

COE241 - Estatística e Modelos Probabilísticos

COS868 - Probabilidade e Estatística para Aprendizado de Máquina

Segundo Semestre de 2022 - Professora: Rosa Maria Meri Leão

Projeto do Curso

1 Dataset

O objetivo deste trabalho é analisar um conjunto de dados aplicando a teoria aprendida em classe. É muito importante que seja realizada uma análise crítica dos resultados encontrados.

O projeto será baseado em dados reais fornecidos por um provedor de Internet de médio porte. Os dados representam a taxa de dados enviados em bps (taxa de upload) e a taxa de dados recebidos em bps (taxa de download) de/por um dispositivo na casa de um usuário do provedor. Dois tipos de dispositivos devem ser analisados: Smart-TV e Chromecast.

Dois arquivos contendo os dados serão disponibilizados: (1) `dataset_chromecast.csv` e (2) `dataset_smart-tv.csv`. Cada arquivo possui os seguintes campos: **device_id**, **date_hour**, **bytes_up**, **bytes_down**. O campo **device_id** é o identificador único do dispositivo que está na casa de um usuário do provedor. O campo **date_hour** é a data e a hora em que a coleta de dados foi realizada com a granularidade de 1 minuto. Os campos **bytes_up** e **bytes_down** contém, respectivamente, as taxas em bps dos dados enviados (taxa de upload) e dos dados recebidos (taxa de download) pelo dispositivo em um minuto.

Atenção reescalonar dados para log 10: como os valores das taxas de upload e de download variam diversas ordens de grandeza, para calcular as estatísticas é necessário reescalonar as taxas para log na base 10. Por exemplo, se o campo **bytes_up** é igual a 1000, o valor que você deve usar para fazer as análises é 3.

2 Estatísticas gerais

O objetivo desse estudo é avaliar os dados sem considerar o horário em que foram gerados, ou seja, você deve considerar todos os dados de cada um dos arquivos para obter as estatísticas descritas a seguir. Para cada tipo de dispositivo, Smart-TV e Chromecast, calcular: Histograma, Função Distribuição Empírica, Box Plot, Média, Variância e Desvio Padrão, para a taxa de upload e taxa de download.

Lembre-se que o tamanho do *bin* deve ser estimado de forma que o histograma represente de forma adequada os dados estudados. Use o método de Sturges apresentado em aula para estimar o tamanho do *bin*.

Comente o que você observou a partir dos gráficos e sobre as diferenças e/ou similaridades entre os dois tipos de dispositivos. É importante interpretar os resultados obtidos.

3 Estatísticas por horário

O objetivo dessa análise é avaliar os dados considerando o horário em que foram gerados independente do dia. Você deve considerar os dados coletados em cada hora para cada tipo de dispositivo para obter as estatísticas descritas a seguir. Para cada tipo de dispositivo, Smart-TV e Chromecast, para cada hora calcular: Box Plot, Média, Variância e Desvio Padrão, para a taxa de upload e taxa de download. Neste item você deve gerar um box plot para cada tipo de dispositivo, para cada taxa coletada (upload e download), para cada hora. Para a média, variância e desvio padrão, você deve fazer 4 gráficos, representando no eixo X a hora e no eixo Y os valores das três estatísticas para cada taxa coletada, para cada tipo de dispositivo.

Faça uma análise dos resultados obtidos. Comente sobre as diferenças ou similaridades entre os dois tipos de dispositivos. O que você pode concluir a respeito das estatísticas obtidas por horário?

4 Caracterizando os horários com maior valor de tráfego

Neste item o objetivo é parametrizar uma ou mais distribuições da literatura para os dois horários com maior valor da mediana/média de cada taxa coletada para cada tipo de dispositivo. O **Passo 1** é escolher, a partir dos gráficos da seção 3, dois horários: um com

maior valor de mediana e o outro com maior valor da média, para a taxa de upload e taxa de download, para cada tipo de dispositivo: Smart-TV e Chromecast.

