

Fundamentals of Machine Learning

Yassine Laguel

Mail : yassine.laguel@univ-cotedazur.fr

Maximal Inequality for bounded losses

Proposition : Upperboudning the max of bounded r.v.

Let X_1, \dots, X_n be n centered random variables ($\mathbb{E}[X_i] = 0$ for all i), such that $X_i \in [a, b]$ almost surely for all $i \in \{1, \dots, n\}$.

Then,

$$\mathbb{E} \left[\max_{i \in \{1, \dots, n\}} X_i \right] \leq \frac{(b - a)}{2} \sqrt{2 \log(n)}.$$

Maximal Inequality for bounded losses

Proposition : Upperboudning the max of bounded r.v.

Let X_1, \dots, X_n be n centered random variables ($\mathbb{E}[X_i] = 0$ for all i), such that $X_i \in [a, b]$ almost surely for all $i \in \{1, \dots, n\}$.

Then,

$$\mathbb{E} \left[\max_{i \in \{1, \dots, n\}} X_i \right] \leq \frac{(b - a)}{2} \sqrt{2 \log(n)}.$$

Proof : on the black board.

Maximal Inequality for bounded losses

Proposition : Uppertbounding the max of bounded r.v.

Let X_1, \dots, X_n be n centered random variables ($\mathbb{E}[X_i] = 0$ for all i), such that $X_i \in [a, b]$ almost surely for all $i \in \{1, \dots, n\}$.

Then,

$$\mathbb{E} \left[\max_{i \in \{1, \dots, n\}} X_i \right] \leq \frac{(b - a)}{2} \sqrt{2 \log(n)}.$$

Proof : on the black board.

Remark :

- This bound is trivial for n large enough, since all variables are almost surely upperbounded by b .

Maximal Inequality for bounded losses

Proposition : Uppertbounding the max of bounded r.v.

Let X_1, \dots, X_n be n centered random variables ($\mathbb{E}[X_i] = 0$ for all i), such that $X_i \in [a, b]$ almost surely for all $i \in \{1, \dots, n\}$.

Then,

$$\mathbb{E} \left[\max_{i \in \{1, \dots, n\}} X_i \right] \leq \frac{(b - a)}{2} \sqrt{2 \log(n)}.$$

Proof : on the black board.

Remark :

- This bound is trivial for n large enough, since all variables are almost surely upperbounded by b .
- Note that we did not need any independence between the X'_i s.

Maximal Inequality for bounded losses

Proposition : Uppertbounding the max of bounded r.v.

Let X_1, \dots, X_n be n centered random variables ($\mathbb{E}[X_i] = 0$ for all i), such that $X_i \in [a, b]$ almost surely for all $i \in \{1, \dots, n\}$.

Then,

$$\mathbb{E} \left[\max_{i \in \{1, \dots, n\}} X_i \right] \leq \frac{(b - a)}{2} \sqrt{2 \log(n)}.$$

Proof : on the black board.

Remark :

- This bound is trivial for n large enough, since all variables are almost surely upperbounded by b .
- Note that we did not need any independence between the X'_i s.

Proposition : Maximal inequality for bounded losses

Assume the loss ℓ satisfies $0 \leq \ell \leq 1$.

Then,

$$\mathbb{E} \left[\max_{f \in \mathcal{H}} r(f) - R(f) \right] \leq \sqrt{\frac{2 \log(|\mathcal{H}|)}{n}}.$$

Maximal Inequality for bounded losses

Proposition : Upperbounding the max of bounded r.v.

Let X_1, \dots, X_n be n centered random variables ($\mathbb{E}[X_i] = 0$ for all i), such that $X_i \in [a, b]$ almost surely for all $i \in \{1, \dots, n\}$.

Then,

$$\mathbb{E} \left[\max_{i \in \{1, \dots, n\}} X_i \right] \leq \frac{(b - a)}{2} \sqrt{2 \log(n)}.$$

Proof : on the black board.

Remark :

- This bound is trivial for n large enough, since all variables are almost surely upperbounded by b .
- Note that we did not need any independence between the X'_i s.

Proposition : Maximal inequality for bounded losses

Assume the loss ℓ satisfies $0 \leq \ell \leq 1$.

Then,

$$\mathbb{E} \left[\max_{f \in \mathcal{H}} r(f) - R(f) \right] \leq \sqrt{\frac{2 \log(|\mathcal{H}|)}{n}}.$$

Proof : on the black board.

Maximal Inequality for bounded losses

Proposition : Upperbounding the max of bounded r.v.

Let X_1, \dots, X_n be n centered random variables ($\mathbb{E}[X_i] = 0$ for all i), such that $X_i \in [a, b]$ almost surely for all $i \in \{1, \dots, n\}$.

Then,

$$\mathbb{E} \left[\max_{i \in \{1, \dots, n\}} X_i \right] \leq \frac{(b - a)}{2} \sqrt{2 \log(n)}.$$

Proof : on the black board.

Remark :

- This bound is trivial for n large enough, since all variables are almost surely upperbounded by b .
- Note that we did not need any independence between the X'_i s.

Proposition : Maximal inequality for bounded losses

Assume the loss ℓ satisfies $0 \leq \ell \leq 1$.

Then,

$$\mathbb{E} \left[\max_{f \in \mathcal{H}} r(f) - R(f) \right] \leq \sqrt{\frac{2 \log(|\mathcal{H}|)}{n}}.$$

Proof : on the black board.

Remarks :

- Taking $\ell \leq c$ with $c > 0$ turns the final bound into $c \sqrt{2 \log(|\mathcal{H}|)/n}$.
→ Exercise

Maximal Inequality for bounded losses

Proposition : Upperbounding the max of bounded r.v.

Let X_1, \dots, X_n be n centered random variables ($\mathbb{E}[X_i] = 0$ for all i), such that $X_i \in [a, b]$ almost surely for all $i \in \{1, \dots, n\}$.

Then,

$$\mathbb{E} \left[\max_{i \in \{1, \dots, n\}} X_i \right] \leq \frac{(b - a)}{2} \sqrt{2 \log(n)}.$$

Proof : on the black board.

Remark :

- This bound is trivial for n large enough, since all variables are almost surely upperbounded by b .
- Note that we did not need any independence between the X'_i s.

Proposition : Maximal inequality for bounded losses

Assume the loss ℓ satisfies $0 \leq \ell \leq 1$.

Then,

$$\mathbb{E} \left[\max_{f \in \mathcal{H}} r(f) - R(f) \right] \leq \sqrt{\frac{2 \log(|\mathcal{H}|)}{n}}.$$

Proof : on the black board.

Remarks :

- Taking $\ell \leq c$ with $c > 0$ turns the final bound into $c \sqrt{2 \log(|\mathcal{H}|)/n}$.
→ Exercise
- Remember that our general goal was to establish a bound of the form

$$\mathbb{E} \left[\sup_{f \in \mathcal{H}} r(f) - R(f) \right] \leq \frac{\psi(\mathcal{H})}{n^\alpha}$$

Maximal Inequality for bounded losses

Proposition : Upperbounding the max of bounded r.v.

Let X_1, \dots, X_n be n centered random variables ($\mathbb{E}[X_i] = 0$ for all i), such that $X_i \in [a, b]$ almost surely for all $i \in \{1, \dots, n\}$.

Then,

$$\mathbb{E} \left[\max_{i \in \{1, \dots, n\}} X_i \right] \leq \frac{(b - a)}{2} \sqrt{2 \log(n)}.$$

Proof : on the black board.

Remark :

- This bound is trivial for n large enough, since all variables are almost surely upperbounded by b .
- Note that we did not need any independence between the X'_i s.

Proposition : Maximal inequality for bounded losses

Assume the loss ℓ satisfies $0 \leq \ell \leq 1$.

Then,

$$\mathbb{E} \left[\max_{f \in \mathcal{H}} r(f) - R(f) \right] \leq \sqrt{\frac{2 \log(|\mathcal{H}|)}{n}}.$$

Proof : on the black board.

Remarks :

- Taking $\ell \leq c$ with $c > 0$ turns the final bound into $c \sqrt{2 \log(|\mathcal{H}|)/n}$.
→ Exercise
- Remember that our general goal was to establish a bound of the form

$$\mathbb{E} \left[\sup_{f \in \mathcal{H}} r(f) - R(f) \right] \leq \frac{\psi(\mathcal{H})}{n^\alpha}$$

- The complexity of the Hypothesis class, captured by $|\mathcal{H}|$ plays a **logarithmic** role in our final statistical bound.

Maximal Inequality for bounded losses

Proposition : Upperboudning the max of bounded r.v.

Let X_1, \dots, X_n be n centered random variables ($\mathbb{E}[X_i] = 0$ for all i), such that $X_i \in [a, b]$ almost surely for all $i \in \{1, \dots, n\}$.

