



Introdução à Inteligência Artificial

Módulo 03: Aprendizado de Máquina (*Machine Learning*, ML)

Parte 1 - Introdução ao Aprendizado de Máquina

Professor: Rafael de Magalhães Dias Frinhani



UNIFEI



Softex



MCTI
FUTURO

FUTURO DO TRABALHO, TRABALHO DO FUTURO



Sumário do Módulo 03

Parte 1 – Introdução ao Aprendizado de Máquina

1. Conceitos
2. Tipos de Aprendizado de Máquina
3. Processo de Aprendizado de Máquina
4. Parâmetros e Hiperparâmetros do Modelo



UNIFEI



Softex



MCTI
FUTURO

FUTURO DO TRABALHO. TRABALHO DO FUTURO



Objetivos do Módulo

- Este módulo tem por objetivo apresentar os conceitos, os tipos e o processo de aprendizado de máquina.
- Entender os conceitos de objetos de dados e tipos de atributos, bem como as estruturas de dados usadas para o seu armazenamento.

Objetivos do Módulo

- Este módulo tem por objetivo apresentar os conceitos, os tipos e o processo de aprendizado de máquina.
- Entender os conceitos de objetos de dados e tipos de atributos, bem como as estruturas de dados usadas para o seu armazenamento.
- Conhecer os **algoritmos** clássicos de aprendizado de máquina **não-supervisionado** e **supervisionado**.
- Apresentar os **conceitos** de **parâmetros**, **hiperparâmetros** e **validação cruzada**.

1. Aprendizado de Máquina – Conceitos

Machine Learning (ML) – campo de pesquisa da Inteligência Artificial, que foca no estudo de como as máquinas computacionais podem obter novos conhecimentos ou habilidades, realizando comportamentos de aprendizagem semelhantes aos de humanos para melhorar seu desempenho.

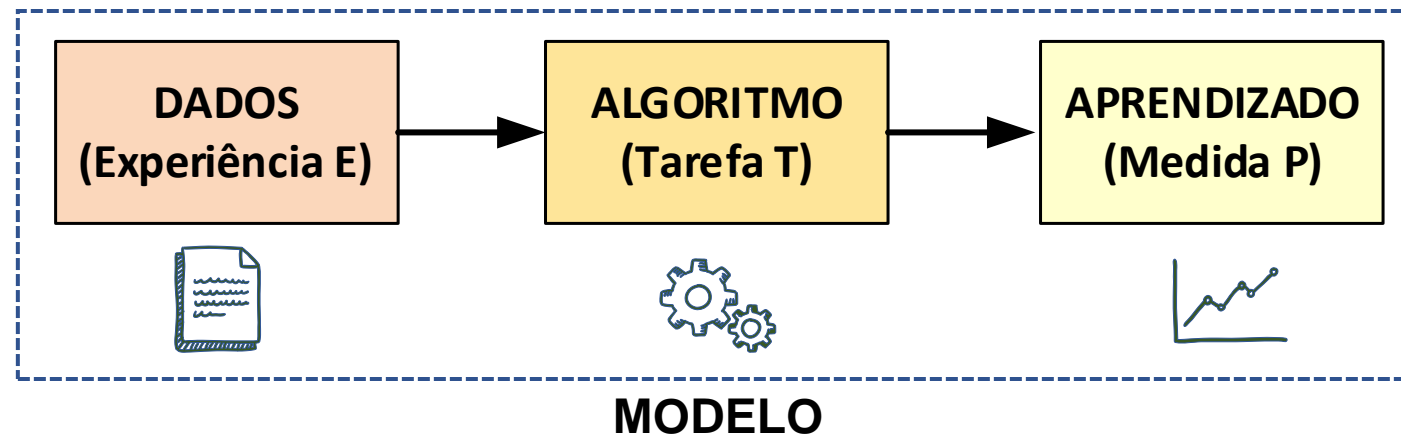
É um subcampo da Ciência da Computação interdisciplinar com a estatística, matemática computacional, teoria da informação e otimização. Em ML é dada maior ênfase aos algoritmos, considerando que a máquina em questão é o computador.

1. Aprendizado de Máquina – Conceitos

Machine Learning (ML) – campo de pesquisa da Inteligência Artificial, que foca no estudo de como as máquinas computacionais podem obter novos conhecimentos ou habilidades, realizando comportamentos de aprendizagem semelhantes aos de humanos para melhorar seu desempenho.

É um subcampo da Ciência da Computação interdisciplinar com a estatística, matemática computacional, teoria da informação e otimização. Em ML é dada maior ênfase aos algoritmos, considerando que a máquina em questão é o computador.

Conceitos centrais do *Machine Learning*:



1. Aprendizado de Máquina – Aplicações

Possui **aplicações em diversas áreas** auxiliando na solução de problemas na manufatura, no comércio e na prestação de serviços, sendo útil para **tarefas como**:

- tradução automática de idiomas
- diagnóstico médico
- taxonomia de espécies, classificação de perfis
- projeções financeiras
- detecção de fraude, filtragem de e-mails
- robôs, carros autônomos
- sistemas de recomendação
- previsões (climática, tráfego, financeira)
- reconhecimento facial e de fala, etc.

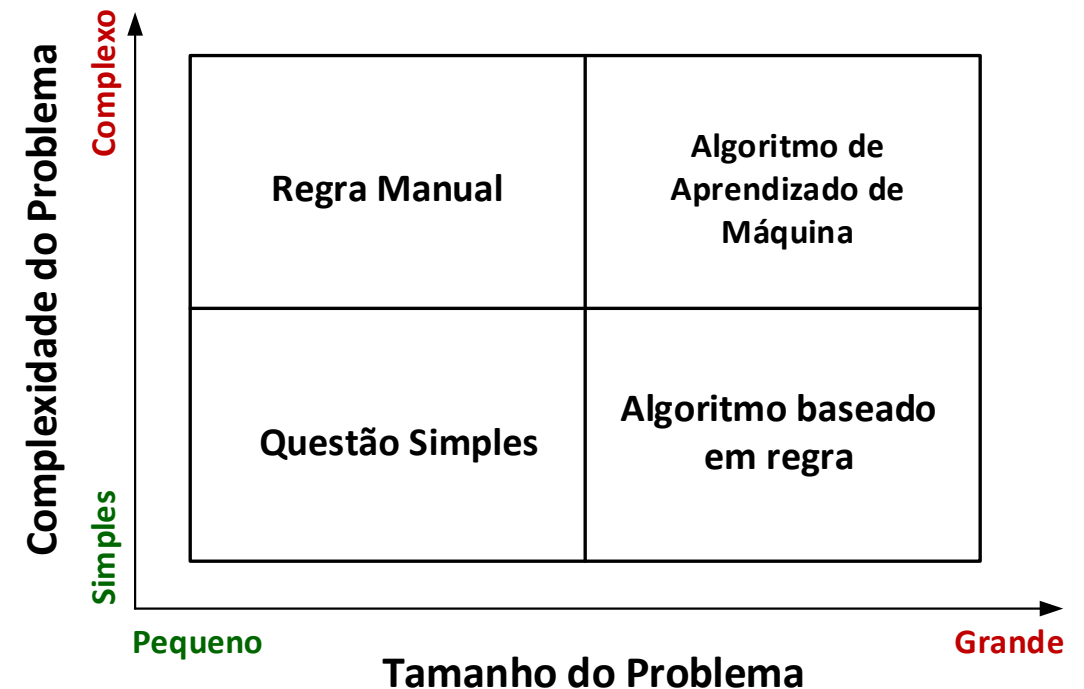
1. Aprendizado de Máquina – Aplicações

Possui aplicações em diversas áreas auxiliando na solução de problemas na manufatura, no comércio e na prestação de serviços, sendo útil para tarefas como:

- tradução automática de idiomas
- diagnóstico médico,
- taxonomia de espécies, classificação de perfis
- projeções financeiras
- detecção de fraude, filtragem de e-mails
- robôs, carros autônomos
- sistemas de recomendação
- previsões (climática, tráfego, financeira)
- reconhecimento facial e de fala, etc.

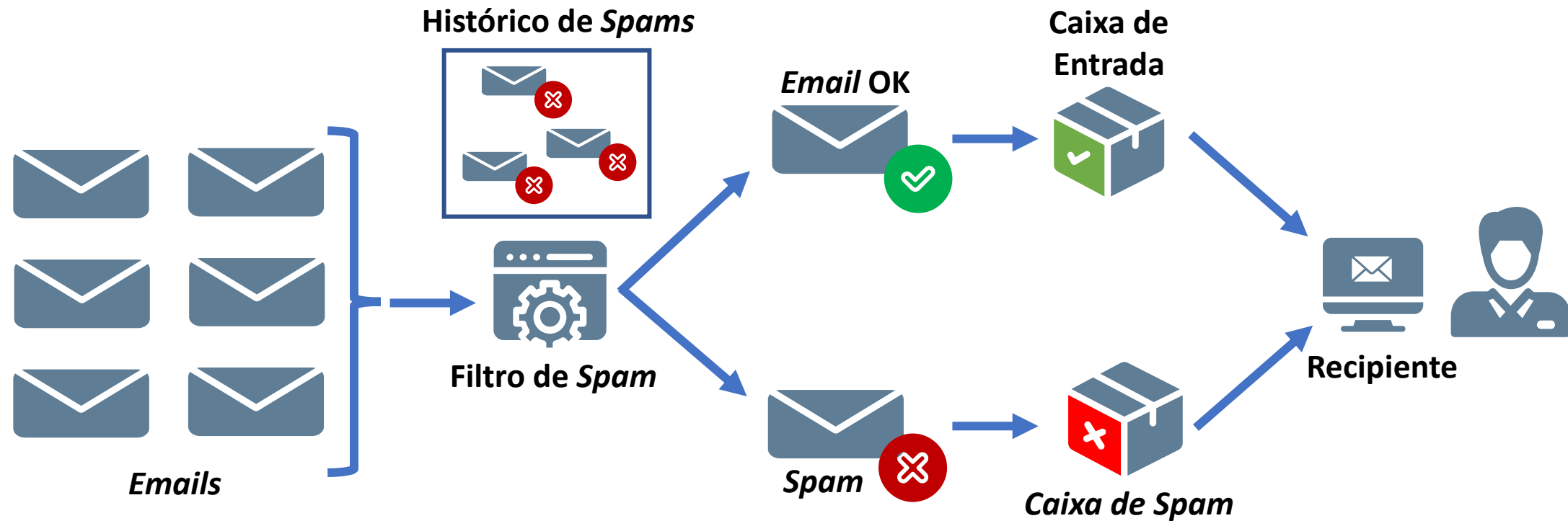
Justificativa para uso de ML

Mesmo sendo adaptável aos mais diversos tipos de problema, **deve-se analisar se seu uso é justificado** considerando a **complexidade** e o **tamanho do problema**.



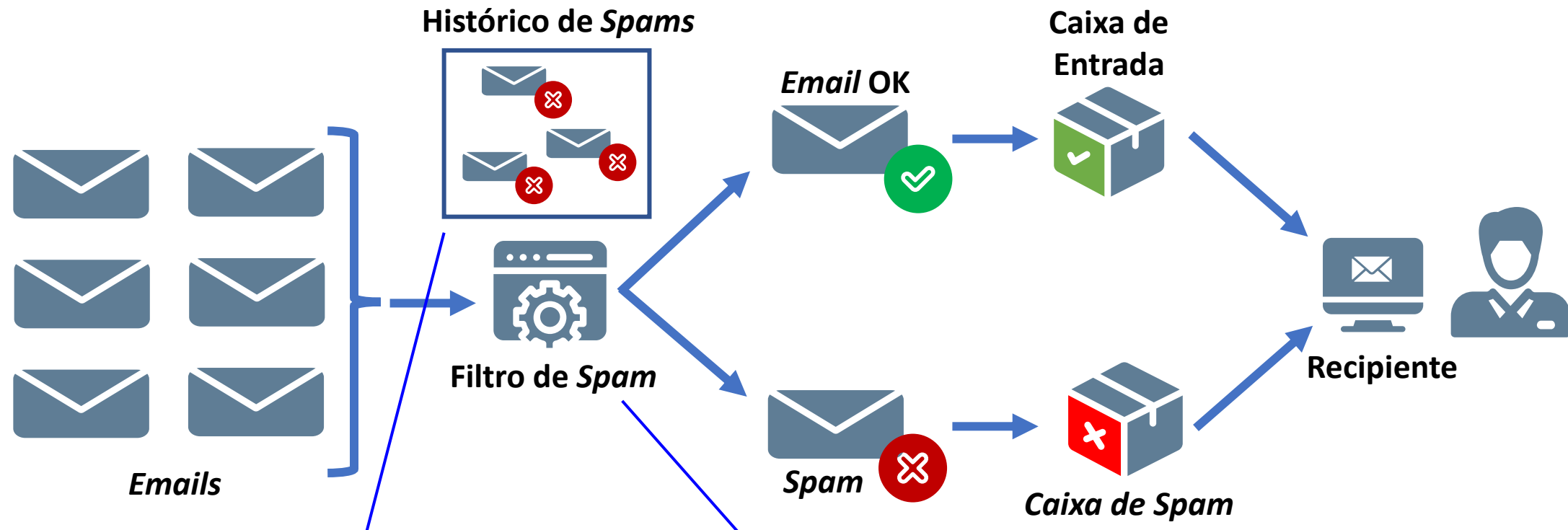
1. Aprendizado de Máquina – Exemplo

Sistema de classificação automática de lixo eletrônico (*spam*) em e-mails.



1. Aprendizado de Máquina – Exemplo

Sistema de classificação automática de lixo eletrônico (*spam*) em e-mails.

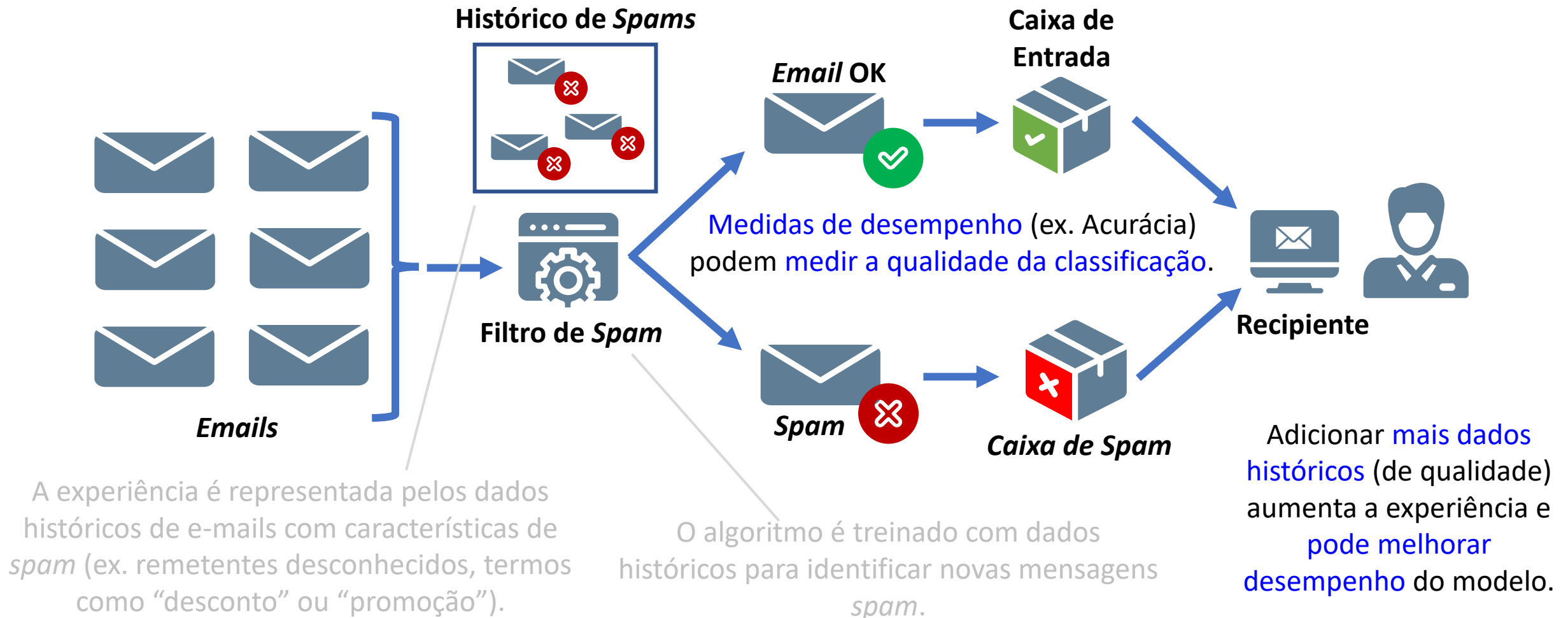


A experiência é representada pelos **dados históricos de e-mails** com características de *spam* (ex. remetentes desconhecidos, termos como “desconto” ou “promoção”).

O algoritmo é treinado com **dados históricos** para identificar novas mensagens *spam*.

1. Aprendizado de Máquina – Exemplo

Sistema de classificação automática de lixo eletrônico (*spam*) em e-mails.



1. Aprendizado de Máquina – Objetos e Atributos

Para ser solucionado por ML o problema precisa ser representado para ser tratado computacionalmente.

Objetos de Dados representam as **entidades** do problema, cujas **características** são representadas pelos **atributos** do objeto.

Um **conjunto de dados reúne** os **objetos de dados** do problema, sendo fornecido com entrada do algoritmo de ML (ex. treinamento, novos dados).

1. Aprendizado de Máquina – Objetos e Atributos

Para ser solucionado por ML o problema precisa ser representado para ser tratado computacionalmente.

Objetos de Dados representam as entidades do problema, cujas características são representadas pelos atributos do objeto.

Um conjunto de dados reúne os objetos de dados do problema, sendo fornecido com entrada do algoritmo de ML (ex. treinamento, novos dados).

Ex. Conjunto de dados de solicitantes de auxílio estudantil.

