

سوال ۱ -

قسمت الف)

ابتدا MPN را تعریف میکنیم. اگر S یک SPN روی X باشد، MPN \hat{S} از S ، شبکه ای با ویژگی های زیر است:

۱- هر گره D distribution D را با یک گره maximizing distribution \hat{D} جایگزین میکنیم به طوری که :

$$\hat{D} : H_{sc(D)} \rightarrow [0, \infty] \quad , \quad \hat{D}(Y) := \max_{y \in Y} D(y)$$

۲- هر گره sum S را با یک گره max \hat{S} جایگزین میکنیم. یال های خروجی گره max همان وزن های یال های خروجی گره sum متناظر را دارند. یک گره max اینگونه تعریف میشود:

$$\hat{S} := \max_{C \in ch(\hat{S})} w_{\hat{S}, C} C$$

۳- هر گره $product$ در SPN یک گره $product$ در MPN است.

حال فرضیه زیر، برای $augmented SPN$ ها، MPN میتواند احتمال MPE را محاسبه کند:

فرضیه: اگر S یک $augmented SPN$ روی X باشد که X از قبل RV های مدل SPN اوریجینال را تشکیل داده و همچنین RV های پنهان را که توسط $augmentation$ بوجود آمده اند، را نیز تشکیل داده است. اگر \hat{S} همان MPN متناظر آن باشد، با فرض اینکه N یک گره دلخواه در S و \hat{N} گره متناظر آن در \hat{S} باشد، برای هر $X \in H_{sc(N)}$ داریم:

$$\hat{N}(X) = \max_{x \in X} N(x)$$

یعنی MPN ، احتمال $augmented\ SPN$ روی X را بیشینه میکند. الگوریتم آن میشود:

Algorithm 5 MPE inference in augmented SPNs

```

1: For all sums  $S$ , let  $\mathbf{ch}'(S) := \{C_S^k \mid k \in \mathcal{X}[Z_S]\}$ 
2: Evaluate  $\mathcal{X}$  in corresponding MPN  $\hat{\mathcal{S}}$  (upwards pass)
3: Initialize  $Q$  as empty queue
4:  $Q \leftarrow$  root node
5: while  $Q$  not empty do
6:    $N \leftarrow Q$ 
7:   if  $N$  is a sum node then
8:      $Q \leftarrow \arg \max_{C \in \mathbf{ch}'(N)} \{w_{\hat{N}, \hat{C}} \hat{C}(\mathcal{X}[\mathbf{sc}(\hat{C})])\}$ 
9:   else if  $N$  is a product node then
10:     $\forall C \in \mathbf{ch}(N) : Q \leftarrow C$ 
11:   else if  $N$  is a distribution node then
12:     $\mathbf{x}^*[\mathbf{sc}(N)] = \arg \max_{\mathbf{x} \in \mathcal{X}[\mathbf{sc}(N)]} N(\mathbf{x})$ 
13:   end if
14: end while
15: return  $\mathbf{x}^*$ 

```

میبینیم که بیشینه کردن یک گره **product** به مسئله بیشینه کردن به صورت مستقل تمام فرزندان آن گره کاهش پیدا میکند. همچنین بیشینه کردن یک گره **sum** ، به مسئله یافتن فرزندی با بیشینه وزن یا خروجی و یافتن انتصاب بیشینه برای این فرزند، کاهش پیدا میکند. الگوریتم بالا با روش **back-tracking** حالت بیشینه $augmented\ SPN$ را میابد.

حال الگوریتم بالا را کمی تغییر میدهیم به طوری که روی یک $arbitrary\ SPN$ اعمال شود اما پاسخ MPE ها $augmented\ SPN$ متناظر را خروجی دهد. البته کمی پیش شرط وجود دارد که بخاطر خلاصه سازی، آن هارا به صورت انگلیسی می آورم:

1. $\forall S \in \mathbf{S}(\mathcal{S}) : \mathcal{X}[Z_S] = \mathbf{val}(Z_S)$, i.e. Z_S is maximized over all its possible states.
2. The states of the RV corresponding to not visited sum nodes are not assigned, i.e. steps 28–30 are not performed.
3. All $\tilde{w}_{S,C} \equiv 1$, which is the case when all twin weights are deterministic, i.e. for each sum

