

Train/Dev/Test Sets.

- Applied ML is iterative progress:

You don't know the value of # layers. # hidden unit ...

You just try \Rightarrow experiment \Rightarrow idea \Rightarrow try \Rightarrow ...

Dataset makes the cycle quickly.



If you have 100/1000/10000 Data. 70/30. OR 60/20/20.

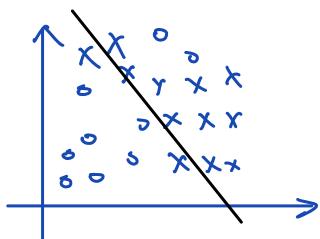
If you have 1,000,000 Data. Maybe 10000 for Dev and Test.
Like 98% train, 1% dev, 1% test.

- Mismatched train/test distribution.

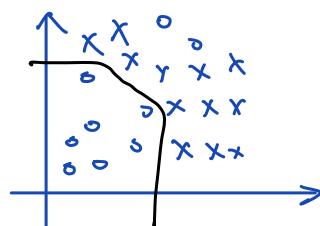
e.g. { Training set from webpage.
Dev/Test sets comes from your phone.

★ Rule of Thumb. Dev and test comes from the same distribution.
It is OK if you only have dev.

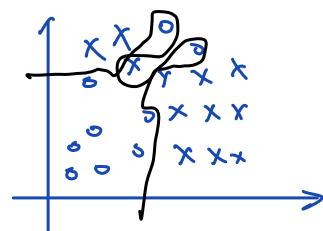
Bias/Variance.



High Bias
Underfitting



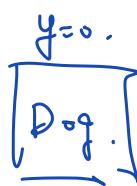
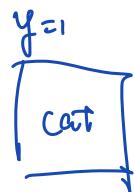
Just Right



High Variance
Overfitting

E.G.:

Cat classification:



Assumption: Train and Dev comes from the same distribution

Trainset Error: 1%

15%

15%

0.5%

Devset Error: 11%

16%

30%

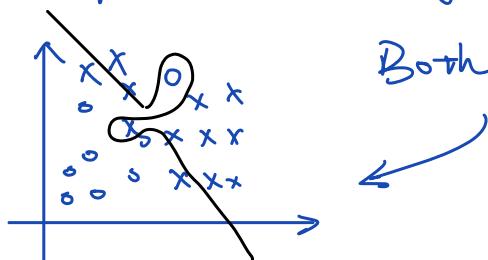
1%

high variance
(Overfitting)

high bias.
(Underfitting)

high bias &
high variance.
low bias
& low variance.

Crucial point: Optal (Bayes) Error: e.g. human classifier error.



Both high bias & variance.

Basic "Recipe" for machine learning.

① High Bias?

(Train data problem)

Yes.

• Bigger Network.

• Train longer.

② High Variance?

(Dev set problem).

Yes.

• More data.

• Regularization.

↓
Fixed ①②. Done.

* High bias problem - "More data" won't help.



* Bias Variance Tradeoff. its not accurate. We can drive down one thing without hurting another for just deep learning.

Regularization

e.g. Logistic regression.

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

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

$$L_2 \text{ regularization. } \|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = w^\top w.$$

L_1 regularization. $\frac{\lambda}{m} \sum_{i=1}^{n_x} |w_i| \rightarrow w$ will be sparse. because some w_i will be 0.

λ : regularization param. determined by cross-validation.

→ is a hyperparam need to be tuned.

Neural Network.

Cost function: $J(w^{[0]}, b^{[0]}, \dots, w^{[L]}, b^{[L]}) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|_F^2$

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

$\Delta w^{[l]}$ from backprop.

$$\frac{\partial J}{\partial w^{[l]}} = \Delta w^{[l]}.$$

$$w^{[l+1]} = w^{[l]} - \alpha \Delta w^{[l]}$$

Adding L_2 regularization,

We call this "weight decay".

$$\Delta w^{[l]} = (\text{from backprop}) + \frac{\lambda}{m} w^{[l]}$$

$$w^{[l+1]} = w^{[l]} - \alpha \left[(\text{from back}) + \frac{\lambda}{m} w^{[l]} \right]$$

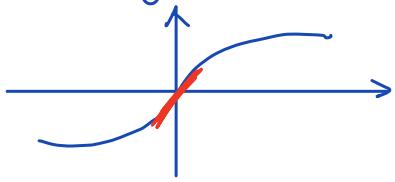
$$= w^{[l]} - \frac{\alpha \lambda}{m} w^{[l]} - \alpha (\text{from BP})$$

Now you are multiply $w^{[l]}$ by $(1 - \frac{\alpha \lambda}{m})$

Why regularization prevent overfitting.

Intuition: when $\lambda \rightarrow \infty$. we are setting some $w \rightarrow 0$. So the network actually become simpler. "Smaller Network".

e.g. Say we have a $\tanh(z)$ activation function.



$$\text{if } \lambda \uparrow, w^{[l]} \downarrow, z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]} \downarrow$$

So z will concentrate with the linear area.

Which we kind simplifies the network.

Dropout Regularization.

We we go through each layer and each node. We have a probability drop some nodes.

Remove all the connections to those nodes.

Do Backprop.

Next iteration (New example input).

Still using the probability to drop some of the nodes.

Implementing Dropout.

Illustrate with layers $l=3$. $\text{keep-prob} = 0.8 \rightarrow 0.2 \text{ prob to drop.}$

vector: $d_3 = \text{np.random.rand}(a_3.\text{shape}^{[0]}, a_3.\text{shape}^{[1]}) < \text{keep-prob}$.

z : represents the 3rd layer.

