# Naive Bayes Classifiers

**Dr. Rahul Kottath**

# Background

- There are three methods to establish a classifier
    a) Model a classification rule directly
        - Examples: k-NN, decision trees, perceptron, SVM
    b) Model the probability of class memberships given input data
        - Example: multi-layered perceptron with the cross-entropy cost
    c) Make a probabilistic model of data within each class
        - Examples: naive Bayes, model based classifiers
- *a*) and *b*) are examples of discriminative classification
- *c*) is an example of generative classification
- *b*) and *c*) are both examples of probabilistic classification

# Things We'd Like to Do

- Spam Classification
  - Given an email, predict whether it is spam or not

- Medical Diagnosis
  - Given a list of symptoms, predict whether a patient has disease X or not

- Weather
  - Based on temperature, humidity, etc... predict if it will rain tomorrow

# Bayesian Classification

$$X = (x_1, x_2, x_3, x_4 \cdots x_n)$$

$-\textcircled{Y}$

- Problem statement:
    - Given features $X_1, X_2, \ldots, X_n$
    - Predict a label Y

# Another Application

**Digit Recognition**



$X_1, \ldots, X_n \in \{0,1\}$ (Black vs. White pixels)

$Y \in \{5,6\}$ (predict whether a digit is a 5 or a 6)

# The Bayes Classifier

- A good strategy is to predict:

$$\arg\max_Y P(Y|\underbrace{X_1,\ldots,X_n}_{\text{pixels}})$$

  - ➢ (for example: what is the probability that the image represents a 5 given its pixels?)

- So … How do we compute that?

# The Bayes Classifier

$$P(A/B) = \frac{P(B/A) \cdot P(A)}{P(B)}$$

- Use Bayes Rule!

Likelihood          Prior

$$P(Y|X_1,\ldots,X_n) = \frac{P(X_1,\ldots,X_n|Y)P(Y)}{P(X_1,\ldots,X_n)}$$

$B$

Normalization Constant

- Why did this help?  Well, we think that we might be able to specify how features are "generated" by the class label

# The Bayes Classifier

$$p + 1-p = 1$$

- Let's expand this for our digit recognition task:

$$\frac{a}{a+b} + \frac{b}{a+b} = \frac{a+b}{a+b} = 1$$

$$P(Y = 5|X_1, \ldots, X_n) = \frac{P(X_1, \ldots, X_n|Y = 5)P(Y = 5)}{P(X_1, \ldots, X_n|Y = 5)P(Y = 5) + P(X_1, \ldots, X_n|Y = 6)P(Y = 6)}$$

$$P(Y = 6|X_1, \ldots, X_n) = \frac{P(X_1, \ldots, X_n|Y = 6)P(Y = 6)}{P(X_1, \ldots, X_n|Y = 5)P(Y = 5) + P(X_1, \ldots, X_n|Y = 6)P(Y = 6)}$$

- To classify, we'll simply compute these two probabilities and predict based on which one is greater

# Probability Basics

- Prior, conditional and joint probability

  - Prior probability: $P(X)$

  - Conditional probability: $P(X_1 \mid X_2), P(X_2 \mid X_1)$

  - Joint probability: $\mathbf{X} = (X_1, X_2), P(\mathbf{X}) = P(X_1, X_2)$

  - Relationship: $P(X_1, X_2) = P(X_2 \mid X_1)P(X_1) = P(X_1 \mid X_2)P(X_2)$

  - Independence: $P(X_2 \mid X_1) = P(X_2), P(X_1 \mid X_2) = P(X_1), P(X_1, X_2) = P(X_1)P(X_2)$

- Bayesian Rule

$$P(C \mid \mathbf{X}) = \frac{P(\mathbf{X} \mid C)P(C)}{P(\mathbf{X})} \qquad Posterior = \frac{Likelihood \times Prior}{Evidence}$$

# Naïve Bayes

$P(X_1/x_2,x_3\cdots x_n,c) \cdot P(x_2,x_3\cdots x_n,c)$

$P(x_1/c) \cdot P(x_2/c) \cdot P(x_3/c) \cdots P(x_n/c) \cdot P(c)$

- Bayes classification

$$P(C \mid \mathbf{X}) \propto P(\mathbf{X} \mid C)P(C) = P(X_1, \cdots, X_n \mid C)P(C)$$

  Difficulty: learning the joint probability

- Naïve Bayes classification

  $P(y_c)_+ = P(x_1/c=y_c) \cdot P(x_2/c=y_c) \cdot P(x_4/c=y_c) \cdots P(x_n/c=y_c) \cdot P(c,y_c)$