حال الگوریتم جدید عبارت است از:

Algorithm 6 MPE inference

```
1: For all sums  $S$ , let  $\text{ch}'(S) := \{C_S^k \mid k \in \mathcal{X}[Z_S]\}$ 

2:  $\forall S \in \mathbf{S}(\mathcal{S}) : \forall C \in \text{ch}(S) : \bar{w}_{S,C} \leftarrow 1$ 
3: for all sum nodes  $S \in \mathbf{S}(\mathcal{S})$  do
4:   for  $S^c \in \mathbf{S}^c(S)$  do
5:     for  $C \in \{C \in \text{ch}(S^c) \mid S \notin \text{desc}(C)\}$  do
6:        $\bar{w}_{S^c,C} \leftarrow \bar{w}_{S^c,C} \times \max_{k \in \mathcal{X}[Z_S]} \bar{w}_{S,C_S^k}$ 
7:     end for
8:   end for
9: end for
10: Equip corresponding MPN  $\hat{\mathcal{S}}$  with weights  $w_{\hat{S},\hat{C}} \leftarrow \bar{w}_{S,C} \times w_{S,C}$ 
11: Evaluate  $\mathcal{X}$  in MPN  $\hat{\mathcal{S}}$  (upwards pass)

12:  $\mathbf{S} \leftarrow \mathbf{S}(\mathcal{S})$ 
13: Initialize  $Q$  as empty queue
14:  $Q \leftarrow$  root node
15: while  $Q$  not empty do
16:    $N \leftarrow Q$ 
17:   if  $N$  is a sum node then
18:      $k^* \leftarrow \arg \max_{k: C_N^k \in \text{ch}'(N)} \left\{ w_{\hat{N},\hat{C}_S^k} \hat{C}_S^k(\mathcal{X}[\text{sc}(\hat{C}_S^k)]) \right\}$ 
19:      $Q \leftarrow C_N^{k^*}$ 
20:      $z_S^* = k^*$ 
21:      $\mathbf{S} \leftarrow \mathbf{S} \setminus \{S\}$ 
22:   else if  $N$  is a product node then
23:      $\forall C \in \text{ch}(N) : Q \leftarrow C$ 
24:   else if  $N$  is a distribution node then
25:      $\mathbf{x}^*[\text{sc}(N)] = \arg \max_{\mathbf{x} \in \mathcal{X}[\text{sc}(N)]} N(\mathbf{x})$ 
26:   end if
27: end while

28: for  $S \in \mathbf{S}$  do
29:    $z_S^* = \arg \max_{k \in \mathcal{X}[Z_S]} \bar{w}_{S,C_S^k}$ 
30: end for

31: return  $\mathbf{x}^*, \mathbf{z}^*$ 
```

Foundations of Sum-Product Networks for probabilistic modeling منبع کمکی: تز دکترا تحت عنوان

نوشته Robert Peharz

قسمت ب) طبق اسلاید های درس، استنتاج marginal برای SPN خطی (tractable) و برای BN به صورت NpHard است. استنتاج conditional هم به صورت مشابه، برای SPN خطی و برای BN به صورت NpHard است. استنتاج MPE برای هر دو NpHard است.

به عنوان منبع اضافی هم برای اطمینان از مقاله learning directed acyclic graph SPNs in sub-quadratic time استفاده کردم.

قسمت ج) تبدیل SPN به bayes net ممکن است: ابتدا SPN را به فرم نرمال در می آوریم. اگر SPN نرمال شده ما S باشد و P یک گره product در این SPN با L فرزند باشد، اگر فرض کنیم k گره sum از ریشه این شبکه تا گره product گفته شده با نام های v_1, v_2, \dots, v_k باشند، داریم:

$$P_S(X_{scope(p)} | H_{v_1} = v_1^*, \dots, H_{v_k} = v_k^*) = \prod_{i=1}^L P_S(x_{scope(p)} | H_{v_1} = v_1^*, \dots, H_{v_k} = v_k^*)$$

حال الگوریتم زیر، نحوه تبدیل bayes net را به SPN نشان میدهد:

Algorithm 3 Build BN Structure

Input: normal SPN \mathcal{S}

Output: BN $\mathcal{B} = (\mathcal{B}_V, \mathcal{B}_E)$

```

1:  $R \leftarrow \text{root of } \mathcal{S}$ 
2: if  $R$  is a terminal node over variable  $X$  then
3:   Create an observable variable  $X$ 
4:    $\mathcal{B}_V \leftarrow \mathcal{B}_V \cup \{X\}$ 
5: else
6:   for each child  $R_i$  of  $R$  do
7:     if BN has not been built for  $\mathcal{S}_{R_i}$  then
8:       Recursively build BN Structure for  $\mathcal{S}_{R_i}$ 
9:     end if
10:  end for
11:  if  $R$  is a sum node then
12:    Create a hidden variable  $H_R$  associated with  $R$ 
13:     $\mathcal{B}_V \leftarrow \mathcal{B}_V \cup \{H_R\}$ 
14:    for each observable variable  $X \in \mathcal{S}_R$  do
15:       $\mathcal{B}_E \leftarrow \mathcal{B}_E \cup \{(H_R, X)\}$ 
16:    end for
17:  end if
18: end if

```

همچنین تبدیل bayes net به SPN نیز امکان پذیر است: ابتدا باید bayes net را به یک درخت junction تبدیل کرد و از روی آن به SPN متناظر برسیم. بخاطر اندازه نمایی SPN تولید شده، اگر bayes net عرض زیادی داشته باشد، SPN تولید

شده مشکلاتی خواهد داشت (قدرت کم در نمایش CPD های محلی). همچنین میتوان یک *bayes net* با *ADDs* را به یک *AC* کامپایل کنیم و از روی آن *SPN* را بدست آوریم. برای اینکار ابتدا باید عملیات های معمول استفاده شده در *VE* را بدست آوریم. دو عملیات معمول، اول ضرب دو *symbolic ADDs* و دومی *summing-out* یک متغیر پنهان است. لذا داریم:

Algorithm 6 Multiplication of two symbolic ADDs, \otimes

Input: Symbolic ADD $\mathcal{A}_{X_1}, \mathcal{A}_{X_2}$
Output: Symbolic ADD $\mathcal{A}_{X_1, X_2} = \mathcal{A}_{X_1} \otimes \mathcal{A}_{X_2}$

- 1: $R_1 \leftarrow \text{root of } \mathcal{A}_{X_1}, R_2 \leftarrow \text{root of } \mathcal{A}_{X_2}$
- 2: **if** R_1 and R_2 are both variable nodes **then**
- 3: **if** $R_1 = R_2$ **then**
- 4: Create a node $R = R_1$ into \mathcal{A}_{X_1, X_2}
- 5: **for each** $r \in \text{dom}(R)$ **do**
- 6: $\mathcal{A}_{X_1}^r \leftarrow \text{Ch}(R_1)|_r$
- 7: $\mathcal{A}_{X_2}^r \leftarrow \text{Ch}(R_2)|_r$
- 8: $\mathcal{A}_{X_1, X_2}^r \leftarrow \mathcal{A}_{X_1}^r \otimes \mathcal{A}_{X_2}^r$
- 9: Link \mathcal{A}_{X_1, X_2}^r to the r th child of R in \mathcal{A}_{X_1, X_2}
- 10: **end for**
- 11: **else**
- 12: $\mathcal{A}_{X_1, X_2} \leftarrow \text{create a symbolic node } \otimes$
- 13: Link \mathcal{A}_{X_1} and \mathcal{A}_{X_2} as two children of \otimes
- 14: **end if**
- 15: **else if** R_1 is a variable node and R_2 is \otimes **then**
- 16: **if** R_1 appears as a child of R_2 **then**
- 17: $\mathcal{A}_{X_1, X_2} \leftarrow \mathcal{A}_{X_2}$
- 18: $\mathcal{A}_{X_1, X_2}^{R_1} \leftarrow \mathcal{A}_{X_1} \otimes \mathcal{A}_{X_2}^{R_1}$
- 19: **else**
- 20: Link \mathcal{A}_{X_1} as a new child of R_2
- 21: $\mathcal{A}_{X_1, X_2} \leftarrow \mathcal{A}_{X_2}$
- 22: **end if**
- 23: **else if** R_1 is \otimes and R_2 is a variable node **then**
- 24: **if** R_2 appears as a child of R_1 **then**
- 25: $\mathcal{A}_{X_1, X_2} \leftarrow \mathcal{A}_{X_1}$
- 26: $\mathcal{A}_{X_1, X_2}^{R_2} \leftarrow \mathcal{A}_{X_2} \otimes \mathcal{A}_{X_1}^{R_2}$
- 27: **else**
- 28: Link \mathcal{A}_{X_2} as a new child of R_1
- 29: $\mathcal{A}_{X_1, X_2} \leftarrow \mathcal{A}_{X_1}$
- 30: **end if**
- 31: **else**
- 32: $\mathcal{A}_{X_1, X_2} \leftarrow \text{create a symbolic node } \otimes$
- 33: Link \mathcal{A}_{X_1} and \mathcal{A}_{X_2} as two children of \otimes
- 34: **end if**
- 35: Merge connected product nodes in \mathcal{A}_{X_1, X_2}

