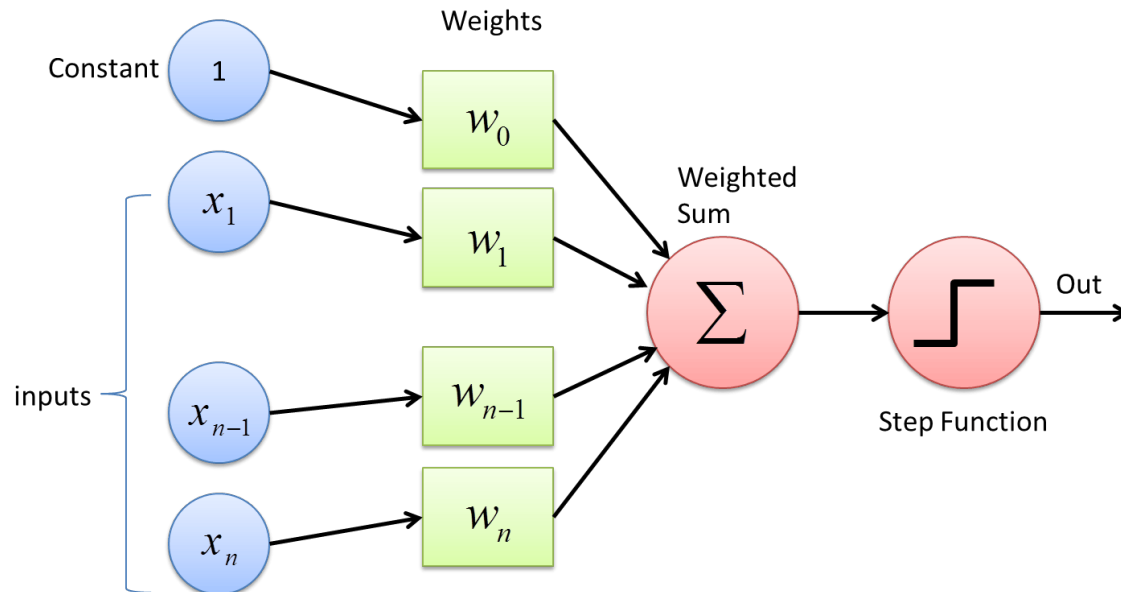
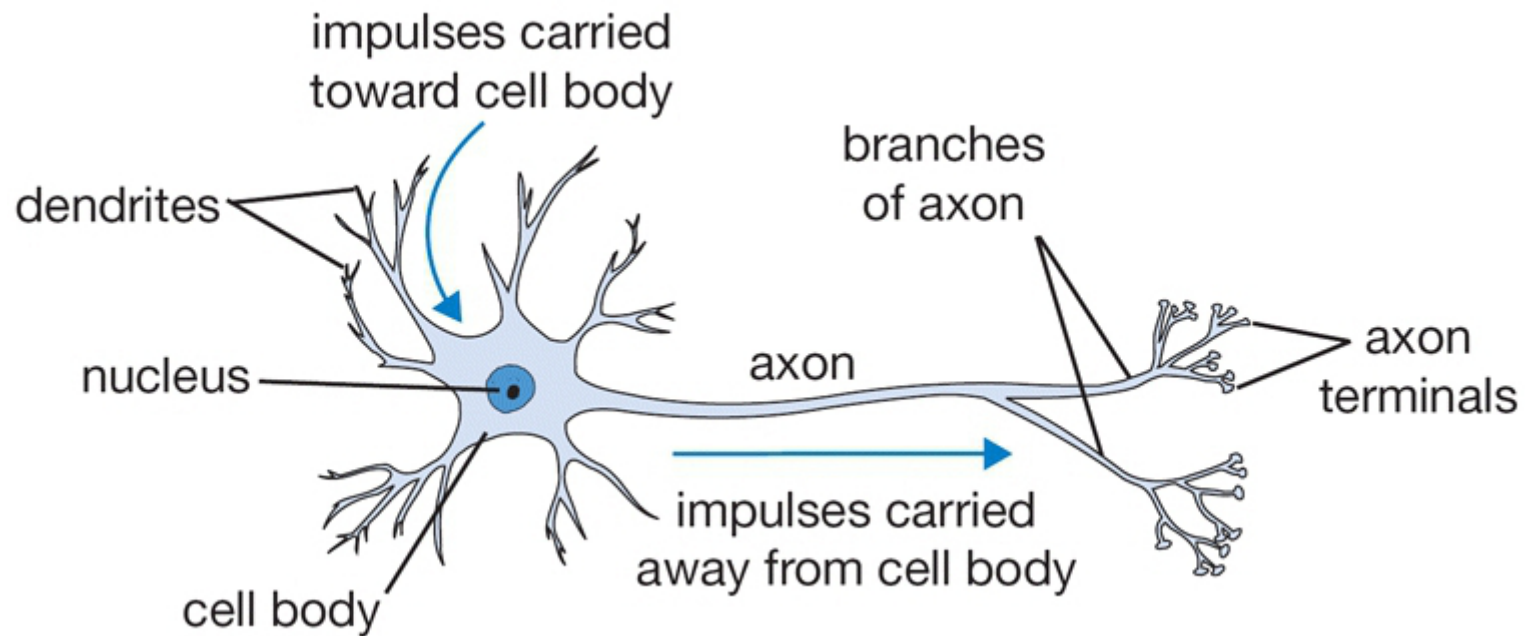


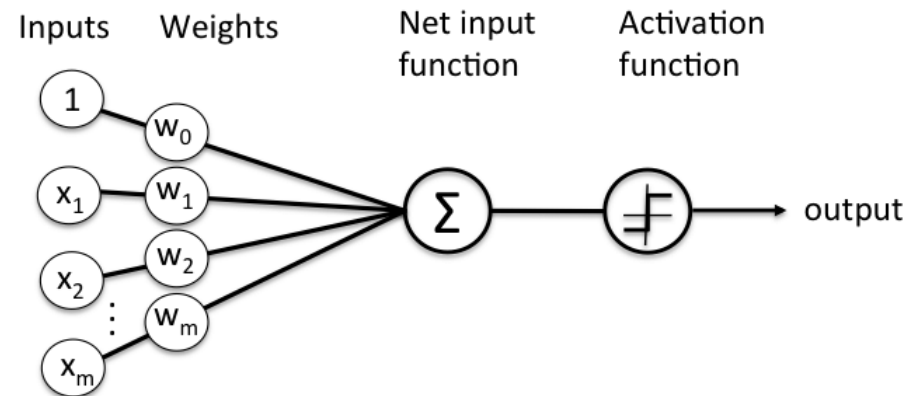
# PERCEPTRON



# Artificial neural networks are inspired in biological neural networks



The **perceptron** computes a weighted sum of its inputs, then applies an **activation function**.

