

# A Bit of Sequence Prediction: Lec 2

# LINEAR BETTING GAME

For  $t = 1$  to  $n$

- Learner has to place a bet of amount  $|\hat{y}_t|$  on either team A or team B (sign of  $\hat{y}_t$  tells us which team we bet on)
- Outcome of the round is revealed as  $y_t \in \{\pm 1\}$
- Learner loses money  $\ell_t = -y_t \cdot \hat{y}_t$

End For

Goal: given  $\phi : \{\pm 1\}^n \mapsto \mathbb{R}$  can we guarantee:

$$\sum_{t=1}^n (-y_t \cdot \hat{y}_t) \leq \phi(y_1, \dots, y_n)$$

# LINEAR BETTING RESULT

## Lemma

For any  $\phi : \{\pm 1\}^n \mapsto \mathbb{R}$ , there exists a strategy with guarantee that

$$\sum_{t=1}^n -(\hat{y}_t \cdot y_t) \leq \phi(y_1, \dots, y_n)$$

If and only if

$$\mathbb{E}[\phi(\epsilon_1, \dots, \epsilon_n)] \geq 0$$

and the strategy achieving this is  $\hat{y}_t = \frac{1}{2} \mathbb{E}[\phi(y_1, \dots, y_{t-1}, -1, \epsilon_{t+1}, \dots, \epsilon_n) - \phi(y_1, \dots, y_{t-1}, +1, \epsilon_{t+1}, \dots, \epsilon_n)]$ .

*Proof:*

Again the reverse direction is easy just plug in random outcomes. For the other direction is similar to Cover's result proof, only this time we don't need restriction on  $\phi$  since we can place any magnitude bet.

# BETTING EXAMPLE

- We are given that  $m$  teams are playing in pairs  $n$  matches against each other.
- Assume that the pairs that are going to play for the  $n$  rounds are announced in advance as  $(i_1, j_1), \dots, (i_n, j_n)$ .
- Benchmark we want to consider is one where, after the fact, we assign scores  $w[1], \dots, w[m]$  to each of the  $m$  teams and when teams  $i, j$  play, the benchmark makes a bet of  $w[i] - w[j]$
- Maximum bet value allowed by the benchmark is some value  $B$

Goal: Strategy to minimize

$$\text{Reg}_n := \sum_{t=1}^n -(\hat{y}_t \cdot y_t) - \min_{w \in \mathbb{R}^m: \max_{i,j} w[i] - w[j] \leq B} \sum_{t=1}^n (-y_t \cdot (w[i_t] - w[j_t]))$$

# BETTING EXAMPLE

- Define

$$\phi(y_1, \dots, y_n) = \min_{w \in \mathbb{R}^m: \max_{i,j} w[i] - w[j] \leq B} \sum_{t=1}^n (-y_t \cdot (w[i_t] - w[j_t])) + C_n$$

- Using the Lemma lets write out  $C_n$

$$\begin{aligned} C_n &= -\mathbb{E}_{\epsilon_1, \dots, \epsilon_n} \left[ \min_{w \in \mathbb{R}^m: \max_{i,j} w[i] - w[j] \leq B} \sum_{t=1}^n (-\epsilon_t \cdot (w[i_t] - w[j_t])) \right] \\ &= \mathbb{E}_{\epsilon_1, \dots, \epsilon_n} \left[ \max_{w \in \mathbb{R}^m: \max_{i,j} w[i] - w[j] \leq B} \sum_{t=1}^n (\epsilon_t \cdot (w[i_t] - w[j_t])) \right] \\ &= \mathbb{E}_{\epsilon_1, \dots, \epsilon_n} \left[ \max_{w \in [0, B]^m} \sum_{t=1}^n (\epsilon_t \cdot (w[i_t] - w[j_t])) \right] \\ &= \mathbb{E}_{\epsilon_1, \dots, \epsilon_n} \left[ \max_{w \in [0, B]^m} \sum_{i=1}^m \sum_{t=1}^n \epsilon_t \cdot w[i] (\mathbf{1}\{i_t = i\} - \mathbf{1}\{j_t = i\}) \right] \\ &= \mathbb{E}_{\epsilon_1, \dots, \epsilon_n} \left[ \sum_{i=1}^m \max_{w[i] \in [0, B]} \sum_{t=1}^n \epsilon_t \cdot w[i] (\mathbf{1}\{i_t = i\} - \mathbf{1}\{j_t = i\}) \right] \\ &= B \mathbb{E}_{\epsilon_1, \dots, \epsilon_n} \left[ \sum_{i=1}^m \max \left\{ \sum_{t=1}^{n_i} \epsilon_t, 0 \right\} \right] \leq \frac{B}{2} \sum_{i=1}^m \sqrt{n_i} \leq \frac{B}{2} \sqrt{mn} \end{aligned}$$

What if we didn't know the which pairs play in advance?

# BETTING GAME WITH ARBITRARY CONTEXTS

For  $t = 1$  to  $n$

- Context  $x_t \in \mathcal{X}$  is provided.
- Learner has to place a bet of amount  $|\hat{y}_t|$  on either team A or team B (sign of  $\hat{y}_t$  tells us which team we bet on)
- Outcome of the round is revealed as  $y_t \in \{\pm 1\}$
- Learner loses money  $\ell_t = -y_t \cdot \hat{y}_t$

End For

Goal: given  $\phi : \mathcal{X}^n \times \{\pm 1\}^n \mapsto \mathbb{R}$  can we guarantee:

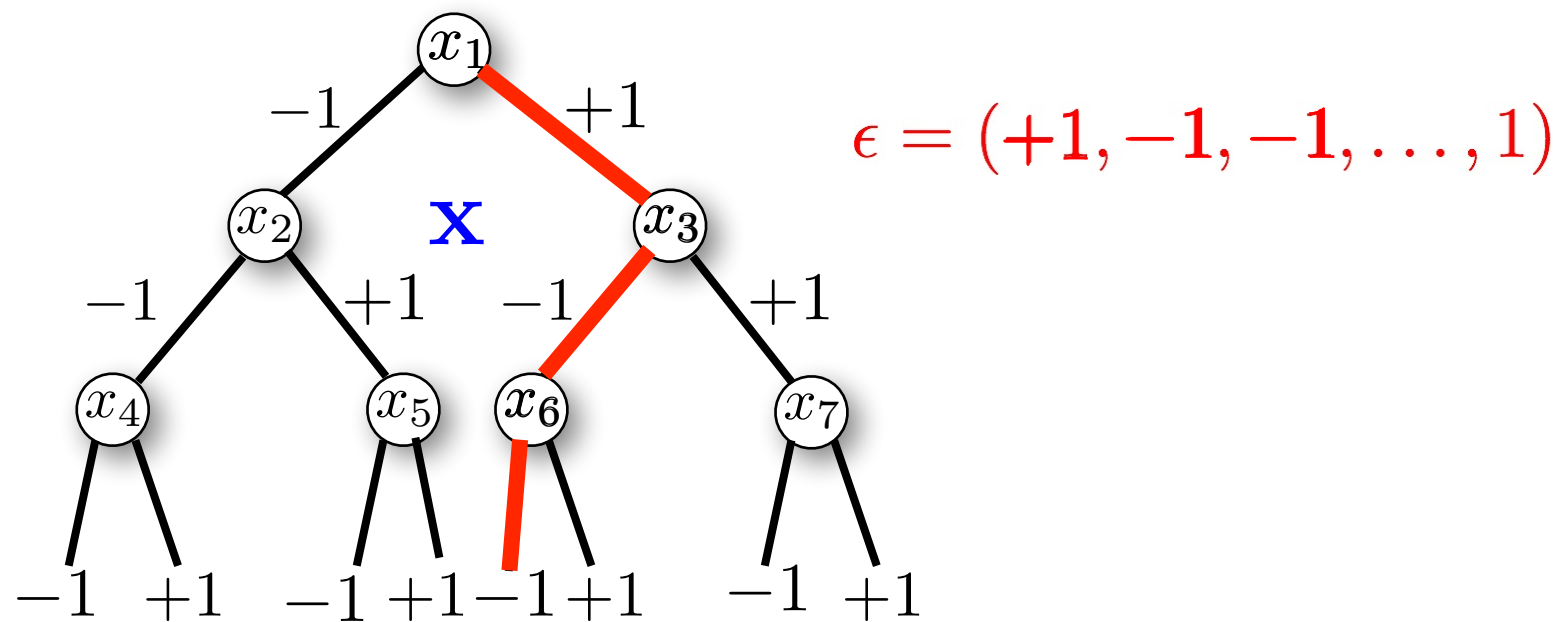
$$\sum_{t=1}^n (-y_t \cdot \hat{y}_t) \leq \phi(x_1, \dots, x_n, y_1, \dots, y_n)$$

Eg.  $x_t = (i_t, j_t)$  teams playing can be chosen arbitrarily as we go.

