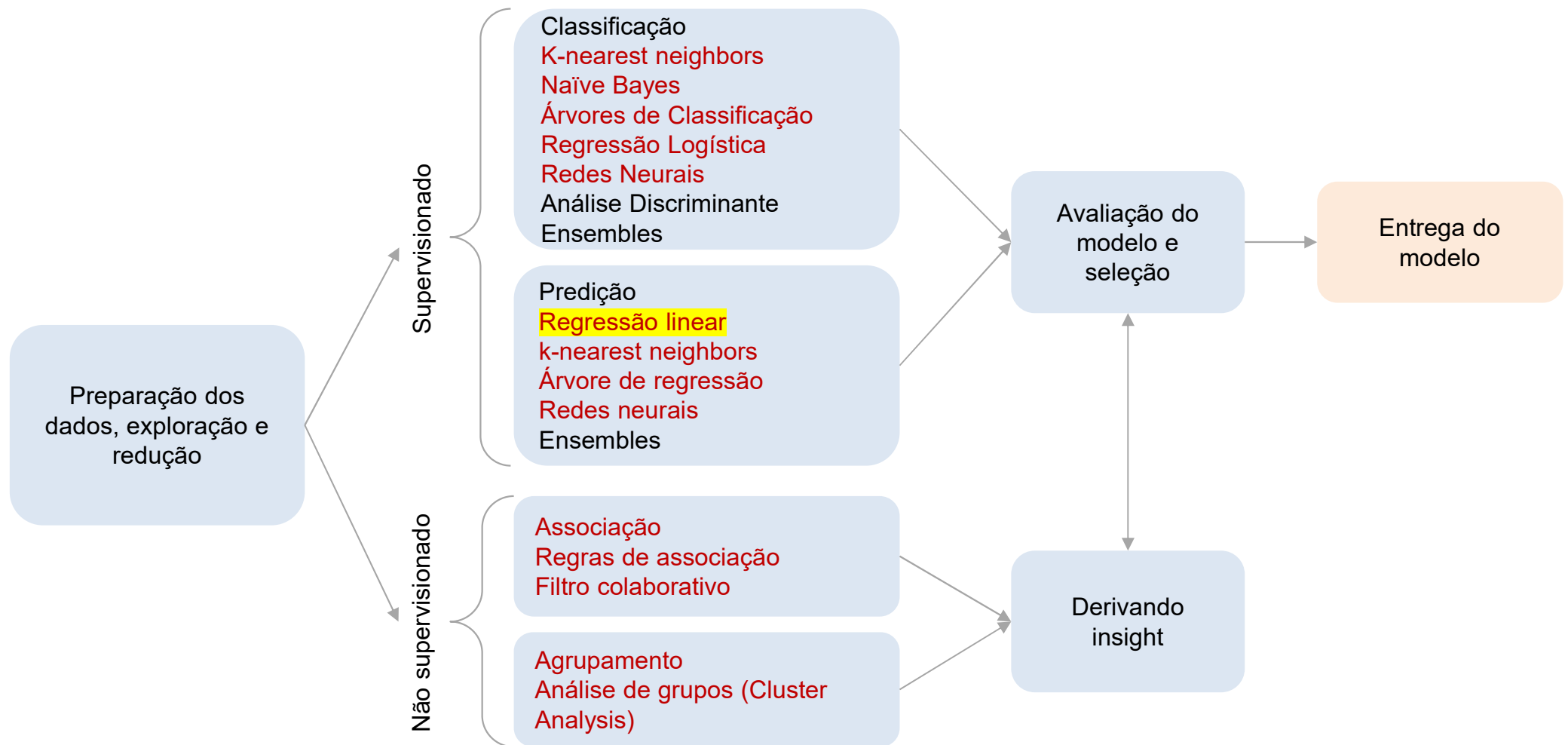


# Agenda

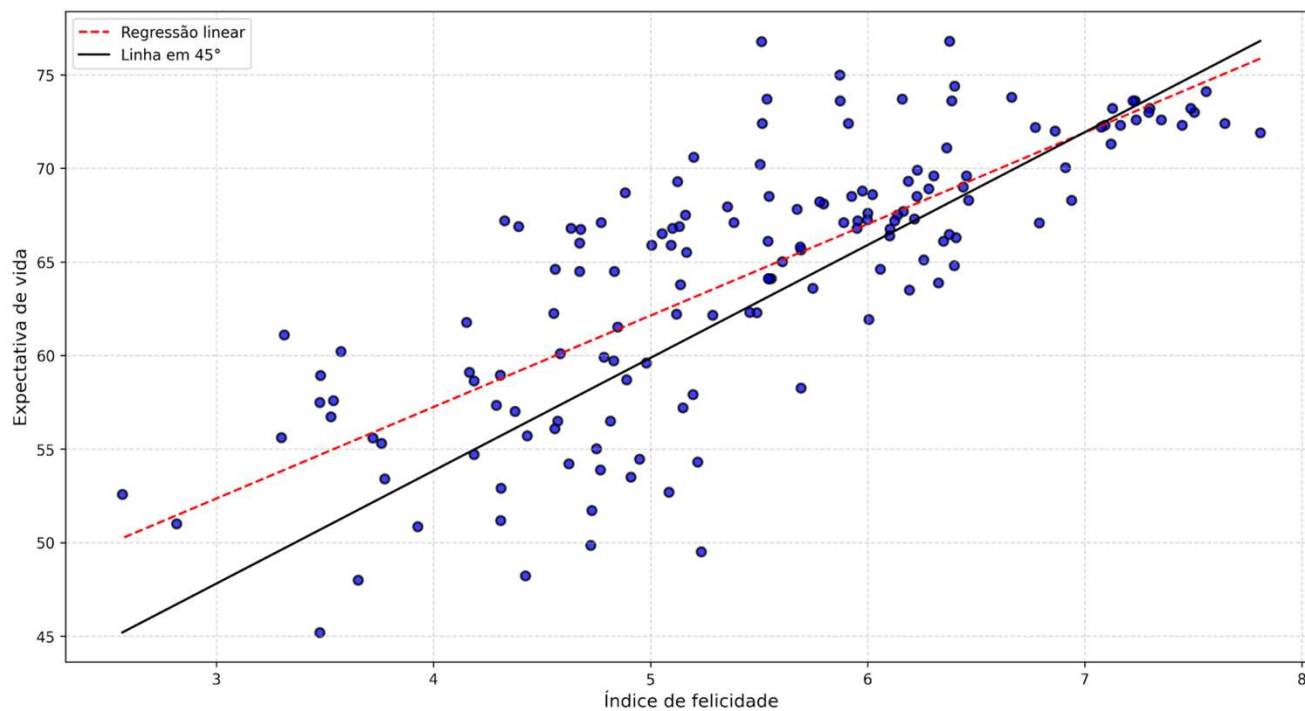
## **UNIDADE 3: Regressão Linear**

- Previsões simples (Regressão linear)
- Previsões complexas (Regressão linear múltipla)



# Relação entre variáveis

- Cenário: **determinar a associação entre duas (ou mais) informações:**
  - Relação entre **índice de felicidade** e **expectativa de vida**
  - Relação entre **número de processos pendentes** e **quantidade de juízes**
- **Variáveis relacionadas** são ditas **correlacionadas**.
- Exemplo de variáveis que aparentemente **não são relacionadas**:
  - Relação entre a **altura de uma criança** e a **de seus pais**
  - Relação entre **número de absolvições** e a **cidade de nascimento dos réus**



Fonte: <https://dataat.github.io/introducao-ao-machine-learning/regressao.html>

É preciso  
estabelecer uma  
relação matemática  
para entender  
como uma variável  
influência outra.

**Para que serve descobrir a relação  
entre variáveis?**

# Regressão x Classificação

- **Regressão:**
  - Verificação se **duas ou mais variáveis estão relacionadas**, como se influenciam.
  - A regressão linear produz uma **previsão numérica**, como base em valores conhecidos.
- **Exemplo de regressão:**
  - Predição de Montantes de Danos em Casos de Indenização
  - Estimativa de Custos Processuais
  - Projeção de Casos de Litígios Trabalhistas
- **Exemplo de classificação**
  - Classificação de Tipos de Casos Jurídicos (Criminal, Civil, Trabalhista, Familiar, etc.)
  - Classificação de Risco de Reincidência Criminal (Alto risco, médio risco, baixo risco)
  - Classificação de Documentos Jurídicos (contratos, petições, sentenças, etc.)

