

Class is 170.

Announcements

Matlab Grader homework, Binary graded.

168 (HW1), 166,165 (HW2) has done the homework. **(If you have not done HW talk to me/TA!)**

Homework 3 due **5 May TODAY**

Homework 4 (SVM +Dictionary Learning) due ~24 May, released soon

Jupiter “GPU” homework. Due 10 May

Today:

- Stanford CNN 12, K-means, EM (Bishop 9) ← *Mike Bincio*
- Play with Tensorflow playground before class <http://playground.tensorflow.org>
Solve the spiral problem

Monday

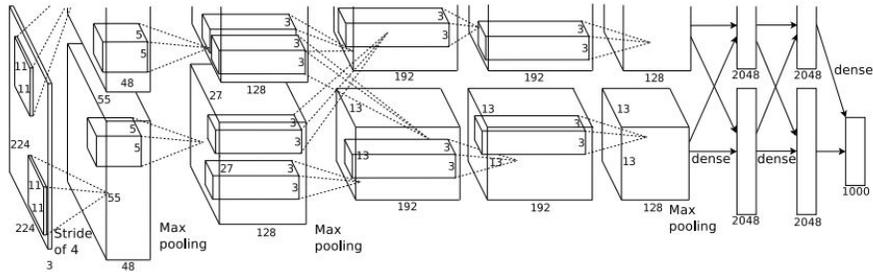
- Stanford CNN 13, Dictionary Learning,

Projects

- **3-4** person groups preferred
- Deliverables: Poster, Report & main code (plus proposal, midterm slide)
- **Topics** your own or chose from suggested topics. Some **physics inspired**.
- **April 26** groups due to TA. **5 students not signed up 44 Groups formed. Guidelines is on Piazza**
- **May 5** proposal due. use dropbox Format "Proposal"+groupNumber
<https://www.dropbox.com/request/XGqCV0qXm9LBZ7J1msS>
- **Wednesday May 22** Midterm slide presentation.
- **Project discussion, 22 May:** We split into 6 sub-classes. The purpose is to make sure your project is on track, good progress and good goals. Each group gives a ~10 min presentation by all members
 - (each person talks for ~2 min, ~1 slide)
 - Motivation & background, which data?
 - small Example,
 - final outcome, (focused on method and data)
 - difficulties,
- There are 7 Groups in each sub-class, thus we have 15 min in total/group. And will use the remaining time for discussion.
- **June 5, 5-8pm** poster. Atkinson Hall with Pizza. Upload June ~3
- Report and code due **Saturday 15 June.**

What's going on inside ConvNets?

This image is CC0 public domain



Class Scores:
1000 numbers

Input Image:
3 x 224 x 224

↑ ↑ ↑ ↑ ↑ ↑ ↑
What are the intermediate features looking for?

Image Processing

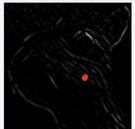
The general expression of a convolution is

$$g(x, y) = \omega * f(x, y) = \sum_{s=-a}^a \sum_{t=-b}^b \omega(s, t) f(x - s, y - t),$$

where $g(x, y)$ is the filtered image, $f(x, y)$ is the original image, ω is the filter kernel. Every element of the filter kernel is considered by $-a \leq s \leq a$ and $-b \leq t \leq b$.

Depending on the element values, a kernel can cause a wide range of effects.

$$\frac{f_{+1} - f_{-1}}{2}$$

		Image result $g(x,y)$
Identity	$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$	
Edge detection	$\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix}$	
	$\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$	
	$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$	
Sharpen	$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$	
Box blur (normalized)	$\frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$	

Feature Inversion

Given a CNN feature vector for an image, find a new image that:

- Matches the given feature vector
- “looks natural” (image prior regularization)

$$\mathbf{x}^* = \underset{\mathbf{x} \in \mathbb{R}^{H \times W \times C}}{\operatorname{argmin}} \ell(\Phi(\mathbf{x}), \Phi_0) + \lambda \mathcal{R}(\mathbf{x})$$

Given feature vector

Features of new image

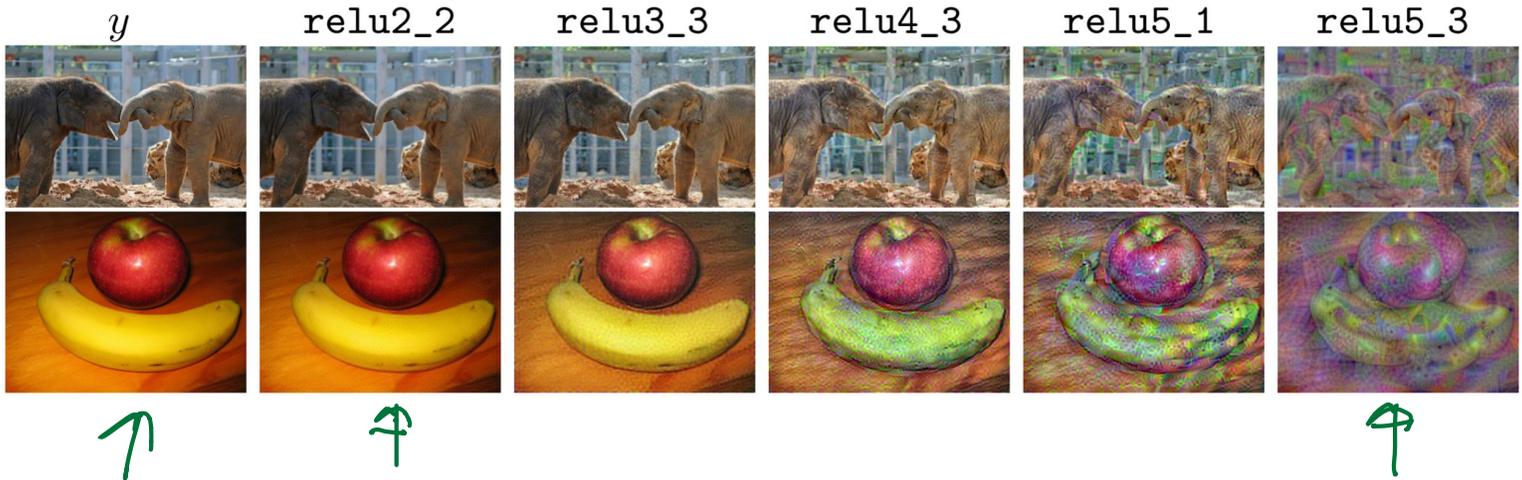
$$\ell(\Phi(\mathbf{x}), \Phi_0) = \|\Phi(\mathbf{x}) - \Phi_0\|^2$$

$$\mathcal{R}_{V^\beta}(\mathbf{x}) = \sum_{i,j} \left((x_{i,j+1} - x_{ij})^2 + (x_{i+1,j} - x_{ij})^2 \right)^{\frac{\beta}{2}}$$

Total Variation regularizer
(encourages spatial smoothness)

