

# Bayes & Uncertainty III

---

*PROF LIM KWAN HUI*

50.021 Artificial Intelligence

*The following notes are compiled from various sources such as textbooks, lecture materials, Web resources and are shared for academic purposes only, intended for use by students registered for a specific course. In the interest of brevity, every source is not cited. The compiler of these notes gratefully acknowledges all such sources.*



# Outline & Objectives

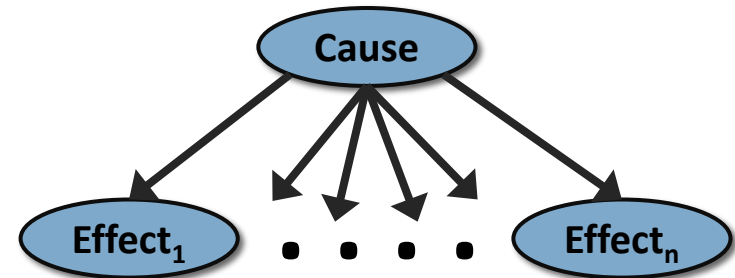
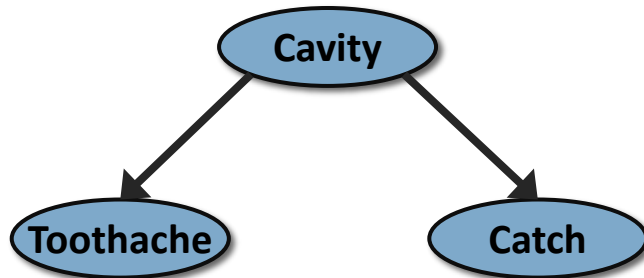
---

- Recap on statistical concepts such as product rule, chain rule, conditional independence, Bayes rules
- Able to represent a problem in terms of a Bayesian network and its corresponding conditional probability table
- Learn about how Bayes net can be used in various scenarios
- Learn about the Naïve Bayes Classifier and its application to text



# Bayes Rule and conditional independence

- $P(\text{Cavity} \mid \text{toothache}, \text{catch})$ 
  - $= \alpha P(\text{toothache}, \text{catch} \mid \text{Cavity}) P(\text{Cavity})$
  - $= \alpha P(\text{toothache} \mid \text{Cavity}) P(\text{catch} \mid \text{Cavity}) P(\text{Cavity})$
- This is an example of a naive Bayes model:
  - $P(\text{Cause}, \text{Effect}_1, \dots, \text{Effect}_n) = P(\text{Cause}) \prod_i P(\text{Effect}_i \mid \text{Cause})$

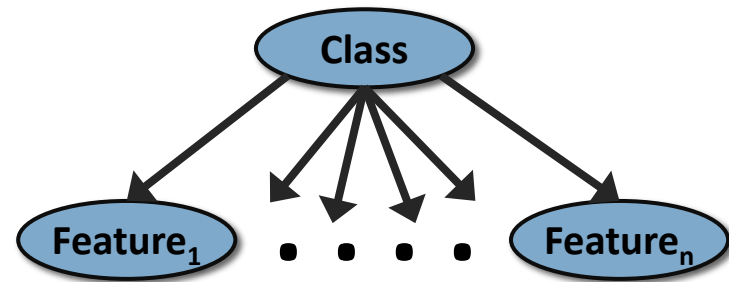
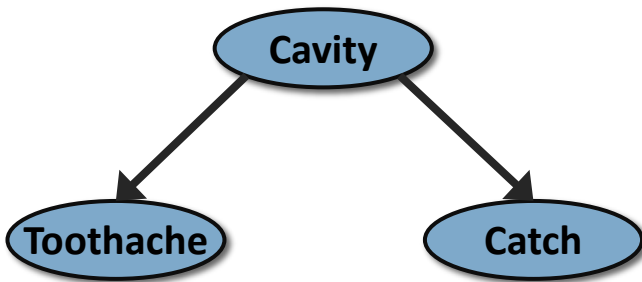


- Total number of parameters is linear in  $n$



# Naïve Bayes Classifier

- $P(\text{Cavity} \mid \text{toothache}, \text{catch})$ 
  - $= \alpha P(\text{toothache}, \text{catch} \mid \text{Cavity}) P(\text{Cavity})$
  - $= \alpha P(\text{toothache} \mid \text{Cavity}) P(\text{catch} \mid \text{Cavity}) P(\text{Cavity})$
- Similar to finding the most likely class given a set of features, e.g.,
  - $P(\text{cavity}=\text{true} \mid \text{toothache}, \text{catch})$  or  $P(\text{cavity}=\text{false} \mid \text{toothache}, \text{catch})$



