

Projeto de Estatística

# Uso do Teorema de Bayes para Classificação de Alunos com Alto Risco de Evasão Escolar



# Nossa Equipe:



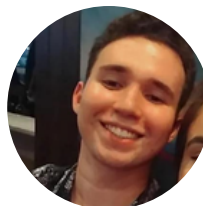
Ênio Henrique Nunes Ribeiro



Filipe Maciel Leicht



Matheus Augusto Monte Silva



Thiago José Grangeiro Costa



Victória Xavier Queiroz

# Fases da Pesquisa

**1**

Entendendo o Problema

**2**

Processamento da Base de Dados

**3**

Análise Estatística dos Dados

**4**

Construção e Treinamento do Modelo

**5**

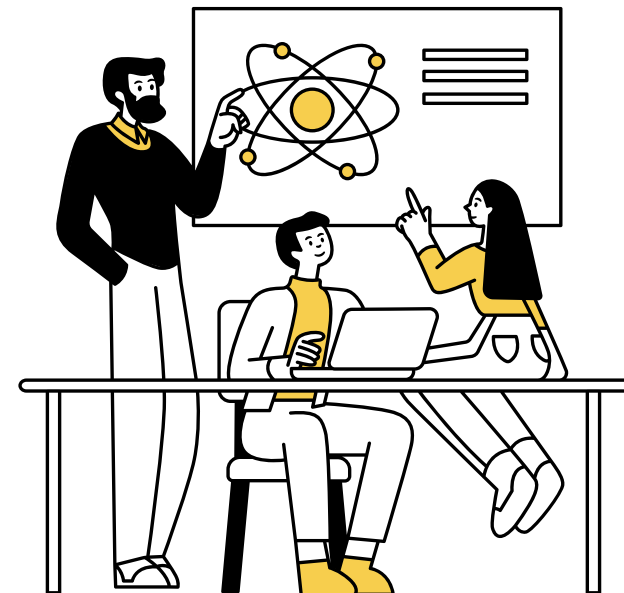
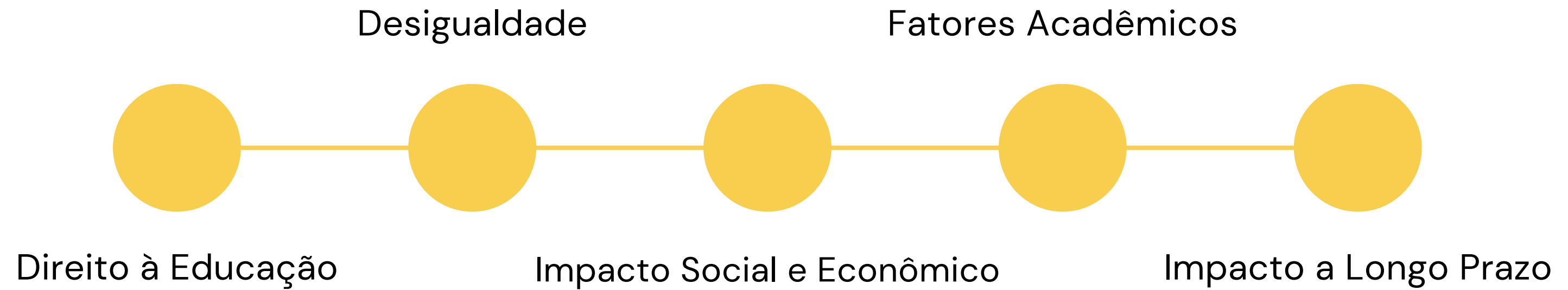
Conclusões

1

# Entendendo o Problema



# Contexto



2

# Processamento da Base de Dados



# Fases de Manipulação da Base de Dados

**1**

Obtenção da Base de Dados

**2**

Filtro de Instâncias

**3**

Engenharia de atributos

**4**

Seleção de Atributos

# Obtenção da Base de Dados



## Predict students' dropout and academic success

Donated on 12/12/2021

A dataset created from a higher education institution (acquired from several disjoint databases) related to students enrolled in different undergraduate degrees, such as agronomy, design, education, nursing, journalism, management, social service, and...

