



WEB SCRAPING E VISUALIZAÇÕES DE NOTÍCIAS RELACIONADAS AO BNDES

MBA Executivo em Business Analytics e Big Data
TBABD-8 – Análise de Mídias Sociais e Mineração de Texto
Professor Fernando Ferreira

RAFAEL HENRIQUE ROCHA DE SOUZA
A56660250

Agenda

- Acessarei links sobre as notícias do BNDES presentes na página de buscas do portal G1;
- Apresentarei o passo a passo do processo de web scraping, criação da nuvem de palavras e desambiguação de entidades;
- Estruturarei um grafo com a coocorrência de entidades e analisarei por comunidades e relevâncias;
- Utilizarei o GEPHI para extrair características, gerar métricas e visualizar as redes.

Reticulando o ambiente Python para ser usado no RStudio

```
# instalando e carregando o pacote reticulate
install.packages("reticulate", repos = "http://cran.us.r-project.org")
library(reticulate)

# verificando a versao reticulada do python e a env
py_config()

# verificando todas as virtual envs
conda_list()

#conda_remove("r-spacy")
#conda_create(envname = "~/.rstudio-desktop/virtualenvs/spacy")

# instalando os pacotes na env
conda_install(envname = "r-reticulate", "spacy")

# instalando e carregando o pacote spacyr
install.packages("spacyr", repos = "http://cran.us.r-project.org")
library(spacyr)

#find_spacy_env()

# instalando os modelos de linguagem na env
spacy_download_langmodel(model = "en", envname = "r-reticulate")
spacy_download_langmodel(model = "pt_core_news_sm", envname = "r-reticulate")
#model = c('pt_core_news_sm', 'en'),
#virtualenv_root = "~/.rstudio-desktop/virtualenvs/spacy" )

# especificando o uso do modelo em portugues
spacy_initialize(model = "pt_core_news_sm",
                 save_profile = TRUE)

#virtualenv = "~/.rstudio-desktop/venv/spacy_virtualenv",
```

Carregando outros pacotes importantes para esse projeto

```
set.seed(123)
```

```
library(rvest) #web scraping
```

```
library(stringr) #manipulacao dos dados
```

```
library(tidyverse) #colecão de pacotes e acessórios
```

```
library(tm) #mineração de texto
```

```
library(igraph) #grafos
```

```
library(wordcloud) #nuvem de palavras
```

```
library(urltools) #manipulação de url
```

```
library(gtools) #funções de suporte
```

Acessando as páginas

Abaixo uma função auxiliar para acessar os links

```
scrape_post_links <- function(site) {  
  # coletando dados da HTML  
  source_html <- read_html(site)  
  # pegando atributos do titulo dos links  
  # buscando tages nos H2 header tags  
  links <- source_html %>%  
    html_nodes("div.widget--info__text-container") %>%  
    html_nodes("a") %>%  
    html_attr("href")  
  # filtrando quaisquer titulo que sejam NA  
  links <- links[!is.na(links)]  
  # retornando vetor com filtros  
  return(links)  
}
```

Fazendo interações em 20 páginas

```
root <- "https://g1.globo.com/busca/?q=BNDES"  
# obtendo o URL da pagina que vamos explorar  
  
all_pages <- c(root, paste0(root, "&page=", 1:20))  
  
# usando a funcao auxiliar que criamos para explorar o titulo de cada postagem  
all_links <- lapply(all_pages, scrape_post_links)  
  
# colapsar em um vetor  
all_links <- unlist(all_links)
```

Criando uma função auxiliar para extrair os links (1/2)

A URL está contida em parâmetro do link “u”

```
extract_urls <- function(raw_url) {  
  params <- urltools::param_get(raw_url)  
  scraped_url <- params$u  
  return (url_decode(scraped_url))  
}  
  
cleaned_links <- lapply(all_links, extract_urls)  
  
# Não estamos interessados em vídeos do globoplay app  
cleaned_links <- Filter(function(x) !any(grepl("globoplay", x)),  
  cleaned_links)
```

Saída com os dados em uma lista, com total de 237 links, com títulos extraídos das 20 páginas do site (2/2)