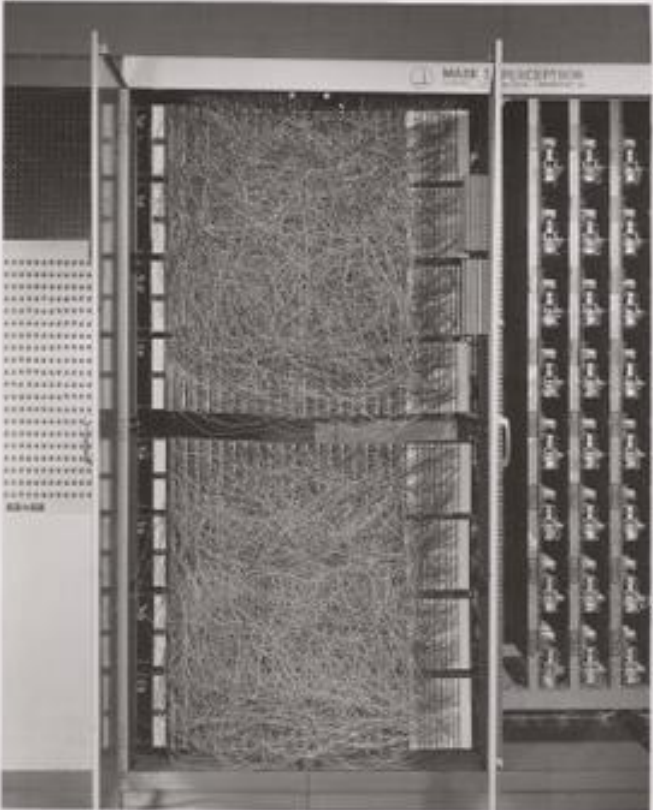


$$\hat{y} = f(\mathbf{x}^T \mathbf{w}) = f(w_0 + w_1 x_1 + \dots + w_m x_m)$$

where  $\hat{y}$  is the output,  $\mathbf{w}$  are the perceptron parameters (also called weights),  $\mathbf{x}$  are the data input features and the  $f$  is the activation function.

# HISTORY

- Invented in 1958 at the Cornell Aeronautical Laboratory by Frank Rosenblatt
- The perceptron was intended to be a machine, rather than a program
- This machine was designed for image recognition: it had an array of 400 photocells, randomly connected to the "neurons". Weights were encoded in potentiometers, and weight updates during learning were performed by electric motors



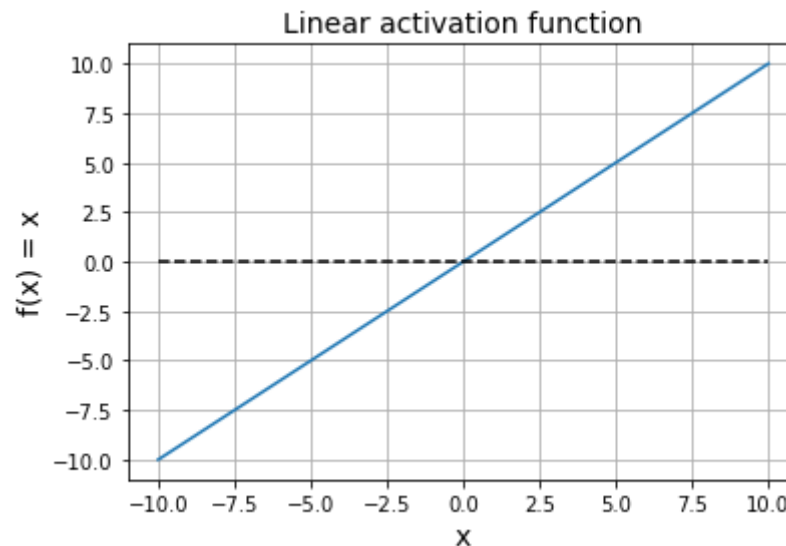
# HISTORY

- Although the perceptron initially seemed promising, it was quickly proved that perceptrons could not be trained to recognise many classes of patterns.
- First AI Winter until the 1980s
- The perceptron is used today as the basic unit of feedforward neural network with two or more layers (also called multilayer perceptrons) with greater processing power

# LINEAR REGRESSION

If we use a **linear** activation function we can use the perceptron for linear regression.

