

Atividade 2: Exploração de dados criminais

Rodney Rick
RA 181.814
Disciplina MO444
Prof. Dr. Anderson Rocha

Resumo—O objetivo deste trabalho é estudar a representação dos dados, exploração e transformação dos mesmos para posterior análise com modelos de Regressão Logística e Rede Neural Multi-Camada (MLP), a fim de prever a categoria do Crime com base em um determinado conjunto de dados entrada. Discussão dos dados de saída e problemas enfrentados para as execuções.

Palavras-chave: Transformação dos dados, Regressão Logística e MLP.

I. INTRODUÇÃO

Nas seções seguintes o conjunto de dados será apresentado, seguido por análise sobre os alguns pontos de vista das características dos atributos e dos modelos abordados, como o número e a distribuição de amostras, a natureza dos atributos e funções de custo de aprendizagem (computacionalmente).

Para o escopo deste trabalho, foi necessário a extração de alguns dados dos atributos a fim de gerar novos campos interessantes para uso, essa é a etapa de engenharia dos dados. Transformar para obter mais recursos para explorar.

Utilizará duas técnicas classificativas, uma de Rede Neural, conhecida como MLP (*Multi-Layer Perceptron*) e outra mais simples para comparação, a Regressão Logística.

Dados as técnicas, avalia-se atributos descritivos e a combinação dos dados sugeridas para análise (Dia da Semana, Distrito Policial, Latitude e Longitude). Assim, caso haja uma entrada criminal, é possível sua classificação conforme esses dados? Outras discussões são levantadas (seção II) e conforme os gráficos existentes no Apêndice A, verifica-se, de modo visual, qual seriam os possíveis dados mais destacados a serem trabalhados. Esse tipo de visualização, normalmente, é considerada a etapa inicial a fase exploratória dos dados.

Ao final do relatório (seção V) faz-se o levantamento dos resultados e discussão das técnicas aplicadas conforme a resolução.

Ainda no Apêndice A, os gráficos apresentados exibem, por exemplo, dados de Distritos Policiais que mais recebem/tratam dos crimes e quais crimes. Na figura 1, um levantamento em formato de distribuição dos crimes dentro da base de dados.

II. A AMOSTRAGEM DOS DADOS

O conjunto de dados referente a base criminal contém 878.049 amostras. Cada amostra está descrita pelo conjunto de informações descritas na tabela I.

As amostras são completas, porém nas 39 Categorias de crimes que a quantidade de dados mostra-se desproporcional conforme a figura 1.

Tabela I: Tabela de descrição dos campos da base de entrada

Atributo	Descrição dos campos
Data	Data e hora do acontecimento do crime
Categoria	categoria do Incidente que o crime é classificado
Descrição	descrição detalhada do incidente crime
Dia da semana	-
Distrito Policial	nome do Departamento de Polícia do Distrito que tratou o crime
Resolução do Crime	Descrição de como foi resolvido o crime
Endereço	Local do crime
Longitude	-
Latitude	-

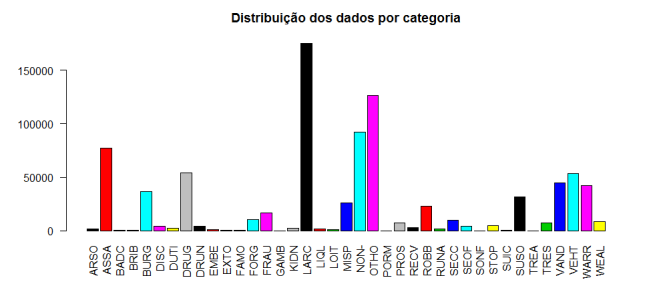


Figura 1: Distribuição dos crimes (os nomes das categorias foram substituídas por abreviações).

Outros levantamentos (citados abaixo) foram feitos e as figuras para demonstração do conceito estão presentes no Apêndice A.

