

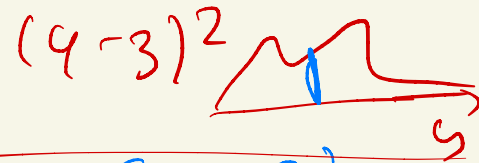
LOCAL AVERAGING | distribution p on $X \times Y$
 data (x_i, y_i) $i = 1, \dots, n$ IID from $p(x, y)$

loss function $\ell(y, z)$

Goal: find f such that $R(f) = \mathbb{E} \ell(y, f(x))$ for $f: X \rightarrow Y$
 Expected risk

Bayes predictor: $f^*(x) = \arg \min_{z \in Y} \mathbb{E}_{p(y|x)} \ell(y, z)$

square loss



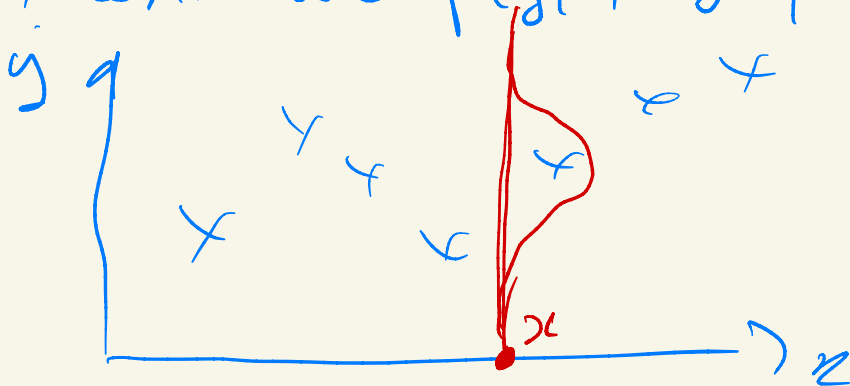
0-1 loss for class $f: f^*(x) = \arg \max_{z \in \{1, \dots, k\}} p(y = z | x)$

$y \in \{1, \dots, k\}$

square loss: $f^*(x) = \mathbb{E}(y|x)$
 $y \in \mathbb{R}$

Principle: estimate $p(y|x)$ by $\hat{p}(y|x)$

plug-in
 estimators

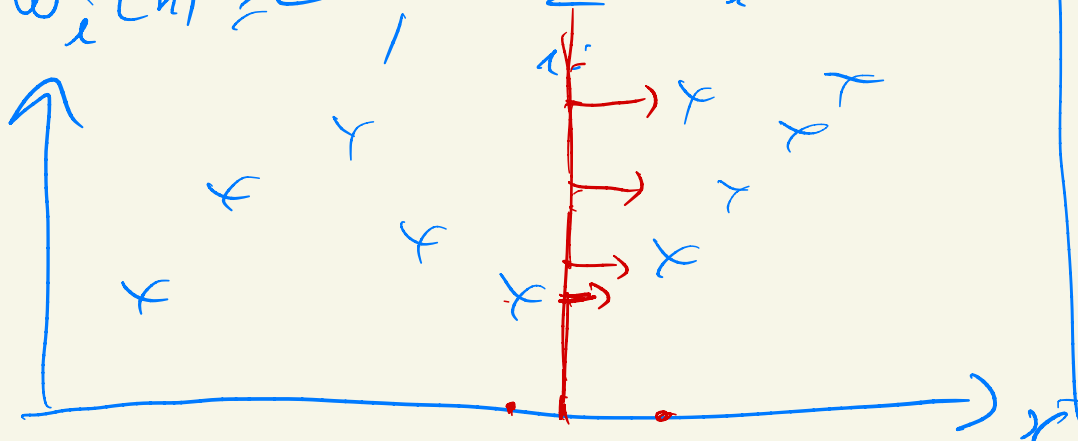


"Linear estimators": $\hat{p}(y|x) = \sum_{i=1}^n w_i(x) \delta_{y_i}(y)$ Dirac at y_i
 weight function $w_i: X \rightarrow \mathbb{R}$

To get a probab. measure:

$\forall x \quad w_i(x) \geq 0$

$\sum_{i=1}^n w_i(x) = 1$



0-1 loss: $\hat{p}(y=3|x)$
 $= \mathbb{E}_{\hat{p}(y|x)} \mathbb{1}_{y=3}$
 $= \sum_{i=1}^n w_i(x) \mathbb{1}_{y_i=3}$
 Majority vote

