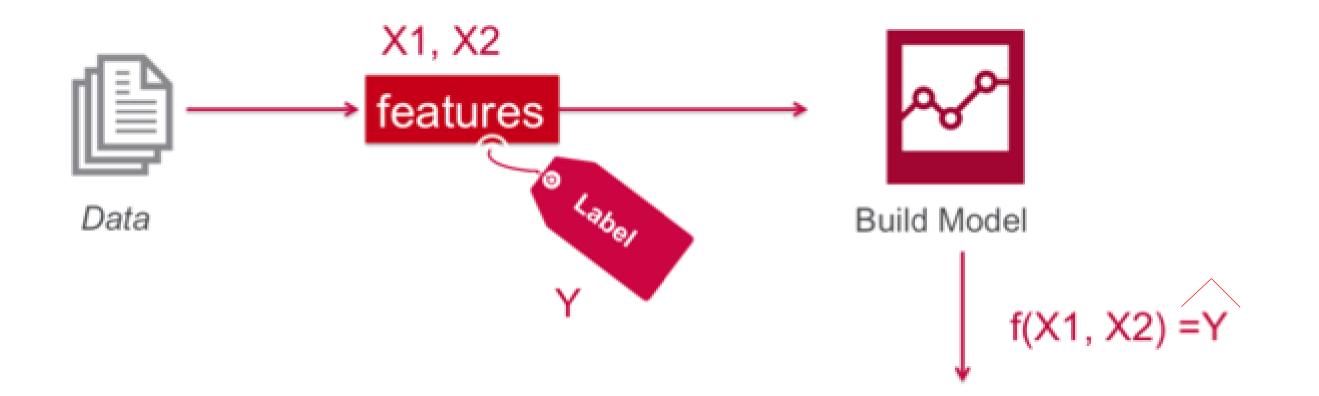
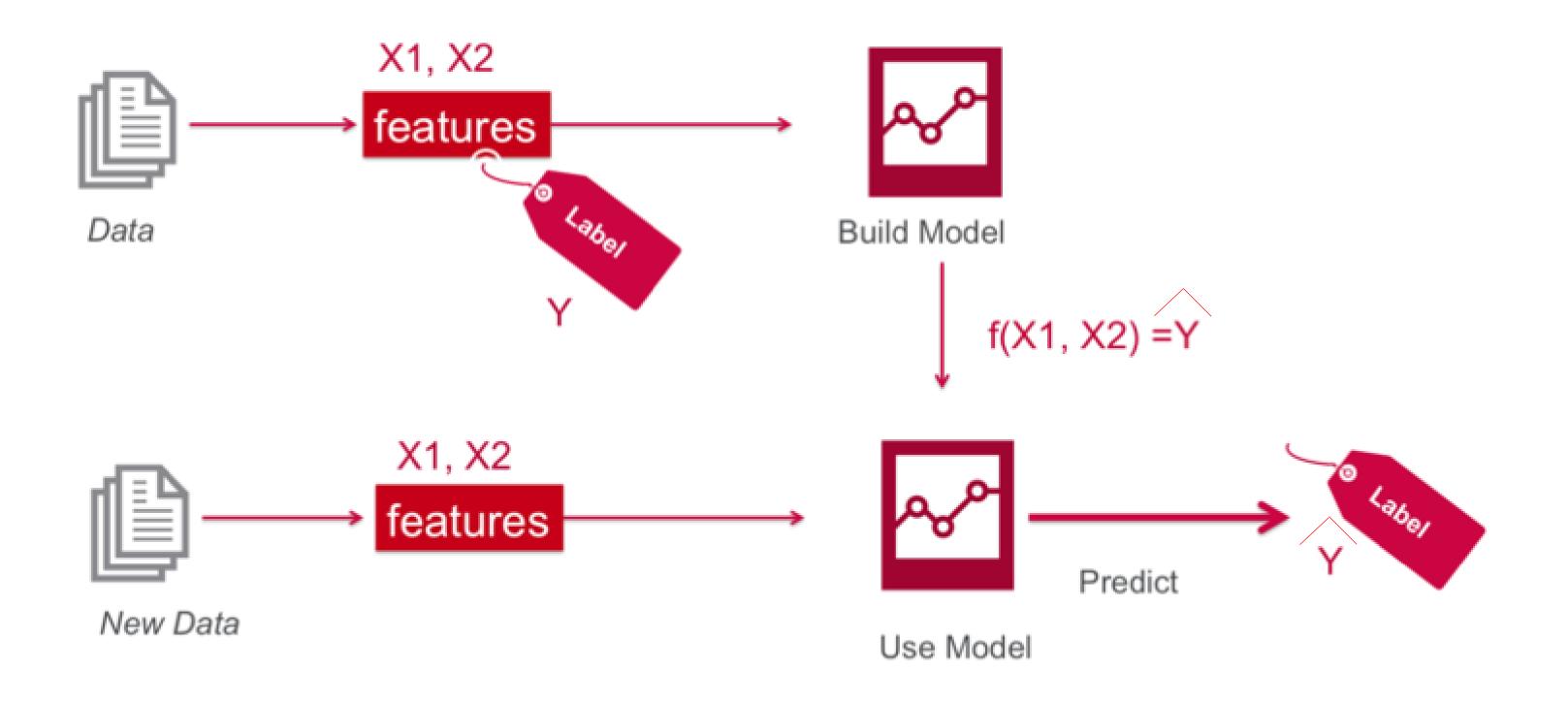




Supervised learning



Supervised learning

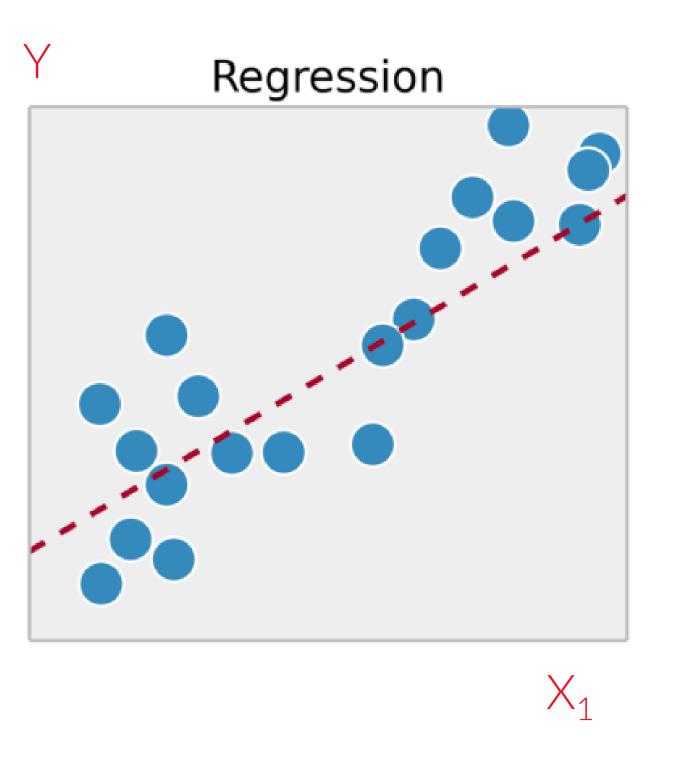




Classification vs Regression

Classification: predict a discrete (or categorical) value based on the input variables

