# Introduction to Machine Learning

**Daniel Sandu** 

Depending on who you ask artificial intelligence is...

Depending on who you ask artificial intelligence is...

• The horseman of the apocalypse

Depending on who you ask artificial intelligence is...

• The horseman of the apocalypse

Going to take your job

Depending on who you ask artificial intelligence is...

- The horseman of the apocalypse
- Going to take your job
- A new era of global prosperity

# **Artificial Intelligence**

But what is (artificial) intelligence?

# **Artificial Intelligence**

But what is (artificial) intelligence?

The property of a biological or artificial entity to become better at solving problems without being told explicitly what to do.

# **Artificial Intelligence**

Why artificial intelligence?

Why not procedural programming?

A couple of examples...

# Stable Diffusion



#### **ChatGPT**

What is artificial intelligence?

Artificial intelligence, or Al, refers to the development of computer systems that can perform tasks that normally require human intelligence, such as visual perception, voice recognition, decision-making, and language translation. These systems use algorithms, mathematical models, and data to learn and improve their performance over time. Al has the potential to revolutionize many aspects of our lives, from the way we work and communicate to the way we diagnose and treat diseases.



# Tesla Autopilot

11

Design machine learning algorithms starting from almost nothing

Design machine learning algorithms starting from almost nothing

What you'll need to succeed...

Design machine learning algorithms starting from almost nothing

What you'll need to succeed...

• High school mathematics – linear algebra and derivatives

Design machine learning algorithms starting from almost nothing

What you'll need to succeed...

- High school mathematics linear algebra and derivatives
- Basic programming keep it simple and always write tests

Design machine learning algorithms starting from almost nothing

What you'll need to succeed...

- High school mathematics linear algebra and derivatives
- Basic programming keep it simple and always write tests
- And most importantly... ask questions if you get stuck

### Tools of the trade

You'll need python version 3.11 (earlier may work but not tested)

#### Tools of the trade

You'll need python version 3.11 (earlier may work but not tested)

Run the following command in the console to install necessary packages

python -m pip install numpy matplotlib

# Types of machine learning

#### **Supervised learning**

The algorithm learns by being supplied labeled data.

Example: learning to provide diagnostics using metrics such as MRI scans and blood tests from previously given diagnostics

### Types of machine learning

#### **Supervised learning**

The algorithm learns by being supplied labeled data.

Example: learning to provide diagnostics using metrics such as MRI scans and blood tests from previously given diagnostics

#### **Unsupervised learning**

The algorithm learns by itself and finds structure in data.

Example: clustering people together using time spent or frequency of various activities to identify athletes, gamers, tv show watchers or other insightful categories

### Types of machine learning

#### **Supervised learning**

The algorithm learns by being supplied labeled data.

Example: learning to provide diagnostics using metrics such as MRI scans and blood tests from previously given diagnostics

#### **Unsupervised learning**

The algorithm learns by itself and finds structure in data.

Example: clustering people together using time spent or frequency of various activities to identify athletes, gamers, tv show watchers or other insightful categories

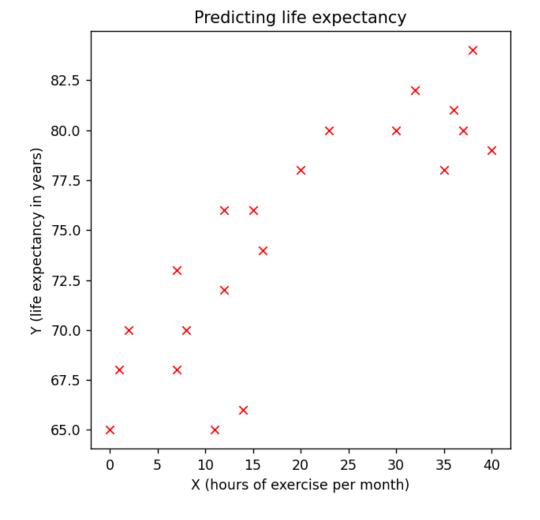
We'll be focusing on supervised learning in this course

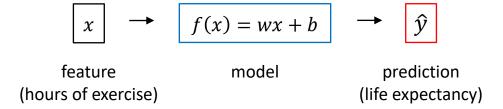
# Supervised learning

#### **Linear Regression**

Given a data set of monthly hours of exercise and life expectancy train a model to make future predictions.

