

UNIP

UNIVERSIDADE PAULISTA

Ciência de Dados

Autora: Profa. Vanessa Santos Lessa

Colaboradora: Profa. Christiane Mazur Doi

Professora conteudista: Vanessa Santos Lessa

É doutora em Ciências e Aplicações Geoespaciais pelo Instituto Presbiteriano Mackenzie, mestre em Engenharia Elétrica com ênfase em Inteligência Artificial Aplicada à Automação pelo Centro Universitário da Fundação Educacional Inaciana (FEI) e bacharel em Engenharia da Computação pela Universidade São Judas Tadeu (USJT). Atua como coordenadora do curso de Bacharelado em Ciência da Computação na Universidade Paulista (UNIP) na modalidade educação a distância (EaD). Atua há mais de 18 anos em cargos técnicos e gerenciais na área de computação.

Dados Internacionais de Catalogação na Publicação (CIP)

L638c Lessa, Vanessa Santos.

 Ciência de Dados / Vanessa Santos Lessa. – São Paulo: Editora Sol, 2023.

 160 p., il.

 Nota: este volume está publicado nos Cadernos de Estudos e Pesquisas da UNIP, Série Didática, ISSN 1517-9230.

 1. Dados. 2. Máquina. 3. Modelos. I. Título.

CDU 004.65

U518.49 – 23

Profa. Sandra Miessa
Reitora

Profa. Dra. Marília Ancona Lopez
Vice-Reitora de Graduação

Profa. Dra. Marina Ancona Lopez Soligo
Vice-Reitora de Pós-Graduação e Pesquisa

Profa. Dra. Claudia Meucci Andreatini
Vice-Reitora de Administração e Finanças

Prof. Dr. Paschoal Laercio Armonia
Vice-Reitor de Extensão

Prof. Fábio Romeu de Carvalho
Vice-Reitor de Planejamento

Profa. Melânia Dalla Torre
Vice-Reitora das Unidades Universitárias

Profa. Silvia Gomes Miessa
Vice-Reitora de Recursos Humanos e de Pessoal

Profa. Laura Ancona Lee
Vice-Reitora de Relações Internacionais

Prof. Marcus Vinícius Mathias
Vice-Reitor de Assuntos da Comunidade Universitária

UNIP EaD

Profa. Elisabete Brihy
Profa. M. Isabel Cristina Satie Yoshida Tonetto
Prof. M. Ivan Daliberto Frugoli
Prof. Dr. Luiz Felipe Scabar

Material Didático

Comissão editorial:

Profa. Dra. Christiane Mazur Doi
Profa. Dra. Ronilda Ribeiro

Apoio:

Profa. Cláudia Regina Baptista
Profa. M. Deise Alcantara Carreiro
Profa. Ana Paula Tôrres de Novaes Menezes

Projeto gráfico:

Prof. Alexandre Ponzetto

Revisão:

Lucas Ricardi
Kleber Souza

Sumário

Ciência de Dados

APRESENTAÇÃO	7
INTRODUÇÃO	8

Unidade I

1 VISÃO GERAL SOBRE CIÊNCIA DE DADOS.....	9
1.1 Big Data e ciência de dados: além do hype.....	12
1.2 Atual cenário data-driven	17
1.3 Problemas e desafios	20
1.4 Soluções baseadas em dados	21
1.5 Conjuntos de habilidades do profissional cientista de dados	24
1.6 Perspectivas de um projeto de ciência de dados.....	25
1.7 Visão geral sobre KDD (do inglês, knowledge discovery in databases)	26
1.8 Representação e extração de conhecimento	31
1.9 Fontes de dados	33
2 VISÃO GERAL SOBRE APRENDIZADO DE MÁQUINA.....	36
2.1 O que é machine learning (ML)?	37
2.2 Pipeline da aprendizagem do modelo	39
2.3 Overfitting e underfitting	40
2.4 Balanço entre viés e variância em modelos de ML	43
2.5 Viés indutivo.....	44
2.6 Sistema de aprendizado	45
2.7 Tipos de aprendizagem	46
2.8 Espaço de hipóteses	47
2.9 Viés de busca: ajuste aos dados	49
3 APRENDIZADO DESCRITIVO E PREDITIVO.....	50
3.1 Aprendizado supervisionado	51
3.2 Classificação	52
3.3 Regressão	57
3.4 Aprendizado não supervisionado	59
3.5 Agrupamento	60
3.6 Associação.....	62
3.7 Sumarização	64
4 MINERAÇÃO DE DADOS	65
4.1 Visão geral	66
4.2 Modelos de ML.....	68

4.3 Árvore de decisão	70
4.4 Naive Bayes	73
4.5 K-vizinhos mais próximos (KNN)	75
4.6 K-médias	77

Unidade II

5 PREPARAÇÃO E PRÉ-PROCESSAMENTO DE DADOS	85
5.1 Principais fontes de dados.....	89
5.2 Coleta, limpeza e organização das informações	90
5.3 Métodos de raspagem.....	93
5.4 Tabulação.....	94
5.5 Seleção de atributos.....	95
5.6 Engenharia de características	96
5.7 Normalização dos dados	97
5.8 Dados ausentes	102
6 MODELOS PREDITIVOS	104
6.1 Regressão linear simples	105
6.2 Ajuste com mínimos quadrados.....	110
6.3 Gradiente descendente	113
6.4 Regressão linear múltipla	117
6.5 Regressão logística	119
7 PLANEJAMENTO DE EXPERIMENTOS	121
7.1 Split de dados – treino, teste e validação	122
7.2 Validação cruzada	124
7.3 Benchmarking.....	127
8 ANÁLISE DE RESULTADOS EXPERIMENTAIS E APLICAÇÕES AVANÇADAS DE ML	129
8.1 Métricas.....	131
8.2 Classificação	135
8.3 Regressão.....	139
8.4 Seleção de modelos.....	141
8.5 Visão computacional.....	143
8.6 Processamento de linguagem natural.....	144
8.7 Reconhecimento de fala	146
8.8 APIs de inteligência artificial.....	147

APRESENTAÇÃO

Como futuro cientista da computação, você encontrará desafios envolvendo a exploração de grande quantidade de dados (Big Data). A ciência de dados é uma área multidisciplinar que utiliza os conceitos das áreas de matemática, estatística, inteligência artificial e engenharia da computação para analisar grandes quantidades de dados. Os dados estão por toda a parte e é necessário utilizá-los de maneira correta.

O objetivo desta disciplina é fornecer os principais conceitos, técnicas e ferramentas referentes à ciência de dados e Big Data, e assim prover teoria básica para que os alunos possam aplicar as novas técnicas e ferramentas estudadas em problemas reais frente à grande quantidade de dados gerados por diferentes fontes.

Neste livro-texto, você terá a oportunidade entender os processos de descoberta de conhecimento em bases de dados, mineração e preparação dos dados, preparação dos dados e pré-processamento, modelagem dos dados, planejamento e análise dos resultados.

Vale acrescentar que este material é escrito em linguagem simples e direta, como se houvesse uma conversa entre a autora e o leitor. Adicionalmente, são inseridas muitas figuras, que auxiliam no entendimento dos tópicos desenvolvidos. Os itens chamados de observação e de lembrete são oportunidades para que você solucione eventuais dúvidas. Os itens chamados de saiba mais possibilitam que você amplie seus conhecimentos. Há, ainda, muitos exemplos, resolvidos em detalhes, o que implica a fixação dos assuntos abordados.

INTRODUÇÃO

As pessoas geram dados a partir de suas redes sociais, de seus e-mails, quando utilizam a internet, criam documentos etc. As máquinas também geram dados a partir dos seus sensores, arquivos de logs, câmeras etc. Por sua vez, as empresas o fazem a partir de suas transações comerciais, seus sistemas de controles administrativos e financeiros, seu comércio eletrônico etc. A quantidade de dados disponíveis cresce e se torna uma parte maior da nossa vida diariamente.

A ciência de dados é uma área que estuda o desenvolvimento de estratégias para analisar dados, preparar dados para análise, explorar, analisar e visualizar, e ainda construir modelos com dados usando linguagens de programação.

A empresa que utiliza a ciência de dados em seu negócio pode conseguir inúmeras vantagens, afinal de contas, os dados estão por toda a parte e é preciso saber utilizá-los de maneira correta. Podemos utilizar os dados no marketing, para melhorar as vendas, no setor de desenvolvimento de produtos, na experiência do cliente com a empresa, no setor financeiro e em qualquer outra área que utilize as informações para suas ações de forma estratégica.

O conteúdo deste livro-texto está dividido em duas unidades (unidade I e unidade II).

Na unidade I, apresentaremos inicialmente os conceitos de ciência de dados e em seguida veremos aprendizado de máquina (aprendizado descritivo e preditivo) e mineração de dados.

Na unidade II, abordaremos a preparação e o pré-processamento de dados, os modelos preditivos, o planejamento de experimentos, a análise de resultados e aplicações práticas de aprendizado de máquina.

Esperamos que você tenha uma boa leitura e se sinta motivado a ler e conhecer mais sobre a disciplina.

Bons estudos!

Unidade I

1 VISÃO GERAL SOBRE CIÊNCIA DE DADOS

A ciência de dados é um campo interdisciplinar que combina técnicas estatísticas, matemáticas e de programação para extrair insights e conhecimentos úteis a partir de conjuntos de dados complexos. Ela envolve a coleta, organização, processamento e análise de grandes volumes de dados, com o objetivo de identificar padrões, fazer previsões e tomar decisões embasadas em evidências.

A visão geral da ciência de dados abrange várias etapas. A primeira delas é a coleta de dados, que pode ser feita por meio de várias fontes, como sensores, bancos de dados, mídias sociais, entre outros. Em seguida, ocorre o processo de limpeza e organização dos dados, em que são removidos ruídos, erros e inconsistências, e os dados são estruturados de forma adequada para análise.

Após a preparação dos dados, a etapa de exploração começa. Aqui, diferentes técnicas estatísticas e de visualização são aplicadas para entender os padrões e relações presentes nos dados. Isso inclui a identificação de tendências, correlações, outliers e insights iniciais.



Observação

Outliers, em estatística e análise de dados, são valores que se diferenciam significativamente do restante dos dados em um conjunto. Esses valores extremos estão longe da média ou dos demais valores do conjunto e podem ser causados por erros de medição, comportamentos anômalos ou eventos raros. Outliers podem distorcer a análise de dados e afetar negativamente a precisão de modelos e estatísticas descritivas.

A detecção e tratamento de outliers são importantes em várias aplicações, pois podem indicar erros de coleta de dados, apontar a presença de eventos incomuns ou fornecer insights valiosos sobre comportamentos excepcionais no conjunto de dados. Existem várias técnicas estatísticas e algoritmos de detecção de outliers para identificar e lidar com esses valores atípicos, a fim de melhorar a qualidade das análises e modelagens.



Observação

Insights são percepções, entendimentos e conclusões significativas e valiosas obtidas a partir da análise de dados ou informações. São descobertas que vão além dos dados brutos, revelando padrões, tendências ou relações ocultas que podem levar a novas ideias, melhorias em processos, estratégias de negócios mais eficientes e tomadas de decisões informadas.

Os insights são uma parte essencial da análise de dados, pois ajudam a transformar dados em conhecimento útil e acionável. Podem ser alcançados por meio de diversas técnicas de análise de dados, como estatísticas descritivas, mineração de dados, aprendizado de máquina e visualização de dados. Eles são valiosos para orientar ações e estratégias de negócios, identificar oportunidades e desafios, prever tendências futuras e entender melhor o comportamento dos clientes e usuários.

A obtenção de insights eficazes requer não apenas a aplicação de técnicas analíticas avançadas, mas também um entendimento claro dos objetivos da análise, a qualidade dos dados e o contexto do problema em questão. Com a capacidade de transformar grandes volumes de dados em informações relevantes, os insights desempenham um papel crítico na tomada de decisões fundamentadas e no desenvolvimento contínuo de negócios e projetos.

Uma vez que os dados foram explorados, é possível aplicar técnicas de modelagem e aprendizado de máquina para criar modelos preditivos e descritivos. Esses modelos são alimentados com os dados disponíveis e treinados para reconhecer padrões e tomar decisões com base neles. Algoritmos de aprendizado de máquina, como regressão, árvores de decisão, redes neurais e algoritmos de agrupamento, são comumente utilizados nessa etapa.

Após a criação dos modelos, eles precisam ser avaliados e validados para garantir sua precisão e eficácia. Métricas apropriadas são definidas e os modelos são testados em dados de validação ou em novos conjuntos de dados. A validação é uma etapa crítica para garantir que os modelos sejam confiáveis e úteis.

Finalmente, os resultados obtidos pela análise dos dados e pelos modelos desenvolvidos são comunicados de forma clara e compreensível para as partes interessadas. Isso geralmente inclui a criação de relatórios, visualizações e apresentações que sintetizam as descobertas e insights.

A ciência de dados tem uma ampla gama de aplicações em diferentes setores e indústrias. Ela é usada para resolver problemas complexos, otimizar processos, prever tendências, melhorar a tomada de decisões, personalizar experiências do usuário, entre outras possibilidades. O seu papel se torna mais importante em um mundo cada vez mais orientado por dados, onde a capacidade de extrair valor de grandes volumes de informação é um diferencial competitivo.

É importante ressaltar que a ciência de dados não se limita apenas à análise de dados, mas também abrange a capacidade de formular perguntas relevantes, coletar dados de maneira adequada, criar modelos robustos e interpretar os resultados de forma crítica. Além disso, a ética e a privacidade dos dados são aspectos fundamentais que devem ser considerados ao realizar projetos de ciência de dados.

Em resumo, a ciência de dados é uma disciplina poderosa que utiliza métodos estatísticos, técnicas de programação e conhecimentos de domínio para extrair informações valiosas e tomar decisões embasadas em dados. Ela desempenha um papel fundamental na Era da Informação, permitindo que organizações e indivíduos transformem dados brutos em insights acionáveis.

Algumas das principais habilidades necessárias na ciência de dados incluem conhecimento de estatística e matemática, habilidades de programação para manipulação e análise de dados, compreensão de algoritmos de aprendizado de máquina e habilidades de visualização de dados para comunicar resultados de forma eficaz.

Além disso, a ciência de dados se beneficia de uma abordagem iterativa e orientada por problemas. Os cientistas de dados frequentemente seguem um ciclo que envolve a formulação de perguntas, a coleta e limpeza dos dados relevantes, a análise exploratória, o desenvolvimento de modelos, a validação e a comunicação dos resultados. Esse ciclo permite um processo contínuo de aprendizado e melhoria.

A ética é um aspecto crítico na ciência de dados. Ao lidar com dados sensíveis, como informações pessoais dos usuários, é essencial garantir a privacidade e a segurança dos dados. Também é importante considerar o viés nos dados e nos modelos, para evitar discriminação ou resultados distorcidos. A transparência e a responsabilidade na utilização dos dados são princípios fundamentais para a prática ética da ciência de dados.

No contexto empresarial, a ciência de dados tem sido adotada para impulsionar a inovação e o crescimento. Ela pode ser aplicada em diversas áreas, como marketing, finanças, saúde, manufatura, transporte, entre outras. Por exemplo, as empresas podem usar a ciência de dados para identificar padrões de comportamento dos clientes, otimizar processos de produção, prever demanda e melhorar a eficiência operacional.

Além disso, a ciência de dados está impulsionando o desenvolvimento de tecnologias avançadas, como a inteligência artificial e a análise de dados em tempo real. Essas tecnologias têm o potencial de transformar setores inteiros, impulsionando a automação, a personalização e a tomada de decisões baseada em dados.

A ciência de dados desempenha um papel crucial na Era Digital, permitindo que organizações e indivíduos transformem dados em insights valiosos. Com sua abordagem interdisciplinar, combinação de técnicas estatísticas e de programação, e foco na resolução de problemas, a ciência de dados continuará a impulsionar a inovação e a tomada de decisões informadas em diversos setores.

1.1 Big Data e ciência de dados: além do hype

Big Data e ciência de dados são termos frequentemente associados e muitas vezes geram um grande hype. No entanto, eles representam conceitos distintos, embora relacionados, que vão além da mera tendência ou moda. Andrea Filatro (2020, p. 14) explica que:

Data Science, ou ciência de dados, pode ser definido como a disciplina que fornece princípios, metodologias e orientações para transformação, validação, análise e criação de significado a partir de dados. O objetivo é extrair conhecimento de conjuntos de dados que, por vezes, podem ser grandes demais (o chamado Big Data) para as análises estatísticas tradicionais. Exemplos incluem a análise de estruturas genômicas complexas, a interpretação de textos e manuscritos e a otimização de estratégias de retenção de alunos.

Big Data

Big Data refere-se à enorme quantidade de dados gerados a partir de várias fontes, como transações comerciais, mídias sociais, sensores, dispositivos móveis e muito mais. Laney (2001) refere-se inicialmente a três características de Big Data utilizando 3 Vs:

- volume (grande quantidade de dados);
- velocidade (gerados em alta velocidade);
- variedade (diversidade de tipos e formatos de dados).

Zikopoulos e Eaton (2011), da mesma forma, destacaram as três características apresentando a evolução que existiu e o que promoveu um novo paradigma. Na figura a seguir podemos observar que a variedade provém da inclusão de dados não estruturados aos dados estruturados existentes, e a velocidade em que os dados passam da sua maioria de batch para dados em stream. Dessa maneira, como resultado das duas características anteriores (variedade e velocidade), passamos dos terabytes para os zettabytes em volume.



Lembrete

Big Data refere-se à enorme quantidade de dados gerados a partir de várias fontes.

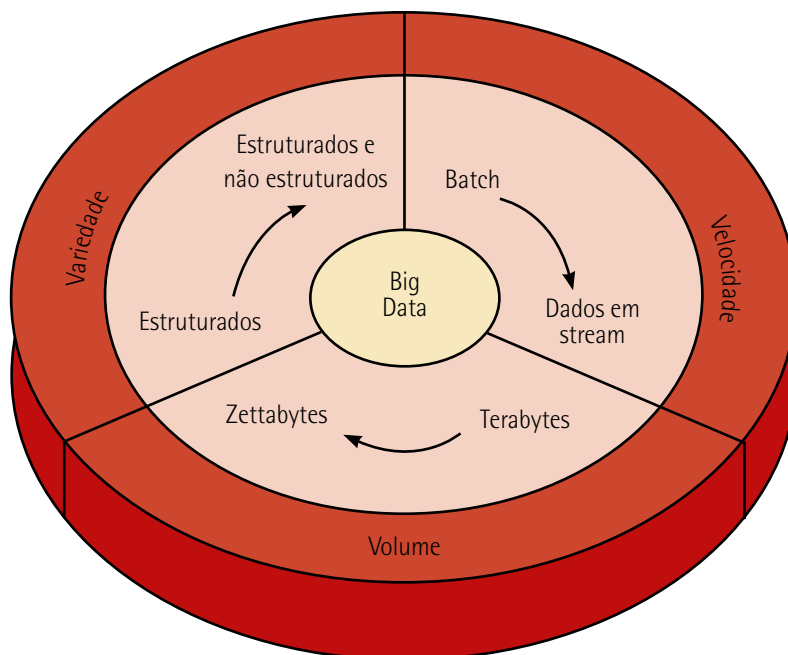


Figura 1 – Caracterização da IBM do modelo 3 Vs

Adaptada de: Zikopoulos e Eaton (2011, p. 5).

Observação

Dados batch referem-se a um conjunto de dados que são coletados, processados e armazenados em blocos ou lotes discretos. Nesse método, os dados são reunidos e processados em grupos, e o processamento ocorre em intervalos específicos, como horas, dias ou semanas. Geralmente, é uma abordagem assíncrona, em que os dados são coletados durante um período e, em seguida, processados em uma operação única ou em vários passos sequenciais.

Dados em stream (ou streaming data) referem-se a um fluxo contínuo de dados que são processados à medida que são gerados ou coletados. Nesse método, os dados são transmitidos em tempo real e são processados em pequenos fragmentos à medida que chegam, em vez de serem agrupados em lotes.

O desafio do Big Data está em como lidar com essas grandes quantidades de informações e extrair valor delas. Em 2013, Demchenko *et al.* publicaram o artigo "Addressing Big Data issues in scientific data infrastructure", no qual descrevem o modelo com as três características e acrescentam mais duas, evoluindo o modelo dos 3 Vs para 5 Vs, como observarmos na figura a seguir. As duas características adicionadas foram:

- veracidade (confiança nos dados);
- valor (dados geram valor e se tornam úteis).

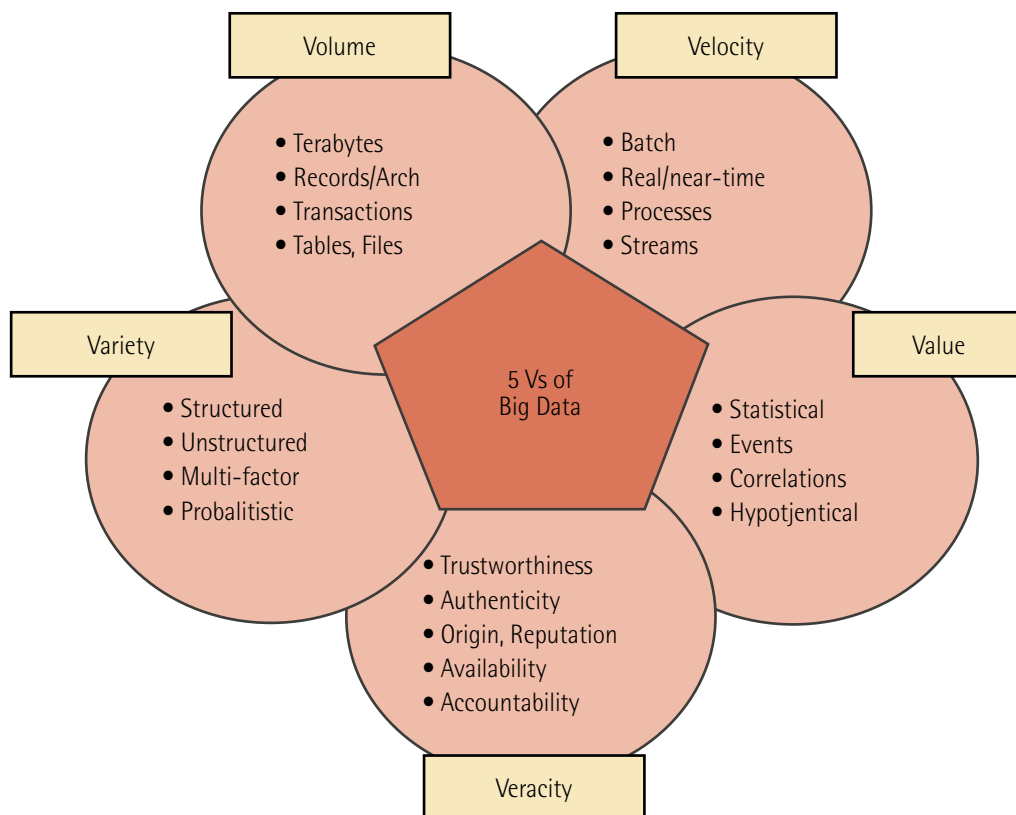


Figura 2 – Modelo dos 5 Vs

Adaptada de: Demchenko *et al.* (2013).

O estudo sobre Big Data é constante e novas características são apresentadas utilizando outras perspectivas. Um exemplo é o trabalho de Khan, Uddin e Gupta, de 2014, que acrescentou mais duas características ao modelo 5 Vs, evoluindo para o modelo 7 Vs. As duas novas características são:

- validade (refere-se à precisão e correção dos dados para o uso pretendido);
- volatilidade (retenção de dados estruturados).

Vamos nos basear na perspectiva do trabalho de Khan, Uddin e Gupta (2014) para definir melhor os 7 Vs.

- **Volume:** refere-se à grande quantidade de dados gerados diariamente. Com o avanço da tecnologia, as organizações coletam e armazenam enormes volumes de dados, muitas vezes em escala petabyte ou exabyte.

- **Velocidade:** representa a taxa na qual os dados são gerados, capturados e processados. Muitas fontes de dados, como sensores IoT (Internet das Coisas), redes sociais e transações financeiras, produzem dados em tempo real, exigindo uma análise rápida para obter insights significativos.



Observação

A Internet das Coisas (IoT) é um conceito tecnológico que se refere à interconexão de dispositivos físicos e objetos cotidianos à internet, permitindo que eles colem e compartilhem dados entre si e com sistemas de computação em nuvem, sem a necessidade de intervenção humana direta.

Esses dispositivos, chamados de "coisas" ou "dispositivos inteligentes", podem variar desde eletrodomésticos, sensores, veículos, câmeras de segurança, wearables (dispositivos vestíveis) até máquinas industriais, entre outros. Cada um desses dispositivos é equipado com sensores, atuadores e conectividade, permitindo que eles colem dados do ambiente ou de si mesmos e interajam com outros dispositivos ou sistemas.

- **Variedade:** refere-se à diversidade de formatos e tipos de dados que são gerados e coletados. Os dados podem ser estruturados (por exemplo, tabelas de bancos de dados), semiestruturados (por exemplo, JSON, XML) ou não estruturados (por exemplo, texto, áudio, vídeo).
- **Veracidade:** refere-se à qualidade e confiabilidade dos dados em um sistema de Big Data. Em outras palavras, diz respeito à certeza de que os dados coletados e utilizados na análise são precisos, completos e verdadeiros. A veracidade dos dados é essencial para garantir que as conclusões e decisões tomadas com base nesses dados sejam confiáveis e sólidas. No ambiente de Big Data, muitas fontes de dados podem conter informações incorretas, desatualizadas ou imprecisas. Por exemplo, dados coletados automaticamente por sensores podem estar sujeitos a erros ou ruídos, e dados extraídos de redes sociais podem conter informações enganosas. É fundamental realizar verificações e processos de validação para garantir a veracidade dos dados, como verificação cruzada com outras fontes confiáveis e aplicação de algoritmos de limpeza e correção de dados.
- **Valor:** refere-se à utilidade e relevância das informações extraídas a partir dos dados. É a capacidade de transformar dados brutos em insights significativos e acionáveis, capazes de gerar benefícios e vantagens para as organizações. O valor do Big Data está em sua capacidade de fornecer informações valiosas que podem ser usadas para tomada de decisões informadas, identificação de tendências, previsão de padrões, otimização de processos e criação de novos produtos ou serviços. Esses insights podem levar a melhorias na eficiência operacional, aumento da satisfação do cliente, descoberta de oportunidades de negócio, entre outros benefícios. Para extrair valor dos dados, as organizações precisam investir em ferramentas e técnicas de análise de dados avançadas, como aprendizado de máquina, mineração de dados e inteligência artificial. Além disso, a capacidade de visualização dos dados é essencial para facilitar a compreensão dos insights obtidos e facilitar a tomada de decisões pelos responsáveis.

