



FACULDADE IMPACTA

CURSOS

BIG DATA - INTRODUÇÃO E OPORTUNIDADES



ALEX SOUSA

SÃO PAULO - 02/2024



SUMÁRIO

SUMÁRIO.....	1
Big Data - Introdução e Oportunidades.....	4
Introdução ao Big Data.....	4
Por que a análise de big data é importante?.....	4
Como Funciona e Principais Tecnologias.....	5
Computação em Nuvem.....	5
Gestão de Dados.....	5
Mineração de Dados.....	5
Armazenamento de Dados.....	5
Hadoop.....	5
Análise na Memória.....	5
Aprendizado de Máquina.....	6
Análise Preditiva.....	6
Mineração de Texto.....	6
Coleta de Dados.....	6
Sensores e Dispositivos IoT.....	6
Websites e Aplicações.....	7
Redes Sociais.....	7
Transações Comerciais.....	7
Logs de Servidores e Dispositivos.....	7
Evolução dos Dados.....	7
Dados Estruturados.....	7
Dados Semi-Estruturados.....	8
Dados Não Estruturados.....	8
Grande Volume de Dados.....	8
Armazenamento.....	8
Transmissão de Dados.....	8
Processamento.....	9
Gerenciamento e Limpeza de Dados.....	9
Qualidade.....	9
Segurança e Privacidade.....	9
Interpretação.....	9
Integração.....	9
Habilidades.....	9
Regulamentação.....	9
Custo.....	10
Ética.....	10
Sistema de Banco de Dados Cassandra.....	10
Principais características do Cassandra.....	10
Aplicações do Cassandra.....	11
Vantagens do Cassandra.....	11
Desvantagens do Cassandra.....	12
Os 5 V's da Big Data.....	12

Volume.....	12
Velocidade.....	12
Variedade.....	12
Veracidade.....	12
Valor.....	13
Como Obter Valor dos Dados.....	13
4 Regras Básica.....	13
Exatidão.....	13
Consistência dos Dados.....	13
Completude dos Dados.....	13
Relevância dos Dados.....	13
Cases de Sucesso - O Mistério da Cerveja e das Fraldas - Walmart.....	13
A Jornada Exploratória.....	13
A Teoria Revelada.....	14
O Impacto da Descoberta.....	14
Resultados Extraordinários.....	14
Lições Aprendidas.....	14
Atividades e áreas de atuação.....	14
Geração de Insights.....	15
Aplicações em Áreas Específicas.....	15
Habilidades Essenciais para o Sucesso.....	15
Ciência de Dados.....	15
1. Análise de Dados.....	15
2. Análise de Tendências.....	15
3. Análise por Algoritmos - Machine Learning.....	16
Engenharias de Dados.....	16
Responsabilidades de um Engenheiro de Dados.....	16
Habilidades de um Engenheiro de Dados.....	16
ETL - Extract Transform Load.....	16
Extract (Extração).....	17
Transform (Transformação).....	17
Load (Carregamento).....	17
Machine Learning.....	17
Tipos de Machine Learning.....	17
Aprendizado Supervisionado.....	18
Aprendizado Não Supervisionado.....	18
Aprendizado por Reforço.....	18
Aplicações de Machine Learning.....	18
Área de Atuação Big Data.....	18
Marketing e Vendas.....	18
Finanças.....	19
Saúde.....	19
Indústria.....	19
Governo.....	19
Outras áreas de atuação.....	19

Big Data - Introdução e Oportunidades

Você será apresentado aos conceitos fundamentais do universo de Big Data, como sua origem e o processo de evolução dos dados. Aprenderá a identificar o valor dos dados por meio da aplicação de 4 regras básicas. Assimilará a importância do uso e investimento em Big Data para melhoria dos resultados de pesquisa, redução de custos e aumento dos lucros das empresas. Por fim, elucidará a respeito das principais atividades e áreas de atuação do profissional de Big Data.

Introdução ao Big Data

Big data é um termo utilizado para descrever conjuntos de dados extremamente grandes e complexos que tradicionais softwares de processamento de dados têm dificuldade em lidar. Esses conjuntos de dados geralmente apresentam uma variedade de formatos e são gerados a uma velocidade muito alta.

Big data são ativos de informações de alto volume, alta velocidade e/ou alta variedade que exigem formas inovadoras e econômicas de processamento de informações que permitem maior percepção, tomada de decisões e automação de processos.

Big data é caracterizado por 5 principais dimensões: volume (a quantidade de dados), velocidade (a taxa de geração e processamento dos dados), variedade (a diversidade de tipos de dados), veracidade (à qualidade dos dados) e o valor (à capacidade de extrair valor útil).

Para lidar com big data, são necessárias técnicas e tecnologias especiais, como computação em nuvem, algoritmos de processamento distribuído, sistemas de armazenamento distribuído, análise de dados avançada e machine learning. O objetivo do processamento de big data é extrair insights significativos, padrões ou conhecimentos úteis a partir desses grandes conjuntos de dados, para tomada de decisões informadas e desenvolvimento de estratégias em diversos campos, como negócios, ciência, saúde, governo, entre outros.

Por que a análise de big data é importante?

