

# 06 - Classification

ml4econ, HUJI 2020

Itamar Caspi

April 26, 2020 (updated: 2020-04-26)

# Packages and setup

Use the `{pacman}` package that automatically loads and installs packages if necessary:

```
if (!require("pacman")) install.packages("pacman")

pacman::p_load(
  tidyverse,    # for data wrangling and visualization
  tidymodels,
  knitr,        # for displaying nice tables
  here,         # for referencing folders and files
  glmnet,       # for estimating lasso and ridge
  ggmosaic,     # for tidy mosaic plots
)
```

Set a theme for `ggplot` (Relevant only for the presentation)

```
theme_set(theme_grey(20))
```

And set a seed for replication

```
set.seed(1203)
```

# Outline

- Binary Classification Problems
- The Confusion Matrix
- The Logistic Regression Model
- Sensitivity Specificity Trade-off
- Multiclass classification (next time)

# Binary Classification Problems

# Bill Gates on Testing for COVID-19

"Basically, there are two critical cases: anyone who is symptomatic, and anyone who has been in contact with someone who tested positive. Ideally both groups would be sent a test they can do at home without going into a medical center. Tests would still be available in medical centers, but the simplest is to have the majority done at home. **To make this work, a government would have to have a website that you go to and enter your circumstances, including your symptoms. You would get a priority ranking, and all of the test providers would be required to make sure they are providing quick results to the highest priority levels.** Depending on how accurately symptoms predict infections, how many people test positive, and how many contacts a person typically has, you can figure out how much capacity is needed to handle these critical cases. For now, most countries will use all of their testing capacity for these cases." - Bill Gates.

Source: "The first modern pandemic by Bill Gates"

# Binary classification

Let  $y_i$  denote the outcome of a COVID-19 test, where

$$y_i = \begin{cases} 1 & \text{if positive,} \\ 0 & \text{if negative,} \end{cases}$$

where the values 1 and 0 are chosen for simplicity.<sup>1</sup>

Two types of questions we might ask:

1. What is the probability of being positive?
2. Can we classify an individual as positive/negative?

[\*] It is common to find a  $\{1, -1\}$  notation for binary outcomes in the ML literature.

# Israeli COVID-19 tests data

The [The Isreali Ministry of Health](#) provides information on more than 100,000 COVID-19 test results. Our aim here is to predict which person will be classified as "positive", i.e. infected by the virus, based on his symptoms and characteristics.

Outcome variable: `corona_result`

Features:

- Symptoms
  - cough
  - fever
  - sore\_throat
  - shortness\_of\_breath
  - head\_ache
- Characteristics
  - age\_60\_and\_above
  - gender

# Read and examine the data

```
covid_raw <- here("06-classification/data", "covid_proc.csv") %>%  
  read_csv()
```

```
covid_raw %>% glimpse()
```

```
## Observations: 107,542  
## Variables: 8  
## $ cough          <dbl> 1, 1, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 0...  
## $ fever          <dbl> 0, 1, 1, 1, 0, 1, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, 1...  
## $ sore_throat    <dbl> 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0...  
## $ shortness_of_breath <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0...  
## $ head_ache      <dbl> 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0...  
## $ corona_result  <chr> "negative", "negative", "negative", "positive", "neg...  
## $ age_60_and_above <chr> "No", "No", "Yes", "Yes", "Yes", "Yes", "No", "No", ...  
## $ gender         <chr> "male", "male", "male", "male", "female", "male", "m...
```

Note that since  $n = 107,542$  and  $p = 7$ , we should not worry much about overfitting.



# Preprocessing

We'll now define all variables, outcome and features, as factors:

```
covid <- covid_raw %>%  
  mutate_all(as_factor)
```

and extract the outcome and features as matrices (for later use with `glmnet`):

```
x <- covid %>%  
  select(-corona_result) %>%  
  model.matrix(~ .-1, data = .)  
  
y <- covid %>% pull(corona_result) %>% as_factor()
```

# Raw detection frequencies

How are test results distributed?

```
covid %>%  
  group_by(corona_result) %>%  
  count()
```

```
## # A tibble: 2 x 2  
## # Groups:   corona_result [2]  
##   corona_result     n  
##   <fct>         <int>  
## 1 negative      98586  
## 2 positive       8956
```

This is an example of **class imbalance** (the distribution of examples across the known classes is skewed), which is a typical feature of classification problems.

