



UNSW
SYDNEY

Australia's
Global
University

COMP9417 – Machine Learning and Data Mining

W12 – Probabilistic Graphical Models

Instructor: Edwin V. Bonilla

School of Computer Science and Engineering

May 23rd, 2017

Acknowledgements

Material derived from:

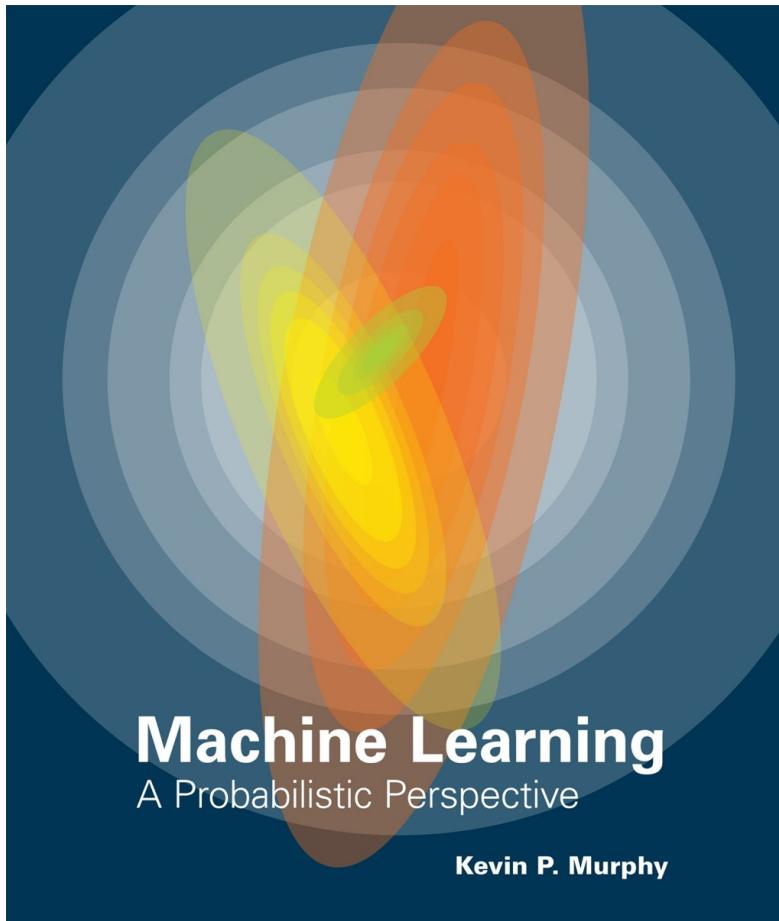
- [Barber, BRML, 2012] Bayesian Reasoning and Machine Learning, David Barber, 2012
- [Bishop, PRML, 2006] Pattern Recognition and Machine Learning, Christopher Bishop, 2006
- [Bonilla, PML, 2017] Slides from course to industry, Edwin V. Bonilla, © 2010 – 2017

Aims

This lecture will reinforce the main concepts in popular machine learning methods for dimensionality reduction, clustering and classification and will explain these methods from a probabilistic graphical model perspective. Following it you should be able to:

- Apply linear dimensionality reduction methods such as PCA and understand them as the limiting case of a linear Gaussian model ([PPCA](#)).
- Apply k-means to clustering problems and its probabilistic counterpart [Gaussian Mixture Models](#).
- Differentiate generative and discriminative approaches to [probabilistic classification](#).
- Apply Naive Bayes and class-conditional Gaussian models to classification problems.
- Estimate the parameters of the above probabilistic models.

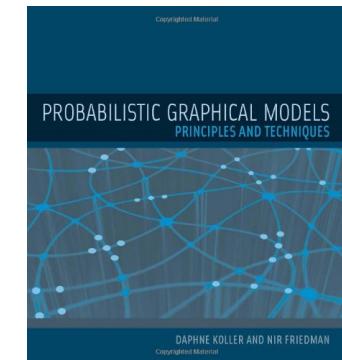
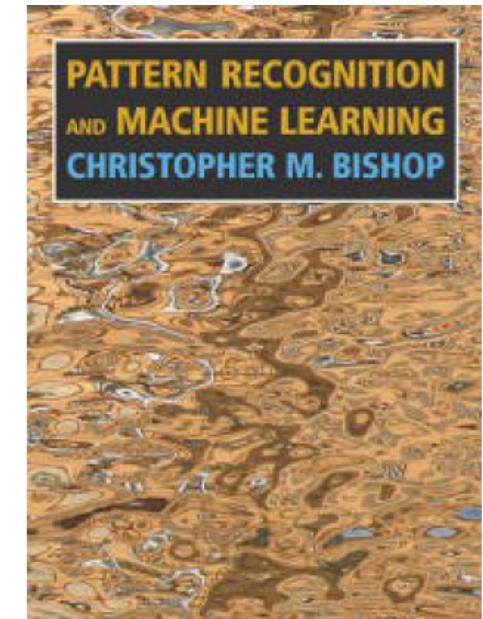
Probabilistic Graphical Models – Best Books Available



[Code available](#)

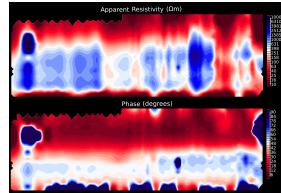


[Fully available online](#)
[Code available](#)

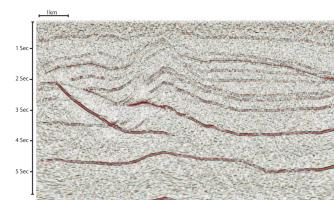


Why Probabilistic Models: Example

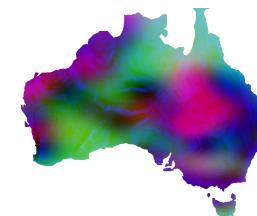
Data Fusion for Geothermal Energy Exploration



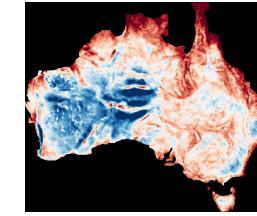
Magnetotellurics



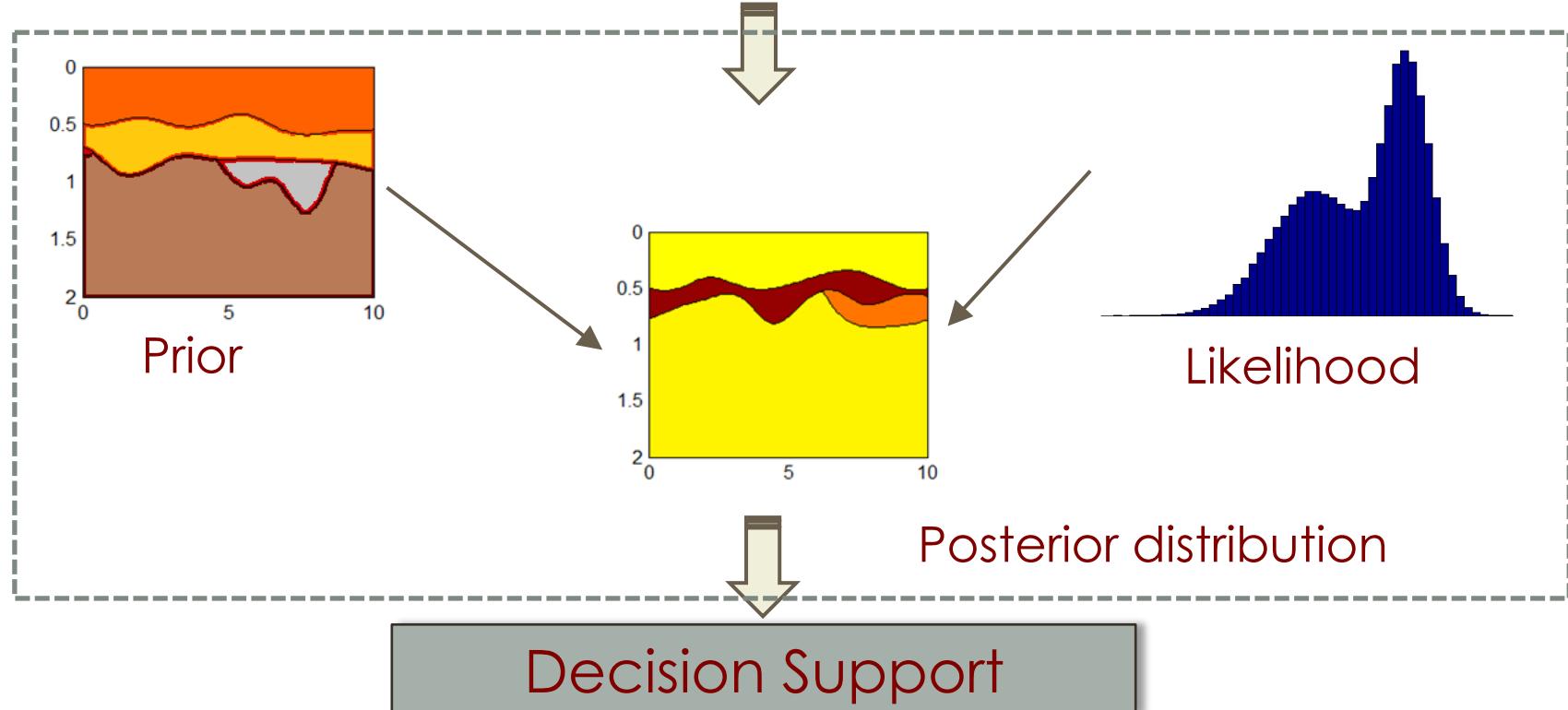
Seismic



Temperature



Gravity



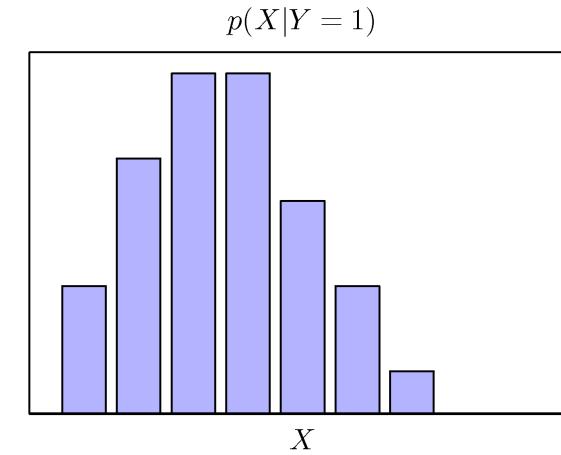
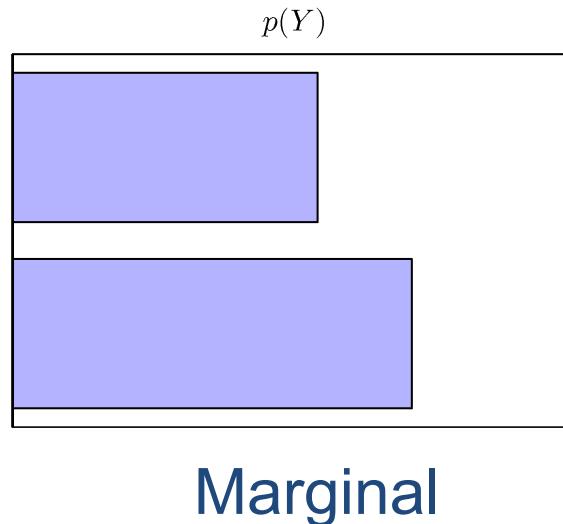
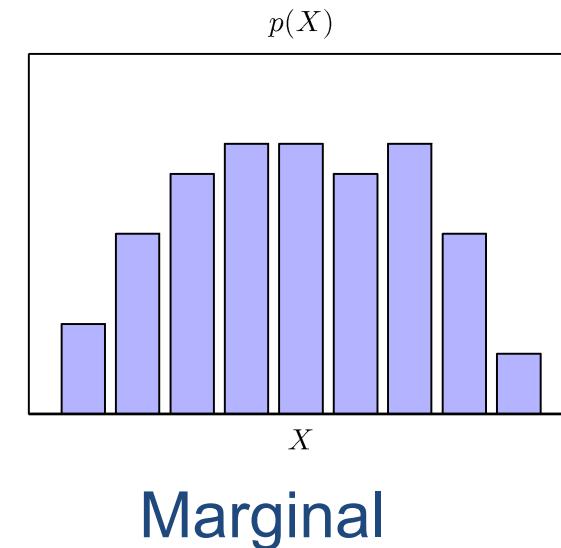
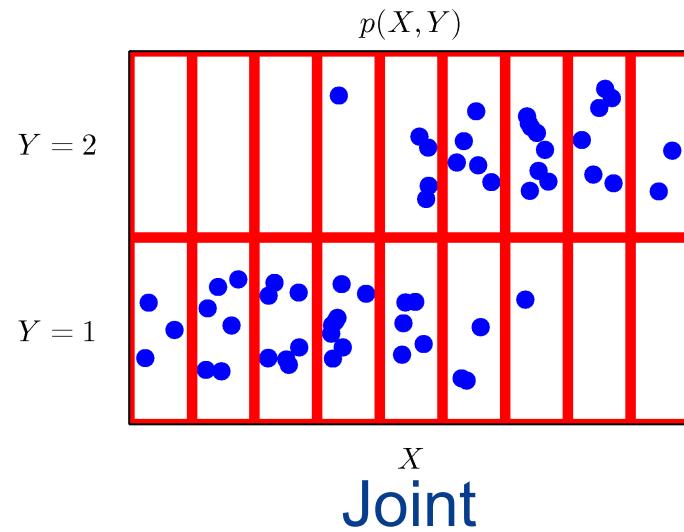
Types of Learning

