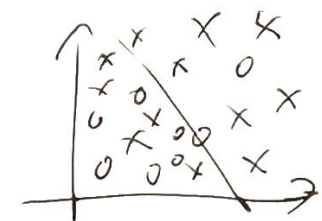
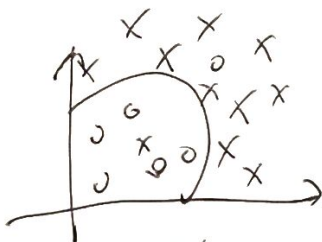


Hyperparameter tuning

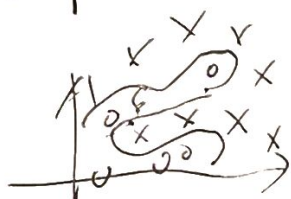
Recipe:



high bias.



just right



high variance.
— overfitting.

High bias? →
try data performance

↓ n

High variance →

Bigger network.

Train longer

(NN architecture search)

More data

Regularization

(NN arch, search)

Train set error 1%

Dev set error 11%

high variance

"Bias variance tradeoff"

↑

modern DNN can reduce both when bigger network & more data.

not doing very well on training set.
TE 15%
DE 16%
high bias.

Regularization:

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

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

single number
omit

$w \in \mathbb{R}^{n \times n}$

L_2 Regularization: $\|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2$

$b \in \mathbb{R}$

$$L_1 \text{ Regn: } \frac{\lambda}{2m} \sum_{j=1}^{n_x} |w_j| = \frac{\lambda}{2m} \|w\|_1 \rightarrow \text{norm}$$

1. Regn \rightarrow sparse W (a lot of 0 in)
easy to store.

2. Regn used much more often.

$\lambda \rightarrow$ Regularization parameter.
(lambda - python)

Neural network:

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

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

$$W \rightarrow \text{to } \mathbb{R}^{n^{[L]} \times n^{[L-1]}}$$

"Frobenius norm"

$dw =$ (from back prop)

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

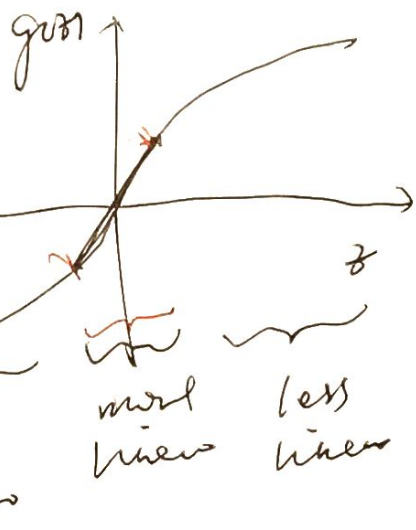
$$dw = (\dots) + \frac{\lambda}{m} w^{[l]}$$

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

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

weight decay

Reasons for using Regularisation



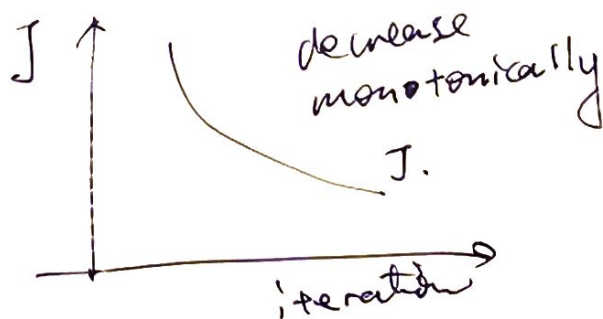
$\lambda \uparrow \rightarrow w^{[l]} \downarrow$
penalised.

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

Small \leftarrow Small

if z is small then cost could be roughly

linear \rightarrow avoid overfit



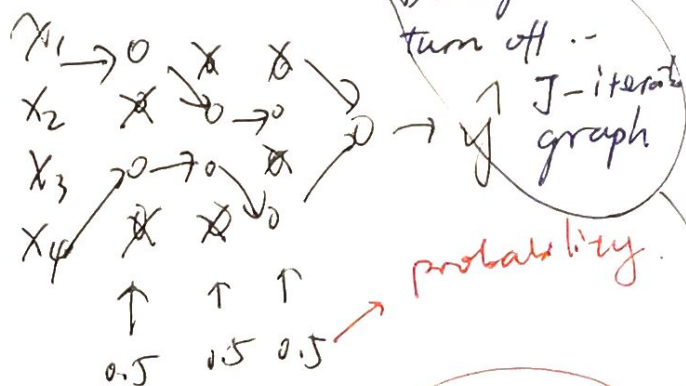
Dropout Regularization

Making predictions at test.

do not use Dropout

→ variety of results

→ add uncertainty / noise



Debug
turn off --
J-iterat
graph

Drop out node randomly.

