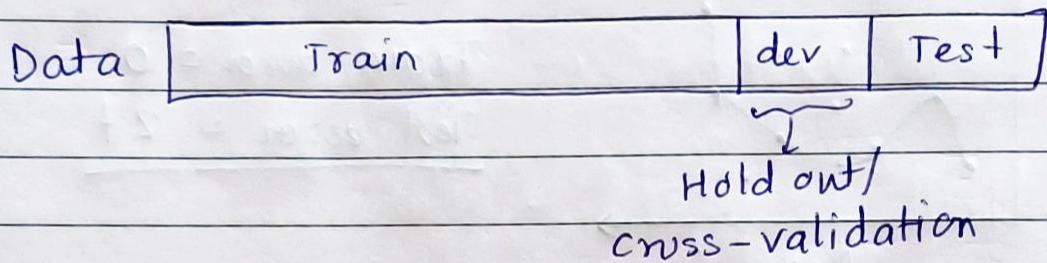


Date : 22/7/24
MON TUE WED THU FRI SAT SUN

Course 2: Improving Deep Neural Networks

Train / dev / test sets

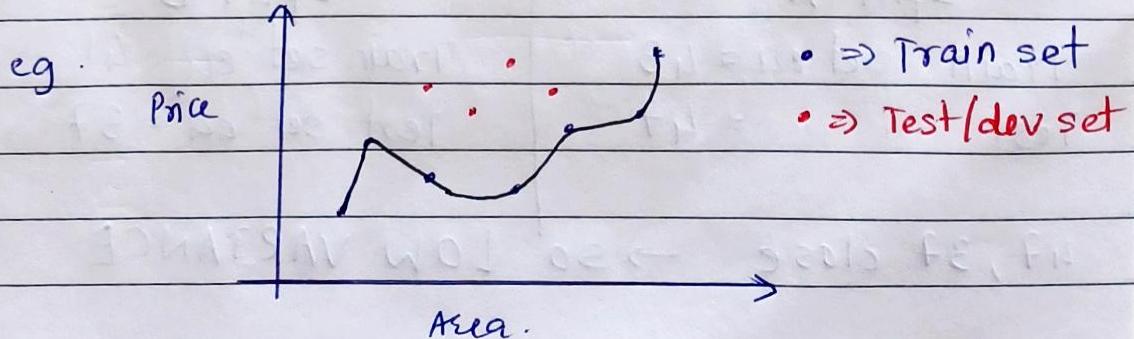


When data very very big: 98% / 1% / 1%
train dev test

- * dev & test sets come from same distribution.

Idea
Experiment
Code

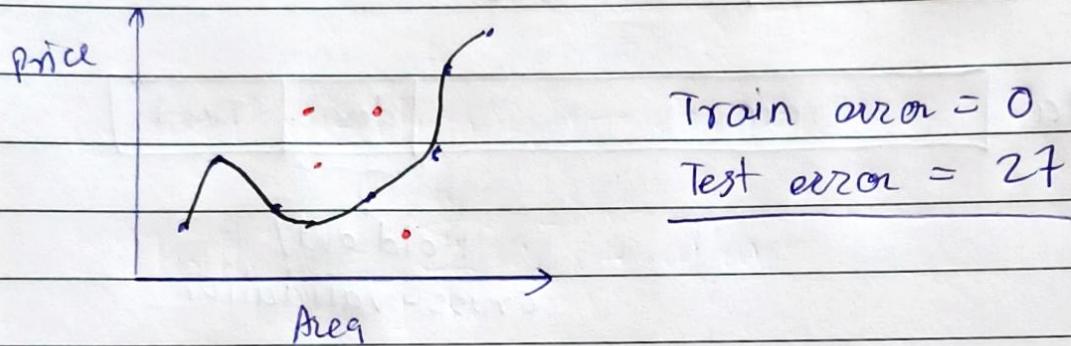
Bias and Variance



Train set error = 0

Test set error = 90 (let)

Now same data but split is different.



→ Test error varies greatly based on your selection of train & test data points.

→ This is an **OVERFIT Model**

→ High variance

If we just do a Linear Regt. in prev. eg.

Split A

Train set err = 43

Test = 47

Split B

Train set err = 41

Test set err = 37

47, 37 close → so LOW VARIANCE

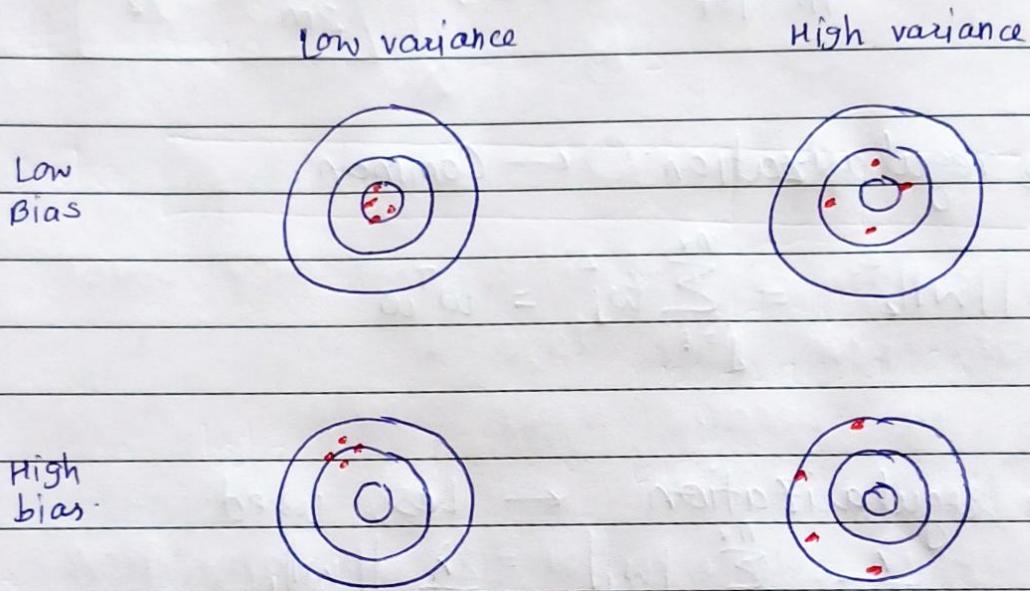
Bayes error = Human's error.

