

Least-square regression

Training data $(x_i, y_i)_{i=1, \dots, n}$

$X \quad \mathbb{R}$

Square loss : $R(\theta) = \mathbb{E}[(y - f(x))^2]$ test error

$$\hat{R}(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2 \quad \text{train error}$$

Model : $f_{\theta}(x) = \varphi(x)^T \theta$ where $\varphi: X \rightarrow \mathbb{R}^d$

linear model w/ parameters

$$\theta \in \mathbb{R}^d$$

feature map

$\in \mathbb{R}^{n \times d}$ design matrix

Matrix notation : $\frac{1}{n} \sum_{i=1}^n (y_i - \varphi(x_i)^T \theta)^2 = \frac{1}{n} \|y - \Phi \theta\|_2^2$

$\in \mathbb{R}^m$ response vector

① Ordinary least-squares regression

OLS:

$$\min_{\theta} \frac{1}{n} \|y - \Phi\theta\|_2^2$$

→

excess risk

$$\frac{\sigma^2}{n}$$

② Ridge regression _____ + $\lambda \|\theta\|_2^2$

③ lower bounds

$$\hat{\theta} \in \arg \min_{\theta} \frac{1}{n} \|y - \Phi\theta\|_2^2 = F(\theta)$$

Assumption $\frac{1}{n} \Phi^T \Phi \in \mathbb{R}^{d \times d}$ empirical (non centered) covariance matrix
 $\hat{\Sigma} \approx \frac{1}{n}$
 invertible $\Rightarrow n \geq d$.

① Alg: $F'(\theta) = \frac{1}{n} \Phi^T (\underbrace{\Phi\theta - y}_{\text{residual vector}}) = 0$

$\Leftrightarrow \Phi^T \Phi \theta = \Phi^T y$ normal equations

$\Leftrightarrow \hat{\theta} = (\Phi^T \Phi)^{-1} \Phi^T y$

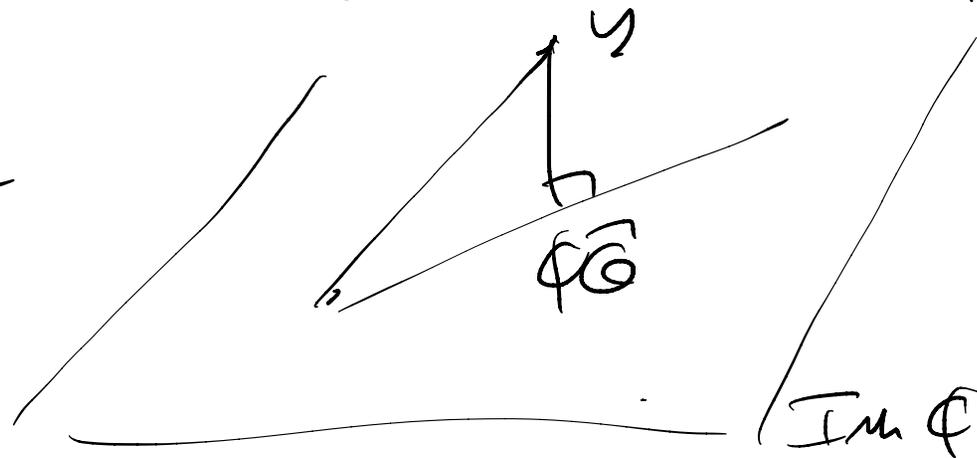
Complexity $O(d^3 + d^2 n)$

$\hookrightarrow d^2 n \cdot c$ using conjugate grad.
 $\hookrightarrow dn$ using SGM

② parametric linear regression

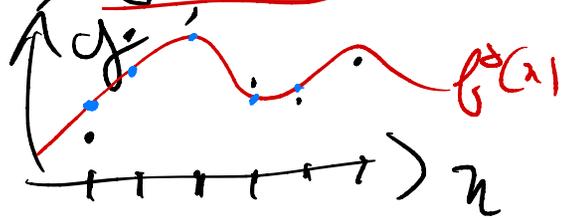
$\min_{\theta \in \mathbb{R}^d} \|y - \Phi\theta\|_2^2 \Leftrightarrow$ projection of $y \in \mathbb{R}^m$
onto the column space of Φ

$$\hat{\theta} = (\Phi^T \Phi)^{-1} \Phi^T y$$



③ Analysis \longrightarrow random design
(x_i, y_i) i.i.d from $p(x, y)$
goal: $R(\hat{\theta})$ small (test risk)

