

# Data Modelling Methods-II

CS771: Introduction to Machine Learning  
Purushottam Kar

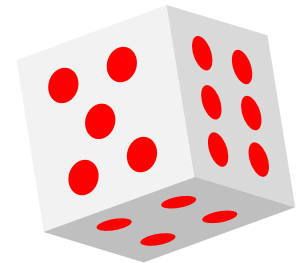
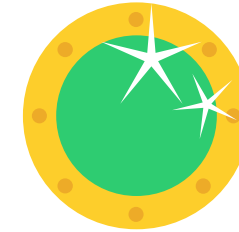


# Outline of today's discussion

- Revise Naïve Bayes method
- Feature modelling methods for unlabelled data
- Gaussian mixture Models (GMMs)
- The alternating optimization approach to learning GMMs
- The k-means approach to learning GMMs

# The Multinoulli Distribution

- Also known as categorical distribution
- Generalizes the Bernoulli distribution
  - Bernoulli models a coin with two outcomes (call them head and tail)
  - ... using one “bias” parameter  $p$  (taken to be the probability of heads)
- Multinoulli models a  $K$ -sided dice
  - ... using a  $K$ -dimensional vector  $\pi$
  - $\pi_k$  is taken to denote the probability of the  $k$ -th side turning up
  - $\pi_k \geq 0$  and  $\sum_{k=1}^K \pi_k = 1$
- Conjugate prior for Bernoulli: Beta  $\mathbb{P}[p; \alpha, \beta] = \frac{1}{B(\alpha, \beta)} p^{\alpha-1} (1-p)^{\beta-1}$
- Conjugate prior for Multinoulli: Dirichlet  $\mathbb{P}[\boldsymbol{\pi}; \boldsymbol{\alpha}] = \frac{1}{B(\boldsymbol{\alpha})} \prod_{i=1}^L \pi_i^{\alpha_i-1}$



# App: Email Categorizer using Naïve Bayes

# App: Email Categorizer using Naïve Bayes

HEC			DoSA	Supervisor
<input type="checkbox"/>	☆	➤	HallPresi	Awesome talk at eSumm
<input type="checkbox"/>	☆	➤	HEC	Urgent! Mess bill overdue
<input type="checkbox"/>	☆	➤	HEC	Urgent! Canteen bill over

# App: Email Categorizer using Naïve Bayes

HEC		DoSA		Supervisor	
<input type="checkbox"/>	☆	➤	HallPresi	Awesome talk at eSumm	
<input type="checkbox"/>	☆	➤	HEC	Urgent! Mess bill overdue	
<input type="checkbox"/>	☆	➤	HEC	Urgent! Canteen bill over	

**X**

talk	meet	project	wake	meeting	wallet	keys	awesome	lost	trip	lost	mess	report

# App: Email Categorizer using Naïve Bayes

HEC			DoSA	Supervisor
<input type="checkbox"/>	☆	➤	HallPresi	Awesome talk at eSumm
<input type="checkbox"/>	☆	➤	HEC	Urgent! Mess bill overdue
<input type="checkbox"/>	☆	➤	HEC	Urgent! Canteen bill over

**X**

--	--	--	--	--	--	--	--	--	--	--	--	--

talk   meet   project   wake   meeting   wallet   keys   awesome   lost   trip   lost   mess   report

$y \in \{ \text{red circle}, \text{green circle}, \text{blue circle} \}$

# App: Email Categorizer using Naïve Bayes

HEC			DoSA	Supervisor
<input type="checkbox"/>	☆	➤	HallPresi	Awesome talk at eSumm
<input type="checkbox"/>	☆	➤	HEC	Urgent! Mess bill overdue
<input type="checkbox"/>	☆	➤	HEC	Urgent! Canteen bill over

*Bag of words* feature  
can record just  
occurrence or count

**X**

talk	meet	project	wake	meeting	wallet	keys	awesome	lost	trip	lost	mess	report

$y \in \{ \text{red circle}, \text{green circle}, \text{blue circle} \}$



# App: Email Categorizer using Naïve Bayes

HEC			DoSA	Supervisor
<input type="checkbox"/>	☆	➤	HallPresi	Awesome talk at eSumm
<input type="checkbox"/>	☆	➤	HEC	Urgent! Mess bill overdue
<input type="checkbox"/>	☆	➤	HEC	Urgent! Canteen bill over

*Bag of words* feature  
can record just  
occurrence or count

<b>x</b>	1	0	0	0	0	0	0	1	0	0	0	0	0
	talk	meet	project	wake	meeting	wallet	keys	awesome	lost	trip	lost	mess	report

$y \in \{ \text{red circle}, \text{green circle}, \text{blue circle} \}$

# App: Email Categorizer using Naïve Bayes

HEC			DoSA	Supervisor
<input type="checkbox"/>	☆	➤	HallPresi	Awesome talk at eSumm
<input type="checkbox"/>	☆	➤	HEC	Urgent! Mess bill overdue
<input type="checkbox"/>	☆	➤	HEC	Urgent! Canteen bill over