Hours of exercise	Life expectancy
20	78
0	65
36	81
12	72

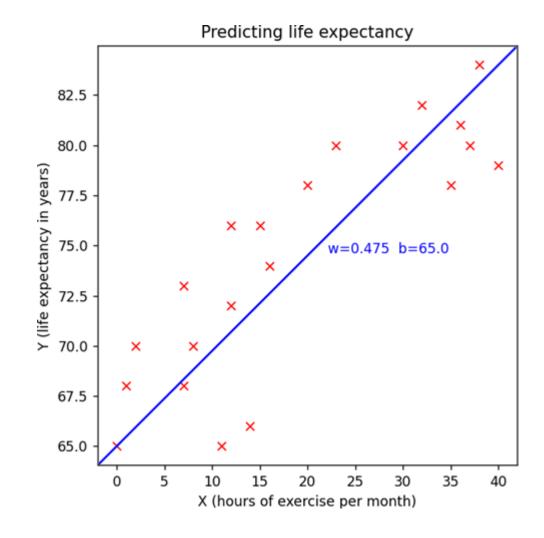




w is the weight, also called parameter, of the model.

b is the bias of the model. It's a flat increase to the prediction.

 $\hat{y}$  is our prediction and y is the ground truth.

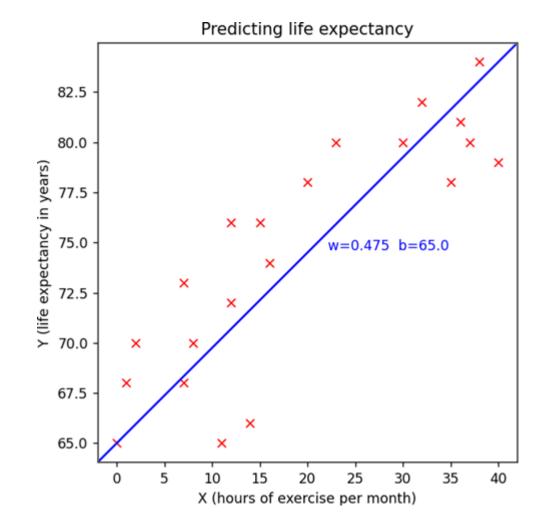


#### **Questions**

What happens if we change w?

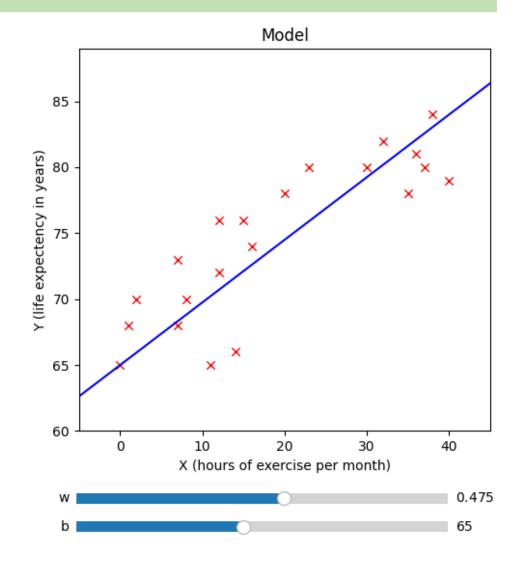
What about *b*?

Can we change w and b automatically to obtain a good model?



#### **Exercise**

Run the **predict\_life.py** script and play with the *w* and *b* values to see how they change the model



#### Can we find the parameters automatically?

We need a "score" system for each model to judge its performance

The score should be good if our predictions are close to the ground truth

We choose the model with the best score to make future predictions

# Linear Regression – Loss function

Previously we mentioned  $\hat{y}$  to be our prediction and y to be the ground truth

Consider for our score, called a loss function, the simple function  $\hat{y} - y$ 

If we are on target the score is 0

The further away we are from the target the further away from 0 we move

What are the problems with this loss function?