Name	Type	Value
cleaned_links	list [237]	List of length 237
[[1]]	character [1]	'https://g1.globo.com/economia/noticia/2021/04/29/bndes-...
[[2]]	character [1]	'https://g1.globo.com/economia/noticia/2021/04/29/bndes-...
[[3]]	character [1]	'https://g1.globo.com/economia/noticia/2021/04/30/no-mel...
[[4]]	character [1]	'https://g1.globo.com/ce/ceara/noticia/2021/04/16/secretari...
[[5]]	character [1]	'https://g1.globo.com/economia/noticia/2021/04/14/preside...
[[6]]	character [1]	'https://g1.globo.com/economia/noticia/2021/01/25/bndes-...
[[7]]	character [1]	'https://g1.globo.com/economia/noticia/2021/01/27/bndes-...
[[8]]	character [1]	'https://g1.globo.com/economia/noticia/2021/01/12/com-r-...
[[9]]	character [1]	'https://g1.globo.com/economia/noticia/2021/05/03/credito...
[[10]]	character [1]	'https://g1.globo.com/economia/noticia/2021/04/29/bndes-...
[[11]]	character [1]	'https://g1.globo.com/economia/noticia/2021/04/29/bndes-...
[[12]]	character [1]	'https://g1.globo.com/economia/noticia/2021/04/30/no-mel...
[[13]]	character [1]	'https://g1.globo.com/ce/ceara/noticia/2021/04/16/secretari...
[[14]]	character [1]	'https://g1.globo.com/economia/noticia/2021/04/14/preside...
[[15]]	character [1]	'https://g1.globo.com/economia/noticia/2021/01/25/bndes-...
[[16]]	character [1]	'https://g1.globo.com/economia/noticia/2021/01/27/bndes-...
[[17]]	character [1]	'https://g1.globo.com/economia/noticia/2021/01/12/com-r-...
[[18]]	character [1]	'https://g1.globo.com/economia/noticia/2021/05/03/credito...
[[19]]	character [1]	'https://g1.globo.com/economia/negocios/noticia/2021/01/...

Acessando de cada link (1/2)

```
scrape_post_body <- function(site) {  
  # Escape 404 Not Found Errors  
  try(  
    text <- site %>%  
      read_html %>%  
      html_nodes("article") %>%  
      html_nodes("p.content-text__container") %>%  
      html_text  
  )  
}  
  
data <- lapply(cleaned_links, scrape_post_body)  
data <- lapply(data,  
  function(item) paste(unlist(item),  
    collapse = ""))
```

Saída com as palavras dos títulos já extraídos dos links das 20 páginas do site (2/2)

Name	Type	Value
data	list [237]	List of length 237
[[1]]	character [9]	' O Banco Nacional de Desenvolvimento Econômico e Social (BNDES) ...
[[2]]	character [15]	' O Banco Nacional de Desenvolvimento Econômico e Social (BNDES) ...
[[3]]	character [30]	' As contas do setor público consolidado registraram superávit prim...
[[4]]	character [18]	' A Secretaria da Saúde do Ceará (Sesa), junto com o Banco Nacional ...
[[5]]	character [7]	' O presidente do Banco Nacional de Desenvolvimento Econômico e S...
[[6]]	character [6]	' O Banco Nacional de Desenvolvimento Econômico e Social (BNDES) i...
[[7]]	character [7]	' A diretoria do Banco Nacional de Desenvolvimento Econômico e Soc...
[[8]]	character [6]	' O Banco Nacional de Desenvolvimento Econômico e Social (BNDES) ...
[[9]]	character [14]	' A concessão de crédito pelos bancos brasileiros atingiu recorde no ...
[[10]]	character [9]	' O Banco Nacional de Desenvolvimento Econômico e Social (BNDES) ...
[[11]]	character [15]	' O Banco Nacional de Desenvolvimento Econômico e Social (BNDES) ...
[[12]]	character [30]	' As contas do setor público consolidado registraram superávit prim...
[[13]]	character [18]	' A Secretaria da Saúde do Ceará (Sesa), junto com o Banco Nacional ...
[[14]]	character [7]	' O presidente do Banco Nacional de Desenvolvimento Econômico e S...
[[15]]	character [6]	' O Banco Nacional de Desenvolvimento Econômico e Social (BNDES) i...
[[16]]	character [7]	' A diretoria do Banco Nacional de Desenvolvimento Econômico e Soc...
[[17]]	character [6]	' O Banco Nacional de Desenvolvimento Econômico e Social (BNDES) ...
[[18]]	character [14]	' A concessão de crédito pelos bancos brasileiros atingiu recorde no ...
[[19]]	character [5]	' O Banco Nacional de Desenvolvimento Econômico e Social (BNDES) ...

Pré-processamento dos textos para criação da nuvem de palavras (1/3)

```
# convertendo todos os titulos em minusculas
```

```
cleaned <- tolower(data)
```

```
# removendo numeros dos titulos
```

```
cleaned <- removeNumbers(cleaned)
```

```
# removendo stopwords e algumas palavras especificas que nao sao relevantes para a nuvem
```

```
cleaned <- removeWords(cleaned, c(stopwords("pt"), "bndes", "banco", "empresa", "empresas",  
  "desenvolvimento", "social", "afirmou", "ano", "dois", "fim", "disse", "rio", "janeiro",  
  "dia", "ter", "antes", "ser", "meses", "segundo", "nacional", "presidente", "ainda",  
  "nesta", "sobre", "parte", "total", "sido", "desde", "agora", "nota", "vai", "pode",  
  "informou", "havia", "outros", "diz", "atual", "junto", "feitos", "grandes", "menos",  
  "quase", "novos", "porque", "apenas", "ativar", "acesso", "abril", "pediu",  
  "acrescentou", "pessoas", "blocos", "apoio", "financiar", "quase", "anos", "afirma",  
  "desta", "capacidade", "pequenas", "maior", "sob", "segunda", "trabalho", "cargo",  
  "outras", "outro", "abaixo", "alta", "todo", "contra", "primeira", "fazer", "cidade",  
  "marco", "onde", "principais", "equipes", "outra", "duas", "nada", "feita", "destacou",  
  "feira", "partir", "forma", "importante", "privado", "longo", "exemplo", "equipe",  
  "aumento", "podem", "novo", "conforme", "entanto", "maiores", "medida", "garantir",  
  "final", "atualmente", "neste", "plano", "conta", "deste", "primeiro", "assim", "deve", "valores",  
  "modelo", "primeiro", "vez", "quanto", "feito", "brasileira", "disso", "meio", "disso", "qualquer",  
  "todas"))
```

