

SGD: Objective $F(\theta) = \mathbb{E}[l(y, f_{\theta}(x))]$ Test error
alternative = $\frac{1}{n} \sum_{i=1}^n l(y_i, f_{\theta}(x_i))$ Training error

Algorithm: $\theta_t = \theta_{t-1} - \gamma g_t \in \mathbb{R}^d$

① Unbiased: $\mathbb{E}[g_t | \theta_{t-1}] = F'(\theta_{t-1})$

(ex: $g_t = \nabla_{\theta} l(y_t, f_{\theta}(x_t)) |_{\theta = \theta_{t-1}}$)
 $t \leq n$ ② $\|g_t\|^2 \leq B^2$ bounded.

Goal: bound $\mathbb{E}[F(\theta_t) - F(\theta_*)]$

minimize of $F(\theta)$

$$\begin{aligned} \|\theta_t - \theta_*\|_2^2 &= \|\theta_{t-1} - \theta_* - \gamma g_t\|_2^2 && \text{because } \theta_t = \theta_{t-1} - \gamma g_t \\ &= \|\theta_{t-1} - \theta_*\|_2^2 - 2\gamma g_t^\top (\theta_{t-1} - \theta_*) + \underbrace{\gamma^2 \|g_t\|_2^2}_{\leq \gamma^2 B^2} \end{aligned}$$

$$\mathbb{E}[\|\theta_t - \theta_*\|_2^2 | \mathcal{G}_{t-1}] \leq \|\theta_{t-1} - \theta_*\|_2^2 - 2\gamma \underbrace{F'(\theta_{t-1})^\top}_{\text{green}} (\theta_{t-1} - \theta_*) + \gamma^2 B^2$$

by convexity $\geq F(\theta_{t-1}) - F(\theta_*)$

towering law of expectation

$$\leq \|\theta_{t-1} - \theta_*\|_2^2 - 2\gamma [F(\theta_{t-1}) - F(\theta_*)] + \gamma^2 B^2$$

$$\mathbb{E}[\|\theta_t - \theta_*\|_2^2] \leq \mathbb{E}[\|\theta_{t-1} - \theta_*\|_2^2] - 2\gamma [\mathbb{E}F(\theta_{t-1}) - F(\theta_*)] + \gamma^2 B^2$$

$$\mathbb{E}[F(\theta_{t-1}) - F(\theta_*)] \leq \frac{\mathbb{E}[\|\theta_{t-1} - \theta_*\|_2^2] - \mathbb{E}[\|\theta_t - \theta_*\|_2^2]}{2\gamma} + \frac{\gamma^2 B^2}{2\gamma}$$

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}[F(\theta_{t-1}) - F(\theta_*)] \leq \frac{1}{2\gamma T} \left[\mathbb{E}[\|\theta_0 - \theta_*\|_2^2] - \underbrace{\mathbb{E}[\|\theta_T - \theta_*\|_2^2]}_{\geq 0} \right] + \frac{\gamma B^2}{2}$$

Jensen for finite sum

$$\mathbb{E}\left[F\left(\frac{1}{T} \sum_{t=1}^T \theta_{t-1}\right)\right] - F(\theta_*) \leq \frac{1}{2\gamma T} \|\theta_0 - \theta_*\|_2^2 + \frac{\gamma B^2}{2}$$

