

Artificial Neural Networks

Lecture 1: Introduction

Dr Juan Neirotti

j.p.neirotti@aston.ac.uk

2024

What is this Module About?

- ▶ In the most general terms: Modeling stationary processes, i.e. probabilistic processes where the density distribution does not changes with time.
- ▶ In the not-that-much general terms: Pattern recognition (speech, face, hand-written characters).
- ▶ All these tasks need a statistics approach to help extract the relevant features (patterns) out of many instances (big data).

Definition

1. Artificial Intelligence (AI) is the area of knowledge that endeavors towards constructing systems (hardware or software) that behave *intelligently*. (Observe that I have not defined what I mean by intelligent just yet).
2. Machine Learning (ML) is a sub-area of AI that tackles problems by extracting the patterns that link questions with correct answers (provided to the AI system during the *training phase*). The ML paradigm differs from the traditional manner to solve problems on the fact that, in the traditional way the rules (patterns) that produce an answer given a question are supposed to be known, whereas in the ML paradigm such rules are unknown and need to be discovered.

Definition

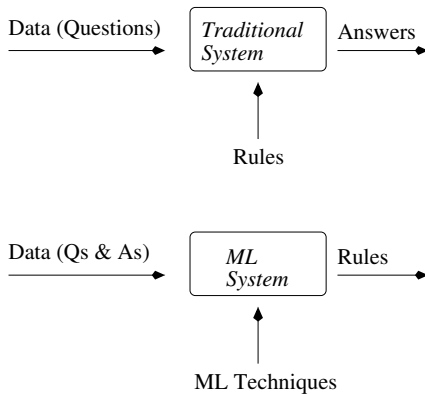


Figure: Different paradigms for solving problems

Definition

1. Statistical ML: The adjective *statistical* acknowledges the nature of the elements in the Data Set to be used $\mathcal{D}_u = \{\mathbf{x}_n\}_{n=1}^M$, which can be described through a probability distribution $\mathcal{P}(\mathbf{x})$, or $\mathcal{D}_s = \{(\mathbf{x}_n, t_n)\}_{n=1}^M$, which can be described through a probability distribution $\mathcal{P}(\mathbf{x}, t)$,
2. Most of the work in SML involves making models of the $\mathcal{P}(\mathbf{x})$ or $\mathcal{P}(\mathbf{x}, t)$ and produce results based on inferences using such a model (the ML techniques in the figure above).

Definition

1. If the type of data set we use to *train* (more on this later) our intelligent system is $\mathcal{D}_u = \{\mathbf{x}_n\}_{n=1}^M$ we say that the training process involves a *unsupervised learning*.
2. If the type of data set we use to *train* (more on this later) our intelligent system is $\mathcal{D}_s = \{(\mathbf{x}_n, t_n)\}_{n=1}^M$, we say that the training process involves a *supervised learning*.

Problems to be Tackled

1. Classification: Given a particular input $\mathbf{x} \in \mathcal{X}$ the intelligent system must produce a classification $t \in \mathcal{T} = \{-1, +1\}$. The map $t = t(\mathbf{x})$ is discrete.
2. Regression: Estimate the functional form of the map $t : \mathcal{X} \rightarrow \mathcal{T}$, i.e. $t = t(\mathbf{x})$, which is, in general, a continuous map.

Why Statistics

1. Both characters are drawn in a 256×256 pixels figure.
2. The total number of possible figures (in black and white) is $2^{256 \times 256} \sim 10^{20000}$. The size of the set with all the possible figures is huge.
3. No all the possible figures are meaningful. There are many (many indeed) figures that wouldn't carry any meaning at all.



Figure: Hand-written characters a and b.

Pattern extraction

1. We can define x_1 and x_2 as the horizontal and vertical dimensions of the characters.
2. a's and b's are of similar width, $x_1(a) \simeq x_2(b)$.
3. a's are expected to be shorter than b's, therefore $x_2(a) < x_2(b)$.
4. Both attributes are distributed variables $x_1(a) \sim \mathcal{P}_{1,a}$, $x_2(a) \sim \mathcal{P}_{2,a}$, $x_1(b) \sim \mathcal{P}_{1,b}$, $x_2(b) \sim \mathcal{P}_{2,b}$.

