

Q1

a) $E(f(x)) = \sum_{x \in X} f(x) p(x) = \underbrace{10 \times \frac{1}{10}}_{\text{for } x=a} + \underbrace{5 \times \frac{2}{10}}_{\text{for } x=b} + \underbrace{\frac{10}{7} \times \frac{7}{10}}_{\text{for } x=c} = [3]$

b) $E\left(\frac{1}{p(x)}\right) = \sum_{x \in X} \frac{1}{p(x)} \times p(x) = \frac{1}{p(a)} \times P(a) + \frac{1}{p(b)} \times P(b) + \frac{1}{p(c)} \times P(c) = [3]$

$\hookrightarrow \frac{1}{p(a)} = f(a), \frac{1}{p(b)} = f(b), \frac{1}{p(c)} = f(c)$

$\Rightarrow \frac{1}{p(x)} = f(x) \Rightarrow E\left(\frac{1}{p(x)}\right) = E(f(x))$

c) $p \neq p \Rightarrow p(a) + p(b) + p(c) + \dots + p(n) = 1$

$\Rightarrow E\left(\frac{1}{p(x)}\right) = \sum_{\substack{x \in X \\ x=a}}^n \frac{1}{p(a)} \times P(a) + \dots + \frac{1}{p(n)} \times P(n)$

$= |X|$ \hookrightarrow number of elements in X or
number of elements which x can get

Q2

a)

$$\begin{aligned}
 E(X) &= \underset{\text{defn}}{\sum_{i=1}^m} E(a_i X_1 + a_2 X_2 + \dots + a_m X_m) \\
 &= E(a_1 X_1) + E(a_2 X_2) + \dots + E(a_m X_m) \\
 &= a_1 E(X_1) + a_2 E(X_2) + \dots + a_m E(X_m) \\
 &= \sum_{i=1}^m a_i E(X_i) \quad \text{(I)}
 \end{aligned}$$

- $E(X)$ for gaussian distribution is μ . Why?
because: (σ^2 is \sum_i , σ^2 feels better as variance)

$$\begin{aligned}
 E(X) &= \int_{-\infty}^{+\infty} \frac{1}{(2\pi\sigma^2)^{\frac{1}{2}}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \nu dx \\
 &= \int_{-\infty}^{+\infty} \frac{1}{(2\pi\sigma^2)^{\frac{1}{2}}} e^{-\frac{(z-\mu)^2}{2\sigma^2}} (z+\mu) dz \\
 &= \frac{1}{(2\pi\sigma^2)^{\frac{1}{2}}} \left[\int_{-\infty}^{+\infty} z e^{-\frac{z^2}{2\sigma^2}} dz + \frac{\mu}{(2\pi\sigma^2)^{\frac{1}{2}}} \int_{-\infty}^{+\infty} e^{-\frac{z^2}{2\sigma^2}} dz \right] \\
 &\quad \text{odd function} = 0 \qquad \text{it is pdf} = \underline{\underline{1}}
 \end{aligned}$$

$$\Rightarrow 0 + \mu = \mu \Rightarrow E(X) = \mu \quad \text{(II)}$$

$$\bullet \text{(I) \& (II)} \Rightarrow E(X) = \sum_{i=1}^m a_i E(X_i) = \boxed{\sum_{i=1}^m a_i \mu_i}$$

$a_i \rightarrow$ is a constant

$\mu_i \rightarrow$ is E/\mathbb{R}^d

$\Rightarrow E(X)$ has the dimension of \mathbb{R}^d

Q2

b)

$$\begin{aligned}\text{Cov}(X) &= \sum_{i,j=1}^m \text{Cov}(a_i X_i, a_j X_j) \\ &= \sum_{i,j=1}^m a_i a_j \text{Cov}(X_i, X_j) \\ &= \sum_{i=1}^m a_i^2 \text{Var}(X_i) + \sum_{i \neq j} a_i a_j \text{Cov}(X_i, X_j) \\ &= \sum_{i=1}^m a_i^2 \text{Var}(X_i) + 2 \sum_{1 \leq i < j \leq m} a_i a_j \text{Cov}(X_i, X_j)\end{aligned}$$

if random variables X_1, \dots, X_m are independent
we will get $\text{Cov}(X_i, X_j) = 0$

$$\Rightarrow \text{Cov}(X) = \sum_{i=1}^m a_i^2 \text{Var}(X_i) = \boxed{\sum_{i=1}^m a_i^2 \Sigma_i}$$

