

الف) $\text{null}(A) = \text{null}(A^T A)$

فرض 1: $x \in \text{null}(A) \Rightarrow A^T (Ax = 0)$

$$\Rightarrow A^T A x = A^T 0 = 0 \Rightarrow x \in \text{null}(A^T A)$$

$$\Rightarrow \boxed{\text{null}(A) \subseteq \text{null}(A^T A)} \quad ①$$

فرض 2: $x \in \text{null}(A^T A) \Rightarrow A^T (A^T A x = 0)$

$$\Rightarrow x^T A^T A x = x^T 0 = 0 \Rightarrow (Ax)^T Ax = 0$$

$$\Rightarrow \|Ax\|_2 = 0 \Rightarrow Ax = 0 \Rightarrow x \in \text{null}(A)$$

$$\Rightarrow \boxed{\text{null}(A^T A) \subseteq \text{null}(A)} \quad ②$$

$$①, ② \Rightarrow \text{null}(A) = \text{null}(A^T A)$$

ب) $\|A\|_2 = \sup_{\|U\|=\|V\|=1} U^T A V, \quad A \in \mathbb{R}^{m \times n}$

این نرم برابر با حداکثر مقداری است که ماتریس A می تواند بردار واحد x را بکشد.

$$\|A\|_2 = \max_{\|x\|_2=1} \|Ax\|_2$$

می دانیم:

$$\sup_{\|U\|=\|V\|=1} U^T A V$$

این عبارت را می توانیم بصورت ضرب داخلی دو بردار U و AV بنویسیم. آنرا:

آنکه اگر بردار v در جهت Av باشد ماکزیم ترین حالت پس می آید چون ضرب داخلی ماکزیم می شود.
 همچنین دو بردار v و u نرم واحد دارند.

$$u = \frac{Av}{\|Av\|}$$

پس v را مانند روبه روی نظر می گیریم:

$$u^T Av = \frac{(Av)^T Av}{\|Av\|_2} = \|Av\|_2$$

$$\sup_{\|v\|_2=1} \|Av\|_2 = \text{Spectral norm of } A$$

$$\max_{\|x\|_2=1} \|Ax\|_2 = \|A\|_2$$

ج) $\|A\|_2 = \sigma_{\max}(A)$

$$\|A\|_2 = \max_{\|x\|_2=1} \|Ax\|_2$$

$$A = U \Sigma V^T$$

یکت تجزیه SVD:

$$\Rightarrow \|Ax\|_2 = \|U \Sigma V^T x\|_2 \xrightarrow{V \text{ is orthogonal}} \|Ax\|_2 = \|\Sigma V^T x\|_2$$

$$\Rightarrow \|Ax\|_2 = \|\Sigma y\|_2 \quad \text{where } y = V^T x$$

زمانی $\|\Sigma y\|_2$ ماکزیم می شود که y با بردار اولیه y_1 align باشد یعنی در جهتش باشد.
 پس $\|Ax\|_2 \leq \sigma_1$. پس بزرگترین مقدار $\|A\|_2$ می تواند دقیقاً σ_1 باشد که بزرگترین
 Singular value است. $\sigma_1 > \sigma_2 > \dots > \sigma_{\min}$
 $\|A\|_2 = \sigma_{\max}(A)$

$$\Rightarrow \|A^T A\|_2 = \|A A^T\|_2 = \|A\|_2^2$$

$$A = U \Sigma V^T$$

باروس SVD :

$$A^T A = V \Sigma^T U^T U \Sigma V^T = V \Sigma^2 V^T$$

- پس Singular value های $A^T A$ نسبت به A به توان 2 رسیده اند.

باتوجه به رابطه بخش "ج" که داشتیم :

$$\|A^T A\|_2 = \sigma_{\max}(A^T A) = \sigma_{\max}(A)^2 = \|A\|_2^2$$

همینطور برای $\|A A^T\|_2$ نیز مانند بالا می توان نوشت .

$$\Rightarrow \text{Tr}(A^T B) = \sum_{i=1}^m \sum_{j=1}^n A_{ij} B_{ij} \quad , \quad A, B \in \mathbb{R}^{m \times n}$$

$$A^T B = K_{n \times n} \Rightarrow \text{Tr}(K)_{n \times n} = \sum_{i=1}^n K_{ii}$$

$$K_{ii} = \sum_{j=1}^m A_{ij}^T B_{j1} = \sum_{j=1}^m A_{j1} \cdot B_{j1}$$

حل کرده element های قطری K را بخواهیم محاسبه کنیم کافی یک Sum روی همان دوم ماتریس نیز بگیریم .

