# Homework 1

## 1) a)

$\mu_{dist}$



$\sigma_{dist}$



(red) $\frac{1}{2}$     (red) $\frac{1}{2}$

## b)

$$z = (x-y)^2$$

$X, Y \sim$ independent r.v. from $\text{Unif}(0,1)$

$p(x)$ — graph from $0$ to $1$, dashed at $\frac{1}{2}$

Mean of $X, Y = \mathbb{E}[X] = \frac{1}{2}$

Variance of $X = \mathbb{E}[(X-\mu)^2] = \frac{1}{12}$

Expected value of $z$

$$\mathbb{E}[z] = \mathbb{E}[(x-y)^2] = \mathbb{E}\left[\left\{(x-\tfrac{1}{2})-(y-\tfrac{1}{2})\right\}^2\right]$$

$$= \mathbb{E}\left[(x-\tfrac{1}{2})^2 + (y-\tfrac{1}{2})^2 + (x-\tfrac{1}{2})(y-\tfrac{1}{2})\right]$$

$$= \mathbb{E}\left[(x-\tfrac{1}{2})^2\right] + \mathbb{E}\left[(y-\tfrac{1}{2})^2\right] + \mathbb{E}\left[(x-\tfrac{1}{2})(y-\tfrac{1}{2})\right]$$

Same    X and Y are independent of each other

Note: $\mathbb{E}\left[(x-\tfrac{1}{2})\right]$

$= \mathbb{E}[X] - \tfrac{1}{2}$

$= \tfrac{1}{2} - \tfrac{1}{2} = 0$

$$= 2\,\mathbb{E}\underbrace{\left[(x-\tfrac{1}{2})^2\right]}_{\tilde{\mu}_x} + \mathbb{E}\left[(x-\tfrac{1}{2})\right]\mathbb{E}\left[(y-\tfrac{1}{2})\right]$$

$$= 2 \times \frac{1}{12} = \frac{1}{6}$$

$\boxed{\mathbb{E}[z] = \frac{1}{6}}$     (red) $0.75$

$$\mathbb{E}\left[(z-\mu_z)^2\right] = \mathbb{E}[z^2] - \mu_z^2$$

$$= \mathbb{E}[(x-y)^4] - (\mathbb{E}[z])^2$$

$$= \mathbb{E}[x^4 + y^4 - 4x^3y - 4xy^3 + 6x^2y^2] - \frac{1}{36}$$

$$= \mathbb{E}[x^4] + \mathbb{E}[y^4] - 4\mathbb{E}[x^3y] - 4\mathbb{E}[xy^3] + 6\mathbb{E}[x^2y^2] - \frac{1}{36}$$

$$\mathbb{E}[X^4] = \int_0^1 x^4 \frac{1}{(1-0)} \, dx = \frac{x^5}{5}\Big|_0^1 = \frac{1}{5}$$

$$\mathbb{E}[X^3 Y] = \mathbb{E}[X^3]\,\mathbb{E}[Y] \quad \text{(due to independence)}$$

$$= \frac{x^4}{4}\Big|_0^1 \times \frac{y^2}{2}\Big|_0^1 = \frac{1}{8}$$

$$\mathbb{E}[XY^3] = \frac{x^2}{2}\Big|_0^1 \times \frac{y^4}{4}\Big|_0^1 = \frac{1}{8}$$

$$\mathbb{E}[Y^4] = \frac{y^5}{5}\Big|_0^1 = \frac{1}{5} \qquad\qquad \mathbb{E}[X^2 Y^2] = \frac{x^3}{3}\Big|_0^1 \times \frac{y^3}{3}\Big|_0^1 = \frac{1}{9}$$

$$\mathbb{E}[(Z-M_Z)^2] = \mathbb{E}[X^4] + \mathbb{E}[Y^4] - 4\mathbb{E}[X^3 Y] - 4\mathbb{E}[XY^3] + 6\mathbb{E}[X^2 Y^2] - \frac{1}{36}$$

