

Lecture 26:

Deriving Randomized Algorithms from Relaxations

RECAP: RECIPE

- 1 Write down sequential Rademacher relaxation for the problem
- 2 Move to upper bound by cutting down the tree
- 3 Ensure that admissibility condition holds
- 4 Solve for the prediction given by relaxation based algorithm

ONLINE VS STATISTICAL LEARNING RATES

- Often optimal Online and statistical learning rates match

ONLINE VS STATISTICAL LEARNING RATES

- Often optimal Online and statistical learning rates match
- Get rid of tree by draw of future from fixed distribution D

$$\mathbf{Rad}_n(x_{1:t}, y_{1:t}) = \sup_{\mathbf{x}} \mathbb{E}_{\epsilon_{t+1:n}} \sup_{f \in \mathcal{F}} \left\{ 2 \sum_{s=t+1}^n \epsilon_s f(\mathbf{x}_{s-t}(\epsilon)) - \sum_{s=1}^t \ell(f(x_s), y_s) \right\}$$

ONLINE VS STATISTICAL LEARNING RATES

- Often optimal Online and statistical learning rates match
- Get rid of tree by draw of future from fixed distribution D

$$\mathbf{Rad}_n(x_{1:t}, y_{1:t}) = \sup_{\mathbf{x}} \mathbb{E}_{x_{t+1:n} \sim D} \sup_{f \in \mathcal{F}} \left\{ 2 \sum_{s=t+1}^n \epsilon_s f(\mathbf{x}_{s-t}(\epsilon)) - \sum_{s=1}^t \ell(f(x_s), y_s) \right\}$$

ONLINE VS STATISTICAL LEARNING RATES

- Often optimal Online and statistical learning rates match
- Get rid of tree by draw of future from fixed distribution D

$$\mathbf{Rad}_n(x_{1:t}, y_{1:t}) = \sup_{\mathbf{x}} \mathbb{E}_{\epsilon_{t+1:n}} \sup_{f \in \mathcal{F}} \left\{ 2C \sum_{s=t+1}^n \epsilon_s f(\mathbf{x}_{s-t}(\epsilon)) - \sum_{s=1}^t \ell(f(x_s), y_s) \right\}$$

- Assume loss ℓ is convex and 1-Lipchitz in first argument

RANDOM PLAYOUT

Define $R_t = x_{t+1:n}, \epsilon_{t+1:n}$ and let $D_t = D^{n-t} \times \text{Unif}\{\pm 1\}^{n-t}$

$$\phi_t(x_{1:t}, y_{1:t}; R_t) = \sup_{f \in \mathcal{F}} \left\{ 2C \sum_{s=t+1}^n \epsilon_s f(x_s) - \sum_{s=1}^t \ell(f(x_s), y_s) \right\}$$

Algorithm : **Draw** $R_t \sim D_t$, **and return**,

$$\tilde{q}_t(R_t) = \operatorname{argmin}_{q \in \Delta(\mathcal{Y})} \sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim q_t} [\ell(\hat{y}_t, y_t)] + \Phi_t(x_{1:t}, y_{1:t}, R_t) \right\}$$

RANDOM PLAYOUT

Define $R_t = x_{t+1:n}, \epsilon_{t+1:n}$ and let $D_t = D^{n-t} \times \text{Unif}\{\pm 1\}^{n-t}$

$$\phi_t(x_{1:t}, y_{1:t}; R_t) = \sup_{f \in \mathcal{F}} \left\{ 2C \sum_{s=t+1}^n \epsilon_s f(x_s) - \sum_{s=1}^t \ell(f(x_s), y_s) \right\}$$

Algorithm : **Draw** $R_t \sim D_t$, **and return,**

$$\tilde{q}_t(R_t) = \operatorname{argmin}_{q \in \Delta(\mathcal{Y})} \sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim q_t} [\ell(\hat{y}_t, y_t)] + \Phi_t(x_{1:t}, y_{1:t}, R_t) \right\}$$

Why/When does this work?

RANDOM PLAYOUT: CONDITION

Sufficient condition for randomized algorithm to work:

$$\begin{aligned} & \sup_{x_t} \mathbb{E}_{\epsilon_t} \left[\sup_{f \in \mathcal{F}} \left\{ 2C \sum_{s=t+1}^n \epsilon_s f(x_s) + 2\epsilon_t f(x_t) - \sum_{s=1}^{t-1} \ell(f(x_s), y_s) \right\} \right] \\ & \leq \mathbb{E}_{x_t \sim D, \epsilon_t} \left[\sup_{f \in \mathcal{F}} \left\{ 2C \sum_{s=t}^n \epsilon_s f(x_s) - \sum_{s=1}^{t-1} \ell(f(x_s), y_s) \right\} \right] \end{aligned}$$

RANDOM PLAYOUT

Initial condition is obvious, as for admissibility,

$$\inf_{q_t} \sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim q_t} [\ell(\hat{y}_t, y_t)] + \mathbf{Rel}_n(x_{1:t}, y_{1:t}) \right\}$$

RANDOM PLAYOUT

Initial condition is obvious, as for admissibility,

$$\inf_{q_t} \sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim q_t} [\ell(\hat{y}_t, y_t)] + \mathbf{Rel}_n(x_{1:t}, y_{1:t}) \right\}$$
$$= \inf_{q_t} \sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim q_t} [\ell(\hat{y}_t, y_t)] + \mathbb{E}_{R_t \sim D_t} [\Phi_t(x_{1:t}, y_{1:t}, R_t)] \right\}$$

RANDOM PLAYOUT

Initial condition is obvious, as for admissibility,

$$\begin{aligned} & \inf_{q_t} \sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim q_t} [\ell(\hat{y}_t, y_t)] + \mathbf{Rel}_n(x_{1:t}, y_{1:t}) \right\} \\ &= \inf_{q_t} \sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim q_t} [\ell(\hat{y}_t, y_t)] + \mathbb{E}_{R_t \sim D_t} [\Phi_t(x_{1:t}, y_{1:t}, R_t)] \right\} \\ &\leq \sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim \tilde{q}_t} [\ell(\hat{y}_t, y_t)] + \mathbb{E}_{R_t \sim D_t} [\Phi_t(x_{1:t}, y_{1:t}, R_t)] \right\} \end{aligned}$$

RANDOM PLAYOUT

Initial condition is obvious, as for admissibility,

$$\begin{aligned} & \inf_{q_t} \sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim q_t} [\ell(\hat{y}_t, y_t)] + \mathbf{Rel}_n(x_{1:t}, y_{1:t}) \right\} \\ &= \inf_{q_t} \sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim q_t} [\ell(\hat{y}_t, y_t)] + \mathbb{E}_{R_t \sim D_t} [\Phi_t(x_{1:t}, y_{1:t}, R_t)] \right\} \\ &\leq \sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim \tilde{q}_t} [\ell(\hat{y}_t, y_t)] + \mathbb{E}_{R_t \sim D_t} [\Phi_t(x_{1:t}, y_{1:t}, R_t)] \right\} \\ &= \sup_{y_t} \left\{ \mathbb{E}_{R_t \sim D_t} [\mathbb{E}_{\hat{y}_t \sim \tilde{q}_t(R_t)} [\ell(\hat{y}_t, y_t)]] + \mathbb{E}_{R_t \sim D_t} [\Phi_t(x_{1:t}, y_{1:t}, R_t)] \right\} \end{aligned}$$

RANDOM PLAYOUT

Initial condition is obvious, as for admissibility,

$$\begin{aligned} & \inf_{q_t} \sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim q_t} [\ell(\hat{y}_t, y_t)] + \mathbf{Rel}_n(x_{1:t}, y_{1:t}) \right\} \\ &= \inf_{q_t} \sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim q_t} [\ell(\hat{y}_t, y_t)] + \mathbb{E}_{R_t \sim D_t} [\Phi_t(x_{1:t}, y_{1:t}, R_t)] \right\} \\ &\leq \sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim \tilde{q}_t} [\ell(\hat{y}_t, y_t)] + \mathbb{E}_{R_t \sim D_t} [\Phi_t(x_{1:t}, y_{1:t}, R_t)] \right\} \\ &= \sup_{y_t} \left\{ \mathbb{E}_{R_t \sim D_t} [\mathbb{E}_{\hat{y}_t \sim \tilde{q}_t(R_t)} [\ell(\hat{y}_t, y_t)]] + \mathbb{E}_{R_t \sim D_t} [\Phi_t(x_{1:t}, y_{1:t}, R_t)] \right\} \\ &\leq \mathbb{E}_{R_t \sim D_t} \left[\sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim \tilde{q}_t(R_t)} [\ell(\hat{y}_t, y_t)] + \Phi_t(x_{1:t}, y_{1:t}, R_t) \right\} \right] \end{aligned}$$

RANDOM PLAYOUT

Initial condition is obvious, as for admissibility,

$$\begin{aligned} & \inf_{q_t} \sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim q_t} [\ell(\hat{y}_t, y_t)] + \mathbf{Rel}_n(x_{1:t}, y_{1:t}) \right\} \\ &= \inf_{q_t} \sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim q_t} [\ell(\hat{y}_t, y_t)] + \mathbb{E}_{R_t \sim D_t} [\Phi_t(x_{1:t}, y_{1:t}, R_t)] \right\} \\ &\leq \sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim \tilde{q}_t} [\ell(\hat{y}_t, y_t)] + \mathbb{E}_{R_t \sim D_t} [\Phi_t(x_{1:t}, y_{1:t}, R_t)] \right\} \\ &= \sup_{y_t} \left\{ \mathbb{E}_{R_t \sim D_t} [\mathbb{E}_{\hat{y}_t \sim \tilde{q}_t(R_t)} [\ell(\hat{y}_t, y_t)]] + \mathbb{E}_{R_t \sim D_t} [\Phi_t(x_{1:t}, y_{1:t}, R_t)] \right\} \\ &\leq \mathbb{E}_{R_t \sim D_t} \left[\sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim \tilde{q}_t(R_t)} [\ell(\hat{y}_t, y_t)] + \Phi_t(x_{1:t}, y_{1:t}, R_t) \right\} \right] \\ &= \mathbb{E}_{R_t \sim D_t} \left[\inf_{q_t} \sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim q_t} [\ell(\hat{y}_t, y_t)] + \Phi_t(x_{1:t}, y_{1:t}, R_t) \right\} \right] \end{aligned}$$

RANDOM PLAYOUT

Initial condition is obvious, as for admissibility,

$$\begin{aligned} & \inf_{q_t} \sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim q_t} [\ell(\hat{y}_t, y_t)] + \mathbf{Rel}_n(x_{1:t}, y_{1:t}) \right\} \\ &= \inf_{q_t} \sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim q_t} [\ell(\hat{y}_t, y_t)] + \mathbb{E}_{R_t \sim D_t} [\Phi_t(x_{1:t}, y_{1:t}, R_t)] \right\} \\ &\leq \sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim \tilde{q}_t} [\ell(\hat{y}_t, y_t)] + \mathbb{E}_{R_t \sim D_t} [\Phi_t(x_{1:t}, y_{1:t}, R_t)] \right\} \\ &= \sup_{y_t} \left\{ \mathbb{E}_{R_t \sim D_t} [\mathbb{E}_{\hat{y}_t \sim \tilde{q}_t(R_t)} [\ell(\hat{y}_t, y_t)]] + \mathbb{E}_{R_t \sim D_t} [\Phi_t(x_{1:t}, y_{1:t}, R_t)] \right\} \\ &\leq \mathbb{E}_{R_t \sim D_t} \left[\sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim \tilde{q}_t(R_t)} [\ell(\hat{y}_t, y_t)] + \Phi_t(x_{1:t}, y_{1:t}, R_t) \right\} \right] \\ &= \mathbb{E}_{R_t \sim D_t} \left[\inf_{q_t} \sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim q_t} [\ell(\hat{y}_t, y_t)] + \Phi_t(x_{1:t}, y_{1:t}, R_t) \right\} \right] \\ &= \mathbb{E}_{R_t \sim D_t} \left[\sup_{p_t} \left\{ \inf_{\hat{y}_t \in \mathcal{Y}} \mathbb{E}_{y_t \sim p_t} [\ell(\hat{y}_t, y_t)] + \mathbb{E}_{y_t \sim p_t} [\Phi_t(x_{1:t}, y_{1:t}, R_t)] \right\} \right] \end{aligned}$$

RANDOM PLAYOUT

To finish admissibility, note that

$$\sup_{p_t} \left\{ \inf_{\hat{y}_t \in \mathcal{Y}} \mathbb{E}_{y_t \sim p_t} [\ell(\hat{y}_t, y_t)] + \mathbb{E}_{y_t \sim p_t} [\Phi_t(x_{1:t}, y_{1:t}, R_t)] \right\}$$

RANDOM PLAYOUT

To finish admissibility, note that

$$\begin{aligned} & \sup_{p_t} \left\{ \inf_{\hat{y}_t \in \mathcal{Y}} \mathbb{E}_{y_t \sim p_t} [\ell(\hat{y}_t, y_t)] + \mathbb{E}_{y_t \sim p_t} [\Phi_t(x_{1:t}, y_{1:t}, R_t)] \right\} \\ &= \sup_{p_t} \left\{ \inf_{\hat{y}_t \in \mathcal{Y}} \mathbb{E}_{y_t \sim p_t} [\ell(\hat{y}_t, y_t)] + \mathbb{E}_{y_t \sim p_t} \left[\sup_{f \in \mathcal{F}} \left\{ 2C \sum_{s=t+1}^n \epsilon_s f(x_s) - \sum_{s=1}^t \ell(f(x_s), y_s) \right\} \right] \right\} \end{aligned}$$

