

**1 Gradient Descent:****Part1:**

Q1. We have given the function:-  

$$\text{argmin}_w f(w^*) + \langle w - w^*, \nabla f(w^*) \rangle + \frac{\lambda}{2} \|w - w^*\|^2$$

To optimize this, let's take the first order derivative and assume

$$w - w^* = x$$

$$f(x) = \text{argmin}_w f(w^*) + \langle x, \nabla f(w^*) \rangle + \frac{\lambda}{2} \|x\|^2$$

$$f'(x) = \nabla f(w^*) + \frac{\lambda}{2} \cdot 2 \|x\|$$

$$f'(x) = \nabla f(w^*) + \lambda \|x\|$$

$$f'(x) = 0$$

$$\lambda = - \frac{\nabla f(w^*)}{\|x\|}$$

$$\|x\| = - \frac{\nabla f(w^*)}{\lambda}$$

$$w - w^* = - \frac{\nabla f(w^*)}{\lambda}$$

$$w = w^* - \frac{\nabla f(w^*)}{\lambda} \quad \text{--- ①}$$

Now, we also know that:

$$w^{t+1} = w^t - \eta \nabla f(w^t)$$

If we compare the eqn ① with above, we can say that

$$\frac{1}{\lambda} = \eta$$

From this we can say that learning rate is inversely proportional to the regularization.

**What does it mean?**

Regularization means controlling overfitting, and one of the ways we can relate regularization with overfitting is by training the model over multiple iterations. The more iterations we have the less overfitting we have.

Now, our relationship tells us that, if we have:

① large  $\eta$ , then in few iterations (or with less regularization) you can optimize your model.

② low  $\eta$ , then in large iterations you can optimize your model.

## Part2:

(2)

We have given:

a sequence of vectors  $v_1, v_2, \dots, v_T$

update eq<sup>t</sup>

$$w^{t+1} = w^t - \eta v_t \quad \text{--- (1)}$$

$$w^{(0)} = 0 \quad \text{--- (2)}$$

To show:-

$$\sum_{t=1}^T \langle w^t - w^*, v_t \rangle \leq \frac{\|w^*\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|v_t\|^2$$

We can write  $\langle w^t - w^*, v_t \rangle$  as

$$\langle w^t - w^*, v_t \rangle = \frac{1}{\eta} \langle w^t - w^*, \eta v_t \rangle \quad \left[ \begin{array}{l} \text{multiply \& divide} \\ \eta \quad 1 \end{array} \right]$$

$$= \frac{1}{2\eta} (-\|w^t - w^* - \eta v_t\|^2 + \|w^t - w^*\|^2 + \eta^2 \|v_t\|^2)$$

$$= \frac{1}{2\eta} (-\|w^{t+1} - w^*\|^2 + \|w^t - w^*\|^2 + \eta^2 \|v_t\|^2)$$

Now, we can sum the equation over  $t$ .

$$\sum_{t=1}^T \langle w^t - w^*, v_t \rangle = \frac{1}{2\eta} \sum_{t=1}^T (-\|w^{t+1} - w^*\|^2 + \|w^t - w^*\|^2 + \eta^2 \|v_t\|^2)$$

$$= \frac{1}{2\eta} \sum_{t=1}^T (-\|w^{t+1} - w^*\|^2 + \|w^t - w^*\|^2) + \sum_{t=1}^T \frac{\eta}{2} \|v_t\|^2$$

The first term on the right hand side collapses to

$$\|w^1 - w^*\|^2 - \|w^{T+1} - w^*\|^2$$

$$\sum_{t=1}^T \langle w^t - w^*, v_t \rangle = \frac{1}{2\eta} (\|w^1 - w^*\|^2 - \|w^{T+1} - w^*\|^2) + \frac{\eta}{2} \sum_{t=1}^T \|v_t\|^2$$

$$\leq \frac{1}{2\eta} \|w^1 - w^*\|^2 + \frac{\eta}{2} \sum_{t=1}^T \|v_t\|^2$$

$$\leq \frac{1}{2\eta} \|w^*\|^2 + \frac{\eta}{2} \sum_{t=1}^T \|v_t\|^2$$

### Part3:

⑤

To show  $\bar{w} = \frac{1}{T} \sum_{t=1}^T w^{(t)}$

we have all the weights as  $w_1, w_2, \dots, w^T$

let calculate their mean  $\bar{w}$

$$\bar{w} = \frac{w_1 + w_2 + w_3 + \dots + w^T}{T}$$

$$\bar{w} = \frac{1}{T} \sum_{t=1}^T w^{(t)}$$

Next:

To prove

$$f(\bar{w}) - f(w^*) \leq \frac{1}{T} \sum_{t=1}^T \langle w^t - w^*, \nabla f(w^t) \rangle$$

let's first Jensen's inequality

it states that:

If  $g(x)$  is a convex function on  $\mathbb{R}_x$ , and  $E(g(x))$  and  $g(E(x))$  are finite, then

$$E(g(x)) \geq g(E(x))$$

Now, do we come up with above eqn

we know that for every random variable  $x$

$$V(x) = E(x^2) - (E(x))^2 \geq 0$$

$$E(x^2) \geq (E(x))^2$$

Now, take any convex function,  $g(x) = x^2$  [for e.g.]

Jensen's inequality states that

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

using the same above inequality, we can write that:

$$\begin{aligned} f(\bar{w}) - f(w^*) &= f\left(\frac{1}{T} \sum_{t=1}^T w^t\right) - f(w^*) \\ &\leq \frac{1}{T} \sum_{t=1}^T (f(w^t) - f(w^*)) \\ &= \frac{1}{T} \sum_{t=1}^T (f(w^t) - f(w^*)) \end{aligned}$$

for every  $t$ , because of convexity, we have

$$f(w^t) - f(w^*) \leq \langle w^t - w^*, \nabla f(w^t) \rangle$$

therefore

$$f(\bar{w}) - f(w^*) \leq \frac{1}{T} \sum_{t=1}^T \langle w^t - w^*, \nabla f(w^t) \rangle$$

We have proved previously

$$\sum_{t=1}^T \langle w^t - w^*, v_t \rangle \leq \frac{\|w^*\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|v_t\|^2$$

upper bounding  $\|w^*\| \rightarrow B$ ,  $\|v_t\| \rightarrow \rho$ ,  $\eta \rightarrow \sqrt{\frac{B^2}{\rho^2 T}}$ , dividing by  $T$

$$\begin{aligned} &\leq \frac{1}{T} \left( \frac{B^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \rho^2 \right) \\ &\leq \frac{1}{T} \left( \frac{B^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \rho^2 \right) \\ &\leq \frac{1}{T} \left( \frac{B^2}{2\eta} + 2\eta^2 \sum_{t=1}^T \rho^2 \right) \\ &\leq \frac{1}{T} \left( \frac{2B^2}{\rho^2 T} + \frac{2B^2}{\rho^2 T} \sum_{t=1}^T \rho^2 \right) \\ &\leq \frac{1}{T} \left( \frac{2B^2}{\rho^2 T} + \frac{2B^2}{\rho^2 T} \cdot T \right) \\ &\propto \frac{1}{\sqrt{T}} \end{aligned}$$

## Part4:

④ Given the objective function :

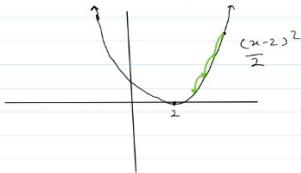
$$f(w) = \frac{1}{2} (w-2)^2 + \frac{1}{2} (w+1)^2$$

$$\text{Term 1} = \frac{1}{2} (w-2)^2$$

$$\text{Term 2} = \frac{1}{2} (w+1)^2$$

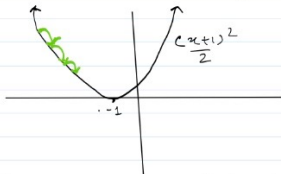
$$\text{Prob. of choosing the term} = \frac{1}{2}$$

Let's say we choose  $\frac{1}{2} (w-2)^2$



Since it is a convex function, and it is given that  $\eta$  is small enough, that every update results in improvement, hence it is bound to converge at some point with small step-size.

Similarly if we choose  $\frac{1}{2} (w+1)^2$



with small step-size this is also bound to converge.

