

Naïve Bayes Classifier & Decision Tree

- Although KNN is used for supervised and K-Means is used for unsupervised, both of them find the exact label for one data point
- Naïve Bayes find the probability that data point will be in each label: $p(y = c | \mathbf{x})$ or $p(c | \mathbf{x})$ in short, then $c = \arg \max_{c \in \{1, \dots, C\}} p(c | \mathbf{x})$
- $p(c | \mathbf{x})$ means that given data point \mathbf{x} , the probability this point is in class c , it is like inference process in test set, but now in training set, we need to apply Bayes Theorem:

$$c = \arg \max_{c \in \{1, \dots, C\}} p(c | \mathbf{x}) = \arg \max_{c \in \{1, \dots, C\}} \frac{p(\mathbf{x} | c) p(c)}{p(\mathbf{x})} = \arg \max_{c \in \{1, \dots, C\}} p(\mathbf{x} | c) p(c)$$

- As mentioned in chapter 4, $p(\mathbf{x} | c)$ hardly finds out. To solve it, we assume independence among the features in one data point

$$p(\mathbf{x} | c) = p(x_1, x_2, \dots, x_d | c) = \prod_{i=1}^d p(x_i | c)$$

- Independence among the features assumption is reason why we call “naive” because in the reality, the features in one dataset are hardly completely independent
- Due to “Naive”, NBC has fast speed while training
- When number of data points (d) is big, $p(x_i | c) \rightarrow 0$, logarithm can solve the problem

$$c = \arg \max_{c \in \{1, \dots, C\}} \left(\log(p(c)) + \sum_{i=1}^d \log(p(x_i | c)) \right)$$

- Finding out $p(x_i | c)$ is based on the distribution of the dataset: *Gaussian Naïve Bayes*, *Multinomial Naïve Bayes*, *Bernoulli Naïve Bayes*

Estimating Probabilities

- In the following dataset, $P(\text{Humidity} = \text{'Low'} | \text{class} = \text{'Yes'}) = P(\text{Humidity} = \text{'Low'} | \text{class} = \text{'No'}) = 0 \rightarrow P(X | \text{class} = \text{'Yes'}) = P(X | \text{class} = \text{'No'}) = 0$, so you cannot categorize the X
- Estimating probabilities: $\frac{n_c + mp}{n + m}$
- Given the dataset, calculate Naïve Bayes to find the class of **last row**

Example	Sky	AirTemp	Humidity	Wind	Water	Forecast	EnjoySport
1	Sunny	Warm	Normal	Strong	Warm	Same	Yes
2	Sunny	Warm	High	Strong	Warm	Same	Yes
3	Rainy	Cold	High	Strong	Warm	Change	No
4	Sunny	Warm	High	Strong	Cool	Change	Yes
5	Cloudy	Warm	High	Weak	Cool	Same	Yes
6	Cloudy	Cold	High	Weak	Cool	Same	No

7	Sunny	Warm	Normal	Strong	Warm	Same	?
8	Sunny	Warm	Low	Strong	Cool	Same	?

$$\begin{aligned} \circ P(C_i): & \begin{cases} P(\text{class} = \text{'yes'}) = \frac{4}{6} = \frac{2}{3} \\ P(\text{class} = \text{'no'}) = \frac{2}{6} = \frac{1}{3} \end{cases} \\ \circ P(\text{Sky} = \text{'Sunny'} | \text{class} = \text{'Yes'}) &= \frac{3}{4} = \frac{3}{4} \\ \circ P(\text{Sky} = \text{'Sunny'} | \text{class} = \text{'No'}) &= 0 \\ \circ P(\text{Sky} = \text{'Rainy'} | \text{class} = \text{'Yes'}) &= 0 \\ \circ P(\text{Sky} = \text{'Rainy'} | \text{class} = \text{'No'}) &= 1 \end{aligned}$$

- Notice: value ‘Low’ is not in ‘Humidity’, so $P(\text{Humidity} = \text{'Low'} | \text{class} = \text{'Yes'}) = P(\text{Humidity} = \text{'Low'} | \text{class} = \text{'No'}) = 0 \rightarrow$ Cannot decide \rightarrow Apply Estimating Probabilities
- $P(\text{Humidity} = \text{'Low'} | \text{class} = \text{'Yes'}) =$

$$\frac{(\text{\# of rows of class 'yes' having 'low'}) + (\text{\# of rows of dataset}) \left(\frac{1}{\text{\# of unique class values}} \right)}{\text{\# of rows of class 'yes' + \# of rows of dataset}} = \frac{0 + 6 * \frac{1}{2}}{4 + 6} = 0.3$$

$$\circ P(\text{Humidity} = \text{'Low'} | \text{class} = \text{'No'}) = \frac{0 + 6 * \frac{1}{2}}{2 + 6} = 0.375$$

