

DESAFIO CIVITAS
EMD

DETECÇÃO DE PLACAS CLONADAS NO
RIO DE JANEIRO

DESAFIO TÉCNICO PARA A VAGA NA PREFEITURA DO RIO DE
JANEIRO

AMANDA AMARO

amandacsamaro@gmail.com

RIO DE JANEIRO, 01 DE JULHO 2024

INTRODUÇÃO

O estudo analisa dados de leitura de radares para identificar possíveis casos de clonagem de placas de veículos no Rio de Janeiro. A clonagem de placas é uma prática ilegal que resulta em multas injustas e fraudes.

Objetivo do Estudo

O objetivo é identificar possíveis placas clonadas, abordando inconsistências nos dados de leitura dos radares. O foco será em:

- Limpeza dos dados
- Verificar a consistência geográfica das leituras.
- Detectar anomalias nas velocidades capturadas.
- Identificar placas associadas a múltiplos tipos de veículos.
- Definição de um modelo para identificar placas clonadas

Foi utilizado SQL para processamento inicial dos dados no BigQuery e Python (Pandas) para análises e visualizações adicionais. A metodologia desenvolvida visa garantir a precisão dos resultados ao obter possíveis placas suspeitas e fornecer insights para melhorar a fiscalização.

METODOLOGIA

Para a análise exploratória dos dados, utilizou-se Python com a biblioteca Pandas para carregar, inspecionar e manipular os dados. Os dados foram analisados em várias etapas para entender sua estrutura e identificar possíveis inconsistências que poderiam influenciar na análise de clonagem de placas.

Limpeza dos Dados

O banco contém 36.358.536 registros ([Código 2.1](#)), com a coluna `datahora_captura` apresentando 1.816.325 valores nulos ([Código 2.2](#)). Considerando que a integridade dos dados é crucial para a confiabilidade da análise e que se trata de 4,99% dos dados, para este trabalho eles foram desconsiderados ([Código 2.3](#)).

```
1 SELECT COUNT(*) AS total_records
2 FROM `rj-cetrio.desafio.readings_2024_06`;
```

Código 2.1: *Análise de quantidade de dados.*

```
1 SELECT
2     COUNTIF(datahora IS NULL) AS datahora_nulls,
3     COUNTIF(datahora_captura IS NULL) AS datahora_captura_nulls,
4     COUNTIF(placa IS NULL) AS placa_nulls,
5     COUNTIF(empresa IS NULL) AS empresa_nulls,
6     COUNTIF(tipoveiculo IS NULL) AS tipoveiculo_nulls,
7     COUNTIF(velocidade IS NULL) AS velocidade_nulls,
8     COUNTIF(camera_numero IS NULL) AS camera_numero_nulls,
9     COUNTIF(camera_latitude IS NULL) AS camera_latitude_nulls,
10    COUNTIF(camera_longitude IS NULL) AS camera_longitude_nulls
11 FROM `rj-cetrio.desafio.readings_2024_06`;
```

Código 2.2: *Análise de dados nulos.*

```
1      SELECT *
2      FROM `rj-cetrio.desafio.readings_2024_06`
3      WHERE
4          datahora IS NOT NULL AND
5          datahora_captura IS NOT NULL AND
6          placa IS NOT NULL AND
7          empresa IS NOT NULL AND
8          tipoveiculo IS NOT NULL AND
9          velocidade IS NOT NULL AND
10         camera_numero IS NOT NULL AND
11         camera_latitude IS NOT NULL AND
12         camera_longitude IS NOT NULL;
```

Código 2.3: Seleção de dados nulos distintos.

Todas essas análises e mapas podem ser vistos com mais detalhes no notebook disponibilizado no repositório do GitHub ou direto pelo Collab.

Análise Exploratória

A análise exploratória é fundamental para entender a estrutura e as possíveis inconsistências nos dados de leitura dos radares. Foram realizadas as seguintes etapas:

- Análise de Outliers Geográficos
- Verificação de Câmeras com Coordenadas Duplicadas
- Análise de Velocidades Inconsistentes

Nos códigos a seguir a utilização do filtro WHERE (usado para excluir itens com nulos) foi omitida para melhorar a legibilidade.

Análise de Outliers Geográficos

Foram identificadas coordenadas de câmeras localizadas em áreas geograficamente improváveis para o escopo do projeto como Petrópolis, no oceano próximo ao Rio de Janeiro e no oceano ao redor da África. Esses pontos foram considerados erros de leitura portanto foram retirados, diminuindo a quantidade de cameras de 1421 para 1405 ao utilizar uma barreira geográfica ([Código 2.4](#)).

