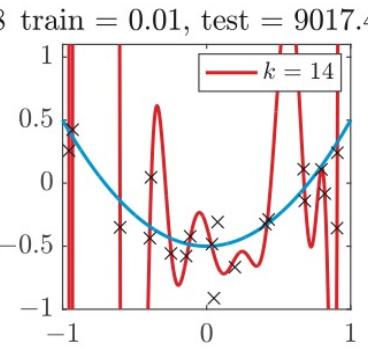
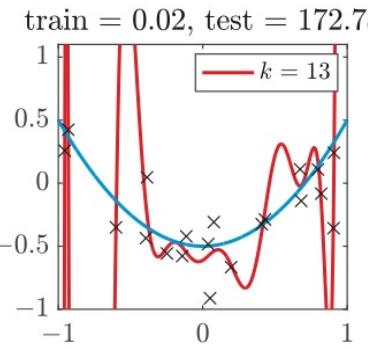
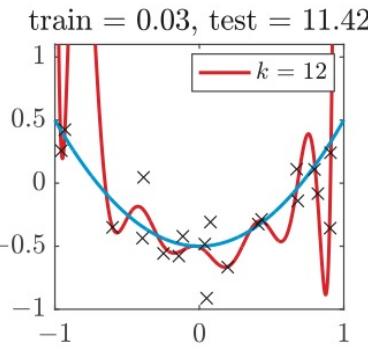
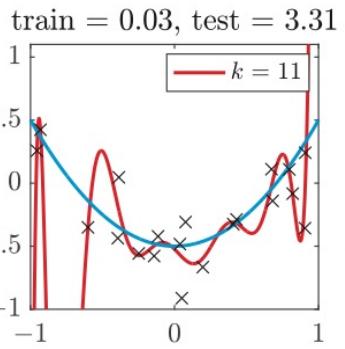
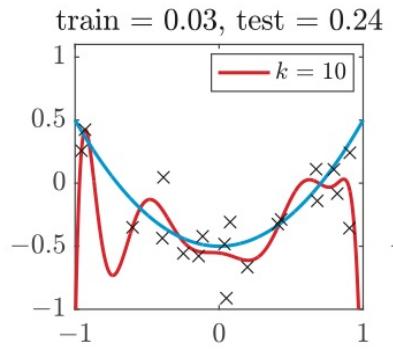
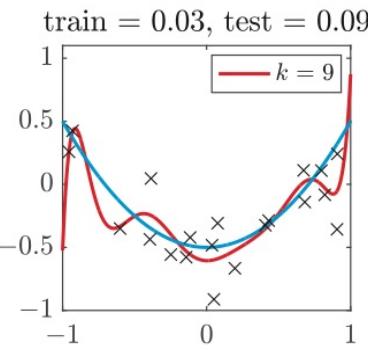
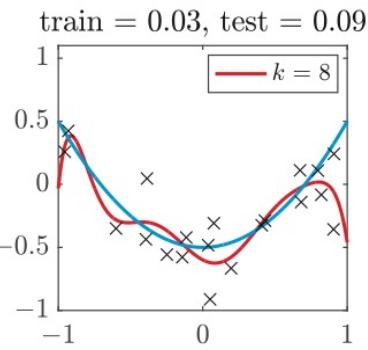
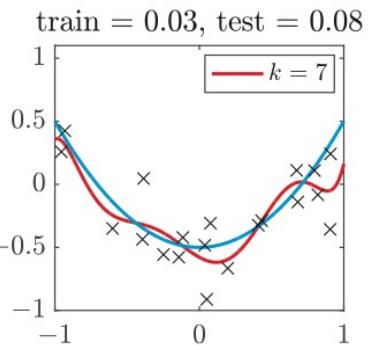