▼

Dataset Characteristics	Subject Area	Associated Tasks
Tabular	Other	Classification
Attribute Type	# Instances	# Attributes
-	4424	36

### Creators

**Valentim Realinho**  
vrealinho@  
Instituto Politécnico de Portalegre

**Mónica Vieira Martins**  
mvmartins@ipportalegre.pt  
Instituto Politécnico de Portalegre

**Jorge Machado**  
jmachado@ipportalegre.pt  
Instituto Politécnico de Portalegre

**Luís Baptista**  
lmbt@ipportalegre.pt  
Instituto Politécnico de Portalegre

**Predict students' dropout and academic success**





# Atributos da Base de Dados

## Classes e seus Atributos

1

**Dados Demográficos**

Status civil, nacionalidade, Gênero, Displaced, Intercambista, Idade no ano de entrada

2

**Dados Socioeconômicos**

Ocupação e Qualificação dos pais, Necessidades Educacionais Especiais, Devedor, Mensalidade em dia, Bolsista

3

**Dados Macroeconômicos**

Produto Interno Bruto, Taxa de Inflação, Taxa de Desemprego

4

**Dados Acadêmicos**

Curso, Turno, Modo de Entrada, Ordem de Escolha, Qualificações Anteriores, Nota da Qualificação Anterior, Dados Curriculares do 1º e 2º Semestres

# Pré-Processamento da Base de Dados

## Filtro de Instâncias

Na base de dados há um atributo chamado "**Target**", referente ao status do estudante em relação à Universidade, o qual poderia ter 3 classificações: **Graduado**, **Desistente** e **Matriculado**. Como o estudo busca saber se o aluno concluiu o curso de fato ou não, a classificação **Matriculado** não é interessante para nós e, por isso, foi deixada de fora.

ANTES → DEPOIS

Target	Target
Dropout	Dropout
Graduate	Graduate
Enrolled	Dropout
Graduate	Graduate
Dropout	Graduate



# Pré-Processamento da Base de Dados

## Engenharia de Atributos

1

**Simplificação do Nome de Alguns Atributos**



Remover caracteres especiais para otimizar o nome das variáveis em geral

2

**Agrupamento de Dados**



Tornar os dados mais familiares para o entendimento do Naive Bayes

3

**Categorizar Variáveis**



Metrificar dados. Neste caso, notas, em 5 categorias.

4

**Unificar Atributos**



Criar um parâmetro único para classificar a performance do aluno em um semestre, a partir da fusão das colunas referentes a desempenho acadêmico

# Pré-Processamento da Base de Dados

## Seleção de Atributos

1	Modo de Entrada	5	Notas (Qualificação Anterior)	9	Gênero
2	Curso	6	Necessidades Educacionais Especiais	10	Idade no Ano de Entrada
3	Turno	7	Devedor	11	Performance no 1º Semestre
4	Qualificações Anteriores	8	Mensalidade em Dia	12	Performance no 2º Semestre

3

# Análise Estatística dos Dados



# Estudo Exploratório

[View in Power BI](#)



