

RWTH AACHEN
UNIVERSITY

Machine Learning – Lecture 4

Probability Density Estimation III

22.10.2017

Bastian Leibe
RWTH Aachen
<http://www.vision.rwth-aachen.de>
leibe@vision.rwth-aachen.de

Machine Learning Winter '18

RWTH AACHEN
UNIVERSITY

Announcements

- Exam dates
 - According to rwth online, the exam dates are
 - 1st try Sat 02.03.2019 10:30 – 12:00h
 - 2nd try Thu 21.03.2019 13:30 – 15:30h
- Exam registration will start in early December...

B. Leibe

RWTH AACHEN
UNIVERSITY

Course Outline

- Fundamentals
 - Bayes Decision Theory
 - Probability Density Estimation
- Classification Approaches
 - Linear Discriminants
 - Support Vector Machines
 - Ensemble Methods & Boosting
 - Randomized Trees, Forests & Ferns
- Deep Learning
 - Foundations
 - Convolutional Neural Networks
 - Recurrent Neural Networks

B. Leibe

RWTH AACHEN
UNIVERSITY

Recap: Maximum Likelihood Approach

- Computation of the likelihood
 - Single data point: $p(x_n|\theta)$
 - Assumption: all data points $X = \{x_1, \dots, x_n\}$ independent
$$L(\theta) = p(X|\theta) = \prod_{n=1}^N p(x_n|\theta)$$
 - Log-likelihood
$$E(\theta) = -\ln L(\theta) = -\sum_{n=1}^N \ln p(x_n|\theta)$$
- Estimation of the parameters θ (Learning)
 - Maximize the likelihood (=minimize the negative log-likelihood)
 - Take the derivative and set it to zero.
$$\frac{\partial}{\partial \theta} E(\theta) = -\sum_{n=1}^N \frac{\partial \ln p(x_n|\theta)}{\partial \theta} \stackrel{!}{=} 0$$

B. Leibe

RWTH AACHEN
UNIVERSITY

Recap: Histograms

- Basic idea:
 - 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, $\Delta_i = \Delta$.
 - This can be done, in principle, for any dimensionality D ...

B. Leibe

Image source: G.M. Bishop, 2006

RWTH AACHEN
UNIVERSITY

Recap: Kernel Density Estimation

- Approximation formula:

$$p(\mathbf{x}) \approx \frac{K}{NV}$$
 - fixed V determine K → Kernel Methods
 - fixed K determine V → K-Nearest Neighbor
- Kernel methods
 - Place a *kernel window* k at location \mathbf{x} and count how many data points fall inside it.
- K-Nearest Neighbor
 - Increase the volume V until the K nearest data points are found.

B. Leibe

Slide adapted from Bernd Schölkopf

RWTH AACHEN UNIVERSITY

Topics of This Lecture

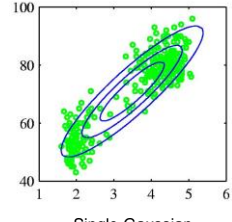
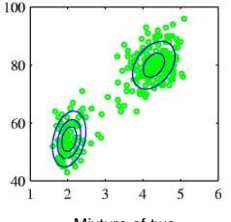
- **Mixture distributions**
 - Mixture of Gaussians (MoG)
 - Maximum Likelihood estimation attempt
- **K-Means Clustering**
 - Algorithm
 - Applications
- **EM Algorithm**
 - Credit assignment problem
 - MoG estimation
 - EM Algorithm
 - Interpretation of K-Means
 - Technical advice
- **Applications**

7

RWTH AACHEN UNIVERSITY

Mixture Distributions

- A single parametric distribution is often not sufficient
 - E.g. for multimodal data

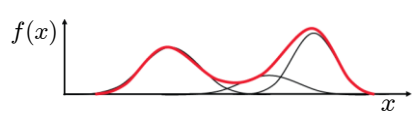



10

RWTH AACHEN UNIVERSITY

Mixture of Gaussians (MoG)

