

Week 1

Monday, October 16, 2017 10:29 AM

Setting Up

Monday, October 16, 2017

10:30 AM

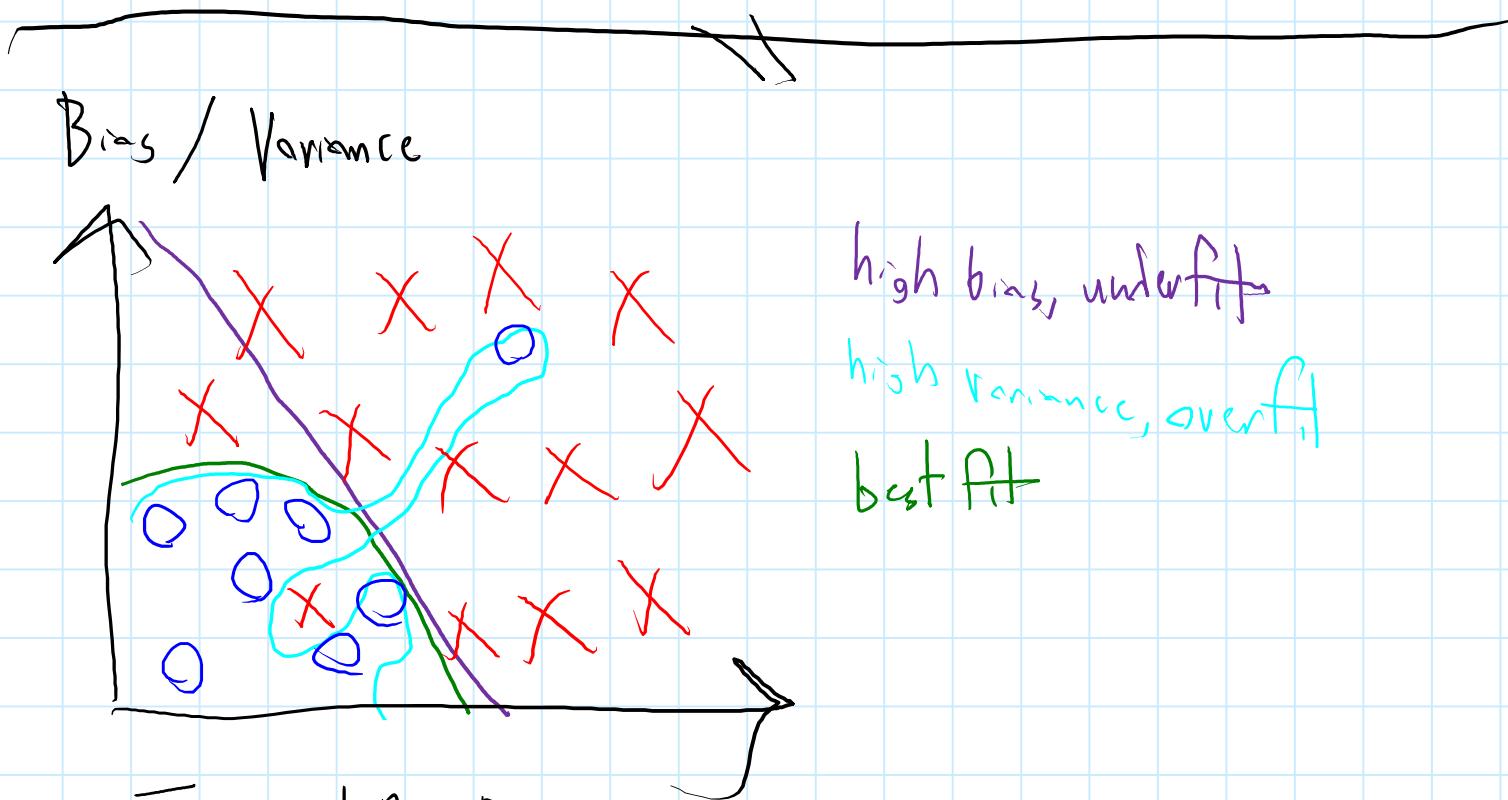
$T_{\text{train}} / \text{Dev} / \text{Test}$ sets

Split data into 3 groups

Train - train NN on bulk of data

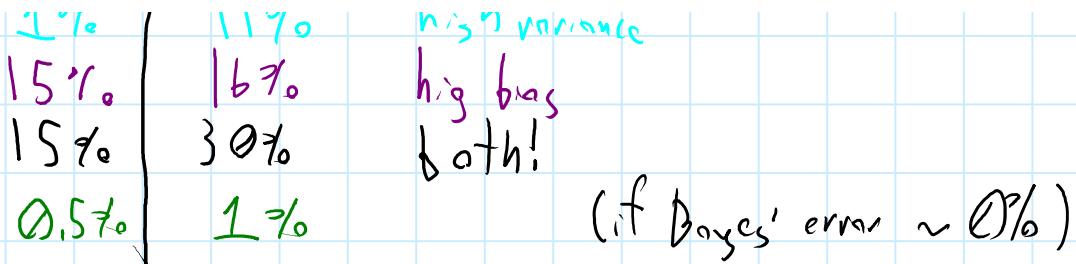
Dev - cross validate multiple training styles to find best

