

Etapas Principais Para Implementar um Projeto de Ciência de Dados

1. Definição do problema

- Entender o objetivo da negócio;
- Definição das objetivos do projeto;
- Identificação das Perguntas de Pesquisa;
- Análise dos dados disponíveis;
- Definição das métricas de sucesso;
- Planejamento do Projeto;

2. Coleta e Armazenamento de Dados

- Dados podem ser coletados por: Pesquisas, Entrevistas, observação, sensores, registros de transações, APIs, Bancos de dados públicos.
- Garantir a qualidade dos dados, técnicas de limpeza e validação dos dados ajudam a manter a integridade dos dados.
- O armazenamento dos dados deve ser eficiente e seguro. Bancos de dados relacionais, NoSQL, Data lakes, Data Lakehouses e Data warehouses são opções.

Data Lakes = É um repositório centralizado que permite armazenar dados estruturados e não estruturados. Podemos realizar qualquer tipo de análise sem a necessidade de estruturar os dados previamente, permite armazenar um grande volume de dados.

→ estão relacionados.

Data warehouses = É um repositório centralizado que só armazena dados estruturados e anteriormente pré-processados, armazena menos tipos de dados, por isso possui uma performance de pesquisa mais rápida.

Data Lakehouses = É um novo conceito que une a flexibilidade, economia e escalabilidade de um Data Lake, e o gerenciamento de dados e os recursos de transações ACID de um Data warehouse.

Transações

Atomicity
Consistency
Isolation
Durability

Um sistema que aplica essas operações é chamado sistema transacional.

- Proteção dos dados coletados. Implementar medidas de segurança, como criptografia, controle de acesso e políticas de privacidade.

3. Preparação e limpeza dos dados

- A preparação de dados é o processo de transformar dados brutos em formato adequado para análise e modelagem;

- Etapas principais:

Limpeza: Envolve a remoção de valores nulos, duplicados e inconsistentes;

Normalização: Padronização dos dados para manter a consistência;

Transformação: Conversão de dados em formatos apropriados para análise. Isso pode envolver agregação, codificação e criação de novas variáveis;

Integração: Combina dados de diferentes fontes para oferecer uma visão unificada. Técnicas de **ETL** são frequentemente utilizadas para consolidar dados de diversos sistemas em um repositório central;

↳ **ETL** = Extração, Transformação e carga.

Codificação: Processo de transformar variáveis de tipo texto em representação numérica. Modificamos os dados sem modificar a informação.

Ex. Variável implicando se um cliente vai ou não fazer uma compra com valores sim/não. Codificamos para 1/0 a fim de treinar um modelo de **Machine Learning**.

Redução: Simplifica conjuntos de dados grandes e complexos, mantendo as informações mais relevantes. Métodos como seleção de características e **PCA** ajudam a reduzir a complexidade e melhorar a eficiência da análise.

↳ **PCA**: Análise de Componentes Principais.

4. Análise Exploratória dos Dados

A exploração e visualização de dados, conhecida como análise exploratória dos dados (EDA), é uma etapa crítica no fluxo de trabalho de ciência de dados.

EDA: envolve análise inicial dos dados para entender suas principais características, padrões e anomalias.

- Esse processo utiliza técnicas estatísticas descritivas e ferramentas de visualização para resumir as distribuições dos dados, identificar relações entre variáveis e detectar valores discrepantes (outliers);

- O objetivo é explorar os dados para compreender seus padrões e detectar eventuais problemas. Também podemos aplicar engenharia de atributos durante ou logo após a EDA.

5. Modelagem preditiva/Estatística

Modelagem Preditiva = Modelagem estatística

- Interesse em utilizar as variáveis para fazer previsões;

Abordagem:

- Inclui uma ampla gama de técnicas, como aprendizado de máquina supervisionado, não supervisionado, etc. Assume que os dados seguem certas distribuições e que há relações lineares entre variáveis. **Métodos: Paramétricos (suposições fortes), Não Paramétricos (suposições fracas).**

Interpretação:

- Resultados mais difíceis de interpretar, especialmente modelos vindos de redes neurais. Foca na precisão das previsões;

Exemplos de uso:

- Usada em negócios para prever vendas futuras, comportamentos de clientes e detecção de fraude. Aplicações em engenharia para prever falha de máquinas.

Ferramentas e técnicas:

Ferramentas: Python com

scikit-learn, pytorch, tensorflow, Linguagem R com caret, randomForest, julia, Rust, C++, Java, Java Script.

Técnicas: Árvores de decisão, Florestas aleatórias, Redes neurais, Deep learning, SVM, Boosting, Bagging, Regressão linear, Regressão logística;

6. Avaliação e Teste

É responsável por avaliar e testar as soluções criadas durante o projeto de data science.

- Métricas de avaliação;
- Validação cruzada;
- Overfitting e Underfitting;
- Curvas de Aprendizado;
- Análise de erros (Resíduos);
- Benchmarking;

7. Entrega de Resultados

A entrega de um projeto de ciência de dados depende dos objetivos, público-alvo e contexto do negócio:

- Relatório técnico ou científico;
- Relatório executivo;
- Dashboard interativo ou infográfico;
- Jupyter Notebook;
- Código fonte e documentação;
- API;
- Aplicação web para deploy do modelo ML;
- Previsões do modelo ML em arquivo CSV;
- Previsões do modelo ML em um banco de dados;
- Podemos também se entregar o arquivo do modelo ML.

Interpretação:

- Resultados interpretáveis e que podem ser usados para inferir relações entre as variáveis. Coeficientes com significados claros, intervalos de confiança, valores p (Resultados amplamente

usados em modelagem estatística);

Exemplos de uso:

- Usada para testar hipóteses, como análise de pesquisa de mercado;

Ferramentas e técnicas:

Ferramentas: Linguagem R, Stata, SAS, SPSS, Python com statsmodels.

Técnicas: Regressão linear, Regressão logística, ANOVA, Análise de sobrevivência, Análise Fatorial, métodos probabilísticos, etc;