$\leq \frac{B}{\sqrt{T}}$ $\sigma \propto \frac{1}{\sqrt{T}}$

Smoothness vs Lipschitz - continuity of $F: \mathbb{R} \rightarrow \mathbb{R}$

$$|F'(a) - F'(b)| \leq L|a - b|$$

\Downarrow

$$|F''(a)| \leq L$$

γ constant is allowed

$$|F(a) - F(b)| \leq B|a - b| \quad \forall a, b$$

\Downarrow

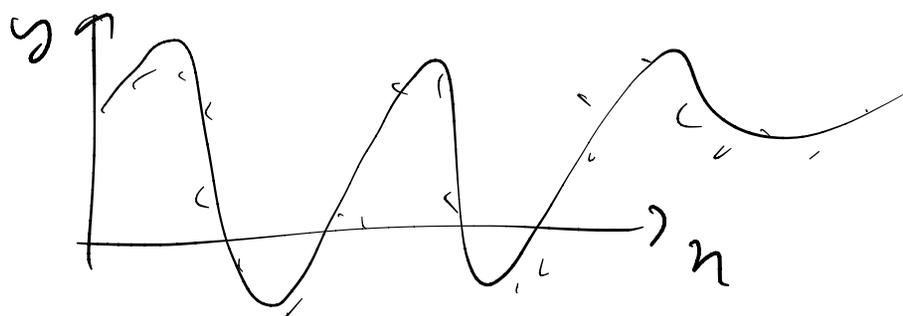
$$|F'(a)| \leq B \quad \forall a$$

γ has to decay with ϵ .

Linear model $\hat{y}_a = \mathbf{Q}^T \varphi(x)$ $\varphi: X \rightarrow \mathbb{R}^d$

Estimation error (Bias/variance)

Optimization (Convex optimization)



LOCAL AVERAGING METHODS

Training data $(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}$, $i = 1, \dots, n$

Expected risk $R(f) = \mathbb{E}[e(y, f(x))]$

Testing error

Optimal predictor) $f_*(x) = \arg \min_{z \in \mathcal{Y}} \mathbb{E}[e(y, z) | x]$

Bayes

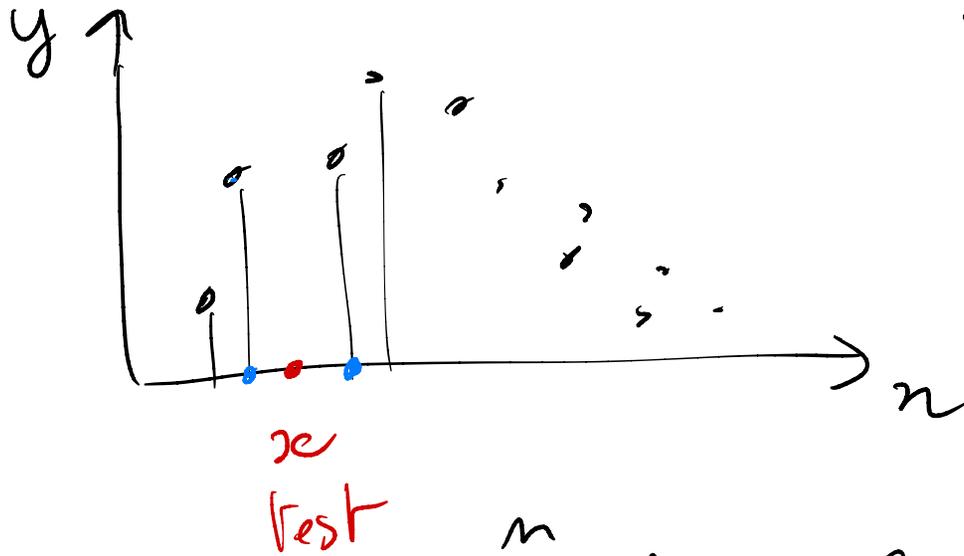
key examples: regression with square loss
 $f_*(x) = \mathbb{E}[y | x]$

classification with 0-1 loss
 $f_*(x) = \arg \max_h P(y = h | x)$

convex surrogates for binary classification
 $e(y, f(x)) = \log(1 + e^{-y f(x)})$
 $y \in \{-1, 1\}$
logistic regression

$$f_{\hat{z}}(x) = \arg \min_{z \in \mathcal{Y}} E[l(y, z) | x] = \int l(y, z) d p(y | x).$$

