

# Edx Capstone: Classification model Project

Luis G. Vazquez de Lara Cisneros.

02/12/2020

## Contents

<b>1</b>	<b>Introduction/overview/executive summary</b>	<b>1</b>
<b>2</b>	<b>Methods/analysis</b>	<b>4</b>
2.1	Data wrangling . . . . .	4
2.2	Data exploration . . . . .	6
2.3	Data preprocessing . . . . .	8
2.4	Selection of the classification algorithms and evaluation metrics. . . . .	10
2.5	Strategy of model generation . . . . .	12
<b>3</b>	<b>Results</b>	<b>12</b>
<b>4</b>	<b>Conclusion</b>	<b>17</b>
	<b>References</b>	<b>18</b>

## 1 Introduction/overview/executive summary

Covid-19 is the name of a disease caused by SARS-CoV-2, a novel type of coronavirus that has spread all over the world, presenting a severe threat to global health. The clinical presentation is very heterogeneous, ranging from an asymptomatic disease, to a life threatening condition with respiratory failure. At present, a specific therapy is lacking. Because the mortality of hospitalized patients with this disease varies from country to country, it is imperative to develop prediction models tailored to the reality of the location. In Mexico, as of 17 December 2020, the health authorities reported 1,289,298 confirmed patients, with 116,487 deaths, around 9.03% of reported cases (Secretaria de Salud 2020). It has also been noted that the mortality of hospitalized patients varies from institution to institution. *The Mexican Institute of Social Security* (IMSS, after the initials in Spanish), carries the biggest burden of public health care and the highest mortality of patients hospitalized with covid-19 (Sanchez Talanquer 2020).

Medical researchers often use prediction models as an aid to estimate the probability of risk for a specific outcome, to inform their decision making. In recent years, an initiative called *The Transparent Reporting of a multivariate prediction model for Individual Prognosis Or Diagnosis* (TRIPOD), made a statement consisting of a checklist with 22 items considered as essential for good reporting of these type of studies (Collins et al. 2015). This work is the final assignment of the capstone course of the *Edx Data Science Professional Certificate*, where we learned the basics of machine learning techniques. Even though the course

is introductory to this topic, I believe it gives the necessary knowledge to fulfill the TRIPOD requirements concerning the statistical analysis (table 1) and presentation of results (table 2).

Table 1: TRIPOD items concerning the statistical analysis (from (Collins et al. 2015)).

Topic	Item	Checklist item
Sample Size	8	Explain how the study size was arrived at.
Missing data	9	Describe how missing data were handled (ag., complete-case analysis, single imputation, multiple imputation) with details of any imputation method.
Statistical analysis methods	10a	Describe how predictors were handled in the analyses.
	10b	Specify type of model, all model-building procedures (including any predictor selection), and method for internal validation.
	10c	For validation, describe how the predictions were calculated.
	10d	Specify all measures used to assess model performance and, if relevant, to compare multiple models.
	10e	Describe any model updating (eg., re-calibration) arising from the validation, if done.
Development vs. validation	12	For validation, identify any differences from the development data in setting, eligibility criteria, outcome, and predictors.

Table 2: TRIPOD items concerning the presentation of results [from (Collins et al. 2015)].

Topic	Item	Checklist item
Participants	13a	Describe the flow of participants through the study, including the number of participants with and without the outcome and, if applicable, a summary of the follow-up time. A diagram may be helpful.
	13b	Describe the characteristics of the participants (basic demographics, clinical features, available predictors), including the number of participants with missing data for predictors and outcome.
	13c	For validation, show a comparison with the development data of the distribution of important variables (demographics, predictors and outcome).
Model development	14a	Specify the number of participants and outcome events in each analysis.
	14b	If done, report the unadjusted association between each candidate predictor and outcome.
Model specification	15a	Present the full prediction model to allow predictions for individuals (ie, all regression coefficients, and model intercept or baseline survival at a given time point).
	15b	Explain how to use the prediction model.
Model performance	16	Report performance measures (with CIs) for the prediction model.

I will use a database of hospitalized patients from a Mexican hospital and try to build an outcome prediction model. The database comprises information from 642 hospital records, with 51 variables. Table 3 shows the column names of the database, with a short description of their meaning and the codification or units employed. I will also use the TRIPOD items as the framework to build the statistical analysis and the presentation of results. Besides, I will try to demonstrate the skills acquired concerning data wrangling. The key steps are summarized as follows:

- Prepare the database for a suitable statistical analysis.
- Explore the data, to decide which variables will be used as predictors (*preprocessing*). In this step, I will use bivariate analysis and imputation methods.
- Split the data set into train and test sets.
- Use logistic regression, k-nearest neighbors, classification and regression trees (CART) and random forest as prediction algorithms. I will use bootstrap as a resampling method to tune-up the parameters. I will use the **caret** package to build the algorithms.
- Compare the performance of the algorithms with the *overall accuracy* in the test set.

Table 3: Description of the variables in the working database.

Column name	Description	Codification/units
sexo	Gender	0 = Female, 1 = Male
ocupacion	Occupation	1 = Health care worker, 2 = Office job, 3 = Outdoor work, 4 = Work in public area, 5 = Work at home, 6 = Unemployed
escolaridad	Schooling	1 = Analphabet, 2 = Primary school, 3 = Junior high school, 4 = High school, 5 = College, 6 = Postgraduate studies
nivsoc	Socioeconomic level	1 = Low, Medium-low, 2 = Medium-high, High
has	History of High blood pressure	1 = Yes, 2 = No
tabaquismo	Smoking	1 = Yes, 2 = No
dm	History of diabetes mellitus	1 = Yes, 2 = No
renal	History of renal failure	1 = Yes, 2 = No
autoinmunidad	History of autoimmune diseases	1 = Yes, 2 = No
edad	Age	Years old
peso	Weight	kg
talla	Height	m
temp	Temperature	Centigrades
fc	Heart rate	Beats per minute
fr	Respiratory rate	Cicles per minute
tas	Arterial pressure, systolic	mm Hg
tad	Arterial pressure, diastolic	mm Hg
score	Severity score	Score points (0-16)
ing_disnea	Short of breath at the moment of hospitalization	1 = Yes, 2 = No
sato2sin	Oxygen saturation	Percentage of Saturated hemoglobin
urea	Blood urea	mg/dL
bun	Blood urea nitrogen	mg/dL
creat	Serum creatinine	mg/dL
colesterol	Serum cholesterol	mg/dL
gluc	Blood glucose	mg/dL
hb	Hemoglobin concentration	mg/dL
leucos	White blood cells	number of cells/ $\hat{A}\mu\text{L}$
plaq	Platelets	number of cells/ $\hat{A}\mu\text{L}$
linfos	Lymphocytes	number of cells/ $\hat{A}\mu\text{L}$
monos	Monocytes	number of cells/ $\hat{A}\mu\text{L}$
eos	Eosinophils	number of cells/ $\hat{A}\mu\text{L}$
basof	Basophils	number of cells/ $\hat{A}\mu\text{L}$
neutros	Neutrophils	number of cells/ $\hat{A}\mu\text{L}$
k	Potassium	mEq/L

