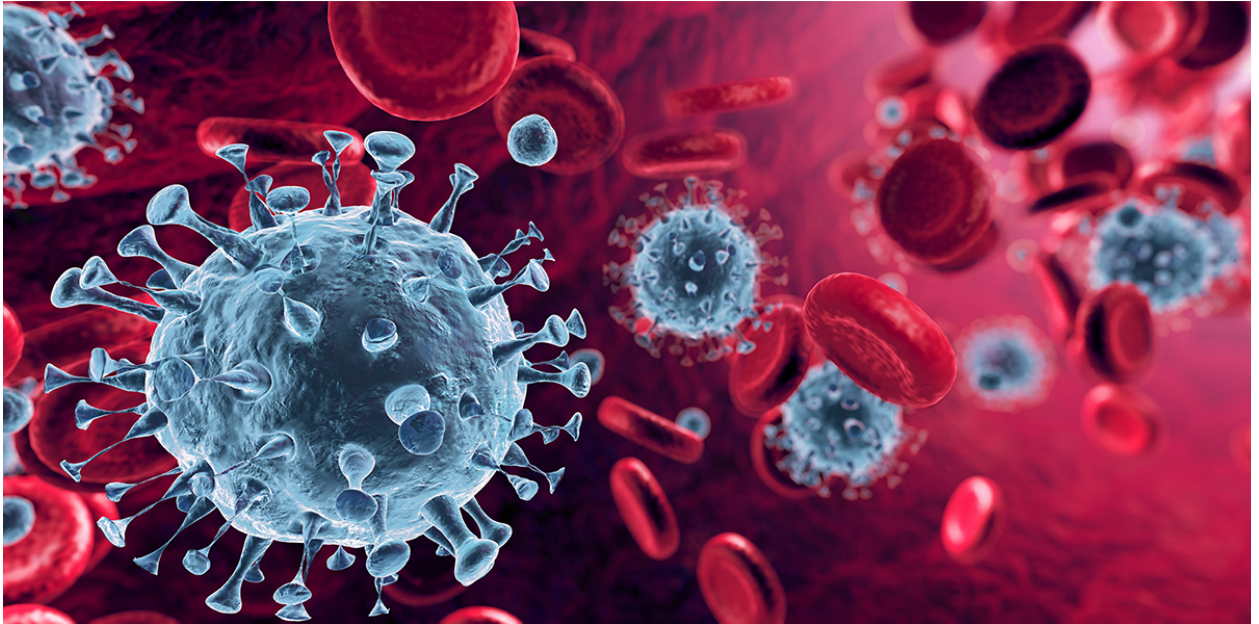


# Logistic Regression Algorithm

Using Breast Cancer Data



## The Intuition

Logistic regression is similar to simple and multiple regression. The difference is it is a classifier of binary data instead of predicting continuous variables. Furthermore, the logistic regression model uses probability to predict the dependent variable.

**Important Formulas :**

$$\hat{y}_i = \sigma(w^T x_i + b)$$

$$L_{ce}(\hat{y}_i, y_i) = -[y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i)]$$

$$Cost(w, b) = \frac{1}{N} \sum L_{ce}(\hat{y}_i, y_i)$$

$$gradient = \frac{\partial l_{ce}(w, b)}{\partial w_i} = [\sigma(w^T x_i + b) - y]x_i$$

$$Bias = [\sigma(w^T x_i + b) - y]$$

## Code

```
In [1]: 1 import pandas as pd
        2 import numpy as np
        3 from sklearn.model_selection import train_test_split
```

```

In [3]: 1 # Read the data file
2 data = pd.read_csv('../_resources/data/cancer.csv')
3
4 # Seperate by dependent and independent variables
5 x = data.iloc[:, [2, 3]].values
6 y = data.iloc[:, [1]].values
7
8 # Change y variables to 1 or 0
9 for i in range (0,len(y)):
10     if(y[i] == 'M'):
11         y[i] = 1
12     else:
13         y[i] = 0
14
15
16 data.head()

```

Out[3]:

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean
0	842302	M	17.99	10.38	122.80	1001.0	0.11840
1	842517	M	20.57	17.77	132.90	1326.0	0.08474
2	84300903	M	19.69	21.25	130.00	1203.0	0.10960
3	84348301	M	11.42	20.38	77.58	386.1	0.14250
4	84358402	M	20.29	14.34	135.10	1297.0	0.10030

5 rows × 33 columns

## Split the data

```

In [ ]: 1 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size = 0

```

## Create Logistic Regression Class

```

In [ ]: 1 class LogisticRegression:
2         def __init__(self, lr=0.01, num_iter=100000, fit_intercept=True, ve
3             self.lr = lr
4             self.num_iter = num_iter
5             self.fit_intercept = fit_intercept
6
7         def __add_intercept(self, X):
8             intercept = np.ones((X.shape[0], 1))
9             return np.concatenate((intercept, X), axis=1)
10
11        def __sigmoid(self, z):
12            return 1 / (1 + np.exp(-z))
13        def __loss(self, h, y):
14            return (-y * np.log(h) - (1 - y) * np.log(1 - h)).mean()
15
16        def fit(self, X, y):
17            if self.fit_intercept:
18                X = self.__add_intercept(X)
19
20            # weights initialization
21            self.theta = np.zeros(X.shape[1])
22
23            for i in range(self.num_iter):
24                z = np.dot(X, self.theta)
25                h = self.__sigmoid(z)
26                gradient = np.dot(X.T, (h - y)) / y.size
27                self.theta -= self.lr * gradient
28
29                if(self.verbose == True and i % 10000 == 0):
30                    z = np.dot(X, self.theta)
31                    h = self.__sigmoid(z)
32                    print(f'loss: {self.__loss(h, y)} \t')
33
34        def predict_prob(self, X):
35            if self.fit_intercept:
36                X = self.__add_intercept(X)
37
38            return self.__sigmoid(np.dot(X, self.theta))
39
40        def predict(self, X, threshold):
41            return self.predict_prob(X) >= threshold

```

```

In [ ]: 1 # Needed Variables
2 n_iter = 1000
3 theta = np.zeros(x.shape[1])
4 lr = 0.01
5
6 model = LogisticRegression(lr=0.1, num_iter=300000)
7 %time model.fit(x_train, y_train)
8
9
10 preds = model.predict(x_test, .5)
11 # accuracy
12 (preds == y_test).mean()
13

```

In [ ]: 1

In [ ]: 1