RANDOM PLAYOUT

To finish admissibility, note that

$$\begin{aligned} & \sup_{p_t} \left\{ \inf_{\hat{y}_t \in \mathcal{Y}} \mathbb{E}_{y_t \sim p_t} [\ell(\hat{y}_t, y_t)] + \mathbb{E}_{y_t \sim p_t} [\Phi_t(x_{1:t}, y_{1:t}, R_t)] \right\} \\ &= \sup_{p_t} \left\{ \inf_{\hat{y}_t \in \mathcal{Y}} \mathbb{E}_{y_t \sim p_t} [\ell(\hat{y}_t, y_t)] + \mathbb{E}_{y_t \sim p_t} \left[\sup_{f \in \mathcal{F}} \left\{ 2C \sum_{s=t+1}^n \epsilon_s f(x_s) - \sum_{s=1}^t \ell(f(x_s), y_s) \right\} \right] \right\} \\ &\leq \sup_{x_t} \mathbb{E}_{\epsilon_t} \left[\sup_{f \in \mathcal{F}} \left\{ 2C \sum_{s=t+1}^n \epsilon_s f(x_s) + 2\epsilon_t f(x_t) - \sum_{s=1}^{t-1} \ell(f(x_s), y_s) \right\} \right] \end{aligned}$$

RANDOM PLAYOUT

To finish admissibility, note that

$$\begin{aligned} & \sup_{p_t} \left\{ \inf_{\hat{y}_t \in \mathcal{Y}} \mathbb{E}_{y_t \sim p_t} [\ell(\hat{y}_t, y_t)] + \mathbb{E}_{y_t \sim p_t} [\Phi_t(x_{1:t}, y_{1:t}, R_t)] \right\} \\ &= \sup_{p_t} \left\{ \inf_{\hat{y}_t \in \mathcal{Y}} \mathbb{E}_{y_t \sim p_t} [\ell(\hat{y}_t, y_t)] + \mathbb{E}_{y_t \sim p_t} \left[\sup_{f \in \mathcal{F}} \left\{ 2C \sum_{s=t+1}^n \epsilon_s f(x_s) - \sum_{s=1}^t \ell(f(x_s), y_s) \right\} \right] \right\} \\ &\leq \sup_{x_t} \mathbb{E}_{\epsilon_t} \left[\sup_{f \in \mathcal{F}} \left\{ 2C \sum_{s=t+1}^n \epsilon_s f(x_s) + 2\epsilon_t f(x_t) - \sum_{s=1}^{t-1} \ell(f(x_s), y_s) \right\} \right] \end{aligned}$$

Condition:

$$\begin{aligned} & \sup_{x_t} \mathbb{E}_{\epsilon_t} \left[\sup_{f \in \mathcal{F}} \left\{ 2C \sum_{s=t+1}^n \epsilon_s f(x_s) + 2\epsilon_t f(x_t) - \sum_{s=1}^{t-1} \ell(f(x_s), y_s) \right\} \right] \\ &\leq \mathbb{E}_{x_t \sim D, \epsilon_t} \left[\sup_{f \in \mathcal{F}} \left\{ 2C \sum_{s=t}^n \epsilon_s f(x_s) - \sum_{s=1}^{t-1} \ell(f(x_s), y_s) \right\} \right] \end{aligned}$$

RANDOM PLAYOUT

Hence,

$$\begin{aligned} & \sup_{x_t} \inf_{q_t} \sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim q_t} [\ell(\hat{y}_t, y_t)] + \mathbf{Rel}_n(x_{1:t}, y_{1:t}) \right\} \\ & \leq \sup_{x_t} \mathbb{E}_{R_t \sim D_t} \left[\sup_{p_t} \left\{ \inf_{\hat{y}_t \in \mathcal{Y}} \mathbb{E}_{y_t \sim p_t} [\ell(\hat{y}_t, y_t)] + \mathbb{E}_{y_t \sim p_t} [\Phi_t(x_{1:t}, y_{1:t}, R_t)] \right\} \right] \end{aligned}$$

RANDOM PLAYOUT

Hence,

$$\begin{aligned} & \sup_{x_t} \inf_{q_t} \sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim q_t} [\ell(\hat{y}_t, y_t)] + \mathbf{Rel}_n(x_{1:t}, y_{1:t}) \right\} \\ & \leq \sup_{x_t} \mathbb{E}_{R_t \sim D_t} \left[\sup_{p_t} \left\{ \inf_{\hat{y}_t \in \mathcal{Y}} \mathbb{E}_{y_t \sim p_t} [\ell(\hat{y}_t, y_t)] + \mathbb{E}_{y_t \sim p_t} [\Phi_t(x_{1:t}, y_{1:t}, R_t)] \right\} \right] \\ & \leq \mathbb{E}_{R_t \sim D_t} \left[\mathbb{E}_{x_t \sim D, \epsilon_t} \left[\sup_{f \in \mathcal{F}} \left\{ 2C \sum_{s=t}^n \epsilon_s f(x_s) - \sum_{s=1}^{t-1} \ell(f(x_s), y_s) \right\} \right] \right] \end{aligned}$$

RANDOM PLAYOUT

Hence,

$$\begin{aligned} & \sup_{x_t} \inf_{q_t} \sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim q_t} [\ell(\hat{y}_t, y_t)] + \mathbf{Rel}_n(x_{1:t}, y_{1:t}) \right\} \\ & \leq \sup_{x_t} \mathbb{E}_{R_t \sim D_t} \left[\sup_{p_t} \left\{ \inf_{\hat{y}_t \in \mathcal{Y}} \mathbb{E}_{y_t \sim p_t} [\ell(\hat{y}_t, y_t)] + \mathbb{E}_{y_t \sim p_t} [\Phi_t(x_{1:t}, y_{1:t}, R_t)] \right\} \right] \\ & \leq \mathbb{E}_{R_t \sim D_t} \left[\mathbb{E}_{x_t \sim D, \epsilon_t} \left[\sup_{f \in \mathcal{F}} \left\{ 2C \sum_{s=t}^n \epsilon_s f(x_s) - \sum_{s=1}^{t-1} \ell(f(x_s), y_s) \right\} \right] \right] \\ & = \mathbb{E}_{R_{t-1} \sim D_{t-1}} [\Phi_t(x_{1:t-1}, y_{1:t-1}, R_{t-1})] \\ & = \mathbf{Rel}_n(x_{1:t-1}, y_{1:t-1}) \end{aligned}$$

EXAMPLE: BIT PREDICTION

- $\mathcal{F} \subset \{\pm 1\}^n$ $\mathcal{X} = \{\}$, $\ell(y', y) = \mathbf{1}\{y \neq y'\} = \frac{1-y \cdot y'}{2}$
- Since there are no x 's the condition is obvious.
- Algorithm : at round t , draw $\epsilon_{t+1:n}$ then play

$$2q_t(\epsilon) - 1$$

$$\begin{aligned} &= \sup_{f \in \mathcal{F}} \left\{ \sum_{s=t+1}^n \epsilon_s f_s - \sum_{s=1}^{t-1} \mathbf{1}\{f_s \neq y_s\} - \mathbf{1}\{f_t \neq 1\} \right\} - \sup_{f \in \mathcal{F}} \left\{ \sum_{s=t+1}^n \epsilon_s f_s - \sum_{s=1}^{t-1} \mathbf{1}\{f_s \neq y_s\} - \mathbf{1}\{f_t \neq -1\} \right\} \\ &= \inf_{f \in \mathcal{F}} \left\{ 2 \sum_{s=t+1}^n \mathbf{1}\{\epsilon_s \neq f_s\} + \sum_{s=1}^{t-1} \mathbf{1}\{f_s \neq y_s\} + \mathbf{1}\{f_t \neq 1\} \right\} \\ &\quad - \inf_{f \in \mathcal{F}} \left\{ 2 \sum_{s=t+1}^n \mathbf{1}\{\epsilon_s \neq f_s\} + \sum_{s=1}^{t-1} \mathbf{1}\{f_s \neq y_s\} + \mathbf{1}\{f_t \neq -1\} \right\} \end{aligned}$$

Solve two ERM's per round.

EXAMPLE: LINEAR PREDICTORS

- Online linear optimization, $\mathcal{F} = \{f : \|f\| \leq 1\}$, $\mathbf{D} = \{\nabla : \|\nabla\|_* \leq 1\}$
- Condition: $\exists D$ and constant C , such that, for any vector w ,

$$\sup_{x_t} \mathbb{E}_{\epsilon_t} [\|w + 2\epsilon_t x_t\|_*] \leq \mathbb{E}_{x_t \sim D} [\|w + Cx_t\|_*]$$

EXAMPLE: LINEAR PREDICTORS

- Online linear optimization, $\mathcal{F} = \{f : \|f\| \leq 1\}$, $\mathbf{D} = \{\nabla : \|\nabla\|_* \leq 1\}$
- Condition: $\exists D$ and constant C , such that, for any vector w ,

$$\sup_{x_t} \mathbb{E}_{\epsilon_t} [\|w + 2\epsilon_t x_t\|_*] \leq \mathbb{E}_{x_t \sim D} [\|w + Cx_t\|_*]$$

- $\ell_1^d / \ell_\infty^d : D = \text{Unif}\{\pm 1\}^d$ or any other symmetric distribution on each coordinate (Eg. normal distribution)

EXAMPLE: LINEAR PREDICTORS

- Online linear optimization, $\mathcal{F} = \{f : \|f\| \leq 1\}$, $\mathbf{D} = \{\nabla : \|\nabla\|_* \leq 1\}$
- Condition: $\exists D$ and constant C , such that, for any vector w ,

$$\sup_{x_t} \mathbb{E}_{\epsilon_t} [\|w + 2\epsilon_t x_t\|_*] \leq \mathbb{E}_{x_t \sim D} [\|w + Cx_t\|_*]$$

- ℓ_1^d / ℓ_∞^d : $D = \text{Unif}\{\pm 1\}^d$ or any other symmetric distribution on each coordinate (Eg. normal distribution)
- Algorithm : Round t draw $R_t \sim N(0, (n - t)I_d)$

$$\hat{y}_t = \operatorname{argmin}_{i \in [d]} \left| \sum_{j=1}^t \nabla_j[i] + R_t[i] \right|$$

EXAMPLE: LINEAR PREDICTORS

- Online linear optimization, $\mathcal{F} = \{f : \|f\| \leq 1\}$, $\mathbf{D} = \{\nabla : \|\nabla\|_* \leq 1\}$
- Condition: $\exists D$ and constant C , such that, for any vector w ,

$$\sup_{x_t} \mathbb{E}_{\epsilon_t} [\|w + 2\epsilon_t x_t\|_*] \leq \mathbb{E}_{x_t \sim D} [\|w + Cx_t\|_*]$$

- ℓ_1^d / ℓ_∞^d : $D = \text{Unif}\{\pm 1\}^d$ or any other symmetric distribution on each coordinate (Eg. normal distribution)
- Algorithm : Round t draw $R_t \sim N(0, (n - t)I_d)$

$$\hat{y}_t = \operatorname{argmin}_{i \in [d]} \left| \sum_{j=1}^t \nabla_j[i] + R_t[i] \right|$$

- Bound : $\mathbb{E}[\text{Reg}_n] \leq \frac{1}{n} \mathbf{Rel}_n(\cdot) = O\left(\sqrt{\frac{\log d}{n}}\right)$

ROUGH SKETCH OF PROOF

- $w = 2C \sum_{s=t+1}^n \nabla_s - \sum_{s=1}^{t-1} \nabla_s$ where $\nabla_{1:t-1}$ are past losses and $\nabla_{t+1:n}$ are drawn from $\text{Unif}\{-1, 1\}^d$

ROUGH SKETCH OF PROOF

- $w = 2C \sum_{s=t+1}^n \nabla_s - \sum_{s=1}^{t-1} \nabla_s$ where $\nabla_{1:t-1}$ are past losses and $\nabla_{t+1:n}$ are drawn from $\text{Unif}\{-1, 1\}^d$
- Assume $t < n - \sqrt{n}$, for last \sqrt{n} rounds even if we are completely off, regret bound does not change

ROUGH SKETCH OF PROOF

- $w = 2C \sum_{s=t+1}^n \nabla_s - \sum_{s=1}^{t-1} \nabla_s$ where $\nabla_{1:t-1}$ are past losses and $\nabla_{t+1:n}$ are drawn from $\text{Unif}\{-1, 1\}^d$
- Assume $t < n - \sqrt{n}$, for last \sqrt{n} rounds even if we are completely off, regret bound does not change
- Hence w can be seen as vector $-\sum_{s=1}^{t-1} \nabla_s$ where each coordinate is perturbed by $2C \sum_{s=t+1}^n \nabla_s$

