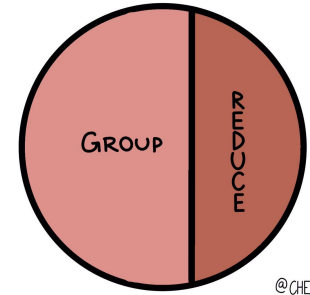


SIMPLIFY

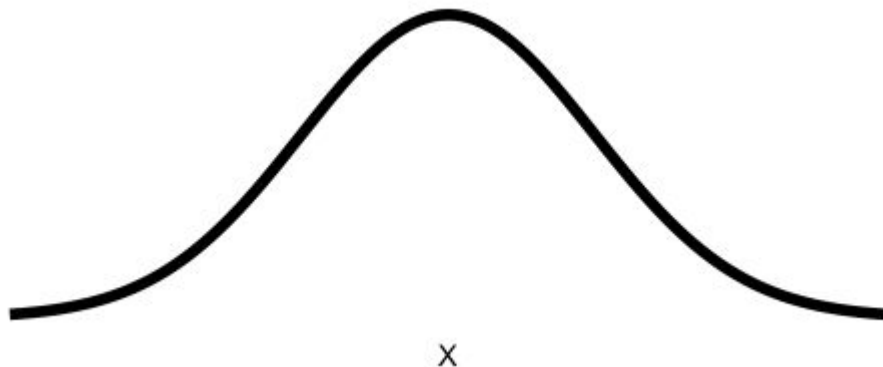


@CHELSEAPARLETT

# Gaussian Mixtures

Dr. Chelsea Parlett-Pelleriti

# Normal (Gaussian) Distribution



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

$\mu$  = Mean

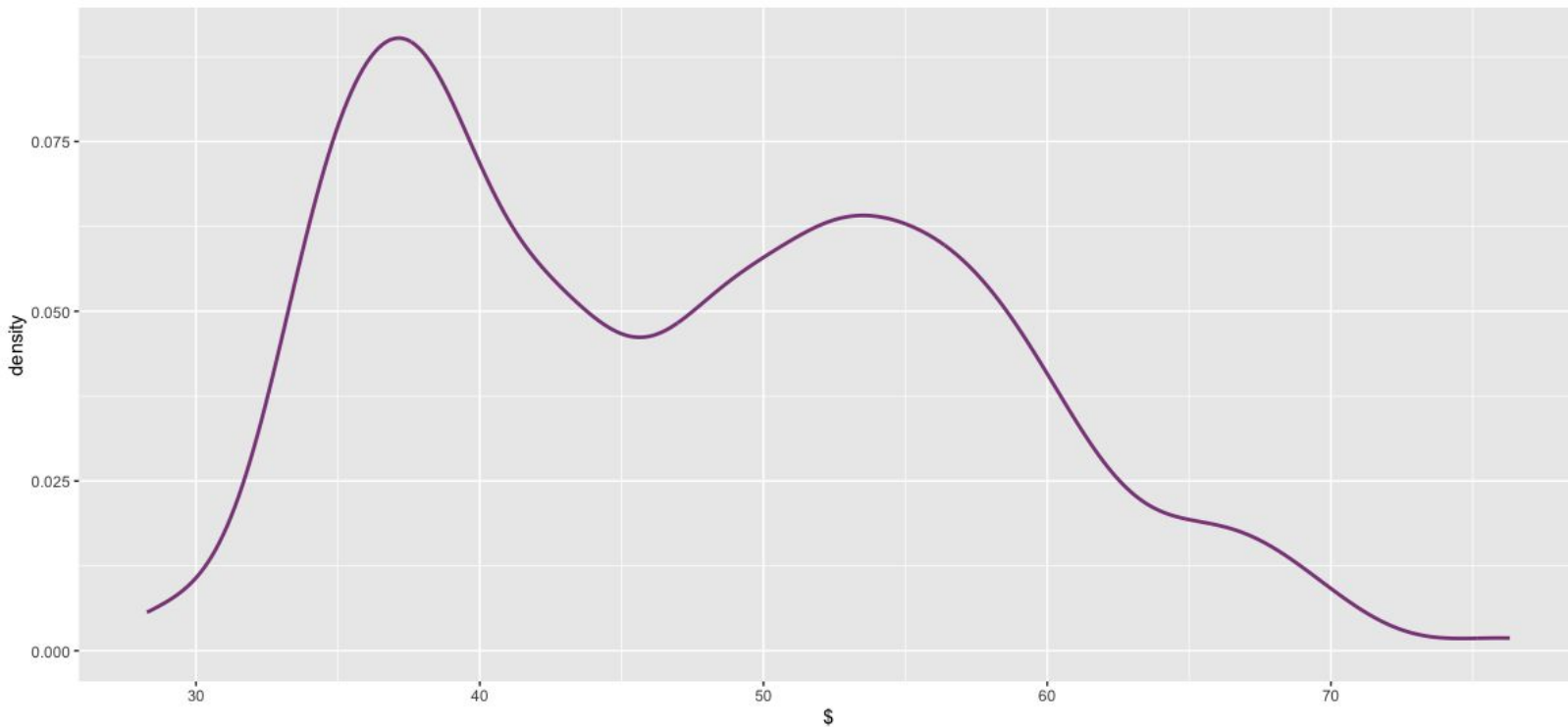
$\sigma$  = Standard Deviation

$\pi \approx 3.14159\dots$

