

# Pythonzão: uma saga do Python no mundo de Big Data

---

Bruno Ábia

[bruno.abia@icomp.ufam.edu.br](mailto:bruno.abia@icomp.ufam.edu.br) 

# Quem sou eu?



- Engenheiro de Computação (Faculdade Fucapi);
- Mestre em Informática (UFAM);
- Doutorando em Informática (UFAM);
  
- Líder de Desenvolvimento (UNASUS/AM) (UEA);
  - Algoritmos de Recomendação;
  
- Pesquisador do Grupo DNS LAB(UFAM);
  - Ciência de dados;
  - Machine Learning;
  - Redes Sociais;
  - Análise de Sentimentos

# Introdução



# Acessórios



# Redes Sociais

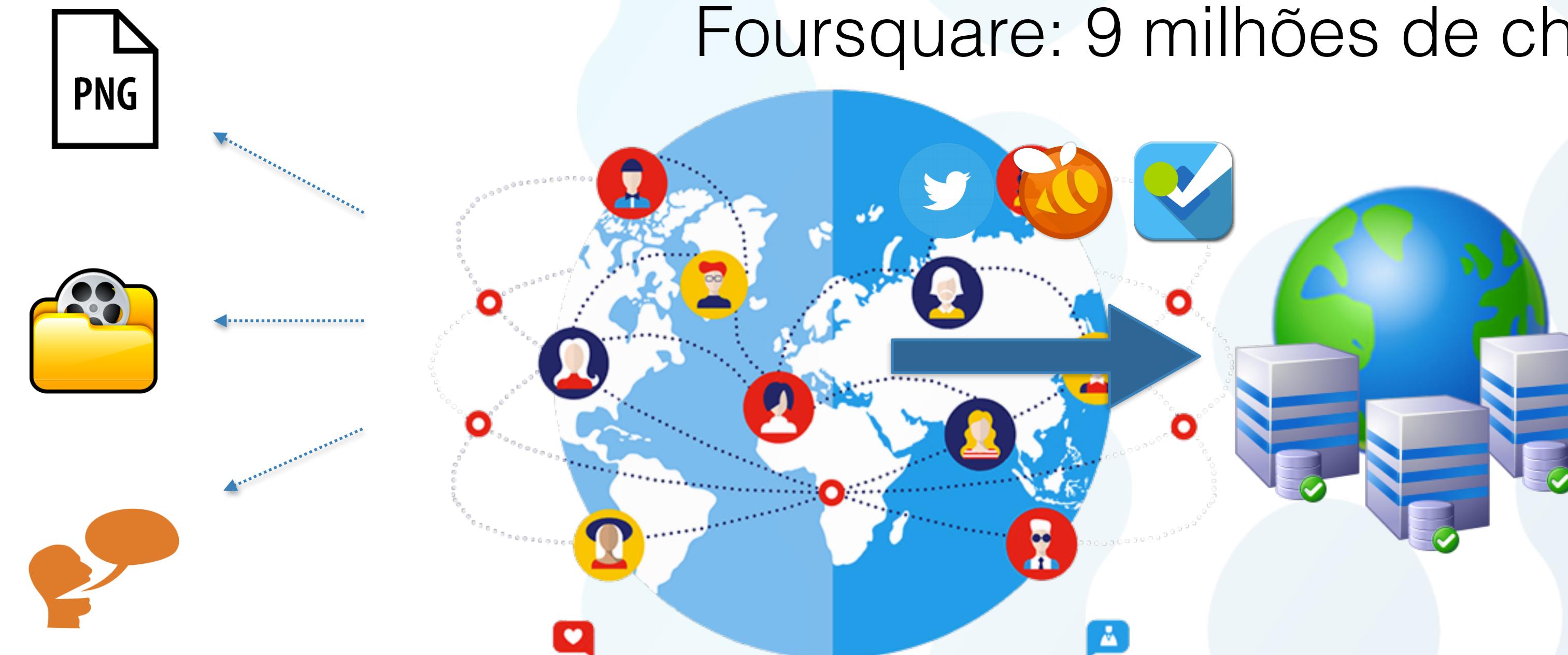


# Smart Cities

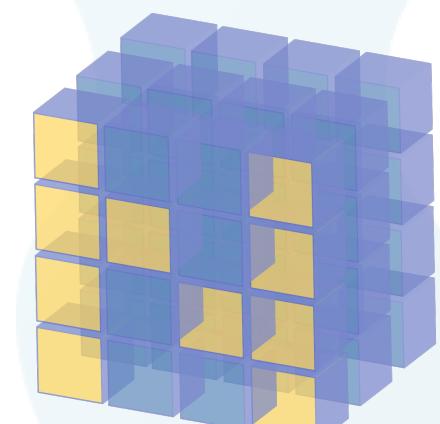
# Introdução

Twitter: 1 Bilhão de tweets mensais

Foursquare: 9 milhões de checkins por mês



# Introdução

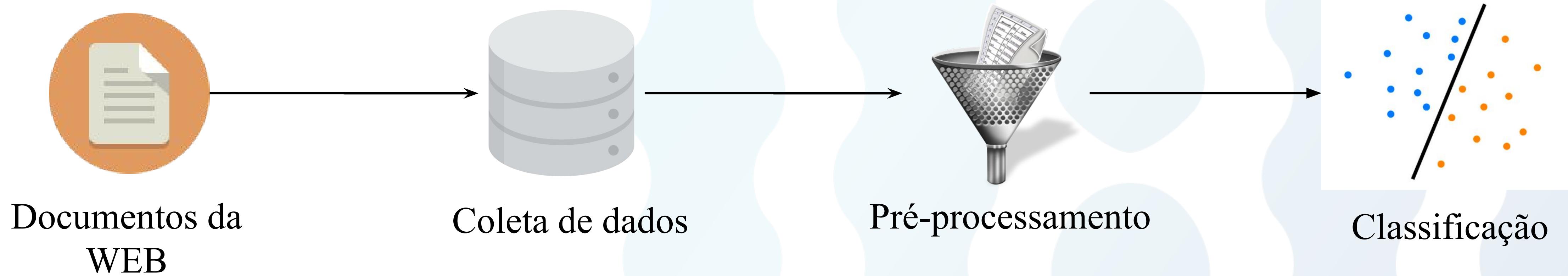


NumPy

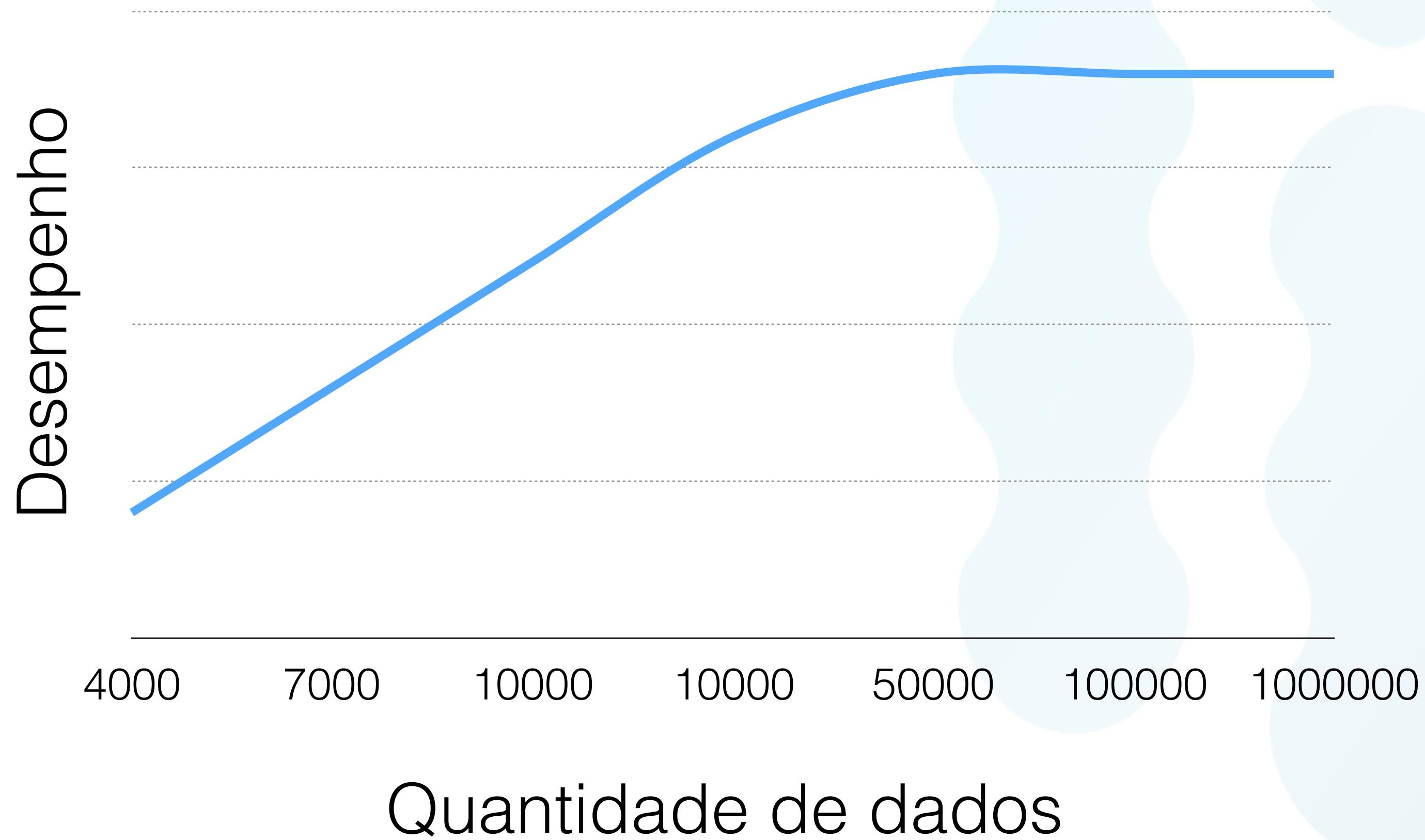


scikit  
learn

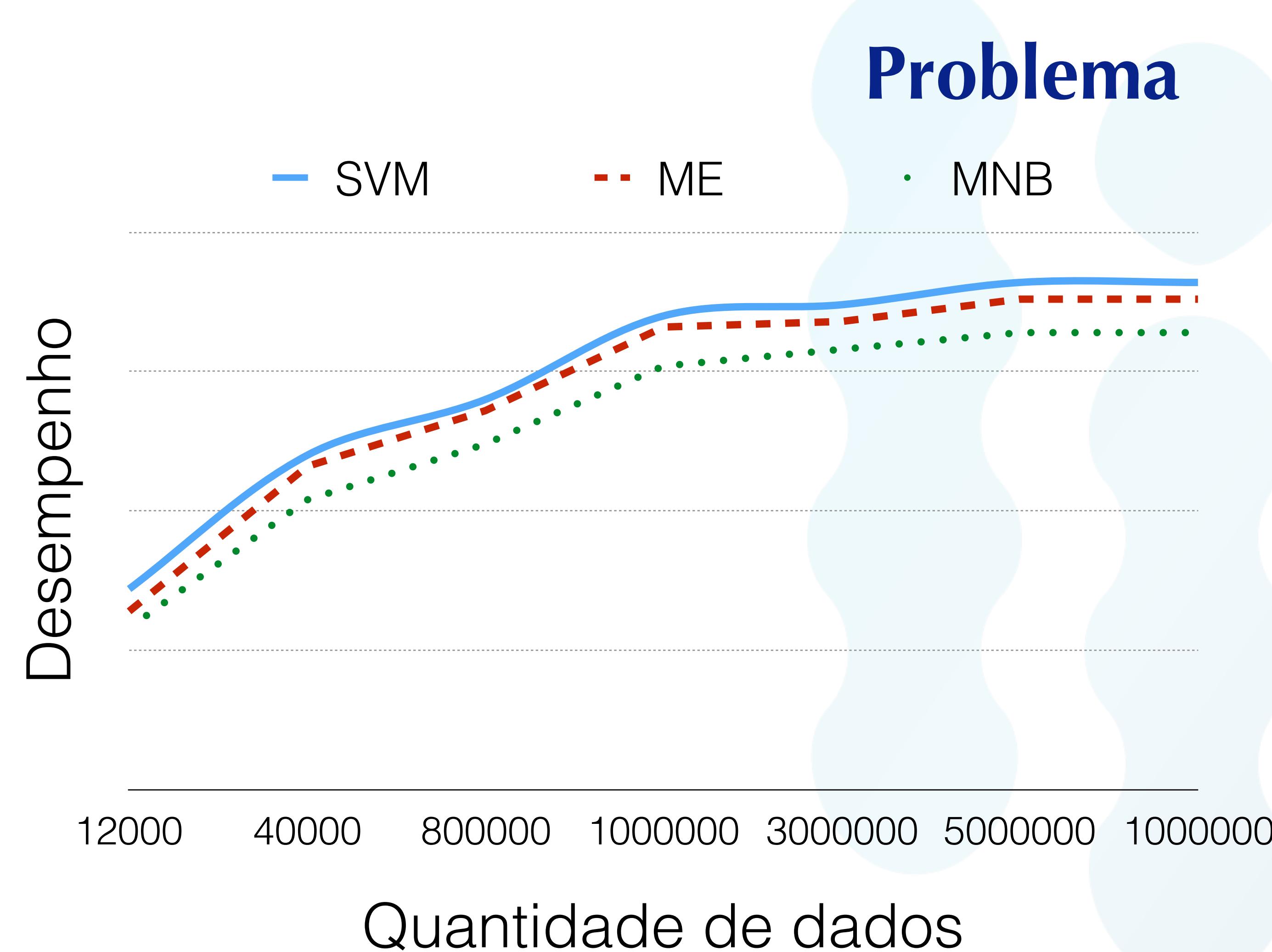