## 2: Automatic Differentiation

### a) Computation Graphs

①  $f_1(w_1, w_2) = e^{e^{w_1 + 2w_2}} + \sin(e^{w_1 + 2w_2})$

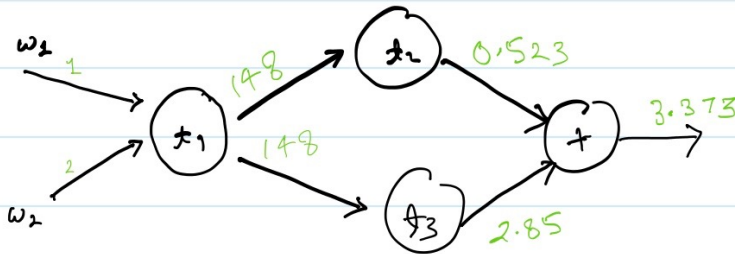
$e^{w_1 + 2w_2} = t_1 \Rightarrow t_1 = e^{1+4} = e^5 = 148.4$

$f_1 = e^{t_1} + \sin(t_1)$

$\sin(t_1) = t_2 \Rightarrow t_2 = \sin(e^5) = 0.523$

$e^{t_1} = t_3 \Rightarrow t_3 = e^{e^5} = 2.85$

$f_1 = t_2 + t_3 \Rightarrow e^{e^5} + \sin(e^5) = 3.373$



$f_2 = w_1 w_2 + \sigma(w_1)$

$w_1 w_2 = t_1$

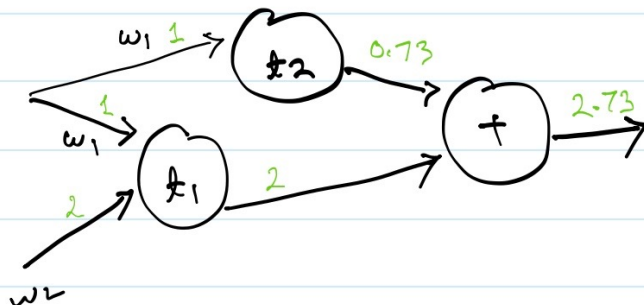
$\sigma(w_1) = t_2$

$f_2 = t_1 + t_2$

$t_1 = 2$

$t_2 = 0.73$

$f_2 = 2.73$



## b) Numerical Differentiation

⑥  $f_1(w_1, w_2) = e^{w_1} \cdot e^{2w_2} + \sin(e^{w_1} \cdot e^{2w_2}) \quad | \quad f_2(w_1, w_2) = w_1 w_2 + \sigma(w_1) \quad | \quad (1, 2)$

$$\frac{\partial f_1}{\partial w_1} = \lim_{\Delta w \rightarrow 0} \frac{f(w_1 + \Delta w, w_2) - f(w_1, w_2)}{\Delta w}$$

$$= \lim_{\Delta w \rightarrow 0} \frac{[e^{w_1 + \Delta w} \cdot e^{2w_2} + \sin(e^{w_1 + \Delta w} \cdot e^{2w_2})] - [e^{w_1} \cdot e^{2w_2} + \sin(e^{w_1} \cdot e^{2w_2})]}{\Delta w}$$

$$= \lim_{\Delta w \rightarrow 0} \frac{[e^{e^{1.01} + e^4} + \sin(e^{1.01} + e^4)] - [e^{e^1 + e^4} + \sin(e^1 + e^4)]}{0.01}$$

$$= 2.1606 \times 10^{25}$$

$$\frac{\partial f_1}{\partial w_2} = \lim_{\Delta w \rightarrow 0} \frac{f(w_1, w_2 + \Delta w) - f(w_1, w_2)}{\Delta w}$$

$$= \lim_{\Delta w \rightarrow 0} \frac{[e^{e^{1.01} + e^{2w_2 + 2\Delta w}} + \sin(e^{e^{1.01} + e^{2w_2 + 2\Delta w}})] - [e^{e^{1.01} + e^{2w_2}} + \sin(e^{e^{1.01} + e^{2w_2}})]}{\Delta w}$$

$$= \lim_{\Delta w \rightarrow 0} \frac{(e^{e^1 + e^{4.02}} + \sin(e^1 + e^{4.02})) - (e^{e^1 + e^4} + \sin(e^1 + e^4))}{0.01}$$

$$= 1.57 \times 10^{27}$$

$$\frac{\partial f_2}{\partial w_1} = \lim_{\Delta w \rightarrow 0} \frac{(w_1 + \Delta w) w_2 + \sigma(w_1 + \Delta w) - (w_1 w_2 + \sigma(w_1))}{\Delta w}$$

$$= \lim_{\Delta w \rightarrow 0} \frac{(1.01)(2) + \sigma(1.01) - (1 \cdot 2 + \sigma(1))}{0.01}$$

$$= \lim_{\Delta w \rightarrow 0} \frac{(0.733 + 2.02 - 2 - 0.731)}{0.01}$$

$$= \frac{(0.002 + 0.02)}{0.01}$$

$$= \frac{0.022}{0.01}$$

$$= 2.2$$

$$\frac{\partial f}{\partial w_2} = \lim_{\Delta w \rightarrow 0} \frac{(w_1 (w_2 + \Delta w) + \sigma(w_1) - w_1 w_2 - \sigma(w_1))}{\Delta w}$$