Estimators: $\hat{p}(x) = \underset{z}{\operatorname{arg\,min}} \mathbb{E}_{\hat{p}(y|x)} \ell(y, z)$
 $\underset{z}{\operatorname{arg\,min}} \sum_{i=1}^n w_i(x) \ell(y_i, z)$

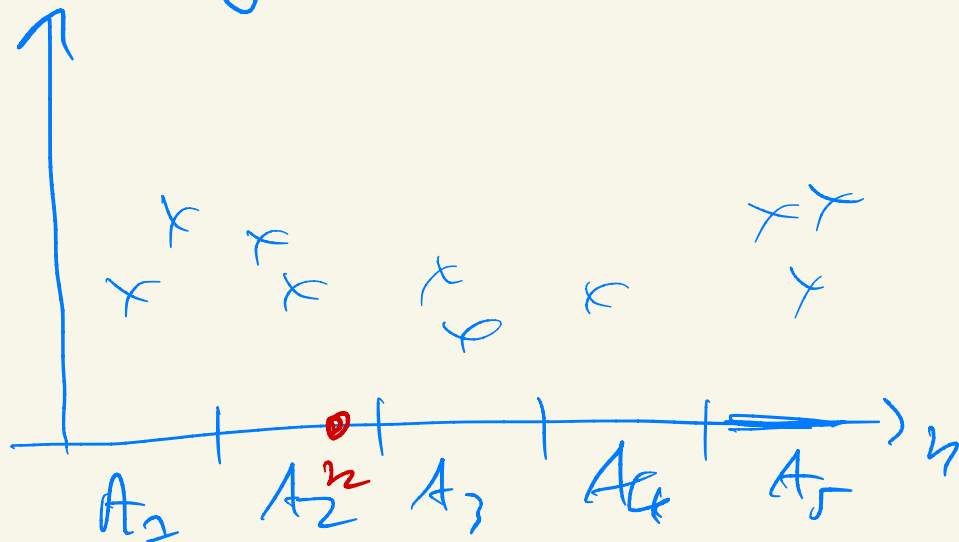
$\mathbb{E}_{\hat{p}(y|x)} \ell(y, z)$
 $\sum_{i=1}^n w_i(x) \ell(y_i, z)$



Square loss: $\hat{p}(x) = \mathbb{E}_{\hat{p}(y|x)} y = \sum_{i=1}^n w_i(x) y_i$

weighted average of all labels

① Partitioning estimate



weight for cell A_j

$v_j =$ number of observations in A_j

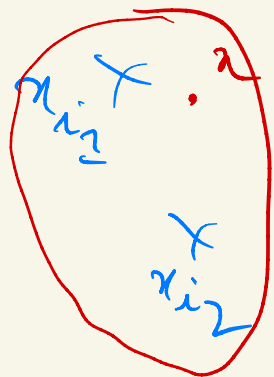
$$w_j(x) = \frac{1}{v_j} \text{ if } x \in A_j$$

1
2 2 2 1 3

② Nearest neighbor: need a metric d on X

given x , order all $d(x, x_i) : d(x, x_{i_1}) \leq \dots \leq d(x, x_{i_k}) \leq \dots \leq d(x, x_{i_n})$

Assign weight $\frac{1}{k}$ to the k k -nearest



x_{i_3}
 x_{i_5}
 x_{i_6}

Running time complexity

kd-trees $\leftarrow O(nd)$ for each test point

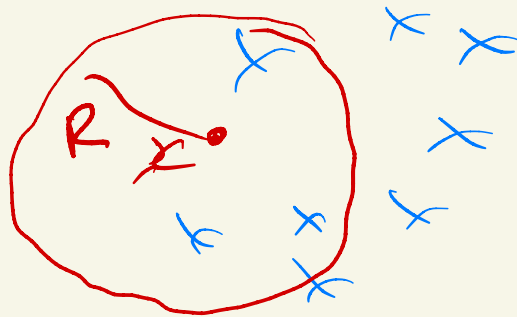
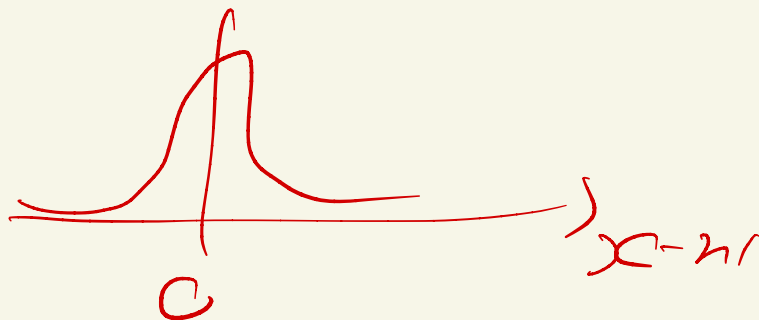
③ kernel regression / Nadaraya-Watson estimator

