

BAYESIAN DECISION THEORY

J. Elder

CSE 4404/5327 Introduction to Machine Learning and Pattern Recognition

Credits

2

Probability & Bayesian Inference

- Some of these slides were sourced and/or modified from:
 - Christopher Bishop, Microsoft UK
 - Simon Prince, University College London
 - Sergios Theodoridis, University of Athens & Konstantinos Koutroumbas, National Observatory of Athens

Bayesian Decision Theory: Topics

3

Probability & Bayesian Inference

1. Probability
2. The Univariate Normal Distribution
3. Bayesian Classifiers
4. Minimizing Risk
5. The Multivariate Normal Distribution
6. Decision Boundaries in Higher Dimensions
7. Parameter Estimation
8. Mixture Models and EM
9. Nonparametric Density Estimation
10. What are Bayes Nets?

Problems for this Meeting

4

Probability & Bayesian Inference

- Problems 2.1-2.4
- Assigned Problem: 2.2

Bayesian Decision Theory: Topics

5

Probability & Bayesian Inference

1. Probability
2. The Univariate Normal Distribution
3. Bayesian Classifiers
4. Minimizing Risk
5. The Multivariate Normal Distribution
6. Decision Boundaries in Higher Dimensions
7. Parameter Estimation
8. Mixture Models and EM
9. Nonparametric Density Estimation
10. Training and Evaluation Methods
11. What are Bayes Nets?

Topic 1. Probability

Random Variables

7

Probability & Bayesian Inference

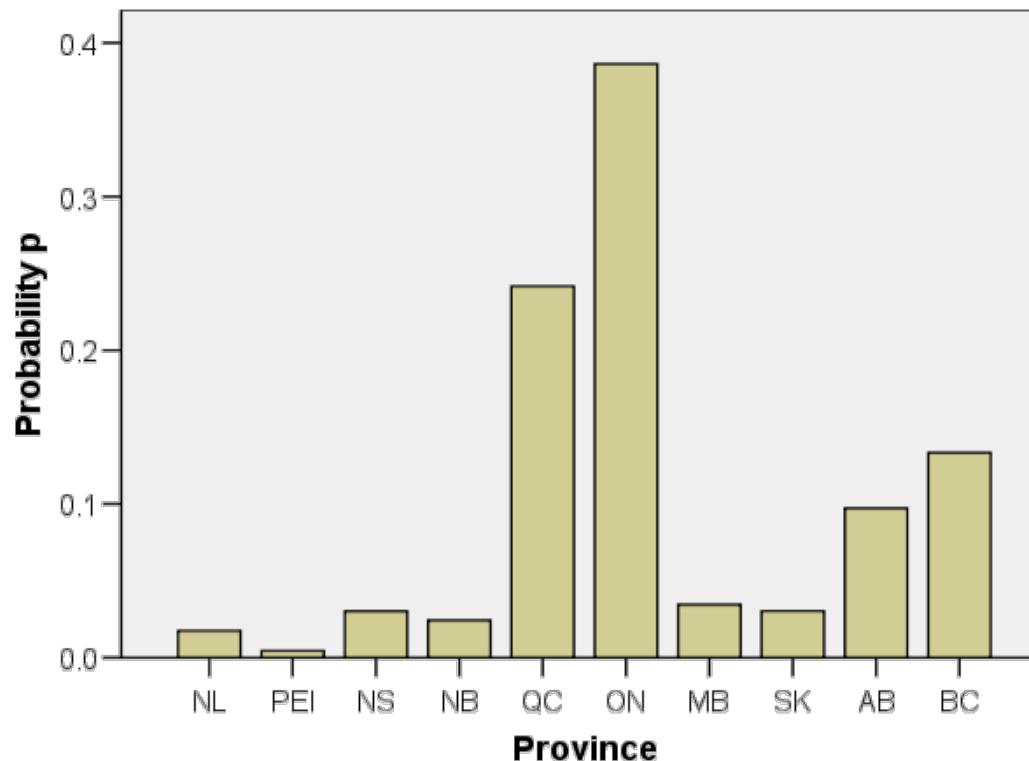
- A **random variable** is a variable whose value is uncertain.
- For example, the height of a randomly selected person in this class is a random variable – I won't know its value until the person is selected.
- Note that we are not completely uncertain about most random variables.
 - For example, we know that height will probably be in the 5'-6' range.
 - In addition, 5'6" is more likely than 5'0" or 6'0".
- The function that describes the probability of each possible value of the random variable is called a **probability distribution**.

Probability Distributions

8

Probability & Bayesian Inference

- For a **discrete** distribution, the probabilities over all possible values of the random variable must **sum** to 1.

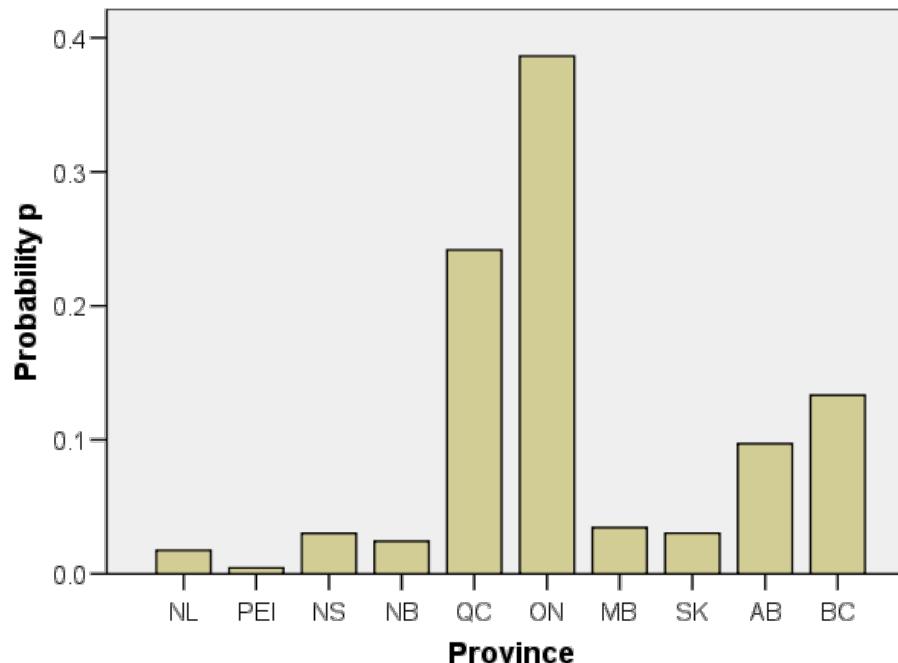


Probability Distributions

9

Probability & Bayesian Inference

- For a **discrete** distribution, we can talk about the probability of a particular score occurring, e.g., $p(\text{Province} = \text{Ontario}) = 0.36$.
- We can also talk about the probability of any one of a subset of scores occurring, e.g., $p(\text{Province} = \text{Ontario or Quebec}) = 0.50$.
- In general, we refer to these occurrences as **events**.



Cases weighted by Sampling weight - master weight

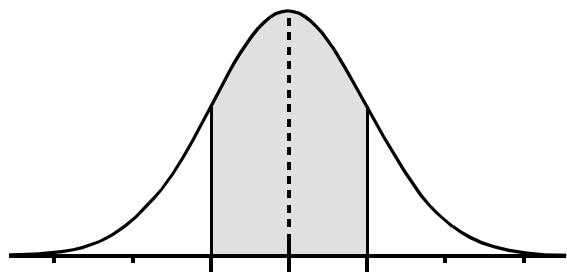
Probability Distributions

10

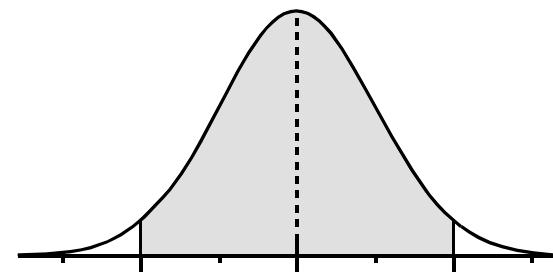
Probability & Bayesian Inference

- For a **continuous** distribution, the probabilities over all possible values of the random variable must **integrate** to 1 (i.e., the area under the curve must be 1).
- Note that the height of a continuous distribution can exceed 1!

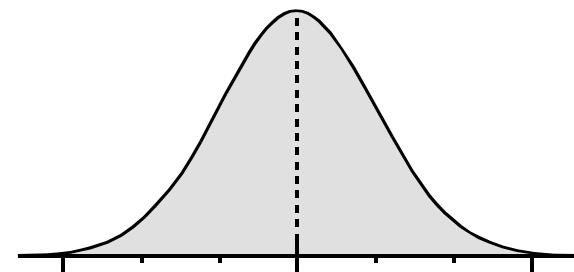
Shaded area = 0.683



Shaded area = 0.954



Shaded area = 0.997

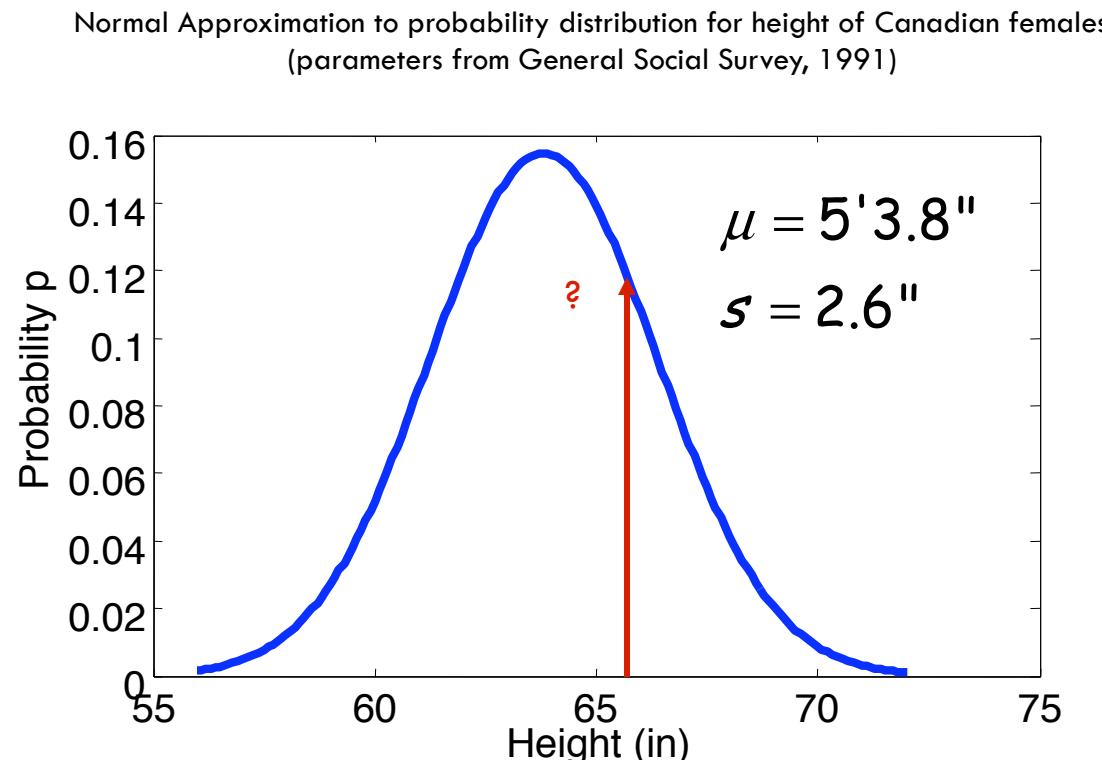


Continuous Distributions

11

Probability & Bayesian Inference

- For continuous distributions, it **does not** make sense to talk about the probability of an exact score.
 - e.g., what is the probability that your height is exactly 65.485948467... inches?



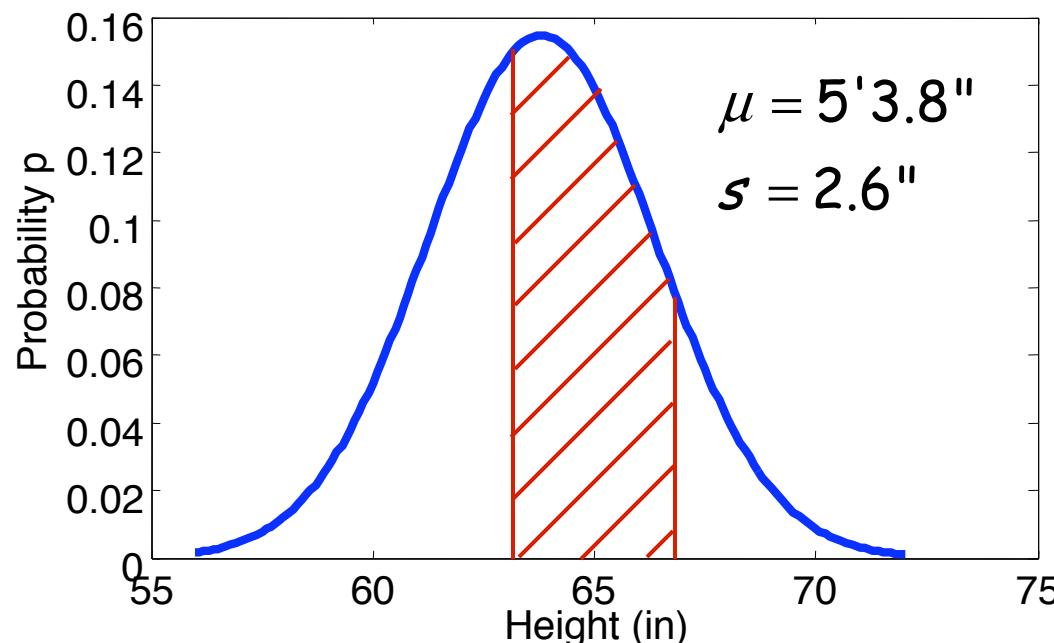
Continuous Distributions

12

Probability & Bayesian Inference

- It **does** make sense to talk about the probability of observing a score that falls within a certain range
 - e.g., what is the probability that you are between 5'3" and 5'7"?
 - e.g., what is the probability that you are less than 5'10"?
- } Valid events

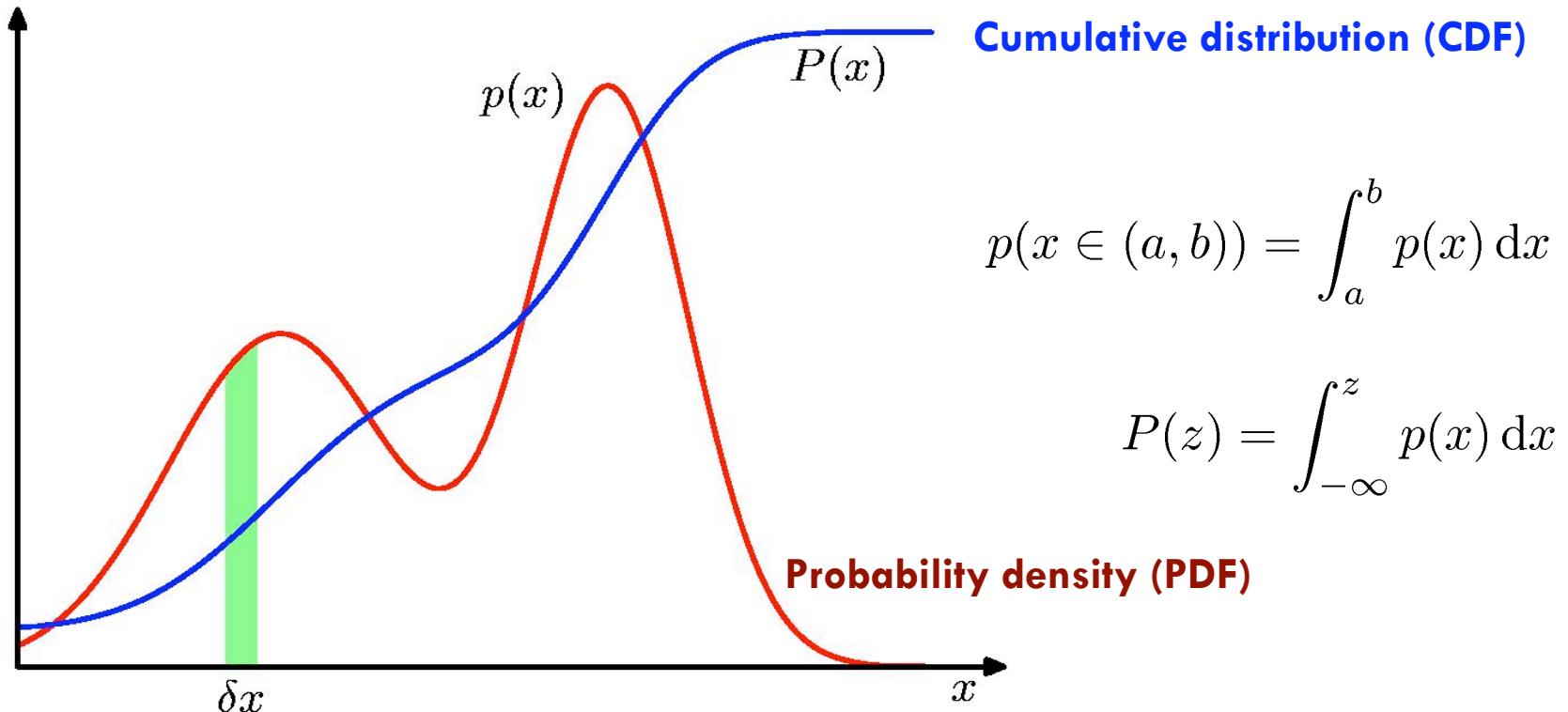
Normal Approximation to probability distribution for height of Canadian females
(parameters from General Social Survey, 1991)



Probability Densities

13

Probability & Bayesian Inference



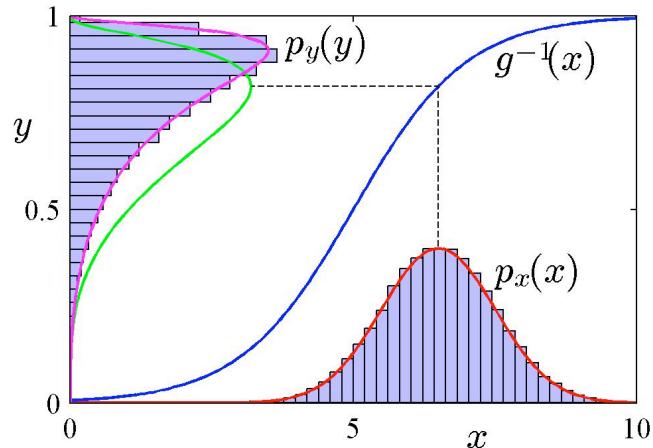
$$p(x) \geq 0$$

$$\int_{-\infty}^{\infty} p(x) dx = 1$$

Transformed Densities

14

Probability & Bayesian Inference



Observations falling within $(x + \delta x)$ transform to the range $(y + \delta y)$

$$\rightarrow p_x(x)|\delta x| = p_y(y)|\delta y|$$

$$\rightarrow p_y(y) \approx p_x(x) \left| \frac{\delta x}{\delta y} \right|$$

Note that in general, $\delta y \neq \delta x$.

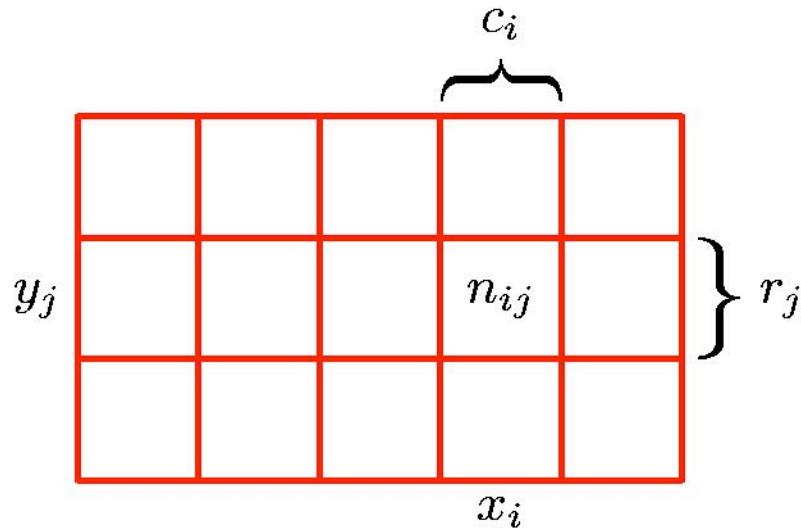
Rather, $\frac{\delta y}{\delta x} \rightarrow \frac{dy}{dx}$ as $\delta x \rightarrow 0$.

$$\text{Thus } p_y(y) = p_x(x) \left| \frac{dx}{dy} \right|$$

Joint Distributions

15

Probability & Bayesian Inference



Marginal Probability

$$p(X = x_i) = \frac{c_i}{N}.$$

Joint Probability

$$p(X = x_i, Y = y_j) = \frac{n_{ij}}{N}$$

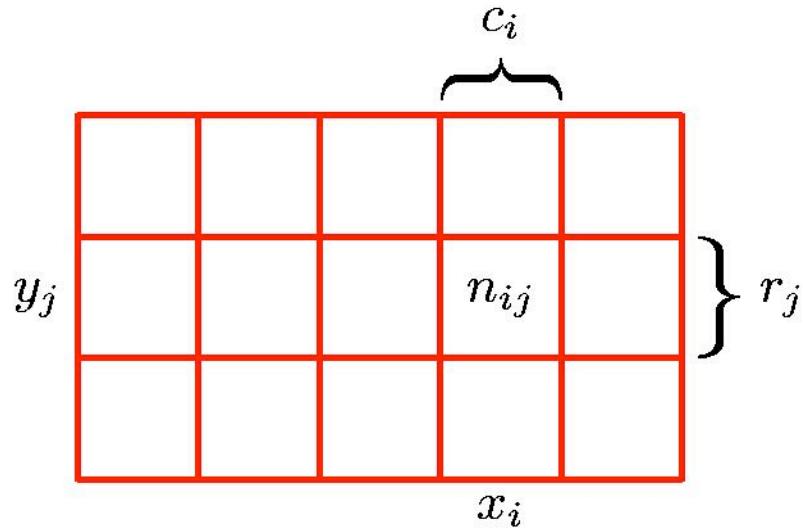
Conditional Probability

$$p(Y = y_j | X = x_i) = \frac{n_{ij}}{c_i}$$

Joint Distributions

16

Probability & Bayesian Inference



Sum Rule

$$\begin{aligned} p(X = x_i) &= \frac{c_i}{N} = \frac{1}{N} \sum_{j=1}^L n_{ij} \\ &= \sum_{j=1}^L p(X = x_i, Y = y_j) \end{aligned}$$

Product Rule

$$\begin{aligned} p(X = x_i, Y = y_j) &= \frac{n_{ij}}{N} = \frac{n_{ij}}{c_i} \cdot \frac{c_i}{N} \\ &= p(Y = y_j | X = x_i) p(X = x_i) \end{aligned}$$

Joint Distributions: The Rules of Probability

17

Probability & Bayesian Inference

- Sum Rule

$$p(X) = \sum_Y p(X, Y)$$

- Product Rule

$$p(X, Y) = p(Y|X)p(X)$$

Marginalization

18

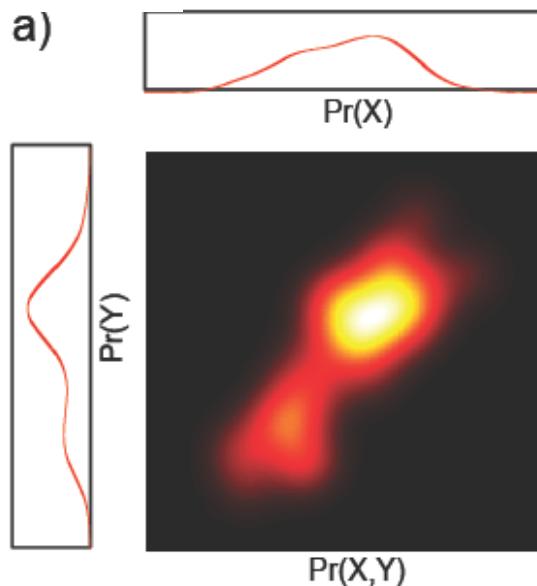
Probability & Bayesian Inference

We can recover probability distribution of any variable in a joint distribution by integrating (or summing) over the other variables

$$Pr(X) = \int Pr(X, Y) dY$$

$$Pr(X, Y) = \sum_W \sum_Z Pr(W, X, Y, Z)$$

$$Pr(Y) = \int Pr(X, Y) dX$$

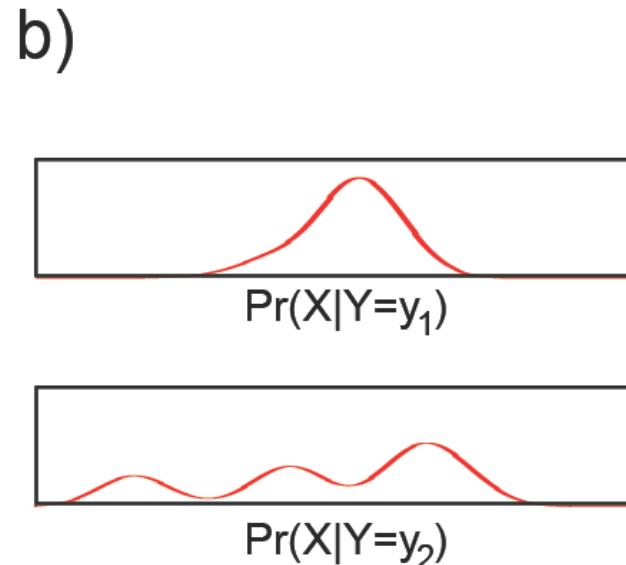
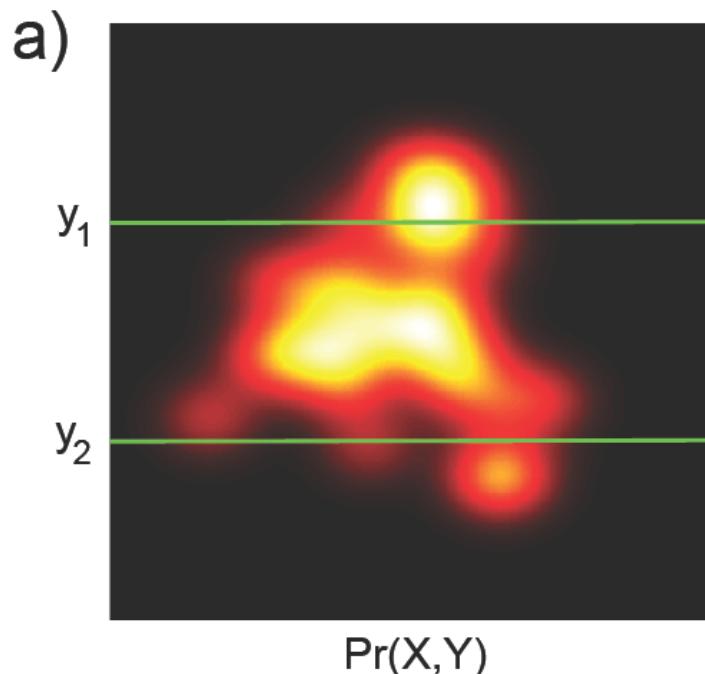


Conditional Probability

19

Probability & Bayesian Inference

- Conditional probability of X given that $Y=y^*$ is relative propensity of variable X to take different outcomes given that Y is fixed to be equal to y^*
- Written as $\Pr(X | Y=y^*)$



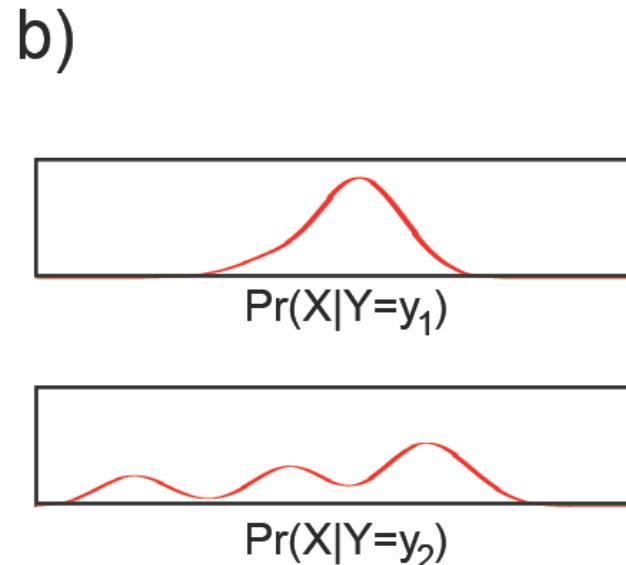
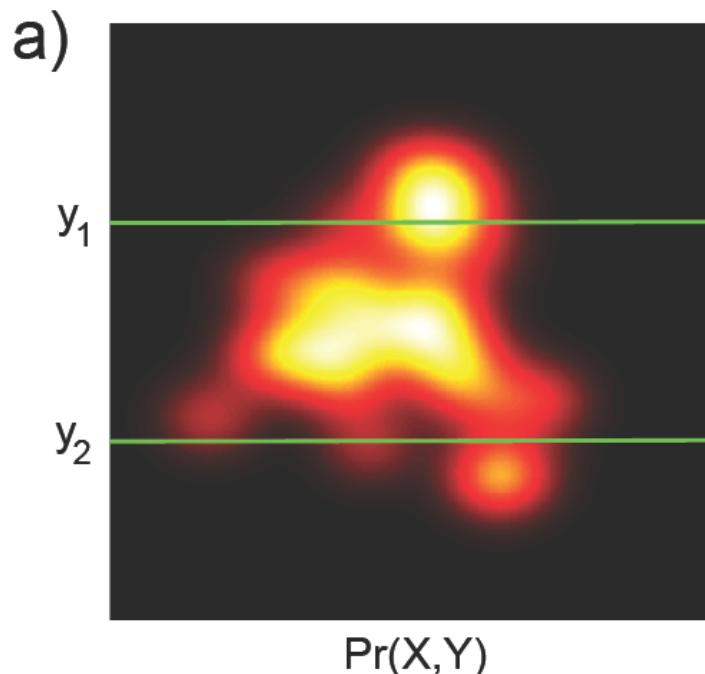
Conditional Probability

20

Probability & Bayesian Inference

- Conditional probability can be extracted from joint probability
- Extract appropriate slice and normalize

$$Pr(X|Y = y^*) = \frac{Pr(X, Y = y^*)}{\int(Pr(X, Y = y^*)dX)} = \frac{Pr(X, Y = y^*)}{Pr(Y = y^*)}$$



Conditional Probability

21

Probability & Bayesian Inference

$$Pr(X|Y = y^*) = \frac{Pr(X, Y = y^*)}{\int(Pr(X, Y = y^*)dX)} = \frac{Pr(X, Y = y^*)}{Pr(Y = y^*)}$$

- More usually written in compact form

$$Pr(X|Y) = \frac{Pr(X, Y)}{Pr(Y)}$$

- Can be re-arranged to give

$$Pr(X, Y) = Pr(X|Y)Pr(Y)$$

Independence

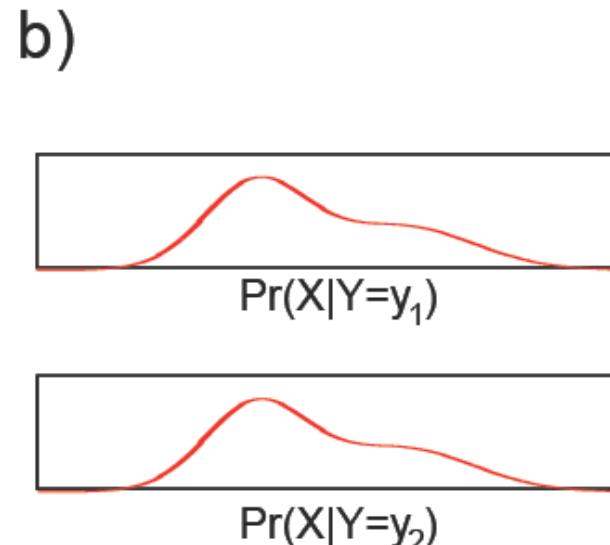
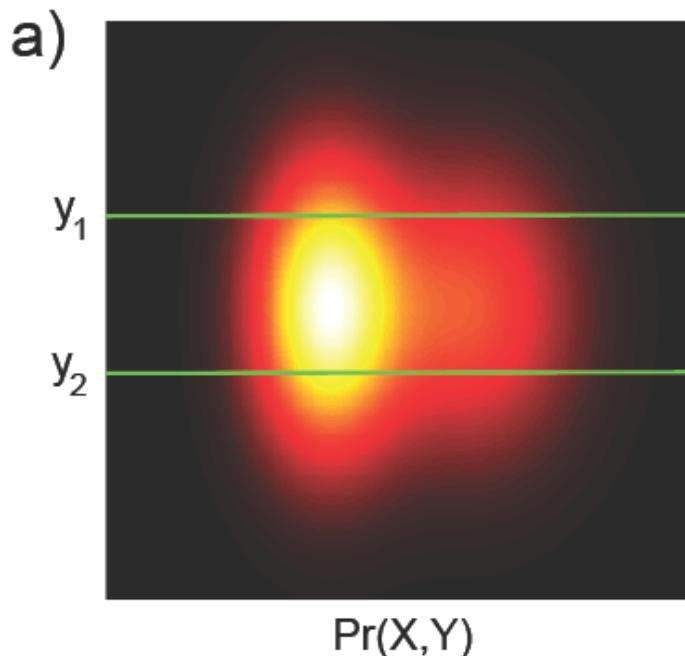
22

Probability & Bayesian Inference

- If two variables X and Y are independent then variable X tells us nothing about variable Y (and vice-versa)

$$Pr(X|Y) = Pr(X)$$

$$Pr(Y|X) = Pr(Y)$$



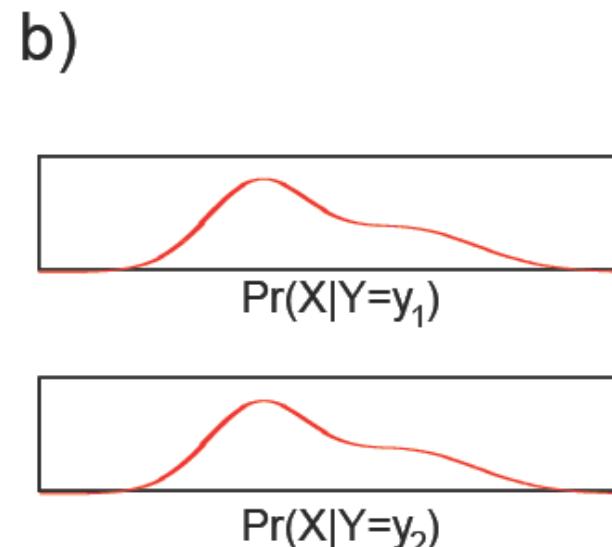
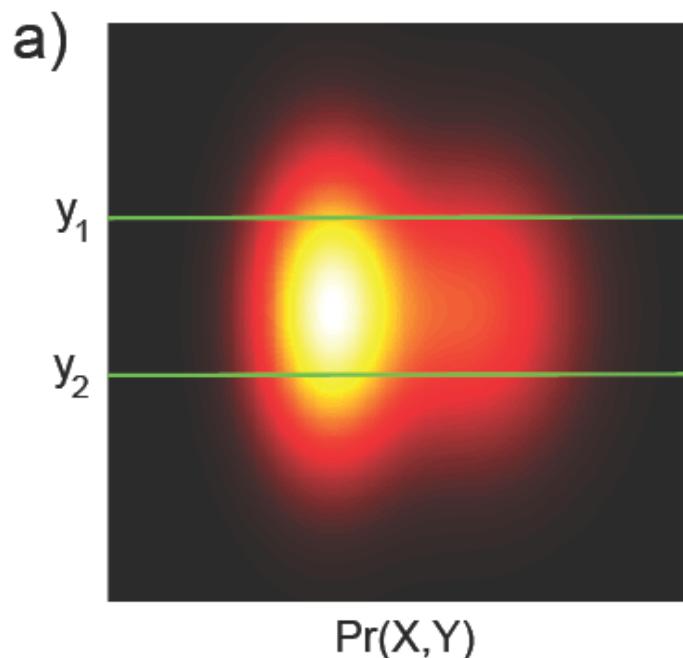
Independence

23

Probability & Bayesian Inference

- When variables are independent, the joint factorizes into a product of the marginals:

$$\begin{aligned} \Pr(X, Y) &= \Pr(X|Y)\Pr(Y) \\ &= \Pr(X)\Pr(Y) \end{aligned}$$



Bayes' Rule

24

Probability & Bayesian Inference

From before:

$$Pr(X, Y) = Pr(X|Y)Pr(Y)$$

$$Pr(X, Y) = Pr(Y|X)Pr(X)$$

Combining:

$$Pr(Y|X)Pr(X) = Pr(X|Y)Pr(Y)$$

Re-arranging:

$$\begin{aligned} Pr(Y|X) &= \frac{Pr(X|Y)Pr(Y)}{Pr(X)} \\ &= \frac{Pr(X|Y)Pr(Y)}{\int Pr(X, Y)dY} \\ &= \frac{Pr(X|Y)Pr(Y)}{\int Pr(X|Y)Pr(Y)dY} \end{aligned}$$

Bayes' Rule Terminology

25

Probability & Bayesian Inference

Likelihood – propensity for observing a certain value of X given a certain value of Y

Prior – what we know about y before seeing x

$$Pr(Y|X) = \frac{Pr(X|Y)Pr(Y)}{Pr(X)}$$

Posterior – what we know about y after seeing x

Evidence – a constant to ensure that the left hand side is a valid distribution



End of Lecture 2

Bayesian Decision Theory: Topics

27

Probability & Bayesian Inference

1. Probability
2. **The Univariate Normal Distribution**
3. Bayesian Classifiers
4. Minimizing Risk
5. The Multivariate Normal Distribution
6. Decision Boundaries in Higher Dimensions
7. Parameter Estimation
8. Mixture Models and EM
9. Nonparametric Density Estimation
10. Training and Evaluation Methods
11. What are Bayes Nets?