# Measuring classification accuracy

What does MSE mean in the context of classification problems?

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{y_i \neq \hat{y}_i\}}$$

In words: In this case, MSE measures the **missclassification rate**, i.e., the ratio between the number of missclassifications and the total number of observations.

**Classification accuracy** is the total number of correct predictions divided by the total number of predictions made for a dataset.

Clearly,

$$accuracy = 1 - missclassification.$$

Are missclassification/accuracy rates useful? Think imbalanced outcome.

# A naive classifier

Our naive "model" says: "classify everyone as being negative"

```
covid %>%  
  mutate(corona_result = as_factor(corona_result)) %>%  
  mutate(.fitted_class = factor("negative", levels = c("negative", "positive"))) %>%  
  conf_mat(corona_result, .fitted_class)
```

```
##           Truth  
## Prediction negative positive  
##   negative    98586     8956  
##   positive         0         0
```

The accuracy of the model is  $98,586/107,542 = 91.67\%$ !

Pretty impressive! Or is it?

This naive classifier lacks the ability to discern one class versus the other, and more importantly, it fails to identify infected individuals - the thing we really care about!

# The Confusion Matrix

# Beyond accuracy – other measures of performance

The **confusion matrix** is a table that categorizes predictions according to whether they match the ground truth.

		Truth	Truth
		Negative	Positive
Prediction	Negative	<i>True negative (TN)</i>	<i>False negative (FN)</i>
Prediction	Positive	<i>False positive (FP)</i>	<i>True positive (TP)</i>

Note that  $TN + FP + TP = N$ , where  $N$  is the number of observations. Accuracy in this case is defined as  $(TN + TP)/N$ .

**Note:** The confusion matrix can be extended to multiclass outcomes.

# Types of classification errors

**False positive rate:** The fraction of negative examples that are classified as positive,  $0/98,586 = 0\%$  in example.

**False negative rate:** The fraction of positive examples that are classified as negative,  $8,956/8,956 = 100\%$  in example.

Can we do better?

# A perfect classifier

Here is a simple example. Let's assume we have a sample of 100 test results, and exactly 20 of them are labeled "positive". If our classifier was perfect, the confusion matrix would look like this:

		Truth	Truth
		Negative	Positive
Prediction	Negative	80	0
Prediction	Positive	0	20

That is, our classifier has a 100% accuracy rate, zero false positive and zero false negative.



# The realistic classifier

Now, here is a classifier that makes some errors:

		Truth	Truth
		Negative	Positive
Prediction	Negative	70	10
Prediction	Positive	5	15

In this example, 10 persons with the pathogen were classified as Negative (not infected), and 5 persons without the pathogen were classified as Positive (infected).

# Logistic Regression Model

# First things first: the linear probability model

Consider a dependent variable  $y_i \in \{0, 1\}$ . Given a vector of features  $\mathbf{x}_i$ , the goal is to predict  $\Pr(y_i = 1|\mathbf{x}_i)$ .

Let  $p_i$  denote the probability of seeing  $y_i = 1$  given  $\mathbf{x}_i$ , i.e.,

$$p_i \equiv \Pr(y_i = 1|\mathbf{x}_i)$$

The linear probability model specifies that

$$p_i = \mathbf{x}_i' \boldsymbol{\beta}$$

However, an OLS regression of  $y_i$  on  $\mathbf{x}_i$  ignores the discreteness of the dependent variable and does not constrain predicted probabilities to be between zero and one.

# Logistic regression model

A more appropriate model is the **logit model** or **logistic regression model** specifies as

$$p_i = \Lambda(\mathbf{x}'_i \boldsymbol{\beta}) = \frac{\exp(\mathbf{x}'_i \boldsymbol{\beta})}{1 + \exp(\mathbf{x}'_i \boldsymbol{\beta})}$$

where  $\Lambda(\cdot)$  is the logistic cdf. As such, the model imposes the restriction that  $0 \leq p_i \leq 1$ .

# Odds-ratio

Note that

$$\frac{p_i}{1 - p_i} = \exp(\mathbf{x}'_i \boldsymbol{\beta})$$

Taking logs yields

$$\ln\left(\frac{p_i}{1 - p_i}\right) = \mathbf{x}'_i \boldsymbol{\beta}$$

The above is useful representation of the logistic regression model. The LHS is called the log **odds ratio** (or relative risk.)