A análise de big data ajuda as organizações a aproveitar seus dados e usá-los para identificar novas oportunidades. Isso, por sua vez, leva a movimentos comerciais mais inteligentes, operações mais eficientes, lucros mais elevados e clientes mais satisfeitos. As empresas que usam big data com análises avançadas ganham valor de várias maneiras, como:

- **Reduzindo custos:** Tecnologias de big data, como análises baseadas em nuvem, podem reduzir significativamente os custos quando se trata de armazenar grandes quantidades de dados (por exemplo, um data lake). Além disso, a análise de big data ajuda as organizações a encontrar formas mais eficientes de fazer negócios.
- **Tomar decisões melhores e mais rápidas:** A velocidade da análise in-memory – combinada com a capacidade de analisar novas fontes de dados, como o streaming de dados da IoT – ajuda as empresas a analisar informações imediatamente e a tomar decisões rápidas e informadas.
- **Desenvolvimento e comercialização de novos produtos e serviços:** Ser capaz de avaliar as necessidades e a satisfação do cliente por meio de análises permite que as empresas ofereçam aos clientes o que eles desejam, quando desejam. Com a análise de big data, mais empresas têm a oportunidade de desenvolver novos produtos inovadores para atender às necessidades em constante mudança dos clientes.

Como Funciona e Principais Tecnologias

Não existe uma tecnologia única que englobe a análise de big data. É claro que existem análises avançadas que podem ser aplicadas a big data, mas, na realidade, vários tipos de tecnologia trabalham em conjunto para ajudá-lo a obter o máximo valor das suas informações. Aqui estão os maiores jogadores:

Computação em Nuvem

Um modelo de entrega baseado em assinatura, a computação em nuvem fornece escalabilidade, entrega rápida e eficiência de TI necessárias para análises eficazes de big data. Como elimina muitas barreiras físicas e financeiras ao alinhamento das necessidades de TI com os objetivos empresariais em evolução, é apelativo para organizações de todas as dimensões.

Gestão de Dados

Os dados precisam ser de alta qualidade e bem controlados antes que possam ser analisados de forma confiável. Com os dados entrando e saindo constantemente de uma organização, é importante estabelecer processos repetíveis para criar e manter padrões de qualidade de dados. Uma vez que os dados sejam confiáveis, as organizações devem estabelecer um programa de gerenciamento de dados mestre que coloque toda a empresa na mesma página.

Mineração de Dados

A tecnologia de mineração de dados ajuda a examinar grandes quantidades de dados para descobrir padrões nos dados – e essas informações podem ser usadas para análises adicionais para ajudar a responder questões comerciais complexas. Com o software de mineração de dados, você pode examinar todo o ruído caótico e repetitivo dos dados, identificar o que é relevante, usar essas informações para avaliar resultados prováveis e, então, acelerar o ritmo da tomada de decisões informadas.

Armazenamento de Dados

incluindo data lake e data warehouse . É vital poder armazenar grandes quantidades de dados estruturados e não estruturados – para que os utilizadores empresariais e os cientistas de dados possam aceder e utilizar os dados conforme necessário. Um data lake ingere rapidamente grandes quantidades de dados brutos em seu formato nativo. É ideal para armazenar big data não estruturado, como conteúdo de mídia social, imagens, voz e dados de streaming. Um data warehouse armazena grandes quantidades de dados estruturados em um banco de dados central. Os dois métodos de armazenamento são complementares; muitas organizações usam ambos.

Hadoop

Essa estrutura de software de código aberto facilita o armazenamento de grandes quantidades de dados e permite a execução de aplicativos paralelos em clusters de hardware comuns. Tornou-se uma tecnologia chave para fazer negócios devido ao aumento constante de volumes e variedades de dados, e seu modelo de computação distribuída processa big data rapidamente. Um benefício adicional é que a estrutura de código aberto do Hadoop é gratuita e utiliza hardware comum para armazenar e processar grandes quantidades de dados.

Análise na Memória

Ao analisar dados da memória do sistema (em vez de da unidade de disco rígido), você pode obter insights imediatos de seus dados e agir rapidamente. Essa tecnologia é capaz de remover latências de preparação de dados e processamento analítico para testar novos cenários e criar modelos; não é apenas uma maneira fácil para as organizações permanecerem ágeis e tomarem melhores decisões de negócios, mas também lhes permite executar cenários analíticos iterativos e interativos.

Aprendizado de Máquina

O aprendizado de máquina, um subconjunto específico de IA que treina uma máquina para aprender, possibilita a produção rápida e automática de modelos que podem analisar dados maiores e mais complexos e fornecer resultados mais rápidos e precisos – mesmo em uma escala muito grande. E ao construir modelos precisos, uma organização tem mais hipóteses de identificar oportunidades lucrativas – ou de evitar riscos desconhecidos.

Análise Preditiva

A tecnologia de análise preditiva utiliza dados, algoritmos estatísticos e técnicas de aprendizado de máquina para identificar a probabilidade de resultados futuros com base em dados históricos. Trata-se de fornecer a melhor avaliação do que acontecerá no futuro, para que as organizações possam se sentir mais confiantes de que estão tomando a melhor decisão de negócios possível. Algumas das aplicações mais comuns de análise preditiva incluem detecção de fraudes, riscos, operações e marketing.

Mineração de Texto

Com a tecnologia de mineração de texto, você pode analisar dados de texto da Web, campos de comentários, livros e outras fontes baseadas em texto para descobrir insights que você não havia notado antes. A mineração de texto usa aprendizado de máquina ou tecnologia de processamento de linguagem natural para vasculhar documentos – e-mails, blogs, feeds do Twitter, pesquisas, inteligência competitiva e muito mais – para ajudá-lo a analisar grandes quantidades de informações e descobrir novos tópicos e relacionamentos de curto prazo.

Coleta de Dados

A coleta de dados é um componente fundamental no processo de gestão e análise de big data. Envolve a captura sistemática e organizada de informações de diversas fontes para serem posteriormente processadas, analisadas e utilizadas para gerar insights valiosos. Existem várias abordagens e técnicas para a coleta de dados, dependendo das fontes de dados, dos objetivos do projeto e das tecnologias disponíveis.