Table 3: Description of the variables in the working database. (*continued*)

Column name	Description	Codification/units
na	Sodium	mEq/L
cl	Chloride	mEq/L
ca	Calcium	mEq/L
ph	pH	Units of pH
pao2	Arterial oxygen partial pressure	mEq/L
paco2	Arterial carbon dioxide partial pressure	mEq/L
hco3	Arterial bicarbonate	mEq/L
dhl	Serum lactate dehydrogenase	IU/L
alat	Serum Alanine aminotransferase	IU/L
aat	Serum Aspartate aminotransferase	IU/L
dimd	D-dimer	µg/mL
fecha	Date of consultation	
fechacov1	Date of covid-19 testing	
fechahosp	Date of hospitalization	
fechalta	Date of hospital discharge	
fechainisint	Date of appearance of first symptoms	
motivoegre	Vital status at the moment of discharge	2 = Alive, 3 = Dead

## 2 Methods/analysis

### 2.1 Data wrangling

I read the data from a `csv` file and created an R object named `dbcovid`. The first thing I noticed is that date variables are of type `character`, thus I used the function `dmy` to transform the date variables into date type:

```
dbcovid <- dbcovid %>% mutate(across(starts_with('fecha'), dmy))
```

The laboratory data must be of `numeric` type, but several are as `character` in `dbcovid`, hence there must be typos. Table 4 shows these errors. To fix this problem, I used a `character` vector to put the names of the problematic variables, employed `regex` patterns and created a “catch-all” function, with the aid of the package `stringr`.

Table 4: Typos in numeric variables.

Typos in numeric variables	
urea	13,9
creat	0,81
bun	4..5
plaq1	160 000
plaq2	203 000
ca	A
aat	BA
alat	NA

```

#Check errors in numeric variables
nomcadenas <- c('urea', 'creat', 'bun', 'plaq', 'ca', 'aat', 'alat')
patron <- '([^\d\\.])|(\.{2,})' #Anything but digits or one decimal point.

fpatron <- function(x){
  x[str_which(x, pattern = patron)]
}

errores <- apply(dbcovid[, nomcadenas], 2, fpatron )

# Catch-all function to detect errors in numeric variables
farreglar <- function(x){
  x = str_trim(x)
  x = case_when(
    str_detect(x, '\\s') ~ str_replace_all(x, '\\s', ''),
    str_detect(x, '[:alpha:]') ~ str_replace_all(x, '[:alpha:]', ''),
    str_detect(x, '|\\.{2,}') ~ str_replace_all(x, '|\\.{2,}', '.'),
    TRUE ~ x
  )
}

#Fix errors in numeric variables
dbcovid <- dbcovid %>% mutate(across(all_of(nomcadenas), farreglar))
#Check if it worked
arregerrores <- apply(dbcovid[, nomcadenas], 2, fpatron )
arregerrores # no mistakes

# Transform numeric variables to numeric type
dbcovid <- dbcovid %>% mutate(across(all_of(nomcadenas), as.numeric))

```

Next, I used the dates to calculate the duration of some events, I also computed other variables such as the body mass index and the presence of obesity. As can be seen in table 3, the codification of dichotomic variables is not uniform; I changed them accordingly using the number 1 to indicate that the feature is present. Finally, I transformed all categorical variables to the **factor** class.

```

#Create new variables with dates, and eliminate the date variables.
dbcovid <- dbcovid %>%
  mutate(bmi = peso/talla^2,
    dayshosp = as.numeric(fechalta - fechahosp),
    duration = as.numeric(fechalta - fechainisint),
    daysdelay = as.numeric(fechahosp - fechainisint),
    obesity = ifelse(bmi >= 30, 1, 2)) %>%
  select(-starts_with('fecha'))

# Change codification of dichotomous categorical variables.
dicot <- c('motivoegre', 'has', 'tabaquismo', 'dm', 'renal', 'autoinmunidad',
  'ing_disnea', 'obesity')
dbcovid <- dbcovid %>%
  mutate(across(all_of(dicot), function(x) ifelse(x == 2, 0, 1)))

# Transform categorical variables to factors.
nomcateg <- c('motivoegre', 'sexo', 'ocupacion', 'nivsoc', dicot[-1])
dbcovid <- dbcovid %>%

```

```
mutate(across(all_of(nomcateg), as.factor))
```

## 2.2 Data exploration

From the 642 patients, 439 survived (68.4%), while 203 died (31.6%); this means that I am facing a relatively unbalanced data. Due to the number of variables, I decided to focus the data exploration on the difference in the outcome, in order to use these results as the first filter to decide which variables are going to be in the prediction models. Tables 5 and 6 show the exploratory analysis of categorical and numeric variables. While some variables look like good candidates, many are not worth including them. Furthermore, looking at the column *valid cases*, I notice that many variables do not add up to the total number of cases, thus there are a lot of missing values.

Table 5: Exploratory data analysis of categorical variables.