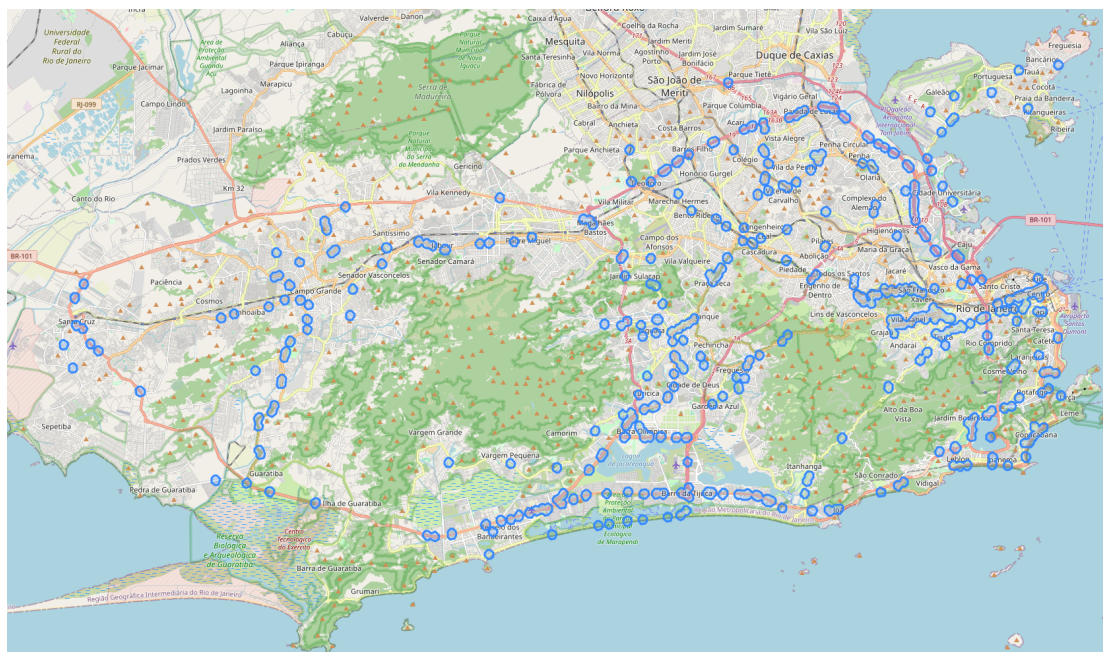


Figura 2.1: Os pontos azuis representam as coordenadas dos radares dentro do limite do Rio de Janeiro

