

$\sqrt{\rho}$  dist

# GMM & EM

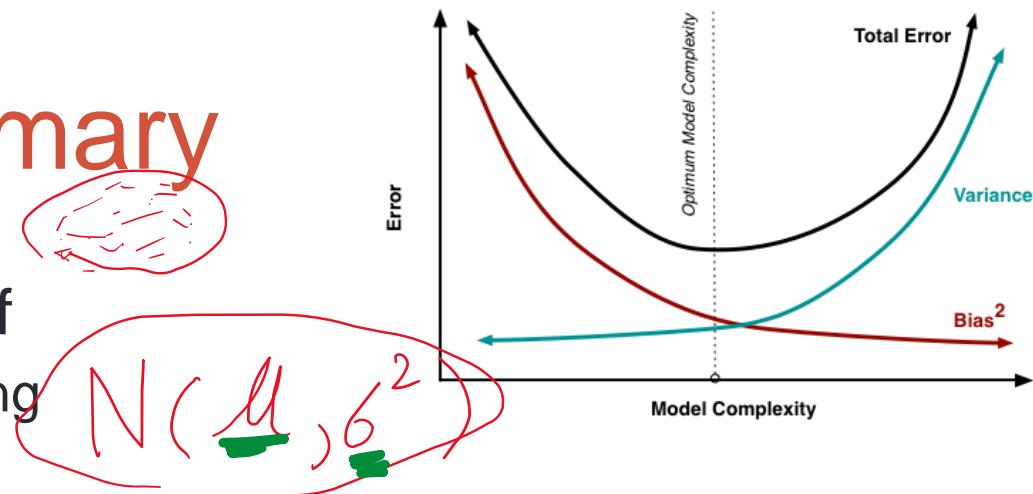
---

Initial parameter

And some RoC

# Last time summary

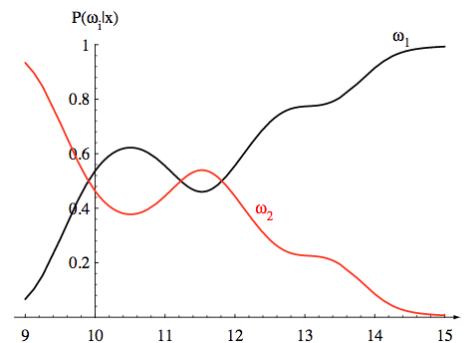
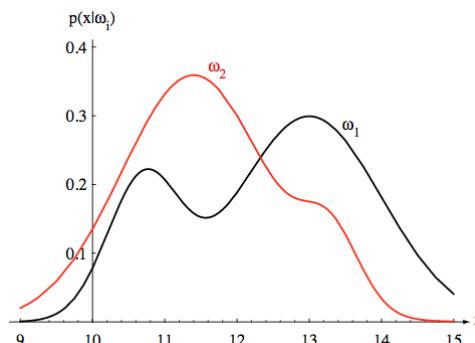
- Bias-Variance trade-off
  - Overfitting and underfitting
- MLE vs MAP estimate
  - How to use the prior
- LRT (Bayes Classifier)
  - Naïve Bayes



Likelihood ratio

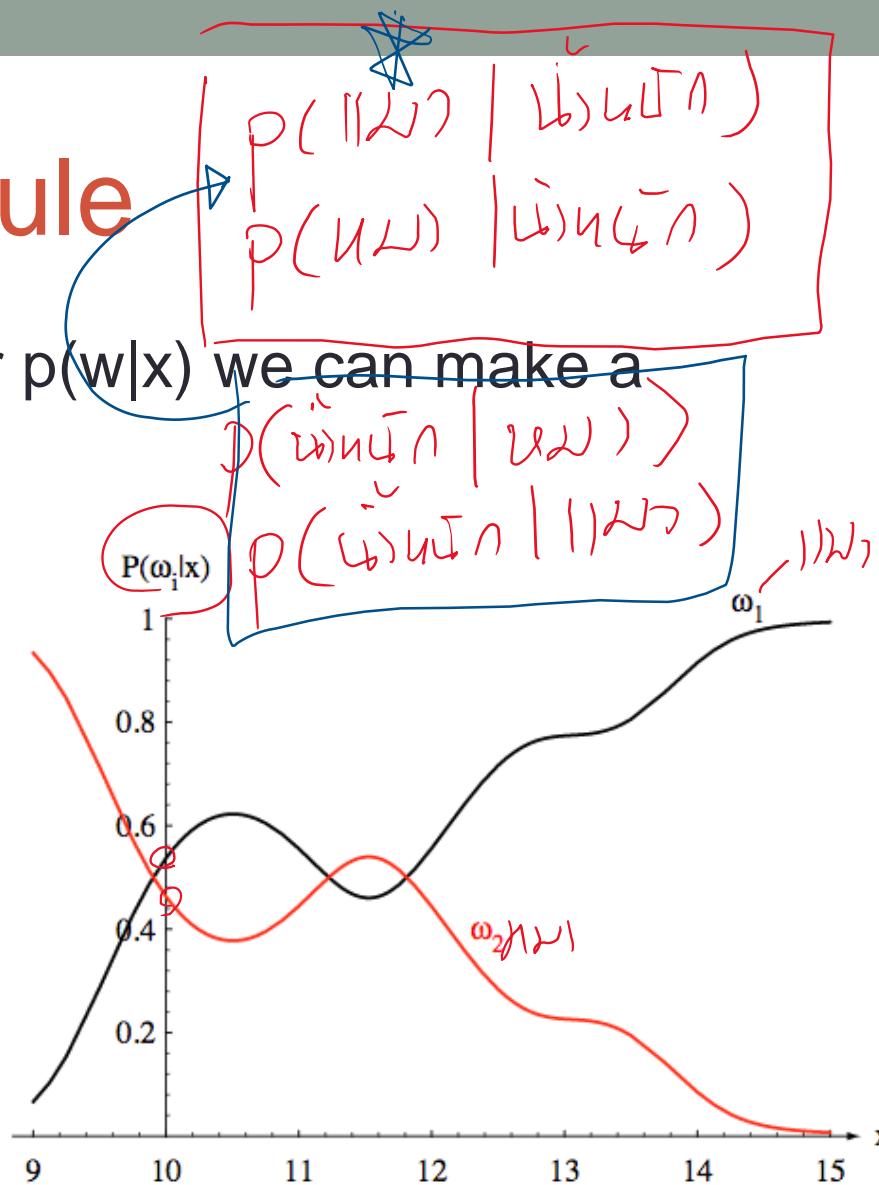
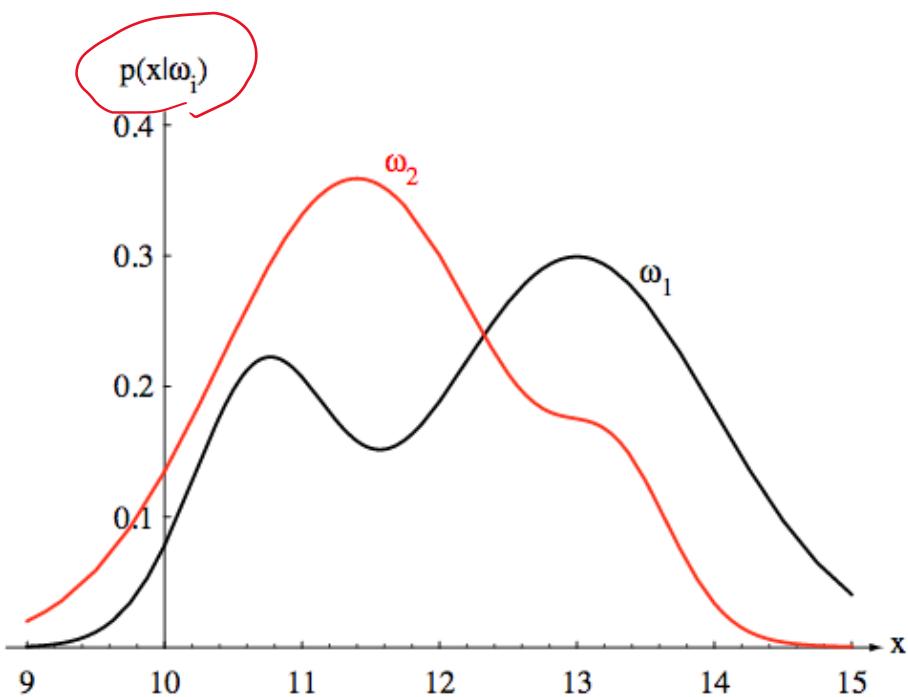
?

Ratio of priors



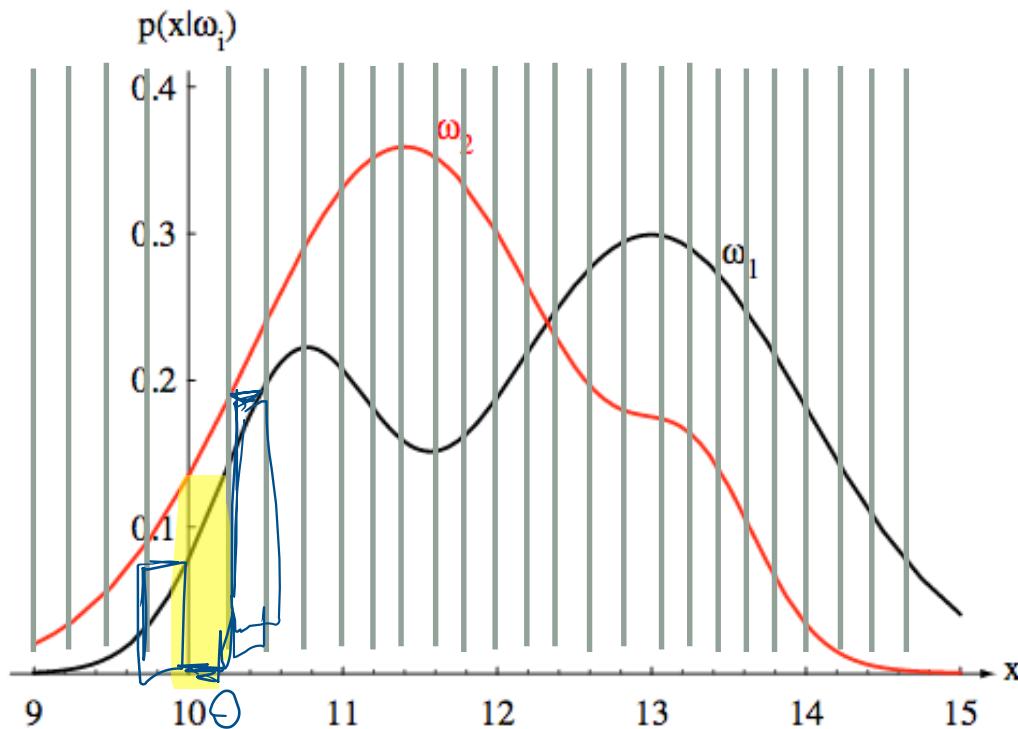
# A simple decision rule

- If we can know either  $p(x|w)$  or  $p(w|x)$  we can make a classification guess

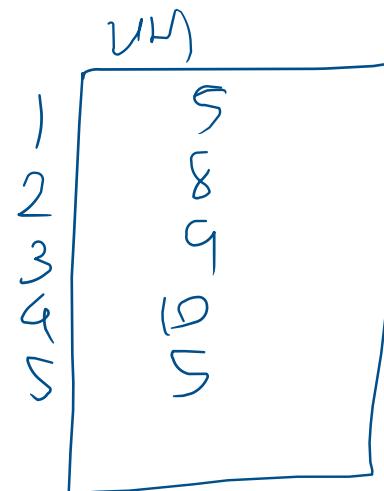


Goal: Find  $p(x|w)$  or  $p(w|x)$  by finding the parameter of the distribution

# A simple way to estimate $p(x|w)$



Make a histogram!



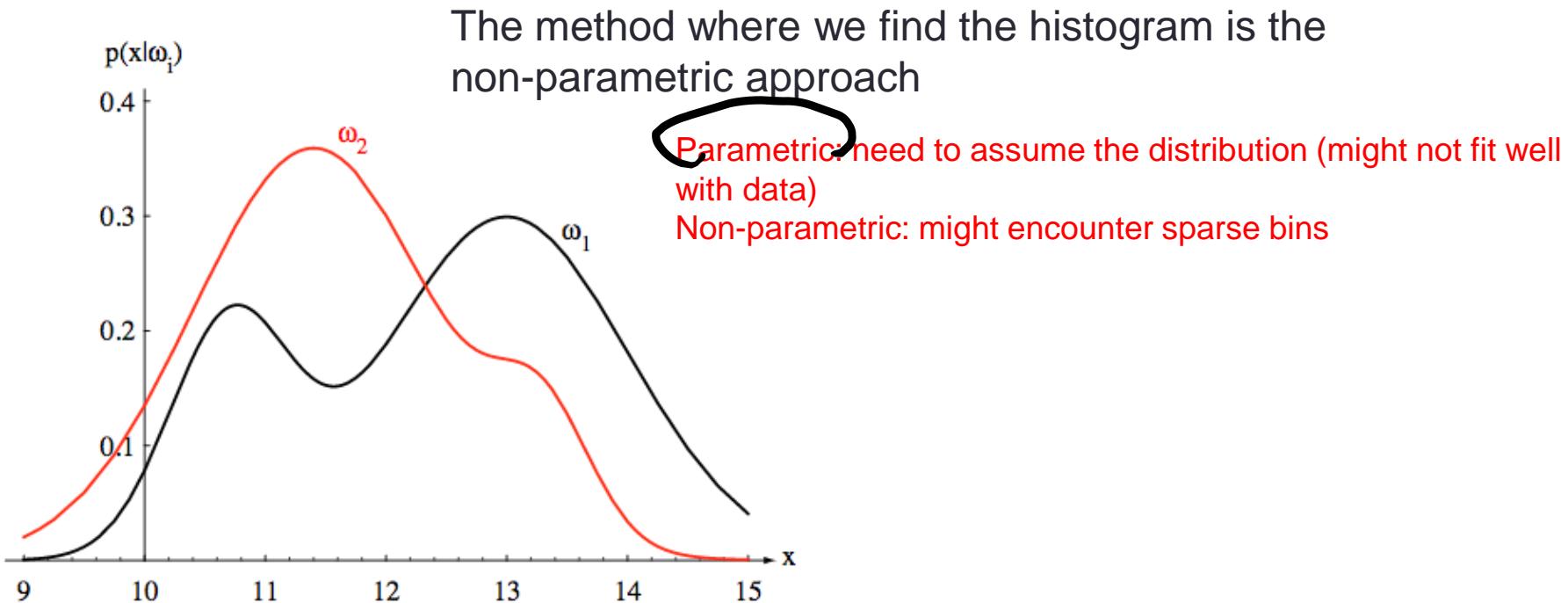
What happens if there is no data in a bin?