# MEET THE TREES

## Definition

An  $\mathcal{X}$  valued binary tree is a sequence of mapping  $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$  where  $\mathbf{x}_t : \{\pm 1\}^{t-1} \mapsto \mathcal{X}$ . Here  $\mathbf{x}_1 \in \mathcal{X}$  is a constant



$$\mathbf{x}_1 = x_1$$

$$\mathbf{x}_2(+1) = x_3$$

$$\mathbf{x}_3(+1, -1) = x_6$$



# ARBITRARY COVARIATES RESULT

## Lemma

For any  $\phi : \mathcal{X}^n \times \{\pm 1\}^n \mapsto \mathbb{R}$ , there exists a strategy with can guarantee that

$$\sum_{t=1}^n -(\hat{y}_t \cdot y_t) \leq \phi(x_1, y_1, \dots, x_n, y_n)$$

If and only if

$$\inf_{\mathbf{x}} \mathbb{E} [\phi(\mathbf{x}_1, \mathbf{x}_2(\epsilon_1), \dots, \mathbf{x}_n(\epsilon_1, \dots, \epsilon_{n-1}), \epsilon_1, \dots, \epsilon_n)] \geq 0$$

# ARBITRARY COVARIATES RESULT

- If we consider the example

$$\phi(x_1, \dots, x_n, y_1, \dots, y_n) = \min_{f \in \mathcal{F}} \left\{ \sum_{t=1}^n -y_t \cdot f(x_t) \right\} + C_n(\mathcal{F})$$

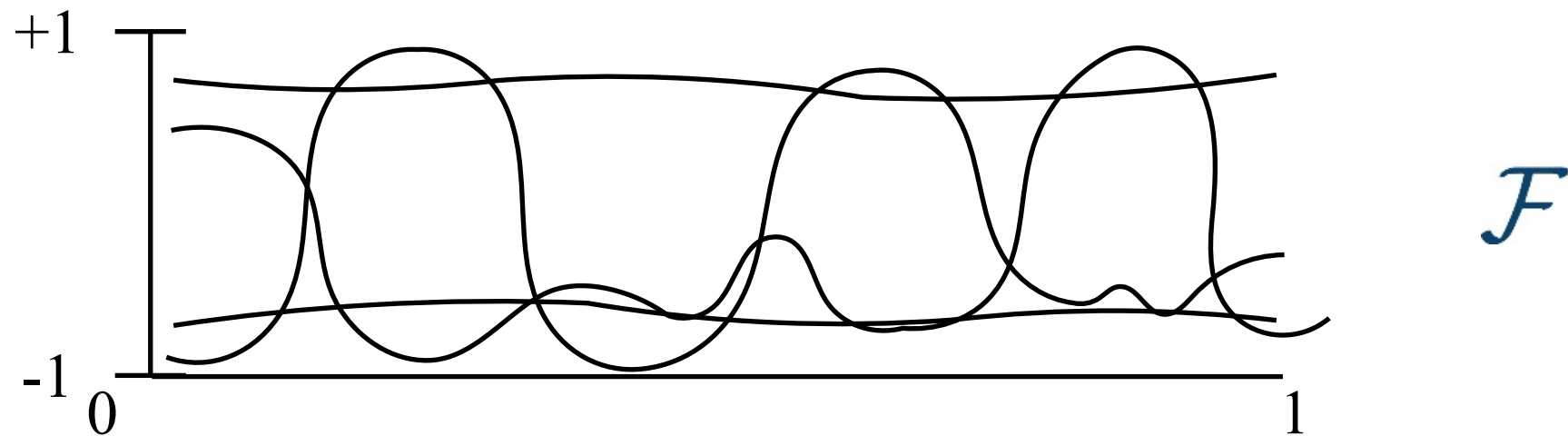
- In this case, using the lemma we get:

$$\begin{aligned} C_n(\mathcal{F}) &= -\mathbb{E}_{\epsilon} \left[ \min_{f \in \mathcal{F}} \left\{ \sum_{t=1}^n -\epsilon_t \cdot f(\mathbf{x}_t(\epsilon_1, \dots, \epsilon_{t-1})) \right\} \right] \\ &= \mathbb{E}_{\epsilon} \left[ \max_{f \in \mathcal{F}} \left\{ \sum_{t=1}^n \epsilon_t \cdot f(\mathbf{x}_t(\epsilon_1, \dots, \epsilon_{t-1})) \right\} \right] := \text{Rad}_n(\mathcal{F}) \end{aligned}$$

