

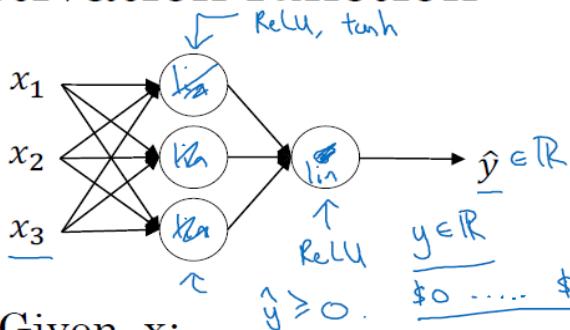
Course 1:

Activation Functions

One place that we use a linear activation function is when we want to do a kind of regression problem (continuous values, not with classes). In this case, we only use linear activation function at the output layer only. And for the rest of layers we use tanh or ReLu activation function.

If we want to use a linear activation function all over the places, we end up with a Linear Regression problem as can be seen in the picture below.

Activation function



$$\begin{aligned}
 & 2 \\
 a^{[1]} &= z^{[1]} = w^{[1]} x + b^{[1]} \\
 a^{[2]} &= z^{[2]} = w^{[2]} a^{[1]} + b^{[2]} \\
 a^{[3]} &= w^{[2]} (w^{[1]} x + b^{[1]}) + b^{[2]} \\
 & 3
 \end{aligned}$$

Given x :

$$\begin{aligned}
 \rightarrow z^{[1]} &= W^{[1]} x + b^{[1]} \\
 \rightarrow a^{[1]} &= g^{[1]}(z^{[1]}) \geq 0 \\
 \rightarrow z^{[2]} &= W^{[2]} a^{[1]} + b^{[2]} \\
 \rightarrow a^{[2]} &= g^{[2]}(z^{[2]}) \geq 0
 \end{aligned}$$

$$\begin{aligned}
 & 1 \\
 g(z) &= z \\
 & \text{"linear activation function"}
 \end{aligned}$$

$$\begin{aligned}
 & 4 \\
 & = (w^{[2]} w^{[1]}) x + (w^{[2]} b^{[1]} + b^{[2]}) \\
 & = \boxed{w' x + b'} \\
 & 5 \\
 g(x) &= z
 \end{aligned}$$

Andrew Ng

For the derivatives of **activation functions** we have: (look at the activation function figures as well)

Sigmoid: It totally makes sense to use sigmoid function for the output layer if the expected output is either 0 or 1 (binary classification), otherwise, better to use tanh function for the middle layers rather than the sigmoid one. The reason is written on tanh description part. (look at the next item)

$$g(z) = \frac{1}{1+e^{-z}} \rightarrow g'(z) = g(z) \times (1 - g(z))$$

Tanh: It is better to use tanh instead of sigmoid since the values are shifted between 1 and -1, so that the mean goes to 0 rather than 0.5 which makes learning for the next layers much easier.

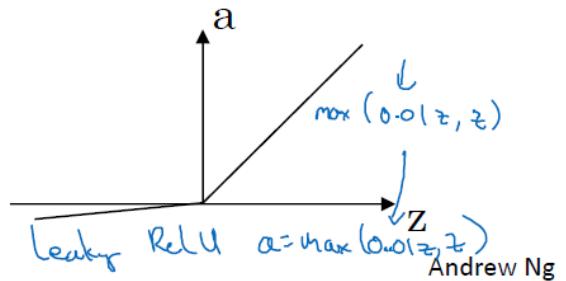
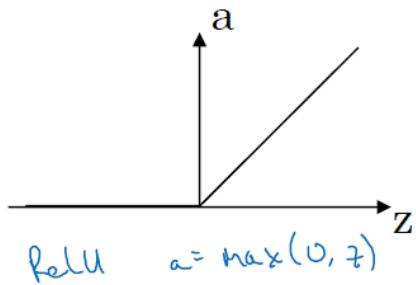
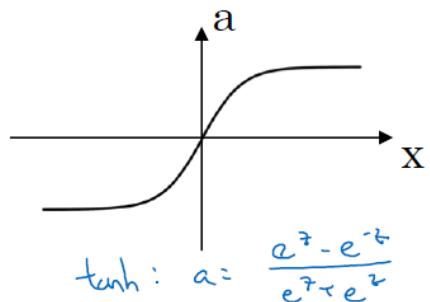
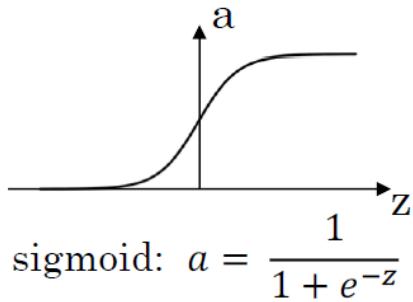
$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} \rightarrow g'(z) = 1 - g(z)^2$$

ReLu: The problem with tanh and sigmoid functions are that, when the values go toward two ends of the function, the derivative goes to 0, which slowing down the gradient descent. In this function only for the negative values we have derivative of 0 but for higher than that the derivative is 1. In most of the cases it is the default choice of activation function as it is faster on gradient descent.

$$g(z) = \max(0, z) \rightarrow g'(z) = 0 \text{ if } z < 0 \text{ else } 1$$

Leaky ReLu: If having a 0 as derivative of negative values is problematic for a problem, this version works.

$$g(z) = \max(0.01 \times z, z) \rightarrow g'(z) = 0.01 \text{ if } z < 0 \text{ else } 1$$



Weight Initialization

What happens when we initialize weights with 0s? It causes all weights to have same rows and causes **symmetric weights** meaning all the weights compute same function as their effect on the output at the very first was the same. So it is useless to do so.

So we initialize it with very small random values. Why small? Because if we use sigmoid or tanh, the weights tend to get large and the gradient descent takes very small steps as it gets so close to 0.

CS230: Lecture 3

Attacking Networks with Adversarial Examples

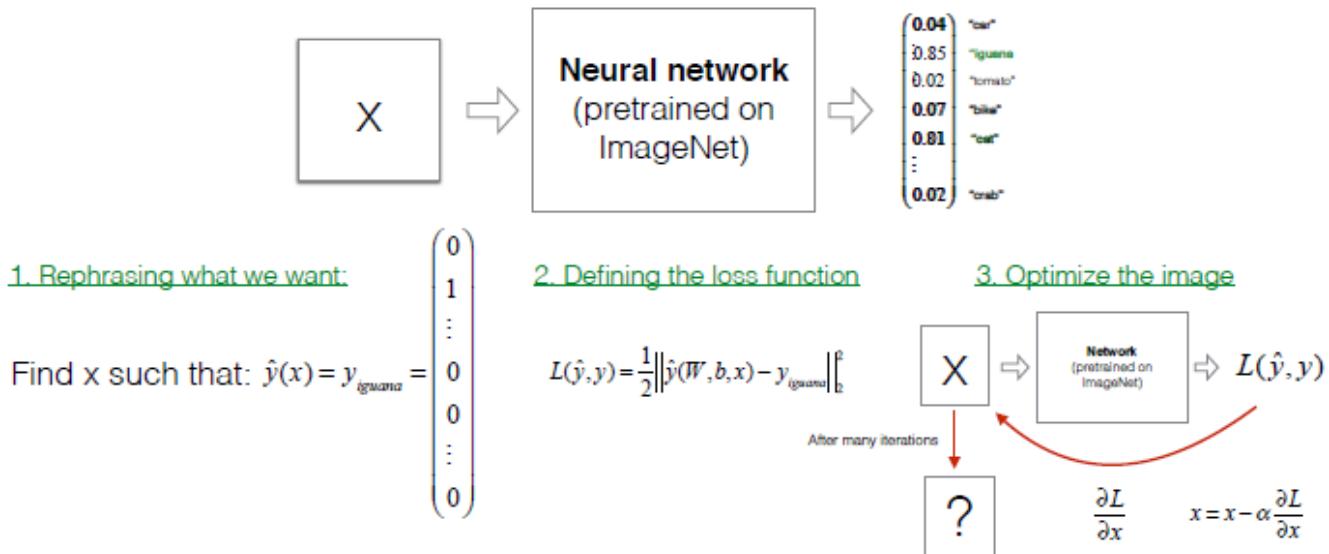
The goal of attacking a network, with adversarial examples is:

Given a network pertained on ImageNet, we want to find an input image, that is not an iguana, but will be classified as an iguana by the network.

For a good network, we expect the output probabilities of an input cat image, have maximum probability for the cat (out of Softmax). But, we want to fool the network in a way to see highest probability for the iguana instead. If we want to formulize the problem, we can see it in the slide below. Note that the loss function could be any other loss like cross entropy.

I. A. Attacking a network with adversarial examples

Goal: Given a network pretrained on ImageNet, find an input image that will be classified as an iguana.

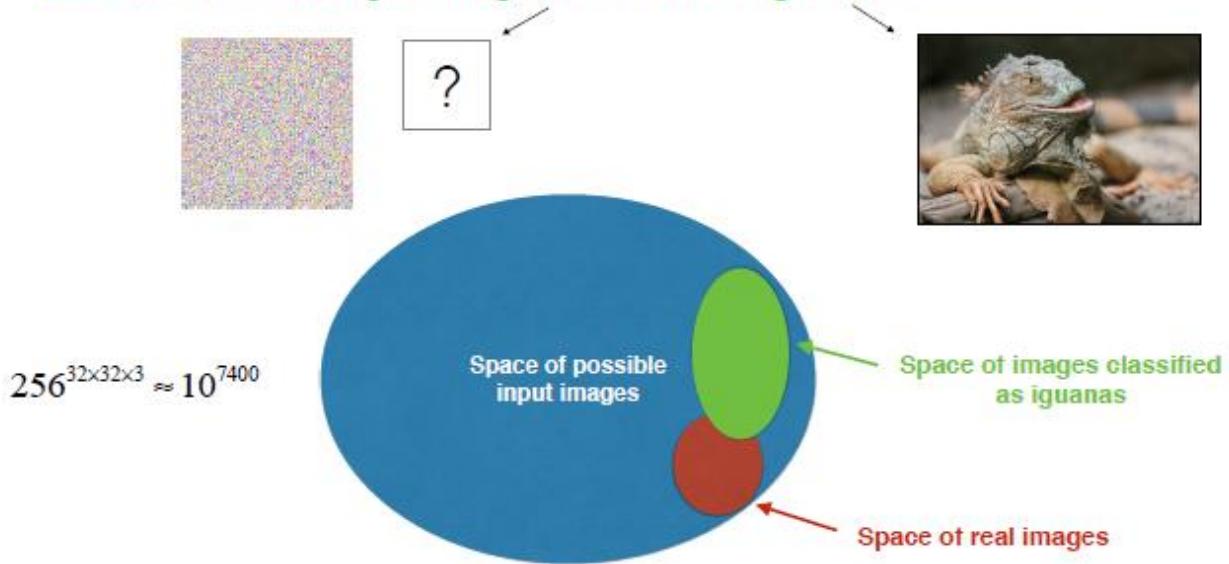


The space of all possible images for a 32×32 image, with 3 RGB channels would be: $256^{32 \times 32 \times 3}$. The space of real and iguana images are shown in red and green respectively. There is a slight overlap between the red and green space. The overlap is where we input an image of a cat and the network finds it as an iguana. Some of its malicious applications could be used in passing face recognition, CAPTCHAs, or violent content detection in social media.

So, in order to fool the network, our loss function needs some constraints which pushes it to the overlapped space.

I. A. Attacking a network with adversarial examples

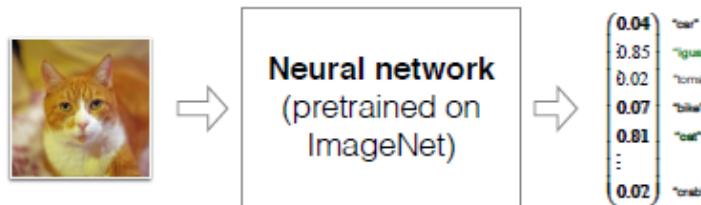
Question: Will the forged image x look like an **iguana**?



In our loss, we want to minimize the difference of the output image and the iguana one as well the difference of the output image and a cat, simultaneously. We can now give the image of the same cat to the network and get an image which is so similar to a cat (human can differentiate easily) but outputted as an iguana.

I. A. Attacking a network with adversarial examples

Goal: Given a network pretrained on ImageNet, find an input image displaying a cat but classified as an iguana.



1. Rephrasing what we want:

$$\text{Find } x \text{ such that: } \hat{y}(x) = y_{\text{iguana}} = \begin{pmatrix} 0 \\ 1 \\ \vdots \\ 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

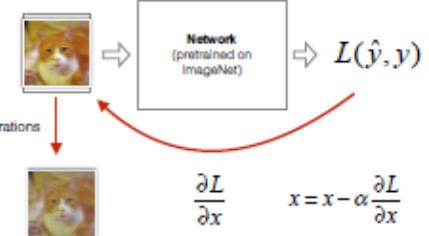
And: $x = x_{\text{cat}}$

2. Defining the loss function

$$L(\hat{y}, y) = \frac{1}{2} \left\| \hat{y}(W, b, x) - y_{\text{iguana}} \right\|_2^2 + \lambda \left\| x - x_{\text{cat}} \right\|_2^2$$

After many iterations

3. Optimize the image

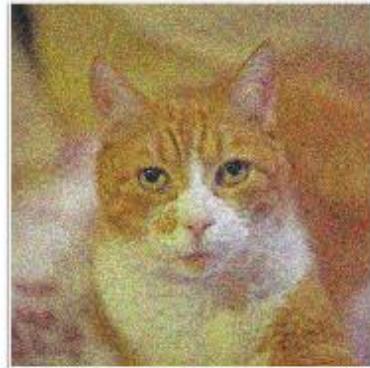


Now let's see how it is going to be:

I. A. Attacking a network with adversarial examples



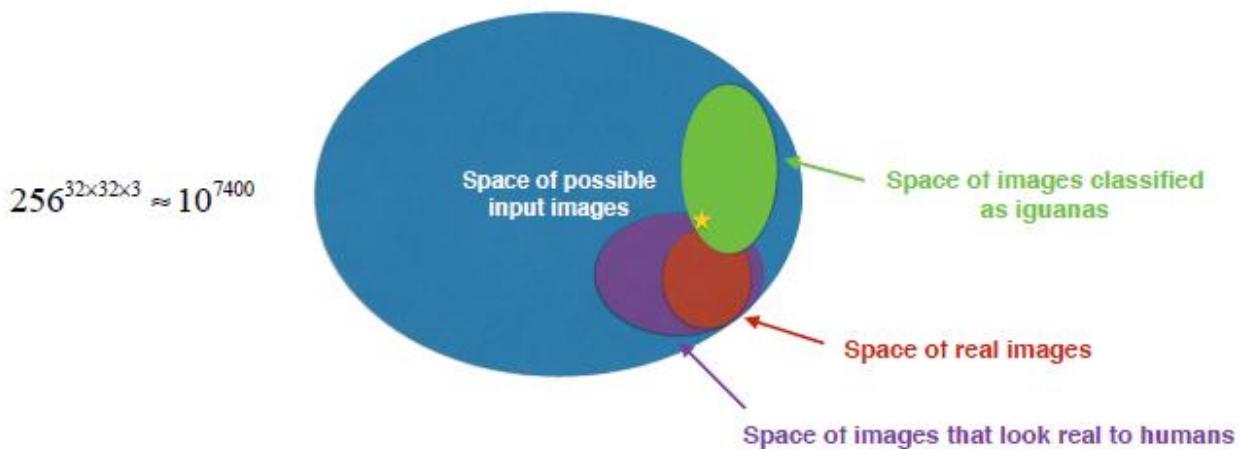
92% Cat



94% Iguana

Why is it happening? Because the space of the images that look real to human is much more than the space of real images, like the right one you see above; it is a blurry image that look real to human.

I. A. Attacking a network with adversarial examples



For defending against these adversarial attacks, we need to first warm up with some knowledge. The attacker may have access to the network and its hyper-parameters as well as parameters. So it is a white-box for the attacker. It is easier to perform white-box attack than black-box one. But how to attack in a black-box attack? One solution is that, if we can query the model, we can use numerical gradient, which means we change the picture and we look into the changes in the gradients. If we do this many times, we can understand how the loss works. But what if we cannot query the model? There is a very complex property of the adversarial attacks. They are **highly translatable**. To this end, if I make my own animal recognition network, then fool it with an adversarial example, it is very likely that the same example work for the black-box network.

One solution to defend against the adversarial attacks is to use a **SafetyNet**. It is a kind of firewall, before any image goes into the actual network, it goes through this firewall layer. But now, the attacker can see this firewall layer as a black-box network to fool again. Yes, it is possible, but now it has two constraints. But maybe if you fool the first one, the second one get fooled as well.

We can also generate adversarial examples and train on them so that they cannot get fooled with similar adversarial examples. But it is super costly and we don't know if we can generalize for the adversarial examples we never saw.

Finally, we can train the model with adversarial examples as well as the actual loss function. Its complexity is still super costly but, at least, it is online.

I. B. Defenses against adversarial examples

Menti

Knowledge of the attacker:

- White-box
- Black-box

Solution 1

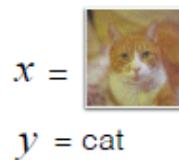
- Create a SafetyNet

Solution 2

- Train on correctly labelled adversarial examples

Solution 3

- Adversarial training $L_{new} = L(W, b, x, y) + \lambda L(W, b, x_{adv}, y)$



[Lu et al. (2017): SafetyNet: Detecting and Rejecting Adversarial Examples Robustly]
[Hanini Kannan et al. (2018): Adversarial Logit Pairing]

Kian Katanforoosh

Previously, researchers thought as we are training our model on true cats, due to high non-linearities of neural networks and overfitting on the training set, the adversarial examples exist. In other words, the network cannot understand what cats are, they only know what we trained on; so cannot generalize well. But it turns out to be wrong. Researchers found out, actually it is the linear part of the network leads to existence of adversarial examples. So, let's use a linear regression for this problem.

For the forward propagation of the linear regression we have: $y_{predicted} = w^T x + b$

So having 6 features, let's for this example say the network has been trained and converged to:

$$w = (1, 3, -1, 2, 2, 3), b = 0, x = [1, -1, 2, 0, 3, -2] \\ \Rightarrow y_{predicted} = -4$$

Now, the question is how to change x into x^* such that $y_{predicted}$ changes radically but $x \approx x^*$?

$$\frac{\delta y_{predicted}}{\delta x} = \text{impact of small changes of } x \text{ on } y_{predicted} \\ \frac{\delta y_{predicted}}{\delta x} = w^T$$

Now, if $x^* = x + \varepsilon w^T$ where ε is the value of the perturbation, then: $y_{predicted}^* = w x^* = w x + \varepsilon w w^T$
And we know: $w w^T = \|w\|^2$ which is always positive and shows that if I push x a little bit, it will pushes the $y_{predicted}^*$ to a larger value.

So let's say $\varepsilon = 0.2$, then the $x^* = [1.2, -0.4, 1.8, 0.4, 3.4, -1.4]$. Now, $y_{predicted}^* = 0.5$ while it was previously was -4. We just pushed x^* by 0.2 and it pushed the result 4.5.

Insights:

1. if w is large then $x \approx x^*$. So, for pushing it to the positive side, we take the $sign(w)$.
2. As x grows in dimension, the impact of $\varepsilon \times sign(w)$ on $y_{predicted}$ increases.

So the general form of generating adversarial, named “fast gradient sign method” is:

$$x^* = x + \varepsilon \times sign(\nabla_x J(w, x, y))$$

Generative Adversarial Networks

Let's now go to the Generative Adversarial Networks (GANs).

II.A - Motivation

Motivation:

- Data synthesis
- Compress and reconstruct data.
- Find a mapping between spaces.
- Image in-painting

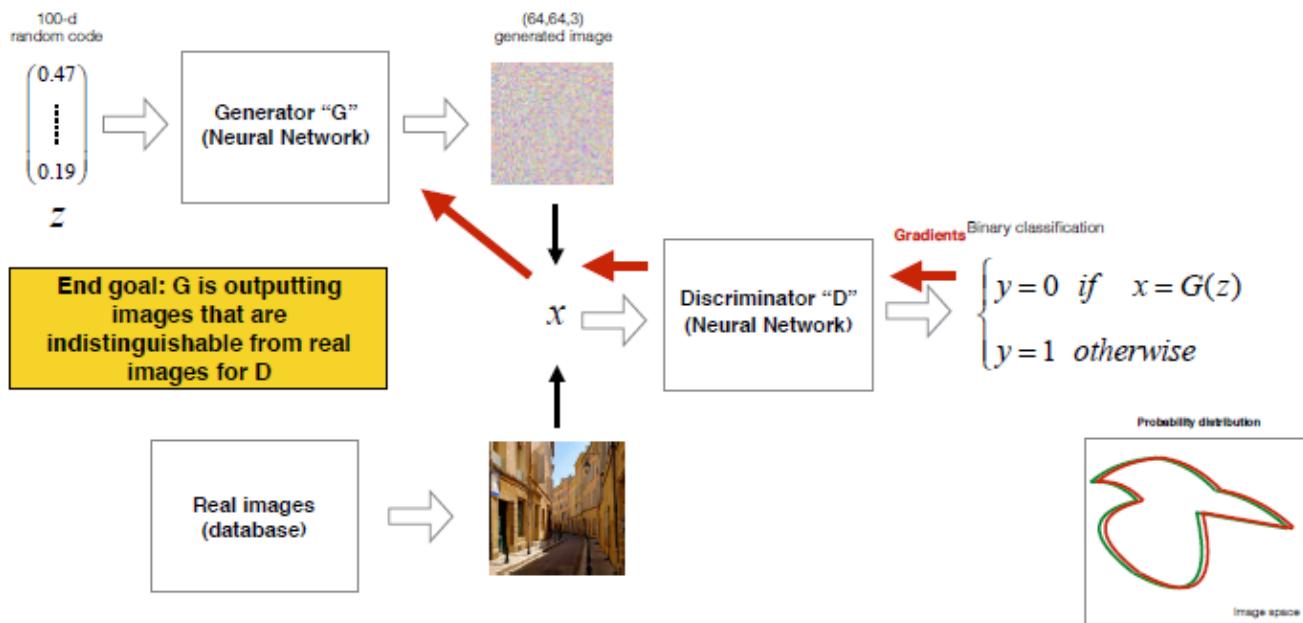
Approach: Collect a lot of data, use it to train a model to generate similar data from scratch.

Intuition: number of parameters of the model << amount of data

The goal here is to not give a cat or overfit to the cat, but train the network in a way that it generates a cat that doesn't exist. For doing so, we have a database of real images. The distribution of the real world images are shown in green and the generated images are shown in red. We want the distribution of generated images be look like the real world images.

We define a discriminator “ D ” whose job is to recognize an image is real or not – a binary classifier. So, it is a simple task to do for the “ D ” network. We want in the backward propagate gradients to the discriminator in order to train generator “ G ”. If the image was generated, then we have $G(z)$ and we can take the gradient and update weights of the generator, but if the image was real, then we do not have any $G(z)$ to take the gradient and it would be 0.

II.B - G/D Game



As for the loss function of the discriminator, we tweak the cross entropy a little bit so that the first term says “D” should correctly label real data as 1 and the second term says “D” should correctly label generated data as 0. So two terms are for two different mini-batches, one on the real images and the other on the “D”.

Now for the loss function of the generator, we want “G” to fool “D”. So, our whole idea would be whether “G” fools “D” or not. If “D” is bad, doesn’t make “G” good. So both should perform well.