# Introdução



# Problema



O desempenho dos métodos de aprendizagem diminui.



Reduzindo o problema  
de variância e  
aumentando generalização

# Soluções

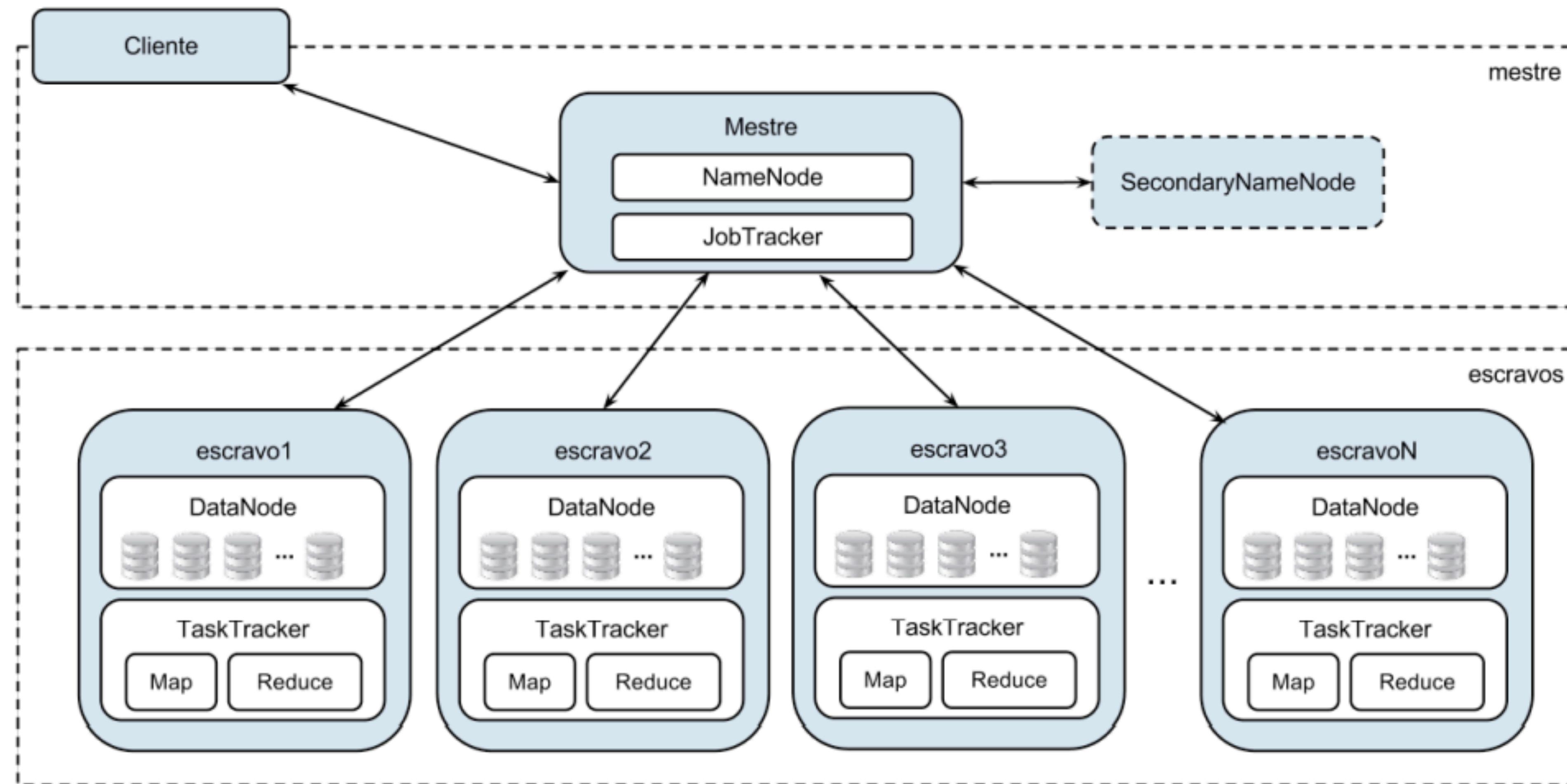


# Soluções



- Plataforma de computação distribuída;
- Custo Baixo e em grande escala;
- Pode analisar dados complexos;

# Soluções





## Soluções

- São variações do Hadoop;
- Se baseiam nos dados salvos no HDFS;
- Tem como objetivo facilitar o acesso aos dados;



