
Instructions:

- This assignment is meant to help you grok certain concepts we will use in the course. Please don't copy solutions from any sources.
 - Avoid verbosity.
 - The assignment needs to be typeset in latex using the attached tex file. The solution for each question should be written in the solution block in space already provided in the tex file. **Handwritten assignments will not be accepted.**
 - Deadline for submission is **11:55PM 1/4/2018**.
-

1. **Independence of Random Variables**

A and B are two random variables which can take values 0 or 1. Two joint probability distributions over A and B are provided in the tables below. For each case, argue whether A and B are independent.

Table 1: (a)

	A=0	A=1
B=0	0.12	0.18
B=1	0.28	0.42

Table 2: (b)

	A=0	A=1
B=0	0.20	0.18
B=1	0.28	0.34

Solution:

1a) From table we get that

$$P(A=0)=0.12+0.28=0.40$$

$$P(A=1)=0.18+0.42=0.60$$

$$P(B=0)=0.12+0.18=0.30$$

$$P(B=1)=0.28+0.42=0.70$$

In case of $A=0, B=0$

by definition of conditional Probability

$$\begin{aligned}
 P(A_0|B_0) &= \frac{P(A_0, B_0)}{P(B_0)} \\
 &= \frac{0.12}{0.3} \\
 &= 0.4
 \end{aligned}$$

also

$$P(A_0) = 0.4 \text{ (calculated)}$$

as

$$P(A_0|B_0) = P(A_0)$$

so we can conclude

$$\boxed{P(A_0) \text{ independent of } P(B_0) \text{ or } P(A_0 \perp B_0)}$$

In case of A=1,B=0

$$\begin{aligned}
 P(A_1|B_0) &= \frac{P(A_1, B_0)}{P(B_0)} \\
 &= \frac{0.18}{0.3} \\
 &= 0.6 \\
 &= P(A_1) \text{ (calculated)}
 \end{aligned}$$

so we can conclude

$$\boxed{P(A_1 \perp B_0)}$$

In case of A=0,B=1

$$\begin{aligned}
 P(A_0|B_1) &= \frac{P(A_0, B_1)}{P(B_1)} \\
 &= \frac{0.28}{0.7} \\
 &= 0.4 \\
 &= P(A_0) \text{ (calculated)}
 \end{aligned}$$

so we can conclude

$$P(A_1 \perp B_0)$$

In case of A=0,B=1

$$\begin{aligned} P(A_1|B_1) &= \frac{P(A_1, B_1)}{P(B_1)} \\ &= \frac{0.42}{0.7} \\ &= 0.6 \\ &= P(A_1) \text{ (calculated)} \end{aligned}$$

so we can conclude

$$P(A_1 \perp B_1)$$

So we conclude $P(A)$ is independent of $P(B)$ or $P(A \perp B)$

1b) From table we get that

$$P(A=0)=0.20+0.28=0.48$$

$$P(A=1)=0.18+0.34=0.52$$

$$P(B=0)=0.20+0.18=0.38$$

$$P(B=1)=0.28+0.34=0.62$$

In case of A=0,B=0

by definition of conditional Probability

$$P(A_0|B_0) = \frac{P(A_0, B_0)}{P(B_0)} = \frac{0.20}{0.38} = 0.5263$$

also

$$P(A_0) = 0.48$$

as

$$P(A_0|B_0) \neq P(A_0)$$

so we can conclude

$$P(A_0) \text{ dependent of } P(B_0) \text{ or } P(A_0 \not\perp B_0)$$

As $P(A_0 \not\perp B_0)$ we conclude $P(A)$ is dependent of $P(B)$ or $P(A \not\perp B)$

2. Ram is trying to study the causes of aggressive behaviour in males. For his initial experiments, he decides to take into account two parameters, namely, the basal level of testosterone in the male (high or low) and the kind of neighbourhood he grew up in (violent/non-violent). Based on a survey of males in a city that he conducted, he estimated that 80% of the males grew up in non-violent neighbourhoods. He also gathered the following posteriors

Neighbourhood	Testosterone		Testosterone	Neighbourhood	Aggression	
	High	Low			High	Low
Violent	0.7	0.3	High	Violent	0.75	0.25
Non-Violent	0.4	0.6	High	Non-Violent	0.22	0.78
			Low	Violent	0.60	0.40
			Low	Non-violent	0.15	0.85

What is the probability that

- (a) A male who grew up in a non-violent neighbourhood is highly aggressive.

Solution: Let following notation are used as abbreviation

Non-Violent=NV , Violent=V Aggressive=A , High=H , Low=L, Testosterone=T

$P(N_{NV}) = P(N = NV)$ etc

$P(N = NV) = 0.8$ and $P(N = V) = 0.2$ (given)

$$\begin{aligned}
 P(A_H|N_{NV}) &= \frac{P(A_H, N_{NV})}{P(N_{NV})} \\
 &= \frac{P(A_H, N_{NV}, T_H) + P(A_H, N_{NV}, T_L)}{P(N_{NV})} \\
 &= \frac{P(A_H|N_{NV}, T_H) * P(T_H|N_{NV}) * P(N_{NV})}{P(N_{NV})} \\
 &\quad + \frac{P(A_H|N_{NV}, T_L) * P(T_L|N_{NV}) * P(N_{NV})}{P(N_{NV})} \\
 &= \frac{0.15 * 0.6 * 0.8 + 0.22 * 0.4 * 0.8}{0.8} \\
 &= \mathbf{0.178}
 \end{aligned}$$

- (b) An arbitrarily chosen male who is highly aggressive, has high levels of testosterone and grew up in a non-violent neighbourhood.