3630

Matriculados



1421

Abandonos

39,15%

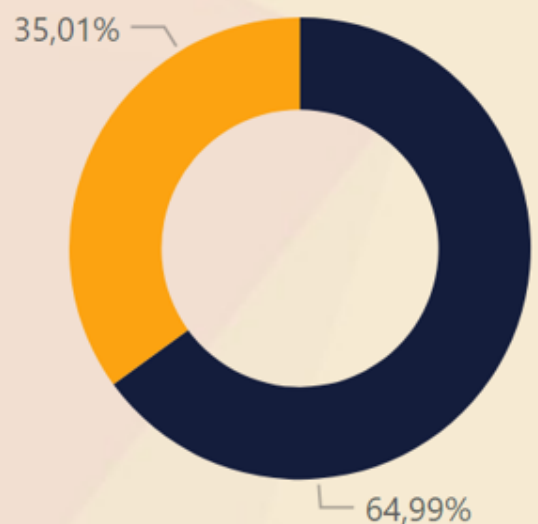


2209

Graduados

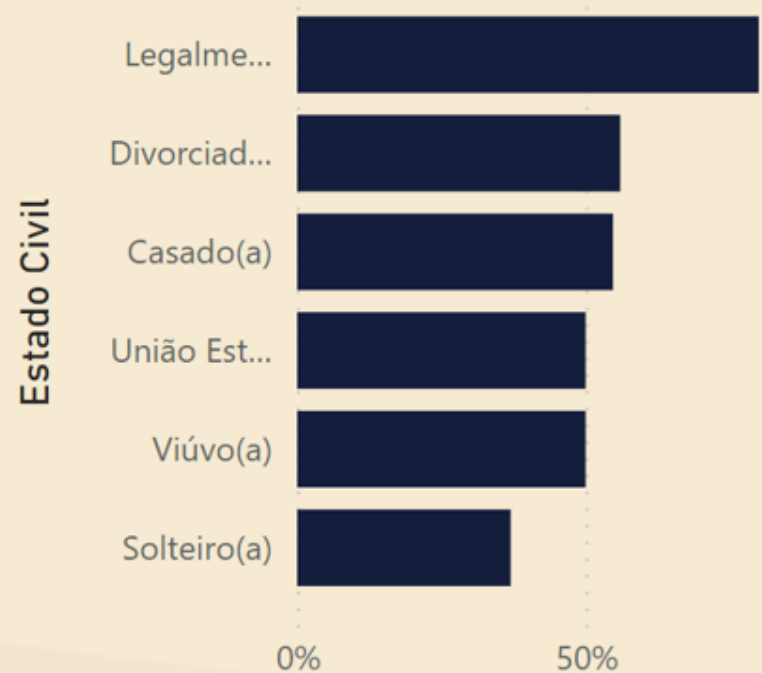
60,85%

(%) Abandonos por Sexo



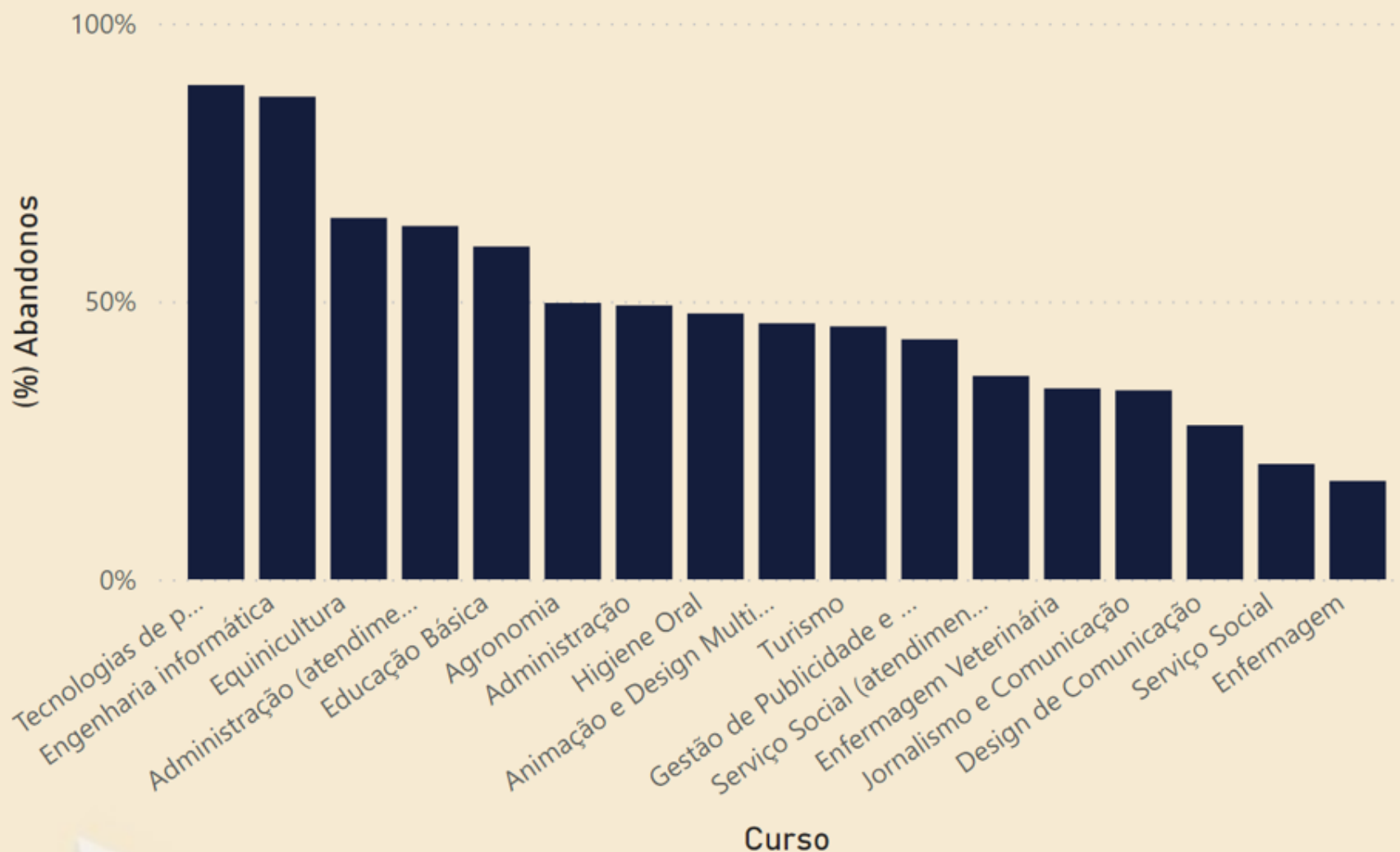
Sexo ● Masculino ● Feminino

(%) Abandonos por Estado Civil

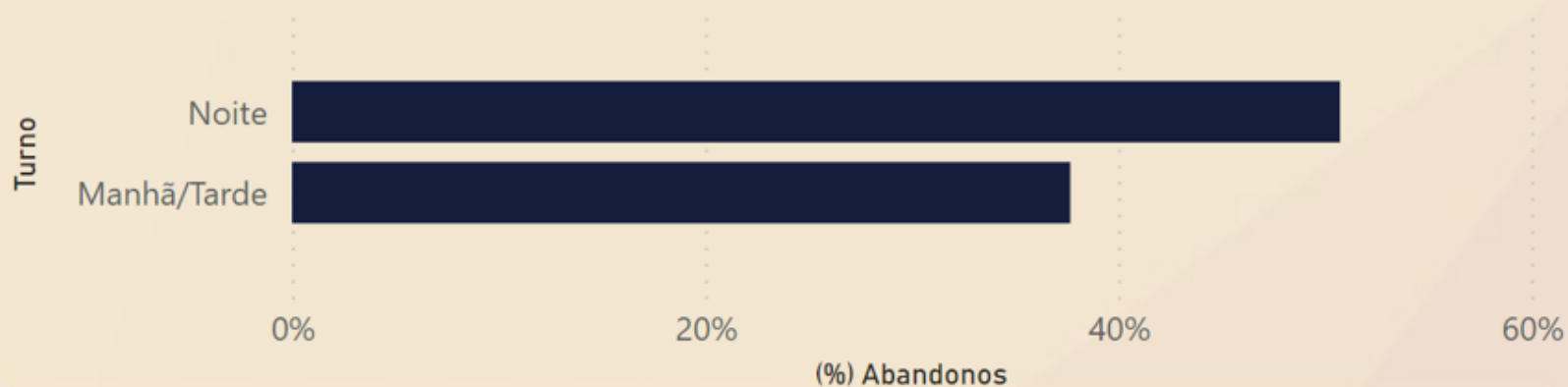


(%) Abandonos

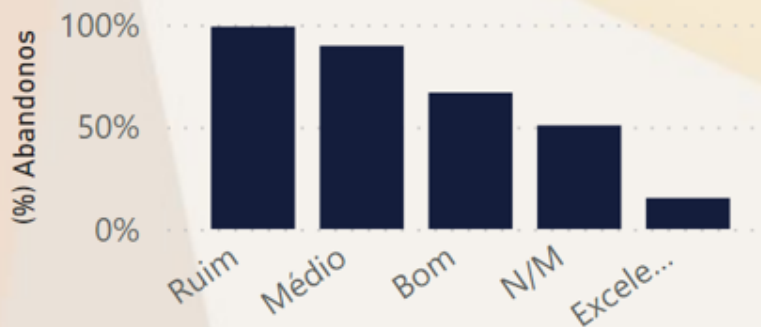
(%) Abandonos por Curso



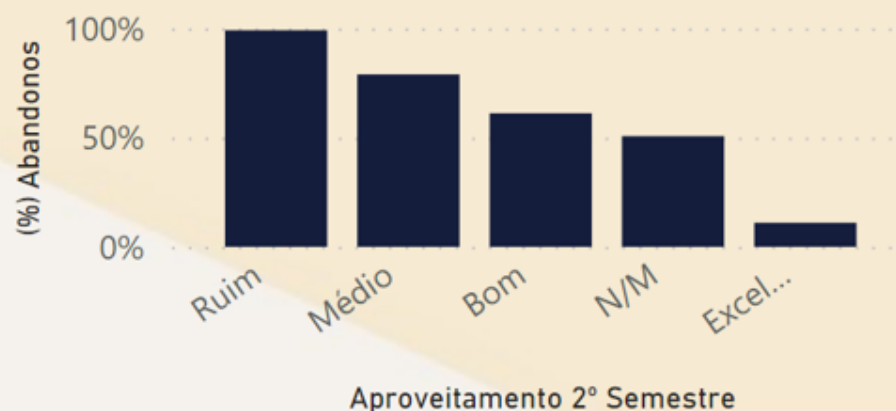
(%) Abandonos por Turno