Algorithm 7 Summing-out a hidden variable H from \mathcal{A} using \mathcal{A}_H, \oplus

Input: Symbolic ADDs \mathcal{A} and \mathcal{A}_H
Output: Symbolic ADD with H summed out

- 1: **if** H appears in \mathcal{A} **then**
- 2: Label each edge emanating from H with weights obtained from \mathcal{A}_H
- 3: Replace H by a symbolic \oplus node
- 4: **end if**

سپس از روی *BN* با الگوریتم زیر، *SPN* را بازیابی میکنیم:

Algorithm 8 Variable Elimination for BN with ADDs

Input: BN \mathcal{B} with ADDs for all observable variables and hidden variables

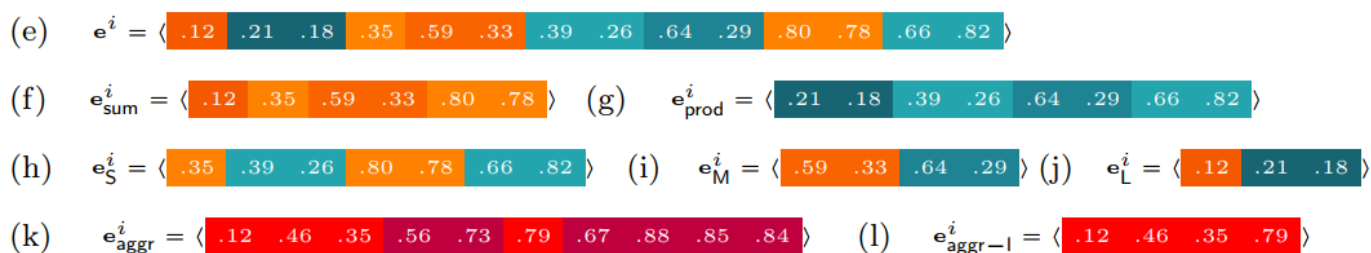
Output: Original SPN \mathcal{S}

- 1: $\pi \leftarrow$ the inverse topological ordering of all the hidden variables present in the ADDs
 - 2: $\Phi \leftarrow \{\mathcal{A}_X \mid X \text{ is an observable variable}\}$
 - 3: **for** each hidden variable H in π **do**
 - 4: $P \leftarrow \{\mathcal{A}_X \mid H \text{ appears in } \mathcal{A}_X\}$
 - 5: $\Phi \leftarrow \Phi \setminus P \cup \{\oplus_H \otimes_{\mathcal{A} \in P} \mathcal{A}\}$
 - 6: **end for**
 - 7: **return** Φ
-

منبع کمکی: مقاله تحت عنوان *on the relationship between sum-product networks and Bayesian networks*

قسمت د) دو مقاله راهکار هایی برای هندل کردن SPN های پیوسته پیدا کردم. در مقاله اول تحت عنوان *on the latent variable interpretation in sum-product networks* یک الگوریتم EM برای یادگیری پارامتر های اینچنین SPN هایی معرفی شده است. همچنین در مقاله *online algorithms for sum-product networks with continuous variables* الگوریتم های *online Bayesian moment matching* برای SPN های با توزیع *gaussian leaf* تعریف شده که میتواند این SPN هارا هندل کند. این الگوریتم ها و محاسبات ریاضی شان بسیار گسترده و خارج سرفصل بوده و به همین دلیل آن هارا در پاسخ ام نیاوردم.