Date :

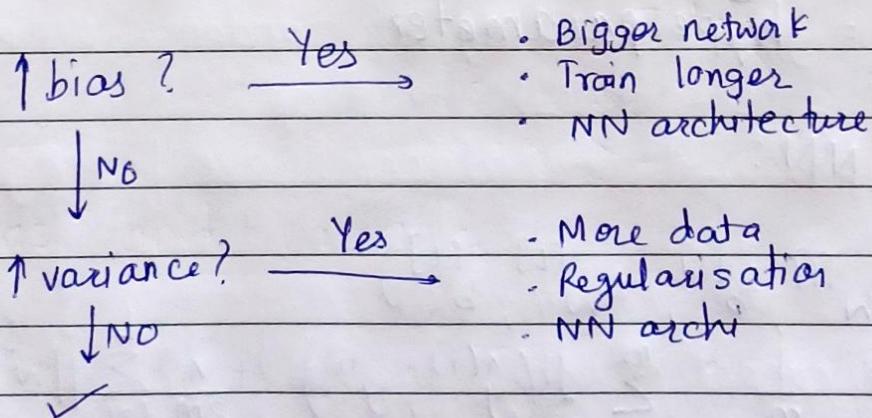
MON TUE WED THU FRI SAT SUN

But huge train err \Rightarrow HIGH BIAS

BIAIS: measurement of how accurately a model can capture a pattern in a training dataset.



Basic recipe for ML



$\|w\| \rightarrow$ norm Frobenius norm

Date :
MON TUE WED THU FRI SAT SUN

Regularization (for \uparrow variance)

e.g. In Logistic Regi:

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

L_2 regularization: ← Common.

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

L_1 regularization ← Less used:

$$\frac{\lambda}{2m} \sum_{j=1}^{n_x} |w_j| = \frac{\lambda}{2m} \|w\|_1$$

λ = regularization parameter.
= hyperparameter.

For NN:

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

Date :

MON TUE WED THU FRI SAT SUN

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

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

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

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

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

Explanation:

By adding the L1 or L2 regul.
we are adding a penalty to the
cost function.

Thus $\lambda \uparrow \rightarrow w^{[l]} \downarrow \rightarrow \text{NN becomes simpler}$

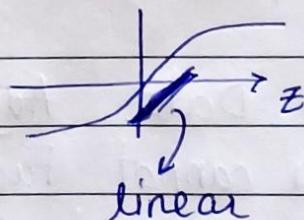
Also, if AF is tanh

$$w^{[l]} \downarrow \Rightarrow z^{[l]} \downarrow$$

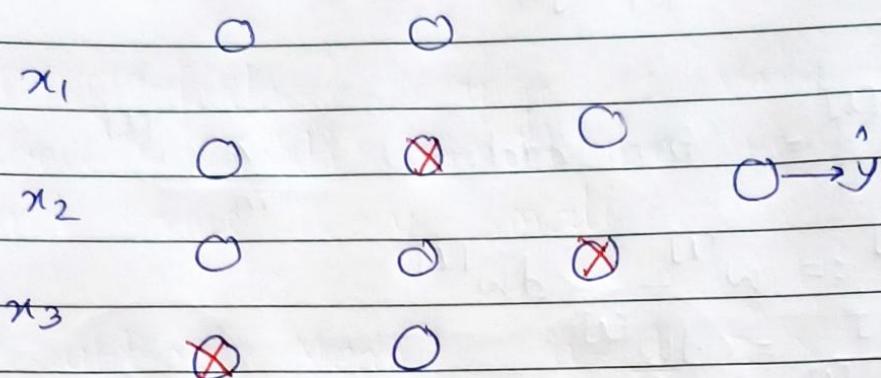
\therefore As if every layer

\approx linear (or linear regression)

\Rightarrow Model will not OVERFIT



→ Dropout Regularization ←



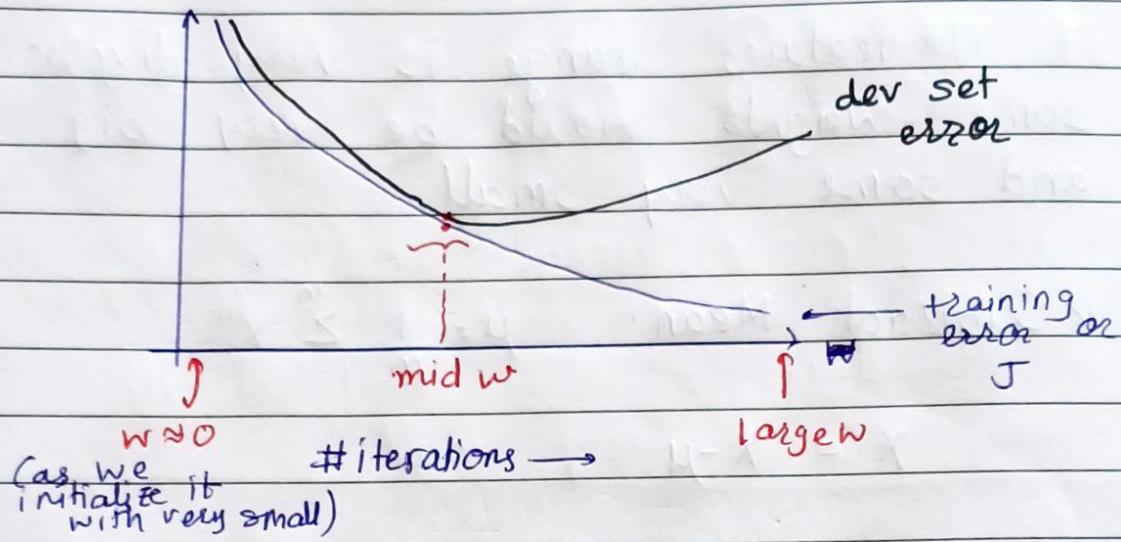
By random probabilities, we drop some nodes from each layer l
So NN became much simpler.

