

Train / Dev / Test Sets

Apply: Mh \Rightarrow # layers, # hidden unit, learning rate, ...	
Data	Test algorithm
training set 80%	Hold out cross-validation
development set	Leave-out test this right to day (left out dev set)
train / test \Rightarrow (70%, 30%)	"Dev" 20% (left out dev set)
more big data = less "dev set" and "test set"	

Mismatched train / test distribution

different picture quality, camera, type make two come from same distribution

- 1) Train pictures from webpage, Dev/Test sets from users

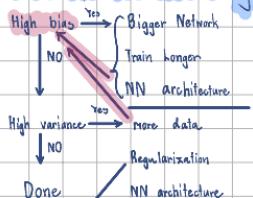
Bias / Variance

1) high bias \Rightarrow underfitting

2) high variance \Rightarrow overfitting



Basic recipe for machine learning



Regularization

$$J(w, b) = \frac{1}{2} \sum_{M=1}^n \|y^{(i)} - \hat{y}^{(i)}\|^2 + \frac{\lambda}{2M} \sum_{j=1}^n \|w_j\|^2 \quad (\text{optional})$$

$$l_2 \text{ regularization} \Rightarrow \sum_{j=1}^n \|w_j\|^2 = w^T w \quad |w: (n-1, n)^T|$$

$$l_1 \text{ regularization} \Rightarrow \sum_{j=1}^n \|w_j\|_1 = \sum_{j=1}^n |w_j|$$

$$\text{Frobenius norm} \Rightarrow \|w\|_F^2 = \sqrt{\sum_{j=1}^n \|w_j\|^2}$$

Weight decay \Rightarrow $w^{(t+1)} = w^{(t)} - \alpha \frac{\partial J}{\partial w} (\text{from backup}) + \frac{\lambda}{M} w^{(t)}$

$$w^{(t+1)} = (1 - \alpha \frac{\lambda}{M}) w^{(t)} - \alpha \frac{\partial J}{\partial w} (\text{from backup})$$

shrink a little each step before the normal gradient descent

Why Regularization Reduces Overfitting

$$J(w, b) = \frac{1}{2} \sum_{M=1}^n \|y^{(i)} - \hat{y}^{(i)}\|^2 + \frac{\lambda}{2M} \sum_{j=1}^n \|w_j\|^2$$

$$W \approx 0$$

$$\text{Example: } \tanh(g(x)) = \tanh(x)$$

$$g \uparrow \quad W \downarrow \quad x \uparrow \quad W = a x + b$$

Every layer \approx linear

Dropout regularization \Rightarrow it is hyperparameter. If chosen properly set some of them to zero proportion of neurons and dropout randomly deactivates

e.g. Illustrate with layer $l = 3$

that as chance

$$a_3 = np.random.rand(a3.shape[0], a3.shape[1]) < keep_prob$$

$$a_3 = np.multiply(a3, a3) \quad |a_3 == a3$$

$a_3 / = \text{keep_prob}$ Inverted dropout \Rightarrow mode expected value
 a_3 running the same

Understanding Dropout

set dropout to lower than another layer because that layer can be overfitting more than another layer
used in computer vision (big data input)

Other Regularization

1) Data augmentation \Rightarrow e.g. flip image in horizon axis, get more data for a little



edit test or train

2) Early stopping \Rightarrow 2 times to know where to stop



Optimize cost function J
e.g. Gradient Descent, etc.
Not overfit
e.g. Regularization, etc.

Normalizing Inputs

$$1. \text{ subtract mean}$$

$$x = x - \mu_x$$

$$2. \text{ Normalize var}$$

$$\sigma^2 = \frac{1}{M} \sum_{m=1}^M (x^{(m)} - \mu_x)^2$$

$$x = \sigma^{-1}(x - \mu_x)$$

$$x_1 \rightarrow x_1'$$

$$x_2 \rightarrow x_2'$$

$$x_3 \rightarrow x_3'$$

$$x_4 \rightarrow x_4'$$

Unnormalize

normalize \Rightarrow make cost function more symmetric
Laplace to local minima more easier

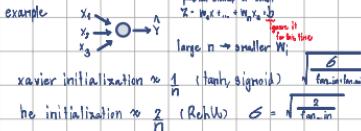
$$x = \sigma^{-1}(x - \mu_x)$$

$$x = \sigma^{-1}(x - \mu_x)$$

- Vanishing/exploding gradients or small (very rare) \Rightarrow big steps, make training difficult
- Vanishing \Rightarrow early layers learn extremely slowly.
 - Exploding \Rightarrow weights become stable (γ very big) \downarrow large layer
- $W^{(k)} \leftarrow$ Identity matrix e.g. 0.9^k (explore)
- $W^{(k)} \leftarrow$ Identity matrix e.g. 0.9^k or 0.01^k (vanish)