Smaller network
→ less overfit.

↑
unstable.

any feature could go away

& current node using

dropout can be independent

eg: Inverted Dropout.

Expected value of $a^{[3]}$ remains

$d_3 = \text{np.random.randn}(a_3.\text{shape}[0], a_3.\text{shape}[1])$

→ layer 3. rand?

motivated by spread out the weights.

keep-prob could vary in different layer.

eg: Big layer could cause overfitting → smaller keep-p.

< keep-prob. has a 0.8 chance of being true & a 0.2 chance of being 0.

$a_3 = a_3 * d_3 \rightarrow$ boolean array.
or (np. multiply)

Some layer → 1.0
use very high value not at input layer

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

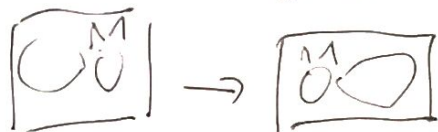
$$z^{[4]} = w^{[4]} a^{[3]} + b^{[4]}$$

↑ reduced by bump back $= 0.8$ 20%

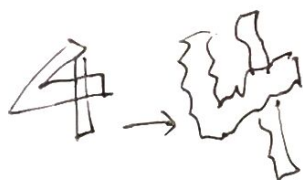
eg: Computer vision → overfitting
large pixel ← successful

Data augmentation

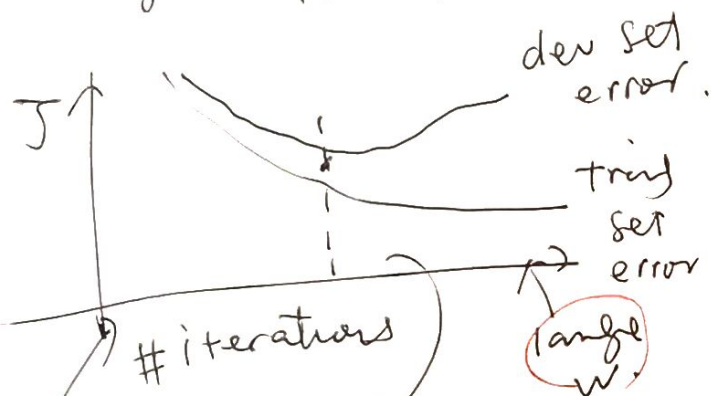
image \rightarrow flipping, distortion



inexpensive way for extra data.



Early Stopping



stop training at this point.

W gets bigger

overfitting less.

pros: run once & try out the different size of W .

Orthogonalization

Coupled $\left\{ \begin{array}{l} \text{Optimize Cost fun } J. \\ - \text{ Grad Des.; momentum} \end{array} \right.$

Regularisation $\left\{ \begin{array}{l} \text{Not overfit} \\ - \text{ Regu.} \end{array} \right.$

cons

J not optimised. mixed two problems.

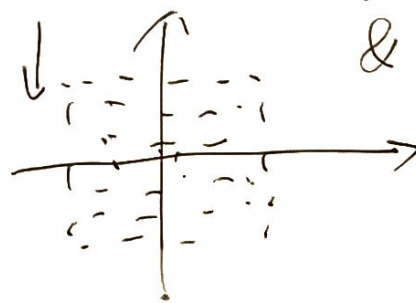
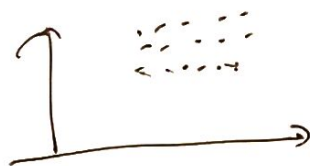
Still trying to avoid overfit

use L_2 Regu

And train as long as possible \rightarrow easier to decompose hyper parameters.

cons: try out \approx computationally expensive

Normalizing training data



centralised & augmented set.

zero out mean number.

