

Machine Learning Engineer Nanodegree

Projeto Capstone

Helton Souza Lima

Junho, 2019

I. Definição

Visão geral do projeto

O Programa Bolsa Família (PBF) é o maior programa de distribuição de renda do Brasil [1], através de um benefício em dinheiro transferido diretamente do governo federal para famílias dentro da linha da pobreza e extrema pobreza, para garantir um alívio mais imediato à pobreza, complementando a renda dessas famílias e condicionando à participação nos serviços de saúde e educação. De acordo com o artigo da Dra. Daniela Dias Kuhn [2], o programa foi efetivo na melhoria dos índices de desenvolvimento humano no Estado do Rio Grande do Sul. Podemos citar outro estudo, realizado em Minas Gerais [3] que aponta a mesma conclusão no âmbito deste estado.

Por outro lado, é recorrente a veiculação de notícias [4] referentes a fraudes nos benefícios do Programa Bolsa Família. Essas fraudes acarretam saques de valores superiores ao necessário para o atingimento do objetivo do programa e precisam ser eliminadas, pois acarretam um custo desnecessário ao governo, chegando ao patamar de bilhões [5] de reais.

A empresa em que trabalho é a DATAPREV [6], empresa de processamento de dados do governo federal. Uma atividade recorrente de nossa empresa é o levantamento e cruzamento de informações entre bases de dados para verificar o correto cumprimento de políticas públicas através de sistemas informatizados. O trabalho com os dados do Bolsa Família permitirá a investigação de situações semelhantes a outras que fazem parte das recorrentes demandas dentro da empresa, e a experiência poderá ser útil dentro de um contexto semelhante ao problema abordado neste trabalho.

Descrição do problema

O público-alvo do PBF são as pessoas que estão dentro da faixa da pobreza ou pobreza extrema. Entende-se que os volumes financeiros disponibilizados para o programa é proporcional à quantidade de pessoas dentro das faixas sociais que são alvo do programa, de forma que, a partir de dados de informações sociais e econômicas, como a população total, esperança de vida ao nascer, taxa de analfabetismo, percentual de crianças na escola, taxa de frequência, renda per capita, percentual de distribuição de renda, proporção de pobres, etc, é possível prever o volume financeiro a ser utilizado para o PBF. Em suma, este trabalho visa verificar se municípios com índices mais baixos (índices que compõem o IDH) recebem mais recursos do PBF, pois correspondem a municípios mais pobres. De forma análoga, em tese, municípios com índices mais altos recebem menos recursos do PBF, considerando a quantidade de pessoas residentes nesses municípios.

Uma das respostas que se desejou responder foi: Será que existem municípios com alto IDHM mas que, mesmo assim, recebem muitos recursos do PBF, em comparação com outros municípios semelhantes?

Sendo assim, de posse dos dados granularizados a nível de município brasileiro, relativos à pesquisa de mapeamento do Índice de Desenvolvimento Humano Municipal (IDHM) no ano de 2010, utilizou-se modelos de machine learning que foram treinados utilizando-se os dados de parte desses municípios e foram capazes de prever o volume financeiro da outra parte desses municípios. Em um momento inicial, a análise dos dados apontou a correlação entre os indicadores sociais e o volume financeiro do PBF associado com cada município. Em seguida, foi possível identificar alguns municípios que apontaram discrepância nessa correlação e foram apontados como municípios onde é possível que tenha sofrido uma maior influência de fraudes.

Métricas

O valor a ser previsto é um valor contínuo, correspondente ao valor, em reais, que é disponibilizado para ser sacado pelos beneficiários do Bolsa Família para cada município. A métrica que foi utilizada é o Root Mean Squared Erros (RMSE), pois é uma métrica que avalia a distância entre o valor previsto e o valor real. Essa métrica é calculada pelo próprio scikit-learn comparando os valores previstos e os valores reais, através da utilização do método "score" dos modelos de regressão. O RMSE é definido como $(1 - u/v)$, onde u é soma das diferenças ao quadrado ($\text{quadrado}(\text{real} - \text{previsto}).\text{sum}()$) e v é o total da soma dos quadrados ($\text{quadrado}(\text{real} - \text{média}(\text{real})).\text{sum}()$). A melhor possibilidade é o valor de "score" ser 1.0 e pode ser negativo se o modelo se comportou de forma muito ruim [11].

