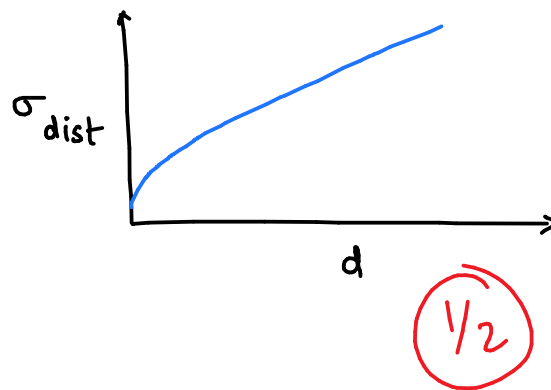
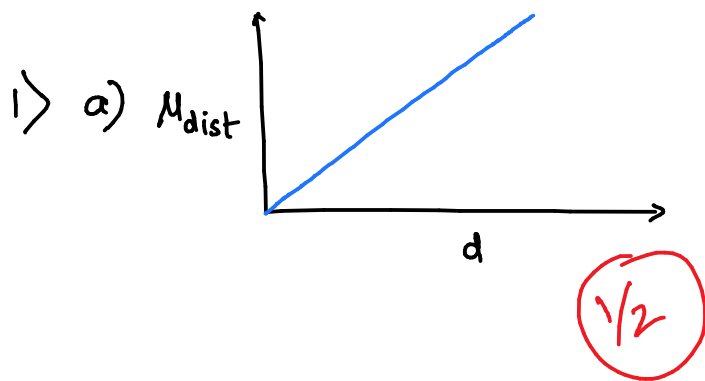


Homework 1

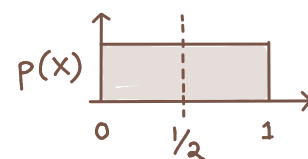


b) $z = (x - y)^2$

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

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

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

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

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

← Same

x and y are independent of each other

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

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

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

$$= \mathbb{E}[x] - \frac{1}{2}$$

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

$\mathbb{E}[z] = \frac{1}{6}$

0.75

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

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

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

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

$$\mathbb{E}[x^4] = \int_0^1 x^4 \frac{1}{(1-0)} dx = \left. \frac{x^5}{5} \right|_0^1 = 1/5$$

$$\begin{aligned}\mathbb{E}[x^3 y] &= \mathbb{E}[x^3] \mathbb{E}[y] \quad (\text{due to independence}) \\ &= \left. \frac{x^4}{4} \right|_0^1 \times \left. \frac{y^2}{2} \right|_0^1 = 1/8\end{aligned}$$

$$\mathbb{E}[x y^3] = \left. \frac{x^2}{2} \right|_0^1 \times \left. \frac{y^4}{4} \right|_0^1 = 1/8$$

$$\mathbb{E}[y^4] = \left. \frac{y^5}{5} \right|_0^1 = 1/5 \quad \mathbb{E}[x^2 y^2] = \left. \frac{x^3}{3} \right|_0^1 \times \left. \frac{y^3}{3} \right|_0^1 = 1/9$$

$$\begin{aligned}\mathbb{E}[(z - \mu_z)^2] &= \mathbb{E}[x^4] + \mathbb{E}[y^4] - 4\mathbb{E}[x^3 y] - 4\mathbb{E}[x y^3] + 6\mathbb{E}[x^2 y^2] - \frac{1}{36} \\ &= \frac{1}{5} + \frac{1}{5} - 4\left(\frac{1}{8}\right) - 4\left(\frac{1}{8}\right) + 6\left(\frac{1}{9}\right) - \frac{1}{36} \\ &= \frac{2}{5} - 1 + \frac{2}{3} - \frac{1}{36} = \frac{7}{180}\end{aligned}$$

$$\boxed{\text{Var}(z) = \mathbb{E}[(z - \mu_z)^2] = \frac{7}{180}}$$

1.25

$$\begin{aligned}c) \quad \mathbb{E}[S] &= \mathbb{E}[z_1 + z_2 + \dots + z_d] \\ &= \mathbb{E}\left[\sum_{i=1}^d z_i\right] = \mathbb{E}\left[\sum_{i=1}^d (x_i - y_i)^2\right] \\ &= \sum_{i=1}^d \mathbb{E}[(x_i - y_i)^2]\end{aligned}$$

