# Empirical Risk Minimization

Vojtěch Franc

March 8, 2022

Prediction task and its solution based on data

Empirical risk minimization

Statistical consistency

Uniform generalization bounds

**XEP33SML – Structured Model Learning, Summer 2022**

## The setting

◆ $\mathcal{X}$ set of input observations

◆ $\mathcal{Y}$ finite set of hidden states, e.g.

- Flat classification: $\mathcal{Y} = \{1, \ldots, K\}$

- Structured classif.: $\mathcal{Y} = \mathcal{Y}_1 \times \cdots \times \mathcal{Y}_{|\mathcal{V}|}$ is a labeling of parts $\mathcal{V}$.

◆ $(x, y) \in \mathcal{X} \times \mathcal{Y}$ randomly drawn from r.v. with p.d.f. $p(x, y)$

◆ $\ell \colon \mathcal{Y} \times \mathcal{Y} \to [0, \infty)$ loss function

## The setting

♦ $\mathcal{X}$ set of input observations

♦ $\mathcal{Y}$ finite set of hidden states, e.g.

- Flat classification: $\mathcal{Y} = \{1, \ldots, K\}$

- Structured classif.: $\mathcal{Y} = \mathcal{Y}_1 \times \cdots \times \mathcal{Y}_{|\mathcal{V}|}$ is a labeling of parts $\mathcal{V}$.

♦ $(x, y) \in \mathcal{X} \times \mathcal{Y}$ randomly drawn from r.v. with p.d.f. $p(x, y)$

♦ $\ell \colon \mathcal{Y} \times \mathcal{Y} \to [0, \infty)$ loss function

**The task:** find a strategy $h \colon \mathcal{X} \to \mathcal{Y}$ with the minimal expected risk

$$R^* = \min_{h \colon \mathcal{X} \to \mathcal{Y}} R(h) \qquad \text{where} \qquad R(h) = \mathbb{E}_{(x,y) \sim p}[\ell(y, h(x))]$$

◆ **Assumption**: we have an access to examples

$$\{(x^1, y^1), (x^2, y^2), \ldots\}$$

drawn from i.i.d. r.v. distributed according to unknown $p(x, y)$.

♦ **Assumption**: we have an access to examples

$$\{(x^1, y^1), (x^2, y^2), \ldots\}$$

drawn from i.i.d. r.v. distributed according to unknown $p(x, y)$.

♦ a) **Evaluation**: Estimate $R(h)$ of a given $h \colon \mathcal{X} \to \mathcal{Y}$ using test set

$$\mathcal{S}^l = \{(x^i, y^i) \in (\mathcal{X} \times \mathcal{Y}) \mid i = 1, \ldots, l\}$$

drawn i.i.d. from $p(x, y)$.

♦ **Assumption**: we have an access to examples

$$\{(x^1, y^1), (x^2, y^2), \ldots\}$$

drawn from i.i.d. r.v. distributed according to unknown $p(x, y)$.

♦ a) **Evaluation**: Estimate $R(h)$ of a given $h\colon \mathcal{X} \to \mathcal{Y}$ using test set

$$\mathcal{S}^l = \{(x^i, y^i) \in (\mathcal{X} \times \mathcal{Y}) \mid i = 1, \ldots, l\}$$

drawn i.i.d. from $p(x, y)$.

♦ b) **Learning**: find $h\colon \mathcal{X} \to \mathcal{Y}$ with small $R(h)$ using training set

$$\mathcal{T}^m = \{(x^i, y^i) \in (\mathcal{X} \times \mathcal{Y}) \mid i = 1, \ldots, m\}$$

drawn i.i.d. from $p(x, y)$.

◆ Given a predictor $h\colon \mathcal{X} \to \mathcal{Y}$, compute the empirical risk

$$R_{\mathcal{S}^l}(h) = \frac{1}{l} \sum_{i=1}^{l} \ell(y^i, h(x^i))$$

and use it as a proxy for $R(h) = \mathbb{E}_{(x,y) \sim p}(\ell(y, h(x)))$.

◆ Given a predictor $h\colon \mathcal{X} \to \mathcal{Y}$, compute the empirical risk

$$R_{\mathcal{S}^l}(h) = \frac{1}{l} \sum_{i=1}^{l} \ell(y^i, h(x^i))$$

and use it as a proxy for $R(h) = \mathbb{E}_{(x,y)\sim p}(\ell(y, h(x)))$.

◆ The value of the empirical risk $R_{\mathcal{S}^l}(h)$ is a random number.

♦ Given a predictor $h \colon \mathcal{X} \to \mathcal{Y}$, compute the empirical risk

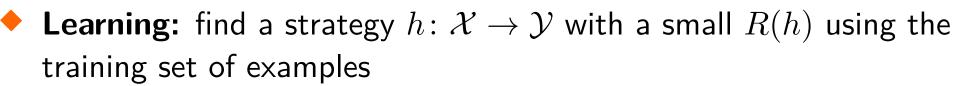$$R_{\mathcal{S}^l}(h) = \frac{1}{l} \sum_{i=1}^{l} \ell(y^i, h(x^i))$$

and use it as a proxy for $R(h) = \mathbb{E}_{(x,y)\sim p}(\ell(y, h(x)))$.

♦ The value of the empirical risk $R_{\mathcal{S}^l}(h)$ is a random number.