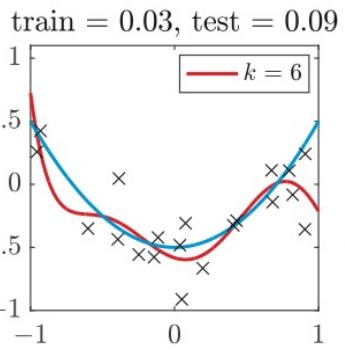
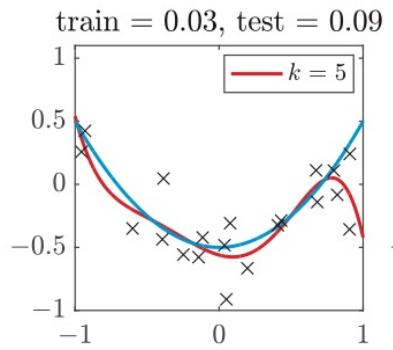
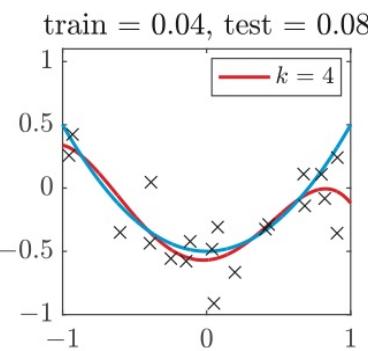
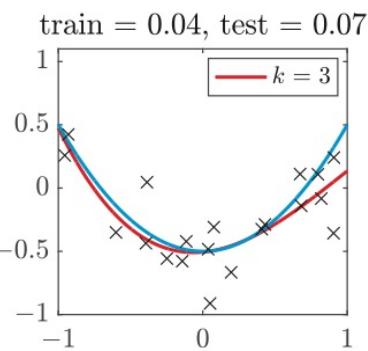
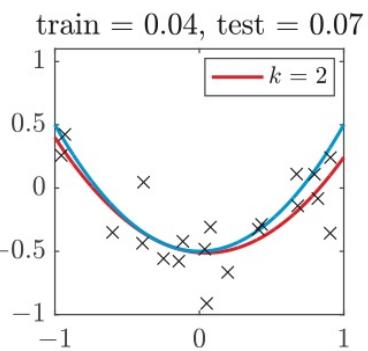
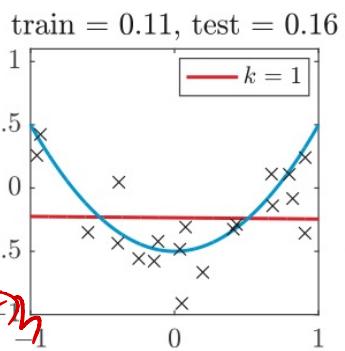
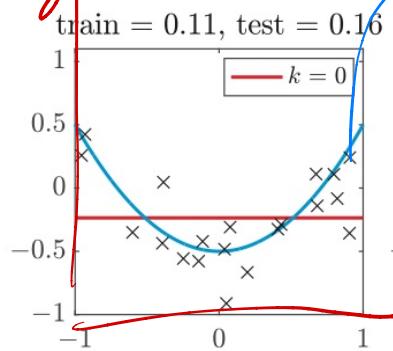


