

2.2 a)

$$H_m = WH_{m-1} + UX_m$$

$$= W(WH_{m-2} + UX_{m-1}) + UX_m$$

$$= W^2 H_{m-2} + WUX_{m-1} + UX_m$$

$$= W^3 H_{m-3} + W^2 UX_{m-2} + WUX_{m-1} + UX_m$$

$$= \dots$$

since  $H_0 = 0$ :

$$H_m = \sum_{t=1}^m W^{m-t} UX_t$$

this means

$$Y_m = V \cdot \sum_{t=1}^m W^{m-t} UX_t$$

$$\left( Y_m = \sum_{t=1}^m V \cdot W^{m-t} UX_t \right)_K \rightarrow \begin{cases} 1 \leq m \leq n \\ 1 \leq K \leq b \end{cases}$$

$\equiv$

$$g(\langle X_t \rangle_{t=1}^n, W, U, V)_{K,m}$$



2.2b)

$$\begin{aligned} \Rightarrow g(\langle X_t \rangle_{t=1}^n, W, U, V)_{km} &= \left( \sum_{t=1}^m V W^{m-t} U X_t \right)_k \\ \Rightarrow f(z)_i &= \sum_{j \in \mathcal{J}} \Theta_{ij} z_j, \quad \forall z \in \mathbb{R}^{\mathcal{J}}, \forall i \in \mathcal{I} \end{aligned}$$

$X \equiv z \Rightarrow$  Given that  $X_t \in \mathbb{R}^d \leadsto X \in \mathbb{R}^{d \times n}$

this implies:

$$\boxed{Z_{(u,v)} = X_{u,v}, \quad \mathcal{J} = \{u,v \mid 1 \leq u \leq d, 1 \leq v \leq n\}}$$

$$g(\langle X_t \rangle_{t=1}^n, W, U, V) \equiv f(z) \Rightarrow$$

$$\boxed{\mathcal{I} = \{(k,m) \mid 1 \leq k \leq b, 1 \leq m \leq n\}}$$

$$\text{so: } g(\langle X_t \rangle_{t=1}^n, W, U, V)_{km} = \sum_{t=1}^m (V W^{m-t} U X_t)_k = (Y_m)_k$$

$$= \sum_{t=1}^m (V W^{m-t} U X_t)_k \quad \begin{cases} 1 \leq k \leq b \\ 1 \leq m \leq n \end{cases}$$

$$= \sum_{v=1}^m V W^{m-v} U X_{v,v}$$

$$= \sum_{v=1}^m V W^{m-v} \left( \sum_{u=1}^d U_{ju} \cdot X_{uv} \right)_j \quad \forall j=1, \dots, d$$

$$= \underbrace{\sum_{j \in \mathcal{J}} \Theta_{ij}}_{\sum_{j \in \mathcal{J}} \Theta_{ij}} \underbrace{\sum_{v=1}^m V W^{m-v} U_{ju} X_{uv}}_{Z_j} \quad \forall j=1, \dots, d$$



2.3  
a) GCN  $\left\{ \begin{array}{l} \hat{B} = 0, (i,j) \in E, \forall i \in V \\ W_{ij} \in \{0,1\}, (i,j) \in E \end{array} \right.$

2.3  
a) GCN  $\left\{ \begin{array}{l} \hat{B} = 0, (i,j) \in E, \forall i \in V \\ W_{ij} \in \{0,1\}, (i,j) \in E \end{array} \right.$

K diff topologies which propagate messages (relational GCN):

$$V' \leftarrow \sigma \left( \sum_{k=1}^K W^k V B^k \right)$$

→ Let  $W^k$  be the matrix representation for all Kernel values defined in the CNN  $\in \mathbb{R}^{(K \times 3) \times (h_y \times h_x)}$  } 2 dimensional equiv

→ Let  $V$  be the vector equivalent ~~the~~ to the ~~the~~ input variables  $X$  with  ~~$\# \# u = Vu \in \mathbb{R}^p$~~   $F = (H - Ky + 1) \times (W_F - K + 1)$

→ let  $B^k$  be the ~~current~~ params to optimize with the adewate optimizer input  $X$

