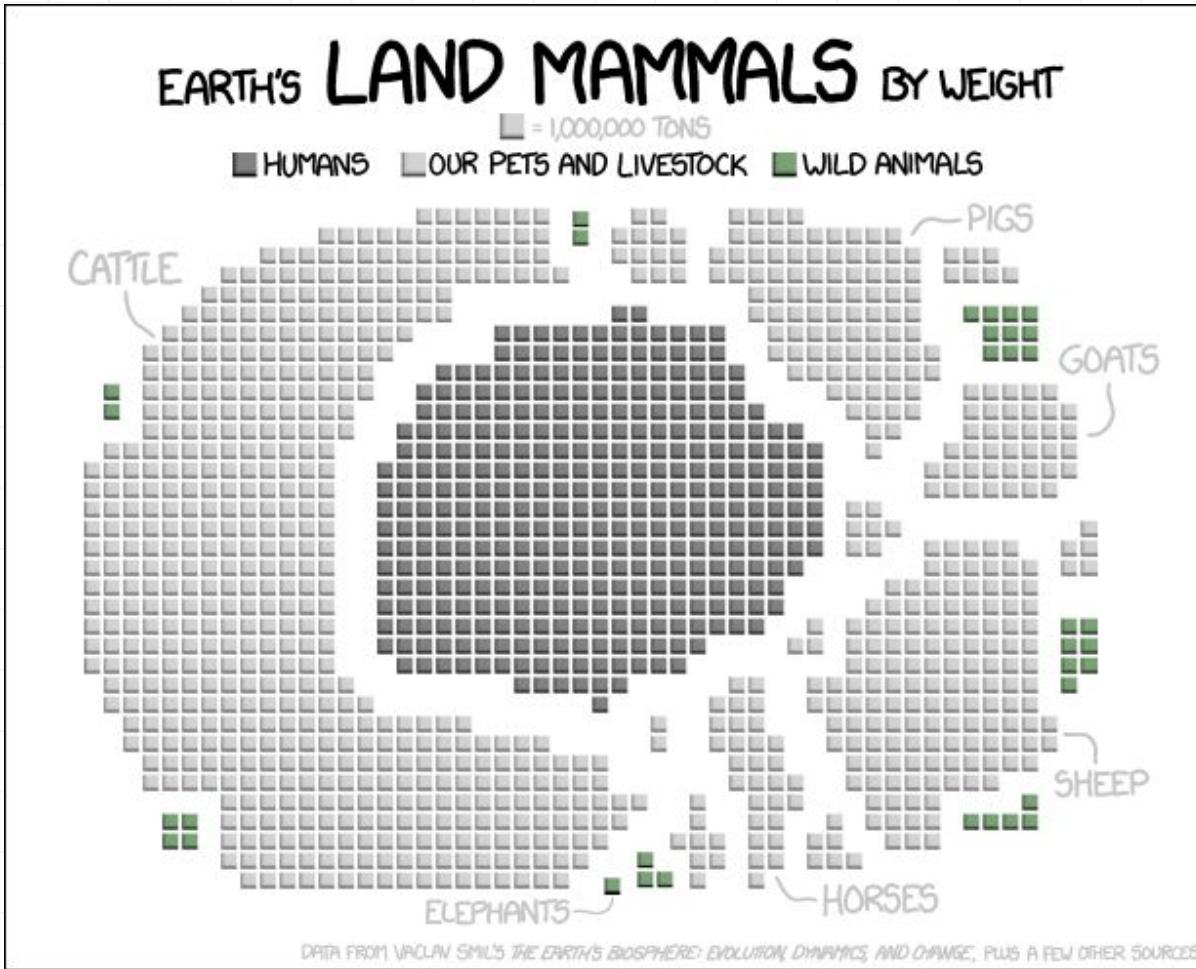


<https://xkcd.com/1338/>



Mixture Models and Expectation Maximisation

Pre-read/watch: [K-Means algorithm \(PRML 9.1\)](#) and update equations of Gaussian Mixture Models ([MML book Ch 11](#))

EM revisited

- An alternative view of EM
- Connections between GMM and K-means
- Bernoulli mixture

EM in general - does it really maximise likelihood, and why?

Practical considerations and other topics

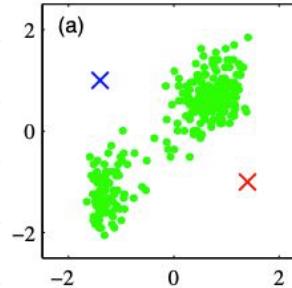
- impossibility of clustering [Kleinberg 2003]
- Kmeans++ [Vassilvitskii and Arthur, 2006]

Recap: what is clustering

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

K-means

- Each cluster contains points close to each other.
- Introduce a prototype $\mu_k \in \mathbb{R}^D$ for each cluster.
- Goal: Find
 - a set prototypes μ_k , $k = 1, \dots, K$, each representing a different cluster.
 - an assignment of each data point to exactly one cluster.