Umidade relativa média do ar

Atributo	Valor
Temperatura	26,4°C
Chuva acumulada	0 mm
Velocidade do vento	2,5 m/s
Radiação solar	2064 kj/m <sup>2</sup>
Sensação térmica	24,7°C
<b>Umidade</b>	<b>54%</b>

$$y = w_0 + w_1 x_1 + w_2 x_2 + \dots$$

## Predição de Cargas de Trabalho para Juízes

Descrição: Prever o volume de trabalho em diferentes tribunais com base em fatores como tamanho da população, tipos de casos comuns na região, etc.

- Tamanho da população da região onde o tribunal está localizado (valores reais entre 30,000 e 500,000)
- Número de juízes ativos no tribunal (valores inteiros entre 1 e 20)
- Número de casos recebidos mensalmente (valores inteiros entre 150 e 800)
- Distribuição dos tipos de casos comuns na região (valores: criminal, civil, trabalhista)
- Taxa de crescimento da população na região (valores reais entre 0.0 e 1.0)
- Taxa de criminalidade na região (valores reais entre 0.0 e 1.0)
- Média de tempo gasto em cada tipo de caso em meses (número inteiro)
- Nível de automação e eficiência do tribunal (valores inteiros entre 0 e 10)
- Nível de congestionamento do sistema judicial na região (baixo, médio, alto)
- Número de advogados atuando na região (valores inteiros entre 1 e 100)
- Número de prédios judiciais na região (valores inteiros entre 1 e 10)
- Nível de urbanização da região (baixo, médio, alto)
- Nível de educação da população na região (valores reais entre 0.0 e 1.0)
- Número de habitantes por juiz (valor fruto da divisão de tamanho da população da região pelo número de juízes ativos no tribunal)
- Índice de Desenvolvimento Humano (valores reais entre 0.0 e 1.0)
- Média de idade dos es (entre 35 e 60)
- **Número de casos pendentes atualmente no tribunal (valores inteiros, coluna que não pode ser nula)**

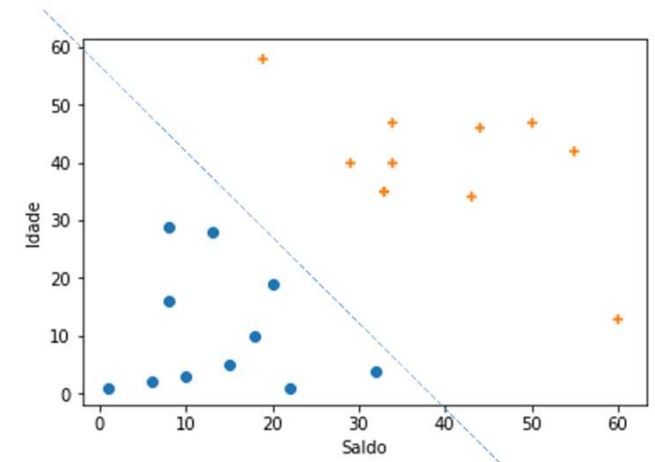
# Modelo linear geral

Ponderação

Característica

$$y = w_0 + w_1x_1 + w_2x_2 + \dots$$

- Vetor de características  $x = (x_1, x_2, \dots, x_n)$
- As ponderações da função linear ( $w_i$ ) são os parâmetros
- A ponderação  $w_0$  é a intersecção



+ Se  $60 - 1.0 \times \text{Idade} - 1.5 \times \text{Saldo} \leq 0$   
• Se  $60 - 1.0 \times \text{Idade} - 1.5 \times \text{Saldo} > 0$

Ajustar o modelo significa encontrar um bom conjunto de ponderações das características.

Quanto maior a magnitude da ponderação mais importante é a característica para a classificação.



# Variáveis Dependentes e Independentes

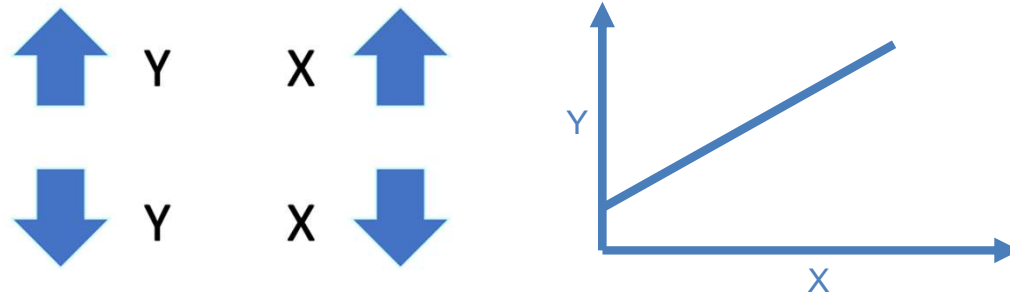
- **Variável dependente** é o valor que estamos **prevendo**
- **Variável independente** é a variável que estamos **usando para prever** uma variável dependente.

