



deeplearning.ai

2 Improving deep neural networks

Andrew Ng

Note by:
Hossein Mamaghanian
amirhm@gmail.com



Improving NNs.

choice of hyper parameters are always highly iterative process
(# layers, learning rates, activation functions, ...)

Train : Model is trained on this

Dev set / hold out / Development set : verify the choice of hyper parameters.

Test : check performance on unseen test sets,

→ before 70/30% ratio was reasonable

but currently with millions of data normally smaller fraction
are dedicated for dev or test (98%, 1%, 1%) exp.

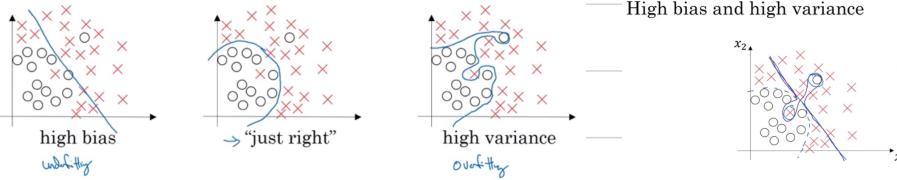
Mismatched train/test distributions,

→ Make sure dev /and test sets come from the same distribution,

→ Not having test set might be ok , (means only dev set could be fine). Just had to be careful not to overfit on the dev set.

→ Bias / Variance trade off

→ example



But in reality its not possible to plot the boundaries for decision

Train set error	10%	15%	15%	0.5%
Dev set error	110%	16%	30%	1%
	high variance	high bias	high bias	low bias
	(*)		≈ variance	low variance

→ This is with assumption that human error is 0% (at least there is solution)

or optimal (Bayes) error is zero. and Train, Dev one from same dist'n distribution,

→ if Bayes error is 15% for example, 2nd column then is not bad actually (Train error 15%, Dev 16%)

→ such case is for example for highly noisy, blurry data which even human cannot fully classify correctly.

Basic Recipe for machine learning:

① does it have high bias?
(training data) →

- | Bigger Net
- | Train longer
- | architecture search.

(at least we should be able to overfit, consider the Bayes error is 0).

② high variance?

→

- | More Data
- | Regularization to reduce overfit
- | architecture search

* Since before all the tools either were reducing Bias or Variance
(But damaging other) it used to be called Bias/Variance trade off.
But in Big Data era and dl, we have such tools
like having More Data will reduce Variance but will not affect
Bias, or Bigger Network (considering regularization is done) will
help ↓ Bias and not ↑ variance

Training Bigger Network never hurts (if regularization are done)
just more computation.

"Regularization"

→ for reducing overfitting.

Imagine the loss function,

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \gamma_2 m \|w\|_2^2$$

$\|w\|_2$ -norm regularization.

→ Most common $\|w\|_2$ -norm

$\|w\|_1$ -norm regularization $\gamma_1 m \sum_{i=1}^m |w_i| \rightarrow w$ will end up being sparse
(May be useful for compressing)

→ in reality this doesn't end up really compressing model,

λ : called regularization Parameter and we use dev train to find right value.

→ (another hyper parameter)

→ In NN

$$J(w^{(l)}, b^{(l)}, -w^{(l)}, b^{(l)}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \gamma_2 m \sum_{j=1}^L \|w_j^{(l)}\|_F^2$$

"Frobenius norm"

BP will change: $d_w^{(l)} = (\text{from BP}) + \gamma_m w^{(l)}$

$$w^{(l)} = w^{(l)} - \alpha d_w^{(l)}$$

→ weight updates becomes

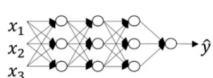
$$w^{[l]} \leftarrow w^{[l]} - \alpha \left[(\text{new}) + \gamma_m w^{[l]} \right]$$

$$\leftarrow w^{[l]} \left(1 - \frac{\alpha \gamma}{m} \right) - \alpha (\text{old})$$

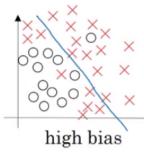
This means that our weights w is always shrinking, for this, ℓ_2 -norm also called "weight decay" regularization.

why regularization helps with overfitting?

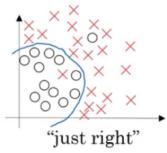
How does regularization prevent overfitting? ①. Intuition: regularization by



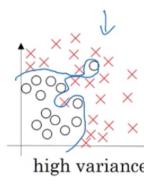
↑
l1



high bias



"just right"

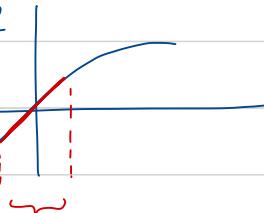


Andrew Ng

penalizing the weights, tries to reduce impact of hidden units, thus simplifying the network if needed.

②. take example of

$\tanh(\cdot)$



$\lambda \uparrow w^{[l]} \downarrow$ then $z^{[l]} \leftarrow w^{[l]} a^{[l-1]} + b^{[l]}$

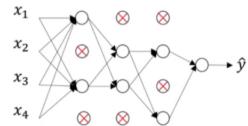
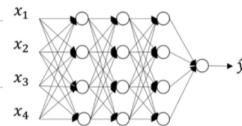
then every layer \rightarrow linear

During the iteration, when we plot the $J(\alpha, \gamma)$, we should consider that now we should add the $(+\frac{\lambda}{2m} \|w\|_F^2)$ term as well otherwise it may seem that it is not decreasing monotonically.

"Dropout" \rightarrow "as a regularization technique"

Dropout is another technique for regularizing which helps prevent overfitting

\rightarrow in each iteration, randomly some % of hidden units set to zero (turned off)



Implementing dropout ("Inverted dropout")

Illustrate with layer $l=3$. $\text{keep-prob} = 0.8$ 0.2

$\rightarrow d3 = np.random.rand(a3.shape[0], a3.shape[1]) < \text{keep-prob}$

$a3' = np.multiply(a3, d3)$ $\# a3' = d3$.