Pattern extraction

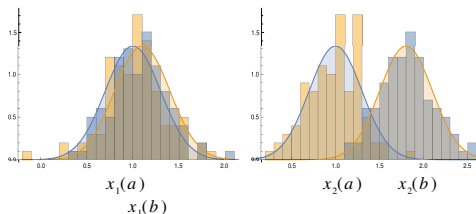
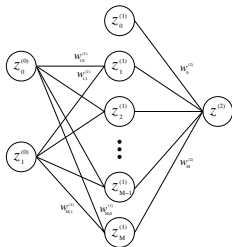


Figure: Histogram and correspondent distribution for x_1 (left) and x_2 (right). Observe that, even for the clear separation between x_2 distributions, there still exists an *overlapping* range of values where no decision can be taken.

Maps

1. $t : \mathcal{X} \rightarrow \mathcal{Y}$, where usually $\mathcal{X} \subset \mathbb{R}^d$
2. d is expected to be large.
3. \mathcal{Y} could be a discrete set (classification problem) or a continuous set (regression)
4. The *estimate* $\hat{t} = \hat{t}(\mathbf{x}, \mathbf{w})$, $\mathbf{x} \in \mathcal{X}$ is the input (or independent variable) and $\mathbf{w} \in \Omega$ are the parameters of the model.
5. Parameters (\mathbf{w}) and the functional form of the map ($\hat{t}(\mathbf{x}, \mathbf{w})$) are part of the model.
6. Neural networks are implementations of maps where the functional form t is fixed by the problem and the parameters \mathbf{w} are *learn* from the available data.

Maps



1.
$$z_k^{(\ell+1)} = \sigma \left(\left\langle z^{(\ell)} \middle| \mathbf{w}_k^{(\ell+1)} \right\rangle + w_{k,0}^{(\ell+1)} \right)$$
2. σ is the activation function that can be: sigmoidal, tanh, sgn, Heaviside, identity, ReLU, etc.
3. The activation variable is always linear:

$$\left\langle z^{(\ell)} \middle| \mathbf{w}_k^{(\ell+1)} \right\rangle + w_{k,0}^{(\ell+1)} = \sum_{j=1}^{\dim(\mathbf{z}^{(\ell)})} z_j^{(\ell)} w_{kj}^{(\ell+1)} + w_{k,0}^{(\ell+1)}$$
4. The number of layers ($0 \leq \ell \leq L$) and the type of units used determine the architecture of the network.

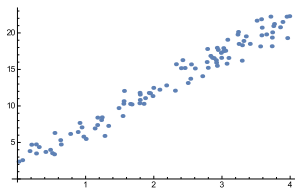
Learning the weights

1. Unsupervised Learning $\mathcal{D} = \{\mathbf{x}_\ell\}_{\ell=1}^N$ the objective is to estimate $\mathcal{P}_X(\mathbf{x}|\mathbf{w})$.
2. Supervised Learning $\mathcal{D} = \{(\mathbf{x}_\ell, t_\ell)\}_{\ell=1}^N$ the objective is to estimate the function $h(\mathbf{x}) = t$ through the map $\hat{t}(\mathbf{x}; \mathbf{w})$ (regression).

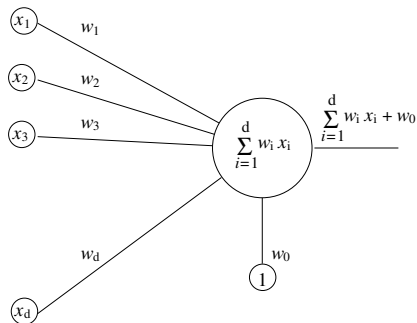
An Example: Regression with a linear function

- Let us suppose we have a set of data presented in the following form:

$$\mathcal{D} = \{(\mathbf{x}_\ell, t_\ell) : t_\ell \in \mathbb{R}, \mathbf{x}_\ell \in \mathbb{R}^d, \ell = 1, \dots, L\}$$



- ▶ The architecture used for this problem is the following:



- ▶ Because the output is the identity function $\sigma(y) = y$ this is known as the linear perceptron.

