

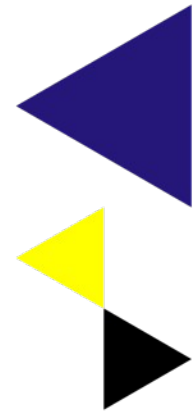


A Tutorial on Language Classification

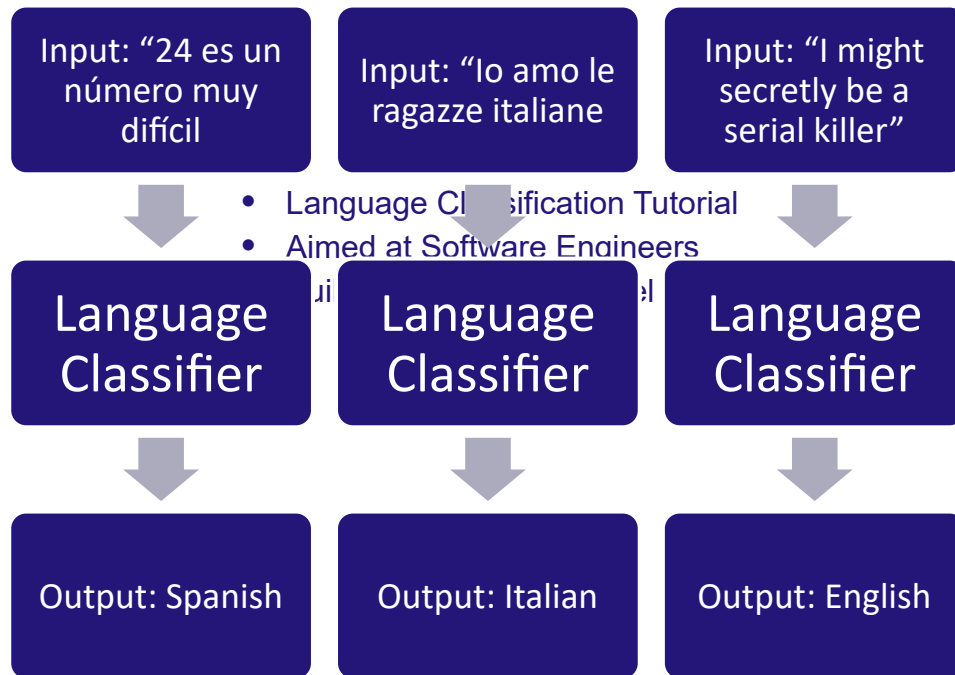
Naam: Kaan Selcuk Gögcay
Datum: 16-01-2026

Table of content

- Context
- Methodology
 - Exploratory Data Analysis (EDA)
 - Modelling
 - Evaluation (including results)
 - Writing Tutorial
- Ideas for Improvements
- Discussions Points

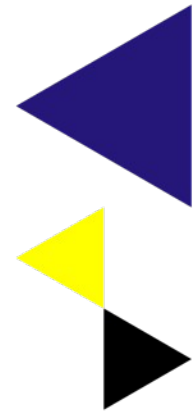


Context



Methodology

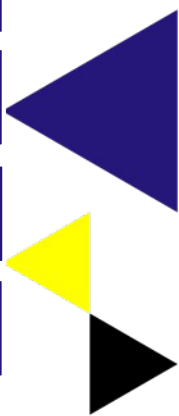
- Exploratory Data Analysis (EDA)
- Modelling
- Evaluation (including results)
- Writing Tutorial



Methodology: EDA

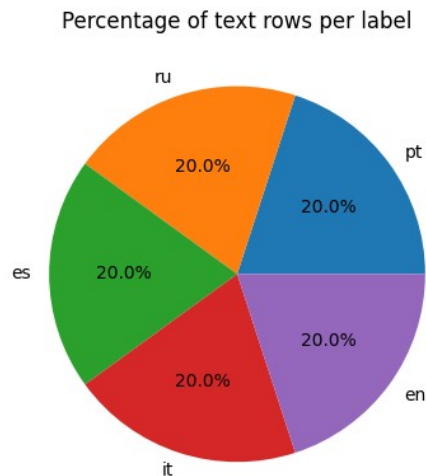
- Language Identification dataset by Papluca
- Amazon reviews
- 20 languages
- Limited to 5 languages (because basic review)
-

Arabic	Bulgarian	German	Greek	English
Spanish	French	Hindi	Italian	Japanese
Dutch	Polish	Portuguese	Russian	Swahili
Thai	Turkish	Urdu	Vietnamese	Chinese



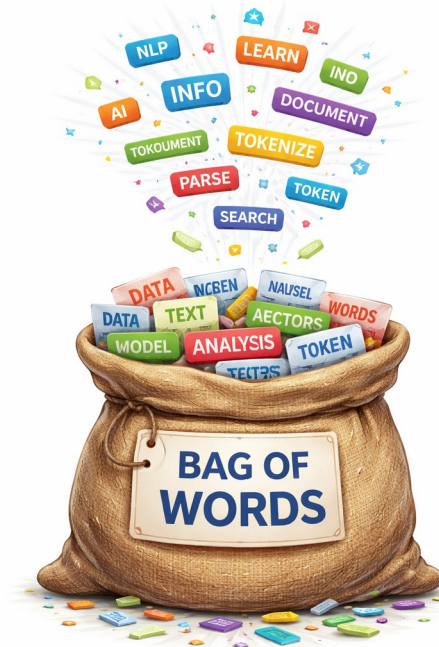
Methodology: EDA

- 5 Languages
- 3500 samples per language
- 17500 samples total
- Evenly split
- Sample: 'A little tight but ok..good quality',
- Sample: 'I cried, oh god it was ugly cries the likes of any I have ever had reading a book. Della and Ren will never be forgotten but their son..... absolutely devastated me. These books will forever be my favorites. They shall never be forgotten!',