- Put $X = (\text{Sunny, Warm, Low, Strong, Cool, Same})$
- $P(X | \text{Class} = \text{'Yes'}) = P(\text{Sky} = \text{'Sunny'} | \text{class} = \text{'Yes'}) * P(\text{AirTemp} = \text{'Warm'} | \text{class} = \text{'Yes'}) * P(\text{Humidity} = \text{'Low'} | \text{class} = \text{'Yes'}) * P(\text{Wind} = \text{'Strong'} | \text{class} = \text{'Yes'}) * P(\text{Water} = \text{'Cool'} | \text{class} = \text{'Yes'}) * P(\text{Forecast} = \text{'Same'} | \text{class} = \text{'Yes'})$
- $P(X | \text{Class} = \text{'No'}) = P(\text{Sky} = \text{'Sunny'} | \text{class} = \text{'No'}) * P(\text{AirTemp} = \text{'Warm'} | \text{class} = \text{'No'}) * P(\text{Humidity} = \text{'Low'} | \text{class} = \text{'No'}) * P(\text{Wind} = \text{'Strong'} | \text{class} = \text{'No'}) * P(\text{Water} = \text{'Cool'} | \text{class} = \text{'No'}) * P(\text{Forecast} = \text{'Same'} | \text{class} = \text{'No'})$
- Compare: $P(X | \text{Class} = \text{'Yes'}) * P(\text{Class} = \text{'Yes'})$ and $P(X | \text{Class} = \text{'No'}) * P(\text{Class} = \text{'No'})$

Common distribution in NBC

Gaussian Naïve Bayes

- For feature i and class c , x_i follows the Gaussian distribution of μ_{ci} and σ_{ci}

$$p(x_i | c) = p(x_i | \mu_{ci}, \sigma_{ci}^2) = \frac{1}{\sqrt{2\pi\sigma_{ci}^2}} \exp\left(-\frac{(x_i - \mu_{ci})^2}{2\sigma_{ci}^2}\right)$$

in which parameter $\theta = \{\mu_{ci}, \sigma_{ci}^2\}$ is determined based on the data points of class c

Multinomial Naïve Bayes (Example in CBD/Week06)

- This model is usually used in text classification in which data point is built on BoW ideas (Bag of words)
- Each data point has number of feature of d (number of words in dictionary). x_i in the data point is number of word i (in the dictionary) appearing in the sentence
- $p(x_i | c)$ ratio of frequency of word i (or feature i in general case) appearing in the sentence of class c

$$p(x_i | c) = \lambda_{ci} = \frac{N_{ci}}{N_c}$$

in which: N_{ci} : total number of word i (feature i) appearing in the sentence of class c , N_c : total number of all words appearing in the sentence of class c

- But if the word does not appear in any sentence of class c , $p(\text{word} | c) = 0$, so if we apply $p(\mathbf{x} | c) = \prod_{i=1}^d p(x_i | c)$ will be 0.

Laplace Smoothing can solve it

$$p(x_i | c) = \hat{\lambda}_{ci} = \frac{N_{ci} + \alpha}{N_c + d\alpha} \text{ in which } \alpha > 0 \text{ and usually } = 1, d: \text{number of words in dictionary, } \sum_{i=1}^d \hat{\lambda}_{ci} = 1, \text{ so each sentence of}$$

class c will have $\hat{\lambda}_c = \{\hat{\lambda}_{c1}, \hat{\lambda}_{c2}, \dots, \hat{\lambda}_{cd}\}$

Bernoulli Naïve Bayes

- This model is applied when each element has value of 0 and 1. E.g. Rather than BoW, we just consider whether the word appears in the sentence or not

$$p(x_i | c) = p(i | c) 1\{x_i = c\} + (1 - p(i | c)) (1 - 1\{x_i \neq c\})$$

in which $p(i | c)$: (the meaning is like $p(x_i | c)$ in Multinomial Naïve Bayes)

Example

North or South

- In training set, we have corpus including 4 documents d1, d2, d3, d4

	Document	Content	Class
Training set	d1	hanoi pho chaolong hanoi	N

	d2	hanoi buncha pho omai	N
	d3	pho banhgio omai	N
	d4	saigon hutiu banhbo pho	S
Test set	d5	hanoi hanoi buncha hutiu	?