Test - Report test results

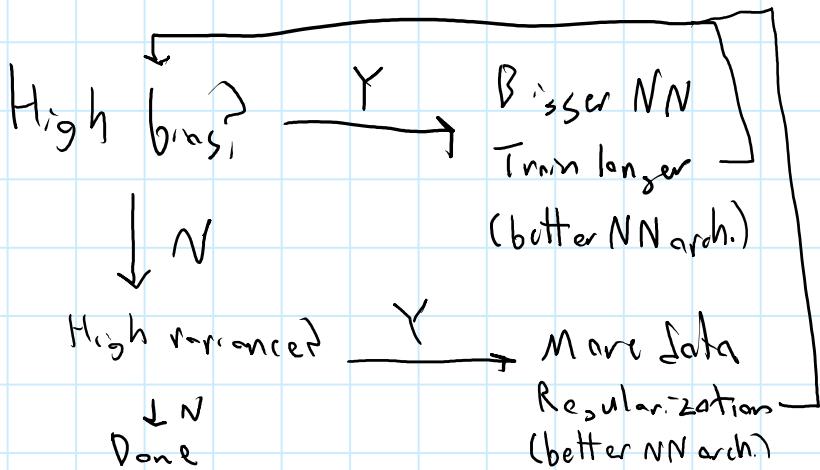


Train Err	Dev Err
1%	11%
15%	16%

high variance
high bias



Basic recipe for ML



Regularization

Monday, October 16, 2017 10:56 AM

Logistic Regression

$$\min_{w,b} J(w,b)$$

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

regularization parameter

L_2 regularization

$$\|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w \quad \text{Used much more}$$

L_1 regularization

$$\frac{\lambda}{2m} \sum_{i=1}^m \|w\|_1 = \frac{\lambda}{2m} \|w\|_1 \quad \text{"sparse" (lots of 0's)}$$

Neural Network

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

$$\|w^{[l]}\|^2 = \sum_{i=1}^n \sum_{j=1}^n (w_{ij}^{[l]})^2 \quad \text{"Frobenius norm" } \|w\|_F^2$$

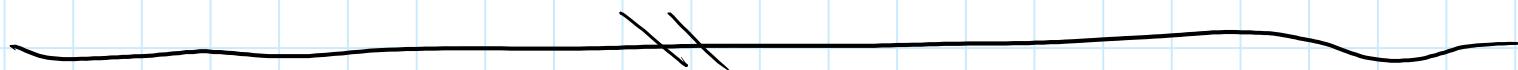
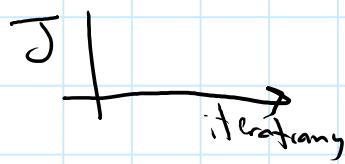
$$\begin{aligned} d_w^{[l]} &= \left(\text{Same from backprop} \right) + \frac{\lambda}{m} w^{[l]} \\ w^{[l]} &= w^{[l]} - \alpha (d_w^{[l]} + \frac{\lambda}{m} w^{[l]}) \end{aligned}$$

Why does regularization help with overfitting?

λ penalizes weights (w) for being too large, reduces impact of hidden units

$\lambda \uparrow \rightarrow w \downarrow \rightarrow z \downarrow$, makes layer more linear
Good de-biasing method!

$x \rightarrow w \downarrow \rightarrow z \downarrow$, more layer more linear
Good debugging method!



Dropout Regularization
Each node has chance of being eliminated

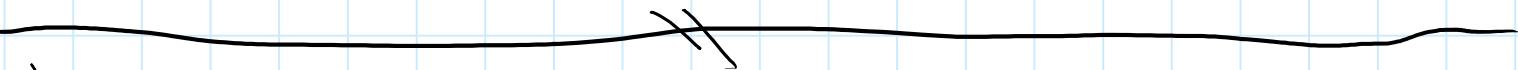
Ex. constraints to drop node in each layer

$$\begin{matrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ y_1 & y_2 & y_3 \end{matrix}$$

Implementing, illustrate w/ layer 3 (inverted dropout)

```
d3 = np.random.random(a3).shape[0], a3.shape[1]) * keep_prob  
a3 = np.multiply(a3, d3)  
a3 /= keep_prob
```

At test time, dropout won't be implemented



Why does dropout work?