# RADEMACHER COMPLEXITY

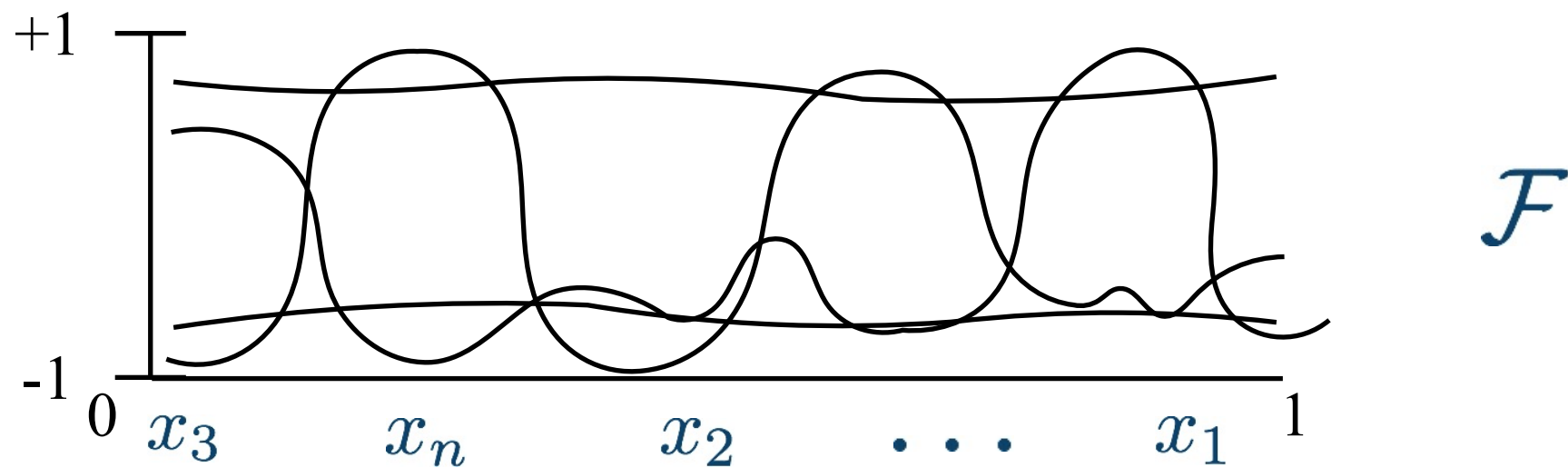
# RADEMACHER COMPLEXITY

Example :  $\mathcal{X} = [0, 1]$ ,  $\mathcal{Y} = [-1, 1]$



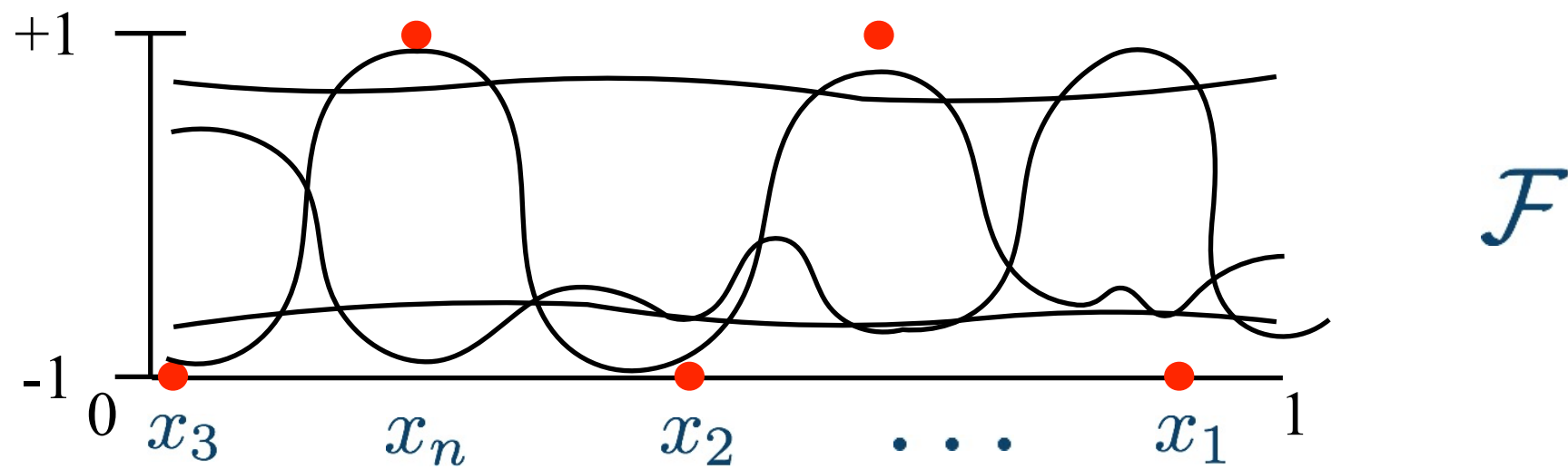
# RADEMACHER COMPLEXITY

Example :  $\mathcal{X} = [0, 1]$ ,  $\mathcal{Y} = [-1, 1]$



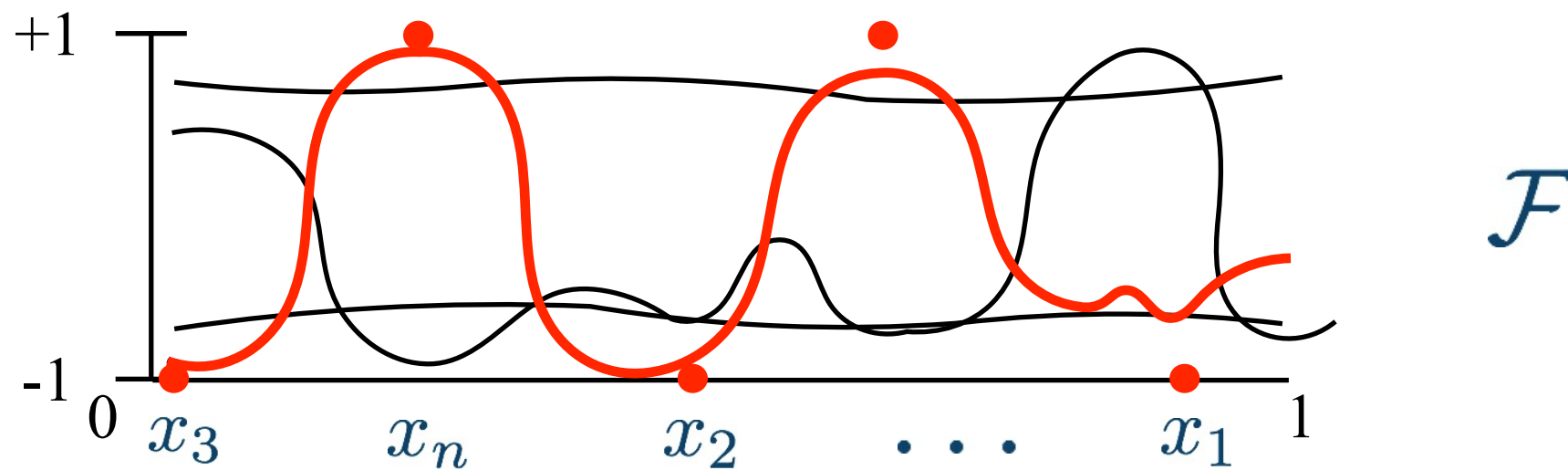
# RADEMACHER COMPLEXITY

Example :  $\mathcal{X} = [0, 1]$ ,  $\mathcal{Y} = [-1, 1]$



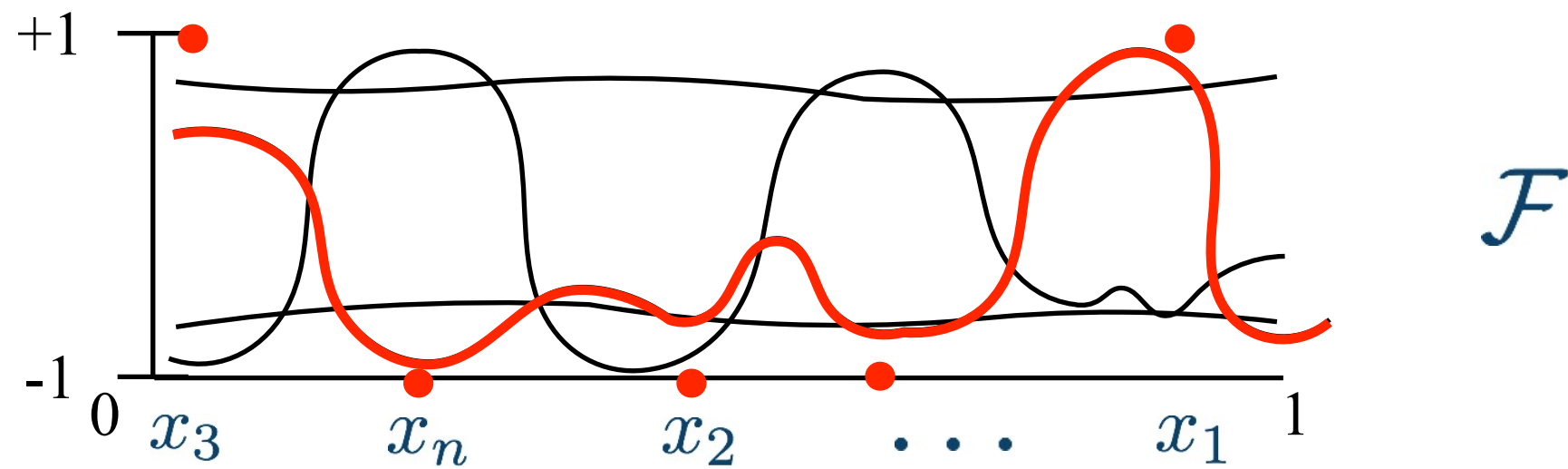
# RADEMACHER COMPLEXITY

Example :  $\mathcal{X} = [0, 1]$ ,  $\mathcal{Y} = [-1, 1]$



# RADEMACHER COMPLEXITY

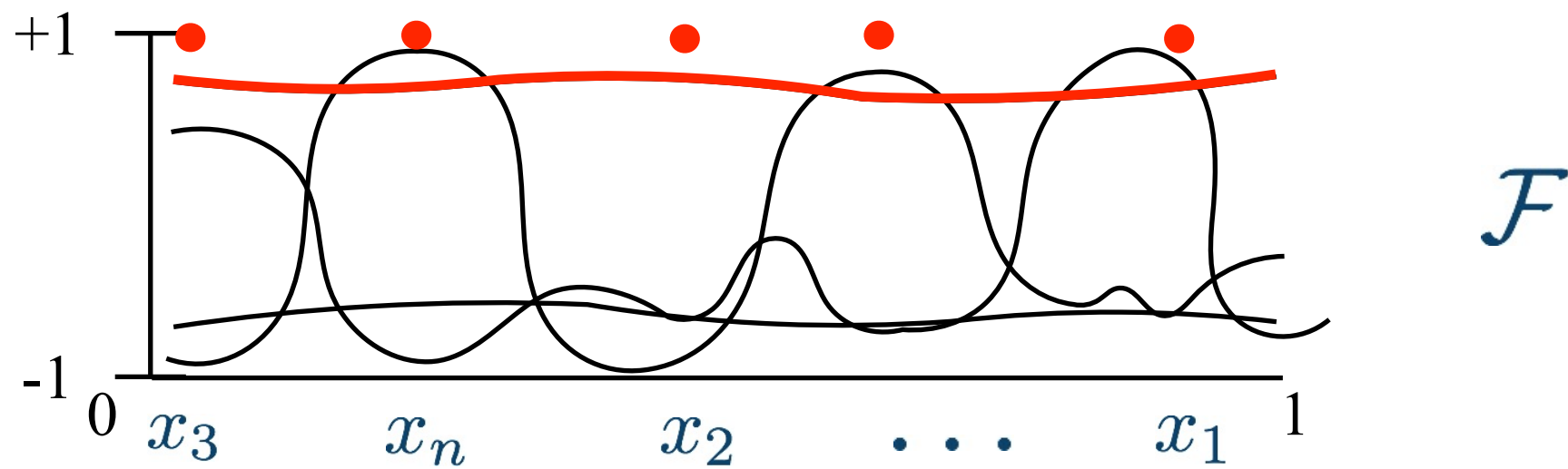
Example :  $\mathcal{X} = [0, 1]$ ,  $\mathcal{Y} = [-1, 1]$





