

part 1

1.1.1)

$$\text{write } f(x) = \dots \underbrace{\sigma(W_3 \sigma(W_2 \sigma(W_1 x + b_1) + b_2) + b_3)}_{n \text{ layers}} \dots$$

$$W_i = I \Rightarrow f(x) = \dots \underbrace{\sigma(\sigma(\sigma(x + b_1) + b_2) + b_3)}_{n \text{ layers}} \dots$$

$$\text{let } y_1 = x + b_1$$

$$y_2 = \sigma(y_1) + b_2$$

$$y_3 = \sigma(y_2) + b_3$$

:

$$f(x) = y_n = \sigma(y_{n-1}) + b_n$$

$$\text{Then } \frac{df(x)}{dx} = \sigma'(y_{n-1}) y'_{n-1}$$

$$= \sigma'(y_{n-1}) \sigma'(y_{n-2}) y'_{n-2}$$

:

$$= \sigma'(y_{n-1}) \sigma'(y_{n-2}) \dots \sigma'(y_1) y'_1$$

$$= \prod_{i=1}^{n-1} \sigma'(y_i)$$

$$= \prod_{i=1}^{n-1} (\sigma(y_i) - \sigma^2(y_i)) \quad \text{since } \sigma'(x) = \sigma(x)(1-\sigma(x)) = \sigma(x) - \sigma^2(x)$$

$$\text{since for any } x \in \mathbb{R}, \quad \frac{d(\sigma(x) - \sigma^2(x))}{d\sigma(x)} = 1 - 2\sigma(x) \stackrel{\text{set}}{=} 0$$

$$\Rightarrow \sigma(x) = \frac{1}{2} \text{ maximizes } \sigma'(x) \text{ at } \sigma'(x) = \frac{1}{4}$$

$$\text{and since } \sigma(x) \in (0, 1) \Rightarrow 0 \leq \sigma(x) - \sigma^2(x) \leq \frac{1}{4}$$

$$\text{Then } 0 \leq \sigma(y_i) - \sigma^2(y_i) \leq \frac{1}{4}$$

$$\Rightarrow 0 \leq \prod_{i=1}^{n-1} (\sigma(y_i) - \sigma^2(y_i)) \leq \left(\frac{1}{4}\right)^{n-1}$$

Note that if there are n hidden layers plus an output layer, i.e. $f(x) = y_{n+1} = \sigma(y_n) + b_{n+1}$, then we can replace all $n-1$ with n to get $0 \leq \prod_{i=1}^n (\sigma(y_i) - \sigma^2(y_i)) \leq \left(\frac{1}{4}\right)^n$ as wanted.

In either case,
the gradients will vanish as n gets large.

1.2.1)

$$x_n = \dots \tanh(W \tanh(W \tanh(x_1))) \dots$$

$$x_2 = \tanh(Wx_1)$$

$$x_3 = \tanh(Wx_2)$$

⋮

$$x_n = \tanh(Wx_{n-1})$$

$$\frac{\partial x_n}{\partial x_1} = \prod_{t=1}^{n-1} \frac{\partial x_{t+1}}{\partial x_t}$$

$$= \prod_{t=1}^{n-1} \left[\frac{\partial \tanh(Wx_t)}{\partial Wx_t} \right] W$$

Note that the Jacobian of \tanh is

$$\begin{bmatrix} \tanh'(x_1) & & & \\ & \tanh'(x_2) & & \\ & & \ddots & \\ & & & \tanh'(x_k) \end{bmatrix}$$

and $\tanh'(x) = \operatorname{sech}^2(x) \in (0, 1]$, for $x \in \mathbb{R}$. So $0 < \sigma_{\max} \left[\frac{\partial \tanh(Wx_t)}{\partial Wx_t} \right] \leq 1$