- **Validade:** refere-se à qualidade e precisão dos dados coletados e armazenados em um sistema de Big Data. Dados válidos são aqueles que são confiáveis, precisos, atualizados e relevantes para os objetivos da análise. A falta de validade nos dados pode levar a análises incorretas e decisões inadequadas. Para garantir a validade dos dados, as organizações devem implementar mecanismos adequados de controle de qualidade durante a coleta, processamento e armazenamento dos dados. Por exemplo, na análise de dados de clientes de um e-commerce, a validade dos dados é essencial para garantir que os insights extraídos reflitam as preferências e o comportamento atual dos clientes, e não informações desatualizadas ou imprecisas.
- **Volatilidade:** refere-se à taxa de mudança dos dados em um sistema de Big Data ao longo do tempo. Em muitos cenários, os dados estão em constante mudança e atualização, e a volatilidade mede a rapidez com que novos dados são gerados e os dados existentes são modificados ou excluídos. A volatilidade dos dados é especialmente relevante para aplicações que requerem análises em tempo real, nas quais as informações precisam ser processadas rapidamente para obter insights atualizados. Por exemplo, em sistemas de monitoramento de tráfego urbano ou redes sociais, a volatilidade é crucial para garantir que os dados utilizados para tomada de decisões sejam os mais recentes disponíveis. Gerenciar a volatilidade é um desafio para as plataformas de Big Data, pois requer infraestruturas e algoritmos eficientes que possam lidar com grandes volumes de dados em constante mudança.

Essas sete características Vs do Big Data destacam os diversos desafios e oportunidades que as organizações enfrentam ao lidar com grandes volumes de dados e demonstram a importância de ter uma estratégia sólida para gerenciar, analisar e interpretar essas informações para obter vantagens competitivas.

Ciência de dados

A ciência de dados é o campo de estudo que se concentra em obter insights, conhecimentos e tomar decisões informadas a partir da análise de dados. Envolve técnicas estatísticas, matemáticas, de programação e conhecimentos de domínio para explorar, limpar, analisar e interpretar os dados, visando a descoberta de padrões, previsões, otimizações e soluções para problemas complexos.

O relacionamento entre Big Data e ciência de dados reside no fato de que a análise de Big Data requer a aplicação de métodos e técnicas da ciência de dados para lidar com as características desse tipo de dado. O Big Data, com seu volume e variedade, apresenta desafios únicos em termos de armazenamento, processamento e análise. A ciência de dados fornece as ferramentas e abordagens necessárias para lidar com esses desafios e aproveitar o potencial dos dados.

No entanto, é importante ressaltar que nem toda análise de dados requer a aplicação do conceito de Big Data. A ciência de dados pode ser aplicada em conjuntos de dados menores, mas igualmente complexos, em que as técnicas e métodos da área são utilizados para obter insights e tomar decisões embasadas.

Além disso, é fundamental compreender que a adoção efetiva da ciência de dados vai além do hype e exige uma abordagem estruturada e orientada a problemas. É necessário ter uma compreensão sólida dos

fundamentos estatísticos, técnicas de programação, matemática e conhecimentos de domínio específico para obter resultados confiáveis e significativos.

A integração bem-sucedida do Big Data e da ciência de dados requer uma combinação de habilidades técnicas, como habilidades de programação, conhecimentos avançados em estatística e algoritmos de aprendizado de máquina, bem como a capacidade de entender o contexto dos dados e formular perguntas relevantes.

Resumindo, Big Data e ciência de dados são conceitos interligados, mas distintos. O Big Data refere-se ao grande volume, velocidade e variedade de dados, enquanto a ciência de dados é o campo que busca extrair insights e conhecimentos a partir desses dados. A adoção efetiva da ciência de dados requer uma abordagem fundamentada e orientada a problemas, além de habilidades técnicas e conhecimentos de domínio específicos.

1.2 Atual cenário data-driven

O atual cenário data-driven, ou orientado por dados, tem sido marcado por um aumento significativo na adoção e valorização da utilização de dados para orientar a tomada de decisões em diversas áreas e setores. As organizações estão reconhecendo cada vez mais o potencial dos dados como um recurso estratégico para impulsionar a inovação, otimizar processos, melhorar a experiência do cliente e alcançar vantagens competitivas.

A figura a seguir apresenta o Data Science no contexto de diversos outros processos relacionados ao tratamento de dados dentro de uma organização. Ela destaca a distinção entre Data Science e outros aspectos do processamento de dados que estão ganhando importância nos negócios atualmente. A tomada de decisão orientada por dados (DOD) é mencionada como uma prática na qual as decisões são embasadas na análise de dados, em vez de dependerem apenas da intuição. Por exemplo, um negociante pode optar por selecionar anúncios com base exclusivamente em sua vasta experiência e intuição no campo ou pode fazer sua escolha com base na análise de dados sobre como os consumidores têm reagido a diferentes anúncios. Alternativamente, ele pode combinar ambas as abordagens. A DOD não é uma prática de "tudo ou nada", e muitas empresas a adotam em diferentes graus, dependendo das suas necessidades e recursos. A figura fornece uma visão geral do papel essencial do Data Science como um processo de análise de dados mais abrangente, envolvendo técnicas avançadas de modelagem e análise para obter insights valiosos a partir dos dados. Isso destaca a crescente importância de usar dados de maneira eficaz para melhorar a tomada de decisões e obter vantagens competitivas no ambiente de negócios em constante mudança.

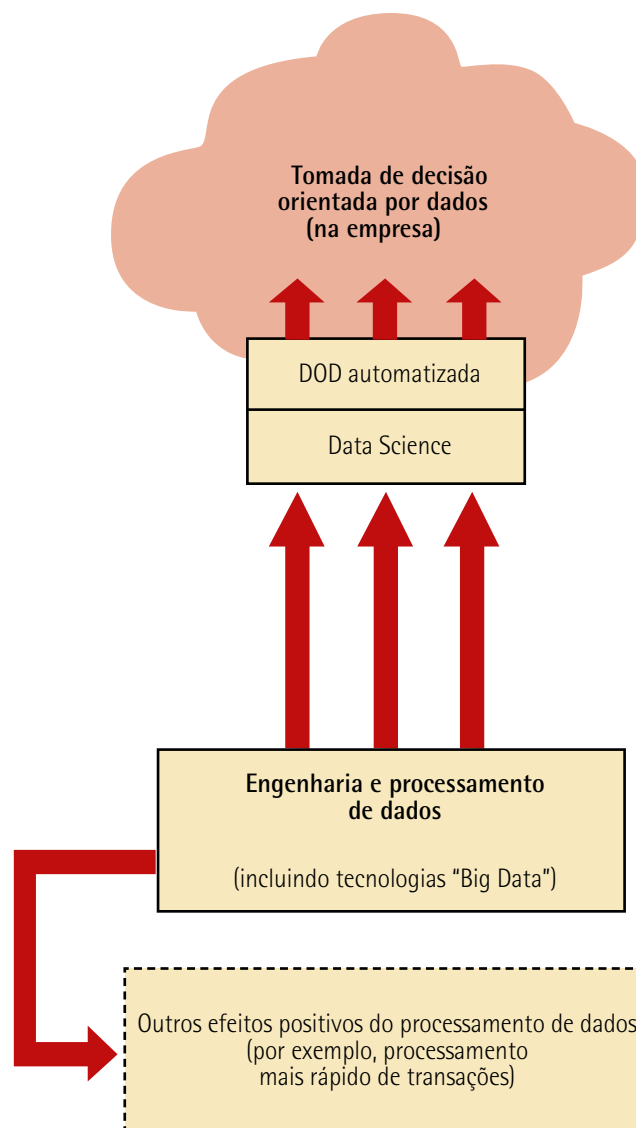


Figura 3 – Data Science no contexto dos diversos processos relacionados a dados na organização

Adaptada de: Provost e Fawcett (2013, p. 40).

Um dos principais impulsionadores desse cenário é o avanço tecnológico, que tornou possível coletar, armazenar e processar grandes volumes de dados de maneira mais acessível e eficiente. O surgimento de tecnologias como computação em nuvem, armazenamento de dados escalável e ferramentas de processamento de dados em tempo real facilitou a manipulação e análise de conjuntos de dados cada vez maiores.

Além disso, o surgimento de dispositivos conectados à internet (Internet das Coisas) e a proliferação de mídias sociais contribuiu para o aumento exponencial da quantidade de dados disponíveis. Essa explosão de dados oferece às organizações uma oportunidade sem precedentes de extrair insights valiosos e tomar decisões mais informadas.

Outro fator importante é a crescente conscientização sobre o valor dos dados. As empresas perceberam que possuem uma riqueza de informações em seus sistemas e processos internos, bem como acesso a dados externos relevantes. Ao aproveitar esses dados de forma inteligente, elas podem identificar padrões, prever tendências, entender o comportamento do cliente, melhorar a eficiência operacional e tomar decisões embasadas em evidências.

Nesse contexto, a cultura data-driven tem ganhado destaque. Ela envolve uma mudança de mentalidade nas organizações, em que a tomada de decisões é baseada em dados e evidências, em vez de intuição ou suposições. Isso requer uma infraestrutura adequada para coletar, armazenar, processar e analisar os dados, bem como a capacitação dos profissionais para lidar com as ferramentas e técnicas de análise de dados.

A análise de dados também se tornou mais acessível com o surgimento de ferramentas de análise de autoatendimento e visualização de dados. Essas ferramentas permitem que os usuários de negócios, mesmo sem habilidades avançadas em programação ou estatística, realizem análises básicas e explorem os dados por conta própria. Isso contribui para uma maior democratização da análise de dados e capacita as equipes a tomar decisões com base em insights diretos dos dados.

No entanto, apesar do crescente interesse e adoção de abordagens data-driven, ainda existem desafios a serem enfrentados. Um dos principais é a qualidade dos dados, já que nem sempre os dados disponíveis são completos, precisos ou confiáveis. A garantia de qualidade dos dados é fundamental para obter resultados precisos e confiáveis.

Outros desafios incluem a privacidade e a segurança dos dados. Com o aumento da coleta e do uso de dados pessoais, as organizações precisam estar em conformidade com regulamentações de privacidade, como a Lei Geral de Proteção de Dados Pessoais (LGPD) que estabelece diretrizes importantes e obrigatórias para a coleta, processamento e armazenamento de dados pessoais. A LGPD foi inspirada na GDPR (Regulamento Geral de Proteção de Dados) que entrou em vigência em 2018 na União Europeia, que garante que os dados sejam protegidos contra acesso não autorizado e ocasionou grandes impactos para empresas e consumidores. No Brasil, a LGPD (Lei n. 13.709, de 14 de agosto de 2018) entrou em vigor em 18 de setembro de 2020, representando um passo importante.

Além disso, a capacitação dos profissionais de dados é um fator crítico. A demanda por cientistas de dados, engenheiros de dados e analistas especializados continua a crescer, mas ainda há escassez de talentos nessa área. É necessário investir em programas de capacitação e desenvolvimento de habilidades para garantir que as equipes tenham o conhecimento e as competências necessárias para lidar com os desafios e extrair o máximo valor dos dados.

Em termos de setores específicos, o cenário data-driven tem impacto em diversas áreas. Na indústria de varejo, por exemplo, os dados são usados para entender o comportamento do consumidor, personalizar ofertas, otimizar estoques e melhorar a experiência de compra. Na área da saúde, os dados são utilizados para análise de históricos médicos, pesquisa clínica, descoberta de medicamentos e monitoramento de pacientes. E no setor financeiro, os dados são empregados na detecção de fraudes, modelagem de riscos e previsões de mercado.

Em resumo, o atual cenário data-driven reflete uma mudança significativa na forma como as organizações abordam a tomada de decisões. O aumento do volume e da variedade de dados, combinado com o avanço tecnológico e a conscientização sobre o valor dos dados, impulsionaram a adoção de abordagens orientadas por dados. No entanto, é fundamental enfrentar desafios como qualidade dos dados, privacidade, segurança e capacitação dos profissionais para aproveitar ao máximo as oportunidades oferecidas pelo cenário data-driven.

1.3 Problemas e desafios

Embora a ciência de dados tenha um grande potencial para impulsionar a inovação e fornecer insights valiosos, também enfrenta uma série de problemas e desafios. Aqui estão alguns dos principais:

- **Qualidade dos dados:** a qualidade dos dados é um desafio fundamental na ciência de dados. Os dados podem conter erros, estar incompletos, ser inconsistentes ou conter viés. A falta de qualidade dos dados pode levar a conclusões errôneas e afetar a confiabilidade dos resultados obtidos.
- **Privacidade e ética:** a coleta e o uso de dados envolvem questões de privacidade e ética. Ao lidar com dados sensíveis, como informações pessoais dos usuários, é essencial garantir a privacidade e a segurança dos dados. Além disso, é necessário considerar o viés nos dados e nos modelos para evitar discriminação ou resultados distorcidos.
- **Escalabilidade:** a ciência de dados lida com conjuntos de dados cada vez maiores, conhecidos como Big Data. A capacidade de lidar com o volume, a velocidade e a variedade desses dados é um desafio em termos de infraestrutura, armazenamento, processamento e análise. A escalabilidade dos sistemas e algoritmos é essencial para lidar eficientemente com grandes volumes de dados.
- **Complexidade dos algoritmos e modelos:** a escolha e implementação de algoritmos adequados para análise e modelagem de dados pode ser um desafio. Existem muitos algoritmos e modelos disponíveis, cada um com suas vantagens e desvantagens, e é necessário entender as características dos dados e os requisitos do problema para selecionar a abordagem mais apropriada.
- **Interpretação e comunicação dos resultados:** a ciência de dados não se trata apenas de analisar os dados, mas de comunicar os resultados de forma clara e compreensível para os usuários finais. A interpretação correta dos resultados e a capacidade de explicar as descobertas de maneira não técnica são habilidades importantes para garantir que os insights sejam compreendidos e utilizados corretamente.
- **Escassez de talentos:** existe uma demanda crescente por profissionais qualificados em ciência de dados, mas há uma escassez de talentos nessa área. Encontrar e contratar cientistas de dados, engenheiros de dados e analistas com habilidades técnicas e conhecimentos de negócios é um desafio enfrentado por muitas organizações.
- **Mudanças rápidas de tecnologia:** a ciência de dados é um campo em constante evolução, com avanços tecnológicos e novas técnicas surgindo regularmente. Manter-se atualizado com

as últimas tecnologias, ferramentas e técnicas requer um esforço contínuo de aprendizado e desenvolvimento profissional.

A ciência de dados enfrenta desafios relacionados à qualidade dos dados, privacidade, escalabilidade, complexidade dos algoritmos, interpretação dos resultados, escassez de talentos e rápidas mudanças de tecnologia. Superar esses desafios requer abordagens rigorosas, conscientização ética, investimento em infraestrutura e capacitação dos profissionais envolvidos.

1.4 Soluções baseadas em dados

Antes de listar algumas técnicas de soluções baseada em dados, vamos apresentar alguns exemplos para contextualizar a utilização dos dados na solução de problemas:

Exemplo 1: Em 2004, antes da chegada do furacão Frances, executivos do Walmart perceberam uma oportunidade de usar sua tecnologia preditiva orientada por dados para prever o impacto da tempestade. Com base nos dados históricos de compras armazenados no banco de dados da empresa, eles tentaram prever os possíveis efeitos do furacão, aprendendo com o que aconteceu em situações similares, como o furacão Charley, que ocorreu semanas antes. A diretora executiva de informação, Linda M. Dillman, liderou a equipe para desenvolver essas previsões e antecipar os eventos em vez de esperar que eles acontecessem. O objetivo era a investigação no aumento da demanda de produtos não usuais nas regiões que seriam afetadas pela aproximação do furacão. Além do aumento esperado do consumo de água mineral, os especialistas do Walmart descobriram um crescimento de sete vezes na demanda de strawberry Pop-Tarts (biscoito americano). Identificaram, também, que na fase pré-furacão, a cerveja foi a campeã de vendas!

Exemplo 2: Em 2012, a Target, concorrente do Walmart, utilizou a tomada de decisão orientada por dados para identificar clientes que estavam esperando um bebê e, assim, antecipar suas necessidades. Com base em dados históricos e técnicas de Data Science, a Target desenvolveu modelos preditivos que identificaram indicadores, como mudanças na dieta e no comportamento de compra, associados a mulheres grávidas. Essas informações permitiram que a Target direcionasse campanhas de marketing personalizadas para esses clientes antes mesmo de outros varejistas, buscando obter uma vantagem competitiva.

Exemplo 3: A empresa de telecomunicações MegaTelCo enfrenta um grande problema de retenção de clientes em seu negócio de produtos e serviços sem fio nos Estados Unidos. Cerca de 20% dos clientes abandonam o serviço quando seus contratos vencem, e a aquisição de novos clientes está se tornando cada vez mais difícil devido à saturação do mercado de telefonia celular. A empresa está buscando encontrar uma solução para reduzir a rotatividade de clientes e evitar o alto custo de atração de novos clientes. A equipe de ciência de dados foi chamada para utilizar os vastos recursos de dados da empresa e desenvolver um plano para decidir quais clientes devem receber uma oferta especial de retenção antes do término de seus contratos. A retenção de clientes tem sido uma das principais aplicações de tecnologias de mineração de dados, especialmente nos setores de telecomunicação e finanças. A principal tarefa da ciência de dados é o desenvolvimento passo a passo, com base no grande volume de dados disponíveis da MegaTelCo, para indicação dos clientes que deve ser oferecido este plano especial de retenção. Como

a MegaTelCo deve definir o público (ou clientes) alvo que irá melhor reter o cliente para um orçamento de incentivo particular? A resposta desta questão é muito mais complicada que parece.

Existem várias soluções baseadas em dados que podem ser implementadas para enfrentar desafios e obter resultados significativos. Aqui estão algumas delas:

- **Limpeza e pré-processamento de dados:** antes de realizar análises ou construir modelos, é fundamental fazer a limpeza e o pré-processamento dos dados. Isso envolve identificar e lidar com dados ausentes, valores discrepantes ou inconsistentes e realizar transformações ou normalizações necessárias. A utilização de técnicas de limpeza e pré-processamento apropriadas garante que os dados estejam prontos para análises mais avançadas.
- **Análise exploratória de dados:** a análise exploratória de dados é uma etapa crítica para compreender as características dos dados e identificar padrões, tendências e insights iniciais. Ela envolve a aplicação de técnicas estatísticas e visualização de dados para explorar a distribuição, a correlação e a estrutura dos dados. A análise exploratória ajuda a definir perguntas e hipóteses mais específicas para análises posteriores.
- **Modelagem preditiva:** a modelagem preditiva envolve a construção de modelos estatísticos ou de aprendizado de máquina para fazer previsões ou classificações com base em dados históricos. Esses modelos podem ajudar a identificar padrões ocultos, fazer previsões futuras ou tomar decisões com base em evidências. A seleção e construção adequada de modelos é crucial para obter resultados precisos e confiáveis.
- **Segmentação e personalização:** com base nos dados do cliente, é possível segmentar o público em grupos com características semelhantes. Essa segmentação permite personalizar estratégias de marketing, ofertas de produtos ou serviços e experiências do cliente. A utilização de algoritmos de agrupamento e técnicas de análise de segmentação ajuda a identificar grupos distintos e entender as necessidades e preferências dos clientes.
- **Deteção de anomalias e fraudes:** a análise de dados pode ser usada para identificar padrões anormais ou suspeitos que podem indicar atividades fraudulentas ou comportamentos inesperados. A detecção de anomalias é útil em várias áreas, como segurança cibernética, detecção de fraudes financeiras e manutenção preditiva de equipamentos. Algoritmos de detecção de anomalias podem ser aplicados para identificar desvios em relação ao comportamento esperado.
- **Otimização de processos:** a análise de dados pode ser aplicada para otimizar processos, identificando gargalos, pontos de melhoria e oportunidades de eficiência. Através da análise de dados de processos, é possível identificar os principais fatores que afetam o desempenho e encontrar soluções para otimizar recursos, reduzir custos e melhorar a qualidade.
- **Visualização de dados:** a visualização de dados é uma solução eficaz para comunicar informações complexas de forma clara e compreensível. Ela permite explorar e apresentar dados de maneira

visual, facilitando a identificação de padrões e insights. A visualização de dados ajuda a contar histórias e transmitir informações importantes para diferentes partes interessadas.

Essas são apenas algumas soluções baseadas em dados que podem ser aplicadas em diversas áreas e setores. É importante ressaltar que a escolha da solução mais adequada depende do contexto e dos objetivos específicos de cada organização. Além das soluções mencionadas anteriormente, outras abordagens podem ser consideradas, como:

- **Text mining e processamento de linguagem natural:** essas técnicas permitem extrair informações valiosas a partir de grandes volumes de texto não estruturado. É possível analisar sentimentos, identificar tópicos relevantes, realizar análise de sentimento em mídias sociais, entre outras aplicações.
- **Aprendizado de máquina interpretável:** o uso de algoritmos de aprendizado de máquina interpretáveis permite compreender o processo de tomada de decisão do modelo. Isso é especialmente importante em setores regulados, em que é necessário explicar as decisões tomadas por um modelo, como no caso de empréstimos, seguros ou sistemas de saúde.
- **Aprendizado de reforço:** o aprendizado de reforço envolve a construção de agentes inteligentes que aprendem a tomar decisões por meio de interações com o ambiente. Essa abordagem é útil em situações em que as regras não são conhecidas antecipadamente e o agente precisa aprender por tentativa e erro.
- **Dados em tempo real:** com o avanço da tecnologia, a capacidade de lidar com dados em tempo real se tornou essencial em muitos setores. Isso permite a detecção de eventos em tempo real, análises preditivas imediatas e tomada de decisões em tempo hábil.
- **Automação de processos:** a automação de processos com base em dados é uma tendência crescente. Por meio da combinação de dados e algoritmos, é possível automatizar tarefas e processos, reduzindo a intervenção humana, melhorando a eficiência e minimizando erros.

É importante destacar que cada solução apresenta desafios específicos, como a disponibilidade e qualidade dos dados, a necessidade de expertise técnica e a interpretação correta dos resultados. Além disso, é essencial considerar questões éticas e legais relacionadas à privacidade, proteção de dados e viés algorítmico ao implementar soluções baseadas em dados.

As soluções baseadas em dados são diversas e têm o potencial de trazer benefícios significativos para as organizações. Por meio da análise, interpretação e aplicação de insights dos dados, é possível impulsionar a inovação, otimizar processos, melhorar a tomada de decisões e obter vantagens competitivas. Cada solução deve ser adaptada às necessidades específicas de cada contexto, considerando-se os desafios e requisitos envolvidos.

1.5 Conjuntos de habilidades do profissional cientista de dados

Um cientista de dados é um profissional altamente qualificado e multidisciplinar, com habilidades técnicas e conhecimentos em várias áreas. Aqui estão alguns conjuntos de habilidades essenciais para um cientista de dados:

- **Conhecimento em programação:** um cientista de dados deve ter habilidades em programação para manipular e analisar dados. Linguagens comuns incluem Python e R, além de conhecimento em SQL para consulta e manipulação de bancos de dados.
- **Estatística e matemática:** um entendimento sólido de estatística é fundamental para a análise de dados. Isso inclui conhecimento em probabilidade, testes de hipóteses, regressão, modelos probabilísticos, entre outros conceitos estatísticos. Além disso, conhecimentos matemáticos, como álgebra linear e cálculo, são importantes para entender algoritmos e técnicas de modelagem.
- **Aprendizado de máquina e mineração de dados:** um cientista de dados precisa ter conhecimento em técnicas de aprendizado de máquina, como regressão, classificação, clustering, árvores de decisão, redes neurais, entre outras. Além disso, habilidades em mineração de dados, incluindo pré-processamento de dados, seleção de recursos e validação de modelos, são essenciais.
- **Conhecimento de bancos de dados:** um cientista de dados deve ter conhecimentos em bancos de dados, tanto em bancos de dados relacionais quanto em bancos de dados não relacionais. Isso inclui habilidades em SQL para consultar, extrair e manipular dados, bem como conhecimento em tecnologias de armazenamento e processamento de Big Data, como Hadoop e Spark.
- **Visualização de dados:** a capacidade de comunicar informações complexas de forma clara e visualmente atraente é fundamental para um cientista de dados. O conhecimento em ferramentas de visualização, como Matplotlib, ggplot, Tableau ou Power BI, permite criar gráficos e visualizações que ajudam a transmitir insights aos stakeholders.
- **Domínio do negócio:** um cientista de dados deve entender o contexto e o domínio do negócio em que está trabalhando. Isso envolve conhecimento sobre as necessidades e desafios específicos do setor, bem como a capacidade de formular perguntas relevantes e aplicar a ciência de dados para solucionar problemas reais.
- **Pensamento analítico e resolução de problemas:** um cientista de dados deve ser capaz de abordar problemas complexos de forma analítica, identificando padrões, formulando hipóteses, testando modelos e tirando conclusões embasadas nos dados. A capacidade de resolver problemas de maneira lógica e criativa é essencial.
- **Comunicação e habilidades interpessoais:** ser capaz de comunicar os resultados e insights de forma clara e eficaz é uma habilidade crucial para um cientista de dados. Isso envolve a capacidade de explicar conceitos técnicos de forma não técnica, trabalhar em equipe, colaborar com outros profissionais e entender as necessidades e expectativas dos stakeholders.