Boolean matrix. 0 or 1.

$a_3 = \text{np.multiply}(a_3, d_3)$.

$a_3 /= \text{keep-prob}$.

explanation: Say you have 10 units. On average, 10 units will get shut down.

$$z^{[l]} = w^{[l]} \cdot a^{[l-1]} + b^{[l]}$$

↑ value will be reduced by 20%.

So dividing 80%, we inverted the dropout effect.

Making prediction at test time.

$$a^{[0]} = x \quad \text{input.}$$

$$z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]}$$

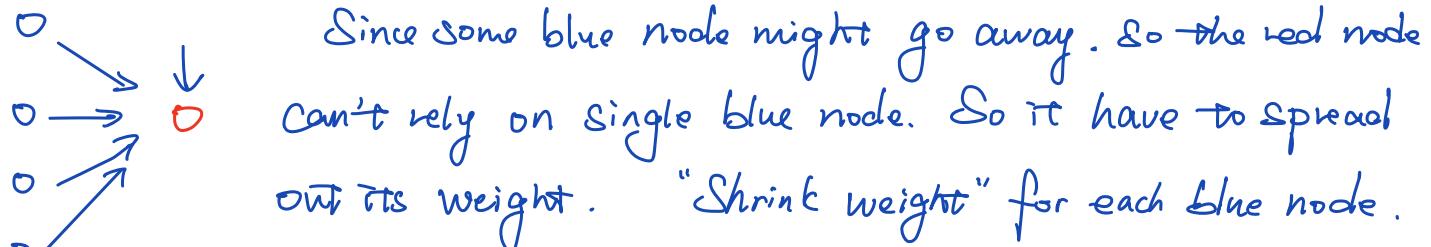
$$a^{[l]} = g^{[l]}(z^{[l]})$$

⋮

In test time, we don't want the dropout layer.

Understanding Dropout.

Intuition: Can't rely on any one feature. So have to spread out weight.



This looks like L_2 regularization.

Remember, no overfitting, no dropout.

Other regularization technique.

Data augmentation.

e.g. flipping images. rotating pics. zoom in/out pics.
rotation, distortion for digits images.

Early stopping.



Stop half way with mid-size weights $\|w\|^2$.

* for DNN, Optimizing J .

- Gradient descent.

Not overfit.

Orthogonalization

By doing early-stopping. You cannot do both things independently.

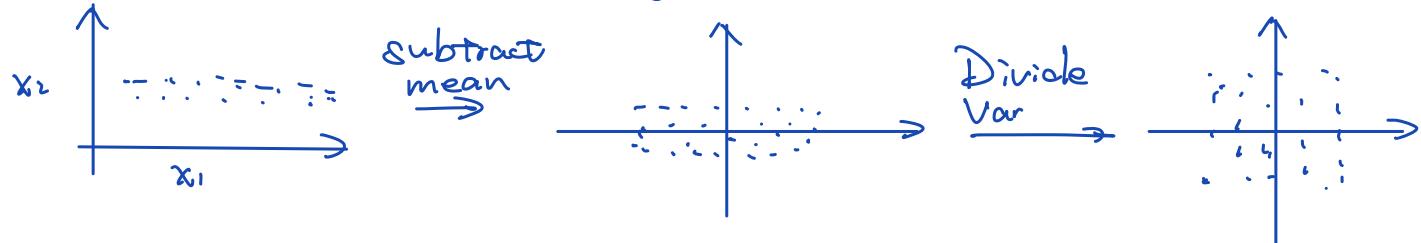
Like doing both things with one tool.

— Regularization.

Normalizing Inputs.

Normalizing Trainig Set. Subtracting mean \Rightarrow Zero mean.

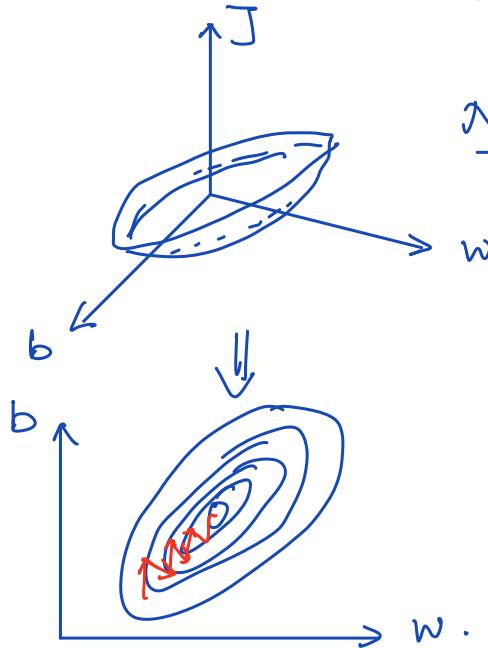
$\sigma^2 = \frac{1}{m} \sum x^{(i)^2}$. Divided by σ^2 . Variance both equal to 1.



