

Charlie Mackie  
May 25th 2020

## Improving Deep Neural Networks : Hyperparameter tuning / Optimization

### Setting up ML app.

#### Data Split:

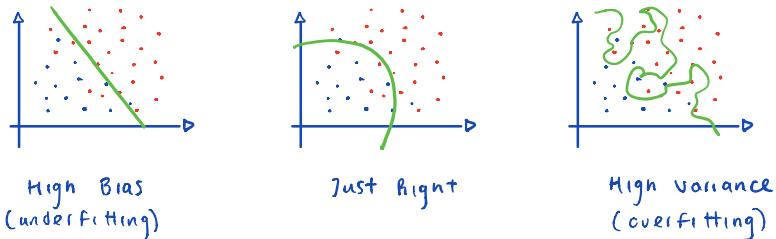
TRAIN	DEV	TEST
98 %	1 %	1 %

- ① will look different with less data, ex 100 examples total:  
70 train, 30 dev/test

#### Distributions:

Test and Dev. sets must come from the same distribution  
(otherwise accuracy could be skewed poorly)

#### Bias / Variance:



#### Data Set Accuracy Scenarios:

- ① Train Set Error: 1% ↗ high variance / overfit  
Test Set Error: 11% ↘ does not generalize well
- ② Train Set Error: 15% ↗ high bias / underfit  
Test Set Error: 16% ↘ error too far from Bayes

Human Error: How well would the most suitable human(s) perform  
Bayes Optimal Error: approx. 0%.

#### Solutions:

- ① High Bias → • Try a bigger network  
• Try a different architecture
- ② High Variance → • Get more data  
• Regularization

## Regularization

Goal: reduce overfitting / variance / accuracy difference

$$J(w, b) = \frac{1}{m} \sum_1^m L(\hat{y}, y) \quad \textcircled{*} \text{ add something to Cost function: this will change the gradient descent}$$

L2 Regularization:  $J(w, b) = \frac{1}{m} \sum_1^m L(\hat{y}, y) + \frac{\lambda}{2m} \|w\|_2^2$

L1 Regularization:  $J(w, b) = \frac{1}{m} \sum_1^m L(\hat{y}, y) + \frac{\lambda}{2m} \|w\|_1$

\textcircled{\*} now have  $\lambda$  as a hyperparameter for Dev. set tuning

In a Neural Network:

$$J(w, b) = \frac{1}{m} \sum_1^m L(\hat{y}, y) + \frac{\lambda}{2m} \sum_1^L \|w\|_F^2$$

$$\|w\|_F^2 = \sum_{i=1}^{n^{(L-1)}} \sum_{j=1}^{n^{(L)}} (w_{ij}^{[L]})^2 \quad \leftarrow \text{Frobenius Norm.}$$

→ This process will effect Gradient descent...

$$\begin{aligned} w^{[L]} &:= w^{[L]} - \alpha \left[ (dw)^{\text{FROM BP}} - \frac{\lambda}{m} w^{[L]} \right] \\ &= w^{[L]} - \frac{\alpha \lambda}{m} w^{[L]} - \alpha (dw)^{\text{FROM BP}} \end{aligned}$$

Why does this work?

- punishes weights with high value / variance
- pushes  $w^{[L]}$  to be smaller  $\therefore$  closer to a linear slope in activation function

## Dropout

Creates a probability parameter that represents if a node will be 'dropped' or not.

Creates more generalized / Lower variance models

Implementation: Inverted Dropout

this could vary by layer

④ construct vector  $d_3$  and multiply  $d_3$  by  $A_3$ , set Keep Prob parameter

- $d_3$  (dropout vector) = `np.random.rand(A3.shape[0], x) < Keep Prob`
- $a_3 = np.multiply(a_3, d_3)$
- $a_3 / \text{Keep Prob} \rightarrow$  need to retain general norm of weights

⑤ Do not apply Dropout @ Test Time

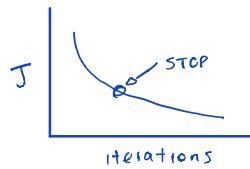
Why does Dropout work?

- Less reliance on a single feature (generalized)
- Smaller network from regularization

Other types of Regularization:

Data Augmentation → crop / rotate images to manipulate / increase the training set

Early Stopping →



Downsides of Some Regularization :

- can reduce orthogonalization (focus on individual tasks)
- want to focus on these separately
  - optimize
  - generalize

## Optimization

Normalize Inputs / Training Set

① subtract mean:  $\mu = \frac{1}{m} \sum x^{(i)}$ ,  $x - \mu$

② normalize variance:  $\sigma^2 = \frac{1}{m} \sum (x^{(i)})^2$

④ cost function is more symmetric and can step through with a higher learning rate

## Vanishing / Exploding Gradients

weights being manipulated millions of times can cause them to reach extrem values

↳ use a different initialization : Glorot / Xavier

$$W^L = np.random.randn(\text{shape}) * \sqrt(2 / n^{L-1})$$

Gradient Checking :  $\frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon}$

essentially random sample and check if the gradient you calculated ( $DW_L$ ) is within  $10^{-3}$  of above formula

## Optimization Algorithms (make your algorithm faster)

### Mini-batch Gradient Descent :

Train Data Set  $X = \{x^{(1)}, x^{(2)} \dots x^{(n)}\}$

$\underbrace{\qquad\qquad\qquad}_{\text{partition your data set and complete}} \rightarrow \text{initialize} \rightarrow F \text{ prop.} \rightarrow J \rightarrow B \text{ prop} \rightarrow \text{update}$   
for each partition.