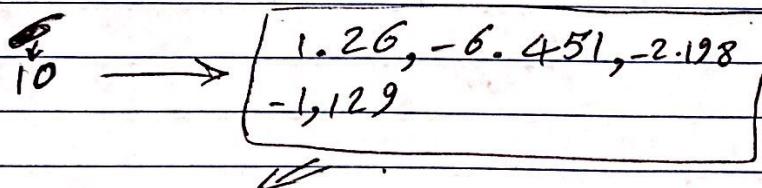
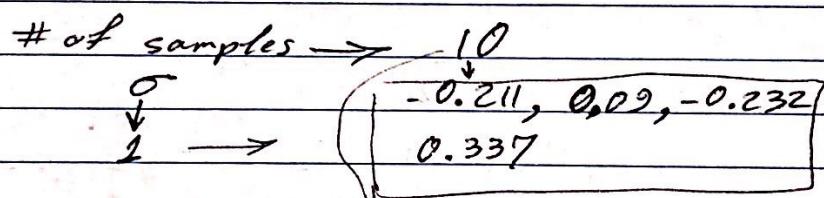
\Rightarrow The dimension of $\text{Cov}(X)$ will be $d \times d$ since
 $\Sigma_i \in \mathbb{R}^{d \times d}$ and a_i is constant

\Rightarrow now if $\text{Cov}(X_1, X_2) = 1$ we will get:

$$\begin{aligned}\text{Cov}(X) &= \sum_{i=1}^m a_i^2 \Sigma_i + 2 \sum_{2 \leq i < j \leq m} a_i a_j \text{Cov}(X_i, X_j) \\ &\quad + 2 \sum_{1 \leq i < j \leq 2} a_i a_j \text{Cov}(X_i, X_j) \\ \Rightarrow \text{Cov}(X) &= \sum_{i=1}^m a_i^2 \Sigma_i + 2 a_1 a_2 \Delta\end{aligned}$$

Q3 a) As I ran the code for this question and with the asked configuration, I saw that as the number of samples increased it tends to get closer to μ which is the 1. Obviously because it should get closer to mean but more samples show the result in the long run.

Now as we changed σ from 1 to 10 it means variance of our samples can be larger, so we saw the same number of samples as $\sigma=1$, for $\sigma=10$ have more space to spread. Therefore, for small number of samples^{in $\sigma=10$} we had more variety of means, in contrast of $\sigma=1$. However, again when we increased number of samples, the samples had less and less space to spread. If we increase σ from 10 to 100 even 1000 samples will find more variety of means.



more variety of numbers because of increase in variance

# of samples $\rightarrow 100$	$\rightarrow 1000$		
$\sigma \rightarrow 1$	0.096, -0.055	-0.013, 0.03, 0.03	
	0.131, 0.212	0.02	
$\sigma \rightarrow 10$	1.017, 0.620 0.203, -0.47	0.16, 0.25, 0.77 0.024	

b) It is basically the covariance matrix of each individual variables:

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} \text{Cov}(X, X) = \text{Var}(X) & \text{Cov}(X, Y) & \text{Cov}(X, Z) \\ \text{Cov}(Y, X) & \text{Cov}(Y, Y) = \text{Var}(Y) & \text{Cov}(Y, Z) \\ \text{Cov}(Z, X) & \text{Cov}(Z, Y) & \text{Cov}(Z, Z) = \text{Var}(Z) \end{bmatrix}$$

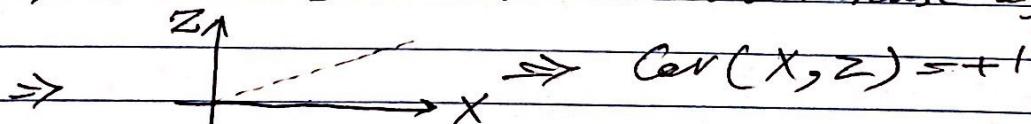
This identity matrix shows X , Y , and Z are independent from each other. Covariance basically shows how two variables vary with each other. All of them are 0. So they don't vary together and they are independent.

c) By this change we will have:

$$\begin{aligned} \text{Cov}(X, X) = \text{Var}(X) &= 1 & \text{Cov}(X, Z) &= 1 & \text{Cov}(Y, Y) = \text{Var}(Y) &= 1 \\ \text{Cov}(Z, X) &= 1 & \text{Cov}(Z, Z) = \text{Var}(Z) &= 1 \end{aligned}$$

