

# Lecture 6: Logistic Regression

## ENEE 691 – Machine Learning and Photonics @ UMBC

Ergun Simsek, Ph.D. and Masoud Soroush, Ph.D.

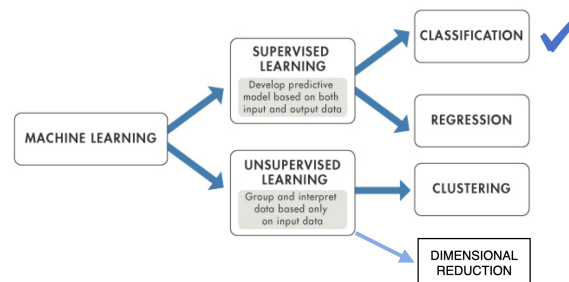
March 3, 2023

### 1 Classification Algorithms

In today's lecture<sup>1</sup>, we introduce the very first classification algorithm, namely the Logistic Regression. During the course of the next few weeks, we will introduce a wide range of classification methods that are commonly used in machine learning. Logistic regression which is the oldest classification algorithm is used to predict the value of a categorical target variable. In addition, in today's lecture, we will introduce the relevant metrics for assessing the performance of classification algorithms.

#### 1.1 Background

Classification algorithms belong to the realm of *supervised* machine learning. As we saw in previous weeks, the target variable for regression methods was a *continuous variable*. In contrast, when the target variable is categorical, we need to apply a different type of algorithm to predict the target. The process of predicting a categorical target variable in machine learning is called **classification**.



There are a wide range of classification techniques in machine learning. The very first classification technique (that we will discuss today) is logistic regression. Note that the term regression in the title of this classification algorithm shouldn't mislead you to think that logistic regression is a type

---

<sup>1</sup> This product is a property of UMBC, and no distribution is allowed. These notes are solely for your own use, as a registered student in this class.

of regression for predicting a continuous target! *Logistic regression is not a regression algorithm! It is rather a classification algorithm!* Other classification algorithms that will be studied during the course of the next few weeks are *Support Vector Machine*, *Naive Bayes*, *Decision Trees*, and some ensemble learning methods such as *Random Forest Classifier*.

## 1.2 Logistic Regression Formulation

Before introducing the formulation of logistic regression, we distinguish two types of classification: the binary and the multinomial classifications.

**Binary Classification:** When the number of possible outcomes for a categorical target variable is exactly two, we refer to the situation as a *binary classification*. Some examples of binary classification problems are:

- target variable: gender (possible outcomes: male, female)
- target variable: party affiliation in US (Possible outcomes: Democratic, Republican)
- target variable: a yes/no survey question (possible outcomes: Yes/Agree, No/Disagree)
- target variable: health condition (possible outcomes: Healthy, Not Healthy)

**Multinomial Classification:** When the number of possible outcomes for a categorical target variable is greater than two, we refer to the situation as a *multinomial classification*. Some examples of multinomial classification problems are:

- target variable: color (possible outcomes: Blue, Red, Green)
- target variable: car types (possible outcomes: SUV, Sedan, Hatchback, Convertible, Coupe)
- target variable: age category (possible outcomes: Children, Youth, Adults, Seniors)
- target variable: house type in US (possible outcomes: Single Family, Condominium, Townhouse, Multi-Family)

We will formulate logistic regression for both binary and multinomial classification problems.

### 1.2.1 Binary Logistic Regression

In this section, we formulate the binary logistic regression. To set the stage, assume that we have  $d$  continuous features collectively represented by  $\vec{x} = (x_1, x_2, \dots, x_d)$ . The target variable  $y$  is a categorical variable with exactly two possible outcomes. It is convenient to label the two outcomes as **Failure** and **Success** associated with  $y = 0$  and  $y = 1$ , respectively.

$$\begin{cases} y = 0 : & \text{Failure} \\ y = 1 : & \text{Success} \end{cases} \quad (1)$$

The above labeling is merely a choice, and you can change the labeling without affecting the predictions of the model. Suppose we have recorded the values of features and the target for  $n$  observations.

Since the target variable is categorical, we cannot assume a relation of the form  $\mathbb{Y} = \mathbb{X} \omega^\top$  (as in the case of linear regression) between the target and features. Such a relation would be meaningless as one equates a categorical quantity with a continuous quantity!. We have to take a different approach. Logistic regression takes a **probabilistic approach** to cure this problem. We notice that although  $y$  is a categorical variable, *the probability for  $y$  to assume a certain class* is continuous! We define

$$P(y^{(i)}|\vec{x}^{(i)}) : \text{probability that the target variable takes value } y^{(i)} \text{ for the } i\text{-th observation provided that the features are } \vec{x}^{(i)} = (x_1^{(i)}, x_2^{(i)}, \dots, x_d^{(i)}) . \quad (3)$$

Note that the probability in (3) is an instance of *conditional probability*. The above probability only concerns the  $i$ -th observation. To cover the whole dataset, we need to consider all  $n$  observations. We define the total probability function  $\mathcal{P}$  by

$$\mathcal{P} = P(y^{(1)}|\vec{x}^{(1)})P(y^{(2)}|\vec{x}^{(2)}) \dots P(y^{(n)}|\vec{x}^{(n)}) = \prod_{i=1}^n P(y^{(i)}|\vec{x}^{(i)}) . \quad (4)$$

**Important Note:** In defining the total probability function, we assume that all  $n$  observations of the dataset are **independent** of one another. Recall from your statistics course that if  $A$  and  $B$  are independent events, then the probability for the event  $A$  and  $B$  is given by  $P(A \cap B) = P(A) \cdot P(B)$ . That is why we multiply all the  $n$  probabilities in equation (4) associated with the  $n$  (independent) observations in the dataset. If you deal with a case where the observations are not independent of one another (*i.e.* the outcome of an earlier observation influences the outcomes of later observations), you cannot use equation (4) for the total probability.