# The parametric approach

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/2\sigma^2}$$

MLE MAP

- We assume  $p(x|w)$  or  $p(w|x)$  follow some distributions with parameter  $\theta$



Goal: Find  $\theta$  so that we can estimate  $p(x|w)$  or  $p(w|x)$

# Maximum Likelihood Estimate (MLE)

$$p(x; \theta)$$

$$p(x|w_i)$$

$$p(w_i|x) = \frac{p(x|w_i)p(w_i)}{p(x)}$$

- Maximizing the likelihood (probability of data given model parameters)

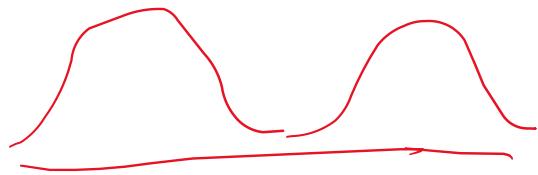
$$\text{Posterior} = \frac{\text{likelihood} * \text{prior}}{\text{evidence}}$$

$$\underline{p(x|\theta)} = L(\theta) \leftarrow \text{This assumes the data is fixed}$$

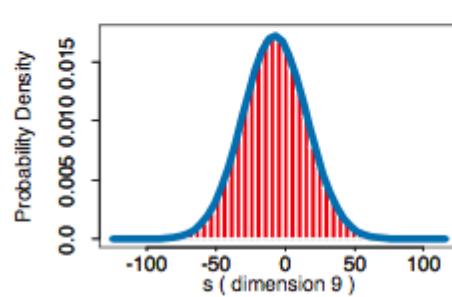
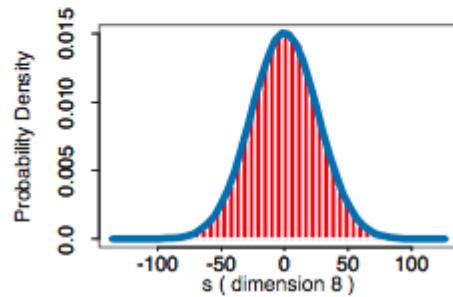
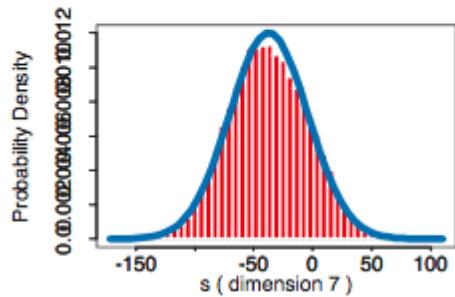
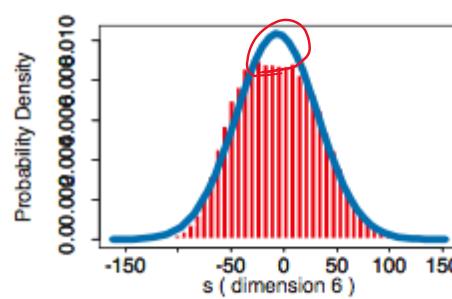
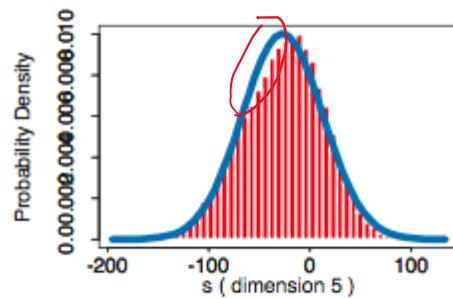
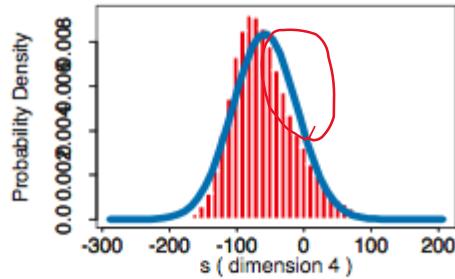
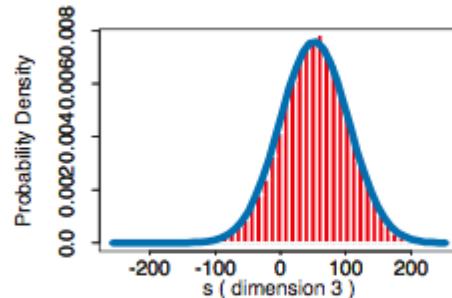
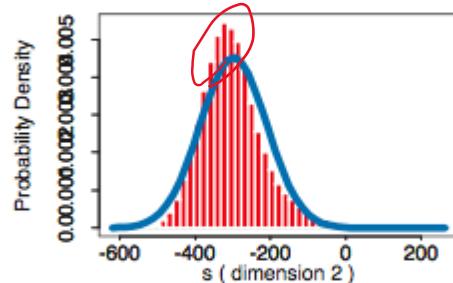
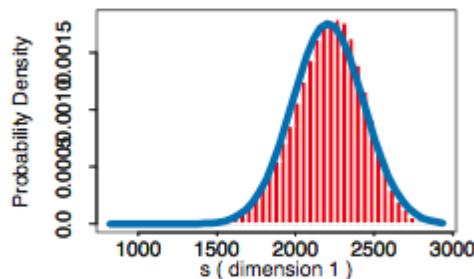
- Usually done on log likelihood

- Take the partial derivative wrt to  $\theta$  and solve for the  $\theta$  that maximizes the likelihood

# Model of one Gaussian

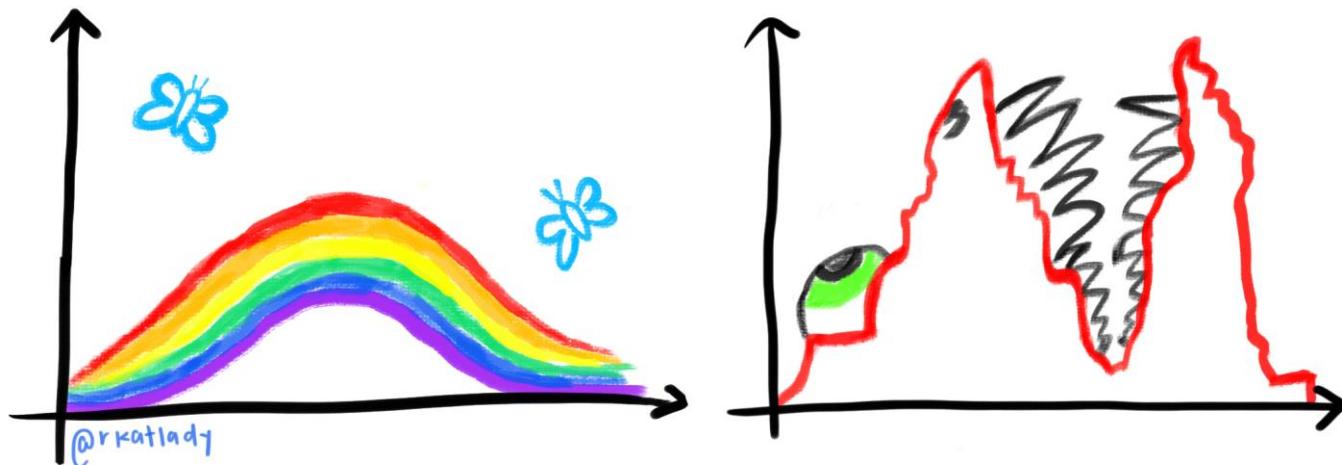


First 9 MFCC's from [s]: Gaussian PDF



# UNDERLYING DISTRIBUTIONS:

~~PARAMETRIC ASSUMPTIONS~~ VS. REALITY

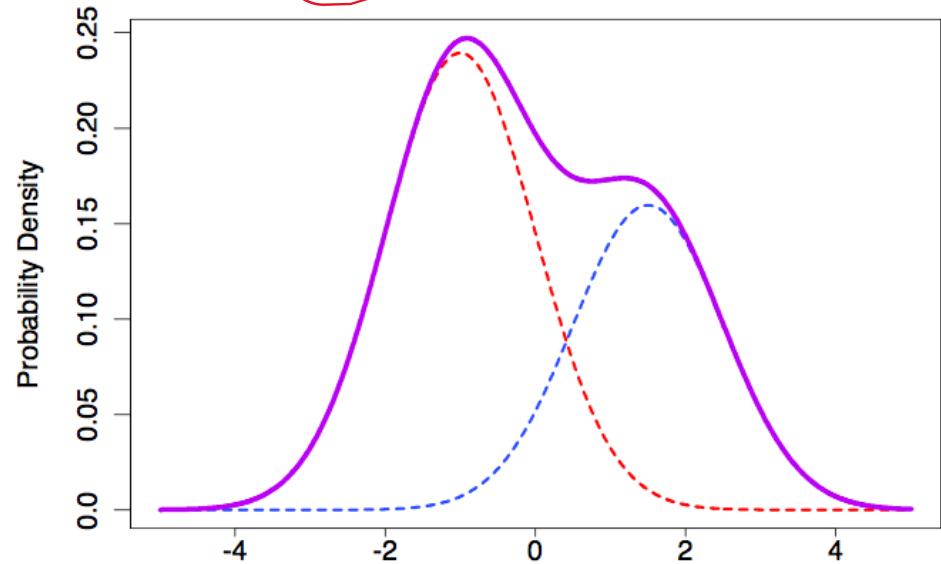


# Gaussian Mixture Models (GMMs)

- Gaussians cannot handle multi-modal data well
- Consider a class can be further divided into additional factors
- Mixing weight makes sure the overall probability sums to 1

$$P(x) \sim \sum_{k=1}^K w_k N(\mu_k, \sigma_k)$$

*K = 3 or 1*



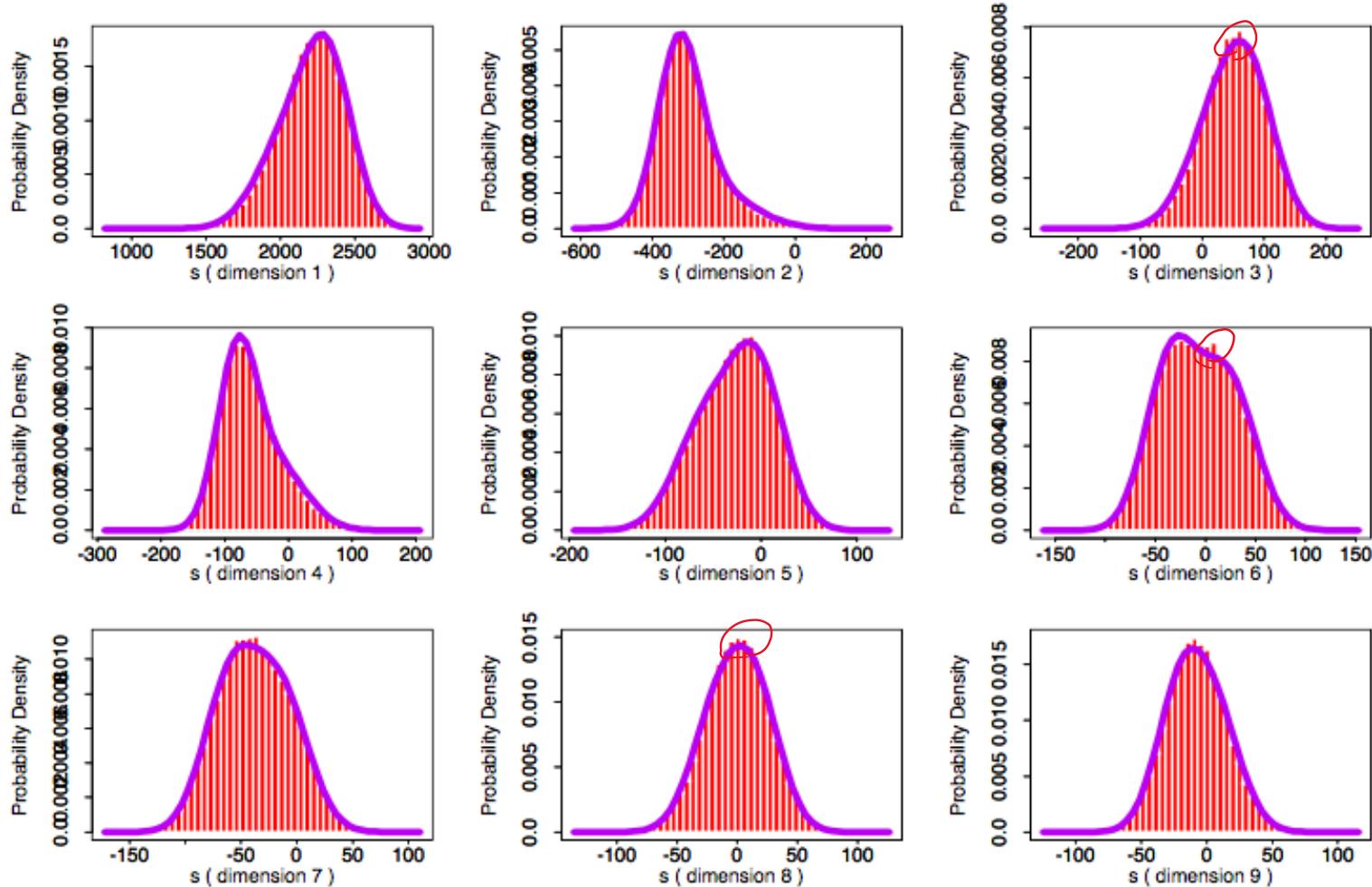
$$p(x) = 0.6 p_1(x) + 0.4 p_2(x)$$

$p_1(x) \sim N(-\sigma, \sigma^2)$        $p_2(x) \sim N(1.5\sigma, \sigma^2)$

$\frac{x}{\sigma}$

# Mixture of two Gaussians

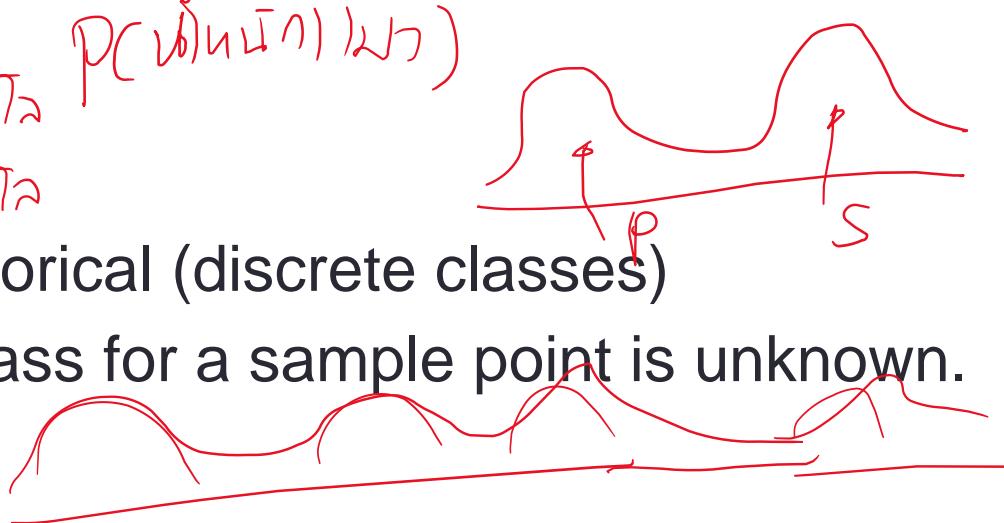
[s]: 2 Gaussian Mixture Components/Dimension



