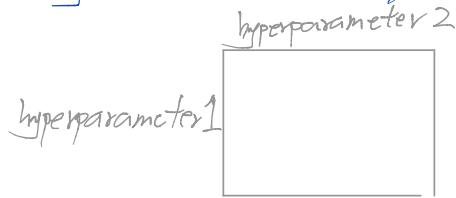


- o Hyper parameter tuning

- o Tuning Process

- △ Hyperparameters:
 - learning rate
 - α , B , β_1 , β_2 , ϵ , # layers, # hidden units
 - Learning rate decay, mini-batch size
 - Top importance
 - Second importance
 - 3rd importance

- △ Try random values = don't use grid



hyperparameter 1 & 2 may not have the same importance

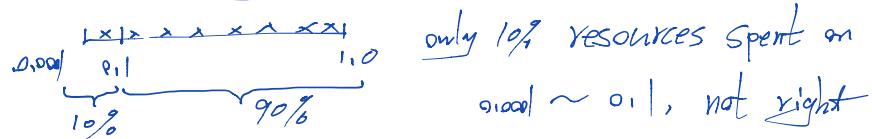
- △ Coarse to fine

Coarse search and then focus on fine search in small region.

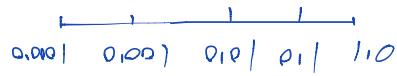
- o Using appropriate scale to pick hyperparameters

- △ E.g. $\alpha = 0.001, \dots, 1$

- if use uniformly chosen random α :



Instead, use log scale



$$\gamma = -4 + \text{np.random().randc()}$$

$$\alpha = 10^\gamma$$

Hyperparameters for exponentially weighted averages

$$\beta = 0.9, \dots, 0.999$$

$$v_i = \beta v_{i-1} + (1-\beta) \theta_i$$

$$1-\beta \rightarrow 0.1, \dots, 0.999$$

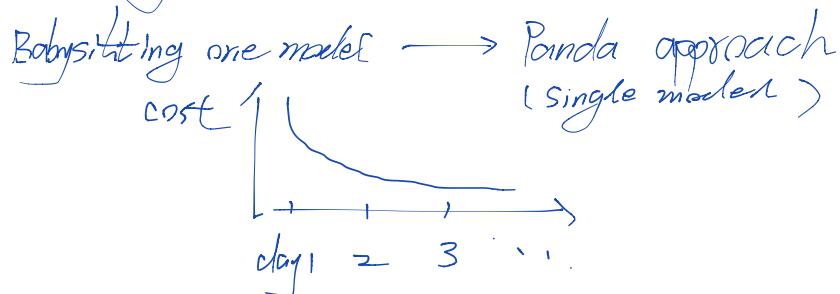
$$\gamma \in [-3, -1]$$

$$1-\beta = 10^\gamma, \quad \beta = 1-10^\gamma$$

$\frac{1}{1-\beta}$ is sensitive to small changes of β when β is close to one

Hyperparameter tuning in practice: Panda VS. Canvar

Model training:

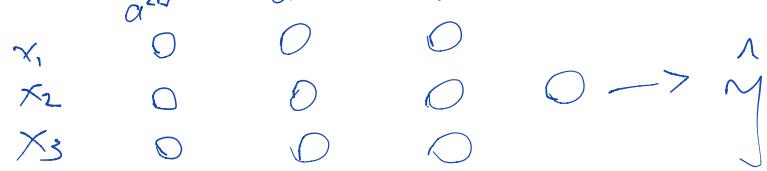


Training many models in parallel \rightarrow canvar approach
(multi-model)

o Batch Normalization

△ Normalizing activations in a network

Normalizing inputs to speed up learning



Can we normalize values of a so that can train
 w, b faster?