- * Dropout efficiency \approx L2 regularization
- * Gen^{ly} we don't drop input layer
- Method → we keep a threshold probability
And generate array of randoms prob.
If prob < thresh \rightarrow drop.
- * Dropout also reduces overfitting.

→ Data Augmentation ←

If u cannot further increase your training data, augment it (rotate, flip, crop, skew, images etc)

→ Early Stopping ←



The point where deflection of dev error starts, we stop training the model there.

But with this method orthogonalization
is violated

do 1 task @ 1 time

- T1 - Optimize J
- T2 - Not overfit

>this method does both @ same time.
∴ J is not reduced much.

Date :

MON TUE WED THU FRI SAT SUN



Normalizing training sets

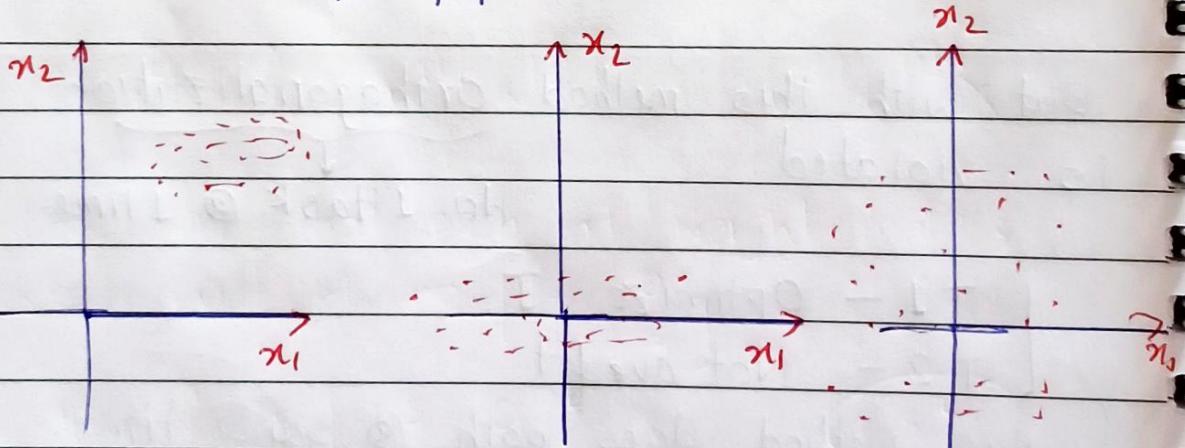
If ILP features range is very diverse:
some weights would be very big
and some very small.

▷ Subtract Mean: $\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$

$$x := x - \mu$$

▷ Normalize variances:

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



org

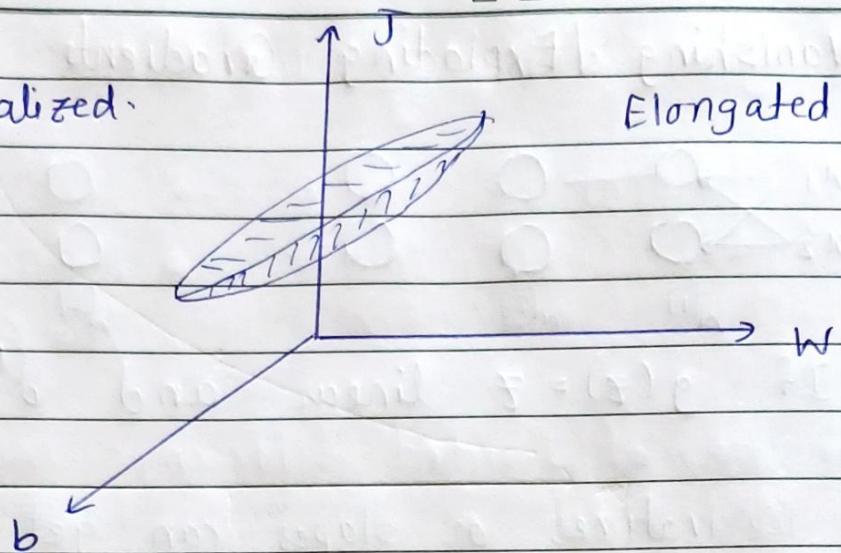
▷

2>

Date :

