# Least-squares regression

last week: $\frac{1}{M} \sum_{i=1}^{M} \ell(y_i, f(x_i))$    empirical risk

output / input

training data

$$R(f) = E\left[ \ell(y, f(x)) \right]$$

losses
square

function classes
linear

① What is the optimal prediction?

$$f^*(x) = \underset{z}{\arg\min}\ E_{p(y|x)}\left[\ell(y, z)\atop (y-z)^2\right]$$

$$= E(y|x)$$



② Model

$$f(x) = f_\Theta(x),\quad \Theta \in$$

$$= \Theta^T \varphi(x)$$

↑ parameters

# Least-squares regression

Data $(x_i, y_i)$, $i = 1, \ldots, n$    feature function $\varphi : \mathcal{X} \to \mathbb{R}^d$

$x_i \in \mathcal{X}$   $y_i \in \mathbb{R}$

Method: $\hat{\Theta} \in \arg\min\limits_{\Theta \in \mathbb{R}^d} \left[ \dfrac{1}{n} \sum\limits_{i=1}^{n} \left( y_i - \Theta^T \varphi(x_i) \right)^2 \right] = F(\Theta)$

Notes: ① often $\varphi(x) = \begin{pmatrix} x \\ 1 \end{pmatrix}$

② Vector / matrix notation

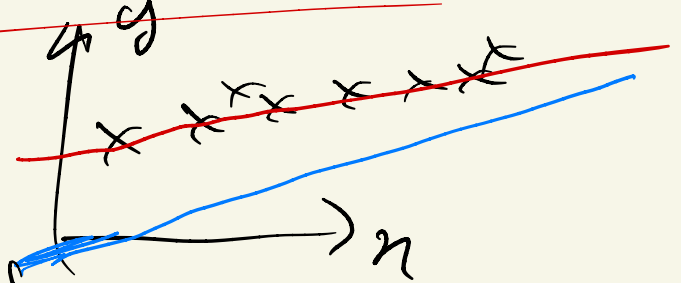$y \in \mathbb{R}^n$, $\quad \overline{\Phi} \in \mathbb{R}^{n \times d}$

$\begin{pmatrix} \varphi(x_1)^T \\ \vdots \\ \varphi(x_n)^T \end{pmatrix}$

$$\boxed{F(\Theta) = \frac{1}{n} \left\| y - \overline{\Phi}\Theta \right\|_2^2}$$ with $\|z\|_2^2 = \sum\limits_{j=1}^{d} z_j^2$ Euclidean norm

$= \dfrac{1}{n} \sum\limits_{i=1}^{n} \left( y_i - (\overline{\Phi}\Theta)_i \right)^2$

Minimizing $F(\theta) = \frac{1}{n} \| y - \phi\theta \|_2^2 \Rightarrow$ Convex $F'(\theta) = 0$

$$F'(\theta) = \frac{2}{n} \overline{\phi}^T (\phi\theta - y)$$

$+ \boxed{\gamma_2 \| \theta \|_2^2}$

$$F(\theta) = \frac{1}{n} \sum_i (y_i - \phi(x_i)^T \theta)^2 = \frac{1}{n} \sum_{i=1}^n \left( y_i - \sum_{j=1}^d \phi(x_i)_j \theta_j \right)^2$$

$$F'(\theta) = -\frac{1}{n} \sum_i 2(y_i - \phi(x_i)^T \theta) \phi(x_i) \qquad = -\frac{2}{n} \sum_i y_i \phi(x_i) + \frac{2}{n} \sum_i \phi(x_i) \phi(x_i)^T \theta$$

$$\frac{\partial F}{\partial \theta_j}(\theta) = -\frac{1}{n} \sum_i 2 \left( y_i - \sum_{j'=1}^d \phi(x_i)_{j'} \theta_{j'} \right) \phi(x_i)_j$$

$$\overline{\phi}^T \phi \in \mathbb{R}^{d \times d} \qquad \phi^T \phi = \sum_{i=1}^n \phi(x_i) \phi(x_i)^T$$