- Distrito Policial → Distribuição crimes por local de atendimento (figura 5);
- Resolução do Crime → Distribuição dos dados conforme a resolução atribuída ao registro (figura 6). Esta parece pouco promissor, por existirem muitos crimes "Sem Resolução" como dados da amostragem.
- Anos → Distribuição dos dados está ao longo de 2003 até início de 2015 e dos crimes realizados (figura 7);
- Mês → Distribuição dos dados está ao longo dos meses do ano, mas que apresenta uma homogeneidade dos crimes (figura 8);
- Dias da semana → Distribuição dos dados está ao longo dos dias da semana, mas que também apresenta uma homogeneidade dos crimes (figura 9);
- Horas → Distribuição dos dados está ao longo das horas dos dias e que demonstra-se relevantes para certos crimes menos cometidos (figura 10);
- Dia da Semana por hora → Mapa de calor da distribuição

A figura 3 (exemplificando somente com 3 classes) demonstra como é o uso da técnica. Necessita-se verificação de uma classe contra todas as outras (em inglês *One-vs-All*), onde torna-se binário o problema para depois os cálculos e cortes no espaço amostral.

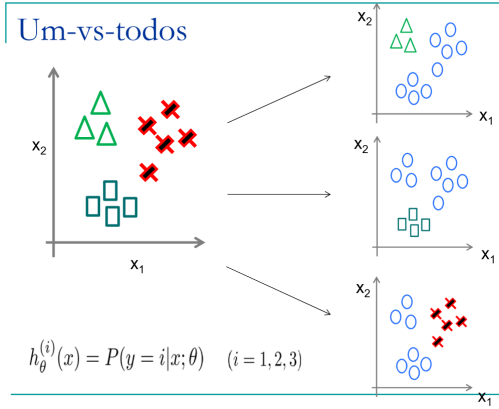


Figura 3: Classificação de Um-vs-Todos (slide da retirada da apresentação da Universidade de Algarve (Portugal) sobre *Machine Learning*)

As equações 1 e 2 representa, respectivamente, as funções *sigmoid* e de custo aplicado para esse processo, esta última estimativa da probabilidade de uma entrada X ser de uma categoria de crime específica. Para o parâmetro θ , inicialmente foram utilizados valores aleatórios, mas que ao decorrer das iterações recebeu forma durante o cálculo da função de gradiente descendente (equação 3).

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad (1)$$

$$h(x) = \frac{1}{1 + e^{\theta^T x}} \quad (2)$$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^i) - y^i)^2 \quad (3)$$

Para o entendimento de como seria elaborado o montagem desta estrutura foram consultados os livros Bishop (capítulo 4) e Wasserman (capítulo 14), Seltman (capítulo 9) e notas das aulas.

B. MLP

Por outro lado, no caso da *MLP*, torna-se necessário trabalhar com um escopo um pouco diferente. Após a binarização dos vetores, tanto para dados descritivos quanto para a combinação (Dia da Semana, Distrito Policial, Latitude, Longitude). Também é feito um vetor binário no qual cada posição representa uma Categoria do crime, atribuindo 1 somente uma vez no vetor para identificação do crime e 0 para todas as outras posições.

A figura 4 representa como entrada a combinação (Dia da Semana, Distrito Policial, Latitude, Longitude) e após a

Tabela II: Comparação dos resultados obtidos para a classificação

	Regressão Logística	MLP
Descrição	92%	65%
Combinação ²	12%	23%

transformação em um vetor binário. Cada neurônio de entrada $I_x, x = [1, 18]$ são os neurônios de entrada. Uma aplicação de uma camada oculta ($H_x, x = [1, 10]$) composta por 10 neurônios. E na camada de saída, 38 categorias ($O_x, x = [1, 38]$) que foram detectadas durante o treino. Lembrando, que mesmo efetuado o balanceamento dos dados, mantendo a aleatoriedade dos dados também para a base de treino, perde-se 1 categoria de crime.

Para o entendimento de como seria feita o montagem desta estrutura foram consultados os livros Bishop (capítulo 5), Carvalho, A. (capítulo 7) e notas das aulas.

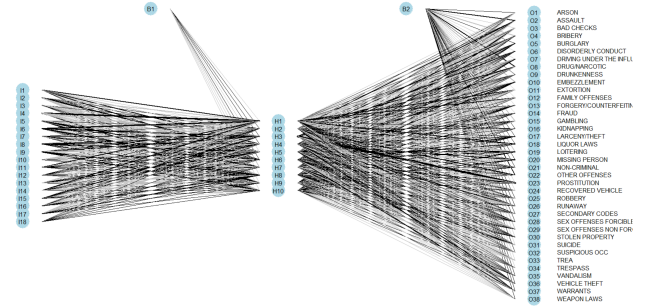


Figura 4: Rede Neural para trabalhar na combinação (Dia da Semana, Distrito Policial, Latitude, Longitude).