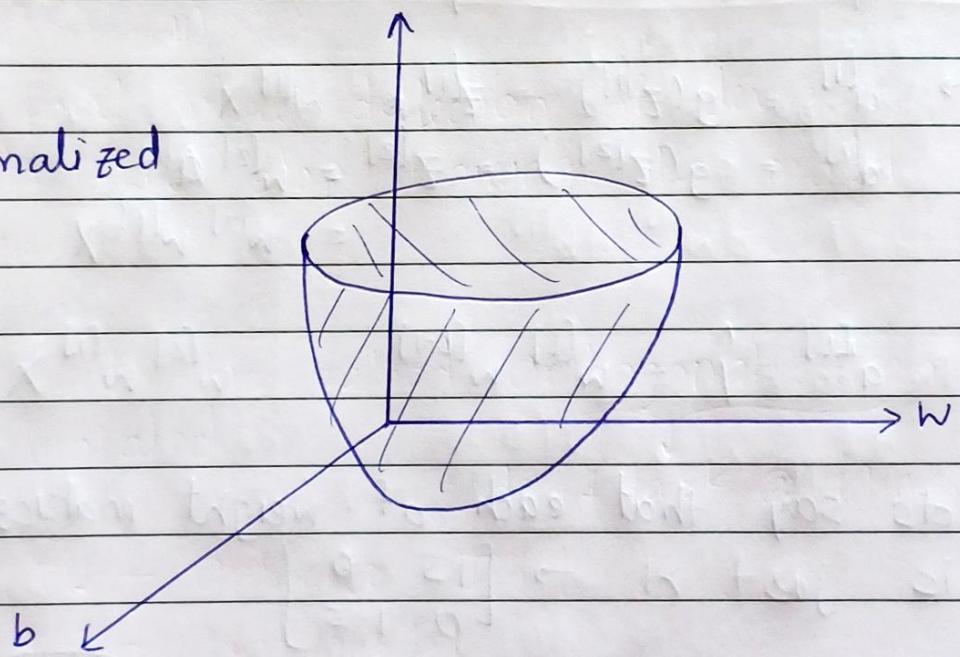
MON TUE WED THU FRI SAT SUN

Unnormalized

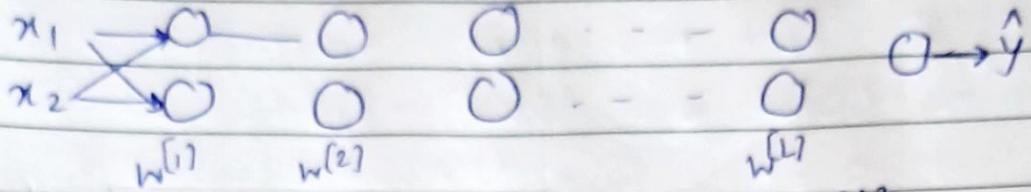


Elongated

Normalized



Vanishing / Exploding Gradients



If $g(z) = z$ linear and $b^{[L]} = 0$

- Derivatives or slopes can get either very big or very small.

$$\begin{aligned} a^{[1]} &= g(z^{[1]}) = \frac{d}{dz} = w^{[1]} x \\ a^{[2]} &= g(z^{[2]}) = \frac{d}{dz} = w^{[2]} a^{[1]} \\ &\quad = w^{[2]} w^{[1]} x \end{aligned}$$

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

Let's say that each of weight matrices is just $\begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$

$$\text{then } \hat{y} \approx \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}^{L-1} x \quad (1.5^{L-1})$$

\Rightarrow Grows exponentially. (Explode)

similarly $\begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}^{L-1}$

\Rightarrow Vanishes.

→ Weight Initialization ← to prevent vanish/explode

⇒ For single neuron:

$$Z = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b^0$$

• If n is large → we want small w_i to prevent exploding

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

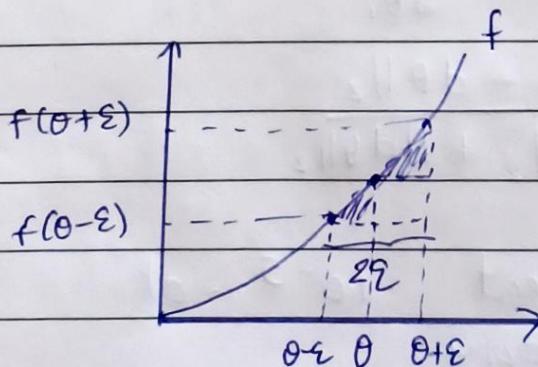
i.e. $w^{[l]} = np \cdot \text{random} \cdot \text{randn}(\dots) * np \cdot \text{sqrt}\left(\frac{1}{n^{[l-1]}}\right)$

If ReLU then it is $\left(\frac{2}{n}\right)$ & $\sqrt{\frac{2}{n^{[l-1]}}}$

If tanh then $\sqrt{\frac{1}{n^{[l-1]}}}$ (xavier initialisation)

Also used are $\sqrt{\frac{2}{n^{[l-1]} + n^{[l]}}}$

Numerical Approx. of Gradients (Gradient checking)



Instead of just moving right, go right then again left so u get a bigger Δ

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

Date :

MON TUE WED THU FRI SAT SUN

$$g(\theta) \propto \frac{f(\theta+\varepsilon) - f(\theta-\varepsilon)}{2\varepsilon}$$

eg.

$f(\theta) = \theta^3$	$\theta = 1 \rightarrow 3(1)^2 = 3$
$(1.01)^3 - (0.99)^3$	$= 3.0001 \Delta$
$2(0.01)$	$\Delta = 0.0001$

In normal 1 sided difference method:

$$\Delta = 0.03$$

Grad. checking:

Reshape $w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}$ into a big vector θ

Simly $dw^{[i]}, db^{[i]}, \dots$ into $d\theta$

For each i :

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

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

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

$\approx 10^{-7}$ - great 10^{-5} - ok
 10^{-3} - worry.

Date :
 MON TUE WED THU FRI SAT SUN

optimization Algo.

Mini-batch gradient descent

e.g. $m = 5 \text{ million}$.