△ Implement batch norm

Given some intermediate value in NN: $z^{(1)}, \dots, z^{(m)}$
 $z^{[l] \in (1)}$

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

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

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

$$\tilde{z}^{(i)} = \gamma z_{norm}^{(i)} + \beta \quad \begin{matrix} \text{learnable parameters} \\ \text{of model} \end{matrix}$$

$$\text{if } \gamma = \sqrt{\sigma^2 + \epsilon} \quad \& \quad \beta = \mu$$

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

\tilde{z} makes hidden units have means & stds other than 0 & 1

Batch Normalization: Not just apply normalization to the input layer, but also deep into the hidden layer in the Neural Network

Fitting Batch Norm Into a NN

→ Adding batch norm to a network

$$x \xrightarrow{w^T b} z \xrightarrow[\text{Batch Norm} \\ (\mu, \sigma^2)} \tilde{z} \rightarrow a = g(\tilde{z})$$

Parameters: $w^{[0]}, b^{[0]}, \dots, w^{[L]}, b^{[L]}$ } SGD
 $\beta^{[0]}, \gamma^{[0]}, \dots, \beta^{[L]}, \gamma^{[L]}$ } mini-batch Adam

Implement batch norm in tensorflow:

ff. nn. batch - normalization

↳ Working with mini-batches

(a common usage in practice)

$$X = \sum_{i=1}^n a_i z^i$$

→ → → → → → → → → →

$$x^{[t]} \rightarrow z^{[t]} \xrightarrow[\text{B.N.}]{\beta^{[t]}, \gamma^{[t]}} \hat{z}^{[t]}$$

$\xrightarrow{\quad}$ within mini-batch

$$z^{[t]} = w^{[t]} a^{[t-1]} + b^{[t]}$$

$$\hat{z}^{[t]} = z_{\text{norm}} \cdot Y + f$$

* $b^{[t]}$ does influence z_{norm} since z_{norm} always has 0 mean.
thus we can remove $b^{[t]}$ term in $\hat{z}^{[t]} = w^{[t]} a^{[t-1]} + b^{[t]}$

$$\hookrightarrow z^{[t]} = w^{[t]} a^{[t-1]}$$

$$\hat{z}^{[t]} = \gamma^{[t]} z_{\text{norm}} + \beta^{[t]}$$

Implement Gradient Descent

for $t=1, \dots, \# \text{ of mini-batches}:$

Compute forward prop on $x^{[t]}$

In each hidden layer, use B.N. to replace $z^{[t]} \rightarrow \hat{z}^{[t]}$

Use backprop to compute $d\hat{z}^{[t]}, d\hat{b}^{[t]}, d\hat{\beta}^{[t]}, d\hat{\gamma}^{[t]}$

Update params $w^{[t]} := w^{[t]} - \alpha d\hat{w}^{[t]}$

$\beta^{[t]} := \beta^{[t]} - \alpha d\hat{\beta}^{[t]}$

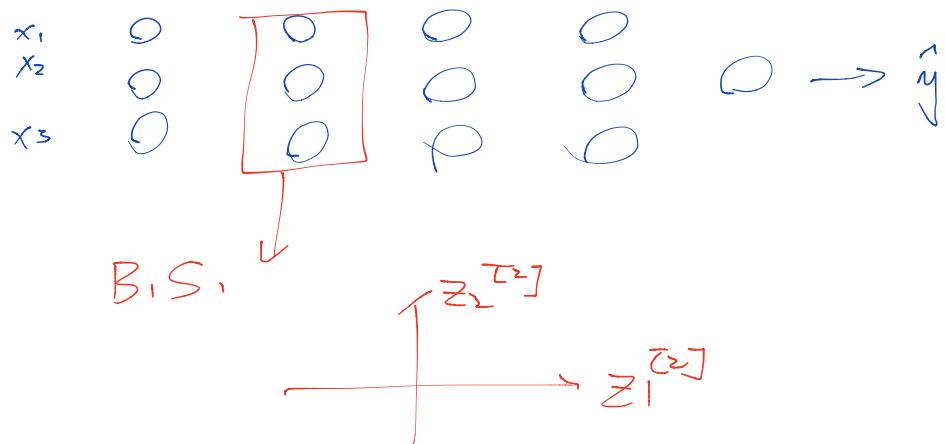
$\gamma^{[t]} := \gamma^{[t]} - \alpha d\hat{\gamma}^{[t]}$

(also works with GD w/ momentum, adam)

o Why does Batch Norm work?