Then,

$$\mathbb{E} \left[\max_{i \in \{1, \dots, n\}} X_i \right] \leq \frac{(b - a)}{2} \sqrt{2 \log(n)}.$$

Proof : on the black board.

Remark :

- This bound is trivial for n large enough, since all variables are almost surely upperbounded by b .
- Note that we did not need any independence between the X'_i s.

Proposition : Maximal inequality for bounded losses

Assume the loss ℓ satisfies $0 \leq \ell \leq 1$.

Then,

$$\mathbb{E} \left[\max_{f \in \mathcal{H}} r(f) - R(f) \right] \leq \sqrt{\frac{2 \log(|\mathcal{H}|)}{n}}.$$

Proof : on the black board.

Remarks :

- Taking $\ell \leq c$ with $c > 0$ turns the final bound into $c \sqrt{2 \log(|\mathcal{H}|)/n}$.
→ Exercise
- Remember that our general goal was to establish a bound of the form

$$\mathbb{E} \left[\sup_{f \in \mathcal{H}} r(f) - R(f) \right] \leq \frac{\psi(\mathcal{H})}{n^\alpha}$$

- The complexity of the Hypothesis class, captured by $|\mathcal{H}|$ plays a **logarithmic** role in our final statistical bound.
- The dependency on n scales as $\frac{1}{n^\alpha}$ with $\alpha = 0.5$.
This will be a common characteristic of our upcoming results.

Maximal Inequality for bounded losses

Proposition : Upperbounding the max of bounded r.v.

Let X_1, \dots, X_n be n centered random variables ($\mathbb{E}[X_i] = 0$ for all i), such that $X_i \in [a, b]$ almost surely for all $i \in \{1, \dots, n\}$.

Then,

$$\mathbb{E} \left[\max_{i \in \{1, \dots, n\}} X_i \right] \leq \frac{(b - a)}{2} \sqrt{2 \log(n)}.$$

Proof : on the black board.

Remark :

- This bound is trivial for n large enough, since all variables are almost surely upperbounded by b .
- Note that we did not need any independence between the X'_i s.

Proposition : Maximal inequality for bounded losses

Assume the loss ℓ satisfies $0 \leq \ell \leq 1$.

Then,

$$\mathbb{E} \left[\max_{f \in \mathcal{H}} r(f) - R(f) \right] \leq \sqrt{\frac{2 \log(|\mathcal{H}|)}{n}}.$$

Proof : on the black board.

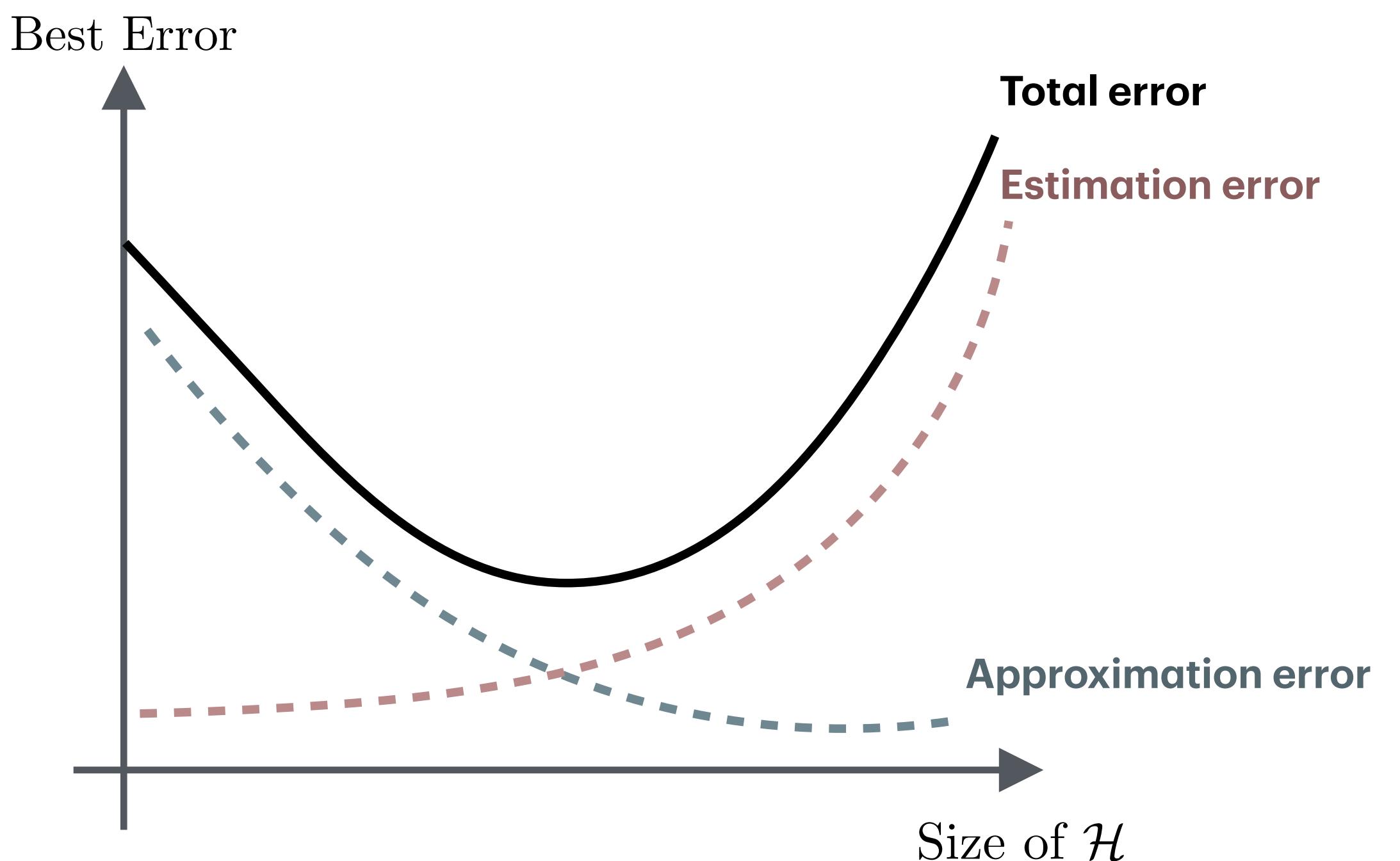
Remarks :

- Taking $\ell \leq c$ with $c > 0$ turns the final bound into $c \sqrt{2 \log(|\mathcal{H}|)/n}$.
→ Exercise
- Remember that our general goal was to establish a bound of the form

$$\mathbb{E} \left[\sup_{f \in \mathcal{H}} r(f) - R(f) \right] \leq \frac{\psi(\mathcal{H})}{n^\alpha}$$

- The complexity of the Hypothesis class, captured by $|\mathcal{H}|$ plays a **logarithmic** role in our final statistical bound.
- The dependency on n scales as $\frac{1}{n^\alpha}$ with $\alpha = 0.5$.
This will be a common characteristic of our upcoming results.
- When $|\mathcal{H}|$ is infinite, this bound is not exploitable.

1. Supervised Learning Setting
2. Estimation vs Approximation
3. Maximal inequalities
4. Rademacher Complexity



Rademacher Complexity of a set

- Rademacher variables and complexity of a set

Rademacher Complexity of a set

■ Rademacher variables and complexity of a set

Definition : Rademacher random variables

A Rademacher random variable is a random variable $X : \Omega \rightarrow \{-1, 1\}$

such that

$$\mathbb{P}[X = 1] = \mathbb{P}[X = -1] = 1/2.$$

Rademacher Complexity of a set

■ Rademacher variables and complexity of a set

Definition : Rademacher random variables

A Rademacher random variable is a random variable $X : \Omega \rightarrow \{-1, 1\}$
such that

$$\mathbb{P}[X = 1] = \mathbb{P}[X = -1] = 1/2.$$

Remark :



Rademacher Complexity of a set

■ Rademacher variables and complexity of a set

Definition : Rademacher random variables

A Rademacher random variable is a random variable $X : \Omega \rightarrow \{-1, 1\}$ such that

$$\mathbb{P}[X = 1] = \mathbb{P}[X = -1] = 1/2.$$

Remark :

- If $U \sim \mathcal{B}(1/2)$, then $X \triangleq 2U - 1$ is a Rademacher variable.

Definition : Rademacher complexity of a set

Let $\Omega_1, \dots, \Omega_n$ be n independent Rademacher variables.

