

*Soluções em Mineração de Dados*

---

# Metodologias de projetos de Data Mining

Prof. Leandro M. Almeida  
[lma3@cin.ufpe.br](mailto:lma3@cin.ufpe.br)

---



---

# Metodologia de Projeto de DM

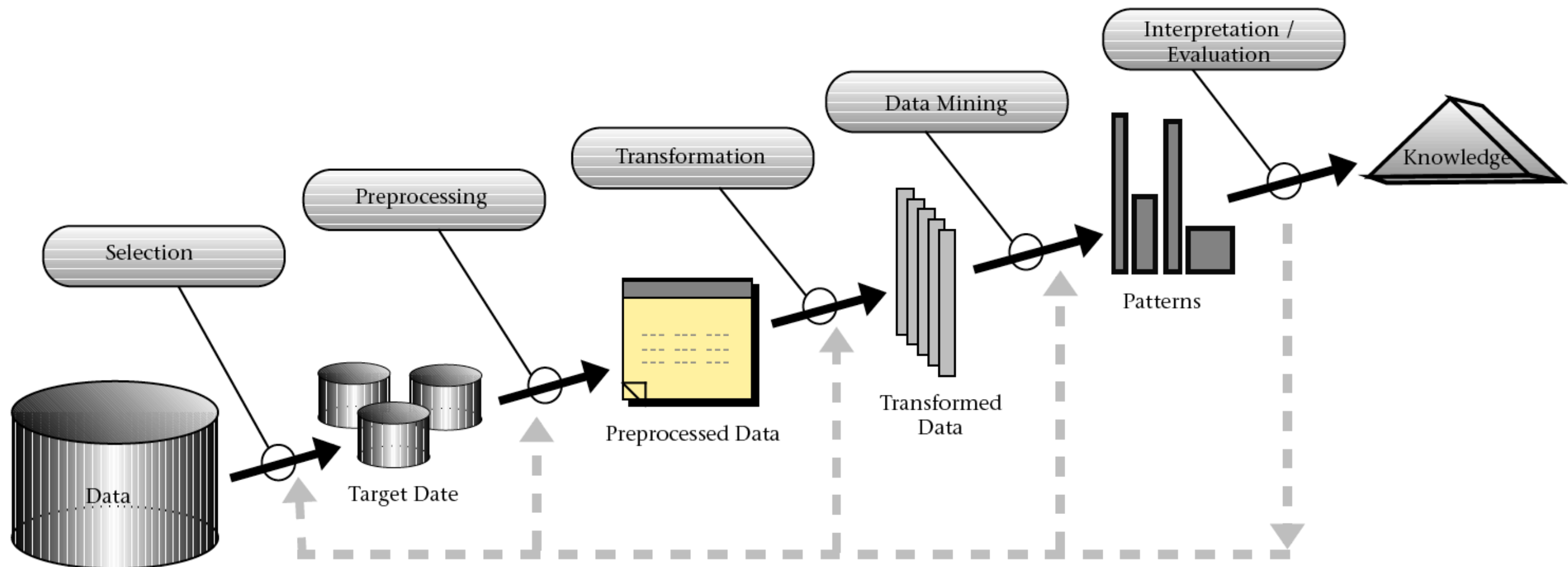
---

- ❖ Na literatura são encontradas metodologias para o desenvolvimento de projetos de mineração com o propósito de guiar os interessados
- ❖ As principais metodologias existentes são:
  - ❖ KDD (knowledge-discovery in databases)
  - ❖ CRISP-DM (Cross Industry Standard Process for Data Mining )
  - ❖ SEMMA (Sample, Explore, Modify, Model, and Assess)



# KDD

- ❖ Processo estabelecido em 1989 com base para a busca por conhecimento em dados, enfatizando a aplicação em alto nível da mineração de dados





---

# KDD

---

- ❖ Seleção

- ❖ possui impacto significativo sobre a qualidade do resultado final
- ❖ Definição do conjunto de dados contendo todas as possíveis variáveis (também chamadas de características ou atributos)
- ❖ Normalmente essa escolha dos dados fica a critério de um especialista do domínio, ou seja, alguém que realmente entende do assunto em questão.