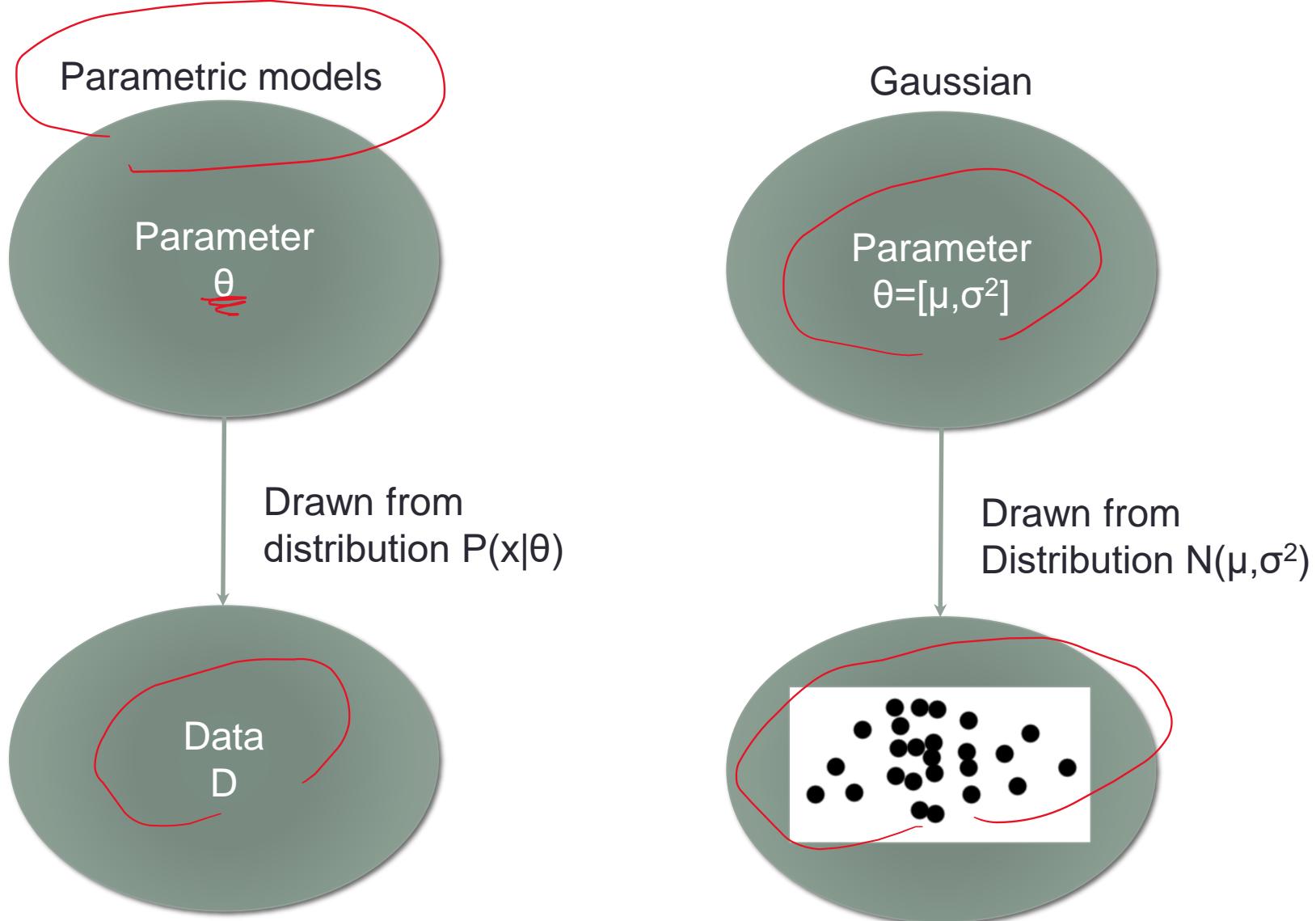
# Mixture models

$$p(x) = \sum_k p(k) p_k(x)$$

- A mixture of models from the same distributions (but with different parameters)
- Different mixtures can come from different sub-class
  - Cat class
    - Siamese cats  $\frac{1}{2}$
    - Persian cats  $\frac{1}{3}$
- $p(k)$  is usually categorical (discrete classes)
- Usually the exact class for a sample point is unknown.
  - Latent variable



# Parametric models



# Maximum A Posteriori (MAP) Estimate

## MLE

- Maximizing the likelihood (probability of data given model parameters)

$$\underset{\theta}{\operatorname{argmax}} p(\mathbf{x}|\theta)$$

$$p(\mathbf{x}|\theta) \\ = L(\theta)$$

- Usually done on log likelihood

- Take the partial derivative wrt to  $\theta$  and solve for the  $\theta$  that maximizes the likelihood

## MAP

- Maximizing the posterior (model parameters given data)

$$\underset{\theta}{\operatorname{argmax}} p(\theta|\mathbf{x})$$

- But we don't know  $p(\theta|\mathbf{x})$

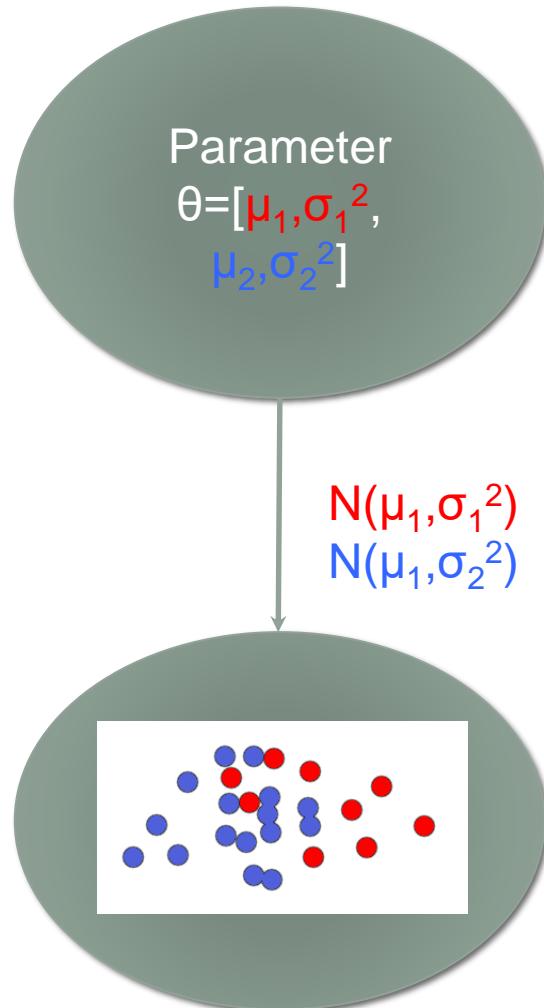
- Use Bayes rule  
$$p(\theta|\mathbf{x}) = \frac{p(\mathbf{x}|\theta)p(\theta)}{p(\mathbf{x})}$$

- Taking the argmax for  $\theta$  we can ignore  $p(\mathbf{x})$

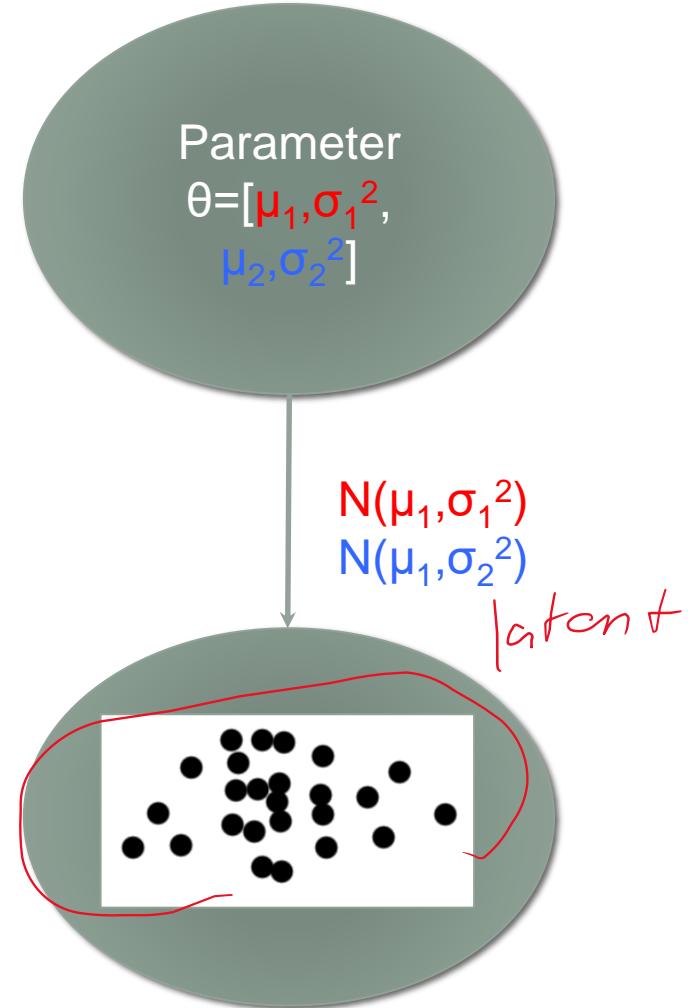
- $$\underset{\theta}{\operatorname{argmax}} p(\mathbf{x}|\theta) p(\theta)$$

# What if some data is missing?

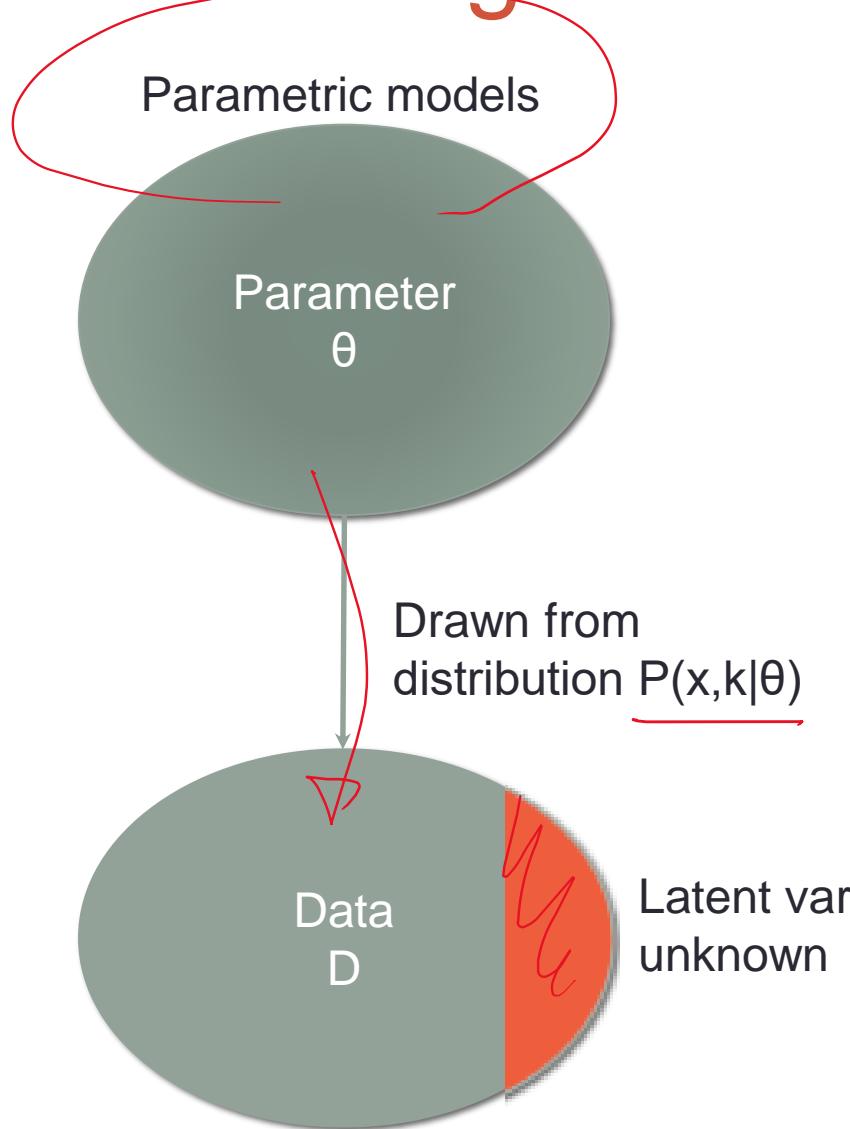
Mixture of Gaussian



Unknown mixture labels



# Estimating missing data



Need to estimate both the latent Variables and the model parameters.

# Slight difference in notation

$p(\mathbf{x}|\theta)$

P

vs  $p(\mathbf{x};\theta)$

$\theta$  as a RV at a fixed value

vs  $\theta$  as a fixed parameter

Most of the time can be used interchangeably

$$P(X | \underline{k; \theta}) = Z \times \underline{k \theta}$$

# Estimating latent variables and model parameters

- GMM
- Observed  $(x_1, x_2, \dots, x_N)$
- Latent  $(k_1, k_2, \dots, k_N)$  from K possible mixtures
- Parameter for  $p(k)$  is  $\phi$ ,  $p(k = 1) = \phi_1$ ,  $p(k = 2) = \phi_2, \dots$

*log likelihood*

$$l(\phi, \mu, \Sigma) = \sum_{n=1}^N \log p(x^{(i)}; \phi, \mu, \sigma)$$

$$= \sum_{n=1}^N \log \left( \sum_{l=1}^K p(x_n | k_{n,l}; \mu, \sigma) p(k_{n,l}; \phi) \right)$$

Make things hard to solve

Cannot be solved by differentiating

~~PS~~

# Assuming k

నేడు అసా వివరచు  $\phi_1$  1, 2, 3, 4

- What if we somehow know  $k_n$ ?
- Maximizing wrt to  $\phi, \mu, \sigma$  gives

$$\phi_j = \frac{1}{N} \sum_{n=1}^N 1(k_n = j) \quad \phi_1 \geq 1 + 0+$$

$$\mu_j = \frac{\sum_{n=1}^N 1(k_n = j) x_n}{\sum_{n=1}^N 1(k_n = j)}$$

$$\sigma_j^2 = \frac{\sum_{n=1}^N 1(k_n = j) (x_n - \mu_j)^2}{\sum_{n=1}^N 1(k_n = j)}$$

1(*condition*)

Indicator function. Equals one if condition is met. Zero otherwise

# Iterative algorithm

- Initialize  $\phi, \mu, \sigma$
- Repeat till convergence
  - Expectation step (E-step) : Estimate the latent labels  $k$
  - Maximization step (M-step) : Estimate the parameters  $\phi, \mu, \sigma$  given the latent labels
- Called Expectation Maximization (EM) Algorithm
- How to estimate the latent labels?  
*ML likelihood*

# Iterative algorithm

- Initialize  $\phi, \mu, \sigma$
- Repeat till convergence
  - Expectation step (E-step) : Estimate the latent labels  $k$  by finding the pdf of  $k$  given everything else  $p(k|x; \phi, \mu, \sigma)$
  - Maximization step (M-step) : Estimate the parameters  $\phi, \mu, \sigma$  given the latent labels by maximizing the expectation of the log likelihood
- Extension of MLE for latent variables
  - MLE :  $\text{argmax } \log p(x;\theta)$
  - EM :  $\text{argmax } \log \sum_k p(x, k;\theta)$   $k$  unknown

How to evaluate  $\log \sum_k p(x, k;\theta)$  when we don't know  $k$ ?