II.B - G/D Game

Training procedure, we want to minimize:

Labels: $\begin{cases} y_{\text{real}} \text{ is always 1} \\ y_{\text{gen}} \text{ is always 0} \end{cases}$

- The cost of the discriminator

$$J^{(D)} = \underbrace{-\frac{1}{m_{\text{real}}} \sum_{i=1}^{m_{\text{real}}} y_{\text{real}}^{(i)} \log(D(x^{(i)}))}_{\text{cross-entropy 1: "D should correctly label real data as 1"}} - \underbrace{\frac{1}{m_{\text{gen}}} \sum_{i=1}^{m_{\text{gen}}} (1 - y_{\text{gen}}^{(i)}) \log(1 - D(G(z^{(i)})))}_{\text{cross-entropy 2: "D should correctly label generated data as 0"}}$$

- The cost of the generator

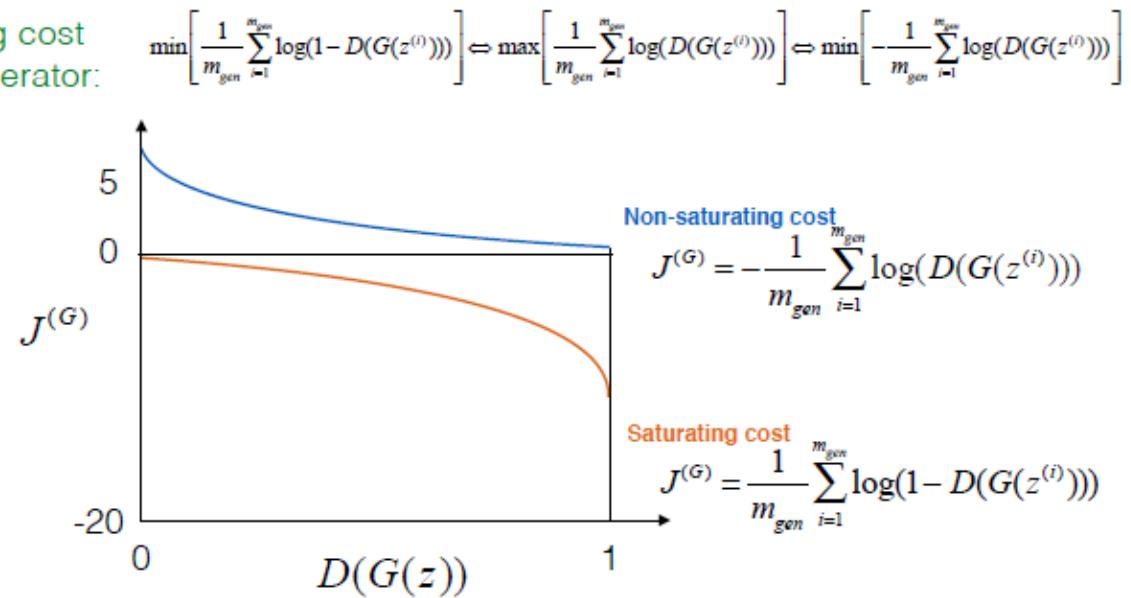
$$J^{(G)} = -J^{(D)} = \frac{1}{m_{\text{gen}}} \sum_{i=1}^{m_{\text{gen}}} \log(1 - D(G(z^{(i)}))) \quad \text{"G should try to fool D: by minimizing the opposite of what D is trying to minimize"}$$

Now, when $D(G(z)) = 1$, then we already fooled the network, but while it is getting closer to fool the network, it goes toward $-\infty$. So a better loss function would be the one with only $\log(D(G(z)))$. Now, the question is, how we should manipulate the loss function above in a way that it uses $\log(D(G(z)))$ instead of $1 - \log(D(G(z)))$? For doing so, we have this formulation:

$$\min \log(1 - x) = \max \log(x) = \min(-\log x)$$

II.C - Training GANs

Saturating cost
for the generator:



[Ian Goodfellow (2014): NIPS Tutorial: GANs]

Kian Katanforoosh

Then we rewrite it as:

II.C - Training GANs

Note that: $\min \left[\frac{1}{m_{gen}} \sum_{i=1}^{m_{gen}} \log(1 - D(G(z^{(i)}))) \right] \Leftrightarrow \max \left[\frac{1}{m_{gen}} \sum_{i=1}^{m_{gen}} \log(D(G(z^{(i)}))) \right] \Leftrightarrow \min \left[-\frac{1}{m_{gen}} \sum_{i=1}^{m_{gen}} \log(D(G(z^{(i)}))) \right]$

New training procedure, we want to minimize:

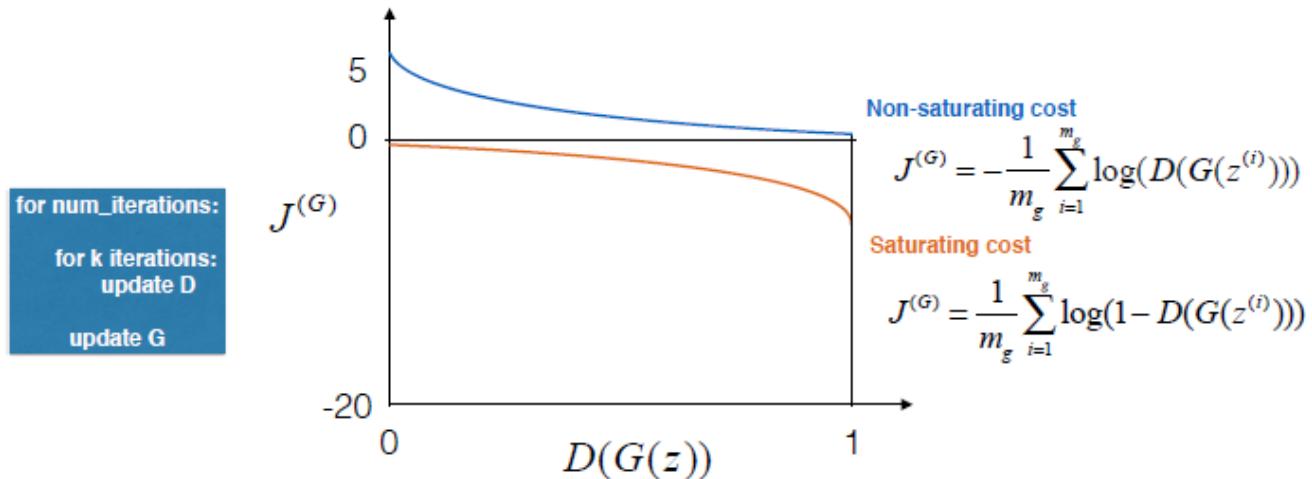
$$J^{(D)} = \underbrace{-\frac{1}{m_{real}} \sum_{i=1}^{m_{real}} y_{real}^{(i)} \cdot \log(D(x^{(i)}))}_{\text{cross-entropy 1: "D should correctly label real data as 1"}} + \underbrace{-\frac{1}{m_{gen}} \sum_{i=1}^{m_{gen}} (1 - y_{gen}^{(i)}) \cdot \log(1 - D(G(z^{(i)})))}_{\text{cross-entropy 2: "D should correctly label generated data as 0"}}$$

$$J^{(G)} = -\frac{1}{m_{gen}} \sum_{i=1}^{m_{gen}} \log(D(G(z^{(i)}))) \quad \text{"G should try to fool D: by minimizing this"}$$

Another trick to train GANs is to use the fact that “D” is usually easier to train than “G”. In addition, as “D” improves, “G” can improve and if “D” cannot improve, “G” cannot as well. So the performance of “D” is an upper bound to what “G” can achieve. So we have to train “D” more than “G”.

II.C - Training GANs

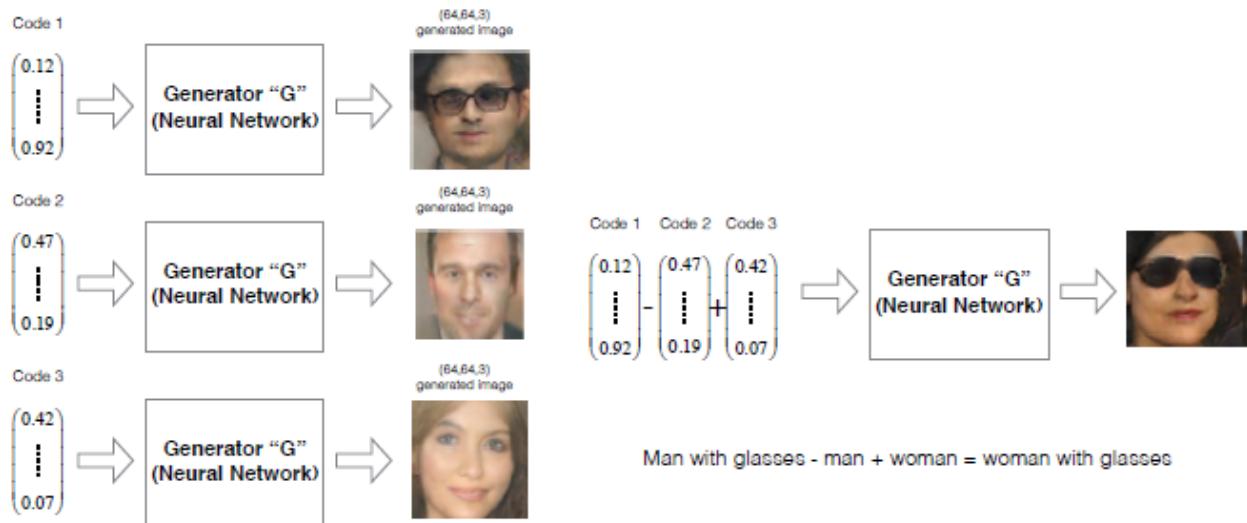
Simultaneously training G/D?



Something funny is that linear operation in latent space of codes have impact directly on the image space.

II.E - Nice results

Operation on codes



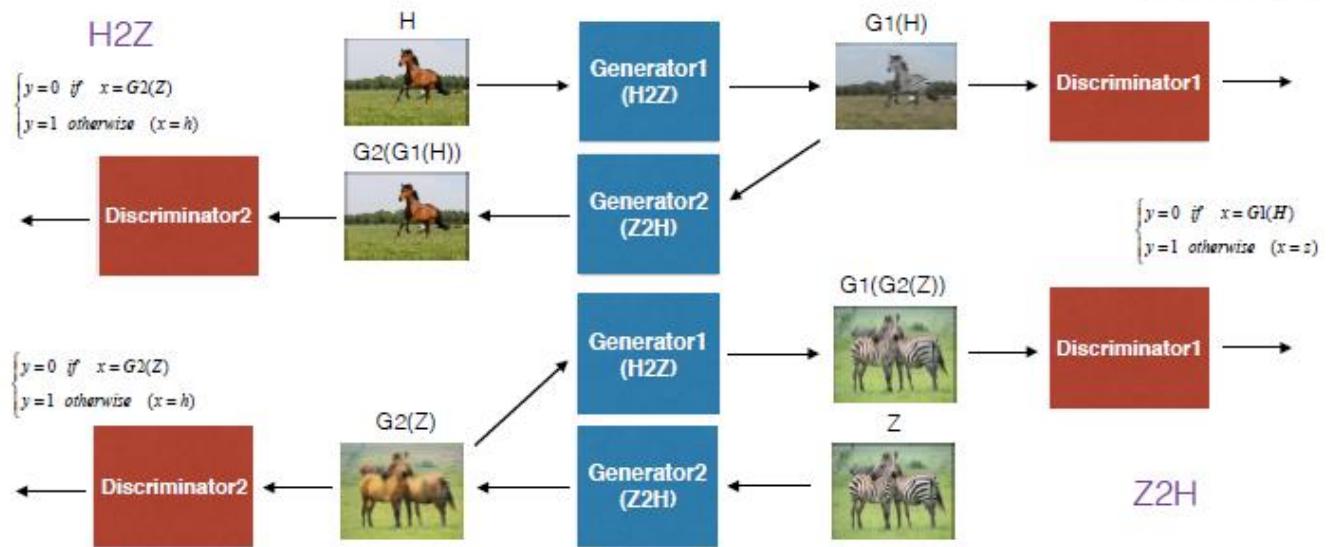
[Radford et al. (2015): UNSUPERVISED REPRESENTATION LEARNING WITH DEEP CONVOLUTIONAL GENERATIVE ADVERSARIAL NETWORKS]

Kian Katanforoosh

And let's see a **CycleGANs** as well:

II.E - Nice results

Architecture?



[Zhu, Park et al. (2017): Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks]

Kian Katanforoosh

II.E - Nice results

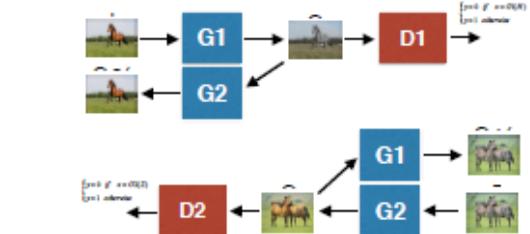
Loss to minimize?

$$J^{(D1)} = -\frac{1}{m_{\text{real}}} \sum_{i=1}^{m_{\text{real}}} \log(D1(z^{(i)})) - \frac{1}{m_{\text{gen}}} \sum_{i=1}^{m_{\text{gen}}} \log(1 - D1(G1(H^{(i)})))$$

$$J^{(G1)} = -\frac{1}{m_{\text{gen}}} \sum_{i=1}^{m_{\text{gen}}} \log(D1(G1(H^{(i)})))$$

$$J^{(D2)} = -\frac{1}{m_{\text{real}}} \sum_{i=1}^{m_{\text{real}}} \log(D2(h^{(i)})) - \frac{1}{m_{\text{gen}}} \sum_{i=1}^{m_{\text{gen}}} \log(1 - D2(G2(Z^{(i)})))$$

$$J^{(G2)} = -\frac{1}{m_{\text{gen}}} \sum_{i=1}^{m_{\text{gen}}} \log(D2(G2(Z^{(i)})))$$



$$J = J^{(D1)} + J^{(G1)} + J^{(D2)} + J^{(G2)} + \lambda J^{\text{cycle}}$$

Course 2

*Improving Deep Neural Networks: Hyperparameter tuning,
Regularization and Optimization*

Week 1

Practical aspects of Deep Learning

Train / Dev / Test sets

Some of the hyper-parameters For NNs are:

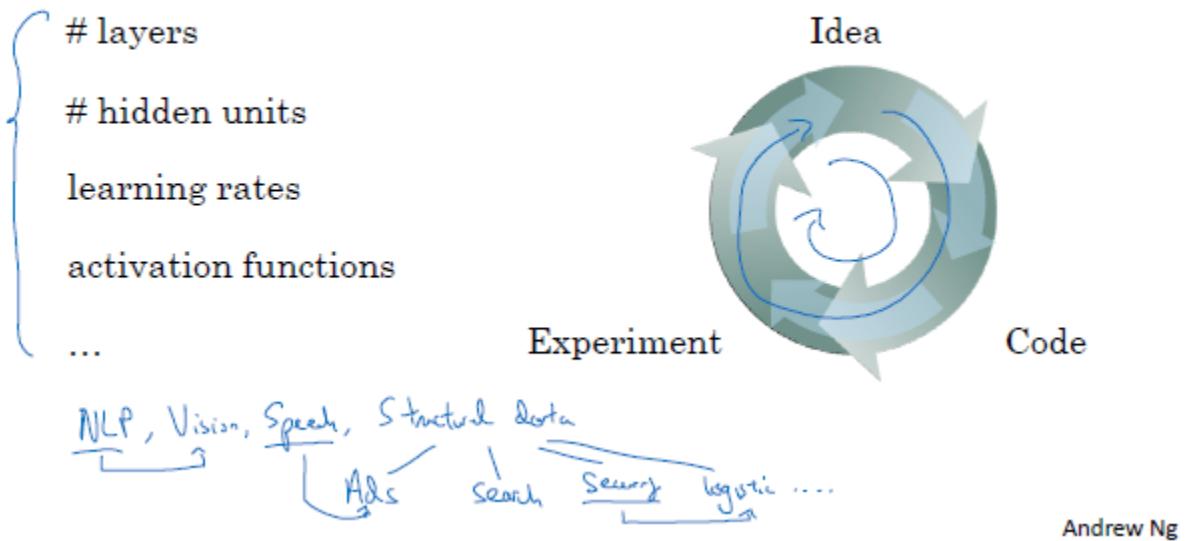
1. # layers
2. # hidden units
3. Value of the learning rate
4. Choice of the activation functions
5. ...

In practical, applied ML is a highly iterative process, so we start with:

1. Idea → number of hidden units, number of layers, activation functions, and so forth
2. Code → we program it and we run it
3. Experiment → we get the results of the experiments and based on the outcome we refine the idea → getting back to step 1

Finding the best choice of network and hyper-parameters is impossible unless test it over and over again

Applied ML is a highly iterative process



To increase the pace of iterating, we divide the dataset into train, dev, and test sets.

A data should have:

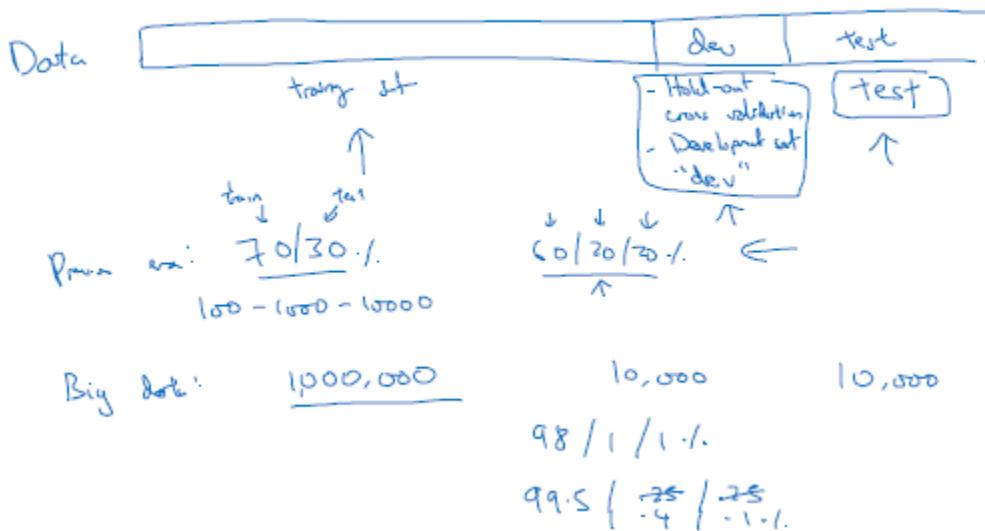
Training Set + Hold-out Cross Validation or Development Set + Test Set

So we run the learning algorithms on the training set and then test it with dev set to find out the best model for the problem. Finally, after finding the best model, we evaluate it with the test set.

Previously: 70/30 train and test sets without dev set or 60/20/20 train/dev/test sets

Modern big data era: dev and test should be so much smaller in percentage since dev is only for testing the models and finding the best and for test set it is only about validating how well the model is doing. So we allocate the most of the data to the training set so that it can generalize better.

Train/dev/test sets



There is a concept about mismatched train and test distributions. For example, we may get cat pictures from web with high resolutions for variety of cat models. Then, our dev and test sets contain user images with low resolution and limited types of cats. So, they are not from the same distributions. What may happen here is that even if our model learns the best on the training set, it cannot perform well on dev and test sets because the type of data is not uniform.

So, we have to make sure the dev/test set and training set are from the same distribution

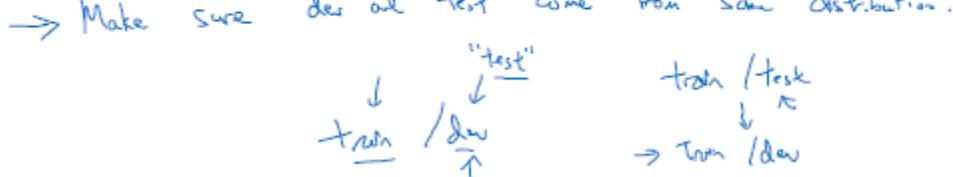
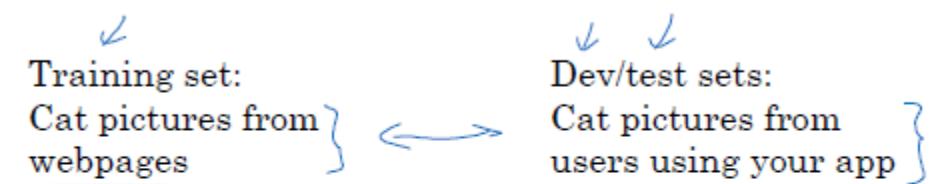
Not having test set might be okay → dev set should be there at any event → evaluation on dev set is mandatory to find the good model

Though in this setting, the result is not unbiased estimate of performance → so it is better to have test set as well

In this setting, sometimes people call dev set as test set.

Mismatched train/test distribution

Certs

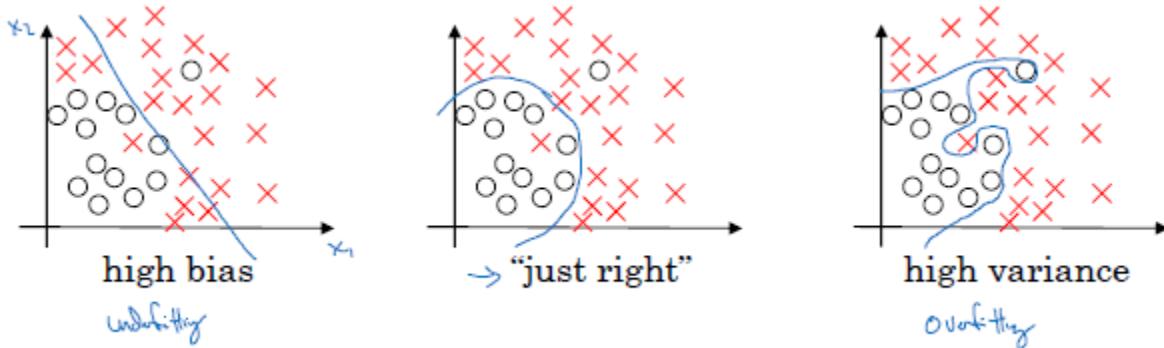


Not having a test set might be okay. (Only dev set.)

Bias / Variance

There is not much about bias/variance trade-off in deep learning.

The leftmost image draws a line as the classifier, so it has no flexibility → high bias → underfitting
The rightmost image is too flexible and is not good for generalization → high variance → overfitting
The middle one is in between of previous two cases so it is right.



Let's see some examples to grasp bias and variance concept better. We use the assumption that human can do the same task with around 0% error.

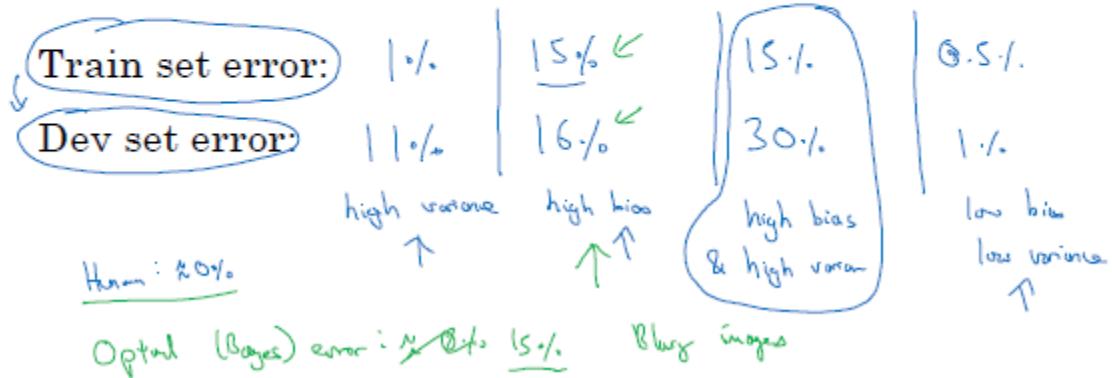
1. High variance
 - a. Training error → 1% → very well
 - b. Dev error → 11% → relatively poorly
 - i. The model is probably overfitting
2. High bias
 - a. Training error → 15% → not fitting
 - b. Dev error → 16% → very well
 - i. The model is probably underfitting
3. High bias and high variance
 - a. Training error → 15% → not fitting
 - b. Dev error → 30% → poor
 - i. Worst of two worlds
4. Low bias and low variance
 - a. Training error → 0.5% → very well fitting
 - b. Dev error → 1% → very well model
 - i. Best of two worlds

We started with the assumption that human level error or even better, the optimal Bayes error, is 0%. If it was 15% we could have said the 2nd model is the best. Then the 4th model is not acceptable because we cannot outperform the optimal Bayes error.

Rule of thumb → compare the error difference of the Train set and Dev set

Bias and Variance

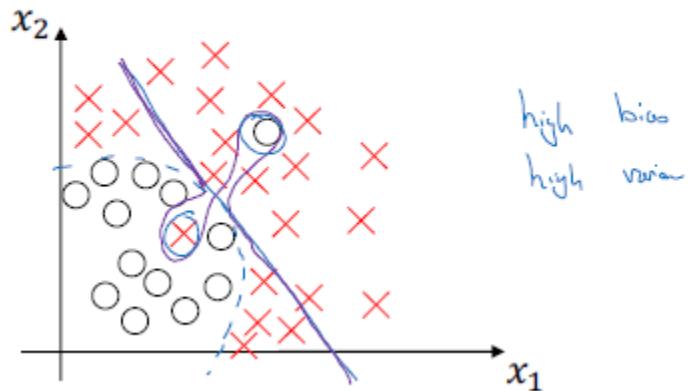
Cat classification



What is High Bias and High Variance look like?

Not fitting the quadratic shape and with too much flexibility in middle or some parts of it
⇒ High bias at some regions and high variance at some other regions.

High bias and high variance

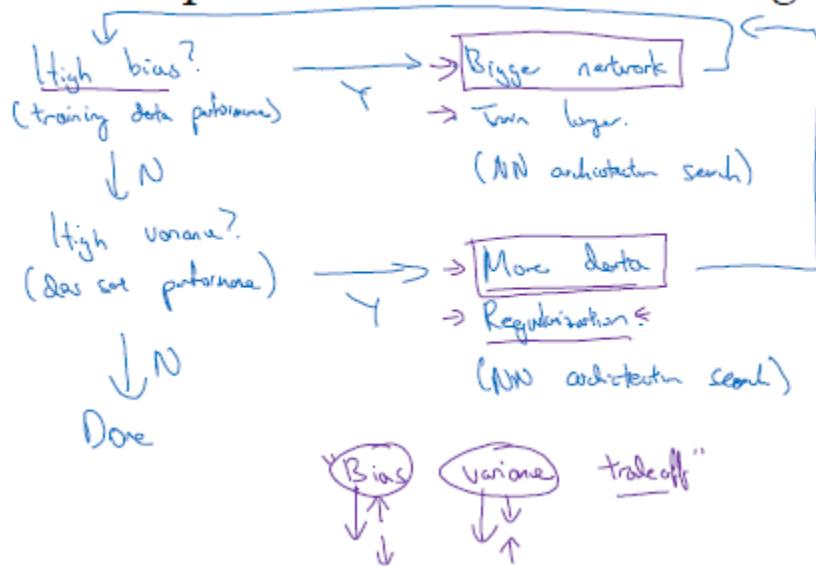


Basic Recipe for Machine Learning

1. Is the algorithm **high bias**? (we look at training data performance)
 - a. Yes?
 - i. Bigger Network
 - ii. Different NN architecture may help
 - iii. Running the model longer
 - b. No?
 - i. Go to next step
2. Is the model high variance? (we look at the error difference of training and dev set)
 - a. Yes?
 - i. Getting more data
 - ii. Trying Regularization
 - iii. Another NN architecture
 1. Going back to step 1
 - b. No?
 - i. We are done!

With bigger network and more data we can reduce bias and variance without hurting the other one.

Basic recipe for machine learning



Regularization

When we have a high variance, we first going to use a regularization, if no more data is available
We can regularize “bias” value as well, but as it is only a single number, it does not change anything much so we can easily omit it. However, “weights” are pretty high dimensional data.

L2 regularizer is the most well-known regularizer, but we also can use L1 regularizer.

L1 regularizer causes the “weights” to become sparse with many zeros → Can be used in compressing the model because of too many zeros

We also use a λ hyper-parameter along with L1 and L2 regularizers.