- ❖ Pré-processamento e Limpeza

- ❖ Realizar tarefas que eliminem dados redundantes e inconsistentes, recuperem dados incompletos e avaliem possíveis dados discrepantes ao conjunto, chamados de outliers

- ❖ Transformação dos Dados

- ❖ Os dados necessitam ser armazenados e formatados adequadamente para que os algoritmos possam ser aplicados (normalização, conversão de categóricos para binário, etc)

- ❖ Data Mining

- ❖ Interpretação e Avaliação



---

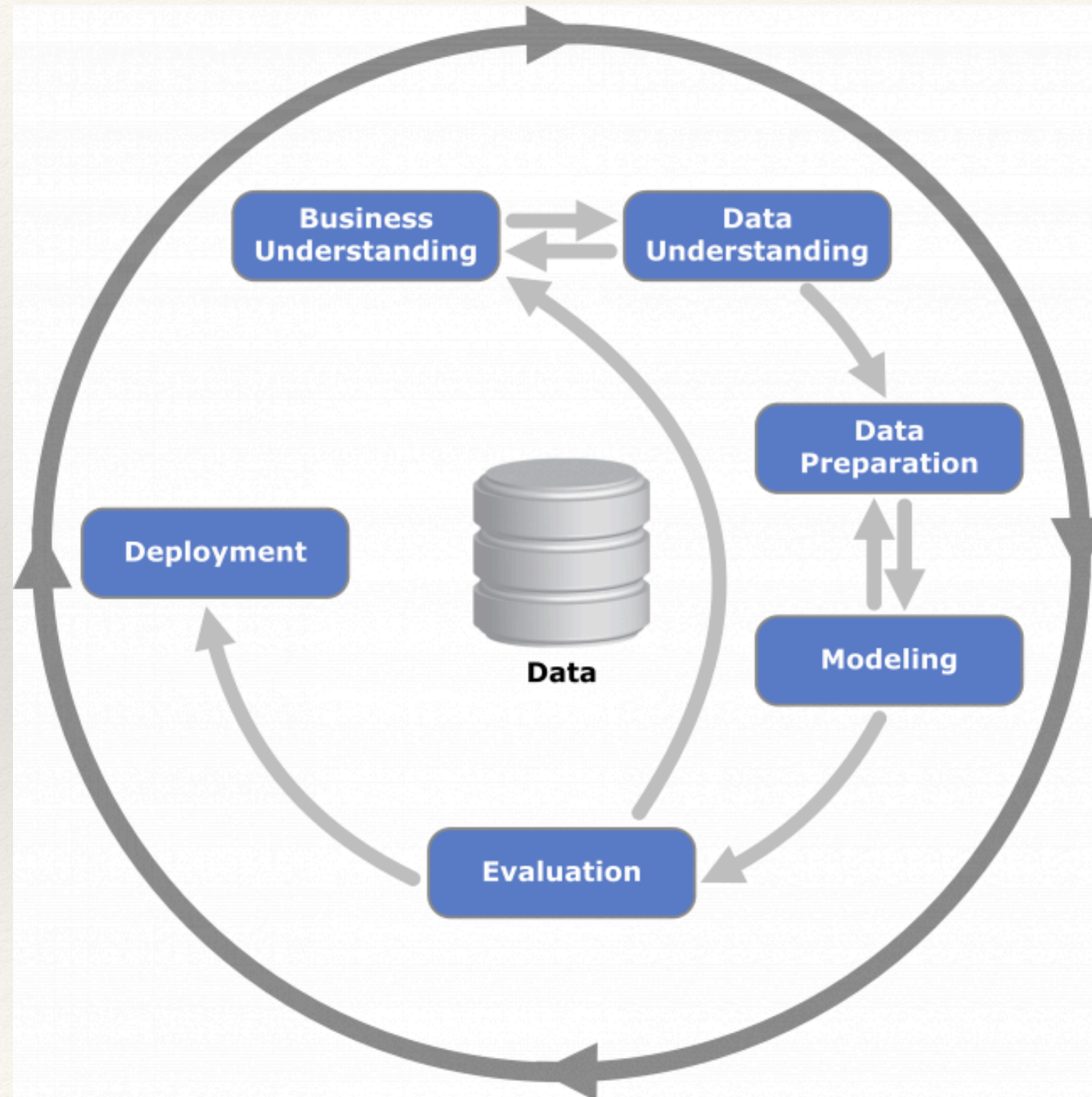
# KDD

---

- ❖ Data Mining
  - ❖ Execução de diferentes algoritmos para a descoberta de padrões de acordo com o propósito do projeto
- ❖ Interpretação e Avaliação
  - ❖ Criar relatórios com gráficos, estatísticas e testes que corroborem o resultado obtido
  - ❖ Apresentar em linguagem não-técnica quais foram os padrões extraídos e quais as possíveis condutas a serem tomadas com o conjunto de informações / conhecimentos obtidos a partir dos dados



# CRISP-DM





---

# CRISP-DM

---

- ❖ Desenvolvida em 1996 com o objetivo de trabalhar com Big Data para descoberta de conhecimento;
- ❖ Consiste em um ciclo com 6 fases:
  1. Entendimento do negócio - buscar uma compreensão adequada do problema que necessita ser resolvido
    - ❖ É preciso buscar detalhes sobre como a questão afeta a organização e quais são os principais objetivos e expectativas em relação ao trabalho como um todo.
  2. Compreensão dos dados
    - ❖ Inspecionar, organizar e descrever todos os dados disponíveis
  3. Preparação dos dados
    - ❖ Preparar todas as databases, definir o formato que será necessário para a análise e ajustar demais questões técnicas