Hence, we can say that the logistic regression model is linear in log odds-ratio.

# The likelihood function

**Likelihood** refers to the probability of seeing the data given parameters.

$$\begin{aligned}\text{Likelihood} &= \prod_{i=1} \Pr(y_i | \mathbf{x}_i) \\ &= \prod_{i=1} p_i^{y_i} (1 - p_i)^{1-y_i} \\ &= \prod_{i=1}^n \left( \frac{\exp(\mathbf{x}_i' \beta)}{1 + \exp(\mathbf{x}_i' \beta)} \right)^{y_i} \left( \frac{1}{1 + \exp(\mathbf{x}_i' \beta)} \right)^{1-y_i}\end{aligned}$$

taking (natural) logs yields the **log likelihood**

$$\log(\text{Likelihood}) = \sum_{i=1}^N \left[ \log \left( 1 + e^{(\beta_0 + x_i' \beta)} \right) - y_i \cdot (\beta_0 + x_i' \beta) \right]$$

In estimation, we want to make the above as big as possible (hence, maximum likelihood estimation, MLE).

# Deviance

Another useful concept is the **deviance**, a generalization of the concept of "least squares" to general linear models (such as logit), and is a measure of the distance between data and fit.

The relationship between deviance and likelihood is given by

$$\text{Deviance} = -2 \times \log(\text{Likelihood}) + \text{Constant}$$

The constant wraps terms that relate to the likelihood of the "perfect" model and we can mostly ignore it.

# Deviance and estimation

In estimation, we want to make deviance as *small* as possible.

$$\begin{aligned}\text{Deviance} &= -2 \sum_{i=1}^N \left[ \log \left( 1 + e^{(\beta_0 + x'_i \beta)} \right) - y_i \cdot (\beta_0 + x'_i \beta) \right] + \text{Constant} \\ &\propto \sum_{i=1}^N \left[ \log \left( 1 + e^{(\beta_0 + x'_i \beta)} \right) - y_i \cdot (\beta_0 + x'_i \beta) \right]\end{aligned}$$

This is the what R's `glm` function minimizes for logistic regressions.

(**NOTE:** In linear models, the deviance is porportional to the RSS)



# Penalized logistic regression

We can also minimize the deviance subject to a standard lasso type ( $\ell_1$  norm) penalty on  $\beta$ :

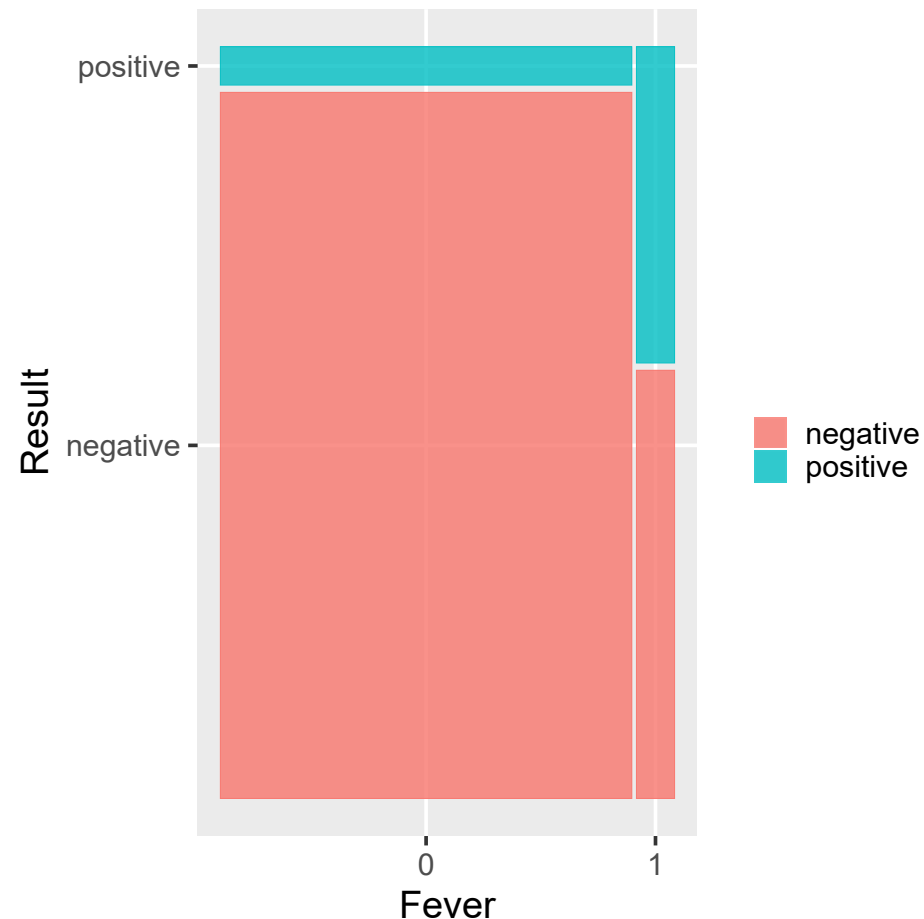
$$\min_{(\beta_0, \beta) \in \mathbb{R}^{p+1}} \left[ \frac{1}{N} \sum_{i=1}^N \log \left( 1 + e^{(\beta_0 + x'_i \beta)} \right) - y_i \cdot (\beta_0 + x'_i \beta) \right] + \lambda \|\beta\|_1$$