$$\text{Tr}(K) = \sum_{i=1}^n \sum_{j=1}^m A_{ij} B_{ij}$$

چون جمع همه ضرب های همان های دوم ماتریس را داریم می توانیم index ها را عوض کنیم :

$$\text{Tr}(K) = \sum_{i=1}^m \sum_{j=1}^n A_{ij} B_{ij} = \text{Tr}(A^T B)$$

$$9) \|A\|_F = \sqrt{\sum \sigma_i^2(A)}$$

$$\|A\|_F = \sqrt{\text{Tr}(A^T A)} = \sqrt{\sum_{i=1}^m \sum_{j=1}^n A_{ij}^2}$$

$$\sqrt{\text{Tr}(A^T A)} \stackrel{\text{SVD}}{=} \sqrt{\text{Tr}(V \Sigma^2 V^T)} \quad \xrightarrow{\text{Tr}(V \Sigma^2 V^T) = \text{Tr}(\Sigma^2)} \\ \text{V is orthogonal}$$

$$\text{Tr}(\Sigma^2) = \sigma_1^2 + \sigma_2^2 + \dots + \sigma_{\min(n,m)}^2 = \sum_{i=1}^{\min(n,m)} \sigma_i^2$$

$$\Rightarrow \|A\|_F = \sqrt{\text{Tr}(\Sigma^2)} = \sqrt{\sum_{i=1}^{\min(n,m)} \sigma_i^2(A)}$$

$$j) \text{Tr}(A) = \sum \lambda_i$$

Eigenvalue decomposition

$$A_{n \times n} = Q \Lambda Q^{-1}$$

$$\text{Tr}(A) = \text{Tr}(Q \Lambda Q^{-1}) = \text{Tr}(\Lambda)$$

تغيير ترتيب التتبع
Similarity transformation

$$\Rightarrow \text{Tr}(A) = \text{Tr}(\Lambda) = \sum_{i=1}^n \lambda_i$$

$$2) \|A + X\|_F^2 = \|A\|_F^2 + \|X\|_F^2 + 2 \text{Tr}(A^T B)$$

$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n A_{ij}^2} \Rightarrow \|A\|_F^2 = \sum_i \sum_j A_{ij}^2$$

$$\Rightarrow \|A + X\|_F^2 = \sum_i \sum_j (A_{ij} + X_{ij})^2 \Rightarrow$$

$$\text{مثلاً) } \|A+X\|_F^2 = \sum \sum A_{ij}^2 + \sum \sum X_{ij}^2 + 2 \sum \sum A_{ij} X_{ij}$$

$$\textcircled{*} \quad \sum \sum A_{ij} X_{ij} = \text{Tr}(A^T X) \quad \text{در بخش "ح" اثبات کردیم که:}$$

دو تا نرم اول هم که تعریف نرم فرنیس هستند:

$$\Rightarrow \quad \|A+X\|_F^2 = \|A\|_F^2 + \|X\|_F^2 + 2 \text{Tr}(A^T X)$$

$$b) \quad \det(A) = \prod_{i=1}^n \lambda_i \quad A \in \mathbb{R}^{n \times n}$$

$$\text{eigenvalue decomposition: } A = Q \Lambda Q^{-1}$$

$$\det(A) = \det(Q \Lambda Q^{-1}) = \det(Q) \cdot \det(\Lambda) \cdot \det(Q^{-1})$$

$$\Rightarrow \det(A) = \det(Q) \cdot \det(\Lambda) \cdot \frac{1}{\det(Q)} = \det(\Lambda)$$

ماتریس قطری $n \times n$ که روی قطرش مقادیر ویژه $\lambda_1, \lambda_2, \dots, \lambda_n$ وجود دارد.

$$\Rightarrow \det(A) = \det(\Lambda) = \prod_{i=1}^n \lambda_i$$

$$c) \quad \|A\|_F \geq \|A\|_2 \geq \frac{1}{\sqrt{n}} \|A\|_F \geq \frac{\text{Tr}(A)}{n} \geq \sqrt[n]{\det(A)}$$

$$\|A\|_F = \sqrt{\sum_i \sigma_i^2(A)}, \quad \|A\|_2 = \sigma_{\max}(A) \quad \text{چون}$$

$$\text{پس: } \boxed{\|A\|_F \geq \|A\|_2}$$

از آنجایی که $\|A\|_2$ بزرگترین Singular value است و $\frac{1}{\sqrt{n}} \|A\|_F$ میانگین Singular value است. پس $\|A\|_2 \geq \frac{\|A\|_F}{\sqrt{n}}$