ROUGH SKETCH OF PROOF

- $w = 2C \sum_{s=t+1}^n \nabla_s - \sum_{s=1}^{t-1} \nabla_s$ where $\nabla_{1:t-1}$ are past losses and $\nabla_{t+1:n}$ are drawn from $\text{Unif}\{-1, 1\}^d$
- Assume $t < n - \sqrt{n}$, for last \sqrt{n} rounds even if we are completely off, regret bound does not change
- Hence w can be seen as vector $-\sum_{s=1}^{t-1} \nabla_s$ where each coordinate is perturbed by $2C \sum_{s=t+1}^n \nabla_s$
- With very high probability, if i^* and j^* are top two coordinates of w , $|w[i^*]| - |w[j^*]| > 4$, hence, with high probability,

$$\begin{aligned} \sup_{x_t \in [-1, 1]^d} \mathbb{E}_{\epsilon_t} [\|w + 2\epsilon_t x_t\|_{\infty}] &= \sup_{x_t \in [-1, 1]^d} \mathbb{E}_{\epsilon_t} [|\mathcal{W}[i^*] + 2\epsilon_t x_t[i^*]|] \\ &= \mathbb{E}_{\epsilon_t} [|\mathcal{W}[i^*] + 2\epsilon_t|] = \mathbb{E}_{x_t \sim D} [\|w + 2\epsilon_t x_t\|_{\infty}] \end{aligned}$$

ROUGH SKETCH OF PROOF

- $w = 2C \sum_{s=t+1}^n \nabla_s - \sum_{s=1}^{t-1} \nabla_s$ where $\nabla_{1:t-1}$ are past losses and $\nabla_{t+1:n}$ are drawn from $\text{Unif}\{-1, 1\}^d$
- Assume $t < n - \sqrt{n}$, for last \sqrt{n} rounds even if we are completely off, regret bound does not change
- Hence w can be seen as vector $-\sum_{s=1}^{t-1} \nabla_s$ where each coordinate is perturbed by $2C \sum_{s=t+1}^n \nabla_s$
- With very high probability, if i^* and j^* are top two coordinates of w , $|w[i^*]| - |w[j^*]| > 4$, hence, with high probability,

$$\begin{aligned} \sup_{x_t \in [-1, 1]^d} \mathbb{E}_{\epsilon_t} [\|w + 2\epsilon_t x_t\|_{\infty}] &= \sup_{x_t \in [-1, 1]^d} \mathbb{E}_{\epsilon_t} [|w[i^*]| + 2\epsilon_t x_t[i^*] |] \\ &= \mathbb{E}_{\epsilon_t} [|w[i^*]| + 2\epsilon_t |] = \mathbb{E}_{x_t \sim D} [\|w + 2\epsilon_t x_t\|_{\infty}] \end{aligned}$$

- In general we don't need this high probability stuff, we can directly prove the condition, just need to check cases.

ROUGH SKETCH OF PROOF

- Why update of form $\hat{y}_t = \operatorname{argmin}_{i \in [d]} |\sum_{j=1}^t \nabla_j[i] + R_t[i]|$
- To see this, note that the algorithm we need is originally of form,

$$\hat{y}_t = \operatorname{argmin}_{\hat{y} \in \mathcal{F}} \sup_{\nabla_t} \left\{ \langle \hat{y}, \nabla_t \rangle + \sup_{f \in \mathcal{F}} \left\{ \langle f, -R_t \rangle - \left\langle f, \sum_{s=1}^t \nabla_s \right\rangle \right\} \right\}$$

ROUGH SKETCH OF PROOF

- Why update of form $\hat{y}_t = \operatorname{argmin}_{i \in [d]} |\sum_{j=1}^t \nabla_j[i] + R_t[i]|$
- To see this, note that the algorithm we need is originally of form,

$$\begin{aligned}\hat{y}_t &= \operatorname{argmin}_{\hat{y} \in \mathcal{F}} \sup_{\nabla_t} \left\{ \langle \hat{y}, \nabla_t \rangle + \sup_{f \in \mathcal{F}} \left\{ \langle f, -R_t \rangle - \left\langle f, \sum_{s=1}^t \nabla_s \right\rangle \right\} \right\} \\ &= \operatorname{argmin}_{\hat{y} \in \mathcal{F}} \sup_{f \in \mathcal{F}} \left\{ \sup_{\nabla_t} \langle \hat{y} - f, \nabla_t \rangle + \left\langle f, -R_t - \sum_{s=1}^{t-1} \nabla_s \right\rangle \right\}\end{aligned}$$

ROUGH SKETCH OF PROOF

- Why update of form $\hat{y}_t = \operatorname{argmin}_{i \in [d]} |\sum_{j=1}^t \nabla_j[i] + R_t[i]|$
- To see this, note that the algorithm we need is originally of form,

$$\begin{aligned}\hat{y}_t &= \operatorname{argmin}_{\hat{y} \in \mathcal{F}} \sup_{\nabla_t} \left\{ \langle \hat{y}, \nabla_t \rangle + \sup_{f \in \mathcal{F}} \left\{ \langle f, -R_t \rangle - \left\langle f, \sum_{s=1}^t \nabla_s \right\rangle \right\} \right\} \\ &= \operatorname{argmin}_{\hat{y} \in \mathcal{F}} \sup_{f \in \mathcal{F}} \left\{ \sup_{\nabla_t} \langle \hat{y} - f, \nabla_t \rangle + \left\langle f, -R_t - \sum_{s=1}^{t-1} \nabla_s \right\rangle \right\} \\ &= \operatorname{argmin}_{\hat{y} \in \mathcal{F}} \sup_{f \in \mathcal{F}} \left\{ \|\hat{y} - f\|_\infty - \left\langle f, R_t + \sum_{s=1}^{t-1} \nabla_s \right\rangle \right\}\end{aligned}$$

EXAMPLE: LINEAR PREDICTORS

- Online linear optimization, $\mathcal{F} = \{f : \|f\| \leq 1\}$, $\mathbf{D} = \{\nabla : \|\nabla\|_* \leq 1\}$
- Condition: $\exists D$ and constant C , such that, for any vector w ,

$$\sup_{x_t} \mathbb{E}_{\epsilon_t} [\|w + 2\epsilon_t x_t\|_*] \leq \mathbb{E}_{x_t \sim D} [\|w + Cx_t\|_*]$$

- ℓ_2/ℓ_2 : $D = \text{Unif}\{\text{unit sphere}\}$ or normalized Gaussian distribution
- Algorithm : Round t draw $R_t \sim N(0, (n - t)I_d)/\sqrt{d}$

$$\hat{y}_t = \operatorname{argmin}_{f: \|f\|_2 \leq 1} \left\langle f, \sum_{j=1}^t \nabla_j + R_t \right\rangle$$

- Bound : $\mathbb{E}[\text{Reg}_n] \leq \frac{1}{n} \mathbf{Rel}_n(\cdot) = O\left(\sqrt{\frac{1}{n}}\right)$

EXAMPLE: FINITE EXPERTS

- Very similar to ℓ_1/ℓ_∞ , think about subtracting -1 from every loss, makes no difference for regret
- But then ℓ_1/ℓ_∞ is same as finite experts
- Algorithm : Round t draw $R_t \sim N(0, (n - t)I_{|\mathcal{F}|})$

$$\hat{y}_t = \operatorname{argmin}_{i \in [d]} \sum_{j=1}^t \ell(i, z_j) + R_t[i]$$

- Bound : $\mathbb{E}[\operatorname{Reg}_n] \leq \frac{1}{n} \mathbf{Rel}_n(\cdot) = O\left(\sqrt{\frac{\log |\mathcal{F}|}{n}}\right)$

EXAMPLE: ONLINE SHORTEST PATH

- Graph $G = (V, E)$, source node S and destination node D .

EXAMPLE: ONLINE SHORTEST PATH

- Graph $G = (V, E)$, source node S and destination node D .
- Every round, we need to pick a path from S to D

EXAMPLE: ONLINE SHORTEST PATH

- Graph $G = (V, E)$, source node S and destination node D .
- Every round, we need to pick a path from S to D
- Adversary picks a delay on every edge $W : E \mapsto [0, 1]$

EXAMPLE: ONLINE SHORTEST PATH

- Graph $G = (V, E)$, source node S and destination node D .
- Every round, we need to pick a path from S to D
- Adversary picks a delay on every edge $W : E \mapsto [0, 1]$
- Learner suffers delay on path chosen which is sum of delays on edges of the path

EXAMPLE: ONLINE SHORTEST PATH

- Graph $G = (V, E)$, source node S and destination node D .
- Every round, we need to pick a path from S to D
- Adversary picks a delay on every edge $W : E \mapsto [0, 1]$
- Learner suffers delay on path chosen which is sum of delays on edges of the path
- Experts bound $|E| \sqrt{\frac{|V| \log |V|}{n}}$

EXAMPLE: ONLINE SHORTEST PATH

- Graph $G = (V, E)$, source node S and destination node D .
- Every round, we need to pick a path from S to D
- Adversary picks a delay on every edge $W : E \mapsto [0, 1]$
- Learner suffers delay on path chosen which is sum of delays on edges of the path
- Experts bound $|E| \sqrt{\frac{|V| \log |V|}{n}}$
- However naive time complexity $O(\#paths)$

EXAMPLE: ONLINE SHORTEST PATH

- Can view it as a different online linear optimization problem
- $\mathcal{F} = \{f \in \{0, 1\}^{|E|} : f \text{ is a path}\}$
- $\mathbf{D} = [0, 1]^{|E|}$ the delays on each edge.

EXAMPLE: ONLINE SHORTEST PATH

- Can view it as a different online linear optimization problem
- $\mathcal{F} = \{f \in \{0, 1\}^{|E|} : f \text{ is a path}\}$
- $\mathbf{D} = [0, 1]^{|E|}$ the delays on each edge.
- Random playout condition satisfied by distribution $D = N(0, 1)$

EXAMPLE: ONLINE SHORTEST PATH

- Can view it as a different online linear optimization problem
- $\mathcal{F} = \{f \in \{0, 1\}^{|E|} : f \text{ is a path}\}$
- $\mathbf{D} = [0, 1]^{|E|}$ the delays on each edge.
- Random playout condition satisfied by distribution $D = N(0, 1)$
- Algorithm: Draw $R_t \sim N(0, (n - t)I_{|E|})$,

$$\text{path}_t = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \left\langle f, \sum_{j=1}^{t-1} \nabla_j + R_t \right\rangle$$

EXAMPLE: ONLINE SHORTEST PATH

- Can view it as a different online linear optimization problem
- $\mathcal{F} = \{f \in \{0, 1\}^{|E|} : f \text{ is a path}\}$
- $\mathbf{D} = [0, 1]^{|E|}$ the delays on each edge.
- Random playout condition satisfied by distribution $D = N(0, 1)$
- Algorithm: Draw $R_t \sim N(0, (n - t)I_{|E|})$,

$$\text{path}_t = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \left\langle f, \sum_{j=1}^{t-1} \nabla_j + R_t \right\rangle$$

- That is solve shortest path algorithm with delay on edge $e \in E$ given by $\sum_{j=1}^{t-1} \nabla_j[e] + R_t[e]$
- Can be solves in poly-time using Bellman-ford algorithm.

LEARNING WITH NON-REPEATED ENTRIES

For $t = 1$ to $|\mathcal{X}|$

Adversary picks $x_t \in \mathcal{X} \setminus \{x_1, \dots, x_{t-1}\}$

Learner predicts $q_t \in \Delta(\mathcal{Y})$

Adversary picks $y_t \in \mathcal{Y}$

Learner draws $\hat{y}_t \sim q_t$ and suffers loss $\ell(\hat{y}_t, y_t)$

End

Regret :

$$\text{Reg}_{|\mathcal{X}|} = \sum_{t=1}^{|\mathcal{X}|} \ell(\hat{y}_t, y_t) - \inf_{f \in \mathcal{F}} \sum_{t=1}^{|\mathcal{X}|} \ell(f(x_t), y_t)$$

LEARNING WITH NON-REPEATED ENTRIES

- For convex Lipschitz loss and binary loss, the symmetrization idea just goes through, only on each path, no node is repeated.
- Sequential Rademacher relaxation:

$$\mathbf{Rad}_{|\mathcal{X}|}(x_{1:t}, y_{1:t}) = \sup_{\mathbf{x}} \mathbb{E}_{\epsilon_{t+1:n}} \sup_{f \in \mathcal{F}} \left\{ \sum_{s=t+1}^{|\mathcal{X}|} \epsilon_s f(\mathbf{x}_{s-t}(\epsilon)) - \sum_{s=1}^t \ell(f(x_s), y_s) \right\}$$

