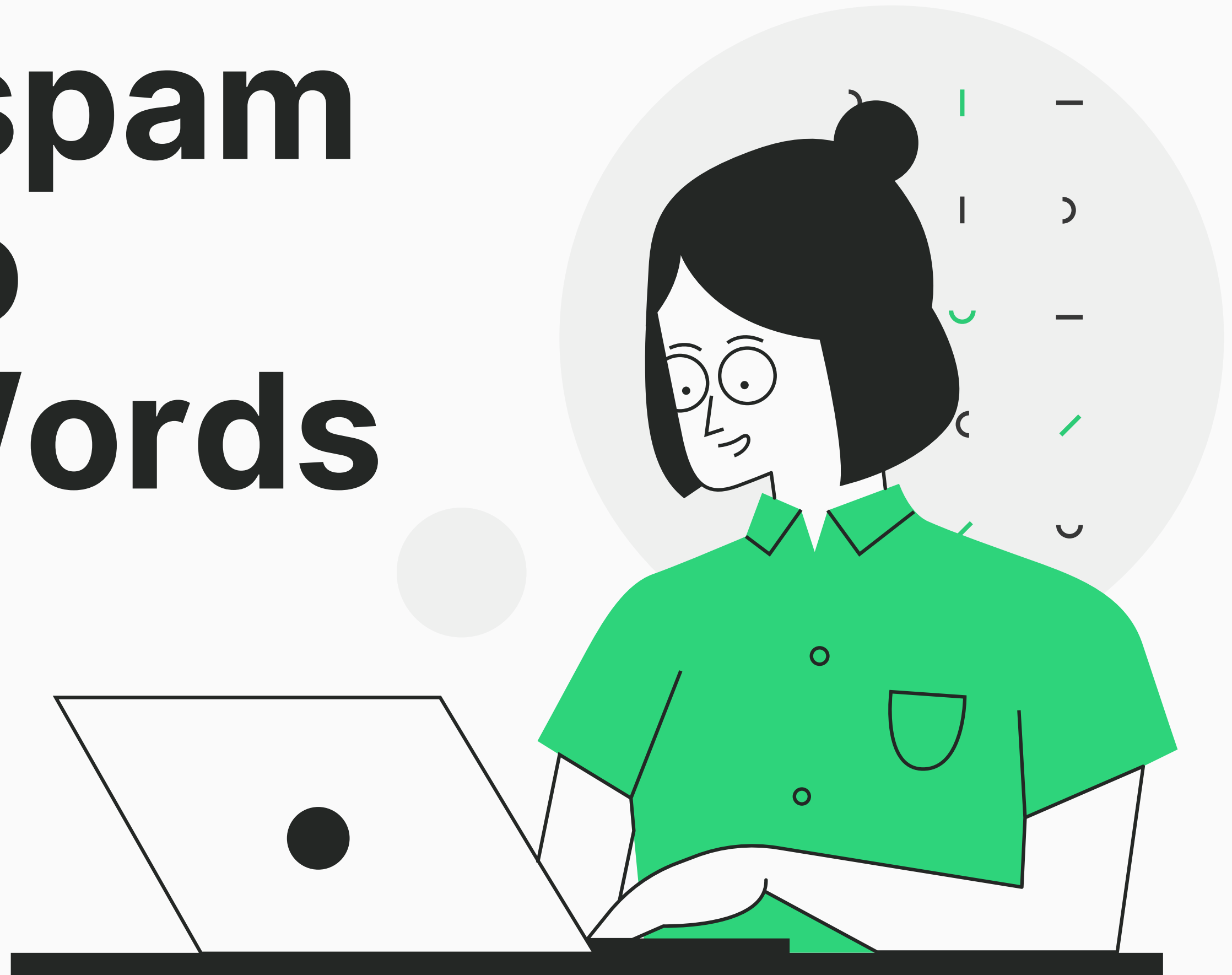


# Filtro de spam utilizando Bag-of-Words

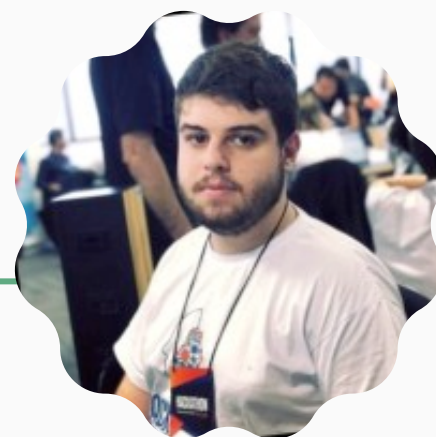
MACHINE LEARNING



# Grupo



**Higor  
Rufino  
Librelato**



**Marcos  
Junior C.  
Sebastião**



**Matheus  
Sant'ana  
Pacheco**



**Max Willian  
Trajano  
Martins**



**Vinicius de  
Lima  
Xavier**

