# Probability Theory Background

## Random Variables $X$ & $Y$

$$x_1, x_2, \cdots x_M \qquad y_1, y_2, \cdots y_L$$

## Joint Distribution $p(X, Y)$



$$p(X = x_i, Y = y_j)$$
$$p(x_i, y_j)$$

$$c_i = \sum_{j=1}^{L} n_{ij}$$

$N$ trials, let $n_{ij}$ be the no. of times we observe $X = x_i$ & $Y = y_j$.

As $N \to \infty$, $\quad p(x_i, y_j) = \dfrac{n_{ij}}{N}$

marginal distribution

$$\boxed{p(X = x_i) = \sum_{j=1}^{L} p(X = x_i, Y = y_j)} = \sum_{j=1}^{L} \frac{n_{ij}}{N} = \frac{c_i}{N}$$

Sum Rule

Conditional Probability of $Y = y_j$ given $X = x_i$

$$p(Y = y_j \mid X = x_i) = \frac{n_{ij}}{c_i}$$

$$\downarrow$$

$$p(y_j \mid x_i)$$

$$p(y_j, x_i) = \frac{n_{ij}}{N} = \left(\frac{n_{ij}}{c_i}\right) \cdot \left(\frac{c_i}{N}\right)$$

**Product Rule**
$$\boxed{p(Y = y_j, X = x_i) = p(Y = y_j \mid X = x_i') \cdot p(X = x_i)}$$

**Sum Rule**
$$\boxed{p(X) = \sum_Y p(X, Y)}$$

**Product Rule**
$$\boxed{p(X, Y) = p(Y \mid X)\, p(X)}$$

$$p(Y \mid X)\, p(X) = p(X, Y) = p(X \mid Y)\, p(Y)$$

**Bayes Rule,**
$$p(Y \mid X) = \frac{p(X \mid Y)\, p(Y)}{p(X)} \longrightarrow \text{prior}$$

$$\downarrow$$

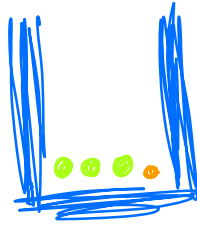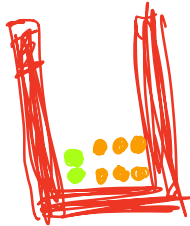posterior          evidence
                "data-likelihood"

$$= \frac{p(X \mid Y)\, p(Y)}{\sum_Y p(X, Y)}$$

$$= \frac{p(X \mid Y)\, p(Y)}{\sum_Y p(X \mid Y)\, p(Y)}$$

Independence : $p(x_i, y_j) = p(x_i) p(y_j)$ , $p(y_j | x_i) = p(y_j)$

$2a \& 6o$

$3a \& 1o$

Two Boxes : Red & Blue

Two kinds of fruit : Apples & Oranges



$p(r) = 0.4 = \frac{2}{5}$

$p(b) = 0.6 = \frac{3}{5}$

$p(a|r) = \frac{1}{4}$ , $p(o|r) = \frac{3}{4}$

$p(a|b) = \frac{3}{4}$ , $p(o|b) = \frac{1}{4}$

|  B | r | b | |
|---|---|---|---|
| F |  |  |  |
| a | 1/10 | 9/20 | 11/20 |
| o | 3/10 | 3/20 | 9/20 |
| | 2/5 | 3/5 | |

$p(a, r) = p(a|r) p(r) = \frac{1}{4} \cdot \frac{2}{5} = \frac{1}{10}$

$p(o, r) = p(o|r) p(r) = \frac{3}{4} \cdot \frac{2}{5} = \frac{3}{10}$

$p(a, b) = p(a|b) p(b) = \frac{3}{4} \cdot \frac{3}{5} = \frac{9}{20}$

$p(o, b) = p(o|b) p(b) = \frac{1}{4} \cdot \frac{3}{5} = \frac{3}{20}$

$p(a) = p(a, r) + p(a, b) = \frac{1}{10} + \frac{9}{20} = \frac{11}{20}$

$p(r|o) = \frac{p(o|r) p(r)}{p(o)} = \frac{p(o, r)}{p(o)} = \frac{3/10}{9/20} = \frac{3}{10} \cdot \frac{20}{9}$

