

# ASSIGNMENT — 3

AKSHAY R  
A20442409

Review questions:

①. Neural networks.

(a).

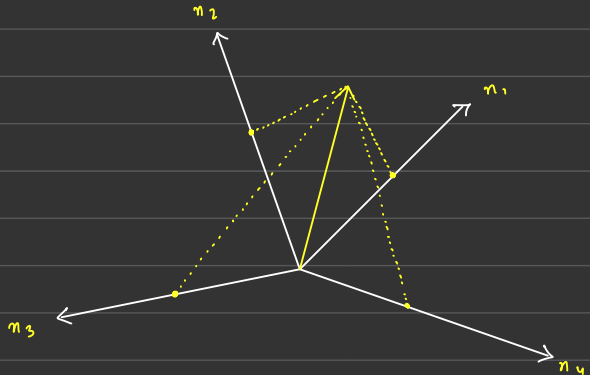
- \* Rows of  $\Theta^T$  (columns of  $\Theta$ ) are templates
- \* With  $k$  rows of  $\Theta^T$  we have  $k$  templates (one template per class)
- \*  $\Theta^T x$  measures how well  $x$  matches with each of the  $k$  templates (dot product of  $x$  with rows of  $\Theta^T$ )
- \* High similarity of a template of a particular class indicates high membership in this class.

$$\hat{y} = \Theta^T x + \theta_0$$

$$\hat{y} = \begin{bmatrix} n_1^T \\ \vdots \\ n_k^T \end{bmatrix} x + \begin{bmatrix} d_1 \\ \vdots \\ d_k \end{bmatrix}$$

$$\hat{y} = \begin{bmatrix} n_1 \cdot x \\ n_2 \cdot x \\ \vdots \\ n_k \cdot x \end{bmatrix} + \begin{bmatrix} d_1 \\ \vdots \\ d_k \end{bmatrix}$$

$n_1, n_2, n_3, n_4$   
are basis vectors.



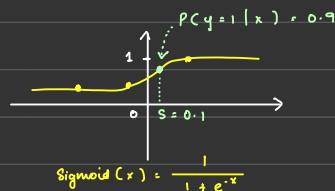
(b)

\* Given similarity score (or decision boundary distance  $s_j$ ) from a linear classifier (the higher the better), compute:

$$\underset{j^{\text{th}} \text{ unit}}{\overset{i^{\text{th}} \text{ example}}{y_j}}(i) \equiv \underset{\text{LABEL}}{p(y=j)} \mid \underset{\text{INPUT}}{x(i)}$$

\* In a 2 class classification case:

$$\hat{y}_j^{(i)} \equiv p(y=1 \mid x^{(i)}) = \text{sigmoid}(s_j^{(i)})$$



\* In a k class classification case:

$$\hat{y}_j^{(i)} \equiv p(y=j \mid x^{(i)}) = \frac{\exp(s_j^{(i)})}{\sum_{l=1}^k \exp(s_l^{(i)})} \leftarrow \text{softmax}$$

\* Used at the last layer

(c)  $L_1$  and  $L_2$  loss

$$L_1: L_i(\theta) = \sum_{j=1}^k |\hat{y}_j^{(i)} - y_j^{(i)}|$$

$$L_2: L_i(\theta) = \sum_{j=1}^k (\hat{y}_j^{(i)} - y_j^{(i)})^2$$

$\hat{y}_j^{(i)}$  is the predicted value :  $y_j^{(i)} \in \mathbb{R}$

Hubber loss:

$$\rho_{\delta}(d) : \begin{cases} \frac{1}{2} d^2 & \text{for normal errors.} \\ \delta(d - \frac{\delta}{2}) & \text{if } |d| \leq \delta \\ 0 & \text{otherwise} \end{cases}$$

$d \rightarrow$  distance b/w pred and actual.  
for outliers

$$L_i(\theta) = \sum_{j=1}^k \rho_{\delta}(\hat{y}_j^{(i)} - y_j^{(i)})$$

Cross entropy loss:

\* Convert similarity scores to probabilities:

$$\hat{y}_j^{(i)} = P(y = j | x^{(i)}) \quad j \in [1, k]$$

\* Likelihood:  $L(\theta) = \prod_{i=1}^m \prod_{j=1}^k (P(y = j | x^{(i)}))^{y_j^{(i)}}$

