

Neural Networks for Data Science Applications

Master's Degree in Data Science

Lecture 4: Fully-connected neural networks

Lecturer: S. Scardapane



SAPIENZA
UNIVERSITÀ DI ROMA

Feedforward neural networks

Limitations of linear models

To understand the limitations of a linear model $f(\mathbf{x})$, consider an input vector $\hat{\mathbf{x}}$, equal to \mathbf{x} but for a single feature $\hat{x}_i = 2x_i$ (e.g., a client with double income).

The output on the two inputs is related by:

$$f(\hat{\mathbf{x}}) = f(\mathbf{x}) + w_i x_i. \quad (1)$$

Effects are linearly super-imposed: it is impossible to model some form of interaction between two features (e.g., 'a client with 10k income is not trustworthy, *unless he is very young*').

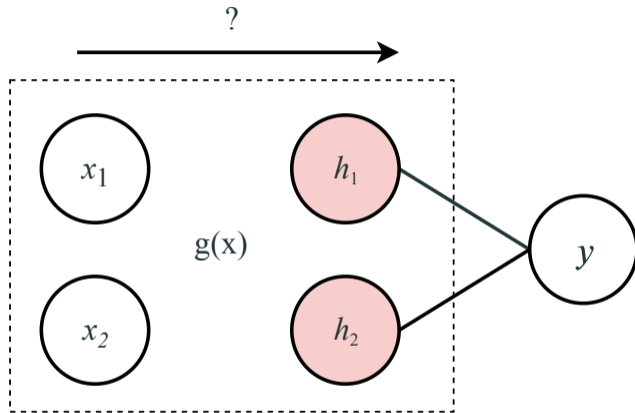


Figure 1: How can we define intermediate computations?

Feedforward neural networks

Fully-connected layers

Can we layer two linear models?

$$\mathbf{h} = \mathbf{W}\mathbf{x}, \quad (2)$$

$$f(\mathbf{h}) = \mathbf{w}^\top \mathbf{h}. \quad (3)$$

We could, but this is equivalent to a single linear model:

$$f(\mathbf{x}) = (\mathbf{w}^\top \mathbf{W}) \mathbf{x}. \quad (4)$$

We need some way of separating the two linear operations.

We can avoid the ‘collapse’ of the two linear models by interleaving them with some element-wise nonlinearity ϕ :

$$\mathbf{h} = \phi(\mathbf{W}\mathbf{x}), \quad (5)$$

$$f(\mathbf{h}) = \mathbf{w}^\top \mathbf{h}. \quad (6)$$

This is the prototype of a **fully-connected** (FC) neural network, sometimes known as a **multilayer perceptron** (MLP).

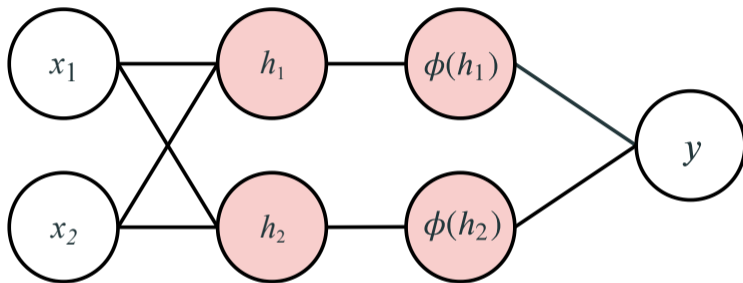


Figure 2: Visualization of a simple feedforward NN. In general, we will never show the nonlinearities explicitly in the graphs from now on.

Feedforward neural networks

Training and complexity considerations

Training of the network can be proceed similarly to a linear model. For example, given a dataset $\{\mathbf{x}_i, y_i\}_{i=1}^n$ for regression, we can minimize the LS loss (note: always remember we are ignoring the biases in the notation):

$$\mathbf{W}^*, \mathbf{w}^* = \arg \min \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2. \quad (7)$$

For classification, we can wrap the output in a sigmoid and minimize a cross-entropy loss.

