

Estruturação de dados do DOU

Tomás Barcellos

22 de maio de 2018

Caminho

1. O problema
2. A solução proposta
3. Os resultados alcançados

O Problema

A fonte

O Diário Oficial da União (DOU) é importante fonte de informações oficiais consultadas por diversos atores sociais. A informações publicadas no DOU são publicadas em formato de textos, tabelas e imagens, dados não-estruturados.

Os desafios

Até dezembro 2017 a Imprensa Nacional disponibilizada o DOU somente em PDF.

A partir de dezembro de 2017 o DOU passa a ser publicado também em HTML (com erros).

Falta de padronização dos textos publicados

A solução proposta

Proposta: um pacote R

1. Em linha com os princípios do tidyverse
2. Usar `regex` para identificar as informações publicadas

rdou: Etapas de processamento

O seguinte fluxo de trabalho foi adotado:

1. Download de todos as páginas do DOU (PDF) do dia;
2. Conversão dos PDFs em TXTs pelo Word;
3. Processamento dos arquivos TXT para estruturar a informação;
e
4. Validação humana da informação processada.

rdou: download

```
# devtools::install_github("tomasbarcellos/rdou")  
library(rdou)  
download_dou("02/03/2017", dest_dir = "pdf")
```

A função `download_dou()` faz o download das páginas do DOU, em PDF. A função é chamada pelo seu efeito colateral (baixar) e retorna a data (invisível).

rdou: conversão

```
paginas <- converter_pdf(  
  data = "02/03/2017", secao = 1,  
  dir_pdf = "pdf", dest_dir = "txt"  
i)
```

A função `converter_pdf()` faz a conversão das páginas do DOU de PDF para TXT. A função é chamada pelo seu efeito colateral e retorna um vetor com o nome dos arquivos TXT criados (invisível).

rdou: processamento

```
agric <- extrair_normas(paginas, "Agricultura")  
faz <- extrair_normas(paginas, "Fazenda")  
str(agric, give.attr = FALSE, vec.len = 1)
```

```
## List of 4  
## $ : chr [1:5] "PORTARIA Nº 27, DE 21 DE FEVEREIRO DE 20  
## $ : chr [1:5] "RESOLUÇÃO Nº 2, DE 24 DE FEVEREIRO DE 20  
## $ : chr [1:12] "RETIFICAÇÃO" ...  
## $ : chr [1:7] "PORTARIA Nº 46, DE 21 DE FEVEREIRO DE 20
```

rdou: processamento

```
# Objetos "norma" possuem alguns atributos  
str(attributes(faz), vec.len = 1)
```

```
## List of 7  
## $ class      : chr "norma"  
## $ orgao      : chr [1:22] "SUPERINTENDÊNCIA DE NORMAS"  
## $ arquivos   : chr [1:19] "inst/doc/txt/DOU1/2017/mar"  
## $ data_dou   : Date[1:1], format: "2017-03-02"  
## $ secao      : num 1  
## $ encodificacao: chr "latin1"  
## $ orgao_alvo  : chr "Ministério da Fazenda"
```

rdou: estruturação das informações

```
df_agric <- estruturar_normas(agric)  
dplyr::glimpse(df_agric)
```

```
## Observations: 4  
## Variables: 9  
## $ numero      <chr> "000000027", "000000002", NA, "000000004"  
## $ tipo        <chr> "PORTARIA N", "RESOLUÇÃO N", "AVISO"  
## $ orgao       <chr> "SECRETARIA DE DEFESA AGROPECUÁRIA"  
## $ texto       <chr> "PORTARIA Nº 27, DE 21 DE FEVEREIRO"  
## $ promulgacao <date> 2017-03-02, 2017-03-02, 2017-03-02  
## $ ementa      <chr> "Credencia a empresa DÍGITOS CERTIFI"  
## $ titulo      <chr> "PORTARIA Nº 27, DE 21 DE FEVEREIRO"  
## $ pagina      <dbl> 5, 5, 6, 6  
## $ secao       <int> 1, 1, 1, 1
```

Usando o pipe

```
library(magrittr)
download_dou("02/03/2017") %>%
  converter_pdf(secao = 1) %>%
  pegar_normas_dou(orgao_alvo = "Agricultura") %>%
  estruturar_normas()
```


Resultados

Resultados

- ▶ Mais de 13.705 normas encontradas entre os meses de abril de 2015 e dezembro de 2017 nas seções 1 e 2 do DOU para o Ministério da Agricultura, Pecuária e Abastecimento.
- ▶ Alimentação semi-automatizada do Sistemas de Consulta a Legislação Agropecuária.