

$$1.a \quad \mathbb{E}_D[P]_{ij} = \mathbb{E}_D[(D \odot X)_{ij}] \stackrel{D_{ij} \sim \text{Bern}(p)}{\text{است}} X_{ij} \mathbb{E}_D[D_{ij}] \stackrel{D_{ij} \sim \text{Bern}(p)}{=} p X_{ij}$$

$$P^T P = \begin{bmatrix} D_{11} x_1^{(1)} & \dots & D_{1n} x_1^{(n)} \\ \vdots & & \vdots \\ D_{1d} x_d^{(1)} & \dots & D_{1d} x_d^{(n)} \end{bmatrix}_{d \times n} \begin{bmatrix} D_{11} x_1^{(1)} & \dots & D_{1d} x_d^{(1)} \\ \vdots & & \vdots \\ D_{n1} x_1^{(n)} & \dots & D_{nd} x_d^{(n)} \end{bmatrix}_{n \times d} = \begin{bmatrix} \sum_{k=1}^n D_{k1}^2 x_1^{(k)2} & \sum_{k=1}^n D_{k1} D_{k2} x_1^{(k)} x_2^{(k)} & \dots \\ \vdots & & \end{bmatrix}_{d \times d}$$

$$\mathbb{E}_D[P^T P]_{ij} = \begin{cases} \mathbb{E}_D \left[\sum_{k=1}^n D_{ki}^2 x_i^{(k)2} \right] = \sum_{k=1}^n \underbrace{\mathbb{E}_D[D_{ki}^2]}_{p \text{ ①}} x_i^{(k)2} = p \sum_{k=1}^n x_i^{(k)} x_i^{(k)} = p (X^T X)_{ij} & i=j \\ \mathbb{E}_D \left[\sum_{k=1}^n D_{ki} D_{kj} x_i^{(k)} x_j^{(k)} \right] = \sum_{k=1}^n \underbrace{\mathbb{E}_D[D_{ki} D_{kj}]}_{p^2 \text{ ②}} x_i^{(k)} x_j^{(k)} = p^2 \sum_{k=1}^n x_i^{(k)} x_j^{(k)} = p^2 (X^T X)_{ij} & i \neq j \end{cases}$$

$$D_{ki} \sim \text{Bern}(p) \rightarrow D_{ki}^2 \sim \text{Bern}(p) \rightarrow \mathbb{E}[D_{ki}^2] = p \text{ ①}$$

$$D_{ki} \perp D_{kj} \rightarrow \mathbb{E}_D[D_{ki} D_{kj}] = \mathbb{E}_D[D_{ki}] \mathbb{E}_D[D_{kj}] = p^2 \text{ ②}$$

1.b

$$\mathcal{L}(\hat{w}) = \mathbb{E}_D[\|y - P \hat{w}\|_2^2] = \mathbb{E}_D[\|y\|^2 + \underbrace{\|P \hat{w}\|^2}_{\hat{w}^T P^T P \hat{w}} - 2 y^T P \hat{w}]$$

$$= \|y\|^2 + \hat{w}^T \mathbb{E}_D[P^T P] \hat{w} - 2 y^T \mathbb{E}_D[P] \hat{w}$$

$$\stackrel{*}{=} \underbrace{\|y\|^2 - 2 y^T p X \hat{w} + p^2 \hat{w}^T X^T X \hat{w}}_{\|y - p X \hat{w}\|^2} + p(1-p) \underbrace{\hat{w}^T \hat{\Gamma}^2 \hat{w}}_{\|\hat{\Gamma} \hat{w}\|^2}$$

طبق بخش 1.a می دانیم که $\mathbb{E}_D[P] = pX$
 همچنین برای $\mathbb{E}_D[P^T P]$ نیز می توان فرم بسته زیر را داشت:

$$\mathbb{E}_D[P^T P] = p^2 X^T X + p \hat{\Gamma}^2 - p^2 \hat{\Gamma}^2$$

برای اینکه عناصر قطر اولی $p^2 X^T X$ را حذف کنیم
 $\hat{\Gamma} = \text{diag}(X^T X)^{\frac{1}{2}}$

تغییر متغیرهای زیر را به ترتیب انجام می دهیم :

1.c $\xrightarrow{p\hat{w}=w} \mathcal{L}(w) = \|y - Xw\|^2 + \frac{1-\rho}{\rho} \|\hat{\Gamma} w\|^2$

$\xrightarrow{\frac{1-\rho}{\rho} \hat{\Gamma} = \Gamma} \mathcal{L}(w) = \|y - Xw\|^2 + \|\Gamma w\|^2$

1.d

ابتدا سعی می کنیم بخش regularization را ایجاد کنیم

$$\Gamma w = \sqrt{\lambda} \tilde{w} \rightarrow \begin{aligned} w &= \sqrt{\lambda} \Gamma^{-1} \tilde{w} \\ \tilde{w} &= \frac{1}{\sqrt{\lambda}} \Gamma w \end{aligned}$$

$$\rightarrow \mathcal{L}(\tilde{w}) = \|y - X \underbrace{\sqrt{\lambda} \Gamma^{-1}}_{\tilde{X}} \tilde{w}\|^2 + \|\sqrt{\lambda} \tilde{w}\|^2 \xrightarrow{\tilde{X} = \sqrt{\lambda} X \Gamma^{-1}} \mathcal{L}(\tilde{w}) = \|y - \tilde{X} \tilde{w}\|^2 + \lambda \|\tilde{w}\|^2$$