♦ Application of Hoeffding inequality: for any $\varepsilon > 0$, the probability of the generalization error being at least $\varepsilon$ can be bound by

$$\mathbb{P}_{\mathcal{S}^l \sim p}\bigg( \underbrace{\Big| R_{\mathcal{S}^l}(h) - R(h) \Big| \geq \varepsilon}_{\text{high generalization error}} \bigg) \leq 2e^{-\frac{2l\,\varepsilon^2}{(\ell_{\min}-\ell_{\max})^2}}$$

◆ **Learning:** find a strategy $h\colon \mathcal{X} \to \mathcal{Y}$ with a small $R(h)$ using the training set of examples

$$\mathcal{T}^m = \{(x^i, y^i) \in (\mathcal{X} \times \mathcal{Y}) \mid i = 1, \ldots, m\}$$

drawn from i.i.d. according to unknown $p(x, y)$.

◆ Use prior knowledge to select hypothesis space

$$\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}} = \{h\colon \mathcal{X} \to \mathcal{Y}\}$$

◆ The learning algorithm

$$A\colon \ \cup_{m=1}^{\infty} (\mathcal{X} \times \mathcal{Y})^m \to \mathcal{H}$$

selects strategy $h_m = A(\mathcal{T}^m)$ based on the training set $\mathcal{T}^m$.

1. Use the training set $\mathcal{T}^m = \{(x^i, y^i) \in \mathcal{X} \times \mathcal{Y} \mid i \in \{1, \ldots, m\}\}$ to approximate $p(x, y)$ by $\hat{p}(x, y)$.

   For example, use the Maximum-Likelihood method:

(a) Guess the shape of the distribution, e.g.

$$\hat{p}_{\boldsymbol{w}}(x, y) = \frac{1}{Z(\boldsymbol{w})} \exp\langle \boldsymbol{w}, \boldsymbol{\phi}(x, y) \rangle, \qquad \boldsymbol{w} \in \mathcal{W}$$

(b) Find the ML estimate

$$\boldsymbol{w}_m \in \underset{\boldsymbol{w} \in \mathcal{W}}{\operatorname{argmax}} \sum_{i=1}^{m} \log \hat{p}_{\boldsymbol{w}}(x^i, y^i)$$

2. Construct a plug-in classifier

$$h_m(x) \in \underset{h \colon \mathcal{X} \to \mathcal{Y}}{\operatorname{argmin}} \mathbb{E}_{(x,y) \sim \hat{p}_{\boldsymbol{w}_m}}[\ell(y, h(x))]$$

◆ Use the training set $\mathcal{T}^m = \{(x^i, y^i) \in \mathcal{X} \times \mathcal{Y} \mid i \in \{1, \ldots, m\}\}$ to approximate the expected risk $R(h)$ by the empirical risk

$$R_{\mathcal{T}^m}(h) = \frac{1}{m} \sum_{i=1}^{m} \ell(y^i, h(x^i))$$

◆ The ERM learning algorithm returns $h_m$ such that

$$h_m \in \underset{h \in \mathcal{H}}{\mathrm{Argmin}}\, R_{\mathcal{T}^m}(h) \tag{1}$$

◆ Depending on the choice of $\mathcal{H}$, $\ell$ and algorithm solving (1) we get individual instances, e.g.: Structured-Output Perceptron, Structured-Output Support Vector Machines, Logistic regression, Neural Networks learned by back-propagation, AdaBoost, . . . .

**The characters of the play:**

- ♦ $R^* = \min_{h \in \mathcal{Y}^{\mathcal{X}}} R(h)$ best attainable (Bayes) risk

- ♦ $R(h_{\mathcal{H}})$ best risk in $\mathcal{H}$; $h_{\mathcal{H}} \in \mathrm{Argmin}_{h \in \mathcal{H}} R(h)$

- ♦ $R(h_m)$ risk of $h_m = A(\mathcal{T}_m)$ learned from $\mathcal{T}^m$

**The characters of the play:**

◆ $R^* = \min_{h \in \mathcal{Y}^{\mathcal{X}}} R(h)$ best attainable (Bayes) risk

◆ $R(h_{\mathcal{H}})$ best risk in $\mathcal{H}$; $h_{\mathcal{H}} \in \mathrm{Argmin}_{h \in \mathcal{H}} R(h)$

◆ $R(h_m)$ risk of $h_m = A(\mathcal{T}_m)$ learned from $\mathcal{T}^m$

**Excess error**: the quantity we want to minimize

$$\underbrace{\left( R(h_m) - R^* \right)}_{\text{excess error}} = \underbrace{\left( R(h_m) - R(h_{\mathcal{H}}) \right)}_{\text{estimation error}} + \underbrace{\left( R(h_{\mathcal{H}}) - R^* \right)}_{\text{approximation error}}$$

**The characters of the play:**

- ◆ $R^* = \min_{h \in \mathcal{Y}^{\mathcal{X}}} R(h)$ best attainable (Bayes) risk

- ◆ $R(h_{\mathcal{H}})$ best risk in $\mathcal{H}$; $h_{\mathcal{H}} \in \operatorname{Argmin}_{h \in \mathcal{H}} R(h)$

- ◆ $R(h_m)$ risk of $h_m = A(\mathcal{T}_m)$ learned from $\mathcal{T}^m$

**Excess error**: the quantity we want to minimize

$$\underbrace{\Big( R(h_m) - R^* \Big)}_{\text{excess error}} = \underbrace{\Big( R(h_m) - R(h_{\mathcal{H}}) \Big)}_{\text{estimation error}} + \underbrace{\Big( R(h_{\mathcal{H}}) - R^* \Big)}_{\text{approximation error}}$$

- ◆ The estimation error is random

- ◆ The estimation error depends on $m$ and $\mathcal{H}$

- ◆ The approximation error depends only on $\mathcal{H}$

◆ The estimation error $R(h_m) - R(h_\mathcal{H})$ is random because it is a function of $h_m = A(\mathcal{T}^m)$ learned on $\mathcal{T}^m$ generated from $p(x, y)$.