Theorem 1.1: Cybenko (1989) - Hornik (1991) - Leshno (1993)

Let $h(\mathbf{x})$ be a continuous function defined on a compact subset $S \subset \mathbb{R}^d$ and $\varepsilon > 0$. For a sufficiently large p , there exists an $f(\mathbf{x})$ as in Eq. (6)-(7) with p hidden units such that:

$$|h(\mathbf{x}) - f(\mathbf{x})| < \varepsilon, \quad \forall \mathbf{x} \in S. \quad (8)$$

This holds for any non-constant, bounded, continuous ϕ .

This is a **universal approximation** theorem. It does not tell us about the *feasibility* for a given problem (e.g., how large should p be?).

Because the NN can be highly non-convex, its optimization problem has multiple local minimum and/or saddle points.

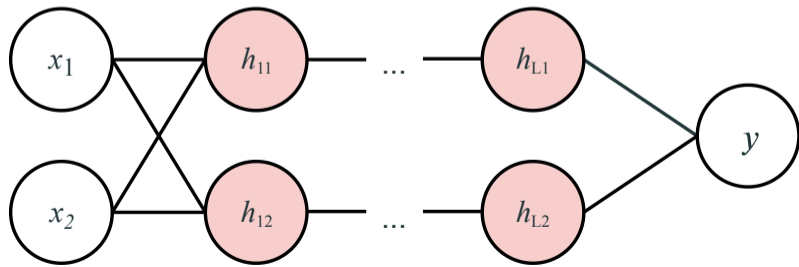
In fact, training a neural network *optimally* is **NP-hard**, even for very simple architectures, and highly dependent on a good initialization.

Finding the global optimum requires running GD from almost *everywhere*: this is similar to an exhaustive search.

Blum, A. and Rivest, R.L., 1989. Training a 3-node neural network is NP-complete. In Advances in neural information processing systems (pp. 494-501).

Nothing prevents us from adding *additional* 'hidden' (intermediate) layers:

$$f(x) = w^T \cdot \overbrace{\phi \left(\underbrace{Z \cdot \phi(W \cdot x) + c}_{h^1} \right)}^{h^2} \tag{9}$$



Differently from a linear model, a NN has several design choices that we have freedom on:

- ▶ The nonlinearity (sometimes called the **activation function**);
- ▶ The dimensionality of the hidden layers;
- ▶ Several others that we will introduce in the next lectures.

We call these **hyper-parameters** to differentiate them from parameters (weights) to be trained via GD.

Choosing the correct set of hyper-parameters is called the **model selection** / **hyper-parameter optimization** problem.

Playing with a neural network

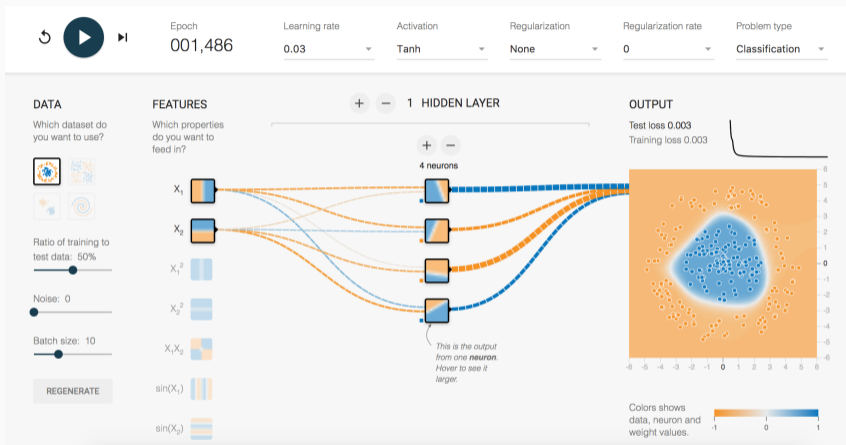


Figure 3: <https://playground.tensorflow.org/>.

Apart from the sigmoid, what kind of non-linearities can we use?

1. Hyperbolic tangent (**tanh**) $\tanh(s) = 2\sigma(s) - 1$, which is just a sigmoid scaled in $[-1, +1]$.
2. Polynomials $\phi(s) = s^p$ are a very bad idea unless a lot of care is taken due to numerical instabilities.
3. We will see that the **rectified linear unit** (ReLU) $\text{ReLU}(s) = \max(0, s)$ is a good default choice. Note that, if we consider the bias of the previous layer part of the activation function, $\text{ReLU}(s) = \max(0, s + b)$ is a threshold function with a trainable threshold.