where again, the penalty is on the sum of the absolute values of  $\beta$  (no including the intercept.)

# Back to the data: can we do better than being "naive"?

There is some evidence that having fever is associated with being "positive".

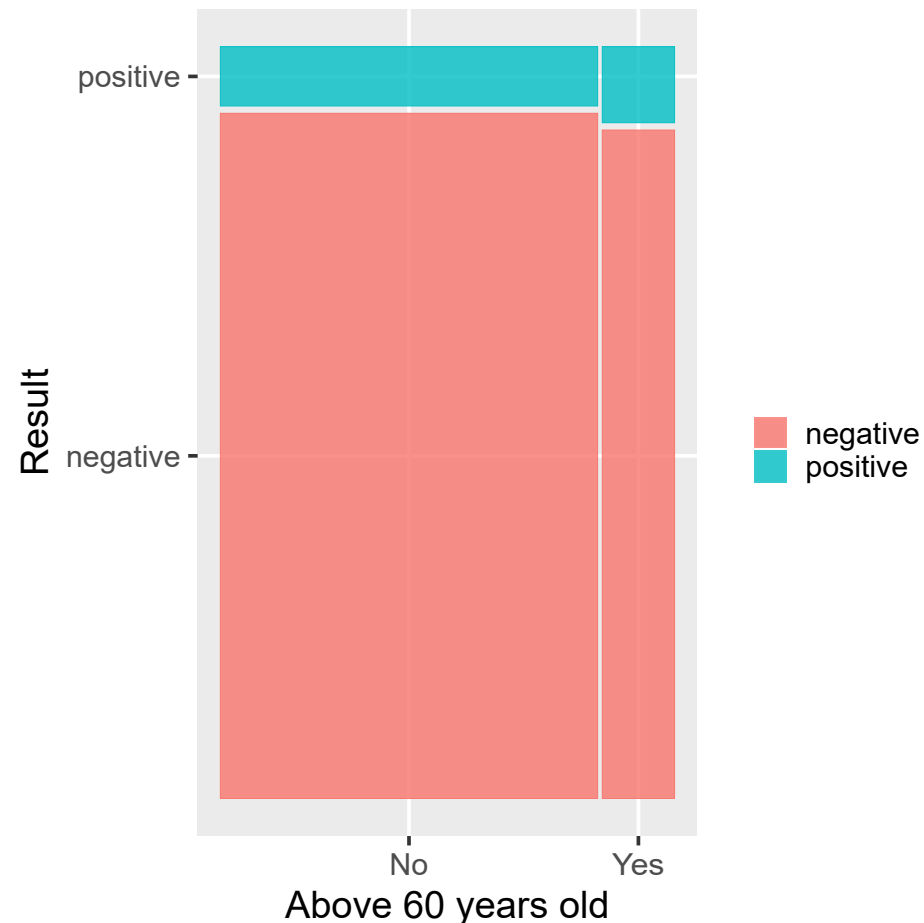
```
covid %>%  
  ggplot() +  
  geom_mosaic(  
    aes(x = product(corona_result, fever),  
        fill = corona_result)  
  ) +  
  labs(  
    x = "Fever",  
    y = "Result",  
    fill = ""  
  )
```



