

Projeto CTTU

Felipe Bormann¹

Recife, Brazil

Abstract

Este é o relatório final para o projeto de Ciências dos dados na Universidade de Federal de Pernambuco. O tema deste trabalho é a análise e exploração de dados de acidentes, infrações e monitoramento de trânsito na cidade do Recife.

Keywords: dados, CTTU, trânsito, data science, geodata

1. Introdução

Dados. Os dados utilizados neste projeto foram todos extraídos do site de dados públicos da Prefeitura do Recife. Busquei utilizar somente os dados da entidade CTTU.

5 *Objetivo.* O objetivo deste projeto é explorar as relações e características dos dados da CTTU para extrair conhecimento, principalmente em relação aos acidentes e infrações cometidas por condutores. Os datasets possuem informações das seguintes entidades:

- Acidentes
- 10 • Infrações
- Equipamentos de Monitoramento
- Semafóros
- Equipamentos de Vigilância

[☆]Fully documented templates are available in the elsarticle package on CTAN.

¹Since 1880.

Tecnologia. A tecnologia usada durante o projeto foi a linguagem Python e as
15 seguintes bibliotecas do seu sistema: Pandas, GeoPandas, Matplotlib, Bokeh e
Shapely. Além da criação de um notebook jupyter para a execução do projeto.

2. Limpeza dos dados

Acidentes. A limpeza dos dados dos acidentes foi relacionado ao fato de que
houve uma transformação na forma como os dados eram armazenados entre o
20 ano de 2014 e 2015, sendo 2016 bem parecido com 2015.

Em relação a datas, enquanto 2014 seguia o formato mês/dia/ano, os anos
de 2015 e 2016 seguiam o formato dia/mês/ano.

Outra diferença gritante estava no nome de colunas que representavam o
mesmo dado, como: "tipo" e "tipo de abertura" ou "quantidade_vitimas" e
25 "quantidade de vitimas".

Nesta variável, "quantidade de vítimas", foi necessário remover dados in-
válidos como "-", "VT" ou "F", como não havia especialista e a quantidade de
dados errados não era tão alta, decidi remover todos estes valores ao invés de
atribuir 0 por pura intuição.

30 Outra ocasião foi a junção de diferentes valores da variável "tipo de ocorrên-
cia" para um valor único, como o conjunto "ATROPELAMENTO", "atropela-
mento" e "Atropelamento", todos se tornando um único valor: "atropelamento",
facilitando o resultado.

Infrações. A primeira ação em relação às infrações foi detectar quais colunas
35 eram do tipo datetime e que podiam ser convertidas para este tipo em Pandas.
Além da extração da criação da coluna "numero semaforo" para todos os campos
que possuíam o número do semáforo associado em sua descrição.

Para as infrações, que não possuíam um ponto geográfico(latitude, longi-
tude) associado, utilizei da informação de ter um semáforo associado à infração,
40 desta forma reduzi o dataset de 1.6 milhões para 890 mil, e associar à posição
geográfica do semáforo à infração, dessa forma concendendo um espaço para

novas inferências baseadas em granularidades maiores e geográficas como Setor censitário, bairro, etc.

Semáforos. A única coluna que possuía dados inválidos em semafóro foi a do
45 tipo de funcionamento, que possuía erros de digitação mas nenhum valor nulo.

3. Análise Exploratória

Acidentes. Em relação à quantidade de acidentes, consegui algumas evidências de que menos acidentes acontecem ao serem instalados equipamentos de monitoramento. Em números absolutos: somente 52 dos 4344 acidentes registrados
50 foram à 100m de distância de um equipamento de monitoramento.

E a figura 1 mostra que há um aumento drástico de acidentes próximo aos equipamentos de monitoramento durante os finais de semana, que são justamente quando estes são desligados.