1. **LeakyReLU** replaces the negative quadrant of ReLU with a small slope (e.g., $\alpha = 0.01$):

$$\text{LeakyReLU}(s) = \begin{cases} s & \text{if } s > 0 \\ \alpha s & \text{otherwise} \end{cases} \quad (10)$$

2. If you want a smoother version of the ReLU, you can use the **exponential linear unit** (ELU):

$$\text{ELU}(s) = \begin{cases} s & \text{if } s > 0 \\ \alpha(\exp s - 1) & \text{otherwise} \end{cases} \quad (11)$$

Viewing the ReLU as a kind of *gating* function $\text{ReLU}(s) = s1_{x \geq 0}$, we can obtain smoother variants as:

$$\phi(s) = s\psi(s), \quad (12)$$

where:

- ▶ If $\psi(s)$ is the standard Gaussian cumulative distribution function, we obtain the **Gaussian Exponential Linear Unit** (GELU).
- ▶ If $\psi(s) = \sigma(s)$ we obtain the **Sigmoid linear unit** (SiLU) or **Swish**.

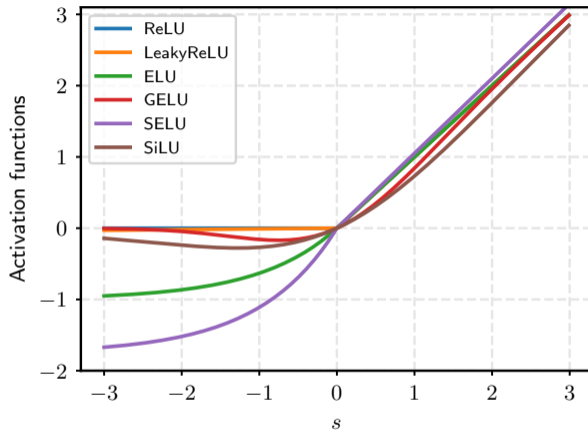


Figure 4: Common activation functions plotted side-by-side.

Activation functions can also have trainable parameters that vary for each unit. For example, the LeakyReLU with trainable α becomes the **parametric ReLU** (PReLU). Or, for a more convoluted example:

$$\text{Trainable-Swish}(s) = \sigma(as + b)(cs + d), \quad (13)$$

with 4 trainable parameters per unit. Variants of this have become popular with the recent LLaMA class of models.

Shazeer, N., 2020. **GLU variants improve transformer**. *arXiv preprint arXiv:2002.05202*.

We call the size of the hidden layers the **width** of the network, and the number of layers its **depth**. *Narrow* networks with potentially infinite layers are also universal approximators.¹

UA theorems are not constructive; there exist families of functions that are efficiently learnable with deep networks as opposed to shallow networks, but it is unclear whether this also holds in practice.²

¹Kidger, P. and Lyons, T., 2020. **Universal approximation with deep narrow networks**. In COLT (pp. 2306-2327). PMLR.

²Nye, M. and Saxe, A., 2018. **Are efficient deep representations learnable?**. arXiv preprint arXiv:1807.06399.

The distribution of NNs with infinite width converges to a mathematical object known as a **Gaussian Process** (GP).

In fact, under a first-order linearization, the predictions of any NN $f(\mathbf{x})$ can be described similarly to what we did in the linear case:

$$(\mathbf{y}_t - \mathbf{y}) = \left(\mathbf{I} - \eta \underbrace{[\partial f(\mathbf{x})]^\top \partial f(\mathbf{x})}_{\mathbf{K} \in \mathbb{R}^{n \times n}} \right) (\mathbf{y}_{t-1} - \mathbf{y}) . \quad (14)$$

The elements $[\mathbf{K}]_{i,j}$ are called the **neural tangent kernel** (NTK), a fundamental recent tool in the analysis of deep neural networks.

Training of the networks

Stochastic optimization

Consider the steps needed for a single iteration of GD:

1. Computing the output of the NN $f(\mathbf{x})$ on *all* examples;
2. Computing the gradient of the cost function with respect to the weights.

Both these operations scale *linearly* in the number of examples, which is unfeasible when we have 10^5 or even more examples.

In practice, to train NNs we use **stochastic** versions of GD, that approximate the real gradient from smaller amounts of data.