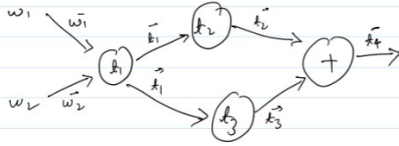
$$= \lim_{\Delta w \rightarrow 0} \frac{(1 \times (2 + 0.01) + \sigma(1) - 1 \times 2 - \sigma(1))}{0.01}$$

$$= \frac{0.01}{0.01}$$

$$= 1$$

### c) Forward Automatic Differentiation

$$\begin{aligned} f_1 &= e^{w_1 + e^{2w_2}} + \sin(e^{w_1} + e^{2w_2}) \\ t_1 &= e^{w_1} + e^{2w_2} \\ f_1 &= \underbrace{e^{t_1}}_{t_2} + \underbrace{\sin t_1}_{t_3} \\ f_1 &= t_2 + t_3 = t_4 \end{aligned}$$



$$\frac{\partial f_1}{\partial \vec{w}} = \frac{\partial (f_1)}{\partial (w_1, w_2)} = \underbrace{\begin{bmatrix} \frac{\partial f_1}{\partial w_1} & \frac{\partial f_1}{\partial w_2} \end{bmatrix}}_{\text{jacobian}}$$

$$\bar{w}_1 = \frac{\partial w_1}{\partial w_1} = 1$$

$$\bar{w}_2 = \frac{\partial w_2}{\partial w_2} = 1$$

$$\begin{aligned} \vec{t}_1 &= \frac{\partial (e^{w_1} + e^{2w_2})}{\partial w_1} = e^{w_1} = 2.713 \\ \frac{\partial (e^{w_1} + e^{2w_2})}{\partial w_2} &= 2e^{2w_2} = 109.19 \end{aligned}$$

$$\bar{t}_2 = \frac{\partial (e^{t_1})}{\partial w_1} = \frac{\partial (e^{t_1})}{\partial t_1} \times \frac{\partial t_1}{\partial w_1} = e^{t_1} \times e^{w_1} = \bar{t}_1 \times e^{w_1} = e^{w_1} \times e^{w_1} = 7.38$$

$$\frac{\partial (e^{t_1})}{\partial w_2} = \frac{\partial (e^{t_1})}{\partial t_1} \times \frac{\partial t_1}{\partial w_2} = e^{t_1} \times 2e^{2w_2} = 2 \bar{t}_1 \times e^{2w_2} = 2 \times 2e^{2w_2} \times e^{2w_2} = 11923.63$$

$$\bar{t}_3 = \frac{\partial (\sin t_1)}{\partial w_1} = \frac{\partial \sin t_1}{\partial t_1} \times \frac{\partial t_1}{\partial w_1} = \cos t_1 \times \frac{\partial (e^{w_1} + e^{2w_2})}{\partial w_1} = \cos t_1 \times e^{w_1} = \cos e^{w_1} \times e^{w_1} = 2.7152$$

$$= \frac{\partial (\sin t_1)}{\partial w_2} = \frac{\partial \sin t_1}{\partial t_1} \times \frac{\partial t_1}{\partial w_2} = \cos t_1 \times \frac{\partial (e^{w_1} + e^{2w_2})}{\partial w_2} = 2 \cos t_1 \times e^{2w_2} = 2 \cos 2e^{2w_2} \times e^{2w_2} = -35.904$$

$$\vec{t}_4 = \frac{\partial (t_2 + t_3)}{\partial w_1} = \frac{\partial t_2}{\partial w_1} + \frac{\partial t_3}{\partial w_1} = \bar{t}_2 + \bar{t}_3 = e^{w_1} \times e^{w_1} + \cos e^{w_1} \times e^{w_1} = 10.10$$

