

# Food Facts

*Rodolfo Viana*

*10-02-2016*

Utilizando os dados do Open Food Facts database (<http://world.openfoodfacts.org/data>), um gratuito, aberto e colaborativo database sobre comida ao redor do mundo, foram feitas análises sobre a concentração de cafeína, proteína e sódio nos alimentos.

Nos dados iniciais temos uma tabela que contém as seguintes colunas:

- code (text)
- url (text)
- creator (text)
- created\_t (text)
- created\_datetime (text)
- last\_modified\_t (text)
- last\_modified\_datetime (text)
- product\_name (text)
- generic\_name (text)
- quantity (text)
- packaging (text)
- packaging\_tags (text)
- brands (text)
- brands\_tags (text)
- categories (text)
- categories\_tags (text)
- categories\_en (text)
- origins (text)
- origins\_tags (text)
- manufacturing\_places (text)
- manufacturing\_places\_tags (text)
- labels (text)
- labels\_tags (text)
- labels\_en (text)
- emb\_codes (text)
- emb\_codes\_tags (text)
- first\_packaging\_code\_geo (text)
- cities (text)
- cities\_tags (text)
- purchase\_places (text)
- stores (text)
- countries (text)
- countries\_tags (text)
- countries\_en (text)
- ingredients\_text (text)
- allergens (text)
- allergens\_en (text)

- traces (text)
- traces\_tags (text)
- traces\_en (text)
- serving\_size (text)
- no\_nutriments (numeric)
- additives\_n (numeric)
- additives (text)
- additives\_tags (text)
- additives\_en (text)
- ingredients\_from\_palm\_oil\_n (numeric)
- ingredients\_from\_palm\_oil (numeric)
- ingredients\_from\_palm\_oil\_tags (text)
- ingredients\_that\_may\_be\_from\_palm\_oil\_n (numeric)
- ingredients\_that\_may\_be\_from\_palm\_oil (numeric)
- ingredients\_that\_may\_be\_from\_palm\_oil\_tags (text)
- nutrition\_grade\_uk (numeric)
- nutrition\_grade\_fr (text)
- pnns\_groups\_1 (text)
- pnns\_groups\_2 (text)
- states (text)
- states\_tags (text)
- states\_en (text)
- main\_category (text)
- main\_category\_en (text)
- image\_url (text)
- image\_small\_url (text)
- energy\_100g (numeric)
- energy\_from\_fat\_100g (numeric)
- fat\_100g (numeric)
- saturated\_fat\_100g (numeric)
- butyric\_acid\_100g (numeric)
- caproic\_acid\_100g (numeric)
- caprylic\_acid\_100g (numeric)
- capric\_acid\_100g (numeric)
- lauric\_acid\_100g (numeric)
- myristic\_acid\_100g (numeric)
- palmitic\_acid\_100g (numeric)
- stearic\_acid\_100g (numeric)
- arachidic\_acid\_100g (numeric)
- behenic\_acid\_100g (numeric)
- lignoceric\_acid\_100g (numeric)
- cerotic\_acid\_100g (numeric)
- montanic\_acid\_100g (numeric)
- melissic\_acid\_100g (numeric)
- monounsaturated\_fat\_100g (numeric)
- polyunsaturated\_fat\_100g (numeric)
- omega\_3\_fat\_100g (numeric)
- alpha\_linolenic\_acid\_100g (numeric)
- eicosapentaenoic\_acid\_100g (numeric)