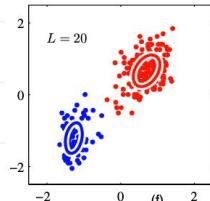
Luxburg, U.V., Williamson, R.C., & Guyon, I. (2012). Clustering: Science or Art? *ICML Unsupervised and Transfer Learning*.



These ambiguities, redundancies, and deficiencies recall those attributed by Dr. Franz Kuhn to a certain Chinese encyclopedia called the [Heavenly Emporium of Benevolent Knowledge](#). In its distant pages it is written that animals are divided into (a) those that belong to the emperor; (b) embalmed ones; (c) those that are trained; (d) sucking pigs; (e) mermaids; (f) fabulous ones; (g) stray dogs; (h) those that are included in this classification; (i) those that tremble as if they were mad; (j) innumerable ones; (k) those drawn with a very fine camel's-hair brush; (l) et cetera; (m) those that have just broken the flower vase; (n) those that at a distance resemble flies. (Borges, 1999)

(initialise a set of cluster centers / component means)

Fig 9.8



Expectation step

K-means

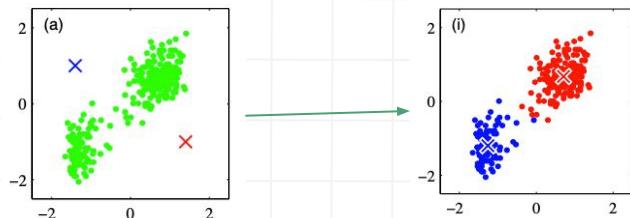
Re-assign data points to clusters,
determine r_{nk}

$$\{0, 1\}$$

Maximisation step

Re-compute the cluster means -
update $\{\mu_k\}$

Fig 9.1



Gaussian Mixture Models

2. **E step.** Evaluate the responsibilities using the current parameter values

$$\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}. \quad (9.23)$$

3. **M step.** Re-estimate the parameters using the current responsibilities

$$\boldsymbol{\mu}_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \quad (9.24)$$

$$\boldsymbol{\Sigma}_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}})^T \quad (9.25)$$

$$\pi_k^{\text{new}} = \frac{N_k}{N} \quad (9.26)$$

where

$$N_k = \sum_{n=1}^N \gamma(z_{nk}). \quad (9.27)$$

K-means and GMM - hard vs soft assignments

\mathbf{x}_n	r_{nk}	
	$\gamma(z_{nk})$	
		θ_k

For k -means clustering,
we have hard assignments

For GMM,
we have soft assignments

Assume a Gaussian mixture model.

Covariance matrices given by $\epsilon \mathbf{I}$, where ϵ is shared by all components.

Then

$$p(\mathbf{x} | \boldsymbol{\mu}_k, \Sigma_k) = \frac{1}{(2\pi\epsilon)^{D/2}} \exp \left\{ -\frac{1}{2\epsilon} \|\mathbf{x} - \boldsymbol{\mu}_k\|^2 \right\}.$$

$$\gamma(z_{nk}) = \frac{\pi_k \exp \left\{ -\frac{\|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2}{2\epsilon} \right\}}{\sum_j \pi_j \exp \left\{ -\frac{\|\mathbf{x}_n - \boldsymbol{\mu}_j\|^2}{2\epsilon} \right\}}$$

$\epsilon \downarrow 0$

$\nearrow \infty$

Taking the limit $\epsilon \rightarrow 0$

$$\gamma(z_{nk}) = \begin{cases} 1 & \text{if } \|\mathbf{x}_n - \boldsymbol{\mu}_k\| < \|\mathbf{x}_n - \boldsymbol{\mu}_j\| \quad \forall j \neq k \\ 0 & \text{otherwise} \end{cases}$$

Wait ... what is being maximised in EM?

$$\ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\} \quad (9.28) \quad \text{also (9.14)}$$

In each maximisation step?

Overall?

$$0 = - \sum_{n=1}^N \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_j \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \underbrace{\boldsymbol{\Sigma}_k (\mathbf{x}_n - \boldsymbol{\mu}_k)}_{\gamma(z_{nk})} \quad (9.16)$$

Cheating :)

M step. Re-estimate the parameters using the current responsibilities

$$\boldsymbol{\mu}_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \quad (9.24)$$

$$\boldsymbol{\Sigma}_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}})^T \quad (9.25)$$

$$\pi_k^{\text{new}} = \frac{N_k}{N} \quad (9.26)$$

Expectation-maximization revisited

X observed, Z “latent”

$$\ln p(\mathbf{X}|\boldsymbol{\theta}) = \ln \left\{ \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) \right\}. \quad (9.29)$$

