



# Judge Panderson: Análise de leads com NLP e Detecção de Anomalias

Patrícia Pampanelli

Data Scientist

[patricia.pampanelli@grupozap.com](mailto:patricia.pampanelli@grupozap.com)

# Apresentação



Patrícia Pampanelli  
Data Scientist @ Grupo ZAP  
[patricia.pampanelli@grupozap.com](mailto:patricia.pampanelli@grupozap.com)



Quem é o grupo **ZAP** ?

**sua**house

O melhor e mais eficiente sistema de gestão, com CRM, plataforma de BI, controle de leads e muito mais.

data **ZAP**

Empresa de inteligência imobiliária que auxilia no processo de decisão para construção, investimento, financiamento e negociação.

fipe **ZAP**

Principal índice de preço de imóveis, de abrangência nacional. Parceria entre Fipe e o Grupo ZAP.

 zap

O portal especialista em imóveis.

 VivaReal

O maior portal de imóveis do Brasil.

conecta imobi

A plataforma que coloca o mercado imobiliário no centro de tudo.

**GEO**IMÓVEL

Plataforma de soluções e pesquisas imobiliárias com as mais diferenciadas ferramentas e instrumentos de análise mercadológica para as empresas que atuam no setor imobiliário

# Motivação



**Mais de 100.000  
mensagens de texto  
por dia!**

*"Sou corretora de Imóveis,  
tenho um cliente perfil para  
esse imóvel, se houver  
interesse em fazer  
intermediação por favor  
entrar em contato."*

*I read your advert and i am  
interested in the purchase. Can we  
discuss more on it? I can be reach  
on this ID: \*\*\*@outlook.com I speak  
only English Please get back to me  
with your Private E-Mail ID: Also  
Maria Olá, tenho interesse neste  
imóvel: Apartamento, 85m<sup>2</sup>, 2  
quartos, ...*

*"Olá, tenho interesse neste  
imóvel e proponho permuta por  
apartamento em Joinville em  
bairro nobre. Apartamento,  
71m<sup>2</sup>, 2 quartos, Rua Roberto  
Sell, 410 - Centro, Palhoça -  
SC, Venda, R\$ 248000.  
Aguardo o contato. Obrigado."*

grupo ZAP

# **Análise de mensagens de texto**



## Pré-processamento

Extração de features:  
frequência que o usuário gera os leads (hora, dia, semana)



Leads de Mensagem

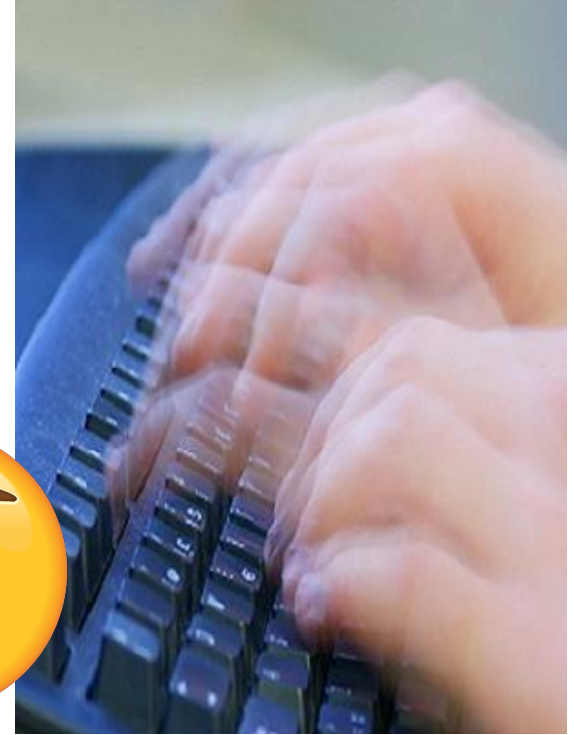
The diagram illustrates a data processing pipeline. On the left, a stack of three overlapping rectangular boxes represents the input data, labeled 'Leads de Mensagem'. A solid teal arrow points from this stack towards the right. On the right, a larger rectangular box with a dashed teal border represents the 'Pré-processamento' (Pre-processing) stage. Inside this box, the text describes the feature extraction process: 'Extração de features: frequência que o usuário gera os leads (hora, dia, semana)'. The overall flow is from left to right, indicating the transformation of raw message leads into processed features.

