

# DS 4400

## Machine Learning and Data Mining I Spring 2024

David Liu

Khoury College of Computer Science  
Northeastern University

February 27 2024

# Announcements

- Final Project
  - Project proposals due this Friday on Gradescope
    - Each member should submit a copy.
  - Expect further details on final project report (6-8 pages) and video (~4 minutes) after spring break.
- Homework 3
  - Released on February 26, due on March 18
- Midterm grades will be released after Spring Break.

# Where We Are

Regression

Simple, Multiple,  
Polynomial

Classification

Logistic Regression,  
kNN

Generative  
Models

LDA, Naïve Bayes

Ensemble  
Learning

Bagging, Boosting

Neural  
Networks

Multi-layer Perceptron,  
Convolutional Neural  
Networks

Tools

Optimization:

- Closed form solution
- Gradient descent
- Backpropagation

Preventing Overfitting:

- Regularization
- Cross Validation

Evaluation of ML:

- Precision, Recall, AUC

# Outline

- Generative classifiers
- Naïve Bayes classifiers [focus for today]
  - Naïve Bayes assumption
  - Laplace smoothing
  - Comparison to LDA
- Decision trees [preview]
  - Information gain / entropy measures
  - Training algorithm

# Recap: ML Concepts

- Supervised vs unsupervised learning
- Classification vs regression
- Linear vs non-linear classifiers
- Generative vs discriminative classifiers
- Loss functions
- Metrics for evaluation

# Generative vs Discriminative

- **Generative model**
  - Given X and Y, learns the joint probability  $P(X, Y)$
  - Can generate more examples from distribution
  - Examples: LDA, Naïve Bayes, language models (GPT-2, GPT-3, BERT)
- **Discriminative model**
  - Given X and Y, learns a decision function for classification

Discriminative  $P(Y|X) \approx P(\text{Cancer} | \text{age, bp, income})$   
Posterior

Generative :  $P(X, Y) \rightarrow P(Y|X) = \frac{P(X, Y)}{P(X)} = \frac{P(X|Y)P(Y)}{P(X)}$

$$P(X \cap Y)$$

Generative Modeling is harder Prior

→ Make Assumptions

LDA →  $P(X|Y)$  is Gaussian

# Generative classifiers based on Bayes Theorem

$$P(A | B) = \frac{P(B | A) \times P(A)}{P(B)}$$

- Exactly the process we just used
- The most important formula in probabilistic machine learning

(Super Easy) Derivation:

$$P(A \wedge B) = P(A | B) \times P(B)$$

$$P(B \wedge A) = P(B | A) \times P(A)$$

these are the same

Just set equal...

$$P(A | B) \times P(B) = P(B | A) \times P(A)$$

and solve...



Bayes, Thomas (1763) An essay towards solving a problem in the doctrine of chances. *Philosophical Transactions of the Royal Society of London*, 53:370-418

# LDA Training and Testing

Given training data  $(x_i, y_i), i = 1, \dots, n, y_i \in \{1, \dots, K\}$

1. Estimate mean  
and variance

$$\begin{aligned}\hat{\mu}_k &= \frac{1}{n_k} \sum_{i:y_i=k} x_i \\ \hat{\sigma}^2 &= \frac{1}{n - K} \sum_{k=1}^K \sum_{i:y_i=k} (x_i - \hat{\mu}_k)^2\end{aligned}$$

2. Estimate prior

$$\hat{\pi}_k = n_k / n.$$

Given testing point  $x$ , predict  $k$  that maximizes:

$$\hat{\delta}_k(x) = x \cdot \frac{\hat{\mu}_k}{\hat{\sigma}^2} - \frac{\hat{\mu}_k^2}{2\hat{\sigma}^2} + \log(\hat{\pi}_k)$$

# Naïve Bayes Classifier

**Idea:** Use the training data to estimate

$$P(X | Y) \text{ and } P(Y) .$$

Then, use Bayes rule to infer  $P(Y|X_{\text{new}})$  for new data

---

$$P[Y = k | X = x] = \frac{P[Y = k] P[X_1 = x_1 \wedge \dots \wedge X_d = x_d | Y = k]}{P[X_1 = x_1 \wedge \dots \wedge X_d = x_d]}$$

Easy to estimate  
from data      Impractical, but necessary

Unnecessary, as it turns out

# Learning Joint Distributions

## Step 1:

Build a JD table for your attributes in which the probabilities are unspecified

A	B	C	Prob
0	0	0	?
0	0	1	?
0	1	0	?
0	1	1	?
1	0	0	?
1	0	1	?
1	1	0	?
1	1	1	?