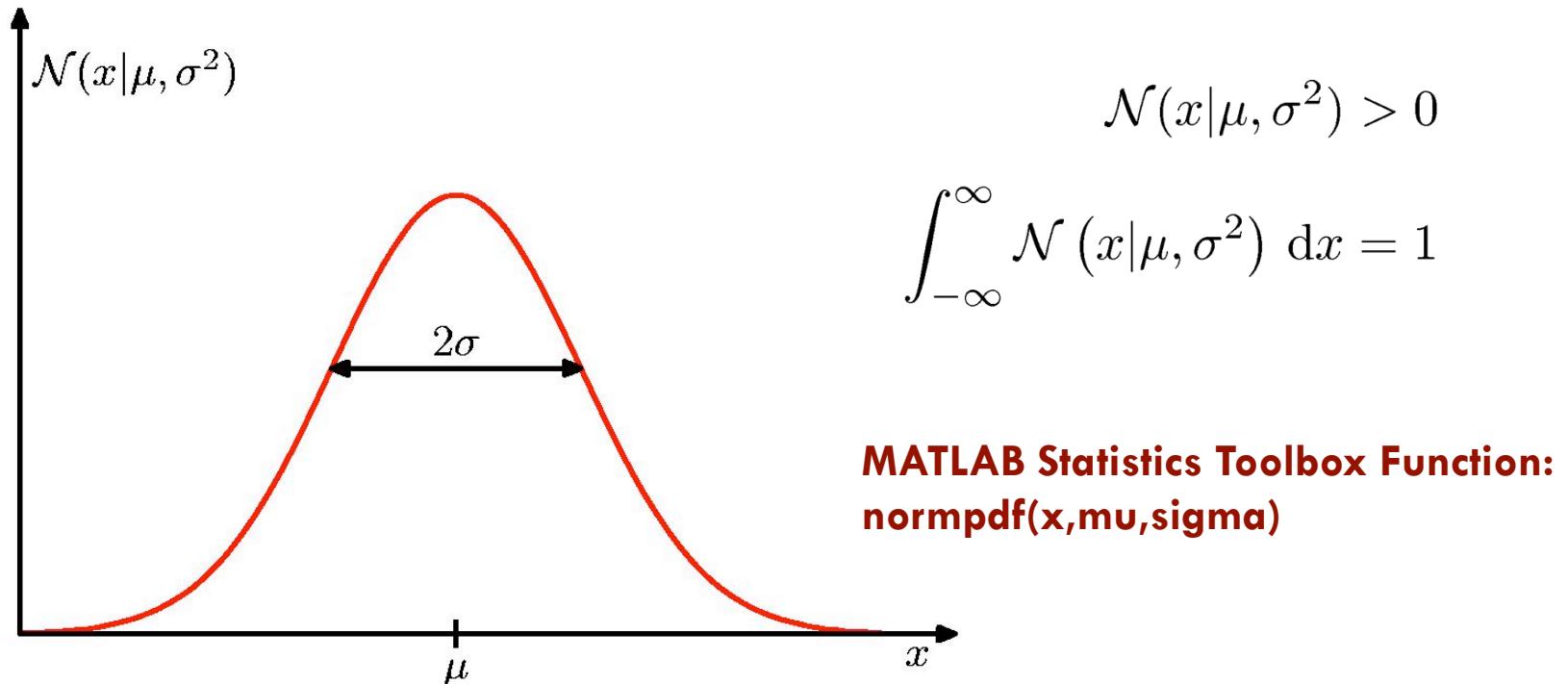
Topic 2. The Univariate Normal Distribution

The Gaussian Distribution

29

Probability & Bayesian Inference

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left\{ -\frac{1}{2\sigma^2}(x - \mu)^2 \right\}$$

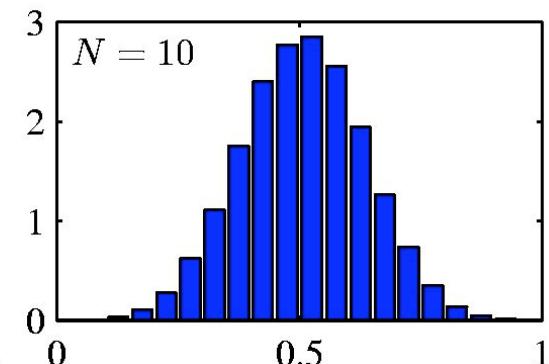
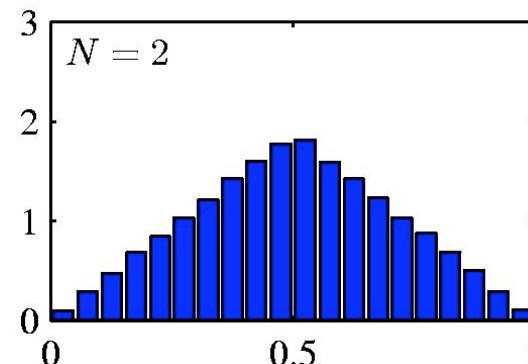
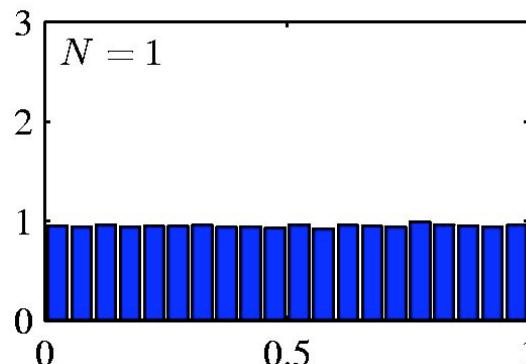


Central Limit Theorem

30

Probability & Bayesian Inference

- The distribution of the sum of N i.i.d. random variables becomes increasingly Gaussian as N grows.
- Example: N uniform $[0, 1]$ random variables.



Expectations

31

Probability & Bayesian Inference

$$\mathbb{E}[f] = \sum_x p(x)f(x)$$

$$\mathbb{E}[f] = \int p(x)f(x) dx$$

$$\mathbb{E}_x[f|y] = \sum_x p(x|y)f(x)$$


Conditional Expectation
(discrete)

$$\mathbb{E}[f] \simeq \frac{1}{N} \sum_{n=1}^N f(x_n)$$

Approximate Expectation
(discrete and continuous)

Variances and Covariances

32

Probability & Bayesian Inference

$$\text{var}[f] = \mathbb{E} \left[(f(x) - \mathbb{E}[f(x)])^2 \right] = \mathbb{E}[f(x)^2] - \mathbb{E}[f(x)]^2$$

$$\begin{aligned}\text{cov}[x, y] &= \mathbb{E}_{x,y} [\{x - \mathbb{E}[x]\} \{y - \mathbb{E}[y]\}] \\ &= \mathbb{E}_{x,y}[xy] - \mathbb{E}[x]\mathbb{E}[y]\end{aligned}$$

$$\begin{aligned}\text{cov}[\mathbf{x}, \mathbf{y}] &= \mathbb{E}_{\mathbf{x},\mathbf{y}} [\{\mathbf{x} - \mathbb{E}[\mathbf{x}]\}\{\mathbf{y}^T - \mathbb{E}[\mathbf{y}^T]\}] \\ &= \mathbb{E}_{\mathbf{x},\mathbf{y}}[\mathbf{x}\mathbf{y}^T] - \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{y}^T]\end{aligned}$$

Gaussian Mean and Variance

33

Probability & Bayesian Inference

$$\mathbb{E}[x] = \int_{-\infty}^{\infty} \mathcal{N}(x|\mu, \sigma^2) x \, dx = \mu$$

$$\mathbb{E}[x^2] = \int_{-\infty}^{\infty} \mathcal{N}(x|\mu, \sigma^2) x^2 \, dx = \mu^2 + \sigma^2$$

$$\text{var}[x] = \mathbb{E}[x^2] - \mathbb{E}[x]^2 = \sigma^2$$

Bayesian Decision Theory: Topics

34

Probability & Bayesian Inference

1. Probability
2. The Univariate Normal Distribution
3. **Bayesian Classifiers**
4. Minimizing Risk
5. The Multivariate Normal Distribution
6. Decision Boundaries in Higher Dimensions
7. Parameter Estimation
8. Mixture Models and EM
9. Nonparametric Density Estimation
10. Training and Evaluation Methods
11. What are Bayes Nets?

Topic 3. Bayesian Classifiers

Bayesian Classification

36

Probability & Bayesian Inference

- Input feature vectors

$$\mathbf{x} = [x_1, x_2, \dots, x_l]^T$$

- Assign the pattern represented by feature vector \mathbf{x} to the **most probable** of the available classes

$$\omega_1, \omega_2, \dots, \omega_M$$

That is, $\mathbf{x} \rightarrow \omega_i : P(\omega_i | \mathbf{x})$ is maximum.

↑
Posterior

Bayesian Classification

37

Probability & Bayesian Inference

- Computation of **posterior** probabilities

- Assume known

- **Prior** probabilities

$$P(\omega_1), P(\omega_2), \dots, P(\omega_M)$$

- **Likelihoods**

$$p(\mathbf{x} | \omega_i), \quad i = 1, 2, \dots, M$$

Bayes' Rule for Classification

38

Probability & Bayesian Inference

$$p(\omega_i | \mathbf{x}) = \frac{p(\mathbf{x} | \omega_i)p(\omega_i)}{p(\mathbf{x})},$$

where

$$p(\mathbf{x}) = \sum_{i=1}^M p(\mathbf{x} | \omega_i)p(\omega_i)$$

M=2 Classes

39

Probability & Bayesian Inference

- Given \mathbf{x} classify it according to the rule

If $P(\omega_1|\mathbf{x}) > P(\omega_2|\mathbf{x}) \rightarrow \omega_1$

If $P(\omega_2|\mathbf{x}) > P(\omega_1|\mathbf{x}) \rightarrow \omega_2$

- Equivalently: classify \mathbf{x} according to the rule

If $p(\mathbf{x}|\omega_1)P(\omega_1) > p(\mathbf{x}|\omega_2)P(\omega_2) \rightarrow \omega_1$

If $p(\mathbf{x}|\omega_2)P(\omega_2) > p(\mathbf{x}|\omega_1)P(\omega_1) \rightarrow \omega_2$

- For equiprobable classes the test becomes

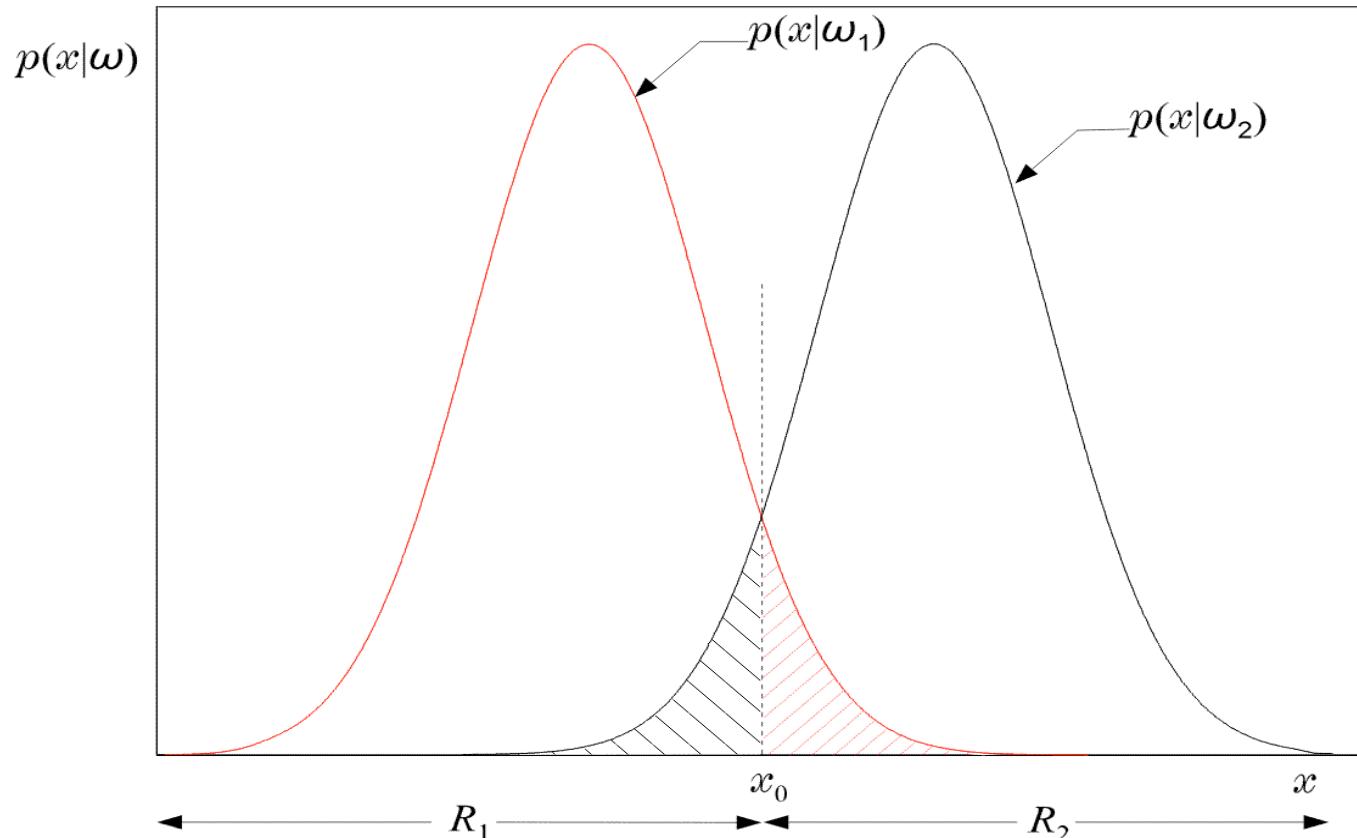
If $p(\mathbf{x}|\omega_1) > p(\mathbf{x}|\omega_2) \rightarrow \omega_1$

If $p(\mathbf{x}|\omega_2) > p(\mathbf{x}|\omega_1) \rightarrow \omega_2$

Example: Equiprobable Classes

40

Probability & Bayesian Inference



$R_1(\rightarrow \omega_1)$ and $R_2(\rightarrow \omega_2)$

Example: Equiprobable Classes

41

Probability & Bayesian Inference

- Probability of error

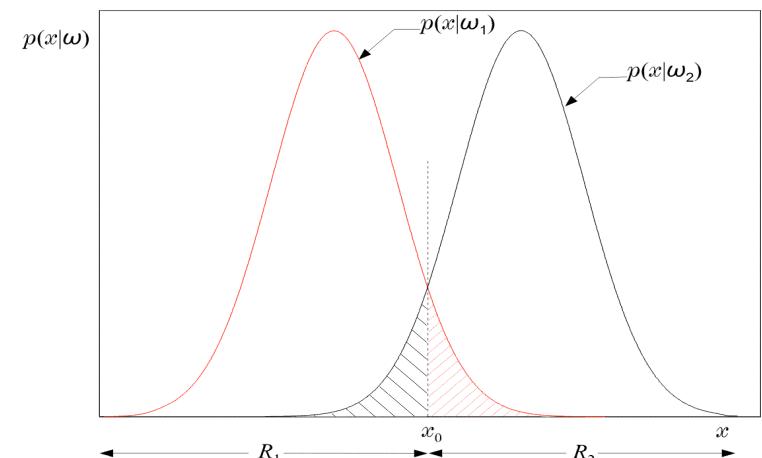
- The black and red shaded areas represent

$$P(\text{error} | \omega_2) = \int_{-\infty}^{x_0} p(x|\omega_2)dx \quad \text{and} \quad P(\text{error} | \omega_1) = \int_{x_0}^{\infty} p(x|\omega_1)dx$$

- Thus

$$\begin{aligned} P_e &\triangleq P(\text{error}) \\ &= P(\omega_2)P(\text{error}|\omega_2) + P(\omega_1)P(\text{error}|\omega_1) \\ &= \frac{1}{2} \int_{-\infty}^{x_0} p(x|\omega_2)dx + \frac{1}{2} \int_{x_0}^{+\infty} p(x|\omega_1)dx \end{aligned}$$

- **Bayesian classifier is OPTIMAL: it minimizes the classification error probability**

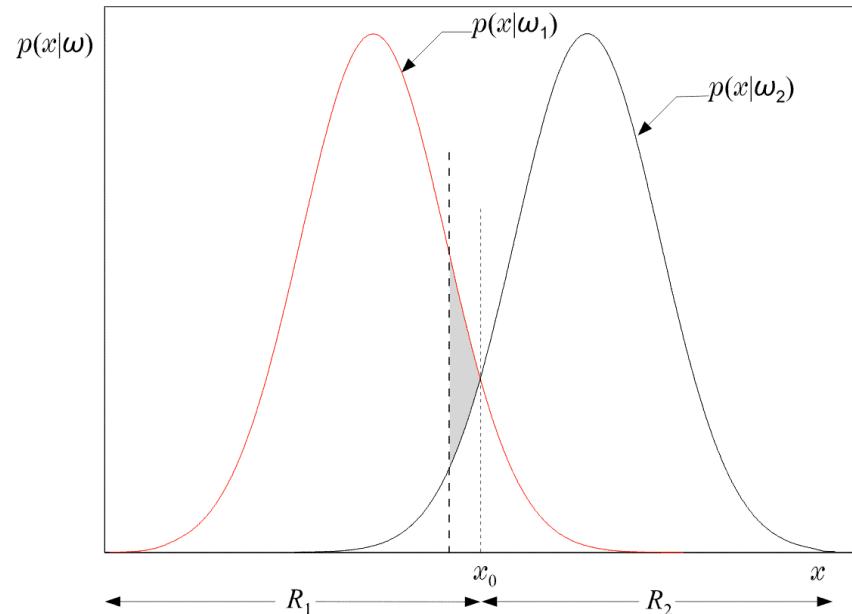


Example: Equiprobable Classes

42

Probability & Bayesian Inference

- To see this, observe that shifting the threshold increases the error rate for one class of patterns more than it decreases the error rate for the other class.



The General Case

43

Probability & Bayesian Inference

- In general, for M classes and unequal priors, the decision rule

$$P(\omega_i | \mathbf{x}) > P(\omega_j | \mathbf{x}) \quad \forall j \neq i \quad \rightarrow \omega_i$$

minimizes the expected error rate.

Types of Error

44

Probability & Bayesian Inference

- Minimizing the expected error rate is a pretty reasonable goal.
- However, it is not always the best thing to do.
- Example:
 - You are designing a pedestrian detection algorithm for an autonomous navigation system.
 - Your algorithm must decide whether there is a pedestrian crossing the street.
 - There are two possible types of error:
 - False positive: there is no pedestrian, but the system thinks there is.
 - Miss: there is a pedestrian, but the system thinks there is not.
 - Should you give equal weight to these 2 types of error?

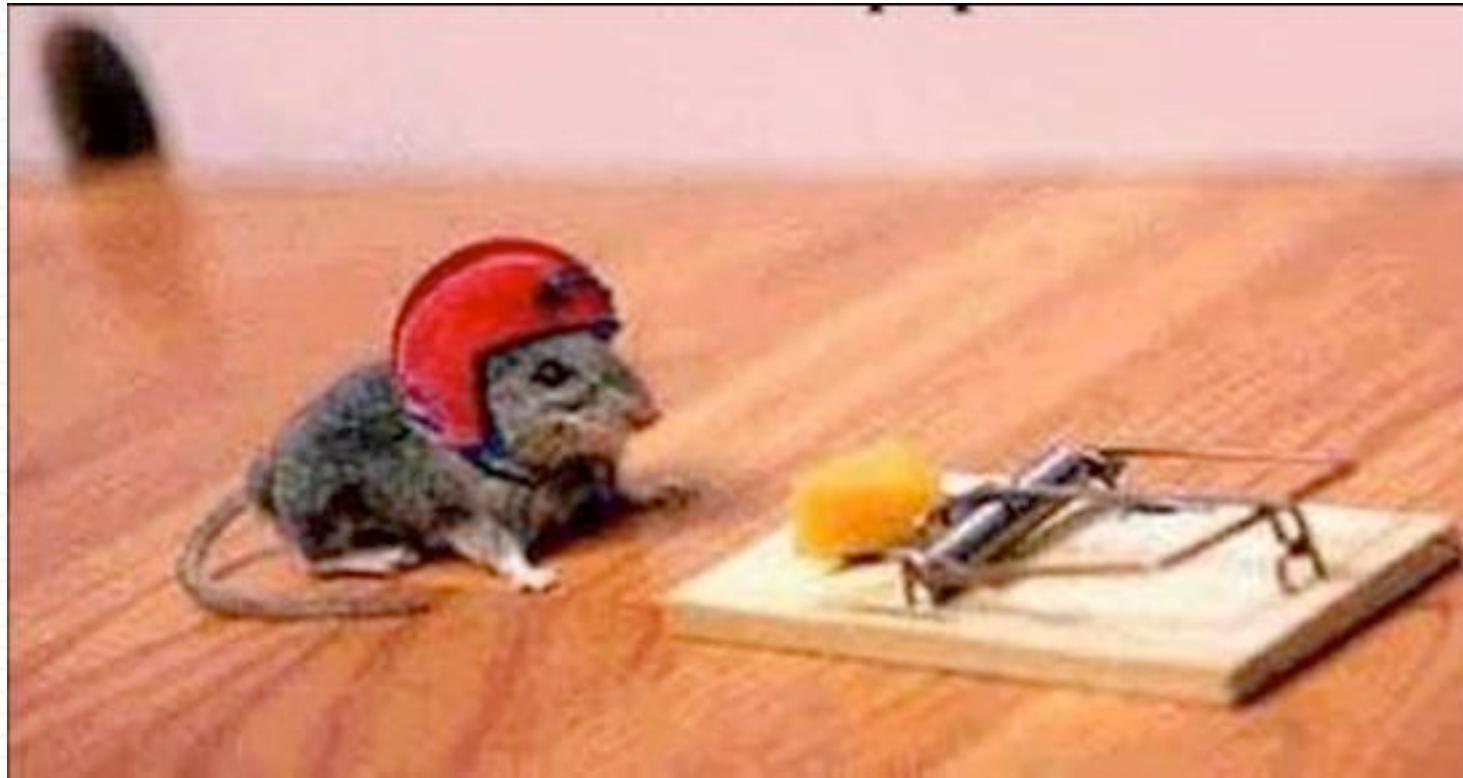
Bayesian Decision Theory: Topics

45

Probability & Bayesian Inference

1. Probability
2. The Univariate Normal Distribution
3. Bayesian Classifiers
4. **Minimizing Risk**
5. The Multivariate Normal Distribution
6. Decision Boundaries in Higher Dimensions
7. Parameter Estimation
8. Mixture Models and EM
9. Nonparametric Density Estimation
10. Training and Evaluation Methods
11. What are Bayes Nets?

Topic 4. Minimizing Risk



The Loss Matrix

47

Probability & Bayesian Inference

- To deal with this problem, instead of minimizing error rate, we minimize something called the **risk**.
- First, we define the **loss matrix L** , which quantifies the cost of making each type of error.
- Element λ_{ij} of the loss matrix specifies the cost of deciding class j when in fact the input is of class i .
- Typically, we set $\lambda_{ii}=0$ for all i .
- Thus a typical loss matrix for the $M = 2$ case would have the form

$$L = \begin{bmatrix} 0 & \lambda_{12} \\ \lambda_{21} & 0 \end{bmatrix}$$

Risk

48

Probability & Bayesian Inference

- Given a loss function, we can now define the risk associated with each class k as:

$$r_k = \sum_{i=1}^M \lambda_{ki} \int_{R_i} p(\mathbf{x} | \omega_k) d\mathbf{x}$$

Probability we will decide Class ω_i , given pattern from Class ω_k

- where R_i is the region of the input space where we will decide ω_i .

Minimizing Risk

49

Probability & Bayesian Inference

- Now the goal is to minimize the expected risk r , given by

$$r = \sum_{k=1}^M r_k P(\omega_k)$$

Minimizing Risk

50

Probability & Bayesian Inference

$$r = \sum_{k=1}^M r_k P(\omega_k) \quad \text{where} \quad r_k = \sum_{i=1}^M \lambda_{ki} \int_{R_i} p(\mathbf{x} | \omega_k) d\mathbf{x}$$

- We need to select the decision regions R_i to minimize the risk r .
- Note that the set of R_i are disjoint and exhaustive.
- Thus we can minimize the risk by ensuring that each input \mathbf{x} falls in the region R_i that minimizes the expected loss for that particular input, i.e.,

$$\text{Letting } l_i = \sum_{k=1}^M \lambda_{ki} p(\mathbf{x} | \omega_k) P(\omega_k),$$

we select the partitioning regions such that

$$\mathbf{x} \in R_i \text{ if } l_i < l_j \quad \forall j \neq i$$

Example: M=2

51

Probability & Bayesian Inference

- For the 2-class case:

$$l_1 = \lambda_{11} p(x | \omega_1) P(\omega_1) + \lambda_{21} p(x | \omega_2) P(\omega_2)$$

and

$$l_2 = \lambda_{12} p(x | \omega_1) P(\omega_1) + \lambda_{22} p(x | \omega_2) P(\omega_2)$$

- Thus we assign x to ω_1 if

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

- i.e., if

Likelihood Ratio Test

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

Likelihood Ratio Test

52

Probability & Bayesian Inference

$$\frac{P(x | \omega_1)}{P(x | \omega_2)} ? \frac{P(\omega_2)(\lambda_{21} - \lambda_{22})}{P(\omega_1)(\lambda_{12} - \lambda_{11})}.$$

- Typically, the loss for a correct decision is 0. Thus the likelihood ratio test becomes

$$\frac{P(x | \omega_1)}{P(x | \omega_2)} ? \frac{P(\omega_2)\lambda_{21}}{P(\omega_1)\lambda_{12}}.$$

- In the case of equal priors and equal loss functions, the test reduces to

$$\frac{P(x | \omega_1)}{P(x | \omega_2)} ? 1.$$

Example

53

Probability & Bayesian Inference

- Consider a one-dimensional input space, with features generated by normal distributions with identical variance:

$$p(x|\omega_1) \sim N(\mu_1, \sigma^2)$$

$$p(x|\omega_2) \sim N(\mu_2, \sigma^2)$$

where $\mu_1 = 0$, $\mu_2 = 1$, and $\sigma^2 = \frac{1}{2}$

- Let's assume equiprobable classes, and higher loss for errors on Class 2, specifically:

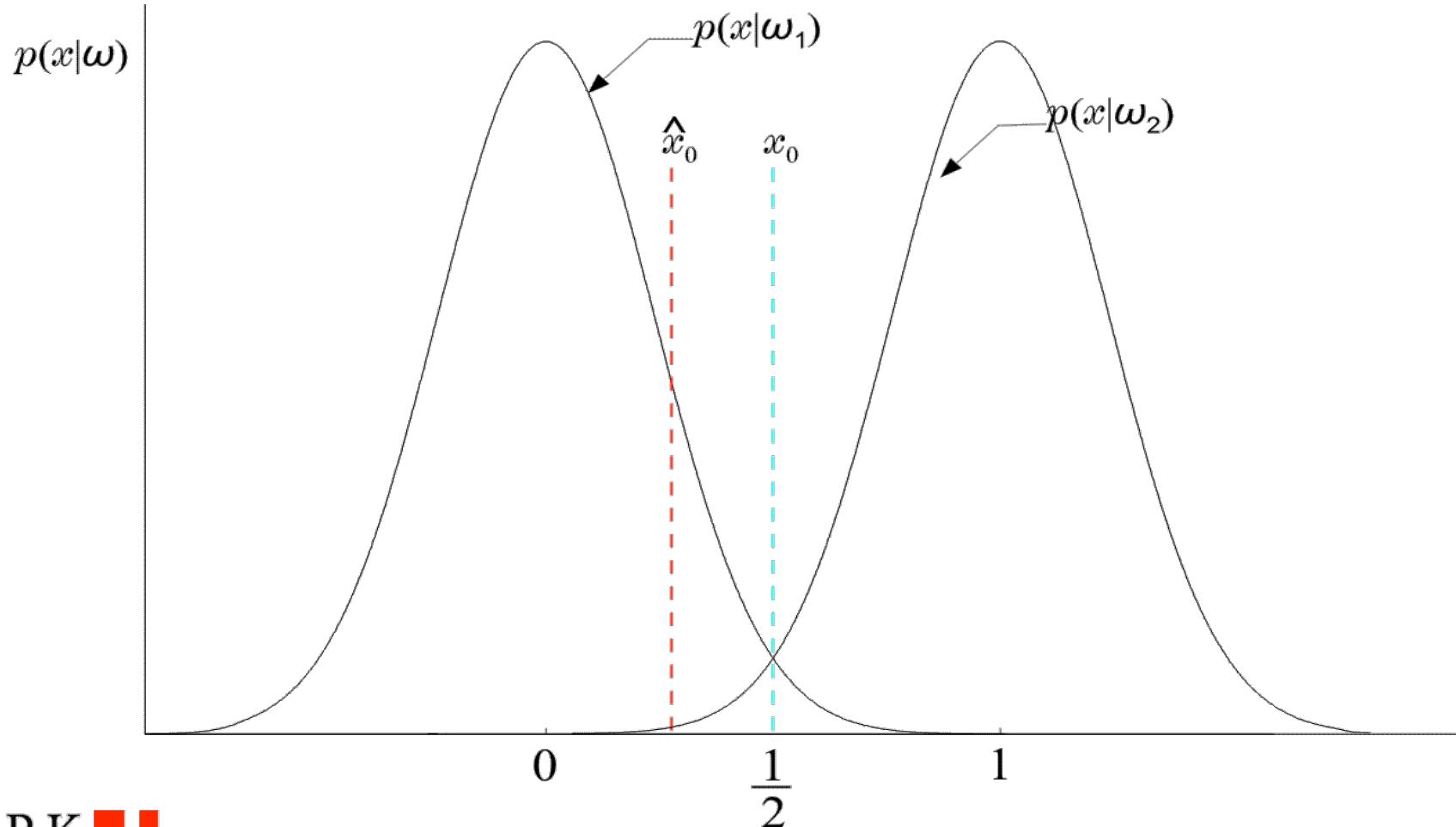
$$\lambda_{21} = 1, \quad \lambda_{12} = \frac{1}{2}.$$

Results

54

Probability & Bayesian Inference

- The threshold has shifted to the left – why?