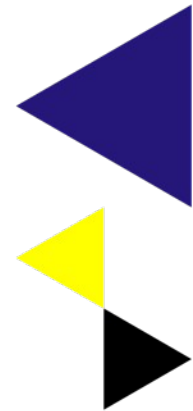
Methodology: Modelling

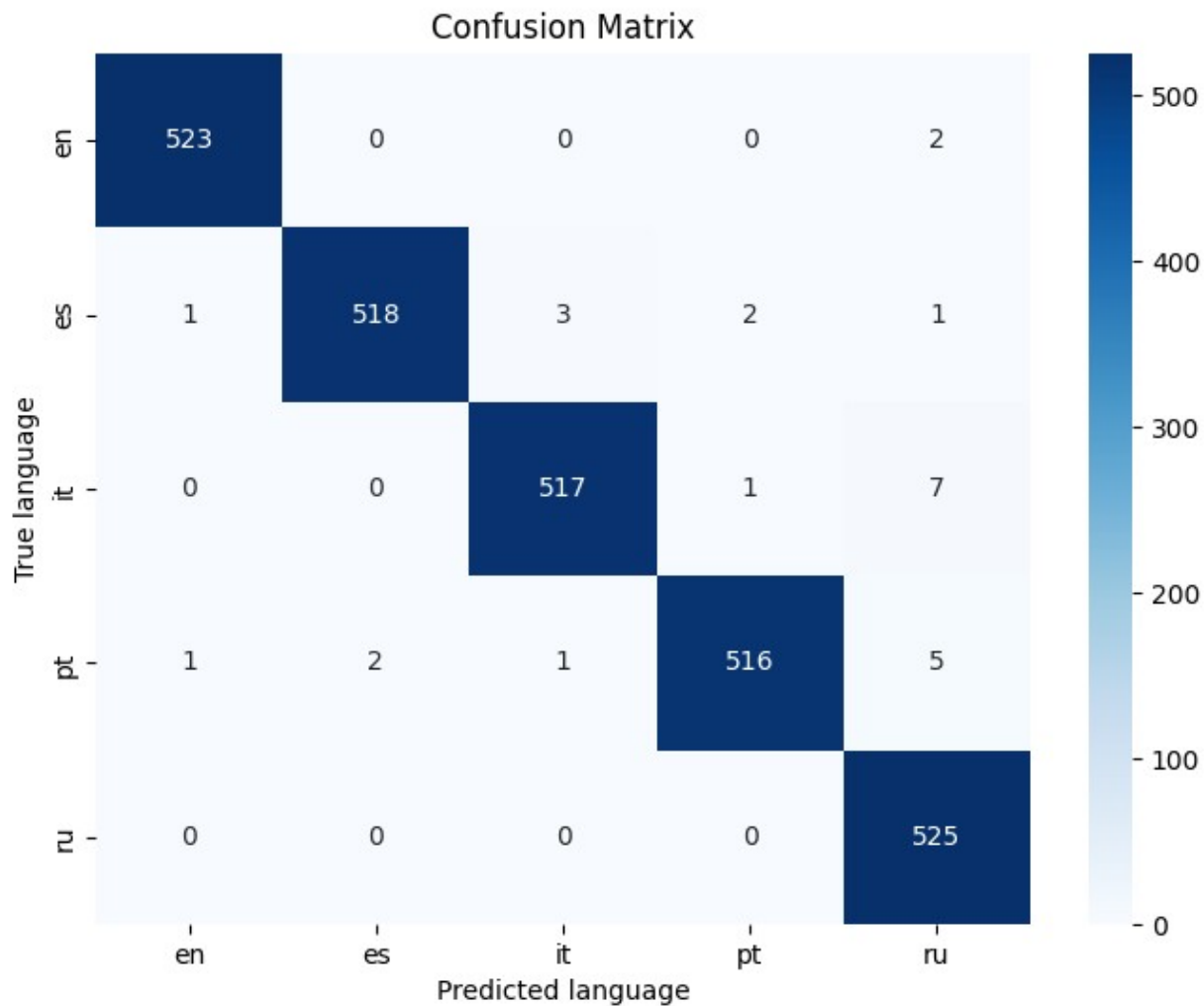
- Deep Learning with Python Chapter 14 Text Classification
- Bag-of-Words model
- Word level tokenization (easiest to understand)
- Vocabulary of 20000 words
-
-



Methodology: Evaluation

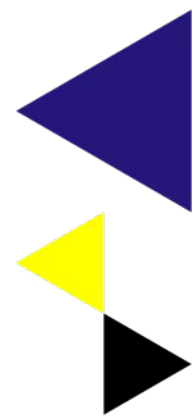
- Train accuracy: 99.94%
- Validation accuracy: 98.78%
- Test accuracy: 99.00%
- Confusion Matrix
-