Weight initialization for deep networks

example

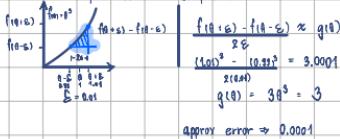


$$\text{xavier initialization} \approx \frac{1}{\sqrt{n}} (\tanh, \text{sigmoid})$$

$$\text{he initialization} \approx \frac{2}{\sqrt{n}} (\text{ReLU}) \quad \sigma = \sqrt{\frac{2}{n}}$$

Numerical Approximations of Gradients (use in gradient checking)

check your derivative computation



Gradient Checking $J(\theta) = J(\theta_0, \theta_1, \dots)$

take all of $W^{(1)}, b^{(1)}$ and reshape a big vector θ

$$\text{for each } i: d\theta_i \text{ approx } J(\theta_0, \theta_1, \dots, \theta_i + \epsilon, \dots) - J(\theta_0, \theta_1, \dots, \theta_i - \epsilon, \dots) / 2\epsilon$$

$$\approx d\theta_i \approx \frac{J(\theta_0, \theta_1, \dots, \theta_i + \epsilon, \dots) - J(\theta_0, \theta_1, \dots, \theta_i - \epsilon, \dots)}{2\epsilon}$$

$$\text{check } \frac{\|d\theta\|_2}{\|d\theta_{\text{approx}}\|_2} \approx 10^{-10} \text{ to } 10^{-8} \text{ Bad}$$

Gradient checking implementation notes

no mode training std

- don't use in training (only for debug).

- if algorithm fails grad check, look at components to identify bug.

look at $d\theta_i$ $d\theta_{i+1}$

- Remember regularization

- Doesn't work with dropout. Instead turn off dropout use diag check.

to double check algorithm is correct without it, and then turn on dropout

- Run at random initialization; perhaps again after some training

W, b close to 0 but sometimes bigger maybe implementation back prop long
is correct only when W, b close to 0, but it gives more accurate when W, b became

Mini Batch Gradient Descent (Hyperparameters)

Batch vs Mini-batch

Vectorization \Rightarrow allows you to efficiently compute

$$X = [x_1, x_2, \dots, x_M] \quad Y = [y_1, y_2, \dots, y_M]$$

$M = 5,000,000$ mini-batches of $1,000$ batch = $5,000$

$$\text{mini-batches} \Rightarrow X_j, Y_j$$

symbol $X^{(i)} \Rightarrow i$ training examples

$$X^{(i)} \Rightarrow z \text{ value at } i \text{ layer NN}$$

$$X^{(i)}, Y^{(i)} \Rightarrow \text{mini-batches}$$

Mini-batch gradient descent

for $t = 1, \dots, 5000$

1) Forward Prop on $X^{(t)}$

$$X^{(t)} = W^{(T)} A^{(T-1)} + b^{(T)}$$

$$A^{(T)} = g(X^{(T)})$$

⋮

$$2) \text{Compute Cost} \Rightarrow J = \frac{1}{1000} \sum_{i=1}^M (y^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot 1000} \|W^{(T)}\|^2$$

3) Backprop to compute gradients $J^{(T)}$ $\frac{\partial J}{\partial w^{(T)}} = (x^{(T)}, y^{(T)})$

$$4) W^{(T+1)} = W^{(T)} - \alpha \frac{\partial J}{\partial w^{(T)}}$$

$$b^{(T+1)} = b^{(T)} - \alpha \frac{\partial J}{\partial b^{(T)}}$$

1 epoch = one full pass through the training set

e.g. 60,000 examples, mini-batch size = 1,000

1 epoch = 60 steps

Understanding Mini-Batch Gradient Descent



Choosing your mini-batch size

$$\text{size} = M / (X^{(T)}, Y^{(T)}) = (X, Y) \text{ too long per iteration}$$

$$\text{size} = 1 \text{ every example is mini-batch base speed for vectorize}$$

In practice value in between 1 - M

small training set ($M < 2000$) \Rightarrow use batch

Typical mini-batch size = 64, 128, 256, 512

$\frac{1}{2} \text{ mode run fast}$

make sure fit in available memory

→ to know it's learn this before
optimization algorithm faster than gradient descent