از آنجایی که نرم فربنیوس جمع Singular value ها است. و تقریباً میانگین آن ها است. و عبارت $\frac{\text{Tr}(A)}{n}$ میانگین eigen value های ماتریس A است، چون eigen value شامل اعداد مثبت و منفی می شود اما Singular value ها همیشه مثبت اند پس:

$$\frac{\|A\|_F}{\sqrt{n}} \geq \frac{\text{Tr}(A)}{n} = \frac{\sum_i \lambda_i}{n}$$

برای بخش آخر آن از eigenvalue استفاده کنیم:
Decomposition

$$A = Q \Lambda Q^{-1}$$

$$\begin{cases} \text{Tr}(A) = \text{Tr}(\Lambda) = \sum_i \lambda_i \\ \det(A) = \det(\Lambda) = \prod_i \lambda_i \end{cases}$$

از آنجایی که میانگین میانگین حسابی بزرگتر از میانگین هندسی است پس:

$$\frac{\sum_i \lambda_i}{n} \geq \sqrt[n]{\prod_i \lambda_i} \Rightarrow \frac{\text{Tr}(A)}{n} \geq \sqrt[n]{\det(A)}$$

(2) با توجه به قضیه مقدار میانگین: $\frac{|f(b) - f(a)|}{|b-a|} = f'(c)$, $c \in (a, b)$

$$f(x) = \log(1 + e^x)$$

$$f'(x) = \frac{1}{1 + e^x} \frac{d}{dx} (1 + e^x) = \frac{e^x}{1 + e^x}$$

$$\Rightarrow x \in (-\infty, +\infty) \Rightarrow \boxed{0 \leq f'(x) \leq 1}$$

So $\Rightarrow \frac{|f(b) - f(a)|}{|b-a|} = f'(c) \leq 1 \Rightarrow$

$$|f(b) - f(a)| \leq \underbrace{1}_K \times |b-a|$$

پس این تابع Lipschitz 1- است.

(3) اگرچه توابع f_i یک تابع Lipschitz p_i - باشند. آنگاه:

$$|f_i(x) - f_i(y)| \leq p_i |x - y|$$

البدلی f_2 این رابطه را بنویسیم و سپس f_1 را به عنوان ورودی به آن بدلیم: $f_2 \circ f_1$

$$|f_2(f_1(x)) - f_2(f_1(y))| \leq p_2 |f_1(x) - f_1(y)|$$

$$|f_1(x) - f_1(y)| \leq p_1 |x - y| \quad \xrightarrow{\text{So}}$$

$$|f_2(f_1(x)) - f_2(f_1(y))| \leq P_2 P_1 |x - y|$$

الذین عبارت را f_n متمم به هم اندازیم:

$$|f_n \circ f_{n-1} \circ \dots \circ f_1(x) - f_n \circ f_{n-1} \circ \dots \circ f_1(y)| \leq P_n P_{n-1} \dots P_1 |x - y|$$

پس ثابت Lipschitz این تابع برابر است با:

$$K = P_n \times P_{n-1} \times \dots \times P_1$$

(4) فرض می‌کنیم دو تابع f و g لیپشیتز هستند. آنگاه:

$$|f(x) - f(y)| \leq L_f |x - y|$$

$$|g(x) - g(y)| \leq L_g |x - y|$$

$$\Rightarrow |f(x) - f(y)| + |g(x) - g(y)| \leq (L_f + L_g) |x - y|$$

از آنجایی که $|A+B| \leq |A| + |B|$ است، پس:

$$|(f(x) + g(x)) - (f(y) + g(y))| \leq |f(x) - f(y)| + |g(x) - g(y)|$$

$$\stackrel{\text{So}}{\Rightarrow} |(f(x) + g(x)) - (f(y) + g(y))| \leq (L_f + L_g) |x - y|$$

پس ثابت برای تابع $f+g$ برابر با $L_f + L_g$ است.

$$\begin{aligned} |f(x)g(x) - f(y)g(y)| &= |f(x)g(x) - f(x)g(y) + f(x)g(y) - f(y)g(y)| \\ &= |f(x)(g(x) - g(y)) + g(y)(f(x) - f(y))| \end{aligned}$$