---

# CRISP-DM

---

## 3. Modelagem

- ❖ São selecionadas e aplicadas as técnicas de mineração de dados mais apropriadas, dependendo dos objetivos identificados na primeira fase

## 4. Avaliação

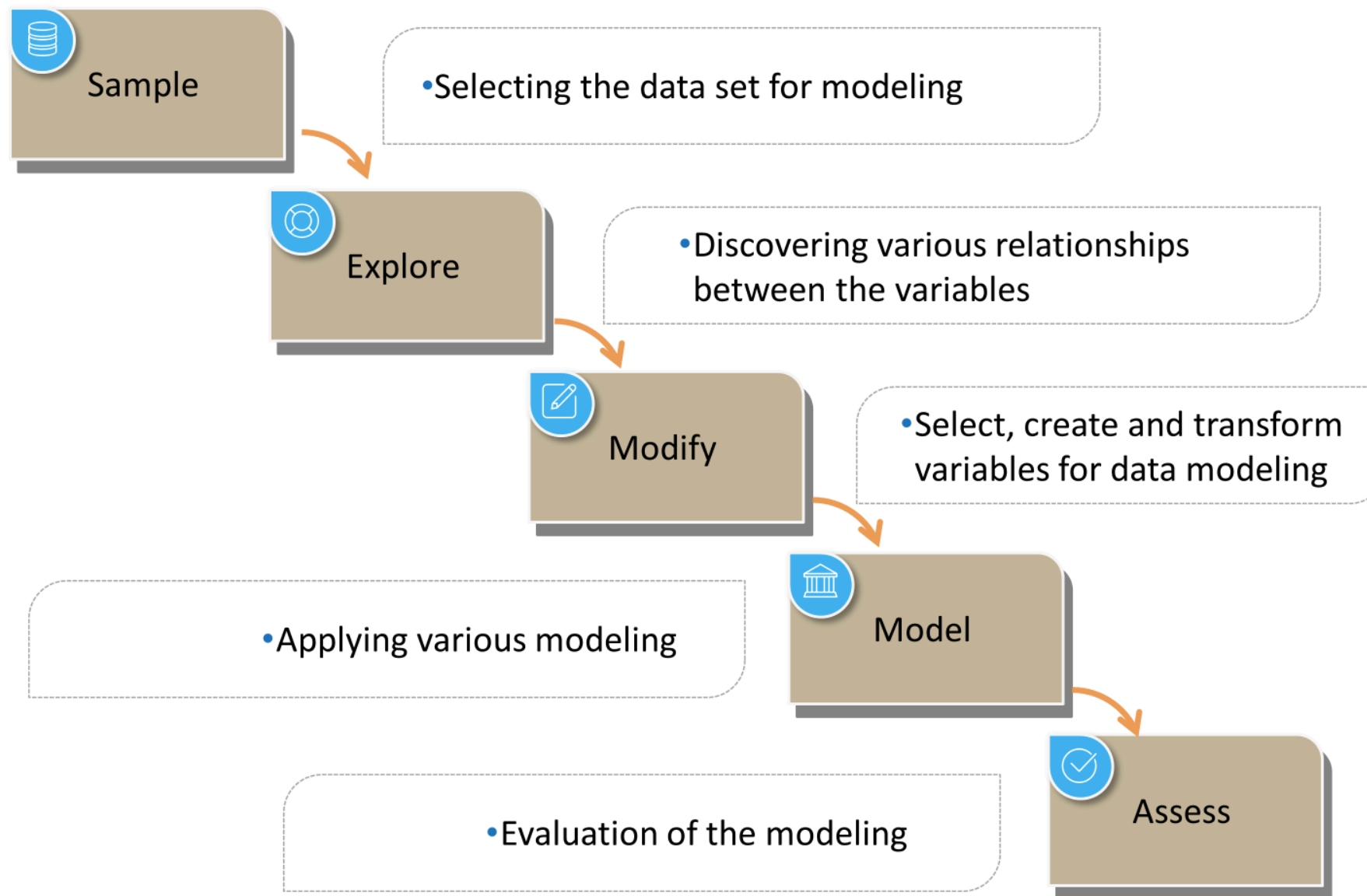
- ❖ Avaliação da aplicabilidade confiável dos insights e conhecimentos obtidos

## 5. Desenvolvimento (deploy)

- ❖ Todo o conhecimento que for obtido por meio do trabalho de mineração e modelagem agora poderá ser aplicado de forma prática. O ideal aqui é dar uma entrega mais palpável e aplicável ao cliente a partir das análises dos dados feitas pela equipe.
- ❖ Algumas das expectativas que se pode ter a partir deste passo é a mudança de processos da empresa ou criação de novos produtos.