# Back to the data: can we do better than being "naive"?

and some evidence for an association with age (above 60)

```
covid %>%  
  ggplot() +  
  geom_mosaic(  
    aes(x = product(corona_result, age_60_and_above)  
        fill = corona_result)  
  ) +  
  labs(  
    x = "Above 60 years old",  
    y = "Result",  
    fill = ""  
  )
```



# Estimating the model using R

We will estimate the model using base R's `glm` (stands for generalized linear model) function:

```
logit_model <- glm(  
  corona_result ~ .,  
  data = covid,  
  family = "binomial"  
)
```

Alternatively, we can estimate the regularized version of the model using `glmnet` with `family = "binomial"`:

```
reg_logit_fit <- cv.glmnet(x, y, family = "binomial")
```

**SPOILER ALERT:** `cv.glmnet` selects all features.

# Model output

The `tidy()` and `glance()` functions from the `{broom}` package provides tidy summary of the output from `glm` objects:

```
logit_model %>% tidy()
```

```
## # A tibble: 8 x 5
##   term                estimate std.error statistic    p.value
##   <chr>              <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept)       -3.23     0.0224   -144.    0.
## 2 cough1             0.656     0.0353    18.6 4.62e- 77
## 3 fever1             1.92      0.0371    51.8 0.
## 4 sore_throat1       4.38      0.119     36.7 2.00e-294
## 5 shortness_of_breath1 4.21      0.138     30.4 1.41e-203
## 6 head_ache1         5.35      0.139     38.6 0.
## 7 age_60_and_aboveYes  0.399     0.0343    11.6 2.83e- 31
## 8 genderfemale      -0.308     0.0279   -11.0 2.34e- 28
```

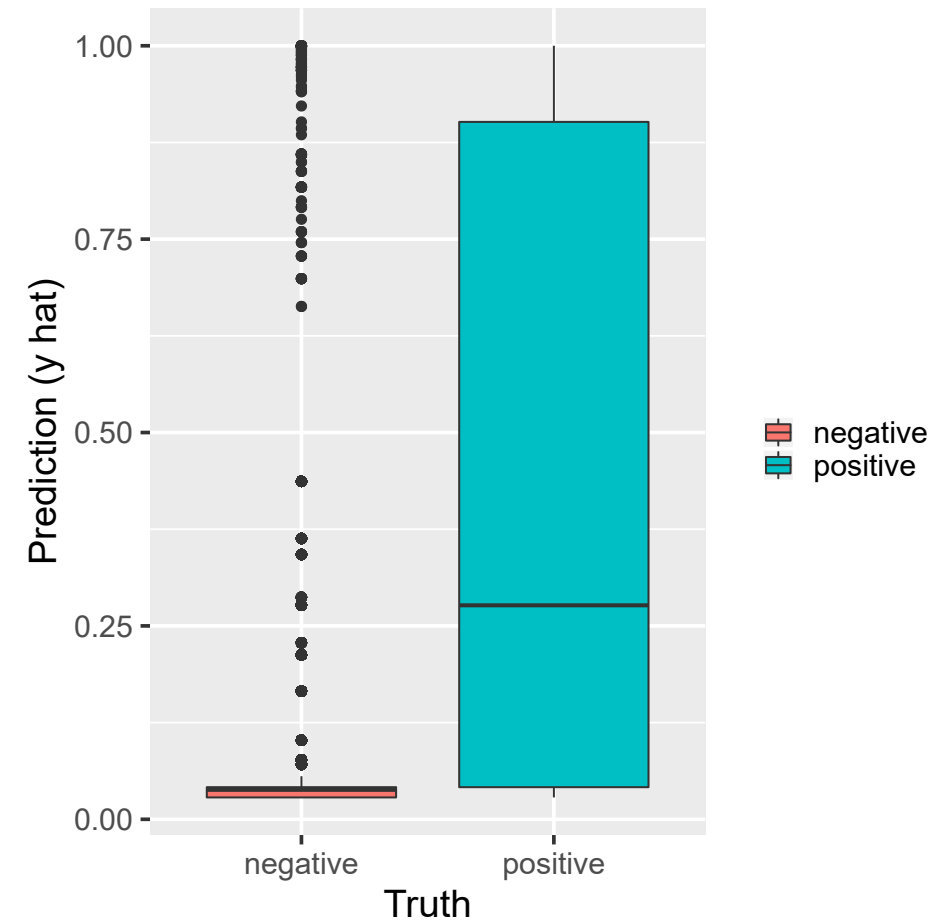
```
logit_model %>% glance()
```

```
## # A tibble: 1 x 7
##   null.deviance df.null  logLik    AIC    BIC deviance df.residual
##   <dbl>    <int>    <dbl>  <dbl>  <dbl>  <dbl>    <int>
## 1      61666.  107541 -20726. 41468. 41544.  41452.    107534
```