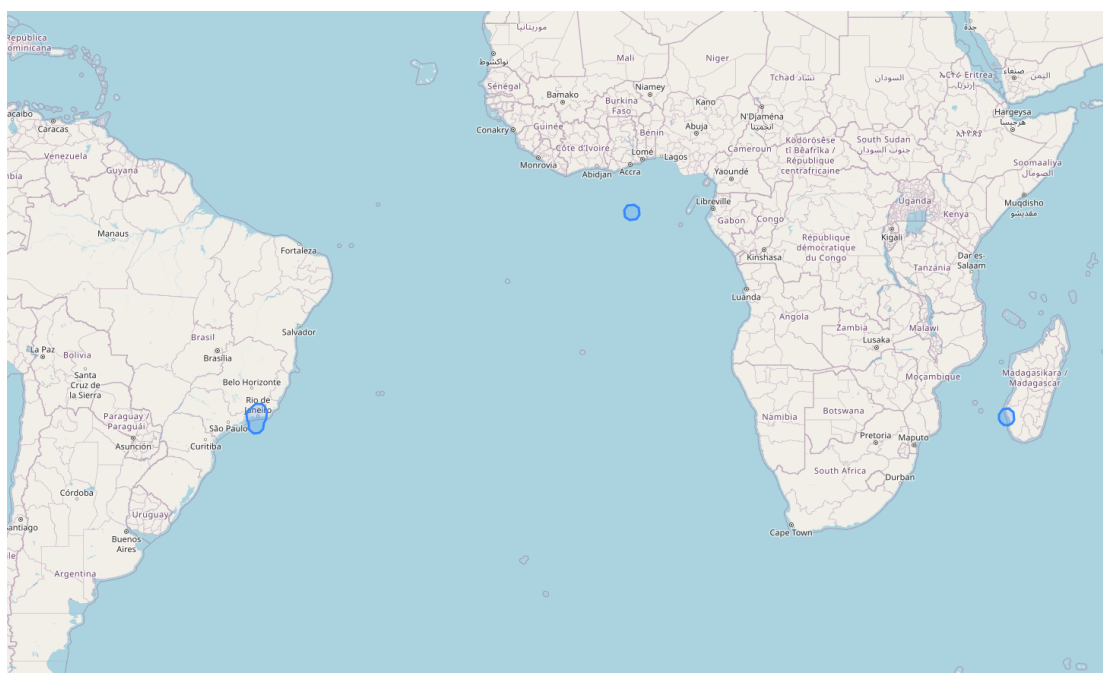


Figura 2.2: Localização dos Outliers nas coordenadas dos radares.

```
1  WITH
2  filtered_valid_data AS (
3      SELECT DISTINCT
4          TO_HEX(camera_numero) AS camera_numero,
5          TO_HEX(placa) AS placa,
6          TO_HEX(empresa) AS empresa,
7          TO_HEX(tipoveiculo) AS tipoveiculo,
8          camera_latitude,
9          camera_longitude,
10         datahora,
11         datahora_captura,
12         velocidade
13     FROM
14         `rj-cetrio.desafio.readings_2024_06`
15
16     cameras_out_of_boundary AS (
17         SELECT
18             camera_numero,
19             camera_latitude,
20             camera_longitude,
21             ST_WITHIN(
22                 ST_GEOPOINT(camera_longitude, camera_latitude),
23                 ST_GEOFROMTEXT('POLYGON((-43.795 -23.082, -43.105 -23.082, -43.105
24                     ↪ -22.738, -43.795 -22.738, -43.795 -23.082))')
25             ) AS within_boundary
26     FROM
27         filtered_valid_data
28     GROUP BY
29         camera_numero, camera_latitude, camera_longitude
30     HAVING
31         within_boundary = FALSE
32 )
33
34     SELECT
35         camera_numero,
36         camera_latitude,
37         camera_longitude
38     FROM
39         cameras_out_of_boundary;
```

Código 2.4: Seleção de dados dentro dos limites do Rio de Janeiro.

Análise de Velocidades Inconsistentes

A tabela abaixo resume as estatísticas descritivas das velocidades capturadas:

Tabela 2.1: Descrição dos dados de velocidade (km/h) por empresa e pelo total.

Empresa	Empresa A	Empresa B	Empresa C	Total
count	8.434.752	18.195.136	2.793.091	34.542.211
mean	34,30	38,52	31,70	37,01
std	8,53	16,46	14,39	15,01
min	1	0	1	0
25%	29	27	22	28
50%	34	38	30	36
75%	39	50	40	46
max	160	255	149	255

Sendo Empresa A = 0891967b413fa4, Empresa B = 1e2545af9d48c6 e Empresa C = 2ce01a80c7f3d0

Velocidades chegando a 255 km/h, o que não é plausível para uma área urbana. Além de que considerando que a velocidade no 75° percentil é de 46 km/h, essas leituras extremas indicam possíveis falhas de medição em alguns radares.

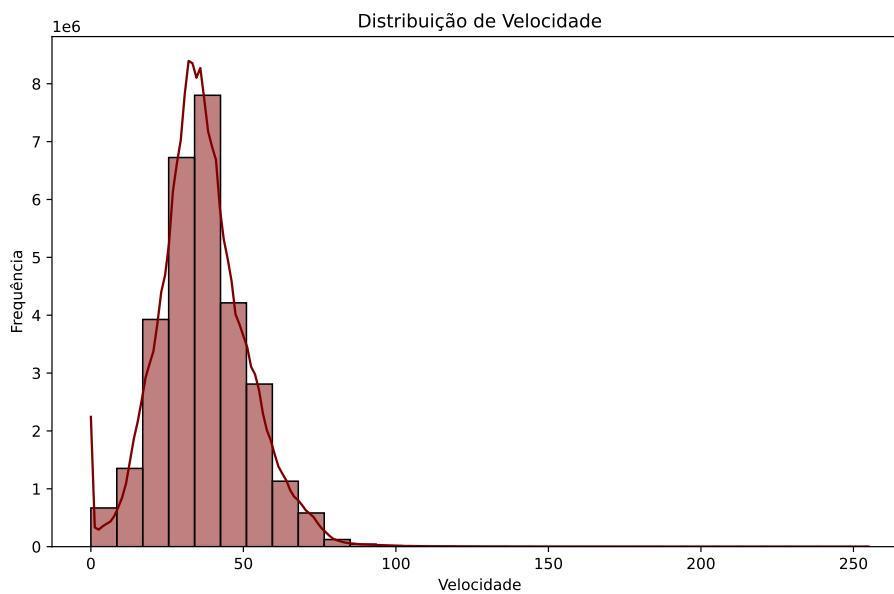


Figura 2.3: Gráfico de distribuição da velocidade.