Variables	Category	Total	Dead <sup>a</sup>	Alive <sup>a</sup>	<i>p</i> <sup>b</sup>
Gender	Female	240	60 (25)	180 (75)	0.007
Gender	Male	402	143 (35.6)	259 (64.4)	
Occupation	Health care worker	97	16 (16.5)	81 (83.5)	< 0.001
Occupation	Office job	113	22 (19.5)	91 (80.5)	
Occupation	Outdoor work	40	20 (50)	20 (50)	
Occupation	Work in public area	149	49 (32.9)	100 (67.1)	
Occupation	Work at home	130	40 (30.8)	90 (69.2)	
Occupation	Unemployed	112	56 (50)	56 (50)	
Socioeconomic level	Low or medium-low	450	164 (36.4)	286 (63.6)	< 0.001
Socioeconomic level	Medium-high or high	192	39 (20.3)	153 (79.7)	
Autoimmunity	No	620	192 (31)	428 (69)	0.098
Autoimmunity	Yes	22	11 (50)	11 (50)	
Diabetes	No	510	146 (28.6)	364 (71.4)	0.002
Diabetes	Yes	132	57 (43.2)	75 (56.8)	
Hypertension	No	419	121 (28.9)	298 (71.1)	0.05
Hypertension	Yes	223	82 (36.8)	141 (63.2)	
Short of breath	No	49	8 (16.3)	41 (83.7)	0.025
Short of breath	Yes	593	195 (32.9)	398 (67.1)	
Obesity	No	383	120 (31.3)	263 (68.7)	1
Obesity	Yes	256	80 (31.2)	176 (68.8)	
Renal disease	No	602	184 (30.6)	418 (69.4)	0.064
Renal disease	Yes	39	18 (46.2)	21 (53.8)	
Smoking	No	596	180 (30.2)	416 (69.8)	0.006
Smoking	Yes	45	23 (51.1)	22 (48.9)	

<sup>a</sup> Values between parentheses represent the percentage for the corresponding category.

<sup>b</sup> *p* value calculated with  $\chi^2$ .

Table 6: Exploratory data analysis of numerical variables.

Variables	Valid cases	Status	Mean $\pm$ SD	<i>p</i> <sup>a</sup>
Blood urea nitrogen	437	Alive	13.78 $\pm$ 12.08	< 0.001
Blood urea nitrogen	196	Dead	23.11 $\pm$ 19.35	< 0.001
Serum creatinine	437	Alive	1.01 $\pm$ 1.1	< 0.001

Table 6: Exploratory data analysis of numerical variables. (*continued*)

Variables	Valid cases	Status	Mean $\pm$ SD	$p^a$
Serum creatinine	196	Dead	1.57 $\pm$ 1.87	< 0.001
Serum lactate dehydrogenase	356	Alive	371.92 $\pm$ 200.29	< 0.001
Serum lactate dehydrogenase	151	Dead	565.43 $\pm$ 250.71	< 0.001
Age	439	Alive	51.56 $\pm$ 13.84	< 0.001
Age	203	Dead	60.31 $\pm$ 13.17	< 0.001
Schooling	439	Alive	3.84 $\pm$ 1.45	< 0.001
Schooling	203	Dead	3.16 $\pm$ 1.5	< 0.001
Respiratory rate	438	Alive	24.47 $\pm$ 4.46	< 0.001
Respiratory rate	200	Dead	26.46 $\pm$ 5.42	< 0.001
Blood glucose	437	Alive	143.34 $\pm$ 83.98	< 0.001
Blood glucose	196	Dead	173.9 $\pm$ 101.57	< 0.001
Potassium	425	Alive	4 $\pm$ 0.55	< 0.001
Potassium	194	Dead	4.34 $\pm$ 0.84	< 0.001
Lymphocytes	437	Alive	1048.73 $\pm$ 704.26	< 0.001
Lymphocytes	196	Dead	829.61 $\pm$ 470.21	< 0.001
Neutrophils	437	Alive	7399.28 $\pm$ 4378.02	< 0.001
Neutrophils	196	Dead	10542.54 $\pm$ 6035.75	< 0.001
Arterial carbon dioxide partial pressure	396	Alive	27.25 $\pm$ 7.35	< 0.001
Arterial carbon dioxide partial pressure	193	Dead	32.26 $\pm$ 15.13	< 0.001
Arterial oxygen partial pressure	396	Alive	74.47 $\pm$ 29.18	< 0.001
Arterial oxygen partial pressure	193	Dead	58.2 $\pm$ 22.73	< 0.001
pH	396	Alive	7.43 $\pm$ 0.07	< 0.001
pH	193	Dead	7.36 $\pm$ 0.15	< 0.001
Oxygen saturation	298	Alive	85.72 $\pm$ 8.75	< 0.001
Oxygen saturation	116	Dead	73.62 $\pm$ 16.71	< 0.001
Severity score	405	Alive	1.22 $\pm$ 1.76	< 0.001
Severity score	193	Dead	2.61 $\pm$ 2.93	< 0.001
Blood urea	437	Alive	29.12 $\pm$ 24.96	< 0.001
Blood urea	196	Dead	49.42 $\pm$ 41.55	< 0.001
Duration	439	Alive	17.87 $\pm$ 8.13	0.001
Duration	203	Dead	15.25 $\pm$ 9.11	0.001
Days hospitalized	439	Alive	10.03 $\pm$ 6.67	0.002
Days hospitalized	203	Dead	8.04 $\pm$ 8.06	0.002
D-dimer	426	Alive	841.63 $\pm$ 2268.33	0.004
D-dimer	187	Dead	2297.73 $\pm$ 6738.31	0.004
Arterial pressure, diastolic	438	Alive	76.19 $\pm$ 10.81	0.007
Arterial pressure, diastolic	200	Dead	73.34 $\pm$ 12.81	0.007
Chloride	425	Alive	103.74 $\pm$ 5.46	0.009
Chloride	194	Dead	102.47 $\pm$ 5.67	0.009
Serum Aspartate aminotransferase	409	Alive	48.47 $\pm$ 37.88	0.01
Serum Aspartate aminotransferase	185	Dead	73.33 $\pm$ 127.89	0.01
White blood cells	437	Alive	10373.05 $\pm$ 14911.91	0.015
White blood cells	196	Dead	12464.77 $\pm$ 6664.19	0.015
Sodium	425	Alive	136.67 $\pm$ 4.21	0.02
Sodium	194	Dead	135.62 $\pm$ 5.54	0.02
Serum cholesterol	357	Alive	149.32 $\pm$ 53.65	0.025

Table 6: Exploratory data analysis of numerical variables. (*continued*)