# Convex functions and Jensen's inequality

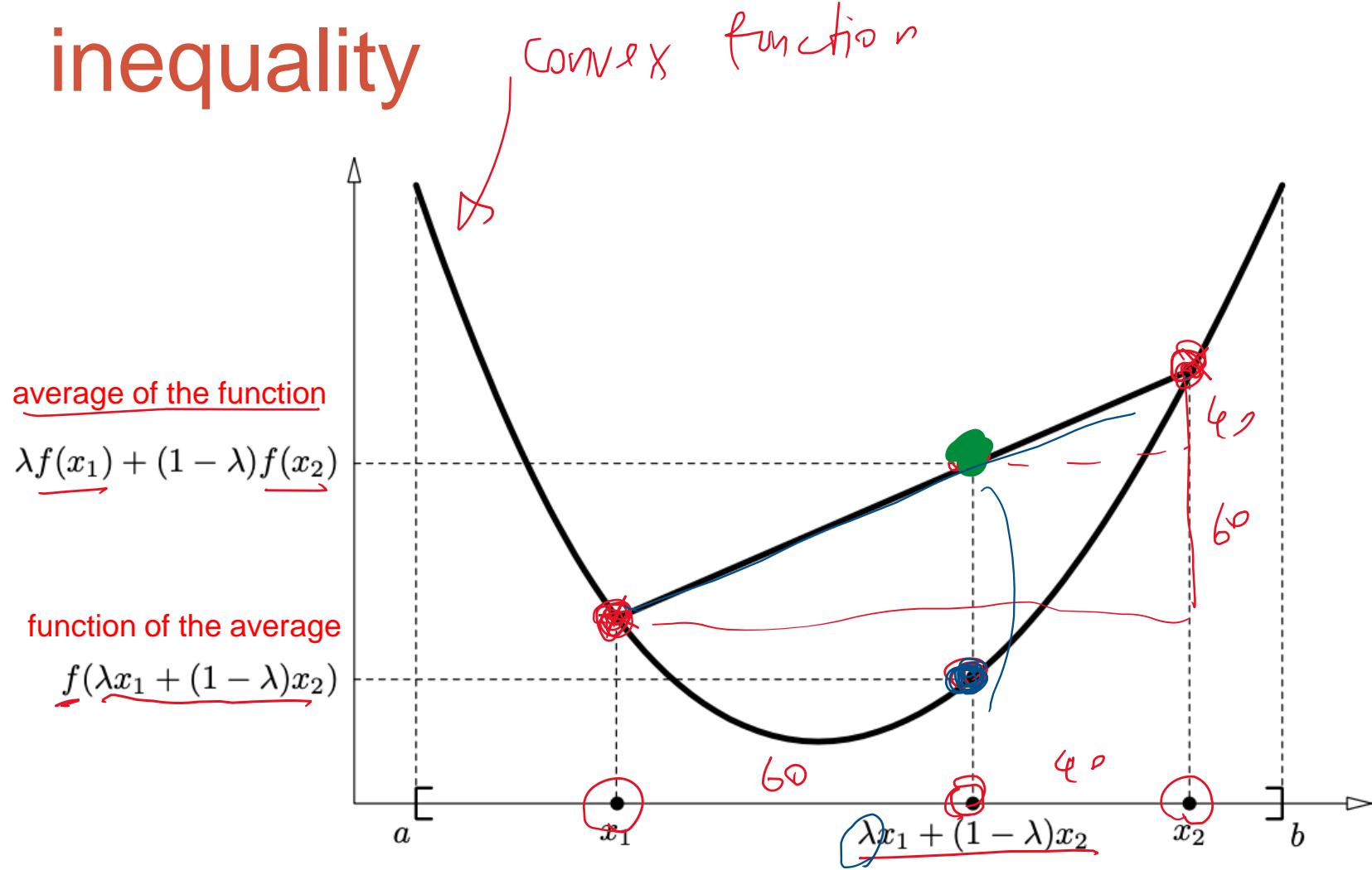


Figure 1:  $f$  is convex on  $[a, b]$  if  $f(\lambda x_1 + (1 - \lambda) x_2) \leq \lambda f(x_1) + (1 - \lambda) f(x_2)$   $\forall x_1, x_2 \in [a, b], \lambda \in [0, 1]$ .

# Jensen's inequality

Let f be a convex function on interval I

If  $x_1, x_2, \dots, x_n$  is in I,  $p(x)$   
 $w_1, \dots, w_n > 0$  and sums to 1  
then,

$$f\left(\sum_i^n w_i x_i\right) \leq \sum_i^n w_i f(x_i)$$

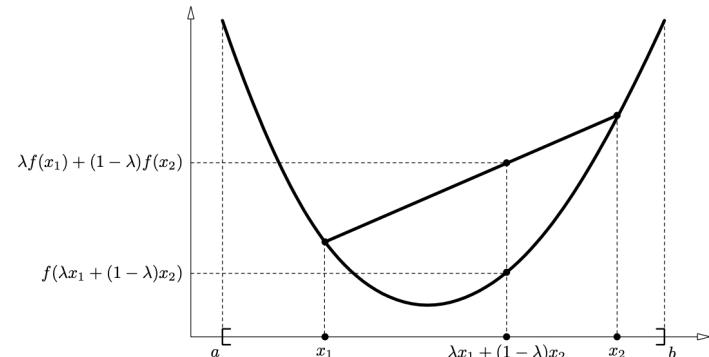
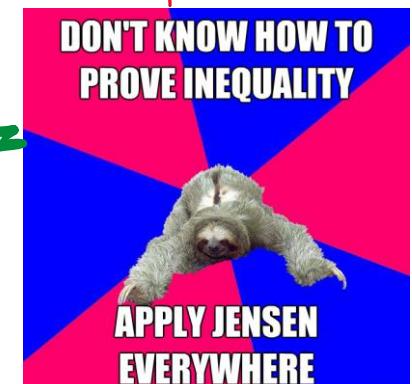


Figure 1:  $f$  is convex on  $[a, b]$  if  $f(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda f(x_1) + (1 - \lambda)f(x_2)$   $\forall x_1, x_2 \in [a, b], \lambda \in [0, 1]$ .

$$E[x] = \sum x p(x)$$
  
$$E[f(x)] =$$



If  $f$  is concave, flip the inequality.  
Can view this as expectation

$$\sum f(x) p(x)$$

$$f(E[X]) \leq E[f(X)]$$

# Jensen's inequality and ELBO



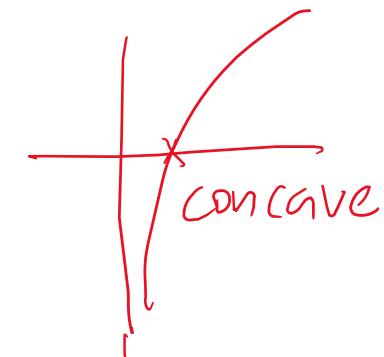
$$\log \sum_k p(x, k | \theta)$$

Let  $Q(k)$  be some dist.  
over  $k$

$$f\left(\sum_i^n w_i x_i\right) \leq \sum_i^n w_i f(x_i)$$

$$= \log \sum_k \left[ p(x, k; \theta) \frac{Q(k)}{Q(k)} \right]$$

$$= \log \sum_k \left[ Q(k) \left[ \frac{p(x, k; \theta)}{Q(k)} \right] \right]$$



$$\geq \sum_k Q(k) \log \frac{P(x, k; \theta)}{Q(k)}$$

$$\triangleleft E_{k \sim Q} \left[ \log \frac{P(x, k; \theta)}{Q(k)} \right]$$

Maximize Evidence Lower Bound (ELBO) =  $\sum_k Q(k) \log (p(x, k; \theta)/Q(k))$

# Making the lower bound tight

We will make the bound tight for fixed  $\theta$   
Jensen's inequality is tight when?

$$f\left(\sum_i^n w_i x_i\right) \leq \sum_i^n w_i f(x_i)$$

$$f(E[X]) \leq E[f(X)]$$

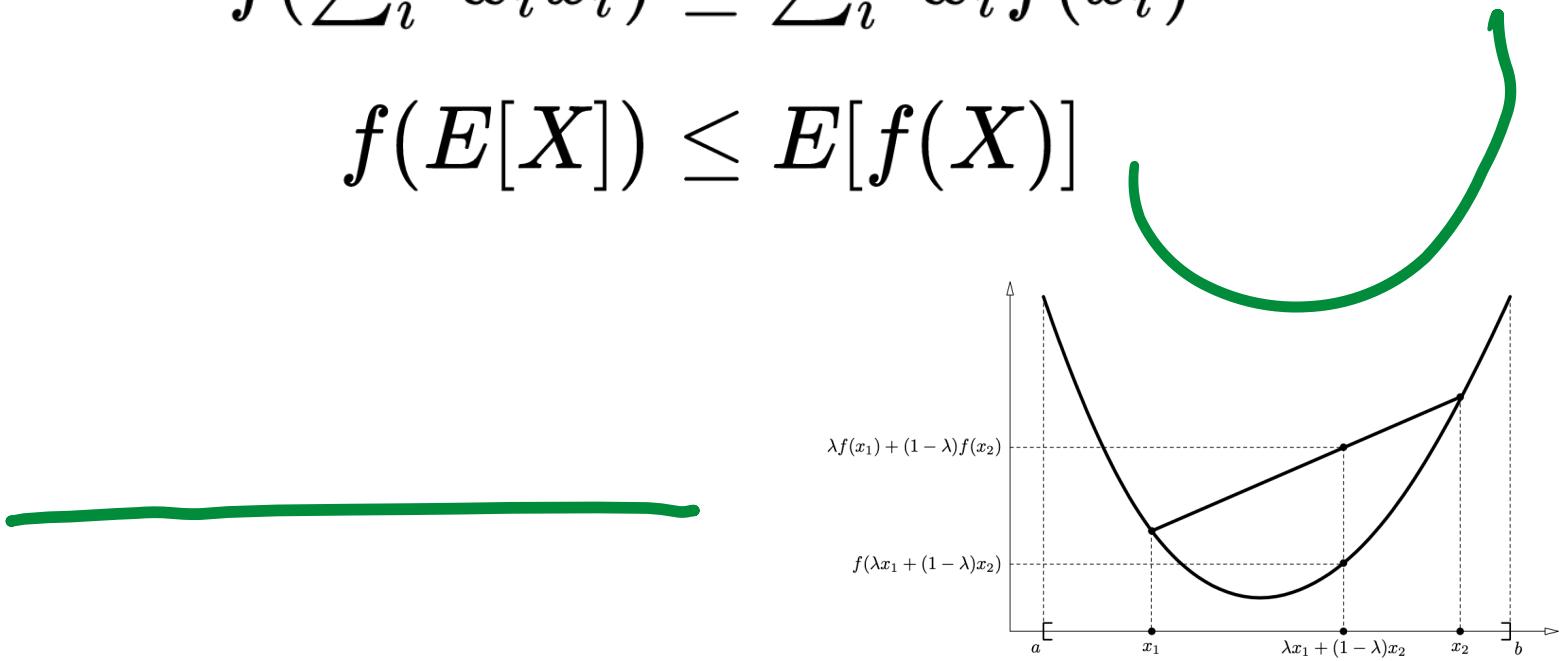


Figure 1:  $f$  is convex on  $[a, b]$  if  $f(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda f(x_1) + (1 - \lambda)f(x_2)$   $\forall x_1, x_2 \in [a, b], \lambda \in [0, 1]$ .

# Making the lower bound tight

We will make the bound tight for fixed  $\theta$

$$f\left(\sum_i^n w_i x_i\right) \leq \sum_i^n w_i f(x_i)$$

$$f(E[X]) \leq E[f(X)]$$

If  $f( )$  is strictly convex, Jensen's inequality is tight ~~iff~~ IFF  
 $x_i$  are all equal

$$E[X] = X = \text{constant}$$

# Making the lower bound tight

We will make the bound tight for fixed  $\theta$

Jensen's inequality is tight when the inside of the expectation is a constant,  $c$  wrt the expectation

$$p(x, k; \theta) / Q(k) = c \quad (\text{constant wrt to } k)$$

Let's pick  $Q(k) \propto p(x, k; \theta)$

We know that  $\sum_k Q(k) = 1$

$$\begin{aligned} Q(k) &= \frac{p(x, k; \theta)}{\sum_k p(x, k; \theta)} \\ &\Rightarrow \frac{p(x, k; \theta)}{p(x; \theta)} = p(k | x; \theta) \end{aligned}$$

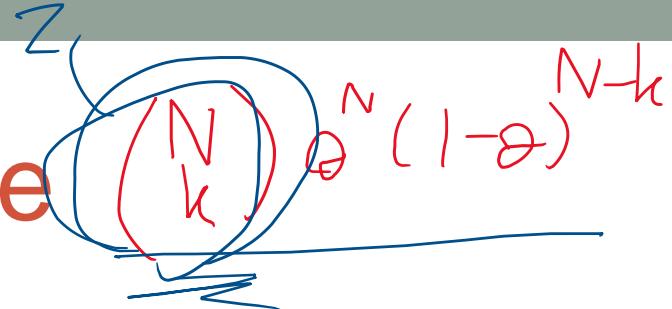
or  $Q(k) = p(k | x; \theta)$

$$\begin{aligned} \sum_k \gamma p(x, k; \theta) &= 1 \\ \gamma &= \frac{1}{\sum_k p(x, k; \theta)} \end{aligned}$$

# Iterative algorithm (general)

- Goal of EM :  $\log \sum_k p(x, k; \theta) \geq \sum_k Q(k) \log (p(x, k; \theta)/Q(k))$
- Maximize the ELBO instead
- Initialize  $\Theta$
- Repeat till convergence
  - Expectation step (E-step) : estimate the conditional expectation  $Q(k) = p(k|x; \theta)$  using the current  $\theta$ .
  - Maximization step (M-step) : Estimate new  $\Theta$  given by maximizing the ELBO given current  $Q(k)$

# EM on a simple example



- Grades in class  $P(A) = 0.5$   $P(B) = 0.5 - \theta$   $P(C) = \theta$
- We want to estimate  $\theta$  from three known numbers  $N_a, N_b, N_c$
- Find the maximum likelihood estimate of  $\theta$