# SEMMA





---

# SEMMA

---

- ❖ Amostragem
  - ❖ Selecionar o conjunto de dados para modelagem.
- ❖ Exploração
  - ❖ Compreensão dos dados, descoberta de relações antecipadas e imprevistas entre as variáveis e também anormalidades, com a ajuda da visualização de dados.
- ❖ Modificação
  - ❖ Usa métodos para selecionar, criar e transformar variáveis na preparação para modelagem de dados.
- ❖ Modelagem
  - ❖ Aplicação de várias técnicas de extração de modelos (mineração de dados) nas variáveis preparadas
- ❖ Avaliação
  - ❖ Avaliação dos resultados da modelagem mostra a confiabilidade e utilidade dos modelos criados.

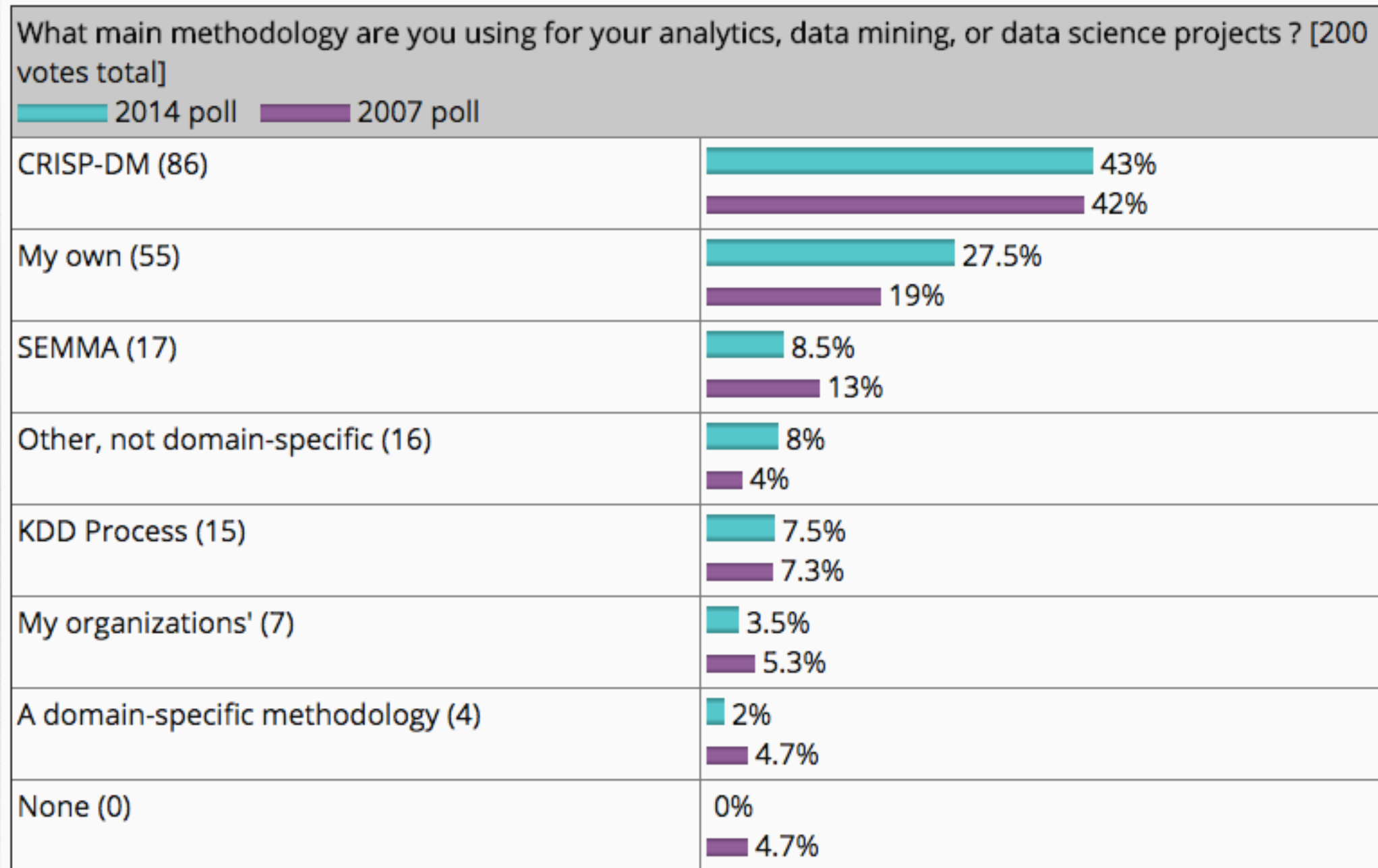


# Comparação dos processos de DM

KDD	SEMMA	CRISP-DM
Pre KDD	-----	Business understanding
Selection	Sample	Data Understanding
Pre processing	Explore	
Transformation	Modify	Data preparation
Data mining	Model	Modeling
Interpretation/Evaluation	Assessment	Evaluation
Post KDD	-----	Deployment



# Comparação dos processos de DM





---

# Comparação dos processos de DM

---

- ❖ O SEMMA e o CRISP-DM são implementações do processo de KDD;
- ❖ CRISP-DM é mais completo que o SEMMA;
- ❖ Muitas empresas vem adotando o CRISP-DM para o desenvolvimento de soluções de mineração devido a sua completude;
- ❖ Os futuros avanços de metodologias de DM estão relacionados a linguagens baseadas em SQL e XML ainda em desenvolvimento.