سوال ①

ادامه بخش 4) با توجه به نامساوی مثلثی :

$$|f(x)g(x) - f(y)g(y)| \leq |f(x)| |g(x) - g(y)| + |g(y)| |f(x) - f(y)|$$

$$|g(x) - g(y)| \leq L_g |x - y|$$

$$|f(x) - f(y)| \leq L_f |x - y|$$

اگر f و g در x وابسته باشند پس $|f(x)| \leq K_f$ ، $|g(y)| \leq K_g$ است پس بصورت کلی خواهیم داشت :

$$|f(x)g(x) - f(y)g(y)| \leq (K_f L_f + K_g L_g) |x - y|$$

← نسبت Lipschitz

سوال ②

تابع شبکه دلتا :

$$y_k(x, w) = \sigma \left(\sum_{j=1}^n w_{kj}^{(2)} h \left(\sum_{i=1}^D w_{ji}^{(1)} x_i + w_{j0} \right) + w_{k0}^{(2)} \right)$$

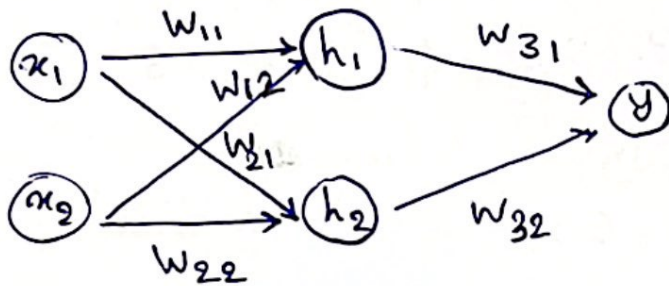
$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad , \quad \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$\sigma(x) = \frac{1}{2} + \frac{1}{2} \tanh\left(\frac{x}{2}\right)$$

اگر \tanh را به جای h قرار دهیم آنگاه :

$$y_k(x, w) = \sigma \left(\sum_{j=1}^n \frac{1}{2} w_{kj}^{(2)} h \left(\sum_{i=1}^D w_{ji}^{(1)} \times \frac{1}{2} + \frac{1}{2} w_{j0} \right) + w_{k0}^{(2)} + \sum_{j=1}^n \frac{1}{2} w_{kj}^{(2)} \right)$$

$$\begin{cases} w_{kj}^{(2)'} = \frac{1}{2} w_{kj}^{(2)} \\ w_{k0}^{(2)'} = w_{k0}^{(2)} + \sum_{j=1}^n \frac{1}{2} w_{kj}^{(2)} \end{cases} \quad \begin{cases} w_{ji}^{(1)'} = \frac{1}{2} w_{ji}^{(1)} \\ w_{j0}' = \frac{1}{2} w_{j0}^{(1)} \end{cases}$$



(1) خطی که بتواند S_2 و S_3 را جدا کند خط $x_1 = 3$ است.

برای اینکه شرط روبه رو را در نظر بگیریم، داریم:

$$\begin{cases} h_1(x) = 0 & S_3 \\ h_1(x) = 1 & S_2 \end{cases}$$

$$z_1(x) = w_{11}x_1 + w_{12}x_2 + b_1$$

$$h_1(x) = \phi(z_1(x)) \quad \text{where} \quad \phi(z) = \begin{cases} 1 & z > 0 \\ 0 & z \leq 0 \end{cases}$$

برای اینکه به ازای $z > 0$ دسته S_2 را تشخیص دهیم و به ازای $z \leq 0$ دسته S_3

باید خط $x_1 - 3 = 0$ را قریباً کنیم پس در این حالت z برابر است با:

$$z_1(x) = -x_1 + 3 \Rightarrow \begin{cases} z > 0 & \phi(z) = 1 \Rightarrow S_2 \\ z \leq 0 & \phi(z) = 0 \Rightarrow S_3 \end{cases}$$

$$\text{So: } \begin{cases} w_{11} = -1 \\ w_{12} = 0 \end{cases}, b_1 = 3$$

(2) برای اینکه دو دسته S_1 و S_2 را جدا کنیم می توانیم از خطی که از دو نقطه

(3 و 1) و (2 و 1) می گذرد را بعنوان جدا کننده در نظر بگیریم که بین این دو دسته

هرامی تولید این خط برابر است با:

$$x_2 + 2x_1 - 5 = 0$$

پس اگر z_2 را بصورت روبه‌یو تعریف کنیم:

$$\begin{aligned} & > 0 & \phi(z) = 1 & S_2 \\ & \leq 0 & \phi(z) = 0 & S_1 \end{aligned}$$