$$z_{nk} \in \{0, 1\}$$

compute posterior $P(Z|X, \boldsymbol{\theta}) \rightarrow \gamma(z_{nk})$

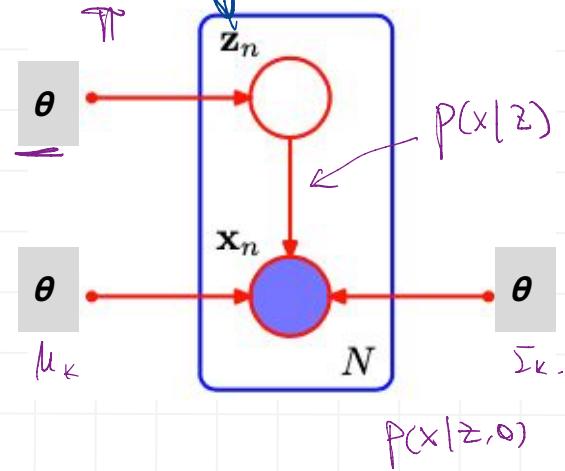
Take expectation of the complete data log-likelihood $P(X, Z|\boldsymbol{\theta})$ w.r.t. Z

$$\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}). \quad (9.30)$$

Maximise this expectation w.r.t. $\boldsymbol{\theta}$

$$\boldsymbol{\theta}^{\text{new}} = \arg \max_{\boldsymbol{\theta}} \mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}). \quad (9.31)$$

plate notation
 N copies represent x_1, \dots, x_N
Fig 9.6



$$P(Z|x, \theta)$$

update $\gamma(z_{nk})$

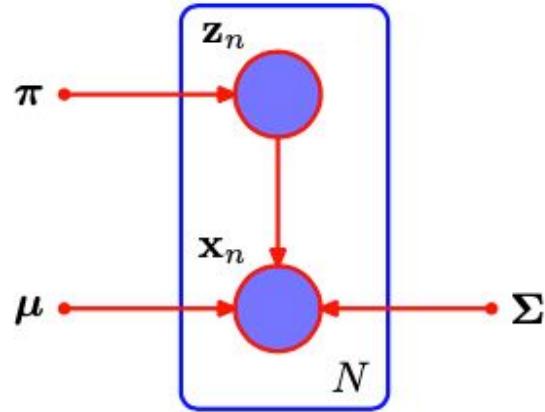
update $\{\pi_i, (\mu_k, \Sigma_k)\}$

Convince ourselves that this is true for GMM M-step

Figure 9.9 This shows the same graph as in Figure 9.6 except that we now suppose that the discrete variables z_n are observed, as well as the data variables x_n .

$$\theta^{\text{new}} = \arg \max_{\theta} Q(\theta, \theta^{\text{old}}). \quad (9.31)$$

$$\ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\} \quad (9.28)$$



$$p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = \prod_{n=1}^N \prod_{k=1}^K \pi_k^{z_{nk}} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_{nk}} \quad (9.35)$$

only \ln in k terms is active

$$\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \{ \ln \pi_k + \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \}. \quad (9.36)$$

$$p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) \propto \prod_{n=1}^N \prod_{k=1}^K [\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]^{z_{nk}}. \quad (9.38)$$

γ

$$\mathbb{E}_{\mathbf{Z}}[\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi})] = \sum_{n=1}^N \sum_{k=1}^K \underbrace{\gamma(z_{nk})}_{\theta^{\text{old}}} \{ \ln \pi_k + \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \}. \quad (9.40)$$

θ^{new}

The General EM Algorithm

Given a joint distribution $p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ over observed variables \mathbf{X} and latent variables \mathbf{Z} , governed by parameters $\boldsymbol{\theta}$, the goal is to maximize the likelihood function $p(\mathbf{X}|\boldsymbol{\theta})$ with respect to $\boldsymbol{\theta}$.

1. Choose an initial setting for the parameters $\boldsymbol{\theta}^{\text{old}}$.

2. **E step** Evaluate $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})$.

3. **M step** Evaluate $\boldsymbol{\theta}^{\text{new}}$ given by

MAP objective

$$\boldsymbol{\theta}^{\text{new}} = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) \quad (9.32) \quad Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) + \ln p(\boldsymbol{\theta})$$

where

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}). \quad (9.33)$$

4. Check for convergence of either the log likelihood or the parameter values.
If the convergence criterion is not satisfied, then let

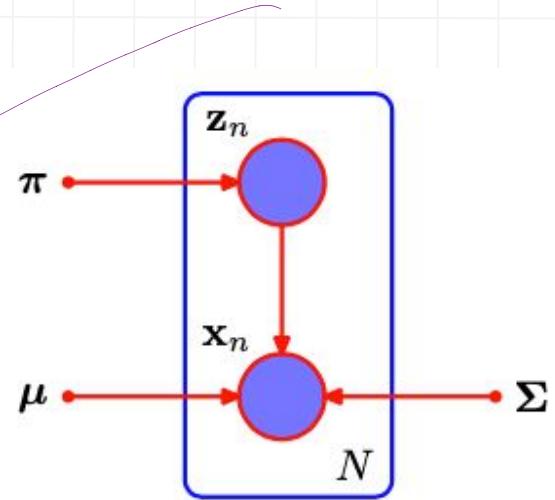
$$\boldsymbol{\theta}^{\text{old}} \leftarrow \boldsymbol{\theta}^{\text{new}} \quad (9.34)$$

and return to step 2.