- ▶ We represent the vectors $|\mathbf{x}\rangle, |\mathbf{z}\rangle \in \mathcal{X} \subseteq \mathbb{R}^d$ and their transposes $\langle \mathbf{x}|, \langle \mathbf{z}|$ using the bra-ket notation.
- ▶ We suppose \mathcal{X} is a vector space over \mathbb{R} or \mathbb{C} with an inner product $\langle \mathbf{z} | \mathbf{x} \rangle$.
- ▶ The quadratic error (or loss) function, given the data set $\mathcal{D} = \{(|\mathbf{y}_m\rangle, t_m)\}_{m=1}^L$ (where $|\mathbf{y}\rangle = |1, \mathbf{x}\rangle$) is defined as

$$\begin{aligned}
 \mathcal{E}(|\mathbf{w}\rangle | \mathcal{D}) &= \frac{1}{L} \sum_{\ell=1}^L \mathcal{E}_{\ell}(|\mathbf{w}\rangle) \\
 &= \frac{1}{2L} \sum_{\ell=1}^L (t_{\ell} - \langle \mathbf{y}_{\ell} | \mathbf{w} \rangle)^2 \\
 &= \frac{1}{2L} \sum_{\ell=1}^L (t_{\ell}^2 - 2t_{\ell} \langle \mathbf{y}_{\ell} | \mathbf{w} \rangle + \langle \mathbf{w} | \mathbf{y}_{\ell} \rangle \langle \mathbf{y}_{\ell} | \mathbf{w} \rangle) \\
 &= \frac{1}{2L} \sum_{\ell=1}^L t_{\ell}^2 - \left(\frac{1}{L} \sum_{\ell=1}^L t_{\ell} \langle \mathbf{y}_{\ell} | \right) |\mathbf{w}\rangle + \frac{1}{2} \left\langle \mathbf{w} \left| \frac{1}{L} \sum_{\ell=1}^L |\mathbf{y}_{\ell}\rangle \langle \mathbf{y}_{\ell}| \right| \mathbf{w} \right\rangle.
 \end{aligned}$$

- ▶ We can identify three terms in the last expression, each term has a component that only depends on the elements of the data set:

$$\mathbb{R}^+ \ni C = \frac{1}{2L} \sum_{\ell=1}^L t_{\ell}^2$$

$$\{1\} \times \mathcal{X} \ni \langle \mathbf{b} | = \frac{1}{L} \sum_{\ell=1}^L t_{\ell} \langle \mathbf{y}_{\ell} |$$

$$\mathcal{X}^{(d+1) \times (d+1)} \ni \mathbf{G} = \frac{1}{L} \sum_{\ell=1}^L |\mathbf{y}_{\ell}\rangle \langle \mathbf{y}_{\ell}|.$$

- ▶ Therefore

$$\mathcal{E}(|\mathbf{w}\rangle | \mathcal{D}) = C - \langle \mathbf{b} | \mathbf{w} \rangle + \frac{1}{2} \langle \mathbf{w} | \mathbf{G} | \mathbf{w} \rangle.$$

- Observe that \mathbf{G} is a symmetric matrix and its eigenvectors and eigenvalues satisfy the following properties:

$$\begin{aligned}\mathbf{G} |\mu\rangle &= \mu |\mu\rangle \\ \langle\mu| \mathbf{G} &= \mu \langle\mu| \\ \langle\mu'| \mathbf{G} |\mu\rangle &= \mu \langle\mu' | \mu\rangle = \mu' \langle\mu' | \mu\rangle \\ (\mu - \mu') \langle\mu' | \mu\rangle &= 0\end{aligned}$$

which implies that $\langle\mu' | \mu\rangle = 0$ for all $\mu \neq \mu'$.

- Also:

$$\begin{aligned}\mu \langle\mu | \mu\rangle &= \langle\mu| \mathbf{G} |\mu\rangle \\ &= \frac{1}{L} \sum_{\ell=1}^L \langle\mu | \mathbf{y}_\ell\rangle \langle\mathbf{y}_\ell | \mu\rangle \\ &= \frac{1}{L} \sum_{\ell=1}^L \langle\mu | \mathbf{y}_\ell\rangle^2 \geq 0\end{aligned}$$

thus $\mu \geq 0$.

Optimization

- ▶ Let us compute the gradient of the loss function:

$$\nabla_{\mathbf{w}} \mathcal{E}(|\mathbf{w}\rangle | \mathcal{D}) = -|\mathbf{b}\rangle + \mathbf{G} |\mathbf{w}\rangle .$$

- ▶ The analytical solution of the problem $\nabla_{\mathbf{w}} \mathcal{E}(|\mathbf{w}\rangle | \mathcal{D}) = \mathbf{0}_{d+1}$ is found if \mathbf{G} admits an inverse (which is true if all its eigenvalues are positive):

$$|\mathbf{w}^*\rangle = \mathbf{G}^{-1} |\mathbf{b}\rangle .$$

- ▶ This solution is not achievable if the condition number is too big $\frac{\mu_{\max}}{\mu_{\min}} \gg 1$, where $\mu_{\max(\min)}$ is the largest (smallest) eigenvalue of \mathbf{G} .