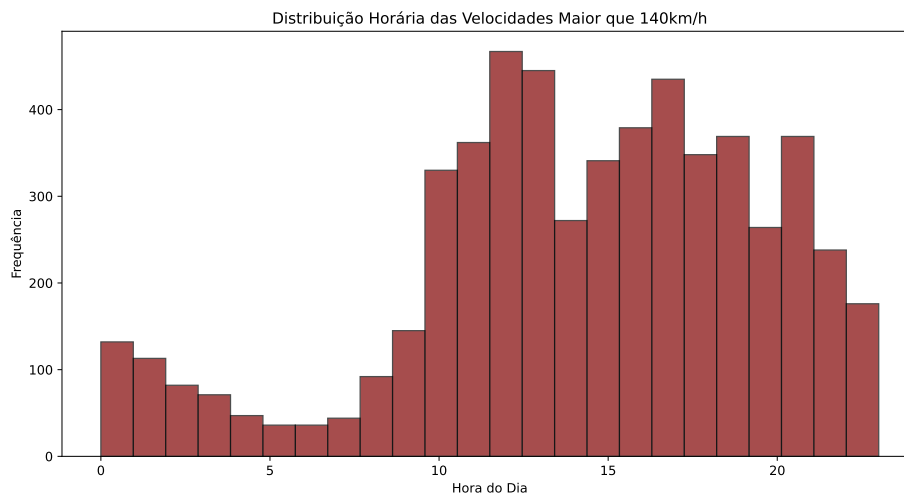


Figura 2.4: Histograma com a frequência de dados com a velocidade acima de 140km/h com relação a hora.

Os gráficos indicam velocidades mais altas durante manhãs e tardes, períodos típicos de maior trânsito. Isso sugere possíveis erros nos radares, uma vez que é improvável que tantas ocorrências de alta velocidade coincidam com horários de pico. Para lidar com esse problema qualquer dado acima de 140km/h entre 6 e 22 horas foi considerado um outlier. Para uma melhor análise seria importante uma auditoria nos equipamentos de medição e revisão dos dados para verificar e corrigir possíveis falhas.

Além disso, essa análise dos dados revelou um total de 306.884 registros com velocidade igual a zero. Este número de leituras não é muito significativo e com as informações dadas não é possível identificar se são problemas na coleta ou no processamento dos dados.

As causas podem incluir:

- **Falhas de Medição:** Problemas técnicos nos dispositivos de monitoramento que resultam em leituras incorretas de velocidade.
- **Paradas de Veículos:** Momentos em que os veículos estavam parados, como em semáforos ou congestionamentos.
- **Erros de Cadastro:** Problemas na configuração dos radares que poderiam levar a registros incorretos.

Detecção de Placas com Múltiplos Tipos de Veículos

Para identificar possíveis casos de clonagem de placas, a metodologia envolveu verificar as placas vinculadas a mais de um tipo de veículo. Inicialmente, foram filtrados os dados inválidos e placas que aparecem apenas uma vez, pois estas não fornecem informações relevantes sobre clonagem (**Código 2.5**). O número de placas que aparecem apenas uma vez é de 4.577.181, representando 12,59% com relação ao total.

```

1  WITH
2  filtered_valid_data AS (
3      SELECT DISTINCT
4          TO_HEX(camera_numero) AS camera_numero,
5          TO_HEX(placa) AS placa,
6          TO_HEX(empresa) AS empresa,
7          TO_HEX(tipoveiculo) AS tipoveiculo,
8          camera_latitude,
9          camera_longitude,
10         datahora,
11         datahora_captura,
12         velocidade
13     FROM
14         `rj-cetrio.desafio.readings_2024_06`
15
16     -- Filtra placas que aparecem mais de uma vez
17     filtered_plates AS (
18         SELECT
19             placa
20         FROM
21             filtered_valid_data
22         GROUP BY
23             placa
24         HAVING
25             COUNT(*) > 1
26     ),
27
28     -- Filtra dados válidos com placas que aparecem mais de uma vez
29     valid_data_with_multiple_plates AS (
30         SELECT
31             *
32         FROM
33             filtered_valid_data
34         WHERE
35             placa IN (SELECT placa FROM filtered_plates)
36     ),
37
38     -- Identifica placas com múltiplos tipos de veículos
39     multiple_vehicle_types AS (
40         SELECT
41             placa,
42             COUNT(DISTINCT tipoveiculo) AS vehicle_type_count
43         FROM
44             valid_data_with_multiple_plates
45         GROUP BY
46             placa
47         HAVING
48             vehicle_type_count > 1
49     )
50     SELECT
51         placa,
52         vehicle_type_count
53     FROM
54         multiple_vehicle_types
55     ORDER BY
56         placa;