Feature Inversion

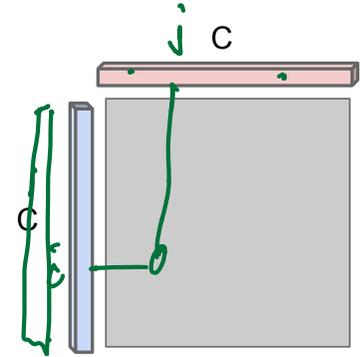
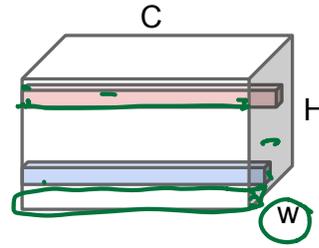
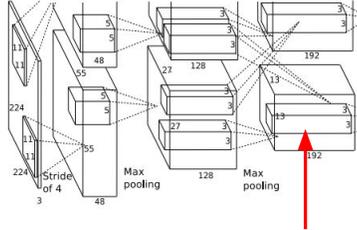
Reconstructing from different layers of VGG-16



Neural Texture Synthesis: Gram Matrix



This image is in the public domain.



Each layer of CNN gives $C \times H \times W$ tensor of features; $H \times W$ grid of C -dimensional vectors

Outer product of two C -dimensional vectors gives $C \times C$ matrix measuring co-occurrence

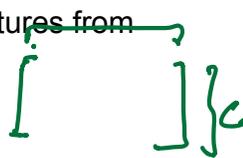
Average over all HW pairs of vectors, giving **Gram matrix** of shape $C \times C$

Efficient to compute; reshape features from $H \times W$ column

$C \times H \times W$ to $= C \times HW$

then compute $G = \underline{F} F^T$

HW column



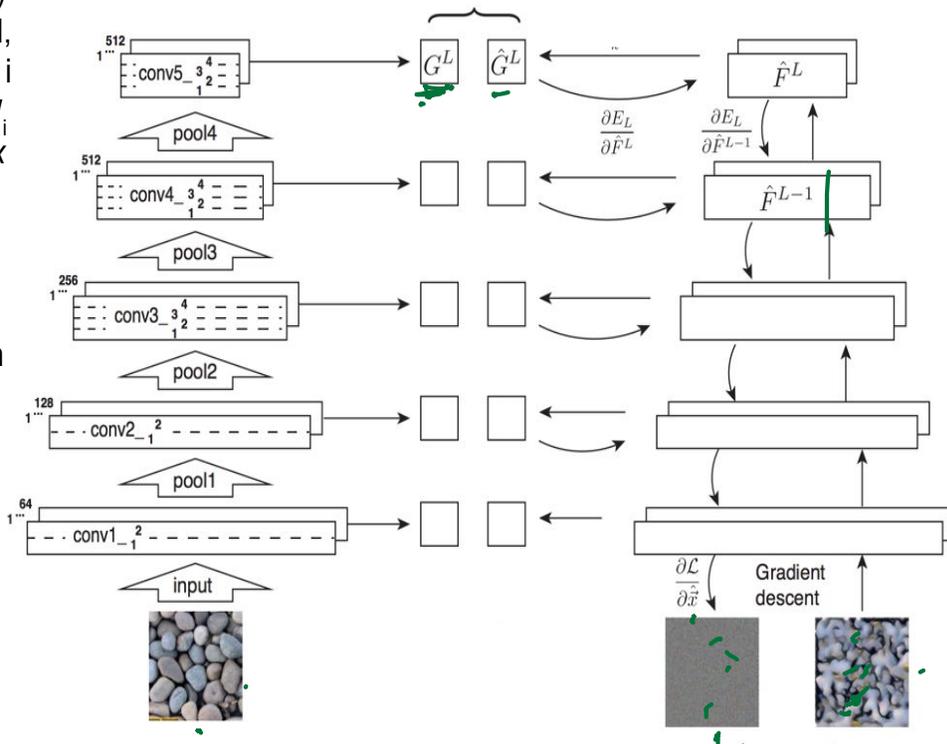
Neural Texture Synthesis

$$E_l = \frac{1}{4N_l^2 M_l^2} \sum_{i,j} (G_{ij}^l - \hat{G}_{ij}^l)^2 \quad \mathcal{L}(\vec{x}, \hat{\vec{x}}) = \sum_{l=0}^L w_l E_l$$

- 1. Pretrain a CNN on ImageNet (VGG-19)
2. Run input texture forward through CNN, record activations on every layer; layer i gives feature map of shape $C_i \times H_i \times W_i$
3. At each layer compute the *Gram matrix* giving outer product of features:

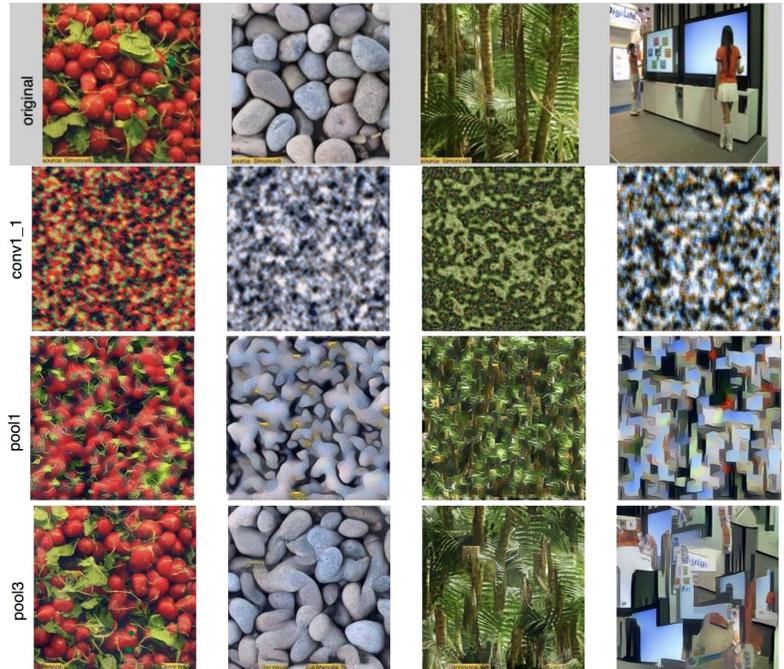
$$G_{ij}^l = \sum_k F_{ik}^l F_{jk}^l \quad (\text{shape } C_i \times C_i)$$

4. Initialize generated image from random noise
5. Pass generated image through CNN, compute Gram matrix on each layer
6. Compute loss: weighted sum of L2 distance between Gram matrices
7. Backprop to get gradient on image
8. Make gradient step on image
9. GOTO 5



Neural Texture Synthesis

Reconstructing texture from higher layers recovers larger features from the input texture



deep

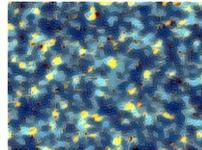
y

relu1_2

relu2_2

relu3_3

relu4_3



Texture synthesis
(Gram reconstruction)