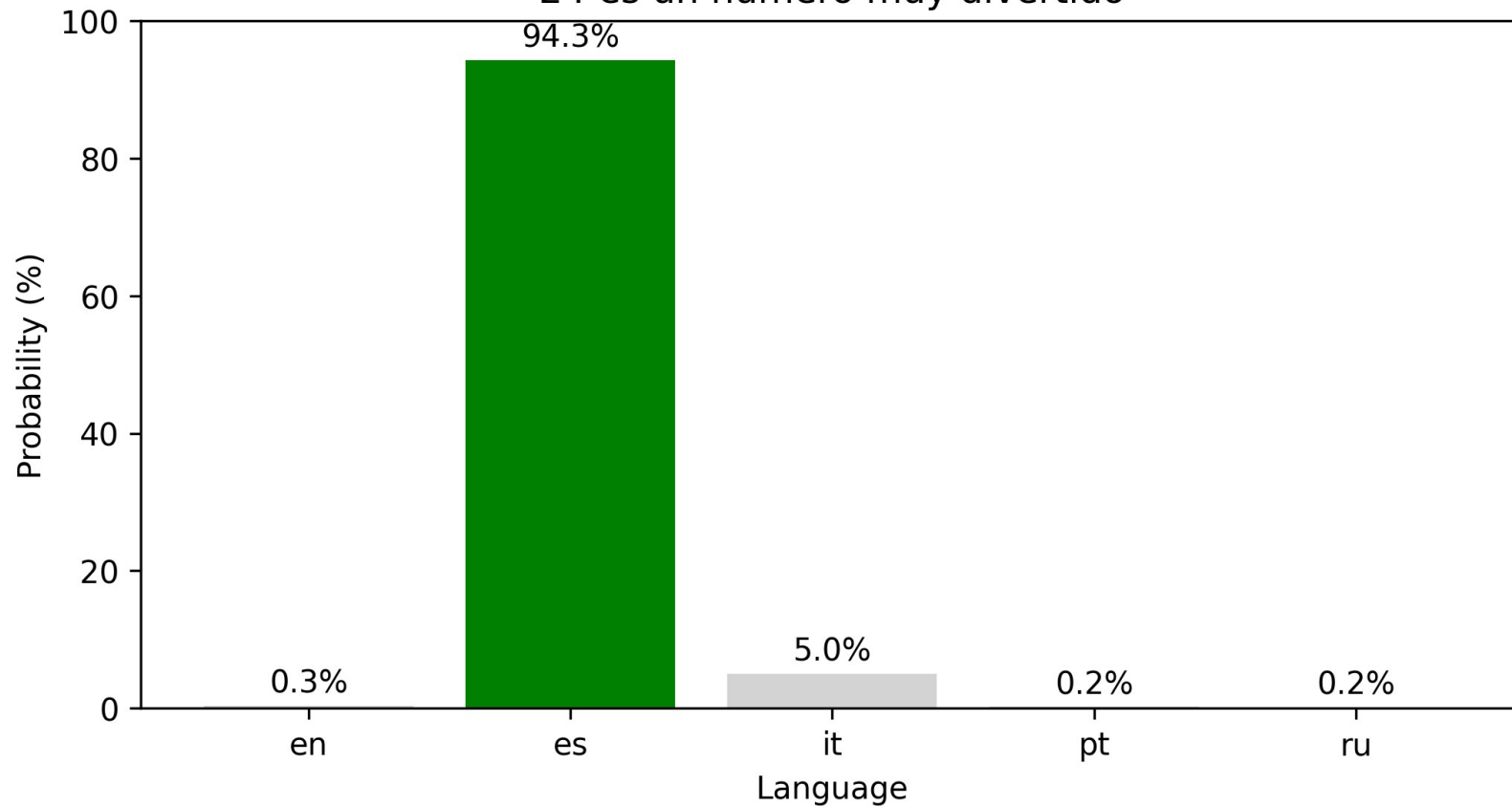


Methodology: Evaluation

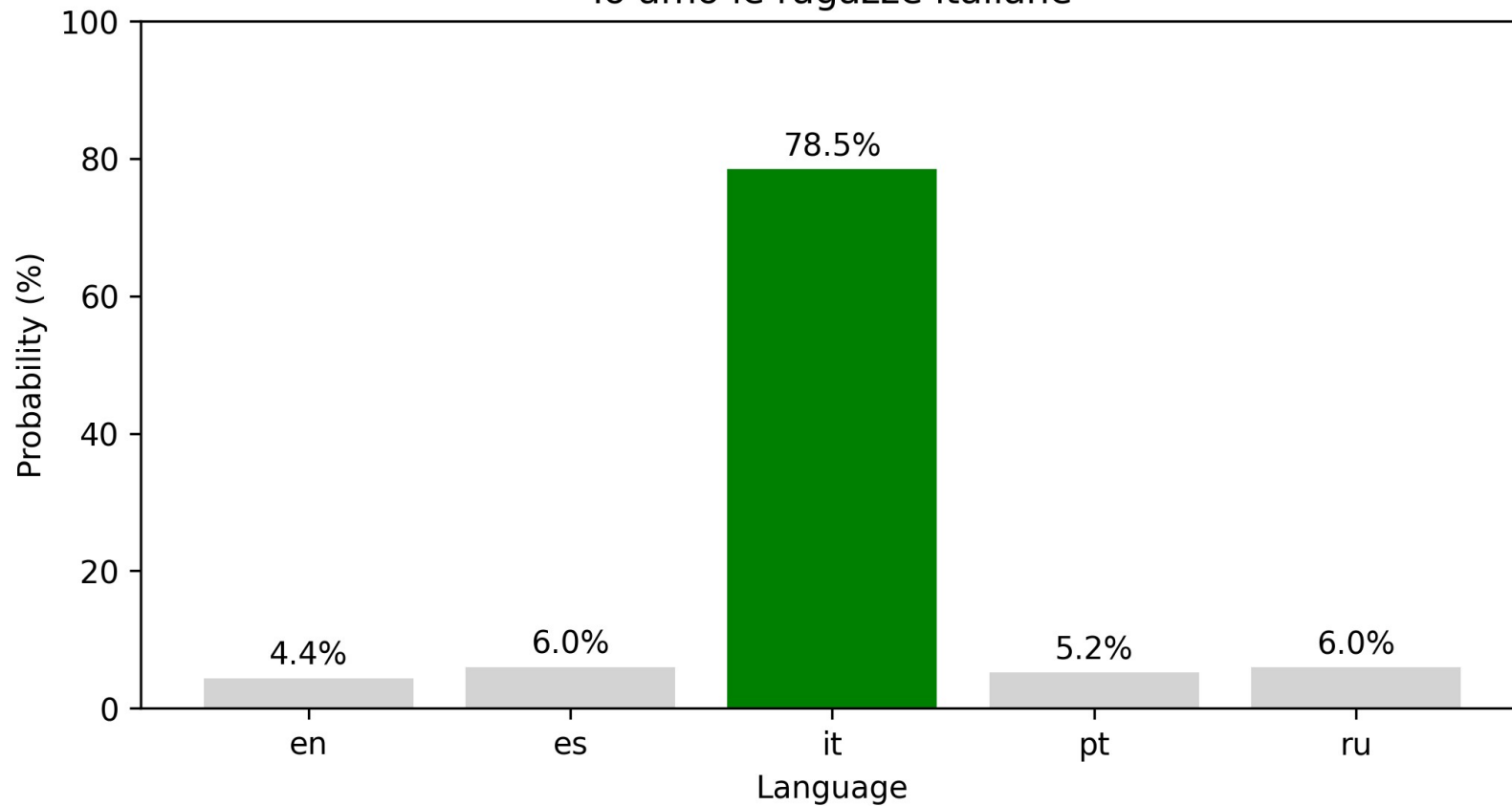
- Train accuracy: 99.94%
- Validation accuracy: 98.78%