$$y = a + bx$$

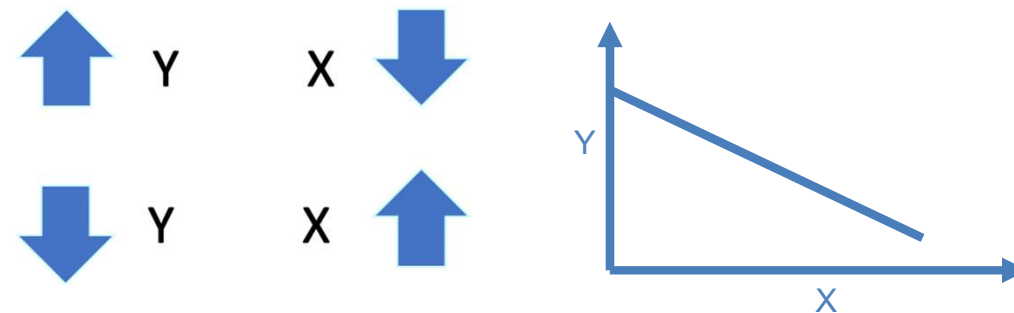
- Uma **variável independente X**, explica a variação em outra variável, que é chamada **variável dependente Y**.
- Esse **relacionamento existe em apenas uma direção**  
**Variável independente (x) → variável dependente (y)**

# Relação entre variáveis Dependentes e Independentes

**Direta (ou positiva)** quando os **valores de Y aumentam** em decorrência do **aumento dos valores de X**.



**Inversa (ou negativa)** quando os **valores de Y variam inversamente** em relação **aos de X**.



# Correlação

- **Correlação**
  - Valor **entre -1 e 1**
  - Quanto mais **próximo de 1 ou -1, mais forte** é a relação
  - Quanto mais **próximo de zero, mais fraca** é a relação
  - O “**valor**” significa a “**força**” da relação
  - O “**sinal**” significa o “**sentido**” da relação
- **1 ou -1 é uma relação perfeita**
- **0 é uma relação inexistente**

# Correlação x Causalidade

- **Correlação**

- Valor **entre -1 e 1**
- Quanto mais **próximo de 1 ou -1, mais forte** é a relação
- Quanto mais **próximo de zero, mais fraca** é a relação
- O “**valor**” significa a “**força**” da relação
- O “**sinal**” significa o “**sentido**” da relação

- **1 ou -1 é uma relação perfeita**

- **0 é uma relação inexistente**

A **correlação nos mostra a força de relacionamento** entre duas variáveis **mas não mostra “como”**.

Se o relacionamento é forte ( $\pm 1$ ) ou não (0)

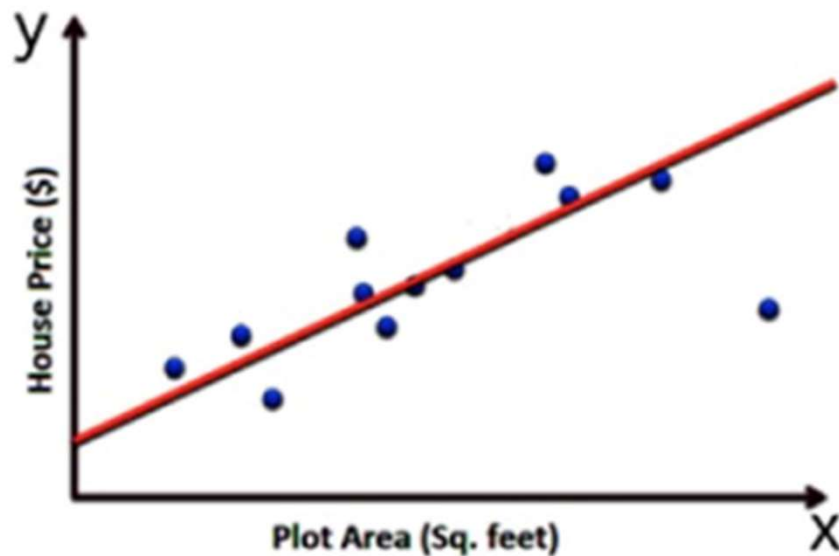
Só porque existe uma **correlação** entre duas variáveis, isso não significa que exista uma **relação causal** entre elas.