The Rademacher complexity of a set $T \subset \mathbb{R}^n$, denoted $\text{Rad}(T)$, is defined as:

$$\text{Rad}(T) = \mathbb{E} \left[\sup_{t \in T} \frac{1}{n} \sum_{i=1}^n \Omega_i t_i \right].$$

Rademacher Complexity of a set

■ Rademacher variables and complexity of a set

Definition : Rademacher random variables

A Rademacher random variable is a random variable $X : \Omega \rightarrow \{-1, 1\}$

such that

$$\mathbb{P}[X = 1] = \mathbb{P}[X = -1] = 1/2.$$

Remark :

- If $U \sim \mathcal{B}(1/2)$, then $X \triangleq 2U - 1$ is a Rademacher variable.

Definition : Rademacher complexity of a set

Let $\Omega_1, \dots, \Omega_n$ be n independent Rademacher variables.

The Rademacher complexity of a set $T \subset \mathbb{R}^n$, denoted $\text{Rad}(T)$, is defined as:

$$\text{Rad}(T) = \mathbb{E} \left[\sup_{t \in T} \frac{1}{n} \sum_{i=1}^n \Omega_i t_i \right].$$

Remarks :

- The quantity $\text{Rad}(T)$ measures the complexity of T , as $\sup_{t \in T} \frac{1}{n} \sum_{i=1}^n \Omega_i t_i$

Rademacher Complexity of a set

■ Rademacher variables and complexity of a set

- describes how well elements in T can replicate the sign pattern of a random signal $(\Omega_1, \dots, \Omega_n)$.

Definition : Rademacher random variables

A Rademacher random variable is a random variable $X : \Omega \rightarrow \{-1, 1\}$ such that

$$\mathbb{P}[X = 1] = \mathbb{P}[X = -1] = 1/2.$$

Remark :

- If $U \sim \mathcal{B}(1/2)$, then $X \triangleq 2U - 1$ is a Rademacher variable.

Definition : Rademacher complexity of a set

Let $\Omega_1, \dots, \Omega_n$ be n independent Rademacher variables.

The Rademacher complexity of a set $T \subset \mathbb{R}^n$, denoted $\text{Rad}(T)$, is defined as:

$$\text{Rad}(T) = \mathbb{E} \left[\sup_{t \in T} \frac{1}{n} \sum_{i=1}^n \Omega_i t_i \right].$$

Remarks :

- The quantity $\text{Rad}(T)$ measures the complexity of T , as $\sup_{t \in T} \frac{1}{n} \sum_{i=1}^n \Omega_i t_i$

Rademacher Complexity of a set

■ Rademacher variables and complexity of a set

Definition : Rademacher random variables

A Rademacher random variable is a random variable $X : \Omega \rightarrow \{-1, 1\}$ such that

$$\mathbb{P}[X = 1] = \mathbb{P}[X = -1] = 1/2.$$

Remark :

- If $U \sim \mathcal{B}(1/2)$, then $X \triangleq 2U - 1$ is a Rademacher variable.

Definition : Rademacher complexity of a set

Let $\Omega_1, \dots, \Omega_n$ be n independent Rademacher variables.

The Rademacher complexity of a set $T \subset \mathbb{R}^n$, denoted $\text{Rad}(T)$, is defined as:

$$\text{Rad}(T) = \mathbb{E} \left[\sup_{t \in T} \frac{1}{n} \sum_{i=1}^n \Omega_i t_i \right].$$

Remarks :

- The quantity $\text{Rad}(T)$ measures the complexity of T , as $\sup_{t \in T} \frac{1}{n} \sum_{i=1}^n \Omega_i t_i$

- describes how well elements in T can replicate the sign pattern of a random signal $(\Omega_1, \dots, \Omega_n)$.

■ First properties

Rademacher Complexity of a set

■ Rademacher variables and complexity of a set

Definition : Rademacher random variables

A Rademacher random variable is a random variable $X : \Omega \rightarrow \{-1, 1\}$ such that

$$\mathbb{P}[X = 1] = \mathbb{P}[X = -1] = 1/2.$$

Remark :

- If $U \sim \mathcal{B}(1/2)$, then $X \triangleq 2U - 1$ is a Rademacher variable.

Definition : Rademacher complexity of a set

Let $\Omega_1, \dots, \Omega_n$ be n independent Rademacher variables.

The Rademacher complexity of a set $T \subset \mathbb{R}^n$, denoted $\text{Rad}(T)$, is defined as:

$$\text{Rad}(T) = \mathbb{E} \left[\sup_{t \in T} \frac{1}{n} \sum_{i=1}^n \Omega_i t_i \right].$$

Remarks :

- The quantity $\text{Rad}(T)$ measures the complexity of T , as $\sup_{t \in T} \frac{1}{n} \sum_{i=1}^n \Omega_i t_i$

- describes how well elements in T can replicate the sign pattern of a random signal $(\Omega_1, \dots, \Omega_n)$.

■ First properties

Proposition : Homotetie and translation

Let $T \subset \mathbb{R}^n, v \in \mathbb{R}^n, c \in \mathbb{R}$, and define $cT + v = \{ct + v, t \in T\}$. Then,

$$\text{Rad}(cT + v) = |c| \text{Rad}(T).$$

Rademacher Complexity of a set

■ Rademacher variables and complexity of a set

Definition : Rademacher random variables

A Rademacher random variable is a random variable $X : \Omega \rightarrow \{-1, 1\}$ such that

$$\mathbb{P}[X = 1] = \mathbb{P}[X = -1] = 1/2.$$

Remark :

- If $U \sim \mathcal{B}(1/2)$, then $X \triangleq 2U - 1$ is a Rademacher variable.

Definition : Rademacher complexity of a set

Let $\Omega_1, \dots, \Omega_n$ be n independent Rademacher variables.

The Rademacher complexity of a set $T \subset \mathbb{R}^n$, denoted $\text{Rad}(T)$, is defined as:

$$\text{Rad}(T) = \mathbb{E} \left[\sup_{t \in T} \frac{1}{n} \sum_{i=1}^n \Omega_i t_i \right].$$

Remarks :

- The quantity $\text{Rad}(T)$ measures the complexity of T , as $\sup_{t \in T} \frac{1}{n} \sum_{i=1}^n \Omega_i t_i$

- describes how well elements in T can replicate the sign pattern of a random signal $(\Omega_1, \dots, \Omega_n)$.

■ First properties

Proposition : Homotetie and translation

Let $T \subset \mathbb{R}^n, v \in \mathbb{R}^n, c \in \mathbb{R}$, and define $cT + v = \{ct + v, t \in T\}$.

Then,

$$\text{Rad}(cT + v) = |c| \text{Rad}(T).$$

Proof : on the black board.

Rademacher Complexity of a set

■ Rademacher variables and complexity of a set

Definition : Rademacher random variables

A Rademacher random variable is a random variable $X : \Omega \rightarrow \{-1, 1\}$ such that

$$\mathbb{P}[X = 1] = \mathbb{P}[X = -1] = 1/2.$$

Remark :

- If $U \sim \mathcal{B}(1/2)$, then $X \triangleq 2U - 1$ is a Rademacher variable.

Definition : Rademacher complexity of a set

Let $\Omega_1, \dots, \Omega_n$ be n independent Rademacher variables.

The Rademacher complexity of a set $T \subset \mathbb{R}^n$, denoted $\text{Rad}(T)$, is defined as:

$$\text{Rad}(T) = \mathbb{E} \left[\sup_{t \in T} \frac{1}{n} \sum_{i=1}^n \Omega_i t_i \right].$$

Remarks :

- The quantity $\text{Rad}(T)$ measures the complexity of T , as $\sup_{t \in T} \frac{1}{n} \sum_{i=1}^n \Omega_i t_i$

- describes how well elements in T can replicate the sign pattern of a random signal $(\Omega_1, \dots, \Omega_n)$.

■ First properties

Proposition : Homotetie and translation

Let $T \subset \mathbb{R}^n, v \in \mathbb{R}^n, c \in \mathbb{R}$, and define $cT + v = \{ct + v, t \in T\}$. Then,

$$\text{Rad}(cT + v) = |c| \text{Rad}(T).$$

Proof : on the black board.

Proposition : Minkowski sum

For any sets $T, T' \subset \mathbb{R}^d$, we have

$$\text{Rad}(T + T') = \text{Rad}(T) + \text{Rad}(T').$$

where $T + T' = \{t + t', t \in T, t' \in T'\}$.

Rademacher Complexity of a set

■ Rademacher variables and complexity of a set

Definition : Rademacher random variables

A Rademacher random variable is a random variable $X : \Omega \rightarrow \{-1, 1\}$ such that

$$\mathbb{P}[X = 1] = \mathbb{P}[X = -1] = 1/2.$$

Remark :

- If $U \sim \mathcal{B}(1/2)$, then $X \triangleq 2U - 1$ is a Rademacher variable.

Definition : Rademacher complexity of a set

Let $\Omega_1, \dots, \Omega_n$ be n independent Rademacher variables.