This speeds up the training process by performing gradient descent more frequently

④ Entire pass through Training Set (call minibatches) = Epoch

choose mini-batch size between 1 - m:

- ① Batch Gradient descent : size = m
- ② Stochastic Gradient Descent (SGD) : size = 1 

Typically powers of 2 are used in correspondance with Comp. memory

⑤ If  $m \% \text{batch size} \neq 0$  then make the last batch smaller

## Gradient Descent with momentum

uses Exponentially weighted Averages: (similar to SMA)

$$V_t = \beta (V_{t-1}) + (1-\beta) \Delta t \quad \frac{1}{1-\beta} = \text{number of } x \text{ (time intervals) average is calculated over}$$

(\*) notice how it's recursive,  
this value accounted for  
all previous values in  $\beta$  range

For GD with mom., change the update step:

↳ calculate  $V_dW = \beta V_dW + (1-\beta) dW$   
 then  $W = W - \alpha V_dW$       ↑ exponentially weighted average

reduces noise:



(\*) Typically use  $\beta=0.9$ , smaller  $\beta$  = smaller average period

## RMS Prop (Reduces vertical movement / noise)

$$w := w - \alpha \frac{dw}{\sqrt{sdw}} \quad b := b - \alpha \frac{db}{\sqrt{sdw}}$$

(\*) where  $sdw = \beta sdw + (1-\beta) dw^2$   
 $sdw = \beta sdw + (1-\beta) dw^2 \leftarrow$  smaller and reduces oscillations in G.D.

(\*) purpose : use a higher learning rate = faster

## Adam Optimization: RMSProp + Momentum (Adaptive Movement Estimation)

$$\begin{aligned} V_dW &= \beta_1 (V_dW) + (1-\beta_1) dW \\ sdw &= \beta_2 (sdw) + (1-\beta_2) dW^2 \end{aligned} \quad \left. \right\} \text{Bias' are also corrected}$$

$$w := w - \alpha \frac{V_dW}{\sqrt{sdw} + \epsilon} \quad b := b - \alpha \frac{V_db}{\sqrt{sdw} + \epsilon}$$

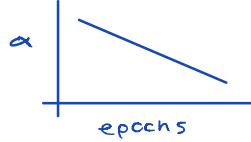
↑  
avoid NAN division

hyperparameters  $\rightarrow \epsilon = 10^{-8}, \alpha = LR, \beta = 0.9, \beta_2 = 0.99$

Learning Rate Decay — slowly reduce LR through Epochs

Intuition: take smaller steps when approaching convergence

$$\alpha = \frac{\alpha_0}{(1 + (\text{decay rate})(\text{epoch num}))}$$



### Hyperparameter Tuning

Things to tune:

- $\alpha$  (learning rate)
- hidden units
- layers
- $\beta_1$  and  $\beta_2$  (Adam)
- mini-batch size
- $\alpha$  decay

④ should tune these on a logarithmic scale

↳ can't just periodically iterate; missing lots of possible results



want to sample more densely at smaller/more precise values

### Batch Normalization (accelerates training)

has a regularizing effect also...

$$\mu = \frac{1}{m} \sum x^{(i)}, \sigma^2 = \frac{1}{m} \sum (x^{(i)} - \mu)^2$$

$$x_{\text{norm}} = \frac{x - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

expressed as:  $\tilde{z}_i = \gamma z_i + \beta$   
where  $\gamma$  and  $\beta$  are params. that we can learn.

④ this is performed within the context of each mini-batch

Benefits:

- input values are more stable while later layers train
- more independence between layers

## Softmax Regression (for multi-class classification)

output a vector of probabilities that sum to 1

$$\vec{z}_L = \begin{bmatrix} 1 \\ z \\ 3 \end{bmatrix} \quad a_L = \begin{bmatrix} e^1 / (e^1 + e^2 + e^3) \\ e^2 / \dots \end{bmatrix} \quad a_L = \frac{e^{(z)}}{\sum_i e^{(z_i)}}$$

could have a 'hard-max' where max is one-hot

$$\text{Loss function} = \mathcal{L}(\hat{y}, y) = -\sum y_j \log(\hat{y}_j)$$