Variables	Valid cases	Status	Mean $\pm$ SD	$p^a$
Serum cholesterol	201	Dead	139.68 $\pm$ 45.73	0.025
Heart rate	438	Alive	97.66 $\pm$ 17.47	0.03
Heart rate	200	Dead	101.21 $\pm$ 19.89	0.03
Basophils	436	Alive	47.21 $\pm$ 96.1	0.069
Basophils	196	Dead	65.49 $\pm$ 124.59	0.069
Hemoglobin concentration	437	Alive	14.47 $\pm$ 2.21	0.078
Hemoglobin concentration	196	Dead	14.14 $\pm$ 2.23	0.078
Monocytes	436	Alive	574.96 $\pm$ 604.01	0.149
Monocytes	196	Dead	645.18 $\pm$ 545.42	0.149
Days of delay	439	Alive	7.84 $\pm$ 5.76	0.171
Days of delay	203	Dead	7.21 $\pm$ 5.33	0.171
Arterial bicarbonate	396	Alive	18.16 $\pm$ 5.18	0.281
Arterial bicarbonate	193	Dead	17.68 $\pm$ 5.05	0.281
Arterial pressure, systolic	438	Alive	123.8 $\pm$ 19.33	0.316
Arterial pressure, systolic	199	Dead	125.86 $\pm$ 25.8	0.316
Body mass index	439	Alive	29.15 $\pm$ 5.11	0.338
Body mass index	200	Dead	29.57 $\pm$ 5.22	0.338
Height	439	Alive	1.63 $\pm$ 0.1	0.367
Height	200	Dead	1.63 $\pm$ 0.09	0.367
Eosinophils	436	Alive	38.36 $\pm$ 133.18	0.397
Eosinophils	196	Dead	49.19 $\pm$ 155.05	0.397
Serum Alanine aminotransferase	407	Alive	52.02 $\pm$ 40.36	0.521
Serum Alanine aminotransferase	179	Dead	56.06 $\pm$ 79.66	0.521
Calcium	33	Alive	7.25 $\pm$ 2.41	0.572
Calcium	12	Dead	7.51 $\pm$ 0.52	0.572
Weight	439	Alive	77.97 $\pm$ 15.47	0.647
Weight	200	Dead	78.6 $\pm$ 16.26	0.647
Platelets	437	Alive	273011.21 $\pm$ 111230.3	0.686
Platelets	196	Dead	268847.55 $\pm$ 123457.44	0.686
Temperature	438	Alive	37.06 $\pm$ 0.82	0.761
Temperature	200	Dead	37.09 $\pm$ 0.93	0.761

<sup>a</sup>  $p$  value calculated with Student t test, except for *Schooling*, where the U Mann-Whitney test was used.

## 2.3 Data preprocessing

First, I excluded the variables that were not significant in the bivariate analysis, choosing (arbitrarily) a cut-off value of  $p < 0.051$ . I also excluded the length of hospital stay and the duration of the disease, as I am interested in the variables recorded the day of hospitalization.

```
finalvarcateg <- varcateg[which(numlistachis < 0.051)]
finalvarnum <- varnum[which(listapesnum < 0.051)]
finalvarnum <- finalvarnum[1:24]
dbcovid <- dbcovid %>%
  select(all_of(c('motivoegre', finalvarcateg, finalvarnum))) %>%
  mutate(motivoegre = factor(motivoegre, labels = c('Alive', 'Dead'))),
```



```

sexo = factor(sexo, labels = c('Female', 'Male')),
nivsoc = factor(nivsoc, labels = c('low', 'High')),
dm = factor(dm, labels = c('No', 'Yes')),
has = factor(has, labels = c('No', 'Yes')),
tabaquismo = factor(tabaquismo, labels = c('No', 'Yes')),
ing_disnea = factor(ing_disnea, labels = c('No', 'Yes'))

```

As part of preprocessing, we also want to know if there are variables close to zero variation. The function `nearZeroVar` allows us to explore this situation. In this case, no variable was found to have this problem.

```
nearZeroVar(dbcovid)
```

```
integer(0)
```

In the next step, I checked the number of missing values in the selected variables (table 7). The simplest form to deal with missing data is to restrict the analysis to individuals with complete data. In fact, this is the default approach in most of the prediction algorithms allowing missingness in `caret`. Considering the selection of variables that are significant, if I were to use *complete case analysis*, I would loose around 62% of the cases, causing a reduction in efficiency. The other form to deal with this problem is using *imputation methods*.

Dealing with missing data is always complicated. I did some research, and found that first we have to think on the reasons of missingness (Madley-Dowd et al. 2019). The missingness mechanism can be classified as:

- **Missing completely at random:** the probability of missingness is independent on observed or missing data.
- **Missing at random:** missingness independent of unobserved data, conditional on the observed data.
- **Missing not at random:** missingness dependent on unobserved data even after conditioning on observed data.

In this study, considering how the data were acquired, it is reasonable to think that missingness is at random, hence, unbiased results can be obtained with large proportions of missing data. In the paper of Madley-Dowd, they simulated data and concluded that when missingness is at random, unbiased results can be obtained even when 90% of the data are missing (Madley-Dowd et al. 2019). Accordingly, the variables in table 7 were used to build the prediction models.

The next consideration is the imputation method. Imputation methods can broadly be classified as *single* and *multiple*. In a single imputation procedure the value of a missing data element is filled by some means, such as the mean or the median of the observed values, without defining an explicit model. In multiple imputation, several draws of the missing elements are made from the posterior distribution of the missing data, given the observed data (Glas 2010).