$$\hat{y} = w_0 + w_1x_1 + w_2x_2 + \dots + w_nx_n$$



# TRAINING A PERCEPTRON

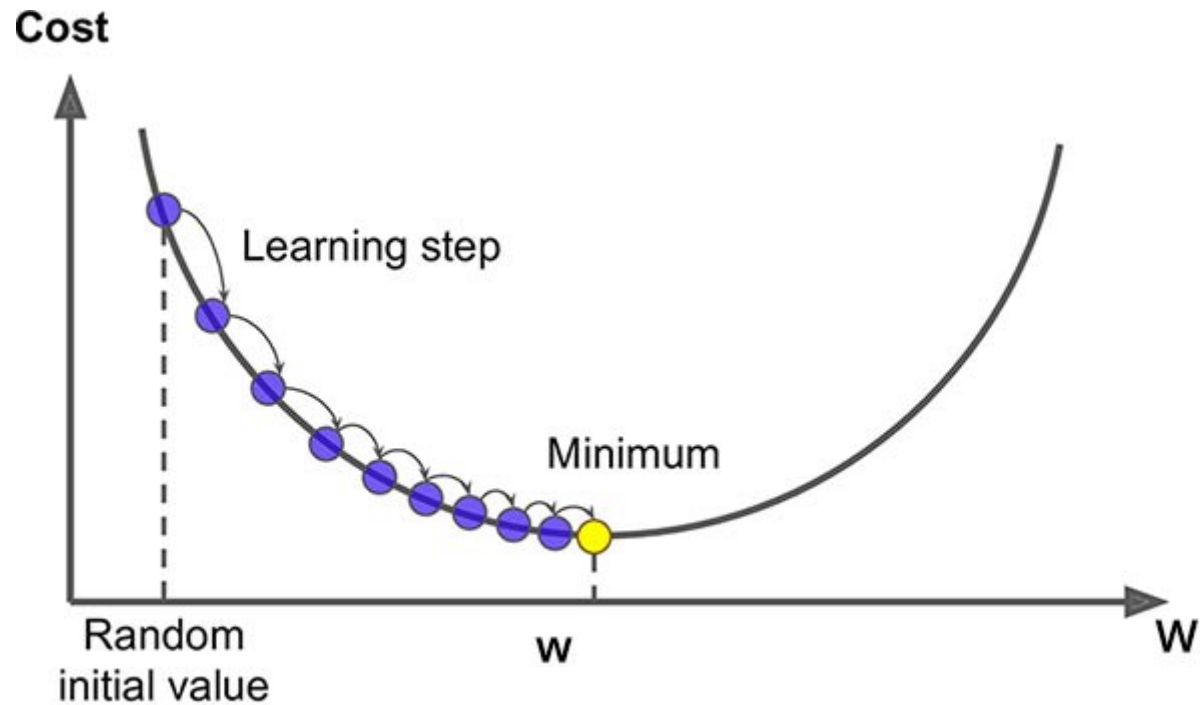
Our objective is to find the values of  $\mathbf{w}$  that minimize a cost function, for regression we usually use the **Mean Square Error** (MSE).

$$MSE(\hat{y}, y) = \frac{1}{N} \sum_{j=1}^N (\hat{y}^{(j)} - y^{(j)})^2$$

where  $N$  is the number of samples and  $y^{(i)}$  is the ground truth of the  $j^{th}$  sample in the dataset.



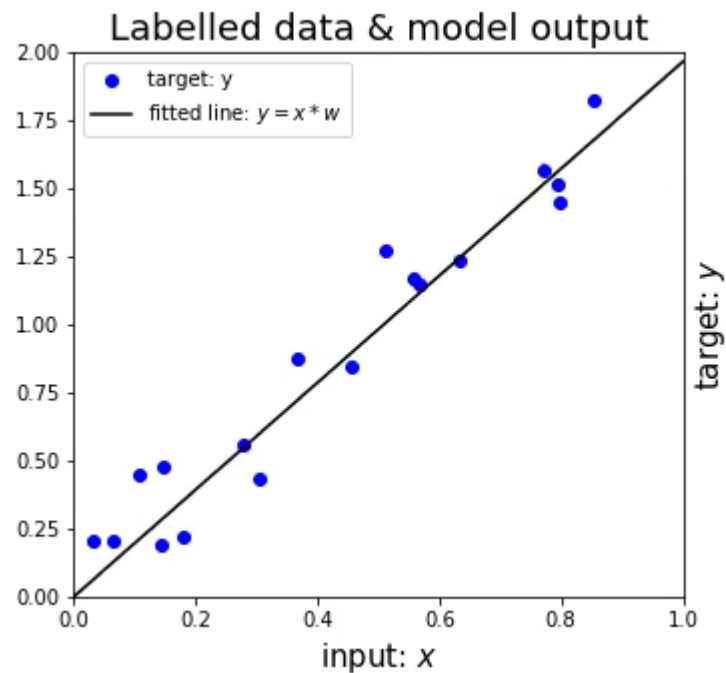
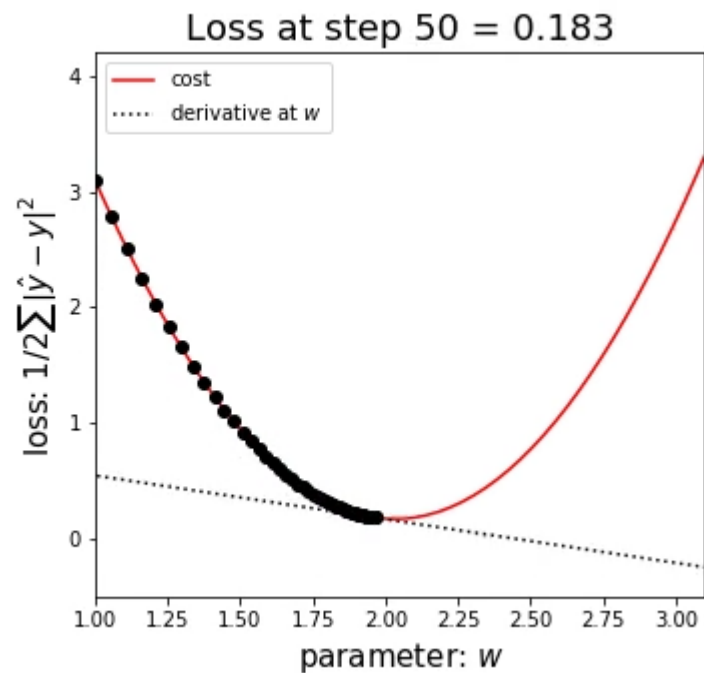
# GRADIENT DESCENT