8 atributos (vetor de atributos)

10 objetos
(registros)

ID	MORADIA ALUNO	MORADIA FAMILIA	ESTADO CIVIL	FILHOS	RENDA/CAPITA	DESPESAS/CAPITA	FAMILIARES DOENÇA	PATRIMÔNIO
11323	Aluguel com Colegas	Heranca	Solteiro	2	166	73	0	0
11289	Aluguel com Colegas	Propria Quitada	Solteiro	0	582	155	0	0
23363	Familia/Terceiros	Propria em Pagamento	Divorciado	0	1118	56	1	15000
22148	Familia/Terceiros	Heranca	Solteiro	1	407	0	0	0
23060	Aluguel com Colegas	Propria Quitada	Casado	0	636	133	1	60000
30394	Familia/Terceiros	Alugada	Solteiro	0	533	53	0	45000
27836	Familia/Terceiros	Heranca	Solteiro	0	281	22	0	0
17099	Aluguel com Colegas	Propria Quitada	Casado	0	1238	54	0	45443
20042	Familia/Terceiros	Alugada	Solteiro	0	400	200	0	0
29757	Aluguel Sozinho	Propria Quitada	Solteiro	0	1227	53	0	7000

1. Aprendizado de Máquina – Objetos e Atributos

A nível computacional um conjunto de dados pode ser armazenado numa Matriz de Dados.

Matriz Objeto–Atributo

Relaciona os objetos (linhas) com seus atributos (colunas), constituindo uma matriz de dimensão $n \times p$.

$$M = \begin{bmatrix} x_{11} & x_{12} & x_{1j} & \dots & x_{1p} \\ x_{21} & x_{22} & x_{2j} & \dots & x_{2p} \\ x_{i1} & x_{i2} & x_{ij} & \dots & x_{ip} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & x_{nj} & \dots & x_{np} \end{bmatrix}$$

1. Aprendizado de Máquina – Objetos e Atributos

A nível computacional um conjunto de dados pode ser armazenado numa Matriz de Dados.

Matriz Objeto–Atributo

Relaciona os objetos (linhas) com seus atributos (colunas), constituindo uma matriz de dimensão $n \times p$.

$$M = \begin{bmatrix} x_{11} & x_{12} & x_{1j} & \dots & x_{1p} \\ x_{21} & x_{22} & x_{2j} & \dots & x_{2p} \\ x_{i1} & x_{i2} & x_{ij} & \dots & x_{ip} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & x_{nj} & \dots & x_{np} \end{bmatrix}$$

Matriz Objeto–Objeto

Relaciona os objetos (linhas) com outros objetos (colunas), constituindo uma matriz de dimensão $n \times n$. Geralmente obtida a partir de uma matriz Objeto–Atributo.

$$D = \begin{bmatrix} 0 & & & & & \\ d(2,1) & 0 & & & & \\ d(3,1) & d(3,2) & 0 & & & \\ d(4,1) & d(4,2) & d(4,3) & 0 & & \\ \vdots & \vdots & \vdots & \vdots & \ddots & \\ d(n,1) & d(n,2) & d(n,3) & d(n,3) & \dots & 0 \end{bmatrix}$$

1. Aprendizado de Máquina – Objetos e Atributos

O tipo do atributo é determinado pelo tipo de dado. Atributos de um mesmo tipo podem possuir limites de valores diferentes. **Tipos de atributos comuns:**

Nominais: nomes (**rótulos**) cujos valores **indicam categoria**, código ou estado. **Ex.** “cor do cabelo” pode ser preto, loiro, ruivo etc. Podem ser números (ex. 0 = preto), mas não aceitam operações matemáticas.

1. Aprendizado de Máquina – Objetos e Atributos

O tipo do atributo é determinado pelo tipo de dado. Atributos de um mesmo tipo podem possuir limites de valores diferentes. **Tipos de atributos comuns:**

QUALITATIVOS

Nominais: nomes (rótulos) cujos valores indicam categoria, código ou estado. **Ex.** “cor do cabelo” pode ser preto, loiro, ruivo etc. Podem ser números (ex. 0 = preto), mas não aceitam operações matemáticas.

Binários: tipo especial de atributo nominal com **apenas duas categorias** (classes): 0 ou 1. **Ex.** Paciente 0 = sem comorbidade, 1 = com comorbidade etc.

Ordinais: **valores possuem ordem significativa entre si**, mas a **magnitude** entre eles **é desconhecida**. Úteis para representar dados subjetivos. **Ex.** 0: muito insatisfeito, 1: um pouco insatisfeito, etc.

1. Aprendizado de Máquina – Objetos e Atributos

O tipo do atributo é determinado pelo tipo de dado. Atributos de um mesmo tipo podem possuir limites de valores diferentes. **Tipos de atributos comuns:**

QUALITATIVOS

Nominais: nomes (rótulos) cujos valores indicam categoria, código ou estado. **Ex.** “cor do cabelo” pode ser preto, loiro, ruivo etc. Podem ser números (ex. 0 = preto), mas não aceitam operações matemáticas.

Binários: tipo especial de atributo nominal com apenas duas categorias (classes): 0 ou 1. **Ex.** Paciente 0 = sem comorbidade, 1 = com comorbidade etc.

Ordinais: valores possuem ordem significativa entre si, mas a magnitude entre eles é desconhecida. Úteis para representar dados subjetivos. **Ex.** 0: muito insatisfeito, 1: um pouco insatisfeito, etc.

QUANTITATIVOS

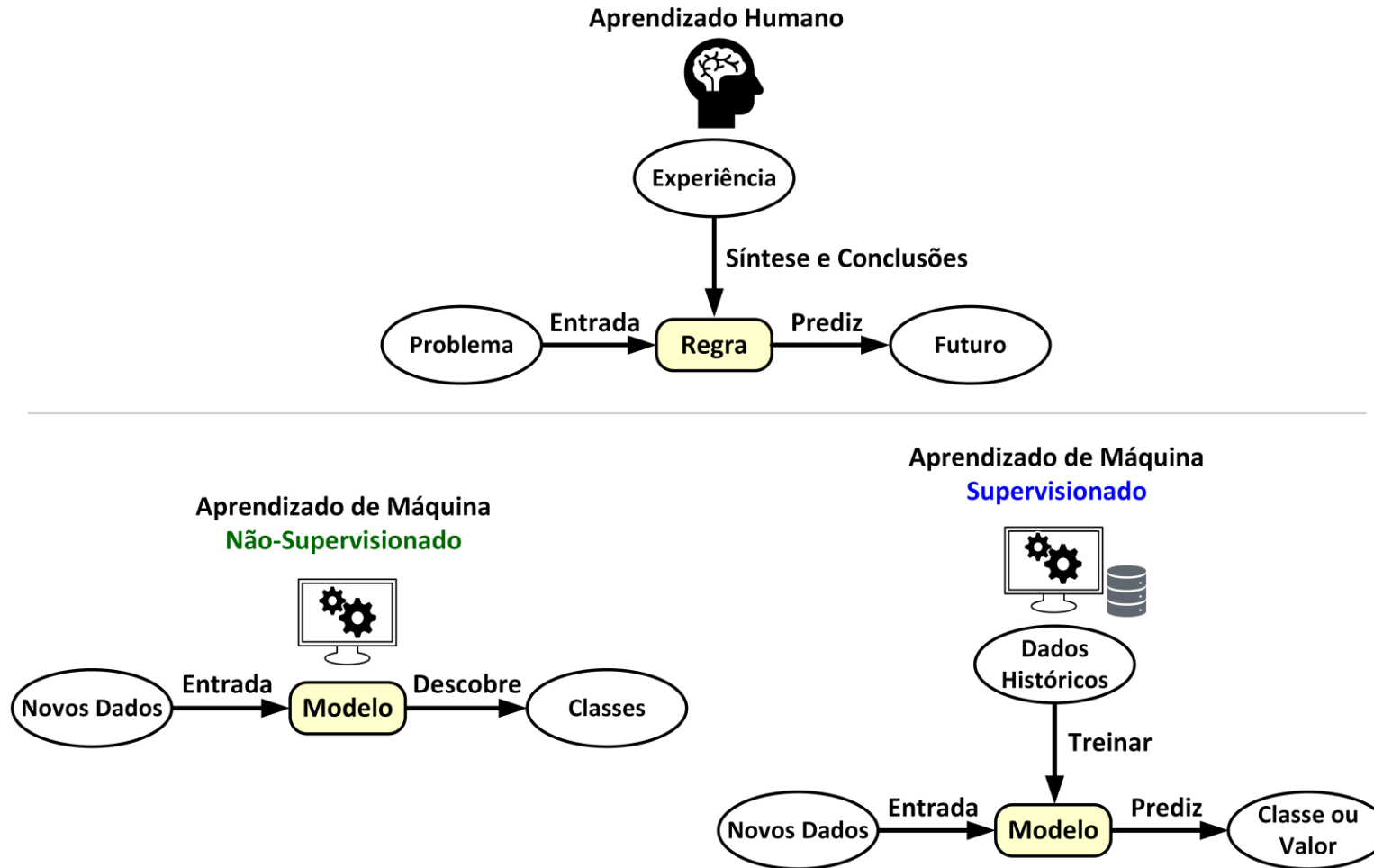
Numéricos: representados por valores inteiros ou reais, de quantidade é mensurável.

Escala de Intervalo: medidos em escala de **unidades de igual tamanho**. Valores positivos ou negativos (**sem zero absoluto**). **Ex.** temperatura.

Escala de Razão: **acomoda o zero absoluto** e **não consideram valores negativos**, permitindo que uma medida seja dimensionada em **proporção**. **Ex.** peso, altura.

2. Tipos de Aprendizado de Máquina

O aprendizado demanda uma interação entre um agente que receberá o conhecimento (aluno) e o ambiente (professor).

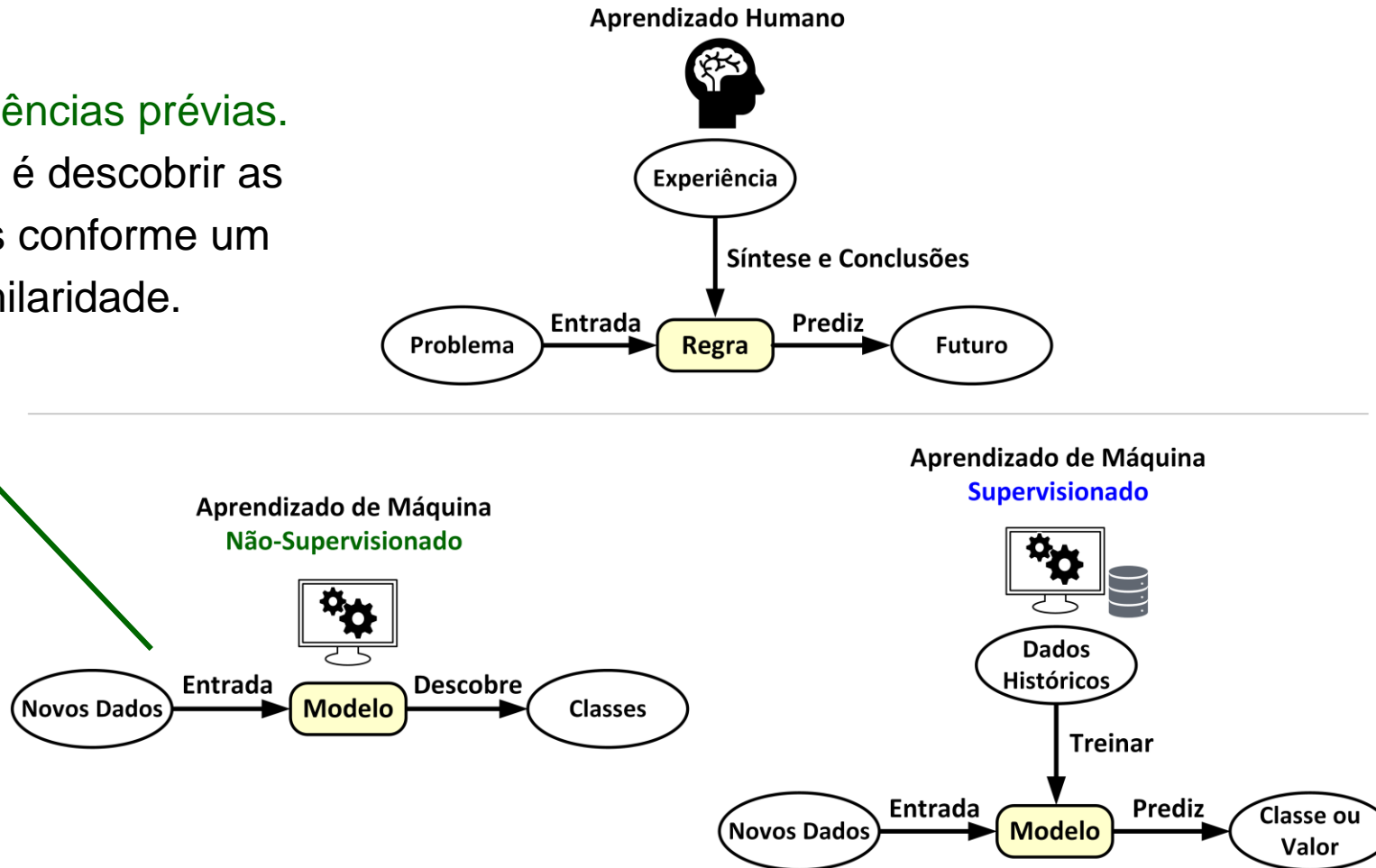


2. Tipos de Aprendizado de Máquina

O aprendizado demanda uma interação entre um agente que receberá o conhecimento (aluno) e o ambiente (professor). Os paradigmas de aprendizagem variam de acordo com o papel do agente.

Não existem experiências prévias.

A tarefa do modelo é descobrir as classes de objetos conforme um critério de similaridade.

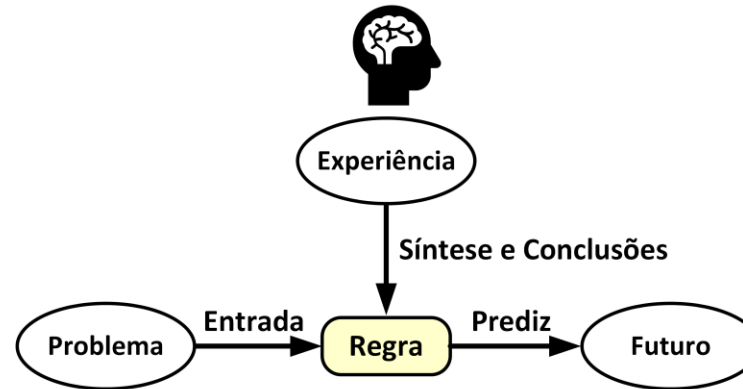


2. Tipos de Aprendizado de Máquina

O aprendizado demanda uma interação entre um agente que receberá o conhecimento (aluno) e o ambiente (professor). **Os paradigmas de aprendizagem variam de acordo com o papel do agente.**

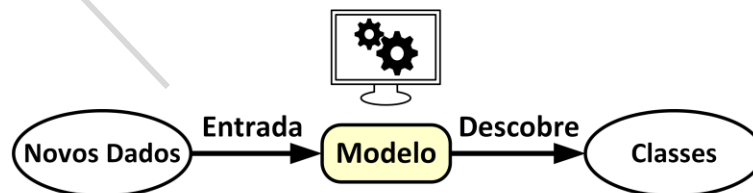
Não existem experiências prévias.
A tarefa do modelo é descobrir as
classes de objetos conforme um
critério de similaridade.

Aprendizado Humano

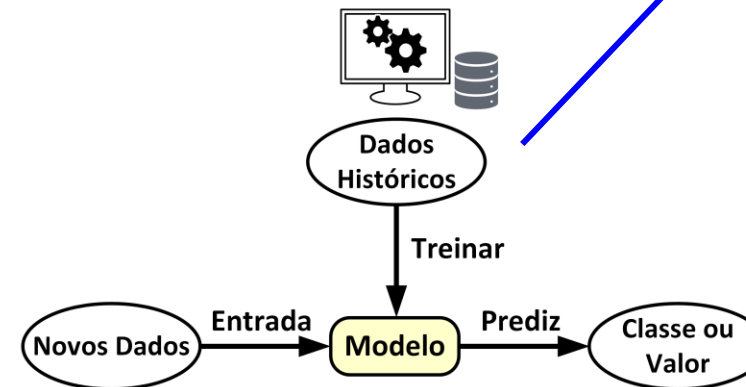


Utiliza dados históricos para o
treinamento do modelo, que
são referência para prever
rótulos de classes ou valores.