Solution:

$$\begin{aligned}
 P(A_H, N_{NV}, T_H) &= P(A_H|N_{NV}, T_H) * P(T_H|N_{NV}) * P(N_{NV}) \\
 &= 0.22 * 0.4 * 0.8 \\
 &= \mathbf{0.0704}
 \end{aligned}$$

3. Game of Diamonds

You are playing a game in which you have an opportunity to win diamonds. You are shown three identical boxes, one of which contains diamonds and the other two boxes are empty. The game proceeds as follows:

- You choose one box, which you think might contain diamonds.
- Among the remaining boxes, either one or both are empty. The game host opens one such empty box.
- Now you have two options: stick to the choice you made earlier, or choose the other box. Depending on the option you choose, you win or lose.

Which option will you choose in the last step and why? (Hint: compute probability of winning in both cases)

Solution: I like to Switch as By **intuition** choosing any box $P(B) = \frac{1}{3}$ and that box has diamond, so if I choose to switch I fail $\frac{1}{3}$ of the time ,but in case of it is empty if I switch to any other box I have chance to win $\frac{2}{3}$ of the time.

Mathematically: Let S be the event of finding the Diamond by switching, and let B_i be the event that the Diamond is in Box i. Then by the Law of Total Probability, $P(S)$ = Probability of Success by switching
 $P(B_i)$ = Probability of Having the diamond in box i where $i \in [1, 2, 3]$

$$\begin{aligned} P(S) &= P(S|D_1) * P(B_1) + P(S|B_2) * P(B_2) + P(S|B_3) * P(B_3) \\ &= P(S|B_1) * \frac{1}{3} + P(S|B_2) * \frac{1}{3} + P(S|B_3) * \frac{1}{3} \end{aligned}$$

assuming that we first choose box B_1

$$= 0 * \frac{1}{3} + P(S|B_2) * \frac{1}{3} + P(S|B_3) * \frac{1}{3}$$

As we choose Box B_1 and switching to again is impossible event so $P(S|B_1) = 0$

$$= 0 + 1 * \frac{1}{3} + 1 * \frac{1}{3}$$

As we choose Box B_1 and switching to B_2 and it has the Diamond is $P(S|B_2) = 1$, similar in case of $P(S|B_3) = 1$

$$= \boxed{\frac{2}{3}}$$

So By switching I have always better chance to win the game

4. Sampling from continuous distributions

You are given a random number generator R , which gives a real number output sampled from the probability distribution $U_{0,1}$. $U_{a,b}$ is defined as:

$$U_{a,b}(x) = \begin{cases} \frac{1}{b-a} & a < x < b \\ 0 & \text{otherwise} \end{cases}$$

(a) Use R to sample from $U_{a,b}$ for $a, b \in \mathbb{R}$ ($a < b$).

Solution: let R_{ab} be random number generator which generate x and $x \in (a, b)$

$$\begin{aligned} a < x < b &\implies 0 < x - a < b - a \\ &\implies 0 < \frac{x - a}{b - a} < 1 \end{aligned}$$

now

$$\begin{aligned} \frac{R_{ab}(x) - a}{b - a} &= R(x) \\ \implies \boxed{R_{ab}(x) &= R(x)(b - a) + a} \end{aligned}$$

(b) Given a probability distribution $P_X(x)$, whose PDF is given by $F_X(x)$ and CDF is given by $C_X(x)$, using a random number generated using R , how will you obtain a sample from the distribution $P_X(x)$?

Solution: Generate a Random No $R \in (0, 1)$.
Find a point x such that

$$F_X(x) = R$$

$$x = F_x^{-1}(R)$$

x is random number sampled from PDF $P_x(x)$

5. Consider the random variables X, Y, Z, W which take 3, 4, 4, 2 values respectively.

- (a) Consider a joint distribution P_1 over these 4 variables. Without any information about the (in)dependencies between the variables, what is the minimum number of parameters you will need to represent this distribution?

Solution: Total No of Parameters = $3 \times 4 \times 4 \times 2 - 1 = 95$ as the 96th can be calculated from other 95 Parameters.

in other way joint probability

$$P(X, Y, W, Z) = P(X|Y, Z, W) * (Y|Z, W) * (Z|W) * (W)$$

so

$$\begin{aligned} Parameters(P(X, Y, W, Z)) &= Parameters(P(X|Y, Z, W)) \\ &+ Parameters(P(Y|Z, W)) \\ &+ Parameters(P(Z|W)) \\ &+ Parameters(P(W)) \end{aligned}$$

$$Parameters(P(X|Y, Z, W)) = (3 - 1) * 4 * 4 * 2 = 64$$

$$Parameters(P(Y|Z, W)) = (4 - 1) * 4 * 2 = 24$$

$$Parameters(P(Z|W)) = (4 - 1) * 2 = 6$$

$$Parameters(P(W)) = 2 - 1 = 1$$

$$\text{Total Parameters} = 64 + 24 + 6 + 1 = 95$$

- (b) An insight into the variables now reveals the information that $(X \perp W|Z)$. What is the minimum number of parameters needed to represent this distribution in this case?