\longrightarrow fixed design: x_1, \dots, x_n are deterministic
 $y_i | x_i$ random
goal: predict well only at x_1, \dots, x_n



Assumptions (1) fixed design (x_1, \dots, x_n deterministic)

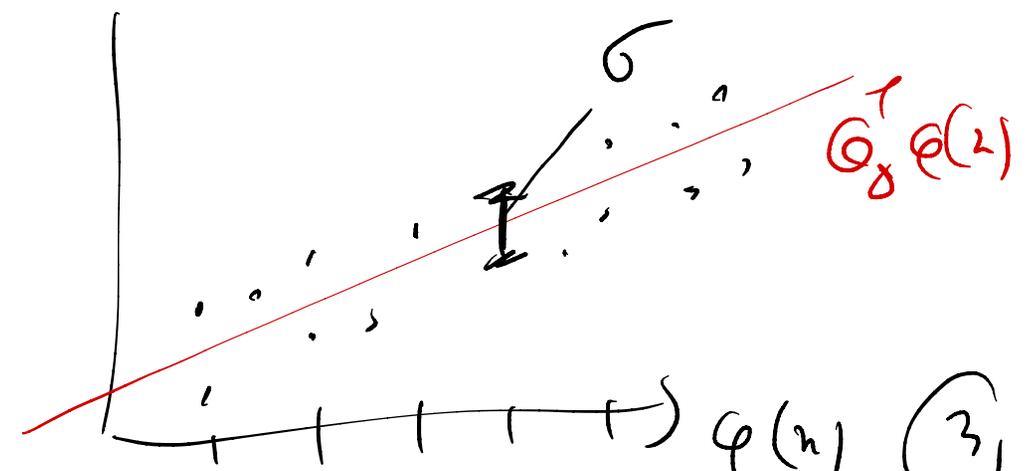
(2) assume model: $y_i = \phi(x_i)^T \Theta_* + \epsilon_i$

where $E \epsilon_i = 0$
 $E \epsilon_i^2 = \sigma^2$
 $\epsilon_1, \dots, \epsilon_n$ independent

"well-specified" model

$$E(y|x) = \phi(x)^T \Theta_* = f^*(x)$$

Bayes predictor
"target" function



~~ϵ_i Gaussian~~

(3) est. method: $\hat{\Theta} = (\Phi^T \Phi)^{-1} \Phi^T y$

Model: $y = \Phi \Theta_* + \Sigma$

"in sample" expected risk. Assume some model Θ

$$R(\theta) = \mathbb{E} \left(\frac{1}{n} \sum_{i=1}^n (y_i - \phi(x_i)^T \theta)^2 \right)$$

determinist

Testing y_i -tes distribution
 $y_i = \phi(x_i)^T \theta_* + \varepsilon_i$

$$= \frac{1}{n} \mathbb{E} \|y - \Phi \theta\|_2^2 = \frac{1}{n} \mathbb{E} \underbrace{\|\Phi \theta_* + \varepsilon - \Phi \theta\|_2^2}_{\text{model}}$$

$$= \frac{1}{n} \mathbb{E} \|\varepsilon + \Phi(\theta_* - \theta)\|_2^2$$

$$= \frac{1}{n} \mathbb{E} \|\varepsilon\|_2^2 + \frac{1}{n} \mathbb{E} \|\Phi(\theta - \theta_*)\|_2^2 + \frac{2}{n} \mathbb{E} (\varepsilon^T \Phi(\theta_* - \theta))$$

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E} \varepsilon_i^2 = \sigma^2$$

$$R(\theta) = \sigma^2 + (\theta - \theta_*)^T \frac{\Phi^T \Phi}{n} (\theta - \theta_*) = \|\theta - \theta_*\|_2^2$$

test risk

$$R(\theta) = \sigma^2 + (\theta - \theta_0)^T \frac{\Phi\Phi^T}{n} (\theta - \theta_0)$$

