# Lecture 6: Bayes Classification

## Richa Singh

# Recap: Bayes' Classification

- Posterior, likelihood, prior, evidence

$$P(\omega_j|x) = \frac{p(x|\omega_j)P(\omega_j)}{p(x)},$$

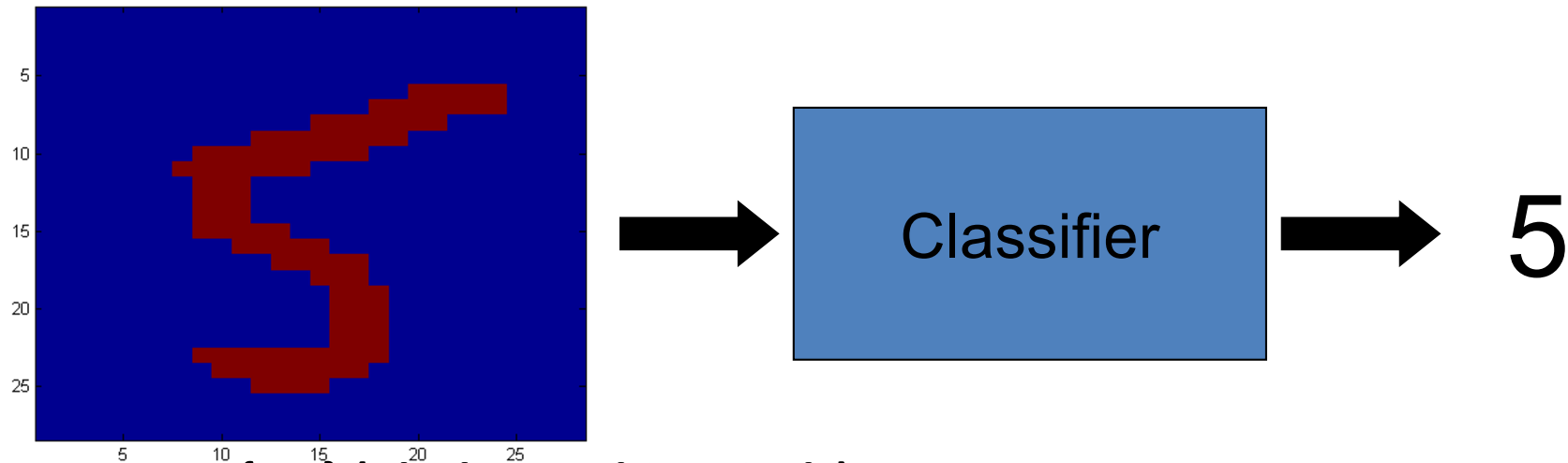— Evidence: In case of two categories

$$p(x) = \sum_{j=1}^{2} p(x|\omega_j)P(\omega_j)$$

$$posterior = \frac{likelihood \times prior}{evidence}$$

# 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:

$$\arg\max_Y P(Y|X_1, \ldots, X_n)$$

  – (for example: what is the probability that the image represents a 5 given its pixels?)
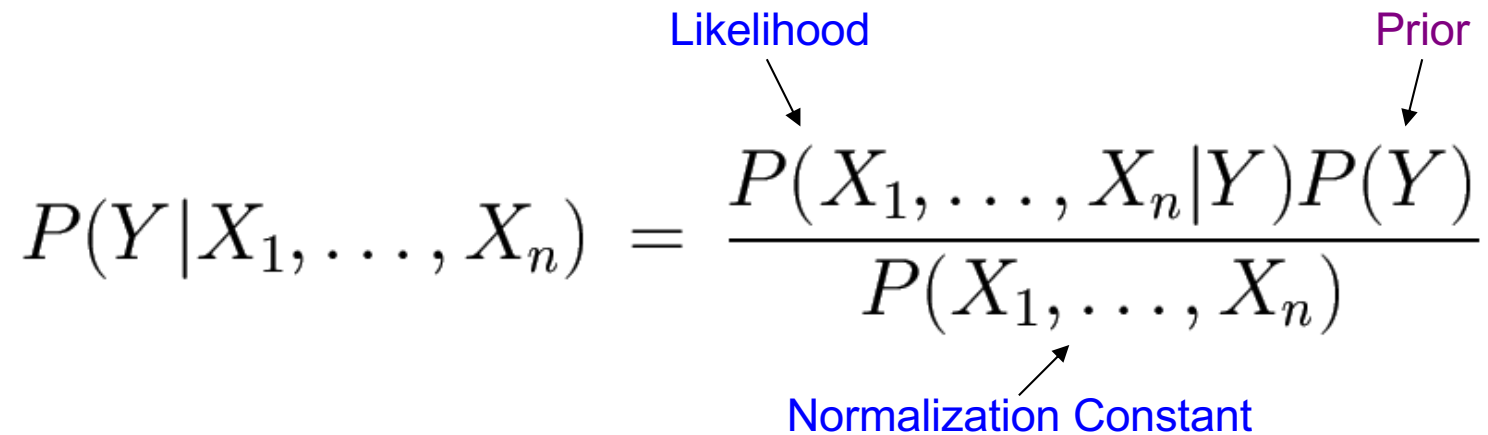
- So … How do we compute that?