End of Lecture 3

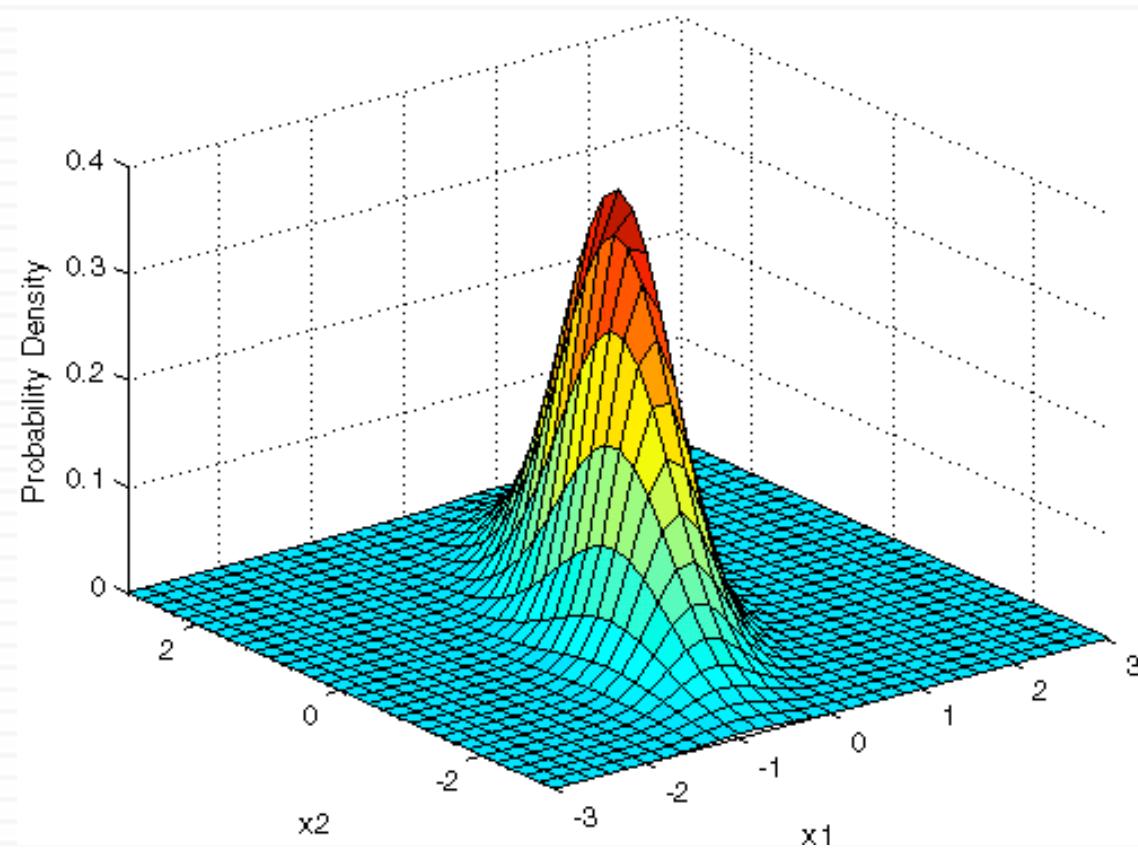
Bayesian Decision Theory: Topics

56

Probability & Bayesian Inference

1. Probability
2. The Univariate Normal Distribution
3. Bayesian Classifiers
4. Minimizing Risk
5. **The Multivariate Normal Distribution**
6. Decision Boundaries in Higher Dimensions
7. Parameter Estimation
8. Mixture Models and EM
9. Nonparametric Density Estimation
10. Training and Evaluation Methods
11. What are Bayes Nets?

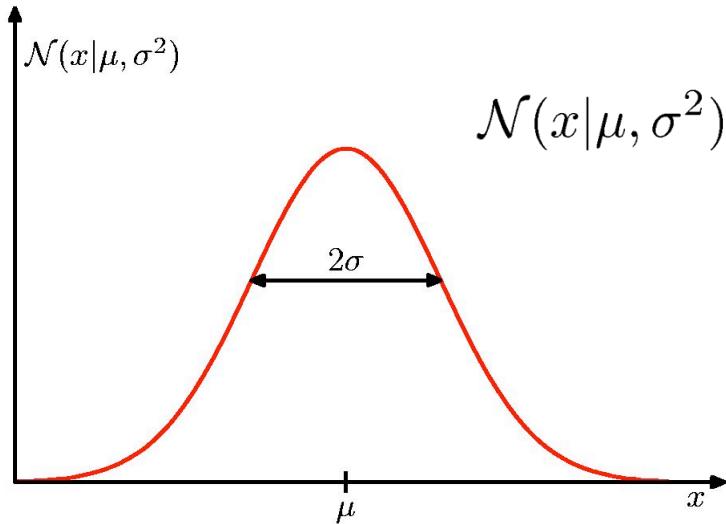
Topic 5 The Multivariate Normal Distribution



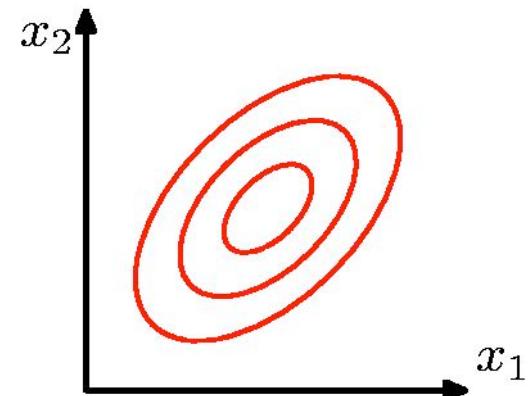
The Multivariate Gaussian

58

Probability & Bayesian Inference



$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left\{ -\frac{1}{2\sigma^2}(x - \mu)^2 \right\}$$



MATLAB Statistics Toolbox Function:
`mvnpdf(x,mu,sigma)`

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

Geometry of the Multivariate Gaussian

59

Probability & Bayesian Inference

$$\Delta^2 = (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \quad \text{where } \Delta \equiv \text{Mahalanobis distance from } \boldsymbol{\mu} \text{ to } \mathbf{x}$$

Eigenvector equation: $\boldsymbol{\Sigma} \mathbf{u}_i = \lambda_i \mathbf{u}_i$

MATLAB Statistics Toolbox Function:
`mahal(x,y)`

where $(\mathbf{u}_i, \lambda_i)$ are the i th eigenvector and eigenvalue of $\boldsymbol{\Sigma}$.

Note that $\boldsymbol{\Sigma}$ real and symmetric $\rightarrow \lambda_i$ real.

**See Appendix B for a review of
matrices and eigenvectors.**

Geometry of the Multivariate Gaussian

60

Probability & Bayesian Inference

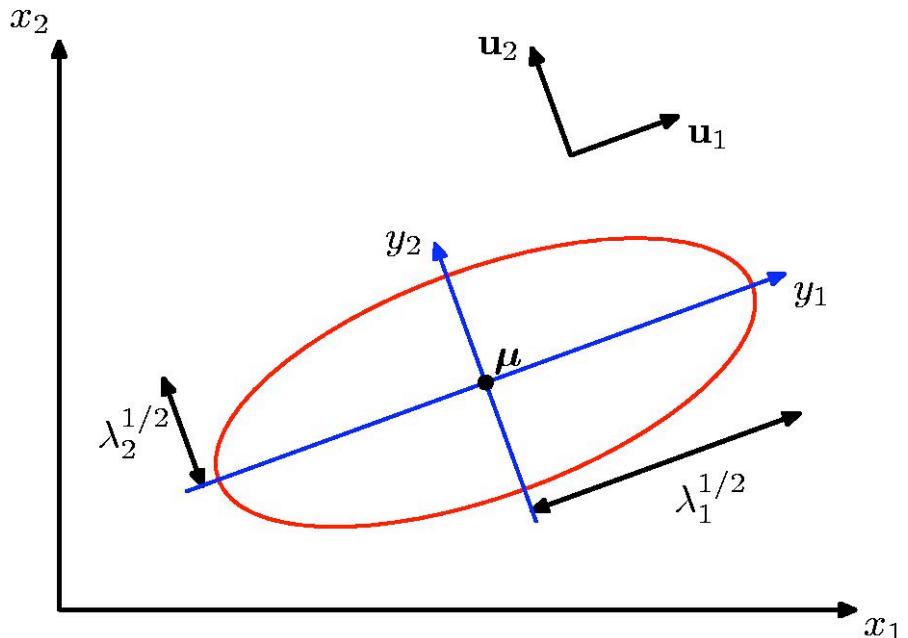
$$\Delta^2 = (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \quad \Delta = \text{Mahalanobis distance from } \boldsymbol{\mu} \text{ to } \mathbf{x}$$

$$\boldsymbol{\Sigma}^{-1} = \sum_{i=1}^D \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T \quad \text{where } (\mathbf{u}_i, \lambda_i) \text{ are the } i\text{th eigenvector and eigenvalue of } \boldsymbol{\Sigma}.$$

$$\Delta^2 = \sum_{i=1}^D \frac{y_i^2}{\lambda_i}$$

$$y_i = \mathbf{u}_i^T (\mathbf{x} - \boldsymbol{\mu})$$

$$\text{or } \mathbf{y} = \mathbf{U}(\mathbf{x} - \boldsymbol{\mu})$$



Moments of the Multivariate Gaussian

61

Probability & Bayesian Inference

$$\begin{aligned}\mathbb{E}[\mathbf{x}] &= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \int \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \mathbf{x} d\mathbf{x} \\ &= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \int \exp \left\{ -\frac{1}{2}\mathbf{z}^T \Sigma^{-1} \mathbf{z} \right\} (\mathbf{z} + \boldsymbol{\mu}) d\mathbf{z}\end{aligned}$$

thanks to anti-symmetry of Σ

$$\mathbb{E}[\mathbf{x}] = \boldsymbol{\mu}$$

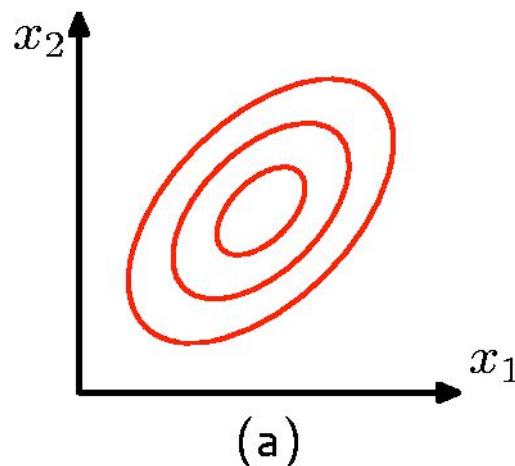
Moments of the Multivariate Gaussian

62

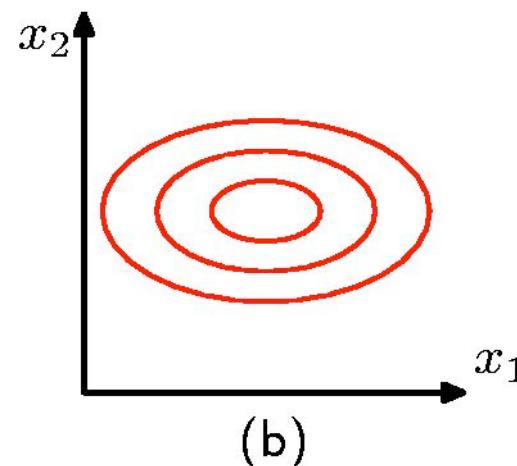
Probability & Bayesian Inference

$$\mathbb{E}[\mathbf{x}\mathbf{x}^T] = \boldsymbol{\mu}\boldsymbol{\mu}^T + \boldsymbol{\Sigma}$$

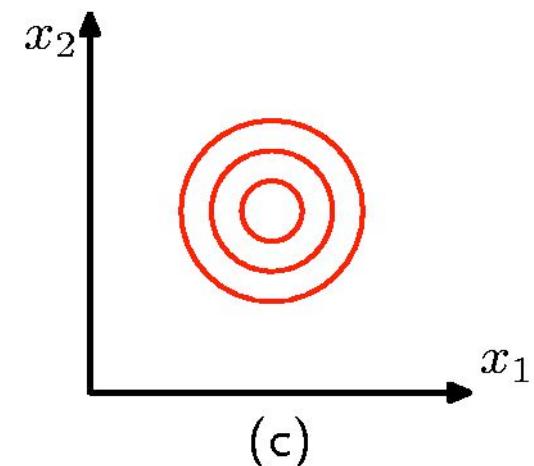
$$\text{cov}[\mathbf{x}] = \mathbb{E} [(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T] = \boldsymbol{\Sigma}$$



(a)



(b)



(c)

5.1 Application: Face Detection



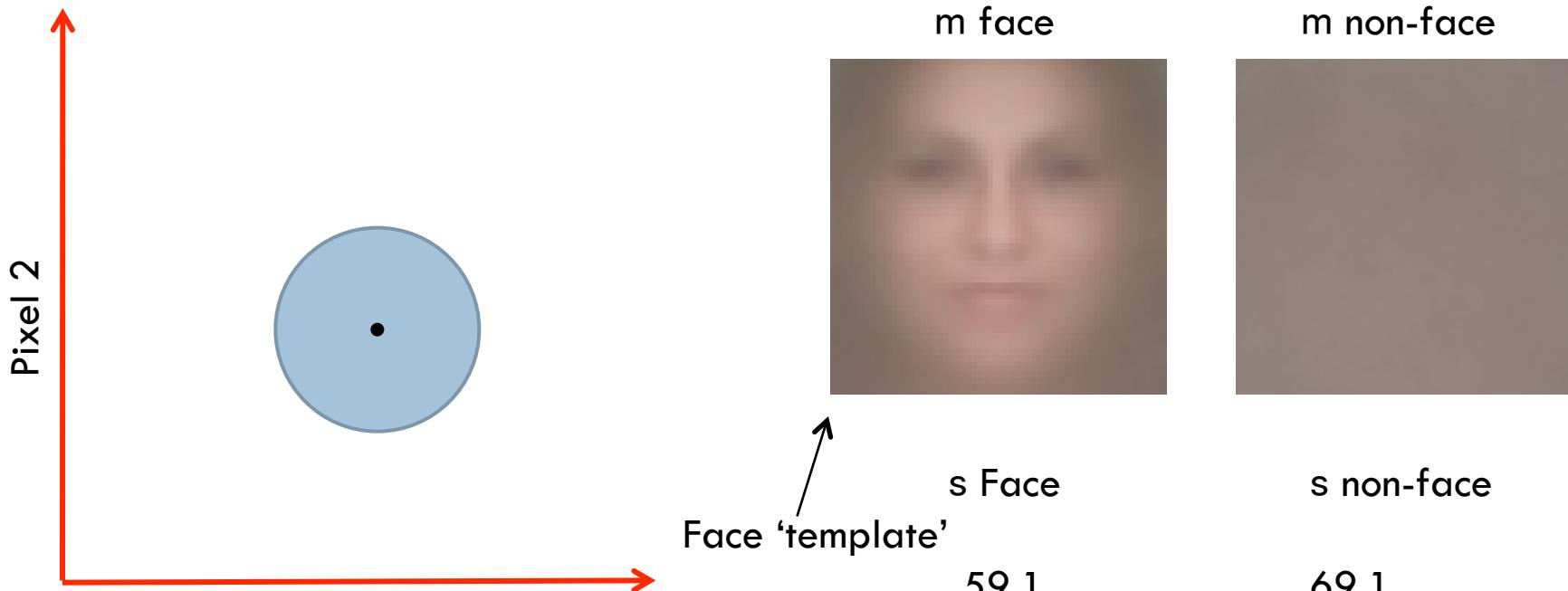
Model # 1: Gaussian, uniform covariance

64

Probability & Bayesian Inference

$$Pr(\mathbf{x}|\text{face}) = \frac{1}{(2\pi)^{d/2}|\Sigma|} \exp \left\{ -0.5(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

Fit model using maximum likelihood criterion

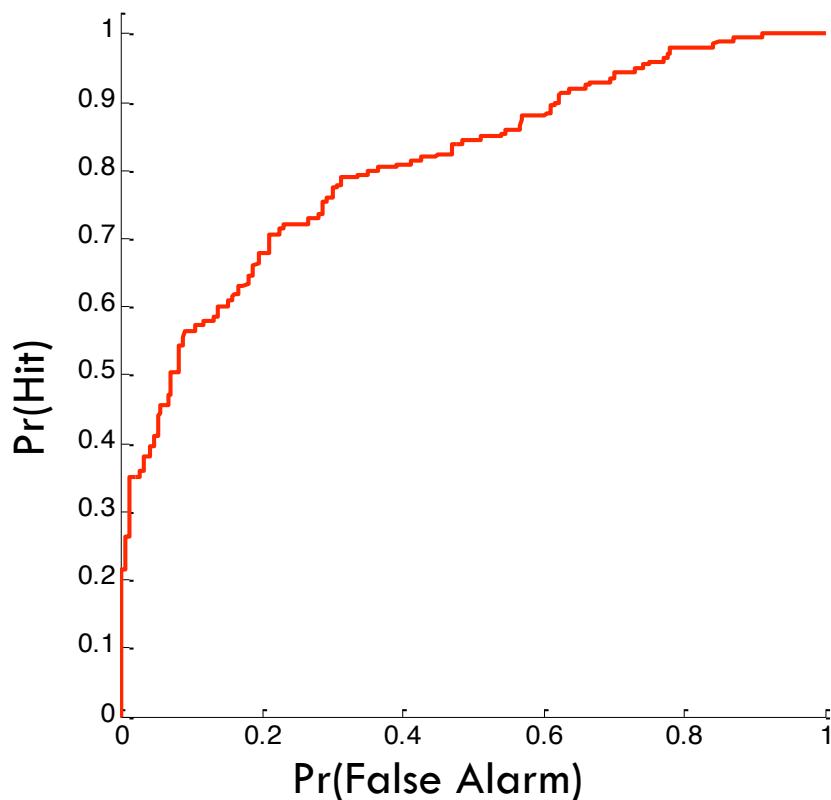


Model 1 Results

65

Probability & Bayesian Inference

Results based on 200 cropped faces and 200 non-faces from the same database.



How does this work with a real image?



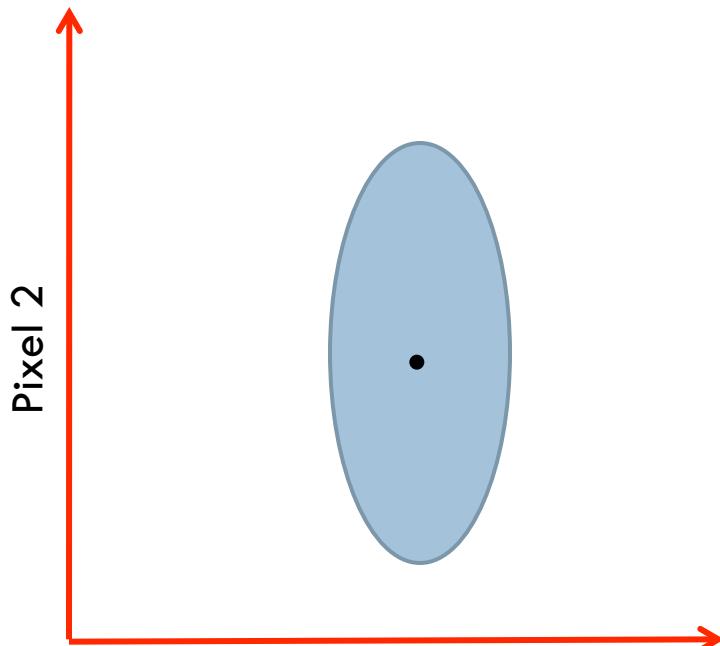
Model # 2: Gaussian, diagonal covariance

66

Probability & Bayesian Inference

$$Pr(\mathbf{x}|\text{face}) = \frac{1}{(2\pi)^{d/2}|\Sigma|} \exp \left\{ -0.5(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

Fit model using maximum likelihood criterion



\mathbf{m} face



\mathbf{m} non-face



\mathbf{s} Face



\mathbf{s} non-face

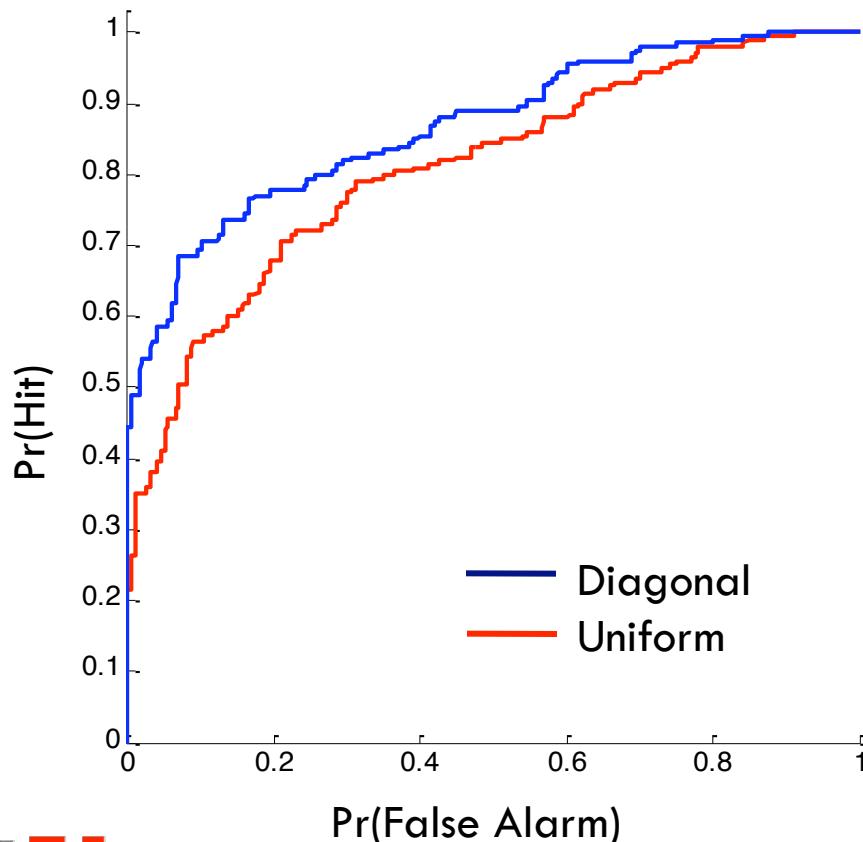


Model 2 Results

67

Probability & Bayesian Inference

Results based on 200 cropped faces and 200 non-faces from the same database.



More sophisticated
model unsurprisingly
classifies new faces
and non-faces better.

Model # 3: Gaussian, full covariance

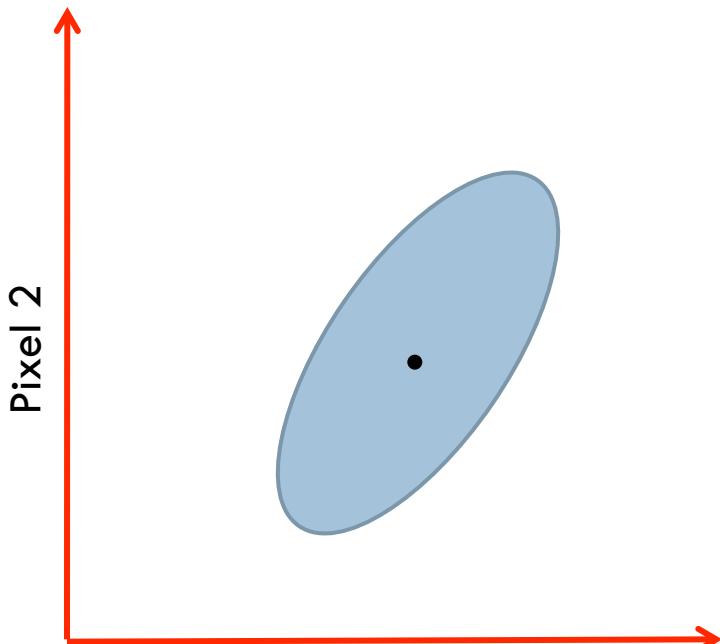
68

Probability & Bayesian Inference

$$Pr(\mathbf{x}|\text{face}) = \frac{1}{(2\pi)^{d/2}|\Sigma|} \exp \left\{ -0.5(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

Fit model using maximum likelihood criterion

PROBLEM: we cannot fit this model. We don't have enough data to estimate the full covariance matrix.



N=800 training images
D=10800 dimensions

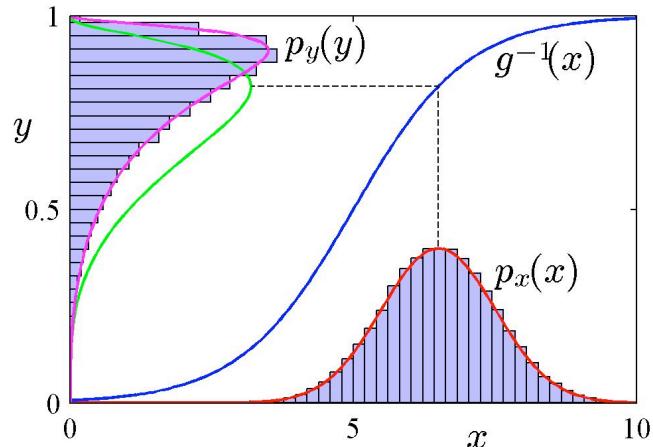
Total number of measured numbers =
ND = $800 \times 10,800 = 8,640,000$

Total number of parameters in cov matrix =
 $(D+1)D/2 = (10,800+1) \times 10,800 / 2 = 58,325,400$

Transformed Densities Revisited

69

Probability & Bayesian Inference



Observations falling within $(x + \delta x)$ transform to the range $(y + \delta y)$

$$\rightarrow p_x(x)|\delta x| = p_y(y)|\delta y|$$

$$\rightarrow p_y(y) \approx p_x(x) \left| \frac{\delta x}{\delta y} \right|$$

Note that in general, $\delta y \neq \delta x$.

Rather, $\frac{\delta y}{\delta x} \rightarrow \frac{dy}{dx}$ as $\delta x \rightarrow 0$.

$$\text{Thus } p_y(y) = p_x(x) \left| \frac{dx}{dy} \right|$$

Problems for this week

70

Probability & Bayesian Inference

- Problems 2.7 – 2.17, 2.19 – 2.21, 2.23 – 2.27 are all good!
 - At least do problem 2.14. We will talk about this Monday. (Hopefully one of you will present a solution!)
- Also, MATLAB exercises up to 1.4.4 are good.
 - At least do Exercise 1.4.2. We will talk about this next week.

Bayesian Decision Theory: Topics

71

Probability & Bayesian Inference

1. Probability
2. The Univariate Normal Distribution
3. Bayesian Classifiers
4. Minimizing Risk
5. The Multivariate Normal Distribution
6. **Decision Boundaries in Higher Dimensions**
7. Parameter Estimation
8. Mixture Models and EM
9. Nonparametric Density Estimation
10. Training and Evaluation Methods
11. What are Bayes Nets?

Topic 6.

Decision Boundaries in Higher Dimensions

Decision Surfaces

73

Probability & Bayesian Inference

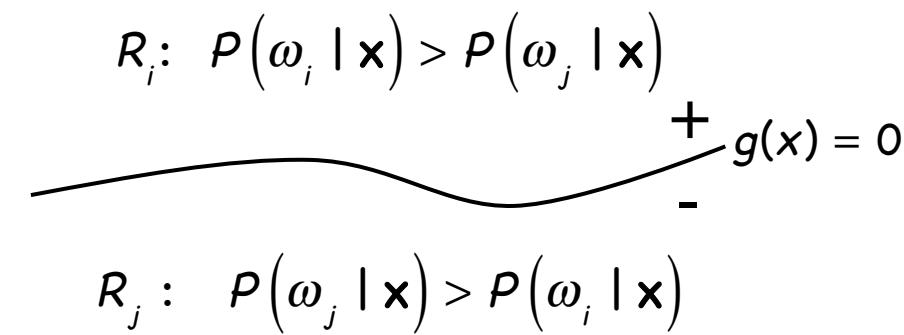
- If decision regions R_i and R_j are contiguous, define

$$g(\mathbf{x}) \equiv P(\omega_i | \mathbf{x}) - P(\omega_j | \mathbf{x})$$

- Then the decision surface

$$g(\mathbf{x}) = 0$$

separates the two decision regions. $g(\mathbf{x})$ is positive on one side and negative on the other.



Discriminant Functions

74

Probability & Bayesian Inference

- If $f(\cdot)$ monotonic, the rule remains the same if we use:

$$\underline{x} \rightarrow \omega_i \text{ if: } f(P(\omega_i | \underline{x})) > f(P(\omega_j | \underline{x})) \quad \forall i \neq j$$

- $g_i(\underline{x}) \equiv f(P(\omega_i | \underline{x}))$ is a **discriminant function**
- In general, discriminant functions can be defined in other ways, independent of Bayes.
- In theory this will lead to a suboptimal solution
- However, non-Bayesian classifiers can have significant advantages:
 - Often a full Bayesian treatment is intractable or computationally prohibitive.
 - Approximations made in a Bayesian treatment may lead to errors avoided by non-Bayesian methods.



End of Lecture 4

Multivariate Normal Likelihoods

76

Probability & Bayesian Inference

- Multivariate Gaussian pdf

$$p(\underline{x} | \omega_i) = \frac{1}{(2\pi)^{\frac{\ell}{2}} |\Sigma_i|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (\underline{x} - \underline{\mu}_i)^T \Sigma_i^{-1} (\underline{x} - \underline{\mu}_i)\right)$$

$$\underline{\mu}_i = E[\underline{x}]$$

$$\Sigma_i = E[(\underline{x} - \underline{\mu}_i)(\underline{x} - \underline{\mu}_i)^T]$$

called the **covariance matrix**

Logarithmic Discriminant Function

77

Probability & Bayesian Inference

$$p(\underline{x} | \omega_i) = \frac{1}{(2\pi)^{\frac{\ell}{2}} |\Sigma_i|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (\underline{x} - \underline{\mu}_i)^T \Sigma_i^{-1} (\underline{x} - \underline{\mu}_i)\right)$$

- $\ln(\cdot)$ is monotonic. Define:

$$g_i(\underline{x}) = \ln(p(\underline{x} | \omega_i) P(\omega_i)) = \ln p(\underline{x} | \omega_i) + \ln P(\omega_i)$$

$$= -\frac{1}{2} (\underline{x} - \underline{\mu}_i)^T \Sigma_i^{-1} (\underline{x} - \underline{\mu}_i) + \ln P(\omega_i) + C_i$$

where

$$C_i = -\frac{\ell}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_i|$$

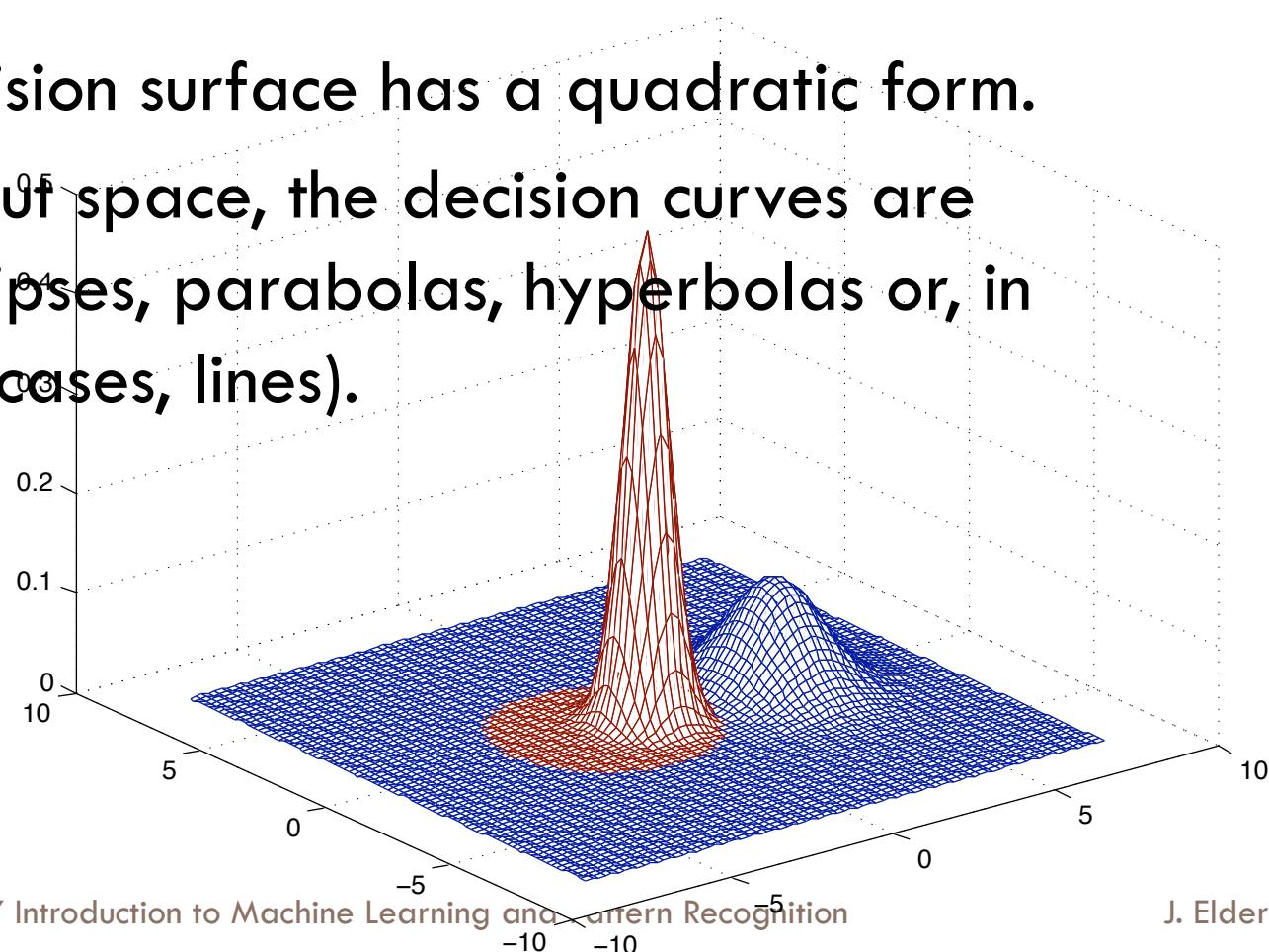
Quadratic Classifiers

78

Probability & Bayesian Inference

$$g_i(\underline{x}) = -\frac{1}{2}(\underline{x} - \underline{\mu}_i)^T \Sigma_i^{-1} (\underline{x} - \underline{\mu}_i) + \ln P(\omega_i) + C_i$$

- Thus the decision surface has a quadratic form.
- For a 2D input space, the decision curves are quadrics (ellipses, parabolas, hyperbolas or, in degenerate cases, lines).