Two types of kernels. TODAY: $h(x, x') \geq 0$
usually, $h(x, x') = q(x - x')$

$$w_i(x) = \frac{h(x, x_i)}{\sum_{j=1}^n h(x, x_j)}$$

\Downarrow

ex: $h(x, x') = 1$ if $|x - x'| \leq R$
 $= 0$ otherwise

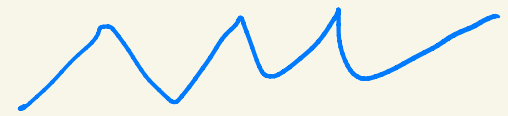


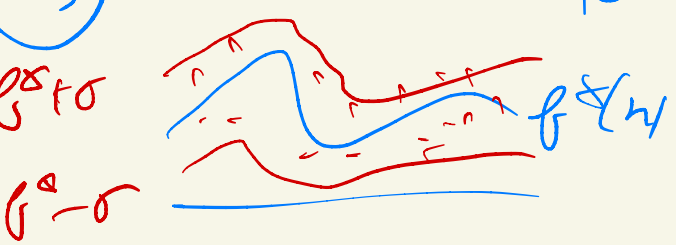
Analysis of local averaging:

linear reg: $f^*(x) = \theta^T x / \frac{\sigma^2 d}{n}$

(t1) $y = \mu + \text{square loss}$, with $f^*(x) = \mathbb{E}(y|x)$

(t2) f^* is B -Lipschitz continuous, $|f^*(x) - f^*(x')| \leq B d(x, x')$

(t3) Bounded noise: $|y - f^*(x)| \leq \sigma$ almost surely \Rightarrow 



w_i only depends on n_1, \dots, n_n

Goal: $\hat{f}(x) = \sum_{i=1}^n w_i(x) y_i$

$R(\hat{f}) - R(f^*)$ is small $\xrightarrow{\text{and du}}$ (test dist)

with $R(f) = \mathbb{E}(y - f(x))^2$

$$R(\hat{f}) - R(f^*) = \mathbb{E} \left[(f(x) - f^*(x))^2 \right]$$

$$= \int (f(x) - f^*(x))^2 p(x) dx$$

~~$\frac{\sigma^2}{n^{1/d}}$~~

$\frac{1}{n^{1/d}} = 0.1$

$\frac{1}{n^{1/d}} = 0.1$

Curse of dimensionality

Goal $\mathbb{E} R(\hat{f}) - R(f^*)$

leaves choice

$$\hat{f}(x) - f^*(x) = \sum_{i=1}^n w_i(x) y_i - f^*(x)$$

n fixed

$$= \sum_{i=1}^n w_i(x) [y_i - \mathbb{E}(y_i|x_i)] + \sum_{i=1}^n w_i(x) [\mathbb{E}(y_i|x_i) - f^*(x)]$$

(because $\sum_{i=1}^n w_i(x) = 1$)

$$\mathbb{E}((\hat{f}(x) - f^*(x))^2 | x_1, \dots, x_n) = \left(\sum_{i=1}^n w_i(x) [f^*(x) - f^*(x_i)] \right)^2 + \sum_{i=1}^n w_i^2(x) \mathbb{E}((y_i - \mathbb{E}(y_i|x_i))^2 | x_i)$$