# RADEMACHER COMPLEXITY

Example :  $\mathcal{X} = [0, 1]$ ,  $\mathcal{Y} = [-1, 1]$



# RADEMACHER COMPLEXITY

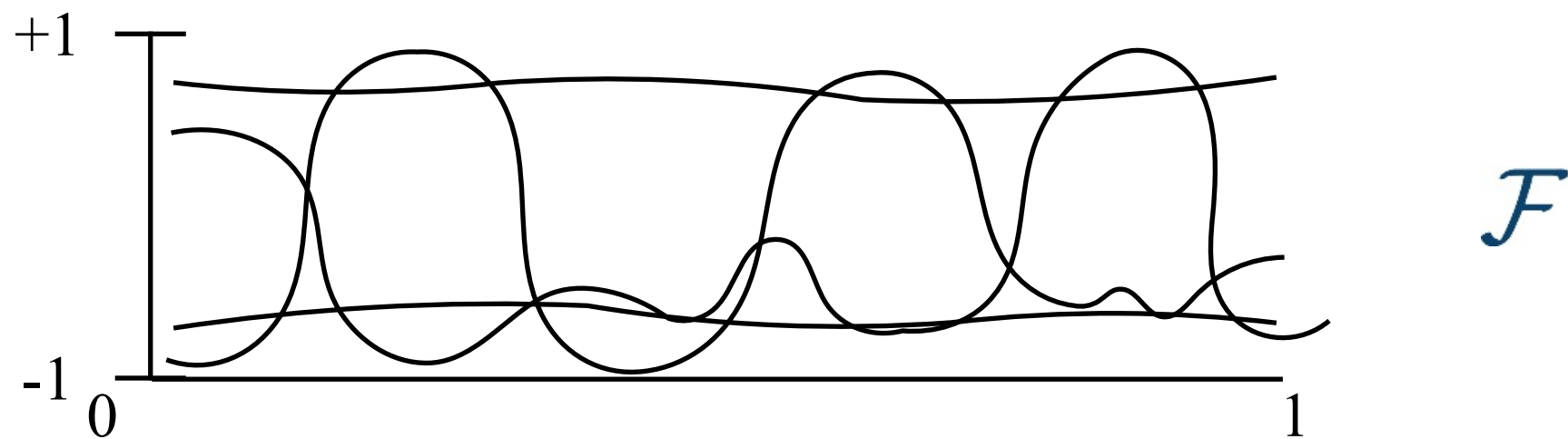
$$\mathcal{R}_n(\mathcal{F}) := \sup_{x_1, \dots, x_n} \mathbb{E}_{\epsilon} \left[ \frac{2}{n} \sup_{f \in \mathcal{F}} \left| \sum_{t=1}^n \epsilon_t f(x_t) \right| \right]$$

# RADEMACHER COMPLEXITY

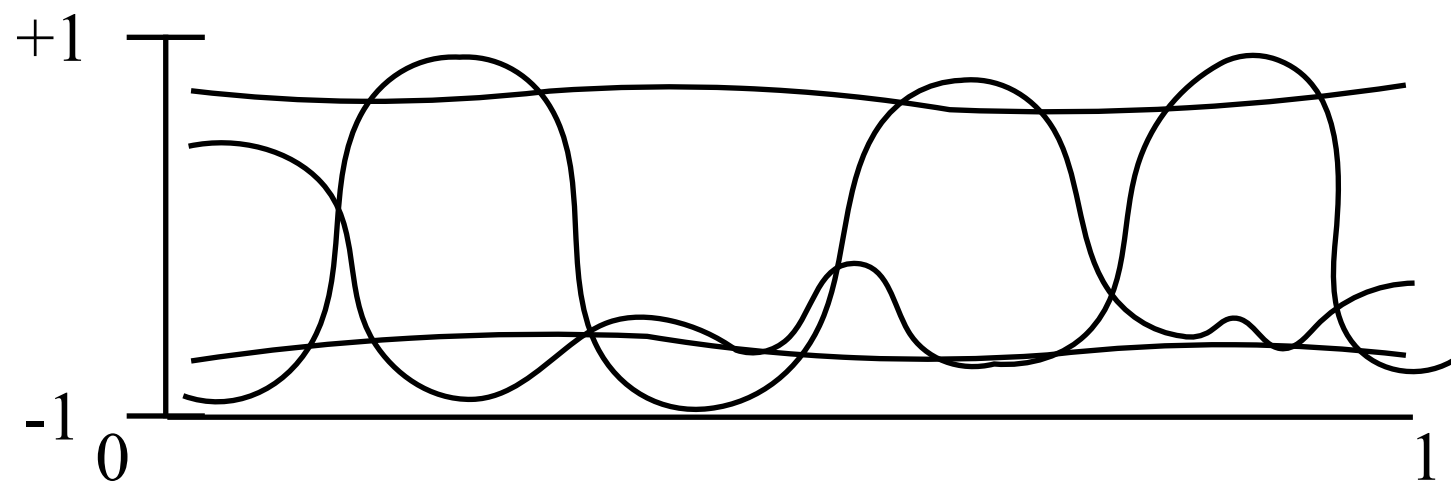
$$\mathcal{R}_n(\mathcal{F}) := \sup_{x_1, \dots, x_n} \mathbb{E}_{\epsilon} \left[ \frac{2}{n} \sup_{f \in \mathcal{F}} \left| \sum_{t=1}^n \epsilon_t f(x_t) \right| \right]$$