The Rademacher complexity of a set $T \subset \mathbb{R}^n$, denoted $\text{Rad}(T)$, is defined as:

$$\text{Rad}(T) = \mathbb{E} \left[\sup_{t \in T} \frac{1}{n} \sum_{i=1}^n \Omega_i t_i \right].$$

Remarks :

- The quantity $\text{Rad}(T)$ measures the complexity of T , as $\sup_{t \in T} \frac{1}{n} \sum_{i=1}^n \Omega_i t_i$

- describes how well elements in T can replicate the sign pattern of a random signal $(\Omega_1, \dots, \Omega_n)$.

■ First properties

Proposition : Homotetie and translation

Let $T \subset \mathbb{R}^n, v \in \mathbb{R}^n, c \in \mathbb{R}$, and define $cT + v = \{ct + v, t \in T\}$. Then,

$$\text{Rad}(cT + v) = |c| \text{Rad}(T).$$

Proof : on the black board.

Proposition : Minkowski sum

For any sets $T, T' \subset \mathbb{R}^d$, we have

$$\text{Rad}(T + T') = \text{Rad}(T) + \text{Rad}(T').$$

where $T + T' = \{t + t', t \in T, t' \in T'\}$.

Proof : on the black board.

More properties

Proposition : Convex hull

Let $T \subset \mathbb{R}^n, v \in \mathbb{R}^n, c \in \mathbb{R}$, and define $cT + v = \{ct + v, t \in T\}$.

Then,

$$\text{Rad}(\text{conv}(T)) = \text{Rad}(T).$$

More properties

Proposition : Convex hull

Let $T \subset \mathbb{R}^n, v \in \mathbb{R}^n, c \in \mathbb{R}$, and define $cT + v = \{ct + v, t \in T\}$.

Then,

$$\text{Rad}(\text{conv}(T)) = \text{Rad}(T).$$

Proof : on the black board.

More properties

Proposition : Convex hull

Let $T \subset \mathbb{R}^n, v \in \mathbb{R}^n, c \in \mathbb{R}$, and define $cT + v = \{ct + v, t \in T\}$.

Then,

$$\text{Rad}(\text{conv}(T)) = \text{Rad}(T).$$

Proof : on the black board.

Proposition : Massart's Lemma - Complexity for finite Sets

Let $T \subset \mathbb{R}^n$ be finite.

Then,

$$\text{Rad}(T) \leq \max_{t \in T} \|t\|_2 \frac{\sqrt{2 \log(|T|)}}{n}.$$

More properties

Proposition : Convex hull

Let $T \subset \mathbb{R}^n, v \in \mathbb{R}^n, c \in \mathbb{R}$, and define $cT + v = \{ct + v, t \in T\}$.

Then,

$$\text{Rad}(\text{conv}(T)) = \text{Rad}(T).$$

Proof : on the black board.

Proposition : Massart's Lemma - Complexity for finite Sets

Let $T \subset \mathbb{R}^n$ be finite.

Then,

$$\text{Rad}(T) \leq \max_{t \in T} \|t\|_2 \frac{\sqrt{2 \log(|T|)}}{n}.$$

Proof : on the black board.

More properties

Proposition : Convex hull

Let $T \subset \mathbb{R}^n, v \in \mathbb{R}^n, c \in \mathbb{R}$, and define $cT + v = \{ct + v, t \in T\}$.

Then,

$$\text{Rad}(\text{conv}(T)) = \text{Rad}(T).$$

Proof : on the black board.

Proposition : Massart's Lemma - Complexity for finite Sets

Let $T \subset \mathbb{R}^n$ be finite.

Then,

$$\text{Rad}(T) \leq \max_{t \in T} \|t\|_2 \frac{\sqrt{2 \log(|T|)}}{n}.$$

Proof : on the black board.

Remark :

- The two previous propositions allow to upper bound the Rademacher complexity of the containing polyhedra.

More properties

Proposition : Convex hull

Let $T \subset \mathbb{R}^n, v \in \mathbb{R}^n, c \in \mathbb{R}$, and define $cT + v = \{ct + v, t \in T\}$.

Then,

$$\text{Rad}(\text{conv}(T)) = \text{Rad}(T).$$

Proof : on the black board.

Proposition : Massart's Lemma - Complexity for finite Sets

Let $T \subset \mathbb{R}^n$ be finite.

Then,

$$\text{Rad}(T) \leq \max_{t \in T} \|t\|_2 \frac{\sqrt{2 \log(|T|)}}{n}.$$

Proof : on the black board.

Remark :

- The two previous propositions allow to upper bound the Rademacher complexity of the containing polyhedra.



More properties

Proposition : Convex hull

Let $T \subset \mathbb{R}^n$, $v \in \mathbb{R}^n$, $c \in \mathbb{R}$, and define $cT + v = \{ct + v, t \in T\}$.

Then,

$$\text{Rad}(\text{conv}(T)) = \text{Rad}(T).$$

Proof : on the black board.

Proposition : Massart's Lemma - Complexity for finite Sets

Let $T \subset \mathbb{R}^n$ be finite.

Then,

$$\text{Rad}(T) \leq \max_{t \in T} \|t\|_2 \frac{\sqrt{2 \log(|T|)}}{n}.$$

Proof : on the black board.

Remark :

- The two previous propositions allow to upper bound the Rademacher complexity of the containing polyhedra.



Proposition : Talagrand's lemma

Let $f_1, \dots, f_n : \mathbb{R} \rightarrow \mathbb{R}$ be n G -Lipchitz functions.

Let $T \subset \mathbb{R}^n$ and define $(f_1, \dots, f_n) \circ T$ as

$$(f_1, \dots, f_n) \circ T = \{(f_1(t_1), \dots, f_n(t_n)), t \in T\}.$$

Then,

$$\text{Rad}((f_1, \dots, f_n) \circ T) \leq G \text{Rad}(T).$$

More properties

Proposition : Convex hull

Let $T \subset \mathbb{R}^n$, $v \in \mathbb{R}^n$, $c \in \mathbb{R}$, and define $cT + v = \{ct + v, t \in T\}$.

Then,

$$\text{Rad}(\text{conv}(T)) = \text{Rad}(T).$$

Proof : on the black board.

Proposition : Massart's Lemma - Complexity for finite Sets

Let $T \subset \mathbb{R}^n$ be finite.

Then,

$$\text{Rad}(T) \leq \max_{t \in T} \|t\|_2 \frac{\sqrt{2 \log(|T|)}}{n}.$$

Proof : on the black board.

Remark :

- The two previous propositions allow to upper bound the Rademacher complexity of the containing polyhedra.



Proposition : Talagrand's lemma

Let $f_1, \dots, f_n : \mathbb{R} \rightarrow \mathbb{R}$ be n G -lipchitz functions.

Let $T \subset \mathbb{R}^n$ and define $(f_1, \dots, f_n) \circ T$ as

$$(f_1, \dots, f_n) \circ T = \{(f_1(t_1), \dots, f_n(t_n)), t \in T\}.$$

Then,

$$\text{Rad}((f_1, \dots, f_n) \circ T) \leq G \text{Rad}(T).$$

Proof : on the black board.

More properties

Proposition : Convex hull

Let $T \subset \mathbb{R}^n$, $v \in \mathbb{R}^n$, $c \in \mathbb{R}$, and define $cT + v = \{ct + v, t \in T\}$.

Then,

$$\text{Rad}(\text{conv}(T)) = \text{Rad}(T).$$

Proof : on the black board.

Proposition : Massart's Lemma - Complexity for finite Sets

Let $T \subset \mathbb{R}^n$ be finite.

Then,

$$\text{Rad}(T) \leq \max_{t \in T} \|t\|_2 \frac{\sqrt{2 \log(|T|)}}{n}.$$

Proof : on the black board.

Remark :

- The two previous propositions allow to upper bound the Rademacher complexity of the containing polyhedra.



Proposition : Talagrand's lemma

Let $f_1, \dots, f_n : \mathbb{R} \rightarrow \mathbb{R}$ be n G -Lipschitz functions.

Let $T \subset \mathbb{R}^n$ and define $(f_1, \dots, f_n) \circ T$ as

$$(f_1, \dots, f_n) \circ T = \{(f_1(t_1), \dots, f_n(t_n)), t \in T\}.$$

Then,

$$\text{Rad}((f_1, \dots, f_n) \circ T) \leq G \text{Rad}(T).$$

Proof : on the black board.

Remark :

- This lemma will later be used to establish statistical bounds when the loss ℓ satisfies a Lipschitz condition.

Rademacher complexity of a hypothesis class

■ Formal definition & Symmetrization Lemma

Definition : Rademacher complexity of a class of functions

Let $(\Omega_1, \dots, \Omega_n)$ be n i.i.d. Rademacher variables.