Example: Isotropic Likelihoods

79

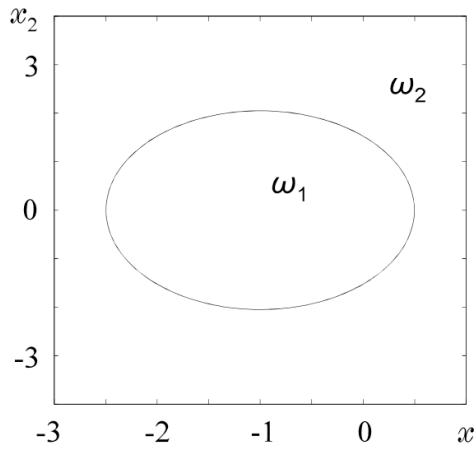
Probability & Bayesian Inference

$$g_i(\underline{x}) = -\frac{1}{2}(\underline{x} - \underline{\mu}_i)^T \Sigma_i^{-1}(\underline{x} - \underline{\mu}_i) + \ln P(\omega_i) + C_i$$

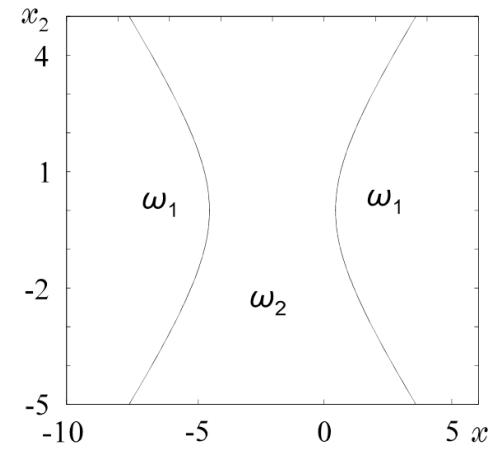
- Suppose that the two likelihoods are both isotropic, but with different means and variances. Then

$$g_i(\underline{x}) = -\frac{1}{2\sigma_i^2}(x_1^2 + x_2^2) + \frac{1}{\sigma_i^2}(\mu_{i1}x_1 + \mu_{i2}x_2) - \frac{1}{2\sigma_i^2}(\mu_{i1}^2 + \mu_{i2}^2) + \ln(P(\omega_i)) + C_i$$

- And $g_i(\underline{x}) - g_j(\underline{x}) = 0$ will be a quadratic equation in 2 variables.



(a)



(b)

Equal Covariances

80

Probability & Bayesian Inference

$$g_i(\underline{x}) = -\frac{1}{2}(\underline{x} - \underline{\mu}_i)^T \Sigma_i^{-1} (\underline{x} - \underline{\mu}_i) + \ln P(\omega_i) + C_i$$

- The quadratic term of the decision boundary is given by

$$\frac{1}{2} \mathbf{x}^T \left(\Sigma_j^{-1} - \Sigma_i^{-1} \right) \mathbf{x}$$

- Thus if the covariance matrices of the two likelihoods are identical, the decision boundary is linear.

Linear Classifier

81

Probability & Bayesian Inference

$$g_i(\underline{x}) = -\frac{1}{2}(\underline{x} - \underline{\mu}_i)^T \Sigma^{-1}(\underline{x} - \underline{\mu}_i) + \ln P(\omega_i) + C_i$$

- In this case, we can drop the quadratic terms and express the discriminant function in linear form:

$$g_i(\underline{x}) = \underline{w}_i^T \underline{x} + w_{i0}$$

$$\underline{w}_i = \Sigma^{-1} \underline{\mu}_i$$

$$w_{i0} = \ln P(\omega_i) - \frac{1}{2} \underline{\mu}_i^T \Sigma^{-1} \underline{\mu}_i$$

Example 1: Isotropic, Identical Variance

82

Probability & Bayesian Inference

$$g_i(\underline{x}) = \underline{w}_i^\top \underline{x} + w_{i0}$$

$$\underline{w}_i = \Sigma^{-1} \underline{\mu}_i$$

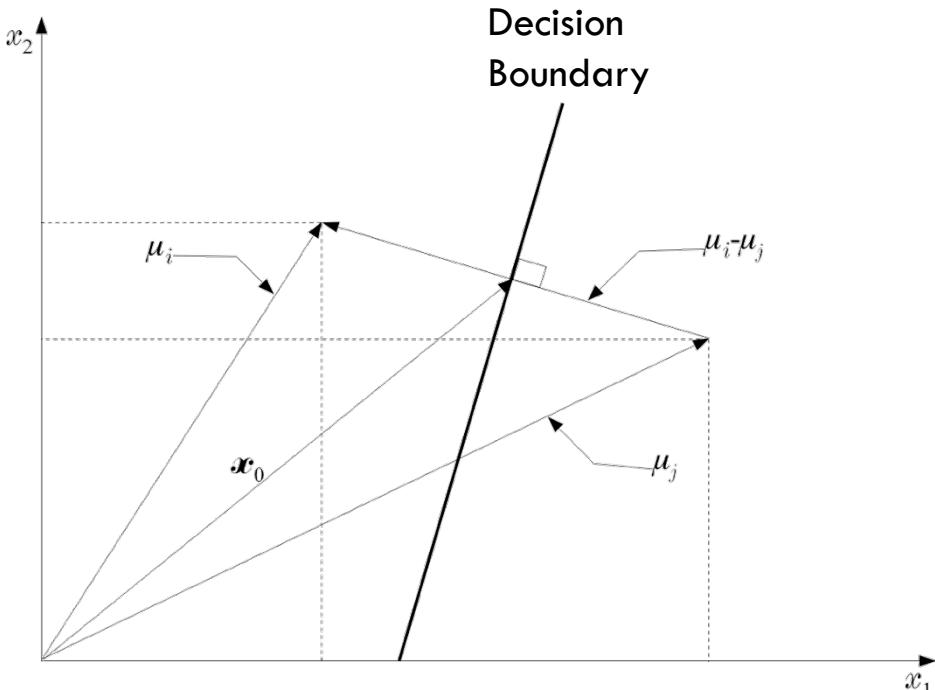
$$w_{i0} = \ln P(\omega_i) - \frac{1}{2} \underline{\mu}_i^\top \Sigma^{-1} \underline{\mu}_i$$

$\Sigma = \sigma^2 I$. Then

$$\underline{w}^\top (\underline{x} - \underline{x}_o) = 0, \text{ where}$$

$\underline{w} = \underline{\mu}_i - \underline{\mu}_j$, and

$$\underline{x}_o = \frac{1}{2} (\underline{\mu}_i + \underline{\mu}_j) - \sigma^2 \ln \frac{P(\omega_i)}{P(\omega_j)} \frac{\underline{\mu}_i - \underline{\mu}_j}{\|\underline{\mu}_i - \underline{\mu}_j\|^2}$$



Example 2: Equal Covariance

83

Probability & Bayesian Inference

$$g_i(\underline{x}) = \underline{w}_i^T \underline{x} + w_{i0}$$

$$\underline{w}_i = \Sigma^{-1} \underline{\mu}_i$$

$$\underline{w}_{i0} = \ln P(\omega_i) - \frac{1}{2} \underline{\mu}_i^T \Sigma^{-1} \underline{\mu}_i$$

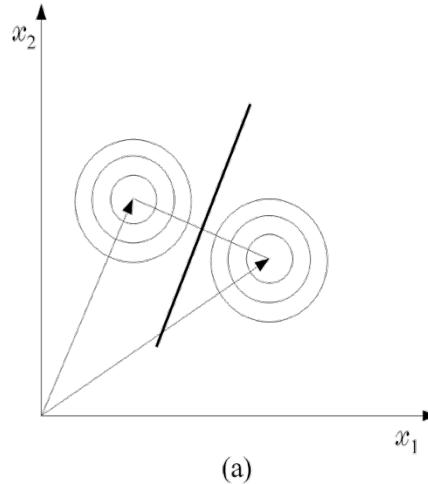
$$g_{ij}(\underline{x}) = \underline{w}^T (\underline{x} - \underline{x}_0) = 0 \text{ where}$$

$$\underline{w} = \Sigma^{-1}(\underline{\mu}_i - \underline{\mu}_j),$$

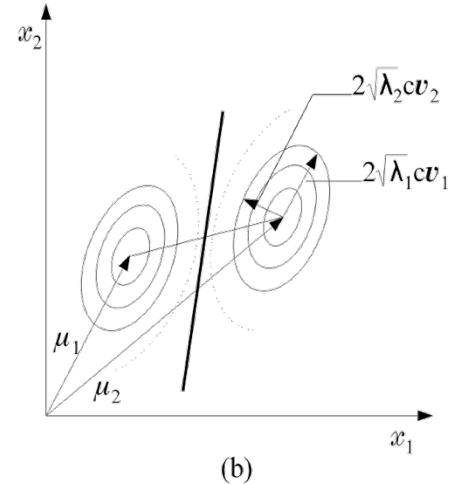
$$\underline{x}_0 = \frac{1}{2}(\underline{\mu}_i + \underline{\mu}_j) - \ln \left(\frac{P(\omega_i)}{P(\omega_j)} \right) \frac{\underline{\mu}_i - \underline{\mu}_j}{\left\| \underline{\mu}_i - \underline{\mu}_j \right\|_{\Sigma^{-1}}^2},$$

and

$$\left\| \underline{x} \right\|_{\Sigma^{-1}} \equiv (\underline{x}^T \Sigma^{-1} \underline{x})^{\frac{1}{2}}$$



(a)



(b)

Minimum Distance Classifiers

84

Probability & Bayesian Inference

- If the two likelihoods have identical covariance AND the two classes are equiprobable, the discrimination function simplifies:

$$g_i(\underline{x}) = -\frac{1}{2}(\underline{x} - \underline{\mu}_i)^T \Sigma_i^{-1}(\underline{x} - \underline{\mu}_i) + \ln P(\omega_i) + C_i$$



$$g_i(\underline{x}) = -\frac{1}{2}(\underline{x} - \underline{\mu}_i)^T \Sigma^{-1}(\underline{x} - \underline{\mu}_i)$$

Isotropic Case

85

Probability & Bayesian Inference

- In the isotropic case,

$$g_i(\underline{x}) = -\frac{1}{2}(\underline{x} - \underline{\mu}_i)^T \Sigma^{-1}(\underline{x} - \underline{\mu}_i) = -\frac{1}{2\sigma^2} \|\underline{x} - \underline{\mu}_i\|^2$$

- Thus the Bayesian classifier simply assigns the class that minimizes the Euclidean distance d_e between the observed feature vector and the class mean.

$$d_e = \|\underline{x} - \underline{\mu}_i\|$$

General Case: Mahalanobis Distance

- To deal with anisotropic distributions, we simply classify according to the Mahalanobis distance, defined as

$$d_m = g_i(\underline{x}) = \left((\underline{x} - \underline{\mu}_i)^T \Sigma^{-1} (\underline{x} - \underline{\mu}_i) \right)^{1/2}$$

- Since the covariance matrix is symmetric, it can be represented as $\Sigma = \Phi \Lambda \Phi^T$

where the columns v_i of Φ are the eigenvectors of Σ and where

Λ is a diagonal matrix whose diagonal elements λ_i are the corresponding eigenvalues.

- Then we have

$$d_m^2 = (\underline{x} - \underline{\mu}_i)^T \Phi^T \Lambda^{-1} \Phi (\underline{x} - \underline{\mu}_i)$$

General Case: Mahalanobis Distance

87

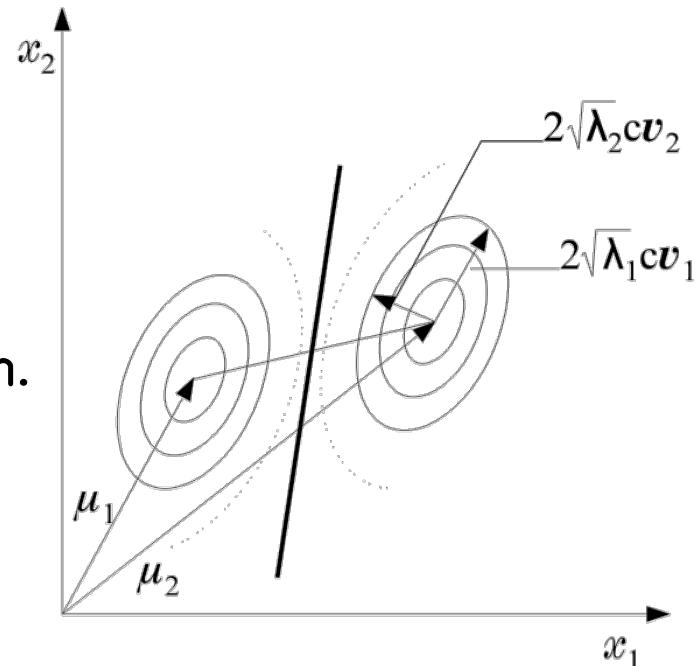
Probability & Bayesian Inference

$$d_m^2 = (\underline{x} - \underline{\mu}_i)^T \Phi^T \Lambda^{-1} \Phi (\underline{x} - \underline{\mu}_i)$$

Let $\mathbf{x}' = \Phi^T \mathbf{x}$. Then the coordinates of \mathbf{x}' are the projections of \mathbf{x} onto the eigenvectors of Σ , and we have:

$$\frac{(\mathbf{x}'_1 - \mu'_{i1})^2}{\lambda_1} + \dots + \frac{(\mathbf{x}'_l - \mu'_{il})^2}{\lambda_l} = d_m^2$$

Thus the curves of constant
Mahalanobis distance c have ellipsoidal form.



Example:

88

Probability & Bayesian Inference

Given ω_1, ω_2 : $P(\omega_1) = P(\omega_2)$ and $p(\underline{x}|\omega_1) = N(\underline{\mu}_1, \Sigma)$, $p(\underline{x}|\omega_2) = N(\underline{\mu}_2, \Sigma)$,

$$\underline{\mu}_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \quad \underline{\mu}_2 = \begin{bmatrix} 3 \\ 3 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} 1.1 & 0.3 \\ 0.3 & 1.9 \end{bmatrix}$$

classify the vector $\underline{x} = \begin{bmatrix} 1.0 \\ 2.2 \end{bmatrix}$ using Bayesian classification:

- $\Sigma^{-1} = \begin{bmatrix} 0.95 & -0.15 \\ -0.15 & 0.55 \end{bmatrix}$

- Compute Mahalanobis d_m from μ_1, μ_2 :

$$d_{m,1}^2 = \left[1.0, \ 2.2 \right] \Sigma^{-1} \begin{bmatrix} 1.0 \\ 2.2 \end{bmatrix} = 2.952, \ d_{m,2}^2 = \left[-2.0, \ -0.8 \right] \Sigma^{-1} \begin{bmatrix} -2.0 \\ -0.8 \end{bmatrix} = 3.672$$

- Classify $\underline{x} \rightarrow \omega_1$. Observe that $d_{E,2} < d_{E,1}$

Bayesian Decision Theory: Topics

89

Probability & Bayesian Inference

1. Probability
2. The Univariate Normal Distribution
3. Bayesian Classifiers
4. Minimizing Risk
5. The Multivariate Normal Distribution
6. Decision Boundaries in Higher Dimensions
7. **Parameter Estimation**
8. Mixture Models and EM
9. Nonparametric Density Estimation
10. Training and Evaluation Methods
11. What are Bayes Nets?

Topic 7. Parameter Estimation

Topic 7.1 Maximum Likelihood Estimation

Maximum Likelihood Parameter Estimation

92

Probability & Bayesian Inference

Suppose we believe input vectors \underline{x} are distributed as $p(\underline{x}) \equiv p(\underline{x}; \underline{\theta})$, where $\underline{\theta}$ is an unknown parameter.

Given independent training input vectors $X = \{\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N\}$

we want to compute the maximum likelihood estimate $\underline{\theta}_{ML}$ for $\underline{\theta}$.

Since the input vectors are independent, we have

$$p(X; \underline{\theta}) \equiv p(\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N; \underline{\theta}) = \prod_{k=1}^N p(\underline{x}_k; \underline{\theta})$$

Maximum Likelihood Parameter Estimation

93

Probability & Bayesian Inference

$$p(\underline{X}; \underline{\theta}) = \prod_{k=1}^N p(\underline{x}_k; \underline{\theta})$$

$$\text{Let } L(\underline{\theta}) \equiv \ln p(\underline{X}; \underline{\theta}) = \sum_{k=1}^N \ln p(\underline{x}_k; \underline{\theta})$$

The general method is to take the derivative of L with respect to $\underline{\theta}$, set it to 0 and solve for $\underline{\theta}$:

$$\hat{\underline{\theta}}_{ML} : \frac{\partial L(\underline{\theta})}{\partial (\underline{\theta})} = \sum_{k=1}^N \frac{\partial \ln p(\underline{x}_k; \underline{\theta})}{\partial (\underline{\theta})} = \underline{0}$$

Properties of the Maximum Likelihood Estimator

94

Probability & Bayesian Inference

Let $\underline{\theta}_0$ be the true value of the unknown parameter vector.
Then

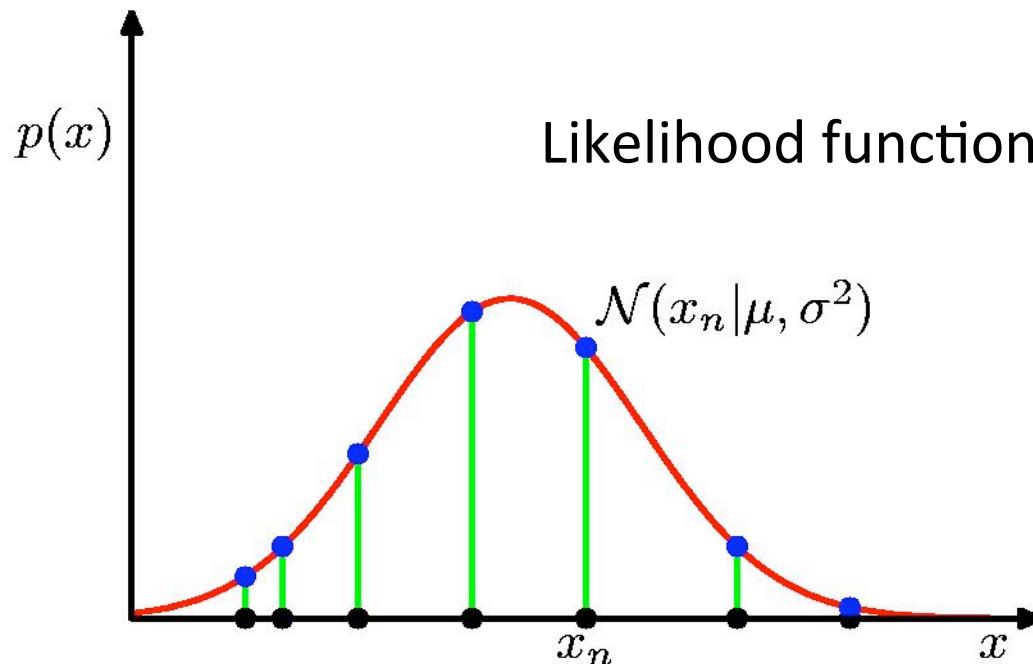
$\underline{\theta}_{ML}$ is **asymptotically unbiased**: $\lim_{N \rightarrow \infty} E[\underline{\theta}_{ML}] = \underline{\theta}_0$

$\underline{\theta}_{ML}$ is **asymptotically consistent**: $\lim_{N \rightarrow \infty} E \left\| \hat{\underline{\theta}}_{ML} - \underline{\theta}_0 \right\|^2 = 0$

Example: Univariate Normal

95

Probability & Bayesian Inference



$$p(\mathbf{x}|\mu, \sigma^2) = \prod_{n=1}^N \mathcal{N}(x_n|\mu, \sigma^2)$$

Example: Univariate Normal

96

Probability & Bayesian Inference

$$\ln p(\mathbf{x}|\mu, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln(2\pi)$$

$$\mu_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N x_n \quad \sigma_{\text{ML}}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{\text{ML}})^2$$

Example: Univariate Normal

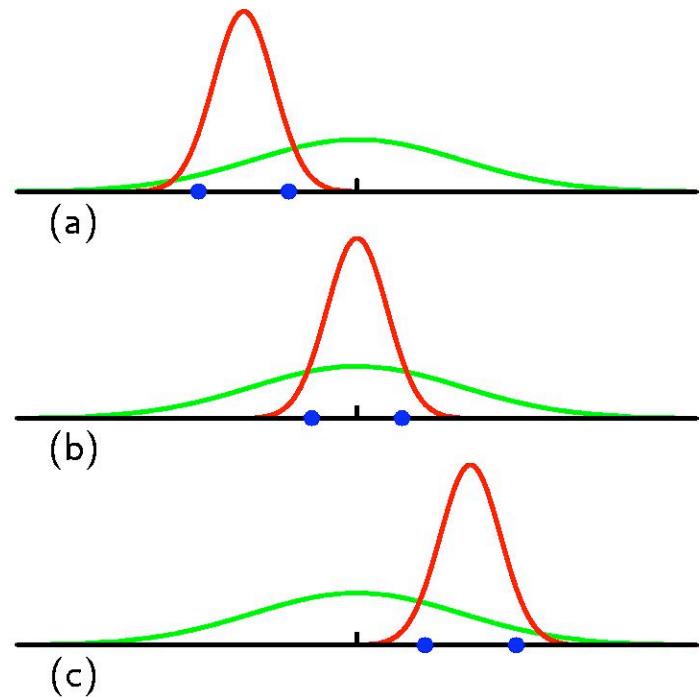
97

Probability & Bayesian Inference

$$\mathbb{E}[\mu_{\text{ML}}] = \mu$$

$$\mathbb{E}[\sigma_{\text{ML}}^2] = \left(\frac{N-1}{N} \right) \sigma^2$$

$$\begin{aligned}\tilde{\sigma}^2 &= \frac{N}{N-1} \sigma_{\text{ML}}^2 \\ &= \frac{1}{N-1} \sum_{n=1}^N (x_n - \mu_{\text{ML}})^2\end{aligned}$$



Thus σ_{ML} is biased (although asymptotically unbiased).



End of Lecture 5

Example: Multivariate Normal

99

Probability & Bayesian Inference

- Given i.i.d. data $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$, the log likelihood function is given by

$$\ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln |\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu})$$

Maximum Likelihood for the Gaussian

100

Probability & Bayesian Inference

- Set the derivative of the log likelihood function to zero,

$$\frac{\partial}{\partial \mu} \ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) = 0$$

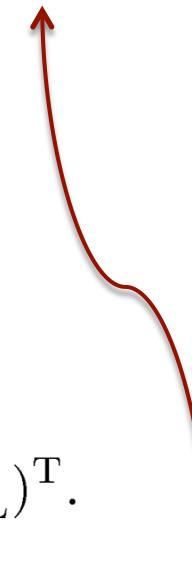
- and solve to obtain

$$\boldsymbol{\mu}_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n.$$

- One can also show that

$$\boldsymbol{\Sigma}_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}_{\text{ML}})(\mathbf{x}_n - \boldsymbol{\mu}_{\text{ML}})^T.$$

Recall: If \mathbf{x} and \mathbf{a} are vectors, then $\frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^T \mathbf{a}) = \frac{\partial}{\partial \mathbf{x}} (\mathbf{a}^T \mathbf{x}) = \mathbf{a}$



Topic 7.1 Bayesian Parameter Estimation

Bayesian Inference for the Gaussian (Univariate Case)

102

Probability & Bayesian Inference

- Assume σ^2 is known. Given i.i.d. data

$\mathbf{x} = \{x_1, \dots, x_N\}$, the likelihood function for μ is given by

$$p(\mathbf{x}|\mu) = \prod_{n=1}^N p(x_n|\mu) = \frac{1}{(2\pi\sigma^2)^{N/2}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2\right\}.$$

- This has a Gaussian shape as a function of μ (but it is *not* a distribution over μ).

Bayesian Inference for the Gaussian (Univariate Case)

103

Probability & Bayesian Inference

- Combined with a Gaussian prior over μ ,

$$p(\mu) = \mathcal{N}(\mu|\mu_0, \sigma_0^2).$$

- this gives the posterior

$$p(\mu|\mathbf{x}) \propto p(\mathbf{x}|\mu)p(\mu).$$

- Completing the square over μ , we see that

$$p(\mu|\mathbf{x}) = \mathcal{N}(\mu|\mu_N, \sigma_N^2)$$

Bayesian Inference for the Gaussian

104

Probability & Bayesian Inference

□ ... where

$$\mu_N = \frac{\sigma^2}{N\sigma_0^2 + \sigma^2} \mu_0 + \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2} \mu_{\text{ML}}, \quad \mu_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N x_n$$

$$\frac{1}{\sigma_N^2} = \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2}.$$

Shortcut: $p(\mu | X)$ has the form $C \exp(-\Delta^2)$.

Get Δ^2 in form $a\mu^2 - 2b\mu + c = a(\mu - b/a)^2 + \text{const}$ and identify

$$\mu_N = b/a$$

$$\frac{1}{\sigma_N^2} = a$$

□ Note:

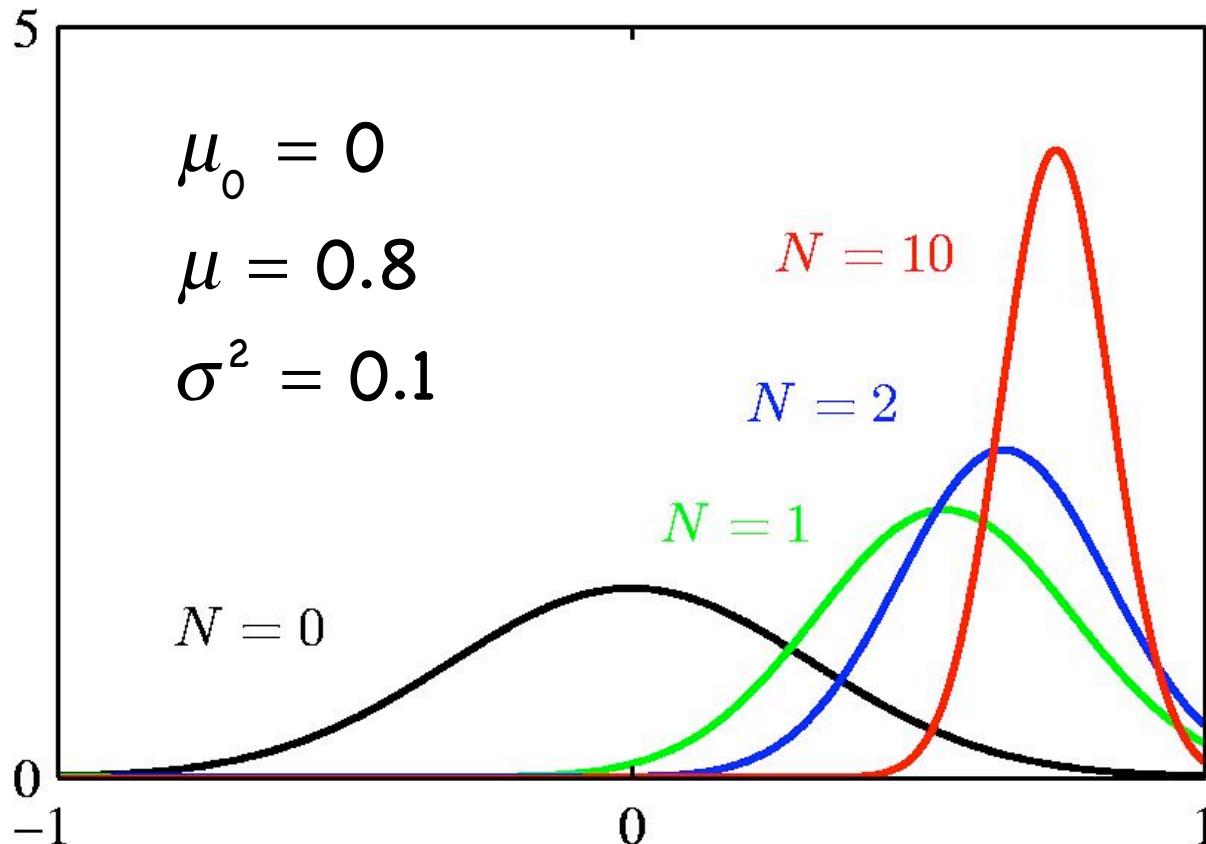
	$N = 0$	$N \rightarrow \infty$
μ_N	μ_0	μ_{ML}
σ_N^2	σ_0^2	0

Bayesian Inference for the Gaussian

105

Probability & Bayesian Inference

□ Example: $p(\mu|x) = \mathcal{N}(\mu|\mu_N, \sigma_N^2)$



Maximum a Posteriori (MAP) Estimation

106

Probability & Bayesian Inference

$$p(\mu|\mathbf{x}) = \mathcal{N}(\mu|\mu_N, \sigma_N^2)$$

$$\mu_N = \frac{\sigma^2}{N\sigma_0^2 + \sigma^2}\mu_0 + \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2}\mu_{ML}, \quad \mu_{ML} = \frac{1}{N} \sum_{n=1}^N x_n$$

$$\frac{1}{\sigma_N^2} = \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2}.$$

In MAP estimation, we use the value of μ that maximizes the posterior $p(\mu | X)$:

$$\mu_{MAP} = \mu_N.$$

Full Bayesian Parameter Estimation

107

Probability & Bayesian Inference

- In both ML and MAP, we use the training data \mathbf{X} to estimate a specific value for the unknown parameter vector $\underline{\theta}$, and then use that value for subsequent inference on new observations \mathbf{x} : $p(\mathbf{x} | \underline{\theta})$
- These methods are suboptimal, because in fact we are always uncertain about the exact value of $\underline{\theta}$, and to be optimal we should take into account the possibility that $\underline{\theta}$ assumes other values.

