

# Neural Networks for Data Science Applications

Master's Degree in Data Science

## Lecture 2: Supervised learning (for vectors)

---

Lecturer: S. Scardapane



**SAPIENZA**  
UNIVERSITÀ DI ROMA

# Supervised learning

---

Setup and examples

A (supervised) dataset is a set of  $n$  **examples**:

$$\mathcal{S} = \{(x_1, y_1), \dots, (x_n, y_n)\} . \quad (1)$$

Informally, given a ‘new’ pair  $(x, y)$  not contained in  $\mathcal{S}$ , we want a function  $f(\cdot)$  such that:

$$f(x) \approx y . \quad (2)$$

More generally, we test the model on a separate dataset  $\mathcal{T}$  never seen during training, i.e.,  $\mathcal{S} \cap \mathcal{T} = \emptyset$ .

We always assume implicitly that the elements in  $\mathcal{S}$  and the elements in  $\mathcal{T}$  are taken from *the same* i.i.d, *unknown* distribution  $p(x, y)$ .

- ▶ **Identically distributed:** the data-generating process is stable (e.g., cats in an apartments).
- ▶ **Independently distributed:** there is no bias in the data collection (e.g., only siamese cats).

If the distribution between  $\mathcal{S}$  and  $\mathcal{T}$  varies, we talk about **domain shift**.

1. **Spam identification:**  $x_i$  is an email, and  $y_i$  describes its probability of being spam.
2. **Robot navigation:**  $x_i$  is a sensory representation of the environment, and  $y_i$  is a motor command.
3. **Text translation:**  $x_i$  is a text and  $y_i$  its corresponding translation.
4. **Product recommendation:**  $x_i$  is a user, and  $y_i$  its affinity w.r.t. a certain catalogue of products.

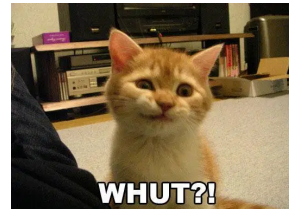
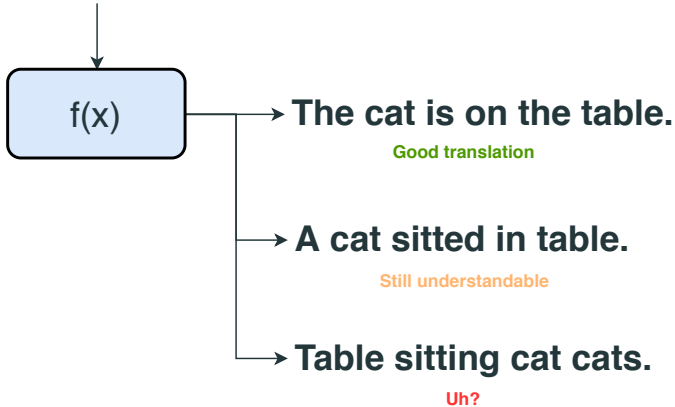
**Note:** ensuring the i.i.d. property sometimes is far from trivial!

# Supervised learning

---

Loss functions

"Il gatto è sul tavolo."



Given a point  $x$ , a desired value  $y$ , and a prediction  $\hat{y} = f(x)$ , we formalize its quality with a **loss function**  $l(y, \hat{y})$ , such that:

1. Low value of loss: good approximation;
2. High value of loss: poor approximation.

In this way, learning becomes a problem of minimizing a certain loss quantity that we designed.



The **expected loss** (risk) of a function  $f$  is:

$$f^*(x) = \arg \min \left\{ \mathbb{E}_{p(x,y)} [l(y, f(x))] \right\} . \quad (3)$$

The expected risk is uncomputable, but can be approximated via **empirical risk minimization**:

$$f^*(x) = \arg \min_f \left\{ \frac{1}{n} \sum_{i=1}^n l(y_i, f(x_i)) \right\} . \quad (4)$$

The gap between the two approaches is called **generalization gap**.

