# *Introduction to Statistical Machine Learning*

## Cheng Soon Ong & Christian Walder

Machine Learning Research Group
Data61 | CSIRO
and
College of Engineering and Computer Science
The Australian National University

Canberra
February – June 2019

(Many figures from C. M. Bishop, "Pattern Recognition and Machine Learning")

# Part VII

## *Mixture Models and EM 1*

# *Marginalisation*

Introduction to Statistical
Machine Learning

©2019
Ong & Walder & Webers
Data61 \ CSIRO
The Australian National
University

- Sum rule $p(A, B) = \sum_C p(A, B, C)$
- Product rule $p(A, B) = p(A|B)p(B)$
- Why do we optimize the log likelihood?

# *Strategy in this course*

- Estimate best predictor = training = learning
  Given data $(x_1, y_1), \ldots, (x_n, y_n)$, find a predictor $f_{\mathbf{w}}(\cdot)$.
    1. Identify the type of input $x$ and output $y$ data
    2. Propose a (linear) mathematical model for $f_{\mathbf{w}}$
    3. Design an objective function or likelihood
    4. Calculate the optimal parameter ($\mathbf{w}$)
    5. Model uncertainty using the Bayesian approach
    6. Implement and compute (the algorithm in python)
    7. Interpret and diagnose results

    We will study unsupervised learning this week

# *Mixture Models and EM*

Introduction to Statistical
Machine Learning

©2019
Ong & Walder & Webers
Data61 \ CSIRO
The Australian National
University

- Complex marginal distributions over observed variables can be expressed via more tractable joint distributions over the expanded space of observed and latent variables.
- Mixture Models can also be used to cluster data.
- General technique for finding maximum likelihood estimators in latent variable models: expectation-maximisation (EM) algorithm.
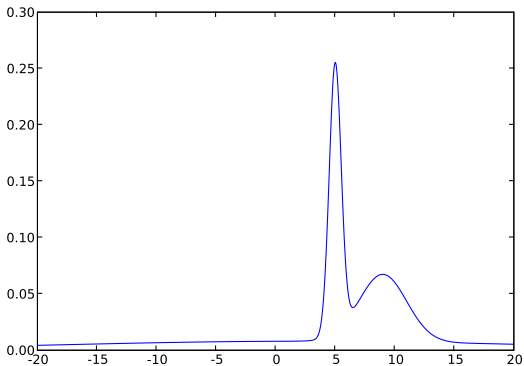
# *Example - Wallaby Distribution*

Introduction to Statistical
Machine Learning

ⓒ2019
Ong & Walder & Webers
Data61 \ CSIRO
The Australian National
University

- Introduced very recently to show . . .

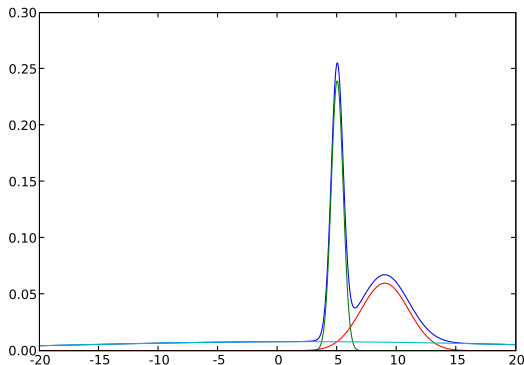# *Example - 'Wallaby' Distribution*

- . . . that already a mixture of three Gaussian can be fun.

$$p(x) = \frac{3}{10}\,\mathcal{N}(x\,|\,5, 0.5) + \frac{3}{10}\,\mathcal{N}(x\,|\,9, 2) + \frac{4}{10}\,\mathcal{N}(x\,|\,2, 20)$$