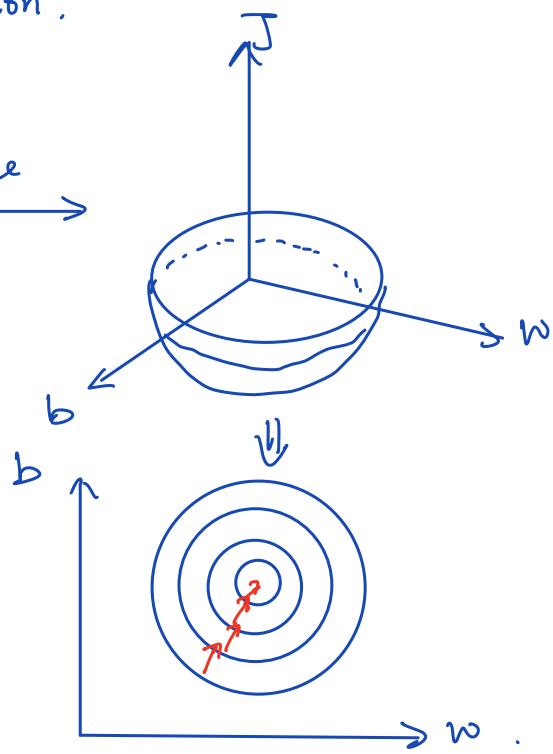
Use same μ and σ^2 to rescale the test set!

Why normalize features?

Unnormalized:



Rough intuition.



The gradient might be going squizy.

The gradient can go straight

Your features are in same scale.

Vanishing / Exploding gradient.

Suppose you are training a really deep NN.



$g(z) = \sigma$ as activation func. $\sigma^{[l]} = 0$.

$$\hat{y} = w^{[L]} w^{[L-1]} w^{[L-2]} \dots w^{[2]} w^{[1]} x$$

$$\underbrace{z^{[1]}}_{z^{[1]} = w^{[1]} x} = w^{[1]} x,$$

$$a^{[1]} = g(z^{[1]}) = \sigma^{[1]}$$

$$\underbrace{a^{[2]}}_{a^{[2]} = g(z^{[2]}) = g(w^{[2]} w^{[1]})} = g(w^{[2]} w^{[1]})$$

Suppose we have. $w^{[L]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$ $\hat{y} = w^{[L]} \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}^{L-1} \cdot x$

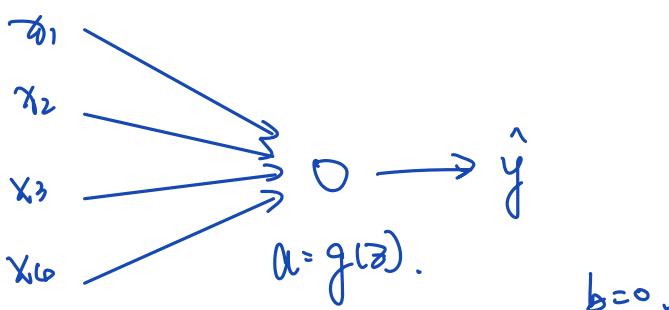
So if $L \rightarrow \infty$. we know this will explode. (exponential). \hat{y} really big.

If $w^{[L]} = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}$. Similarly, with $L \rightarrow \infty$. $\hat{y} \rightarrow 0$.

So. with really deep neural network. it will take really long time to train

Weight initialization in a deep network.

Single neuron example:



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

We don't want z blow up or too small.

So we want large n and small w_i .

One reasonable thing to do:

$$\text{Var}(w_i) = \frac{1}{n}$$

n : Number of features fed into one neuron.

So we can set: $W^{[l]} = \text{np.random.rand}(\text{shape},) * \text{np.sqrt}(\frac{1}{n^{[l-1]}})$
* If use ReLU. $(\frac{1}{n^{[l-1]}})$ sometimes perform better.

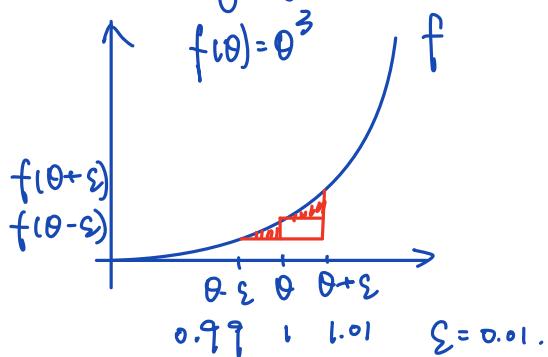
e.g. Xavier initialization. Random uniform distribution between:

$$\pm \sqrt{\frac{6}{n^{[l-1]} + n^{[l]}}} . \quad n^{[l]}: \text{Connections into the neuron.}$$

$n^{[l+1]}$ connection going out of the neuron.

Numerical Approximation of Gradient.

Checking your derivative computation.



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

$$\frac{(1.01)^3 - (0.99)^3}{0.02} = 3.0001.$$

$$g(\theta) = 3\theta^2 = 3. \quad \text{Approximation err: } 0.0001.$$

Gradient Checking. (Don't use in training. Only for debug).

Take $w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}$ and reshape it into big vector θ .

$$J(w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}) = J(\theta).$$

Take $dw^{[1]}, db^{[1]}, \dots$ and reshape it into big vector $d\theta$.

Q: is $d\theta$ the gradient of $J(\theta)$. $J(\theta) = J(\theta_0, \theta_1, \theta_2, \dots)$.

for each i:

$$\frac{d\theta_{\text{approx}}^{[i]}}{d\theta^{[i]}} = \frac{J(\theta_0, \theta_1, \dots, \theta_i + \epsilon, \dots) - J(\theta_0, \theta_1, \dots, \theta_i - \epsilon, \dots)}{2\epsilon}. \quad \epsilon = 10^{-7}$$

We know this should approximately equal to $\frac{d\theta^{[i]}}{d\theta^{[i]}} = \frac{dJ}{d\theta^{[i]}}$

Check $\frac{\|d\theta_{\text{approx}} - d\theta\|_2}{\|d\theta_{\text{approx}}\|_2 + \|d\theta\|_2} \approx 10^{-7}$ Great!

$$\approx 10^{-5}$$

10^{-3} worry.