Full Bayesian Parameter Estimation

108

Probability & Bayesian Inference

- In full Bayesian parameter estimation, we do not estimate a specific value for $\underline{\theta}$.
- Instead, we compute the posterior over $\underline{\theta}$, and then integrate it out when computing $p(\underline{x} | \mathcal{X})$:

$$p(\underline{x} | \mathcal{X}) = \int p(\underline{x} | \underline{\theta}) p(\underline{\theta} | \mathcal{X}) d\underline{\theta}$$

$$p(\underline{\theta} | \mathcal{X}) = \frac{p(\mathcal{X} | \underline{\theta}) p(\underline{\theta})}{\int p(\mathcal{X} | \underline{\theta}) p(\underline{\theta}) d\underline{\theta}}$$

$$p(\mathcal{X} | \underline{\theta}) = \prod_{k=1}^N p(\underline{x}_k | \underline{\theta})$$

Example: Univariate Normal with Unknown Mean

109

Probability & Bayesian Inference

Consider again the case $p(\underline{x}|\mu) \sim N(\mu, \sigma^2)$ where σ is known and $\mu \sim N(\mu_0, \sigma_0^2)$

We showed that $p(\mu|\underline{X}) \sim N(\mu_N, \sigma_N^2)$, where

$$\begin{aligned}\mu_N &= \frac{\sigma^2}{N\sigma_0^2 + \sigma^2}\mu_0 + \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2}\mu_{\text{ML}}, & \mu_{\text{ML}} &= \frac{1}{N} \sum_{n=1}^N x_n \\ \frac{1}{\sigma_N^2} &= \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2}.\end{aligned}$$

In the MAP approach, we approximate $p(\underline{x}|\underline{X}) \sim N(\mu_N, \sigma^2)$

In the full Bayesian approach, we calculate $p(\underline{x}|\underline{X}) = \int p(\underline{x}|\mu)p(\mu|\underline{X})d\mu$

which can be shown to yield $p(\underline{x}|\underline{X}) \sim N(\mu_N, \sigma^2 + \sigma_N^2)$

Hints for Exercise 1.4.2

110

Probability & Bayesian Inference

- Here are some MATLAB functions you may find useful in solving Exercise 1.4.2
 - `mnrnd`
 - `mvnrnd`
 - `mvnpdf`
 - `mean`
 - `cov`
 - `squeeze`
 - `sum`
 - `repmat`
 - `inv`
 - `min`
 - `max`
 - `zeros`
 - `ones`

Problem 1.4.2

111

Probability & Bayesian Inference

```
function pe=pr142
%Exercise 1.4.2 from PR Matlab Manual

m(:,1)=[0 0 0]';
m(:,2)=[1 2 2]';
m(:,3)=[3 3 4]';
S=[0.8 0.2 0.1;0.2 0.8 0.2; 0.1 0.2 0.8];
N=1000;

%Part 1
ntrain=mrnd(N,ones(3,1)/3); %Number of training pts generated by each
class
ntest=mrnd(N,ones(3,1)/3); %Number of test pts generated by each class
test=[];
mml=zeros(3,3);
Smli=zeros(3,3,3);
for i=1:3
    train=mvnrnd(m(:,i),S,ntrain(i)); %training vectors from class i
    test=[test;mvnrnd(m(:,i),S,ntest(i))]; %test vectors from class i

    mml(i,:)=mean(train); %ML estimate of mean for class i
    Smli(i,:,:)=ntest(i)*cov(train,1);%weighted ML estimate of
covariance for class i
end

Sml=squeeze(mean(Smli)/N); %ML estimate of common covariance

%Part 2: Euclidean distance
for i=1:3
    dsq(:,i)=sum((test-repmat(mml(i,:),N,1)).^2,2);
end
[m,idx(:,1)]=min(dsq');

%Part 3: Mahalanobis distance
for i=1:3
    y=test-repmat(mml(i,:),N,1);
    dsq(:,i)=sum((y*inv(Sml)).*y,2);
end
[m,idx(:,2)]=min(dsq');

%Part 4: Maximum likelihood classifier
for i=1:3
    p(:,i)=mvnpdf(test,mml(i,:),Sml);
end
[m,idx(:,3)]=max(p');

%Ground truth classes
idxgt=[ones(ntest(1),1);2*ones(ntest(2),1);3*ones(ntest(3),1)];

for i=1:3
    pe(i)=mean(idx(:,i)~=idxgt); %Error rate for class i
end
```

Bayesian Decision Theory: Topics

112

Probability & Bayesian Inference

1. Probability
2. The Univariate Normal Distribution
3. Bayesian Classifiers
4. Minimizing Risk
5. The Multivariate Normal Distribution
6. Decision Boundaries in Higher Dimensions
7. Parameter Estimation
8. **Mixture Models and EM**
9. Nonparametric Density Estimation
10. Training and Evaluation Methods
11. What are Bayes Nets?

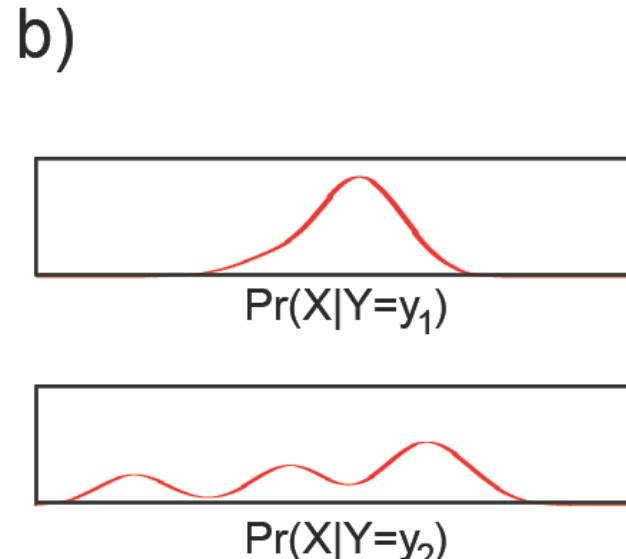
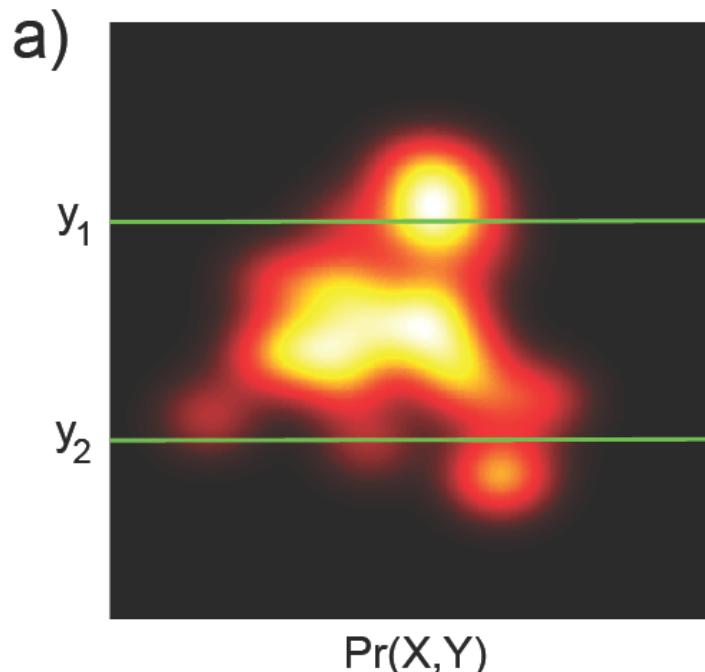
Topic 8 Mixture Models and EM

Motivation

114

Probability & Bayesian Inference

- What do we do if a distribution is not well-approximated by a standard parametric model?



8.1 Intuition

Mixtures of Gaussians

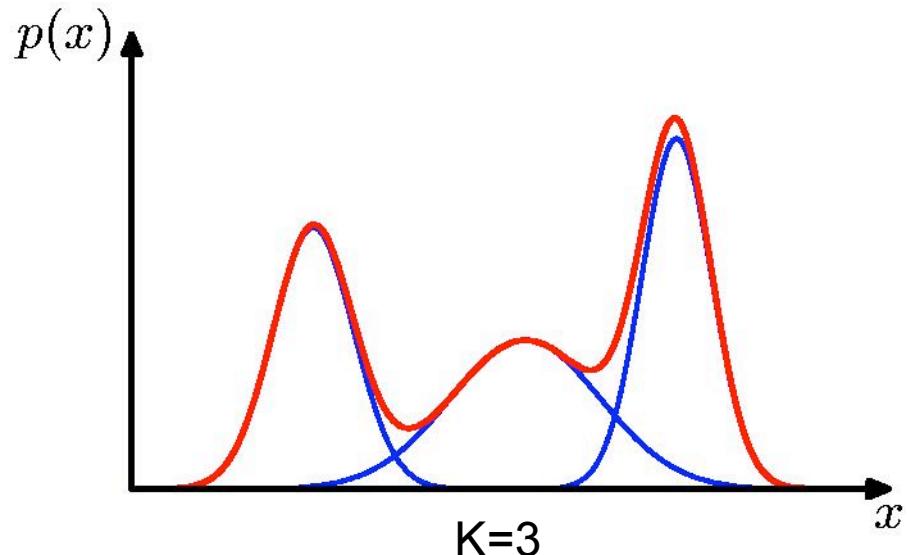
116

Probability & Bayesian Inference

- Combine simple models into a complex model:

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

↑
Component
Mixing coefficient

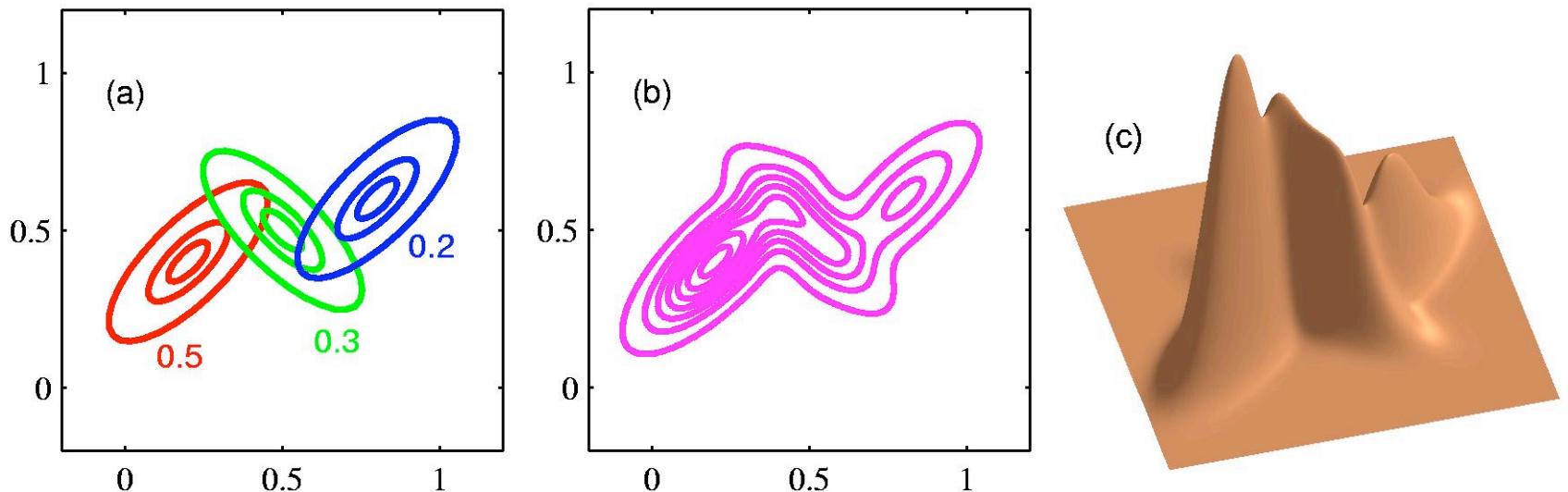


$$\forall k : \pi_k \geq 0 \quad \sum_{k=1}^K \pi_k = 1$$

Mixtures of Gaussians

117

Probability & Bayesian Inference



Mixtures of Gaussians

118

Probability & Bayesian Inference

- Determining parameters μ , σ and π using maximum log likelihood

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left\{ \underbrace{\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}_{\text{Log of a sum; no closed form maximum.}} \right\}$$

Log of a sum; no closed form maximum.

- Solution: use standard, iterative, numeric optimization methods or the **expectation maximization** algorithm.



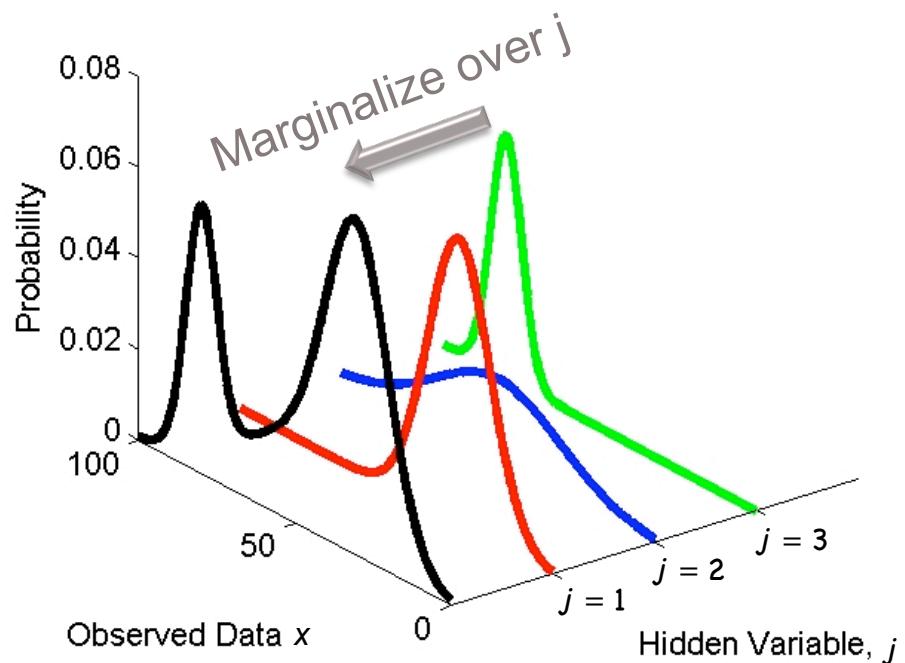
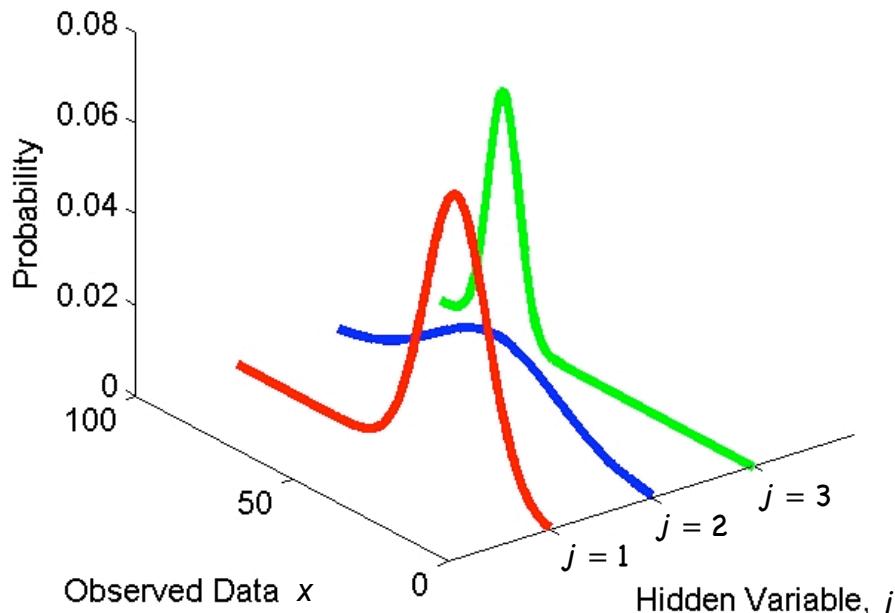
End of Lecture 6

Hidden Variable Interpretation

120

Probability & Bayesian Inference

$$p(x | P_1, \dots, P_J, \mu_1, \dots, \mu_J, \sigma_1, \dots, \sigma_J) = \sum_{j=1}^J P_j N(x; \mu_j, \sigma_j^2) = \sum_{j=1}^J p(j) p(x | j)$$



Hidden Variable Interpretation

121

Probability & Bayesian Inference

$$p(x | P_1, \dots, P_J, \mu_1, \dots, \mu_J, \sigma_1^2, \dots, \sigma_J^2) = \sum_{j=1}^J P_j N(x; \mu_j, \sigma_j^2) = \sum_{j=1}^J p(j) p(x | j)$$

ASSUMPTIONS

- for each training datum x_i there is a hidden variable j_i .
- j_i represents which Gaussian x_i came from
- hence j_i takes discrete values

OUR GOAL:

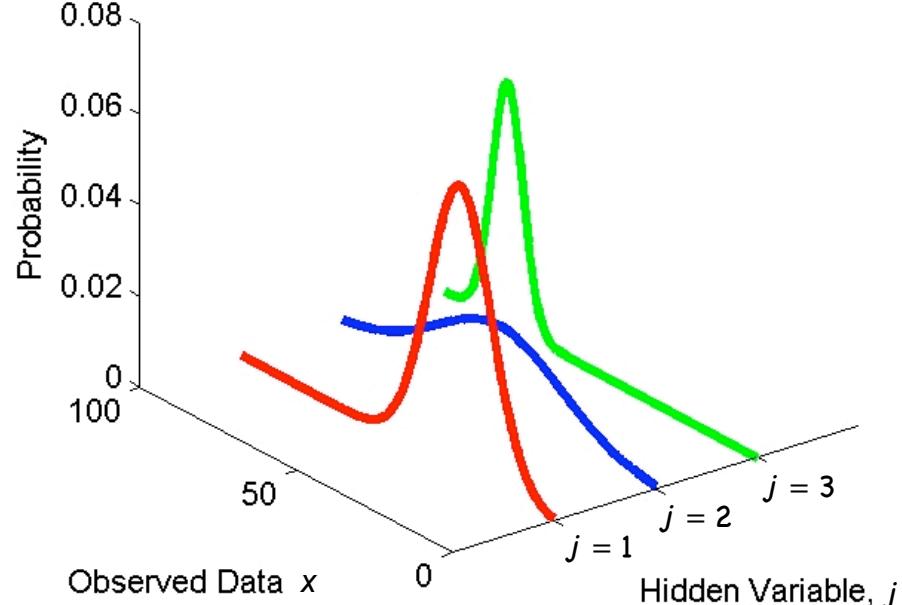
To estimate the parameters θ :

The means μ_j

The covariances Σ_j

The weights (mixing coefficients) P_j

for all J components of the model.



THING TO NOTICE:

If we knew the hidden variables j_i for the training data it would be easy to estimate parameters θ – just estimate individual Gaussians separately.

Hidden Variable Interpretation

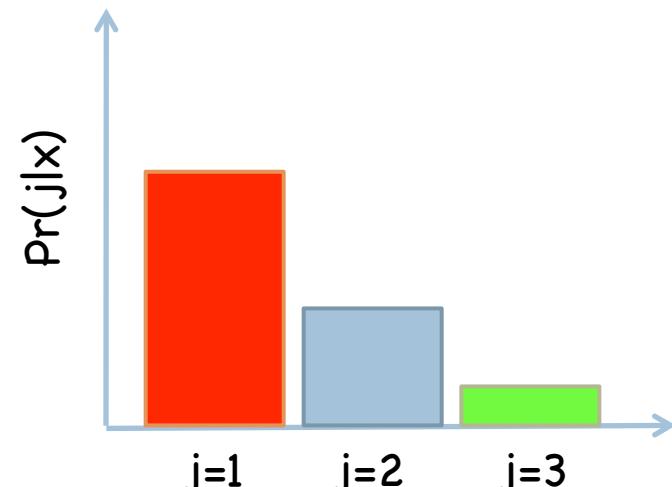
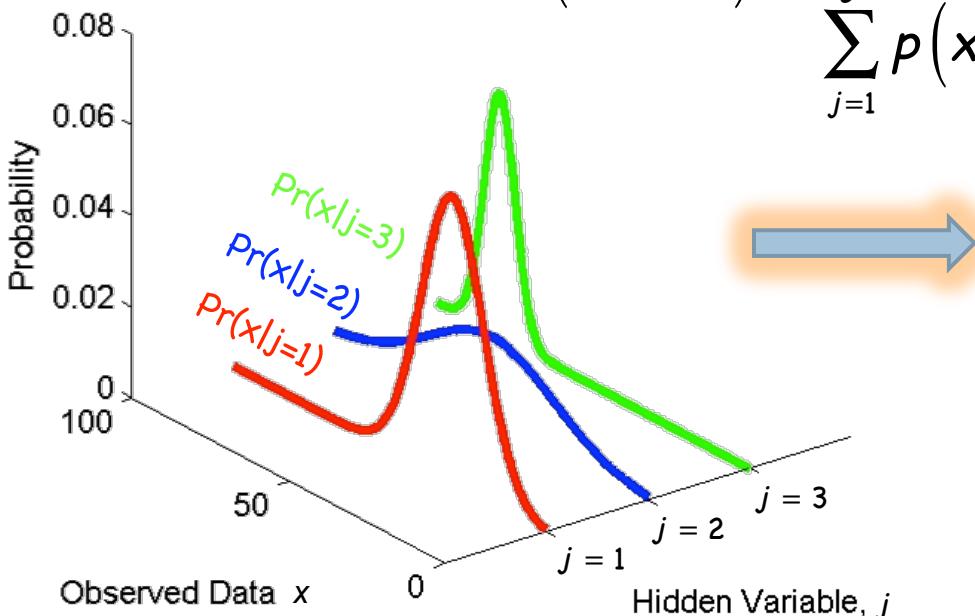
122

Probability & Bayesian Inference

THING TO NOTICE #2:

If we knew the parameters θ it would very easy to estimate the posterior distribution over each hidden variable j_i using Bayes' rule:

$$p(j|x, \theta) = \frac{p(x|j, \theta)P_j}{\sum_{j=1}^J p(x|j, \theta)P_j}$$



Expectation Maximization

123

Probability & Bayesian Inference

Chicken and egg problem:

- could find $j_{1\dots N}$ if we knew θ
- could find θ if we knew $j_{1\dots N}$

Solution: Expectation Maximization (EM) algorithm

(Dempster, Laird and Rubin 1977)

EM for Gaussian mixtures can be viewed as alternation between 2 steps:

1. Expectation Step (E-Step)

- For fixed θ find posterior distribution over responsibilities $j_{1\dots N}$

2. Maximization Step (M-Step)

- Now use these posterior probabilities to re-estimate θ

8.2 Math

Mixture Model

125

Probability & Bayesian Inference

Let

$x_k, k = 1, \dots, N$ denote the training input observations, assumed to be independent
 $j_k \in [1, \dots, J]$ indicate the component of the mixture from which the observation was drawn
(Note that x_k is observable but j_k is hidden.)

Let

$\Theta^t = (\theta^t, P^t)$ represent the unknown parameters we are trying to estimate, where
 θ represents the vector of coefficients for the component distributions and
 P represents the mixing coefficients.

Our mixture model is $p(x_k | \Theta) = \sum_{j=1}^J P_{j_k} p(x_k | j_k; \Theta)$

Q Function

126

Probability & Bayesian Inference

We will iteratively estimate Θ , starting with an initial guess $\Theta(0)$ and monotonically improving our estimate $\Theta(t)$ at successive time steps t .

For this purpose, we define a **Q function**

$$Q(\Theta; \Theta(t)) = E \left[\sum_{k=1}^N P_{j_k} \log p(x_k | j_k; \theta) \right]$$

The Q function represents the **expected log likelihood** of the training data, given our most recent estimate of the parameters $\Theta(t)$, where the expectation is taken over the possible values of the hidden labels j_k .

Expectation Step

127

Probability & Bayesian Inference

- In the E-Step, we calculate the (expected) log probability over the possible parameter values:

$$\begin{aligned} Q(\Theta; \Theta(t)) &= E \left[\sum_{k=1}^N P_{j_k} \log p(x_k | j_k; \theta) \right] \\ &= \sum_{k=1}^N E \left[P_{j_k} \log p(x_k | j_k; \theta) \right] \\ &= \sum_{k=1}^N \sum_{j_k=1}^J P(j_k | x_k; \Theta) P_{j_k} \log p(x_k | j_k; \theta) \end{aligned}$$

Maximization Step

128

Probability & Bayesian Inference

- In the M-Step, we select for our new parameter estimate the value that maximizes this expected log probability:

$$\Theta(t + 1) = \arg \max_{\Theta} Q(\Theta; \Theta(t))$$

Example: Mixture of Isotropic Gaussians

129

Probability & Bayesian Inference

$$p(\mathbf{x}_k \mid j; \theta) = \frac{1}{(2\pi\sigma_j^2)^{\frac{d}{2}}} \exp\left(-\frac{\|\mathbf{x}_k - \mu_j\|^2}{2\sigma_j^2}\right)$$

□ E-Step:

$$Q(\Theta; \Theta(t)) = \sum_{k=1}^N \sum_{j=1}^J p(j \mid \mathbf{x}_k; \Theta) \left(-\frac{l}{2} \log \sigma_j^2 - \frac{1}{2\sigma_j^2} \|\mathbf{x}_k - \mu_j\|^2 + \log P_j \right)$$

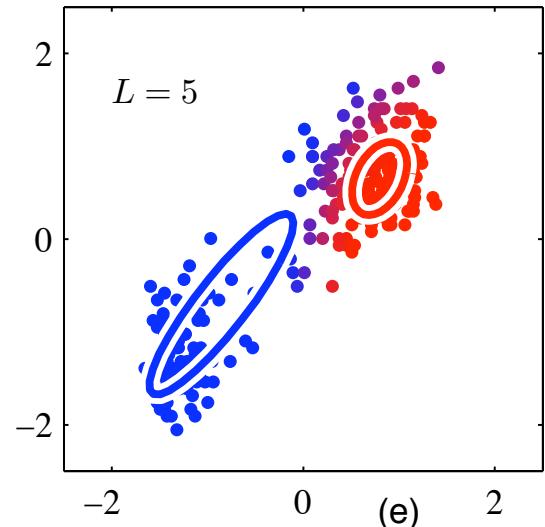
Example: Mixture of Isotropic Gaussians

130

Probability & Bayesian Inference

□ Responsibilities Update Equation:

$$P(j | \mathbf{x}_k; \Theta(t)) = \frac{P(\mathbf{x}_k | j; \theta(t)) P_j(t)}{\sum_{j=1}^J P(\mathbf{x}_k | j; \theta(t)) P_j(t)}$$



□ Parameter Update Equations:

$$\mu_j(t+1) = \frac{\sum_{k=1}^N P(j | \mathbf{x}_k; \Theta(t)) \mathbf{x}_k}{\sum_{k=1}^N P(j | \mathbf{x}_k; \Theta(t))}$$

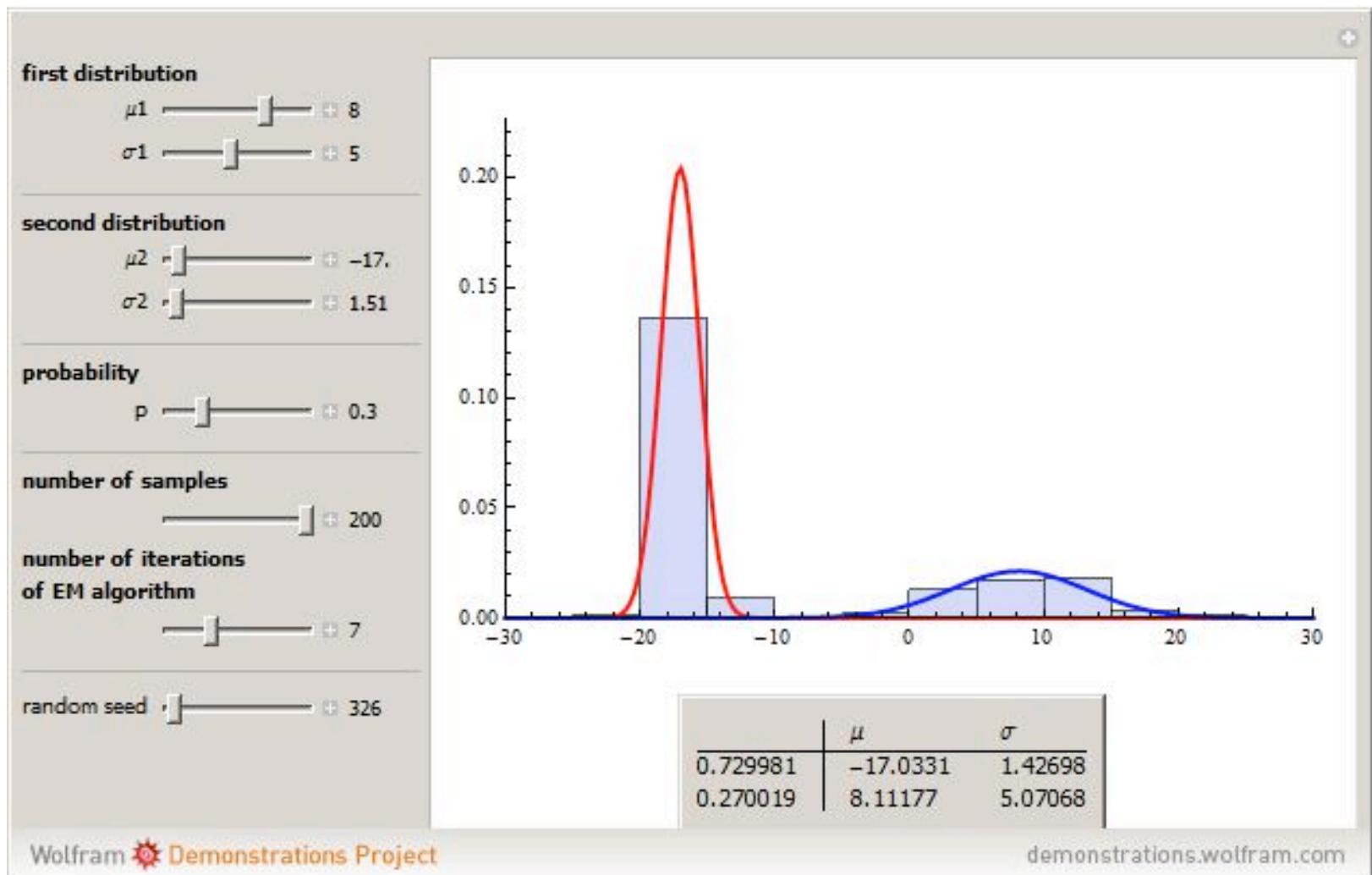
$$\sigma_j^2(t+1) = \frac{\sum_{k=1}^N P(j | \mathbf{x}_k; \Theta(t)) \|\mathbf{x}_k - \mu_j(t+1)\|^2}{\sum_{k=1}^N P(j | \mathbf{x}_k; \Theta(t))}$$

$$P_j(t+1) = \frac{1}{N} \sum_{k=1}^N P(j | \mathbf{x}_k; \Theta(t))$$

Univariate Gaussian Mixture Example

131

Probability & Bayesian Inference

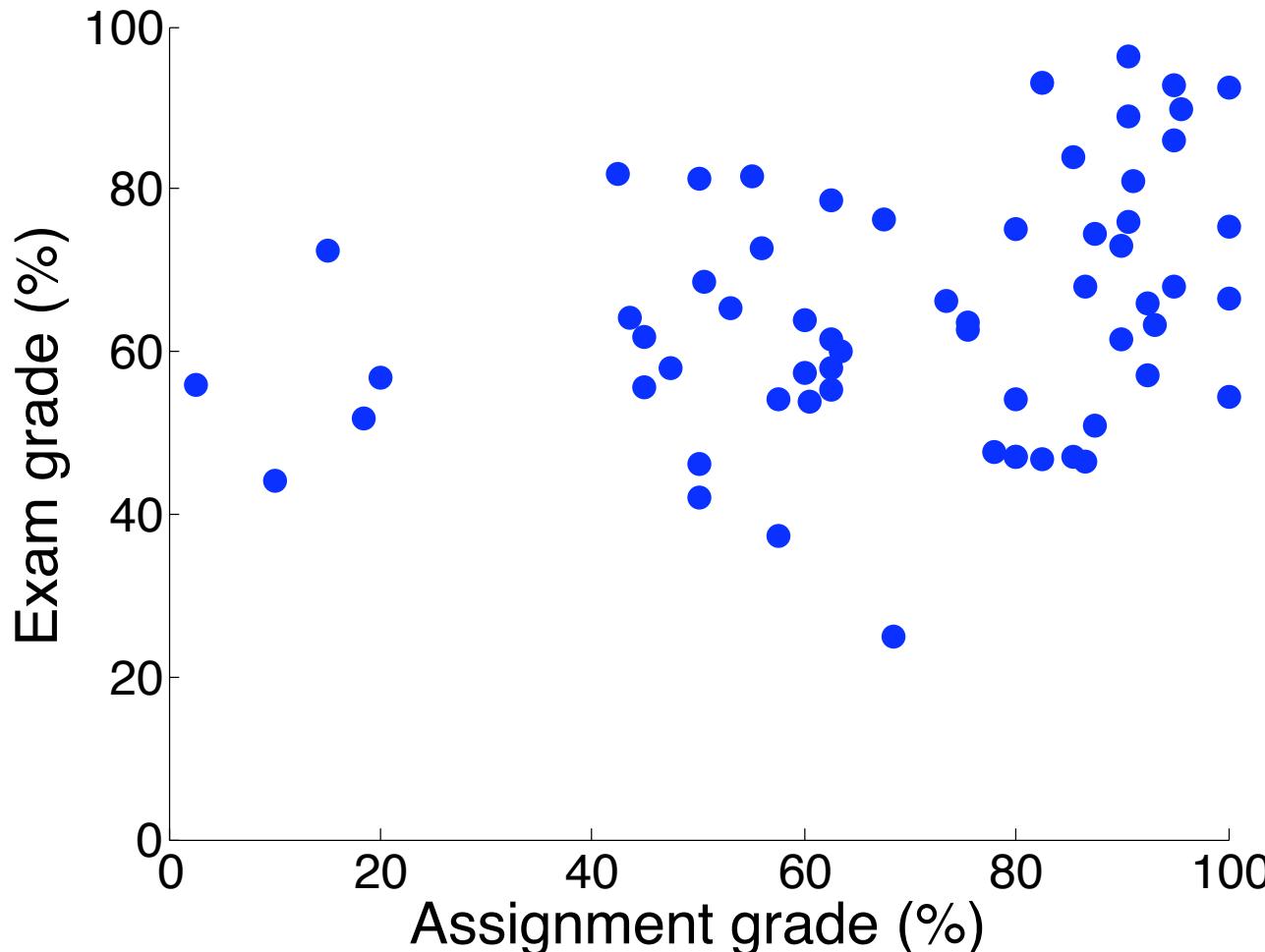


2-Component Bivariate MATLAB Example

132

Probability & Bayesian Inference

CSE 2011Z 2010W



2-Component Bivariate MATLAB Example

133

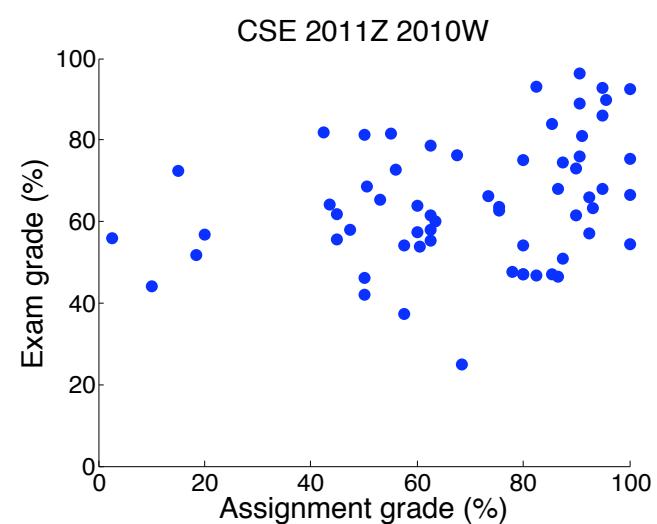
Probability & Bayesian Inference

```
%update responsibilities
```

```
for i=1:k  
    p(:,i)=alphas(i).*mvnpdf(x,mu(:,i),squeeze(S(:,:,i)));  
end  
p=p./repmat((sum(p,2)),1,k);
```

```
%update parameters
```

