

# Covid v/s Food

Jose Antonio Castro <sup>1,†</sup>, Juan Vicuña <sup>2,†</sup>, Emily Skog <sup>3,†</sup> and Clemente Cambara <sup>4,†</sup>

<sup>1</sup> Departamento de Ciencias de la Computación, Escuela de Ingeniería, Pontificia Universidad Católica de Chile, Santiago, Chile; jacastr18@uc.cl

<sup>2</sup> Departamento de Ciencias de la Computación, Escuela de Ingeniería, Pontificia Universidad Católica de Chile, Santiago, Chile; juan.vg@uc.cl

<sup>3</sup> Instituto de Ingeniería Biológica y Biomédica, Escuela de Ingeniería, Pontificia Universidad Católica de Chile, Santiago, Chile; emily.skog@uc.cl

<sup>4</sup> Departamento de Ciencias de la Computación, Escuela de Ingeniería, Pontificia Universidad Católica de Chile, Santiago, Chile; ccambara.d@uc.cl

† These authors contributed equally to this work.

## 1. Introducción

Este trabajo busca responder la interrogante: ¿Cómo se relacionan los hábitos alimenticios de distintos países con el impacto biológico que ha tenido el virus SARS-CoV-2 en cada uno de ellos? Para esto, se han usado datasets de números actualizados de casos totales confirmados, fallecidos y activos para 218 países [1] junto con datasets de porcentaje de ingesta alimentaria para estos países [2]. El plan es unir los datasets más relevantes utilizando los países como llave y realizar diversos análisis mediante el uso de técnicas de procesamiento de datos vistas en clase.

### Impacto del problema

Debido al impacto que el virus SARS-CoV-2 ha tenido en el mundo y considerando lo poco que se sabe de su comportamiento, es relevante estudiar distintos factores que tradicionalmente no son relacionados con la pandemia de Covid-19 y analizarlos a nivel estadístico para determinar la relación entre ellos. De esta manera, se espera entender más las implicancias que distintas circunstancias tienen sobre el virus y con esta información poder combatirlo mejor.

Se decidió investigar específicamente la alimentación dado que es un pilar fundamental en la salud y no suele ser correctamente dimensionada al momento de relacionarla con distintas complicaciones biológicas.

## 2. Metodología

El primer paso fue limpiar los datos y agruparlos en una matriz de tal forma que cada entrada representara un país y cada columna tuviera la siguiente información respecto a los hábitos alimenticios:

**Citation:** Castro, J.A.; Vicuña, J.; Skog, E.; Cambara, C. Covid v/s Food. *Journal Not Specified* **2021**, *1*, 0. <https://doi.org/>

Received:

Accepted:

Published:

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.

**Copyright:** © 2021 by the authors. Submitted to *Journal Not Specified* for possible open access publication under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

Alcoholic Beverages
Animal Products
Animal fats
Aquatic Products, Other
Cereals - Excluding Beer
Eggs
Fish, Seafood
Fruits - Excluding Wine
Meat
Milk - Excluding Butter
Miscellaneous
Offals
Oilcrops
Spices
Starchy Roots
Stimulants
Sugar Crops
Sugar & Sweeteners
Treenuts
Vegetal Products
Vegetable Oils
Vegetables
Obesity

**Tabla 1:** Esta tabla representa los datos de alimentación utilizados para determinar si existe una relación entre estos y el avance del coronavirus en el mundo. Es importante notar que cada columna representando un tipo de alimento es una proporción calculada como  $\frac{\text{kcal totales consumidas}}{\text{cantidad de habitantes en el país}}$  mientras que la columna **Obesity** es el porcentaje de gente con obesidad respecto a la cantidad total de habitantes en ese país.