Solution:

$$P(X, Y, W, Z) = P(X|Y, Z, W) * (Y|Z, W) * (Z|W) * (W)$$

which Simplifies given the information $(X \perp W|Z)$

$$\begin{aligned}
 &= P(X|Y, Z, \cancel{W}) * (Y|Z, W) * (Z|W) * (W) \\
 &= P(X|Y, Z) * (Y|Z, W) * (Z|W) * (W)
 \end{aligned}$$

SO Total No Of Parameters

$$\begin{aligned}
 Parameters(P(X, Y, W, Z)) &= Parameters(P(X|Y, Z)) \\
 &+ Parameters(P(Y|Z, W)) \\
 &+ Parameters(P(Z|W)) \\
 &+ Parameters(P(W))
 \end{aligned}$$

$$Parameters(P(X|Y, Z)) = (3 - 1) * 4 * 4 = 32$$

$$Parameters(P(Y|Z, W)) = (4 - 1) * 4 * 2 = 24$$

$$Parameters(P(Z|W)) = (4 - 1) * 2 = 6$$

$$Parameters(P(W)) = 2 - 1 = 1$$

$$\mathbf{Total\ Parameters} = 32 + 24 + 6 + 1 = \boxed{63}$$

- (c) An oracle further tells you that $(Y \perp X|Z, W)$. What is the minimum number of parameters needed to represent this distribution in this case?

Solution:

$$P(X, Y, W, Z) = P(X|Y, Z, W) * (Y|Z, W) * (Z|W) * (W)$$

which Simplifies given the information $(Y \perp X|Z, W)$

$$= P(X|\cancel{Y}, Z, W) * (Y|Z, W) * (Z|W) * (W)$$

we all ready have the information $(X \perp W|Z)$ which further simplifies

$$\begin{aligned}
 &= P(X|Z, \cancel{W}) * (Y|Z, W) * (Z|W) * (W) \\
 &= P(X|Z) * (Y|Z, W) * (Z|W) * (W)
 \end{aligned}$$

SO Total No Of Parameters

$$\begin{aligned} \text{Parameters}(P(X, Y, W, Z)) &= \text{Parameters}(P(X|Z)) \\ &\quad + \text{Parameters}(P(Y|Z, W)) \\ &\quad + \text{Parameters}(P(Z|W)) \\ &\quad + \text{Parameters}(P(W)) \end{aligned}$$

$$\begin{aligned} \text{Parameters}(P(X|Y, Z)) &= (3 - 1) * 4 = 8 \\ \text{Parameters}(P(Y|Z, W)) &= (4 - 1) * 4 * 2 = 24 \\ \text{Parameters}(P(Z|W)) &= (4 - 1) * 2 = 6 \\ \text{Parameters}(P(W)) &= 2 - 1 = 1 \\ \text{Total Parameters} &= 8 + 24 + 6 + 1 = \boxed{39} \end{aligned}$$

6. The students of a college have the option of choosing between two mess caterers, namely, Fake Foods (FF) and Terrible Taste (TT). At the end of each month, each student needs to pick his/her caterer of choice. Based on past experience, we know that 80% of the students choosing FF opt to continue eating in FF for the next month, while 40% of the students eating in TT choose to opt for FF in the subsequent month.
- (a) If 50% of the students are assigned to FF in the first month and the rest for TT, what fraction of students are in FF at the start of the 4th month?

Solution: Let us assume that starting of the month K_1 , student assign to Fake Foods (FF) and K_2 to Terrible Taste (TT), Let after First month fraction stay with FF and 1-p goes with TT, and similar case q fraction stay with TT and 1-q, so after one month students in FF = $pk_1 + (1 - q)k_2$ and in TT = $qk_2 + (1 - p)k_1$. We can write it as

$$\begin{bmatrix} p & 1 - q \\ 1 - p & q \end{bmatrix} \cdot \begin{bmatrix} k_1 \\ k_2 \end{bmatrix}$$

so after three month this will become

$$\begin{bmatrix} p & 1 - q \\ 1 - p & q \end{bmatrix}^3 \cdot \begin{bmatrix} k_1 \\ k_2 \end{bmatrix}$$

so putting $p = 0.8$ $q = 0.6$ $k_1 = 0.5$ and $k_2 = 0.5$ in the above matrix we get

that

$$\begin{aligned}
 \begin{bmatrix} 0.8 & 0.4 \\ 0.2 & 0.6 \end{bmatrix}^3 \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} &= \begin{bmatrix} 0.8 \cdot 0.8 + 0.4 \cdot 0.2 & 0.8 \cdot 0.4 + 0.4 \cdot 0.6 \\ 0.2 \cdot 0.8 + 0.6 \cdot 0.2 & 0.2 \cdot 0.4 + 0.6 \cdot 0.6 \end{bmatrix} \begin{bmatrix} 0.8 & 0.4 \\ 0.2 & 0.6 \end{bmatrix} \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} \\
 &= \begin{bmatrix} 0.72 & 0.56 \\ 0.28 & 0.44 \end{bmatrix} \begin{bmatrix} 0.8 & 0.4 \\ 0.2 & 0.6 \end{bmatrix} \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} \\
 &= \begin{bmatrix} 0.688 & 0.624 \\ 0.312 & 0.376 \end{bmatrix} \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} \\
 &= \begin{bmatrix} 0.656 \\ 0.344 \end{bmatrix}
 \end{aligned}$$

So at the start of the 4th month **65.6%** students will be in FF and **34.4%** students will be in TT mess

- (b) Does the fraction of students eating in FF converge to a certain value? If yes, what is the value?