Aprendizado de Máquina
Não-Supervisionado



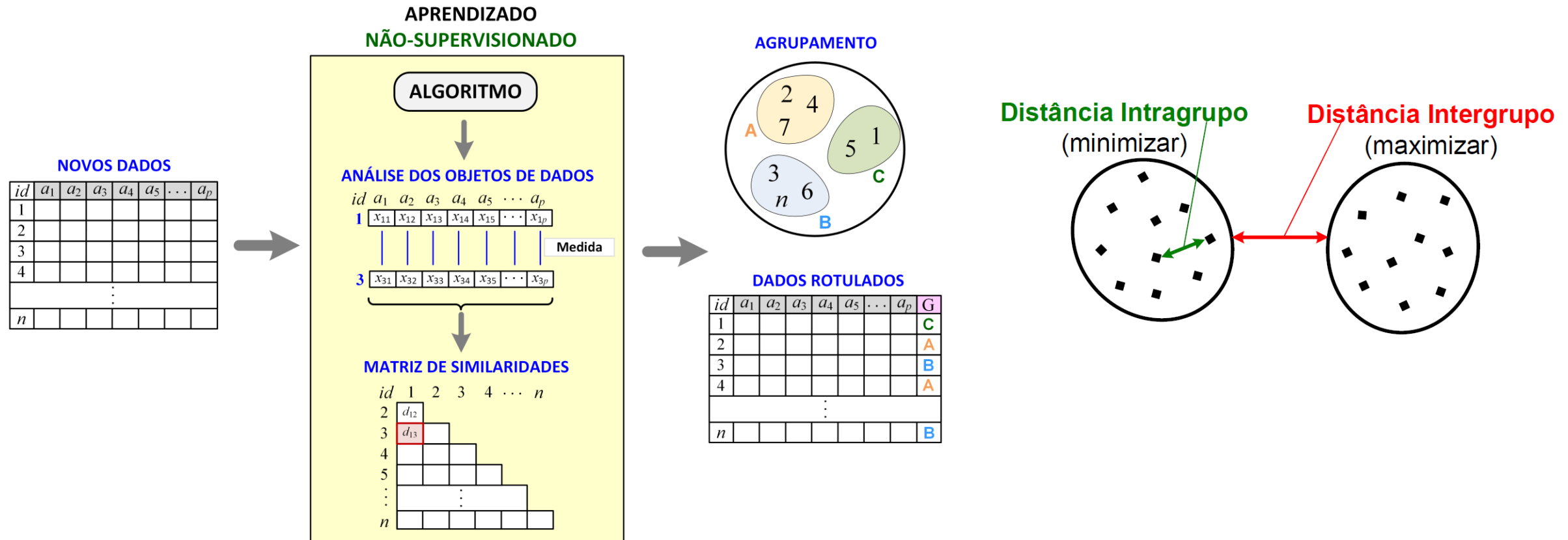
Aprendizado de Máquina
Supervisionado



2. Tipos de Aprendizado de Máquina

Aprendizado Não-Supervisionado

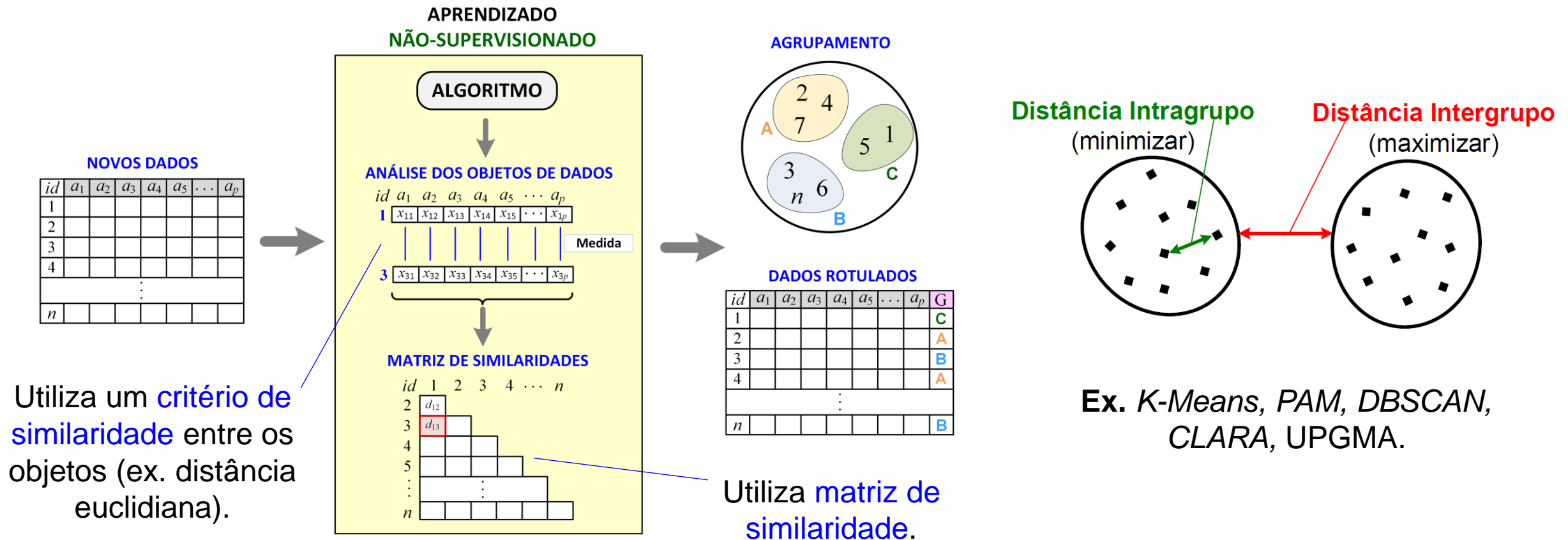
Abrange os métodos que realizam a tarefa de agrupamento (*clustering*) de dados, nas situações em que não se sabe de previamente quais são as classes de objetos ou mesmo sua quantidade.



2. Tipos de Aprendizado de Máquina

Aprendizado Não-Supervisionado

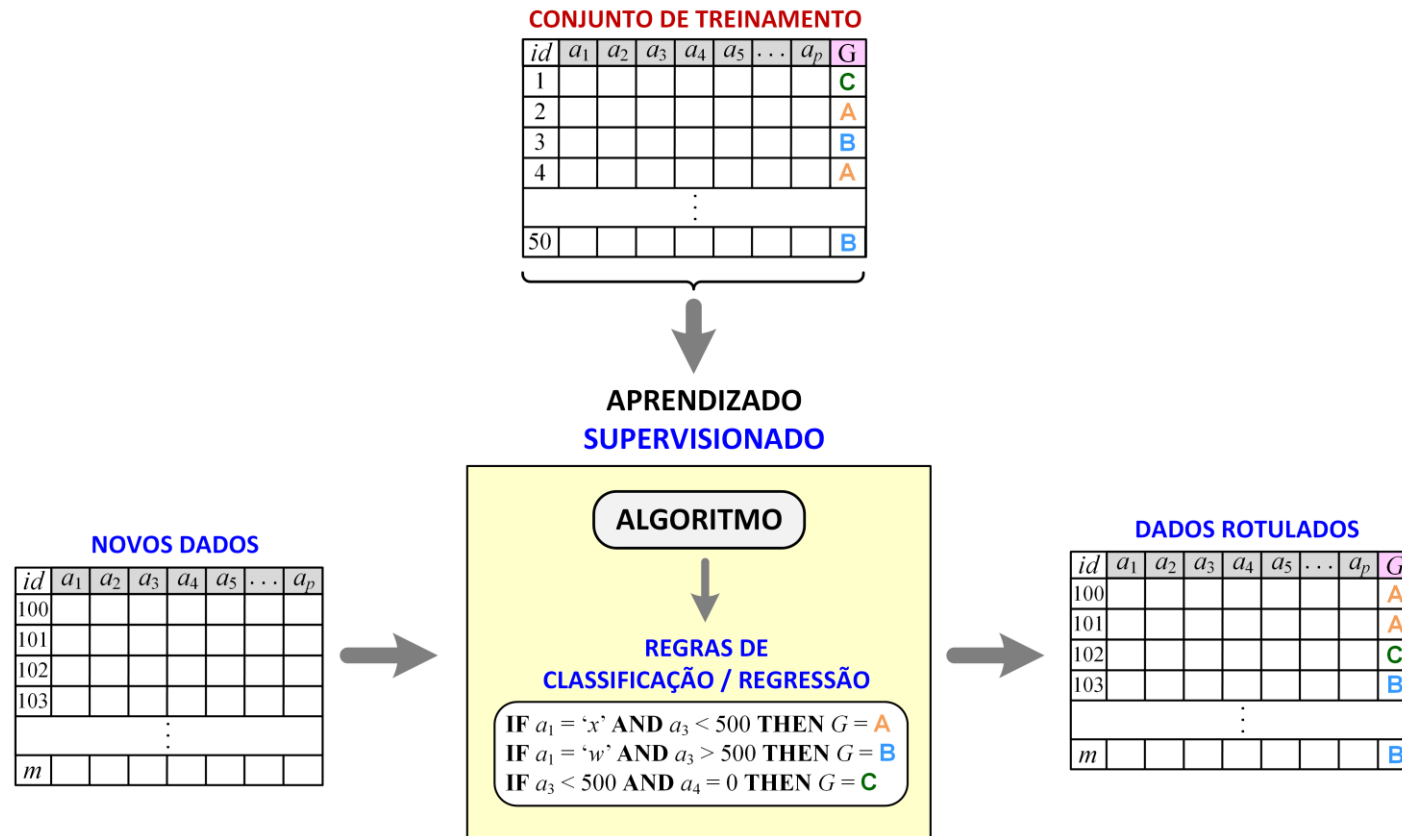
Abrange os métodos que realizam a tarefa de agrupamento (*clustering*) de dados, nas situações em que não se sabe de previamente quais são as classes de objetos ou mesmo sua quantidade.



2. Tipos de Aprendizado de Máquina

Aprendizado Supervisionado

Reúne os métodos preditivos, que são capazes de deduzir valores a partir de dados anteriores, através de algoritmos que realizam a tarefa de **classificação** (rótulos de classe) e **regressão** (valores).

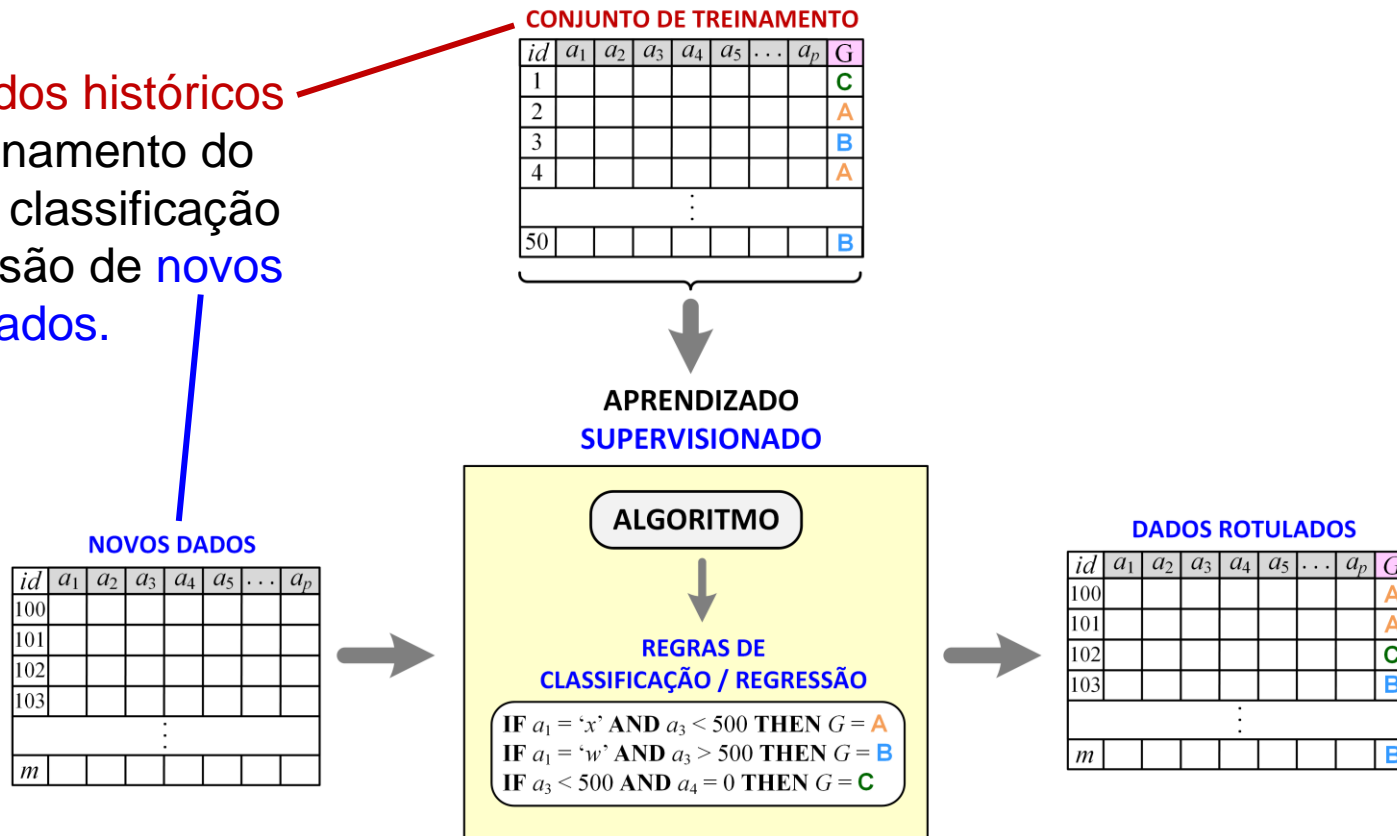


2. Tipos de Aprendizado de Máquina

Aprendizado Supervisionado

Reúne os métodos preditivos, que são capazes de deduzir valores a partir de dados anteriores, através de algoritmos que realizam a tarefa de **classificação** (rótulos de classe) e **regressão** (valores).

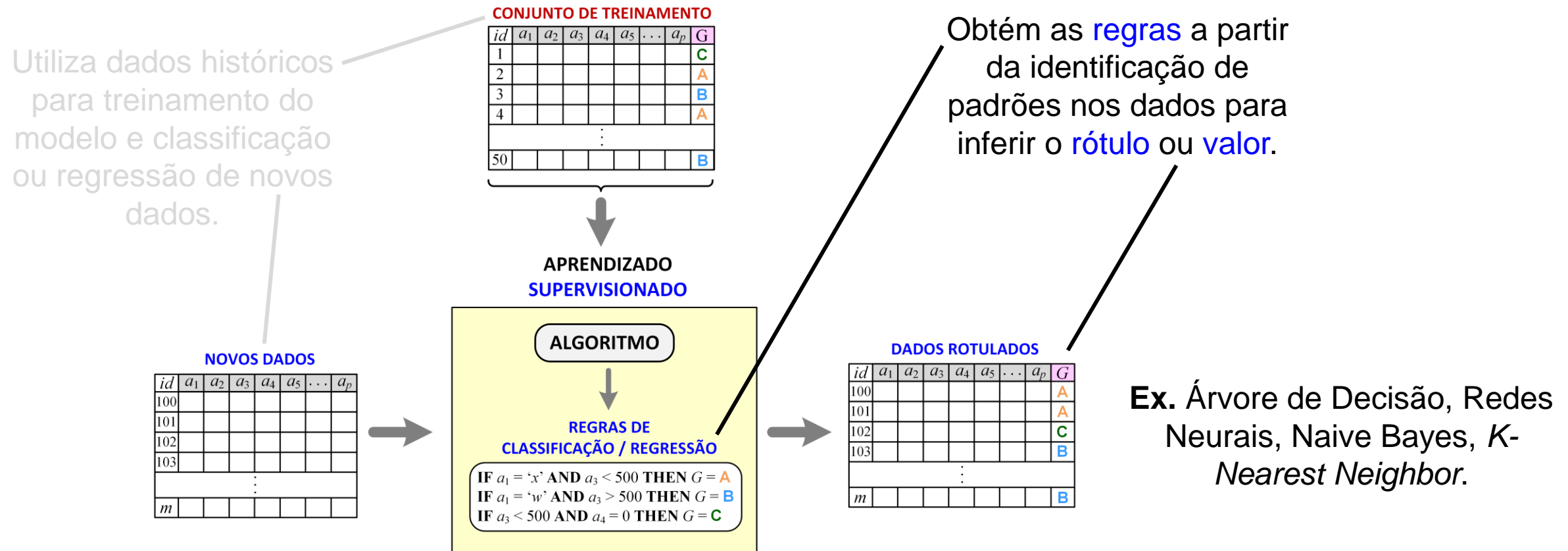
Utiliza **dados históricos** para treinamento do modelo e classificação ou regressão de **novos dados**.



2. Tipos de Aprendizado de Máquina

Aprendizado Supervisionado

Reúne os métodos preditivos, que são capazes de deduzir valores a partir de dados anteriores, através de algoritmos que realizam a tarefa de **classificação** (rótulos de classe) e **regressão** (valores).

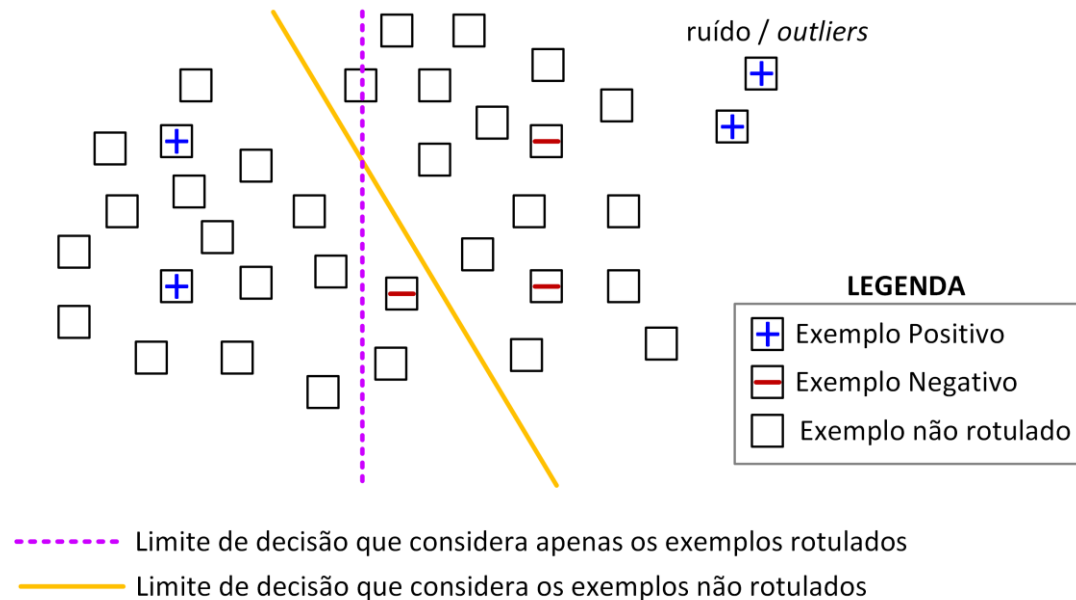


2. Tipos de Aprendizado de Máquina

Aprendizado Semi-Supervisionado

Meio termo entre aprendizado não-supervisionado e supervisionado.

Objetos rotulados são usados para aprender modelos de classe e os não-rotulados para refinar os limites entre as classes.

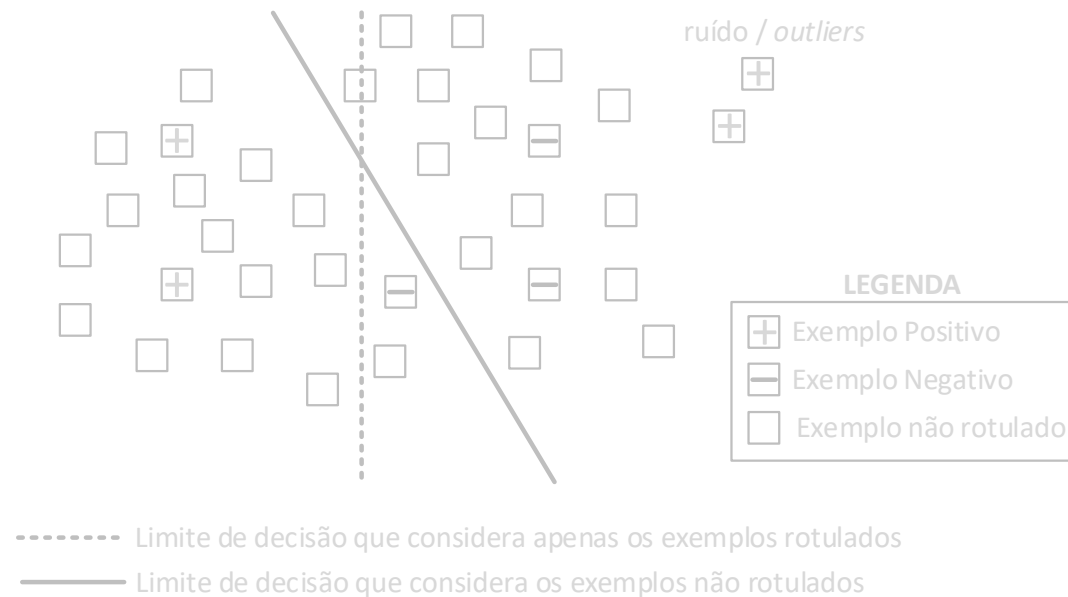


2. Tipos de Aprendizado de Máquina

Aprendizado Semi-Supervisionado

Meio termo entre aprendizado não-supervisionado e supervisionado.