$\leq B d(x, x_i)$ $\leq \sigma^2$

$$\leq \left(\sum_{i=1}^n w_i(x) B d(x, x_i) \right)^2$$

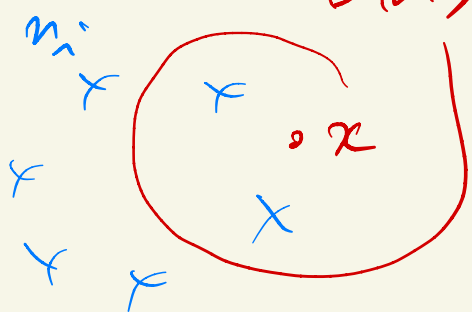
$$\leq B^2 \sum_{i=1}^n w_i(x) d(x, x_i)^2$$

bias

$$+ \underbrace{\sigma^2 \sum_{i=1}^n w_i^2(x)}_{\text{variance}}$$

$\sum_{i=1}^n w_i(x) = 1$

$$NW = \frac{1}{h}$$



$$\sigma^2 \left[\sum_{i=1}^n \left(w_i(x) - \frac{1}{n} \right)^2 - n \times \frac{1}{n^2} + \frac{2}{n} \right]$$

measure of non uniformity

Task 1: compute variance for k-NN $\Rightarrow \boxed{\sigma^2/h}$

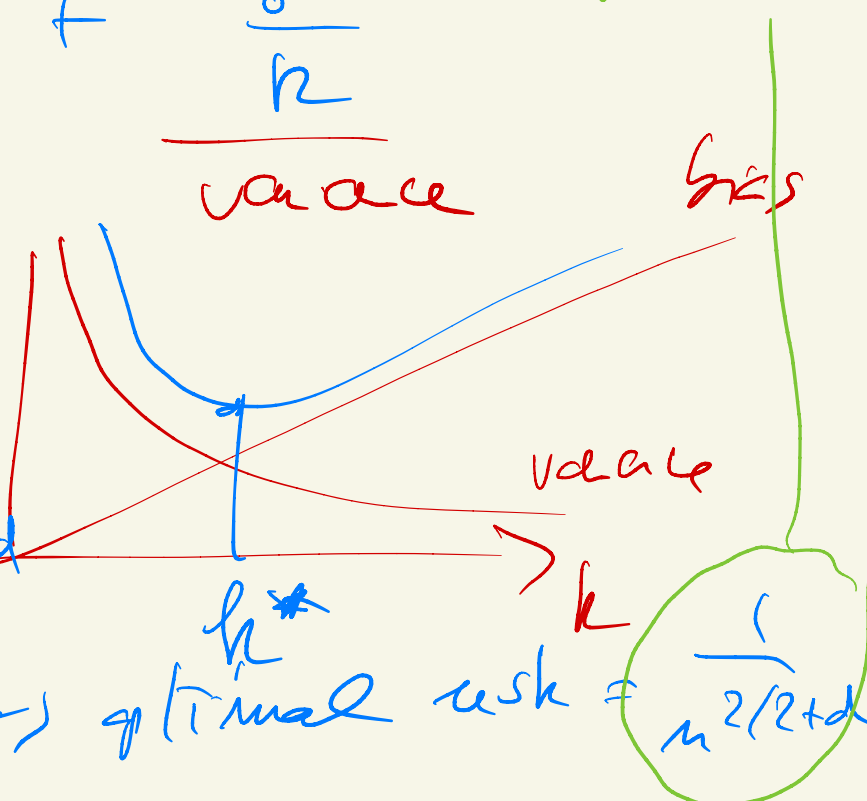
Task 2: $\int \sum_{i=1}^m w_i(x) d(x, x_i)^2 p(x) dx$ for k-NN
 $\leq \text{diam}(Z)^2 \left(\frac{h}{n}\right)^{2/d}$ if $d \geq 2$
 $\frac{h}{n}$ if $d=1$

$w_i(x) = \frac{1}{h}$ for $x_i \rightarrow x_i$
 0 otherwise

Trade-off: $B^2 \text{diam}(Z)^2 \left(\frac{h}{n}\right)^{2/d} + \frac{\sigma^2}{h}$ optimal

