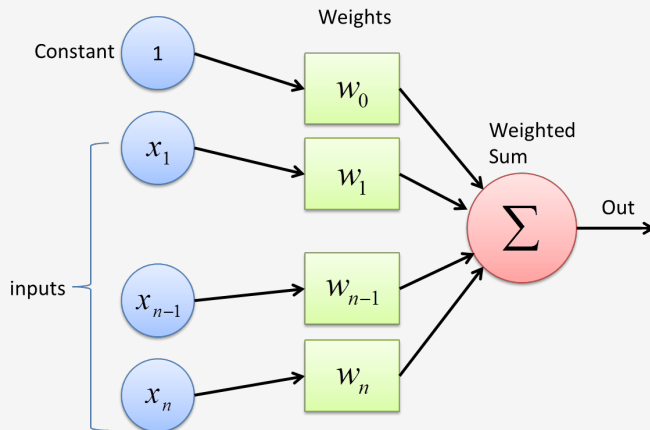# Supervised learning

Perceptron

Section 1

The perceptron (regression)

# Where we see a first neuron



$$h_w(x) = w_0 + w_1 \times x_1 + w_2 \times x_2 + \cdots + w_n \times x_n$$

## Characterisation

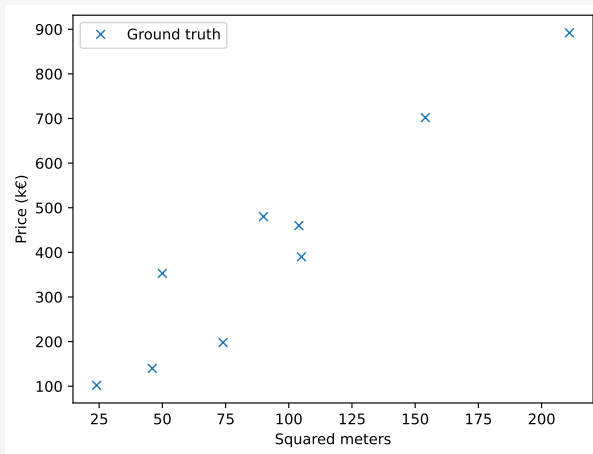The elements of the **w** vector are called the weights of the perceptron.

Given an input vector $x$, the output of the perceptron is a linear combination of its input.

$$h_w(x) = w_0 + w_1 \times x_1 + w_2 \times x_2 + \cdots + w_n \times x_n$$

The set of functions representable by a perceptron is the set of linear functions. That's our hypothesis space $\mathcal{H}_{perc}$.

# An even simpler dataset

| X | Y |
|-------|-----------|
| $m^2$ | Price (€) |
| 24 | 102 000 |
| 46 | 140 000 |
| 50 | 353 600 |
| 211 | 892 000 |
| 74 | 198 000 |

Perceptron for predicting the price

We have a single feature ($x_1$: squared meters), thus a function representable by a perceptron would have the form:

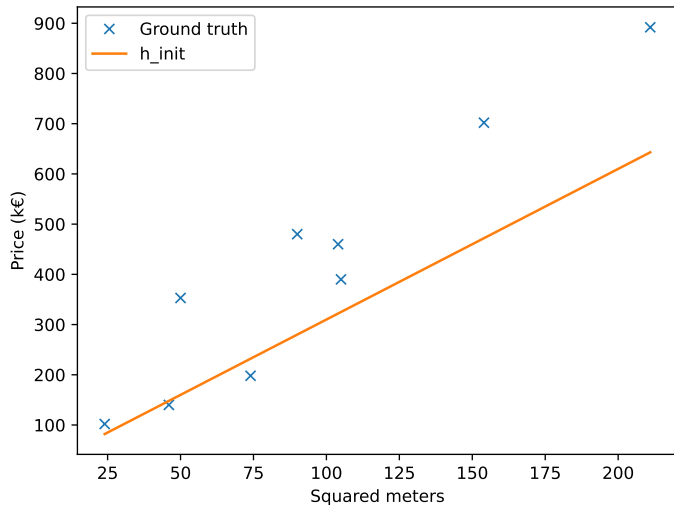$$h_w(sqm) = w_0 + w_1 \times x_1$$

where we can interpret:

- $w_0$ as the base price
- $w_1$ as the price per squared meters

# Perceptron for predicting the price

Making an educated guess we could set:

$$w_0 = 10000 \quad (\texteuro)$$
$$w_1 = 3000 \quad (\texteuro/m^2)$$
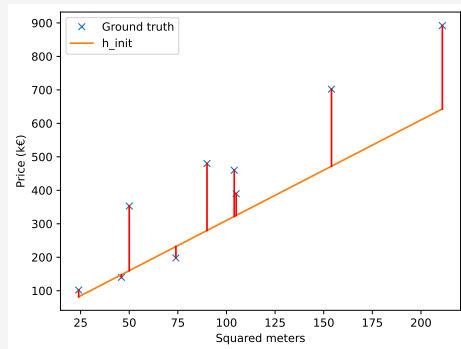
# Perceptron: prediction error

For an example $(x, y)$ and predictor $h_w$

prediction error $= |y - h_w(x)|$ (length of red segments)

$L_2$ loss: $L_2(y, \hat{y}) = (y - h_w(x))^2$ (squared error)

This leads to the empirical loss (for $L_2$) over the entire dataset $E$:

$$EmpLoss_{L_2, E}(h_w) = \sum_{(x,y) \in E} \frac{(y - h_w(x))^2}{|E|}$$

# Finding the best hypothesis

I want the hypothesis $\hat{h}^*$, with the minimum empirical loss:

$$\hat{h}^* = \underset{h_w \in \mathcal{H}_{perc}}{\arg\min} \, EmpLoss_{L_2, E}(h_w)$$

For the perceptron, this means finding the best weights $\hat{w}^*$ in the weight space.

$$\hat{w}^* = \underset{w}{\arg\min} \, EmpLoss_{L_2, E}(h_w)$$

Posing $Loss(w) = EmpLoss_{L_2, E}(h_w)$, we obtain:

$$\hat{w}^* = \underset{w}{\arg\min} \, Loss(w)$$

# Gradient: direction of steepest ascent

In any point $w$ of the function, the gradient defines the direction of steepest ascent:

$$\vec{\nabla} g(w)$$

It can be computed from the partial derivatives:

$$\vec{\nabla} g(w) = \begin{bmatrix} \frac{\delta}{\delta w_0} g(w) \\ \frac{\delta}{\delta w_1} g(w) \\ \vdots \\ \frac{\delta}{\delta w_m} g(w) \end{bmatrix}$$
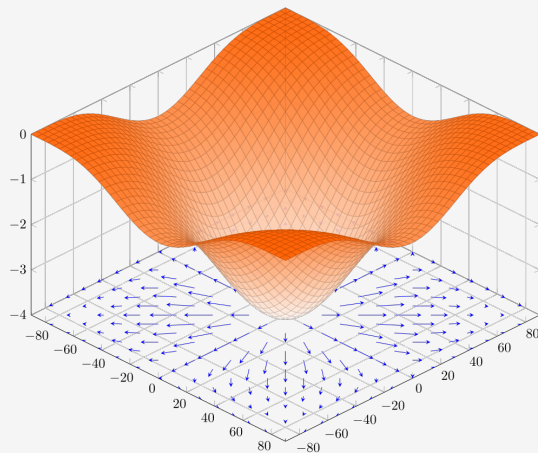


Figure: Gradient of $f(x, y) = -(cos^2(x) + cos^2(y))^2$

## Gradient descent

From a point $w$, compute a new candidate $w'$ by following the direction of steepest descent (opposite of the gradient).

$$w' = w - \alpha \times \vec{\nabla} g(w)$$

The distance traveled is parameterized by the **step size** $\alpha$.

Since the function is decreasing in this direction, there is a good chance that

$$g(w') < g(w)$$

## Gradient descent

Applying this repeatedly, we get the gradient descent algorithm:[1]

$w \leftarrow$ any value in the parameter space
**while** not converged **do**
$\quad w \leftarrow w - \alpha \times \vec{\nabla} Loss(w)$

Typical convergence criteria: stop when the update did not provide an improvement for the last $k$ iterations (e.g. $k = 5$).