Intuition: can't rely on any one feature, spread out weights

i.e. - n 0 1 .

i.e.



if any input disappears,
none favored

(can vary keep-prob by layer (lower w/ more nodes, higher w/ less))

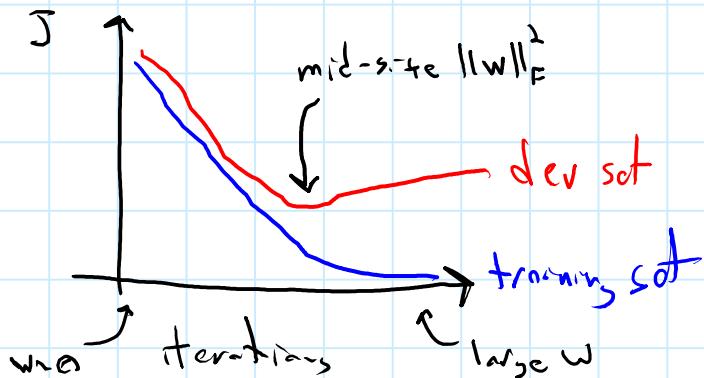


Other techniques

Augmentation

Flip images horizontally, random zooms, crops, rotations, distortions,

early stopping, (alt to L_2)

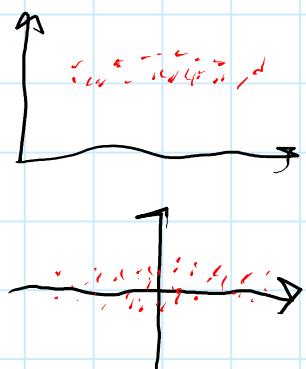


Setting Up Optimization

Wednesday, October 18, 2017 10:28 AM

Normalizing Training Sets

Subtract mean: $\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$, $X := x - \mu e$



Normalize variance: $\sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)^2$ (elementwise square)

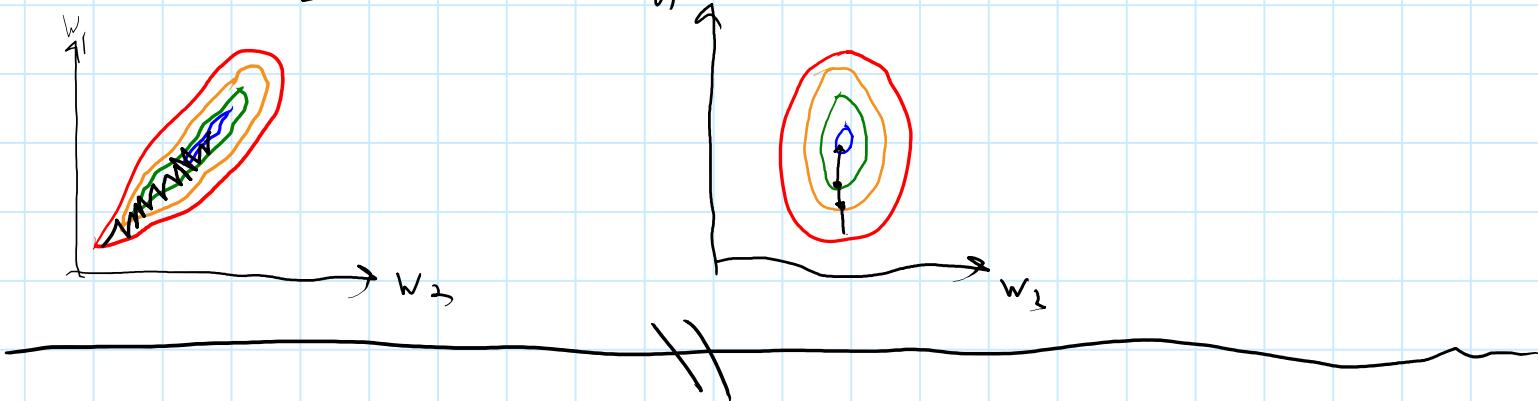


Make sure to use same σ, μ for test set

Why? Faster gradient descent

Unnormalize

Normalized



Vanishing / Exploding Gradients

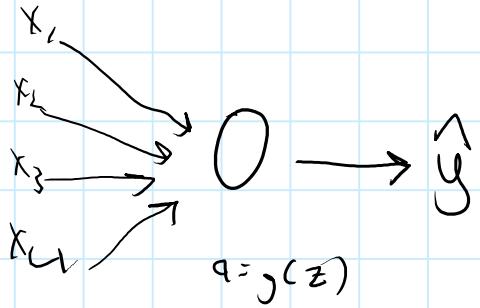
Given $w^{[L]}$, $g(z) = z$, $b^{[L+1]} = \emptyset$, $w^{[L]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$

$\hat{y} = \left(\prod_{i=1}^L w^{[i]} \right) x \rightarrow \hat{y} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}^{L-1} x \rightarrow$ Causes exponential change