Abaixo apenas algumas das muitas fontes e métodos de coleta de dados disponíveis. É importante escolher as abordagens mais adequadas às necessidades específicas do projeto, garantindo que os dados coletados sejam relevantes, precisos e legalmente obtidos. Além disso, é crucial considerar questões de privacidade e segurança durante todo o processo de coleta e armazenamento de dados.

Sensores e Dispositivos IoT

Muitos dados são gerados por dispositivos IoT (Internet das Coisas) e sensores em tempo real. Esses dispositivos podem incluir sensores de temperatura, umidade, GPS, acelerômetros, câmeras, entre outros. Eles geram continuamente dados que podem ser coletados e usados para uma variedade de finalidades, como monitoramento ambiental, análise de tráfego, saúde digital, entre outros.

Websites e Aplicações

Dados podem ser coletados de websites, aplicativos móveis e plataformas online por meio de ferramentas de análise da web, como Google Analytics. Essas ferramentas rastreiam o comportamento dos usuários, como cliques, visualizações de página, tempo gasto em cada página, entre outros, fornecendo insights valiosos sobre o engajamento do usuário e o desempenho do site ou aplicativo.

Redes Sociais

As redes sociais geram enormes volumes de dados constantemente, incluindo posts, curtidas, compartilhamentos, comentários e mensagens. Ferramentas de monitoramento de redes sociais podem ser usadas para coletar e analisar esses dados, permitindo às empresas entender melhor o sentimento do cliente, identificar tendências, realizar campanhas de marketing direcionadas, entre outros.

Transações Comerciais

Dados transacionais, como registros de vendas, compras, pagamentos e transações financeiras, são uma fonte valiosa de informações para as empresas. Sistemas de gestão de banco de dados (DBMS) são usados para coletar, armazenar e gerenciar esses dados, que podem ser analisados para identificar padrões de compra, detectar fraudes, otimizar operações, entre outros.

Logs de Servidores e Dispositivos

Logs de servidores, aplicativos e dispositivos de rede contêm informações detalhadas sobre atividades, eventos e falhas. Esses logs podem ser coletados e analisados para monitorar a integridade e o desempenho de sistemas de TI, identificar problemas de segurança, rastrear atividades de usuários, entre outros.

Evolução dos Dados

A evolução dos dados reflete a crescente complexidade e diversidade das informações geradas e armazenadas nas últimas décadas. Essa evolução reflete a necessidade de lidar com uma variedade cada vez maior de tipos e formatos de dados à medida que as tecnologias de informação avançam.

Enquanto os dados estruturados continuam sendo uma parte importante das operações comerciais, os dados semi-estruturados e não estruturados oferecem oportunidades adicionais para análise avançada, como processamento de linguagem natural, reconhecimento de padrões em imagens e áudio, e análise de sentimento em texto não estruturado.

Essa evolução dos dados reflete o ritmo acelerado das mudanças tecnológicas e sociais nas últimas décadas. À medida que entramos em uma era cada vez mais digital, a capacidade de entender, gerenciar e aproveitar o poder dos dados se torna essencial para o sucesso das organizações em todos os setores. As organizações estão cada vez mais explorando todas essas formas de dados para obter insights mais profundos e valiosos.

Dados Estruturados

- Os dados estruturados são organizados em um formato tabular com campos predefinidos e tipos de dados específicos.
- Eles são altamente organizados e seguem um esquema fixo, geralmente representado em bancos de dados relacionais.

- Exemplos comuns de dados estruturados incluem tabelas em bancos de dados SQL, onde cada linha representa uma entrada individual e cada coluna representa um atributo específico.

Dados Semi-Estruturados

- Os dados semi-estruturados não seguem um esquema rígido como os dados estruturados, mas ainda possuem alguma estrutura básica.
- Eles podem ser representados em diferentes formatos, como JSON (JavaScript Object Notation), XML (Extensible Markup Language) ou YAML (YAML Ain't Markup Language).
- Embora os dados semi-estruturados possam conter campos repetidos e aninhados, eles geralmente não têm um esquema formal definido.

Dados Não Estruturados

- Os dados não estruturados não seguem uma organização pré definida e não se encaixam em um formato tabular.
- Eles podem incluir texto livre, imagens, áudio, vídeo, e-mails, posts de redes sociais, entre outros.
- Diferentemente dos dados estruturados e semi-estruturados, os dados não estruturados podem ser mais desafiadores de analisar devido à falta de estrutura formal.

Grande Volume de Dados

A era do Big Data nos presenteou com um oceano de informações, mas navegar nesse mar tempestuoso não é tarefa fácil. Lidar com grandes volumes de dados apresenta diversos desafios que exigem soluções inovadoras e estratégias robustas. Vamos mergulhar em alguns dos principais problemas que podem surgir nesse contexto.

Armazenamento

Armazenar grandes volumes de dados requer uma infraestrutura de armazenamento escalável e robusta. O aumento do volume de dados pode levar a desafios relacionados ao gerenciamento de espaço em disco, custos de armazenamento e escalabilidade.

Processar grandes volumes de dados pode ser demorado e exigir recursos computacionais significativos. Algoritmos de processamento em massa precisam ser desenvolvidos para lidar com a análise eficiente de grandes conjuntos de dados.

O processamento paralelo e distribuído é muitas vezes necessário para lidar com a carga de trabalho, utilizando clusters de computadores ou sistemas distribuídos para processar dados em paralelo e acelerar o tempo de processamento.

Transmissão de Dados

Movimentar grandes volumes de dados entre sistemas e redes pode ser um desafio. A largura de banda limitada e a latência podem afetar a velocidade e a eficiência da transferência de dados, especialmente em ambientes distribuídos.

Soluções como compressão de dados, otimização de rede e protocolos de transferência de dados eficientes são necessárias para reduzir o tempo e os custos associados à transmissão de grandes volumes de dados.