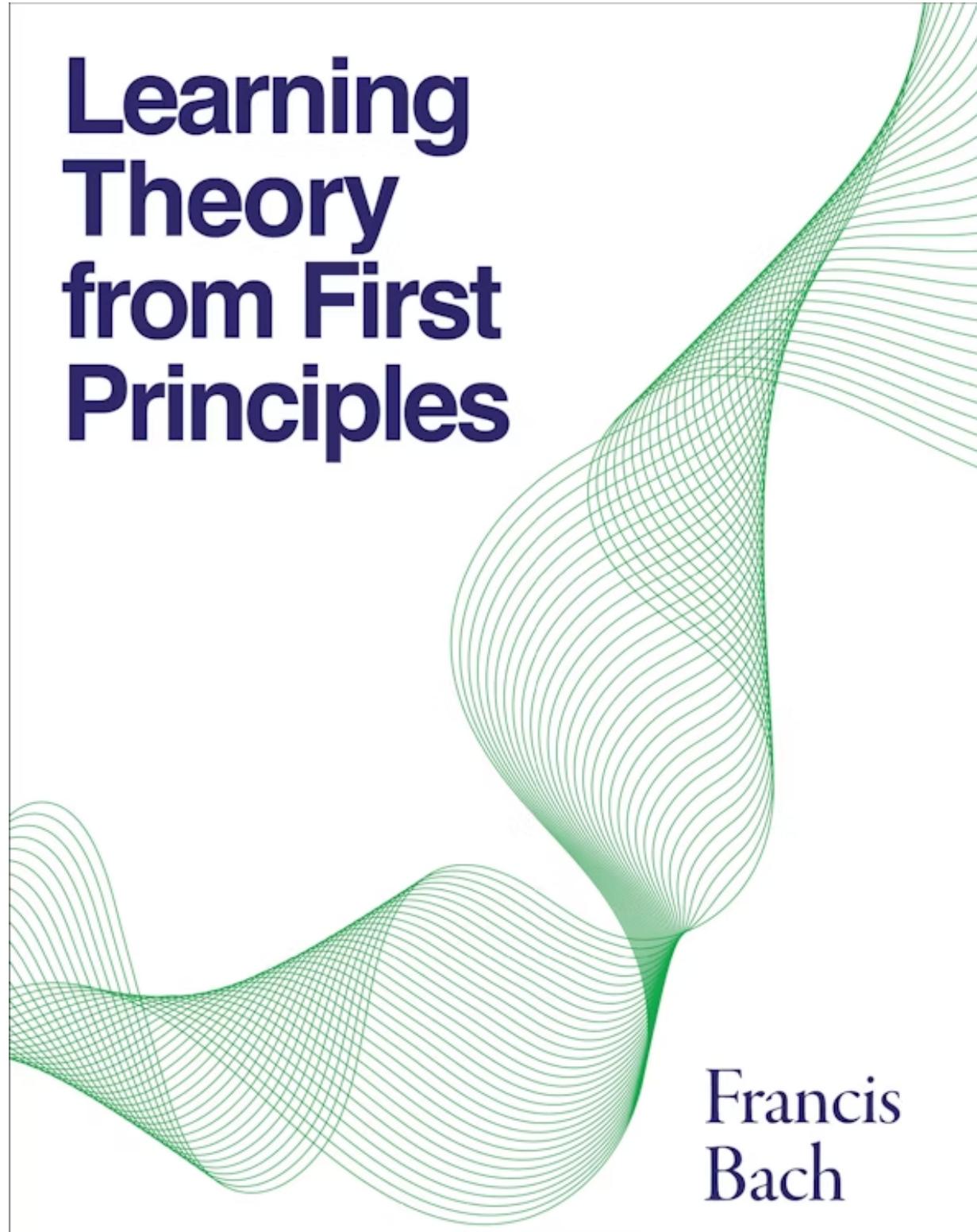
✓

$$f^k(x) = F(S^k)$$

$f_0(x)$ = polynomials of order 0

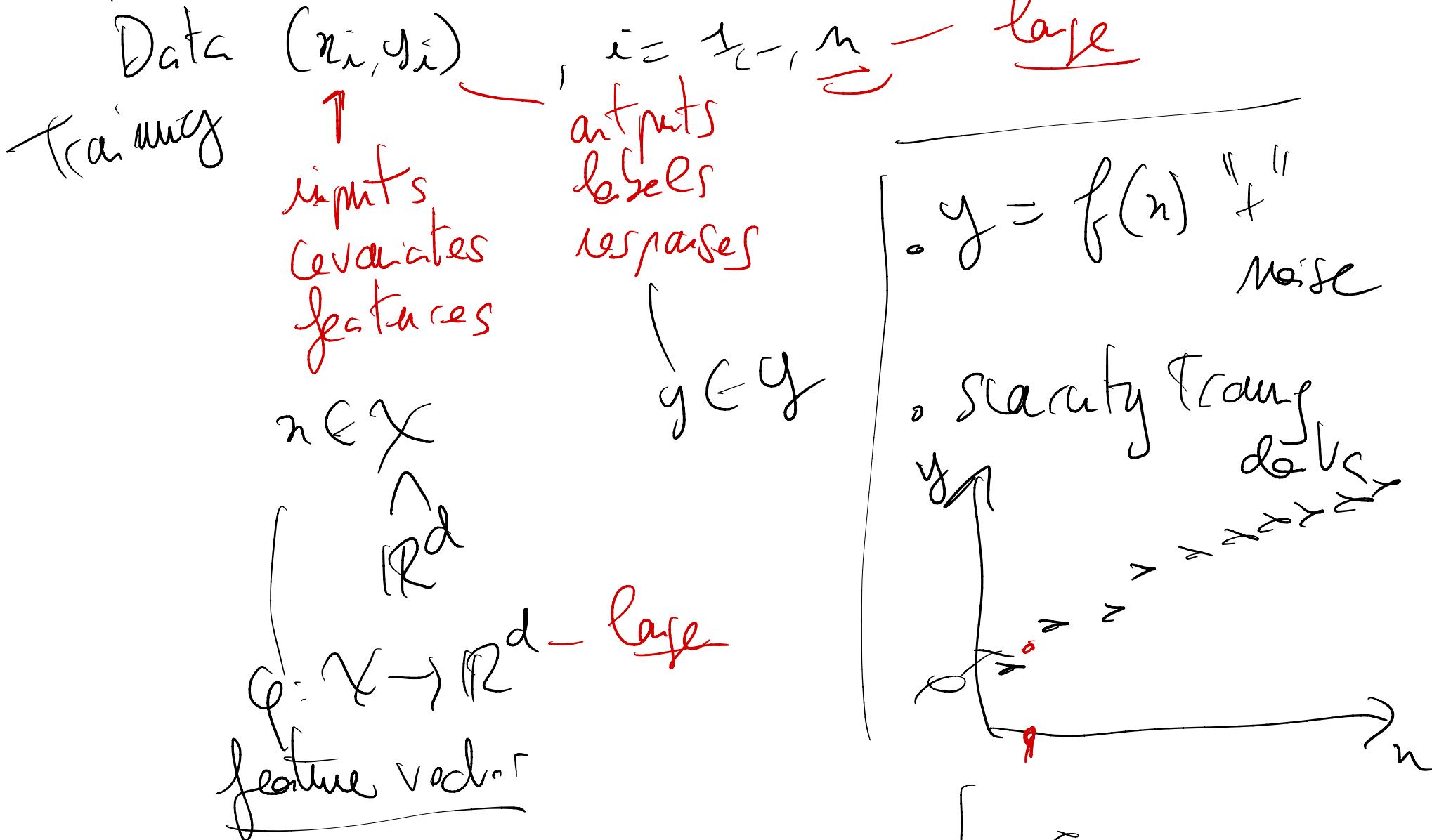


Learning Theory from First Principles

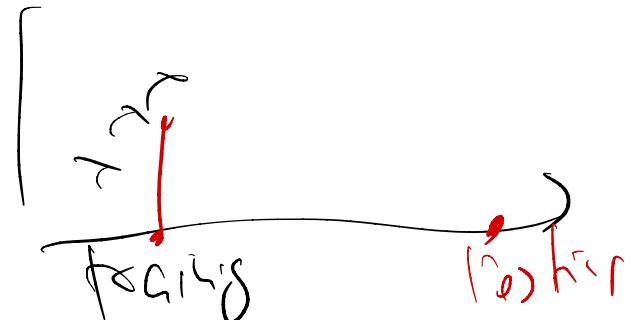


Francis
Bach

Supervised Learning



Test performance: $x \in \mathcal{X}$



Probabilistic formulation

(x_i, y_i) IID $i = 1 \rightarrow n$
 independent and identically distributed
 from distribution $p(x, y) = \text{test dist. func.}$