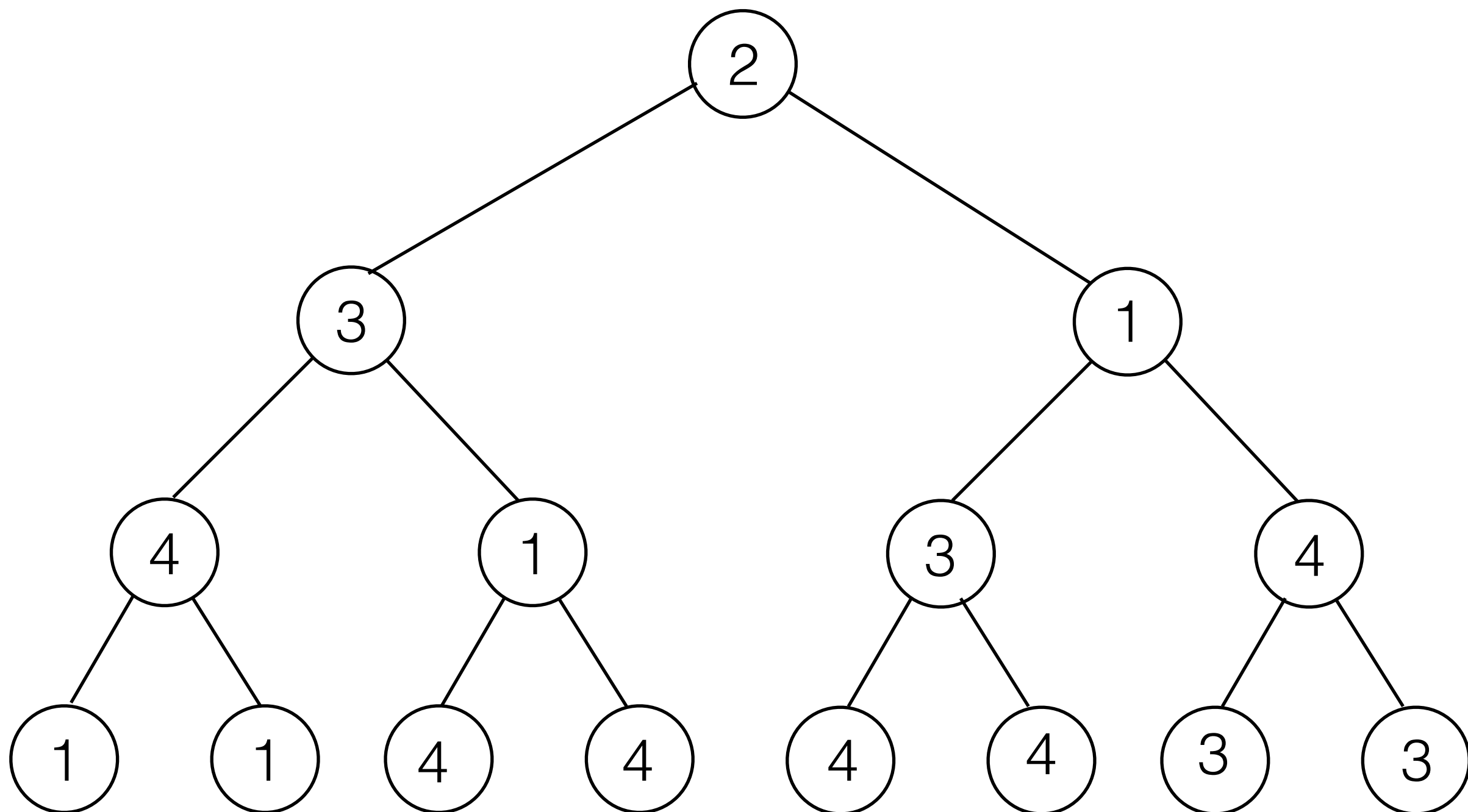
where \mathbf{x} is a tree with values in $\mathcal{X} \setminus \{x_1, \dots, x_t\}$ with no node repeated on any path.

LEARNING WITH NON-REPEATED ENTRIES

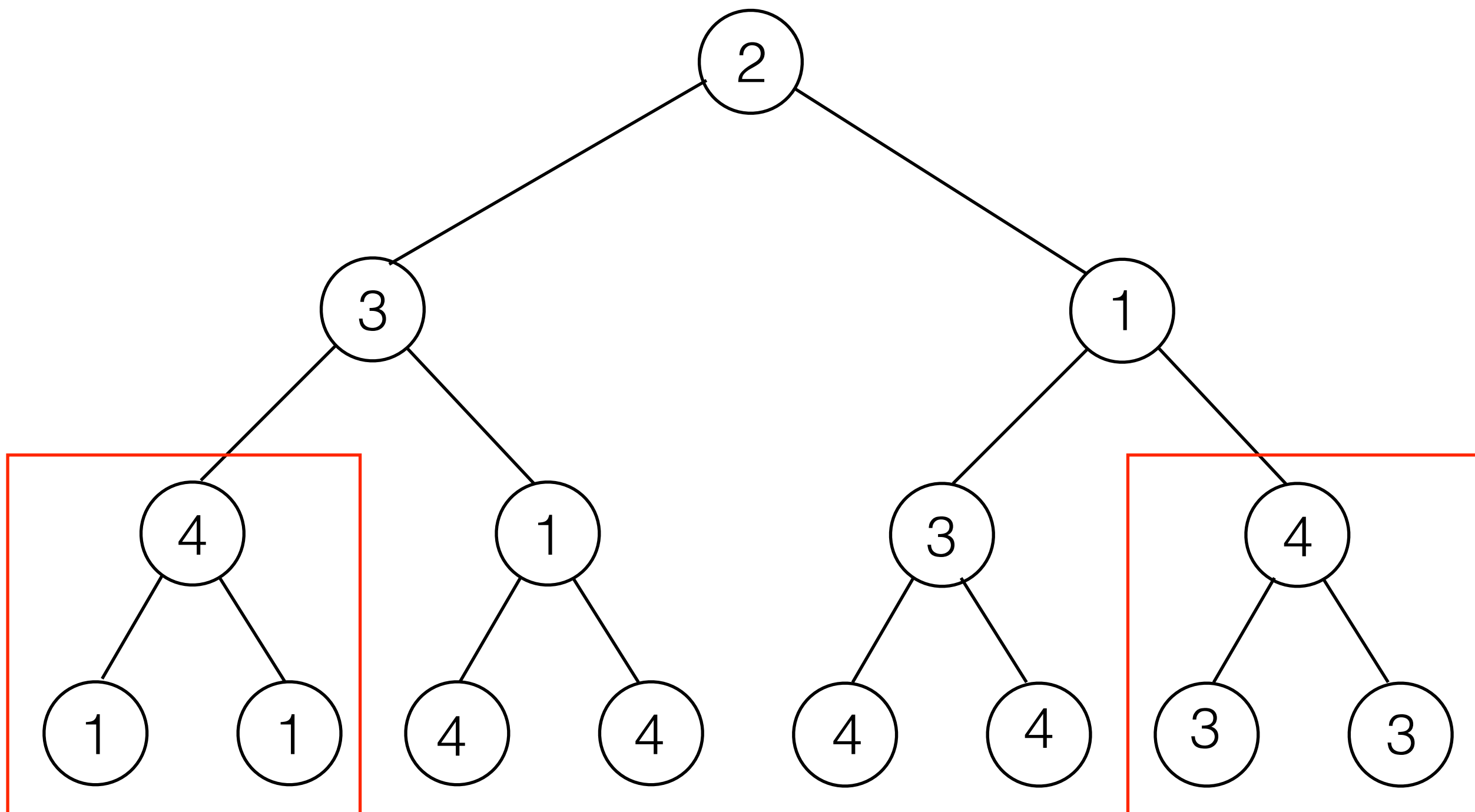
- Inductively we can show that:

$$\mathbf{Rad}_{|\mathcal{X}|}(x_{1:t}, y_{1:t}) = \mathbb{E}_{\epsilon_{t+1:n}} \sup_{f \in \mathcal{F}} \left\{ \sum_{s=t+1}^{|\mathcal{X}|} \epsilon_s f(x_s) - \sum_{s=1}^t \ell(f(x_s), y_s) \right\}$$

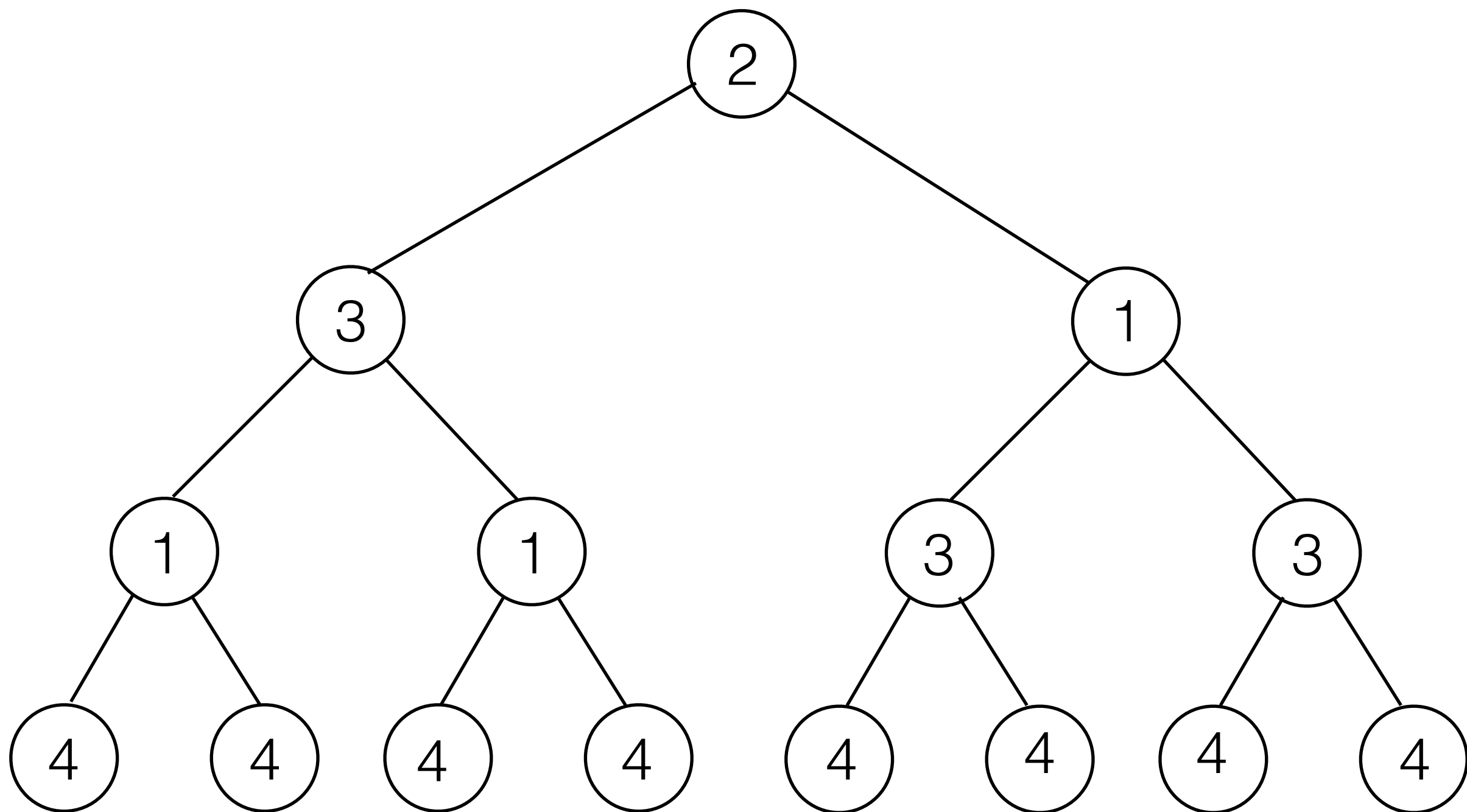
where $x_{t+1}, \dots, x_{|\mathcal{X}|}$ are elements from $\mathcal{X} \setminus \{x_1, \dots, x_t\}$ in any order non-repeated.



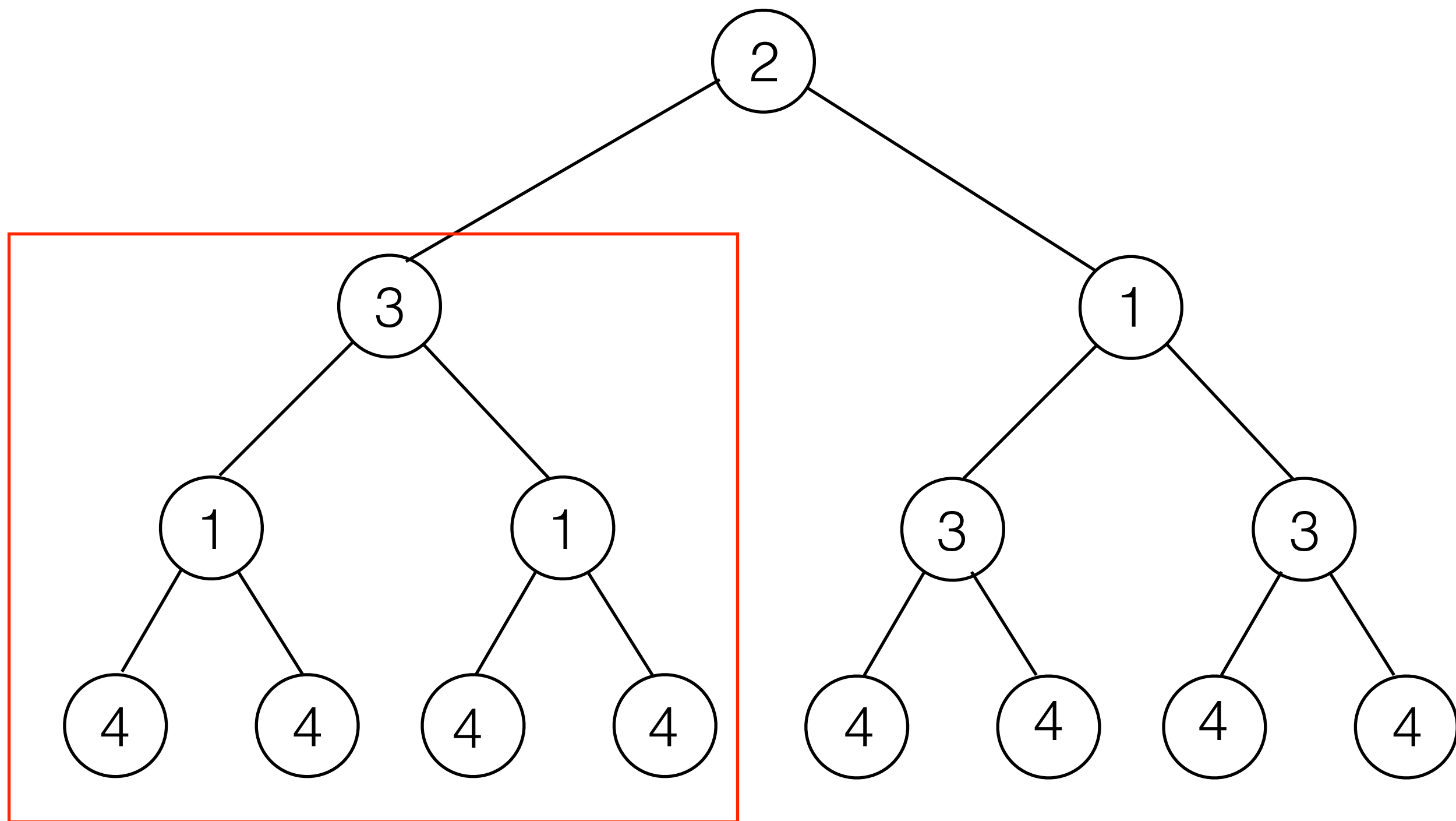
$$\mathbb{E} \sup_{\epsilon_{t+1:n} f \in \mathcal{F}} \left\{ \sum_{s=t+1}^{|\mathcal{X}|} \epsilon_s f(\mathbf{x}_{s-t}(\epsilon)) - \sum_{s=1}^t \ell(f(x_s), y_s) \right\}$$