# Model predictions (in sample)

The figure on the right shows the resulting in-sample fit. There appears to be little overlap between probabilities for the true positives and the true negatives.

```
covid_pred %>%  
  ggplot(aes(x = corona_result,  
             y = .fitted,  
             fill = corona_result)) +  
  geom_boxplot() +  
  labs(  
    x = "Truth",  
    y = "Prediction (y hat)",  
    fill = ""  
  )
```



# Sensitivity Specificity Trade-off

# Classification rule

To classify individuals as positive/negative we first need to set a **classification rule** (cut-off), i.e., a probability  $p^*$  above which we classify an individual as positive.

For illustration, we'll set  $p^* = 0.8$ :

```
class_rule <- 0.8
```

This means that whenever  $\hat{y}_i > 0.8$ , we would classify individual  $i$  as positive.

**QUESTION:** Is this rule overly aggressive or passive?



# Classification under the rule

```
covid_pred <- logit_model %>%  
  augment(type.predict = "response") %>%  
  mutate(  
    .fitted_class = if_else(.fitted < class_rule, "negative", "positive"),  
    .fitted_class = as_factor(.fitted_class)  
  ) %>%  
  select(corona_result, .fitted, .fitted_class)  
  
covid_pred
```

```
## # A tibble: 107,542 x 3  
##   corona_result .fitted .fitted_class  
##   <fct>         <dbl> <fct>  
## 1 negative      0.0709 negative  
## 2 negative      0.342  negative  
## 3 negative      0.287  negative  
## 4 positive      0.437  negative  
## 5 negative      0.0770 negative  
## 6 positive      0.437  negative  
## 7 negative      0.0381 negative  
## 8 negative      1.00   positive  
## 9 negative      0.817  positive  
## 10 negative     1.00   positive  
## # ... with 107,532 more rows
```

# Sensitivity specificity trade-off

As we've seen, classifying everyone as "negative" ( $p^* = 1$ ), fails to be specific, i.e., it fails to identify any positive results (what we really care about!):

**Sensitivity:** The fraction of positive examples that are classified as positive ("true positive rate"),  $98,586/98,586 = 100\%$  in example.

**Specificity:** The fraction of negative examples (Yes) that are classified as negative ("true negative rate"),  $0/8,956 = 0\%$  in example.

Note that in general,

$$\text{false negative rate} = 1 - \text{specificity}$$

$$\text{false positive rate} = 1 - \text{sensitivity}$$

# Our model's confusion matrix

The function `cnf_mat()` from the `{yardstick}` package provides easy access to a model's confusion matrix and the implied performance statistics.

```
covid_conf_mat <-  
  covid_pred %>%  
  conf_mat(corona_result, .fitted_class)  
  
covid_conf_mat
```

```
##           Truth  
## Prediction negative positive  
##   negative    98455     6179  
##   positive      131     2777
```

