

# Machine Learning

Important points and calculations

- Kathirmani Sukumar

# Disclaimer

- Please use LMS notes for introduction and theory about each algorithm
- This notes is mainly to cover important points discussed in class

# Naïve Bayes

Outlook	Temperature	Humidity	Wind	Play Tennis
Sunny	Hot	High	Weak	No
Sunny	Hot	High	Strong	No
Overcast	Hot	High	Weak	Yes
Rain	Mild	High	Weak	Yes
Rain	Cool	Normal	Weak	Yes
Rain	Cool	Normal	Strong	No
Overcast	Cool	Normal	Strong	Yes
Sunny	Mild	High	Weak	No
Sunny	Cool	Normal	Weak	Yes
Rain	Mild	Normal	Weak	Yes
Sunny	Mild	Normal	Strong	Yes
Overcast	Mild	High	Strong	Yes
Overcast	Hot	Normal	Weak	Yes
Rain	Mild	High	Strong	No

# Frequency Tables

Outlook	Play=Yes	Play=No	Temperature	Play=Yes	Play=No
<i>Sunny</i>	2/9	3/5	<i>Hot</i>	2/9	2/5
<i>Overcast</i>	4/9	0/5	<i>Mild</i>	4/9	2/5
<i>Rain</i>	3/9	2/5	<i>Cool</i>	3/9	1/5

Humidity	Play=Yes	Play=No
<i>High</i>	3/9	4/5
<i>Normal</i>	6/9	1/5

Wind	Play=Yes	Play=No
<i>Strong</i>	3/9	3/5
<i>Weak</i>	6/9	2/5

$P(\text{Play=Yes}) = 9/14$   
 $P(\text{Play=No}) = 5/14$

# Predictions

Given a new instance, predict its label

$X = (\text{Outlook}=\text{Sunny}, \text{Temperature}=\text{Cool}, \text{Humidity}=\text{High}, \text{Wind}=\text{Strong})$

$$P(\text{Outlook}=\text{Sunny} \mid \text{Play}=\text{Yes}) = (2 + 1) / (9 + 3) = 0.25$$

$$P(\text{Outlook}=\text{Sunny} \mid \text{Play}=\text{No}) = (3 + 1) / (5 + 3) = 0.5$$

$$P(\text{Temperature}=\text{Cool} \mid \text{Play}=\text{Yes}) = (3 + 1) / (9 + 3) =$$

$$P(\text{Temperature}=\text{Cool} \mid \text{Play}=\text{No}) = (1 + 1) / (5 + 3) = 0.25$$

0.33

$$P(\text{Humidity}=\text{High} \mid \text{Play}=\text{No}) = (4 + 1) / (5 + 2) = 0.71$$

$$P(\text{Humidity}=\text{High} \mid \text{Play}=\text{Yes}) = (3 + 1) / (9 + 2) = 0.36$$

$$P(\text{Wind}=\text{Strong} \mid \text{Play}=\text{No}) = (3 + 1) / (5 + 2) = 0.57$$

$$P(\text{Wind}=\text{Strong} \mid \text{Play}=\text{Yes}) = (3 + 1) / (9 + 2) = 0.36$$

$$P(\text{Play}=\text{No}) = 5 / 14 = 0.35$$

$$P(\text{Play}=\text{Yes}) = 9 / 14 = 0.64$$

$$P(\text{Yes} \mid X) = [P(\text{Sunny} \mid \text{Yes}) P(\text{Cool} \mid \text{Yes}) P(\text{High} \mid \text{Yes}) P(\text{Strong} \mid \text{Yes})] P(\text{Play}=\text{Yes}) = 0.0068$$

$$P(\text{Yes} \mid X) = [P(\text{Sunny} \mid \text{No}) P(\text{Cool} \mid \text{No}) P(\text{High} \mid \text{No}) P(\text{Strong} \mid \text{No})] P(\text{Play}=\text{No}) = 0.0177$$

Given the fact  $P(\text{Yes} \mid X) < P(\text{No} \mid X)$ . We label  $X$  to No

# K Nearest Neighbor

X1	X2	Y
7	7	Bad
7	4	Bad
3	4	Good
1	4	Good

Predict Y if X1=3, X2=7; K= 3

# K Nearest Neighbor

X1	X2	Euclidian Distance	Neighbor	Y
7	7	$\text{Sqrt}((7-3)^2 + (7-7)^2) = 4$	3	Bad
7	4	$\text{Sqrt}((7-3)^2 + (4-7)^2) = 5$	4	Bad
3	4	$\text{Sqrt}((3-3)^2 + (4-7)^2) = 3$	1	Good
1	4	$\text{Sqrt}((1-3)^2 + (4-7)^2) = 3.6$	2	Good

