



MINERAÇÃO DE DADOS COMPLEXOS

Curso de Extensão



INF-0619 – PROJETO FINAL

Lista de Projetos

A avaliação da disciplina de Projeto Final (INF-0619) será constituída por um projeto alocado para cada grupo, considerando a lista de preferência de projetos que cada grupo deverá fornecer através do seu coordenador. Lembrando que não há garantia de que o projeto com maior preferência escolhido por um grupo será alocado para o mesmo. Os projetos disponíveis são:

- Análise de Avaliações de Livros
- Análise de Sentimento de Memes
- Análise do Sono
- Classificação de Atributos Faciais
- Classificação de Notícias Brasileiras
- Classificação de Sons Urbanos
- Detecção de Ações Humanas
- Indicadores de Doenças Cardíacas
- Predição de Atributos Visuais
- Predição de Seguro de Veículos
- Reconhecimento de Emoções em Diálogos

Observações

- A **Apresentação Parcial** será feita de forma oral (8 a 10 minutos), com auxílio de slides.
- A **Apresentação Final** será feita por meio de um vídeo de 8 a 10 minutos, preparado por cada grupo.
- Pontos importantes a serem abordados nas apresentações:
 - Descrição do problema e análise exploratória dos dados.
 - Preparação e normalização dos dados.
 - Técnicas utilizadas (baseline e soluções mais complexas).
 - Avaliação dos métodos.
 - Dificuldades encontradas, análise de casos de erro.
- Após as apresentações, os membros da banca avaliadora podem realizar perguntas ao grupo sobre o trabalho realizado.
- Complementando a Apresentação Final de cada projeto, um relatório de até 12 páginas (em formato PDF) deverá ser submetido no Moodle, pelo coordenador do grupo, junto com todos os arquivos fontes utilizados na realização do projeto.
- Os vídeos das Apresentações Finais serão publicadas no canal YouTube do curso (<http://bit.ly/youtube-mdc>).
- Qualquer tentativa de fraude implicará em nota zero na disciplina para todos os membros do grupo, sem prejuízo de outras sanções.
- O projeto final não pode ser compartilhado entre os grupos, o que caracteriza tentativa de fraude.
- Projeto final realizado com ajuda externa será considerado como uma tentativa de fraude.



ANÁLISE DE AVALIAÇÕES DE LIVROS

Descrição

Os comentários e avaliações (*reviews*) auxiliam outros usuários na escolha e compra de produtos em lojas *online*. Além do texto dos *reviews*, os usuários geralmente podem indicar uma nota, normalmente entre 0 e 5, da qualidade do produto. Neste projeto, a partir do comentário de um avaliador da plataforma de livros Goodreads, deve-se estimar qual a nota dada ao livro pela avaliação.

O dataset deste projeto está associado a uma competição do Kaggle. A base possui dados de avaliações, cada um com o texto do comentário, a nota da avaliação (entre 0 e 5) e os metadados do *review*. No conjunto de treinamento, existem cerca de 900.000 amostras, enquanto no teste há cerca de 478.000.

Objetivo Principal

Este projeto envolve um problema multi-classe, cujo objetivo é prever corretamente qual é a nota dada ao livro a partir da avaliação feita por um usuário. A avaliação oficial da competição é realizada através da métrica F1-score.

Técnicas Envolvidas

- Processamento de linguagem natural.
- Análise de dados.
- Classificação multi-classe.

Desafios

Para esse projeto, alguns desafios são:

- Analisar as palavras e termos chaves de cada classe.
- Analisar os meta-dados de cada *review*.
- Explorar técnicas de processamento de linguagem natural.
- Lidar com desbalanceamento das classes.
- Classificar corretamente a nota do livro a partir dos dados textuais.

Conjunto de Dados

- Kaggle da competição.
- Base de dados.



MINERAÇÃO DE DADOS COMPLEXOS

Curso de Extensão



ANÁLISE DE SENTIMENTO DE MEMES

Descrição

Com o crescente uso das redes sociais na última década, novos métodos de comunicação surgiram e se tornaram comuns. Entre eles, *Memes* são imagens ou vídeos, geralmente acompanhados de alguma frase curta. Apesar de frequentemente trazerem mensagens irônicas ou cômicas, algumas vezes carregam algum tipo de discurso de ódio que é replicado pelas redes. Identificar e reduzir o alcance deste tipo de mensagem é essencial para combater esses discursos e diminuir seu impacto na nossa sociedade.