Se midió la relevancia de esta matriz **X** de datos para predecir vectores relacionados al avance del virus SARS-CoV-2 en cada país de distintas maneras. Primero, dado un vector **y** que representa la proporción por país de  $\frac{\text{casos confirmados de coronavirus}}{\text{cantidad de tests totales realizados}}$ , se utilizó un método de regresión lineal entrenado con esta matriz para que predijera dos variaciones del vector **y**

1. Si la proporción de casos confirmados para un país particular esta sobre la media mundial (vector binario con valor 0 si esta proporción está bajo la media y 1 si está por sobre).
2. El cuartil en el que se encuentra esta proporción respecto al resto de países

Luego, se realizó esta misma metodología con un clasificador de Árbol de Decisiones, esta vez haciendo un testeo de hiper-parámetros en ambas versiones del vector **y** para encontrar aquellos que maximizaran el `accuracy_score`. En particular, se probó con las combinaciones de criterios `gini` y `entropy` y 100 profundidades distintas para maximizar la eficacia del clasificador.

De manera alterna, se decidió probar un clasificador SVC para modelar el problema, pero aun con un estudio pertinente de hiperparametros, no se logró un precisión lo adecuadamente alta como para consiferar útil. En más detalles, la mejor precisión (0.4) fué obtenida con kernel `poly` de con un valor de `C` de 7.

### 3. Materiales y Métodos

#### *Materiales*

Para obtener los resultados del proyecto, es necesario usar datasets de números actualizados de casos totales confirmados, fallecidos y activos para 218 países junto con datasets de porcentaje de ingesta alimentaria para estos países. Estos, como se explico

anteriormente, son utilizados para originar un nuevo dataset contenedor donde se pueda realizar diversos análisis. Para obtener los datos se pueden utilizar nuestra fuente de datos, la página web más grande de la comunidad de la ciencia de datos y análisis del tráfico web: Kaggle. [1] [2].

### Métodos

Diversas técnicas de procesamiento de datos fueron utilizadas para desarrollo del proyecto, ya sea para la limpieza de datos como para la clasificación de datos.

En primer lugar, para la unión de los datasets antes mencionados y la limpieza de datos, se hizo uso de la librería Pandas [3] junto con sus funciones integradas (como *df.dropna()* o *df.merge()*). Luego se llevó a cabo una función formulada por el equipo que calcula la correlación entre las columnas del dataset, para así quitarlas de este mismo.

Con el dataset ya preprocesado, se utilizó un método agrupador para formular el vector y según el criterio determinado (según la media mundial o según los cuartiles).

Teniendo los datos listos para entrenar (como un vector **X** y otro **y**), se eligió la librería SkLearn [4] para realizar el entrenamiento de las regresiones lineales, los árboles de decisión y el clasificador SVC. Para lograr lo anterior, las siguientes funciones más relevantes de la librería que fueron utilizadas son:

63

Funciones SKLearn utilizadas	
Función	Explicación
<i>train_test_split()</i>	Divide el vector <b>X</b> en <i>X_train</i> , <i>X_test</i> y el vector <b>y</b> en <i>y_train</i> , <i>y_test</i> respecto a una proporción deseada (0.2 para el proyecto).
<i>linear_model.LinearRegression()</i>	Crea un modelo de Regresión lineal a ser entrenado.
<i>tree.DecisionTreeClassifier()</i>	Crea un árbol de decisión de datos que clasifica los datos con nodos internos con atributos.
<i>metrics.accuracy_score</i>	Retorna la precisión de una predicción respecto a sus valores verdaderos.
<i>svm.SVC()</i>	Crea un modelo de SVM a ser entrenado.

64

Para el estudio de hiperparametros se utilizó una simple función iterativa que encontrara la mejor precisión para un conjunto de valores.

65

66

## 67 4. Resultados

68 Se mostraran los resultados generales del proyecto realizado en base a dos defini-  
69 ciones distintas de las muestras procesadas, es decir, en base a las definiciones del vector  
70  $y$  y de la matriz  $X$ .

