

# Statistical Machine Learning (BE4M33SSU)

## Lecture 3: Empirical Risk Minimization II

Czech Technical University in Prague

## Linear classifier with minimal classification error

- ◆  $\mathcal{X}$  is a set of observations and  $\mathcal{Y} = \{+1, -1\}$  is a set of hidden labels
- ◆  $\phi: \mathcal{X} \rightarrow \mathbb{R}^n$  is fixed feature map embedding observations from  $\mathcal{X}$  to  $\mathbb{R}^n$
- ◆ Task: we search for a linear classification strategy  $h: \mathcal{X} \rightarrow \mathcal{Y}$

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

with minimal expected risk

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

- ◆ We are given a set of training examples

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

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

## ERM learning for linear classifiers

- ◆ The Empirical Risk Minimization principle leads to solving

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

where the empirical risk is

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

In this lecture we address the following issues:

1. Algorithmic issues: In the general case there is no known algorithm solving the task (1) in time polynomial in  $m$ .
2. Is the ERM algorithm for hypothesis space containing linear classifiers statistically consistent? ... yes.

## Vapnik-Chervonenkis (VC) dimension

**Definition 1.** Let  $\mathcal{H} \subseteq \{-1, +1\}^{\mathcal{X}}$  and  $\{x^1, \dots, x^m\} \in \mathcal{X}^m$  be a set of  $m$  input observations. The set  $\{x^1, \dots, x^m\}$  is said to be shattered by  $\mathcal{H}$  if for all  $\mathbf{y} \in \{+1, -1\}^m$  there exists  $h \in \mathcal{H}$  such that  $h(x^i) = y^i$ ,  $i \in \{1, \dots, m\}$ .

**Definition 2.** 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 1.** The VC-dimension of the hypothesis space of all linear classifiers operating in  $n$ -dimensional feature space  $\mathcal{H} = \{h(x; \mathbf{w}, b) = \text{sign}(\langle \mathbf{w}, \phi(x) \rangle + b) \mid (\mathbf{w}, b) \in (\mathbb{R}^n \times \mathbb{R})\}$  is  $n + 1$ .

**Theorem 2.** Let  $\mathcal{H} \subseteq \{+1, -1\}^{\mathcal{X}}$  be a hypothesis space with VC dimension  $d < \infty$  and  $\mathcal{T}^m = \{(x^1, y^1), \dots, (x^m, y^m)\} \in (\mathcal{X} \times \mathcal{Y})^m$  a training set drawn from i.i.d. random variables with distribution  $p(x, y)$ . Then, for any  $\varepsilon > 0$  it holds

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

**Corollary 1.** Let  $\mathcal{H} \subseteq \{+1, -1\}^{\mathcal{X}}$  be a hypothesis space with VC dimension  $d < \infty$ . Then ERM is statistically consistent in  $\mathcal{H}$  w.r.t.  $\ell^{0/1}$  loss function.

**Corollary 2.** Let  $\mathcal{H} \subseteq \{+1, -1\}^{\mathcal{X}}$  be a hypothesis space with VC dimension  $d < \infty$ . Then, for any  $0 < \delta < 1$  the inequality

$$R^{0/1}(h) \leq R_{\mathcal{T}^m}^{0/1}(h) + \sqrt{\frac{8\left(d \log \frac{2em}{d} + \log \frac{4}{\delta}\right)}{m}}$$

holds for any  $h \in \mathcal{H}$  with probability  $1 - \delta$  at least.

## Training linear classifier from separable examples

**Definition 3.** The examples  $\mathcal{T}^m = \{(x^i, y^i) \in (\mathcal{X} \times \mathcal{Y}) \mid i = 1, \dots, m\}$  are linearly separable w.r.t. feature map  $\phi: \mathcal{X} \rightarrow \mathbb{R}^n$  if there exists  $(\mathbf{w}, b) \in \mathbb{R}^{n+1}$  such that

$$y^i(\langle \mathbf{w}, \phi(x^i) \rangle + b) > 0, \quad i \in \{1, \dots, m\} \quad (2)$$

- ◆ Implementation of the ERM for linearly separable examples  $\mathcal{T}^m$  leads to solving (2) which yields  $h(x; \mathbf{w}, b)$  with  $R_{\mathcal{T}^m}^{0/1}(h(\cdot; \mathbf{w}, b)) = 0$ .

