

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

Os desafios

Nº 37 sexta-feira, 26 de abril de 2013

Ministério da Agricultura, Pecuária e Abastecimento

INSTITUTO NACIONAL DE METEOROLOGIA

PORTARIA Nº 22, DE 22 DE ABRIL DE 2013

O Diretor do Instituto Nacional de Meteorologia, no uso de suas atribuições legais, faz saber que o Edital nº 01/2013, publicado em 20 de abril de 2013 no DOU nº 80, de 22 de abril de 2013, publicado no Diário Oficial da União nº 22 de abril de 2013, publicado no Diário Oficial da União nº 22 de abril de 2013, para a contratação de serviços de consultoria técnica e elaboração de projeto de lei de organização de serviços, encontra-se em fase de julgamento e o prazo para apresentação de propostas é de 15 dias, contados a partir da publicação desta Portaria, no Diário Oficial da União nº 26 de abril de 2013, publicado no Diário Oficial da União nº 26 de abril de 2013.

JOSE MAURO DE BEZINNE

SUPERINTENDÊNCIA FEDERAL DO ESTADO DE ALAGOAS

PORTARIA Nº 01, DE 22 DE ABRIL DE 2013

O Superintendente Federal de Agricultura, Registro e Abastecimento, no uso de suas atribuições legais, faz saber que o Edital nº 01/2013, publicado em 20 de abril de 2013 no DOU nº 80, de 22 de abril de 2013, publicado no Diário Oficial da União nº 22 de abril de 2013, para a contratação de serviços de consultoria técnica e elaboração de projeto de lei de organização de serviços, encontra-se em fase de julgamento e o prazo para apresentação de propostas é de 15 dias, contados a partir da publicação desta Portaria, no Diário Oficial da União nº 26 de abril de 2013, publicado no Diário Oficial da União nº 26 de abril de 2013.

ALAN CORREIA DE ARAÚJO

SUPERINTENDÊNCIA FEDERAL DO ESTADO DA BAHIA

PORTARIA Nº 46, DE 22 DE ABRIL DE 2013

A SUPERINTENDÊNCIA FEDERAL DE AGRICULTURA DO ESTADO DA BAHIA, no uso das competências que lhe são conferidas pelo Decreto nº 11.218, de 22 de junho de 2010, e pelo Edital nº 01/2013, publicado em 20 de abril de 2013 no DOU nº 80, de 22 de abril de 2013, publicado no Diário Oficial da União nº 22 de abril de 2013, para a contratação de serviços de consultoria técnica e elaboração de projeto de lei de organização de serviços, encontra-se em fase de julgamento e o prazo para apresentação de propostas é de 15 dias, contados a partir da publicação desta Portaria, no Diário Oficial da União nº 26 de abril de 2013, publicado no Diário Oficial da União nº 26 de abril de 2013.

VIVIANA ALVES ALMEIDA RANGEL

SUPERINTENDÊNCIA FEDERAL DO ESTADO DO RIO GRANDE DO SUL

PORTARIAS Nº 12 DE ABRIL DE 2013

O SUPERINTENDENTE FEDERAL DE AGRICULTURA DO ESTADO DO RIO GRANDE DO SUL, no uso das atribuições legais, faz saber que o Edital nº 01/2013, publicado em 20 de abril de 2013 no DOU nº 80, de 22 de abril de 2013, publicado no Diário Oficial da União nº 22 de abril de 2013, para a contratação de serviços de consultoria técnica e elaboração de projeto de lei de organização de serviços, encontra-se em fase de julgamento e o prazo para apresentação de propostas é de 15 dias, contados a partir da publicação desta Portaria, no Diário Oficial da União nº 26 de abril de 2013, publicado no Diário Oficial da União nº 26 de abril de 2013.

Nº 101 - DEFENSIVA e, ainda, o prazo de 15 de abril de 2013, e o prazo para apresentação de propostas é de 15 dias, contados a partir da publicação desta Portaria, no Diário Oficial da União nº 26 de abril de 2013, publicado no Diário Oficial da União nº 26 de abril de 2013.

Nº 102 - DEFENSIVA e, ainda, o prazo de 15 de abril de 2013, e o prazo para apresentação de propostas é de 15 dias, contados a partir da publicação desta Portaria, no Diário Oficial da União nº 26 de abril de 2013, publicado no Diário Oficial da União nº 26 de abril de 2013.

Nº 103 - DEFENSIVA e, ainda, o prazo de 15 de abril de 2013, e o prazo para apresentação de propostas é de 15 dias, contados a partir da publicação desta Portaria, no Diário Oficial da União nº 26 de abril de 2013, publicado no Diário Oficial da União nº 26 de abril de 2013.