Pré-processamento dos textos para criação da nuvem de palavras (2/3)

```
# removendo pontuacao
```

```
cleaned <- removePunctuation(cleaned)
```

```
# removendo espacos no comeco e final de cada titulo
```

```
cleaned <- str_trim(cleaned)
```

Pré-processamento dos textos para criação da nuvem de palavras (3/3)

```
# convertendo o vetor dos titulos em corpus
cleaned_corpus <- Corpus(VectorSource(cleaned))

# vaporizando cada palavra em cada titulo
# cleaned_corpus <- tm_map(cleaned_corpus, stemDocument)
doc_object <- TermDocumentMatrix(cleaned_corpus)
doc_matrix <- as.matrix(doc_object)

# contando aparicoes de cada palavra
counts <- sort(rowSums(doc_matrix),decreasing=TRUE)

# filtrando quaisquer palavras que nao contenham letras
counts <- counts[grepl("^[a-z]+$", names(counts))]

# criando quadro de dados a partir de informacoes de frequencia de palavras
frame_counts <- data.frame(word = names(counts), freq = counts)
```

Saída com as palavras apresentadas em um matriz e a frequência de aparições de cada uma delas em cada texto:

	1	2	3	4	5	6	7	8	9	10	11	12
micro	2	3	0	0	0	0	0	0	1	2	3	0
mil	2	1	0	2	0	0	0	0	0	2	1	0
milhões	1	1	0	1	0	1	2	2	1	1	1	0
modalidade	2	0	0	0	0	0	0	0	0	2	0	0
negociada	1	0	0	0	0	0	0	0	0	1	0	0
negócios	1	1	0	0	2	0	0	0	0	1	1	0
nessa	1	0	0	0	0	0	0	0	1	1	0	0
nova	2	2	0	0	0	0	0	0	0	2	2	0
novamente	1	0	0	0	0	0	0	0	0	1	0	0
oferecer	1	0	0	0	0	0	0	0	0	1	0	0
oferecida	1	0	0	0	0	0	0	0	0	1	0	0
operações	1	0	0	0	0	0	0	1	0	1	0	0
original	1	0	0	0	0	0	0	0	0	1	0	0
pagamento	1	1	2	0	0	0	1	1	1	1	1	2
pagamentos	5	3	0	0	0	0	0	0	0	5	3	0
pandemia	1	2	1	1	0	1	0	0	4	1	2	1
parcelas	1	3	0	0	0	0	0	0	3	1	3	0
participações	1	1	0	0	0	0	0	0	0	1	1	0
passado	2	2	8	0	1	0	0	0	0	2	2	8
pausa	1	3	0	0	0	0	0	0	0	1	3	0

Showing 63 to 82 of 6,563 entries, 237 total columns

Criação da nuvem de palavras em wordcloud e wordcloud2

```
wordcloud(words = frame_counts$word,  
          freq = frame_counts$freq,  
          scale=c(4, .2),  
          min.freq = 4,  
          max.words=150, random.order=FALSE,  
          rot.per=0,  
          colors=brewer.pal(8, "Dark2"))
```

```
install.packages("wordcloud2")  
library(wordcloud2)  
wordcloud2(data = frame_counts)
```


Resultados

No nível mais superficial de observação, é evidente o destaque à palavra GOVERNO, já que o Banco Nacional de Desenvolvimento Econômico e Social (BNDES) é uma empresa pública federal, cujo principal objetivo é o financiamento de longo prazo e investimento em todos os segmentos da economia brasileira, palavras que também apareceram em nossa nuvem. Obviamente, JAIR BOLSONARO e os nomes do atual presidente e seu antecessor GUSTAVO MONTEZANO e JOAQUIM LEVY também são apontadas.

Se nota algumas postagens ligadas a comparações sobre a situação atual e passada do banco, sobre seu novo modelo econômico e relacionadas a pipeline de modelagem e desenvolvimento de seus projetos. E também aparecem palavras direcionadas a prováveis reportagens sobre aprovação de financiamento para empresas diversas, escolha de consórcios para estudos de desestatizações dos Correios e Cedae, e ações emergenciais contra efeitos da pandemia de Covid-19.

