

Course 2:

Week 1 Train/dev/Test Sets

- Development set is often used to test the trained model!
- For example, If we have two models and we need to determine which will work better, we test on dev set
- the dev and test sets must come from the same distribution

Bias/Variance

Train set error: 1%

Dev set error: 11%

in this case, high variance

overfitting

15% cannot fit the training set

16% as well as dev set

high bias

underfitting

Train set error: 15%

0.5%

Assuming optimal
error is 0%.

Dev set error: 30%

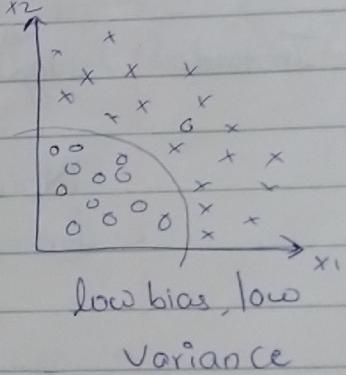
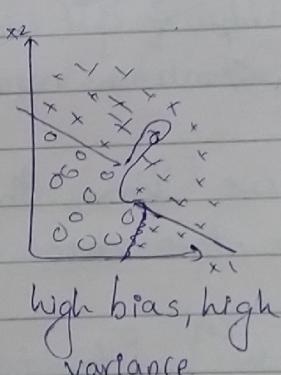
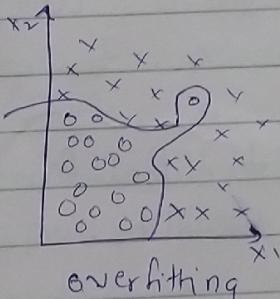
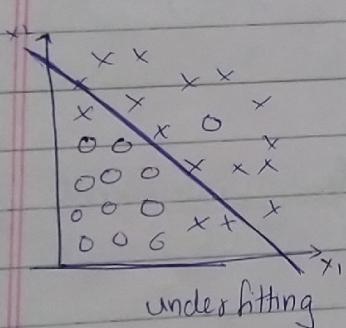
1%

high bias

low bias

high variance

low variance



High bias?
(training data problem)

- Bigger n/ω
- Train longer
- NN architecture

High variance
(dev set problem)

- More data
- Regularization
- NN architecture

previously → bias-variance trade-offs

modernly → big data & big n/ω \Rightarrow is possible
low bias low variance

Regularization

(helps over fitting)

In logistic Regression,

we tried to minimize the cost function i.e $\min_{w,b} J(w,b)$
i.e

$$J(w,b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \lambda \frac{\|w\|_2^2}{2m} \quad \rightarrow \lambda = \text{regularizat'n parameter}$$

$$L_2 \text{ Regularizat'n } \|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w$$

$$\text{for } L_1 \text{ Regularization } \frac{\lambda}{2m} \sum_{j=1}^{n_x} |w_j| = \frac{\lambda}{2m} \|w\|_1 \quad \rightarrow w \text{ will be sparse}$$

λ will be Regularizat'n parameter

In neural Network

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

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

Instead of L_2 regularization "Forbenius norm"

Now we know that to calculate d_w , we need $\frac{\partial J}{\partial w}$,
for in this case, it will be:

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

$$\therefore w^{[l]} := w^{[l]} - \alpha d_w^{[l]}$$

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

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

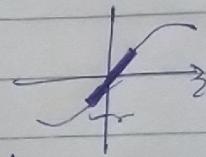
because $w^{[l]}$ is multiplied by $\frac{\alpha \lambda}{m} I$ and we subtract it

Why does Regularization help in less overfitting?

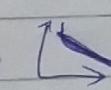
intuition-1: It lessens the impact (weights) of some hidden units.

So if they are hampering in the o/p it reduces

Thus regularization penalizes weights being too large.

So, in intuition, it can be seen that if $\lambda \uparrow$, then $w^{[l]} \downarrow$ which means $z^{[l]} \downarrow$ ($z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]}$).
 For function ~~tanh~~, $g(z)$ will roughly be linear 

Hence, the whole network will become linear. If such is the case, the NN will not be able to compute complex decisions.
 So, if regularization becomes large, $w \nmid$, activation=tanh \rightarrow overfit

A debugging tool is to plot J vs #iterations. It should be  if not, regularization is needed

Dropout Regularization

Consider that each layer has 0.5 chances of keeping the nodes and eliminating the nodes. For a particular example, randomly $n^{[l]}/2$ nodes ~~will~~ will be chosen from each layer L . All the incoming and outgoing ~~nodes~~ ^{weights} from these nodes will be eliminated. Thus, we will be training a small network. For the next set of examples, another set of nodes will be removed.