## Pré-processamento

Leads de  
Mensagem



Extração de  
features:  
frequência que  
o usuário gera  
os leads (hora,  
dia, semana)



## Pré-processamento

**Leads de  
Mensagem**



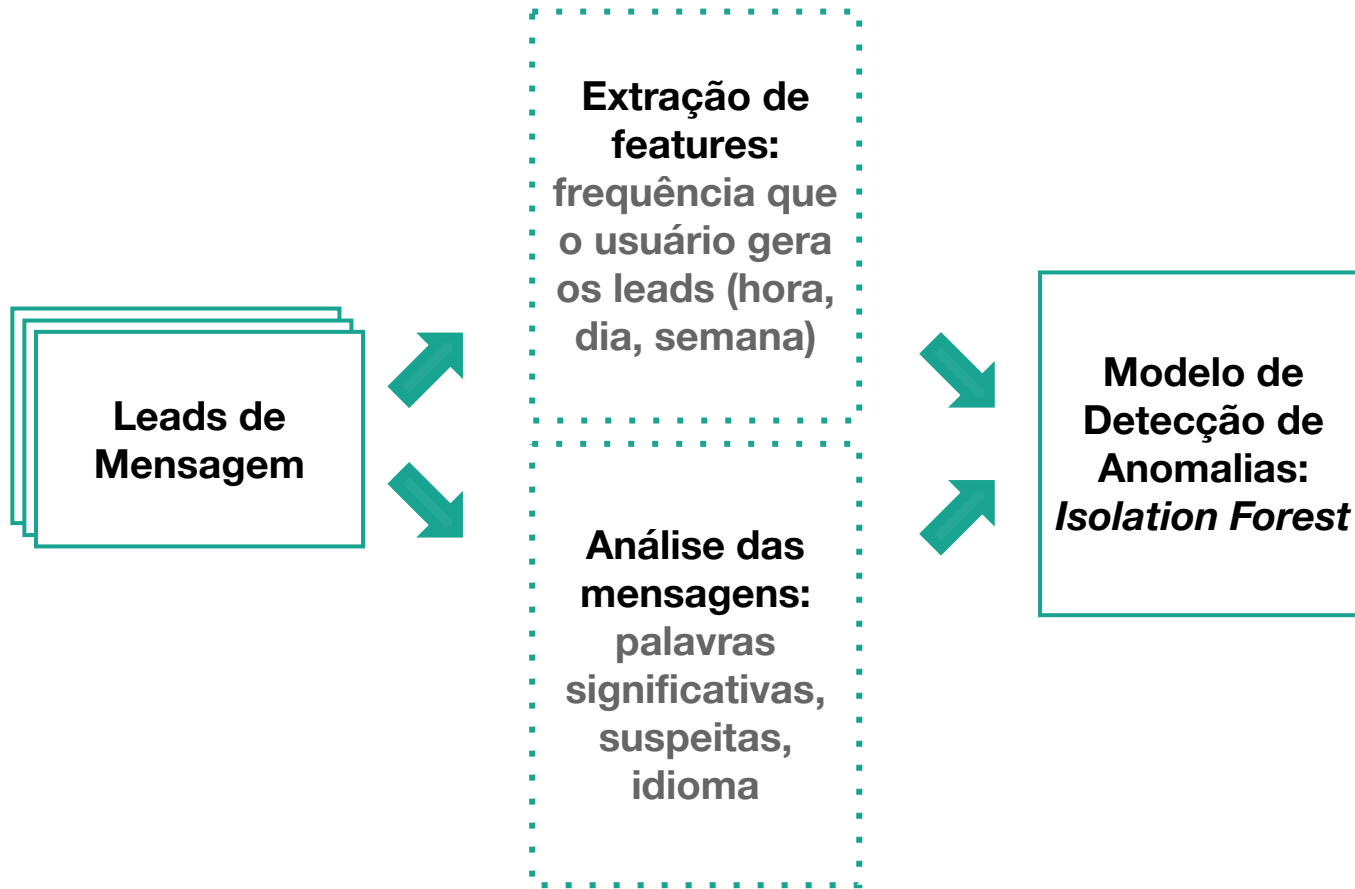
**Extração de  
features:**  
frequência que  
o usuário gera  
os leads (hora,  
dia, semana)