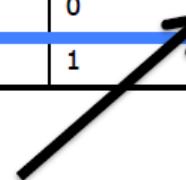
## Step 2:

Then, fill in each row with:

$$\hat{P}(\text{row}) = \frac{\text{records matching row}}{\text{total number of records}}$$

A	B	C	Prob
0	0	0	0.30
0	0	1	0.05
0	1	0	0.10
0	1	1	0.05
1	0	0	0.05
1	0	1	0.10
1	1	0	0.25
1	1	1	0.10

Fraction of all records in which  
A and B are true but C is false



# Example – Learning Joint Probability Distribution

This Joint PD was obtained by learning from three attributes in the UCI “Adult” Census Database [Kohavi 1995]



# Naïve Bayes Classifier

**Problem:** estimating the joint density isn't practical

However, if we make the assumption that the attributes are independent given the class label, estimation is easy!

# Naïve Bayes Classifier

**Problem:** estimating the joint PD or CPD isn't practical

- Severely overfits, as we saw before

However, if we make the assumption that the attributes are independent given the class label, estimation is easy!

$$P[X_1 = x_1 \wedge \dots \wedge X_d = x_d | Y = k] = \prod_{j=1}^d P[X_j = x_j | Y = k]$$

- In other words, we assume all attributes are *conditionally independent* given  $Y$
- Often this assumption is violated in practice, but more on that later...

# Using the Naïve Bayes Classifier

- Now, we have

$$P[Y = k | X = x] = \frac{P[Y = k]P[X_1 = x_1 \wedge \dots \wedge X_d = x_d | Y = k]}{P[X_1 = x_1 \wedge \dots \wedge X_d = x_d]}$$

This is constant for a given instance,  
and so irrelevant to our prediction

# Using the Naïve Bayes Classifier

- Now, we have

$$P[Y = k | X = x] = \frac{P[Y = k] P[X_1 = x_1 \wedge \dots \wedge X_d = x_d | Y = k]}{P[X_1 = x_1 \wedge \dots \wedge X_d = x_d]}$$

This is constant for a given instance, and so irrelevant to our prediction

- In practice, we use log-probabilities to prevent underflow
  - To classify a new point  $\mathbf{x}$ ,
- $$h(\mathbf{x}) = \arg \max_{y_k} P(Y = k) \prod_{j=1}^d P(X_j = x_j | Y = k)$$
- j<sup>th</sup> attribute value of  $\mathbf{x}$*

# Naïve Bayes Classifier

## TRAIN

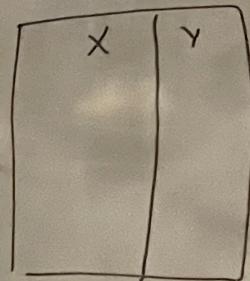
- For each class label  $k$ 
  1. Estimate prior  $\pi_k = P[Y = k]$  from the data
  2. For each value  $v$  of attribute  $X_j$ 
    - Estimate  $P[X_j = v | Y = k]$

## TEST on INPUT $x = (x_1, \dots, x_d)$

- For every  $k$ , compute the probabilities
  - $p_k = P[Y = k] \prod_{j=1}^d P[X_j = x_j | Y = k]$
  - Classify  $x$  to the class  $k$  that maximizes  $p_k$

Assume

$P(x, y) \rightarrow \text{Data} \rightarrow \text{Choose distribution to Model} \rightarrow \text{MLE}$



$$\hat{P}(x, y)$$

LDA

$$[\mu, \omega, \pi] \quad P(y|x)$$

One-hot

[ 0 0 0 1 0 0 ]  
Dog Cat fish bat

$$\left[ \pi_k, P(x_i | Y=Y_k) \right] \text{ NAIVE BAYES}$$

given

# Training Naïve Bayes

Estimate  $P[X_j = x_j | Y = k]$  and  $P[Y = k]$  directly from the training data by counting!

<u>Sky</u>	<u>Temp</u>	<u>Humid</u>
sunny	warm	normal
sunny	cold	high
rainy	cold	high
sunny	warm	high

<u>Play?</u>
yes
yes
no
yes

Prior:  $P(\text{Play} = \text{Yes}) =$

$P(\text{Play} = \text{No}) =$

Conditional feature distributions

$P(\text{Sky} = \text{sunny} | \text{Play} = \text{Yes}) =$

$P(\text{Sky} = \text{sunny} | \text{Play} = \text{No}) =$

$P(\text{Temp} = \text{warm} | \text{Play} = \text{Yes}) =$

$P(\text{Temp} = \text{warm} | \text{Play} = \text{No}) =$

$P(\text{Humid} = \text{high} | \text{Play} = \text{Yes}) =$

$P(\text{Humid} = \text{high} | \text{Play} = \text{No}) =$

# Training Naïve Bayes

Estimate  $P[X_j = x_j | Y = k]$  and  $P[Y = k]$  directly from the training data by counting!