Essas são apenas algumas das principais habilidades necessárias para um cientista de dados. É importante ressaltar que a ciência de dados é um campo em constante evolução, portanto, a capacidade de aprendizado contínuo e adaptação às novas tecnologias e técnicas também é fundamental para um cientista de dados. Além disso, habilidades de pensamento crítico, curiosidade, resiliência e capacidade de lidar com grandes volumes de dados são características valorizadas nessa área.

É importante destacar que um cientista de dados raramente possui todas essas habilidades mencionadas em um nível avançado desde o início de sua carreira. É um processo contínuo de aprendizado e aprimoramento, no qual profissionais podem se especializar em áreas específicas de acordo com seus interesses e necessidades do mercado.

Além disso, equipes de ciência de dados frequentemente são compostas de profissionais com habilidades complementares, como engenheiros de dados, analistas de negócios e especialistas em visualização de dados, que colaboram para alcançar os objetivos de análise e insights em uma organização.

No geral, um cientista de dados deve ter um conjunto diversificado de habilidades técnicas, conhecimento em estatística e matemática, compreensão do domínio do negócio e habilidades interpessoais. A combinação dessas habilidades permite que eles aproveitem o poder dos dados para tomar decisões informadas, identificar oportunidades de negócios e solucionar problemas complexos em diversas áreas e setores.

1.6 Perspectivas de um projeto de ciência de dados

Um projeto de ciência de dados apresenta várias perspectivas e possibilidades. A primeira delas é identificar e compreender claramente o problema ou desafio que o projeto de ciência de dados busca resolver. Isso envolve definir os objetivos, as metas e as expectativas do projeto, bem como entender o contexto e as restrições envolvidas.

A coleta e a preparação de dados são perspectivas cruciais em um projeto de ciência de dados. Isso inclui identificar as fontes de dados relevantes, coletar os dados necessários e realizar a limpeza, o pré-processamento e a transformação dos dados, garantindo que eles estejam prontos para análises posteriores.

A análise exploratória de dados é essencial para compreender a estrutura, a distribuição e as características dos dados. Isso envolve a aplicação de técnicas estatísticas e visualização de dados para identificar padrões, tendências e insights iniciais, ajudando a definir perguntas e hipóteses mais específicas.

A perspectiva de modelagem inclui a seleção e a construção de modelos estatísticos ou de aprendizado de máquina adequados para o problema em questão. Isso requer a escolha dos algoritmos apropriados, o ajuste dos parâmetros do modelo e a validação do desempenho do modelo por meio de métricas relevantes. Uma vez que os modelos são treinados e avaliados, é necessário interpretar e analisar os resultados. Essa perspectiva envolve a avaliação do desempenho do modelo, a interpretação dos principais fatores que afetam os resultados e a comunicação dos insights obtidos para as partes interessadas.

Após a construção do modelo, a perspectiva de implantação envolve a integração do modelo em sistemas ou processos existentes. É importante acompanhar e monitorar continuamente o desempenho do modelo em produção, garantindo que ele forneça resultados precisos e atualizados. Uma perspectiva fundamental em projetos de ciência de dados é o aprendizado contínuo e a busca pela melhoria. Isso envolve a análise dos resultados obtidos, a identificação de oportunidades de aprimoramento, a exploração de novas técnicas e abordagens, e a atualização do modelo à medida que novos dados e informações se tornam disponíveis.

Cada perspectiva que acabamos de mencionar é interligada e contribui para o sucesso geral do projeto de ciência de dados. É importante abordar cada uma delas de forma cuidadosa e sistemática para garantir resultados relevantes, confiáveis e aplicáveis para a organização ou problema em questão.

1.7 Visão geral sobre KDD (do inglês, *knowledge discovery in databases*)

Mineração de dados

A análise exploratória de dados é uma subdivisão da estatística ligada à mineração de dados, que por sua vez surgiu da convergência de estatística clássica, inteligência artificial e aprendizado de máquina. Além disso, a mineração de dados se relaciona com descoberta de conhecimento e aprendizado de máquina na área de inteligência artificial. O termo "mineração de dados" (data mining) refere-se aos estágios de descoberta do processo de KDD, integrando-se a esse processo. A descoberta de conhecimento em bancos de dados (KDD) está relacionada com a mineração de dados, conforme a figura a seguir:

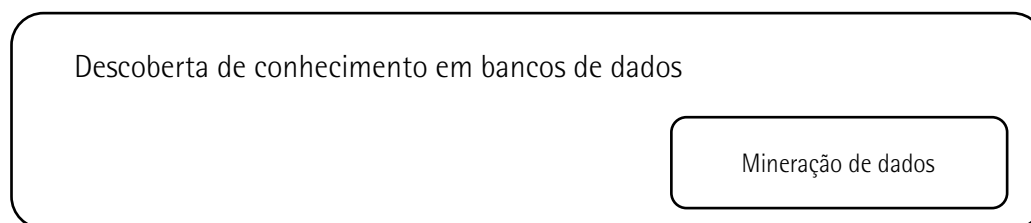


Figura 4 – Relação entre a descoberta de conhecimento em bancos de dados e a mineração de dados

O processo de descoberta de conhecimento

A descoberta de conhecimento em bancos de dados é um processo que envolve a extração de informações valiosas, padrões ocultos, conhecimento e insights úteis a partir de grandes volumes de dados. Esse processo faz parte do campo mais amplo da ciência de dados e é essencial para transformar dados brutos em informações de fácil acesso.

A descoberta de conhecimento em bancos de dados geralmente segue uma sequência de etapas, que podem variar em detalhes dependendo do contexto. A figura a seguir apresenta as etapas do KDD:

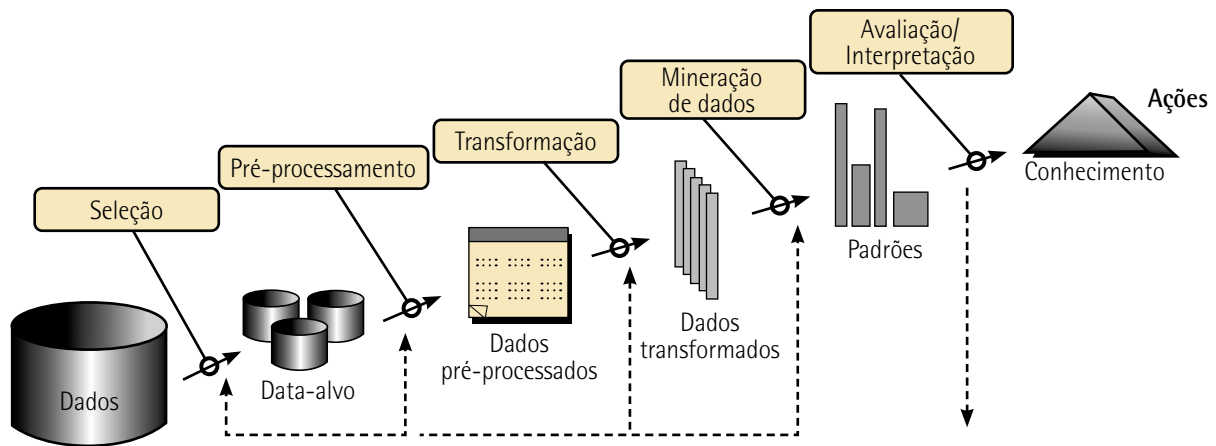


Figura 5 – Etapas do KDD

Adaptada de: Fayyad, Piatetsky-Shapiro e Smyth (1996).

A seguir descrevemos melhor as etapas do KDD:

- **Seleção de dados:** nesta etapa, os dados relevantes são identificados e selecionados para a análise. Isso envolve a definição de critérios de inclusão e exclusão e a obtenção dos conjuntos de dados adequados para o problema em questão.
- **Pré-processamento de dados:** os dados brutos podem ser complexos, inconsistentes ou conter ruído. Nesta etapa, ocorre a limpeza e a transformação dos dados, incluindo a remoção de dados ausentes ou duplicados, normalização, discretização e outras técnicas de preparação dos dados para análise.
- **Transformação de dados:** a transformação de dados é feita para representar os dados em uma forma mais adequada para análise. Isso pode envolver a agregação de dados, a criação de novos atributos ou a redução da dimensionalidade através de técnicas como análise de componentes principais (PCA) ou seleção de recursos.
- **Mineração de dados:** a mineração de dados é a etapa central do processo de descoberta de conhecimento. Nessa etapa, são aplicadas técnicas e algoritmos de aprendizado de máquina, estatística e visualização de dados para identificar padrões, tendências, associações ou relações interessantes nos dados. Isso pode incluir técnicas como classificação, regressão, clusterização, regras de associação, redes neurais, entre outras.
- **Avaliação e interpretação dos resultados:** após a aplicação das técnicas de mineração de dados, os resultados obtidos são avaliados e interpretados. Isso envolve a análise dos padrões descobertos, a validação dos modelos construídos e a interpretação dos insights obtidos em termos do problema ou domínio específico em questão.

- **Utilização e aplicação dos conhecimentos:** os conhecimentos e insights descobertos durante o processo são utilizados para tomar decisões informadas, desenvolver estratégias, resolver problemas e promover melhorias nos negócios ou em outras áreas de aplicação. A utilização efetiva dos conhecimentos é a finalidade última da descoberta de conhecimento em bancos de dados.

É importante ressaltar que a descoberta de conhecimento em bancos de dados não é um processo único e estático. É um ciclo contínuo de exploração, análise, interpretação e utilização dos dados. À medida que novos dados são coletados e o contexto evolui, o processo de descoberta de conhecimento pode ser iterado e refinado para capturar novos insights e informações valiosas.

A descoberta de conhecimento em bancos de dados desempenha um papel fundamental na ciência de dados e na tomada de decisões baseada em dados. Ela permite que as organizações aproveitem o poder dos dados para obter vantagens competitivas, identificar oportunidades de negócios, otimizar processos, personalizar experiências de usuário, prever tendências e comportamentos, entre outras aplicações.

Os desafios do processo de descoberta de conhecimento

Devemos destacar que a descoberta de conhecimento em bancos de dados também apresenta desafios. Com o aumento da disponibilidade de dados, lidar com grandes volumes de dados e a complexidade inerente a eles pode ser um desafio. A eficiência computacional e o uso de técnicas de processamento distribuído são necessários para lidar com esses dados em tempo hábil.

Considerando a qualidade dos dados pode ser um problema significativo na descoberta de conhecimento. Os dados podem conter erros, ruído, valores ausentes ou inconsistentes, o que pode afetar a precisão e a confiabilidade dos resultados obtidos. É importante realizar uma limpeza e uma validação adequadas dos dados antes de prosseguir com a análise.

Outro desafio importante é a interpretação dos resultados da descoberta de conhecimento que pode ser desafiadora, pois é importante considerar possíveis vieses nos dados e nos algoritmos. Além disso, a interpretação dos padrões descobertos requer um entendimento profundo do contexto e do domínio em que os dados estão inseridos. Embora a análise de dados seja automatizada, o conhecimento humano desempenha um papel crucial na descoberta de conhecimento. A integração de especialistas do domínio e de profissionais de análise de dados é necessária para interpretar os resultados, validar as descobertas e traduzir os insights em ações práticas.

A descoberta de conhecimento em bancos de dados também apresenta preocupações em relação à privacidade dos dados e ao uso ético das informações. É necessário garantir a conformidade com regulamentações de privacidade, como o LGPD, e adotar práticas responsáveis de proteção de dados ao lidar com informações sensíveis.

Superar esses desafios requer uma abordagem cuidadosa, combinando o conhecimento técnico, o domínio do negócio e as habilidades analíticas. Além disso, é fundamental manter uma mentalidade crítica e estar disposto a iterar e ajustar o processo de descoberta de conhecimento à medida que novos desafios surgem.

No geral, a descoberta de conhecimento em bancos de dados desempenha um papel vital na ciência de dados, permitindo a transformação de dados brutos em insights valiosos. Ao enfrentar os desafios e aplicar as melhores práticas, as organizações podem aproveitar ao máximo seus dados para obter uma vantagem competitiva e impulsionar a inovação em diversos setores.

O processo de KDD é um ciclo contínuo, pois novos dados podem estar disponíveis ou os objetivos podem mudar ao longo do tempo. Além disso, é importante ressaltar que o KDD não é uma tarefa única, mas sim um processo que requer uma abordagem sistemática e colaborativa.

Os benefícios do processo de descoberta de conhecimento

O objetivo final do processo de KDD é transformar dados em conhecimento valioso, que possa ser usado para tomar decisões mais informadas, gerar insights e impulsionar a inovação. A descoberta de conhecimento em bancos de dados desempenha um papel fundamental em diversos setores e aplicações, como negócios, medicina, finanças, marketing, ciências sociais, entre outros.

A descoberta de conhecimento permite que as organizações baseiem suas decisões em dados e evidências sólidas. Isso reduz a dependência de suposições ou intuições subjetivas, levando a decisões mais fundamentadas e de melhor qualidade. Através da mineração de dados, é possível identificar padrões, associações e tendências ocultas nos dados. Essas informações podem fornecer insights valiosos sobre comportamentos de clientes, preferências do mercado, tendências de vendas, riscos financeiros, entre outros aspectos relevantes para a organização.

A análise dos dados pode revelar ineficiências e gargalos em processos organizacionais. Com base nesse conhecimento, as empresas podem otimizar seus processos, reduzir custos, melhorar a produtividade e a eficiência operacional. Com a descoberta de conhecimento, é possível entender melhor os clientes e suas preferências. Isso permite personalizar ofertas, serviços e experiências, aumentando a satisfação do cliente e a fidelidade à marca.

A análise dos dados históricos e atuais pode ser utilizada para prever eventos futuros e antecipar tendências. Isso é especialmente útil em áreas como previsão de demanda, detecção de fraudes, análise de riscos e tomada de medidas preventivas. A descoberta de conhecimento em bancos de dados pode estimular a inovação e o desenvolvimento de novos produtos. Os insights obtidos a partir da análise dos dados podem revelar oportunidades de mercado, necessidades não atendidas e novas abordagens para a criação de produtos ou serviços.

Uma metodologia utilizada em projetos de ciência dos dados

O CRISP-DM (Cross-Industry Standard Process for Data Mining), conforme Chapman *et al.* (2000), é um modelo de processo muito utilizado na área de mineração de dados para guiar projetos de análise de dados (Piatetsky, 2014). Ele fornece uma estrutura flexível e abrangente para a condução de projetos de mineração de dados, permitindo que as equipes enfrentem desafios complexos e tomem decisões informadas ao longo de todo o ciclo de vida do projeto. O Crisp-DM é ilustrado na figura a seguir e cada uma das seis fases será explicada adiante:

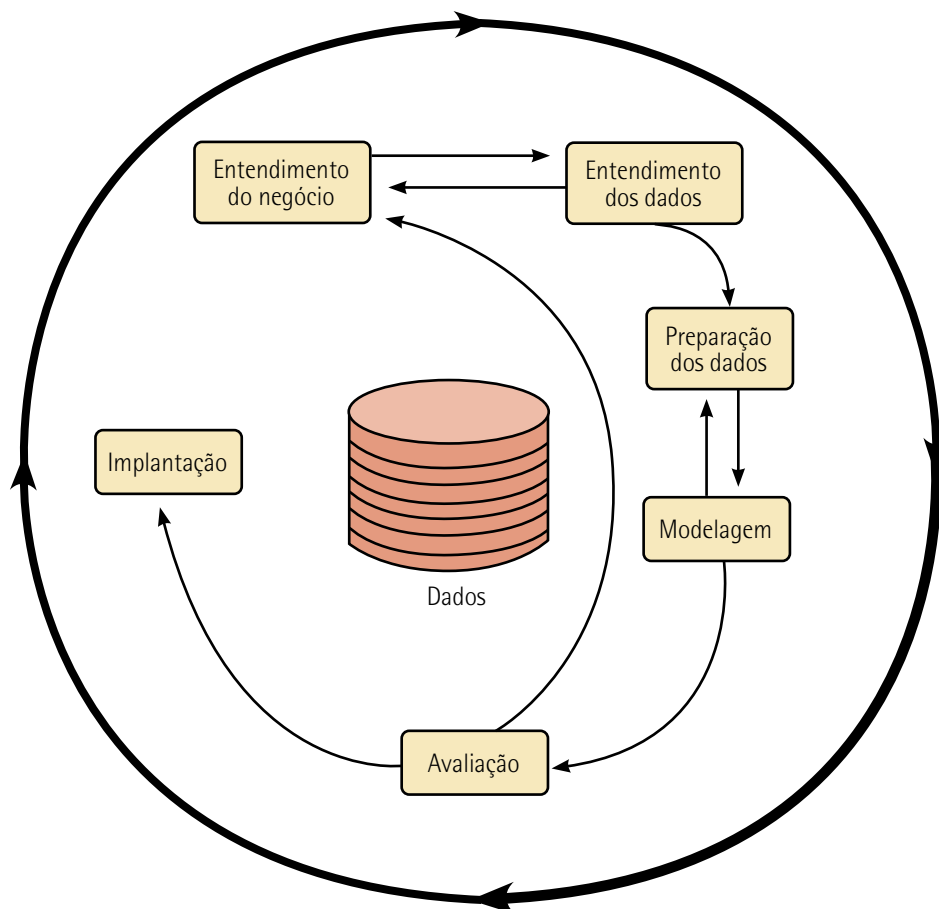


Figura 6 – Fases do modelo de referência CRISP-DM

Adaptada de: Chapman *et al.* (2000, p. 12).

- **Entendimento do negócio:** nesta fase inicial, a equipe trabalha para compreender os objetivos e requisitos do projeto, identificando como a mineração de dados pode contribuir para as metas de negócios.
- **Entendimento dos dados:** nesta fase, os dados disponíveis são explorados e analisados para identificar sua qualidade, relevância e potencial para atender aos objetivos do projeto. Isso envolve a realização de análises exploratórias e a compreensão das características dos dados.
- **Preparação dos dados:** aqui, os dados são limpos, transformados e preparados para análise. Isso inclui lidar com valores ausentes, normalização, seleção de atributos relevantes e outras tarefas de preparação.
- **Modelagem:** nesta fase, são desenvolvidos modelos de mineração de dados, como algoritmos de aprendizado de máquina, para explorar os padrões e relacionamentos nos dados. Diferentes abordagens são testadas e avaliadas para encontrar a mais adequada.

- **Avaliação:** os modelos construídos na fase de modelagem são avaliados para garantir que eles atendam aos critérios de sucesso do projeto. Isso pode envolver testes de desempenho, validação cruzada e outros métodos para garantir que os modelos sejam robustos e generalizáveis.
- **Implantação:** nesta última fase, os resultados da análise são apresentados aos stakeholders e são tomadas medidas para implementar os insights obtidos no ambiente de negócios. Isso pode envolver a criação de relatórios, integração com sistemas existentes ou outras formas de utilização prática.

É importante notar que o CRISP-DM é um processo iterativo e cíclico, o que significa que as fases não são necessariamente executadas em uma única passagem linear. À medida que a equipe ganha insights ao longo do processo, é comum voltar a fases anteriores para refinar abordagens ou ajustar estratégias. O CRISP-DM oferece uma estrutura sólida para guiar projetos de mineração de dados, promovendo uma abordagem sistemática e informada para a análise de dados em contextos variados.

1.8 Representação e extração de conhecimento

Representação e extração de conhecimento são etapas cruciais no processo de descoberta de conhecimento em bases de dados (KDD). Essas etapas envolvem a conversão dos dados brutos em uma forma adequada para análise e a extração de informações valiosas e conhecimento útil a partir desses dados. Vamos explorar esses conceitos em detalhes:

Representação de dados

A representação de dados refere-se à forma como os dados são organizados e estruturados para análise. É essencial escolher uma representação adequada para capturar as características relevantes dos dados e permitir a aplicação de técnicas de descoberta de conhecimento. Algumas técnicas comuns de representação de dados são:

- **Estruturas de dados:** os dados podem ser representados em diferentes estruturas, como tabelas, matrizes, grafos ou objetos. A escolha da estrutura depende da natureza dos dados e do tipo de análise a ser realizada.
- **Codificação:** os dados podem ser codificados usando diferentes formatos, como numérico, categórico, binário, texto ou outros formatos específicos. A codificação adequada dos dados é essencial para garantir que eles sejam interpretados corretamente pelos algoritmos de descoberta de conhecimento.
- **Vetorização:** em muitos casos, é necessário representar dados não estruturados, como texto ou imagens, em uma forma vetorial. A vetorização envolve a extração de características relevantes dos dados e a representação dessas características como vetores numéricos, permitindo a aplicação de técnicas de análise.

Extração de conhecimento

A extração de conhecimento envolve a aplicação de algoritmos e técnicas de descoberta de padrões, associações e tendências nos dados para identificar informações valiosas e conhecimento útil. Algumas técnicas comuns de extração de conhecimento incluem:

- **Mineração de dados:** a mineração de dados é um campo que abrange várias técnicas, como classificação, regressão, clusterização, regras de associação e detecção de anomalias. Essas técnicas são aplicadas para identificar padrões, relações e tendências nos dados que podem ser interpretados como conhecimento valioso.
- **Aprendizado de máquina:** o aprendizado de máquina é uma abordagem que permite aos sistemas aprenderem automaticamente com os dados e melhorarem seu desempenho ao longo do tempo. Os algoritmos de aprendizado de máquina são usados para treinar modelos que podem fazer previsões, tomar decisões ou identificar padrões nos dados.
- **Processamento de linguagem natural:** em casos em que os dados são textuais, técnicas de processamento de linguagem natural (NLP) podem ser usadas para extrair informações e conhecimentos significativos dos documentos de texto. Isso pode incluir a identificação de tópicos, a extração de entidades, a análise de sentimento e a sumarização automática.
- **Visualização de dados:** a visualização de dados desempenha um papel importante na extração de conhecimento, permitindo a representação gráfica dos padrões e relações encontrados nos dados. Gráficos, gráficos de dispersão, mapas de calor e outras técnicas visuais são utilizados para facilitar a compreensão e interpretação dos resultados.

Ao combinar essas técnicas de representação de dados e extração de conhecimento, é possível obter insights valiosos e conhecimento útil a partir dos dados. Essas informações podem ser usadas para tomar decisões informadas, identificar oportunidades de negócio, otimizar processos, prever tendências, personalizar experiências de usuário, entre outras aplicações.

Desafios da representação e extração de conhecimento

À medida que a quantidade de dados aumenta, a dimensionalidade dos dados também pode se tornar um desafio. Dados de alta dimensionalidade podem levar a problemas de escalabilidade e complexidade na representação e extração de conhecimento. Técnicas de redução de dimensionalidade, como análise de componentes principais (PCA) ou seleção de características, podem ser aplicadas para lidar com esse desafio. Os dados também podem conter ruído, erros ou valores ausentes, o que pode afetar a qualidade dos resultados obtidos. A limpeza e o tratamento adequado dos dados são essenciais para garantir a confiabilidade dos insights e do conhecimento extraído.

Devemos considerar possíveis tendências nos dados e nos algoritmos utilizados na representação e extração de conhecimento. Além disso, a interpretação dos resultados requer um entendimento profundo do contexto e do domínio dos dados, para evitar conclusões errôneas ou interpretações equivocadas.

Ao lidar com dados sensíveis, como informações pessoais dos clientes, é essencial garantir a privacidade e a conformidade com regulamentações de proteção de dados. As práticas éticas no uso dos dados devem ser seguidas, garantindo a anonimização, segurança e consentimento adequados dos indivíduos envolvidos.

Em algumas situações, os resultados da representação e extração de conhecimento podem ser complexos e difíceis de serem interpretados. É importante garantir que os resultados sejam explicáveis e compreensíveis, especialmente quando se trata de tomada de decisões críticas.

Apesar dos desafios, a representação e extração de conhecimento são componentes essenciais na jornada de descoberta de conhecimento em bancos de dados. Com as técnicas adequadas, os profissionais de ciência de dados podem extrair informações valiosas e conhecimento útil dos dados, impulsionando a inovação e tomada de decisões informadas em uma ampla gama de setores e aplicações.

1.9 Fontes de dados

As fontes de dados são os pontos de origem dos dados que serão analisados e processados para extrair informações valiosas e conhecimento útil. Essas fontes podem ser diversas e variar dependendo do domínio, dos objetivos do projeto e dos tipos de dados necessários. Vamos explorar algumas das principais fontes de dados comumente utilizadas em ciência de dados:

- **Bases de dados estruturadas:** as bases de dados estruturadas são um tipo comum de fonte de dados, que armazena informações em tabelas com colunas e linhas. Esses dados são organizados em um formato predefinido e são altamente acessíveis. Exemplos de bases de dados estruturadas incluem bancos de dados relacionais, como MySQL, Oracle e SQL Server.
- **Bases de dados não estruturadas:** as bases de dados não estruturadas armazenam dados em formatos que não seguem uma estrutura rígida, como documentos de texto, arquivos de áudio, imagens, vídeos e dados de mídia social. Essas fontes de dados são caracterizadas pela sua variedade e complexidade, e exigem técnicas especiais para extração e processamento, como processamento de linguagem natural (PLN) e análise de imagens.
- **Dados de sensores e dispositivos IoT:** com o crescimento da Internet das Coisas (IoT), os dispositivos estão gerando uma enorme quantidade de dados em tempo real. Sensores incorporados em veículos, edifícios, aparelhos domésticos e outras fontes capturam informações como temperatura, umidade, localização, movimento e muito mais. Esses dados são usados para monitorar e controlar sistemas, otimizar operações e tomar decisões com base em informações em tempo real.