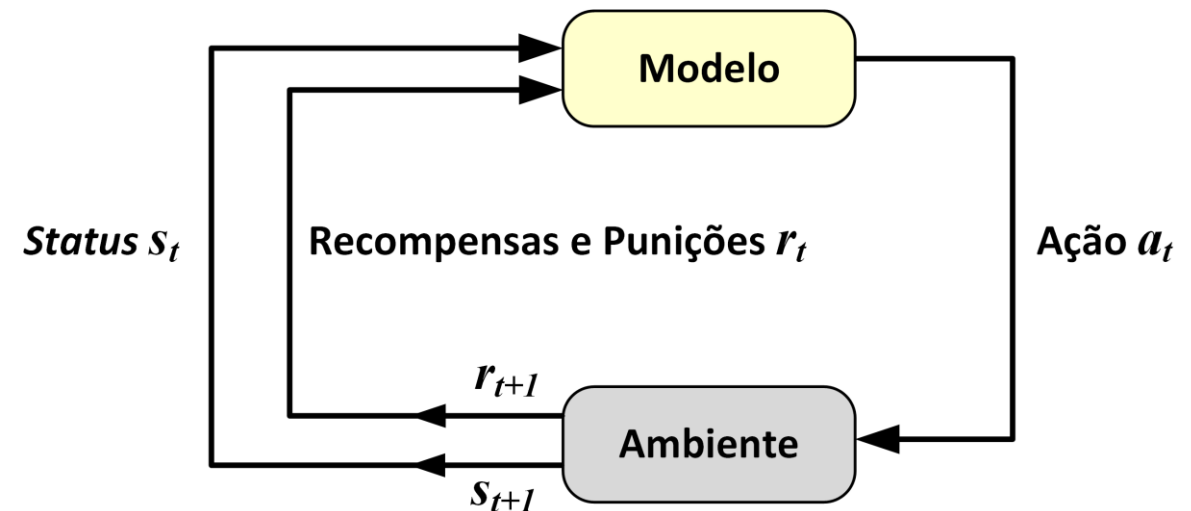
Objetos rotulados são usados para aprender modelos de classe e os não-rotulados para refinar os limites entre as classes.



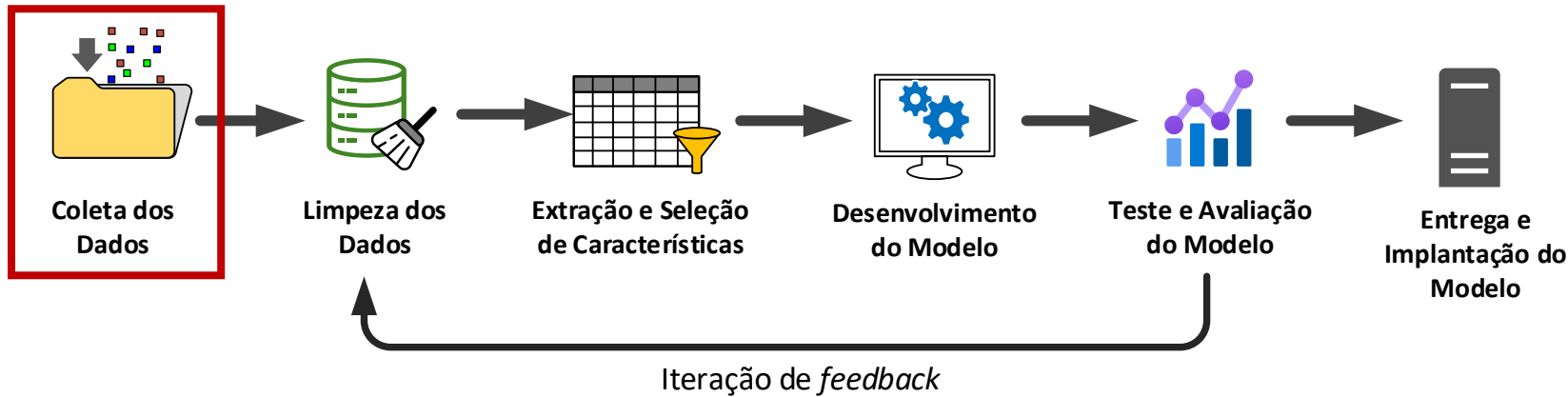
Aprendizado por Reforço

Utiliza uma estratégia de tentativa e erro, que **recompensa boas ações e penaliza as ruins** com objetivo de maximizar a recompensa.

Todos os dados e sinais de recompensa são gerados dinamicamente. O **aprendizado ativo** tem **participação humana** na rotulagem.

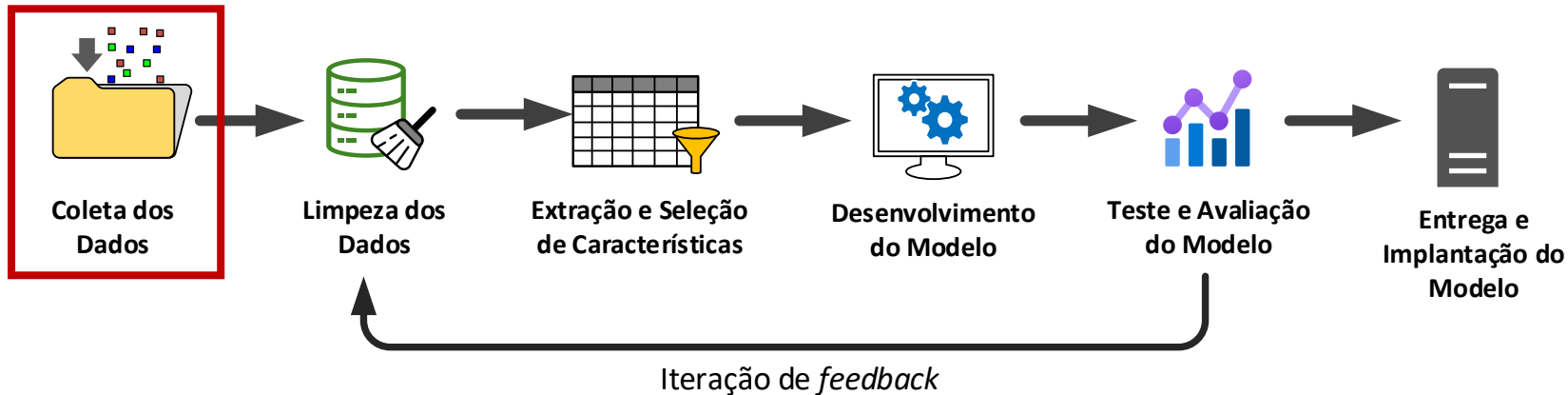


3. Processo de Aprendizado de Máquina



- Visa o **conhecer a fonte dos dados** (ex. arquivo cru, banco de dados, *data cube*, arquivo formatado etc.), bem como o **modo como estão organizados** (ex. não-estruturados, estruturados, semi-estruturados, híbrido).

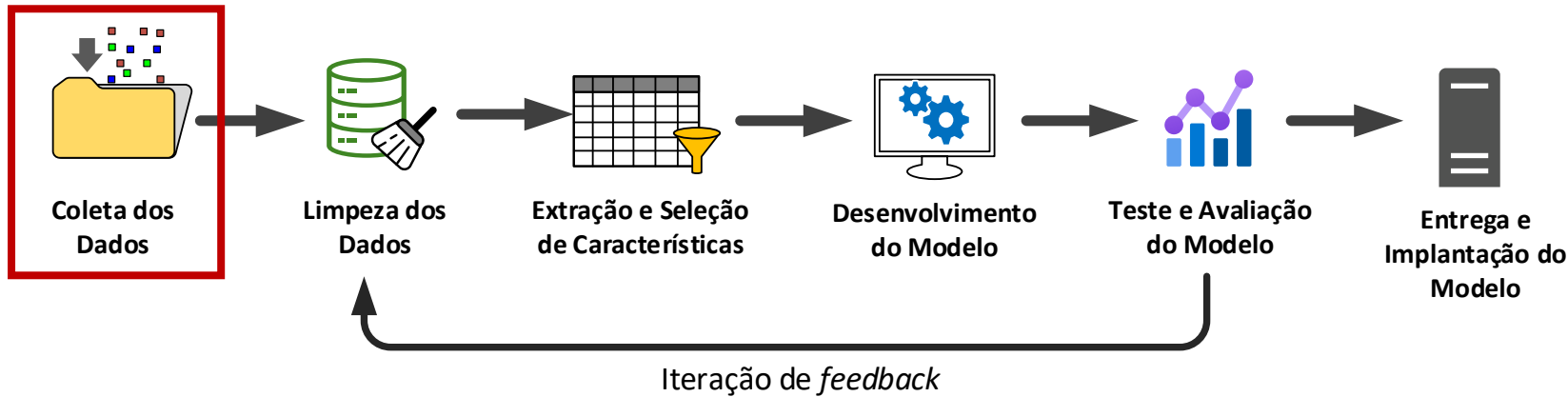
3. Processo de Aprendizado de Máquina



- Os dados coletados constituem o **conjunto de dados** que é dividido em **conjuntos** de **treinamento** e de **teste**. Uma parte do conjunto de treinamento pode constituir um conjunto de **validação** (ajustes de hiperparâmetros).

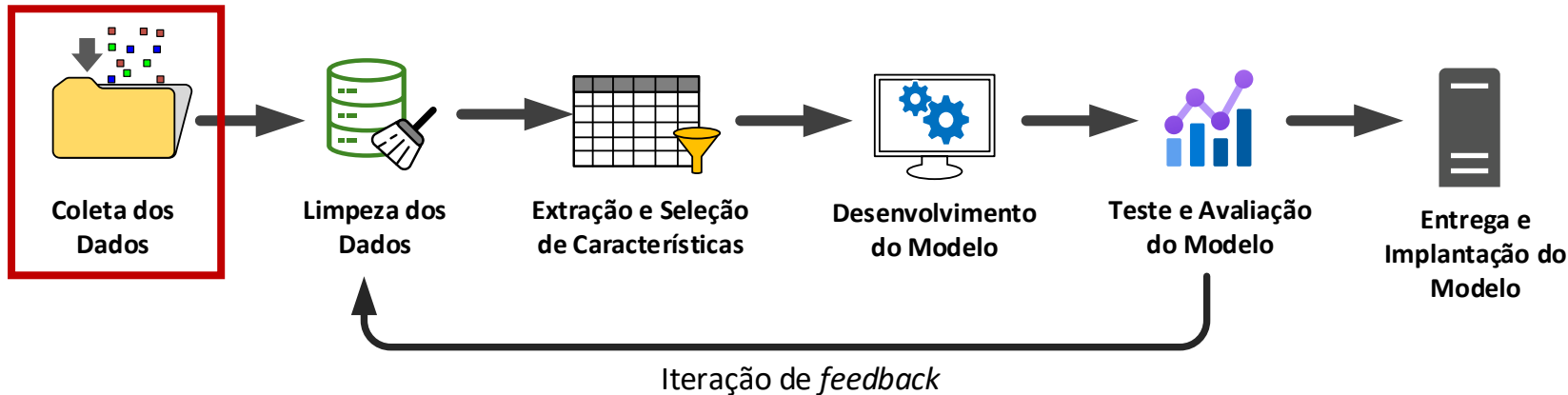
		Características								
		ID	MORADIA ALUNO	MORADIA FAMILIA	ESTADO CIVIL	FILHOS	RENDIA	DESPESAS	DOENTES	PATRIMÔNIO
Conjunto de Treinamento		11323	Aluguel com Colegas	Heranca	Solteiro	2	166	73	0	0
		11289	Aluguel com Colegas	Propria Quitada	Solteiro	0	582	155	0	0
		23363	Familia/Terceiros	Propria em Pagamento	Divorciado	0	1118	56	1	15000
		22148	Familia/Terceiros	Heranca	Solteiro	1	407	0	0	0
		23060	Aluguel com Colegas	Propria Quitada	Casado	0	636	133	1	60000
		30394	Familia/Terceiros	Alugada	Solteiro	0	533	53	0	45000
		27836	Familia/Terceiros	Heranca	Solteiro	0	281	22	0	0
		17099	Aluguel com Colegas	Propria Quitada	Casado	0	1238	54	0	45443
Conjunto de Teste		20042	Familia/Terceiros	Alugada	Solteiro	0	400	200	0	0
		29757	Aluguel Sozinho	Propria Quitada	Solteiro	0	1227	53	0	7000

3. Processo de Aprendizado de Máquina



- Uma vez definido, o conjunto de dados é analisado para uma maior familiaridade com os dados e identificar quais são os objetos de dados e tipos de atributos.

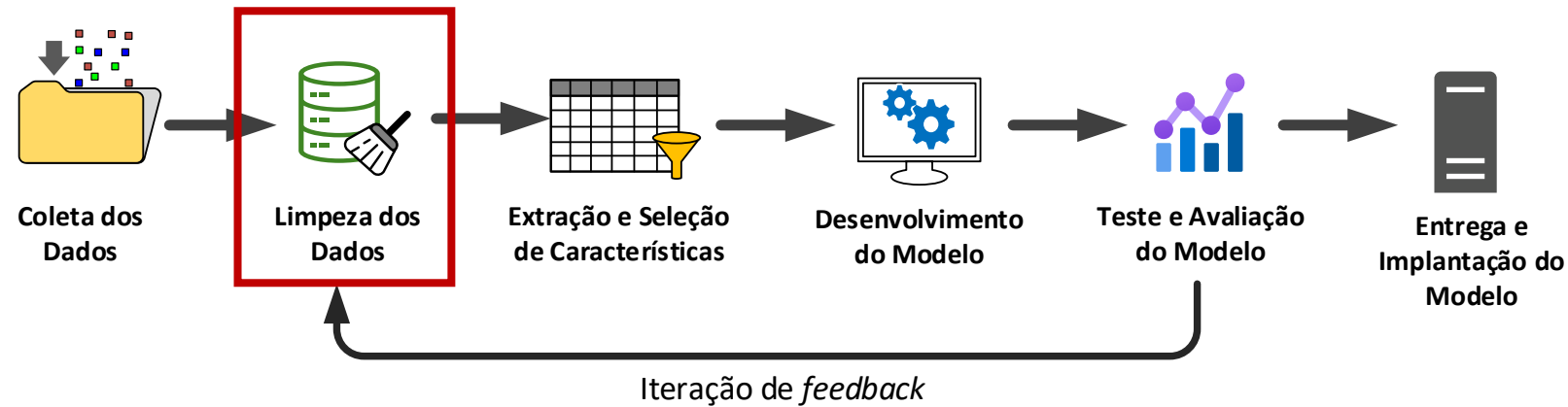
3. Processo de Aprendizado de Máquina



- Uma vez definido, o conjunto de dados é analisado para uma maior familiaridade com os dados e identificar quais são os objetos de dados e tipos de atributos.
- Métodos úteis nesta fase: estatística descritiva, análise exploratória de dados, **dicionário de dados** (contém nome do atributo, descrição, tipo, rótulos de classes, limites de valores por atributo).

Atributo	Tipo	Limite de Valores
proc_escolar	String	Procedência escolar do solicitante. Classes: "Pública"; "Filantrópica"; "Particular bolsa ≥50%"; "Particular bolsa <50%"; "Particular sem bolsa"
sit_moradia	String	Situação da moradia do solicitante. Classes: "Família, Parentes, Terceiros"; "Divide Aluguel"; "Sozinho Imóvel Alugado ou Financiado"; "Sozinho Imóvel Próprio"
sit_moradia_fam	String	Situação da moradia da família do solicitante. Classes = "Alugada"; "Cedida, Herança"; "Própria e Quitada"; "Própria em Pagamento"
filhos	Inteiro	Número de filhos que o solicitante possui. Valores inteiros iguais ou maiores que 0.
renda	Inteiro	Renda per capita familiar. Classes: 1 = [0, 200]; 2 = (200, 286]; 3 = (286, 360]; 4 = (360, 433]; 5 = (433, 500]; 6 = (500, 579]; 7 = (579, 676]; 8 = (676, 767]; 9 = (767, 988]; 10 = (988, 2902]
despesa	Inteiro	Despesa per capita familiar. Classes: 1 = [0, 20]; 2 = (20, 45]; 3 = (45, 66]; 4 = (66, 86]; 5 = (86, 111]; 6 = (111, 139]; 7 = (139, 179]; 8 = (179, 247]; 9 = (247, 347]; 10 = (347, 1659]
doenca	Inteiro	Número de indivíduos com doença grave no grupo familiar. Valores inteiros iguais ou maiores que 0.
fam_ensino_completo	Inteiro	Número de familiares com Ensino Superior Completo ou Pós-Graduação. Valores inteiros iguais ou maiores que 0.
bens_familiares	Inteiro	Valor total dos bens familiares. Classes: 1 = [0, 3.000]; 2 = (3.000, 7.918]; 3 = (7.918, 11.415]; 4 = (11.415, 15.917]; 5 = (15.917, 22.775]; 6 = (22.775, 36.122]; 7 = (36.122, 60.000]; 8 = (60.000, 100.000]; 9 = (10.000, 156.687]; 10 = (156.687, 700.000]

3. Processo de Aprendizado de Máquina



Os dados são preparados para estarem aptos para o treinamento do modelo, solucionando problemas de qualidade como valores ausentes, suavizar o ruído, redução e corrigir inconsistências.