The package `caret` has some functions that can be used for preprocessing the data. the `preProcess` function has several imputation methods. In this work, I used two methods of imputation: `medianImpute`, an example of single imputation which takes the median of each predictor in the training set, and `bagImpute`, an example of multiple imputation which fits a bagged tree model for each predictor<sup>1</sup>. These methods create an object that has to be used in the `predict` function to finally have the imputed data sets (see code).

```

#Impute with the median
preprocmd <- preProcess(dbcovid, method = 'medianImpute')
dbcovidmd <- predict(preprocmd, dbcovid)
sum(is.na(dbcovidmd)) # Show the missings in the categorical variables

```

<sup>1</sup>The methods used do not impute categorical variables. In this particular case, there were only two cases with a missing element in these type of variables, thus I decided to omit them.

```
## [1] 2
```

```
dbcovidmd <- na.omit(dbcovidmd)
sum(is.na(dbcovidmd)) #check we have no missing data
```

```
## [1] 0
```

```
#Imputation via bagging (if I leave the factor variables, it throws an error).
set.seed(271220)
preprocbag <- dbcovid %>%
  select(where(is.numeric)) %>%
  preProcess(., method = 'bagImpute')

dbcovidbag <- dbcovid %>%
  select(where(is.numeric)) %>%
  predict(preprocbag, .)

dbcovidbag <- dbcovid %>%
  select(where(is.factor)) %>%
  bind_cols(dbcovidbag) %>%
  na.omit
sum(is.na(dbcovidbag)) # check we have no missing data
```

```
## [1] 0
```

## 2.4 Selection of the classification algorithms and evaluation metrics.

The `caret` package embodies more than 200 classification and regression models. In this work, I used the following algorithms:

- Logistic regression, with the `glm` method.
- k-Nearest neighbors (kNN), with the `knn` method.
- Classification and regression trees (CART), with the `rpart` method.
- Random forests, with the `rf` method.

When the output is continuous, the main goal of the machine learning algorithm is referred as *prediction*; when the output is categorical, then we use the word *classification* (Irizarry 2019a). Hence, as the outcome of the dataset is categorical, I will use this term going forward.

In general, machine learning algorithms estimate conditional probabilities as function of the predictors. For classification models, we can represent the main task with the following expression:

$$p_k(\mathbf{x}) = \Pr(Y = k \mid \mathbf{X} = \mathbf{x}), \text{ for } k = 1, \dots, K \quad (1)$$

where  $\mathbf{X}$  represents the set of predictors, and  $K$  the total number of classes of the categorical variable.

Logistic regression is a specific case of a set of *generalized linear models*, it can be regarded as the simplest prediction model when the outcome is categorical. The *method kNN* is similar to bin smoothing, but easier to adapt to multiple dimensions; the general idea is to define a distance between all observations based on the values of the predictors, then we look for the  $k$  nearest points and calculate the proportion of the outcome associated with these points (Irizarry 2019b). CART models predict an outcome variable  $Y$  by partitioning

Table 7: Missing values in the putative predictive variables.

Variable	Valid	Total missing	Percentage missing	Total
sato2sin	414	228	35.5	642
dhl	507	135	21.0	642
colesterol	558	84	13.1	642
ph	589	53	8.3	642
pao2	589	53	8.3	642
paco2	589	53	8.3	642
aat	594	48	7.5	642
score	598	44	6.9	642
dimd	613	29	4.5	642
k	619	23	3.6	642
na	619	23	3.6	642
cl	619	23	3.6	642
urea	633	9	1.4	642
bun	633	9	1.4	642
creat	633	9	1.4	642
gluc	633	9	1.4	642
leucos	633	9	1.4	642
linfos	633	9	1.4	642
neutros	633	9	1.4	642
fc	638	4	0.6	642
fr	638	4	0.6	642
tad	638	4	0.6	642
ocupacion	641	1	0.2	642
tabaquismo	641	1	0.2	642
motivoegre	642	0	0.0	642
sexo	642	0	0.0	642
nivsoc	642	0	0.0	642
dm	642	0	0.0	642
has	642	0	0.0	642
ing_disnea	642	0	0.0	642
escolaridad	642	0	0.0	642
edad	642	0	0.0	642

the predictors, defining decision rules and creating *decision trees* with prediction at the ends, referred to as *nodes* (Irizarry 2019c). *Random forests* extend the concept of classification trees. The goal is to improve prediction performance with the construction of a number of individual trees randomly different through bootstrapping (Irizarry 2019c).

## 2.5 Strategy of model generation

When we build classification or regression models, because the data are random, the parameters and results obtained are also random variables. If we repeat the experiment the results are going to be different. According to the law of big numbers, after a lot of repetitions, results will tend to converge. Thus, one of the most important concepts in machine learning is the use of a resampling technique to stabilize the results. The `caret` package has a number of functions that help us to develop a good strategy for model building, including resampling. Once we preprocessed the data, the next step is to create a training set and a test set using the function `createDataPartition`. The text book of the course recommends the training set to be the biggest. Because the final data has 640 cases, to keep a decent amount of cases in the test set, I decided to partition the data into 85% of the cases for the training set and the remaining for the test set. Next, I chose the resampling method embedded in the function `train`. In this work I decided to use *bootstrap*; it is computationally more expensive than *k-fold cross validation*, but the size of the data set makes it feasible.

Before training any algorithm, we must decide on which evaluation metric to use. During the training step, the `train` function uses the *accuracy* as the default when the outcome is binary. *Overall accuracy* is defined as the proportion of cases that were correctly predicted in the test set (Irizarry 2019b). The function `confusionMatrix` calculates all the metrics for binary outcomes, and I will use it to calculate the overall accuracy.

Many algorithms have parameters that have to be calibrated by means of resampling. For the algorithms selected in this work we have the following parameters:

- With `knn` we need to find the best value for the number of nearest neighbors with the argument `k`.
- With `rpart` we have to tune up the *complexity parameter* with the argument `cp`, this parameter sets a minimum for how much the evaluation metric must improve for another partition to be added, because if we split until every point is its own partition, we would be overtraining the model. Although not used as tuning parameters, we can also change, using customized functions, the minimum number of observations required in a partition before deciding the next partition (`minsplit`), as well as the number of observations in each node (`minbucket`) .
- With `rf` we tune up the number of variables randomly sampled as candidates at each split with `mtry`, but it is also possible to optimize other parameters such as the minimum node size with customized functions.

In summary, for the purpose of this work, I trained the models with the bootstrap method, varying only the tuning parameters and keeping the rest with the default values. The evaluation metric was accuracy.

## 3 Results

The four models tested were trained with both imputation methods [<sup>2</sup>]. The summary of the accuracy obtained with each method is presented in table 8, together with the 95% confidence interval. While the imputation method had little influence in the accuracy of the logistic regression and de kNN models, the CART and random forests algorithms tended to perform better with the median imputation, although the confidence intervals overlap. The algorithm that renders the best accuracy is the random forests with median imputation (0.8210526).

[<sup>2</sup>] For the sake of clarity in the manuscript, I decided to omit the code. Please [us this link](#) to open the corresponding R Script file.

Table 8: Accuracy of the models tested.