قسمت ه) ابتدا یک نمونه Xi به SPN خود وارد میکنیم (a) و activation های گره های آن را جمع آوری میکنیم (b). سپس فیلتر های مختلفی را اعمال کرده و embedding تشکیل میدهم برای مثال جمع آوری activation های تمام گره های داخلی (e) یا فیلتر کردن بر اساس نوع گره مثلا جمع آوری activation های گره های sum (f) یا گره های product (g) یا فیلتر کردن گره ها با معیار طول scope کوچک (h) یا متوسط (i) یا بزرگ (j) یا aggregate کردن گره های با scope برابر (k) (همان گره های sum قرمز رنگ در c و ارزیابی شده در d) و یا در آخر با استفاده از فقط گره های داخلی (l):



سوال ۲ -

$$\min_{\theta_{XY}} - \frac{1}{|\mathcal{D}|} \sum_{(x,y) \in \mathcal{D}} \log P(y|x; \theta_{XY});$$

$$\min_{\theta_{YX}} -\frac{1}{|\mathcal{D}|} \sum_{(x,y) \in \mathcal{D}} \log P(x|y; \theta_{YX}).$$

متعاقبا، میتوانیم دو تابع پیش بینی زیر را برای فرآیند $primal$ و $dual$ ارائه دهیم:

$$f(x; \theta_{XY}) \triangleq \arg \max_{y' \in \mathcal{Y}} P(y'|x; \theta_{XY}),$$

$$g(y; \theta_{YX}) \triangleq \arg \max_{x' \in \mathcal{X}} P(x'|y; \theta_{YX}).$$

در DSL ، به صورت مشترک آموزش هم برای $primal$ هم برای $dual$ انجام میشود که بهینه سازی زیر نتیجه میشود:

$$\text{objective 1: } \min_{\theta_{XY}} \frac{1}{|\mathcal{D}|} \sum_{(x,y) \in \mathcal{D}} \ell_1(f(x; \theta_{XY}), y),$$

$$\text{objective 2: } \min_{\theta_{YX}} \frac{1}{|\mathcal{D}|} \sum_{(x,y) \in \mathcal{D}} \ell_2(g(y; \theta_{YX}), x),$$

$$\text{s.t. } P(x)P(y|x; \theta_{XY}) = P(y)P(x|y; \theta_{YX}), \forall (x, y) \in \mathcal{D},$$

اگر فرموله سازی لاگرانژ آن را بنویسیم و هر دو $objective function$ ها را تلفیق کنیم، داریم:

$$\ell_{dsl} = (\log \hat{P}(x) + \log P(y|x; \theta_{XY}) - \log \hat{P}(y) - \log P(x|y; \theta_{YX}))^2.$$

الگوریتم کلی DSL برابر است با:

Algorithm 1 Dual supervise learning algorithm

Require: : Marginal distributions $\hat{P}(x)$ and $\hat{P}(y)$; Lagrange parameters λ_{XY} and λ_{YX} ; optimizers Opt_1 and Opt_2 ;

repeat

Sample a minibatch of m pairs $\{(x_j, y_j)\}_{j=1}^m$;

Calculate the gradients as follows:

$$G_f = \nabla_{\theta_{XY}} (1/m) \sum_{j=1}^m [\ell_1(f(x_j; \theta_{XY}), y_j) + \lambda_{XY} \ell_{dsl}(x_j, y_j; \theta_{XY}, \theta_{YX})];$$

$$G_g = \nabla_{\theta_{YX}} (1/m) \sum_{j=1}^m [\ell_2(g(y_j; \theta_{YX}), x_j) + \lambda_{YX} \ell_{dsl}(x_j, y_j; \theta_{XY}, \theta_{YX})];$$