Now that we have the total probability function (4), we should maximize it in order to find the predictions of the model for the target variable  $y$ . However, before we get into the maximization of  $\mathcal{P}$ , we need to decide how we model the individual probabilities  $P(y^{(i)}|\vec{x}^{(i)})$ .

**Question:** Since  $P(y^{(i)}|\vec{x}^{(i)})$  is a continuous quantity, can we model it via a linear relation (*i.e.*  $P(y^{(i)}|\vec{x}^{(i)}) = \omega_0 + \omega_1 x_1^{(i)} + \dots + \omega_d x_d^{(i)}$ )?

**Answer:** No! The linear (and polynomial) function is unbounded (*i.e.* the value of the function can vary from  $-\infty$  to  $+\infty$  when at least one of the weights  $\omega_i \neq 0$ ), whereas  $P(y^{(i)}|\vec{x}^{(i)})$  is a probability, and hence it belongs to the interval  $(0, 1)$ ! The probability function cannot assume any negative values or any positive values greater than 1.

To address the above problem, we need to model the probability function  $P(y^{(i)}|\vec{x}^{(i)})$  with a *bounded function*. In logistic regression, the *sigmoid function* is chosen to take care of this job. The **sigmoid function**  $\sigma$  is defined by

$$\sigma : \mathbb{R} \rightarrow (0, 1) \\ \sigma(z) = \frac{1}{1 + e^{-z}} . \quad (5)$$

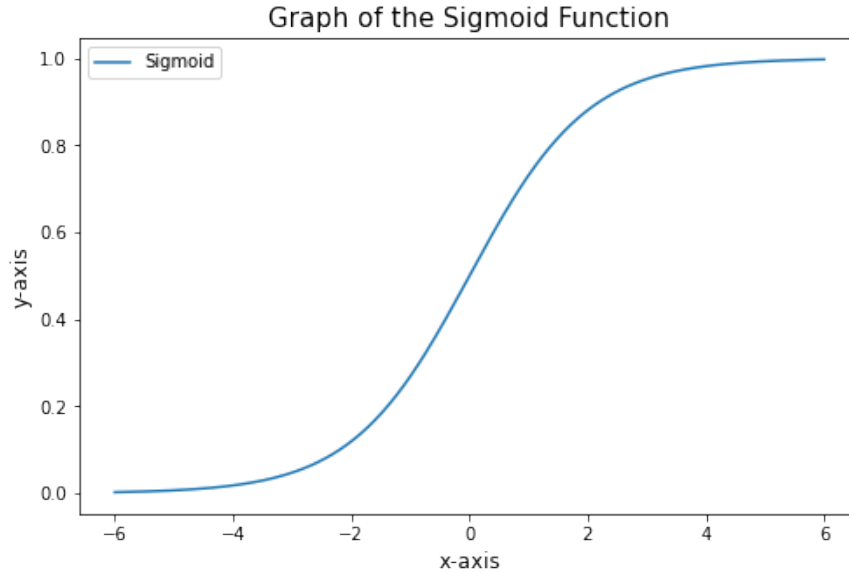
Notice that the sigmoid function (5) maps the entire real line  $\mathbb{R} = (-\infty, +\infty)$  to the interval  $(0, 1)$ . To obtain a better sense about the sigmoid function, let us plot the function in (5) and find its graph. We can easily accomplish this task in python:

```
[1]: import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

def sigmoid(z):                                # Define sigmoid function
    return 1/(1 + np.exp(-z))

x_vals = np.linspace(-6, 6, 200)
y_vals = sigmoid(x_vals)

sig_plot = sns.lineplot(x=x_vals, y=y_vals, label='Sigmoid')
sig_plot.set_xlabel("x-axis", fontsize=12)
sig_plot.set_ylabel("y-axis", fontsize=12)
sig_plot.set_title("Graph of the Sigmoid Function", fontsize=15)
plt.gcf().set_size_inches(8, 5)
plt.show()
```



As is clear from the above graph, sigmoid function  $\sigma$  is a continuous and smooth bounded function that varies in the interval  $(0, 1)$  (This is exactly what we need for a probability function). It is important to note that the sigmoid function is not the only function that can be modeled for the probability function. One can choose a different bounded function for this purpose. However, the sigmoid function is a particularly simple choice for this purpose.