# Filtro de Spam

O filtro de spam possui o comportamento como o de um **funil**. É um funil inteligente, **pois é composto por um conjunto de regras, fatores e algoritmos** que classificam a mensagem entrante como legítima ou não. O **alarme de um ou mais desses fatores** combinados definirá se a mensagem é de fato um spam.

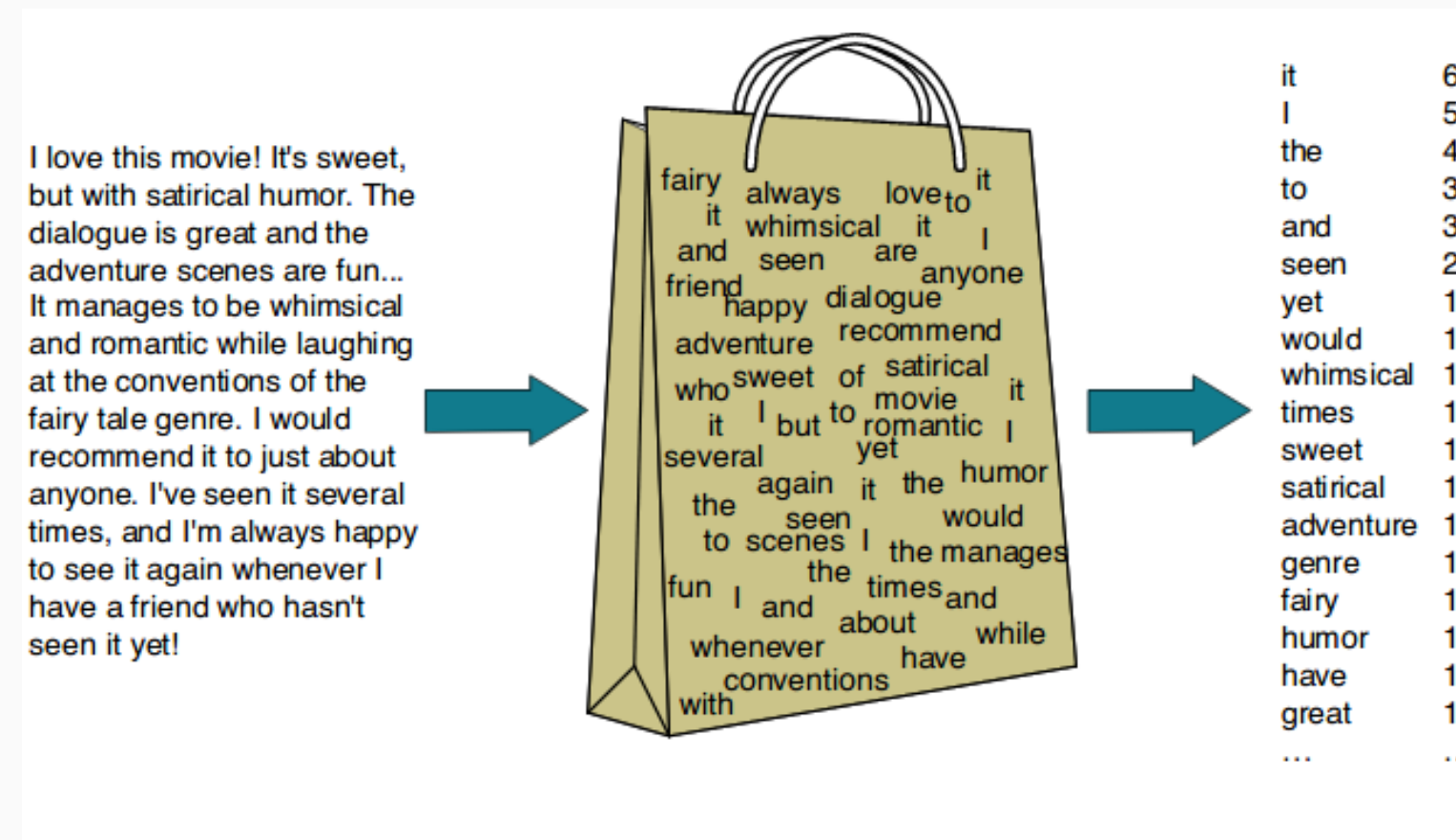
## Os tipos de filtro:

- Filtro de palavras;
- Filtro heurístico;
- Filtro bayesiano;
- Filtro de blacklist;
- Filtro de whitelist;
- Filtro de greylist.

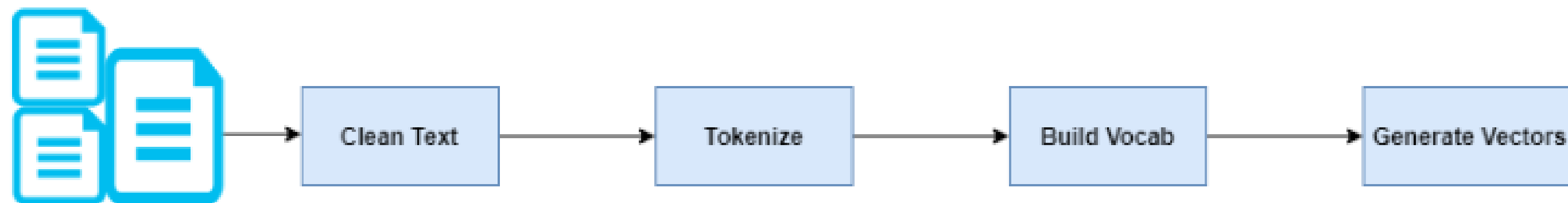


# Bag-of-words

Bag of Words (BOW – ou, em português, sacola de palavras) é um método para extrair recursos de documentos de texto. Esses recursos podem ser usados para treinar algoritmos de aprendizado de máquina. Ele cria um vocabulário de todas as palavras únicas que ocorrem em todos os documentos do conjunto de treinamento.



# Bag-of-words



Limpeza do texto -> Tokenizar -> Criar o vocabulário -> Gerar vetores

# Modelos de Teste

O resultado de um modelo de machine learning é dependente dos dados de entrada que o alimentam. Esses dados serão inseridos em uma equação, que retornará o resultado do modelo.

Assim, podemos ver que o que determina o desempenho de um modelo é o desenvolvimento desta equação, que acontece através do algoritmo utilizado.

- **KNN** (Aprendizagem supervisionada);
- **SVM** (Aprendizagem supervisionada);
- **Random Forest** (Aprendizagem supervisionada).