- Supervised Learning: There is a “teaching signal”
 - Classification, regression, time series prediction
- Unsupervised Learning: Descriptive modelling
 - Dimensionality reduction, visualisation, clustering, latent variable modelling, density estimation, time series modelling
- Reinforcement Learning: Delayed reward (indirect supervisory signal)

We will study **these topics** from a (probabilistic) graphical models perspective

Preliminaries

Probability Distributions: Illustration



The Rules of Probability and Terminology

Discrete Case

- Sum Rule: $P(X) = \sum_Y P(X, Y)$
- Product Rule: $P(X, Y) = P(Y|X)P(X)$
- By symmetry: $P(Y, X) = P(X, Y)$
- Therefore: $P(X) = \sum_Y P(X, Y) = \sum_Y P(X|Y)P(Y)$
- Normalisation, marginalisation

Bayes' Theorem

Bayesian Inference

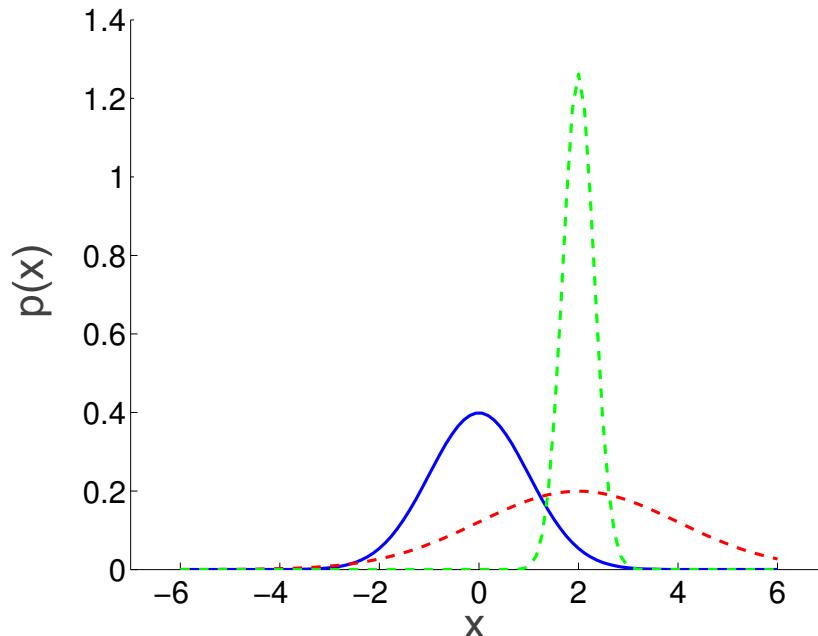
Bayesian inference provides us with a mathematical framework explaining how to change our (prior) beliefs in the light of new evidence.

$$\underbrace{p(Z|X)}_{\text{posterior}} = \frac{\overbrace{p(X|Z)}^{\text{likelihood}} \overbrace{p(Z)}^{\text{prior}}}{\underbrace{p(X)}_{\text{evidence}}}$$

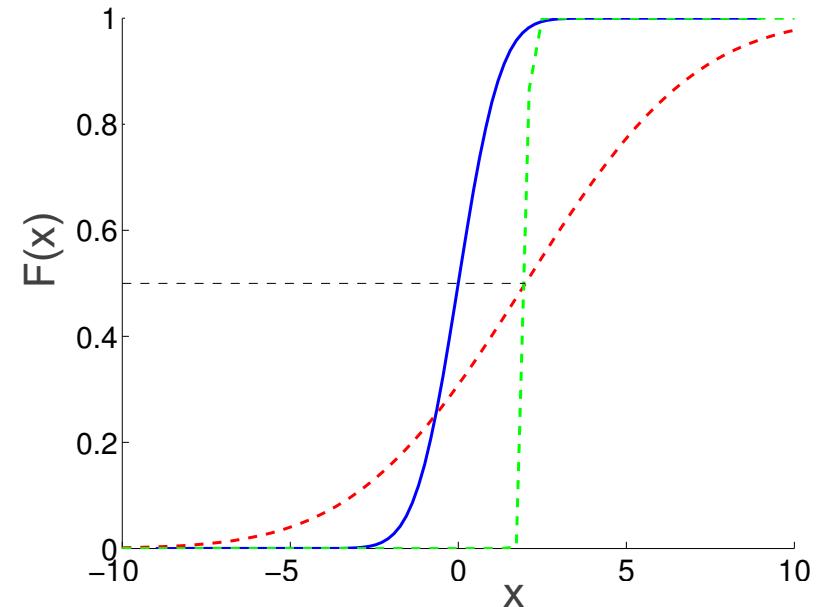
$$p(Z|X) = \frac{p(X|Z)p(Z)}{\sum_{Z'} p(X|Z')p(Z')}$$

Continuous Random Variables

Probability Density Function (pdf)



Cumulative Distribution Function (cdf)



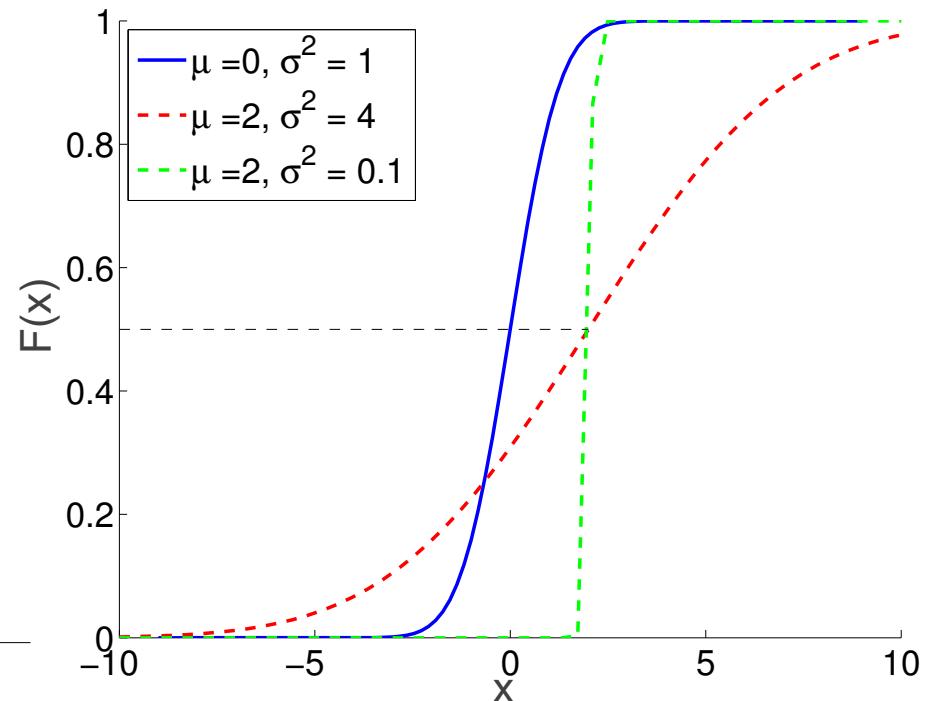
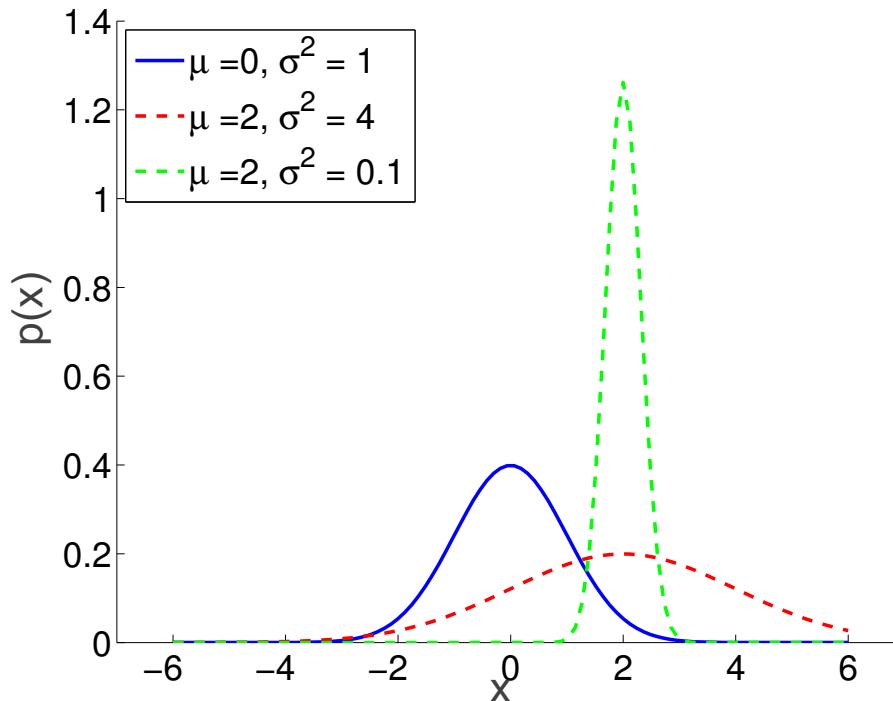
$$p(x) \geq 0, \quad \int_{-\infty}^{\infty} p(x)dx = 1$$

$$p(a < x < b) = \int_a^b p(x)dx$$

$$\begin{aligned} F(x) &= p(X \leq x) \\ &= \int_{-\infty}^x p(z)dz \end{aligned}$$

All rules of probability apply where we replace sums with integrals

The Gaussian Distribution – 1D Example



$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right)$$

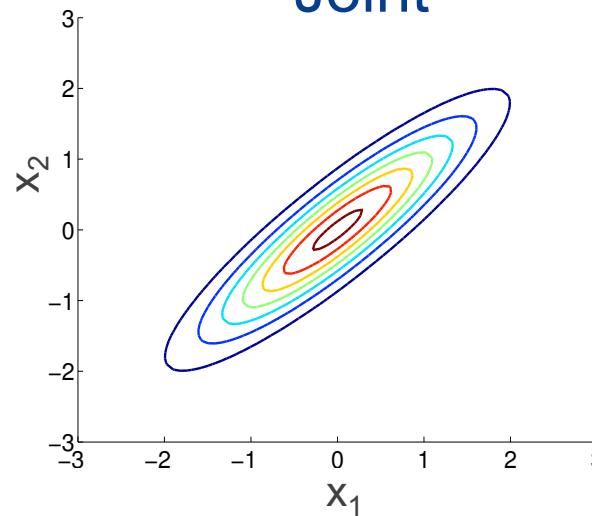
Parameters

In $D > 1$ dimensions:

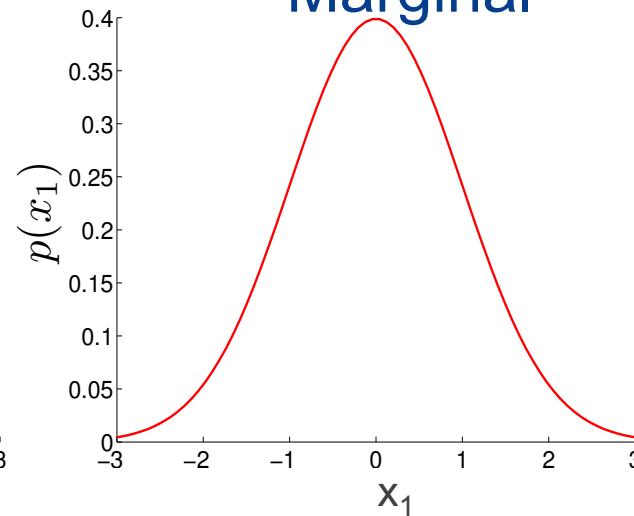
$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{|2\pi\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right)$$