We can now define a linear function  $z$

$$z^{(i)} = \omega_0 + \omega_1 x_1^{(i)} + \cdots + \omega_d x_d^{(i)} = \omega_0 + \vec{\omega} \cdot \vec{x}^{(i)}, \quad (6)$$

where  $\omega_0$  as usual represents the bias term and the weight vector  $\vec{\omega} = (\omega_1, \omega_2, \cdots, \omega_d)$  contains

the weights that correlate with the features  $\vec{x} = (x_1, x_2, \dots, x_d)$ . In order to fulfill the boundedness property of the probability function  $P(y^{(i)}|\vec{x}^{(i)})$ , we equate  $\sigma(z^{(i)})$  with the probability to observe success (*i.e.*  $y^{(i)} = 1$ ) for the  $i$ -th observation. Since success and failure are two disjoint outcomes, the probability for observing failure (*i.e.*  $y^{(i)} = 0$ ) for the  $i$ -th observation will be  $1 - \sigma(z^{(i)})$ . In summary

$$P(y^{(i)}|\vec{x}^{(i)}) = \begin{cases} \sigma(z^{(i)}), & y^{(i)} = 1 \\ 1 - \sigma(z^{(i)}), & y^{(i)} = 0 \end{cases}. \quad (7)$$

Let us denote the probability of finding success in the  $i$ -th observation by  $p(z^{(i)}) = P(y^{(i)} = 1|\vec{x}^{(i)})$ . Then, the first case of (7) implies that  $p(z^{(i)}) = \frac{1}{1 + e^{-z^{(i)}}}$ . One can then show that

$$e^{z^{(i)}} = e^{\omega_0 + \vec{\omega} \cdot \vec{x}^{(i)}} = \frac{p(z^{(i)})}{1 - p(z^{(i)})}. \quad (8)$$

The quantity on the right hand side of (8) is called the **odds** for the  $i$ -th observation. Note that odds belong to the interval  $(0, \infty)$ . Values of odds close to zero indicate a low probability of success for the  $i$ -th observation. On the other hand, large values of odds indicate a high probability of success for the  $i$ -th observation. Odds introduce a measure for the probability of success for each individual observation of the dataset. *Note that each observation has its own odds.* The logarithm of odds for the  $i$ -th observation is known as the **log-odds** or **logit**:

$$z^{(i)} = \omega_0 + \vec{\omega} \cdot \vec{x}^{(i)} = \log \left( \frac{p(z^{(i)})}{1 - p(z^{(i)})} \right), \quad (9)$$

and as is clear from equation (9), log-odds can vary from  $-\infty$  to  $+\infty$ .

Equivalently, we can combine the two branches of (7) into one single expression. Using a simple trick, this can be done and both branches of (7) can be simultaneously described as follows:

$$P(y^{(i)}|\vec{x}^{(i)}) = \sigma(z^{(i)})^{y^{(i)}} (1 - \sigma(z^{(i)}))^{1-y^{(i)}}. \quad (10)$$

Note that when  $y^{(i)} = 1$  (*i.e.* success takes place), then  $1 - y^{(i)} = 0$ , and hence  $P(y^{(i)}|\vec{x}^{(i)}) = \sigma(z^{(i)})$ . On the other hand, when  $y^{(i)} = 0$  (*i.e.* failure takes place), then  $1 - y^{(i)} = 1$  and  $P(y^{(i)}|\vec{x}^{(i)}) = 1 - \sigma(z^{(i)})$ . As usual, to present the total probability function  $\mathcal{P}$ , we need to multiply the probabilities associated with all  $n$  observations of the dataset. Hence, we arrive at

$$\boxed{\mathcal{P} = \prod_{i=1}^n \sigma(z^{(i)})^{y^{(i)}} (1 - \sigma(z^{(i)}))^{1-y^{(i)}}}. \quad (11)$$

**Goal:** We aim to find the **optimal values** for the weights  $\omega_0$  and  $\vec{\omega}$  that maximize the total probability function (11).

In maximizing the total probability function, we encounter a technical difficulty. Recall from lecture 1 that in order to maximize a function, you need to take the derivatives with regard to parameters (in this case  $\omega_0, \omega_1, \dots, \omega_d$ ) and set them to zero. From the product rule for derivatives, you recall that taking the derivative of a product is a tedious task! On the other hand, taking the derivative of a sum is very straightforward: the derivative of the sum is equal to the sum of individual derivatives! Therefore computationally, it would be highly desirable to convert the product in (11) into a sum:

$$\prod_{i=1}^n \rightarrow \sum_{i=1}^n . \quad (12)$$

We can easily do that by defining the **likelihood function** to be the logarithm of the total probability function

$$\ell(\omega_0, \vec{\omega}) = \log(\mathcal{P}) . \quad (13)$$

The logarithm on the right hand side of (13) automatically converts the product into a sum due to the property  $\log(ab) = \log(a) + \log(b)$ . Since logarithm is a monotonically increasing function, maximization of the total probability function results in the maximization of the likelihood function. In machine learning, however, people prefer to minimize a cost function (rather than maximizing a function). We can easily fulfill the desire by introducing the cost function to be

$$J(\omega_0, \vec{\omega}) = -\ell(\omega_0, \vec{\omega}) = -\sum_{i=1}^n \left[ y^{(i)} \log(\sigma(z^{(i)})) + (1 - y^{(i)}) \log(1 - \sigma(z^{(i)})) \right] . \quad (14)$$

Because of the minus sign on the right hand side of (14), the maxima of the likelihood function correspond to the minima of the cost function of the logistic regression. Note that  $J(\omega_0, \vec{\omega})$  is referred to as the **cost function of the logistic regression algorithm**. To minimize  $J(\omega_0, \vec{\omega})$ , we set  $\frac{\partial J(\omega_0, \vec{\omega})}{\partial \omega_j} = 0$  which, after some simple algebraic manipulations, lead to

$$\sum_{i=1}^n (y^{(i)} - \sigma(z^{(i)})) x_j^{(i)} = 0 , \quad (15)$$

for  $j = 0, 1, \dots, d$  with  $x_0^{(i)} = 1$ . Note that (15) represents  $d + 1$  equations, as  $j$  is a free index ranging from 0 to  $d$  (Do not take (15) as a single equation!). This system of equations need to be solved in order to determine the optimal values for  $\omega_0$  and  $\vec{\omega}$ . However, unlike the regression case, this cannot be done analytically! scikit-learn internally solves this system of equations using a gradient descent approach through the substitution  $\vec{\omega} \rightarrow \vec{\omega} - \eta \vec{\nabla} J(\vec{\omega})$ .