Model	Imputed with the median		Imputed with bagging	
	Accuracy	CI <sub>95%</sub>	Accuracy	CI <sub>95%</sub>
Generalized linear model	0.779	0.682-0.858	0.789	0.694-0.866
k nearest neighbors	0.642	0.537-0.738	0.653	0.548-0.747
CART	0.705	0.603-0.794	0.621	0.516-0.719
Random forests	0.821	0.729-0.892	0.789	0.694-0.866

Table 9: Results of the logistic regression model (bagging imputation).

Variable	Estimate	Std. error	z value	p
sato2sin	-0.0488	0.0136	-3.588	0.0003
neutros	0.0001	0.0000	3.276	0.0011
pao2	-0.0191	0.0059	-3.221	0.0013
paco2	0.0432	0.0165	2.621	0.0088
dhl	0.0018	0.0008	2.146	0.0319
linfos	-0.0006	0.0003	-2.104	0.0353
edad	0.0255	0.0134	1.899	0.0575
tabaquismoYes	0.8471	0.4886	1.734	0.0830
ph	-2.7400	1.6350	-1.675	0.0939
colesterol	-0.0049	0.0030	-1.654	0.0982
nivsocHigh	-0.6515	0.3957	-1.646	0.0997
escolaridad	-0.1669	0.1058	-1.578	0.1146
ocupacion5	-0.8407	0.5442	-1.545	0.1224
score	0.1123	0.0802	1.401	0.1613
cl	-0.0411	0.0310	-1.324	0.1854
ing_disneaYes	0.8315	0.6382	1.303	0.1926
urea	0.0340	0.0297	1.143	0.2531
gluc	-0.0018	0.0018	-0.991	0.3219
leucos	0.0000	0.0000	-0.951	0.3415
tad	-0.0115	0.0122	-0.943	0.3457

Table 9 shows the coefficients of the logistic regression model (bagging imputation), ordered by the  $p$  value. Unfortunately, the `glm` function does not provide the standardized coefficients, thus their absolute values are not comparable; however, we notice that some values are negative, meaning that the values are inversely related with the outcome. We also see that only 6 variables out of the 31 predictors analyzed had a  $p < 0.05$ .

The kNN model had the lowest accuracy of all the models tested (0.6526316 using `bagImpute`). The best number of neighbors was 30 (figure 1).

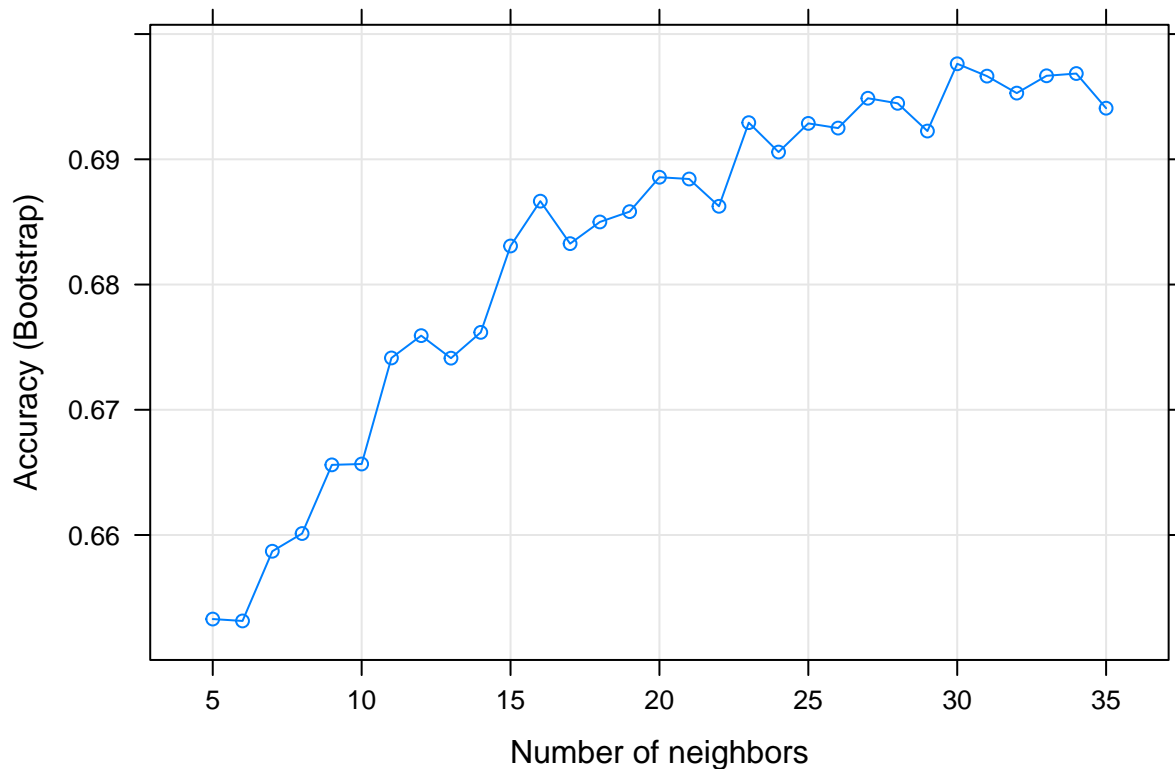


Figure 1: kNN model. Effect of the number of neighbors on accuracy.

The CART model, when imputation of the data with the median was used, had an accuracy of 0.7052632. The best value of the complexity parameter was 0.026 (figure 2). This algorithm allows to visualize a decision tree, hence is highly interpretable (figure 3). The tree begins with the concentration of d-dimers, closely related with clot formation inside the vessels (intravascular coagulation). We have to follow the left side of the branch when the condition is met. This model tells us that a patient with d-dimer above 321.5  $\mu\text{g/mL}$ , with white blood cells above 8899 cell/mL, arterial partial oxygen pressure below 71.5 mm Hg and heart rate above 109, will probably die. On the contrary, a patient with a d-dimer concentration below 321.5  $\mu\text{g/mL}$ , has more chances of survival. Albeit very useful, the accuracy is modest, and due to the recursive partitioning, it is very easy to overtrain (Irizarry 2019c).