The Gaussian Distribution – 2D Example

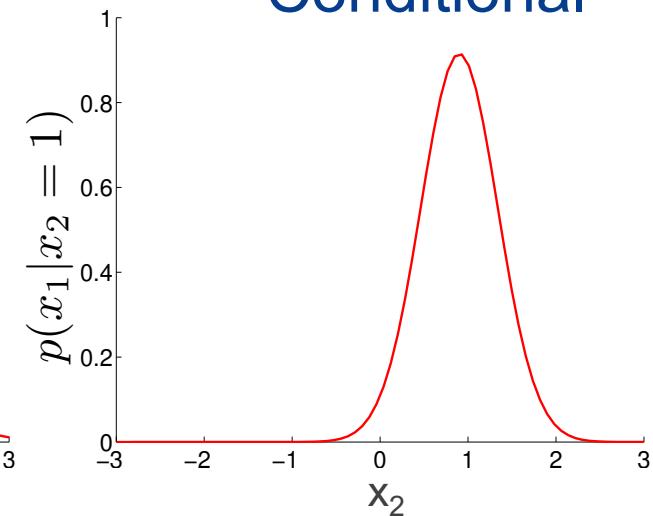
Joint



Marginal



Conditional



$$p(x_1, x_2) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \rightarrow \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{12}^T & \boldsymbol{\Sigma}_{22} \end{bmatrix}\right)$$

The marginal and the conditional distributions are also Gaussians:

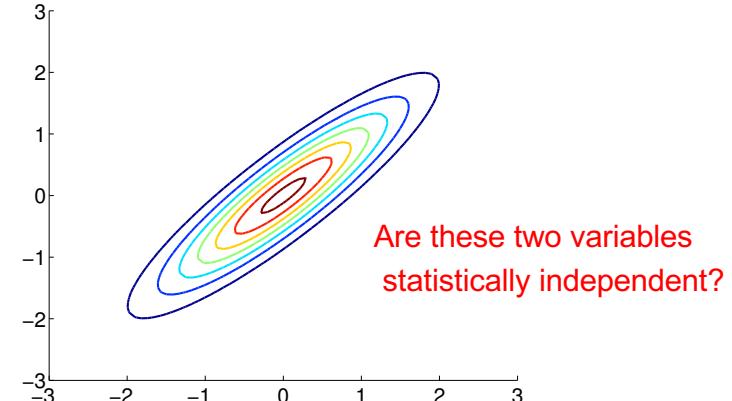
$$p(\mathbf{x}_1) = \mathcal{N}(\mathbf{x}_1|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11})$$

$$p(\mathbf{x}_1|\mathbf{x}_2) = \mathcal{N}(\mathbf{x}_1|\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2), \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{12}^T)$$

Covariance and Precision Matrices

Σ : is the covariance matrix

Σ^{-1} : is the precision matrix



- An entry $\Sigma_{ij} = 0$ indicates that the variables i and j are marginally independent
- An entry $(\Sigma^{-1})_{ij} = 0$ indicates that the variables i and j are conditionally independent given all the other variables
- Marginalizing out a variable leaves Σ unchanged but changes Σ^{-1}

In general, there is a difference between statistical independence and correlation

Parameter Estimation for The Gaussian Distribution – Maximum Likelihood Training (1)

We are given a dataset $\mathcal{D} = \{\mathbf{x}^{(n)}\}_{n=1}^N$ with $\mathbf{x} \in \mathbb{R}^D$ and we want to fit a multivariate Gaussian $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$

Goal: Find the best parameters $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ that fit the data.

Assuming iid observations, the data log-likelihood can be written as:

$$L(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \prod_{n=1}^N p(\mathbf{x}^n | \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$\begin{aligned} \mathcal{L}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \log L(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \log p(\mathbf{x}^n | \boldsymbol{\mu}, \boldsymbol{\Sigma}) \\ &= -\frac{1}{2} \sum_{n=1}^N (\mathbf{x}^n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}^n - \boldsymbol{\mu}) - \frac{N}{2} \log |\boldsymbol{\Sigma}| \end{aligned}$$

Parameter Estimation for The Gaussian Distribution – Maximum Likelihood Training (2)

The **optimal** μ and Σ can be obtained by computing the corresponding derivatives of $L(\mu, \Sigma)$ and setting them to zero:

- For the **mean** we have:

$$\frac{\partial \mathcal{L}}{\partial \mu} = \sum_{n=1}^N \Sigma^{-1} (\mathbf{x}^n - \mu)$$

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

- Similarly, for the **covariance**:

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

These are simply the sample mean and sample covariance of the data!

Outline

I. Dimensionality Reduction

- PCA and PPCA

II. Clustering

- K-means and Gaussian Mixture Models

III. Probabilistic Classification – Generative Approaches

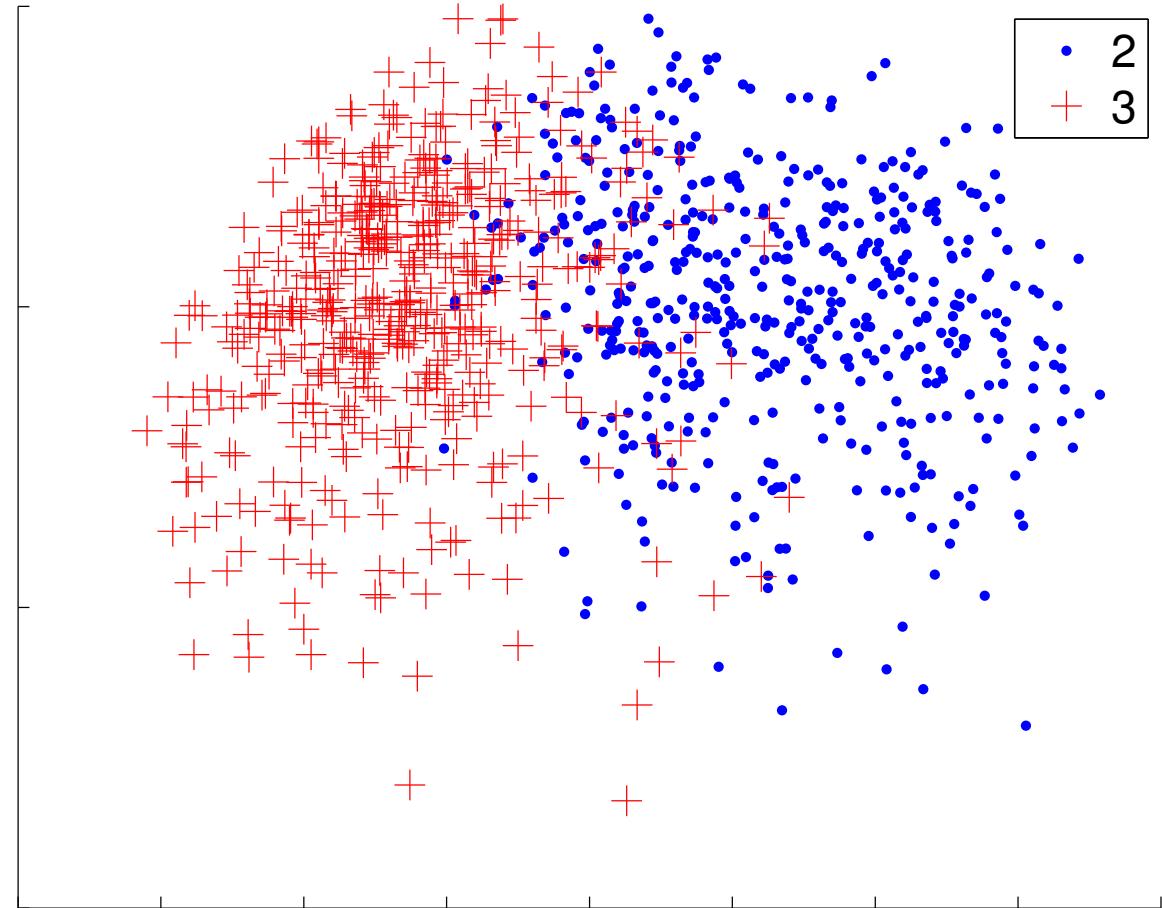
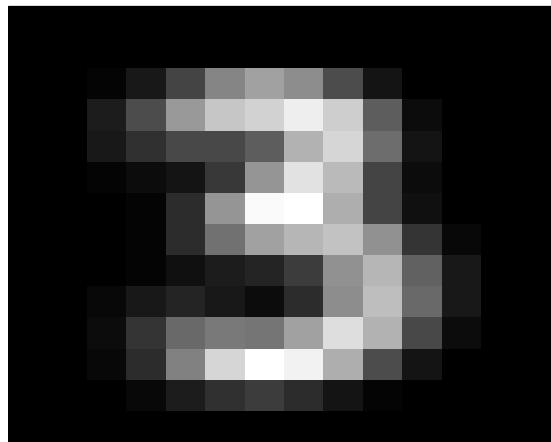
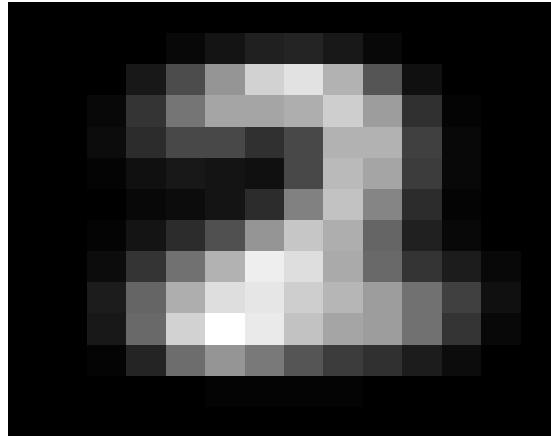
- Naïve Bayes and class-conditional Gaussian models

I. Dimensionality Reduction: Probabilistic Principal Component Analysis (PPCA)

Low Dimensional Manifolds

High-dimensional data may lie in low-dimensional manifolds

- Only a small number of dimensions (degrees of freedom) may be sufficient to explain well the data



Why Dimensionality Reduction?

- Learning in high-dimensional spaces may be very hard
- We can use low-dimensional representations for visualisation
- These representations can be more discriminative
- We will focus on **unsupervised** and **linear** techniques for dimensionality reduction
 - Note that we may be “throwing” away useful information

Principal Component Analysis (PCA)

Problem Definition

Consider a set of N datapoints $\{\mathbf{x}^{(n)}\}_{n=1}^N$ with $\mathbf{x} \in \mathbb{R}^D$.