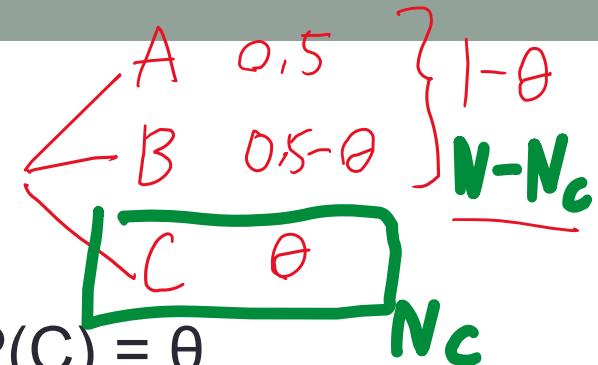
$$P(D; \underline{\theta}) = Z \cdot 0.5^{N_a} (0.5 - \theta)^{N_b} \theta^{N_c}$$

$$\log P(D; \theta) = \log Z + N_a \log 0.5 + N_b \log(0.5 - \theta) + N_c \log \theta$$

$$\frac{\partial L}{\partial \theta} = 0 + 0 \cdot N_b \frac{1}{0.5 - \theta} (-1) + N_c \frac{1}{\theta} = 0$$

$$\Rightarrow \theta = \frac{1}{2} \frac{N_c}{N_c + N_B}$$

# EM on a simple example



- Grades in class  $P(A) = 0.5$   $P(B) = 0.5 - \theta$   $P(C) = \theta$
- We want to estimate  $\theta$  from ONE known number
  - $N_c$  (we also know  $N$  the total number of students)
- Find  $\theta$  using EM

E-Step

$$Q(N_b) \sim \text{Binomial}(N - N_c, \frac{0.5 - \theta}{1 - \theta})$$

$$= \binom{N - N_c}{N_b} \left( \frac{0.5 - \theta}{1 - \theta} \right)^{N_b} \left( \frac{0.5}{1 - \theta} \right)^{1 - N_b}$$

$$Q(N_a) \sim \text{Binomial}(N - N_c, \frac{0.5}{1 - \theta})$$

M-Step

$$\sum_{N_B} Q(N_B; \theta') \log \frac{P(N_c, N_B | N; \theta)}{Q(N_B; \theta')}$$

$$\Rightarrow \sum_{N_B} Q(N_B; \theta') \left[ \log P(N_c, N_B | N; \theta) - \log Q(N_B; \theta') \right]$$

$$\frac{\partial \ell}{\partial \theta} = \sum_{N_B} Q(N_B; \theta') \frac{P'(N_c, N_B | N; \theta)}{P(N_c, N_B | N; \theta)} - 0$$

$$P(N_c, N_B | N; \theta) = Z_{0.5}^{N-N_B-N_c} (0.5-\theta)^{N_B} \theta^{N_c}$$

$$P'(N_c, N_B | N; \theta) = Z_{0.5}^{N-N_B-N_c} \left[ (0.5-\theta)^{N_B} N_c \theta^{N_c-1} - \theta^N N_B (0.5-\theta)^{N_B-1} \right]$$

$$\frac{P'}{P} = \frac{N_c}{Q} - \frac{N_B}{0.5-\theta}$$

$$\sum_{N_B} Q(N_B; \underline{\theta}') \left[ \frac{N_c}{\theta} - \frac{N_B}{0.5 - \theta} \right] = Q$$

$$\frac{\sum Q(N_B; \underline{\theta}') \frac{N_c}{\theta}}{N_B} = \sum_{N_B} Q(N_B; \underline{\theta}') \frac{N_B}{0.5 - \theta}$$

$$\frac{N_c}{\theta} \underbrace{\sum Q(N_B; \underline{\theta}')}_{N_B}$$

$$\frac{N_c}{\theta}$$

$$\Rightarrow \theta = \frac{1}{2} \frac{N_c}{N_c + N_b}$$

$$\frac{\sum Q(N_B; \underline{\theta}') N_B}{N_B} = \frac{(N - N_c) Q_{0.5 - \underline{\theta}'}}{1 - \underline{\theta}'} \boxed{\frac{N_B \sim Q(N_B; \underline{\theta})}{1 - \underline{\theta}'}}$$

$\hat{N}_b$   
estimate von  $N_b$

# Will this work?

For iteration  $i$ , with  $\theta^{(i)}$

$$\log \sum_k p(x, k; \theta^{(i)}) \geq \underbrace{\sum_k Q(k) \log (p(x, k; \theta^{(i)}) / Q(k))}_{\text{ELBO}}$$

E-step, making the bound tight by picking  $Q'(k)$  yields

$$\log \sum_k p(x, k; \theta^{(i)}) = \sum_k Q'(k) \log (p(x, k; \theta^{(i)}) / Q'(k))$$

M-step, maximize ELBO by finding  $\theta^{(i+1)}$

$$\sum_k Q'(k) \log (p(x, k; \theta^{(i)}) / Q'(k)) \leq \sum_k Q'(k) \log (p(x, k; \theta^{(i+1)}) / Q'(k))$$

For iteration  $i+1$ , with  $\theta^{(i+1)}$

$$\log \sum_k p(x, k; \theta^{(i+1)}) \geq \sum_k Q(k) \log (p(x, k; \theta^{(i+1)}) / Q(k))$$

Thus,

$$\log \sum_k p(x, k; \theta^{(i+1)}) \geq \log \sum_k p(x, k; \theta^{(i)})$$

So EM improves the likelihood at every step!

# Notes on ELBO

We set  $Q(k) = p(k | x; \theta)$  to make the inequality tight.

What if we cannot compute  $p(k | x; \theta)$  ?

Use a looser bound by picking any  $Q(k)$

Estimate  $p(k | x; \theta)$  with  $q(k | x; \theta)$  that we can compute

This is called **Variational Inference**

We will revisit this.

15'00

# Estimating latent variables and model parameters

- GMM  $p(x) = \sum_k p(k)N(\mu_k, \sigma_k)$
- Observed  $(x_1, x_2, \dots, x_N)$
- Latent  $(k_1, k_2, \dots, k_N)$  from K possible mixtures
- Parameter for  $p(k)$  is  $\phi$ ,  $p(k = 1) = \phi_1$ ,  $p(k = 2) = \phi_2 \dots$

$$l(\phi, \mu, \Sigma) = \sum_{n=1}^N \log p(x^{(i)}; \phi, \mu, \sigma)$$

$$= \sum_{n=1}^N \log \sum_{l=1}^K p(x_n | k_{n,l}; \mu, \sigma) p(k_{n,l}; \phi)$$

Make things hard to solve

Cannot be solved by differentiating

# EM on GMM

- E-step
  - Set soft labels:  $w_{n,j}$  = probability that nth sample comes from jth mixture p
  - Using Bayes rule
    - $p(k|x ; \mu, \sigma, \phi) = \frac{p(x|k ; \mu, \sigma, \phi) p(k; \mu, \sigma, \phi)}{p(x; \mu, \sigma, \phi)}$
    - $p(k|x ; \mu, \sigma, \phi)$  is proportional to  $p(x|k ; \mu, \sigma, \phi) p(k; \phi)$

$$p(k_n = j | x_n; \phi, \mu, \Sigma) = \frac{p(x_n; \mu_j, \sigma_j) p(k_n = j; \phi)}{\sum_l p(x_n; \mu_l, \sigma_l) p(k_n = l; \phi)}$$

Annotations:

- $w_{n,j}$  is circled in red.
- $p(k_n = j | x_n; \phi, \mu, \Sigma)$  is circled in red.
- $N(\mu_j, \sigma_j)$  has a red arrow pointing to it from the term  $p(x_n; \mu_j, \sigma_j)$ .
- $\phi_j$  has a red arrow pointing to it from the term  $p(k_n = j; \phi)$ .

# EM on GMM

ເຕີກວ່າວິນ

- M-step (hard labels)

$$\phi_j = \frac{1}{N} \sum_{n=1}^N \underbrace{1(k_n = j)}_{\text{hard labels}}$$

$$\mu_j = \frac{\sum_{n=1}^N 1(k_n = j) x_n}{\sum_{n=1}^N 1(k_n = j)}$$

$$\sigma_j^2 = \frac{\sum_{n=1}^N 1(k_n = j) (x_n - \mu_j)^2}{\sum_{n=1}^N 1(k_n = j)}$$

# EM on GMM

- M-step (soft labels)

Summe:

$$\phi_j = \frac{1}{N} \sum_{n=1}^N w_{n,j}$$

$$\mu_j = \frac{\sum_{n=1}^N w_{n,j} \underline{x_n}}{\sum_{n=1}^N w_{n,j}}$$

0,3

0,5

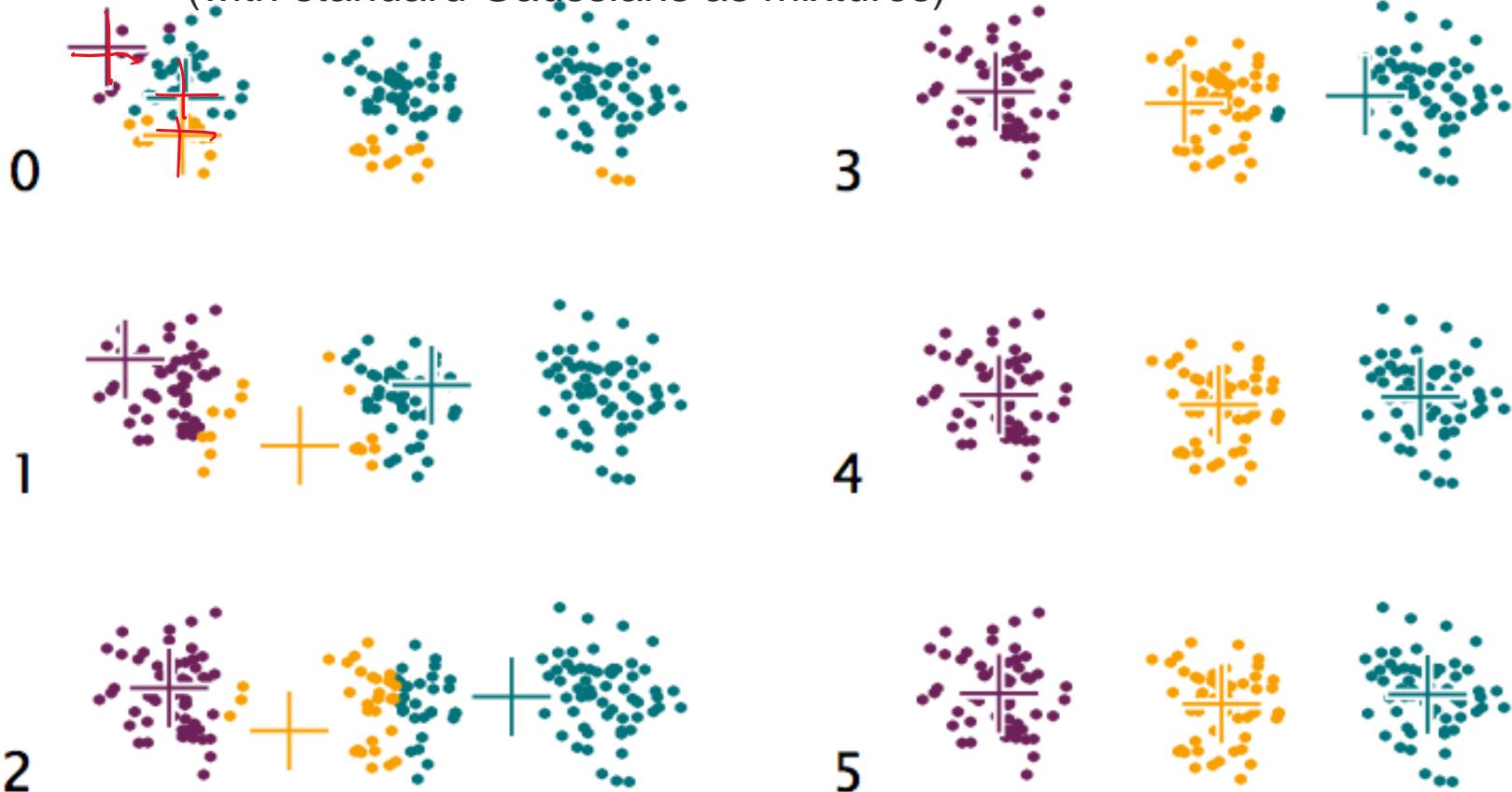
0,25

$$\sigma_j^2 = \frac{\sum_{n=1}^N w_{n,j} (x_n - \mu_j)^2}{\sum_{n=1}^N w_{n,j}}$$

# K-mean vs EM

*assign*

EM on GMM can be considered as EM with soft labels  
(with standard Gaussians as mixtures)



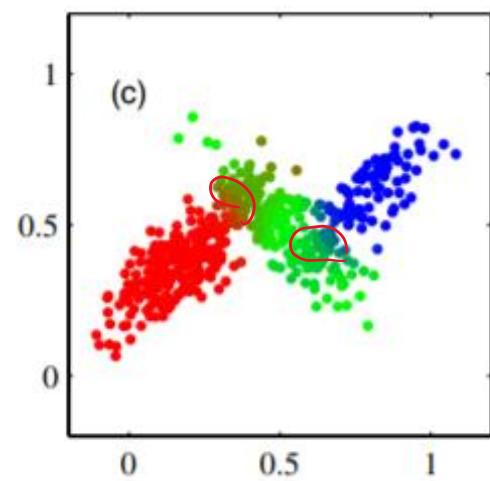
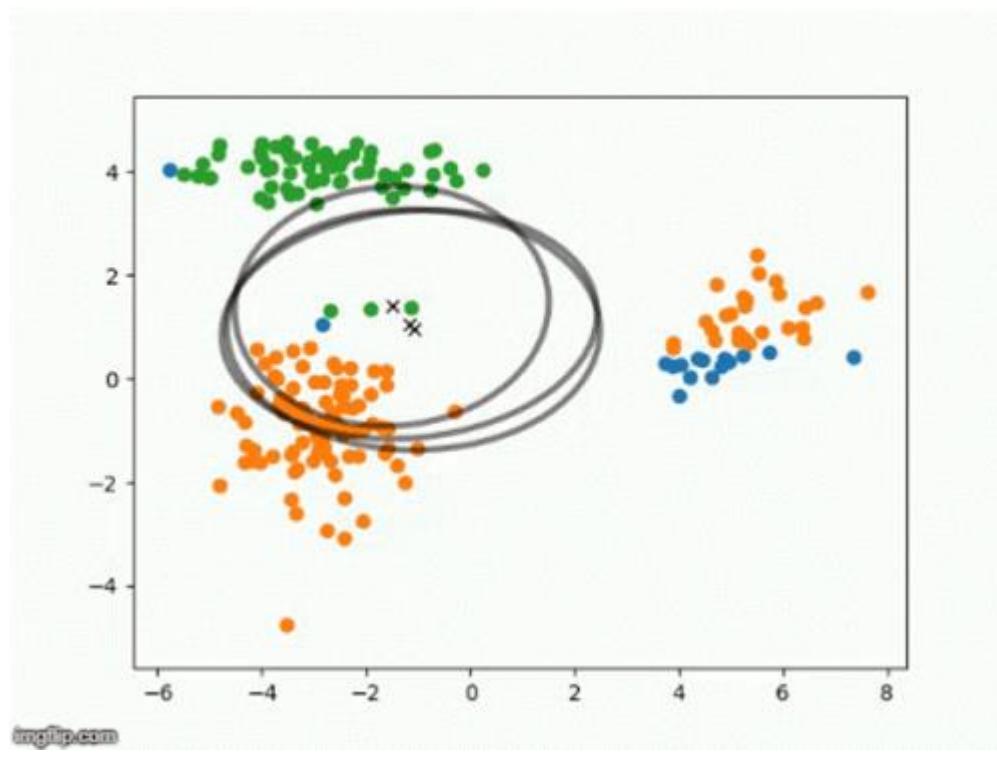
# K-mean clustering