We define the empirical Rademacher complexity of the hypothesis class \mathcal{H} , denoted $\widetilde{\text{Rad}}(\mathcal{H})$, as

$$\widetilde{\text{Rad}}(\mathcal{H}) = \mathbb{E} \left[\sup_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^N \Omega_i \ell(f(X_i), Y_i) \mid (X_1, Y_1), \dots, (X_n, Y_n) \right],$$

and the Rademacher complexity of the class \mathcal{H} , denoted $\text{Rad}(\mathcal{H})$, as:

$$\text{Rad}(\mathcal{H}) = \mathbb{E}[\widetilde{\text{Rad}}(\mathcal{H})] = \mathbb{E} \left[\sup_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^N \Omega_i \ell(f(X_i), Y_i) \right].$$

Rademacher complexity of a hypothesis class

■ Formal definition & Symmetrization Lemma

Definition : Rademacher complexity of a class of functions

Let $(\Omega_1, \dots, \Omega_n)$ be n i.i.d. Rademacher variables.

We define the empirical Rademacher complexity of the hypothesis class \mathcal{H} , denoted $\widetilde{\text{Rad}}(\mathcal{H})$, as

$$\widetilde{\text{Rad}}(\mathcal{H}) = \mathbb{E} \left[\sup_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^N \Omega_i \ell(f(X_i), Y_i) \mid (X_1, Y_1), \dots, (X_n, Y_n) \right],$$

and the Rademacher complexity of the class \mathcal{H} , denoted $\text{Rad}(\mathcal{H})$, as:

$$\text{Rad}(\mathcal{H}) = \mathbb{E}[\widetilde{\text{Rad}}(\mathcal{H})] = \mathbb{E} \left[\sup_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^N \Omega_i \ell(f(X_i), Y_i) \right].$$

Remarks :

- We will use from now on the Rademacher complexity to establish as a measure of complexity of the hypothesis class \mathcal{H} .

Rademacher complexity of a hypothesis class

■ Formal definition & Symmetrization Lemma

Definition : Rademacher complexity of a class of functions

Let $(\Omega_1, \dots, \Omega_n)$ be n i.i.d. Rademacher variables.

We define the empirical Rademacher complexity of the hypothesis class \mathcal{H} , denoted $\widetilde{\text{Rad}}(\mathcal{H})$, as

$$\widetilde{\text{Rad}}(\mathcal{H}) = \mathbb{E} \left[\sup_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^N \Omega_i \ell(f(X_i), Y_i) \mid (X_1, Y_1), \dots, (X_n, Y_n) \right],$$

and the Rademacher complexity of the class \mathcal{H} , denoted $\text{Rad}(\mathcal{H})$, as:

$$\text{Rad}(\mathcal{H}) = \mathbb{E}[\widetilde{\text{Rad}}(\mathcal{H})] = \mathbb{E} \left[\sup_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^N \Omega_i \ell(f(X_i), Y_i) \right].$$

Remarks :

- We will use from now on the Rademacher complexity to establish as a measure of complexity of the hypothesis class \mathcal{H} .
- A fundamental result in this direction is the symmetrisation lemma.

Rademacher complexity of a hypothesis class

■ Formal definition & Symmetrization Lemma

Definition : Rademacher complexity of a class of functions

Let $(\Omega_1, \dots, \Omega_n)$ be n i.i.d. Rademacher variables.

We define the empirical Rademacher complexity of the hypothesis class \mathcal{H} , denoted $\widetilde{\text{Rad}}(\mathcal{H})$, as

$$\widetilde{\text{Rad}}(\mathcal{H}) = \mathbb{E} \left[\sup_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^N \Omega_i \ell(f(X_i), Y_i) \mid (X_1, Y_1), \dots, (X_n, Y_n) \right],$$

and the Rademacher complexity of the class \mathcal{H} , denoted $\text{Rad}(\mathcal{H})$, as:

$$\text{Rad}(\mathcal{H}) = \mathbb{E}[\widetilde{\text{Rad}}(\mathcal{H})] = \mathbb{E} \left[\sup_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^N \Omega_i \ell(f(X_i), Y_i) \right].$$

Remarks :

- We will use from now on the Rademacher complexity to establish as a measure of complexity of the hypothesis class \mathcal{H} .
- A fundamental result in this direction is the symmetrisation lemma.

Proposition : Symmetrization lemma

Under the Empirical Risk Minimization Framework, we have :

$$\mathbb{E} \left[\sup_{f \in \mathcal{H}} \{r(f) - R(f)\} \right] \leq 2 \text{Rad}(\mathcal{H}).$$

Rademacher complexity of a hypothesis class

■ Formal definition & Symmetrization Lemma

Definition : Rademacher complexity of a class of functions

Let $(\Omega_1, \dots, \Omega_n)$ be n i.i.d. Rademacher variables.

We define the empirical Rademacher complexity of the hypothesis class \mathcal{H} , denoted $\widetilde{\text{Rad}}(\mathcal{H})$, as

$$\widetilde{\text{Rad}}(\mathcal{H}) = \mathbb{E} \left[\sup_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^N \Omega_i \ell(f(X_i), Y_i) \mid (X_1, Y_1), \dots, (X_n, Y_n) \right],$$

and the Rademacher complexity of the class \mathcal{H} , denoted $\text{Rad}(\mathcal{H})$, as:

$$\text{Rad}(\mathcal{H}) = \mathbb{E}[\widetilde{\text{Rad}}(\mathcal{H})] = \mathbb{E} \left[\sup_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^N \Omega_i \ell(f(X_i), Y_i) \right].$$

Remarks :

- We will use from now on the Rademacher complexity to establish as a measure of complexity of the hypothesis class \mathcal{H} .
- A fundamental result in this direction is the symmetrisation lemma.

Proposition : Symmetrization lemma

Under the Empirical Risk Minimization Framework, we have :

$$\mathbb{E} \left[\sup_{f \in \mathcal{H}} \{r(f) - R(f)\} \right] \leq 2 \text{Rad}(\mathcal{H}).$$

Proof : on the black board.

Rademacher complexity of a hypothesis class

■ Formal definition & Symmetrization Lemma

Definition : Rademacher complexity of a class of functions

Let $(\Omega_1, \dots, \Omega_n)$ be n i.i.d. Rademacher variables.

We define the empirical Rademacher complexity of the hypothesis class \mathcal{H} , denoted $\widetilde{\text{Rad}}(\mathcal{H})$, as

$$\widetilde{\text{Rad}}(\mathcal{H}) = \mathbb{E} \left[\sup_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^N \Omega_i \ell(f(X_i), Y_i) \mid (X_1, Y_1), \dots, (X_n, Y_n) \right],$$

and the Rademacher complexity of the class \mathcal{H} , denoted $\text{Rad}(\mathcal{H})$, as:

$$\text{Rad}(\mathcal{H}) = \mathbb{E}[\widetilde{\text{Rad}}(\mathcal{H})] = \mathbb{E} \left[\sup_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^N \Omega_i \ell(f(X_i), Y_i) \right].$$

Remarks :

- We will use from now on the Rademacher complexity to establish as a measure of complexity of the hypothesis class \mathcal{H} .
- A fundamental result in this direction is the symmetrisation lemma.

Proposition : Symmetrization lemma

Under the Empirical Risk Minimization Framework, we have :

$$\mathbb{E} \left[\sup_{f \in \mathcal{H}} \{r(f) - R(f)\} \right] \leq 2 \text{Rad}(\mathcal{H}).$$

Proof : on the black board.

■ Applications in the regression setting

Proposition : Complexity for Lipschitz Losses

Under the Empirical Risk Minimization Framework, if $\hat{y} \mapsto \ell(\hat{y}, y)$ is G -lipschitz, then

$$\widetilde{\text{Rad}}(\mathcal{H}) \leq G \mathbb{E} \left[\sup_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \Omega_i f(X_i) \mid X_1, \dots, X_n \right].$$

Rademacher complexity of a hypothesis class

■ Formal definition & Symmetrization Lemma

Definition : Rademacher complexity of a class of functions

Let $(\Omega_1, \dots, \Omega_n)$ be n i.i.d. Rademacher variables.

We define the empirical Rademacher complexity of the hypothesis class \mathcal{H} , denoted $\widetilde{\text{Rad}}(\mathcal{H})$, as

$$\widetilde{\text{Rad}}(\mathcal{H}) = \mathbb{E} \left[\sup_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^N \Omega_i \ell(f(X_i), Y_i) \mid (X_1, Y_1), \dots, (X_n, Y_n) \right],$$

and the Rademacher complexity of the class \mathcal{H} , denoted $\text{Rad}(\mathcal{H})$, as:

$$\text{Rad}(\mathcal{H}) = \mathbb{E}[\widetilde{\text{Rad}}(\mathcal{H})] = \mathbb{E} \left[\sup_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^N \Omega_i \ell(f(X_i), Y_i) \right].$$

Remarks :

- We will use from now on the Rademacher complexity to establish as a measure of complexity of the hypothesis class \mathcal{H} .
- A fundamental result in this direction is the symmetrisation lemma.