- Our goal is to find a low-dimensional linear representation $\mathbf{z}^{(n)} \in \mathbb{R}^K$ with $K < D$ that **explains the data well**.
- We can have (at least) two criteria:
 1. Find an orthogonal projection such as the **variance** of the projected data is maximized.
 2. Find a linear projection that minimizes the **reconstruction** error (squared error).
- These criteria lead to the same solution given by PCA

Principal Component Analysis (PCA)

The Algorithm

Input: D-dimensional data $\{\mathbf{x}^{(n)}\}_{n=1}^N, K$

1. Compute sample mean $\bar{\mathbf{x}}$ and covariance \mathbf{S}
2. Find K eigenvectors $\mathbf{e}^1, \dots, \mathbf{e}^K$ corresponding to the largest K eigenvalues $\lambda_1, \dots, \lambda_K$ and construct the matrix

$$\mathbf{E} = (\mathbf{e}^1, \dots, \mathbf{e}^K)$$

3. **Output:** Lower dimensional representation is given by:

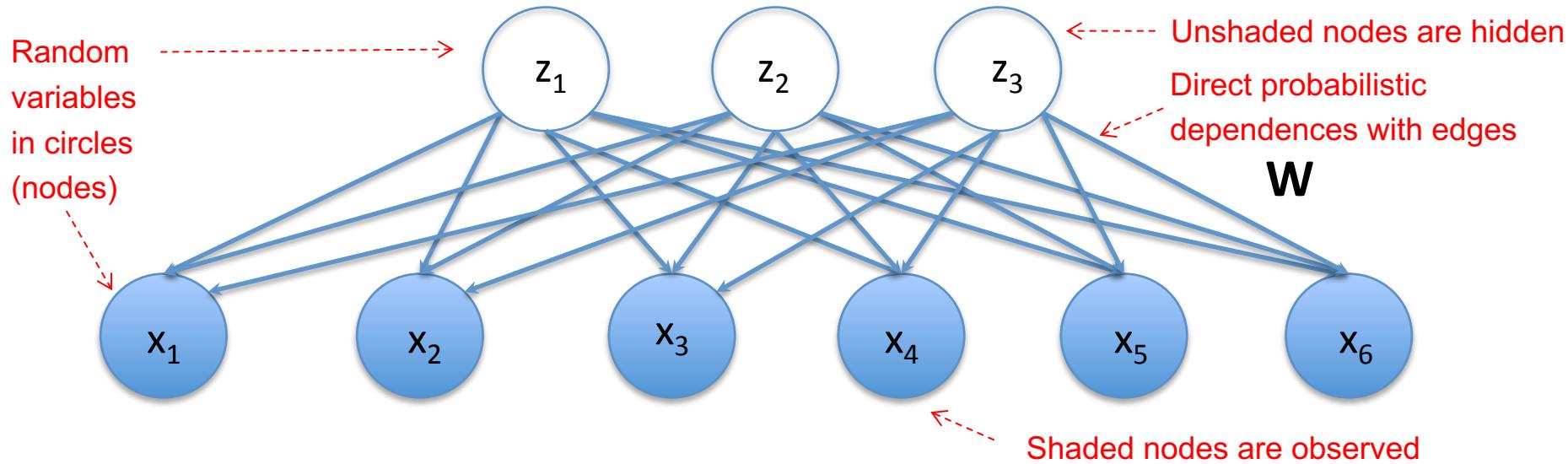
$$\mathbf{z}^n = \mathbf{E}^T(\mathbf{x}^n - \bar{\mathbf{x}})$$

How much variance have we explained by selecting K components?

- HINT: We can use the fact that the total variance is given by $\sum_{i=1}^D \lambda_i$
- We can use this to establish a criterion to select K in an unsupervised setting.

Probabilistic Principal Component Analysis (PPCA)

In fact, PCA can be derived from a proper probabilistic latent variable model:



We observe $\mathbf{x} \in \mathbb{R}^D$ and assume that it is generated by a latent low-dimensional vector $\mathbf{z} \in \mathbb{R}^K$:
 $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I}), \quad p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I})$

- The $D \times K$ matrix \mathbf{W} , the D -dimensional bias vector $\boldsymbol{\mu}$ and the noise variance σ^2 are the model parameters
 - How do we learn these parameters?

PPCA: Parameter Estimation

Maximum Likelihood Parameters

PPCA defines $p(\mathbf{z}, \mathbf{x})$ but \mathbf{z} are latent as we only observe $\{\mathbf{x}^{(n)}\}_{n=1}^N$

- Need to marginalise \mathbf{z} and get $p(\mathbf{x}) \rightarrow$ Analytical solution!

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}),$$

Assuming iid observations, the data log-likelihood is:

$$\mathcal{L}(\mathbf{W}, \boldsymbol{\mu}, \sigma^2) = \sum_{n=1}^N \log p(\mathbf{x}^n | \mathbf{W}, \boldsymbol{\mu}, \sigma^2)$$

and Maximizing \mathcal{L} wrt $\boldsymbol{\mu}$, \mathbf{W} and σ^2 we obtain:

$$\hat{\boldsymbol{\mu}}_{\text{ML}} = \bar{\mathbf{x}} \quad \hat{\mathbf{W}}_{\text{ML}} = \mathbf{E}(\Lambda - \sigma_{\text{ML}}^2 \mathbf{I})^{1/2} \mathbf{R}$$

↑ ↑ ↑
 Ordered eigenvalues and
eigenvectors of data covariance Rotation matrix, e.g. $\mathbf{R} = \mathbf{I}$
↑ ↑
 $\hat{\sigma}_{\text{ML}}^2 = \frac{1}{D-K} \sum_{i=K+1}^D \lambda_i$
 Average variance of
Discarded dimensions

PPCA

How to Do Dimensionality Reduction

Probabilistic inference problem

- Compute $p(\mathbf{z} | \mathbf{x}) \rightarrow$ Analytical solution!

$$p(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\mathbf{z} | \mathbf{M}^{-1} \mathbf{W}^T (\mathbf{x} - \boldsymbol{\mu}), \sigma^{-2} \mathbf{M}), \quad \mathbf{M} = \mathbf{W}^T \mathbf{W} + \sigma^2 \mathbf{I}.$$

Using the ML estimators, the posterior mean is:

$$\mathbb{E}[\mathbf{z} | \mathbf{x}] = (\widehat{\mathbf{W}}_{\text{ML}}^T \widehat{\mathbf{W}}_{\text{ML}} + \widehat{\sigma}_{\text{ML}}^2 \mathbf{I})^{-1} \widehat{\mathbf{W}}_{\text{ML}}^T (\mathbf{x} - \bar{\mathbf{x}})$$

- and taking the limit $\sigma^2 \rightarrow 0$ we obtain PCA, an **orthogonal projection** of the datapoints into the latent space.

Advantages of probabilistic approach:

- Proper likelihood function
 - Mixture models and Bayesian extensions
- Can find dimensionality automatically

PCA & PPCA

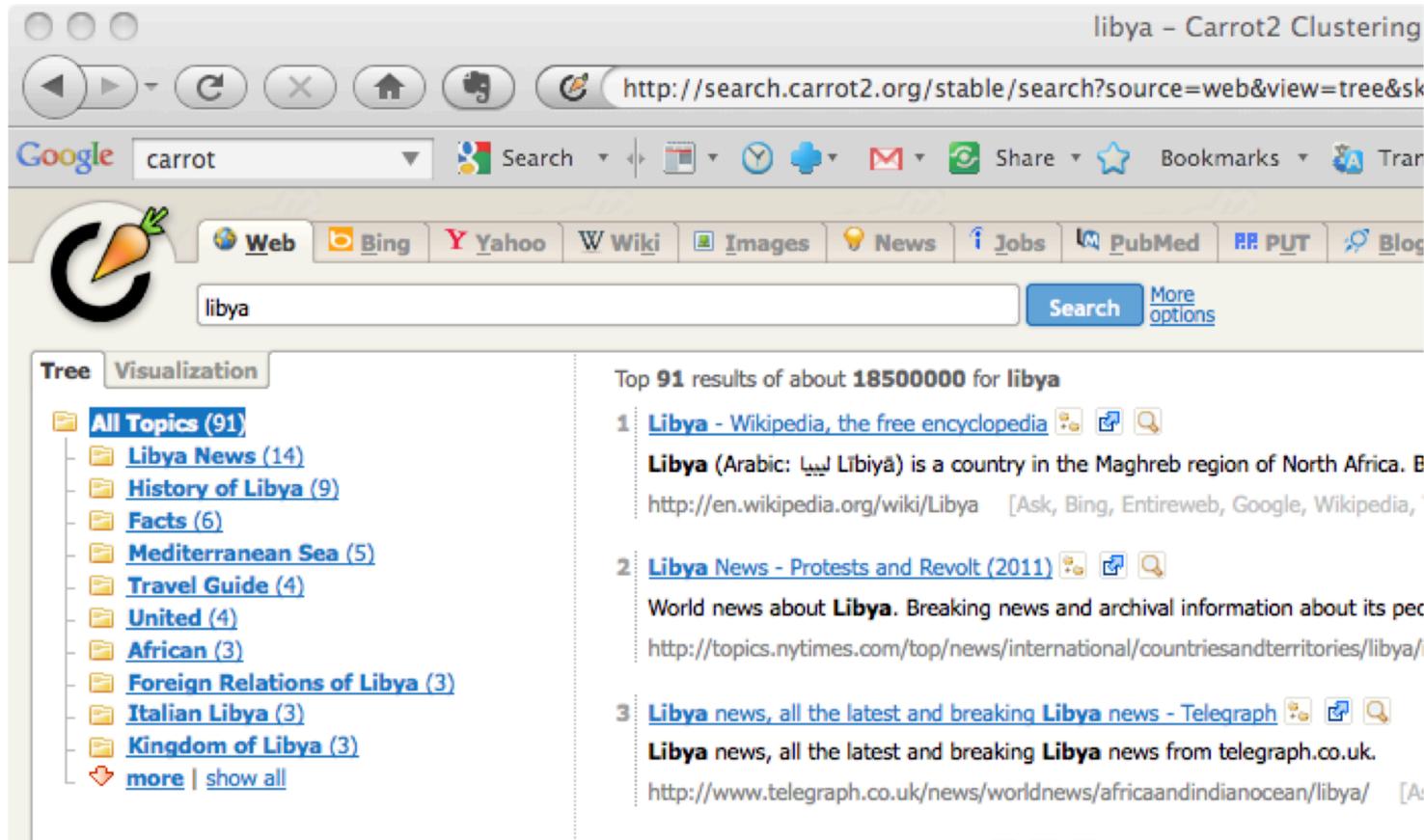
Final Remarks

- PPCA is a constrained Gaussian model
 - $\Sigma = \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I}$
 - if $K = D$ then $\Sigma = \mathbf{S}$, the sample covariance
 - The number of independent parameters only grows linearly with D (cf. full Gaussian, isotropic Gaussian)
 - But we still capture the K most significant correlations
- Caution in supervised scenarios: We may be throwing away the important information
- We are assuming the data lies close to a hyperplane (linear technique)
 - Manifold may be curved or data can be clustered
- PCA is a very old technique with many different names, e.g. the Karhunen-Loève transform
 - Most commonly used dimensionality reduced method

II. Clustering and Mixture Models: K-Means and Gaussian Mixtures

Clustering

Goal: Given (e.g. documents), assign “similar” ones to one of K groups (clusters)



The screenshot shows a web browser window titled "libya - Carrot2 Clustering". The address bar displays the URL <http://search.carrot2.org/stable/search?source=web&view=tree&sk>. The search query "libya" is entered in the search bar. The interface includes a navigation bar with icons for back, forward, search, and refresh, and a toolbar with links for Google, carrot, Search, Share, Bookmarks, and Translate.

The main content area is divided into two tabs: "Tree" and "Visualization". The "Tree" tab is selected, showing a hierarchical list of topics under "All Topics (91)". The topics include "Libya News (14)", "History of Libya (9)", "Facts (6)", "Mediterranean Sea (5)", "Travel Guide (4)", "United (4)", "African (3)", "Foreign Relations of Libya (3)", "Italian Libya (3)", and "Kingdom of Libya (3)". There are also links for "more" and "show all".

The "Visualization" tab is visible but not active. To the right, the results for the search query "libya" are displayed. The results are titled "Top 91 results of about 18500000 for libya". The first three results are listed:

- 1 Libya - Wikipedia, the free encyclopedia**
Libya (Arabic: ليبية Libyā) is a country in the Maghreb region of North Africa. It is bordered by Egypt to the east, Sudan to the west, and the Mediterranean Sea to the north. The capital city is Tripoli. The official language is Arabic, and the primary religion is Islam. Libya has a population of approximately 6.5 million people.
- 2 Libya News - Protests and Revolt (2011)**
World news about Libya. Breaking news and archival information about its recent political crisis and revolution.
- 3 Libya news, all the latest and breaking Libya news - Telegraph**
Libya news, all the latest and breaking Libya news from telegraph.co.uk.

Clustering with K-means

Each cluster \mathcal{C}_k is associated with a prototype vector (or centre) $\boldsymbol{\mu}_k$ and each \mathbf{x}^n with \mathbf{z}^n where $z_k^n = 1$ iff \mathbf{x}^n belongs to cluster \mathcal{C}_k .