Marjani *et al.* (2017) dividem em três etapas para permitir o gerenciamento de dados da IoT. O processo começa gerenciando as fontes de dados da IoT, em que dispositivos conectados interagem. Isso gera diversas fontes de dados com diferentes formatos, armazenadas em nuvem de baixo custo. Em seguida, os dados gerados são considerados Big Data devido ao seu volume, velocidade e variedade, sendo armazenados em bancos de dados compartilhados. Na etapa final, ferramentas como MapReduce, Spark, Splunk e Skytree são aplicadas para analisar os conjuntos de dados de IoT. A análise começa com dados de treinamento e há quatro níveis de análise.

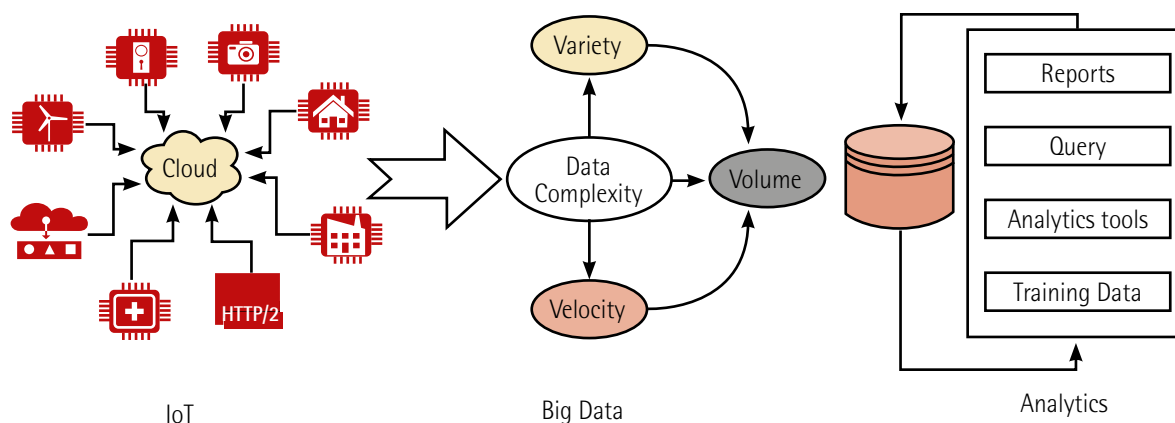


Figura 7 – Relação entre IoT e análise de Big Data

Adaptada de: Marjani et al. (2017, p. 5).



Saiba mais

Com o objetivo de compreender mais sobre Big Data e Internet das Coisas (IoT), recomenda-se a leitura dos capítulos "Big Data e Internet das Coisas (IoT)" e "Integrando Big Data e IoT" do livro a seguir:

MORAIS, I. S. et al. *Introdução a Big Data e Internet das Coisas (IoT)*. Porto Alegre: Grupo A, 2018. E-book.

Disponível em: <https://tinyurl.com/3ybwsed6>. Acesso em: 18 ago. 2023.

Dados de mídias sociais

As mídias sociais são uma fonte rica de dados para análise. Plataformas como Facebook, Twitter, Instagram e LinkedIn geram enormes volumes de dados relacionados a atividades de usuários, postagens, interações sociais, avaliações, comentários e muito mais. Esses dados são valiosos para entender o comportamento do usuário, sentimentos do cliente e tendências de mercado e também para realizar análises de redes sociais.

Dados de texto e documentos

Dados textuais, como documentos, relatórios, e-mails, artigos e registros, são fontes importantes de informações em muitos domínios. A análise de texto permite extrair conhecimento de grandes volumes de texto, identificar tópicos e realizar análise de sentimento, detecção de spam e categorização de documentos.

Dados de streaming

Dados de streaming são gerados em tempo real e requerem análise em tempo real. Eles podem incluir transações financeiras, feeds de sensores, dados de tráfego, dados de mídias sociais em tempo real e muito mais. Esses dados são processados e analisados em tempo real para identificar eventos, padrões ou anomalias em tempo hábil.

Dados de fontes externas

Além das fontes internas, os cientistas de dados podem aproveitar fontes externas de dados, como bancos de dados públicos, conjuntos de dados abertos, dados governamentais, dados de pesquisa acadêmica, feeds de dados climáticos e outras fontes disponíveis publicamente. Essas fontes externas podem enriquecer os dados internos da organização, fornecer contexto adicional e ampliar as possibilidades de análise e descoberta de conhecimento.

Dados transacionais

Dados transacionais são registros de transações e atividades comerciais, como vendas, compras, registros de clientes, registros de inventário e transações financeiras. Esses dados são fundamentais para análises de desempenho de negócios, identificação de tendências de vendas, detecção de fraudes e tomada de decisões estratégicas.

Dados geoespaciais

Dados geoespaciais incluem informações relacionadas à localização geográfica, como coordenadas GPS, endereços, mapas, imagens de satélite e dados de sensores geoespaciais. Esses dados são usados em aplicações de geolocalização, logística, planejamento urbano, monitoramento ambiental e outras análises que envolvem aspectos espaciais.

Dados históricos

Os dados históricos são registros passados de eventos, transações e atividades. Esses dados são usados para análise retrospectiva, identificação de padrões históricos, previsão de tendências futuras e modelagem preditiva. Os dados históricos são especialmente úteis em cenários financeiros, de mercado, clima e saúde, entre outros.

É importante ressaltar que a escolha e a combinação adequada das fontes de dados são essenciais para garantir a representatividade, qualidade e relevância dos dados utilizados na análise. Além disso, é necessário considerar a legalidade, privacidade e ética no acesso e uso das fontes de dados, especialmente quando se trata de dados sensíveis ou restritos.

A ciência de dados utiliza uma variedade de fontes de dados para extrair informações valiosas e conhecimento útil. A combinação correta das fontes de dados permite análises mais abrangentes e insights mais significativos, impulsionando a tomada de decisões informadas e a obtenção de vantagens competitivas nas organizações.

2 VISÃO GERAL SOBRE APRENDIZADO DE MÁQUINA

O aprendizado de máquina, também conhecido como machine learning, é uma subárea da inteligência artificial que se concentra no desenvolvimento de algoritmos e modelos capazes de aprender e tomar decisões a partir dos dados, sem serem explicitamente programados. O objetivo principal do aprendizado de máquina é permitir que os sistemas "aprendam" automaticamente a partir dos dados e melhorem seu desempenho ao longo do tempo, sem a necessidade de regras ou instruções específicas.

Em uma visão geral, o aprendizado de máquina pode ser dividido em três categorias principais:

- **Aprendizado supervisionado:** no aprendizado supervisionado, os algoritmos são treinados utilizando um conjunto de dados de entrada pré-rotulados, chamados de conjunto de treinamento. O objetivo é aprender a relação entre as entradas e as saídas correspondentes, para que o modelo seja capaz de fazer previsões ou tomar decisões em novos dados não vistos anteriormente. Exemplos de algoritmos de aprendizado de máquina supervisionado incluem regressão linear, regressão logística, árvores de decisão e redes neurais.
- **Aprendizado não supervisionado:** no aprendizado não supervisionado, os algoritmos são aplicados a conjuntos de dados não rotulados, ou seja, não há informações prévias sobre as saídas desejadas. O objetivo é descobrir estruturas, padrões ou grupos intrínsecos nos dados, fornecendo uma visão mais profunda dos dados e insights sobre seu comportamento. Algoritmos de clusterização, como K-means, e algoritmos de redução de dimensionalidade, como análise de componentes principais (PCA), são exemplos comuns de aprendizado de máquina não supervisionado.
- **Aprendizado semissupervisionado:** o aprendizado semissupervisionado é usado para as mesmas finalidades que o aprendizado supervisionado, porém envolve tanto dados com rótulos como sem rótulos para treinamento. Geralmente, uma pequena quantidade de dados tem rótulos e uma grande quantidade não tem (devido à menor complexidade e custo na aquisição). Ele pode ser aplicado a tarefas como classificação, regressão e previsão. É vantajoso quando rotular todos os dados é caro demais.

Aprendizado por reforço

O aprendizado por reforço envolve o treinamento de algoritmos através de interações com um ambiente. O agente de aprendizado toma ações em um ambiente e recebe recompensas ou punições com base no desempenho de suas ações. O objetivo é maximizar as recompensas ao longo do tempo, aprendendo a melhor política de ação. Isso é frequentemente aplicado em jogos, robótica e otimização de processos. Algoritmos de aprendizado por reforço incluem a Q-Learning, Sarsa e Deep Q-Networks (DQN).

O aprendizado de máquina é utilizado em diversas áreas, como reconhecimento de padrões, processamento de linguagem natural, visão computacional, recomendação de produtos, detecção de fraudes, diagnóstico médico, análise de mercado e muito mais. Ele oferece a capacidade de extrair conhecimento valioso dos dados, automatizar tarefas complexas e tomar decisões baseadas em informações, impulsionando a inovação e a transformação em diversos setores e indústrias.

Desafios do aprendizado de máquina

O desempenho dos modelos de aprendizado de máquina depende da disponibilidade de dados relevantes e de qualidade. É necessário ter conjuntos de dados representativos, com quantidade suficiente de exemplos e características adequadas para treinar os modelos de forma eficaz. Antes de aplicar técnicas de aprendizado de máquina, é preciso realizar a etapa de pré-processamento de dados, que envolve tratamento de valores ausentes, normalização, codificação de variáveis categóricas, entre outras tarefas. Essa etapa é crucial para garantir a qualidade dos dados e melhorar o desempenho dos modelos.

Nem todas as características dos dados podem ser relevantes para o problema em questão, e algumas podem até introduzir ruído. Portanto, é importante realizar a seleção adequada de recursos, identificando as características mais informativas e descartando as redundantes ou irrelevantes.

O overfitting ocorre quando um modelo se ajusta muito bem aos dados de treinamento, mas não consegue generalizar para novos dados. O underfitting, por sua vez, ocorre quando o modelo não é capaz de capturar os padrões dos dados de treinamento. Encontrar o equilíbrio entre esses dois problemas é fundamental para obter modelos com bom desempenho.

Alguns modelos de aprendizado de máquina, como redes neurais profundas, podem ser complexos e difíceis de interpretar. A interpretabilidade dos modelos é um desafio, especialmente em casos em que é preciso explicar as decisões tomadas pelo modelo ou quando são necessárias justificativas para questões éticas e legais.

À medida que a quantidade de dados aumenta, os modelos de aprendizado de máquina precisam ser escaláveis para lidar com grandes volumes de informações. Algoritmos eficientes e técnicas de processamento distribuído são necessários para lidar com conjuntos de dados massivos.

Apesar desses desafios, o aprendizado de máquina continua a evoluir e se tornar uma parte essencial das soluções baseadas em dados. Com avanços contínuos em algoritmos, infraestrutura de computação e técnicas de processamento de dados, o aprendizado de máquina tem o potencial de revolucionar várias áreas, impulsionando a automação, a tomada de decisões inteligentes e a inovação em um mundo cada vez mais orientado a dados.

2.1 O que é machine learning (ML)?

Machine learning, ou aprendizado de máquina, é um campo da inteligência artificial que envolve o desenvolvimento de algoritmos e modelos capazes de aprender e tomar decisões a partir dos dados, sem serem explicitamente programados. É uma abordagem que permite que as máquinas sejam treinadas para aprender com exemplos e experiências passadas, com o objetivo de automatizar tarefas complexas, fazer previsões, identificar padrões e tomar decisões com base em informações.

Em vez de programar regras específicas para cada cenário, o machine learning permite que os modelos aprendam a partir dos dados, detectem padrões e generalizem esse conhecimento para casos futuros. Isso é útil quando lidamos com grandes volumes de dados e situações complexas, nas quais seria difícil ou inviável criar um conjunto de regras manuais.

O processo de machine learning geralmente envolve as seguintes etapas:

- **Coleta de dados:** é necessário obter um conjunto de dados que seja representativo e relevante para o problema em questão. Esses dados podem ser coletados de diversas fontes, como bancos de dados, sensores, mídias sociais, entre outros.
- **Pré-processamento de dados:** os dados coletados podem precisar passar por uma etapa de pré-processamento, na qual são tratados valores ausentes e realizada normalização, codificação de variáveis categóricas e outras transformações para garantir que os dados estejam prontos para serem usados no treinamento dos modelos.
- **Seleção e engenharia de recursos:** nesta etapa, são selecionadas as características relevantes dos dados que serão utilizadas para treinar o modelo. Também é possível criar características, por meio de técnicas de engenharia de recursos, que possam melhor representar os padrões e informações nos dados.
- **Treinamento do modelo:** os dados preparados são utilizados para treinar o modelo de machine learning. Durante o treinamento, o modelo ajusta seus parâmetros de acordo com os exemplos fornecidos, buscando otimizar uma função de perda ou maximizar uma função de recompensa, dependendo do tipo de aprendizado.
- **Avaliação e ajuste do modelo:** após o treinamento, o modelo é avaliado para verificar seu desempenho em dados não utilizados no treinamento. São utilizadas métricas de avaliação adequadas ao problema em questão. Se necessário, o modelo pode ser ajustado, e o processo de treinamento e avaliação pode ser repetido até que se obtenha um desempenho satisfatório.
- **Aplicação do modelo:** após o treinamento e ajuste, o modelo está pronto para ser aplicado em novos dados. Ele é capaz de fazer previsões, classificações ou tomar decisões com base nas informações fornecidas, utilizando os padrões e conhecimentos aprendidos durante o treinamento.

Machine learning é aplicado em várias áreas, como reconhecimento de padrões, processamento de linguagem natural, visão computacional, recomendação de produtos, detecção de fraudes, diagnóstico médico, análise de mercado e muitas outras. É uma tecnologia poderosa que impulsiona a automação, a tomada de decisões inteligentes e a inovação em diversos setores.

2.2 Pipeline da aprendizagem do modelo

O pipeline da aprendizagem do modelo, também conhecido como pipeline de aprendizado de máquina, é uma sequência de etapas que são seguidas para desenvolver e treinar um modelo de aprendizado de máquina de forma estruturada. Esse pipeline abrange desde o pré-processamento dos dados até a avaliação e implantação do modelo em um ambiente de produção. Cada etapa do pipeline desempenha um papel fundamental no desenvolvimento e no treinamento bem-sucedido do modelo. As principais etapas do pipeline da aprendizagem do modelo são:

- **Coleta e pré-processamento dos dados:** nesta etapa, os dados são coletados a partir de fontes relevantes e podem passar por um processo de limpeza e pré-processamento. Isso inclui a remoção de valores ausentes, tratamento de outliers, normalização ou padronização dos dados, além de possíveis transformações para melhorar a qualidade e a representação dos dados.
- **Divisão dos dados:** os dados coletados são divididos em conjuntos de treinamento, validação e teste. O conjunto de treinamento é usado para treinar o modelo, o conjunto de validação é usado para ajustar os hiperparâmetros do modelo e o conjunto de teste é usado para avaliar o desempenho final do modelo em dados não vistos anteriormente.
- **Seleção de recursos:** nesta etapa, são selecionadas as características (ou atributos) mais relevantes dos dados que serão utilizados para treinar o modelo. A seleção adequada de recursos pode influenciar diretamente o desempenho do modelo, pois características irrelevantes podem introduzir ruído e prejudicar a capacidade de generalização.
- **Escolha do modelo:** com base no problema em questão e nas características dos dados, é selecionado o modelo de aprendizado de máquina mais adequado. Existem vários algoritmos disponíveis, como regressão linear, árvores de decisão, redes neurais, SVM (support vector machines) e muitos outros. A escolha do modelo depende do tipo de problema (classificação, regressão, agrupamento etc.) e das características dos dados.
- **Treinamento do modelo:** nesta etapa, o modelo é treinado utilizando o conjunto de treinamento. O objetivo é ajustar os parâmetros do modelo para minimizar a função de perda ou maximizar a função de recompensa, dependendo do tipo de aprendizado. Isso é feito através de algoritmos de otimização que buscam encontrar os melhores valores para os parâmetros do modelo.
- **Avaliação do modelo:** após o treinamento, o modelo é avaliado utilizando o conjunto de validação. Métricas adequadas ao problema são utilizadas para medir o desempenho do modelo, como acurácia, precisão, recall, F1-score, erro quadrático médio (MSE) etc. Essas métricas fornecem uma medida objetiva do desempenho do modelo e permitem a comparação entre diferentes modelos ou configurações.
- **Ajuste de hiperparâmetros:** além dos parâmetros aprendidos durante o treinamento, existem hiperparâmetros que controlam o comportamento do modelo. Esses hiperparâmetros precisam ser ajustados manualmente para otimizar o desempenho do modelo. Isso pode ser feito utilizando

técnicas como validação cruzada, busca em grade (grid search) ou otimização bayesiana para encontrar a combinação ideal de hiperparâmetros que leve a um desempenho máximo do modelo.

- **Teste final:** após o ajuste do modelo e a obtenção de resultados satisfatórios na etapa de validação, é hora de testar o modelo final no conjunto de teste, que contém dados não utilizados anteriormente. Essa etapa é fundamental para verificar o desempenho do modelo em dados "não viciados" e ter uma avaliação final do seu poder de generalização.
- **Implantação e monitoramento:** uma vez que o modelo tenha sido treinado e testado, ele pode ser implantado em um ambiente de produção, onde estará pronto para fazer previsões ou tomar decisões em tempo real. Durante a implantação, é importante monitorar o desempenho contínuo do modelo, avaliando sua precisão e eficácia. Isso pode envolver o acompanhamento das métricas de desempenho, detecção de anomalias e recalibração do modelo, se necessário.
- **Melhoria contínua:** o pipeline da aprendizagem do modelo é um processo iterativo e contínuo. À medida que novos dados são coletados e o modelo é usado em produção, é possível identificar áreas de melhoria e refinamento. Isso pode envolver a inclusão de novos recursos, a exploração de técnicas mais avançadas de modelagem, a atualização do modelo com novos dados ou o desenvolvimento de modelos mais complexos. A melhoria contínua garante que o modelo esteja atualizado, preciso e relevante ao longo do tempo.

O pipeline da aprendizagem do modelo fornece uma estrutura para guiar o processo de desenvolvimento de modelos de aprendizado de máquina, garantindo a utilização adequada dos dados, a escolha do modelo correto, o treinamento eficiente e a avaliação precisa. Seguir um pipeline bem definido ajuda a otimizar o tempo e os recursos gastos no desenvolvimento de modelos e aumenta as chances de obter resultados confiáveis e de alta qualidade.

2.3 Overfitting e underfitting

O overfitting e o underfitting são dois fenômenos indesejáveis no aprendizado de máquina que afetam a capacidade do modelo de generalizar corretamente a partir dos dados de treinamento. Eles ocorrem quando o modelo se ajusta demais ou não se ajusta o suficiente aos dados, resultando em um desempenho insatisfatório na etapa de teste. Vamos entender melhor cada um desses conceitos:

Overfitting (sobreajuste)

O overfitting ocorre quando o modelo se torna excessivamente complexo e se ajusta perfeitamente aos dados de treinamento, capturando até mesmo o ruído presente nesses dados. Como resultado, o modelo memoriza os exemplos de treinamento em vez de aprender os padrões subjacentes que permitem generalizar para novos dados. Isso pode levar a um desempenho pobre na etapa de teste, em que o modelo falha em fazer previsões precisas em dados não vistos.

Veja alguns sinais de overfitting:

- O desempenho do modelo é excelente nos dados de treinamento, mas ruim nos dados de teste.
- O modelo possui uma complexidade excessiva em relação ao tamanho dos dados disponíveis.
- O modelo captura o ruído presente nos dados de treinamento, resultando em uma precisão excessivamente alta nesses dados, mas não em novos dados.

Algumas técnicas para combater o overfitting incluem:

- Aumentar o tamanho do conjunto de treinamento, se possível.
- Utilizar técnicas de regularização, como L1 ou L2, para penalizar pesos excessivamente grandes.
- Utilizar técnicas de seleção de recursos para reduzir a dimensionalidade dos dados.
- Utilizar validação cruzada para avaliar o desempenho do modelo em conjuntos de dados diferentes.

Underfitting (subajuste)

O underfitting ocorre quando o modelo é muito simples ou não é capaz de capturar os padrões presentes nos dados de treinamento. Nesse caso, o modelo não consegue se ajustar adequadamente aos dados e acaba subestimando a complexidade do problema. O resultado é um desempenho insatisfatório tanto nos dados de treinamento quanto nos dados de teste.

Veja alguns sinais de underfitting:

- O desempenho do modelo é ruim tanto nos dados de treinamento quanto nos dados de teste.
- O modelo não consegue capturar os padrões e relações importantes presentes nos dados.
- O modelo é muito simples em relação à complexidade do problema.

Algumas técnicas para combater o underfitting incluem:

- Aumentar a complexidade do modelo, adicionando mais camadas ou neurônios em redes neurais, por exemplo.
- Adicionar mais recursos ou características relevantes aos dados.
- Utilizar algoritmos mais avançados ou complexos.

Equilibrar entre o overfitting e o underfitting é um desafio importante no aprendizado de máquina. O objetivo é encontrar um modelo que generalize bem para dados não vistos, capturando os padrões e relações essenciais sem se ajustar demais aos dados de treinamento. Essa busca pelo equilíbrio é fundamental para obter um modelo com bom desempenho e capacidade de generalização. A figura a seguir apresenta um exemplo de underfitting, de overfitting e de como poderia ficar o equilíbrio entre os dados ajustando o hiperplano.

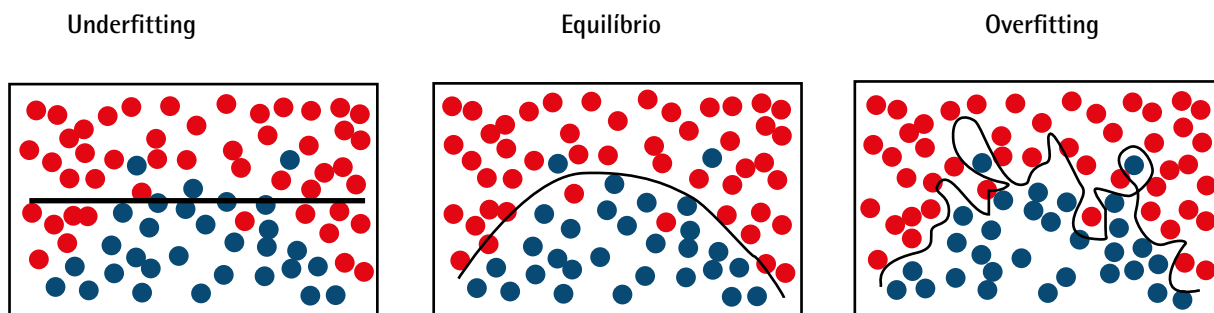


Figura 8 – Exemplo de underfitting, de overfitting e do equilíbrio entre os dados

Além das técnicas mencionadas anteriormente para combater o overfitting e o underfitting, existem outras abordagens que podem ser úteis:

- **Ajuste de hiperparâmetros:** os hiperparâmetros são parâmetros que controlam o comportamento do modelo de aprendizado de máquina, como a taxa de aprendizado, o número de camadas em uma rede neural ou a profundidade de uma árvore de decisão. Ajustar adequadamente esses hiperparâmetros pode ajudar a encontrar um equilíbrio entre o overfitting e o underfitting. Isso pode ser feito manualmente, através de tentativa e erro, ou utilizando técnicas de otimização automática, como busca em grade (grid search) ou otimização bayesiana.
- **Ensemble learning:** o ensemble learning é uma abordagem em que vários modelos são combinados para melhorar o desempenho geral. Essa técnica é útil para lidar tanto com o overfitting quanto com o underfitting. Alguns métodos populares de ensemble learning incluem bagging, boosting e random forests. Essas técnicas combinam as previsões de vários modelos individuais para obter uma previsão mais robusta e geralmente mais precisa.
- **Cross-validation:** a validação cruzada é uma técnica que permite avaliar o desempenho do modelo de forma mais robusta, fornecendo uma estimativa mais confiável de como o modelo se sairá em dados não vistos. O k-fold cross-validation é um exemplo comum, em que os dados são divididos em k subconjuntos. O modelo é treinado e avaliado k vezes, cada vez utilizando um subconjunto diferente como conjunto de teste e o restante como conjunto de treinamento. Essa abordagem ajuda a avaliar o desempenho médio do modelo e a detectar problemas de overfitting ou underfitting.

- **Aumento de dados (data augmentation):** esta técnica é utilizada em problemas de visão computacional e processamento de linguagem natural. Consiste em gerar dados sintéticos adicionais, aplicando transformações aos dados de treinamento existentes. Por exemplo, ao treinar um modelo de reconhecimento de imagens, é possível aplicar rotações, zooms, espelhamentos ou outras transformações nas imagens para aumentar a diversidade dos dados de treinamento. Isso pode ajudar a evitar o overfitting, fornecendo mais exemplos variados para o modelo aprender.

Combater o overfitting e o underfitting é essencial para desenvolver modelos de aprendizado de máquina robustos e de alta qualidade. É importante encontrar um equilíbrio entre a complexidade do modelo e a capacidade de generalização, além de utilizar técnicas como ajuste de hiperparâmetros, ensemble learning, validação cruzada e aumento de dados para melhorar o desempenho do modelo. A escolha correta dessas técnicas dependerá do problema em questão e das características dos dados disponíveis.

2.4 Balanço entre viés e variância em modelos de ML

O balanço entre viés e variância é um conceito fundamental no desenvolvimento de modelos de aprendizado de máquina. Ele está relacionado à capacidade do modelo de generalizar corretamente a partir dos dados de treinamento. Vamos entender melhor esses dois conceitos:

- **Viés (bias):** o viés de um modelo é a simplificação ou suposições errôneas que ele faz sobre os dados. Um modelo com alto viés tende a fazer suposições simplistas e ignorar detalhes complexos dos dados. Isso pode resultar em um modelo subajustado, que não é capaz de capturar os padrões e relações importantes presentes nos dados. Um modelo com alto viés pode ter um desempenho ruim tanto nos dados de treinamento quanto nos dados de teste.
- **Variância (variance):** a variância de um modelo refere-se à sensibilidade excessiva do modelo aos dados de treinamento. Um modelo com alta variância é altamente complexo e se ajusta perfeitamente aos dados de treinamento, mas não generaliza bem para novos dados. Isso pode resultar em um modelo sobreajustado, que memoriza os exemplos de treinamento e não é capaz de capturar os padrões subjacentes que permitem fazer previsões precisas em novos dados. Um modelo com alta variância pode ter um desempenho excelente nos dados de treinamento, mas um desempenho ruim nos dados de teste.