Neural Style Transfer

Content Image



This image is licensed under CC-BY 3.0

+

Style Image



Starry Night by Van Gogh is in the public domain

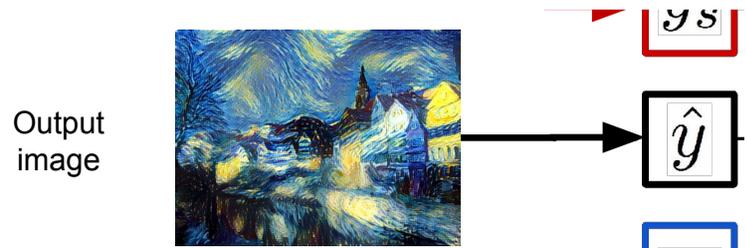
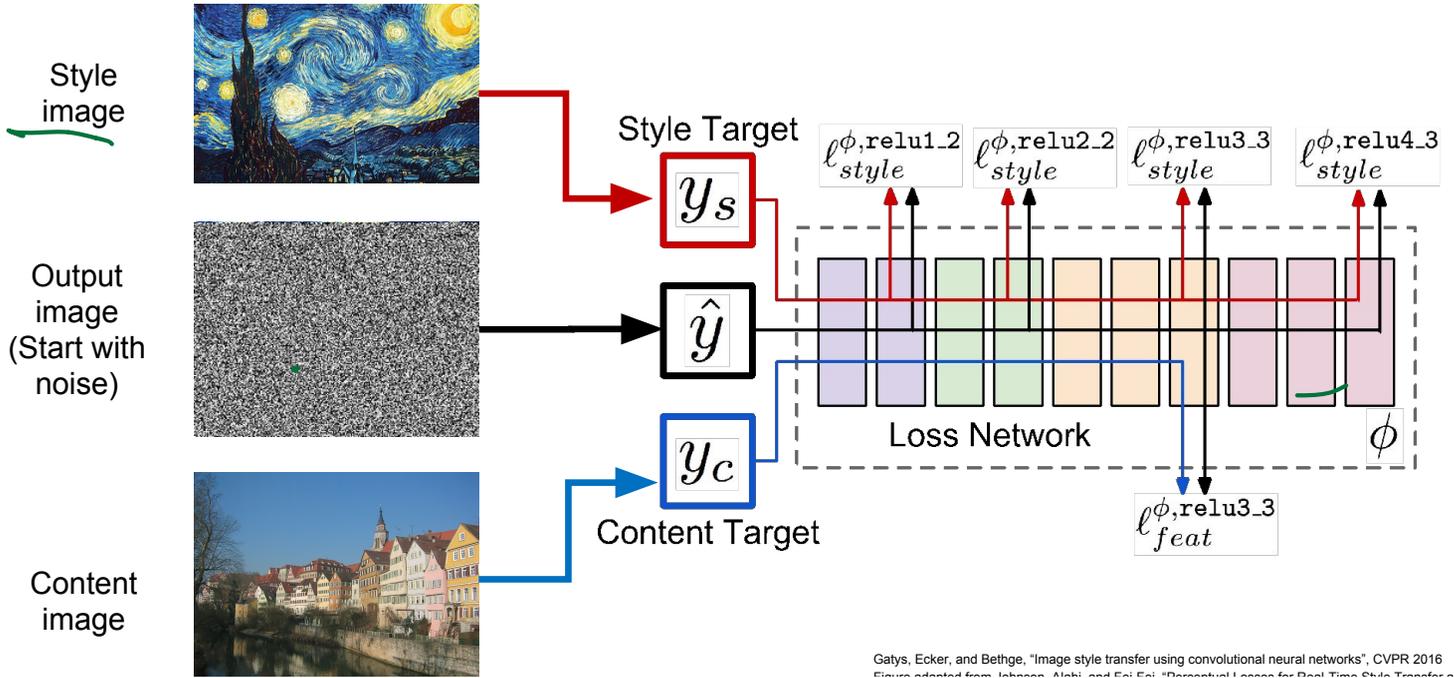
=

Style Transfer!



This image copyright Justin Johnson, 2015. Reproduced with permission.

Style transfer = Feature reconstruction loss + Gram matrix



K-means and Expectation Maximization

Mike Bianco and Peter Gerstoft

ECE228

5/6/2019

K-means and expectation maximization (EM) can be considered unsupervised learning

- In *supervised learning*, we have desired machine learning (ML) model output or ‘action’ \mathbf{y} based on inputs \mathbf{x} (features), and model parameters θ
 - Probabilities of the form: $p(\mathbf{y}|\mathbf{x},\theta)$
 - Linear regression and classification, support vector machines, etc.
- In *unsupervised learning*, we are interested in discovering useful patterns in the features. This can be for discovering latent data ‘causes’ or significant ‘groups’
 - Probabilities of the form: $p(\mathbf{x}|\theta)$
 - Principal components analysis (PCA), K-means, dictionary learning, etc.

Unsupervised learning

Unsupervised machine learning is inferring a function to describe hidden structure from "unlabeled" data (a classification or categorization is not included in the observations). Since the examples given to the learner are unlabeled, there is no evaluation of the accuracy of the structure that is output by the relevant algorithm—which is one way of distinguishing unsupervised learning from supervised learning.

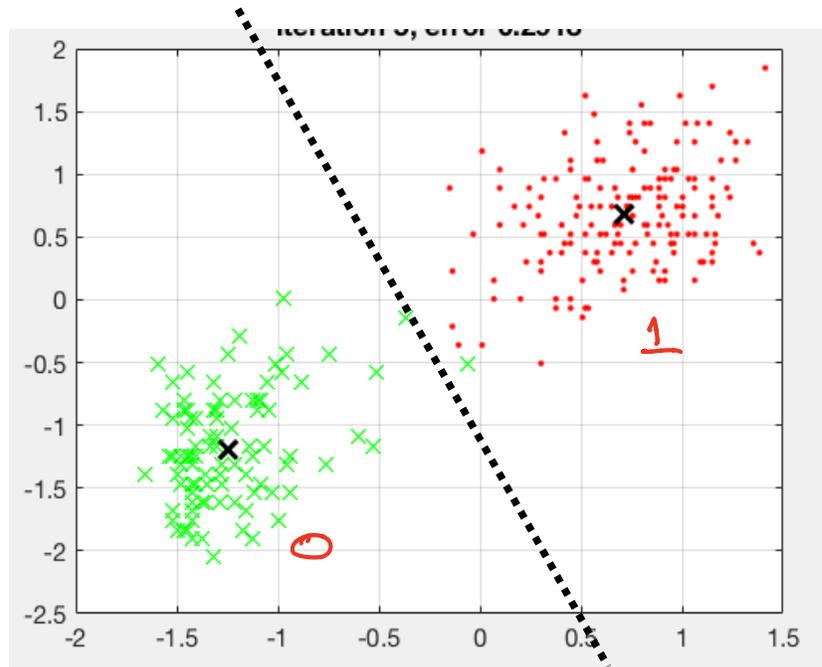
We are not interested in prediction

Supervised learning: all classification and regression.

$$Y = \mathbf{w}^T \mathbf{X}$$

Prediction is important.