---

[1]Recal that $Loss(w)$ is shortcut for $EmpLoss_{L,E}(h_w)$

Computing the gradient ($L_2$ loss, single example)

Partial derivative of the $L_2$ loss for a single example $(x, y)$:

$$\frac{\delta}{\delta w_i} Loss(w) = \frac{\delta}{\delta w_i}(y - h_w(x))^2$$
$$= 2(y - h_w(x)) \times \frac{\delta}{\delta w_i}(y - h_w(x))$$

Applied to our system with a single feature $(x_1)$ we obtain:

$$\frac{\delta}{\delta w_0} Loss(w) = -2(y - h_w(x))$$
$$\frac{\delta}{\delta w_1} Loss(w) = -2(y - h_w(x)) \times x_1$$

## Update rules

Updating the weights based on a single example:[2]

$$w_0 \leftarrow w_0 + \alpha \times (y - h_w(x))$$
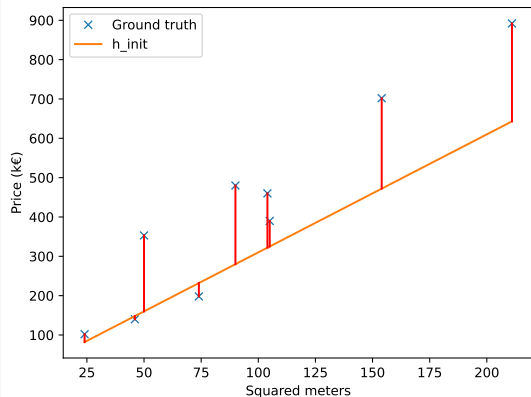$$w_1 \leftarrow w_1 + \alpha \times (y - h_w(x)) \times x_1$$

Updating the weights based on the entire training set $E$:

$$w_0 \leftarrow w_0 + \alpha \times \sum_{(x,y) \in E} (y - h_w(x))$$
$$w_1 \leftarrow w_1 + \alpha \times \sum_{(x,y) \in E} (y - h_w(x)) \times x_1$$

---

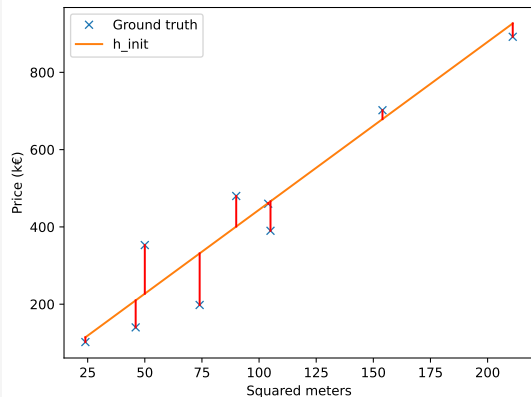[2]Note that the $-2$ factor from the previous equation is included in $\alpha$ term.

# Updating our initial guess[3]



$$w = [10000, 3000]$$

$$w' = [10010.53, 4343.82]$$
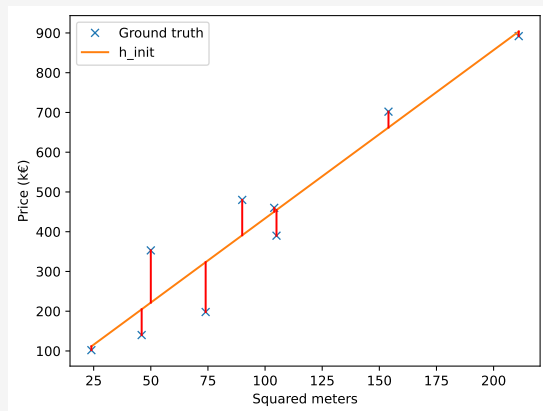
[3]With $alpha = 10^{-5}$

# Result of the gradient descent

Repeating the gradient descent step, we eventually converge to our best solution.

$$\hat{w}^* = [9947.29, 4235.02]$$

I should sell my apartment for:

$$9947.29 + 4235.02 \times 165 = 708725 €$$

Gradient descent (applied to a perceptron with $L_2$ loss)

$w \leftarrow$ any value in the parameter space
**while** not converged **do**
    $w_0 \leftarrow w_0 + \alpha \times \sum_{(x,y) \in E}(y - h_w(x))$
    **for** $i \in 1 \ldots n$ **do**
        $w_i \leftarrow w_i + \alpha \times \sum_{(x,y) \in E}(y - h_w(x)) \times x_i$

## Representation trick

In the previous slides, we always had to deal with the $w_0$ weight specially because it has no corresponding feature.