  – Making the assumption that (all input attributes are independent)

$$P(X_1, X_2, \cdots, X_n \mid C) = P(X_1 \mid X_2, \cdots, X_n ; C)P(X_2, \cdots, X_n \mid C)$$
$$= P(X_1 \mid C)P(X_2, \cdots, X_n \mid C)$$
$$= P(X_1 \mid C)P(X_2 \mid C) \cdots P(X_n \mid C)$$

$P(C/x) \qquad = \prod_{i=1}^{n} P(x_i/c) \cdot P(c)$

# Example

## Example: Play Tennis

$P(Yes) = 9/14$

$P(No) = 5/14$

*PlayTennis*: training examples

| Day | Outlook | Temperature | Humidity | Wind | PlayTennis |
|-----|---------|-------------|----------|------|------------|
| D1 | Sunny | Hot | High | Weak | No |
| D2 | Sunny | Hot | High | Strong | No |
| D3 | Overcast | Hot | High | Weak | Yes |
| D4 | Rain | Mild | High | Weak | Yes |
| D5 | Rain | Cool | Normal | Weak | Yes |
| D6 | Rain | Cool | Normal | Strong | No |
| D7 | Overcast | Cool | Normal | Strong | Yes |
| D8 | Sunny | Mild | High | Weak | No |
| D9 | Sunny | Cool | Normal | Weak | Yes |
| D10 | Rain | Mild | Normal | Weak | Yes |
| D11 | Sunny | Mild | Normal | Strong | Yes |
| D12 | Overcast | Mild | High | Strong | Yes |
| D13 | Overcast | Hot | Normal | Weak | Yes |
| D14 | Rain | Mild | High | Strong | No |

$x_1$  $x_2$  $x_3$  $x_4$

4 YES
0 NO

→ 9 YES

5 NO

D15

# Example

2 (YES)   3 (NO)

| Outlook | Play=*Yes* | Play=*No* |
|---------|------------|-----------|
| *Sunny* | 2/9 | 3/5 |
| *Overcast* | 4/9 | 0/5 |
| *Rain* | 3/9 | 2/5 |

| Temperature | Play=*Yes* | Play=*No* |
|-------------|------------|-----------|
| *Hot* | 2/9 | 2/5 |
| *Mild* | 4/9 | 2/5 |
| *Cool* | 3/9 | 1/5 |

# Example

| Humidity | Play=*Yes* | Play=N*o* |
|----------|:----------:|:---------:|
| *High* | 3/9 | 4/5 |
| *Normal* | 6/9 | 1/5 |

| Wind | Play=*Yes* | Play=N*o* |
|------|:----------:|:---------:|
| *Strong* | 3/9 | 3/5 |
| *Weak* | 6/9 | 2/5 |

*P*(Play=*Yes)* = 9/14

*P*(Play=*No)* = 5/14

# Example

$$P(Yes/x') = P(Outlook = Sunny/Yes) \cdot P(Temp = Cool/Yes) \cdot P(Hum = High/Yes) \cdot P(Wind = Strong/Yes) \cdot P(Yes)$$

$$P(No/x') =$$

- Test Phase

  play = ?

  - Given a new instance,

    x'=(Outlook=*Sunny*, Temperature=*Cool*, Humidity=*High*, Wind=*Strong*)

  - Look up tables

    P(Outlook=*Sunny*|Play=*Yes*) = 2/9          P(Outlook=S*unny*|Play=*No*) = 3/5

    P(Temperature=*Cool*|Play=*Yes*) = 3/9       P(Temperature=*Cool*|Play==*No*) = 1/5

    P(Huminity=*High*|Play=*Yes*) = 3/9          P(Huminity=*High*|Play=*No*) = 4/5

    P(Wind=*Strong*|Play=*Yes*) = 3/9            P(Wind=*Strong*|Play=*No*) = 3/5

    P(Play=*Yes*) = 9/14                         P(Play=*No*) = 5/14

    Yes                                          No

$$P(Yes/x')$$

# Example

## MAP rule

P($Yes$ | **x′**): [P($Sunny$ | Yes)P($Cool$ | Yes)P($High$ | Yes)P($Strong$ | Yes)]P(Play=$Yes$) = 0.0053

P($No$ | **x′**): [P($Sunny$ | No) P($Cool$ | No)P($High$ | No)P($Strong$ | No)]P(Play=$No$) = 0.0206