Processamento

Analisar um tsunami de dados exige computadores poderosos e algoritmos eficientes. Processar e transformar dados brutos em informações úteis é um desafio técnico que requer expertise e infraestrutura robusta.

Gerenciamento e Limpeza de Dados

Com grandes volumes de dados, a qualidade dos dados pode se tornar um problema. É comum lidar com dados incompletos, duplicados, inconsistentes ou incorretos, o que pode afetar a precisão e a confiabilidade das análises.

Estratégias de limpeza e integração de dados precisam ser implementadas para garantir que os dados sejam limpos, padronizados e prontos para análise antes de serem usados para tomada de decisão.

Qualidade

Nem todo dado é um diamante. A qualidade dos dados coletados é fundamental para a confiabilidade das análises e insights. Garantir a qualidade e a consistência dos dados exige processos rigorosos de limpeza, organização e validação.

Segurança e Privacidade

Aumentar o volume de dados também aumenta os riscos de segurança e privacidade. Grandes conjuntos de dados contêm informações confidenciais e sensíveis que precisam ser protegidas contra acesso não autorizado, roubo ou violação.

É necessário implementar medidas robustas de segurança, como criptografia, controle de acesso e monitoramento de atividades, para garantir a integridade e a confidencialidade dos dados.

Interpretação

Transformar dados em insights acionáveis é como decifrar um mapa do tesouro. Extrair significado e valor dos dados requer habilidades analíticas avançadas e ferramentas de visualização que possibilitem a compreensão das informações por diferentes públicos.

Integração

Combinar dados de diferentes fontes, como rios que se juntam em um delta, pode ser um desafio. Integrar dados de forma eficiente e transparente é essencial para obter uma visão completa e holística do cenário.

Habilidades

Navegar no mar do Big Data exige uma tripulação experiente. Encontrar profissionais com habilidades técnicas e analíticas para lidar com os desafios dessa área é crucial para o sucesso de qualquer projeto.

Regulamentação

As leis que regem o uso de dados, como bússolas que guiam a navegação, precisam ser claras e eficazes. Garantir a conformidade com as regulamentações de proteção de dados e privacidade é essencial para evitar riscos e multas.

Custo

Construir e manter um navio para navegar no Big Data pode ser caro. O alto custo de infraestrutura, software e profissionais especializados pode ser um obstáculo para empresas e organizações com recursos limitados.

Ética

A ética deve ser a bússola que guia a navegação no mar de dados. É fundamental garantir que o uso do Big Data seja ético e responsável, respeitando a privacidade dos indivíduos e evitando vieses e discriminações.

Sistema de Banco de Dados Cassandra

Apache Cassandra é um sistema de gerenciamento de banco de dados distribuído, altamente escalável e altamente disponível, projetado para lidar com grandes volumes de dados em vários data centers e na nuvem.

O Cassandra é uma poderosa ferramenta para lidar com grandes volumes de dados distribuídos de forma eficiente e escalável. Sua arquitetura distribuída, flexibilidade de esquema e alto desempenho o tornam uma escolha popular para uma ampla gama de aplicativos e casos de uso em grandes empresas e organizações.

Cassandra é um sistema de banco de dados NoSQL¹ poderoso e versátil, ideal para gerenciar grandes volumes de dados em diversos cenários.

Aqui estão algumas características e aspectos importantes do Cassandra.

Principais características do Cassandra

- **Escalabilidade:** Horizontalmente escalável, o Cassandra pode ser facilmente expandido adicionando mais servidores à sua infraestrutura, sem comprometer o desempenho. É como construir um armazém cada vez maior para atender às suas necessidades crescentes.
- **Alta disponibilidade:** O Cassandra garante que seus dados estejam sempre disponíveis, mesmo em caso de falhas em um ou mais servidores. Sua arquitetura resiliente replica seus dados em vários servidores, garantindo a continuidade das operações.
- **Consistência ajustável:** O Cassandra oferece diferentes níveis de consistência, permitindo que você escolha o equilíbrio ideal entre performance e confiabilidade para suas necessidades específicas. É como ajustar a temperatura do seu armazém para garantir a preservação dos seus produtos.
- **Modelo de dados flexível:** O Cassandra utiliza um modelo de dados de coluna ampla, ideal para armazenar grandes quantidades de dados semiestruturados, como logs de sensores, dados de IoT e eventos de clickstream. Imagine um armazém versátil, capaz de armazenar diferentes tipos de produtos de forma organizada e eficiente.

¹ NoSQL, ou "Not Only SQL" (Não Apenas SQL), é um termo genérico utilizado para descrever um conjunto de sistemas de gerenciamento de banco de dados (SGBDs) que diferem dos tradicionais bancos de dados relacionais baseados em SQL. NoSQL surgiu como uma resposta às limitações dos bancos de dados relacionais em lidar com grandes volumes de dados não estruturados ou semiestruturados, alta escalabilidade, distribuição geográfica e flexibilidade de esquema.