$\rightarrow a3' / \cancel{=} \text{keep-prob}$

So units \rightsquigarrow 10 units shut off

$$z^{(l+1)} = W^{(l+1)} \cdot a^{(l)} + b^{(l)}$$

\downarrow reduced by 20%.

$$/ = 0.8$$

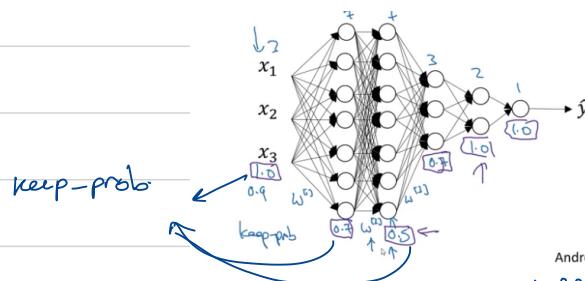
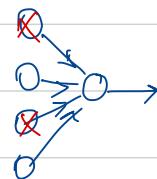
in each iteration the random dropout should regenerate not to keep the same

in implementation, it's important that if dropouts is used in each iteration a should be scaled to keep the same expected value.

dropout obviously is not implemented on the test time. (if done, it shouldn't change a lot but it will be possibly random). or could be done few times and averaged the results (not good idea still).

"why dropout works as regularization?" (intuition: can't rely on any feature, so have to spread the weights).

by knocking randomly the units, it force not to rely just on some and distribute the weight \rightarrow shrinking weight: similar to L_2 -norm



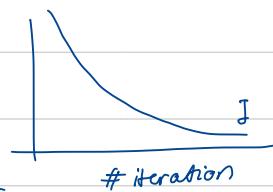
one essential parameter is keep-prob. and it can be different for each layer.
for example for bigger w we can set smaller keep-prob (0.5)
and 1 for input/output layers.

bad: More hyper parameters.

\rightarrow More is used in 'cv' which data is huge but not enough. very frequently used.

bad: cost function J is not well defined anymore

\rightarrow set keep-prob to 1 and check J is monotonically decreasing and then turn on dropout.



plotting J is like
debugging..

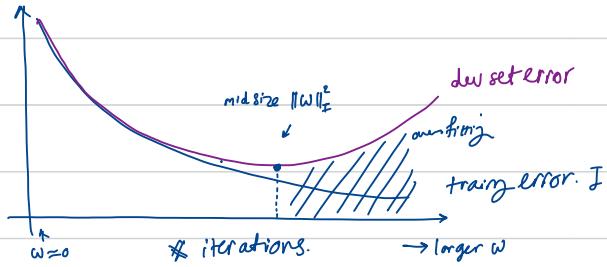
→ other techniques:

* Data augmentation

1. flipping.
2. cropping
3. zooming.

For digits we can add also random rotation and distortions as well.

* early stopping:



1. optimize cost function J (gradient descent, adam, ...)
2. Not overfit : use regularization

→ these two step should (better) to be look two different task. First Focus on 1 to choose the related hp, optimization and so, later 2. (Orthogonalization).
(one task at a time).

bad: early stopping breaks this rule (rather both are affected). alternative is to use l_2 which only focus to solve the overfitting (More h-p choice).

good: by running one gradient decent process we will check all the possible options for w .

better to use l_2 -norm (with different λ) if possible.

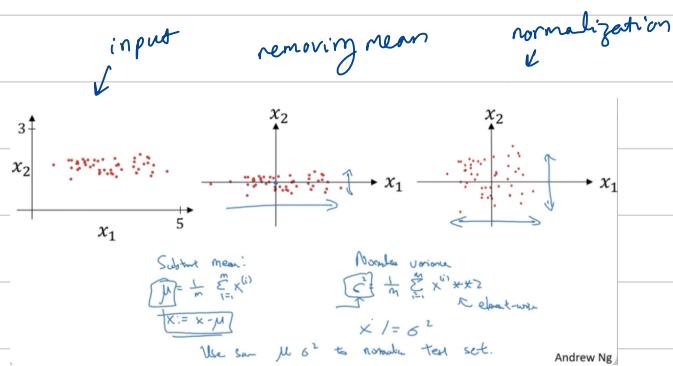
Setting up optimization problem

Normalizing inputs

(speed up training)!

→ test set should also use the same

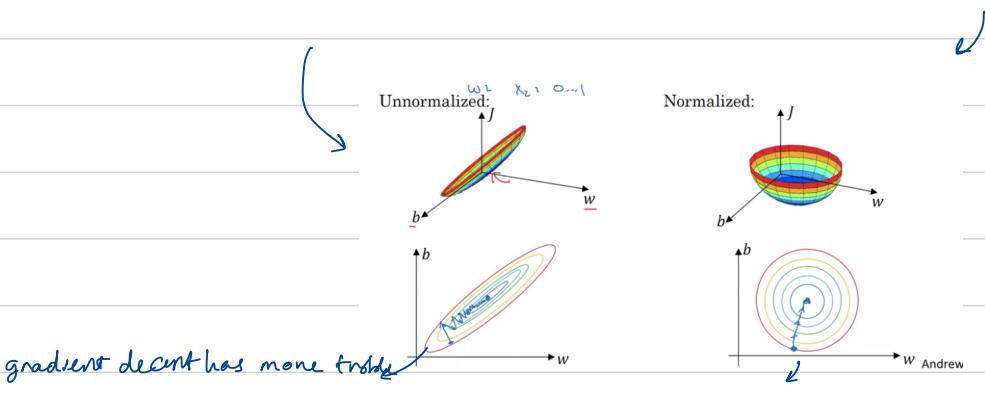
$$\mu, \sigma^2$$



why?

→ cost function could be very elongated.

more symmetric

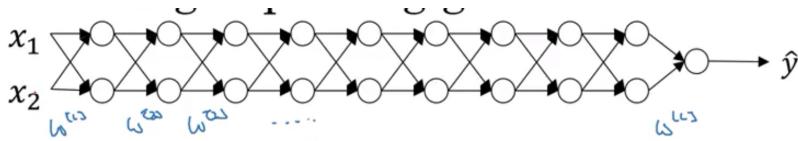


* at least it never hurts!

Vanishing or exploding gradient.

↳ gradient becomes very small both makes training difficult!
 ↳ very big.

take example of very deep NN. (take bsp and linear activation).



$$y = w^{(c)} \underbrace{w^{(c-1)} \cdots w^{(1)}}_{\text{---}} x$$

→ imagine w are same $\begin{bmatrix} a_0 & 0 \\ 0 & b_0 \end{bmatrix} = w^{(c)} \begin{bmatrix} a_0 & b_0 \end{bmatrix}^{L-1} x$

then if $a_0, b_0 \rightarrow > 1 \rightarrow$ activation will explode

$< 1 \rightarrow$ vanishes

→ same arguments are true for gradients as well.