**Análise das  
mensagens:**  
palavras  
significativas,  
suspeitas,  
idioma





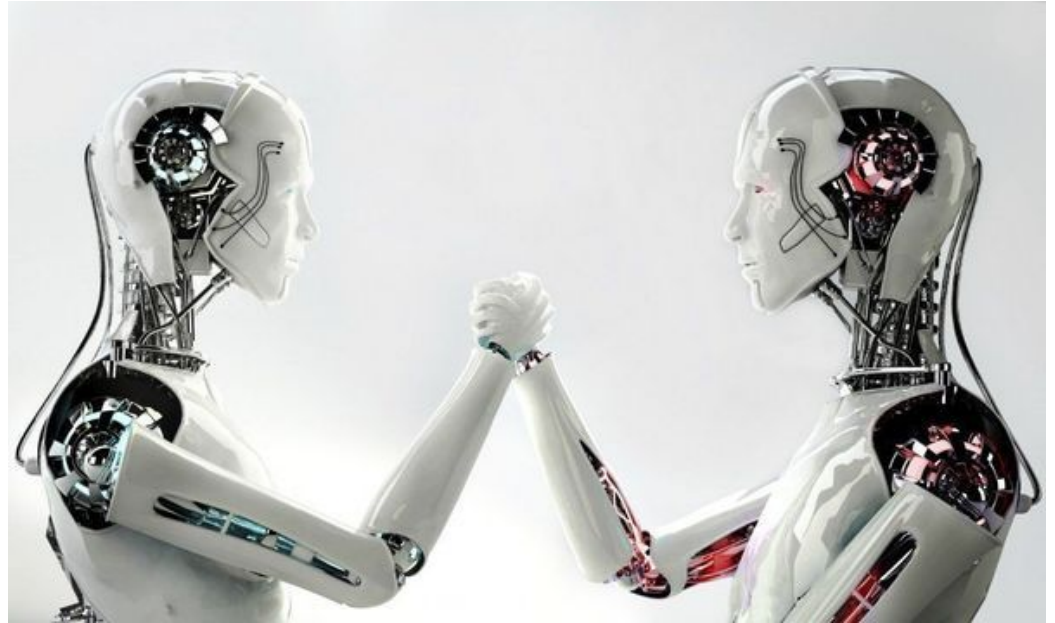
## Pré-processamento



# Detecção de anomalias com Isolation Forest

# Isolation Forest

- Método *ensemble* que funciona “isolando” as observações anômalas
- Métodos *ensembles*: conjunto de classificadores treinados para os quais os outputs são combinados



# Isolation Forest

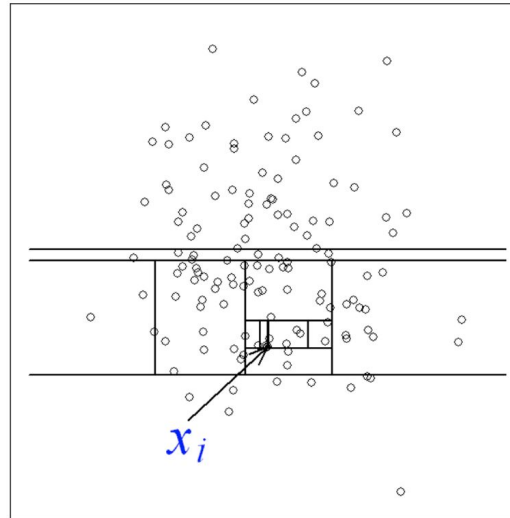
- Tem como base o método de árvores de decisão
- Define explicitamente quais são as amostras anômalas





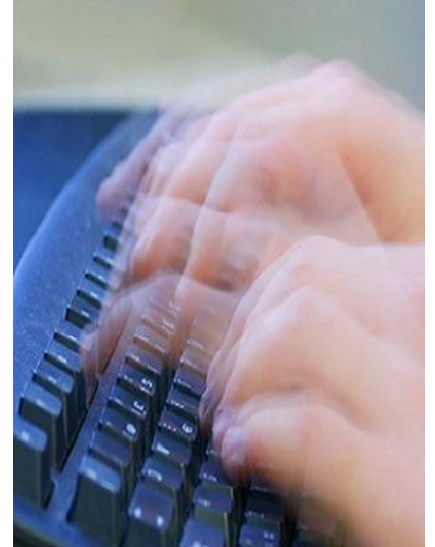
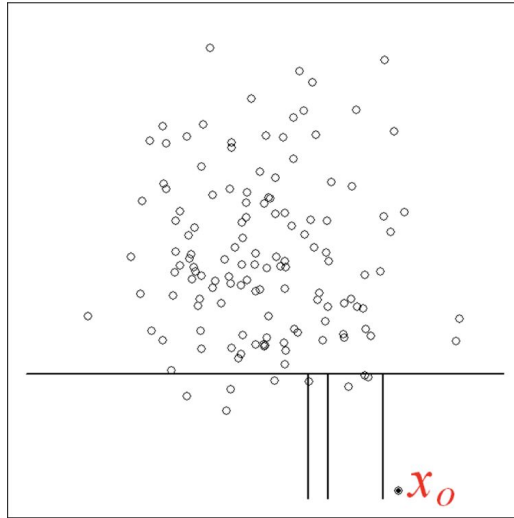
# Isolation Forest