◆ The estimation error $R(h_m) - R(h_{\mathcal{H}})$ is random because it is a function of $h_m = A(\mathcal{T}^m)$ learned on $\mathcal{T}^m$ generated from $p(x, y)$.

◆ We can derive bounds on the probability that the estimation error is above $\varepsilon > 0$, that is,

$$\mathbb{P}\left( R(h_m) - R(h_{\mathcal{H}}) \geq \varepsilon \right) \leq U(m, \varepsilon, \mathcal{H})$$

◆ The estimation error $R(h_m) - R(h_{\mathcal{H}})$ is random because it is a function of $h_m = A(\mathcal{T}^m)$ learned on $\mathcal{T}^m$ generated from $p(x, y)$.

◆ We can derive bounds on the probability that the estimation error is above $\varepsilon > 0$, that is,
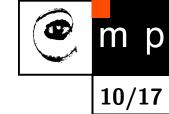
$$\mathbb{P}\left( R(h_m) - R(h_{\mathcal{H}}) \geq \varepsilon \right) \leq U(m, \varepsilon, \mathcal{H})$$

**Definition 1.** *The algorithm $A\colon \cup_{m=1}^{\infty} (\mathcal{X} \times \mathcal{Y})^m \to \mathcal{H}$ is statistically consistent in $\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}}$ if for any $p(x, y)$ it holds that*

$$\forall \varepsilon > 0\colon \lim_{m \to \infty} \mathbb{P}\left( R(h_m) - R(h_{\mathcal{H}}) \geq \varepsilon \right) = 0$$

*where $h_m = A(\mathcal{T}^m)$ is learned from $\mathcal{T}^m$ generated from $p(x, y)$.*

◆ Let $\mathcal{X} = [a, b] \subset \mathbb{R}$, $\mathcal{Y} = \{+1, -1\}$, $\ell(y, y') = [y \neq y']$, $p(x \mid y = +1)$ and $p(x \mid y = -1)$ be uniform distributions on $\mathcal{X}$ and $p(y = +1) = 0.8$.

◆ Let $\mathcal{X} = [a, b] \subset \mathbb{R}$, $\mathcal{Y} = \{+1, -1\}$, $\ell(y, y') = [y \neq y']$, $p(x \mid y = +1)$ and $p(x \mid y = -1)$ be uniform distributions on $\mathcal{X}$ and $p(y = +1) = 0.8$.

◆ The optimal strategy is $h(x) = +1$ with the Bayes risk $R^* = 0.2$.

◆ Let $\mathcal{X} = [a, b] \subset \mathbb{R}$, $\mathcal{Y} = \{+1, -1\}$, $\ell(y, y') = [y \neq y']$, $p(x \mid y = +1)$ and $p(x \mid y = -1)$ be uniform distributions on $\mathcal{X}$ and $p(y = +1) = 0.8$.

◆ The optimal strategy is $h(x) = +1$ with the Bayes risk $R^* = 0.2$.

◆ Consider learning algorithm which for a given training set $\mathcal{T}^m = \{(x^1, y^1), \ldots, (x^m, y^m)\}$ returns strategy

$$h_m(x) = \begin{cases} y^j & \text{if } x = x^j \text{ for some } j \in \{1, \ldots, m\} \\ -1 & \text{otherwise} \end{cases}$$

◆ Let $\mathcal{X} = [a, b] \subset \mathbb{R}$, $\mathcal{Y} = \{+1, -1\}$, $\ell(y, y') = [y \neq y']$, $p(x \mid y = +1)$ and $p(x \mid y = -1)$ be uniform distributions on $\mathcal{X}$ and $p(y = +1) = 0.8$.

◆ The optimal strategy is $h(x) = +1$ with the Bayes risk $R^* = 0.2$.

◆ Consider learning algorithm which for a given training set $\mathcal{T}^m = \{(x^1, y^1), \ldots, (x^m, y^m)\}$ returns strategy

$$h_m(x) = \begin{cases} y^j & \text{if } x = x^j \text{ for some } j \in \{1, \ldots, m\} \\ -1 & \text{otherwise} \end{cases}$$

◆ The empirical risk is $R_{\mathcal{T}^m}(h_m) = 0$ with probability 1 for any $m$.