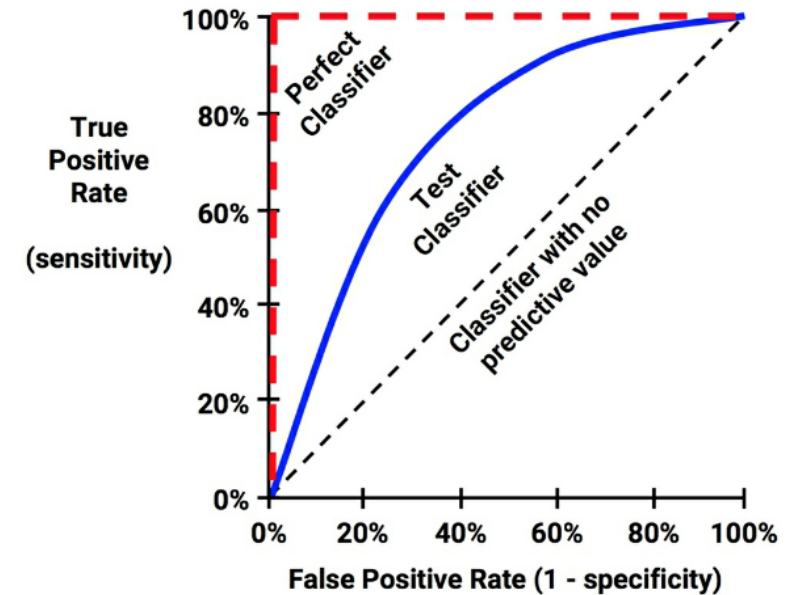
```
covid_conf_mat%>%  
  summary() %>%  
  filter(.metric %in% c("accuracy", "sens", "spec"))  
  mutate("1-estimate" = 1 - .estimate)
```

```
## # A tibble: 3 x 4  
##   .metric .estimator .estimate `1-.estimate`  
##   <chr>   <chr>      <dbl>      <dbl>  
## 1 accuracy binary      0.941      0.0587  
## 2 sens    binary      0.999      0.00133  
## 3 spec    binary      0.310      0.690
```

As we can see, for `class_rule = 0.8`, the model is highly sensitive but not so sensitive. Clearly, changing the rule would change the model's classification properties.

# Visualizing the sens-spec trade-off with ROC curves

A receiver **operating characteristic (ROC) curve**, plots sensitivity against 1-specificity. By doing so, it highlights the trade-off between false-positive and true-positive error rates as the classifier threshold is varied.



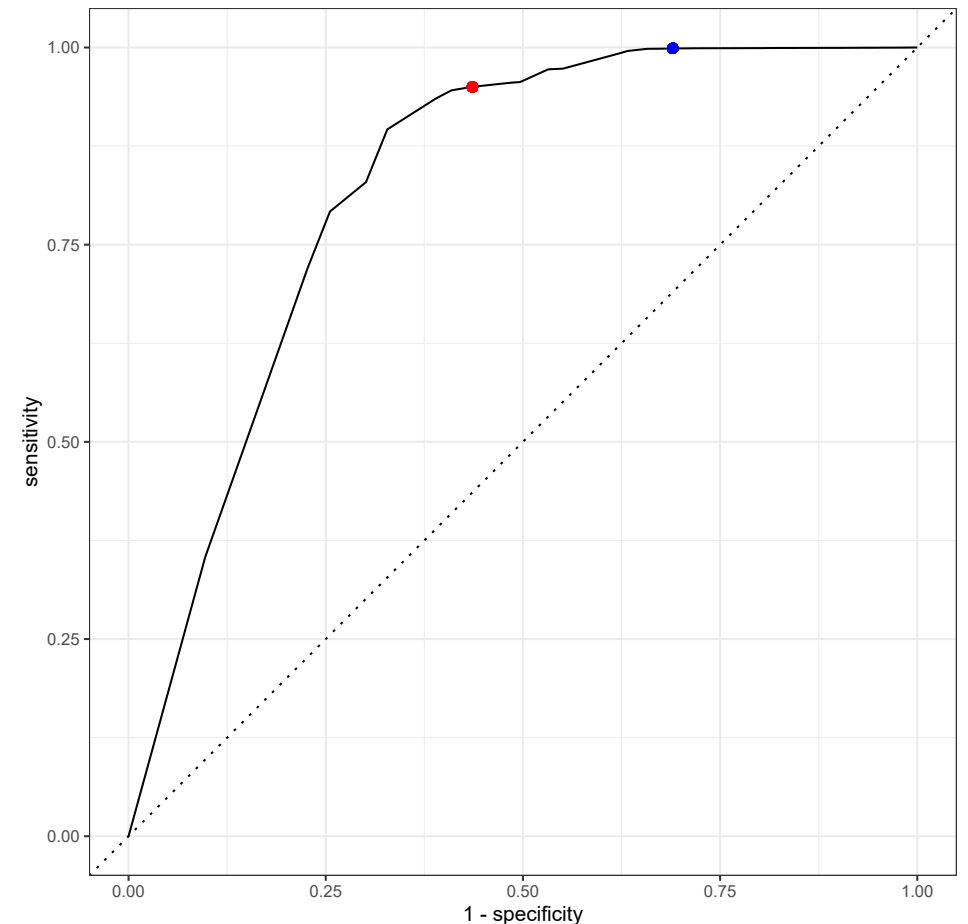
Source: "Machine Learning with R: Expert techniques for predictive modeling"

