

מגשים:

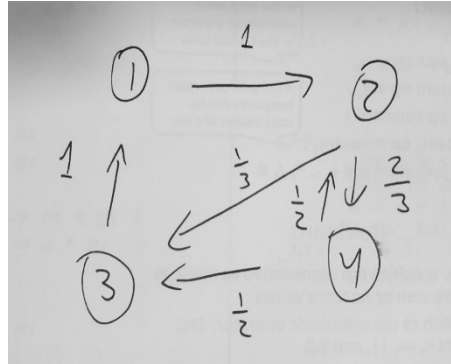
דן מלכא 304773591

אופק ברנסקי 311170138

דן ברוך 207042391

שאלה 1

1. מכונת המצבים הינה:



2. הגרף מחלקת קשירות אחת, לכן הגרף *Irreducible*, וקיימת מחלקת קשירות אחת

3. נתבונן בגרף שקיבלנו, נחשב את הזמן מחזור של כל מצב:

$$p_{1,1}^3, p_{1,1}^5 > 0 \rightarrow \gcd(1 \rightarrow 1) = d_1 = 1$$

$$p_{2,2}^2, p_{2,2}^3 > 0 \rightarrow \gcd(2 \rightarrow 2) = d_2 = 1$$

$$p_{3,3}^3, p_{3,3}^4 > 0 \rightarrow \gcd(3 \rightarrow 3) = d_3 = 1$$

$$p_{4,4}^2, p_{4,4}^5 > 0 \rightarrow \gcd(4 \rightarrow 4) = d_4 = 1$$

לפיכך ניתן לקבוע כי השרשרת היא *aperiodic*.

4. נחשב את ההתפלגות הסטוציונרית:

$$[x_1, x_2, x_3, x_4] \cdot \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1/3 & 2/3 \\ 1 & 0 & 0 & 0 \\ 0 & 1/2 & 1/2 & 0 \end{bmatrix} = [x_1, x_2, x_3, x_4]$$

↓

$$\begin{cases} x_3 = x_1 \\ x_1 + \frac{x_4}{2} = x_2 \\ \frac{x_2}{3} + \frac{x_4}{2} = x_3 \\ \frac{2x_2}{3} = x_4 \end{cases} \rightarrow \begin{cases} x_1 + x_2 + x_3 + x_4 = 1 \\ x_3 = x_1 \\ x_1 + \frac{x_4}{2} = x_2 \\ \frac{x_2}{3} + \frac{x_4}{2} = x_3 \\ \frac{2x_2}{3} = x_4 \end{cases} \rightarrow \begin{cases} x_1 = \frac{2}{3}x_2 \\ x_3 = \frac{2}{3}x_2 = x_1 \\ x_1 + x_2 + x_3 + x_4 = 1 \\ \frac{2x_2}{3} = \left(\frac{2}{3}\right)\left(\frac{3}{2}\right)x_1 = x_4 \end{cases}$$

$$\rightarrow x_1 + \frac{3}{2}x_2 + x_1 + \left(\frac{2}{3}\right)\left(\frac{3}{2}\right)x_1 = 1 \rightarrow 4.5x_1 = 1 \rightarrow x_1 = \frac{2}{9}$$

↓

$$x_2 = 1.5\left(\frac{2}{9}\right)x_1 = \frac{1}{3}, \quad x_3 = \frac{2}{9}, \quad x_4 = \left(\frac{2}{3}\right)\left(\frac{1}{3}\right) = \frac{2}{9}$$

$$\pi = \left[\frac{2}{9}, \frac{1}{3}, \frac{2}{9}, \frac{2}{9}\right]$$

5. נחשב את ה *expected return time* לכל מצב :
 הראנו כי השרשרת אי פריקה ומחזורית, לכן המצבים הינם *Positive recurrent*, לפי הנלמד בשיעור,
 לכל מצב מתקיים :

$$\pi_i = \frac{1}{\mathbb{E}[T_i]}$$

התפלגות
 סטוציונרית

לכן :

