



**Servicio de vigilancia y detección de insights y oportunidades para la oferta exportable de alimentos con soporte en la gastronomía y tendencias post-COVID19'**

**Lima, Enero 2021**

# TEXT MINING & SENTIMENT ANALYSIS

Promperú - 2021



- 
- A close-up photograph of many blueberries, which are covered in a fine layer of white frost or dew. The berries are a deep blue color and are densely packed together, filling the entire frame.
1. Software utilizado
  2. Modelos analíticos implementados
  3. Metodología
  4. Proceso de análisis de Insights



1. **Procesador: 3.2 GHz Core I7**
2. **Memoria RAM: 32 GB**
3. **R version: 3.6.2**
4. **R studio version: 1.3**

### Razones para usar R:

1. R es gratuito
2. Es intuitivo
3. Esta avalado por una comunidad científica mundial.
4. Constantemente aparecen nuevos paquetes gratuitos que expande la capacidad de r
5. Crea gráficos de calidad superior a otros paquetes
6. Es compatible con 'todos' los formatos de datos

### Razones para usar Power BI:

1. Es intuitivo
2. Compatible con el modelo de datos actual
3. Control y seguridad de información
4. Actualizaciones periódicas
5. Múltiple visibilidad (PC, móvil)

- 
- A close-up photograph of many blueberries, which are covered in a fine layer of white powder (bloom) and some water droplets. The berries are a deep blue color and are densely packed together.
1. Software utilizado
  2. Modelos analíticos implementados
  3. Metodología
  4. Proceso de análisis de Insights

1

Evolución de las menciones en las semanas de estudio

Las líneas de tiempo que permiten ver la estacionalidad del producto en el periodo de estudio

2

Keyword Cloud/ frequency

Una nube de palabras que permite ver cuáles se citan más, qué otras palabras acompañan, etc.

3

Análisis de bigramas

Grafo de palabras que se combinan de dos en dos para dar significado a los mensajes.

4

Topic Modeling

Cluster de términos, para ver cómo se agrupan, qué resúmenes se podrían hacer, etc.

5

Sentiment Analysis

Análisis del sentimiento, tanto positivo y negativos, así como por emociones básicas de la persona (alegría, tristeza, miedo, rechazo, etc.).



- 
- A close-up photograph of many blueberries, which are covered in a fine layer of white frost or sugar. The berries are a deep blue color and are densely packed together, filling the entire frame. The lighting is soft, highlighting the texture of the berries and the frost.
1. Software utilizado
  2. Modelos analíticos implementados
  3. Metodología
  4. Proceso de análisis de Insights

## Proceso Promperú

### Análisis y Exploración

### Limpieza de datos y tokenización

### Construcción de marco de datos Tidy

### Modelamiento

### Visualización de datos

- Construcción de matriz de datos
- Filtro de negocio.
- Estudio de variables

- Limpieza de textos.
- Filtro de stopwords
- Tokenización
- Lematizado de palabras

*Se realiza la limpieza de la data conservando textos relevantes para el estudio*

- Construcción de marco de datos considerando palabras tokenizadas.