- **Comunidade vibrante:** O Cassandra possui uma comunidade ativa e global de desenvolvedores e usuários que contribuem para o desenvolvimento e aprimoramento do sistema. É como ter uma equipe de especialistas à sua disposição para ajudá-lo a tirar o máximo proveito do Cassandra.
- **Modelo de Dados:** Cassandra é baseado em um modelo de dados NoSQL do tipo chave-valor. Os dados são organizados em tabelas, onde cada linha é identificada por uma chave primária única. O esquema das tabelas é flexível, permitindo a adição e remoção de colunas sem afetar o restante dos dados na tabela.
- **Distribuição de Dados:** Cassandra distribui automaticamente os dados através de vários nós em um cluster. Cada nó é responsável por uma parte do conjunto de dados, determinado pelo algoritmo de particionamento. Isso permite que Cassandra escale horizontalmente, adicionando novos nós ao cluster conforme necessário para lidar com mais dados e tráfego.
- **Consistência e Tolerância a Falhas:** Cassandra é projetado para ser altamente disponível e tolerante a falhas. Ele usa um modelo de consistência baseado em escrita, permitindo que os dados sejam gravados em vários nós de forma assíncrona. Os dados são replicados entre nós para garantir redundância e disponibilidade. Os administradores podem configurar o número de réplicas e a estratégia de replicação de acordo com os requisitos de disponibilidade e tolerância a falhas.
- **Leitura e Gravação de Dados:** Cassandra oferece suporte a leitura e gravação de dados de baixa latência, mesmo em grandes volumes de dados. Os dados podem ser lidos e gravados em qualquer nó do cluster, eliminando gargalos de desempenho e pontos únicos de falha.
- **Consultas:** Cassandra suporta consultas através de sua linguagem de consulta CQL (Cassandra Query Language), que é semelhante ao SQL. As consultas podem ser realizadas em chaves primárias, índices secundários e através de operações de filtragem. No entanto, Cassandra não suporta operações de junção entre tabelas, pois não segue um modelo relacional tradicional.
- **Escalabilidade Linear:** Uma das principais vantagens do Cassandra é sua capacidade de escalar linearmente com o número de nós no cluster. Adicionar mais nós ao cluster aumenta a capacidade de armazenamento e o desempenho de leitura/gravação de forma proporcional.
- **Uso em Casos de Uso:** Cassandra é amplamente utilizado em aplicativos que exigem alta disponibilidade, escalabilidade e tolerância a falhas, como redes sociais, análise de logs, sistemas de mensagens, aplicativos da Internet das Coisas (IoT) e muito mais.

Aplicações do Cassandra

- **Armazenamento de logs:** O Cassandra é ideal para armazenar grandes volumes de logs gerados por aplicativos e sistemas.
- **Análise de Big Data:** O Cassandra é uma excelente opção para analisar grandes conjuntos de dados em tempo real.
- **IoT:** O Cassandra é uma plataforma robusta para armazenar e gerenciar dados gerados por dispositivos IoT.
- **Comércio eletrônico:** O Cassandra é utilizado por grandes plataformas de comércio eletrônico para armazenar dados de produtos, clientes e pedidos.
- **Mídias sociais:** O Cassandra é usado por plataformas de mídias sociais para armazenar dados de usuários, posts e atividades.

Vantagens do Cassandra

- Altamente escalável
- Alta disponibilidade
- Consistência ajustável
- Modelo de dados flexível

- Comunidade vibrante
- Código aberto
- Desempenho superior
- Baixo custo

Desvantagens do Cassandra

- Complexidade de configuração e gerenciamento
- Não é ideal para transações ACID
- Curva de aprendizado mais acentuada

Os 5 V's da Big Data

Big data é caracterizado por cinco principais dimensões, os cinco "V's" do big data são volume, velocidade, variedade, veracidade e valor. Esses conceitos são fundamentais para compreender a natureza e os desafios associados ao processamento de grandes conjuntos de dados.

Volume

Refere-se à quantidade de dados gerados, armazenados e processados em um determinado período de tempo. O volume de dados em um ambiente de big data é geralmente extremamente grande, indo de gigabytes a petabytes e além. Esses dados podem ser provenientes de várias fontes, como transações comerciais, registros de atividades em redes sociais, dispositivos IoT (Internet das Coisas), sensores, dispositivos móveis, entre outros. O grande volume de dados apresenta desafios em termos de armazenamento, gerenciamento e processamento eficiente desses dados.

Velocidade

Refere-se à taxa na qual os dados são gerados, capturados, processados e analisados. Em um ambiente de big data, a velocidade com que os dados são produzidos e precisam ser processados pode ser extremamente alta. Por exemplo, dados gerados em tempo real, como transmissões de redes sociais, transações financeiras ou sensores industriais, exigem sistemas capazes de lidar com o processamento em tempo real para fornecer insights relevantes e acionáveis. A capacidade de lidar com grandes volumes de dados em tempo real é fundamental para muitos aplicativos de big data, como detecção de fraudes, monitoramento de desempenho em tempo real e análise de sentimentos em mídias sociais.

Variedade

Refere-se à diversidade de tipos e formatos de dados disponíveis em um ambiente de big data. Os dados podem ser estruturados, semi estruturados ou não estruturados e podem incluir texto, imagens, áudio, vídeo, logs de servidor, dados de sensores, entre outros. Essa diversidade de dados apresenta desafios adicionais para a análise e interpretação, pois os dados podem ser provenientes de fontes heterogêneas e apresentar diferentes estruturas e semântica. Lidar com essa variedade requer ferramentas e técnicas flexíveis de análise de dados que possam trabalhar com diferentes tipos e formatos de dados, permitindo a extração de insights valiosos independentemente de como os dados são estruturados.

Veracidade

Refere-se à qualidade dos dados, incluindo sua precisão, confiabilidade e integridade. Em um ambiente de big data, onde os conjuntos de dados podem ser vastos e heterogêneos, garantir a veracidade dos dados é crucial para garantir que as decisões tomadas com base nesses dados sejam precisas e confiáveis.

Valor

Refere-se à capacidade de extrair valor útil e insights significativos dos dados. Embora os dados em si sejam abundantes, o verdadeiro benefício do big data reside na capacidade de analisar esses dados de maneira eficaz para identificar padrões, tendências e correlações que possam ser usados para melhorar processos, desenvolver produtos e serviços, otimizar operações e tomar decisões estratégicas.

Como Obter Valor dos Dados

O valor real de um dado em Big Data se traduz na sua capacidade de gerar insights acionáveis que impactam positivamente a sua empresa. Para determinar se um dado possui valor real, considere os aspectos abaixo.