Solution: The fraction of student in FF converge to certain value after 12 month as

$$\begin{bmatrix} 0.8 & 0.4 \\ 0.2 & 0.6 \end{bmatrix}^{12} \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} = \begin{bmatrix} 0.66666\dots \\ 0.33334\dots \end{bmatrix}$$

so after 12 month **66.66%** people will be in FF mess

- (c) Repeat part (a) with 60% of students assigned to TT for the first month. What is the answer to part (b) in this case? If it converges, does it converge to the same value or is it different? Justify your answer.

Solution: Part (a)

$$\begin{bmatrix} 0.688 & 0.624 \\ 0.312 & 0.376 \end{bmatrix} \begin{bmatrix} 0.4 \\ 0.6 \end{bmatrix} = \begin{bmatrix} 0.6496 \\ 0.3504 \end{bmatrix}$$

So at the start of the 4th month **64.96%** students will be in FF and **35.04%** students will be in TT mess.

Part (b) In case **60%-40%** split of student after 12 month as

$$\begin{bmatrix} 0.8 & 0.4 \\ 0.2 & 0.6 \end{bmatrix}^{12} \begin{bmatrix} 0.4 \\ 0.6 \end{bmatrix} = \begin{bmatrix} 0.66667\dots \\ 0.33333\dots \end{bmatrix}$$

will be **66.66%** in FF mess and **33.33%** in TT mess.

Fact is that when we perform repeated multiplication by the same matrix, we are

really repeatedly scaling the coefficients corresponding to the eigenvector basis by the size of the eigenvalue. When we do this a large number of times, the largest eigenvalue/eigenvector will come to dominate irrespective of the input vector. So what will be the starting student split after 12 month it will always converges to **66.66%** in FF and **33.33%** in TT.

7. A Markov Chain is a discrete time stochastic process. It consists of N states and is characterized a $N \times N$ transition probability matrix P , whose entries lie in the interval $[0, 1]$, entries in each row adding up to 1. The entry P_{ij} is the probability of the state in the next time step being j , given that the state at the current time step is i .

Consider a Markov Chain with x states and the following transition probability matrix:

$$P = \begin{bmatrix} 0 & 0.7 & 0.3 & 0 & 0 & 0 & 0 & 0 \\ 0.4 & 0.6 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0.8 & 0 & 0 & 0 & 0.2 \\ 0 & 0 & 0.5 & 0 & 0.5 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0.1 & 0 & 0.9 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \end{bmatrix}$$

- (a) Define $p_{ij}^{(n)}$ as the probability of reaching state j in n steps starting from state i . Calculate $p_{12}^{(3)}$ and $p_{22}^{(3)}$.

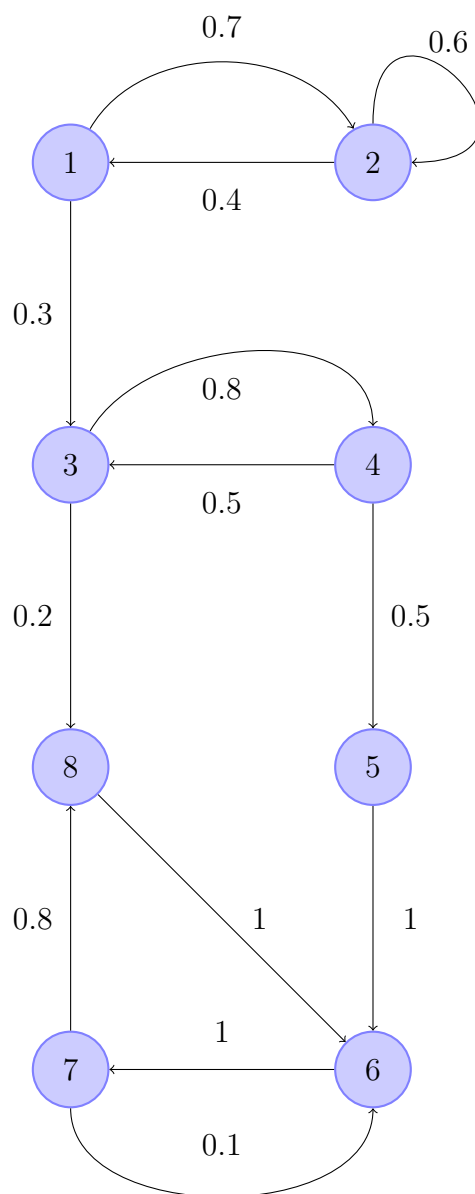
Solution:

$$P^3 = \begin{bmatrix} 0.168 & 0.448 & 0.204 & 0 & 0.12 & 0.06 & 0 & 0 \\ 0.256 & 0.552 & 0.072 & 0.096 & 0 & 0 & 0 & 0.024 \\ 0 & 0 & 0 & 0.32 & 0 & 0.4 & 0.2 & 0.08 \\ 0 & 0 & 0.2 & 0 & 0.2 & 0.1 & 0.5 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0.1 & 0 & 0.9 \\ 0 & 0 & 0 & 0 & 0 & 0.9 & 0.1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0.01 & 0.9 & 0.09 \\ 0 & 0 & 0 & 0 & 0 & 0.1 & 0 & 0.9 \end{bmatrix}$$

so From above Matrix we get $p_{12}^{(3)} = \mathbf{0.448}$ and $p_{22}^{(3)} = \mathbf{0.552}$.

- (b) The period d_i of a state i is defined as $d_i = \gcd\{n \geq 1 : p_{ii}^{(n)} > 0\}$. If $p_{ii}^{(n)} = 0 \forall n \geq 1$, $d_i = \inf$. Find the period of each state in the Markov Chain.

Solution:



using the equation $d_i = \gcd\{n \geq 1 : p_{ii}^{(n)} > 0\}$ we calculate the periods as

following

$$\begin{aligned}
 d_1 &= \gcd(2, 3, 4, \dots) = 1 \\
 d_2 &= \gcd(1, 2, 3, 4, \dots) = 1 \\
 d_3 &= \gcd(2, 4, 6, \dots) = 2 \\
 d_4 &= \gcd(2, 4, 6, \dots) = 2 \\
 d_5 &= \inf \text{ because } P_{55}^n = 0 \\
 d_6 &= \gcd(2, 3, 4, 5, \dots) = 1 \\
 d_7 &= \gcd(2, 3, 4, 5, \dots) = 1 \\
 d_8 &= \gcd(3, 5, 7, \dots) = 1
 \end{aligned}$$

- (c) The state j is said to be accessible from state i if $p_{ij}^{(n)} > 0$ for some n . States i and j are said to communicate if they are accessible from each other. Show that communication is an equivalence relation.

Solution: Reflexive : All the states communicate with themselves : $P_{ii}^0 = 1 > 0$

Symmetry: If $P_{ij}^n > 0$, then $P_{ji}^n > 0$

Transitivity: If $P_{ik}^n > 0$ and $P_{kj}^m > 0$, then $P_{ij}^{(m+n)} \geq P_{ik}^n P_{kj}^m > 0$
as all the above three relations satisfy then the communication is an equivalence relation

- (d) A Markov Chain can be partitioned into classes based on the communication relation defined previously. For the given Markov Chain, find all the equivalence classes. Out of these, which classes are aperiodic (i.e. $d_i = 1 \forall i$ in the class)?

Solution: From a state all those states which are reachable from that state form an equivalence class.

The following are the equivalence classes we get as follows :

$\{1, 2\}$ $\{3, 4\}$ $\{6, 7, 8\}$

among those classes aperiodic $\{d_i = 1 \forall i\}$ are $\{1, 2\}$ and $\{6, 7, 8\}$

8. $\mathbb{X} = \{X_1, \dots, X_N\} \in \Lambda^N$ is a multivariate random variable, with $x_i \in \Lambda, \forall i \in \{1, 2, \dots, N\}$. The possible values taken by the samples (\mathbf{x}) of \mathbb{X} can be thought of as a state of a Markov Chain (refer Question 7) with Λ^N states. Consider such a Markov Chain with transition probability between states defined as

$$p_{\mathbf{xy}} = \begin{cases} q(i)\pi(y_i|(x_v)_{v \in \{1, \dots, N\} \setminus \{i\}}), & \text{if } \exists i \in \{1, \dots, N\} \text{ such that } \forall j \in 1, \dots, N \text{ with } j \neq i, x_j = y_j \\ 0, & \text{otherwise} \end{cases}$$

where q is a density function over the indices $\{1, \dots, N\}$ and π is a joint distribution over $\{X_1, X_2, \dots, X_N\}$ (you can think of π as the current state probabilities). Show that for this Markov Chain, the following condition (called the detailed balance condition) is satisfied

$$\pi(\mathbf{x})p_{\mathbf{x}\mathbf{y}} = \pi(\mathbf{y})p_{\mathbf{y}\mathbf{x}}$$

$\forall \mathbf{x}, \mathbf{y} \in \Lambda^N$.

Hint: Prove it separately for the cases where i) $\mathbf{x} = \mathbf{y}$, ii) \mathbf{x} and iii) \mathbf{y} differ in only one variable and ii) \mathbf{x} and \mathbf{y} differ in more than one variables.