sequence random signs max correlation

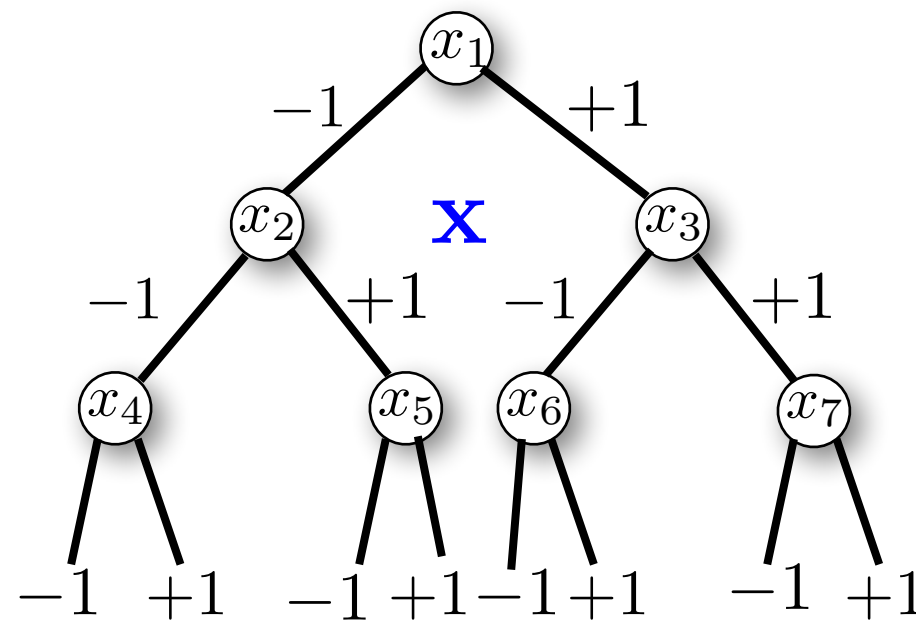
# SEQUENTIAL RADEMACHER COMPLEXITY



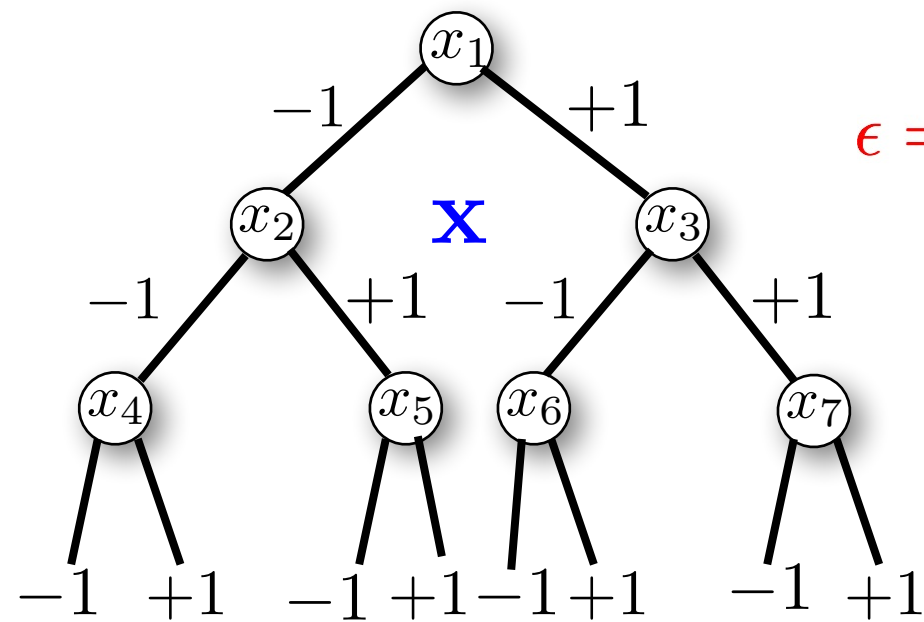
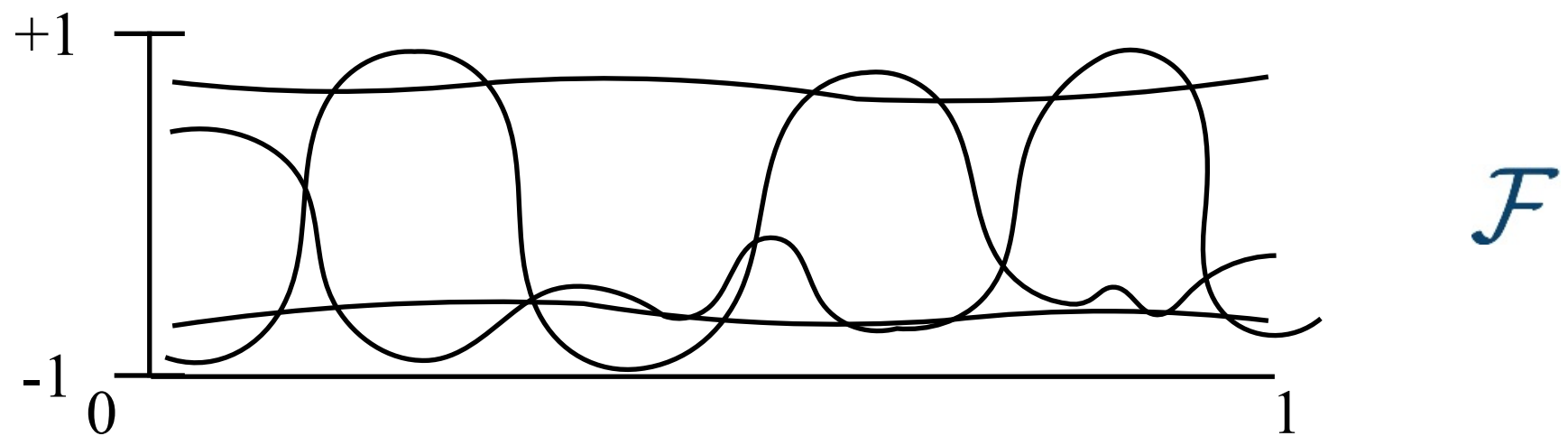
# SEQUENTIAL RADEMACHER COMPLEXITY



$\mathcal{F}$

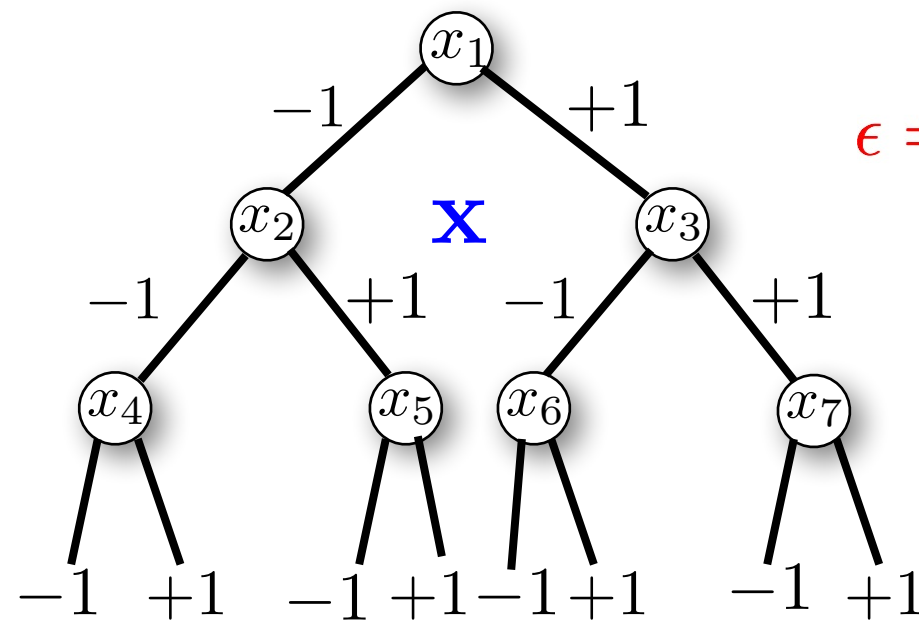
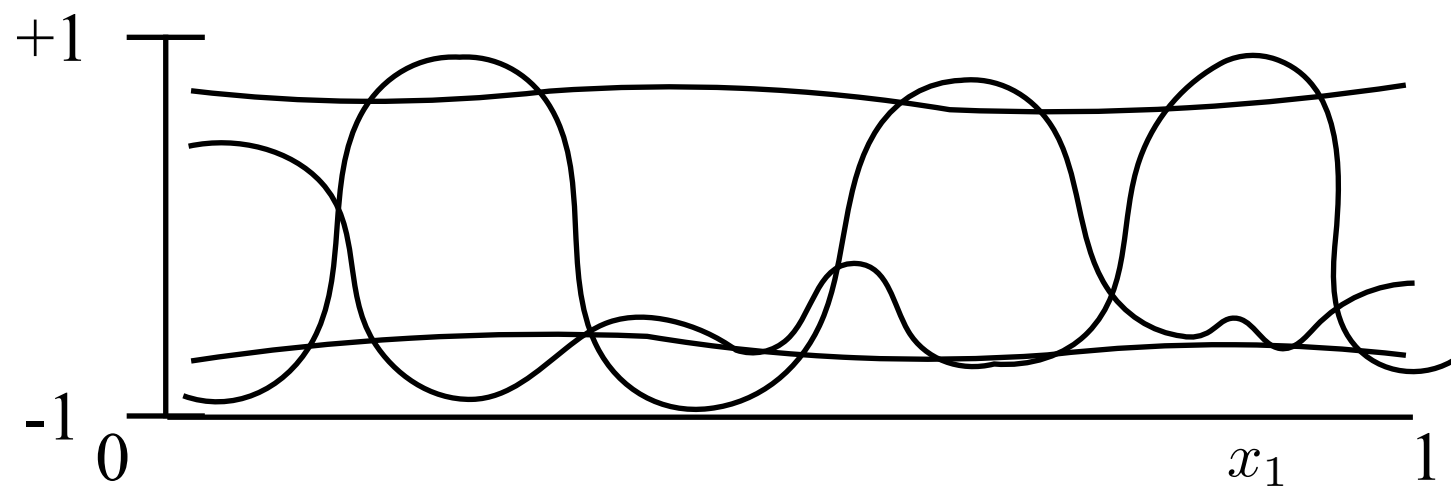