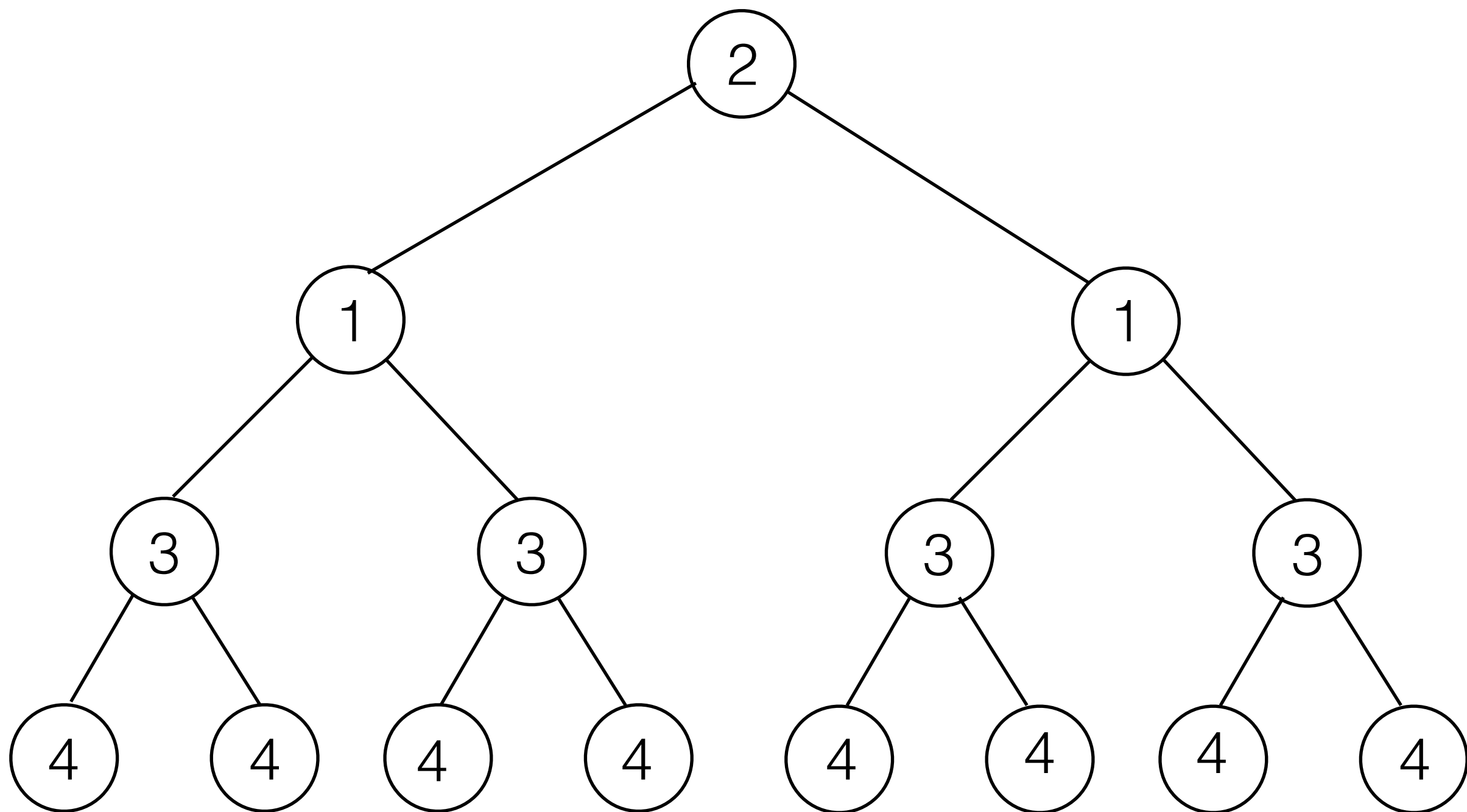
$$\mathbb{E} \sup_{\epsilon_{t+1:n} f \in \mathcal{F}} \left\{ \sum_{s=t+1}^{|\mathcal{X}|} \epsilon_s f(\mathbf{x}_{s-t}(\epsilon)) - \sum_{s=1}^t \ell(f(x_s), y_s) \right\}$$



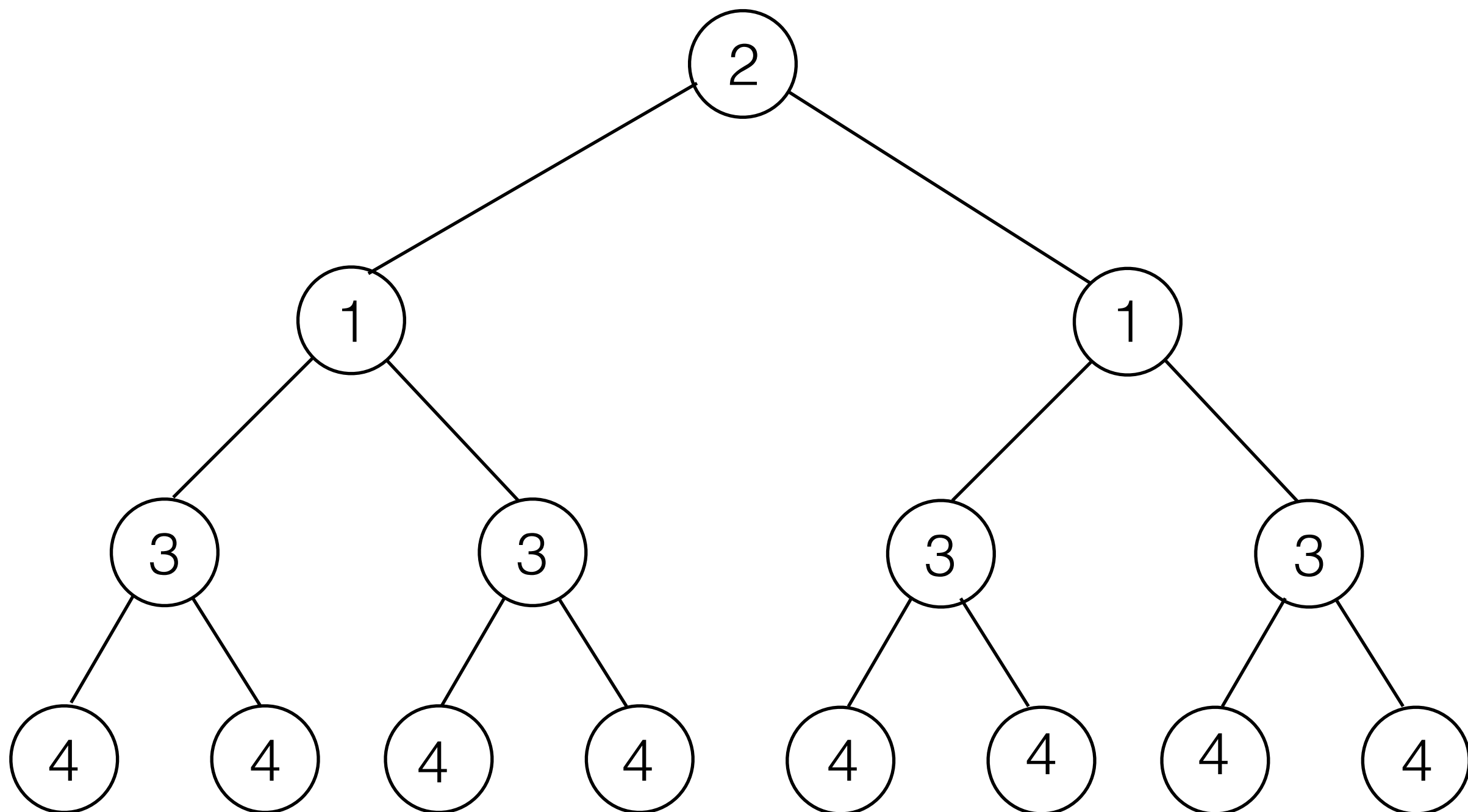
$$\mathbb{E} \sup_{\epsilon_{t+1:n}} \sup_{f \in \mathcal{F}} \left\{ \sum_{s=t+1}^{|\mathcal{X}|} \epsilon_s f(\mathbf{x}_{s-t}(\epsilon)) - \sum_{s=1}^t \ell(f(x_s), y_s) \right\}$$



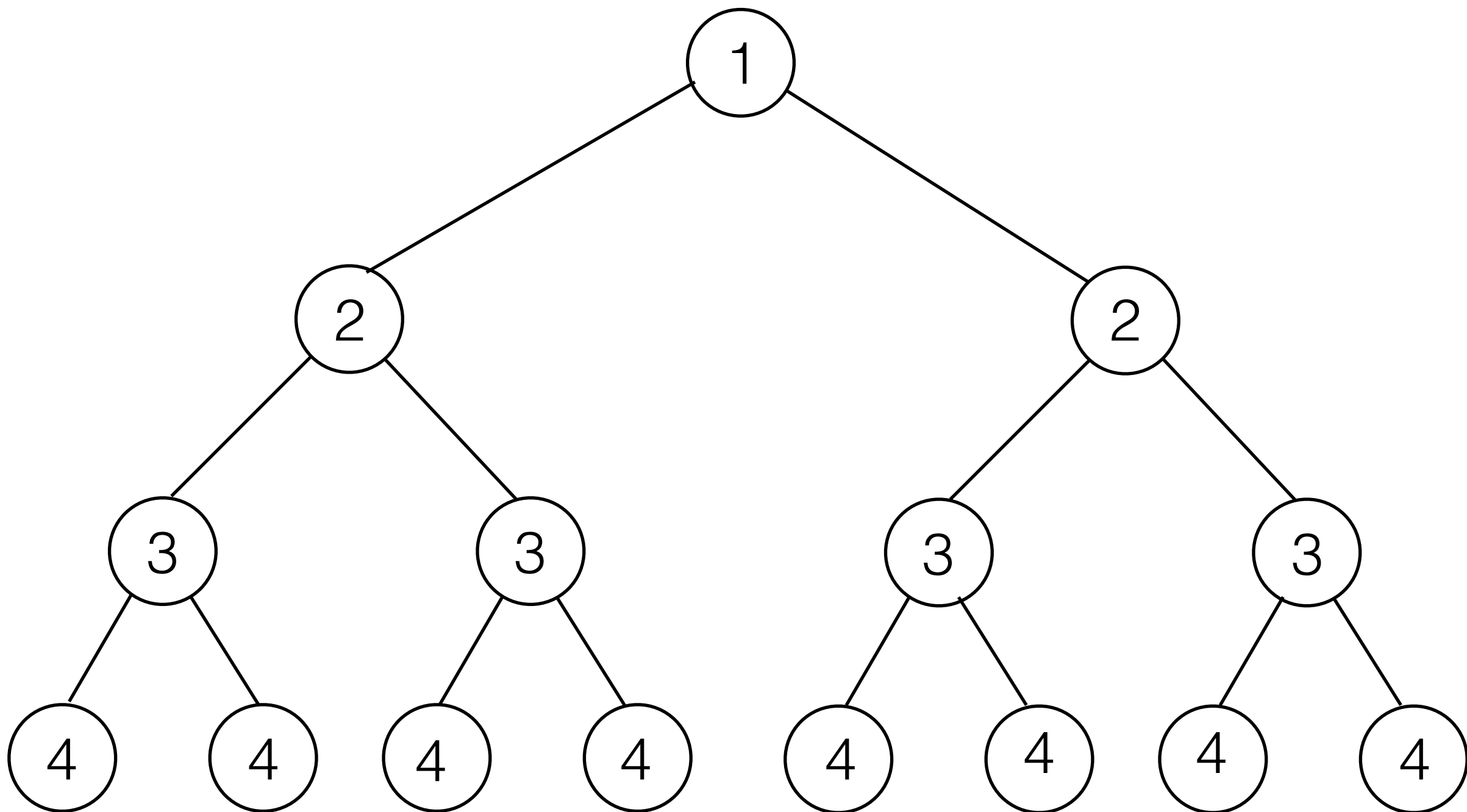
$$\mathbb{E} \sup_{\epsilon_{t+1:n} f \in \mathcal{F}} \left\{ \sum_{s=t+1}^{|\mathcal{X}|} \epsilon_s f(\mathbf{x}_{s-t}(\epsilon)) - \sum_{s=1}^t \ell(f(x_s), y_s) \right\}$$



