

# Laboratorio di Apprendimento Automatico

Fabio Aioli

Università di Padova

# Probabilistic Classifiers

- Compute  $P(c_j|d_i)$  by means of the Bayes' theorem
  - $P(c_j|d_i) = P(d_i|c_j)P(c_j)/P(d_i)$
  - Maximum a posteriori Hypothesis (MAP)  $\operatorname{argmax} P(c_j|d_i)$
- Classes are viewed as generators of documents
- The prior probability  $P(c_j)$  is the probability that a document  $d$  is in  $c_j$

# Naive Bayes Classifiers

Task: Classify a new instance  $D$  based on a tuple of attribute values  $D = \langle x_1, x_2, \dots, x_n \rangle$  into one of the classes  $c_j \in \mathcal{C}$

$$c_{MAP} = \operatorname{argmax}_{c_j \in \mathcal{C}} P(c_j | x_1, x_2, \dots, x_n)$$

$$= \operatorname{argmax}_{c_j \in \mathcal{C}} \frac{P(x_1, x_2, \dots, x_n | c_j) P(c_j)}{P(x_1, x_2, \dots, x_n)}$$

$$= \operatorname{argmax}_{c_j \in \mathcal{C}} P(x_1, x_2, \dots, x_n | c_j) P(c_j)$$

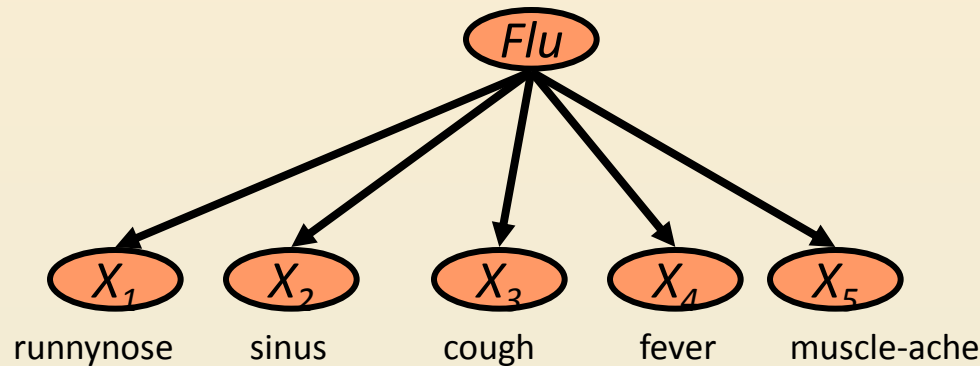
# Naive Bayes Classifier: Assumption

- $P(c_j)$ 
  - Can be estimated from the frequency of classes in the training examples.
- $P(x_1, x_2, \dots, x_n | c_j)$ 
  - $O(|X|^n / |C|)$  parameters
  - Could only be estimated if a very, very large number of training examples was available.

## Naive Bayes Conditional Independence Assumption:

- Assume that the probability of observing the conjunction of attributes is equal to the product of the individual probabilities  $P(x_i | c_j)$ .

# The Naïve Bayes Classifier

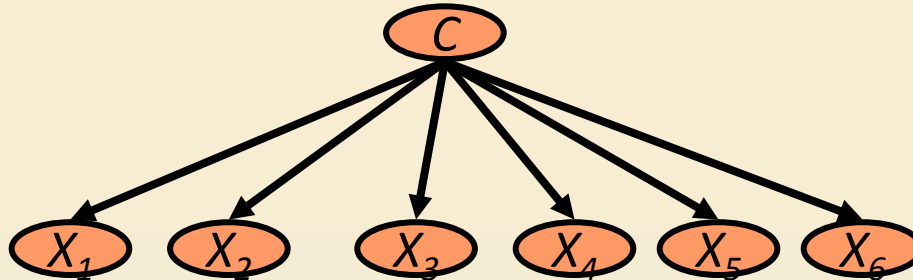


- **Conditional Independence Assumption:** features are independent of each other given the class:

$$P(X_1, \dots, X_5 | C) = P(X_1 | C) \cdot P(X_2 | C) \cdot \dots \cdot P(X_5 | C)$$

- Only  $n|C|$  parameters ( $+|C|$ ) to estimate

# Learning the Model



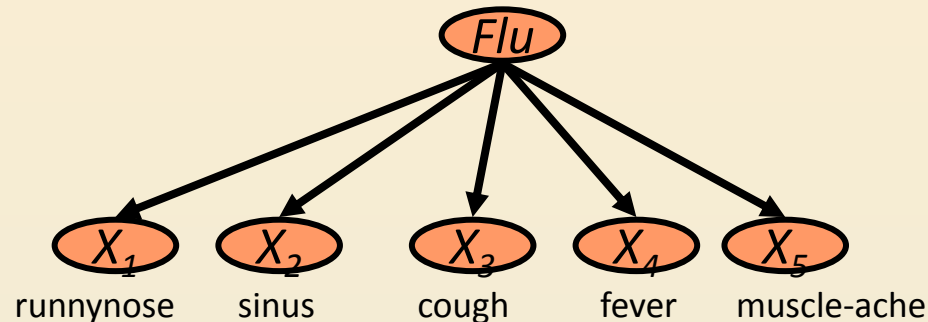
maximum likelihood estimates: most likely value of each parameter given the training data

– i.e. simply use the frequencies in the data

$$\hat{P}(c_j) = \frac{N(C = c_j)}{N}$$

$$\hat{P}(x_i | c_j) = \frac{N(X_i = x_i, C = c_j)}{N(C = c_j)}$$

# Problem with Max Likelihood



$$P(X_1, \dots, X_5 | C) = P(X_1 | C) \cdot P(X_2 | C) \cdot \dots \cdot P(X_5 | C)$$

- What if we have seen no training cases where patient had no flu and muscle aches?

$$\hat{P}(X_5 = t | C = nf) = \frac{N(X_5 = t, C = nf)}{N(C = nf)} = 0$$

- Zero probabilities cannot be conditioned away, no matter the other evidence!

$$\ell = \arg \max_c \hat{P}(c) \prod_i \hat{P}(x_i | c)$$

# Smoothing to Avoid Overfitting

$$\hat{P}(x_i | c_j) = \frac{N(X_i = x_i, C = c_j) + 1}{N(C = c_j) + k}$$

# of values of  $X_i$