Normally, you process entire dataset
 then take a small gradient descent step
 Then again process 5 million.

solution: Mini-batches.

$$X = [\underbrace{x^{(1)} \quad x^{(2)} \quad \dots \quad x^{(1000)}}_{X \{1\}} \mid \underbrace{x^{(1001)} \quad \dots \quad x^{(2000)} \quad \dots \quad x^{(m)}}_{X \{2\}}]$$

Batch Gradient Descent: Algo till now.

Process entire training set at same time

for $t = 1 \dots 5000$:

Forward prop:

$$z^{[1]} = w^{[1]} x^{[t]} + b^{[1]}$$

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

$$A^{[L]} = g^{[L]}(z^{[L]})$$

} rect
implen.

Date :
 MON TUE WED THU FRI SAT SUN

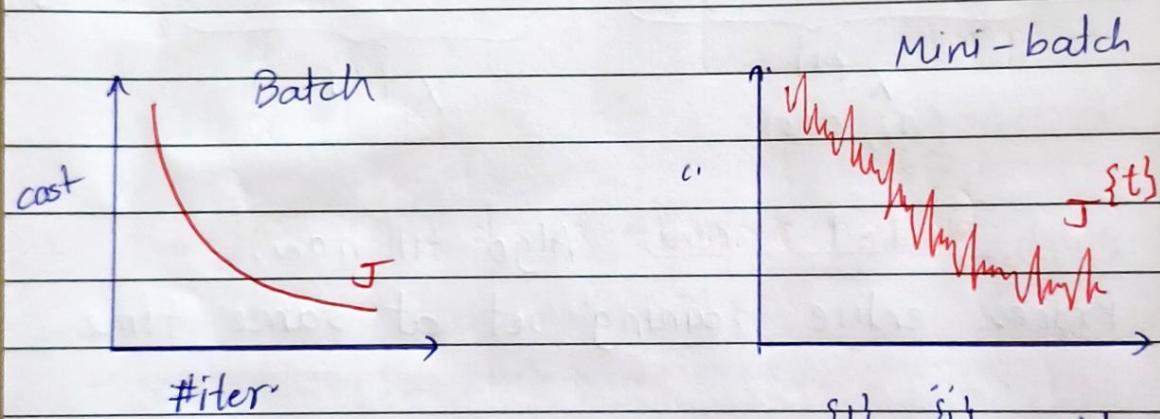
$$\text{Cost } J^{(t)} = \frac{1}{1000} \sum_{i=1}^{\frac{1}{2}} l(y^{(i)}, \hat{y}^{(i)}) + \lambda \sum_{j=0}^{2-10} \|w^{(j)}\|^2$$

Backprop:

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

$$b^{(l)} := b^{(l)} - \alpha d b^{(l)}$$

This is a single pass (or 1 epoch)
thru training set



$b \in x^{(1)}, y^{(1)}$ might
be easier minibatch but
 $x^{(2)}, y^{(2)}$ might be
harder.

→ Mini-batch size ←

minibatch size = m: Batch gradient descent

" " = 1: STOCHASTIC GRADIENT DESCENT

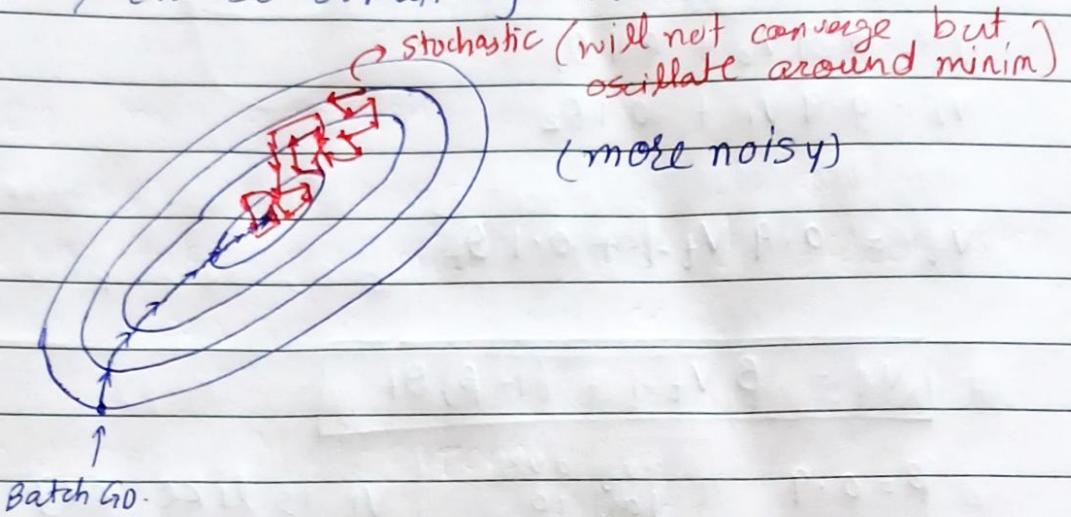
Every eg = own mini-batch.

Mini-batch size = HyperParam.

Date :

MON TUE WED THU FRI SAT SUN

So, choose something in betw.



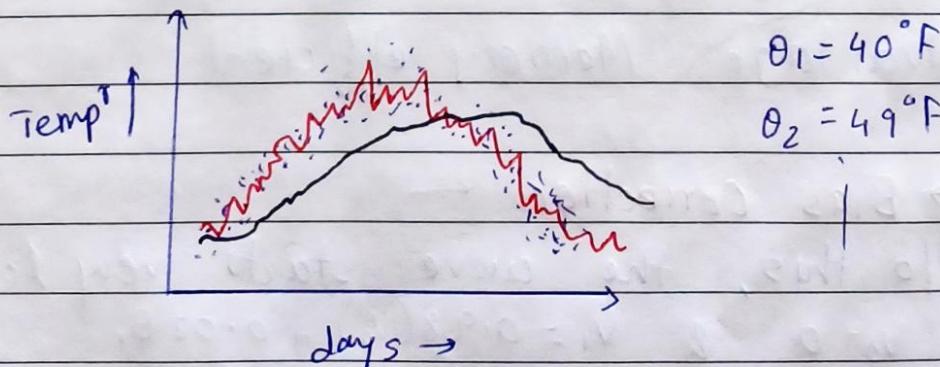
Small train set ($m \leq 2000$) \rightarrow Batch.