Note that  $y^i(\langle \mathbf{w}, \phi(x^i) \rangle + b) > 0$  implies

$$h(x^i) = \text{sign}(\langle \mathbf{w}, \phi(x^i) \rangle + b) = y^i$$

- ◆ The linear programming task (2) can be solved by the Perceptron algorithm.

## Auxiliary prediction problem leading to tractable ERM

- ◆  $\mathcal{X}, \mathcal{Y} = \{+1, -1\}$  and  $\phi: \mathcal{X} \rightarrow \mathbb{R}^n$  defined as before.
- ◆ Auxiliary prediction problem: find a decision function  $f: \mathcal{X} \rightarrow \mathbb{R}$  minimizing the expectation of the hinge loss  $\psi: \mathcal{Y} \times \mathbb{R} \rightarrow \mathbb{R}_+$ :

$$R^\psi(f) = \mathbb{E}_{(x,y) \sim p}(\psi(y, f(x))) \quad \text{where} \quad \psi(y, t) = \max\{0, 1 - y t\}$$

- ◆ Assuming the hypothesis space which contains the linear functions

$$\mathcal{F} = \{f(x) = \langle \phi(x), \mathbf{w} \rangle + b \mid (\mathbf{w}, b) \in \mathbb{R}^{n+1}\}$$

the ERM principle leads to solving

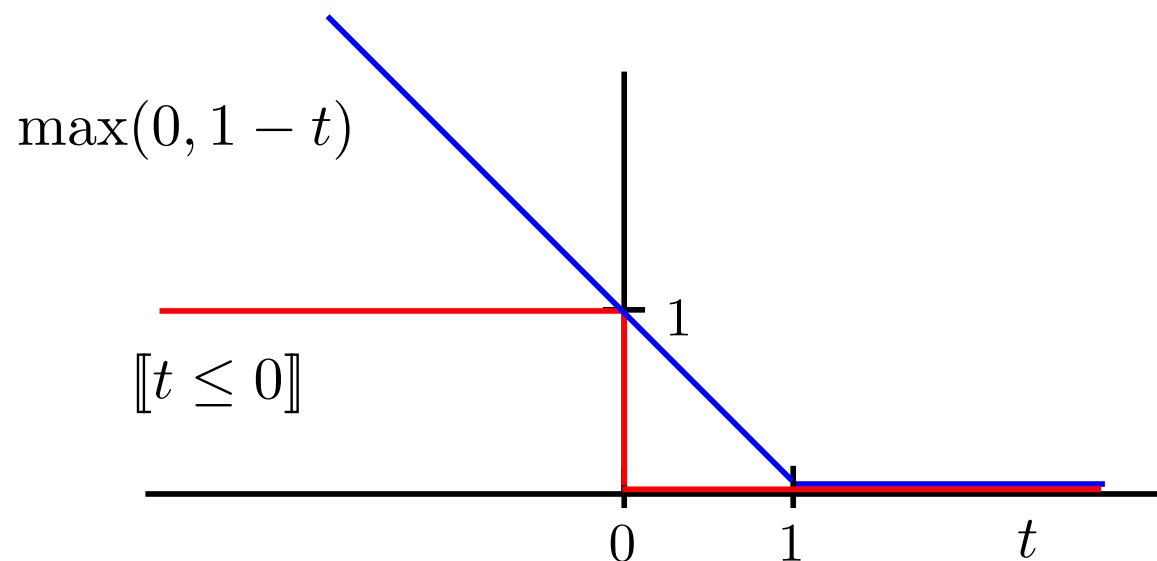
$$f^* = \underset{f \in \mathcal{F}}{\text{Argmin}} R_{\mathcal{T}^m}^\psi(f) \quad \text{where} \quad R_{\mathcal{T}^m}^\psi(f) = \frac{1}{m} \sum_{i=1}^m \psi(y^i, f(x^i))$$

- ◆ How is this task related to minimization of the classification error?