- Intuitively, we can predict d5 is class of N
- The problem can be solved by *Multinomial Naïve Bayes* or *Bernoulli Naïve Bayes*. We should test 2 models to choose the best one. Now *Multinomial Naïve Bayes* is implemented

- First we find $p(c)$:
$$\begin{cases} p(N) = \frac{3}{4} \\ p(S) = \frac{1}{4} \end{cases}$$

- We got dictionary $V = \{\text{hanoi, pho, chaolong, buncha, omai, banhgio, saigon, hutiu, banhbo}\} \rightarrow d = |V| = 9$
- TRAINING
 - class = N

	hanoi	pho	chaolong	buncha	omai	banhgio	saigon	huti	banhbo	
d1: \mathbf{x}_1	2	1	1	0	0	0	0	0	0	
d2: \mathbf{x}_2	1	1	0	1	1	0	0	0	0	
d3: \mathbf{x}_3	0	1	0	0	1	1	0	0	0	
Total	3	3	1	1	2	1	0	0	0	$N_N = 11$
$\rightarrow \hat{\lambda}_N$	4/20	4/20	2/20	2/20	3/20	2/20	1/20	1/20	1/20	$N_N + V \alpha = 20$

class = S

d4: \mathbf{x}_4	0	1	0	0	0	0	1	1	1	$N_S = 4$
$\rightarrow \hat{\lambda}_S$	1/13	1/13	1/13	1/13	1/13	1/13	2/13	2/13	2/13	$N_S + V \alpha = 13$

- TEST

d5: \mathbf{x}_5	2	0	0	1	0	0	0	1	0	
--------------------	---	---	---	---	---	---	---	---	---	--

$d5 = [\text{hanoi, hanoi, buncha, hutiu}]$

$$\begin{cases} p(c = N | d5) \propto p(c = N) \prod_{i=1}^{d5} p(x_i | N) = \frac{3}{4} \left(\frac{4}{20} \right)^2 \frac{2}{20} \frac{1}{20} \approx 1.5 \times 10^{-4} \\ p(c = S | d5) \propto p(c = S) \prod_{i=1}^{d5} p(x_i | S) = \frac{1}{4} \left(\frac{1}{13} \right)^2 \frac{1}{13} \frac{2}{13} \approx 1.75 \times 10^{-5} \end{cases}$$

$$\rightarrow p(c = N | d5) > p(c = S | d5) \rightarrow d5 \in \text{class}(N)$$

- In above example, we use Laplace Smoothing $\alpha = 1$. 1.5×10^{-4} and 1.75×10^{-5} just helps you to find the class.

Note: $p(d5) = p(c = N | d5) + p(c = S | d5)$

$$\begin{cases} p(c = N | d5) = \frac{p(c = N) \prod_{i=1}^{d5} p(x_i | N)}{p(d5)} = \frac{1.5 \times 10^{-4}}{1.5 \times 10^{-4} + 1.75 \times 10^{-5}} = 0.8955 \\ p(c = S | d5) = 1 - p(c = N | d5) = 0.1045 \end{cases}$$

- So, probability that d5 is in class of N: 89.55% and in class of S: 10.45%

Summary

- NBC is usually applied in text classification
- Due to the independence assumption, NBC usually have fast training and testing
- If independence assumption is correct with the dataset we are considering, NBC will have better result than SVM and logistics regression
- Laplace Smoothing is used in MultinomialNB to avoid the absence of word in the dictionary in the class

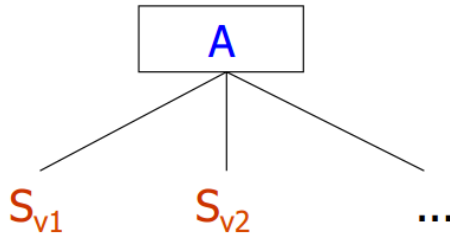
Decision Tree

ID3

- Entropy S: đo sự hỗn loạn của S (Sample).