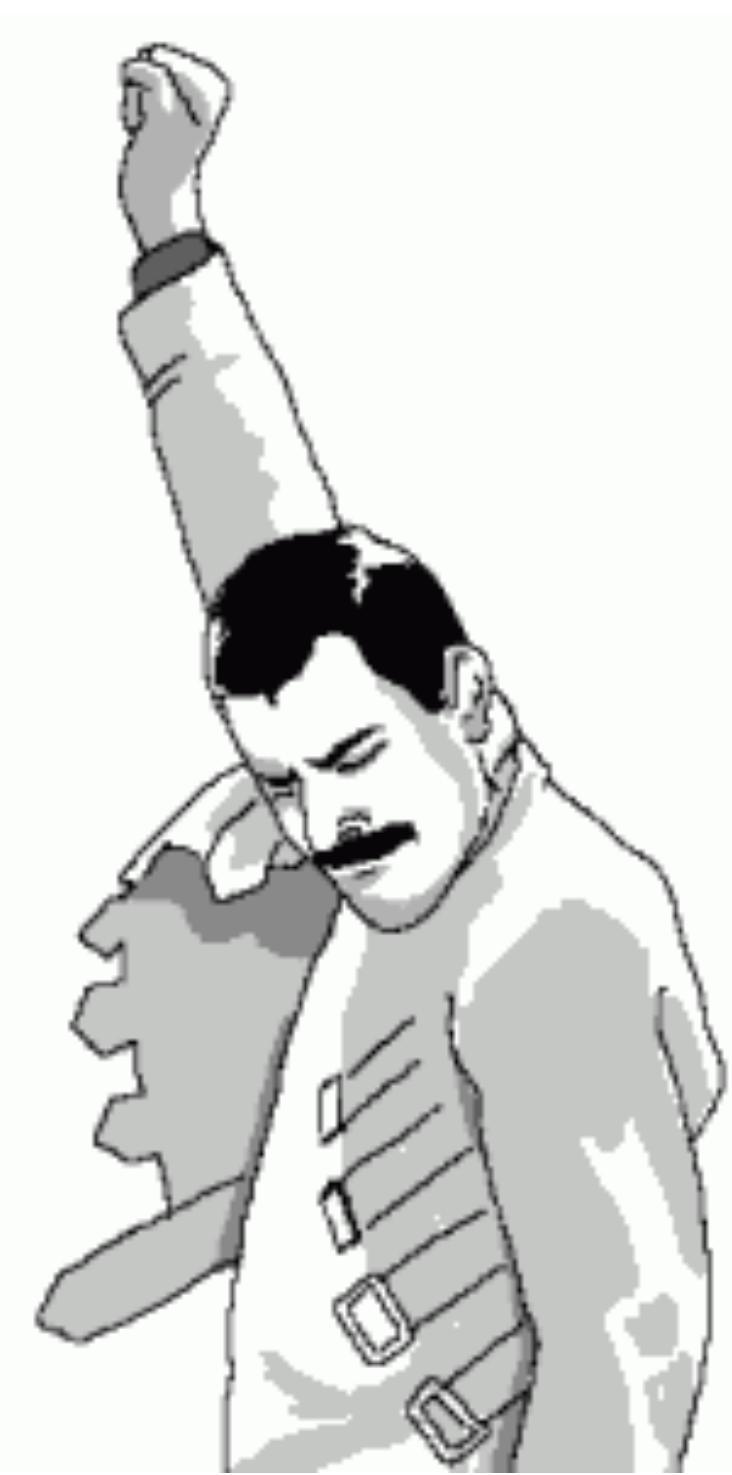
### Desvantagem:

- Só funciona para dados estáticos;
- Necessário conhecer Java;
- Podemos escrever códigos em Python, porém precisamos da Interface **Jython**;



## Soluções

- Framework que utiliza Resilient Distributed Datasets (RDDs);
- Opera o conjunto de dados de uma só vez;
- 100 vezes mais rápido que o Hadoop para análises in-memory;