Optimization

- Observation. Assume that the formal expression $|\mathbf{w}^*\rangle = \mathbf{G}^{-1} |\mathbf{b}\rangle$ is given (not computed). Then:

$$\begin{aligned}\mathcal{E}(|\mathbf{w}\rangle|\mathcal{D}) &= C - \langle \mathbf{b} | \mathbf{w} - \mathbf{w}^* + \mathbf{w}^* \rangle + \frac{1}{2} \langle \mathbf{w} - \mathbf{w}^* + \mathbf{w}^* | \mathbf{G} | \mathbf{w} - \mathbf{w}^* + \mathbf{w}^* \rangle \\ &= C - \langle \mathbf{b} | \mathbf{w} - \mathbf{w}^* \rangle - \langle \mathbf{b} | \mathbf{w}^* \rangle + \langle \mathbf{w} - \mathbf{w}^* | \mathbf{G} | \mathbf{w}^* \rangle + \\ &\quad + \frac{1}{2} \langle \mathbf{w} - \mathbf{w}^* | \mathbf{G} | \mathbf{w} - \mathbf{w}^* \rangle \\ &= C - \langle \mathbf{b} | \mathbf{w}^* \rangle + \frac{1}{2} \langle \mathbf{w} - \mathbf{w}^* | \mathbf{G} | \mathbf{w} - \mathbf{w}^* \rangle \\ &= \mathcal{E}(|\mathbf{w}^*\rangle|\mathcal{D}) + \frac{1}{2} \langle \mathbf{w} - \mathbf{w}^* | \mathbf{G} | \mathbf{w} - \mathbf{w}^* \rangle.\end{aligned}$$

where $\mathcal{E}(|\mathbf{w}^*\rangle|\mathcal{D})$ is the residual error of the optimal model.

Gradient Descent Method

- ▶ Let us propose the following iterative method to estimate the solution ($|\mathbf{w}^*\rangle$) of our problem:

$$|\mathbf{w}_{k+1}\rangle = |\mathbf{w}_k\rangle - \alpha_k |\mathbf{g}_k\rangle$$

where

$$|\mathbf{g}_k\rangle = \nabla_{\mathbf{w}} \mathcal{E}(|\mathbf{w}_k\rangle | \mathcal{D}) = -|\mathbf{b}\rangle + \mathbf{G} |\mathbf{w}_k\rangle$$
$$\alpha_k \geq 0$$

are the gradient and the learning rate (or step-size) respectively.

- ▶ We want to find $\alpha \geq 0$ that improves the estimates of the solution, i.e. $\mathcal{E}(|\mathbf{w}_{k+1}\rangle|\mathcal{D}) \leq \mathcal{E}(|\mathbf{w}_k\rangle|\mathcal{D})$.
- ▶ Let us define $\phi(\alpha) = \mathcal{E}(|\mathbf{w}_k\rangle - \alpha|\mathbf{g}_k\rangle|\mathcal{D})$, thus

$$\begin{aligned}
 \phi(\alpha) &= C - \langle \mathbf{b} | (|\mathbf{w}_k\rangle - \alpha|\mathbf{g}_k\rangle) + \\
 &\quad + \frac{1}{2} ((\langle \mathbf{w}_k | - \alpha \langle \mathbf{g}_k |) \mathbf{G} (|\mathbf{w}_k\rangle - \alpha|\mathbf{g}_k\rangle)) \\
 &= \mathcal{E}(|\mathbf{w}_k\rangle|\mathcal{D}) - \alpha \langle \mathbf{g}_k | (\langle \mathbf{b} | - \mathbf{G} |\mathbf{w}_k\rangle) + \frac{\alpha^2}{2} \langle \mathbf{g}_k | \mathbf{G} |\mathbf{g}_k\rangle \\
 &= \phi(0) - \alpha \langle \mathbf{g}_k | \mathbf{g}_k\rangle + \frac{\alpha^2}{2} \langle \mathbf{g}_k | \mathbf{G} |\mathbf{g}_k\rangle
 \end{aligned}$$

► If

$$\frac{2 \langle \mathbf{g}_k | \mathbf{g}_k \rangle}{\langle \mathbf{g}_k | \mathbf{G} | \mathbf{g}_k \rangle} > \alpha > 0$$

then the error function decreases.