Regression: predict a continuous value based on the input variables

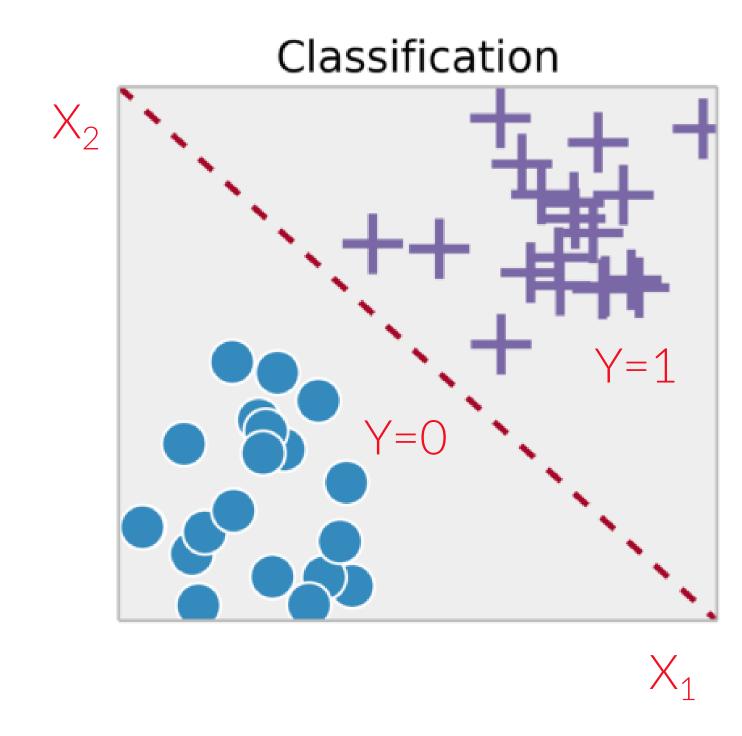




Classification vs Regression

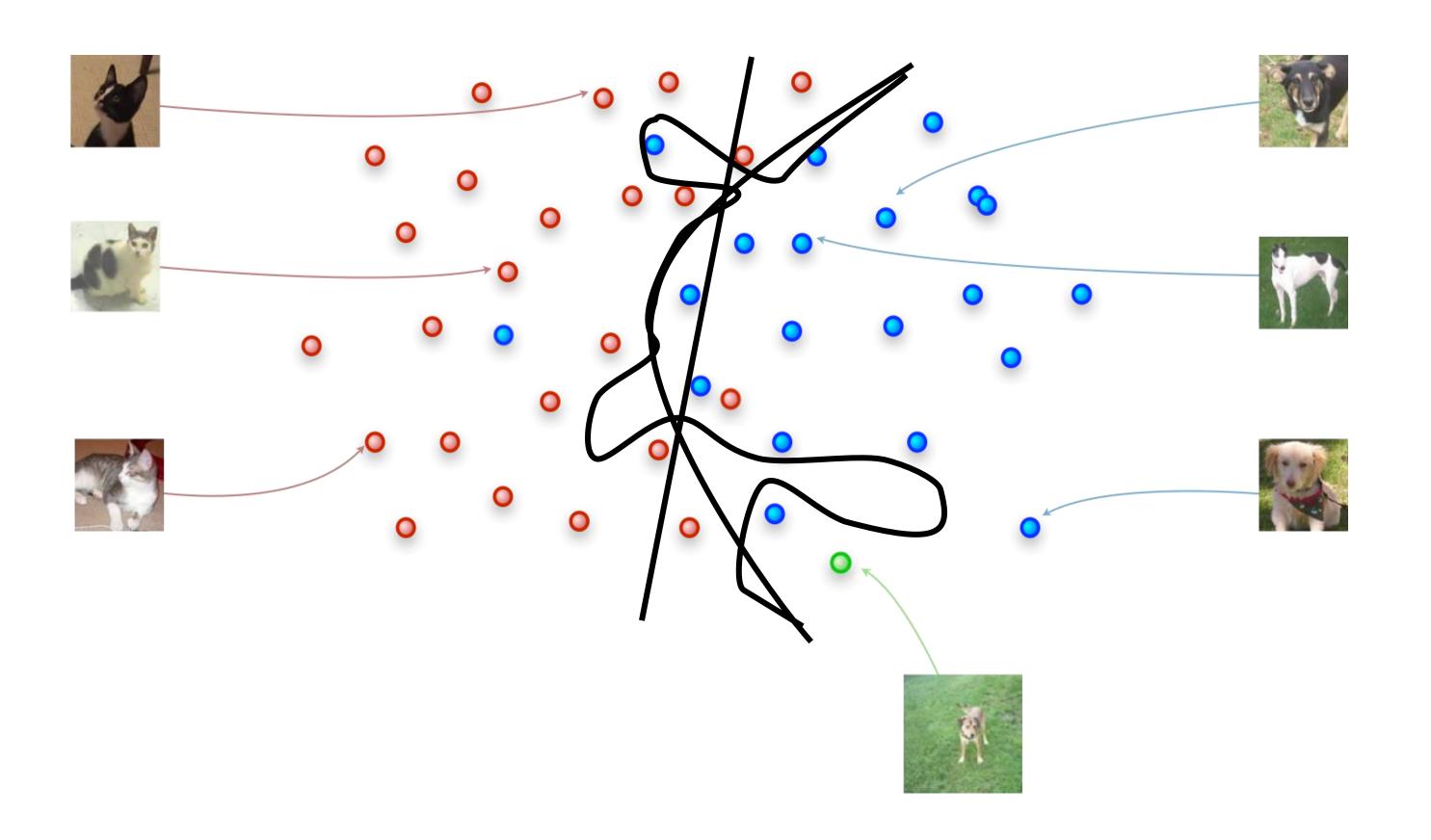
Classification: predict a discrete (or categorical) value based on the input variables