# SEQUENTIAL RADEMACHER COMPLEXITY



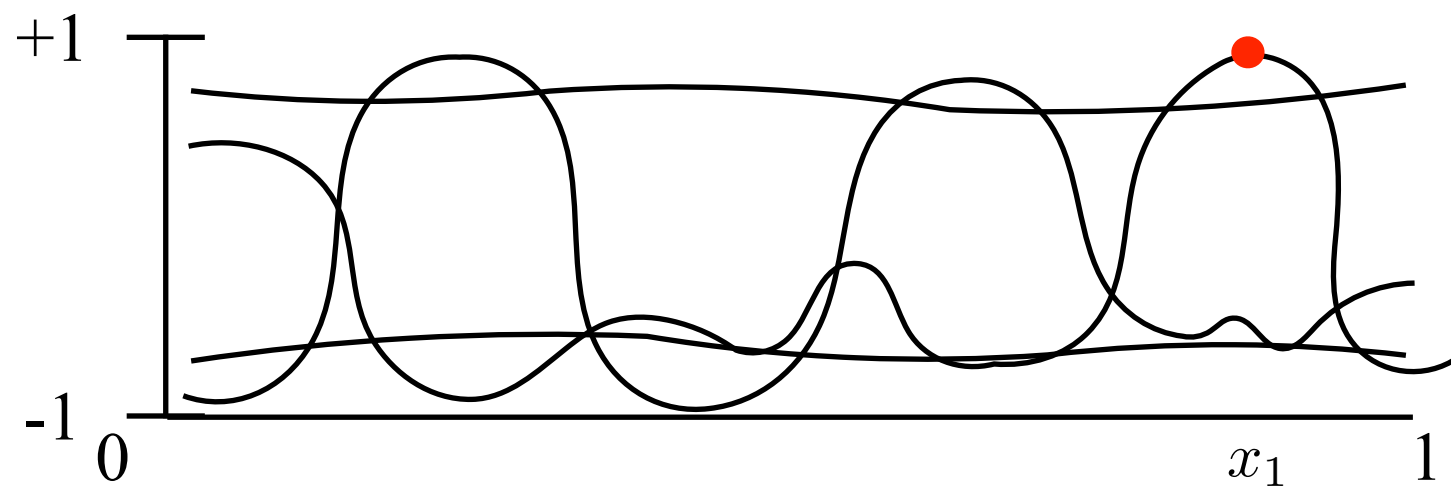
$$\epsilon = (+1, -1, -1, \dots, 1)$$

# SEQUENTIAL RADEMACHER COMPLEXITY

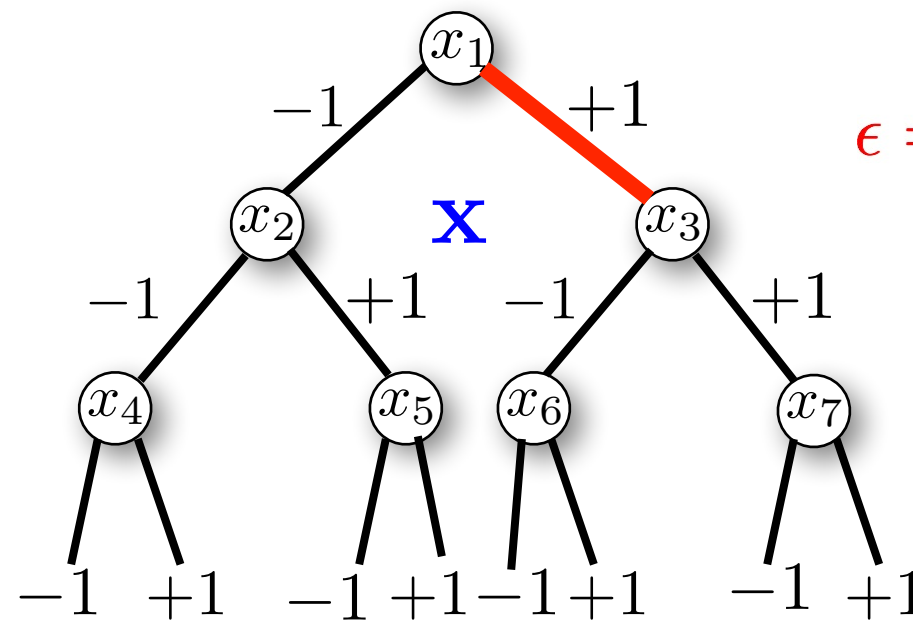


$$\epsilon = (+1, -1, -1, \dots, 1)$$

# SEQUENTIAL RADEMACHER COMPLEXITY



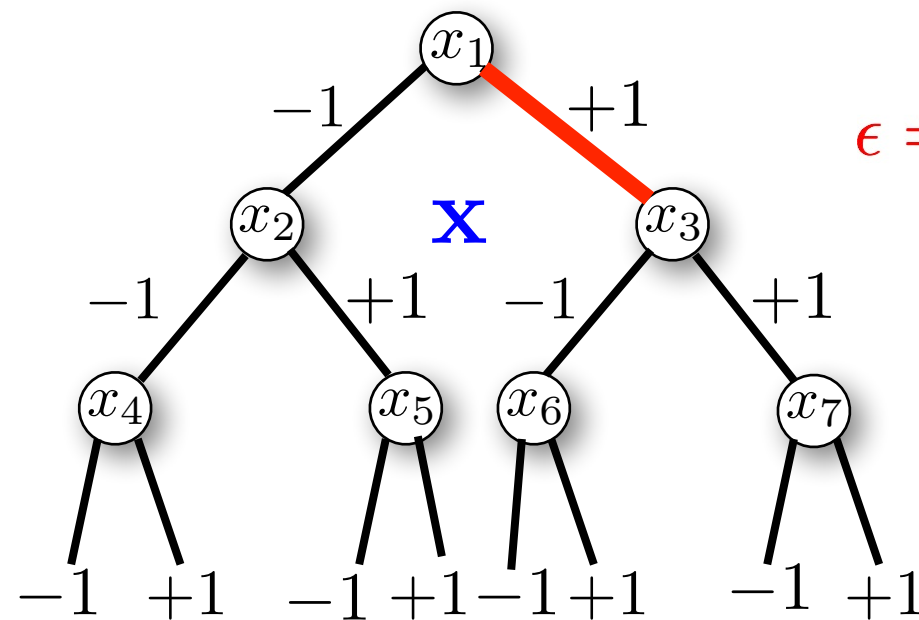
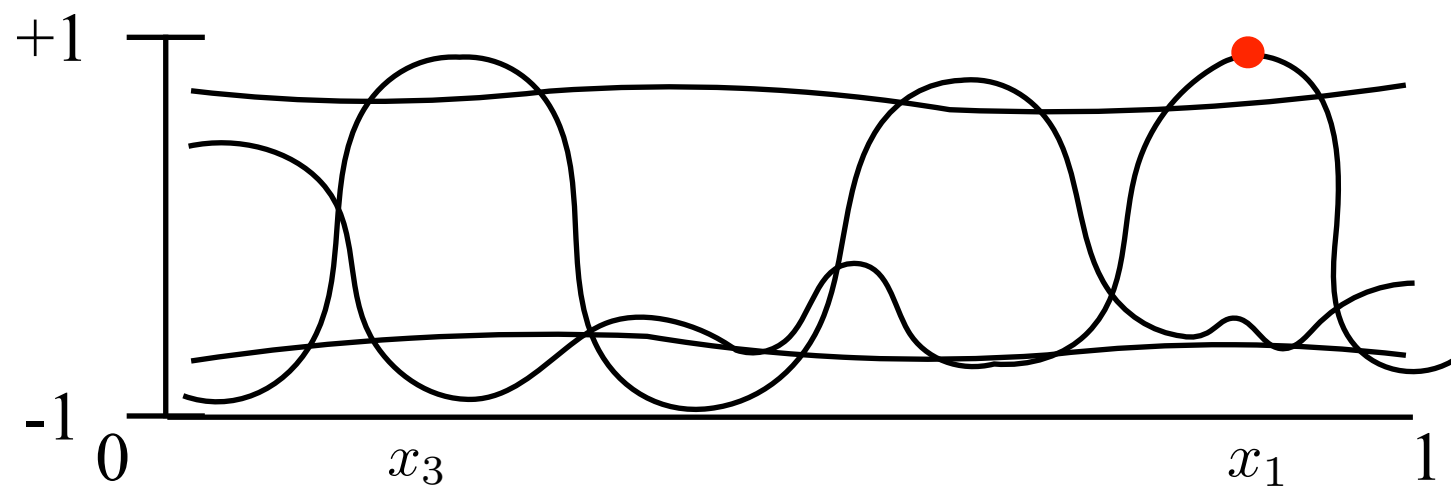
$\mathcal{F}$