3. Processo de Aprendizado de Máquina

Exemplos de Problemas de Qualidade em Dados

#	ID	NOME	ANIVERSÁRIO	GENERO	PROFESSOR	ALUNOS	ESTADO	CIDADE
1	111	João	31/12/2000	M	0	0	São Paulo	Piracicaba
2	222	Maria	19/11/1977	F	1	15	Bahia	
3	333	Alice	19/04/2001	F	0	0	Rio de Janeiro	Resende
4	444	Marcos	03/08/1968	M	0	0	Minas erais	Ipatinga
5	555	Alex	14/05/1980	A	1	23	São Paulo	Barretos
6	555	Pedro	25/12/2001	M	1	10	Amazonas	Amazonas
7	777	Caio	1986/12/01	M	0	0	Piauí	Cocal
8	888	Rosana	31/05/1997	F	0	0	Maranhão	São Luís
9	999	Ana	05/08/2004	F	0	5	Minas Gerais	Itajubá
10	101010	Paulo	14/07/1993	M	1	26	Espírito Santo	Vitória

Valor Inválido

Erro ortográfico

Valor Faltante

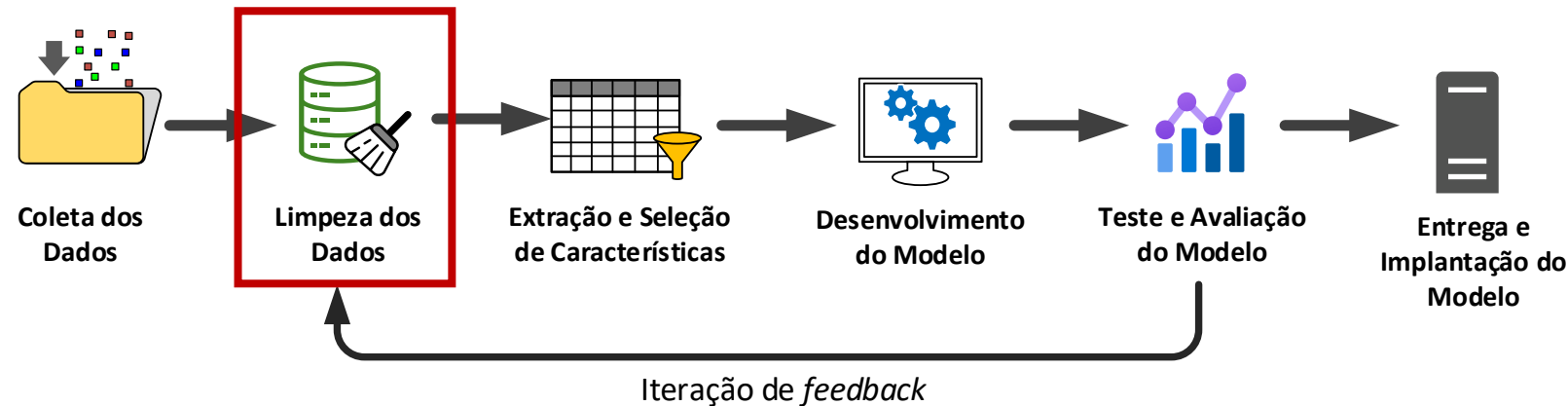
Item duplicado inválido

Formato Incorreto

Dependência de Atributo

Valor que deveria estar em outra coluna

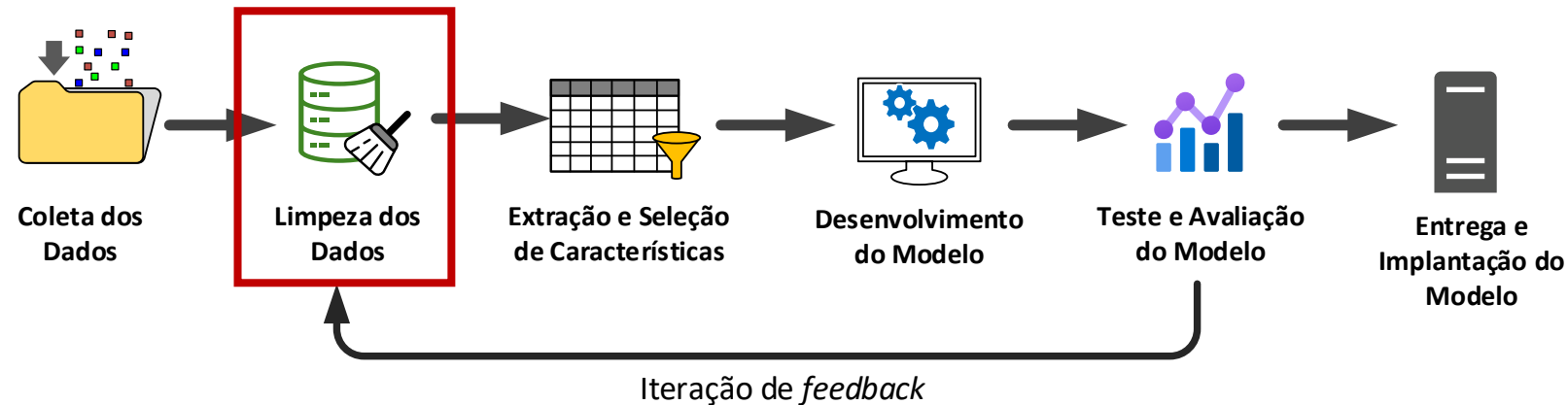
3. Processo de Aprendizado de Máquina



Os dados são preparados para estarem aptos para o treinamento do modelo, **solucionando problemas** de qualidade como **valores ausentes**, suavizar o ruído, redução e corrigir inconsistências. **Métodos:**

- **Ignorar o objeto:** Geralmente é usado quando o rótulo da classe está ausente. Recomendado quando existem vários atributos com valores ausentes.
- **Preenchimento Manual:** Quando se conhece o valor. Pode ser inviável em grandes conjuntos de dados.
- **Constante global:** Substituir os valores de atributo ausentes pela mesma constante (ex. NULL).
- **Medida de tendência central:** Preencher com a média dos valores do atributo em distribuições normais de dados (simétricas) ou a mediana (assimétricas).
- **Valor mais provável:** determinado com regressão ou outras ferramentas baseadas em inferência.

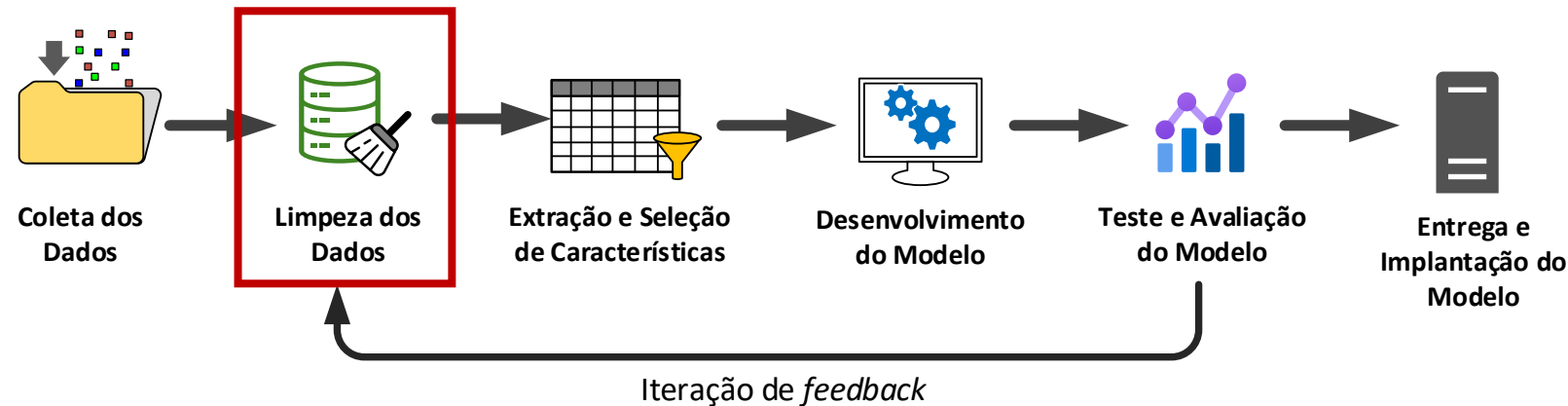
3. Processo de Aprendizado de Máquina



Os dados são preparados para estarem aptos para o treinamento do modelo, **solucionando problemas** de qualidade como valores ausentes, **suavizar o ruído**, redução e corrigir inconsistências. **Métodos:**

- **Análise de *Outliers*:** Técnicas básicas de **estatística descritiva** (ex. *boxplots* e *scatter plots*) ou **agrupamento** podem ser usadas para identificar outliers, que podem representar ruído.
- **Categorização (*Binning*):** suavizam valores de dados ordenados distribuindo-os em compartimentos (Bins) com a mesma quantidade de valores. Ex. um atributo cujos valores foram particionados em: Bin 1 (4, 8 e 15), Bin 2 (20, 22, 24). Usando a média, os valores originais de cada Bin são substituídos por: Bin 1 (9, 9, 9), Bin 2 (22, 22, 22).
- **Regressão:** A regressão linear envolve **encontrar a melhor função para ajustar dois atributos** de modo que um atributo possa ser usado para prever outro.

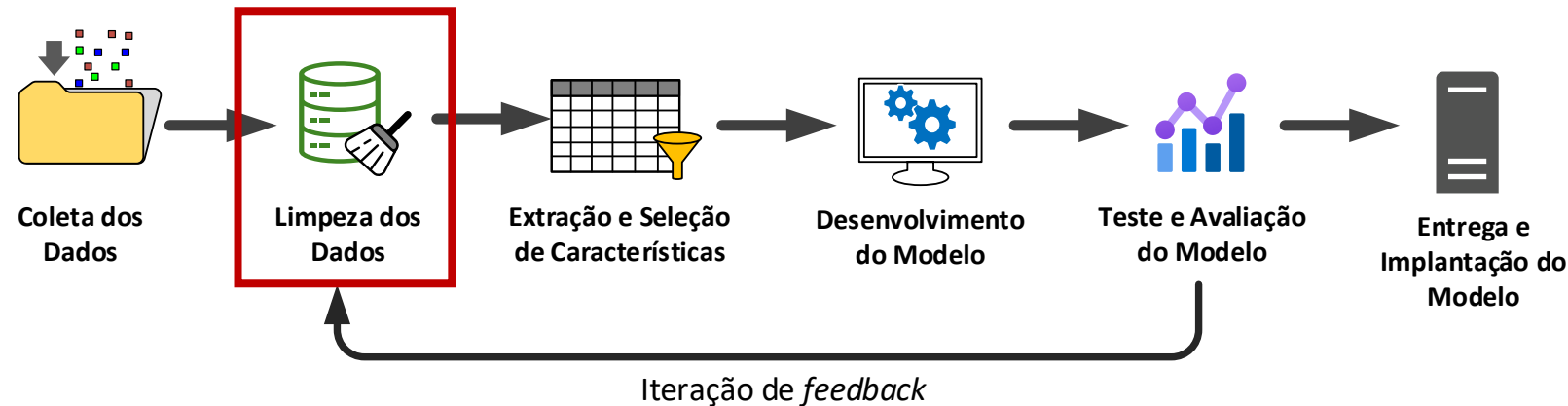
3. Processo de Aprendizado de Máquina



Os dados são preparados para estarem aptos para o treinamento do modelo, **solucionando problemas** de qualidade como valores ausentes, suavizar o ruído, **redução** e corrigir inconsistências. **Métodos:**

- **Análise de *Outliers*:** Redução de dimensionalidade: métodos que diminuem a quantidade variáveis aleatórias ou atributos a serem considerados. Ex. Transformada Wavelet e Análise de Componentes Principais (*Principal Component Analysis*, PCA).
- **Redução de numerosidade:** substitui o volume de dados original por formas alternativas e menores de representação de dados. Inclui métodos paramétricos (ex. Regressão, Log-Linear) e os não paramétricos (ex. histogramas, agrupamento, amostragem e agregação de cubo de dados).

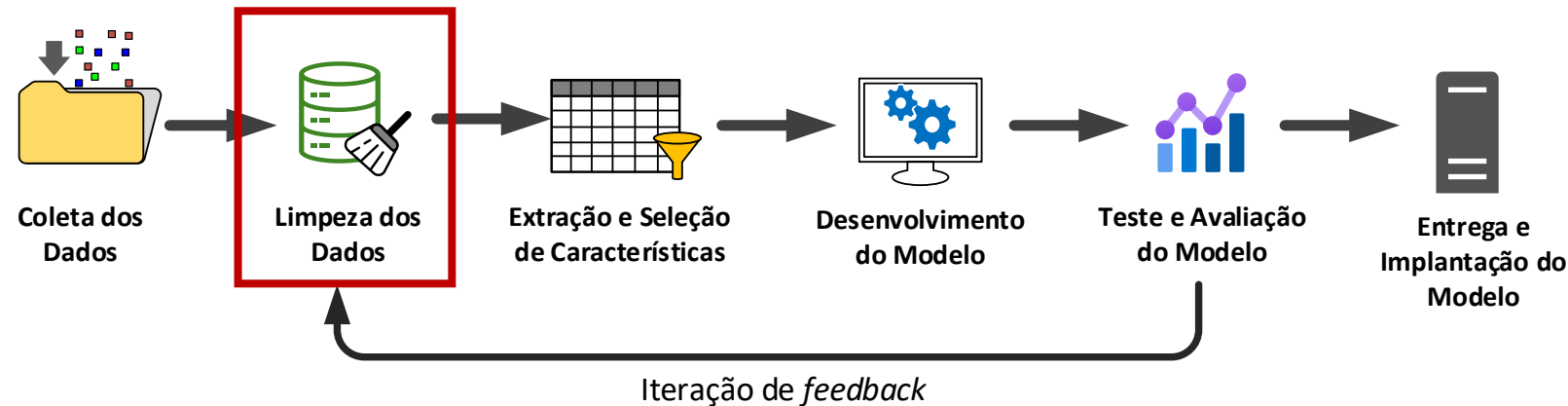
3. Processo de Aprendizado de Máquina



Os dados são preparados para estarem aptos para o treinamento do modelo, **solucionando problemas** de qualidade como valores ausentes, suavizar o ruído, redução e **corrigir inconsistências**. **Métodos:**

- **Suavização:** mesmo princípio usado na redução de ruídos em dados (ex. *binning*, regressão e agrupamento).
- **Construção de atributo:** **novos atributos a partir da combinação** de um conjunto **de atributos**.
- **Agregação:** **síntese** ou agregação de dados (ex. dados de vendas diárias agregado a venda mensal).

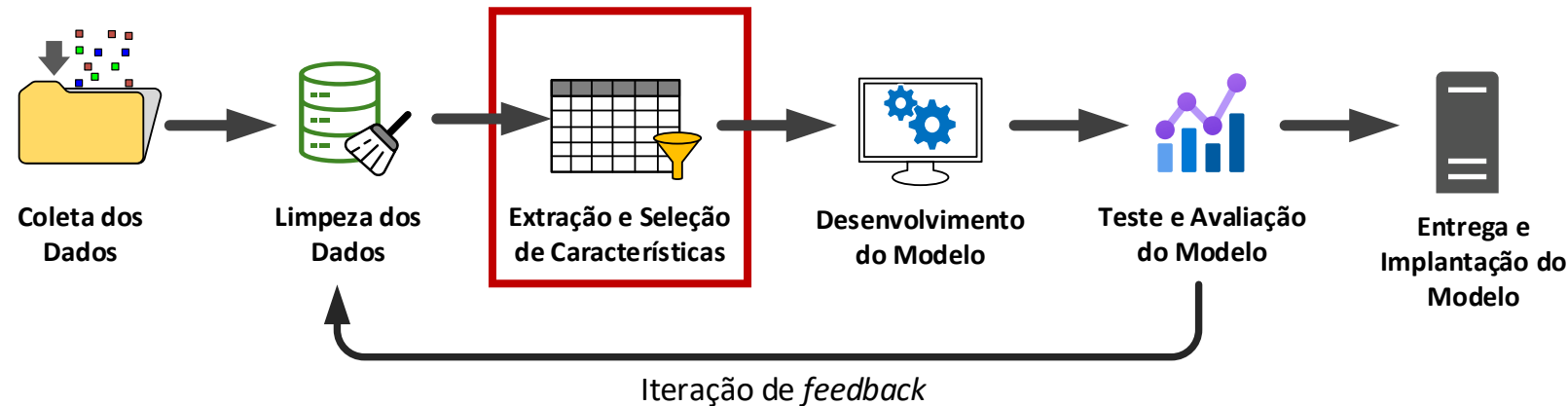
3. Processo de Aprendizado de Máquina



Os dados são preparados para estarem aptos para o treinamento do modelo, **solucionando problemas** de qualidade como valores ausentes, suavizar o ruído, redução e **corrigir inconsistências**. **Métodos:**

- **Suavização:** mesmo princípio usado na redução de ruídos em dados (ex. *binning*, regressão e agrupamento).
- **Construção de atributo:** novos atributos a partir da combinação de um conjunto de atributos.
- **Agregação:** síntese ou agregação de dados (ex. dados de vendas diárias agregado a venda mensal).
- **Normalização:** dados são redimensionados para um menor intervalo [0, 1]. Ex. Normalização z-score, Min-Max.
- **Discretização:** **valores** de um atributo **numérico** (ex. idade) **são substituídos por rótulos de intervalo** (ex. 0–10, 11–20, etc.) ou **rótulos conceituais** (ex. jovem, adulto). Ex. *Binning*, Histograma, agrupamento, análise de Correlação.
- **Hierarquia de Conceitos para Dados Nominais:** uso de **conceitos de nível superior** (ex. bairro ao invés das ruas).