V. RESULTADOS

Os resultados obtidos são demonstrados na tabela II. Na qual representa a porcentagem de acerto dado uma entrada X_{ij} onde i representa a amostra dos dados de validação ou teste e j a posição no vetor da amostra (sendo esse dado 0 ou 1).

Os resultados na tabela estimam valores condizentes ao escopo da execução de 30.000 amostras para treino (mas por questão de tempo e hardware, a MLP necessitou uma amostragem menor para processamento), 20.000 amostras de validação e 40.000 para teste. A amostragem inicialmente sugerida, de 700.000 exemplares para treino, e 178.000 para teste, não foi possível ser feita devido os recursos computacionais.

O tempo computacional usado para a execução dos algoritmos, trabalhando com 1000 iterações, MLP (6 horas, lembrando com um escopo menor de dados) foi extraordinariamente maior que da Regressão Logística (40 minutos). Porém para uso de dados formados pela combinação dos atributos, a MLP obteve um valor melhor (23%) que a Regressão Logística (12%).

Foram feitos alguns testes para a utilização da hora e minutos (10), combinado com dia da semana (9) e também

combinados junto os atributos de Latitude e Longitude. Mas não foi obtido valores acima entre 5% e 10%. Independente de trabalhar com uma amostragem maior dos dados, ou no caso da MLP, incluir mais camadas ocultas de neurônios (no caso de 8, 10 e 6 neurônios). Como os resultados não foram promissores, o modelo este modelo foi abandonado.

No geral, pode-se concluir (através da comparação da base de treino, e posterior testes e validações) que o pré-processamento, transformação ajudou para otimização e exploração das técnicas usadas e classificação dos próximos dados de entrada, diminuindo o erro médio e aumento a porcentagem de probabilidade da classificação de saída. Trabalhar com o aumento do número de iterações ou de uma amostragem não afetou de forma substancial os resultados (e algumas vezes variou na margem de 1.5% para erro da classificação, ou para menos ou para mais). Para melhorar ainda mais o resultado, seria necessário para tanto adicionar mais atributos (através da combinação dos já existentes ou obtê-los a partir de outras fontes) ou aplicar modelos mais sofisticados, com funções de custo mais relevantes, combinação de modelos e os modelos não lineares.

Resumindo, não há nenhuma única para melhor solução, no entanto, dentre as propostas neste trabalho, seleciona-se as mais viáveis em resultados promissores, e continua-se lapidando os resultados trabalhando as técnicas e/ou explorar mais formas de extrair conteúdos dos atributos disponibilizados

REFERÊNCIAS

- [1] C. M. Bishop, *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Secaucus, NJ, USA: Springer-Verlag New York, Inc., 2006.
- [2] L. Wasserman, *All of Statistics: A Concise Course in Statistical Inference*. Springer Publishing Company, Incorporated, 2010.
- [3] H. J. Seltman, "Experimental Design and Analysis," 2010. [Online]. Available: <http://www.stat.cmu.edu/~hjseltman/309/Book/>
- [4] Carvalho, A., *Inteligência Artificial: Uma Abordagem de Aprendizagem de Máquina*. LTC. [Online]. Available: <https://books.google.com.br/books?id=4DwelAEACAAJ>

APÊNDICE

Abaixo seguem figuras 5, 6, 7, 8, 9, 11 comentadas anteriormente e separadas por alguns critérios conforme a classificação dos crimes.

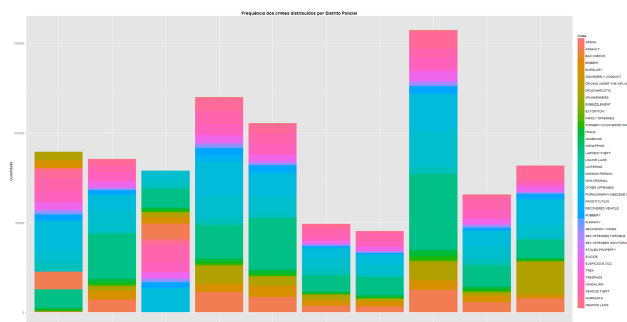


Figura 5: Distribuição dos crimes pelos distritos policiais.