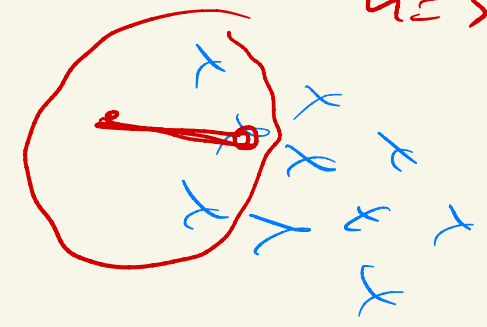
bias variance bias

"optimal" choice of $h \Rightarrow \left(\frac{h}{n}\right)^{2/d} = \frac{1}{h}$
 $h^{\frac{2}{d}+1} = n^{2/d} \Rightarrow h = \frac{n^{2/d}}{d} = \frac{2}{2+d}$
 $d n^{2/2+d} \Rightarrow$ optimal risk = $\frac{1}{n^{2/2+d}}$



$\frac{1}{n^{2/2+d}}$

$h=3$



goal: $\int p(x) \sum_{i=1}^m w_i(x) d(x, x_i)^2$

$\leq d(x, h-h \text{ nearest neighbors})$

$m+1$ points: x_1, \dots, x_m, x_{m+1}

training data test point

Lemma: let x_1, \dots, x_m, x_{m+1} $m+1$ (i.i.d) points from $p(x)$

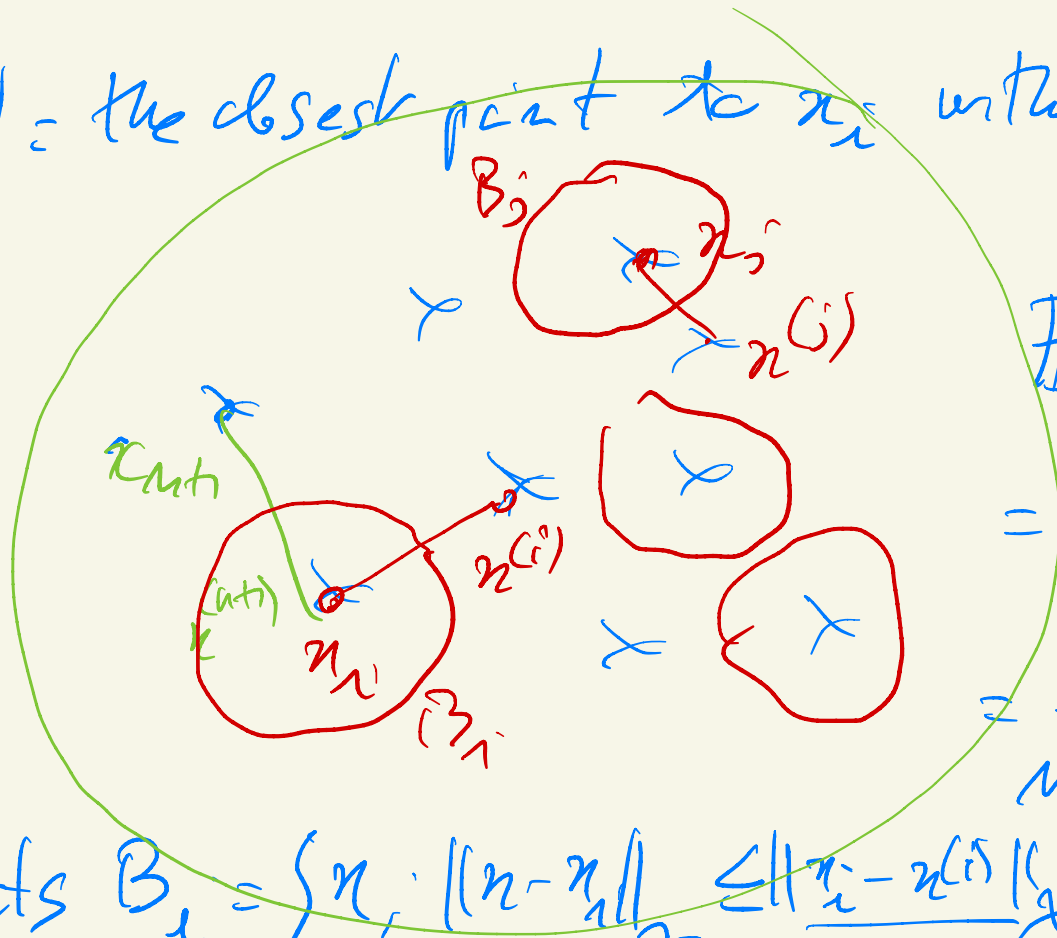
then $E d(x_{m+1}, \text{1st nearest neighbor within } (x_1, \dots, x_m)) \leq \frac{1}{m^{2/d}}$ for $d \geq 2$

