# Statistical Machine Learning (BE4M33SSU)
# Lecture 4: Empirical Risk Minimization II

Czech Technical University in Prague

V. Franc

**BE4M33SSU – Statistical Machine Learning, Winter 2021**

◆ $\mathcal{X}$ is a set of observations and $\mathcal{Y} = \{+1, -1\}$ a set of hidden labels

◆ $\phi \colon \mathcal{X} \to \mathbb{R}^n$ is fixed feature map embedding $\mathcal{X}$ to $\mathbb{R}^n$

◆ Task: find linear classification strategy $h \colon \mathcal{X} \to \mathcal{Y}$

$$h(x; \boldsymbol{w}, b) = \text{sign}(\langle \boldsymbol{w}, \boldsymbol{\phi}(x) \rangle + b) = \begin{cases} +1 & \text{if} \quad \langle \boldsymbol{w}, \boldsymbol{\phi}(x) \rangle + b \geq 0 \\ -1 & \text{if} \quad \langle \boldsymbol{w}, \boldsymbol{\phi}(x) \rangle + b < 0 \end{cases}$$

with minimal expected risk

$$R^{0/1}(h) = \mathbb{E}_{(x,y) \sim p}\Big( \ell^{0/1}(y, h(x)) \Big) \quad \text{where} \quad \ell^{0/1}(y, y') = [y \neq y']$$

◆ We are given a set of training examples

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

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

◆ The Empirical Risk Minimization principle leads to solving

$$(\boldsymbol{w}^*, b^*) \in \underset{(\boldsymbol{w}, b) \in (\mathbb{R}^n \times \mathbb{R})}{\operatorname{Argmin}} R_{\mathcal{T}^m}^{0/1}(h(\cdot; \boldsymbol{w}, b)) \tag{1}$$

where the empirical risk is

$$R_{\mathcal{T}^m}^{0/1}(h(\cdot; \boldsymbol{w}, b)) = \frac{1}{m} \sum_{i=1}^{m} [y^i \neq h(x^i; \boldsymbol{w}, b)]$$

◆ Algorithmic issues (next lecture): in general, there is no known algorithm solving the task (1) in time polynomial in $m$.

◆ The uniform bound on the generalization error (this lecture):

$$\mathbb{P}\left( \sup_{h \in \mathcal{H}} \left| R^{0/1}(h) - R_{\mathcal{T}^m}^{0/1}(h) \right| \geq \varepsilon \right) \leq B(m, \mathcal{H}, \varepsilon)$$

**Definition:** Let $\mathcal{H} \subseteq \{-1,+1\}^{\mathcal{X}}$ and $\{x^1,\ldots,x^m\} \in \mathcal{X}^m$ be a set of $m$ input observations. The set $\{x^1,\ldots,x^m\}$ is said to be shattered by $\mathcal{H}$ if for all $\boldsymbol{y} \in \{+1,-1\}^m$ there exists $h \in \mathcal{H}$ such that $h(x^i) = y^i$, $i \in \{1,\ldots,m\}$.
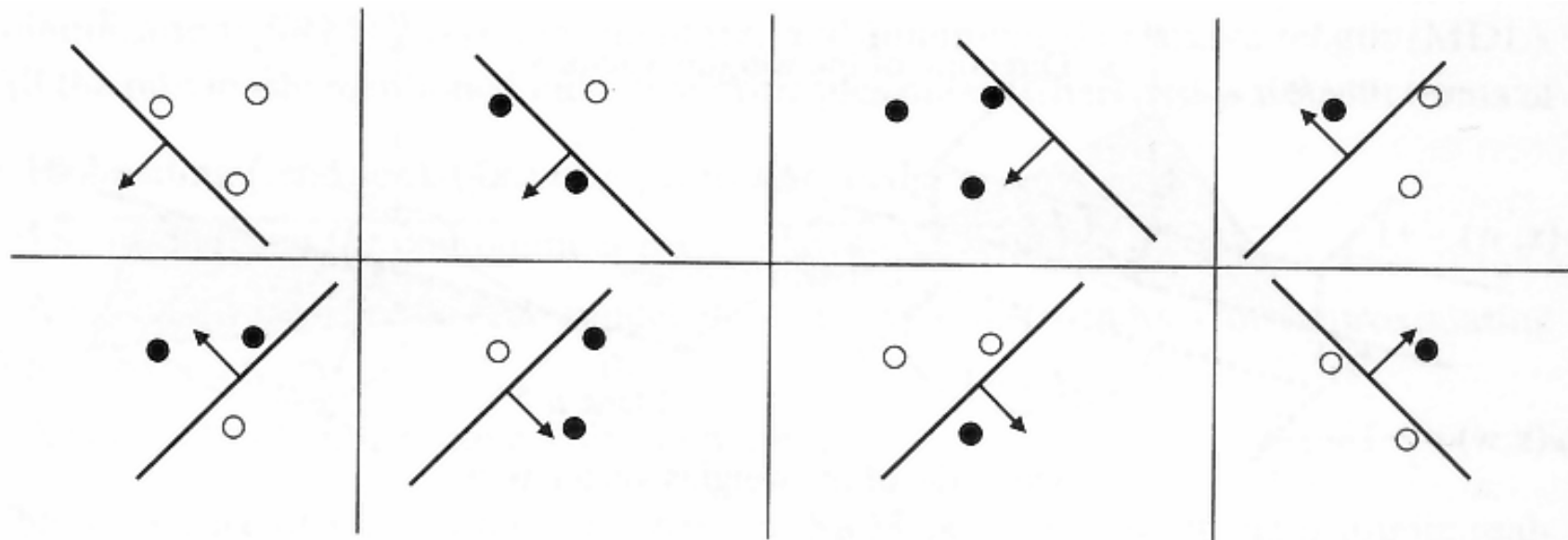
**Definition:** Let $\mathcal{H} \subseteq \{-1,+1\}^{\mathcal{X}}$. The Vapnik-Chervonenkis dimension of $\mathcal{H}$ is the cardinality of the largest set of points from $\mathcal{X}$ which can be shattered by $\mathcal{H}$.