We can define an artificial feature $x_0$ that always has the value 1. And reformulate $h_w$:

$$h_w(x) = \sum_{i \in [0,m]} w_i \times x_i$$

$$h_w(x) = w \cdot x \quad \text{(dot product)}$$

and the update rule (for $L_2$ loss):

$$w_i \leftarrow w_i + \alpha \times \sum_{(x,y) \in E} (y - h_w(x)) \times x_i$$

Section 2

A perceptron for classification

## A classification problem

I now want to buy a new apartment to replace the one I just sold. To be reactive I built an automated system that sends me any new announce of an apartment for sale.
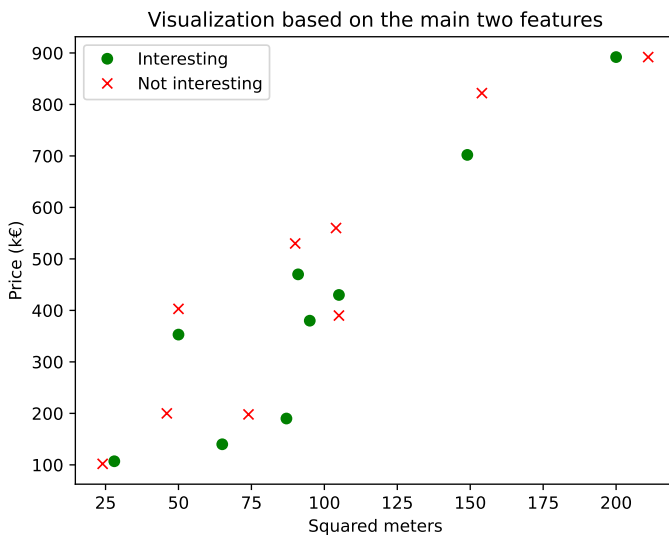
- Problem: there are dozens of announces every day and I don't have time to look at them all.

- Solution: build an AI system that will predict whether I will be interested in a particular apartment based on a few of its features. If it predicts that I am not interested, it will discard the announce.
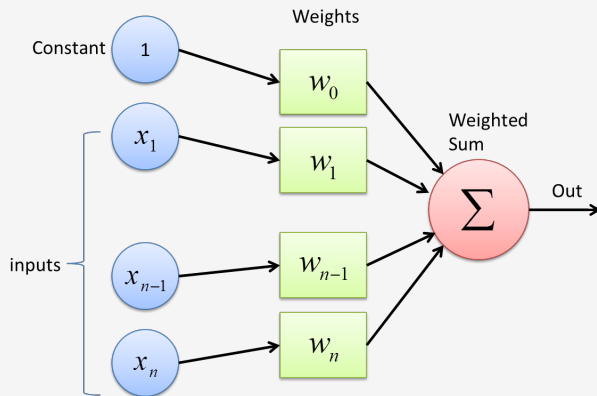
A classification problem: dataset

So far, I collect the following information stating whether an announce that I previously saw was interesting.

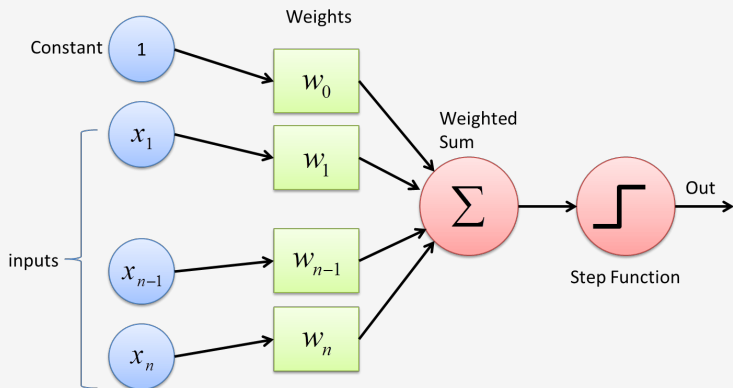| X | | | | Y |
|---|---|---|---|---|
| $m^2$ | Num Rooms | Floor | Price ($\in$) | Interesting |
| 24 | 1 | 4 | 102 000 | true |
| 46 | 3 | 2 | 140 000 | false |
| 50 | 3 | 6 | 353 600 | false |
| 211 | 5 | 3 | 892 000 | true |
| 74 | 3 | 1 | 198 000 | true |