II. Análise

Exploração dos dados

Os dados utilizados foram obtidos de duas fontes. A primeira fonte são os dados relacionados ao Índice de Desenvolvimento Humano Municipal (IDHM), disponibilizado pelo site Atlas do Desenvolvimento Humano no Brasil [7] ou no site da Kaggle [8]. Os dados do IDHM são disponibilizados para cada um dos 5565 municípios brasileiros, sendo composto por dados que podem ser agrupados em 3 dimensões: dados sobre longevidade, dados sobre o nível de acesso ao conhecimento e dados sobre a renda. O cálculo do IDHM foi realizado a partir das informações dos 3 últimos Censos Demográficos do IBGE (1991, 2000 e 2010). Neste trabalho foram utilizados os dados do IDHM de 2010.



Expectativa
de vida
ao nascer

Escolaridade
da população
adulta

Fluxo escolar
da população
jovem

Renda
per capita

$$\sqrt[3]{(\text{Ícone 1}) \times (\text{Ícone 2})}$$

MÉDIA GEOMÉTRICA
RAIZ CÚBICA DA MULTIPLICAÇÃO DOS
SUBÍNDICES COM PESOS 1 E 2



IDHM
longevidade



IDHM
educação



IDHM
renda

$$\sqrt[3]{(\text{Círculo 1}) \times (\text{Círculo 2}) \times (\text{Círculo 3})}$$

MÉDIA GEOMÉTRICA
RAIZ CÚBICA DA MULTIPLICAÇÃO DOS 3 IDHMS

=

IDHM

A segunda fonte são os dados relacionados à quantidade de famílias beneficiárias e o total de pagamentos disponibilizados pelo PBF para cada município brasileiro. Os dados são disponibilizados pelo Ministério da Cidadania [9]. Os dados utilizados são de janeiro de 2010, ou seja, 7 anos após o ano de lançamento do PBF, que pode ser considerado como suficiente para o programa ter atingido uma maturidade em sua operacionalização e gestão e os dados serem considerados consolidados. Também são dados que coincidem com o ano da realização do Censo, em 2010, como forma de aproximar o levantamento social realizado pelo Censo dos dados de recursos disponibilizados pelo Bolsa Família.

União dos dados

A primeira etapa do trabalho foi a união de ambas as fontes de dados para formar um único conjunto de dados. O resultado final é composto de 5565 linhas (correspondentes a cada município) e 241 colunas (4 colunas dos dados do Bolsa Família, incluindo a *Quantidade de Famílias Beneficiárias do Bolsa Família* e o *Valor Repassado para Bolsa Família*, e 237 colunas dos dados para composição do IDHM).

Tratamento de variáveis categóricas

Foram identificadas 7 colunas com valores não-numéricos e de códigos pertencentes a domínios:

ano: Sempre o mesmo ano em todas as linhas (2010)

codmun6, ibge, codmun7: Códigos identificadores do município

município: Nome do município

anomes: Competência (mês + ano) do valor disponibilizado pelo Bolsa Família, sempre com valor "201001", que significa janeiro de 2010.

uf: Código do IBGE identificador da Unidade de Federação ao qual o município pertence.

Todas essas variáveis foram removidas para a continuação da análise exploratória e alimentação dos modelos de predição.

Transformação de valores

Através da exploração inicial dos dados, verificou-se que a variável **idhm** possuía alguns registros entre 0 e 1 e o restante dos registros entre 400 e 900. Em verificações individuais destes casos, percebeu-se que os registros estavam apenas transformados para valores entre 0 e 1. Por exemplo, para o município de Cabixi, em Rondônia, o valor que se verificou foi 0,65. Entretanto, após pesquisa no portal Atlas Brasil, este município foi avaliado com IDHM 650.



Portanto, decidiu-se realizar a transformação destes casos para que todos ficassem com a mesma base. O mesmo procedimento foi realizado para **idhm_e, idhm_l, idhm_r, i_freq_prop e i_escolaridade**.