Regression: predict a continuous value based on the input variables





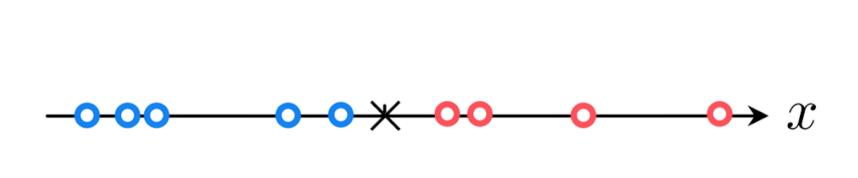
Binary classification



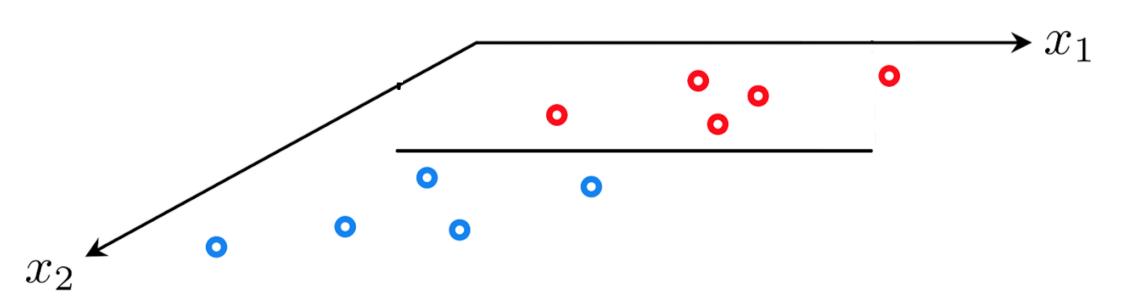




Two perspectives on classification



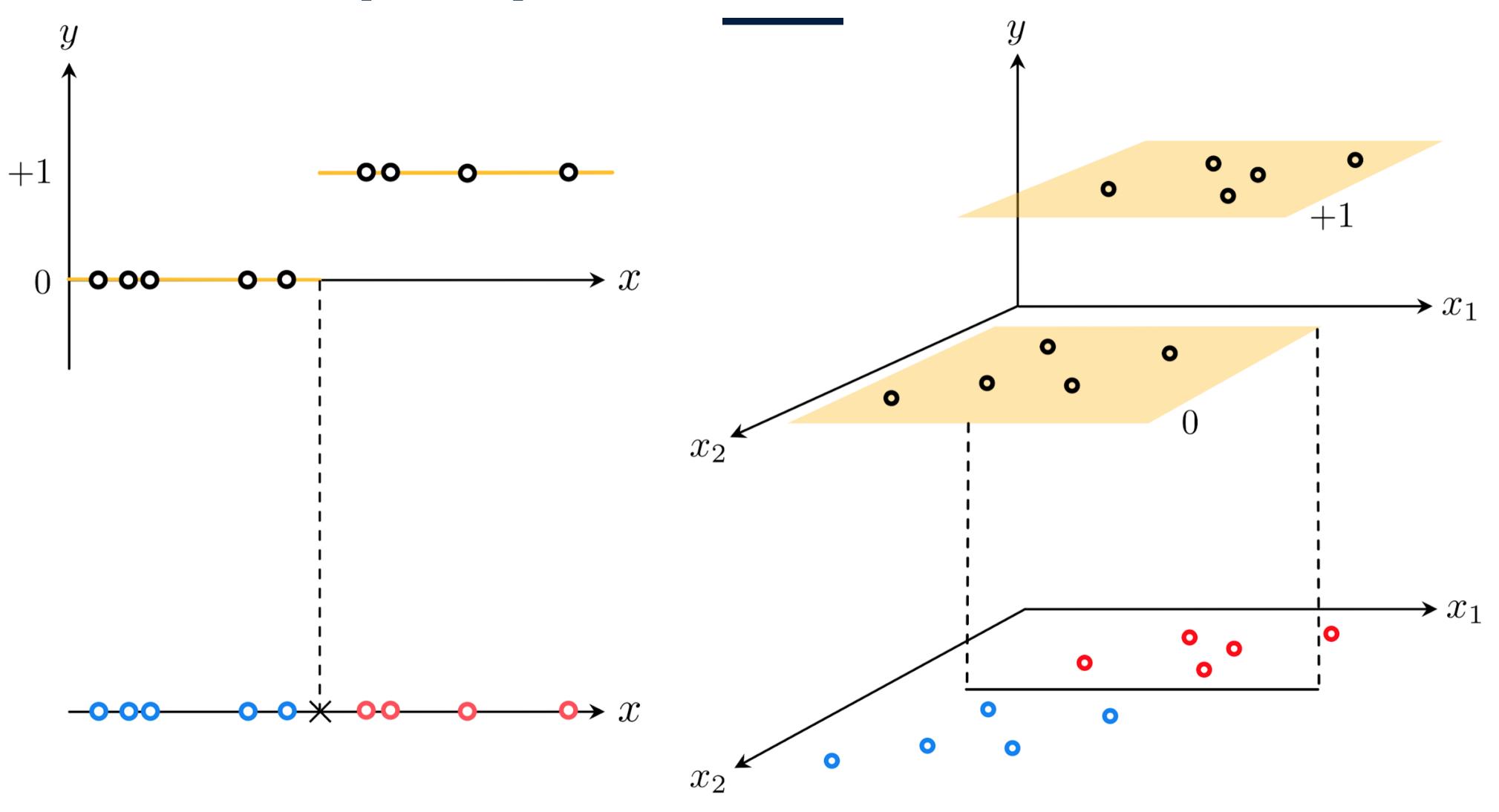
1d input: decision boundary is a point



2d input: decision boundary is a line



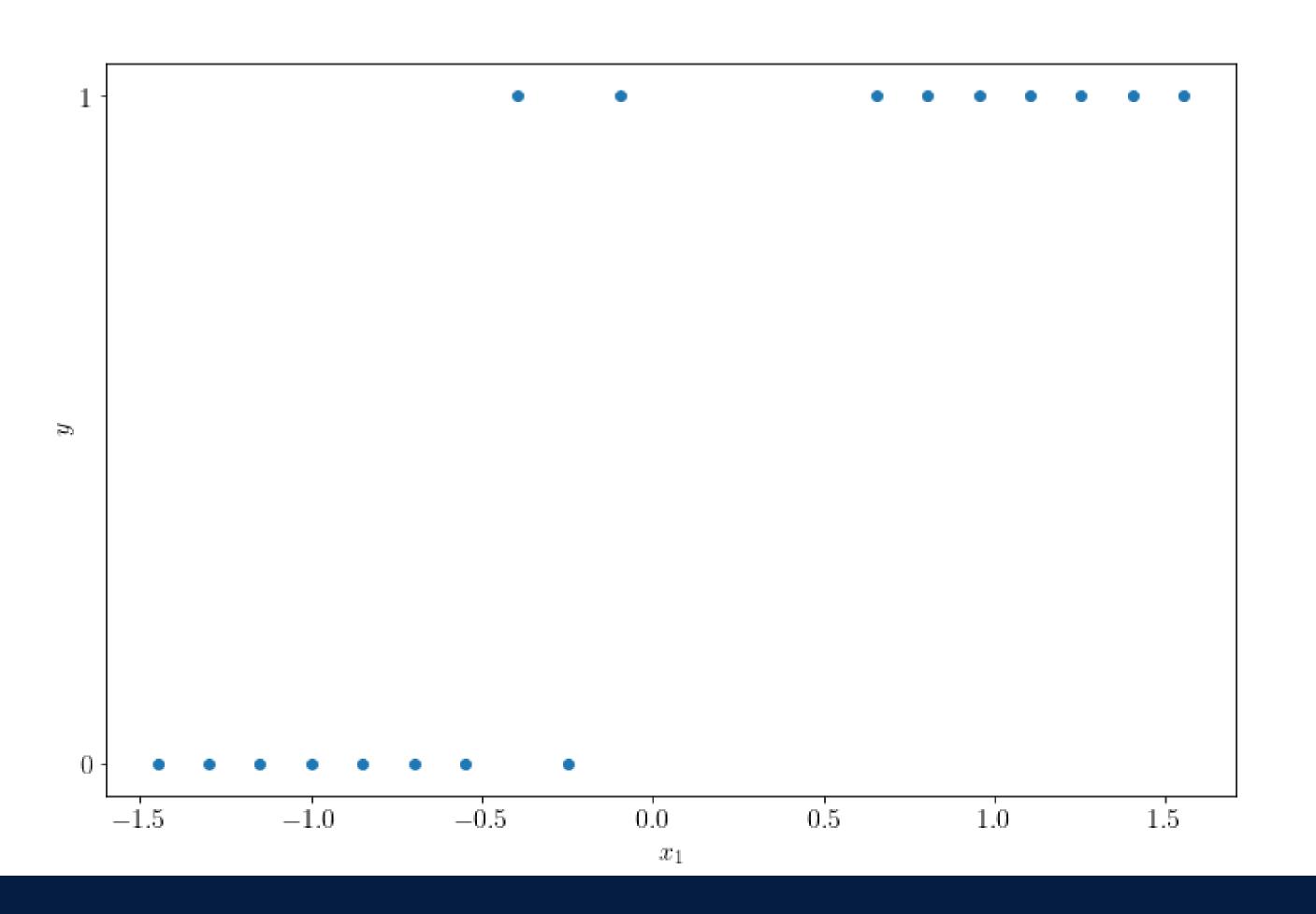
Two perspectives on classification



1d input: decision boundary is a point

2d input: decision boundary is a line

Illustrative 1D example

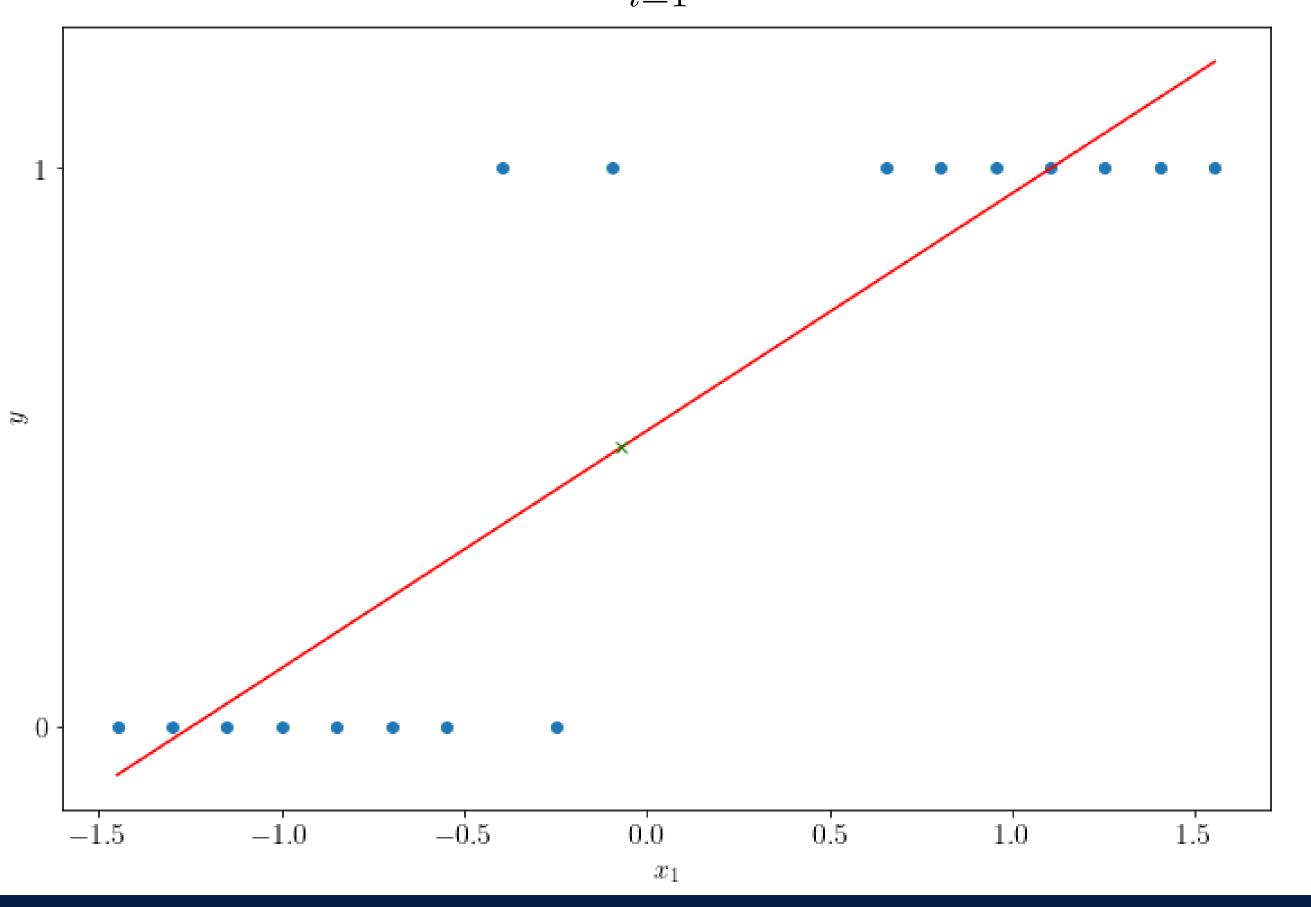






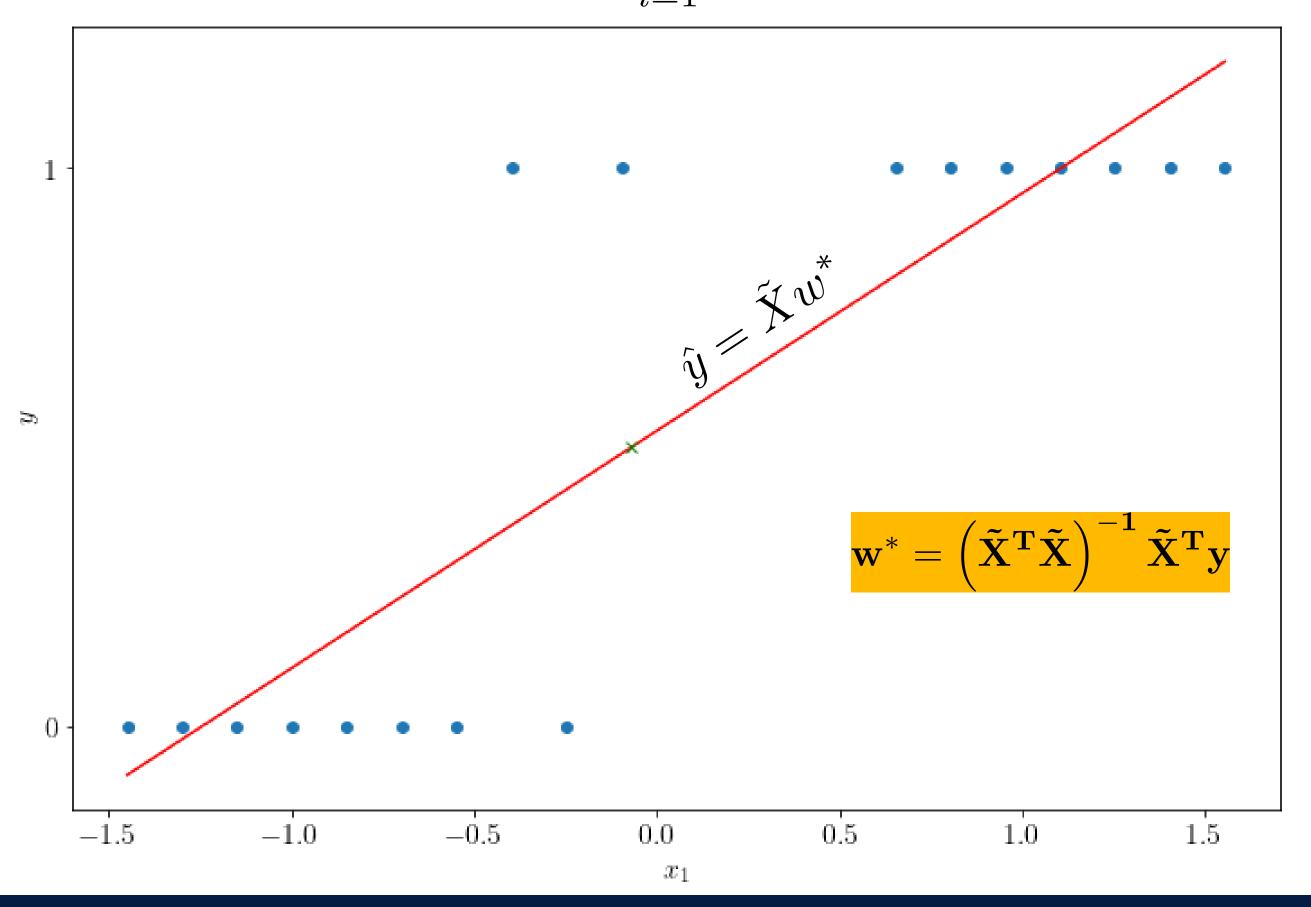
How about applying the linear regression framework seen earlier to fit such data?