O objetivo é encontrar um equilíbrio entre viés e variância, em que o modelo seja suficientemente complexo para capturar os padrões importantes nos dados, mas não seja excessivamente complexo a ponto de se ajustar ao ruído presente nos dados de treinamento. Esse equilíbrio garante que o modelo seja capaz de generalizar bem para novos dados.

Existem algumas abordagens para encontrar esse equilíbrio.

- **Regularização:** a regularização é uma técnica que permite controlar a complexidade do modelo, adicionando um termo de penalidade aos parâmetros do modelo durante o processo de treinamento. Isso ajuda a evitar o overfitting, reduzindo a variância do modelo.

- **Ajuste de hiperparâmetros:** os hiperparâmetros do modelo, como a taxa de aprendizado, o número de camadas em uma rede neural ou a profundidade de uma árvore de decisão, podem ser ajustados para encontrar o equilíbrio entre viés e variância. Por exemplo, aumentar a complexidade do modelo pode reduzir o viés, mas aumentar a variância, enquanto diminuir a complexidade pode reduzir a variância, mas aumentar o viés.
- **Ensemble learning:** o ensemble learning é uma técnica em que vários modelos são combinados para melhorar o desempenho geral. Essa abordagem ajuda a reduzir a variância do modelo, combinando as previsões de vários modelos individuais. Exemplos de ensemble learning incluem bagging, boosting e random forests.
- **Aumento de dados (data augmentation):** aumentar a diversidade dos dados de treinamento por meio de técnicas de aumento de dados pode ajudar a reduzir o viés do modelo, fornecendo mais exemplos e variações para o modelo aprender.

Encontrar o equilíbrio certo entre viés e variância é uma tarefa desafiadora no desenvolvimento de modelos de aprendizado de máquina. Isso requer experimentação, ajuste fino e compreensão do problema em questão. É importante lembrar que não existe uma solução única para todos os problemas. O equilíbrio entre viés e variância dependerá da complexidade do problema, da quantidade e qualidade dos dados disponíveis, bem como das características específicas do modelo e dos algoritmos utilizados.

É comum que, ao ajustar um modelo para reduzir o viés, a variância aumente e vice-versa. Portanto, encontrar o ponto ideal de equilíbrio é uma tarefa delicada. É essencial realizar uma avaliação cuidadosa do desempenho do modelo em dados de treinamento e teste, bem como utilizar técnicas como validação cruzada para obter estimativas mais confiáveis do desempenho geral.

Um erro comum é focar demais na redução do erro nos dados de treinamento, levando a um modelo sobreajustado. É fundamental considerar também o desempenho em dados não vistos, para garantir que o modelo generalize bem e seja capaz de fazer previsões precisas em situações reais.

O balanço entre viés e variância é um aspecto importante na modelagem de aprendizado de máquina. Buscar um modelo com baixo viés e baixa variância é o objetivo, mas encontrar esse equilíbrio requer uma compreensão profunda dos dados, técnicas apropriadas de ajuste de modelo e uma avaliação cuidadosa do desempenho. Encontrar o ponto ideal entre viés e variância ajudará a construir modelos mais precisos, confiáveis e generalizáveis.

2.5 Viés indutivo

Viés indutivo é um conceito importante no campo do aprendizado de máquina e refere-se às suposições ou tendências que um algoritmo de aprendizado faz ao inferir padrões a partir dos dados de treinamento. O viés indutivo pode ser entendido como um conjunto de preconceitos incorporados ao algoritmo, que moldam a maneira como ele generaliza a partir dos exemplos de treinamento.

Cada algoritmo de aprendizado de máquina possui um viés indutivo, que é uma consequência das escolhas feitas na definição do espaço de hipóteses ou das restrições impostas ao modelo. O viés indutivo é necessário, pois permite ao algoritmo realizar inferências e generalizar além dos dados de treinamento. No entanto, o viés também pode levar a erros se as suposições subjacentes ao algoritmo não estiverem de acordo com a natureza dos dados.

O viés indutivo pode assumir diferentes formas, dependendo do algoritmo e do contexto do problema. Alguns algoritmos de aprendizado de máquina têm um viés que favorece determinados tipos de funções. Por exemplo, em algoritmos de regressão linear, há um viés indutivo para soluções que se ajustam a uma linha reta. Isso significa que o algoritmo tende a buscar soluções lineares mesmo quando o relacionamento entre as variáveis não é linear. Outros algoritmos assumem independência entre as variáveis, como na técnica de Naive Bayes. Isso significa que o algoritmo supõe que as variáveis são independentes umas das outras ao fazer suas previsões. Essa suposição nem sempre é verdadeira na prática, mas é feita para simplificar o problema e facilitar o cálculo. Alguns algoritmos de aprendizado que aplicam técnicas de regularização, como regressão ridge ou lasso, possuem um viés indutivo para preferir modelos mais simples. A regularização adiciona uma penalidade para evitar coeficientes excessivamente grandes, o que ajuda a evitar o overfitting e a favorecer soluções mais parcimoniosas.

O viés indutivo é uma escolha projetada pelos desenvolvedores do algoritmo e pode ser adaptado às características e aos requisitos específicos do problema em questão. É importante compreender o viés indutivo do algoritmo selecionado, pois ele influenciará a maneira como o modelo aprende e generaliza a partir dos dados de treinamento. O viés indutivo pode levar a resultados melhores ou piores, dependendo da adequação entre as suposições do algoritmo e a realidade dos dados.

O viés indutivo é um componente essencial dos algoritmos de aprendizado de máquina e desempenha um papel importante na generalização e inferência a partir dos dados de treinamento. Compreender o viés indutivo do algoritmo escolhido é crucial para tomar decisões adequadas na modelagem e no ajuste do modelo de aprendizado de máquina.

2.6 Sistema de aprendizado

Um sistema de aprendizado, também conhecido como sistema de aprendizado de máquina, é um conjunto de componentes e algoritmos que permitem que uma máquina aprenda a partir dos dados e faça previsões ou tomadas de decisão com base nesse aprendizado. O objetivo é capacitar a máquina a reconhecer padrões, extrair informações úteis e melhorar seu desempenho ao longo do tempo, sem ser explicitamente programada para cada tarefa específica.

O sistema de aprendizado requer um conjunto de dados de treinamento, que são exemplos de entrada juntamente com suas saídas correspondentes. Os algoritmos de aprendizado são os principais componentes do sistema de aprendizado. Com base nos algoritmos de aprendizado, é construído o modelo, que é a representação aprendida a partir dos dados. O modelo captura os padrões e relações descobertos durante o treinamento e é capaz de fazer previsões ou tomar decisões com base nesse aprendizado.

O desempenho do modelo utilizado no sistema de aprendizado é avaliado usando dados de teste ou validação, que são separados dos dados de treinamento. A avaliação do modelo permite verificar sua capacidade de generalização e medir sua precisão e desempenho em dados não vistos. Após a avaliação e validação do modelo, ele pode ser implantado em um ambiente de produção, no qual é utilizado para fazer inferências ou previsões em dados novos. A implantação do modelo requer a integração com outros sistemas e a implementação de um pipeline de inferência para processar os dados de entrada e gerar as saídas desejadas.

Um sistema de aprendizado requer monitoramento contínuo para garantir seu bom desempenho ao longo do tempo. É importante analisar regularmente os resultados das previsões ou decisões do modelo e atualizá-lo conforme novos dados se tornem disponíveis. Isso permite aprimorar o modelo e garantir sua relevância e precisão contínuas.

Os sistemas de aprendizado têm uma ampla gama de aplicações em diversos setores, como saúde, finanças, varejo, automação industrial e muito mais. Eles são capazes de lidar com grandes volumes de dados, descobrir padrões complexos e fornecer insights valiosos para melhorar a tomada de decisões e automatizar tarefas. À medida que os algoritmos e as técnicas de aprendizado de máquina avançam, os sistemas de aprendizado se tornam cada vez mais poderosos e capazes de enfrentar desafios complexos do mundo real.

Um sistema de aprendizado é composto de vários componentes inter-relacionados, desde a preparação e treinamento dos dados até a implantação e monitoramento contínuo do modelo. A combinação adequada de técnicas, algoritmos e considerações éticas é essencial para construir sistemas de aprendizado eficazes e confiáveis, capazes de extrair informações valiosas e melhorar a tomada de decisões em uma ampla variedade de domínios.

2.7 Tipos de aprendizagem

Existem vários tipos de aprendizado de máquina que podem ser aplicados em diferentes cenários, dependendo das características dos dados e dos objetivos do problema. Aqui estão alguns dos principais tipos de aprendizado:

- **Aprendizado supervisionado:** no aprendizado supervisionado, o modelo é treinado usando pares de entrada e saída rotulados. O objetivo é aprender a mapear os dados de entrada para as saídas desejadas com base nos exemplos fornecidos. O modelo é treinado usando algoritmos como regressão linear, árvores de decisão, redes neurais, entre outros. Esse tipo de aprendizado é usado para tarefas de classificação, em que o objetivo é prever uma classe ou categoria, e para tarefas de regressão, nas quais o objetivo é prever um valor contínuo.
- **Aprendizado não supervisionado:** no aprendizado não supervisionado, o modelo é treinado em dados de entrada não rotulados. O objetivo é descobrir padrões, estruturas ou relações intrínsecas nos dados. Os algoritmos de aprendizado não supervisionado incluem técnicas de clusterização, como o K-means e agrupamento hierárquico, e técnicas de redução de dimensionalidade, como análise de componentes principais (PCA) e análise de fatores.

- **Aprendizado por reforço:** no aprendizado por reforço, o modelo aprende a tomar decisões sequenciais com base em interações com um ambiente. O modelo recebe feedback em forma de recompensas ou penalidades, dependendo das ações tomadas. O objetivo é maximizar a recompensa cumulativa ao longo do tempo, aprendendo a tomar ações que levem ao melhor resultado. Algoritmos populares de aprendizado por reforço incluem Q-Learning e Monte Carlo.
- **Aprendizado semissupervisionado:** no aprendizado semissupervisionado, o modelo é treinado com uma combinação de dados rotulados e não rotulados. A ideia é utilizar a informação disponível nos dados rotulados e a estrutura dos dados não rotulados para melhorar o desempenho do modelo. Isso é especialmente útil quando a obtenção de rótulos para todos os dados é difícil ou custosa.

Esses são apenas alguns exemplos dos tipos de aprendizado de máquina. Cada tipo tem suas próprias características e é adequado para diferentes problemas e cenários. Em muitos casos, combinações desses tipos de aprendizado são usadas para abordar problemas mais complexos e desafiadores. Além disso, é importante destacar que a escolha do tipo de aprendizado depende das características dos dados, dos recursos disponíveis e dos objetivos específicos do projeto.

É válido mencionar também que alguns problemas de aprendizado de máquina podem exigir técnicas de aprendizado híbridas, que combinam elementos de diferentes tipos de aprendizado. Por exemplo, pode-se utilizar aprendizado não supervisionado para realizar um pré-agrupamento dos dados e, em seguida, aplicar aprendizado supervisionado em cada grupo individualmente.

É importante ressaltar que o tipo de aprendizado escolhido pode ter impacto na performance e na eficácia do modelo. Portanto, é essencial compreender as características e limitações de cada tipo de aprendizado, além de explorar abordagens diferentes para determinar a mais adequada para cada caso.

A escolha do tipo de aprendizado de máquina depende do tipo de dado disponível, da natureza do problema e dos objetivos do projeto. A compreensão das diferenças entre os tipos de aprendizado permite selecionar as técnicas e algoritmos mais apropriados para resolver problemas específicos e alcançar resultados de alta qualidade.

2.8 Espaço de hipóteses

No contexto do aprendizado de máquina, o espaço de hipóteses se refere ao conjunto de todas as possíveis funções ou modelos que um algoritmo de aprendizado pode escolher como solução para um determinado problema. Essas hipóteses são expressas através de parâmetros, pesos ou estruturas específicas, dependendo do algoritmo e do tipo de aprendizado utilizado.

O espaço de hipóteses define as restrições sobre o conjunto de soluções possíveis que o algoritmo de aprendizado pode explorar durante o processo de treinamento. É um espaço multidimensional que representa as diferentes combinações e configurações que o modelo pode assumir para mapear os dados de entrada para as saídas desejadas.

O tamanho e a complexidade do espaço de hipóteses podem variar dependendo do tipo de modelo e do conjunto de características consideradas. Por exemplo, em um modelo linear simples, o espaço de hipóteses seria definido pelos diferentes valores dos pesos atribuídos a cada característica. Já em um modelo de redes neurais com várias camadas e muitos neurônios, o espaço de hipóteses seria muito maior e mais complexo.

A busca pelo melhor modelo dentro do espaço de hipóteses é realizada pelo algoritmo de aprendizado, que tenta encontrar a hipótese que melhor se ajusta aos dados de treinamento e que generaliza bem para dados não vistos. A qualidade da solução encontrada depende da representatividade do espaço de hipóteses em relação ao verdadeiro relacionamento entre os dados de entrada e saída.

Uma abordagem comum é utilizar uma função de perda ou critério de desempenho que quantifica o quão bem a hipótese se ajusta aos dados de treinamento. O algoritmo de aprendizado então busca minimizar essa função de perda para encontrar a melhor hipótese dentro do espaço de hipóteses.

É importante destacar que a escolha do espaço de hipóteses pode afetar o desempenho do modelo. Se o espaço de hipóteses for muito restrito, o modelo pode não ter capacidade suficiente para capturar a complexidade dos dados, levando a um subajuste (underfitting). Por outro lado, se o espaço de hipóteses for muito amplo, o modelo pode se tornar excessivamente complexo e se ajustar demais aos dados de treinamento, resultando em um sobreajuste (overfitting).

Portanto, encontrar um equilíbrio adequado no espaço de hipóteses é um desafio no aprendizado de máquina. É necessário considerar a complexidade dos dados, a quantidade de dados de treinamento disponíveis e as restrições computacionais para determinar um espaço de hipóteses que permita um bom desempenho e generalização do modelo.

Dentro do espaço de hipóteses, é comum encontrar diferentes subclasses de hipóteses, que representam configurações específicas do modelo ou restrições impostas sobre as soluções possíveis. Essas subclasses podem ser definidas com base em suposições prévias sobre o problema, restrições de recursos computacionais ou conhecimento prévio do domínio.

Além disso, o espaço de hipóteses não precisa ser estático. Em alguns casos, é possível adaptar ou expandir o espaço de hipóteses durante o processo de aprendizado. Isso pode ser feito usando técnicas como seleção de características, seleção de modelos, ajuste de hiperparâmetros ou até mesmo combinação de diferentes modelos.

A seleção adequada do espaço de hipóteses é crucial para obter um bom desempenho do modelo. Se o espaço de hipóteses for muito limitado e não incluir a verdadeira função subjacente, o modelo não será capaz de aprender a relação entre os dados de forma precisa. Por outro lado, se o espaço de hipóteses for muito amplo, o modelo pode se ajustar demasiadamente aos dados de treinamento e ter dificuldade em generalizar para novos dados.

Encontrar o melhor espaço de hipóteses é um desafio importante no aprendizado de máquina. Isso requer uma combinação de conhecimento do domínio, experiência, exploração de diferentes abordagens e técnicas, além de experimentação e avaliação criteriosa do desempenho do modelo. A capacidade de definir e explorar adequadamente o espaço de hipóteses é fundamental para construir modelos de aprendizado de máquina eficazes e com boa capacidade de generalização.

2.9 Viés de busca: ajuste aos dados

O viés de busca refere-se à tendência de um algoritmo de aprendizado de máquina preferir certos tipos de modelos ou hipóteses em relação a outros. Esse viés afeta a maneira como o algoritmo realiza a busca no espaço de hipóteses à procura da melhor solução para um problema.

O viés de busca é influenciado por várias razões, incluindo a escolha do algoritmo de aprendizado, a definição do espaço de hipóteses e as suposições subjacentes ao modelo. Essas suposições podem ser explícitas, como a escolha de um modelo linear em vez de um modelo não linear, ou implícitas, como a suposição de que os dados seguem uma distribuição gaussiana.

O ajuste aos dados é uma consequência do viés de busca. Quando o modelo é treinado, ele tenta encontrar a melhor hipótese dentro do espaço de hipóteses disponível. No entanto, o modelo pode se ajustar demasiadamente aos dados de treinamento, capturando ruídos e padrões irrelevantes. Isso é conhecido como overfitting ou sobreajuste.

O overfitting ocorre quando o modelo se torna muito complexo em relação à quantidade de dados de treinamento disponíveis. O modelo pode memorizar os exemplos de treinamento em vez de generalizar corretamente para novos dados. Como resultado, o desempenho do modelo pode ser excelente nos dados de treinamento, mas se deteriorar significativamente em dados não vistos.

Para evitar o overfitting, é importante considerar o viés de busca e encontrar um equilíbrio adequado entre a capacidade do modelo e a quantidade de dados disponíveis. Reduzir a complexidade do modelo, aumentar a quantidade de dados de treinamento ou usar técnicas de regularização, como a adição de termos de penalização nos pesos do modelo, são abordagens comuns para mitigar o overfitting.

Por outro lado, o ajuste insuficiente aos dados, conhecido como underfitting ou subajuste, ocorre quando o modelo é muito simples em relação à complexidade dos dados. Nesse caso, o modelo não consegue capturar corretamente a estrutura dos dados e apresenta um desempenho insatisfatório tanto nos dados de treinamento quanto nos dados não vistos.

3 APRENDIZADO DESCRITIVO E PREDITIVO

O **aprendizado descritivo** e o **aprendizado preditivo** são duas abordagens distintas dentro do campo do aprendizado de máquina, cada uma com objetivos e métodos específicos.

O **aprendizado descritivo**, também conhecido como aprendizado não supervisionado, tem como objetivo explorar os dados e descobrir padrões, estruturas ou relações ocultas neles, sem a necessidade de rótulos ou saídas conhecidas. Nessa abordagem, o algoritmo de aprendizado analisa os dados de entrada e tenta encontrar agrupamentos naturais, similaridades, anomalias ou outras características relevantes.

Os algoritmos de aprendizado descritivo incluem técnicas como agrupamento (clusterização), análise de componentes principais (PCA), regras de associação e detecção de anomalias. Esses algoritmos são úteis para explorar grandes conjuntos de dados, identificar grupos semelhantes de exemplos e entender a estrutura subjacente dos dados. Eles são utilizados em áreas como segmentação de mercado, detecção de fraudes, análise de redes sociais e processamento de linguagem natural.

Já o **aprendizado preditivo**, também conhecido como aprendizado supervisionado, tem como objetivo construir um modelo capaz de fazer previsões ou classificar novos exemplos com base em dados de treinamento rotulados. Nessa abordagem, o algoritmo de aprendizado recebe um conjunto de exemplos de entrada, com suas respectivas saídas conhecidas, e aprende a mapear corretamente as entradas para as saídas desejadas.

Os algoritmos de aprendizado preditivo incluem técnicas como árvores de decisão, regressão linear, redes neurais, máquinas de vetores de suporte (SVM) e algoritmos de classificação como k-vizinhos mais próximos (k-NN) e Naive Bayes. Esses algoritmos são utilizados em problemas de classificação, regressão e previsão, em que é necessário inferir uma função que mapeia os dados de entrada para as saídas desejadas.

O aprendizado preditivo é usado em uma variedade de aplicações, como previsão de vendas, diagnóstico médico, detecção de spam, reconhecimento de padrões, entre outros. Ele permite que os modelos sejam treinados com base em exemplos rotulados e usem essa informação para fazer previsões precisas em novos dados não rotulados.

A figura a seguir apresenta uma separação didática entre aprendizado descritivo e aprendizado preditivo, pois ambas as abordagens têm suas próprias técnicas e métodos específicos, mas podem ser complementares em muitos casos. Por exemplo, o aprendizado descritivo pode ser usado como uma etapa de pré-processamento para identificar padrões nos dados e gerar características relevantes para alimentar os algoritmos de aprendizado preditivo. Veremos a seguir as técnicas de classificação, regressão, agrupamento, associação e sumarização.

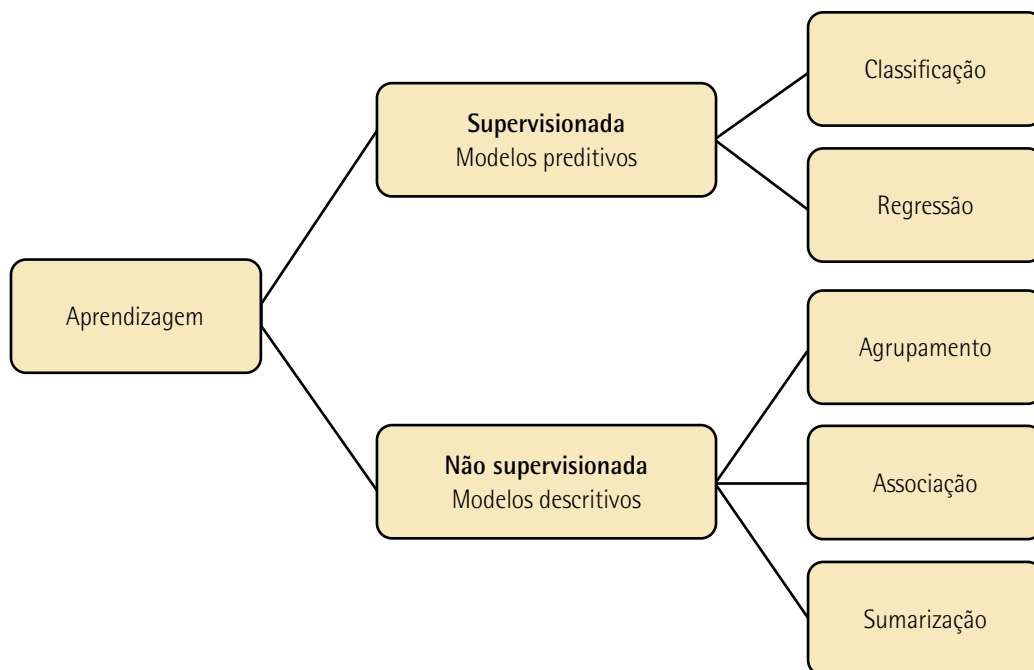


Figura 9 – Organização para o aprendizado de máquina

O aprendizado descritivo busca explorar e entender os dados, descobrindo estruturas e padrões ocultos, enquanto o aprendizado preditivo visa construir modelos que possam fazer previsões ou classificações com base nos dados disponíveis. Ambas as abordagens desempenham um papel fundamental no campo do aprendizado de máquina e são aplicadas em uma ampla variedade de problemas e setores.

3.1 Aprendizado supervisionado

O aprendizado supervisionado é uma abordagem do campo de aprendizado de máquina em que um algoritmo é treinado para aprender a mapear um conjunto de dados de entrada para suas respectivas saídas conhecidas. Nessa abordagem, o algoritmo recebe um conjunto de exemplos de treinamento, que consiste em pares de entrada e saída rotulados, e utiliza essas informações para aprender um modelo capaz de prever as saídas corretas para novos dados não vistos.

O processo de aprendizado supervisionado envolve os seguintes elementos:

- **Dados de treinamento:** são exemplos de entrada e saída rotulados que são fornecidos ao algoritmo de aprendizado. Os dados de treinamento são usados para estimar os parâmetros do modelo.
- **Características (features):** são as variáveis ou atributos que descrevem as entradas dos dados de treinamento. As características são usadas como informações de entrada para o modelo de aprendizado.

- **Saídas desejadas:** são os rótulos ou valores corretos associados às entradas nos dados de treinamento. Essas saídas são fornecidas ao algoritmo de aprendizado para que ele possa aprender a relação entre as entradas e as saídas corretas.
- **Modelo de aprendizado:** é a representação matemática ou estatística que o algoritmo de aprendizado cria a partir dos dados de treinamento. O modelo pode ser um conjunto de regras, uma árvore de decisão, uma rede neural, uma função matemática ou qualquer outra forma que possa mapear as entradas para as saídas.
- **Algoritmo de aprendizado:** é o algoritmo que utiliza os dados de treinamento para ajustar os parâmetros do modelo de aprendizado. O algoritmo realiza uma busca no espaço de hipóteses para encontrar o modelo que melhor se ajusta aos dados de treinamento.
- **Avaliação do modelo:** após o treinamento, é importante avaliar o desempenho do modelo em dados não vistos. Isso é feito usando um conjunto de dados de teste separado, que não foi usado durante o treinamento. A avaliação permite verificar se o modelo é capaz de generalizar bem e fazer previsões precisas em novos dados.

O aprendizado supervisionado é aplicado em uma variedade de problemas, como classificação, regressão e previsão. Na classificação, o objetivo é atribuir rótulos ou categorias às entradas. Exemplos incluem identificação de spam, reconhecimento de imagens e detecção de fraudes. Na regressão, o objetivo é prever um valor numérico contínuo com base nas entradas. Exemplos incluem previsão de preços de imóveis, estimativa de vendas e análise de tendências.

O aprendizado supervisionado é uma abordagem poderosa para o desenvolvimento de modelos preditivos e oferece uma ampla gama de algoritmos e técnicas para diferentes tipos de problemas. É importante escolher cuidadosamente as características relevantes, o algoritmo de aprendizado apropriado e avaliar adequadamente o desempenho do modelo para garantir resultados precisos e confiáveis.

3.2 Classificação

A classificação é uma técnica muito importante no campo do aprendizado de máquina. Seu objetivo é atribuir rótulos ou categorias a diferentes instâncias de dados com base em suas características. É uma forma de aprendizado supervisionado em que o algoritmo de aprendizado é treinado usando um conjunto de dados rotulados, em que as classes ou categorias já são conhecidas.