- Test accuracy: 99.00%
- Confusion Matrix
- Classifying sentences
-



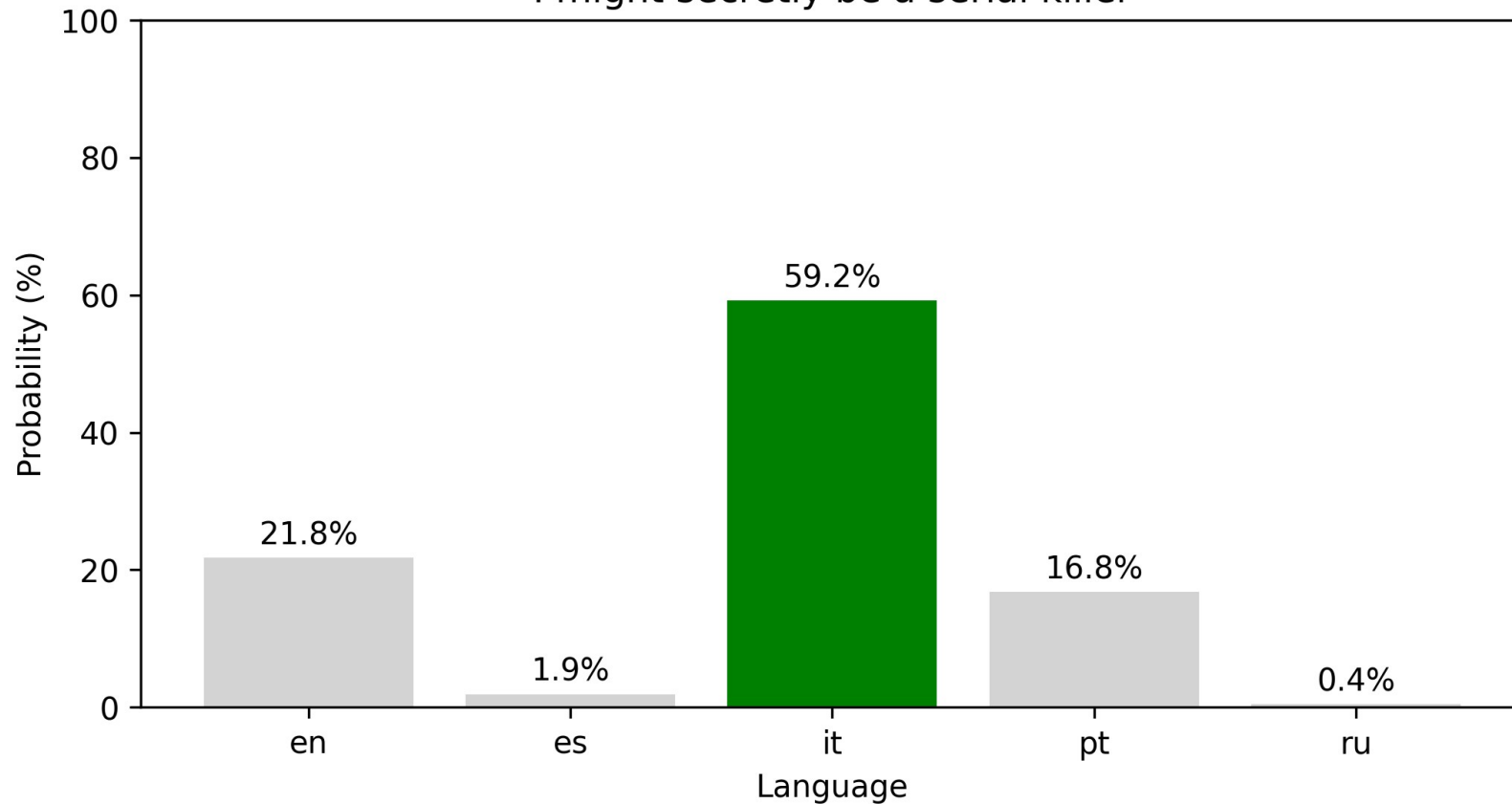
Language prediction for:
"24 es un número muy divertido"