Among the first 3 Neighbors, 2 samples are Good and 1 sample is Bad; Prediction is Good

X1	X2	Euclidian Distance	Neighbor	Y
7	7	$\text{Sqrt}((7-3)^2 + (7-7)^2) = 4$	3	Bad
7	4	$\text{Sqrt}((7-3)^2 + (4-7)^2) = 5$	4	Bad
3	4	$\text{Sqrt}((3-3)^2 + (4-7)^2) = 3$	1	Good
1	4	$\text{Sqrt}((1-3)^2 + (4-7)^2) = 3.6$	2	Good

# Random Forest

- Each tree is grown fully to reduce bias; Due to this variance will be high
- To reduce variance, multiple such trees are grown
- For each tree, 60% of samples are randomly selected
- While randomly selecting samples, each sample is given equal weight
- For each tree,  $\sqrt{\text{No. of cols}}$  are randomly selected
- Prediction is based on polling of predictions from all trees



# Adaboost

- In Random Forest, each sample is given equal weight while sampling
- In Adaboost, while building first decision tree, each sample is given equal weight
- While building second decision tree, for each sample new weights are calculated
- Samples which are misclassified in first decision tree are given higher weights in second decision tree

# Sample Data

Outlook	Temperature	Humidity	Wind	Target Original	Target New
Sunny	Hot	High	Weak	No	-1
Sunny	Hot	High	Strong	No	-1
Overcast	Hot	High	Weak	Yes	1
Rain	Mild	High	Weak	Yes	1
Rain	Cool	Normal	Weak	Yes	1
Rain	Cool	Normal	Strong	No	-1
Overcast	Cool	Normal	Strong	Yes	1
Sunny	Mild	High	Weak	No	-1
Sunny	Cool	Normal	Weak	Yes	1

N = 9

# First Tree Prediction on Training Data

Sample Weight	Actual	Predicted	Correct Prediction?
1/N = 1 / 9	-1	1	Wrong
1 / 9	-1	-1	Correct
1 / 9	1	-1	Wrong
1 / 9	1	1	Correct
1 / 9	1	1	Correct
1 / 9	-1	1	Wrong
1 / 9	1	1	Correct
1 / 9	-1	-1	Correct
1 / 9	1	1	Correct

## Error Rate

Error Rate = Wrong / Total Predictions

Error Rate = 3 / 9 = 0.33

## Classifier's weight

$$\alpha_t = \frac{1}{2} \ln \left( \frac{1 - \epsilon_t}{\epsilon_t} \right)$$

```
(1 / 2) * np.log((1-0.33) / 0.33)
```

```
0.3540925289622428
```

# Samples new weight

$$D_{t+1}(i) = \frac{D_t(i) \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$

- $Y_i$  = Actual value of  $i^{\text{th}}$  sample
- $H_t(x_i)$  = Prediction value of  $i^{\text{th}}$  sample
- $\text{Alpha}_t$  = Classifiers weight

Sample Weight ( $D_1$ )	Actual	Predicted	Correct Prediction?	Sample Weight ( $D_2$ ) (numerator)
$1/N = 1/9$	-1	1	Wrong	$(1/9) \exp(-0.35*(-1)*(1)) = 0.15$
$1/9$	-1	-1	Correct	$(1/9) \exp(-0.35*(-1)*(-1)) = 0.07$
$1/9$	1	-1	Wrong	$(1/9) \exp(-0.35*(1)*(-1)) = 0.15$
$1/9$	1	1	Correct	$(1/9) \exp(-0.35*(1)*(1)) = 0.07$
$1/9$	1	1	Correct	...
$1/9$	-1	1	Wrong	...
$1/9$	1	1	Correct	...
$1/9$	-1	-1	Correct	...
$1/9$	1	1	Correct	...

- Samples which are misclassified given weight as 0.15
- Samples which are correctly classified given weight as 0.07
- Samples which are misclassified are given more weight
- $Z_t$  is nothing but summation of all new weights
- Mainly used to make values add up to 1
- Second decision tree gives more weights to misclassified samples

# Weighted Polling

$$H(x) = \text{sign} \left( \sum_{t=1}^T \alpha_t h_t(x) \right)$$

- While polling prediction from each tree is added together along with classifier weight
- It is called as weighted polling
- If sum comes out to be positive, final prediction is +1
- If sum comes out to be negative, final prediction is -1