

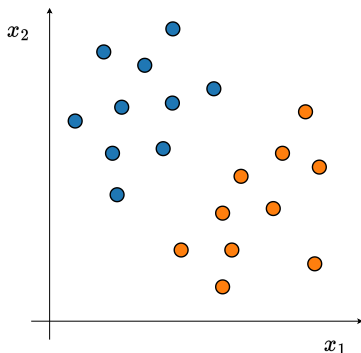
Linear classifiers

Classification problem

For the **classification** task, the target variable y to predict is **discrete**, *i.e.* for K classes (or labels), y can only take K distinct values.

For example, we can have $y \in \{0, \dots, K-1\}$ or $y \in \{-1, +1\}$.

We will first focus on the **binary classification** problem, where $K = 2$.



Linear model for classification

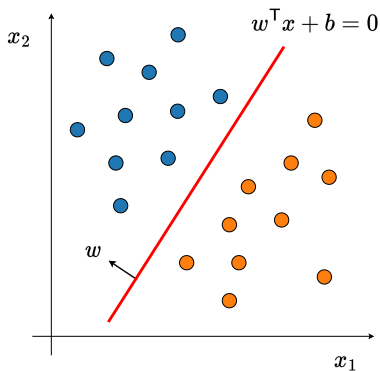
Our model will still be a **linear function**, but also with an **activation function** (also called *discriminant function*):

$$y = f(x; w, b) = \sigma(w^T x + b)$$

For example with:

$$\sigma(x) = \text{sign}(x) = \begin{cases} +1 & \text{if } x \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

this is called the **perceptron** model.



Logistic regression

If we use the **sigmoid** function (or *logistic* function) as the activation function, the model is called **logistic regression**.

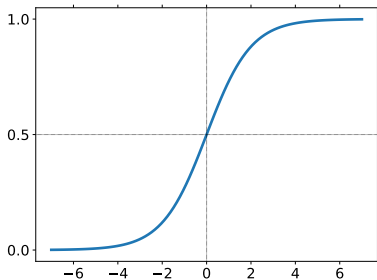
$$y = \sigma(z) \quad \text{with} \quad z = \mathbf{w}^T \mathbf{x} + b \quad \text{and} \quad \sigma(z) = \frac{1}{1 + e^{-z}}$$

The scalar z before the activation is called a **logit**.

The output of the model can be interpreted as a **probability**.

We can take the final classification decision by thresholding:

$$\begin{cases} \text{class 1} & \text{if } y \geq 0.5 \\ \text{class 0} & \text{otherwise} \end{cases}$$



Maximum likelihood for logistic regression

With logistic regression formulation, we can write:

$$p(y_i | \mathbf{x}_i, \mathbf{w}, b) = \underbrace{\sigma(\mathbf{w}^\top \mathbf{x} + b)}_{\hat{y}_i}^{y_i} (1 - \underbrace{\sigma(\mathbf{w}^\top \mathbf{x} + b)}_{\hat{y}_i})^{(1-y_i)}$$

We can apply maximum likelihood (ML) estimation on the training data:

$$\mathbf{w}_{\text{ML}}, b_{\text{ML}} = \arg \max_{\mathbf{w}, b} \log p(\mathbf{y} | \mathbf{X}, \mathbf{w}, b)$$

Assuming *i.i.d.* samples, we obtain:

$$\begin{aligned} \log p(\mathbf{y} | \mathbf{X}, \mathbf{w}, b) &= \log \prod_{i=1}^N p(y_i | \mathbf{x}_i, \mathbf{w}, b) = \sum_{i=1}^N \log p(y_i | \mathbf{x}_i, \mathbf{w}, b) \\ &= \sum_{i=1}^N \log \left((\hat{y}_i)^{y_i} (1 - \hat{y}_i)^{(1-y_i)} \right) \\ &= \sum_{i=1}^N (y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i)) \end{aligned}$$

Binary cross-entropy

Following the maximum likelihood principle, we typically want to **minimize** the loss function:

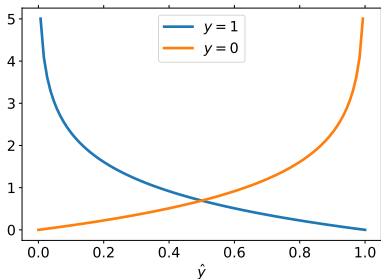
$$\mathcal{L}_{\text{CE}}(\mathbf{w}, b) = -\frac{1}{N} \sum_{i=1}^N (y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i))$$

where $\hat{y}_i = \sigma(\mathbf{w}^\top \mathbf{x}_i + b)$ is the prediction of the model.

This loss function is called **binary cross-entropy**.

This loss function is **convex**. But unlike linear regression, there is **no closed-form analytical solution** to this minimization problem.

→ Use **gradient descent**.



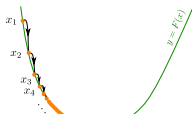
Gradient descent

Principle

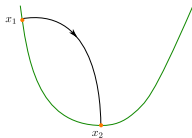
- Update the parameters by taking steps in the opposite direction of the gradient of the function to minimize:

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \gamma \nabla \mathcal{L}(\mathbf{w}^{(t)}), \gamma > 0.$$

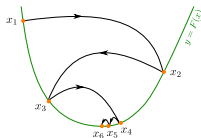
- The hyperparameter γ is called the **learning rate** (or *step size*).



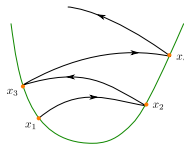
Small step size
Slow convergence



Good step size
Fast convergence



Large step size
Slow convergence



Too large step size
Divergence

Gradient descent for binary cross-entropy

Recall the binary cross-entropy loss:

$$\mathcal{L}_{\text{CE}}(\mathbf{w}, b) = -\frac{1}{N} \sum_{i=1}^N (y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i))$$

where $\hat{y}_i = \sigma(\mathbf{w}^\top \mathbf{x}_i + b)$ is the prediction of the model.

By successive applications of the chain rule, we find that:

$$\begin{aligned}\nabla_{\mathbf{w}} \mathcal{L}_{\text{CE}} &= \frac{1}{N} \sum_{i=1}^N (\hat{y}_i - y_i) \mathbf{x}_i \\ \frac{\partial \mathcal{L}_{\text{CE}}}{\partial b} &= \frac{1}{N} \sum_{i=1}^N (\hat{y}_i - y_i).\end{aligned}$$

Extension to multiclass classification

Logistic regression can be extended to the **multiclass** setting (K classes):

- Output targets are encoded into one-hot vectors $\mathbf{y} \in \mathbb{R}^K$
(for example: $\mathbf{y} = [0, 1, 0, 0]^\top$ for class 2 out of 4)
- Learnable parameters are $\mathbf{W} \in \mathbb{R}^{D \times K}$ and $\mathbf{b} \in \mathbb{R}^K$
- Predictions are $\hat{\mathbf{y}} = \sigma(\mathbf{W}^\top \mathbf{x} + \mathbf{b})$, where σ is applied component-wise
- σ is the softmax function: $\sigma(\mathbf{x})_i = \frac{\exp x_i}{\sum_{j=1}^K \exp x_j}$
- The loss is the **categorical cross-entropy** function:

$$\mathcal{L}_{\text{CCE}}(\mathbf{W}, \mathbf{b}) = -\frac{1}{N} \sum_{i=1}^N \ell(\mathbf{y}_i, \hat{\mathbf{y}}_i) \quad \text{with} \quad \ell(\mathbf{y}, \hat{\mathbf{y}}) = \sum_{j=1}^K y_j \log \hat{y}_j$$