### 1.2.2 Multinomial Logistic Regression

In this section, we formulate the logistic regression algorithm for the multinomial classification problems. In the literature, the multinomial logistic regression is sometimes known as the *Softmax Regression* as well. The first observation to make is that the formulation presented in the previous section (*i.e.* binary logistic regression) is exclusively for the case where the number of possible

outcomes for the categorical target variable is exactly two. There is no immediate generalization of the binary case which would work for the multinomial case.

Before we formally generalize the logistic regression to the multinomial case, we introduce an important function, namely the **softmax function**. The softmax function appears very often in machine learning (and in particular in deep learning), and you need to have a solid understanding of this important function. The softmax function is defined by

$$\begin{aligned} \text{softmax} : \mathbb{R}^n &\rightarrow (0, 1)^n \\ \text{softmax}(z_1, z_2, \dots, z_n) &= \left( \frac{e^{z_1}}{\sum_{i=1}^n e^{z_i}}, \frac{e^{z_2}}{\sum_{i=1}^n e^{z_i}}, \dots, \frac{e^{z_n}}{\sum_{i=1}^n e^{z_i}} \right). \end{aligned} \quad (16)$$

Note that the input of the softmax function is an  $n$ -tuple and its output is another  $n$ -tuple. The output of the softmax function possesses two crucial properties:

- Note that every element of the output  $n$ -tuple  $\frac{e^{z_j}}{\sum_{i=1}^n e^{z_i}} = \frac{e^{z_j}}{e^{z_1} + \dots + e^{z_n}}$  belongs to the interval  $(0, 1)$ . Therefore, every element of the output  $n$ -tuple can be interpreted as a **probability**!
- The sum of the  $n$  elements of the output  $n$ -tuple is 1:  $\sum_{j=1}^n \frac{e^{z_j}}{\sum_{i=1}^n e^{z_i}} = \frac{e^{z_1} + e^{z_2} + \dots + e^{z_n}}{e^{z_1} + e^{z_2} + \dots + e^{z_n}} = 1$ .

This shows that the output of the softmax function is a **discrete probability distribution**!

The essence of the above argument can be summarized as follows: Suppose we have  $n$  entities, and we assign score to these entities. By passing the scores to the softmax function, we convert the scores into a discrete probability distribution. The individual entries of the output of the softmax function correspond to the probabilities of occurrences. The fact that the probabilities add up to 1 indicates that one of the cases will certainly take place.

Let us define the softmax function in python using numpy, and see how the function works in practice.

```
[2]: def softmax(x):                                # Define softmax function
      e = np.exp(x)
      return e/e.sum()

examples = [[100, 100], [7, 5], [1, 2, 0], [4, 2, -100], [1, 2, 3, 4, 5], [-1, 2, 0, 1, -1, 3]]

for x in examples:
    print('Softmax(%s)=%s\n' % (x, softmax(x)))
```

```
Softmax([100, 100])=[0.5 0.5]
```

```
Softmax([7, 5])=[0.88079708 0.11920292]
```

```
Softmax([1, 2, 0])=[0.24472847 0.66524096 0.09003057]
```

```
Softmax([4, 2, -100])=[8.80797078e-01 1.19202922e-01 6.00136094e-46]
```

```
Softmax([1, 2, 3, 4, 5])=[0.01165623 0.03168492 0.08612854 0.23412166
0.63640865]
```

```
Softmax([-1, 2, 0, 1, -1, 3])=[0.01152193 0.23142412 0.03131985 0.08513618
0.01152193 0.62907599]
```

**Exercise:** Calculate the softmax of (3,5,6,8), (3.5,5.5,6.5,8.5), (-2,0,1,3), (5,7,8,10). What do you conclude from this exercise?

```
[3]: four_tuples = [[3, 5, 6, 8], [3.5, 5.5, 6.5, 8.5], [-2, 0, 1, 3], [5, 7, 8, 10]]

for x in four_tuples:    # Calculating the softmax for above 4-tuples
    print('Softmax(%s)=%s\n' % (x, softmax(x)))
```

```
Softmax([3, 5, 6, 8])=[0.0056533 0.04177257 0.11354962 0.83902451]
```

```
Softmax([3.5, 5.5, 6.5, 8.5])=[0.0056533 0.04177257 0.11354962 0.83902451]
```

```
Softmax([-2, 0, 1, 3])=[0.0056533 0.04177257 0.11354962 0.83902451]
```

```
Softmax([5, 7, 8, 10])=[0.0056533 0.04177257 0.11354962 0.83902451]
```

The above exercise indicates that the softmax function enjoys an important symmetry! If we shift all elements of the input  $n$ -tuple by the same constant, the softmax function stays constant:

$$\text{softmax}(z_1 + c, z_2 + c, \dots, z_n + c) = \text{softmax}(z_1, z_2, \dots, z_n) . \quad (17)$$