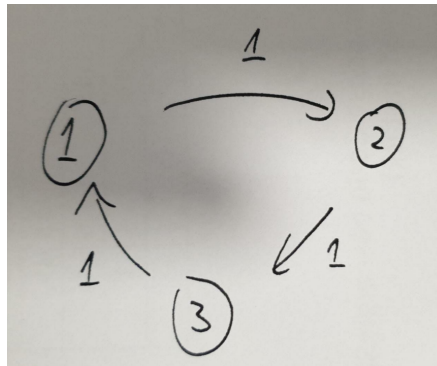
$$\mathbb{E}[T_1] = \mathbb{E}[T_3] = \mathbb{E}[T_4] = \frac{1}{\frac{2}{9}} = 4.5$$

$$\mathbb{E}[T_2] = \frac{1}{\frac{1}{3}} = 3$$

6. נגדיר מטריצה P' באופן הבא :

$$\begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}$$

נקבל :



נחשב את ההתפלגות הסטוציונרית :

$$[x_1, x_2, x_3] \cdot \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix} = [x_1, x_2, x_3] \rightarrow \pi = \left[\frac{1}{3}, \frac{1}{3}, \frac{1}{3} \right]$$

נקבל :

$$\mathbb{E}[T_1] = 3$$

וכמובן, מתקיים :

$$p_{1,1}^m = \begin{cases} 1 & \forall m \text{ s.t. } m \bmod 3 = 0 \\ 0 & \text{else} \end{cases}$$

שאלה 2

1. נגדיר MDP לבעיה:

$$\begin{aligned} S &= \{0, 1, \dots, 2k-1\}, S_0 = k \\ A &= \{CW, CCW\} \\ P(s'|s, CW) &= \begin{cases} 1, & s' = s+1 \bmod 2k-1 \\ 0, & \text{otherwise} \end{cases} \\ P(s'|s, CCW) &= \begin{cases} 1, & s' = s-1 \bmod 2k-1 \\ 0, & \text{otherwise} \end{cases} \\ r(s, a) &= \begin{cases} 1, & s = 0 \\ 0, & \text{otherwise} \end{cases} \end{aligned}$$

2. נבחין שעל מנת למקסם את הסיכוי לסיים ב-0, אנו רוצים להיות במצבים הקרובים ל-0 עד כמה שאפשר. לכן, המדיניות האופטימלית תהיה כזו ה"שואפת" להתקרב ל-0 עד כמה שאפשר, ולהישאר סביב 0 כאשר נמצאים במצבים הקרובים ל-0. נוכל להגדיר באופן הבא:

$$\pi^*(s) = \begin{cases} CW, & s > k \\ CCW, & s \leq k \end{cases}$$

עבור מדיניות זו, כל עוד אנו רחוקים ממצב 0, נתקדם לכיוון 0 במסלול הקצר ביותר, ואם אנו נמצאים במצב 0, נתקדם למצב $2k-1$, ובצעד הבא נחזור למצב 0. נבצע איטרציה אחת של value iteration:

$$\begin{aligned} V_0(s) &= 0, \forall s \\ V_1(s) &= \max_{a \in A} \left\{ r(s, a) + \gamma \sum_{\{s' \in S\}} p(s'|s, a) V_0(s') \right\} = \max_{a \in A} \{r(s, a)\} \\ \text{מכיוון שלכל } s \neq 0 \text{ מתקיים } r(s, a) = 0 \text{ לכל } a, \text{ ועבור } s = 0 \text{ מתקיים } r(s, a) = 1, \text{ נקבל:} \\ V_1(s) &= \begin{cases} 1, & s = 0 \\ 0, & \text{otherwise} \end{cases} \end{aligned}$$

4. נבצע איטרציה נוספת:

$$\begin{aligned} V_2(s) &= \max_{a \in A} \left\{ r(s, a) + \gamma \sum_{\{s' \in S\}} p(s'|s, a) V_1(s') \right\} \\ \text{עבור } s \neq 2k-1, 0, 1, & \text{ מתקיים } p(s'|s, a) = 0 \text{ או } V_1(s') = 0 \text{ לכל } s' \text{ (הם לא סמוכים למצב 0, לכן} \\ & \text{ההסתברות לעבור למצב } s' \text{ עבורו } V_1(s') \text{ היא 0).} \\ & \text{נתבונן במצבים שנותרו:} \end{aligned}$$

$$\begin{aligned} V_2(2k-1) &= \max_{a \in A} \left\{ r(2k-1, a) + \gamma \sum_{\{s' \in S\}} p(s'|2k-1, a) V_1(s') \right\} \\ &= \max_{a \in A} \left\{ \gamma \sum_{\{s' \in S\}} p(s'|2k-1, a) V_1(s') \right\} = \gamma p(s'|2k-1, CW) = \gamma \end{aligned}$$