Logistic regression

$$\min_{w,b} J(w,b) \quad w \in \mathbb{R}^n, b \in \mathbb{R} \quad \lambda = \text{regularization parameter}$$

$$J(w,b) = \underbrace{\frac{1}{m} \sum_{i=1}^m l(y^{(i)}, h^{(i)})}_{\text{L2 regularization}} + \underbrace{\frac{\lambda}{2m} \|w\|_2^2}_{\text{L1 regularization}}$$

~~+ $b^2 / 2m$~~
~~omit~~

$$\|w\|_2^2 = \sum_{j=1}^n w_j^2 = w^T w \leftarrow$$

$$\|w\|_1 = \frac{\lambda}{2m} \sum_{j=1}^n |w_j| = \frac{\lambda}{2m} \|w\|_1, \quad w \text{ will be sparse}$$

Now for the Neural Network, the concept is still the same. Note that, the weight update rule changes as we add regularization term since its derivative also comes in the formula. L2 regularizer sometimes being called “weight decay” as it makes the “weight” update gets a little smaller by multiplying by a number less than 1.

Neural network

$$\rightarrow J(w^{(0)}, b^{(0)}, \dots, w^{(L)}, b^{(L)}) = \underbrace{\frac{1}{m} \sum_{i=1}^m l(y^{(i)}, h^{(i)})}_{\text{“Frobenius norm”}} + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{(l)}\|_F^2$$

$$\|w^{(l)}\|_F^2 = \sum_{i=1}^{n^{(l)}} \sum_{j=1}^{n^{(l+1)}} (w_{ij}^{(l)})^2 \quad w^{(l)}: (n^{(l)}, n^{(l+1)})$$

$$\frac{\partial J}{\partial w^{(l)}} = \underbrace{(\text{from backprop}) + \frac{\lambda}{m} w^{(l)}}_{\|w^{(l)}\|_F^2} \quad \frac{\partial J}{\partial w^{(l)}} = \partial w^{(l)}$$

$$\rightarrow w^{(l)} := w^{(l)} - \alpha \frac{\partial J}{\partial w^{(l)}} \quad w^{(l)} := w^{(l)} - \alpha \left[(\text{from backprop}) + \frac{\lambda}{m} w^{(l)} \right]$$

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

$$= \underbrace{\left(1 - \frac{\alpha \lambda}{m}\right) w^{(l)}}_{\leq 1} - \alpha (\text{from backprop})$$

“Weight decay”

Andrew Ng

Additional Readings:

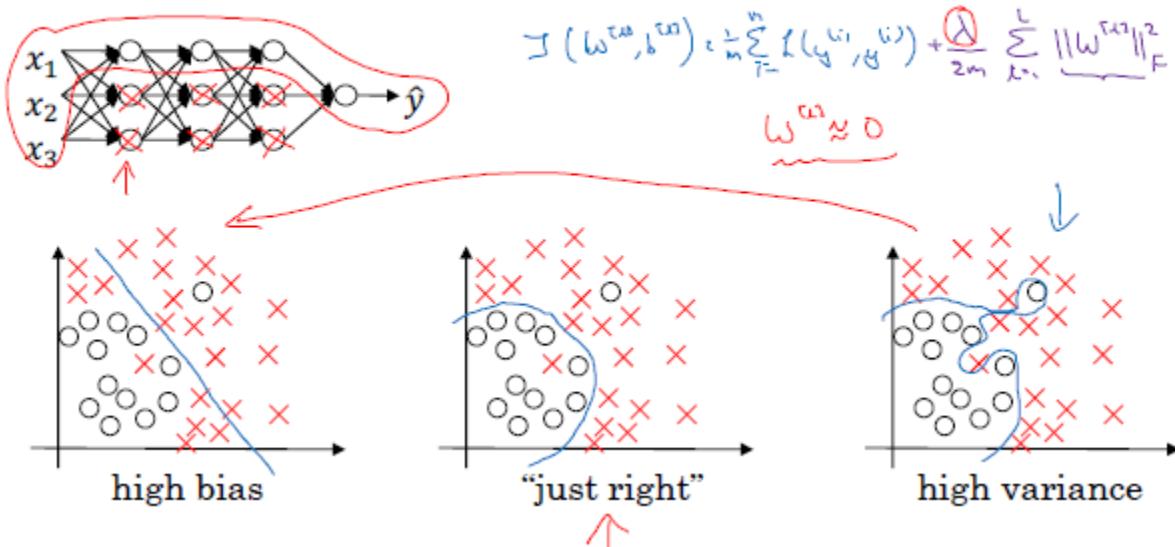
<https://www.quora.com/What-is-the-difference-between-L1-and-L2-regularization-How-does-it-solve-the-problem-of-overfitting-Which-regularizer-to-use-and-when>

<https://www.quora.com/What-are-the-Frobenius-L2-and-Euclidian-norms-How-are-they-different-from-one-another>

Why regularization reduces overfitting?

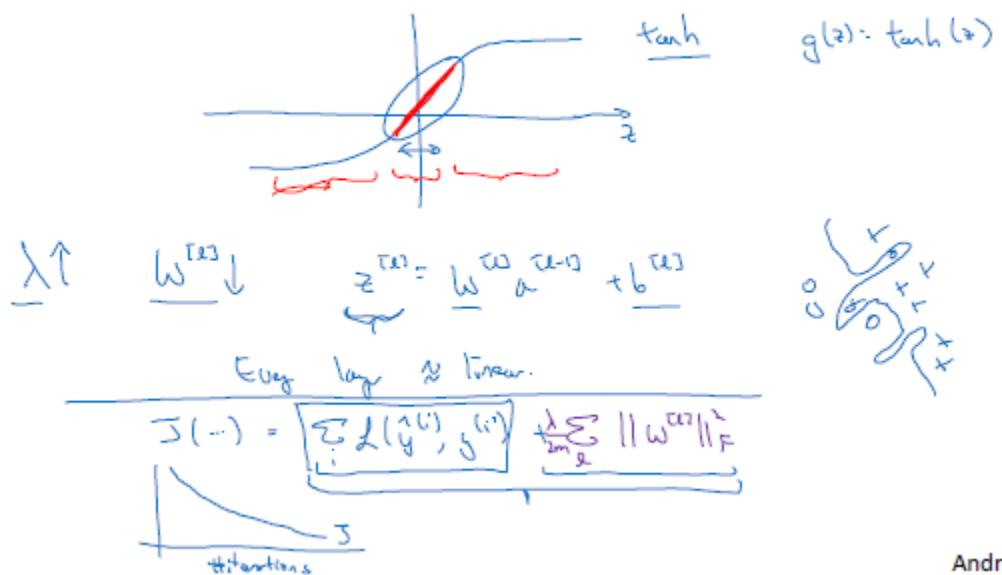
One intuition for using regularization is that it zeroing out many hidden units and make it so much simpler so it goes from high variance to lower variance. However, it is not quite right. It actually reduces the impact of some of these hidden units and prevents from overfitting.

How does regularization prevent overfitting?



For another intuition we use tanh. In tanh activation function, with a large value of λ , we are going to have smaller weights. With smaller weights, the "Z" gets smaller as well. Now, small "Z" values have $G(Z)$ close to 0, the part that would be roughly linear and as we go far from the zero it does not grow very fast. So it is calculating a linear like function which is so simple and less prone to overfitting.

How does regularization prevent overfitting?

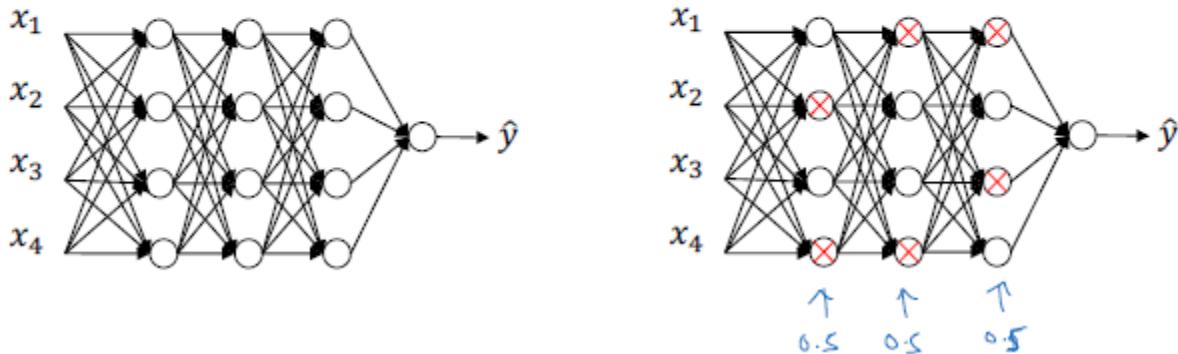


Andrew Ng

Dropout Regularization

In the dropout, we give a probability to remove or keep some nodes in NN \rightarrow we end up with much smaller and diminished network. We then do it over and over again to create other diminished networks.

Dropout regularization



A concrete approach of dropout implementation is called “inverted dropout”. We define a “keep-prob” as a value which is the chance of keeping a node. We randomly initialize a matrix for each layer which is either 0 for eliminating a node or 1 for keeping the node. Then we multiply this matrix with the activation matrix. Then we need to scale up the result by “keep-prob” to keep the expected value of the matrix as before. The last step is for “inverted dropout” which helps us a lot in the test time.

Implementing dropout (“Inverted dropout”)

Illustrate with layer $l=3$. $\text{keep-prob} = \frac{0.8}{5} = 0.2$

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

$$\underbrace{a_3}_{\uparrow \uparrow} = \text{np.multiply}(a_3, d_3) \quad \# a_3 \neq d_3.$$

$$\rightarrow a_3 \leftarrow \underbrace{a_3 / \text{keep-prob}}_{\substack{50 \text{ units. } \rightsquigarrow 10 \text{ units shut off}}}$$

$$z^{(4)} = w^{(4)} \cdot \underbrace{\frac{a_3}{5}}_{\substack{\text{Test} \\ \downarrow \text{reduced by } 20\%}} + b^{(4)}$$

$$/ \approx 0.8$$

In the test time, we never use dropout. Otherwise we need to run dropout many times. If we don't use inverted dropout, we need to scale down the results we get in the test time for each epoch which is much harder.

Making predictions at test time

$$a^{(0)} = X$$

No drop out.

$$\uparrow z^{(1)} = w^{(1)} \underline{a^{(0)}} + b^{(1)}$$

$$a^{(1)} = g^{(1)}(z^{(1)})$$

$$z^{(2)} = w^{(2)} \underline{a^{(1)}} + b^{(2)}$$

$$a^{(2)} = \dots$$

$$\downarrow$$

$$\hat{y}$$

λ = keep-prob

Understanding Dropout

The intuition is that we can't rely on any one of the features as it might fade away, so have to spread out weights. It is like shrinking the weights as we previously did with L2 regularizer.

One note is that, we can choose different keep-prob for each layer. Mainly we use lower keep-prob values for those layers with many units and parameters that may cause the network to overfit.

For example, for layers having nodes more than 7×7 , it is more likely that this matrix tend to overfit, so we can dropout more here but in somewhere else with less possibility of overfitting we can eliminate less units.

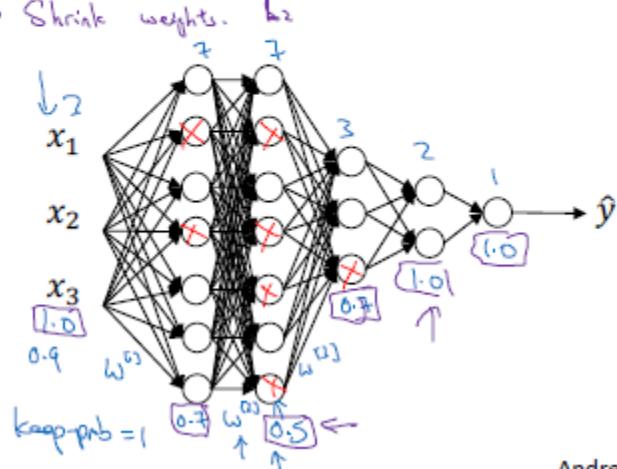
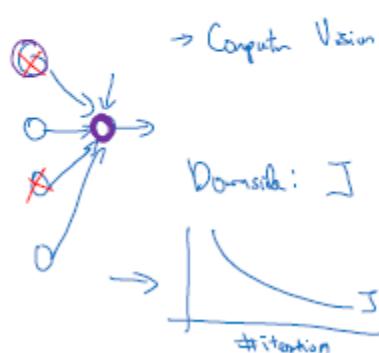
Note that for the input and output layers it is better to keep the keep-prob as 1.

We mainly use the dropout in computer vision as we have many layers in this application.

The downside of the dropout is that the cost function is no longer well-defined as we are eliminating some nodes randomly at each iteration. So for making sure that the cost function monotonically decreases we have to turn off the dropout.

Why does drop-out work?

Intuition: Can't rely on any one feature, so have to spread out weights. \rightarrow Shrink weights.

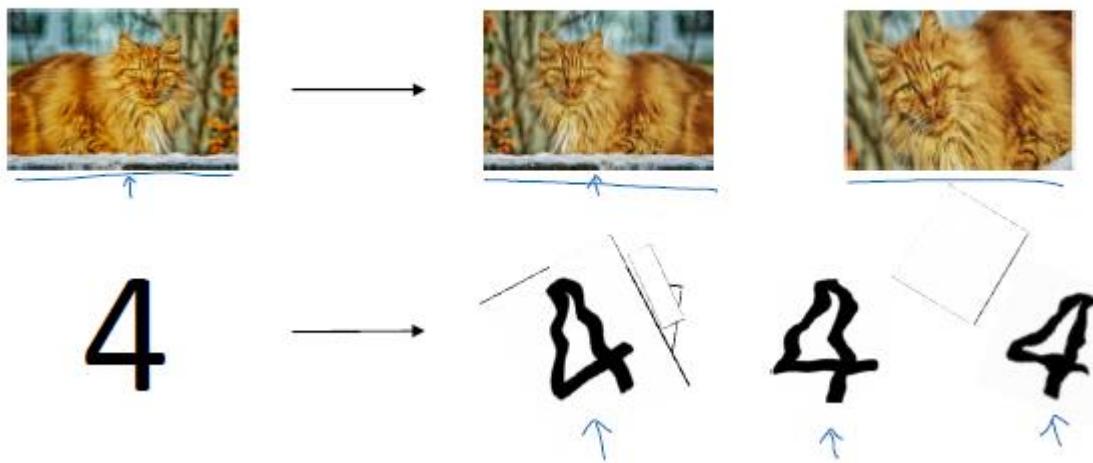


Other regularization methods

Data augmentation: We may use dropout or any kind of regularizer as we do not have data to enlarge training set. In this case, we can use data augmentation to gain more data. Although the new data set is not as good as having brand new independent examples.

1. flipping a picture → gaining 2X samples
2. Zooming into part of the image
3. Adding distortion is also practical

Data augmentation



Early stopping: As iterations go on, the training error and cost function decreases, but for dev set error at some point starts to increase. At very first the weights are close to zero and very small. As we train, the weights get bigger and bigger. So, similar to regularization, the early stopping keep the weights at the mid size. like a L2 that scale the values, it keeps the early values of W.

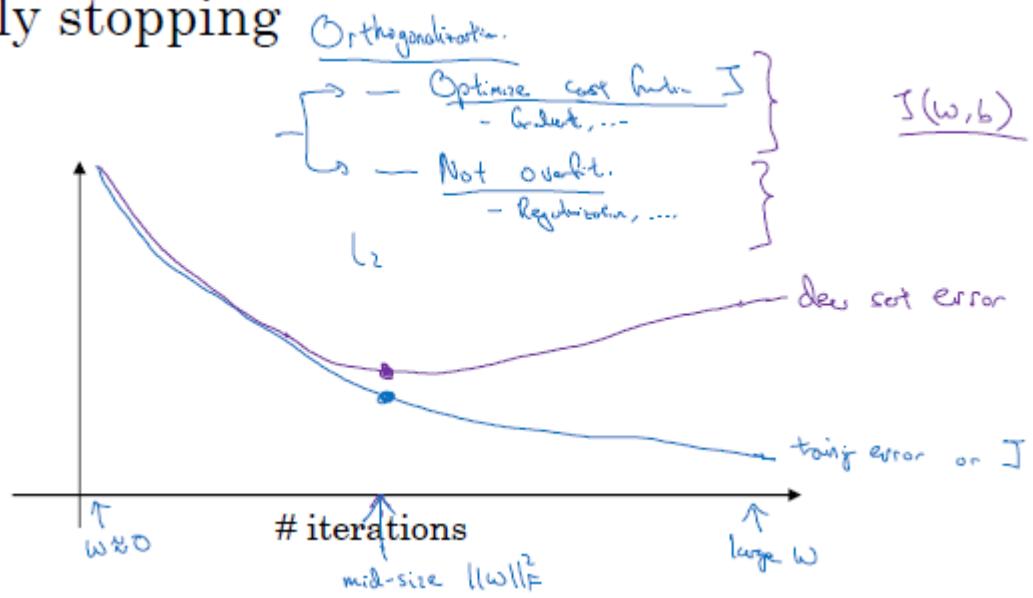
Its downside:

We want to optimize the cost function with an optimizer like Adam or gradient descent. We also don't want to overfit by using a technique like regularizer. So they are many hyper-parameters to tune and each one of them are for independent tasks, namely orthogonalizing. Early stopping bounds these tasks into one so we cannot understand exactly which of the chosen parameters are not quite right.

Its benefit:

It reduces the search space because we don't need to check all the possible hyper-parameter values. Moreover, by running it once, we check weights with small, mid, and large size.

Early stopping



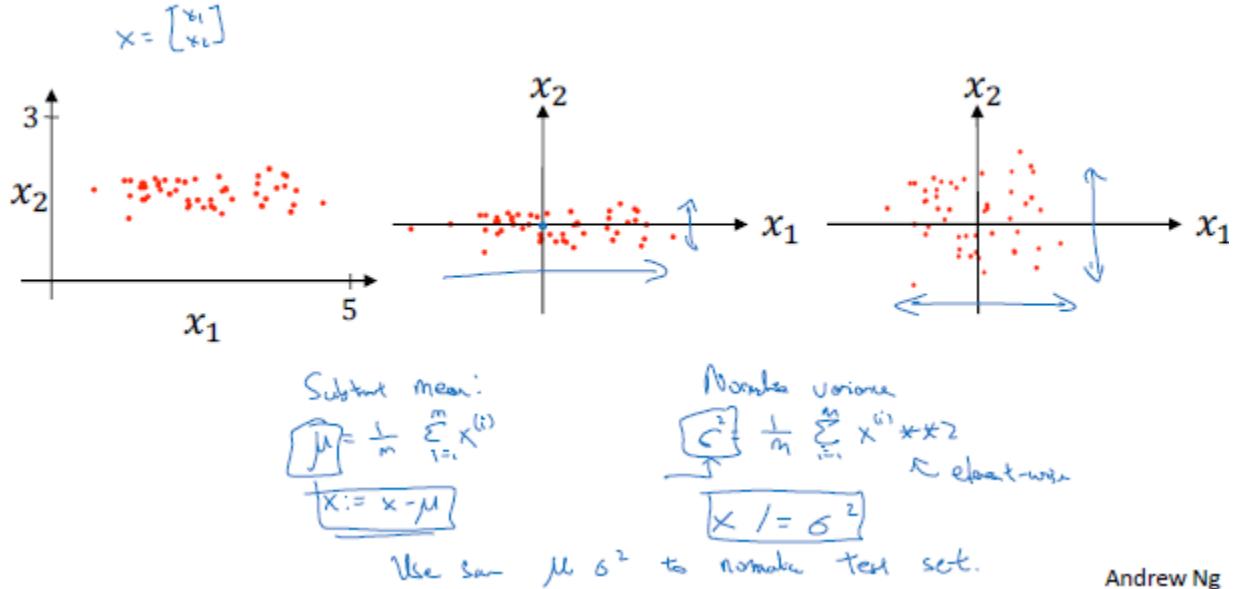
Normalizing inputs

One technique to increase the pace of the iteration while learning algorithm is working is to normalize the input data. For normalizing input data, we have to take 2 steps:

- 1- zeroing out the mean → subtracting data from the mean value of the data
- 2- normalizing the variance → so that data concentrated around the 0,0

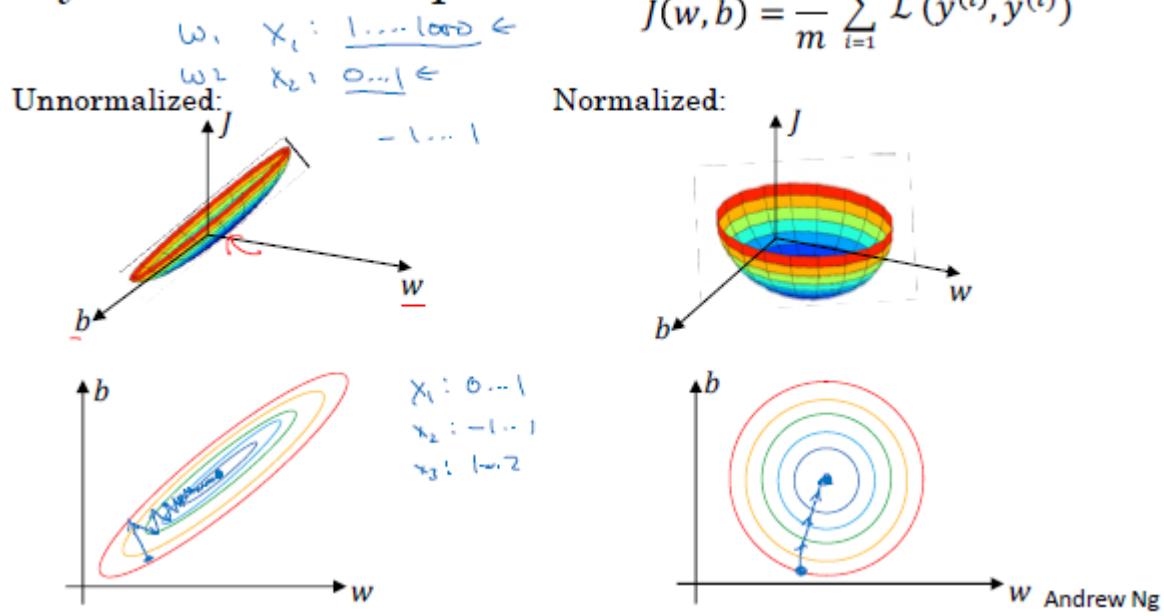
Note that we use the same variance and mean for the test set as well.

Normalizing training sets



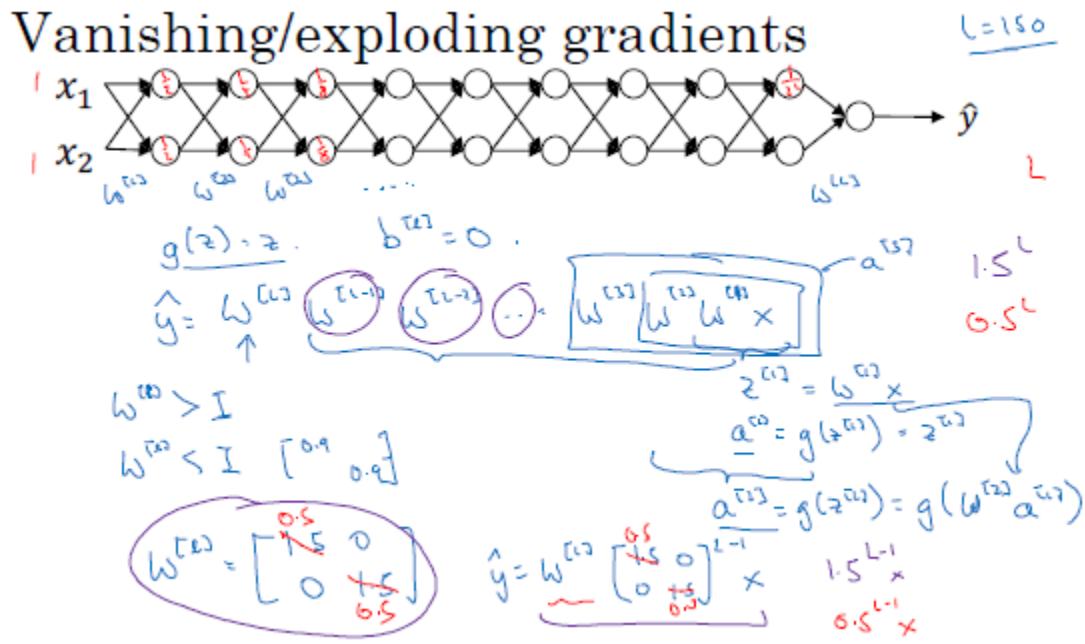
Now the question is, how does it affect the training pace. Non-normalized data is elongated and normalized is more symmetric. For the gradient descent in the former case takes small steps because need many steps to reach the minimum but in the latter case larger steps works very well and learning algorithm runs faster. So we aim to bring them to similar scale.

Why normalize inputs?