Partial remedy:

Weight Initialization for Deep Networks

Single Neuron Example



$$z = \sum_{i=1}^n w_i x_i$$

larger $n \rightarrow$ want smaller w_i , variance $w_i \approx \frac{1}{n}$

$$w^{[l]} \approx \sqrt{\frac{c}{n^{[l-1]}}}$$

ReLU : $c=1$
tanh : $c=1$ (Xavier init)

Numerical Approx of Gradients

Homework Notes

Saturday, October 21, 2017 11:13 AM

Week 2

Thursday, October 26, 2017 10:29 AM

Notations:

~~X~~ [layer] {mini-batches} (node / example)

Optimization Algorithms

Thursday, October 26, 2017 10:29 AM

Mini-batch Gradient Descent $\{t\}$

$$\text{batch: } X = [x^{(1)}, \dots, x^{(m)}]; Y = [y^{(1)}, \dots, y^{(m)}]$$

mini batch $x^{\{1\}}, x^{\{2\}}, \dots, x^{\{m\}}$; $y^{\{1\}}, y^{\{2\}}, \dots, y^{\{m\}}$

Split into smaller chunks & learn/update sequentially

ex. $t \in \{1, \dots, 5000\}$ $m = 50000$

for $t = [1, 5000]$

Forward pass

$$Z^{[1]} = W^{[1]} X^{\{t\}} + b^{[1]}$$

$$A^{[1]} = g^{[1]}(Z^{[1]})$$

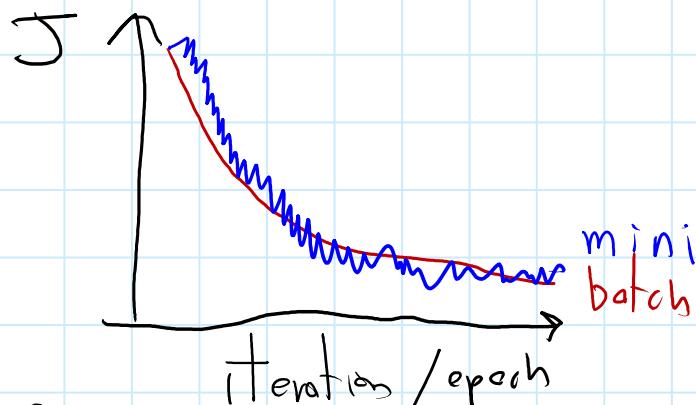
Compute cost $J^{\{t\}} = \frac{1}{1000} \sum_{i=1}^m \ell(y^{(i)}, A^{[1]}) + \frac{\lambda}{2 \cdot 1000} \sum_k \|W^{[k]}\|_F^2$

Backprop to get gradients
 $W^{[k]} = \alpha dW^{[k]}, b^{[k]} = \alpha db^{[k]}$

All this constitutes 1 "epoch" (pass through training set)