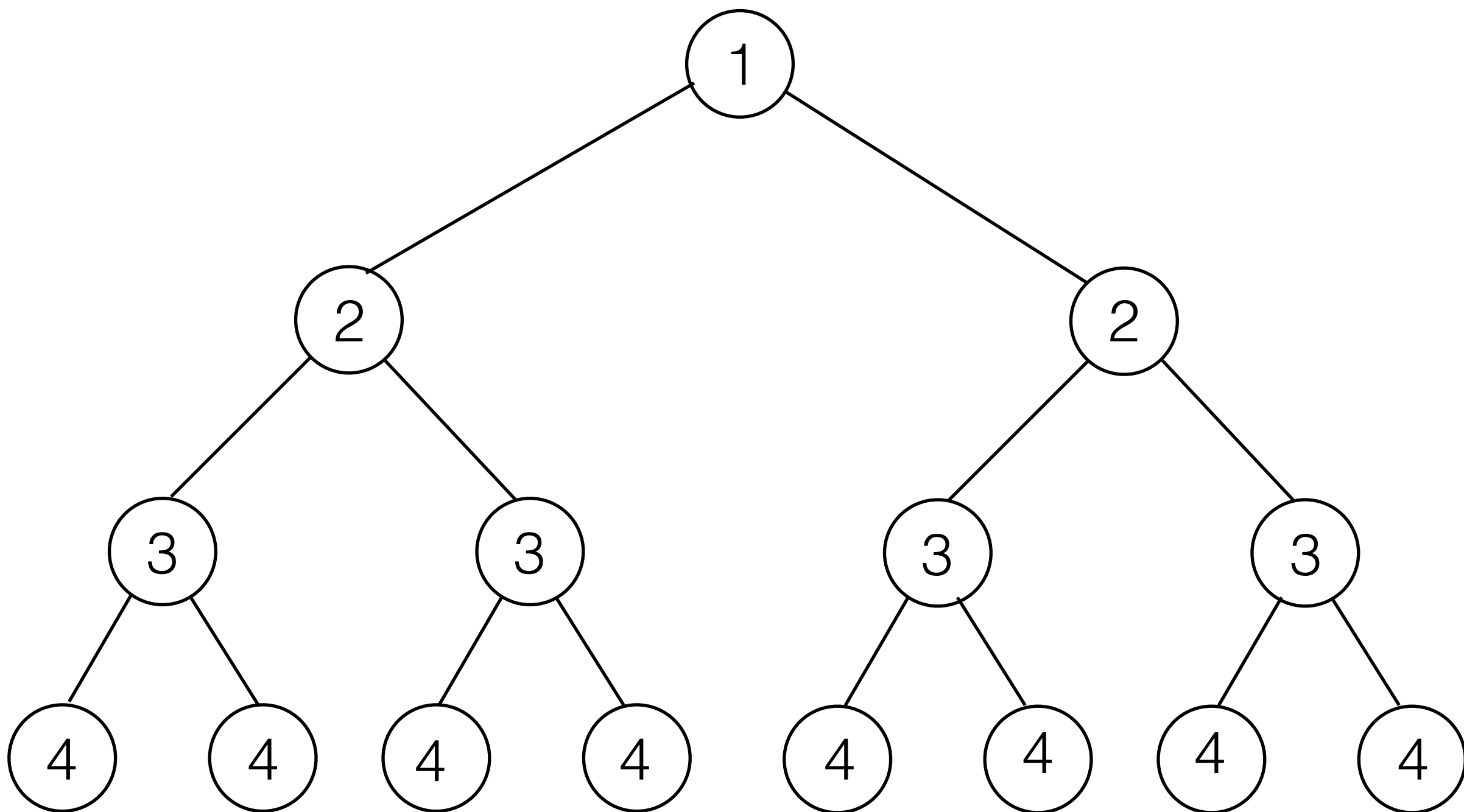
$$\mathbb{E} \sup_{\epsilon_{t+1:n}} \sup_{f \in \mathcal{F}} \left\{ \sum_{s=t+1}^{|\mathcal{X}|} \epsilon_s f(\mathbf{x}_{s-t}(\epsilon)) - \sum_{s=1}^t \ell(f(x_s), y_s) \right\}$$



$$\mathbb{E} \sup_{\epsilon_{t+1:n} f \in \mathcal{F}} \left\{ \sum_{s=t+1}^{|\mathcal{X}|} \epsilon_s f(\mathbf{x}_{s-t}(\epsilon)) - \sum_{s=1}^t \ell(f(x_s), y_s) \right\}$$



$$\mathbb{E} \sup_{\epsilon_{t+1:n}} \sup_{f \in \mathcal{F}} \left\{ \sum_{s=t+1}^{|\mathcal{X}|} \epsilon_s f(\mathbf{x}_{s-t}(\epsilon)) - \sum_{s=1}^t \ell(f(x_s), y_s) \right\}$$



$$= \mathbb{E} \sup_{\epsilon_{t+1:n} f \in \mathcal{F}} \left\{ \sum_{s=t+1}^{|\mathcal{X}|} \epsilon_s f(x_s) - \sum_{s=1}^t \ell(f(x_s), y_s) \right\}$$

LEARNING WITH NON-REPEATED ENTRIES

- Inductively we can show that:

$$\mathbf{Rad}_{|\mathcal{X}|}(x_{1:t}, y_{1:t}) = \mathbb{E}_{\epsilon_{t+1:n}} \sup_{f \in \mathcal{F}} \left\{ \sum_{s=t+1}^{|\mathcal{X}|} \epsilon_s f(x_s) - \sum_{s=1}^t \ell(f(x_s), y_s) \right\}$$

where $x_{t+1}, \dots, x_{|\mathcal{X}|}$ are elements from $\mathcal{X} \setminus \{x_1, \dots, x_t\}$ in any order non-repeated.

LEARNING WITH NON-REPEATED ENTRIES

- Inductively we can show that:

$$\mathbf{Rad}_{|\mathcal{X}|}(x_{1:t}, y_{1:t}) = \mathbb{E}_{\epsilon_{t+1:n}} \sup_{f \in \mathcal{F}} \left\{ \sum_{s=t+1}^{|\mathcal{X}|} \epsilon_s f(x_s) - \sum_{s=1}^t \ell(f(x_s), y_s) \right\}$$

where $x_{t+1}, \dots, x_{|\mathcal{X}|}$ are elements from $\mathcal{X} \setminus \{x_1, \dots, x_t\}$ in any order non-repeated.

- We can use $\mathbf{Rel}_{|\mathcal{X}|}(x_{1:t}, y_{1:t}) = \mathbf{Rad}_{|\mathcal{X}|}(x_{1:t}, y_{1:t})$ as a relaxation

LEARNING WITH NON-REPEATED ENTRIES

- Inductively we can show that:

$$\mathbf{Rad}_{|\mathcal{X}|}(x_{1:t}, y_{1:t}) = \mathbb{E}_{\epsilon_{t+1:n}} \sup_{f \in \mathcal{F}} \left\{ \sum_{s=t+1}^{|\mathcal{X}|} \epsilon_s f(x_s) - \sum_{s=1}^t \ell(f(x_s), y_s) \right\}$$

where $x_{t+1}, \dots, x_{|\mathcal{X}|}$ are elements from $\mathcal{X} \setminus \{x_1, \dots, x_t\}$ in any order non-repeated.

- We can use $\mathbf{Rel}_{|\mathcal{X}|}(x_{1:t}, y_{1:t}) = \mathbf{Rad}_{|\mathcal{X}|}(x_{1:t}, y_{1:t})$ as a relaxation
- Condition satisfied trivially, with constant 1,

$$\begin{aligned} \sup_{x_t \in \mathcal{X} \setminus \{x_1, \dots, x_{t-1}\}} \mathbb{E}_{\epsilon_t} \left[\sup_{f \in \mathcal{F}} \left\{ 2 \sum_{s=t+1}^{|\mathcal{X}|} \epsilon_s f(x_s) + 2\epsilon_t f(x_t) - \sum_{s=1}^{t-1} \ell(f(x_s), y_s) \right\} \right] \\ = \mathbb{E}_{\epsilon_t} \left[\sup_{f \in \mathcal{F}} \left\{ 2 \sum_{s=t}^{|\mathcal{X}|} \epsilon_s f(x_s) - \sum_{s=1}^{t-1} \ell(f(x_s), y_s) \right\} \right] \end{aligned}$$

because the sum $2 \sum_{s=t}^{|\mathcal{X}|} \epsilon_s f(x_s)$ is independent of order.

LEARNING WITH NON-REPEATED ENTRIES

- Algorithm: Fix some order over elements of \mathcal{X} . On each round t , draw $\epsilon_{t+1}, \dots, \epsilon_{|\mathcal{X}|}$.

- Solve

$$q_t = \operatorname{argmin}_{q \in \Delta(\mathcal{Y})} \sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim q_t} [\ell(\hat{y}_t, y_t)] + \sup_{f \in \mathcal{F}} \left\{ 2 \sum_{s=t+1}^n \epsilon_s f(x_s) - \sum_{s=1}^t \ell(f(x_s), y_s) \right\} \right\}$$

- Bound : $\mathbb{E}[\operatorname{Reg}_n] \leq \mathcal{R}_n^{\text{stat}}(\mathcal{F})$

LEARNING WITH NON-REPEATED ENTRIES

- Algorithm: Fix some order over elements of \mathcal{X} . On each round t , draw $\epsilon_{t+1}, \dots, \epsilon_{|\mathcal{X}|}$.


















- Solve

$$q_t = \operatorname{argmin}_{q \in \Delta(\mathcal{Y})} \sup_{y_t} \left\{ \mathbb{E}_{\hat{y}_t \sim q_t} [\ell(\hat{y}_t, y_t)] + \sup_{f \in \mathcal{F}} \left\{ 2 \sum_{s=t+1}^n \epsilon_s f(x_s) - \sum_{s=1}^t \ell(f(x_s), y_s) \right\} \right\}$$

- Bound : $\mathbb{E}[\operatorname{Reg}_n] \leq \mathcal{R}_n^{\operatorname{stat}}(\mathcal{F})$
- Example: binary classification



















$$q_t = \frac{1}{2} + \frac{1}{2} \sup_{f \in \mathcal{F}} \left\{ \sum_{s=t+1}^n \epsilon_s f(x_s) + \frac{1}{2} \sum_{s=1}^{t-1} y_s f(x_s) + \frac{1}{2} f(x_t) \right\} \\ - \frac{1}{2} \sup_{f \in \mathcal{F}} \left\{ \sum_{s=t+1}^n \epsilon_s f(x_s) + \frac{1}{2} \sum_{s=1}^{t-1} y_s f(x_s) - \frac{1}{2} f(x_t) \right\}$$

ONLINE COLLABORATIVE FILTERING

for $t = 1$ to n



















ONLINE COLLABORATIVE FILTERING

for $t = 1$ to n

Entry to predict $x_t = (\text{User}, \text{Wii})$

ONLINE COLLABORATIVE FILTERING

















					
					
					
					
					

for $t = 1$ to n

Entry to predict $x_t = (\text{User}, \text{Item})$

Learner picks $\hat{y}_t \in [-1, 1]$

ONLINE COLLABORATIVE FILTERING

for $t = 1$ to n


















Entry to predict $x_t = (\text{woman}, \text{Wii})$

Learner picks $\hat{y}_t \in [-1, 1]$

True rating $y_t \in \{\pm 1\}$ revealed

Learner suffers loss $|\hat{y}_t - y_t|$

ONLINE COLLABORATIVE FILTERING

for $t = 1$ to n

Entry to predict $x_t = (\text{User}, \text{Item})$

Learner picks $\hat{y}_t \in [-1, 1]$


















True rating $y_t \in \{\pm 1\}$ revealed

Learner suffers loss $|\hat{y}_t - y_t|$



















$$\text{Reg}_n := \frac{1}{n} \sum_{t=1}^n |\hat{y}_t - y_t| - \inf_{M: \|M\| \leq B} \frac{1}{n} \sum_{t=1}^n |M[x_t] - y_t|$$

($\|\cdot\|$: trace norm)

























ONLINE COLLABORATIVE FILTERING

























ONLINE COLLABORATIVE FILTERING

























ONLINE COLLABORATIVE FILTERING

























					
					
					
					
					