4 Regras Básicas

Exatidão

Os dados são valiosos se estiverem alinhados com os objetivos estratégicos e operacionais da organização. Se os dados podem ser utilizados para informar decisões, melhorar processos, otimizar operações ou impulsionar inovações que ajudam a alcançar metas de negócios, então eles têm valor real.

Consistência dos Dados

O dado é confiável e preciso? A qualidade e integridade dos dados são essenciais para determinar seu valor real. Dados precisos, completos e atualizados são mais confiáveis e, portanto, mais valiosos do que dados imprecisos, incompletos ou desatualizados.

Completeness dos Dados

Os dados devem fornecer informações completas sobre o tema em questão.

Relevância dos Dados

Os dados são valiosos se forem relevantes e úteis para resolver problemas ou responder a perguntas específicas. Se os dados fornecem informações significativas que podem levar a insights acionáveis ou ações tangíveis, então eles têm valor real.

Cases de Sucesso - O Mistério da Cerveja e das Fraldas - Walmart

Em 2008, a rede americana de supermercados Walmart se deparou com um enigma intrigante: um aumento inexplicável nas vendas de cerveja e fraldas nas noites de quinta-feira. Intrigados com essa peculiaridade, os analistas de dados da empresa mergulharam em um mar de informações para desvendar o mistério.

A Jornada Exploratória

Utilizando técnicas avançadas de análise de Big Data, os analistas vasculharam dados de compras, históricos de clientes, pesquisas online e até mesmo dados demográficos. A investigação revelou uma correlação surpreendente: nas noites de quinta-feira, os pais solteiros compravam cerveja e fraldas juntos.

A Teoria Revelada

A explicação para essa correlação reside no comportamento dos pais solteiros. Na noite de quinta-feira, após uma longa semana de trabalho, esses pais buscavam um momento de descontração com uma cerveja. Ao mesmo tempo, precisavam cuidar dos seus filhos, o que implicava em comprar fraldas.

O Impacto da Descoberta

Essa descoberta proporcionou à Walmart insights valiosos sobre o comportamento dos seus clientes. A empresa:

- Adaptou as promoções e ofertas para atender às necessidades dos pais solteiros.
- Otimizou o layout das lojas para facilitar a compra de cerveja e fraldas.
- Personalizou as campanhas de marketing para esse público específico.

Resultados Extraordinários

As ações tomadas pela Walmart resultaram em:

- Aumento significativo nas vendas de cerveja e fraldas nas noites de quinta-feira.
- Aumento da fidelidade dos pais solteiros à marca.
- Melhoria da percepção da empresa como um parceiro que entende as necessidades dos seus clientes.

Lições Aprendidas

O caso da Walmart demonstra o poder do Big Data para:

- Identificar padrões de comportamento: Descobrir insights valiosos sobre os clientes.
- Tomar decisões estratégicas: Implementar ações que impactam positivamente as vendas e a fidelização.
- Criar uma experiência personalizada: Atender às necessidades específicas de cada público.

A história da cerveja e das fraldas é um exemplo inspirador de como o Big Data pode ser utilizado para desvendar mistérios, gerar insights valiosos e impulsionar o sucesso de uma empresa.

Este case demonstra o potencial do Big Data para o varejo. É importante que as empresas se inspirem no sucesso do Walmart e invistam em soluções de Big Data para se manterem competitivas.

Atividades e áreas de atuação

O Big Data se configura como um universo em constante expansão, abrindo um leque de oportunidades para diversas áreas de atuação. As atividades e áreas de atuação do Big Data podem ser categorizadas em três pilares principais.

Geração de Insights

- **Análise de dados:** Extrair informações valiosas de grandes conjuntos de dados para entender o comportamento do cliente, identificar tendências de mercado e otimizar processos.
- **Machine learning:** Criar modelos preditivos para antecipar cenários futuros, auxiliar na tomada de decisões e personalizar ofertas.
- **Visualização de dados:** Transformar dados complexos em representações visuais intuitivas para facilitar a comunicação e o compartilhamento de insights.

Aplicações em Áreas Específicas

- **Varejo:** Personalizar ofertas, otimizar preços, prever demanda e melhorar a experiência do cliente.
- **Finanças:** Detectar fraudes, gerenciar riscos, otimizar investimentos e tomar decisões mais assertivas.
- **Saúde:** Diagnosticar doenças, personalizar tratamentos, prever epidemias e otimizar a gestão de recursos.
- **Indústria:** Otimizar processos de produção, prever falhas de equipamentos, reduzir custos e aumentar a eficiência.
- **Governo:** Combater crimes, melhorar a segurança pública, otimizar serviços públicos e tomar decisões baseadas em dados.

Habilidades Essenciais para o Sucesso

As áreas de atuação do Big Data estão em constante crescimento e novas aplicações são descobertas continuamente. É importante que as empresas invistam em profissionais qualificados para aproveitar ao máximo o potencial do Big Data.

Ciência de Dados

A ciência de dados é um campo interdisciplinar em franca expansão que combina matemática, estatística, computação e conhecimento de domínio para extrair insights valiosos de grandes conjuntos de dados. Essa área oferece um leque de oportunidades para diversos setores, desde o varejo e finanças até a saúde e indústria.

A ciência de dados é um campo em constante evolução, com novas oportunidades surgindo a cada dia. Profissionais qualificados nesta área são altamente requisitados no mercado de trabalho e podem contribuir significativamente para o sucesso de qualquer organização.

As principais atividades da ciência de dados se dividem em três etapas principais.

1. Análise de Dados

A análise de dados é um processo que transforma dados brutos em informações valiosas e acionáveis. Ela envolve a coleta, organização, limpeza, análise e visualização de dados para identificar padrões, tendências e insights que podem ser utilizados para tomar decisões mais inteligentes e estratégicas.