- Amostras consideradas normais precisam de muitas partições para serem identificadas



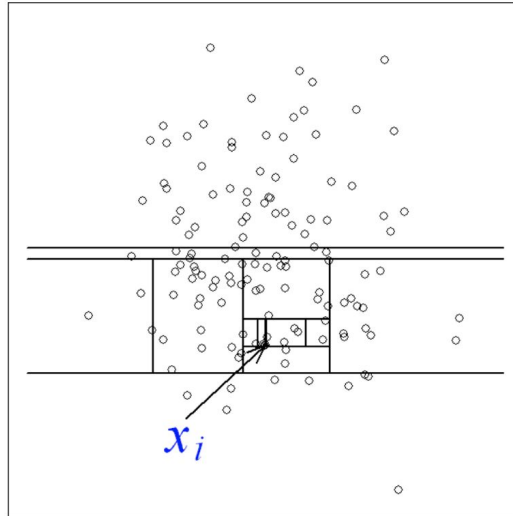
# Isolation Forest

- Amostras anômalas são identificadas por nós mais próximos da raiz da árvore

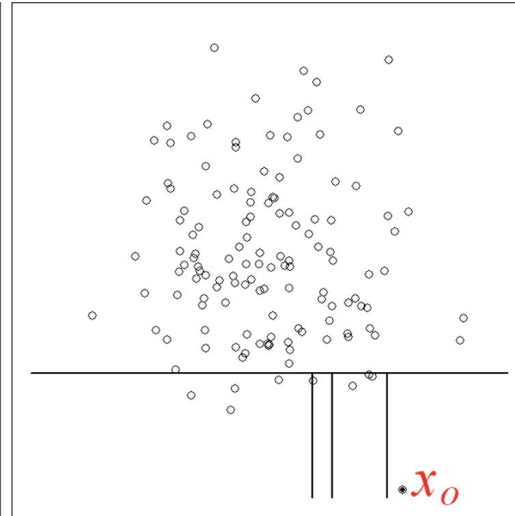


# Isolation Forest

Score  $\approx 0$ : normais

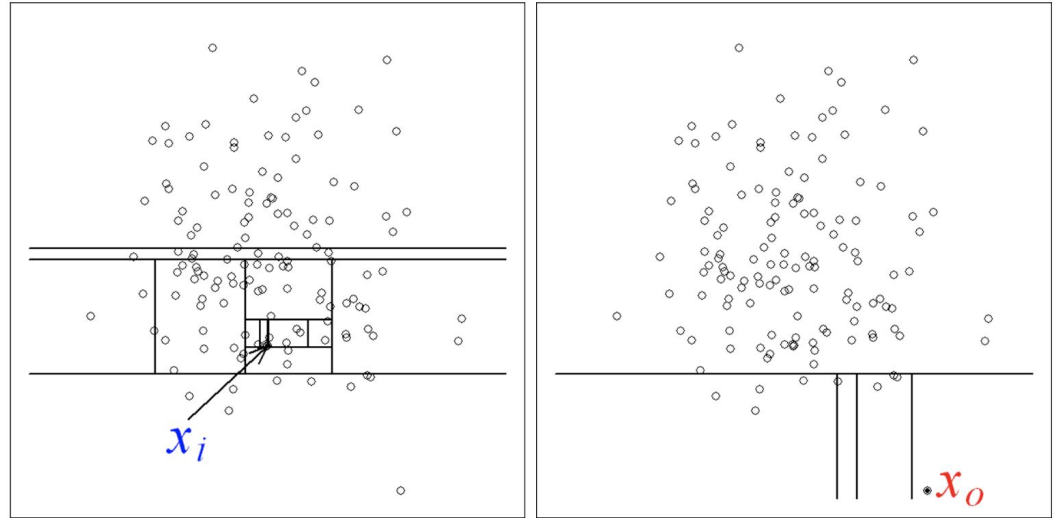


Score  $\approx 1$ : anomalias



# Isolation Forest

- **Parâmetro de contaminação:**
  - entre 0 e 0.5
  - representa a porcentagem de amostras que possivelmente são *outliers*



```
from sklearn.ensemble import IsolationForest
```