1.e

$$\tilde{X} = \sqrt{\lambda} X \Gamma^{-1} = \sqrt{\lambda} \begin{bmatrix} x_1^{(1)} & \dots & x_d^{(1)} \\ \vdots & & \vdots \\ x_1^{(N)} & \dots & x_d^{(N)} \end{bmatrix} \begin{bmatrix} \Gamma_{1,1}^{-1} & & 0 \\ & \ddots & \\ 0 & & \Gamma_{d,d}^{-1} \end{bmatrix} = \sqrt{\lambda} \begin{bmatrix} \Gamma_{1,1}^{-1} x_1^{(1)} & \dots & \Gamma_{d,d}^{-1} x_d^{(1)} \\ \vdots & & \vdots \\ \Gamma_{1,1}^{-1} x_1^{(N)} & \dots & \Gamma_{d,d}^{-1} x_d^{(N)} \end{bmatrix}$$

زیر Γ متناسب با نورم ستون d ام X است. بنابراین مشاهده می شود که ماتریس به دست آمده، هر ستونش نرمالایز شده است یعنی نورم ستون های \tilde{X} یکسان شده است. همچنین توسط λ می توان مقدار آنرا کنترل کرد.

این مسئله در حقیقت همان batch normalization است زیرا آنجا نیز تغییرها به صورت جدا در هر batch نرمالایز می شوند. می توان گفت اگر نیز به نوعی تغییر λ در batch norm عمل می کند

1.f

$$\frac{\partial \mathcal{J}}{\partial w_i} = \left(y_d - \sum_{k=1}^n (w_k + \delta_k) x_k \right) (-x_i)$$

$$\rightarrow \mathbb{E}_{\delta} \left[\frac{\partial \mathcal{J}}{\partial w_i} \right] = \mathbb{E}_{\delta} \left[-y_d x_i + x_i \sum_{k=1}^n w_k x_k + x_i \sum_{k=1}^n \delta_k x_k \right]$$

$$= -y_d x_i + x_i \sum_{k=1}^n w_k x_k + x_i \sum_{k=1}^n \underbrace{\mathbb{E}_{\delta_k} [\delta_k]}_{0} x_k = x_i \sum_{k=1}^n w_k x_k - y_d x_i$$

* مشاهده می شود که چون δ_k از توزیع نرمال که مقدار آن صفر است پیروی می کند $\mathbb{E} \left[\frac{\partial \mathcal{J}}{\partial w_i} \right]$ با $\mathbb{E} \left[\frac{\partial \mathcal{J}_n}{\partial w_i} \right]$ تابع هزینه بدون dropout است، یکی می شود

1.g

$$\mathcal{J}_n = \frac{1}{2} \left(y_d - \sum_{k=1}^n w_k x_k \right)^2$$

تابع هزینه بدون dropout را به صورت زیر تعریف می کنیم:

حال امید ریاضی تابع با dropout را حساب می کنیم:

$$\mathbb{E}[\mathcal{J}] = \mathbb{E} \left[\frac{1}{2} \left(y_d - \underbrace{\sum_{k=1}^n w_k x_k}_H - \sum_{k=1}^n \delta_k x_k \right)^2 \right]$$

$$= \frac{1}{2} \mathbb{E} \left[H^2 + \underbrace{\left(\sum_{k=1}^n x_k \delta_k \right)^2}_{\text{این عبارت، z نامیده می شود}} - 2H \sum_{k=1}^n \delta_k x_k \right] = \frac{1}{2} \mathbb{E}[H^2] + \frac{1}{2} \mathbb{E}[Z^2] - 2H \sum_{k=1}^n \underbrace{\mathbb{E}[\delta_k]}_{0} x_k$$

می دانیم که $\mathbb{E}[Z]$ برابر 0 است بنابراین $\mathbb{E}[Z^2]$ در حقیقت همان واریانس Z است و داریم:

$$\text{Var}(Z) = \text{Var} \left(\sum_{k=1}^n x_k \delta_k \right) \stackrel{\text{از استقلال } \delta_k}{=} \sum_{k=1}^n \text{Var}(x_k \delta_k) = \sum_{k=1}^n x_k^2 \text{Var}(\delta_k) = \sum_{k=1}^n x_k^2 \alpha w_k^2$$

$$\mathbb{E}[\mathcal{J}] = \mathcal{J}_n + \frac{\alpha}{2} \sum_{k=1}^n x_k^2 w_k^2$$

همان طور که مشاهده می شود، فرمی شبیه L2 regularization داریم.

البته در حالت کلی می توان نشان داد اگر از تابع هزینه least squares استفاده شود، صرف نظر از توزیع که در dropout استفاده شود، به تقریبی از L2 regularization می رسیم!

1.h

additive gaussian dropout

- یک نویز گاوسی به وزن هر نورون اضافه می‌کند و به این صورت تأثیر نورون‌ها را کم یا زیاد می‌کند. دقت شود که در dropout عادی یک متغیر تصادفی بزرگی در وزن اضافه می‌شود و dropout گاوسی یک جبر تخمینی است که از توزیع بیوسه گاوسی استفاده می‌کند.
- به واسطه اضافه کردن نویز به activation ها، قدرت generalization بیشتری به مدل می‌دهد و مدل را robust تر می‌کند.

spatial dropout

- در شبکه‌های CNN کاربرد دارد. به جای اعمال روی نورون‌ها، روی تمام feature map اعمال می‌شود. یعنی در حالت عادی که از توزیع بزرگی استفاده می‌شود با احتمال $1-p$ ممکن است هر کدام از feature map ها حذف شوند.
- مدل را مجبور می‌کند که از تمام feature map ها استفاده کند و به یک پیرن خاص تکی نباشد \rightarrow مدل robust تر.