↑  
Maximize

\* Log-likelihood:  $-\log L(\theta) = - \sum_{i=1}^m \sum_{j=1}^k y_j^{(i)} \log(P(y = j | x^{(i)}))$

↑  
Minimize  
(negative LL)

$$= - \sum_{i=1}^m \sum_{j=1}^k y_j^{(i)} \log(\hat{y}_j^{(i)})$$

\* Sample loss:  $L_i(\theta) = - \sum_{j=1}^k y_j^{(i)} \log(\hat{y}_j^{(i)})$

(d) \* Occam's razor: simpler explanations are better.

\* A simpler explanation is when the weights  $\theta$  are lower (eg: when  $\theta_{ij} = 0$  we remove one coefficient)

\* Smaller coefficients  $\Rightarrow$  more stable solution that will generalize better.

$$* \quad L(\theta) = \frac{1}{m} \sum_{i=1}^m L_i(f(x_{ij}, w), y_i) + \lambda R(\theta)$$

ERROR TERM
Regularization term

↳ Weight of Regularization  
HYPER-PARAMETER

$$R(\theta) = \sum_i \sum_j \theta_{ij}^2 \quad \text{or} \quad \sum_i \sum_j |\theta_{ij}|$$

$L_2$ 
 $L_1$

\*  $L_2$  regularization minimizes weights while spreading term  
(eg:  $0.5^2 + 0.5^2 < 1^2 + 0^2$ )

\*  $L_1$  regularization doesn't have this property and may concentrate terms  
(eg:  $10.51 + 10.51 = 1 + 0$ )

(e) We move in the opposite direction of the gradient because we are trying to minimize the loss function and hence in a -ve direction with respect to the gradient.

(f) Stochastic gradient descent:

\* Randomly order examples.

\* for  $j: 1 \dots m$

↳ 1 EPOCH: Traverse through entire dataset.

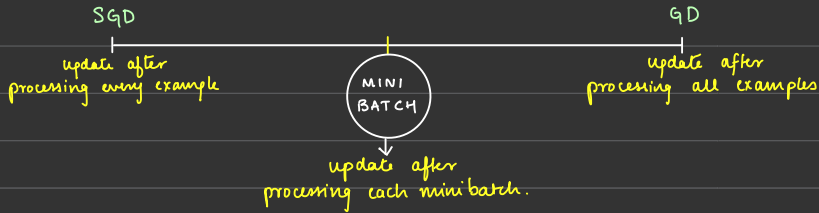
$$\theta^{(i+1)} \leftarrow \theta^{(i)} - \eta \nabla L_j(\theta^i)$$

↳ Based on example  $j$

\* Example for linear regression

$$\theta^{(i+1)} \leftarrow \theta^{(i)} - \eta (x^{(i)T} \theta - y^{(i)}) x^{(i)}$$

update after every example.



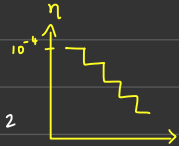
(g)

\* Learning rate need not be fixed!

→ make learning rate smaller as iterations progress

\* Strategies :

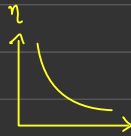
① Step decay : every  $k$  iterations  $\eta \leftarrow \eta/2$



② Exponential decay :

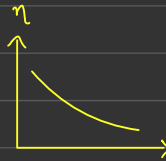
$$\eta_t = \eta \cdot e^{-\frac{k}{t}}$$

$\nwarrow$  Decay rate  
 $\nwarrow$  Iteration index



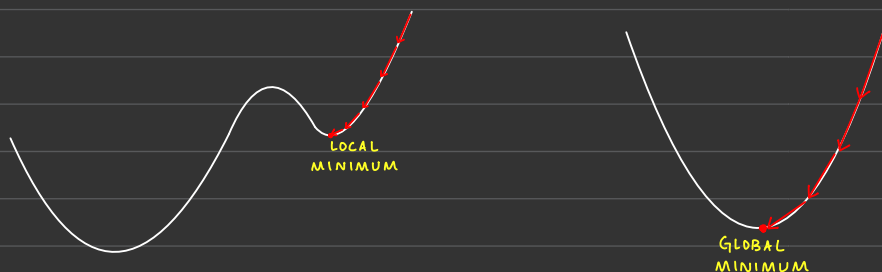
③ Fractional decay :

$$\eta = \eta_0 / (1 + kt)$$



→ Use when loss is not converging.