EM for GMM, revisited

Figure 9.9 This shows the same graph as in Figure 9.6 except that we now suppose that the discrete variables z_n are observed, as well as the data variables x_n .

$$\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \{ \ln \pi_k + \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \}. \quad (9.36)$$



$$\mathbb{E}_{\mathbf{Z}} [\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi})] = \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \{ \ln \pi_k + \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \}. \quad (9.40)$$

$\partial \mathbb{E}[\cdot]$
 $\partial \boldsymbol{\mu}_k, \pi_k, \dots$

K-means and GMM - differences & connections

MLE

$$\text{GMM} \quad \ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}. \quad \begin{matrix} \dots \exp \{ (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \} \\ \vdots \\ \mathbf{x}_n \\ \boldsymbol{\mu}_k \\ \boldsymbol{\Sigma}_k \end{matrix} \quad (9.14)$$

K-Means

$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2 \quad \begin{matrix} \nearrow \text{only one term active} \\ \text{data} \\ \downarrow \\ \text{cluster centers} \end{matrix} \quad (9.1)$$

K-means and GMM - differences & connections

$$\rightarrow (\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k)$$

	$\boldsymbol{\mu}_k$	$\boldsymbol{\Sigma}_k$	π_k	N_k
GMM	$\frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n$	$\frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \underbrace{(\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T}_{}$	$\frac{N_k}{N}$	$\sum_{n=1}^N \gamma(z_{nk}).$
K-Means	$\frac{\sum_n r_{nk} \mathbf{x}_n}{\sum_n r_{nk}}.$	ϵI	Does not matter if > 0	$\sum_n r_{nk}$

$$\mathbf{x}_n \in \mathbb{R}^D \quad \text{GMM: } K + KD + K \cdot \frac{D(D+1)}{2}$$

kmean K.D

Complexity: how many parameters for each?

Implementation:

Can a cluster center “die” in each model, how to handle? \rightarrow Initialize another center,

What if some Gaussian components are singular?



D separate Bernoulli Distributions

a set of D binary variables x_i , each governed by Bernoulli distribution with mean μ_i

$$\underbrace{\mathbf{x} = (x_1, \dots, x_D)^T}_{\text{and } \underline{\boldsymbol{\mu}} = (\mu_1, \dots, \mu_D)^T} \quad \text{and } \boldsymbol{\mu} = (\mu_1, \dots, \mu_D)^T$$

x_i turn the i -th component on/off.

$$p(\mathbf{x}|\boldsymbol{\mu}) = \prod_{i=1}^D \mu_i^{x_i} (1 - \mu_i)^{(1-x_i)} \quad (9.44)$$

$$\mathbb{E}[\mathbf{x}] = \boldsymbol{\mu} \quad (9.45)$$

$$\text{cov}[\mathbf{x}] = \text{diag}\{\mu_i(1 - \mu_i)\}. \quad (9.46)$$

Mixture of Bernoulli Distributions

$$\mathbf{x} = (x_1, \dots, x_D)^T$$

$$\boldsymbol{\mu} = \{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K\}, \boldsymbol{\pi} = \{\pi_1, \dots, \pi_K\}$$

$$p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\pi}) = \sum_{k=1}^K \pi_k p(\mathbf{x}|\boldsymbol{\mu}_k) \quad (9.47)$$

$$p(\mathbf{x}|\boldsymbol{\mu}_k) = \prod_{i=1}^D \mu_{ki}^{x_i} (1 - \mu_{ki})^{(1-x_i)}. \quad (9.48)$$

$$\mathbb{E}[\mathbf{x}] = \sum_{k=1}^K \pi_k \boldsymbol{\mu}_k \quad (9.49)$$

$$\text{cov}[\mathbf{x}] = \sum_{k=1}^K \pi_k \left\{ \boldsymbol{\Sigma}_k + \boldsymbol{\mu}_k \boldsymbol{\mu}_k^T \right\} - \underline{\mathbb{E}[\mathbf{x}] \mathbb{E}[\mathbf{x}]^T} \quad (9.50)$$

$$\mathbf{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_D \end{pmatrix}$$

$$\boldsymbol{\mu} = \begin{bmatrix} \mu_{11} & \cdots & \mu_{1D} \\ \vdots & \ddots & \vdots \\ \mu_{D1} & \cdots & \mu_{DD} \end{bmatrix}^k$$

$$\mathbb{E}(x^4) - \underline{(\mathbb{E}(x))^4}$$

$$\boldsymbol{\Sigma}_k = \text{diag} \{ \mu_{ki} (1 - \mu_{ki}) \}$$

Data likelihood : (9.51)

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

$$p(\mathbf{x}|\mathbf{z}, \boldsymbol{\mu}) = \prod_{k=1}^K p(\mathbf{x}|\boldsymbol{\mu}_k)^{z_k} \quad p(\mathbf{z}|\boldsymbol{\pi}) = \prod_{k=1}^K \pi_k^{z_k}.$$