$$\phi = \begin{pmatrix} \phi(x_1)^T \\ \vdots \\ \phi(x_n)^T \end{pmatrix} \in \mathbb{R}^{n \times d}$$

$$\left. \begin{array}{c} \hat{\Sigma}_{jj} \\ = \frac{1}{n} \sum_{i=1}^n \phi(x_i)_j^2 \end{array} \right.$$

$$\hat{\Sigma} = \frac{1}{n} \phi^T \phi = \frac{1}{n} \sum_{i=1}^n \phi(x_i) \phi(x_i)^T$$

matrix of second order moment

(non centered covariance matrix)

$$F(\Theta) = \frac{1}{n} \Phi^T (\Phi\Theta - y) = 0 \quad \Longrightarrow \quad \frac{1}{n}\Phi^T\Phi \cdot \Theta = \frac{1}{n}\Phi^T y$$

<span style="color:red">(normal equations)</span>

<span style="color:red">$\in \mathbb{R}^d$</span> $\quad$ <span style="color:red">$\in \mathbb{R}^d$</span>

linear system

Assumption: $\frac{1}{n}\Phi^T\Phi$ invertible <span style="color:red">$\mathbb{R}^{d \times d}$</span>

$$\Longrightarrow \quad \text{rank}\,\Phi = d \quad \left( \begin{array}{c} \text{remember that} \\ \Phi \in \mathbb{R}^{n \times d} \end{array} \right)$$

This imposes that $n \geq d$

Solution: $\Theta = \left(\frac{1}{n}\Phi^T\Phi\right)^{-1} \frac{\Phi^T y}{n} \quad$ $\underline{1 \text{ line of code}}$

<span style="color:blue">$+ \delta I$</span>

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

<span style="color:red">$\frac{1}{n}\Phi^T\Phi = \hat{\Sigma}$</span>

<span style="color:red">dominant</span>

complexity? $O(d^3)$ only for small $d$ $\underbrace{+ O(d^2 n)}_{\text{(circled in red)}}$
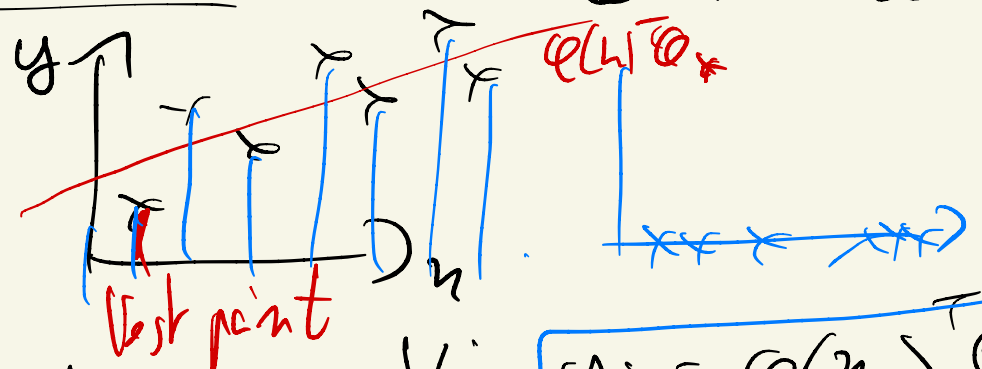
gradient descent / stochastic $\to O(dn)$
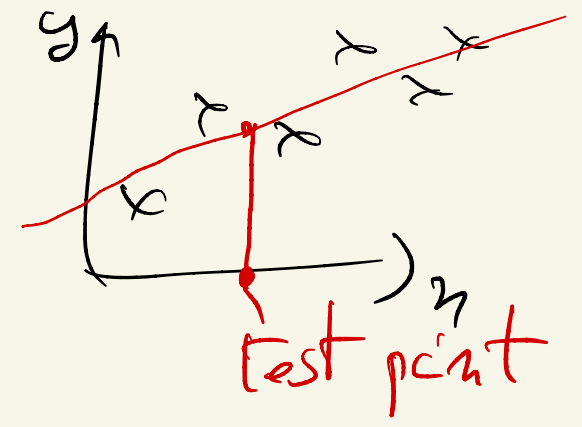
Summary: $F(\theta) = \frac{1}{n} \| y - \Phi\theta \|_2^2$ and $\boxed{\hat{\theta} = \left(\frac{1}{n}\Phi^\top\Phi\right)^{-1}\frac{1}{n}\Phi^\top y}$