Typical size: 64, 128, 256, 512.

More Optimization Algo other than
Grad. Desc.

\rightarrow Exponentially Weighted Avg.

eg.



β also a hyperparam.

Date :.....

MON TUE WED THU FRI SAT SUN

$$V_0 = 0$$

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

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

|

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

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

$$\beta = 0.9 \quad (\text{Avg over last 10 days})$$

$$\beta = 0.98$$

lot of wt to prev. value

V_t is like avg over $\approx \frac{1}{1-\beta}$ days
temp.

$$\beta = 0.5 \approx 2 \text{ days (more noisy).}$$

$$\begin{aligned} \Rightarrow V_{100} &= 0.1 \theta_{100} + 0.9 V_{99} (0.1 \theta_{99} + 0.9 \theta_{98}) + \dots \\ &= 0.1 \theta_{100} + 0.1 \times 0.9 \theta_{99} + 0.1 (0.9)^2 \theta_{98} \\ &\quad + \dots \end{aligned}$$

Advantage: Memory efficient

→ Bias Correction ←

w/o this, the curve starts very low.

$$\therefore V_0 = 0 \quad \& \quad V_1 = 0.98(V_0) + \underbrace{0.02 \theta_1}_{\text{very low}}$$

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

Instead of $v_t \rightarrow \frac{v_t}{1-\beta_t}$ (t = current day)

$$\text{eg. } t=2 \quad 1-\beta^t = 1-(0.98)^2 = 0.0396$$

$$\frac{v_2}{1-\beta^t} = \frac{0.0196 \theta_1 + 0.02 \theta_2}{0.0396}$$

$$\uparrow \text{sum} = D^r$$

Gradient Descent with Momentum

↳ faster than just G.D.

↳ calc. exponentially weighted avg of gradients & then use it to update wts.

On iteration t :

compute dW, db on current minibatch/batch.

$$VdW = \beta VdW + (1-\beta) dW$$

$$Vdb = \beta Vdb + (1-\beta) db$$

$$w := w - \alpha VdW \quad b := b - \alpha Vdb$$

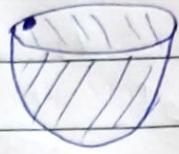


Most common $\beta = 0.9$

Date :

MON TUE WED THU FRI SAT SUN

Bowl analogy:

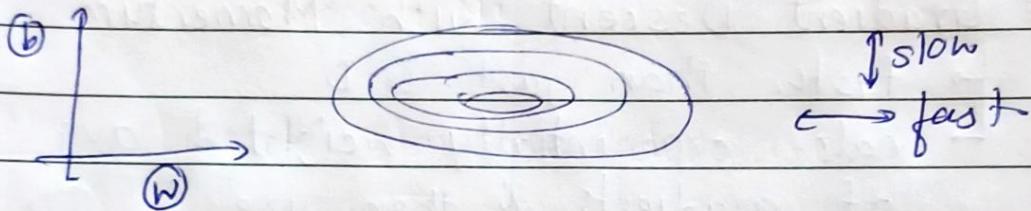


d_w, d_b are like acc given to ball

v_{dw}, v_{db} are velocities & β is friction.

Gen'ly we don't do bias corr here.

RMSprop [Root Mean square Prop]



On iter. t:

calc. d_w, d_b on current mini-batch / bat

$$\begin{aligned} Sd_w &= \beta Sd_w + (1-\beta) d_w^2 && \left. \right\} \text{elem-wise} \\ Sd_b &= \beta Sd_b + (1-\beta) d_b^2 && \left. \right\} \text{squaring.} \end{aligned}$$

$$w := w - \frac{\alpha}{\sqrt{Sd_w}}$$

$$b := b - \frac{\alpha}{\sqrt{Sd_b}}$$

Now, d_w will be small & $d_w = 1$ orge.

* damps out oscillations in Gradient desc.

Adaptive Moment Estimation.

Date :

MON	TUE	WED	THU	FRI	SAT	SUN
<input type="checkbox"/>						

Adam optimization Algo: (RMSProp Momentum)

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

on iter. t:

compute dw, db using current mini-batch

$$V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) dw,$$

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) dw^2$$

$$V_{db} = \beta_1 V_{db} + (1 - \beta_1) db,$$

$$S_{db} = \beta_2 S_{db} + (1 - \beta_2) db^2$$

$$\text{Bias } \left\{ \begin{array}{l} V_{dw}^{\text{corr}} = V_{dw} / (1 - \beta_1^t), \quad V_{db}^{\text{corr}} = V_{db} / (1 - \beta_1^t) \\ S_{dw}^{\text{corr}} = S_{dw} / (1 - \beta_2^t), \quad \dots \end{array} \right.$$

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

Hyperparam. choice:

$\alpha \rightarrow$ to be tuned

$$\beta_1 = 0.9 \quad \beta_2 = 0.999$$

$$\epsilon = 10^{-8} \text{ (keep anything doesn't matter)}$$

Learning Rate Decay:

slowly reduce α o/w it will keep oscillating around the minimum.