# *Example - 'Wallaby' Distribution*

Introduction to Statistical
Machine Learning

© 2019
Ong & Walder & Webers
Data61 | CSIRO
The Australian National
University

- Use $\mu, \sigma$ as latent variables and define a distribution

$$p(\mu, \sigma) = \begin{cases} \frac{3}{10} & \text{if } (\mu, \sigma) = (5, 0.5) \\ \frac{3}{10} & \text{if } (\mu, \sigma) = (9, 2) \\ \frac{4}{10} & \text{if } (\mu, \sigma) = (2, 20) \\ 0 & \text{otherwise.} \end{cases}$$

$$\begin{aligned} p(x) &= \int_{-\infty}^{\infty} \int_0^{\infty} p(x, \mu, \sigma) \; \mathrm{d}\mu \; \mathrm{d}\sigma \\ &= \int_{-\infty}^{\infty} \int_0^{\infty} p(x \,|\, \mu, \sigma) \, p(\mu, \sigma) \; \mathrm{d}\mu \; \mathrm{d}\sigma \\ &= \frac{3}{10} \, \mathcal{N}(x \,|\, 5, 0.5) + \frac{3}{10} \, \mathcal{N}(x \,|\, 9, 2) + \frac{4}{10} \, \mathcal{N}(x \,|\, 2, 20) \end{aligned}$$

# *K-means Clustering*

Introduction to Statistical
Machine Learning

©2019
Ong & Walder & Webers
Data61 \ CSIRO
The Australian National
University

- Given a set of data $\{\mathbf{x}_1, \ldots, \mathbf{x}_N\}$ where $\mathbf{x}_n \in \mathbb{R}^D$, $n = 1, \ldots, N$.
- Goal: Partition the data into $K$ clusters.

# K-means Clustering

Introduction to Statistical
Machine Learning

©2019
Ong & Walder & Webers
Data61 \ CSIRO
The Australian National
University

- Given a set of data $\{\mathbf{x}_1, \ldots, \mathbf{x}_N\}$ where $\mathbf{x}_n \in \mathbb{R}^D$, $n = 1, \ldots, N$.
- Goal: Partition the data into $K$ clusters.
- Each cluster contains points close to each other.
- Introduce a prototype $\boldsymbol{\mu}_k \in \mathbb{R}^D$ for each cluster.
- Goal: Find
  1. a set prototypes $\boldsymbol{\mu}_k$, $k = 1, \ldots, K$, each representing a different cluster.
  2. an assignment of each data point to exactly one cluster.

# *K-means Clustering - The Algorithm*

Introduction to Statistical
Machine Learning

© 2019
Ong & Walder & Webers
Data61 \ CSIRO
The Australian National
University

- Start with arbitrary chosen prototypes $\boldsymbol{\mu}_k$, $k = 1, \ldots, K$.
  1. Assign each data point to the closest prototype.
  2. Calculate new prototypes as the mean of all data points assigned to each of them.
- In the following, we will formalise this introducing a notation which will be useful later.

# *K-means Clustering - Notation*

- Binary indicator variables

$$r_{nk} = \begin{cases} 1, & \text{if data point } \mathbf{x}_n \text{ belongs to cluster } k \\ 0, & \text{otherwise} \end{cases}$$

  using the $1$-of-$K$ coding scheme.

- Define a distortion measure

$$J = \sum_{n=1}^{N} \sum_{k=1}^{K} r_{nk} \left\| \mathbf{x}_n - \boldsymbol{\mu}_k \right\|^2$$

- Find the values for $\{r_{nk}\}$ and $\{\boldsymbol{\mu}_k\}$ so as to minimise $J$.

# *K-means Clustering - Notation*

- Find the values for $\{r_{nk}\}$ and $\{\boldsymbol{\mu}_k\}$ so as to minimise $J$.
- But $\{r_{nk}\}$ depends on $\{\boldsymbol{\mu}_k\}$, and $\{\boldsymbol{\mu}_k\}$ depends on $\{r_{nk}\}$.

