

MAT 180 Notes

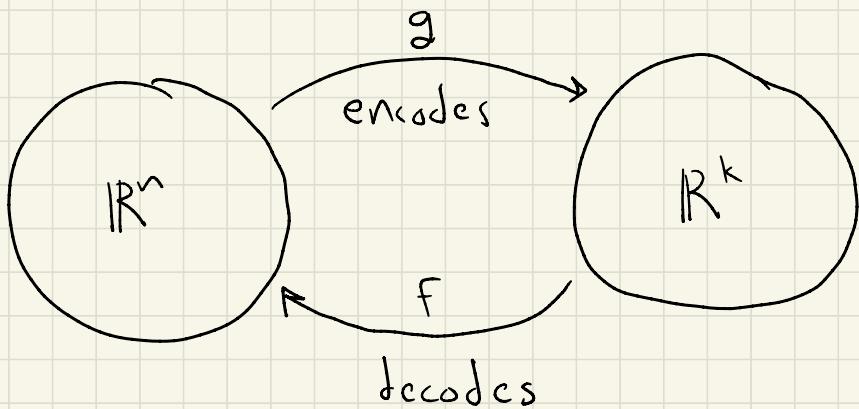


Principal Component Analysis

Dataset $\{x^{(1)}, \dots, x^{(m)}\} \in \mathbb{R}^n$.

Goal: compress data to \mathbb{R}^k

Want functions



$$\text{s.t. } x \approx f(g(x))$$

$$\text{Fix } f(y) = Dy \quad D \in \mathbb{R}^{n \times k}$$

D orthonormal columns

$$(\text{i.e. } D^T D = I)$$

To be optimal, we want

$$g(x) = \underset{y}{\operatorname{argmin}} \|x - f(y)\|_2^2$$

$$= \underset{y}{\operatorname{argmin}} (x - f(y))^T (x - f(y))$$

$$= \underset{y}{\operatorname{argmin}} (x^T x - f(y)^T x - \underline{x^T f(y)} + f(y)^T f(y))$$

$$= \underset{y}{\operatorname{argmin}} (-2x^T f(y) + f(y)^T f(y))$$

$$= " (-2x^T D_y + y^T D^T D_y)$$

$$= " (-2x^T D_y + y^T y)$$

Solve the optimization problem

$$\nabla_y (-2x^T D_y + y^T y) = 0$$

Rule: $\nabla_y (M y) = y^T \nabla_y (M^T) + M \cdot \nabla_y$

$$\Rightarrow -2x^T D + y^T \nabla(y) + y^T \nabla(y) = 0$$

argmin
means the
argument y
that
minimizes

$$\Rightarrow -2x^T D + 2y^+ = 0 \Rightarrow y^+ = x^T D$$

$$\Rightarrow \boxed{y = D^T x}$$

The optimal solution is

$$g(x) = D^T x.$$

* still need D

Recall: we want $f(g(x))$ as close as possible to x for all points in the data set. That is,

$$f(g(x^{(i)}))$$

$$D = \underset{\substack{C: C^T C = I}}{\operatorname{argmin}} \sum_{i=1}^m \|x^{(i)} - C(C^T x^{(i)})\|_2^2$$

$$X = \underbrace{\begin{pmatrix} -x^{(1)} \\ \vdots \\ -x^{(m)} \end{pmatrix}}_{\in \mathbb{R}^{m \times n}} = \operatorname{argmin}_{\|X\|_F} \|X - X C C^T\|_F^2$$

$$= \text{Tr}((X - XCCT^T)^T(X - XCCT^T))$$

$$= \text{Tr}\left(\cancel{X^T X} - CCT^T X^T X - X^T X CCT^T + CCT^T X^T X CCT^T\right)$$

$$= \text{Tr}\left(-\text{Tr}(CCT^T X^T X) - \text{Tr}(X^T X CCT^T) + \underbrace{\text{Tr}(CCT^T X^T X CCT^T)}_{\uparrow}\right)$$

$$= \text{Tr}\left(-2\text{Tr}(\underbrace{CCT^T X^T X}_{\uparrow}) + \text{Tr}(X^T X CCT^T)\right)$$

$$= \underset{C: C^T C = I}{\text{argmin}} - \text{Tr}\left(\underbrace{X^T X CCT^T}_{\uparrow}\right)$$

$$= \underset{C: C^T C = I}{\text{argmax}} \text{Tr}(C^T X^T X C)$$

Lemma The optimal encoding function

$$\text{is } g(x) = D^T x \quad \text{where}$$

$$D = \underset{C: C^T C = I}{\text{argmax}} \text{Tr}(C^T X^T X C)$$

To proceed, consider $k=1$, so

C is a unit vector.