Local averaging: estimate directly $p(y | x)$



$$\delta y \quad \hat{p}(y | x) \quad \left(\int \delta y_i y = y_i \right)$$

① Regression

$$\hat{p}(x) = \hat{E}(y | x)$$

$$= \sum_{i \rightarrow} \hat{w}_i(x) y_i$$

Majority vote

② Classification

$$\hat{p}(x) = \arg \max_k \hat{p}(y = k | x)$$

$$\hat{p}(y | x) = \sum_{i=1}^n \hat{w}_i(x) \delta_{y_i}$$

weight functions $\forall n$

Dirac at y_i

$$\sum_{i \rightarrow} \hat{w}_i(x) \mathbb{1}_{y_i = k}$$

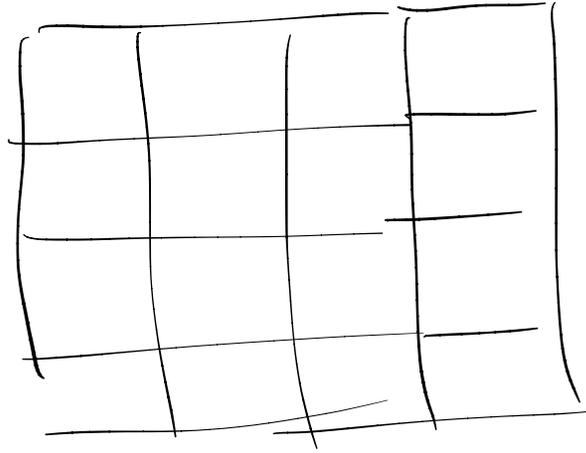
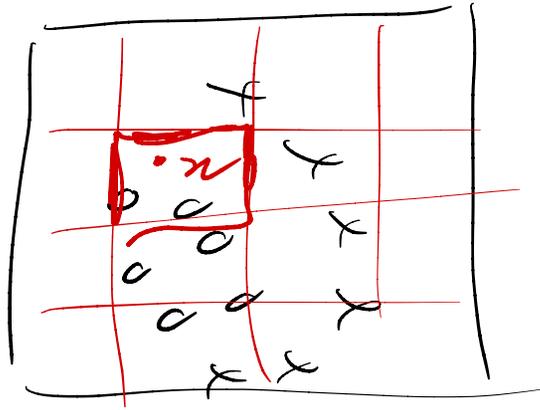
$$\hat{w}_i(x) \geq 0$$

$$\sum_{i \rightarrow} \hat{w}_i(x) = 1$$

Estimator

$$\hat{p}(x) = \arg \min_{z \in \mathcal{Y}} \int l(y, z) d \hat{p}(y | x)$$

Example 1, partition estimator

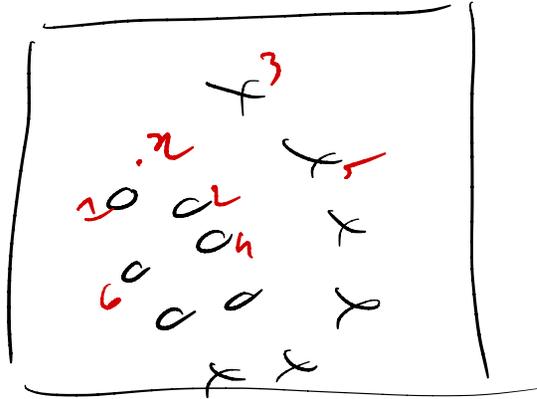


$$X = \mathbb{R}^2$$

$$Y = \{0, x\}$$

NB: simplified decision tree

Example 2: h -nearest neighbors



$$X = \mathbb{R}^2$$

$$y = \{0, x\}$$

$x_1(x)$ nearest-neighbor

$x_h(x)$ h -nearest neighbor

$$\Delta(x, x_{\underline{i_1}(x)}) \leq \Delta(x, x_{\underline{i_2}(x)}) \leq \dots$$

L_1 distance over X .