Você terá 8 datasets, cada um contendo os dados com as seguintes características:

- Dataset 1: Horário com a maior mediana da taxa de upload em uma hora, Smart-TV
- Dataset 2: Horário com a maior média da taxa de upload em uma hora, Smart-TV
- Dataset 3: Horário com a maior mediana da taxa de download em uma hora, Smart-TV
- Dataset 4: Horário com a maior média da taxa de download em uma hora, Smart-TV
- Dataset 5: Horário com a maior mediana da taxa de upload em uma hora, Chromecast
- Dataset 6: Horário com a maior média da taxa de upload em uma hora, Chromecast
- Dataset 7: Horário com a maior mediana da taxa de download em uma hora, Chromecast
- Dataset 8: Horário com a maior média da taxa de download em uma hora, Chromecast

No **Passo 2**, faça um histograma para cada um dos 8 datasets. Lembre-se que você deve escolher o tamanho do *bin* usando o método de Sturges.

No **Passo 3** calcule o maximum likelihood estimator (MLE) para estimar os parâmetros das seguintes distribuições: Gaussiana e Gamma, para cada um dos 8 datasets. Explique como você calculou o MLE.

No **Passo 4** você deve fazer um gráfico para cada um dos 8 datasets contendo 3 curvas: o histograma, a função densidade Gaussiana com os parâmetros obtidos com o MLE e a função densidade Gamma com os parâmetros obtidos com o MLE. Observando os gráficos você deve comentar se existe ou não uma variável aleatória da literatura que possivelmente possa ser usada para representar os dados de cada um dos 8 datasets.

O **Passo 5** consiste em fazer o gráfico *Probability Plot* comparando os dados de cada dataset com as distribuições parametrizadas. No total são 16 gráficos, comparando os dados reais dos 8 datasets com cada uma das duas distribuições parametrizadas.

A partir dos resultados dessa seção você deve analisar as seguintes questões:

1. Quais foram os horários escolhidos para cada dataset?

2. O que você pôde observar a partir dos histogramas dos datasets?
3. Comente sobre as diferenças e/ou similaridades entre os datasets 1 e 2, 3 e 4, 5 e 6, 7 e 8. O objetivo é comparar as características dos datasets com a maior mediana em um determinado horário com os datasets com a maior média em um determinado horário.
4. É possível caracterizar os datasets acima por uma variável aleatória da literatura?
5. Se a resposta for não, qual o motivo?
6. O que você pôde observar a partir dos gráficos *Probability Plot*?

5 Análise da correlação entre as taxas de upload e download para os horários com o maior valor de tráfego

Você deve calcular o coeficiente de correlação amostral e fazer o gráfico scatter plot comparando as taxas dos seguintes datasets: **dataset 1 e dataset 3, dataset 2 e dataset 4, dataset 5 e dataset 7, dataset 6 e dataset 8**. Você deve analisar os resultados e indicar se existe alguma correlação entre as taxas de download e upload dos dispositivos.

6 Comparação dos dados gerados pelos dispositivos Smart-TV e Chromecast

O objetivo deste estudo é avaliar se os dois dispositivos que são usados prioritariamente para assistir vídeo, possuem distribuição de probabilidade das taxas de upload e download semelhante nos horários de maior tráfego. Você deve usar a estatística G do teste Chi-Square for goodness of fit (apresentado em classe) para fazer essa análise. Esse teste também é chamado de G-test (<https://en.wikipedia.org/wiki/G-test>).

Você deve realizar o G-test comparando as taxas dos seguintes datasets: **dataset 1 e dataset 5, dataset 2 e dataset 6, dataset 3 e dataset 7, dataset 4 e dataset 8**.

Note que para realizar o G-test entre cada par de datasets, o número de bins e os valores dos bins dos histogramas dos datasets a serem comparados devem ser os mesmos. Portanto

se os bins dos histogramas que você obteve na seção 4 para um determinado par de datasets for diferente, escolha o número de bins e os valores dos bins de um deles e use como referência para aquele par de datasets.

7 Relatório

Você deve fazer um relatório contendo todos os resultados que você obteve e explicando como você os obteve. É importante comentar cada um dos resultados e explicar se o resultado que você obteve poderá auxiliar o provedor de serviço de Internet a entender os dados que passam pela sua rede. A avaliação do projeto será feita com base na qualidade do relatório.

Você deve fazer upload do seu relatório (arquivo pdf) na plataforma do Google Classroom na atividade Projeto do Curso.

No relatório deve estar indicado um link para o código que você usou para obter os resultados do trabalho.