Solution: i) if $\mathbf{x} = \mathbf{y}$

$$p_{xx} = \sum_{i \in V} q(i) \pi(x_i | (x_v)_{v \in (V-i)})$$

ii) \mathbf{x} and \mathbf{y} differ in one variables

$$\begin{aligned} \pi(x)p_{xy} &= \pi(x)q(i)\pi(y_i | (x_v)_{v \in V_1}) \\ &= \pi(x_i, (x_v)_{v \in V_1})q(i) \frac{\pi(y_i, (x_v)_{v \in V_1})}{\pi((x_v)_{v \in V_1})} \\ &= \pi(y_i, (x_v)_{v \in V_1}) \\ &= \pi(y_i, (x_v)_{v \in V_1})q(i) \frac{\pi(x_i, (x_v)_{v \in V_1})}{\pi((x_v)_{v \in V_1})} \\ &= \pi(y)q(i)\pi(x_i | (x_v)_{v \in V_1}) \\ &= \pi(y)p_{yx} \end{aligned}$$

iii) \mathbf{x} and \mathbf{y} differ in more than one variables then

$$p_{xx} = p_{yx} = 0$$

9. Consider binary random variables $\mathbf{V} = \{V_1, V_2, \dots, V_m\}$, $\mathbf{H} = \{H_1, H_2, \dots, H_n\}$ taking values $(\mathbf{v}, \mathbf{h}) \in \{0, 1\}^{m+n}$. The joint probability distribution is given by

$$p(\mathbf{v}, \mathbf{h}) = \frac{1}{Z} e^{-E(\mathbf{v}, \mathbf{h})}$$

where E is an energy function defined as

$$E(\mathbf{v}, \mathbf{h}) = - \sum_{i=1}^n \sum_{j=1}^m w_{ij} h_i v_j - \sum_{j=1}^m b_j v_j - \sum_{i=1}^n c_i h_i$$

and Z is the normalizing constant.

$$Z = \sum_{\mathbf{v}, \mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})}$$

(a) Show that

$$P(V_l = 1|\mathbf{h}) = \sigma\left(\sum_{i=1}^n w_{il}h_i + b_l\right)$$

Solution:

$$P(V_1|h) = \sigma\left(\sum_{i=1}^n W_{i1}h_{hi} + b_i\right)$$

Representing all variables as X, except V_1

$$\begin{aligned} P(V_1|Xh) &= \frac{P(V_1|Xh)}{P(X, h)} \\ &= \frac{e^{-E(V=1, X, h)}}{e^{-E(V_1=1, X, h)} + e^{E(V_1=0, X, h)}} \\ &= \frac{e^{-\epsilon(X, h) - 1\eta(h)}}{e^{-\epsilon(X, h) - 1\eta(h)} + e^{-\epsilon(X, h) - 0\eta(h)}} \\ &= \frac{e^{-\epsilon(X, h)} \cdot e^{-\eta(h)}}{e^{-\epsilon(X, h)}e^{-\eta(h)} + e^{-\epsilon(X, h)}} \\ &= \frac{e^{-\eta(h)}}{e^{-\eta(h)} + 1} \\ &= \frac{1}{1 + e^{\eta h}} \\ &= \sigma\left(-e^{\eta(h)}\right) \\ &= \sigma\left(\sum_{i=1}^n w_{il}h_i + b_i\right) \end{aligned}$$

(b) Show that

$$P(h_k = 1|\mathbf{v}) = \sigma\left(\sum_{j=1}^n w_{kj}v_j + c_k\right)$$

Solution:

$$\begin{aligned} P(h_k = 1|V) &= \sigma\left(\sum_{j=1}^n w_{kj}v_j + C_k\right) \\ X_{h_k} &= -\sum_{j=1}^m w_{kj}h_k - C_k \end{aligned}$$

With, $Y_{i \neq k}$

$$P(h_k = 1|h', V) = \frac{P(h_k = 1|h', V)}{P(h', V)}$$

From above part of solution, we have:

$$\begin{aligned} \frac{e^{-E(h_k=1, h', V)}}{e^{-E(h_k=1, h', V)} + e^{-E(h_k=0, h', V)}} &= \frac{e^{-y(h', V) - x(V)}}{e^{-y(h', V) - 1x(V)} + e^{-y(h', V) - 0 \cdot x(V)}} \\ &= \frac{e^{-x(V)}}{e^{-x(V)} + 1} \\ &= \frac{1}{1 + e^{x(V)}} \\ &= \sigma \left(\sum_{j=1}^m w_{kj} V_j + C_k \right) \end{aligned}$$

(c) Show that the marginal is given as below

$$p(\mathbf{v}) = \frac{1}{Z} \prod_{j=1}^m e^{b_j v_j} \prod_{i=1}^n (1 + e^{c_i + \sum_{j=1}^m w_{ij} v_j})$$

Solution:

$$P(V) = \frac{1}{Z} \prod_{j=1}^m e^{b_j V_j} \prod_{i=1}^n \left(1 + e^{c_j + \sum_{j=1}^m w_{ij} V_j} \right)$$

also,

$$\begin{aligned}
P(V) &= \frac{1}{Z} \sum_n p(V, h) \\
&= \frac{1}{Z} \sum_n e^{-E(V, h)} \\
&= \frac{1}{Z} \sum_{n_1} \sum_{n_2} \dots \sum_{n_n} e^{\sum_{j=1}^m b_j V_j} \prod_{i=1}^n e^{h_i (C_i + \sum_{j=1}^m w_{ij} V_j)} \\
&= \frac{1}{Z} e^{\sum_{j=1}^m b_j V_j} \sum_{h_1} e^{h_1 (C_1 + \sum_{j=1}^m w_{1j} V_j)} \dots \sum_{h_n} e^{h_n (C_n + \sum_{j=1}^m w_{nj} V_j)} \\
&= \frac{1}{Z} e^{\sum_{j=1}^m b_j V_j} \prod_{i=1}^n \sum_{h_i} e^{h_i (C_i + \sum_{j=1}^m w_{ij} V_j)} \\
&= \frac{1}{Z} \sum_{j=1}^m e^{b_j V_j} \prod_{i=1}^n (1 + e^{C_i + \sum_{j=1}^m w_{ij} V_j})
\end{aligned}$$