# *K-means Clustering - Notation*

- Find the values for $\{r_{nk}\}$ and $\{\boldsymbol{\mu}_k\}$ so as to minimise $J$.
- But $\{r_{nk}\}$ depends on $\{\boldsymbol{\mu}_k\}$, and $\{\boldsymbol{\mu}_k\}$ depends on $\{r_{nk}\}$.
- Iterate until no further change
  1. Minimise $J$ w.r.t. $r_{nk}$ while keeping $\{\boldsymbol{\mu}_k\}$ fixed,

  $$r_{nk} = \begin{cases} 1, & \text{if } k = \arg\min_j \|\mathbf{x}_n - \boldsymbol{\mu}_j\|^2 \\ 0, & \text{otherwise.} \end{cases} \qquad \forall n = 1, \ldots, N$$

  Expectation step
  2. Minimise $J$ w.r.t. $\{\boldsymbol{\mu}_k\}$ while keeping $r_{nk}$ fixed,

  $$0 = 2 \sum_{n=1}^{N} r_{nk}(\mathbf{x}_n - \boldsymbol{\mu}_k)$$

  $$\boldsymbol{\mu}_k = \frac{\sum_{n=1}^{N} r_{nk}\mathbf{x}_n}{\sum_{n=1}^{N} r_{nk}}$$

  Maximisation step

# K-means Clustering - Example

Introduction to Statistical
Machine Learning

ⓒ2019
Ong & Walder & Webers
Data61 \ CSIRO
The Australian National
University

# K-means Clustering - Example

Introduction to Statistical
Machine Learning

ⓒ2019
Ong & Walder & Webers
Data61 \ CSIRO
The Australian National
University

# *K-means Clustering - Cost Function*

Introduction to Statistical
Machine Learning

ⓒ2019
Ong & Walder & Webers
Data61 \ CSIRO
The Australian National
University

Cost function $J$ after each E step (blue points)
and M step (red points).

# K-means Clustering - Notes

- Initial condition crucial for convergence.
- What happens, if at least one cluster centre is too far from all data points?
- Complex step: Finding the nearest neighbour. (Use triangle inequality; build K-D trees, ...)
- Generalise to non-Euclidean dissimilarity measures $\mathcal{V}(\mathbf{x}_n, \boldsymbol{\mu}_k)$ (called $K$-medoids algorithm),

$$\widetilde{J} = \sum_{n=1}^{N} \sum_{k=1}^{K} r_{nk} \mathcal{V}(\mathbf{x}_n, \boldsymbol{\mu}_k).$$

- Online stochastic algorithm
  1. Draw data point $\mathbf{x}_n$ and locate nearest prototype $\boldsymbol{\mu}_k$.
  2. Update only $\boldsymbol{\mu}_k$ using decreasing learning rate $\eta_n$

$$\boldsymbol{\mu}_k^{\mathsf{new}} = \boldsymbol{\mu}_k^{\mathsf{old}} + \eta_n(\mathbf{x}_n - \boldsymbol{\mu}_k^{\mathsf{old}}).$$

# *K-means Clustering - Image Segmentation*

*Introduction to Statistical Machine Learning*

©2019
*Ong & Walder & Webers
Data61 \ CSIRO
The Australian National
University*

- Segment an image into regions of reasonable homogeneous appearance.
- Each pixel is a point in $\mathbb{R}^3$ (red, blue, green). (Note that the pixel intensities are bounded in the range $[0, 1]$ and therefore this space is strictly speaking not Euclidean).
- Run $K$-means on all points of the image until convergence. Replace all pixels with the corresponding mean $\boldsymbol{\mu}_k$.
- Results in an image with a palette only $K$ different colours.
- There are much better approaches to image segmentation (but it is an active research topic), this here serves only to illustrate $K$-means.