◆ Let $\mathcal{X} = [a, b] \subset \mathbb{R}$, $\mathcal{Y} = \{+1, -1\}$, $\ell(y, y') = [y \neq y']$, $p(x \mid y = +1)$ and $p(x \mid y = -1)$ be uniform distributions on $\mathcal{X}$ and $p(y = +1) = 0.8$.
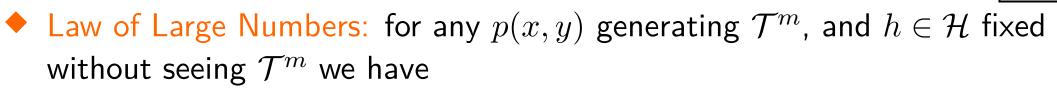
◆ The optimal strategy is $h(x) = +1$ with the Bayes risk $R^* = 0.2$.

◆ Consider learning algorithm which for a given training set $\mathcal{T}^m = \{(x^1, y^1), \ldots, (x^m, y^m)\}$ returns strategy

$$h_m(x) = \begin{cases} y^j & \text{if } x = x^j \text{ for some } j \in \{1, \ldots, m\} \\ -1 & \text{otherwise} \end{cases}$$

◆ The empirical risk is $R_{\mathcal{T}^m}(h_m) = 0$ with probability 1 for any $m$.

◆ The expected risk is $R(h_m) = 0.8$ for any $m$.

◆ Law of Large Numbers: for any $p(x, y)$ generating $\mathcal{T}^m$, and $h \in \mathcal{H}$ fixed without seeing $\mathcal{T}^m$ we have

$$\forall \varepsilon > 0 \colon \quad \lim_{m \to \infty} \mathbb{P}\Big( \underbrace{\big| R(h) - R_{\mathcal{T}^m}(h) \big| \geq \varepsilon}_{\text{high generalization error}} \Big) = 0$$

◆ Law of Large Numbers: for any $p(x,y)$ generating $\mathcal{T}^m$, and $h \in \mathcal{H}$ fixed without seeing $\mathcal{T}^m$ we have

$$\forall \varepsilon > 0: \lim_{m \to \infty} \mathbb{P}\left( \underbrace{\left| R(h) - R_{\mathcal{T}^m}(h) \right| \geq \varepsilon}_{\text{high generalization error}} \right) = 0$$

◆ Uniform Law of Large Numbers: if for any $p(x,y)$ generating $\mathcal{T}^m$ it holds that

$$\forall \varepsilon > 0: \lim_{m \to \infty} \mathbb{P}\Bigg( \begin{array}{l} \left| R(h_1) - R_{\mathcal{T}^m}(h_1) \right| \geq \varepsilon \quad \text{or} \\[1mm] \left| R(h_2) - R_{\mathcal{T}^m}(h_2) \right| \geq \varepsilon \quad \text{or} \\[1mm] \vdots \\[1mm] \underbrace{\left| R(h_{|\mathcal{H}|}) - R_{\mathcal{T}^m}(h_{|\mathcal{H}|}) \right| \geq \varepsilon}_{\substack{\text{high generalization error at least} \\ \text{for one strategy}}} \end{array} \Bigg) = 0$$

we say that ULLN applies for $\mathcal{H}$.

◆ Law of Large Numbers: for any $p(x, y)$ generating $\mathcal{T}^m$, and $h \in \mathcal{H}$ fixed without seeing $\mathcal{T}^m$ we have

$$\forall \varepsilon > 0: \quad \lim_{m \to \infty} \mathbb{P}\Big( \underbrace{\big| R(h) - R_{\mathcal{T}^m}(h) \big| \geq \varepsilon}_{\text{high generalization error}} \Big) = 0$$

◆ Uniform Law of Large Numbers: if for any $p(x, y)$ generating $\mathcal{T}^m$ it holds that

$$\forall \varepsilon > 0: \quad \lim_{m \to \infty} \mathbb{P}\Big( \underbrace{\sup_{h \in \mathcal{H}} \big| R(h) - R_{\mathcal{T}^m}(h) \big| \geq \varepsilon}_{\substack{\text{high generalization error at least} \\ \text{for one strategy}}} \Big) = 0$$

we say that ULLN applies for $\mathcal{H}$.

◆ Law of Large Numbers: for any $p(x, y)$ generating $\mathcal{T}^m$, and $h \in \mathcal{H}$ fixed without seeing $\mathcal{T}^m$ we have

$$\forall \varepsilon > 0: \lim_{m \to \infty} \mathbb{P}\Big( \underbrace{\big| R(h) - R_{\mathcal{T}^m}(h) \big| \geq \varepsilon}_{\text{high generalization error}} \Big) = 0$$
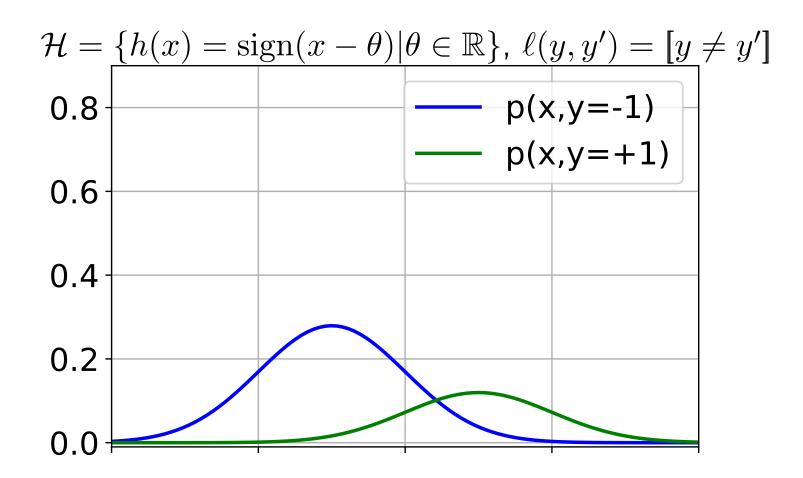
◆ Uniform Law of Large Numbers: if for any $p(x, y)$ generating $\mathcal{T}^m$ it holds that