# The Bayes Classifier

- Use Bayes Rule!

Likelihood             Prior

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

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

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

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

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

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

# Model Parameters

- For the Bayes classifier, we need to "learn" two functions, the likelihood and the prior

- How many parameters are required to specify the prior for our digit recognition example?

# Model Parameters

- How many parameters are required to specify the likelihood?
  - (Supposing that each image is 30x30 pixels)

?

# Model Parameters

- The problem with explicitly modeling $P(X_1,…,X_n|Y)$ is that there are usually way too many parameters:

  - We'll run out of space

  - We'll run out of time

  - And we'll need tons of training data (which is usually not available)

# The Naïve Bayes Model

- The *Naïve Bayes Assumption*: Assume that all features are independent **given the class label Y**
- Equationally speaking:
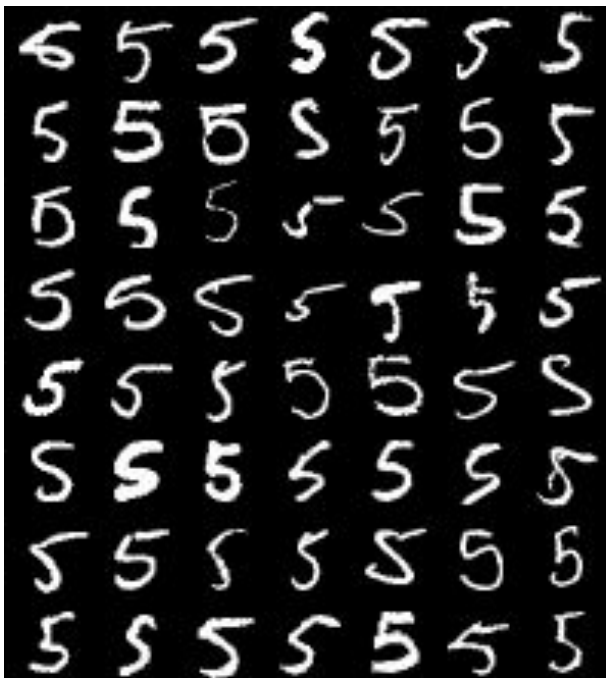
$$P(X_1, \ldots, X_n | Y) = \prod_{i=1}^{n} P(X_i | Y)$$

- (We will discuss the validity of this assumption later)

# Why is this useful?

- # of parameters for modeling $P(X_1,\ldots,X_n|Y)$:
- - Given each x_i is a binary attribute and y is boolean
  - $2(2^n-1)$


- # of parameters for modeling $P(X_1|Y),\ldots,P(X_n|Y)$

  - $2n$

# Naïve Bayes Training

- Now that we've decided to use a Naïve Bayes classifier, we need to train it with some data:



MNIST Training Data

# Naïve Bayes Training

- Training in Naïve Bayes is **easy**:
  - Estimate P(Y=v) as the fraction of records with Y=v

$$P(Y = v) = \frac{Count(Y = v)}{\# \ records}$$

  - Estimate P(X$_i$=u|Y=v) as the fraction of records with Y=v for which X$_i$=u

$$P(X_i = u | Y = v) = \frac{Count(X_i = u \wedge Y = v)}{Count(Y = v)}$$

- (This corresponds to Maximum Likelihood estimation of model parameters)

# Naïve Bayes Training

- In practice, some of these counts can be zero

- Fix this by adding "virtual" counts:

$$P(X_i = u | Y = v) = \frac{Count(X_i = u \wedge Y = v) + 1}{Count(Y = v) + 2}$$

  - (This is like putting a prior on parameters and doing MAP estimation instead of MLE)
  - This is called *Smoothing*

# Naïve Bayes Training

- For binary digits, training amounts to averaging all of the training fives together and all of the training sixes together.