$$= \frac{1}{5} + \frac{1}{5} - 4\left(\frac{1}{8}\right) - 4\left(\frac{1}{8}\right) + 6\left(\frac{1}{9}\right) - \frac{1}{36}$$

$$= \frac{2}{5} - 1 + \frac{2}{3} - \frac{1}{36} = \frac{7}{180}$$

$$\boxed{Var(Z) = \mathbb{E}[(Z-M_Z)^2] = \frac{7}{180}}$$

<span style="color:red">**0.75**</span>

c)
$$\mathbb{E}[S] = \mathbb{E}[Z_1 + Z_2 + \cdots + Z_d]$$

$$= \mathbb{E}\left[\sum_{i=1}^{d} Z_i\right] = \mathbb{E}\left[\sum_{i=1}^{d} (X_i - Y_i)^2\right]$$

$$= \sum_{i=1}^{d} \mathbb{E}\left[(X_i - Y_i)^2\right]$$

We know that $\mathbb{E}[(X_i - Y_i)^2] = \frac{1}{6}$ and that $\mathbb{E}[(X_i - Y_i)^2] = \mathbb{E}[(X_j - Y_j)^2]$
$$\forall \; i,j \in d$$

$$\therefore \; \mathbb{E}[S] = \sum_{i=1}^{d} \mathbb{E}[(X_i - Y_i)^2] = \sum_{i=1}^{d}\left(\frac{1}{6}\right) = \frac{d}{6}$$

$$\boxed{\mathbb{E}[S] = d\,\mathbb{E}[Z] = \frac{d}{6}} \qquad \text{<span style=\"color:red\">0.5</span>}$$

Similarly, we calculate variance of $S$

$$Var[S] = Var[Z_1 + Z_2 + \cdots + Z_d]$$

$$= \sum_{i=1}^{d} Var(Z_i) \quad \begin{bmatrix}\text{since } Z_i\text{'s are} \\ \text{independent}\end{bmatrix}$$

<span style="color:green">(Note $Var(X+Y) = Var(X) + Var(Y)$ $+ Cov(X,Y)$</span>

<span style="color:green">If $X$ and $Y$ are independent $\Rightarrow Var(X+Y) = Var(X) + Var(Y)$)</span>

$$= d\,Var(Z) = \frac{7d}{180} \qquad \text{<span style=\"color:red\">0.5</span>}$$

d) Markov's inequality says

$$P\left(|z - \mathbb{E}[z]| \geq a\right) \leq \frac{\text{Var}[z]}{a^2}$$

or,

$$P\left(|s - \mathbb{E}[s]| \geq a\right) \leq \frac{\text{Var}[s]}{a^2}$$

$$\Rightarrow \quad P\left(|s - d/6| \geq a\right) \leq \frac{7d}{180\, a^2}$$

Note $S = \|\underline{X} - \underline{Y}\|_2^2$
represents the distance
between two points lying
in $d$-dimensional space

$\left|s - d/6\right| \rightarrow$ is also a distance

probability that this distance $\left|s - \frac{d}{6}\right|$ is greater than 'a'

(0.5)

Say $a = 1$

For $d = 1$

$$P\left(|s - 1/6| \geq 1\right) \leq \frac{7}{180}$$

In 1-D, the chances of the distance between 2 points exceeding a certain values is less