$$\forall \varepsilon > 0: \lim_{m \to \infty} \mathbb{P}\Big( \underbrace{\sup_{h \in \mathcal{H}} \big| R(h) - R_{\mathcal{T}^m}(h) \big| \geq \varepsilon}_{\substack{\text{high generalization error at least} \\ \text{for one strategy}}} \Big) = 0$$
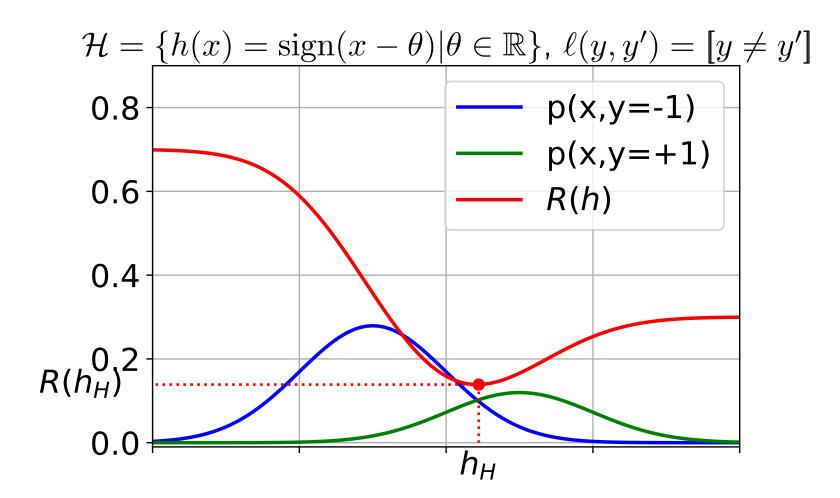
we say that ULLN applies for $\mathcal{H}$.

**Theorem 1.** *If ULLN applies for $\mathcal{H}$ then ERM is statistically consistent in $\mathcal{H}$.*
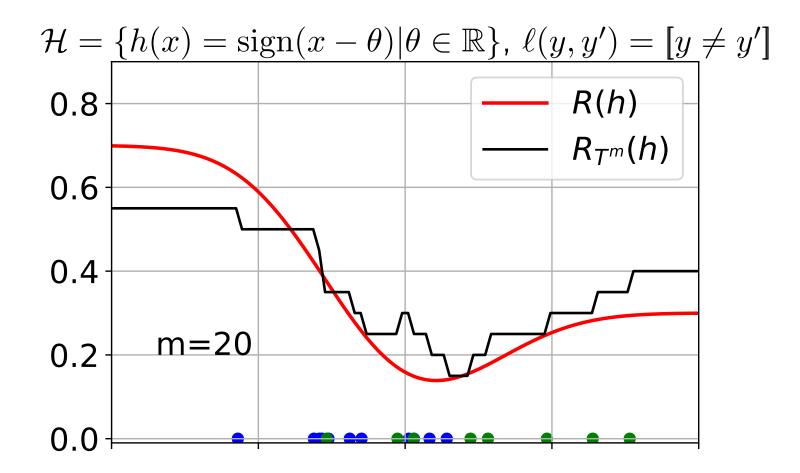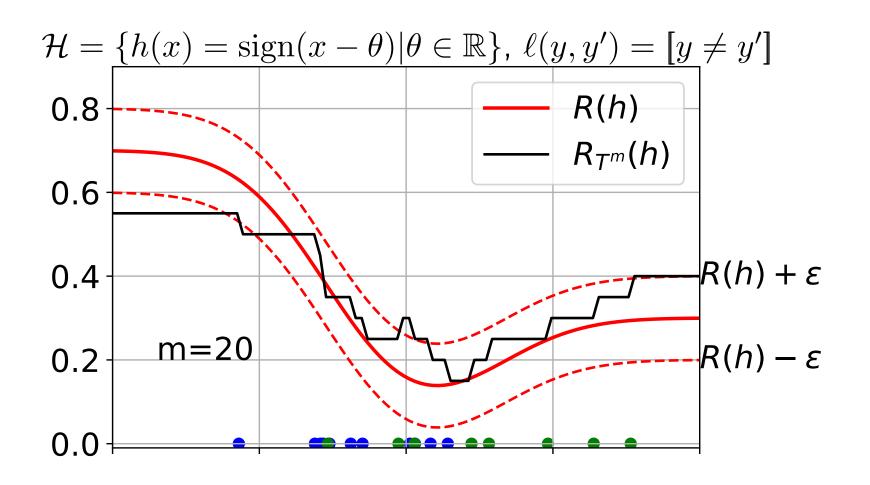
$$\mathcal{H} = \{h(x) = \mathrm{sign}(x - \theta) | \theta \in \mathbb{R}\}, \ \ell(y, y') = [y \neq y']$$

$$\mathcal{H} = \{h(x) = \mathrm{sign}(x - \theta) | \theta \in \mathbb{R}\}, \ \ell(y, y') = [y \neq y']$$

$$\mathbb{P}\left( \underbrace{\sup_{h \in \mathcal{H}} \left| R(h) - R_{\mathcal{T}^m}(h) \right|}_{\text{worst generalization error}} \geq \varepsilon \right) \leq B(m, \mathcal{H}, \varepsilon)$$