Date :

MON TUE WED THU FRI SAT SUN

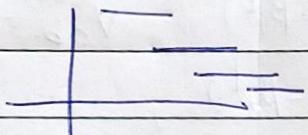
M-1: $\alpha = \frac{1}{1 + \text{decay} \cdot \text{epoch No.}} \times \alpha_0$

M-2: Exponential decay:

$$\alpha = 0.95^{\text{epoch-No.}} \cdot \alpha_0$$

M-3: $\alpha = \frac{\text{const}}{\sqrt{\text{epoch No.}}} \cdot \alpha_0$

M4: Discrete staircase decrease-



M-5: Manually decr. α .

Local optima:

In a very high-dimensional model, you are most likely to get stuck in a saddle point rather than local optimum

Problem = Plateau \rightarrow slows down learning very much
 ADAM & RMSprop help to solve this problem.

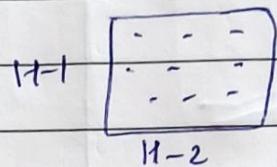
Date : 25/07/24
MON TUE WED THU FRI SAT SUN

Order of tuning of Hyperparameters

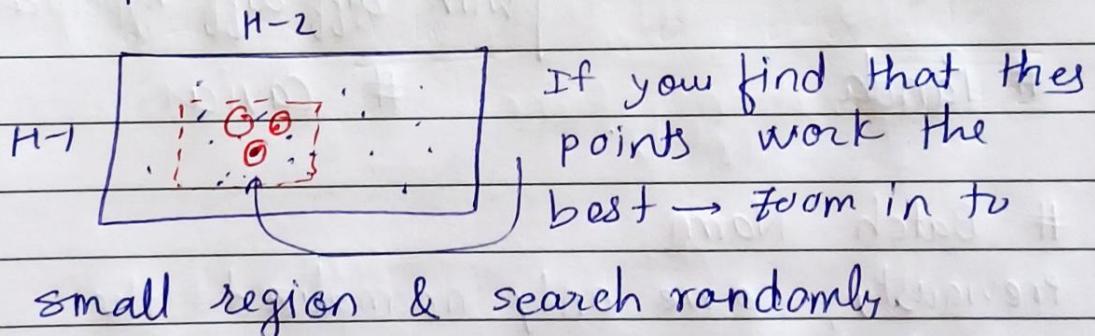
- 1) α , Most imp.
- 2) β , mini-batch size
- 3) $\#$ layers, learning rate decay.

Sampling Hyperparams.

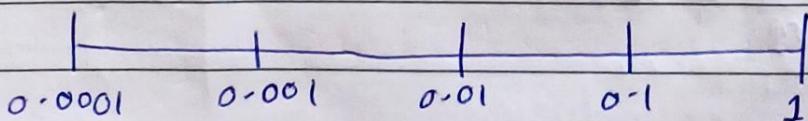
Try Random combinations - don't use grid like



Also, we can COARSE to fine the sampling scheme.



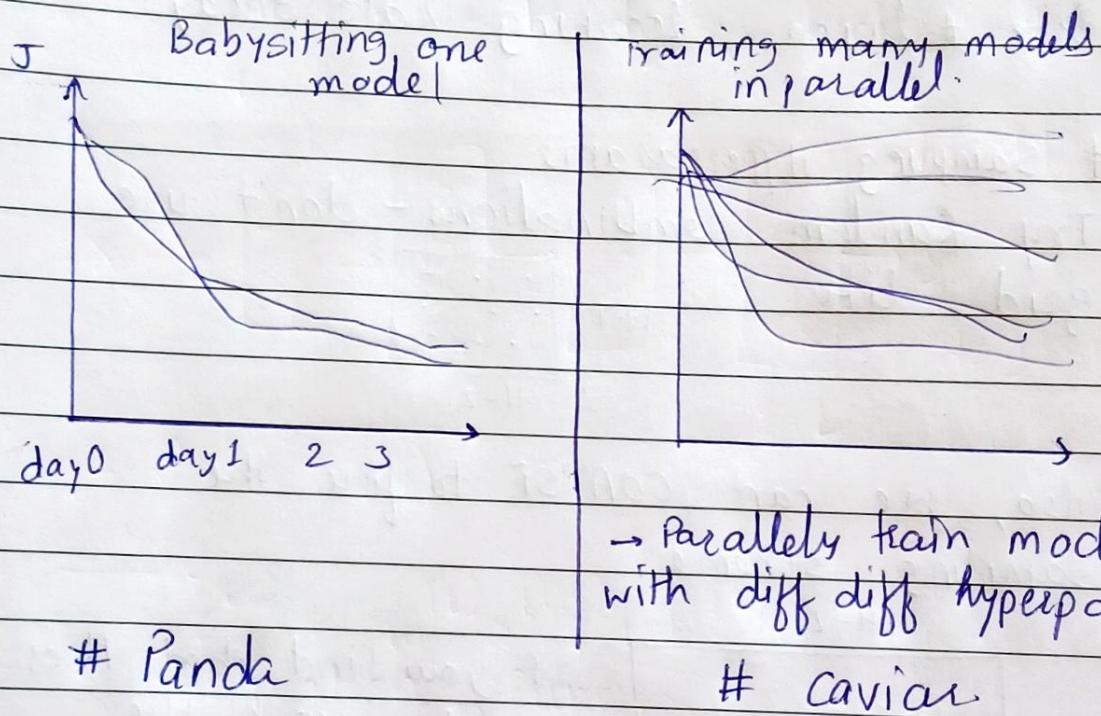
Instead of using linear scale to pick random, use log scale.



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

$\alpha = 10^r$

Retest hyperparams once every several months.



Batch Norm

Previously, we only normalized the IIP features. Here we apply this in some hidden layer as well.