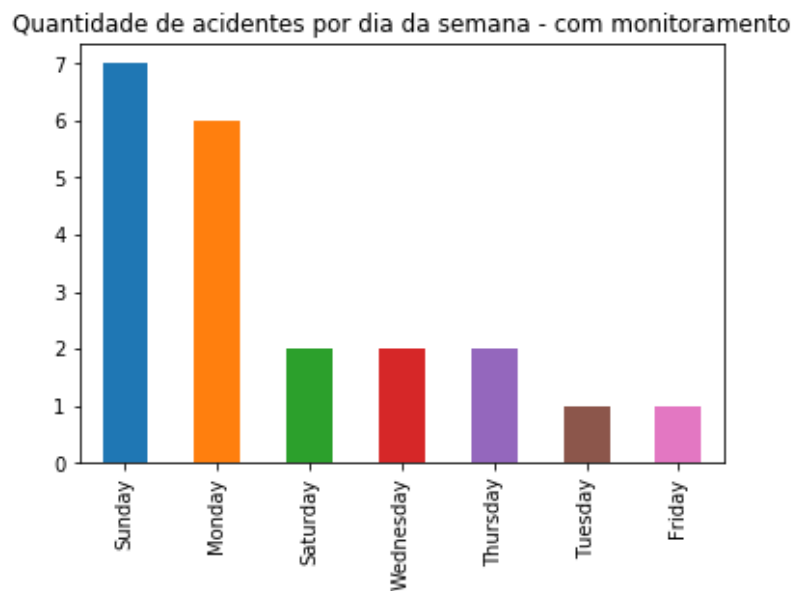


Figure 1: Quantidade de acidentes próximos à equipamentos de monitoramento

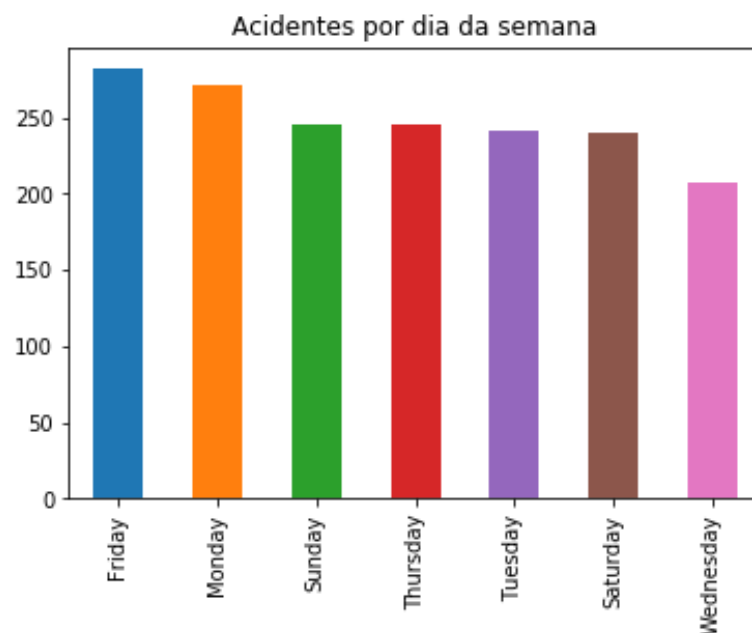


Figure 2: Quantidade de acidentes por dia da semana

Infrações. Procurei por variáveis que fossem interessantes e tentei relacioná-las com as variáveis de tempo em alguns casos ou explorar informações e insights simples que o dataframe poderia dar, como na Figura 3, onde fui capaz de tentar começar a investigar, apesar de sem sucesso, se onde haviam mais multas onde a velocidade estava acima entre 20 e 50 por cento acima do permitido causava mais mortes do que "somente" 20 por cento acima.

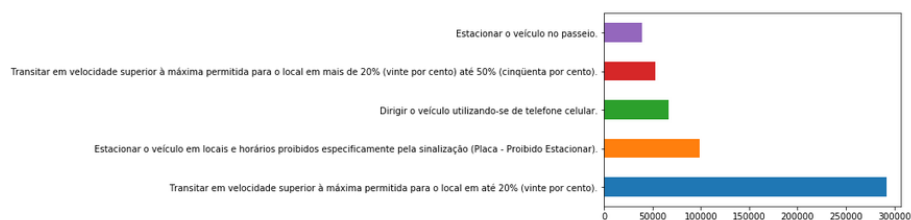


Figure 3: Tipo de Infrações mais frequentes

60 *Infrações por semana.* Outra descoberta foi, ao tentar correlacionar a semana do ano entre acidentes e infrações, não cheguei à resposta desejada, dado que não há correlação entre as duas curvas, como mostra, os dois gráficos 4 e 5 mas consegui descobrir que há um aumento crescente até à 43^a semana do ano, em outubro, onde o número de multas vem de uma crescente desde da 39^a semana.

