

CS 472 Midterm #1

1. Single Layer perceptron.

$$\Delta w_i = c(t-z)x_i$$

$$c = 0.6 \quad w = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

First Epoch:

$$net_0 = 0 \quad z_0 = 0$$

$$\Delta w_0 = 0.6(1-0)1 = 0.6 \quad net_1 = 0.5(0.6) + 0.7(-0.3) + 0.5(0.6) + 0.6 = .99 \quad z_1 = 1$$

$$\Delta w_1 = 0.6(1-0) - 0.5 = -0.3$$

$$\Delta w_2 = 0.6(1-0)1 = 0.6$$

$$\Delta w_3 = 0.6(1-0)1 = 0.6$$

$$\Delta w = \begin{bmatrix} -0.3 \\ -0.42 \\ 0.3 \\ -0.6 \end{bmatrix} \quad w = \begin{bmatrix} 0.3 \\ -0.72 \\ 0.3 \\ 0 \end{bmatrix}$$

$$net_2 = 0.3 + (-0.5)(-0.72) + 0.5(0.3) + 0 = 0.21 \quad z_2 = 1$$

$$\Delta w = [0 \ 0 \ 0 \ 0] \quad w = [0.3 \ -0.72 \ 0.3 \ 0]$$

$$net_3 = 0 + (-0.72) + 0.3^2 = -0.63 \quad z_3 = 0$$

$$\Delta w = [0 \ 0 \ 0 \ 0] \quad w = [0.3 \ -0.72 \ 0.3 \ 0]$$

~~convergence approached the initial value
of epoch.~~

Second Epoch

$$net_0 = 0.3 + 0.5(0.72) + 0.3 = 0.96 \quad z_0 = 1$$

No weight changes necessary.

$$net_1 = 0.5(0.3) - 0.7(0.72) + 0.5(0.3) = -0.204 \quad z_1 = 0$$

No weight changes necessary. We have gone through all 4 datapoints with no weight changes, therefore convergence has been achieved.

$$w = [0.3 \ -0.72 \ 0.3 \ 0]$$

2. Feature Reduction / Selection

A. PCA Algorithm :

$X = \text{dataset}$

$\text{Means} = \text{mean}(\text{for each column of } X)$ # find mean of each
feature

$X - \text{Means}$ # subtract the mean of the corresponding
feature from each value

$\text{covariance}(X)$ # create the covariance matrix for X

$\text{eigenvalues}(X)$ # find the eigenvalues & eigenvectors for
each column of X (feature)

$p = 2$ # how many of the ~~extra~~ eigenvectors with
the highest eigenvalues we will select

sort (eigenvectors by eigenvalues)
eigenvectors = top p rows

$T = \text{fancy_matrix_multiply}(\text{eigenvectors}, \text{stuff})$

bar-graph (T) # This will display how much information
we get from these particular features

B. Greedy Wrapper Algorithm: Backwards Search

$S = \text{set of features}$

~~then consider features~~ while ($S.\text{size} > \text{desired_size}$):

for each feature in S :

train (feature)

$a = \text{accuracy}(\text{feature})$

5. remove (feature with lowest accuracy)

C. Pros & Cons.

The backwards search can be very thorough, but also quite expensive since we are training our model many times.

Principal Component Analysis does not have many cons except that there is no set threshold of information at which a feature should be removed, which could cause us to remove a feature that the model actually did need. A pro is that it does not require training to use this algorithm.

3. Backpropagation.

$$w_{21} = -1.5 \quad w_{52} = -1 \quad w_{53} = -1.2 \\ w_{31} = 2 \quad w_{62} = 0.4 \quad w_{63} = -0.5 \quad x = \begin{bmatrix} 1 \\ 0.6 \end{bmatrix} \quad t = 1 \\ w_{41} = 0.5 \quad w_{72} = 1.2 \quad w_{73} = 1.6 \quad c = 0.5$$

$$\Delta w_{ij} = c a_i \delta_j \quad a_i = \frac{1}{1+e^{-\text{net}_i}} \quad f'(\text{net}_i) = a_i(1-a_i)$$

output: $\delta_i = (t - a_i) f'(\text{net}_i)$

hidden: $\delta_i = \sum_k (\delta_k w_{ik}) f'(\text{net}_i)$

$$\text{net}_2 = -0.1 + 0.6(0.4) + 1.2 = 1.34 \quad a_2 = \frac{1}{1+e^{-1.34}} = 0.79$$