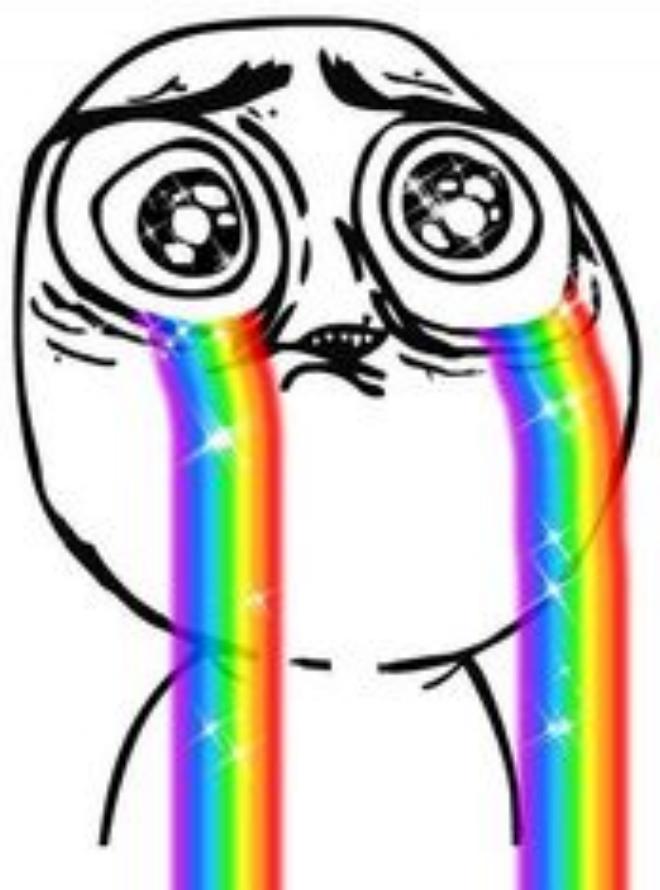


**PySpark**

## Soluções

**E sim!!!! Temos um versão pra Python!!.**

# Recursos Disponíveis



Machine Learning

Deep Learning

Complex Networks

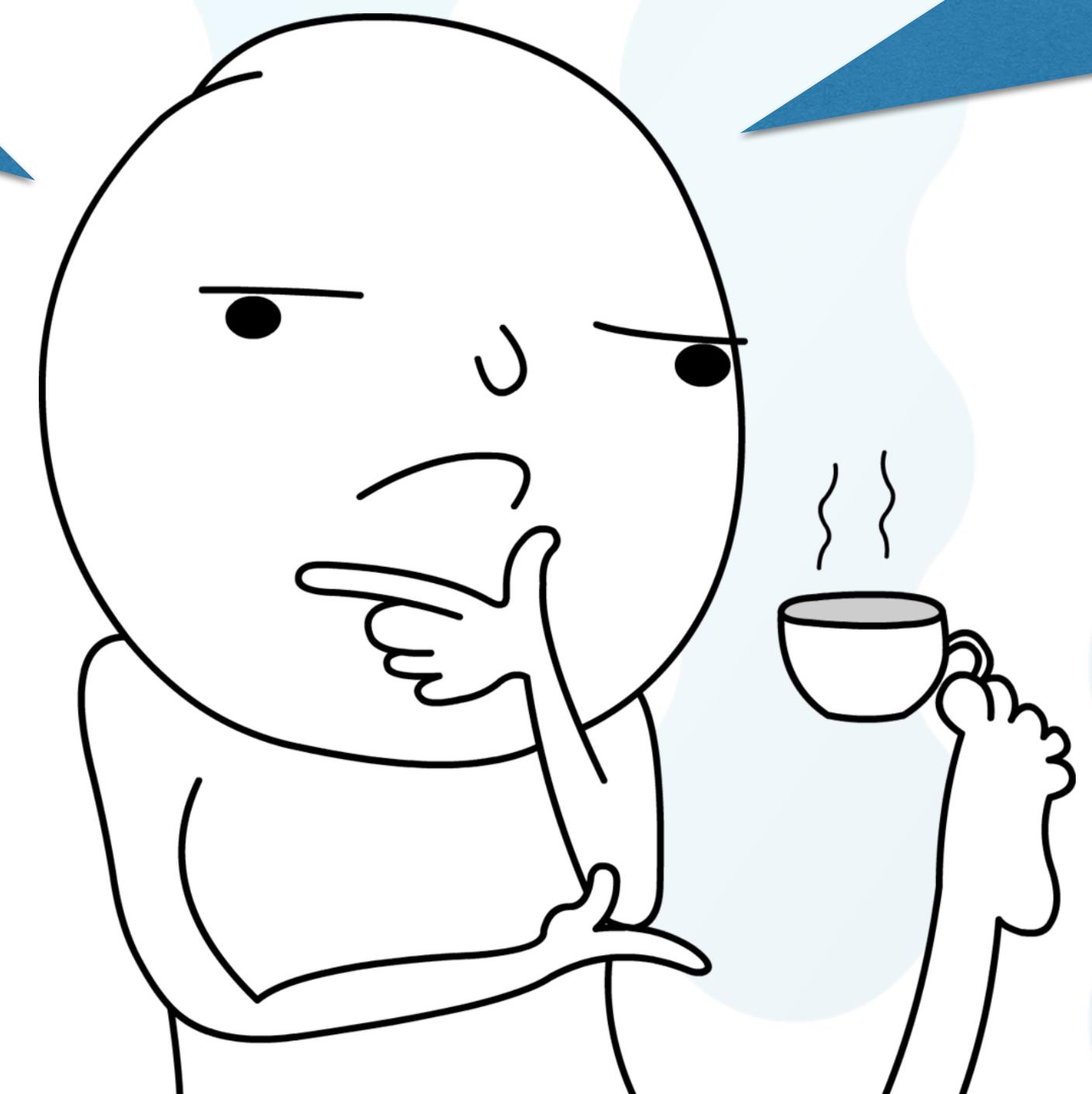
Natural Process Language

Data Mining

# Dúvidas Comuns

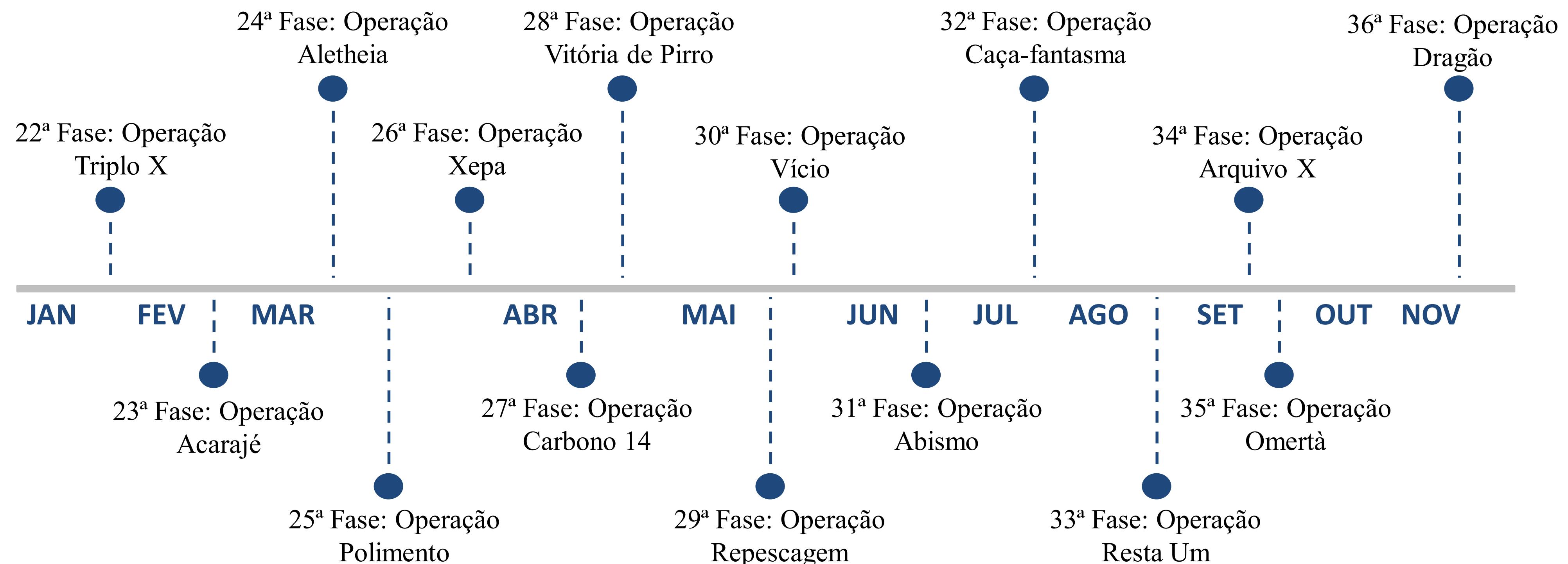
Onde posso aplicar?