Goal: maximize $\text{Tr}(C^T X^T X C) = \underbrace{C^T X^T X}_{} C$

Note: $X^T X$ is positive semi-definite, symmetric

Spectral theorem $\Rightarrow X^T X$ is diagonalizable with
orthonormal eigenvectors v_1, \dots, v_n
eigenvalues $\lambda_1 \geq \dots \geq \lambda_n \geq 0$

$$\text{so } C \in \mathbb{R}^n \Rightarrow C = \sum_i a_i v_i \quad \sum a_i^2 = 1$$

$$\Rightarrow C^T X^T X C = \left(\sum_i a_i v_i \right)^T X^T X \left(\sum_j a_j v_j \right)$$

$$X^T X v_j = \lambda_j v_j$$

$$= \sum_{ij} a_i a_j v_i^T \lambda_j v_j$$

$$v_i^T v_j = \begin{cases} 1 & i=j \\ 0 & \text{else} \end{cases}$$

$$= \sum_i a_i^2 \lambda_i$$

largest value is λ_1

maximized when $a_1 = 1$ $a_2 = \dots = a_n = 0$

$$S_0 \quad C = V_1 \quad \begin{pmatrix} \text{First principal} \\ \text{component} \end{pmatrix}$$

the eigenvector of $X^T X$
corresponding to the
largest eigenvalue

Claim In general, the optimal
solution is

$$D = C = \min \left[\begin{matrix} 1 & | \\ V_1 & \dots & V_k \\ 1 & | \end{matrix} \right] \quad \begin{matrix} \text{eigenvectors of} \\ X^T X \text{ corresponding} \\ \text{to largest } k \\ \text{eigenvalues} \end{matrix}$$

↑
principal components
of X

Proof: (induction) base case above
(HW)

§ 3.1 Probability & information Theory

9/28/22

(non-rigorous) (for rigour need)
Measure theory

§ 3.2

A random variable (RV) is a variable that takes on values from a specified set Ω of "states"

E.g. Outcomes of some random experiment

Notation X RV

x_1, x_2, x_3 values that X takes

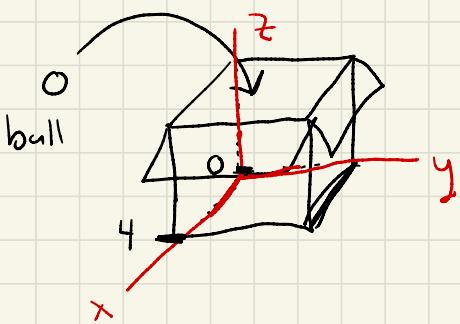
(discrete)

E.g. X measures outcomes of rolling a 6-sided die. Maybe $x_1 = 2 \quad x_2 = 6 \quad x_3 = 1, \dots$

$$\Omega = \{1, 2, 3, 4, 5, 6\}$$

(continuous)

E.g. $X =$ resting x -coord



$$\Omega = [0, 4]$$

$$x_1 = 1.213$$

$$x_2 = \pi$$

§ 3.3 Probability Distributions

How likely is each state of a RV?