$(n_i, g_i)$ sampled iid from $dp(x, y)$ $\Rightarrow$ generalization performance
Random design $= \uparrow$ goal is to minimize $R(\theta) = \mathbb{E}(y - \varphi(\theta)^\top x)^2$

Fixed design: $(x_i)$'s are deterministic

$\varphi(n_i)^\top\theta_*$

Test point

Test point

Assumptions: $\forall i,$ $\boxed{y_i = \varphi(n_i)^\top\theta_* + \varepsilon_i}$ only source of randomness

$\mathbb{E}\varepsilon_i = 0$
$\mathbb{E}\varepsilon_i^2 = \sigma^2$
noise variance independent

goal: find $\theta$ such that
$\mathbb{E}_y(F(\theta))$ is as small possible

$\boxed{y = \Phi\theta_* + \varepsilon}$
$n_0$

$$F(\theta) = \hat{R}(\theta) = \frac{1}{n} \|y - \phi\theta\|_2^2 \quad \text{with model} \quad y = \phi\theta^* + \varepsilon \in \mathbb{R}^n$$

$$\text{with} \quad \mathbb{E}\,\varepsilon = c \qquad \left(\mathbb{E}\varepsilon_i^2 = \sigma^2\right.$$

$$\boxed{\mathbb{E}(\varepsilon\varepsilon^T) \in \mathbb{R}^{n \times n}}$$
$$\| \quad \sigma^2 I$$

$$\left(\mathbb{E}(\varepsilon\varepsilon^T)\right)_{ij} = \mathbb{E}\varepsilon_i\varepsilon_j = \begin{cases} \sigma^2 & \text{if } i = j \\ c & \text{if } i \neq j \end{cases}$$

goal: $\min R(\theta) = \mathbb{E}\hat{R}(\theta)$

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

deterministic

$$= \frac{1}{n}\mathbb{E}\|\varepsilon\|_2^2 + \frac{1}{n}\mathbb{E}\|\phi(\theta^* - \theta)\|_2^2 + \frac{2}{n}\mathbb{E}\varepsilon^T\phi(\theta^* - \theta) \qquad \boxed{\|z\|_2^2 = z^Tz}$$

$$= \frac{1}{n}\sum_{i=1}^{n}\underbrace{\mathbb{E}\varepsilon_i^2}_{\sigma^2} + \frac{1}{n}\underbrace{(\theta^* - \theta)^T\phi^T\phi(\theta^* - \theta)}_{} + c$$

minimized for $\theta = \theta^*$

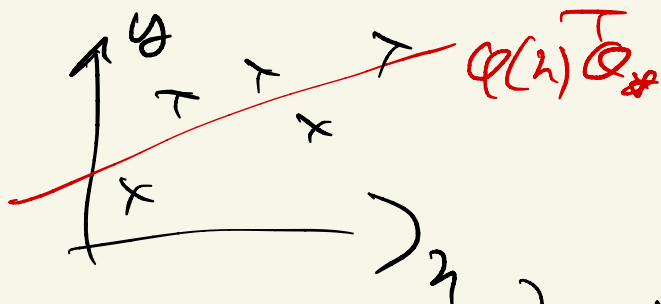$$\boxed{R(\theta) = \sigma^2 + (\theta^* - \theta)^T\hat{\Sigma}(\theta^* - \theta)}$$
$$R^* $$

$$\Rightarrow \boxed{R(\theta) - R^* = (\theta^* - \theta)^T\hat{\Sigma}(\theta^* - \theta)} \quad \text{excess risk}$$