Supervised learning: least square classifier (binary)



$$\|y - Xw\|_2$$

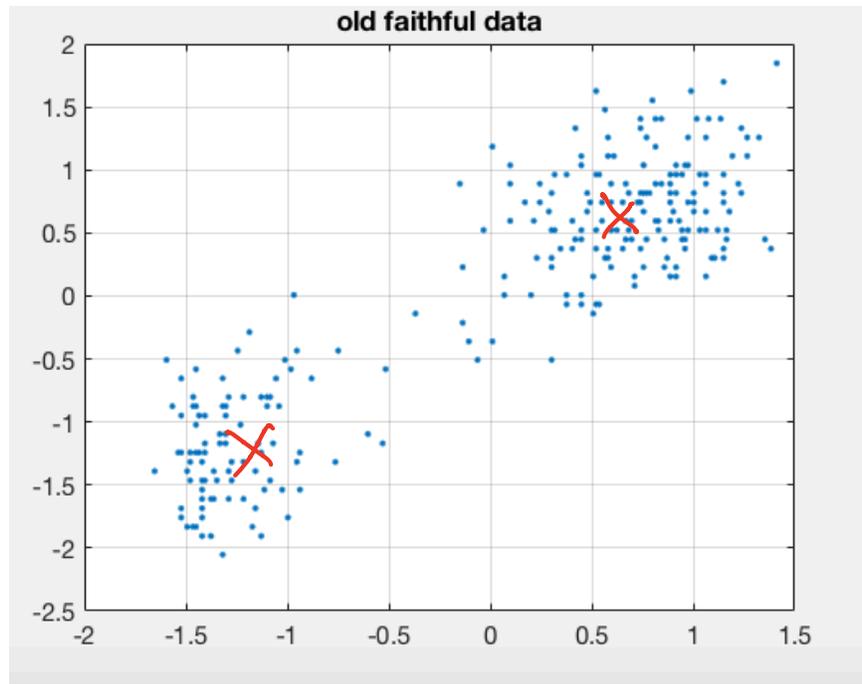
$$y = Xw$$

Training set $\{(x^1, y^1), (x^2, y^2), (x^3, y^3)\}$

We are given the two classes

(green = 0, red = 1)

Unsupervised learning: how are features best divided?



K-means

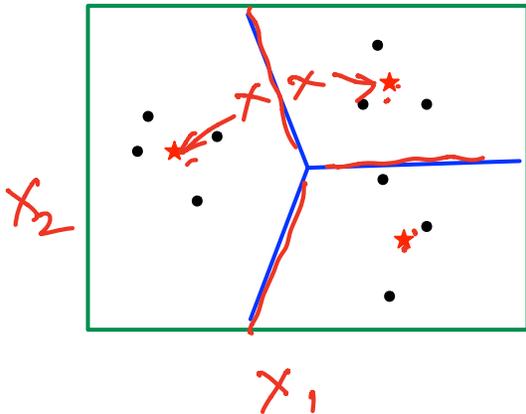
Just have features $\{(x_1^1, x_2^1), (x_1^2, x_2^2), (x_1^3, x_2^3)\}$

K-means

- **Input:** Points $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^p$; integer K
- **Output:** “Centers”, or representatives, $\mu_1, \dots, \mu_K \in \mathbb{R}^p$
- Output also $\mathbf{z}_1, \dots, \mathbf{z}_N \in \mathbb{R}^K$

Goal: Minimize average squared distance between points and their nearest representatives:

- $cost(\mu_1, \dots, \mu_K) = \sum_{n=1}^N \min_j \|x_n - \mu_j\|^2$



The centers carve \mathbb{R}^p up into k convex regions: μ_j 's region consists of points for which it is the closest center.

K-means

$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2 \quad (9.1)$$

① assign ptr. centroids.

Solving for r_{nk}

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

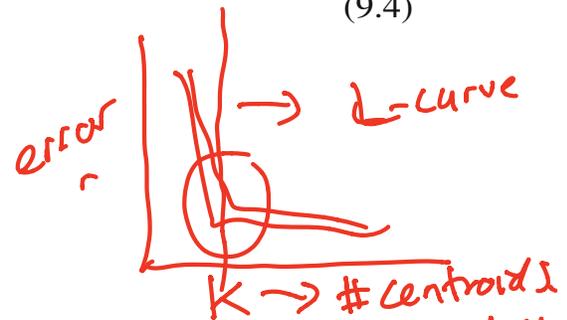
② update centroids based on ptr. (9.2)

Differentiating for $\boldsymbol{\mu}_k$

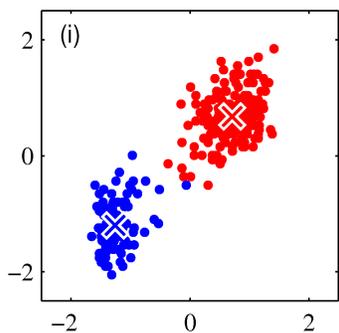
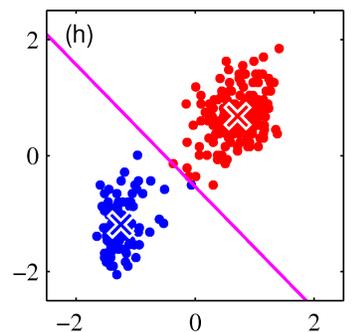
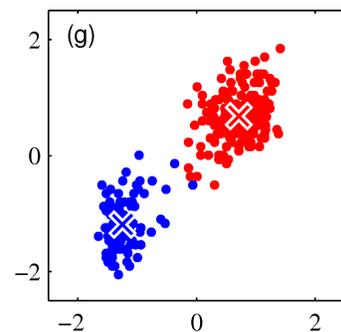
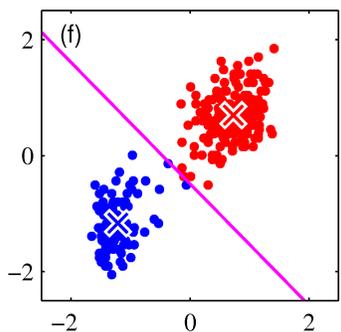
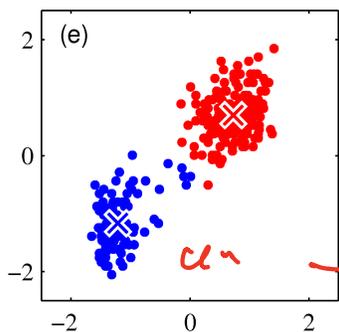
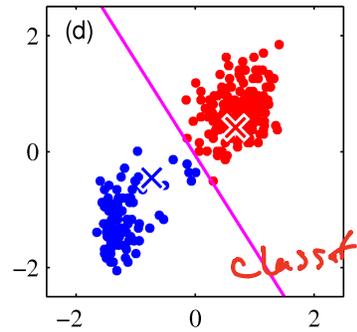
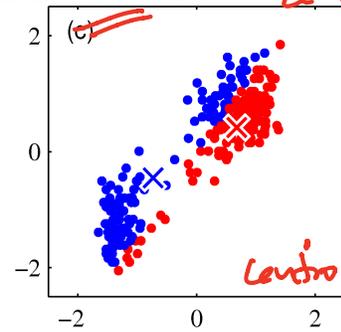
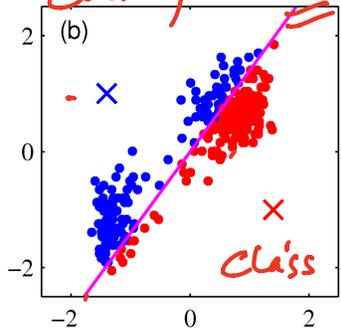
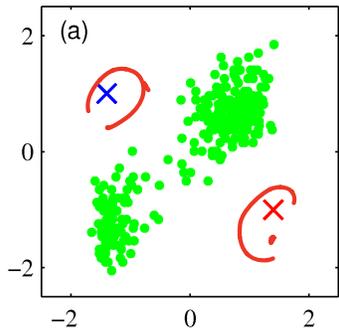
$$2 \sum_{n=1}^N r_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k) = 0 \quad (9.3)$$