---

# Avaliação de classificadores

---

- ❖ A construção de classificadores de dados usa um conjunto de dados com rótulos conhecidos;
- ❖ A validação cruzada é uma técnica para avaliar a capacidade de generalização de um modelo;
- ❖ Amplamente empregada em problemas onde o objetivo da modelagem é a classificação;
- ❖ Verificar o seu desempenho para um novo conjunto de dados



---

# Validação cruzada

---

- ❖ Particiona o conjunto de dados em subconjuntos mutualmente exclusivos
- ❖ Utiliza alguns subconjuntos para a estimação dos parâmetros do modelo (dados de treinamento)
- ❖ O restante dos subconjuntos (dados de validação ou de teste) são empregados na validação do modelo.
- ❖ Diversas formas de realizar a divisão dos dados foram sugeridas, sendo as três mais utilizadas: o método *holdout*, o *k-fold* e o *leave-one-out*



---

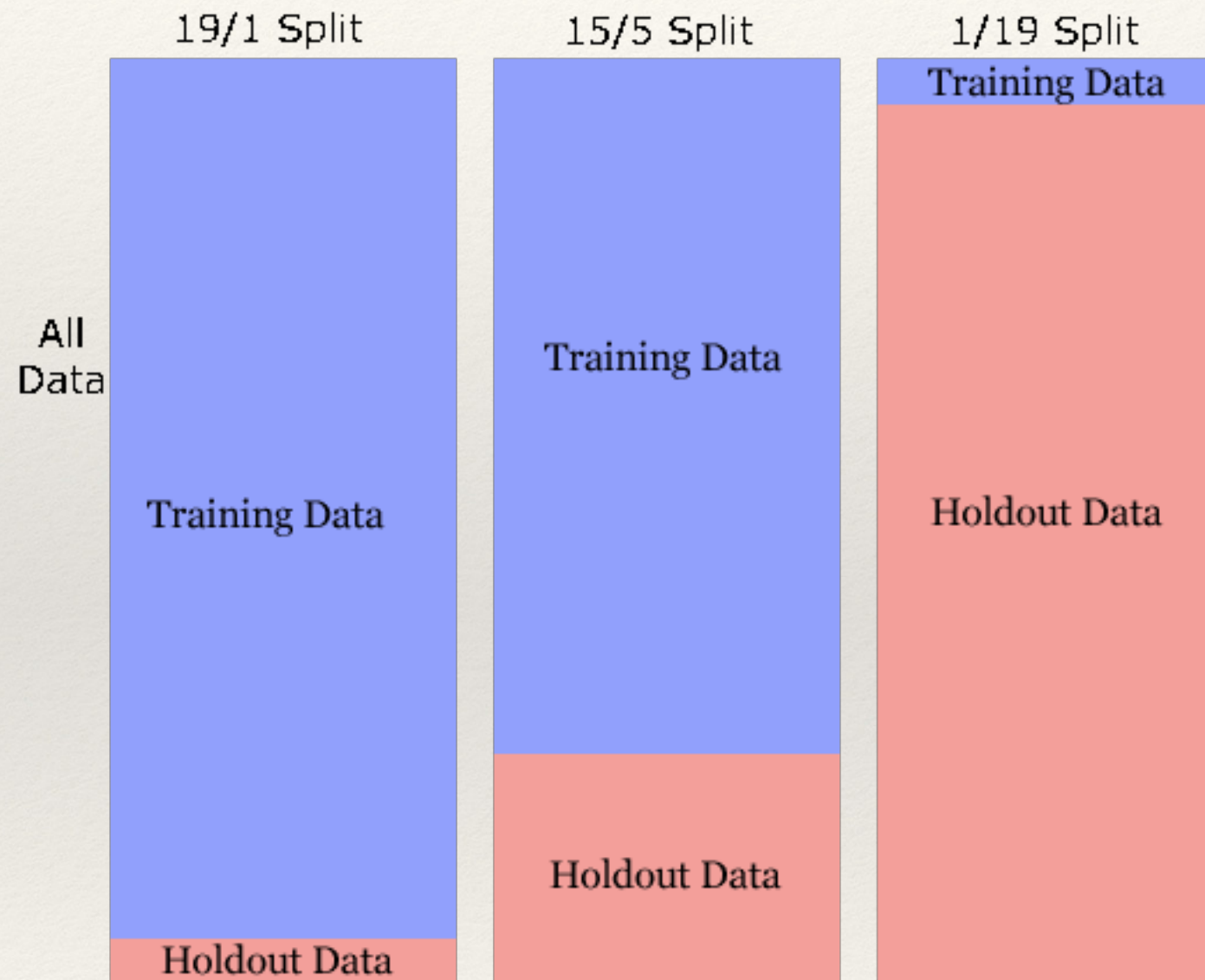
# Validação cruzada

---

- ❖ Método *holdout*:
  - ❖ Divide o conjunto total de dados em dois subconjuntos mutuamente exclusivo
  - ❖ O conjunto de dados pode ser separado em quantidades iguais ou não
  - ❖ Após o particionamento, a estimação do modelo é realizada e, posteriormente, os dados de teste são aplicados
  - ❖ Esta abordagem é indicada quando está disponível uma grande quantidade de dados.
  - ❖ Caso o conjunto total de dados seja pequeno, o erro calculado na predição pode sofrer muita variação



# Validação cruzada





---

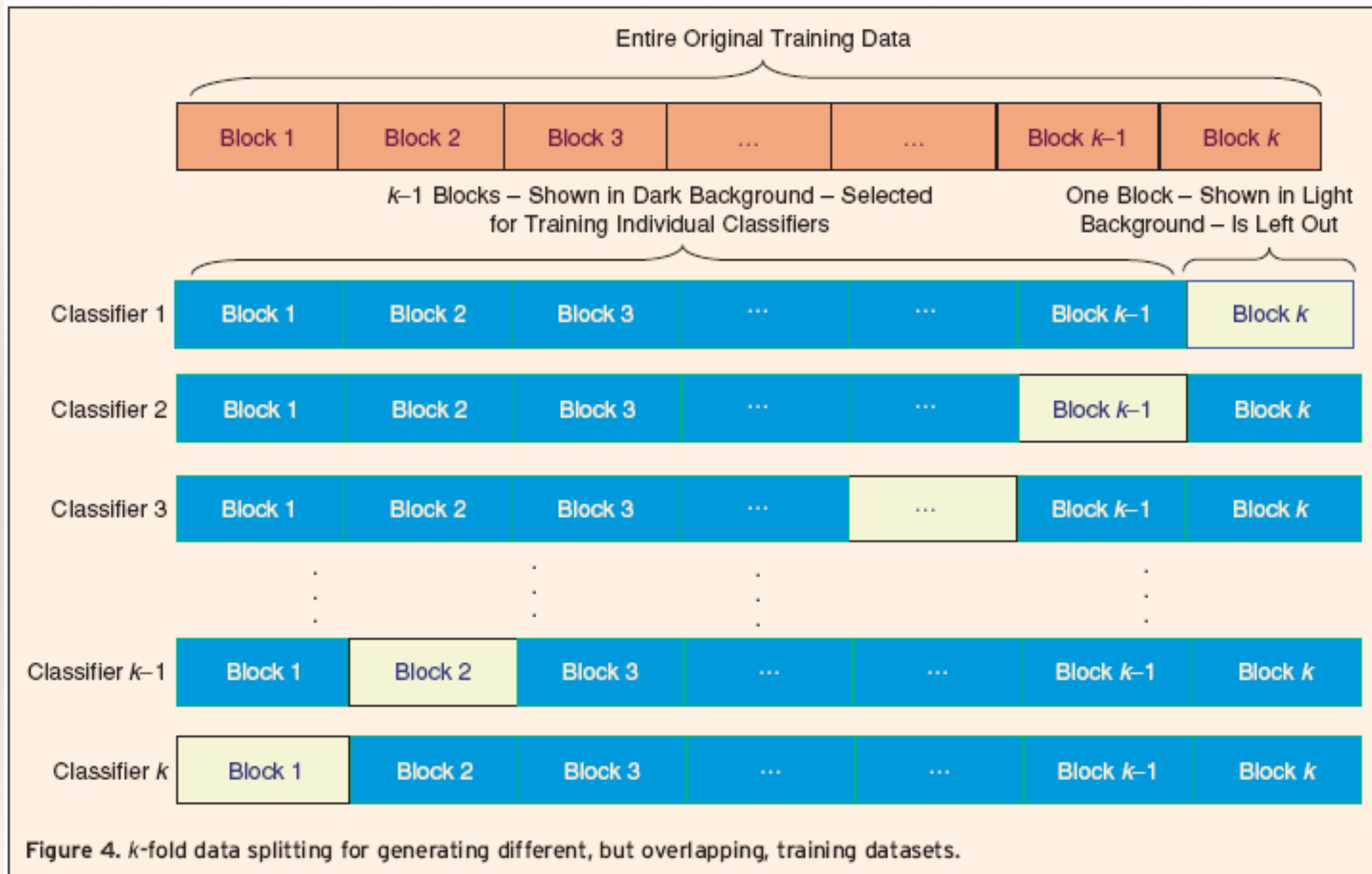
# Validação cruzada

---

- ❖ Método *k-fold*:
  - ❖ Dividir o conjunto total de dados em  $k$  subconjuntos mutuamente exclusivos do mesmo tamanho
  - ❖ Um subconjunto é utilizado para teste e os  $k-1$  restantes são utilizados para estimação dos parâmetros (treinamento). Ao final das  $k$  iterações calcula-se a acurácia sobre os erros encontrados
  - ❖ Ao final tem-se uma medida mais confiável sobre a capacidade do modelo de representar o processo gerador dos dados