**Theorem:** The VC-dimension of the hypothesis class of all two-class linear classifiers operating in $n$-dimensional feature space
$\mathcal{H} = \{h(x; \boldsymbol{w}, b) = \text{sign}(\langle \boldsymbol{w}, \boldsymbol{\phi}(x) \rangle + b) \mid (\boldsymbol{w}, b) \in (\mathbb{R}^n \times \mathbb{R})\}$ is $n + 1$.

Example for $n = 2$-dimensional feature class

**Theorem:** Let $\mathcal{H} \subseteq \{+1, -1\}^{\mathcal{X}}$ be a hypothesis class with VC dimension $d < \infty$ and $\mathcal{T}^m = \{(x^1, y^1), \ldots, (x^m, y^m)\} \in (\mathcal{X} \times \mathcal{Y})^m$ a training set draw from i.i.d. rand vars with distribution $p(x, y)$. Then, for any $\varepsilon > 0$ it holds

$$\mathbb{P}\left( \sup_{h \in \mathcal{H}} \left| R^{0/1}(h) - R^{0/1}_{\mathcal{T}^m}(h) \right| \geq \varepsilon \right) \leq 4 \left( \frac{2\,e\,m}{d} \right)^d e^{-\frac{m\,\varepsilon^2}{8}}$$

**Corollary:** Let $\mathcal{H} \subseteq \{+1, -1\}^{\mathcal{X}}$ be a hypothesis class with VC dimension $d < \infty$. Then ULLN applies.

◆ We learned how to bound the generalization error uniformly for:

- Finite hypothesis class $\mathcal{H} = \{h_1, \ldots, h_K\}$:

$$\mathbb{P}\left( \max_{h \in \mathcal{H}} |R_{\mathcal{T}^m}(h) - R(h)| \geq \varepsilon \right) \leq 2|\mathcal{H}|e^{-\frac{2m\,\varepsilon^2}{(b-a)^2}}$$

- Two-class classifiers $\mathcal{H} \subseteq \{+1, -1\}^{\mathcal{X}}$ a finite VC-dimensions $d$:

$$\mathbb{P}\left( \sup_{h \in \mathcal{H}} \left| R^{0/1}(h) - R_{\mathcal{T}^m}^{0/1}(h) \right| \geq \varepsilon \right) \leq 4\left( \frac{2\,e\,m}{d} \right)^d e^{-\frac{m\,\varepsilon^2}{8}}$$

In both cases the bound goes to zero, i.e., ULLN applies.

◆ Does ERM algorithm $h_m \in \underset{h \in \mathcal{H}}{\mathrm{Argmin}}\, R_{\mathcal{T}^m}(h)$ finds strategy with the minimal risk $R(h)$?

# Statistically consistent learning algorithm

◆ $h_{\mathcal{H}} \in \text{Argmin}_{h \in \mathcal{H}} R(h)$ the best strategy in $\mathcal{H}$ has the risk $R(h_{\mathcal{H}})$

◆ $h_m = A(\mathcal{T}_m)$ strategy learned from $\mathcal{T}_m$ with has risk $R(h_m)$

◆ $R(h_m) - R(h_{\mathcal{H}})$ is the estimation error

◆ The statistically consistent algorithm can make the estimation error arbitrarily small if it has enough examples.