For $d = 5$

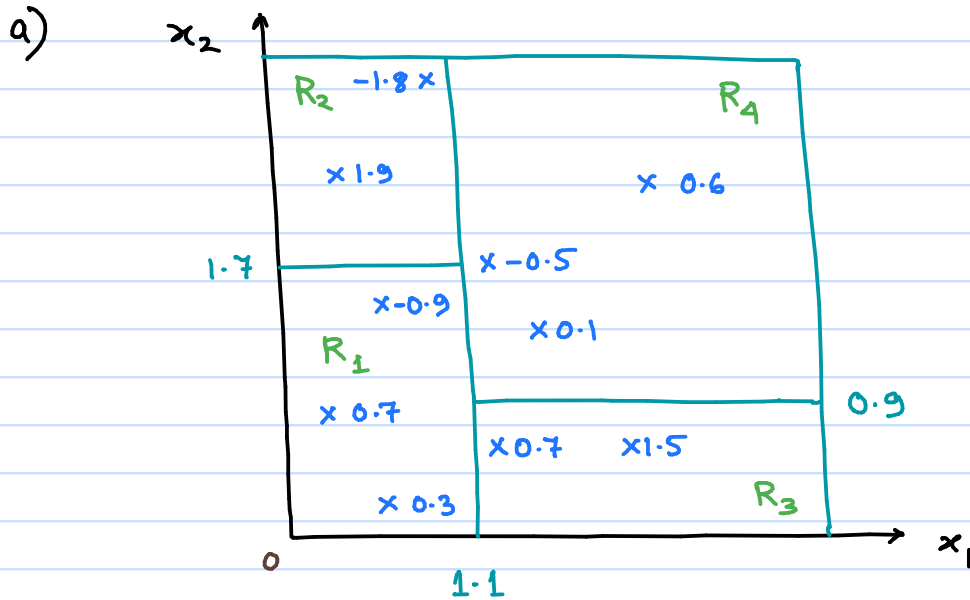$$P\left(|s - 5/6| \geq 1\right) \leq \frac{35}{180}$$

For $d = 10$

$$P\left(|s - 10/6| \geq 1\right) \leq \frac{70}{180}$$

In 10-D, the chances of the distance between 2 points exceeding a certain values is much more

Hence, we find that with increasing dimension, the distance between points increases, and most points in higher dimensions are quite far apart!
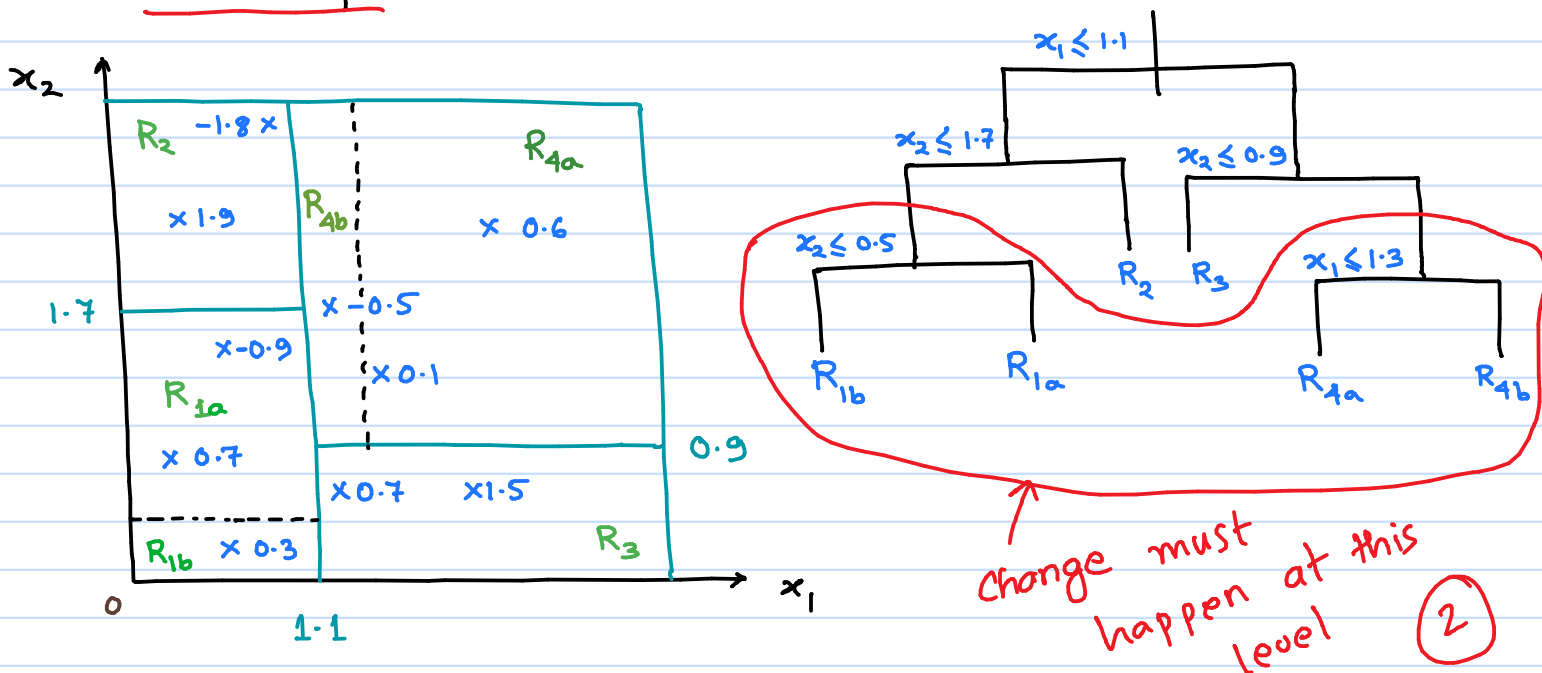
# 3) Regression Tree