Remember Regularization
Doesn't work with dropout
Run at random init.
Again after some training

Mini-Batch Gradient Descent. Speed up.

Vectorization allows you to efficiently compute on m examples.

$$X = [x^{(1)}, x^{(2)}, \dots, x^{(m)}]$$

$(n \times m)$

What if $m = 1,000,000$?

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

$(1 \times m)$

Split examples into "baby" training set. e.g. mini-batches of 1,000 each.

$$\text{e.g. } X^{(1)} = [x^{(1)}, \dots, x^{(1000)}]. \quad X^{(2)} = [x^{(1001)}, \dots, x^{(2000)}].$$

We have up to $X^{(1000)}$.

Split y accordingly.

Mini-batch t : $X^{(t)}, Y^{(t)}$. Dimension of $X^{(t)}$: $(n \times 1000)$. $Y^{(t)}$: 1×1000

for $t=1, \dots, 1000$.

Forward prop on $X^{(t)}$.

$$\left. \begin{array}{l} z^{(t)} = w^{(t)} x^{(t)} + b^{(t)} \\ A^{(t)} = g^{(1)}(z^{(t)}) \\ \vdots \\ A^{(t)} = g^{(L)}(z^{(t)}) \end{array} \right\}$$

1 step of gradient descent using $X^{(t)}, Y^{(t)}$.

You can vectorize this.

$$\text{Compute Cost } J = \frac{1}{1000} \sum_{i=1}^L \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum L ||w^{(t)}||_F^2$$

Back prop to compute gradient respect to $J^{(t)}$ (using $(X^{(t)}, Y^{(t)})$)

$$\text{Update weights: } w^{(t)} = w^{(t)} - \alpha \cdot \Delta w^{(t)}, \quad b^{(t)} = b^{(t)} - \alpha \cdot \Delta b^{(t)}.$$

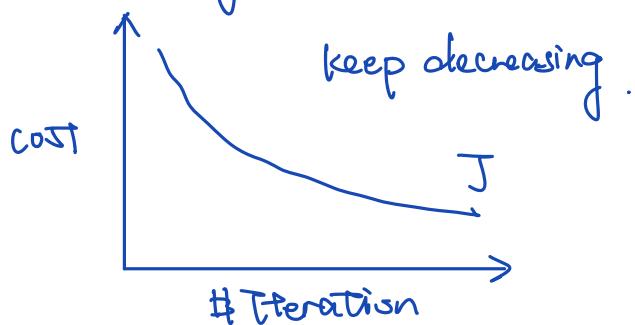
} → This entire for loop finish is 1 epoch of training.

"Single pass through the training set".

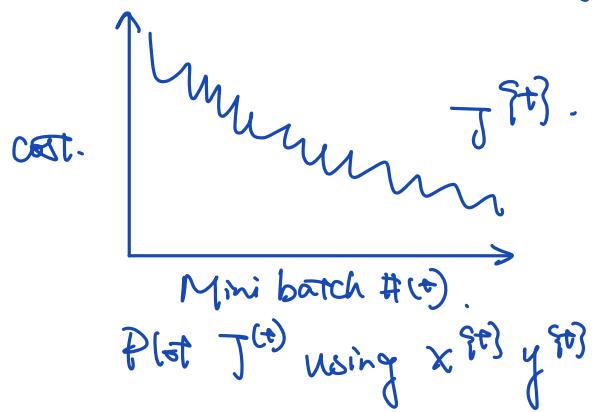
e.g. 100 epochs: Repeat 100 times for loop.

Understand Mini-batch Gradient Descent.

Batch gradient descent.



Trend: Going down.



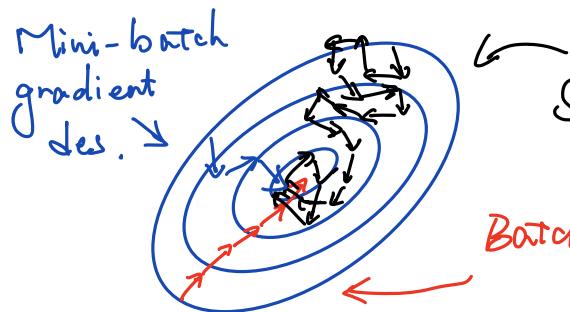
Each time you are throwing a different training

Choosing your size:

① If mini-batch size = m. Batch gradient descent.

② If ... size = 1. Stochastic gradient descent:

Every-example itself is a mini batch.



Stochastic Gradient Des.

Noise may be reduced by decreasing learning rate. Lossing adv. for vectorization.

Batch gradient des.

Processing a huge dataset in every iteration.

Time of each iteration is long.

So we want to take mini-batch size not too large nor too small.

How to choose batch-size?

If small training set. (e.g. < 2000). Use batch gradient descent.

Typical mini-batch size: 64, 128, 256, 512.

(Reason: the structure of memory is binary. Maybe faster).

Make sure minibatch fit the CPU/GPU memory.