- Task: cluster data into groups
- K-mean algorithm
  - Initialization: Pick K data points as cluster centers
  - Assign: Assign data points to the closest centers
  - Update: Re-compute cluster center
  - Repeat: Assign and Update

# EM algorithm for GMM

- Task: cluster data into Gaussians
- EM algorithm
  - Initialization: Randomly initialize parameters Gaussians
  - Expectation: Assign data points to the closest Gaussians
  - Maximization: Re-compute Gaussians parameters according to assigned data points
  - **Repeat**: Expectation and Maximization
- Note: assigning data points is actually a soft assignment (with probability)

# K-mean vs EM



<https://towardsdatascience.com/gaussian-mixture-models-vs-k-means-which-one-to-choose-62f2736025f0>

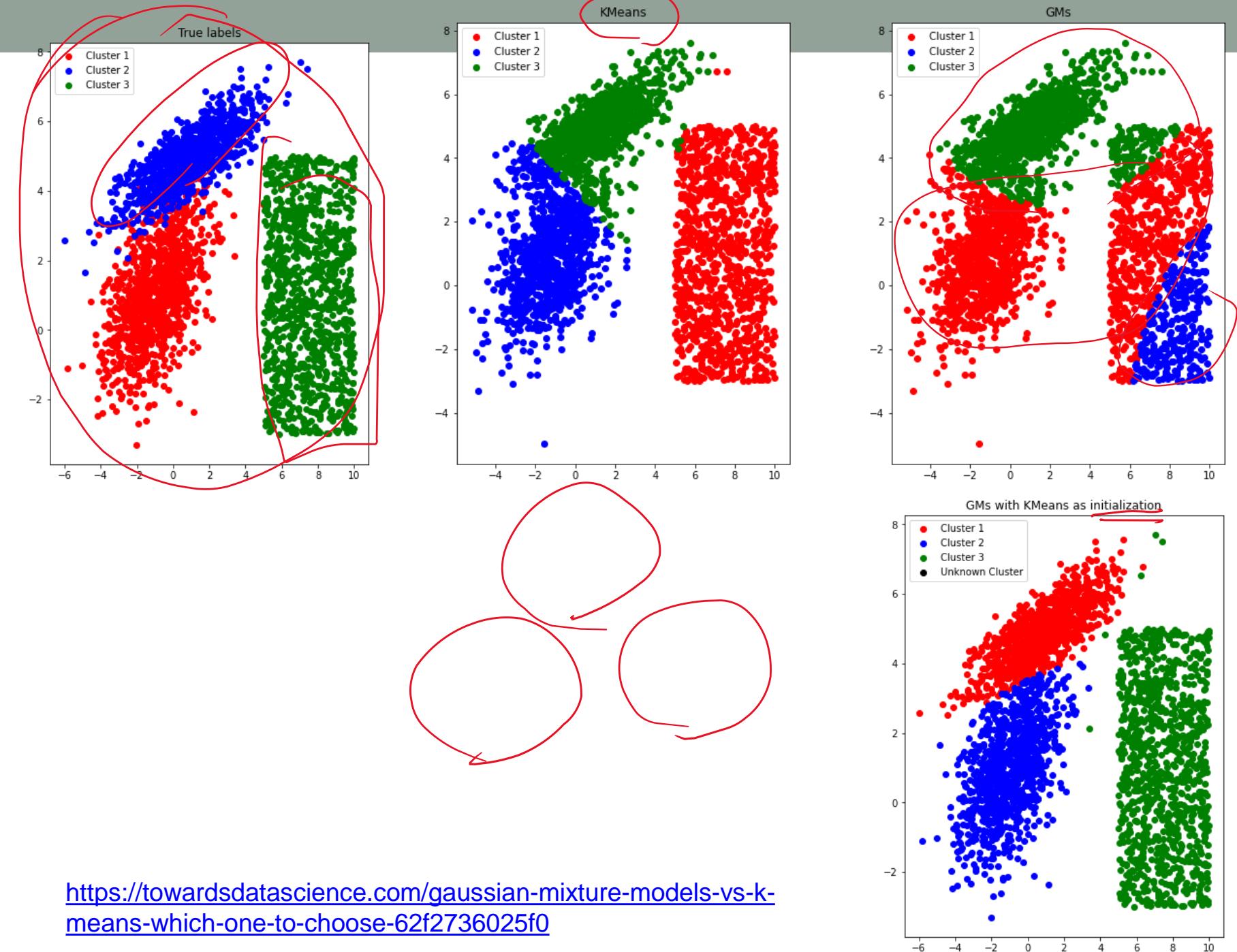
# EM/GMM notes

0, 05

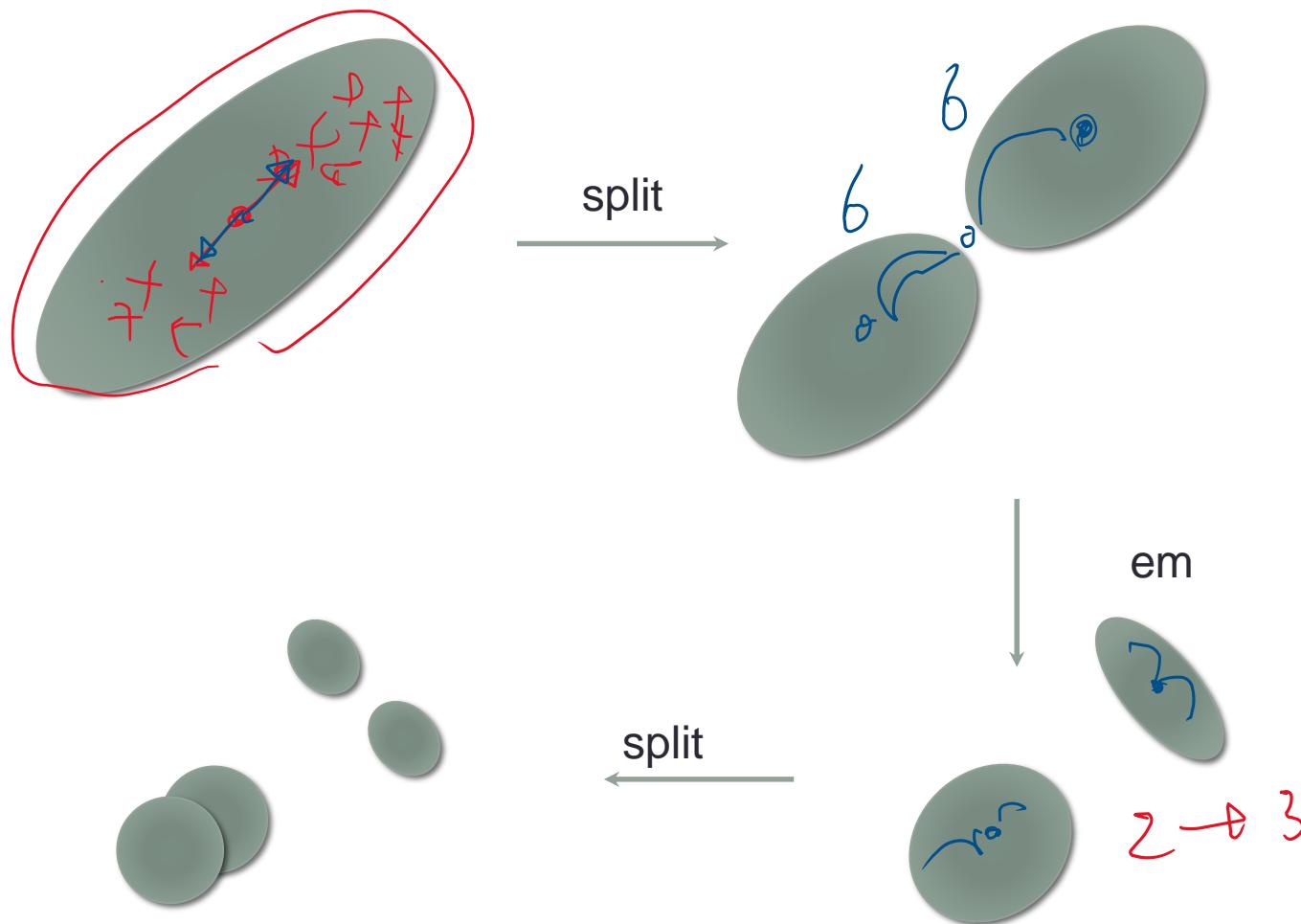
0, 07

k  
10<sup>-4</sup>? 10<sup>-6</sup>

- Converges to local maxima (maximizing likelihood)
  - Just like k-means, need to try different initialization points
- EM always improve the likelihood for each iteration
  - Stops EM when likelihood changes < threshold 10e-6
- Just like k-means some centroid can get stuck with one sample point and no longer moves
  - For EM on GMM this cause variance to go to 0...
    - Introduce variance floor (minimum variance a Gaussian can have)
- Tricks to avoid bad local maxima
  - Starts with 1 Gaussian
  - Split the Gaussians according to the direction of maximum variance
  - Repeat until arrive at k Gaussians
  - Does not guarantee global maxima but works well in practice

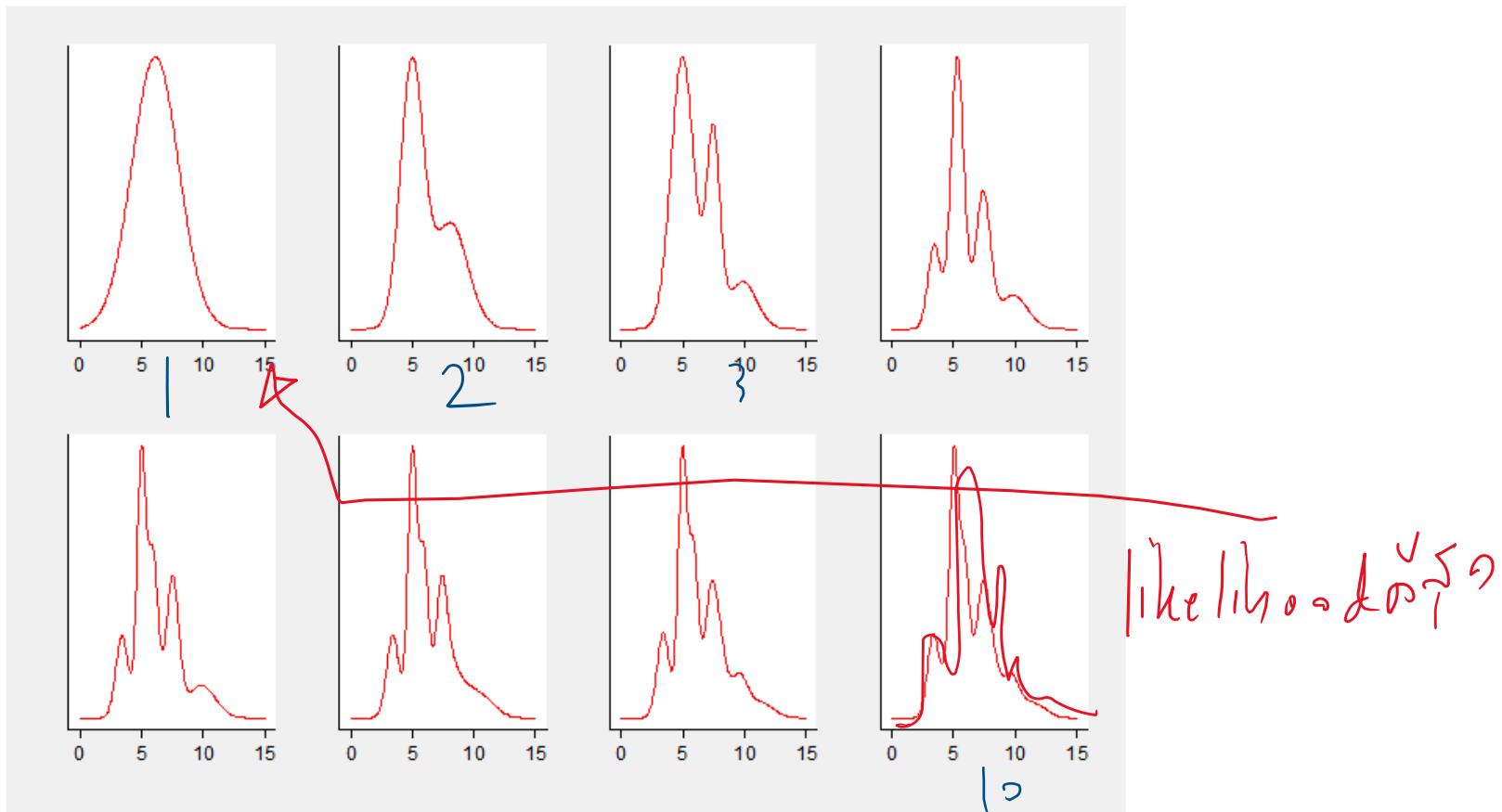


# Gaussian splitting



# Picking the amount of Gaussians

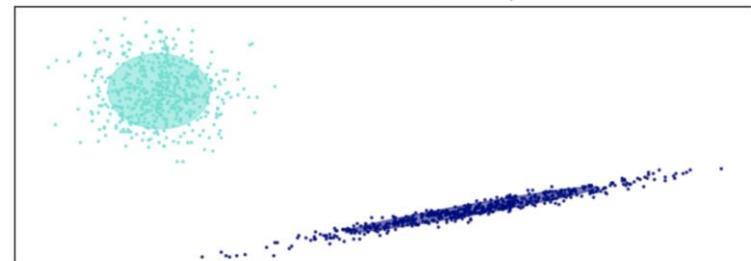
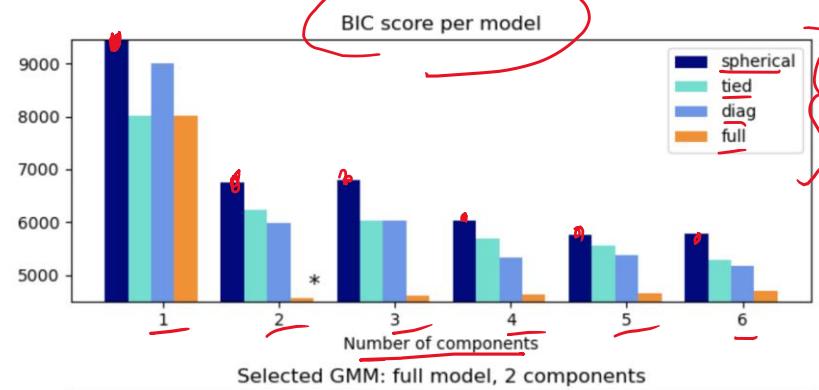
- As we increase K, the likelihood will keep increasing
- More mixtures -> more parameters -> overfits



# Picking the amount of Gaussians

*Tied*

- Need a measure of goodness (like Elbow method in k-mean)
- Bayesian Information Criterion (BIC)
- Penalize the log likelihood from the data by the number of parameters in the model
  - $-2 \log L + t \log (n)$
  - $t$  = number of parameters in the model
  - $n$  = number of data points
- We want to minimize BIC



# BIC is bad use cross validation!

- Just like how I don't recommend using elbow method for clustering
- BIC is bad use cross validation!
- Test on the goal of your model

$$p(w_i | x)$$

# Latent variables?

EM is all about problem formulation. You can solve the same task with different formulations.