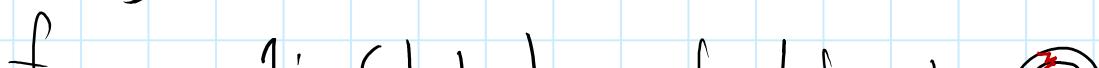
II

Understanding Minibatch



oscillates since smaller sets of training data trained

Choosing mini batch size

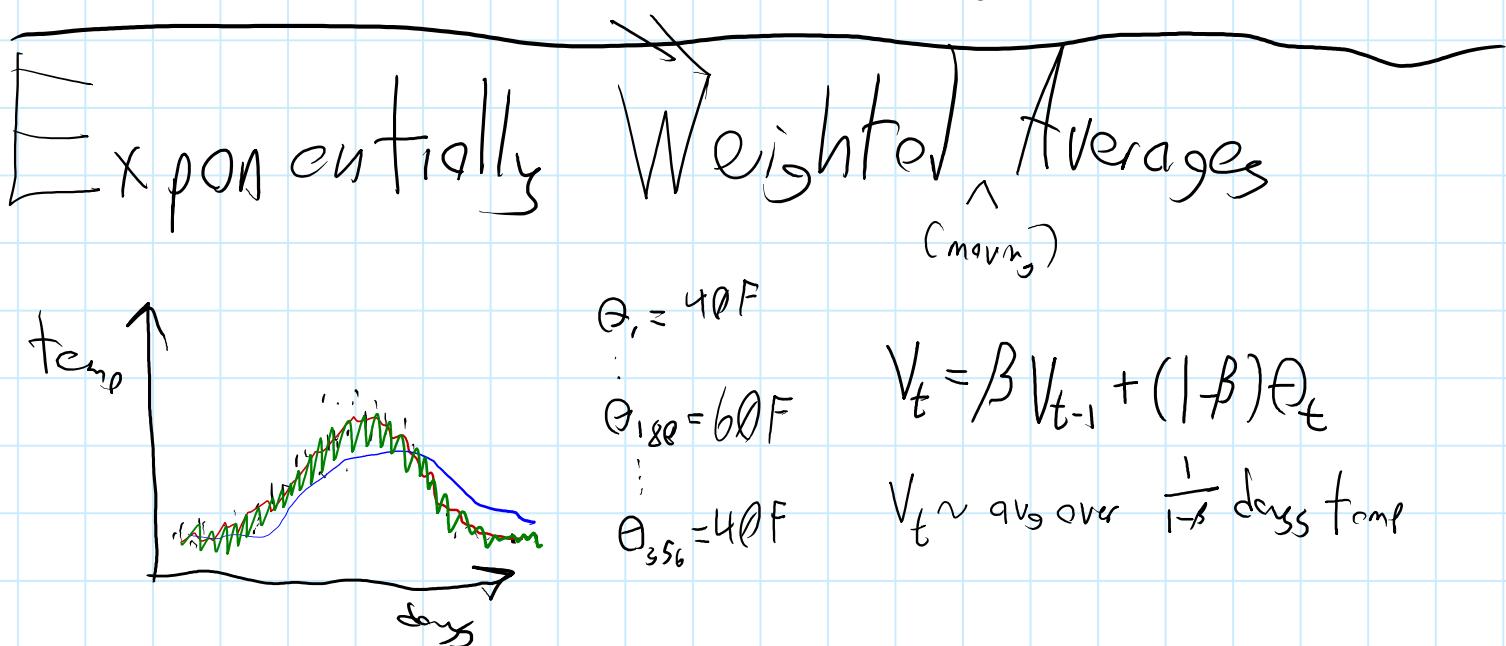


if size = 1: Stochastic gradient descent

→ Less speed gain from vectorization

Small training set (≤ 2500): Just use batch $\nabla \rightarrow$

Typical mini-batch size: 2^n (64, 128, 256, 512)
 - make sure it fits in (C/GPU) memory!



$$\beta = 0.9 \sim 10 \text{ days long}$$

$$\beta = 0.98 \sim 50 \text{ days long (adopts slower)}$$

$$\beta = 0.5 \sim 2 \text{ days (adopts rapidly, very sensitive)}$$

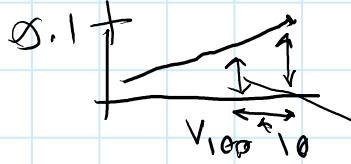
||

Understanding

$$\beta = 0.9$$

$$V_{100} = 0.1\theta_{100} + 0.9(0.1\theta_{99} + 0.9(0.1\theta_{98} + 0.9(\dots)))$$

$$= 0.1\theta_{100} + 0.1 \times 0.9 \cdot \theta_{99} + 0.1 \cdot (0.9)^2 \theta_{98} + 0.1 \cdot (0.9)^3 \theta_{97} + \dots$$



$$[0.9^{10} \sim \frac{1}{e}]$$

$$v_0 = \emptyset$$

repeat for all θ_i {

$$v_t := \beta v_{t-1} + (1-\beta) \theta_i$$

// Bias Correction

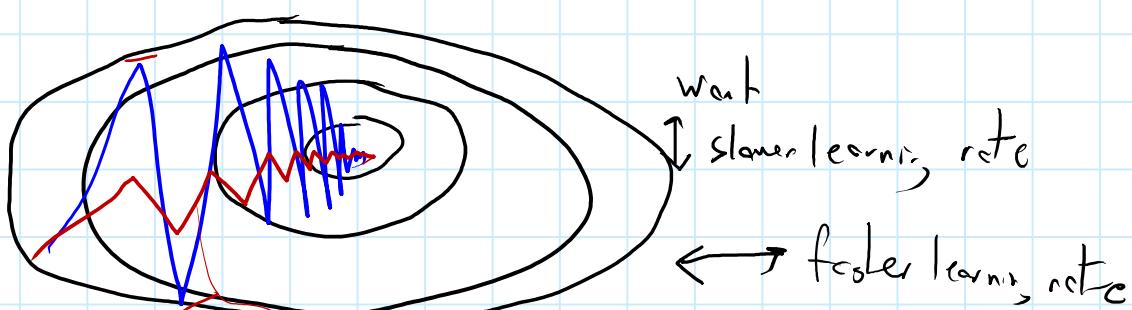
First handful of values will be abnormally low

$$\text{Correction: } \frac{v_t}{1-\beta^t}$$

$$v_p = \emptyset$$

$$v_t = 0.18 v_0 + 0.03 \theta_i$$

Gradient Descent w/ Momentum



avg of oscillations much closer to 0

On iteration t'