$$z_2(x) = 2x_1 + x_2 - 5$$

$$h_2(x) = \phi(z_2(x))$$

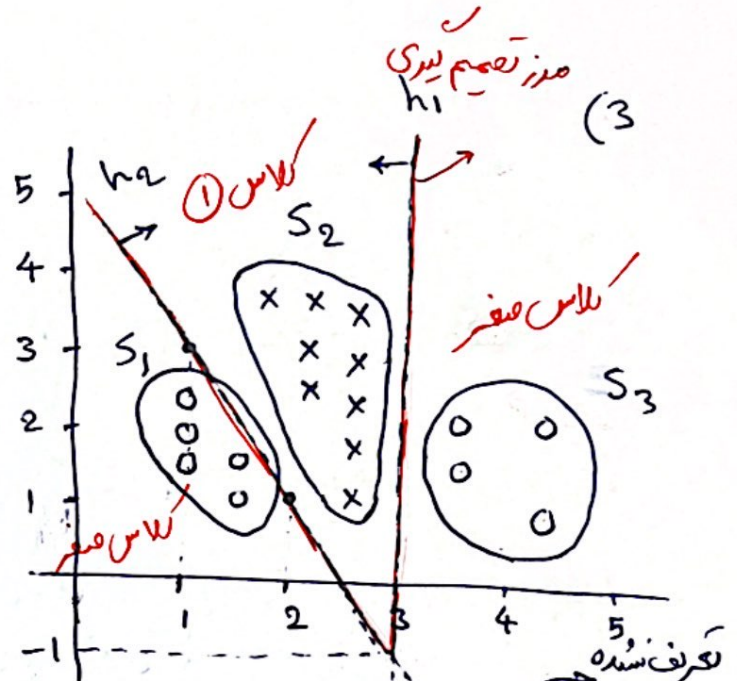
$$\begin{cases} w_{21} = 2 \\ w_{22} = 1 \end{cases}, b_2 = -5$$

$S_1, S_3 \Rightarrow$ کلاس رایبه

$S_2 \Rightarrow$ کلاس x

$$h_1(x) = \begin{cases} 0 & S_3 \\ 1 & S_2 \end{cases}$$

$$h_2(x) = \begin{cases} 0 & S_1 \\ 1 & S_2 \end{cases}$$



تحلیل: برای کلاس فربدر همواره هردو نوروں h_1 و h_2 یک می‌شود و برای کلاس رایبه نیز هردو صفر می‌شوند. اگر جمع این دو نوروں را حساب کنیم:

$$h_1(x) + h_2(x) = \begin{cases} 2 & \text{کلاس فربدر} \\ 1 & \text{کلاس رایبه} \\ 0 & \text{تعریف شده} \end{cases}$$

نقطه ای که هم سمت راست h_1 قرار بگیرد و همزمان سمت چپ h_2 قرار بگیرد مربوط به هیچ کلاسی نیست.

اگر یک بایاس $-\frac{3}{2}$ به جمع بالا بدیم کلاس بندی بعد از اعمال تابع ϕ درست انجام می‌شود

$$h_1(x) + h_2(x) - \frac{3}{2} = \begin{cases} 0.5 & \text{کلاس فربدر} \\ -0.5 & \text{کلاس رایبه} \end{cases} \quad \phi(z) = \begin{cases} 1 & \text{کلاس فربدر} \\ 0 & \text{کلاس رایبه} \end{cases}$$

$$w_{31} = 1, w_{32} = 1, b_3 = -\frac{3}{2}$$

(4) الف) دلیل استفاده از SGD متعذر است:

1) وقتی از mini batch استفاده می‌کنیم بار محاسباتی در یک iteration کاهش می‌یابد یعنی به جای اینکه ماتریس ورودی $[m \times d]$ باشد که بار محاسباتی فنرب و محاسبه بردارین را زیاد می‌کند. می‌توانیم در هر mini batch از K Sample از m استفاده کنیم و ماتریس ورودی $[K \times d]$ باشد در هر iteration $(K \ll m)$

2) از لحاظ معموری و GPU به صرفه است. چون نیاز است تعداد داده کمتری در GPU برای محاسبات در هر iteration ذخیره کنیم و در iteration بعدی داده‌ها قبل از نیاز نداریم را می‌توانیم پاک کنیم.