O fato de as pessoas usarem guarda-chuvas quando chove não significa que os guarda-chuvas façam a chuva cair.

Uma **análise de regressão** nos permite **começar a ver como**.

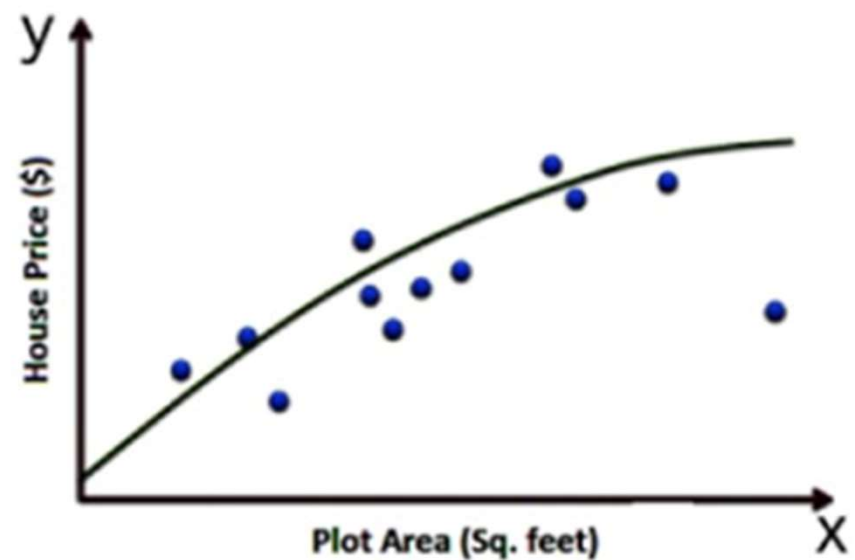
# Modelos de Regressão

Linear Regression



$$y = a + bx$$

Multiple Regression (Polynomial)



$$y = w_0 + w_1x_1 + w_2x_2 + \dots$$

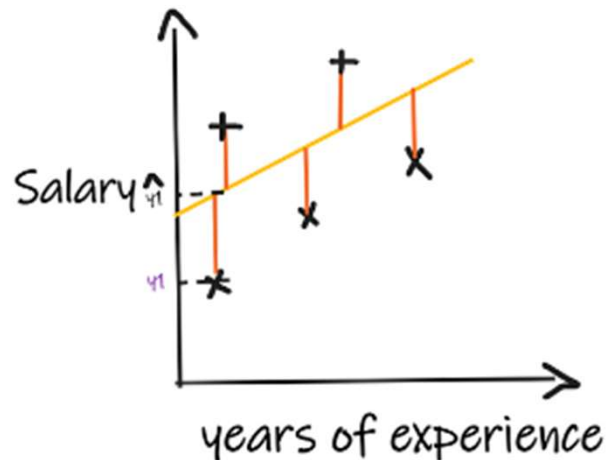


# Como avaliar um modelo

- Como avaliar a qualidade do modelo?
  - Os erros têm distribuição normal?
  - Existem “outliers” no conjunto de dados?
  - O modelo gerado é adequado?

# Métricas de avaliação: regressão

- **Mean Absolute Error (MAE):** Erro Absoluto Médio é a **média do valor absoluto dos erros (diferenças absolutas entre os valores preditos e os valores observados)**.