- docosaehaenoic\_acid\_100g (numeric)
- omega\_6\_fat\_100g (numeric)
- linoleic\_acid\_100g (numeric)
- arachidonic\_acid\_100g (numeric)
- gamma\_linolenic\_acid\_100g (numeric)
- dihomogamma\_linolenic\_acid\_100g (numeric)
- omega\_9\_fat\_100g (numeric)
- oleic\_acid\_100g (numeric)
- elaidic\_acid\_100g (numeric)
- gondoic\_acid\_100g (numeric)
- mead\_acid\_100g (numeric)
- erucic\_acid\_100g (numeric)
- nervonic\_acid\_100g (numeric)
- trans\_fat\_100g (numeric)
- cholesterol\_100g (numeric)
- carbohydrates\_100g (numeric)
- sugars\_100g (numeric)
- sucrose\_100g (numeric)
- glucose\_100g (numeric)
- fructose\_100g (numeric)
- lactose\_100g (numeric)
- maltose\_100g (numeric)
- maltodextrins\_100g (numeric)
- starch\_100g (numeric)
- polyols\_100g (numeric)
- fiber\_100g (numeric)
- proteins\_100g (numeric)
- casein\_100g (numeric)
- serum\_proteins\_100g (numeric)
- nucleotides\_100g (numeric)
- salt\_100g (numeric)
- sodium\_100g (numeric)
- alcohol\_100g (numeric)
- vitamin\_a\_100g (numeric)
- beta\_carotene\_100g (numeric)
- vitamin\_d\_100g (numeric)
- vitamin\_e\_100g (numeric)
- vitamin\_k\_100g (numeric)
- vitamin\_c\_100g (numeric)
- vitamin\_b1\_100g (numeric)
- vitamin\_b2\_100g (numeric)
- vitamin\_pp\_100g (numeric)
- vitamin\_b6\_100g (numeric)
- vitamin\_b9\_100g (numeric)
- vitamin\_b12\_100g (numeric)
- biotin\_100g (numeric)
- pantothenic\_acid\_100g (numeric)
- silica\_100g (numeric)
- bicarbonate\_100g (numeric)

- potassium\_100g (numeric)
- chloride\_100g (numeric)
- calcium\_100g (numeric)
- phosphorus\_100g (numeric)
- iron\_100g (numeric)
- magnesium\_100g (numeric)
- zinc\_100g (numeric)
- copper\_100g (numeric)
- manganese\_100g (numeric)
- fluoride\_100g (numeric)
- selenium\_100g (numeric)
- chromium\_100g (numeric)
- molybdenum\_100g (numeric)
- iodine\_100g (numeric)
- caffeine\_100g (numeric)
- taurine\_100g (numeric)
- ph\_100g (numeric)
- fruits\_vegetables\_nuts\_100g (numeric)
- collagen\_meat\_protein\_ratio\_100g (numeric)
- cocoa\_100g (numeric)
- chlorophyll\_100g (numeric)
- carbon\_footprint\_100g (numeric)
- nutrition\_score\_fr\_100g (numeric)
- nutrition\_score\_uk\_100g (numeric)

Por se tratar de mais de 50M de dados, foi decidido que iriamos trabalhar apenas com as seguintes colunas:

- product\_name
- generic\_name
- quantity
- brands\_tags
- categories\_en
- origins\_tags
- stores
- countries\_en
- serving\_size
- no\_nutriments
- additives\_n
- main\_category\_en
- trans\_fat\_100g
- fiber\_100g
- proteins\_100g
- salt\_100g
- sodium\_100g
- calcium\_100g
- caffeine\_100g
- image\_url
- ingredients\_text

```
library(dplyr)
library(ggplot2)
```

```
file <- read.csv("food_facts.csv", sep=";")

file$sodium_100g <- as.character(file$sodium_100g)
file$sodium_100g[is.na(file$sodium_100g)] <- 0
file$sodium_100g <- gsub(",", ".", file$sodium_100g)
file$sodium_100g <- as.numeric(file$sodium_100g)

file$caffeine_100g <- as.character(file$caffeine_100g)
file$caffeine_100g[is.na(file$caffeine_100g)] <- 0
file$caffeine_100g <- gsub(",", ".", file$caffeine_100g)
file$caffeine_100g <- as.numeric(file$caffeine_100g)

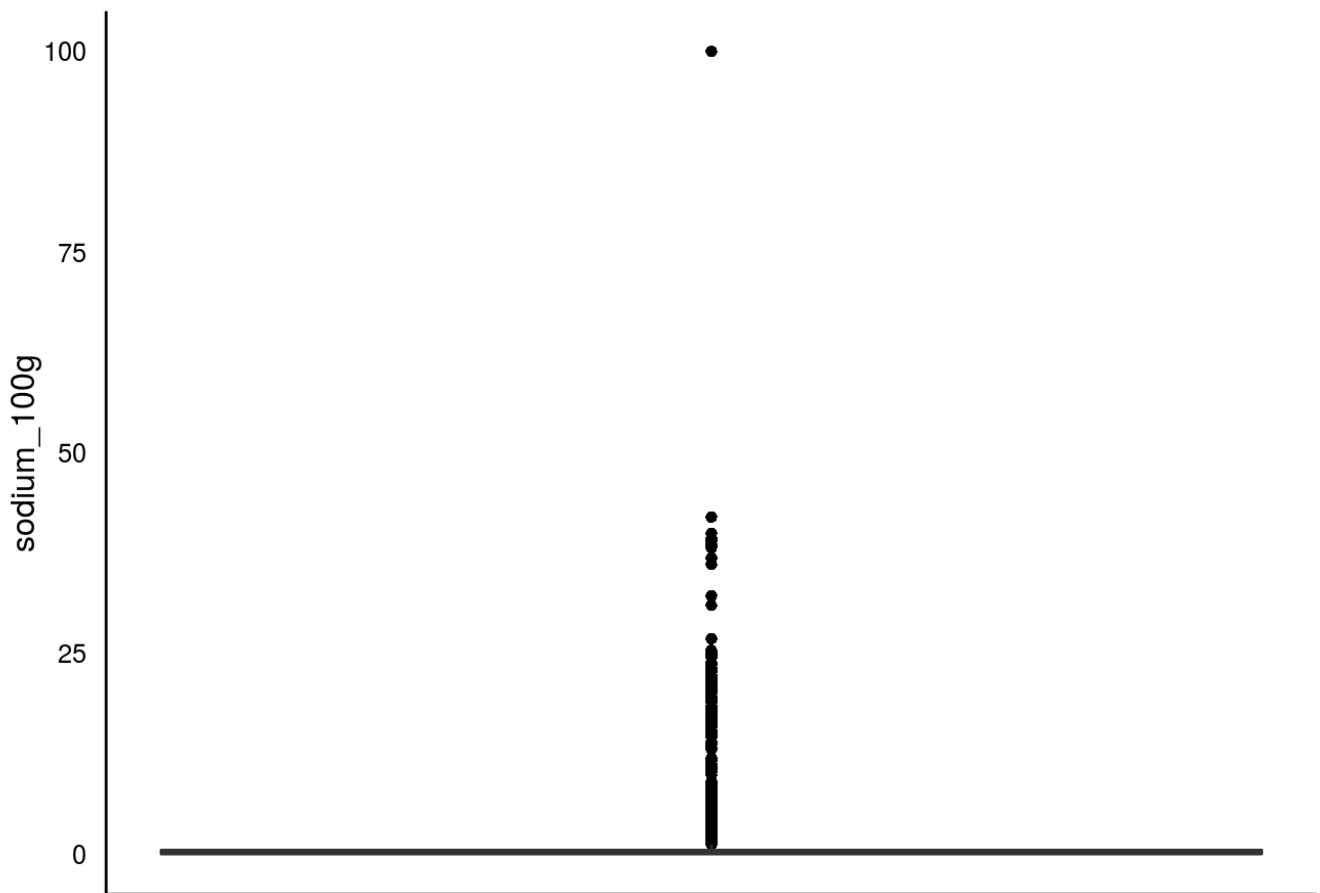
file$proteins_100g <- as.character(file$proteins_100g)
file$proteins_100g[is.na(file$proteins_100g)] <- 0
file$proteins_100g <- gsub(",", ".", file$proteins_100g)
file$proteins_100g <- as.numeric(file$proteins_100g)

file_sodium <- filter(file, sodium_100g != 0)
file_caffeine <- filter(file, caffeine_100g != 0)
file_proteins <- filter(file, proteins_100g != 0)
```

A nossa primeira curiosidade foi descobrir como que era a distribuição de Sódio, Proteína e Cafeína em todos os alimentos. Para uma melhor visualização da distribuição dos dados utilizamos os boxplot. Por se tratar de um banco de dados colaborativo, existe um grande número de valores em branco. Por causa disso, filtramos e retiramos os alimentos com valores NA para proteína, cafeína e sódio.