$e \approx 2.71828\dots$

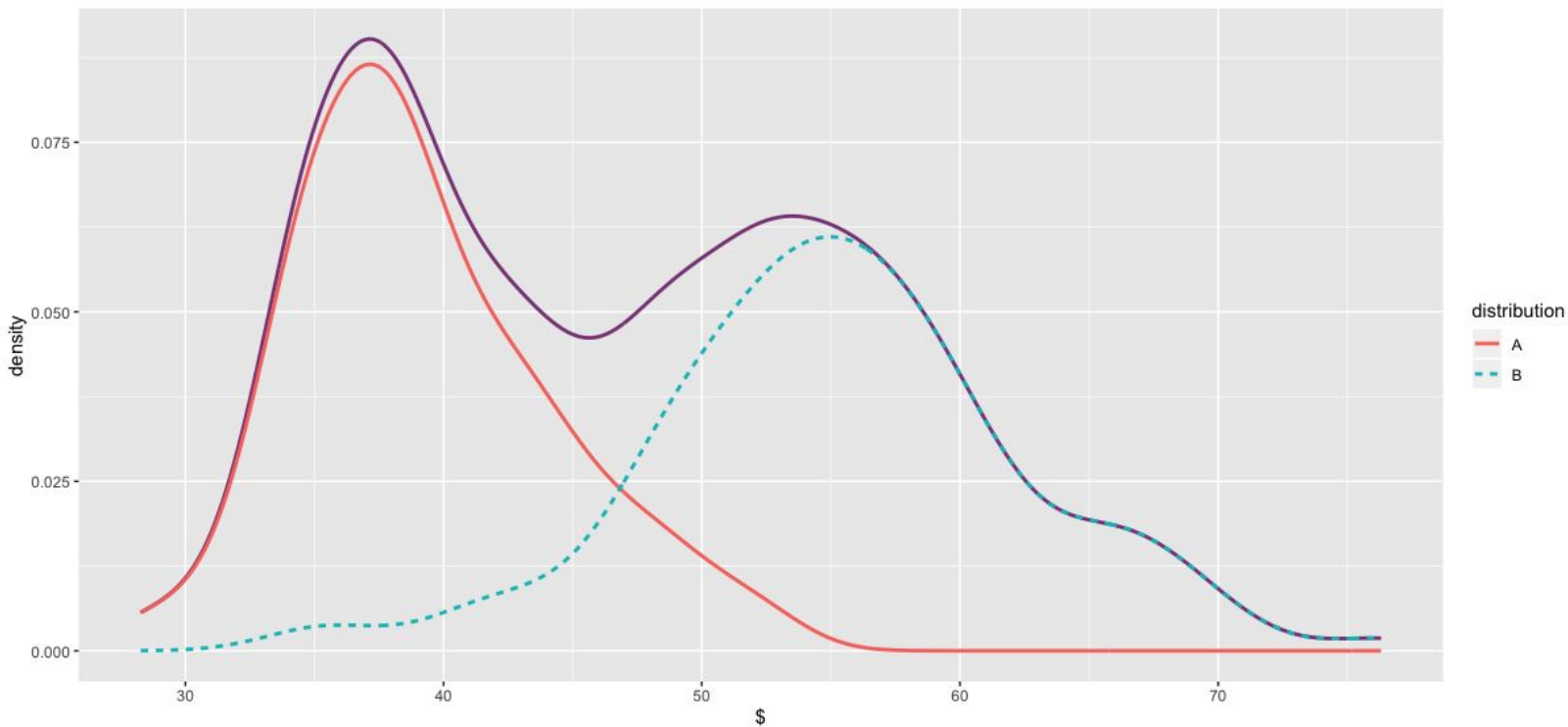
# GMM

\$ Spent on Fast Food per Week



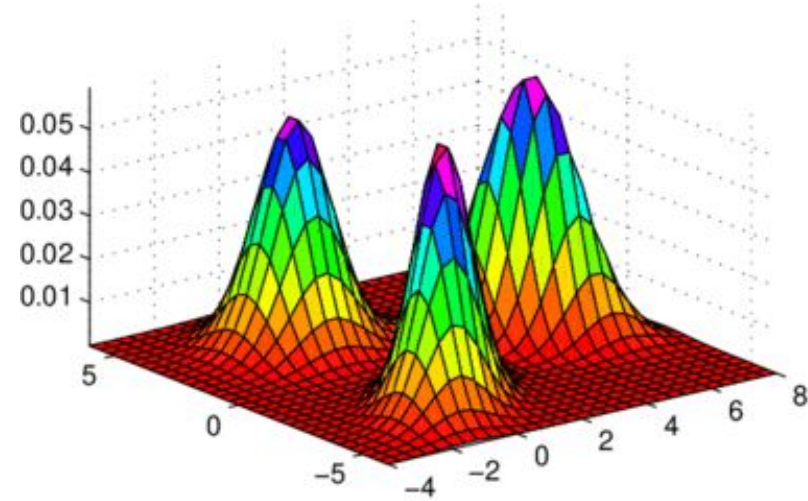
# GMM

\$ Spent on Fast Food per Week

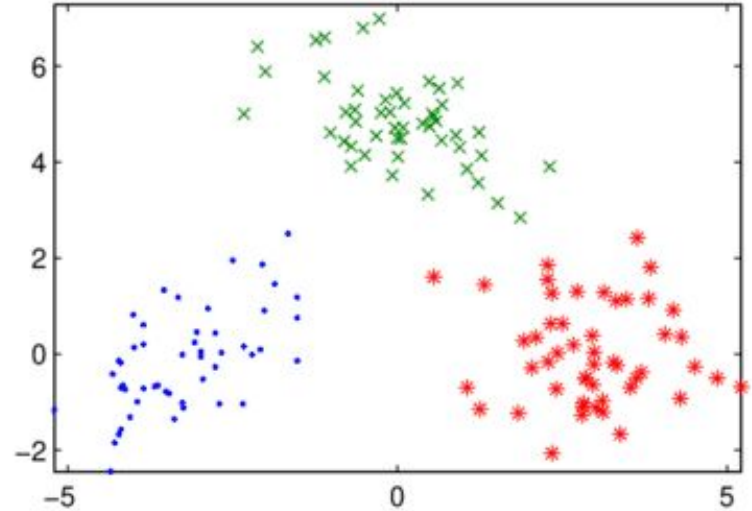


# Multivariate Normal Distributions

( a ) 3 components Gaussian mixture density



( b ) Data from 3 components Gaussian mixture density



Assume:

- Latent Groups
- Groups are Gaussian

# GMM

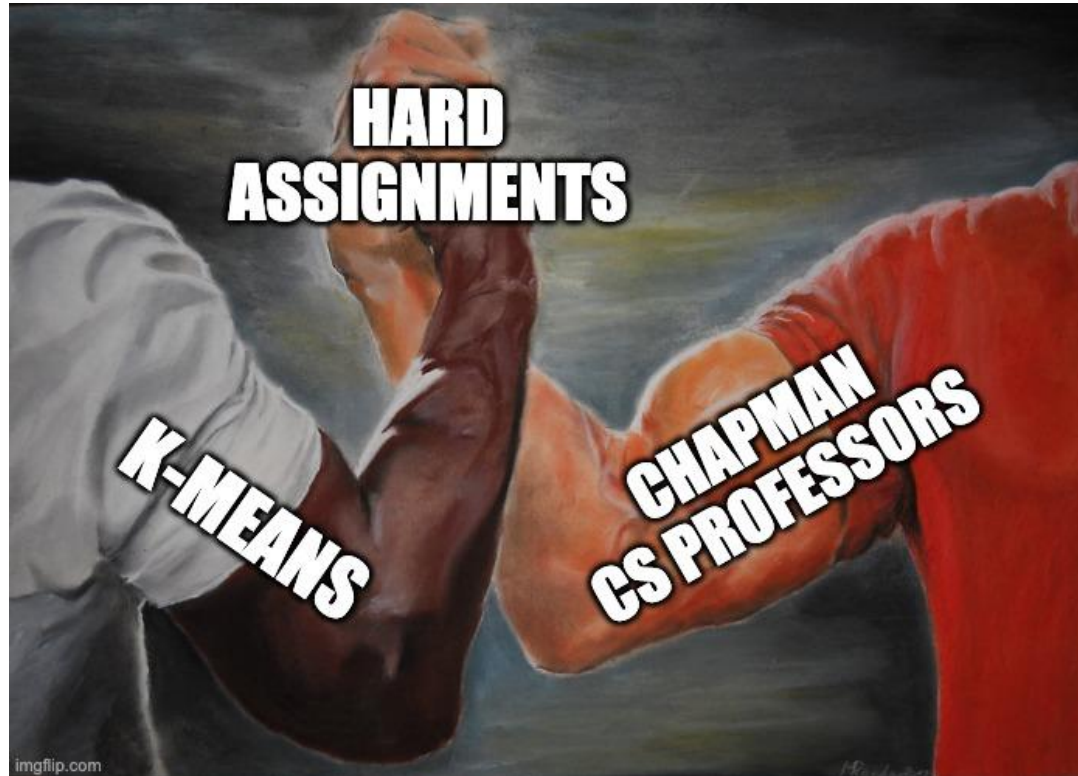
## K means

- Hard Assignment
- All Variances the Same
- Roughly same # of datapoints

## GMM

- Soft (probabilistic) Assignment
- Variances can be different
- Explicitly model # of data points

GMM



# GMM

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

Goal: choose  $r_{nk}$  and  $\mu_k$  that minimize  $J$



# GMM

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

$$r_{nk} = \begin{cases} 1 & \text{if } k = \arg \min_j \|x_n - \mu_j\|^2 \\ 0 & \text{otherwise} \end{cases}$$

# GMM

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

$$\underbrace{2 \sum_{n=1}^N r_{nk} (x_n - \mu_k)}_{\text{derivative of } J} = 0$$

# GMM

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

$$\underbrace{2 \sum_{n=1}^N r_{nk} (x_n - \mu_k)}_{\text{derivative of } J} = 0$$

$$\mu_k = \frac{\sum_n r_{nk} x_n}{\sum_n r_{nk}}$$

GMM

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

derivative of  $J$

$$\mu_k = \frac{\sum_n r_{nk} x_n}{\sum_n r_{nk}}$$

# GMM

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

when  $r_{nk}$  can  
only be 0 or 1  
(**hard assignment**)

$$\underbrace{2 \sum_{n=1}^N r_{nk} (x_n - \mu_k)}_{\text{derivative of } J} = 0$$

$$\mu_k = \frac{\sum_n r_{nk} x_n}{\sum_n r_{nk}}$$

# GMM

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

when  $r_{nk}$  can  
only be 0 or 1  
(**hard assignment**)

$$2 \sum_{n=1}^N r_{nk} (x_n - \mu_k) = 0$$

derivative of  $J$

The value of  $\mu_k$   
that minimizes our  
loss is the mean of  
all data points  
belonging to a  
cluster

$$\mu_k = \frac{\sum_n r_{nk} x_n}{\sum_n r_{nk}}$$

# GMM

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

when  $r_{nk}$  can  
only be 0 or 1  
(**hard assignment**)

$$2 \sum_{n=1}^N r_{nk} (x_n - \mu_k) = 0$$

derivative of  $J$

The value of  $\mu_k$   
that minimizes our  
loss is the mean of  
all data points  
belonging to a  
cluster

$$\mu_k = \frac{\sum_n r_{nk} x_n}{N_k} = \frac{1}{N_k} \sum_n r_{nk} x_n$$

# K-Means Algorithm

1. Choose  $k$  random points to be cluster centers
2. For each data point, assign it to the cluster whose center is closest
3. Using these assignments, recalculate the centers
4. Repeat 2 and 3 until either:
  - a. Cluster membership does not change
  - b. Centers change only a tiny amount



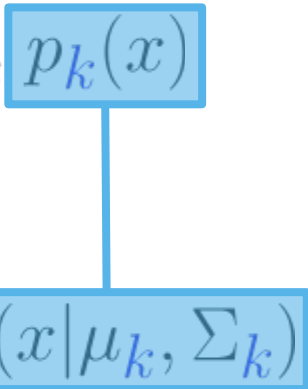
# Gaussian Mixture Model (EM Algorithm)

1. Choose **k** random points to be cluster centers (or estimate using **k-means...etc**)
2. For each data point, calculate the **probability** of belonging to each cluster
3. Using these probability weights, recalculate the **means + variances** (and weights)
4. Repeat 2 and 3 until **distributions converge**.