△ Learning on shifting input distribution

Covariate shift: if the input distribution changes, then one may retrain the model



B.S. made sure no matter how input/early layers change, z_1^{L2} , z_2^{L2} will always have the same mean and variance

★ B.S. reduces influence of input value changing

o Batch Norm as regularization

- Each mini-batch is scaled by the mean/variance computed on just that mini-batch
- This adds some noise to the values $\tilde{z}^{[t]}$ within that mini-batch. So similar to dropout, it adds some noise to

each hidden layer's activations.

- This has a slight regularization effect

- You can turn B.N. as a regularization method

- Batch Norm at test time

- At test time, there may not be enough examples to take mean/variance to do Batch Norm
- Have to come up with some separate estimate of μ, σ^2
- μ, σ^2 : estimate using exponentially weighted average across mini-batch

e.g. $x^{(1)}[i]$ $x^{(2)}[i]$ $x^{(3)}[i]$ in training
 \downarrow \downarrow \downarrow
 $\mu^{(1)}[i]$ $\mu^{(2)}[i]$ $\mu^{(3)}[i]$ $\rightarrow \mu, \sigma^2$
 θ_1 θ_2 θ_3

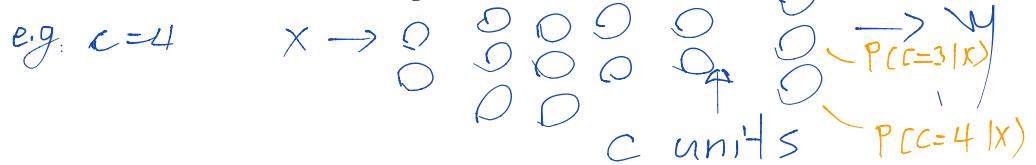
in testing, use μ, σ from above training result to do the B.S.

- Multi-class classification

- Softmax regression

$$C = \# \text{ of classes}$$

e.g. $C=4$



\hat{y} is (4×1) matrix

△ Softmax layer

$$z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]}$$

activation function =

$$t^{(i)} = e^{z^{[l]}(i)}$$

$$a^{[l]} = \frac{t^{(i)}}{\sum_{i=1}^n t^{(i)}}$$

e.g.: $z^{[l]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \rightarrow t = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 2.1 \end{bmatrix}$

$$\hookrightarrow \frac{1}{\sum_i} t^{(i)} = 16.3$$

$$a^{[l](1)} = 0.842 \quad a^{[l](2)} = 0.042$$

$$a^{[l](3)} = 0.002 \quad a^{[l](4)} = 0.114$$

$$\hookrightarrow \hat{y} = \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix}$$

Softmax + no hidden layers \rightarrow linear decision boundary classification

△ Training a Softmax Classifier

- Understanding Softmax

"hard max": $\begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$

"Softmax": $\begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \rightarrow \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix}$ "more gentle"

Softmax regression generalizes logistic regression to C classes. If $C=2$, softmax \Rightarrow logistic regression.

- loss function

$$L(y, \hat{y}) = - \sum_{j=1}^C y_j \log \hat{y}_j$$

- Cost

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

- $Y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}]_{C \times m}$

- Back prop. $\Delta z^{(i)} = \hat{y}^{(i)} - y^{(i)}$

- Deep Learning Frameworks

- Caffe/Caffe2
- Lasagne
- Theano
- CNTK
- mxnet
- Torch
- DL4J
- Paddle/Paddle
- TensorFlow
- Keras

- TensorFlow