$$\Rightarrow \sigma_{\max} \left(\frac{\partial x_n}{\partial x_1} \right) = \sigma_{\max} \left[\prod_{t=1}^{n-1} \frac{\partial x_{t+1}}{\partial x_t} \right] \leq \prod_{t=1}^{n-1} \sigma_{\max} \left[\frac{\partial x_{t+1}}{\partial x_t} \right] \text{ by hint.}$$

$$\begin{aligned} \text{and } \sigma_{\max} \left[\frac{\partial x_{t+1}}{\partial x_t} \right] &= \sigma_{\max} \left[\frac{\partial \tanh(Wx_t)}{\partial Wx_t} \right] W \\ &\leq \sigma_{\max} \left(\frac{\partial \tanh(Wx_t)}{\partial Wx_t} \right) \sigma_{\max}(W) \\ &\leq 1 \cdot \frac{1}{2} = \frac{1}{2} \end{aligned}$$

$$\Rightarrow \sigma_{\max} \left(\frac{\partial x_n}{\partial x_1} \right) = \sigma_{\max} \left[\prod_{t=1}^{n-1} \frac{\partial x_{t+1}}{\partial x_t} \right] \leq \left(\frac{1}{2}\right)^{n-1}$$

Also, since singular values are non-negative,

$$\Rightarrow 0 \leq \sigma_{\max} \left[\prod_{t=1}^{n-1} \frac{\partial x_{t+1}}{\partial x_t} \right] \leq \left(\frac{1}{2}\right)^{n-1}$$

$$\Rightarrow 0 \leq \sigma_{\max} \left(\frac{\partial x_n}{\partial x_1} \right) \leq \left(\frac{1}{2}\right)^{n-1} \text{ as wanted.}$$

1.3.1)

$$\text{sub in kernel function: } \alpha_i = \frac{\sum_{j=1}^n \phi(Q_i)^T \phi(k_j) V_j}{\sum_{j=1}^n \phi(Q_i)^T \phi(k_j)}$$
$$= \frac{\phi(Q_i)^T \sum_{j=1}^n \phi(k_j) V_j^T}{\phi(Q_i)^T \sum_{j=1}^n \phi(k_j)}$$

we can compute $\sum_{j=1}^n \phi(k_j) V_j^T$ and $\sum_{j=1}^n \phi(k_j)$ once and reuse them for every query. So attention can be calculated in $O(n)$.

(3.2)

Assume we have SVD of $P = U\Sigma W^T$, $U, W \in \mathbb{R}^{n \times k}$ are semi-unitary s.t. $U^T U = W^T W = I$.

$$\text{Then } PV = U\Sigma W^T V$$

step ① computing $U\Sigma$ is multiplying the k n -dimensional vectors by the k non-zero singular values, which is $O(nk)$, and results in a $n \times k$ matrix

step ② computing $W^T V$ is multiplying $k \times n$ matrix by $n \times d$ matrix, which is $O(knd)$. and results in a $k \times d$ matrix

step ③ computing $U\Sigma W^T V$ is multiplying the $n \times k$ matrix by the $k \times d$ matrix which is $O(nkd)$

\Rightarrow so the entire 3 steps combined is $O(nkd)$ as wanted.

(1.33)

P is rank $k \Rightarrow$ we can write $P = U\Sigma W^T$, $U, W \in \mathbb{R}^{n \times k}$ by SVD.
where U, W are semi-unitary.

Then we can take $C = D = W^T \in \mathbb{R}^{k \times n}$.

$$\begin{aligned}\Rightarrow Q(Ck)^T DV &= Q(W^T k)^T W^T V \\ &= Q k^T W W^T V \quad \downarrow \text{since } WW^T = I, W \text{ is orthogonal.} \\ &= Q k^T V \\ &= P V \quad \text{since } P = Q k^T\end{aligned}$$

complexity:
Step ①: compute Ck ($k \times n$ times $n \times d$), $O(knd)$. output $k \times d$ matrix.
Step ②: compute DV ($k \times n$ times $n \times d$), $O(knd)$. output $k \times d$ matrix.
Step ③: $Q(Ck)^T$ ($n \times d$ times $d \times k$), $O(ndk)$. output $n \times k$ matrix.
Step ④: $Q(Ck)^T DV$ ($n \times k$ times $k \times d$), $O(nkd)$. output $n \times d$ matrix.