2

$$\frac{\partial z}{\partial w} = \begin{bmatrix} \frac{\partial z_1}{\partial w_1} & \dots & \frac{\partial z_1}{\partial w_k} \\ \vdots & & \vdots \\ \frac{\partial z_l}{\partial w_1} & \dots & \frac{\partial z_l}{\partial w_k} \end{bmatrix}_{l \times k}$$

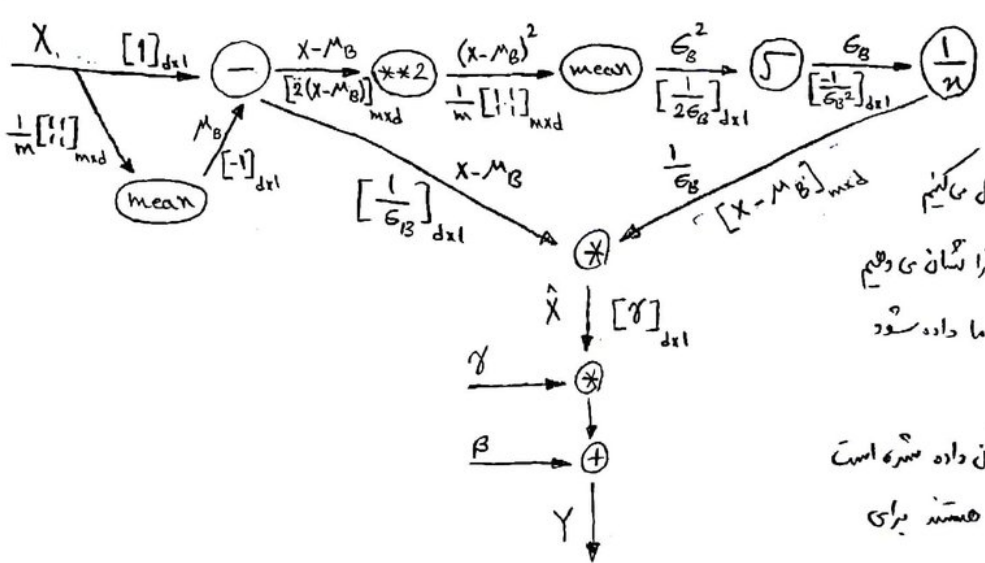


$$\frac{\partial L}{\partial w} = \left[\frac{\partial z}{\partial w} \right]_{k \times l}^T @ \left[\frac{\partial L}{\partial z} \right]_{l \times 1} \rightarrow \frac{\partial L}{\partial w_j} = \left[\frac{\partial z}{\partial w_j} \right]_{1 \times l} @ \left[\frac{\partial L}{\partial z} \right]_{l \times 1} = \begin{bmatrix} \frac{\partial z_1}{\partial w_j} \\ \vdots \\ \frac{\partial z_l}{\partial w_j} \end{bmatrix} \cdot \begin{bmatrix} \frac{\partial L}{\partial z_1} \\ \vdots \\ \frac{\partial L}{\partial z_l} \end{bmatrix}$$

$$\rightarrow \frac{\partial L}{\partial w_j} = \sum_{i=1}^l \underbrace{\frac{\partial z_i}{\partial w_j}}_{X_{i+j}} \times \underbrace{\frac{\partial L}{\partial z_i}}_{w'_i} \rightarrow \frac{\partial L}{\partial w_j} = \sum_{i=1}^l w'_i X_{j+i}$$

فرض سطر مقدار w'_i را برای همه مقادیر i می‌دانیم
و مشاهده می‌شود که دقیقاً به نرم اعمال کانولوشن مانند
صورت سوال رسیدیم

3.a



⊕ مسئله را در کلی‌ترین حالت و به صورت ماتریسی حل می‌کنیم
⊕ تمام ضرب‌ها element wise است و با * آنرا نشان می‌دهیم
⊕ $\frac{\partial L}{\partial Y}$ گرادیان upstream است که باید به ما داده شود

⊕ مقادیر forward, backward روی گراف نشان داده شده است
مقادیری که ایجاد شده ذکر شده و درون پرانتز هستند برای
backward

$$\frac{\partial L}{\partial \gamma} = \text{sum} \left(\left[\frac{\partial L}{\partial Y} \right]_{m \times d} * \left[\hat{X} \right]_{m \times d}, \text{axis} = 0 \right) = \sum_{i=1}^m \frac{\partial L}{\partial y_i} * \hat{x}_i$$

$$\frac{\partial L}{\partial \beta} = \text{sum} \left(\left[\frac{\partial L}{\partial Y} \right]_{m \times d}, \text{axis} = 0 \right) = \sum_{i=1}^m \frac{\partial L}{\partial y_i}$$