Running-time complexity

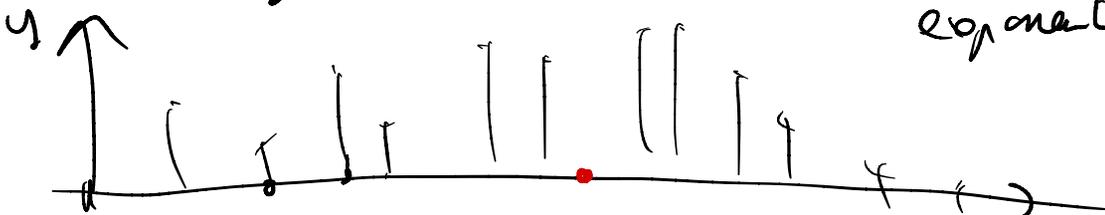
n points $x_i \in \mathbb{R}^d$

$O(nd)$ per test point (time)

$O(nd)$ memory \Rightarrow store the data

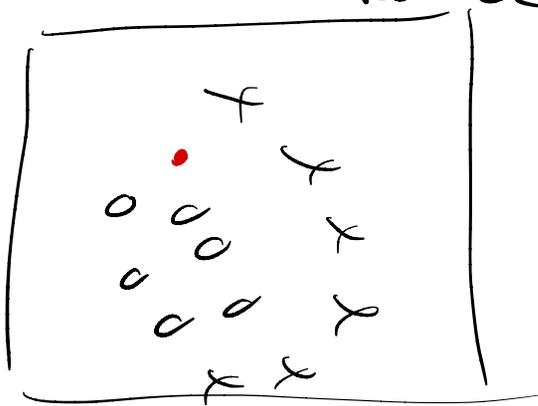
"indexing" technique

$\log(n)$ complexity
exponential in d



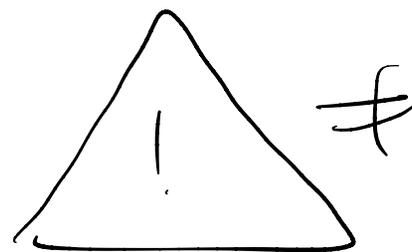
$$w_i(x) = \begin{cases} \frac{1}{h} & \text{if } x \in \{x_1(x), x_2(x), \dots, x_h(x)\} \\ 0 & \text{otherwise} \end{cases}$$

Example 3: Nadaraya-Watson estimator
kernel regression



$$X = \mathbb{R}^2$$

$$Y = \{0, x\}$$



≠ kernel methods

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

$h: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}_+$

$$k(x, x') = \exp\left(-\frac{\|x - x'\|_2^2}{h^2}\right)$$

10^4	28
10^4	48

Estimator = $\hat{f}(x) = \sum_{i=1}^n \hat{w}_i(x) y_i$ (Regression) | Criterion: excess expected risk
 $R(\hat{f}) - R(\beta_*) = \mathbb{E} | \hat{f}(x) - \beta_*(x) |^2$

Assumption: (1) banded noise:
 $|y - \beta_*(x)| \leq \sigma$ almost surely

(2) Regularity = β_* is Lipschitz-continuous
with constant B w.r.t distance Δ

$$\hat{\beta}(x) - \beta_0(x) = \sum_{i \rightarrow}^n \hat{w}_i(x) y_i - \left(\sum_{i \rightarrow}^n \hat{w}_i(x) \right) \beta_0(x) \quad \text{for } x \in \mathcal{X} \text{ a test point}$$

$$= \sum_{i \rightarrow}^n \hat{w}_i(x) (y_i - \beta_0(x)) + \sum_{i \rightarrow}^n \hat{w}_i(x) (\beta_0(x_i) - \beta_0(x))$$

$$\mathbb{E} \left[(\hat{\beta}(x) - \beta_0(x))^2 \mid x_1, \dots, x_n \right] = \underbrace{\left(\mathbb{E}[\hat{\beta}(x) - \beta_0(x) \mid x_1, \dots, x_n] \right)^2}_{\text{weights sum to 1}} + \text{var}(\hat{\beta}(x) \mid x_1, \dots, x_n)$$

$$= \left(\sum_i \hat{w}_i(x) (\beta_0(x_i) - \beta_0(x)) \right)^2 + \text{var} \left(\sum_{i \rightarrow}^n \hat{w}_i(x) (y_i - \beta_0(x)) \mid x_1, \dots, x_n \right)$$

$$= \underbrace{\left(\sum_i \hat{w}_i(x) (\beta_0(x_i) - \beta_0(x)) \right)^2}_{\text{bias}} + \underbrace{\sum_{i \rightarrow}^n \hat{w}_i(x)^2 \text{var}(y_i - \beta_0(x) \mid x)}_{\text{variance}}$$

$$\leq \left(\sum_i \hat{w}_i(x) |\beta_0(x_i) - \beta_0(x)| \right)^2 + \sigma^2 \sum_{i \rightarrow}^n \hat{w}_i(x)^2$$