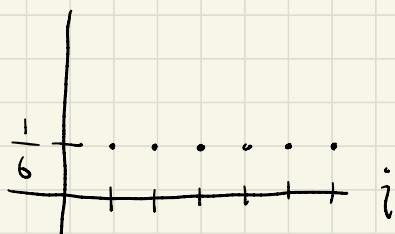
Discrete Probability Mass Function

$$P: \Omega \rightarrow [0,1] \text{ s.t. } \sum_{x \in \Omega} P(x) = 1$$

Notation $P(x) = P(x=x)$ probability of x occurring

E.g. $x = 6$ -sided die roll

$$i \in \{1, \dots, 6\}, P(i) = \frac{1}{6} \quad \sum_{i=1}^6 P(i) = 1.$$



E.g. $x = \text{two die rolls } (6\text{-sided, } 8\text{-sided})$

$$\Omega = \{(i,j) \mid 1 \leq i \leq 6, 1 \leq j \leq 8\}$$

$$P(i,j) = \frac{1}{48}$$

Joint distribution

$X = \text{outcome of 6-sided die}$
 $Y = \text{" " " 8-sided die}$

$$\Omega_X = \{1, \dots, 6\}$$

$$\Omega_Y = \{1, \dots, 8\}$$

$$P_X(i) = \frac{1}{6}$$

$$P_Y(j) = \frac{1}{8}$$

$$P(X=i, Y=j) = \frac{1}{48} = \frac{1}{6} \cdot \frac{1}{8} = P_X(i) P_Y(j)$$

↑
independent
variables

E.g. Two 6-sided d.e rolls (indistinguishable)

$$\Omega = \left\{ \{i, j\} \mid 1 \leq i, j \leq 6 \right\}$$

$$P(\{i, i\}) = \frac{1}{36} \quad 1 \leq i \leq 6$$

12 13 14 15 16

23 24 25 26

31 33 36

45 46

56

$$P(\{i, j\}) = \frac{2}{36} \quad 1 \leq i \neq j \leq 6$$

$$6 \cdot \frac{1}{36} + 15 \cdot \frac{2}{36} = \frac{6+30}{36} = 1$$

Continuous \cap uncountable (\mathbb{R} or $[0,1]$)

probability density function

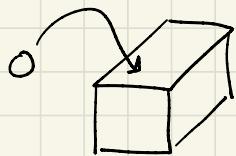
$$p: \Omega \rightarrow \mathbb{R}_{\geq 0} \text{ s.t.}$$

$$\underbrace{\int_{\Omega} p}_{\text{Lebesgue integral}} = 1$$

Interpretation: the probability
that the outcome is in

$$S \subseteq \Omega \text{ is } \int_S p$$

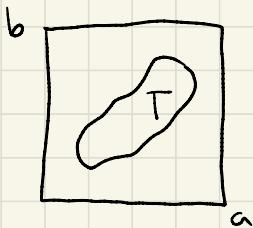
E.g.



$$\Omega = [0, a] \times [0, b] \quad x \quad y$$

$$P(x=x, y=y) = c$$

$$\int_{\Omega} p = 1 \implies \int_0^a \int_0^b p \, dy \, dx = cab \implies c = \frac{1}{ab}$$



$$P(T) = \int_T p = \iint_T \frac{1}{ab} dA = \frac{\text{area}(T)}{ab}$$

§ 3.4 Marginal Probability

Joint dist on x, y

$$p(x=x) = \sum_y P(x=x, y=y) \quad \text{discrete}$$

$$p(x=x) = \int p(x=x, y=y) dy \quad \text{continuous}$$

§ 3.5 Conditional Probability

Probability that $x=x$ given that $y=y$

$$p(x=x | y=y) = \frac{p(x=x, y=y)}{p(y=y)}$$

E.g. 2 - 6 sided die $x = \text{outcome } \{i, j\}$
 $y = \text{value of } i+j$

$$p(x=\{3, 4\} | y=7) = \frac{p(x=\{3, 4\}, y=7)}{p(y=7)}$$

$$= \frac{p(x=\{3, 4\})}{p(y=7)} = \frac{\frac{2}{36}}{\frac{1}{6}} = \boxed{\frac{1}{3}}$$

Chain Rule

Bayes Rule

} read
these

Independence of RV

Expectation value of $f(x)$ w.r.t dist P

$$f: \Omega \rightarrow \mathbb{R}$$

Expectation value: most likely outcome $f(x)$

$$E_p[f] = \sum_{x \in \Omega} p(x)f(x) \quad \text{discrete}$$

$$= \int_{\Omega} p f \quad \text{continuous}$$