Repeat until convergence:

1. Initialize centres $\{\boldsymbol{\mu}_k\}_{k=1}^K$ (e.g. using K data vectors)
2. Set $z_k^n = 1$ iff $k = \operatorname{argmin}_j \|\mathbf{x}^n - \boldsymbol{\mu}^j\|^2$ and $z_k^n = 0$ otherwise
3. Set $\boldsymbol{\mu}_k = \frac{1}{|\mathcal{C}_k|} \sum_{n=1}^N z_k^n \mathbf{x}^n$

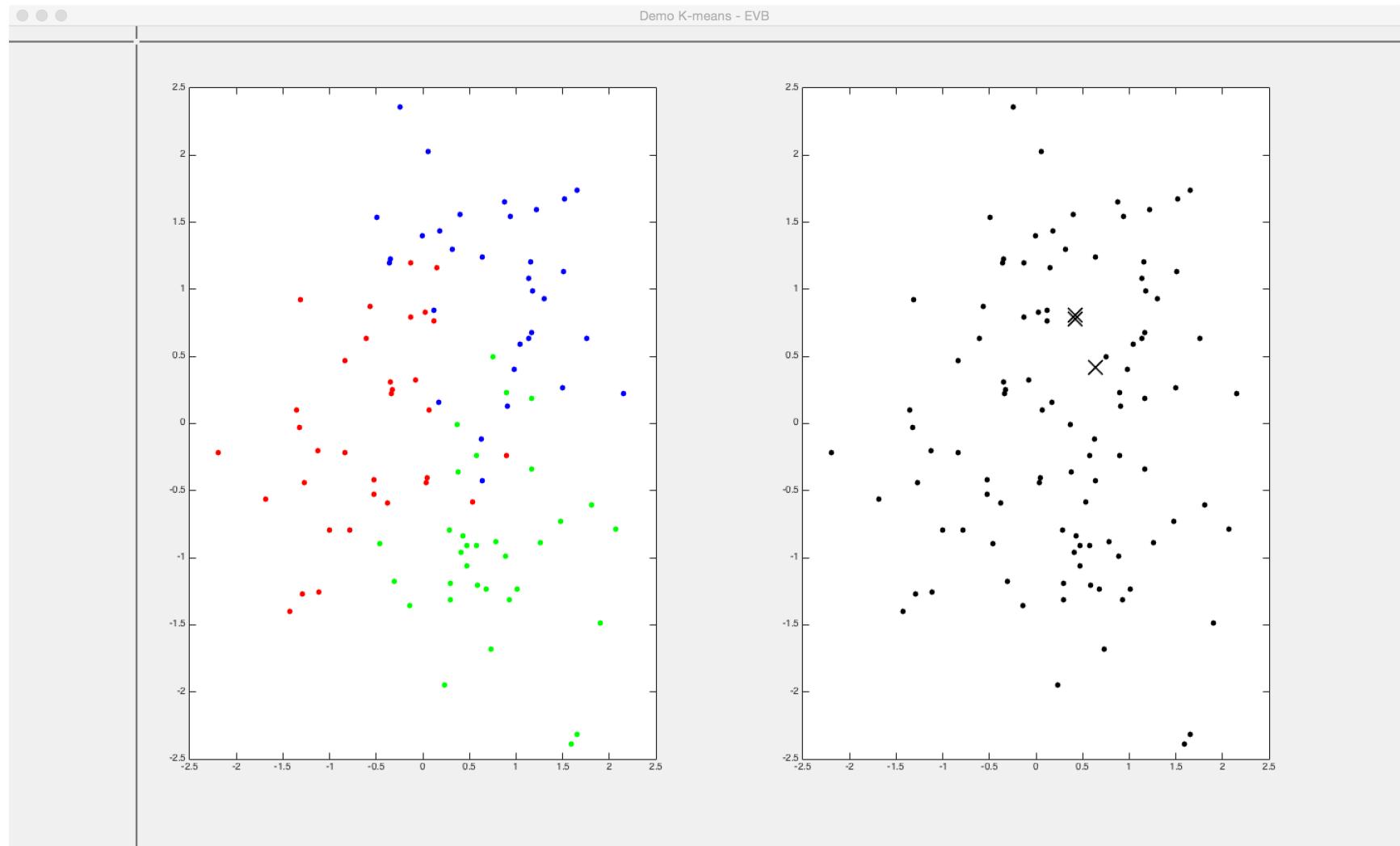
- This procedure can be shown to minimize the error function:

$$\mathcal{E} = \sum_{k=1}^K \sum_{n=1}^N z_k^n \|\mathbf{x}^n - \boldsymbol{\mu}^k\|^2$$

- Convergence only to a local minimum

K-means Clustering Demo

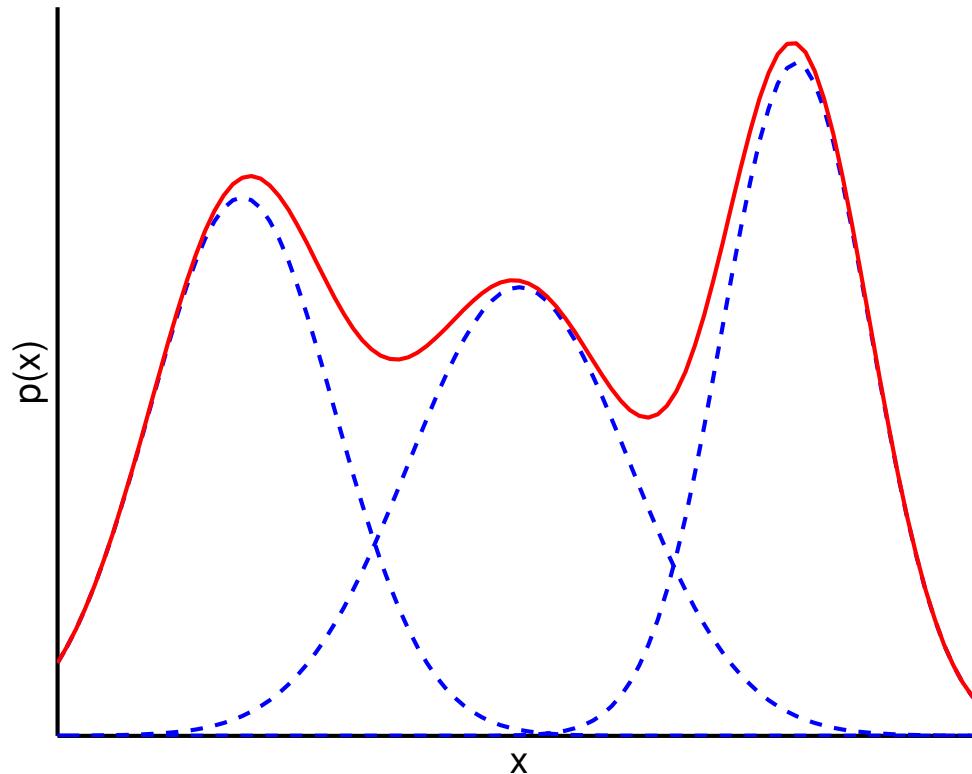
29



demo_kmeans.m

Gaussian Mixture Models (GMMs)

Modelling Complex Densities



In blue: original Gaussians scaled by their priors $p(z_k = 1)$ with $k = 1, 2, 3$

In red: the resulting mixture

The final density is multimodal

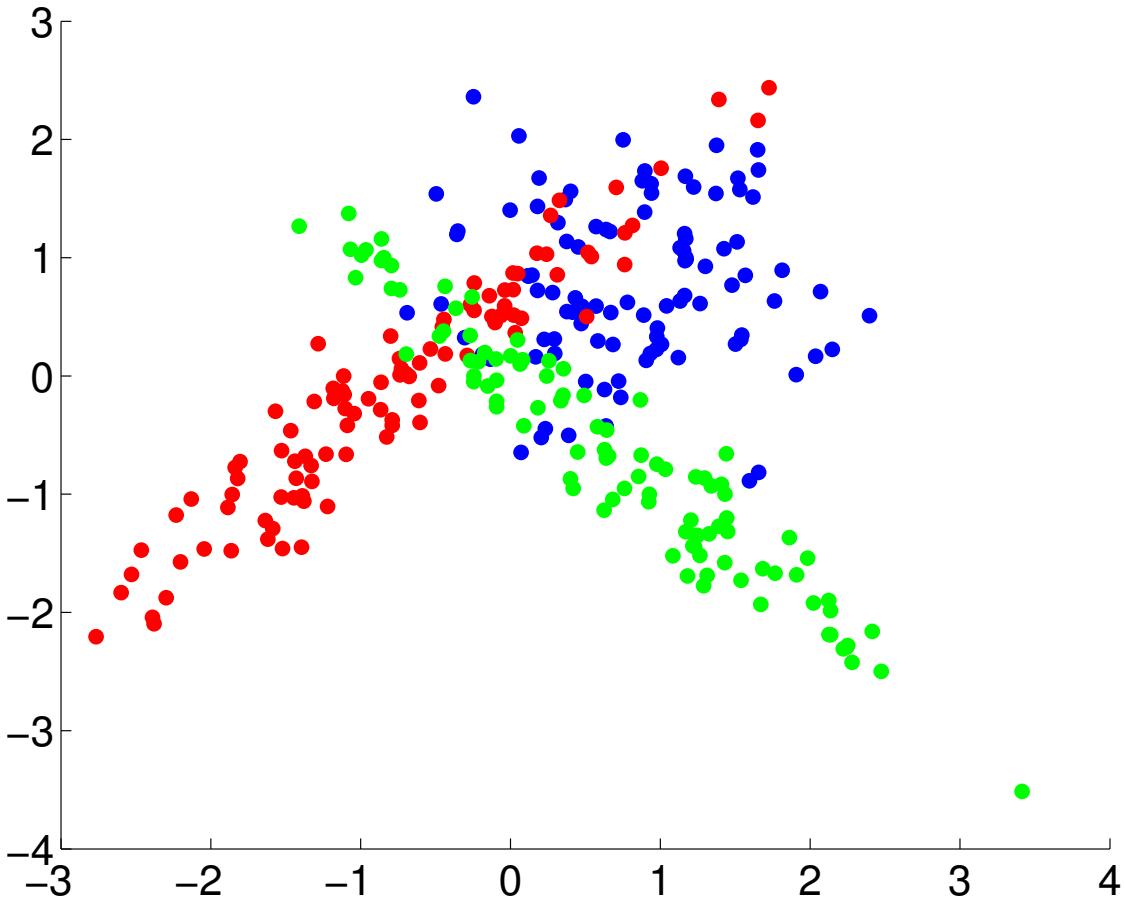
Generative Process: For $k = 1, 2, 3$

1. Pick component (cluster) k with probability π_k
2. Draw $\mathbf{x} \sim N(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$

Gaussian Mixture Models (GMMs)

31

Data Generation – How can mixtures be used for clustering?



$$p_k = \frac{1}{3}, k = 1, \dots, 3$$

$$\Sigma_1 = \begin{pmatrix} 0.5 & 0 \\ 0 & 0.5 \end{pmatrix}$$

$$\Sigma_2 = \begin{pmatrix} 1 & 0.95 \\ 0.95 & 1 \end{pmatrix}$$

$$\Sigma_3 = \begin{pmatrix} 1 & -0.95 \\ -0.95 & 1 \end{pmatrix}$$

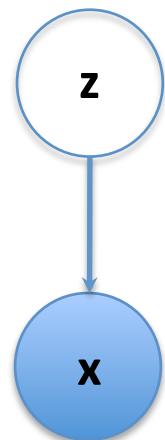
How can we represent this mixture of Gaussians with a graphical model?

Gaussian Mixture Models (GMMs)

Probabilistic Latent Variable model

In fact, K-means can be seen as a limiting case of a Gaussian mixture model:

- Given $z_k \in \{0, 1\}$ with $\sum_{k=1}^K z_k = 1$



$$p(z_k = 1) = \pi_k \text{ with } 0 \leq \pi_k \leq 1, \quad \sum_{k=1}^K \pi_k = 1, \quad p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k}$$

$$p(\mathbf{x}|z_k = 1) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k), \quad p(\mathbf{x}|\mathbf{z}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k}$$

Marginal:

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

How can we learn $\{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}$?

Gaussian Mixture Models (GMMs)

Parameter Estimation – Direct Likelihood Maximisation

- We can write down the data log-likelihood:

$$\mathcal{L}(\{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K) = \sum_{n=1}^N \log \left(\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}^n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right)$$

- We can do gradient-based optimization directly but:
 - Note logarithm of the sum, no closed form solution
 - Constraints on π_k and on $\boldsymbol{\Sigma}_k$
- We can address these by having suitable parameterizations:
 - Softmax for π_k and Cholesky decomposition for $\boldsymbol{\Sigma}_k$
- An alternative solution is obtained via the EM algorithm
 - “closed-form” updates
 - The above constraints are automatically satisfied

Gaussian Mixture Models (GMMs)

Parameter Estimation – The Expectation-Maximisation (EM) Algorithm

Consider the case that \mathbf{z} had been observed. Therefore:

$$\mathcal{L}^{\text{comp}}(\boldsymbol{\theta}) = \sum_{n=1}^N \log p(\mathbf{x}^n, \mathbf{z}^n | \boldsymbol{\theta}),$$

where $\mathcal{L}^{\text{comp}}$ refers to the **complete data log-likelihood** and $\boldsymbol{\theta}$ are the model parameters.

