Tuan Tran

A20357808

CS 584 - Summer 2020.

1,

- ~~Gene~~ In generative learning, we model the feature distribution and class priors —> allow generating new examples and also do classification.

- In discriminative learning, we model $P(Y|X)$ —> strictly doing prediction.

- Given nD feature vector:
$$g(x) = \theta^T x = \theta_0 + \theta_1 x_1 + \ldots + \theta_n x_n.$$

- Parameters $\theta_1 \ldots \theta_n$ can be thought of as the components of the normal vector to the decision boundary

  $\theta_0$ can be thought of as the negative distance from decision boundary to origin.

2, We have 2 options
- One vs all, ~~we~~ where we model the decision ~~boundary~~ boundaries that separates each class against ~~eve~~ all remaining classes
  └> This will produce K discriminant functions (one for each class)

- One vs each other where we separates each class against each of the remaining classes —> this will produce $\frac{K(K-1)}{2}$ discriminant functions, (one for each pair of classes)

→ One vs each other produces more discriminant Functions and thus more parameters
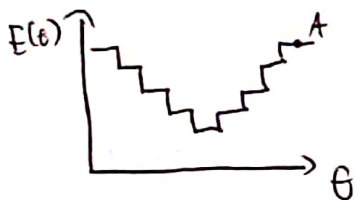
3,

- Empirical Error:

$$E(\theta) = \# x^\ast + \# x^{\ast\ast}$$

where $x^\ast = \{x^{(i)} \mid \mathbb{1}(y^{(i)}=0) \wedge \theta^T x > 0\}$

$x^{\ast\ast} = \{x^{(i)} \mid \mathbb{1}(y^{(i)}=1) \wedge \theta^T x < 0\}$.

↳ The problem is that $E(\theta)$ is a piecewise Function:



↳ we can't compute the gradient since we can easily stop at point A (gradient $=0$).

→ Can't optimize using GD.

4,

PC: $E(\theta) = \sum\limits_{x^{(i)} \in x^\ast} \theta^T x^{(i)} - \sum\limits_{x^{(i)} \in x^{\ast\ast}} \theta^T x^{(i)}$

⌣) ✗ $\nabla E(\theta) = \sum\limits_{x^{(i)} \in x^\ast} x^{(i)} - \sum\limits_{x^{(i)} \in x^{\ast\ast}} x^{(i)}$

↳ Update:

$$\theta \leftarrow \theta - \eta\left(\sum\limits_{x^{(i)} \in x^\ast} x^{(i)} - \sum\limits_{x^{(i)} \in x^{\ast\ast}} x^{(i)}\right)$$

**5,**

- In order to arrive at the logistic function as the hypothesis, we assume that the log of probabilities ratio can be modelled as a linear function:

$$\log \frac{P(y=1|x)}{P(y=0|x)} = \theta^T x .$$

- We simply equate $P(y=1|x) = \text{Sigmoid} = \dfrac{1}{1+e^{-\theta^T x}}$

➤ the result is in the range of $(0, 1)$ where $0.5$ represents Complete uncertainty.

**6,**

- $\ell(\theta) = \log\left(\overline{\prod_{x^{(i)} \in C_1} P(y=1|x^{(i)})}\; \overline{\prod_{x^{(i)} \in C_0} P(y=0|x^{(i)})}\right)$

$$= \log \prod_{i=1}^{m} P(y=1|x^{(i)})^{y^{(i)}} \; P(y=0|x^{(i)})^{y^{(i)}}$$

$$= \sum_{i=1}^{m} y^{(i)} \log P(y=1|x^{(i)}) + (1-y^{(i)}) \log P(y=0|x^{(i)})$$

$$\underbrace{\qquad\qquad}_{\log \overset{||}{P}(1-P(y=1|x^{(i)}))}$$

where $h_\theta(x) = P(y=1|x)$

➤ $\ell(\theta) = \displaystyle\sum_{i=1}^{m} y^{(i)} \log h_\theta(x^{(i)}) + (1-y^{(i)}) \log(1-h_\theta(x))$

- $\dfrac{d}{d\theta}\ell(\theta) = \dfrac{d}{d\theta} \displaystyle\sum_{i=1}^{m} y^{(i)} \log h_\theta(x^{(i)}) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)}))$

$$= \sum_{i=1}^{m} y^{(i)}(1-h_\theta(x^{(i)}))x^{(i)} + (1-y^{(i)})(-h_\theta(x^{(i)}))x^{(i)}$$

$$= \sum_{i=1}^{m} x^{(i)}\left(y^{(i)} - y^{(i)} h_\theta(x^{(i)}) - h_\theta(x^{(i)}) + y^{(i)} h_\theta(x^{(i)})\right)$$

$$= \sum_{i=1}^{m} (y^{(i)} - h_\theta(x^{(i)})) x^{(i)}$$

$\rightarrow$ ~~negative log like~~ gradient of negative log likelihood:

$$\frac{d}{d\theta}(-\ell(\theta)) = \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)}$$