# Linear Regression – Loss function

We can change the loss function to  $|\hat{y} - y|$  to obtain positive numbers

Can we do better?

Now consider the loss function  $(\hat{y} - y)^2$ 

What are the benefits of this loss function?

Now consider the loss function  $(\hat{y} - y)^2$ 

What are the benefits of this loss function?

- Values are always positive
- When we are on target the result is zero
- The further way from the target the larger the loss
- Big mispredictions are taxed more than small mispredictions
- Under-predictions are symmetrical to over-predictions
- It's differentiable on its entire domain

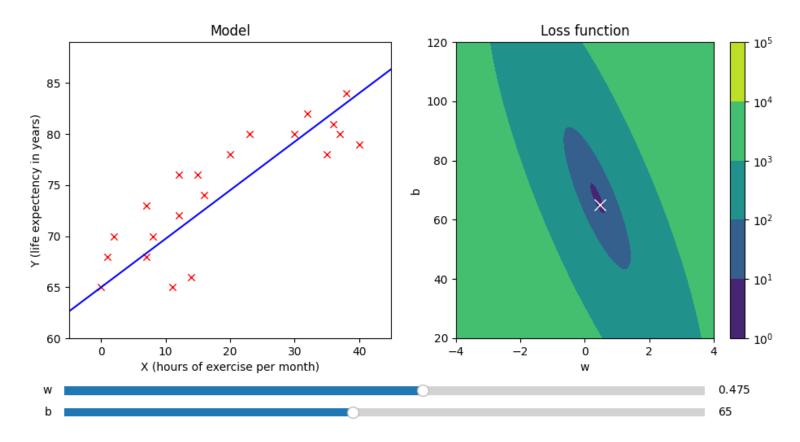
To get the overall performance of a model with parameters  $\boldsymbol{w}$  and  $\boldsymbol{b}$  we generalize the loss function to all samples in the data set

$$loss(w,b) = \frac{1}{2m} \sum_{i=1}^{m} (\hat{y}^{(i)} - y^{(i)})^2$$

m is the number of samples in the data set  $\hat{y}^{(i)}$  is the prediction for the  $i^{th}$  sample in the data set  $y^{(i)}$  is the ground truth for the  $i^{th}$  sample in the data set

#### **Exercise**

Run the **loss\_function.py** script and play with the w and b values to see the loss function behavior



### Linear Regression – Finding a good model

Now we have a way to measure the performance of our models

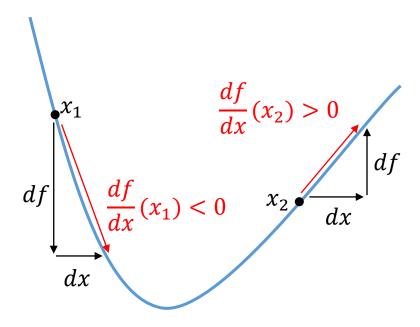
How can we find the parameters w and b automatically?

We can test values in a grid and pick up the best model, but this is slow and doesn't scale well

We can use some mathematical magic to find good parameters

# Linear Regression – Derivatives

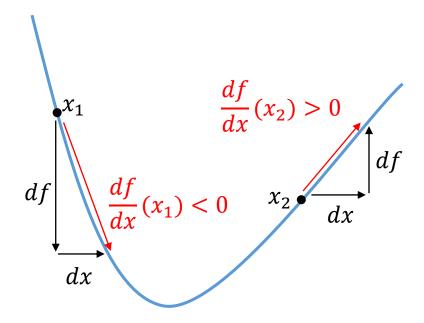
The derivative of a function f tells use how much f(x) changes if we change x by a tiny amount



# Linear Regression – Derivatives

The derivative of a function f tells use how much f(x) changes if we change x by a tiny amount

We can use the derivative of f to change x such that f(x) is decreasing



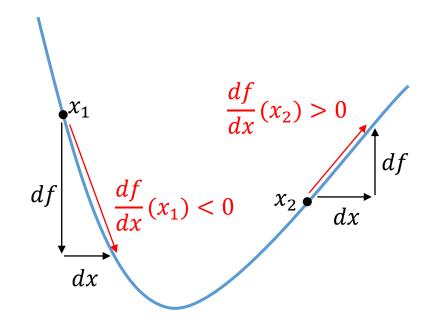
# Linear Regression – Derivatives

The derivative of a function f tells use how much f(x) changes if we change x by a tiny amount