Research: dependence through Markov chains
 domain adaptation

Practical performance evaluation

Test data $(x_i^{test}, y_i^{test})_{i \in \mathcal{I}^T test}$

$$\left(E \ell(y, f(x)) - \frac{1}{n} \sum_{i=1}^n \ell(y_i^{test}, f(x_i^{test})) \right) = O\left(\frac{1}{\sqrt{n}}\right)$$

log ~~linear~~

[Training { validate }
 } test]

Decision theory : learning when best distribution is known

distribution $p(x, s)$ on $X \times Y$

$$\mathbb{E}[f(z)] = \int f(z) dp(z)$$

$$\mathbb{E}[g(z, s)] = \int_X \int_{Z \times Y} g(z, s) dp(z, s)$$

Loss function : $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$

$\ell(y, z)$ = loss of predicting z while the observed output is y

Binary classifier : 0-1 loss $\ell(y, z) = \mathbf{1}_{y \neq z}$

$$y = \{0, 1\}, \quad \boxed{\{-1, 1\}}, \quad \{1, 2\}$$

Multicategory : $y = \{1, \dots, k\}$

Regression : $y = \mathbb{R}$ $\ell(y, z) = (y - z)^2 \text{ or } \frac{1}{2}(y - z)^2$
 $\log(1 - e^{-y f(z)})$

Expected risk
generalization prob
test error
testing error

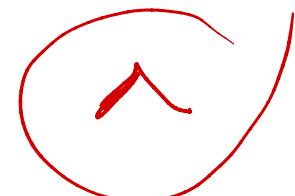
$$f: X \rightarrow Y$$

$$R(f) = E[e(y, f(x))] \quad \text{testing data}$$

"function"

two sources of randomness

training data
testing data



Empirical risk
training error

$$\hat{R}(f) = \frac{1}{n} \sum_{i=1}^n e(y_i, f(x_i)) \quad \xrightarrow{\text{to training data}}$$

Examples: classification $f: X \rightarrow Y$ finite with 0-1 loss
 $R(f) = E(1_{y \neq f(x)}) = P(y \neq f(x))$ error rate

Regression: $R(f) = E(y - f(x))^2$ "least-squares" \leftarrow accuracy

Bayes Risk and Bayes predictor

$$R^* = R(\hat{f}^*)$$

(is optimal $\hat{f} \sim f^*$)

$$S(f) = \mathbb{E}\{e(g, f(z))\} = \mathbb{E} \left[\underbrace{\mathbb{E}(e(y, f(z))|z)}_{\text{funct. a. of } z} \right]$$

To simplify: $X = \text{finite}$

$$= \sum_{n' \in X} \mathbb{E} \left[e(g, f(z')) | z = z' \right] \cdot p(z = z')$$

minimize separately.

$$f: X \rightarrow Y$$

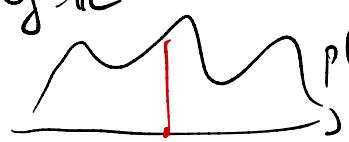
$$f(z') \forall z' \in X$$

(classif.
y finite)

$$\begin{aligned} f^*(z') &= \arg \min_{g \in \mathcal{G}} P(g \neq y | z = z') \\ &= \arg \max_{g \in \mathcal{G}} P(g = y | z = z') \end{aligned}$$

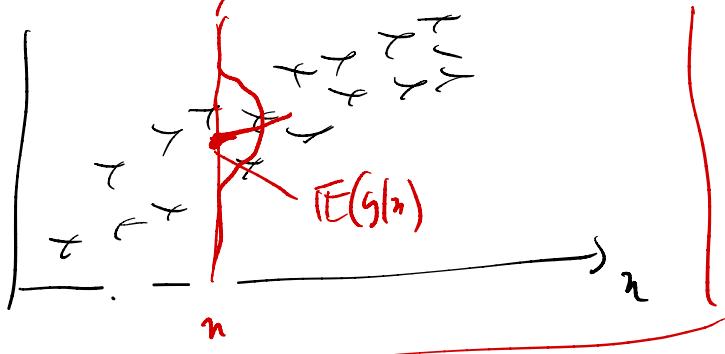
$$f^*(z') \in \arg \min_{g \in \mathcal{G}} \mathbb{E}(e(y, g) | z = z')$$