We minimize the following cost function $g(w) = \frac{1}{N} \sum_{i=1}^{N} \left(\tilde{X}^{i} w - y^{i} \right)^{2}$

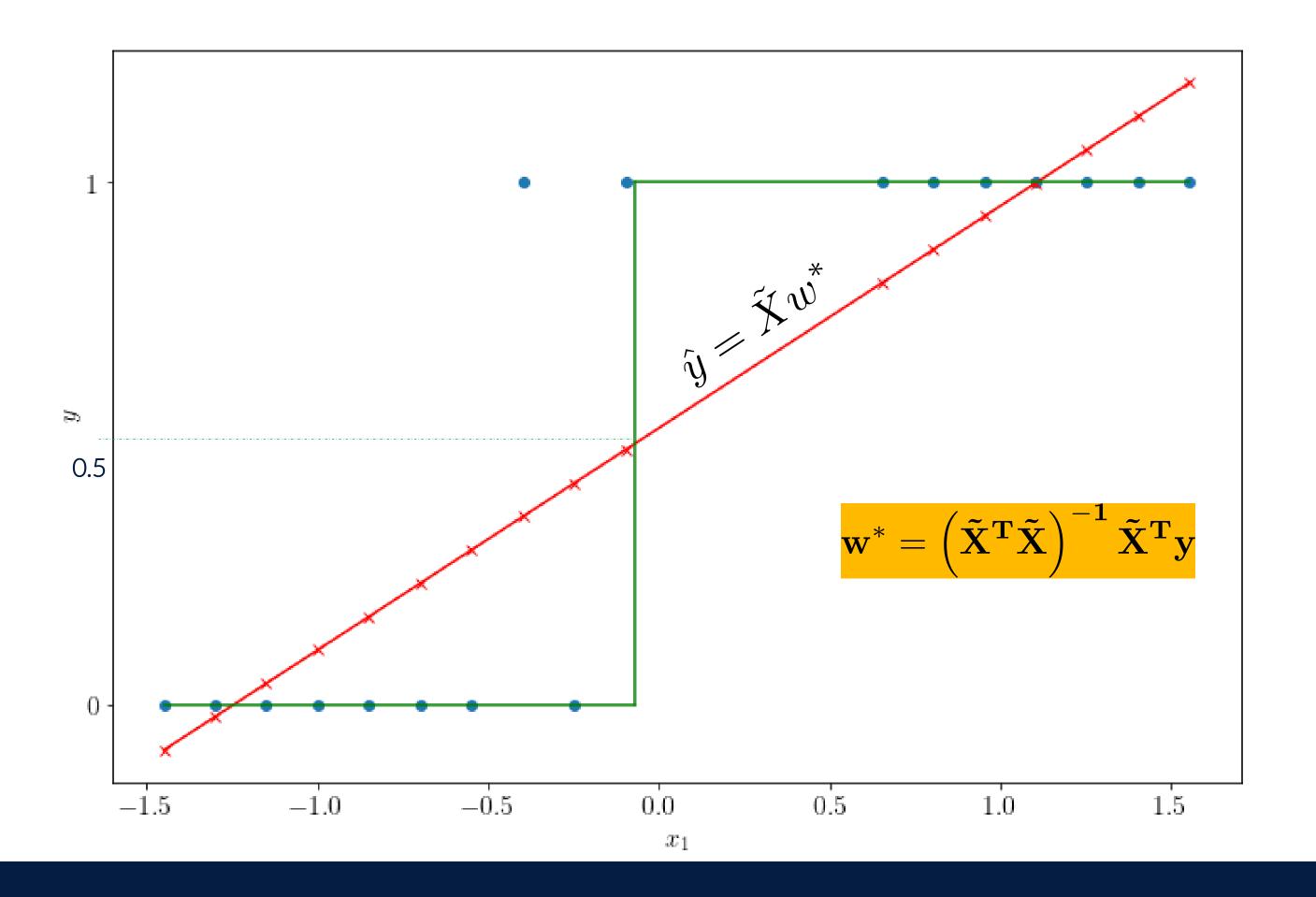


How about applying the linear regression framework seen earlier to fit such data?

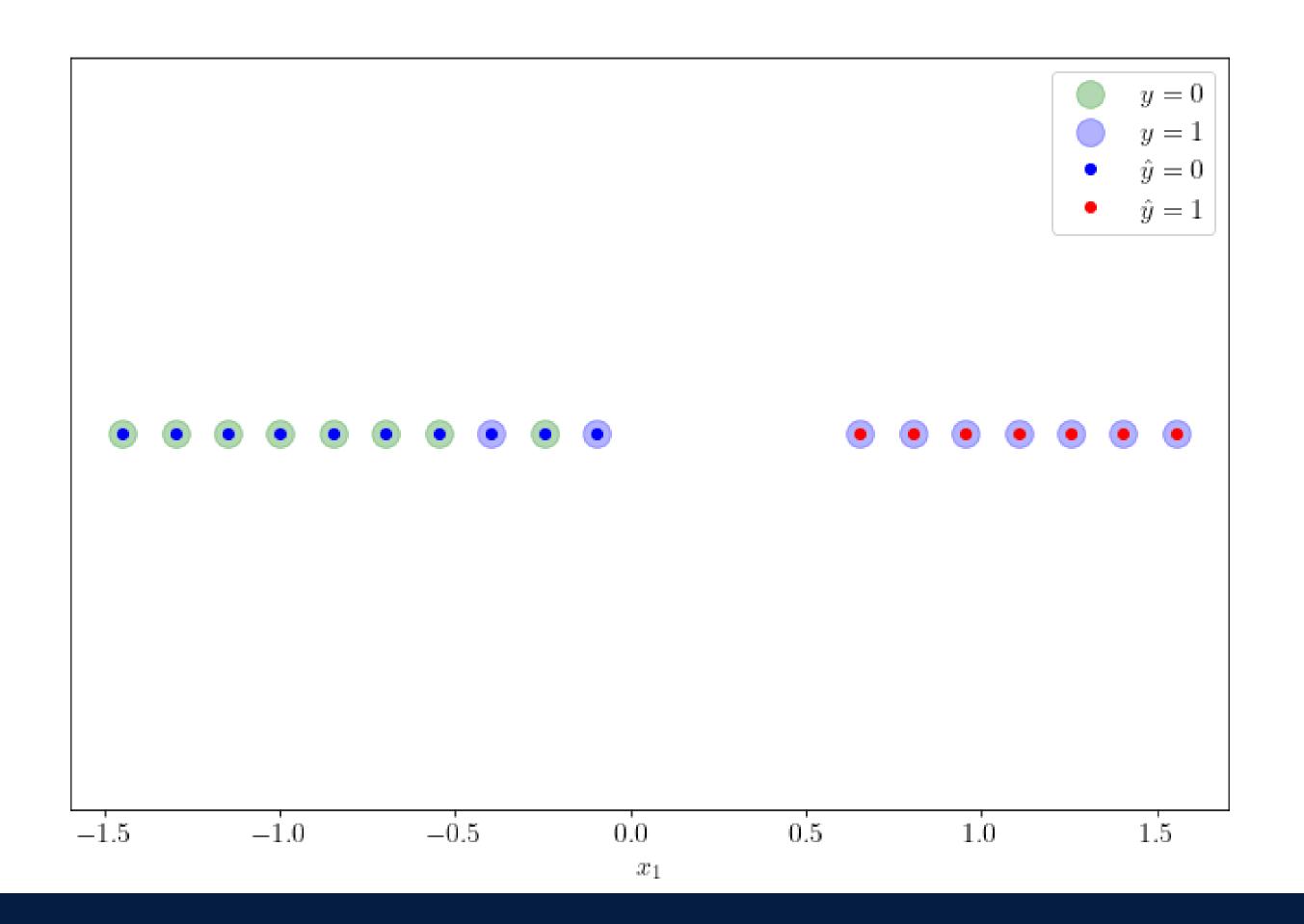
We minimize the following cost function $g(w) = \frac{1}{N} \sum_{i=1}^{N} \left(\tilde{X}^{i} w - y^{i} \right)^{2}$



And we apply a step function $\operatorname{step}(\tilde{X}w^*-0.5)$



Results from the other perspective: 2 classification errors. In this case, 2 false negatives.