# *Illustrating K-means Clustering - Segmentation*

Introduction to Statistical
Machine Learning

ⓒ2019
Ong & Walder & Webers
Data61 \ CSIRO
The Australian National
University

$K = 2$

$K = 10$

$K = 3$

Original image

# Illustrating K-means Clustering - Segmentation

# K-means Clustering - Compression

- Lossy data compression: accept some errors in the reconstruction as trade-off for higher compression.
- Apply $K$-means to the data.
- Store the code-book vectors $\boldsymbol{\mu}_k$.
- Store the data in the form of references (labels) to the code-book. Each data point has a label in the range $[1, \ldots, K]$.
- New data points are also compressed by finding the closest code-book vector and then storing only the label.
- This technique is also called vector quantisation.

# *Illustrating K-means Clustering - Compression*

$K = 2$

4.2%

$K = 3$

8.3%

$K = 10$

16.7%

Original image

100 %

# *Latent variable modeling*

- We have already seen a mixture of two Gaussians for linear classification
- However in the clustering scenario, we do not observe the class membership
- Strategy (this is vague)
  - We have a difficult distribution $p(\mathbf{x})$
  - We introduce a new variable $\mathbf{z}$ to get $p(\mathbf{x}, \mathbf{z})$
  - Model $p(\mathbf{x}, \mathbf{z}) = p(\mathbf{x}|\mathbf{z})p(\mathbf{z})$ with easy distributions

# *Mixture of Gaussians*

- A Gaussian mixture distribution is a linear superposition of Gaussians of the form

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

- As $\int p(\mathbf{x}) \, d\mathbf{x} = 1$, if follows $\sum_{k=1}^{K} \pi_k = 1$.
- Let us write this with the help of a latent variable $\mathbf{z}$.

### Definition (Latent variables)

Latent variables (as opposed to observable variables), are variables that are not directly observed but are rather inferred (through a mathematical model) from other variables that are observed and directly measured. They are also sometimes called hidden variables, model parameters, or hypothetical variables.

# *Mixture of Gaussians*

- Let $\mathbf{z} \in \{0, 1\}^K$ and $\sum_{k=1}^{K} z_k = 1$. In words, $\mathbf{z}$ is a $K$-dimensional vector in 1-of-$K$ representation.
- There are exactly $K$ different possible vectors $\mathbf{z}$ depending on which of the $K$ entries is 1.
- Define the joint distribution $p(\mathbf{x}, \mathbf{z})$ in terms of a marginal distribution $p(\mathbf{z})$ and a conditional distribution $p(\mathbf{x} \,|\, \mathbf{z})$ as

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

# *Mixture of Gaussians*

- Set the marginal distribution to

$$p(z_k = 1) = \pi_k$$

  where $0 \leq \pi_k \leq 1$ together with $\sum_{k=1}^{K} \pi_k = 1$.

- Because $\mathbf{z}$ uses $1$-of-$K$ coding, we can also write

$$p(\mathbf{z}) = \prod_{k=1}^{K} \pi_k^{z_k}.$$

- Set the conditional distribution of $\mathbf{x}$ given a particular $\mathbf{z}$ to

$$p(\mathbf{x} \,|\, z_k = 1) = \mathcal{N}(\mathbf{x} \,|\, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k),$$

  or

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

# *Mixture of Gaussians*

- The marginal distribution over $\mathbf{x}$ is now found by summing the joint distribution over all possible states of $\mathbf{z}$

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

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

- The marginal distribution of $\mathbf{x}$ is a Gaussian mixture.
- For several observations $\mathbf{x}_1, \ldots, \mathbf{x}_N$ we need one latent variable $\mathbf{z}_n$ per observation.
- What have we gained? Can now work with the joint distribution $p(\mathbf{x}, \mathbf{z})$. Will lead to significant simplification later, especially for EM algorithm.