Para o sódio temos o seguinte boxplot:

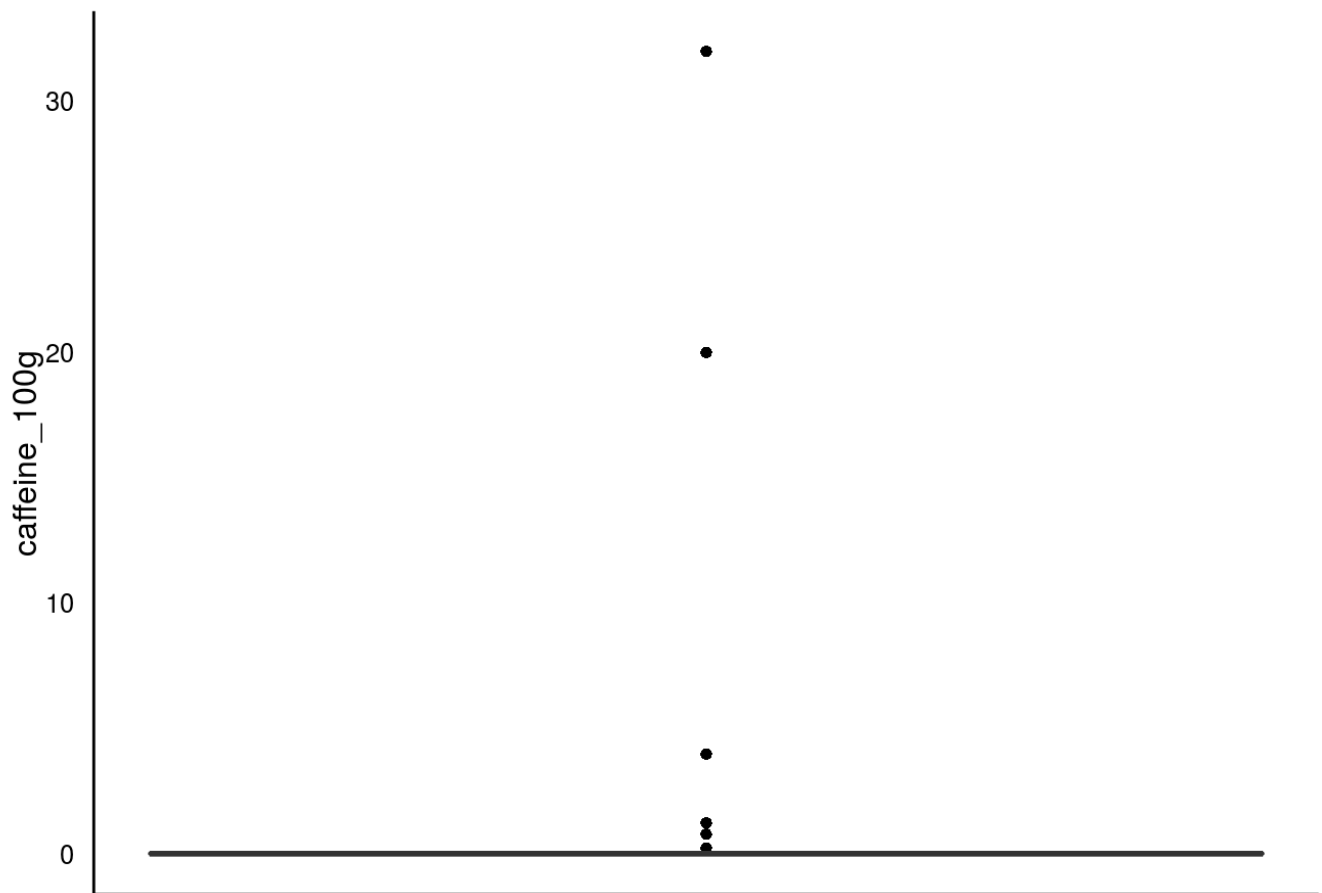
```
ggplot(file_sodium, aes(1, sodium_100g)) +
  geom_boxplot() +
  theme_classic() +
  theme(axis.ticks = element_blank(),
        axis.text.x=element_blank(),
        axis.title.x=element_blank(),
        legend.position="none")
```