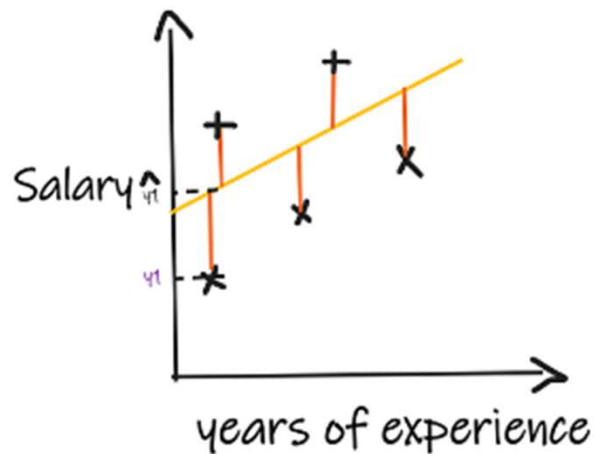


$$MAE = \frac{\sum_{i=1}^n |y_i - \hat{y}_i|}{n}$$



# Métricas de avaliação: regressão

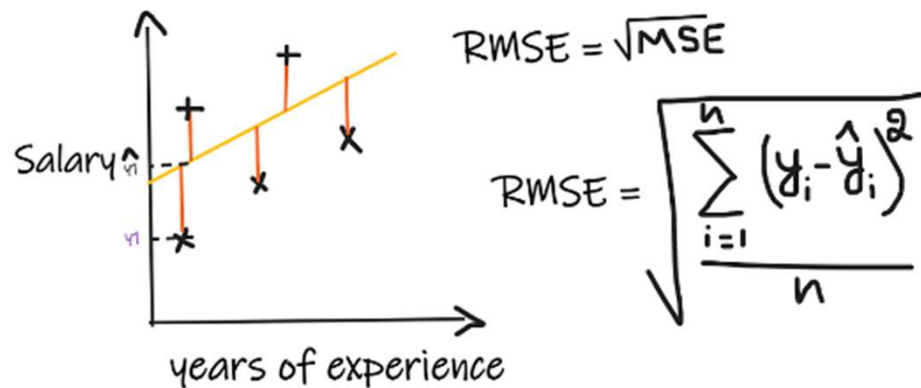
- **Mean Squared Error (MSE):** Erro Médio Quadrático é a **média dos erros quadrados, pune erros maiores**. MSE é preferido quando se deseja **penalizar mais fortemente os erros maiores**, o que é comum em modelos onde grandes desvios são particularmente indesejáveis.



$$MSE = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}$$

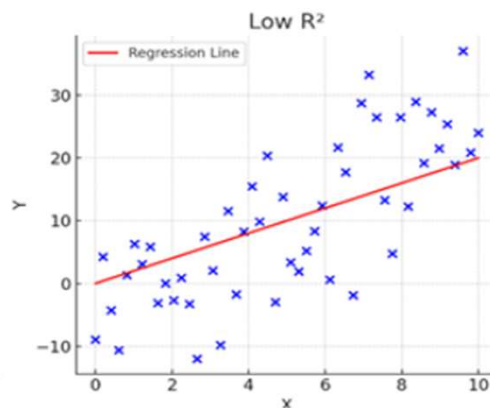
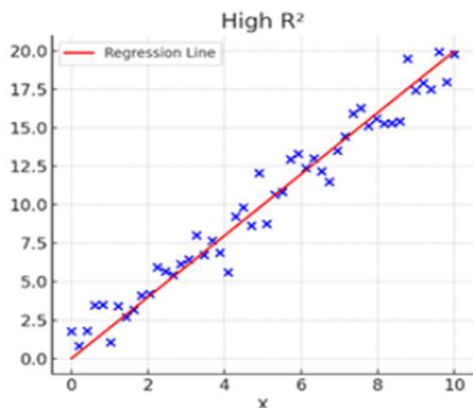
# Métricas de avaliação: regressão

- **Root Mean Square Error (RMSE):** Raiz do Erro Quadrático Médio é a **raiz quadrada da média dos erros quadrados**. RMSE é frequentemente usado em contextos onde é importante **manter as unidades do erro comparáveis com os dados originais** e ao mesmo tempo penalizar erros maiores. O RMSE, em particular, **nos dá uma ideia do erro médio em relação às unidades dos dados originais**.



# Métricas de avaliação: regressão

- **Coeficiente de Determinação ( $R^2$ ):** Ele é uma medida de **quão bem os valores preditos se ajustam aos valores reais**. O  $R^2$  varia de 0 a 1. Um valor de  $R^2$  de 1 indica que o modelo explica 100% da variância nos dados; um valor de 0.75 que explica 75% dos dados, um valor de 0 indica que o modelo não explica nenhuma variância em relação aos valores reais.



