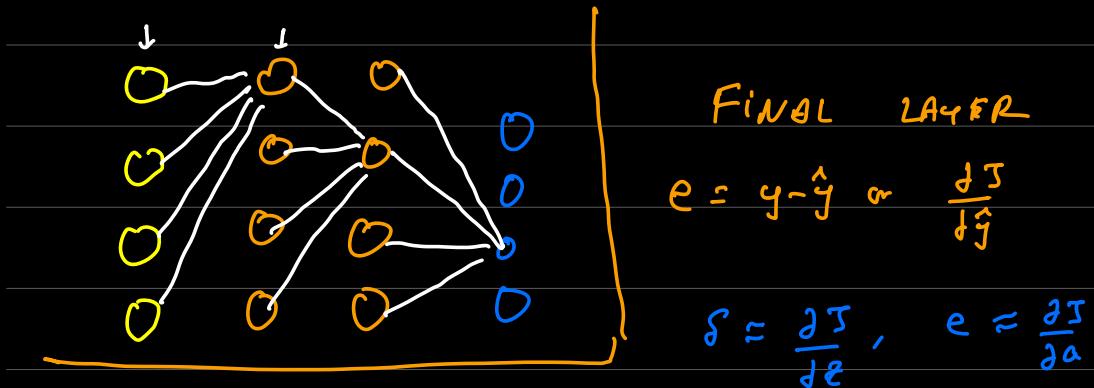


- 1) Summarize backprop algorithm for MLP
- 2) Variants of gradient descent

Σ

RECAP

Backprop expressions for an MLP



$$\delta^{(k)} = g'(z^{(k)}) \odot e^{(k)}$$

$$e^{(k+1)} = w^T \delta^{(k+1)}$$

L transpose of wd used in forward pass

FULL ALGORITHM

1) Initialize weights randomly → Another class!

2) Do forward propagation at 'forward pass'.

$$z^{(k+1)} = w^{(k)} a^{(k)}$$

$$a^{(k+1)} = g(z^{(k)}) \rightarrow \text{pointwise}$$

Calculate e - at final layer $\rightarrow \frac{\partial J}{\partial \theta} \therefore y - \hat{y}$

3) Backprop

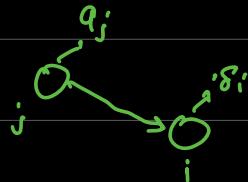
$$\delta^{(k)} = g'(z^{(m)}) \odot e^{(k)}$$

$$e^{(k)} = (w^{(k+1)})^T \delta^{(k+1)}$$

calculate this till the first hidden layer

4)

$$\Delta w_{ij} = +\delta_i q_j \frac{\partial}{\partial w_{ij}}$$



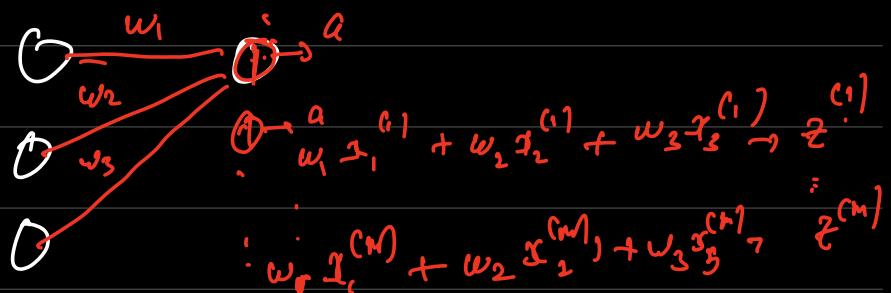
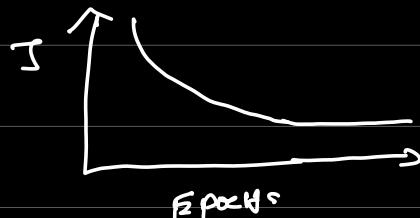
$$w_{ij} = w_{ij} + \Delta w_{ij} \leftarrow \begin{cases} \text{stochastic gradient} \\ \text{descent} \end{cases}$$

5)

Repeat 2-4 for all data pts (SGD)

6)

Repeat 2-5 till Convergence



Variants of gradient descent

$$\text{vector } \leftarrow w_{t+1} \leftarrow w_t - \alpha \nabla_w J^{(t)} \quad \nabla_w \\ \text{gradient wrt } w$$

t - time or iteration number

{ one data point at a time }

L_s SGD

1) what is α ? How to choose α

*

2) Should we change alpha with time?