Finally, we analyze the performance of the random forests algorithm. It addresses the shortcomings of the CART models, but we lose interpretability. As expected, this machine learning technique gave the best results concerning accuracy. When the data imputed with the median was used, random forests had an accuracy of 0.8210526 8. The best number of variables sampled as candidates at each split was 5 (figure 4, left panel). Although I did not tune up the number of trees, the right panel of figure 4 shows that the number used (150) was enough to stabilize the error.

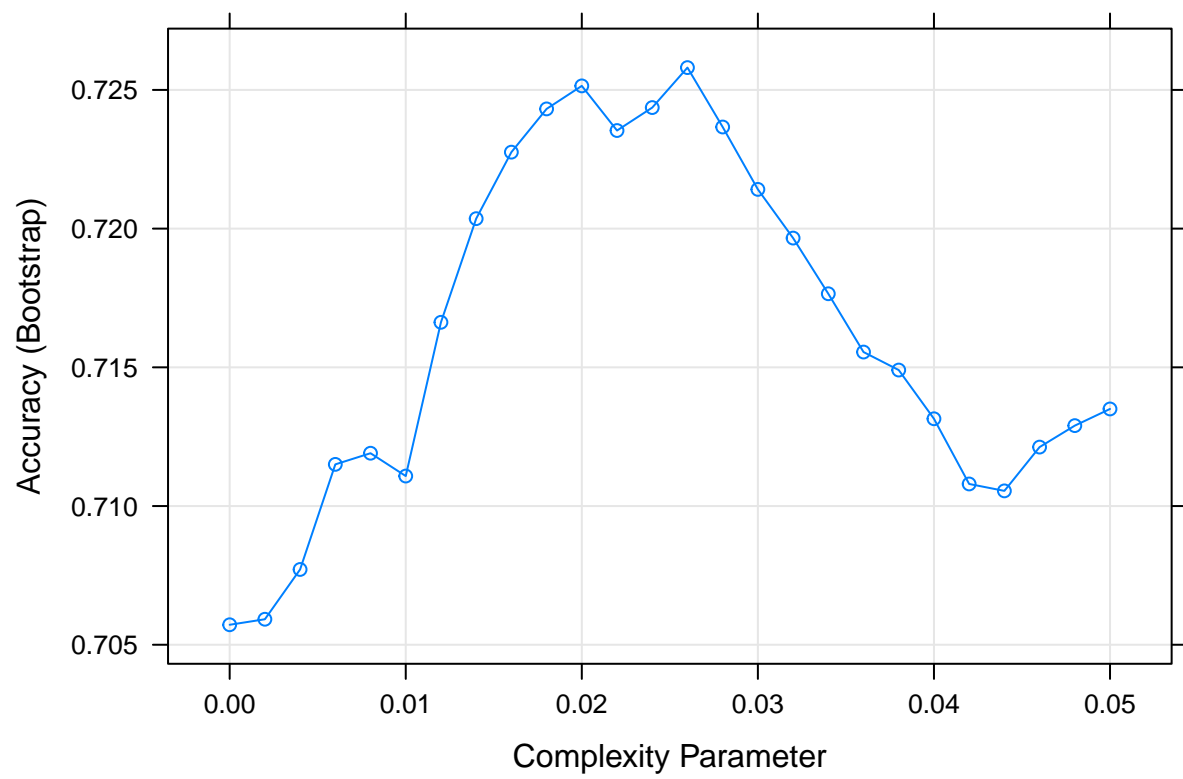


Figure 2: Classification tree model. Effect of different values of the complexity parameter on model accuracy.

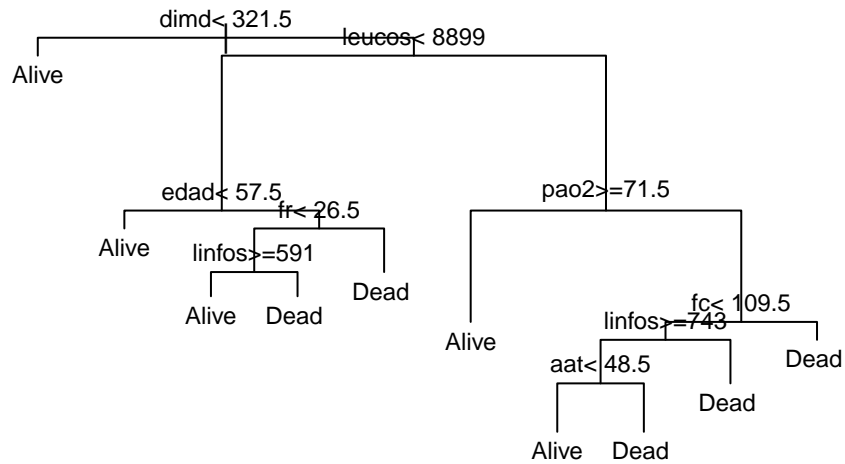


Figure 3: Classification tree model. Final decision tree.

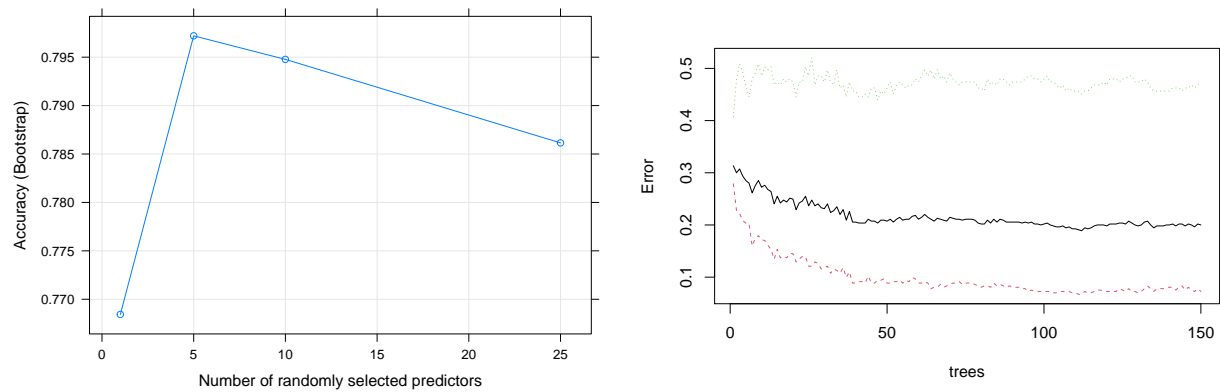


Figure 4: Random forest model. *Left panel*: effect of the number of randomly selected predictors on accuracy. *Right panel*: effect of the number of trees on the stabilization of the error.