<u>Sky</u>	<u>Temp</u>	<u>Humid</u>
sunny	warm	normal
sunny	cold	high
rainy	cold	high
sunny	warm	high

<u>Play?</u>
yes
yes
no
yes

Classify new point:

X: Sky = sunny, Temp = cold, Humid = high

Priors ( $\pi_k$ )

$$\left. \begin{array}{l} P(\text{Play} = \text{Yes}) = \frac{3}{4} \\ P(\text{Play} = \text{No}) = \frac{1}{4} \end{array} \right\} \text{Joint}$$

$$\nearrow \frac{1}{1+2} = \frac{1}{3}$$

Prediction

New day: Sunny, cold,  
(x) low humidity

$$P(\text{Sky} = \text{Sunny} \mid \text{Play} = \text{Yes}) = \frac{3}{3}$$

$$P(\text{Sky} = \text{Sunny} \mid \text{Play} = \text{No}) = \frac{0}{1}$$

$$P(\text{Play} = \text{Yes} \mid x)$$

$$\boxed{P(\text{Temp} = \text{Warm} \mid \text{Play} = \text{Yes}) = \frac{2}{3}}$$

$$\boxed{P(\text{Temp} = \text{Warm} \mid \text{Play} = \text{No}) = \frac{0}{1}}$$

Play = Yes

$$(\frac{3}{4})(1)(\frac{1}{3})(\frac{1}{3}) = \frac{1}{12}$$

$$P(\text{Humid} = \text{high} \mid \text{Play} = \text{Yes}) = \frac{2}{3}$$

$$P(\text{Humid} = \text{high} \mid \text{Play} = \text{No}) = \frac{1}{1}$$

Play = No

$$(\frac{1}{4}) \cdot \underline{0} \cdot \underline{1} \cdot \underline{0} = 0$$

# Laplace Smoothing

- Notice that some probabilities estimated by counting might be zero
  - Possible overfitting!

# Laplace Smoothing

- Notice that some probabilities estimated by counting might be zero
  - Possible overfitting!
- Fix by using Laplace smoothing:

- Adds 1 to each count

$$P(X_j = v \mid Y = k) = \frac{c_v + 1}{\sum_{v' \in \text{values}(X_j)} c_{v'} + |\text{values}(X_j)|}$$

where

- $c_v$  is the count of training instances with a value of  $v$  for attribute  $j$  and class label  $k$
  - $|\text{values}(X_j)|$  is the number of values  $X_j$  can take on

# Naïve Bayes Classifier

## TRAIN

- For each class label  $k$ 
  1. Estimate prior  $\pi_k = P[Y = k]$  from the data
  2. For each value  $v$  of attribute  $X_j$ 
    - Estimate  $P[X_j = v | Y = k]$  with Laplace smoothing

## TEST on INPUT $x = (x_1, \dots, x_d)$

- For every  $k$ , compute the probabilities
  - $P[Y = k] \prod_{j=1}^d P[X_j = x_j | Y = k]$
  - Classify  $x$  to the class  $k$  that maximizes the above product

# Continuous Features

- Naïve Bayes can be extended to continuous features
- Gaussian Naïve Bayes
  - Here an additional assumption is that each distribution  $P[X_j|Y = k]$  is Gaussian  $N(\mu_j, \sigma_j)$
  - It estimates the mean and standard deviation from training data
- This leads to a linear classifier

# Comparison to LDA

- **Similarity to LDA**
  - Both are generative models
  - They both estimate:
$$P[X = x \text{ and } Y = k] = P[X = x|Y = k]P[Y = k]$$
Using Bayes Theorem
- **Difference from LDA**
  - Naïve Bayes can handle discrete data
  - LDA uses multi-variate normal
  - LDA assumes same variances for all classes
  - Naïve Bayes make the conditional independence assumption
  - LDA is linear, while Naïve Bayes is usually not linear

# Naïve Bayes Summary

## Advantages:

- Fast to train (single scan through data)
- Fast to classify
- Not sensitive to irrelevant features
- Handles real and discrete data
- Handles streaming data well

## Disadvantages:

- Assumes independence of features