Vanishing / Exploding gradients

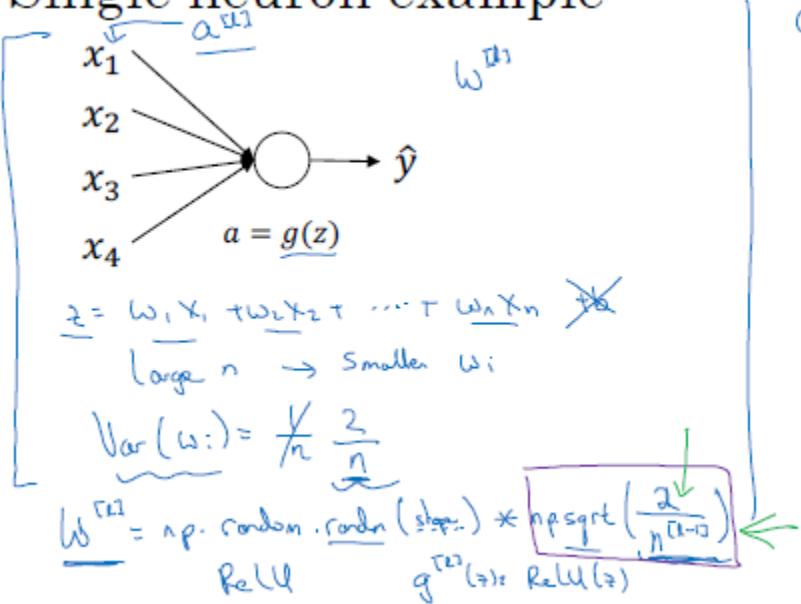
A very huge problem in very deep networks is vanishing and exploding gradients. Let's see its math below to understand how vanishing and exploding gradients is happening. Now, starting with a weight value larger than 1, either using an activation function or assuming it linear, in the long run, they the value larger than 1 goes to the power of a big value which ends up with a very huge value. In the contrast, using weight with value less than 1, when it goes to the power of a large value, it ends up with a very small value, very close to 0. Specially with very small values, the optimizer should take very small steps which reduces the pace of training.



Weight Initialization for Deep Networks

Initializing the weights showed that helps a lot in preventing vanishing and exploding gradients. For weight initialization, we see that as the number of features increase, we need smaller weights because as they are summing up, it doesn't cause explosion of gradient. So, we set the variance of the weight equal to 1 or 2 over number of features at each layer. We can see different initialization methods for different activation functions.

Single neuron example



Other variants:

Tanh

Xavier initialization

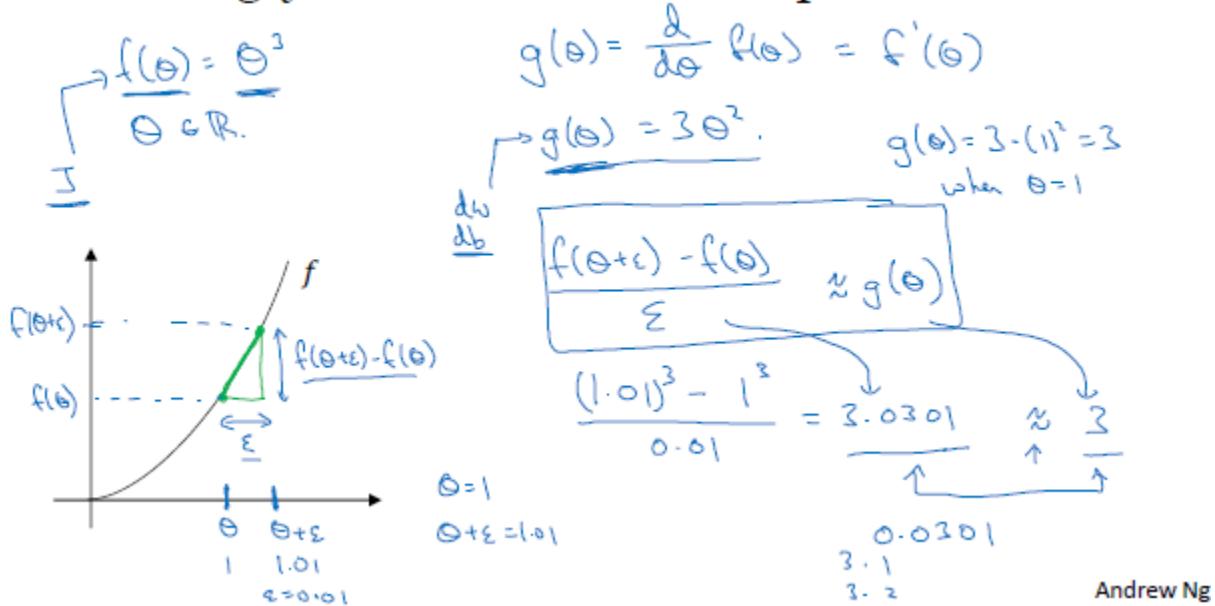
$$\frac{2}{n^{(L-1)} + n}$$

Numerical approximation of gradients

When we implement back propagation, we need to make sure the implementation works fine. To this end, we use a tool, named “numerical checking”. Let’s see how it works.

First of all, for taking derivative of a function, we define a very small value, ε , in which while we are calculating the slope of the function, we push the ε to 0 that is the mathematical definition of derivative. Let’s see its result:

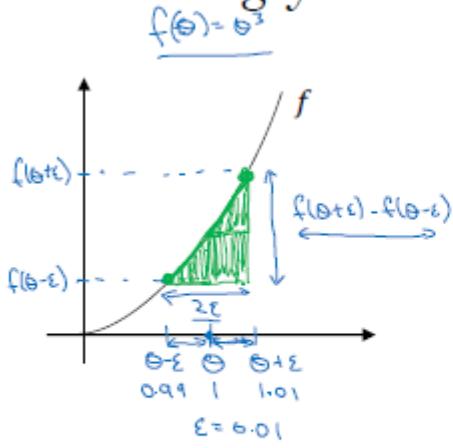
Checking your derivative computation



Now we use ε both on right and left side of the actual θ . By doing so, we make a bigger triangle. Now if take the derivative for $\theta \pm \varepsilon$ equation, we should get a value very closer to actual derivative of function than the one side we saw before. Although it is much more accurate, but twice slower than the previous one.

We can see accuracy of using two side difference is in the order of $O(\varepsilon^2)$ while the one side difference gives the order of $O(\varepsilon)$.

Checking your derivative computation



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

$$\frac{(1.01)^3 - (0.99)^3}{2(0.01)} = 3.0001 \approx 3$$

$$g(\theta) = 3\theta^2 = 3$$

approx error: 0.0001
(prev slide: 3.0301, error: 0.03)

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

$\frac{\cancel{0}(\epsilon^2)}{\cancel{0.01}} = \underline{\underline{0.0001}}$	$\frac{f(\theta+\epsilon) - f(\theta)}{\epsilon}$ error: $\mathcal{O}(\epsilon)$ $\uparrow \quad \uparrow$ 0.01
---	--

Andrew Ng

Gradient checking

We use gradient checking to make sure our back-propagation implementation is right. For using it first we concatenate all the weights and biases into a very big matrix. Then we take the derivative of all weights and biases, then again concatenate them into a very giant matrix. Now let's see whether or not the derivative we performed is the gradient of cost function $J(\theta)$.

Gradient check for a neural network

Take $\underline{W^{[1]}, b^{[1]}, \dots, W^{[L]}, b^{[L]}}$ and reshape into a big vector $\underline{\theta}$.

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

Take $\underline{dW^{[1]}, db^{[1]}, \dots, dW^{[L]}, db^{[L]}}$ and reshape into a big vector $\underline{d\theta}$.

concatenate

Is $d\theta$ the gradient of $J(\underline{\theta})$?

Now we apply the idea we got from numerical approximation, we take two side difference of the of the giant matrix $d(\theta)$. Then we take the distance between this approximated and actual result. It should get a value less than or around 10^{-7} , it is great. But if it is larger than or around 10^{-3} , we should get worried.

Gradient checking (Grad check)

$$J(\underline{\theta}) = J(\underline{\theta_0, \theta_1, \dots})$$

for each i :

$$\rightarrow \underline{d\theta_{approx}[i]} = \frac{J(\underline{\theta_0, \theta_1, \dots, \theta_i + \varepsilon, \dots}) - J(\underline{\theta_0, \theta_1, \dots, \theta_i - \varepsilon, \dots})}{2\varepsilon}$$

$$\approx \underline{d\theta[i]} = \frac{\partial J}{\partial \theta_i}$$

$$d\theta_{approx} \approx d\theta$$

Check

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

$$\varepsilon = 10^{-7}$$

$$\approx \boxed{10^{-7} - \text{great!}} \leftarrow$$

$$\rightarrow 10^{-3} - \text{worry.} \leftarrow$$

Gradient Checking Implementation Notes

The first bullet, suggest that we should not use grad check in the whole training process as it is so slow. The second bullet, provides the fact that with grad checking, we can identify whether the problem is from the gradient of the bias terms or weights.

The third bullet, highlights that if we want to use grad check and we have a regularizer, we should take it into account in the grad check formulation.

The fourth bullet, warns that we should not use grad check while we are using dropout. We can do each one of them separately.

The last bullet is a hint about when to use grad check. It is the best to use at the random initialization step and after some epochs because if the weights and biases get large values, even if we implemented everything correctly, the difference of approximation and actual gradient will be larger than expected.

Gradient checking implementation notes

- Don't use in training – only to debug

$$\frac{\partial \text{approx}[t]}{\partial \theta} \leftrightarrow \frac{\partial \text{gt}}{\partial \theta}$$

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

$$\frac{\partial b^{(l)}}{\partial \theta} \quad \frac{\partial w^{(l)}}{\partial \theta}$$

$$J(\theta) = \frac{1}{n} \sum_i \ell(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2n} \sum_j \|w^{(j)}\|_F^2$$

$\frac{\partial J}{\partial \theta} = \text{grad of } J \text{ wrt. } \theta$

- Remember regularization.

- Doesn't work with dropout.

J

keep-prob = 1.0

- Run at random initialization; perhaps again after some training.

$$w, b \approx 0$$

Andrew Ng

Assignment Notes

In general, initializing all the weights to zero results in the network failing to break symmetry. This means that every neuron in each layer will learn the same thing, and you might as well be training a neural network with $n[l]=1$ for every layer, and the network is no more powerful than a linear classifier such as logistic regression.

Poor initialization can lead to vanishing/exploding gradients, which also slows down the optimization algorithm.

If you train this network longer you will see better results, but initializing with overly large random numbers slows down the optimization.

The cost starts very high. This is because with large random-valued weights, the last activation (sigmoid) outputs results that are very close to 0 or 1 for some examples, and when it gets that example wrong it incurs a very high loss for that example. Indeed, when $\log(a[3])=\log(0)$, the loss goes to infinity.

****In summary**:**

- Initializing weights to very large random values does not work well.
- Hopefully initializing with small random values does better. The important question is: how small should be these random values be? Lets find out in the next part!

****Observations**:**

- The value of λ is a hyperparameter that you can tune using a dev set.
- L2 regularization makes your decision boundary smoother. If λ is too large, it is also possible to "oversmooth", resulting in a model with high bias.

****What is L2-regularization actually doing?**:**

L2-regularization relies on the assumption that a model with small weights is simpler than a model with large weights. Thus, by penalizing the square values of the weights in the cost function you drive all the weights to smaller values. It becomes too costly for the cost to have large weights! This leads to a smoother model in which the output changes more slowly as the input changes.

****What you should remember** -- the implications of L2-regularization on:**

- The cost computation:
 - A regularization term is added to the cost
- The backpropagation function:
 - There are extra terms in the gradients with respect to weight matrices
- Weights end up smaller ("weight decay"):
- Weights are pushed to smaller values.
- Dropout is a regularization technique.
- You only use dropout during training. Don't use dropout (randomly eliminate nodes) during test time.
- Apply dropout both during forward and backward propagation.
- During training time, divide each dropout layer by `keep_prob` to keep the same expected value for the activations. For example, if `keep_prob` is 0.5, then we will on average shut down half the nodes, so the output will be scaled by 0.5 since only the remaining half are contributing to the solution. Dividing by 0.5 is equivalent to multiplying by 2. Hence, the output now has the same expected value. You can check that this works even when `keep_prob` is other values than 0.5.

It seems that there were errors in the `backward_propagation_n` code we gave you! Good that you've implemented the gradient check. Go back to `backward_propagation` and try to find/correct the errors *(Hint: check dW2 and db1)*. Rerun the gradient check when you think you've fixed it. Remember you'll need to re-execute the cell defining `backward_propagation_n()` if you modify the code.

Can you get gradient check to declare your derivative computation correct? Even though this part of the assignment isn't graded, we strongly urge you to try to find the bug and re-run gradient check until you're convinced backprop is now correctly implemented.

****Note****

- Gradient Checking is slow! Approximating the gradient with $\frac{\partial J}{\partial \theta} \approx \frac{J(\theta + \varepsilon) - J(\theta - \varepsilon)}{2\varepsilon}$ is computationally costly. For this reason, we don't run gradient checking at every iteration during training. Just a few times to check if the gradient is correct.
- Gradient Checking, at least as we've presented it, doesn't work with dropout. You would usually run the gradient check algorithm without dropout to make sure your backprop is correct, then add dropout.

Course 2

Week 2

Mini-batch gradient descent

Vectorization allows you to efficiently compute on m examples.

Processing the whole training set and take its gradient descent to find the optimal solution. However, processing the entire training set is time consuming. Get a faster algorithm is possible if we start to make some progress even before finishing processing entire dataset. That is why we jump from batch gradient descent into mini-batch gradient descent.

We have to take mini-batches (baby training set). We take the mini-batch gradient descent which is a gradient descent for each batch.

Mini-batch gradient descent

repeat $\{$
for $t = 1, \dots, 5000 \}$

Forward prop on X^{t+1} .

$$z^{(t)} = w^{(t)} X^{t+1} + b^{(t)}$$

$$A^{(t)} = g^{(t)}(z^{(t)})$$

$$\vdots$$

$$A^{(t)} = g^{(t)}(z^{(t)})$$

$$\left. \begin{array}{l} \text{Compute cost } J^{t+1} = \frac{1}{m} \sum_{i=1}^m L(A^{(t)}, y^{(i)}) \\ \quad + \frac{\lambda}{2 \cdot m} \sum_{j=1}^n \|w^{(t)}\|_F^2 \end{array} \right\} \text{Vectorized implementation (1000 examples)}$$

Backprop to comput gradients wrt J^{t+1} (using (X^{t+1}, Y^{t+1}))

$$w^{(t+1)} = w^{(t)} - \alpha \nabla J^{t+1}, \quad b^{(t+1)} = b^{(t)} - \alpha \nabla b^{(t)}$$

3 3

"1 epoch"
pass through training set.

1 step of gradient desc
using $\underline{X^{t+1}}, \underline{Y^{t+1}}$
(as if $m=1000$)

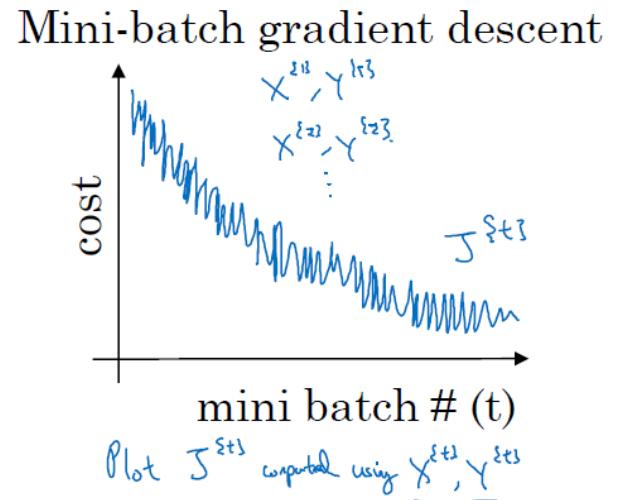
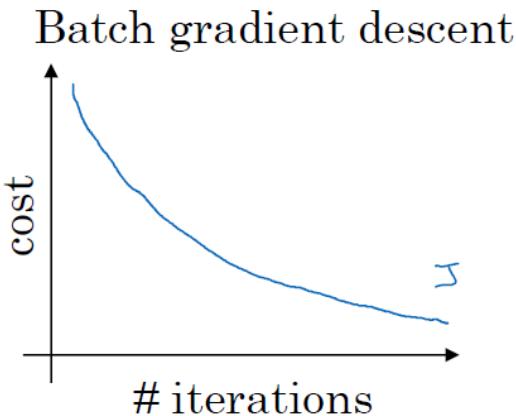
X, Y

Andrew Ng

Understanding mini-batch gradient descent

In batch gradient descent the cost decreases continuously (if everything goes well), however, in the mini-batch gradient descent we are oscillating as the trend goes down in overall. The reason is that maybe in one batch we are doing well in training and the cost goes down; but, on the next batch we may see that the training is harder due to many reasons such as incorrect labels, so the cost may not go down and instead may go a little worse than the previous one.

Training with mini batch gradient descent



Two extremes are:

- 1) If mini-batch size = m = batch size: Batch gradient descent \rightarrow too long per iteration
- 2) If mini-batch size = 1: stochastic gradient descent \rightarrow lose speed up from not having vectorization

In the first extreme, we know that we have a little noise and we end-up to the minimum; however takes longer. (the blue in the figure below)

In the second extreme, we are always oscillating because updating the rule based on one sample is not ideal and causes oscillation, specifically around the global minimum. (the purple in the figure below)

Something in between is the best option for mini-batch size which is not too big or too small \rightarrow fastest learning since it:

- supports vectorization
- does not need to wait for the whole training set to update weights

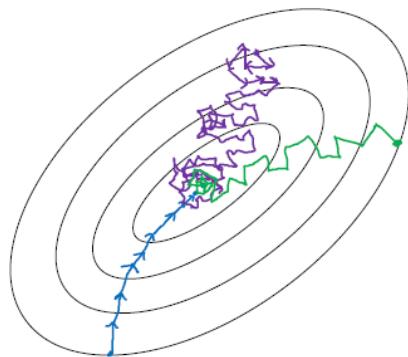
We can see it as the green in the figure below.

Choosing your mini-batch size

→ If mini-batch size = m : Batch gradient descent. $(X^{(1)}, Y^{(1)}) = (X, Y)$.

→ If mini-batch size = 1 : Stochastic gradient descent. Every example is its own $(X^{(1)}, Y^{(1)}) = (x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)})$ mini-batch.

In practice: Somewhere in-between 1 and m



Stochastic
gradient
descent
↓
Use speedup
from vectorization

In-between
(mini-batch size
not too big/small)

Fastest learning.

- Vectorization.
 $(n \times n)$
- Make passes without
processing entire training set.

Batch
gradient descent
(mini-batch size = m)
↓
Too long
per iteration

Andrew Ng

But how to choose its size? There are some rule of thumbs:

- Small training set: use BGD because we can go through the whole dataset very fast ($m \leq 2000$)
- Typical mini-batch size: better to be power of 2 like 64, 128, 256, 512, or even 1024, etc.
- Make sure mini-batch size fits in CPU/GPU memory → mini-batch size is a hyper-parameter, try some of the power of 2 and check which ones converge faster.

Choosing your mini-batch size

If small training set : Use batch gradient descent.
 $(m \leq 2000)$

Typical mini-batch sizes:

→ $64, 128, 256, 512$
 $\underbrace{2^6, 2^7, 2^8, 2^9}_{\frac{1024}{2^{10}}}$

Make sure mini-batch fits in CPU/GPU memory.

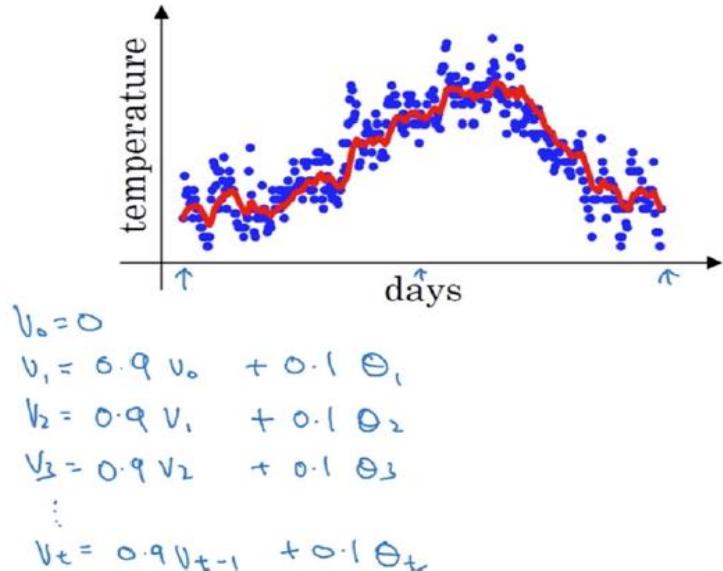
$$X^{(1)}, Y^{(1)}$$

Exponentially weighted averages

For this one, let's start with an example of a moving average which gives 0.9 weight to the previous value and 0.1 to the current one. Its plot is in red:

Temperature in London

$$\begin{aligned}\theta_1 &= 40^{\circ}\text{F} & 4^{\circ}\text{C} \leftarrow \\ \theta_2 &= 49^{\circ}\text{F} & 9^{\circ}\text{C} \\ \theta_3 &= 45^{\circ}\text{F} & \vdots \\ &\vdots \\ \theta_{180} &= 60^{\circ}\text{F} & 15^{\circ}\text{C} \\ \theta_{181} &= 56^{\circ}\text{F} & \vdots \\ &\vdots\end{aligned}$$

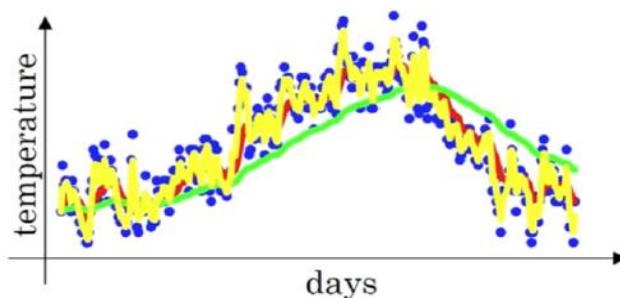


In the concept of exponentially weighted average, the information from the past $1 / (1 - \beta)$ days are preserved. For this reason, if we increase β , the number of preserved days in the moving average is higher so the green one tends to be smoother and shifted to the right as adapts slower since the current value get less importance. On the flip side, having β as 0.5 (the yellow one) oscillates more and converges sooner (adapts itself sooner to the current value).

Exponentially weighted averages

$$\begin{aligned}V_t &= \beta V_{t-1} + (1-\beta) \theta_t \leftarrow \\ \beta = 0.9 &: \approx 10 \text{ day}^{\text{'}\text{}} \text{ temp.} \\ \beta = 0.98 &: \approx 50 \text{ days} \\ \beta = 0.5 &: \approx 2 \text{ days}\end{aligned}$$

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



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

Understanding exponentially weighted averages

if the beta is 0.9, it takes around 10 days to get around one third of the 0.1. If the beta is 0.98 it takes 50 days to go to the 1/3 of the first value which is 0.02. In other words if the very first value is epsilon we have $(1-\epsilon)^{(1/\epsilon)} = 1/e$

Exponentially weighted averages

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

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

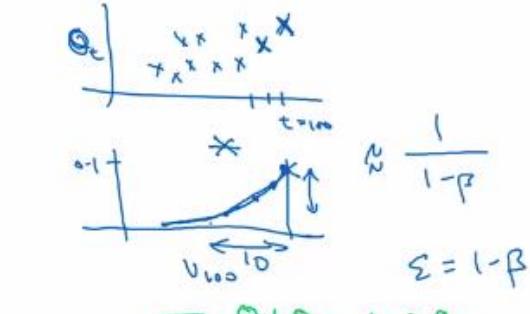
$$v_{99} = 0.9 v_{98} + 0.1 \theta_{99}$$

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

...

$$\begin{aligned} \rightarrow v_{100} &= 0.1 \theta_{100} + 0.9 (0.1 \theta_{99} + 0.9 v_{99}) \\ &= 0.1 \theta_{100} + 0.1 \cdot 0.9 \cdot \theta_{99} + 0.1 (0.9)^2 \theta_{98} + 0.1 (0.9)^3 \theta_{97} + 0.1 (0.9)^4 \theta_{96} \end{aligned}$$

$$0.9^{10} \approx 0.35 \approx \frac{1}{e}$$



$$v_{100} = (1 - \beta)^{10} v_0$$

$$\frac{(1-\beta)^{10}}{0.9} = \frac{1}{e} \quad ?$$

$$\beta = 0.98 \quad \epsilon = 0.02 \rightarrow 0.98^{50} \approx \frac{1}{e}$$

Andrew Ng

$$\begin{aligned} v_t &= (1 - \beta) \theta_t + \beta v_{t-1} = (1 - \beta) \theta_t + \beta(1 - \beta) \theta_{t-1} + \beta^2 v_{t-2} \\ &= (1 - \beta) \theta_t + \beta(1 - \beta) \theta_{t-1} + \beta^2(1 - \beta) \theta_{t-2} + \dots + \beta^t (1 - \beta) \theta_{t-t} \\ &= (1 - \beta)[\theta_t + \beta \theta_{t-1} + \beta^2 \theta_{t-2} + \dots + \beta^t \theta_{t-t}] \end{aligned}$$

We know that:

$$\begin{aligned} \frac{1}{(1 - \beta)} &= 1 + (1 - (1 - \beta)) + (1 - (1 - \beta))^2 + \dots \\ &= 1 + (\beta) + (\beta)^2 + \dots \end{aligned}$$

Now, back to our problem, we have:

$$v_t = \frac{[\theta_t + \beta \theta_{t-1} + \beta^2 \theta_{t-2} + \dots + \beta^t \theta_{t-t}]}{1 + (\beta) + (\beta)^2 + \dots}$$

If we assume the θ_t 's variation is small, we can omit this variable. Moreover, we assume that after k_{th} averaging, we can neglect the previous ones. So we re-write the above equation as:

$$\frac{(\beta^k + \beta^{k+1} + \beta^{k+2} + \dots)}{1 + (\beta) + (\beta)^2 + \dots} = \frac{\beta^k (1 + \beta^1 + \beta^2 + \dots)}{1 + (\beta) + (\beta)^2 + \dots} = \beta^k$$

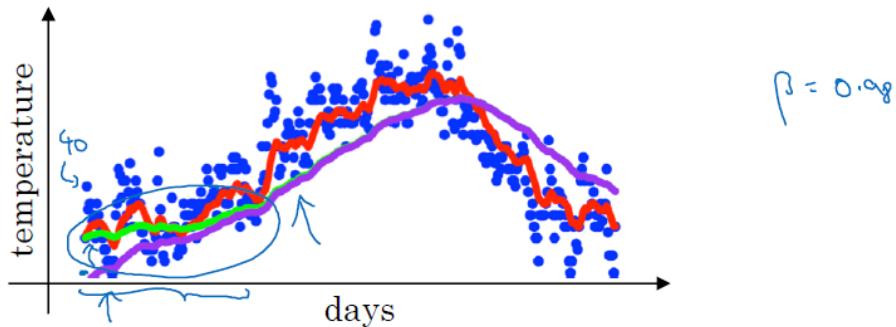
Now let's say that a value less than $1/e$ which is around 0.37 is negligible. Then we can write it down as:

$$\beta^k = \frac{1}{e} \Rightarrow k \log(\beta) = -\log(e) \Rightarrow k = -\frac{\log(e)}{\log(\beta)}$$

For example, if $\beta = 0.9$ as the above example, we get $k = 9.5$ in which means that after around 10 days the importance of the first day goes less than $1/e$ and we can ignore it.