OLS

$$\hat{\theta} = (\Phi\Phi^T)^{-1} \Phi^T y = (\Phi\Phi^T)^{-1} \Phi^T (\Phi\theta_0 + \varepsilon)$$

model

$$y = \Phi\theta_0 + \varepsilon = \theta_0 + \underbrace{(\Phi\Phi^T)^{-1} \Phi^T \varepsilon}$$

goal: $\mathbb{E}_\varepsilon R(\hat{\theta}) = \sigma^2$

since $\mathbb{E}\varepsilon = 0 \Rightarrow$ "unbiased"

$$= \mathbb{E} \left[\left[(\Phi\Phi^T)^{-1} \Phi^T \varepsilon \right]^T \frac{\Phi\Phi^T}{n} \left[(\Phi\Phi^T)^{-1} \Phi^T \varepsilon \right] \right]$$

$$= \mathbb{E} \varepsilon^T \left[\Phi (\Phi\Phi^T)^{-1} \frac{\Phi\Phi^T}{n} (\Phi\Phi^T)^{-1} \Phi^T \right] \varepsilon$$

$$= \mathbb{E} \frac{1}{n} \text{tr} \varepsilon^T \Phi (\Phi\Phi^T)^{-1} \Phi^T \varepsilon = \mathbb{E} \text{tr} \Phi (\Phi\Phi^T)^{-1} \Phi^T \varepsilon \varepsilon^T$$

$$= \frac{\sigma^2}{n} \mathbb{E} \text{tr} \Phi (\Phi\Phi^T)^{-1} \Phi^T = \frac{\sigma^2}{n} \text{tr} \mathbb{I}_d = \frac{\sigma^2 d}{n}$$

$\mathbb{E}(\varepsilon\varepsilon^T)$ is
 $= \sigma^2 \mathbb{I}_{n \times n}$

Ridge regression: $\frac{1}{n} \|y - \phi\theta\|_2^2 + \lambda \|\theta\|_2^2$ MAP

Random design
 Low bias

$\text{Det} < 0$ $\lambda \|\theta\|_2$

$\hat{\theta} = \arg \min_{\theta} \|y - \phi\theta\|_2$

$\|y - \phi\theta\|_2$

$\text{tr} AB = \text{tr} BA$
 $\text{tr} \lambda = \lambda \quad \forall \lambda \in \mathbb{R}$

$2^5 5^8$
 $3^4 2^6$

Ridge regression: $\frac{1}{n} \|y - \Phi\theta\|_2^2 + \lambda \|\theta\|_2^2 = F(\theta)$

empirical risk

regularization

exp ()
 exp ()
 label loss

exp ($-\lambda \|\theta\|_2^2$)
 prior

Maximum A Posteriori
 estimate
 with $y | \theta \sim \text{Gaussian}$
 $\theta \sim \text{Gaussian}$

optimize

Estimator: no assumption on $\Phi \in \mathbb{R}^{n \times d}$

$$\frac{1}{2} F'(\theta) = \frac{1}{n} \Phi^T (\Phi\theta - y) + \lambda \theta = 0$$

$$\Rightarrow \hat{\theta}_\lambda = (\Phi^T \Phi + n\lambda I)^{-1} \Phi^T y$$

Estimator: $\hat{\Theta}_\lambda = (\Phi^T \Phi + n\lambda I)^{-1} \Phi^T y = (\Phi^T \Phi + n\lambda I)^{-1} \Phi^T \Theta_0 + (\Phi^T \Phi + n\lambda I)^{-1} \Phi^T \varepsilon$

Model: $y = \Phi \Theta_0 + \varepsilon$

excess risk: $\frac{1}{n} \|\Phi(\Theta - \Theta_0)\|_2^2 = R(\Theta) - \sigma^2$

deviate $\hat{\Theta}_\lambda - \Theta_0 = (\Phi^T \Phi + n\lambda I)^{-1} \Phi^T \varepsilon - n\lambda (\Phi^T \Phi + n\lambda I)^{-1} \Theta_0$

$E[R(\hat{\Theta}_\lambda) - \sigma^2] = \frac{1}{n} E \left[\underbrace{\|\Phi (\Phi^T \Phi + n\lambda I)^{-1} \Phi^T \varepsilon\|_2^2}_V + \underbrace{n\lambda^2 \|\Phi (\Phi^T \Phi + n\lambda I)^{-1} \Theta_0\|_2^2}_B \right] + 0$