```

Código 2.5: Código para selecionar veículos com múltiplos tipos.

```

1      -- Calcula diferença de tempo entre registros consecutivos
2      base_time_diff_data AS (
3          SELECT
4              *,
5              LEAD(datahora) OVER (PARTITION BY placa ORDER BY datahora) AS
                ↳ next_datahora,
6              LEAD(camera_numero) OVER (PARTITION BY placa ORDER BY datahora) AS
                ↳ next_camera_numero,
7              LEAD(camera_latitude) OVER (PARTITION BY placa ORDER BY datahora) AS
                ↳ next_latitude,
8              LEAD(camera_longitude) OVER (PARTITION BY placa ORDER BY datahora) AS
                ↳ next_longitude,
9              TIMESTAMP_DIFF(LEAD(datahora) OVER (PARTITION BY placa ORDER BY datahora),
                ↳ datahora, SECOND) AS time_diff
10         FROM
11             valid_data_with_multiple_plates
12     ),
13
14     -- Calcula a distância geodésica entre pontos consecutivos
15     geo_distance_calculation AS (
16         SELECT
17             *,
18             ST_DISTANCE(
19                 ST_GEOGPOINT(camera_longitude, camera_latitude),
20                 ST_GEOGPOINT(next_longitude, next_latitude)
21             ) / 1000 AS distance_km
22         FROM
23             base_time_diff_data
24     ),
25
26     -- Calcula velocidade baseada na distância e tempo
27     velocity_calculation AS (
28         SELECT
29             *,
30             CASE
31                 WHEN time_diff > 0 THEN distance_km / (time_diff / 3600)
32                 ELSE NULL
33             END AS velocity_kmh
34         FROM
35             geo_distance_calculation
36     ),

```

Código 2.6: Parte do código final que foi utilizada para fazer os cálculos de distância, tempo e velocidade.

Cálculo de Distância entre Pontos e Velocidade

Para detectar possíveis casos de clonagem de placas, uma parte crucial da metodologia envolveu o cálculo da distância entre os pontos de leitura e a velocidade dos veículos (**Código 2.6**). Esse processo foi realizado em várias etapas:

- **Agrupamento por Placa e Datahora:** Inicialmente, os dados foram agrupados por placa e ordenados por data e hora. Esse agrupamento permitiu analisar a sequência de leituras para cada placa de forma cronológica.
- **Cálculo da Diferença de Tempo entre Leituras Consecutivas:** Utilizando a função LEAD, foi possível calcular a diferença de tempo entre leituras consecutivas de uma mesma placa. A função LEAD gera uma nova coluna com o valor da próxima linha dentro do grupo, permitindo calcular o tempo transcorrido entre duas leituras consecutivas.
- **Cálculo da Distância Geodésica:** A distância geodésica entre as coordenadas das leituras consecutivas foi calculada utilizando a função ST_DISTANCE. Essa função computa a distância entre dois pontos geográficos (latitude e longitude), que foi posteriormente convertida de metros para quilômetros.
- **Cálculo da Velocidade:** Com a distância e a diferença de tempo calculadas, foi possível determinar a velocidade média do veículo entre duas leituras consecutivas.

ABORDAGEM E RESULTADOS

A abordagem utilizada no **Código 3.1** para detectar possíveis casos de clonagem de placas resultou na identificação de 583.885 placas suspeitas, sendo 427.574 placas que se enquadraram no caso de ter diversos tipos de veículos.

As principais métricas e critérios empregados são descritos nessa abordagem foram:

- **Escolha da Velocidade Média Máxima:** Para determinar anomalias de velocidade, foi escolhido um valor conservador para a velocidade média máxima. A análise dos dados indicou uma média de aproximadamente 44 km/h. Considerando essa média, foi definido um limite de 70 km/h para identificar possíveis casos de clonagem. Este valor é significativamente mais alto que a média, proporcionando uma margem de segurança que minimiza falsos positivos.
- **Diferença de Tempo:** Foi aplicado um critério conservador ao definir a diferença de tempo entre leituras consecutivas de uma mesma placa. Estabeleceu-se que diferenças de tempo menores que 3600 segundos (1 hora) deveriam ser analisadas, garantindo que apenas deslocamentos plausíveis em períodos razoáveis fossem considerados.
- **Análise das Velocidades Elevadas:** Foi imposto um limite de 140 km/h para velocidades capturadas entre 6h e 22h, excluindo qualquer leitura acima desse valor. Esta escolha baseou-se no fato de que, na cidade do Rio de Janeiro, não há rodovias com limites de velocidade superiores a 120 km/h. Assim, qualquer registro acima de 140 km/h foi considerado uma anomalia e excluído da análise para garantir a integridade dos dados.

```

1      WITH
2      filtered_valid_data AS (
3          SELECT DISTINCT
4              TO_HEX(camera_numero) AS camera_numero,
5              TO_HEX(placa) AS placa,
6              TO_HEX(empresa) AS empresa,
7              TO_HEX(tipoveiculo) AS tipoveiculo,
8              camera_latitude,
9              camera_longitude,
10             datahora,
11             datahora_captura,
12             velocidade
13      FROM
14          `rj-cetrio.desafio.readings_2024_06`
15      WHERE
16          datahora IS NOT NULL
17          AND datahora_captura IS NOT NULL
18          AND placa IS NOT NULL
19          AND empresa IS NOT NULL
20          AND tipoveiculo IS NOT NULL
21          AND velocidade IS NOT NULL
22          AND camera_numero IS NOT NULL
23          AND camera_latitude IS NOT NULL
24          AND camera_longitude IS NOT NULL
25          AND velocidade > 0
26          AND NOT (velocidade > 140 AND EXTRACT(HOUR FROM datahora) BETWEEN 6 AND 22)
27          AND ST_WITHIN(
28              ST_GEOPOINT(camera_longitude, camera_latitude),
29              ST_GEOFROMTEXT('POLYGON((-43.795 -23.082, -43.105 -23.082, -43.105
30                  ↵ -22.738, -43.795 -22.738, -43.795 -23.082))')
31          ),
32
33      -- Filtra placas que aparecem mais de uma vez
34      filtered_plates AS (
35          SELECT
36              placa
37          FROM
38              filtered_valid_data
39          GROUP BY
40              placa
41          HAVING
42              COUNT(*) > 1
43      ),
44
45      -- Filtra dados válidos com placas que aparecem mais de uma vez
46      valid_data_with_multiple_plates AS (
47          SELECT
48              *
49          FROM
50              filtered_valid_data
51          WHERE
52              placa IN (SELECT placa FROM filtered_plates)
53      ),

```

Código 3.1: Abordagem final para os resultados de placas clonadas (pt.1).

```

1      -- Calcula diferença de tempo entre registros consecutivos
2      base_time_diff_data AS (
3          SELECT
4              *,
5              LEAD(datahora) OVER (PARTITION BY placa ORDER BY datahora) AS
6                  ↪ next_datahora,
7              LEAD(camera_numero) OVER (PARTITION BY placa ORDER BY datahora) AS
8                  ↪ next_camera_numero,
9              LEAD(camera_latitude) OVER (PARTITION BY placa ORDER BY datahora) AS
10                 ↪ next_latitude,
11              LEAD(camera_longitude) OVER (PARTITION BY placa ORDER BY datahora) AS
12                 ↪ next_longitude,
13              TIMESTAMP_DIFF(LEAD(datahora) OVER (PARTITION BY placa ORDER BY datahora),
14                 ↪ datahora, SECOND) AS time_diff
15          FROM
16              valid_data_with_multiple_plates
17      ),
18
19      -- Calcula a distância geodésica entre pontos consecutivos
20      geo_distance_calculation AS (
21          SELECT
22              *,
23              ST_DISTANCE(
24                  ST_GEOGPOINT(camera_longitude, camera_latitude),
25                  ST_GEOGPOINT(next_longitude, next_latitude)
26              ) / 1000 AS distance_km
27          FROM
28              base_time_diff_data
29      ),
30
31      -- Calcula velocidade baseada na distância e tempo
32      velocity_calculation AS (
33          SELECT
34              *,
35              CASE
36                  WHEN time_diff > 0 THEN distance_km / (time_diff / 3600)
37                  ELSE NULL
38              END AS velocity_kmh
39          FROM
40              geo_distance_calculation
41      ),