$$\begin{aligned} V_2(1) &= \max_{a \in A} \left\{ r(1, a) + \gamma \sum_{\{s' \in S\}} p(s'|1, a) V_1(s') \right\} = \max_{a \in A} \left\{ \gamma \sum_{\{s' \in S\}} p(s'|1, a) V_1(s') \right\} \\ &= \gamma p(s'|2k-1, CCW) = \gamma \\ V_2(0) &= \max_{a \in A} \left\{ r(0, a) + \gamma \sum_{\{s' \in S\}} p(s'|0, a) V_1(s') \right\} = r(0, CW) = 1 \end{aligned}$$

סה"כ:

$$V_2(s) = \begin{cases} 1, & s = 0 \\ \gamma, & s \in \{2k-1, 1\} \\ 0, & \text{otherwise} \end{cases}$$

5. נחשב את $V^*(s)$ לכל $s \in \{0, 1, 2, 3\}$

חישבנו את V_0, V_1, V_2 לכל s , נמשיך בחישוב (אבחנה, משיקולי סימטריה, $V_i(1) = V_i(3)$, לכן אחשב רק את $V_i(1)$):

$$V_3(0) = \max_{a \in A} \left\{ r(0, a) + \gamma \sum_{\{s' \in S\}} p(s'|0, a) V_2(s') \right\} = 1 + \gamma \cdot 1 \cdot \gamma = 1 + \gamma^2$$

$$V_3(1) = \max_{a \in A} \left\{ r(1, a) + \gamma \sum_{\{s' \in S\}} p(s'|0, a) V_2(s') \right\} = 0 + \gamma \cdot 1 \cdot 1 = \gamma$$

$$V_3(2) = \max_{a \in A} \left\{ r(2, a) + \gamma \sum_{\{s' \in S\}} p(s'|0, a) V_2(s') \right\} = 0 + \gamma \cdot 1 \cdot \gamma = \gamma^2$$

$$V_4(0) = \max_{a \in A} \left\{ 1 + \gamma \sum_{\{s' \in S\}} p(s'|0, a) V_3(s') \right\} = 1 + \gamma \cdot 1 \cdot \gamma = 1 + \gamma^2$$

$$V_4(1) = \max_{a \in A} \left\{ 0 + \gamma \sum_{\{s' \in S\}} p(s'|0, a) V_3(s') \right\} = \gamma \cdot 1 \cdot (1 + \gamma^2) = \gamma(1 + \gamma^2)$$

$$V_4(2) = \max_{a \in A} \left\{ 0 + \gamma \sum_{\{s' \in S\}} p(s'|0, a) V_3(s') \right\} = \gamma \cdot 1 \cdot \gamma = \gamma^2$$

$$V_5(0) = \max_{a \in A} \left\{ 1 + \gamma \sum_{\{s' \in S\}} p(s'|0, a) V_4(s') \right\} = 1 + \gamma \cdot 1 \cdot (\gamma(1 + \gamma^2)) = 1 + \gamma^2(1 + \gamma^2)$$

$$V_5(1) = \max_{a \in A} \left\{ 0 + \gamma \sum_{\{s' \in S\}} p(s'|0, a) V_4(s') \right\} = \gamma \cdot 1 \cdot (1 + \gamma^2) = \gamma(1 + \gamma^2)$$

$$V_5(2) = \max_{a \in A} \left\{ 0 + \gamma \sum_{\{s' \in S\}} p(s'|0, a) V_4(s') \right\} = \gamma \cdot 1 \cdot \gamma(1 + \gamma^2) = \gamma^2(1 + \gamma^2)$$

ניתן לראות דפוס, ולחשב במקרה הכללי:

$$V_i(0) = 1 + \gamma \cdot V_{i-1}(1)$$

$$V_i(1) = \gamma \cdot V_{i-1}(0)$$

כאשר $V_{i+1} = V_i$ עבור i זוגי.