If $z^{(1)}, z^{(2)}, \dots, z^{(m)}$:

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

Date:

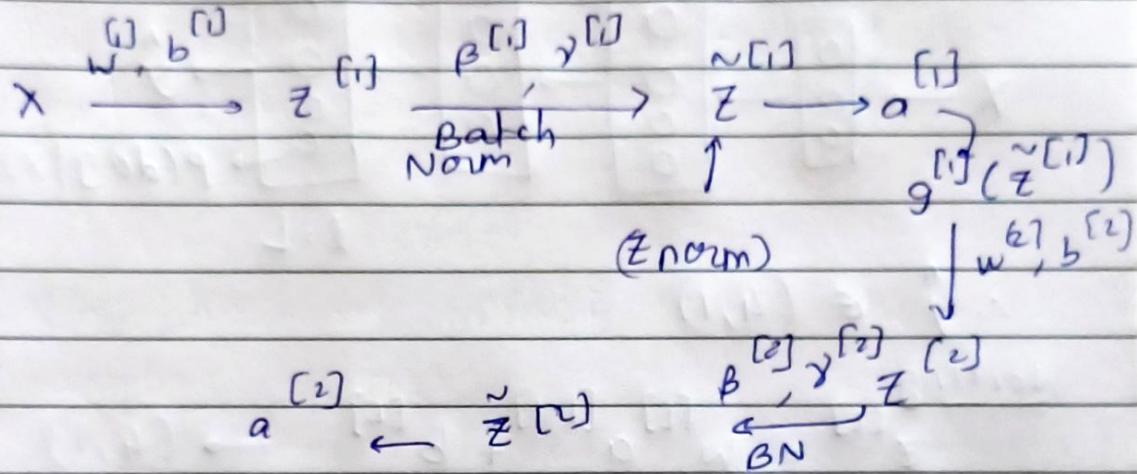
MON TUE WED THU FRI SAT SUN



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

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

$$\gamma = \sqrt{\sigma^2 + \epsilon} \quad \& \quad \beta = \mu \quad \left. \begin{array}{l} \text{Learnable} \\ \text{parameters} \end{array} \right.$$



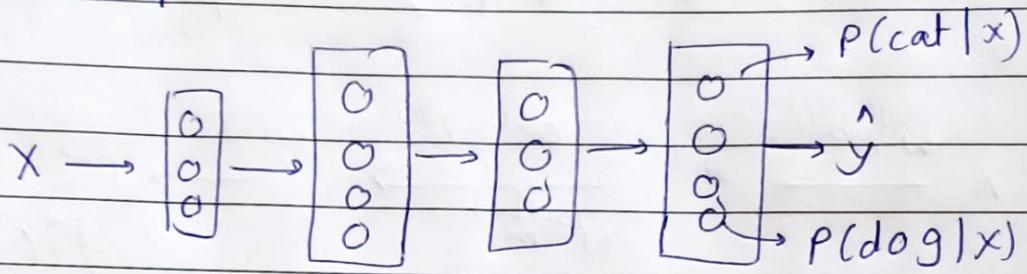
$$\beta^{[1]} = \beta^{[1]} - d\beta^{[1]} \quad \left. \begin{array}{l} \text{Upgrading } \beta \end{array} \right.$$

Softmax Regression .

- Generalizatⁿ of logistic reg.
- Predict to recognize one of c classes

eg. Recognizing cats, dogs, chicks, other

$$C = 4$$



$$\therefore \hat{y} \in (4, 1)$$

Working:

$$z^{[L]} = w^{[L]} a^{[L-1]} + b^{[L]} \rightarrow (4, 1)$$

Activation func:

$$t = e^{(z^{[L]})} \rightarrow (4, 1)$$

$$a^{[L]} = \frac{e^{z^{[L]}}}{\sum_{j=1}^4 t_i}, \quad a_i^{[L]} = \frac{t_i}{\sum t_i}$$

↙

(4,1)

eg. with zero hidden layers

with hidden layers

get a non-linear
decision boundary.



Date :

MON TUE WED THU FRI SAT SUN

[Hard max used to O/P as softmax is more gentle.]

If $C = 2$: softmax = logistic reg.

Loss func. in softmax:

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

$$\hat{y} = a^{[L]} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$$

Binary class.

$$L(\hat{y}, y) = -\sum_{j=1}^4 y_j \log \hat{y}_j$$

In this eq. $L(\hat{y}, y) = -y_2 \log \hat{y}_2 = -\log(\hat{y}_2)$
To make L small, we need \hat{y}_2 big.

vectorized impl:

$$y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}] = \begin{bmatrix} 0 & 1 & \dots \\ 0 & 0 & \dots \\ 0 & 0 & \dots \\ 0 & 0 & \dots \end{bmatrix}$$

$$\hat{y} = [\hat{y}^{(1)}, \hat{y}^{(2)}, \dots, \hat{y}^{(m)}] = \underbrace{\begin{bmatrix} 0.3 \\ 0.2 \\ 0.05 \\ 0.05 \end{bmatrix}}_{(4, m)}$$

In Tensorflow, we only need to implement forward prop & can use GradientTape. for easy backprop.

Date :

MON TUE WED THU FRI SAT SUN

<input type="checkbox"/>						
--------------------------	--------------------------	--------------------------	--------------------------	--------------------------	--------------------------	--------------------------

Grad. Desc. in softmax.

Backprop:

$$\begin{aligned} d\mathbf{z}^{[l]} &= \hat{\mathbf{y}} - \mathbf{y} \\ (4,1) & \\ \hookrightarrow \frac{\partial J}{\partial \mathbf{z}^{[l]}}. \end{aligned}$$

- * Usually prog. frameworks take care of backprop.