65 Não consegui confirmar a partir de fontes oficiais o motivo de tal incentivo mas é notório o crescimento como mostra a figura 4.

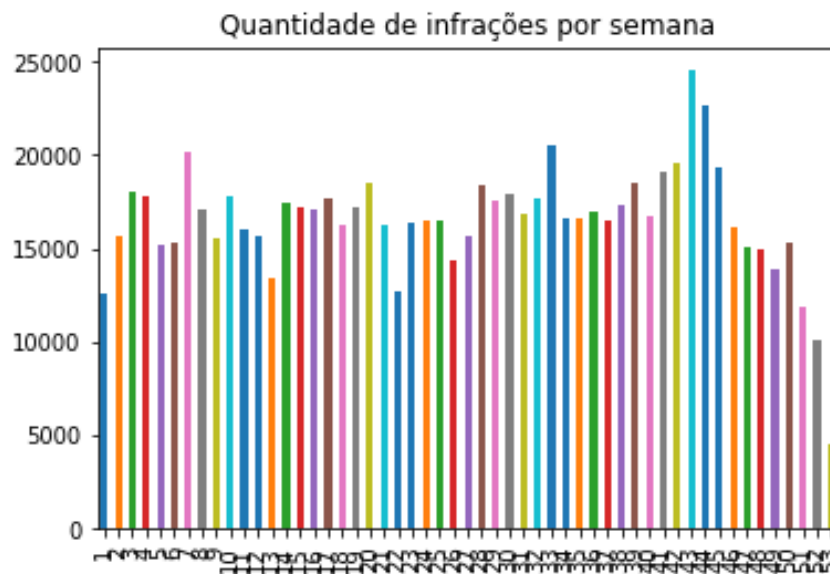


Figure 4: Quantidade de infrações por semana no ano

4. Conclusão

Consegui alcançar o objetivo de gerar alguns insights e visualizações que extraíam alguma informação, contudo, a tarefa se tornou mais complicada do

70 ponto de vista de pré-processamento e análise de dados geográficos, ao ter de aprender sobre visões geográficas, distâncias que fizessem sentido e à busca de informações não tão fáceis do governo, como o sistema de bonificação dos funcionários da CTTU, por exemplo.

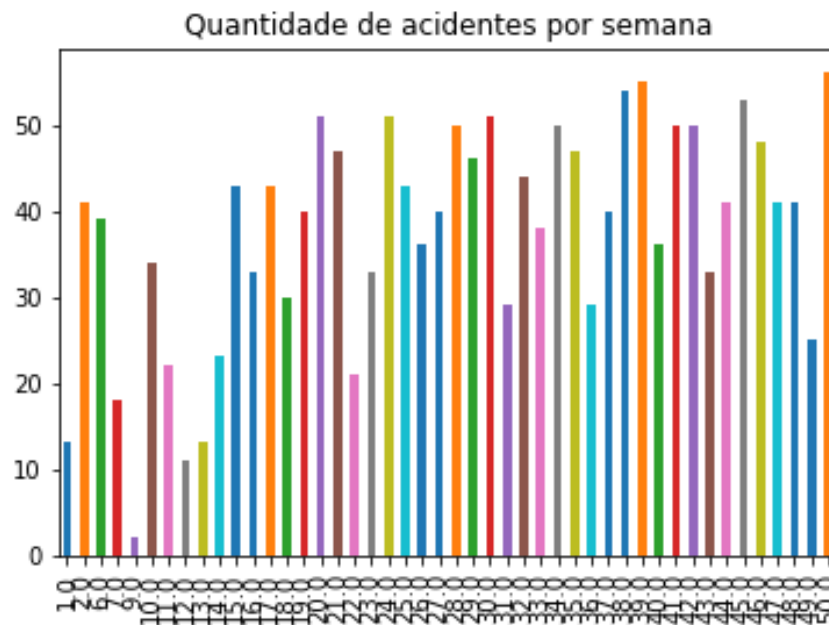


Figure 5: Quantidade de Acidentes por semana do ano

Objetivos Futuros.

- 75
- Gerar novas visualizações a partir dos dados de acidentes
 - Gerar modelo de predição para a variável "houve um vítima?" em um acidente
 - Correlacionar novas informações sobre a influência de infrações na ocorrência de acidentes