(%) Abandonos por Aproveitamento 1º Semestre



(%) Abandonos por Aproveitamento 2º Semestre



4

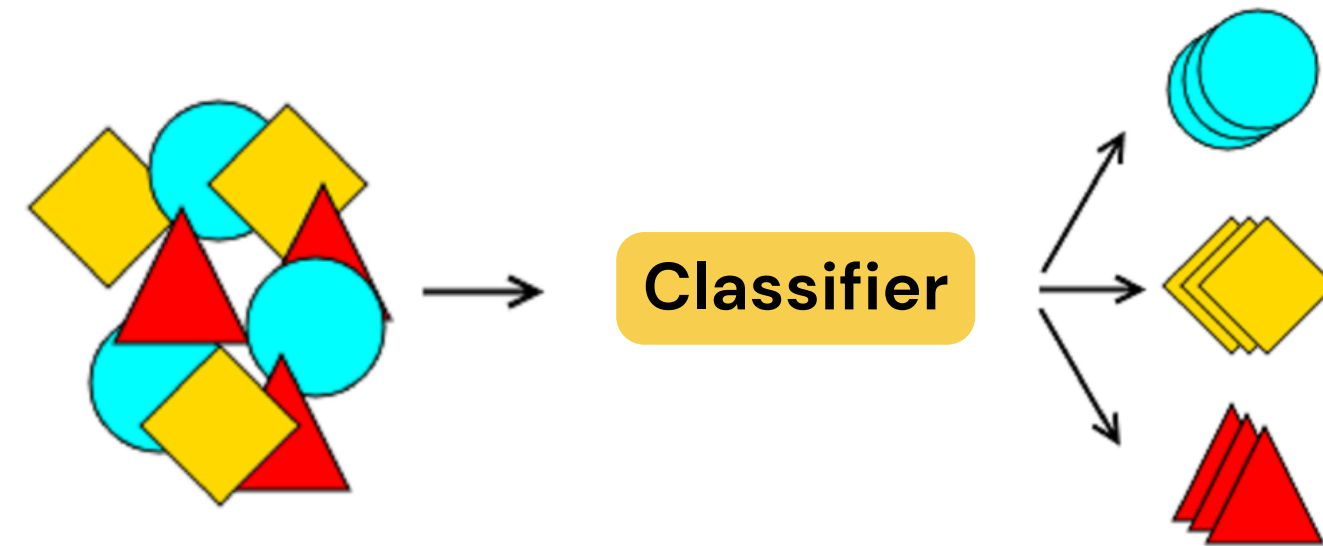
# Construção e Treinamento do Modelo



# Naive Bayes

Algoritmo de classificação probabilístico baseado no **Teorema de Bayes**

- 1 Gaussian Naive Bayes
- 2 Multinomial Naive Bayes
- 3 Complement Naive Bayes
- 4 Bernoulli Naive Bayes
- 5 Categorical Naive Bayes



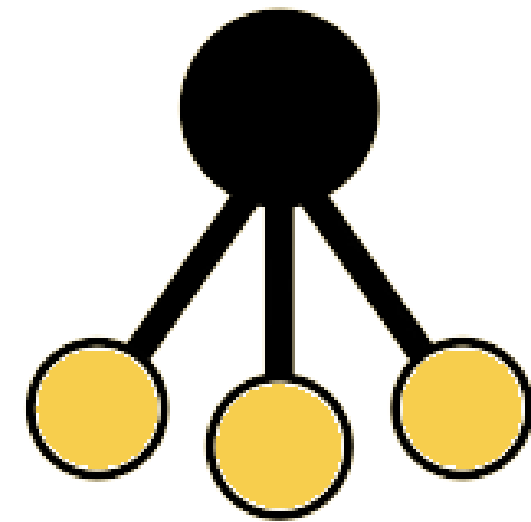