which we can easily solve for $\boldsymbol{\mu}_k$ to give

$$\boldsymbol{\mu}_k = \frac{\sum_n r_{nk} \mathbf{x}_n}{\sum_n r_{nk}}. \quad (9.4)$$

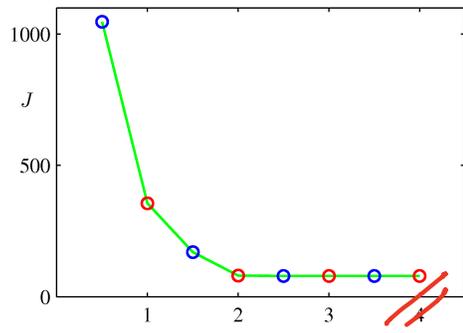


K-means Gap statistic

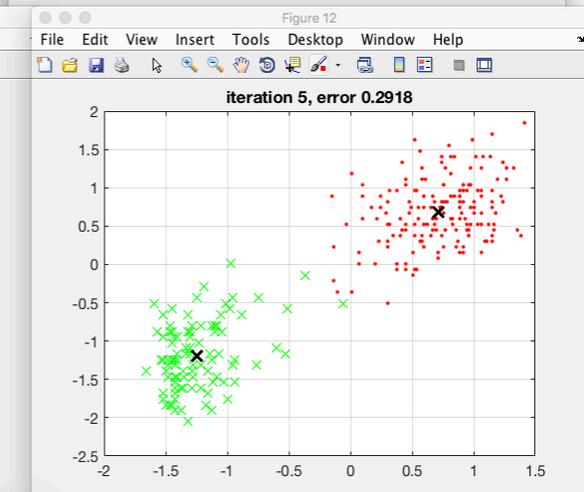
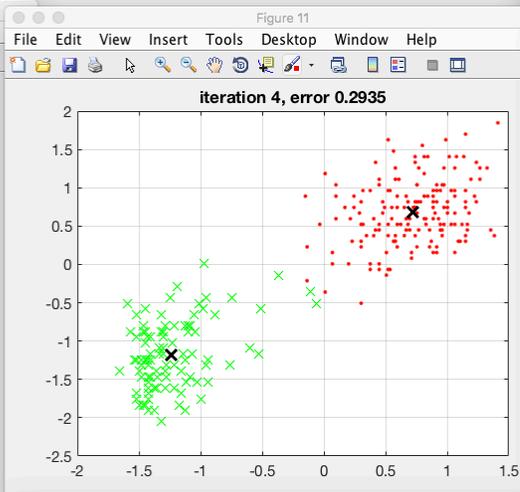
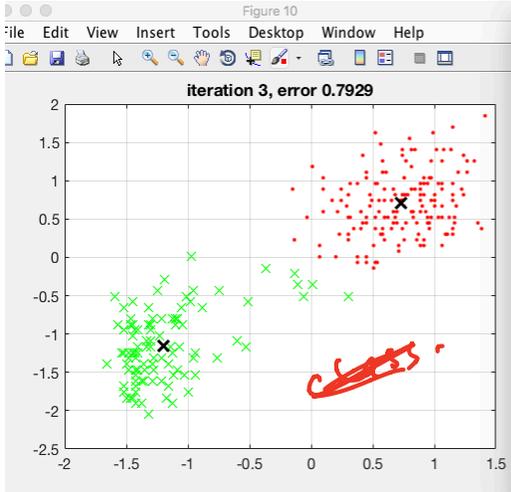
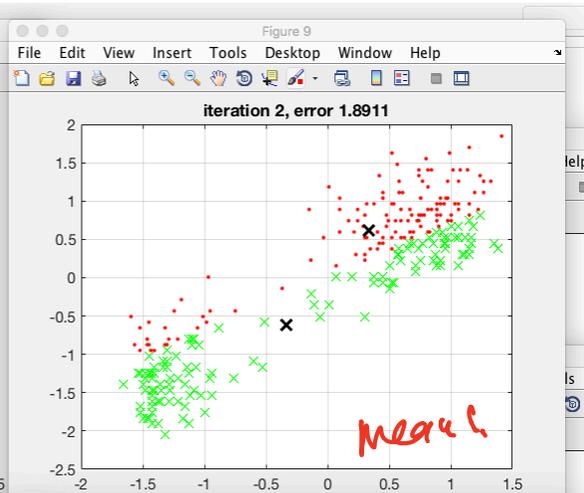
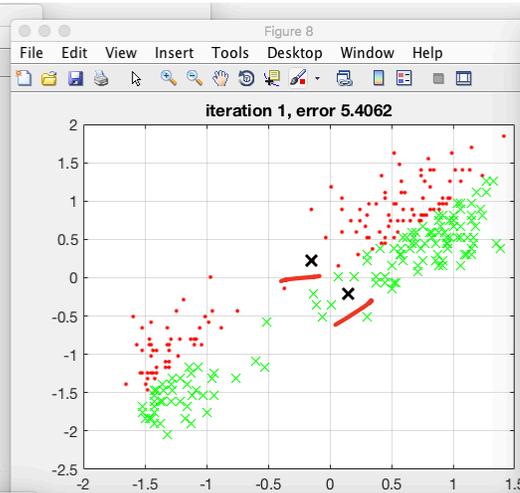
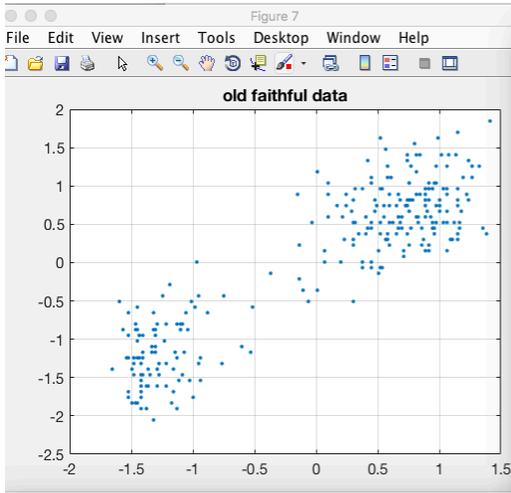