To begin our exploration of supervised learning, we will make a few simplifying assumptions:

- ▶ The input  $\mathbf{x}$  is a vector of shape  $d$ .
- ▶ The output  $y \in \mathbb{R}$  is a single real number.

In this case, a common loss is the squared norm:

$$l(y, \hat{y}) = (y - \hat{y})^2 . \quad (5)$$

Alternative losses are the absolute value  $|y - \hat{y}|$  and the Huber loss.

# Linear models

---

## Linear models for regression

A linear model  $f$  is defined as:

$$f(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle = \mathbf{w}^T \mathbf{x} = \sum_j w_j x_j, \quad (6)$$

where  $\mathbf{w}$  is a vector of adaptable parameters.

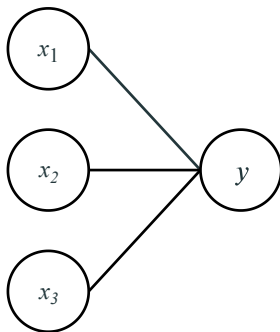
This model is fundamental in many disciplines, ranging from econometrics to statistics.

A more general formulation considers the inclusion of an **offset** (bias)  $b \in \mathbb{R}$ :

$$f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b. \quad (7)$$

Because we can always rewrite this as  $\mathbf{w}^\top \bar{\mathbf{x}}$ , with  $\bar{\mathbf{x}} = [\mathbf{x}; 1]$ , we can avoid writing the bias explicitly to simplify the notation.

**Hint:** Everytime we write a linear model, mentally add an offset term whenever needed.



**Figure 1:** Each arrow represents a *linear* influence on the destination, which sums the results.

Combining the squared loss with a linear model results in the **least-squares** optimization problem:

$$\text{LS}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2. \quad (8)$$

We can vectorize LS as:

$$\text{LS}(\mathbf{w}) = \frac{1}{n} \left\| \underset{(n)}{\mathbf{y}} - \underset{(n,d)(d)}{\mathbf{X} \mathbf{w}} \right\|^2, \quad (9)$$

where  $[\mathbf{X}]_i = \mathbf{x}_i$  and  $[\mathbf{y}]_i = y_i$ .

LS is a convex problem, with a simple gradient (**normal equations**):

$$\nabla \text{LS}(\mathbf{w}) = -\frac{2}{n} \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\mathbf{w}) . \quad (10)$$

However, LS is special in the sense that  $\nabla \text{LS}(\mathbf{w}) = 0$  is a linear equation that can be solved explicitly:

$$\mathbf{w}^* = \underbrace{(\mathbf{X}^\top \mathbf{X})^{-1}}_{(d,d)} \underbrace{\mathbf{X}^\top}_{(d,n)} \underbrace{\mathbf{y}}_{(n)} . \quad (11)$$



Numerical problems in the inversion of  $(\mathbf{X}^\top \mathbf{X})$  can be solved by adding a small amount of  $\ell_2$  **regularization** (**ridge regression**):

$$\text{LS-REG}(\mathbf{w}) = \text{LS}(\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|^2, \quad (12)$$

for some  $\lambda > 0$ . This makes the problem *strictly* convex and forces the solution to be contained in a ball of given radius, modifying the gradient and the explicit solution as:

$$\nabla \text{LS-REG}(\mathbf{w}) = \nabla \text{LS}(\mathbf{w}) + \lambda \mathbf{w}. \quad (13)$$

$$\mathbf{w}^* = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}. \quad (14)$$

where  $\mathbf{I}$  is the identity matrix of appropriate shape.