Neste projeto, para cada meme composto por imagem e texto, deve-se classificar qual é o sentimento associado. O dataset deste projeto consiste em cerca de 7.000 memes, representados por uma imagem e um texto associado. Cada meme possui uma descrição entre cinco sentimentos possíveis, sendo eles muito negativo, negativo, neutro, positivo e muito positivo.

Objetivo Principal

Este projeto envolve um problema multi-classe, cujo objetivo é prever corretamente qual é o sentimento relacionado ao meme, considerando a imagem e o texto. A métrica de avaliação que deverá ser utilizada neste projeto é a Acurácia Balanceada.

Técnicas Envolvidas

- Processamento de linguagem natural.
- Processamento de imagens.
- Classificação multi-classe.

Desafios

Para esse projeto, alguns desafios são:

- Explorar técnicas de processamento de linguagem natural.
- Explorar técnicas de processamento de imagens.
- Analisar modelos multi-modais para a classificação multi-classe.
- Classificar corretamente o sentimento do meme.
- Explorar técnicas que permitam explicar quais elementos textuais e visuais são responsáveis pela classificação.

Conjunto de Dados

- Descrição da base de dados.
- Base de dados.



MINERAÇÃO DE DADOS COMPLEXOS

Curso de Extensão



ANÁLISE DO SONO

Descrição

A popularização de dispositivos *wearables*, como *smartwatches* e pulseiras *fitness*, permitiu a captura em tempo real de sinais biológicos que auxiliam no monitoramento da saúde de um usuário. Estes dispositivos são equipados com sensores que capazes de medir batimento cardíaco, níveis de oxigenação no sangue, intensidade de atividades físicas, e até a qualidade do sono de um indivíduo. O dataset deste projeto consiste nas medições feitas com sensores de Apple Watches de 31 pacientes durante exames PSG (*Polysomnography*) que avaliam as fases e a qualidade do sono. Para cada paciente, foram medidos os batimentos cardíacos, a aceleração e uma estimativa de seu relógio biológico através dos sensores do *smartwatch* durante todo o procedimento, pareando-os com a fase do sono (acordado, N1, N2, N3 ou REM) apontada pelo exame PSG.

Objetivo Principal

Este projeto envolve dois problemas multi-classe. O objetivo principal será treinar um modelo que indique se um paciente está acordado ou dormindo, a partir das medições do dispositivo. O objetivo secundário consiste em determinar a fase do sono de um paciente (acordado, N1, N2, N3 ou REM). A métrica que deverá ser utilizada neste projeto é a Acurácia Balanceada.

Técnicas Envolvidas

- Processamento de dados tabulares.
- Análise de séries temporais.
- Engenharia de características.
- Classificação binária e multi-classe.

Desafios

Para esse projeto, alguns desafios são:

- Explorar técnicas de processamento de dados tabulares.
- Combinar os dados pré-processados com os dados temporais.
- Lidar com o desbalanceamento das classes.
- Classificar corretamente a fase do sono e o estado do paciente.

Conjunto de Dados

- Descrição e informações gerais da base.
- Base de dados.



CLASSIFICAÇÃO DE ATRIBUTOS FACIAIS

Descrição

Nos últimos anos, modelos de *Deep Learning* obtiveram resultados impressionantes em diversos problemas que até recentemente eram abordados de forma manual por especialistas. Apesar das diversas aplicações que tais métodos apresentaram bons resultados — por exemplo, em contextos médicos, financeiros, forenses e vigilância — há uma crescente preocupação em garantir que estes modelos sejam justos e não reproduzam vieses indesejáveis em cenários cujas decisões afetem a vida de pessoas. Podemos perceber isso, por exemplo, em métodos que apresentam acurácias mais altas para um grupo étnico específico em detrimento de outros. O dataset deste projeto consiste em 108.501 imagens de rostos anotadas em relação a etnia, grupo etário e gênero. A base de dados foi criada com o intuito de avaliar os vieses existentes em modelos de aprendizado de máquina treinados para tarefas de detecção, reconhecimento, e verificação facial. Devido a representatividade de seus atributos, ela também permite o treinamento de modelos que reconheçam atributos faciais.