## Latent variable considerations

- Imaginary quantity meant to provide a simplified view of the process
  - GMM mixtures. Speech recognizer states. Customer segmentation.
- Real-world thing, but impossible to directly measure
  - Cause of a disease. Temperature of a star.
- Real-world thing, that is not measured because of noise/faulty sensors

$$y = x + \text{noise}$$

*latent*

# Latent variables?

- Discrete latent variables: clusters/partitions data into  
subgroups
- Continuous latent variables: can be used for  
dimensionality reduction (factor analysis, etc)

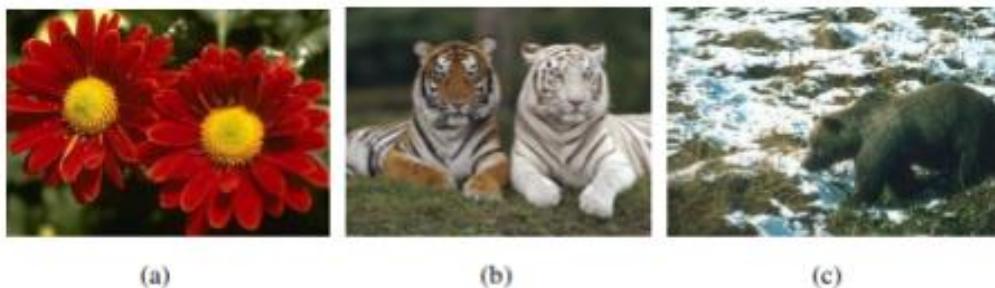
# EM usage examples

# Image segmentation with GMM EM

- $D - \{r,g,b\}$  value at each pixel
- Latent : segment where each pixel comes from
- Hyperparameters: number of mixtures ( $K$ ), initial values



# Image segmentation with GMM EM



**Fig. 1.** Original images: (a) flower, (b) tiger, (c) bear

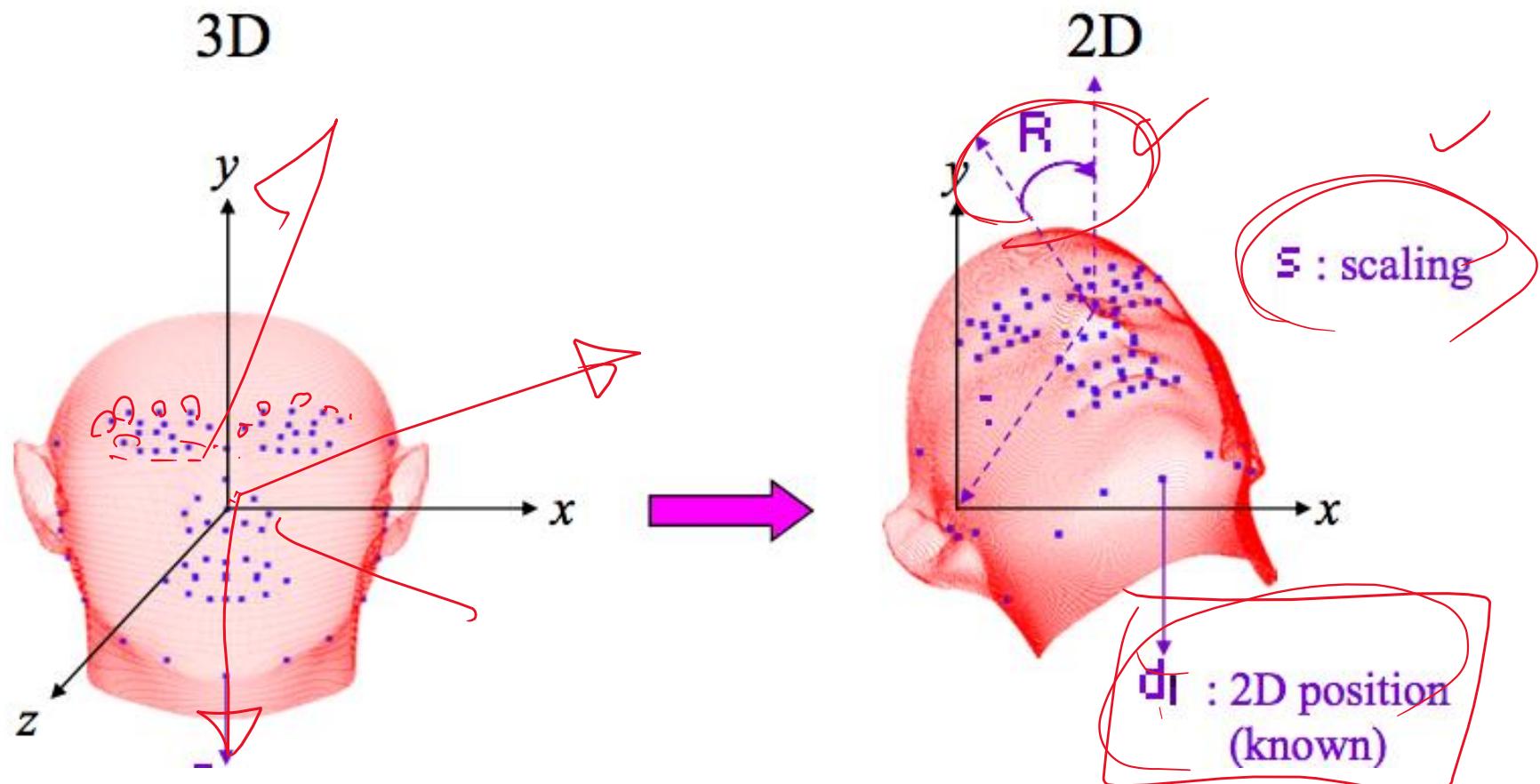


**Fig. 2.** Segmentation results ( $M = 2$ )



**Fig. 3.** Segmentation results ( $M = 5$ )

# Face pose estimation (estimate 3d coordinates from 2d picture)



# Language modeling

## THE UNITED STATES CONSTITUTION

We the People of the United States, in Order to form a more perfect Union, establish Justice, insure domestic Tranquility, provide for the common Defense, promote the general Welfare, and ensure the Blessings of Liberty to ourselves and our Posterity, do ordain and establish this Constitution for the United States of America.

### ARTICLE I.

#### Section 1.

All legislative Powers herein granted shall be vested in a Congress of the United States, which shall consist of a Senate and House of Representatives.

#### Section 2.

Chusec 1: The House of Representatives shall be composed of Members chosen every second Year by the People of the several States, and the Electors in each State shall have the Qualifications requisite for Electors of the most numerous Branch of the State Legislature.

Chusec 2: No Person shall be a Representative who shall not have attained to the Age of twenty-five Years, and been seven Years a Citizen of the United States, and who shall not, when elected, be an Inhabitant of that State in which he shall be chosen.

Chusec 3: Representations and direct Taxes shall be apportioned among the several States which may be included within that Union, according to their respective Numbers, which shall be determined by adding to the whole Number of free Persons, including those bound to Service for a Term of Years, and excluding Indians not taxed, three-fifths of all other Persons. The actual Enumeration shall be made within three Years after the First Meeting of the Convention of the United

Latent variable:  
Topic  
 $P(\text{word}|\text{topic})$

For examples: see Probabilistic latent semantic analysis

# MEME

Know  
Your  
Meme

## Multiple EM for Motif Elicitation

From Wikipedia, the free encyclopedia

ATCG

For other uses, see [MEME \(disambiguation\)](#).

**Multiple Expectation maximizations for Motif Elicitation (MEME)** is a tool for discovering motifs in a group of related [DNA](#) or [protein](#) sequences.<sup>[1]</sup>

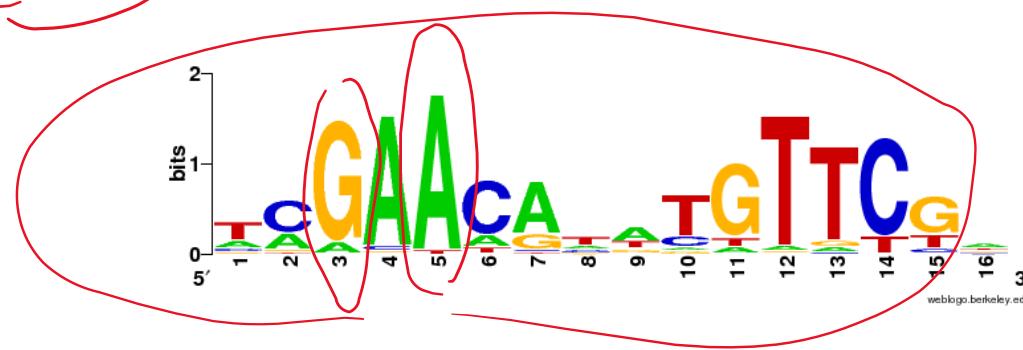
A **motif** is a sequence pattern that occurs repeatedly in a group of related protein or DNA sequences and is often associated with some biological function. MEME represents motifs as [position-dependent letter-probability matrices](#) which describe the probability of each possible letter at each position in the pattern. Individual MEME motifs do not contain gaps. Patterns with variable-length gaps are split by MEME into two or more separate motifs.

MEME takes as input a group of DNA or protein sequences (the training set) and outputs as many motifs as requested. It uses statistical modeling techniques to automatically choose the best width, number of occurrences, and description for each motif.

MEME is the first of a collection of tools for analyzing motifs called the [MEME suite](#).

### Contents [hide]

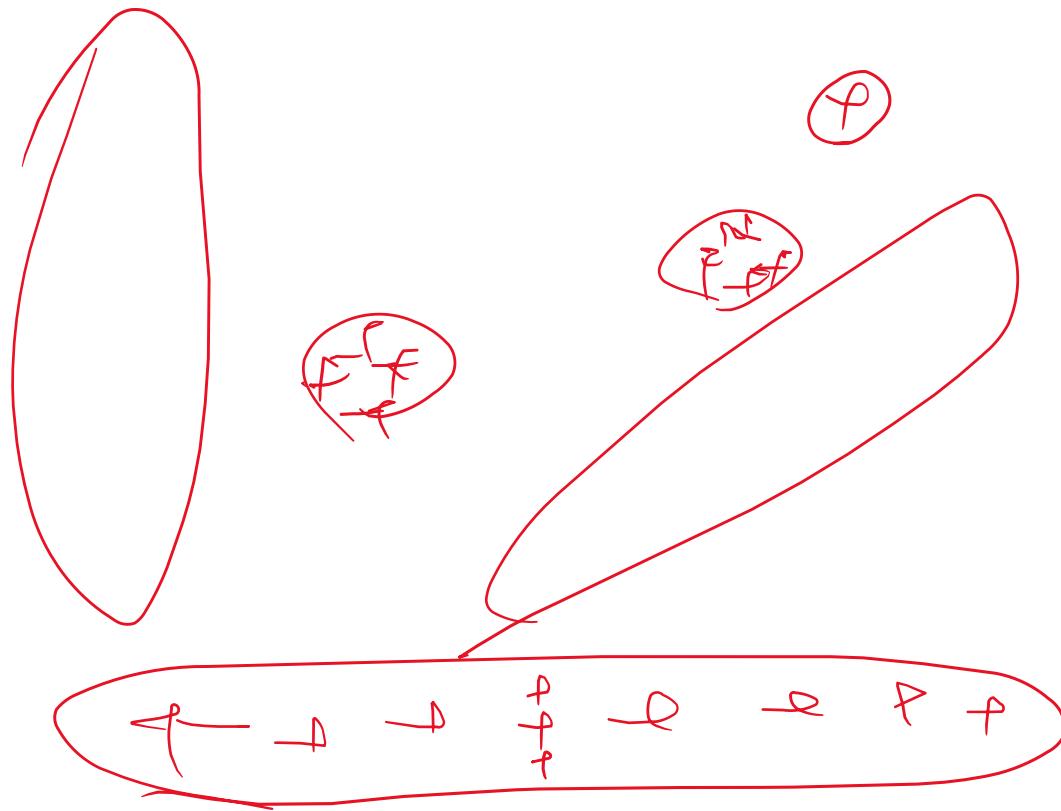
- 1 Definition
- 2 Use
- 3 Algorithm components
- 4 See also
- 5 References
- 6 External links



[https://en.wikipedia.org/wiki/Multiple\\_EM\\_for\\_Motif\\_Elicitation](https://en.wikipedia.org/wiki/Multiple_EM_for_Motif_Elicitation)  
[https://en.wikipedia.org/wiki/Position\\_weight\\_matrix](https://en.wikipedia.org/wiki/Position_weight_matrix)

# Summary

- GMM
  - Mixture of Gaussians
- EM
  - Expectation
  - Maximization



More info and exact proofs

- ✓ [https://www.cs.utah.edu/~piyush/teaching/EM\\_algorithm.pdf](https://www.cs.utah.edu/~piyush/teaching/EM_algorithm.pdf)
- ✓ <http://cs229.stanford.edu/summer2019/cs229-notes8.pdf>

# Homework notes

- T8
- RoC



# Prediction and thresholds

0 - 1

Beer	Grass	Rice	Flood	Prediction
100	3	3	Yes	0.8
20	1	1	Yes	0.3
80	3	2	No	0.6
40	1	1	No	0.2
40	1	1	No	0.1

What happens if I set my threshold at 0.5?

z

# Prediction and thresholds

Beer	Grass	Rice	Flood	Prediction	Metric
100	3	3	Yes	0.8 <i>flood</i>	TP ✓
20	1	1	Yes	0.3 <i>nf</i>	FN
80	3	2	No	0.6 <i>f</i>	FA
40	1	1	No	0.2 <i>nf R</i>	TN
40	1	1	No	0.1 <i>nf</i>	TN

What happens if I set my threshold at 0.5?

True positive rate =

$$\frac{1}{2}$$

False alarm rate =

$$\frac{1}{3}$$

Precision =

Recall =

# Prediction and thresholds

Beer	Grass	Rice	Flood	Prediction	Metric
100	3	3	Yes	0.8	TP
20	1	1	Yes	0.3	FN
80	3	2	No	0.6	FA
40	1	1	No	0.2	TN
40	1	1	No	0.1	TN

What happens if I set my threshold at 0.5?

True positive rate =  $\frac{1}{2}$

False alarm rate =  $\frac{1}{3}$

Precision =  $\frac{1}{2}$

Recall =  $\frac{1}{2}$

# Prediction and thresholds

Beer	Grass	Rice	Flood	Prediction	Metric
100	3	3	Yes	0.8 <i>l</i>	
20	1	1	Yes	0.3 <i>#</i>	
80	3	2	No	0.6 <i>f</i>	
40	1	1	No	0.2 <i>f</i>	
40	1	1	No	0.1 <i>nt</i>	

What happens if I set my threshold at 0.15?

