

Supervised learning: data $(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}$

$\ell(y, f(x))$ loss function

\hookrightarrow prediction at x : $f(x) \in \mathcal{Y}$

goal: $\min \mathbb{E} \ell(y, f(x)) = R(f)$ -

\hookrightarrow testing distribution

- $f^*(x)$ Bayes predictor for square loss $\mathbb{E}(y|x)$
for 0-1 loss: $\arg \max p(y|x)$

• local averaging (k-NN)

• empirical risk minimization

$$- R(f) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i))$$

empirical risk

x dealing with discrete
x Estimation error

Binary classification: $y = \{-1, 1\}$ $\xrightarrow{g(x)}$

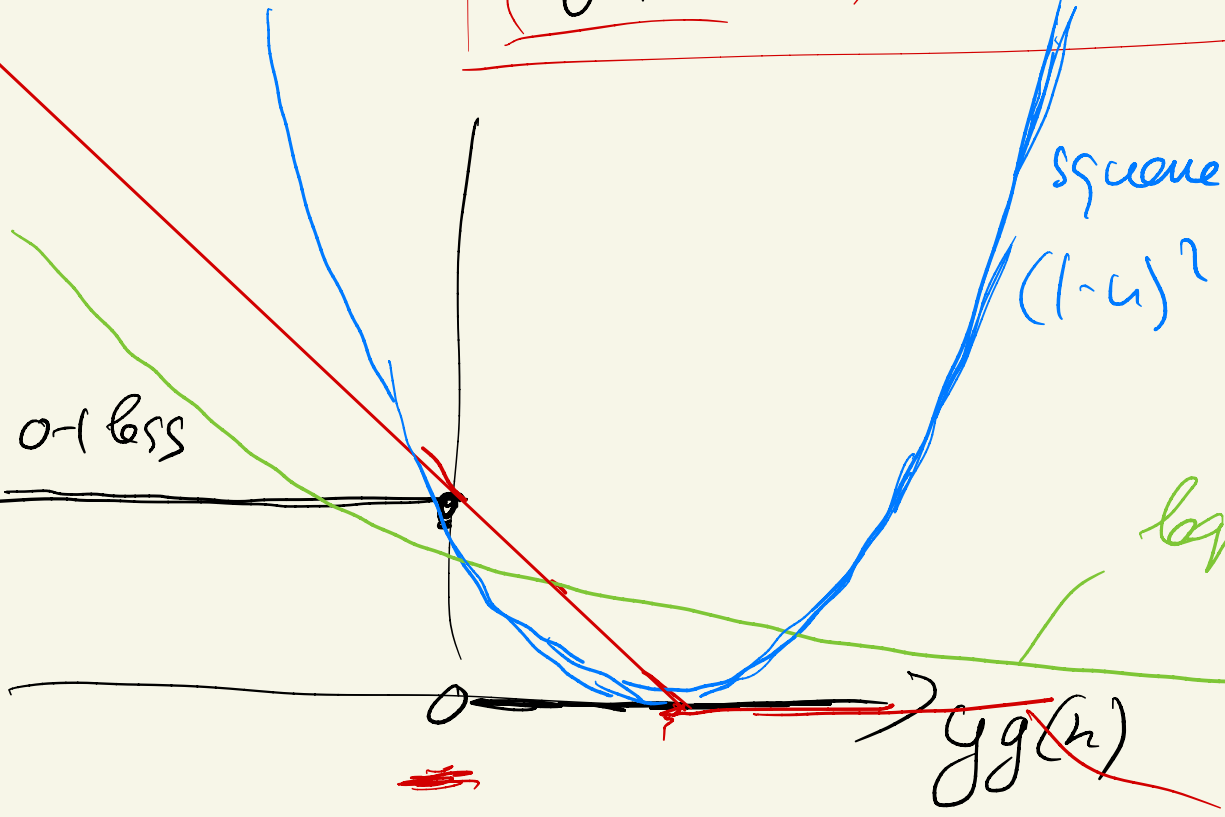
$f: X \rightarrow \{-1, 1\}$

$g: X \rightarrow \mathbb{R}$ with $f(x) = \text{sign}(g(x))$ $\begin{cases} = 1 & \text{if } g(x) > 0 \\ = 0 & \text{if } g(x) = 0 \\ = -1 & \text{if } g(x) < 0 \end{cases}$