Given the fact P($Yes$ | **x′**) < P($No$ | **x′**), we label **x′** to be "$No$".

$$P(Yes/x') = \frac{0.0053}{0.0053 + 0.0206}$$

$$P(No/x') = \frac{0.0206}{0.0053 + 0.0206}$$

$$P(Yes/x') + P(No/x') = 1$$

# Conclusions

- Naïve Bayes based on the independence assumption
  - Training is very easy and fast; just requiring considering each attribute in each class separately
  - Test is straightforward; just looking up tables or calculating conditional probabilities with normal distributions
- A popular generative model
  - Performance competitive to most of state-of-the-art classifiers even in presence of violating independence assumption
  - Many successful applications, e.g., spam mail filtering
  - A good candidate of a base learner in ensemble learning
  - Apart from classification, naïve Bayes can do more…

# Evaluating a Classification model:
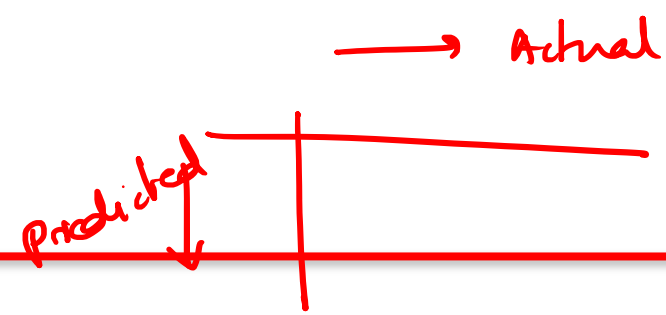
## 1. Log Loss or Cross-Entropy Loss:

- It is used for evaluating the performance of a classifier, whose output is a probability value between the 0 and 1.

- For a good binary Classification model, the value of log loss should be near to 0.

- The value of log loss increases if the predicted value deviates from the actual value.

- The lower log loss represents the higher accuracy of the model.

- For Binary classification, cross-entropy can be calculated as:

$$Loss = -(y\log(p)+(1-y)\log(1-p))$$

# Evaluating a Classification model:

_Actual_

_Predicted_

## 2. Confusion Matrix:

- The confusion matrix provides us a matrix/table as output and describes the performance of the model.

- It is also known as the error matrix.

- The matrix consists of predictions result in a summarized form, which has a total number of correct predictions and incorrect predictions. The matrix looks like as below table:

_Act_

|  | **Actual Positive** | **Actual Negative** |
|---|---|---|
| Predicted Positive | True Positive | False Positive |
| Predicted Negative | False Negative | True Negative |

$$Accuracy = \frac{TP+TN}{Total\ Population}$$

# Example

example confusion matrix for a binary classifier

| n=165 | Predicted: NO | Predicted: YES |
|---|---|---|
| Actual: NO | 50 | 10 |
| Actual: YES | 5 | 100 |

$$\frac{150}{165}$$

# Example

|        | Predicted: NO | Predicted: YES |     |
|--------|---------------|----------------|-----|
| n=165  |               |                |     |
| Actual: NO  | TN = 50  | FP = 10   | 60  |
| Actual: YES | FN = 5   | TP = 100  | 105 |
|        | 55            | 110            |     |

# Example

- Accuracy: Overall, how often is the classifier correct?
  - (TP+TN)/total = (100+50)/165 = 0.91

- Misclassification Rate: Overall, how often is it wrong?
  - (FP+FN)/total = (10+5)/165 = 0.09
  - equivalent to 1 minus Accuracy
  - also known as "Error Rate"

- True Positive Rate: When it's actually yes, how often does it predict yes?
  - TP/actual yes = 100/105 = 0.95
  - also known as "Sensitivity" or "Recall"

# Example

- False Positive Rate: When it's actually no, how often does it predict yes?

  - FP/actual no = 10/60 = 0.17

- True Negative Rate: When it's actually no, how often does it predict no?

  - TN/actual no = 50/60 = 0.83

  - equivalent to 1 minus False Positive Rate

  - also known as "Specificity"

  F1 score

- Precision: When it predicts yes, how often is it correct?

  - TP/predicted yes = 100/110 = 0.91

# Python Packages needed

- pandas
  - Data Analytics
- numpy
  - Numerical Computing
- matplotlib.pyplot
  - Plotting graphs
- sklearn
  - Classification and Regression Classes

# Implementation Using sklearn

Let's go to Jupyter Notebook!