

Módulo 5 Tarefa 1

Base de nascidos vivos do DataSUS

O DataSUS disponibiliza diversos arquivos de dados com relação a seus segurados, conforme a [lei da transparência de informações públicas \(https://www.sisgov.com/transparencia-acesso-informacao/#:~:text=A%20Lei%20da%20Transpar%C3%Aancia%20\(LC,em%20um%20site%20na%20intern](https://www.sisgov.com/transparencia-acesso-informacao/#:~:text=A%20Lei%20da%20Transpar%C3%Aancia%20(LC,em%20um%20site%20na%20intern)

Essas informações podem ser obtidas pela internet [aqui \(http://www2.datasus.gov.br/DATASUS/index.php?area=0901&item=1\)](http://www2.datasus.gov.br/DATASUS/index.php?area=0901&item=1). Como o processo de obtenção desses arquivos foge um pouco do nosso escopo, deixamos o arquivo `SINASC_RO_2019.csv` já como vai ser encontrado no DataSUS. O dicionário de dados está no arquivo `estrutura_sinasc_para_CD.pdf` (o nome do arquivo tal qual no portal do DataSUS).

Nosso objetivo

Queremos deixar uma base organizada para podermos estudar a relação entre partos com risco para o bebê e algumas condições como tempo de parto, consultas de pré-natal etc.

Preparação da base

1. Carregue a base 'SINASC_RO_2019.csv'. Conte o número de registros e o número de registros não duplicados da base. Dica: você aprendeu um método que remove duplicados, encadeie este método com um outro método que conta o número de linhas. **Há linhas duplicadas?**
2. Conte o número de valores *missing* por variável.
3. Ok, no item anterior você deve ter achado pouco prático ler a informação de tantas variáveis, muitas delas nem devem ser interessantes. Então crie uma seleção dessa base somente com as colunas que interessam. São elas:

```
['LOCNASC', 'IDADEMAE', 'ESTCIVMAE', 'ESMAE', 'QTDFILVIVO',  
 'GESTACAO', 'GRAVIDEZ', 'CONSULTAS', 'APGAR5']
```

Refaça a contagem de valores *missings*.

4. Apgar é uma *nota* que o pediatra dá ao bebê quando nasce de acordo com algumas características associadas principalmente à respiração. Apgar 1 e Apgar 5 são as notas 1 e 5 minutos do nascimento. Apgar5 será a nossa variável de interesse principal. Então remova todos os registros com Apgar5 não preenchido. Para esta seleção, conte novamente o número de linhas e o número de *missings*.
5. observe que as variáveis ['ESTCIVMAE', 'CONSULTAS'] possuem o código 9, que significa *ignorado*. Vamos assumir que o não preenchido é o mesmo que o código 9.
6. Substitua os valores faltantes da quantitativa (QTDFILVIVO) por zero.
7. Das restantes, decida que valor te parece mais adequado (um 'não preenchido' ou um valor 'mais provável' como no item anterior) e preencha. Justifique. Lembre-se de que tratamento de dados é trabalho do cientista, e que estamos tomando decisões a todo o momento - não há necessariamente certo e errado aqui.
8. O Apgar possui uma classificação indicando se o bebê passou por asfixia:
 - Entre 8 e 10 está em uma faixa 'normal'.
 - Entre 6 e 7, significa que o recém-nascido passou por 'asfixia leve'.
 - Entre 4 e 5 significa 'asfixia moderada'.

- Entre 0 e 3 significa 'asfixia severa'.

Crie uma categorização dessa variável com essa codificação e calcule as frequências dessa categorização.

9. Renomeie as variáveis para que fiquem no *snake case*, ou seja, em letras minúsculas, com um *underscore* entre as palavras. Dica: repare que se você não quiser criar um *dataframe* novo, você vai precisar usar a opção `inplace = True`.