Now that we know the softmax function and its basic properties, we can formulate logistic regression for the multinomial classification. Suppose we have a classification problem with  $K$  classes. This means that there are  $K$  possible values that the categorical target variable can assume. In the case of binary classification, we modeled the probability function with the sigmoid function. In the multinomial case, we model the probabilities of for different outcomes with the softmax function. We need to specify the probability functions  $P(y^{(i)} = k | \vec{x}^{(i)})$  (i.e. the probability that the target  $y$  is observed in class  $k$  in the  $i$ -th observation, given that the observed features are  $\vec{x}^{(i)}$ ). We can specify all probabilities at once by

$$(P(y^{(i)} = 1 | \vec{x}^{(i)}), P(y^{(i)} = 2 | \vec{x}^{(i)}), \dots, P(y^{(i)} = K | \vec{x}^{(i)})) = \text{softmax}(z_1^{(i)}, z_2^{(i)}, \dots, z_K^{(i)}) , \quad (18)$$

where  $z_j^{(i)} = \omega_{0,j} + \vec{\omega}_j \cdot \vec{x}^{(i)} = \omega_{0,j} + \omega_{1,j}x_1^{(i)} + \omega_{2,j}x_2^{(i)} + \dots + \omega_{d,j}x_d^{(i)}$  for  $j = 1, 2, \dots, K$ .

**Important Note:** In the case of multinomial classification, a *new set of weights* ( $\omega_{0,j}$  and  $\vec{\omega}_j$  for  $j = 1, 2, \dots, K$ ) is introduced for *each class*. The total number of parameters of the multinomial logistic regression model is  $K(d + 1)$ , with  $K$  being the number of classes and  $d$  the number of features. The parameters of the model can be accommodated in a  $(d + 1) \times K$  matrix:



$$\omega = \begin{pmatrix} \omega_{0,1} & \omega_{0,2} & \cdots & \omega_{0,K} \\ \omega_{1,1} & \omega_{1,2} & \cdots & \omega_{1,K} \\ \vdots & \vdots & \ddots & \vdots \\ \omega_{d,1} & \omega_{d,2} & \cdots & \omega_{d,K} \end{pmatrix}. \quad (19)$$

Each column of (19) casts the weights associated with a class of the target variable. The properties of the softmax function guarantees that through (18), we get a discrete probability distribution with  $K$  classes. As usual, in order to find the optimal values of the parameters (weights  $\omega$  in (19)), we need to minimize an appropriate cost function. The cost function for the multinomial case is a generalization of the cost function of the binary logistic regression. The multinomial logistic regression cost function is defined by

$$J(\omega) = - \sum_{i=1}^n \sum_{k=1}^K \log (P(y^{(i)} = k | \vec{x}^{(i)})) = - \sum_{i=1}^n \sum_{k=1}^K \log \left[ \frac{e^{z_k^{(i)}}}{\sum_{j=1}^K e^{z_j^{(i)}}} \right]. \quad (20)$$

Similar to the case of binary classification, the optimization equations (*i.e.* derivatives of  $J(\omega)$  with respect to the weights  $\omega$  being set to zero) cannot be solved analytically. scikit-learn library uses an iterative gradient descent approach to solve the optimization equations numerically.

We close the discussion on the multinomial logistic regression by two side remarks.

**Side Remark:** Due to property (17), the multinomial logistic model is *overparameterized*! We can use the symmetry property (17) to reduce the number of parameters to  $(d+1)(K-1)$ . In other words, property (17) can be used to get rid of one of the columns of (19).

**Side Remark:** It can be shown that the multinomial logistic regression with two classes (*i.e.*  $K=2$ ) is equivalent to the binary logistic regression when one uses property (17).

### 1.3 Metrics to Evaluate Classifiers

In the previous section, we introduced the very first classifier, namely the logistic regression. It is now time to illustrate how we should evaluate classifiers, and measure how good they perform. To achieve this goal, we need to introduce appropriate criteria by which classifiers can be assessed. Recall that when we introduced regression models in the previous sessions, we saw how one was supposed to assess regressors. In this section, we would like to do the same thing for classifiers. Note that the same metrics (*e.g.*  $RSE$ ,  $R^2$ -score,  $F$ -statistic, Adjusted  $R^2$ -score, etc.) introduced for assessing regressors *cannot be applied* to classifiers, and we first need to come up with the right metrics that can be applied to classifiers.

In order to make comparisons between classifiers, we first introduce two boring classifiers. These two classifiers help you understand why your constructed classifier works the way it does. In every real world classification problem, you need to make sure that your classifier performs better than these two classifiers. Moreover, to interpret your model's performance, it is important to establish by how much your model beats both of these boring classifiers.

- **Sharp Classifier:** A sharp classifier labels *all observations* with one class regardless of the values the features assume. Suppose you are working on a binary classification problem with

the goal to predict the gender based on the height. A sharp classifier labels *all observations* with one class (*i.e.* male or female) regardless of the height.

- **Monkey Classifier:** A monkey classifier labels each observation randomly (based on pure chance). In other words, if there are  $K$  classes for the target, the probability of choosing each of the  $K$  classes for the monkey classifier is  $\frac{1}{K}$ , regardless of the values the features assume. For instance, if we are dealing with a binary classification problem, the monkey classifier picks either the Success ( $y = 1$ ) or the Failure ( $y = 0$ ) with a 50-50 chance for each instance (*i.e.* monkey classifier tosses a fair coin for each observation), regardless of the values the features assume in each instance.