# Our model's ROC curve

On the left, you can see our model's ROC curve, plotted using the `roc_curve()` function. The red and blue dots correspond to two cut-offs, 0.8 and 0.2, respectively.

```
covid_pred %>%  
  roc_curve(corona_result, .fitted) %>%  
  autoplot() +  
  geom_point(  
    aes(x = 0.690, y = 0.999),  
    color = "blue"  
  ) + # 0.8 threshold  
  geom_point(  
    aes(x = 0.436, y = 0.950),  
    color = "red"  
  ) # 0.2 threshold
```

Note that we've used `.fitted` instead of `.fitted_class`.

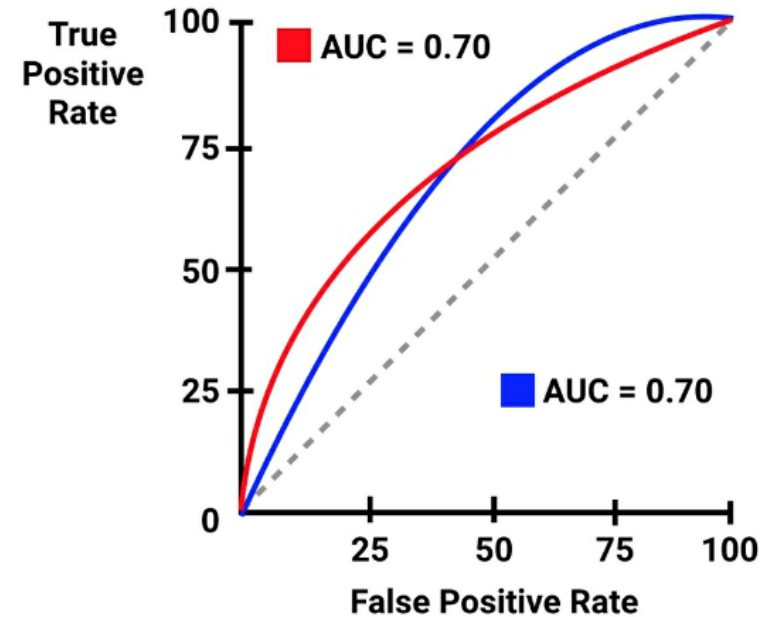


# Area under the curve (AUC)

- Ranking of classifiers can be made based on the area under the ROC curve (AUC).
- For example, a perfect classifier has  $\text{auc}=1$  and a classifier with no discriminate value has  $\text{auc}=0.5$ .
- Nevertheless, identical auc values can result from two different ROC curves. Thus, qualitative examination is warrant.

```
covid_pred %>% roc_auc(corona_result, .fitted)
```

```
## # A tibble: 1 x 3  
##   .metric .estimator .estimate  
##   <chr>   <chr>       <dbl>  
## 1 roc_auc binary      0.827
```



Source: "Machine Learning with R: Expert techniques for predictive modeling"

# AUC and cross-validation

When it comes to classification tasks, it is sometimes more reasonable to tune the penalty parameter based on classification performance metrics (and not on, say, deviance.)

For example, we can use the `cv.glmnet()` function while setting the `type.measure = "auc"` in order to tune based on auc values

```
cvfit <- cv.glmnet(  
  x, y,  
  family = "binomial",  
  type.measure = "auc"  
)
```

or set `type.measure = "class"` to tune based on the misclassification rate.

# Multiclass Classification

(Next time)



```
slides::end()
```

 [Source code](#)

# References

Lantz, Brett. Machine Learning with R: Expert techniques for predictive modeling, 3rd Edition (p. 333). Packt Publishing.