# Categorical Naive Bayes

É uma variação do algoritmo **Naive Bayes** projetada especificamente para lidar com dados categóricos

$$P(x_i = t \mid y = c; \alpha) = \frac{N_{tic} + \alpha}{N_c + \alpha n_i}$$

Por que escolhemos o  
Categorical Naive Bayes?



# Implementação

Nossa implementação do algoritmo aplicado à base de dados estudada:



[Google Colaboratory](#)



5

# Conclusões

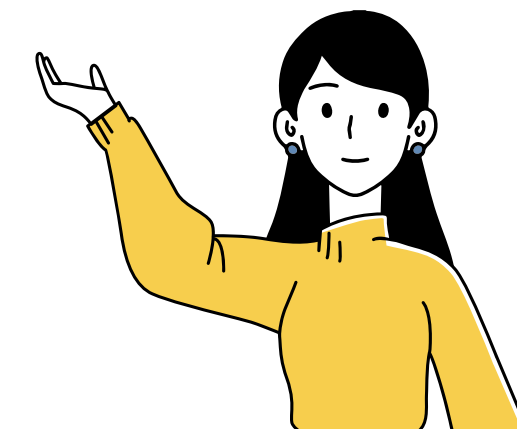
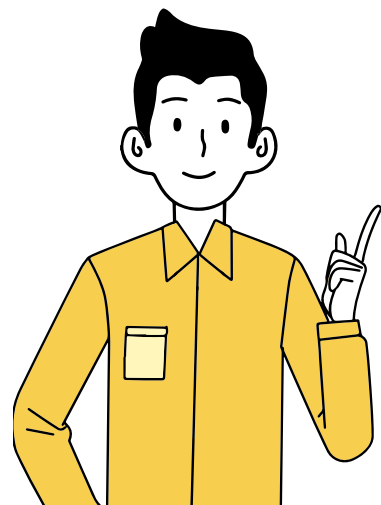


# Nossas Conclusões

Notamos a necessidade de alterar a forma como os dados eram apresentados para **maximizar a precisão** do modelo utilizado

Percebemos que analisar os dados relativos à vida de um estudante para prever se ele vai ou não abandonar o curso nos leva a **conclusões imprecisas**

O dado mais impactante no nosso modelo é o **desempenho do semestre do estudante**, o resto não apresentou impacto significativo



# Obrigado pela Atenção!