Começamos o processo de classificação com o conjunto de dados, que é um conjunto de exemplos de treinamento, no qual cada exemplo é descrito por um conjunto de características (features) e possui um rótulo ou categoria associada. Os exemplos de treinamento são usados para treinar o modelo de classificação.

As características são as variáveis que descrevem as instâncias de dados e são usadas como entrada para o modelo de classificação. As características podem ser numéricas, categóricas ou binárias, dependendo do tipo de dados. Os rótulos são as classes ou categorias às quais as instâncias de dados

pertencem. Por exemplo: diagnósticos de pessoas (saudáveis ou doentes), perfil de pagamento das pessoas (bom ou mau) ou qual o tipo de animal (cachorro, gato ou rato).

O modelo de classificação é uma representação matemática ou estatística que o algoritmo de aprendizado cria a partir dos dados de treinamento. O modelo é construído com base nas características dos dados e é usado para fazer previsões sobre a classe de novos exemplos. O algoritmo de classificação aprende a relação entre as características dos dados de treinamento e seus rótulos correspondentes. Existem diferentes algoritmos de classificação disponíveis, como árvores de decisão, máquinas de vetores de suporte (SVM), k-vizinhos mais próximos (KNN), redes neurais, entre outros. Cada algoritmo tem suas próprias características, vantagens e limitações.

Após o treinamento é importante avaliar o desempenho do modelo em dados não vistos. Isso é feito usando um conjunto de dados de teste separado, que não foi usado durante o treinamento. A avaliação permite medir a precisão do modelo na classificação correta das instâncias de dados e identificar possíveis problemas, como overfitting ou underfitting.

A figura a seguir representa um exemplo de classificação. De um lado estão classificadas bolinhas azuis; do outro, as bolinhas vermelhas, aí separamos os dois objetos, cada um com sua classificação. Por exemplo, bolinhas vermelhas poderiam ser pessoas doentes e bolinhas azuis, pessoas saudáveis.

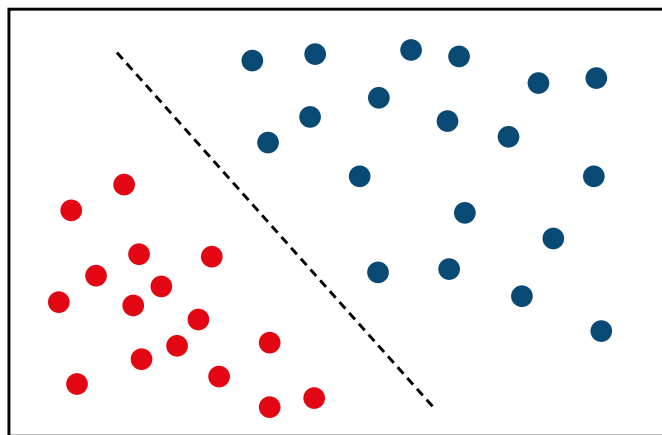


Figura 10 – Exemplo de classificação

Podemos utilizar a classificação em aplicações de várias áreas, como medicina, finanças, marketing, processamento de linguagem natural, detecção de fraudes, entre outras. Por exemplo, um modelo de classificação pode ser usado para prever se um cliente fará uma compra ou não, classificar e-mails como spam ou não spam, identificar doenças com base em sintomas ou classificar imagens em diferentes categorias.

Para obter bons resultados na classificação, é importante selecionar cuidadosamente as características relevantes, escolher o algoritmo de classificação adequado, ajustar os hiperparâmetros do modelo e realizar uma avaliação rigorosa. Além disso, devemos ter um conjunto de dados representativo e suficientemente grande para treinar o modelo de classificação e evitar vieses ou problemas de generalização.

A classificação é uma técnica fundamental no aprendizado de máquina e é utilizada para resolver problemas de tomada de decisão. Através da identificação de padrões nos dados, os modelos de classificação podem automatizar processos de categorização, filtragem, previsão e suporte à decisão, trazendo eficiência e precisão para diversas aplicações. No entanto, é importante ressaltar que a escolha adequada das características, a qualidade dos dados de treinamento e a seleção do algoritmo de classificação são fatores cruciais para o desempenho e a acurácia do modelo.

Existem diferentes tipos de algoritmos de classificação, cada um com suas próprias suposições e métodos de classificação. Podemos citar alguns algoritmos como exemplo: árvore de decisão, máquinas de vetores de suporte (SVM), k-vizinhos mais próximos (KNN) e redes neurais.

- **Árvores de decisão:** o algoritmo cria uma estrutura em forma de árvore que representa decisões e condições baseadas nas características dos dados. Cada nó da árvore representa uma característica e cada ramo representa uma decisão. As folhas da árvore correspondem às classes ou categorias finais.
- **Máquinas de vetores de suporte (SVM):** esse algoritmo busca encontrar um hiperplano de separação que maximize a margem entre as classes. Ele mapeia as características de entrada em um espaço de dimensões superiores e realiza a classificação com base na posição dos exemplos neste espaço.
- **k-vizinhos mais próximos (KNN):** esse algoritmo classifica as instâncias de acordo com a classe das instâncias vizinhas mais próximas. A distância entre as instâncias é calculada com base nas características, e os k-vizinhos mais próximos são considerados para determinar a classe do exemplo em questão.
- **Redes neurais:** são modelos inspirados no funcionamento do cérebro humano, compostos de camadas de neurônios interconectados. Cada neurônio processa as características de entrada e a rede aprende os pesos das conexões para realizar a classificação.

Esses são apenas alguns exemplos de algoritmos de classificação, e existem muitos outros disponíveis, cada um com suas próprias características e adequado para diferentes tipos de problemas. A escolha do algoritmo correto depende das características dos dados, do tamanho do conjunto de dados, da dimensionalidade, da interpretabilidade e de outros fatores relevantes.

Vamos exemplificar a classificação utilizando SVM (support vector machine) em Python. Isso pode ser realizado usando a biblioteca scikit-learn, que é uma das mais populares para aprendizado de máquina.



Saiba mais

Para saber sobre a biblioteca scikit-learn, acesse:

<https://scikit-learn.org/>. Acesso em: 21 ago. 2023.

Vamos mostrar um exemplo simples de como usar SVM para classificar dados fictícios em duas classes. Certifique-se de ter o scikit-learn instalado em seu ambiente antes de executar este código. Você pode instalar com o comando `pip install scikit-learn`.

```
1. import numpy as np
2. import matplotlib.pyplot as plt
3. from sklearn import datasets
4. from sklearn.model_selection import train_test_split
5. from sklearn import svm
6. from sklearn.metrics import accuracy_score

7. # Gerar dados fictícios para classificação
8. X, y = datasets.make_classification(n_samples=100, n_features=2, n_classes=2, n_clusters_per_
   class=1, n_redundant=0, random_state=42)

9. # Dividir os dados em conjunto de treinamento e teste
10. X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

11. # Criar um classificador SVM
12. clf = svm.SVC(kernel='linear') # SVM com kernel linear

13. # Treinar o modelo
14. clf.fit(X_train, y_train)

15. # Fazer previsões no conjunto de teste
16. y_pred = clf.predict(X_test)

17. # Calcular a acurácia das previsões
18. accuracy = accuracy_score(y_test, y_pred)
19. print("Acurácia: {:.2f}%".format(accuracy * 100))

20. # Plotar os dados e o hiperplano de separação
21. plt.scatter(X[:, 0], X[:, 1], c=y, cmap=plt.cm.Paired)
22. ax = plt.gca()
23. xlim = ax.get_xlim()
24. ylim = ax.get_ylim()
```

```
25. # Criar grid para plotar o hiperplano
26. xx, yy = np.meshgrid(np.linspace(xlim[0], xlim[1], 50), np.linspace(ylim[0], ylim[1], 50))
27. Z = clf.decision_function(np.c_[xx.ravel(), yy.ravel()])
28. Z = Z.reshape(xx.shape)

29. # Plotar o hiperplano e margens
30. plt.contour(xx, yy, Z, colors='k', levels=[-1, 0, 1], alpha=0.5, linestyles=['--', '-', '--'])
31. plt.scatter(clf.support_vectors_[0], clf.support_vectors_[1], s=100, linewidth=1,
               facecolors='none', edgecolors='k')
32. plt.xlabel('Feature 1')
33. plt.ylabel('Feature 2')
34. plt.title('SVM com Kernel Linear')
35. plt.show()
```

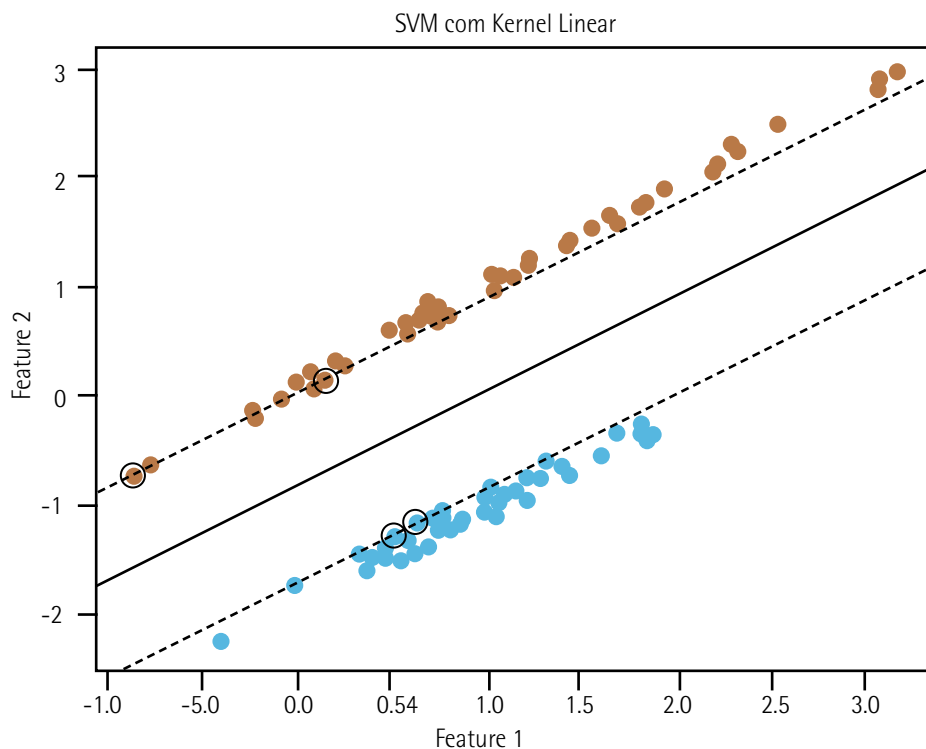


Figura 11

Neste exemplo, estamos gerando dados fictícios usando `make_classification` da biblioteca `scikit-learn`, dividindo os dados em conjuntos de treinamento e teste, criando um classificador SVM com kernel linear, treinando-o, fazendo previsões e calculando a acurácia. Também estamos plotando os dados e o hiperplano de separação no gráfico.

Lembre-se de que este é um exemplo básico para entender como usar SVM para classificação em Python. Em cenários reais, você deve adaptar o código para seus próprios dados e explorar diferentes configurações de SVM (como diferentes kernels e parâmetros de regularização) para obter melhores resultados.

A classificação é uma técnica importante no campo do aprendizado de máquina que permite automatizar a categorização de dados com base em características relevantes. Com a capacidade de prever a classe ou categoria de novos exemplos, os modelos de classificação são utilizados em diversas áreas para melhorar a tomada de decisões, a eficiência operacional e o entendimento dos dados.

3.3 Regressão

A regressão é uma técnica estatística utilizada no campo do aprendizado de máquina para modelar e prever relações entre variáveis. Ao contrário da classificação, em que o objetivo é atribuir rótulos ou categorias a instâncias de dados, a regressão busca prever um valor numérico contínuo com base em um conjunto de variáveis independentes.

O objetivo da regressão é encontrar uma função matemática ou estatística que relacione as variáveis independentes (também chamadas de características ou variáveis de entrada) a uma variável dependente (também conhecida como variável de saída ou variável-alvo). Essa função é chamada de modelo de regressão e é utilizada para fazer previsões sobre o valor da variável dependente para novos exemplos de dados.

Existem vários tipos de regressão, cada um adequado para diferentes tipos de problemas e dados. Alguns dos principais tipos de regressão incluem:

- **Regressão linear:** é um tipo de regressão que assume uma relação linear entre as variáveis independentes e dependentes. O modelo de regressão linear encontra a melhor linha reta que representa essa relação, minimizando a soma dos quadrados dos erros entre os valores reais e os valores previstos.
- **Regressão logística:** é usada quando a variável dependente é binária, ou seja, possui apenas duas classes. O modelo de regressão logística utiliza uma função logística para estimar a probabilidade de um exemplo pertencer a uma das classes.
- **Regressão polinomial:** é uma extensão da regressão linear, em que a relação entre as variáveis é modelada usando um polinômio. Isso permite capturar relações não lineares entre as variáveis e aumentar a flexibilidade do modelo.
- **Regressão de séries temporais:** é usada para prever valores futuros com base em padrões temporais nos dados. A regressão de séries temporais leva em consideração a dependência temporal dos dados e utiliza métodos como médias móveis, Arima (AutoRegressive Integrated Moving Average) e modelos baseados em suavização exponencial.
- **Regressão Ridge e Lasso:** são técnicas de regressão que ajudam a lidar com problemas de multicolinearidade, em que as variáveis independentes estão altamente correlacionadas entre si. Essas técnicas adicionam termos de regularização ao modelo de regressão para controlar a complexidade e evitar overfitting.

Avaliar a qualidade do modelo de regressão é essencial para garantir sua eficácia. Métricas comuns de avaliação incluem o erro quadrático médio (mean square error – MSE), o coeficiente de determinação (R^2) e o erro absoluto médio (mean absolute error – MAE). Além disso, é importante realizar validação cruzada e dividir os dados em conjuntos de treinamento, validação e teste para garantir que o modelo seja capaz de generalizar bem para novos dados.

A regressão tem amplas aplicações em diversas áreas, como previsão de vendas, análise financeira, análise de risco, previsão de preços, modelagem de demanda, entre outras. Ela permite extrair insights valiosos dos dados e tomar decisões informadas com base nas relações encontradas.

Vamos utilizar a Python para exemplificar. Consideramos um conjunto de dados com duas variáveis: o número de horas estudadas (X) e o resultado de um teste (Y). Vamos usar a regressão linear para tentar prever os resultados dos testes com base nas horas estudadas. Vamos começar importando as bibliotecas necessárias e criando alguns dados fictícios:

```
1. import numpy as np
2. import matplotlib.pyplot as plt

3. # Dados fictícios de horas estudadas e resultados do teste
4. horas_estudadas = np.array([1, 2, 3, 4, 5, 6, 7, 8, 9, 10])
5. resultados_teste = np.array([35, 47, 60, 72, 80, 92, 100, 112, 120, 130])

6. # Calculando a média das horas estudadas e dos resultados do teste
7. media_horas = np.mean(horas_estudadas)
8. media_resultados = np.mean(resultados_teste)

9. # Calculando os coeficientes da regressão
10. numerator = np.sum((horas_estudadas - media_horas) * (resultados_teste - media_resultados))
11. denominator = np.sum((horas_estudadas - media_horas) ** 2)
12. slope = numerator / denominator
13. intercept = media_resultados - slope * media_horas

14. # Função de regressão linear
15. def regressao_linear(x):
16.     return slope * x + intercept

17. # Fazendo a previsão para algumas horas estudadas
18. horas_para_previsao = np.array([11, 12, 13, 14])
19. previsao_resultados = regressao_linear(horas_para_previsao)
```

20. # Plotando os dados e a linha de regressão
21. `plt.scatter(horas_estudadas, resultados_teste, label="Dados reais")`
22. `plt.plot(horas_estudadas, regressao_linear(horas_estudadas), color='red', label="Regressão linear")`
23. `plt.scatter(horas_para_previsao, previsao_resultados, color='green', label="Previsão")`
24. `plt.xlabel("Horas Estudadas")`
25. `plt.ylabel("Resultados do Teste")`
26. `plt.legend()`
27. `plt.show()`

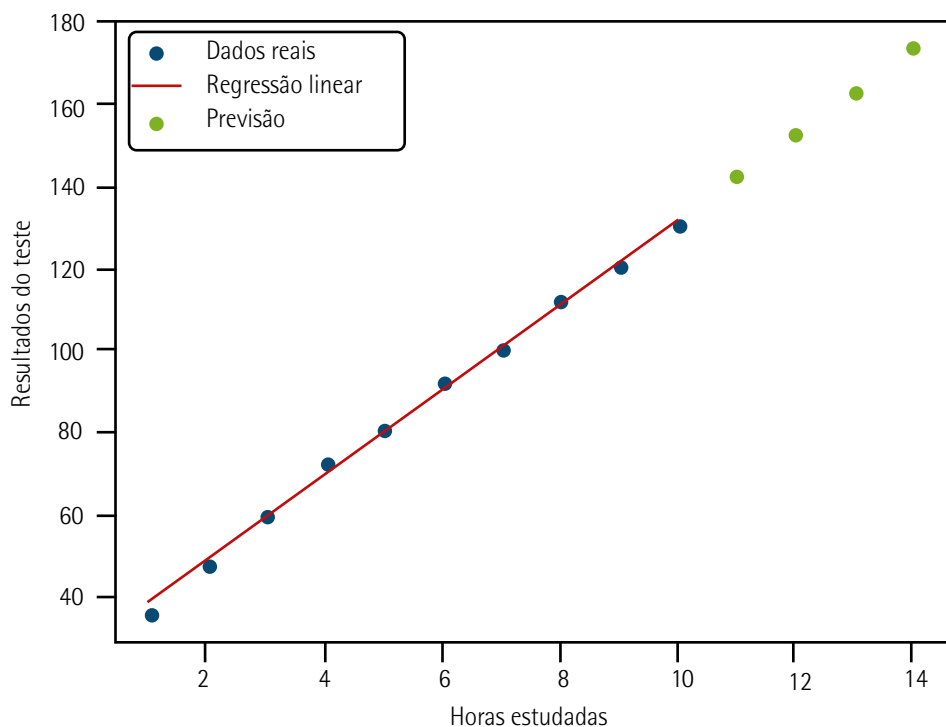


Figura 12

Este código cria um gráfico com os pontos representando os dados reais, a linha de regressão linear e alguns pontos verdes que representam as previsões para novas horas estudadas. A regressão linear é usada para encontrar a linha que melhor se ajusta aos dados, minimizando a soma dos erros quadráticos.

Lembrando que este é um exemplo simples e a qualidade da previsão pode variar dependendo dos dados reais e do modelo estatístico utilizado.

3.4 Aprendizado não supervisionado

O aprendizado não supervisionado é uma das principais abordagens de aprendizado de máquina, em que o objetivo é extrair informações úteis a partir de dados não rotulados. Diferentemente do aprendizado supervisionado, em que os dados são rotulados, o aprendizado não supervisionado busca identificar padrões e estruturas nos dados sem uma definição prévia das classes.

O aprendizado não supervisionado é utilizado em várias aplicações, como análise de dados de clientes, agrupamento de dados em bases de dados, detecção de anomalias, identificação de padrões em imagens, reconhecimento de voz, entre outras. O sucesso do aprendizado não supervisionado depende de uma escolha cuidadosa das técnicas e algoritmos utilizados, bem como da qualidade dos dados de entrada.

3.5 Agrupamento

O agrupamento, também conhecido como clustering, é uma técnica de aprendizado não supervisionado que busca identificar grupos ou clusters de objetos similares em um conjunto de dados. O objetivo do agrupamento é encontrar estruturas e padrões nos dados sem a necessidade de rótulos ou categorias predefinidas.

O processo de agrupamento envolve a divisão dos dados em grupos de tal forma que objetos dentro do mesmo grupo sejam mais semelhantes entre si do que com objetos de outros grupos. A semelhança é geralmente medida com base nas características ou atributos dos objetos. Existem várias abordagens e algoritmos de agrupamento, cada um com suas próprias suposições e critérios de similaridade.

Alguns dos algoritmos de agrupamento mais comuns incluem:

- **K-means:** é um algoritmo de particionamento que divide os dados em K grupos, onde K é um valor predefinido. Ele inicializa os centroides dos grupos de forma aleatória e, em seguida, itera alternando entre atribuir objetos ao grupo mais próximo e atualizar os centroides com base nos objetos atribuídos.
- **Hierárquico:** é uma abordagem que constrói uma estrutura hierárquica de clusters. Existem dois tipos principais: aglomerativo, em que cada objeto começa como um cluster e os clusters são combinados de forma iterativa, e divisivo, no qual todos os objetos começam em um único cluster e são divididos em subclusters.
- **DBSCAN:** é um algoritmo baseado em densidade que agrupa os objetos com base na densidade local. Ele identifica regiões densas de objetos conectados e atribui esses grupos como clusters, enquanto objetos isolados são considerados ruídos.
- **Mean shift:** é um algoritmo que busca iterativamente o centro de massa dos pontos em uma vizinhança definida por uma janela de busca. Ele move os pontos em direção aos centros de massa até atingir uma convergência, formando assim os clusters.

O agrupamento tem uma ampla gama de aplicações em diferentes áreas, como segmentação de clientes, análise de mercado, agrupamento de documentos, análise de redes sociais, reconhecimento de padrões em imagens, entre outros. É importante escolher o algoritmo de agrupamento adequado com base na natureza dos dados, na estrutura desejada e nos requisitos específicos do problema.

A avaliação dos resultados do agrupamento também é um desafio importante. Existem várias métricas, como índice de Silhueta e índice de Davies-Bouldin, que podem ser usadas para avaliar a qualidade e a coerência dos clusters obtidos.

Vamos exemplificar o agrupamento utilizando o algoritmo K-means, que é um método de clustering popular para agrupar dados não rotulados. Aqui está um exemplo simples de como usar o K-means em Python usando a biblioteca scikit-learn:

```
1. import numpy as np
2. import matplotlib.pyplot as plt
3. from sklearn.cluster import KMeans
4. from sklearn.datasets import make_blobs

5. # Gerar dados fictícios
6. X, _ = make_blobs(n_samples=300, centers=4, random_state=42)

7. # Criar um objeto KMeans com 4 clusters
8. kmeans = KMeans(n_clusters=4)

9. # Treinar o modelo KMeans
10. kmeans.fit(X)

11. # Obter os centros dos clusters e os rótulos para cada ponto
12. centroids = kmeans.cluster_centers_
13. labels = kmeans.labels_

14. # Plotar os dados e os centros dos clusters
15. plt.scatter(X[:, 0], X[:, 1], c=labels, cmap='viridis')
16. plt.scatter(centroids[:, 0], centroids[:, 1], marker='x', s=200, linewidths=3, color='r')

17. plt.xlabel("Feature 1")
18. plt.ylabel("Feature 2")
19. plt.title("K-Means Clustering")
20. plt.show()
```

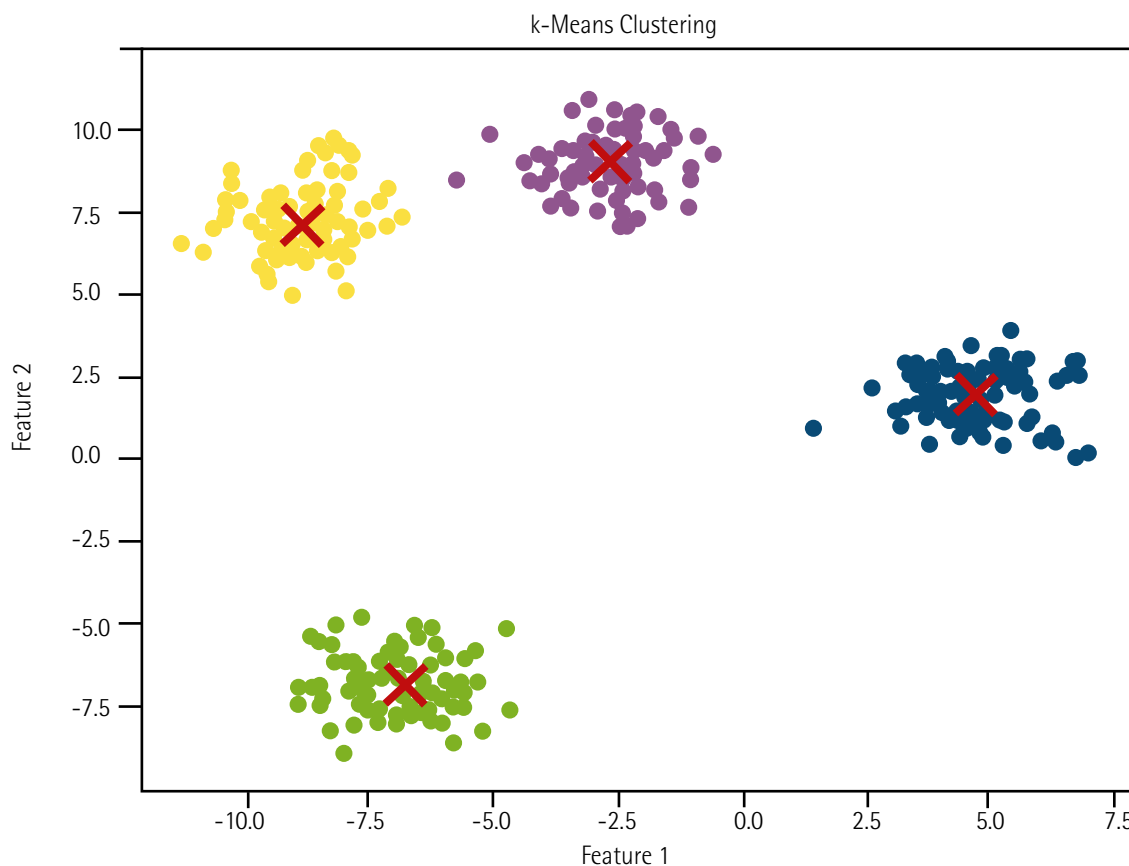


Figura 13

Neste exemplo, estamos gerando dados fictícios usando `make_blobs` da biblioteca `scikit-learn`, criando um objeto K-means com quatro clusters, treinando o modelo, obtendo os centros dos clusters e os rótulos para cada ponto e, finalmente, plotando os dados e os centros dos clusters no gráfico.