# Bayes for Classification

---

- Given the following:
  - An observation  $\mathbf{o}$  represented by a feature set  $X_o = \{x_1, x_2, \dots, x_m\}$
  - A fixed set of classes  $Y = \{y_1, y_2, \dots, y_n\}$
- We are interested to classify the class  $Y$  that observation  $\mathbf{o}$  belongs to given its feature set  $X_o$

$$y_{MAP} = \operatorname{argmax} P(Y | X_o)$$



# Bayes for Classification

---

- We are interested to classify the class  $Y$  that observation  $o$  belongs to given its feature set  $X_o$ .
- Applying Bayes Theorem, we have:

$$y_{MAP} = \operatorname{argmax} P(Y | X_o)$$

$$y_{MAP} = \operatorname{argmax} \frac{P(X_o | Y) P(Y)}{P(X_o)}$$

$$y_{MAP} = \operatorname{argmax} P(X_o | Y) P(Y)$$

$$y_{MAP} = \operatorname{argmax} P(x_1, x_2, \dots, x_m | Y) P(Y)$$



# Example: Tax Avoidance

- Given an observation  $o$  with feature set  $X_o = \{\text{Industry}=\text{IT}, \text{Status}=\text{Married}, \text{Income}=\text{low}\}$ , how to estimate class  $Y$ ?

$$y_{MAP} = \operatorname{argmax} P(Y | X_o)$$

- Find probabilities by counting:
  - $P(Y) = N_c / N$ 
    - E.g.,  $P(\text{AvoidTax}=\text{Yes}) = 3/10$
  - $P(X_i | Y_k) = |X_{ik}| / N_c$ 
    - E.g.,  $P(\text{Industry}=\text{Sales} | \text{AvoidTax}=\text{Yes}) = 2/3$

ID	Industry	Marital Status	Income Level	Avoid Tax
1	IT	Single	High	No
2	Sales	Married	Medium	No
3	Sales	Single	Low	No
4	IT	Married	High	No
5	Sales	Divorced	Low	Yes
6	Sales	Married	Low	No
7	IT	Divorced	High	No
8	IT	Married	Medium	Yes
9	Sales	Married	Low	No
10	Sales	Single	Medium	Yes



# Example: Tax Avoidance

- Given an observation  $o$  with feature set  $X_o = \{\text{Industry}=\text{IT}, \text{Status}=\text{Married}, \text{Income}=\text{low}\}$ , how to estimate class  $Y$ ?

- $P(Y=\text{Yes} \mid X) = P(X \mid Y=\text{Yes}) P(Y=\text{Yes})$   
 $= P(\text{Industry}=\text{IT} \mid \text{Yes}) \times$   
 $P(\text{Status}=\text{Married} \mid \text{Yes}) \times$   
 $P(\text{Income}=\text{Low} \mid \text{Yes}) \times P(\text{Yes})$   
 $= 1/3 \times 1/3 \times 1/3 \times 3/10$
- $P(Y=\text{No} \mid X) = P(X \mid Y=\text{No}) P(Y=\text{No})$   
 $= P(\text{Industry}=\text{IT} \mid \text{No}) \times$   
 $P(\text{Status}=\text{Married} \mid \text{No}) \times$   
 $P(\text{Income}=\text{Low} \mid \text{No}) \times P(\text{No})$   
 $= 3/7 \times 4/7 \times 3/7 \times 7/10$

ID	Industry	Marital Status	Income Level	Avoid Tax
1	IT	Single	High	No
2	Sales	Married	Medium	No
3	Sales	Single	Low	No
4	IT	Married	High	No
5	Sales	Divorced	Low	Yes
6	Sales	Married	Low	No
7	IT	Divorced	High	No
8	IT	Married	Medium	Yes
9	Sales	Married	Low	No
10	Sales	Single	Medium	Yes