# *Mixture of Gaussians*

Introduction to Statistical
Machine Learning

©2019
Ong & Walder & Webers
Data61 | CSIRO
The Australian National
University

- Conditional probability of $\mathbf{z}$ given $\mathbf{x}$ by Bayes' theorem

$$\gamma(z_k) = p(z_k = 1 \,|\, \mathbf{x}) = \frac{p(z_k = 1) \, p(\mathbf{x} \,|\, z_k = 1)}{\sum_{j=1}^{K} p(z_j = 1) \, p(\mathbf{x} \,|\, z_j = 1)}$$
$$= \frac{\pi_k \, \mathcal{N}(\mathbf{x} \,|\, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^{K} \pi_j \, \mathcal{N}(\mathbf{x} \,|\, \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

- $\gamma(z_k)$ is the responsibility of component $k$ to 'explain' the observation $\mathbf{x}$.

# *Mixture of Gaussians - Ancestral Sampling*

Introduction to Statistical
Machine Learning

©2019
Ong & Walder & Webers
Data61 \ CSIRO
The Australian National
University

- Goal: Generate random samples distributed according to the mixture model.
    1. Generate a sample $\hat{\mathbf{z}}$ from the distribution $p(\mathbf{z})$.
    2. Generate a value $\hat{\mathbf{x}}$ from the conditional distribution $p(\mathbf{x} \mid \hat{\mathbf{z}})$.
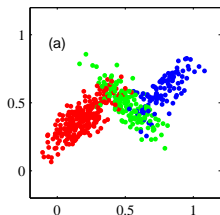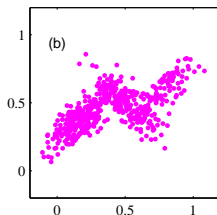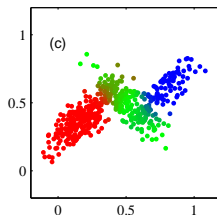- Example: Mixture of 3 Gaussians, 500 points.



Original states of $\mathbf{z}$.  Marginal $p(\mathbf{x})$.  (R, G, B) - colours mixed according to $\gamma(z_{nk})$.

# *Mixture of Gaussians - Maximum Likelihood*

Introduction to Statistical
Machine Learning

©2019
Ong & Walder & Webers
Data61 \ CSIRO
The Australian National
University

- Given $N$ data points, each of dimension $D$, we have the data matrix $\mathbf{X} \in \mathbb{R}^{N \times D}$ where each row contains one data point.
- Similarly, we have the matrix of latent variables $\mathbf{Z} \in \mathbb{R}^{N \times K}$ with rows $\mathbf{z}_n^T$.
- Assume the data are drawn i.i.d., the distribution for the data can be represented by a graphical model.

# *Mixture of Gaussians - Maximum Likelihood*

- The log of the likelihood function is then

$$\ln p(\mathbf{X} \mid \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^{N} \ln \left\{ \sum_{k=1}^{K} \pi_k \, \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

- Significant problem: If a mean $\boldsymbol{\mu}_j$ 'sits' directly on a data point $\mathbf{x}_n$ then

$$\mathcal{N}(\mathbf{x}_n \mid \mathbf{x}_n, \sigma_j^2 \mathbf{I}) = \frac{1}{(2\pi)^{1/2}} \frac{1}{\sigma_j}.$$

- Here we assumed $\boldsymbol{\Sigma}_k = \sigma_k^2 \mathbf{I}$. But problem is general, just think of a main axis transformation for $\boldsymbol{\Sigma}_k$.
- Overfitting (in disguise) occuring again with the maximum likelihood approach.
- Use heuristics to detect this situation and reset the mean of the corresponding component of the mixture.

# *Mixture of Gaussians - Maximum Likelihood*

- A $K$ component mixture has a total of $K!$ equivalent solutions corresponding to the $K!$ ways of assigning $K$ sets of parameters to $K$ solutions.
- Also called identifiability problem. Needs to be considered when the parameters discovered by a model are interpreted.
- Maximising the log likelihood of a Gaussian mixture is more complex then for a single Gaussian. Summation over all $K$ components inside of the logarithm make it harder.
- Setting the derivatives of the log likelihood to zero does not longer result in a closed form.
- May use gradient-based optimisation.
- Or EM algorithm. Stay tuned.