- The log of sum does not appear!
 - We **would not have** to marginalize the latent variable \mathbf{z}

But life is hard and \mathbf{z} has not been observed. What can we do?

1. Compute its **expectation** over the posterior $\langle \mathcal{L}^{\text{comp}} \rangle_{p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})}$
2. **Maximize** this expectation wrt $\boldsymbol{\theta}$

The EM algorithm iterates the above steps and converges to a local minimum of the (incomplete) data log-likelihood

Gaussian Mixture Models (GMMs)

The EM Algorithm: E-step

In the E-step we need to compute the posterior $p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})$

$$\begin{aligned} p(z_k = 1|\mathbf{x}, \boldsymbol{\theta}^{\text{old}}) &= \frac{p(z_k = 1|\boldsymbol{\theta}^{\text{old}})p(\mathbf{x}|z_k = 1, \boldsymbol{\theta}^{\text{old}})}{\sum_{j=1}^K p(z_j = 1|\boldsymbol{\theta}^{\text{old}})p(\mathbf{x}|z_j = 1, \boldsymbol{\theta}^{\text{old}})} \\ &= \frac{\pi_k^{\text{old}} \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k^{\text{old}}, \boldsymbol{\Sigma}_k^{\text{old}})}{\sum_{j=1}^K \pi_j^{\text{old}} \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j^{\text{old}}, \boldsymbol{\Sigma}_j^{\text{old}})} \stackrel{\text{def}}{=} \gamma_k \end{aligned}$$

π_k is the **prior** probability of component k and γ_k is the **posterior** probability of component k given that we have observed \mathbf{x} .

- γ_k can be viewed as the **responsibility** that component k takes for explaining observation \mathbf{x} .

Gaussian Mixture Models (GMMs)

The EM Algorithm: M-step

In the M-step we need to maximise the objective function:

$$\begin{aligned} \mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) &= \left\langle \sum_{n=1}^N \log p(\mathbf{x}^n, \mathbf{z}^n | \boldsymbol{\theta}) \right\rangle_{p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta}^{\text{old}})} \\ &= \sum_{n=1}^N \sum_{k=1}^K \gamma_k^n [\log \pi_k + \log \mathcal{N}(\mathbf{x}^n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)] \end{aligned}$$

Optimising wrt $\boldsymbol{\theta}$ and defining $N_k = \sum_{n=1}^N \gamma_k^n$ we have:

$$\pi_k = \frac{N_k}{N} \quad \boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_k^n \mathbf{x}^n \quad \boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_k^n (\mathbf{x}^n - \boldsymbol{\mu}_k)(\mathbf{x}^n - \boldsymbol{\mu}_k)^T$$

- Interpretation of $N_k, \pi_k, \boldsymbol{\mu}_k$ and $\boldsymbol{\Sigma}_k$?

EM for Gaussian Mixture Models (GMMs)

The Algorithm

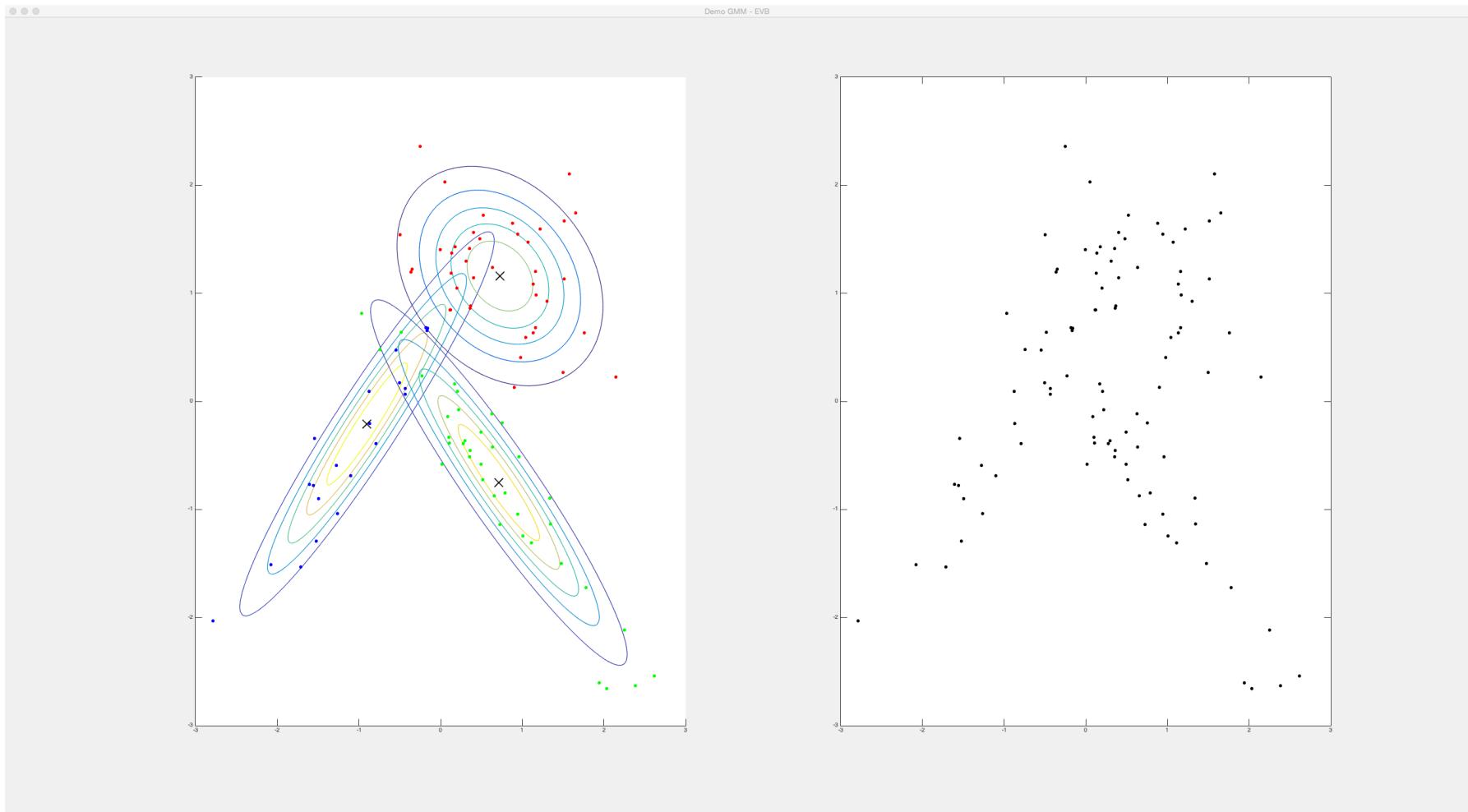
Goal: Given $\mathcal{D} = \{\mathbf{x}^n\}_{n=1}^N$ (e.g. documents) learn a GMM

Input: $\mathcal{D} = \{\mathbf{x}^n\}_{n=1}^N, K$

1. Initialise $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$ (**how?**) and evaluate data log-likelihood
2. Repeat until convergence of parameters or data log-likelihood
 - a. **E-step:** Compute responsibilities γ_k^n using current parameters $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$
 - b. **M-step:** Re-estimate $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$ using current (fixed) responsibilities γ_k^n
 - c. Compute data log-likelihood with current parameters and evaluate convergence

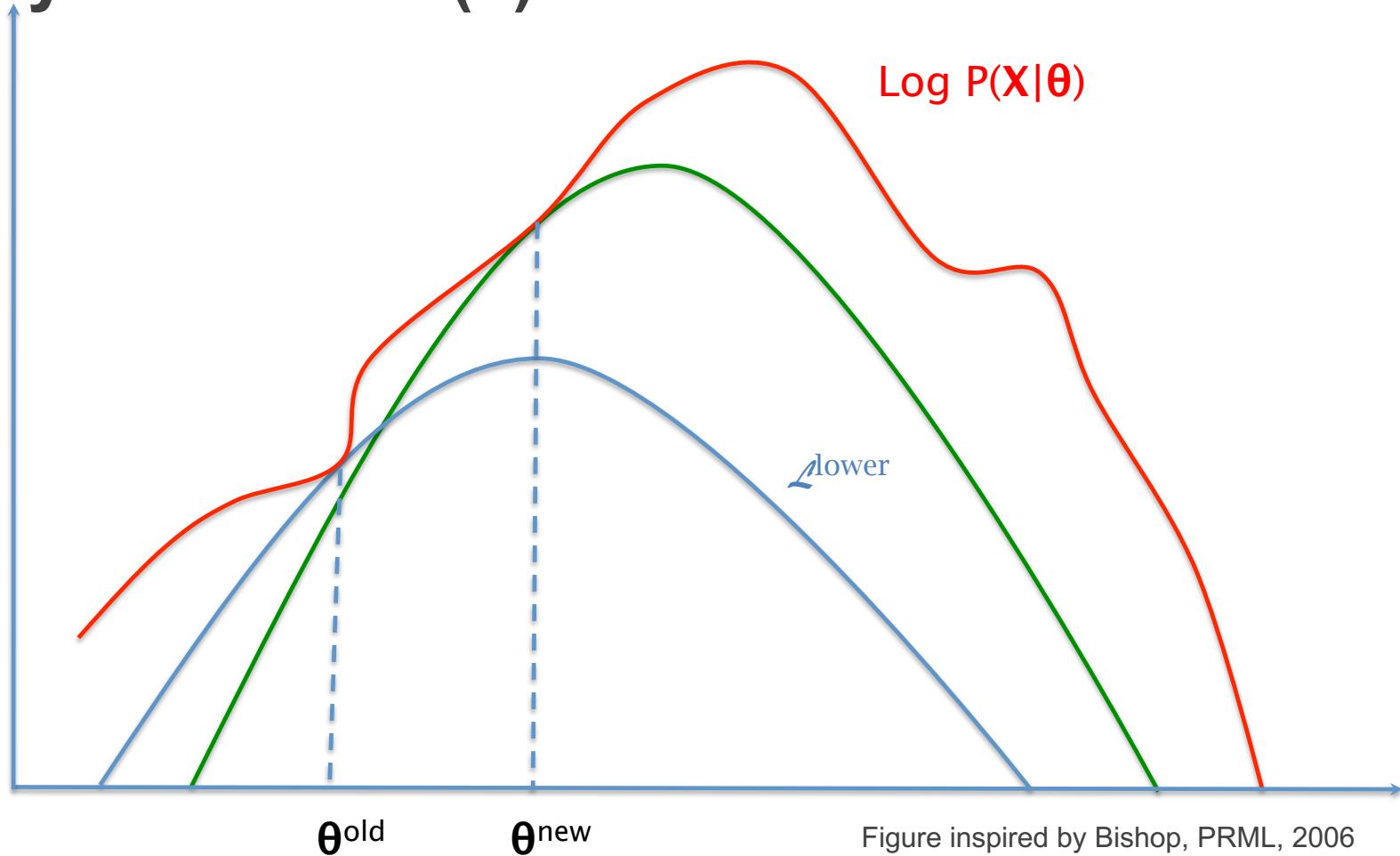
Learning GMMs: Demo

38



demo_gmm.m

Why EM Works (2) – Illustration



Iteratively, compute lower bound for fixed parameters θ^{old} and then maximize this lower bound obtaining θ^{new}

Relation Between GMMs and K-means

40

- GMM does soft clustering - we have posterior $p(z_k^n = 1 | \mathbf{x})$
- Recall documents example
- K-means does hard clustering
- Consider a GMM model with $\Sigma_k = \sigma^2 \mathbf{I}$
 - We obtain a **softmax** function for the responsibilities:

$$\gamma_k^n = \frac{\pi_k \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x}^n - \boldsymbol{\mu}_k\|^2\right)}{\sum_{j=1}^K \pi_j \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x}^n - \boldsymbol{\mu}_j\|^2\right)}$$

- The (hard) max is obtained by making $\sigma^2 \rightarrow 0$
 - Then $\gamma_k^n \rightarrow z_k^n$ (in the K-means algorithm, i.e. hard assignment)
 - Same update for the centres
 - Optimizing complete data log-likelihood same as optimizing error function in K-means