# Naïve Bayes on Text

---

- Given a document  $d$ , we want to find the most likely class  $c$ :

$$c_{NB} = \operatorname{argmax}_{c_i \in C} P(c_i | d)$$



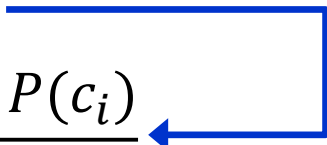
# Naïve Bayes on Text

---

- Given a document  $d$ , we want to find the most likely class  $c$ :

$$c_{NB} = \operatorname{argmax}_{c_i \in \mathcal{C}} P(c_i | d)$$
$$c_{NB} = \operatorname{argmax} \frac{P(d | c_i) P(c_i)}{P(d)}$$

**Bayes Rule**





# Naïve Bayes on Text

---

- Given a document  $d$ , we want to find the most likely class  $c$ :

$$c_{NB} = \operatorname{argmax}_{c_i \in \mathcal{C}} P(c_i | d)$$

$$c_{NB} = \operatorname{argmax} \frac{P(d | c_i) P(c_i)}{P(d)}$$

$$c_{NB} = \operatorname{argmax} P(d | c_i) P(c_i)$$

**Normalization**



# Naïve Bayes on Text

- Given a document  $d$ , we want to find the most likely class  $c$ :

$$c_{NB} = \operatorname{argmax}_{c_i \in \mathcal{C}} P(c_i | d)$$

$$c_{NB} = \operatorname{argmax} \frac{P(d | c_i) P(c_i)}{P(d)}$$

$$c_{NB} = \operatorname{argmax} P(d | c_i) P(c_i)$$

$$c_{NB} = \operatorname{argmax} P(w_1, w_2, \dots, w_m | c_i) P(c_i)$$

**Document as  
words**



# Naïve Bayes on Text

- Given a document  $d$ , we want to find the most likely class  $c$ :

$$c_{NB} = \operatorname{argmax}_{c_i \in C} P(c_i | d)$$

$$c_{NB} = \operatorname{argmax} \frac{P(d | c_i) P(c_i)}{P(d)}$$

$$c_{NB} = \operatorname{argmax} P(d | c_i) P(c_i)$$

$$c_{NB} = \operatorname{argmax} P(w_1, w_2, \dots, w_m | c_i) P(c_i)$$

$$c_{NB} = \operatorname{argmax} P(w_1 | c_i) P(w_2 | c_i) \dots P(w_m | c_i) P(c_i)$$

**Conditional  
Independence**



# Naïve Bayes on Text

- Given a document  $d$ , we want to find the most likely class  $c$ :

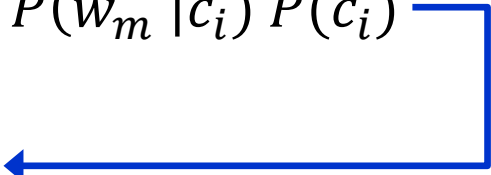
$$c_{NB} = \operatorname{argmax}_{c_i \in C} P(c_i | d)$$

$$c_{NB} = \operatorname{argmax} \frac{P(d | c_i) P(c_i)}{P(d)}$$

$$c_{NB} = \operatorname{argmax} P(d | c_i) P(c_i)$$

$$c_{NB} = \operatorname{argmax} P(w_1, w_2, \dots, w_m | c_i) P(c_i)$$

$$c_{NB} = \operatorname{argmax} P(w_1 | c_i) P(w_2 | c_i) \dots P(w_m | c_i) P(c_i)$$

$$c_{NB} = \operatorname{argmax}_{c_i \in C} P(c_i) \prod_{w_j \in W} P(w_j | c_i)$$




# Naïve Bayes on Text

---

- How do we estimate the different probabilities?

$$c_{NB} = \operatorname{argmax}_{c_i \in C} P(c_i) \prod_{w_j \in W} P(w_j \mid c_i)$$

- For a set of documents  $D$ , vocabulary of words  $W$ , and classes  $C$ ,

$$P(c_i) = \frac{|c_i|}{|D|}$$

$$P(w_1 \mid c_i) = \frac{\text{count}(w_1, c_i)}{\sum_{w_j \in W} \text{count}(w_j, c_i)}$$



# Naïve Bayes on Text

- How do we estimate the different probabilities?

$$c_{NB} = \operatorname{argmax}_{c_i \in C} P(c_i) \prod_{w_j \in W} P(w_j \mid c_i)$$

