## Opções de Conjuntos de Dados

Escolha um dos conjuntos de dados abaixo ou use um seu próprio (veja abaixo considerações para este caso).

| Conjunto de Dados   | Visão Geral  | Questão Guia   | Estimativ<br>a de<br>Tempo |
|---|--|--|----------------------------|
| Red Wine Quality¹  Leia este arquivo texto que descreve as variáveis e como os dados foram coletados.             | Este conjunto de dados contém<br>1.599 vinhos tintos com 11<br>variáveis de propriedades<br>químicas do vinho. Ao menos 3<br>especialistas em vinhos<br>avaliaram cada vinho,<br>fornecendo uma nota entre 0<br>(muito ruim) e 10 (muito<br>excelente).  | Quais propriedades<br>químicas influenciam a<br>qualidade dos vinhos<br>tintos?  | 10-20 horas                |
| White Wine Quality <sup>2</sup> Leia este arquivo texto que descreve as variáveis e como os dados foram coletados | Este conjunto de dados contém<br>4.898 vinhos brancos com 11<br>variáveis de propriedades<br>químicas do vinho. Ao menos 3<br>especialistas em vinhos<br>avaliaram cada vinho,<br>fornecendo uma nota entre 0<br>(muito ruim) e 10 (muito<br>excelente). | Quais propriedades<br>químicas influenciam a<br>qualidade dos vinhos<br>brancos?   | 10-20 horas                |
| Financial Contributions to Presidential Campaigns by State  | Selecione uma eleição usando os botões e clique em "Export Contributor Data" para obter os conjuntos de dados. Escolha UM estado e explore as contribuições feitas para um candidato em um ano de eleições.  | perguntas sobre este conjunto de dados. Você pode adicionar variáveis a este conjunto de dados como sexo ou partido político do candidato. |                            |
| Loan Data from<br>Prosper   | Este conjunto de dados possui<br>113.937 empréstimos com 81<br>variáveis em cada um, incluindo<br>o valor, taxa de juros, status do  | Faça suas próprias<br>perguntas sobre este<br>conjunto de dados.<br>Existem MUITAS variáveis   | 15-30 horas                |

<sup>&</sup>lt;sup>1</sup> P. Cortez, A. Cerdeira, F. Almeida, T. Matos and J. Reis. Modeling wine preferences by data mining from physicochemical properties. In Decision Support Systems, Elsevier, 47(4):547-553. ISSN: 0167-9236. Available at: [@Elsevier] <a href="http://dx.doi.org/10.1016/j.dss.2009.05.016">http://dx.doi.org/10.1016/j.dss.2009.05.016</a> [Pre-press (pdf)] <a href="http://www3.dsi.uminho.pt/pcortez/winequality09.pdf">http://www3.dsi.uminho.pt/pcortez/winequality09.pdf</a> [bib] <a href="http://www3.dsi.uminho.pt/pcortez/dss09.bib">http://www3.dsi.uminho.pt/pcortez/dss09.bib</a>

| Última atualização em 11/03/2014 Este <u>dicionário de</u> variáveis explica as variáveis do conjunto de dados. | pagamento, receita do<br>mutuário, seu emprego atual,<br>histórico do cartão de crédito e<br>informações sobre seu último<br>pagamento.  | neste conjunto de dados<br>e você não deverá<br>explorar todas. Escolha<br>entre 10 a 15 variáveis<br>para sua análise. |           |
|---|--|---|-----------|
| Encontre seu<br>próprio conjunto<br>de dados!   | Lembre-se que encontrar e limpar o conjunto de dados é uma tarefa que demanda tempo e esforço significativos! Veja a lista abaixo caso você deseje utilizar seu próprio conjunto de dados. | Faça suas próprias<br>perguntas sobre o<br>conjunto de dados!   | 30+ horas |

## Caso você esteja usando um conjunto de dados próprio...

| Sau | coni | iunto  | d۸ | dados | dovo. |
|-----|------|--------|----|-------|-------|
| seu | COIL | iuiito | ue | uauos | ueve: |

- ☐ possuir ao menos 1.000 observações
- ☐ conter ao menos uma variável categórica (você pode criar uma)
- ☐ conter ao menos 8 variáveis diferentes
- estar em um formato "limpo"¹ (você pode ter que realizar a limpeza e formatação dos seus dados como parte da exploração)
- estar em um formato usual, como .csv, .tsv, .txt, ou .xls

## Aqui estão algumas fontes para encontrar conjuntos de dados:

- <a href="http://www.inside-r.org/howto/finding-data-internet">http://www.inside-r.org/howto/finding-data-internet</a> (não utilize o conjunto de dados do Titanic)
- <a href="http://opendata.stackexchange.com/">http://opendata.stackexchange.com/</a>
- http://www.data.gov/

<sup>&</sup>lt;sup>1</sup> Conjuntos de dados "limpos" (tidy) são aqueles que possuem uma estrutura particular. Leia mais sobre este tipo de conjunto de dados no artigo de Hadley Wickham's, <a href="http://vita.had.co.nz/papers/tidy-data.pdf">http://vita.had.co.nz/papers/tidy-data.pdf</a>