3.a

برای محاسبه $\frac{\partial L}{\partial X}$ به جای γ می‌نویسیم و حساب می‌کنیم:

$$\begin{aligned} \frac{\partial L}{\partial X} &= \left[\frac{\partial L}{\partial Y} \right]_{m \times d} * \left[\gamma \right]_{d \times 1} * \left[\frac{1}{\sigma_B} \right]_{d \times 1} \\ &+ \text{Sum} \left(\left[\frac{\partial L}{\partial Y} \right]_{m \times d} * \left[\gamma \right]_{d \times 1} * \left[\frac{1}{\sigma_B} \right]_{d \times 1} * [-1]_{d \times 1}, \text{axis}=0 \right) * \frac{1}{m} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}_{m \times d} \\ &+ \text{Sum} \left(\left[\frac{\partial L}{\partial Y} \right]_{m \times d} * \left[\gamma \right]_{d \times 1} * [X - \mu_B]_{m \times d}, \text{axis}=0 \right) * \left[\frac{-1}{\sigma_B^2} \right]_{d \times 1} * \left[\frac{1}{2\sigma_B} \right]_{d \times 1} * \frac{1}{m} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}_{m \times d} * [2(X - \mu_B)]_{m \times d} \\ &+ \text{Sum} \left(\left[\frac{\partial L}{\partial Y} \right]_{m \times d} * \left[\gamma \right]_{d \times 1} * [X - \mu_B]_{m \times d}, \text{axis}=0 \right) * \left[\frac{-1}{\sigma_B^2} \right]_{d \times 1} * \left[\frac{1}{2\sigma_B} \right]_{d \times 1} * \frac{1}{m} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}_{m \times d} \\ &* \text{Sum} \left([2(X - \mu_B)]_{m \times d} * [-1]_{d \times 1}, \text{axis}=0 \right) * \frac{1}{m} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}_{m \times d} \end{aligned}$$

برای کوتاه‌تر شدن جواب به جای γ از $\frac{\partial L}{\partial \hat{X}}$ استفاده می‌کنیم.
 همچنین برای 2 عبارت آخر σ_B ها را زیر $X - \mu_B$ می‌آوریم تا به \hat{X} تبدیل شوند.

$$\begin{aligned} \frac{\partial L}{\partial X} &= \left[\frac{\partial L}{\partial \hat{X}} \right]_{m \times d} * \left[\frac{1}{\sigma_B} \right]_{d \times 1} \\ &- \left[\frac{1}{m\sigma_B} * \sum_{i=1}^m \frac{\partial L}{\partial \hat{n}_i} \right]_{d \times 1} * \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}_{m \times d} \\ &- \left[\sum_{i=1}^m \frac{\partial L}{\partial \hat{n}_i} * \hat{n}_i \right]_{d \times 1} * \left[\frac{1}{m\sigma_B} \right]_{d \times 1} * [\hat{X}]_{m \times d} \\ &- \left[\sum_{i=1}^m \frac{\partial L}{\partial \hat{n}_i} * \hat{n}_i \right]_{d \times 1} * \left[\frac{1}{m\sigma_B} \right]_{d \times 1} * \underbrace{\left[\sum_{j=1}^m \hat{n}_j \right]_{d \times 1}}_{\text{مجموع ردیف‌های آن 0 می‌شود}} * \frac{1}{m} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}_{m \times d} \end{aligned}$$

چون \hat{X} ماتریس standardize شده است
 پس جمع ردیف‌های آن 0 می‌شود و عبارت آخر بی‌تأثیر است

بنابراین برای $\frac{\partial L}{\partial \hat{n}_i}$ کافی است به سطرهای ماتریس $\frac{\partial L}{\partial X}$ نگاه کنیم و داریم:

$$\frac{\partial L}{\partial \hat{n}_i} = \left(\frac{\partial L}{\partial \hat{n}_i} * \frac{1}{\sigma_B} \right) - \left(\frac{1}{m\sigma_B} * \sum_{j=1}^m \frac{\partial L}{\partial \hat{n}_j} \right) - \left(\frac{1}{m\sigma_B} * \hat{n}_i * \sum_{j=1}^m \frac{\partial L}{\partial \hat{n}_j} * \hat{n}_j \right)$$

3.b

در حالت کلی در شبکه های مصنوعی به واسطه تغییر وزن های لایه های قبل (در back prop) توزیع ورودی لایه های بعدی تغییر می کند و این مسئله روی گزاردان جدیری که برای آن لایه محاسبه می شود تأثیر گذاشته و باعث می شود گزاردان، سیگنال sub optimal ای برای آپدیت وزن های آن لایه باشد. batch norm با محدود کردن توزیع pre activation های هر لایه، این مسئله را از بین می برد و باعث می شود گزاردان های بعدی برای آپدیت وزن ها داشته باشیم. بنابراین مدل سریع تر و stable تر آموزش می بیند. (internal covariate shift) به بیان دیگر coupling بین لایه ها کمتر می شود.

جدای از پایدار تر شدن مدل نسبت به تغییرات جزئی در وزن ها و به طور کلی ورودی هر لایه، batch norm باعث می شود هر لایه بتواند pre activation هایش را مستقل از لایه های دیگر، در رنج مناسب activation func بیارد. این مسئله به خصوص در activation func های مانند sigmoid مهم است که اگر ورودی شان در رنج مناسبی نباشد saturate می شوند و گزاردان های بسیار کمی (تدریج منفی) خواهند داشت.

بنابراین batch norm به شکل vanishing gradients هم کمک می کند.