► Rayleigh's inequality:

$$\mu_{\min} \langle \boldsymbol{\mu}_{\min} | \boldsymbol{\mu}_{\min} \rangle \leq \langle \mathbf{g}_k | \mathbf{G} | \mathbf{g}_k \rangle \leq \mu_{\max} \langle \boldsymbol{\mu}_{\max} | \boldsymbol{\mu}_{\max} \rangle$$

thus

$$\frac{2}{\mu_{\max}} > \alpha > 0$$

is the condition to ensure convergence.

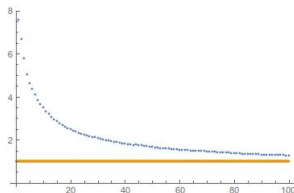
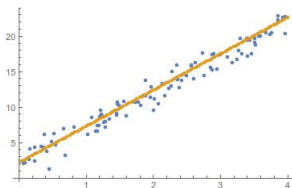
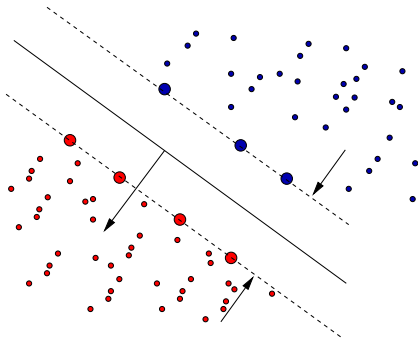


Figure: Linear regression of data points generated by the function $y(x) = 5x + 2 + \varepsilon$, where ε is a random Gaussian noise with zero mean and unit variance (left) and error $\sqrt{\frac{1}{L} \sum_{\ell=1}^L (y_{\ell} - w_0 - w_1 x_{\ell})^2}$ vs iterations (right). Observe that the error converges to the noise's standard deviation. The eigenvalues of $\frac{1}{L} \sum_{\ell=1}^L \mathbf{Y}_{\ell} \mathbf{Y}_{\ell}^T$ are $\mu_{\max} = 6.603 \pm 0.001$ and $\mu_{\min} = 0.196 \pm 0.001$. Thus $\eta \approx 0.03 \pm 0.01$.

Second Example: Classification with a binary perceptron

- Let us suppose we have a set of data presented in the following form:

$$\mathcal{D} = \{(t_\ell; \mathbf{x}_\ell) : t_\ell \in \{-1, +1\}, \mathbf{x}_\ell \in \mathbb{R}^d, \ell = 1, \dots, L\}.$$



Linear separability

Perceptron Learning

- ▶ The choice of the function $\Theta(x)$, where $\Theta(x > 0) = 1$ and 0 otherwise, is well suited for the classification problem.
- ▶ Suppose we have a set $\mathcal{D} = \{(t_\ell; \mathbf{x}_\ell) : t_\ell \in \{-1, +1\}, \mathbf{x}_\ell \in \mathbb{R}^d, \ell = 1, \dots, L\}$, thus there are two classes of vectors \mathbf{x} , those which carry the label $t = +1$ and those with $t = -1$.
- ▶ The error is defined as the total number of elements in \mathcal{D} that have been wrongly classified by \mathbf{w} :

$$E(\mathbf{w}) = \sum_{\ell=1}^L \Theta(-t_\ell \langle \mathbf{x}_\ell | \mathbf{w} \rangle)$$

- ▶ Suppose that there exists $\mathbf{w}^* \in \mathbb{R}^d$ such that $\text{sgn}(t_\ell \langle \mathbf{x}_\ell | \mathbf{w} \rangle) = 1$ for all $1 \leq \ell \leq L$.

Update:

$$|\mathbf{w}_{\ell+1}\rangle = |\mathbf{w}_\ell\rangle + \eta \Theta(-t_\ell \langle \mathbf{x}_\ell | \mathbf{w}_\ell \rangle) t_\ell |\mathbf{x}_\ell\rangle. \quad (1)$$