different cost: $\|\theta^* - \theta\|_2^2$ not used

$$R(\Theta) - R^* = (\Theta - \Theta_*)^T \hat{\Sigma} (\Theta - \Theta_*) \quad \text{with model } y = \phi \Theta_* + \varepsilon$$

$$\hat{\Theta} = \frac{\hat{\Sigma}^{-1} \phi^T y}{m}$$

(normal equations)



$$\phi(x)^T \Theta_*$$

$$\mathbb{E}\varepsilon = 0$$
$$\mathbb{E}\varepsilon\varepsilon^T = \sigma^2 I$$

$$\mathbb{E}(R(\hat{\Theta}) - R_*) = \mathbb{E}\left[ (\hat{\Theta} - \Theta_*)^T \hat{\Sigma} (\hat{\Theta} - \Theta_*) \right]$$

$$\text{with } \hat{\Theta} = \frac{\hat{\Sigma}^{-1} \phi^T y}{m}$$
$$= \hat{\Sigma}^{-1} \frac{\phi^T (\phi \Theta_* + \varepsilon)}{m}$$

Consequence) : $\mathbb{E}\hat{\Theta} = \Theta_*$ unbiased

bias : $\mathbb{E}\hat{\Theta} - \Theta_*$

$$\hat{\Theta} = \Theta_* + \frac{\hat{\Sigma}^{-1} \phi^T \varepsilon}{m}$$



$$\hat{\Theta}$$

Covariance matrix $\Theta_*$

$$\mathbb{E}(\hat{\Theta} - \Theta_*)(\hat{\Theta} - \Theta_*)^T = \mathbb{E}\left[ \hat{\Sigma}^{-1} \phi^T \frac{\varepsilon\varepsilon^T}{} \phi \hat{\Sigma}^{-1} \right] \quad \sigma^2 I$$

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

$$= \frac{\sigma^2}{m^2} \hat{\Sigma}^{-1} \phi^T \phi \hat{\Sigma}^{-1} = \boxed{\frac{\sigma^2}{m} \hat{\Sigma}^{-1}}$$

Consequences

$$\mathbb{E}\, R(\hat{\Theta}) - R^*$$
$$= \text{tr}\left( \hat{\Sigma}\, \mathbb{E}(\hat{\Theta} - \Theta_*)(\hat{\Theta} - \Theta_*)^T \right)$$
$$= \text{tr}\, \hat{\Sigma} \cdot \frac{\sigma^2}{m} \hat{\Sigma}^{-1}$$
$$= \text{tr}\, \frac{\sigma^2}{m} I = \boxed{\frac{\sigma^2 d}{m}}$$

Excess risk = $\dfrac{\sigma^2 d}{n}$ — dimension

$R(\hat{\theta}) - R^* \leq$ — pessimistic

- it is an equality for any $n$
- it is only for fixed design $\Rightarrow$ see book for analysis for random design
- it is "optimal" $\Rightarrow$

  if $y = \varphi(x)^T \theta_* + \varepsilon$ $\Rightarrow$ see precise statements in book

- it is disappointing $\Rightarrow$ need for regularization
- Why random design has err?

Regularization: replace $F(\theta) = \frac{1}{n}\|y - \phi\theta\|_2^2$ $\quad \frac{\lambda}{2}\theta\theta^T$

$\text{by} \quad \boxed{F(\theta) + \frac{\lambda}{2}\|\theta\|_2^2} \Rightarrow \text{ridge regression}$

$\underbrace{\phantom{F_\lambda(\theta)}}_{F_\lambda(\theta)} \quad + \lambda\|\theta\|_1 \Rightarrow \text{Lasso}$

gradient: $F_\lambda'(\theta) = F'(\theta) + \lambda\theta$

$\qquad = \frac{1}{n}\phi^T(\phi\theta - y) + \lambda\theta = (\hat\Sigma + \lambda I)\theta - \frac{1}{n}\phi^T y$

$\qquad\qquad \underbrace{\frac{1}{n}\phi^T\phi}_{\hat\Sigma}$

normal equations: $(\hat\Sigma + \lambda I)\theta = \frac{1}{n}\phi^T y$. Always unique solution when $\lambda > 0$. $\quad \hat\theta_\lambda = (\hat\Sigma + \lambda I)^{-1}\frac{\phi^T y}{n}$

goal: compute bias and variance of excess risk.