### 71 4.1. Resultados Muestra Casos Confirmados / Test Totales

72 Para esta sección, el vector  $y$  es definido como  $\frac{\text{casos confirmados de coronavirus}}{\text{cantidad de tests totales realizados}}$  y la  
73 matriz  $X$  como las columnas restantes del dataset.

74 Mostraremos dos tabla con todos los clasificadores modelados y los resultados  
75 pertinentes. Cada tabala estara basada según cada variación del vector  $y$  que exista.

76 Siendo mas especificos, mostraremos los resultados para cuando el vector  $y$  esta basado  
77 en la media mundial de la proporción (asignando 0 y 1) y mostraremos otra tabla para  
78 cuando el vector  $y$  esta basado en los cuartiles de la proporción (asignando 0, 1, 2 y 3).

79

Vector y Media Mundial	
Clasificador	Resultados
Regresión Lineal	Error de media cuadratica: 0.14 y coeficiente de error cuadratico $r^2$ de 0.34.
Árbol de decisión	Mejor resultado posible con hiperparametros: $criterion = gini$ y $max\_depth = 20$ . Precisión obtenida de un 74.07%
SVM	Mejor resultado posible con hiperparametros: $kernel = poly$ , $C = 1$ y con una presición del 0.64.
Vector y Cuartiles	
Clasificador	Resultados
Regresión Lineal	Error de media cuadratica: 0.38 y coeficiente de error cuadratico $r^2$ de 0.78.
Regresión Lineal Entrenada	Error de media cuadratica: 0.14 y coeficiente de error cuadratico $r^2$ de 0.34.
Árbol de decisión	Mejor resultado posible con hiperparametros: $criterion = gini$ y $max\_depth = 2$ . Precisión obtenida de un 77.7%
SVM	Mejor resultado posible con hiperparametros: $kernel = poly$ , $C = 7$ y con una presición del 0.4.

### 82 4.2. Resultados Muestra Muertes / Tamaño población

83 Para esta sección, el vector  $y$  es definido como  $\frac{\text{muertes totales}}{\text{cantidad de población}}$  y la matriz  $X$  como  
84 las columnas restantes del dataset.

85 Mostraremos dos tabla con todos los clasificadores modelados y los resultados  
86 pertinentes. Cada tabala estara basada según cada variación del vector  $y$  que exista.

87 Siendo mas especificos, mostraremos los resultados para cuando el vector  $y$  esta basado  
88 en la media mundial de la proporción (asignando 0 y 1) y mostraremos otra tabla para  
89 cuando el vector  $y$  esta basado en los cuartiles de la proporción (asignando 0, 1, 2 y 3).

90

Vector y Media Mundial	
Clasificador	Resultados
Regresión Lineal	Error de media cuadrática: 0.20 y coeficiente de error cuadrático $r^2$ de 0.6.
Árbol de decisión	Mejor resultado posible con hiperparámetros: <i>criterion</i> = <i>entropy</i> y <i>max_depth</i> = 14. Precisión obtenida de un 51.85%
SVM	Mejor resultado posible con hiperparámetros: <i>kernel</i> = <i>rbf</i> , <i>C</i> = 3 y con una precisión del 0.76.
Vector y Cuartiles	
Clasificador	Resultados
Regresión Lineal	Error de media cuadrática: 0.664 y coeficiente de error cuadrático $r^2$ de 0.27.
Árbol de decisión	Mejor resultado posible con hiperparámetros: <i>criterion</i> = <i>gini</i> y <i>max_depth</i> = 2. Precisión obtenida de un 77.7%
SVM	Mejor resultado posible con hiperparámetros: <i>kernel</i> = <i>poly</i> , <i>C</i> = 3 y con una precisión del 0.55.