Proposition : Symmetrization lemma

Under the Empirical Risk Minimization Framework, we have :

$$\mathbb{E} \left[\sup_{f \in \mathcal{H}} \{r(f) - R(f)\} \right] \leq 2 \text{Rad}(\mathcal{H}).$$

Proof : on the black board.

■ Applications in the regression setting

Proposition : Complexity for Lipschitz Losses

Under the Empirical Risk Minimization Framework, if $\hat{y} \mapsto \ell(\hat{y}, y)$ is G -lipschitz, then

$$\widetilde{\text{Rad}}(\mathcal{H}) \leq G \mathbb{E} \left[\sup_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \Omega_i f(X_i) \mid X_1, \dots, X_n \right].$$

Proof : on the black board.

Rademacher complexity of a hypothesis class

■ Formal definition & Symmetrization Lemma

Definition : Rademacher complexity of a class of functions

Let $(\Omega_1, \dots, \Omega_n)$ be n i.i.d. Rademacher variables.

We define the empirical Rademacher complexity of the hypothesis class \mathcal{H} , denoted $\widetilde{\text{Rad}}(\mathcal{H})$, as

$$\widetilde{\text{Rad}}(\mathcal{H}) = \mathbb{E} \left[\sup_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^N \Omega_i \ell(f(X_i), Y_i) \mid (X_1, Y_1), \dots, (X_n, Y_n) \right],$$

and the Rademacher complexity of the class \mathcal{H} , denoted $\text{Rad}(\mathcal{H})$, as:

$$\text{Rad}(\mathcal{H}) = \mathbb{E}[\widetilde{\text{Rad}}(\mathcal{H})] = \mathbb{E} \left[\sup_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^N \Omega_i \ell(f(X_i), Y_i) \right].$$

Remarks :

- We will use from now on the Rademacher complexity to establish as a measure of complexity of the hypothesis class \mathcal{H} .
- A fundamental result in this direction is the symmetrisation lemma.

Proposition : Symmetrization lemma

Under the Empirical Risk Minimization Framework, we have :

$$\mathbb{E} \left[\sup_{f \in \mathcal{H}} \{r(f) - R(f)\} \right] \leq 2 \text{Rad}(\mathcal{H}).$$

Proof : on the black board.

■ Applications in the regression setting

Proposition : Complexity for Lipschitz Losses

Under the Empirical Risk Minimization Framework, if $\hat{y} \mapsto \ell(\hat{y}, y)$ is G -lipschitz, then

$$\widetilde{\text{Rad}}(\mathcal{H}) \leq G \mathbb{E} \left[\sup_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \Omega_i f(X_i) \mid X_1, \dots, X_n \right].$$

Proof : on the black board.

Remarks :

- Hence, assuming a Lipschitz error function allows us to focus on the complexity spanned by the set of estimators rather than the composition

$$\ell(f(X), Y).$$

Linear Regression Examples

Proposition : Linear prediction with ℓ_2 -constraints, and ℓ_2 bounded data.

Let $\mathcal{H} \triangleq \{x \mapsto \omega^\top x, \|\omega\|_2 \leq c\}$, for $c > 0$ fixed.

Let $x = (x_1, \dots, x_n) \in \mathcal{X}^n$ be fixed and $T_x \triangleq \{(f(x_1), \dots, f(x_n)), f \in \mathcal{H}\}$

Then,

$$\text{Rad}(T_x) \leq c \cdot \frac{\max_{1 \leq i \leq n} \|x_i\|_2}{\sqrt{n}}.$$

Linear Regression Examples

Proposition : Linear prediction with ℓ_2 -constraints, and ℓ_2 bounded data.

Let $\mathcal{H} \triangleq \{x \mapsto \omega^\top x, \|\omega\|_2 \leq c\}$, for $c > 0$ fixed.

Let $x = (x_1, \dots, x_n) \in \mathcal{X}^n$ be fixed and $T_x \triangleq \{(f(x_1), \dots, f(x_n)), f \in \mathcal{H}\}$

Then,

$$\text{Rad}(T_x) \leq c \cdot \frac{\max_{1 \leq i \leq n} \|x_i\|_2}{\sqrt{n}}.$$

Proof : on the black board.

Linear Regression Examples

Proposition : Linear prediction with ℓ_2 -constraints, and ℓ_2 bounded data.

Let $\mathcal{H} \triangleq \{x \mapsto \omega^\top x, \|\omega\|_2 \leq c\}$, for $c > 0$ fixed.

Let $x = (x_1, \dots, x_n) \in \mathcal{X}^n$ be fixed and $T_x \triangleq \{(f(x_1), \dots, f(x_n)), f \in \mathcal{H}\}$

Then,

$$\text{Rad}(T_x) \leq c \cdot \frac{\max_{1 \leq i \leq n} \|x_i\|_2}{\sqrt{n}}.$$

Proof : on the black board.

Remark :

- We directly deduce

$$\text{Rad}(T_x) \leq c \cdot \max_{1 \leq i \leq n} \|x_i\|_\infty \sqrt{\frac{d}{n}}.$$

Linear Regression Examples

Proposition : Linear prediction with ℓ_2 -constraints, and ℓ_2 bounded data.

Let $\mathcal{H} \triangleq \{x \mapsto \omega^\top x, \|\omega\|_2 \leq c\}$, for $c > 0$ fixed.

Let $x = (x_1, \dots, x_n) \in \mathcal{X}^n$ be fixed and $T_x \triangleq \{(f(x_1), \dots, f(x_n)), f \in \mathcal{H}\}$

Then,

$$\text{Rad}(T_x) \leq c \cdot \frac{\max_{1 \leq i \leq n} \|x_i\|_2}{\sqrt{n}}.$$

Proof : on the black board.

Remark :

- We directly deduce

$$\text{Rad}(T_x) \leq c \cdot \max_{1 \leq i \leq n} \|x_i\|_\infty \sqrt{\frac{d}{n}}.$$

Proposition : Linear prediction with ℓ_1 -constraints, and ℓ_∞ bounded data.

Let $\mathcal{H} \triangleq \{x \mapsto \omega^\top x, \|\omega\|_1 \leq c\}$, for $c > 0$ fixed.

Let $x = (x_1, \dots, x_n) \in \mathcal{X}^n$ be fixed and $T_x \triangleq \{(f(x_1), \dots, f(x_n)), f \in \mathcal{H}\}$

Then,

$$\text{Rad}(T_x) \leq c \cdot \max_{1 \leq i \leq n} \|x_i\|_\infty \sqrt{\frac{2 \log(2d)}{n}}.$$

Linear Regression Examples

Proposition : Linear prediction with ℓ_2 -constraints, and ℓ_2 bounded data.

Let $\mathcal{H} \triangleq \{x \mapsto \omega^\top x, \|\omega\|_2 \leq c\}$, for $c > 0$ fixed.

Let $x = (x_1, \dots, x_n) \in \mathcal{X}^n$ be fixed and $T_x \triangleq \{(f(x_1), \dots, f(x_n)), f \in \mathcal{H}\}$

Then,

$$\text{Rad}(T_x) \leq c \cdot \frac{\max_{1 \leq i \leq n} \|x_i\|_2}{\sqrt{n}}.$$

Proof : on the black board.

Remark :

- We directly deduce

$$\text{Rad}(T_x) \leq c \cdot \max_{1 \leq i \leq n} \|x_i\|_\infty \sqrt{\frac{d}{n}}.$$

Proposition : Linear prediction with ℓ_1 -constraints, and ℓ_∞ bounded data.

Let $\mathcal{H} \triangleq \{x \mapsto \omega^\top x, \|\omega\|_1 \leq c\}$, for $c > 0$ fixed.

Let $x = (x_1, \dots, x_n) \in \mathcal{X}^n$ be fixed and $T_x \triangleq \{(f(x_1), \dots, f(x_n)), f \in \mathcal{H}\}$

Then,

$$\text{Rad}(T_x) \leq c \cdot \max_{1 \leq i \leq n} \|x_i\|_\infty \sqrt{\frac{2 \log(2d)}{n}}.$$

Remark :

- Note the logarithmic dependency in the dimension!

Linear Regression Examples

Proposition : Linear prediction with ℓ_2 -constraints, and ℓ_2 bounded data.

Let $\mathcal{H} \triangleq \{x \mapsto \omega^\top x, \|w\|_2 \leq c\}$, for $c > 0$ fixed.