3) آپدیت بردارین در نتیجه محاسبه لاس زودتر انجام می‌شود. یعنی شبکه را سریع‌تر آپدیت می‌کنیم و می‌توانیم کاهش لاس و عملکرد شبکه را ببینیم.

ب)

$$h(x) = \sigma(w_1 x + b_1)$$

$$p(x) = \sigma(w_2 h(x) + b_2)$$

$$\mathcal{L} = -\frac{1}{m} \sum_i y_i \log p_i + (1-y_i) \log (1-p_i)$$

$$h_i = h(x_i)$$

$$p_i = p(x_i)$$

$$\frac{\partial \mathcal{L}}{\partial p_i} = -\frac{1}{m} \left(\frac{y_i}{p_i} - \frac{1-y_i}{1-p_i} \right)$$

$$\frac{\partial p_i}{\partial w_2} = p_i (1-p_i) h_i$$

$$\frac{\partial \mathcal{L}}{\partial w_2} = \sum_{i=1}^m \frac{\partial \mathcal{L}}{\partial p_i} \cdot \frac{\partial p_i}{\partial w_2} = \sum_{i=1}^m \frac{\partial \mathcal{L}}{\partial p_i} p_i (1-p_i) h_i$$

$$\frac{\partial \mathcal{L}}{\partial b_2} = \sum_{i=1}^m \frac{\partial \mathcal{L}}{\partial P_i} \underbrace{P_i(1-P_i)}_{\frac{\partial P_i}{\partial b_1} = P_i(1-P_i)} \quad \frac{\partial P_i}{\partial h_i} = P_i(1-P_i) \cdot W_2$$

$$\frac{\partial \mathcal{L}}{\partial h_i} = \frac{\partial \mathcal{L}}{\partial P_i} \frac{\partial P_i}{\partial h_i} = \frac{\partial \mathcal{L}}{\partial P_i} P_i(1-P_i) \cdot W_2$$

$$\frac{\partial \mathcal{L}}{\partial w_1} = \sum_{i=1}^m \frac{\partial \mathcal{L}}{\partial h_i} \cdot \frac{\partial h_i}{\partial w_1} = \sum_{i=1}^m \frac{\partial \mathcal{L}}{\partial h_i} \cdot h_i(1-h_i) \cdot x_i^T$$

$$\frac{\partial \mathcal{L}}{\partial b_1} = \sum_{i=1}^m \frac{\partial \mathcal{L}}{\partial h_i} \cdot \frac{\partial h_i}{\partial b_1} = \sum_{i=1}^m \frac{\partial \mathcal{L}}{\partial h_i} \cdot h_i(1-h_i)$$

سوال (4)

Input image: $i \times i \times 3$

بخش اول

First layer: $K_1 \times K_1 \times 3$ فیلتر d_1

$$\left(\frac{H - K + 2P}{S} + 1 \right) \cdot$$

output: $(i - K_1 + 1) \times (i - K_1 + 1) \times d_1$

$$\left(\frac{W - K + 2P}{S} + 1 \right)$$

Second layer: $K_2 \times K_2 \times d_1$ فیلتر d_2

output: $(i - K_1 - K_2 + 2) \times (i - K_1 - K_2 + 2) \times d_2$

Third layer: $K_3 \times K_3 \times d_2$ فیلتر d_3 dilation = n

$$\text{dilation} = n \rightarrow K_{\text{effective}} = K_3 + (K_3 - 1)(n - 1)$$

Third layer:

$$\text{output: } i - K_1 - K_2 + 2 - (K_3 + (K_3 - 1)(n - 1)) + 1$$

$$\Rightarrow (i - K_1 - K_2 - nK_3 + n + 2) \times (i - K_1 - K_2 - nK_3 + n + 2) \times d_3$$

بخش دوم

$$\begin{cases} \text{First layer} & K \times K & \text{dilation} = d_1 \\ \text{Second layer} & K \times K & \text{dilation} = d_2 \\ \text{Third layer} & K \times K & \text{dilation} = d_3 \end{cases}$$

فرمول receptive field اگر $\text{Stride} = 1$ باشد برابر است با:

$$RF_K = 1 + \sum_{i=1}^K (F-1) \times d$$

\downarrow \downarrow \downarrow
 تعداد لایه ها سائز فیلتر dilation

$$RF_3 = 1 + (K-1)d_1 + (K-1)d_2 + (K-1)d_3$$

$$\Rightarrow RF_3 = RF_{\text{output}} = 1 + (K-1)(d_1 + d_2 + d_3)$$

مقدار (K) برای عناصرهای ورودی خروجی