Exponentially weighted averages

$$V_t = \beta(V_{t-1}) + (1-\beta)A_t$$

e.g. $\beta=0.3 \approx 10$ days temp

$\beta=0.3 \approx 50$ days temp

$\beta=0.5 \approx 2$ days temp

V_t as approximately average over

$$\alpha \frac{1}{1-\beta} \text{ day temp}$$

Understanding exponentially weighted averages

example

$$V_{100} = 0.3 V_{99} + 0.1 A_{100} = 0.10 \frac{1}{100} + 0.9 V_{99}$$

$$= 0.3 V_{99} + 0.3 A_{99}$$

$$= 0.1 A_{100} + 0.3(0.1 A_{100} + 0.9(0.1 A_{99} + 0.9 V_{99}))$$

$$= 0.1 A_{100} + 0.9(0.1 A_{100} + 0.9(0.1 A_{99} + 0.9^2 V_{99}))$$

$$= 0.1 A_{100} + (0.1)(0.3) A_{99} + (0.1)(0.9)^2 A_{99} + (0.9)^3 V_{99}$$



in contrast $0.98^{10} \approx \frac{1}{e}$

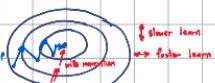
Bias correction of Exponentially Weighted Averages

$$\frac{V_t}{1-\beta^t} = \frac{1}{2}; \quad 1 - \beta^t = 1 - (0.9)^t = 0.036$$

$$\frac{V_t}{0.036} = \frac{0.036 A_t + 0.9 V_{t-1}}{0.036}$$

Gradient Descent with momentum

Example:



on iteration t :

Compute dW, dB on current mini-batch
from yellow acceleration

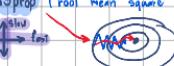
$V_{dW} = \beta V_{dW} + (1-\beta)dW$ imagine as bowl down from hill

$V_{dB} = \beta V_{dB} + (1-\beta)dB$

$W = W - \alpha V_{dW}$

$b = b - \alpha V_{dB}$

RMSprop (root mean square prop)
also known as gradient descent



On iteration t :

Compute dW, dB on gradient mini-batch

$$S_{dW} = \beta^t S_{dW} + (1-\beta) dW^2 \rightarrow \text{slowly decaying operations}$$

$$S_{dB} = \beta^t S_{dB} + (1-\beta) dB^2 \rightarrow \text{longer decaying operations}$$

$$W = W - \alpha \frac{dW}{\sqrt{S_{dW}}} \rightarrow \text{use small learning rate}$$

$$b = b - \alpha \frac{dB}{\sqrt{S_{dB}}} \rightarrow \text{use very small learning rate}$$

Adam Optimization Algorithm (just taking 'momentum'+'RMSprop')

$$V_{dW} = 0, V_{dB} = 0, V_{ab} = 0, V_{ab} = 0$$

on iteration t :

Compute dW, dB on gradient mini-batch

$$\text{momentum } V_{dW} = \beta_1 V_{dW} + (1-\beta_1) dW; V_{ab} = \beta_1 V_{ab} + (1-\beta_1) dB$$

$$\text{RMSprop } S_{dW} = \beta_2 S_{dW} + (1-\beta_2) dW^2; S_{ab} = \beta_2 S_{ab} + (1-\beta_2) dB^2$$

$$V_{dW} = \frac{V_{dW}}{1-(\beta_1)^t}; V_{ab} = \frac{V_{ab}}{1-(\beta_1)^t}$$

$$S_{dW} = \frac{S_{dW}}{1-(\beta_2)^t}; S_{ab} = \frac{S_{ab}}{1-(\beta_2)^t}$$

$$W = W - \alpha V_{dW} \rightarrow \text{stuck to local minima, bad reduction in gradient norm}$$

$$b = b - \alpha \frac{V_{ab}}{\sqrt{S_{ab}}} \rightarrow \text{more it easier to reach local minima}$$

hyperparameters

learning rate decay (reduces learning rate over time)

1 epoch = 1 pass through data

Epoch

1	0.1
2	0.07
3	0.05
4	0.04

other methods

$$\alpha = \frac{0.95}{k} \cdot \alpha_0$$

$$\alpha = \frac{1}{N_{epoch}} \cdot \alpha_0 \text{ or } \frac{k}{N_{epoch}} \cdot \alpha_0$$

or just manual decay

Tuning Process

Hyperparameters

α most important

$$\beta_1 \approx 0.9$$

$$\beta_2 \approx 0.999$$

$$\beta_1, \beta_2, \epsilon$$

layers

hidden units

learning rate decay

mini-batch size

The problem of local optima



Tensorflow

Multiple problem

cost function $w=5$
 $J(w) = w^2 - 10w + 25 = (w-5)^2$

Code example

```
import numpy as np
import tensorflow as tf
coefficients = np.array([1.0, 1.0-20j, 25j])
w = tf.Variable([0], dtype=tf.float32)
x = tf.placeholder(tf.float32, [3, 1])
cost = x[0][0]*w**2 + x[1][0]*w + x[2][0] # (w-5)**2
train = tf.train.GradientDescentOptimizer(0.01).minimize(cost)
init = tf.global_variables_initializer()
session = tf.Session()
with tf.Session() as session:
    session.run(init)
    print(session.run(w))
for i in range(1000):
    session.run(train, feed_dict={x: coefficients})
print(session.run(w))
```



Andrew Ng