Both the sharp and the monkey classifiers are very boring, and in most of the real life problems, they perform poorly. However, when we construct a genuine classifier, it is important to figure out *by how much* our classifier does a better job than the two boring classifiers defined above!

### 1.3.1 Metrics to Evaluate Binary Classifiers

Now, we are ready to define the metrics by which we assess classifiers. We first introduce these metrics in the context of binary classifications, and then generalize them to the multinomial case. To be consistent with the machine learning literature, instead of labeling the two outcomes of the categorical target variable by Success and Failure, we label them by **positive** and **negative**. In this labeling the negative label is used for the bigger class, and the positive label is used for the smaller class. If the sizes of the two classes are equal, you can label either of the them with the positive and the other one with the negative label. In assessing different classifiers, it is important to make a distinction between the two following types of qualitatively different classification problems:

- **Balanced Case:** In a balanced binary classification problem, the size of the positive class is close to the size of the negative class (Note that still the size of the negative class is, by definition, greater than the size of the positive class, but in the balanced case, not much bigger). In other words, the sizes of the two classes are comparable. An example of a balanced classification is to predict the gender of fair sample of people based on their height. The size of the positive class (say males) are more or less the same as the size of the negative class (say females).
- **Unbalanced Case:** In an unbalanced binary classification problem, the size of the negative class is much bigger (orders of magnitude bigger) than the size of the positive class. A prominent example of an unbalanced binary classification is to predict whether individuals chosen from a fair sample of participants are diagnosed with cancer. In this example, the size of the positive class (*i.e.* cancer positive) is much much smaller than the size of the negative class (*i.e.* cancer negative).

In each binary classification problem, in order to understand the behavior of your classifier, you should ask yourself which type of classification (balanced or unbalanced) you are dealing with.

With the terminology introduced so far, we can imagine four possible cases to take place for each instance in a binary classification problem:

- **True Positive:** A classifier labels a positive instance as positive, resulting in a win for the classifier.
- **True Negative:** A classifier labels a negative item as negative, resulting in a win for the classifier.

- **False Positive:** A classifier mistakenly labels a negative item as a positive, resulting in a **Type I** classification error.
- **False Negative:** A classifier mistakenly labels a positive item as a negative, resulting in a **Type II** classification error.

We can conveniently represent the sizes of the above four distinct cases through a matrix. This matrix is called the **Confusion Matrix**, as is defined by

$$CM = \begin{pmatrix} TP & FN \\ FP & TN \end{pmatrix}, \quad (21)$$

where  $TP$ ,  $FN$ ,  $FP$ , and  $TN$  represent the total number of true positives, false negatives, false positives, and true negatives, respectively.

**Accuracy Statistic** The first metric defined to measure the performance of a classifier is the **accuracy**. Accuracy is defined as the ratio of the number of correct predictions over total predictions:

$$accuracy = \frac{TP + TN}{TP + TN + FN + FP}. \quad (22)$$

Multiplying the fraction in (22) by 100, accuracy is usually expressed in percentages. Also note that accuracy is a metric that is defined for the classification altogether (*i.e.* It is not defined for classes individually). The rest of the metrics that we will study are defined for each class of the classification problem. Accuracy is a metric that is relatively easy to explain, and it is worth providing in any classification evaluation.

Now, let us contemplate about the accuracy of the sharp and the monkey classifiers in a balanced binary classification problem. In a balanced binary classification, the monkey classifier is expected to achieve an accuracy around 50% by (pure) random guessing. The sharp classifier would achieve the same accuracy because for the sharp all predictions are positive (they could be all negative as well), and hence:  $TN = FN = 0$ , and  $TP = FP = \frac{n}{2}$  with  $n$  being the total number of instances. This implies that a poor classifier (like the monkey or the sharp) can be 50% accurate in a balanced binary classification problem, and hence, a genuine classifier must considerably be more accurate than 50%!

**Limitation of Accuracy Alone:** Accuracy *alone* has, however, limitations as an evaluation metric. To show the limitations of accuracy (as the only evaluation metric), we consider an unbalanced binary classification problem. Consider the cancer classification problem from a fair sample. In a fair sample of people, the majority of the participants are healthy. Let us assume that the chance of finding people with cancer is 5%. Now suppose a sharp classifier labels all participants as healthy (*i.e.* the sharp classifier puts everybody in the negative class). What would be the accuracy of this sharp classifier? Well, there is no positive class predictions, hence  $TP = FP = 0$ . Thus,  $accuracy = \frac{TN}{TN + FN} = \frac{0.95n}{0.95n + 0.05n} = 0.95$ . Therefore, the accuracy of the sharp classifier is 95%! We know that a sharp classifier is a very weak classifier. Nonetheless, it achieves a remarkably high accuracy! This shows that accuracy cannot be trusted as the only evaluation metric. In particular, in dealing with unbalanced classification problems, a high accuracy alone does not mean anything!

In order to offer a wider range of metrics to evaluate classifiers, we define three more metrics, namely *precision*, *recall*, and *F1-score*. Note that unlike the accuracy, these three metrics are defined for each class of the target variable!