É possível notar que existe um grande número de outliers e que tanto a média, 1 e 3 quartil estão próximos a zero. Interessante observar que existe um alimento com mais da metade da concentração de sódio dos demais.

Para a cafeína temos o seguinte boxplot:

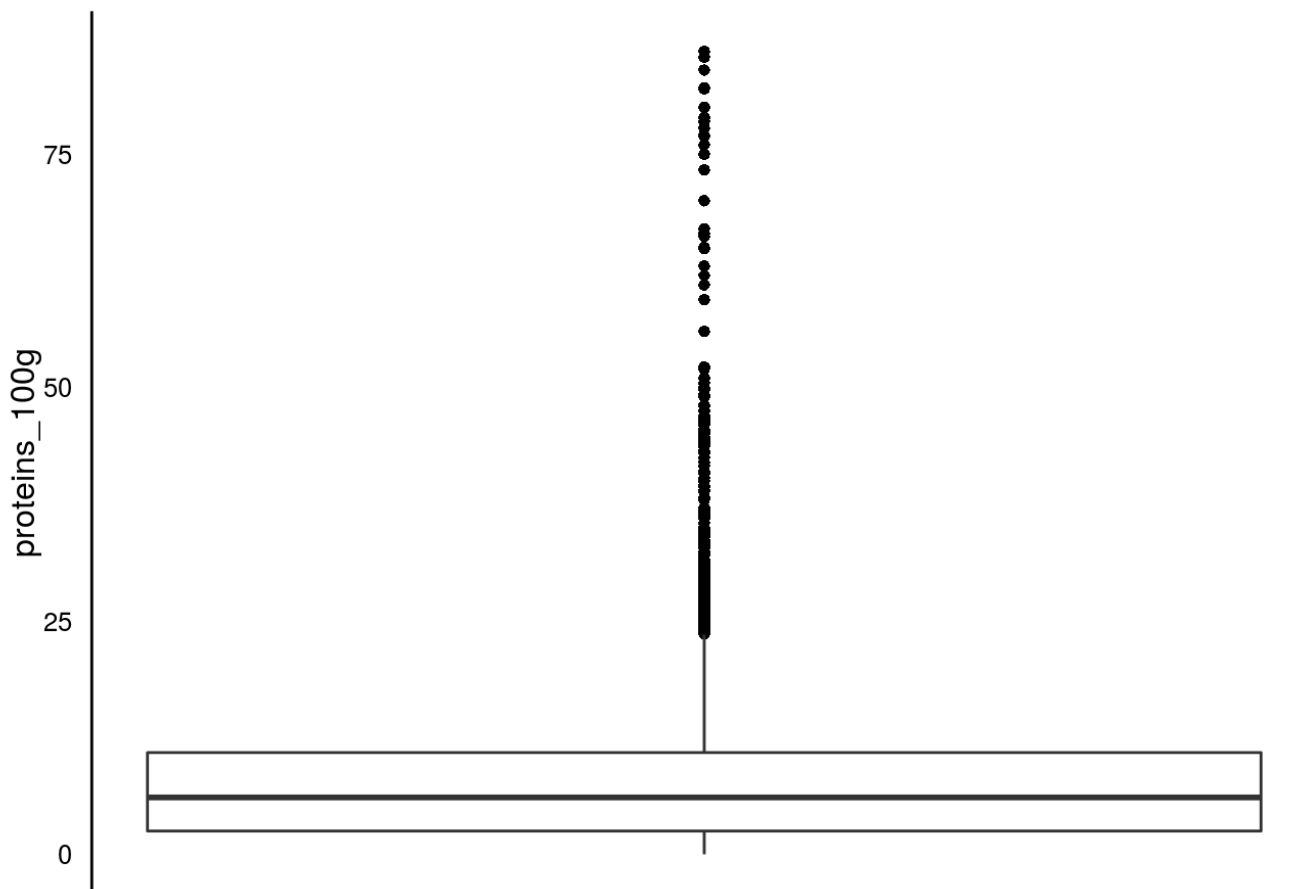
```
ggplot(file_caffeine, aes(1, caffeine_100g)) +  
  geom_boxplot() +  
  theme_classic() +  
  theme(axis.ticks = element_blank(),  
        axis.text.x = element_blank(),  
        axis.title.x = element_blank())
```