We can use the derivative of f to change x such that f(x) is decreasing

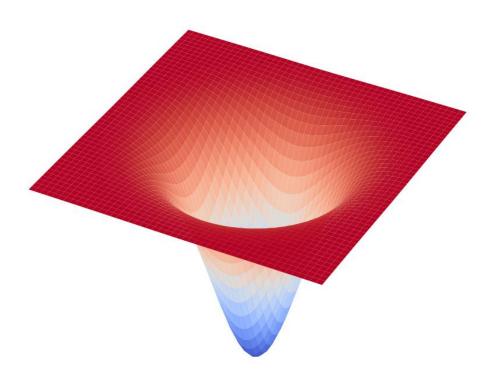
Changing 
$$x$$
 to  $x + \frac{df}{dx}(x)$  will increase  $f(x)$ 

Changing 
$$x$$
 to  $x - \frac{df}{dx}(x)$  will decrease  $f(x)$ 



# Linear Regression – Finding a good model

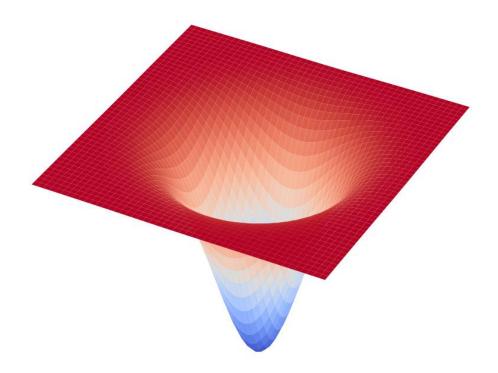
We can generalize the derivative from single to multiple function variables with partial derivatives



#### Linear Regression – Finding a good model

We can generalize the derivative from single to multiple function variables with partial derivatives

Calculating the derivative of the loss function and subtracting it from each variable will decrease the loss function



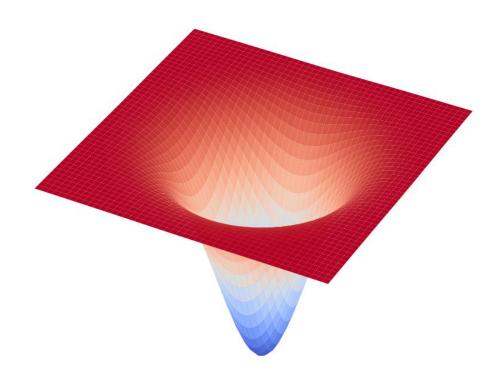
### Linear Regression – Finding a good model

We can generalize the derivative from single to multiple function variables with partial derivatives

Calculating the derivative of the loss function and subtracting it from each variable will decrease the loss function

$$w \coloneqq w - \frac{\partial}{\partial w} loss(w, b)$$

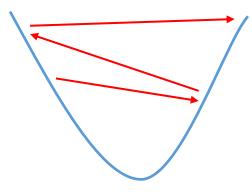
$$b \coloneqq b - \frac{\partial}{\partial b} loss(w, b)$$



This iterative process of updating the variables of the loss function is called **Gradient Descent** 

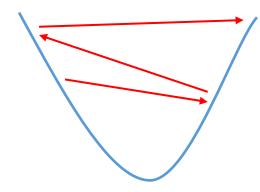
This iterative process of updating the variables of the loss function is called **Gradient Descent** 

The gradient can be very large causing the algorithm to overshoot the minimum and to not converge



This iterative process of updating the variables of the loss function is called **Gradient Descent** 

The gradient can be very large causing the algorithm to overshoot the minimum and to not converge