```

Código 3.2: Continuação da abordagem final para os resultados de placas clonadas (pt.2).

```
1      -- Identifica placas com múltiplos tipos de veículos
2      multiple_vehicle_types AS (
3          SELECT
4              placa,
5              COUNT(DISTINCT tipoveiculo) AS vehicle_type_count
6          FROM
7              valid_data_with_multiple_plates
8          GROUP BY
9              placa
10         HAVING
11             vehicle_type_count > 1
12     ),
13
14     -- Identifica possíveis placas clonadas
15     possible_clones AS (
16         SELECT
17             v.*,
18             mv.placa IS NOT NULL AS multiple_vehicle_types
19         FROM
20             velocity_calculation v
21         LEFT JOIN
22             multiple_vehicle_types mv
23         ON
24             v.placa = mv.placa
25         WHERE
26             (v.velocity_kmh > 70 AND v.time_diff < 3600)
27             OR mv.placa IS NOT NULL
28     )
29
30     SELECT
31         *
32     FROM
33         possible_clones
34     ORDER BY
35         placa, datahora;
```

Código 3.3: Continuação da abordagem final para os resultados de placas clonadas (pt.3).

DISCUSSÃO E CONCLUSÃO

Considerações Conservadoras

Adotou-se uma abordagem conservadora ao selecionar os valores para a velocidade média máxima e a diferença de tempo entre leituras. Este cuidado minimiza a possibilidade de falsos positivos, assegurando que apenas os casos mais prováveis de clonagem sejam destacados.

Análise dos Locais de Maior Ocorrência

Durante a análise dos dados de leitura dos radares, identificou-se que havia várias câmeras registradas com as mesmas coordenadas geográficas. Para a análise no Rio de Janeiro, todas as câmeras foram agrupadas de acordo com a localização (Código 4.1).

```
1      SELECT
2          camera_latitude,
3          camera_longitude,
4          COUNT(DISTINCT TO_HEX(camera_numero)) AS num_cameras
5      FROM `rj-cetrio.desafio.readings_2024_06`
6      GROUP BY camera_latitude, camera_longitude
7      HAVING num_cameras > 1;
```

Código 4.1: Agrupamento de cameras de acordo com a localizacao.

Dessa foram identificados 10 locais com maior número de aparições de casos suspeitos, ajudando a identificar pontos críticos na cidade onde há maior incidência de possíveis veículos com placas clonadas passarem. A identificação desses hotspots é crucial para orientar ações de fiscalização e investigação.