- Sum of M individual Normal distributions



- In the limit, every smooth distribution can be approximated this way (if M is large enough)

$$p(x|\theta) = \sum_{j=1}^M p(x|\theta_j)p(j)$$

11

RWTH AACHEN UNIVERSITY

Mixture of Gaussians

$$p(x|\theta) = \sum_{j=1}^M p(x|\theta_j)p(j)$$

Likelihood of measurement x given mixture component j

$$p(x|\theta_j) = \mathcal{N}(x|\mu_j, \sigma_j^2) = \frac{1}{\sqrt{2\pi}\sigma_j} \exp\left\{-\frac{(x - \mu_j)^2}{2\sigma_j^2}\right\}$$

$$p(j) = \pi_j \text{ with } 0 \leq \pi_j \leq 1 \text{ and } \sum_{j=1}^M \pi_j = 1$$

Prior of component j

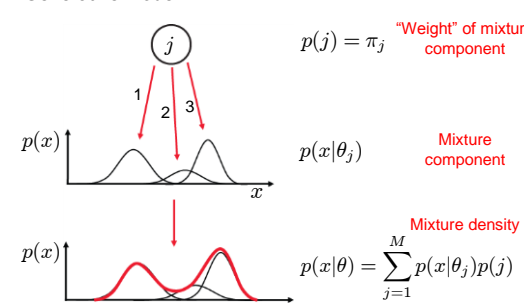
- **Notes**
 - The mixture density integrates to 1: $\int p(x)dx = 1$
 - The mixture parameters are $\theta = (\pi_1, \mu_1, \sigma_1, \dots, \pi_M, \mu_M, \sigma_M)$

12

RWTH AACHEN UNIVERSITY

Mixture of Gaussians (MoG)

- "Generative model"



"Weight" of mixture component

Mixture component

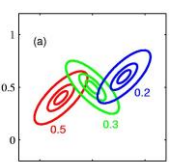
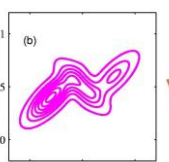
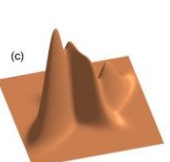
Mixture density

$$p(x|\theta) = \sum_{j=1}^M p(x|\theta_j)p(j)$$

13

RWTH AACHEN UNIVERSITY

Mixture of Multivariate Gaussians

14

Machine Learning Winter '18

Mixture of Multivariate Gaussians

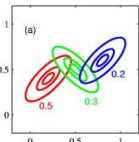
- Multivariate Gaussians

$$p(\mathbf{x}|\theta) = \sum_{j=1}^M p(\mathbf{x}|\theta_j)p(j)$$

$$p(\mathbf{x}|\theta_j) = \frac{1}{(2\pi)^{D/2}|\Sigma_j|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \mu_j)^T \Sigma_j^{-1}(\mathbf{x} - \mu_j)\right\}$$
- Mixture weights / mixture coefficients:

$$p(j) = \pi_j \text{ with } 0 \leq \pi_j \leq 1 \text{ and } \sum_{j=1}^M \pi_j = 1$$
- Parameters:

$$\theta = (\pi_1, \mu_1, \Sigma_1, \dots, \pi_M, \mu_M, \Sigma_M)$$



Slide credit: Bernt Schiele B. Leibe Image source: C.M. Bishop, 2006 15

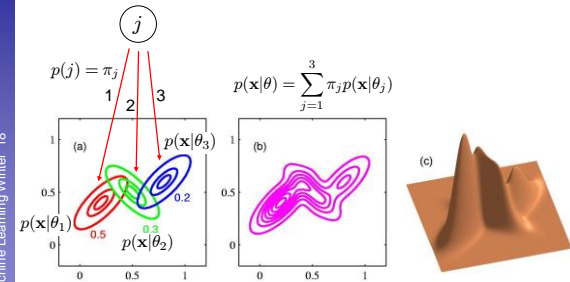
Machine Learning Winter '18

Mixture of Multivariate Gaussians

- "Generative model"