# Desafios

## ESCOLHA DO TEMA

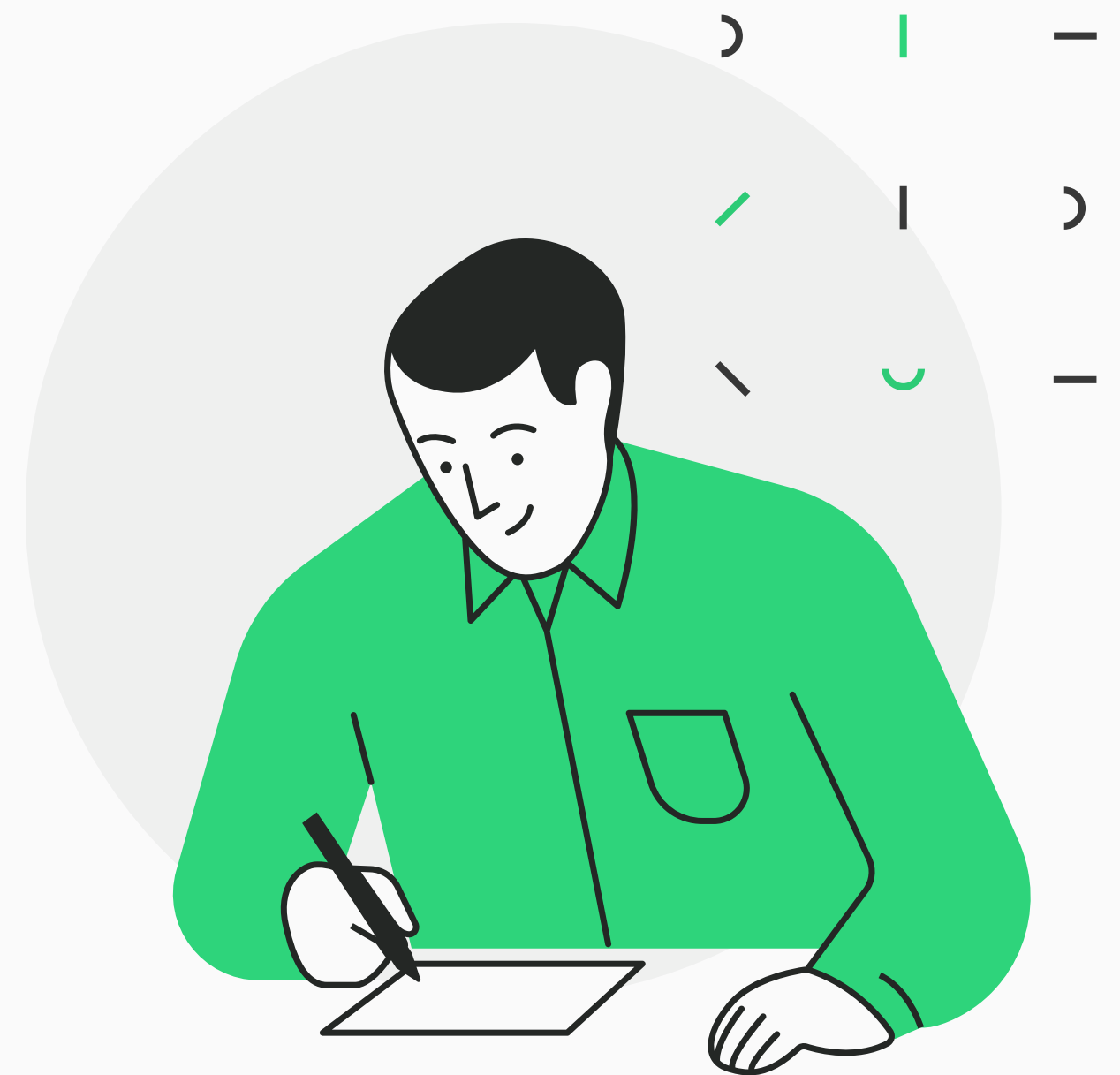
A escolha de um tema divergente do resto da turma e que agradasse a todos os participantes do grupo, foi um fator de divergência e debate.

## ENCONTRAR UM BOM DATASET PARA TREINO

Existem vários datasets para baixar, porem grande parte não estava sendo interessantes para a realização do treino e evolução do algoritmo.

## ESCOLHA DOS MODELOS

Encontramos diversos tipos de modelos para aplicação do modelo de treinamento, mas levamos um tempo para definir qual modelo aplicar.





# Obrigado!!

Dúvidas?



) | —

/ | )

\ ∪ —