É complexo usar o Pyspark?



# Exemplo de Aplicação

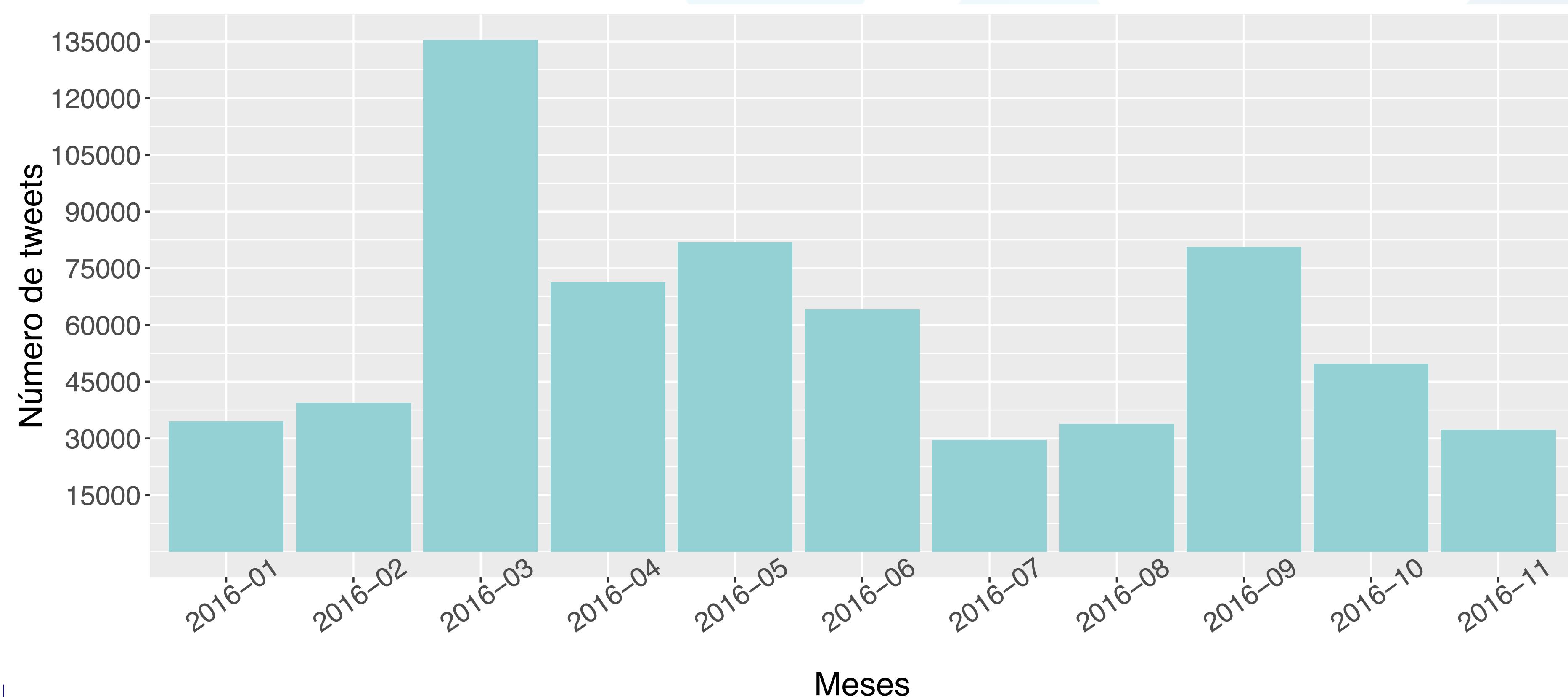
## Operação Lava Jato da Polícia Federal



# Exemplo de Aplicação

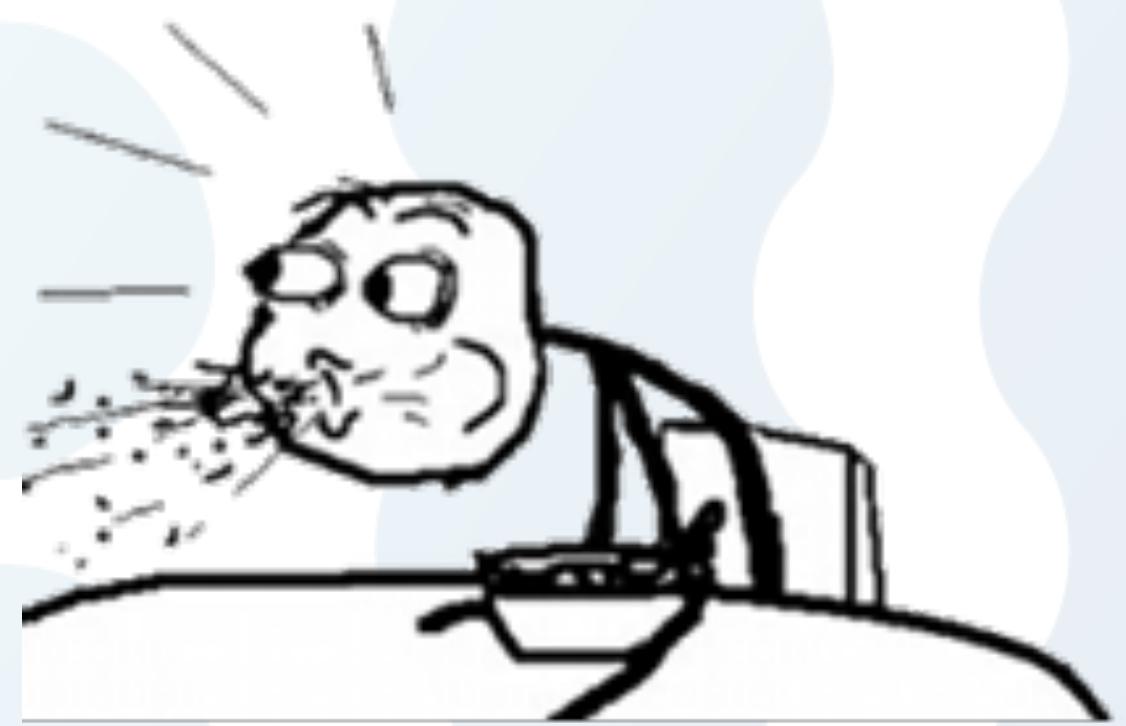
## Base de Dados

- Período Coletado 01/01/2016 a 20/11/2016;
- Total de 652.210 tweets;