# A classification problem: dataset visualization

# Our previous perceptron



$$h_w(x) = w_0 + w_1 \times x_1 + w_2 \times x_2 + \cdots + w_n \times x_n$$

$$h_w(x) = w \cdot x$$

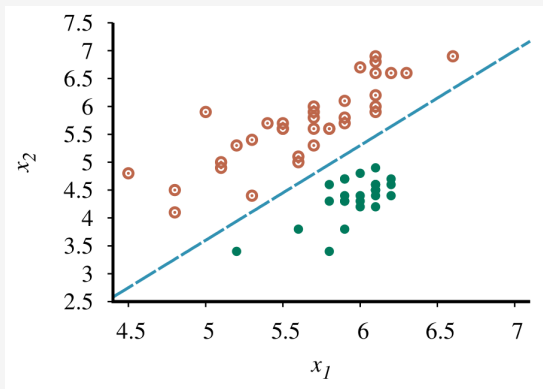# Perceptron for classification



$$h_w(x) = Step(w \cdot x)$$

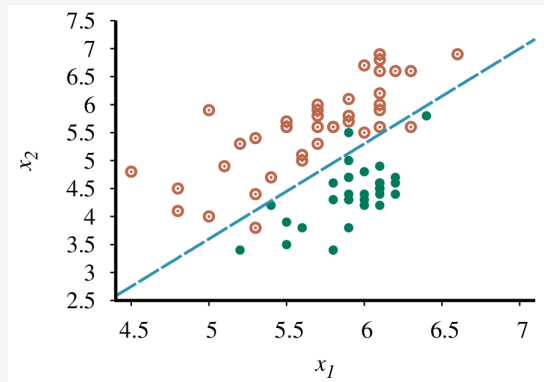where $Step(z) = 1$ if $z \geq 0$ and $0$ otherwise

# The perceptron as a linear classifier

The perceptron defines a **decision boundary** that separates two classes.



Linearly separable
Perfectly classifiable by a perceptron

**Not** linearly separable
**Not** perfectly classifiable by a perceptron
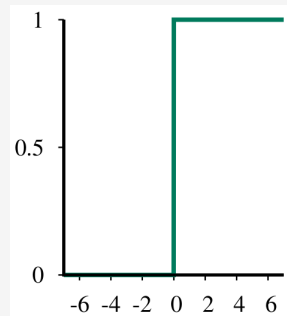
## The step function

$$h_w(x) = Step(w \cdot x)$$

where $Step(z) = 1$ if $z \geq 0$ and $0$ otherwise

We now have a function that we could train in order to output:

- 1 if the example is in the class (interesting)
- 0 otherwise (not interesting)

Problem: the function is:

- non-differentiable in 0
- the gradient is 0 everywhere else



$Step(z)$

## The perceptron learning rule

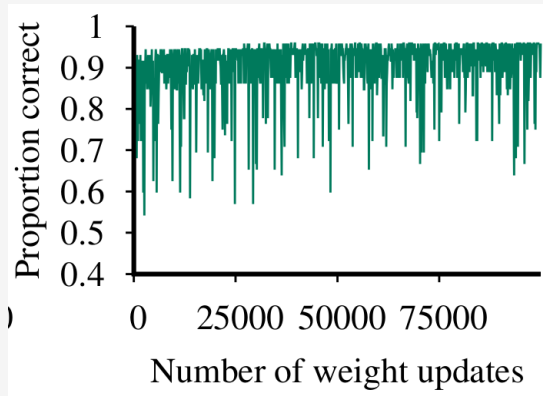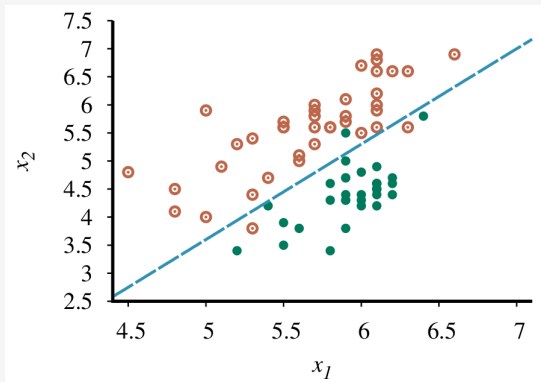Nevertheless, an rule was proposed the **perceptron update rule** (here for a single example $(x, y)$):

$$w_i \leftarrow w_i + \alpha \times (y - h_w(x)) \times x_i$$

which is identical to the update rule for linear regression (for $L_2$).

The rule is show to converge to a solution when the data is linearly separable.