Exponentially weighted averages.

e.g. Temperature in London.

$$\theta_1 = 40^\circ\text{F. } 4^\circ\text{C.}$$

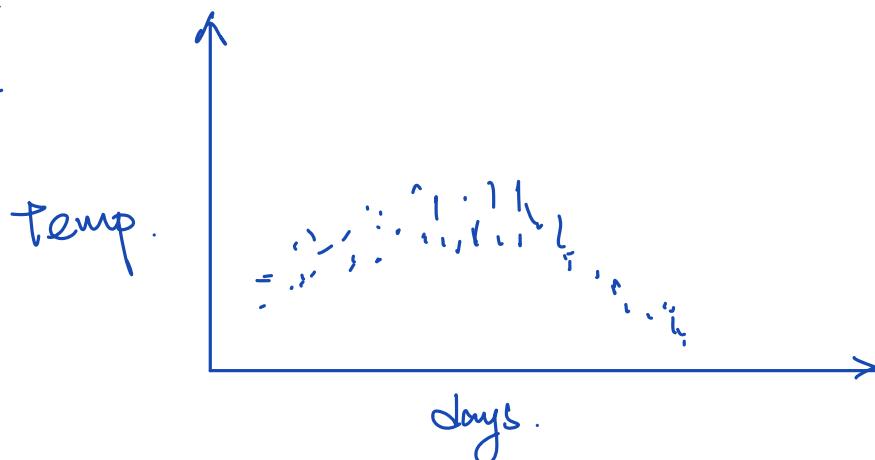
$$\theta_2 = 49^\circ\text{F. } 9^\circ\text{C.}$$

$$\theta_3 = 41^\circ\text{F}$$

:

$$\theta_{180} = 60^\circ\text{F.}$$

$$\theta_{181} = 56^\circ\text{F.}$$



$$V_0 = 0$$

$$V_1 = 0.9 V_0 + 0.1 \theta_1$$

$$V_2 = 0.9 V_1 + 0.1 \theta_2$$

$$V_3 = 0.9 V_2 + 0.1 \theta_3$$

:

$$V_t = 0.9 V_{t-1} + 0.1 \theta_t$$

Weighted average:

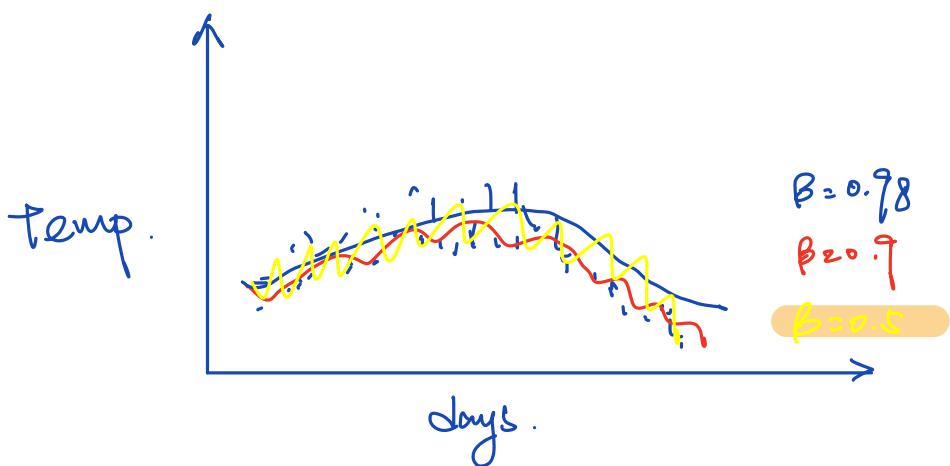
$$V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

V_t as approx. average over $\approx \frac{1}{1-\beta}$ days temperature.

$\beta = 0.9 \approx 10$ day temperature.

$\beta = 0.98 \approx 50$ days temperature.

$\beta = 0.1 \approx 2$ days temperature



Understanding exponentially weighted averages

Exponentially weighted average:

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

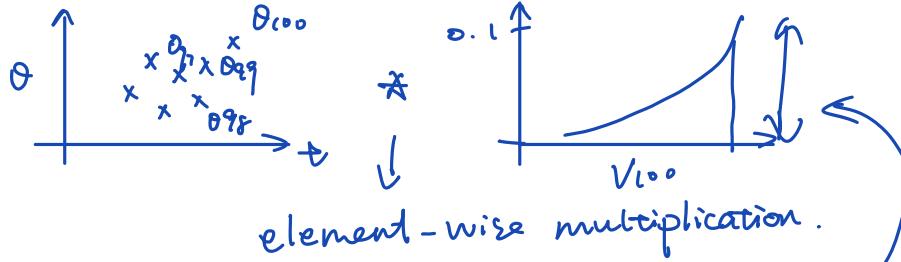
Take an example:

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

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

$$= 0.1 \theta_{100} + 0.1 \times 0.9 \theta_{99} + 0.1 \times 0.9 \times 0.9 \theta_{98} + \dots + 0.1 \times (0.9)^{99} \theta_1 + \dots$$

Suppose we have:



Wonder: how many days you are averaging?

Notice $0.9^{10} \approx \frac{1}{e}$. Generally: $(1-\varepsilon)^{\frac{1}{\varepsilon}} = \frac{1}{e}$. When $\varepsilon \rightarrow 0$. $\varepsilon = t - \beta$.

So in here, we take 10 days height of this decrease $\frac{2}{3}$

After 10 days, the weight decrease to $\frac{1}{3}$ approx.

Other e.g. $\varepsilon = 0.2$. $0.78^{10} \approx \frac{1}{e}$. Averaging over 10 days.

Bias Correction in exponentially weighted average.

We implement $V_0 = 0$. $V_1 = 0.98 V_0 + 0.02 \theta_1 = 0.02 \theta_1$.

$$V_2 = 0.98 V_1 + 0.02 \theta_2$$

$$= 0.9608 \theta_1 + 0.02 \theta_2$$

This is not a good estimate of
the first day temperature.