طبق موارد ذکر شده batch norm نیاز به weight initialization را ساده تر کرده و اجازه استفاده از learning rate های بزرگتری را می دهد.

اجازه استفاده از activation func های پیچیده تر را داده و می توان شبکه های عمیقتری داشت زیرا وابستگی لایه ها کمتر شده است و جریان گزاردان بهتری داریم.

3.c

استفاده از \hat{x} قدرت مدل را به سمت کاهش می دهد زیرا pre activation ما همگی مدل منفی و با وابستگی 1 خواهند بود. به این صورت تنها از بخشی از activation func استفاده می شود. مثلاً در sigmoid تنها از قسمت شبه خطی آن استفاده می شود و یا در relu دیگر آن expressive بودن نورون ها در بخش مثبت را نداریم.

هدف اصلی آن است که $y = \gamma \hat{x} + \beta$ هر میانگین و واریانس بتواند داشته باشد اما میانگین و واریانس آن از لایه های قبل decouple شده باشد.

3.d

یکی از دغدغه های که در استفاده از batch norm داریم آن است که moving average به دست آمده در زمان آموزش تدریجی به میانگین هر batch در زمان تست است. به واسطه مشکل ذکر شده میانگین به دست آمده برای هر batch در زمان آموزش، به سمت یکی از کلاس های سنگ یا کمر به بایاس می شود و بنابراین توزیع داده های مقادیر را در زمان آموزش و تست (برای هر batch) فاصله می دهد. بنابراین pre activation ها به صورت مطلوب نخواهند بود.

بنابراین دین آموزش و تست discrepancy خواهیم داشت.

4.a
 خنیر ایجاد شکل نمی‌کند. وقتی پس از لایه conv از batch norm استفاده می‌شود، مقادیر آن لایه standardize می‌شوند و بایاس که به همه فیلترهای هر feature map اضافه شده است، از بین می‌رود. ضمناً خود لایه batch norm دارای پارامتر β است که حکم بایاس را دارد. بنابراین معمولاً وقتی از batch norm استفاده می‌شود، بایاس لایه قبل را False می‌کنند.

4.b

وزن می‌گیریم منظور از وزن‌ها، وزن‌های conv layer باشد
 در این صورت مملکت شبکه در زمان تست دچار مشکل می‌شود. چه وزن‌ها در α ضرب شوند چه ورودی، ورودی BN α برابر می‌شود. حال چون μ_B و σ_B از moving average زمان آموزش به دست آمده پس داریم:

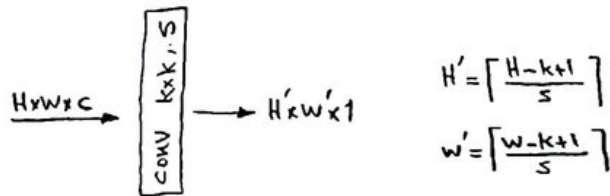
$$BN_{\gamma, \beta}(\alpha XW) = \gamma * \frac{\alpha XW - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} + \beta$$

که با حالت مورد انتظار برای تست متفاوت است.

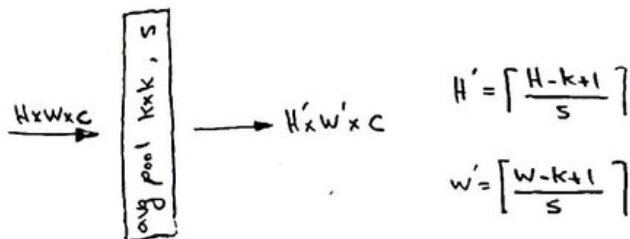
④ نمایان ذکر است که این مورد در زمان آموزش تأثیری ندارد زیرا standardization اگر آنرا خنیر می‌کند
 ⑤ اگر منظور از وزن‌ها μ_B ، σ_B ، γ و β نیز باشد، در حالتی که وزن‌ها در α ضرب شوند ضربی BN نیز تقریباً α برابر می‌شود

4.c

برای اعمال یک لایه:

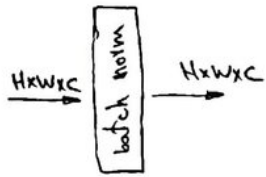


$$\begin{aligned} \text{تعداد ضرب} &= H' \times W' \times k \times k \times C \\ \text{تعداد جمع} &= H' \times W' \times (k \times k \times C - 1) \end{aligned} \rightarrow \text{مجموع} = (H-k)(W-k) \frac{k^2 C}{s^2}$$



$$\text{تعداد جمع و تقسیم} = H' \times W' \times C \times (k \times k - 1 + 1) \rightarrow \text{مجموع} = (H-k)(W-k) \frac{k^2 C}{s^2}$$

ادامه
4.c



در حالت کلی اگر در فاز train باسیم m ، batch size داریم:

$$\text{calc } \mu_B : C \left(\overbrace{m-1+1}^{\text{تقسیم جمع}} \right) HW$$

$$\text{calc } \sigma_B : C \left(\underbrace{1}_{\text{تفریق}} + \underbrace{1}_{\text{توان 2}} + \underbrace{m-1}_{\text{جمع}} + \underbrace{1}_{\text{تقسیم}} + \underbrace{1}_{5} \right) HW$$

$$\text{calc } \hat{X} : C \left(\underbrace{1}_{\text{تفریق}} + \underbrace{1}_{\text{تقسیم}} \right) HW$$

$$\text{calc } Y : C \left(\underbrace{1}_{\text{جمع ضرب}} + \underbrace{1}_{\text{جمع}} \right) HW$$