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

$$X := X - \mu$$

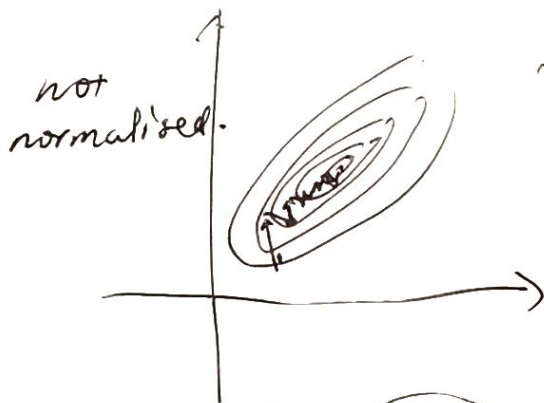
Normalize variance

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m X^{(i)} ** 2.$$

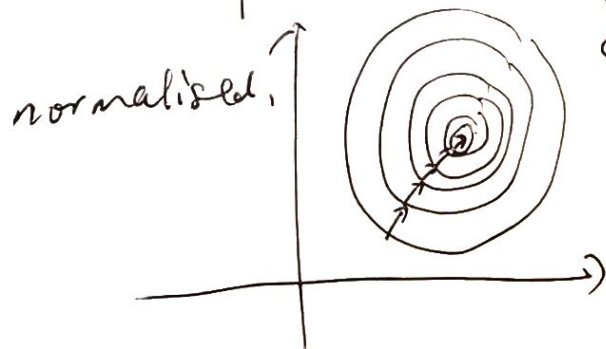
$$X := X / \sigma$$

Reasons for normalising

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



$x_1: 1 \dots 1000$
 $x_2: 0 \dots 1$
 range differ
 more symmetric



normalise test
 Set with same
 μ & σ^2 paras.

Vanishing gradient.
 exploding

For deep neural network.

$w^{[l]} > 1 \rightarrow$ exponentially
 $w^{[l]} < 1$ explodes.

DNN

$$\hat{y} = (w^{[l]})^{(n-1)} x$$

Partial solution:

Initialise, assume ReLU

$$\begin{matrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{matrix} > 0 \rightarrow \hat{y} \quad \frac{1}{n} \text{ for ReLU}$$

more $n \uparrow \rightarrow$ less weight
 $w_i \downarrow$ wanted.

$$\text{Var}(w_i) = \frac{2}{n}$$

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

$$\tanh \rightarrow \sqrt{\frac{1}{n^{[l-1]}}}$$

Xavier initialization

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

Reason?

tuning choices.

Gradient check for NN

Concatenate $W^{[1]} b^{[1]} \dots W^{[L]} b^{[L]}$

* $J(W, b, \dots) = J(\theta)$ into single vector θ .

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

* $J(dW^{[1]} db^{[1]} \dots dW^{[L]} db^{[L]}) = J(d\theta)$

$$J(\theta) = J(\theta_1, \theta_2, \dots, \theta_m)$$

For each i : (component of θ)

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

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

check relative distance

$$\frac{\|d\theta_{\text{app}} - d\theta\|_2}{\|d\theta_{\text{app}}\|_2 + \|d\theta\|_2} < e^{-7}$$

$\epsilon = 10^{-7}$ valid.

if $\begin{cases} 10e^{-5} \rightarrow \text{maybe components too large.} \\ 10e^{-3} \rightarrow \text{highly likely bug.} \end{cases}$

Practical implementation notes:

① Do not use in training — only to debug.

$$d\theta_{\text{approx}}[i]$$

→ slow computationally expensive.

only when debug.

⑤ Run at random initialization; again after some train $w, b \approx 0 \rightarrow w, b \uparrow$ debug. could error

② Turn to different components to identify bug

$$db^{[L]}, dW^{[L]}$$

$$③ J(\theta) = \frac{1}{m} \sum_i L(y^{(i)}, \hat{y}^{(i)})$$

Remember Regularization term.