We know that $\mathbb{E}[(x_i - y_i)^2] = 1/6$ and that $\mathbb{E}[(x_i - y_i)^2] = \mathbb{E}[(x_j - y_j)^2] \quad \forall i, j \in d$

$$\therefore \mathbb{E}[S] = \sum_{i=1}^d \mathbb{E}[(x_i - y_i)^2] = \sum_{i=1}^d (1/6) = d/6$$

$$\boxed{\mathbb{E}[S] = d \mathbb{E}[z] = d/6}$$

0.5

Similarly, we calculate variance of S

$$\begin{aligned}\text{Var}[S] &= \text{Var}[z_1 + z_2 + \dots + z_d] \\ &= \sum_{i=1}^d \text{Var}(z_i) \quad \left[\text{Since } z_i \text{'s are independent} \right] \\ &= d \text{Var}(z) = 7d/180\end{aligned}$$

(Note $\text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y) + \text{Cov}(X, Y)$)

If X and Y are independent
 $\Rightarrow \text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y)$

0.5

You can further use Markov's inequality to prove that points in higher dimensions are far apart

Markov's inequality says

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

or,

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

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

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

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

$|S - d/6| \rightarrow$ is also a distance

Say $a = 1$

For $d = 1$

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

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

For $d = 5$

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

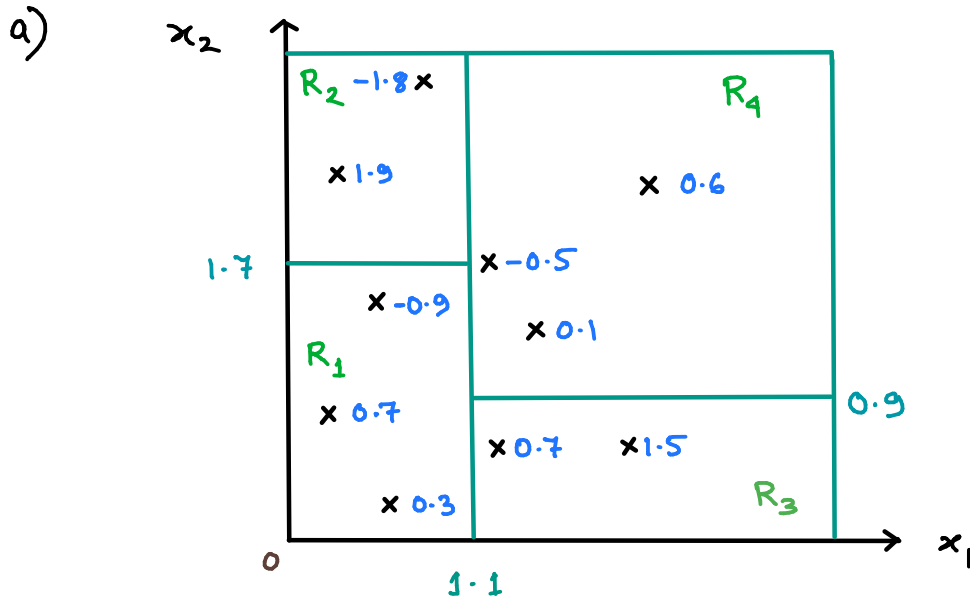
For $d = 10$

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

In 10-D, the chances of
the distance between 2 points
exceeding a certain value
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