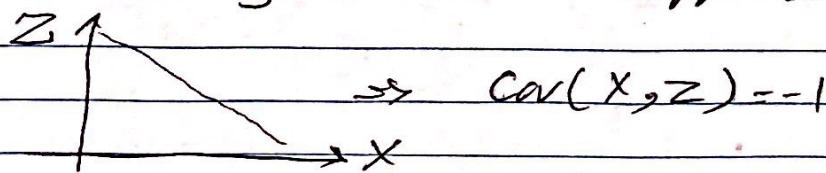
This indicates: $\text{Cov}(X, Z) = \frac{\text{Cov}(X, Z)}{\sqrt{\text{Var}(X)} \sqrt{\text{Var}(Z)}} = 1$

\Rightarrow It means X and Z will increase together



$\text{Cov}(Z, X)$ will do the same as $\text{Cov}(X, Z)$, so the result won't change.

If this $\text{Cov}(X, Z) = -1$, the exact opposite of this diagram would happen.



(24)

a) based on our initial feelings, since it is an exponential distribution with $P(x|λ) = \frac{1}{2} e^{-\frac{1}{2}λ}$, we can assume

$$P(λ) = \frac{1}{2} e^{-\frac{1}{2}λ}. \text{ It looks like } \begin{array}{c} \nearrow \\ \curvearrowleft \end{array} \rightarrow$$

which easily has the highest point at $λ=0$.
 Now, even with maximum, $\frac{d f(λ)}{d λ} = -\frac{1}{4} e^{-\frac{1}{2}λ} = 0$
 $\Rightarrow f(0) = \frac{1}{2}$

b) $λ_{ML} = \operatorname{argmax} \{ P(D|X) \}$

$$\rightarrow P(D|X) = P(X_1, n_2, n_3, \dots, n_n | λ)$$

$$= \prod_{i=1}^n P(n_i | λ) = \prod_{i=1}^n \frac{λ^{n_i} e^{-λ}}{n_i!}$$

$$= \prod_{i=1}^n \frac{λ^{n_i} e^{-λ}}{n_i!} = \frac{λ^{\sum_{i=1}^n n_i} e^{-nλ}}{\prod_{i=1}^n n_i!}$$

$$\Rightarrow -\ln(P(D|λ)) = -\ln \frac{λ^{\sum_{i=1}^n n_i} e^{-nλ}}{\prod_{i=1}^n n_i!} = \ln \lambda \sum_{i=1}^n n_i - nλ - \sum_{i=1}^n \ln(n_i!)$$

$$\Rightarrow \frac{\partial \ln P(D|λ)}{\partial λ} = \frac{1}{λ} \sum_{i=1}^n n_i - n = 0$$

$$\Rightarrow λ = \frac{\sum_{i=1}^n n_i}{n} \rightarrow 9 \Rightarrow λ = \frac{79}{9}$$

$$c) \lambda_{MAP} = \arg \max \{ P(D|\lambda) P(\lambda) \}$$

we have showed that $\textcircled{I} P(D|\lambda) = \lambda^{\sum_{i=1}^n x_i - \lambda n} e^{-\lambda n}$

$$P(\lambda) = \theta e^{-\theta \lambda} \xrightarrow{\theta = \frac{1}{2}} \frac{1}{2} e^{-\frac{\lambda}{2}}$$

$$(P(D|\lambda) P(\lambda)) \textcircled{I} \times \textcircled{II} = \frac{1}{2} \frac{\lambda^{\sum_{i=1}^n x_i - \lambda n} e^{-\lambda/2 - \lambda n}}{\prod_{i=1}^n x_i!}$$

$$\ln(P(D|\lambda) P(\lambda)) \rightarrow \frac{1}{2} (\ln \lambda \sum_{i=1}^n x_i - \frac{\lambda}{2} - \lambda n - \sum_{i=0}^n \ln(x_i!))$$

$$\Rightarrow \frac{\partial (\ln(P(D|\lambda) P(\lambda)))}{\partial \lambda} = \frac{1}{2} \left(\frac{1}{\lambda} \sum_{i=1}^n x_i - \frac{1}{2} \right) = 0$$

$$\Rightarrow \lambda = \frac{2 \sum_{i=1}^n x_i}{2n+1} = \boxed{\frac{158}{19}}$$

d) by using ML and MAP we get $\arg \max_{ML} \lambda$

and $\arg \max_{MAP} \lambda$. As these are coming

from Poisson distribution, the arg max will be

the mode of these distributions which is 2.1.