# Mixtures

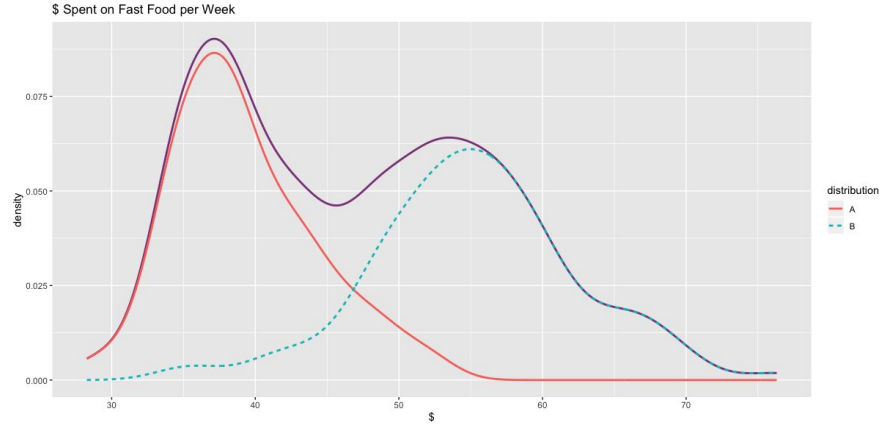
$$p(x) = \sum_{k=1}^K w_k p_k(x)$$

# Mixtures

$$p(x) = \sum_{k=1}^K w_k p_k(x)$$


$$p(x) = \sum_{k=1}^K w_k \mathcal{N}(x | \mu_k, \Sigma_k)$$

# Mixtures



$$p(x) = \sum_{k=1}^K w_k \mathcal{N}(x | \mu_k, \Sigma_k)$$

# Mixtures

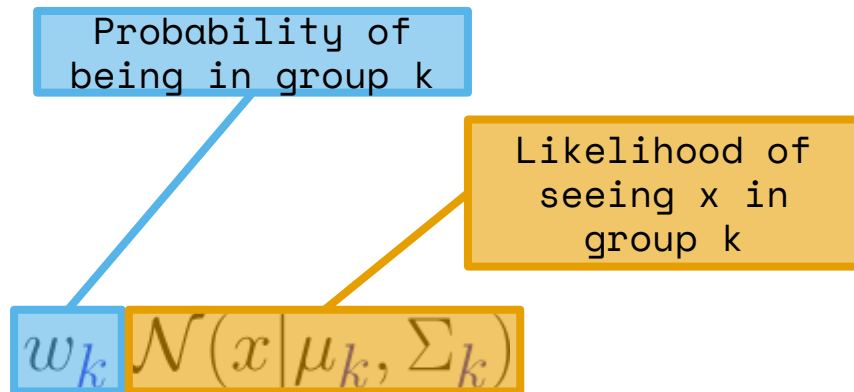
The diagram illustrates the components of a mixture model equation. The equation is  $p(x) = \sum_{k=1}^K w_k \mathcal{N}(x | \mu_k, \Sigma_k)$ . The term  $w_k$  is highlighted in a blue box, with a blue line connecting it to a text box above that reads "Probability of being in group k". The term  $\mathcal{N}(x | \mu_k, \Sigma_k)$  is highlighted in an orange box, with an orange line connecting it to a text box to its right that reads "Likelihood of seeing x in group k".

$$p(x) = \sum_{k=1}^K w_k \mathcal{N}(x | \mu_k, \Sigma_k)$$

Probability of being in group k

Likelihood of seeing x in group k

# Mixtures



# Posterior Probabilities

Prior probability of  
being in cluster  $k$

Likelihood of seeing  
 $x$  in cluster  $k$

$p(\text{cluster } k|x)$

$$= \frac{w_k \mathcal{N}(x|\mu_k, \Sigma_k)}{\sum_{j=1}^K w_j \mathcal{N}(x|\mu_j, \Sigma_j)}$$

Posterior  
probability of  
being in cluster  $k$

# Maximum Likelihood Estimation

$$p(x) = \sum_{k=1}^K w_k \mathcal{N}(x | \mu_k, \Sigma_k)$$

Probability of being in group k

Likelihood of seeing x in group k

$$p(\mathbf{X} | \mathbf{w}, \mu, \Sigma) = p(x_1, x_2, \dots, x_n | \mathbf{w}, \mu, \Sigma) =$$

$$\prod_{n=1}^N \sum_{k=1}^K w_k \mathcal{N}(x_n | \mu_k, \Sigma_k)$$



# Maximum Likelihood Estimation

$$p(x) = \sum_{k=1}^K w_k \mathcal{N}(x|\mu_k, \Sigma_k)$$

Probability of being in group k

Likelihood of seeing x in group k

$$p(\mathbf{X}|\mathbf{w}, \mu, \Sigma) = p(x_1, x_2, \dots, x_n|\mathbf{w}, \mu, \Sigma) =$$

$$\prod_{n=1}^N \sum_{k=1}^K w_k \mathcal{N}(x_n|\mu_k, \Sigma_k)$$

$$\log(p(\mathbf{X}|\mathbf{w}, \mu, \Sigma)) = \sum_{n=1}^N \log \left\{ \sum_{k=1}^K w_k \mathcal{N}(x_n|\mu_k, \Sigma_k) \right\}$$