a)

$x_2$ axis plot:

$R_2$  $-1.8 \times$    $R_4$
$\times 1.9$    $\times 0.6$
1.7
$\times -0.5$
$\times -0.9$
$\times 0.1$
$R_1$
$\times 0.7$    0.9
$\times 0.7$  $\times 1.5$
$\times 0.3$    $R_3$
0    1.1    $x_1$

①

b) Since $x_1^* = 1.5 > 1.1$ and $x_2^* = 1.8 > 0.9$, the test points belongs to region $R_4$.

The mean of the training point output in $R_4$ is $\hat{y}_{R_4} = 0.0667$

Therefore, prediction becomes $\hat{y}^* = 0.067$    ⑤ 0.5

c) There could be many possibilities of creating a deeper tree. One example could be

$x_2$ axis plot:

$R_2$  $-1.8 \times$    $R_{4a}$
$\times 1.9$  $R_{4b}$  $\times 0.6$
1.7
$\times -0.5$
$\times -0.9$
$\times 0.1$
$R_{1a}$
$\times 0.7$    0.9
$\times 0.7$  $\times 1.5$
$R_{1b}$  $\times 0.3$    $R_3$
0    1.1    $x_1$

Tree:

$x_1 \leq 1.1$
$x_2 \leq 1.7$        $x_2 \leq 0.9$
$x_2 \leq 0.5$    $R_2$   $R_3$   $x_1 \leq 1.3$
$R_{1b}$   $R_{1a}$    $R_{4a}$   $R_{4b}$

change must happen at this level    ②

d) Based on the above tree, $x^*$ belongs to region $R_{4a}$,

0.5    thus $\hat{y}^* = \hat{y}_{R_{4a}} = 0.35$

4)

a) $$\underline{y} = \underline{\underline{X}}\,\underline{\theta} + \underline{\epsilon}, \qquad \underline{\epsilon} \sim N(\underline{0}, \sigma^2 I_N)$$

- The likelihood turns out to be Gaussian

$I_N$ is an identity matrix of size $N \times N$

$N$ — size of training data

$$P(\underline{y} \mid \underline{\underline{X}}\,\underline{\theta}) = N(\underline{\underline{X}}\,\underline{\theta}, \sigma^2 I_N)$$

$$= \frac{1}{(2\pi)^{N/2}\,|\sigma^2 I_N|^{1/2}}\,\exp\left(-\frac{1}{2\sigma^2}(\underline{y} - \underline{\underline{X}}\,\underline{\theta})^T(\underline{y} - \underline{\underline{X}}\,\underline{\theta})\right)$$