Partial Solution

→ random initialization of weights. if $g \leq$ ReLU

$$z = w_1 x_1 + w_2 x_2 + \dots + w_n x_n$$

$$\begin{array}{c} z_1 \\ z_2 \\ \vdots \\ z_n \end{array} \Rightarrow \sum \hat{w}_i \rightarrow \hat{z} \quad n \uparrow \rightarrow w_i \text{ should be smaller}$$

↳ $\text{Var}(w_i) \leq \frac{1}{n}$ (for ReLU).

Python: $w^{(c)} \leftarrow \text{np.random.randn(} \right) \leftarrow \text{np.sqrt}\left(\frac{2}{n^{(c-1)}}\right)$

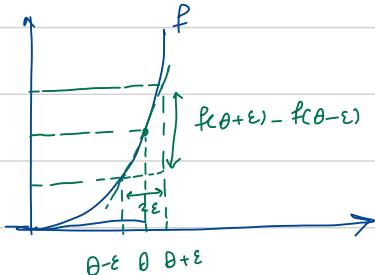
if tanh Xavier initialization $\sqrt{\frac{1}{n^{(c-1)}}}$

$$\sqrt{\frac{2}{n^{(c-1)} + n^{(c)}}}$$

Numerical approximation of gradients.

To double check if gd implementation is correct.

How to check derivative computation?



$$\frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \approx g(\theta)$$

$\theta = 1$
 $\epsilon = 0.01$

$$= 3.0001 \quad (\text{calculated})$$

$$= 3 \rightarrow (\text{real})$$

approximation error = 0.0001

(if done just 1 side error = 0.03)

$$f'(\theta) = \lim_{\epsilon \rightarrow 0} \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \rightarrow \text{error } O(\epsilon^2)$$

If we use $\lim_{\epsilon} \frac{f(\theta + \epsilon) - f(\theta)}{\epsilon}$

error $\rightarrow O(\epsilon)$

→ Now this numerical gradient calculation could be used to verify if our implementation of Back prop is correct.

→ Implementation *

Take all the $W^{[l]}, b^{[l]}$ and reshape into a big vector θ .

Now the cost $J(W^{[0]}, b^{[0]}, \dots, W^{[L]}, b^{[L]}) = J(\theta)$

Now all the $dW^{[l]}, db^{[l]}, \dots$ into a big vector $d\theta$

Gradient checking:

for each i :

$$d\theta_{approx}^{[i]} = \frac{J(\theta_1, \theta_2, \dots, \theta_i + \epsilon, \dots) - J(\theta_1, \dots, \theta_i - \epsilon, \dots)}{2\epsilon}$$

$$\approx d\theta^{[i]} \cdot \frac{\partial J}{\partial \theta^i}$$

check if they are close enough

$$\frac{\|d\theta_{approx} - d\theta\|_2}{\|d\theta_{approx}\|_2 + \|d\theta\|_2} \approx$$

for $\epsilon = 10^{-7}$ \sim should be $\ll 10^{-7}$

$\ll 10^{-5}$ Maybe fine but should be checked

$> 10^{-3}$ problematic

check which value is creating such diff \rightarrow

and track the problem.

Implementation Notes!

\rightarrow grad check only for debug not in training (.it's very slow)

\rightarrow if grad check fails, look at components which is problematic.

↳ For example $\mathbf{d}b^{(e)}$ are very different but not $\mathbf{d}W^{(e)}, \dots$
→ if regularization is used some remember to use it

$$\text{if } J(\theta) = \frac{1}{m} \sum_i L(\hat{y}_i^{(e)}, y_i^{(e)}) + \frac{\gamma}{2m} \sum_e \|\mathbf{w}_e^{(e)}\|^2$$
$$\mathbf{d}\theta \text{ is grad of } J \text{ w.r.t. } \theta.$$

→ grad check doesn't work with dropout. (obviously).

implement without dropout (keep_prob=1.0) and check
use grad check and if pass use drop-out.

→ Run grad check at random initialization and after few training
as well

if $\mathbf{w}, \mathbf{b} \approx 0$ is maybe looks correct but not.

"Optimization algorithms"

Batch vs mini-batch gradient descent.

Normally size of training sets are huge, if we want to process all once and then run just one step of gradient descent, it is not at all efficient.

Training set are divided in smaller mini-training set (minibatches).

$$X \leftarrow \begin{bmatrix} x^{(1)} & x^{(2)} \\ \vdots & \vdots \\ x^{(1)} & x^{(2)} \\ \vdots & \vdots \\ x^{(1)} & x^{(2)} \end{bmatrix}$$

similar for the Y as well,

in each iteration one mini-batch is given and gradient descent is applied

(i) \rightarrow training example index

$[]$ \rightarrow layer number

$\{ \}$ \rightarrow mini-batch index

Batch g.d.: if all the training set is used in each step of g.d.

mini-batch: if g.d. ran on mini-batches.

For $t = 1, \dots, 5000$

Forward prop on $X^{\{t\}}$

$$\begin{aligned} Z^{[0]} &= W^{[1]} X^{[t]} + b^{[1]} \\ A^{[0]} &= g^{[0]}(Z^{[0]}) \\ A^{[1]} &= g^{[1]}(Z^{[1]}) \end{aligned} \quad \left. \begin{array}{l} \\ \\ \end{array} \right\}$$

compute cost $J^{\{t\}} = \frac{1}{1000} \sum_{i=1}^l L(\hat{y}^i, y^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_w \|W^{[w]}\|_F^2$

For $t = 1, \dots, 5000$

Forward prop on $X^{[t]}$

$$\begin{aligned} Z^{[0]} &= W^{[0]} X^{[t]} + b^{[1]} \\ A^{[0]} &= g^{[0]}(Z^{[0]}) \\ \vdots \\ A^{[L]} &= g^{[L]}(Z^{[L]}) \end{aligned}$$

compute cost $J = \frac{1}{1000} \sum_{i=1}^l L(\hat{y}_i, y_i) + \frac{\lambda}{2 \cdot 5000} \sum_k \|W^{[k]}\|_F^2$