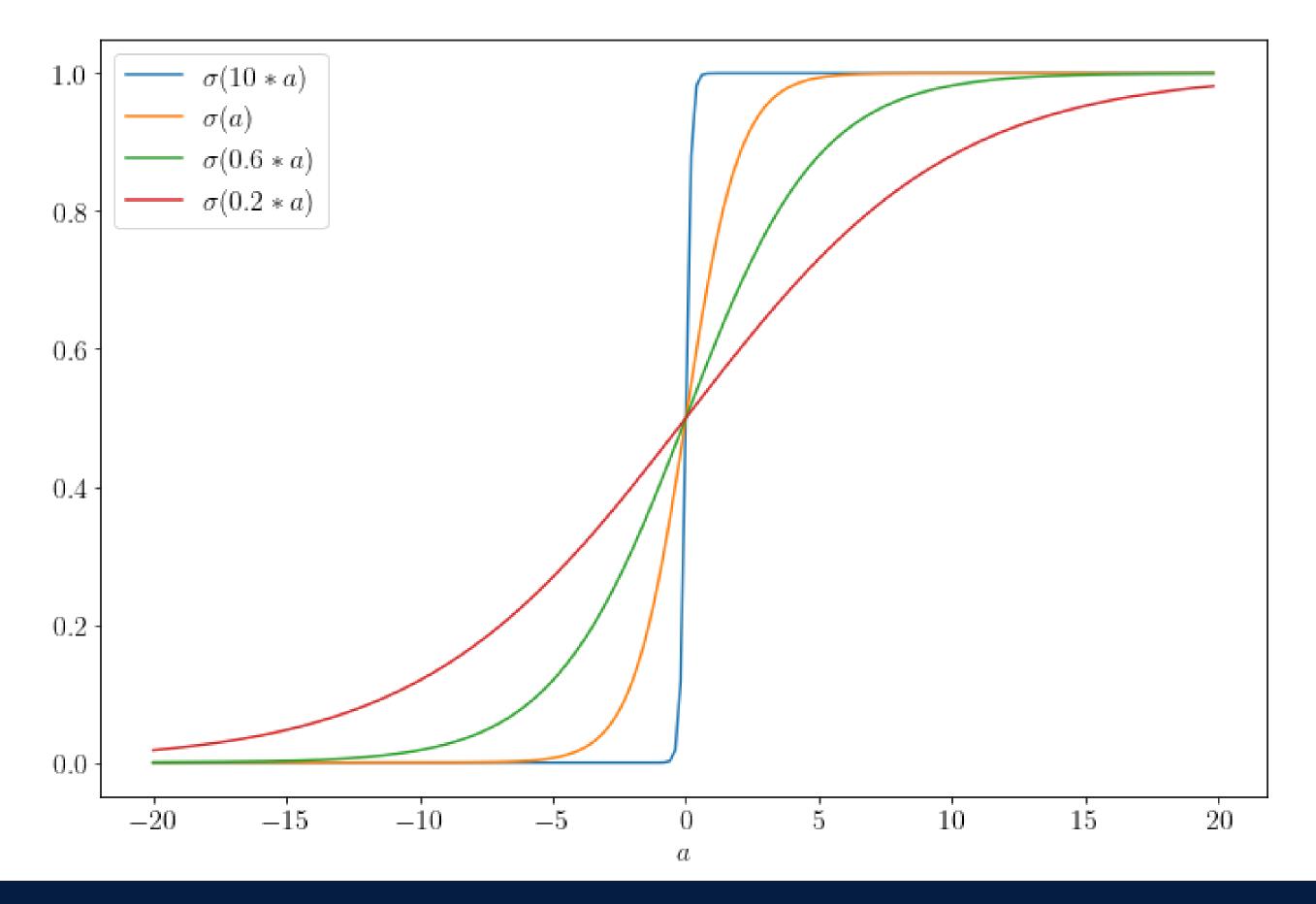




Proposition #2 – Intuition

How about integrating the step function in the cost function to minimize?

Or, better yet, a smooth approximation of the step function, the logistic sigmoï $\phi(x) = \frac{1}{1 + e^{-x}}$



Proposition #2 – Formally

We'll explicitly model the probability $p := P(\hat{y} = 1|X)$

The probability takes values in [0, 1]. The range of a linear model is not bounded.

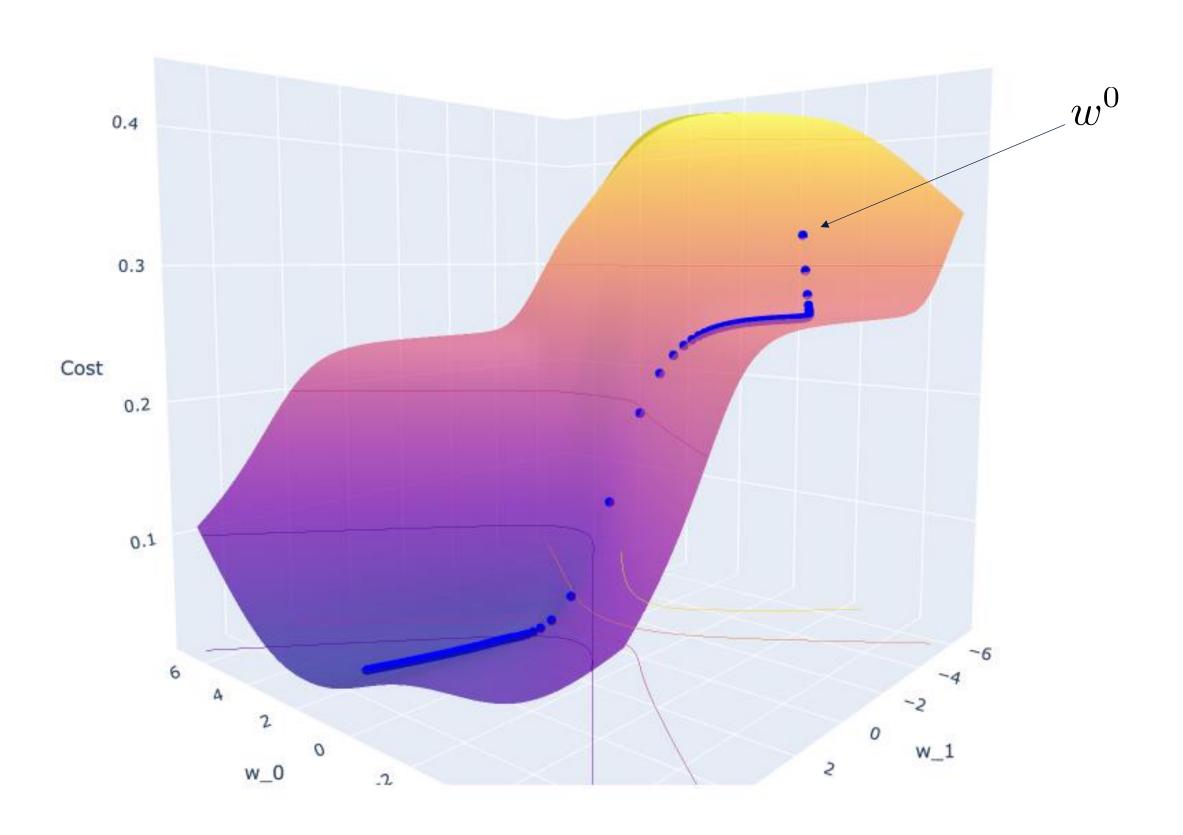
We'll extend the range to
$$[-\infty, +\infty]$$
 , and approxima $\left(e\frac{p}{1-p}\right)$ by a linear model Hence, $\ln\left(\frac{p}{1-p}\right) = \tilde{X}w$, which can be written as:

Hence,
$$\ln\left(rac{p}{1-p}
ight)= ilde{X}w$$
 , which can be written as:

$$p := P(\hat{y} = 1|X) = \frac{e^{\tilde{X}w}}{1 + e^{\tilde{X}w}} = \frac{1}{1 + e^{-\tilde{X}w}} = \sigma(\tilde{X}w)$$

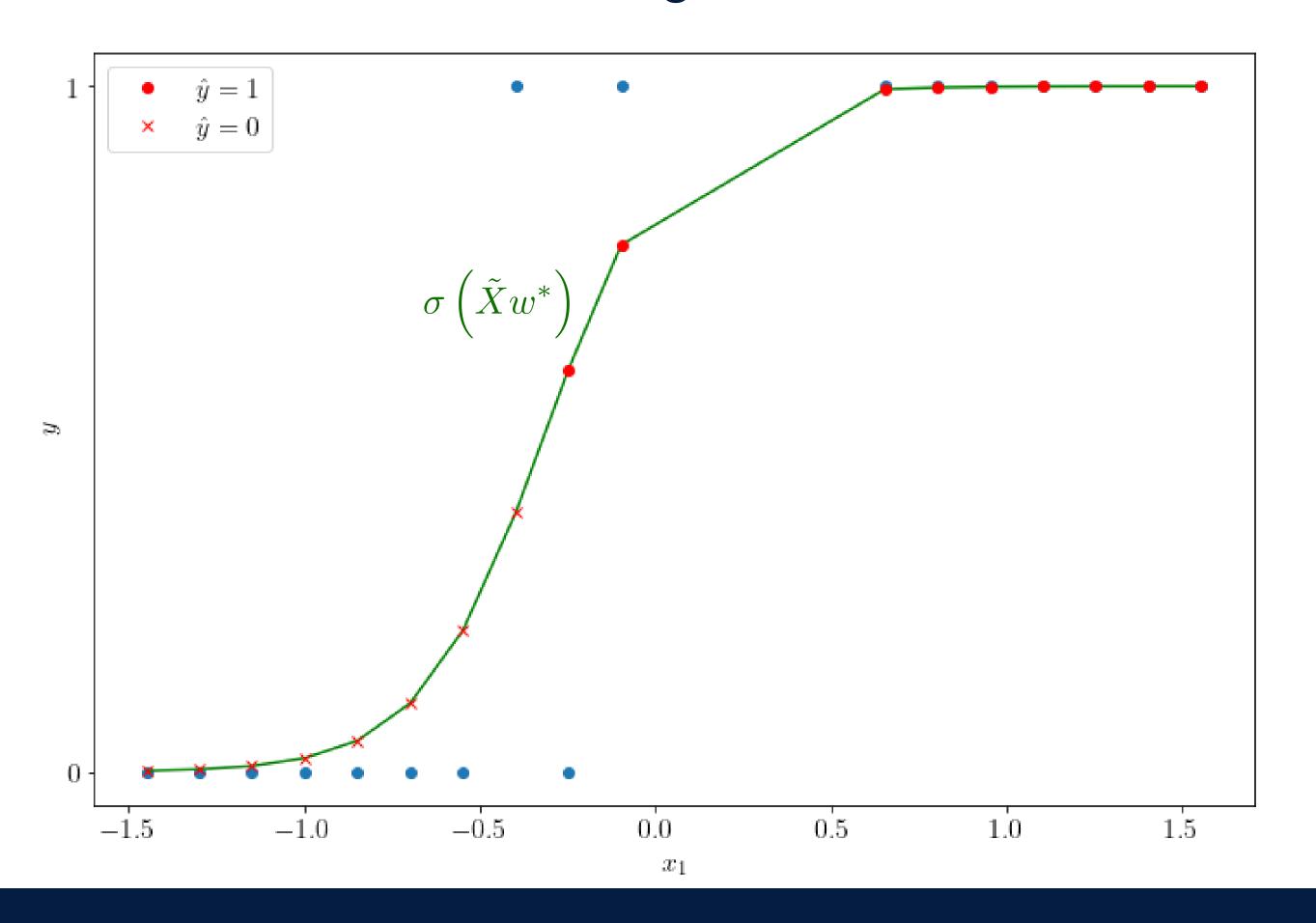
So we'll now minimize the following cost function $g_{\rm sl}(w) = \frac{1}{2N} \sum_{i=1}^N \left(\sigma\left(\tilde{X}^i w\right) - y^i \right)^2$

There is no closed-form solution for w^* . We'll use $\operatorname{\textit{gradient descent}}$ (cf. notebook).