Dados ausentes e outliers

Não foram identificados dados ausentes no conjunto de dados, após a realização de busca por lacunas. Em relação aos *outliers*, nenhum caso foi interpretado como *outlier*. Todas as variáveis analisadas individualmente apresentaram distribuição normal ou distribuição normal mista, com dois picos. Foram avaliadas individualmente as variáveis

valor_repassado_bolsa_familia, qtd_familias_beneficiarias_bolsa_familia, idhm, idhm_e, idhm_l, idhm_r, i_freq_prop, i_escolaridade, theil, gini, pmpob, pind e pesotot.

Visualização Exploratória

Variáveis avaliadas individualmente através de gráficos

Conforme relatado na seção anterior, as variáveis avaliadas individualmente foram analisadas através de gráficos que estão a seguir:

valor_repassado_bolsa_familia

- Grande parte dos municípios recebe até 100 mil reais. A quantidade de municípios com valor maior que 200 mil reais está em torno de 25%



qtd_familias_beneficiarias_bolsa_familia

- Grande parte dos municípios possui até mil famílias beneficiadas. A quantidade de municípios com mais de 2 mil famílias beneficiárias está em torno de 25%. O maior valor é 181531 família beneficiárias.



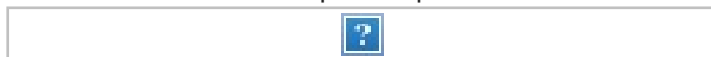
idhm

- Índice de Desenvolvimento Humano do Município. É uma distribuição normal mista, com dois pontos de picos, próximo dos valores 600 e 720. O menor valor é 418 e o maior é 862.



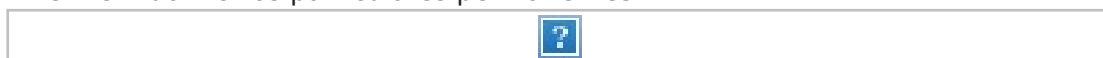
gini

- Mede o grau de desigualdade existente na distribuição de indivíduos segundo a renda domiciliar per capita. Seu valor varia de 0, quando não há desigualdade (a renda domiciliar per capita de todos os indivíduos tem o mesmo valor), a 1, quando a desigualdade é máxima (apenas um indivíduo detém toda a renda). O universo de indivíduos é limitado àqueles que vivem em domicílios particulares permanentes.



pmpob

- Proporção dos indivíduos com renda domiciliar per capita igual ou inferior a R\$ 140,00 mensais, em reais de agosto de 2010. O universo de indivíduos é limitado àqueles que vivem em domicílios particulares permanentes.



pind

- Proporção dos indivíduos com renda domiciliar per capita igual ou inferior a R\$ 70,00 mensais, em reais de agosto de 2010. O universo de indivíduos é limitado àqueles que

vivem em domicílios particulares permanentes.

pesotot

- População total de cada município



Análise de variáveis correlacionadas

In this section, you will need to provide some form of visualization that summarizes or extracts a relevant characteristic or feature about the data. The visualization should adequately support the data being used. Discuss why this visualization was chosen and how it is relevant.

Questions to ask yourself when writing this section:

- *Have you visualized a relevant characteristic or feature about the dataset or input data?*
- *Is the visualization thoroughly analyzed and discussed?*
- *If a plot is provided, are the axes, title, and datum clearly defined?*

Algoritmos e técnicas

Modelo de referência

III. Metodologia

Pré-processamento dos dados

O que mais chamou a atenção na análise exploratória dos dados foi a quantidade de variáveis existentes no conjunto de dados referentes ao cálculo do IDHM: 237 variáveis. Uma hipótese levantada no início do trabalho e que norteou a preparação dos dados foi a possibilidade de eliminar variáveis que fossem redundantes para alimentação de modelos de machine learning. Sendo assim, buscou-se inicialmente eliminar as variáveis categóricas e, em seguida, aquelas com forte relação e que agregariam muito pouco aos modelos em relação à capacidade de predição, sendo apenas informações que deixam o processamento mais lento.