$$r_{nk} = \begin{cases} 1 & \text{if } k = \arg \min_j \|\mathbf{x}_n - \boldsymbol{\mu}_j\|^2 \\ 0 & \text{otherwise.} \end{cases}$$

$$\boldsymbol{\mu}_k = \frac{\sum_n r_{nk} \mathbf{x}_n}{\sum_n r_{nk}}$$



Old Faithful, Kmeans from Murphy



Application of K-means to data compression: Vector Quantization

Each pixel x_i is represented
By codebook of K entries μ_k

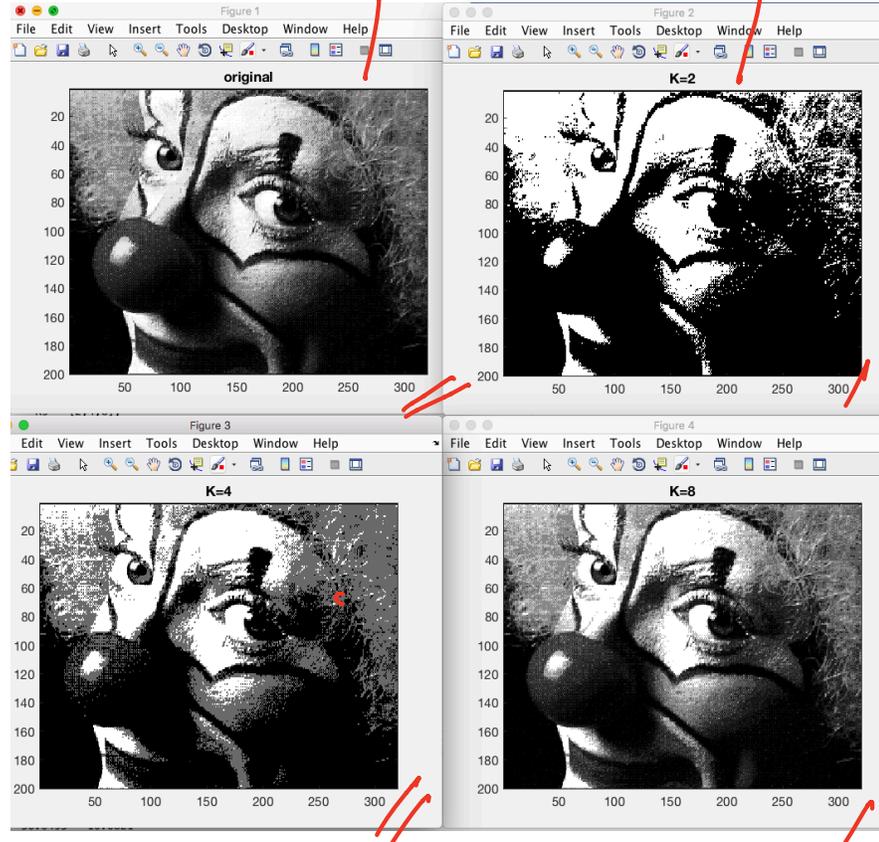
$$\text{Encode}(x_i) = \underset{k}{\text{argmin}} \|x_i - \mu_k\|$$

Consider $N=64k$ observations, of
 $D=1$ (b/w) dimension, $C=8$ bit

$$NC=513k$$

$N \log_2 K + KC$ bits is needed
 $K=4$ gives 128k a factor 4.

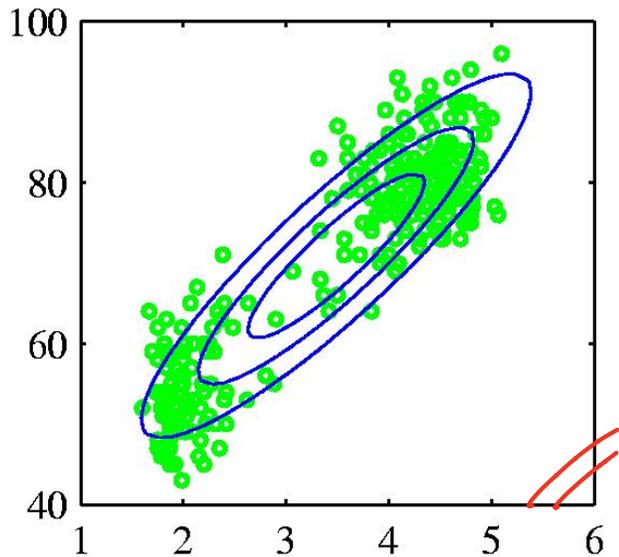
"codebook"



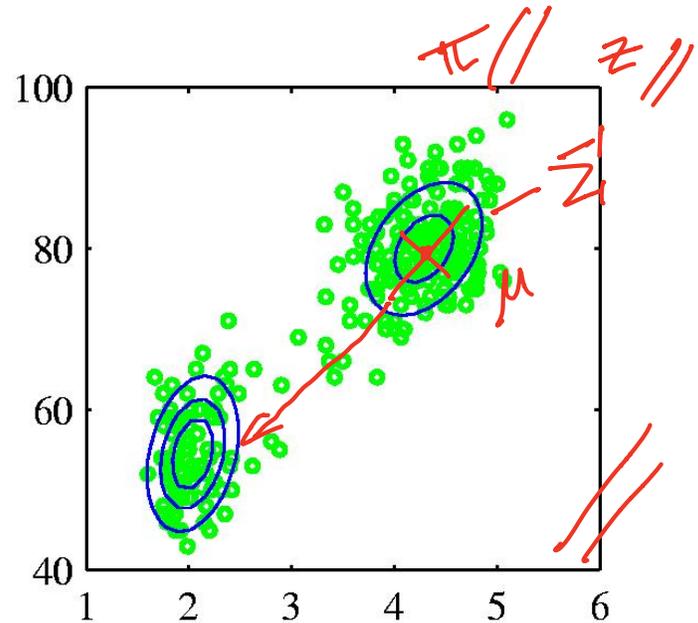
Mixtures of Gaussians (1)

Old Faithful geyser:

The time between eruptions has a [bimodal distribution](#), with the mean interval being either 65 or 91 minutes, and is dependent on the length of the prior eruption. Within a margin of error of ± 10 minutes, Old Faithful will erupt either 65 minutes after an eruption lasting less than $2\frac{1}{2}$ minutes, or 91 minutes after an eruption lasting more than $2\frac{1}{2}$ minutes.