$$p(j) = \pi_j$$

$$p(\mathbf{x}|\theta) = \sum_{j=1}^3 \pi_j p(\mathbf{x}|\theta_j)$$

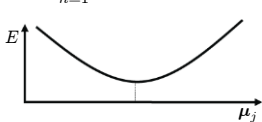


Slide credit: Bernt Schiele B. Leibe Image source: C.M. Bishop, 2006 16

Machine Learning Winter '18

Mixture of Gaussians – 1st Estimation Attempt

- Maximum Likelihood
 - Minimize $E = -\ln L(\theta) = -\sum_{n=1}^N \ln p(\mathbf{x}_n|\theta)$
 - Let's first look at μ_j :

$$\frac{\partial E}{\partial \mu_j} = 0$$

 - We can already see that this will be difficult, since

$$\ln p(\mathbf{X}|\pi, \mu, \Sigma) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n|\mu_k, \Sigma_k) \right\}$$

This will cause problems!

Slide adapted from Bernt Schiele B. Leibe 17

Machine Learning Winter '18

Mixture of Gaussians – 1st Estimation Attempt

- Minimization:

$$\frac{\partial E}{\partial \mu_j} = -\sum_{n=1}^N \frac{\frac{\partial}{\partial \mu_j} p(\mathbf{x}_n|\theta_j)}{\sum_{k=1}^K p(\mathbf{x}_n|\theta_k)}$$

$$= -\sum_{n=1}^N \left(\Sigma^{-1}(\mathbf{x}_n - \mu_j) \frac{p(\mathbf{x}_n|\theta_j)}{\sum_{k=1}^K p(\mathbf{x}_n|\theta_k)} \right)$$

$$= -\sum_{n=1}^N (\mathbf{x}_n - \mu_j) \frac{\pi_j \mathcal{N}(\mathbf{x}_n|\mu_j, \Sigma_j)}{\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n|\mu_k, \Sigma_k)} \stackrel{!}{=} 0$$

= $\gamma_j(\mathbf{x}_n)$
"responsibility" of component j for \mathbf{x}_n
- We thus obtain

$$\Rightarrow \mu_j = \frac{\sum_{n=1}^N \gamma_j(\mathbf{x}_n) \mathbf{x}_n}{\sum_{n=1}^N \gamma_j(\mathbf{x}_n)}$$

Slide credit: Bernt Schiele B. Leibe 18

Machine Learning Winter '18

Mixture of Gaussians – 1st Estimation Attempt

- But...

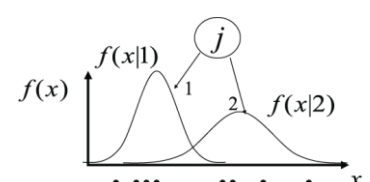
$$\mu_j = \frac{\sum_{n=1}^N \gamma_j(\mathbf{x}_n) \mathbf{x}_n}{\sum_{n=1}^N \gamma_j(\mathbf{x}_n)} \quad \gamma_j(\mathbf{x}_n) = \frac{\pi_j \mathcal{N}(\mathbf{x}_n|\mu_j, \Sigma_j)}{\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n|\mu_k, \Sigma_k)}$$
- I.e. there is no direct analytical solution!

$$\frac{\partial E}{\partial \mu_j} = f(\pi_1, \mu_1, \Sigma_1, \dots, \pi_M, \mu_M, \Sigma_M)$$
 - Complex gradient function (non-linear mutual dependencies)
 - Optimization of one Gaussian depends on all other Gaussians!
 - It is possible to apply iterative numerical optimization here, but in the following, we will see a simpler method.