Table 10 describes the overall variable importance calculated for the models with the best imputation method using the `varImp` function of the `caret` package. Only the 20 most important variables are shown for each model. It is worth noting the absence of d-dimer (`dimd`) in the logistic regression model, but with the other variables, although there is no agreement with respect to the order, the algorithms agree in general with the 5 most important variables. Interestingly, neither diabetes mellitus nor high blood pressure were considered risk factors by these models.

Table 10: Variable importance of the models with the best imputation performance.

Generalized linear model		k nearest neighbors		CART		Random forests	
Predictor	Overall importance	Predictor	Overall importance	Predictor	Overall importance	Predictor	Overall importance
sato2sin	100.00	sato2sin	100.00	dhl	100.00	dhl	100.00
neutros	91.21	dhl	98.91	dimd	69.39	dimd	99.56
pao2	89.65	dimd	82.79	sato2sin	65.50	pao2	81.80
paco2	72.71	pao2	69.59	neutros	65.39	neutros	64.94
dhl	59.30	urea	63.44	linfos	61.33	bun	63.50
linfos	58.13	bun	61.48	pao2	50.94	leucos	51.06
edad	52.34	edad	60.12	urea	36.27	sato2sin	51.03
tabaquismoYes	47.66	neutros	58.59	leucos	35.92	paco2	50.24
ph	46.01	score	55.96	ph	32.22	ph	49.70
colesterol	45.40	leucos	54.35	k	28.46	edad	47.55
nivsocHigh	45.19	creat	52.88	score	24.55	creat	45.05
escolaridad	43.26	escolaridad	46.35	fc	13.22	urea	43.05
ocupacion5	42.33	ph	45.92	edad	11.51	gluc	39.88
score	38.26	fr	42.00	aat	9.48	linfos	38.99
cl	36.10	gluc	36.01	fr	8.26	fc	37.83
ing_disneaYes	35.50	k	35.95	nivsocHigh	6.83	fr	37.69
urea	30.98	ocupacion	33.83	colesterol	5.72	score	34.77
gluc	26.68	aat	29.77	bun	4.60	colesterol	33.27
leucos	25.56	linfos	29.07	creat	4.47	aat	30.75
tad	25.33	nivsoc	24.62	sexoMale	0.00	tad	29.20

## 4 Conclusion

The random forests model with data imputed with the median gave the highest accuracy. Lactic dehydrogenase, d-dimer, partial arterial oxygen pressure, number of neutrophils (a type of white blood cell), oxygen saturation breathing room air, age, number of white blood cells, arterial pH, serum urea and partial arterial carbon dioxide were considered the 10 most important variables by this supervised machine learning model.

Machine learning algorithms are powerful tools and a promising area of research in the clinical setting. They will help in better decision making with complex diseases such as covid-19, but as stated in the TRIPOD requirements, the study design is also important. In this particular case, the data was gathered from the everyday hospital records, hence the quality during data acquisition is compromised; the typos and data missingness are examples of this problem. Another limitation could be that the sample size was not calculated *a priori*, I only used the amount of patients hospitalized in a certain time frame. The estimation of the sample size for classification algorithms in an open question (Brownlee 2017). However, given the number of predictors and the type of classification models used, we do not expect that the sample size used could bias the results obtained (Vabalas et al. 2019).

Finally, the strategy used in this work did not include a validation data set. As future work, we plan to use the random forests prediction algorithm, and test it prospectively, with a completely different data set and check its performance.

## References

- Brownlee, Jason. 2017. “How Much Training Data Is Required for Machine Learning?” <https://machinelearningmastery.com/much-training-data-required-machine-learning/>.
- Collins, Gary S., Johannes B. Reitsma, Douglas G. Altman, and Karel G. M. Moons. 2015. “Transparent Reporting of a Multivariable Prediction Model for Individual Prognosis or Diagnosis (TRIPOD): The TRIPOD Statement.” *Journal of Clinical Epidemiology* 68 (2): 112–21. <https://doi.org/10.1016/j.jclinepi.2014.11.010>.
- Glas, C. A. W. 2010. “Missing Data.” In, edited by Penelope Peterson, Eva Baker, and Barry McGaw, 283–88. Oxford: Elsevier. <https://doi.org/10.1016/B978-0-08-044894-7.01346-4>.
- Irizarry, Rafael A. 2019a. *Chapter 27 Introduction to Machine Learning / Introduction to Data Science*. <https://rafalab.github.io/dsbook/introduction-to-machine-learning.html#conditional-probabilities-and-expectations>.
- . 2019b. *Chapter 29 Cross Validation / Introduction to Data Science*. <https://rafalab.github.io/dsbook/cross-validation.html>.
- . 2019c. *Chapter 31 Examples of Algorithms / Introduction to Data Science*. <https://rafalab.github.io/dsbook/examples-of-algorithms.html#classification-and-regression-trees-cart>.
- Madley-Dowd, Paul, Rachael Hughes, Kate Tilling, and Jon Heron. 2019. “The Proportion of Missing Data Should Not Be Used to Guide Decisions on Multiple Imputation.” *Journal of Clinical Epidemiology* 110 (June): 63–73. <https://doi.org/10.1016/j.jclinepi.2019.02.016>.
- Sanchez Talanquer, Mariano. 2020. “La letalidad hospitalaria por covid-19 en México: desigualdades institucionales.” <https://datos.nexos.com.mx/?p=1625>.
- Secretaria de Salud. 2020. “Informe técnico Diario Covid-19 México,” December. [https://www.gob.mx/cms/uploads/attachment/file/601308/Comunicado\\_Tecnico\\_Diario\\_COVID-19\\_2020.12.17.pdf](https://www.gob.mx/cms/uploads/attachment/file/601308/Comunicado_Tecnico_Diario_COVID-19_2020.12.17.pdf).
- Vabalas, Andrius, Emma Gowen, Ellen Poliakoff, and Alexander J. Casson. 2019. “Machine Learning Algorithm Validation with a Limited Sample Size.” *PLOS ONE* 14 (11): e0224365. <https://doi.org/10.1371/journal.pone.0224365>.