FRANCISCO NIVAL SINGER

Diário Oficial da União - Seção 2

ISSN 1677-7030



Ministério da Ciência, Tecnologia e Inovação

AGÊNCIA ESPECIAL BRASILEIRA

BRTE/AC-01

No Parecer nº 04/ABR, de 15 de abril de 2013, publicado no DOU nº 80 de 22 de abril de 2013, Seção 2, página 1, foi o conteúdo das informações.

ROBERTO DE ANDRADE OLIVEIRA BARBOSA

ROBERTO DE ANDRADE OLIVEIRA BARBOSA

INSTITUTO NACIONAL DE PESQUISAS ESPACIAIS

PORTARIA Nº 1.574, DE 14 DE ABRIL DE 2013

O Diretor do Instituto Nacional de Pesquisas Espaciais, no uso das atribuições legais que lhe são conferidas pelo Decreto nº 11.218, de 22 de junho de 2010, e pelo Edital nº 01/2013, publicado em 20 de abril de 2013 no DOU nº 80, de 22 de abril de 2013, publicado no Diário Oficial da União nº 22 de abril de 2013, para a contratação de serviços de consultoria técnica e elaboração de projeto de lei de organização de serviços, encontra-se em fase de julgamento e o prazo para apresentação de propostas é de 15 dias, contados a partir da publicação desta Portaria, no Diário Oficial da União nº 26 de abril de 2013, publicado no Diário Oficial da União nº 26 de abril de 2013.

LEONEL FERNANDO PERDINI

Ministério da Cultura

SECRETARIA EXECUTIVA

PORTARIA Nº 104, DE 22 DE ABRIL DE 2013

O SECRETÁRIO EXECUTIVO DO MINISTÉRIO DA CULTURA, no uso das atribuições que lhe são conferidas pelo Decreto nº 11.218, de 22 de junho de 2010, e pelo Edital nº 01/2013, publicado em 20 de abril de 2013 no DOU nº 80, de 22 de abril de 2013, publicado no Diário Oficial da União nº 22 de abril de 2013, para a contratação de serviços de consultoria técnica e elaboração de projeto de lei de organização de serviços, encontra-se em fase de julgamento e o prazo para apresentação de propostas é de 15 dias, contados a partir da publicação desta Portaria, no Diário Oficial da União nº 26 de abril de 2013, publicado no Diário Oficial da União nº 26 de abril de 2013.

ROSELI CARVALHO DE MOURA

ROSELI CARVALHO DE MOURA

FUNDAÇÃO CULTURAL PALMARES

PORTARIA Nº 45, DE 22 DE ABRIL DE 2013

O PRESIDENTE DA FUNDAÇÃO CULTURAL PALMARES, no uso das atribuições legais que lhe são conferidas pelo Decreto nº 11.218, de 22 de junho de 2010, e pelo Edital nº 01/2013, publicado em 20 de abril de 2013 no DOU nº 80, de 22 de abril de 2013, publicado no Diário Oficial da União nº 22 de abril de 2013, para a contratação de serviços de consultoria técnica e elaboração de projeto de lei de organização de serviços, encontra-se em fase de julgamento e o prazo para apresentação de propostas é de 15 dias, contados a partir da publicação desta Portaria, no Diário Oficial da União nº 26 de abril de 2013, publicado no Diário Oficial da União nº 26 de abril de 2013.

JOSE HERTON SANTANA ALMEIDA

JOSE HERTON SANTANA ALMEIDA

INSTITUTO BRASILEIRO DE MUSEUS

PORTARIAS Nº 22 DE ABRIL DE 2013

O SUPERINTENDENTE DO INSTITUTO BRASILEIRO DE MUSEUS, no uso das atribuições legais que lhe são conferidas pelo Decreto nº 11.218, de 22 de junho de 2010, e pelo Edital nº 01/2013, publicado em 20 de abril de 2013 no DOU nº 80, de 22 de abril de 2013, publicado no Diário Oficial da União nº 22 de abril de 2013, para a contratação de serviços de consultoria técnica e elaboração de projeto de lei de organização de serviços, encontra-se em fase de julgamento e o prazo para apresentação de propostas é de 15 dias, contados a partir da publicação desta Portaria, no Diário Oficial da União nº 26 de abril de 2013, publicado no Diário Oficial da União nº 26 de abril de 2013.