10. (a) Show that for random variable $x \in \mathbf{R}^n$ drawn from the distribution $\mathcal{N}(\mu, \Sigma)$, the random variable $y = Ax + b$ follows the distribution $\mathcal{N}(A\mu + b, A\Sigma A^T)$.

Solution: We find the mean of y by using the fact that E is a linear operator.

$$\bar{y} = \mathbb{E}\{y\} = \mathbb{E}\{Ax + b\} = A\mathbb{E}\{x\} + b = A\bar{x} + b$$

we find the covariance as such follows

$$\begin{aligned}
\mathbf{C}_y &\triangleq \mathbb{E}\{(y - \bar{y})(y - \bar{y})^\top\} \\
&= \mathbb{E}\left\{\left[(Ax + b) - (A\bar{x} + b)\right]\left[(Ax + b) - (A\bar{x} + b)\right]^\top\right\} \\
&= \mathbb{E}\left\{\left[A(x - \bar{x})\right]\left[A(x - \bar{x})\right]^\top\right\} \\
&= \mathbb{E}\left\{A(x - \bar{x})(x - \bar{x})^\top A^\top\right\} \\
&= A\mathbb{E}\left\{(x - \bar{x})(x - \bar{x})^\top\right\}A^\top \\
&= A\Sigma A^\top
\end{aligned}$$

Hence $y = Ax + b$ follows the distribution $\mathcal{N}(A\mu + b, A\Sigma A^T)$ (**Proved**)

- (b) Using the result from part (a), give a method for sampling from an arbitrary normal distribution with mean μ and variance Σ , given that you have the means to sample from the standard normal distribution.

Solution: $Y = \sum x + \mu$ where $Y \sim \mathcal{N}(\mu, \sigma^2)$. and $X \sim \mathcal{N}(0, 1)$.

11. We use the notion of “distance” to measure the difference between 2 quantities. One standard measure of distance is the Euclidean norm for points in d -dimensional space. There are problems where we would like to measure distances between different mathematical objects such as probability distributions. Consider the situation where the model outputs a probability distribution (for example: softmax) and we want the loss function to capture the distance between the predicted output and the true output. One of the commonly used functions to measure distance between 2 probability distributions is the KL-divergence. The KL-divergence between 2 distributions $p(x)$ and $q(x)$ is defined as

$$KL(p||q) = \sum_x p(x) \log \frac{p(x)}{q(x)}$$

The summation is replaced by integration for continuous distributions.

Based on the above definition, answer the following questions

- (a) Under what condition(s) is $KL(p||q) = 0$?

Solution: If a function $f(x)$ is convex, then we can say the following:

$$E[f(x)] \geq f(E[x])$$

A function is convex if $\forall \lambda \in [0, 1]$ and

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y)$$

A convex function between two points is always lower than the straight line between those points. Now, if we rearrange the KL divergence formula,

$$\begin{aligned} KL(p||q) &= \int p(x) \log \frac{p(x)}{q(x)} dx \\ &= E_p \left[\log \frac{p(x)}{q(x)} \right] \\ &= -E_p \left[\log \frac{q(x)}{p(x)} \right] \\ &\geq -\log \left(E_p \left[\frac{q(x)}{p(x)} \right] \right) \\ &= -\log \left(\int p(x) \frac{q(x)}{p(x)} dx \right) \\ &= -\log \left(\int q(x) dx \right) \\ &= -\log(1) \\ &= 0 \end{aligned}$$

If $p(x) = q(x)$ that means $p(x)$ and $q(x)$ have similar behaviour on that condition $KL(p||q) = 0$

(b) Is the function symmetric?

Solution: Although the KL divergence measures the distance between two distributions, it is not a distance measure. This is because that the KL divergence is not a metric measure. **It is not symmetric**

$$KL(p||q) - KL(q||p) = \sum_{i=1}^n \ln \left(\frac{p(i)}{q(i)} \right) (p(i) + q(i))$$

there is no reason that R.H.S could become 0.

as i is not a random variable, it just a dummy index. However, there can be a random variable that takes value i with probability $p(i)$, and another that takes the value i with probability $q(i)$.

(c) One necessary property of a distance function is that it needs to be ≥ 0 . Prove that $KL(p||q) \geq 0$.

Solution: for convex f and random variable Y , $\mathbb{E}f(Y) \geq f(\mathbb{E}Y)$

$f(x) = -\log x$ is convex

We can let $Y = \frac{p(X)}{q(X)}$. As X is a random variable, Y is also just a random variable.