\Rightarrow so entire steps ① \rightarrow ④ takes $O(nkd)$.

Part 2

$$2.1) \quad g[\theta, \tilde{\alpha}] = f(\tilde{\alpha}) \frac{\partial}{\partial \theta} \log p(a=\tilde{\alpha} | \theta), \quad \tilde{\alpha} \sim p(a | \theta)$$

notice $\tilde{\alpha}$ is either 1 or 0

$$\text{if } \tilde{\alpha} = 0: \quad g[\theta, \tilde{\alpha}] = f(0) \frac{\partial}{\partial \theta} \log p(a=0 | \theta) \\ = 0 \cdot \frac{\partial}{\partial \theta} \log p(a=0 | \theta) = 0.$$

$$\begin{aligned} \text{if } \tilde{\alpha} = 1: \quad g[\theta, \tilde{\alpha}] &= f(1) \frac{\partial}{\partial \theta} \log p(a=1 | \theta) \\ &= 1 \cdot \frac{\partial}{\partial \theta} \log p(a=1 | \theta) \\ &= \frac{\partial}{\partial \theta} \log u \\ &= \frac{1}{u} \frac{\partial u}{\partial \theta} \\ &= \frac{1}{u} x \sigma' \left(\sum_{d=1}^D \theta_d x_d \right) \\ &= \frac{1}{u} x u (1-u) \quad \text{since } \sigma'(x) = \sigma(x)(1-\sigma(x)) \\ &= x(1-u) \end{aligned}$$

Combining them \Rightarrow
$$g[\theta, \tilde{\alpha}] = \tilde{\alpha} x (1-u)$$

$$\begin{aligned} 2.2) \quad \text{Var}[g[\theta, \tilde{\alpha}],] &= \text{Var}[\tilde{\alpha} x, (1-u)] \\ &= x_1^2 (1-u)^2 \text{Var}[\tilde{\alpha}] \\ &= x_1^2 (1-u)^2 u (1-u) \\ &= x_1^2 u (1-u)^3 \end{aligned}$$

Part 3

3.1.1

$$\begin{aligned} \mathbf{x}^T L \mathbf{x} &= \mathbf{x}^T (\mathbf{D} - \mathbf{A}) \mathbf{x} \\ &= \mathbf{x}^T \mathbf{D} \mathbf{x} - \mathbf{x}^T \mathbf{A} \mathbf{x} \end{aligned}$$

$$\begin{aligned} &= [\mathbf{x}_1 \dots \mathbf{x}_n] \begin{bmatrix} \deg(1) & & & \\ & 0 & & \\ & & \ddots & \\ & & & \deg(n) \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{bmatrix} - [\mathbf{x}_1 \dots \mathbf{x}_n] \begin{bmatrix} A_{11} & A_{12} & \cdots & A_{1n} \\ A_{21} & A_{22} & \cdots & A_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ A_{n1} & A_{n2} & \cdots & A_{nn} \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{bmatrix} \\ &= \sum_{u=1}^n x_u \deg(u) x_u - [\mathbf{x}_1 \dots \mathbf{x}_n] \begin{bmatrix} A_{11} & A_{12} & \cdots & A_{1n} \\ A_{21} & A_{22} & \cdots & A_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ A_{n1} & A_{n2} & \cdots & A_{nn} \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{bmatrix} \end{aligned}$$

since $A_{uv} = 1$ if there is an edge between vertex u, v . and $A_{uv} = 0$ otherwise,

$$\text{so } [\mathbf{x}_1 \dots \mathbf{x}_n] \begin{bmatrix} A_{11} & A_{12} & \cdots & A_{1n} \\ A_{21} & A_{22} & \cdots & A_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ A_{n1} & A_{n2} & \cdots & A_{nn} \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{bmatrix} = \sum_{(u,v) \in E} 2 x_u x_v \quad \text{since } (u,v) \in E \\ \Rightarrow A_{uv} \text{ and } A_{vu} \text{ are both 1}$$