④ Doesn't work with dropout. (turn off)

Mini-batch gradient descent

$V_{ec} \rightarrow m$ examples.

$$\begin{cases} X = [X^{(1)} X^{(2)} \dots X^{(m)}] \\ Y = [Y^{(1)} Y^{(2)} \dots Y^{(m)}] \end{cases}$$

(n_x, m)
 $(1, m)$

if $m = 50,000,000$

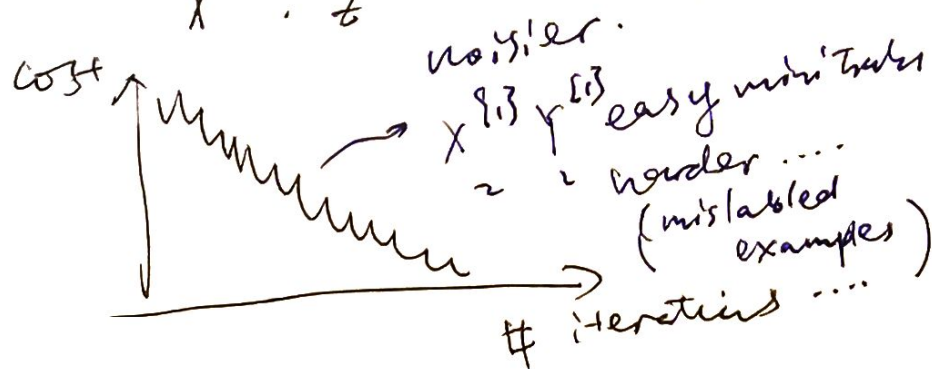
have to go through (process)
the entire training set
before the next step
— slow

mini batches of 1,000 each.

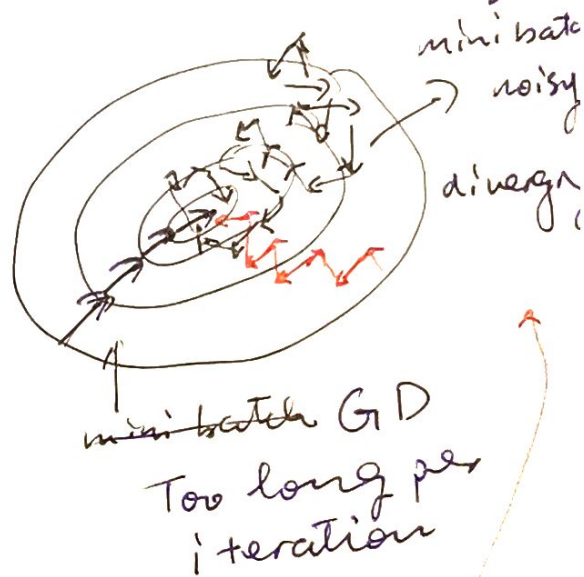
$$\begin{bmatrix} X^{(1)} & \dots & X^{(1000)} & \dots & X^{(m)} \\ \hline \underbrace{X^{(1)} \dots X^{(1000)}}_{X^{(1)}} & \dots & \underbrace{X^{(50000)}}_{X^{(50000)}} \end{bmatrix}$$

mini batch t : $X^{(t)}$, $Y^{(t)}$.
 $(n_x, 1000)$ $(1, 1000)$.

$$X^{(i)}, Z^{(i)}$$



mini-batch size = m :
batch gradient descent
mini-batch size = 1:
Stochastic GD.
loss full velocity from
vectorization



mini batches
trade off

$$m \leq 2000$$

→ BD

Typical mini-batches:

$$64, 128, 256, 512$$

$$2^6, 2^7, 2^8, 2^9$$

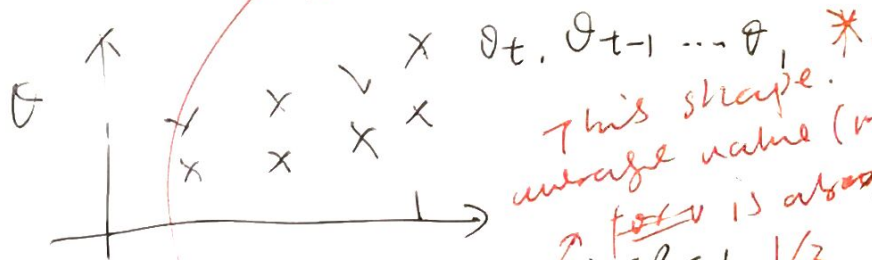
Exponentially Weighted Average (moving).