$$\epsilon = (+1, -1, -1, \dots, 1)$$

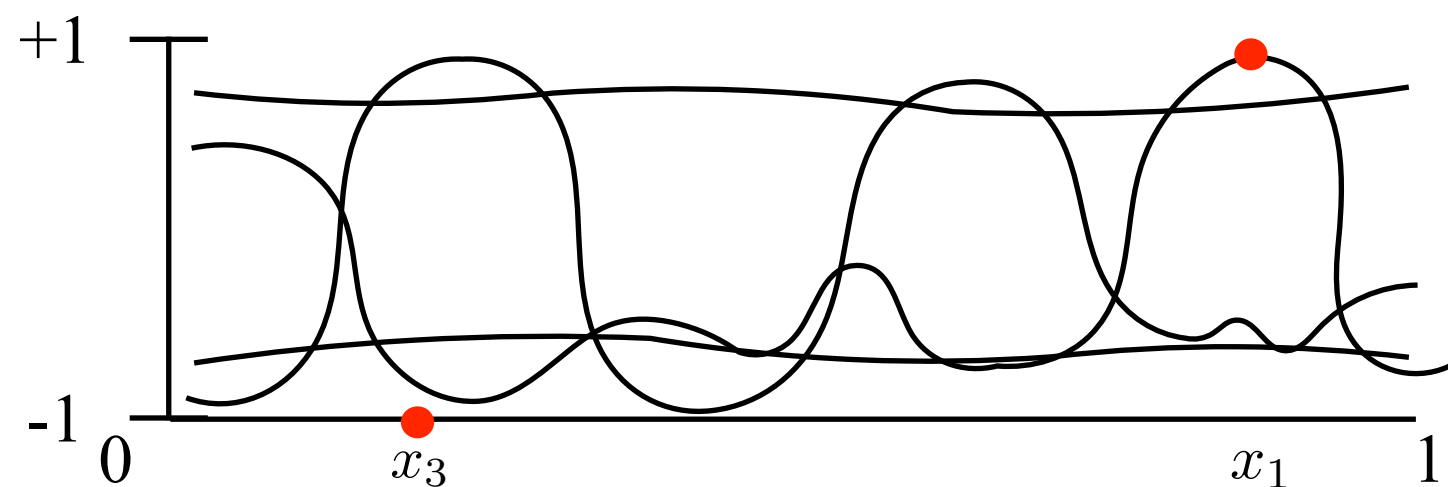


# SEQUENTIAL RADEMACHER COMPLEXITY

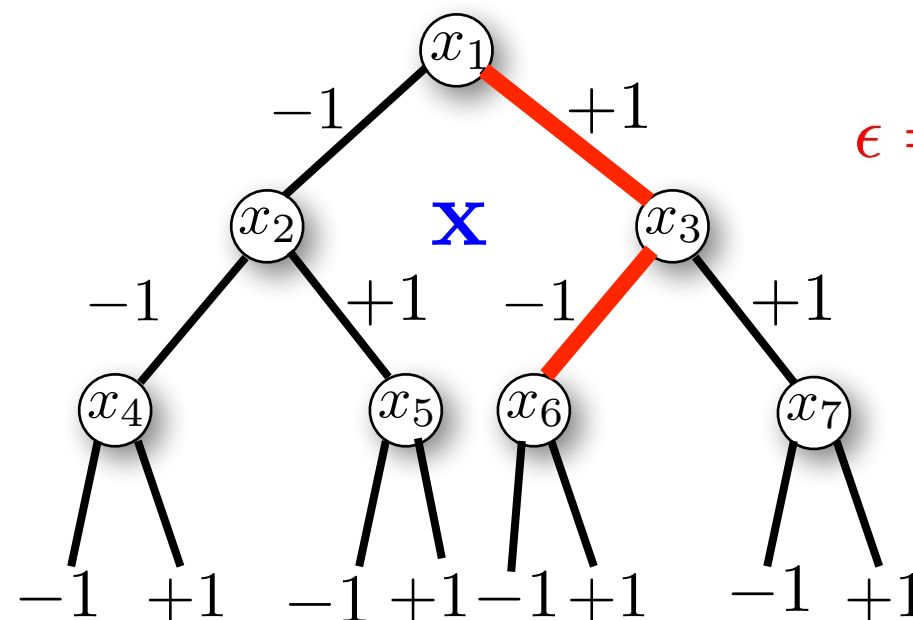


$$\epsilon = (+1, -1, -1, \dots, 1)$$

# SEQUENTIAL RADEMACHER COMPLEXITY

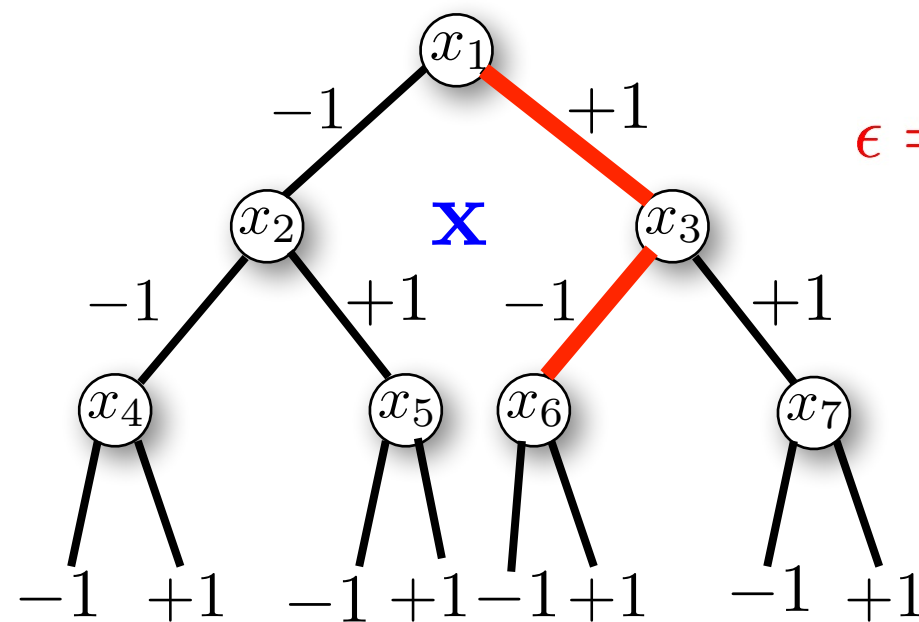
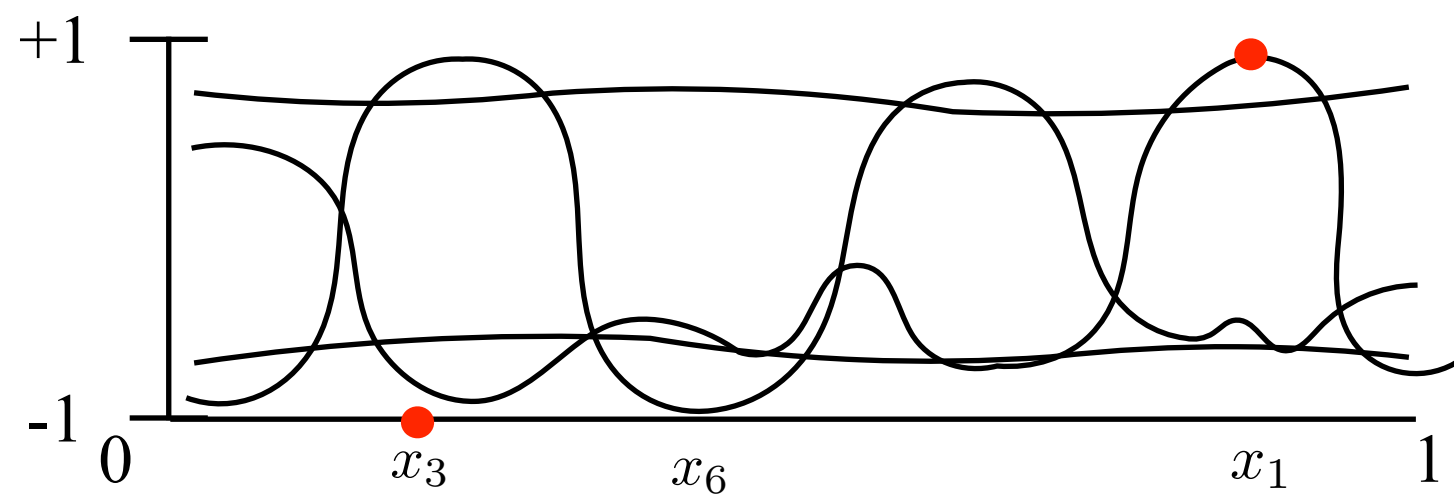


$\mathcal{F}$



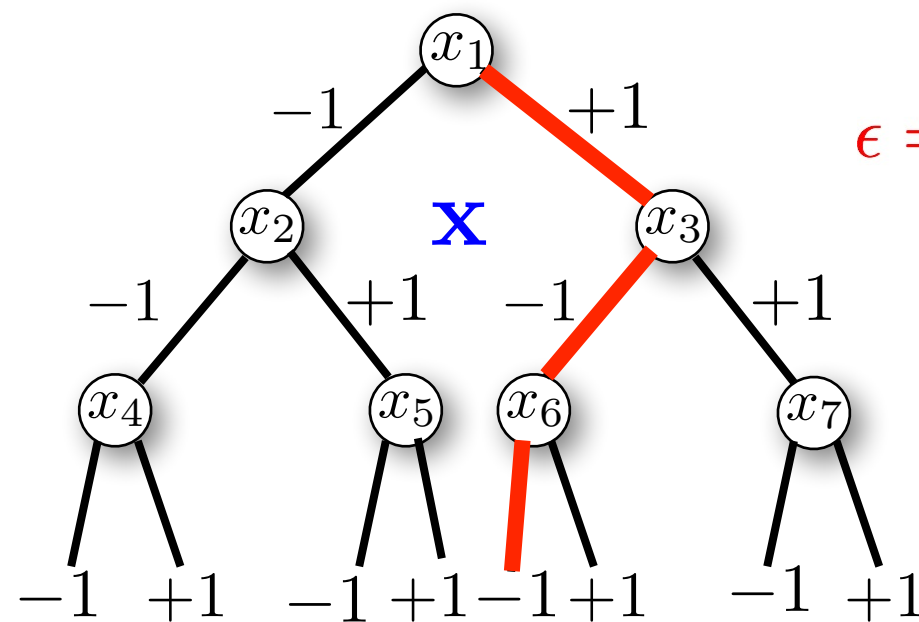
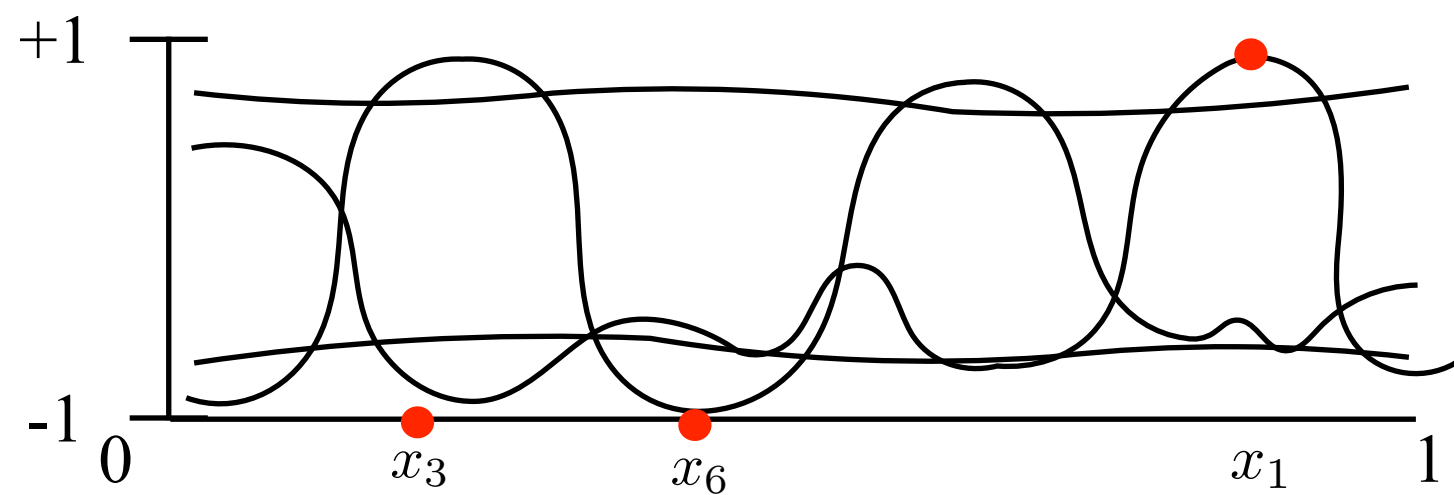
$\epsilon = (+1, -1, -1, \dots, 1)$