\Rightarrow Therefore

$$\begin{aligned} \mathbf{x}^T L \mathbf{x} &= \sum_{u=1}^n x_u^2 \deg(u) - \sum_{(u,v) \in E} 2 x_u x_v \\ &= \sum_{u=1}^n \sum_{(u,v) \in E} x_u^2 - \sum_{(u,v) \in E} 2 x_u x_v \\ &= \sum_{(u,v) \in E} (x_u^2 + x_v^2) - \sum_{(u,v) \in E} 2 x_u x_v \\ &= \sum_{(u,v) \in E} (x_u^2 - 2 x_u x_v + x_v^2) \\ &= \sum_{(u,v) \in E} (x(u) - x(v))^2 \quad \text{as wanted,} \end{aligned}$$

3.1.2)

Since L is a symmetric matrix (because both A and D are symmetric) and $x^T L x = \sum_{(u,v) \in E} (x(u) - x(v))^2 \geq 0$, L must be PSD.

by 3.1.1.

next consider the vector $v = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$

$$\begin{aligned}
 Lv &= (D - A)v = Dv - Av \\
 &= \begin{bmatrix} \deg(1) & 0 & \dots & 0 \\ \vdots & & & \\ 0 & \deg(n) & & \end{bmatrix} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} - \begin{bmatrix} A_{11} & A_{12} & \dots & A_{1n} \\ A_{21} & A_{22} & \dots & A_{2n} \\ \vdots & \ddots & \ddots & \vdots \\ A_{n1} & A_{n2} & \dots & A_{nn} \end{bmatrix} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \\
 &= \begin{bmatrix} \deg(1) \\ \deg(2) \\ \vdots \\ \deg(n) \end{bmatrix} - \begin{bmatrix} \sum_{i=1}^n A_{1i} \\ \sum_{i=1}^n A_{2i} \\ \vdots \\ \sum_{i=1}^n A_{ni} \end{bmatrix} \\
 &= 0 \quad \left(\text{since } \sum_{i=1}^n A_{ji} = \deg(j) \right)
 \end{aligned}$$

so the smallest eigenvalue $\lambda_{\min}(L) = 0$ with eigenvector $v = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$

3.2.1)

$$\lambda_{\max}(\tilde{L}) = \max_{\|x\|=1} x^T \tilde{L} x, \quad x \in \mathbb{R}^n.$$

so it is sufficient to show that $x^T \tilde{L} x \leq 2$. for any $x \in \mathbb{R}^n$ st. $\|x\|=1 \Leftrightarrow x^T x = 1$

First, $x^T (I + \tilde{A}) x$

$$\begin{aligned} &= x^T x + x^T \tilde{A} x \\ &= \sum_{i=1}^n x_{(i)}^2 + \sum_{(i,j) \in E} \frac{2x_{(i)}x_{(j)}}{\sqrt{d(i)d(j)}} \\ &= \sum_{(i,j) \in E} \left(\frac{x_{(i)}}{\sqrt{d(i)}} + \frac{x_{(j)}}{\sqrt{d(j)}} \right)^2 \end{aligned}$$

same reasoning as 3.1.1.

$$\geq 0.$$

$$\Rightarrow x^T (I + \tilde{A}) x \geq 0$$

$$\Rightarrow x^T I x + x^T \tilde{A} x \geq 0$$

$$\Rightarrow -x^T \tilde{A} x \leq x^T x$$

$$\Rightarrow x^T x - x^T \tilde{A} x \leq 2x^T x$$

also, $x^T \tilde{L} x = x^T (I - \tilde{A}) x$

$$= x^T x - x^T \tilde{A} x$$

$$\Rightarrow x^T \tilde{L} x \leq 2x^T x$$

$$\Rightarrow x^T \tilde{L} x \leq 2 \quad \text{since } \|x\| = x^T x = 1$$

as wanted.