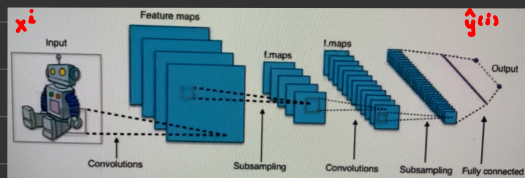
(h)



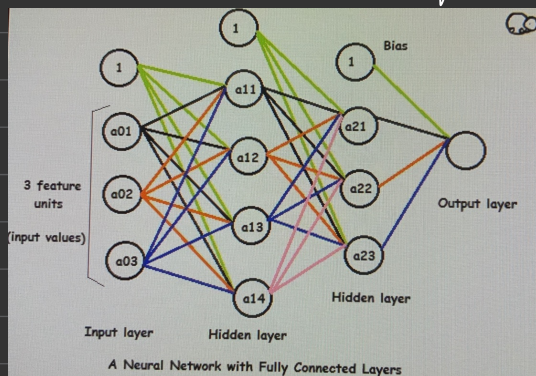
Not always will we have a loss function that will yield a global minimum while gradient descent. Sometimes we can get stuck in a local minima during gradient descent. Momentum is used to overcome the local minima problem.

(j) As the name states a fully connected layer is one where every neuron is connected to every other neurons in the preceding or the previous layer. Weights are calculated through forward and backward passes.

A CNN layer is one where weights are calculated by convolving a filter over an image and then using gradient descent.

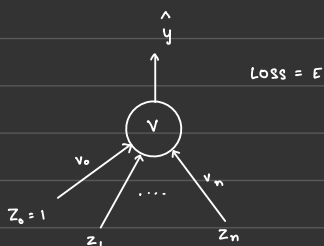


Convolution layers



Fully connected layers.

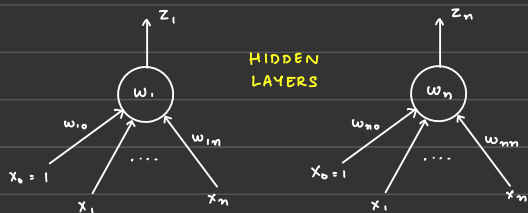
(i)



chain rule:

$$\frac{\partial E}{\partial v} = \frac{\partial E}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial v} \rightarrow \text{find } v$$

Loss  $\nearrow$   
parameters  $\nwarrow$



HIDDEN LAYERS

$$v^{(i+1)} \leftarrow v^{(i)} - \eta \nabla \frac{\partial E}{\partial v}$$

$$\frac{\partial E}{\partial w_j} = \frac{\partial E}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_j} \frac{\partial z_j}{\partial w_j} \rightarrow \text{find } w_j$$

- \* Start with a guess of parameters :  $v, \{w_j\}$  (at random close to zero).
- \* use  $x^{(i)}$  and current parameters to compute outputs  $z^{(i)}, \hat{y}^{(i)}$  : forward pass
- \* use outputs  $z^{(i)}, \hat{y}^{(i)}$  to update parameters :  $v, \{w_j\}$  : backward pass (push back gradients)
- \* Continue while loss changes.
- \* Represent the network using a computational graph.
- \* Because each node is simple it has a simple explicit expression for its derivatives.
- \* Forward pass: Push input to compute all intermediate node values
- \* Backward pass: Starting with end nodes push gradients towards the beginning.
- \* Multiply backpropagated gradients (from back) by current gradient and propagate this

Gradients are propagated in the backward pass.

(k) Regularization measures:

\* Large number of parameters tend to overfit

Methods to prevent overfitting:

① Dropout:

- At each training stage drop our units in **fully connected** layers with probability of  $(1-p)$  ← HYPER-PARAMETER
  - Removed nodes are reinitiated with original weights in the subsequent stage.
-