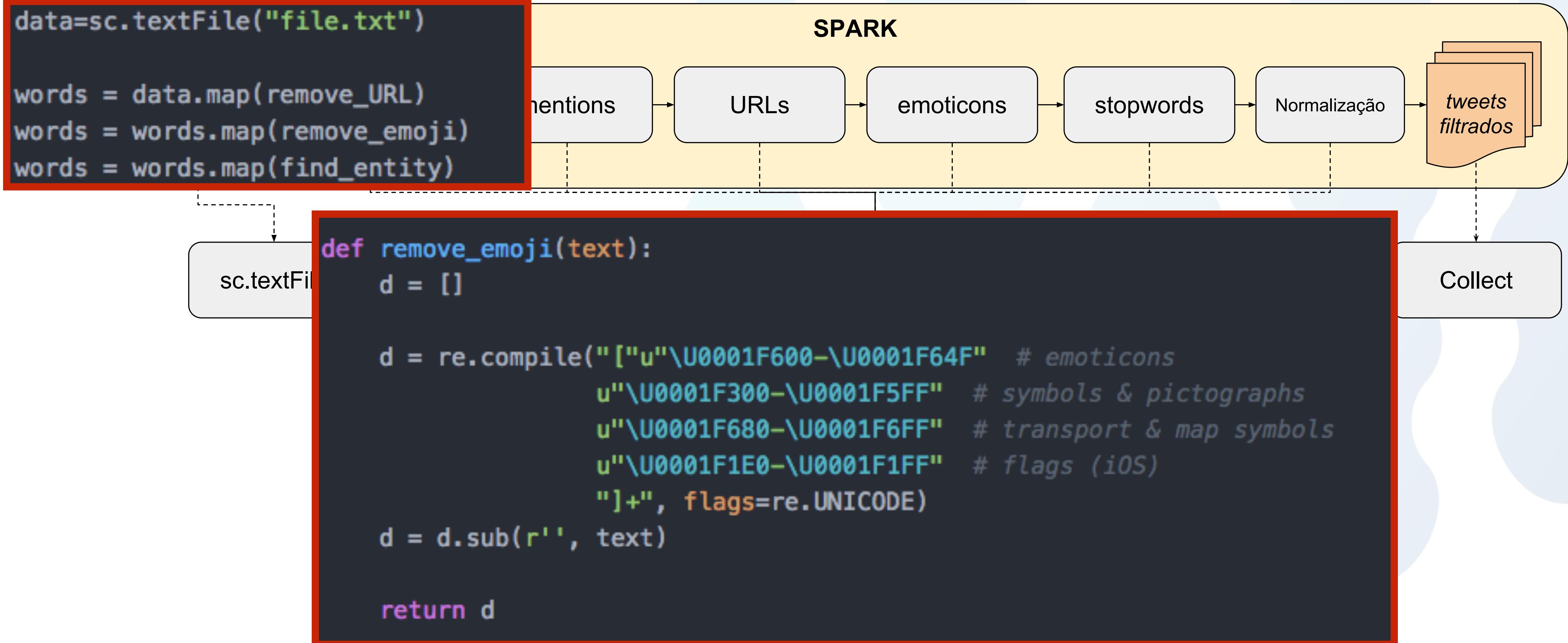


# Como iniciar?

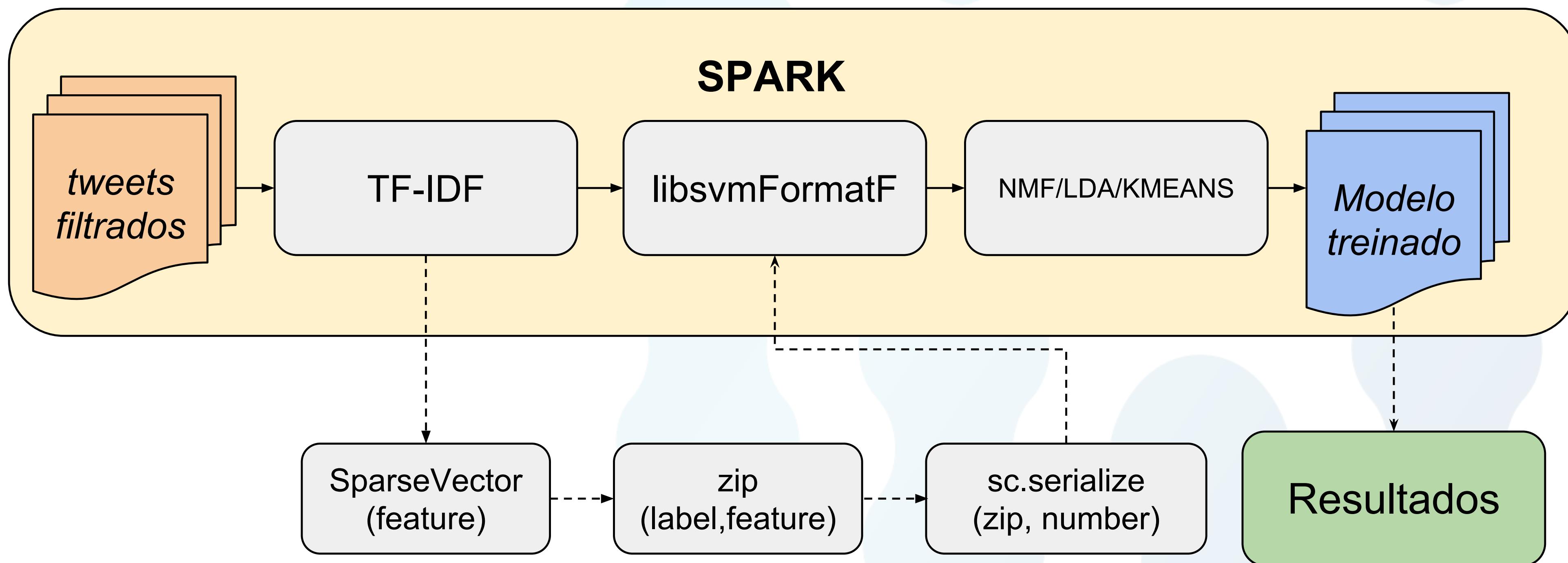
```
os.environ['SPARK_HOME'] = "/home/aluno/spark-1.6.1-bin-hadoop1"  
|  
# # Append pyspark to Python Path  
sys.path.append("/home/aluno/spark-1.6.1-bin-hadoop1/python")  
  
#Create Pyspark import  
try:  
  
    from pyspark import SparkContext  
    from pyspark import SparkConf  
  
    print("Successfully imported Spark Modules")  
    sc = SparkContext("local", "Simple Language Model Computing")  
  
except ImportError as e:  
    print("Can not import Spark Modules", e)  
  
    sys.exit(1)
```