- ~~&~~ SGD update ~~s~~ *: For batch $X_k$:

$$\theta \leftarrow \theta - \eta \underset{\cancel{x \in X_k}}{\overset{}{\sum}} \underset{i \in X_k}{\sum} (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)}$$

- Should initialize parameters randomly ~~&~~ ~~closely~~ close to 0

  Since as $\theta \rightarrow 0$, $h_\theta(x) \rightarrow 0.5$ $\rightarrow$ represents ~~great~~ greatest uncertainty

7,

- For $k$ classes, we use the softmax function since it
  outputs a distribution over all the ~~cla~~ classes and adds up ~~s~~ to 1.
- Softmax is essentially a "one vs all" version of ~~logis~~ sigmoid.

- Update equation:

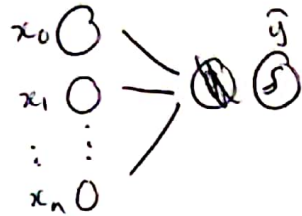$$\theta_j \leftarrow \theta_j - \eta (h_{\theta_j}(x) - 1(y=j)) x^{(i)}$$

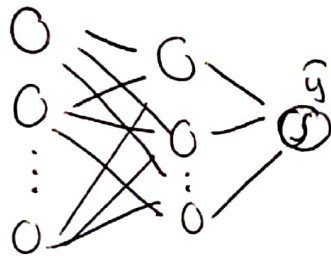  where $j = 1 \dots k$ and $\theta_j$ = parameter vector for class $j$

# Ⓐ Neural Net

※ 1,

- Basic Structure:

  A logistic unit:

  

  $x_0$, $x_1$, $\vdots$, $x_n$ → Ⓝ Ⓢ → $\hat{y}$

  → We can stack mutiple logistic unit on top of each other to Form a neural net.

  Ex:

  

- Activation function is important as it introduces non-linearity into the neural net model → give it ability to learn complex patterns.

- Several types of activation function: sigmoid, softmax, Relu, leaky Relu, linear (no activation), etc

2,

- ~~In~~ ※ Feedforward net pushes the input ~~allth~~ forward all the way to output, Feedback net also does the forward push, but the out put ~~loops~~ loops back to the input layer

- In the case of single output, we have exactly 1 parameter vector $\theta$ to plug into activation function

  In K-class case, we have $\underline{\underline{K}}$ parameter vectors $\theta_1 \dots \theta_k$

  In terms of classification, we use sigmoid for single output and softmax for K-class outputs.

3,

- number of hidden units < number of imp inputs :

  -> we are basically doing dimensionality reduction
  Since we are condensing a larger number of features
  into a smaller dimensional space
    -> Extracting the more important features

- number of hidden units > no of inputs :

  -> We are doing non-linear mapping of features to higher
  dimensional space
  -> Allows for more complex models .

4,

- Difficulty is that the update equations require knowing
  if $z^{(i)} = j$ or not where $z_j^{(i)}$ is the output of the unit $j$
  in hidden layer .

- Chain rule is used to derive update equations by moving backward
  one layer at a time, calculating gradient at each step. Since
  Since neural net is essentially a chain of functions, consider
  parameter $v$ into one output unit $\hat{g}$ with error $E$
  -> $\dfrac{\partial E}{\partial v} = \dfrac{\partial E}{\partial \hat{g}} \cdot \dfrac{\partial \hat{g}}{\partial v}$ .

5, Assume 2-layer

- Single output regression:

$$V = \begin{bmatrix} v_0 \\ \vdots \\ v_n \end{bmatrix} = \text{parameter vector of output layer}$$

$$W_j = \begin{bmatrix} w_{j0} \\ \vdots \\ w_{jn} \end{bmatrix} = \text{parameter vector of hidden unit } z_j.$$

$$E = \text{loss} = \frac{1}{2} \sum_{i=1}^{m} (\hat{y}^{(i)} - y^{(i)})^2$$

$$\rightarrow \frac{\partial E}{\partial v} = \frac{\partial E}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial v} = \frac{1}{2} \cdot 2 \sum_{i=1}^{m} (\hat{y}^{(i)} - y^{(i)}) z^{(i)}$$

$$= \sum_{i=1}^{m} (\hat{y}^{(i)} - y^{(i)}) z^{(i)}$$

+, Update for $V$:

$$V \leftarrow V - \eta \sum_{i=1}^{m} (\hat{y}^{(i)} - wy^{(i)}) z^{(i)}$$

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

$$= \sum_{i=1}^{m} (\hat{y}^{(i)} - y^{(i)}) v_j z_j^{(i)} (1 - z_j^{(i)}) x^{(i)}$$

