**Question 1.** [20 MARKS]

**Part 1: (10)**

In this question, we ask you to give an extension of the law of total probability. The law of total probability applied to an unconditional probability $P(B)$ is given by:

$$P(B) = \sum_k P(B|A_k)P(A_k).$$

Here we are writing the probability of event $B$ in an expanded form with conditional probabilities where the conditional events $A_k$ is a partition of the sample space $\Omega$, which amounts to the following three conditions:

$$A_i \neq \{\} \, \forall i; \text{ each event is non-empty,}$$
$$A_i \cap A_j = \{\}, i \neq j, \forall i, j; \text{ mutually exclusive,}$$
$$\cup_i A_i = \Omega; \text{ collectively exhaustive.}$$

Now, instead of $P(B)$, if we want apply the law of total probability to conditional probability $P(B|C)$, what will be the formula? Write all the correct options and the corresponding formulas in your answer.

(a) $P(B|C) = \sum_k P(B|A_k)P(A_k)$

(b) $P(B|C) = \sum_k P(B|A_k \cap C)P(A_k)$ *⟸ only true when A and C are independent.*

(c) ✓ $P(B|C) = \sum_k P(B|A_k \cap C)P(A_k|C).$ *→ both have to given C.*

*A与C协3独立.*

The conditional events $A_k$ are still a partition of the sample space here.

(c). ∵ $P(A \cap B | C) = P(A|B \cap C) P(B|C)$  ∴ c对.
↓
这儿需相好了 ∩

∵ $P(B) = \sum_j P(B|A) P(A)$.  ∴ a 对对.

**Part 2: (10)**

If two non-empty events $D$ and $F$ are *independent*, then their probabilities do not change if the other event is given as a condition: $P(D) = P(D|F)$ and $P(F|D) = P(F)$.

In part 1, if the events $A_k$ in the partition are all independent of event $C$, which of the above options (a), (b) and (c) are correct? Write all the correct options and the corresponding formulas in your answer.

A 与 c independent. ∴ $P(A_k) = P(A_k | c)$

∴ (b) 与 (c) 在这情况下相同.

∴ 都对 (b). (c).

**Question 2.** [20 MARKS]

In this question, we ask you to derive a formula related to the Bellman equation for action value $q_\pi$. Recall that $q_\pi(x, c)$ is the action value of state action pair $x$ and $c$ under policy $\pi$ defined as the expected return:

$$q_\pi(x, c) \doteq E_\pi\left[G_t | S_t = x, A_t = c\right],$$

and $v_\pi(x)$ is the state value of state $x$ under policy $\pi$ defined as the expected return:

$$v_\pi(x) \doteq E_\pi\left[G_t | S_t = x\right].$$

If $g(x, c)$ is the expected reward:

$$g(x, c) \doteq E\left[R_{t+1} | S_t = x, A_t = c\right],$$

then derive the following identity:

*[handwritten: 因是離散的 一個季所以用LOTUS, 所以 $E[g_i|x)] = \sum_{x \in X} g(x) \cdot P(X=x)$]*

$$q_\pi(x, c) = g(x, c) + \gamma \boxed{\sum_{x'} p(x'|x, c) v_\pi(x'),} \quad \forall x, \forall c,$$

*[handwritten under box: $P(X=x)$ 的？是 $\sum_{x}$. $g(x)$.]*

where $p(x'|x, c) = P(S_{t+1} = x'|S_t = x, A_t = c)$ is the probability of next state $x'$ given the current state-action pair $x, c$.

Use the linearity of expectation (LE), the law of total expectation (LOTE), the law of the unconscious statistician (LOTUS) and the Markov property (MP) in your derivation. For each step where you use one of these rules, write the name of the rule beside that step as (LE), (LOTE), (LOTUS), and (MP). *[handwritten: ⇒ 為 $q_\pi(x, c)$ 化簡]*

*[handwritten derivation on left side:]*

$q_\pi(x, c) \doteq E_\pi[G_t | S_t = x, A_t = c]$

$= E_\pi[R_{t+1} + \gamma G_{t+1} | S_t = x, A_t = c]$  （增加$\gamma$）

$= E_\pi[R_{t+1} | S_t = x, A_t = c] + E_\pi[\gamma G_{t+1} | S_t = x, A_t = c]$  by (LE: Linearity $E[X+Y] = E[X] + E[Y]$ ✓)

$= E_\pi[R_{t+1} | S_t = x, A_t = c] + \gamma E_\pi[G_{t+1} | S_t = x, A_t = c]$  by (LE: Linearity $E[aX] = aE[X]$ ✓)
*[因為$R_{t+1}$的分布 不涉及到 policy$\pi$, 所以把 $R_{t+1}$... 把 expectation中的 constant 提出来. 把 expectation中的constant 提出来. 因为$\gamma$是常数.]*