Regression: $f^*(z') \in \arg \min_{g \in \mathcal{G}} \mathbb{E}((y - g)^2 | z = z')$
 $y \sim \mathcal{N}$



law of total variance
 $\mathbb{E}(y - g)^2 = (\mathbb{E}y - g)^2 + \text{var}(y)$

Lemmas: $\arg \min_a \mathbb{E}_{y \sim \mathcal{G}} ((y - a)^2) = \mathbb{E}y$
 $\mathbb{E}((y - a)^2) = \mathbb{E}(y^2 - 2ay + a^2) = (\mathbb{E}y)^2 + s^2 - 2a\mathbb{E}y + a^2$
 $\frac{\partial}{\partial a} = -2\mathbb{E}y + 2 = 0$



Boges Modell der $f^*(x) = \underset{g}{\operatorname{argmin}} \mathbb{E}[e(g, \xi) | n=x]$

Boges risk: $R(f^*) = \mathbb{E}_{n \mid} \min \mathbb{E}[e(g, \xi) | n=x] = R^*$

Coverian : $\mathcal{R}(f) - R^*$ excess risk

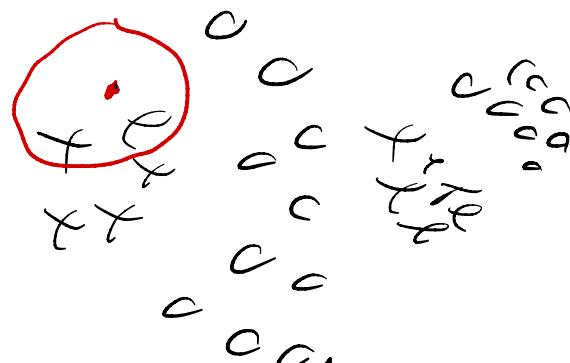
$(C_{BS}) \rightarrow (C^L SR)$

Two classes:

- ✗ local averaging
- ✗ Empirical risk min.

local averaging

$$\hat{f}(x) = \arg \min E\{e(y, g)\}_{n=x} = \int e(y, g) \underbrace{dp(y|x)}_{\downarrow}$$