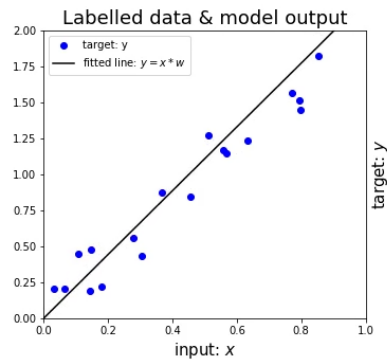
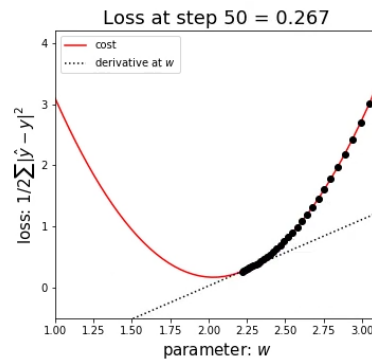
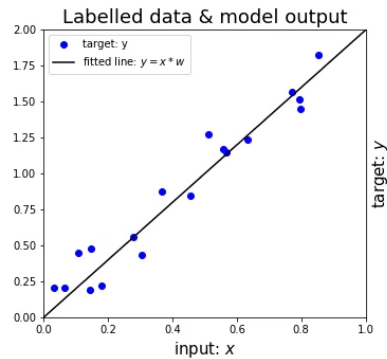
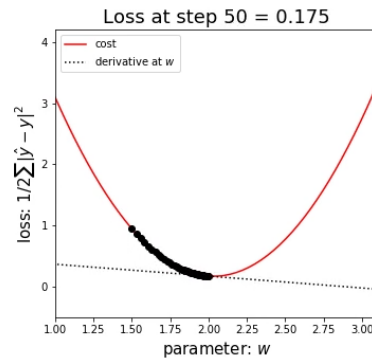
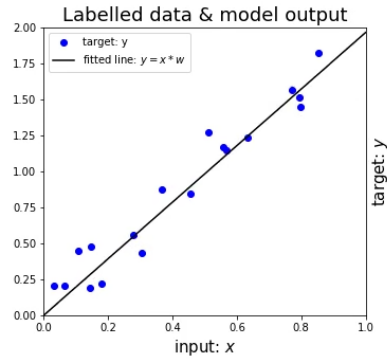
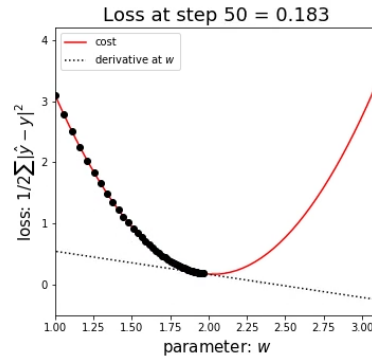
# GRADIENT DESCENT

The algorithm goes as follows:

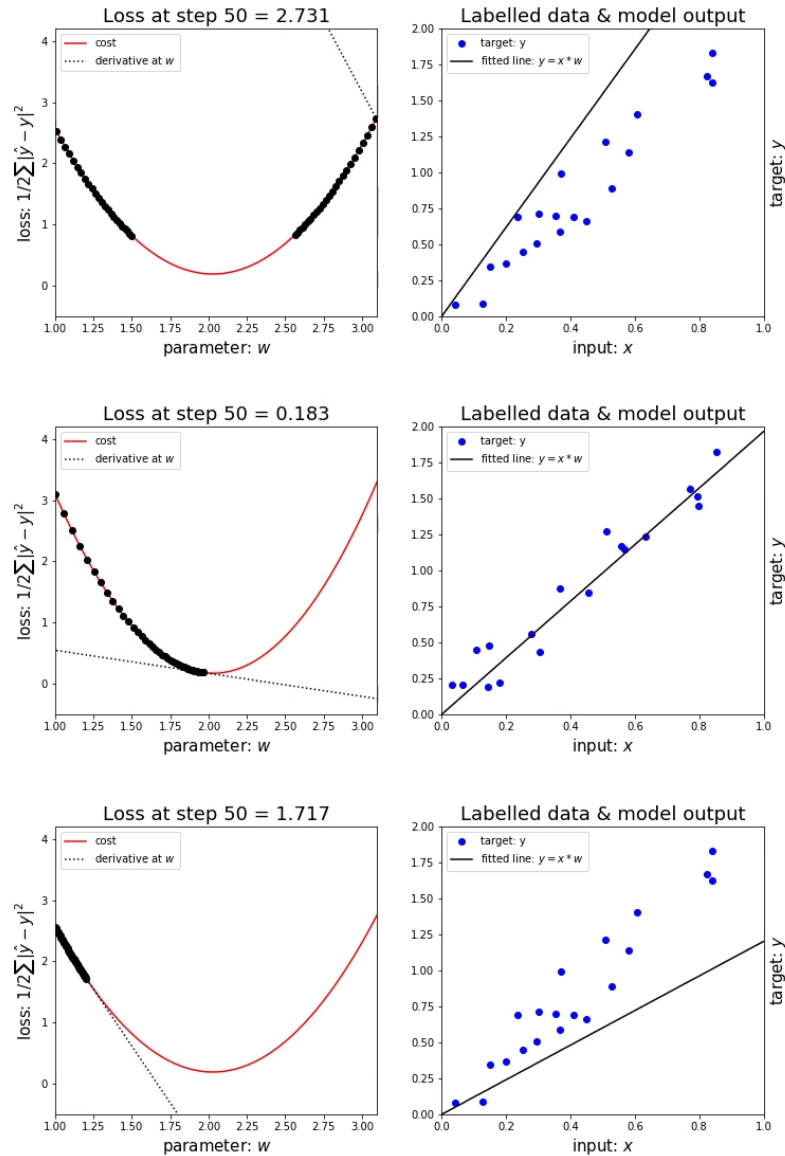
- Compute the output of the model,  $\hat{y}$ .
- Compute the gradient of the loss function with respect to the parameters,  $\frac{\partial MSE}{\partial w} = \frac{2}{N} \frac{\partial \hat{y}}{\partial w} (\hat{y} - y)$  where  $\frac{\partial \hat{y}}{\partial w} = x$ .
- Update the parameters,  $w \leftarrow w - \eta \frac{\partial MSE}{\partial w}$ , where  $\eta$  is the learning rate.
- Repeat until convergence.