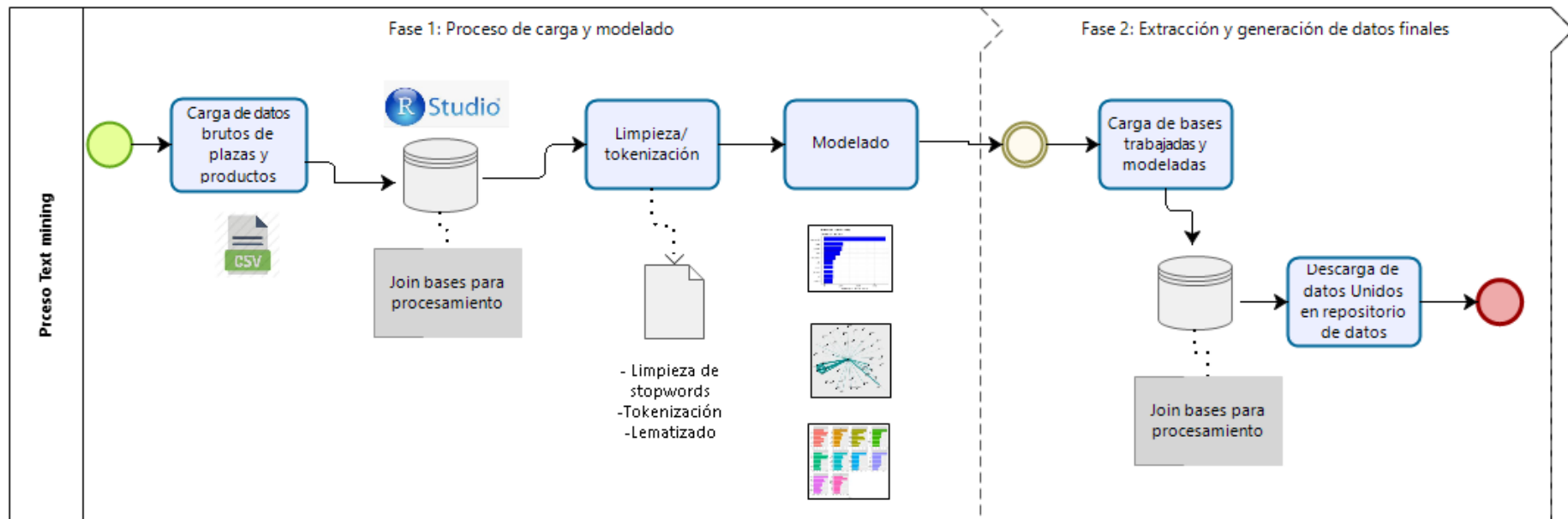
- Análisis de tendencias
- Frecuencia de palabras mas recurrentes
- Análisis clustering
- Análisis de bigramas
- Análisis de sentimientos
- Generación de tablas de consumo

- Carga de tablas de consumo.
- Creación de variables adicionales.
- Actualización y publicación de reporte.

**Retroalimentación**



## Proceso Promperú



- 
- A close-up photograph of many blueberries, which are covered in a fine layer of white frost or sugar. The berries are a deep blue color and are densely packed together, filling the entire frame. The lighting is soft, highlighting the texture of the berries.
1. Software utilizado
  2. Modelos analíticos implementados
  3. Metodología
  4. Proceso de análisis de Insights

A.

## Diccionario de variables requeridas

Variable	Descripción
url	La URL donde está el comentario de texto.
sentimiento	Obtenido por la herramienta de escucha (quizás nos aporte algo, lo he dejado por si las moscas, aunque prefiero hacerlo a partir del texto disponible)
autor	Quién lo escribió, por si quisieran localizar algún blog o identidad de redes sociales de referencia (influencia)
texto	El contenido de texto en sí. El campos más relevante, lógicamente.
hashtags	Quizás sea interesante sacar cuáles son los más citados, para ver tendencias, modas, etc.
impacto	Cuánto gente ha reaccionado al texto (mide a ver el éxito de una publicación)
impresiones	Cuánto gente lo ha llegado a leer (es lógicamente una cifra mayor al impacto)
localización	Desde dónde se ha escrito en el mercado de referencia, por si quisiera analizar también por ubicaciones.
día	Día del periodo de estudio (octubre a diciembre de 2020)
hora	hora
plaza	País en análisis
producto	Producto en análisis

## A. Preparando el entorno de trabajo

### *Limpieza del entorno de trabajo*

```
7 # clean the workspace
8 rm(list = ls())
9 cat("\014")
10 # fijamos a UTF-8
11 options(encoding = "utf-8")
```

### *Carga de funciones diseñadas*

Functions	
calcular...	function (dtm)
cvLDA	function (Ntopics, dtm, ...)

### *Carga de librerías necesarias para trabajar*

```
26 # Cargamos las librerías que vamos a necesitar
27 library(readxl) # Para leer archivos excel
28 library(tidyverse) # Para las operaciones con datos
29 library(syuzhet) # Libreria para emociones
30 library(tidytext)
31 library(stringr) # Para operar con datos de tipo St
32 library(stopwords) # Para poder quitar las stopwords
33 library(ggplot2) # Librería de visualización gráficas
34 library(lubridate) # Para el formateo de fechas y s
35 library(scales) # Para trabajar con datos de coma
36 library(igraph) # Para el análisis de bigramas en fo
```