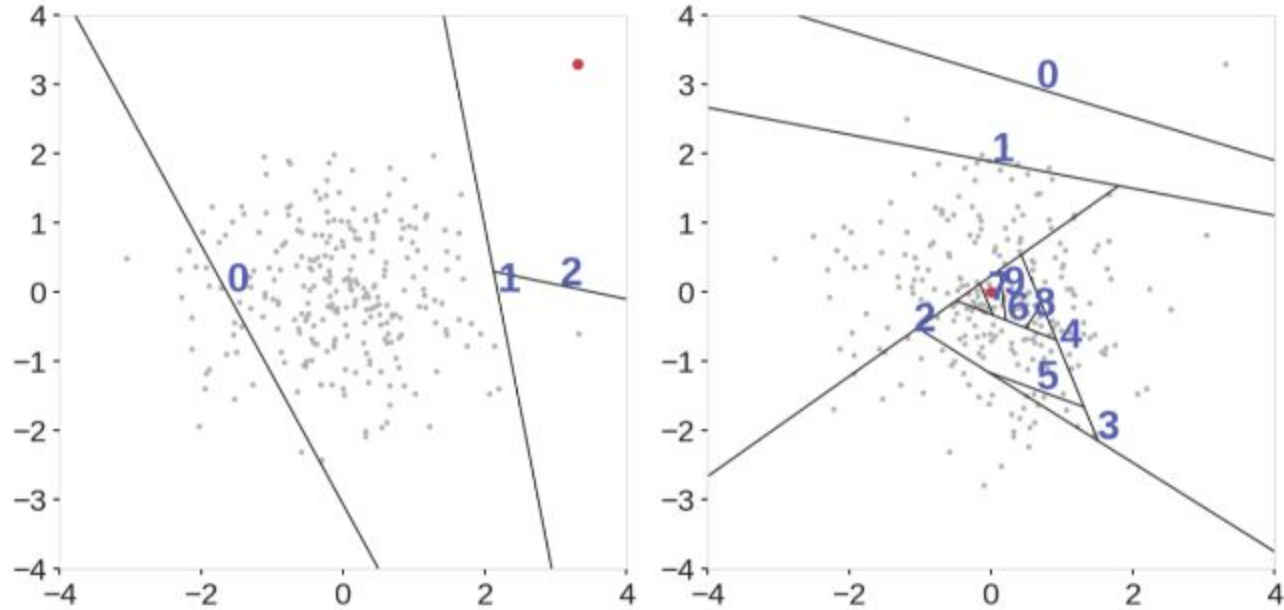


```
from sklearn.ensemble import IsolationForest
```

```
model = IsolationForest(n_estimators=300, max_features=0.7, contamination=0.015)
```



# Extended Isolation Forest

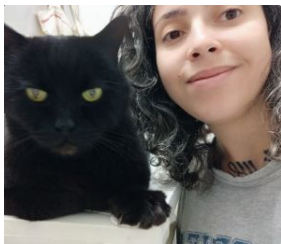


## Pré-processamento





# Ninguém faz nada sozinho...



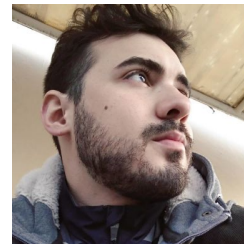
**Ana Carolina**



**Audrey Ferreira**



**Daniel Vainsencher**



**Vinícius Silva**



**Fábio Zacarias**



**Guilherme Bertola**



**Renan Barbioni**

# Classificação de leads

- **BLOCKED:** leads que contém palavras ou links proibidos
- **AGENT:** leads de corretor
- **NOT\_IN\_PORTUGUESE:** leads que não estão em português
- **UNDEFINED:** menos de X leads do mesmo usuário e em menos de X minutos
- **SUSPECT:** leads gerados em uma frequência que o modelo identifica como sendo bot;
- **EXCHANGE:** leads gerados com a intenção de fazer troca/permuta de imóveis
- **NORMAL:** leads considerados normais
- **ETC...**

# Obrigada!

**Estamos contratando!**  
**[jobs.kenoby.com/grupozap](https://jobs.kenoby.com/grupozap)**