Bias correction in exponentially weighted averages

Bias correction



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

$$v_0 = 0$$

$$v_1 = 0.98 v_0 + 0.02 \theta_1$$

$$v_2 = 0.98 v_1 + 0.02 \theta_2$$

$$= 0.98 \times 0.02 \times \theta_1 + 0.02 \theta_2$$

$$= 0.0196 \theta_1 + 0.02 \theta_2$$

$$\left| \begin{array}{l} \frac{v_t}{1 - \beta^t} \\ t=2: 1 - \beta^t = 1 - (0.98)^2 = 0.0396 \\ \frac{v_2}{0.0396} = \frac{0.0196 \theta_1 + 0.02 \theta_2}{0.0396} \end{array} \right.$$

Andrew Ng

If we look at the example on the left side of the figure above, we see that when we start from $v_0 = 0$, it takes some time for the system to warm up. In this case we need to give it a bias which helps the system to warm up sooner. For doing so, we want the coefficient of the first data in the first averaging to be 1, instead of $1 - \beta$ (in the example above it is 0.02). For doing so, the whole equation should be divided by $1 - \beta$. However, we also want the effect of bias to converge to 0 in the long run. This happens when we use β^t which means when $t \rightarrow \infty$ and $0 < \beta < 1$, $\beta^t \rightarrow 0$; so $1 - \beta^t \rightarrow 1$. By doing so, the purple line becomes the green one.

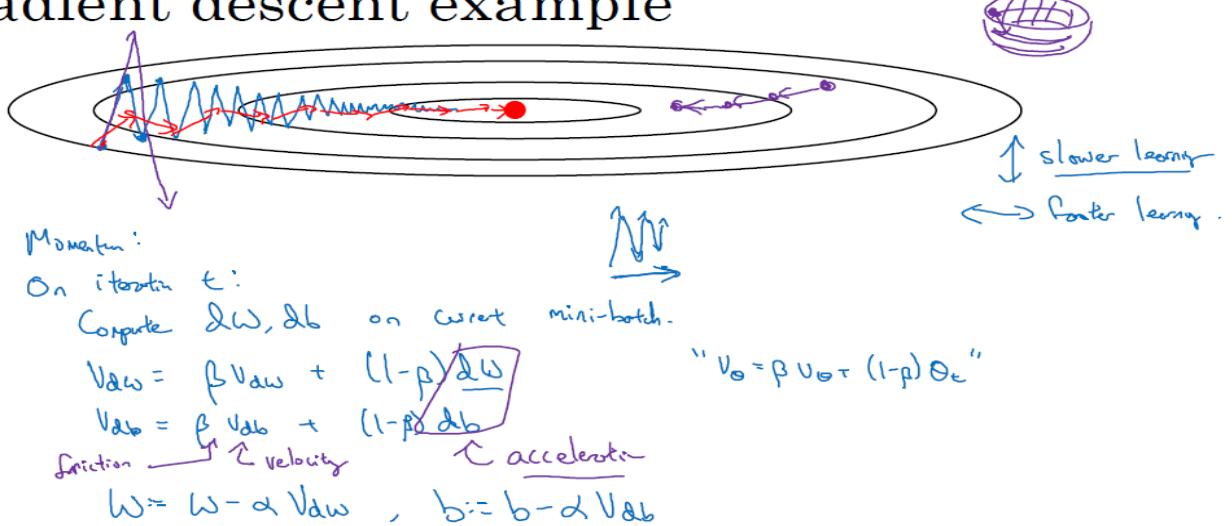
Gradient Descent with momentum

Oscillations toward the center, most of the times, takes many steps as sometimes steps causes divergence from the correct path. Moreover, like the example below, learning rate decreases faster vertically than horizontally.

Momentum optimizer smooth up each iteration in the gradient descent by manipulating learning rate for each dimension. If you take the average of the vertical oscillation, it will tends to zero. So it slows down the vertical learning rate. But for the horizontal one it is much faster. → more straight forward path to the minimum.

In these equations, v_{dw} and v_{db} are acting as velocity or momentum toward directions by taking the moving average of the previous steps and dw and db play the role of accelerator as always.

Gradient descent example



α and β control exponentially moving weighted average and both of them are hyper-parameters. You can use beta as hyper parameters. Bias correction is not needed in practice; however, using it does not hurt and actually is useful.

Also, sometimes $1 - \beta$ terms are omitted. Using it or not using it are both acceptable but omitting them makes it less intuitive; then we need to retune α again.

Implementation details

$$V_{aw} = 0, \quad V_{ab} = 0$$

On iteration t :

Compute dW, db on the current mini-batch

$$\begin{aligned} \rightarrow v_{dW} &= \beta v_{dW} + (1-\beta) dW \\ \rightarrow v_{db} &= \beta v_{db} + (1-\beta) db \end{aligned} \quad \left. \begin{array}{l} v_{dW} = \beta v_{dW} + dW \\ v_{db} = \beta v_{db} + db \end{array} \right\} \quad \left. \begin{array}{l} v_{dW} = \beta v_{dW} + dW \\ v_{db} = \beta v_{db} + db \end{array} \right\} \quad \left. \begin{array}{l} v_{dW} = \beta v_{dW} + dW \\ v_{db} = \beta v_{db} + db \end{array} \right\} \quad \left. \begin{array}{l} v_{dW} = \beta v_{dW} + dW \\ v_{db} = \beta v_{db} + db \end{array} \right\}$$

Hyperparameters: α, β

$$\beta = 0.9$$

average over last 10 graduates

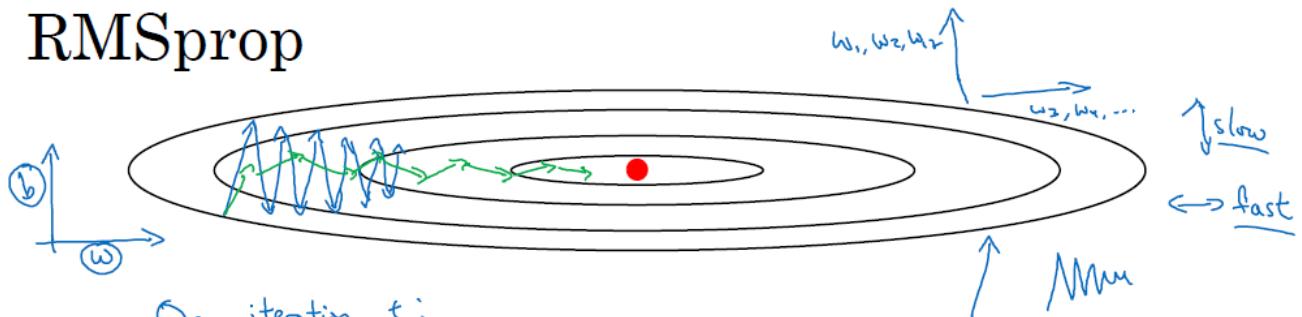
Andrew Ng

RMSprop

To increase the momentum we have RMSprop. For the sake of the example shown in the slide below, if we assume b is the horizontal movement and w is the vertical movement, our goal is to move toward the goal by decreasing b such way that it oscillates less vertically. The same applies on w . In this example, we want to slow down learning on vertical, speed up on horizontal. The power of 2 causes the dw to become smaller (if it is less than 1) and db becomes larger so that the Sdw and Sdb are getting smaller and bigger respectively. Then in the update rule, by dividing the gradient on a squared root of these values the path toward the global minima would be less oscillated. In this way we can use larger learning rate. In practice dW is a high dimensional value.

One more point is that we need to add epsilon to denominator as it might get so close to 0 and we may encounter a division by 0 problem.

RMSprop



On iteration t :

Compute dw, db on current mini-batch,

$$S_{dw} = \beta_2 S_{dw} + (1-\beta_2) dw^2 \quad \leftarrow \text{element-wise}$$

$$\rightarrow S_{db} = \beta_2 S_{db} + (1-\beta_2) db^2 \quad \leftarrow \text{large}$$

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

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

Adam Optimization Algorithm

Momentum + RMSprop = Adam (Adaptive Moment Estimation)

It has bias correction in it. → highly popular and useful

In the slide below, line 1 is the same as what we did in the momentum. Line 2 is the same as what we did in the RMSprop. Lines 3 and 4 are only about bias correction. Finally, line 5 is the mixture of momentum and RMSprop algorithms for updating rules.

Adam optimization algorithm

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

On iteration t :

Compute $\delta w, \delta b$ using current mini-batch

$$1 \quad V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) \delta w, \quad V_{db} = \beta_1 V_{db} + (1 - \beta_1) \delta b \quad \leftarrow \text{"momentum"} \beta_1$$

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

$$3 \quad V_{dw}^{\text{corr}} = V_{dw} / (1 - \beta_1^t), \quad V_{db}^{\text{corr}} = V_{db} / (1 - \beta_1^t)$$

$$4 \quad S_{dw}^{\text{corr}} = S_{dw} / (1 - \beta_2^t), \quad S_{db}^{\text{corr}} = S_{db} / (1 - \beta_2^t)$$

$$5 \quad W := W - \alpha \frac{V_{dw}^{\text{corr}}}{\sqrt{S_{dw}^{\text{corr}}} + \epsilon} \quad b := b - \alpha \frac{V_{db}^{\text{corr}}}{\sqrt{S_{db}^{\text{corr}}} + \epsilon}$$

For the hyper-parameters:

- α : needs to be tuned
- β_1 : 0.9 (dw) → default value and does not need to be tuned most of the times
- β_2 : 0.999 (δw^2) → default value and does not need to be tuned most of the times
- ϵ : 10^{-8} → default value and does not need to be tuned at all

Learning rate decay

Sometimes with fix value alpha and noise we wander around min but never meet it because the learning rate is large. So we need to implement an adaptive learning rate.

Learning rate decay

1 epoch = 1 pass through data.

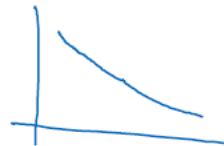
$$\alpha = \frac{1}{1 + \text{decay-rate} * \text{epoch-num}} \cdot \alpha_0$$

Epoch	α
1	0.1
2	0.062
3	0.05
4	0.04
:	i



$$\alpha_0 = 0.2$$

$$\text{decay-rate} = 1$$



Decay rate is another hyper-parameter. We have also exponentially decay, discrete stair case decay, manual decay, etc.

Other learning rate decay methods

formulas

$$\left\{ \begin{array}{l} \alpha = 0.95^{\text{epoch-num}} \cdot \alpha_0 \quad - \text{exponentially decay} \\ \alpha = \frac{k}{\sqrt{\text{epoch-num}}} \cdot \alpha_0 \quad \text{or } \frac{k}{\sqrt{t}} \cdot \alpha_0 \end{array} \right.$$

α | t

discrete staircase

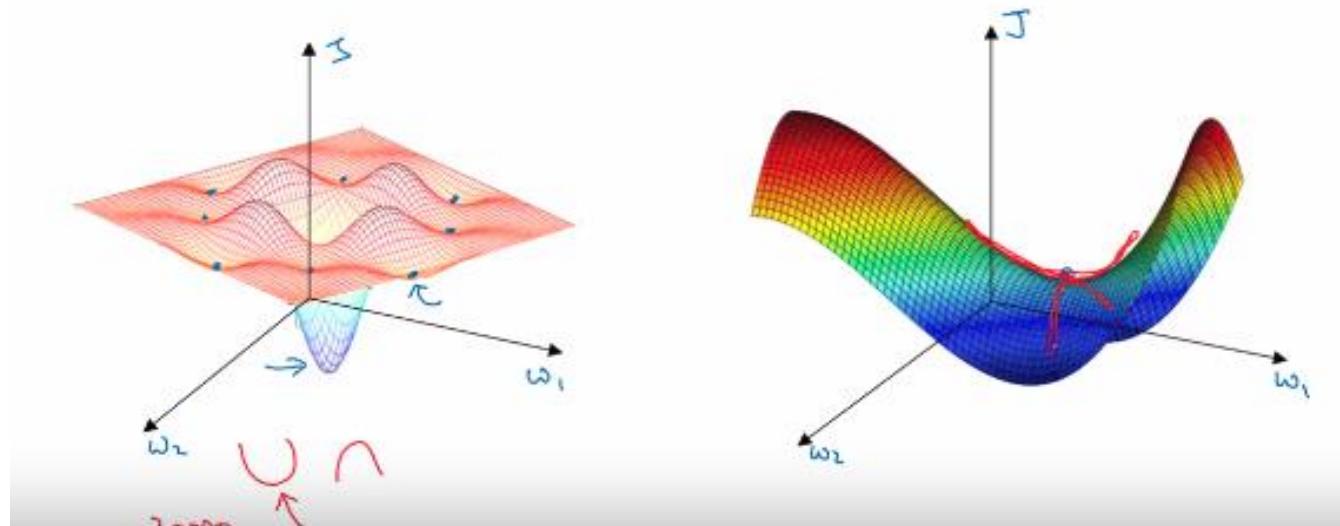
Manual decay.

The problem of local optima

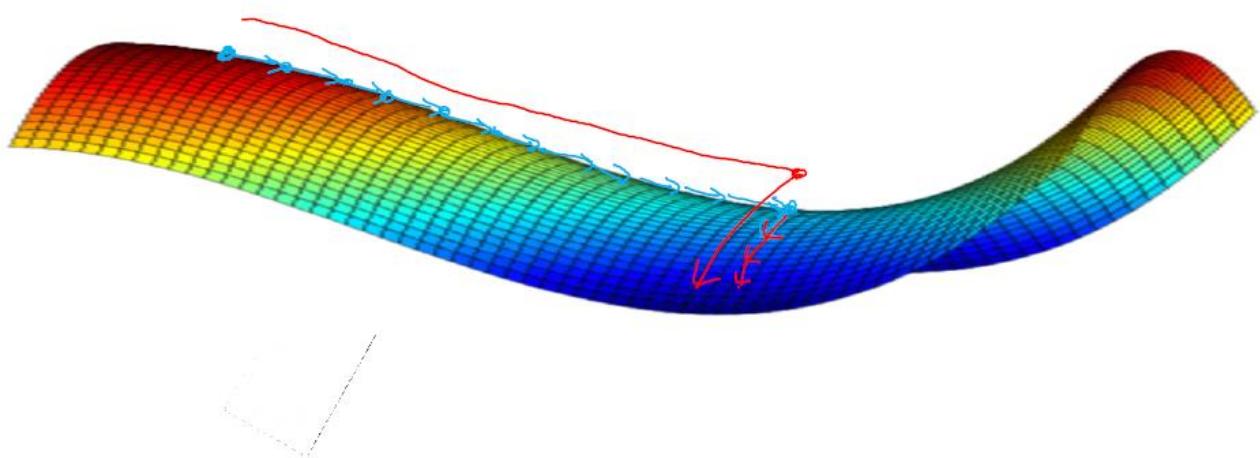
Local optima problem: gradient descent might end up in local optima. However, in neural network, most points that were assumed as local optima are only saddle points (derivative of zero).

The plateau is a region where the derivative is close to zero for a long time and this is the problem that really slow down the learning. So not only local optima is a problem, but plateau is also a problem. So it is unlikely to get stuck in a bad local optima. Plateaus are the problems which adam or RMSprop can be helpful here to move down the plateau.

Local optima in neural networks



Problem of plateaus



- Unlikely to get stuck in a bad local optima
- Plateaus can make learning slow

Assignment

A variant of this is Stochastic Gradient Descent (SGD), which is equivalent to mini-batch gradient descent where each mini-batch has just 1 example. The update rule that you have just implemented does not change. What changes is that you would be computing gradients on just one training example at a time, rather than on the whole training set. The code examples below illustrate the difference between stochastic gradient descent and (batch) gradient descent.

- **(Batch) Gradient Descent**:

```
''' python
X = data_input
Y = labels
parameters = initialize_parameters(layers_dims)
for i in range(0, num_iterations):
    # Forward propagation
    a, caches = forward_propagation(X, parameters)
    # Compute cost.
    cost += compute_cost(a, Y)
    # Backward propagation.
    grads = backward_propagation(a, caches, parameters)
    # Update parameters.
    parameters = update_parameters(parameters, grads)
```

'''

- **Stochastic Gradient Descent**:

```
'''python
X = data_input
Y = labels
parameters = initialize_parameters(layers_dims)
for i in range(0, num_iterations):
    for j in range(0, m):
        # Forward propagation
        a, caches = forward_propagation(X[:,j], parameters)
        # Compute cost
        cost += compute_cost(a, Y[:,j])
        # Backward propagation
        grads = backward_propagation(a, caches, parameters)
        # Update parameters.
        parameters = update_parameters(parameters, grads)
'''
```

- The difference between gradient descent, mini-batch gradient descent and stochastic gradient descent is the number of examples you use to perform one update step.
- You have to tune a learning rate hyperparameter α .

- With a well-turned mini-batch size, usually it outperforms either gradient descent or stochastic gradient descent (particularly when the training set is large).

In Stochastic Gradient Descent, you use only 1 training example before updating the gradients. When the training set is large, SGD can be faster. But the parameters will "oscillate" toward the minimum rather than converge smoothly. Here is an illustration of this:

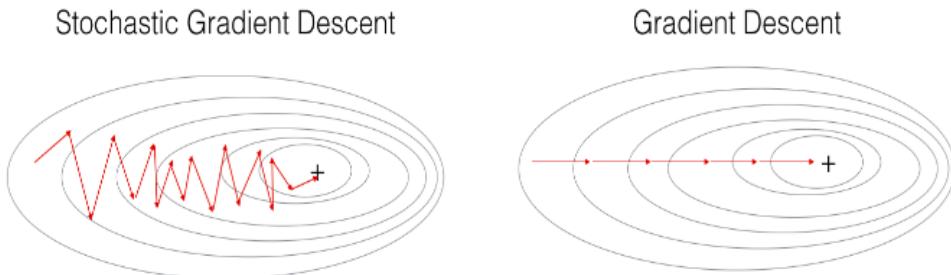


Figure 1 : SGD vs GD

"+" denotes a minimum of the cost. SGD leads to many oscillations to reach convergence. But each step is a lot faster to compute for SGD than for GD, as it uses only one training example (vs. the whole batch for GD).

Note also that implementing SGD requires 3 for-loops in total:

1. Over the number of iterations
2. Over the m training examples
3. Over the layers (to update all parameters, from $(W^{[1]}, b^{[1]})$ to $(W^{[L]}, b^{[L]})$)

In practice, you'll often get faster results if you do not use neither the whole training set, nor only one training example, to perform each update. Mini-batch gradient descent uses an intermediate number of examples for each step. With mini-batch gradient descent, you loop over the mini-batches instead of looping over individual training examples.

In practice, you'll often get faster results if you do not use neither the whole training set, nor only one training example, to perform each update. Mini-batch gradient descent uses an intermediate number of examples for each step. With mini-batch gradient descent, you loop over the mini-batches instead of looping over individual training examples.

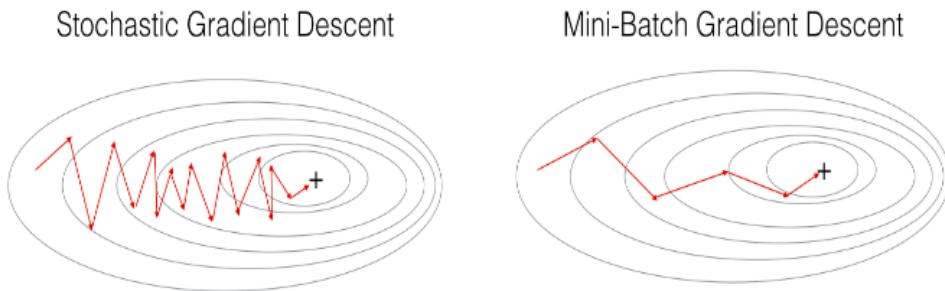


Figure 2 : SGD vs Mini-Batch GD

"+" denotes a minimum of the cost. Using mini-batches in your optimization algorithm often leads to faster optimization.

What you should remember:

- The difference between gradient descent, mini-batch gradient descent and stochastic gradient descent is the number of examples you use to perform one update step.
- You have to tune a learning rate hyperparameter α .
- With a well-turned mini-batch size, usually it outperforms either gradient descent or stochastic gradient descent (particularly when the training set is large).

2 - Mini-Batch Gradient descent

Let's learn how to build mini-batches from the training set (X, Y).

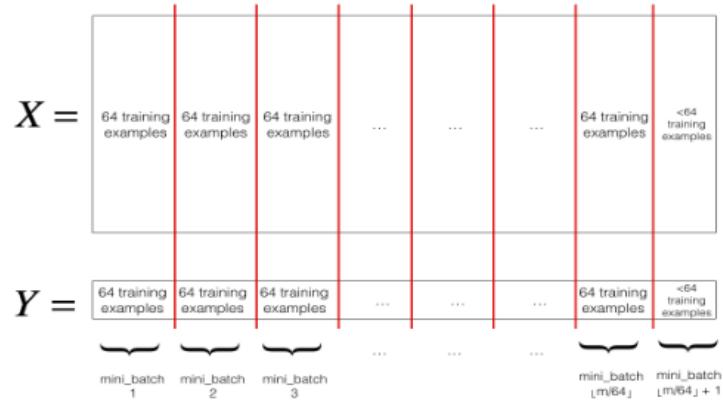
There are two steps:

- **Shuffle:** Create a shuffled version of the training set (X, Y) as shown below. Each column of X and Y represents a training example. Note that the random shuffling is done synchronously between X and Y . Such that after the shuffling the i^{th} column of X is the example corresponding to the i^{th} label in Y . The shuffling step ensures that examples will be split randomly into different mini-batches.

$$X = \begin{pmatrix} X_0^{(1)} & X_0^{(2)} & \dots & X_0^{(m-1)} & X_0^{(m)} \\ X_1^{(1)} & X_1^{(2)} & \dots & X_1^{(m-1)} & X_1^{(m)} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ X_{12286}^{(1)} & X_{12286}^{(2)} & \dots & X_{12286}^{(m-1)} & X_{12286}^{(m)} \\ X_{12287}^{(1)} & X_{12287}^{(2)} & \dots & X_{12287}^{(m-1)} & X_{12287}^{(m)} \end{pmatrix} \quad Y = \begin{pmatrix} y^{(1)} & y^{(2)} & \dots & y^{(m-1)} & y^{(m)} \end{pmatrix}$$

$$X = \begin{pmatrix} X_0^{(1)} & X_0^{(2)} & \dots & X_0^{(m-1)} & X_0^{(m)} \\ X_1^{(1)} & X_1^{(2)} & \dots & X_1^{(m-1)} & X_1^{(m)} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ X_{12286}^{(1)} & X_{12286}^{(2)} & \dots & X_{12286}^{(m-1)} & X_{12286}^{(m)} \\ X_{12287}^{(1)} & X_{12287}^{(2)} & \dots & X_{12287}^{(m-1)} & X_{12287}^{(m)} \end{pmatrix} \quad Y = \begin{pmatrix} y^{(1)} & y^{(2)} & \dots & y^{(m-1)} & y^{(m)} \end{pmatrix}$$

- **Partition:** Partition the shuffled (X, Y) into mini-batches of size `mini_batch_size` (here 64). Note that the number of training examples is not always divisible by `mini_batch_size`. The last mini batch might be smaller, but you don't need to worry about this. When the final mini-batch is smaller than the full `mini_batch_size`, it will look like this:



What you should remember:

- Shuffling and Partitioning are the two steps required to build mini-batches
- Powers of two are often chosen to be the mini-batch size, e.g., 16, 32, 64, 128.

3 - Momentum

Because mini-batch gradient descent makes a parameter update after seeing just a subset of examples, the direction of the update has some variance, and so the path taken by mini-batch gradient descent will "oscillate" toward convergence. Using momentum can reduce these oscillations.

Momentum takes into account the past gradients to smooth out the update. We will store the 'direction' of the previous gradients in the variable v . Formally, this will be the exponentially weighted average of the gradient on previous steps. You can also think of v as the "velocity" of a ball rolling downhill, building up speed (and momentum) according to the direction of the gradient/slope of the hill.

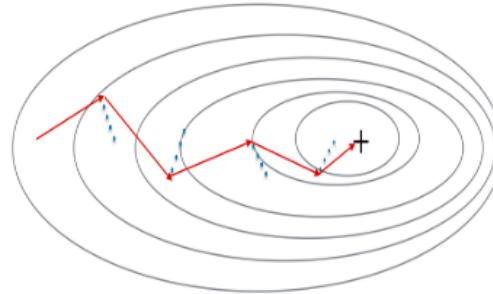


Figure 3: The red arrows shows the direction taken by one step of mini-batch gradient descent with momentum. The blue points show the direction of the gradient (with respect to the current mini-batch) on each step. Rather than just following the gradient, we let the gradient influence v and then take a step in the direction of v .

Note that:

- The velocity is initialized with zeros. So the algorithm will take a few iterations to "build up" velocity and start to take bigger steps.
- If $\beta = 0$, then this just becomes standard gradient descent without momentum.

How do you choose β ?

- The larger the momentum β is, the smoother the update because the more we take the past gradients into account. But if β is too big, it could also smooth out the updates too much.
- Common values for β range from 0.8 to 0.999. If you don't feel inclined to tune this, $\beta = 0.9$ is often a reasonable default.
- Tuning the optimal β for your model might need trying several values to see what works best in term of reducing the value of the cost function J .