در صورت صبرال چون گفته شده برای یک تصویر ورودی پس مرتبه محاسباتی CHW می شود

4.d

از 1×1 conv می توان برای کنترل تعداد چنل ها استفاده کرد. اکثراً کاربرد آن، آن است که قبل از لایه های conv دیگر با ساینر کنترل بالاتر، از 1×1 conv استفاده می کنیم تا عمق را کاهش دهیم و هزینه محاسبات را کاهش دهیم. به عنوان مثال inception architecture گلول قبل از لایه های کانولوشنی 3×3 و 5×5 از آن استفاده شده. (به طور مستقیم هم می توان عمق را توسط لایه های conv با اندازه کنترل بالاتر تعیین کرد اما بسیار کمتر خواهد بود)

4.e

لایه های pooling مزایای زیر را دارند:

- با کاهش ابعاد feature map ها هزینه محاسبات را پایین می آورند
- با حذف پنجره های کم اهمیت تر مدل را تا حدی در برابر تغییر ورودی robust می کنند و تا حدی translation invariance ایجاد می کنند
- generalization مدل را افزایش می دهند و تا حدی overfitting را بهبود می دهند
- به واسطه کاهش ابعاد، receptive field را افزایش می دهند

نکته مهم آن است که این موارد را توسط روش های دیگری نیز می توان به دست آورد اما pooling پارامتری برای learn شدن ندارد و محاسبات را سریع تر و efficient تر می کند.

4. P

اعداد مختلفی را برای سائز ورودی امتحان کردیم به طوری که به سائز کرنل تمام مراحل بخش پذیر باشد. عدد 21 و مضارب آن این ویژگی را دارند:

$$21 \times 21 \xrightarrow{\text{conv } 3 \times 3 \text{ } s=2} 10 \times 10 \xrightarrow{\text{pool } 2 \times 2} 5 \times 5 \xrightarrow{\text{conv } 3 \times 3 \text{ } s=2} 2 \xrightarrow{\text{pool } 2 \times 2} 1$$

پس هر پیکسل خروجی تحت تأثیر 21×21 پیکسل ورودی است

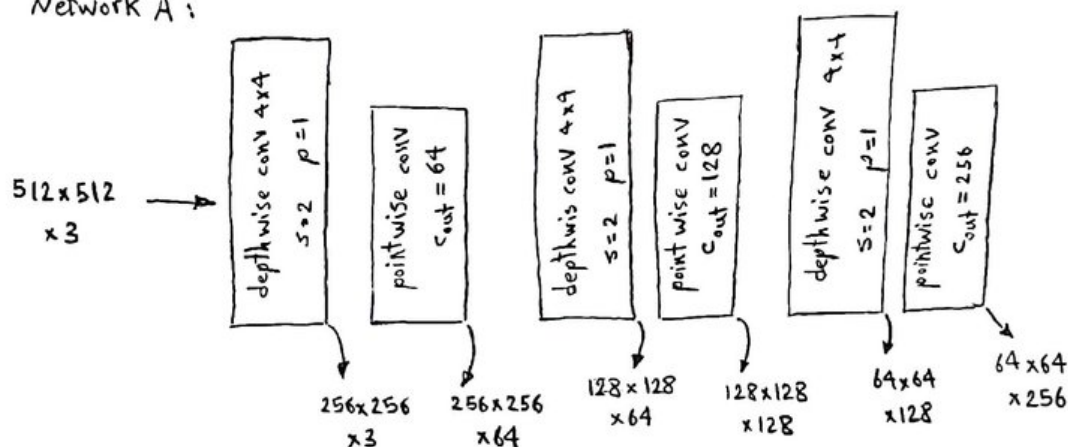
$$\text{Receptive Field} = 21^2 = 441$$

5.a

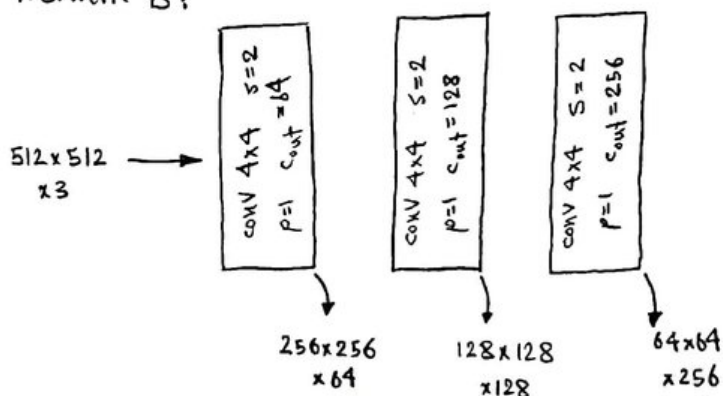
نیزدیک به دست آوردن
مقیاس H, W

$$O = \left\lceil \frac{i+2p-k+1}{s} \right\rceil$$

Network A:



Network B:



5.b

Network A parameters:

$$\begin{aligned}
 &(4 \times 4 \times 3) + (3 \times 64) = 240 \\
 &+ (4 \times 4 \times 64) + (64 \times 128) = 9,216 \\
 &+ (4 \times 4 \times 128) + (128 \times 256) = 34,816 \\
 &= 44,272
 \end{aligned}$$