Generating some data:

```
1 # Linear model with unknown coefficients
2 X = tf.random.normal((10, 5))
3 y = X @ tf.random.normal((5, 1))
```

Computing a linear model:

```
1 w = tf.random.normal((5, 1))
2 yhat = X @ w # (10, 1)
```

Computing the objective function:

```
1 mse = tf.reduce_sum((y - yhat)**2)
```

Explicit solution (numerically unstable):

```
1 wopt = tf.linalg.inv(tf.transpose(X) @ X) @ tf.transpose(X) @ y
```

Explicit solution (better numerical conditioning):

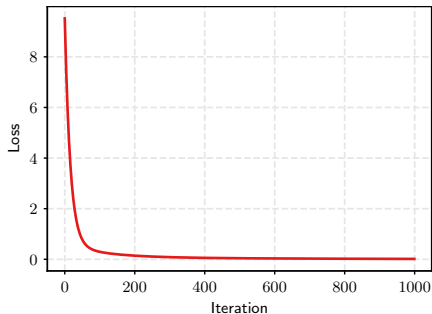
```
1 wopt = tf.linalg.solve(tf.transpose(X) @ X, tf.transpose(X) @ y)
```

Simple implementation of gradient descent:

---

```
1 for i in range(15000):  
2     # Note the sign: the derivative has a minus!  
3     w = w + 0.001 * tf.transpose(X) @ (y - X @ w)
```

---



# Linear models

---

## Linear models for classification

Another important class of supervised learning problems is **classification**, where  $y$  is an integer  $\{1, \dots, c\}$ , such that  $y_i = j$  means that  $\mathbf{x}_i$  is of class  $j$ .

For example, with  $c = 3$  we might have:

- ▶  $y = 1$ : the email is spam;
- ▶  $y = 2$ : the email is legit;
- ▶  $y = 3$ : the email is dubious.

Solving these as regression tasks is generally not an optimal choice: among other things, it is not guaranteed that classes have a definite ordering.

A common solution is to predict a **probability distribution** over the classes.

A vector **a** belongs to the **probability simplex**  $\Delta_c$  if:  
(c)

$$\sum_i [a]_i = 1, \quad [a]_i \geq 0. \quad (15)$$

If  $f(\mathbf{x}) = \hat{\mathbf{y}} \in \Delta_c$ , we can interpret it as a probability distribution, e.g., we can select the class with highest probability as:

$$\text{class} = \arg \max_i [\hat{\mathbf{y}}]_i. \quad (16)$$

Note that we *cannot* directly predict an integer with our models, because it would require some form of threshold operation which is not compatible with gradient descent (gradient zero almost everywhere).

Predicting a probability distribution can be seen as a *soft* approximation to this problem.



The **softmax** function maps any vector to the probability simplex:

$$[\text{softmax}(\mathbf{a})]_i = \frac{\exp(a_i)}{\sum_j \exp(a_j)} \quad (17)$$

The numerator ensures that all outputs are positive, while the denominator ensures that the final vector sums to 1.

Our linear model for classification becomes:

$$f(\mathbf{x}) = \text{softmax}(\underbrace{\mathbf{W}}_{(c,d)} \cdot \underbrace{\mathbf{x}}_{(d)}) \quad (18)$$

The pre-softmax values  $\mathbf{W}\mathbf{x}$  are called the **logits** of the model.

In order to compare the predictions with the ground truth, we encode our targets using a **one-hot encoding**. Given a pair  $(\mathbf{x}, \mathbf{y})$ :

$$y_i = \begin{cases} 1 & \text{if } \mathbf{x} \text{ is of class } i, \\ 0 & \text{otherwise.} \end{cases} \quad (19)$$

For example, with 3 classes  $\{\text{cat}, \text{dog}, \text{other}\}$ :

$$\text{cat} = [1, 0, 0] \quad \text{dog} = [0, 1, 0] \quad \text{other} = [0, 0, 1]. \quad (20)$$

This is a probability distribution putting all the **mass** on a single class.

Finally, we need a loss function  $l$  to compare two probability distributions.

The **cross-entropy** loss is defined for two vectors  $\mathbf{y}, \hat{\mathbf{y}} \in \Delta_c$  as:

$$\text{CE}(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_i y_i \log(\hat{y}_i) . \quad (21)$$

The cross-entropy can be interpreted as the Kullback–Leibler divergence (a common distance measure between probability distributions) between  $\mathbf{y}$  and  $\hat{\mathbf{y}}$ .