# Effect of the initialization



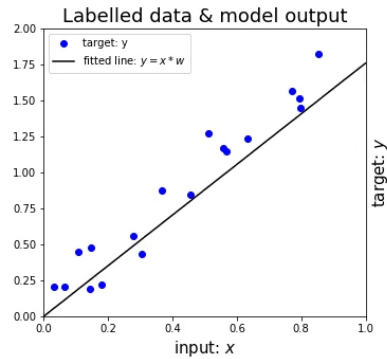
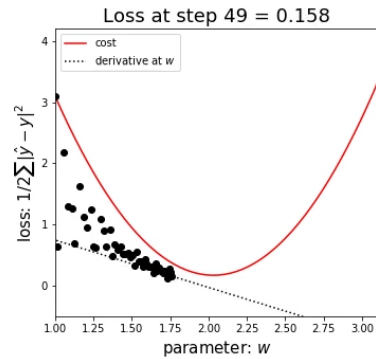
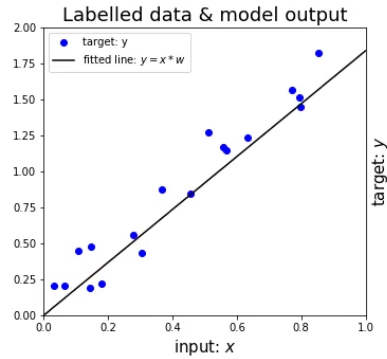
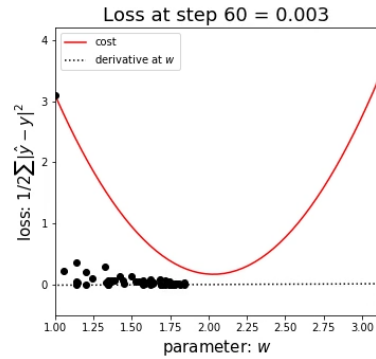
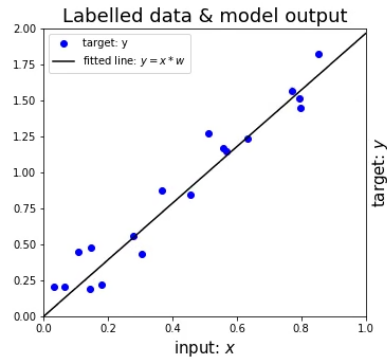
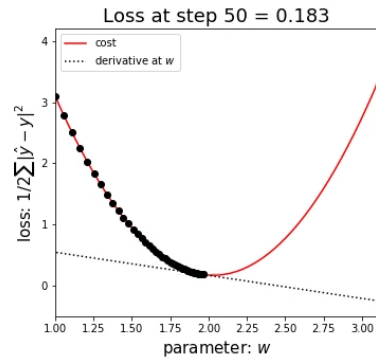
# Effect of the learning rate



# GRADIENT DESCENT FLAVOURS

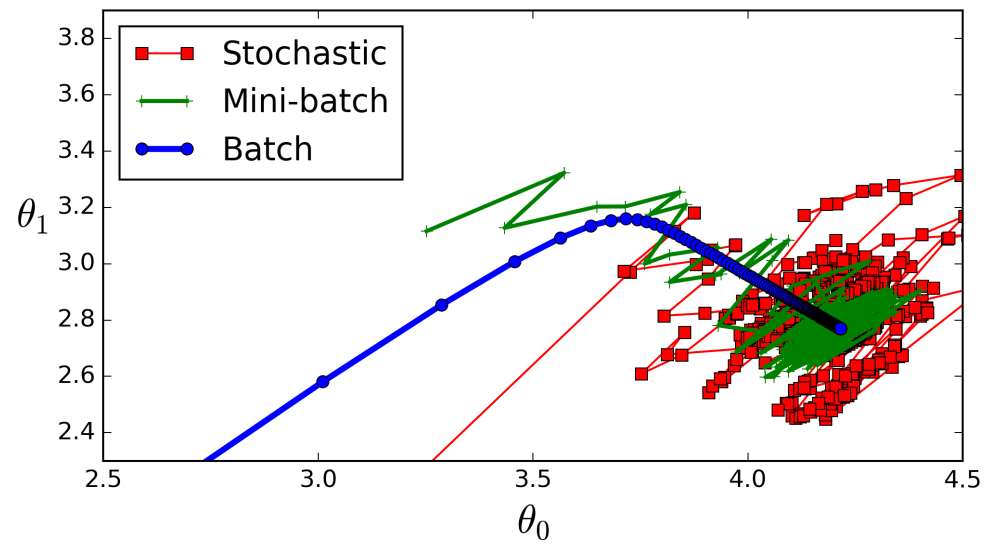
Depending on the data used to evaluate the gradients we distinguish between:

- **Batch Gradient Descent:** We use the entire dataset to compute gradients.
- **Stochastic Gradient Descent:** We use a random instance to evaluate the gradient.
- **Mini-Batch Gradient Descent:** We use a small random set of instances to evaluate gradient (default choice).



- Batch gradient descent can converge nicely to the minimum, but since it requires the entire dataset it is very slow (and sometimes impracticable).
- The alternatives converge more erratically, but faster.
- To address this issue we normally use learning rate schedulers (decrease-increase learning rate during training).