What you should remember:

- Momentum takes past gradients into account to smooth out the steps of gradient descent. It can be applied with batch gradient descent, mini-batch gradient descent or stochastic gradient descent.
- You have to tune a momentum hyperparameter β and a learning rate α .

4 - Adam

Adam is one of the most effective optimization algorithms for training neural networks. It combines ideas from RMSProp (described in lecture) and Momentum.

How does Adam work?

1. It calculates an exponentially weighted average of past gradients, and stores it in variables v (before bias correction) and $v^{corrected}$ (with bias correction).
2. It calculates an exponentially weighted average of the squares of the past gradients, and stores it in variables s (before bias correction) and $s^{corrected}$ (with bias correction).
3. It updates parameters in a direction based on combining information from "1" and "2".

The update rule is, for $l = 1, \dots, L$:

$$\begin{cases} v_{dW^{[l]}} = \beta_1 v_{dW^{[l]}} + (1 - \beta_1) \frac{\partial J}{\partial W^{[l]}} \\ v_{dW^{[l]}}^{corrected} = \frac{v_{dW^{[l]}}}{1 - (\beta_1)^t} \\ s_{dW^{[l]}} = \beta_2 s_{dW^{[l]}} + (1 - \beta_2) \left(\frac{\partial J}{\partial W^{[l]}} \right)^2 \\ s_{dW^{[l]}}^{corrected} = \frac{s_{dW^{[l]}}}{1 - (\beta_2)^t} \\ W^{[l]} = W^{[l]} - \alpha \frac{v_{dW^{[l]}}^{corrected}}{\sqrt{s_{dW^{[l]}}^{corrected} + \epsilon}} \end{cases}$$

where:

- t counts the number of steps taken of Adam
- L is the number of layers
- β_1 and β_2 are hyperparameters that control the two exponentially weighted averages.
- α is the learning rate
- ϵ is a very small number to avoid dividing by zero

As usual, we will store all parameters in the `parameters` dictionary

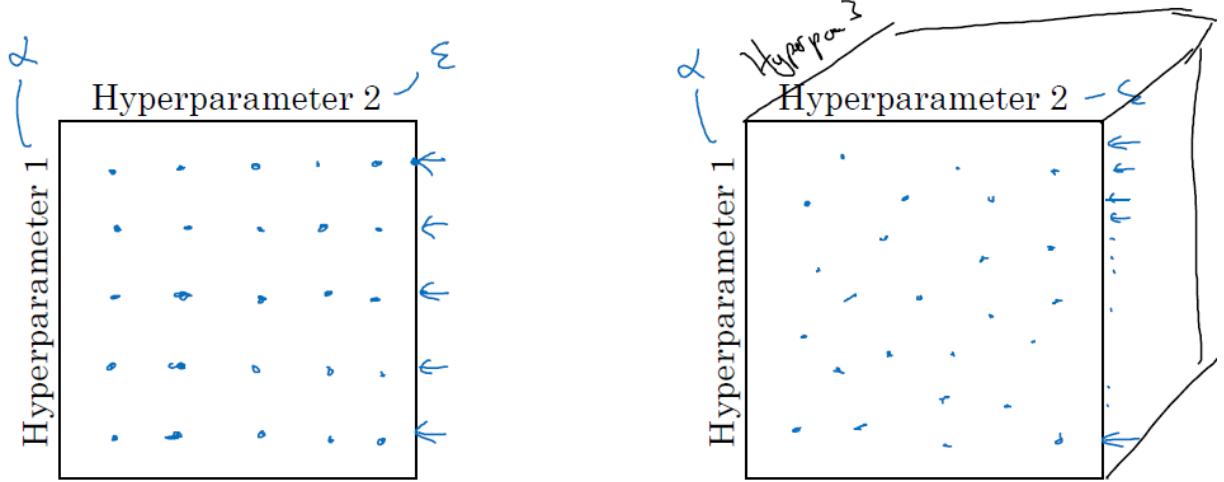
Course 2

WEEK 3

Tunning process

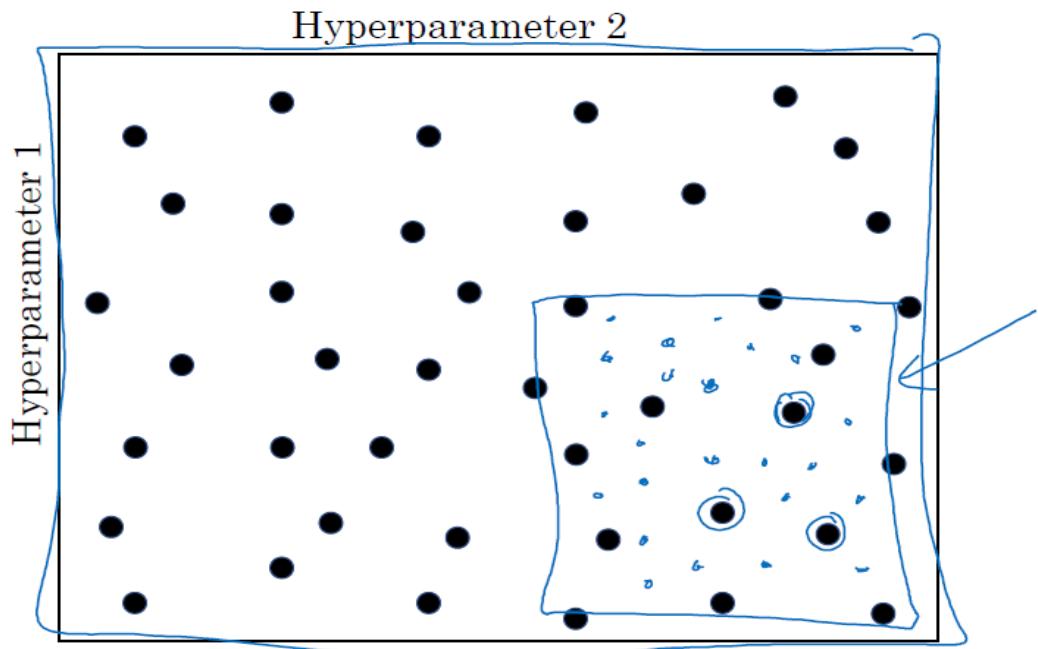
Sampling the points in the grid is practical. In an era before the deep learning, choosing the points such that is shown in left figure of the below slide was practical. In deep learning era, points must be random. Because in advance it is hard to know which hyper-parameter is the best for tunning. Because some are more important than the others. The reason is that, in the left figure, for the hyper-parameter 2, we only have 5 values to check. However, in the right figure, there are 25 of them which gives us more flexibility to choose the value wisely.

Try random values: Don't use a grid



In the coarse to fine scheme, we want to find which values for hyper-parameters in our random grid works the best. So we find it out roughly by search different point of the grid. Then for the one we found the best, we zoom into that part of the grid and try different values randomly again.

Coarse to fine



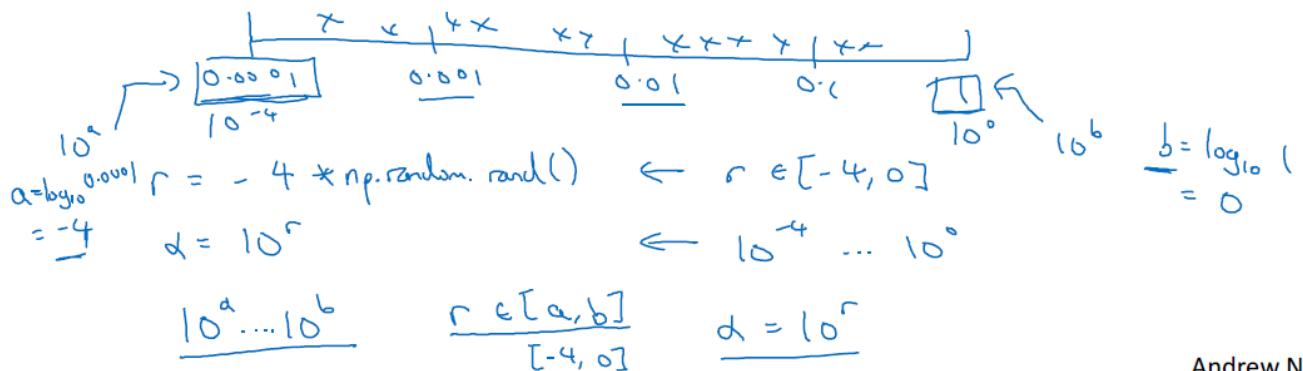
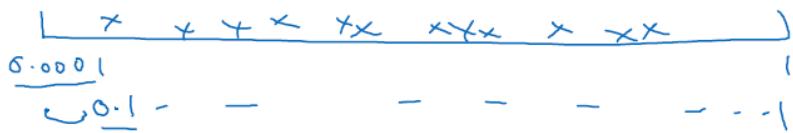
Using an appropriate scale to pick hyper-parameters

Sometimes, for some hyper-parameter, choosing values uniformly is the best way. Like for the time when we want to choose number of hidden layers which we assume the best would be between 2 and 5. So we try 2, 3, 4, and 5 uniformly. But it is not always the case.

Sometimes picking an appropriate value in an original scale does not work well like the example below since it is uniformly spread. In the example below if we want to find the learning rate which is between 0.0001 and 1, 90% of times we are checking values between 0.1 and 1 and only 10% of times we are checking between 0.0001 and 0.1. For this reason, we use a logarithmic scale for searching the learning rate. So we divide the hyper-parameter environment to different scaled environment so that it would be way more useful in searching for the actual hyper-parameter.

Appropriate scale for hyperparameters

$$\alpha = 0.0001, \dots, 1$$



Andrew Ng

For weighted averages we know that if β was 0.9, we are averaging for the past 10 days while for 0.999 we are averaging over the past 1000 days. So again, linear environment is not useful at all. So instead we use $1 - \beta$ for 0.1 to 0.001 in the logarithmic scale which is between 10^{-1} and 10^{-3} .

Hyperparameters for exponentially weighted averages

$$\beta = 0.9 \dots 0.999$$

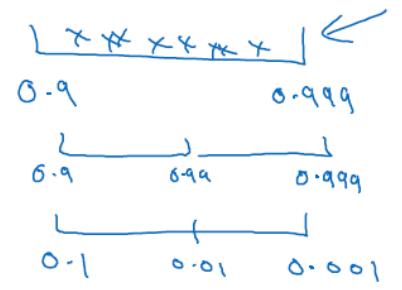
\downarrow \downarrow
 10 1000

$$1-\beta = 0.1 \dots 0.001$$

$$\beta: 0.900 \rightarrow 0.9005 \quad \} \sim 10$$

$$\beta: 0.999 \rightarrow 0.9995$$

~ 1000 ~ 2000



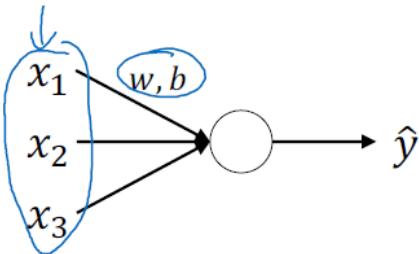
$$r \in [-3, -1]$$
$$1-\beta = 10^r$$
$$\beta = 1 - 10^r$$

$$\frac{1}{1-\beta}$$

Normalizing activations in a network

We saw it before that normalizing inputs are important to speed up the learning. So we compute the variance and mean and we normalize the input features based on them. But in deep network we have many input to layers. So two of our options is to normalize data before going into the activation function or after that. The latter is more widely used in practice.

Normalizing inputs to speed up learning

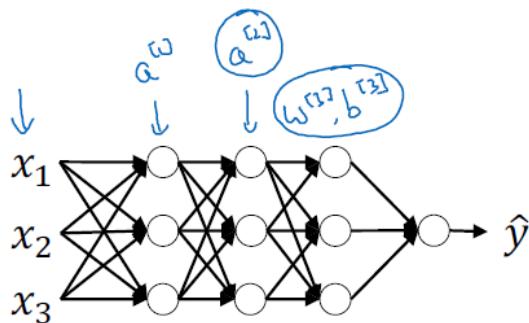


$$\mu = \frac{1}{m} \sum_i x^{(i)}$$

$$X = X - \mu$$

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

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



Can we normalize $\frac{\alpha^{[2]}}{\sigma^{[2]}}$ so
as to tan $w^{[2]}, b^{[2]}$ faster

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

How to implement batch norm?

First we get the mean and variance for values in each layer independently and we will normalize them. It makes the model to have mean 1 and variance zero. However, we do not always want that to happen since maybe we want to use different distributions. So we will come with different γ and β hyper-parameters in order to control the distribution shape. It allows the data to be whatever mean and variance we want so γ is the same as variance and β same as mean. So for example for sigmoid, we don't want our values get too large or small because it will end up with 0 and 1. So we use a distribution which prevents that problem. They may get the same values of mean and variance of the actual values of that layer which means that the normalized values would be the same as the actual values.

Implementing Batch Norm

Given some intermediate values in NN

$$\begin{aligned}\mu &= \frac{1}{m} \sum z^{(i)} \\ \sigma^2 &= \frac{1}{m} \sum (z_i - \mu)^2 \\ z_{\text{norm}}^{(i)} &= \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}} \\ z^{(i)} &= \gamma z_{\text{norm}}^{(i)} + \beta\end{aligned}$$

Use $\hat{z}^{(i)}$ instead of $\tilde{z}^{(i)}$.

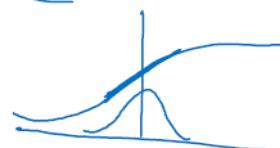
If $\gamma = \sqrt{\sigma^2 + \epsilon}$ ←
 $\beta = \mu$ ←
 then $\hat{z}^{(i)} = z^{(i)}$

learnable parameters
of model.

$$z^{(1)}, \dots, z^{(m)} \downarrow \downarrow \quad z^{(i)}$$

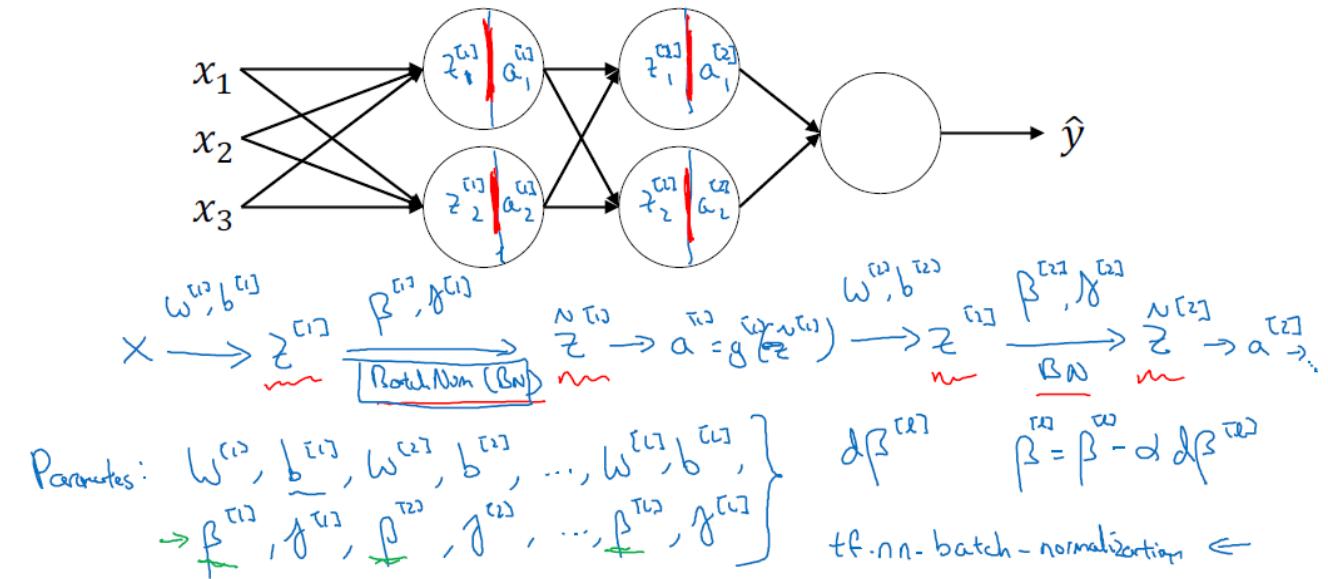
$$x \leftarrow$$

$$\hat{z}^{(i)} \leftarrow$$



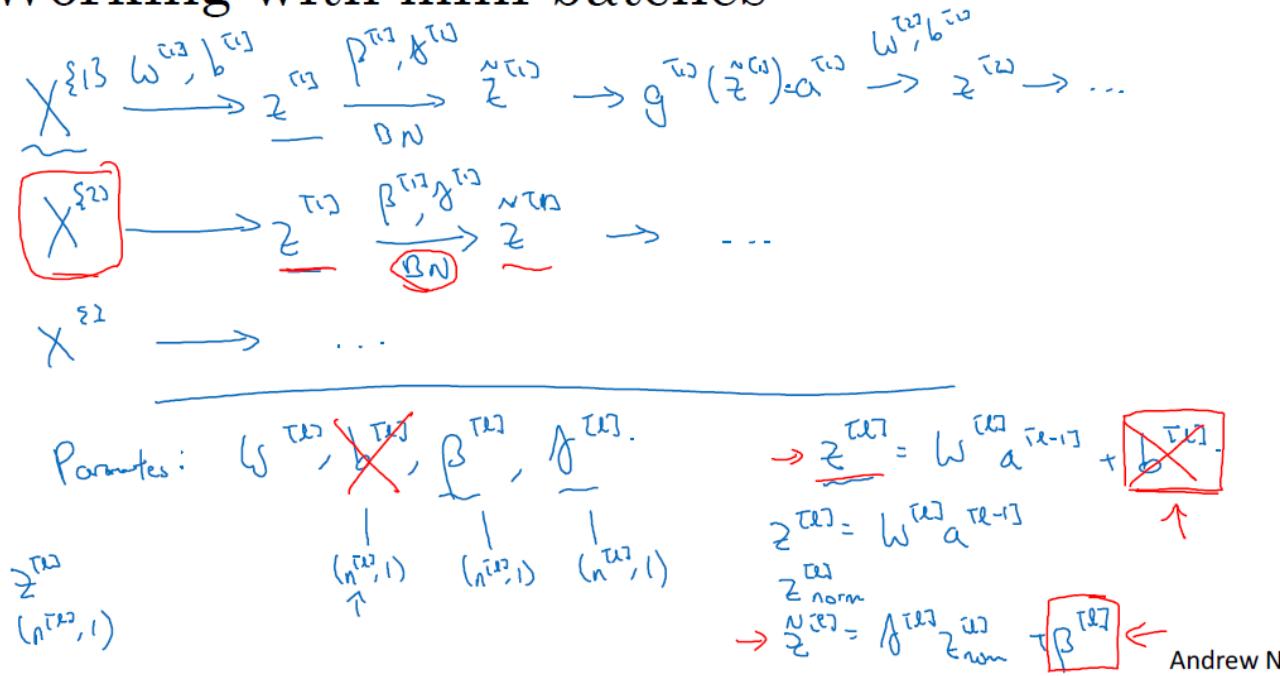
Fitting Batch Norm into a neural network

Adding Batch Norm to a network



The batch norm is different for different mini-batches for each layer. The bias constant will get cancelled out when it comes to normalization. So "b" is not needed as parameters for batch normalization because of zeroing out the mean.

Working with mini-batches



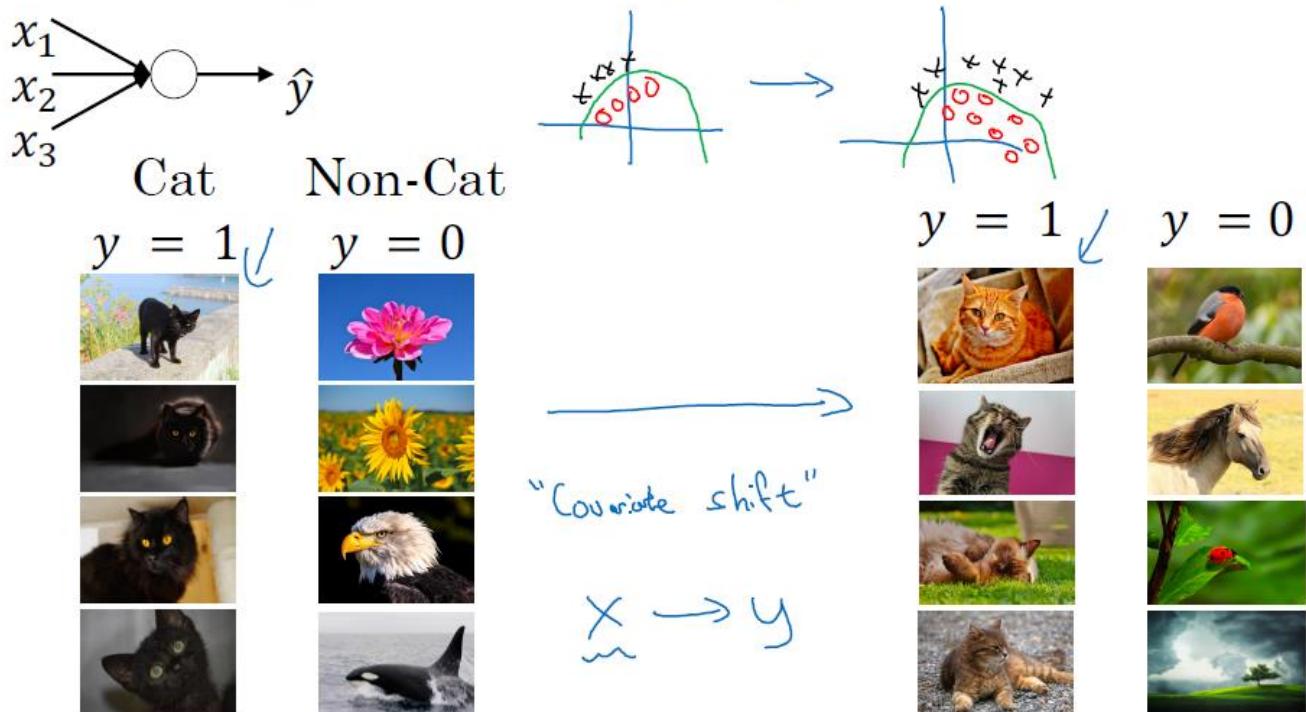
Any update algorithm works for updating parameters.

Why does batch norm work?

One reason is the same reason why we normalized input data.

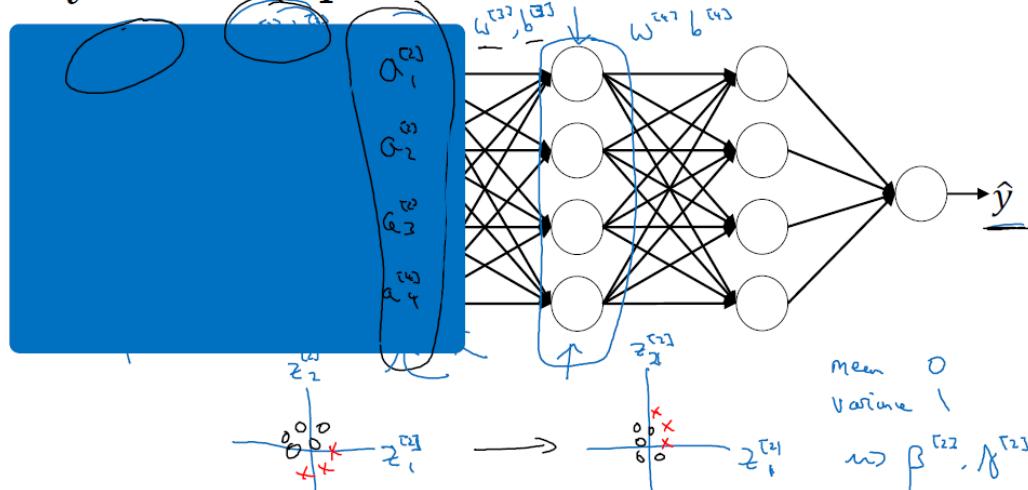
We also need to learn the shifting input distribution without re-learning. So one thing that may happen is that the data distribution may change so that the previously learned decision boundary cannot work on the data with new distribution. This incident has been called *covariant shift*. So we do not want to retrain our network.

Learning on shifting input distribution



It is obvious that the weights and biases in the previous layers are changing so that the values of activation layers will change as well. When the values of an activation layer changes all the time, it suffers from covariant shifting. So batch norm tries to use a single distribution for that layer so that the values in that layer will not change a lot as their mean and variance will stay the same.

Why this is a problem with neural networks?



Batch norm can also be used as regularization but practically it is not known for that and try to not use that as a regularizer because its effect is so less than expected.

Batch Norm as regularization

X

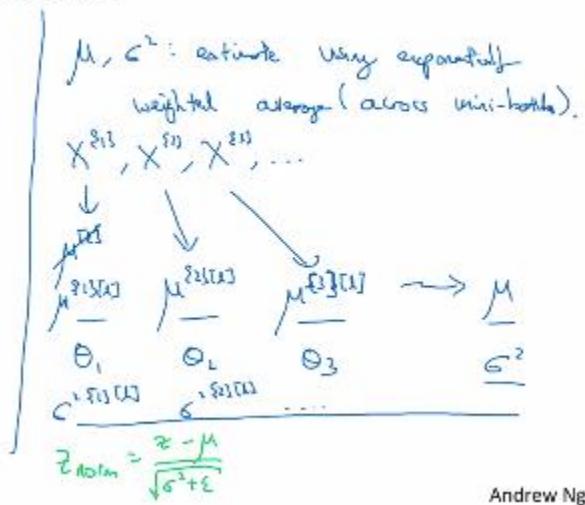
- Each mini-batch is scaled by the mean/variance computed on just that mini-batch.
 $\xrightarrow{z^{[l]}}$ $\underline{\mu}, \underline{\sigma^2}$ $\underline{z^{[l]}}$ $\times^{\{+\}}$
- This adds some noise to the values $z^{[l]}$ within that minibatch. So similar to dropout, it adds some noise to each hidden layer's activations.
 μ, σ^2
- This has a slight regularization effect.