- Log-likelihood

(0.5)

$$\ln P(\underline{y} \mid \underline{\underline{X}}\,\underline{\theta}) = -\frac{N}{2}\log 2\pi - \frac{N}{2}\log \sigma^2 - \frac{1}{2\sigma^2}\underbrace{(\underline{y} - \underline{\underline{X}}\,\underline{\theta})^T(\underline{y} - \underline{\underline{X}}\,\underline{\theta})}_{\text{dependence on } \underline{\theta}}$$

To maximize the log-likelihood, we take derivative w.r.t. $\underline{\theta}$ and set it to zero

$$\frac{\partial}{\partial \underline{\theta}}\ln P(\underline{y} \mid \underline{\underline{X}}\,\underline{\theta}) = \frac{1}{2\sigma^2}\,2\,\underline{\underline{X}}^T(\underline{y} - \underline{\underline{X}}\,\underline{\theta}) = 0$$

$$\underline{\underline{X}}^T(\underline{y} - \underline{\underline{X}}\,\underline{\theta}) = 0$$

$$\Rightarrow \underline{\underline{X}}^T\underline{\underline{X}}\,\underline{\theta} = \underline{\underline{X}}^T\underline{y}$$

If $\underline{\underline{X}}^T\underline{\underline{X}}$ is invertible, then

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

(0.5)

b) In practice, $\underline{\underline{X}}$ is a tall matrix with more rows than columns

The columns of matrix $\underline{\underline{X}}$ denote the different input features

If $\underline{\underline{X}}^T\underline{\underline{X}}$ is not invertible $\rightarrow$ $\underline{\underline{X}}$ is rank-deficient

(0.5) $\Rightarrow$ In practice, it means some input features are redundant

5) Logistic function, $h(x) = \dfrac{e^x}{1 + e^x}$

(0.5)

a) $\dfrac{dh(x)}{dx} = \dfrac{e^x(1+e^x) - e^x \cdot e^x}{(1+e^x)^2} = \dfrac{e^x(1+e^x-e^x)}{(1+e^x)^2}$

$$= \dfrac{e^x}{1+e^x} \cdot \dfrac{1}{1+e^x}$$

$$= \left(\dfrac{e^x}{1+e^x}\right) \cdot \left(1 - \dfrac{e^x}{1+e^x}\right)$$

$$= h(x) \cdot (1 - h(x))$$

b) We will now consider the two classes as $\{0,1\}$ (instead of $\{-1,1\}$)

Treat

$$p(y=1 \mid \underline{x}; \underline{\theta}) = h(\underline{x}^T\underline{\theta}) \quad ; \quad p(y=0 \mid \underline{x}; \underline{\theta}) = 1 - h(\underline{x}^T\underline{\theta})$$

$$= \dfrac{e^{\underline{x}^T\underline{\theta}}}{1 + e^{\underline{x}^T\underline{\theta}}} \qquad\qquad = \dfrac{1}{1 + e^{\underline{x}^T\underline{\theta}}}$$

Log-Likelihood for a data pair $\{\underline{x}^{(i)}, y^{(i)}\}$

$$\ln p\left(y^{(i)} \mid \underline{x}^{(i)}; \underline{\theta}\right) = \begin{cases} \ln h(\underline{x}^T\underline{\theta}) & \text{if } y^{(i)} = 1 \\ \ln(1 - h(\underline{x}^T\underline{\theta})) & \text{if } y^{(i)} = 0 \end{cases}$$

To make the expression more compact, we write

$$\ln p\left(y^{(i)} \mid \underline{x}^{(i)}; \underline{\theta}\right) = y^{(i)} \ln h(\underline{x}^{(i)T}\underline{\theta}) + (1 - y^{(i)}) \ln(1 - h(\underline{x}^{(i)T}\underline{\theta}))$$