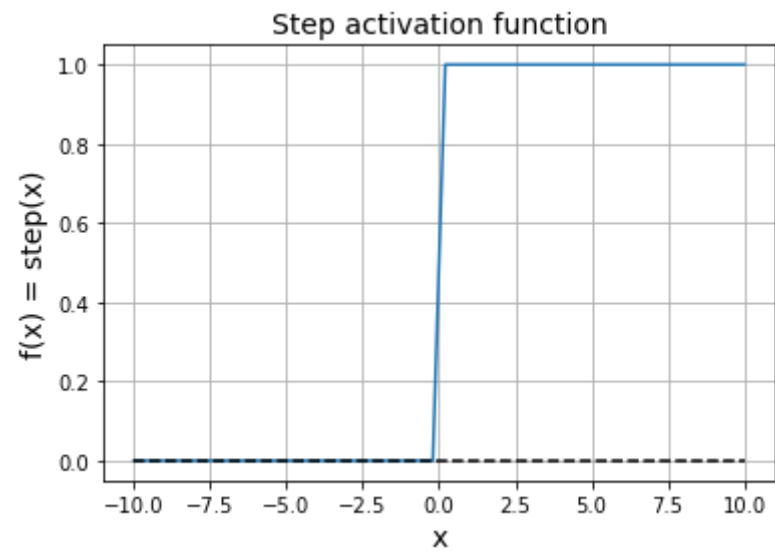


# BINARY CLASSIFICATION

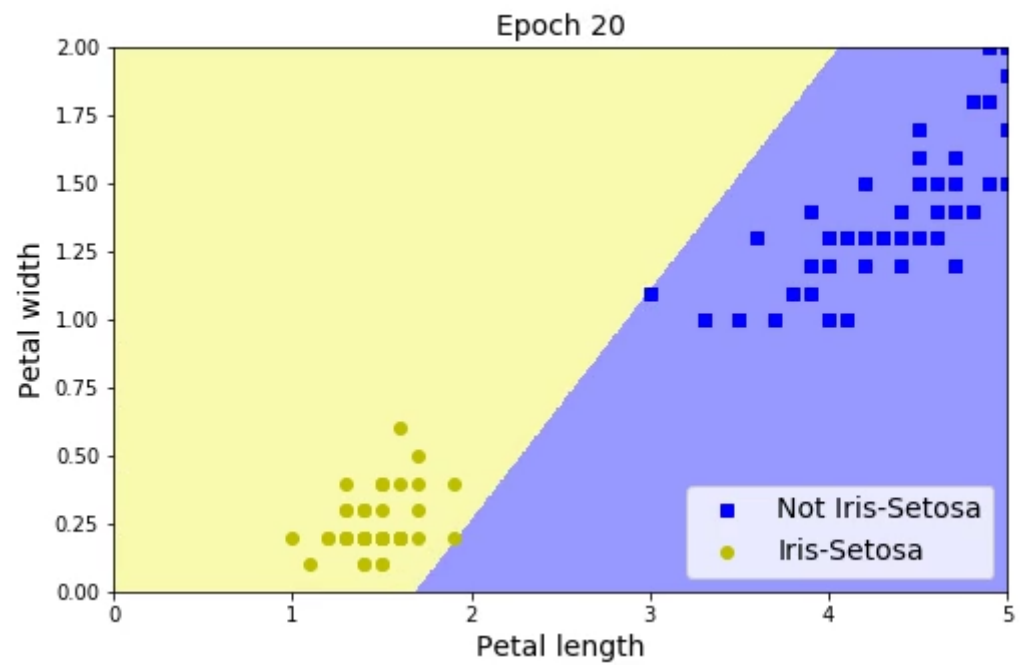
We can use a perceptron to perform binary classification, say if one particular instance belongs to a particular class or not. For that we use a **step** activation function.

$$\hat{y} = \text{step}(\mathbf{x}^T \mathbf{w}) = \text{step}(w_0 + w_1 x_1 + \dots + w_m x_m)$$

where the step function sets the output to 0 if  $\mathbf{x}^T \mathbf{w} \leq 0$  or 1 if  $\mathbf{x}^T \mathbf{w} > 0$ .

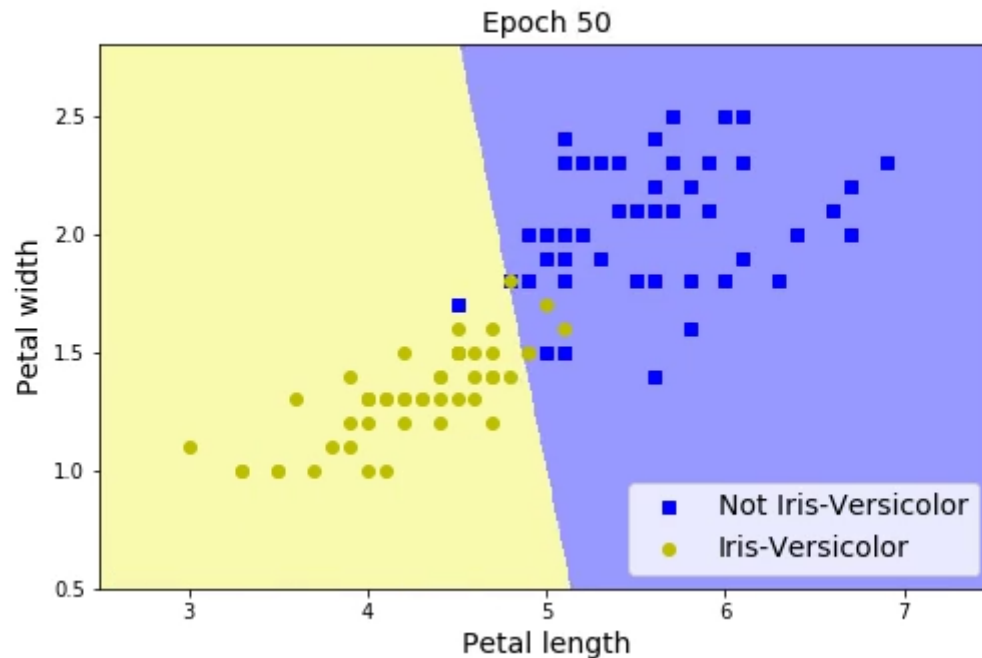


```
1 class Perceptron():
2     def __init__(self, size):
3         self.w = np.random.randn(size)
4
5     def __call__(self, w, x):
6         return w.dot(x.T) > 0
7
8     def fit(self, X, Y, epochs, lr):
9         X = np.c_[np.ones(len(X)), X]
10        for epoch in range(epochs):
11            for x, y in zip(X, Y):
12                y_hat = self(self.w, x)
13                # loss function (MSE)
14                l = 0.5*(y_hat - y)**2
15                # derivatives
16                dldh = (y_hat - y)
17                dhdw = x
18                dldw = dldh*dhdw
19                # update
20                self.w = self.w - lr*dldw
```



# LIMITATIONS

The perceptron is a **linear classifier**, therefore it will never get to the state with all the input vectors classified correctly if the training set is not **linearly separable**.