Single Gaussian



Mixture of two Gaussians

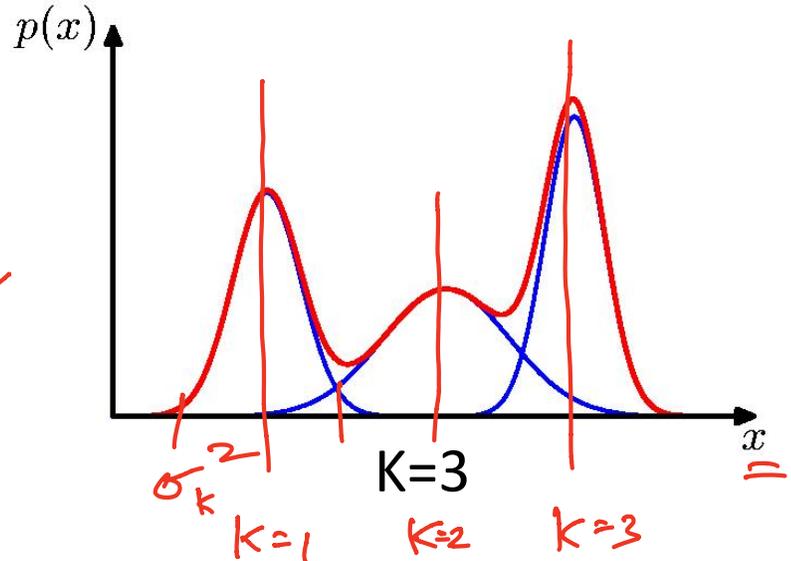
Mixtures of Gaussians (2)

Combine simple models into a complex model:

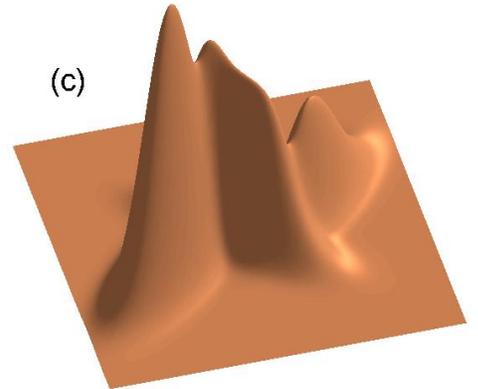
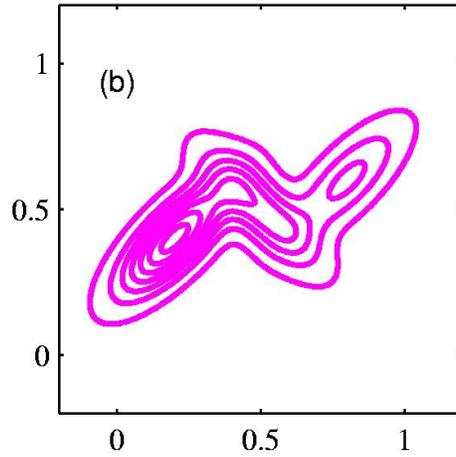
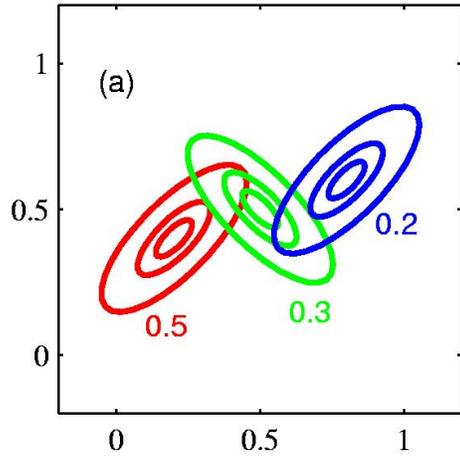
$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Component
Mixing coefficient

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



Mixtures of Gaussians (3)



Mixture of Gaussians

- Mixtures of Gaussians

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

- Expressed with latent variable \mathbf{z}

"marginalization"

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

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

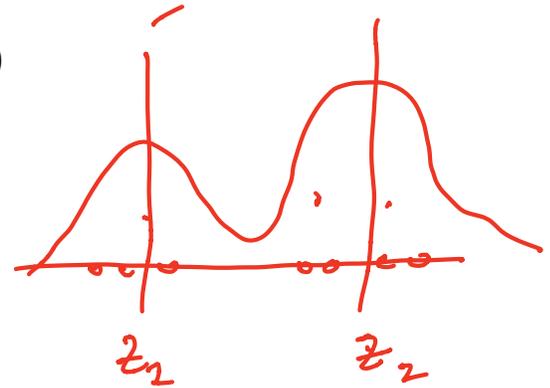
$$p(z_k = 1) = \pi_k //$$

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

$$p(\mathbf{z}) = \prod_k \pi_k^{z_k} \quad z_k = \{0, 1\}$$

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

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



Want to estimate the latent variables for data X

- Probability of data given latent representation

$$p(\mathbf{X} | \underbrace{\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}}_{\boldsymbol{\theta}}) = \prod_N \sum_k \pi_k \mathcal{N}(x | \mu_k, \Sigma_k)$$

- Log likelihood

$$\ln p(\mathbf{X} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_N \ln \sum_k \pi_k \mathcal{N}(x | \mu_k, \Sigma_k)$$
$$= \ln \sum_z p(x, z | \boldsymbol{\theta})$$

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

Can't we just solve for the latent variables by maximizing log likelihood?

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

- Take derivative w.r.t. μ_k : $\rightarrow \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k) = \frac{1}{2\pi^{n/2} |\Sigma_k|^{1/2}}$

$$\frac{\partial}{\partial \mu_k} \mathcal{N}(\cdot) = \Sigma_k^{-1} (\mathbf{x}_n - \mu_k) \mathcal{N}(\cdot) \cdot \exp \left\{ -\frac{1}{2} (\mathbf{x}_n - \mu_k)^T \Sigma_k^{-1} (\mathbf{x}_n - \mu_k) \right\}$$

$$\frac{\partial}{\partial \mu_k} \ln p(\mathbf{X} | \pi, \mu, \Sigma) = \sum_n \underbrace{\pi_k \Sigma_k^{-1} (\mathbf{x}_n - \mu_k) \mathcal{N}(\cdot)}_{\text{responsibility } r_{nk}} \underbrace{\left(\frac{\Sigma_k^{-1} (\mathbf{x}_n - \mu_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\cdot)} \right)}_{\cdot \Sigma_k^{-1} (\mathbf{x}_n - \mu_k)}$$

$$\mu_k = \frac{\sum_N \delta(z_{nk}) \mathbf{x}_n}{\sum_N \delta(z_{nk})} \quad \text{K-means.}$$

Can't we just solve for the latent variables by maximizing log likelihood?