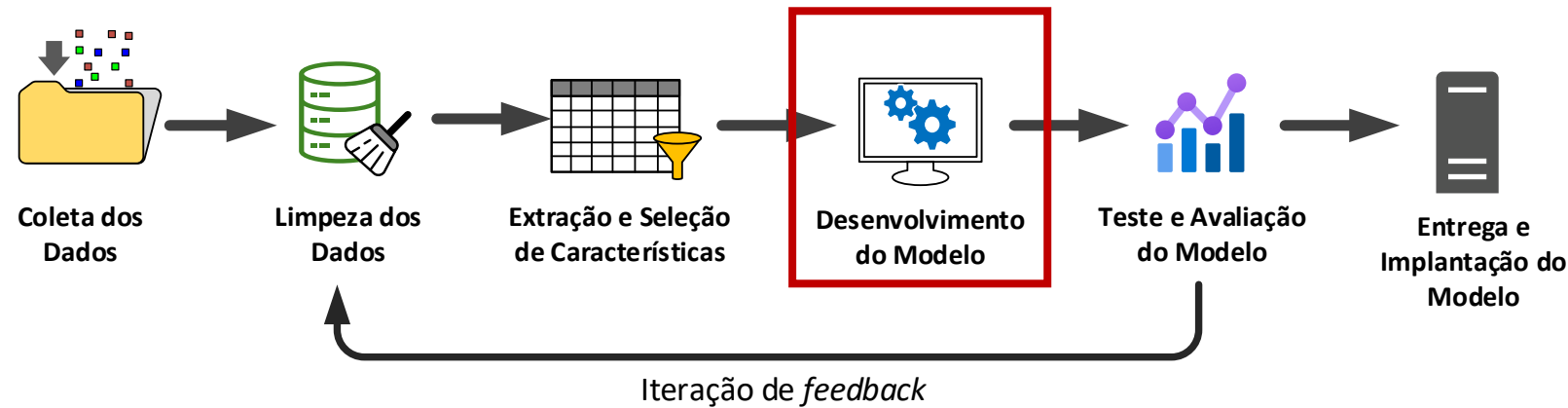
3. Processo de Aprendizado de Máquina



Tem por objetivo **extrair e selecionar os atributos mais relevantes** para compor o conjunto de treinamento e de testes, desconsiderando os atributos que não contribuem para solução do problema.

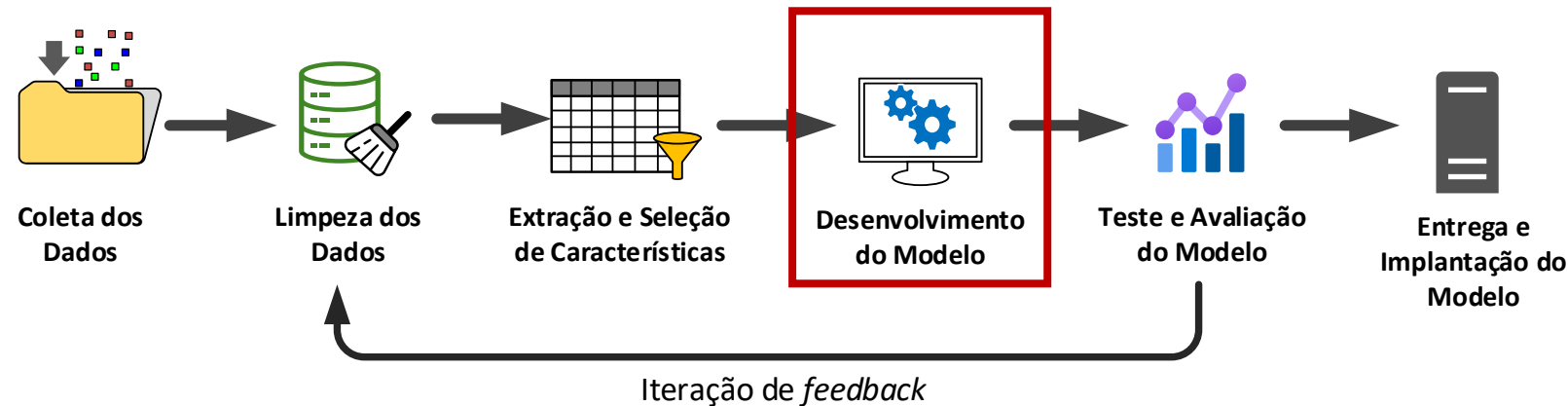
- **Filtro:** método que seleciona o atributo de forma independente do modelo em si. **Uso da correlação entre atributos para decidir quais manter e quais desconsiderar.** Ex. correlação de Pearson, coeficiente qui-quadrado e Informações Mútuas.
- **Wrapper (Invólucro):** usa um modelo preditivo para pontuar um subconjunto de atributos e considera o problema de seleção como uma busca, na qual o **wrapper avalia e compara diferentes combinações de atributos**, com o modelo preditivo sendo usado para avaliar estas combinações.
- **Embedded (Embutido):** usa a **seleção de atributos como parte da construção do modelo**. Ao contrário dos métodos supracitados, **aprende a seleção de atributos durante o treinamento**. Ex. Regressão Lasso e Regressão Ridge.

3. Processo de Aprendizado de Máquina



Esta fase visa a **construção e o treinamento do modelo de aprendizado**. Dependendo da estratégia de verificação adotada uma parte dos **dados de treinamento poderão constituir o conjunto de validação**, o qual será usado para analisar o funcionamento do modelo e ajustes na sua parametrização.

3. Processo de Aprendizado de Máquina



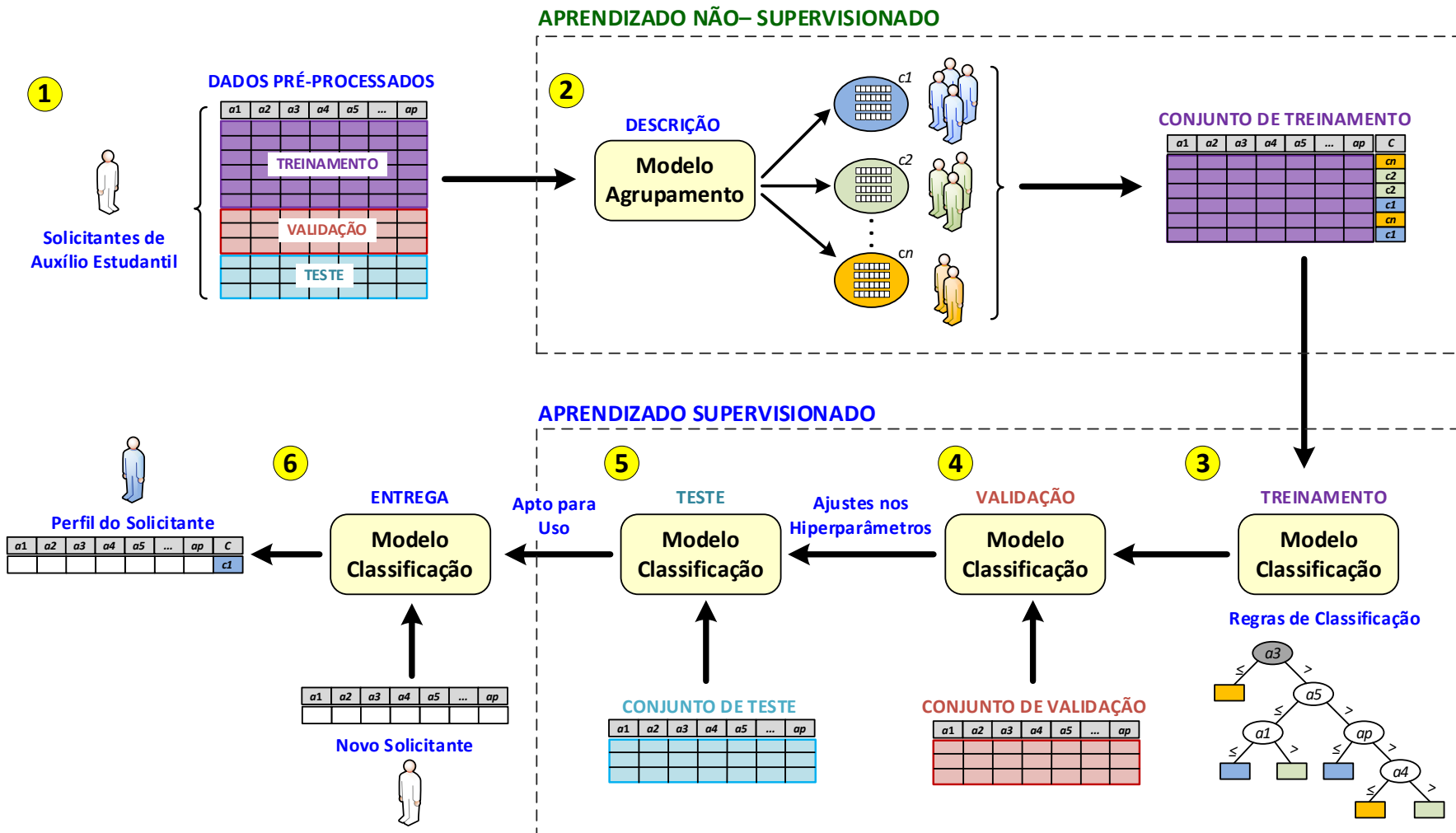
Esta fase visa a construção e o treinamento do modelo de aprendizado. Dependendo da estratégia de verificação adotada uma parte dos dados de treinamento poderão constituir o conjunto de validação, o qual será usado para analisar o funcionamento do modelo e ajustes na sua parametrização.

A partir do entendimento do problema e dos dados disponíveis, **são definidos os métodos de aprendizado** que sejam mais adequados para constituir o modelo.

A modelagem varia conforme o propósito do projeto, sendo que o **modelo final pode ser um único método** (não-supervisionado ou supervisionado), **ou um modelo híbrido** dado pela combinação de ambos os tipos de aprendizado.

3. Processo de Aprendizado de Máquina

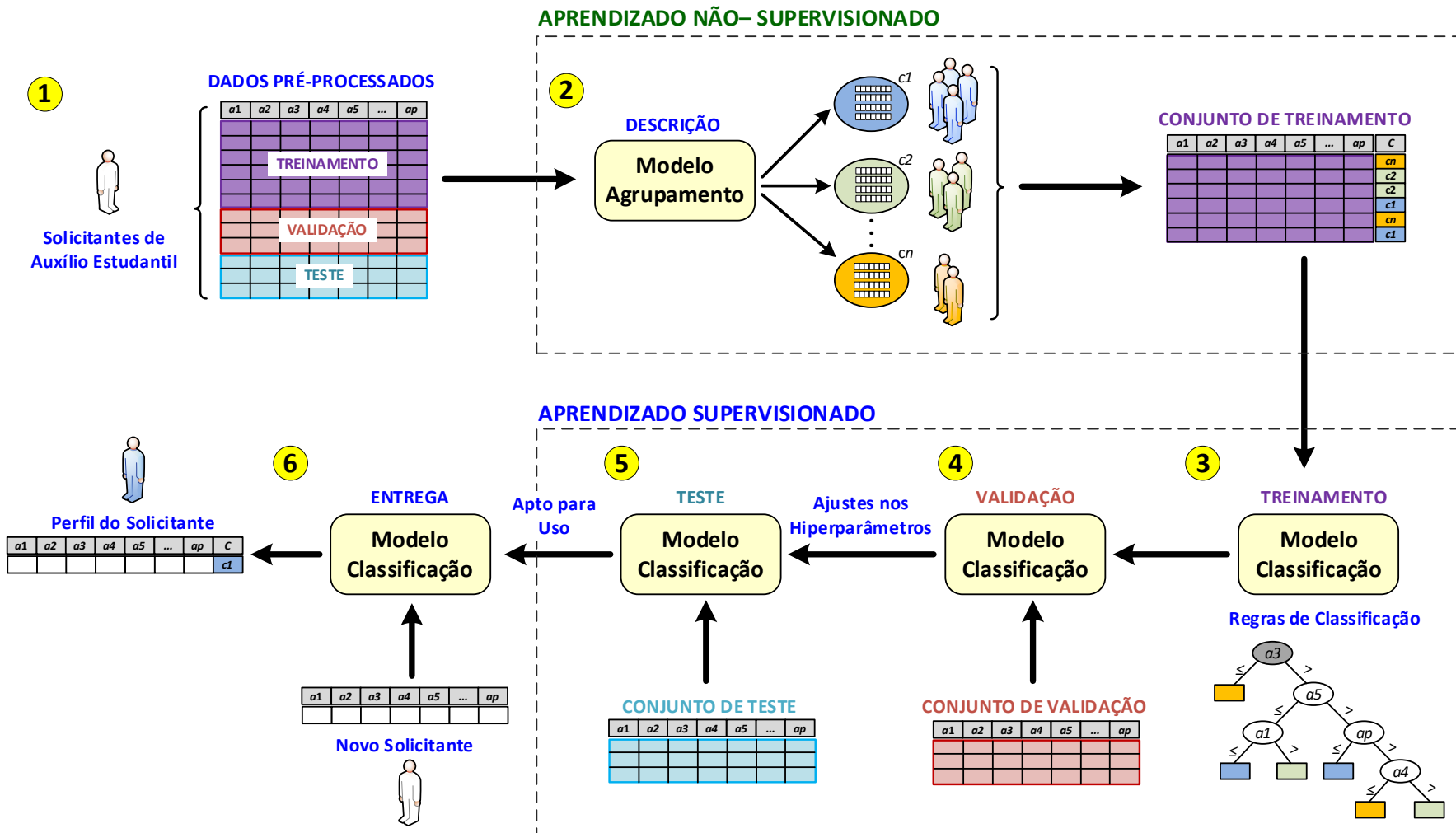
Exemplo de Modelo Híbrido de Aprendizado de Máquina



- (1) o conjunto de dados é dividido em treinamento, validação e teste.
- (2) método não-supervisionado para identificar as classes do conjunto
- (3) treinamento do método supervisionado com os dados rotulados.

3. Processo de Aprendizado de Máquina

Exemplo de Modelo Híbrido de Aprendizado de Máquina



(1) o conjunto de dados é dividido em treinamento, validação e teste.

(2) método não-supervisionado para identificar as classes do conjunto

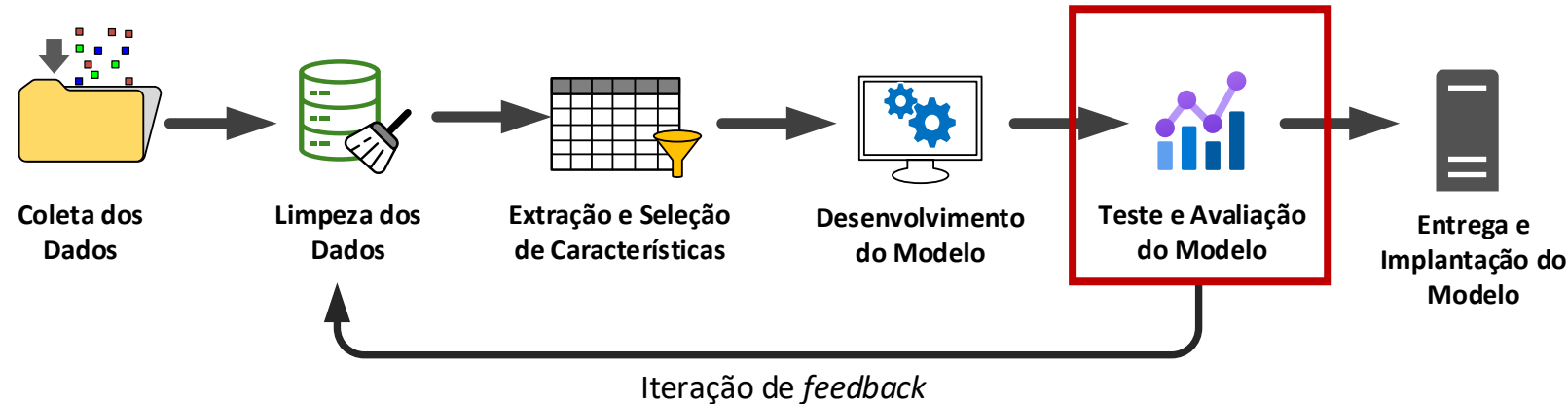
(3) treinamento do método supervisionado com os dados rotulados.

(4) verificação e ajustes no modelo com o conjunto de validação.

(5) treino do modelo e teste da sua capacidade de generalização

(6) Modelo final é disponibilizado para uso.

3. Processo de Aprendizado de Máquina



Os resultados são interpretados e a qualidade do modelo é verificada através de análise dos dados e métricas de desempenho.

CARACTERÍSTICAS DE UM MODELO DE QUALIDADE:

- **Capacidade de Generalização:** ele pode prever com precisão os dados reais (não apenas treinamento)?
- **Interpretabilidade:** o resultado da previsão é fácil de interpretar?
- **Velocidade de Previsão:** quanto tempo leva para prever cada dado?
- **Praticabilidade:** a taxa de previsão ainda é aceitável com um grande volume de dados?

3. Processo de Aprendizado de Máquina

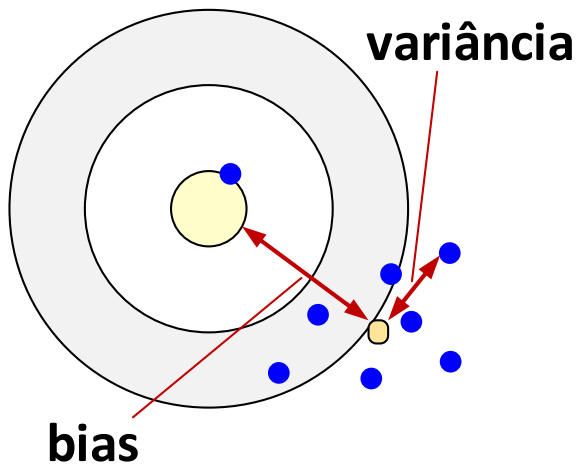
Validade do Modelo

Erro: diferença entre o resultado previsto pelo modelo após o aprendizado e o resultado amostral real.

- **de treinamento:** obtido ao executar o modelo com os dados de treinamento.
- **de generalização:** obtido ao executar o modelo com novas amostras.

viés (bias): diferença entre o valor esperado e o valor correto (ou média) que se tenta prever.

variância: desvio em relação a média dos resultados da predição.



3. Processo de Aprendizado de Máquina

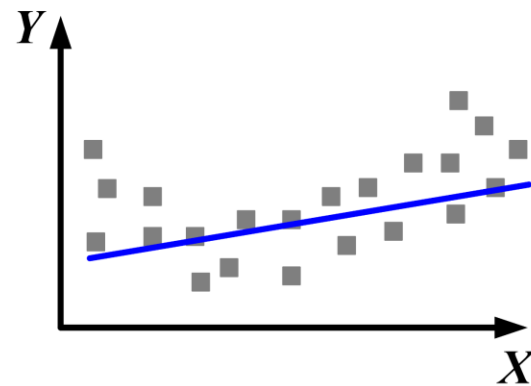
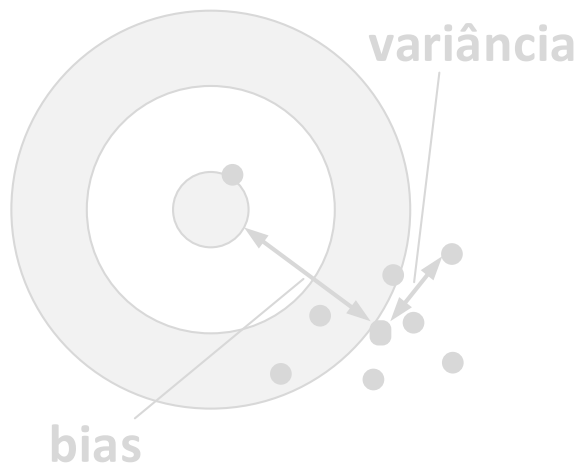
Validade do Modelo

Erro: diferença entre o resultado previsto pelo modelo após o aprendizado e o resultado amostral real.