$\mathcal{H} = \{ h(x) = \text{sign}(x - \theta) | \theta \in \mathbb{R} \},\ \ell(y, y') = [y \neq y']$

$$\mathbb{P}\Big( \underbrace{\sup_{h \in \mathcal{H}} \big| R(h) - R_{\mathcal{T}^m}(h) \big|}_{\text{worst generalization error}} \geq \varepsilon \Big) \leq B(m, \mathcal{H}, \varepsilon)$$

$$\mathcal{H} = \{h(x) = \mathrm{sign}(x - \theta) | \theta \in \mathbb{R}\},\ \ell(y, y') = [y \neq y']$$

$$\mathbb{P}\bigg(\underbrace{\sup_{h\in\mathcal{H}}\big|R(h)-R_{\mathcal{T}^m}(h)\big|}_{\text{worst generalization error}}\geq\varepsilon\bigg)\leq B(m,\mathcal{H},\varepsilon)$$



$\mathcal{H}=\{h(x)=\operatorname{sign}(x-\theta)\,|\,\theta\in\mathbb{R}\},\ \ell(y,y')=[y\neq y']$

$$\mathbb{P}\left(\underbrace{\sup_{h\in\mathcal{H}}\big|R(h)-R_{\mathcal{T}^m}(h)\big|}_{\text{worst generalization error}}\geq\varepsilon\right)\leq B(m,\mathcal{H},\varepsilon)$$

$$\mathcal{H}=\{h(x)=\operatorname{sign}(x-\theta)|\theta\in\mathbb{R}\},\ \ell(y,y')=[y\neq y']$$

$$\mathbb{P}\Big(\underbrace{\sup_{h\in\mathcal{H}}\big|R(h)-R_{\mathcal{T}^m}(h)\big|}_{\text{worst generalization error}} \geq \varepsilon\Big) \leq B(m,\mathcal{H},\varepsilon)$$



$\mathcal{H}=\{h(x)=\mathrm{sign}(x-\theta)|\theta\in\mathbb{R}\},\ \ell(y,y')=[y\neq y']$

$$\mathbb{P}\left( \underbrace{\sup_{h \in \mathcal{H}} \left| R(h) - R_{\mathcal{T}^m}(h) \right|}_{\text{worst generalization error}} \geq \varepsilon \right) \leq B(m, \mathcal{H}, \varepsilon)$$

$$\underbrace{R(h_m) - R(h_{\mathcal{H}})}_{\text{estimation error}}$$

$$\mathcal{H} = \{h(x) = \text{sign}(x - \theta) | \theta \in \mathbb{R}\}, \ \ell(y, y') = [y \neq y']$$

$$\mathbb{P}\left(\underbrace{\sup_{h\in\mathcal{H}}\left|R(h)-R_{\mathcal{T}^m}(h)\right|}_{\text{worst generalization error}} \geq \varepsilon\right) \leq B(m,\mathcal{H},\varepsilon)$$

$$\underbrace{R(h_m)-R(h_\mathcal{H})}_{\text{estimation error}} \leq 2\underbrace{\sup_{h\in\mathcal{H}}\left|R(h)-R_{\mathcal{T}^m}(h)\right|}_{\text{worst generalizaton error}}$$

$$\mathcal{H} = \{h(x) = \operatorname{sign}(x-\theta)|\theta\in\mathbb{R}\},\ \ell(y,y') = [y\neq y']$$

$$\mathbb{P}\left( \underbrace{\sup_{h \in \mathcal{H}} \left| R(h) - R_{\mathcal{T}^m}(h) \right|}_{\text{worst generalization error}} \geq \varepsilon \right) \leq B(m, \mathcal{H}, \varepsilon)$$

$$\mathbb{P}\left( \underbrace{R(h_m) - R(h_{\mathcal{H}})}_{\text{estimation error}} \geq \varepsilon \right) \leq \mathbb{P}\left( \underbrace{\sup_{h \in \mathcal{H}} \left| R(h) - R_{\mathcal{T}^m}(h) \right|}_{\text{worst generalizaton error}} \geq \frac{\varepsilon}{2} \right)$$

$$\mathcal{H} = \{ h(x) = \text{sign}(x - \theta) | \theta \in \mathbb{R} \}, \ \ell(y, y') = [y \neq y']$$

For fixed $\mathcal{T}^m$ and $h_m \in \operatorname{Argmin}_{h \in \mathcal{H}} R_{\mathcal{T}^m}(h)$ we have:

$$R(h_m) - R(h_{\mathcal{H}}) = \left( R(h_m) - R_{\mathcal{T}^m}(h_m) \right) + \left( R_{\mathcal{T}^m}(h_m) - R(h_{\mathcal{H}}) \right)$$

$$\leq \left( R(h_m) - R_{\mathcal{T}^m}(h_m) \right) + \left( R_{\mathcal{T}^m}(h_{\mathcal{H}}) - R(h_{\mathcal{H}}) \right)$$

