

Estruturação de dados do DOU

Tomás Barcellos

22 de maio de 2018

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.

A fonte



DIÁRIO OFICIAL DA UNIÃO

República Federativa do Brasil - Imprensa Nacional

Em circulação desde 1º de outubro de 1862

Ano CLIII Nº 80

Brasília - DF, terça-feira, 10 de maio de 2016

ISSN 1677-7042



SEÇÃO 1

Sumário

Assunto	Página
Atos do Poder Executivo	1
Presidência da República	7
Ministério da Agricultura, Pecuária e Abastecimento	9
Ministério da Cultura, Patrimônio e Esportes	10
Ministério da Cultura	11
Ministério da Defesa	18
Ministério da Educação	19
Ministério da Fazenda	30
Ministério da Saúde	41
Ministério da Justiça	46
Ministério da Planeta	47
Ministério da Relações Exteriores	48
Ministério da Segurança Pública	49
Ministério da Saneamento	50
Ministério da Saúde	51
Ministério da Segurança Pública	52
Ministério da Saneamento	53
Ministério da Saúde	54
Ministério da Segurança Pública	55
Ministério da Saneamento	56
Ministério da Saúde	57
Ministério da Segurança Pública	58
Ministério da Saneamento	59
Ministério da Saúde	60
Ministério da Segurança Pública	61
Ministério da Saneamento	62
Ministério da Saúde	63
Ministério da Segurança Pública	64
Ministério da Saneamento	65
Ministério da Saúde	66
Ministério da Segurança Pública	67
Ministério da Saneamento	68
Ministério da Saúde	69
Ministério da Segurança Pública	70
Ministério da Saneamento	71
Ministério da Saúde	72
Ministério da Segurança Pública	73
Ministério da Saneamento	74
Ministério da Saúde	75
Ministério da Segurança Pública	76
Ministério da Saneamento	77
Ministério da Saúde	78
Ministério da Segurança Pública	79
Ministério da Saneamento	80
Ministério da Saúde	81
Ministério da Segurança Pública	82
Ministério da Saneamento	83
Ministério da Saúde	84
Ministério da Segurança Pública	85
Ministério da Saneamento	86
Ministério da Saúde	87
Ministério da Segurança Pública	88
Ministério da Saneamento	89
Ministério da Saúde	90
Ministério da Segurança Pública	91
Ministério da Saneamento	92
Ministério da Saúde	93
Ministério da Segurança Pública	94
Ministério da Saneamento	95
Ministério da Saúde	96
Ministério da Segurança Pública	97
Ministério da Saneamento	98
Ministério da Saúde	99
Ministério da Segurança Pública	100

Atos do Poder Executivo

DECRETO Nº 8.706, DE 9 DE MAIO DE 2016

Altera o Decreto nº 6.464, de 27 de maio de 2010, que dispõe sobre a organização e o funcionamento dos órgãos e entidades da administração pública federal direta, com o seguinte texto:

A PRESIDENTE DA REPÚBLICA, no uso da atribuição que lhe confere o art. 84, inciso V, alínea "f", da Constituição,

DECRETA:

Art. 1º O Decreto nº 6.464, de 27 de maio de 2010, passa a vigorar com as seguintes alterações:

TABELA DE PREÇOS DE JORNAL AVULSOS

Página	Valor	Valor	Valor
de 1 a 10	R\$ 0,50	R\$ 0,50	R\$ 0,50
de 11 a 20	R\$ 0,50	R\$ 0,50	R\$ 0,50
de 21 a 30	R\$ 0,50	R\$ 0,50	R\$ 0,50
de 31 a 40	R\$ 0,50	R\$ 0,50	R\$ 0,50
de 41 a 50	R\$ 0,50	R\$ 0,50	R\$ 0,50
de 51 a 60	R\$ 0,50	R\$ 0,50	R\$ 0,50
de 61 a 70	R\$ 0,50	R\$ 0,50	R\$ 0,50
de 71 a 80	R\$ 0,50	R\$ 0,50	R\$ 0,50
de 81 a 90	R\$ 0,50	R\$ 0,50	R\$ 0,50
de 91 a 100	R\$ 0,50	R\$ 0,50	R\$ 0,50

Brasília, 9 de maio de 2016; 17ª de Independência e 170ª da República.