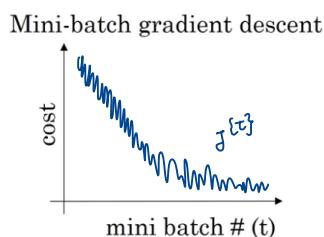
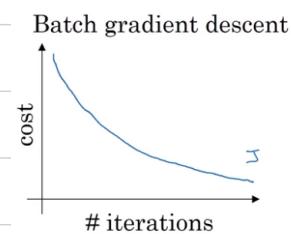
Back-prop to compute gradient w.r.t $J^{[t]}$ (using $(X^{[t]}, Y^{[t]})$)

$$w^{[k]} \leftarrow w^{[k]} - \alpha d w^{[k]}, b^{[k]} \rightarrow \text{update}$$

→ One full pass over full training set: "1 epoch"

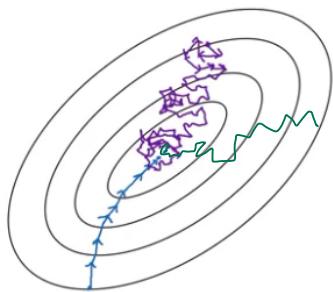
in batch gd in each epoch 1 gd step is done while in mini-batch gd, 5000 (ex.)

Normally many epochs has to be done (always mini-batch is used obviously).



this noisier trend is that
due to the fact that some
mini-batches one maybe
harder or easier than others.

{ mini batch size $\sim m$, Batch gradient descent (Just 1 min-batch)
 $\sim n \sim n = 1$, stochastic $n \sim n$ (any example is mini-batch).



batch gd trajectory (Too long)

stochastic gradient decent trajectory

always very noisy and will never reach to exact local minimum (will oscillate around)

(could be improved by using smaller learning rate)

advantage: gd step in each example of t.s.

→ Main disadvantage: losing speedup from vectorization.

* something between should be selected, (not too small / not too big).

→ Faster due to vectorization

→ faster progress on gd.

guide lines:

smaller training set : batch gradient descent. ($m \leq 2000$)

otherwise: mini-batch size: (64, 128, ... 512)

* Make sure mini-batch fits in cpu/gpu memory. (otherwise will be slow).

mini-batch size is another hyper-parameter.

Exponentially weighted averages (Moving averages).

faster than batch or mini-batch optimization.

applying a MA over data

with parameter β (< 1)

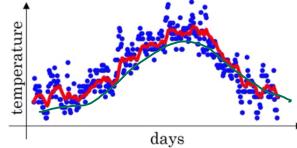
$$V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

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

- $\beta = 0.9$: ≈ 10 days' lag.
- $\beta = 0.98$: ≈ 50 days

V_t is approximately
average over
 $\approx \frac{1}{1-\beta}$ days'
temperature.

$$\frac{1}{1-0.98} = 50$$



Exponentially weighted averages

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

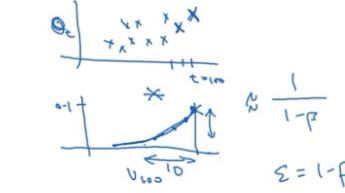
$$v_{100} = 0.9 v_{99} + 0.1 \theta_{100}$$

$$v_{99} = 0.9 v_{98} + 0.1 \theta_{99}$$

$$v_{98} = 0.9 v_{97} + 0.1 \theta_{98}$$

...

$$\begin{aligned} v_{100} &= 0.1 \theta_{100} + 0.9 \underbrace{(0.1 \theta_{99} + 0.9 \underbrace{\dots}_{\text{...}})}_{v_{99}} \\ &= 0.1 \theta_{100} + 0.1 \times 0.9 \cdot \theta_{99} + 0.1 (0.9)^2 \theta_{98} + 0.1 (0.9)^3 \theta_{97} + 0.1 (0.9)^4 \theta_{96} \\ &\quad + \dots \\ &\approx 0.9^{10} \approx 0.35 \approx \frac{1}{e} \end{aligned}$$



$$\sum = 1 - \beta$$

it is roughly we could
say that it is
averaging (weighted)

over $\frac{1}{1-\beta}$ past samples.

→ Implementation:

$$VA = \beta V_B + (1-\beta) \theta_B$$

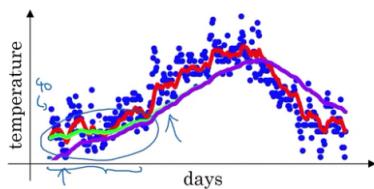
$$\frac{(1-\beta)^{1/e}}{\beta} = \frac{1}{e} \quad \beta=0.98?$$

$$\beta=0.98 \rightarrow \frac{0.98^{50}}{e} \approx \frac{1}{e}$$

Andrew Ng

in comparison with really averaging has the advantages of speed / Memory compared to real averaging.

"Bias Correction"



$$\begin{aligned} \rightarrow v_t &= \beta v_{t-1} + (1 - \beta) \theta_t \\ v_0 &= 0 \\ v_1 &= 0.98 v_0 + 0.02 \theta_1 \\ v_2 &= 0.98 v_1 + 0.02 \theta_2 \\ &= 0.98^2 v_0 + 0.02 \theta_1 + 0.02 \theta_2 \\ &= 0.0196 \theta_1 + 0.02 \theta_2 \end{aligned}$$

$$\beta = 0.98$$

$$\left| \begin{array}{l} \frac{v_t}{1 - \beta^t} \\ t=2: \quad 1 - \beta^t = 1 - (0.98)^2 = 0.0396 \\ \frac{v_1}{0.0396} = \frac{0.0196 \theta_1 + 0.02 \theta_2}{0.0396} \end{array} \right.$$

Andrew Ng

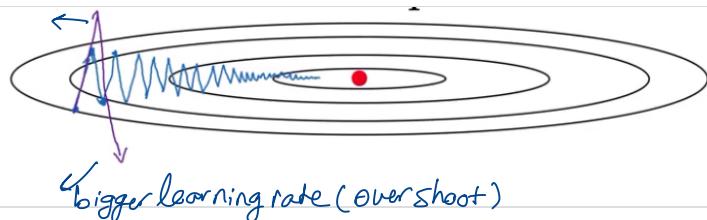
if β is close to 1 and $v_0 = 0$

Normally purple line is what we get, we apply bias correction term to reach better estimate of green line

if $v_0 = 0$ it takes more sample to reach better estimates.

adding bias correction term $\frac{1}{1 - \beta^t}$ (t = sample number)

→ "gradient descent with Momentum"



* blue line is the trajectory of gd(mini or batch)

what we want slower gradient on vertical axis

bigger ν \Rightarrow horizontal ν

"gradient descent with Momentum"

→ on each iteration +

compute d_w, d_b on current mini-batch (or batch).