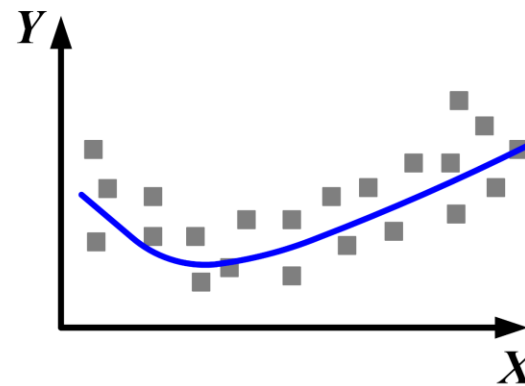
- **de treinamento:** obtido ao executar o modelo com os dados de treinamento.
- **de generalização:** obtido ao executar o modelo com novas amostras.

viés (bias): diferença entre o valor esperado e o valor correto (ou média) que se tenta prever.

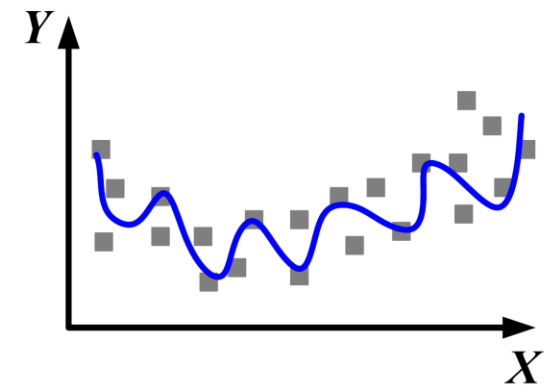
variância: desvio em relação a média dos resultados da predição.



(a) *Underfitting*
falha no aprendizado das características



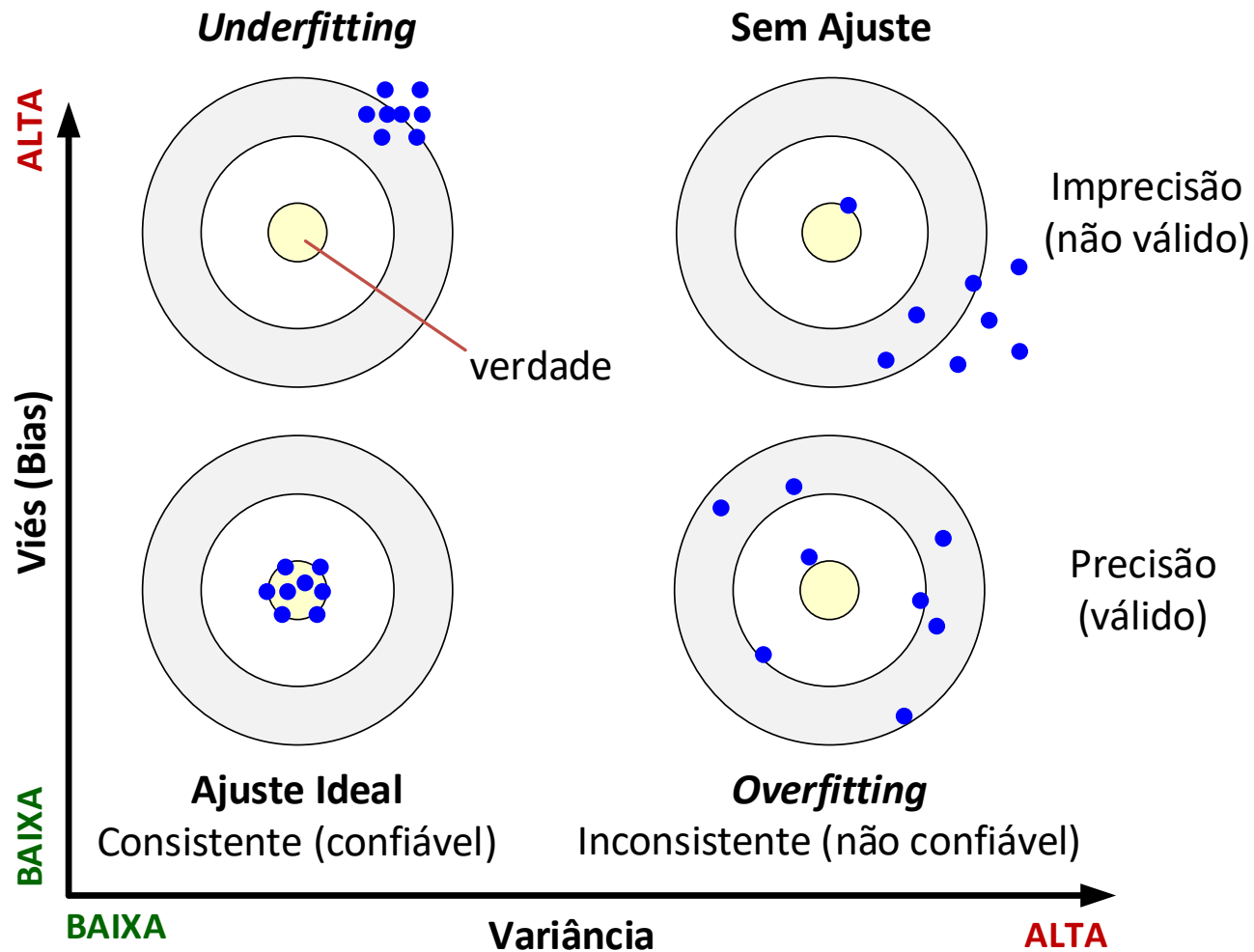
(b) Ajuste adequado



(c) *Overfitting*,
aprendizado do ruído

3. Processo de Aprendizado de Máquina

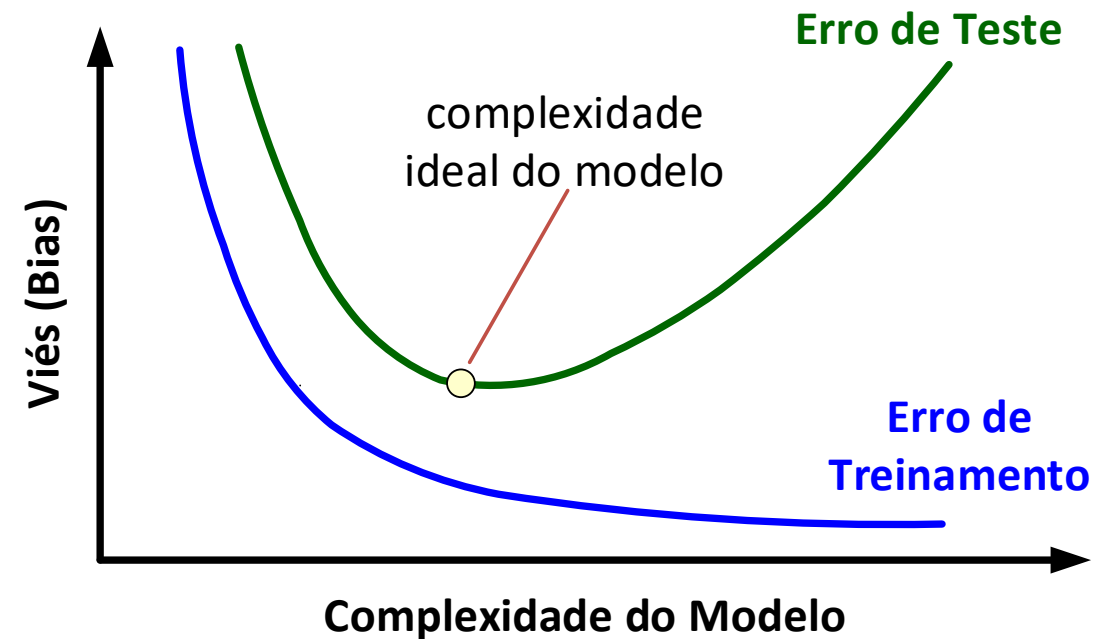
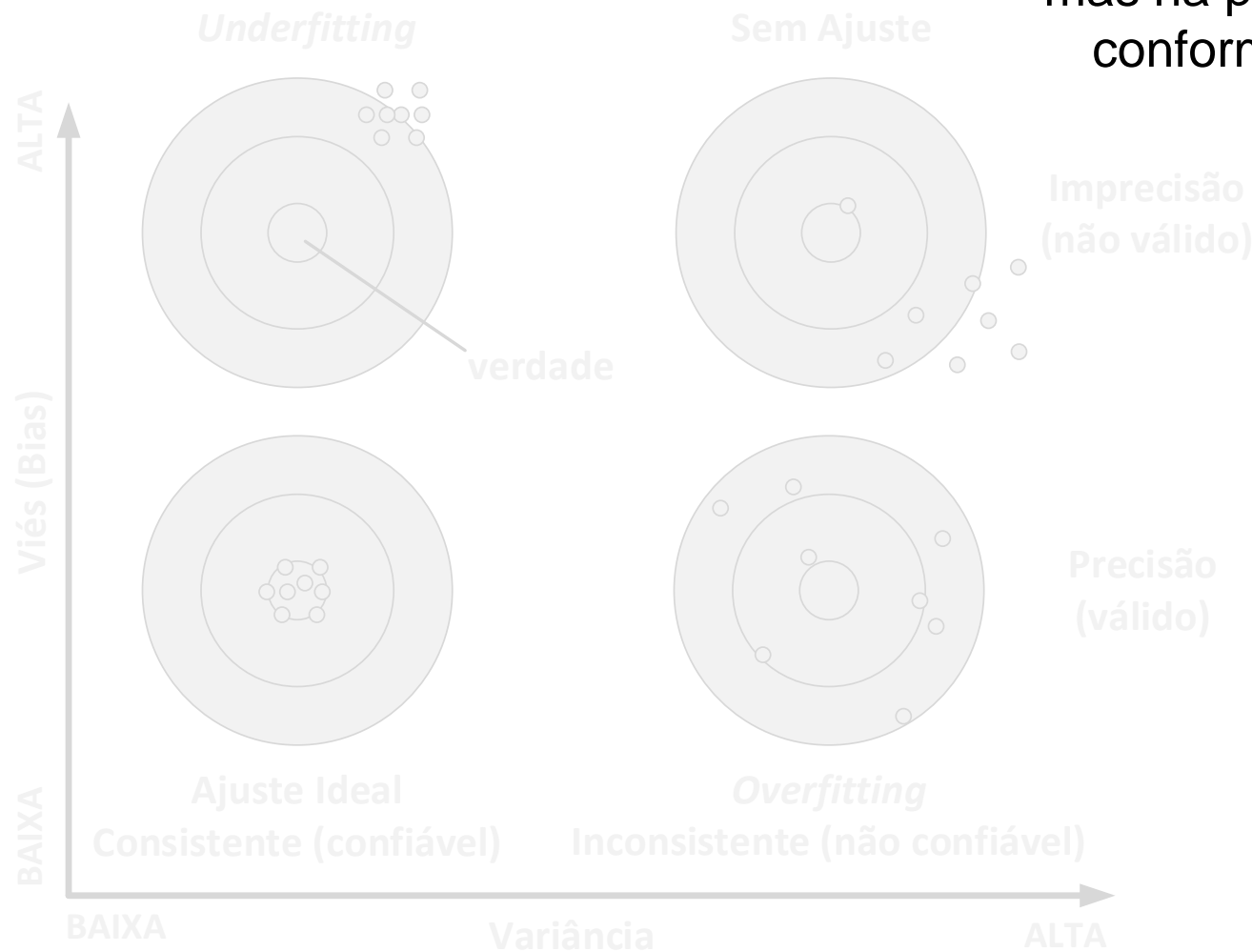
Validade do Modelo



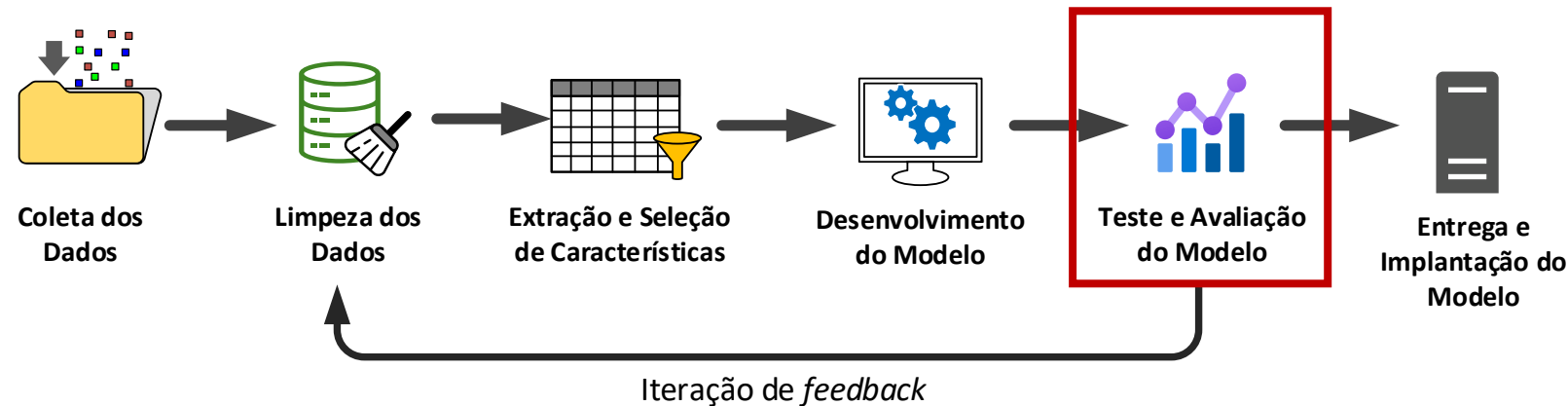
3. Processo de Aprendizado de Máquina

Validade do Modelo

O objetivo é obter um modelo com baixo viés e variância, mas na prática é impossível atender as duas condições conforme se aumenta a complexidade do modelo.



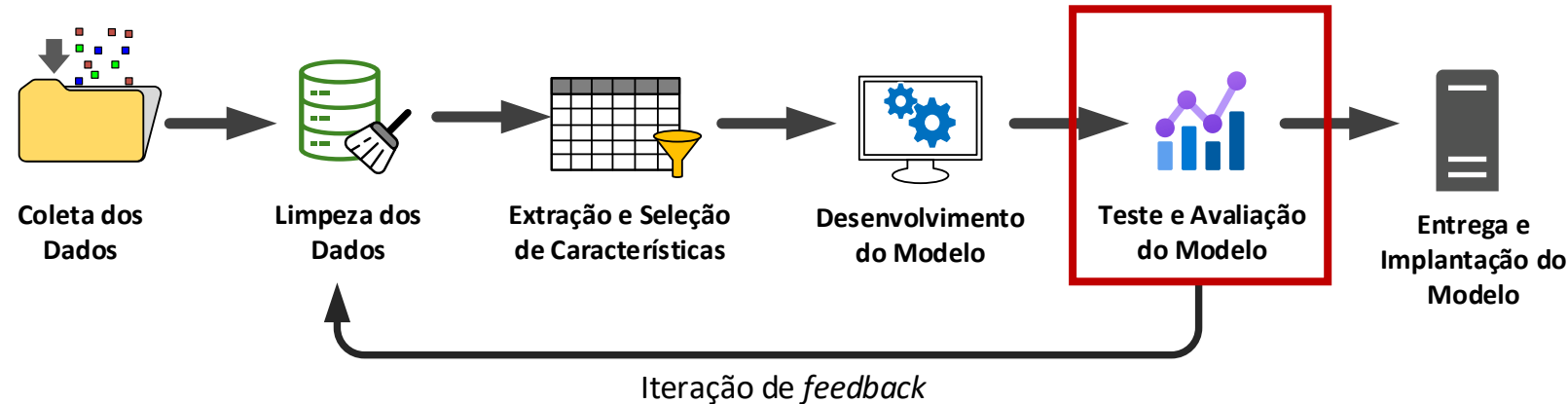
3. Processo de Aprendizado de Máquina



Métodos Não-Supervisionados:

- **Índices Internos:** não se conhece os rótulos das classes. Usados por métodos de agrupamento em tempo de execução como critério de melhoria (ex. Davies-Bouldin, Dunn, Coeficiente de Silhouette, Calinski-Harabasz, Matriz de Contingência).

3. Processo de Aprendizado de Máquina



Métodos Não-Supervisionados:

- **Índices Internos:** não se conhece os rótulos das classes. Usados por métodos de agrupamento em tempo de execução como critério de melhoria (ex. Davies-Bouldin, Dunn, Coeficiente de Silhouette, Calinski-Harabasz, Matriz de Contingência).
- **Índices Externos:** são conhecidos os rótulos das classes. Usados após a execução do método de agrupamento para comparar o desempenho entre métodos (ex. ARI, CRand, Informação Mútua, Pureza, Jaccard, Dice, Cardinalidade.)

3. Processo de Aprendizado de Máquina

Métodos **Supervisionados** – Classificação

Matriz de confusão: Útil para analisar o quanto um modelo reconhece tuplas diferentes. Na forma básica tem duas classes (interesse e outros casos). Cada **quadrante indica** uma das **possibilidades de classificação**.

		VALORES PREDITOS		
		POSITIVO	NEGATIVO	
VALORES REAIS	POSITIVO	Verdadeiros Positivos (VP)	Falsos Negativos (FN)	$P = VP + FN$
	NEGATIVO	Falsos Positivos (FP)	Verdadeiros Negativos (VN)	$N = FP + VN$

3. Processo de Aprendizado de Máquina

Métodos Supervisionados – Classificação

Matriz de confusão: Útil para analisar o quanto um modelo reconhece tuplas diferentes. Na forma básica tem duas classes (interesse e outros casos). Cada quadrante indica uma das possibilidades de classificação.