# Maximum Likelihood Estimation

$$p(x) = \sum_{k=1}^K w_k \mathcal{N}(x | \mu_k, \Sigma_k)$$

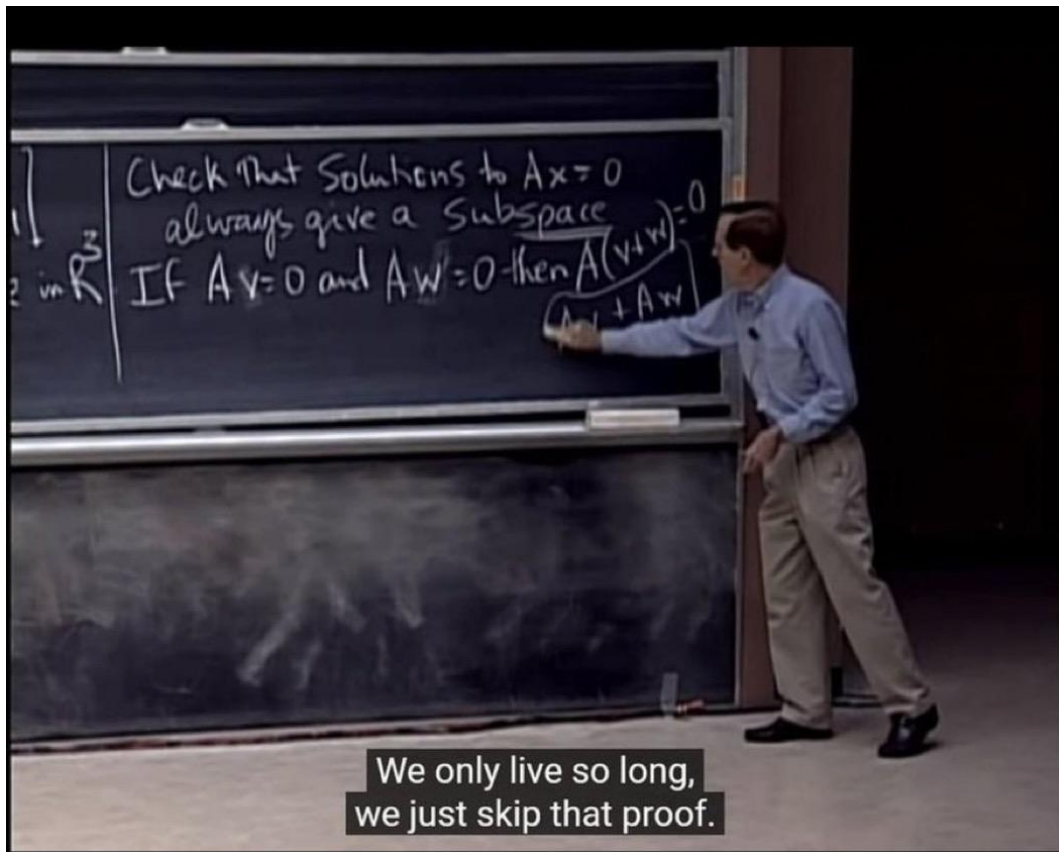
Probability of being in group k

Likelihood of seeing x in group k

$$\log(p(\mathbf{X} | \mathbf{w}, \mu, \Sigma)) = \sum_{n=1}^N \log \left\{ \sum_{k=1}^K w_k \mathcal{N}(x_n | \mu_k, \Sigma_k) \right\}$$

Goal: choose  $w$ ,  $\mu$ ,  $\Sigma$  that maximize the log likelihood

# Maximum Likelihood Estimation



$$p(x) = \sum_{k=1}^K w_k \mathcal{N}(x | \mu_k, \Sigma_k)$$

Probability of being in group  $k$

Likelihood of seeing  $x$  in group  $k$

$$\left. w_k \mathcal{N}(x | \mu_k, \Sigma_k) \right\}$$

## Formulas (E-Step)

$$r_{nk} = \frac{w_k N(x_n | \mu_k, \Sigma_k)}{\sum_j w_j N(x_n | \mu_j, \Sigma_j)}$$

$$w_k \mathcal{N}(x|\mu_k, \Sigma_k)$$

## Formulas (E-Step)

$$r_{nk} = \frac{w_k \mathcal{N}(x_n|\mu_k, \Sigma_k)}{\sum_j w_j \mathcal{N}(x_n|\mu_j, \Sigma_j)}$$

Responsibilities are  
the posterior  
probability of a data  
point being in cluster  
k

$$p(\text{cluster } k|x)$$

Posterior  
probability of being  
in cluster k

Prior probability of  
being in cluster k

Likelihood of seeing  
x in cluster k

$$= \frac{w_k \mathcal{N}(x|\mu_k, \Sigma_k)}{\sum_{j=1}^K w_j \mathcal{N}(x|\mu_j, \Sigma_j)}$$

## Formulas (M-Step)

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} x_n$$

$$N_k = \sum_{n=1}^N r_{nk}$$

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} (x_n - \mu_k)(x_n - \mu_k)^T$$

## Formulas (M-Step)

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} x_n$$

$$N_k = \sum_{n=1}^N r_{nk}$$

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} (x_n - \mu_k)(x_n - \mu_k)^T$$

## Formulas (M-Step)

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N \boxed{r_{nk}} \boxed{x_n}$$

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N \boxed{r_{nk}} \boxed{(x_n - \mu_k)(x_n - \mu_k)^T}$$

$$N_k = \sum_{n=1}^N r_{nk}$$



## Formulas (M-Step)

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} x_n$$

$$N_k = \sum_{n=1}^N r_{nk}$$

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} (x_n - \mu_k)(x_n - \mu_k)^T$$

## Formulas (M-Step)

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} x_n$$

$$N_k = \sum_{n=1}^N r_{nk}$$

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} (x_n - \mu_k)(x_n - \mu_k)^T$$