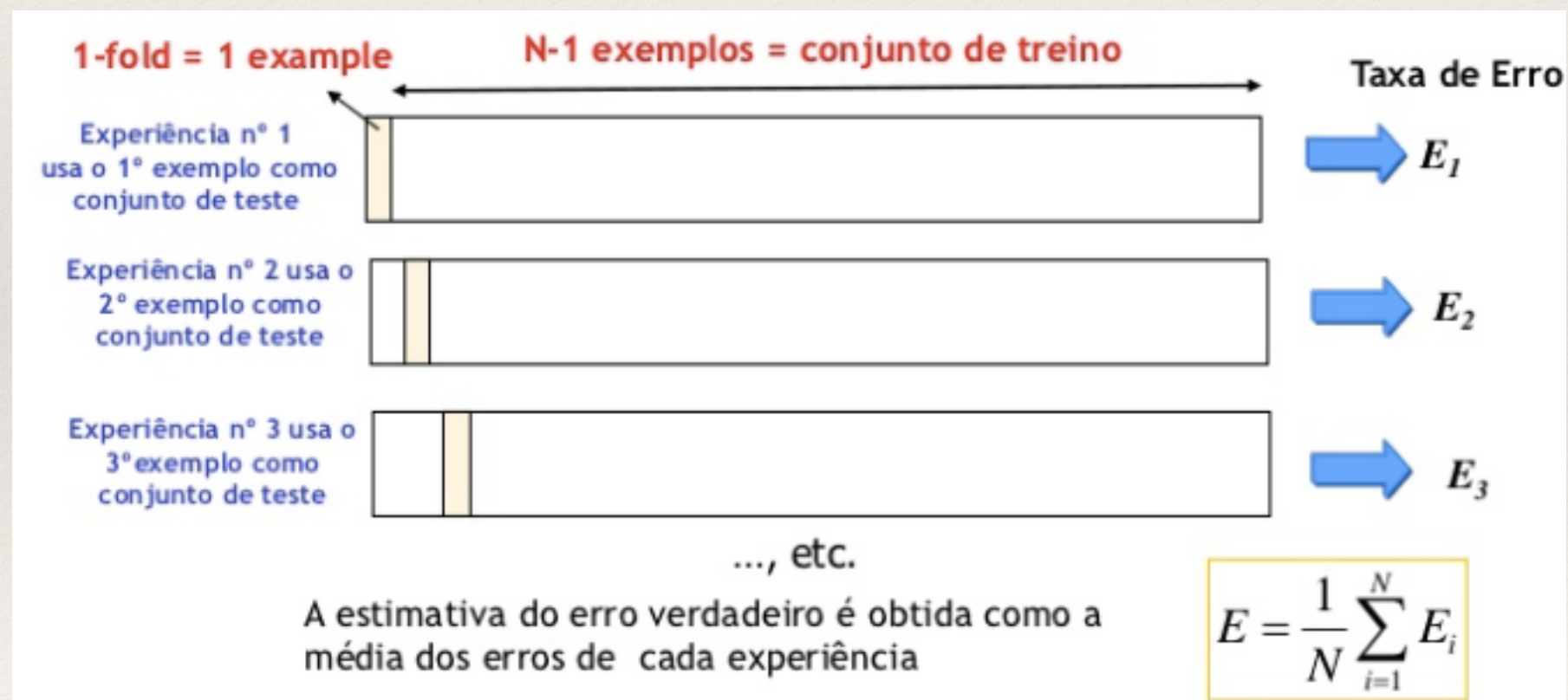
# Validação cruzada





# Validação cruzada

- ❖ Método *leave one out*:
  - ❖ Caso específico do k-fold, com k igual ao número total de dados N. Nesta abordagem são realizados N cálculos de erro, um para cada dado.
  - ❖ Alto custo computacional





---

Avaliação sensível à  
distribuição das classes e  
ao custo

---



# Tomada de decisão. Podemos errar?

Um classificador permite auxiliar à tomada de decisões entre diferentes ações. Podemos permitirmos tomar decisões erradas?