De forma geral, ganharam destaque pela coocorrência em textos e notícias, conteúdo referente a investigações pelo TCU por desvios e fraudes em contratos e programas de pagamentos de Estados ao banco.

Nosso próximo passo, será submeter esse banco de palavras a um processo de desambiguação de entidades e uma limpeza mais aprofundada antes de podermos transformá-las em um grafo.

Extraindo todas as entidades dos textos

```
entities <- spacy_extract_entity(unlist(data))  
head(entities)  
tail(entities)
```

```
> head(entities)  
  doc_id      text ent_type start_id length  
1  text1 Banco Nacional de Desenvolvimento Econômico e Social      ORG         3       7  
2  text1      BNDES      MISC        11       1  
3  text4      BNDES      MISC        17       1  
4  text4      Bruno Laskowski      PER        37       2  
5  text4      Participações      LOC        42       1  
6  text4      Mercado de Capitais      LOC        44       3  
  
> tail(entities)  
  doc_id      text ent_type start_id length  
4983 text1993 Comissões Parlamentares de      PER        46       3  
4984 text1994  Supremo Tribunal Federal      LOC        24       3  
4985 text1994      Poder Legislativo      ORG        41       2  
4986 text1994      Congresso Nacional      LOC        52       2  
4987 text1995      Celso de Mello      PER         2       3  
4988 text1995      STF      LOC         9       1
```

Entidades extraídas com base nos atributos de tag, analisadas pela biblioteca spaCy. O modelo entendeu, por exemplo, que as entidades: Celso de Mello e Bruno Laskowski são pessoas (PER) e Banco Nacional e Desenvolvimento Econômico e Social e Poder Legislativo seriam organizações (ORG).

Criando a lista de adjacências

```
# Criando funcao auxiliar para dar suporte a geracao da lista de adjacencias
```

```
get_adjacent_list <- function(edge_list) {  
  if(length(edge_list)>2)  
    adjacent_matrix <- combinations(  
      length(edge_list), 2, edge_list)  
  #return(adjacent_matrix)}  
}
```

```
adjacent_matrix <- edges %>%  
  lapply(get_adjacent_list) %>%  
  reduce(rbind)
```

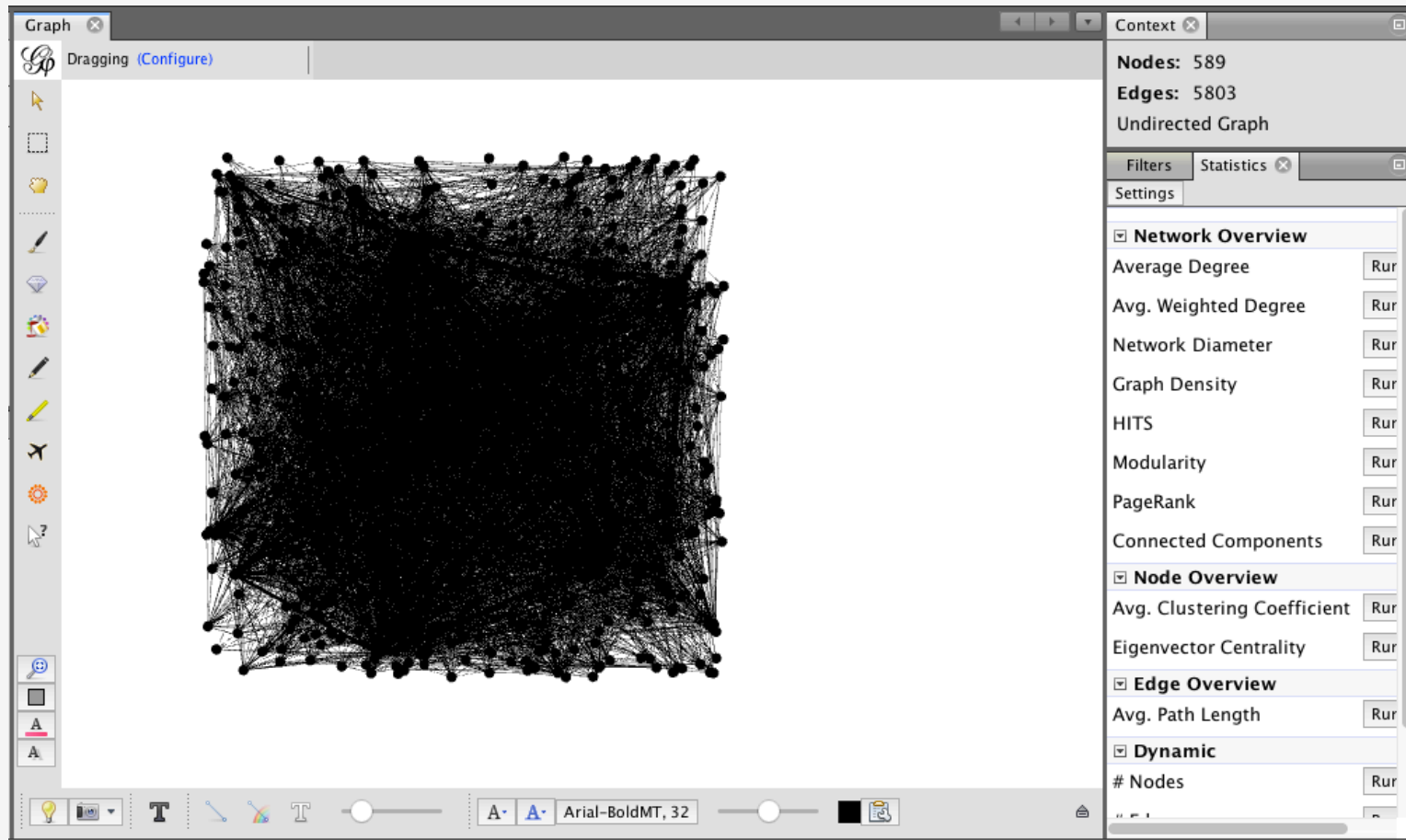
Criando o objeto grafo que será exportado para o GEPHI.

```
df <- as_tibble(adjacent_matrix, colnames=c('source', 'target'))  
  
weighted_edgelist <- df %>%  
  group_by(V1, V2) %>%  
  summarise(weight=n())  
  
news_graph <- weighted_edgelist %>% graph_from_data_frame(directed=F)  
write_graph(news_graph, 'news_graph.graphml', 'graphml')
```