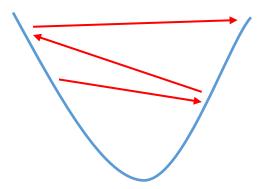
We can add the learning rate  $\alpha$  to reduce the gradient

$$w \coloneqq w - \alpha \frac{\partial}{\partial w} loss(w, b)$$

$$b \coloneqq b - \alpha \frac{\partial}{\partial b} loss(w, b)$$

We can also normalize the data set by mapping all features to the interval [-1,1] centering all them at 0

$$X_n^{(i)} = \frac{X^{(i)} - \frac{\max X + \min X}{2}}{\frac{\max X - \min X}{2}}$$



The loss function and gradient will be drastically reduced for data sets with large features

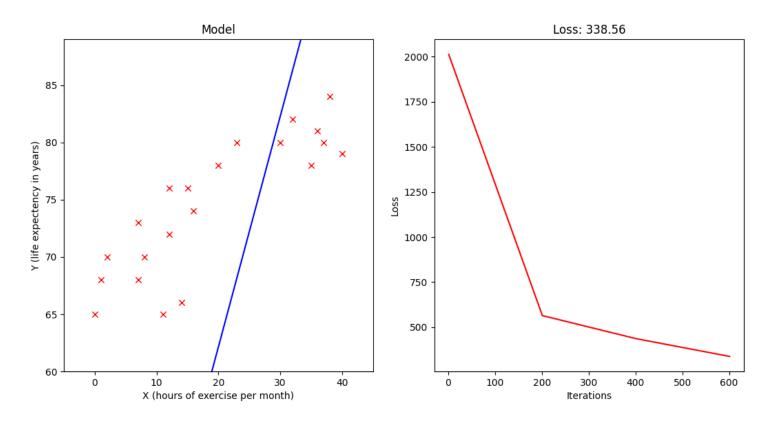
The algorithm's convergence speed will also be increased

$$w \coloneqq w - \alpha \frac{\partial}{\partial w} loss(w, b)$$

$$b \coloneqq b - \alpha \frac{\partial}{\partial b} loss(w, b)$$

#### **Exercise**

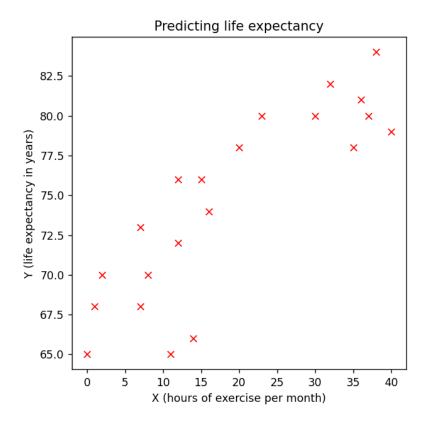
Implement gradient descent to automatically find good w and b parameters. Start with the code written in the **gradient\_descent.py** script.



# Lunch break

### Discrete prediction

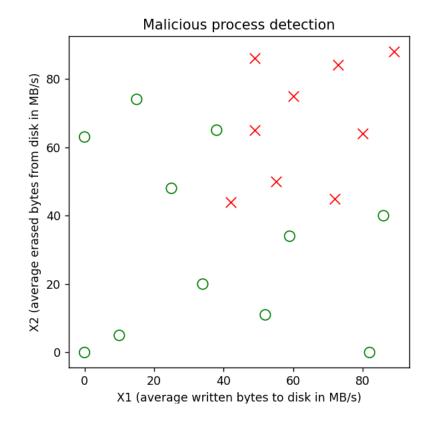
Linear Regression allows us to map features to continuous outputs



### Discrete prediction

Linear Regression allows us to map features to continuous outputs

What if we want to map to the discrete outputs 0 and 1?

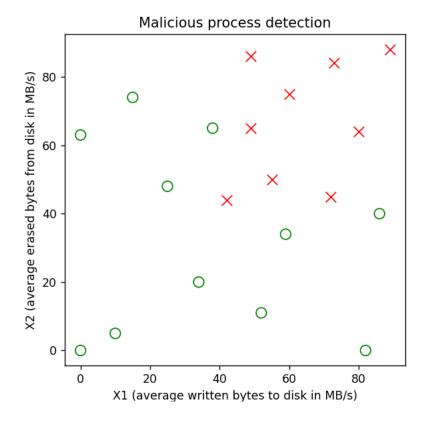


### Discrete prediction

Linear Regression allows us to map features to continuous outputs

What if we want to map to the discrete outputs 0 and 1?