So we'll now minimize the following cost function $g_{\rm sl}(w) = \frac{1}{2N} \sum_{i=1}^N \left(\sigma\left(\tilde{X}^i w\right) - y^i \right)^2$

There is no closed-form solution for w^* . We'll use gradient descent (cf. notebook).



Proposition #3 - Cross entropy

A random variable follows a Bernoulli distribution if it only has two possible outcomes: 0 or 1.

$$P(Y = y \mid X) = p^y (1 - p)^{1-y}, \quad y = \{0, 1\}.$$

The likelihood of a Bernouilli distribution over N observations:

$$L(p) = \prod_{i=1}^{N} p^{y^i} (1-p)^{1-y^i}$$

which we wish to maximize.

Recall that
$$p^i = P\left(\hat{y}^i = 1 | \tilde{X}^i\right) = \sigma\left(\tilde{X}^i w\right)$$

We want to minimize the negative log-likelihood:

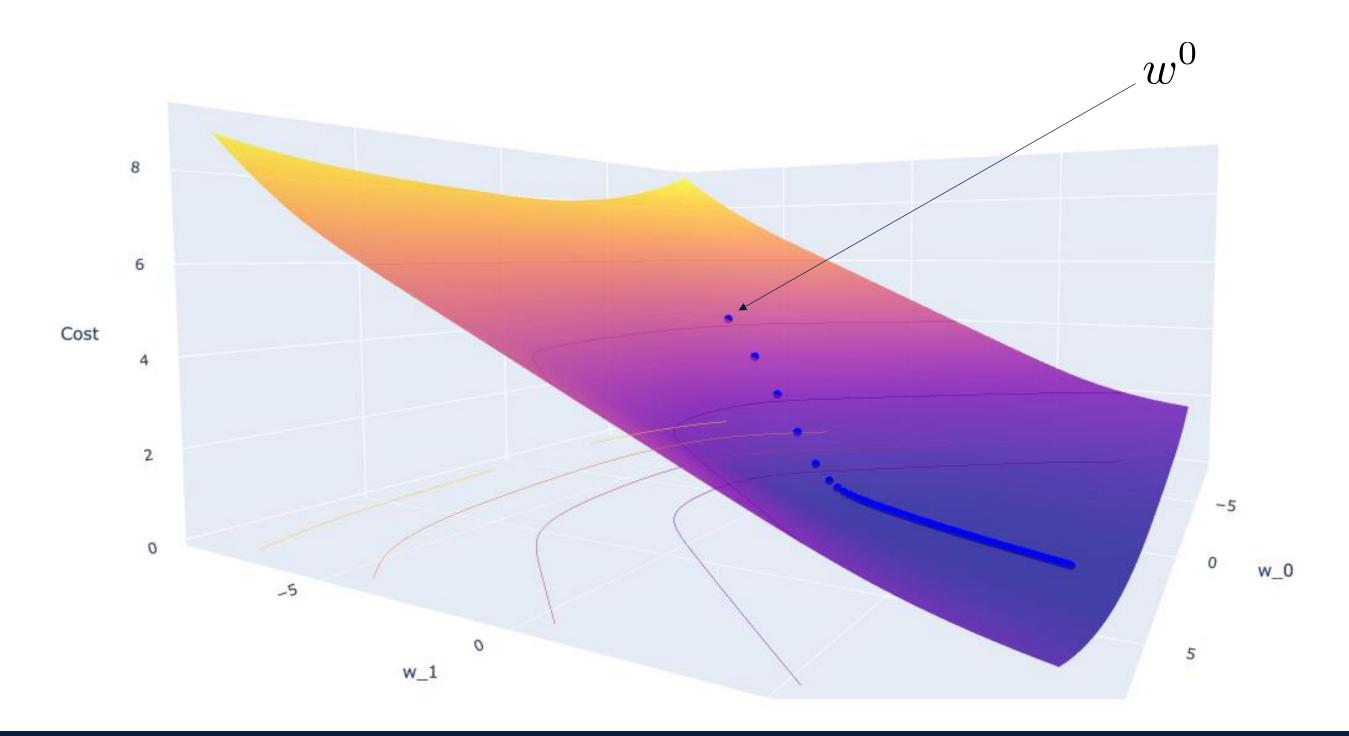
$$g_{\text{cel}} = -\frac{1}{N} \sum_{i=1}^{N} \left(y^{i} \ln \left(\sigma \left(\tilde{X}^{i} w \right) \right) + (1 - y^{i}) \ln \left(1 - \sigma \left(\tilde{X}^{i} w \right) \right) \right)$$



We want to minimize the negative of the log-likelihood (or cross entropy):

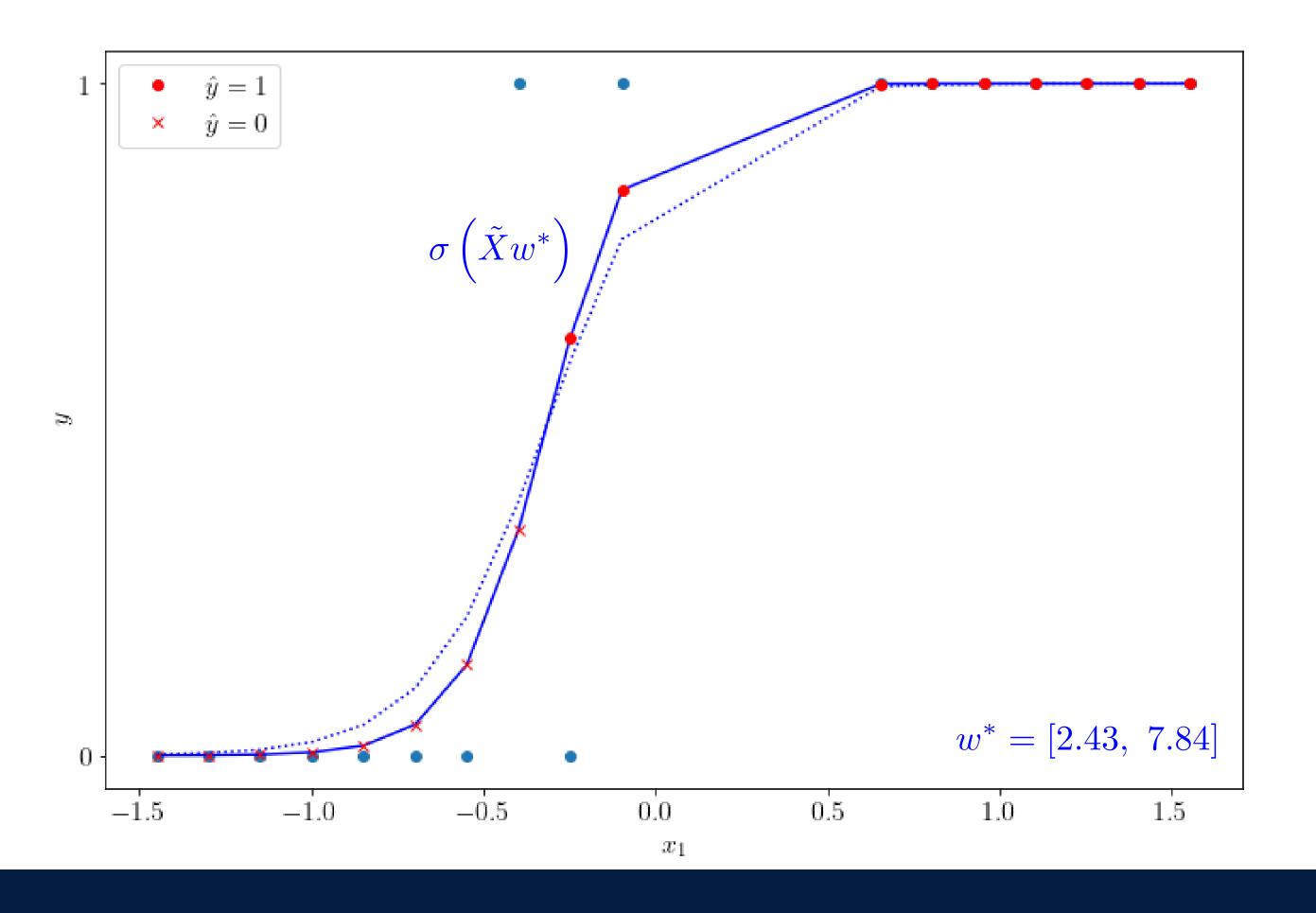
$$g_{\text{cel}} = -\frac{1}{N} \sum_{i=1}^{N} \left(y^{i} \ln \left(\sigma \left(\tilde{X}^{i} w \right) \right) + (1 - y^{i}) \ln \left(1 - \sigma \left(\tilde{X}^{i} w \right) \right) \right)$$

There is no closed-form solution fo w^* . We'll use **gradient descent** (cf. notebook).