## B.

## ANÁLISIS Y EXPLORACIÓN

### Carga y concatenación de datos según productos y plazas

```
206 arandano_esp <- read.csv("2. Datos brutos - Data 2020/productos/arandano_esp.csv",encoding = "Latin1",sep = ";")
207 arandano_fr <- read.csv("2. Datos brutos - Data 2020/productos/arandano_fr.csv",encoding = "Latin1",sep = ";")
208 arandano_uk <- read.csv("2. Datos brutos - Data 2020/productos/arandano_uk.csv",encoding = "Latin1",sep = ";")
209
```

### Construcción de variables

```
210 # Creando variables plaza y producto
211 arandano_esp$plaza <- "Espana"
212 arandano_fr$plaza <- "Francia"
213 arandano_uk$plaza <- "Reino_Unido"
214 arandano <- rbind(arandano_esp,arandano_fr,arandano_uk)
215
216 rm(arandano_esp,arandano_fr,arandano_uk)
217 arandano$producto <- "arandano"
218
219 names(arandano) <- c("x","url","sentimiento","autor","texto","hashtags","impacto",
220 "impresiones","localizacion","fecha","hora","plaza","producto")
```

### Joins bases

```
390 dfpromperu2020 <- rbind(arandano,cafe,palta,pisco,quinua,uva,superfood)
391 dfpromperu2020$X <- NULL
392 names(dfpromperu2020)
```

Número de mensajes producto - mercado

Cantidades en Miles

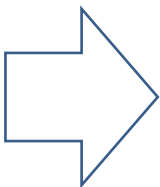


## C.

## LIMPIEZA DE DATOS

### 1 Seleccionar el producto y plaza que se quiere analizar

```
432 # corremos considerando base genérica
433 # -----
434 # Carga de datos para un producto en específico (todos los países)
435 datos <- arandano
436 # datos <- cafe
437 # datos <- palta
438 # datos <- pisco
439 # datos <- quinua
440 # datos <- uva
441 # datos <- superfood
442
443 # Filtro por plaza (país específico)
444 # -----
445 # -----
446 datos <- datos %>% filter(plaza == 'Espana')
447 # datos <- datos %>% filter(plaza == 'Francia')
448 # datos <- datos %>% filter(plaza == 'Reino_Unido')
449
450 table(datos$producto, datos$plaza)
451 # -----
```



### 2 Limpieza de textos

```
453 # Vamos a hacer un poco de limpieza de texto
454 # -----
455 # espana$texto<-gsub("#[A-Za-z0-9]+|[A-Za-z0-9]+|\\w+(?:\\.\\w+)*\\/\\S+", "",
456 datos$texto<-gsub("@[A-Za-z0-9]+|\\w+(?:\\.\\w+)*\\/\\S+", "", datos$texto)
457 datos$texto <- gsub("RT", "", datos$texto)
458 datos$texto <- gsub("https://", "", datos$texto)
459 datos$texto <- gsub("http://", "", datos$texto)
460
461 datos$texto <- tolower(datos$texto) # Se convierte todo el texto a minúsculas
462 datos$texto <- str_replace_all(datos$texto, "http\\S*", "") # Eliminaci?n de http
463 datos$texto <- str_replace_all(datos$texto, "[[:punct:]]", " ") # Eliminaci?n de puntuaci?n
464 datos$texto <- str_replace_all(datos$texto, "[[:digit:]]", " ") # Eliminaci?n de n?meros
465 datos$texto <- str_replace_all(datos$texto, "[\\s]+", " ") # Eliminaci?n de espacios
466 datos$texto <- chartr('áéíóüñ', 'aeiouñ', datos$texto)
```