Logistic regression creates a boundary line to separate positives from negatives

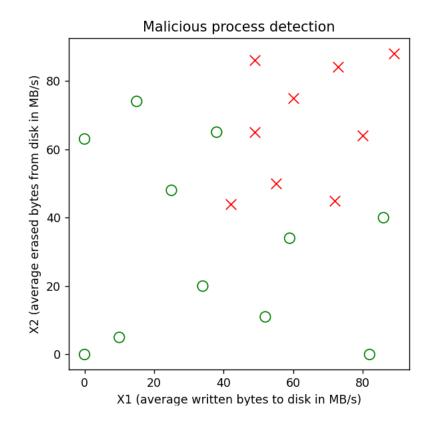


#### **Logistic Regression**

Calculate the "distance" between  $\boldsymbol{X}$  and the decision boundary  $\boldsymbol{f}$ 

Map said "distance" to the interval (0,1) using the activation function a

If a(f(X)) > 0.5 then the prediction is positive  $\hat{y} = 1$  otherwise negative  $\hat{y} = 0$ 



### **Logistic Regression**

$$X \longrightarrow f(X) = w_1 x_1 + w_2 x_2 + b \longrightarrow a(f(X)) \longrightarrow \hat{y}$$

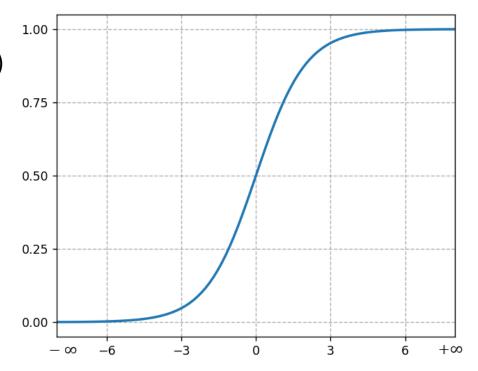
features (disk activity)

model

prediction (malicious or not)

A common activation function is the sigmoid activation function

$$a(z) = \frac{1}{e^{-z} + 1}$$



#### Logistic Regression – Loss function

#### Linear regression loss function

$$loss(w,b) = \frac{1}{2m} \sum_{i=1}^{m} (\hat{y}^{(i)} - y^{(i)})^2$$

### Logistic Regression – Loss function

#### Linear regression loss function

$$loss(w,b) = \frac{1}{2m} \sum_{i=1}^{m} (\hat{y}^{(i)} - y^{(i)})^2$$

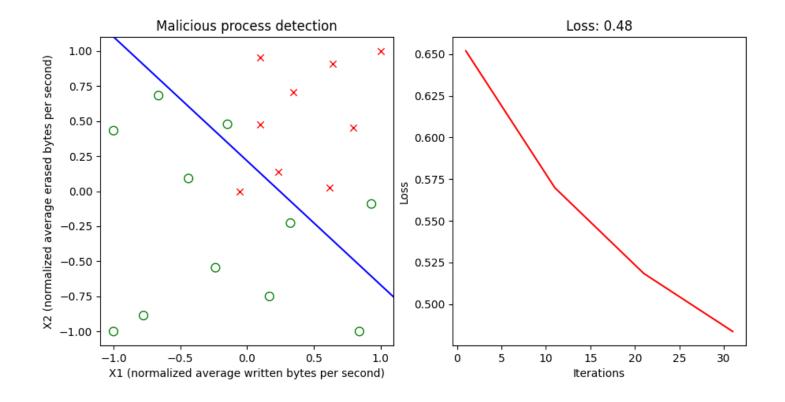
#### Logistic regression loss function

$$loss(w,b) = -\frac{1}{m} \sum_{i=1}^{m} (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}) + y^{(i)} \log(\hat{y})$$