both of MAP and ML are predicting 8 incidences

see next page

(~~$\frac{158}{19} \approx 8$ and $2.1 \approx 2$~~) for tomorrow. ML and

MAP are not always the same, but this time

both are predicting 8 accidents. (continue next page)

more precise way next page

We put the numbers into the distribution for lot of ML and MAP to find the most probability of accidents.

$$R_{ML} = \frac{79}{9} \times e^{-\frac{79}{9}} / 8! \quad \begin{matrix} x=8 \\ x=9 \end{matrix}$$

$$\Rightarrow \left(\frac{79}{9} \right)^8 \times e^{-\frac{79}{9}} / 8! = 0.1347 \rightarrow \text{more probable}$$

$$\Rightarrow \left(\frac{79}{9} \right)^9 \times e^{-\frac{79}{9}} / 9! = 0.1313$$

ML
→ number of accidents: 8

$$R_{MAP} = \frac{158}{19}$$

$$\Rightarrow \left(\frac{158}{19} \right)^8 \times e^{-\frac{(158)}{19}} / 8! = 0.1387$$

$$\Rightarrow \left(\frac{158}{19} \right)^9 \times e^{-\frac{(158)}{19}} / 9! = 0.1281$$

MAP
→ number of accidents: 8

e) We mainly use prior knowledge to enhance our predictions estimation.

In this example, we used this information that λ is exponential distribution.

The prior information will be less and less effective when our samples increase.

In fact ML is as MAP with uniform distribution for λ , an information which is not always true.

If we don't have enough data, these information helps us to increase our accuracy of estimation to avoid overfitting.

f) As our goal is to decrease the number of

accidents, easily we have to increase θ in the prior. There are two ways to prove it.

First is, as we know, increasing θ means decreasing the mean of the distribution

\Rightarrow estimated λ will decrease as well

Second, you can see by $\frac{\sum k_i}{n+\theta}$ in MAP

that increasing θ means decreasing the statement(λ) and vice versa.

Q25 a)

We have 2 random variables. Let's call "is sunny" as X and "is table free" as Y . For writing the ML for maximizing $P(D|X)$ we need to come up with a distribution. Obviously for "D" we are dealing with Bayes rule. So the problem will be

$$\lambda_{ML} = \arg \max p(D|\lambda)$$

$$\hookrightarrow p(X_i, Y_i | \lambda) = P(Y_i | X_i, \lambda) P(X_i)$$

For distribution we can use Bernoulli distribution for being sunny and not sunny as win and lose respectively. The same is applicable for vacancy of table.

We also could use ^{marginal} pmf of for example "being sunny" was ^{being} related to "not empty table" and find another distribution related to that. But we don't have those information to compute exactly right now.

b) Now that we have Maximum likelihood for some number of data (10 days), we can use " λ " as sunny parameter. It means we can put " $\lambda=1$ " in our distribution and find its value (mode) from ML technique. It means we have to put our problem equal to zero after derivation of y (emptiness of a table). Now, there ~~is~~ is an approach to estimate y since y can only be 0 and 1. We can use a threshold and say if the value that we take is more than the threshold or not. This threshold can be guessed by using the ~~mean~~ ^{threshold} value (~~in this case~~ if it is not biased the ^{it} mean would be 0.5%). In any case with any threshold, we can easily say if $|y|$ is more than threshold, y is equal to 1 which means table is available, otherwise it is not available.

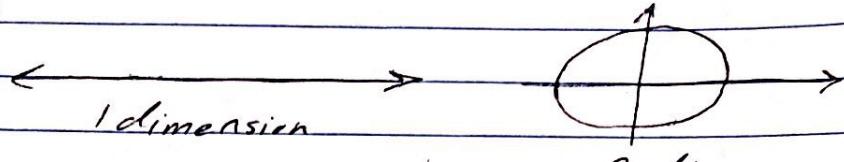
C) Now we have to add another variable with name Z . We can use three different variables for different times of day, but it is not suggested since it cannot be "morning", "afternoon", or "evening" at the same time. So we use a Z with values of 0 for morning, 1 for afternoon, and 2 for evening. Now it means based on Z and λ we want to estimate Y (vacancy of table). Again we have to write down the formula of Bayes which will be used for $P(D|Z)$:

$$P(D|Z) = P(X_i, Y_i, Z_i | \lambda) = P(X_i, Y_i | Z_i, \lambda) P(Z_i)$$

$$= P(Y_i | X_i, Z_i, \lambda) P(X_i) P(Z_i)$$

Again like part "a" and b we have to solve ML problem by using derivative and equalizing with 0.