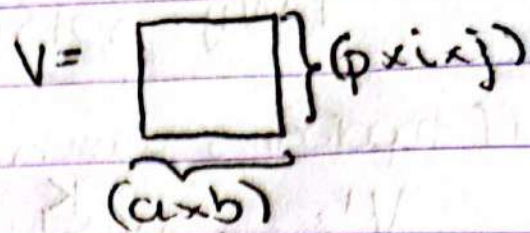
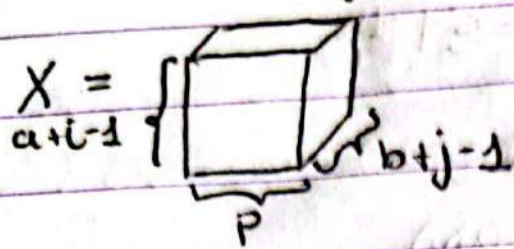
$\in \mathbb{R} \in \underbrace{(H-K_y+1) \times (W-K+1)}_{1 \text{ dimensional}}^T$

1 dimensional equivalent to

This GCN representation would be equivalent to CNN presented



$$2.3 \text{ b) } V \begin{pmatrix} a \\ b \end{pmatrix} \begin{pmatrix} p \\ i \\ j \end{pmatrix} = X_{p, a+i-1, b+j-1}$$



Weight matrix  $W = \begin{pmatrix} p \\ i \\ j \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} =$

Deriving the matrix with a single-head attention layer ~~is~~ would imply that we don't know the exact topology to use for the problem suggested.

Given that:

$$V' \leftarrow \sigma \left( \sum_{k=1}^K W^k V B^k \right)$$

$$\text{if } K=1$$

$$V' \leftarrow \sigma(WVB)$$

So

$$W \leftarrow (B^{-1} V^{-1} V')$$



2.4 a) We know that:  $V^*(s_0) = r(s_0) + \gamma V^*(s')$  so:

Best policy without memory:

Given that  $p(c=\text{green}) = 1/3$  the best policy is the one that goes to  $s_8$  [ $p(c=\text{red}) = 2/3$ ] and the policy value is

$$V^*(s_0) = \overset{+1 \text{ reward}}{2/3} \cdot \gamma^4 - \overset{-1 \text{ reward}}{1/3} \gamma^4 \quad (\text{since we do 4 steps})$$
$$\boxed{= \gamma^4/3}$$

Best policy with memory: Given that we have a memory we should take ~~advantage~~ advantage of state  $s_2$  so:

$$\boxed{V^*(s_0) = \gamma^8}$$

Since you first go to  $s_2$   
In order to know what is the value of  $c$  and be sure were is the possible value

The range of values for which  $\gamma$  is better a memory-less policy is  $[0 \dots \sqrt[4]{1/3}]$

since:  $\gamma^4/3 = \gamma^8 \rightarrow \text{equality}$

$$1 = 3\gamma^4$$

$$\boxed{\gamma = \sqrt[4]{1/3}}$$



2.4 b) The smallest set of past observations would be a set of 4. since the decision ~~is~~ the decision ~~is~~ in which the policy is going to influence ~~is~~  $s_4$  is related to state  $s_4$  and from  $s_2$  to  $s_4$  there are 4 steps. When you are in ~~all~~ other states  $s_5, s_7$  the policy is clear always. ( $s_5 \rightarrow s_6$  &  $s_7 \rightarrow s_8$ )

2.4 c)

Without memory

There is no effect in the policy since the action associated into a state always remains the same. Analyze both cases.

50% go in desired direction,  $\rightarrow$  the policy remains the same

50% chance you stay in current state,  $\rightarrow$  you are not in the desired ~~stay~~ state but the action to go there remains the same

With memory

The effect on the policy is the same, it only affects the capacity of the memory that should increase in order to provide  $V^*$ .