$= \frac{2}{3}$

Probabilities wrt continuous variables, $x \in \mathbb{R}$, $x \in \mathbb{R}^d$

pdf = probability density function $p(x)$

$p(x) \geqslant 0$ , $\int_{-\infty}^{\infty} p(x) \, dx = 1$

Sum Rule : $\quad p(x) = \int_{-\infty}^{\infty} p(x,y) \, dy$

Product Rule : $\quad p(x,y) = p(y|x)p(x) = p(x|y)p(y)$

## Expectation (Mean)

$$E[f(x)] = \sum_x p(x)f(x) \qquad \int p(x)f(x) \, dx$$

## Variance

$$\left( E[x+y] = E[x] + E[y] \right)$$

$$Var[f(x)] = E\left[ \left( f(x) - E[f(x)] \right)^2 \right]$$

$$= E\left[ (f(x))^2 + (E[f(x)])^2 - 2f(x)E[f(x)] \right]$$

$$= E\left[ f(x)^2 \right] + E\left[ (E[f(x)])^2 \right] - 2E[f(x)]E[f(x)]$$

$$= E[(f(x))]^2 - \left( E[f(x)] \right)^2$$

$$Var(x) = E\left[ (x - E[x])^2 \right] = E[x^2] - (E[x])^2$$

$$cov(x,y) = E\left[ \{x - E[x]\}\{y - E[y]\} \right]$$

$$= E[xy] - E[x]E[y]$$

What if $x$ & $y$ are independent? $\quad p(x,y) = p(x)p(y)$

$$cov(x,y) = E[xy] - E[x]E(y)$$

$$\int p(x,y) \, xy \, dx \, dy = \int p(x)p(y) \, x \, y \, dx \, dy$$

$$= \int p(x)x \, dx \cdot \int p(y) \cdot y \cdot dy$$

$$= E[x] E[y]$$

$$\text{cov}(x, y) = 0 \quad \text{if} \quad x \& y \text{ are independent.}$$

## Gaussian Distribution / Normal Distribution

$$x \in \mathbb{R}$$

$$p(x \mid \mu, \sigma^2) = p(x) = \frac{1}{\sqrt{2\pi}\, \sigma} \cdot e^{-\frac{1}{2}(x-\mu)^2 / \sigma^2}, \quad \begin{array}{l} \mu = \text{mean} \\ \quad \text{or} \\ \quad \text{expectation} \\ \sigma^2 = \text{variance} \end{array}$$

$$\int_{-\infty}^{\infty} x\, p(x)\, dx = \mu = E[x]$$

$$E[(x-\mu)^2] = E[x^2] - \mu^2 = \sigma^2$$

$$x \in \mathbb{R}^d, \quad x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} \qquad E[x] = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_d \end{bmatrix} = \mu, \quad E[(x_i - \mu_i)^2] = \sigma_i^2$$

$$p(x) = p(x_1, x_2, \ldots, x_d)$$

Suppose $x_i$ is independent of $x_j$, $\forall \, i \neq j$

$$p(x) = \prod_{i=1}^{d} p(x_i)$$

$$= \frac{1}{(\sqrt{2\pi})^d \sigma_1 \sigma_2 \cdots \sigma_d} \prod_{i=1}^{d} e^{-\frac{1}{2}(x_i - \mu_i)^2 / \sigma_i^2}$$

$$= \frac{1}{(2\pi)^{d/2} \underbrace{\sigma_1 \sigma_2 \cdots \sigma_d}_{\to \, \det(\Sigma)^{1/2}}} \underbrace{e^{-\frac{1}{2} \sum_i (x_i - \mu_i)^2 / \sigma_i^2}}_{} \quad \to e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$$

$$(x-\mu)^T(x-\mu) = \sum_i (x_i - \mu_i)^2$$

$$(x-\mu)^T \Sigma^{-1}(x-\mu) = \sum_i (x_i - \mu_i)^2 / \sigma_i^2 \qquad \Sigma = \begin{bmatrix} \sigma_1^2 \ \sigma_2^2 & & O \\ & \ddots & \\ O & & \sigma_d^2 \end{bmatrix}$$