Este é um exemplo básico para entender como usar o K-means. Em aplicações reais, você pode precisar ajustar o número de clusters (`n_clusters`) com base em algum critério (por exemplo, usando o método Elbow) e também explorar outros parâmetros para obter resultados adequados ao seu conjunto de dados.

O agrupamento é uma técnica de aprendizado não supervisionado que permite descobrir padrões e estruturas nos dados, agrupando objetos similares em clusters. É uma ferramenta valiosa para explorar grandes conjuntos de dados e encontrar insights úteis sem a necessidade de rótulos prévios.

3.6 Associação

Associação é uma técnica de aprendizado não supervisionado usada para descobrir padrões interessantes ou relações entre itens em um conjunto de dados. Ela é frequentemente utilizada em aplicações como análise de cesta de compras, análise de comportamento de usuários e análise de interações em redes sociais.

A técnica de associação envolve encontrar conjuntos de itens que ocorrem juntos com frequência em um conjunto de dados. Esses conjuntos são chamados de **itemsets frequentes**. Para encontrar esses itemsets frequentes, algoritmos de associação como o Apriori são utilizados. O algoritmo Apriori funciona através da geração iterativa de itemsets candidatos de tamanhos crescentes, seguido pela verificação de sua frequência no conjunto de dados.

Além de encontrar os itemsets frequentes, a técnica de associação pode ser usada para gerar regras de associação a partir desses itemsets. As regras de associação são declarações que indicam uma relação entre dois ou mais itens. Por exemplo, uma regra de associação pode ser: "se um cliente comprar pão e leite, então é provável que ele também compre manteiga".

As regras de associação são avaliadas com base em métricas como suporte, confiança e lift. O suporte mede a frequência com que um itemset aparece no conjunto de dados. A confiança mede a proporção de vezes em que a regra é verdadeira em relação ao número total de vezes que a condição antecedente ocorre. O lift é uma métrica que mede a força da relação entre os itens, levando em conta a frequência do itemset e a frequência dos itens individualmente.

A técnica de associação pode ser útil em muitas aplicações, como recomendações de produtos, análise de comportamento de usuários em sites de comércio eletrônico, análise de redes sociais e muito mais. No entanto, é importante notar que o uso de associação pode ser desafiador em conjuntos de dados com muitos itens, pois o espaço de busca pode se tornar muito grande e os algoritmos de associação podem levar muito tempo para serem executados.

Vamos exemplificar a técnica de associação de uma forma bem simples. Vamos considerar um cenário hipotético de uma loja de produtos eletrônicos e tentar identificar algumas regras de associação com base nas compras dos clientes.

Suponha que temos os seguintes dados de transações, onde cada linha representa uma compra de um cliente:

- **Compra:** smartphone, fone de ouvido, carregador.
- **Compra:** smartphone, capa protetora.
- **Compra:** smartwatch, fone de ouvido.
- **Compra:** smartphone, carregador, capa protetora.
- **Compra:** fone de ouvido, carregador.

Agora, podemos usar esses dados para identificar algumas regras de associação:

- **Regra 1:** se um cliente comprar um smartphone, é provável que ele também compre um carregador.
- **Regra 2:** se um cliente comprar um fone de ouvido, é provável que ele também compre um carregador.
- **Regra 3:** se um cliente comprar um smartphone e um carregador, é provável que ele também compre uma capa protetora.
- **Regra 4:** se um cliente comprar um smartwatch, é provável que ele também compre um fone de ouvido.

Essas regras de associação podem ajudar a loja a entender os padrões de compra dos clientes e otimizar a disposição dos produtos na loja ou criar promoções direcionadas. Esses exemplos são hipotéticos e criados apenas para ilustrar o conceito de regras de associação sem uso de código. Em cenários reais, você usaria algoritmos como o Apriori ou FP-Growth para descobrir automaticamente essas regras a partir dos dados.

3.7 Sumarização

A sumarização é uma técnica que visa resumir informações em um conjunto de dados de forma concisa, mas informativa. Ela é usada para extrair as principais ideias, características ou padrões de um conjunto de dados extenso, reduzindo-o para um resumo mais compacto e fácil de entender.

Existem dois tipos principais de sumarização: a **sumarização automática** e a **sumarização manual**.

A **sumarização automática** é realizada por algoritmos e técnicas de processamento de linguagem natural (PLN) que analisam o texto ou os dados brutos e extraem informações relevantes para gerar o resumo. Existem várias abordagens para a sumarização automática, incluindo:

- **Sumarização extrativa:** nesta abordagem, as frases ou trechos mais importantes do texto são identificados e selecionados para formar o resumo. Geralmente, são considerados critérios como relevância, importância e coerência para determinar quais partes do texto devem ser incluídas no resumo.
- **Sumarização abstrativa:** nesta abordagem, o sistema de sumarização gera frases sinteticamente que capturam o significado do texto original, em vez de simplesmente extrair frases do texto original. Isso envolve a compreensão do texto, a interpretação do significado e a geração de frases com base nesse entendimento.

Já a **sumarização manual** é feita por seres humanos, que leem, analisam e selecionam as informações mais importantes do conjunto de dados para criar um resumo. Isso é comum em áreas como jornalismo,

em que os profissionais resumem e sintetizam informações de várias fontes para criar notícias ou artigos resumidos.

A sumarização tem diversas aplicações, como:

- Sumarização de textos longos, como artigos de notícias, documentos acadêmicos e relatórios, para fornecer uma visão geral rápida do conteúdo.
- Sumarização de grandes conjuntos de dados para identificar os principais padrões ou insights.
- Sumarização de conversas ou discussões em fóruns ou redes sociais para extrair as principais ideias ou tópicos discutidos.
- Sumarização de informações de mercado, como análise de tendências, estatísticas ou resultados de pesquisa.

A sumarização eficaz requer a compreensão do contexto, a identificação das informações mais relevantes e a habilidade de comunicar essas informações de forma clara e concisa. Tanto a sumarização automática quanto a sumarização manual desempenham papéis importantes em diferentes contextos, e a escolha entre elas depende do objetivo, da natureza dos dados e das necessidades do usuário.

4 MINERAÇÃO DE DADOS

A mineração de dados, também conhecida como descoberta de conhecimento em bancos de dados (KDD – knowledge discovery in databases), surgiu como uma resposta à crescente quantidade de dados disponíveis nas empresas e instituições. O termo "mineração de dados" foi cunhado no final da década de 1980 e o campo evoluiu a partir dos campos da estatística, inteligência artificial e banco de dados.

O rápido avanço da tecnologia da informação resultou na geração de grandes volumes de dados em diversos setores, como finanças, saúde, varejo, telecomunicações, entre outros. No entanto, muitas vezes, esses dados eram armazenados sem serem utilizados de maneira eficiente, devido à falta de ferramentas e técnicas para extrair insights e conhecimentos úteis deles.

A mineração de dados surgiu como uma abordagem para lidar com esse desafio. Ela visa descobrir padrões, tendências, relações e conhecimentos ocultos nos dados, permitindo que as organizações tomem decisões mais informadas e estratégicas. A ideia central da mineração de dados é explorar os dados em busca de informações valiosas e previamente desconhecidas.

Inicialmente, a mineração de dados concentrou-se principalmente em técnicas estatísticas, como análise de regressão, análise de agrupamento e análise de correlação. No entanto, à medida que o campo evoluiu, foram desenvolvidas e aplicadas várias técnicas e algoritmos mais avançados, como árvores de decisão, redes neurais, algoritmos genéticos, máquinas de vetores de suporte e algoritmos de aprendizado de máquina.

A mineração de dados é um processo iterativo que envolve várias etapas, incluindo seleção de dados, pré-processamento, transformação, mineração, interpretação e avaliação dos resultados. Essas etapas visam extrair conhecimento útil e significativo dos dados, que pode ser usado para tomar decisões melhores, identificar oportunidades de negócios, otimizar processos, detectar fraudes, entre outras aplicações.

Com o avanço da tecnologia e o surgimento de grandes quantidades de dados estruturados e não estruturados, como dados de redes sociais, textos, imagens e vídeos, a mineração de dados continua evoluindo. Novas técnicas e algoritmos são desenvolvidos para lidar com esses desafios e explorar o potencial dos dados em diversas áreas, como medicina, marketing, ciências sociais, segurança, entre outras.

A mineração de dados desempenha um papel importante na Era da Informação, permitindo que as organizações aproveitem o poder dos dados para tomar decisões estratégicas, obter vantagem competitiva e impulsionar a inovação. Ela continua a evoluir à medida que surgem novas fontes de dados e desafios, tornando-se uma disciplina interdisciplinar e essencial para a compreensão e o uso eficaz dos dados em nosso mundo moderno.

4.1 Visão geral

A mineração de dados é uma área de estudo que visa extrair informações valiosas e significativas a partir de grandes conjuntos de dados. É um processo iterativo e interdisciplinar que envolve a aplicação de técnicas e algoritmos para descobrir padrões, tendências, relações e conhecimentos ocultos nos dados. Para *Morais et al.* (2018, p. 82):

Dentre as características mais importantes da mineração de dados, está o grande volume de dados e a capacidade de mudança de escala com relação ao tamanho dos dados. Algoritmos têm a capacidade de mudança de escala, mas a mineração é muito mais do que aplicar algoritmos, pois, geralmente, os dados contêm ruído ou estão incompletos, sendo provável que padrões sejam perdidos, de modo que a confiabilidade será baixa. Logo, o analista precisa tomar a decisão sobre quais tipos de algoritmos de mineração serão necessários, aplicando-os em um conjunto de amostra de dados específico, sintetizando os resultados, aplicando ferramentas de apoio à decisão e mineração, iterando o processo.

A mineração de dados parte do pressuposto de que os dados contêm informações valiosas que podem ser usadas para tomar decisões mais informadas, identificar oportunidades de negócios, resolver problemas complexos, fazer previsões e melhorar processos. Através da aplicação de técnicas estatísticas, algoritmos de aprendizado de máquina, análise de dados e outras abordagens, a mineração de dados busca transformar os dados em conhecimento acionável.

O processo de mineração de dados geralmente envolve as seguintes etapas:

1. **Seleção de dados:** envolve a identificação dos dados relevantes para o problema em questão. Esta etapa inclui a definição dos critérios de seleção e a obtenção dos dados necessários.

2. **Pré-processamento:** é a fase em que os dados brutos são limpos, organizados e preparados para análise. Isso pode envolver a remoção de dados ausentes, correção de erros, normalização e transformação dos dados.

3. **Transformação:** nesta etapa, os dados são convertidos em uma forma adequada para análise. Isso pode incluir a redução de dimensionalidade, a extração de características relevantes e a aplicação de técnicas estatísticas ou algoritmos de processamento de dados.

4. **Mineração de dados:** é a fase central do processo, em que são aplicados algoritmos e técnicas de mineração de dados para descobrir padrões, relações e conhecimentos nos dados. Isso pode envolver a aplicação de técnicas de aprendizado de máquina, análise estatística, visualização de dados e outras abordagens.

5. **Avaliação e interpretação:** os resultados da mineração de dados são avaliados quanto à sua relevância, qualidade e utilidade. Os padrões e conhecimentos descobertos são interpretados para extrair informações significativas e compreender seu impacto no problema em questão.

6. **Utilização do conhecimento:** os resultados e insights obtidos são utilizados para tomar decisões informadas, desenvolver estratégias, resolver problemas e gerar valor para a organização ou área de estudo.

A mineração de dados tem uma ampla gama de aplicações em diferentes setores, como finanças, saúde, varejo, marketing, telecomunicações, ciências sociais e muitos outros. Ela é utilizada para análise de mercado, detecção de fraudes, previsão de demanda, recomendação de produtos, segmentação de clientes, diagnóstico médico, análise de sentimentos, entre outras finalidades.

No entanto, a mineração de dados também apresenta desafios, como a qualidade dos dados, a privacidade e segurança, o gerenciamento de grandes volumes de dados e a interpretação correta dos resultados. Portanto, é importante aplicar técnicas adequadas, considerar questões éticas e ter uma compreensão sólida dos dados e do contexto em que a mineração de dados está sendo aplicada.

Com o avanço da tecnologia e o surgimento de técnicas mais sofisticadas, como aprendizado de máquina e inteligência artificial, a mineração de dados continua evoluindo. Novos algoritmos e abordagens estão sendo desenvolvidos para lidar com conjuntos de dados cada vez maiores, mais complexos e variados. Isso inclui técnicas de aprendizado profundo (deep learning), processamento de linguagem natural (NLP), análise de redes sociais, entre outras.

Além disso, a mineração de dados está se beneficiando do uso de ferramentas e plataformas específicas, como linguagens de programação para análise de dados (R, Python), bibliotecas de aprendizado de máquina (scikit-learn, TensorFlow, Keras), ferramentas de visualização de dados (Tableau, Power BI) e plataformas de Big Data (Hadoop, Spark).

A aplicação efetiva da mineração de dados requer uma combinação de conhecimentos em estatística, aprendizado de máquina, programação, matemática e conhecimento de domínio específico. Profissionais qualificados em mineração de dados são capazes de identificar padrões relevantes, criar modelos preditivos e interpretar os resultados de forma adequada. A figura a seguir apresenta a multidisciplinaridade da mineração de dados:

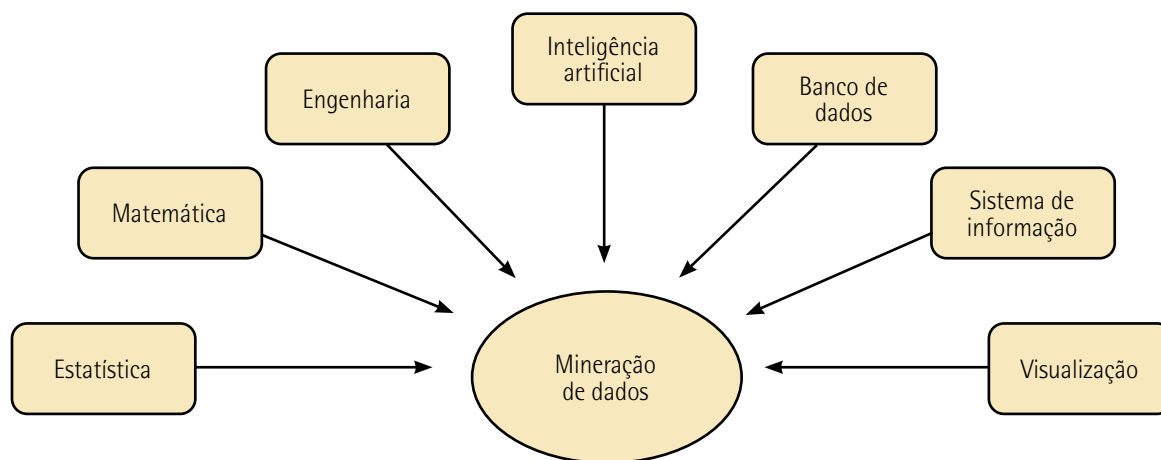


Figura 14 – Multidisciplinaridade da mineração de dados

Adaptada de: Castro e Ferrari (2016, p. 49).

Embora a mineração de dados ofereça grandes oportunidades para melhorar a tomada de decisões e gerar insights valiosos, é importante também considerar os desafios e limitações. Isso inclui a necessidade de garantir a qualidade dos dados, lidar com a dimensionalidade e complexidade dos conjuntos de dados, evitar o viés nos resultados e garantir a privacidade e segurança das informações.

A mineração de dados é uma disciplina multidisciplinar que visa extrair conhecimento valioso e útil a partir de grandes volumes de dados. Ela oferece inúmeras oportunidades para melhorar a eficiência, a precisão e a compreensão em diversas áreas. No entanto, seu sucesso depende do uso adequado de técnicas e algoritmos, bem como da compreensão correta do contexto e dos desafios específicos de cada problema de mineração de dados.

4.2 Modelos de ML

Os modelos de aprendizado de máquina (machine learning – ML) são representações matemáticas ou estatísticas que capturam os padrões e relações nos dados e são usados para fazer previsões, classificações ou tomar decisões com base nesses padrões. Eles são a base do processo de aprendizado de máquina, permitindo que os algoritmos aprendam a partir dos dados e gerem resultados.

Existem diversos tipos de modelos de ML, cada um adequado para diferentes tarefas e tipos de dados. Vamos explorar alguns dos modelos mais comuns:

- **Regressão linear:** é um modelo que estabelece uma relação linear entre variáveis de entrada e uma variável de saída contínua. É usado para fazer previsões numéricas, estimando valores em uma escala contínua.
- **Regressão logística:** é um modelo usado para problemas de classificação binária, em que a variável de saída é categórica com duas classes. Ele estima a probabilidade de um evento ocorrer, mapeando a entrada para uma função logística.
- **Árvores de decisão:** são modelos que dividem o conjunto de dados em subconjuntos menores, criando uma árvore de regras de decisão. Cada nó da árvore representa uma característica do dado, e as folhas são as classes ou valores de saída. As árvores de decisão são populares devido à sua interpretabilidade.
- **Random forest:** é um conjunto de árvores de decisão que trabalham em conjunto para realizar classificação ou regressão. Cada árvore é treinada em um subconjunto aleatório dos dados e os resultados são combinados para obter uma predição final.
- **Support vector machines (SVM):** são modelos que mapeiam os dados em um espaço dimensional superior, em que as classes podem ser separadas por um hiperplano. O SVM busca encontrar o hiperplano de melhor separação entre as classes.
- **Redes neurais artificiais:** são modelos inspirados no funcionamento do cérebro humano. Consistem em camadas de neurônios interconectados que processam os dados de entrada para gerar uma saída. Redes neurais profundas, como as redes neurais convolucionais (CNNs) e redes neurais recorrentes (RNNs), são comumente usadas para tarefas de processamento de imagem e sequência, respectivamente.
- **Máquinas de vetores de suporte (SVM):** são modelos que buscam encontrar o hiperplano que melhor separa os dados em um espaço dimensional superior, em que a separação é linear. Podem ser estendidos para realizar classificação não linear usando truques de kernel.
- **Naive Bayes:** é um modelo baseado no Teorema de Bayes e na suposição de independência condicional entre os recursos. É frequentemente usado para classificação de texto e análise de sentimentos.
- **Algoritmos de agrupamento (clustering):** são modelos usados para identificar grupos ou clusters nos dados, com base em sua similaridade. Algoritmos populares incluem o K-means, hierarchical clustering e DBSCAN.

Esses são apenas alguns exemplos de modelos de aprendizado de máquina. A escolha do modelo correto depende da natureza dos dados, do problema a ser resolvido e da disponibilidade de recursos computacionais. Cada modelo tem suas próprias suposições e propriedades matemáticas, o que os torna mais adequados para determinados tipos de dados ou tarefas específicas. Além dos modelos mencionados, existem muitos outros algoritmos e abordagens, como redes bayesianas, algoritmos genéticos, máquinas de aprendizado extremo e muitos mais.

Ao selecionar um modelo de aprendizado de máquina, é importante considerar algumas características, como a capacidade de generalização do modelo, interpretabilidade, complexidade computacional, requisitos de dados e a quantidade de dados disponíveis para treinamento. Cada modelo tem suas vantagens e limitações, e a escolha adequada dependerá do contexto e dos objetivos do projeto.

Além disso, é importante mencionar que a construção de um modelo de aprendizado de máquina não se resume à seleção do algoritmo. Envolve também a etapa de treinamento, em que o modelo é ajustado aos dados de treinamento para aprender os padrões e relações. Isso requer a definição adequada das características e atributos relevantes, a escolha de uma função de perda apropriada e a otimização dos parâmetros do modelo para minimizar o erro.

Após o treinamento, é necessário avaliar o desempenho do modelo usando dados de validação ou teste. Isso permite verificar a capacidade do modelo de generalizar e fazer previsões precisas em dados não vistos anteriormente. Durante a avaliação, métricas como acurácia, precisão, recall e F1-score são frequentemente utilizadas para medir o desempenho do modelo.

É importante destacar que a escolha e o desempenho do modelo não são garantia de sucesso. Outros fatores, como a qualidade dos dados, o pré-processamento adequado, a seleção de recursos relevantes e a interpretação correta dos resultados, também desempenham um papel fundamental no sucesso de um projeto de aprendizado de máquina.

Os modelos de aprendizado de máquina são representações matemáticas ou estatísticas que capturam os padrões e relações nos dados. Existem diversos tipos de modelos, cada um com suas características e aplicações específicas. A escolha do modelo correto depende do problema em questão, dos dados disponíveis e dos objetivos do projeto. A construção de um modelo envolve também o treinamento adequado, a avaliação de desempenho e a consideração de outros fatores importantes para o sucesso do projeto.

4.3 Árvore de decisão

A árvore de decisão é um modelo de aprendizado de máquina que representa uma estrutura hierárquica de decisões e suas possíveis consequências. Essa técnica é utilizada para problemas de classificação e regressão, em que o objetivo é tomar decisões ou prever valores com base em características ou atributos dos dados.

Uma árvore de decisão é composta de nós que representam as características dos dados, arestas que representam as decisões tomadas com base nessas características e folhas que representam as classes ou

valores de saída. A construção da árvore envolve a seleção de características relevantes e a definição de regras de decisão que dividem os dados em subconjuntos mais puros ou que melhor separam as classes.

A figura a seguir apresenta uma possível árvore de decisão da "adivinhação" de animais:

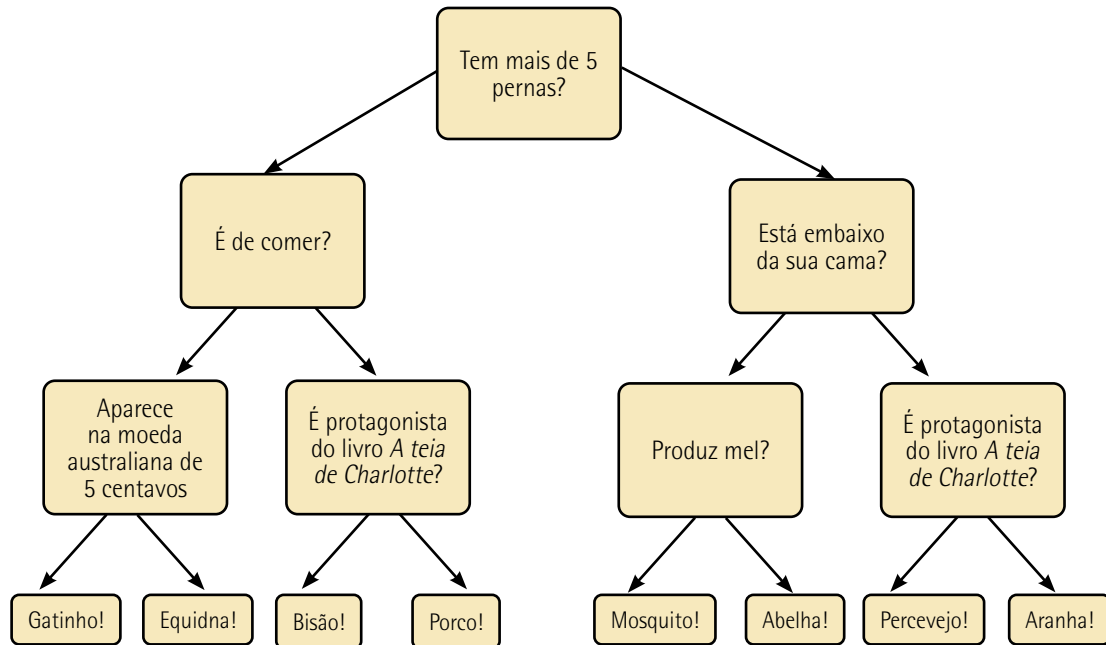


Figura 15 – A árvore de decisão da "adivinhação" de animais

Fonte: Grus (2016, p. 231).

O processo de construção de uma árvore de decisão é geralmente baseado em algoritmos como o ID3 (Iterative Dichotomiser 3), C4.5, CART (Classification and Regression Trees) ou random forest. Esses algoritmos buscam encontrar a melhor divisão dos dados em cada nó da árvore, considerando critérios como ganho de informação, índice de Gini ou erro quadrático médio.

Uma das principais vantagens das árvores de decisão é a sua interpretabilidade. As regras de decisão representadas pela árvore são facilmente compreensíveis e podem fornecer insights sobre os padrões presentes nos dados. Além disso, as árvores de decisão podem lidar com dados numéricos e categóricos, e ser aplicadas em problemas com várias classes.

No entanto, as árvores de decisão também apresentam algumas limitações. Elas podem ser sensíveis a pequenas variações nos dados de treinamento, o que pode levar a uma alta variância e overfitting. Além disso, a construção de árvores de decisão com muitos nós ou níveis pode levar à complexidade excessiva e à perda de generalização.

Para mitigar esses problemas, podem ser aplicadas técnicas como a poda da árvore (pruning), que remove ramos desnecessários, e a combinação de várias árvores em um conjunto, como no caso do random forest. Essas abordagens visam melhorar o desempenho e a robustez das árvores de decisão.

Veremos na sequência exemplo simples de como construir uma árvore de decisão usando a biblioteca scikit-learn em Python. Vamos usar um conjunto de dados de flores íris para criar a árvore de decisão e prever a classe das flores com base em suas características. As características utilizadas são comprimento da sépala, largura da sépala, comprimento da pétala e largura da pétala. As classes das flores são *Iris setosa*, *Iris virginica* e *Iris versicolor*.