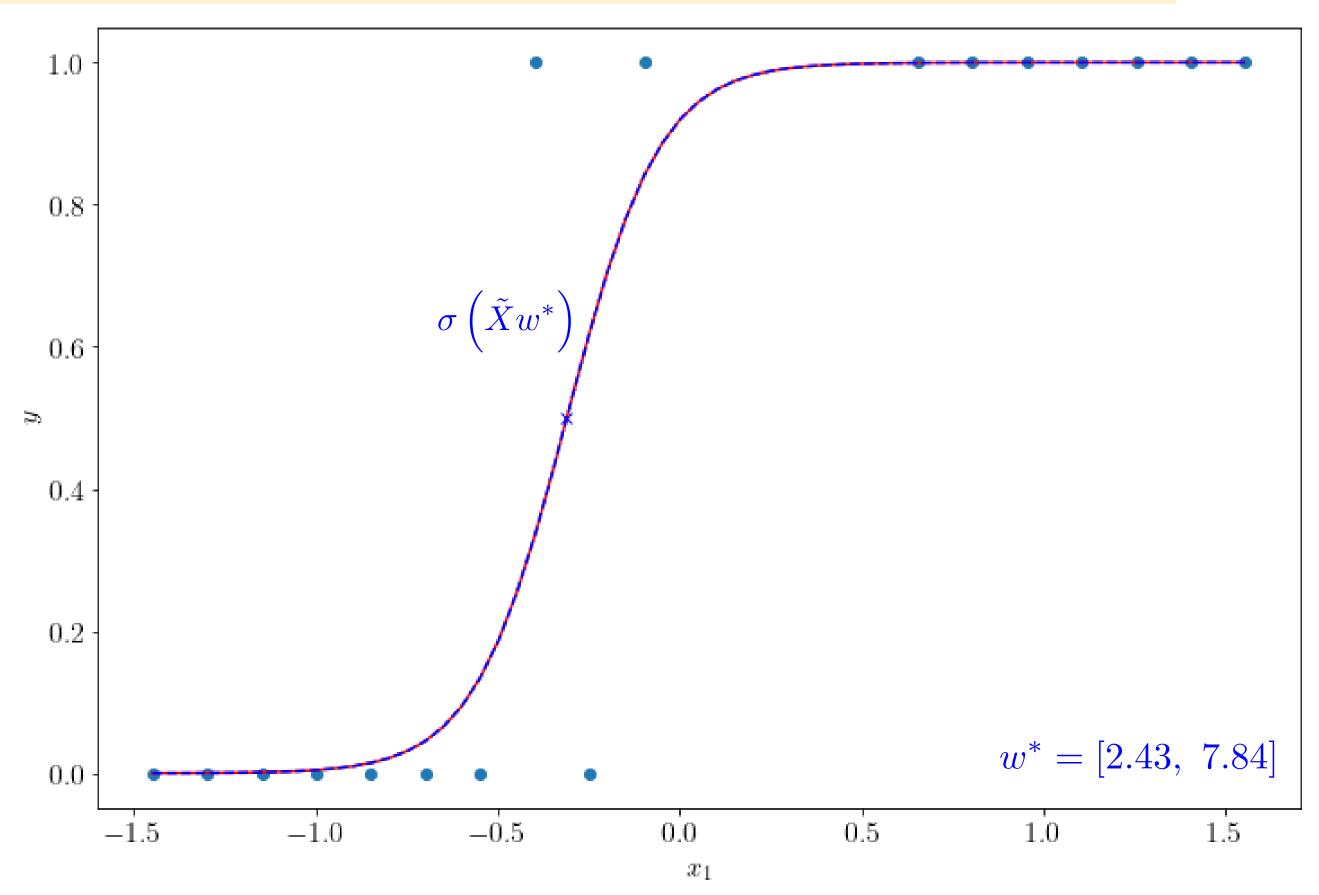
Minimizing the Cross Entropy

$$g_{\text{cel}} = -\frac{1}{N} \sum_{i=1}^{N} \left(y^{i} \ln \left(\sigma \left(\tilde{X}^{i} w \right) \right) + (1 - y^{i}) \ln \left(1 - \sigma \left(\tilde{X}^{i} w \right) \right) \right)$$



Now using scikit-learn

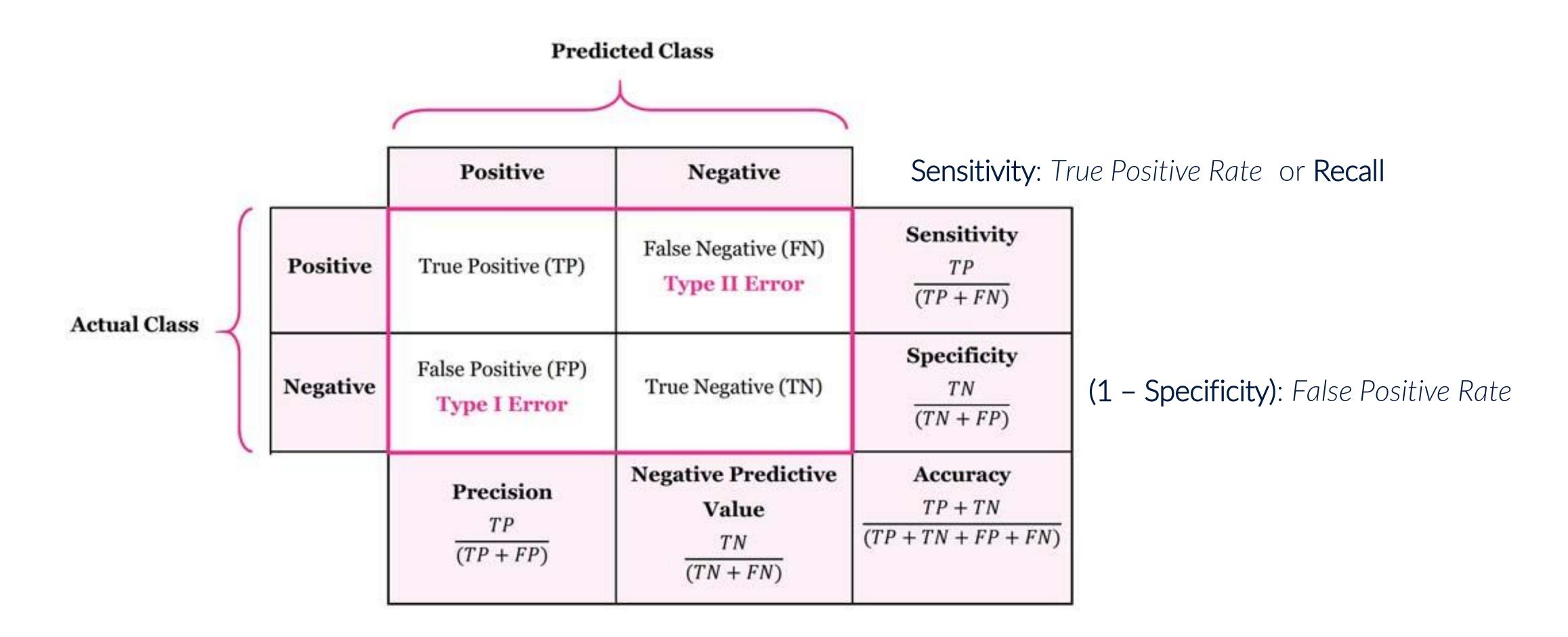
```
# Régularisation: 'none' ou '12'
clf = sklearn.linear_model.LogisticRegression(solver='lbfgs',penalty='none')
clf.fit(X, y.reshape(N,))
w_0, w_1 = clf.intercept_[0], clf.coef_[0][0]
```







Confusion matrix and statistics



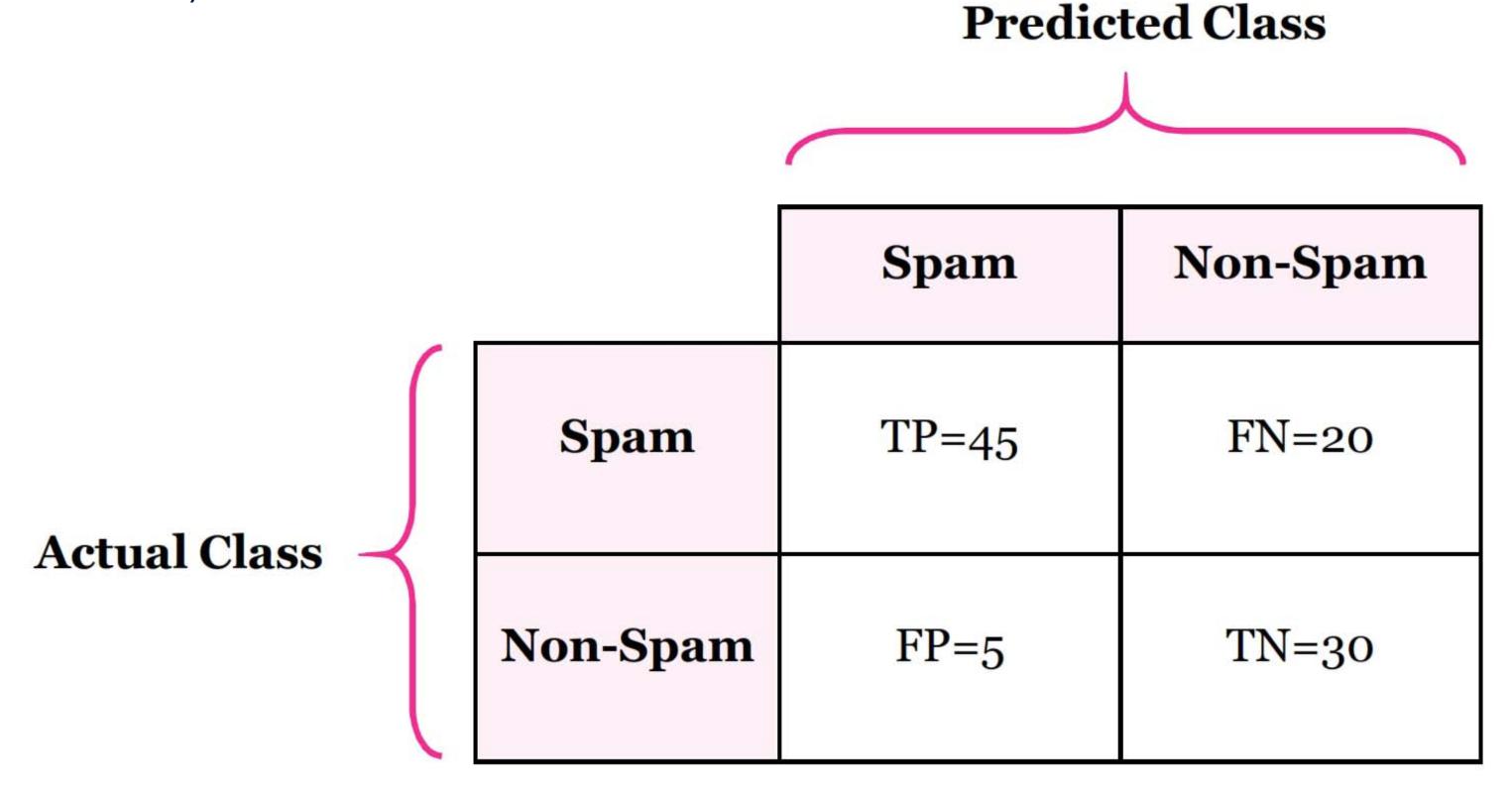
F1 Score =
$$2 \times \frac{Precision \times Recall}{Precision + Recall}$$