Perceptron Convergence Theorem

- ▶ Theorem: For any linearly separable set $\mathcal{D} = \{(t_\ell; \mathbf{x}_\ell) : t_\ell \in \{-1, +1\}, \mathbf{x}_\ell \in \mathbb{R}^d, \ell = 1, \dots, L\}$, the learning rule (1) is guaranteed to find a solution in a finite number of steps.
- ▶ Proof: Suppose τ is the total number of real updates

$$\begin{aligned} |\mathbf{w}_{\ell+1}\rangle &= |\mathbf{w}_\ell\rangle + \eta \Theta(-t_\ell \langle \mathbf{x}_\ell | \mathbf{w}_\ell \rangle) t_\ell |\mathbf{x}_\ell\rangle \\ &\vdots \\ &= |\mathbf{w}_0\rangle + \eta \sum_{j=1}^{\ell} \Theta(-t_j \langle \mathbf{x}_j | \mathbf{w}_j \rangle) t_j |\mathbf{x}_j\rangle \\ &= \eta \sum_{j=1}^L \tau_j t_j |\mathbf{x}_j\rangle \end{aligned}$$

where $0 \leq \tau_j$ is the number of times the vector \mathbf{x}_j has been misclassified and $|\mathbf{w}_0\rangle = |\mathbf{0}\rangle$.

- Then

$$\langle \mathbf{w}^* | \mathbf{w}_{\ell+1} \rangle = \eta \sum_{j=1}^L \tau_j t_j \langle \mathbf{w}^* | \mathbf{x}_j \rangle$$

where $t_j \langle \mathbf{w}^* | \mathbf{x}_j \rangle > 0$ because $|\mathbf{w}^*\rangle$ is the solution to the problem (and, therefore, classifies correctly all examples from the data set \mathcal{D}).

- Then

$$\langle \mathbf{w}^* | \mathbf{w}_{\ell+1} \rangle \geq \eta \tau \min_{(\mathbf{x}_j, t_j) \in \mathcal{D}} \{t_j \langle \mathbf{w}^* | \mathbf{x}_j \rangle\},$$

where $\tau = \sum_{j=1}^L \tau_j$.

- Thus

$$\|\mathbf{w}_{\ell+1}\|_2 \geq \eta \tau \min_{(\mathbf{x}_j, t_j) \in \mathcal{D}} \{t_j \langle \mathbf{w}^* | \mathbf{x}_j \rangle\}$$

► Also

$$\begin{aligned}\langle \mathbf{w}_{\ell+1} | \mathbf{w}_{\ell+1} \rangle &= \langle \mathbf{w}_{\ell} | \mathbf{w}_{\ell} \rangle + \eta^2 \Theta(-t_{\ell} \langle \mathbf{x}_{\ell} | \mathbf{w}_{\ell} \rangle) \langle \mathbf{x}_{\ell} | \mathbf{x}_{\ell} \rangle + \\ &\quad + 2\eta \Theta(-t_{\ell} \langle \mathbf{x}_{\ell} | \mathbf{w}_{\ell} \rangle) t_{\ell} \langle \mathbf{w}_{\ell} | \mathbf{x}_{\ell} \rangle,\end{aligned}$$

and observe that the last term is always less than zero.

► Thus

$$\begin{aligned}\langle \mathbf{w}_{\ell+1} | \mathbf{w}_{\ell+1} \rangle &\leq \langle \mathbf{w}_{\ell} | \mathbf{w}_{\ell} \rangle + \eta^2 \Theta(-t_{\ell} \langle \mathbf{x}_{\ell} | \mathbf{w}_{\ell} \rangle) \langle \mathbf{x}_{\ell} | \mathbf{x}_{\ell} \rangle \\ &\leq \langle \mathbf{w}_0 | \mathbf{w}_0 \rangle + \eta^2 \sum_{j=1}^{\ell} \Theta(-t_j \langle \mathbf{x}_j | \mathbf{w}_j \rangle) \langle \mathbf{x}_j | \mathbf{x}_j \rangle \\ &\leq \eta^2 \tau \max_{\mathbf{x}_j \in \mathcal{D}} \{ \langle \mathbf{x}_j | \mathbf{x}_j \rangle \}.\end{aligned}$$

- ▶ Thus

$$\|\mathbf{w}_{\ell+1}\|_2 \leq \eta \sqrt{\tau \max_{\mathbf{x}_j \in \mathcal{D}} \{\langle \mathbf{x}_j | \mathbf{x}_j \rangle\}}.$$

- ▶ Both conditions are satisfied if

$$0 \leq \tau \leq \frac{\max_{\mathbf{x}_j \in \mathcal{D}} \{\langle \mathbf{x}_j | \mathbf{x}_j \rangle\}}{\left(\min_{(\mathbf{x}_j, t_j) \in \mathcal{D}} \{t_j \langle \mathbf{w}^* | \mathbf{x}_j \rangle\} \right)^2} < \infty.$$

- ▶ That indicates that there is a maximum number of mistakes the algorithm could make.