DEMA RICUSSEFF
Diretor Geral

DECRETO Nº 8.706, DE 9 DE MAIO DE 2016

Altera o Decreto nº 6.464, de 27 de maio de 2010, que dispõe sobre a organização e o funcionamento dos órgãos e entidades da administração pública federal direta, com o seguinte texto:

A PRESIDENTE DA REPÚBLICA, no uso da atribuição que lhe confere o art. 84, inciso V, alínea "f", da Constituição,

DECRETA:

Art. 1º O Decreto nº 6.464, de 27 de maio de 2010, passa a vigorar com as seguintes alterações:

TABELA DE PREÇOS DE JORNAL AVULSOS

Página	Valor	Valor	Valor
de 1 a 10	R\$ 0,50	R\$ 0,50	R\$ 0,50
de 11 a 20	R\$ 0,50	R\$ 0,50	R\$ 0,50
de 21 a 30	R\$ 0,50	R\$ 0,50	R\$ 0,50
de 31 a 40	R\$ 0,50	R\$ 0,50	R\$ 0,50
de 41 a 50	R\$ 0,50	R\$ 0,50	R\$ 0,50
de 51 a 60	R\$ 0,50	R\$ 0,50	R\$ 0,50
de 61 a 70	R\$ 0,50	R\$ 0,50	R\$ 0,50
de 71 a 80	R\$ 0,50	R\$ 0,50	R\$ 0,50
de 81 a 90	R\$ 0,50	R\$ 0,50	R\$ 0,50
de 91 a 100	R\$ 0,50	R\$ 0,50	R\$ 0,50

Brasília, 9 de maio de 2016; 17ª de Independência e 170ª da República.

DEMA RICUSSEFF
Diretor Geral

Este documento pode ser verificado no endereço eletrônico <http://www.in.gov.br/imprensa/leitura>.

Documento autêntico digitalmente assinado por 10.738.626 de 20/05/2016, que insere a

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).

A solução proposta

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("projdiario/projdiario", subdir  
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 `dest_dir` (invisível).

rDOU: conversão

```
paginas <- converter_pdf(dir_pdf = "pdf", secao = 1,  
                          data = "02/03/2017",  
                          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")  
mcti <- extrair_normas(paginas, "Ciência")  
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(mcti))
```

```
## List of 7  
## $ class      : chr "norma"  
## $ orgao      : chr [1:12] "GABINETE DO MINISTRO" "GA  
## $ 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 Ciência, Tecnologia
```

rDOU: estruturação das informações

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

```
## Observations: 4
## Variables: 11
## $ NR_ATO          <chr> "000000027", "000000002", NA, "0000
## $ SG_TIPO         <chr> "POR", "RES", "ART", "POR"
## $ AN_ATO          <chr> "2017", "2017", "2017", "2017"
## $ SG_ORGAO        <chr> "SDA/MAPA", "SDA/MAPA", "SPA/MAPA
## $ CD_TIPO_ATO     <chr> "A", "A", "A", "A"
## $ TX_TEXTO        <chr> "<p>Ministério da Agricultura, Pe
## $ DT_PROMULGACAO  <date> 2017-03-02, 2017-03-02, 2017-03-
## $ TX_EMENTA       <chr> "Credencia a empresa DÍGITOS CERT
## $ DS_TITULO        <chr> "PORTARIA Nº 27, DE 21 DE FEVERE
## $ NM_PAGINA       <dbl> 5, 5, 6, 6
```

Usando o pipe

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