3.2.2).

By 3.2.1, $\lambda_{\max}(\tilde{L}) \leq 2$.

$$x^T \tilde{L} x \leq 2$$

$$\Rightarrow x^T(I - \tilde{A})x \leq 2$$

$$\Rightarrow x^T x - x^T \tilde{A} x \leq 2$$

$$\Rightarrow 1 - x^T \tilde{A} x \leq 2 \quad \text{if we take } \|x\| = x^T x = 1$$

$$\Rightarrow 1 - 2 \leq x^T \tilde{A} x$$

$$\Rightarrow -1 \leq x^T \tilde{A} x$$

$$\text{so } -1 \leq \lambda_i(\tilde{A}). \quad \textcircled{1}$$

$$\text{Also, } x^T \tilde{L} x = x^T(I - \tilde{A})x$$

$$= x^T x - x^T \tilde{A} x \\ = \sum_{i=1}^n x(i)^2 - \sum_{(i,j) \in E} \frac{2x(i)x(j)}{\sqrt{d(i)d(j)}}$$

$$= \sum_{(i,j) \in E} \left(\frac{x(i)}{\sqrt{d(i)}} - \frac{x(j)}{\sqrt{d(j)}} \right)^2 \quad \text{same reasoning as 3.1.1.}$$

$$\geq 0.$$

$$\Rightarrow x^T(I - \tilde{A})x \geq 0$$

$$\Rightarrow x^T x - x^T \tilde{A} x \geq 0$$

$$\Rightarrow x^T x \geq x^T \tilde{A} x$$

$$\Rightarrow 1 \geq x^T \tilde{A} x \quad \text{by hint 3.2.1.}$$

$$\text{so } \lambda_i(\tilde{A}) \leq 1 \quad \textcircled{2}$$

$$\text{combining } \textcircled{1} + \textcircled{2} \Rightarrow -1 \leq \lambda_i(\tilde{A}) \leq 1$$

As wanted.

3.3.1)

we can take the eigendecomposition of $\tilde{A} = Q \Lambda Q^T$ where Q is an orthonormal matrix.
since \tilde{A} is a square symmetric matrix.

$$\begin{aligned}\text{Then } \tilde{A}^k &= (Q \Lambda Q^T)^k \\ &= Q \underbrace{\Lambda Q^T Q \Lambda Q^T \dots Q \Lambda Q^T}_{k \text{ times}} \\ &= Q \Lambda^k I \Lambda^k \dots I \Lambda^k Q^T \quad \text{since } Q^T Q = I \\ &= Q \Lambda^k Q^T\end{aligned}$$

since Λ is diagonal, Λ^k is equal to element wise power across the diagonal entries, which is much more efficient.

3.3.2)

No. It is not always beneficial to use a deeper model since the eigenvalues of \tilde{A} are all between -1 and 1, so taking large number of layers k would result in Λ^k approaching 0 and thus \tilde{A}^k approaching 0 (vanishing gradient) as well as the problem of over smoothing, and hence $H^{(k)}$.

3.4.1)

$$\alpha_{ij} = \frac{\exp(\text{LeakyReLU}(\vec{a}^\top [\vec{w_h}_i || \vec{w_h}_j]))}{\sum_{k \in N_i} \exp(\text{LeakyReLU}(\vec{a}^\top [\vec{w_h}_i || \vec{w_h}_k]))}$$

→ GAT allows for (implicitly) assigning different importance to nodes of a same neighbourhood, enabling greater model capacity.

→ In GAT, analyzing learned attention weights may result in better model interpretability.