## Sinais de Overfitting

- Treino: BAIXO erro (MAE, MSE, RMSE) e ALTO  $R^2$
- Teste: ALTO erro (MAE, MSE, RMSE) e BAIXO  $R^2$

## Sinais de Underfitting

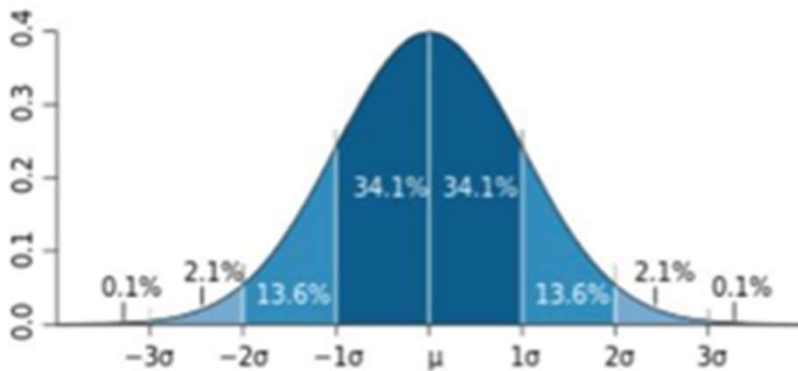
- Treino: ALTO erro (MAE, MSE, RMSE) e BAIXO  $R^2$
- Teste: ALTO erro (MAE, MSE, RMSE) e BAIXO  $R^2$

# Diagnóstico de Resíduos (erro)

- Analise a distribuição dos resíduos (**diferença entre valores previstos e reais**)  $e_i = y_i - \hat{y}_i$ 
  - Overfitting:** **PODE** apresentar resíduos sistematicamente **distribuídos em uma direção específica ou grandes resíduos em dados de teste**.
  - Underfitting:** **Resíduos altos e de distribuição uniforme** **PODE** indicar que o modelo não está capturando bem a variação nos dados.

Quando o erro segue a distribuição normal, podemos garantir que nossos parâmetros beta estimados com mínimos quadrados são iguais a uma estimativa gerada via máxima verossimilhança.

A recíproca não é verdadeira, ou seja, a não normalidade dos resíduos não nos permite afirmar que, com certeza, nossos parâmetros estão errados



O histograma de resíduos deve ser semelhante a uma normal

- Se os erros tiverem uma **distribuição normal**,
- Aproximadamente, 95% dos resíduos estarão no intervalo de um desvio padrão da média**
- Caso contrário, deve existir a presença de “outlier”



# Alguns outros modelos de regressão linear

## Lasso

A regressão lasso adiciona uma penalização do tipo L1 ao termo de erro da regressão linear, que é a soma das magnitudes dos coeficientes (elimina variáveis irrelevantes)

$$RSS_{\text{lasso}} = \sum_{i=1}^n [y_i - (\mathbf{w} \cdot \mathbf{x}_i + b)]^2 + \alpha \sum_{j=1}^p |w_j|$$

regularização  $\ell_1$

Soma dos quadrados dos resíduos + penalidade \* |inclinação|

## Ridge

A regressão ridge adiciona uma penalização do tipo L2 ao termo de erro da regressão linear, que é o quadrado da magnitude dos coeficientes.

$$RSS_{\text{ridge}} = \sum_{i=1}^n [y_i - (\mathbf{w} \cdot \mathbf{x}_i + b)]^2 + \alpha \sum_{j=1}^p w_j^2$$

regularização  $\ell_2$

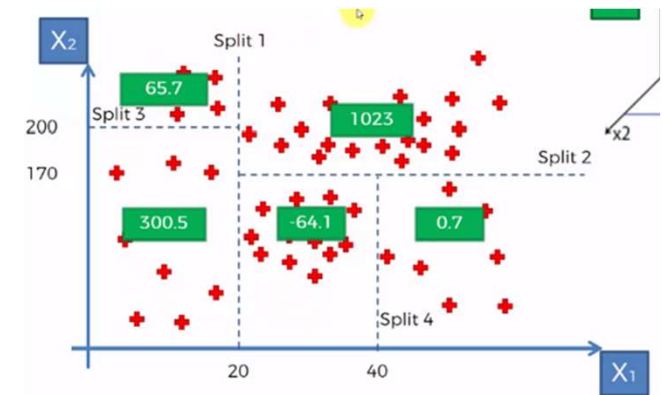
Soma dos quadrados dos resíduos + penalidade \* (inclinação)<sup>2</sup>

## Elastic-Net Regression

Combinação das técnicas de regularização L1 (Lasso) e L2 (Ridge).