### (1) Evolución de las menciones en las semanas de estudio



## E.

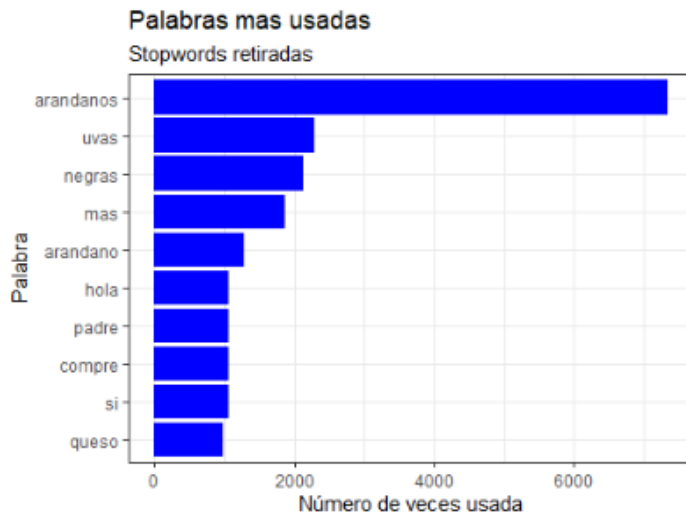
## MODELAMIENTO DE DATOS

### (2) Keyword Cloud/ Frequency

```
661 # Representacion grafica de las frecuencias
662 df_tidy %>% group_by(token) %>% summarise(n1 = n()) %>%
663   arrange(desc(n1)) %>% slice_max(order_by = n1, n = 10) %>%
664   ggplot(aes(x = reorder(token,n1), y = n1)) +
665   geom_col(show.legend = TRUE,fill = "blue") +
666   theme_bw() +
667   labs(y = "", x = "") +
668   theme(legend.position = "none") +
669   coord_flip() +
670   labs(title = "Palabras mas usadas",
671        subtitle = "Stopwords retiradas",
672        x = "Palabra",
673        y = "Número de veces usada")
```

### Backup de base de datos

```
677 # Para guardar la bases
678 # -----
679 df_backup <- df_tidy %>% group_by(token) %>% summarise(n1 = n()) %>%
680   arrange(desc(n1)) %>% slice_max(order_by = n1, n = 10)
681 |
682 df_backup$producto <- "arandano"
683 df_backup$plaza <- "españa"
684
685 df_backup <- df_backup[c(3,4,1,2)]
686 write.csv(df_backup,'3. Datos Procesados - 2020/df_tidy(uva_UK).csv',row.names = F)
```



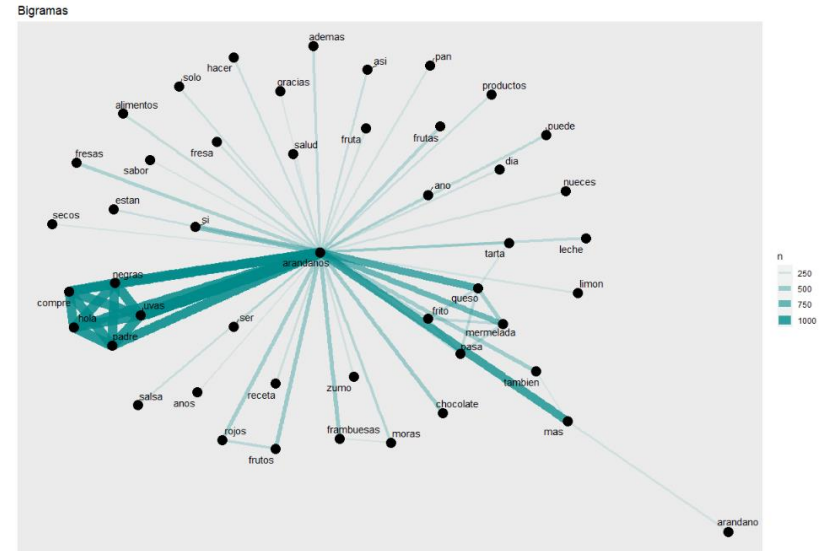
# MODELAMIENTO DE DATOS

### (3) Análisis de bigramas

```
757 # A veces nos interesa entender la relacion entre palabras en una opini?n.
758 review_bigrams <- df_bigrams %>%
759   unnest_tokens(bigram, frase, token = "ngrams", n = 2) # separamos token 2 - grams
760 bigrams_separated <- review_bigrams %>%
761   separate(bigram, c("word1", "word2"), sep = " ") # separamos word por bigrama
762 bigrams_filtered <- bigrams_separated %>%
763   filter(!word1 %in% lexicon$word) %>%
764   filter(!word2 %in% lexicon$word) # eliminamos stop words por bigrama
765 bigram_counts <- bigrams_filtered %>%
766   dplyr::count(word1, word2, sort = TRUE) # contamos la cantidad de words por bigrama
767 bigrams_united <- bigrams_filtered %>%
768   unite(bigram, word1, word2, sep = " ") # count bigrams cleaning
769 bigrams_united %>%
770   dplyr::count(bigram, sort = TRUE)
```

## Backup de base de datos

```
821 analisis_bigrams <- title_word_pairs %>% filter(n >= 500)
822
823 analisis_bigrams$producto <- "arandano"
824 analisis_bigrams$plaza <- "España"
825
826 analisis_bigrams <- analisis_bigrams[c(4,5,1,2,3)]
827 write.csv(analisis_bigrams, '3. Datos Procesados - 2020/analisis_bigrams(uva_uk).csv')
```





E.