The key observation is that the training problem in NNs is an expectation with respect to all data (θ is the set of weights):

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2. \quad (15)$$

To limit the complexity, we can use only a **mini-batch** \mathcal{B} of M examples from the full dataset:

$$J(\theta) \approx \tilde{J}(\theta) = \frac{1}{M} \sum_{i \in \mathcal{B}} (y_i - f(\mathbf{x}_i))^2. \quad (16)$$

We can use the gradient from the approximated function to perform an iteration of GD.

The previous algorithm is called **stochastic gradient descent** (SGD).

The computational complexity of an iteration of SGD is fixed with respect to M (the **batch size**) and does not depend on the size of the dataset.

Because we assumed that samples are i.i.d., we can prove SGD also converges to a stationary point in average, albeit *with noisy steps*.

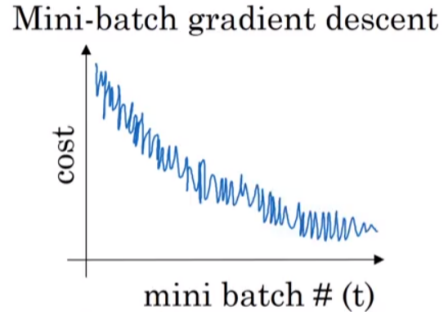
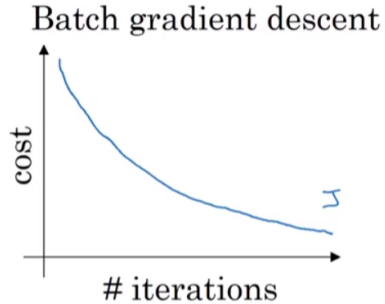


Figure 5: SGD will converge on average to a stationary point, although in an apparently noisy fashion (Source: EngMRK).

Instead of randomly sampling batches from the training dataset at each iteration, we generally apply the following procedure:

1. Shuffle the full dataset;
2. Split the dataset into blocks of M elements and process them sequentially, batch-by-batch;
3. After the last block, return to point (1) and iterate.

A full pass over the dataset is called an **epoch**. This is efficient because the majority of the time we deal with data stored *contiguously* in memory. In practice, even the shuffling operation can be approximated.

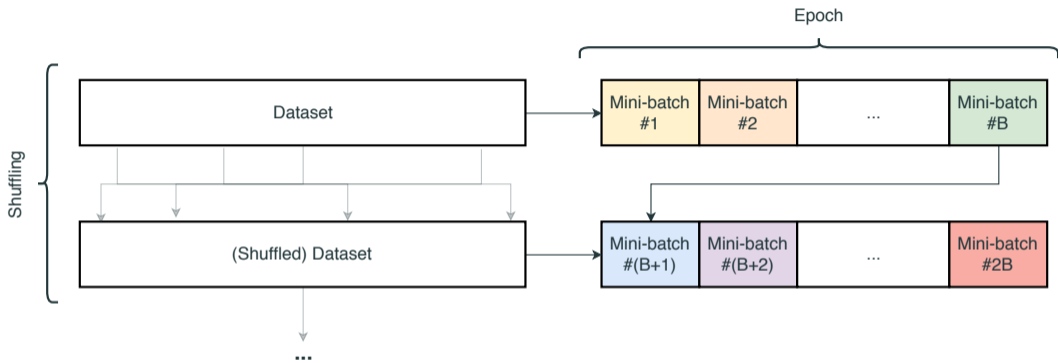


Figure 6: Dividing the optimization process into epochs is also helpful to define metrics and stopping criteria.

Running SGD requires selecting the size of the mini-batch, which is an additional hyper-parameter:

- ▶ Smaller mini-batches are faster, but the network might require more iterations to converge because the gradient is noisier.
- ▶ Larger mini-batches provide a more reliable estimation of the gradient, but are slower.

In general, it is typical to choose power-of-two sizes (32, 64, 128, ...) depending on the hardware configuration and the total memory available.

In the limit $M = 1$ we obtain an **online** (streaming) optimization.

- ▶ **D2L**: Chapter 5; **UDL**: Chapters 3-6³; **PPA**: Chapters 5 and 7.

³Of interest: the discussion on linear regions induced by a ReLU network.