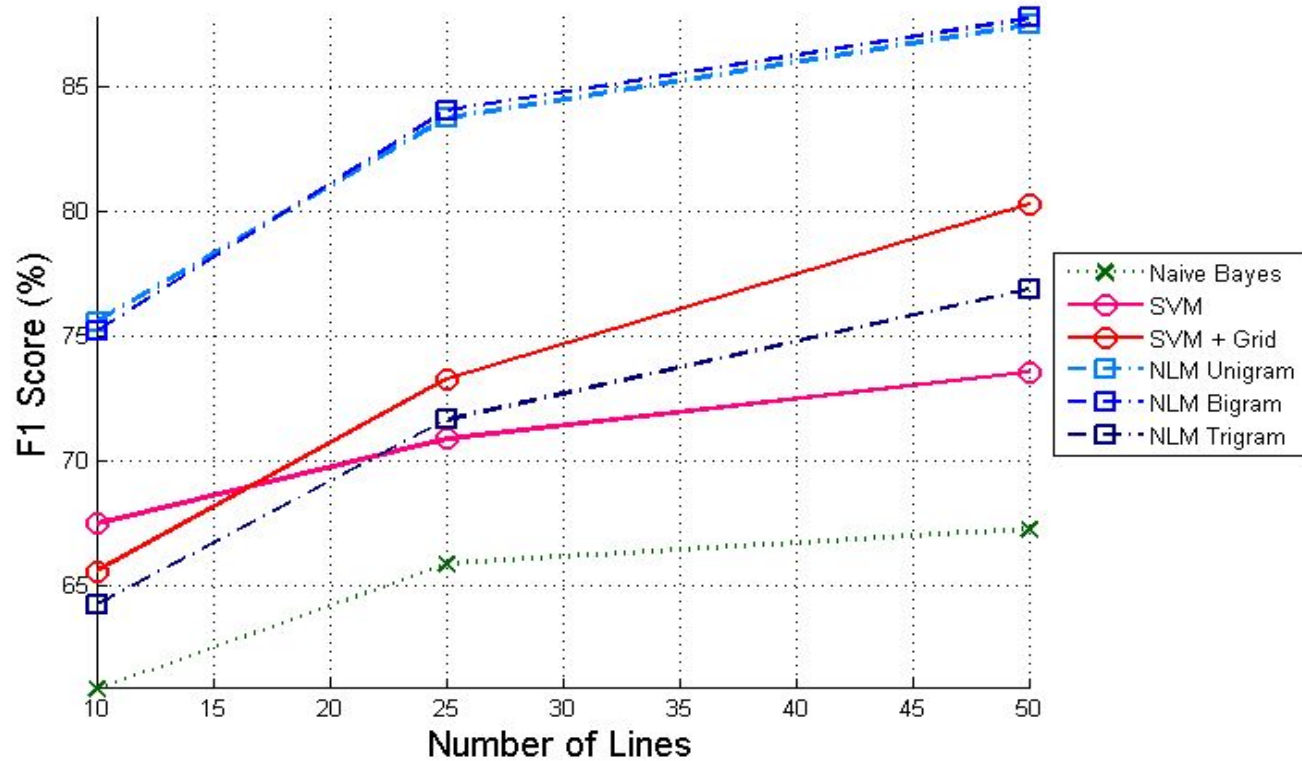
4

# Results

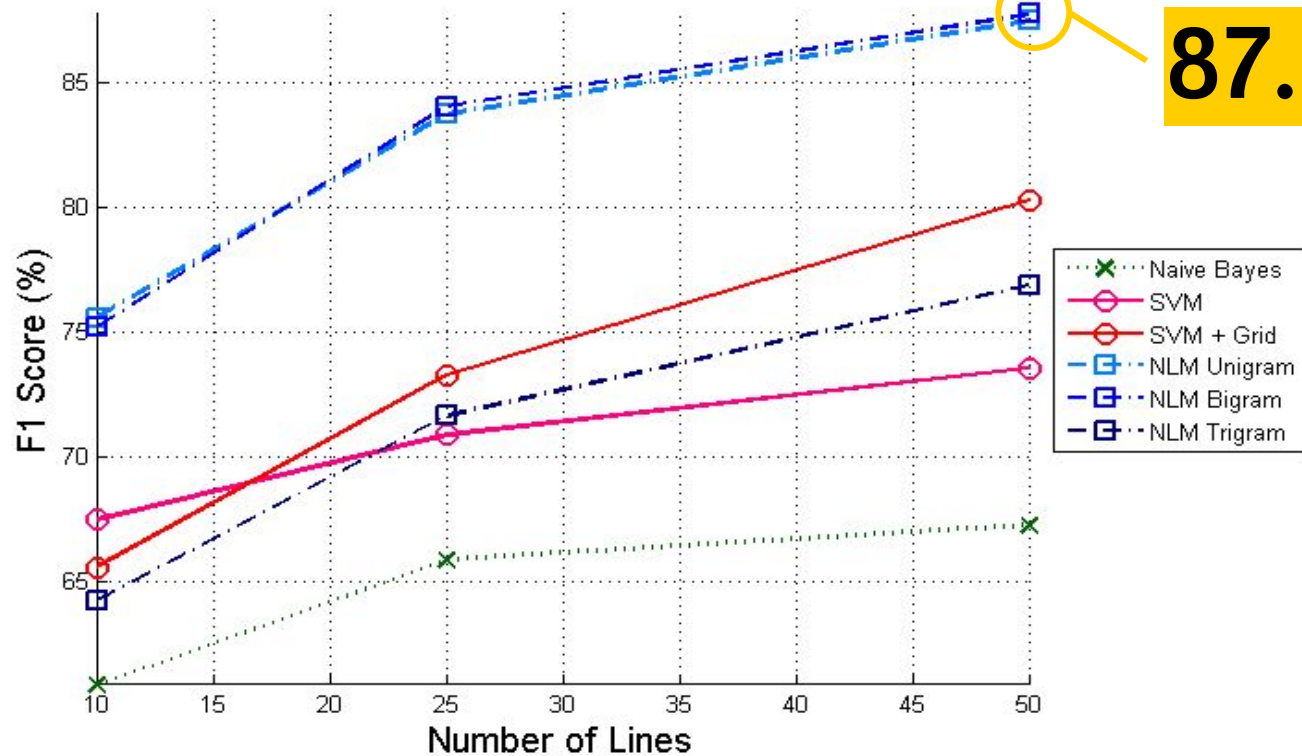
---



Classification Results



Classification Results



87.6%



---

5

# **Conclusions. Future Work**

---

- ◎ **NLM** is, by far, the **best option**
- ◎ **Larger chunks** (more lines) lead to **better results**
- ◎ **Correlation** represents a **problem** in simple approaches

- **Increase** database (training and test)
- Further analyze **relation** between **chunk size** and **accuracy**
- **Equilibrate** number of samples per author
- **Explore** other **smoothing** methods of NLM
- Use of **Neural Networks**



---

# Thanks!

---



## Credits

---

Special thanks to all the people who made and released these awesome resources for free:

- Presentation template by SlidesCarnival
- Photographs by Unsplash

To evaluate the system

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

More Compact way:

$$F_1 \text{ Score} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$