Complete data likelihood

$$\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\mu}, \boldsymbol{\pi}) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \left\{ \ln \pi_k + \sum_{i=1}^D [x_{ni} \ln \mu_{ki} + (1 - x_{ni}) \ln(1 - \mu_{ki})] \right\} \quad (9.54)$$

Expectations of complete data likelihood

$$\mathbb{E}_{\mathbf{Z}}[\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\mu}, \boldsymbol{\pi})] = \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \left\{ \ln \pi_k + \sum_{i=1}^D [x_{ni} \ln \mu_{ki} + (1 - x_{ni}) \ln(1 - \mu_{ki})] \right\} \quad (9.55)$$

where

$$\gamma(z_{nk}) = \mathbb{E}[z_{nk}] = \frac{\pi_k p(\mathbf{x}_n|\boldsymbol{\mu}_k)}{\sum_{j=1}^K \pi_j p(\mathbf{x}_n|\boldsymbol{\mu}_j)}. \quad (9.56)$$

Similar calculation as with mixture of Gaussian

$$\gamma(z_{nk}) = \frac{\pi_k p(\mathbf{x}_n | \boldsymbol{\mu}_k)}{\sum_{j=1}^K \pi_j p(\mathbf{x}_n | \boldsymbol{\mu}_j)}$$

$$N_k = \sum_{n=1}^N \gamma(z_{nk})$$

$$\bar{\mathbf{x}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \quad \mu_k = \bar{\mathbf{x}}$$

$$\pi_k = \frac{N_k}{N}$$

16×16
 $D=256$

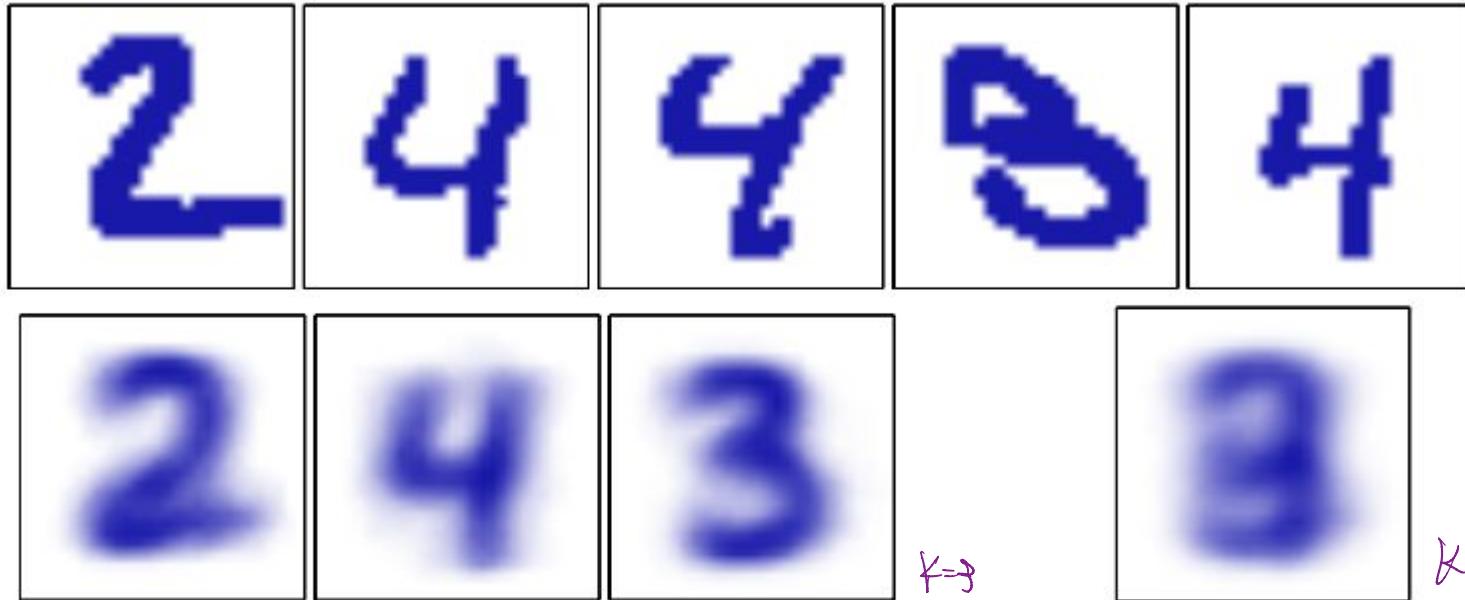


Figure 9.10 Illustration of the Bernoulli mixture model in which the top row shows examples from the digits data set after converting the pixel values from grey scale to binary using a threshold of 0.5. On the bottom row the first three images show the parameters μ_{ki} for each of the three components in the mixture model. As a comparison, we also fit the same data set using a single multivariate Bernoulli distribution, again using maximum likelihood. This amounts to simply averaging the counts in each pixel and is shown by the right-most image on the bottom row.

Outline

EM revisited

- Connections between GMM and K-means
- Bernoulli mixture

EM in general - does it really maximise likelihood, and why?