2. Análise de Tendências

A análise de tendências é um processo crucial para identificar padrões e previsões em dados, permitindo que empresas e indivíduos se preparem para o futuro com mais assertividade. Através da análise de dados

históricos, eventos atuais e projeções futuras, é possível tomar decisões estratégicas e aproveitar oportunidades antes que se tornem mainstream.

3. Análise por Algoritmos - Machine Learning

Os cientistas de dados aplicam técnicas estatísticas avançadas e algoritmos de machine learning para construir modelos que possam fazer previsões, classificações, agrupamentos e outras análises preditivas com base nos dados disponíveis.

Engenharias de Dados

A engenharia de dados é um campo interdisciplinar que combina conhecimentos de engenharia de software, ciência da computação e matemática para construir e gerenciar sistemas de dados escaláveis e confiáveis. A principal função da engenharia de dados é criar a infraestrutura necessária para que os cientistas de dados possam realizar seus trabalhos de forma eficiente e eficaz.

O engenheiro de dados deve possuir habilidades para construir e gerenciar infraestruturas de Big Data.

Responsabilidades de um Engenheiro de Dados

- **Desenvolver e gerenciar pipelines de dados:** Criar sistemas para coletar, limpar, transformar e integrar dados de diversas fontes.
- **Construir e gerenciar armazéns de dados:** Criar repositórios de dados centralizados e escaláveis para facilitar o acesso e a análise.
- **Desenvolver e gerenciar ferramentas de análise:** Criar ferramentas para que os cientistas de dados possam analisar os dados de forma eficiente.
- **Garantir a segurança e a qualidade dos dados:** Implementar medidas para garantir que os dados sejam seguros, confiáveis e precisos.
- **Colaborar com stakeholders:** Trabalhar em conjunto com cientistas de dados, analistas de negócios e outros stakeholders para entender suas necessidades e fornecer soluções adequadas.

Habilidades de um Engenheiro de Dados

- **Programação:** Linguagens de programação como Python, Java, Scala e SQL.
- **Engenharia de software:** Conhecimentos sobre princípios de engenharia de software para construir sistemas escaláveis e confiáveis.
- **Ciência da computação:** Conhecimentos sobre algoritmos, estruturas de dados e bancos de dados.
- **Matemática:** Conhecimentos básicos de matemática e estatística para entender e aplicar técnicas de análise de dados.
- **Comunicação:** Habilidade para comunicar ideias complexas de forma clara e concisa para stakeholders.

diferença entre cientista de dados e engenharia de dados

ETL - Extract Transform Load

ETL (Extract, Transform, Load) é um processo fundamental na engenharia de dados usado para mover e transformar dados de várias fontes para um sistema de destino, como um data warehouse, data lake ou banco de dados, onde podem ser analisados. Aqui está uma explicação de cada etapa do processo ETL.

Extract (Extração)

Nesta etapa, os dados são extraídos de uma ou várias fontes de dados. Isso pode incluir bancos de dados, sistemas de arquivos, APIs, feeds de streaming, entre outros. A extração pode envolver a leitura de grandes volumes de dados de fontes heterogêneas, como bancos de dados relacionais, bancos de dados NoSQL, arquivos CSV, JSON, XML, entre outros. Os dados extraídos são frequentemente em bruto e podem estar em diferentes formatos e estruturas.

Transform (Transformação)

Após a extração, os dados são transformados para garantir que estejam em um formato adequado para análise. Isso pode envolver várias operações de transformação, como limpeza de dados, padronização de formatos, conversão de tipos de dados, agregação, filtragem, enriquecimento com dados adicionais, e muito mais. O objetivo desta etapa é garantir que os dados sejam consistentes, completos e corretos antes de serem carregados no sistema de destino.

Load (Carregamento)

Na etapa final do processo ETL, os dados transformados são carregados no sistema de destino, onde podem ser armazenados e analisados. Isso geralmente envolve a inserção dos dados em tabelas de um data warehouse, data lake ou banco de dados. O carregamento pode ser feito de várias maneiras, incluindo inserção em lote, inserção em tempo real (streaming), ou utilizando técnicas de particionamento e paralelismo para otimizar o desempenho do carregamento.

O processo ETL é essencial para preparar e integrar dados de várias fontes para análise. Ele permite que as organizações obtenham insights valiosos a partir de grandes volumes de dados, garantindo que os dados sejam precisos, consistentes e prontos para análise. Nos últimos anos, com o aumento do volume e da complexidade dos dados, surgiram técnicas e ferramentas alternativas, como ELT (Extract, Load, Transform), que invertem a ordem das etapas de transformação e carregamento para melhor lidar com grandes volumes de dados brutos.

Machine Learning

O Machine Learning (ML) revolucionou a forma como as organizações lidam com dados, permitindo a automatização de tarefas complexas, a tomada de decisões baseada em dados e a criação de sistemas inteligentes capazes de aprender e melhorar com o tempo.

Machine Learning é um subcampo da inteligência artificial (IA) que se concentra no desenvolvimento de algoritmos e modelos computacionais que permitem aos sistemas aprenderem padrões a partir de dados e tomarem decisões sem intervenção humana. Em vez de serem explicitamente programados para executar uma tarefa específica, os sistemas de ML são treinados com exemplos de dados para aprender a realizar a tarefa de forma autônoma.

Através de algoritmos e modelos estatísticos, o ML processa grandes volumes de dados para identificar padrões e tomar decisões autonomamente. Essa tecnologia revolucionária transforma a maneira como vivemos, trabalhamos e tomamos decisões em diversos setores da sociedade.