Language prediction for:
"lo amo le ragazze italiane"

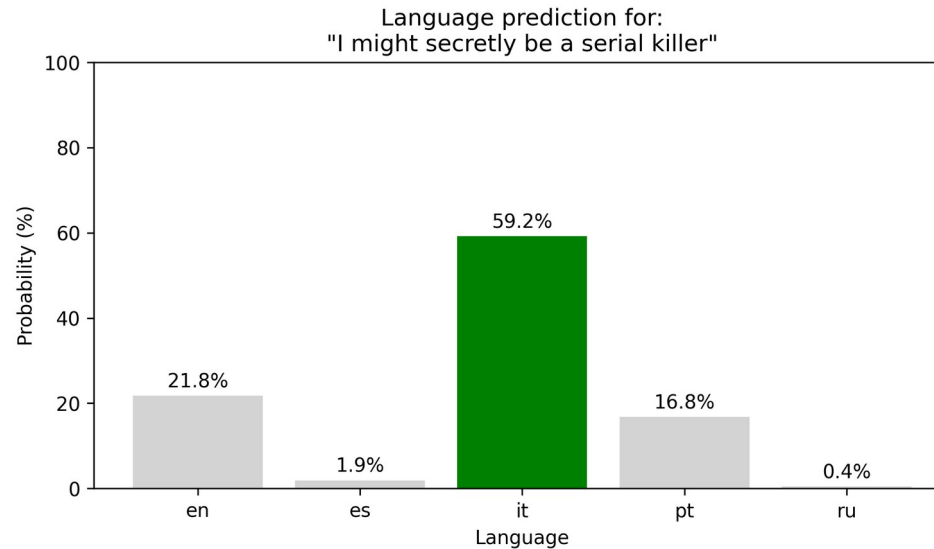


Language prediction for:
"I might secretly be a serial killer"



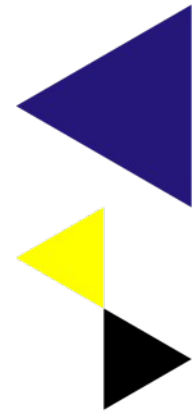
Methodology: Evaluation

- Vocabulary
 - i - 6
 - might - 1223
 - secretly - 0
 - be - 77
 - a - 3
 - serial - 13297
 - killer - 0



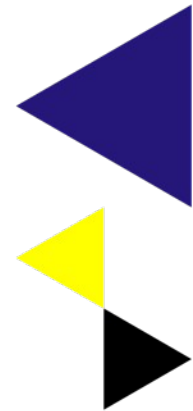
Methodology: Writing Tutorial

- What you will learn
- Introduction
- Approach Overview
- Scope and Limitations
- Coding
 - Step 1: Setting up the environment
 - Step 2: Loading the data
 - Step 3: Filtering languages
 - Step 4: Data Preperation
 - Step 5: Tokenization
 - Step 6: Bag-of-Words model
 - Step 7: Training the Classifier
 - Step 8: Evaluating the model
 -