- Log likelihood $\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$.
- Take derivative w.r.t. $\boldsymbol{\mu}_k$:

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

“responsibility”, from Bayes’s rule:

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

$E \gamma(z_{nk})$

using parameter initialization

Solving for μ_k, Σ_k

Take derivative w.r.t. μ_k :

Solving for π_k

Use Lagrange multipliers with constraint

$$\sum_{k=1}^K \pi_k = 1$$

EM Gauss Mix

1. Initialize the means $\underline{\mu}_k$, covariances $\underline{\Sigma}_k$ and mixing coefficients $\underline{\pi}_k$, and evaluate the initial value of the log likelihood.
2. **E step.** Evaluate the responsibilities using the current parameter values

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

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

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

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

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

where

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

4. Evaluate the log likelihood

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

and check for convergence of either the parameters or the log likelihood. If the convergence criterion is not satisfied return to step 2.

calculate class

like K-means

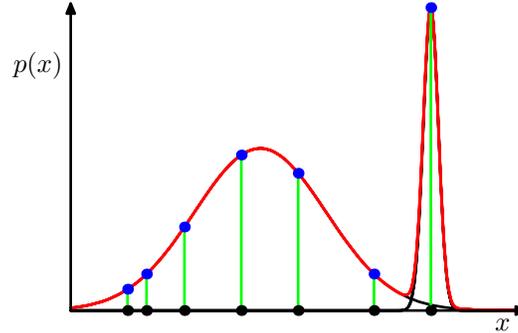
calculate means

heuristic EM

$\log E[\cdot] \geq E \log(\cdot)$
"Jensen's inequality"

Important not to have singularities

maximum likelihood



General EM

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

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

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

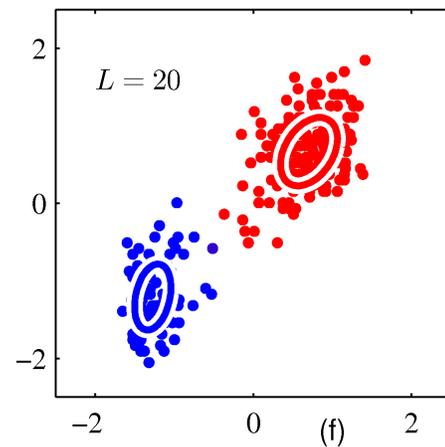
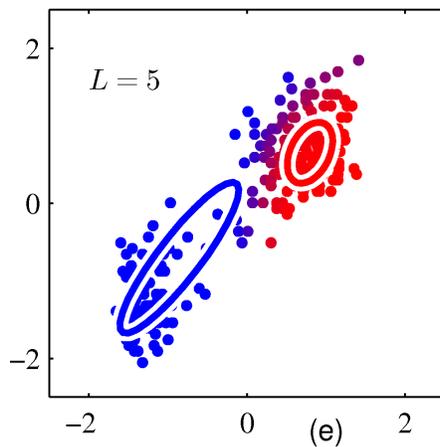
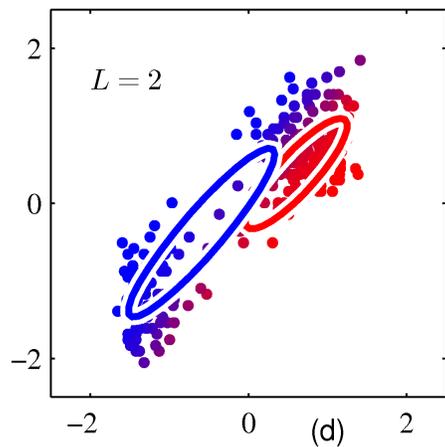
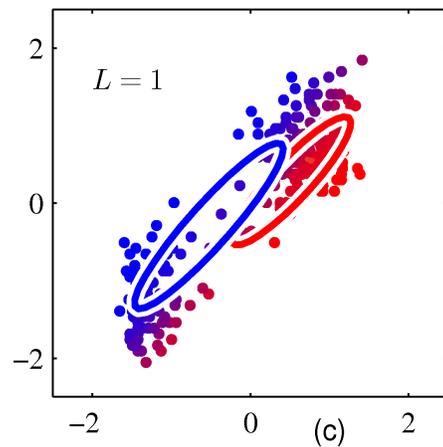
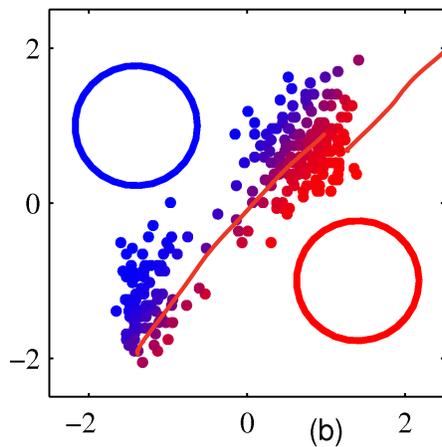
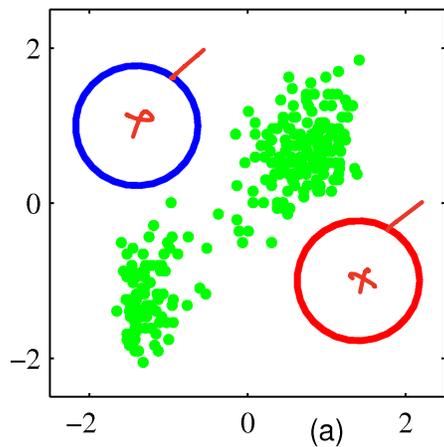
$$\boldsymbol{\theta}^{\text{old}} \leftarrow \boldsymbol{\theta}^{\text{new}} \quad (9.34)$$

and return to step 2.

Gaussian Mixtures

« local method »

$\gamma(z_{kn})$



Kmeans and EM (9.3.2)

K-means:
"hard" assignment /
EM "soft" /

$$\underline{\Sigma_k = \epsilon I.}$$

$$p(\mathbf{x} | \boldsymbol{\mu}_k, \Sigma_k) = \frac{1}{(2\pi\epsilon)^{1/2}} \exp \left\{ -\frac{1}{2\epsilon} \|\mathbf{x} - \boldsymbol{\mu}_k\|^2 \right\}. \quad (9.41)$$

$\exp\left(\frac{1}{2\epsilon} d_{\min}\right)$

Whereby the responsibilities

$$\gamma(z_{nk}) = \frac{\pi_k \exp \left\{ -\frac{\|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2}{2\epsilon} \right\}}{\sum_j \pi_j \exp \left\{ -\frac{\|\mathbf{x}_n - \boldsymbol{\mu}_j\|^2}{2\epsilon} \right\}} \quad (9.42)$$

Handwritten notes: $1/k$ (pointing to π_k), 0 (above $\boldsymbol{\mu}_k$), 0 (below $\boldsymbol{\mu}_j$), $0 \rightarrow 0$ (above denominator), 1 (below denominator), $1/2$ (below 2ϵ)

Becomes delta functions.

$\min \|x_n - \mu_k\|, \gamma = 1$
else 0

And the EM means approach the Kmeans

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \quad (9.17)$$

// //