$\hat\theta_\lambda = (\hat\Sigma + \lambda I)^{-1}\frac{\phi^T y}{n} = (\hat\Sigma + \lambda I)^{-1}\frac{\phi^T}{n}\underbrace{(\phi\theta_* + \varepsilon)}_{\text{model}} = (\hat\Sigma + \lambda I)^{-1}\hat\Sigma\theta_* + (\hat\Sigma + \lambda I)^{-1}\frac{\phi^T\varepsilon}{n}$

Bias: $\mathbb{E}\hat\theta_\lambda = (\hat\Sigma + \lambda I)^{-1}\hat\Sigma\theta_* = (\hat\Sigma + \lambda I)^{-1}(\hat\Sigma + \lambda I - \lambda I)\theta_*$

$\qquad\qquad = \theta_* - \lambda(\hat\Sigma + \lambda I)^{-1}\theta_*$

$$\boxed{\mathbb{E}\hat\theta_\lambda - \theta_* = -\lambda(\hat\Sigma + \lambda I)^{-1}\theta_*}$$

$$\hat{Q}_\delta = (\hat{\Sigma}+\lambda I)^{-1}\hat{\Sigma}\hat{Q}_* + (\hat{\Sigma}+\lambda I)^{-1}\frac{1}{n}\Phi^T \mathcal{E}$$

<span style="color:red">$\mathbb{E}\hat{Q}_\lambda$</span>

<span style="color:red">Covariance matrix</span>

$$\text{var}(\hat{Q}_\delta) = \mathbb{E}\left[\left(\hat{Q}_\delta - \mathbb{E}\hat{Q}_\delta\right)\left(\hat{Q}_\delta - \mathbb{E}\hat{Q}_\delta\right)^T\right] = \mathbb{E}\left[(\hat{\Sigma}+\lambda I)^{-1}\Phi^T\underbrace{\mathcal{E}\mathcal{E}^T}_{\sigma^2 I}\Phi(\hat{\Sigma}+\lambda I)^{-1}\right]\Big/ n^2$$

$$= \frac{(\hat{\Sigma}+\lambda I)^{-1}\hat{\Sigma}(\hat{\Sigma}+\lambda I)^{-1}}{n} = \boxed{\frac{\sigma^2\hat{\Sigma}(\hat{\Sigma}+\lambda I)^{-2}}{n}}$$

<span style="color:blue">$\Rightarrow$ when $\lambda = 0$</span>

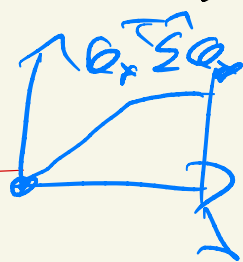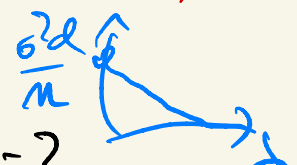<span style="color:blue">$\dfrac{\hat{\Sigma}^{-1}}{n}$</span>

**Main result:**

$$\mathbb{E} R(\hat{Q}_\delta) - R_* = \mathbb{E}(\hat{Q}_\delta - Q_*)^T\hat{\Sigma}(\hat{Q}_\delta - Q_*)$$