B. Leibe 19

Machine Learning Winter '18

Mixture of Gaussians – Other Strategy

- Other strategy:
 
 - Observed data: $\bullet \dots \bullet$
 - Unobserved data: $1 \ 111 \quad 22 \ 2 \ 2$
 - Unobserved = "hidden variable": $j|x$
$$h(j=1|x_n) = \begin{matrix} 1 & 111 & 00 & 0 & 0 \end{matrix}$$

$$h(j=2|x_n) = \begin{matrix} 0 & 000 & 11 & 1 & 1 \end{matrix}$$

Slide credit: Bernt Schiele B. Leibe 20

RWTH AACHEN UNIVERSITY

Mixture of Gaussians – Other Strategy

- Assuming we knew the values of the hidden variable...

ML for Gaussian #1 \uparrow \uparrow ML for Gaussian #2

assumed known \rightarrow 1 111 22 2 2 j

$h(j=1|x_n) =$ 1 111 00 0 0

$h(j=2|x_n) =$ 0 000 11 1 1

$$\mu_1 = \frac{\sum_{n=1}^N h(j=1|x_n)x_n}{\sum_{i=1}^N h(j=1|x_n)} \quad \mu_2 = \frac{\sum_{n=1}^N h(j=2|x_n)x_n}{\sum_{i=1}^N h(j=2|x_n)}$$

Machine Learning Winter '18 21

Slide credit: Bernt Schiele B. Leibe

RWTH AACHEN UNIVERSITY

Mixture of Gaussians – Other Strategy

- Assuming we knew the mixture components...

$p(j=1|x)$ \downarrow \downarrow $p(j=2|x)$

1 111 22 2 2 j

- Bayes decision rule: Decide $j=1$ if $p(j=1|x_n) > p(j=2|x_n)$

Machine Learning Winter '18 22

Slide credit: Bernt Schiele B. Leibe

RWTH AACHEN UNIVERSITY

Clustering with Hard Assignments

- Let's first look at clustering with "hard assignments"

1 111 22 2 2 j

Machine Learning Winter '18 23

Slide credit: Bernt Schiele B. Leibe

RWTH AACHEN UNIVERSITY

Topics of This Lecture

- Mixture distributions
 - Mixture of Gaussians (MoG)
 - Maximum Likelihood estimation attempt
- K-Means Clustering
 - Algorithm
 - Applications
- EM Algorithm
 - Credit assignment problem
 - MoG estimation
 - EM Algorithm
 - Interpretation of K-Means
 - Technical advice
- Applications

Machine Learning Winter '18 24

Slide credit: Bernt Schiele B. Leibe

RWTH AACHEN UNIVERSITY

K-Means Clustering

- Iterative procedure
 - Initialization: pick K arbitrary centroids (cluster means)
 - Assign each sample to the closest centroid.
 - Adjust the centroids to be the means of the samples assigned to them.
 - Go to step 2 (until no change)
- Algorithm is guaranteed to converge after finite #iterations.
 - Local optimum
 - Final result depends on initialization.

Machine Learning Winter '18 25

Slide credit: Bernt Schiele B. Leibe

RWTH AACHEN UNIVERSITY

K-Means – Example with K=2

Machine Learning Winter '18 26

Slide credit: Bernt Schiele B. Leibe Image source: C.M. Bishop, 2006

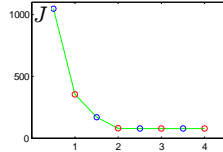
RWTH AACHEN UNIVERSITY

K-Means Clustering

- K-Means optimizes the following objective function:

$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2$$
- where

$$r_{nk} = \begin{cases} 1 & \text{if } k = \arg \min_j \|x_n - \mu_j\|^2 \\ 0 & \text{otherwise.} \end{cases}$$
- i.e., r_{nk} is an indicator variable that checks whether μ_k is the nearest cluster center to point x_n .
- In practice, this procedure usually converges quickly to a local optimum.

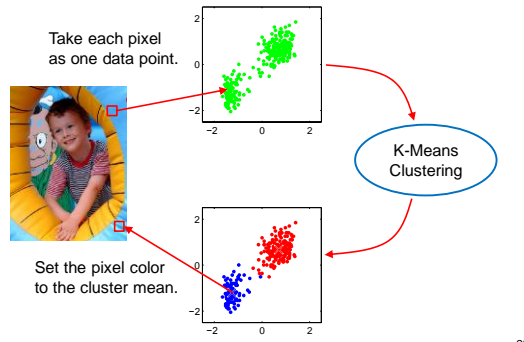


27
B. Leibe Image source: C.M. Bishop, 2006

RWTH AACHEN UNIVERSITY

Example Application: Image Compression

Take each pixel as one data point.