A **logistic regression** is a linear model  $f(\mathbf{x}) = \text{softmax}(\mathbf{W}\mathbf{x})$  trained by optimizing the cross-entropy:

$$\text{LR}(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \text{CE}(\mathbf{y}_i, f(\mathbf{x}_i)) . \quad (22)$$

It is not possible to solve the logistic regression problem explicitly. A linear model for classification has  $dc$  parameters.

A special case is **binary classification**, where  $c = 2$ . In this case, we can predict a single scalar value  $f(\mathbf{x}) \in [0, 1]$  since:

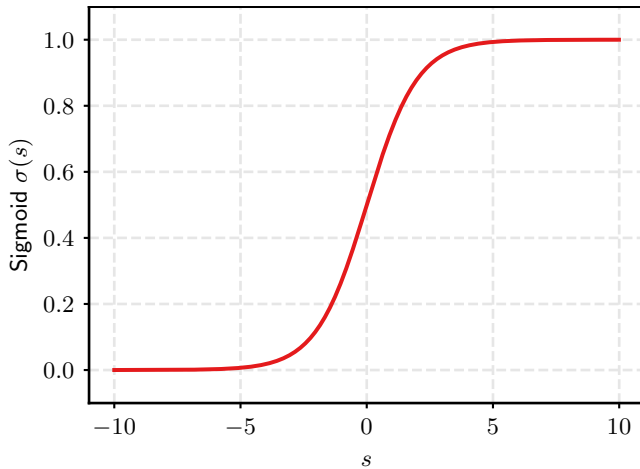
$$f(\mathbf{x}) \quad \text{probability of class 1 ,} \quad (23)$$

$$1 - f(\mathbf{x}) \quad \text{probability of class 2 .} \quad (24)$$

In this case, the softmax simplifies to the **sigmoid function**:

The sigmoid  $\sigma(s) \in [0, 1]$  is defined as:

$$\sigma(s) = \frac{1}{1 + \exp(-s)} . \quad (25)$$



**Figure 2:** A visualization of the sigmoid function. Note that 0 and 1 are only approached asymptotically.

Combining everything, we obtain a binary version of the logistic regression algorithm:

$$\text{BIN-LR}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \left[ \underbrace{-y_i \log(\sigma(\mathbf{w}^\top \mathbf{x}))}_{\text{Class 1}} - \underbrace{(1 - y_i) \log(1 - \sigma(\mathbf{w}^\top \mathbf{x}))}_{\text{Class 2}} \right] \quad (26)$$

In this case, we can obtain the most probable class from the model as:

$$\text{class} = \begin{cases} 1 & \text{if } \sigma(\mathbf{w}^\top \mathbf{x}) > 0.5, \\ 2 & \text{otherwise .} \end{cases} \quad (27)$$

By manually differentiating we obtain:

$$\sigma'(s) = \sigma(s)(1 - \sigma(s)). \quad (28)$$

Plugging this into the gradient computation we obtain:

$$\nabla \text{BIN-LR}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n (\sigma(\mathbf{w}^\top \mathbf{x}_i) - y_i) \mathbf{x}_i, \quad (29)$$

showing its similarity to the regression case.



```
1 from tensorflow.keras.metrics import *
2
3 # The one we have described up to now.
4 categorical_crossentropy(ytrue, yhat)
5
6 # ytrue should contain the indexes of the classes instead of the
7 # one-hot encodings.
8 sparse_categorical_crossentropy(ytrue, yhat)
9
10 # Numerically-stable versions requiring the logits as inputs
11 # (see the LogSumExp trick).
12 categorical_crossentropy(ytrue, yhat, from_logits=True)
13 sparse_categorical_crossentropy(ytrue, yhat, from_logits=True)
```

# Linear models

---

Calibration and a probabilistic  
formulation

More formally, for the classification setting we assumed  $f(x)$  parameterizes a *probability distribution*  $\hat{p}(y | f(x))$  over the output  $y$ . We can always do this, e.g., for regression:

$$\hat{p}(y | f(x)) = \mathcal{N}(y | f(x), \sigma^2), \quad (30)$$

where the model predicts the center of a Gaussian distribution with fixed variance (hyper-parameter). This probabilistic formulation can be more flexible or useful in many contexts.