2.5 a)

$$\begin{aligned}\nabla_{\theta} E[f(s, a) | a \sim \pi_{\theta}(a|s)] &= \int \nabla_{\theta} [f(s, a) p(s, a)] da \\ &= \int f(s, a) \nabla_{\theta} p(s, a) da \\ &\text{since } \nabla_{\theta} \ln(p(s, a)) = \frac{1}{p(s, a)} \nabla_{\theta} p(s, a) \\ &= \int f(s, a) \nabla_{\theta} \ln p(s, a) p(s, a) da \\ &= E[f(s, a) \nabla_{\theta} \ln \pi_{\theta}(s|a) | a \sim \pi_{\theta}(s, a)]\end{aligned}$$

$$\begin{aligned}2.6 b) \nabla_{\theta} E_{\pi_{\theta}}[R_t | s_t^{\theta}] &= \nabla_{\theta} E_{\pi_{\theta}}[r(s_t, a_t) + \gamma R_{t+1} | s_t^{\theta}] \\ &= \nabla_{\theta} (E_{\pi_{\theta}}[r(s_t, a_t) | s_t^{\theta}] + \gamma E[R_{t+1} | s_t^{\theta}]) \\ &= \nabla_{\theta} E_{\pi_{\theta}}[r(s_t, a_t) | s_t^{\theta}] + \nabla_{\theta} \gamma E[R_{t+1} | s_t^{\theta}] \\ &= \nabla_{\theta} E_{\pi_{\theta}}[r(s_t, a_t) | s_t^{\theta}] + \gamma E[R_{t+1} \nabla_{\theta} \ln \pi_{\theta}(a_t | s_t) | s_t^{\theta}]\end{aligned}$$

$$\begin{aligned}\nabla_{\theta} E_{\pi_{\theta}}[r(s_t, a_t) | s_t^{\theta}] &= \nabla_{\theta} \int r(s_t, a_t) p(s_t, a_t) da \\ &= \int r(s_t, a_t) \nabla_{\theta} p(s_t, a_t) da \\ &= \int r(s_t, a_t) \nabla_{\theta} p(s_t, a_t) p(s_t, a_t) da \\ &= E_{\pi_{\theta}}[r(s_t, a_t) \nabla_{\theta} \ln \pi_{\theta}(a_t | s_t) | s_t^{\theta}]\end{aligned}$$

$$\begin{aligned}&= E_{\pi_{\theta}}[r(s_t, a_t) \nabla_{\theta} \ln \pi_{\theta}(a_t | s_t) | s_t^{\theta}] + \\ &\gamma E_{\pi_{\theta}}[R_{t+1} \nabla_{\theta} \ln \pi_{\theta}(a_t | s_t) | s_t^{\theta}]\end{aligned}$$



$$2.5c) \nabla_{\theta} J[\pi_{\theta}] = \nabla_{\theta} E_{\pi_{\theta}}[R_0]$$

$$\left( \begin{aligned} &= \nabla_{\theta} E_{\pi_{\theta}} \left[ \sum_{k=1}^{n-1} \gamma^{k-t} r(s_0, a_0) \right] \\ &= \nabla_{\theta} E \left[ \sum_{k=1}^{n-1} \gamma^{k-t} r(s_0 | a_0) \pi_{\theta}(a_0 | s_0) \right] \\ &= \nabla_{\theta} \sum_{k=1}^{n-1} E[r^{k-t} | r(s_0 | a_0) \pi_{\theta}(a_0 | s_0)] \end{aligned} \right)$$

$$= \nabla_{\theta} E_{\pi_{\theta}} [R_0 + \gamma R_1 + \gamma^2 R_2 + \dots]$$

$$= E_{\pi_{\theta}} [\nabla_{\theta} R_0 + \nabla_{\theta} \gamma R_1 + \nabla_{\theta} \gamma^2 R_2 + \dots]$$

$\Rightarrow$  this can be also expressed as

$$= E_{\pi_{\theta}} \nabla_{\theta} E_{\pi_{\theta}} \left[ \sum_{t=0}^{n-1} \gamma^t R_t \right]$$

$$= E_{\pi_{\theta}} \left[ \sum_{t=0}^{n-1} \nabla_{\theta} \gamma^t R_t \right]$$

$$\boxed{= E_{\pi_{\theta}} \left[ \sum_{t=0}^{n-1} \gamma^t R_t \nabla_{\theta} \ln \pi_{\theta}(a_t | s_t) \right]}$$

Given previously  
exercise 2.5a)

2.5a)