

# Mineração de Dados – Classificação

## Preparação da base de microdados do ENEM 2012

Jonathan Coutinho Luz de Queiroz<sup>1</sup>  
Guilherme Lima Bernal<sup>1</sup>

<sup>1</sup>Instituto de Matemática – Universidade Federal da Bahia

jonathanqueiroz@dcc.ufba.br, ufba@lbguilherme.com

### 1. Introdução

Nesta atividade, preparamos a base de microdados do ENEM 2012 para a futura utilização em um algoritmo de classificação. O objetivo final do trabalho consiste em classificar os candidatos em quatro grupos, de acordo com a nota obtida na prova de redação: “ruim”, “regular”, “boa” e “ótima”. O intervalo de valores para a nota da redação que está associado a cada um desses grupos é baseado na divisão do dataset original em quartis, sendo “ruim” o primeiro quartil e “ótima” o último. Para passar pela qualificação, todas as colunas precisam ser normalizadas para números de 0 a 1.

Optamos, nesta etapa, por remover colunas majoritariamente por razões técnicas (número de identificação do candidato, por exemplo, além de outros campos descritos a seguir). Consequentemente, colunas que simplesmente suspeitamos serem poucos relevantes para a classificação, devido à análise estatística previamente realizada, foram mantidas. Essa abordagem permite que a análise estatística sirva como um guia durante o desenvolvimento do algoritmo de classificação, ao invés de determinar completamente quais dados serão utilizados nele.

### 2. Normalização da base de dados

As seguintes operações foram realizadas com o intuito de normalizar a base de dados:

- Foi criada a coluna **UF\_INSC\_IGUAL\_PROVA**, que guarda “1” se **UF\_INSC** for igual ao **UF\_MUNICIPIO\_PROVA**, indicando que o candidato não mudou de estado para fazer a prova. “0” caso contrário.
- **UF\_INSC** foi expandido para um total de 26 colunas, cada uma contendo “1” ou “0” para indicar se o candidato fez a prova nesse estado ou não. Nenhum candidato poder ter mais de um estado marcado como “1”. Informações adicionais sobre a localização do candidato (nome do município, código do município, etc.) foram removidas.
- De forma similar, **TP\_SEXO** foi expandido para duas colunas: **GENERO\_HOMEM** e **GENERO\_MULHER**, com apenas uma delas sendo marcada como “1”.
- As colunas **NOTA\_CN**, **NOTA\_CH**, **NOTA\_LC** e **NOTA\_MT** passaram a guardar a nota do candidato em cada uma das provas. Para que o valor seja consistente de 0 a 1, essa nota foi dividida por 1000.

- Várias colunas que guardavam as opções marcadas pelo candidato, gabarito, identificação da prova e afins foram removidas por serem irrelevantes para obtenção do resultado.
- A coluna ***TP\_ESCOLA*** foi expandida para duas colunas: ***ESCOLA\_PUBLICA*** e ***ESCOLA\_PRIVADA***. Para a maioria dos candidatos, o tipo de escola não foi informado, e portanto ambas as colunas apresentam valor falso. Demais informações sobre a escola e a entidade certificadora do candidato (nome, município, etc.) foram removidas.
- Analogamente, a coluna ***ID\_LOCALIZACAO*** foi expandida para duas colunas: ***LOCALIZACAO\_URBANA*** e ***LOCALIZACAO\_RURAL***. Para a maioria dos candidatos, o tipo de localização não foi informado, e portanto ambas as colunas apresentam valor falso.
- A idade foi normalizada para valores de 0 a 1 através de uma divisão por 100. Dessa forma, a menor idade passou a ser 0.12 e a maior idade passou a ser 0.70.
- ***TP\_COR\_RACA*** foi dividido em 5 colunas, uma com cada cor/raça. Cada uma possui “1” se ela é verdadeira para o candidato. Nenhum candidato tem duas dessas colunas como verdadeiras.
- Analogamente, ***TP\_ESTADO\_CIVIL*** foi dividido em 4 colunas booleanas: ***ESTADO\_CIVIL\_SOLTEIRO***, ***ESTADO\_CIVIL\_CASADO***, ***ESTADO\_CIVIL\_DIVORCIADO*** e ***ESTADO\_CIVIL\_VIUVO***.
- ***ST\_CONCLUSAO*** foi dividido em 4 colunas seguindo o mesmo princípio: ***ENSINO\_MEDIO\_CONCLUIDO*** (alunos que já haviam concluído o Ensino Médio quando a prova foi aplicada), ***ENSINO\_MEDIO\_EM\_2012*** (alunos que estavam concluindo o Ensino Médio no mesmo ano de aplicação da prova), ***ENSINO\_MEDIO\_DEPOIS\_2012*** (alunos que só concluiriam o Ensino Médio em anos posteriores à aplicação da prova) e ***ENSINO\_MEDIO\_NAO\_FAZ*** (alunos que nem haviam concluído previamente nem estavam cursando o Ensino Médio). O ano exato de conclusão do Ensino Médio, quando informado, foi descartado.
- Por fim, a nota da redação foi dividida nos 4 grupos anteriormente citados: “ruim”, “regular”, “boa” e “ótima”.