$$V_d w = \beta V_d w + (1-\beta) d w$$

$$V_d b = \beta V_d b + (1-\beta) d b$$

rolling ball downhill analogy

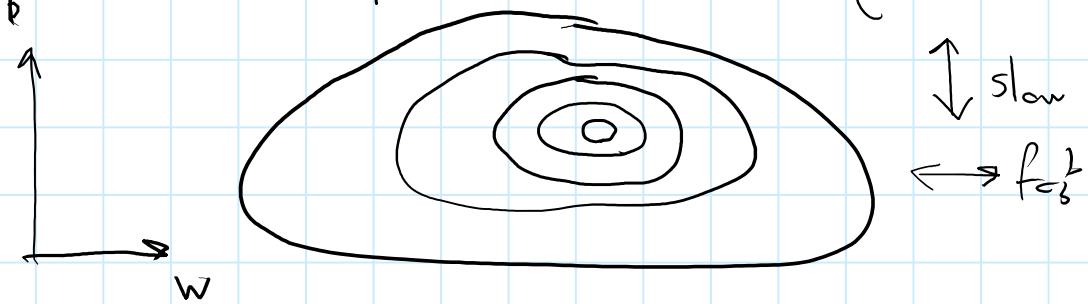
friction velocity "acceleration"

$$w := w - \alpha V_d w; \quad b := b - \alpha V_d b$$

(most people don't bias correct)

(most people don't bias correct)

RMS Prop (Root Mean Squared)



On iteration t:

Compute ∇_w, ∇_b

$$\begin{aligned} S_{dw} &= \beta S_{dw} + (1-\beta) \nabla_w^2 \xrightarrow{\text{element wise}} \text{will be small} \\ S_{db} &= \beta S_{db} + (1-\beta) \nabla_b^2 \xrightarrow{\text{will be large}} \end{aligned}$$

$$w := w - \alpha \frac{\nabla_w}{\sqrt{S_{dw}} (+\epsilon)} \quad b := b - \alpha \frac{\nabla_b}{\sqrt{S_{db}}} \quad \text{axis division by } \theta$$

Adam Optimization Algorithm (ADaptive Moment estimation)

$$V_{dw}, S_{dw}, V_{db}, S_{db} = \emptyset$$

On iter t:

Compute ∇_w, ∇_b

$$\begin{aligned} V_{dw} &= \beta_1 V_{dw} + (1-\beta_1) \nabla_w \quad V_{db} = \beta_1 V_{db} + (1-\beta_1) \nabla_b \\ S_{dw} &= \beta_2 S_{dw} + (1-\beta_2) \nabla_w^2 \quad S_{db} = \beta_2 S_{db} + (1-\beta_2) \nabla_b^2 \end{aligned}$$

$$V_{dw}^{\text{corrected}} = \sqrt{V_{dw}} / (1 - \beta_1 t), \dots$$
$$S_{dw}^{\text{corrected}} = \sqrt{S_{dw}} / (1 - \beta_2 t), \dots$$

$$w := w - \alpha \frac{V_{dw}^{\text{corr}}}{\sqrt{S_{dw}^{\text{corr}}} + \epsilon}, b := b - \alpha \frac{V_{db}^{\text{corr}}}{\sqrt{S_{db}^{\text{corr}}} + \epsilon}$$

Choice of hyperparams:

α : must be tuned

β_1 : 0.4; β_2 : 0.999

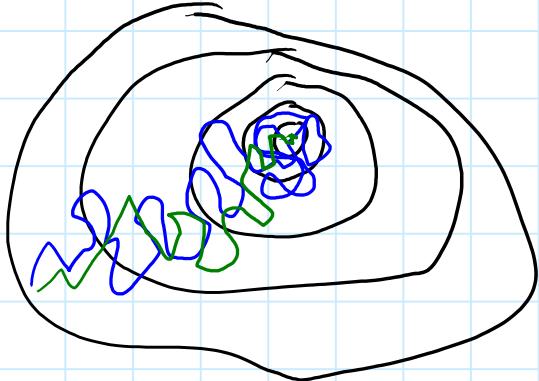
ϵ : 10^{-8}

Hyperparameter Tuning

Sunday, October 29, 2017

2:35 PM

Learning Rate Decay (Learn on Dr. Ng's priority)



lower α as training goes on

$$\alpha = \frac{\alpha_0}{1 + \text{decay_rate} * \text{epoch_num}}$$

$$\alpha = 0.95^{(\text{epoch})} \cdot \alpha_0, \sqrt{\frac{k \cdot \alpha_0}{\text{epoch}}}, \text{ discrete staircase!}$$



(can also control manually)

The Problem of Local Optima

While easy to imagine issue of local optima as "wells", due to high dimensionality, they're more like saddle points, real issue is plateaus causing learning to slow, but not stall (ADAM helps here!)