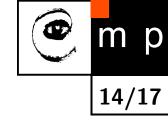$$\leq 2 \sup_{h \in \mathcal{H}} \left| R(h) - R_{\mathcal{T}^m}(h) \right|$$

Therefore $\varepsilon \leq R(h_m) - R(h_{\mathcal{H}})$ implies $\frac{\varepsilon}{2} \leq \sup_{h \in \mathcal{H}} \left| R(h) - R_{\mathcal{T}^m}(h) \right|$ and

$$\mathbb{P} \left( R(h_m) - R(h_{\mathcal{H}}) \geq \varepsilon \right) \leq \mathbb{P} \left( \sup_{h \in \mathcal{H}} \left| R(h) - R_{\mathcal{T}^m}(h) \right| \geq \frac{\varepsilon}{2} \right)$$

1. $\mathcal{H}$ is a finite set and $\ell \colon \mathcal{Y} \times \mathcal{Y} \to [\ell_{min}, \ell_{max}]$. Then,

$$
\mathbb{P}_{\mathcal{T} \sim p^m} \left( \max_{h \in \mathcal{H}} \left| R(h) - R_{\mathcal{T}^m}(h) \right| \geq \varepsilon \right) \leq 2 |\mathcal{H}| \exp \left( \frac{-2m\, \varepsilon^2}{(\ell_{max} - \ell_{min})^2} \right)
$$

holds for any $\varepsilon > 0$ and $m \in \mathcal{N}$.

1. $\mathcal{H}$ is a finite set and $\ell \colon \mathcal{Y} \times \mathcal{Y} \to [\ell_{min}, \ell_{max}]$. Then,

$$\mathbb{P}_{\mathcal{T} \sim p^m} \left( \max_{h \in \mathcal{H}} \left| R(h) - R_{\mathcal{T}^m}(h) \right| \geq \varepsilon \right) \leq 2|\mathcal{H}| \exp \left( \frac{-2m\,\varepsilon^2}{(\ell_{max} - \ell_{min})^2} \right)$$

holds for any $\varepsilon > 0$ and $m \in \mathcal{N}$.

2. $\ell(y, y') = [y \neq y']$, $\mathcal{Y} = \{+1, -1\}$ and VC-dimension of $\mathcal{H}$ is finite. VC-dimension $d$ of $\mathcal{H}$ is the maximal number of inputs which can be classified by strategies from $\mathcal{H}$ in all possible (that is $2^d$) ways. Then,

$$\mathbb{P}_{\mathcal{T} \sim p^m} \left( \sup_{h \in \mathcal{H}} \left| R(h) - R_{\mathcal{T}^m}(h) \right| \geq \varepsilon \right) \leq 4 \left( \frac{2\,e\,m}{d} \right)^d e^{-\frac{m\,\varepsilon^2}{8}}$$

♦ Let $z = (x, y) \in \mathcal{Z} = \mathcal{X} \times \mathcal{Y}$, $p(z) = p(x, y)$ and $g(z) = \ell(y, h(x))$.

◆ Let $z = (x, y) \in \mathcal{Z} = \mathcal{X} \times \mathcal{Y}$, $p(z) = p(x, y)$ and $g(z) = \ell(y, h(x))$.

**Definition 2.** *Let $\mathcal{G} \subseteq [a, b]^{\mathcal{Z}}$ be a set of functions $g \colon \mathcal{Z} \to [a, b]$ where $a, b \in \mathbb{R}$ and $a < b$. Let $\mathcal{U}^m = \{z^1, \dots, z^m\} \in \mathcal{Z}^m$ be drawn i.i.d. from $p(z)$.*

◆ Let $z = (x, y) \in \mathcal{Z} = \mathcal{X} \times \mathcal{Y}$, $p(z) = p(x, y)$ and $g(z) = \ell(y, h(x))$.

**Definition 2.** *Let $\mathcal{G} \subseteq [a, b]^{\mathcal{Z}}$ be a set of functions $g \colon \mathcal{Z} \to [a, b]$ where $a, b \in \mathbb{R}$ and $a < b$. Let $\mathcal{U}^m = \{z^1, \ldots, z^m\} \in \mathcal{Z}^m$ be drawn i.i.d. from $p(z)$.*

*The empirical Rademacher complexity of $\mathcal{G}$ w.r.t. to the sample $\mathcal{U}^m$ is*

$$\hat{\mathfrak{R}}_m(\mathcal{G}, \mathcal{U}^m) = \mathbb{E}_{\sigma \sim \mathrm{Unif}\{-1, +1\}} \left[ \sup_{g \in \mathcal{G}} \frac{1}{m} \sum_{i=1}^{m} \sigma_i \, g(z_i) \right]$$

◆ Let $z = (x, y) \in \mathcal{Z} = \mathcal{X} \times \mathcal{Y}$, $p(z) = p(x, y)$ and $g(z) = \ell(y, h(x))$.