$$= \frac{\partial (t_2 + t_3)}{\partial w_2} = \frac{\partial t_2}{\partial w_2} + \frac{\partial t_3}{\partial w_2} = \bar{t}_2 \times \bar{t}_3 = 2 \times 2e^{2w_2} \times e^{2w_2} + 2 \cos 2e^{2w_2} \times e^{2w_2} = 11887.9$$

$$f_2 = w_1 w_2 + \sigma(w_1)$$

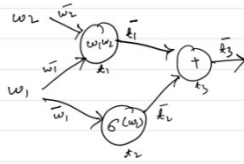
$$x_1 = w_1 w_2$$

$$x_2 = \sigma(w_1)$$

$$x_3 = x_1 + x_2 = f_2$$

$$\bar{w}_1 = \frac{\partial w_1}{\partial w_1} = 1, \frac{\partial w_1}{\partial w_2} = 0$$

$$\bar{w}_2 = \frac{\partial w_2}{\partial w_2} = 1, \frac{\partial w_2}{\partial w_1} = 0$$



$$\bar{x}_1 = \frac{\partial w_1 w_2}{\partial w_1} = w_2$$

$$\frac{\partial w_1 w_2}{\partial w_2} = w_1$$

$$\bar{x}_2 = \frac{\partial \sigma(w_1)}{\partial w_1} = \sigma(w_1) (1 - \sigma(w_1))$$

$$= \frac{\partial \sigma(w_1)}{\partial w_2} = 0$$

$$\bar{x}_3 = \frac{\partial (x_1 + x_2)}{\partial w_1} = \frac{\partial x_1}{\partial w_1} + \frac{\partial x_2}{\partial w_1} = \bar{x}_1 + \bar{x}_2 = w_2 + \sigma(w_1) (1 - \sigma(w_1))$$

$$= \frac{\partial (x_1 + x_2)}{\partial w_2} = \frac{\partial x_1}{\partial w_2} + \frac{\partial x_2}{\partial w_2} = \bar{x}_1 + \bar{x}_2 = w_1 = 1$$



## d) Backward Automatic Differentiation

③-

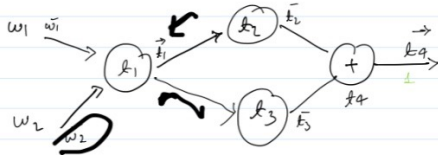
$$f_1 = e^{w_1 + e^{2w_2}} \times \sin(e^{w_1} + e^{2w_2})$$

$$k_1 = e^{w_1 + 2w_2}$$

$$k_2 = e^{k_1}$$

$$k_3 = \sin(k_1)$$

$$k_4 = k_2 + k_3$$



In reverse mode, we go from the right

$$\bar{k}_4 = \frac{\partial f}{\partial k_4} = 1$$

$$\bar{k}_2 = \frac{\partial f}{\partial k_2} = \frac{\partial f}{\partial k_4} \times \frac{\partial k_4}{\partial k_2} = 1 \times \frac{\partial}{\partial k_2} (k_2 + k_3) = 1$$

$$\bar{k}_3 = \frac{\partial f}{\partial k_3} = \frac{\partial f}{\partial k_4} \times \frac{\partial k_4}{\partial k_3} = 1 \times \frac{\partial}{\partial k_3} (k_2 + k_3) = 1$$

$$\bar{k}_1 = \frac{\partial f}{\partial k_1} = \frac{\partial f}{\partial k_3} \times \frac{\partial k_3}{\partial k_1} = 1 \times \frac{\partial}{\partial k_1} (\sin(k_1)) = \cos k_1$$

$$\bar{k}_1 = \frac{\partial f}{\partial k_1} = \frac{\partial f}{\partial k_2} \times \frac{\partial k_2}{\partial k_1} = 1 \times \frac{\partial}{\partial k_1} e^{k_1} = e^{k_1}$$

$$\bar{w}_1 = \frac{\partial f}{\partial w_1} = \frac{\partial f}{\partial k_1} \times \frac{\partial k_1}{\partial w_1} = \cos k_1 \times \frac{\partial}{\partial w_1} (e^{w_1 + 2w_2}) = e^{w_1} \cos k_1 = e^{w_1} \cos e^{w_1 + 2w_2}$$

$$= \frac{\partial f}{\partial k_1} \times \frac{\partial k_1}{\partial w_1} = e^{k_1} \times \frac{\partial}{\partial w_1} (e^{w_1 + 2w_2}) = e^{w_1} \times e^{k_1} = e^{w_1} \times e^{e^{w_1 + 2w_2}}$$

$$\bar{w}_2 = \frac{\partial f}{\partial w_2} = \frac{\partial f}{\partial k_1} \times \frac{\partial k_1}{\partial w_2} = \cos k_1 \times \frac{\partial}{\partial w_2} (e^{w_1 + 2w_2}) = \cos k_1 \times 2e^{2w_2} = \cos(e^{w_1 + 2w_2}) \times 2e^{2w_2}$$