$$\text{net}_3 = -1.2(0.1) - 0.6(0.5) + 1.6 = 1.18 \quad a_3 = \frac{1}{1+e^{-1.18}} = 0.77$$

$$\text{net}_1 = -0.79(1.5) + 0.77(2) + 0.5 = 0.86 \quad a_1 = \frac{1}{1+e^{-0.86}} = 0.70$$

$$\delta_1 = (1 - 0.7) 0.7(1 - 0.7) = 0.063$$

$$\Delta w_{21} = 0.5(0.79)(0.063) = \Delta w_{21} = 0.025$$

$$\Delta w_{31} = 0.5(0.77)(0.063) = \Delta w_{31} = 0.024$$

$$\Delta w_{41} = 0.5(1)(0.063) = \Delta w_{41} = 0.032$$

$$\delta_2 = 0.063(-1.5) 0.79(1 - 0.79) = \delta_2 = -0.016$$

$$\delta_3 = 0.063(2) 0.77(1 - 0.77) = \delta_3 = 0.022$$

$$\Delta w_{52} = 0.5(0.1)(-0.016) = \Delta w_{52} = -0.0008$$

$$\Delta w_{62} = 0.5(0.6)(-0.016) = \Delta w_{62} = -0.0048$$

$$\Delta w_{72} = 0.5(1)(-0.016) = \Delta w_{72} = -0.008$$

$$\Delta w_{53} = 0.5(0.1)(0.022) = \Delta w_{53} = 0.0011$$

$$\Delta w_{63} = 0.5(0.6)(0.022) = \Delta w_{63} = 0.066$$

$$\Delta w_{73} = 0.5(1)(0.022) = \Delta w_{73} = 0.011$$

$w_{21} = -1.48$	$w_{52} = -1.0008$	$w_{53} = -1.19$
$w_{31} = 2.024$	$w_{62} = 0.39$	$w_{63} = -0.43$
$w_{41} = 0.53$	$w_{72} = 1.19$	$w_{73} = 1.61$

4. Quadratic Machine.

I think

- A. I don't know what you mean by layers, but[^] there will be 2, one for the normal inputs & another for the quadratic inputs.
- B. We would use the perceptron rule since it is for binary classification.
- C. It does well with binary classification problems with data that isn't necessarily linearly separable.
- D. $\text{net} = w_x x + w_y y + w_{x^2} x^2 + w_{y^2} y^2 + w_{xy} xy$
 $= 0.1(-0.2) + 0.1(1.1) + 0.1(0.04) + 0.1(1.21) + 0.1(-0.22)$
 $= 0.193$
 $z = 1$

No weight changes. All weights still = 0.1

5. Overfitting

- A. Overfitting can happen due to too much training on the training set, or because we've added in duplicates of datapoints to give our set a more normal distribution. It can also be caused by using the same training set throughout all training.
- B. 1. N-fold cross-validation. This involves dividing a dataset into N equally sized subsets, then training the model, starting with $k=1$ & going to N , with dataset minus subset $_k$, then testing the model's accuracy on subset $_k$. This will work for any algorithm.
2. Random split training & validation. ~~separate~~ Train the model multiple times each time using a different split of the dataset for the training & validation sets. Ex: first training, 90% training set, 10% validation set, second training, 70% training set, 30% validation set. ~~the~~ The data should be shuffled each time to reduce overlap across training sessions. This will work for any algorithm.

3. Make sure a feature isn't too highly correlated with the output. For example, if you collected all the data for democrats on a Wednesday, & all the data for republicans on a Tuesday, & included the day of data collection as a feature, the model would overfit a ton & be helpless when presented w/ the real data. This applies for any algorithm.

6. Decision Tree

$$Info(S) = \sum_{i=1}^{|S|} \frac{|S_i|}{|S|} Info(S_i) \quad Info(S_i) =$$

~~Info_{avg}~~ Info_{avg} (Ent)

comprob

$$Prob_N = 4/7$$

$$Prob_P = 3/7$$

$$Prob_{single} = 4/7$$

$$Prob_{S,N} = 2/4$$

$$Prob_{S,P} = 2/4$$

$$Info_{single} = -\frac{4}{7} \left(\frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2} \right)$$

$$= 0.571$$

$$Prob_{Married} = 3/7$$

$$Prob_{M,N} = 2/3$$

$$Prob_{M,P} = 1/3$$

$$Info_{Married} = -\frac{3}{7} \left(\frac{2}{3} \log_2 \frac{2}{3} + \frac{1}{3} \log_2 \frac{1}{3} \right)$$