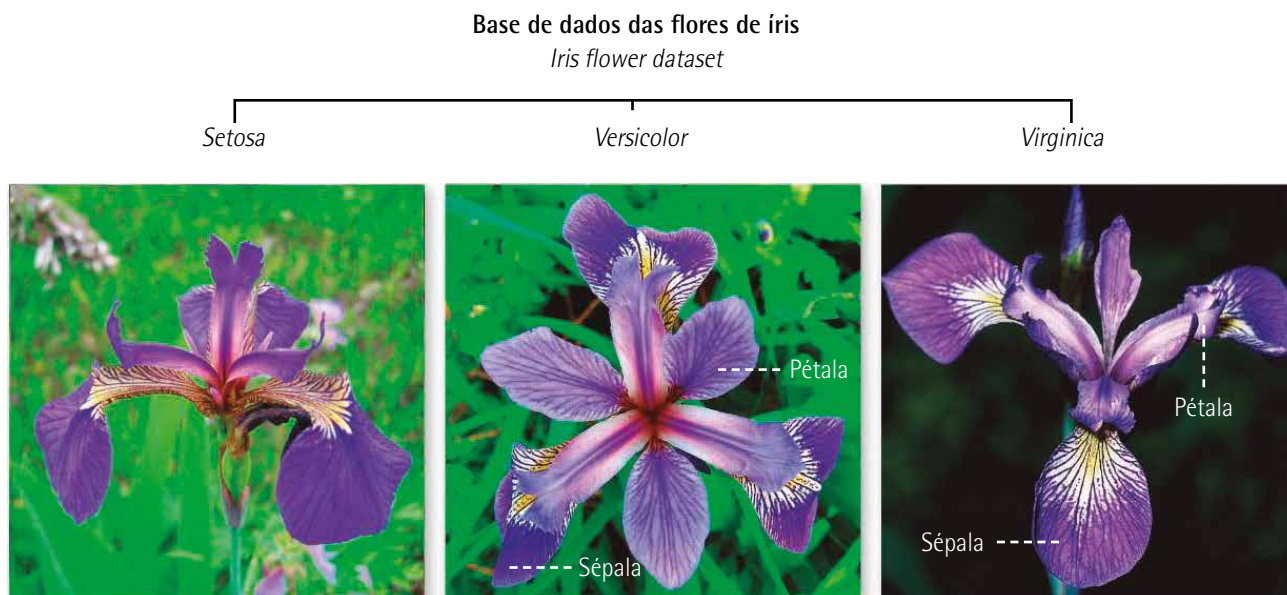


Figura 16 – Diferenças entre os três tipos de flores de íris usadas no Iris Flower Dataset



Lembrete

Disponível em: <https://tinyurl.com/yks2dd37>. Acesso em: 10 jul. 2023.

Você deve estar com o scikit-learn instalado. É possível instalá-lo usando o seguinte comando: `pip install scikit-learn`

1. # Importando as bibliotecas necessárias
2. `from sklearn.datasets import load_iris`
3. `from sklearn.tree import DecisionTreeClassifier`
4. `from sklearn.model_selection import train_test_split`
5. `from sklearn.metrics import accuracy_score`
6. # Carregando o conjunto de dados Iris
7. `iris = load_iris()`
8. `X = iris.data` # Características das flores
9. `y = iris.target` # Classes das flores
10. # Dividindo o conjunto de dados em treinamento e teste
11. `X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)`
12. # Criando o classificador de árvore de decisão
13. `clf = DecisionTreeClassifier()`


```
14. # Treinando o classificador com o conjunto de treinamento
15. clf.fit(X_train, y_train)

16. # Fazendo previsões no conjunto de teste
17. y_pred = clf.predict(X_test)

18. # Calculando a precisão do classificador
19. accuracy = accuracy_score(y_test, y_pred)
20. print(f"Precisão do classificador: {accuracy:.2f}")
```

Saída: precisão do classificador: 1.00

O código carrega o conjunto de dados Iris, divide-o em conjuntos de treinamento e teste, cria um classificador de árvore de decisão, treina o classificador com o conjunto de treinamento e, em seguida, avalia a precisão do classificador no conjunto de teste. Lembre-se de que a precisão pode variar dependendo da divisão aleatória dos dados em treinamento e teste.

A acurácia de 1 (ou 100%) significa que o classificador de árvore de decisão fez previsões perfeitamente corretas para todas as amostras do conjunto de teste. Em outras palavras, todas as flores no conjunto de teste foram classificadas corretamente em suas respectivas classes. Uma acurácia de 1 é um resultado ideal e indica que o modelo de árvore de decisão conseguiu aprender com sucesso as relações entre as características das flores e suas classes no conjunto de treinamento e aplicou esse conhecimento para fazer previsões precisas no conjunto de teste.

A árvore de decisão é um modelo de aprendizado de máquina que utiliza uma estrutura hierárquica de decisões para realizar classificação ou regressão. Elas são facilmente interpretáveis e adequadas para uma variedade de problemas. No entanto, é importante considerar as limitações e aplicar técnicas adicionais para obter resultados mais precisos e robustos.

4.4 Naive Bayes

O Naive Bayes é um modelo de aprendizado de máquina baseado no Teorema de Bayes e na suposição de independência condicional entre os recursos. Ele é comumente usado para problemas de classificação, especialmente em tarefas de processamento de linguagem natural, como análise de sentimentos, detecção de spam e categorização de documentos.

O modelo Naive Bayes é chamado "ingênuo" porque assume que todas as características são independentes entre si, ou seja, não há correlação entre elas. Essa é uma suposição simplificada, mas que permite uma implementação eficiente e resultados satisfatórios na prática.

O funcionamento do Naive Bayes é baseado na aplicação do Teorema de Bayes, que descreve a probabilidade de um evento ocorrer dado o conhecimento prévio sobre o evento. No caso do Naive Bayes, a probabilidade de uma classe dadas as características é calculada usando a seguinte fórmula:

Onde:

- $P(\text{Classe} \mid \text{Características})$ é a probabilidade da classe dada as características observadas.
- $P(\text{Classe})$ é a probabilidade de a classe ocorrer independentemente das características.
- $P(\text{Características} \mid \text{Classe})$ é a probabilidade de as características ocorrerem dada a classe.
- $P(\text{Características})$ é a probabilidade de as características ocorrerem independentemente da classe.

Para realizar a classificação, o modelo Naive Bayes calcula a probabilidade de cada classe para um determinado conjunto de características e seleciona a classe com a maior probabilidade como a predição.

Existem diferentes variantes do Naive Bayes, como o Naive Bayes Gaussiano, Bernoulli e Multinomial, que se diferenciam na forma como as distribuições de probabilidade são modeladas para diferentes tipos de características (contínuas, binárias ou discretas, respectivamente).

Uma das principais vantagens do Naive Bayes é a sua simplicidade e eficiência computacional. Ele requer uma quantidade relativamente pequena de dados para treinamento e é menos sensível a overfitting em comparação com outros modelos mais complexos. Além disso, o Naive Bayes é robusto em relação a recursos irrelevantes ou redundantes, pois eles são considerados independentes.

No entanto, a suposição de independência condicional entre as características pode ser inadequada em alguns casos, especialmente quando as características estão correlacionadas. Além disso, o Naive Bayes pode ter dificuldade em lidar com características ausentes ou raros eventos que não foram observados durante o treinamento.

Apesar dessas limitações, o Naive Bayes é utilizado e fornece resultados satisfatórios em muitos problemas de classificação. É uma escolha popular, especialmente quando se lida com grandes volumes de dados e problemas de processamento de linguagem natural. Vamos dar um exemplo de código em Python utilizando o classificador Naive Bayes.

```
1. # Importando as bibliotecas necessárias
2. from sklearn.datasets import load_iris
3. from sklearn.model_selection import train_test_split
4. from sklearn.naive_bayes import GaussianNB
5. from sklearn.metrics import accuracy_score

6. # Carregando o conjunto de dados Iris
7. iris = load_iris()
8. X = iris.data # Características das flores
9. y = iris.target # Classes das flores
```

```
10. # Dividindo o conjunto de dados em treinamento e teste
11. X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

12. # Criando o classificador Naive Bayes
13. clf = GaussianNB()

14. # Treinando o classificador com o conjunto de treinamento
15. clf.fit(X_train, y_train)

16. # Fazendo previsões no conjunto de teste
17. y_pred = clf.predict(X_test)

18. # Calculando a precisão do classificador
19. accuracy = accuracy_score(y_test, y_pred)
20. print(f"Precisão do classificador Naive Bayes: {accuracy:.2f}")
```

Saída: precisão do classificador Naive Bayes: 1.00

O código carrega o conjunto de dados Iris, divide-o em conjuntos de treinamento e teste, cria um classificador Naive Bayes usando a classe 'GaussianNB' do 'scikit-learn', treina o classificador com o conjunto de treinamento e, em seguida, avalia a precisão do classificador no conjunto de teste. O classificador Naive Bayes assume independência entre as características, o que pode não ser estritamente verdadeiro para todos os conjuntos de dados. Portanto, o desempenho do Naive Bayes pode variar dependendo da natureza do conjunto de dados.

4.5 K-vizinhos mais próximos (KNN)

O algoritmo dos k-vizinhos mais próximos, conhecido como KNN (do inglês, K-nearest neighbors), é um método de aprendizado de máquina utilizado para classificação e regressão. Ele é baseado no princípio de que amostras com características semelhantes tendem a ter rótulos ou valores de saída semelhantes.

No KNN, o objetivo é classificar uma nova amostra ou prever seu valor de saída com base nas informações dos k-vizinhos mais próximos presentes no conjunto de treinamento. A distância entre as amostras é geralmente medida usando métricas como a distância euclidiana ou a distância de Manhattan.

O funcionamento do KNN é relativamente simples. Ao receber uma nova amostra, o algoritmo calcula sua distância em relação a todas as outras amostras do conjunto de treinamento. Em seguida, seleciona os k-vizinhos mais próximos com base nessa distância. A classe mais frequente entre esses vizinhos é atribuída à nova amostra no caso de classificação, ou a média ou mediana dos valores de saída dos vizinhos é usada como uma previsão no caso de regressão.

A escolha do valor de k é um aspecto crítico no KNN. Valores pequenos de k tornam o modelo mais sensível ao ruído nos dados, enquanto valores grandes de k podem levar a uma suavização excessiva e à perda de detalhes importantes. Portanto, é importante selecionar o valor de k adequado para cada problema específico.

O KNN é um algoritmo não paramétrico, o que significa que não faz suposições sobre a distribuição subjacente dos dados. Ele é capaz de lidar com dados numéricos e categóricos, e não requer treinamento prévio extenso, pois os dados de treinamento são armazenados diretamente para a classificação ou regressão.

No entanto, o KNN tem algumas limitações. Um dos principais desafios é o cálculo da distância em conjuntos de dados grandes, o que pode tornar o processo computacionalmente custoso. Além disso, o desempenho do KNN pode ser afetado por desequilíbrios de classe nos dados, bem como pela presença de atributos irrelevantes ou de alta dimensionalidade.

Vamos exemplificar como podemos usar o algoritmo k-nearest neighbors (KNN) com a base de dados Iris usando a biblioteca scikit-learn em Python:

```
1. # Importando as bibliotecas necessárias
2. from sklearn.datasets import load_iris
3. from sklearn.model_selection import train_test_split
4. from sklearn.neighbors import KNeighborsClassifier
5. from sklearn.metrics import accuracy_score

6. # Carregando o conjunto de dados Iris
7. iris = load_iris()
8. X = iris.data # Características das flores
9. y = iris.target # Classes das flores

10. # Dividindo o conjunto de dados em treinamento e teste
11. X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

12. # Criando o classificador k-NN com k=3 (vamos usar 3 vizinhos mais próximos)
13. k = 3
14. clf = KNeighborsClassifier(n_neighbors=k)

15. # Treinando o classificador com o conjunto de treinamento
16. clf.fit(X_train, y_train)

17. # Fazendo previsões no conjunto de teste
18. y_pred = clf.predict(X_test)
```

```
19. # Calculando a precisão do classificador
20. accuracy = accuracy_score(y_test, y_pred)
21. print(f"Precisão do classificador k-NN (k={k}): {accuracy:.2f}")
```

Saída: precisão do classificador k-NN (k=3): 1.00

O código carrega o conjunto de dados Iris, divide-o em conjuntos de treinamento e teste, cria um classificador KNN usando a classe `KNeighborsClassifier` do `scikit-learn`, treina o classificador com o conjunto de treinamento, especificando que queremos usar três vizinhos mais próximos (você pode ajustar o valor de `k` para experimentar diferentes números de vizinhos) e, em seguida, avalia a precisão do classificador no conjunto de teste.

O KNN é um algoritmo simples de aprendizado de máquina baseado em instâncias que faz previsões com base nos vizinhos mais próximos no espaço de características. A eficácia do KNN pode ser influenciada pelo valor de `k` escolhido e pela natureza do conjunto de dados. Experimente diferentes valores de `k` para ver como isso afeta a precisão do classificador.

4.6 K-médias

O algoritmo K-médias, também conhecido como K-means, é um método de aprendizado de máquina não supervisionado utilizado para realizar agrupamento de dados. Ele é utilizado em problemas de mineração de dados e análise exploratória, em que o objetivo é encontrar grupos ou clusters de amostras que sejam similares entre si.

O algoritmo K-médias funciona da seguinte maneira:

1. Escolha o número de clusters `K` que você deseja criar. Isso representa o número de grupos em que os dados serão divididos.
2. Selecione aleatoriamente `K` pontos iniciais como centros dos clusters. Esses pontos são chamados de centroides.
3. Atribua cada amostra do conjunto de dados ao cluster cujo centroide está mais próximo. A distância entre as amostras e os centroides pode ser medida usando a distância euclidiana.
4. Recalcule os centroides dos clusters, que são os novos pontos centrais baseados nas amostras atribuídas a cada cluster.
5. Repita os passos 3 e 4 até que haja convergência, ou seja, até que não ocorram mais mudanças na atribuição das amostras aos clusters ou nos centroides.

O objetivo do algoritmo K-médias é minimizar a soma dos quadrados das distâncias entre as amostras e seus centroides correspondentes. Essa métrica é chamada de "inércia" e representa a coesão dentro de cada cluster. Quanto menor a inércia, mais compactos e bem-definidos são os clusters.

O algoritmo K-médias é rápido e escalável, tornando-o adequado para grandes conjuntos de dados, e é um método de agrupamento rígido em que cada amostra é atribuída exclusivamente a um cluster. O resultado do K-médias pode variar dependendo da inicialização aleatória dos centroides. Portanto, é comum executar o algoritmo várias vezes com diferentes inicializações e selecionar a solução com a menor inércia. O número de clusters K deve ser fornecido pelo usuário. Determinar o valor adequado de K pode ser um desafio e requer conhecimento do domínio e análise dos dados. O K-médias pressupõe que os clusters sejam esféricos e de tamanho aproximadamente igual. Portanto, ele pode não ser eficaz em dados com formas irregulares ou com tamanhos de cluster muito diferentes.

Neste exemplo, vamos usar a biblioteca scikit-learn para aplicar o K-médias na base de dados Iris, mas lembre-se de que o K-means é um algoritmo de clustering, e a base Iris é frequentemente usada para classificação, então estamos aplicando o K-means para encontrar grupos não rotulados nos dados.

```
1. # Importando as bibliotecas necessárias
2. from sklearn.datasets import load_iris
3. from sklearn.cluster import KMeans
4. import matplotlib.pyplot as plt

5. # Carregando o conjunto de dados Iris
6. iris = load_iris()
7. X = iris.data # Características das flores

8. # Criando o modelo K-Means com 3 clusters (pois sabemos que existem 3 classes de flores
   na base Iris)
9. n_clusters = 3
10. kmeans = KMeans(n_clusters=n_clusters, random_state=42)

11. # Aplicando o K-Means aos dados
12. kmeans.fit(X)

13. # Obtendo as etiquetas dos clusters (rótulos dos grupos) para cada amostra
14. cluster_labels = kmeans.labels_

15. # Plotando os resultados
16. plt.scatter(X[:, 0], X[:, 1], c=cluster_labels, cmap='viridis')
17. centers = kmeans.cluster_centers_
18. plt.scatter(centers[:, 0], centers[:, 1], c='red', marker='X', s=200)
19. plt.xlabel('Comprimento da Sépala')
20. plt.ylabel('Largura da Sépala')
21. plt.title('K-Means Clustering - Base Iris')
22. plt.show()
```

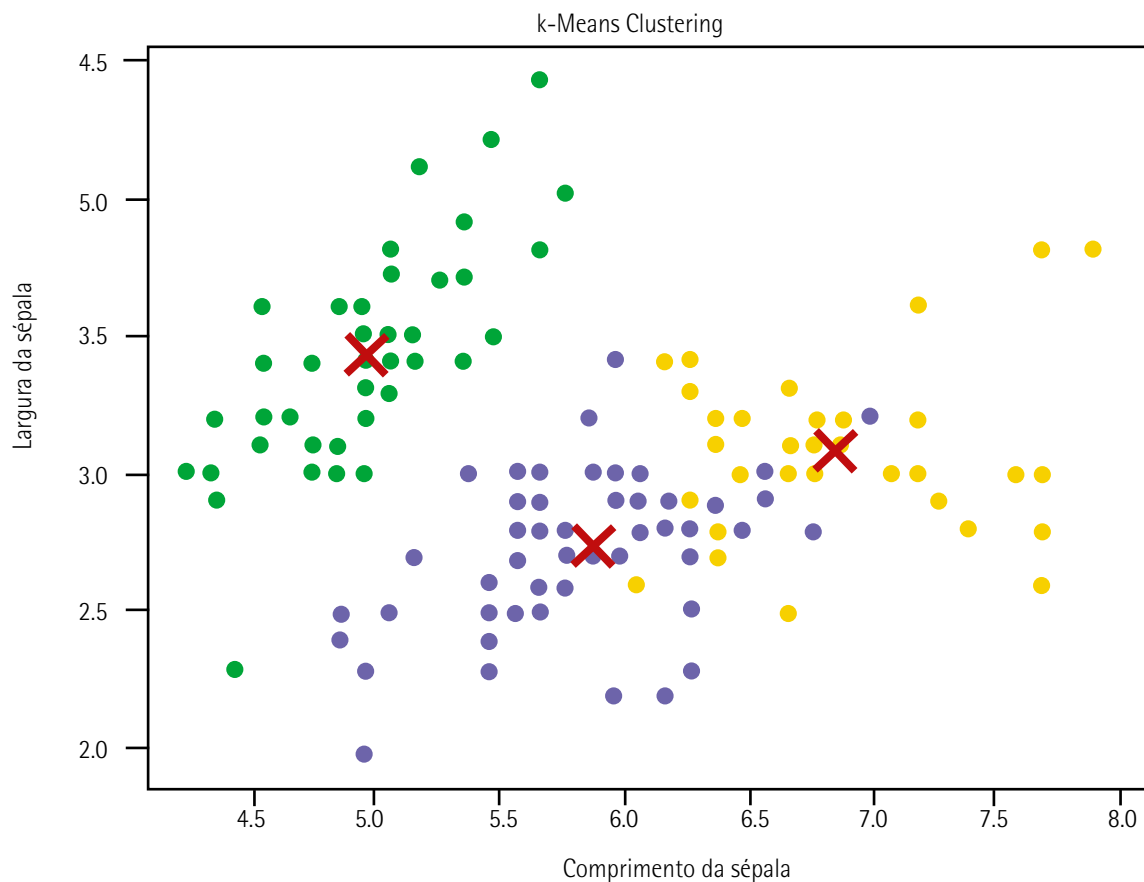


Figura 17

Neste exemplo, carregamos a base de dados Iris, criamos um modelo K-means com três clusters (pois sabemos que existem três classes de flores na base Iris), aplicamos o K-means aos dados e, em seguida, plotamos as amostras de acordo com os rótulos dos clusters encontrados. Além disso, marcamos os centroides dos clusters com marcadores vermelhos (X).

É importante observar que, como mencionado anteriormente, o K-means é um algoritmo de clustering, e a base de dados Iris é rotulada, o que significa que os rótulos reais das classes são conhecidos. Portanto, você pode comparar visualmente os resultados do K-means com os rótulos reais da base Iris. Neste caso, esperamos ver uma boa separação entre os clusters, mas a correspondência direta com as classes reais pode não ser perfeita, já que o K-means não tem conhecimento dos rótulos de classe na base Iris.

O K-médias é utilizado em aplicações como segmentação de clientes, agrupamento de documentos, análise de imagens e identificação de padrões em dados. Ele é uma técnica poderosa para explorar e entender a estrutura de conjuntos de dados não rotulados, permitindo a descoberta de grupos sem a necessidade de rótulos prévios.



Resumo

Vimos nesta unidade que a ciência de dados é um campo interdisciplinar que envolve a coleta, processamento, análise e interpretação de dados para obter insights, conhecimentos e apoio à tomada de decisões. Ela combina técnicas de programação, estatística, matemática, conhecimento de domínio e habilidades de comunicação para extrair informações significativas de conjuntos de dados complexos. O ciclo de vida típico de um projeto de ciência de dados envolve a definição de objetivos, coleta e preparação de dados, análise exploratória, modelagem estatística ou de aprendizado de máquina, avaliação e interpretação dos resultados.

Exibimos o tema aprendizado de máquina, que é uma subárea da inteligência artificial que se concentra no desenvolvimento de algoritmos e modelos que permitem que um sistema "aprenda" a partir de dados e melhore seu desempenho em tarefas específicas ao longo do tempo. Ele é dividido em três categorias principais: aprendizado supervisionado (em que os modelos são treinados em dados rotulados), aprendizado não supervisionado (em que os modelos identificam padrões em dados não rotulados) e aprendizado por reforço (em que os modelos aprendem a tomar ações para maximizar recompensas em um ambiente).

Mostramos a diferença entre aprendizado descritivo e aprendizado preditivo. O aprendizado descritivo refere-se à análise exploratória de dados para entender padrões, tendências e relações entre variáveis. Ele não visa fazer previsões, mas sim fornecer insights sobre os dados. Em contraste, o aprendizado preditivo envolve a construção de modelos que podem fazer previsões com base em dados históricos. Esses modelos são treinados em dados passados e são usados para fazer previsões futuras. A combinação desses dois tipos de aprendizado pode fornecer uma compreensão abrangente e útil dos dados.

Exploramos o tema mineração de dados, que é o processo de descobrir padrões, informações relevantes e conhecimentos implícitos em grandes volumes de dados. Envolve a aplicação de técnicas de análise de dados, estatísticas, aprendizado de máquina e visualização para identificar padrões ocultos, tendências e relacionamentos em conjuntos de dados. A mineração de dados pode ser usada em várias áreas, como marketing, finanças, saúde e muito mais, para extrair informações valiosas e tomar decisões informadas.

A ciência de dados é a disciplina que trata da extração de insights de dados, o aprendizado de máquina é uma técnica-chave dentro dessa disciplina para construir modelos de previsão e análise, o aprendizado descritivo e o preditivo diferenciam a análise exploratória da previsão, e a mineração de dados é um processo para descobrir padrões ocultos em grandes conjuntos de dados.



Exercícios

Questão 1. (FGV 2022, adaptada) Um time de ciência de dados utilizou um modelo linear para resolver uma tarefa de análise de dados financeiros provenientes de diferentes unidades de uma organização. Um membro do time, que não participou da modelagem, testa o modelo e verifica que ele apresenta um péssimo resultado. Preocupado, ele busca os resultados apresentados no treino para verificar o problema. A respeito desse cenário, avalie as afirmativas.

- I – Se o resultado do treino foi ótimo, ocorreu underfitting. Uma possível solução corresponde à utilização de um modelo mais complexo e à redução do tempo de treinamento.
- II – Se o resultado do treino também foi péssimo, ocorreu overfitting. Uma possível solução corresponde à utilização de técnicas de regularização e de métodos de validação cruzada.
- III – Se o resultado do treino foi ótimo, ocorreu overfitting. Uma possível solução corresponde à utilização de um modelo menos complexo e de métodos de validação cruzada.

É correto o que se afirma em:

- A) I, apenas.
- B) III, apenas.
- C) I e III, apenas.
- D) II e III, apenas.
- E) I, II e III.

Resposta correta: alternativa B.

Análise das afirmativas

I – Afirmativa incorreta.

Justificativa: o underfitting ocorre quando o modelo de aprendizado de máquina é muito simples ou não é capaz de capturar os padrões presentes nos dados de treinamento. Um dos sinais do underfitting é o desempenho ruim, tanto no treinamento quanto nos testes. Portanto, a apresentação de resultados ótimos de treino não indica a ocorrência de underfitting.

II – Afirmativa incorreta.

Justificativa: o overfitting ocorre quando o modelo se torna excessivamente complexo e se ajusta perfeitamente aos dados de treinamento, capturando até mesmo o ruído presente nesses dados. Um dos sinais do overfitting é o desempenho ótimo no treino, mas ruim nos testes. Logo, se o resultado do treino também foi péssimo, não há evidência de overfitting no modelo.

III – Afirmativa correta.

Justificativa: o cenário em que o resultado de teste foi péssimo, mas o resultado de treino foi ótimo, aponta para o problema de overfitting. Nesse caso, o modelo é muito complexo e se ajusta perfeitamente aos dados de treino, capturando até mesmo o ruído presente nesses dados. Em tal cenário, duas atitudes que podem ajudar a resolver o problema são reduzir a complexidade do modelo e utilizar validação cruzada, para avaliar o desempenho do modelo em conjuntos de dados distintos.

Questão 2. (IDIB 2020, adaptada) A mineração de dados, também conhecida como data mining, é uma das muitas áreas da computação e tem como objetivo identificar correlações e padrões em um grande conjunto de dados, com o intuito de prever resultados. A respeito dos conceitos que fazem parte da mineração de dados, avalie as afirmativas.

I – Redes neurais e árvores de decisão são dois conhecidos exemplos de ramificações da mineração de dados.

II – São exemplos de etapas da mineração de dados: seleção de dados, preparação de dados e utilização do conhecimento.

III – A mineração de dados faz uso de fundamentos pertencentes a outras três grandes áreas de conhecimento: matemática, estatística e data warehouse.

É correto o que se afirma em:

A) I, apenas.

B) III, apenas.

C) I e II, apenas.

D) II e III, apenas.

E) I, II e III.

Resposta correta: alternativa C.