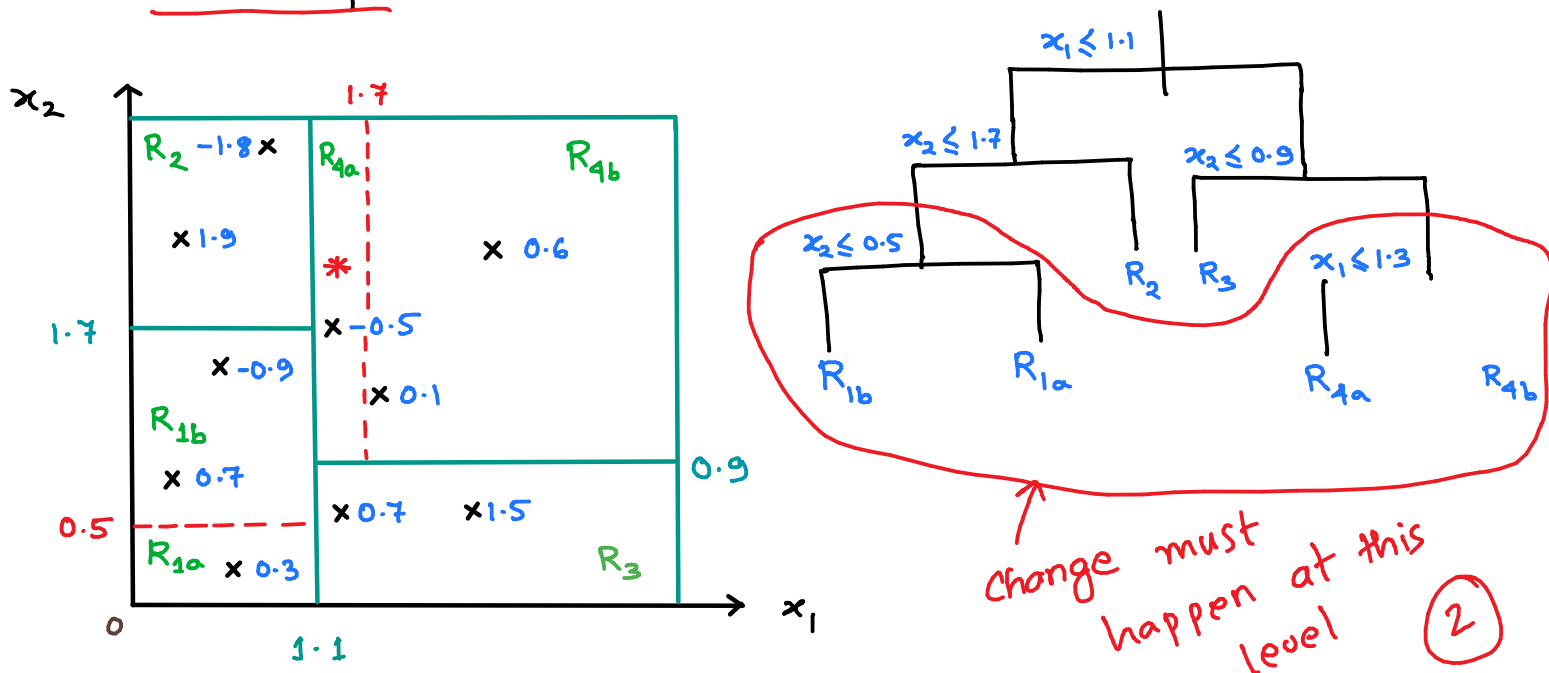
b) Since $x_{1,*} = 1.5 > 1.1$ and $x_{2,*} = 1.8 > 0.9$, the test point belongs to region R_4 .

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

Therefore, the prediction becomes $\hat{y}_* = 0.067$ 0.5

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

One example could be



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

0.5

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

5)

$$a) \quad \underline{y} = \underline{X} \underline{\theta} + \underline{\epsilon}, \quad \underline{\epsilon} \sim \mathcal{N}(\underline{0}, \sigma^2 \mathbf{I}_N)$$

- The likelihood turns out to be Gaussian

\mathbf{I}_N is an identity matrix of size $N \times N$
 N - size of training data

$$P(\underline{y} | \underline{X} \underline{\theta}) = \mathcal{N}(\underline{X} \underline{\theta}, \sigma^2 \mathbf{I}_N)$$