- For a set of documents  $D$ , vocabulary of words  $W$ , and classes  $C$ ,

$$P(c_i) = \frac{|c_i|}{|D|}$$

**Proportion of class  $c_i$  in the entire dataset**

$$P(w_1 \mid c_i) = \frac{\text{count}(w_1, c_i)}{\sum_{w_j \in W} \text{count}(w_j, c_i)}$$

**# times word  $w_1$  is used in all documents of class  $c_i$**

**Total words in all documents of class  $c_i$**





# Naïve Bayes on Text

- How do we use this to test if an email is a scam or real?

$$P(c_i) = \frac{|c_i|}{|D|}$$

$$P(w_1 | c_i) = \frac{\text{count}(w_1, c_i)}{\sum_{w_j \in W} \text{count}(w_j, c_i)}$$

ID	Email Text	Class
1	iphone, free, password	Scam
2	job, easy, password	Scam
3	lunch, restaurant, discount, iphone	Real
4	restaurant, reservation, discount	Real
5	iphone, discount, reservation	?

$$P(\text{Real} | \text{iphone, discount, reservation})$$

$$= P(\text{iphone, discount, reservation} | \text{Real}) P(\text{Real})$$

$$= P(\text{iphone} | \text{Real}) P(\text{discount} | \text{Real}) P(\text{reservation} | \text{Real}) P(\text{Real})$$

$$= (1/7) \times (2/7) \times (1/7) \times (2/4) = 0.0029$$



# Exercise: Naïve Bayes on Text

- How do we use this to test if an email is a scam or real?

$$P(c_i) = \frac{|c_i|}{|D|}$$

$$P(w_1 | c_i) = \frac{\text{count}(w_1, c_i)}{\sum_{w_j \in W} \text{count}(w_j, c_i)}$$

ID	Email Text	Class
1	iphone, free, password	Scam
2	job, easy, password	Scam
3	lunch, restaurant, discount, iphone	Real
4	restaurant, reservation, discount	Real
5	iphone, discount, reservation	?

$$P(\text{Scam} | \text{iphone, discount, reservation}) = ?$$

Which is the most likely class?



# Exercise: Naïve Bayes on Text

- How do we use this to test if an email is a scam or real?

$$P(c_i) = \frac{|c_i|}{|D|}$$

$$P(w_1 | c_i) = \frac{\text{count}(w_1, c_i)}{\sum_{w_j \in W} \text{count}(w_j, c_i)}$$

ID	Email Text	Class
1	iphone, free, password	Scam
2	job, easy, password	Scam
3	lunch, restaurant, discount, iphone	Real
4	restaurant, reservation, discount	Real
5	iphone, discount, reservation	?

Which is the most likely class?

$$P(\text{Scam} | \text{iphone, discount, reservation})$$

$$= P(\text{iphone, discount, reservation} | \text{Scam}) P(\text{Scam})$$

$$= P(\text{iphone} | \text{Scam}) P(\text{discount} | \text{Scam}) P(\text{reservation} | \text{Scam}) P(\text{Scam})$$

$$= (1/6) \times (0/6) \times (0/6) \times (2/4) = 0$$

$$P(\text{Real} | i, d, r) = 0.0029$$

$$P(\text{Scam} | i, d, r) = 0$$

**What is the issue here?**



# Laplace Smoothing

- Problem with new words, or unseen words for a specific class
  - E.g., Consider a new word  $w_1$ , which gives  $P(w_1 | c) = \frac{\text{count}(w_1, c_i)}{\sum \text{count}(w, c_i)} = 0$
  - The probability  $P(w_1 | c) = 0$  affects the entire equation

$$c_{NB} = \operatorname{argmax} P(c | d)$$
$$\vdots$$

$$c_{NB} = \operatorname{argmax} P(w_1 | c) P(w_2 | c) P(w_3 | c) P(w_4 | c) P(c)$$

Usually  $\alpha = 1$

- Laplace Smoothing: Add a constant  $\alpha$  to our counts

$$P(w_1 | c_i) = \frac{\text{count}(w_1, c_i) + \alpha}{\sum_{w_j \in W} (\text{count}(w_j, c_i) + \alpha)} = \frac{\text{count}(w_1, c_i) + \alpha}{\sum_{w_j \in W} (\text{count}(w_j, c_i)) + \alpha |W|}$$



# Exercise: Naïve Bayes on Text

- How do we use this to test if an email is a scam or real?