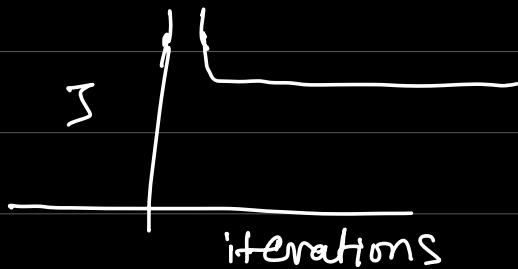
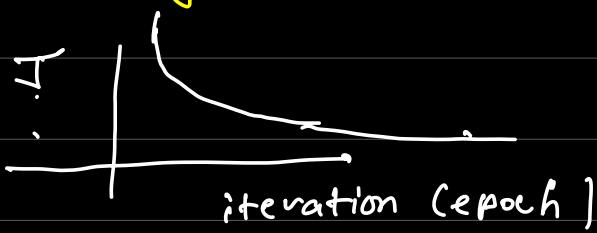
Meaning \rightarrow Can we change with each iteration

Scheduling techniques \rightarrow Preimposed

e.g. Reduce α by a factor of 10 every 100 epochs

3) why is ~~it~~ same for all parameters?

4) How do we get out of local minima, saddle points and regions of small gradients?



Solutions: Various adaptive approaches

Momentum, NAG: speeds up convergence }
even when gradient is low }

Adagrad, RMS prop: Adapt \rightarrow over time i.e.
iterations, over
parameters

Adam : Hybrid \rightarrow goto method ~~at~~

Vanilla GD

$$w \leftarrow w + \Delta w$$

$$\Delta w = -\alpha g$$

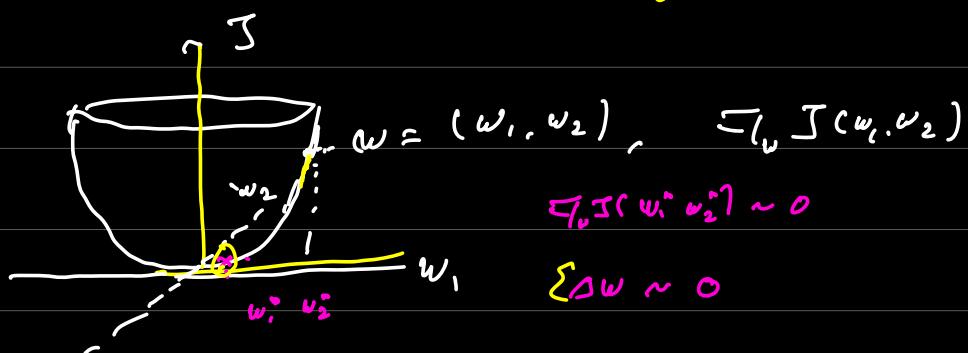
$$g = \nabla_w J$$

Momentum

$$w \leftarrow w + m$$

$$m = -\alpha g + \beta m_{t-1}$$

$$\begin{cases} m_t = -\alpha g_t + \beta m_{t-1} \\ w_t \leftarrow w_{t-1} + m_t \end{cases}$$



Momentum \leftarrow Earlier steps should also effect
where we will move in the current

step

$$\boxed{\begin{cases} m_t \leftarrow \frac{(\Delta w_t)}{\alpha} + \\ + m_t = -\alpha g_t + \beta m_{t-1} \\ w_t \leftarrow w_{t-1} + m_t \end{cases}}$$

Initially $m_0 = 0$

$\beta < 1$

$m_1 = -\alpha g_1 + \beta m_0 \rightarrow$ gradient descent
 $\hookrightarrow (\Delta w_1)$

$$M_2 = -\tilde{\Delta}g_2 + \beta M_1 \\ \Delta w_2 + \beta (-\Delta g_1)$$

$\Delta w_2 + \beta \Delta w_1 \rightarrow$ first gradient update
→ vector

$$M_3 = \Delta w_3 + \beta M_2$$

$$= \Delta w_3 + \underbrace{\beta \Delta w_2}_{\text{Exponential smoothing}} + \beta^2 \Delta w_1$$

$$\nabla J(w) = \begin{pmatrix} \frac{\partial J}{\partial w_1} \\ \frac{\partial J}{\partial w_2} \\ \vdots \\ \frac{\partial J}{\partial w_n} \end{pmatrix}$$

Exponential smoothing

$$\left[M_t = \Delta w_t + \beta \Delta w_{t-1} + \beta^2 \Delta w_{t-2} + \dots + \beta^{t-1} \Delta w_1 \right]$$