		VALORES PREDITOS		
		POSITIVO	NEGATIVO	
VALORES REAIS	POSITIVO	Verdadeiros Positivos (VP)	Falsos Negativos (FN)	$P = VP + FN$
	NEGATIVO	Falsos Positivos (FP)	Verdadeiros Negativos (VN)	$N = FP + VN$

Diversas medidas são derivadas a partir da matriz de confusão. 

MEDIDA	EQUAÇÃO
Acurácia, Taxa de Reconhecimento (Acurácia = Sensitividade + Especificidade)	$\frac{VP + VN}{P + N}$
Taxa de Erro, Taxa de Classificação Incorreta	$\frac{FP + FN}{P + N}$
Sensitividade, Taxa de Verdadeiros Positivos, Recall	$\frac{VP}{P}$
Especificidade, Taxa de Verdadeiros Negativos	$\frac{VN}{N}$
Precisão	$\frac{VP}{VP + FP}$
F, F ₁ , F-score (média harmônica da precisão e recall)	$\frac{2 \times \text{Precisão} \times \text{Recall}}{\text{Precisão} + \text{Recall}}$

3. Processo de Aprendizado de Máquina

Métodos Supervisionados – Classificação

Matriz de confusão: Útil para analisar o quanto um modelo reconhece tuplas diferentes. Na forma básica tem duas classes (interesse e outros casos). Cada quadrante indica uma das possibilidades de classificação.

Ex. **tuplas positivas** = imagens com um gato.

tuplas negativas = todas as demais que não possui.

		VALORES PREDITOS		TOTAL
		POSITIVO	NEGATIVO	
VALORES REAIS	POSITIVO	140	30	170
	NEGATIVO	20	10	30
		160	40	200

Acurácia = $(140 + 10) / (170 + 30) = 75\%$

Recall = $140 / 170 = 82,4\%$

Precisão = $140 / (140 + 20) = 87,5\%$

MEDIDA	EQUAÇÃO
Acurácia, Taxa de Reconhecimento (Acurácia = Sensitividade + Especificidade)	$\frac{VP + VN}{P + N}$
Taxa de Erro, Taxa de Classificação Incorreta	$\frac{FP + FN}{P + N}$
Sensitividade, Taxa de Verdadeiros Positivos, Recall	$\frac{VP}{P}$
Especificidade, Taxa de Verdadeiros Negativos	$\frac{VN}{N}$
Precisão	$\frac{VP}{VP + FP}$
F, F ₁ , F-score (média harmônica da precisão e recall)	$\frac{2 \times \text{Precisão} \times \text{Recall}}{\text{Precisão} + \text{Recall}}$

3. Processo de Aprendizado de Máquina

Métodos **Supervisionados** – Regressão

Modelos de qualidade apresentam MAE e MSE próximos a 0, ou igual a 1 no caso do coeficiente de determinação.

Erro Médio Absoluto $MAE = \frac{1}{m} \sum_{i=1}^m |y_i - \hat{y}_i|$

verdadeiros valores-alvo de teste

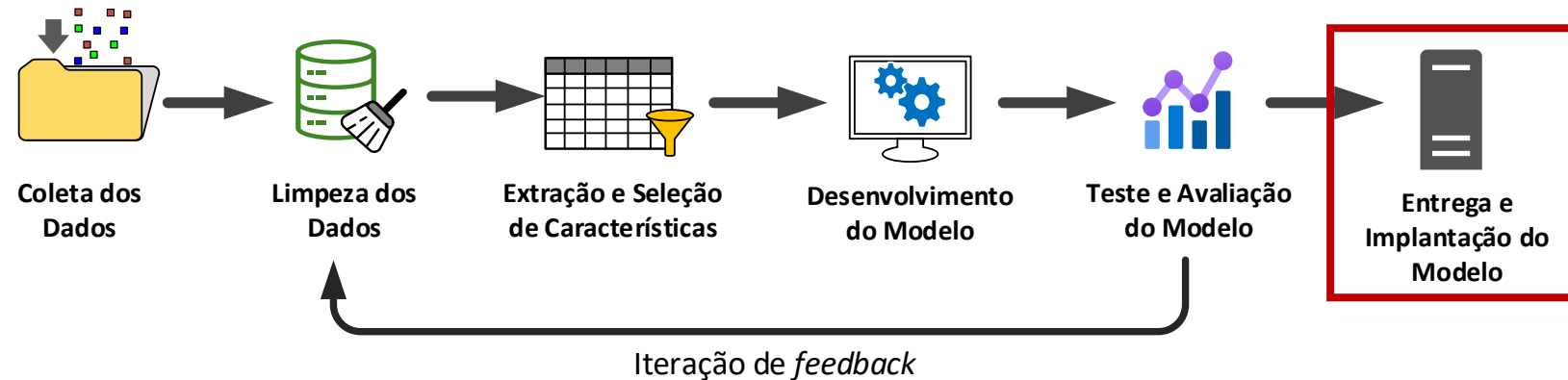
valores previstos correspondentes aos valores-alvo

Erro Médio Quadrático $MSE = \frac{1}{m} \sum_{i=1}^m (y_i - \hat{y}_i)^2$

Coeficiente de Determinação $R^2 = 1 - \frac{RSS}{TSS} = 1 - \frac{\sum_{i=1}^m (y_i - \hat{y}_i)^2}{\sum_{i=1}^m (y_i - \bar{y})^2}$

média dos valores previstos

3. Processo de Aprendizado de Máquina



Visa **disponibilizar o modelo para uso** em um sistema automatizado ou auxílio à tomada de decisão.

Os **resultados** da aplicação do modelo **são analisados por especialistas e usuários finais**, que podem fornecer **feedbacks** para sua melhoria.

Dependendo dos requisitos, esta fase pode ser tão simples como gerar um relatório ou tão complexo como a implementação do modelo como um módulo de um sistema especialista.

4. Parâmetros e Hiperparâmetros

Modelos de aprendizado possuem estratégias e critérios distintos para realização de uma tarefa.

O **parâmetro** é uma **variável de configuração interna ao modelo** cujo valor é definido ou estimado a partir dos dados ou por aprendizado de máquina (ex. pesos em uma RNA, coeficientes de regressão linear ou logarítmica).

4. Parâmetros e Hiperparâmetros

Modelos de aprendizado possuem estratégias e critérios distintos para realização de uma tarefa.

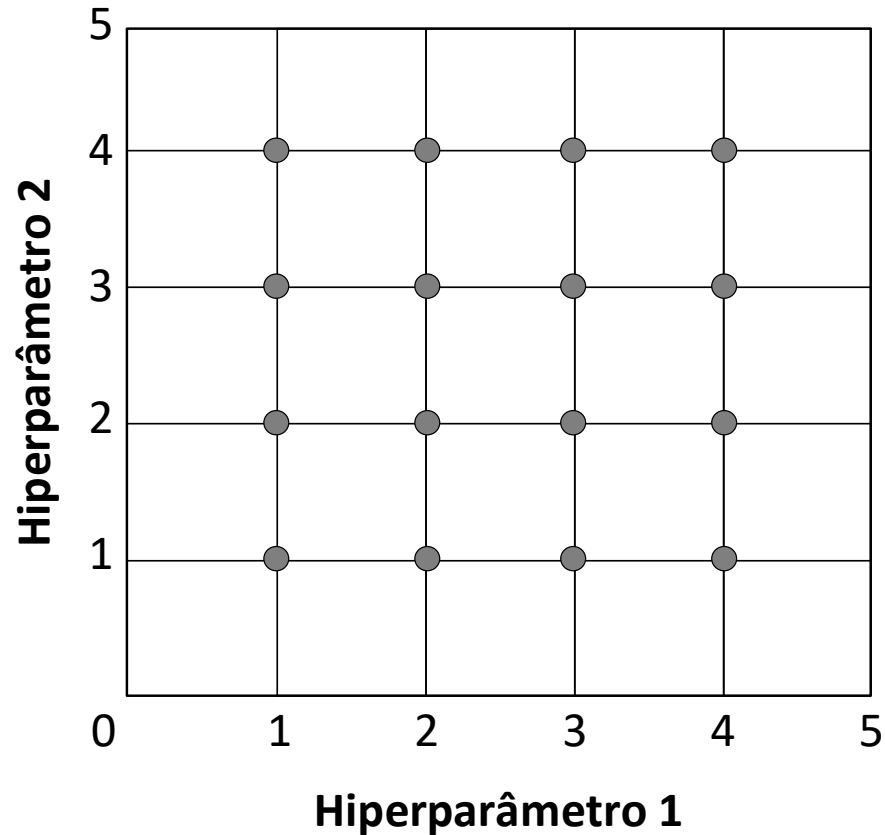
O parâmetro é uma variável de configuração interna ao modelo cujo valor é definido ou estimado a partir dos dados ou por aprendizado de máquina (ex. pesos em uma RNA, coeficientes de regressão linear ou logarítmica).

O treinamento do modelo geralmente se refere à otimização dos seus parâmetros, sendo esse processo concluído por um algoritmo de *Gradient Descent* (ex. BGD, SGD)

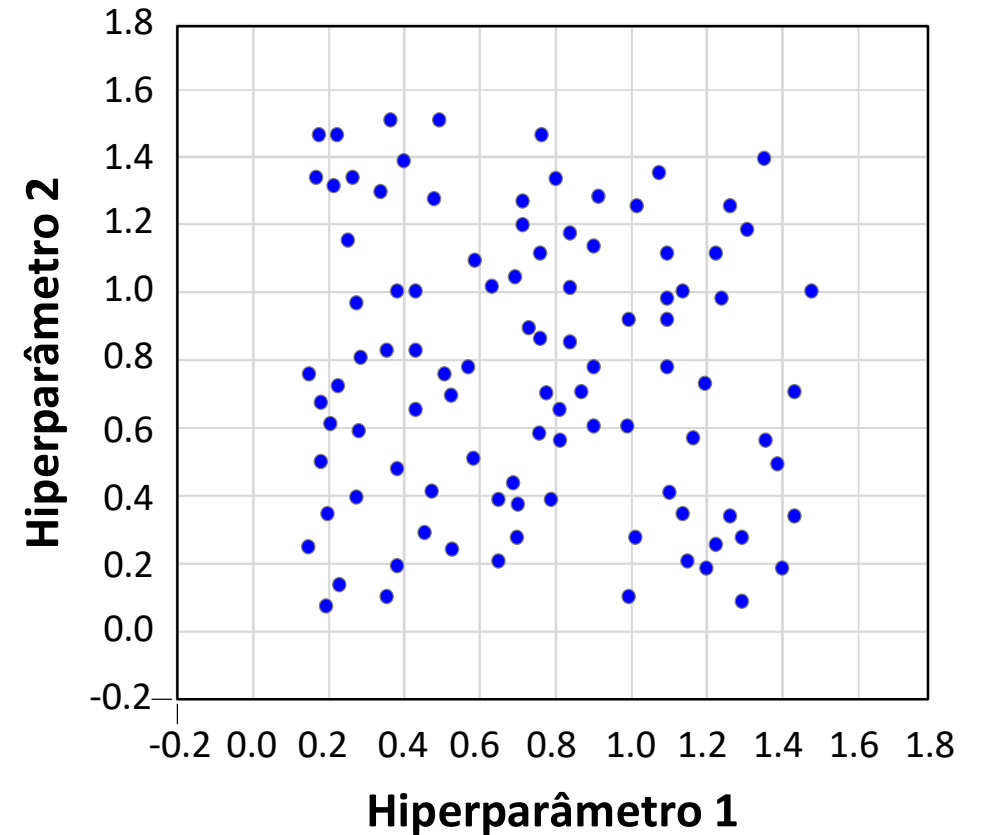
Hiperparâmetro é uma **configuração externa** de um algoritmo de ML, não sendo afetado pelo próprio algoritmo. Definido antes do treinamento, sendo constante durante sua execução (ex. quantidade de camadas de uma RNA, níveis de uma Árvore de Decisão, valor K do KNN).

4. Parâmetros e Hiperparâmetros

Métodos como **Busca em Grade** e **Busca Aleatória** podem auxiliar na descoberta dos hiperparâmetros mais adequados para um modelo.



(a) Busca em Grade

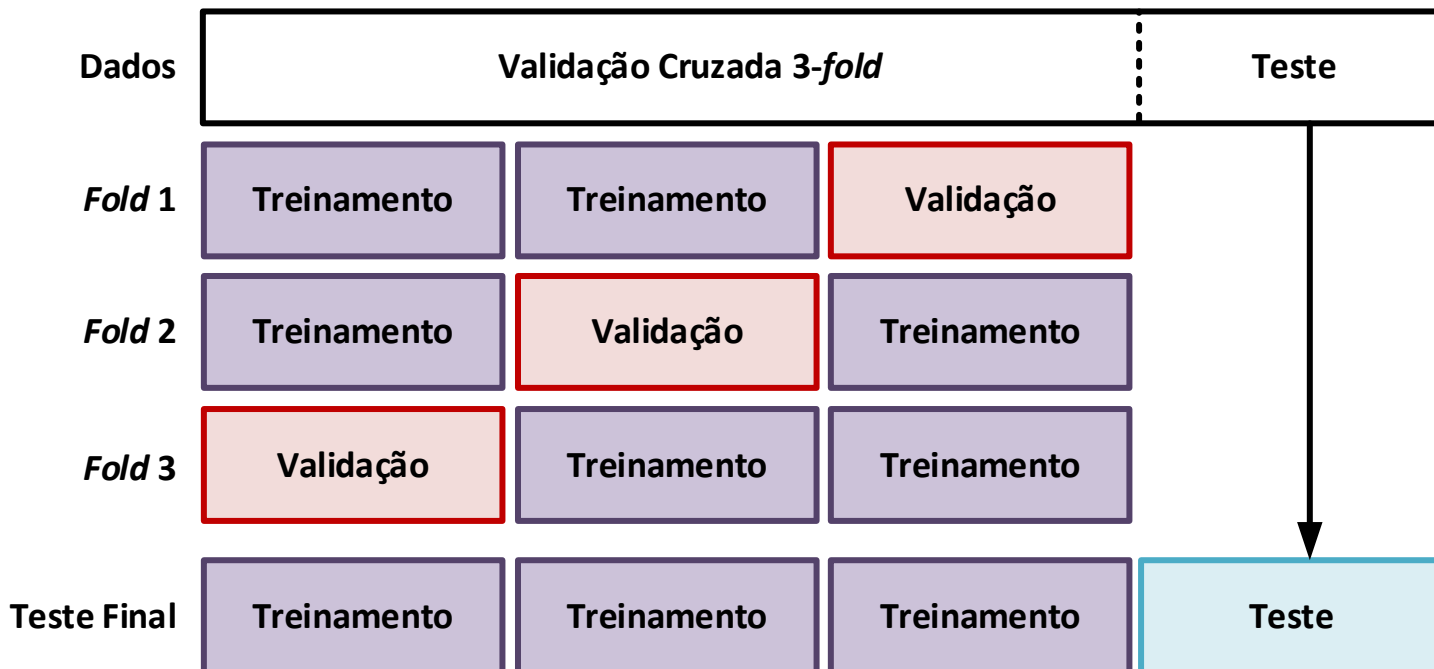


(b) Busca Aleatória

4. Parâmetros e Hiperparâmetros

Validação Cruzada é uma técnica considerada mais robusta que as anteriores. O conjunto de dados é dividido em treinamento e teste.

O conjunto de treinamento é dividido em k pastas. A cada iteração é escolhida uma **partição diferente para validação** do modelo.



Cada amostra é usada o mesmo número de vezes para treinamento, e apenas uma vez para validação.

A **precisão** é dada pelo número total de classificações corretas das k iterações, dividido pelo número total de objetos dos dados originais.

Resumo

Foi vista a **definição de Aprendizado de Máquina** e seus **conceitos centrais** (dados, algoritmo e aprendizado), suas **aplicações** e em quais situações seu uso é justificado.

Foram apresentadas as **estruturas de dados comumente usadas para o ML** e os **tipos de atributos** observados em conjuntos de dados (nominais, binários, categóricos, numéricos).

Descritos os **tipos de Aprendizado de Máquina** (Não-Supervisionados, Supervisionados, Semi-Supervisionados e Aprendizado por Reforço).

Resumo

Foi vista a definição de Aprendizado de Máquina e seus conceitos centrais (dados, algoritmo e aprendizado), suas aplicações e em quais situações seu uso é justificado.

Foram apresentadas as estruturas de dados comumente usadas para o ML e os tipos de atributos observados em conjuntos de dados (nominais, binários, categóricos, numéricos).

Descritos os tipos de Aprendizado de Máquina (Não-Supervisionados, Supervisionados, Semi-Supervisionados e Aprendizado por Reforço).

Foram detalhadas as fases do **Processo de Aprendizado de Máquina** (1. Coleta de Dados, 2. Limpeza de Dados, 3. Extração e Seleção de Características, 4. Desenvolvimento, 5. Teste e Avaliação, 6. Implantação do Modelo) e principais técnicas para sua respectiva execução.

Apresentados os aspectos relacionados a **parâmetros** e **hiperparâmetros** de algoritmos de aprendizado de máquina.

Referências Bibliográficas

- [1] Huawei Technologies Co. Artificial Intelligence Technology - Official Textbooks for Huawei ICT Academy. Morgan Kaufmann, 3rd edition, 2023.
- [2] Marc Peter Deisenroth, A Aldo Faisal, and Cheng Soon Ong. Mathematics for Machine Learning. Cambridge University Press, 2020.
- [3] Aurélien Géron. Hands-on machine learning with Scikit-Learn, Keras, and TensorFlow. O'Reilly Media, Inc., 3rd edition, 2022.
- [4] Trevor Hastie, Robert Tibshirani, JeromeHFriedman, and JeromeHFriedman. The elements of statistical learning: data mining, inference, and prediction, volume 2. Springer, 2009.
- [5] Han Jiawei, Kamber Micheline, and Pei Jian. Data Mining: Concepts and Techniques. Morgan Kaufmann, 3rd edition, 2012.
- [6] Andreas C Müller and Sarah Guido. Introduction to Machine Learning with Python: a guide for data scientists. O'Reilly Media, Inc., 2016.
- [7] Shai Shalev-Shwartz and Shai Ben-David. Understanding Machine Learning: From theory to algorithms. Cambridge university press, 3rd edition, 2014.

Apoio

Projeto Residência 8 (TPA N° 068/SOFTEX/UNIFEI)

Este projeto é apoiado pelo Ministério da Ciência, Tecnologia e Inovações, com recursos da Lei nº 8.248, de 23 de outubro de 1991, no âmbito do [PPI-Softex| PNM-Design], coordenado pela Softex.