## General Case: $x \in \mathbb{R}^d$

## Multivariate Gaussian Distribution:

$$p(x) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$$

$\underbrace{\qquad}$ → determinant of $\Sigma$

$$\Sigma = \text{Covariance Matrix} = E[(x-\mu)(x-\mu)^T]$$

$\Sigma_{ij}$ is covariance between $x_i$ & $x_j$

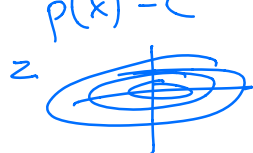$\Sigma$ is $d \times d$, symmetric, positive definite

$$\Sigma = V \Lambda V^T \qquad (\Lambda \text{ is diagonal}, \Lambda_{ii} > 0$$

$$\Sigma^{-1} = V \Lambda^{-1} V^T \qquad V^T V = I, \ V V^T = I)$$

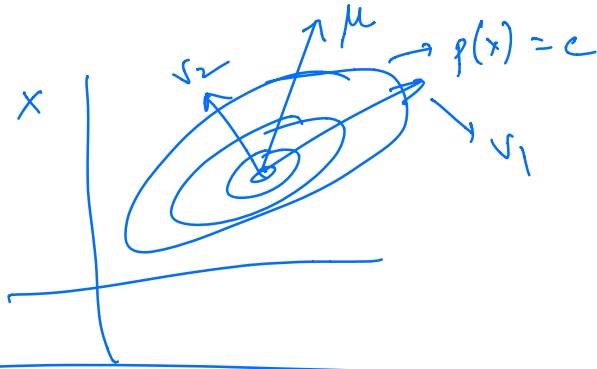$$\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu) = \frac{1}{2}(x-\mu)^T V \Lambda^{-1} V^T (x-\mu)$$

$$z = V^T(x-\mu)$$

$$= \frac{1}{2} z^T \Lambda^{-1} z$$

$$p(x) = c \quad \Rightarrow \quad \frac{1}{2} z^T \Lambda^{-1} z = c$$


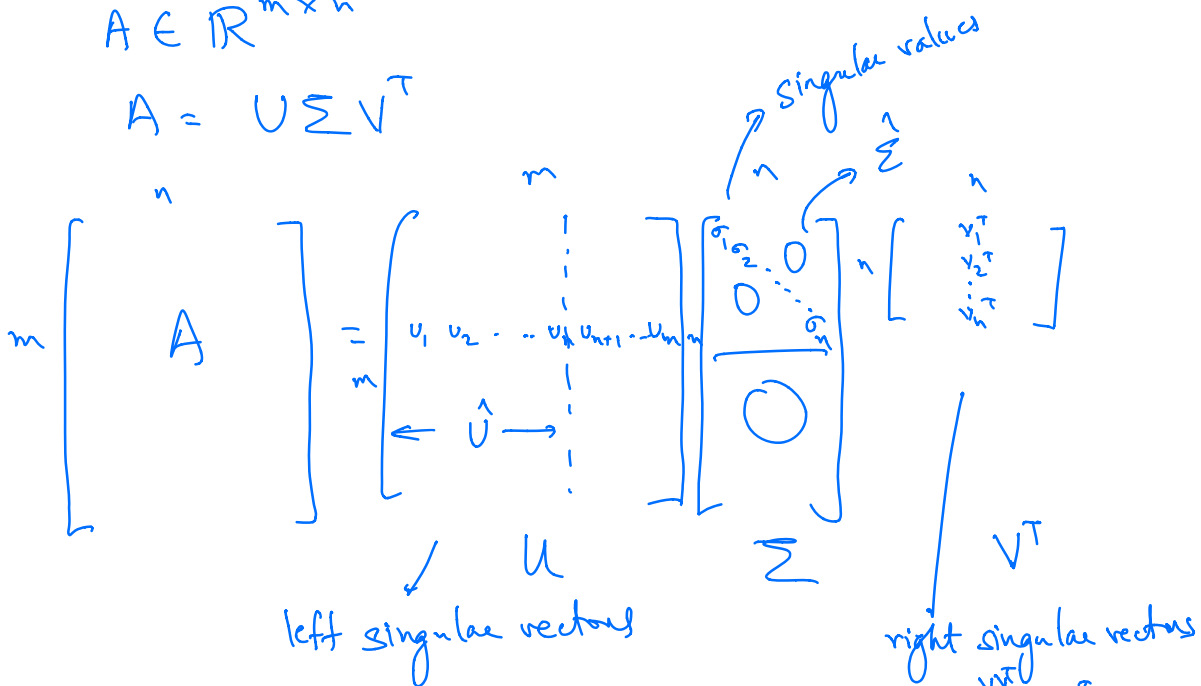
$$\boxed{\frac{1}{2} \sum \frac{z_i^2}{\sigma_i^2} = c}$$ → Equation of ellipse