Set the pixel color to the cluster mean.

28
B. Leibe Image source: C.M. Bishop, 2006


RWTH AACHEN UNIVERSITY

Example Application: Image Compression


$K = 2$




$K = 3$




$K = 10$



Original image



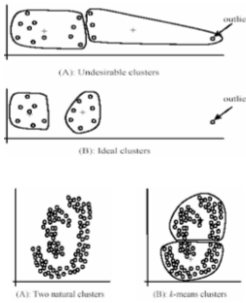


29
B. Leibe Image source: C.M. Bishop, 2006

RWTH AACHEN UNIVERSITY

Summary K-Means

- Pros**
 - Simple, fast to compute
 - Converges to local minimum of within-cluster squared error
- Problem cases**
 - Setting k ?
 - Sensitive to initial centers
 - Sensitive to outliers
 - Detects spherical clusters only
- Extensions**
 - Speed-ups possible through efficient search structures
 - General distance measures: k-medoids



30
Slide credit: Kristen Grauman B. Leibe

RWTH AACHEN UNIVERSITY

Topics of This Lecture

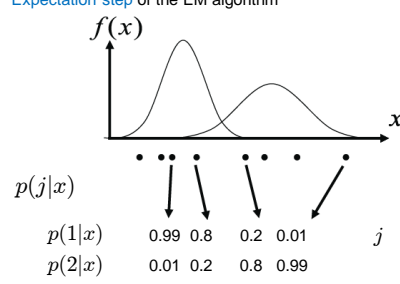
- Mixture distributions
 - Recap: Mixture of Gaussians (MoG)
 - Maximum Likelihood estimation attempt
- K-Means Clustering
 - Algorithm
 - Applications
- EM Algorithm**
 - Credit assignment problem
 - MoG estimation
 - EM Algorithm
 - Interpretation of K-Means
 - Technical advice
- Applications

31
B. Leibe

RWTH AACHEN UNIVERSITY

EM Clustering

- Clustering with "soft assignments"
 - Expectation step of the EM algorithm



$p(1 x)$	0.99	0.8	0.2	0.01	j
$p(2 x)$	0.01	0.2	0.8	0.99	

32
Slide credit: Bernt Schiele B. Leibe

RWTH AACHEN UNIVERSITY

EM Clustering

- Clustering with "soft assignments"
 - Maximization step of the EM algorithm

$p(1 x)$	0.99	0.8	0.2	0.01
$p(2 x)$	0.01	0.2	0.8	0.99

$$\mu_j = \frac{\sum_{n=1}^N p(j|\mathbf{x}_n) \mathbf{x}_n}{\sum_{n=1}^N p(j|\mathbf{x}_n)}$$

Maximum Likelihood estimate

Machine Learning Winter '18 | Slide credit: Bernt Schiele | B. Leibe | 33

RWTH AACHEN UNIVERSITY

Credit Assignment Problem

- "Credit Assignment Problem"
 - If we are just given \mathbf{x} , we don't know which mixture component this example came from

$$p(\mathbf{x}|\theta) = \sum_{j=1}^2 \pi_j p(\mathbf{x}|\theta_j)$$
 - We can however evaluate the posterior probability that an observed \mathbf{x} was generated from the first mixture component.

$$p(j=1|\mathbf{x}, \theta) = \frac{p(j=1, \mathbf{x}|\theta)}{p(\mathbf{x}|\theta)}$$

$$p(j=1, \mathbf{x}|\theta) = p(\mathbf{x}|j=1, \theta)p(j=1) = p(\mathbf{x}|\theta_1)p(j=1)$$

$$p(j=1|\mathbf{x}, \theta) = \frac{p(\mathbf{x}|\theta_1)p(j=1)}{\sum_{j=1}^2 p(\mathbf{x}|\theta_j)p(j)} = \gamma_j(\mathbf{x})$$

"responsibility" of component j for \mathbf{x} .