Applying Linearity of Expectation we can say $\mathbb{E}[-\log Y] \geq -\log \mathbb{E}Y$

Taking expectations with respect to q , this becomes

$$\begin{aligned} \mathbb{E} \left[-\log \frac{p(X)}{q(X)} \right] &= \mathbb{E} \left[\log \frac{q(X)}{p(X)} \right] \\ \implies KL(p||q) &\geq -\log \mathbb{E}_q \left(\frac{p(X)}{q(X)} \right) = -\log \mathbb{E}_p(1) = 0 \quad \textbf{(Proved)} \end{aligned}$$

12. In future classes we will encounter situations where we need to compute the KL-divergence between 2 gaussians. As a warm up, derive the expression for the KL-divergence between 2 univariate gaussians p and q .

$p \sim \mathcal{N}(\mu_1, \sigma_1)$ and $q \sim \mathcal{N}(\mu_2, \sigma_2)$.

Understand the effect of each of the terms to the value of the divergence and try to convince yourself of these findings.

Solution:

$$\begin{aligned}
KL(p, q) &= - \int p(x) \log q(x) dx + \int p(x) \log p(x) dx \\
&= - \int p(x) \log \frac{1}{(2\pi\sigma_2^2)^{(1/2)}} e^{-\frac{(x-\mu_2)^2}{2\sigma_2^2}} dx \\
&\quad + \int p(x) \log \frac{1}{(2\pi\sigma_1^2)^{(1/2)}} e^{-\frac{(x-\mu_1)^2}{2\sigma_1^2}} dx \\
&= \frac{1}{2} \log(2\pi\sigma_2^2) - \int p(x) \log e^{-\frac{(x-\mu_2)^2}{2\sigma_2^2}} dx \\
&\quad - \frac{1}{2} \log(2\pi\sigma_1^2) + \int p(x) \log e^{-\frac{(x-\mu_1)^2}{2\sigma_1^2}} dx \\
&= \frac{1}{2} \log(2\pi\sigma_2^2) - \int p(x) \left(-\frac{(x-\mu_2)^2}{2\sigma_2^2}\right) dx \\
&\quad - \frac{1}{2} \log(2\pi\sigma_1^2) + \int p(x) \left(-\frac{(x-\mu_1)^2}{2\sigma_1^2}\right) dx \\
&= \frac{1}{2} \log(2\pi\sigma_2^2) + \frac{\int p(x)x^2 dx - \int p(x)2x\mu_2 dx + \int p(x)\mu_2^2 dx}{2\sigma_2^2} \\
&\quad - \frac{1}{2} \log(2\pi\sigma_1^2) - \frac{\int p(x)x^2 dx - \int p(x)2x\mu_1 dx + \int p(x)\mu_1^2 dx}{2\sigma_1^2}
\end{aligned}$$

Letting $\langle \rangle$ denote the expectation operator under p , I can rewrite this as

$$\begin{aligned}
&= \frac{1}{2} \log(2\pi\sigma_2^2) + \frac{\langle x^2 \rangle - 2\langle x \rangle \mu_2 + \mu_2^2}{2\sigma_2^2} \\
&\quad - \frac{1}{2} \log(2\pi\sigma_1^2) - \frac{\langle x^2 \rangle - 2\langle x \rangle \mu_1 + \mu_1^2}{2\sigma_1^2}
\end{aligned}$$

can replace x with $x + \mu_1$. The expected value of x^2 is σ_1^2 by simplifying we get

$$\begin{aligned}
&= \frac{1}{2} \log(2\pi\sigma_2^2) + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2} (1 + \log 2\pi\sigma_1^2) \\
&= \log \frac{\sigma_2}{\sigma_1} + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}
\end{aligned}$$

Bonus question Derive the expression for the general multivariate case.

Solution: Let, p and q be the pdfs of normal distributions with mean and variances as μ_1, μ_2 and Σ_1, Σ_2 , respectively. Such distributions are two multivariate gaussians. KL divergence for them is given as:

$$\begin{aligned}
 KL &= \int \left[\frac{1}{2} \log \frac{|\Sigma_2|}{|\Sigma_1|} - \frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) + \frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) \right] \times p(x) dx \\
 &= \frac{1}{2} \log \frac{|\Sigma_2|}{|\Sigma_1|} - \frac{1}{2} \text{tr} \{ E[(x - \mu_1)(x - \mu_1)^T] \Sigma_1^{-1} \} + \frac{1}{2} E[(x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2)] \\
 &= \frac{1}{2} \log \frac{|\Sigma_2|}{|\Sigma_1|} - \frac{1}{2} \text{tr} \{ I_d \} + \frac{1}{2} (\mu_1 - \mu_2)^T \Sigma_2^{-1} (\mu_1 - \mu_2) + \frac{1}{2} \text{tr} \{ \Sigma_2^{-1} \Sigma_1 \} \\
 &= \frac{1}{2} \left[\log \frac{|\Sigma_2|}{|\Sigma_1|} - d + \text{tr} \{ \Sigma_2^{-1} \Sigma_1 \} + (\mu_2 - \mu_1)^T \Sigma_2^{-1} (\mu_2 - \mu_1) \right].
 \end{aligned}$$