$= g(x, c) + \gamma E_\pi[G_{t+1} | S_t = x, A_t = c]$  （即 expected reward.）

$= g(x, c) + \gamma E_\pi[E_\pi[G_{t+1} | (S_{t+1}), S_t = x, A_t = c] | S_t = x, A_t = c]$  （加入新状态。相当$Y$.）

$= g(x, c) + \gamma E_\pi[E_\pi[G_{t+1} | (S_{t+1})] | S_t = x, A_t = c]$  by (MP.)
*[根据新状态. 根据新状态 MP, we don't need the previous state and action... 因为 只有 $S_{t+1}$ ]*

$= g(x, c) + \gamma E_\pi[v_\pi(S_{t+1}) | S_t = x, A_t = c]$
*[涉及新状态 $S_{t+1}$ 的分布。]*

$= g(x, c) + \gamma \sum_{x'} P(x'|x, c) v_\pi(x')$  by (LOTUS.)

*[handwritten: probability, 涉及$\sum$的, input 必须是事件!! 新的$S_{t+1}, S_t, A_t$的分布. $g$的分布. 所以, $g$里:$E[v_\pi(S_{t+1}) | S_t = x, A_t = c]$ 应为( $x'$ | $x$, $c$ )]*

*[middle/right handwritten annotations:]*

Linearity: $E[X + Y] = E[X] + E[Y]$ ✓

Linearity: $E[aX] = aE[X]$ ✓

This is known as the **law of total expectation** and also the tower rule.

Again, the conditional **expectation** $E[X|Y]$ integrates only over the random variable $X$ but not $Y$. On the other hand, in $E[E[X|Y]]$, the outer **expectation** integrates over the random variable $Y$ but not $X$ which is already integrated out.
*[条件期望×随机变量, 为另一个$Y$的... 是随机变量× 随机变量, 为另一个$Y$的... 在$E[X|Y]$中是随机变量了.]*

**Interesting results**
*[注: $E[X|Z] = E[E[X|Y,Z]|Z]$ ⇒ LOTE for **conditional expectation**. 加入的新变量.]*

$E[E[f(X)|Y] g(Y)] = E[E[f(X)g(Y)|Y]].$

$E[f(X, Y)g(Y)] = E[E[f(X, Y)|Y] g(Y)].$

*[expectation of an conditional expectation.]*

Show that $E[X] = E[E[X | Y]]$ ⇒ law of total expectation.
*[R.V.]*

*[comment box: D — Daniel Asimakwini 对所有人说 下午 1:45 — for the MP step, we can drop the variables in the inner expectation because we don't need the previous state and action to determine the current one? Yes!!]*

Law of the unconscious statistician: $E[g(X)] = \sum_{x \in \mathcal{X}} g(x) P(X=x)$
*[X is the value that is a function of another R.V. 所以 $E[f(Y)], Y = g(X)$. 所以 $= \sum_x P(y=g(x))$ ... $= \sum g(x) P(y=g(x) = g(x))$ ... $= \sum_{x \in X} g(x) P(X=x)$]*

**Question 3.**  [20 MARKS]

(10+10)

In this question, we ask you to give the Bellman optimality equation and value iteration update rule for action values. The Bellman optimality equation for $v_*$ is given by:
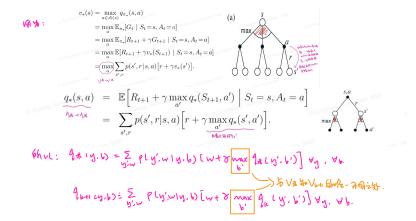
$$v_*(y) = \max_b \sum_{y',w} p(y',w|y,b) \left[ w + \gamma v_*(y') \right], \forall y, \tag{1}$$

where $p(y',w|y,b)$ is the joint probability of next state $y'$ and reward $w$ given the current state-action pair $y,b$. Then the value iteration method can be directly obtained based on the optimality equation by replacing the optimal state value $v_*$ with estimate $v_k$ for the $k$th iteration and replacing equality $=$ with assignment $\doteq$:

$$v_{k+1}(y) \doteq \max_b \sum_{y',w} p(y',w|y,b) \left[ w + \gamma v_k(y') \right], \forall y. \tag{2}$$

Provide the Bellman optimality equation for $q_*$ and the corresponding value iteration update rule for action value using the state-action pair $y, b$.