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

$$= (1-\beta) \theta_t + \beta ((1-\beta) \theta_{t-1} + \beta v_{t-2})$$

$$= (1-\beta) \theta_t + \beta (1-\beta) \theta_{t-1} + \beta^2 (1-\beta) \theta_{t-2}$$

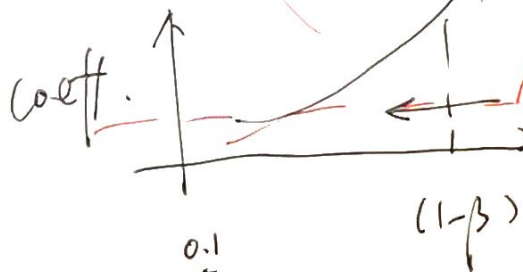
+ ... values about 10 days for value to decay to $1/e$.



This shape. average value (mean) for v is about $1/3$

$$0 < \beta < 1 \quad 1/3 \quad 1/e$$

exponentially decay



$$(1-\epsilon)^{1/\epsilon} = 1/e$$

$$\epsilon = 1-\beta \rightarrow \begin{cases} 0.9^{10} \approx 1/e \\ 0.98^{50} \approx 1/e \end{cases}$$

$$\approx \frac{1}{1-\beta} \text{ days average.}$$

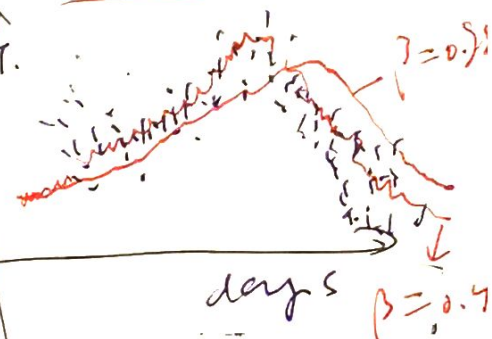
temperature is v_t . approximately

$\beta = 0.9 \Rightarrow v_t$ average over the last 10 days of temperature.

$\beta = 0.98 \rightarrow 50.$ $\beta = 0.5 \rightarrow$ averaging over 2 days.

$\beta \uparrow$ plot is much smoother averaging over more days of temp.

Curve shift to right. larger window of temperature Adapts slowly.



$\beta \uparrow \rightarrow (1-\beta) \downarrow$. adapt slowly to current θ_t . stronger simulation with former v_t . \rightarrow smoother

more noisy adapt much quicker

To implement.

Smooth out the steps of GD

just iteratively rewrite

$$V_0 = \beta V + (1-\beta) \sigma,$$

Just a ^{rough} coarse estimation

→ accurate mean requires storage of previous X days of temperature whereas this method only takes last value.

* { efficiency ✓
computationally cheap

Bias correction

large β . more damp to neutralise its oscillation

initialise with $V_0 = 0$.
which influence the performance
to produce averaging. ^{unavoidable}

During initial phase: ^{bad start}

* Do $\frac{V_t}{1-\beta^t}$

mini batch

$$t=2: 1-\beta^t = 1-(0.98)^2 = 0.396$$

$$\begin{cases} V_0 = 0 \\ V_1 = 0.98 V_0 + 0.02 \sigma \\ V_2 = 0.98 V_1 + 0.02 \sigma \end{cases}$$

$$\frac{V_2}{0.396} = \frac{0.196 \sigma_1 + 0.02 \sigma_2}{0.396} \approx \text{averages}$$

$t \uparrow \rightarrow$ no effect as $(1-\beta^t) \sim 1$



Gradient descent oscillates a lot & this prevents a large learning rate α as it may overshoot ^{diverging}.

momentum = (mv)

slower learning \uparrow

faster learning \rightarrow

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

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

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

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

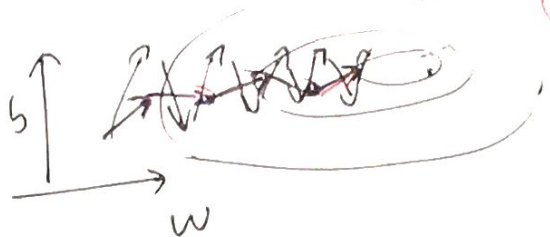
average $\uparrow \uparrow \uparrow \uparrow$ out
the oscillation

vertically slow down
horizontally - no big influence



$dW \rightarrow$ acceleration
 $W \rightarrow$ velocity
 $\beta \rightarrow$ friction.

RMSprop: Root Mean Square



want to speed up at w.
slow down at b.

On minibatches t :

Compute $dw, db \dots$

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

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

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

slope is very large on b direction $db \rightarrow$ large.
 $\rightarrow S_{db}$ large. $\rightarrow b$ updates vertically more slowly.