$$P(c_i) = \frac{|c_i|}{|D|}$$

*Laplace  
Smoothing*

$$P(w_1 | c_i) = \frac{\text{count}(w_1, c_i) + 1}{\sum_{w_j \in W} \text{count}(w_j, c_i) + |W|}$$

ID	Email Text	Class
1	iphone, free, password	Scam
2	job, easy, password	Scam
3	lunch, restaurant, discount, iphone	Real
4	restaurant, reservation, discount	Real
5	iphone, discount, reservation	?

$$P(\text{Real} | \text{iphone, discount, reservation}) = ?$$



# Exercise: Naïve Bayes on Text

- How do we use this to test if an email is a scam or real?

$$P(c_i) = \frac{|c_i|}{|D|}$$

*Laplace Smoothing*

$$P(w_1 | c_i) = \frac{\text{count}(w_1, c_i) + 1}{\sum_{w_j \in W} \text{count}(w_j, c_i) + |W|}$$

ID	Email Text	Class
1	iphone, free, password	Scam
2	job, easy, password	Scam
3	lunch, restaurant, discount, iphone	Real
4	restaurant, reservation, discount	Real
5	iphone, discount, reservation	?

$$P(\text{Real} | \text{iphone, discount, reservation})$$

$$= P(\text{iphone, discount, reservation} | \text{Real}) P(\text{Real})$$

$$= P(\text{iphone} | \text{Real}) P(\text{discount} | \text{Real}) P(\text{reservation} | \text{Real}) P(\text{Real})$$

$$= (1 + 1)/(7 + 9) \times (2 + 1)/(7 + 9) \times (1 + 1)/(7 + 9) \times (2/4) \approx 0.0015$$



# Exercise: Naïve Bayes on Text

- How do we use this to test if an email is a scam or real?

$$P(c_i) = \frac{|c_i|}{|D|}$$

*Laplace  
Smoothing*

$$P(w_1 | c_i) = \frac{\text{count}(w_1, c_i) + 1}{\sum_{w_j \in W} \text{count}(w_j, c_i) + |W|}$$

ID	Email Text	Class
1	iphone, free, password	Scam
2	job, easy, password	Scam
3	lunch, restaurant, discount, iphone	Real
4	restaurant, reservation, discount	Real
5	iphone, discount, reservation	?

$$P(\text{Scam} | \text{iphone, discount, reservation})$$
$$= ?$$



# Exercise: Naïve Bayes on Text

- How do we use this to test if an email is a scam or real?

$$P(c_i) = \frac{|c_i|}{|D|}$$

*Laplace Smoothing*

$$P(w_1 | c_i) = \frac{\text{count}(w_1, c_i) + 1}{\sum_{w_j \in W} \text{count}(w_j, c_i) + |W|}$$

ID	Email Text	Class
1	iphone, free, password	Scam
2	job, easy, password	Scam
3	lunch, restaurant, discount, iphone	Real
4	restaurant, reservation, discount	Real
5	iphone, discount, reservation	?

$$P(\text{Scam} | \text{iphone, discount, reservation})$$

$$= P(\text{iphone, discount, reservation} | \text{Scam}) P(\text{Scam})$$

$$= P(\text{iphone} | \text{Scam}) P(\text{discount} | \text{Scam}) P(\text{reservation} | \text{Scam}) P(\text{Scam})$$

$$= (1 + 1)/(6 + 9) \times (0 + 1)/(6 + 9) \times (0 + 1)/(6 + 9) \times (2/4) \approx 0.0003$$





# Naïve Bayes on Text

---

- Naïve Bayes is a decent baseline classifier but has its own limitations
- Assumes there is conditional independence (given a class)
  - $P(w_1, w_2, w_3, w_4 | c) = P(w_1 | c) P(w_2 | c) P(w_3 | c) P(w_4 | c)$
  - Is this always true?
- Assumes that order does not matter
  - “I like burger but dislike fruits” and “I like fruits but dislike burgers”
  - Does the NB classifier treat the above two sentences differently?



# Summary

---

- Recap on statistical concepts such as product rule, chain rule, conditional independence, Bayes rules
- Problem representation in terms of a Bayesian network and its corresponding conditional probability table
- Learn about how Bayes net can be used in various scenarios
- Learn about the Naïve Bayes Classifier and its application to text