K-means is the hard limit of GMMs

GMMs and EM

Final Remarks

Singularities In GMMs: Consider the simple case $\Sigma_k = \sigma_k^2 I$

- Centre one of the components to a single datapoint $\mu_k = x^n$ and Make $\sigma_k \rightarrow 0$
- **The log-likelihood will tend to infinity!**
- Can use heuristics or be Bayesian (priors over variances)

Identifiability: $K!$ ways to assign K parameters to K components

- These solutions will have the same likelihood
- Problems with interpretability (not with density estimation)
- Possible initial confusion as components try to explain same data and symmetry eventually is broken

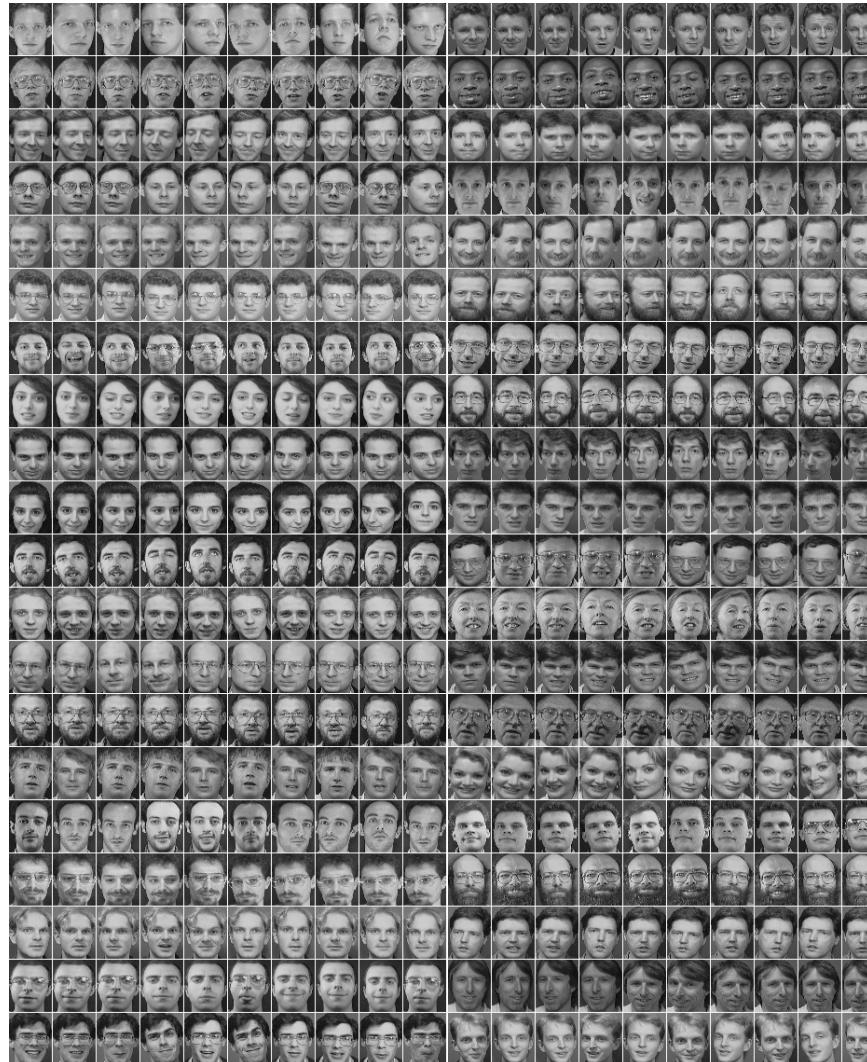
Local optima: EM can get stuck in local optima

- susceptible to different initializations (we can use K-means)

III. Probabilistic Classification: Generative Approaches – Naïve Bayes

Classification Problems

Facial Recognition



Examples of faces and their identity



What is her
identity?

Classification Problems

Handwritten Digit Recognition

7210414959
0690159734
9665407401
3134727121
1742351244

2

What is the number in the image?

Examples of images and their corresponding digit

Classification Problems

Supervised Document Classification

ODP – Open Directory Project

<http://www.dmoz.org/>

Google dmoz Search Share Bookmarks Translate AutoFill Sign in

d[m]o[z] open directory project In partnership with AOL Search.

[about dmoz](#) | [dmoz blog](#) | [suggest URL](#) | [help](#) | [link](#) | [editor login](#)

[Search](#) [advanced](#)

Arts Movies, Television, Music...	Business Jobs, Real Estate, Investing...	Computers Internet, Software, Hardware...
Games Video Games, RPGs, Gambling...	Health Fitness, Medicine, Alternative...	Home Family, Consumers, Cooking...
Kids and Teens Arts, School Time, Teen Life...	News Media, Newspapers, Weather...	Recreation Travel, Food, Outdoors, Humor...
Reference Maps, Education, Libraries...	Regional US, Canada, UK, Europe...	Science Biology, Psychology, Physics...
Shopping Clothing, Food, Gifts...	Society People, Religion, Issues...	Sports Baseball, Soccer, Basketball...
World Català, Dansk, Deutsch, Español, Français, Italiano, 日本語, Nederlands, Polski, Русский, Svenska...		

[Become an Editor](#) Help build the largest human-edited directory of the web

Copyright © 2011 Netscape

4,868,290 sites - 90,614 editors - over 1,006,417 categories



Difference with clustering search engines example?

The Classification Problem

Problem Definition (Refresher)

In all previous examples we are dealing with **discrete** targets.

Given a set of input-output pairs $\mathcal{D} = \{\mathbf{x}^n, y^n\}_{n=1}^N$

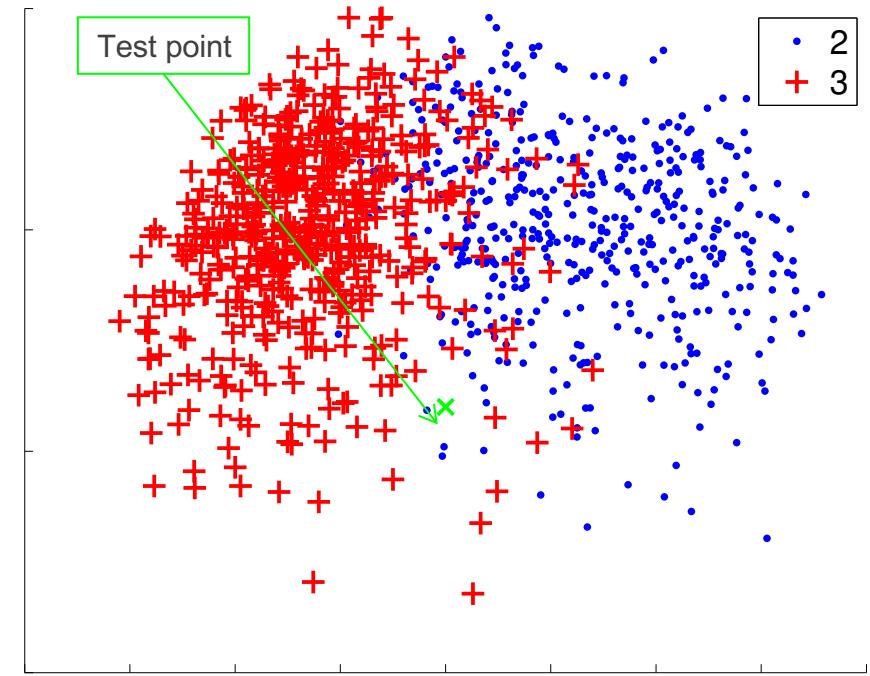
- \mathbf{x} : D-dimensional feature vector
- $y \in \{\mathcal{C}_1, \dots, \mathcal{C}_M\}$

Goal: Learn a mapping

- $f(\mathbf{x}) : \mathbf{x} \rightarrow y$
- Make predictions at unseen datapoints \mathbf{x}^* .

Probabilistic methods

- Confidence on our predictions
 - Decision theory
 - Reject option



- A probabilistic classifier provides an elegant framework for decision theory
- Suppose that we have a predictive probability $p(y = c | \mathbf{x})$
- Let $\underline{\mathcal{L}}(c, c')$ be the loss incurred by making a decision c' when the true class is c .

Then for a new \mathbf{x}^* we predict c^* that solves:

$$c^* = \operatorname{argmin}_{c'} \sum_c \underline{\mathcal{L}}(c, c') p(c | \mathbf{x}^*)$$

Expected loss

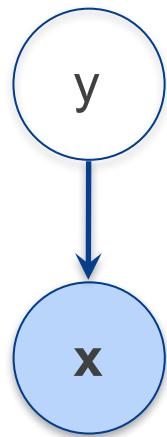
- For the case of the 0-1 loss (a unit penalty is paid for a misclassification), the optimal decision maximizes $p(c | \mathbf{x}^*)$
 - » This optimal classifier is known as **Bayes classifier**

How would we make predictions in probabilistic binary classification?

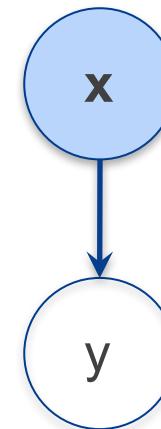
Generative vs. Discriminative Approaches

48

Generative



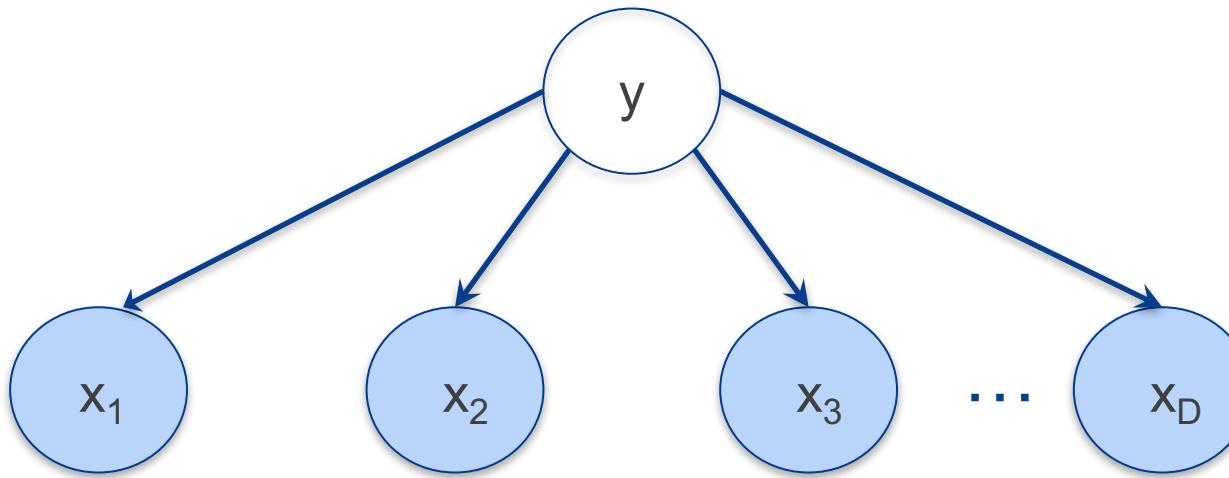
Discriminative



- Model the joint $p(\mathbf{x}, \mathbf{y})$ via models for $p(\mathbf{x}|\mathbf{y})$ and $p(\mathbf{y})$.
- Predictions $p(\mathbf{y}|\mathbf{x})$ via Bayes rule
 - Amenable to incorporation of prior knowledge
 - Indirect approach
- Model $p(\mathbf{y}|\mathbf{x})$ directly
 - Direct prediction
 - Difficult to incorporate prior information
 - Focus on the task at hand

Generative Approaches

Naïve Bayes



Generative model: $p(\mathbf{x}, y) = p(\mathbf{x} | y) p(y)$

- Where the features are **conditionally independent** given the class label:

$$p(\mathbf{x}|y) = \prod_{i=1}^D p(x_i|y)$$

- Inference problem: Predict the class for a new test point (use Baye's rule):

$$p(y|\mathbf{x}^*) = \frac{p(\mathbf{x}^*|y)p(y)}{\sum_{y'} p(\mathbf{x}^*|y')p(y')} \quad \text{← Marginal } p(\mathbf{x}^*)$$

Naïve Bayes

Binary Classification and Binary Features: Example (1)