A cafeína possui menos outliers que o sódio, porém a média, 1 e 3 quartil também se encontra perto de zero.

Para a proteína temos o seguinte boxplot:

```
ggplot(file_proteins, aes(1, proteins_100g)) +  
  geom_boxplot() +  
  theme_classic() +  
  theme(axis.ticks = element_blank(),  
        axis.text.x=element_blank(),  
        axis.title.x=element_blank(),  
        legend.position="none")
```



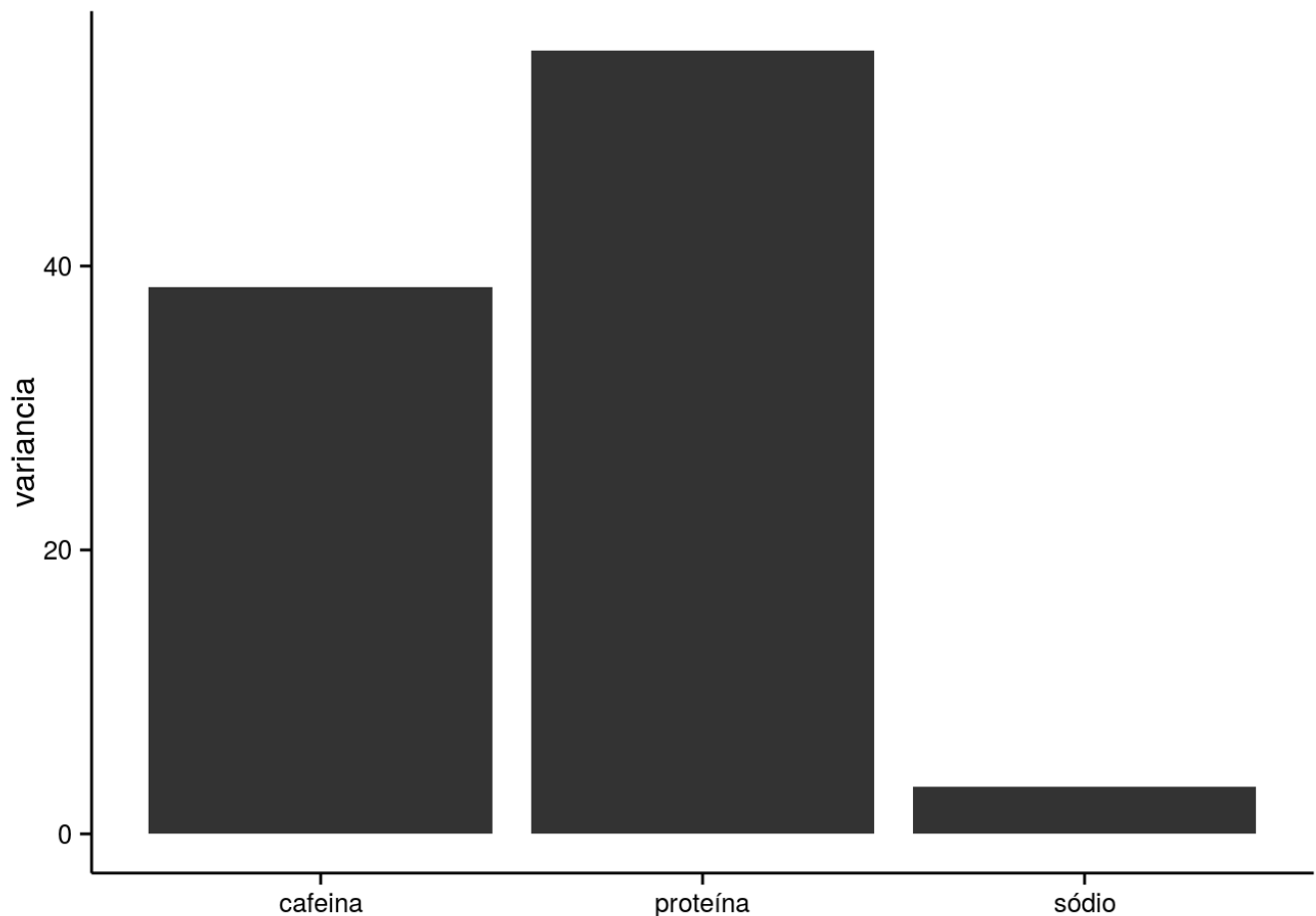
A proteína, diferente das outras substâncias, possui um boxplot com 1 e 3 quartil com valores diferentes. Mostrando dessa forma que a proteína possui diversos alimentos com diferentes grau de concentração