ONLINE COLLABORATIVE FILTERING

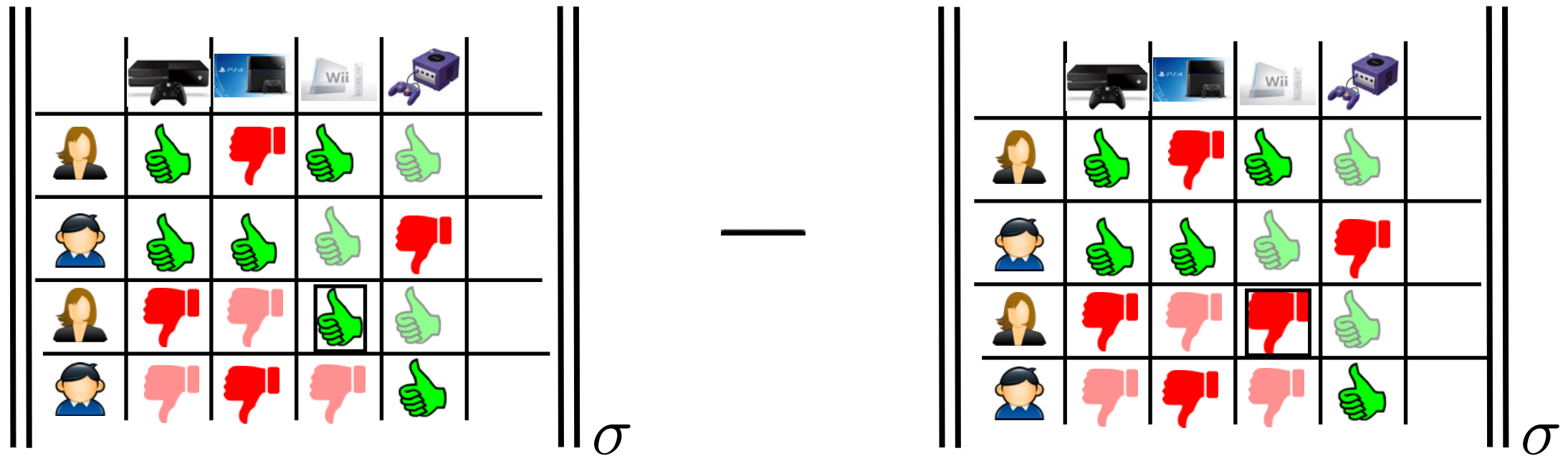
					
					
					
					
					

ONLINE COLLABORATIVE FILTERING

ONLINE COLLABORATIVE FILTERING



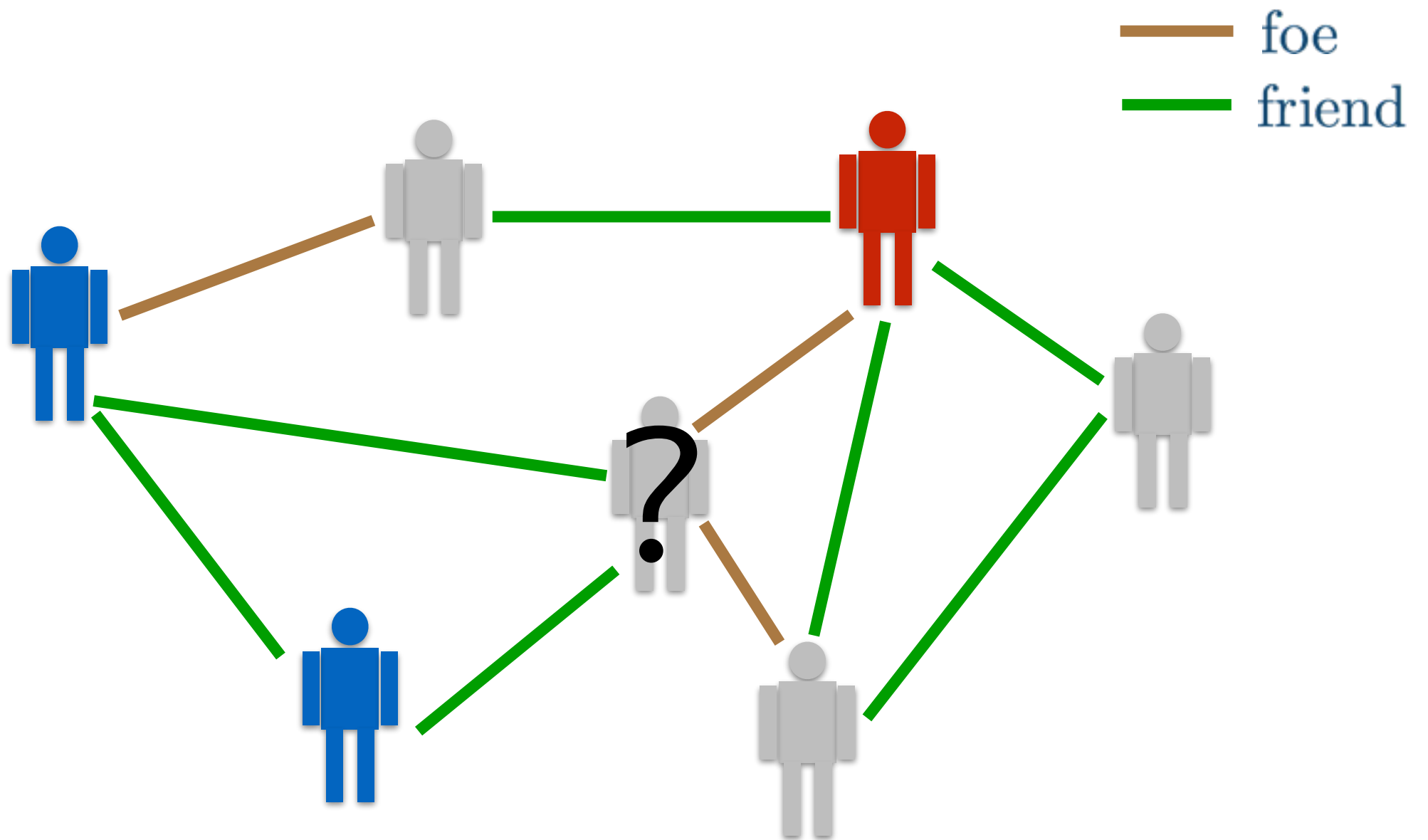
ONLINE COLLABORATIVE FILTERING

- M users and N products, regret bound :

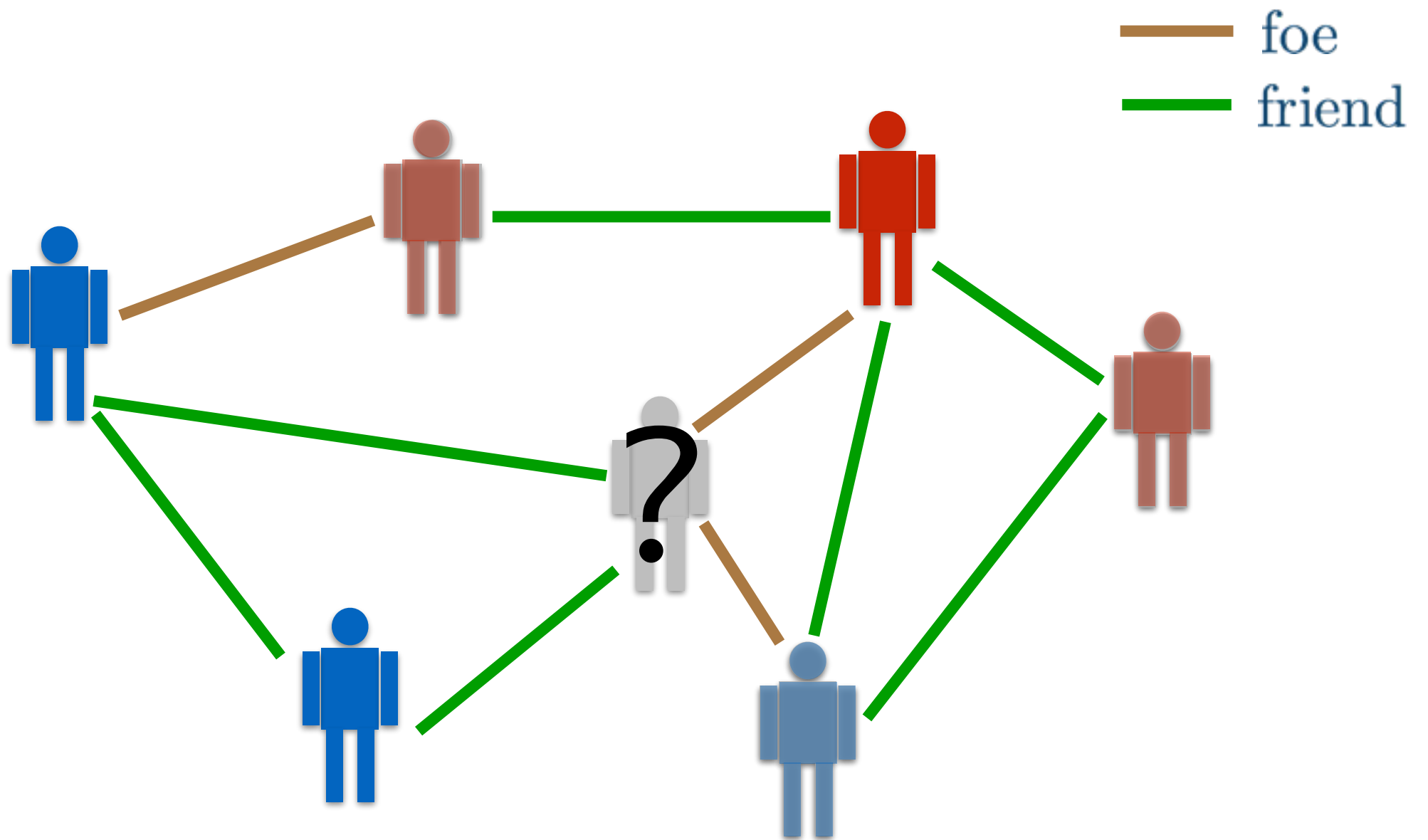
$$\mathbb{E} [\text{Reg}_n] \leq \frac{B\sqrt{M+N}}{n}$$

- Statistical learning : for same rate, require assumption that user product pair is uniformly distributed [Srebro & Shraibman'05]
- Improves over [Cesa-Bianchi & Shamir'11], [Hazan et al'12] both in terms of regret bound and time complexity.
- Algorithm for online edge prediction and link classification in social networks (adjacency matrix)