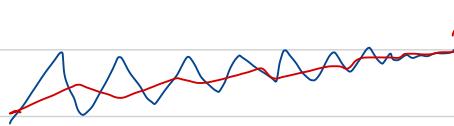
$$V_{dw} = \beta V_{dw} + (1 - \beta) d_w$$

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

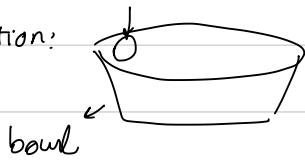
$$w = w - \alpha V_{dw}, b = b - \alpha V_{db}$$

This smooths out steps of gradient descent.

smooth out



Intuition:



d_w or d_b (acceleration)
 V_{dw} or V_{db} (velocity)
 β Friction

helps get the momentum in correct direction.

Implementation details

On iteration t :

Compute dW, db on the current mini-batch

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

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

$$W = W - \alpha v_{dw}, b = b - \alpha v_{db}$$

In practice we don't usually add

bias correction term $\frac{V_{dw}}{1 - \beta}$

(it just need to to reach, why bother)

V_{db}, V_{dw} one initialize so

Hyperparameters: α, β

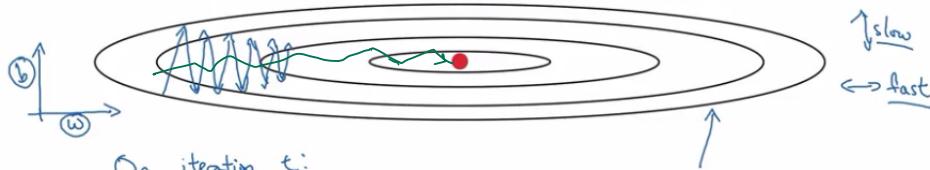
$$\beta = 0.9$$

some even omits the $(1-\beta)$ in updates: $V_{dw} \leftarrow \beta V_{dw} + dw$
 like integrating with α , it would mean the α should be modified
 properly (better to keep with original version).

* almost always work better than g.d. without momentum.

→ RMS prop

we saw the problem of g.d. when the steps are not right, for example we want gradient steps to be small in b , and larger for w (This could be for set of parameters, not necessarily all).



On iteration t :

Compute dW, db on current mini-batch,

$$S_{dw} = \beta S_{dw} + (1-\beta) \frac{dW^2}{\text{element-wise}} \quad \rightarrow \text{element wise power 2.}$$

similar to {
 exponentially averaging }
 $S_{db} = \beta S_{db} + (1-\beta) \frac{db^2}{\text{element-wise}}$ ← large
 (momentum)
 $w := w - \alpha \frac{dW}{\sqrt{S_{dw} + \epsilon}}$ ← ϵ : for Numerical stability
 $b := b - \alpha \frac{db}{\sqrt{S_{db} + \epsilon}}$ ←

S_{dw} is small then boost
 the steps for w → is big then reduces the
 steps for update of b .

This results in faster learning without diverging. The green line will be new trajectory. β is new hyperparameter, do not mix with the one from momentum (shown by β_2)

for Numerically stability we add ϵ to prevent dividing to 0

→ Adam Optimization (Adaptive moment estimation).

→ gradient decent with momentum + RMS prop

Adam optimization algorithm

$$V_{dw} = 0, S_{dw} = 0, V_{db} = 0, S_{db} = 0$$

On iteratn t :

Compute dW, db using current mini-batch

$$V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) dW, \quad V_{db} = \beta_1 V_{db} + (1 - \beta_1) db \leftarrow \text{"moment"} \beta_1$$

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) dW^2, \quad S_{db} = \beta_2 S_{db} + (1 - \beta_2) db \leftarrow \text{"RMSprop"} \beta_2$$

$$V_{dw}^{\text{corrected}} = V_{dw} / (1 - \beta_1^t), \quad V_{db}^{\text{corrected}} = V_{db} / (1 - \beta_1^t)$$

$$S_{dw}^{\text{corrected}} = S_{dw} / (1 - \beta_2^t), \quad S_{db}^{\text{corrected}} = S_{db} / (1 - \beta_2^t)$$

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

Andrew Ng

α : need to be tuned

$$\beta_1 : 0.9$$

$$(dW)$$

$$\beta_2 : 0.999$$

$$(dW^2)$$

$$\epsilon : 10^{-8}$$

, normally these default values
are used.

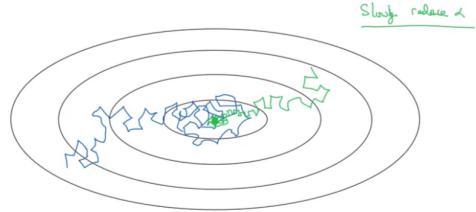
"learning rate decay"

Instead of using a fixed value for α , we reduce the α as we iterate to get closer to optimal point.

at earlier steps it is better to have $\uparrow \alpha$ to go faster but as it get closer $\downarrow \alpha$

$$\rightarrow \alpha_s = \frac{1}{1 + \text{decay rate} * \text{epoch_num}} \alpha_0$$

→ New hyper parameter



exponential decay. $\alpha_s = 0.95^{\text{epoch_num}} \cdot \alpha_0$

$$\alpha_s = \frac{K}{\sqrt{\text{epoch_num}}} \alpha_0 \quad \text{or} \quad \frac{K}{\sqrt{t}} \alpha_0$$

discrete staircase decay



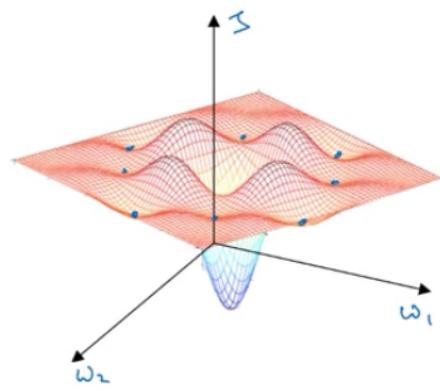
Manual decay : if it is just one model it is possible just to look at the error rate and adjust the α as it goes