Objetivo Principal

Este projeto envolve três problemas multi-classe. O objetivo principal será treinar um ou mais modelos que classifiquem o gênero (*Male* ou *Female*), grupo étnico (*White*, *Black*, *Indian*, *East Asian*, *Southeast Asian*, *Middle Eastern*, ou *Latino*) e o grupo etário (*0-2*, *3-9*, *10-19*, *20-29*, *30-39*, *40-49*, *50-59*, *60-69*, ou *70+*) a partir de uma imagem de rosto de uma pessoa. Para as três tarefas, a métrica de avaliação que deverá ser utilizada é a Acurácia Balanceada.

Técnicas Envolvidas

- Processamento de imagens.
- Classificação multi-classe.

Desafios

Para esse projeto, alguns desafios são:

- Explorar técnicas de processamento de imagem.
- Lidar com o desbalanceamento de dados.
- Classificar corretamente cada atributo facial, reduzindo possíveis vieses entre os grupos.
- Explorar técnicas que permitam explicar quais regiões da face foram mais importantes na identificação de cada atributo facial.

Conjunto de Dados

- Informações gerais sobre o dataset.
- Descrição da base.
- Base de dados.



MINERAÇÃO DE DADOS COMPLEXOS

Curso de Extensão



CLASSIFICAÇÃO DE NOTÍCIAS BRASILEIRAS

Descrição

A classificação automática de assuntos de texto possui grande impacto na organização de notícias. A indicação de qual assunto a notícia se refere auxilia tanto na organização do veículo de entrega de notícias quanto na seleção de notícias de interesse para um determinado leitor. Neste projeto, a partir do título e do corpo da notícia, deve-se classificar o assunto associado.

O dataset consiste em cerca de 10.000 dados textuais de notícias do grupo Globo. Para cada uma, estão disponíveis o título, o texto e o assunto associado. Os assuntos são divididos entre *esporte*, *economia*, *política*, *tecnologia* e *famosos*.

Objetivo Principal

Este projeto envolve um problema multi-classe, cujo objetivo é prever corretamente qual o assunto referente ao título e texto da notícia. A métrica de avaliação que deverá ser utilizada neste projeto é a Acurácia Balanceada.

Técnicas Envolvidas

- Processamento de linguagem natural.
- Sanitização, organização e representação de dados textuais.
- Classificação multi-classe.

Desafios

Para esse projeto, alguns desafios são:

- Analisar as palavras e termos chaves de cada classe.
- Explorar técnicas de processamento de linguagem natural.
- Lidar com o desbalanceamento dos dados.
- Classificar corretamente o assunto da notícia a partir dos dados textuais.

Conjunto de Dados

- Descrição da base de dados.
- Base de dados.



MINERAÇÃO DE DADOS COMPLEXOS

Curso de Extensão



CLASSIFICAÇÃO DE SONS URBANOS

Descrição

A detecção e identificação automática de elementos sonoros urbanos é uma importante atividade com aplicações em diversos contextos, como recuperação multimídia e análise urbanística. Compreender quais são as principais fontes de sons em um bairro pode, por exemplo, nos ajudar a entender as características urbanas daquela região e auxiliar autoridades a planejar políticas de redução de poluição sonora. O dataset deste projeto consiste em 8.732 trechos de áudios de até 4 segundos em formato WAV, coletados em ambientes urbanos e disponíveis na plataforma www.freesound.org.

Objetivo Principal

Este projeto envolve um problema multi-classe, cujo objetivo é classificar corretamente qual a fonte do som em um trecho de áudio urbano. Há um total de 10 classes: ar condicionado, buzina de automóvel, crianças brincando, latido de cachorros, furadeiras, ruído de motor, disparos de armas de fogo, britadeiras, sirenes e música. A métrica de avaliação que deverá ser utilizada neste projeto é a Acurácia Balanceada.

Técnicas Envolvidas

- Processamento e representação de conteúdo de áudio.
- Classificação multi-classe.

Desafios

Para esse projeto, alguns desafios são:

- Explorar formas de representação do conteúdo de áudio (e.g., espectrograma de Mel) que explicitem as diferenças entre cada fonte. Um possível ponto de partida é a biblioteca python Librosa.
- Lidar com o desbalanceamento entre as classes.
- Investigar técnicas de aumento de dados para áudio (inserção de ruído, som ambiente, etc).
- Classificar corretamente as múltiplas fontes presentes nos trechos de áudio.