- Consider the case of $y \in \{C_1, C_2\}$ and each $x_i \in \{0, 1\}$
 - We want to classify documents as being about sports (C_1) or politics (C_2).
- Bag-of-words
 - Describe each document as a binary vector indicating the presence/absence of a word in a vocabulary V .
 - $V = \{\text{football}, \text{defence}, \text{strategy}, \text{goal}, \text{office}\}$

- The conditional probability tables(CPTs) are given by:

$p(\mathcal{C}_1) = 0.5$ $p(f = 1 \mathcal{C}_1) = 0.8$ $p(d = 1 \mathcal{C}_1) = 0.7$ $p(s = 1 \mathcal{C}_1) = 0.2$ $p(g = 1 \mathcal{C}_1) = 0.7$ $p(o = 1 \mathcal{C}_1) = 0.2$	$p(f = 1 \mathcal{C}_2) = 0.1$ $p(d = 1 \mathcal{C}_2) = 0.7$ $p(s = 1 \mathcal{C}_2) = 0.8$ $p(g = 1 \mathcal{C}_2) = 0.3$ $p(o = 1 \mathcal{C}_2) = 0.7$
--	--

Naïve Bayes

Binary Classification and Binary Features: Example (2)

A new document arrives and is described by $\mathbf{x}^* = (0, 1, 1, 1, 0)$

- Is it about sports or politics? any guess?

$$p(C_1 | \mathbf{x}^*) = \frac{(0.2)(0.7)(0.2)(0.7)(0.8)}{(0.2)(0.7)(0.2)(0.7)(0.8) + (0.9)(0.7)(0.8)(0.3)(0.3)} \\ \approx 0.26$$

- What happened with the priors $p(C_1)$ and $p(C_2)$?
- If were to minimize the 0-1 loss we would classify this document as politics.
- How can we learn the model parameters (i.e. the CPTs) from data?
 - $\theta_i^y \stackrel{\text{def}}{=} p(x_i = 1 | y)$ and $\theta^y \stackrel{\text{def}}{=} p(y)$ for $y \in \{C_1, C_2\}$ and $i = 1, \dots, D$
- Intuitively, what would be reasonable estimates for these parameters?

Naïve Bayes

Parameter Estimation via Maximum Likelihood (1)

At training time we observe $\{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N$

- Assuming iid observations, the data log-likelihood is:

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_{n=1}^N \log p(\mathbf{x}^n, y^n) = \sum_{n=1}^N \log p(y^n) + \sum_{n=1}^N \log p(\mathbf{x}^n | y^n)$$

To optimise $\mathcal{L}(\boldsymbol{\theta})$ wrt $\theta^y \stackrel{\text{def}}{=} p(y)$ we need to maximise:

$$\mathcal{L}(\boldsymbol{\theta}^y) = \sum_{n=1}^N \log p(y^n) = N^{\mathcal{C}_1} \log \theta^{\mathcal{C}_1} + N^{\mathcal{C}_2} \log \theta^{\mathcal{C}_2} \quad \text{s.t.} \quad \sum_{y \in \{\mathcal{C}_1, \mathcal{C}_2\}} \theta^y = 1$$

- Where $N^{\mathcal{C}_j}$ is the number of times class \mathcal{C}_j occurs in the training data
- The solution to the above optimisation problem is: $\hat{\theta}^{\mathcal{C}_1} = \frac{N^{\mathcal{C}_1}}{N}$

Naïve Bayes

Parameter Estimation via Maximum Likelihood (2)

To optimize $\mathcal{L}(\theta)$ wrt $\theta_i^y \stackrel{\text{def}}{=} p(x_i = 1|y)$ we need to consider only:

$$\mathcal{L}(\theta_i^y) = \sum_{n=1}^N \log p(\mathbf{x}^n | y^n)$$

- Due to the independence assumption and binary nature of x_i :

$$p(\mathbf{x}|y) = \prod_{i=1}^D (\theta_i^y)^{x_i} (1 - \theta_i^y)^{1-x_i}$$

- Therefore, we want to maximise (wrt each θ_i^y):

$$\sum_n \sum_i \left(x_i^n \log \theta_i^{y^n} + (1 - x_i^n) \log (1 - \theta_i^{y^n}) \right)$$

- Which yields $\theta_i^{\mathcal{C}_j} = p(x_i = 1|y = \mathcal{C}_j) = \frac{\#(x_i = 1, \mathcal{C}_j)}{N_{\mathcal{C}_j}}$
- Where $\#(x_i = 1, \mathcal{C}_j)$ is the number of times $x_i=1$ for class \mathcal{C}_j

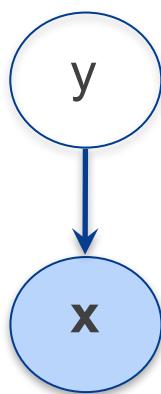
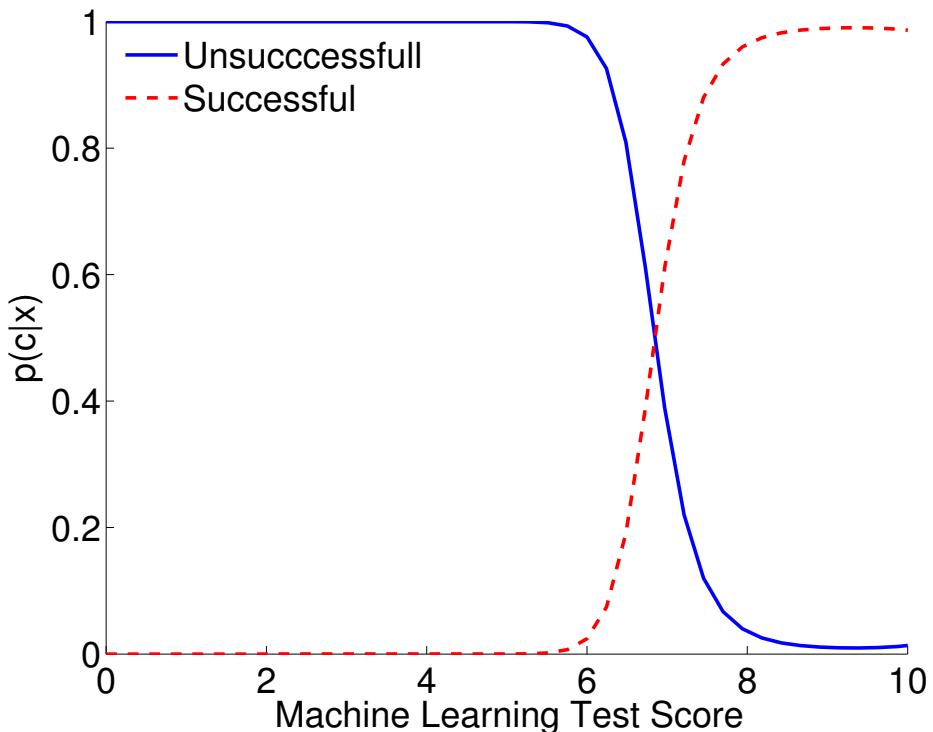
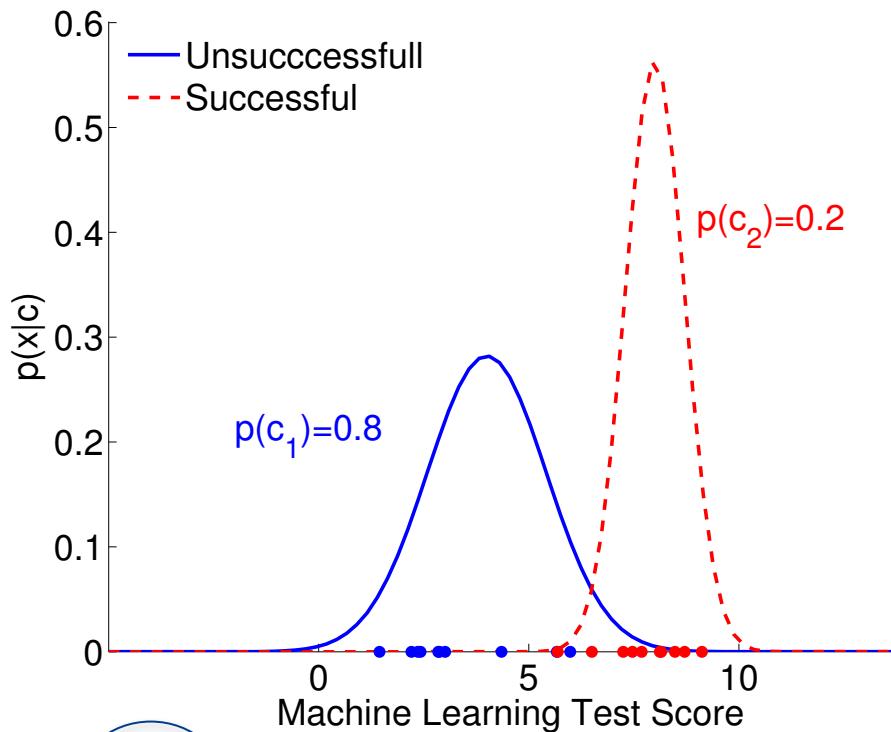
Naïve Bayes

Additional Remarks

- Naive Bayes is a simple, computationally efficient yet effective method in many practical applications
 - e.g. document classification
 - Conditional independence (not marginal independence)
 - [Linear decision boundary](#)
- Maximum likelihood estimates may be poor if data is very sparse
 - What is $p(\text{sport} | \mathbf{x}^*)$ if the word *goal* is never seen in training (and the new document does contain the word *goal*)?
 - Need to “smooth” estimates or be Bayesian
- What about features with more than two states?
 - Need to consider a categorical distribution
 - What parameter estimates would we get?
- What about continuous features?
 - Need a model for $p(x_i | C_j)$, Gaussian?

Generative Approaches

Class-conditional Gaussian Models



Generative: Modelling $p(x,y)$ via $p(x|y)$ and $p(y)$

Predictions via Bayes rule $p(y | x) \propto p(x|y) p(y)$

- Need to model the density $p(x|y)$

Class-conditional Gaussian Models

Inference and Learning

- Consider the binary case of $y \in \{C_1, C_2\}$:

- Learning:** Fitting a Gaussian to each class: How?

$$p(\mathbf{x}|\mathcal{C}_1) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) \quad p(\mathbf{x}|\mathcal{C}_2) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$$

- And estimating $p(C_1), p(C_2)$ How?

- Predictions:** For the 0-1 loss, we can compare the log of the (un-normalised posteriors):

$$\log p(\mathbf{x}^*|\mathcal{C}_1) + \log p(\mathcal{C}_1) > \log p(\mathbf{x}^*|\mathcal{C}_2) + \log p(\mathcal{C}_2)$$

- This gives us a quadratic decision boundary

Class-conditional Gaussian Models

Additional Remarks

- High-dimensional data Need to evaluate $\log|\Sigma_j|$ and Σ_j^{-1}
 - What can we do?
- How can we model more complex densities? i.e. to address multi-modality?
- Relation to Naïve Bayes: For continuous variables in Naïve Bayes
 - you can make: $p(x_i | C_j) = N(x_i | \mu_{ji}, \sigma_{ji}^2)$.
- This is equivalent to having a class-conditional Gaussian model with:
 - $p(\mathbf{x} | C_j) = N(\mathbf{x} | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$ where $\boldsymbol{\Sigma}_j = \text{diag}(\sigma_{j1}, \dots, \sigma_{jD})$

Summary & Conclusions

- Popular machine learning methods can be seen as probabilistic graphical models
- PCA vs. PPCA
- K-means vs GMMs
- Naïve Bayes
- Discriminative vs generative
- Inference → prediction
- Learning → parameter estimation
 - All data visible (maximum likelihood)
 - Hidden variables (EM)

Research Projects at all Levels

- I'm looking for **outstanding** students to work with me in research projects (PhD, Masters, Honours, fun)
 - Creative
 - Strong math and stats background (or willing to climb the hill)
 - Excellent programming skills
- Project on *structured prediction with deep learning*, selected for the prestigious [2018 UNSW Scientia Scholarship Scheme](#)
 - Stipend of \$40,000 per annum for 4 years and a support package of up to \$10,000 per annum
- Other topics:
 - Bayesian deep learning, probabilistic programming, reinforcement learning
 - Spatio-temporal modelling, causal inference, inverse problems
 - Privacy-preserving machine learning
- For a sample of my research look at <http://ebonilla.github.io/>
 - Drop me an email if interested (e.bonilla@unsw.edu.au)