## 5. Discusión

Al analizar los resultados de los casos confirmados divididos por el total de tests, se puede notar que ciertos clasificadores generan una precisión más alta. Al segmentar el vector y según la media se puede apreciar la mejor precisión con el clasificador de árbol de decisión, variando los hiperparámetros (*criterion*=gini y *max\_depth*=20). Mientras que los peores resultados se dan con la regresión lineal. Esto tiene sentido, en el ámbito que la información del data set es no lineal, no sigue un patrón evidente. También, en el árbol de decisión al usar un "max\_depth" de valor moderado, no se genera "overfitting".

Luego, al cambiar la información del vector y en cuartiles, se puede evidenciar resultados similares: la regresión lineal entrenada genera precisión baja y el árbol de decisión la más alta (*criterion*=gini y *max\_depth*=2).

Por otro lado, al analizar los resultados de los muertos divididos por la población, nuevamente se puede notar que ciertos clasificadores generan una precisión más alta. Al segmentar el vector y según la media se puede apreciar la mejor precisión con *Support Vector Machines* (SVM), variando los hiperparámetros (*kernel*=rbf y *C*=3). Mientras que los peores resultados se dan con el árbol de decisión. Esto tiene sentido, en el ámbito que la información del data set es no lineal, no sigue un patrón evidente. Por lo que el kernel del SVM, que es no lineal, se ajusta de la mejor manera a esta información. También, sus valores de *slack* son un poco mayor que "1", por lo que da un poco de flexibilidad al segmentar la información.

Luego, al cambiar la información del vector y en cuartiles, se puede evidenciar resultados distintos: la regresión lineal entrenada genera precisión baja y el árbol de decisión la más alta (*criterion*=gini y *max\_depth*=2), notando que actúa de manera similar a la tasa de incidencia de casos.

Todos estos resultados han explicitado que hay cierta correlación entre la manera de alimentación de los países con la tasa de muertos y tasa de incidencia a raíz de COVID-19, ya que se genera una precisión para todos los casos mayor al 70%, lo cual es

124 suficiente para decir que no es mera coincidencia.

125

126 Sin embargo, la manera de segmentar el vector  $y$ , aunque demuestra una cor-  
127 relación, no es muy útil al momento de la aplicación, ya que solo nos puede decir si  
128 un país (a partir de su alimentación) tiene un riesgo leve, moderado, alto, muy alto, de  
129 contraer la enfermedad o morir de ella. No es muy específico.

130

131 Tampoco se analizó ciertas alimentaciones en particular, como consumo de alcohol  
132 o verduras. Lo cual habría sido interesante ver la tendencia de los países solo analizando  
133 ciertas columnas del data set.

## 134 6. Conclusión

135 Se puede concluir que se llegó a clasificar de manera exitosamente la incidencia  
136 y muertes a raíz del COVID-19 a partir de la alimentación de cada país. Se llegó para  
137 todos los casos a precisiones mayores al 70%, con el más alto rondando el 78%. En todos  
138 los casos esto fue con clasificadores no lineales, SVM y árbol de decisión, variando los  
139 hiperparámetros.

140

141 No obstante, a futuro se desea de una manera segmentar el vector  $y$ , para dar  
142 información más clara de la gravedad de la enfermedad en el país en cuestión, no solo en  
143 2 o 4 categorías. Además, se podría entrenar con otros métodos, como redes neuronales,  
144 que abarcan la no linealidad del problema.

145

146 También, sería interesante analizar como ciertas alimentaciones afecta la incidencia  
147 del virus en cada país, ¿cómo si se alimentan con más carne o más frutos tienen diferentes  
148 resultados?. Por lo que todavía hay tantas interrogantes que se pueden responder de  
149 este dataset, independiente de sus fallencias.

## 150 7. References

151

- 152 1. Joseph Assaker, [Covid-19 Global Dataset: Up to date numbers of daily Confirmed, Death](#)  
153 [and Active cases for 218 countries](#)
- 154 2. [Food and Agriculture Organization of the United Nations](#)
- 155 3. [Documentación de Pandas](#)
- 156 4. [Documentación de Scikit Learn](#)