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

- Log-likelihood

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

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

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

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

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

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

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

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

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

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

0.5

\rightarrow In practice, it means some input features are redundant

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

a)
$$\begin{aligned}\frac{dh(x)}{dx} &= \frac{e^x(1+e^x) - e^x \cdot e^x}{(1+e^x)^2} = \frac{e^x(1+e^x - e^x)}{(1+e^x)^2} \\ &= \frac{e^x}{1+e^x} \cdot \frac{1}{1+e^x} \\ &= \left(\frac{e^x}{1+e^x}\right) \cdot \left(1 - \frac{e^x}{1+e^x}\right) \\ &= h(x) \cdot (1 - h(x))\end{aligned}$$

0.5

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

Treat

$$\begin{aligned}p(y=1 | \underline{x}; \underline{\theta}) &= h(\underline{x}^T \underline{\theta}) ; \quad p(y=0 | \underline{x}; \underline{\theta}) = 1 - h(\underline{x}^T \underline{\theta}) \\ &= \frac{e^{\underline{x}^T \underline{\theta}}}{1 + e^{\underline{x}^T \underline{\theta}}} \quad \quad \quad = \frac{1}{1 + e^{\underline{x}^T \underline{\theta}}}\end{aligned}$$

Log-Likelihood for a data pair $\{\underline{x}_i, y_i\}$

$$\ln p(y_i | \underline{x}_i; \underline{\theta}) = \begin{cases} \ln h(\underline{x}_i^T \underline{\theta}) & \text{if } y_i = 1 \\ \ln (1 - h(\underline{x}_i^T \underline{\theta})) & \text{if } y_i = 0 \end{cases}$$

To make the expression more compact, we write

$$\ln p(y_i | \underline{x}_i; \underline{\theta}) = y_i \ln h(\underline{x}_i^T \underline{\theta}) + (1 - y_i) \ln (1 - h(\underline{x}_i^T \underline{\theta}))$$

The log-likelihood for entire training data is

$$\ln p(y_1, \dots, y_N | \underline{x}_1, \dots, \underline{x}_N) = \sum_{i=1}^N y_i \ln h(\underline{x}_i^T \underline{\theta}) + (1 - y_i) \ln (1 - h(\underline{x}_i^T \underline{\theta}))$$

1

$$c) \ln p(y_1, \dots, y_N | x_1, \dots, x_N) = \sum_{i=1}^N y_i \ln h(x_i^T \underline{\theta}) + (1-y_i) \ln (1-h(x_i^T \underline{\theta}))$$

$$\frac{d}{d\underline{\theta}} \left[y_i \ln \underbrace{h(x_i^T \underline{\theta})}_h + (1-y_i) \ln \underbrace{(1-h(x_i^T \underline{\theta}))}_h \right]$$

$$= y_i \frac{1}{h} \left(\frac{dh}{d\underline{\theta}} \right) x_i + (1-y_i) \frac{1}{1-h} \left(-\frac{dh}{d\underline{\theta}} \right) x_i$$

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

$$= y_i (1-h) x_i - (1-y_i) h x_i$$

$$= y_i x_i - \cancel{y_i h x_i} - h x_i + \cancel{y_i h x_i}$$

$$= (y_i - h) x_i$$

$$= \underbrace{(y_i - h(x_i^T \underline{\theta}))}_{\text{residual}} x_i^{(i)}$$

Therefore,

$$\frac{dL}{d\underline{\theta}} = \frac{d}{d\underline{\theta}} \ln p(y | \underline{x}; \underline{\theta}) = \sum_{i=1}^N (y_i - h(x_i^T \underline{\theta})) x_i^T$$

①

d) Differentiating further,

$$\frac{d^2 \ln p(y_i | x_i; \underline{\theta})}{d\underline{\theta} d\underline{\theta}^T} = \frac{d}{d\underline{\theta}^T} (y_i - h(x_i^T \underline{\theta})) x_i$$

$P \times P$
matrix

$$= - \frac{dh}{d\underline{\theta}^T} x_i x_i^T$$

$$\underline{\theta} \in \mathbb{R}^P$$

$$= - h(1-h) \underbrace{x_i x_i^T}_{P \times P}$$

①

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