(Q5) d) as number of dimensions increase, obviously the mean of distance also increases. It is because of the fact that with the same variance and no covariance between each of the two dimensions, we have way more space to spread our point. for example:



for 1 dimension which is a line, we have to spread 1000 points. Now for a circle (2d) we have more space to spread 1000 points. for a sphere (3d) we have again more space to spread the points. As we can see, from each point to center of our shape we have same variance and more space causes more chance for points to spread far away from center. (also $\sqrt{x^2} < \sqrt{x^2+y^2} < \sqrt{x^2+y^2+z^2} < \dots$) Consequently, for K-means, if number of classes stay the same we have to increase its range because for example for 1D we have 0.78 as mean and for 256D we have 16.008 as mean. We cannot cover all numbers (if the number of classes and their range stays the same) for 256D as it will be covered for 1D. So, either the number of classes should increase or the range of the classes that will cover the points should increase or both. It is happening while we have same number of points. If now we increase number of points as we increase the dimensions, still we have to do the same because even if we increase the samples and afterward the mean decreases for higher dimension to the exact mean of lower dimensions, the biggest number of each sample has been increased as dimension increased. So increasing range of classes or number of classes for higher dimensions is the only way. (PICTURE)

now based on K-means definition, we have ~~some~~
an specified number of clusters each with a center
which is the mid point. When we introduce another dimension
it will increase the distance of each point to the
center of the cluster because obviously $\sqrt{x^2 + y^2}$ ~~for x^2 + y^2~~
which means distance will increase as D increases.
Later it will make the points to find the nearest
center of cluster and then the distance of the
center of cluster to origin point will increase. It
causes, at the end, ~~to~~ the ^{distance of} center of clusters
to the origin point to increase, the range of distances
that a cluster supports to increase, and some
points change their clusters as the number of
dimensions increase.

Q5

e) This question means we have two d-dimensional hypercubes inside of each other. One with $\frac{1}{\epsilon}$ side length and one with $(1-\epsilon)$ side length which means the smaller one is inside of the larger one. The volume of the bigger hypercube is 1^d and its π for the smaller one is $(1-\epsilon)^d \rightarrow \frac{1^d}{(1-\epsilon)^d}$ is the ratio of the ~~second~~ first hypercube to the volume of the second one or in the inverse order it will be $\frac{(1-\epsilon)^d}{1^d}$.

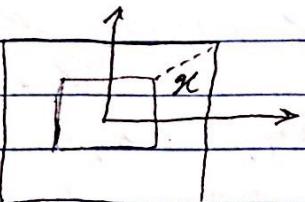
As d gets larger 1^d will ~~stay~~ 1 and $(1-\epsilon)^d$ will get be smaller since $0 < 1-\epsilon < 1$. So the volume of the π hypercube decreases.

It means that $\sqrt[d]{\text{distance of each point on the corner of bigger hyper cube to the corresponding corner of the smaller hypercube}}$ will increase. It causes some points in the smaller hypercube falls into the bigger hypercube.

If the points are on the surface of the hypercube the distance between each two of them from bigger hypercube to the smaller one will increase. (if we assume each hypercube is our samples)

Now if the points should spread inside of the hypercube as the volume decreases ~~distance~~ of each angle to the center will decrease (for the smaller hypercube). But for the larger it will stay the same. At the same time by adding a dimension we are adding a new variable to our distance (e.g. $\sqrt{x^2} < \sqrt{x^2+y^2} < \sqrt{x^2+y^2+z^2} < \dots$) so distance will increase. It will make the density to increase. ~~however cannot say much about this~~ We can see it mathematically in the next page.

for example for 2D we have


$$x = \sqrt{\left(\frac{1-(1-\alpha)}{2}\right)^2 + \left(\frac{1-(1-\alpha)}{2}\right)^2}$$
$$x = \sqrt{2\alpha^2} = \sqrt{2} \alpha$$

for 3D:

$$x = \sqrt{3\alpha^2} = \sqrt{3} \alpha$$

for d dimensions:

$$x = \sqrt{d\alpha^2} = \text{vol } \sqrt{\alpha^2} = \sqrt{d} \alpha$$

→ as you can see, as d increases
it gets larger and the volume between
larger hypercube and the smaller one increases
so it proves that as dimensions increase,
each two points' ~~dist~~ distance will also
increase.