# Naïve Bayes Classification

# Another Example of the Naïve Bayes Classifier

**The weather data, with counts and probabilities**

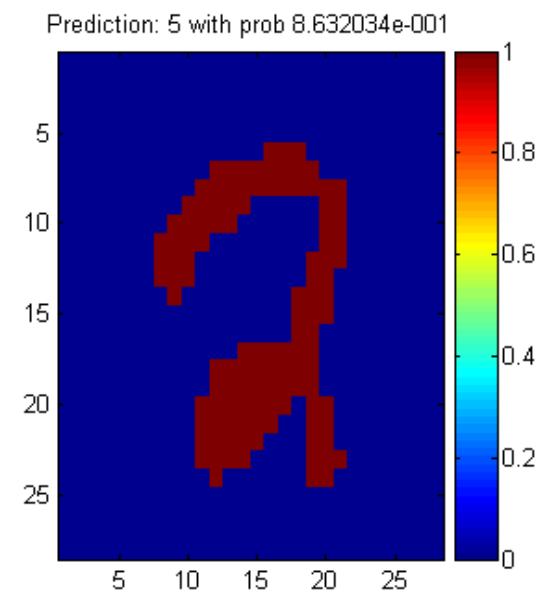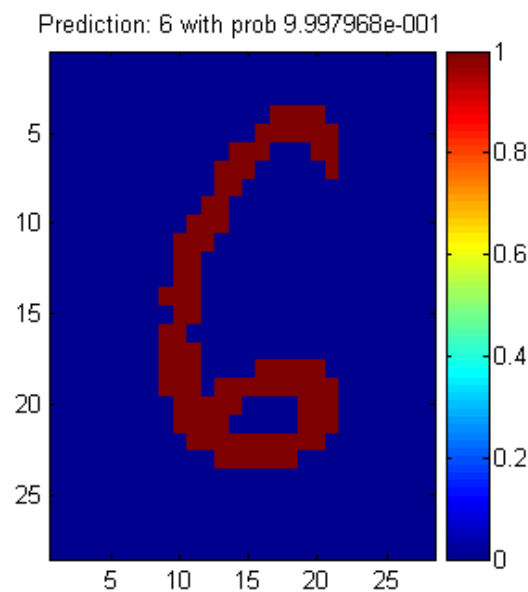| outlook | yes | no | temperature | yes | no | humidity | yes | no | windy | yes | no | play | yes | no |
|---------|-----|----|-------------|-----|----|----------|-----|----|-------|-----|----|------|-----|----|
| sunny | 2 | 3 | hot | 2 | 2 | high | 3 | 4 | false | 6 | 2 | | 9 | 5 |
| overcast | 4 | 0 | mild | 4 | 2 | normal | 6 | 1 | true | 3 | 3 | | | |
| rainy | 3 | 2 | cool | 3 | 1 | | | | | | | | | |
| sunny | 2/9 | 3/5 | hot | 2/9 | 2/5 | high | 3/9 | 4/5 | false | 6/9 | 2/5 | | 9/14 | 5/14 |
| overcast | 4/9 | 0/5 | mild | 4/9 | 2/5 | normal | 6/9 | 1/5 | true | 3/9 | 3/5 | | | |
| rainy | 3/9 | 2/5 | cool | 3/9 | 1/5 | | | | | | | | | |

**A new day**

| outlook | temperature | humidity | windy | play |
|---------|-------------|----------|-------|------|
| sunny | cool | high | true | ? |

- Likelihood of yes

$$= \frac{2}{9} \times \frac{3}{9} \times \frac{3}{9} \times \frac{3}{9} \times \frac{9}{14} = 0.0053$$

- Likelihood of no

$$= \frac{3}{5} \times \frac{1}{5} \times \frac{4}{5} \times \frac{3}{5} \times \frac{5}{14} = 0.0206$$

- Therefore, the prediction is No

# The Naive Bayes Classifier for Data Sets with Numerical Attribute Values

- One common practice to handle numerical attribute values is to assume normal distributions for numerical attributes.