Machine Learning Winter '18 | Slide credit: Bernt Schiele | B. Leibe | 34

RWTH AACHEN UNIVERSITY

EM Algorithm

- Expectation-Maximization (EM) Algorithm
 - E-Step: softly assign samples to mixture components

$$\gamma_j(\mathbf{x}_n) \leftarrow \frac{\pi_j \mathcal{N}(\mathbf{x}_n|\mu_j, \Sigma_j)}{\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n|\mu_k, \Sigma_k)} \quad \forall j=1, \dots, K, \quad n=1, \dots, N$$
 - M-Step: re-estimate the parameters (separately for each mixture component) based on the soft assignments

$$\tilde{N}_j \leftarrow \sum_{n=1}^N \gamma_j(\mathbf{x}_n) = \text{soft number of samples labeled } j$$

$$\hat{\pi}_j^{\text{new}} \leftarrow \frac{\tilde{N}_j}{N}$$

$$\hat{\mu}_j^{\text{new}} \leftarrow \frac{1}{\tilde{N}_j} \sum_{n=1}^N \gamma_j(\mathbf{x}_n) \mathbf{x}_n$$

$$\hat{\Sigma}_j^{\text{new}} \leftarrow \frac{1}{\tilde{N}_j} \sum_{n=1}^N \gamma_j(\mathbf{x}_n) (\mathbf{x}_n - \hat{\mu}_j^{\text{new}})(\mathbf{x}_n - \hat{\mu}_j^{\text{new}})^T$$

Machine Learning Winter '18 | Slide adapted from Bernt Schiele | B. Leibe | 38

RWTH AACHEN UNIVERSITY

EM Algorithm – An Example

Machine Learning Winter '18 | B. Leibe | Image source: G.M. Bishop, 2006 | 39

RWTH AACHEN UNIVERSITY

EM – Technical Advice

- When implementing EM, we need to take care to avoid singularities in the estimation!
 - Mixture components may collapse on single data points.
 - E.g. consider the case $\Sigma_k = \sigma_k^2 \mathbf{I}$ (this also holds in general)
 - Assume component j is exactly centered on data point \mathbf{x}_n . This data point will then contribute a term in the likelihood function

$$\mathcal{N}(\mathbf{x}_n|\mathbf{x}_n, \sigma_j^2 \mathbf{I}) = \frac{1}{\sqrt{2\pi}\sigma_j}$$
 - For $\sigma_j \rightarrow 0$, this term goes to infinity!

⇒ Need to introduce regularization

- Enforce minimum width for the Gaussians
- E.g., instead of Σ^{-1} , use $(\Sigma + \sigma_{\min} \mathbf{I})^{-1}$

Machine Learning Winter '18 | B. Leibe | Image source: G.M. Bishop, 2006 | 40

RWTH AACHEN UNIVERSITY

EM – Technical Advice (2)

- EM is very sensitive to the initialization
 - Will converge to a local optimum of E .
 - Convergence is relatively slow.
- ⇒ Initialize with k-Means to get better results!
 - k-Means is itself initialized randomly, will also only find a local optimum.
 - But convergence is much faster.
- Typical procedure
 - Run k-Means M times (e.g. $M = 10-100$).
 - Pick the best result (lowest error J).
 - Use this result to initialize EM
 - Set μ_j to the corresponding cluster mean from k-Means.
 - Initialize Σ_j to the sample covariance of the associated data points.