Source: Data Science and machine learning: The confusion matrix



Exercise

- Calculate precision, sensitivity and specificity
- Calculate accuracy and F1 score

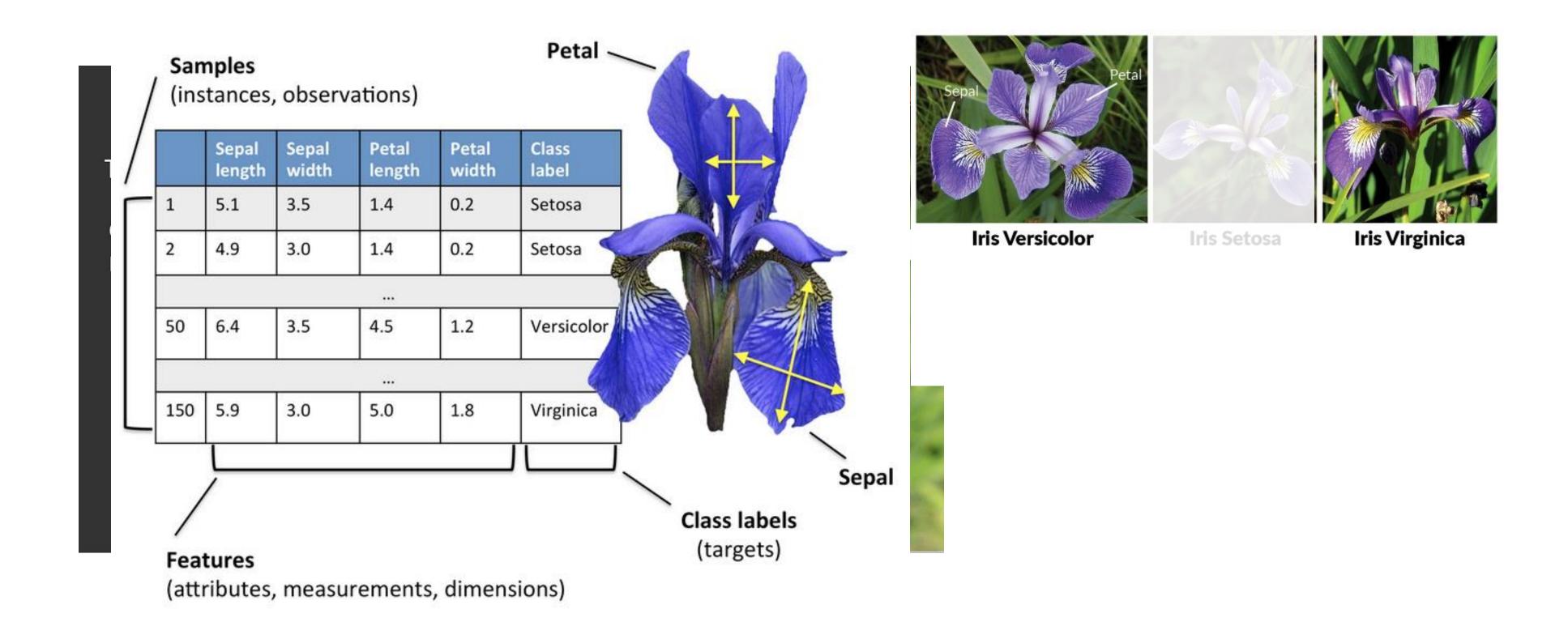


Source: <u>Data Science and machine learning: The confusion matrix</u>



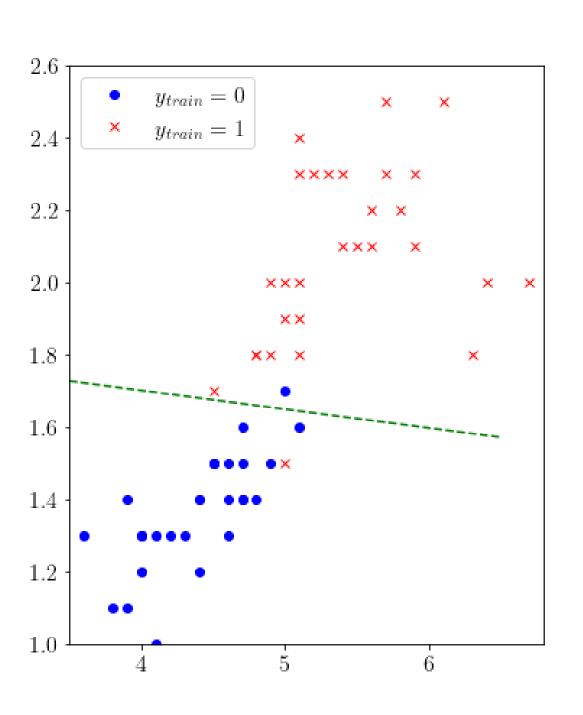
Real example using scikit-learn

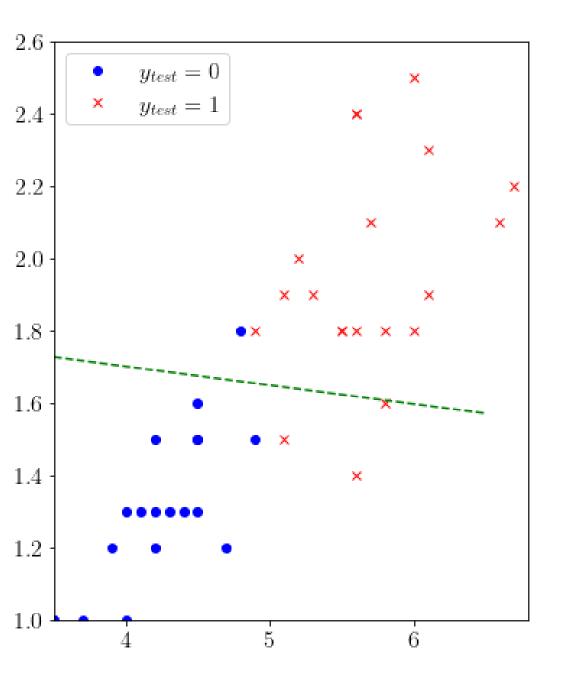
UCI Machine Learning Repository

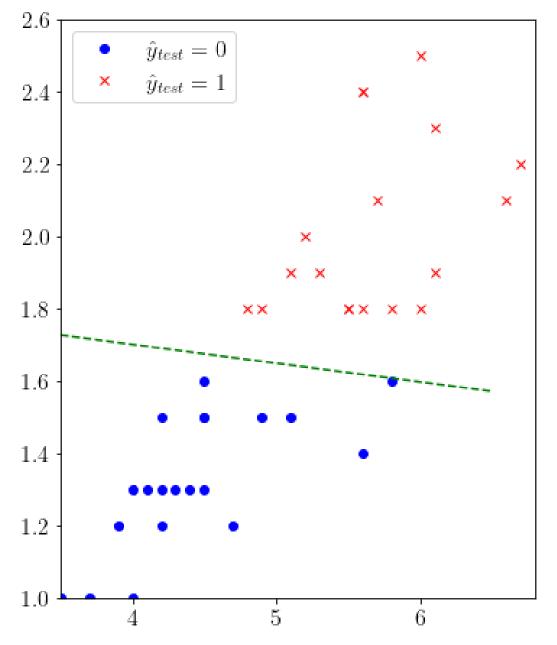


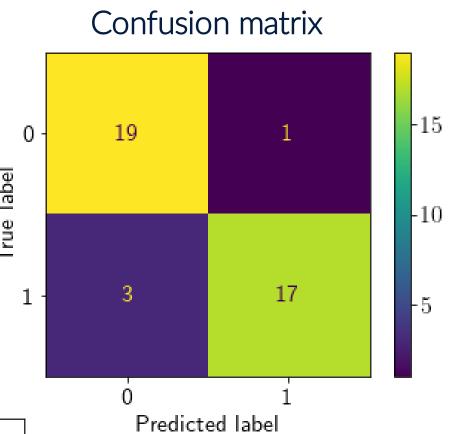
Iris dataset {versicolor: 0, virginica: 1}

- Using train-test-split, 30% test
- Using logistic regression (no regularization)









What did I learn?

- Definition of (binary) classification (vs regression)
- Logistic regression using various cost functions
- In particular, the cross entropy cost (or negative log-likelihood)
- Confusion matrix for classification performance