**Precision Statistic** Precisions for the positive and negative classes in a binary classification problem are defined as follows:

$$\begin{aligned} \text{precision}[P] &= \frac{TP}{TP + FP} , \\ \text{precision}[N] &= \frac{TN}{TN + FN} . \end{aligned} \tag{23}$$

If we now go back to the unbalanced classification example (*i.e.* predicting cancer patients) at the end of the previous section, we realize that since for the sharp classifier  $TP = FP = 0$ , it is not even possible to define  $\text{precision}[P]$ ! This is a bad sign, as a strong classifier should produce precision close to 1. Although one gets a high accuracy in the cancer classification problem, it should become evident that the sharp classifier has a problem in producing a viable result for the precision of the positive class (even though the  $\text{precision}[N]$  is still high), and this is a sign of trouble! It would be impossible for a sharp classifier to achieve high precision on both classes ( $[P]$  and  $[N]$ ) in any real classification problem.

**Recall Statistic** Recalls for the positive and negative classes in a binary classification problem are defined as follows:

$$\begin{aligned} \text{recall}[P] &= \frac{TP}{TP + FN} , \\ \text{recall}[N] &= \frac{TN}{TN + FP} . \end{aligned} \tag{24}$$

**Note:** The recall of the positive class ( $\text{recall}[P] = \frac{TP}{TP + FN}$ ) is sometimes referred to as **sensitivity** in the machine learning literature.

It is straightforward to see that recall is also able to detect the deficiency of the sharp classifier in the cancer classification problem. In that example, the recall of the negative class is high ( $\text{recall}[N] = 1$  as  $FP = 0$ ), but the recall of the positive class is 0 ( $\text{recall}[P] = \frac{0}{0+0.05n} = 0$ )! This means that the sharp classifier performs very poor on the positive class (a fact we already knew!).

**Question:** Why do we need both precision and recall? Can't we live with one of them?

**Answer:** Sometimes we might be willing to tolerate false positives more than false negatives. For instance, in the context of the cancer classification problem, a false positive case is a participant who does not actually have cancer, but has mistakenly been labeled as a cancer patient. On the other hand, a false negative case is a cancer patient who has been labeled as healthy. Due to the severe consequences, it is evident that we can tolerate false positive cases, but our tolerance for false negative is very low! From (24), it is straightforward to see that a high  $\text{recall}[P]$  (close to 1) guarantees that your classifier has a very few false negatives!

**Exercise:** Why can't we use  $\text{precision}[N]$  in the same manner to make sure that we have very few false negatives?

**Important Remark:** In training classifiers, increasing precision of a class decreases the recall of the same class, and vice versa. This phenomenon is referred to as the **precision/recall trade-off**.

**F1-Score** F1-score is defined as the *harmonic mean* of the precision and recall for each class. First, let us quickly explain what a harmonic mean is. There are three types of mean defined in statistics: the *arithmetic mean*, the *geometric mean*, and the *harmonic mean*. Arithmetic mean is the most famous type of mean that you have worked with. Suppose  $x_1, x_2, \dots, x_n$  are  $n$  real numbers. The harmonic mean  $H$  of  $x_1, x_2, \dots, x_n$  is given by

$$H = \frac{n}{\frac{1}{x_1} + \frac{1}{x_2} + \dots + \frac{1}{x_n}} = \frac{n}{\sum_{i=1}^n \frac{1}{x_i}}. \quad (25)$$

Now, the F1-scores for the positive and negative classes in a binary classification are defined by

$$\begin{aligned} F1[P] &= \frac{2 \text{precision}[P] \text{recall}[P]}{\text{precision}[P] + \text{recall}[P]}, \\ F1[N] &= \frac{2 \text{precision}[N] \text{recall}[N]}{\text{precision}[N] + \text{recall}[N]}. \end{aligned} \quad (26)$$

If one is looking for a single metric to make a good judgement about the performance of a classifier, F1-score would be that metric. F1-score is a very tough measure to beat! That is essentially because the harmonic mean is always less than or equal to the arithmetic mean (Note that F1-score is the harmonic mean of the precision and recall for each class), and the smaller number has a larger effect. Achieving a *high F1-score requires both high precision and high recall*! None of the poor classifiers (e.g. the monkey or the sharp) can produce a decent F1-score despite high accuracy and recall.

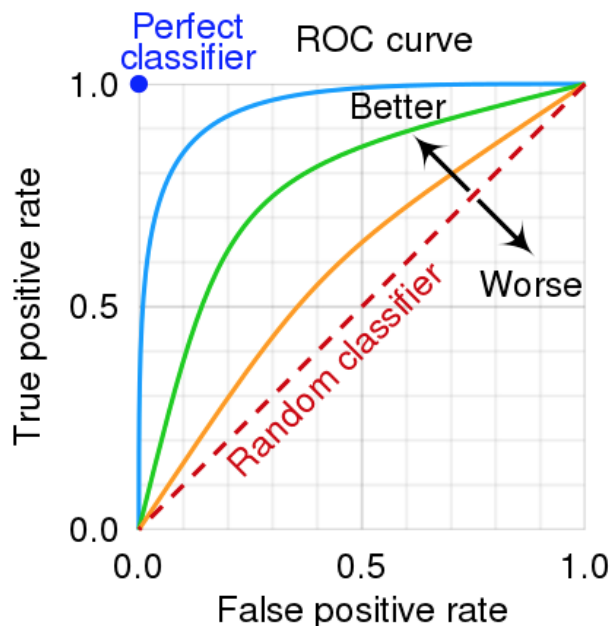
#### Lessons to Take Away:

1. Accuracy is a misleading metric in unbalanced classification problems!
2. High precision is very hard to achieve in unbalanced classification problems!
3. F1-score does the best job of any single metric, but all four of them work together to evaluate the performance of a classifier!