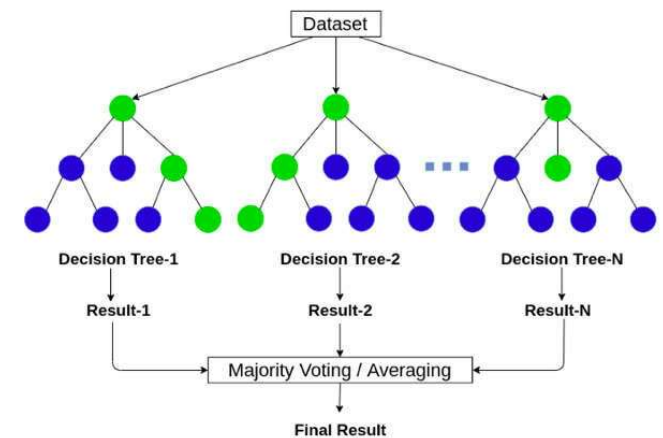
## DecisionTreeRegressor

O critério para dividir um nó pode ser baseado em diferentes métricas de erro, como o mean squared error (MSE) ou mean absolute error (MAE)



## RandomForestRegressor

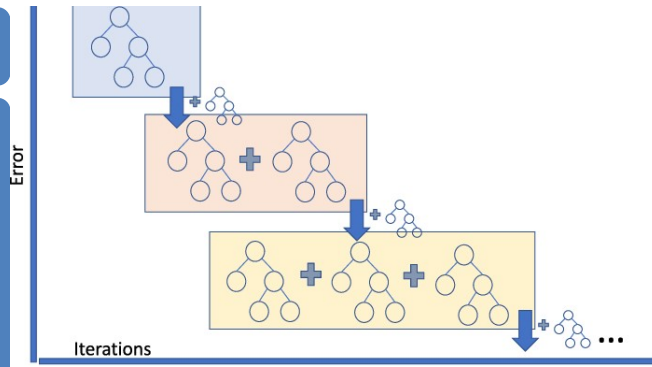
Baseado em florestas aleatórias que são métodos de ensemble que combina múltiplas árvores de decisão fazendo a média das previsões de todas as árvores para produzir uma previsão final.



# Alguns outros modelos de regressão linear

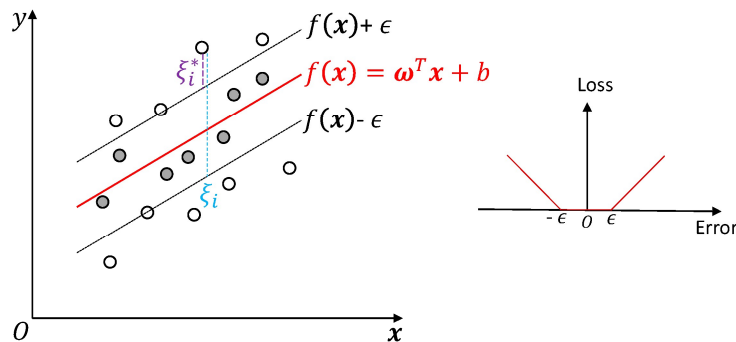
## GradientBoostingRegressor

Baseado em gradiente para problemas de regressão, adicionando iterativamente modelos fracos (geralmente árvores de decisão) para corrigir os erros dos modelos anteriores.



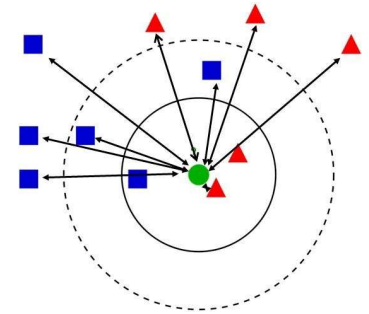
## Support Vector Regression

Implementação do algoritmo de regressão baseado em SVM, tenta encontrar uma função que tenha no máximo uma margem de erro epsilon para todos os pontos de treinamento.



## KNeighborsRegressor

Implementação do algoritmo de regressão baseado nos vizinhos mais próximos (k-Nearest Neighbors, k-NN).



## BayesianRidge

Combina princípios de regressão ridge com probabilidade bayesiana para fornecer estimativas dos coeficientes de regressão.

$$p(y|\lambda) = N(w|0, \lambda^{-1}I_p)$$

## MLPRegressor

Implementação de uma rede neural feedforward (perceptron multicamada) para problemas de regressão.