# Singular Value Decomposition (SVD)

$$A \in \mathbb{R}^{m \times n}$$

$$A = U \Sigma V^T$$



$U$ — left singular vectors

$\Sigma$ — singular values

$V^T$ — right singular vectors

$$UU^T = U^TU = I, \qquad \sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n \geq 0, \quad V^TV = I$$

$$A = U\Sigma V^T \qquad \left( U \in \mathbb{R}^{m \times m}, \ \Sigma \in \mathbb{R}^{m \times n} \right)$$

$$AV = U\Sigma \qquad \qquad V \in \mathbb{R}^{n \times n}$$

$$A v_i = u_i \sigma_i$$

"Thin or reduced SVD": $A = \hat{U} \hat{\Sigma} \hat{V}^T, \quad \hat{U} \in \mathbb{R}^{m \times n}$

$m \geq n \qquad \hat{\Sigma} \in \mathbb{R}^{n \times n}$ diagonal matrix $\hat{V} \in \mathbb{R}^{n \times n}$
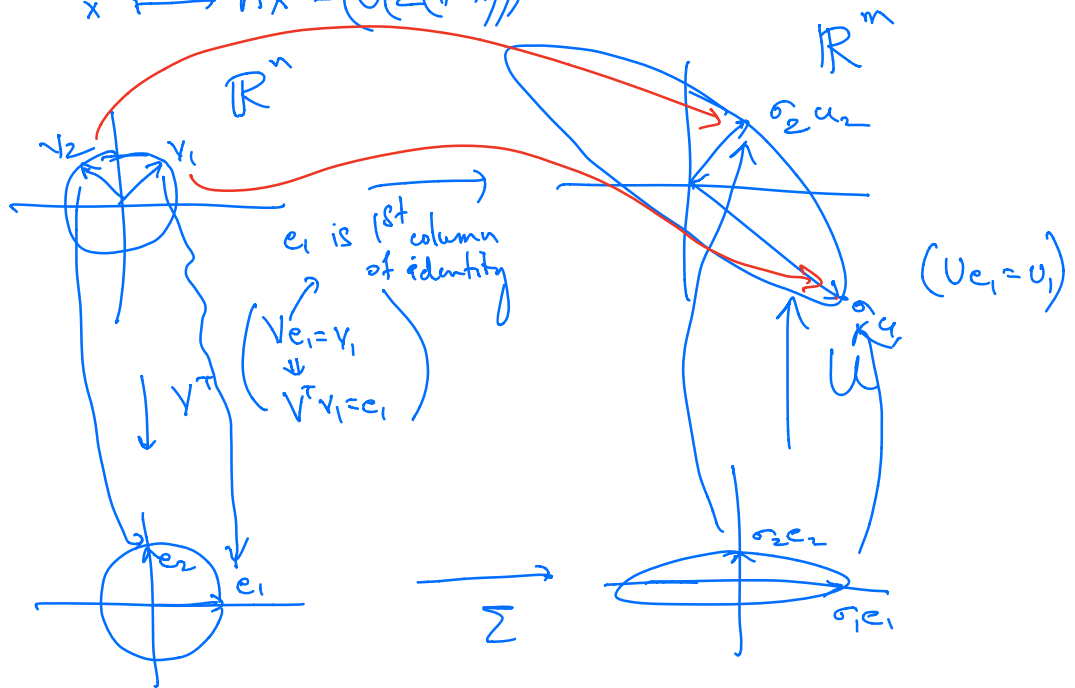
$$\hat{U}^T\hat{U} = I \quad , \quad \hat{U}\hat{U}^T \neq I$$

$$\hat{V}^T\hat{V} = I$$

$$n \quad \overset{n}{A^T A} = \left(\hat{V}\hat{\Sigma}\hat{U}^T\right)\left(\hat{U}\hat{\Sigma}\hat{V}^T\right) = \hat{V}\hat{\Sigma}^2\hat{V}^T \longrightarrow \text{eigenvalue decomposition of } A^T A$$

$$m \quad \overset{m}{AA^T} = \left(\hat{U}\hat{\Sigma}\hat{V}^T\right)\left(\hat{V}\hat{\Sigma}\hat{U}^T\right) = \hat{U}\hat{\Sigma}^2\hat{U}^T$$

$$A = U\Sigma V^T \qquad , \qquad A^T = V\Sigma^T U^T$$

$$A^T U = V\Sigma^T$$

$$AV = U\Sigma \qquad\qquad\qquad A^T u_i = v_i \sigma_i \quad , \quad i = 1, 2, \dots n$$