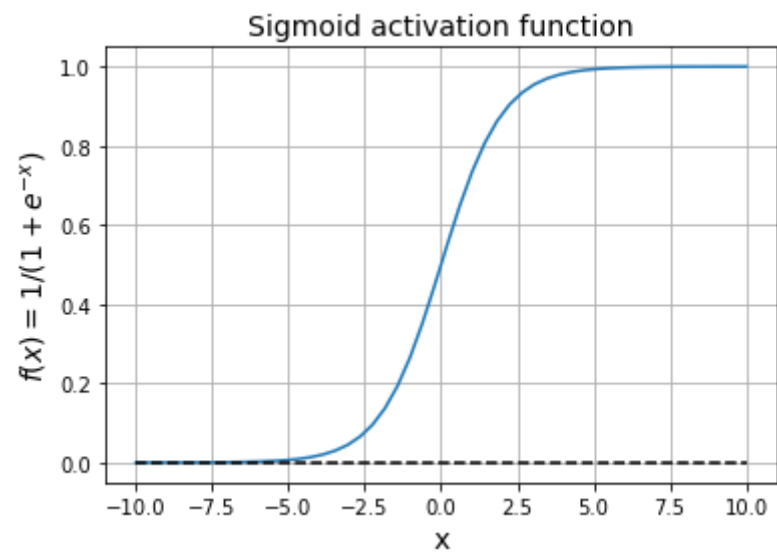


# LOGISTIC REGRESSION

We can use a perceptron to predict the probability of an instance to belong to a particular class. For that we use a **sigmoid** activation function:

$$\hat{y} = \sigma(\mathbf{x}^T \mathbf{w})$$

Where  $\sigma(x) = \frac{1}{1+e^{-x}}$  is the sigmoid function, that will limit the value of  $\hat{y}$  between 0 and 1.



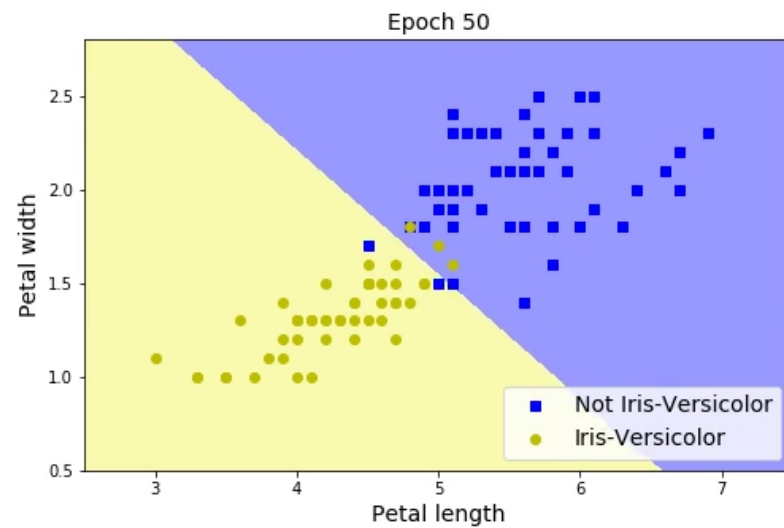
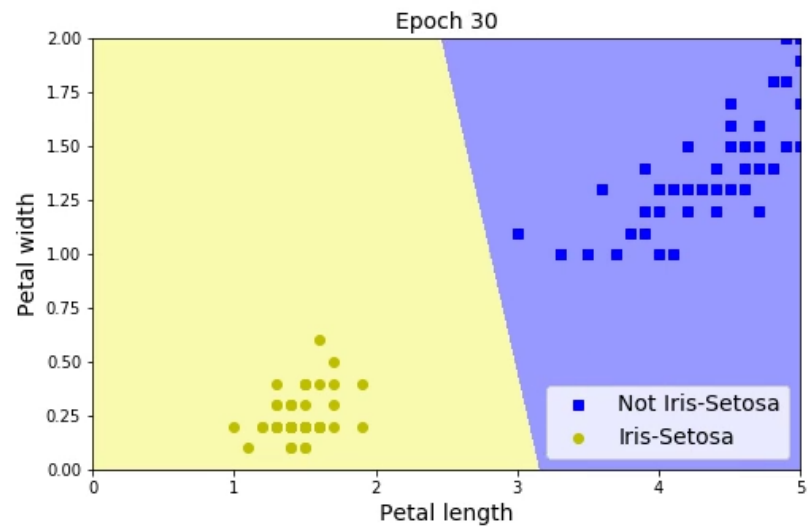


To train a logistic regression model we need a different loss function, called the **log loss** function.

$$J(\mathbf{w}) = -\frac{1}{N} \sum_{j=1}^N \left[ y^{(j)} \log(\hat{y}^{(j)}) + (1 - y^{(j)}) \log(1 - \hat{y}^{(j)}) \right]$$

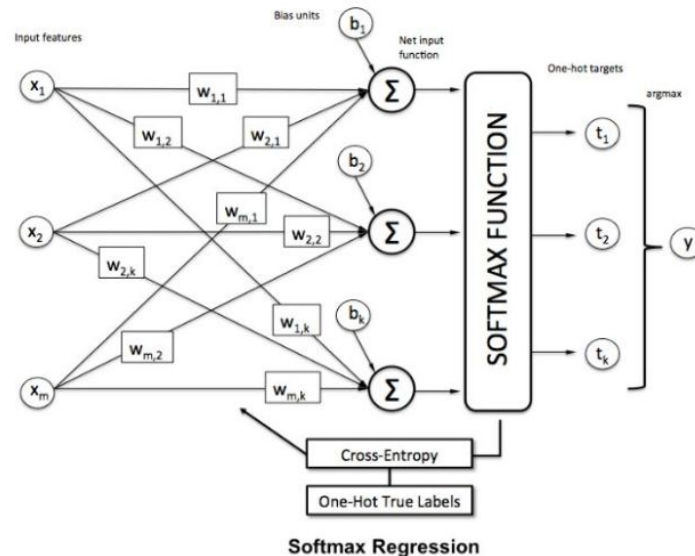
There is no known close-form to minimize this function, but gradient descent is guaranteed to find the global minimum if the learning rate is not too large. For that we need the derivative of the cost function with respect to the model parameters

```
1 class LogisticRegression():
2     def __init__(self, size):
3         self.w = np.random.randn(size)
4
5     def __call__(self, w, x):
6         return self.sigmoid(w.dot(x.T))
7
8     def sigmoid(self, x):
9         return 1. / (1. + np.exp(-x))
10
11    def fit(self, X, Y, epochs, lr):
12        X = np.c_[np.ones(len(X)), X]
13        for epoch in range(epochs):
14            for x, y in zip(X, Y):
15                y_hat = self(self.w, x)
16                self.w = self.w - lr*(y_hat - y)*x
```