# SEQUENTIAL RADEMACHER COMPLEXITY



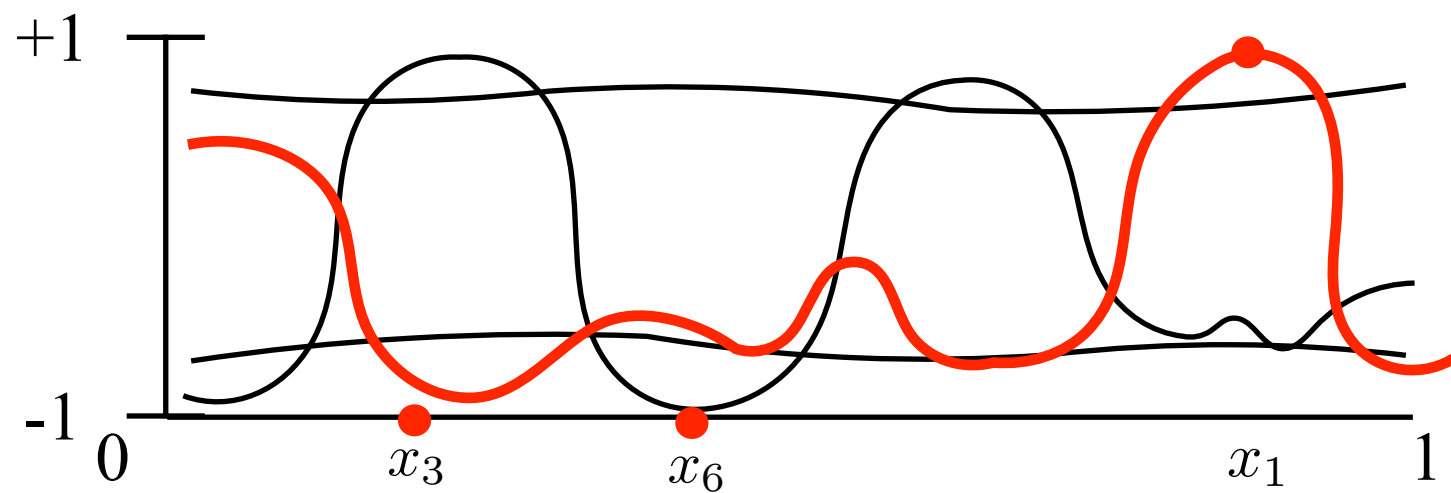
$$\epsilon = (+1, -1, -1, \dots, 1)$$

# SEQUENTIAL RADEMACHER COMPLEXITY

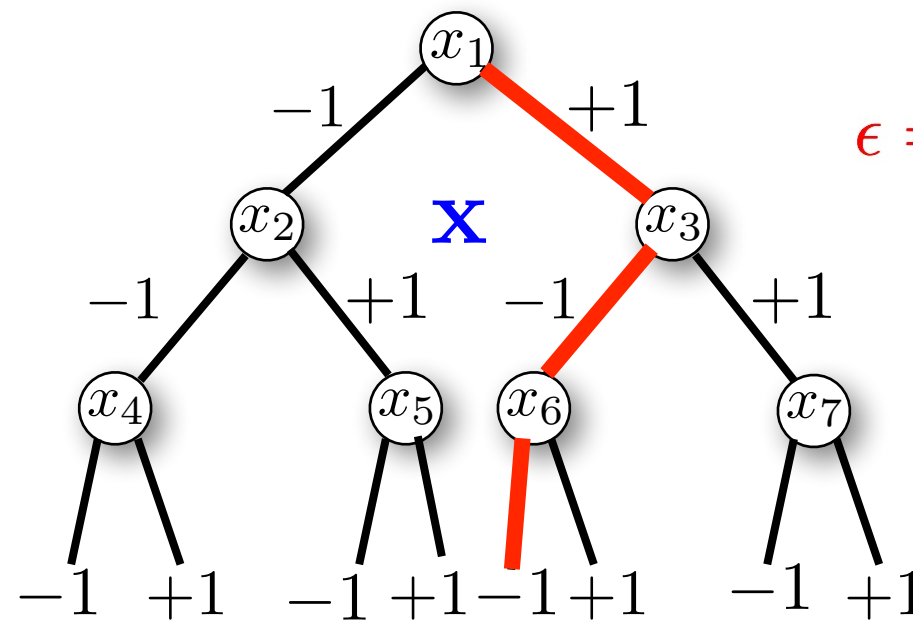


$$\epsilon = (+1, -1, -1, \dots, 1)$$

# SEQUENTIAL RADEMACHER COMPLEXITY



$\mathcal{F}$



$$\epsilon = (+1, -1, -1, \dots, 1)$$

# ONLINE SUPERVISED LEARNING

Given: model class  $\mathcal{F} \subset \mathbb{R}^{\mathcal{X}}$ , convex  $L$  Lipschitz loss  $\ell : \mathbb{R} \times \mathcal{Z} \mapsto \mathbb{R}$

For  $t = 1$  to  $n$

- Context  $x_t \in \mathcal{X}$  is provided.
- Learner picks prediction  $\hat{y}_t \in \mathbb{R}$
- Outcome of the round  $z_t \in \mathcal{Z}$  is revealed
- Learner suffers loss  $\ell(\hat{y}_t, z_t)$

End For

Goal: minimize regret

$$\text{Reg}_n(\mathcal{F}) = \sum_{t=1}^n \ell(\hat{y}_t, z_t) - \inf_{f \in \mathcal{F}} \sum_{t=1}^n \ell(f(x_t), z_t)$$

# ONLINE SUPERVISED LEARNING: REDUCTION

$$\begin{aligned}\text{Reg}_n(\mathcal{F}) &= \sum_{t=1}^n \ell(\hat{y}_t, z_t) - \inf_{f \in \mathcal{F}} \sum_{t=1}^n \ell(f(x_t), z_t) \\&= \sup_{f \in \mathcal{F}} \sum_{t=1}^n (\ell(\hat{y}_t, z_t) - \ell(f(x_t), z_t)) \\&\leq \sup_{f \in \mathcal{F}} \sum_{t=1}^n \partial \ell(\hat{y}_t, z_t) \cdot (\hat{y}_t - f(x_t)) \\&= L \sup_{f \in \mathcal{F}} \sum_{t=1}^n \mathbb{E}_{b_t \sim \frac{1 + \partial \ell(\hat{y}_t, z_t)/L}{2}} [b_t \cdot (\hat{y}_t - f(x_t))] \\&\leq L \mathbb{E} \left[ \sup_{f \in \mathcal{F}} \sum_{t=1}^n b_t \cdot (\hat{y}_t - f(x_t)) \right] \\&= L \mathbb{E} \left[ \sum_{t=1}^n b_t \cdot \hat{y}_t - \inf_{f \in \mathcal{F}} \sum_{t=1}^n b_t \cdot f(x_t) \right]\end{aligned}$$

# ONLINE SUPERVISED LEARNING: REDUCTION

- So online supervised learning with any convex, L-Lipschitz loss we get guarantee that

$$\text{Reg}_n(\mathcal{F}) \leq L \text{Rad}_n(\mathcal{F})$$

- With a more complicated proof technique, one can show the same result but without requiring convexity of loss  $\ell$ .