Homework Notes

Sunday, October 29, 2017 2:55 PM

Week 3

Tuesday, October 31, 2017 11:01 AM

Hyperparameter Tuning

Tuesday, October 31, 2017

11:01 AM

α

$\beta \sim 0.9$

$$\beta_1, \beta_2, \epsilon = [0.9, 0.999, 10^{-8}]$$

layers

hidden units

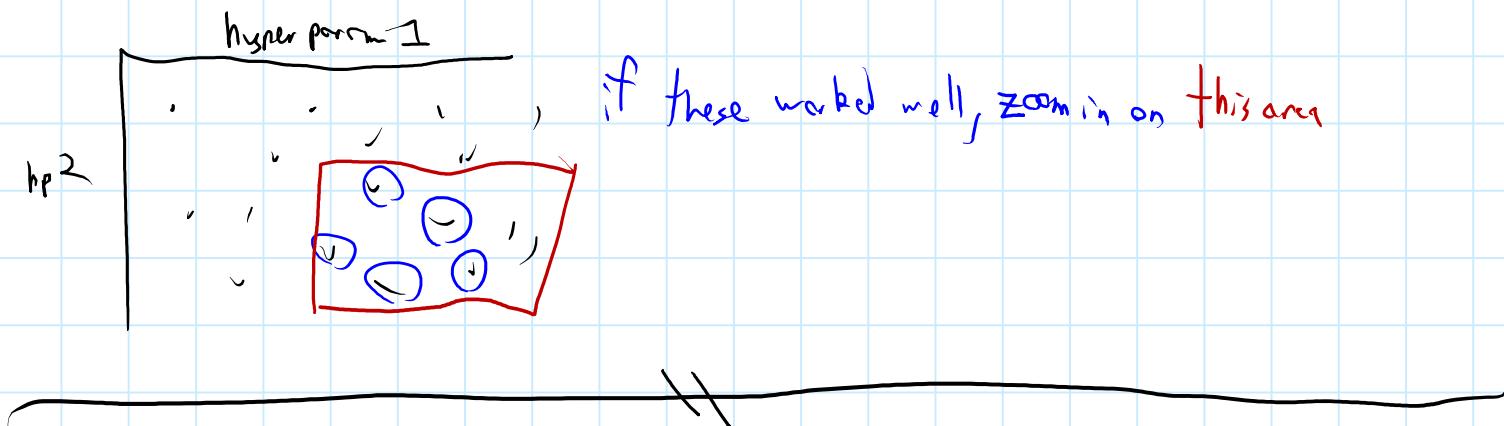
Learning rate decay

mini batch size

rough
importance
1, 2, 3

Try random values, don't use grid of values ($P^1 \times P^2$)

Use coarse to fine method



Picking randomly doesn't mean uniformly random!

e.g. $\alpha = 10^{-4} \dots 1$, sample along log scale

$r = -4 \Rightarrow np.random.rand() \leftarrow r \in [-4, 0]$

$\alpha = 10^{r}$

$$\beta = 0.9 \dots 0.999 ; (1-\beta) = 10^{-1} \dots 10^{-3}$$

implement same as α

Note $\lim_{\beta \rightarrow 1} \frac{\beta}{1-\beta} = \infty$



Organizing Hyperparameter searches
"panda vs corvo" parenting

panda': low resources, "babysit" one model & modify regularly, looking for changes that significantly improve J

corvo': many resources, train many models in parallel

Batch Normalization

Tuesday, October 31, 2017 10:22 PM

Normalizing features improves learning, the same applies to hidden layers!

Typically normalize $Z^{[l]}$, not $g^{[l]}$

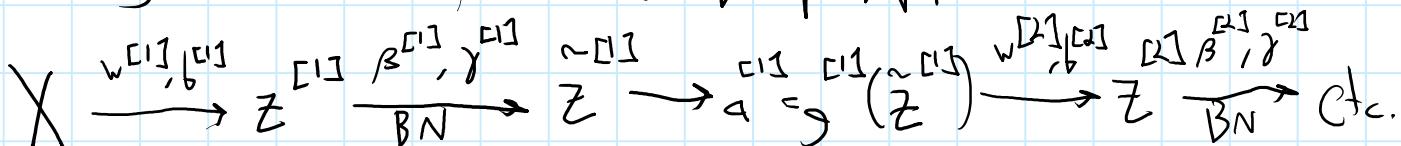
Ex, Given some values $Z^{[l]}$

$$\mu = \frac{1}{m} \sum_i^m Z^{[l]}_i$$
$$\sigma^2 = \frac{1}{m} \sum_i (Z_i - \mu)^2$$
$$Z_{\text{norm}}^{[l]} = \frac{(Z^{[l]} - \mu)}{\sqrt{\sigma^2 + \epsilon}}$$

$$\tilde{Z}^{[l]} = \gamma Z_{\text{norm}}^{[l]} + \beta$$

Don't always want $\mu=0$ & $\sigma=1$ for hidden units
learnable parameters

Using Batch Norm in a Deep NN



Typically used w/ mini-batches (Same as above but X^{batch} is input, \tilde{Z} is specific to batch)