$$A v_i = u_i \sigma_i \qquad\qquad\qquad A^T u_i = 0 \quad , \quad i = n+1, \dots m$$

$$A : \mathbb{R}^n \to \mathbb{R}^m$$

$$x \longmapsto Ax = \left(U\left(\Sigma\left(V^T x\right)\right)\right)$$



$\mathbb{R}^n$

$\mathbb{R}^m$

$v_2 \quad v_1$

$\sigma_2 u_2$

$(U e_1 = u_1)$

$e_1$ is 1st column of identity

$$\begin{pmatrix} V e_1 = v_1 \\ \Downarrow \\ V^T v_1 = e_1 \end{pmatrix}$$

$V^T$

$e_2 \quad e_1$

$\Sigma$

$\sigma_2 e_2$

$\sigma_1 e_1$

If A has rank $r$

$$\sigma_1 \geq \sigma_2 \geq \ldots \geq \sigma_r > 0, \quad \sigma_{r+1} = \sigma_{r+2} \cdots = \sigma_n = 0$$

$$A : \mathbb{R}^n \to \mathbb{R}^m \qquad\qquad A^T : \mathbb{R}^m \to \mathbb{R}^n$$

$$v_1 \longmapsto u_1 \sigma_1 \qquad\qquad u_1 \longmapsto v_1 \sigma_1$$
$$v_2 \longmapsto u_2 \sigma_2 \qquad\qquad u_2 \longmapsto v_2 \sigma_2$$
$$\vdots \qquad\qquad\qquad \vdots$$
$$v_r \longmapsto u_r \sigma_r \qquad\qquad u_r \longmapsto v_r \sigma_r$$
$$v_{r+1} \longmapsto 0 \qquad\qquad u_{r+1} \longmapsto 0$$
$$\vdots \qquad\qquad\qquad \vdots$$
$$v_n \longmapsto 0 \qquad\qquad u_m \longmapsto 0$$

SVD provides orthogonal basis for the four fundamental subspaces of A

Column Space $= R(A) = \langle u_1, u_2 \ldots u_r \rangle$

Row Space $= R(A^T) = \langle v_1, v_2, \ldots v_r \rangle$

Null Space(A) $= N(A) = \langle v_{r+1}, \ldots v_n \rangle$

Null Space $(A^T) = N(A^T) = \langle u_{r+1}, \ldots \ldots, u_m \rangle$

$A \in \mathbb{R}^{m \times n}$

Truncated SVD, $A_k = U_k \Sigma_k V_k^T$, $\quad U_k \in \mathbb{R}^{m \times k}$, $\Sigma_k \in \mathbb{R}^{k \ast k}$

$U_k^T U_k = I$
$V_k^T V_k = I \longleftarrow V_k \in \mathbb{R}^{n \times k}$

Among all rank-k approximations of A, $A_k$ is the "best"

$$A_k = \arg\min_{B \text{ of rank } k} \|A - B\|_2 \quad , \quad A_k = \arg\min_{B \text{ of rank } k} \|A - B\|_F$$

---

## Regression

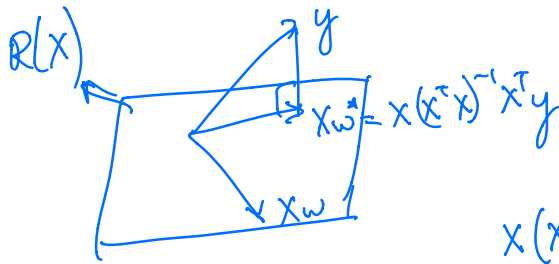$$\boxed{\min_w \|y - Xw\|_2^2} \quad , \quad X = \begin{array}{c} \\ N \end{array} \overset{d+1}{\begin{bmatrix} 1 & x_1^T \\ 1 & x_2^T \\ \vdots \\ 1 & x_N^T \end{bmatrix}}$$