#### Logistic Regression – Gradient Descent

#### **Exercise**

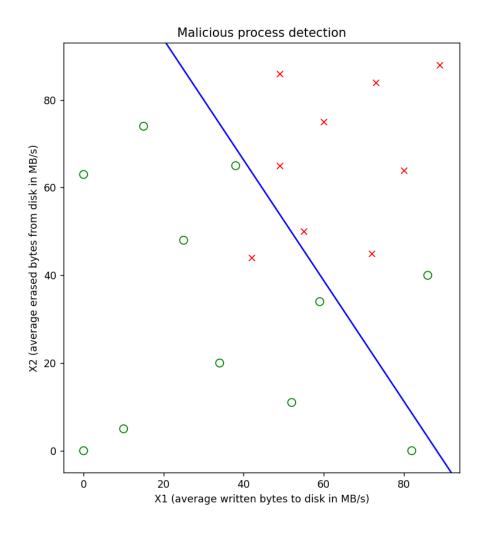
Implement gradient descent for Logistic Regression. Start with the code written in the logistic\_regression.py script as base. Note that the data has been normalized for you.



# Logistic Regression – Can we do better?

Our simple model doesn't fit the data well.

We can add more features to obtain a better fit.

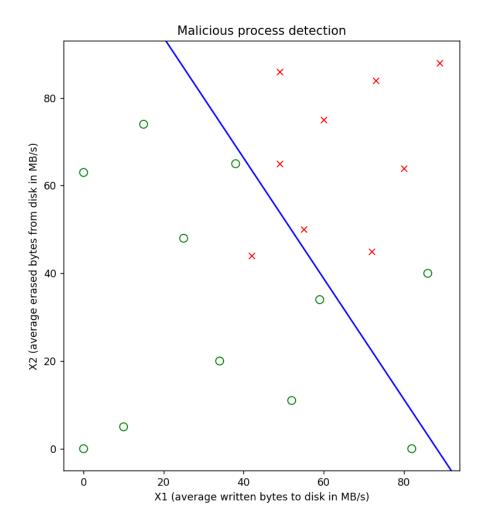


#### Logistic Regression – Can we do better?

Our simple model doesn't fit the data well

We can add more features to obtain a better fit

We are not adding new features such as the number of files created or erased but instead reuse the features  $x_1$  and  $x_2$  to keep the dimensionality to 2D to better visualize the data.



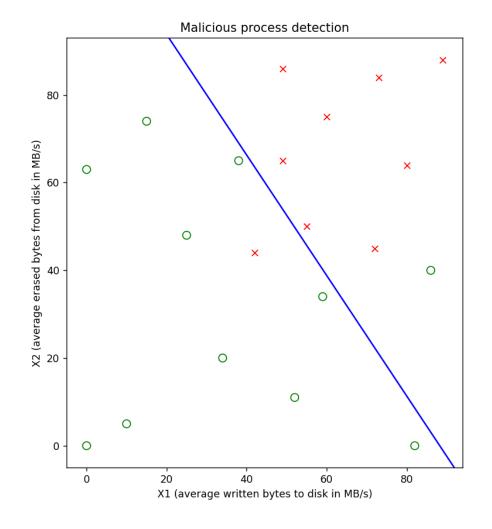
#### Logistic Regression – Can we do better?

Our simple model doesn't fit the data well

We can add more features to obtain a better fit

We are not adding new features such as the number of files created or erased but instead reuse the features  $x_1$  and  $x_2$  to keep the dimensionality to 2D to better visualize the data.

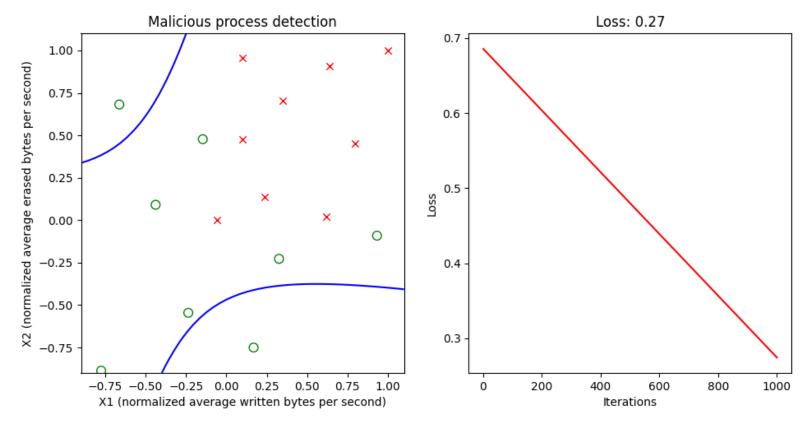
We are adding the features  $x_1^2$ ,  $x_1x_2$  and  $x_2^2$ . Our boundary will be a polynomial and thus be more flexible than a simple line.