0-1 loss:

$l(y, f(x)) = \begin{cases} 1 & y \neq f(x) \\ 0 & y = f(x) \end{cases} = \begin{cases} 1 & yg(x) \leq 0 \\ 0 & yg(x) > 0 \end{cases} = \varphi_{0-1}(yg(x))$

convex surrogate



square loss

$(1-u)^2 \quad (1-yg(x))^2 = (y-g(x))^2$

because $y \in \{-1, 1\}$

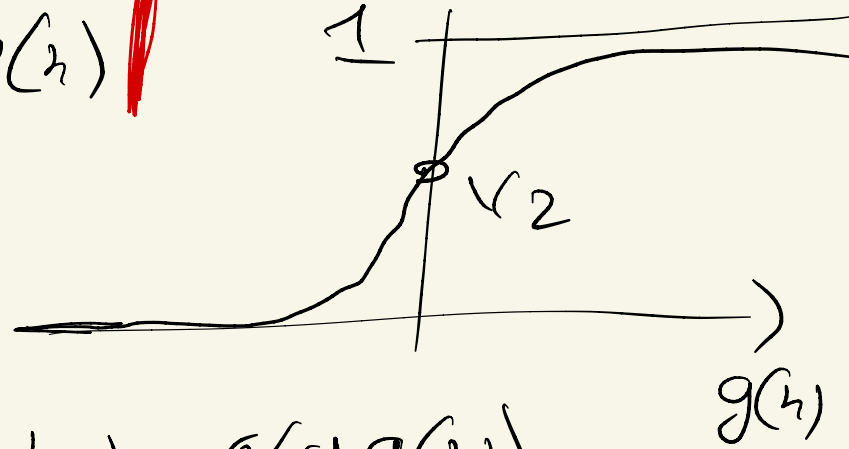
logistic loss $\log(1+e^{-u})$

$\max(1-u, 0) = (1-u)_+$
hinge loss \Rightarrow SVM

logistic loss: $\log(1+e^{-u})$!!

probabilistische Interpretation mit conditional probabilistic model

$$p(y=1|x) = \sigma(g(x)) = \frac{1}{1+e^{-g(x)}} \quad \text{Sigmoid}$$



$$p(y=-1|x) = 1 - \sigma(g(x))$$

$$= \sigma(-g(x)) \Rightarrow p(y|x) = \sigma(yg(x))$$

Maximum likelihood with independent data

$$\max_{\hat{x}} \sum_i \log p(y_i|x_i) = \max_{\hat{x}} \sum_i \log \sigma(y_i g(x_i))$$

$$= \max_{\hat{x}} \sum_i -\log(1+e^{-y_i g(x_i)})$$

"cross-entropy loss"

Analysis of convex surrogate:

0-1 error $R_{0-1}(g) = \mathbb{E} \phi_{0-1}(yg(h))$ -

$$R_{\phi}(g) = \mathbb{E} \phi(yg(h)) -$$

square, hinge, logistic

Goal: then $R_{0-1}(g) - R_{0-1}^*(g)$ is small \leftarrow
if $R_{\phi}(g) - R_{\phi}^*(g)$ is small \leftarrow

Square loss (1) check that optimal predictions are the same

$$(2) R_{0-1}(g) - R_{0-1}^*(g) \leq \sqrt{R_{\text{square}}(g) - R_{\text{square}}^*(g)}$$

optimal prediction:

0-1 loss: $f^*(x) = 1$ if $P(y=1|x) > 1/2 = \text{sign}(\mathbb{E}(y|x))$
 -1 if $P(y=-1|x) > 1/2$

$$\mathbb{E}(y|x) = 1 \cdot P(y=1|x) + (-1) \cdot P(y=-1|x) = 2P(y=1|x) - 1 = 2q(x) - 1$$

Lemma: $R_{0-1}(g) - R_{0-1}^*(g) \leq E |2y(n)-1 - g(n)|$

consequence

$$\leq E |g_{square}^*(n) - g(n)|$$

$$\leq \sqrt{E (g_{square}^*(n) - g(n))^2}$$

(Jensen)

$$= \sqrt{R_{square}(g) - R_{square}^*(g)}$$

proof: $R_{0-1}(g) - R_{0-1}^*(g)$

$$= E \left[\mathbb{1}_{yg(n) \leq 0} - \mathbb{1}_{yg^*(n) \leq 0} \right]$$

$$= E \left(E \left[\mathbb{1}_{yg(n) \leq 0} - \mathbb{1}_{yg^*(n) \leq 0} \mid n \right] \right)$$

non zero only when $g(n)$ and $g^*(n)$ have different signs

two cases: ① $g(n) > 0, g^*(n) < 0$; ② $g(n) < 0, g^*(n) > 0$

$$E \left[\mathbb{1}_{yg(n) \leq 0} - \mathbb{1}_{yg^*(n) \leq 0} \mid n \right] = \frac{1}{y=-1} - \frac{1}{y=+1} = p(y=-1(n)) - p(y=+1(n))$$

$$|0^{n(5-)} - 0^{n(3)}|$$

$$= 1 - 2q(n) \leq 1 - 2q(n) + g(n) = |(-2q(n) + g(n))|$$

Decomposition of the risk: $\ell(y, f(x))$ where $f(x) \in \mathbb{R}$
 impl. = $\ell(y, f(x))$ φ logistic σ

$R(\beta) = \mathbb{E} \ell(y, f(x))$ expected risk | if $y \in \{-1, 1\}$

$\hat{R}(\beta) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i))$

"generalization error"
 class of models
 $\mathcal{F} = \{ \beta: X \rightarrow \mathbb{R} \}$
 $R(\beta)$
 $\hat{R}(\beta)$

Empirical risk minimization

$\hat{\beta} \in \arg \min_{\beta \in \mathcal{F}} \hat{R}(\beta)$

$R(\beta)$ "Bayes error rate"
 $\mathcal{F} \subset \mathcal{F}^*$

goal:

$R(\hat{\beta}) - R^* = \underbrace{R(\hat{\beta}) - \inf_{\beta \in \mathcal{F}} R(\beta)}_{\text{estimation error}} + \underbrace{\inf_{\beta \in \mathcal{F}} R(\beta) - R^*}_{\text{approximation error}}$

estimation error $\rightarrow 0$
 random $\downarrow n$

approximation error $\rightarrow 0$
 det. must \subset under of \mathcal{F}

Estimation error =

$$R(\hat{\beta}) - \inf_{f \in \mathcal{F}} R(f) = R(\hat{\beta}) - R(\beta_{\mathcal{F}}^*)$$

arg min $R(f)$
 $f \in \mathcal{F}$

$$R(\hat{\beta}) - \hat{R}(\hat{\beta}) + \hat{R}(\hat{\beta}) - \hat{R}(\beta_{\mathcal{F}}^*) + \hat{R}(\beta_{\mathcal{F}}^*) - R(\beta_{\mathcal{F}}^*)$$

$$\leq 0$$

because $\hat{\beta} \in \text{arg min}_{f \in \mathcal{F}} \hat{R}(f)$

$$\leq \sup_{f \in \mathcal{F}} R(f) - \hat{R}(f)$$

$$+ \sup_{f \in \mathcal{F}} \hat{R}(f) - R(f)$$

$$\text{estimation error} \leq 2 \sup_{f \in \mathcal{F}} |R(f) - \hat{R}(f)|$$

(x_i, y_i) iid

if f is fixed =
Central limit theorem

$$\hat{R}(f) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i))$$

Hoeffding

$$\frac{1}{n} \sum Z_i \sim \text{Normal}\left(\mathbb{E}Z, \frac{\sigma^2}{n}\right)$$