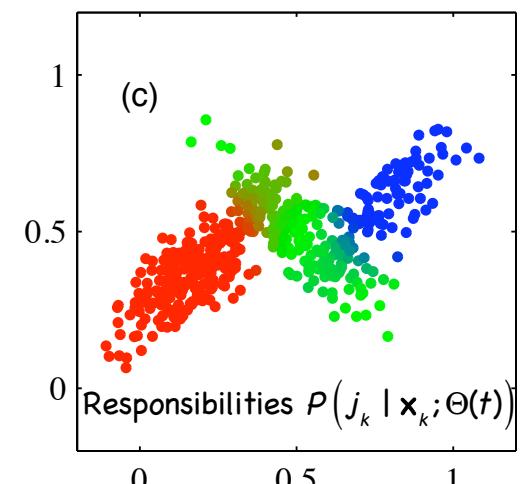
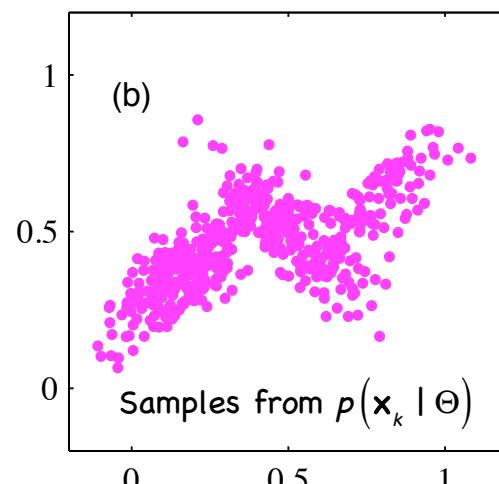
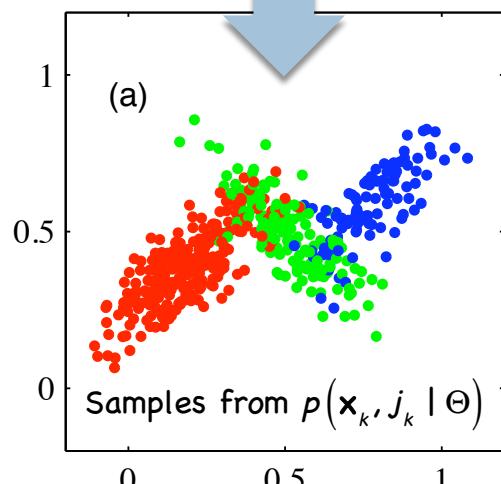
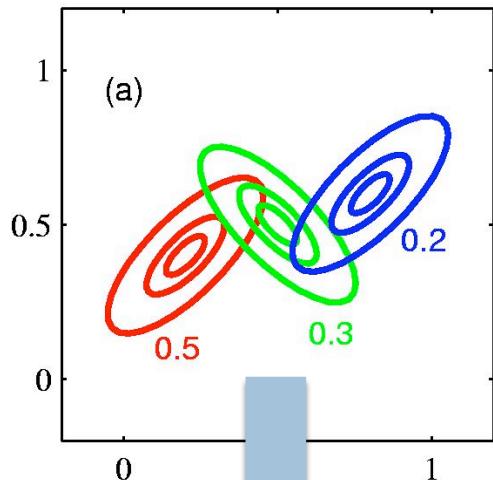
```
for i=1:k  
    Nk=sum(p(:,i));  
    mu(:,i)=p(:,i)'*x/Nk;  
    dev=x-repmat(mu(:,i),N,1);  
    S(:,:,i)=(repmat(p(:,i),1,D).*dev)'.*dev/Nk;  
    alphas(i)=Nk/N;  
end
```



Bivariate Gaussian Mixture Example

134

Probability & Bayesian Inference



Expectation Maximization

135

Probability & Bayesian Inference

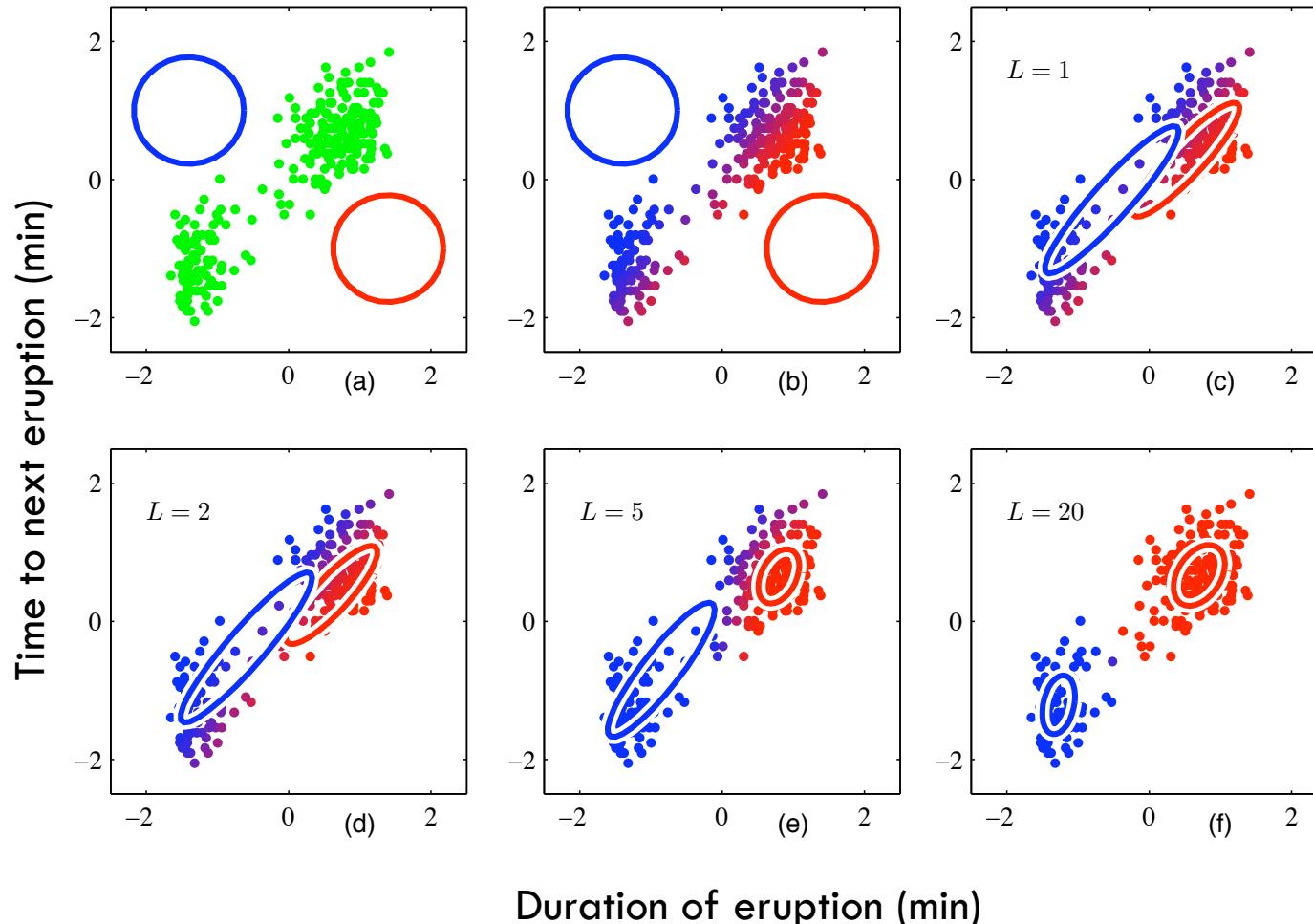
- EM is guaranteed to monotonically increase the likelihood.
- However, since in general the likelihood is non-convex, we are not guaranteed to find the globally optimal parameters.

8.3 Applications

Old Faithful Example

137

Probability & Bayesian Inference



Face Detection Example: 2 Components

138

Probability & Bayesian Inference

Face Model
Parameters

Face Model Parameters

0.4999

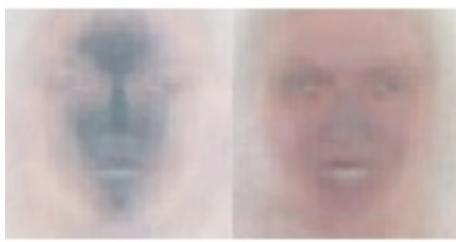
0.5001

Prior

Mean



Each component is still assumed to have diagonal covariance.



Standard deviation

The face model and non-face model have divided the data into two clusters. In each case, these clusters have roughly equal weights.

0.4675

0.5325

Prior

Mean



The primary thing that these seem to have captured is the photometric (luminance) variation.

Non-Face
Model
Parameters

Standard deviation

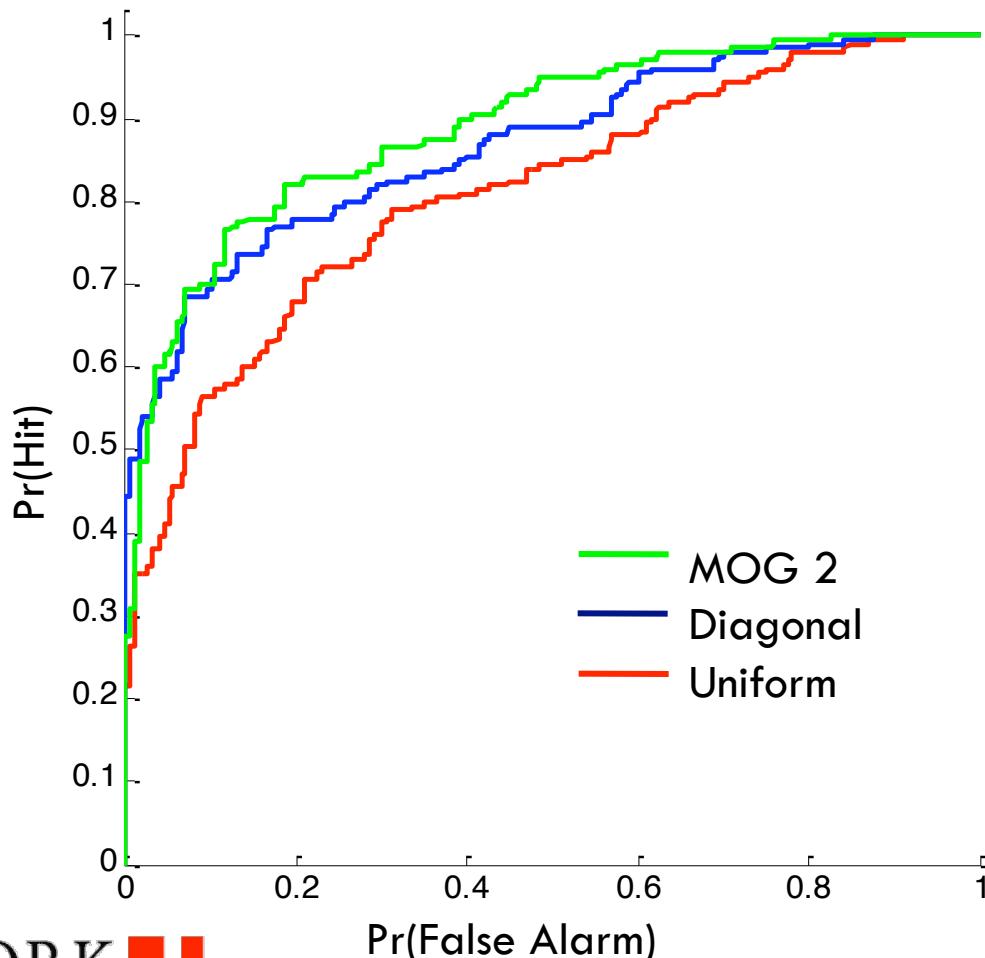


Note that the standard deviations have become smaller than for the single Gaussian model as any given data point is likely to be close to one mean or the other.

Results for MOG 2 Model

139

Probability & Bayesian Inference



Performance improves relative to a single Gaussian model, although it is not dramatic.

We have a better description of the data likelihood.

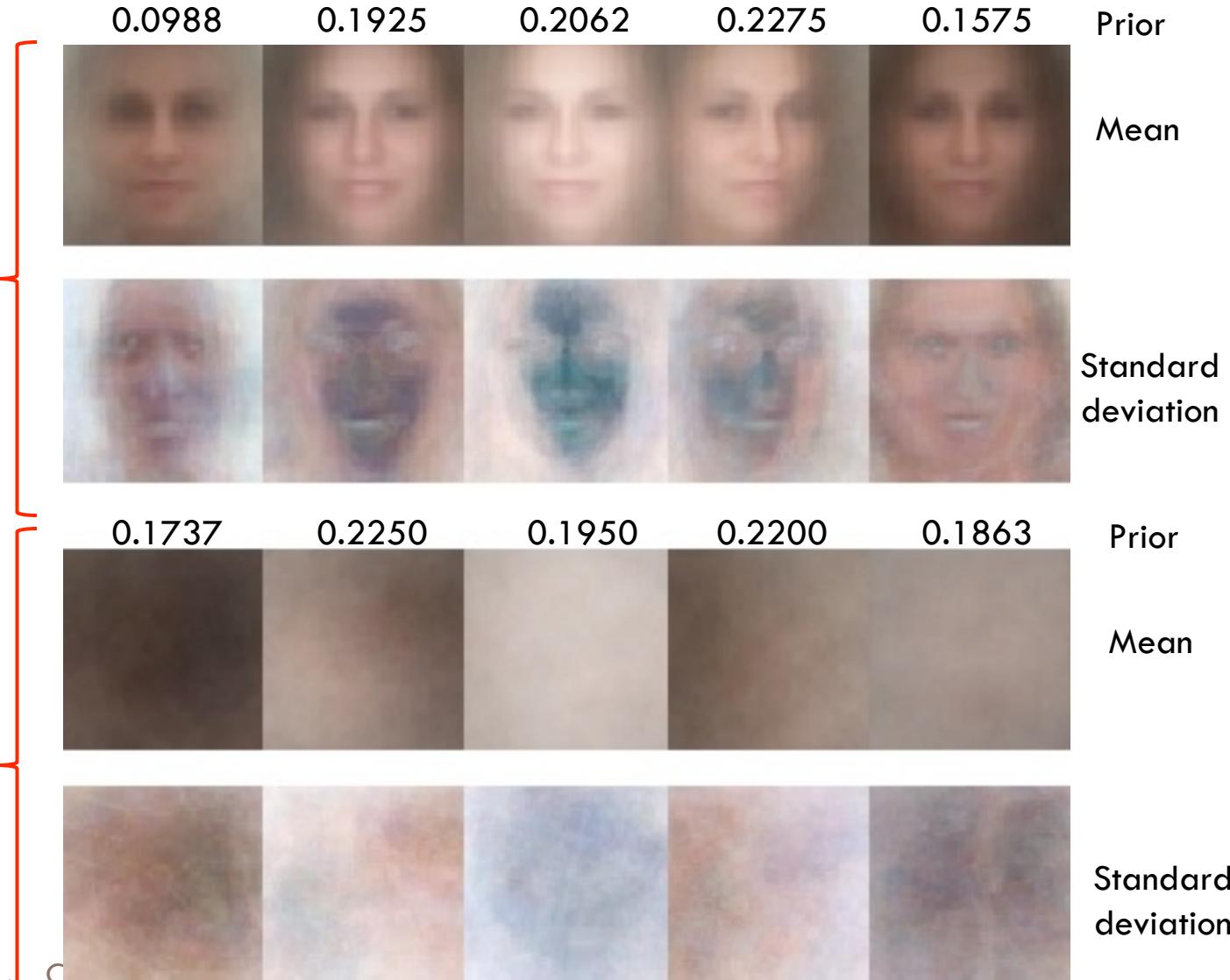
MOG 5 Components

140

Probability & Bayesian Inference

Face Model
Parameters

Non-Face
Model
Parameters



MOG 10 Components

141

Probability & Bayesian Inference

0.0075 0.1425 0.1437 0.0988 0.1038 0.1187 0.1638 0.1175 0.1038 0.0000



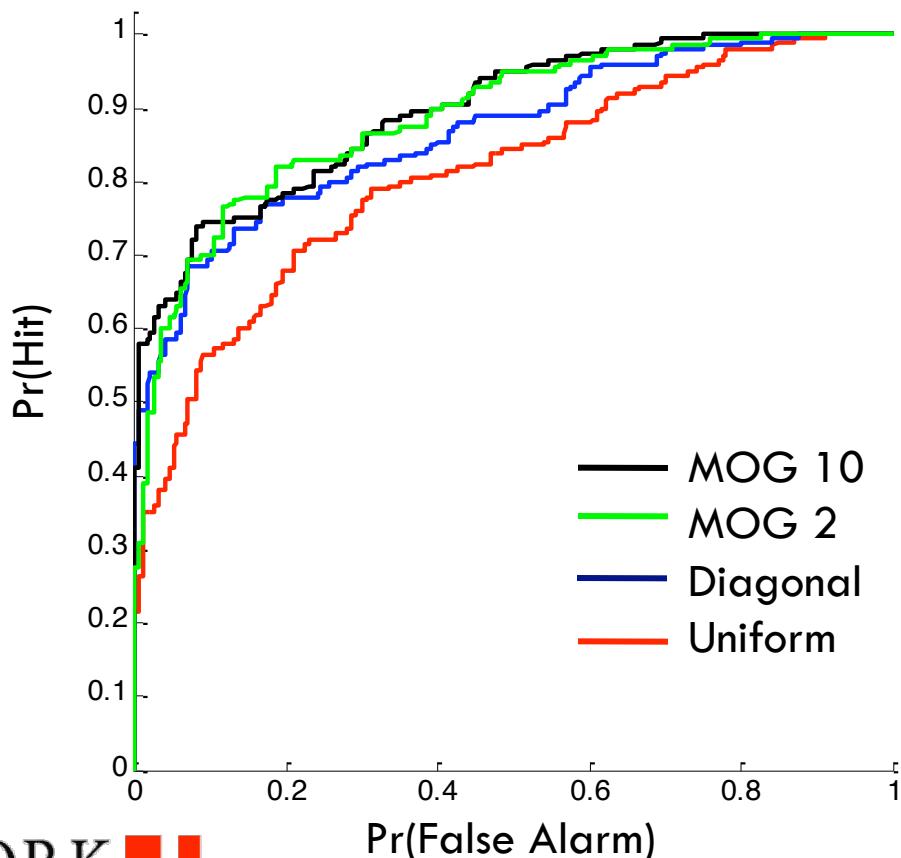
0.1137 0.0688 0.0763 0.0800 0.1338 0.1063 0.1063 0.1263 0.0900 0.0988



Results for Mog 10 Model

142

Probability & Bayesian Inference



Performance improves
slightly more,
particularly at low false
alarm rates.

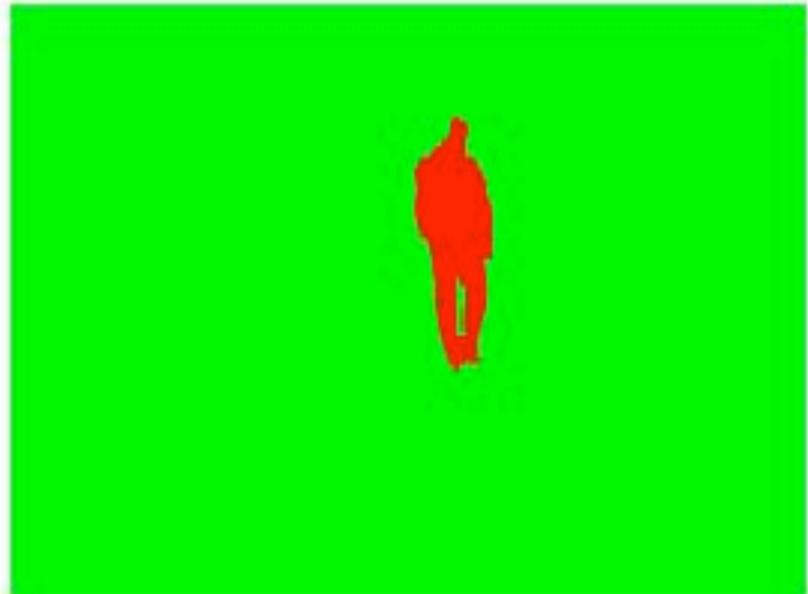
Background Subtraction

143

Probability & Bayesian Inference



Test Image



Desired Segmentation

GOAL : (i) Learn background model (ii) use this to segment regions where the background has been occluded

What if the scene isn't static?

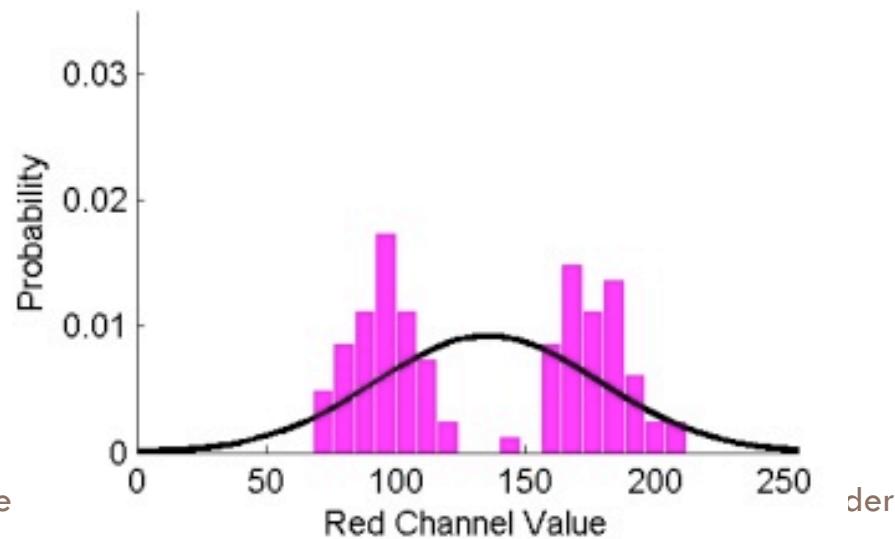
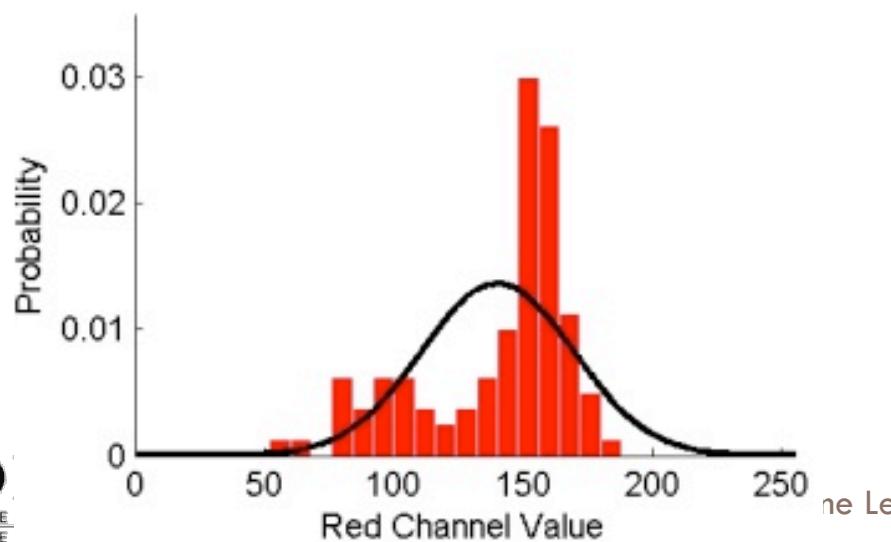
144

Probability & Bayesian Inference



Gaussian is no longer a good fit to the data.

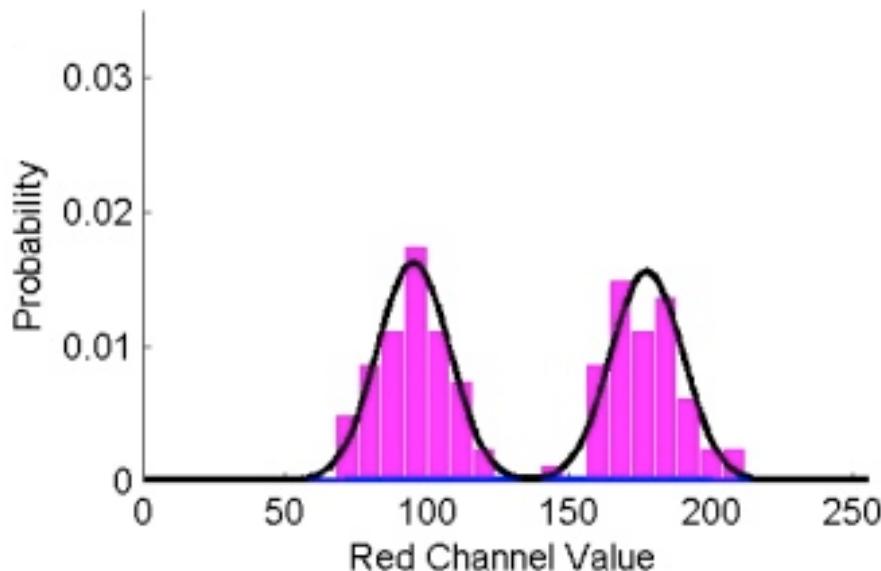
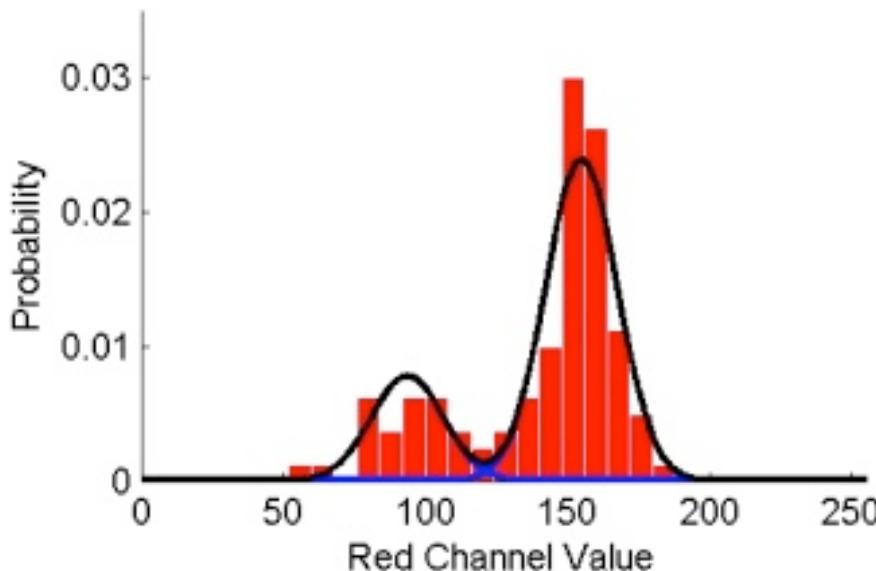
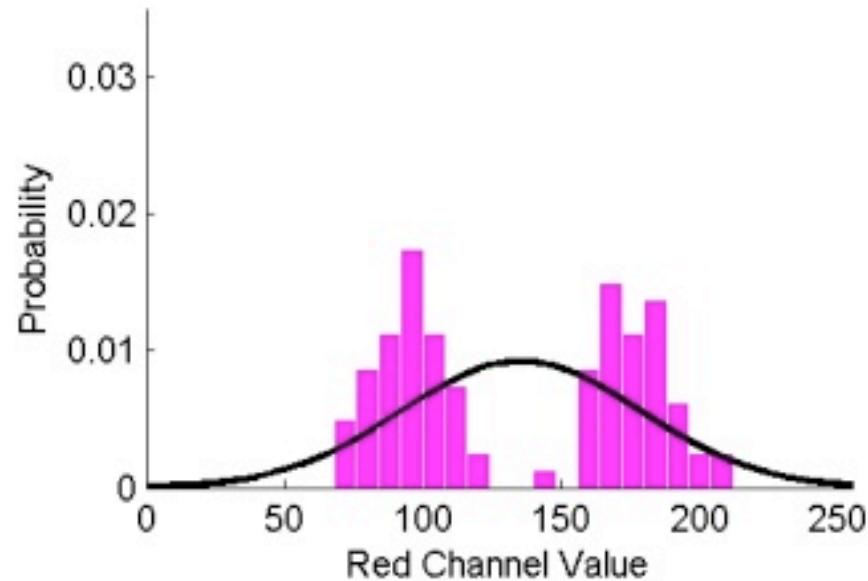
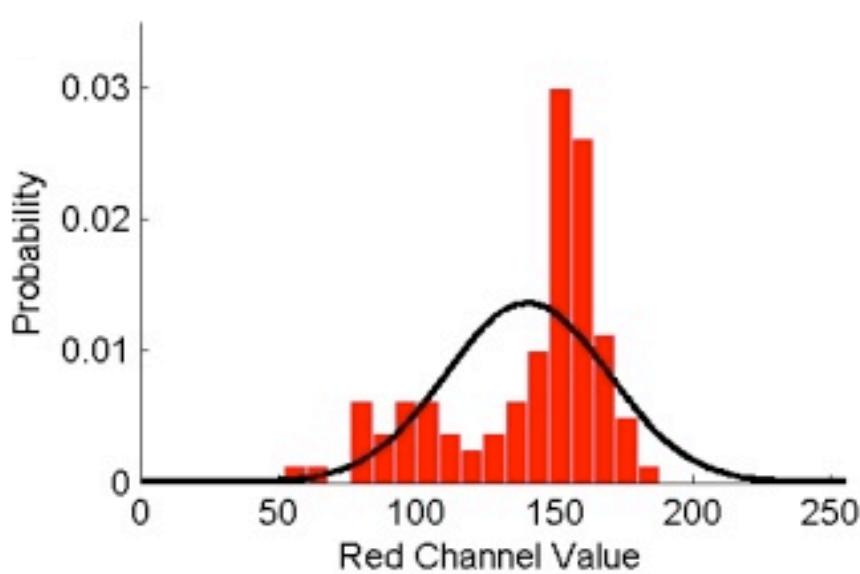
Not obvious exactly what probability model would fit better.



Background Mixture Model

145

Probability & Bayesian Inference



Background Subtraction Example

146

Probability & Bayesian Inference





End of Lecture 7

Bayesian Decision Theory: Topics

148

Probability & Bayesian Inference

1. Probability
2. The Univariate Normal Distribution
3. Bayesian Classifiers
4. Minimizing Risk
5. The Multivariate Normal Distribution
6. Decision Boundaries in Higher Dimensions
7. Parameter Estimation
8. Mixture Models and EM
9. **Nonparametric Density Estimation**
10. Training and Evaluation Methods
11. What are Bayes Nets?

9. Nonparametric Methods

Nonparametric Methods

150

Probability & Bayesian Inference

- Parametric distribution models are restricted to specific forms, which may not always be suitable; for example, consider modelling a multimodal distribution with a single, unimodal model.
- You can use a mixture model, but then you have to decide on the number of components, and hope that your parameter estimation algorithm (e.g., EM) converges to a global optimum!
- Nonparametric approaches make few assumptions about the overall shape of the distribution being modelled, and in some cases may be simpler than using a mixture model.

Histogramming

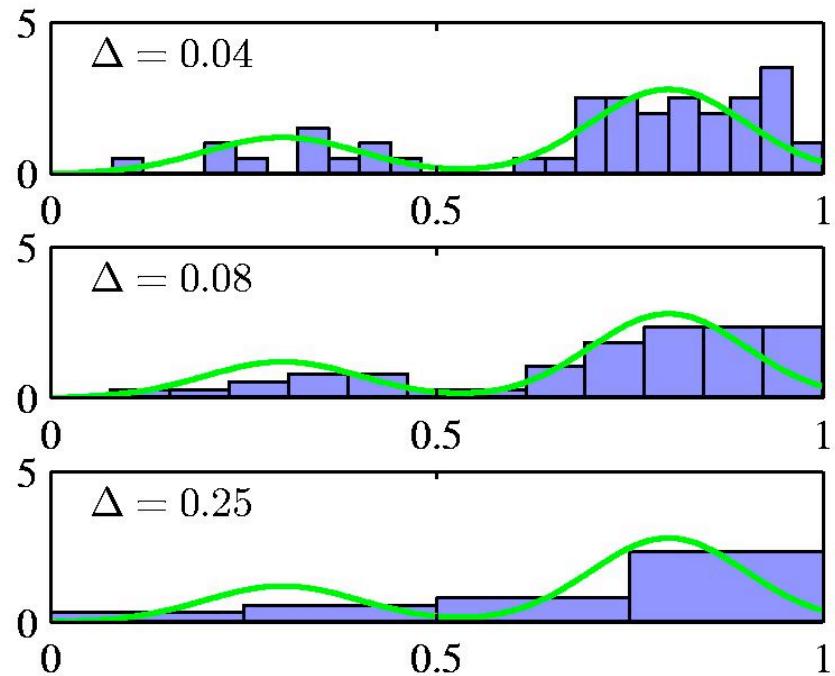
151

Probability & Bayesian Inference

- **Histogram methods** partition the data space into distinct bins with widths Δ_i and count the number of observations, n_i , in each bin.

$$p_i = \frac{n_i}{N\Delta_i}$$

- Often, the same width is used for all bins, $\Delta_i = \Delta$.
- Δ acts as a smoothing parameter.