# MULTI-CLASS CLASSIFICATION

The logistic regression model can be generalized to support multiple classes. For that we use a Perceptron with multiple outputs (one for each class) and then apply a **softmax** activation function.



The output can be computed as

$$\hat{y} = \arg \max_k \sigma(\mathbf{w} \cdot \mathbf{x})_k$$

where now  $\sigma(t)_k = \frac{e^{t_k}}{\sum_{k=1}^K e^{t_k}}$  is the softmax function,  $K$  is the number of classes and  $t_k$  is the score of each class for the instance  $\mathbf{x}$ .

The loss function used for training softmax regression model is called  
Cross Entropy

$$J(\mathbf{w}) = -\frac{1}{N} \sum_{j=1}^N \sum_{k=1}^K y_k^{(j)} \log(\hat{y}_k^{(j)})$$

where  $y_k^{(j)}$  is the probability that the  $j^{th}$  instance belongs to class  $k$  (usually 1 or 0). If  $K = 2$  the cost function is equivalent to the log loss function. The gradient of the cross entropy cost function is

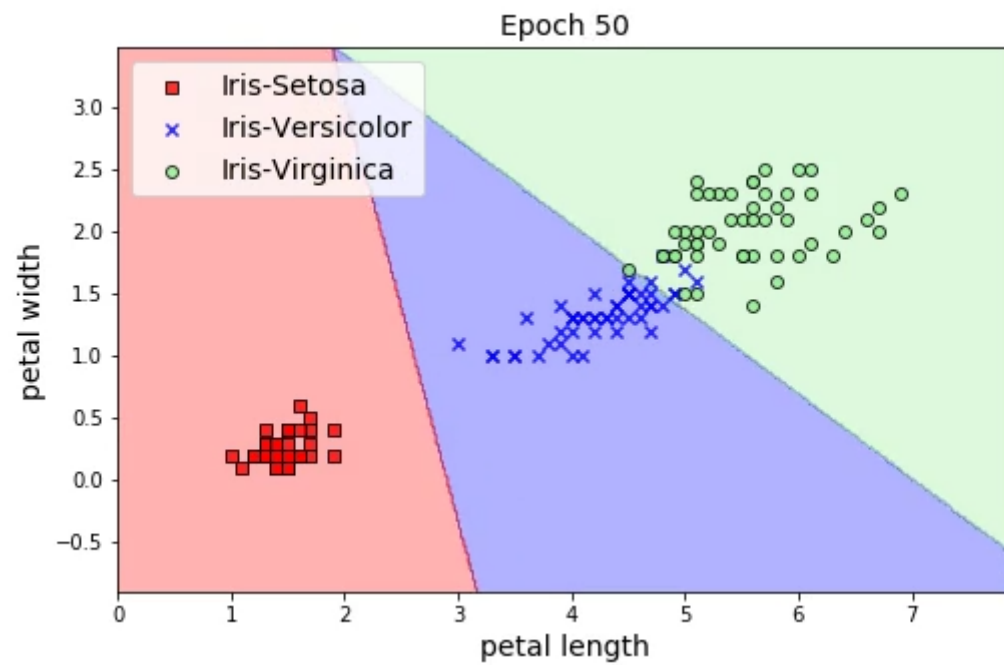
```
1 class SoftmaxRegression():
2     def __init__(self, inputs, outputs):
3         self.w = np.random.randn(inputs, outputs)
4
5     def fit(self, X, Y, epochs, lr):
6         X = np.c_[np.ones(len(X)), X]
7         for epoch in range(epochs):
8             for x, y in zip(X, Y):
9                 x = x[None, :]
10                y_hat = np.dot(x, self.w)
11                loss = crossentropy(y_hat, y).mean()
12                dldo = grad_crossentropy(y_hat, y)
13                grad_w = np.dot(x.T, dldo)
14                self.w = self.w - lr * grad_w
15
16    def predict(self, x):
17        x = np.c_[np.ones(len(x)), x]
18        y_hat = np.dot(x, self.w)
19        return np.argmax(softmax(y_hat), axis=1)
```

```
1 def crossentropy(output, target):
2     logits = output[np.arange(len(output)), target]
3     entropy = - logits + np.log(np.sum(np.exp(output), axis=-1))
4     return entropy
```

```
1 def softmax(x):
2     return np.exp(x) / np.exp(x).sum(axis=-1, keepdims=True)
```

```
1 def grad_crossentropy(output, target):
2     answers = np.zeros_like(output)
3     answers[np.arange(len(output)), target] = 1
4     return (- answers + softmax(output)) / output.shape[0]
```





# SUMMARY

- ANNs are inspired in biological NNs
- The perceptron computes a weighted sum of its inputs and applies an activation function
- We can make different tasks depending on the activation and loss functions used
- We use Gradient Descent to find the set of parameters that minimize the loss function
- Linear models have limitations with data that is not linearly separable

Task	Activation	Loss
Regression	Linear	MSE
Binary Classification	Step	MSE
	Sigmoid	LogLoss
Multi-class Classification	Softmax	CE
Multi-label Classification	Sigmoid	LogLoss

See the code



Practice with this exercise