$$= \frac{1}{n} \sum Z_i \rightarrow \text{FLZ}$$

law of large #'s

Hoeffding's inequality: $Z_i \text{ IID}$, $Z_i \in [a, b]$ almost surely

$$P\left(\left|\frac{1}{n} \sum_i Z_i - \mathbb{E}Z\right| \geq t\right) \leq 2e^{-2nt^2} = \delta$$

with probability greater than $1 - \delta$

$$\left|\frac{1}{n} \sum_i Z_i - \mathbb{E}Z\right| \leq \sqrt{\frac{\log \frac{2}{\delta}}{2n}}$$

$$2nt^2 = \log \frac{2}{\delta}$$

$$t = \sqrt{\frac{\log \frac{2}{\delta}}{2n}}$$

Proof: look at chapter 1

$$P(S) \leq \frac{\mathbb{E}(e^{\lambda S})}{e^{\lambda t}}$$

Markov / inequalities
Chebyshev

$$Z \in (a, b) \Rightarrow 2e^{-\frac{2nt^2}{(b-a)^2}}$$

Finite number of models: $\mathcal{F} = \{f_1, \dots, f_m\}$, $|\mathcal{F}| = m$

$\sup_{f \in \mathcal{F}} |R(f) - \hat{R}(f)|$

$\forall f, P(|R(f) - \hat{R}(f)| \geq t) \leq 2 \exp(-\frac{nt^2}{4\epsilon_0^2})$

if $R(f)$ almost surely $\in [a, b]$ $4\epsilon_0^2$

$R(f) \in [-\epsilon_0, \epsilon_0]$

$P(\sup_{f \in \mathcal{F}} |R(f) - \hat{R}(f)| \geq t) = P(\cup_{f \in \mathcal{F}} \{|R(f) - \hat{R}(f)| \geq t\})$

union bound $\leq \sum_{f \in \mathcal{F}} P(|\hat{R}(f) - R(f)| \geq t)$

$\leq 2e^{-nt^2/2\epsilon_0^2}$

With proba $1 - \delta$

$\sup_{f \in \mathcal{F}} |R(f) - \hat{R}(f)| \leq \sqrt{\frac{\log 2/|\mathcal{F}|/\delta}{n} \frac{2\epsilon_0^2}{2}} = \sqrt{\frac{2\epsilon_0^2}{n}} \sqrt{\log |\mathcal{F}| + \log \frac{2}{\delta}} = \delta$

Rademacher complexity:

definition

$$\mathcal{H} = \left\{ h: (x, y) \mapsto \ell(y, h(x)) \right\}$$
$$\sup_{f \in \mathcal{F}} \bar{R}(f) - R(f) = \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i)) - \mathbb{E} \ell(y_i, f(x_i))$$
$$= \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n h(z_i) - \mathbb{E} h(z) \quad (z = (x, y))$$

goal: get a bound on

$$\mathbb{E} \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n h(z_i) - \mathbb{E} h(z) \leq 2R_n(\mathcal{H})$$

Def. Rademacher complexity of \mathcal{H} .

$$R_n(\mathcal{H}) = \mathbb{E} \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \varepsilon_i h(z_i)$$

with ε a vector of Rademacher random variables

$$P(\varepsilon_i = \pm 1) = 1/2 \quad \forall i \in \{1, \dots, n\}$$

Symmetrization lemma:

Proof: $D = \{z_1, \dots, z_n\}$, independent copy $D' = \{z'_1, \dots, z'_n\}$

$$E_{D, h \in \mathcal{H}} \sup \frac{1}{n} \sum_i h(z_i) - \underbrace{E h(z)} = E_{D, D', h \in \mathcal{H}} \sup \frac{1}{n} \sum_i \left[h(z_i) - \underbrace{E(h(z'_i) | D)} \right]$$

$$= E_{D, D'} \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_i E(h(z_i) - h(z'_i) | D) \quad \text{sup } E \leq E \text{ sup}$$

$$= E \left(\sup_{h \in \mathcal{H}} \frac{1}{n} \sum_i h(z_i) - h(z'_i) \mid D \right)$$

$$\leq E_{D, D'} \left(\sup_{h \in \mathcal{H}} \frac{1}{n} \sum_i h(z_i) - h(z'_i) \mid D \right)$$

$$E E(\cdot | D) = E \cdot$$

$$= E_{D, D'} \left(\sup_{h \in \mathcal{H}} \frac{1}{n} \sum_i \varepsilon_i (h(z_i) - h(z'_i)) \right)$$

$$\frac{1}{n} \sum_i \varepsilon_i h(z_i) - \frac{1}{n} \sum_i \varepsilon_i h(z'_i) \quad \text{sup } A \in \mathbb{R} \leq \text{sup } A + \text{sup } B$$

$$\leq E_{D, D', \varepsilon} \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_i \varepsilon_i h(z_i) + E_{D, D', \varepsilon} \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_i (-\varepsilon_i) h(z'_i)$$

$$R_n(f) \leq R_n(g) = 2R_n(H)$$

$$h(z) = h(x, y) = \ell(y, f(x))$$

$$\begin{aligned} & \sum \varepsilon_i \ell(\overline{y_i}, \overline{f(x_i)}) \quad \text{difficult} \\ & = \sum \varepsilon_i \varphi_i(\overline{f(x_i)}) \end{aligned}$$

see proof in book

Contraction principle

Lipschitz continuous with constant G

$$|\varphi_i(a) - \varphi_i(b)| \leq G|a - b|$$

Prop: $\forall \varepsilon \sup_{f \in \mathcal{F}} \sum_{i=1}^n \varepsilon_i \varphi_i(\overline{f(x_i)}) \leq G \cdot \varepsilon \sup_{f \in \mathcal{F}} \sum_{i=1}^n \varepsilon_i \overline{f(x_i)}$

Consequence: φ Mebss function are G -Lipschitz

$$\|R_n(H)\| \leq G \cdot \|R_n(F)\|$$

Example: linear functions: $\mathcal{F} = \{ \theta: x \rightarrow \theta^T \varphi(x) \mid \|\theta\|_2 \leq D \}$

$$R_n(F) = \mathbb{E} \sup_{\|\theta\|_2 \leq D} \frac{1}{n} \sum_{i=1}^n \varepsilon_i \theta^T \varphi(x_i)$$

$$\|\theta\|_2 \leq D$$

$$R(F) = \mathbb{E} \sup_{\| \theta \|_2 \leq D} \left(\frac{1}{n} \sum_i \varepsilon_i \underbrace{\phi(x_i)^T}_{f(x_i)} \theta \right)$$

$$= \mathbb{E} \left\| \frac{1}{n} \sum_i \varepsilon_i \phi(x_i) \right\|_2 \times D$$

$$\leq \sqrt{\mathbb{E} \left\| \frac{1}{n} \sum_i \varepsilon_i \phi(x_i) \right\|_2^2} \times D$$

$$= \sqrt{\frac{1}{n^2} \sum_i \mathbb{E} \varepsilon_i^2 \|\phi(x_i)\|_2^2} \times D$$

$$R(F) = \frac{1}{\sqrt{n}} D \times \mathbb{E} \|\phi(x)\|_2^2$$

$$\sup_{\| \theta \|_2 \leq D} \theta^T \theta = \| \theta \|_2 \cdot D$$

Jensen's inequality

because

$$\mathbb{E} \varepsilon_i = 0$$

+ i.i.d

Next week:

• Summary chapter 4

• chapter 5:
optimization