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

## Comparison to K-Means

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} x_n$$

$$N_k = \sum_{n=1}^N r_{nk}$$

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} (x_n - \mu_k)(x_n - \mu_k)^T$$

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

# Math Summary

- GMM does **soft assignment**, every data point belongs to every cluster with some probability
- Data points that are more likely to be in a cluster have **more influence** over its parameters
- GMM uses the EM algorithm to iteratively update the cluster distributions. It first assigning a responsibility to each data point (**E-step**), and then using them to calculate weighted means and variances for each cluster (**M-step**)
- Responsibilities measure the **probability of a data point being in each cluster** (technically the **posterior probability**).
- Responsibilities contain information about **how common a cluster is** as well as the **likelihood of a data point belonging to that cluster**

# Math Summary

- GMM does **soft assignment**, every data point belongs to every cluster with some probability
- Data points that are more likely to be in a cluster have **more influence** over its parameters
- GMM uses the EM algorithm to iteratively update the cluster distributions. It first assigning a responsibility to each data point (**E-step**), and then using them to calculate weighted means and variances for each cluster (**M-step**)
- Responsibilities measure the **probability of a data point being in each cluster** (technically the **posterior probability**).
- Responsibilities contain information about **how common a cluster is** as well as the **likelihood of a data point belonging to that cluster**

# Math Summary

- GMM does **soft assignment**, every data point belongs to every cluster with some probability
- Data points that are more likely to be in a cluster have **more influence** over its parameters
- GMM uses the EM algorithm to iteratively update the cluster distributions. It first assigning a responsibility to each data point (**E-step**), and then using them to calculate weighted means and variances for each cluster (**M-step**)
- Responsibilities measure the **probability of a data point being in each cluster** (technically the **posterior probability**).
- Responsibilities contain information about **how common a cluster is** as well as the **likelihood of a data point belonging to that cluster**

# Math Summary

- GMM does **soft assignment**, every data point belongs to every cluster with some probability
- Data points that are more likely to be in a cluster have **more influence** over its parameters
- GMM uses the EM algorithm to iteratively update the cluster distributions. It first assigning a responsibility to each data point (**E-step**), and then using them to calculate weighted means and variances for each cluster (**M-step**)
- Responsibilities measure the **probability of a data point being in each cluster** (technically the **posterior probability**).
- Responsibilities contain information about **how common a cluster is** as well as the **likelihood of a data point belonging to that cluster**

# Math Summary

- GMM does **soft assignment**, every data point belongs to every cluster with some probability
- Data points that are more likely to be in a cluster have **more influence** over its parameters
- GMM uses the EM algorithm to iteratively update the cluster distributions. It first assigning a responsibility to each data point (**E-step**), and then using them to calculate weighted means and variances for each cluster (**M-step**)
- Responsibilities measure the **probability of a data point being in each cluster** (technically the **posterior probability**).
- Responsibilities contain information about **how common a cluster is** as well as the **likelihood of a data point belonging to that cluster**