"local optima"

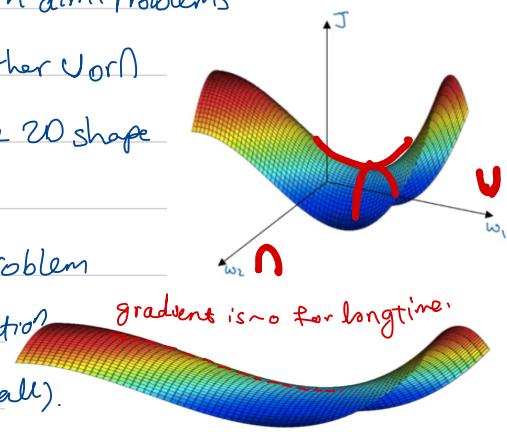
Previously, the intuition about the optimization problem is considered like this, One global optima with many local optima with zero gradient which could act like a trap.

(This is correct for smaller problems like this 2D visualization)

But for high dimensional problems, most of local optima points are actually saddle points, for high dim. problems at each point each param should be either U or V the chance all be U to be like the 2D shape is low (2^{-n} , n : # Parameters)



In higher dim. problem plateaus are problem
→ (slows the opt.) since for a large iteration the gradient is close to zero (very small).



→ local optima are not problem is DN, rather plateaus are problem in slowing the optimization
→ momentum, RMS prop or Adam work well to solve this

"hyper Parameter tuning process"

→ we saw that there are so many hyper parameter to tune always in DNN

α → the most important hyperparameter to tune.

B

B_1, B_2, ϵ
0.9 0.999 10^{-8}

* layers

hidden units

(for adam)

learning rate decay

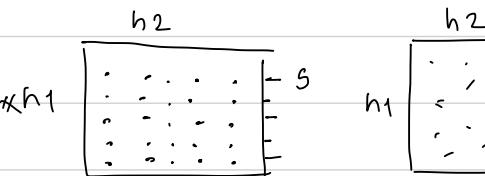
mini-batch size

Most important

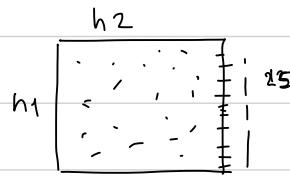
second priority

3rd priority

→ In earlier version it was common to do a grid over hyper parameters and check. but now we more tend to sample randomly over the hyper parameter of interest.

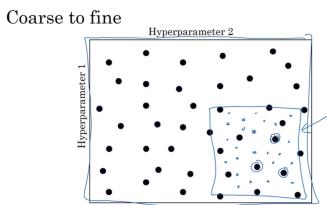


(sampled uniquely)



(randomly)

→ if we do randomly at same number of training we have checked more possibility for each (search more efficiently our h.p. space).



while searching for best value of hp. one approach could be to go in a coarse to fine grid (start randomly more coarse and zoom in region which looks better).

→ Using an appropriate scale to pick h.p.

→ Sampling uniformly random over h.p. space doesn't always
on efficient way.

for example for # layer [1 4] }
$n^{[4]}$ [50 ---- 100]

is not bad but for example consider α

$$[\alpha = 0.0001 \text{ ---- } 1]$$

if we sample uniformly 90% are [0.1 1] and only 10% [0.0001 0.1].
for such h.p. it is more reasonable to sample in log scale.

$$\begin{cases} r \sim 4 * np.random.rand() \\ n \sim 10^r \end{cases}$$

h.p for exponentially weighted averages.

$\beta \approx 0.9 \text{ ---- } 0.999$ (uniform random does not make sense).

$$1 - \beta \approx 0.1 \text{ ---- } 0.001 \quad [10^1 \quad 10^{-3}] \quad r \in [-3 \quad 1] \rightarrow 1 - \beta \approx 10^r \rightarrow \beta \approx 1 - 10^r$$

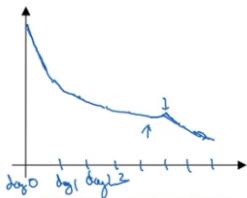
→ This helps to sample efficiently in ranges which are important.

$\beta \approx 0.9000 \rightarrow 0.9005$] roughly same average over 10

$\beta \approx 0.9990 \rightarrow 0.9995$] change from 1000 to 2000 huge change.

Re-test hyperparameters occasionally

Babysitting one model

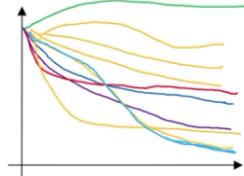


Panda ↪

when there are a lot of data but relatively less computation power.

Take one model and baby sit model but observing error (dev) and change h.p. accordingly.

Training many models in parallel



Caviar ↪

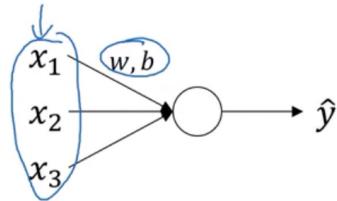
Andrew Ng

When there are enough computation power or less data we can train multiple model in parallel and chose the best one!

Batch Normalization.

→ search of h.p easier
robust for range of h.p
more easy to train even more deep networks.

Previously, we used the idea of normalizing inputs, to get more uniformed learning curves and to speedup algorithms like gradient-decent to optimize.



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

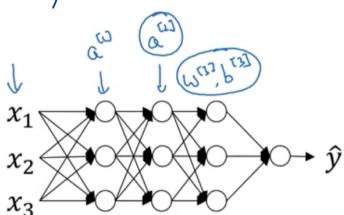
$$X = X - \mu$$

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

$$X = X / \sigma^2$$



Same idea in batch-normalization could be used for intermediate layers as well. (There is a debate whether to use Z or A for batch-Norm, but more we use Z (before activation)).