Implementing ~~dropout~~ ("Inverted dropout")

layer $d=3$, keep_prob=0.8 (80% prob of keeping neurons)

$d3 = np.random.rand(a3.shape[0], a3.shape[1]) < keep_prob$

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

$a3/ = keep_prob$ \leftarrow is used to keep the expected value of $a3$ as it is

e.g. 50 units \rightarrow 10 units shut off

$$z^{[4]} = w^{[4]} a^{[3]} + b^{[4]}$$

\uparrow reduced by 20%

disadvantage: The size of $a3$ remains same during training, no need to add $a3/ = ?$.

The results are random. You don't want random results at test time

Trying different combinations of neurons is computationally inefficient and will also lead to similar o/p.

Why does drop-out work?

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

L₂-regularization shrinks weights. But it should not only rely on particular features.

Also, you may apply different keep-prob for different layers.

Generally used in computer vision where, there are a lot of features.

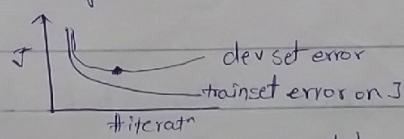
Should be used only when the model overfits

disadvantage: J is not well defined.

Cannot draw J vs #iteration's plot.

Other Regularizⁿ methods

- 1) Data Augmentation: change orientation of image & add data
- 2) Early stopping: To stop a neural n/w early when you get optimized value. that is min. J(Cost).



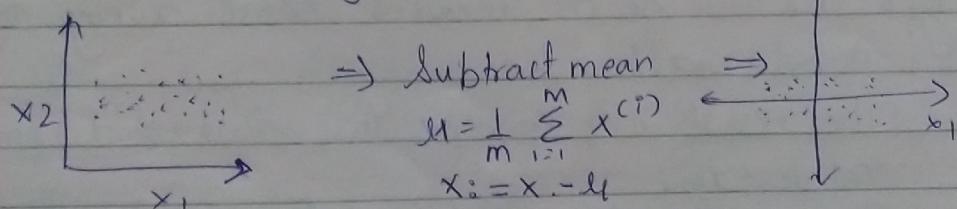
disadvantages taking 2 ml problems: optimize J & not overfit, to solve which adds complexity.

Advantage: With a single go over the gradient descent, we can try small, med. & large vals. of w.

Normalizⁿ of Inputs

Adv.: Speed up training of model

$$X = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$



Normalizing variance

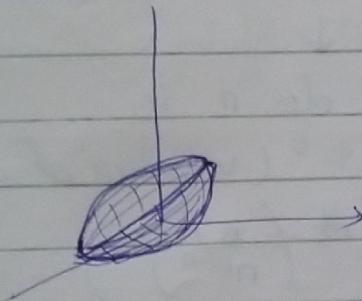
$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m x^{(i)} \cdot x^{(i)}$$

elementwise

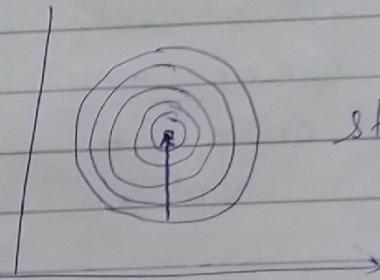
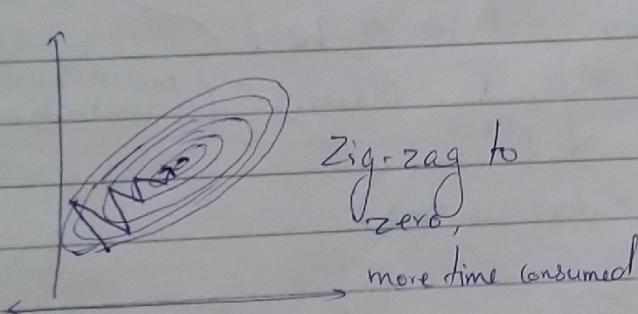
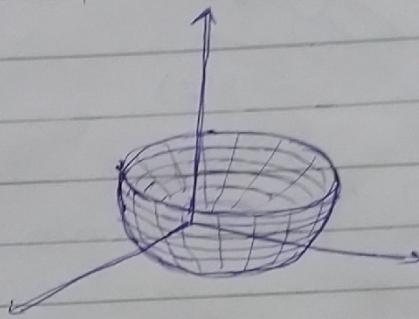
$$x / \sigma^2$$

Use the same σ & σ^2 to normalize your test data, do not use a new one

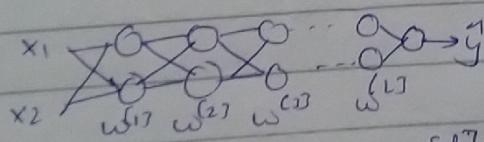
Unnormalized:



Normalized:



Exploding/Vanishing gradient



$$g(z) = z \quad b^{[L]} = 0$$

$$\hat{y} = w^{[L]} w^{[L-1]} \dots w^{[2]} w^{[1]} w^{[0]} x$$

$f = g(z^{[0]}) = z^{[0]}$

$w^{[0]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$

$\hat{y} = w^{[L]} \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}^{L-1} x \rightarrow 1.5^{L-1} x$

if it were 0.5, it were vanishing

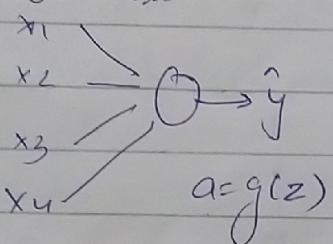
expanded

If the gradients are exponentially small or big, it is difficult to train based on gradient descent, as it will take tiny steps.