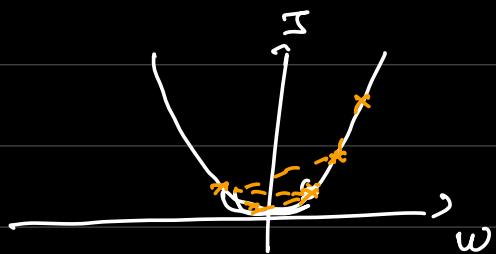
$\beta = 0$, gradient descent

$$\beta = 0.1 \rightarrow (10^{-1})^{10}, 10^{-10}$$

Nesterov Accelerated Gradient

(NAG)

Reduces the oscillations in momentum



Idea: Look-ahead

Use information about the update you are going to make

Momentum

$$w \leftarrow w + m$$

$$m \leftarrow \Delta w + (\beta m)$$

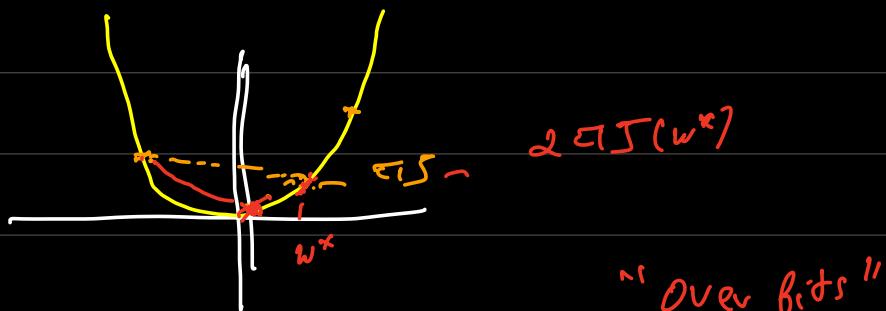
\swarrow \downarrow
Current Previous steps

NAG

$$\begin{aligned} w &\leftarrow w + m \\ m &\leftarrow \Delta w + (\beta m) \\ \Delta w^* &= -\nabla J(w^*) \\ w^* &= w + \beta m_{t-1} \end{aligned} \quad \left| \begin{array}{l} w \leftarrow w_{t-1} \\ w^* \leftarrow w + \beta m_{t-1} \xrightarrow{\text{previous momentum}} \\ \text{look ahead} \\ \Delta w^* = -\alpha \nabla J(w^*) \end{array} \right.$$

$$\left\{ \begin{array}{l} m_t = \Delta w^* + \beta m_{t-1} \\ w = w + m_t \end{array} \right.$$

Same as momentum



Adagrad

2- ideally depend on the parameter

→ take into account the gradient wrt every parameter

~ for a particular $w_i^{(t)}$ if it's gradient ↑
then its alpha-hd should ↓



$$w_{t+1} \leftarrow w_t - \frac{\alpha \nabla J_t}{\sqrt{V_i} + \epsilon}$$

Adaptive update → Component wise

$$V_i = \sum_{i=0}^t g_i^2, \quad g_i = (\nabla J)_i$$

↳ sum of all gradients

$$\sim \begin{pmatrix} \left(\frac{\partial J}{\partial w_1}\right)^2 \\ \left(\frac{\partial J}{\partial w_L}\right)^2 \\ \left(\frac{\partial J}{\partial w_3}\right)^2 \end{pmatrix} + \begin{pmatrix} \left(\frac{\partial J}{\partial w_1}\right)^2 \\ \left(\frac{\partial J}{\partial w_2}\right)^2 \\ \left(\frac{\partial J}{\partial w_3}\right)^2 \end{pmatrix} + \dots + \begin{pmatrix} \left(\frac{\partial J}{\partial w_1}\right)^2 \\ \left(\frac{\partial J}{\partial w_2}\right)^2 \\ \left(\frac{\partial J}{\partial w_3}\right)^2 \end{pmatrix}$$

$$\begin{pmatrix} \left(\frac{\partial J}{\partial w_1}\right)^2 \\ \vdots \\ \left(\frac{\partial J}{\partial w_n}\right)^2 \end{pmatrix}$$