JOSE HERTON SANTANA ALMEIDA

JOSE HERTON SANTANA ALMEIDA

Nº 171 - DEFENSIVA e, ainda, o prazo de 15 de abril de 2013, e o prazo para apresentação de propostas é de 15 dias, contados a partir da publicação desta Portaria, no Diário Oficial da União nº 26 de abril de 2013, publicado no Diário Oficial da União nº 26 de abril de 2013.

Nº 172 - DEFENSIVA e, ainda, o prazo de 15 de abril de 2013, e o prazo para apresentação de propostas é de 15 dias, contados a partir da publicação desta Portaria, no Diário Oficial da União nº 26 de abril de 2013, publicado no Diário Oficial da União nº 26 de abril de 2013.

Nº 173 - DEFENSIVA e, ainda, o prazo de 15 de abril de 2013, e o prazo para apresentação de propostas é de 15 dias, contados a partir da publicação desta Portaria, no Diário Oficial da União nº 26 de abril de 2013, publicado no Diário Oficial da União nº 26 de abril de 2013.

Nº 174 - DEFENSIVA e, ainda, o prazo de 15 de abril de 2013, e o prazo para apresentação de propostas é de 15 dias, contados a partir da publicação desta Portaria, no Diário Oficial da União nº 26 de abril de 2013, publicado no Diário Oficial da União nº 26 de abril de 2013.

Nº 175 - DEFENSIVA e, ainda, o prazo de 15 de abril de 2013, e o prazo para apresentação de propostas é de 15 dias, contados a partir da publicação desta Portaria, no Diário Oficial da União nº 26 de abril de 2013, publicado no Diário Oficial da União nº 26 de abril de 2013.

Nº 176 - DEFENSIVA e, ainda, o prazo de 15 de abril de 2013, e o prazo para apresentação de propostas é de 15 dias, contados a partir da publicação desta Portaria, no Diário Oficial da União nº 26 de abril de 2013, publicado no Diário Oficial da União nº 26 de abril de 2013.

Nº 177 - DEFENSIVA e, ainda, o prazo de 15 de abril de 2013, e o prazo para apresentação de propostas é de 15 dias, contados a partir da publicação desta Portaria, no Diário Oficial da União nº 26 de abril de 2013, publicado no Diário Oficial da União nº 26 de abril de 2013.

Nº 178 - DEFENSIVA e, ainda, o prazo de 15 de abril de 2013, e o prazo para apresentação de propostas é de 15 dias, contados a partir da publicação desta Portaria, no Diário Oficial da União nº 26 de abril de 2013, publicado no Diário Oficial da União nº 26 de abril de 2013.

Nº 179 - DEFENSIVA e, ainda, o prazo de 15 de abril de 2013, e o prazo para apresentação de propostas é de 15 dias, contados a partir da publicação desta Portaria, no Diário Oficial da União nº 26 de abril de 2013, publicado no Diário Oficial da União nº 26 de abril de 2013.

Nº 180 - DEFENSIVA e, ainda, o prazo de 15 de abril de 2013, e o prazo para apresentação de propostas é de 15 dias, contados a partir da publicação desta Portaria, no Diário Oficial da União nº 26 de abril de 2013, publicado no Diário Oficial da União nº 26 de abril de 2013.

Nº 181 - DEFENSIVA e, ainda, o prazo de 15 de abril de 2013, e o prazo para apresentação de propostas é de 15 dias, contados a partir da publicação desta Portaria, no Diário Oficial da União nº 26 de abril de 2013, publicado no Diário Oficial da União nº 26 de abril de 2013.

Nº 182 - DEFENSIVA e, ainda, o prazo de 15 de abril de 2013, e o prazo para apresentação de propostas é de 15 dias, contados a partir da publicação desta Portaria, no Diário Oficial da União nº 26 de abril de 2013, publicado no Diário Oficial da União nº 26 de abril de 2013.

para acesso: 00224-33000000

Documento assinado digitalmente conforme MP nº 2.204-2 de 2009/06, que institui a

assinatura de Chave Pública Brasileira - ICP-Brasil.

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 todas 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"  
)
```

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/ma"  
## $ data_dou   : Date[1:1], format: "2017-03-02"  
## $ secao      : num 1  
## $ encodificacao: chr "latin1"  
## $ orgao_alvo  : chr "Ministério da Fazenda"
```

udou: estruturação das informações

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

```
## Observations: 4  
## Variables: 9  
## $ numero      <int> 27, 2, NA, 46  
## $ tipo        <chr> "PORTARIA", "RESOLUÇÃO", "AVISO DE P  
## $ 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 CERTIF  
## $ titulo      <chr> "PORTARIA Nº 27, DE 21 DE FEVEREIRO  
## $ pagina      <int> 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.