In general, lower estimate of early days.

$$\text{Better estimate: } \frac{V_t}{1-\beta^t}$$

$$t=2: 1-\beta^2 = 1-(0.98)^2 = 0.0396.$$

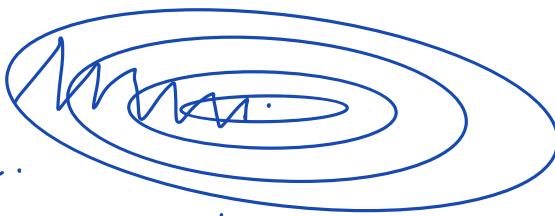
$$\text{New estimate: } \frac{V_2}{0.0396} = \frac{0.9608 \theta_1 + 0.02 \theta_2}{0.0396}.$$

Remove the bias.

When t goes large, the correction come to 1.

Gradient Descent with Momentum.

GD e.g.



keep oscillating.

But we cannot increase learning rate because it will overshoot.

Vertically: slower. Horizontally: faster.

Momentum:

On iteration t :

Compute $\Delta w \cdot \Delta b$ on current minibatch.

$$V_{dw} = \beta V_{dw} + (1-\beta) \Delta w \quad (\text{Similar to } V_b = \beta V_b + (1-\beta) \Delta b)$$

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

$$w = w - \alpha V_{dw}$$

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

velocity

acceleration

Hyperparam: α, β .

Normally: $\beta = 0.9$. Average across past 10 iterations.

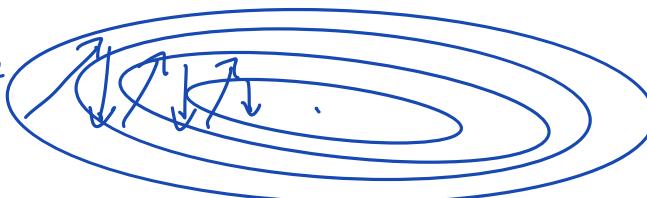
Now say we are doing descent:

Vertical average will go to the middle.

Horizontally, we kept going to the right. So we'll go much quicker.

RMSprop.

Root Mean Square



in practical:

some w_i will be vertical

some w_j will be horizontal

Say b : speed of learning vertically.

w : horizontally

On iteration t :

Compute $\Delta w \cdot \Delta b$ on current mini-batch.

$$Sdw = \beta Sdw + (1-\beta) dw^2 \rightarrow \text{elementwise square.}$$

$$Sdb = \beta Sdb + (1-\beta) db^2$$

$$w = w - \alpha \frac{dw}{\sqrt{Sdw} + \epsilon}$$

$$b = b - \alpha \cdot \frac{db}{\sqrt{Sdb} + \epsilon}$$

Notice: Sdw relative small. and Sdb is relative large.

So dw divide by small number and db divide by large value.

Note. we don't want Sdb be zero. So we add some epsilon to bottom.

Adam Optimization Algorithm.

Adam: Adaptive moment estimation.

Initialize: $Vdw = 0$, $Sdw = 0$, $Vdb = 0$, $Sdb = 0$.

On iteration t :

Compute dw , db using current mini-batch.

$$Vdw = \beta_1 Vdw + (1-\beta_1) dw. \quad \left. \right\} \text{Momentum } \beta_1.$$

$$Vdb = \beta_1 Vdb + (1-\beta_1) db \quad \left. \right\}$$

$$Sdw = \beta_2 Sdw + (1-\beta_2) dw^2 \quad \left. \right\} \text{RMSprop } \beta_2.$$

$$Sdb = \beta_2 Sdb + (1-\beta_2) db^2.$$

$$Vdw^{\text{corrected}} = Vdw / (1 - \beta_1^t).$$

$$Vdb^{\text{corrected}} = Vdb / (1 - \beta_1^t)$$

$$Sdw^{\text{corrected}} = Sdw / (1 - \beta_2^t)$$

$$Sdb^{\text{corrected}} = Sdb / (1 - \beta_2^t)$$

$$w = w - \alpha \frac{Vdw^{\text{corrected}}}{\sqrt{Sdw^{\text{corrected}}} + \epsilon}$$

$$b = b - \alpha \frac{Vdb^{\text{corrected}}}{\sqrt{Sdb^{\text{corrected}}} + \epsilon}$$

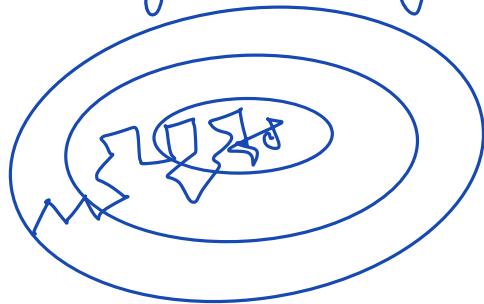
Bias
Correction

Hyperparam Choice:

α : Need to be tune. (Learning rate).

$\beta_1: 0.9$, $\beta_2: 0.999$, $\epsilon: 10^{-8}$. (Default value).

Learning Rate Decay.



Say your learning at the end just wander around.
Maybe just reduce α .
"Beginning, take large steps. later smaller step"

1 epoch = 1 pass through data.

Set $\alpha = \frac{1}{1 + \text{decay-rate} * \text{epoch-num}}$ α_0 . Decay-rate need to be tuned.

$\alpha_0 = 0.2$. Decay-rate : 1.

Epoch	α
1	0.1
2	0.67
3	0.5
4	0.4
:	:

Other learning rate decay method:

$$\alpha = \alpha_0 \cdot \frac{k}{\text{epoch-num}} - \text{exp. decay.}$$

$$\alpha = \frac{k}{\sqrt{\text{epoch-num}}} \cdot \alpha_0 \quad \text{or} \quad \frac{k}{\sqrt{t}} \alpha_0.$$

Hyperparam Tuning.