# The perceptron learning rule (under non separable data)

However the perceptron learning rule is unstable when the data is not linearly separable:
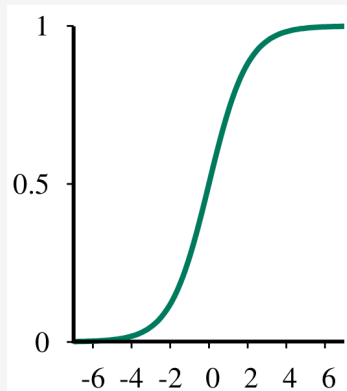
## Replacing the step function

Turns out we can replace the step function with one with nicer properties.

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

and redefine our hypothesis function:

$$h_w(x) = Logistic(w \cdot x) = \frac{1}{1 + e^{-w \cdot x}}$$

Often called the **logistic regression**.

# Back on track: gradient descent

$Logistic(x)$ is differentiable on $]-\infty, \infty[$. This allows us to reuse gradient descent for training:
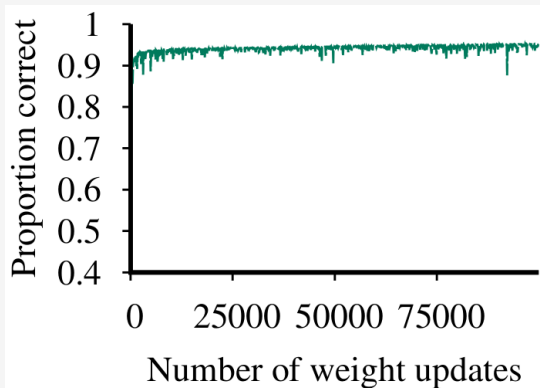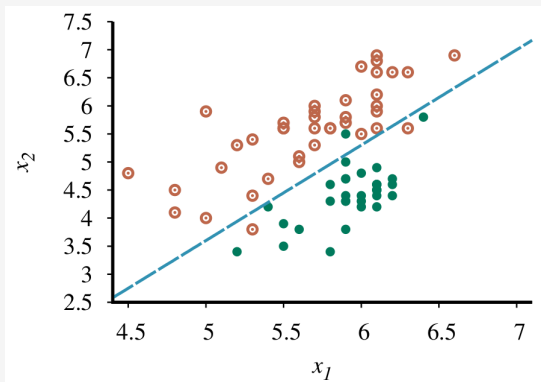
$$w \leftarrow w - \alpha \times \vec{\nabla} Loss(w)$$

For an $L_2$ loss we obtain the update rule:

$$w_i \leftarrow w_i + \alpha(y - h_w(x)) \times h_w(x) \times (1 - h_w(x)) \times x_i$$

# Training the logistic regression (under non-separable data)

Compared to the step function, the logistic regression tends to converge more quickly and reliably in the presence of noisy and non-separable data.
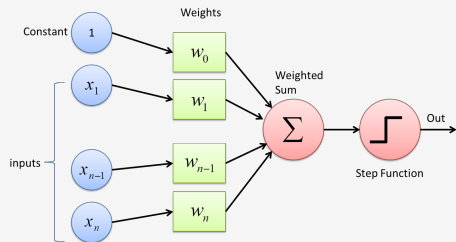
Section 3

Synthesis

## Synthesis

We saw two classes of perceptrons:
- linear regressor
- linear classifier

Both can be trained with **gradient descent** in attempt to **minimize the loss**.

In the next course, the perceptron will be a **neural unit** in a **neural network**.

Exercises

- For each of the following datasets, propose a perceptron that would correctly classify it.

| X | Y |
|----|--------|
| -1 | Flower |
| 2  | Cat    |
| 1  | Flower |
| 0  | Flower |

| X | Y |
|----|--------|
| -1 | Flower |
| 2  | Cat    |
| 1  | Cat    |
| 0  | Cat    |

| X | Y |
|----|--------|
| -1 | Flower |
| 2  | Flower |
| 1  | Cat    |
| 0  | Cat    |

## Exercises

1. You have $N$ examples in your dataset, each with $M$ features. Give an estimate of the computational cost of a single update step. Does it scale to large-scale datasets (e.g. $N = 10^5, M = 10^4$)

2. For the linear regression (regression perceptron), are we guaranteed to find the optimal weights with gradient descent?

3. What's a reasonable loss for spam detection? (you can make some assumptions about the proportion of spam)

4. What's the update formula of a regression perceptron using the $L_1$ loss?