אז:

$$V_i(0) = 1 + \gamma^2 \cdot V_{i-2}(0) = 1 + \gamma^2 + \gamma^4 \cdot V_{i-4}(0) \dots$$

ולכן בגבול לאינסוף:

$$V^*(0) = \dots = 1 + \gamma^2 + \gamma^4 + \dots = \sum_{i=0}^{\infty} (\gamma^2)^i = \sum_{i=0}^{\infty} (\gamma^2)^i$$

זהו טור הנדסי ו- $\gamma^2 < 1$ לכן $V^*(0) = \frac{1}{1-\gamma^2}$

באופן דומה, נקבל:

$$V_i(1) = \gamma \cdot (1 + \gamma \cdot V_{i-2}(1)) = \gamma + \gamma^2 \cdot V_{i-2}(1) = \gamma + \gamma^3 + \gamma^5 + \dots$$

$$V^*(1) = \sum_{i=0}^{\infty} \gamma^{2i+1} = \gamma \sum_{i=0}^{\infty} \gamma^{2i} = \frac{\gamma}{1-\gamma^2}$$

$$V^*(3) = \frac{\gamma}{1-\gamma^2}$$

עבור $V_i(2)$ הנוסחה הכללית היא:

$$V_i(2) = \gamma \cdot V_i(1) = \gamma \cdot (\gamma + \gamma^3 + \gamma^5 + \dots) = \gamma^2 + \gamma^4 + \gamma^6 \dots$$

$$V^*(2) = \sum_{i=0}^{\infty} \gamma^{2i+1} = \gamma^2 \sum_{i=0}^{\infty} \gamma^{2i} = \frac{\gamma^2}{1-\gamma^2}$$

שאלה 3

1. נגדיר MDP:

$$\mathcal{S} = \{(x, d_i) \mid x \subseteq \{0, 1, 2, \dots, N\}, d_i \in \{0, 1, \dots, 9\}\}$$

כאשר:

- x הינו קבוצת slots הפנויים, שלא שובצו מספרים בהם (יכול לקבל כל תת קבוצה של $\{0, 1, \dots, N\}$)
- d_i הינו המספר שהוגרל, ונדרש לשיבוץ באחד ה slots

בנוסף:

$$\mathcal{A} = \{a_{x,j} \mid x \subseteq \{0, 1, \dots, N\}, j \in x\}$$

ז"א, בהינתן מצב $s' := (x', d'_i)$, קבוצת הפעולות האפשריות (ז"א הצבת מספר d_i שהוגרל באחד מהתאים הריקים מתוך s), הינו האוסף:

$$\mathcal{A}_{s'} = \{a_{x',j} \mid x' = x, j \in x'\}$$

ונגדיר את פונקציית הרווח עבור (x, d_i) ו $a_{x,j}$

$$\mathcal{R}(s', a_{x,j}) = d_i \cdot 10^j$$

ז"א, עבור הצבת המספר d_i בתא j , נקבל "תוספת" רווח למספר הסופי מהצורה הנ"ל.

נגדיר את מטריצת המעברים-

עבור פעולה $a_{x,j}$ המציבה את המספר d_i בתא j , יש לנו 10 מצבים שונים שניתן לעבור אליהם (בהתאם להגרלת המספר הבא) כל אחד בסיכוי $1/10$:

אז, עבור מעבר ממצב $s_t := (x, d_i)$ ביצענו את הפעולה $a_{x,j}$, נעבור למצב $s_{t+1} := (x_{t+1}, d'_i)$ באופן הבא:

$$p_{s_t, s_{t+1}} = \begin{cases} 1/10 & x' := x/\{j\} \\ 0 & \text{else} \end{cases}$$

לדוגמה, עבור:

$$s_t = (x_t, d_i \mid d_i \in \{0, \dots, 9\}), \quad a_t = a_{x_t, j}$$

↓

$$s_{t+1} := (x_{t+1}, d'_i) := ((x_t/\{j\}), d'_i)$$

ומצבי ההתחלה יהיו:

$$s_0 = \{\{0, 1, \dots, N\}, d_i \mid d_i \in \{0, 1, \dots, 9\}\}$$

בנוסף הבעיה היא *finite horizon*, לכן $\gamma = 1$

2. למדנו בשיעור כי:

the tail of an optimal policy is optimal for the tail problem

עבור MDP עם $finite horizon$, לפיכך המדיניות האופטימלית עבור מצב s בזמן t איננו מושפע מהגרלת המספרים קודמת, ואיפה הם שובצו אלא רק מהמצב הנתון – קבוצת ה $slots$ הריקים והמספר הרנדומי שדורש שיבוץ.