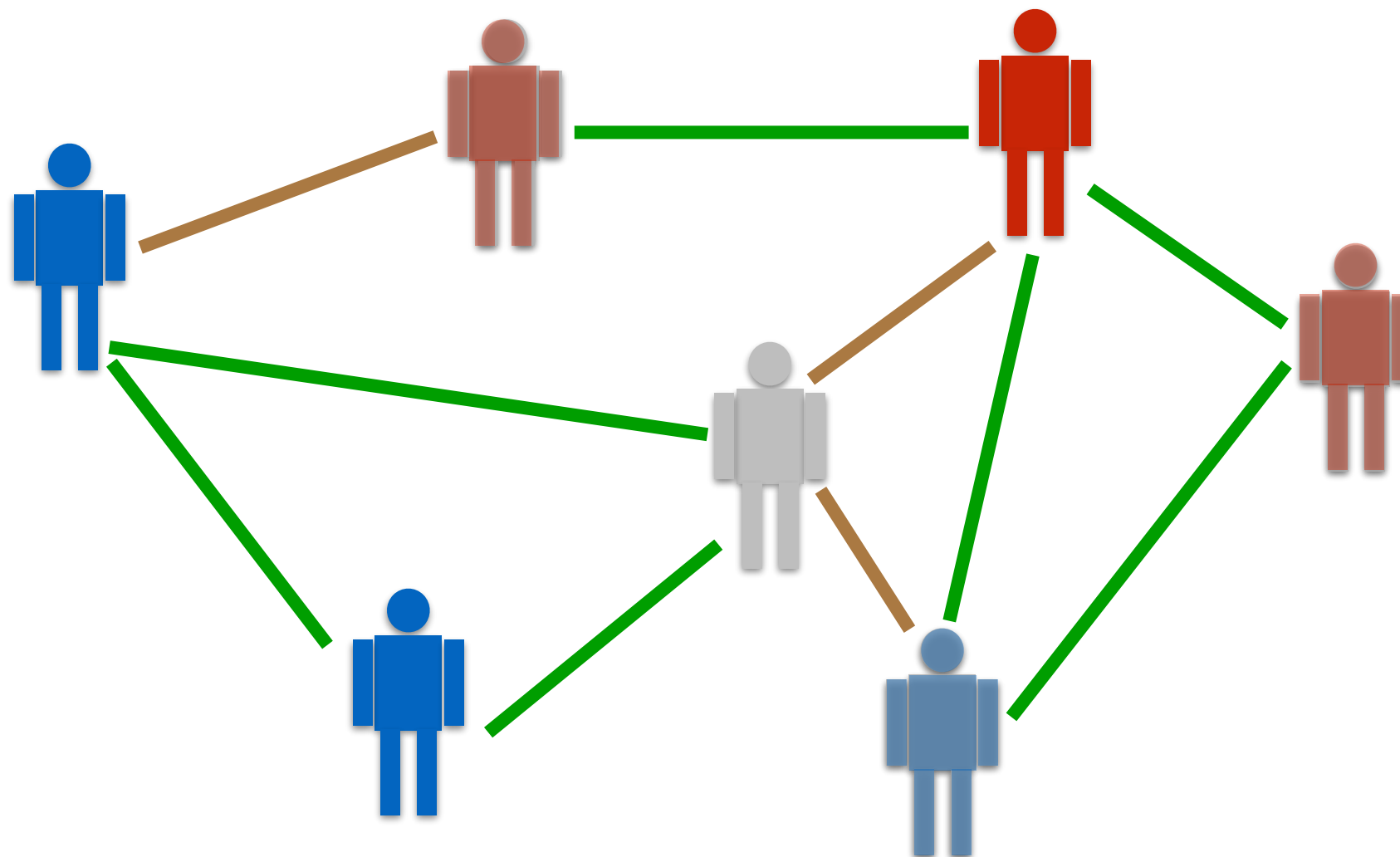
ONLINE NODE CLASSIFICATION



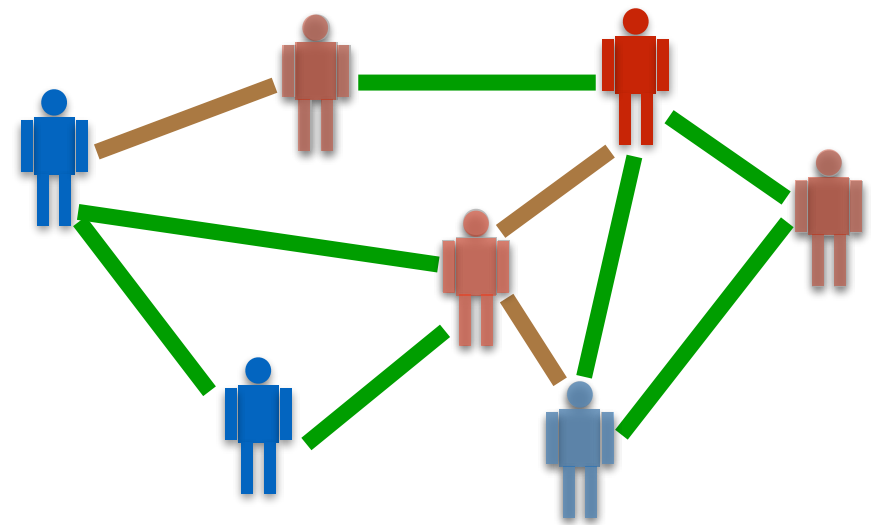
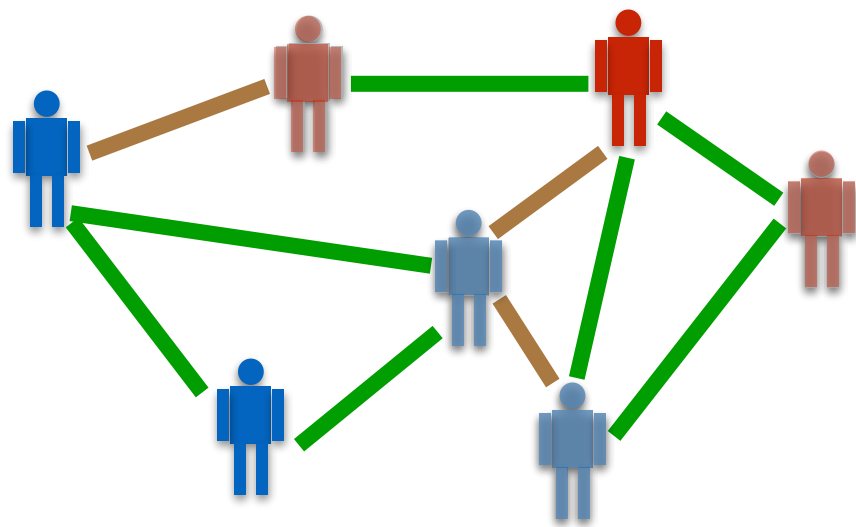
ONLINE NODE CLASSIFICATION



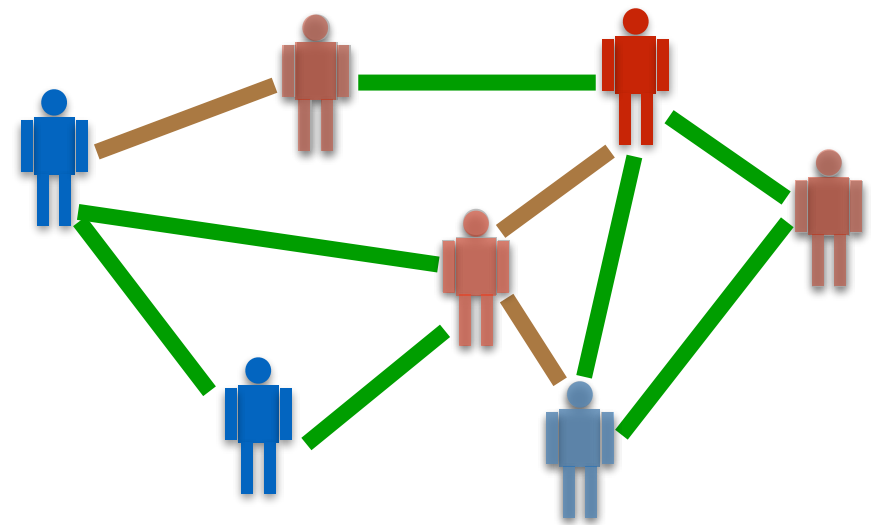
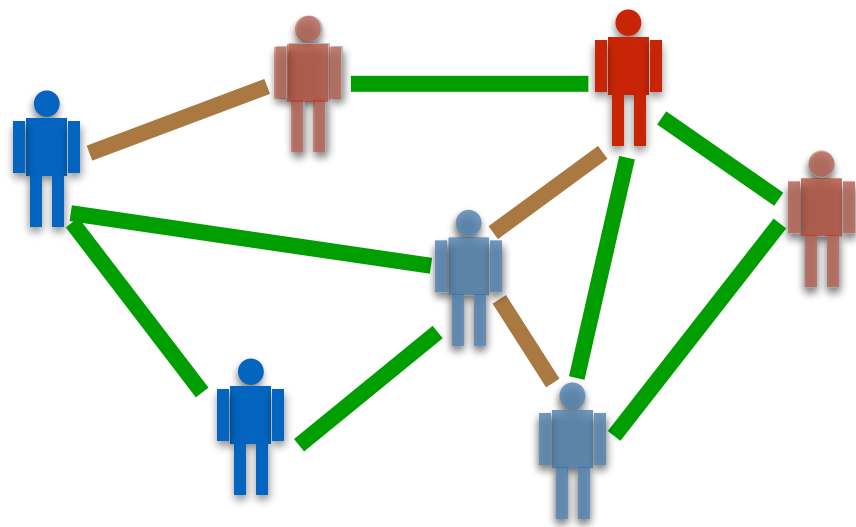
ONLINE NODE CLASSIFICATION



ONLINE NODE CLASSIFICATION



ONLINE NODE CLASSIFICATION



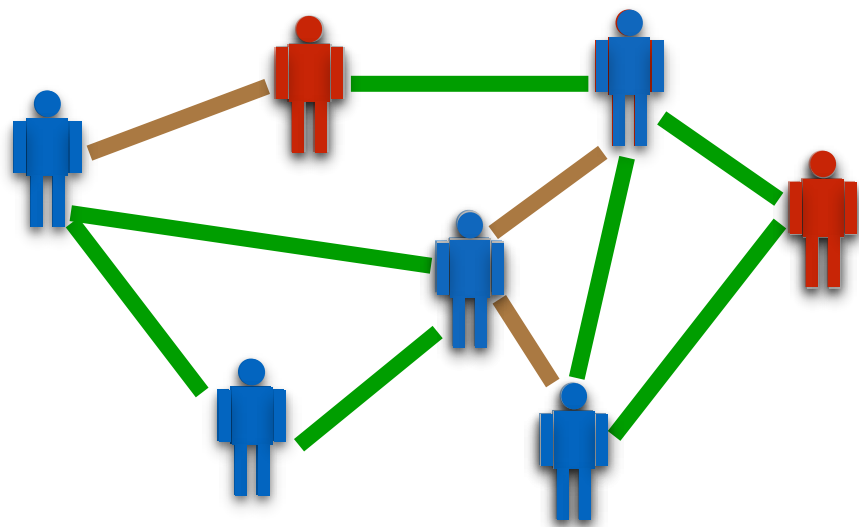
\mathcal{F}

ONLINE NODE CLASSIFICATION

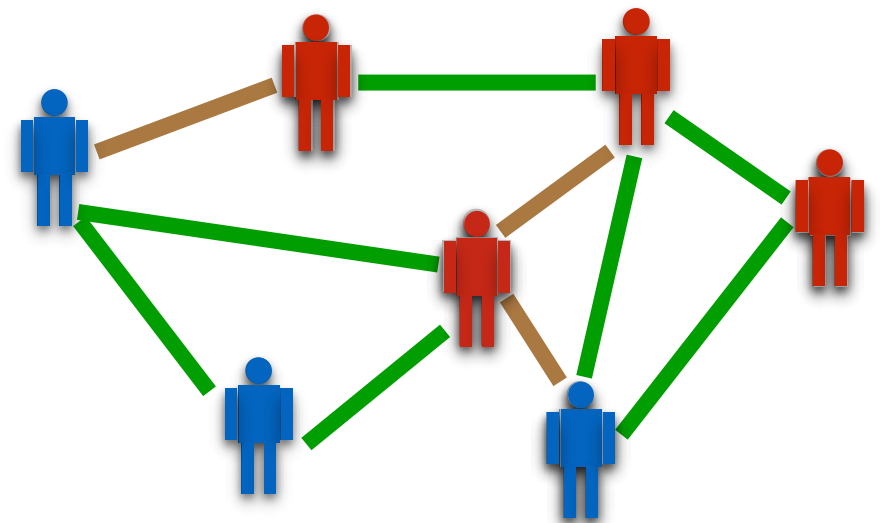


\mathcal{F}

ONLINE NODE CLASSIFICATION



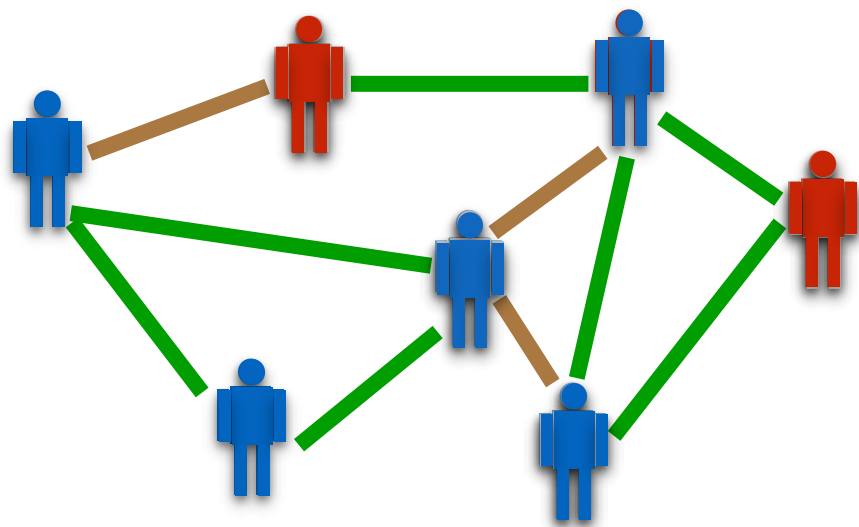
—



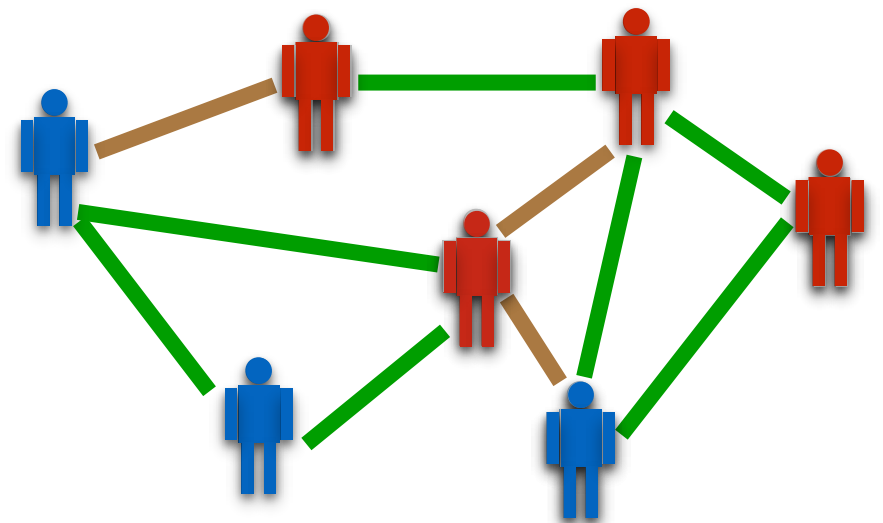
\mathcal{F}

Computationally hard!

ONLINE NODE CLASSIFICATION



—



\mathcal{F}
Relax

Computationally hard!

ONLINE NODE CLASSIFICATION



Regret bound : $\mathbb{E} [\text{Reg}_n] \leq n^{-1} \sqrt{|V| \log |\mathcal{F}|}$