Usando a GEPHI

Após a criação do objeto grafo, exportaremos esse arquivo para o GEPHI, que é uma ferramenta para visualização, análise e manipulação de redes e grafos.

O uso de técnicas de manipulação e apresentação de dados levam a novas possibilidades de análises. Nesse projeto iremos construir uma representação, e naturalmente será proposto um debate sobre a formação de comunidades e alguns “nós” principais, que disseminam e moldam a circulação dos sentidos.



Nosso arquivo abriu com 589 nós e 5803 arestas. A partir daí vamos começar os ajustes.

Na barra de filtros, em topologia, foi utilizado o filtro “componente gigante”, o que consequentemente diminuiu sutilmente a quantidade de nós e arestas. Na barra distribuição foi selecionado o filtro “ForceAtlas2”, após esse passo já se nota o aparecimento de comunidades no grafo. Foi utilizado no menu alternativas de comportamento o tópico “dissuadir hubs” com a intenção de espalhá-los.

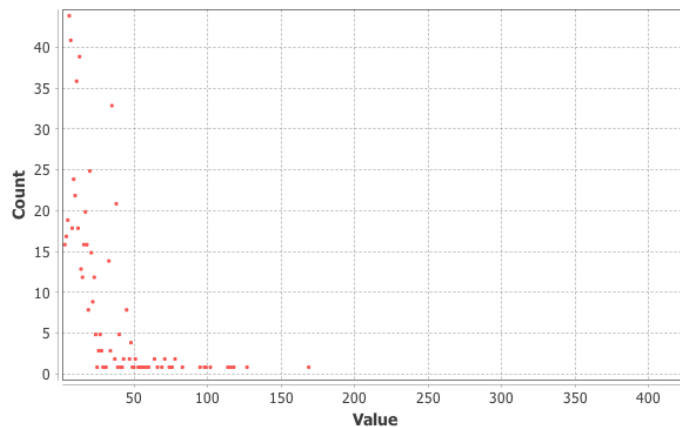
Logo após, partimos para geração de métricas no menu de estatísticas, na figura 1 o grau médio e na figura 2 o grau médio ponderado, ambos apresentaram uma curva com formato em cauda longa.

Degree Report

Results:

Average Degree: 19.705

Degree Distribution

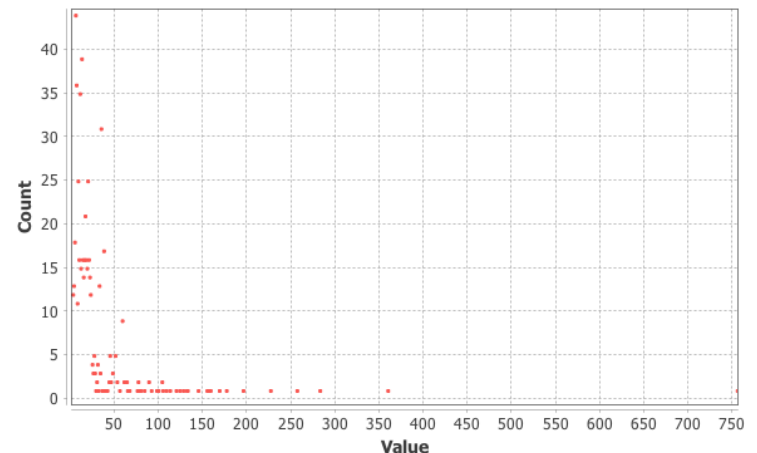


Weighted Degree Report

Results:

Average Weighted Degree: 25.205

Degree Distribution



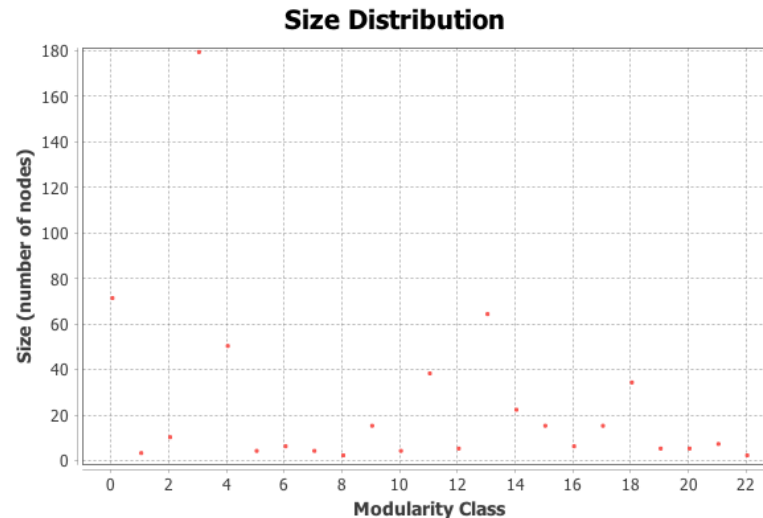
Em seguida, foi ajustado o tamanho dos nós usando o parâmetro grau e no menu distribuição foi selecionada a opção que possibilitar evitar sobreposição, foram afinados também o dimensionamento e a gravidade para trazer equilíbrio. Feito isso foi gerado um relatório de modularidade, que como ilustrado na figura abaixo, classificou a nossa base em 23 comunidades, utilizamos essa feature para colorir nosso grafo com esse parâmetro.

Parameters:

Randomize: On
Use edge weights: On
Resolution: 1.0

Results:

Modularity: 0.564
Modularity with resolution: 0.564
Number of Communities: 23

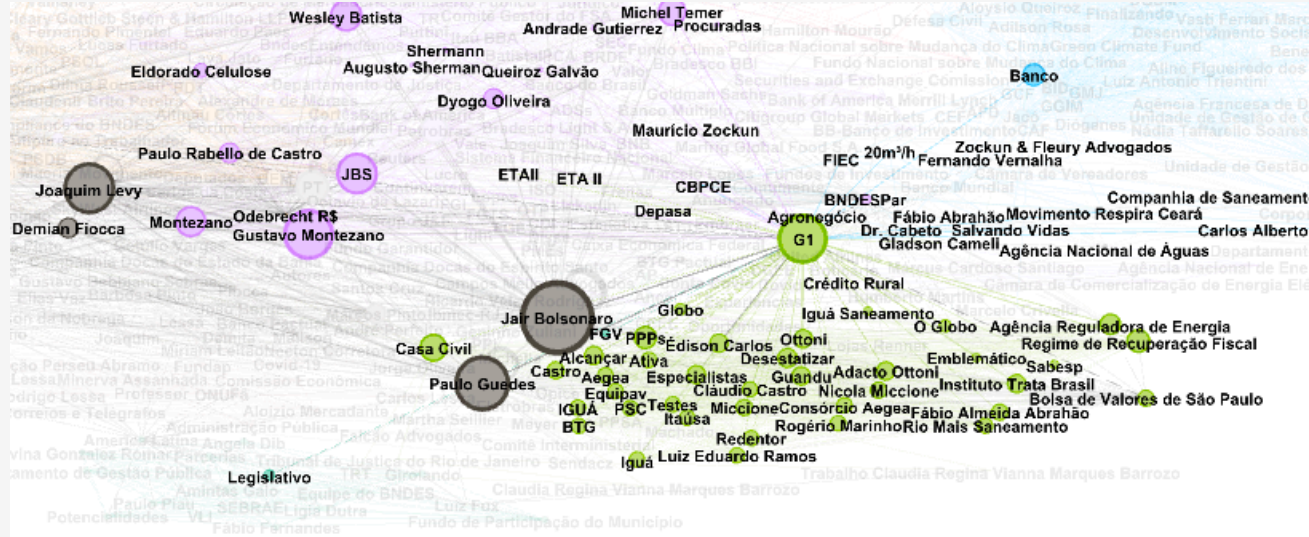


O foco agora está em analisar algumas comunidades, para facilitar esse processo foram ajustados no laboratório de dados os nomes e labels; logo após foi selecionada a opção “ajustar rótulos”, desta forma a visualização ficou mais sofisticada nos permitindo extrair suas características. Em seguida foi feita uma análise detalhada no laboratório de dados e após a identificação de alguns nós irrelevantes, optamos por simplesmente excluir alguns deles, mesclamos também algumas entidades que nitidamente se referiam ao mesmo nó, levando em consideração os graus de frequências e quantidade de conexões.