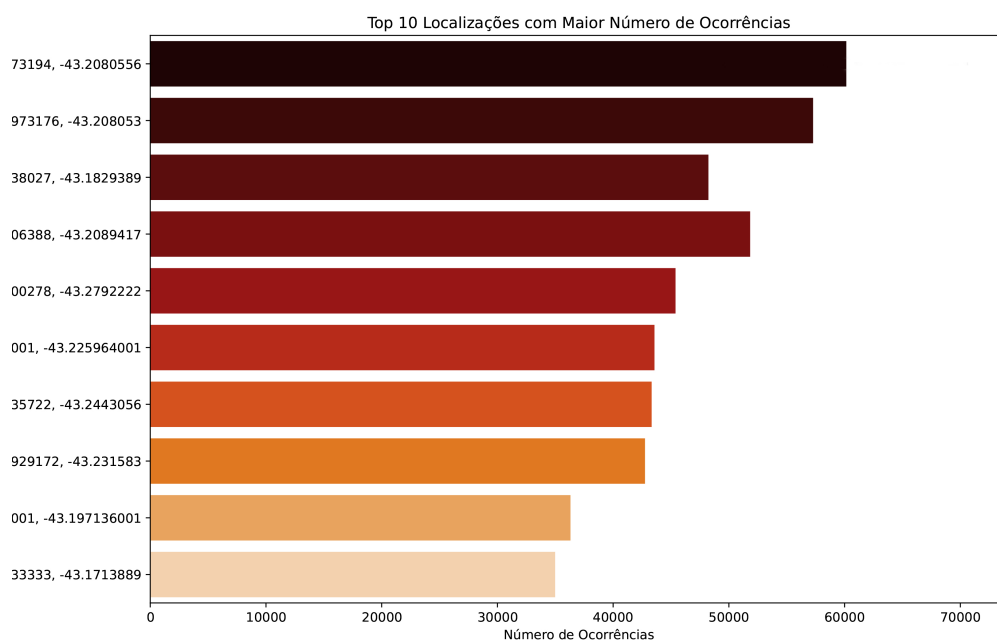


Figura 4.1: Número de ocorrências por coordenada

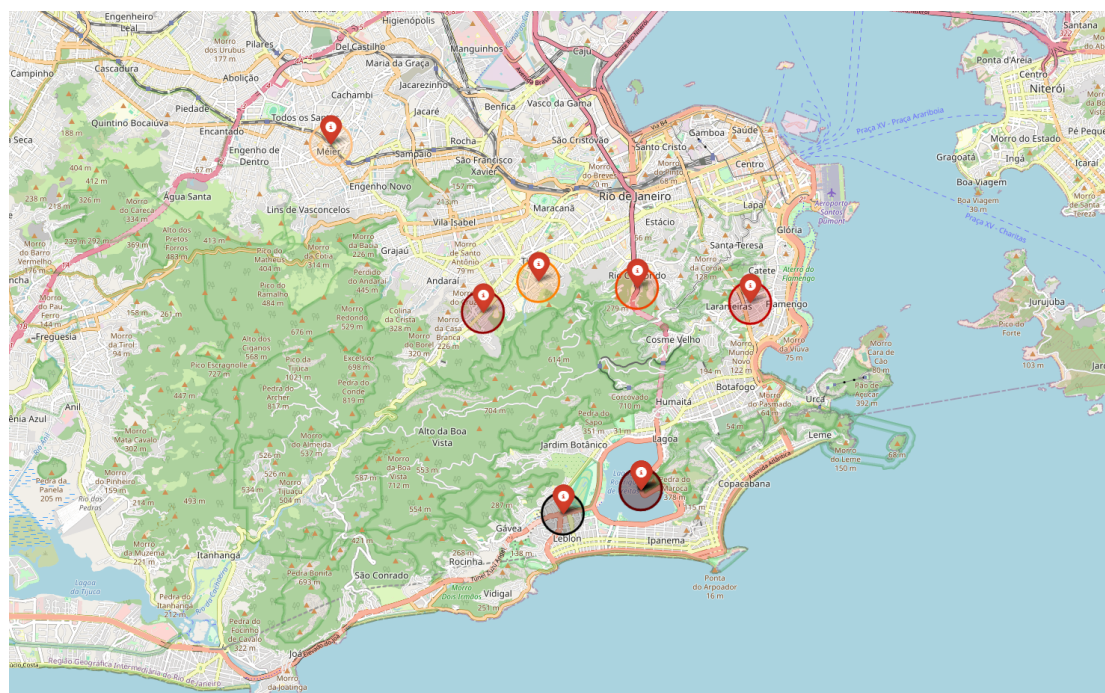


Figura 4.2: Localização das coordenadas com maior número de ocorrência de suspeita de clonagem. Cada ponto está envolto em um círculo de cor referente à Figura 4.1. O número de tags parece menor por conta de radares muito próximos uns dos outros.

Tipos de Veículos

A análise também revelou a distribuição dos tipos de veículos envolvidos nas leituras de placas suspeitas. A [Tabela 4.1](#) apresenta a contagem dos diferentes tipos de veículos detectados e percebe-se que o tipo de veículo identificado como Tipo D representa a maioria das leituras suspeitas, com 8.721.661 ocorrências, seguido pelos outros tipos com menor frequência.

Tabela 4.1: Contagem de cada tipo de veículo dentre o resultado de placas suspeitas.

Tipo Veículo	Count
Tipo A	465.113
Tipo B	785.076
Tipo C	309.443
Tipo D	8.721.661
Total	10.281.293

Sendo Tipo A = 031cc0037e816d, Tipo B = 7a6376f47ca915, Tipo C = b88652111099ed e Tipo D = e2e0029fc0d3e5

Conclusão

A análise realizada foi conservadora ao identificar possíveis casos de clonagem de placas de veículos no Rio de Janeiro, mas atende bem ao propósito de ser um estudo inicial. Utilizando técnicas de análise de dados e SQL, foi possível filtrar, processar e detectar anomalias nos dados de leitura dos radares.

A abordagem apresentada pode ser utilizada como uma base para futuras análises e aprimoramentos. Além disso, a análise dos locais com maior incidência de casos suspeitos e a distribuição dos tipos de veículos fornecem informações valiosas para direcionar os esforços de combate à clonagem de placas.

Para aprimorar ainda mais a detecção de placas clonadas, recomenda-se integrar outras fontes de dados e utilizar técnicas de machine learning para identificar padrões mais complexos.

Trabalhos Futuros

A duplicidade de registros em uma mesma localidade pode introduzir redundâncias nos dados, afetando a precisão das análises de velocidade e de detecção de anomalias. Como melhoria futura, recomenda-se a implementação de um mecanismo para identificar e consolidar os dados provenientes de câmeras sobrepostas, garantindo assim uma base de dados mais limpa e precisa para estudos subsequentes.

DESAFIO CIVITAS
EMD

DETECÇÃO DE PLACAS CLONADAS NO
RIO DE JANEIRO

AMANDA AMARO
amandacsamaro@gmail.com

RIO DE JANEIRO, 01 DE JULHO 2024