Let $x = (x_1, \dots, x_n) \in \mathcal{X}^n$ be fixed and $T_x \triangleq \{(f(x_1), \dots, f(x_n)), f \in \mathcal{H}\}$

Then,

$$\text{Rad}(T_x) \leq c \cdot \frac{\max_{1 \leq i \leq n} \|x_i\|_2}{\sqrt{n}}.$$

Proof : on the black board.

Remark :

- We directly deduce

$$\text{Rad}(T_x) \leq c \cdot \max_{1 \leq i \leq n} \|x_i\|_\infty \sqrt{\frac{d}{n}}.$$

Proposition : Linear prediction with ℓ_1 -constraints, and ℓ_∞ bounded data.

Let $\mathcal{H} \triangleq \{x \mapsto \omega^\top x, \|w\|_1 \leq c\}$, for $c > 0$ fixed.

Let $x = (x_1, \dots, x_n) \in \mathcal{X}^n$ be fixed and $T_x \triangleq \{(f(x_1), \dots, f(x_n)), f \in \mathcal{H}\}$

Then,

$$\text{Rad}(T_x) \leq c \cdot \max_{1 \leq i \leq n} \|x_i\|_\infty \sqrt{\frac{2 \log(2d)}{n}}.$$

Remark :

- Note the logarithmic dependency in the dimension!

Proposition : Linear prediction with simplex-constraints, and ℓ_∞ bounded data.

Let $\mathcal{H} \triangleq \{x \mapsto \omega^\top x, w \in \Delta_d\}$, where $\Delta_d = \{w \in (\mathbb{R}_+)^d, \|w\|_1 = 1\}$.

Let $x = (x_1, \dots, x_n) \in \mathcal{X}^n$ be fixed and $T_x \triangleq \{(f(x_1), \dots, f(x_n)), f \in \mathcal{H}\}$

Then,

$$\text{Rad}(T_x) \leq \max_{1 \leq i \leq n} \|x_i\|_\infty \sqrt{\frac{2 \log(d)}{n}}.$$

Linear Regression Examples

Proposition : Linear prediction with ℓ_2 -constraints, and ℓ_2 bounded data.

Let $\mathcal{H} \triangleq \{x \mapsto \omega^\top x, \|w\|_2 \leq c\}$, for $c > 0$ fixed.

Let $x = (x_1, \dots, x_n) \in \mathcal{X}^n$ be fixed and $T_x \triangleq \{(f(x_1), \dots, f(x_n)), f \in \mathcal{H}\}$

Then,

$$\text{Rad}(T_x) \leq c \cdot \frac{\max_{1 \leq i \leq n} \|x_i\|_2}{\sqrt{n}}.$$

Proof : on the black board.

Remark :

- We directly deduce

$$\text{Rad}(T_x) \leq c \cdot \max_{1 \leq i \leq n} \|x_i\|_\infty \sqrt{\frac{d}{n}}.$$

Proposition : Linear prediction with ℓ_1 -constraints, and ℓ_∞ bounded data.

Let $\mathcal{H} \triangleq \{x \mapsto \omega^\top x, \|w\|_1 \leq c\}$, for $c > 0$ fixed.

Let $x = (x_1, \dots, x_n) \in \mathcal{X}^n$ be fixed and $T_x \triangleq \{(f(x_1), \dots, f(x_n)), f \in \mathcal{H}\}$

Then,

$$\text{Rad}(T_x) \leq c \cdot \max_{1 \leq i \leq n} \|x_i\|_\infty \sqrt{\frac{2 \log(2d)}{n}}.$$

Remark :

- Note the logarithmic dependency in the dimension!

Proposition : Linear prediction with simplex-constraints, and ℓ_∞ bounded data.

Let $\mathcal{H} \triangleq \{x \mapsto \omega^\top x, w \in \Delta_d\}$, where $\Delta_d = \{w \in (\mathbb{R}_+)^d, \|w\|_1 = 1\}$.

Let $x = (x_1, \dots, x_n) \in \mathcal{X}^n$ be fixed and $T_x \triangleq \{(f(x_1), \dots, f(x_n)), f \in \mathcal{H}\}$

Then,

$$\text{Rad}(T_x) \leq \max_{1 \leq i \leq n} \|x_i\|_\infty \sqrt{\frac{2 \log(d)}{n}}.$$

Proof : Almost identical to the previous one.

A Nonlinear Example - Feed-forward neural nets

- Definition

A Nonlinear Example - Feed-forward neural nets

■ Definition

Definition

Neural networks are non-linear estimators of the form

$$f(w, x) = \sigma(w_s^\top \sigma(\dots, \sigma(w_1^\top x)))$$

where for all $i \in \{1, \dots, s\}$, $w_i \in \mathbb{R}^{n_i \times n_{i-1}}$

$\sigma : u \mapsto (\sigma(u_1), \dots, \sigma(u_p))^\top$ applies the scalar function σ

to each coefficient of the input vector.

A Nonlinear Example - Feed-forward neural nets

■ Definition

Definition

Neural networks are non-linear estimators of the form

$$f(w, x) = \sigma(w_s^\top \sigma(\dots, \sigma(w_1^\top x)))$$

where for all $i \in \{1, \dots, s\}$, $w_i \in \mathbb{R}^{n_i \times n_{i-1}}$

$\sigma : u \mapsto (\sigma(u_1), \dots, \sigma(u_p))^\top$ applies the scalar function σ

to each coefficient of the input vector.

Remarks :

- n_0 must match the dimension d of the input x .

A Nonlinear Example - Feed-forward neural nets

■ Definition

Definition

Neural networks are non-linear estimators of the form

$$f(w, x) = \sigma(w_s^\top \sigma(\dots, \sigma(w_1^\top x)))$$

where for all $i \in \{1, \dots, s\}$, $w_i \in \mathbb{R}^{n_i \times n_{i-1}}$

$\sigma : u \mapsto (\sigma(u_1), \dots, \sigma(u_p))^\top$ applies the scalar function σ

to each coefficient of the input vector.

Remarks :

- n_0 must match the dimension d of the input x .
- σ is often called an **activation function**.

A Nonlinear Example - Feed-forward neural nets

■ Definition

Definition

Neural networks are non-linear estimators of the form

$$f(w, x) = \sigma(w_s^\top \sigma(\dots, \sigma(w_1^\top x)))$$

where for all $i \in \{1, \dots, s\}$, $w_i \in \mathbb{R}^{n_i \times n_{i-1}}$

$\sigma : u \mapsto (\sigma(u_1), \dots, \sigma(u_p))^\top$ applies the scalar function σ

to each coefficient of the input vector.

Remarks :

■ n_0 must match the dimension d of the input x .

■ σ is often called an **activation function**.

■ Examples of standard activation functions include

$$\text{ReLU} : t \mapsto \max(t, 0) \quad \text{Sigmoid} : t \mapsto \frac{1}{1 + \exp(-t)} \quad \text{Tanh} : t \mapsto \frac{e^t - e^{-t}}{e^t + e^{-t}}$$

A Nonlinear Example - Feed-forward neural nets

■ Definition

Definition

Neural networks are non-linear estimators of the form

$$f(w, x) = \sigma(w_s^\top \sigma(\dots, \sigma(w_1^\top x)))$$

where for all $i \in \{1, \dots, s\}$, $w_i \in \mathbb{R}^{n_i \times n_{i-1}}$

$\sigma : u \mapsto (\sigma(u_1), \dots, \sigma(u_p))^\top$ applies the scalar function σ

to each coefficient of the input vector.

Remarks :

■ n_0 must match the dimension d of the input x .

■ σ is often called an **activation function**.

■ Examples of standard activation functions include

$$\text{ReLU} : t \mapsto \max(t, 0) \quad \text{Sigmoid} : t \mapsto \frac{1}{1 + \exp(-t)} \quad \text{Tanh} : t \mapsto \frac{e^t - e^{-t}}{e^t + e^{-t}}$$

■ The integer s is often called the **depth** of the neural net. A neural network is called a deep neural net as soon as $s \geq 2$

A Nonlinear Example - Feed-forward neural nets

■ Definition

Definition

Neural networks are non-linear estimators of the form

$$f(w, x) = \sigma(w_s^\top \sigma(\dots, \sigma(w_1^\top x)))$$

where for all $i \in \{1, \dots, s\}$, $w_i \in \mathbb{R}^{n_i \times n_{i-1}}$

$\sigma : u \mapsto (\sigma(u_1), \dots, \sigma(u_p))^\top$ applies the scalar function σ to each coefficient of the input vector.