Least Squares Solution $\quad X^T X w^* = X^T y$

$$w^* = (X^T X)^{-1} X^T y$$

Prediction on training set: $\quad Xw^* = X(X^T X)^{-1} X^T y$



$$Xw^* = X(X^T X)^{-1} X^T y$$

$$U \in \mathbb{R}^{N \times (d+1)}$$

$$X = U\Sigma V^T \quad - \text{ reduced SVD}$$

$$X(X^T X)^{-1} X^T \longrightarrow \text{ "hat" matrix}$$

$$(U\Sigma V^T)(V\Sigma U^T U\Sigma V^T)^{-1} V\Sigma U^T$$

$$\underbrace{\phantom{U\Sigma U^T}}_{I}$$

$$U\Sigma V^T (V\Sigma^2 V^T)^{-1} V\Sigma U^T$$

$$U\Sigma V^T \underbrace{(V\Sigma^{-2} V^T)}_{I} V\underbrace{\Sigma U^T}_{I} = UU^T$$

$$U = [v_1 \; v_2 \; \dots \; v_{d+1}]$$

$$U^T U = I \quad , \quad UU^T \neq I \quad (\text{orthogonal projector})$$

$$UU^T = [v_1 \; v_2 \dots v_{d+1}] \begin{bmatrix} v_1^T \\ v_2^T \\ \vdots \\ v_{d+1}^T \end{bmatrix} = v_1 v_1^T + v_2 v_2^T + \dots + v_{d+1} v_{d+1}^T$$

Least Squares Prediction $= \boxed{UU^T y} \longrightarrow \Sigma u_i u_i^T$

$$= \sum_{i=1}^{n} u_i (u_i^T y)$$

Least Squares Regression: $\min_{w} \|y - Xw\|_2^2$

Ridge Regression: $\min_{w} \|y - Xw\|_2^2 + \lambda \|w\|_2^2$, $\lambda \geq 0$

Solution: $(X^T X + \lambda I) w^* = X^T y$

$\Rightarrow \quad w^* = (X^T X + \lambda I)^{-1} X^T y$ $\qquad N \geq d+1$

Prediction: $Xw^* = X(X^T X + \lambda I)^{-1} X^T y$

$X = U\Sigma V^T$ — reduced SVD

$X^T X = V \Sigma^2 V^T$

$X^T X + \lambda I = V(\Sigma^2 + \lambda I) V^T$ $\qquad (VV^T = I)$

$(X^T X + \lambda I)^{-1} = V(\Sigma^2 + \lambda I)^{-1} V^T$

$$U\Sigma \underbrace{V^T \cdot V}_{I} (\Sigma^2 + \lambda I)^{-1} \underbrace{V^T V}_{I} \Sigma U^T$$

$$U\Sigma \underbrace{(\Sigma^2 + \lambda I)^{-1}}_{} \Sigma U^T$$

$$\begin{bmatrix} \frac{\sigma_1^2}{\sigma_1^2 + \lambda} & & & \\ & \frac{\sigma_2^2}{\sigma_2^2 + \lambda} & & \\ & & \ddots & \\ & & & \frac{\sigma_{d+1}^2}{\sigma_{d+1}^2 + \lambda} \end{bmatrix}$$

Ridge Regression Solution $= \sum u_i \left( \dfrac{\sigma_i^2}{\sigma_i^2 + \lambda} \right) v_i^T y$

$\sigma_i^2 \gg \lambda$ , $\dfrac{\sigma_i^2}{\sigma_i^2 + \lambda} \approx 1$

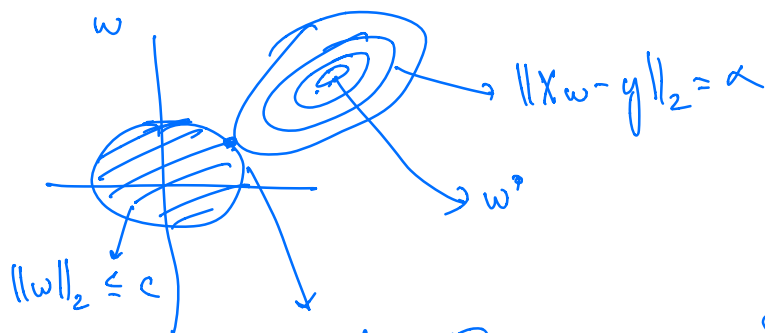$\sigma_i^2 \ll \lambda$ , $\dfrac{\sigma_i^2}{\sigma_i^2 + \lambda} \approx 0$

$\downarrow$ Shrinkage

Ridge Regression can equivalently be thought of

as :

$$\min_{w} \|Xw - y\|_2$$

$$st \quad \|w\|_2 \le c$$



$w$

$\|Xw - y\|_2 = \alpha$

$w^*$

$\|w\|_2 \le c$

Ridge Regression Solution