The probabilistic formulation also provides a principled way to interpret training by maximizing the **likelihood** of the model (assuming the elements of the dataset are i.i.d.):

$$f^*(x) = \arg \max \prod_{(x_i, y_i)} \hat{p}(y_i | f(x_i)) . \quad (31)$$

For (30), this is **equivalent** to training with a squared loss. Similarly, training with cross-entropy is **equivalent** to assuming a categorical distribution over the output (can you prove it?).

A common misconception when doing classification is that  $[f(x)]_i$  can be immediately interpreted as *the probability of pattern  $x$  being of class  $i$* .

However, this is only true whenever the trained model satisfies:

$$p(y = i | x) = [f(x)]_i. \quad (32)$$

We say the model is well **calibrated**, but this must be checked manually.

---

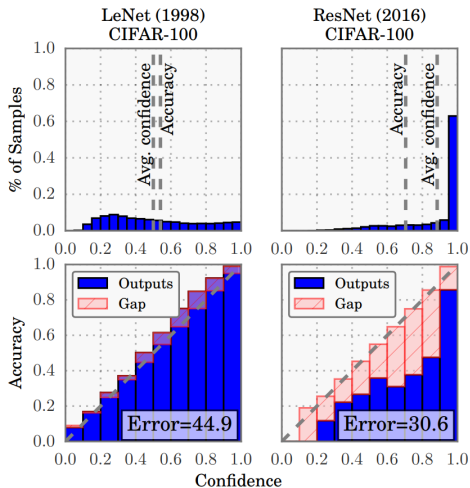
Guo, C., et al.. On calibration of modern neural networks. ICML 2017.

To measure the calibration of a model, we keep a separate validation set, and we split the interval  $[0, 1]$  into  $m$  equispaced bins (each of size  $1/m$ ). Define:

- ▶  $B_m$  the number of samples from the validation set, whose predicted confidence falls in bin  $m$ .
- ▶  $p_m$  the average confidence of the network for that bin.
- ▶  $a_m$  the average accuracy of the network for these elements.

Then, the **expected calibration error** (ECE) is given by:

$$\text{ECE} = \sum_m \frac{B_m}{n} |a_m - p_m|. \quad (33)$$



**Figure 3:** Plotting  $a_m$  against  $p_m$  for every bin gives us a **reliability plot** (from Guo et al., 2017).

This topic is important because more complex networks may be highly over (or under) confident, with many methods to improve it (temperature scaling, logit normalization, ...).

A simple (and popular) option is to decrease the weight given to ‘easy’ samples using a variant of cross-entropy call the **focal loss**:

$$\text{FL}_\alpha(\mathbf{y}, \hat{\mathbf{y}}) = -(1 - \hat{y}_c)^\alpha \log \hat{y}_c, \quad (34)$$

where  $c = \arg \max \mathbf{y}$ .

---

Mukhoti, J., et al., 2020. **Calibrating deep neural networks using focal loss**. *Advances in Neural Information Processing Systems*, 33, pp. 15288-15299.



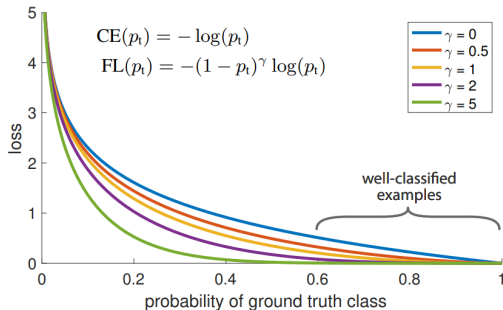


Figure 1. We propose a novel loss we term the *Focal Loss* that adds a factor  $(1 - p_t)^\gamma$  to the standard cross entropy criterion. Setting  $\gamma > 0$  reduces the relative loss for well-classified examples ( $p_t > .5$ ), putting more focus on hard, misclassified examples. As

- ▶ **(Required)** Chapter 3 and 4 from the book.