**Definition 2.** *Let $\mathcal{G} \subseteq [a, b]^{\mathcal{Z}}$ be a set of functions $g \colon \mathcal{Z} \to [a, b]$ where $a, b \in \mathbb{R}$ and $a < b$. Let $\mathcal{U}^m = \{z^1, \ldots, z^m\} \in \mathcal{Z}^m$ be drawn i.i.d. from $p(z)$.*

*The empirical Rademacher complexity of $\mathcal{G}$ w.r.t. to the sample $\mathcal{U}^m$ is*

$$\hat{\mathfrak{R}}_m(\mathcal{G}, \mathcal{U}^m) = \mathbb{E}_{\sigma \sim \mathrm{Unif}\{-1, +1\}} \left[ \sup_{g \in \mathcal{G}} \frac{1}{m} \sum_{i=1}^{m} \sigma_i \, g(z_i) \right]$$

*The Rademacher complexity of $\mathcal{G}$ w.r.t. distribution $p(z)$ is*

$$\mathfrak{R}_m(\mathcal{G}) = \mathbb{E}_{\mathcal{U}^m \sim p^m(z)} \left[ \hat{\mathfrak{R}}_m(\mathcal{G}, \mathcal{U}^m) \right]$$

◆ Let $\mathcal{G} \subseteq [a,b]^{\mathcal{Z}}$ be a set of functions. Then, for every $\delta \in (0,1)$

$$\sup_{g \in \mathcal{G}} \left| \mathbb{E}_{z \sim p}(g(z)) - \frac{1}{m} \sum_{i=1}^{m} g(z_i) \right| \leq 2\,\mathfrak{R}_m(\mathcal{G}) + (b-a)\sqrt{\frac{\log 2/\delta}{2\,m}}$$

holds with probability $1 - \delta$ at least, w.r.t. $\mathcal{U}^m = \{z^1, \ldots, z^m\} \sim p^m(z)$.

◆ Let $\mathcal{G} \subseteq [a,b]^{\mathcal{Z}}$ be a set of functions. Then, for every $\delta \in (0,1)$

$$\sup_{g \in \mathcal{G}} \left| \mathbb{E}_{z \sim p}(g(z)) - \frac{1}{m} \sum_{i=1}^{m} g(z_i) \right| \leq 2\,\mathfrak{R}_m(\mathcal{G}) + (b-a)\sqrt{\frac{\log 2/\delta}{2\,m}}$$

holds with probability $1 - \delta$ at least, w.r.t. $\mathcal{U}^m = \{z^1, \ldots, z^m\} \sim p^m(z)$.

◆ For every $\delta \in (0,1)$

$$\sup_{g \in \mathcal{G}} \left| \mathbb{E}_{z \sim p}(g(z)) - \frac{1}{m} \sum_{i=1}^{m} g(z_i) \right| \leq 3\,\hat{\mathfrak{R}}_m(\mathcal{G}, \mathcal{U}^m) + (b-a)\sqrt{\frac{\log 4/\delta}{2\,m}}$$

holds with probability $1 - \delta$ at least, w.r.t. $\mathcal{U}^m = \{z^1, \ldots, z^m\} \sim p^m(z)$.

◆ Assume that $\mathcal{X} \subseteq \mathbb{R}^n$ and $p(\boldsymbol{x}, y)$ is such that $\|\boldsymbol{x}\| \leq R$.

◆ Assume that

$$\mathcal{G} = \left\{ \psi(\langle \boldsymbol{w}, \boldsymbol{x} \rangle, y) \mid \|\boldsymbol{w}\|_2 \leq B \right\}$$

where $\psi \colon \mathbb{R} \times \mathcal{Y} \to \mathbb{R}$ is such that $f(t) = \psi(t, y)$ is $\rho$-Lipschitz continuous for all $y \in \mathcal{Y}$.

♦ Assume that $\mathcal{X} \subseteq \mathbb{R}^n$ and $p(\boldsymbol{x}, y)$ is such that $\|\boldsymbol{x}\| \leq R$.

♦ Assume that

$$\mathcal{G} = \left\{ \psi(\langle \boldsymbol{w}, \boldsymbol{x} \rangle, y) \mid \|\boldsymbol{w}\|_2 \leq B \right\}$$

where $\psi \colon \mathbb{R} \times \mathcal{Y} \to \mathbb{R}$ is such that $f(t) = \psi(t, y)$ is $\rho$-Lipschitz continuous for all $y \in \mathcal{Y}$.

E.g. $\psi(t, y) = \max\{0, 1 - t\,y\}$ and $\psi(t) = |t - y|$ are 1-Lipschitz.

◆ Assume that $\mathcal{X} \subseteq \mathbb{R}^n$ and $p(\boldsymbol{x}, y)$ is such that $\|\boldsymbol{x}\| \leq R$.

◆ Assume that

$$\mathcal{G} = \left\{ \psi(\langle \boldsymbol{w}, \boldsymbol{x} \rangle, y) \mid \|\boldsymbol{w}\|_2 \leq B \right\}$$

where $\psi \colon \mathbb{R} \times \mathcal{Y} \to \mathbb{R}$ is such that $f(t) = \psi(t, y)$ is $\rho$-Lipschitz continuous for all $y \in \mathcal{Y}$.

E.g. $\psi(t, y) = \max\{0, 1 - t\,y\}$ and $\psi(t) = |t - y|$ are 1-Lipschitz.

◆ Then,

$$\hat{\mathfrak{R}}_m(\mathcal{G}) \leq \frac{\rho\,B\,R}{\sqrt{m}}$$

◆ We can also compute

$$b = \max_{t \in [-BR, BR]} \psi(t, y) \qquad \text{and} \qquad a = \min_{t \in [-BR, BR]} \psi(t, y)$$