

Course 2: Improving Deep Neural Networks :

Hyperparameter Tuning, Regularization and
Optimization

Content :

Week 1: Practical aspects of Deep learning

Week 2: Optimization algorithms

Week 3: Hyperparameter tuning, Batch normalization
and programming Frameworks

Coursera

Deep learning specialization

Train / Dev / Test sets

① dataset size : $100 \approx 1000000 \rightarrow 60/20/20$
 ↓
 small data

② size : $1000000 \approx \infty \rightarrow 98/1/1$
 ↓
 large data
 $\underline{=}$ $99.5/0.25/0.25$

Bias / variance

	overfitting	underfitting	Underfit	overfit	
error	high var	high bias	high bias, high var	low bias, low var	
Train	1%	15%	15%	0.5%	
Dev	11%	16%	30%	1%	

" high bias \rightarrow more training WON'T help !



get bigger network,
more data

assumption:

human has 0% error

Best

if high bias:

- // Try to make m bigger
- // Try different model
- // Train longer
- // Try different optimization algorithms

if high var:

- // more data
 - // Try regularization
 - // Try different model
- * Train bigger network never hurts

Reg for logistic
Reg. D

Regularization → reduce variance (overfitting)

$$J(w, b) = \frac{1}{m} \sum_i L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|w\|_2^2$$

(L2) : most common type $\|w\|_2^2 = w^T w$

(L1) : sparse! $\|w\|_1 = \sum_{j=1}^{n_x} |w_j|$

Req. for N.N.

$$J(w^{[l]}, b^{[l]}, \dots, w^{[L]}, b^{[L]}) = \frac{1}{m} \sum_i^L L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|w^{[L]}\|_F^2$$
$$\|w^{[L]}\|_F^2 = \sum_{i=1}^L \sum_{j=1}^{n^{[L-1]}} (w_{i,j}^{[L]})^2$$

"Frobenius" norm.
Not L2 norm.

why Regularization?

as $\lambda \rightarrow \infty$, weight matrix $\rightarrow 0$.



zero out a lot of hidden units

↓
simplify nn.

use DNN in training!

↑ if use in test, it
would add noise.

↓
reduce overfitting.

Dropout regularization → apply both during forward &
backward propagation

→ "drop" some neurons on each iteration based on
a probability. e.g. on each iteration, work with a
smaller neural network.

* spread out weights

Optimization algorithms

notation: i → i th training sample

$[l]$ → l th layer

$\{i\}$ → t th mini batch

Batch gradient descent: run gradient descent on whole data

Mini batch gradient descent: run gradient descent on mini data

↙ of small training set < 2000

= m → Batch gradient descent → cost func decreases on each iter.

Too long / iter

= 1 → Stochastic gradient descent → cost func oscillates, can

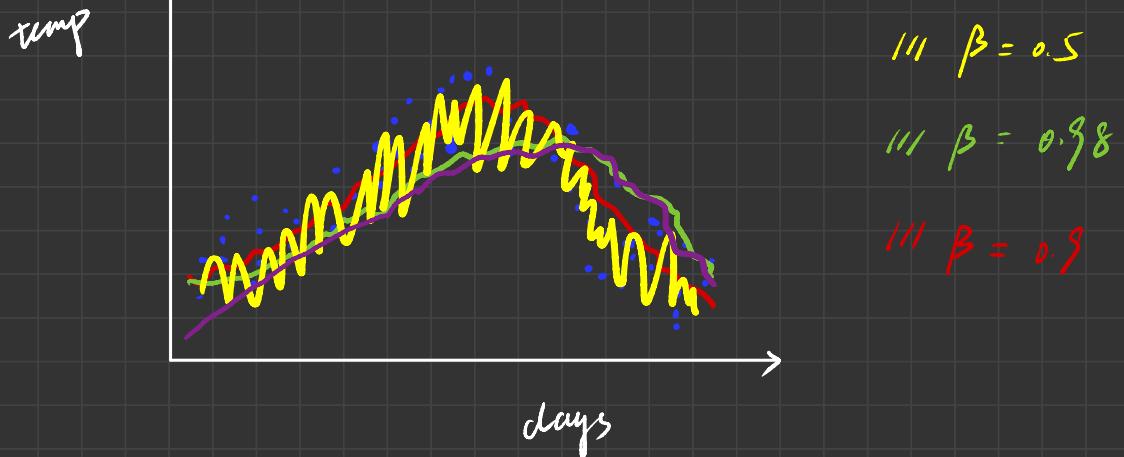
use smaller learning rate be extremely noisy, won't converge, lose
speedup from vectorization

batch $1 \times m$ → mini batch gradient descent → fast, don't always

has to be ↑ converge

power of 2
X GPU / GPU dev.

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



Note: when $\beta = 0.98$. Won't actually have green curve (see purple curve instead). A lower starting.

Bias correction

$$V_1 = 0.02 \theta_1 \rightarrow 0.02 \theta_1 / (1 - 0.98) = \theta_1$$

$$V_2 = 0.0196 \theta_1 + 0.02 \theta_2 \rightarrow (0.0196 \theta_1 + 0.02 \theta_2) / (1 - 0.98^2)$$

Gradient descent with momentum

on iteration t : β often used

compute d_w , db on current mini batch

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

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

$$w = w - \alpha V_{dw}, \quad b = b - \alpha V_{db}$$

RMSprop

computes the exponentially weighted averages of the squared gradients

Adam optimization algorithm → Adaptive moment estimation

Putting momentum & RMSprop together & combines the effect of gradient descent with momentum together with gradient descent with RMSprop.

CODE see lab

hyper parameter	:	α . (learning rate)		tune
β_1				0.9
β_2				0.999
ϵ				10^{-8}

$$\alpha = \left(1 / (1 + \text{decay-rate} * \text{epoch_num}) \right) \alpha_0 \leftarrow *$$

$$\alpha = 0.95 \text{ epoch_num} \quad \alpha_0$$

$$\alpha = K / \text{epoch_num} \quad \alpha_0$$

Tuning

1. d .
2. β . mini batch size, # of hidden units
3. # of layers, α decay, Adam β_1 , β_2 , ϵ .

choose points at random.

in practice: Pandas vs. Caviar

not enough computational capacity: baby sit one model
→ training models in parallel

Batch Normalization

softmax regression

TensorFlow

see labs