Practical considerations and other topics

- impossibility of clustering [Kleinberg 2003]
- Kmeans++ [Vassilvitskii and Arthur, 2006]

KL divergence (reminder)

$$\begin{aligned} \text{KL}(p\|q) &= - \int p(\mathbf{x}) \ln q(\mathbf{x}) d\mathbf{x} - \left(- \int p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x} \right) \\ &= - \int p(\mathbf{x}) \ln \left\{ \frac{q(\mathbf{x})}{p(\mathbf{x})} \right\} d\mathbf{x}. \end{aligned} \quad (1.113)$$

If we have “incorrectly” represented $p(\mathbf{x})$ with $q(\mathbf{x})$, how much more *information* do we need to recover $p(\mathbf{x})$?

↓
Apply Jensen’s inequality, $-\ln()$ is convex

$$f \left(\int \mathbf{x} p(\mathbf{x}) d\mathbf{x} \right) \leq \int f(\mathbf{x}) p(\mathbf{x}) d\mathbf{x}. \quad (1.117)$$

$$\text{KL}(p\|q) = - \int p(\mathbf{x}) \ln \left\{ \frac{q(\mathbf{x})}{p(\mathbf{x})} \right\} d\mathbf{x} \geq - \ln \int q(\mathbf{x}) d\mathbf{x} = 0 \quad (1.118)$$

Equality will only hold iff. $p(\mathbf{x}) = q(\mathbf{x})$ for all \mathbf{x}

The EM algorithm in general

Goal: show that the EM algorithm maximises the likelihood function.

$$p(\mathbf{X}|\boldsymbol{\theta}) = \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}). \quad (9.69)$$

$$\ln p(\mathbf{X}|\boldsymbol{\theta}) = \mathcal{L}(q, \boldsymbol{\theta}) + \text{KL}(q||p) \quad (9.70)$$

lower bound ↗
 ↓ ↗
 ↗

where we have defined

$$\mathcal{L}(q, \boldsymbol{\theta}) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{q(\mathbf{Z})} \right\} \quad (9.71)$$

$$\text{KL}(q||p) = - \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})}{q(\mathbf{Z})} \right\}. \quad (9.72)$$

Derive (9.70)

Likelihood

complete data likelihood

$$P(X|\theta) = \sum_z P(X, z|\theta)$$

sum rule

$$\ln P(X|\theta) + \ln P(z|x,\theta) = \ln P(X, z|\theta)$$

$$\ln P(X|\theta) = \ln P(X, z|\theta) - \ln P(z|x,\theta)$$

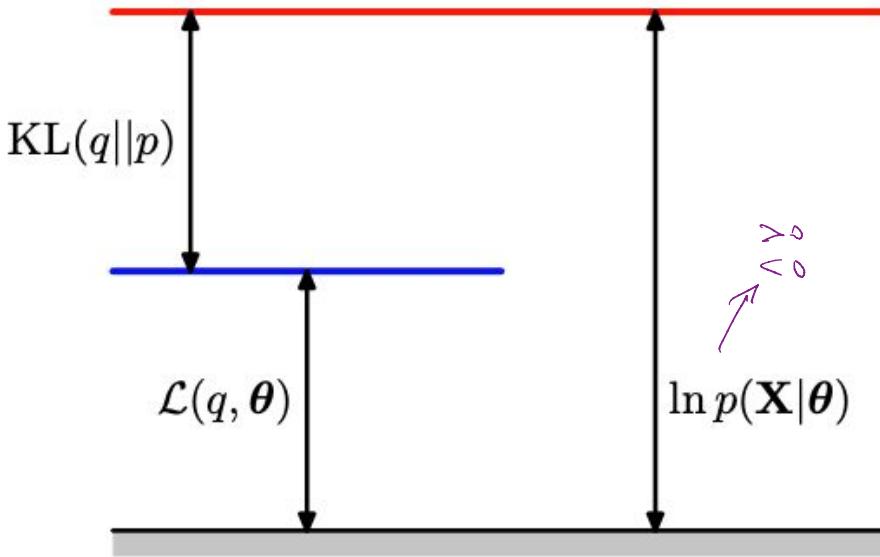
introduce $q(z)$

$$\begin{aligned} \sum_z q(z) &= 1 \\ &= \sum_z q(z) \ln P(X|\theta) \\ &= \sum_z q(z) [\ln P(X, z|\theta) - \ln P(z|x,\theta) + \ln q(z) - \ln q(z)] \\ &= \sum_z q(z) [\ln P(X, z|\theta) - \ln q(z)] - \sum_z q(z) \frac{\ln \frac{P(z|x,\theta)}{q(z)}}{\ln q(z)} \\ &\quad \text{L}(q, \theta) \quad + \quad \text{KL}(q(z) \parallel P(z|x,\theta)) \end{aligned}$$

Illustrating the decomposition

$$\ln p(\mathbf{X}|\boldsymbol{\theta}) = \mathcal{L}(q, \boldsymbol{\theta}) + \text{KL}(q\|p) \quad (9.70)$$

Figure 9.11 Illustration of the decomposition given by (9.70), which holds for any choice of distribution $q(\mathbf{Z})$. Because the Kullback-Leibler divergence satisfies $\text{KL}(q\|p) \geq 0$, we see that the quantity $\mathcal{L}(q, \boldsymbol{\theta})$ is a lower bound on the log likelihood function $\ln p(\mathbf{X}|\boldsymbol{\theta})$.