Tomada de decisão numa central nuclear: Um classificador  $h$  prediz se **abrir** ou **fechar** a válvula do módulo de refrigeração num dado momento

- Avaliamos desempenho num conjunto de teste = **100 000 dados** acumulados no último mês; a classe é o resultado da decisão tomada por um operário (esperto) em cada momento
  - Número de exemplos da classe “**fechar**”: **99 500**
  - Número de exemplos da classe “**abrir**”: **500**
- Suponhamos  $h$  prediz sempre “**fechar**” (classe maioritária). A taxa de erro é muito pequena:

$$\text{Err} = \frac{500}{100000} \times 100 = 0.5\%$$

É  $h$  um bom clasificador?



# Problema de Decisão Central Nuclear

## Matriz de Confusão

CLASSE ATUAL	CLASSE PREDITA	
	abrir	fechar
	abrir	fechar
	TP	FN
	FP	TN

Taxa de acerto (accuracy):

$$\frac{TP + TN}{TP + TN + FP + FN}$$

Diagonal dos acertos

✓ TP (true positive) - positivos verdadeiros

nº de exemplos classificados “abrir” que são “abrir”

(correctamente classificados)

✓ FP (false positive) - positivos falsos

nº de exemplos classificados “abrir” que são “fechar”

(incorrectamente classificados)

✓ TN (true negative) - negativos verdadeiros

nº de exemplos classificados “fechar” que são “fechar”

(correctamente classificados)

✓ FN (false negative) - negativos falsos

nº de exemplos classificados “fechar” que são “abrir”

(incorrectamente classificados)



# Matriz de Confusão

## Problema de Classificação Binária

CLASSE ACTUAL	CLASSE PREDITA		
		Yes (+)	No (-)
	Yes (+)	TP	FN
	No (-)	FP	TN

Taxa de acerto (accuracy):

$$\frac{TP + TN}{TP + TN + FP + FN}$$

✓ TP (true positive) - positivos verdadeiros

nº de exemplos classificados positivos que são positivos (correctamente classificados)

✓ FP (false positive) - positivos falsos

nº de exemplos classificados positivos que são negativos (incorrectamente classificados)

✓ TN (true negative) - negativos verdadeiros

nº de exemplos classificados negativos que são negativos (correctamente classificados)

✓ FN (false negative) - negativos falsos

nº de exemplos classificados negativos que são positivos (incorrectamente classificados)



# Medidas de Avaliação

- ✓ **True Positive Rate** = recall (sensibility):  
proporção de positivos verdadeiros do total de positivos

$$\text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

- ✓ **False Positive Rate**: proporção positivos falsos (incorrectamente classificados como positivos) do total de negativos

$$\text{FPR} = \frac{\text{FP}}{\text{FP} + \text{TN}}$$

- ✓ **True Negative Rate**:  
proporção de negativos verdadeiros do total de negativos

$$\text{TNR} = \frac{\text{TN}}{\text{FP} + \text{TN}}$$

- ✓ **False Negative Rate**: proporção de negativos falsos (incorrectamente classificados como negativos) do total de positivos

$$\text{FNR} = \frac{\text{FN}}{\text{TP} + \text{FN}}$$

		predita	
		+	-
actual	+	TP	FN
	-	FP	TN

		predita	
		+	-
actual	+	TP	FN
	-	FP	TN

		predita	
		+	-
actual	+	TP	FN
	-	FP	TN

		predita	
		+	-
actual	+	TP	FN
	-	FP	TN



# Precisão e Sensibilidade

- ✓ Precision (precisão): proporção de positivos verdadeiros do total dos exemplos classificados como positivos

$$\text{precision} = \frac{TP}{TP + FP}$$

		predita	
		+	-
actual	+	TP	FN
	-	FP	TN

- ✓ Recall (sensibilidade) (true positive rate): proporção de exemplos positivos que foram correctamente classificados

$$\text{recall} = \text{TPR} = \frac{TP}{TP + FN}$$

		predita	
		+	-
actual	+	TP	FN
	-	FP	TN

$$\text{F - measure} = \frac{2}{\frac{1}{\text{precision}} + \frac{1}{\text{recall}}}$$



# Precisão e Sensibilidade

Duas medidas de desempenho muito usadas nos sistemas de recuperação de informação (information retrieval systems). Os documentos de uma base de dados podem ser **recuperados** (classificados como relevantes) ou **rejeitados** a partir de uma “**query**” à base de dados realizado por um utilizador

- ✓ **Precision (precisão)**: mede a proporção dos documentos recuperados que são realmente relevantes do total de documentos recuperados.

$$\text{precisão} = \frac{\text{documentos relevantes recuperados}}{\text{documentos recuperados}}$$

- ✓ **Recall (sensibilidade)** (true positive rate): reflete a probabilidade de que um documento realmente relevante seja recuperado pelo sistema

$$\text{sensibilidade} = \frac{\text{documentos relevantes recuperados}}{\text{documentos relevantes}}$$