# GMM

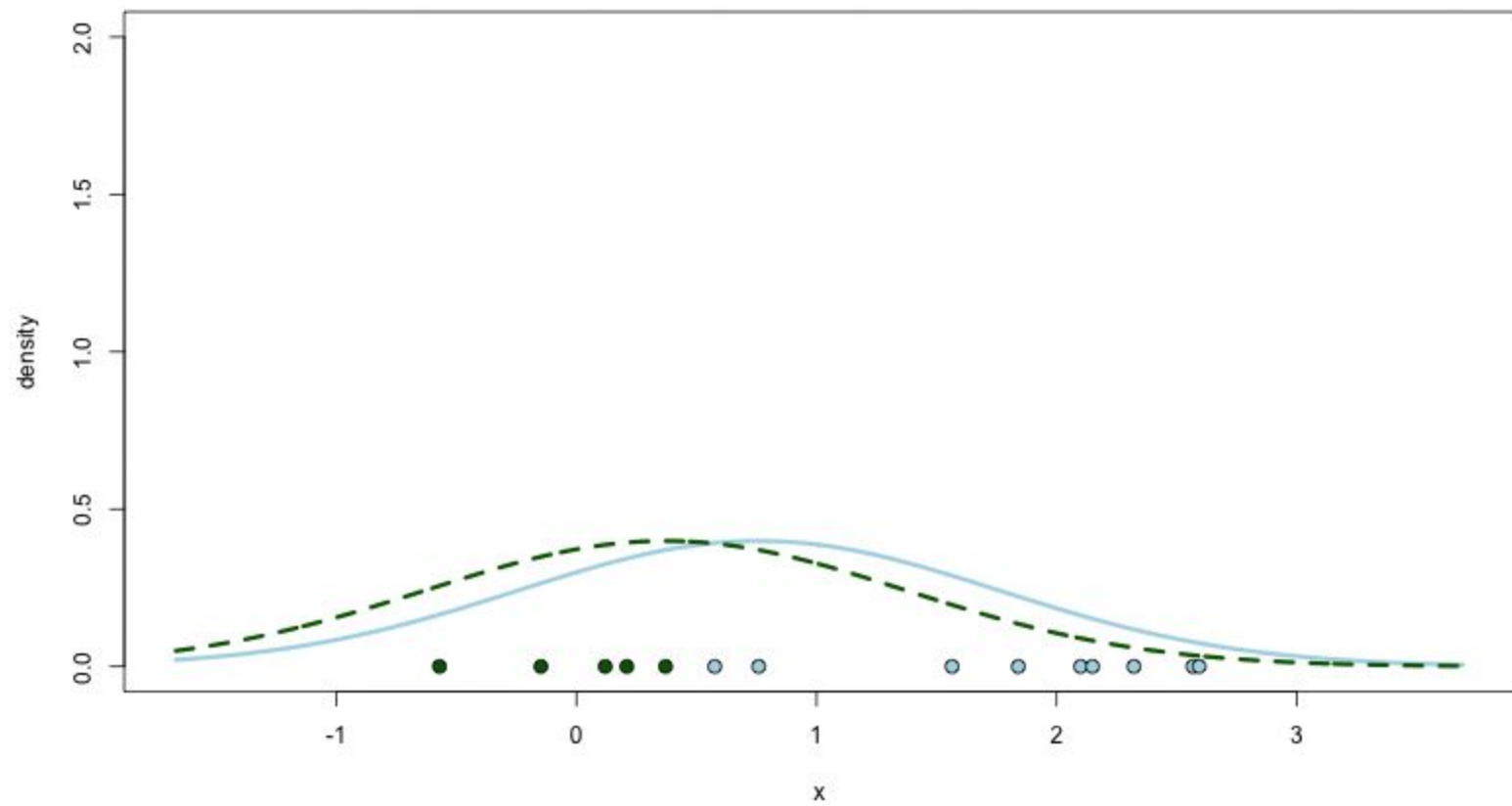
## K means

- Hard Assignment
- All Variances the Same
- Roughly same # of datapoints

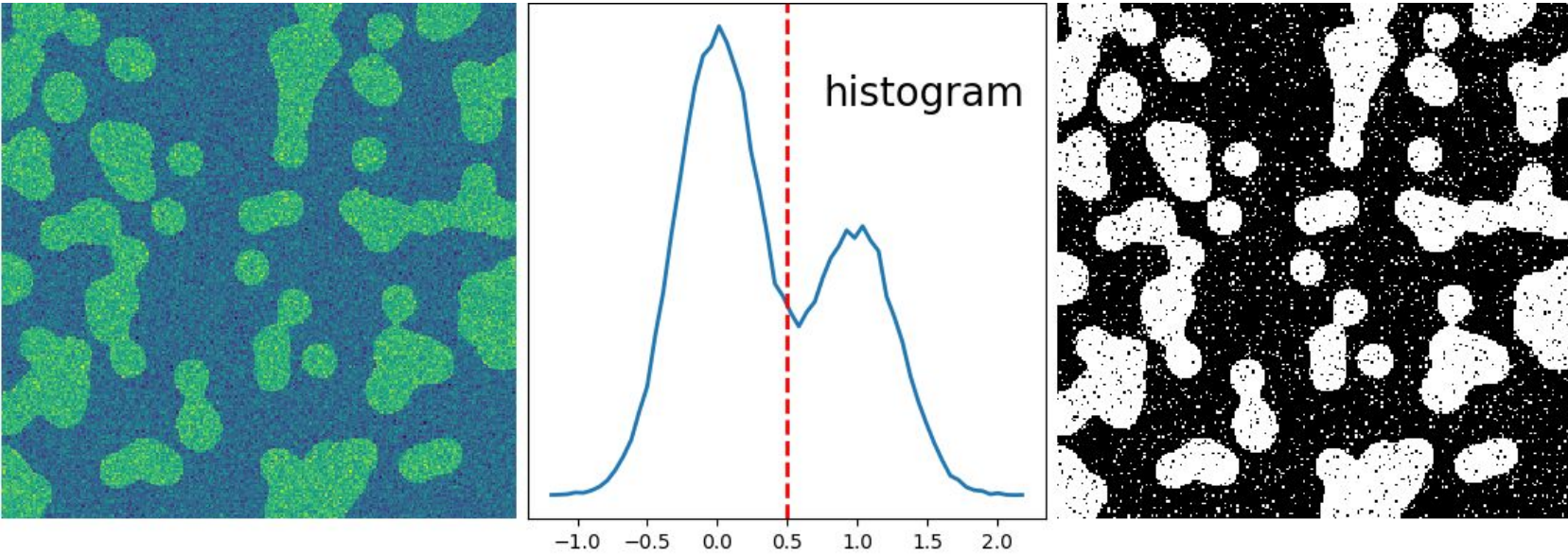
## GMM

- Soft (probabilistic) Assignment
- Variances can be different
- Explicitly model # of data points

E1



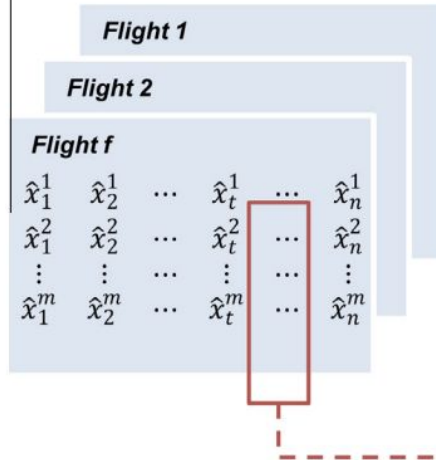
# Applications



# Applications

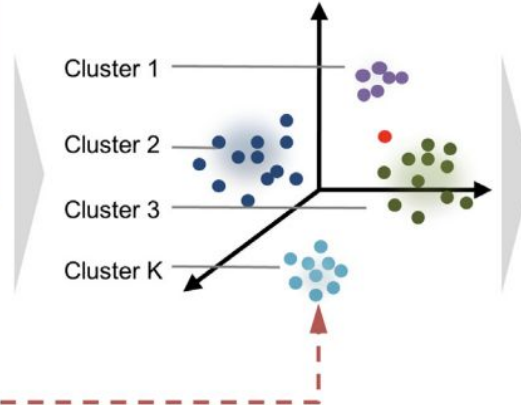
## Normalized vectors

Every flight parameter is normalized to have “zero mean and unit variance”