Hyperparameters: ① α : learning rate. ② β : Momentum. β_1, β_2, γ .
 $0.9, 0.999, 10^{-8}$.

Barely Tuned

② # layers, ① # hidden units, ③ learning rate decay, ② Mini-batch size.

Try random values: Don't use a grid. Use random sampling.

e.g. Hyperparam ①

Hyper
②

x	x	x	x	x
x	x	x	x	x
x	x	x	x	x
x	x	x	x	x
x	x	x	x	x

e.g. ②. Hyper ①

Hyper
②

x	x	x
x	x	
x	y	y

Randomly sample.

Grid Sample.

● Coarse to fine.

When sampling. Shrink sample range. Sample more densely.

Use an appropriate scale to pick hyperparams.

Picking HP at random.

$$n^{[l]} = 50 \text{ or } 100$$

Number of hidden units.

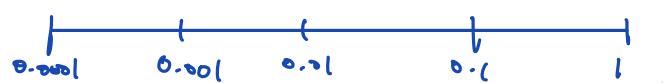
You think it is good to pick 10 - 100 units.

layers: 2-4.

Appropriate scale for HP.

Say for $\lambda = 0.0001 \text{ n.l.}$. If sample uniformly. 90% will between 0.1 to 1. Only 10% in 0.0001 to 0.1.

Instead, you should do log.



in python: $r = -4 * \text{np.random.rand}()$. $\leftarrow r \in [-4, 0]$.

$$\alpha = 10^r : 10^{-4}, \dots, 10^0.$$

HP for exp. weighted averages.

$\beta = 0.9, 0.99, 0.999$. So it is un-reasonable to search from 0.9

Averaging: $\frac{\downarrow}{10}, \frac{\downarrow}{100}, \frac{\downarrow}{1000}$ to 0.999 linearly.

We can use the previous method to sample on



$$10^{-1} \quad 10^{-3}$$

Reason behind this.

$$r \in [-3, -1]. \quad 1 - \beta = 10^r. \quad \beta = 1 - 10^{-r}$$

When β change from 0.9 to 0.999. (Average change 10 values).

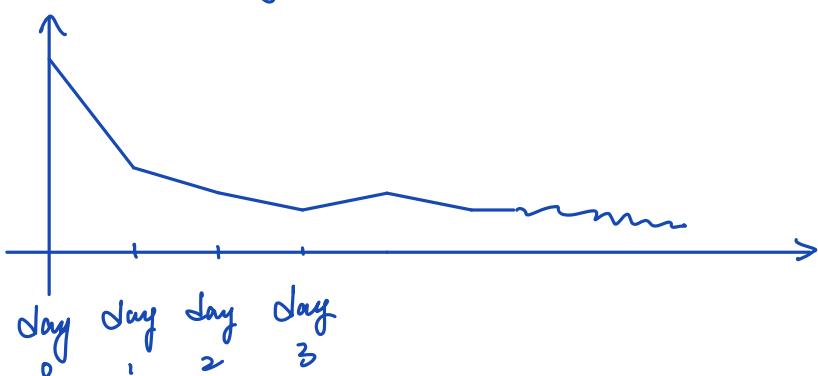
Not a big deal.

When β change from 0.999 to 0.9999. (Average change 1000 values).

It is a huge change.

HP tuning in practice: Pandas VS Caviar.

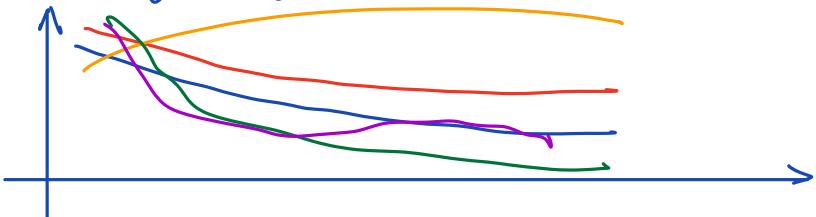
Baby-sitting one model.



Panda model.

Only one baby. Try lots of things.

Training many models in parallel.



Basically trying a lot of models.

Caviar. fishing. Didn't pay a lot of attention on specific one.

Normalizing activations in a network.

Normalizing inputs can speed up learning.



What about multiple layers.

Can we normalize $a^{[i]}$ so to train $w^{[i+1]} b^{[i+1]}$ faster.

We normalizing $\tilde{z}^{[i]}$. (There are debate whether normalizing before or after activation function).
 (\tilde{z}) (a) .

Implementing Batch-Norm.

Given some intermediate value in NN.

Hidden units: $\tilde{z}^{(i)}, \dots, \tilde{z}^{(m)}$. $\rightarrow (\tilde{z}^{[i]}, \dots)$
 $\mu = \frac{1}{m} \sum_i \tilde{z}^{(i)}$ which layer.