נתבונן על החלטה של מדיניות האופטימלית π^* עבור מצב s_t , אז לפי הטענה שצינו לעיל, המדיניות האופטימלית עבור מצב s_t תהיה זהה למדיניות כאילו s_t היה המצב ההתחלתי, עבור בעיה קטנה יותר. נראה כי קיימת מדיניות אופטימלית יחידה-

עבור 2 מצבים שונים s_1, s_2 (כאשר הוצבו מספרים שהוגרלו בעבר, בתאים שונים, אך בשני המצבים נשארו m תאים ריקים), אנו נדרשים להציב מספר שהוגרל d_i באחד מהתאים הריקים מתוך m תאים. הרי כי מצבים אלו שקולים, עד כדי הפעלת פונקציה מונוטונית על הרווח שהצטבר עד כה- נניח כי הרווח המקסימלי עבור s_2 יתבצע ע"י הצבת המספר d_i שהוגרל בתא j , אז המדיניות עבור s_1 תניב רווח מקסימלי ע"י השמת המספר d_i בתא j' , כמובן ברווח יהיה כפולה של המספר 10.

נראה את המדיניות בעזרת אינדוקציה - על כמות המקומות הריקים – $slots$ שנמצאים ב"משחק":

מקרה בסיס-

נותר $slot$ אחד פנוי, אז ברור כי המדיניות האופטימלית תבצע השמה של המספר i שהוגרל בתא הריק.

הנחה-

עבור m תאים פנויים, נניח כי בוצע מדיניות אופטימלית של השמת מספרים בתאים $N - M$, נסמן אותה כ a' , והרווח שלה כ $r'(s', a')$ (כאשר s' מתאר את המצב בו אנו נמצאים, ואיזה מספר הוגרל).

צעד-

עבור $m+1$ תאים פנויים, נראה כי המדיניות האופטימלית איננה תלויה בעבר, אלא רק בתאים הריקים והמספר שהוגרל:

$$\pi^* = \underset{a_{\{...,j\}}}{\operatorname{argmax}} \left[\left\{ j \cdot d_i + \frac{1}{10} \cdot \left(\sum_0^9 r'(s', a') \right) \right\} \mid \forall j \in \text{empty slots} \right]$$

לפיכך ניתן לראות כי עבור פונקציית הרווח, היא איננה תלויה בשיבוצים הקודמים, כל עוד בוצע השמה אופטימלית, אלא רק בכמה תאים פנויים יש, ומה המספר שהוגרל.

3. נשתמש ב $backwards recursion$ –

a. עבור הספרה האחרונה d_i , יש לנו בדיוק מקום אחד לשים אותה, לפיכך:

$$\forall d_i \in \{0,1,..9\}, \quad j \in \{0,1,2\}$$

↓

$$\pi^*({j}, d_i) = a_{j,j}, \quad r = d_i \cdot 10^j$$

b. עבור הספרה לפני אחרונה, ישנם 3 מצבים אפשריים

$$_ _ X, _ X _, X _ _$$

נחשב את הפעולה האופטימלית עבור כל מצב:

$$: _ _ X$$

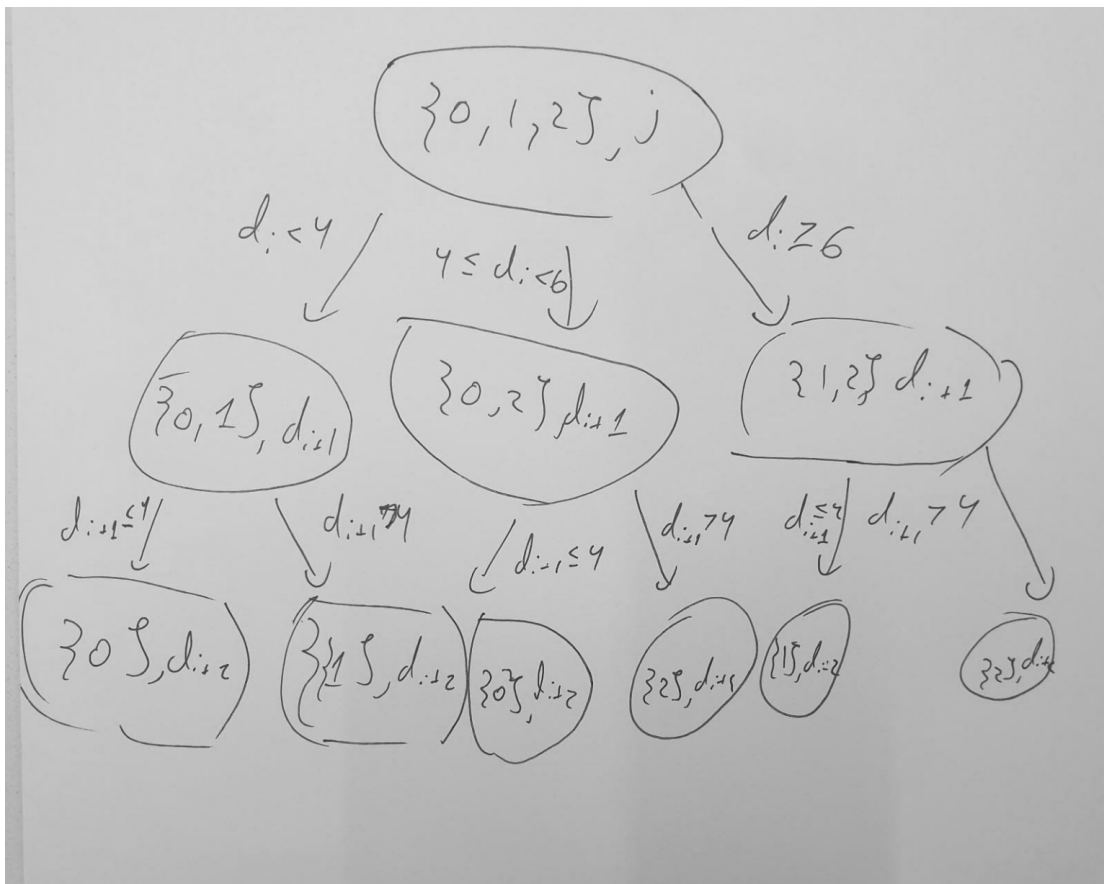
$$\begin{aligned} \pi^*({0,1}, d_i) &= \underset{a}{\operatorname{argmax}} \left\{ 100d_i + \frac{1}{10} \cdot \sum_{k=0}^9 10k, \quad 10d_i + \frac{1}{10} \cdot \sum_{k=0}^9 100k \right\} \\ &= \underset{a}{\operatorname{argmax}} \{ 100d_i + 45, 10d_i + 450 \} \end{aligned}$$

↓

$$a_{\{0,1\},0} \text{ if } d_i > 4, \quad a_{\{0,1\},1} \text{ if } d_i \leq 4$$

ז"א, אם המספר d_i שהוגרל קטן מ-5, נשים אותו בתא האמצעי (אינדקס 1), אחרת נשים אותו בתא השמאלי (אינדקס 0).

$$: _ _ X$$



$$V \in \mathbb{R}^{|S|}; \quad TV(s) = \frac{1}{|A|} \sum_{a \in A} \left(r(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) V(s') \right)$$

Show that T is γ -contracting w.r.t Max-Norm

Which is equivalent to:

$$\begin{aligned} \|T(V_1) - T(V_2)\|_\infty &\leq \gamma \cdot \|V_1 - V_2\|_\infty \\ \left\| \frac{1}{|A|} \sum_{a \in A} \left(\cancel{r(s, a)} + \gamma \sum_{s' \in S} p(s'|s, a) V_1(s') \right) - \frac{1}{|A|} \sum_{a \in A} \left(\cancel{r(s, a)} + \gamma \sum_{s' \in S} p(s'|s, a) V_2(s') \right) \right\|_\infty &= \\ \left\| \frac{1}{|A|} \gamma \sum_{a \in A} \left(\sum_{s' \in S} p(s'|s, a) (V_1(s') - V_2(s')) \right) \right\|_\infty &\stackrel{\leq}{=} \\ &\stackrel{\substack{\text{Triangle ineq.} \\ \text{Max-Norm} \\ \text{not affecting} \\ \text{left term}}}{=} \\ &\leq \frac{1}{|A|} \gamma \cdot \sum_{a \in A} \left(\underbrace{\sum_{s' \in S} p(s'|s, a)}_{\sum_p \text{prob} = 1} \right) \|V_1(s') - V_2(s')\|_\infty = \frac{1}{|A|} \gamma \cdot \sum_{a \in A} \underbrace{1}_{|A|} \cdot \|V_1(s') - V_2(s')\|_\infty = \\ &= \gamma \cdot \|V_1 - V_2\|_\infty; \quad Q.E.D \end{aligned}$$