+, Update for $w_j$:

$$w_j \leftarrow w_j - \eta \sum_{i=1}^{m} (\hat{y}^{(i)} - y^{(i)}) v_j z_j^{(i)} (1 - z_j^{(i)}) x^{(i)}$$

- Mutiple outputs: & assume $K$ outputs

→ $v_j$ = parameter vector of $\hat{q}_j$, $j = 1...k$.

$w_j$ = same definition as before.

$$E = \text{loss} = \frac{1}{2} \sum_{i=1}^{m} \sum_{\ell=1}^{K} \left( \hat{q}_\ell^{(i)} - y_\ell^{(i)} \right)^2$$

$\hookrightarrow$ sum over all $K$ classes.

+, $\dfrac{\partial E}{\partial v_j} = \dfrac{\partial E}{\partial \hat{q}_j} \cdot \dfrac{\partial \hat{q}_j}{\partial v_j}$

$$= \sum_{i=1}^{m} \left( \hat{q}_j^{(i)} - y_j^{(i)} \right) z^{(i)}$$

→ Update for $v_j$ where $j = 1...k$:

$$v_j \leftarrow v_j - \eta \sum_{i=1}^{m} \left( \hat{q}_j^{(i)} - y_j^{(i)} \right) z^{(i)}$$

+, $\dfrac{\partial E}{\partial w_j} = \sum_{\ell=1}^{K} \dfrac{\partial E}{\partial \hat{q}_\ell} \cdot \dfrac{\partial \hat{q}_\ell}{\partial z_j} \cdot \dfrac{\partial z_j}{\partial w_j}$

$\hookrightarrow$ summation because $w_j$ leads to $z_j$ which is input to all of ~~$k$ classes~~

Output units.

$$= \sum_{i=1}^{m} \sum_{\ell=1}^{K} \left( \hat{q}_\ell^{(i)} - y_\ell^{(i)} \right) v_{\ell j} z_j^{(i)} \left( 1 - z_j^{(i)} \right) x^{(i)}$$

→ update for $w_{ji}$

$$w_j \leftarrow v_j - \eta \sum_{i=1}^{m} \sum_{\ell=1}^{K} \left( \hat{q}_\ell^{(i)} - y_\ell^{(i)} \right) v_{\ell j} z_j^{(i)} \left( 1 - z_j^{(i)} \right) x^{(i)}$$

- In case of regression, objective function is MSE, for multiclass, it's MSE over all classes.

7, Assume 2-layer:

- Single output Classification:

  Using the sigmoid derivative rules

  and $\dfrac{\partial E}{\partial v} = \dfrac{\partial E}{\partial \hat{q}} \cdot \dfrac{\partial \hat{q}}{\partial v}$

  $\dfrac{\partial E}{\partial w_j} = \dfrac{\partial E}{\partial \hat{q}} \cdot \dfrac{\partial \hat{q}}{\partial z_j} \cdot \dfrac{\partial z_j}{\partial w_j}$

  where $E = \displaystyle\sum_{i=1}^{m} y^{(i)} \log \hat{q}^{(i)} + (1-y^{(i)}) \log(1-\hat{q}^{(i)})$

  $\llcorner$ we observe that $\dfrac{\partial E}{\partial v}$ and $\dfrac{\partial E}{\partial w_j}$ has same value as their regression

  Counterpart.

  $\rightarrow \quad \therefore v \leftarrow v - \eta \displaystyle\sum_{i=1}^{m} (\hat{q}^{(i)} - y^{(i)}) z^{(i)}$

  $w_j \leftarrow w_j - \eta \displaystyle\sum_{i=1}^{m} (\hat{q}^{(i)} - y^{(i)}) v_j z_j^{(i)} (1 - z_j^{(i)}) x^{(i)}$

- ~~class~~ K-class ~~&~~ Classification:

  Similarly, we also ~~obt~~ obtain same gradient and update equations

  as the regression ~~counter~~ Counterpart in case of multiple

  Outputs.