Conjunto de Dados

- Descrição da base de dados.
- Base de dados.



MINERAÇÃO DE DADOS COMPLEXOS

Curso de Extensão



DETECÇÃO DE AÇÕES HUMANAS

Descrição

Com a popularização de câmeras, *smartphones* e sistemas de vigilância, a detecção de ações humanas se tornou importante em diversas aplicações, tais como monitoramento, reconhecimento de gestos e ações. Neste projeto, para cada imagem da base de dados, deve-se classificar qual ação está sendo representada na imagem.

O dataset deste projeto consiste em 15.000 imagens para treinamento e 3.000 imagens para teste. Cada imagem representa uma de quinze tipos de ações, como, por exemplo, comendo, bebendo, lutando, correndo e dormindo.

Objetivo Principal

Este projeto envolve um problema multi-classe, cujo objetivo é prever corretamente qual é a ação realizada em cada imagem da base de dados. A métrica de avaliação que deverá ser utilizada neste projeto é a Acurácia Balanceada.

Técnicas Envolvidas

- Processamento de imagens.
- Classificação multi-classe.

Desafios

Para esse projeto, alguns desafios são:

- Explorar técnicas de processamento de imagem.
- Lidar com o desbalanceamento de dados.
- Classificar corretamente a ação em cada imagem.
- Explorar técnicas que permitam explicar quais elementos visuais são responsáveis pela classificação.

Conjunto de Dados

- Base de dados.



MINERAÇÃO DE DADOS COMPLEXOS

Curso de Extensão



INDICADORES DE DOENÇAS CARDÍACAS

Descrição

A análise de indicadores de problemas cardíacos possui grande importância no diagnóstico precoce e na prevenção de doenças. Segundo o Centro de Controle e Prevenção de Doenças (em inglês, *Centers for Disease Control and Prevention*, CDC), a agência de controle e prevenção de doenças do Departamento de Saúde dos Estados Unidos, problemas cardíacos são a principal causa de mortes nos Estados Unidos, com cerca de 700.000 casos de morte em 2020. Neste projeto, para cada amostra, deve-se classificar se o paciente é propenso a desenvolver problemas cardíacos.

O dataset deste projeto consiste em cerca de 320.000 dados que possuem informações e anotações reportando se um paciente tem doenças cardíacas. As características disponíveis estão relacionadas à saúde do paciente (tais como IMC, presença de diabetes, se já teve um infarto e dificuldades de mobilidade), comportamentais (se bebe bebidas alcoólicas, se fuma e seu tempo de sono) e sobre a pessoa (idade, sexo, etnia, entre outras).

Objetivo Principal

Este projeto envolve um problema binário, cujo objetivo é prever corretamente se uma pessoa terá ou não problemas cardíacos. A métrica que deverá ser utilizada neste projeto é a Acurácia Balanceada.

Técnicas Envolvidas

- Processamento de dados tabulares.
- Análise de dados.
- Engenharia de características.
- Classificação binária.

Desafios

Para esse projeto, alguns desafios são:

- Explorar técnicas de processamento de dados tabulares.
- Lidar com o desbalanceamento das classes.
- Identificar as características associadas aos problemas cardíacos.
- Classificar corretamente se uma pessoa terá ou não problemas cardíacos.

Conjunto de Dados

- Descrição da base de dados.
- Base de dados.



PREDIÇÃO DE ATRIBUTOS VISUAIS

Descrição

Vivemos em um mundo visualmente dinâmico, no qual a aparência de um local pode mudar drasticamente em algumas horas ou de uma estação do ano para outra. Por exemplo, um parque ao meio dia de um verão apresenta características visuais muito diferentes do mesmo local capturado ao final de uma tarde de inverno. Identificar tais alterações visuais em uma cena permite adaptar um sistema automaticamente à passagem do tempo ou a mudanças climáticas. O dataset deste projeto consiste em 8.571 imagens de cenas ao ar livre, anotadas com 40 atributos visuais em relação à aparência do local. Os atributos visuais se dividem entre aqueles que representam o período do dia (como manhã, noite, nascer/por do sol), condições climáticas (por exemplo, frio, ensolarado, com neblina), iluminação (tais como claro, escuro), vegetação, estação do ano, cores, e até características subjetivas (como se a cena é estressante, calma, bonita).