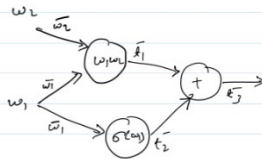
$$= e^{k_1} \times \frac{\partial}{\partial w_2} (e^{w_1 + 2w_2}) = e^{k_1} \times 2e^{2w_2} = e^{e^{w_1 + 2w_2}} \times 2e^{2w_2}$$

$$f_2 = w_1 w_2 + \sigma(w_1)$$

$$k_1 = w_1 w_2$$

$$k_2 = \sigma(w_1)$$

$$k_3 = k_1 + k_2$$



$$\bar{k}_3 = \frac{\partial f}{\partial k_3} = \frac{\partial k_3}{\partial k_3} = 1$$

$$\bar{k}_1 = \frac{\partial f}{\partial k_1} = \frac{\partial f}{\partial k_3} \times \frac{\partial k_3}{\partial k_1} = 1 \times \frac{\partial (k_1 + k_2)}{\partial k_1} = 1$$

$$\bar{k}_2 = \frac{\partial f}{\partial k_2} = \frac{\partial f}{\partial k_3} \times \frac{\partial k_3}{\partial k_2} = 1 \times \frac{\partial (k_1 + k_2)}{\partial k_2} = 1$$

$$\bar{w}_1 = \frac{\partial f}{\partial w_1} = \frac{\partial f}{\partial k_2} \times \frac{\partial k_2}{\partial w_1} = 1 \times \frac{\partial \sigma(w_1)}{\partial w_1} = \sigma(w_1) (1 - \sigma(w_1))$$

$$= \frac{\partial f}{\partial k_1} \times \frac{\partial k_1}{\partial w_1} = 1 \times \frac{\partial (w_1 w_2)}{\partial w_1} = w_2 = 2$$

$$\bar{w}_2 = \frac{\partial f}{\partial w_2} = \frac{\partial f}{\partial k_1} \times \frac{\partial k_1}{\partial w_2} = 1 \times \frac{\partial (w_1 w_2)}{\partial w_2} = w_1 = 1$$

Ⓢ Yes, it is pretty easy with software libraries.

e): Yes it is pretty easy with software libraries

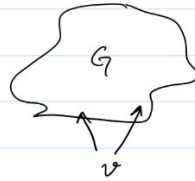
### 3 Directed Acyclic Graphs:

④: If a graph  $G$  is a DAG, then  $G$  has a topological ordering

Proof: (By induction)

Base Case: If we have only 1 node,

then  $G$  has a topological ordering.



Hypothesis: If  $G'$  has  $\leq n$  nodes, then  $G'$  has topological ordering.

Consider  $G$  with  $n+1$  nodes,

we know  $G$  is a DAG, & we added one extra node  $v$ , we add  $v$  with no incoming edges.

$G - \{v\}$  is a DAG, since deleting  $v$  cannot create any cycles.

So, in short:

\* Place  $v$  first; then add topological ordering of  $G - \{v\}$ .

\*  $v$  has no incoming edges, hence first point is valid.

By induction the lemma is proven.

⑤: If  $G$  has a topological ordering, then  $G$  is a DAG.

Proof: By contradiction

Let's suppose  $G$  has topological ordering  $x_1, x_2, x_3, \dots, x_n$

Suppos  $G$  has a directed cycle like below



from the above figure  $x_2$  has lower index than  $x_3$  since  $2 < 3$

Does Node  $x_2$  should come before  $x_3$ , if they are topological order.

But the edge  $x_3 \rightarrow x_2$ , and  $x_1, x_2, \dots, x_n$  is a topological order, we must have  $3 < 2$ , which is a contradiction.

So  $G$  has no cycle.  $G$  is a DAG.