# Pré-Processamento



# Treinamento dos Algoritmos de M.L.



# Resultados

## Notícia compartilhada (SET/2016)



27/09/2016 16h21 - Atualizado em 27/09/2016 22h21

### STF aceita denúncia e torna Gleisi e Paulo Bernardo réus na Lava Jato

Casal é acusado de pedir e receber R\$ 1 milhão desviados da Petrobras.  
Defesa nega repasse e aponta divergências entre delações premiadas.

$$\text{Percentual} = \frac{\text{Tópicos correspondentes}}{\text{Total de Notícias}}$$



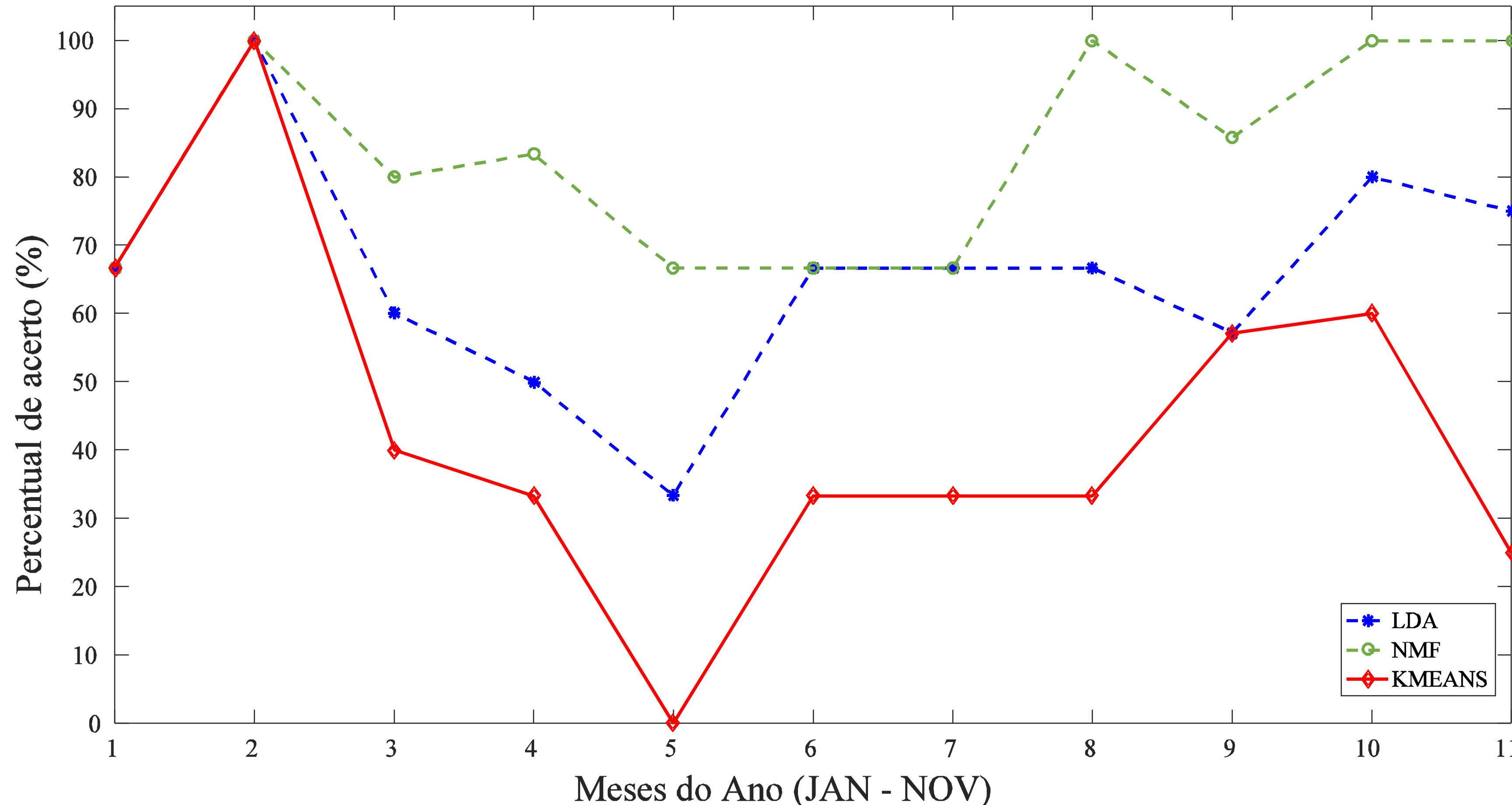
Exemplo de *tweet*:  
O **Paulo Bernardo** e **Gleisi Hoffman** são reús na Lava Jato.. Justiça sendo feita!!! Fora PT!!

Petrobras Mantega STF  
Lula Marisa MPF  
Sergio Moro Bumlai  
Palocci Antonio Dilma  
PT Paulo Gleisi

(Zhao et al., 2011)

# Resultados

## Percentual de Acerto por mês



# Conclusão

- Big Data é uma realidade em nosso dia a dia;
- Precisamos pensar como cientistas de dados;
- Python é realmente a linguagem para resolver esses problemas;
- O céu é o limite;



# Obrigado!



Bruno Ábia Souza

[bruno.abia@icomp.ufam.edu.br](mailto:bruno.abia@icomp.ufam.edu.br)