The log-likelihood for entire training data is

(1)

$$\ln p\left(y^{(1)}, \ldots, y^{(N)} \mid \underline{x}^{(1)}, \ldots, \underline{x}^{(N)}\right) = \sum_{i=1}^{N} y^{(i)} \ln h(\underline{x}^{(i)T}\underline{\theta})$$
$$+ (1 - y^{(i)}) \ln(1 - h(\underline{x}^{(i)T}\underline{\theta}))$$

c)
$$\ln p\left(y^{(1)}, \ldots, y^{(N)} \mid \underline{x}^{(1)}, \ldots, \underline{x}^{(N)}\right) = \sum_{i=1}^{N} y^{(i)} \ln h(\underline{x}^{(i)^T} \underline{\theta})$$
$$+ (1-y^{(i)}) \ln (1 - h(\underline{x}^{(i)^T} \underline{\theta}))$$

$$\frac{d}{d\underline{\theta}} \left[ y^{(i)} \ln \underbrace{h(\underline{x}^{(i)^T} \underline{\theta})}_{h} + (1-y^{(i)}) \ln (1 - \underbrace{h(\underline{x}^{(i)^T} \underline{\theta})}_{h}) \right]$$

$$= \quad y^{(i)} \frac{1}{h} \left(\frac{dh}{d\underline{\theta}}\right) \underline{x}^{(i)} + (1-y^{(i)}) \frac{1}{1-h} \left(-\frac{dh}{d\underline{\theta}}\right) \underline{x}^{(i)}$$

Using the relation $\frac{dh}{d\underline{\theta}} = h(1-h)$

$$= \quad y^{(i)} (1-h) \underline{x}^{(i)} - (1-y^{(i)}) h \underline{x}^{(i)}$$

$$= \quad y^{(i)} \underline{x}^{(i)} - y^{(i)} h \underline{x}^{(i)} - h \underline{x}^{(i)} + y^{(i)} h \underline{x}^{(i)}$$

$$= \quad \left(y^{(i)} - h\right) \underline{x}^{(i)}$$

$$= \quad \left(y^{(i)} - h(\underline{x}^{(i)^T} \underline{\theta})\right) \underline{x}^{(i)}$$

Therefore,

$$\frac{dL}{d\underline{\theta}} = \frac{d}{d\underline{\theta}} \ln p(\underline{y} \mid \underline{\underline{x}}; \underline{\theta}) = \sum_{i=1}^{N} \left(y^{(i)} - h(\underline{x}^{(i)^T} \underline{\theta})\right) \underline{x}^{(i)^T} \qquad \textcircled{1}$$

d) Differentiating further,

$$\underbrace{\frac{d^2 \ln p(y^{(i)} \mid \underline{x}^{(i)}; \underline{\theta})}{d\underline{\theta} \, d\underline{\theta}^T}}_{\substack{P \times P \\ \text{matrix} \\ \underline{\theta} \in \mathbb{R}^P}} = \frac{d}{d\underline{\theta}^T} \left(y^{(i)} - h(\underline{x}^{(i)^T} \underline{\theta})\right) \underline{x}^{(i)}$$

$$= - \frac{dh}{d\underline{\theta}^T} \underline{x}^{(i)} \underline{x}^{(i)^T}$$

$$= - h(1-h) \underbrace{\underline{x}^{(i)} \underline{x}^{(i)^T}}_{P \times P} \qquad \textcircled{1}$$

$$\frac{d^2 L}{d\underline{\theta} \, d\underline{\theta}^T} = - \sum_{i=1}^{N} h(\underline{x}^{(i)^T} \underline{\theta}) \left(1 - h(\underline{x}^{(i)^T} \underline{\theta})\right) \underline{x}^{(i)} \underline{x}^{(i)^T}$$