Note : $Z^{[l]} = W^{[l]} g^{[l-1]} + b^{[l]}$ $b^{[l]}$ will be cancelled in normalization process due to mean subtraction, gets replaced by $\beta^{[l]}$

$$\tilde{Z}^{[l]} = \gamma^{[l]} Z^{[l]} + \beta^{[l]}$$

Implementing: (w/ $\nabla \downarrow$)

for $t=1 \dots \text{num_mini_batches}$

for $t = 1 \dots \text{num_mini_batches}$

Compute FP

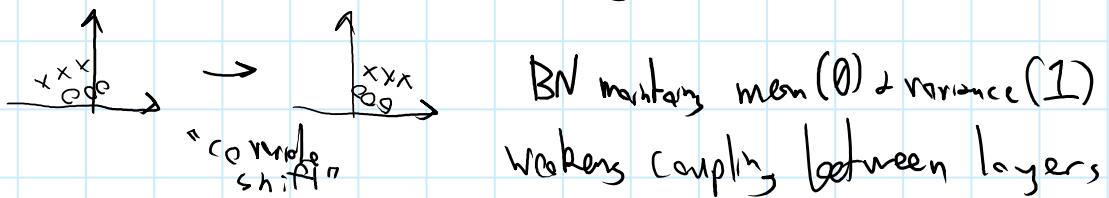
Use BN to replace $\mathbb{E}^{[l]}$ w/ $\hat{\mathbb{E}}^{[l]}$
Use BP to find $\frac{\partial w^{[l]}}{\partial b}, \frac{\partial \beta^{[l]}}{\partial b}, \frac{\partial \gamma^{[l]}}$

Update params

Why Does BN Work?

Learns on shifting input distribution

e.g. teaching cat classifier on only black cats may not recognize non-black



BN slight as regularization

Since each minibatch has its own μ/σ computed & normalized, some noise is introduced, forcing downstream neurons not to rely on one in particular (like dropout)
[Not a standin for actual regularization]

BN @ Test Time

BN @ test time

Not using mini batches during test

μ, σ^2 estimated w/ exp. weighted avg. across mini batches
i.e.

$$\begin{aligned} & X \{t\} \\ & \downarrow \\ & \mu^{\{t\}[x]} \quad \sigma^2 \{t\}[x] \quad \left(\text{combine over } t \right) \\ & Z_{\text{norm}} = \frac{Z - \mu}{\sqrt{\sigma^2 + \epsilon}} \end{aligned}$$

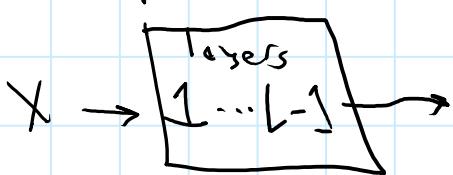
Multi-class Classification

Wednesday, November 1, 2017 10:18 AM

Softmax Regression

$C = \# \text{ classes trying to classify}, n^{[L]} = C$

ex. 4 possibilities



- $P(\text{none} | X)$
- $P(\text{class 1} | X)$
- $P(\text{class 2} | X)$
- $P(\text{class 3} | X)$

$$\rightarrow y$$

Softmax Layer ($Z^{[L]}$)

$$Z^{[L]} = W^{[L]} a^{[L-1]} + b^{[L]}$$

Activation function:
 $t = e^{\frac{z}{\sum t}}$, $a^{[L]} = \frac{t}{\sum t_i}$, $a_i^{[L]} = \frac{t_i}{\sum t_i}$

Ex. $Z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \rightarrow t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix}, \sum_{j=1}^4 t_j \approx 176.3 \rightarrow a = \frac{t}{176.3}$

Training a Softmax Classifier

Loss Function

$$-\ln \left[\frac{1}{n} \sum P(y_i) \right] \text{ for } i = 1, \dots, n$$

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \quad \hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix} \quad L(\hat{y}, y) = - \sum_{j=1}^J y_j \log(\hat{y}_j)$$

here, $L(\hat{y}, y) = -y_2 \log \hat{y}_2$, only way to rectify is to increase \hat{y}_2

"Maximum likelihood estimation" $J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i, y_i)$

"Gradient Descent w/ Softmax"

Backprop : $\underbrace{\frac{dL}{dz}}_{\frac{dJ}{dz}} = \hat{y} - y$