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

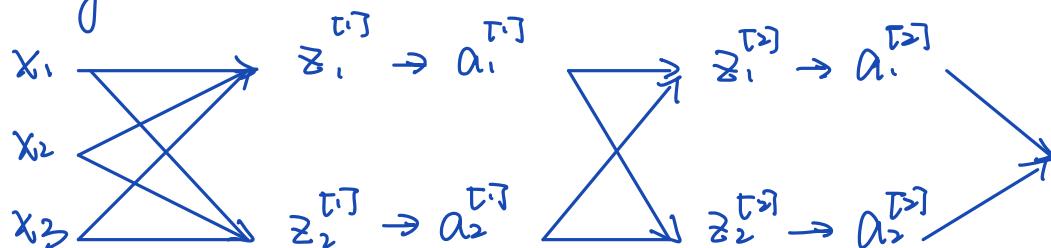
$$\tilde{z}_{\text{norm}}^{(i)} = \frac{\tilde{z}^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$\tilde{z}^{(i)} = \gamma \tilde{z}_{\text{norm}}^{(i)} + \beta. \quad (\gamma \text{ and } \beta \text{ are learnable params})$$

This allow you to set mean and variance of $\tilde{z}_{\text{norm}}^{(i)}$

$$\text{If } \gamma = \sqrt{\sigma^2 + \epsilon}, \beta = \mu. \text{ Then } \tilde{z}^{(i)} = \tilde{z}^{(i)}.$$

Fitting BN in to a NN.



$$x \xrightarrow{W^{[t]}, b^{[t]}} z^{[t]} \xrightarrow[\text{Batch-Norm (BN)}}{B^{[t]}, \gamma^{[t]}} \tilde{z}^{[t]} \rightarrow \alpha^{[t]} \xrightarrow{W^{[t+1]}, b^{[t+1]}} z^{[t+1]} \xrightarrow[\text{BN}}{B^{[t+1]}, \gamma^{[t+1]}} \tilde{z}^{[t+1]} \rightarrow \alpha^{[t+1]}$$

Parameters: $\left. \begin{array}{l} W^{[t]}, b^{[t]}, W^{[t+1]}, b^{[t+1]}, \dots, W^{[L]}, b^{[L]} \\ B^{[t]}, \gamma^{[t]}, \beta^{[t]}, \gamma^{[t]}, \dots, \beta^{[L]}, \gamma^{[L]} \end{array} \right\} \frac{d\beta^{[L]}}{\beta^{[L]} - d\beta^{[L]}}$

e.g. tensorflow: `tf.nn.batch_normalization.` (One-line code).

Working with mini-batch.

$$x^{[i]} \xrightarrow{W^{[t]}, b^{[t]}} z^{[i]} \xrightarrow[\text{BN}}{B^{[t]}, \gamma^{[t]}} \tilde{z}^{[i]} \dots$$

Normalizing only inside the mini-batch.

Params: $W^{[L]}, b^{[L]}, \beta^{[L]}, \gamma^{[L]}$.

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

When compute the norm of $z^{[L]}$. We actually cancel the constant $b^{[L]}$.

So we only do $z^{[L]} = W^{[L]} a^{[L-1]}$. then compute $\tilde{z}^{[L]}$.

$\Rightarrow \tilde{z}^{[i]} = \gamma^{[L]} z_{\text{norm}}^{[L]} + \beta^{[L]}$. ($\beta^{[L]}$ is actually performing the "bias")

Implement GD.

for $t=1, \dots, \# \text{Mini-batch}$.

Compute forward prop on $x^{[t]}$.

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

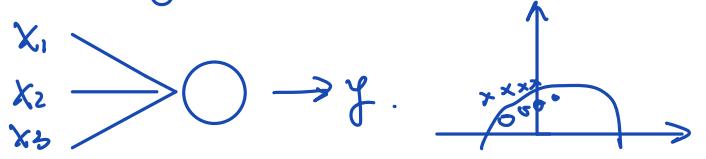
Use back prop to compute $dW^{[L]}, d\beta^{[L]}, d\gamma^{[L]}$.

Update: $W^{[L]} = W^{[L]} - \lambda dW^{[L]}$.
 $\beta^{[L]} \dots, \gamma^{[L]} = \dots$

Work with momentum. RMSprop. Adam. ...

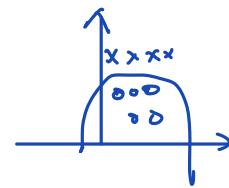
Why does BN work?

Learning on shifting input distribution.



e.g. Cat. Non-Cat.
 $y=1$ $y=0$.

Cat Non-Cat.
 $y=1$ $y=0$



Dr Pos e.g.
x: Neg e.g.

Black cats. \longrightarrow Colored - Cat

So your classifier won't work
really well.

"Covariate Shift".

Why this is a prob. for NN?

e.g. for w^{T3}, b^{T3} ,
 w^{T4}, b^{T4} , w^{T5}, b^{T5}

a_1^{T2}

a_2^{T2}

a_3^{T2}

a_4^{T2}

x_1
 $x_2 \Rightarrow 0$
 x_3



As $w^{T1}, b^{T1}, w^{T2}, b^{T2}$, change.

a_i^{T2} keep changing in w^{T2}, b^{T2} perspective.
Distribution

What BN does. is "fixing" the shift of the hidden units distribution.

No matter how it changes. Mean and variance governed by γ and β .

Make hidden units value more STABLE. (Learn indep. for each layer)

- BN as regularization.

which has some noise \downarrow Similar for β and γ .

- Each mini-batch is scaled by the mean/var computed on just that minibatch.

- This adds some noise to the values z^{T2} within that minibatch.

Some similar to dropout. it adds some noise to each hidden layer's activation.

- This has a slight regularization effect. (unintended effect).

Batch Norm in Test time. in test time, you may not have "mini-batch".

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

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

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

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

μ, σ^2 : estimated using exp. weighted average (across mini-batch).

$$x_1^{(1)}, x_2^{(1)}, \dots$$

$$\downarrow \quad \downarrow$$

$$\mu^{[1]}, \sigma^{[1]} \dots$$

$$\theta_1, \theta_2, \dots$$

Compute $\theta_1, \theta_2, \dots$ using exp. weight average. And finally get the μ .

Keep a running average. Similar for σ^2 .

Compute $z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$. $z^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta$. Use these μ, σ^2 for the hidden units in your network during test time.