**The numeric weather data with summary statistics**

| outlook | yes | no | temperature | yes | no | humidity | yes | no | windy | yes | no | play | yes | no |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| sunny | 2 | 3 | | 83 | 85 | | 86 | 85 | false | 6 | 2 | | 9 | 5 |
| overcast | 4 | 0 | | 70 | 80 | | 96 | 90 | true | 3 | 3 | | | |
| rainy | 3 | 2 | | 68 | 65 | | 80 | 70 | | | | | | |
| | | | | 64 | 72 | | 65 | 95 | | | | | | |
| | | | | 69 | 71 | | 70 | 91 | | | | | | |
| | | | | 75 | | | 80 | | | | | | | |
| | | | | 75 | | | 70 | | | | | | | |
| | | | | 72 | | | 90 | | | | | | | |
| | | | | 81 | | | 75 | | | | | | | |
| sunny | 2/9 | 3/5 | mean | 73 | 74.6 | mean | 79.1 | 86.2 | false | 6/9 | 2/5 | | 9/14 | 5/14 |
| overcast | 4/9 | 0/5 | std dev | 6.2 | 7.9 | std dev | 10.2 | 9.7 | true | 3/9 | 3/5 | | | |
| rainy | 3/9 | 2/5 | | | | | | | | | | | | |

- Let $x_1$, $x_2$, …, $x_n$ be the values of a numerical attribute in the training data set.

$$\mu = \frac{1}{n}\sum_{i=1}^{n} x_i$$

$$\sigma = \frac{1}{n-1}\sum_{i=1}^{n}(x_i - \mu)^2$$

$$f(w) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(w-\mu)^2}{\sigma^2}}$$

- For examples,

$$f\left(\text{temperature} = 66 \mid \text{Yes}\right) = \frac{1}{\sqrt{2\pi}\,(6.2)}\,e^{-\frac{(66-73)^2}{2(6.2)^2}} = 0.0340$$

- Likelihood of Yes = $\dfrac{2}{9} \times 0.0340 \times 0.0221 \times \dfrac{3}{9} \times \dfrac{9}{14} = 0.000036$

- Likelihood of No = $\dfrac{3}{5} \times 0.0291 \times 0.038 \times \dfrac{3}{5} \times \dfrac{5}{14} = 0.000136$

# Bayesian Decision Theory

- Generalization of the preceding ideas
  - Use of more than one feature
  - Use more than two states of nature
  - Allowing actions other than decide on the state of nature
    - Allowing actions other than classification primarily allows the possibility of rejection
    - Refusing to make a decision in close or bad cases!
  - Introduce a loss function which is more general than the probability of error
    - The loss function states how costly each action taken is

# Bayesian Decision Theory – Continuous Features…

- Let $\{\omega_1, \omega_2,…, \omega_c\}$ be the set of c states of nature (or "categories")

- Let $\{\alpha_1, \alpha_2,…, \alpha_a\}$ be the set of possible actions

- Let $\lambda(\alpha_i \mid \omega_j)$ be the loss incurred for taking action $\alpha_i$ when the true state of nature is $\omega_j$

# Two-category Classification

- $\alpha_1$ : deciding $\omega_1$
- $\alpha_2$ : deciding $\omega_2$
- $\lambda_{ij} = \lambda(\alpha_i \mid \omega_j)$
- Loss incurred for deciding $\alpha_i$ when the true state of nature is $\omega_j$

# Two-category Classification

- $\alpha 1$: deciding $\omega 1$
- $\alpha 2$: deciding $\omega 2$
- $\lambda ij = \lambda(\alpha i \mid \omega j)$
- Loss incurred for deciding $\alpha i$ when the true state of nature is $\omega j$
- Conditional risk:

$$
\begin{aligned}
R(\alpha_1|\mathbf{x}) &= \lambda_{11} P(\omega_1|\mathbf{x}) + \lambda_{12} P(\omega_2|\mathbf{x}) \\
R(\alpha_2|\mathbf{x}) &= \lambda_{21} P(\omega_1|\mathbf{x}) + \lambda_{22} P(\omega_2|\mathbf{x}).
\end{aligned}
$$

# Two-category Classification

- Our rule is the following:

  if R($\alpha$1 | x) < R($\alpha$2 | x)

- Action $\alpha$1: "decide $\omega$1" is taken

- This results in the equivalent rule :

- Decide $\omega$1 if:

$$(\lambda_{21} - \lambda_{11})p(\mathbf{x}|\omega_1)P(\omega_1) > (\lambda_{12} - \lambda_{22})p(\mathbf{x}|\omega_2)P(\omega_2)$$

- and decide $\omega$2 otherwise

# Bayesian Decision Theory – Continuous Features…

- Overall risk

  R = Sum of all R($\alpha$i | x) for i = 1,…,a

  **Conditional risk**

- Minimizing R $\Longleftrightarrow$ Minimizing R($\alpha$i | x) for i = 1,…, a

- $R(\alpha_i \mid x) = \sum_{j=1}^{j=c} \lambda(\alpha_i \mid \omega_j) P(\omega_j \mid x)$  for i = 1,…,a