Objetivo Principal

Este projeto envolve um problema de classificação multi-rótulo, cujo objetivo é estimar a presença ou ausência de cada um dos 40 atributos visuais relacionados com a aparência de uma cena. A métrica de avaliação que deverá ser utilizada neste projeto é o erro quadrático médio (*Mean Squared Error*, MSE).

Técnicas Envolvidas

- Processamento de imagens.
- Classificação multi-rótulo.

Desafios

Para esse projeto, alguns desafios são:

- Explorar técnicas de processamento de imagens.
- Explorar relações entre os atributos visuais.
- Explorar técnicas que permitam explicar quais elementos na cena são responsáveis pela classificação.

Conjunto de Dados

- Descrição da base.
- Exemplos visuais de cada atributo.
- Base de dados.



MINERAÇÃO DE DADOS COMPLEXOS

Curso de Extensão



PREDIÇÃO DE SEGURO DE VEÍCULOS

Descrição

A análise de clientes é uma importante área de estudo para empresas de seguro. Este tipo de técnica tem como objetivo identificar características importantes para direcionar propagandas e ações de captação de potenciais novos clientes. O dataset deste projeto consiste cerca de 381.000 dados com informações de possíveis novos clientes. Para cada amostra da base de dados, existem informações sobre o cliente (tais como idade, gênero, se ele possui carteira de habilitação e o código da região do cliente), veículo (idade do veículo, se o veículo já teve algum acidente anteriormente) e em relação ao seguro.

Objetivo Principal

Este projeto envolve um problema binário, cujo objetivo é prever corretamente se um cliente irá querer ou não um seguro veicular. A métrica que deverá ser utilizada neste projeto é a Acurácia Balanceada.

Técnicas Envolvidas

- Processamento de dados tabulares.
- Análise de dados.
- Engenharia de características.
- Classificação binária.

Desafios

Para esse projeto, alguns desafios são:

- Explorar técnicas de processamento de dados tabulares.
- Analisar os dados da base de dados.
- Lidar com o desbalanceamento das classes.
- Identificar as características importantes para o cliente querer comprar o seguro veicular.
- Classificar corretamente se um cliente irá querer o seguro veicular.

Conjunto de Dados

- Descrição da base de dados.
- Base de dados.



MINERAÇÃO DE DADOS COMPLEXOS

Curso de Extensão



RECONHECIMENTO DE EMOÇÕES EM DIÁLOGOS

Descrição

A identificação da emoção de um interlocutor é uma habilidade importante para assistentes pessoais (como Alexa, Google Assistant e Siri), que permite que seus comportamentos sejam adaptados dependendo do estado emocional do usuário. A medida que novos tipos de dispositivos são integrados aos assistentes, podemos explorar o que é dito, a forma como o interlocutor fala e também suas expressões faciais para reconhecer seu estado emocional. O dataset deste projeto consiste de 13.000 frases coletadas do seriado *Friends*, anotadas em relação à emoção da fala (neutro, raiva, nojo, medo, felicidade, tristeza e surpresa) e ao seu sentimento (positivo, neutro ou negativo). Além da transcrição da fala, há também o trecho em vídeo da fala retirado do seriado.

Objetivo Principal

Este projeto envolve dois problemas multi-classe. O objetivo principal é prever corretamente qual é a emoção relacionada à frase a partir do texto, imagem e áudio. O objetivo secundário consiste em prever o sentimento da frase. A métrica de avaliação que deverá ser utilizada em ambas as tarefas é a Acurácia Balanceada.

Técnicas Envolvidas

- Processamento de linguagem natural.
- Processamento de imagens.
- Processamento de áudio.
- Classificação multi-classe.

Desafios

Para esse projeto, alguns desafios são:

- Analisar o texto das frases e sua relação com as emoções.
- Explorar métodos de representação e treinamento de modelos para linguagem natural.
- Explorar técnicas de reconhecimento de emoções através das imagens faciais.
- Explorar métodos de representação e treinamento de modelos a partir do sinal de áudio.
- Explorar abordagens multi-modais para a classificação multi-classe.
- Lidar com o desbalanceamento entre as classes.

Conjunto de Dados

- Descrição da base de dados.
- Base de dados.