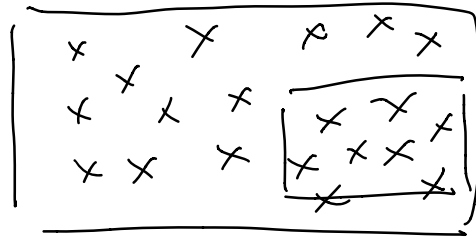


Tuning process

¹ α , ² β , $\beta_1, \beta_2, \epsilon$, # layers, ²# hidden units, ²mini-batch size
²learning rate decay

* Try random values: don't use a grid

* Coarse to fine



Use appropriate scale to pick up hyperparameters

eg. weighted smooth average

β search between $0.9 \sim 0.999$

$$\Downarrow$$
$$[10^{-1} \sim 10^{-3}]$$

`np.power(10, 3 * np.random.randn())`

batch normalization

$$\mu = \frac{1}{m} \sum_i z^{(i)}$$

$$\sigma^2 = \frac{1}{m} \sum_i (z_i - \mu)^2$$

$$z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$\tilde{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta$$

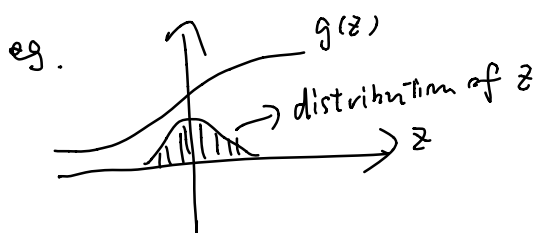
if $\gamma = 1/\sigma^2 + \epsilon$, $\beta = \mu$

then $\tilde{z}^{(i)} = z^{(i)}$

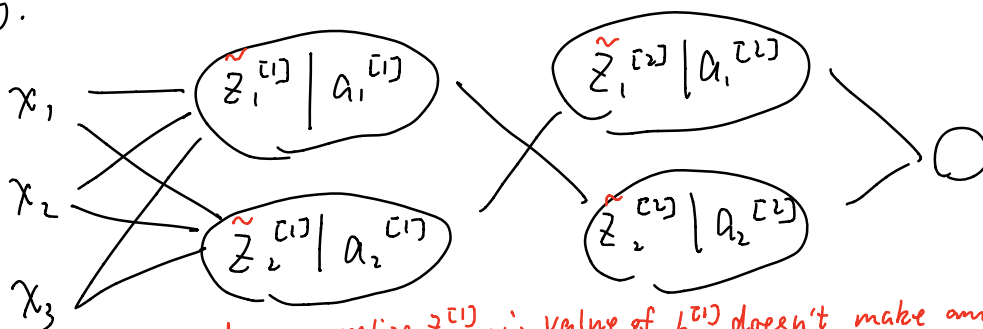
∴ 我们可以将 \tilde{z} 理解为将 z 调整 mean 和 variance 之后得到的点

learnable parameter

intuition: control the mean and variance of $z^{(i)}$, so that its distribution will fit in the activation function better



eg.



∴ we need to normalize $z^{(i)}$, ∵ value of $b^{(i)}$ doesn't make any difference
∴ set it to zero

$$x \xrightarrow{w^{(i)}, b^{(i)}} z^{(i)} \xrightarrow{\text{normalize}} \tilde{z}^{(i)} \xrightarrow{\gamma^{(i)}, \beta^{(i)}} \tilde{z}^{(i)} \longrightarrow a^{(i)} = g^{(i)}(\tilde{z}^{(i)})$$

∴ we can tune β and γ just like other parameter

$$\beta^{[L]} = \beta^{[L]} - \alpha d\beta^{[L]}$$

eg.

apply mini-batch with batch normalization

for $t = 1, \dots$ # mini batches

forward prop on $X^{\{t\}}$

in hidden layer, use $\tilde{z}^{[l]}$ to replace $z^{[l]}$

backward prop:

$dw^{[l]}, d\beta^{[l]}, dz^{[l]}$

update parameter:

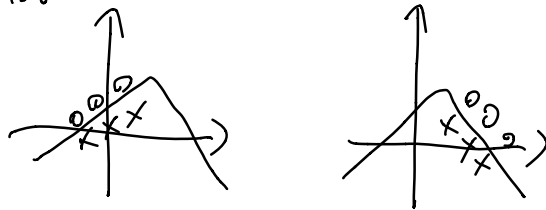
$w^{[l]} = w^{[l]} - \alpha dw^{[l]}$
 $\beta^{[l]} = \dots$
 $\gamma^{[l]} = \dots$

} can also try
RMS prop,
weighted momentum
or Adam

Why Batch Norm work in training?

1. force to control mean and variance of z ,
thus make it easier to train. (standardize hidden
layers)

反例:



hard to use the same model when distributions
are different.

2. because every mini-batch has dif mean & variance,
we batch-norm can add noise to training
(so that we can't rely just on

Thus it is similar to regularization (any nodes)

how to use batch norm at test time?

$$z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$\tilde{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta$$

$$\begin{array}{ccc} x^{\{1\}} & \dots & x^{\{T\}} \\ \downarrow & & \downarrow \\ \mu^{\{1\}}[L] & & \mu^{\{T\}}[L] \end{array} \xrightarrow{\text{exponentially moving average}} \mu = [\mu^{[1]} \dots \mu^{[L]}]$$

} training
(one batch at a time)

$$z_{\text{norm}}^{(i)[L]} = \frac{z^{(i)[L]} - \mu^{[L]}}{\sqrt{\sigma^{[L]2} + \epsilon}}, \quad \tilde{z}_{\text{norm}}^{(i)[L]} = \gamma z_{\text{norm}}^{(i)[L]} + \beta$$

} testing
(one instance at a time)

training时每个 mini-batch 单独根据各自的 mean 和 variance 校正,

在 testing 时,沿用之前得到的 γ 和 β ,

但由于是新数据,需要重新计算 $\hat{\mu}_{\text{test}}$ 和 $\hat{\sigma}_{\text{test}}^2$

① 若只有一个 sample, 则用之前各个 mini-batch 的 $\mu^{[1]\{1\}} \dots \mu^{[L]\{T\}}$ 计算 EWA 从而得到 $\hat{\mu}_{\text{test}}$, $\hat{\sigma}_{\text{test}}^2$ 同理

② 若有多个 samples, 直接 $\hat{\mu}_{\text{test}} \leftarrow \hat{\mu}_{\text{population}}, \hat{\sigma}_{\text{test}}^2 \leftarrow \hat{\sigma}_{\text{population}}^2$

softmax is a generalization of logistic regression to multiple classes

Multi-Classification: softmax

$$X \rightarrow \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \begin{matrix} a_1^{[L]} \\ a_2^{[L]} \\ a_3^{[L]} \\ \vdots \\ a_k^{[L]} \end{matrix} \rightarrow L(a^{[L]}, y)$$

k classes

$$\begin{bmatrix} 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 0 \end{bmatrix}$$

m points

eg. "softmax" with 4 classes

L denotes last layer

$$z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \rightarrow t^{[L]} = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix}$$

probability for each class

$$a^{[L]} = \frac{1}{\sum_{j=1}^4 t_j^{[L]}} t^{[L]}$$

"hardmax"

$$z^{[L]} \rightarrow \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

loss function of softmax

$$L(a^{[L]}, y) = \sum_{i=1}^m \left(\sum_{k=1}^K -y_k^{(m)} \log a_k^{[L](m)} \right)$$

↓

$-\log a_k^{[L]}$

backprop of softmax

$$dz^{(1)} = \hat{y} - y$$

Deep learning Framework

$$\begin{aligned} J(w) &= w^2 - 10w + 25 \\ &= (w - 5)^2 \end{aligned}$$

session.run (init)

session.run (train, feed_dict {x: coefficients})