To solve this, careful weight initialization provided as a partial solutⁿ

Weight Initialization for Deep Networks

Overcomes



$$z = w_1x_1 + w_2x_2 + \dots + w_nx_n$$

in order to keep z in a normal range, if $n \uparrow \Rightarrow w_i \downarrow$.

One way to do it is,

$$\text{Var}(w_i) = \frac{1}{n} \quad (\text{Variance of weights is set})$$

$$w^{[l]} = np.random.randn(\text{shape}) * np.sqrt\left(\frac{1}{n^{[l-1]}}\right)$$

in case of relu function, it is 2 here

in case of tanh functⁿ $\sqrt{\frac{1}{n^{[l-1]}}}$ is used. (called xavier initialization)

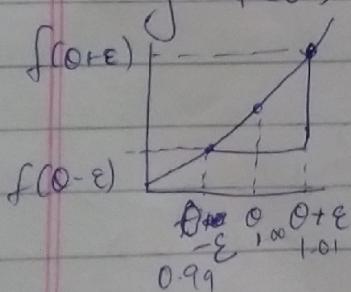
$$\text{or } \sqrt{\frac{1}{n^{[l-1]} * n^{[l]}}} \quad (\text{some authors use this})$$

You could also add another parameter to this variance, for hyper parameter surge. It might have a side effect too

Train NN's faster

Checking your derivative computatⁿ

$$\text{e.g. } f(x) = x^3$$



$$g(x) = \text{slope} = \frac{f(x+\varepsilon) - f(x-\varepsilon)}{2\varepsilon} = \frac{3.001 - 0.999}{2 \cdot 0.01} \approx 3$$

$$\text{Actual } g(x) = 3x^2 = 3 \quad \text{correct}$$

Gradient Checking: Tool to find bugs in back propagation.

Take $w^{(1)}, b^{(1)}$... $w^{(L)}, b^{(L)}$ & reshape into a big vector θ
Concatenate

$$J(w^{(1)}, b^{(1)} \dots w^{(L)}, b^{(L)}) = J(\theta) = J(\theta_1, \theta_2, \theta_3)$$

Similarly, $d\theta^{(1)}, db^{(1)} \dots d\theta^{(L)}, db^{(L)} \Rightarrow d\theta$

Grad check

For $d\theta_{approx}^{(l)} = \frac{J(\theta_1, \theta_2, \dots, \theta_l + \epsilon, \dots) - J(\theta_1, \theta_2, \dots, \theta_l - \epsilon, \dots)}{2\epsilon}$

should be $\approx d\theta^{(l)} = \frac{\partial J}{\partial \theta_l}$

To check this, we take euclidean vector b/w.

$$d\theta_{approx}^{(l)} \& d\theta^{(l)} = \frac{\|d\theta_{approx}^{(l)} - d\theta^{(l)}\|_2}{\|d\theta_{approx}^{(l)}\|_2 + \|d\theta^{(l)}\|_2}$$

Consider we take ϵ of order 10^{-7}
in that case, if o/p is $\approx 10^{-7}$ great!
 $\approx 10^{-5}$ ok
 $\approx 10^{-3}$ wrong

Tips

- 1 Don't use it in training - only to debug
- 2 If algo. fails grad check, look at components to try to identify bug
- 3 Remember Regularization $d\theta$ = gradient of J (add λ thing)
- 4 Doesn't work w/ dropout \rightarrow in such cases, keep, keep-prob=1.0 do grad check and then put \div anything
- 5 Run at random initialization; perhaps again after some training. (Not use oftenly)

Optimization Algorithms & train your NN faster

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(1000)} & x^{(1001)} & x^{(2000)} & \dots & x^{(m)} \\ (n_x, m) & x_{\Sigma 1^3} & & x_{\Sigma 2^3} & & & & \end{bmatrix}$$

Now, in case if $m = 5\text{million}$, then it will take so long to train.

In such a case, we make baby training sets. we call them mini batches.

For example, lets make size of mini batch as 1000, we name it $x^{\Sigma 1^3}$. Here, 5000 such mini batches will be made. Simillarly we do for $y, (y^{\Sigma 1^3}, y^{\Sigma 2^3}, \dots)$

Size of $x^{\Sigma 1^3}$ is $(n_x, 1000)$