<span style="color:red">Lemma:</span>

$$\mathbb{E}(Z-a)^T M(Z-a) = (\mathbb{E}Z-a)^T M(\mathbb{E}Z-a) + \text{tr } M\,\text{var}(Z)$$

$$= \text{tr}\,\frac{\hat{\Sigma}^2(\hat{\Sigma}+\lambda I)^{-2}}{n}\sigma^2 + \left(-\lambda(\hat{\Sigma}+\lambda I)^{-1}Q_*\right)^T\hat{\Sigma}\left(-\lambda(\hat{\Sigma}+\lambda I)^{-1}Q_*\right)$$

$$\lambda^2 Q_*^T(\hat{\Sigma}+\lambda I)^{-2}\hat{\Sigma}Q_*$$

<span style="color:red">variance term</span>

<span style="color:red">bias term</span>

goal: getting upper bounds and optimize over $\lambda$

variance: $\frac{\sigma^2}{n} \, \text{tr} \, \hat{\Sigma}^2 (\hat{\Sigma} + \lambda I)^{-2}$    bias $\lambda^2 \, \theta_*^T \hat{\Sigma} (\hat{\Sigma} + \lambda I)^{-2} \theta_*$

$\text{tr} \, \hat{\Sigma} \times \text{matrix}$

requirement : $\hat{\Sigma}$ may not be invertible

main tool : $\underbrace{\hat{\Sigma}(\hat{\Sigma} + \lambda I)^{-2}}_{\text{matrix with eigenvalues } \mu(\mu+\lambda)^{-2}} \preceq \frac{1}{2}\lambda^{-1} I$  where $\mu$ is an eigenvalue of $\hat{\Sigma}$.

Lemma : $\mu(\mu+\lambda)^{-2} \le \frac{1}{2\lambda} \quad \forall \lambda, \mu \ge 0$

$\implies (\mu+\lambda)^{-2} \le \frac{1}{2\lambda\mu} \implies (\mu+\lambda)^2 \ge 2\lambda\mu \quad \underline{\text{true}}$

variance $\le \frac{\sigma^2}{2n} \lambda^{-1} (\text{tr}\,\hat{\Sigma})$    bias $\le \frac{1}{2}\lambda \|\theta_*\|^2$

$\subset$ • not optimal
$\subset$ • useless in practice

"optimal" tradeoff : $\boxed{\lambda^2 = \frac{\sigma^2}{n} \frac{\text{tr}\,\hat{\Sigma}}{\|\theta_*\|^2}} \implies \text{Risk} = \frac{\sigma}{\sqrt{n}} \|\theta_*\| \sqrt{\text{tr}\,\hat{\Sigma}}$

compare to $\sigma^2 d$

## Homogeneity :

$$y = \varphi(x)^\top \Theta + \varepsilon$$

$$\mathbb{E}\,\varepsilon^2 = \sigma^2$$

red annotations: $m$, $hg \cdot m^{-1}$, $\dfrac{\text{Rish} + \lambda \|\Theta\|^2}{hg^{1/2}}$, $\dfrac{}{hg^2 \cdot m^{-2}}$, $\sigma \sim hg$

$hg$

$$\|\varphi(x)\| \le R \quad, \quad \mathrm{tr}\,\hat{\Sigma} = \mathrm{tr}\, \frac{1}{m} \sum_i \varphi(x_i)\varphi(x_i)^\top$$

$$\le R^2 - m^2$$

$hg$

$$\text{Expected rish} \sim \frac{\sigma}{\sqrt{m}} \sqrt{\frac{\mathrm{tr}\,\hat{\Sigma}}{m}} \|\Theta_\sigma\| \sim hg \cdot m^{-1} \sim hg^2$$

$hg^{1/2}$

$$\delta^2 = \frac{\sigma^2}{m} \frac{\mathrm{tr}\,\hat{\Sigma}}{\|\Theta_\sigma\|^2} \sim \frac{hg^2 m^2}{hg^2 m^{-2}} = m^4$$

$$\delta \sim m^2$$

× Least squares ⌐ ordinary $\frac{\sigma^2 d}{m}$

⌐ regularized

$$\frac{1}{n_b} + \lambda \underline{\quad}$$

variance     bias

× Extensions ⌐ square loss ⌐ stat.
                linear parameters ─
                beyond $L^2$ ─ $L_1$