$$w_{t+1} \leftarrow w_t - \frac{\alpha \nabla J_t}{\sqrt{V_i + \epsilon}}.$$

$$\rightarrow V_i = \sum_{i=0}^t G_i^2, \quad G_i = (\nabla J)_i$$

\nwarrow Sum of all gradients

$\mathcal{L} \rightarrow$ for each parameter

\rightarrow scales differently - based on its gradient history

$$w_0 \rightarrow [0, 0.1] \quad [0, 0.01]$$

$$V_1 \sim [0.5, 0.98]^{-1/10}$$

RHS prop: ADAGRAD :: NAG: Momentum

\uparrow smoothed version of Adagrad

\rightarrow ADAM \rightarrow ADAGRAD + MOMENTUM



Adapt α for
parameters



Time adaptation
of α

ADAM first

\rightarrow RMSprop - Adagrad

SGD

ADAM UPDATES

STANDARD

$$m_t = \beta_1 m_{t-1} - \alpha \nabla_w J(w_t) \rightarrow \text{momentum update}$$

$$v_t = \beta_2 v_{t-1} + (\nabla_w J(w_t))^2 \rightarrow \text{adaptive gradient}$$

$$w_{t+1} := w_t - \frac{\alpha m_t}{\sqrt{v_t} + \epsilon}$$

{ ADAGRAD + MOMENTUM }

INPUT DATA PRE-PROCESSING

1) NORMALIZATION

2) STANDARDIZATION \rightarrow Z-SCORE

3) PCA \rightarrow Discuss this briefly

1) NORMALISATION

no assumptions

INPUT - X

$$\hat{X} = \frac{X - X_{\min}}{X_{\max} - X_{\min}} \rightarrow [0, 1]$$

↳ what does this transformation accomplish?

X_{\min} and X_{\max} how are they estimated?

what data is used to estimate X_{\min} , X_{\max} ?

- ↳ 1) Use training data? \rightarrow Yes!
- 2) Use training for normalization during training
- 3) In test data - normalize using X_{\min} and X_{\max} estimated from test data
- 4) clip test data

Standardization

- 1) Assume that features are from a Gaussian distribution

- 2) Make the distribution 0-mean and 1-standard deviation

8) features are independent

$$\hat{x}_i = x_i - \bar{x}_i, \quad \bar{x}_i = \frac{1}{m} \sum_{j=1}^m x_i^{(j)}$$

↳ training set

$$\sigma_i^2 = \frac{1}{m} \sum_{j=1}^m (\hat{x}_i^{(j)})^2$$

$$z_i \leftarrow \left[\frac{x_i - \bar{x}_i}{\sigma_i} \right] \rightarrow Z\text{-score}$$

$$\sim N(0, 1)$$

(LBS) BW ↑ Ht (feet)	
x_1	\bar{x}_2
100	5.
200	6
150	7
300	4.5
	6.5

$\frac{|x_1|}{|\bar{x}_2|} \approx 20$

$y \sim \mathcal{BNI}$

$$\hat{y} = w_1 \hat{x}_1 + w_2 \hat{x}_2$$

w_1 and w_2 cannot

$$w_1 \leftarrow w_1 + \Delta w_1 \approx 10$$

Change the "same way"

$$w_2 \leftarrow w_2 + \Delta w_2 \approx 10$$

$$\hat{y} = w_1 \underbrace{100}_{\text{This term dominates}} + w_2 \underbrace{5}_{\text{can be ignored}}$$

$$w_1, w_2 \approx 0.5$$

Faster Convergence

$$\begin{aligned} \mathbf{x}_{\text{train}} &\rightarrow [-25, \dots, 125] \\ x_{\min} &= -25 & x_{\max} &= 125 \\ \mathbf{x}_{\text{test}} &\rightarrow \begin{bmatrix} -25 & -25 & -25 \\ -28, -27, -26, \dots, 120, 125 \end{bmatrix} \\ x_{\min} &\geq -45, & x_{\max} &\leq 150 \\ &\xrightarrow{\text{To normalize}} \end{aligned}$$

$$\frac{x^{(\text{test})} - x_{\min}^{(\text{test})}}{x_{\max}^{(\text{test})} - x_{\min}^{(\text{test})}} \leftarrow x$$