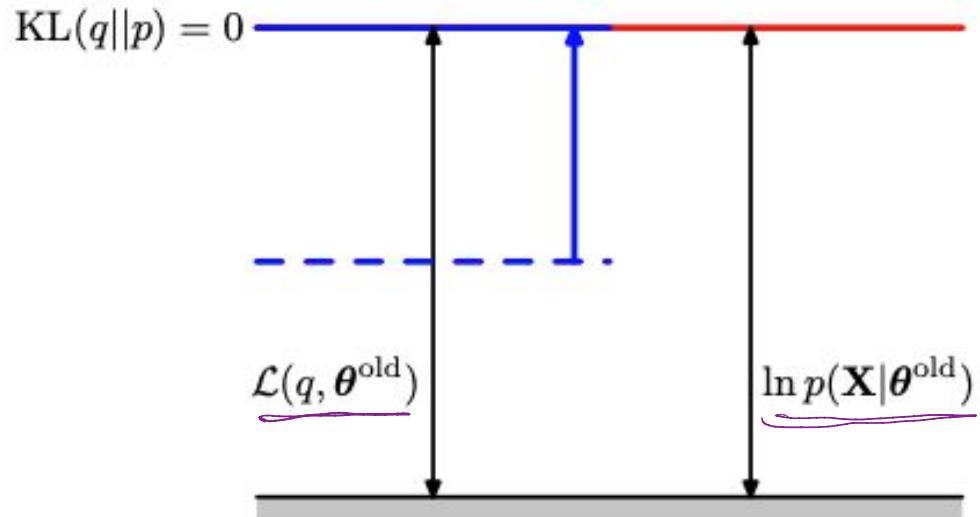
E step

set $q(Z) = p(Z | X, \theta)$

$$KL(q||p) = - \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{Z}|\mathbf{X}, \theta)}{q(\mathbf{Z})} \right\}. \quad (9.72)$$

Figure 9.12

Illustration of the E step of the EM algorithm. The q distribution is set equal to the posterior distribution for the current parameter values θ^{old} , causing the lower bound to move up to the same value as the log likelihood function, with the KL divergence vanishing.



3. M step Evaluate θ^{new} given by

$$\theta^{\text{new}} = \arg \max_{\theta} Q(\theta, \theta^{\text{old}}) \quad (9.32)$$

where

$$Q(\theta, \theta^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\theta). \quad (9.33)$$

$$\mathcal{L}(q, \theta) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{X}, \mathbf{Z}|\theta)}{q(\mathbf{Z})} \right\}$$

(9.71)

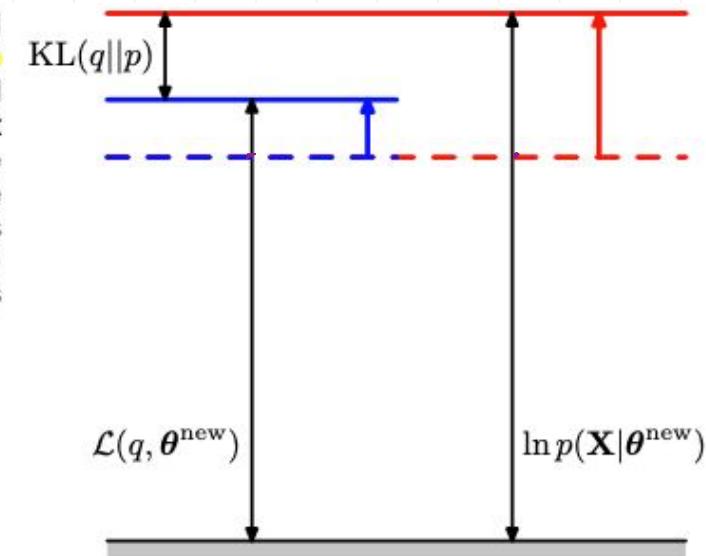
updated.

$H_q(\mathbf{Z}) \rightarrow \text{fixed}$.