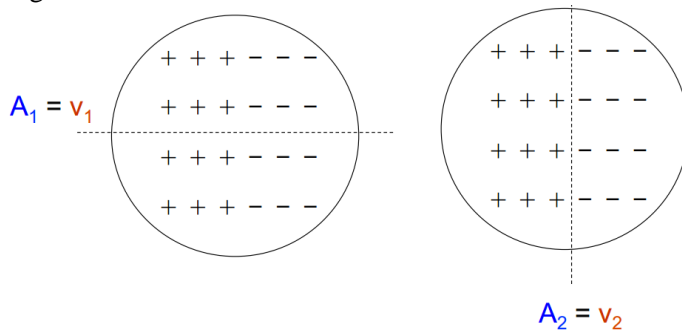
- E.g. If your dataset has 2 classes, $\text{Entropy}(S) = -p_+ \log_2(p_+) - p_- \log_2(p_-) = \sum_{i \in \{1, c\}} -p_i \log_2(p_i)$

- Gain(S, A) trong đó A là feature: lượng thông tin A chứa để giải thích cho S, so you choose feature which has maximum Gain



$$\text{Gain}(S, A) = \text{Entropy}(S) - \sum_{v \in \text{Values}(A)} (S_v / S) \cdot \text{Entropy}(S_v)$$

- E.g.



- Should choose $A_2 = v_2$ to minimize the entropy S

- Given the dataset, calculate ID3:

No.	Sky	AirTemp	Humidity	Wind	Water	Forecast	Enjoy
1	Sunny	Warm	Normal	Strong	Warm	Same	Yes
2	Sunny	Warm	High	Strong	Warm	Same	Yes
3	Rainy	Cold	High	Strong	Warm	Change	No
4	Sunny	Warm	High	Strong	Cool	Change	Yes
5	Cloudy	Warm	High	Weak	Cool	Same	Yes
6	Cloudy	Cold	High	Weak	Cool	Same	No

$$\text{Entropy}(S) =$$

$$-p_{\text{yes}} \log_2(p_{\text{yes}}) - p_{\text{no}} \log_2(p_{\text{no}}) = -\frac{4}{6} \log_2 \frac{4}{6} - \frac{2}{6} \log_2 \frac{2}{6} = 0.917$$

$$\text{Gain}(S, \text{Sky})$$

$$= \text{Entropy}(S) - \sum_{v \in \{\text{Sunny}, \text{Rainy}, \text{Cloudy}\}} (S_v / S) \text{Entropy}(S_v)$$

$$= \text{Entropy}(S) - \left[\frac{3}{6} \text{Entropy}(S_{\text{Sunny}}) + \frac{1}{6} \text{Entropy}(S_{\text{Rainy}}) + \frac{2}{6} \text{Entropy}(S_{\text{Cloudy}}) \right]$$

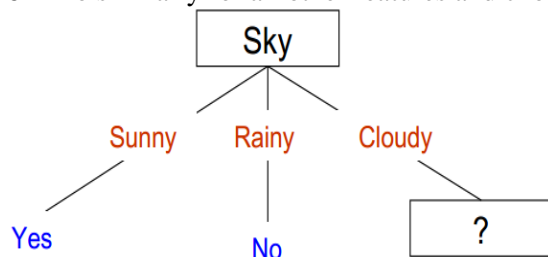
(Intuitively, you see that $\text{Entropy}(S_{\text{Sunny}}) = 0$ because Sunny just include one class: 'Yes' and $\text{Entropy}(S_{\text{Rainy}}) = 0$ because Rainy just include one class: 'No')

$$= \text{Entropy}(S) - \frac{2}{6} \text{Entropy}(S_{\text{Cloudy}})$$

$$= \text{Entropy}(S) - \frac{2}{6} \left(-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right)$$

$$= 0.584$$

- Do similarly for all other features and choose maximum Gain, and the largest Gain is feature Sky



- Delete all rows which does not include Cloudy in Sky

$$\text{Entropy}(S_{\text{Cloudy}}) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$$

$$\text{Gain}(S_{\text{Cloudy}}, \text{AirTemp})$$

$$\begin{aligned}
&= \text{Entropy}(S_{\text{Cloudy}}) - \sum_{v \in (\text{Warm}, \text{Cold})} (S_v / S) \text{Entropy}(S_v) \\
&= 1 - \left[\frac{1}{2} \text{Entropy}(S_{\text{Warm}}) + \frac{1}{2} \text{Entropy}(S_{\text{Cold}}) \right] \\
&= 1 - (0 + 0) \\
&= 1
\end{aligned}$$

C4.5

- One weakness of ID3: if one feature in your dataset is continuous, the ID3 is highly likely to get overfitting when every edge in Decision Tree of this continuous feature will be one number
- $\text{SplitInfo}_A D = - \sum_{v \in A} p_v \log_2 p_v$
- $\text{Gain}(A) = \text{Gain}(S, A)$
- $\text{GainRatio}(A) = \frac{\text{Gain}(A)}{\text{SplitInfo}_A D}$
- Choose feature max GainRatio
- Given the dataset, Calculate C4.5