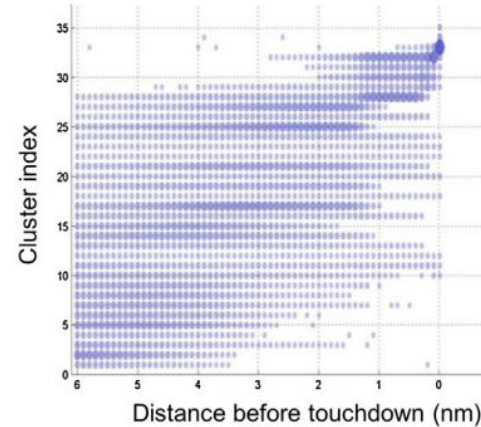
## Clusters

GMM clustering is performed on normalized vectors; each cluster represent a typical operation of aircraft



## Temporal distribution of clusters

The temporal distribution of clusters is summarized by observation frequency of each cluster along the temporal reference



Larger circle size and darker color indicates a higher observation frequency

**Fig. 3.** Cluster analysis: identify typical operations and temporal distribution.

# Applications

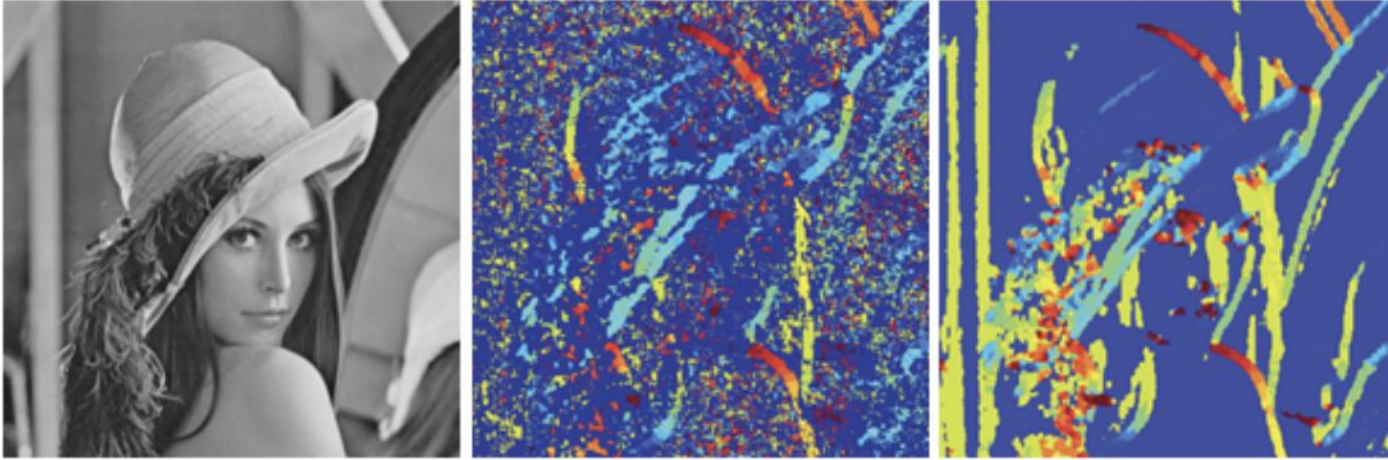
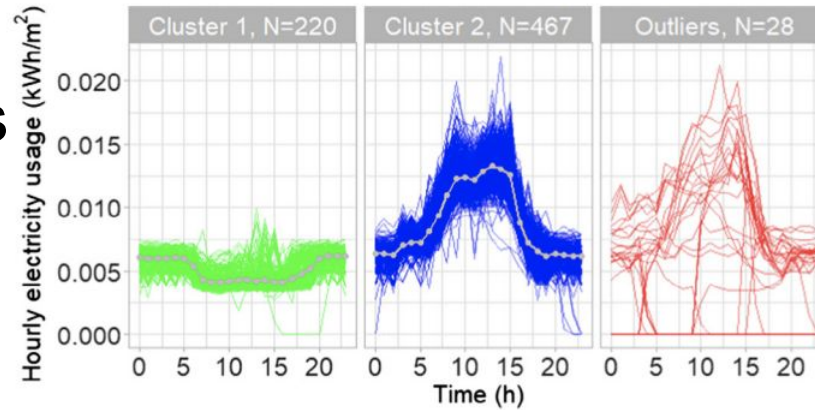
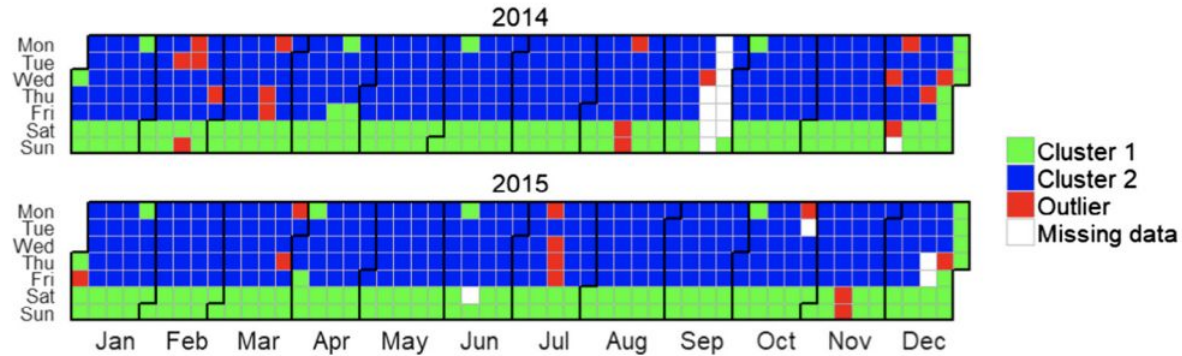


Fig. 1. Illustration of clustering of patches in the PLE method for the Lena image. LEFT: Original image; RIGHT: Clustered image; The pixels in the same color indicate that  $8 \times 8$  patches around them are in the same cluster. It can be seen that patches from different parts of image are grouped into one cluster [17].

# Applications



a) TDEU profiles and outliers



b) Distribution of the TDEU profiles

**Fig. 7.** Visualisation of the intra-building clustering result of Building #16.