- In a D -dimensional space, using M bins in each dimension will require M^D bins!

The curse of dimensionality

Kernel Density Estimation

152

Probability & Bayesian Inference

- Assume observations drawn from a density $p(\mathbf{x})$ and consider a small region R containing \mathbf{x} such that
- If the volume V of R is sufficiently small, $p(\mathbf{x})$ is approximately constant over R and

$$P = \int_{\mathcal{R}} p(\mathbf{x}) d\mathbf{x}.$$

$$P \simeq p(\mathbf{x})V$$

- The expected number K out of N observations that will lie inside R is given by

- Thus

$$p(\mathbf{x}) = \frac{K}{NV}.$$

$$K \simeq NP.$$

Kernel Density Estimation

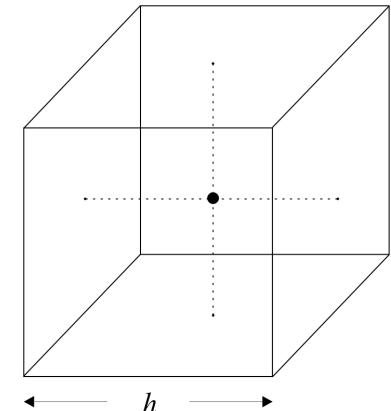
153

Probability & Bayesian Inference

Kernel Density Estimation: fix V , estimate K from the data. Let R be a hypercube centred on x and define the kernel function (Parzen window)

$$k((\mathbf{x} - \mathbf{x}_n)/h) = \begin{cases} 1, & |(x_i - x_{ni})/h| \leq 1/2, \\ 0, & \text{otherwise.} \end{cases} \quad i = 1, \dots, D,$$

$$p(\mathbf{x}) = \frac{K}{NV}.$$



It follows that

$$K = \sum_{n=1}^N k\left(\frac{\mathbf{x} - \mathbf{x}_n}{h}\right)$$

and hence

$$p(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N \frac{1}{h^D} k\left(\frac{\mathbf{x} - \mathbf{x}_n}{h}\right).$$

Kernel Density Estimation

154

Probability & Bayesian Inference

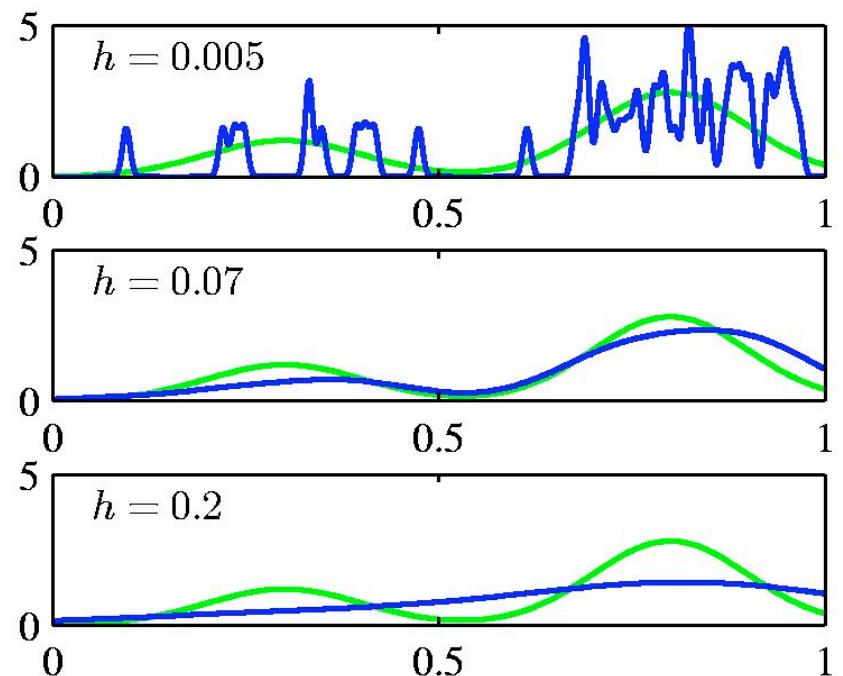
To avoid discontinuities in $p(x)$, use a smooth kernel, e.g. a Gaussian

$$p(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N \frac{1}{(2\pi h^2)^{D/2}} \exp \left\{ -\frac{\|\mathbf{x} - \mathbf{x}_n\|^2}{2h^2} \right\}$$

(Any kernel $k(u)$ such that

$$\begin{aligned} k(\mathbf{u}) &\geq 0, \\ \int k(\mathbf{u}) d\mathbf{u} &= 1 \end{aligned}$$

will work.)

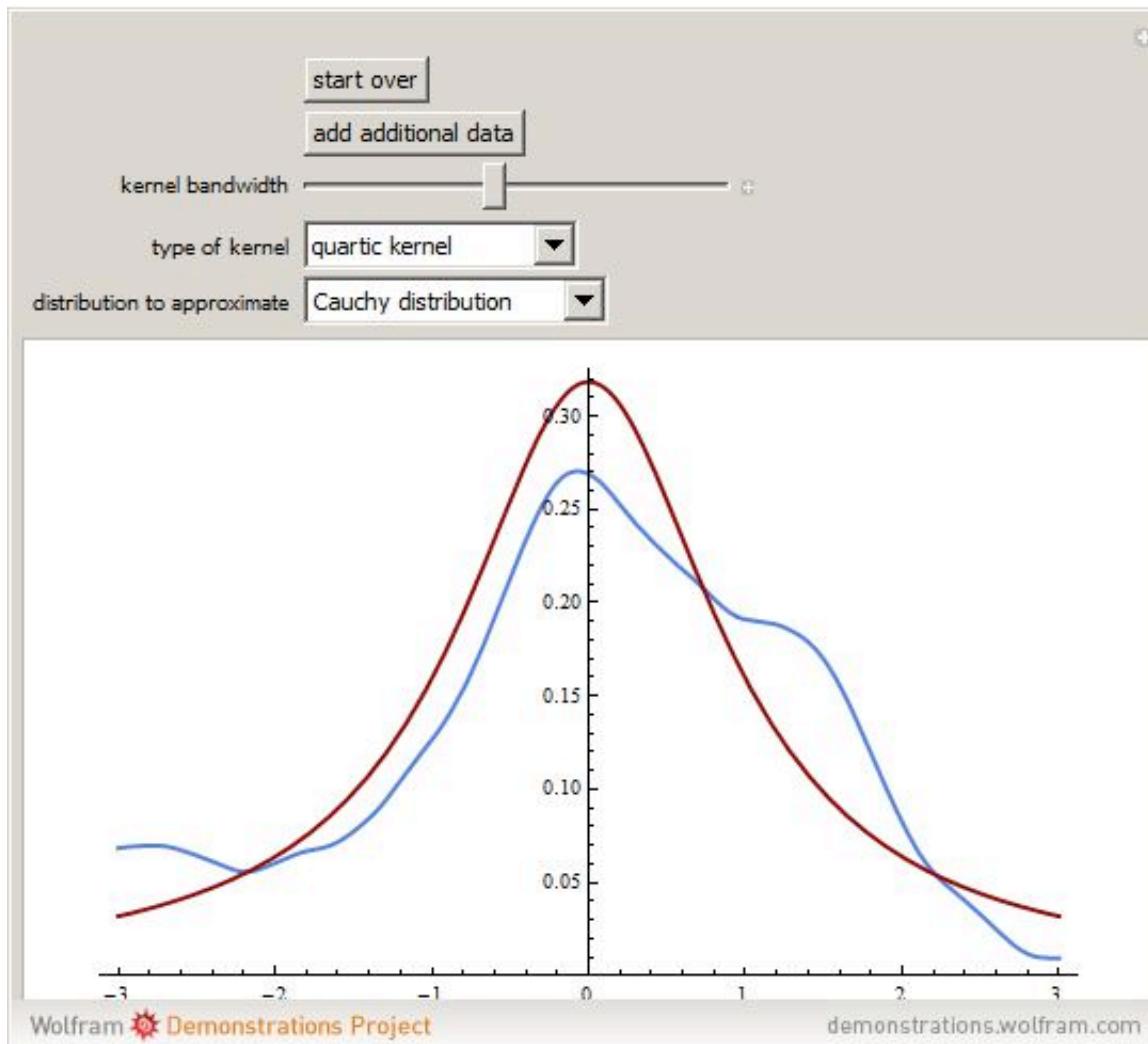


h acts as a smoother.

KDE Example

155

Probability & Bayesian Inference



Wolfram Demonstrations Project

demonstrations.wolfram.com

Kernel Density Estimation

156

Probability & Bayesian Inference

- Problem: if V is fixed, there may be too few points in some regions to get an accurate estimate.

Nearest Neighbour Density Estimation

157

Probability & Bayesian Inference

Nearest Neighbour

Density Estimation: fix K ,
estimate V from the data.

Consider a hypersphere
centred on x and let it
grow to a volume V^* that
includes K of the given N
data points. Then

$$p(x) \simeq \frac{K}{NV^*}.$$

```
for j=1:np
    d=sort(abs(x(j)-xi));
    V=2*d(K(i));
    phat(j)=K(i)/(N*V);
end
```

Nearest Neighbour Density Estimation

158

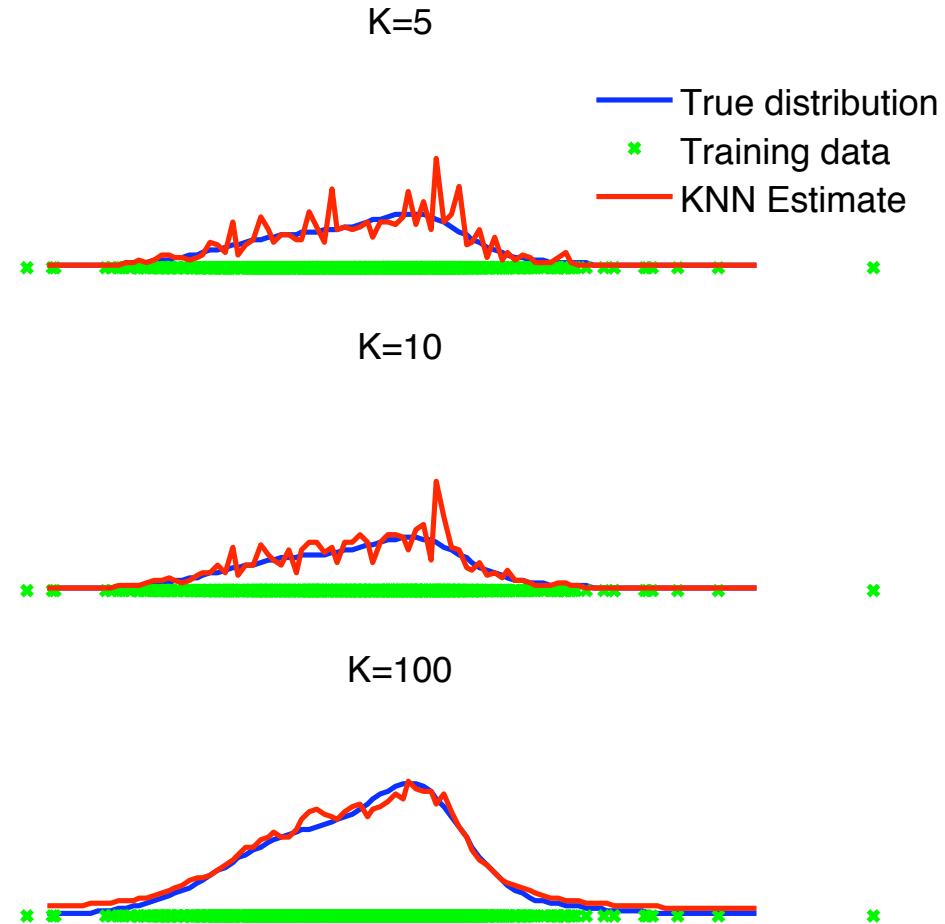
Probability & Bayesian Inference

Nearest Neighbour

Density Estimation: fix K ,
estimate V from the data.

Consider a hypersphere
centred on x and let it
grow to a volume V^* that
includes K of the given N
data points. Then

$$p(x) \simeq \frac{K}{NV^*}.$$



Nearest Neighbour Density Estimation

159

Probability & Bayesian Inference

- Problem: does not generate a proper density (for example, integral is unbounded on \mathbb{R}^D)
- In practice, on finite domains, can normalize.
- But makes strong assumption on tails $\left(\propto \frac{1}{x}\right)$

Nonparametric Methods

160

Probability & Bayesian Inference

- Nonparametric models (not histograms) require storing and computing with the entire data set.
- Parametric models, once fitted, are much more efficient in terms of storage and computation.

K-Nearest-Neighbours for Classification

161

Probability & Bayesian Inference

- Given a data set with N_k data points from class C_k and $\sum_k N_k = N$, we have

$$p(\mathbf{x}) = \frac{K}{NV}$$

- and correspondingly

$$p(\mathbf{x}|\mathcal{C}_k) = \frac{K_k}{N_k V}.$$

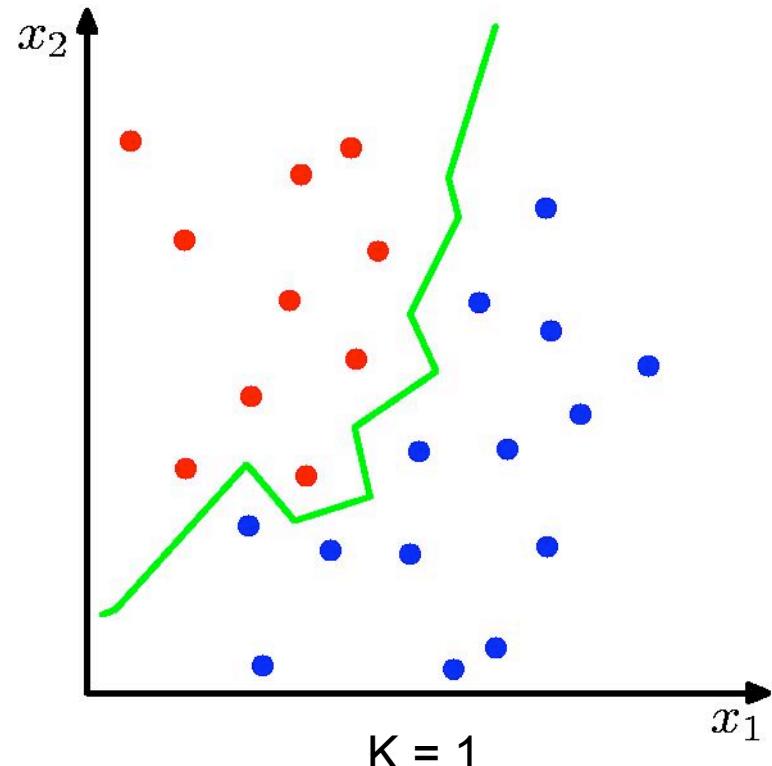
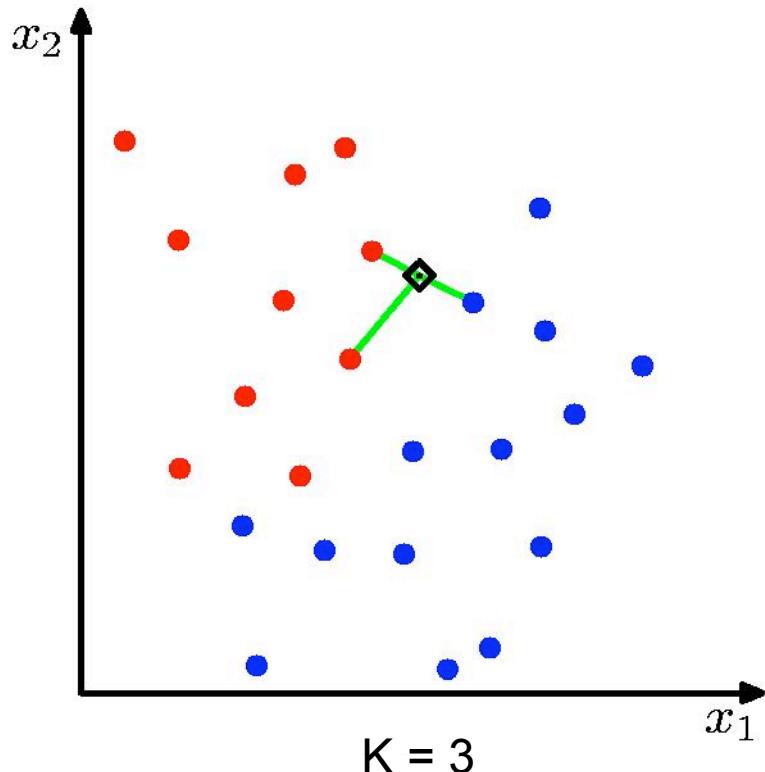
- Since $p(\mathcal{C}_k) = N_k/N$, Bayes' theorem gives

$$p(\mathcal{C}_k|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{p(\mathbf{x})} = \frac{K_k}{K}.$$

K-Nearest-Neighbours for Classification

162

Probability & Bayesian Inference

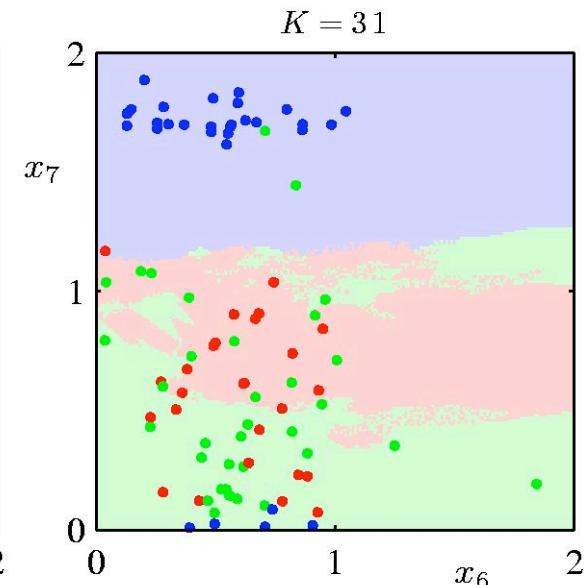
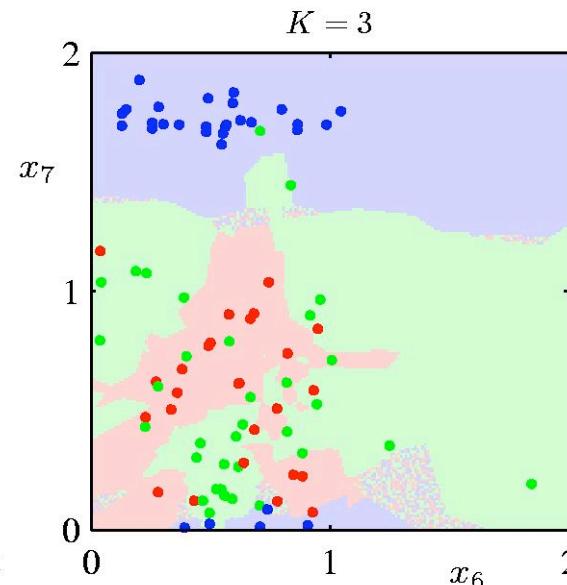
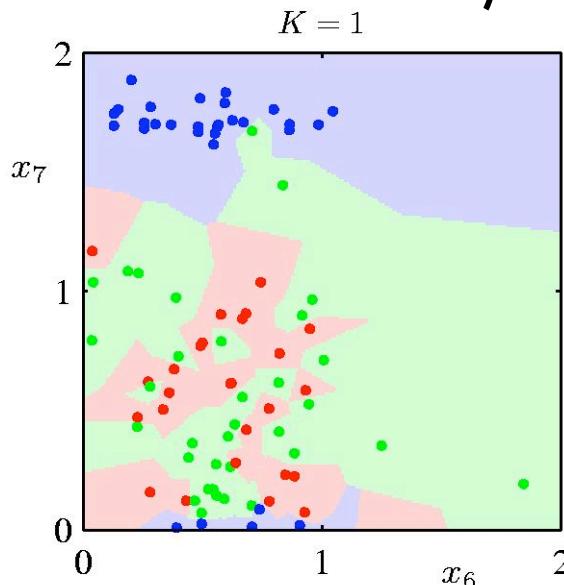


K-Nearest-Neighbours for Classification

163

Probability & Bayesian Inference

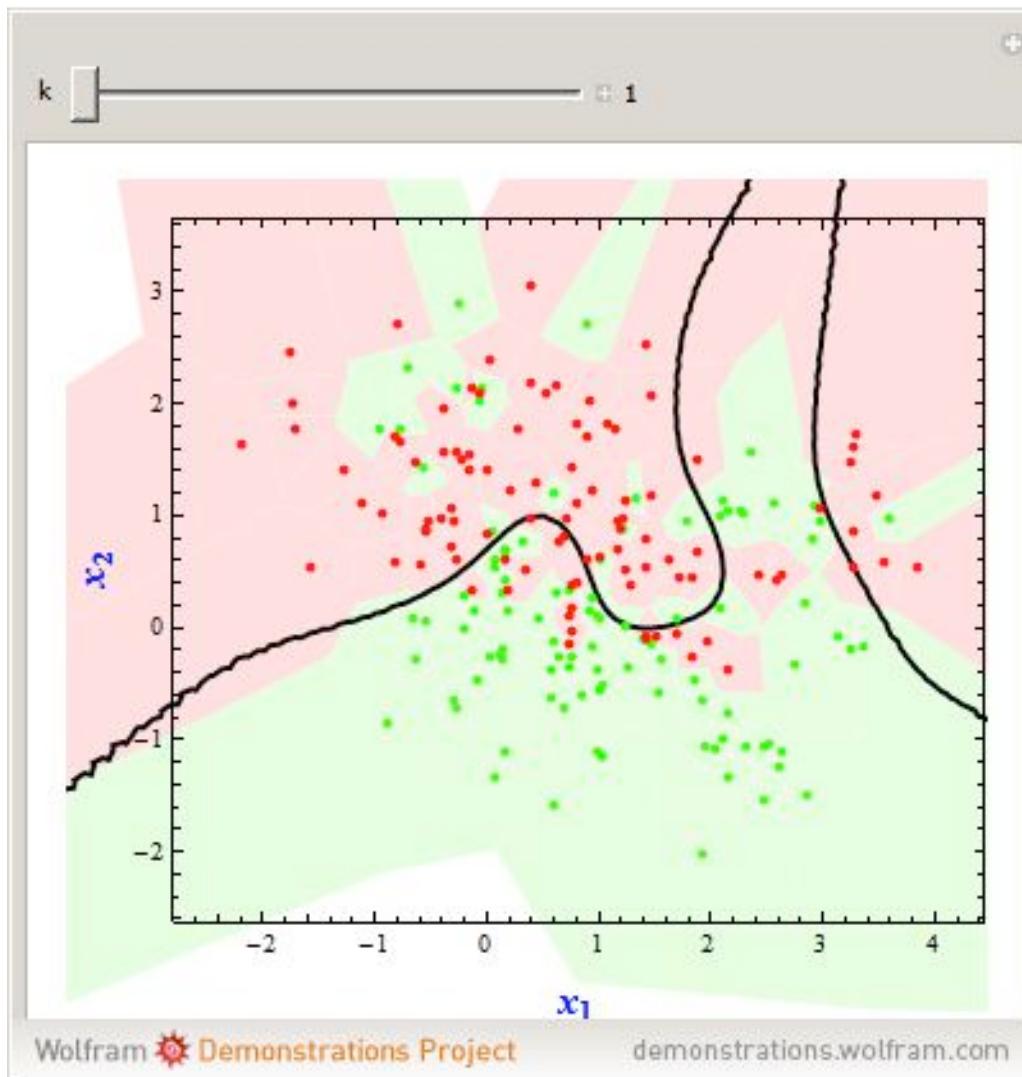
- K acts as a smother
- As $N \rightarrow \infty$, the error rate of the 1-nearest-neighbour classifier is never more than twice the optimal error (obtained from the true conditional class distributions).



KNN Example

164

Probability & Bayesian Inference



Naïve Bayes Classifiers

165

Probability & Bayesian Inference

- All of these nonparametric methods require lots of data to work. If $\mathcal{O}(N)$ training points are required for accurate estimation in 1 dimension, then $\mathcal{O}(N^D)$ points are required for D -dimensional input vectors.
- It may sometimes be possible to assume that the individual dimensions of the feature vector are conditionally independent. Then we have

$$p(\underline{x} \mid \omega_i) = \prod_{j=1}^D p(x_j \mid \omega_i)$$

- This reduces the data requirements to $\mathcal{O}(DN)$.



End of Lecture 8

Bayesian Decision Theory: Topics

167

Probability & Bayesian Inference

1. Probability
2. The Univariate Normal Distribution
3. Bayesian Classifiers
4. Minimizing Risk
5. The Multivariate Normal Distribution
6. Decision Boundaries in Higher Dimensions
7. Parameter Estimation
8. Mixture Models and EM
9. Nonparametric Density Estimation
10. **Training and Evaluation Methods**
11. What are Bayes Nets?

10. Training and Evaluation Methods

Machine Learning System Design

169

Probability & Bayesian Inference

- The process of solving a particular classification or regression problem typically involves the following sequence of steps:
 1. **Design and code** promising candidate systems
 2. **Train** each of the candidate systems (i.e., learn the parameters)
 3. **Evaluate** each of the candidate systems
 4. **Select and deploy** the best of these candidate systems

Using Your Training Data

170

Probability & Bayesian Inference

- You will always have a finite amount of data on which to train and evaluate your systems.
- The performance of a classification system is often **data-limited**: if we only had more data, we could make the system better.
- Thus it is important to use your finite data set wisely.

Overfitting

171

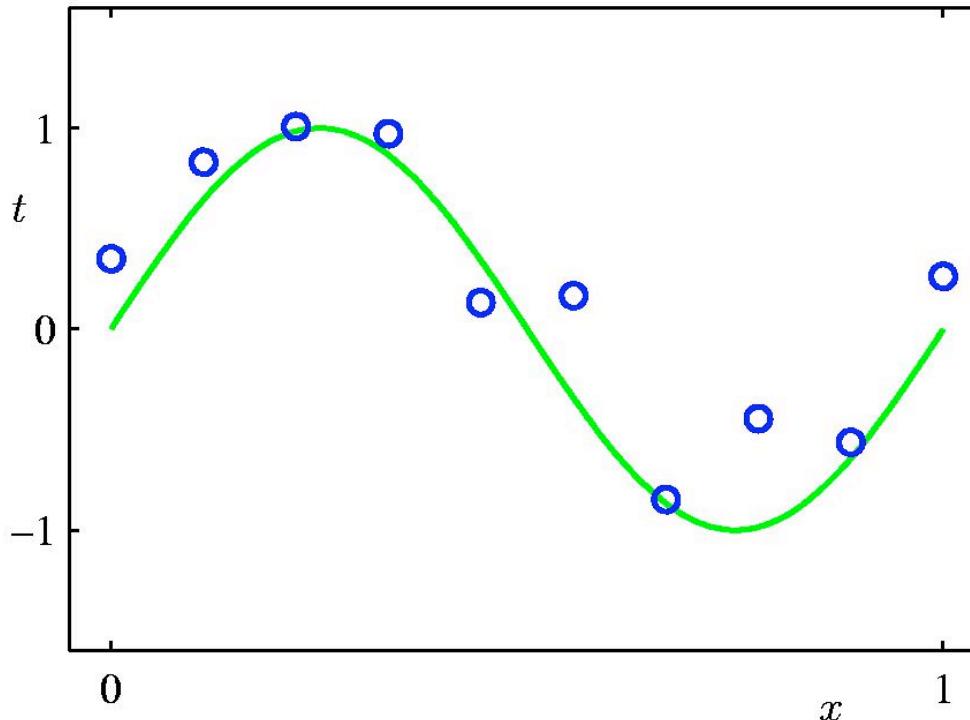
Probability & Bayesian Inference

- Given that learning is often data-limited, it is tempting to use all of your data to estimate the parameters of your models, and then select the model with the lowest error on your training data.
- Unfortunately, this leads to a notorious problem called **over-fitting**.

Example: Polynomial Curve Fitting

172

Probability & Bayesian Inference

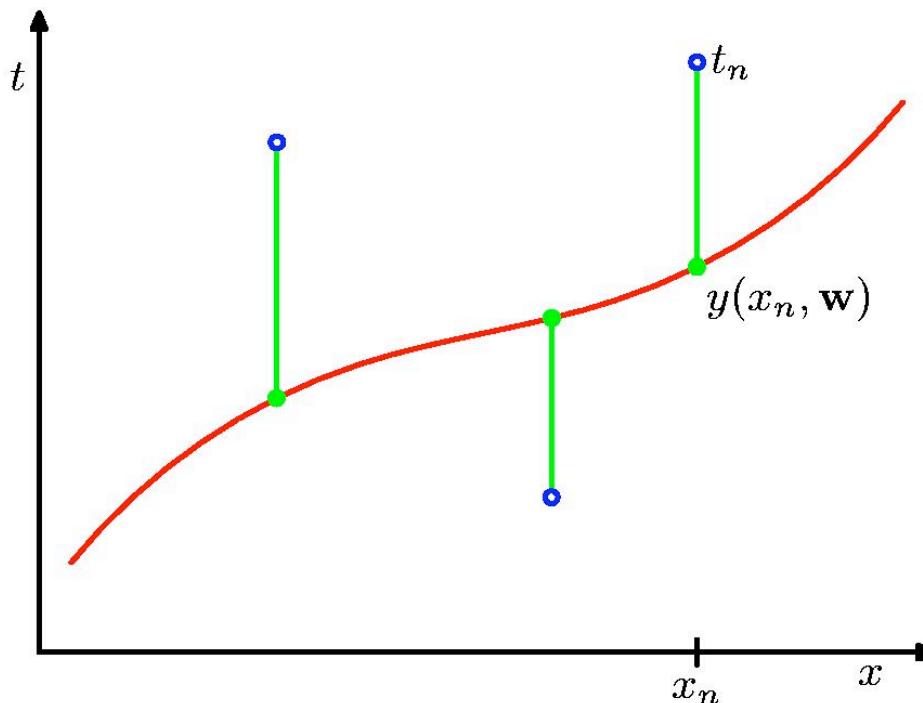


$$y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \dots + w_M x^M = \sum_{j=0}^M w_j x^j$$

Sum-of-Squares Error Function

173

Probability & Bayesian Inference

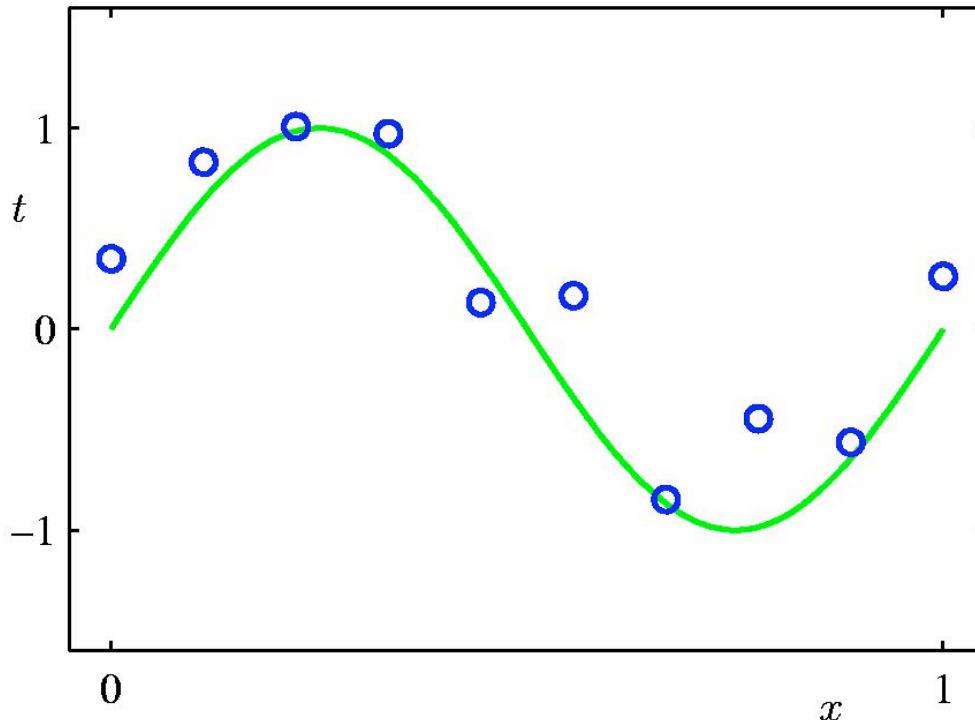


$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2$$

How do we choose M , the order of the model?