# Prediction and thresholds

Beer	Grass	Rice	Flood	Prediction	Metric
100	3	3	Yes	0.8	TP
20	1	1	Yes	0.3	TP
80	3	2	No	0.6	FA
40	1	1	No	0.2	FA
40	1	1	No	0.1	TN

What happens if I set my threshold at 0.15?

True positive rate = 1 X

False alarm rate = 2/3 X

Precision = 2/4 ✓

Recall = 2/2 ✓

# Prediction and thresholds

Beer	Grass	Rice	Flood	Prediction	Metric
100	3	3	Yes	0.8	
20	1	1	Yes	0.3	
80	3	2	No	0.6	
40	1	1	No	0.2	
40	1	1	No	0.1	

What happens if I set my threshold at 0.5?

$$\text{True positive rate} = \frac{1}{2}$$

$$\text{False alarm rate} = \frac{1}{3}$$

$$\text{Precision} = \frac{1}{2}$$

$$\text{Recall} = \frac{1}{2}$$

What happens if I set my threshold at 0.15?

$$\text{True positive rate} = 1$$

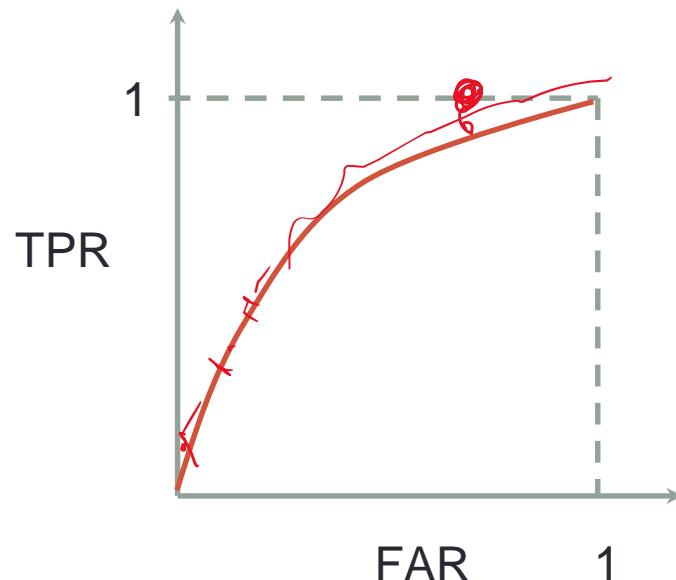
$$\text{False alarm rate} = 2/3$$

$$\text{Precision} = 2/4$$

$$\text{Recall} = 2/2$$

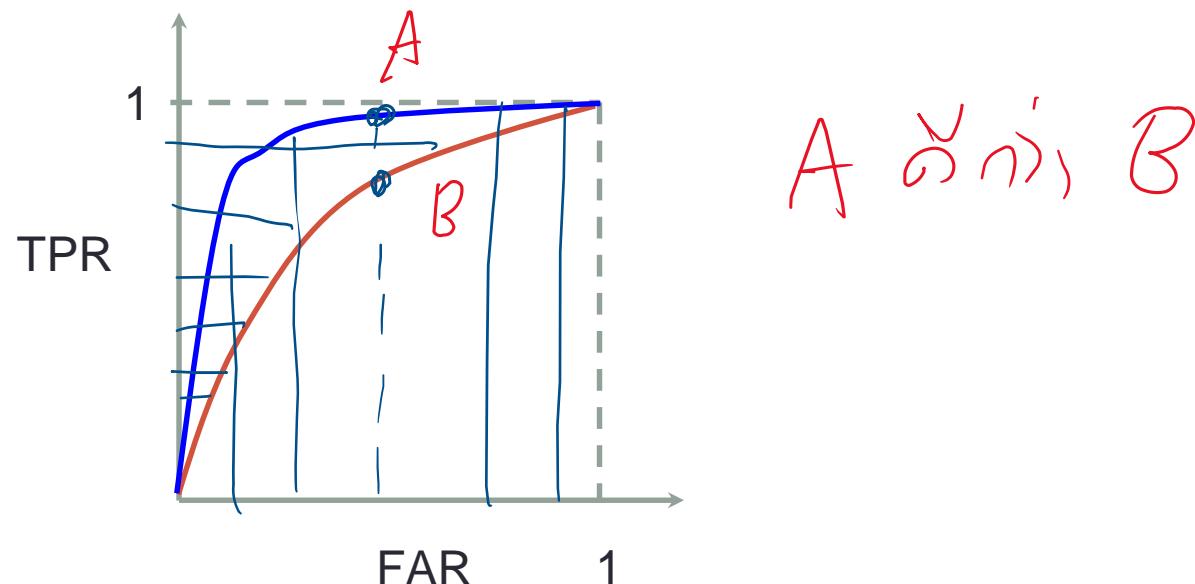
# Receiver operating Characteristic (RoC) curve

- What if we change the threshold
- FA TP is a tradeoff
  - This is why we need to think of the application when thinking of metrics.
- Plot FA rate and TP rate as threshold changes



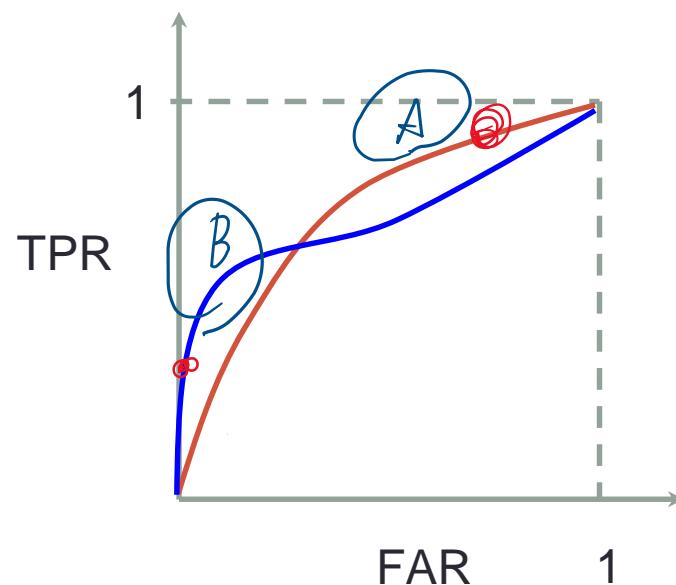
# Comparing detectors

- Which is better?



# Comparing detectors

- Which is better?

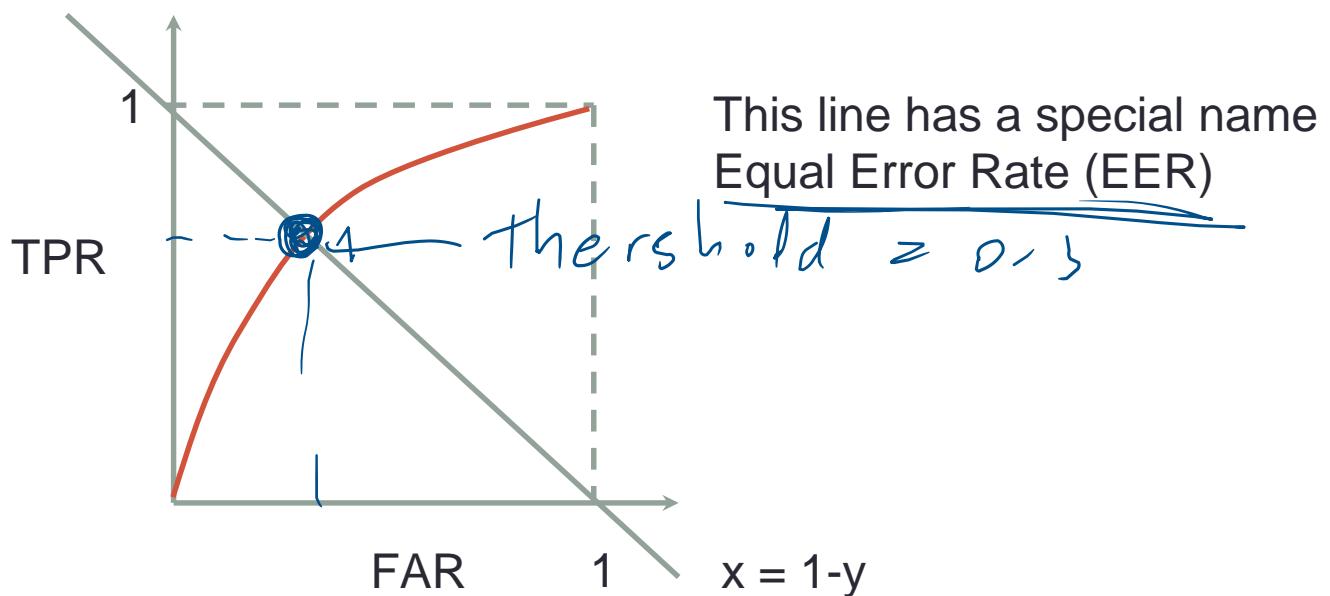


# Selecting the threshold

FN FA

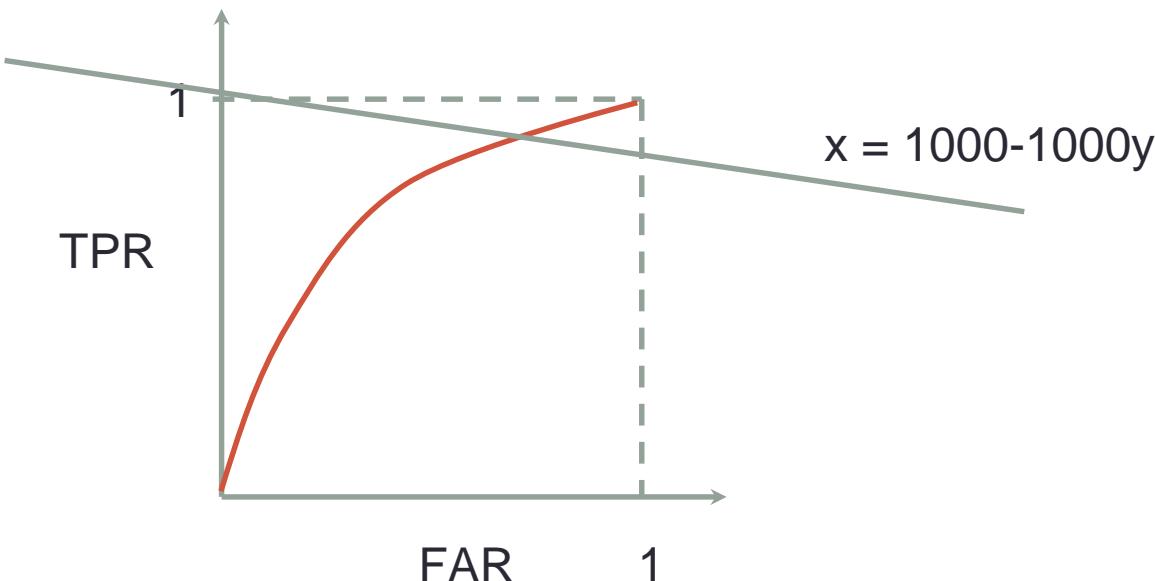
- Select based on the application
- Trade off between TP and FA. Know your application, know your users.
  - A miss is as bad as a false alarm

$$\text{FAR} = \frac{\text{FN}}{\text{TPR}} \Rightarrow x = 1-y$$



# Selecting the threshold

- Select based on the application
- Trade off between TP and FA. Know your application, know your users. Is the application about safety?
  - A miss is 1000 times more costly than false alarm.
    - $\text{FAR} = 1000(1-\text{TPR}) \Rightarrow x = 1000-1000y$



# Churn prediction

TP

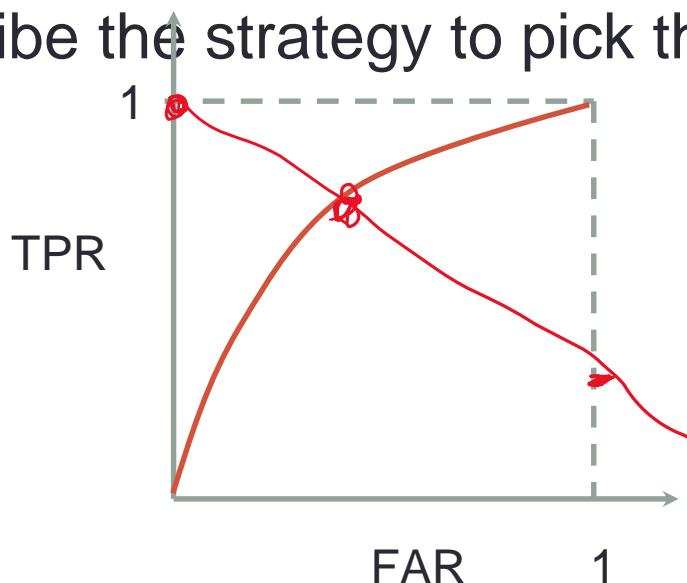
Predict whether a customer will stop subscription, so we can send a promotional ad.

Usual subscription fee 50

Cost of calling the customer 5

Promotional subscription fee 25

Describe the strategy to pick the threshold



FN

$\rightarrow$

FA

$$5 + (50 - 25)$$

30

$$30 \text{ FAR} = 50 \text{ FN}$$

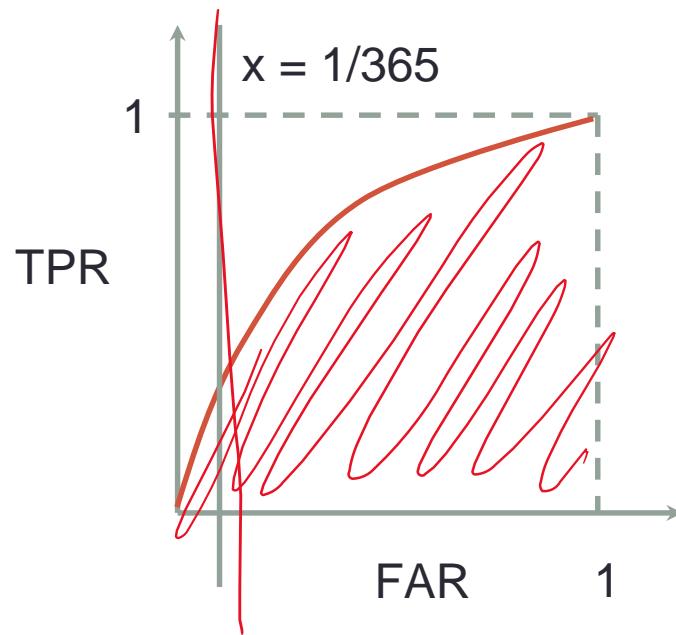
$$30 \text{ FAR} = 50(1 - \text{TPR})$$

$$30X = 50 - 50Y$$

$$y = -\frac{30}{50}x + 1$$

# Selecting the threshold

- Select based on the application
- Trade off between TP and FA.
  - Regulation or hard threshold
  - Cannot exceed 1 False alarm per year
    - If 1 decision is made everyday,  $\text{FAR} = 1/365$



# Notes about RoC

INFERENCE

- Ways to compress RoC to just a number for easier comparison -- use with care!!
  - EER
  - Area under the curve
  - F score
- Other similar curve - Detection Error Tradeoff (DET) curve
  - Plot False alarm vs Miss rate
- Can plot on log scale for clarity