Através de gráficos que demonstram a correlação

In this section, all of your preprocessing steps will need to be clearly documented, if any were necessary. From the previous section, any of the abnormalities or characteristics that you identified about the dataset will be addressed and corrected here. Questions to ask yourself when writing this section:

- *If the algorithms chosen require preprocessing steps like feature selection or feature transformations, have they been properly documented?*
- *Based on the **Data Exploration** section, if there were abnormalities or characteristics that needed to be addressed, have they been properly corrected?*
- *If no preprocessing is needed, has it been made clear why?*

Implementação

Refinamento

IV. Resultados

Avaliação e validação do modelo

Justificativa

V. Conclusão

Visualização de forma livre

Reflexão

Aperfeiçoamento

Referências

- [1] [Portal do Programa Bolsa Família. Ministério da Cidadania. \(http://mds.gov.br/assuntos/bolsa-familia\)](http://mds.gov.br/assuntos/bolsa-familia)
- [2] [Kuhn, Daniela Dias. Tonetto, Elci da Silva. O Programa Bolsa Família e os indicadores sociais no Rio Grande do Sul. Desenvolvimento em Questão \(https://www.revistas.unijui.edu.br/index.php/desenvolvimentoemquestao/article/view/5799/530\)](https://www.revistas.unijui.edu.br/index.php/desenvolvimentoemquestao/article/view/5799/530)
- [3] [Denubila, Lais Atanaka. Ferreira, Marco Aurelio Marques. Monteiro, Doraliza Auxiliadora Abranches. Programa Bolsa Família: Análise Da Trajetória Dos Indicadores Sociais Em Minas Gerais. Associação Nacional de Pós-Graduação e Pesquisa em Administração \(http://www.anpad.org.br/admin/pdf/apb1239.pdf\)](http://www.anpad.org.br/admin/pdf/apb1239.pdf)
- [4] [Busca no Google sobre fraudes no Bolsa Família \(https://www.google.com/search?q=bolsa+fam%C3%ADlia+fraudes&rlz=1C1GCEU_pt-brBR835BR835&source=lnms&tbn=nws&sa=X&ved=0ahUKEwiz_MzgsLbhAhU7KLkGHcQzCmq\)](https://www.google.com/search?q=bolsa+fam%C3%ADlia+fraudes&rlz=1C1GCEU_pt-brBR835BR835&source=lnms&tbn=nws&sa=X&ved=0ahUKEwiz_MzgsLbhAhU7KLkGHcQzCmq)
- [5] ["Controladoria-Geral acha R\\$ 1,3 bi em fraudes no Bolsa Família", Revista Exame Online, 4 de janeiro de 2018 \(https://exame.abril.com.br/brasil/controladoria-geral-acha-r-13-bi-em-fraudes-no-bolsa-familia/\)](https://exame.abril.com.br/brasil/controladoria-geral-acha-r-13-bi-em-fraudes-no-bolsa-familia/)
- [6] [Portal da Dataprev. Empresa de Tecnologia e Informações da Previdência Social \(http://www.dataprev.gov.br/\)](http://www.dataprev.gov.br/)
- [7] [Portal do Atlas do Desenvolvimento Humano no Brasil \(http://www.atlasbrasil.org.br/2013/pt/o_atlas/idhm/\)](http://www.atlasbrasil.org.br/2013/pt/o_atlas/idhm/)
- [8] [Human Development Indexes and Census data for Brazilian municipalities. Portal Kaggle \(https://www.kaggle.com/pauloeduneves/hdi-brazil-idh-brasil\)](https://www.kaggle.com/pauloeduneves/hdi-brazil-idh-brasil)
- [9] [Visualizador de Dados Sociais. Um portal do Ministério da Cidadania \(https://aplicacoes.mds.gov.br/saqi/vis/data/data-table.php\)](https://aplicacoes.mds.gov.br/saqi/vis/data/data-table.php)
- [10] [Human Development Indexes and Census data for Brazilian municipalities. Kaggle DataSet. Setembro/2018 \(https://www.kaggle.com/kerneler/starter-hdi-brazil-idh-brasil-80f68b4b-6\)](https://www.kaggle.com/kerneler/starter-hdi-brazil-idh-brasil-80f68b4b-6)
- [11] [Método "score" do modelo Linear Regression. Biblioteca scikit-learn v0.21.2 \(https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LinearRegression.html#sklearn.linear\)](https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LinearRegression.html#sklearn.linear)