# The hinge-loss upper bounds the 0/1-loss

- ◆ The hinge-loss is an upper bound of the 0/1-loss evaluated for the predictor  $h(x) = \text{sign}(f(x))$ :

$$\underbrace{[\text{sign}(f(x)) \neq y]}_{\ell^{0/1}(y, f(x))} = [y f(x) \leq 0] \leq \underbrace{\max\{0, 1 - y f(x)\}}_{\psi(y, f(x))}$$



- ◆ Therefore 0/1-risk of  $h(x) = \text{sign}(f(x))$  is upper-bounded by  $\psi$ -risk:

$$R^{0/1}(\text{sign}(f)) \leq R^\psi(f) \quad \text{for any } f: \mathcal{X} \rightarrow \mathbb{R}$$



## Excess error of $\psi$ -risk upper bounds excess risk of 0/1-risk

- ◆ The best attainable 0/1-risk is  $R_*^{0/1} = \inf_{h \in \mathcal{Y}^{\mathcal{X}}} R^{0/1}(h)$ .
- ◆ The best attainable  $\psi$ -risk is  $R_*^\psi = \inf_{f \in \mathbb{R}^{\mathcal{X}}} R^\psi(f)$
- ◆ The best predictor in  $\mathcal{F}$  is  $f_{\mathcal{F}} \in \text{Argmin}_{f \in \mathcal{F}} R^\psi(f)$ .

**Theorem 3.** For any  $f: \mathcal{X} \rightarrow \mathbb{R}$  the following inequality holds

$$\underbrace{R^{0/1}(\text{sign}(f)) - R_*^{0/1}}_{\text{excess error of original task}} \leq \underbrace{R^\psi(f) - R_*^\psi}_{\text{excess error of auxiliary task}}$$

**Corollary 3.** Let  $A': \cup_{m=1}^{\infty} (\mathcal{X} \times \mathcal{Y}) \rightarrow \mathcal{F}$  be a learning algorithm statistically consistent in  $\mathcal{F} \subseteq \mathbb{R}^{\mathcal{X}}$  w.r.t.  $\psi$ -risk. In addition, let  $R^\psi(f_{\mathcal{F}}) = R_*^\psi$ . Then, the learning algorithm  $A(\mathcal{T}^m) = \text{sign}(A'(\mathcal{T}^m))$  is statistically consistent in  $\mathcal{H} = \{\text{sign}(f) \mid f \in \mathcal{F}\}$  w.r.t. 0/1-risk.

## Solving ERM problem of the auxiliary prediction task

- ◆ Let us consider a space of linear score functions with parameter vector inside a ball of radius  $r$ , that is,

$$\mathcal{F}_r = \{f(x) = \langle \phi(x), \mathbf{w} \rangle + b \mid (\mathbf{w}, b) \in \mathbb{R}^{n+1}, \|\mathbf{w}\| \leq r\}$$

- ◆ The ERM problem for  $\psi(y, t) = \max\{0, 1 - y t\}$  loss reads

$$f^* = \underset{f \in \mathcal{F}_r}{\text{Argmin}} R_{\mathcal{T}^m}^\psi(f) \quad \text{where} \quad R_{\mathcal{T}^m}^\psi(f) = \frac{1}{m} \sum_{i=1}^m \psi(y^i, f(x^i))$$

- ◆ The ERM problem is a convex unconstrained optimization task

$$(\mathbf{w}^*, b^*) = \underset{\|\mathbf{w}\| \leq r, b \in \mathbb{R}}{\text{argmin}} \left( \frac{1}{m} \sum_{i=1}^m \max\{0, 1 - y^i(\langle \mathbf{w}, \phi(x^i) \rangle + b)\} \right)$$

## Topics covered in the lecture

- ◆ Linear classifier
- ◆ Vapnik-Chervonenkis dimension
- ◆ Consistency and generalization bound for two-class prediction and 0/1-loss
- ◆ ERM problem for linear classifiers
- ◆ Auxiliary prediction problem ERM of which is tractable
- ◆ Excess error of the auxiliary problem upper bounds the excess error of the original problem