Após esses ajustes, foram recalculadas algumas estatísticas e na figura ilustrada no próximo slide já podemos ver nossas redes tomarem forma.

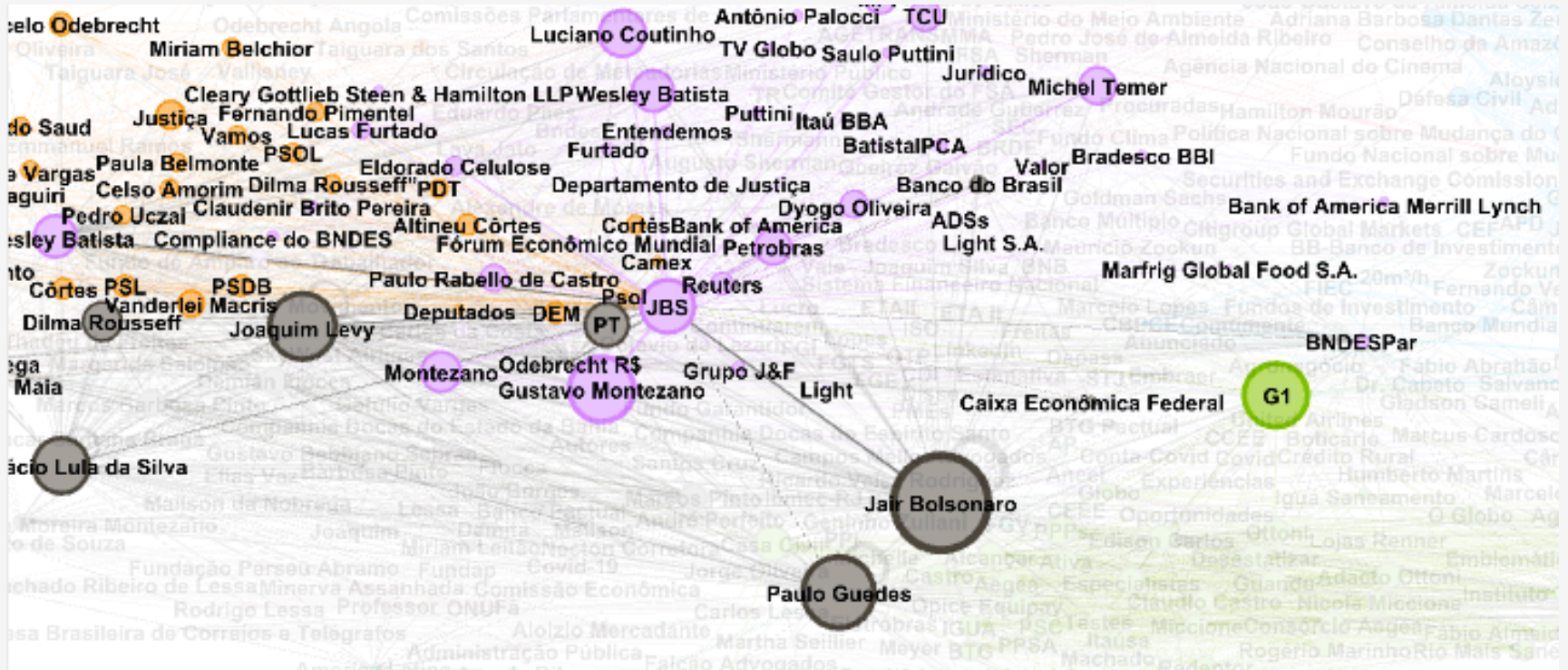


CONCLUSÕES – Comunidade em verde



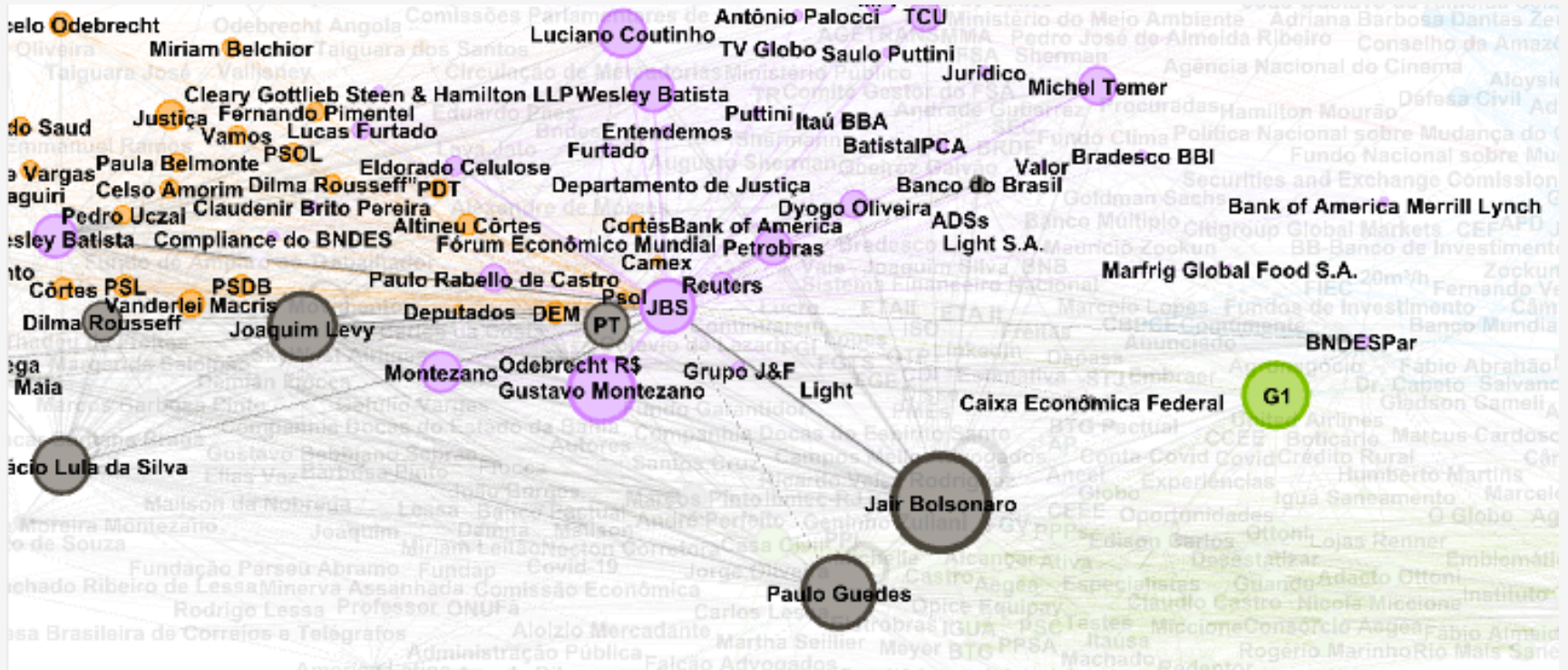
- Na comunidade ilustrada em verde se nota links entre o governo federal, o presidente Jair Bolsonaro, o ministro Paulo Guedes e entidades relacionadas a linhas de financiamento e crédito (Crédito Rural) e estruturação de projetos de concessões e parcerias público privadas (PPPs), com essas ligações obviamente relacionadas ao atual presidente do BNDES (Gustavo Montezano) e seu antecessor (Joaquim Levy) e Nicola Miccione (Secretário da Casa Civil do RJ) que participou das negociações sobre a concessão da CEDAE. Consequentemente o então governador em exercício do RJ, Claudio Castro também apresenta links.
- Se nota também movimentações direcionadas a venda de ações (Bolsa de Valores de SP).
- Diante mão, se nota uma relevante população relacionada ao setor de saneamento, mais especificamente, as vendas de participações em consórcios habilitados e na implantação de investimentos e desestatização nesse setor, como no caso da Iguá Saneamento, Trata Brasil e da aparição de seu presidente Edison Carlos e de Adacto Ottoni, que é Professor Associado do Departamento de Engenharia Sanitária e Ambiental da UERJ.
- Rogério Marinho (Ministro do Desenvolvimento Social), que foi recentemente acusado por Guedes por furar o teto de gastos; e Luiz Eduardo Ramos (Ministro Chefe da Secretaria de Governo), também aparecem bem relacionados na comunidade.