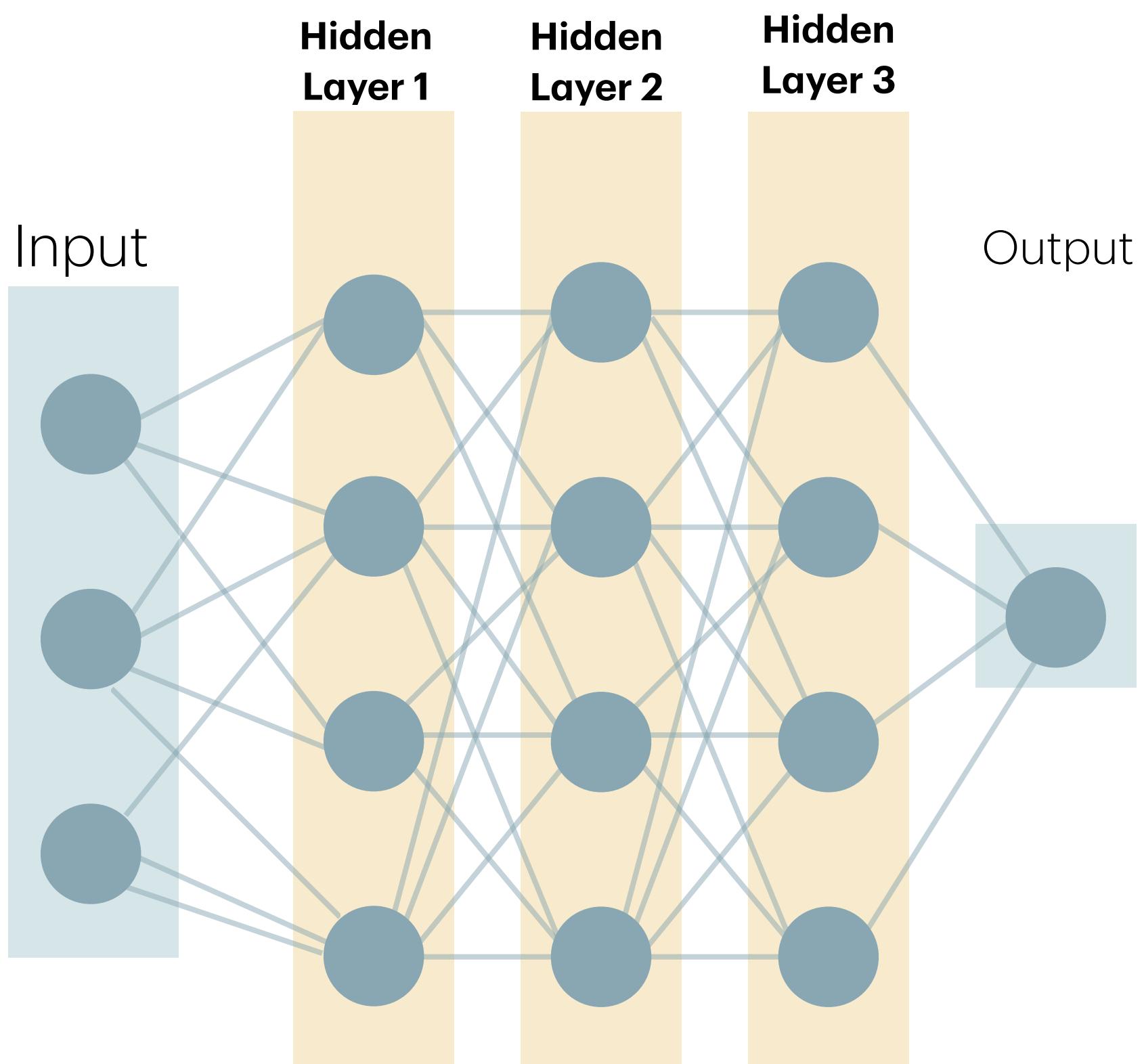
Remarks :

- n_0 must match the dimension d of the input x .
- σ is often called an **activation function**.
- Examples of standard activation functions include

$$\text{ReLU} : t \mapsto \max(t, 0) \quad \text{Sigmoid} : t \mapsto \frac{1}{1 + \exp(-t)} \quad \text{Tanh} : t \mapsto \frac{e^t - e^{-t}}{e^t + e^{-t}}$$

- The integer s is often called the **depth** of the neural net. A neural network is called a deep neural net as soon as $s \geq 2$

■ Visualisation



A complexity result for feed-forward neural nets

■ A recursive bound

Proposition : A recursive bound

Let \mathcal{L} be a class of functions from \mathbb{R}^d to \mathbb{R} that includes the zero function.

Let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ be γ -Lipschitz.

Define

$$\mathcal{L}' := \left\{ x \in \mathbb{R}^d \rightarrow \sigma \left(\sum_{j=1}^m w_j l_j(x) + b \right) \in \mathbb{R} : |b| \leq \beta, \|w\|_1 \leq \omega, l_1, \dots, l_m \in \mathcal{L} \right\}.$$

Then, for $x = (x_1, \dots, x_n) \in \mathcal{X}^n$,

$$\text{Rad}(\mathcal{L}'_x) \leq \gamma \left(\frac{\beta}{\sqrt{n}} + 2\omega \text{Rad}(\mathcal{L}_x) \right),$$

where

$$\mathcal{L} = \{\ell(x_1), \dots, \ell(x_n), \ell \in \mathcal{L}\},$$

$$\mathcal{L}'_x = \{\ell(x_1), \dots, \ell(x_n), \ell \in \mathcal{L}'\}.$$

A complexity result for feed-forward neural nets

■ A recursive bound

Proposition : A recursive bound

Let \mathcal{L} be a class of functions from \mathbb{R}^d to \mathbb{R} that includes the zero function.

Let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ be γ -Lipschitz.

Define

$$\mathcal{L}' := \left\{ x \in \mathbb{R}^d \rightarrow \sigma \left(\sum_{j=1}^m w_j l_j(x) + b \right) \in \mathbb{R} : |b| \leq \beta, \|w\|_1 \leq \omega, l_1, \dots, l_m \in \mathcal{L} \right\}.$$

Then, for $x = (x_1, \dots, x_n) \in \mathcal{X}^n$,

$$\text{Rad}(\mathcal{L}'_x) \leq \gamma \left(\frac{\beta}{\sqrt{n}} + 2\omega \text{Rad}(\mathcal{L}_x) \right),$$

where

$$\begin{aligned} \mathcal{L} &= \{\ell(x_1), \dots, \ell(x_n), \ell \in \mathcal{L}\}, \\ \mathcal{L}'_x &= \{\ell(x_1), \dots, \ell(x_n), \ell \in \mathcal{L}'\}. \end{aligned}$$

Proof : on the black board.

A complexity result for feed-forward neural nets

■ A recursive bound

Proposition : A recursive bound

Let \mathcal{L} be a class of functions from \mathbb{R}^d to \mathbb{R} that includes the zero function.

Let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ be γ -Lipschitz.

Define

$$\mathcal{L}' := \left\{ x \in \mathbb{R}^d \rightarrow \sigma \left(\sum_{j=1}^m w_j l_j(x) + b \right) \in \mathbb{R} : |b| \leq \beta, \|w\|_1 \leq \omega, l_1, \dots, l_m \in \mathcal{L} \right\}.$$

Then, for $x = (x_1, \dots, x_n) \in \mathcal{X}^n$,

$$\text{Rad}(\mathcal{L}'_x) \leq \gamma \left(\frac{\beta}{\sqrt{n}} + 2\omega \text{Rad}(\mathcal{L}_x) \right),$$

where

$$\mathcal{L} = \{ \ell(x_1), \dots, \ell(x_n), \ell \in \mathcal{L} \},$$

$$\mathcal{L}'_x = \{ \ell(x_1), \dots, \ell(x_n), \ell \in \mathcal{L}' \}.$$

Proof : on the black board.

■ Complexity result

Proposition : A recursive bound

Let σ be γ -Lipschitz.

Consider the class \mathcal{L} of neural networks losses with n layers

$$\begin{aligned} \ell^{(k)}(x) &:= \begin{cases} \sigma(\mathbf{w}^{(k)}x + b^{(k)}) & \text{if } k \leq n-1 \\ \mathbf{w}^{(n)}x + b^{(n)} & \text{if } k = n \end{cases} \text{ such that} \\ \|w^{(k)}\|_\infty &\leq \omega, b^{(k)} \leq \beta, \forall k \in \{1, \dots, n\}. \end{aligned}$$

Then, for any $x = (x_1, \dots, x_n) \in \mathbb{R}^{d \times n}$, the set \mathcal{L}_x satisfies

$$\text{Rad}(\mathcal{L}_x) \leq \frac{1}{\sqrt{n}} \left(\beta + 2\omega\beta\lambda \sum_{k=0}^{n-3} (2\omega\lambda)^k + 2\omega(2\omega\lambda)^{n-2} \max_i \|x_i\|_\infty \sqrt{2 \log(2d)} \right).$$

Proof : exercise!