# Logistic Polynomial Regression

#### **Exercise**

Implement logistic polynomial regression using the **polynomial.py** script to obtain a better fit of the data.



#### Logistic Polynomial Regression

#### **Exercise**

Implement logistic polynomial regression using the **polynomial.py** script to obtain a better fit of the data.

We are not adding new features such as the number of files created or erased but instead reuse the features  $x_1$  and  $x_2$  to keep the dimensionality to 2D to better visualize the data.

We are adding the features  $x_1^2$ ,  $x_1x_2$  and  $x_2^2$ . Our boundary will be a polynomial and thus be more flexible than a simple line.

### **Polynomial Regression**

#### Homework

Implement gradient descent with polynomial regression using the **gradient\_descent.py** script file as the base to obtain a better fit of the life expectancy data set.

#### **Congratulations!**

You just finished day #1 of this course!

See you tomorrow for day #2!

### Disadvantages of Linear and Logistic regression

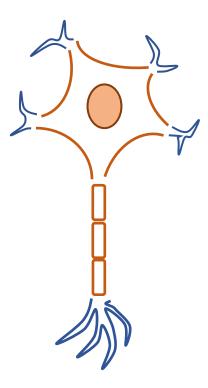
We must choose polynomial features manually for complex data sets ( $x_1^2$ ,  $x_1x_2$  ...)

It would be nice if the algorithm would automatically scale for complex data sets

This shortcoming is addressed by Artificial Neural Networks

#### **Artificial Neural Networks**

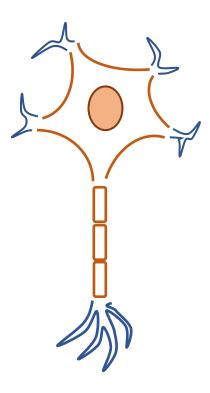
Artificial neural networks are inspired from biological neurons



#### **Artificial Neural Networks**

Artificial neural networks are inspired from biological neurons

A nerve cell can receive impulses from other nerve cells and send its own impulses



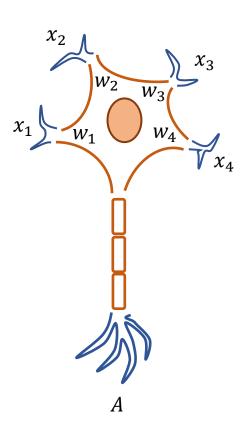
#### **Artificial Neural Networks**

Artificial neural networks are inspired from biological neurons

A nerve cell can receive impulses from other nerve cells and send its own impulses

Can be modeled mathematically as the dot product between the output of other neurons and the weights of the current neuron and its activation

$$A = a(w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + b)$$



#### Artificial Neural Networks - Layering

We can model multiple neurons using m by n matrices in  $\mathbb R$  set

Below is the activation of an entire layer of m neurons taking as input the output of the previous layer of n neurons using the activation function a

$$A = a \begin{pmatrix} \begin{pmatrix} w_{1,1} & \cdots & w_{1,n} \\ \vdots & \ddots & \vdots \\ w_{m,1} & \cdots & w_{m,n} \end{pmatrix} \begin{pmatrix} a_1 \\ \vdots \\ a_n \end{pmatrix} + \begin{pmatrix} b_1 \\ \vdots \\ b_m \end{pmatrix} = a \begin{pmatrix} \begin{pmatrix} w_{1,1}a_1 + \cdots + w_{1,n}a_n + b_1 \\ \vdots \\ w_{m,1}a_1 + \cdots + w_{m,n}a_n + b_m \end{pmatrix} \end{pmatrix}$$