CONCLUSÕES – Comunidade em laranja



- Na comunidade laranja se observa uma concentração relevante de partidos políticos e de nomes relacionados a câmara dos deputados.
- Nomes como os de Altineu Cortês e Vanderlei Macris, ambos relatores da CPI do BNDES, que pediram o indiciamento de dezenas de pessoas e depoimentos como por exemplo do ex presidente do BNDES, Paulo Rabello de Castro, após acusações de possíveis ilícitos nas operações financeiras do banco, enquanto presidia.
- Links relacionados a Claudenir Brito Pereiram Diretor de Compliance e Riscos do BNDES, também demonstram relevância.

CONCLUSÕES – Comunidade em lilás



- Em lilás se observa nomes diretamente relacionados a escândalos e a investigações, como indicado pela forte presença de empresas como Odebrecht, JBS e seu presidente executivo Wesley Batista, no caso dessa última, foi noticiado diversas vezes que ela foi a empresa que mais recebeu dinheiro do BNDES. A Petrobrás também aparece com relevância, provavelmente por ter tido seus escândalos de corrupção noticiados com frequência.
- Nomes como os de Dyogo Oliveira, que assumiu a presidência do banco estatal no Governo Temer, e de Luciano Coutinho, presidente do BNDES por 10 anos também apresentam grau de relevância por terem participado das investigações; da mesma forma aparece Saulo Puttini, que já foi diretor jurídico do banco e veio a público reforçar a expressão “caixa-preta”, que se tornou popular por se reforçar que alguns fatos contribuíram para que o BNDES passasse a ser visto como uma “caixa-preta” do governo.