RID	age	income	student	credit_rating	Class: buys_computer
1	youth	high	no	fair	no
2	youth	high	no	excellent	no
3	middle_aged	high	no	fair	yes
4	senior	medium	no	fair	yes
5	senior	low	yes	fair	yes
6	senior	low	yes	excellent	no
7	middle_aged	low	yes	excellent	yes
8	youth	medium	no	fair	no
9	youth	low	yes	fair	yes
10	senior	medium	yes	fair	yes
11	youth	medium	yes	excellent	yes
12	middle_aged	medium	no	excellent	yes
13	middle_aged	high	yes	fair	yes
14	senior	medium	no	excellent	no

○ $\text{SplitInfo}_{\text{Income}} S$

$$= -\frac{4}{14} \log_2 \frac{4}{14} - \frac{6}{14} \log_2 \frac{6}{14} - \frac{4}{14} \log_2 \frac{4}{14} = 0.926$$

○ $\text{Gain}(\text{Income})$

$$= \text{Entropy}(S) - \sum_{v \in (\text{low}, \text{medium}, \text{high})} (S_v / S) \text{Entropy}(S_v) = 0.94 - 0.91 = 0.03$$

$$\text{GainRatio}(\text{Income}) = \frac{0.03}{0.926} = 0.032$$

○ Do similarly for all other features and choose feature max GainRatio

Cart

- Binary Split for feature A. S_A is subset of A which has 1 or v-1 unique values of A
- $\text{Gini}(D) = 1 - \sum_{i \in \text{Values}(D)} p_i^2$
- $\text{Gini}_A D = \sum_{i \in \text{Values}(A)} p_i \text{Gini}(D_i)$
- $\Delta \text{Gini}(A) = \text{Gini}(D) - \text{Gini}_A D$
- Choose feature has min $\text{Gini}_A D$ or max $\Delta \text{Gini}(A)$
- Given the dataset, Calculate C4.5

RID	age	income	student	credit_rating	Class: buys_computer
1	youth	high	no	fair	no
2	youth	high	no	excellent	no
3	middle_aged	high	no	fair	yes
4	senior	medium	no	fair	yes
5	senior	low	yes	fair	yes
6	senior	low	yes	excellent	no
7	middle_aged	low	yes	excellent	yes
8	youth	medium	no	fair	no
9	youth	low	yes	fair	yes
10	senior	medium	yes	fair	yes
11	youth	medium	yes	excellent	yes
12	middle_aged	medium	no	excellent	yes
13	middle_aged	high	yes	fair	yes
14	senior	medium	no	excellent	no

$$\text{Gini}(S) = 1 - \left(\frac{9}{14} \right)^2 - \left(\frac{5}{14} \right)^2 = 0.46$$

○

$$\text{Gini}_{i \in (\text{Low}, \text{medium}) \in \text{income}}(S) = \sum_{i \in (\text{Low}, \text{medium})} p_i \text{Gini}(D_i)$$

$$= \frac{10}{14} \left(1 - \frac{6^2}{10} - \frac{4^2}{10} \right) + \frac{4}{14} \left(1 - \frac{1^2}{4} - \frac{3^2}{4} \right)$$

$$= 0.45 = \text{Gini}_{i \in (\text{High})}$$

$$\text{Gini}_{\text{income} \in (\text{Low}, \text{High})} = \text{Gini}_{\text{income} \in (\text{Medium})} = 0.315$$

$$\text{Gini}_{\text{income} \in (\text{Medium}, \text{High})} = \text{Gini}_{\text{income} \in (\text{High})} = 0.3$$

$$\text{You choose } \text{Gini}_{\text{income} \in (\text{Medium}, \text{High})} = 0.3$$

$$\text{Gini}_{\text{age} \in \{\text{youth}, \text{senior}\} / \{\text{middle_aged}\}} = 0.375$$

$$\text{Gini}_{\text{student}} = 0.367$$

$$\text{Gini}_{\text{credit rating}} = 0.429$$

○ Choose Income to split, $\text{Gini}_A D$ is min

	Splitting Criteria	Attribute type	Missing values	Pruning Strategy	Outlier Detection
ID3	Information Gain	Handles only Categorical value	Do not handle missing values.	No pruning is done	Susceptible to outliers
CART	Towing Criteria	Handles both Categorical & Numeric value	Handle missing values.	Cost-Complexity pruning is used	Can handle Outliers
C4.5	Gain Ratio	Handles both Categorical & Numeric value	Handle missing values.	Error Based pruning is used	Susceptible to outliers