Batch norm for testing

Mean and variance is getting calculated for every single mini-batch in training time. But for testing we do not have mini-batches so we use exponential weighted average across all mini-batches in the training time and we will use it for the test time.

Batch Norm at test time

$$\begin{aligned} \rightarrow \mu &= \frac{1}{m} \sum_i z^{(i)} \\ \rightarrow \sigma^2 &= \frac{1}{m} \sum_i (z^{(i)} - \mu)^2 \\ \rightarrow z_{\text{norm}}^{(i)} &= \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}} \leftarrow \\ \rightarrow \tilde{z}^{(i)} &= \gamma z_{\text{norm}}^{(i)} + \beta \end{aligned}$$



Andrew Ng

COURSE 3

WEEK1

Why ML strategy?

There are many options we can try so that the model may improve. But which ones shall we try and why some of them worth to try and some don't. For example, we may put a 6 months effort on collecting more data while it may not help at all in increasing the accuracy.

Motivating example



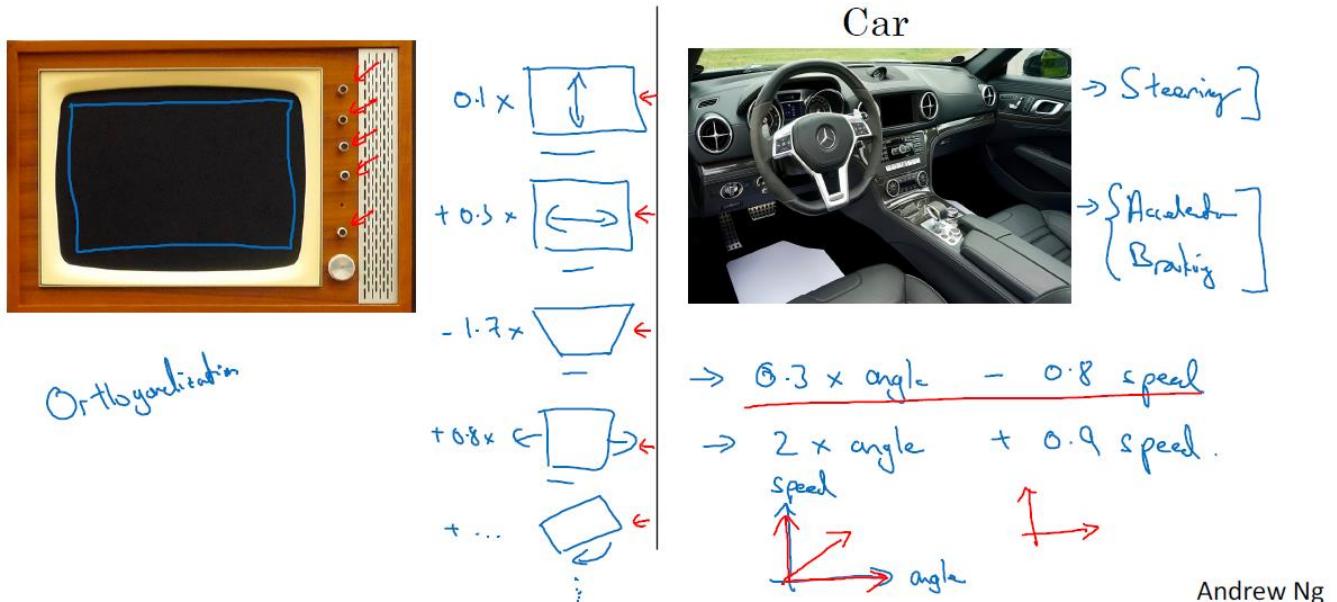
Ideas:

- Collect more data
- Collect more diverse training set
- Train algorithm longer with gradient descent
- Try Adam instead of gradient descent
- Try bigger network
- Try smaller network
- Try dropout
- Add L_2 regularization
- Network architecture
 - Activation functions
 - # hidden units
 - ...

Andrew Ng

Orthogonalization

One of the things happening is about the most effective machine learning people which are very clear eyed about what to tune in order to try to achieve one effect. This is a process we call orthogonalization. Like old TV knobs for tuning the video. In this context, orthogonalization is existence of different knobs in which each one is responsible for changing the width, height, and many other things. Each knob does not interfere with work of the other one. In the figures below, we can see orthogonalization for car and TV.



We expect for ideal case:

- training set fits well on cost function
- → then we hope dev set fits well on cost function
- → then we hope test set fits well on cost function
- → then we hope the result can perform well in real world

However, for each state, it may not work as expected, so we can play with its knob:

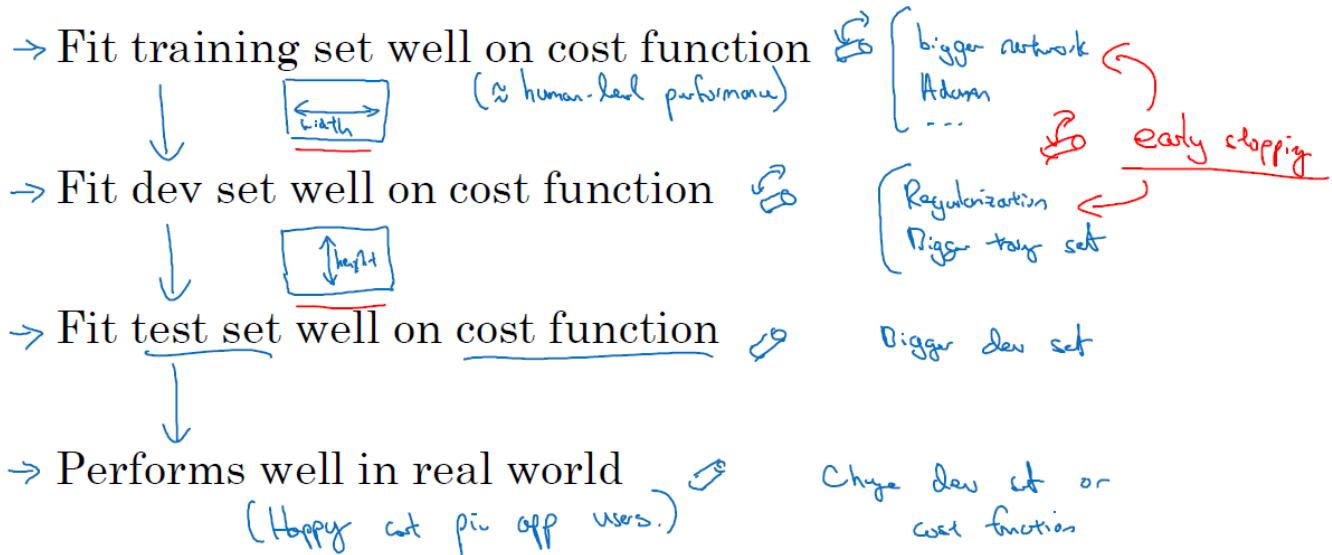
- training set does NOT fit well on cost function
 - 1- maybe bigger network
 - 2- better optimization algorithm
 - 3- ...
- dev set does NOT fit well on cost function but training set fits
 - 1- using regularization
 - 2- bigger training set to make sure the algorithm generalize better
 - 3- ...
- test set does NOT fit well on cost function but dev set fits
 1. bigger dev set might be needed due to overfitting on dev set
 2. ...
- does not perform well on real world but works great on test set
 - 1- changing dev set
 - 2- changing cost function
 - 3- ...

Note:

Try not to use early stopping. Why?

Because it both works on how well we are doing on the training set and dev set. The former happens since we stop early on training set and we are not letting the model to be trained on the training set very well. The latter happens since it improves the dev set performance. It is less orthogonal since it is a knob to control two things simultaneously.

Chain of assumptions in ML



Single number evaluation metric

For choosing the best model, we have to choose the correct metric which is the best if it can be a single valued number.

As an example, in the slide below, we are looking at *precision* and *recall*:

- precision: Of the examples that our classifier recognizes as cat, what percentage are actually cats
- recall: from all of cats images, what percentage of actual cats are correctly recognized

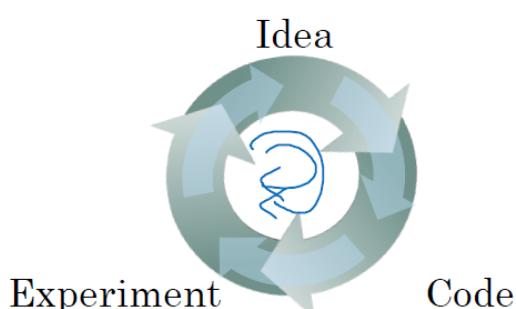
We can see in this example, recall of the classifier A and precision of the classifier B outperform the others. So which one now is better? Classifier A or B?

For this example, we can use *F1 score* which is the *harmonic mean* of precision and recall. Then we will end up with one single metric, with one we can iterate faster as for our dev set we know what we want to increase and what is exactly important for us. In this case it is F1 score. The metrics should always be a single real number for evaluation.

For the case of F1 score with harmonic mean we have:

$$F_1 \text{ Score} = \frac{2}{\frac{1}{P} + \frac{1}{R}} = \frac{2PR}{P + R}$$

Using a single number evaluation metric



→ of examples recognized as cat,
what % actually are cats?
→ what % of actual cats
are correctly recognized

Classifier	Precision	Recall
A	95%	90%
B	98%	85%

F_1 score = "Average" of P and R.

$$\left(\frac{2}{\frac{1}{P} + \frac{1}{R}} \right) \cdot \text{"Harmonic mean"}$$

Dev set + Single number evaluation metric
real speak up iterating

Andrew Ng

As another example we have several classifiers, each with values of error percentage for different countries. It is really hard to say which one of the classifiers works the best. A technique would be taking the error average of each classifier.

Algorithm	US	China	India	Other
A	<u>3%</u>	7%	5%	9%
B	5%	6%	5%	10%
C	2%	3%	4%	5%
D	5%	8%	7%	2%
E	4%	5%	2%	4%
F	7%	11%	8%	12%

Satisficing and Optimizing metric

Sometimes it is hard to reach a single real number value when there are many metrics involved in the problem.

In the example below, we have two metrics, one is the accuracy of the model and the other is the running time. It is hard to come up with a single value for accuracy and running time, writing it in a linear fashion. So instead we take another approach:

- We want to maximize accuracy (\rightarrow optimizing) while the *running time is less than X* (\rightarrow satisficing metric).

So for N metrics, one of them will be optimization variable and N-1 will be satisficing variables. Instead of writing a formula to reach a real number, we write it like this as the optimization problem.

Another cat classification example

Classifier	Accuracy	Running time
A	90%	80ms
B	92%	95ms
C	95%	1,500ms

optimizing ↓ ↓ Satisficing

Cost = accuracy - $0.5 \times \underline{\text{running Time}}$

maximize Accuracy
subject to Running Time $\leq 100\text{ ms}$.

N metrics : | optimizing
 N-1 satisficing

Wakewords / Trigger words
 Alexa, OK Google,
 Hey Siri, nihao baidu
 你好 百度

accuracy.
 #false positive

Maximize accuracy.
 s.t. ≤ 1 false positive
 every 24 hours.

Another example would be accuracy of wake up detection: when someone says the trigger word, accuracy of false positive is also important \rightarrow maximizing accuracy + at most one false positive in every 24 hours of operation.

In this example, the first one is optimizing metric and second is satisfice metric.

Train/dev/test distributions

How to setup dev and test set?

The problem here is that, we have data from different countries. We may spend months on training on the dev set so that we can target the bull's eye. However, we may finally hit it but when we go to the test set, we see that the distribution is different and the result will not well generalized. In this scenario we have to change the bull's eye completely which means months of waste of time.
So, they must be from the same distribution, not each one from different parts of the world. A way is to shuffle the dev and test set which are from different distribution.

Cat classification dev/test sets

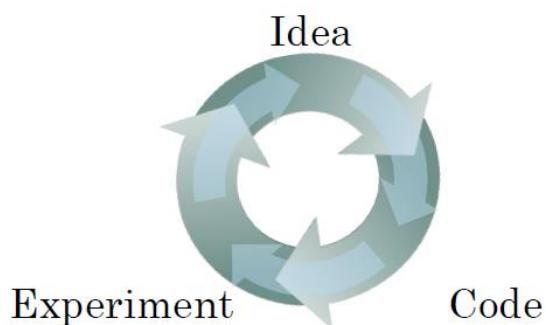
↳ development set, hold out cross validation corp

Regions:

- US
 - UK
 - Other Europe
 - South America
 - India
 - China
 - Other Asia
 - Australia
- Randomly shuffle into dev/test



dev set
+
Metric



Another example is shown in the slide below.

Optimizing on dev set on loan approvals for
medium income zip codes

\uparrow $x \rightarrow y$ (repay loan?)

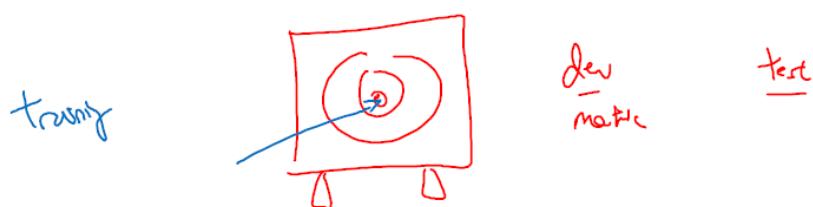


Tested on low income zip codes

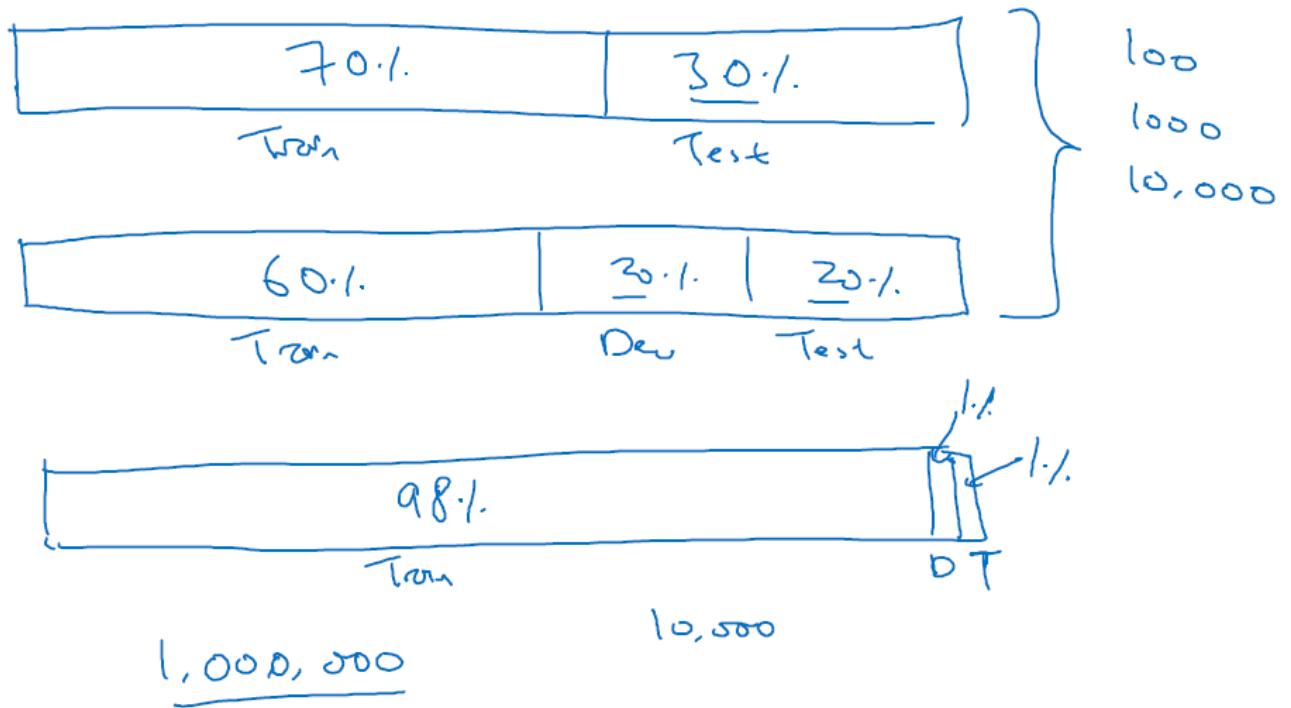
~ 3 month

Guideline

Choose a dev set and test set to reflect data you expect to get in the future and consider important to do well on.



Size of the dev and test sets



For the size of test set:

- set your test set to be big enough to give high confidence in the overall performance of your system

Not having test set is ok while having train and dev sets. But it is always recommended to have all three instead of only train and dev sets.

When to change dev/test sets and metrics

Algorithm A has some pornographic images which is not acceptable while Algorithm B having higher error with no pornographic images in. So in this example, the metric and dev set prefer algorithm A, while we prefer B. Consequently, we have to change the bull's eye.

Changing metric by giving weights to pornographic images and normal ones.

Cat dataset examples

Metric + Dev : Prefer A
You/users : Prefer B.

→ Metric: classification error

Algorithm A: 3% error → pornographic

✓ Algorithm B: 5% error

$$\left\{ \begin{array}{l} \text{Error: } \frac{1}{\sum_{i=1}^{m_{\text{dev}}} w^{(i)}} \sum_{i=1}^{m_{\text{dev}}} w^{(i)} \cdot \mathbb{I}\{y_{\text{pred}}^{(i)} \neq y^{(i)}\} \\ \rightarrow w^{(i)} = \begin{cases} 1 & \text{if } x^{(i)} \text{ is non-porn} \\ 10 & \text{if } x^{(i)} \text{ is porn} \end{cases} \end{array} \right.$$

l predicted value (0/1)

So the whole goal here with respect to orthogonalization would be:

1. define a metric to evaluate classifiers → placing the target
2. doing well on this metric → aim to shot at the target (maybe by changing the cost function)

As another example, shown in the slide below, sometimes, we are doing well on dev/test set, but the goal is to do well on user images which might be blurry and from different distribution of what we trained our model based on.

Another example

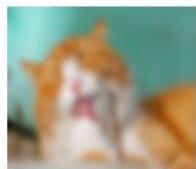
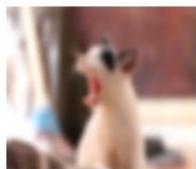
Algorithm A: 3% error

✓ Algorithm B: 5% error ←

→ Dev/test



→ User images



If doing well on your metric + dev/test set does not correspond to doing well on your application, change your metric and/or dev/test set.

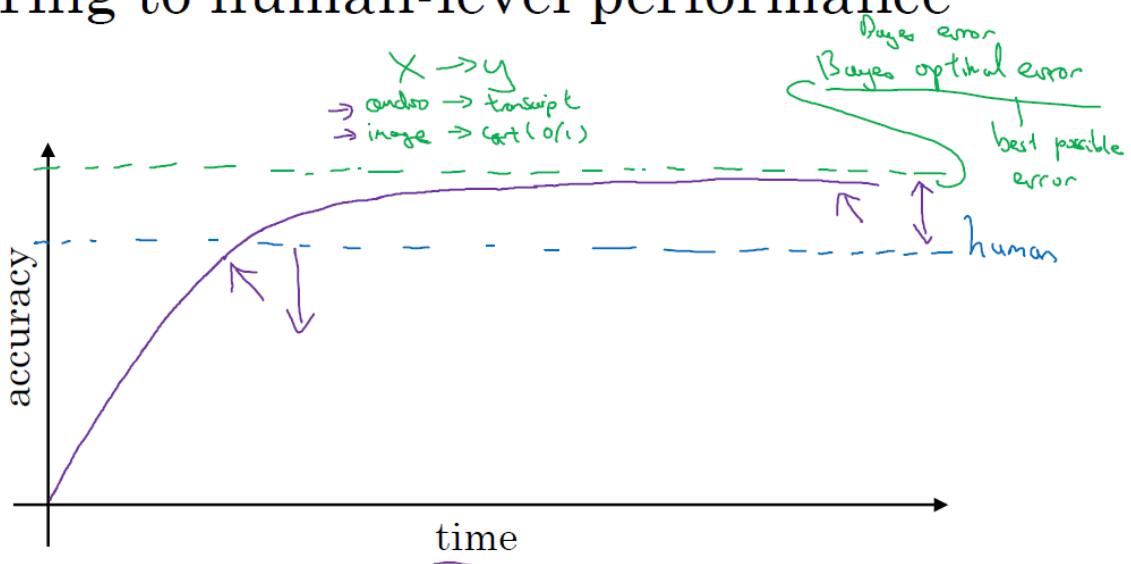
Why human-level performance?

In big data era, the accuracy of machine learning models can reach the human-level performance. Moreover, these models try to mimic the human way of solving the problems. This is the reason recently models are being compared with human-level performance.

The accuracy of machine learning models can beat the human accuracy but they cannot surpass a theoretical limit, "Bayes Optimal Error" which has the highest accuracy.

In the below diagram, we can see that when the model reaches the accuracy of human-level performance, its momentum decreases highly or maybe never reaches to that. Because the difference between the human-level performance is not much different than the Bayes error. Moreover, as long as the model didn't reach the human-level performance, there are some ways and tools to improve them, but it is not the case when the model passes the human-level performance.

Comparing to human-level performance



Why compare to human-level performance

Humans are quite good at a lot of tasks. So long as ML is worse than humans, you can:

- - Get labeled data from humans. (x, y)
- - Gain insight from manual error analysis: Why did a person get this right?
- - Better analysis of bias/variance.

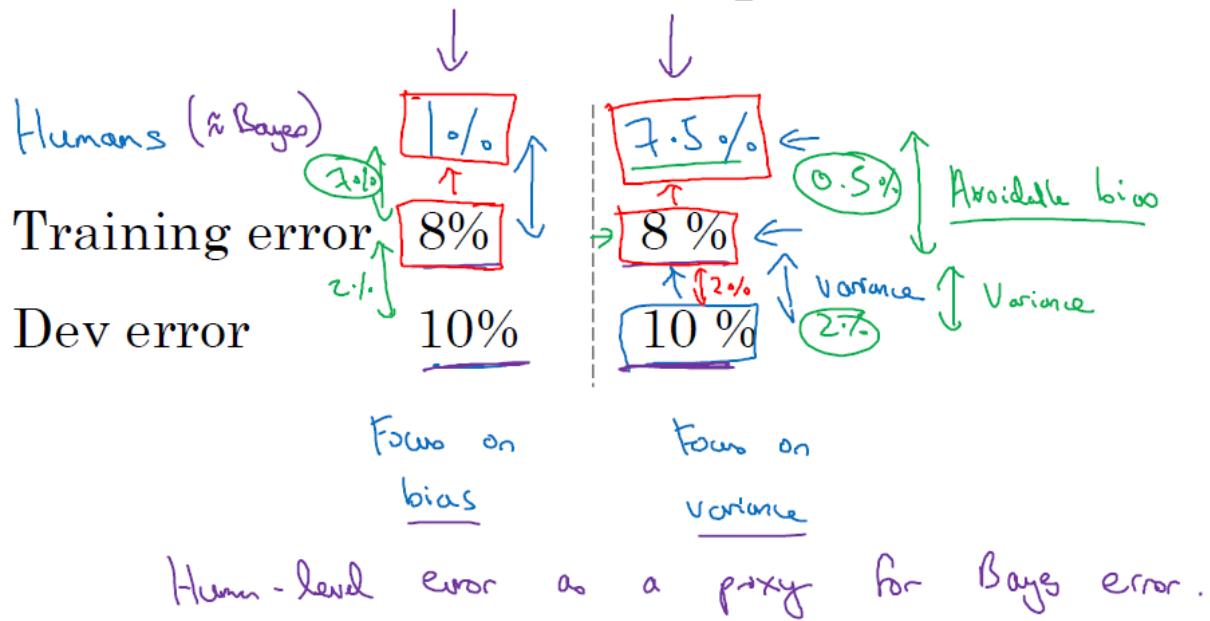
Avoidable bias

Previously, we assumed the human error is better than the Bayes error which is 0%, however it is not true. → we have an avoidable bias.

Human-level error is a proxy for Bayes error in most of the cases.

Difference between the training error and Bayes error is “avoidable bias” since we cannot do better than Bayes error. We need to get as close as possible to the Bayes Error and then try to reduce difference error between training and dev error which is variance.

Cat classification example



Human-level error as a proxy for Bayes error

Medical image classification example:



Suppose:

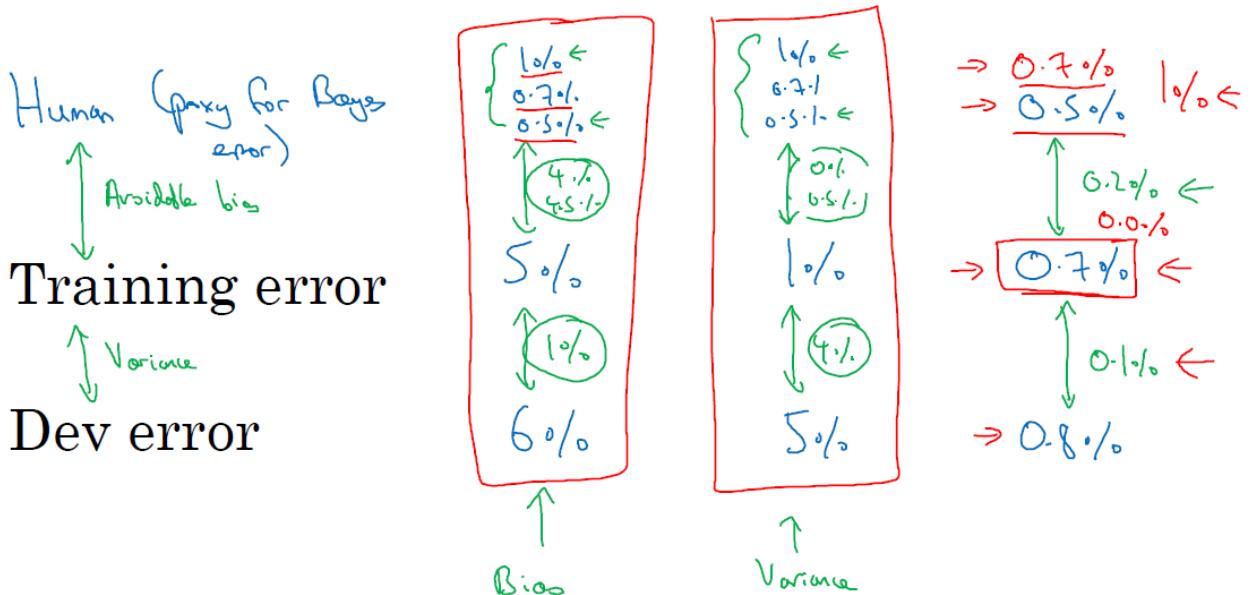
- (a) Typical human 3 % error
- (b) Typical doctor 1 % error
- (c) Experienced doctor 0.7 % error
- (d) Team of experienced doctors .. 0.5 % error

What is “human-level” error?