174

Probability & Bayesian Inference

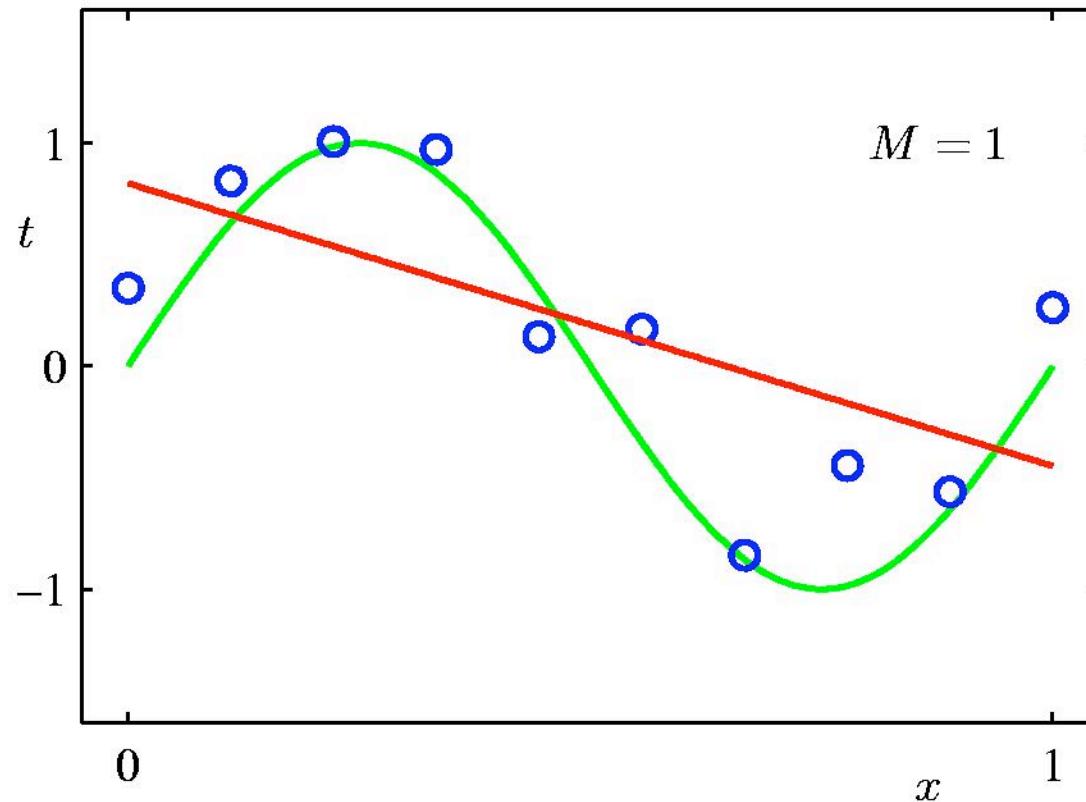


$$y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \dots + w_M x^M = \sum_{j=0}^M w_j x^j$$

1st Order Polynomial

175

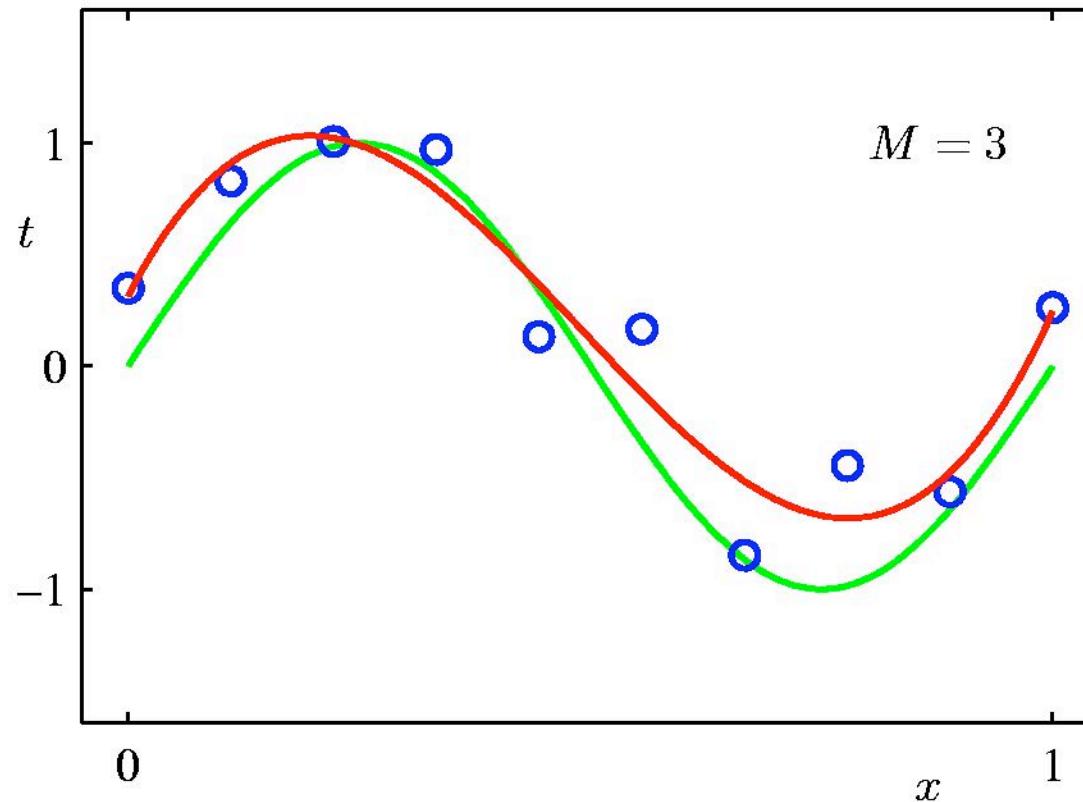
Probability & Bayesian Inference



3rd Order Polynomial

176

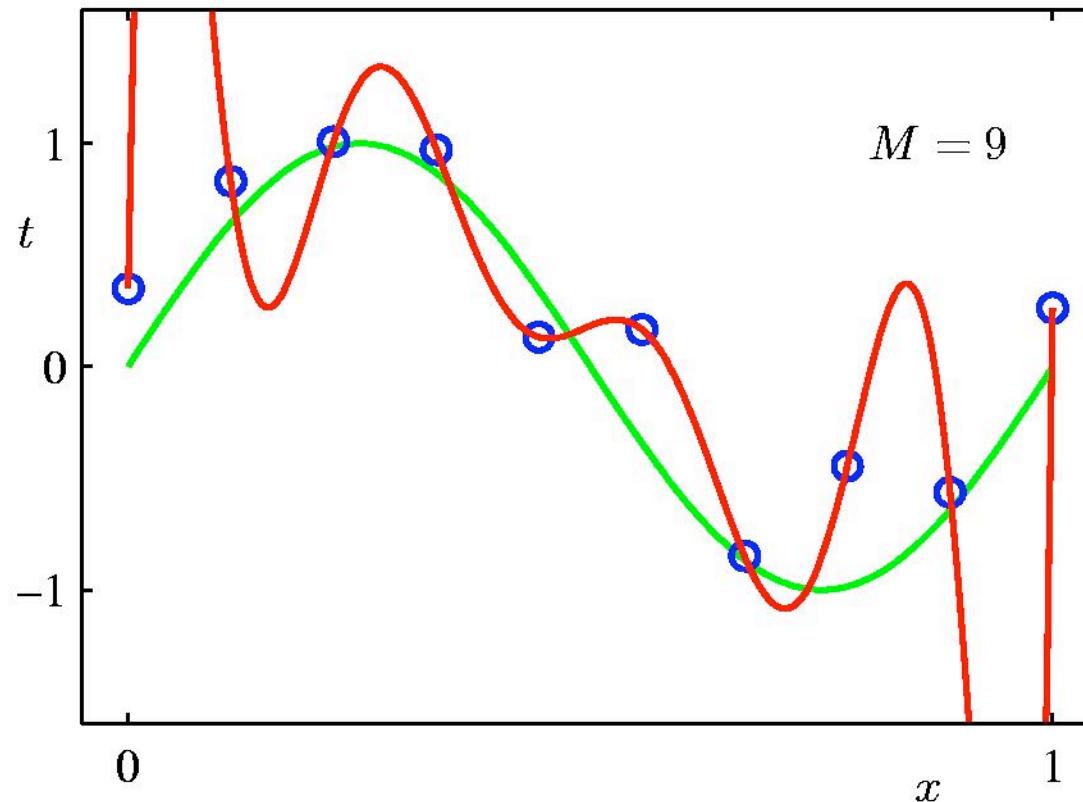
Probability & Bayesian Inference



9th Order Polynomial

177

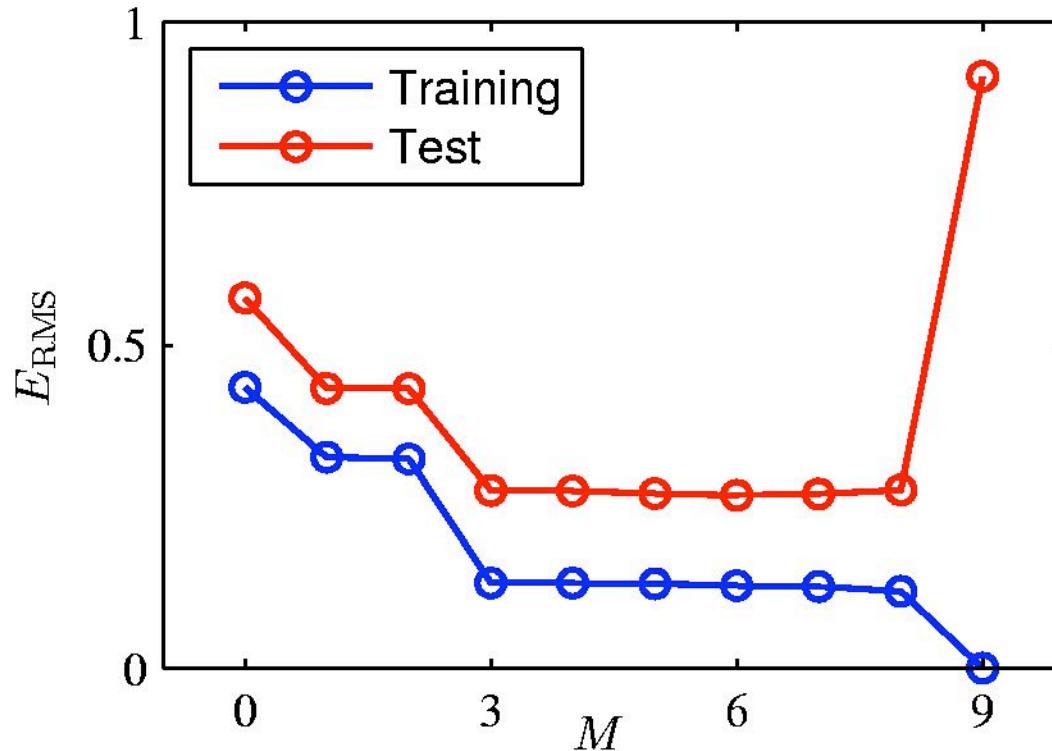
Probability & Bayesian Inference



Over-fitting

178

Probability & Bayesian Inference



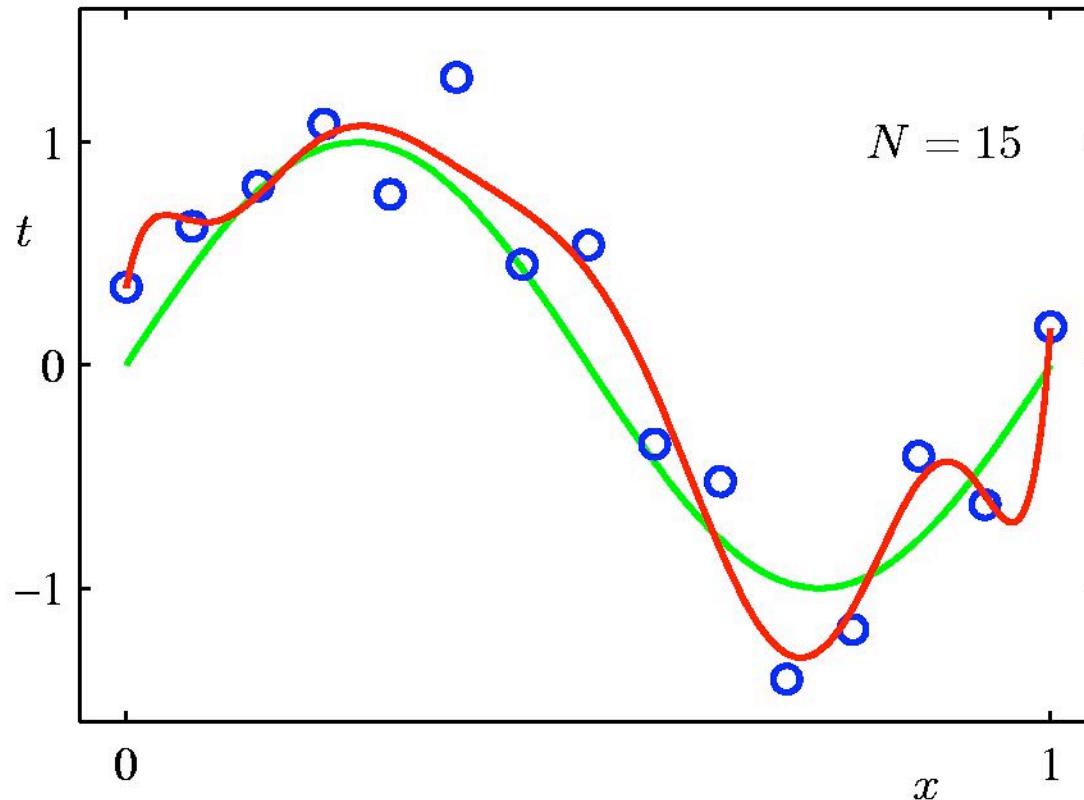
Root-Mean-Square (RMS) Error: $E_{\text{RMS}} = \sqrt{2E(\mathbf{w}^*)/N}$

Overfitting and Sample Size

179

Probability & Bayesian Inference

9th Order Polynomial

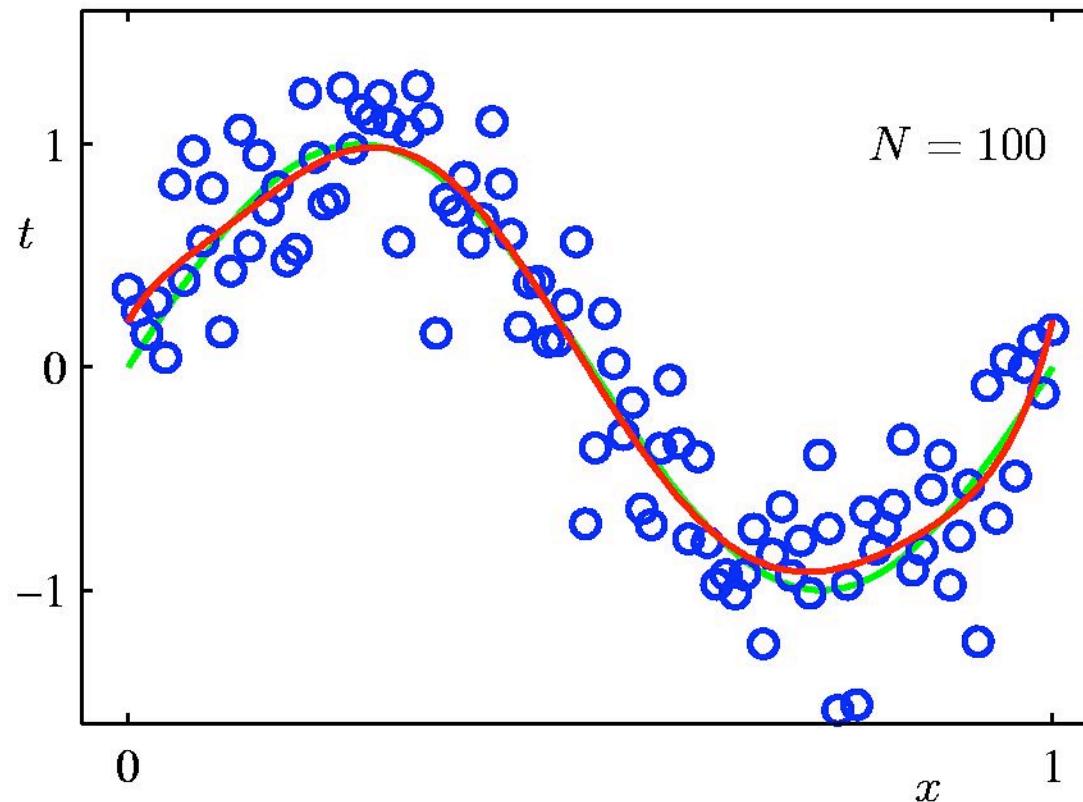


Over-fitting and Sample Size

180

Probability & Bayesian Inference

9th Order Polynomial



Methods for Preventing Over-Fitting

181

Probability & Bayesian Inference

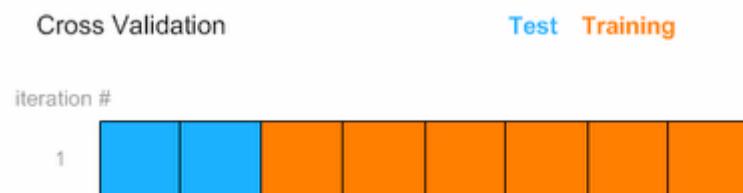
- Bayesian parameter estimation
 - ▣ Application of prior knowledge regarding the probable values of unknown parameters can often limit over-fitting of a model
- Model selection criteria
 - ▣ Methods exist for comparing models of differing complexity (i.e., with different types and numbers of parameters)
 - Bayesian Information Criterion (BIC)
 - Akaike Information Criterion (AIC)
- Cross-validation
 - ▣ This is a very simple method that is universally applicable.

Cross-Validation

182

Probability & Bayesian Inference

- The available data are partitioned into disjoint training and test subsets.
- Parameters are learned on the training sets.
- Performance of the model is then evaluated on the test set.
- Since the test set is independent of the training set, the evaluation is fair: models that overlearn the noise in the training set will perform poorly on the test set.

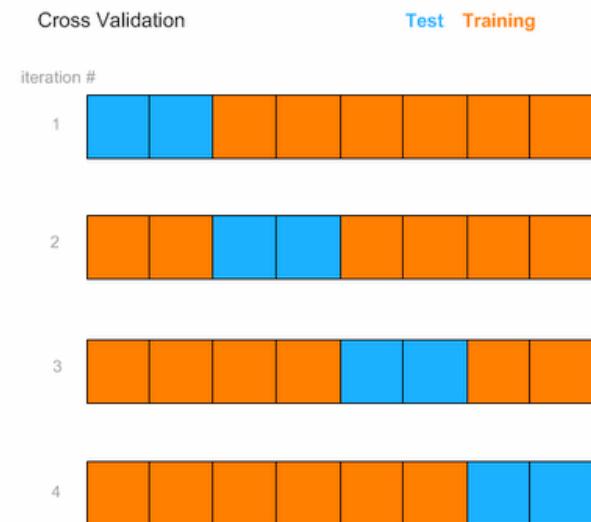


Cross-Validation: Choosing the Partition

183

Probability & Bayesian Inference

- What is the best way to partition the data?
 - A larger training set will lead to more accurate parameter estimation.
 - However a small test set will lead to a noisy performance score.
 - If you can afford the computation time, repeat the training/test cycle on complementary partitions and then average the results. This gives you the best of all worlds: accurate parameter estimation and accurate evaluation.
 - In the limit: the **leave-one-out method**



Bayesian Decision Theory: Topics

184

Probability & Bayesian Inference

1. Probability
2. The Univariate Normal Distribution
3. Bayesian Classifiers
4. Minimizing Risk
5. The Multivariate Normal Distribution
6. Decision Boundaries in Higher Dimensions
7. Parameter Estimation
8. Mixture Models and EM
9. Nonparametric Density Estimation
10. Training and Evaluation Methods
11. **What are Bayes Nets?**

10. What are Bayes Nets?

Directed Graphical Models and the Role of Causality

186

Probability & Bayesian Inference

- Bayes nets are directed acyclic graphs in which each node represents a random variable.
- Arcs signify the existence of direct causal influences between linked variables.
- Strengths of influences are quantified by conditional probabilities

$$p(\mathbf{x}) = \prod_{k=1}^K p(x_k | \text{pa}_k)$$

where pa_k is the set of 'parent' nodes of node k .

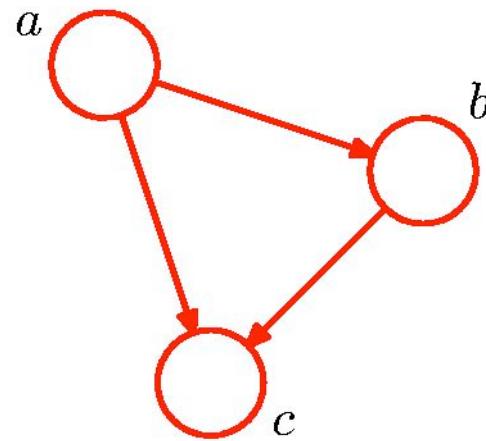
- NB: For this to hold it is critical that the graph be acyclic.

Bayesian Networks

187

Probability & Bayesian Inference

□ Directed Acyclic Graph (DAG)



From the definition of conditional probabilities (product rule):

$$p(a, b, c) = p(c|a, b)p(a, b) = p(c|a, b)p(b|a)p(a)$$

In general:

$$p(x_1, \dots, x_K) = p(x_K|x_1, \dots, x_{K-1}) \dots p(x_2|x_1)p(x_1)$$

This corresponds to a complete graph.

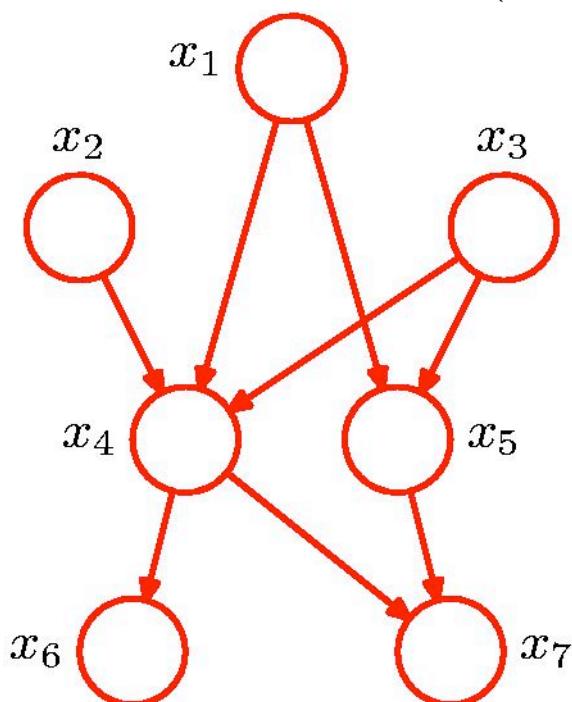
Bayesian Networks

188

Probability & Bayesian Inference

- However, many systems have sparser causal relationships between their variables.

$$p(x_1, \dots, x_7) = p(x_1)p(x_2)p(x_3)p(x_4|x_1, x_2, x_3) \\ p(x_5|x_1, x_3)p(x_6|x_4)p(x_7|x_4, x_5)$$



General Factorization

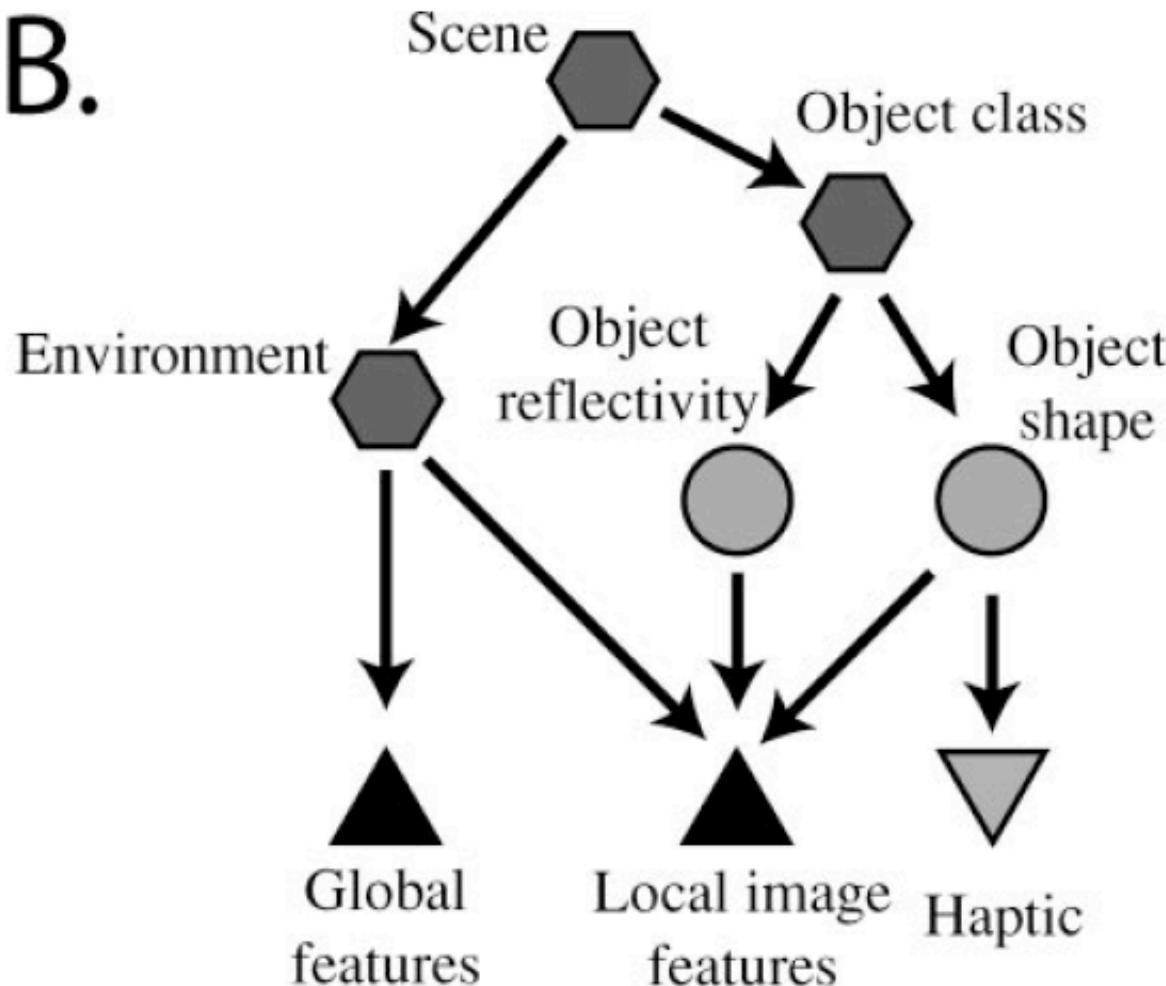
$$p(\mathbf{x}) = \prod_{k=1}^K p(x_k | \text{pa}_k)$$

Generative Models of Perception

189

Probability & Bayesian Inference

B.



Discrete Variables

190

Probability & Bayesian Inference

- General joint distribution: $K^2 - 1$ parameters



$$p(\mathbf{x}_1, \mathbf{x}_2 | \boldsymbol{\mu}) = \prod_{k=1}^K \prod_{l=1}^K \mu_{kl}^{x_{1k} x_{2l}}$$

- Independent joint distribution: $2(K - 1)$ parameters



$$\hat{p}(\mathbf{x}_1, \mathbf{x}_2 | \boldsymbol{\mu}) = \prod_{k=1}^K \mu_{1k}^{x_{1k}} \prod_{l=1}^K \mu_{2l}^{x_{2l}}$$

Discrete Variables

191

Probability & Bayesian Inference

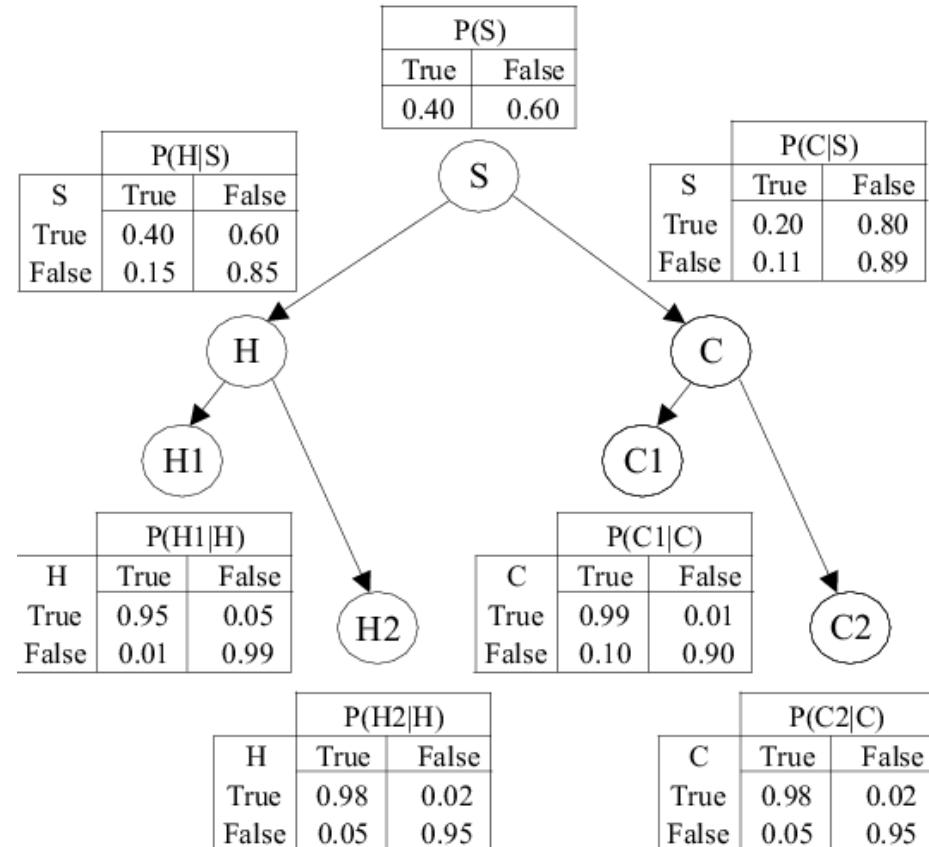
- General distributions require many parameters.
- General joint distribution over M variables:
 $K^M - 1$ parameters
- It is thus extremely important to identify structure in the system that corresponds to a sparser graphical model and hence fewer parameters.

Binary Variable Example

192

Probability & Bayesian Inference

- S: Smoker?
- C: Cancer?
- H: Heart Disease?
- (H₁, H₂): Results of medical tests for heart disease
- (C₁, C₂): Results of medical tests for cancer



Discrete Variables

193

Probability & Bayesian Inference

- Example: M -node Markov chain
 - $K - 1 + (M - 1) K(K - 1)$ parameters



Using Bayes Nets

194

Probability & Bayesian Inference

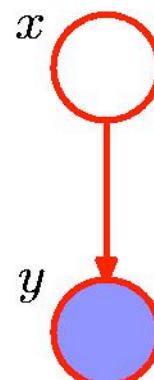
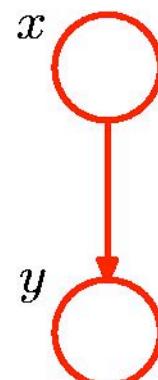
- Once a DAG has been constructed, the joint probability can be obtained by multiplying the marginal (root nodes) and the conditional (non-root nodes) probabilities.
- Training:** Once a topology is given, probabilities are estimated via the training data set. There are also methods that learn the topology.
- Probability Inference:** This is the most common task that Bayesian networks help us to solve efficiently. Given the values of some of the variables in the graph, known as evidence, the goal is to compute the conditional probabilities for some of the other variables, given the evidence.

Inference in Bayes Nets

195

Probability & Bayesian Inference

- In inference, we clamp some of the variables to observed values, and then compute the posterior over other, unobserved variables.
- Simple example:



$$p(y) = \sum_{x'} p(y|x')p(x')$$

$$p(x|y) = \frac{p(y|x)p(x)}{p(y)}$$

Example

196

Probability & Bayesian Inference

$$P(x1)=0.60 \quad P(y1|x1)=0.40 \quad P(z1|y1)=0.25 \quad P(w1|z1)=0.45$$

$$P(y1|x0)=0.30 \quad P(z1|y0)=0.60 \quad P(w1|z0)=0.30$$



$$P(x0)=0.40 \quad P(y0|x1)=0.60 \quad P(z0|y1)=0.75 \quad P(w0|z1)=0.55$$

$$P(y0|x0)=0.70 \quad P(z0|y0)=0.40 \quad P(w0|z0)=0.70$$

$$P(y1)=0.36 \quad P(z1)=0.47 \quad P(w1)=0.37$$

$$P(y0)=0.64 \quad P(z0)=0.53 \quad P(w0)=0.63$$

a) Suppose x has been measured and its value is 1. What is the probability that w is 0?

b) Suppose w is measured and its value is 1. What is the probability that x is 0?

Message Passing

197

Probability & Bayesian Inference

- For a), computation **propagates** from node x to node w , resulting in $P(w0|x1) = 0.63$.
- For b), computation **propagates** in the opposite direction, resulting in $P(x0|w1) = 0.4$.
- In general, the required inference information is computed via a combined process of “**message passing**” among the nodes of the DAG.
- **Complexity:**
 - For singly connected graphs, message passing algorithms amount to a complexity **linear** in the **number of nodes**.

Bayesian Decision Theory: Topics

198

Probability & Bayesian Inference

1. Probability
2. The Univariate Normal Distribution
3. Bayesian Classifiers
4. Minimizing Risk
5. The Multivariate Normal Distribution
6. Decision Boundaries in Higher Dimensions
7. Parameter Estimation
8. Mixture Models and EM
9. Nonparametric Density Estimation
10. What are Bayes Nets?