$\leq B \Delta(x; n)$

$$\leq \sum_i \hat{w}_i(x) B^2 \Delta(x; n)^2 + \sigma^2 \sum_{i \rightarrow}^n \hat{w}_i(x)^2 \quad \text{by Jensen's inequality}$$

$$\mathbb{E}(|\hat{f}(x) - \hat{f}(n)|^2) \leq B^2 \sum_{i=1}^n \mathbb{E}[\hat{w}_i(x) \Delta(x, x_i)^2] + \sigma^2 \sum_{i=1}^n \mathbb{E}[\hat{w}_i(x)^2]$$

expected test error at $x \in \mathcal{X}$

$$\int \mathbb{E}[(\hat{f}(x) - \hat{f}(n))^2] d\mathcal{P}(x) \leq B^2 \sum_{i=1}^n \int \mathbb{E}[\hat{w}_i(x) \Delta(x, x_i)^2] d\mathcal{P}(x) + \sigma^2 \sum_{i=1}^n \int \mathbb{E}[\hat{w}_i(x)^2] d\mathcal{P}(x)$$

↑ training data ↑ test point

excess risk

minimized
for $w_i(x)$ very localized.

$$\boxed{R_{nn}} = \sum_{i=1}^n \hat{w}_i(x)^2 = h \cdot \frac{1}{h^2} = \frac{1}{h}$$

\Rightarrow variance term = σ^2/h

↓

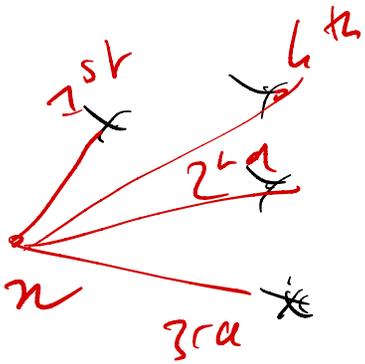
$$\sum_{i=1}^n \hat{w}_i(x)^2$$

$$= \sum_{i=1}^n \left(\hat{w}_i(x) - \frac{1}{n} \right)^2 + \frac{1}{n}$$

minimized for uniform weights

maximized for all weights on single x_i

goal: bound $\mathbb{E} \int \sum_{i=1}^n \hat{w}_i(x) \Delta(x, x_i)^2 d\rho(x)$



$$\mathbb{E} \int \frac{1}{k} \sum_{j=1}^k \Delta^2(x, x_{i_j}(x)) d\rho(x)$$

j -th nearest neighbor

Lemma: $X \subset \mathbb{R}^d$, $\Delta = \|\cdot\|_2$ distance, X bounded
 x_1, \dots, x_{n+1} sampled iid from a distribution

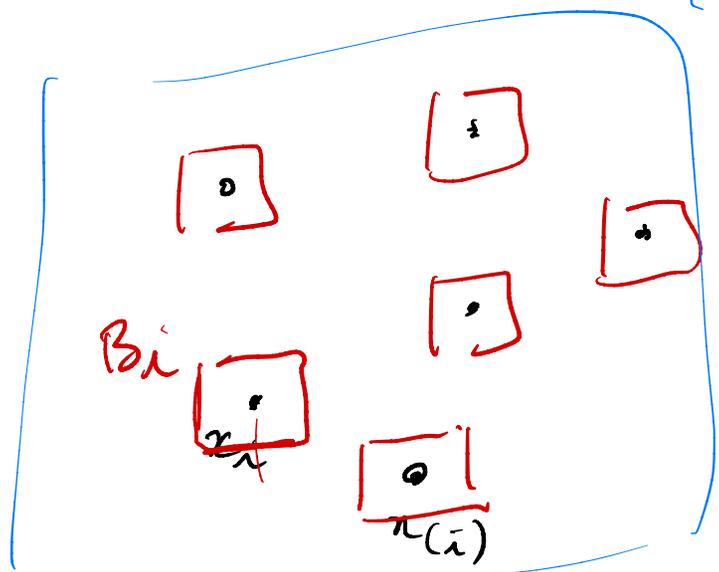
Expected square distance between x_{n+1} and
its nearest neighbor ($k=1$) $\leq \frac{4}{n^{2/d}} \text{diam}(X)^2$
(among x_1, \dots, x_n)