(handwritten annotations)

$$v_*(s) = \max_{a \in \mathcal{A}(s)} q_{\pi_*}(s,a)$$
$$= \max_a \mathbb{E}_{\pi_*}[G_t \mid S_t = s, A_t = a]$$
$$= \max_a \mathbb{E}_{\pi_*}[R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t = a]$$
$$= \max_a \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a]$$
$$= \max_a \sum_{s',r} p(s',r|s,a)[r + \gamma v_*(s')].$$

$$q_*(s,a) = \mathbb{E}\left[ R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') \mid S_t = s, A_t = a \right]$$
$$= \sum_{s',r} p(s',r|s,a)\left[ r + \gamma \max_{a'} q_*(s', a') \right].$$

$\text{BhvL:}\quad q_*(y,b) = \sum_{y',w} p(y',w|y,b)\left[ w + \gamma \boxed{\max_{b'} q_*(y',b')} \right] \forall y, \forall b.$

$q_{k+1}(y,b) \doteq \sum_{y',w} p(y',w|y,b)\left[ w + \gamma \boxed{\max_{b'} q_k(y',b')} \right] \forall y, \forall b.$

## Question 4.  [20 MARKS]

Prove that the discounted sum of rewards is always finite if the rewards are bounded: $|R_{t+1}| \leq R_{\max}$ for all $t$ for some finite $0 < R_{\max} < \infty$:

$$\left| \sum_{i=0}^{\infty} \gamma^i R_{t+1+i} \right| < \infty; \qquad \text{for } \gamma \in [0,1).$$

$\because |R_{t+1}| \leq R_{\max}$

$\therefore R_{t+1} \leq R_{\max}$. ← 和求3个条件倒到

$\because$ 在有 $\gamma$ 的情况下, $\sum\limits_{i=0}^{\infty} \gamma^i R_{t+1+i} = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots$

$\textcircled{\leq} R_{\max} + \gamma R_{\max} + \gamma^2 R_{\max} + \cdots$

↳ 提取公因子 $R_{\max}$.

$= R_{\max} (1 + \gamma + \gamma^2 + \cdots)$

$= \boxed{\dfrac{R_{\max}}{1-\gamma}}$ → 如果R一直等于 $R_{\max}$.

$< \infty$

$\therefore \left| \sum\limits_{i=0}^{\infty} \gamma^i R_{t+1+i} \right| < \infty$

• Continuing Task 有 $\gamma$ 下面的如果都比私有无穷无尽. 当时候 $\gamma$ 为一0 1的时 极限之无私无尽)

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \qquad 0 \leq \gamma \leq 1$$

$$G_t \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \cdots$$
$$= R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \gamma^2 R_{t+4} + \cdots)$$
$$= R_{t+1} + \gamma G_{t+1}$$

如果R一直等于1: $\quad G = \dfrac{1}{1-\gamma}$

**Question 5.** [20 MARKS]

In this question, we ask you to amend incorrect statements. In the following, there are two incorrect statements. Write mathematical equations or expressions that will make the statements correct and briefly list the changes you made on the provided equations or expressions.

**Part 1: (10)**

The Bellman equation for action value can be written as:

$$q_\pi(s,a) = \sum_{s',r,s,a'} p(s',r|s,a)\pi(a'|s')\left[r + \gamma q_\pi(s',a')\right], \forall s, \forall a.$$

*(handwritten annotations in Chinese around the equation)*

**Part 2: (10)**

The optimal state value can be related to the optimal action value in the following way:

$$v_*(s) \le \max_a q_\pi(s,a), \forall s, \forall a.$$

Part 1. 应为
$$q_\pi(s,a) \doteq \mathbb{E}_\pi[G_t | S_t = s, A_t = a]$$
$$= \sum_{s'}\sum_r p(s',r|s,a)\left[r + \gamma\mathbb{E}_\pi[G_{t+1} | S_{t+1} = s']\right]$$
$$= \sum_{s'}\sum_r p(s',r|s,a)\left[r + \gamma\sum_{a'}\pi(a'|s')\mathbb{E}_\pi[G_{t+1} | S_{t+1} = s', A_{t+1} = a']\right]$$
$$= \sum_{s'}\sum_r p(s',r|s,a)\left[r + \gamma\sum_{a'}\pi(a'|s')q_\pi(s',a')\right]$$

$$v_\pi(s) = \sum_a \pi(a|s)\sum_{s'}\sum_r p(s',r|s,a)[r + \gamma v_\pi(s')]$$
$$q_\pi(s,a) = \sum_{s'}\sum_r p(s',r|s,a)\left[r + \gamma\sum_{a'}\pi(a'|s')q_\pi(s',a')\right]$$

所以 $q_\pi(s,a) = \sum_{s',r} p(s',r|s,a)\sum_{a'}\pi(a'|s')\left[r + \gamma q_\pi(s',a')\right], \forall s, \forall a.$

Part 2. 应为:

$$v_*(s) = \max_a q_*(s,a)$$

所以 L: $v_*(s) = \max_a q_*(s,a)$