Can we normalize $\frac{a^{[2]}}{w^{[2]}, b^{[2]}}$ so fast?

$$\text{Normalizing } \frac{z^{[2]}}{\uparrow}$$

Andrew Ng

Implementation

given some intermediate values in Network :

$$\underbrace{z^{(1)}, \dots, z^{(n)}}_{\text{all from layer [l]}}$$

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

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

$$z_{\text{norm}}^{(i)} = \frac{z^{(i)} - m}{\sqrt{\sigma^2 + \epsilon}}$$

ϵ : for numerical stability

$$\tilde{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta \quad (\gamma, \beta \text{ learnable p. of model}).$$

Now use $\tilde{z}^{(i)}$ instead of $z^{(i)}$

→ This helps to control the mean and variance of values.

$$\begin{aligned} \gamma &= \sqrt{\sigma^2 + \epsilon} \\ \beta &= m \end{aligned} \quad \left. \begin{array}{l} \text{reverse the operation.} \\ \tilde{z} = z \end{array} \right\}$$

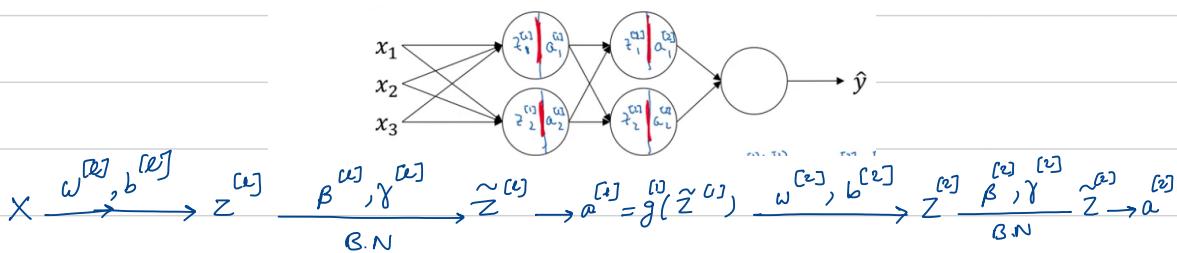
for example if we are using sigmoid activation, 
we don't want to distribute all around the linear part
and we use γ, β to use more of non-linearity of sig.

→ by Normalizing the intermediate layers b.n. helps speedup optimization similar to normalizing the inputs.

"Adding b.n. to a Network"

→ adding the B.N is as follow

in every layer z (output before activation) is normalized.

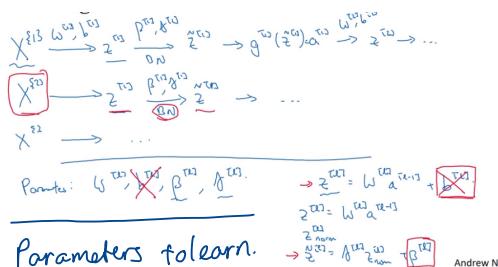


Parameter to be updated $w^{[0]}, b^{[0]}, w^{[1]}, b^{[1]}, \dots$

and new b.n ones $\beta^{[0]}, \gamma^{[0]}, \beta^{[1]}, \gamma^{[1]}, \dots \rightarrow$ these are also updated

similar to w, b $\beta^{[0]} \leftarrow \beta^{[0]} - \alpha \frac{\partial L}{\partial \beta^{[0]}}$

batch normalization works with mini-batches as well.



→ Now, having b is really not necessary

since in. b.n step we are re-adjusting
the bias. we can remove b

Implementing gradient descent

for $t = 1 \dots \text{num MiniBatchs}$
Compute forward pass on $X^{[t]}$.

In each hidden layer, use BN to map $\tilde{z}^{[t]}$ with $\tilde{z}^{[t]}$.

Use backprop & compute $\frac{\partial L^{[t]}}{\partial z^{[t]}}$, $\frac{\partial L^{[t]}}{\partial w^{[t]}}$, $\frac{\partial L^{[t]}}{\partial b^{[t]}}$

Update points $w^{[t]} := w^{[t]} - \alpha \frac{\partial L^{[t]}}{\partial w^{[t]}}$
 $\beta^{[t]} := \beta^{[t]} - \alpha \frac{\partial L^{[t]}}{\partial \beta^{[t]}}$
 $\gamma^{[t]} := \dots$

Works w/ momentum, RMSprop, Adam.

B.N work with / momentum,

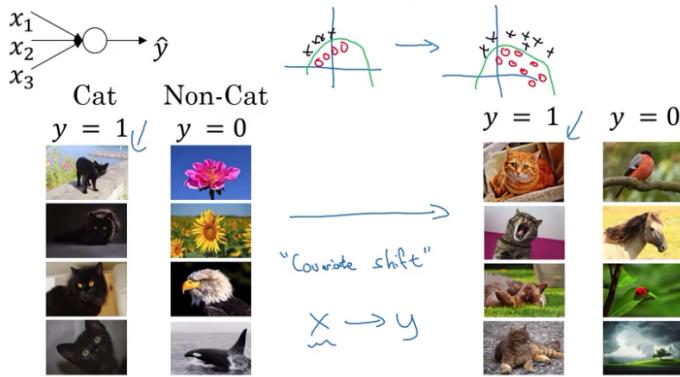
RMSprop, Adam,..

→ why P.N works? How speed ups optimization?

*→ similar to normalization of input could speed up learning (all values are similar scale and learning contour is more uniform).

*2: It make the weights in deeper layers more robust to changes in weights in earlier layers of N.N.

Learning on shifting input distribution



for example if we have trained our net on left data (black cats), Now will have problem with new data of changed distribution (other color).

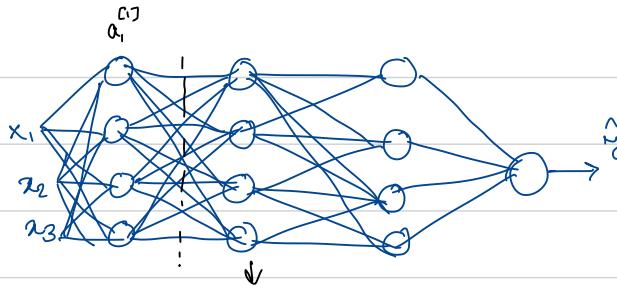
Andrew Ni

Like a model learned on left could not generalize to right even though it is same function.



→ This is called "covariate shift": if you learned a mapping $X \rightarrow Y$ if the ds's of data changed, You should re-learn the mapping.