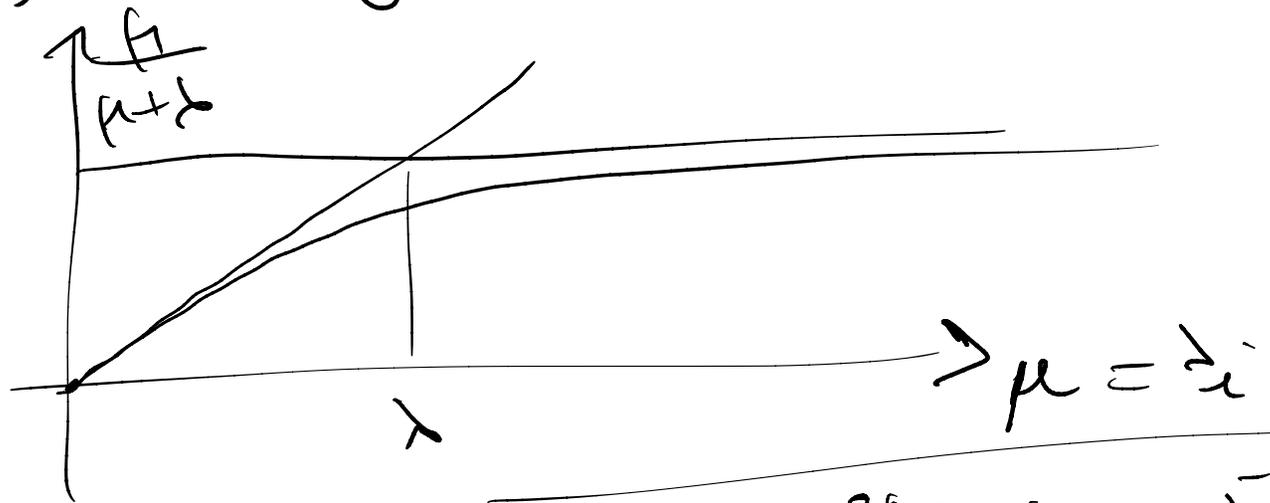
$V = \frac{1}{n} E \left[\text{tr} \left(\Phi (\Phi^T \Phi + n\lambda I)^{-1} \Phi^T \varepsilon \varepsilon^T \Phi (\Phi^T \Phi + n\lambda I)^{-1} \Phi^T \right) \right]$

$= \frac{\sigma^2}{n} \text{tr} \left(\underbrace{\Phi^T \Phi}_{n \Sigma} \right) (\Phi^T \Phi + n\lambda I)^{-2} = \frac{\sigma^2}{n} \text{tr} \hat{\Sigma}^2 (\hat{\Sigma} + \lambda I)^{-2}$ (in place of $\frac{\sigma^2}{n}$)

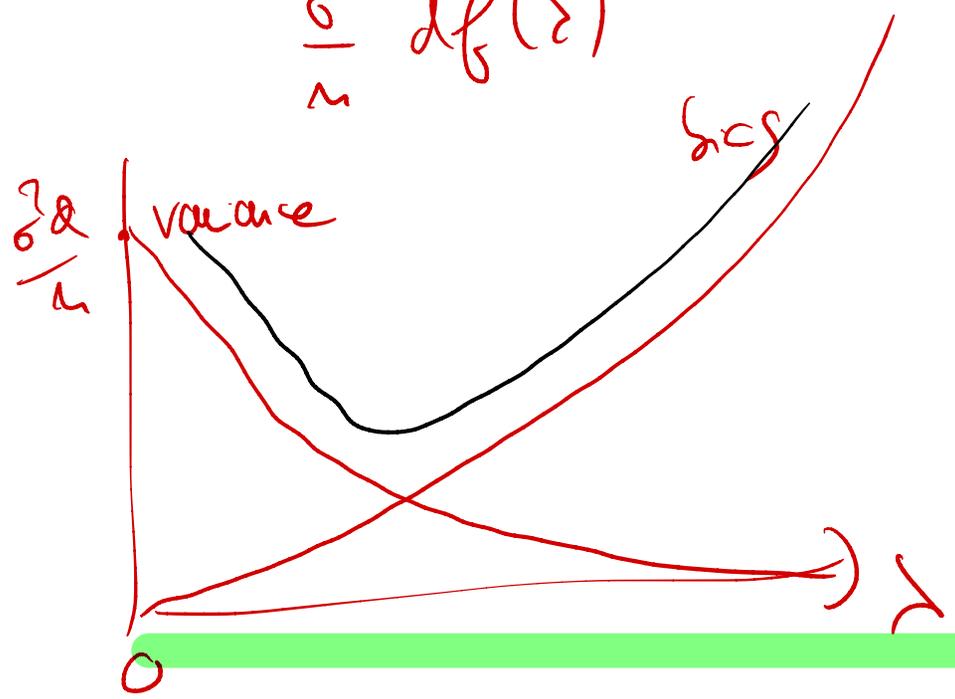
= degrees of freedom

Variance $\frac{\sigma^2}{n} \text{tr} \hat{\Sigma}^2 (\hat{\Sigma} + \lambda I)^{-2} = \sum_{i=1}^d \frac{\lambda_i^2}{(\lambda_i + \lambda)^2} \approx$ "counting" the eigenvalues above λ