$$= 0.394$$

$$Prob_T = 4/7$$

$$Prob_{T,N} = 2/4$$

$$Prob_{T,P} = 2/4$$

$$Info_{Married, short} = 0.965$$

$$Info_T = 0.571$$

$$Prob_{short} = 3/7$$

$$Prob_{S,N} = 2/3$$

$$Prob_{S,P} = 1/3$$

$$Info_{short} = 0.394$$

$$Info_{height} = 0.965$$

$$\text{Prob}_{H_i} = 3/7$$

$$\text{Prob}_{H_i, N} = 1/3$$

$$\text{Prob}_{H_i, P} = 2/3$$

$$\text{Info}_{H_i} = 0.394$$

$$\text{Prob}_{M_{\text{red}}} = 2/7$$

$$\text{Prob}_{M_i, N} = 1/2$$

$$\text{Prob}_{M_i, P} = 1/2$$

$$\begin{aligned}\text{Info}_{M_{\text{red}}} &= -\frac{2}{7} \left(\frac{1}{2} \log_2 \left(\frac{1}{2} \right) + \frac{1}{2} \log_2 \left(\frac{1}{2} \right) \right) \\ &= 0.286\end{aligned}$$

$$\text{Prob}_{L_0} = 2/7$$

$$\text{Prob}_{L_i, N} = 2/2$$

$$\text{Prob}_{L_i, P} = 0/1$$

$$\begin{aligned}\text{Info}_{L_0} &= -2 \left(\log_2(1) \right) \\ &= 0\end{aligned}$$

$$\text{Info}_{GPA} = 0.68$$

Due to having the least information (entropy), we should split on GPA first.

7

10. Dealing with unknown/missing data.

1. Set the missing value equal to the mean or mode of the rest of the datapoints for that feature.
2. Simply include "missing" as another possible value for the feature. Works best for nominal data.
3. Train another model without that feature, & when the original model encounters a datapoint with that missing value, it will send that datapoint to the new model to be classified.

perception rule: $y=1$ if net > 0, bias input always 1. $\Delta w_i = C(t-y)x_i$, output, good for binary classification

delta rule: $\Delta w_i = C(t-y)x_i$. Not good for classification, but good for regression.

stochastic update/gradient descent: update after each iteration

batch update/gradient descent: update after each epoch.

Quadratic Machine: uses perceptron rule to update weights but also includes quadratics of inputs multiplied together. x_1^2, y^2, xy

Decision Boundary equation: $X_2 = (-w_1/w_2)X_1 + \theta/w_2$

$$w_1x_1 + w_2x_2 = \theta \rightarrow X_2 + \frac{w_1}{w_2}x_1 = \frac{\theta}{w_2}$$

Delta rules wants to minimize SSE of a regression line.

Perceptron rule wants to maximize correct classification.

Backprop: $\Delta w_{ij} = C O_i \delta_j$ "output" of input nodes is just their input value

$$O_i = \frac{1}{1+e^{-net_i}}$$

$$f'(net_i) = O_i(1-O_i)$$

$$\delta_i = (t_i - O_i) f'(net_i) * \text{output node}$$

$$\delta_i = \sum_k (\delta_k w_{ik}) f'(net_i) * \text{hidden node}$$

N -Fold cross-validation: use all data. For training & validation.

momentum in backprop: $\Delta w_{ij}(t+1) = \alpha \Delta w_{ij}(t) + \beta \Delta w_{ij}(t)$

$$\alpha \text{ usually } = 0.9, \text{ range of values for hidden nodes } = \text{Uniform}(5, 50)$$

Grid search: choose potential hyperparameters & try all combinations. Only useful if not a lot of hyperparameter options.

$L_1: E|t_i - z_i|$ Info(S) = data with predicted $E \frac{|z_i|}{|t_i|}$ Info(S_i) otherwise

$L_2: E(t_i - z_i)^2$ Info(S) = data with predicted $E \frac{1}{|t_i|}$

MSE: $\frac{SSE}{\# \text{ of datapoints}}$ all depends on current set

RMSSE: \sqrt{MSE} 1. calculate baseline entropy of S

If multiple outputs 2. calculate entropy of each target, $L_1, L_2, \text{ & MSE}$ 3. split S based on target

then δ if output layer for each target, t_i , output on previous layer.

then δ if output layer for each target, t_i , output on previous layer.

then weight changes for be $\sqrt{MSE_{total}}$ 5. repeat until new S.

values leading to output nodes.