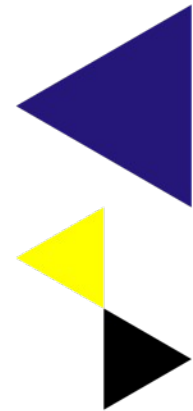
Methodology: Writing Tutorial

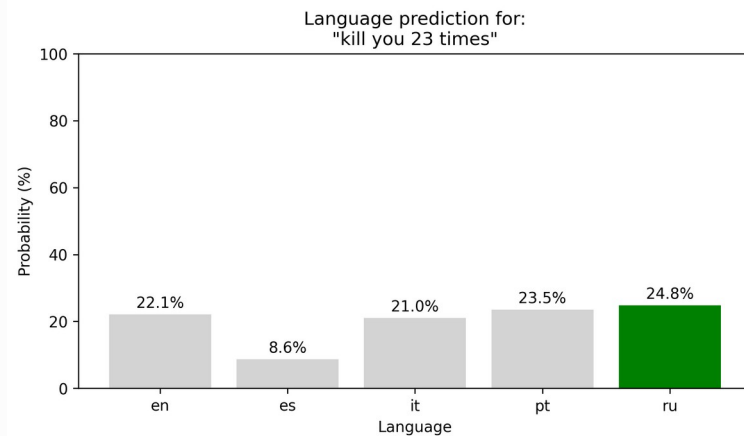
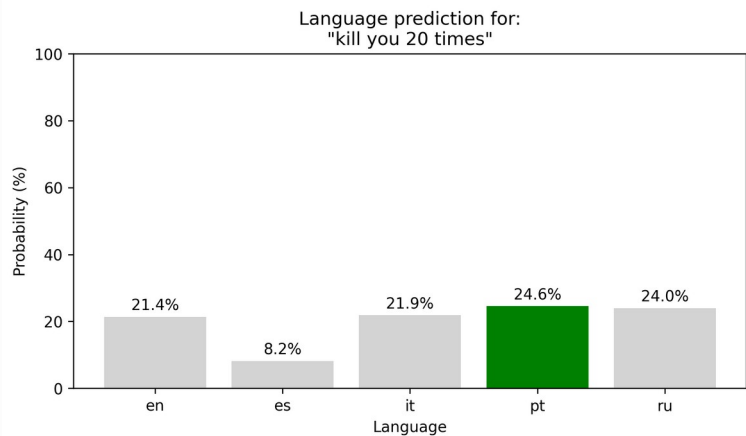
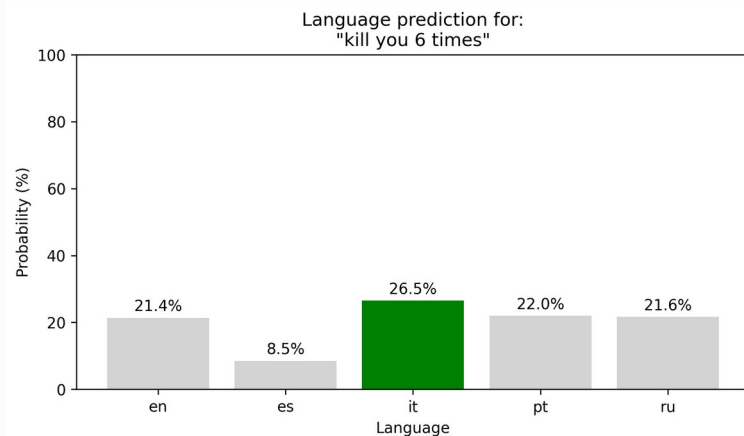
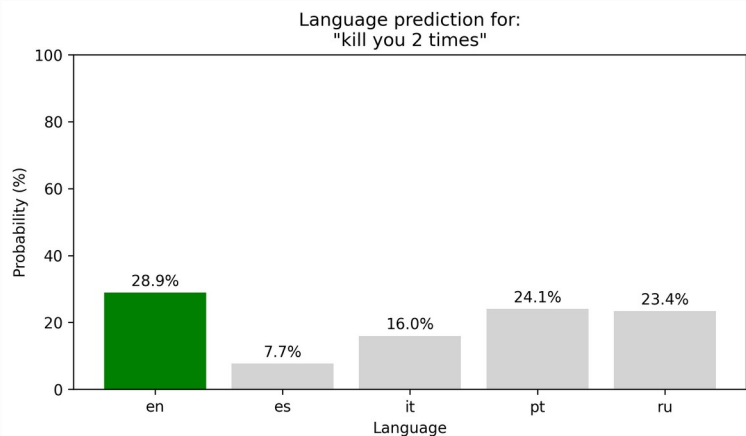
- What you will learn
- Introduction
- Approach Overview
- Scope and Limitations
- Coding
 - Step 1: Setting up the environment
 - Step 2: Loading the data
 - Step 3: Filtering languages
 - Step 4: Data Preperation
 - Step 5: Tokenization
 - Step 6: Bag-of-Words model
 - Step 7: Training the Classifier
 - Step 8: Evaluating the model
 -