λ_i are eigenvalues of $\hat{\Sigma}$



$\frac{\sigma^2}{n} df(\lambda)$



Bias = $n d^2 \|\phi (\phi \phi^T + n \lambda I)^{-1} \phi_0\|_2^2$
 $= \phi_0^T (\phi \phi^T + n \lambda I)^{-2} \phi \phi^T \phi_0 n d^2$
 $= \sum \phi_0^T (\hat{\Sigma} + \lambda I)^{-2} \hat{\Sigma} \phi_0$
 increasing in λ

$$\text{Risk} = \text{variance} + \text{Bias}^2$$

$$\frac{\sigma^2}{m} \text{tr} \hat{\Sigma}^2 (\hat{\Sigma} + \lambda I)^{-2} + \lambda^2 \mathbf{Q}_*^T \hat{\Sigma} (\hat{\Sigma} + \lambda I)^{-2} \mathbf{Q}_*$$

lemma: $\forall \mu > 0, \forall \lambda \geq 0, (\mu + \lambda)^{-2} \lambda \mu \leq \frac{1}{2} \Leftrightarrow (\mu + \lambda)^2 \geq 2\lambda\mu$

applied to $\mu = \text{eigenvalues of } \hat{\Sigma}$

$$\Rightarrow (\hat{\Sigma} + \lambda I)^{-2} \hat{\Sigma} \preceq \frac{1}{2} I$$

$$V \leq \frac{\sigma^2}{2m} \frac{\text{tr} \hat{\Sigma}}{\lambda}$$

variance

$$B \leq \frac{\lambda}{2} \|\mathbf{Q}_*\|^2$$

bias

optimal $\lambda^0 = \frac{\sigma}{\sqrt{m}} \frac{\sqrt{\text{tr} \hat{\Sigma}}}{\|\mathbf{Q}_*\|_2}$ and optimal performance

$$\frac{\sigma}{\sqrt{m}} \frac{\sqrt{\text{tr} \hat{\Sigma}}}{R} \|\mathbf{Q}_*\|_2$$

dimension independent bound $\neq \frac{\sigma^2}{m}$

Key: adding regularizer

① explicitly $\hat{R}(\beta) + \frac{\lambda}{2} \|\beta\|_2^2$
 $+ \lambda \|\beta\|_1$

$\hat{R}(\beta)$ such that $\|\beta\|_2 \leq D$

② Computational regularization.

single pass SGD

Random design = $(x_1, y_1), \dots, (x_n, y_n) \text{ IID}$

≠ fixed design

$$\hat{\theta}_{OLS} = (\Phi^T \Phi)^{-1} \Phi^T y = \sum_{i=1}^n \left(\frac{\Phi_i^T y_i}{n} \right)$$

Model: $y_i = \varphi(x_i)^T \theta_* + \varepsilon_i \Rightarrow$ exposed risk

$$\mathbb{E} \| y - \varphi(x_i)^T \theta \|^2$$

(x, y) testing data

$$= \sigma^2 + (\theta - \theta_*)^T \mathbb{E} \varphi(x) \varphi(x)^T (\theta - \theta_*)$$

$$= \sigma^2 + (\theta - \theta_*)^T \Sigma (\theta - \theta_*)$$

excess risk of OLS = $\frac{\sigma^2}{n} \mathbb{E} \left[\text{tr} \hat{\Sigma}^{-1} \Sigma \right]$

$\varphi(x)$ is gaussian $\Rightarrow \frac{\sigma^2 d}{n \left(1 - \frac{d+1}{n} \right)}$

Morning : ERM with square loss. / Statistics
Rademacher complexity

Afternoon : convex optimization / optimization
SGD \longrightarrow combination

\implies linear models
in finite dim.

Morning : kernel methods

after : neural networks