A nossa segunda curiosidade foi descobrir qual das três substancias possui maior variância nos alimentos

```
df <- data.frame(tipo = c("sódio", "cafeína", "proteína"),
                  variancia = c(var(file_sodium$sodium_100g), var(file_caffeine
                    $caffeine_100g), var(file_proteins$proteins_100g)))

ggplot(df, aes(x = tipo, y = variancia)) +
  geom_bar(stat = "identity") +
  theme_classic() +
  theme(axis.title.x = element_blank())
```





É possível notar que a proteína é a substância que possui maior variância e o sódio menor variância de concentração nos alimentos

A nossa última curiosidade foi descobrir qual alimento possui a maior concentração de Sódio, Proteína e Cafeína

Para a cafeína temos o alimento:

```
file_caffeine <- file_caffeine[order(-file_caffeine$caffeine_100g),]
file_caffeine[1,]$product_name
```

```
## [1] Red Bull energy drink
## 50932 Levels:   알 \U0001f37a 통깨 짜왕 ゆず 미역 울무차 콩두유 ... 黑瓶眼药水
```

foto do alimento (<http://en.openfoodfacts.org/images/products/22220768/front.8.400.jpg>)

O que não foi nenhuma surpresa, já que o RedBull é conhecido por ser um alimento que pode te proporcionar uma “energia extra”

Para a proteína temos:

```
file_proteins <- file_proteins[order(-file_proteins$proteins_100g),]
file_proteins[1,]$product_name
```

```
## [1] Blattgelatine weiss
## 50932 Levels:   알 \U0001f37a 통깨 짜왕 ゆず 미역 울무차 콩두유 ... 黑瓶眼药水
```

```
file_proteins[1,]$image_url
```

```
## [1] http://en.openfoodfacts.org/images/products/20153465/front.8.400.jpg  
## 61200 Levels: ...
```

foto do alimento (<http://en.openfoodfacts.org/images/products/20153465/front.8.400.jpg>)

Uma especie de gelatina vendida na Alemanha

Para o sódio tivemos uma surpresa, pois o alimento cadastrado como tendo a maior concentração de sódio não possuía praticamente nenhum registro (nome, localidade, foto, etc). Por se tratar de um conjunto de dados aberto e colaborativo isso é considerado “normal”. Por essa razão resolvemos observar o alimento na segunda posição:

```
file_sodium <- file_sodium[order(-file_sodium$sodium_100g),]  
file_sodium[2,]$product_name
```

```
## [1] Himalayan Pink Salt  
## 50932 Levels: 알 \U0001f37a 통깨 짜왕 ゆず 미역 울무차 콩두유 ... 黑瓶眼药水
```

Foto do alimento (<http://en.openfoodfacts.org/images/products/009/661/991/1936/front.6.400.jpg>)

O alimento encontrado foi um sal vendido nos Estados Unidos