## MODELAMIENTO DE DATOS

### (4) Topic Modeling

```
849 corpus <- corpus(datos$texto)
850 cdfm <- dfm(corpus, remove=c(lista_stopwords),
851             verbose=TRUE, remove_punct=TRUE, remove_numbers=TRUE)
852
853 # Quitamos palabras que solo salgan 1 vez
854 cdfm <- dfm_trim(cdfm, min_docfreq = 2, verbose=TRUE)
855
856 # Ahora lo exportamos a un formato para procesar los Topic Models.
857 dtm <- convert(cdfm, to="topicmodels")
858
859 # Calculamos ahora los topics óptimos
860 # calculartopics(dtm)
861 # Estimamos el LDA con el número óptimo de topics que nos haya salido
862 lda <- LDA(dtm, k = 10, method = "Gibbs",
863           control = list(verbose=25L, seed = 123, burnin = 100, iter = 500))
```

### Backup de base de datos

```
889 terminos$producto <- "arandano"
890 terminos$plaza <- "española"
891 terminos <- terminos[c(4,5,1,2,3)]
892
893 write.csv(terminos,'3. Datos Procesados - 2020/cluster_terminos(uva_uk).csv',row.names = F)
```



E.

## MODELAMIENTO DE DATOS

### (4) Topic Modeling

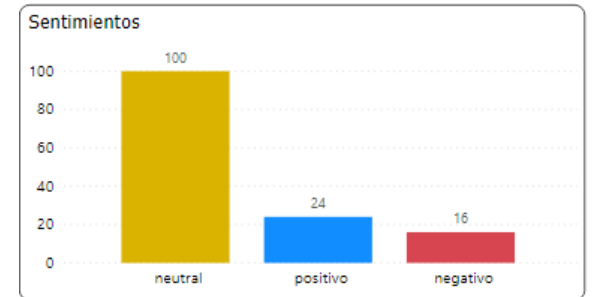
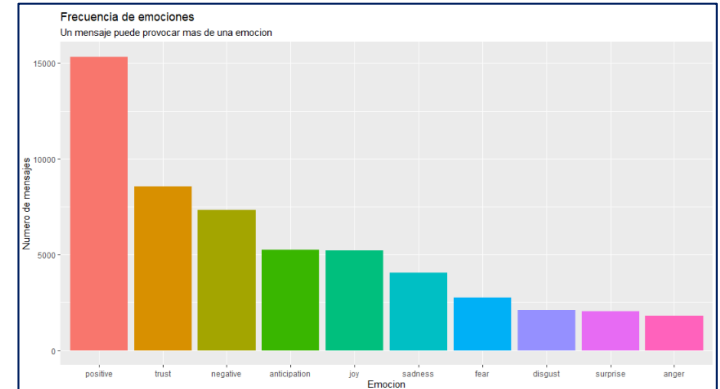
```

920 # Vamos a coger las palabras en español del diccionario NRC
921 nrc<- get_sentiment(datos$texto, method="nrc",lang="spanish")
922 #nrc<- get_sentiment(datos$texto, method="nrc",lang="french")
923 #nrc<- get_sentiment(datos$texto, method="nrc",lang="english")
924
925 # Obtenemos las emociones
926 emotions <- get_nrc_sentiment(datos$texto,lang="spanish")
927 #emotions <- get_nrc_sentiment(datos$texto,lang="french")
928 #emotions <- get_nrc_sentiment(datos$texto,lang="english")
929
930 emo_bar = colSums(emotions)
    
```

### Backup de base de datos

```

949 emo_sum$producto <- "quinua"
950 emo_sum$plaza <- "Reino_Unido"
951
952 emo_sum <- emo_sum[c(3,4,1,2)]
953
954 write.csv(emo_sum,'3. Datos Procesados - 2020/sentimientos (quinua_uk).csv',row.names = F)
    
```





**Servicio de vigilancia y detección de insights y oportunidades para la oferta exportable de alimentos con soporte en la gastronomía y tendencias post-COVID19'**

**Lima, Enero 2021**