$p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})$

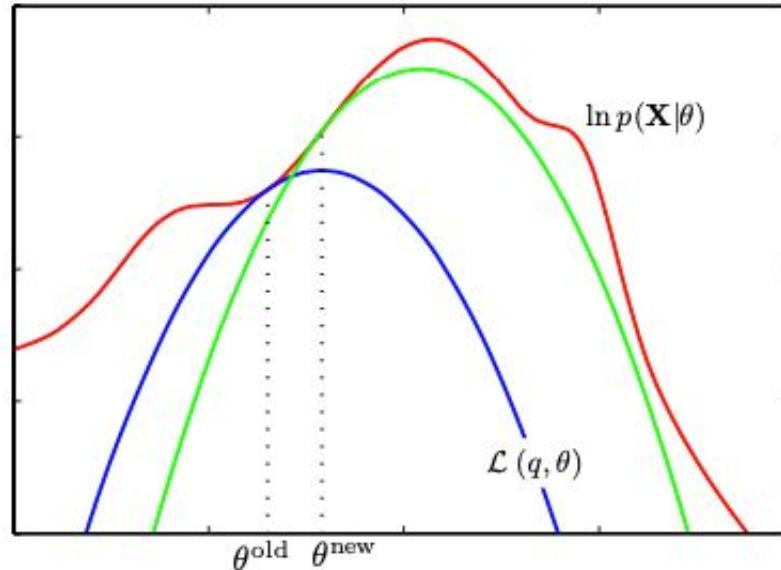
Figure 9.13

Illustration of the M step of the EM algorithm. The distribution $q(\mathbf{Z})$ is held fixed and the lower bound $\mathcal{L}(q, \theta)$ is maximized with respect to the parameter vector θ to give a revised value θ^{new} . Because the KL divergence is nonnegative, this causes the log likelihood $\ln p(\mathbf{X}|\theta)$ to increase by at least as much as the lower bound does.



EM as alternating maximization

Figure 9.14 The EM algorithm involves alternately computing a lower bound on the log likelihood for the current parameter values and then maximizing this bound to obtain the new parameter values. See the text for a full discussion.



Lower bound $L(q, \theta)$ is a convex function having a unique maximum (for mixture components from the exponential family).

Extensions: Generalised EM seeks to improve rather than maximise $L(q, \theta)$; expectation conditional maximisation seeks to maximise $L(q, \theta)$ for a subset of the parameters; Incremental algorithms also exist.

Outline

EM revisited

- Connections between GMM and K-means
- Bernoulli mixture

EM in general - does it really maximise likelihood, and why?

Practical considerations and other topics

- impossibility of clustering [Kleinberg 2003]
- Kmeans++ [Vassilvitskii and Arthur, 2006]

An Impossibility Theorem for Clustering

[NeuRIPS 2003]



Jon Kleinberg

Professor of Computer Science, Cornell University
Verified email at cs.cornell.edu - [Homepage](#)

algorithms data mining information networks

Theorem 2.1 For each $n \geq 2$, there is no clustering function f that satisfies Scale-Invariance, Richness, and Consistency.

SCALE-INVARIANCE. For any distance function d and any $\alpha > 0$, we have $f(d) = f(\alpha \cdot d)$.

RICHNESS. Range(f) is equal to the set of all partitions of S .

CONSISTENCY. Let d and d' be two distance functions. If $f(d) = \Gamma$, and d' is a Γ -transformation of d , then $f(d') = \Gamma$.

shrink or not but, points in the same cluster

Single-linkage operates by initializing each point as its own cluster, and then repeatedly merging the pair of clusters whose distance to one another (as measured from their closest points of approach) is minimum.

Stopping conditions:

K-clusters	Distance- r	Scale- α	K-means
✓		✓	
	✓		✓
		✓	X

K-means ++

Vassilvitskii, Sergei, and David Arthur. "k-means++: The advantages of careful seeding." In *Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*, pp. 1027-1035. 2006.



Sergei Vassilvitskii

I am a Research Scientist at [Google](#) New York. Previously I was a Research Scientist at [Yahoo! Research](#) and an Adjunct Assistant Professor at [Columbia University](#). I completed my PhD at [Stanford University](#) under the supervision of [Rajeev Motwani](#). Prior to that I was an undergraduate at [Cornell University](#).

sergei@cs.stanford.edu

Problem with K-means

Finds a local optimum that *could be arbitrarily worse than the global optimum.*

Algorithm summary:

1. Choose one center uniformly at random among the data points.
2. For each data point x not chosen yet, compute $D(x)$, the distance between x and the nearest center that has already been chosen.
3. Choose one new data point at random as a new center, using a weighted probability distribution where a point x is chosen with probability proportional to $D(x)^2$.
4. Repeat Steps 2 and 3 until k centers have been chosen.
5. Now that the initial centers have been chosen, proceed using standard [k-means clustering](#).

Theorem: k-means++ is $\Theta(\log k)$ approximate in expectation.

What do you do in practice? Normalise data. Use Kmeans++ to initialise Kmeans (`sklearn.cluster.kmeans_plusplus`). Use Kmeans (and spherical covariances) to initialise GMM. Try a number of different initialisations `sklearn.cluster.KMeans(n_clusters=8, *, init='k-means++', n_init=10, ...)`

Summary for today

EM revisited

- Connections between GMM and K-means
- Bernoulli mixture

EM in general - does it really maximise likelihood, and why?

Practical considerations and other topics

- impossibility of clustering [Kleinberg 2003]
- Kmeans++ [Vassilvitskii and Arthur, 2006]