Damp out those noisy parameters with large slope.

in case the denominator is too small that leads to exponential explosion: add ϵ epsilon

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

combine momentum & RMSprop.

ADAM *

α = needs to be tune

$$\beta_1: 0.9 \quad (dw)$$

$$\beta_2: 0.999 \quad (dw^2)$$

$$\epsilon: 10^{-8} \rightarrow \text{each mini batch.}$$

For iteration t :

$$\begin{cases} V_{dw} = \beta_1 V_{dw} + (1-\beta_1)dw \\ V_{db} = \beta_1 V_{db} + (1-\beta_1)db \end{cases}$$

$$\begin{cases} S_{dw} = \beta_2 S_{dw} + (1-\beta_2)dw^2 \\ S_{db} = \beta_2 S_{db} + (1-\beta_2)db^2 \end{cases}$$

Combine.
correct bias correction

$$V_{dw}^{co} = V_{dw} / (1-\beta_1^t)$$

$$V_{db}^{co} = V_{db} / (1-\beta_1^t)$$

$$S_{dw} = S_{dw} / (1-\beta_2^t)$$

$$S_{db} = S_{db} / (1-\beta_2^t)$$

$$\begin{cases} w := w - \alpha \frac{V_{dw}}{\sqrt{S_{dw}} + \epsilon} \\ b := b - \alpha \frac{V_{db}}{\sqrt{S_{db}} + \epsilon} \end{cases}$$

Mini Batches won't converge.

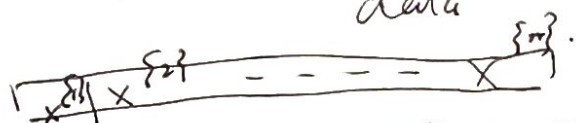


oscillates in a tight region.

oscillates around the optimised point & may not converge.

Learning rate decay.

1 epoch = 1 pass through data



epoch 1

epoch 2.

set $\alpha = \frac{1}{1 + \text{decayrate} * \text{epoch-number}}$ α_0

Epoch	α
1	0.1
2	0.67
3	0.5
4	0.4
...	...

2 for tune
gradually decrease.

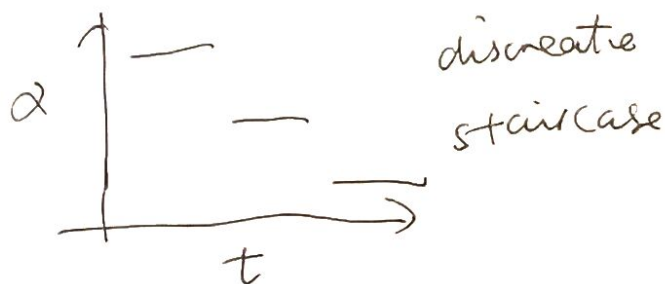
$\alpha = 0.2$

decay-rate = 1

$\alpha = 0.55 \cdot \alpha_0$

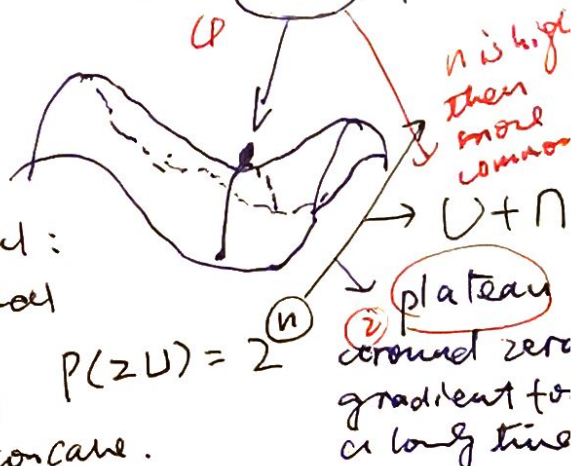
$\alpha = \frac{k}{\sqrt{ep}} \cdot \alpha_0$

or $\frac{k}{\sqrt{t}} \cdot \alpha_0$



OR manual decay.
→ for small number of models.
training

local optimal
advanced perceivition
most zero gradients
are saddle points.



③ local optimal:
n dimensional
space.

need. $U \cap$
convex concave.

$2U \rightarrow$ impossible nearby

$P(ZU) = 2$

② plateau
around zero
gradient for
a long time