Update the parameters of f and g :

$$\theta_{XY} \leftarrow Opt_1(\theta_{XY}, G_f), \theta_{YX} \leftarrow Opt_2(\theta_{YX}, G_g).$$

until models converge

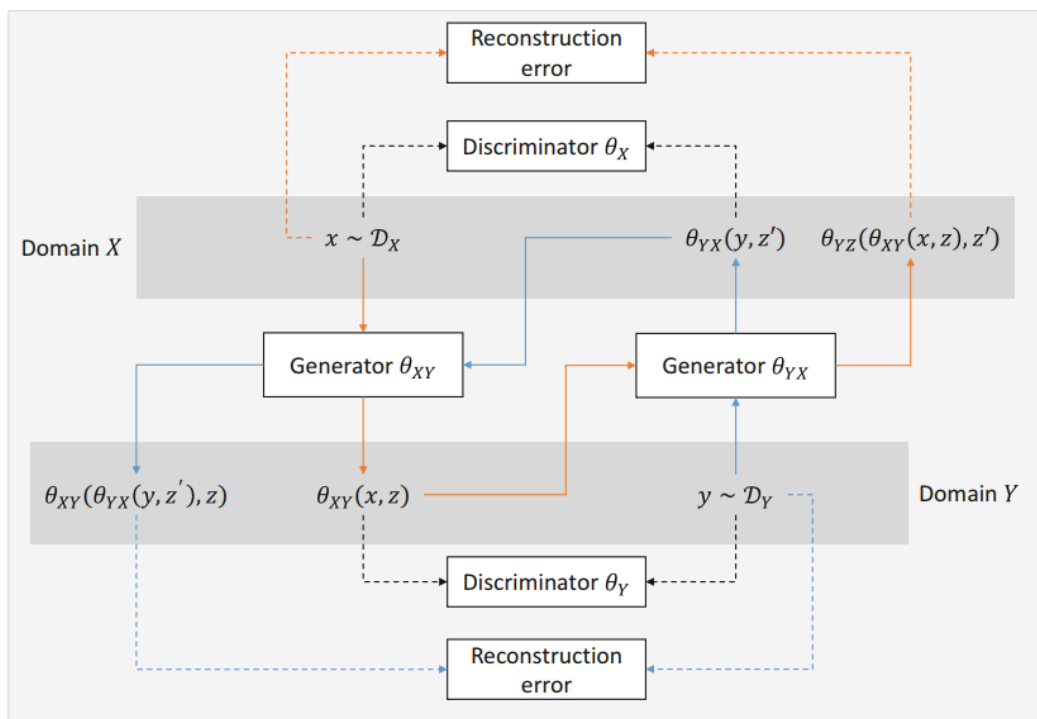
در *Dual unsupervised learning* معماری *encoder-decoder* لازم است زیرا برچسب داده ها یا جمله های دوزبانه وجود ندارد و باید برای مثال از زبان مبدا، جمله . برای مثال در کاربرد ترجمه ماشینی، این نوع یادگیری بسیار پرکاربرد است. معماری آن به این صورت است که نسخه نویزی جمله مبدا توسط *encoder* به یک *embedding* ثابت تبدیل شده و توسط *decoder* به جمله مقصد تبدیل میشود. لذا تابع هزینه برابر است با :

$$\ell_{dae}(\theta_{en}^l, \theta_{de}^l) = \frac{1}{|\mathcal{M}_l|} \sum_{x \in \mathcal{M}_l} \Delta(x, \theta_{de}^l(\theta_{en}^l(c(x)))),$$

لذا در فرآیند *dual* آن تابع هزینه میشود:

$$\begin{aligned} \ell_{dual}(\theta_{en}^l, \theta_{de}^l) &= \frac{1}{|\mathcal{M}_A|} \sum_{x \in \mathcal{M}_A} \Delta(x, \theta_{de}^A(\theta_{en}^B(c(\hat{y})))) \\ &+ \frac{1}{|\mathcal{M}_B|} \sum_{x \in \mathcal{M}_B} \Delta(x, \theta_{de}^B(\theta_{en}^A(c(\hat{y}))))), \end{aligned}$$

معماری کلی *dual reconstruction* برابر است با:



هر دو مدل ترجمه تصویری (*primal, dual*) برای کمینه کردن خطای بازسازی تصویر اوريجینال، آموزش داده میشوند. *DualGan* از *L1 distance* استفاده میکند زیرا *L2 distance* باعث تار شدن تصویر بازسازی شده میشود. تابع هزینه آن برابر است با:

$$\begin{aligned} \ell(x, y; \theta_{YX}, \theta_{XY}) &= \lambda_X \|x - \theta_{YX}(\theta_{XY}(x, z), z')\| \\ &\quad + \lambda_Y \|y - \theta_{XY}(\theta_{YX}(y, z'), z)\| \\ &\quad - \theta_Y(\theta_{XY}(x, z)) - \theta_X(\theta_{YX}(y, z')). \end{aligned}$$