Network B parameters:

$$\begin{aligned}
 &(4 \times 4 \times 3) \times 64 = 3,072 \\
 &+ (4 \times 4 \times 64) \times 128 = 131,072 \\
 &+ (4 \times 4 \times 128) \times 256 = 524,288 \\
 &= 658,432
 \end{aligned}$$

5.c

Network A mul operations:

$$\begin{aligned}
 & (4 \times 4)(256 \times 256 \times 3) + (3 \times 256 \times 256 \times 64) = 15,728,640 \\
 & + \\
 & (4 \times 4)(128 \times 128 \times 64) + (64 \times 128 \times 128 \times 128) = 150,994,944 \\
 & + \\
 & (4 \times 4)(64 \times 64 \times 128) + (128 \times 64 \times 64 \times 256) = 142,606,336 \\
 & = 309,329,920
 \end{aligned}$$

Network B mul operations:

$$\begin{aligned}
 & (4 \times 4 \times 3)(256 \times 256 \times 64) = 201,326,592 \\
 & + \\
 & (4 \times 4 \times 64)(128 \times 128 \times 128) = 2,147,483,650 \\
 & + \\
 & (4 \times 4 \times 128)(64 \times 64 \times 256) = 2,147,483,650 \\
 & = 4,496,293,888
 \end{aligned}$$

5.d

با اینکه ابعاد Feature map ها در نهایت یکسان است اما شبکه A پارامتر و محاسبات بسیار کمتری دارد
از این نظر می توان از این روش در شرایطی که منابع محدودی داریم استفاده کنیم

5.e

هدف از معماری MobileNet آن است که حجم مدل، پیچیدگی محاسبات را پایین بیاورد تا بتوان از لایه های کانولوشنی در دستگاه های با منابع محدود، mobile نیز استفاده کرد.

به این منظور نشان داده می شود که هر لایه کانولوشن عادی را می توان به 2 لایه depthwise conv و pointwise conv تقسیم کرد به طوری که ابعاد ورودی و خروجی این لایه ها با کانولوشن عادی یکسان است اما تعداد پارامتر ها و مرتبه محاسبات کمتر می شود (مانند چیزی که در بحث قبل دیدیم) در لایه depthwise conv کانولوشن تنها بر روی یک چنل درونی (عمق=1) انجام می شود و بنابراین به تعداد چنل های ورودی فیلتر $k \times k$ داریم پس در لایه pointwise conv به اندازه ی تعداد چنل های خروجی مورد نیاز، فیلتر های 1×1 استفاده می کنیم. در حقیقت لایه های depthwise در هر چنل به صورت مجزا دنبال پترن می گردند اما لایه های pointwise از ترکیب چنل ها استفاده کرده اما از همسایگی های داخل هر Feature map یعنی ویژگی های spatial استفاده نمی کنند.

به طور کلی در حالتی که ابعاد ورودی و خروجی تغییر نکند و فیلتر $k \times k$ داشته باشیم ($stride=1$). درونی conv و depthwise separable را، به صورت زیر می شود مقایسه کرد:

	#params	#operations
depthwise separable conv	$NM + Mk^2$	$NMHW + MHWk^2$
conv	NMk^2	$NMHWk^2$

همان طور که مشاهده می شود مرتبه محاسبات و تعداد پارامتر ها بسیار کمتر شده و فضا نشان داده شده که دقت کار تا حد خوبی حفظ شده است

5.f

width multiplier (α)

مقدار α بین 0 و 1 است که در عمق (تعداد چنل) ضرب شده و آنرا کاهش می دهد یعنی: $N \rightarrow \alpha N$ و $M \rightarrow \alpha M$. بنابراین تعداد عملیات های محاسباتی و تعداد پارامتر ها کاهش می یابد که می توان مقدار آنرا توسط جدول بحث قبل حساب کرد. کاربرد این هایپر پارامتر نازک تر کردن لایه ها در مواردی است که منابع کمتری داریم. البته قطعاً دقت نیز با کاهش α کاهش می یابد.

resolution multiplier (ϕ)

مقدار ϕ بین 0 و 1 است که در ابعاد (H, W) ورودی ضرب شده و آنرا کاهش می دهد یعنی: $W \rightarrow \phi W$ و $H \rightarrow \phi H$. این هایپر پارامتر تعداد پارامتر ها را تغییر می دهد (از جدول معلوم است) اما تعداد عملیات ها را ϕ^2 برابر می کند. به صورت فنی با تعیین resolution ورودی این هایپر پارامتر ست می شود.

6.a

به طور معمول از این تابع ضرر به عنوان regularization term در کنار توابع دیگری مانند cross entropy loss استفاده می شود.

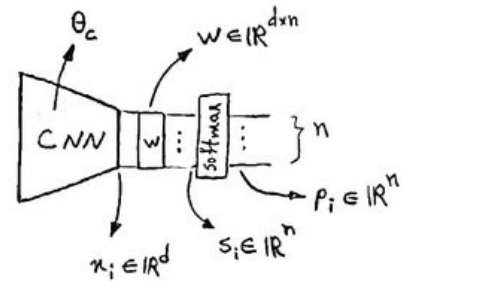
یکی از مشکلاتی که در دسته بندی و به خصوص در حوزه تشخیص چهره با آن مواجه هستیم کم بودن فاصله کلاس های مختلف از هم است. به این صورت که مدل دسته بندی را به درستی انجام می دهد اما representation های به دست آمده برای سیپل های کلاس های مختلف به هم نزدیک هستند. اکثر دیتای تست با اعمال تغییر کوچکی از سیپل درست دسته بندی شده به دست آید، ممکن است مدل آن را به کلاس دیگری نسبت دهد.