Standardization

Features \sim with different σ^2

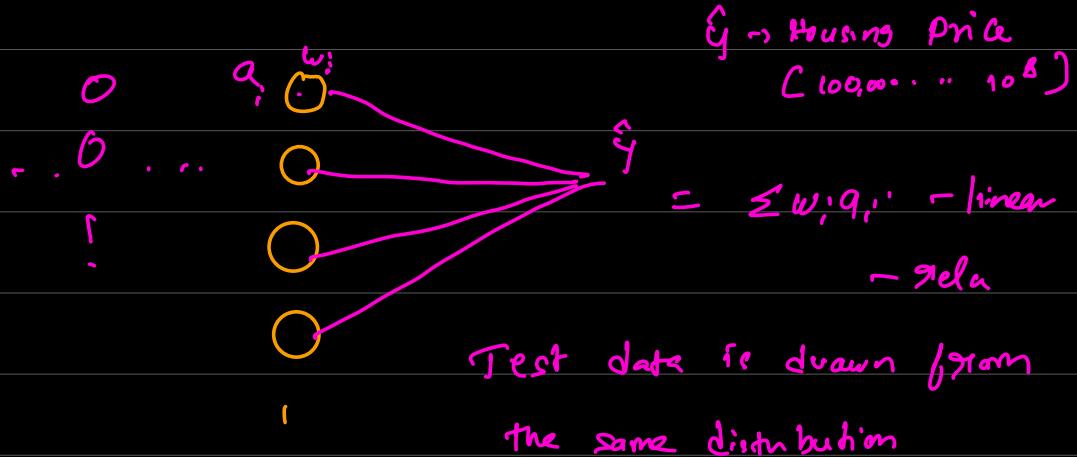
$$\begin{array}{ll} \sigma^2 = 100 \\ \sigma^2 = 10 \\ \vdots \end{array}$$

$X \rightarrow Y \sim$ Regression \rightarrow Det. change of Y
Classification \rightarrow one hot vector

If there are outliers in $Y(x)$ \rightarrow remove them from

Training

Final layer Fully connected neural network



Trained $y \in [10^4 \dots 10^6]$ - might fail
testing $y \in [10^6 \dots 10^9]$

problems in generalization

lack of generalization

Principal Component analysis

In many real world cases, features might not be independent

\Rightarrow Take an image \rightarrow pixels are correlated!

$X \rightarrow n$ features, m -dat points, you first

construct the covariance matrix

size of the covariance matrix $(n \times n)$

$$x = x_1 \dots x_n$$

$$\Sigma = \mathbf{U} \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix} \mathbf{U}^{-1}$$

zero out non diagonal elements

diagonal \rightarrow variance of individual features

Represents covariance between individual features

Diagonalization of the covariance matrix

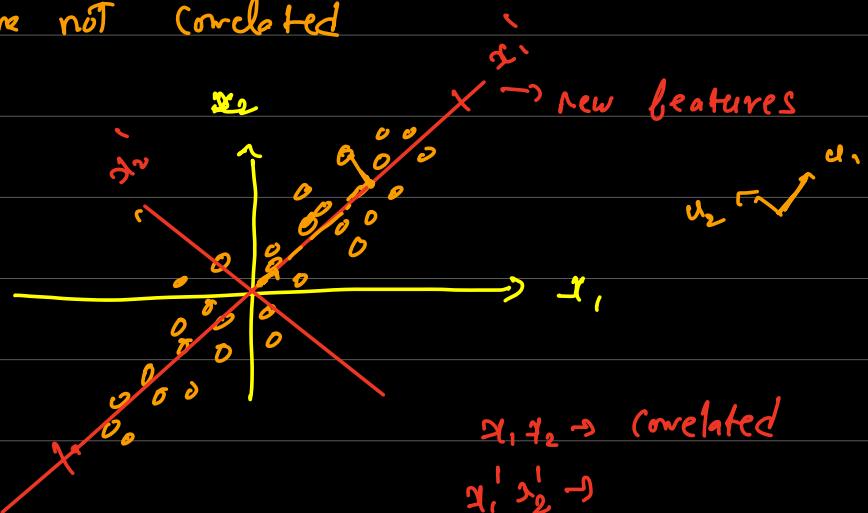
Solve the eigen value problem

$$\Sigma \mathbf{U} = \mathbf{U} \Lambda$$

\mathbf{U} eigen vectors Λ eigen values

\mathbf{U}^{-1} can be used to project the features so that

they are not correlated



You can just use \mathbf{x}' to represent your data

Varianc is captured by $\lambda \rightarrow n$

Throwing out features (transformed) corresponding
to smallest $\lambda \rightarrow$ Dimensionality reduction

$$\Sigma = \frac{1}{n} (x - \bar{x})^T (x - \bar{x}) \in \mathbb{R}^{m \times m}$$

$$x' = x \cup$$

Weight initialization }
Batch normalization } Dropouts -
 L1 and L2 regularization
 ↓
Convolutional optimization problems - $w=0$.
Neural network ← what happens?

Look at weight update equations

$w \sim N(0, I)$ → own set of problems

$\left\{ \begin{array}{l} w \sim N(0, \frac{2}{n} I) \sim \text{he initialization} \\ \frac{1}{n} \cdot I \quad n \text{ is the number of input} \\ \text{features in a layer} \end{array} \right\}$