Ideas for improvement

- Remove numbers from data
 -

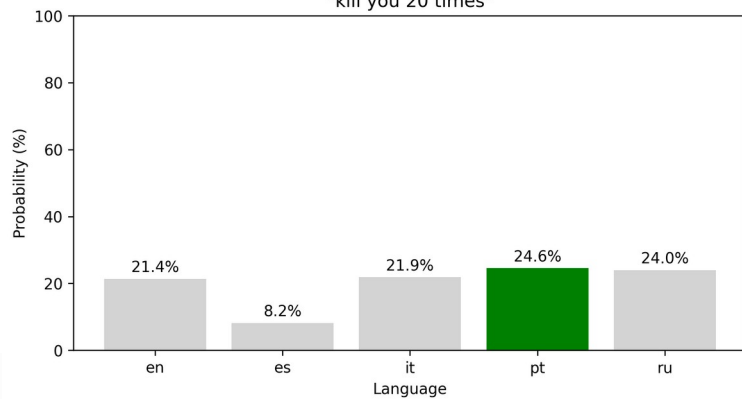




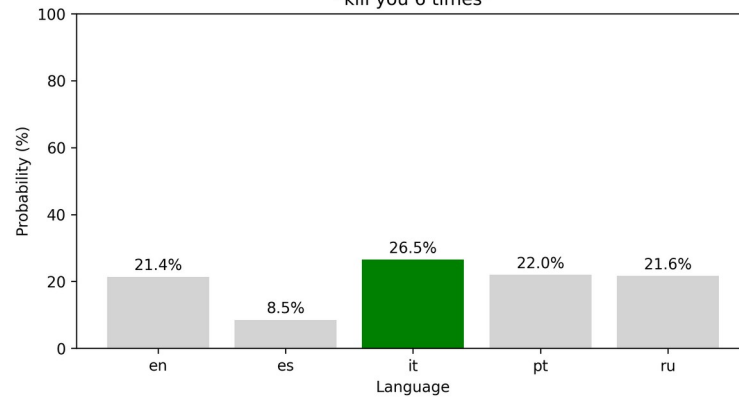
Vocabulary

- kill - 10661
- you – 69
- 2 – 100

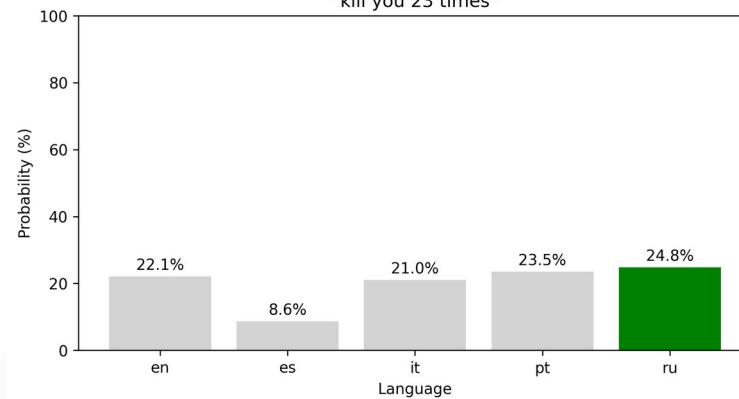
Language prediction for:
"kill you 20 times"



Language prediction for:
"kill you 6 times"



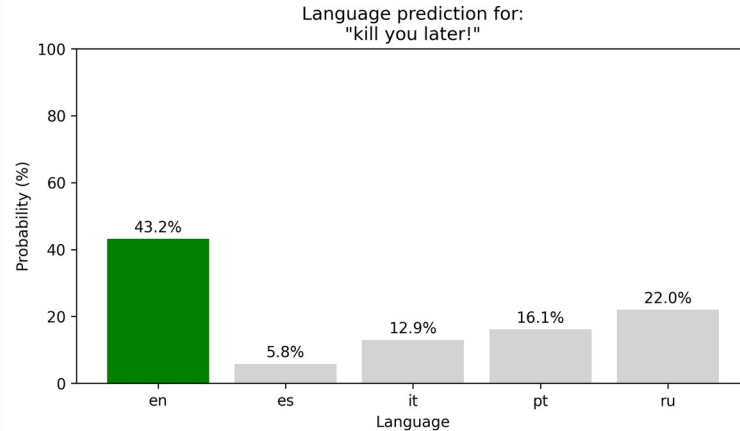
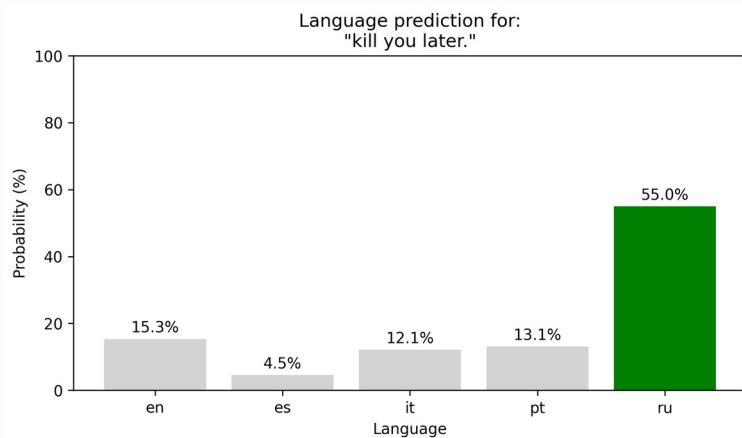
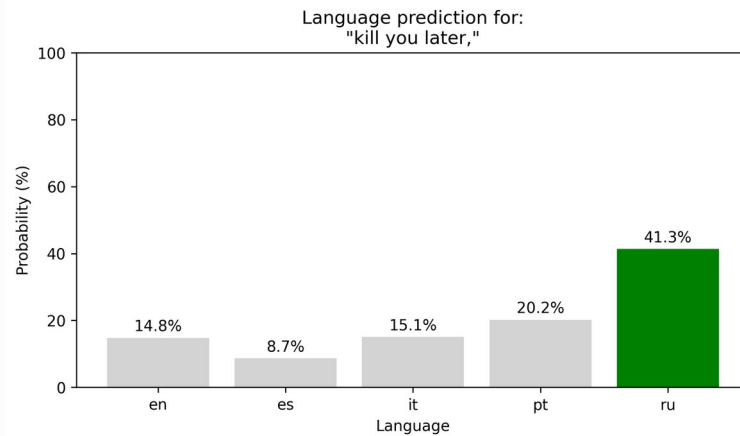
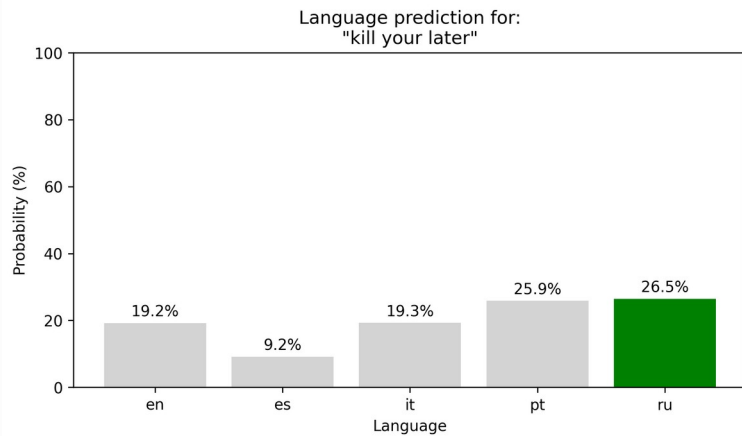
Language prediction for:
"kill you 23 times"