Choice D is close to Bayes error so we can assume it as human-level error. We have also this concept that if we pass the average doctor diagnosis accuracy, it is human-level error. In that case choice B is human-level error but it is not a proxy to Bayes error.

In the example below, in the left column, it doesn't matter which proxy we choose as far as its difference with training error is so much higher; so we want to reduce bias. The next one same applies for variance. Finally, the last one, we should be concerned about the proxy we use as it affects which one of bias or variance we should focus on.

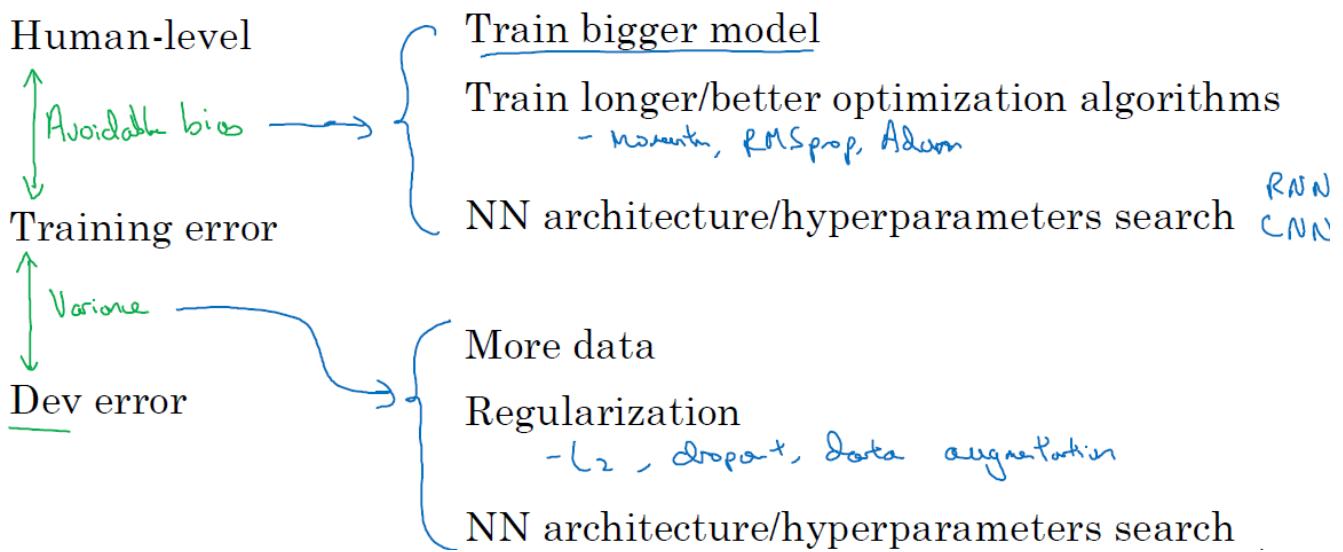
Error analysis example



There are many cases that computer can surpass human-level performance due to limitation of human in learning patterns. Like voice recognition.

Improving your model performance

Reducing (avoidable) bias and variance



Course 3

Week 2

Carrying out error analysis

Does working on something new worth the effort and time we put on it?

So let's say that we want to come up with a dog classifier. In the *error analysis* phase, we see around 100 pictures of dogs are mislabeled in the dev set. Then we count up how many of the mislabeled pictures are dog pictures. If it is only 5% of them, it means that if the error of dog classifier was 10%, after fixing the mislabeled data in dev set, the error goes down to 9.5% (it is relative error). On the other hand, if 50% of mislabeled dev set examples are dog, from 10% error we can reduce it to 5%. In other words, we can find the ceiling in this way and understand whether it is worth our time or not.

Look at dev examples to evaluate ideas



90% accuracy
→ 10% error

Should you try to make your cat classifier do better on dogs? ↩

- Error analysis: → 5-10 min
- {
- Get ~100 mislabeled dev set examples. "ceiling"
 - Count up how many are dogs.
- 5% 10%
5/100 9.5% | → 50%
 50/100 10%
 ↓ ↓
 5%.

Evaluate multiple ideas in parallel

Ideas for cat detection:

- Fix pictures of dogs being recognized as cats ↩
- Fix great cats (lions, panthers, etc..) being misrecognized
- Improve performance on blurry images ↩

Image	Dog	Great Cats	Blurry	Instagram	Comments
1	✓				✓ Pitbull
2			✓	✓	
3		✓	✓		Rainy day at zoo
:	:	:	:		
% of total	8%	43%	61%	12%	

Cleaning up incorrectly labeled data

As far as the number of samples are too much and the errors are not much, for the training set we can say that as DL algorithms are quite robust to random errors in the training set, we do not need to correct the mislabelled data. However, they are not robust at all to the systematic errors.

For dev set and test set → in the error analysis table we add a column of incorrectly labeled. If it makes a significant difference to your ability to evaluate algorithms on your dev set, so fix it. But if its difference is not significant, just let it go.

For the fixing of mislabelled data, as it is time consuming, in the left example does not help much; but, in the right one it does help.

Error analysis

Image	Dog	Great Cat	Blurry	Incorrectly labeled	Comments
...					
98				✓	Labeler missed cat in background
99		✓			
100				✓	Drawing of a cat: Not a real cat.
% of total	8%	43%	61%	6%	

Overall dev set error 10%

Errors due incorrect labels 0.6% ←

Errors due to other causes 9.4% ←

↑

$\begin{array}{r} \checkmark \\ 2\% \\ \hline 0.6\% \\ \hline 1.4\% \\ \hline 2.1\% \end{array}$ 1.9%

Goal of dev set is to help you select between two classifiers A & B.

Now some guidelines in correcting the mislabelled data:

1. First bullet: dev and test must come from the same distribution. So everything we do with dev set, same must be done on test set.
2. Third bullet: training set may come from slightly different distribution, but as far as dev and test sets are from the same distribution, it is ok as the DL models are quite robust to slight changes in distributions.
3. Second bullet: the examples of correct prediction might get wrong by changing the correcting a label. But if the model is accurate enough, it may not happen very often or if it happens, it is not much prevalent.

Correcting incorrect dev/test set examples

- Apply same process to your dev and test sets to make sure they continue to come from the same distribution
- Consider examining examples your algorithm got right as well as ones it got wrong. {
- Train and dev/test data may now come from slightly different distributions.

Build your first system quickly, then iterate

1. Quickly setup dev/set and metric
2. Build initial system quickly
3. Use bias/variance analysis and error analysis to prioritize next steps.

This is for systems that we do not have much literature for it. If we have some, then it is ok to make a more complex system.

Training and testing on different distributions

We have some pictures from web and some from mobile app. We want our model works the best for the mobile app. If we combine them and shuffle them, then split it to train, dev, and test sets, our data in dev and test might be more about web rather than mobile, so the target is not right due to different distributions of dev, test, and train sets although shuffling them urge them to come from the same distributions. If dev and test are from different distributions that is also wrong.

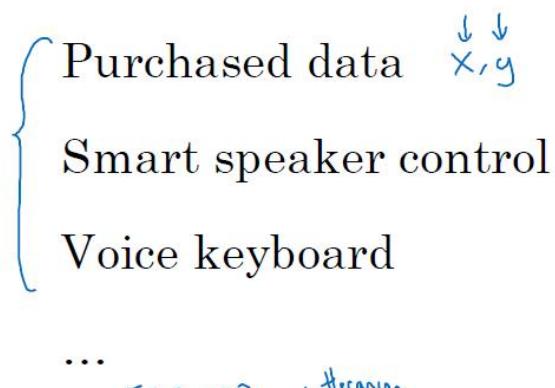
ANSWER is to use mobile app pictures for dev and test and put the remaining for the training set. → in the long run it works better although the distributions of training with the rest is different but we come up with techniques so that the model works the best on the dev set.

Speech recognition example

Speech activated rearview mirror



Training



Dev/test

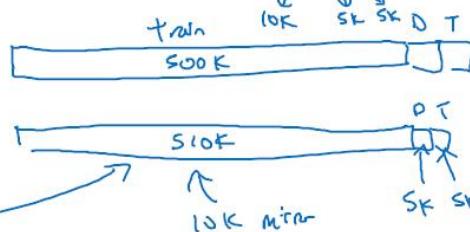
Speech activated
rearview mirror

→ 20,000

train
500K

510K
10K mirr

PT
SK SK



Bias and Variance with mismatched data distributions

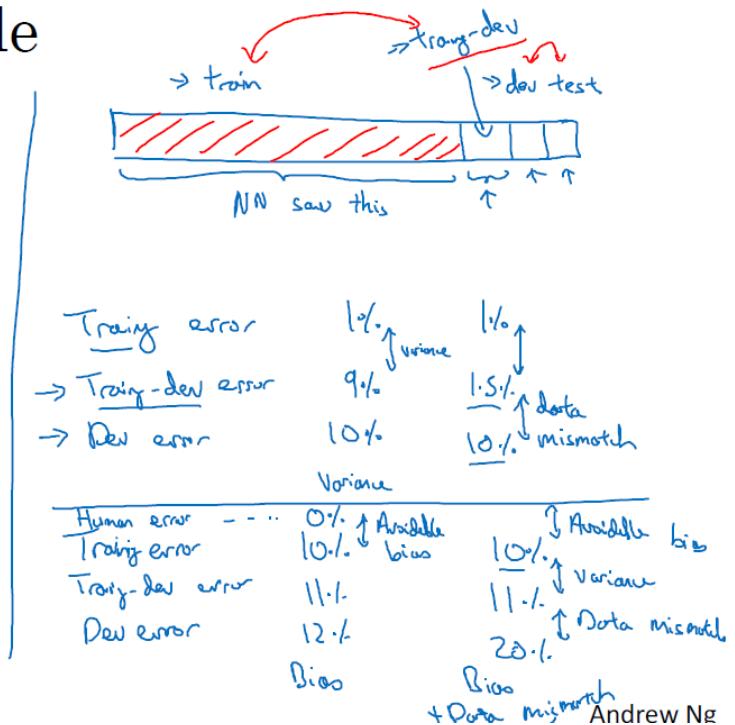
If training error is low, but dev error is high, the reason might be due to different distributions of training and dev or maybe the dev set is harder for the model to recognize. So we use a training-dev set from the training distribution as the new dev set. Then we compare the results on dev set and training-dev set. Now, if the difference of training and train-dev set is high, it means the problem is only variance because there is no difference in distributions anymore. If the difference of train-dev set and dev set is high, we call it data-mismatch problem and happening due to different distributions.

Cat classifier example

Assume humans get $\approx 0\%$ error.

Training error 1% $\downarrow 9\%$
 Dev error 10%

Training-dev set: Same distribution as training set, but not used for training



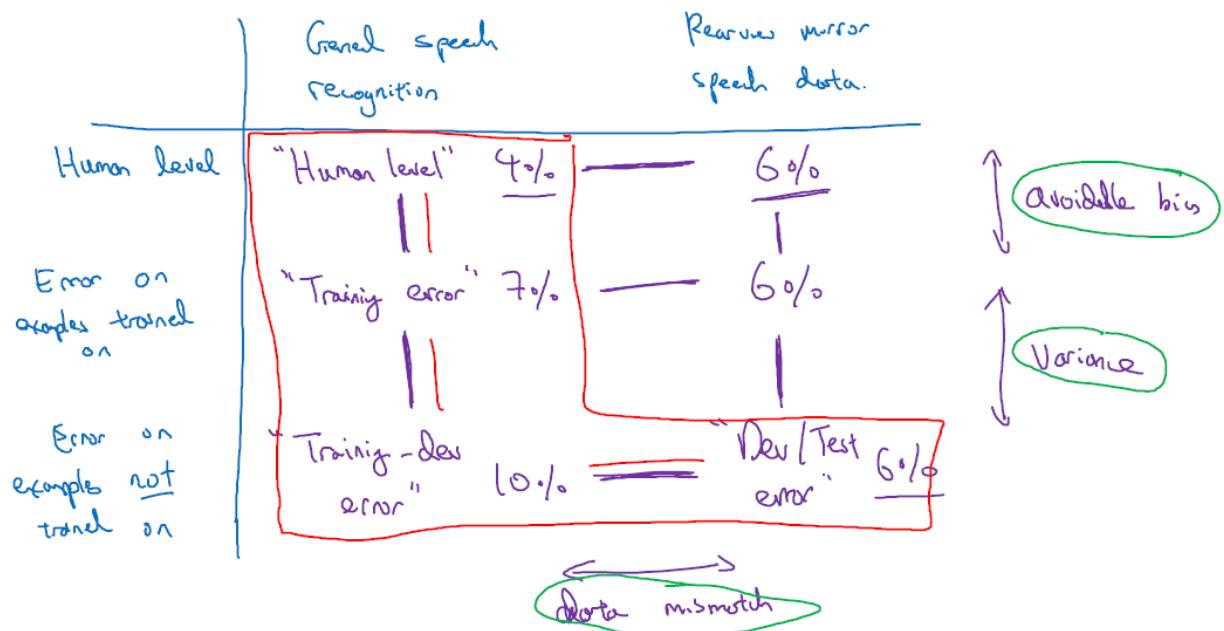
The errors on dev and test set can go down which means their examples are so much easier for the learner to recognize than the training set examples.

Bias/variance on mismatched training and dev/test sets

Human level	
Training set error	4% ↑ available bias
Training-dev set error	7% ↑ Variance
→ Dev error	10% ↑ Data mismatch
→ Test error	12% ↑ degree of difficulty to dev set.

4%	
7%	}
10%	
6%	
6%	

More general formulation



Addressing data mismatch

There is no systematic way to handle it but what we can do is:

Addressing data mismatch

- • Carry out manual error analysis to try to understand difference between training and dev/test sets

E.g. noisy - car noise street numbers

- • Make training data more similar; or collect more data similar to dev/test sets

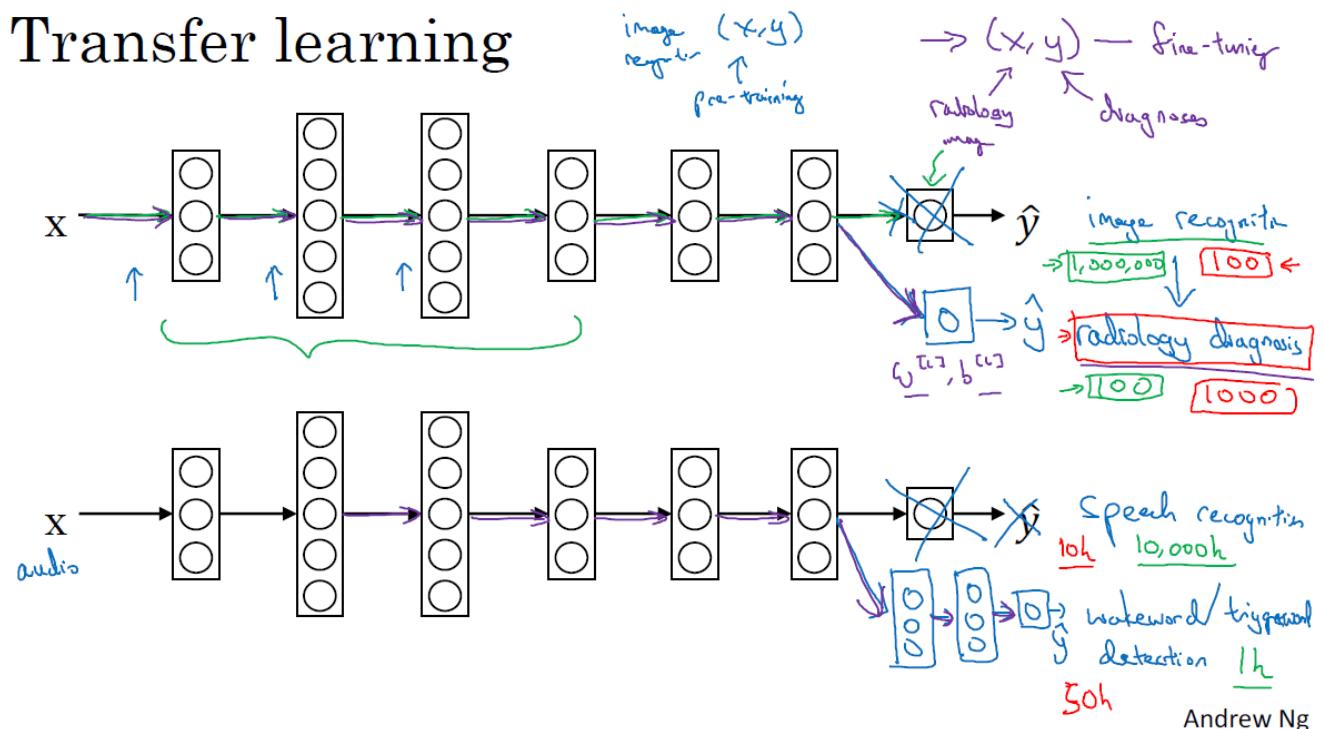
E.g. Simulate noisy in-car data

One of the ways to make it more similar is creating artificial data synthesis. E.G. adding car noise to a sentence to make it more similar to dev data with car noise.

Transfer learning

We might want to take advantage of trained weights to use it on another similar task that we do not have much data on it. So we reinitialize the last layer (for example) of the NN weights and run it for the new data set. Some of its knowledge like detecting edges and some shapes are the same, so we can take advantage of them.

Transfer learning



Transfer learning makes sense when:

When transfer learning makes sense

Transfer from A \rightarrow B

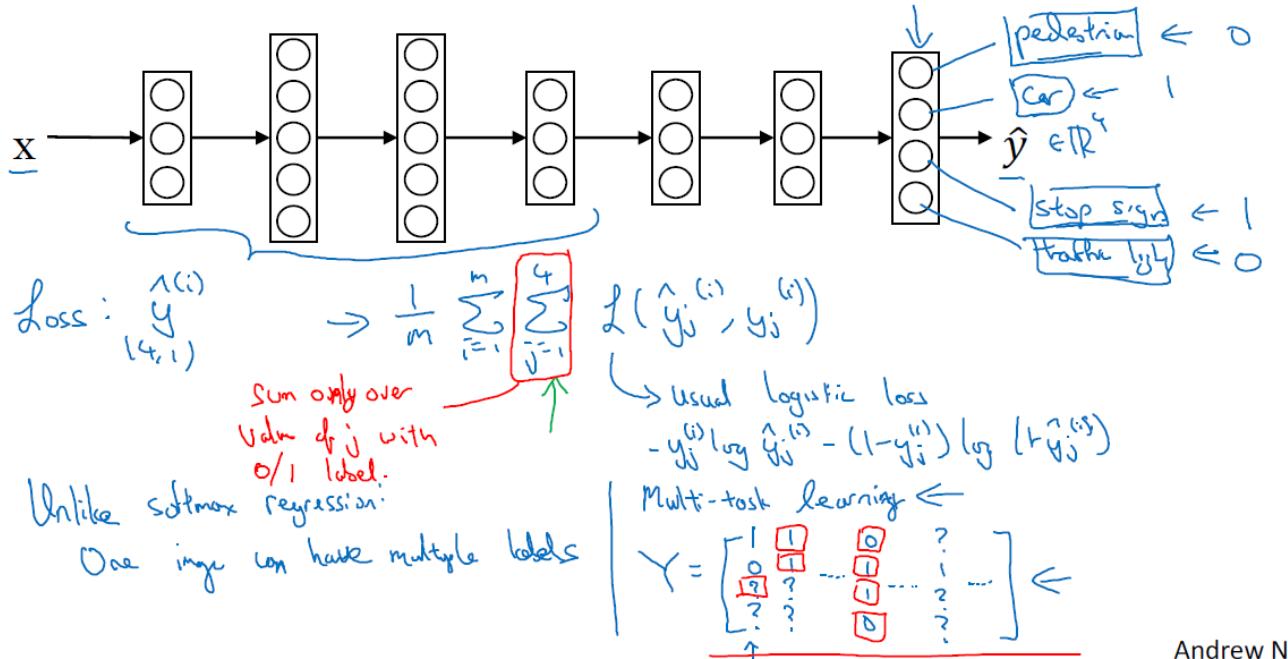
- Task A and B have the same input x .
- You have a lot more data for Task A than Task B.
- Low level features from A could be helpful for learning B.

Multi-task learning

Solving several problems in a single architecture NN like a multi-label task → 4 different NNs
 Note that the last layer can have multiple labels, so it cannot be Softmax.

Note that it is also acceptable to not have all the labels

Neural network architecture



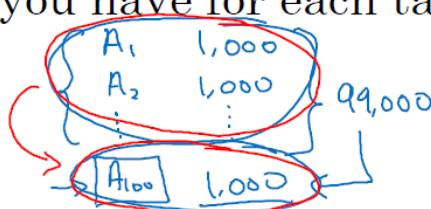
Andrew Ng

It makes sense when:

When multi-task learning makes sense

- Training on a set of tasks that could benefit from having shared lower-level features.
- Usually: Amount of data you have for each task is quite similar.

$$\begin{array}{ll} A & 1,000,000 \\ \downarrow & \downarrow \\ B & 1,000 \end{array}$$

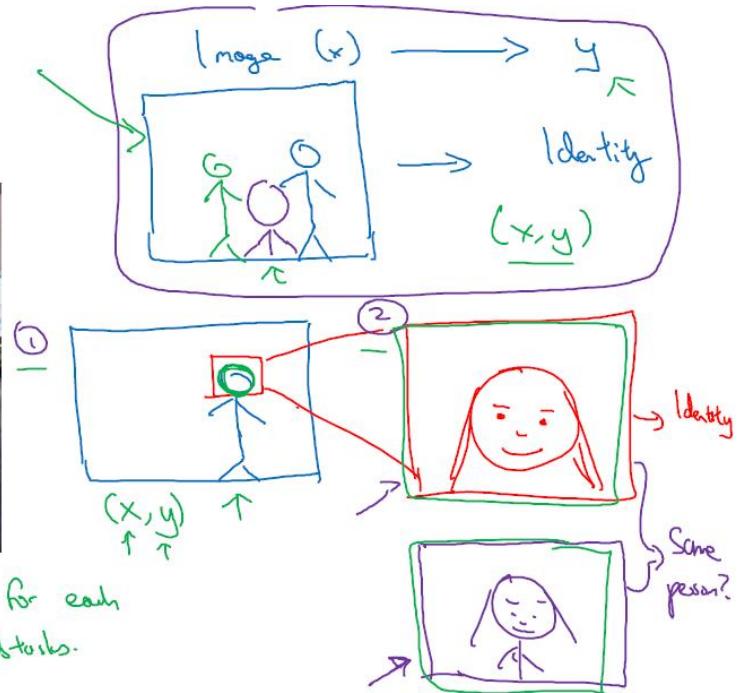


- Can train a big enough neural network to do well on all the tasks.

What is end-to-end deep learning?

Breaking down the problem into two simplified problems which ends into better overall performance.

Face recognition

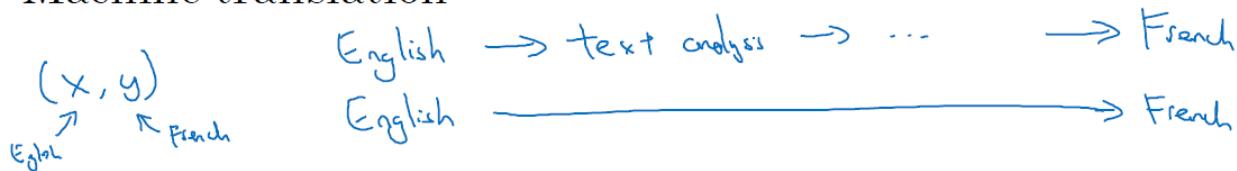


We have a lot of data for each task here (I) recognizing face from an image (II) identifying person, but we have not much data for the mixed version.

For machine translation but as we have a lot of data for translating a language to other one, so we do not need to simplify the model and we can have an end-to-end translation. It does not work well for estimating child's age by the x-ray.

More examples

Machine translation



Estimating child's age:

Image $\xrightarrow{①}$ bones $\xrightarrow{②}$ age

Image \longrightarrow age \leftarrow



Andrew Ng

Whether to use end-to-end deep learning

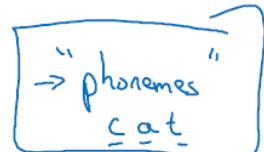
First bullet: we do not want to think about phonemes when instead we can pronounce the whole word in once. While if we have less data it might be more useful to let the model learn on its own so it might find better representative of data.

The last bullet is about injecting manual information into network and can be so much helpful in many places.

Pros and cons of end-to-end deep learning

Pros:

- Let the data speak $x \rightarrow y$
- Less hand-designing of components needed



Cons:

- May need large amount of data $\underline{x} \rightarrow \underline{y}$ $(\underline{x}, \underline{y})$
- Excludes potentially useful hand-designed components

Data. Hand-design.

$$x - - - - \rightarrow y$$