**Receiver-Operator Characteristic (ROC)** Receiver-Operator Characteristic (ROC) is a metric to solely evaluate *binary classifiers* (ROC is not applicable to multinomial classifiers). ROC is a graphic method of evaluating binary classifiers in which one plots the true positive rate versus the false positive rate at various thresholds. The method was originally developed for operators of military radar receivers starting in 1941, which led to its name. The **true positive rate (TPR)** and the **false positive rate (FPR)** for a binary classifier are defined by

$$\begin{aligned} TPR &= \frac{TP}{TP + FN}, \\ FPR &= \frac{FP}{FP + TN}. \end{aligned} \quad (27)$$

Note that  $TPR$  is the same as the *sensitivity* (or equivalently  $recall[P]$ ), and  $FPR$  is equal to  $FPR = 1 - specificity$ , where **specificity** is defined as  $specificity = \frac{TN}{TN+FP}$ . Since  $TPR$  and  $FPR$  are both rates, and hence the plot only concerns the square  $[0, 1] \times [0, 1]$ . The diagonal (*i.e.* the line connecting the origin  $(0, 0)$  to the point  $(1, 1)$ ) in this plot describes the performance of the monkey classifier. To evaluate your classifier, you need to calculate its corresponding  $TPR$  and  $FPR$ . Then plot the point  $(TPR, FPR)$  in the ROC plot. If the point stands above the diagonal, the performance of your classifier is better than the monkey classifier. On the other hand, if the point stands below the diagonal, your classifier performs worse than a monkey!



In order to measure by how much a classifier performs better than a monkey classifier, we connect the point  $(TPR, FPR)$  to the origin  $(0, 0)$  and to  $(1, 1)$ . The Area Under the Curve (AUC) is a metric that shows by how much a given classifier performs better than the monkey classifier. Note that AUC for the monkey classifier is  $AUC = \frac{1}{2}$  (*i.e.* the area of a half square). For a **perfect classifier**, the  $(TPR, FPR)$  point is  $(0, 1)$ , and hence  $AUC = 1$ . It is then evident that the AUC of a classifier which performs better than the monkey, but is not perfect, is  $\frac{1}{2} < AUC < 1$ .

**Note:** One can easily show that AUC (*i.e.* the area under the ROC curve) is equal to the arithmetic mean of sensitivity and specificity (*i.e.*  $AUC = \frac{sensitivity + specificity}{2}$ ).

**Note:** scikit-learn library can easily plot the ROC curve of any binary classifier, and automatically shows its corresponding AUC in the ROC plot.

### 1.3.2 Metrics to Evaluate Multinomial Classifiers

We now generalize the definitions for the evaluation metrics offered for the binary classification in the previous section to the multinomial case. Suppose, we are dealing with a classification problem with  $K$  classes. Let  $C[i, j]$  denote the number of instances of class  $i$  that received label  $j$  by the classifier. In other words,  $i$  (the first index) in  $C[i, j]$  denotes the true (actual) label and  $j$  (the second index) denotes the predicted label. We can then define a  $K \times K$  **confusion matrix** as follows:

$$CM = \begin{pmatrix} C[1,1] & C[1,2] & \cdots & C[1,K] \\ C[2,1] & C[2,2] & \cdots & C[2,K] \\ \vdots & \vdots & \ddots & \vdots \\ C[K,1] & C[K,2] & \cdots & C[K,K] \end{pmatrix}. \quad (28)$$

Note that the diagonal elements ( $C[i, i]$ ) of the confusion matrix (28) represent the number of true predictions for each class, and the off diagonal elements ( $C[i, j]$  with  $i \neq j$ ) of (28) represent the number false predictions of different pairs of classes. Equation (28) is the generalization of the confusion matrix (21) in the binary case.

The accuracy in the multinomial classification is then defined by

$$accuracy = \frac{\sum_{i=1}^K C[i, i]}{\sum_{j=1}^K \sum_{k=1}^K C[j, k]}. \quad (29)$$

The above equation is the generalization of (22) in the binary case. The precision for each class  $[i]$  of the classification is then defined by

$$precision[i] = \frac{C[i, i]}{\sum_{j=1}^K C[j, i]}. \quad (30)$$

In the denominator of (30), we consider all instances with predicted label  $[i]$ . This formula, (30), is consistent with the formula (23) in the binary case. Similarly, the recall for each class  $[i]$  is defined by

$$recall[i] = \frac{C[i, i]}{\sum_{j=1}^K C[i, j]}. \quad (31)$$

In the denominator of (31), we consider all instances with actual label  $[i]$ . The above formula, (31), is consistent with the definition (24) in the binary case. Finally, the  $F1$ -score for each class  $[i]$  is calculated as the harmonic mean of the precision and recall of the same class:

$$F1[i] = \frac{2 \, precision[i] \, recall[i]}{precision[i] + recall[i]}. \quad (32)$$

The above formula is the same as (26).

**Note:** scikit-learn library calculates the confusion matrix and all the above mentioned metrics for a classifier and represents the summary in a nice manner.