و همچنین الگوریتم آموزش این شبکه برابر است با:

Algorithm 1 DualGAN training procedure

Require: Two image collections \mathcal{D}_X and \mathcal{D}_Y , clipping parameter c , batch size m , and d

- 1: Randomly initialize θ_{YX} , θ_{XY} , θ_X , and θ_Y
 - 2: **repeat**
 - 3: **for** $t = 1, \dots, d$ **do**
 - 4: Sample images $\{x^{(k)}\}_{k=1}^m \subseteq \mathcal{D}_X$, $\{y^{(k)}\}_{k=1}^m \subseteq \mathcal{D}_Y$
 - 5: Update θ_X to minimize $\frac{1}{m} \sum_{k=1}^m \ell(x^{(k)}, y^{(k)}; \theta_X)$
 - 6: Update θ_Y to minimize $\frac{1}{m} \sum_{k=1}^m \ell(x^{(k)}, y^{(k)}; \theta_Y)$
 - 7: $\text{clip}(\theta_X, -c, c)$, $\text{clip}(\theta_Y, -c, c)$
 - 8: **end for**
 - 9: Sample images $\{x^{(k)}\}_{k=1}^m \subseteq \mathcal{D}_X$, $\{y^{(k)}\}_{k=1}^m \subseteq \mathcal{D}_Y$
 - 10: Update θ_{YX} and θ_{XY} to minimize $\frac{1}{m} \sum_{k=1}^m \ell(x^{(k)}, y^{(k)}; \theta_{YX}, \theta_{XY})$
 - 11: **until** convergence
-

منبع کمکی هر دو قسمت: کتاب *Dual learning* نوشته Tao Qin

سوال ۳ -

قسمت الف) در *Q-learning*، *updated policy* با *behavior policy* متفاوت است. لذا میتوانیم برای مثال ابتدا اپیزود هایی از محیط را اجرا کرده و مدلی بدست آوریم و سپس در آینده این مدل را استفاده کنیم (طبق گفته های کلاس دکتر رهبان درس یادگیری ماشین). از طرفی در *Q-learning*، برای عمل های آینده ما یک پاداش تخمین میزنیم و مقدار جدیدی به حالت های جدید میدهیم بدون اینکه از *greedy policy* استفاده کرده باشیم. زیرا به این دلایل، *off policy* است زیرا *off policy* یعنی یادگیرنده از این نوع، مقادیر *optimal policy* را به صورت مستقل از عمل های عامل یاد میگیرد.

یکی از مزایای *off policy* این است که *updated policy* و *behavior policy* میتوانند یکی *deterministic* باشد و دیگری نباشد. یعنی مستقل بودن و جدا بودن این دو این کمک را به ما میکند. از طرف دیگر اگر عامل در محیطی *offline*

بخواید یادگیری انجام دهد و *explore* زیادی نداشته باشد، *off policy* بهتر عمل میکند. در دنیای واقعی هم از لحاظ هزینه، به صرفه تر است.

قسمت ب) در Q-learning ما Q-table داریم اما در DQN، به جای آن، یک شبکه عصبی عمیق داریم. در Q-learning ما علاقه داریم Q-function را بیابیم اما از طرفی با بزرگ شدن اعمال و فضای مسئله، سایز Q-table بزرگ شده و آموزش با Q-learning تقریباً غیر ممکن میشود که در این موقعیت، DQN بخاطر اینکه با شبکه عصبی کار میکند، مشکل را حل میکند. از طرفی Q-learning با حالت های پیوسته کار نمیکند اما DQN مشکلی ندارد. اما Q-learning همیشه همگرا میشود اما DQN ممکن است همگرا نشود. از طرفی Q-learning نوعی الگوریتم online است اما DQN اینطور نیست.

قسمت ج)

استفاده از **experience replay buffer**: آموزش روی داده های ترتیبی مرتبط (correlated) میتواند منجر به ناپایداری مدل شود. لذا از این بافر برای ذخیره نمونه برداری (transition tuple) ها استفاده میشود و سپس به صورت رندوم از این بافر نمونه برداری میکنیم به صورت minibatch و اینگونه ناپایداری کاهش میابد.