Machine Learning Winter '18 | B. Leibe | 41

Machine Learning Winter '18

K-Means Clustering Revisited

- Interpreting the procedure
 - Initialization: pick K arbitrary centroids (cluster means)
 - Assign each sample to the closest centroid. (E-Step)
 - Adjust the centroids to be the means of the samples assigned to them. (M-Step)
 - Go to step 2 (until no change)

B. Leibe 42

Machine Learning Winter '18

K-Means Clustering Revisited

- K-Means clustering essentially corresponds to a Gaussian Mixture Model (MoG or GMM) estimation with EM whenever
 - The covariances are of the K Gaussians are set to $\Sigma_j = \sigma^2 I$
 - For some small, fixed σ^2

Slide credit: Bernt Schiele B. Leibe 43

Machine Learning Winter '18

Summary: Gaussian Mixture Models

- Properties
 - Very general, can represent any (continuous) distribution.
 - Once trained, very fast to evaluate.
 - Can be updated online.
- Problems / Caveats
 - Some numerical issues in the implementation
 - Need to apply regularization in order to avoid singularities.
 - EM for MoG is computationally expensive
 - Especially for high-dimensional problems!
 - More computational overhead and slower convergence than k-Means
 - Results very sensitive to initialization
 - Run k-Means for some iterations as initialization!
 - Need to select the number of mixture components K .
 - Model selection problem (see later lecture)

B. Leibe 44

Machine Learning Winter '18

Topics of This Lecture

- Mixture distributions
 - Recap: Mixture of Gaussians (MoG)
 - Maximum Likelihood estimation attempt
- K-Means Clustering
 - Algorithm
 - Applications
- EM Algorithm
 - Cluster assignment problem
 - MoG estimation
 - EM Algorithm
 - Interpretation of K-Means
 - Technical advice
- Applications

B. Leibe 45

Machine Learning Winter '18

Applications

- Mixture models are used in many practical applications.
 - Wherever distributions with complex or unknown shapes need to be represented...
- Popular application in Computer Vision
 - Model distributions of pixel colors.
 - Each pixel is one data point in, e.g., RGB space.
 - Learn a MoG to represent the class-conditional densities.
 - Use the learned models to classify other pixels.

B. Leibe Image source: C.M. Bishop, 2006 46

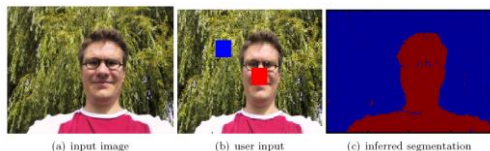
Machine Learning Winter '18

Application: Background Model for Tracking

- Train background MoG for each pixel
 - Model "common" appearance variation for each background pixel.
 - Initialization with an empty scene.
 - Update the mixtures over time
 - Adapt to lighting changes, etc.
- Used in many vision-based tracking applications
 - Anything that cannot be explained by the background model is labeled as foreground (=object).
 - Easy segmentation if camera is fixed.

C. Stauffer, E. Grimson, [Learning Patterns of Activity Using Real-Time Tracking](#), IEEE Trans. PAMI, 22(8):747-757, 2000. B. Leibe Image Source: Daniel Roth, Tobias Jaeger 47

Application: Image Segmentation



- User assisted image segmentation
 - User marks two regions for foreground and background.
 - Learn a MoG model for the color values in each region.
 - Use those models to classify all other pixels.
- ⇒ Simple segmentation procedure
(building block for more complex applications)

B. Leibe

48

References and Further Reading

- More information about EM and MoG estimation is available in Chapter 2.3.9 and the entire Chapter 9 of Bishop's book (recommendable to read).

Christopher M. Bishop
Pattern Recognition and Machine Learning
Springer, 2006



- Additional information
 - Original EM paper:
 - A.P. Dempster, N.M. Laird, D.B. Rubin, „[Maximum-Likelihood from incomplete data via EM algorithm](#)“, In Journal Royal Statistical Society, Series B. Vol 39, 1977
 - EM tutorial:
 - J.A. Bilmes, "[A Gentle Tutorial of the EM Algorithm and its Application to Parameter Estimation for Gaussian Mixture and Hidden Markov Models](#)", TR-97-021, ICSI, U.C. Berkeley, CA, USA

B. Leibe

51