استفاده از center loss به همراه cross entropy loss باعث می شود که intra class compactness و inter class dispersion را هم در نظر بگیرد. این دو معیار را به هم داشته باشیم و مرکز بین representation های سیپل های هر کلاس بیست شود.

برای این منظور از فرمول $L_c = \frac{1}{2} \sum_{i=1}^m \|x_i - c_j\|^2$ برای center loss استفاده می شود که فاصله representation سیپل های هر کلاس (یعنی x_i) را از مرکز آن دسته کند می کند.

⊛ نکته سائز اصیت آن است که استفاده هر دو لاس ذکر شده به صورت joint است که بهترین بهبود را ایجاد می کند.
• اگر فقط از cross entropy loss استفاده شود که همان شکل اولیه را داریم.
• اگر فقط از center loss استفاده شود مدل ترجیح می شود که همه representation ها را به دست آورد.

6.b $L = L_s + \lambda L_c$
cross entropy loss
center loss



update $w \rightarrow w^{t+1} = w^t - \eta^t \left(\frac{\partial L}{\partial w^t} \right)$

update $\theta_c \rightarrow \theta_c^{t+1} = \theta_c^t - \eta^t \sum_{i=1}^m \frac{\partial L}{\partial x_i^t} \cdot \frac{\partial x_i^t}{\partial \theta_c^t}$
داخل CNN $\left[\frac{\partial L}{\partial x^t} \right] \cdot \left[\frac{\partial x^t}{\partial \theta_c^t} \right]$

$\frac{\partial L}{\partial w} = \frac{\partial L_s}{\partial w} = X^T @ \frac{\partial L_s}{\partial s} = [X^T] @ \begin{bmatrix} p_{j-1} \\ p_j \end{bmatrix}$
where $\begin{cases} p_{j-1} & j = \text{true class} \\ p_j & j \neq \text{true class} \end{cases}$

update $c_j \rightarrow c_j^{t+1} = c_j^t - \alpha \Delta c_j^t$

$\Delta c_j = \frac{\sum_{i=1}^m 1_{y_i=j} \times (c_j - x_i)}{1 + \sum_{i=1}^m 1_{y_i=j}}$

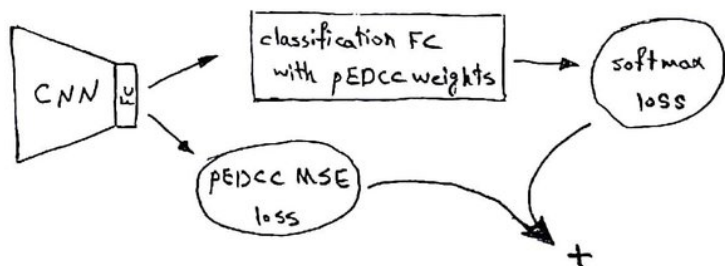
$\frac{\partial L}{\partial x} = \frac{\partial L_s}{\partial x} + \lambda \frac{\partial L_c}{\partial x}$
 $\begin{bmatrix} x_1 - c_{y_1} \\ \vdots \\ x_m - c_{y_m} \end{bmatrix}_{m \times d}$

⊛ در حالت کلی باید برای به دست آوردن مراکز دسته ها، میانگین را به ازای تمام سیپل ها بگیریم اما این کار هزینه بسیار زیادی دارد. بنابراین Δc_j^t را به صورت رویه او به عنوان جایگزین این کار تعریف می کنیم.
ماند که از فرمول معلوم است، میانگین فاصله تا مرکز برای سیپل های هر دسته (به ازای mini batch) حساب می شود و در جهت عکس آن آپدیت می کنیم.

6.c

در این مقاله از PEDCC استفاده شده است تا به صورت evenly distributed مراکز دسته‌ها چسبند.

در روش center loss برای مراکز دسته‌ها از یک تقریب استفاده کردیم. در اینجا نیز تفاوت آن است که بارهای دیگری این مراکز دسته‌ها را مدیریت می‌کنیم. اگر hyper sphere واحد را در فضای representation تصور کنیم، در این روش مراکز دسته‌ها در ابتدا به صورت evenly distributed روی سطح این ابرکره قرار می‌گیرند و فیلکس می‌شوند. در ادامه به صورت تقریباً مشابه از توابع خطای softmax و MSE به صورت زیر استفاده می‌شود:



نکته حائز اهمیت آن است که چون مراکز دسته‌ها فیلکس هستند یک لایه FC جدید قبل softmax اضافه می‌شود که وزن‌های آن فریز شده و در فرآیند آموزش آپدیت نمی‌شوند.

بنابراین با جاگذاری و فیلکس کردن وزن‌های PEDCC در لایه قبل softmax، مراکز دسته‌ها فیلکس می‌شوند و در تئوری بیشترین فاصله را از هم دارند پس 2 تابع لاس مانند قبل intra class compactness و inter class dispersion را ایجاد می‌کنند اما به واسطه مراکز دسته‌هایی که از هم فاصله مناسبی دارند تقسیم بهتری در تشکیل representation سمبل‌ها به وجود می‌آید.

* در برخی سوال‌ها با عملی پایایی همفکری داریم.