Tipos de Machine Learning

Existem três tipos principais de abordagens de Machine Learning:

Aprendizado Supervisionado

Nesse tipo de aprendizado, os algoritmos são treinados com pares de entrada e saída, permitindo que o modelo aprenda a mapear entradas para saídas. Exemplos incluem classificação, regressão e detecção de anomalias.

Aprendizado Não Supervisionado

Aqui, os algoritmos são treinados com dados não rotulados e são deixados para encontrar padrões e estruturas nos dados por conta própria. Exemplos incluem clustering, redução de dimensionalidade e associação.

Aprendizado por Reforço

Nesse tipo de aprendizado, os algoritmos aprendem a partir de feedback em um ambiente dinâmico, ajustando seu comportamento para maximizar uma recompensa. Exemplos incluem jogos, robótica e controle de processos.

Aplicações de Machine Learning

As aplicações de Machine Learning são vastas e estão presentes em diversas indústrias e domínios, incluindo:

- **Saúde:** Diagnóstico médico, descoberta de medicamentos, monitoramento de pacientes.
- **Finanças:** Detecção de fraudes, análise de crédito, previsão de mercado.
- **Varejo:** Recomendação de produtos, previsão de demanda, otimização de preços.
- **Automotivo:** Condução autônoma, manutenção preditiva, personalização de veículos.
- **Marketing:** Segmentação de clientes, personalização de campanhas, análise de sentimentos.
- **Manufatura:** Controle de qualidade, otimização de processos e manutenção preventiva.

O Machine Learning continua a ser uma área emocionante e em rápido desenvolvimento, com o potencial de transformar radicalmente a forma como as organizações operam e os serviços que oferecem. No entanto, é importante abordar os desafios éticos e técnicos associados ao uso de ML para garantir que ele seja aplicado de maneira responsável e ética. Com o avanço da tecnologia e a conscientização sobre questões éticas, o Machine Learning tem o potencial de beneficiar a sociedade de maneiras profundas e significativas.

Área de Atuação Big Data

O Big Data, que se refere à coleta e análise de grandes conjuntos de dados, possui um amplo leque de aplicações em diversas áreas, impactando significativamente os negócios e a sociedade. As principais áreas de atuação onde o Big Data pode ser utilizado são:

Marketing e Vendas

- **Análise de comportamento do consumidor:** Compreender as preferências e hábitos dos clientes para personalizar ofertas, campanhas e produtos.
- **Segmentação de mercado:** Identificar grupos específicos de clientes com características e necessidades em comum para direcionar campanhas de marketing.
- **Otimização de preços:** Definir preços dinâmicos que se ajustam à demanda e ao perfil do cliente.
- **Detecção de fraudes:** Identificar transações fraudulentas em tempo real.

- **Análise de campanhas:** Avaliar a efetividade das campanhas de marketing e identificar áreas de otimização.

Finanças

- **Gestão de risco:** Identificar e mitigar riscos de crédito, mercado e operacional.
- **Deteção de fraudes:** Identificar transações fraudulentas em tempo real.
- **Análise de mercado:** Monitorar tendências do mercado e identificar oportunidades de investimento.
- **Gestão de investimentos:** Otimizar carteiras de investimento e tomar decisões mais inteligentes.
- **Desenvolvimento de novos produtos:** Criar produtos financeiros personalizados para atender às necessidades dos clientes.

Saúde

- **Diagnóstico de doenças:** Analisar dados de pacientes para identificar padrões e auxiliar no diagnóstico de doenças.
- **Desenvolvimento de novos medicamentos:** Identificar novos targets para o desenvolvimento de medicamentos.
- **Personalização de tratamentos:** Criar planos de tratamento personalizados para cada paciente.
- **Prevenção de doenças:** Identificar fatores de risco para doenças e desenvolver medidas preventivas.
- **Monitoramento de pacientes:** Monitorar a saúde dos pacientes em tempo real e identificar sinais de alerta.

Indústria

- **Otimização da produção:** Identificar gargalos na produção e otimizar processos para aumentar a eficiência.
- **Previsão de falhas de equipamentos:** Identificar sinais de falha em equipamentos e realizar manutenções preventivas.
- **Gestão da cadeia de suprimentos:** Otimizar a cadeia de suprimentos para reduzir custos e aumentar a eficiência.
- **Desenvolvimento de novos produtos:** Criar novos produtos que atendam às necessidades dos clientes e do mercado.
- **Melhoria da qualidade dos produtos:** Identificar e corrigir problemas de qualidade dos produtos.

Governo

- **Combate à criminalidade:** Analisar dados para identificar padrões de crime e desenvolver estratégias de combate.
- **Prevenção de fraudes:** Identificar fraudes em programas sociais e benefícios públicos.
- **Gestão de recursos públicos:** Otimizar a gestão de recursos públicos para melhorar a qualidade dos serviços prestados à população.
- **Tomada de decisões:** Tomar decisões mais inteligentes baseadas em dados concretos.
- **Melhoria da qualidade de vida da população:** Desenvolver políticas públicas que melhorem a qualidade de vida da população.

Outras áreas de atuação

- **Recursos Humanos:** Análise de dados para recrutamento, seleção, treinamento e desenvolvimento de funcionários.
- **Logística:** Otimização de rotas e entregas.
- **Educação:** Personalização do aprendizado e avaliação do desempenho dos alunos.
- **Energia:** Otimização do consumo de energia e desenvolvimento de fontes de energia renovável.
- **Telecomunicações:** Análise de dados para melhorar a qualidade dos serviços e identificar oportunidades de negócio.

O Big Data é uma ferramenta poderosa que pode ser utilizada para melhorar a eficiência, a produtividade e a tomada de decisões em diversas áreas. Ao utilizar o Big Data de forma inteligente, as empresas e organizações podem obter uma vantagem competitiva significativa e contribuir para o desenvolvimento da sociedade.