**Definition:** 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

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

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

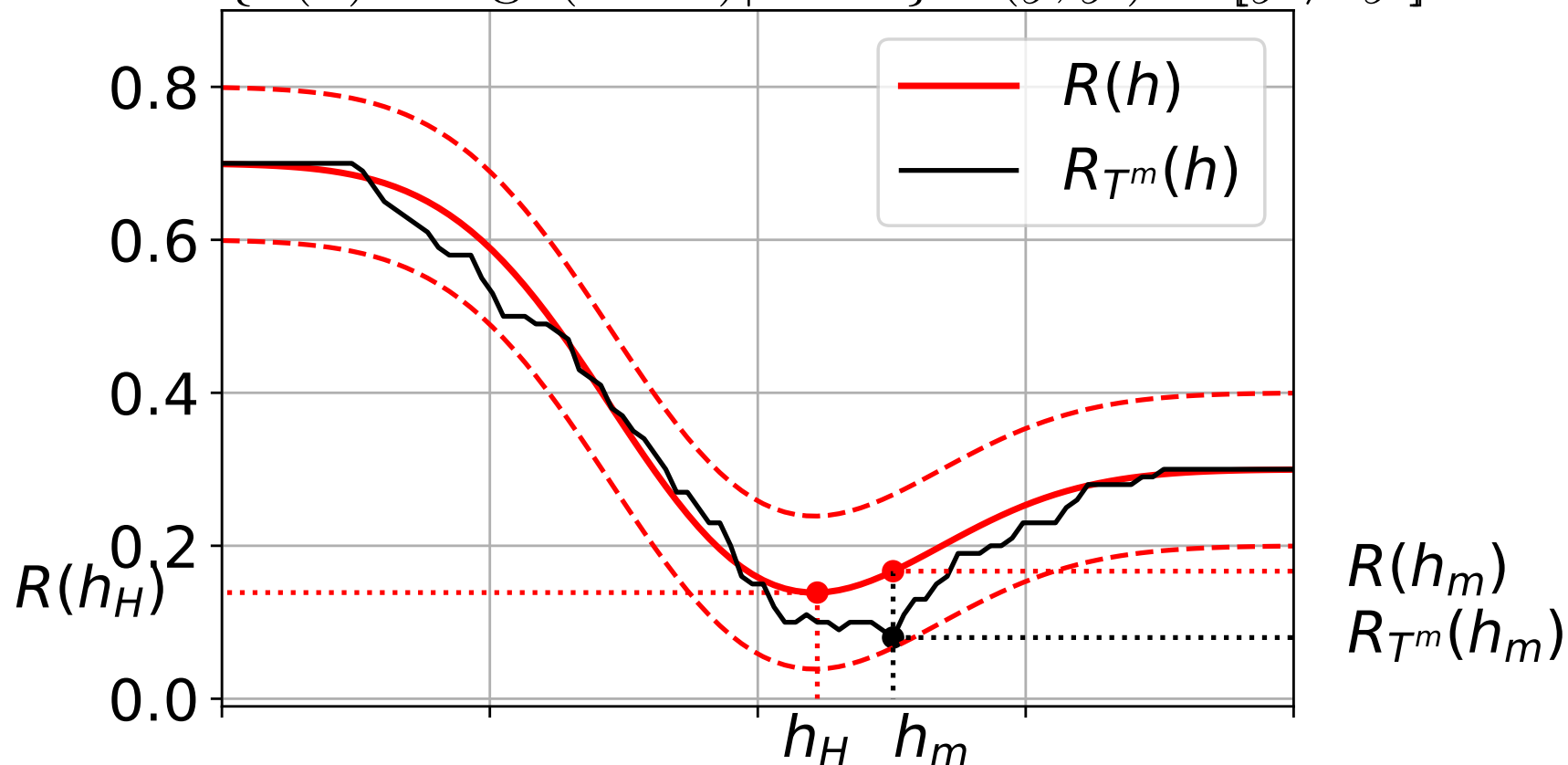$$\mathbb{P}\left(\underbrace{\sup_{h \in \mathcal{H}} \left| R(h) - R_{\mathcal{T}^m}(h) \right|}_{\text{highest generalization error}} \geq \varepsilon \right) \leq B(m, \mathcal{H}, \varepsilon)$$

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

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

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

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

so if converges the RHS to zero (ULLN) so does the LHS (estimation error).

◆ Let $\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}}$ be the hypothesis class with the best predictor

$$h_{\mathcal{H}} \in \underset{h \in \mathcal{H}}{\mathrm{Argmin}} \, R(h)$$

◆ We learn $h_m$ from $\mathcal{T}^m \sim p(x, y)$ with the ERM algorithm

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

◆ Assume we have for $\mathcal{H}$ the uniform bound on the generalization error

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

◆ Then, for any $\varepsilon > 0$ the inequality

$$R(h_m) \leq R(h_{\mathcal{H}}) + \varepsilon$$

holds with the probability $1 - B(m, \mathcal{H}, \varepsilon/2)$ at least.

**The characters of the play:**

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

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

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

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

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

Questions:

◆ What causes individual errors ?

◆ How do the errors depend on $\mathcal{H}$ and $m$?