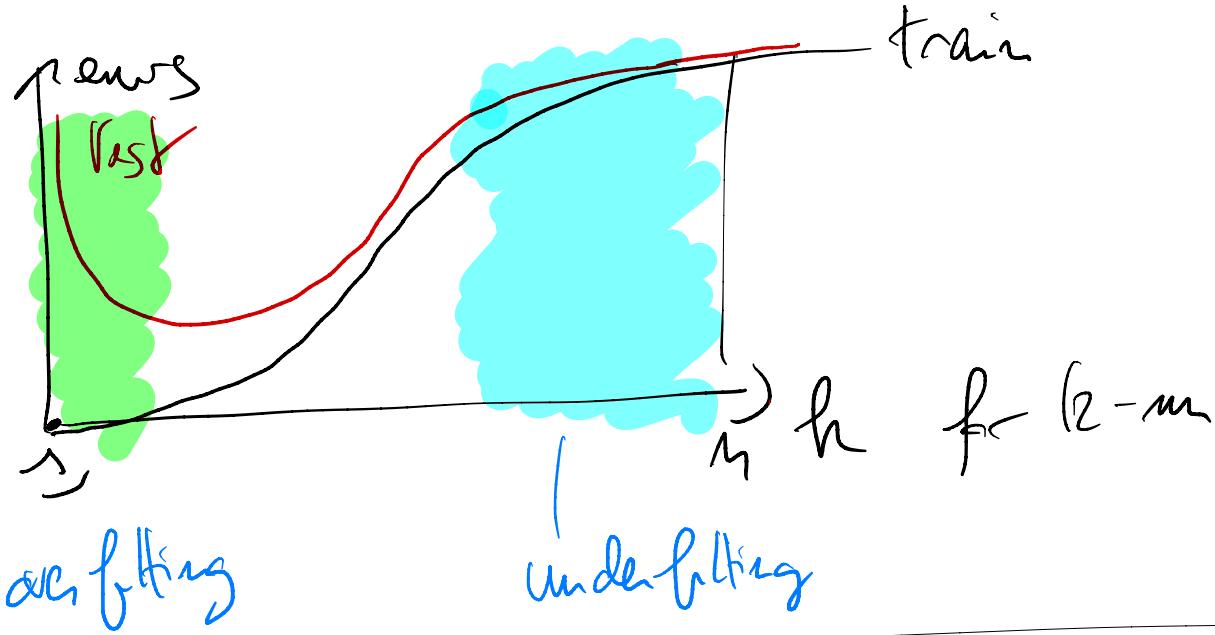


$$dp(y|x) = \sum_{i=1}^n w_i(x) \text{Diac}(y|y_i) \rightarrow \begin{array}{l} K\text{-nearest neighbors} \\ \text{Nadaraya-Watson} \end{array}$$

Ex: square loss: $\hat{f}(x) = \sum w_n(x) y_n$

0-1 loss: majority rule weighted by w_n

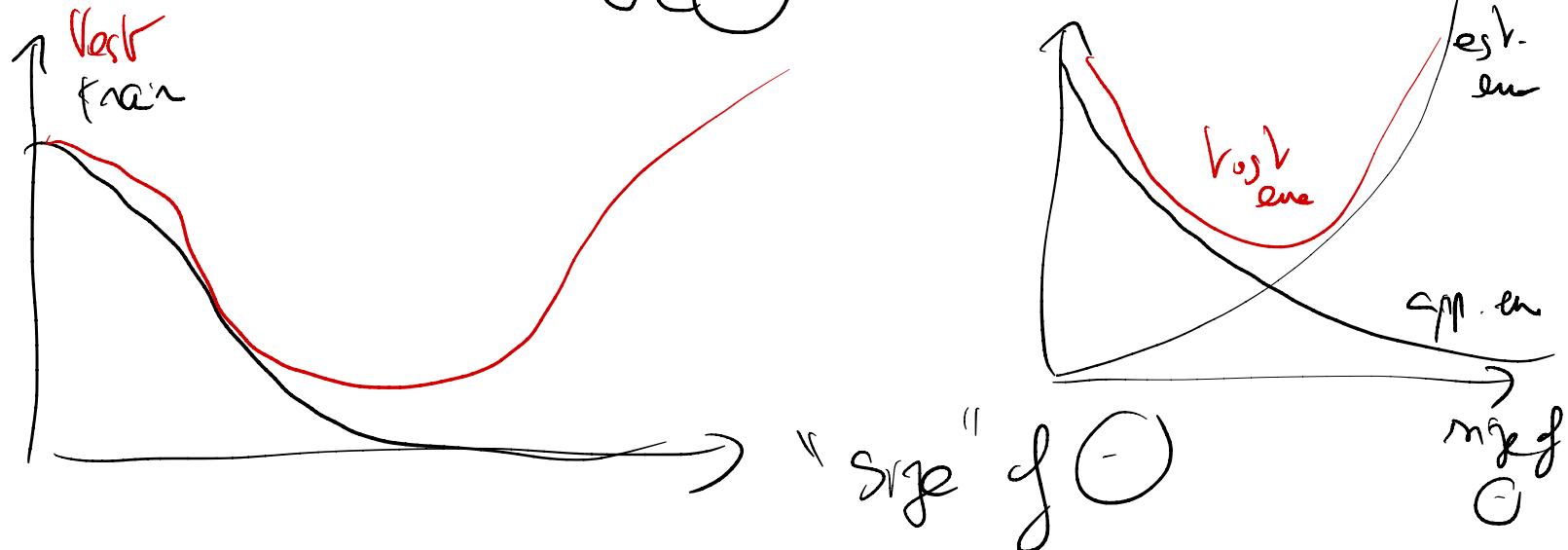
$\hat{dp}(s|x)$
training data dependent
smoothing