استفاده از تکنیک **freeze the target network**: در Q-learning ما با استفاده از یک حدس، حدس بعدی را میزنیم. لذا اینکار باعث correlation مضری میتواند بشود و شبکه عصبی نمیتواند بین دو چیز خیلی شبیه به هم تمایز قائل شود. وقتی ما میخواهیم Q value را برای یک حالت با آپدیت کردن پارامتر های شبکه، بهینه کنیم، ناخودآگاه به مقادیر حالت های دیگر و نزدیک به آن هم دست میزنیم و این ناپایداری ایجاد میکند. برای حل این مشکل یک کپی از کل شبکه عصبی گرفته و آن را target network مینامیم. از Q value پیش بینی شده توسط target network، برای backpropagation و یادگیری شبکه اصلی و اولیه، استفاده میشود. توجه شود پارامتر های target network آموزش داده نمیشوند یعنی freeze هستند اما به صورت متناوب و هر چند مدت یک بار، با شبکه اصلی sync میشوند.

قسمت د) در DDQN، از دو مدل شبکه عصبی یکسان استفاده میکنند. اولی هنگام experience replay یادگیری انجام میدهد و دیگری یک کپی از آخرین اپیزود اولین مدل است که Q value با این شبکه محاسبه میشود. در DQN، Q value محاسبه شده برابر است با بیشترین Q value حالت بعدی به اضافه پاداش کسب شده. لذا اگر هر دفعه Q value برای یک حالت خاص مقدار زیادی باشد، مقداری که خروجی شبکه عصبی به ازای آن حالت خاص به ما میدهد، بزرگ و بزرگ تر میشود. حال اگر برای یک حالتی، مقدار عمل a از عمل b بیشتر باشد، پس همیشه a انتخاب میشود و اگر برای بعضی از memory experience ها b بهتر از a باشد، شبکه نمیتواند b را انتخاب کند زیرا اختلاف بین مقادیر خروجی زیاد است. لذا باید این

اختلاف را کم کنیم. در DDQN در شبکه دوم، اختلاف خروجی ها پایین است و به همین دلیل از این شبکه برای بدست آوردن Q value استفاده میکنیم. لذا در DQN ما تخمین بالاتر پاداش واقعی را داریم که در DDQN با جدا کردن action selection و action evaluation، این مشکل را حل میکنیم.

تفاوت شبکه دوم در DDQN با target network این است که target network کپی کامل از شبکه اصلی است و برای محاسبه Q value از آن استفاده نمیشود به صورت مستقیم اما در DDQN شبکه دوم، کپی آخرین اپیزود از شبکه اول است و برای محاسبه Q value استفاده میشود.

سوال ۴ - در نوتبوک آورده شده است.

سوال ۵ -

قسمت الف) الگوریتم DDPG ترکیبی از الگوریتم های DPG و DQN است. به این صورت که از دو شبکه target استفاده میکند که با یکی یک سری target را تخمین زده و با دیگری از target های تخمین زده شده، یادگیری را انجام میدهیم. در این الگوریتم از experience replay و slow-learning ای که در DQN بود نیز استفاده میشود. برای آپدیت شبکه عامل، از sampled policy gradient استفاده میشود به این صورت:

$$\nabla_{\theta^{\mu}} J \approx \frac{1}{N} \sum_i \nabla_a Q(s, a | \theta^Q) |_{s=s_i, a=\mu(s_i)} \nabla_{\theta^{\mu}} \mu(s | \theta^{\mu}) |_{s_i}$$

قسمت ب) برای حل مشکل exploration از یک فرآیند نويز به صورت correlated (بخاطر explore خوب در محیط فیزیکی) استفاده میشود. برای مثال میتوان از فرآیند Ornstein-uhlenbeck استفاده کرد. این فرآیند velocity را در حرکت براونی مدل میکند. از طرف دیگر الگوریتم ما off-policy است یعنی بحث exploration از الگوریتم یادگیری مستقل است. لذا یک exploration policy به شکل زیر تعریف میکنیم که حاصل جمع actor policy و همان نويز گفته شده است:

$$\mu'(s_t) = \mu(s_t | \theta_t^{\mu}) + \mathcal{N}$$

سوال ۶ - در نوتبوک آورده شده است.