Ideas for improvement

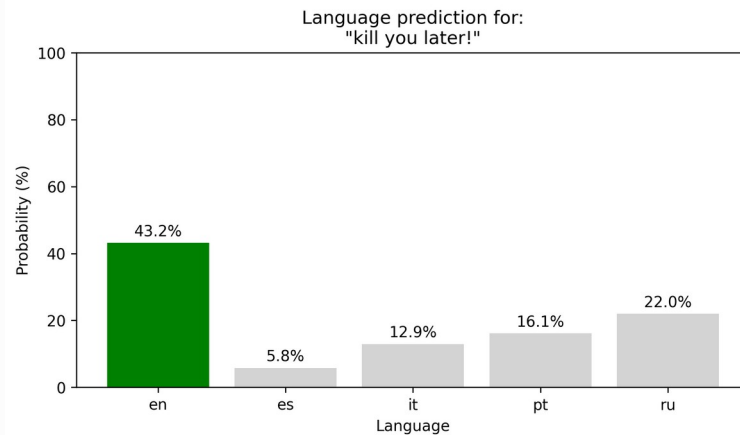
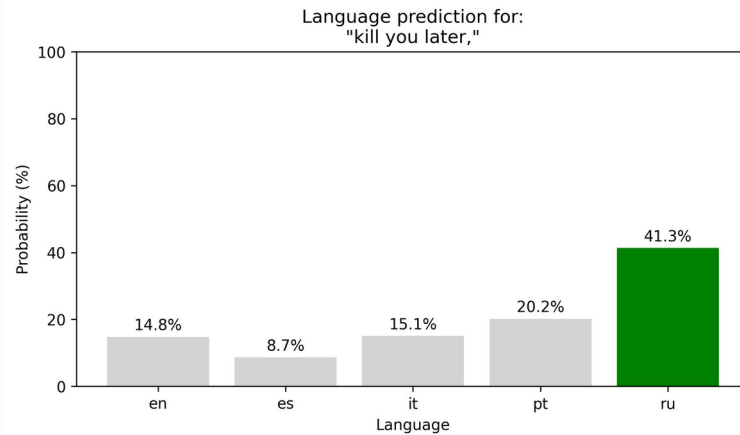
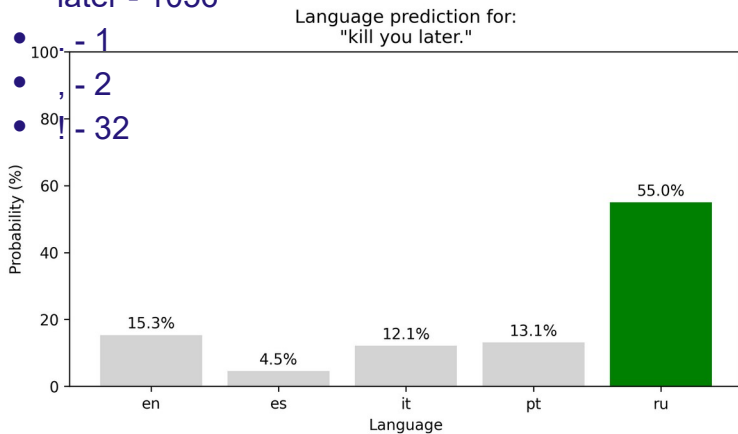
- Remove numbers from data
- Consider removing punctuation
 -





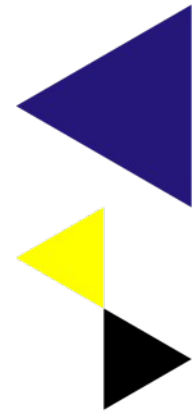
Vocabulary

- kill - 10661
- you – 69
- later - 1056



Ideas for improvement

- Remove numbers from data
- Consider removing punctuation
- Add more languages
-



Discussion Points

- Punctuation. Keep? Remove? Partially?
- Every language contains some english words
 - В блант резюме (у меня мало времени) ,
возможно , просить нас (в массы) **Shell** 20
баксов в год...
 - На веб-сайте агентства по вопросам
международного развития , на сайте " **Care**
" или...
 - Обсуждаем статью в журнале " **The New York
Times** " в журнале " **New York Times**
автор...
 - Lo utilice con unos guantes de boxeo **everlast** de
12 onzas y también...
 - Per l'anno in corso, si prevede un fatturato di 94
milioni di dollari e un utile di 26 centesimi per
azione da operazioni **continue**.