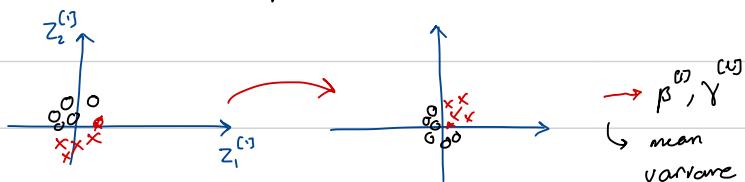
imagine a NN



from a prospective of this hidden layer, it receives

$a^{[L]}$ as input and tries to map it to output.

but indeed since weight in previous layers are changing then, distribution of $a^{[L]}$ are also changing. b.n. ensure that no matter how weight are changing at least it keeps the distribution same.



it limits the amount of effects of changing weights in previous layers, could change distribution of output.

→ (Weakening the dependency of layers learning).

Batch Norm as regularization

- Each mini-batch is scaled by the mean/variance computed on just that mini-batch. X
- This adds some noise to the values $z^{[L]}$ within that minibatch. So similar to dropout, it adds some noise to each hidden layer's activations. → drop out has multiplicative noise, b.n. has additive noise.
- This has a slight regularization effect. (regularization effect is not that much, better to be used with dropout. smaller b.s has more regularization but don't use b.n as regularization.)

applyin B.N at test time:

→ B.N is done as follow on forward path. on training:

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

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

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

$$\tilde{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta$$

μ, σ^2 are calculated on mini-batch.

but on test time we may not have a mini-batch, we tend to estimate μ, σ^2 over mini-batch one time on traing time.

→ μ, σ^2 : estimate using exponentially weighted average (across mini-batch)

$$x^{\{1\}}, x^{\{2\}}, \dots$$

↓

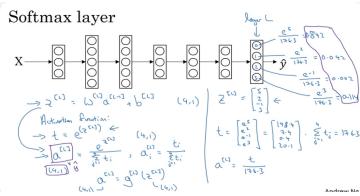
$$\begin{array}{ccccccccc} \mu^{\{1\}[2]} & \mu^{\{2\}[2]} & \mu^{\{3\}[2]} & \rightarrow & \text{average} & \mu^{[2]} \\ & & & & & & \\ & & & & \rightarrow & & \sigma^2 \end{array}$$

then at test time instead of using formula in \star use:

$$z_{\text{norm}} = \frac{z - \mu}{\sqrt{\sigma^2 + \epsilon}} , \tilde{z} = \gamma z_{\text{norm}} + \beta$$

softmax layer

in multi class classification, we could use an output vector representing each class



→ in this case we use softmax layer at output.

$$\rightarrow Z^{[c]} = w^{[c]} a^{[L-1]} + b^{[c]}$$

Activation function:

$$\rightarrow t = e^{z^{[c]}}$$

$$a^{[c]} = \frac{e^{z^{[c]}}}{\sum_{j=1}^C t_j}$$

$$Z^{[L]} = \begin{bmatrix} 5 \\ -2 \\ -1 \\ 3 \end{bmatrix}$$

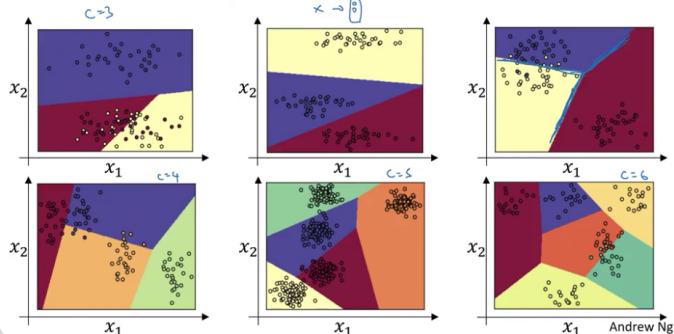
$$t = \begin{bmatrix} e^5 \\ e^{-2} \\ e^{-1} \\ e^3 \end{bmatrix}$$

$$\sum_j t_j = 176.3$$

$$a^{[c]} = \frac{t}{176.3} = \begin{bmatrix} 0.84 \\ 0.04 \\ 0.002 \\ 0.114 \end{bmatrix}$$

→ decision boundaries softmax represents:

Softmax examples



a NN with no hidden layer

$$Z^{[c]} = w^{[c]} x + b^{[c]}$$

$$a^{[c]} = \hat{y} = g(Z^{[c]})$$

Multi-class classification

Understanding softmax

$$z^{(l)} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix}$$

"soft row"

$$\alpha^{(l)} = g^{(l)}(z^{(l)}) = \begin{bmatrix} e^5/(e^5 + e^2 + e^{-1} + e^3) \\ e^2/(e^5 + e^2 + e^{-1} + e^3) \\ e^{-1}/(e^5 + e^2 + e^{-1} + e^3) \\ e^3/(e^5 + e^2 + e^{-1} + e^3) \end{bmatrix} = \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix}$$

"hard row"

Softmax regression generalizes logistic regression to C classes.

If $C=2$, softmax reduces to logistic regress. $\alpha^{(l)} = \begin{bmatrix} 0.842 \\ 0.158 \end{bmatrix}$

Andrew Ng

Loss function

$$\hat{y} = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} \quad \text{calc } y_2 = 1$$

$$y = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} \quad \hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix} \quad C=4$$

$$\mathcal{L}(\hat{y}, y) = -\sum_{j=1}^m y_j \log \hat{y}_j \quad \left| \mathcal{J}(w^{(1)}, b^{(1)}, \dots) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) \right.$$

small $-y_2 \log \hat{y}_2 = -\log \hat{y}_2$. make \hat{y}_2 big.

$$Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}] \quad \hat{Y} = [\hat{y}^{(1)} \ \dots \ \hat{y}^{(m)}]$$

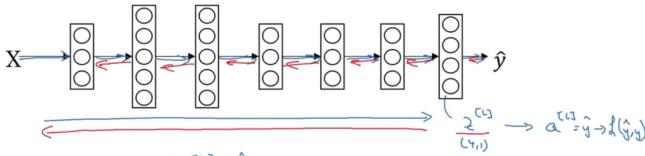
$$= \begin{bmatrix} 0 & 0 & 1 & \dots \\ 1 & 0 & 0 & \dots \\ 0 & 1 & 0 & \dots \\ 0 & 0 & 0 & \dots \end{bmatrix} \quad = \begin{bmatrix} 0 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix} \ \dots \ .$$

Andrew Ng

Maximum likelihood

loss function issues

Gradient descent with softmax



Backprop: $\frac{\partial \mathcal{L}}{\partial z^{(l)}} = \hat{y} - y$

drive the equation top over!

Andrew Ng