$\text{diam}(X) = \sup_{x, x' \in X} \Delta(x, x')$ (if $d \geq 2$)

$\{x, \|x - a\|_2 \leq r\} \Rightarrow \text{diameter: } 2r$

$k > 1: 8 \text{diam}(X)^2 \left(\frac{2k}{n}\right)^{2/d}$

Prob: $x_{(i)}$ the nearest neighbor of x_i among the n data points.



goal: $\mathbb{E} \Delta(x_{(n+1)}, x_{(n+1)})^2$
 $= \mathbb{E} \Delta(x_i, x_{(i)})^2$

$= \frac{1}{n+1} \sum_{i=1}^{n+1} \mathbb{E} [\Delta(x_i, x_{(i)})^2] = \frac{1}{n+1} \sum_{i=1}^{n+1} R_i^2$

$R_i = \Delta(x_i, x_{(i)}) > 0$

sets $B_i = \{x \in \mathbb{R}^d, \|x - x_i\|_\infty < \frac{R_i}{2}\}$

Note that $\Delta(x_i, x_j) \geq \max(R_i, R_j) \Rightarrow$ the balls B_i are disjoint

$\text{volume}(\cup B_i) \leq \text{volume}(\{x + y, x \in X, \|y\|_\infty \leq \text{diam}(X)\})$

|| because disjoint

$= \sum_i \text{vol}(B_i) = \sum_i R_i^d$

$\Rightarrow \left[\sum_i R_i^d \leq (2 \text{diam}(X))^d \right]$

good bound on $\frac{1}{n} \sum R_i^2$

we know $\sum_i R_i^d \leq (2 \operatorname{diam}(\mathcal{X}))^d$

$$\frac{1}{n} \sum_i \underbrace{(R_i^2)^{\frac{d}{2}}}_{R_i^d} \geq \left(\frac{1}{n} \sum_i R_i^2 \right)^{d/2}$$

$$\frac{d}{2} \geq 1$$

Jensen's applied to $n \mapsto n^{d/2}$

$$\begin{aligned} \Rightarrow \left(\frac{1}{n} \sum_i R_i^2 \right) &\leq \left(\frac{1}{n} \left([2 \operatorname{diam}(\mathcal{X})]^d \right)^{2/d} \right)^{2/d} \\ &\leq \frac{1}{n^{2/d}} 4 \operatorname{diam}(\mathcal{X})^2 \end{aligned}$$

Summary

$$\text{Excess risk of } k\text{-nn} \leq \underbrace{\frac{\sigma^2}{h}}_{\text{variance}} + \underbrace{B^2 \delta \text{diam}(\mathcal{X})^2 \left(\frac{2h}{m}\right)^{2/d}}_{\text{bias}}$$

variance
decreasing in h

bias
increasing in h

$$\frac{\sigma^2}{h} + B^2 \left(\frac{h}{m}\right)^{2/d}$$

$$\text{optimal balance} = \frac{1}{h} \sim \left(\frac{h}{m}\right)^{2/d} \quad (\Rightarrow) \quad h^{d/2} \sim \frac{m}{h}$$

$$(\Rightarrow) \quad h \sim m^{\frac{2}{2+d}}$$

$$h \sim m^{\frac{2}{2+d}}$$

increasing in h

$$\text{optimal value} \sim \frac{1}{m^{2/(2+d)}}$$

curse of
dimensionality