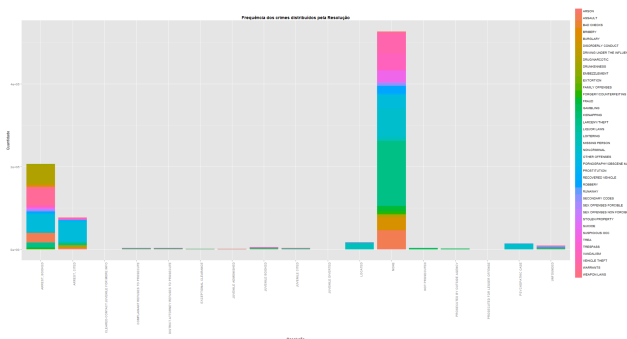


Figura 6: Distribuição dos crimes conforme a resolução.

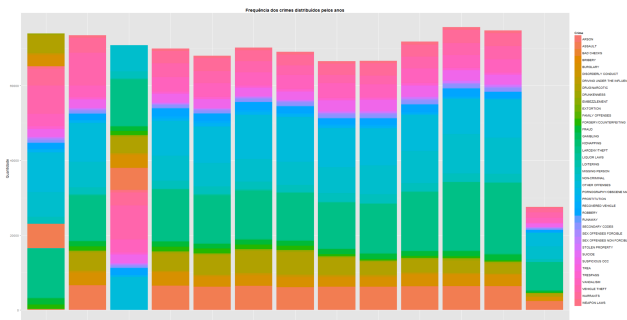


Figura 7: Distribuição dos crimes baseado pelos anos.

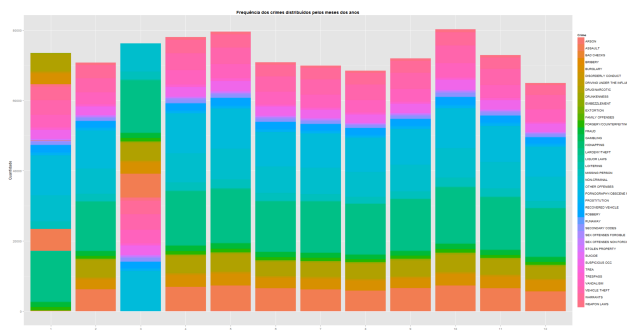


Figura 8: Distribuição dos crimes baseado no mês.

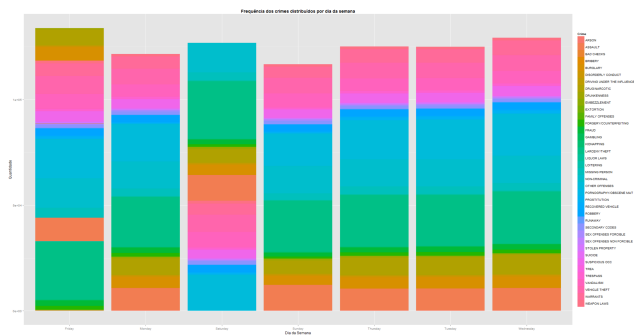


Figura 9: Distribuição dos crimes baseado no dia da semana.

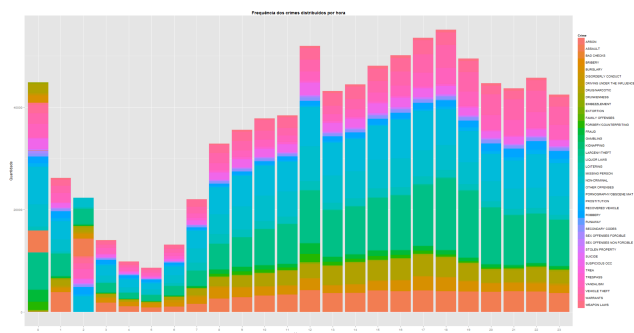


Figura 10: Distribuição dos crimes baseado no horário.

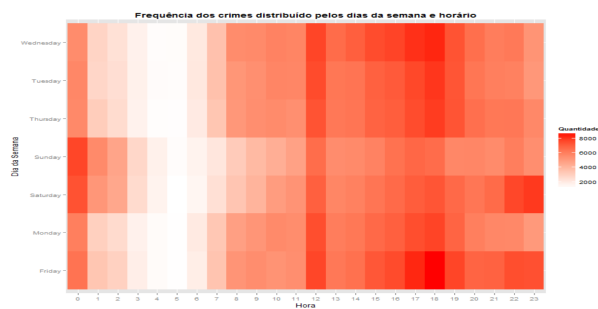


Figura 11: Heatmap da distribuição dos dados pela semana e hora.