Proof in \mathbb{R}^d with $d(x, x') = \|x - x'\|_\infty$

Notation: $x^{(i)}$ the nearest-neighbor of x_i in x_1, \dots, x_{m+1}

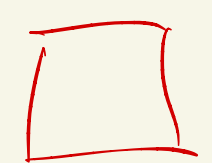
$(x_i$

$x^{(i)}$ = the closest point to x_i within $n_1, \dots, n_{i-1}, n_{i+1}, \dots, n_{n+1}$



$$\begin{aligned} & \mathbb{E} \|x_{n+1} - n^{(n+1)}\|_\infty^2 \\ &= \mathbb{E} \|n_k - n^{(k)}\|_\infty^2 \quad \forall k \\ &= \frac{1}{n+1} \sum_{k=1}^{n+1} \mathbb{E} \|n_k - n^{(k)}\|_\infty^2 \end{aligned}$$

Sets $B_i = \left\{ n_i : \|n - n_i\|_\infty \leq \frac{\|n_i - n^{(i)}\|_\infty}{2} \right\}$



The balls are disjoint

$$\sum_i \text{volume}(B_i) \leq \text{volume}(Z) \quad \text{where } Z = \{a, \|a - y\|_\infty \leq \text{diam}(X)\}$$

$$\frac{1}{n+1} \sum_i \left[\frac{\|n_i - n^{(i)}\|_\infty}{2} \right]^d \leq \left(2 \text{diam}(X) \right)^{\frac{1}{n+1} d} \text{diam}(Z) \leq 2 \text{diam}(X)$$

Then, with Jensen's inequality

$$\sum_{i=1}^{m+1} \frac{\|x_i - x^{(i)}\|^d}{2^d} \leq 2^d \text{diam}(X)^d$$

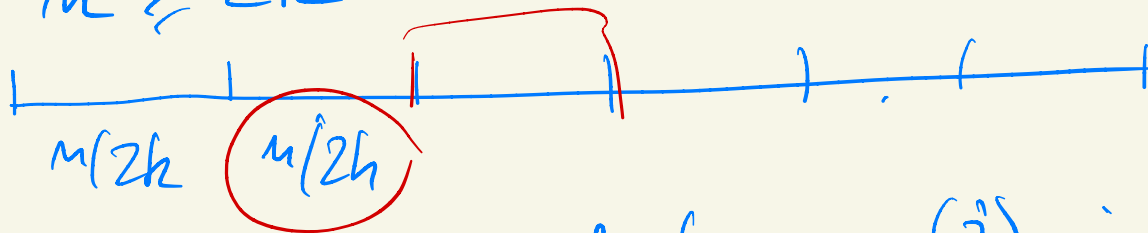
Jensen: $\frac{1}{m+1} \sum_{i=1}^{m+1} \left(\|x_i - x^{(i)}\|^2 \right)^{d/2} \geq \left(\frac{1}{m+1} \sum_{i=1}^{m+1} \|x_i - x^{(i)}\|^2 \right)^{d/2}$

$$\frac{1}{m+1} \sum_{i=1}^{m+1} \|x_i - x^{(i)}\|^2 \leq \left(\frac{4^d \text{diam}(X)^d}{m+1} \right)^{2/d}$$

$$\leq \frac{4}{(m+1)^{2/d}} \text{diam}(X)^2$$

Lemma: the average squared distance to the h -wall is less than $\frac{8}{n} \left(\frac{h}{n}\right)^{2/d} \text{diam}(\mathcal{X})^2$ if $d \geq 2$

Proof: $n \geq 2h$



$\rightarrow 2h$ pieces, notation: $n^{(i)}$ is the nearest neighbor of n_{n+1} within block i

$$\|n_{n+1} - h\text{-wall}\|^2 \leq \frac{2h}{n} \sum_{i=1}^{n/(2h)} \|n_{n+1} - n^{(i)}\|^2$$

$$\leq \frac{1}{n/(2h)} \sum_{i=1}^{n/(2h)} \frac{8 \text{diam}(\mathcal{X})^2}{(n/(2h))^{2/d}} \leq \frac{(2h)^{2/d}}{n} 8 \text{diam}(\mathcal{X})^2$$

Next weeks = Empirical risk minimization