Empirical Risk Minimization : ERM

Model $f_\theta: \mathcal{X} \rightarrow \mathcal{Y}$ of functions parameterized by $\theta \in \Theta$

$$\hat{f} = f_{\hat{\theta}} \text{ s.t. } \hat{\theta} \in \arg \min_{\theta \in \Theta} \hat{R}(f_\theta) = \frac{1}{n} \sum_{i=1}^n e(y_i, f_\theta(x_i))$$



Decomposition between estimator error & approximator error
Model $\{f_{\theta}, \mathcal{O}(\Theta)\}$.

$$\text{Criteria: } R(f_{\theta}) - R_x = R(f_{\theta}) - \underbrace{\inf_{\mathcal{O}(\Theta)} R(f_{\theta})}_{\geq 0} + \inf_{\mathcal{O}(\Theta)} R(f_{\theta}) - R_x$$

estimator error

approximator error

≥ 0

random

decreasing in
increasing of Θ

deterministic

independent in
decreasing Θ

goal of supervised learning

Data $D_n(p) = (x_i, y_i)_{i=1, \dots, n}$ $\xrightarrow{\text{pred. dist. } p}$ $\hat{f}: X \rightarrow Y$
 Mlese
 f

Prediction function
 $\hat{f}: X \rightarrow Y$

$A: \text{datasets} \rightarrow \text{functions}$

Given $R_p(\hat{f}) - R_p^* = R_p(A(D_n(p))) - R_p^*$

2 ways of removing randomness \rightarrow the expectation

$E_{\text{training data}} R_p(A(D_n(p))) - R^* \leq R_p^*(p)$

$\rightarrow E_{\text{testing data}}$

\rightarrow (high pred. $HSE(G_1)$, with $mse \geq -s$)
 $R_p(A(D_n(p))) - R^* \leq R_p^*(s)$

Markov inequality: $X \geq c$

$$\forall \epsilon, P(X > \epsilon) \leq \frac{E[X]}{\epsilon}$$

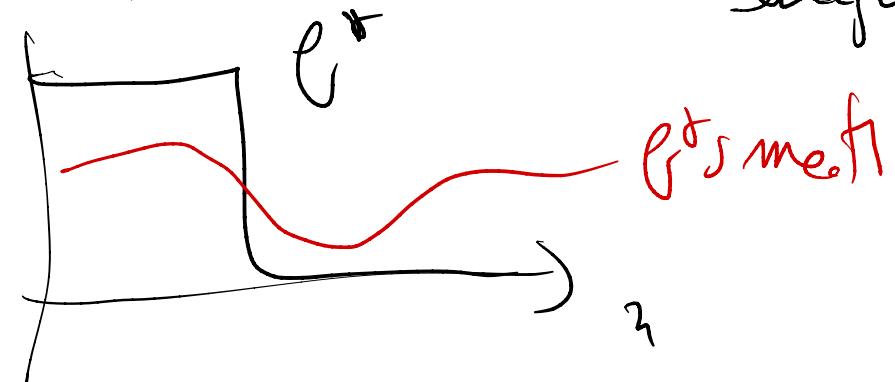
$\exists \delta$ with prob $[-\delta, X < \frac{E[X]}{\delta}]$

$= \delta$

$$\sup_{P \in \mathcal{P}} \mathbb{E}_{\text{training data}} R_p(\hat{A}(D_n)) - R_p^*$$

No free lunch

$\mathcal{P} = \{g(x, \gamma), \quad g^\alpha(x) \text{ smooth enough}\}$



- least squares
- ERM
- optimization
- local or
kernel method
- sparse ps
- NN