- Objective function used is the cross-entropy / negative log likelihood

  Obtained from negating the ~~likelihood~~ likelihood.

6,

8,

- ~~Overfitt~~ Prevent overfitting by:
  - +, Early stopping
  - +, L1 , L2 regularization ~~and~~
  - +, Dropout .

- We can add regularization term to loss function:

Ex of L2:

$$E = \frac{1}{2} \sum_{i=1}^{m} (\bar{y}^{(i)} - y^{(i)})^2 + \lambda \|G\|^2$$

$\downarrow$ update equation:

$$w_j \leftarrow w_j - \eta \frac{\partial E}{\partial w_j} \ast - \lambda \|w_j\|$$

$\downarrow$ Basically weight decay.

— Scaling input is important ~~some~~ since some input may have way <u>bigger</u> values than the other, which will have more influence on output  even though ~~that may not~~ they may not be important features.

— One possible approach to adapt learning rate is learning rate decay, for ex:

$$\eta^{(i)} \ast = \frac{\eta^{(i-1)}}{2}.$$

$\downarrow$ · Useful since as we approach the minima, we ~~nee~~ may  need smaller and smaller learning rate to avoid overshooting the minima.

— ~~Moment u~~ Update equation with momentum added:

$$w_j^{(t+1)} \leftarrow w_j^{(t)} - \eta \frac{\partial E}{\partial w_j} + \beta(w_j^{(t)} - w_j^{(t-1)})$$

$\downarrow$ useful to help us ~~avoid~~ overcome local minimum.

— Dropout is performed by randomly "shutdown" units in hidden layers during ~~their~~ training

$\rightarrow$ A·Prevent ~~adapted~~ Co-adaptation

9,

CNN is used to process image

The idea is that we arrange the parameters in an ~~n to~~ n×n matrix
(dotproduct)
called filter and we perform convolution between ~~the~~ the Filter
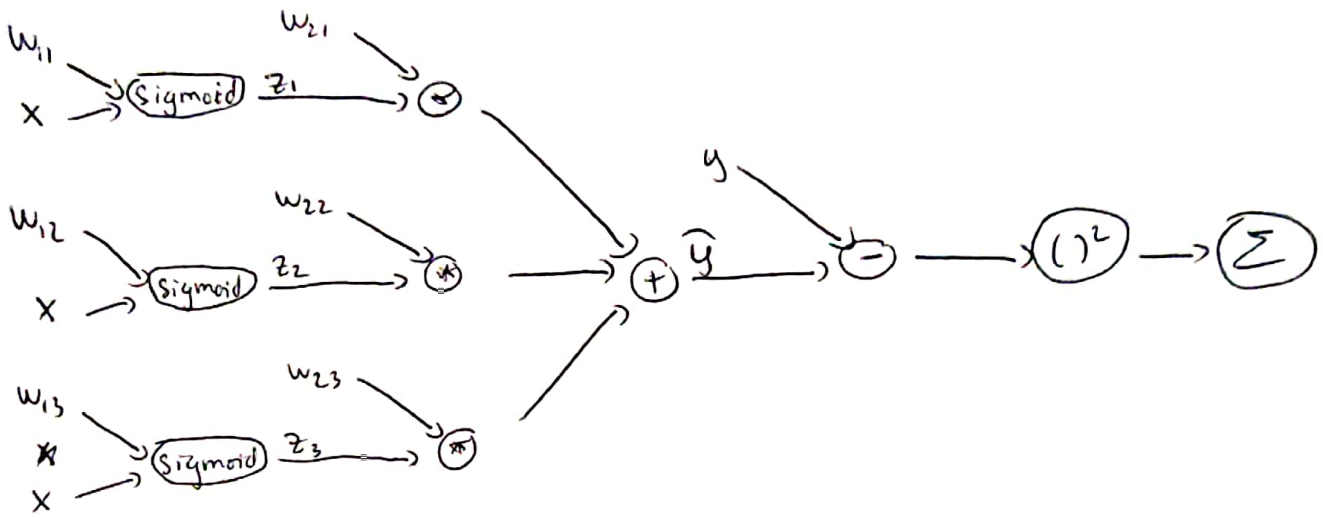
and each n × n ~~ve~~ region of the image

⌐ ~~Advan~~ ~~A tot~~ Advantages of doing this include ~~we get~~:

+, We can keep the ~~image~~ input image in its Original Form
and not have to vectorize the image

+, There are ~~mo~~ ~~much~~ Fewer parameter Compared to vectorizing
the image and feed into MLP.

(A) Deep learning.

1,