Choice of words crucial  
– stemming, throw  
away articles etc

*Bag of words* feature  
can record just  
occurrence or count

<b>x</b>	1	0	0	0	0	0	0	1	0	0	0	0	0
	talk	meet	project	wake	meeting	wallet	keys	awesome	lost	trip	lost	mess	report

$y \in \{ \text{red circle}, \text{green circle}, \text{blue circle} \}$

# App: Email Categorizer using Naïve Bayes

HEC			DoSA	Supervisor
<input type="checkbox"/>	☆	➤	HallPresi	Awesome talk at eSumm
<input type="checkbox"/>	☆	➤	HEC	Urgent! Mess bill overdue
<input type="checkbox"/>	☆	➤	HEC	Urgent! Canteen bill over

Commonly used in NLP

Choice of words crucial  
– stemming, throw away articles etc

*Bag of words* feature  
can record just  
occurrence or count

**x**

1	0	0	0	0	0	0	1	0	0	0	0	0
talk	meet	project	wake	meeting	wallet	keys	awesome	lost	trip	lost	mess	report

$y \in \{ \text{red circle}, \text{green circle}, \text{blue circle} \}$

# App: Email Categorizer using Naïve Bayes

HEC	DoSA	Supervisor
<input type="checkbox"/> <input checked="" type="checkbox"/> <input type="checkbox"/>	HallPresi	Awesome talk at eSumm
<input type="checkbox"/> <input checked="" type="checkbox"/> <input type="checkbox"/>	HEC	Urgent! Mess bill overdue
<input type="checkbox"/> <input checked="" type="checkbox"/> <input type="checkbox"/>	HEC	Urgent! Canteen bill over

Commonly used in NLP

Choice of words crucial  
– stemming, throw away articles etc

*Bag of words* feature  
can record just  
occurrence or count

**x**

1	0	0	0	0	0	0	1	0	0	0	0	0
talk	meet	project	wake	meeting	wallet	keys	awesome	lost	trip	lost	mess	report

$y \in \{ \text{red circle}, \text{green circle}, \text{blue circle} \}$

Usually very high dimensional

# App: Email Categorizer using Naïve Bayes

HEC	DoSA	Supervisor
<input type="checkbox"/> <input checked="" type="checkbox"/> <input type="checkbox"/> HallPresi		Awesome talk at eSumm
<input type="checkbox"/> <input checked="" type="checkbox"/> <input type="checkbox"/> HEC		Urgent! Mess bill overdue
<input type="checkbox"/> <input checked="" type="checkbox"/> <input type="checkbox"/> HEC		Urgent! Canteen bill over

Commonly used in NLP

Choice of words crucial  
– stemming, throw away articles etc

*Bag of words* feature can record just occurrence or count

**x**

1	0	0	0	0	0	0	1	0	0	0	0	0
talk	meet	project	wake	meeting	wallet	keys	awesome	lost	trip	lost	mess	report

$y \in \{ \text{red circle}, \text{green circle}, \text{blue circle} \}$

Usually very high dimensional

Usually very very sparse

# App: Email Categorizer using Naïve Bayes

Sept 6, 2017



# App: Email Categorizer using Naïve Bayes

$$\mathbb{P}[\bullet] = \frac{|\text{\#emails from supervisor}|}{|\text{\#total emails}|}$$

# App: Email Categorizer using Naïve Bayes

$$\mathbb{P}[\bullet] = \frac{|\text{\#emails from supervisor}|}{|\text{\#total emails}|}$$

$$\mathbb{P}[\text{awesome} = 1 \mid \bullet] = \frac{|\text{\#emails from sup. with "awesome"}|}{|\text{\#total emails from supervisor}|}$$



# App: Email Categorizer using Naïve Bayes

$$\mathbb{P}[\bullet] = \frac{|\text{\#emails from supervisor}|}{|\text{\#total emails}|}$$

$$\mathbb{P}[\text{awesome} = 1 \mid \bullet] = \frac{|\text{\#emails from sup. with "awesome"}|}{|\text{\#total emails from supervisor}|}$$

At test time ...

# App: Email Categorizer using Naïve Bayes

$$\mathbb{P}[\bullet] = \frac{|\text{\#emails from supervisor}|}{|\text{\#total emails}|}$$

$$\mathbb{P}[\text{awesome} = 1 \mid \bullet] = \frac{|\text{\#emails from sup. with "awesome"}|}{|\text{\#total emails from supervisor}|}$$

At test time ...

$$\mathbb{P}[\mathbf{x}^t, \bullet] = \mathbb{P}[\bullet] \cdot \prod_{j=1}^d \mathbb{P}[\mathbf{x}_i^t \mid \bullet]$$

# App: Email Categorizer using Naïve Bayes

$$\mathbb{P}[\bullet] = \frac{|\text{\#emails from supervisor}|}{|\text{\#total emails}|}$$

$$\mathbb{P}[\text{awesome} = 1 \mid \bullet] = \frac{|\text{\#emails from sup. with "awesome"}|}{|\text{\#total emails from supervisor}|}$$

At test time ...

$$\begin{aligned}\mathbb{P}[\text{awesome} = 0 \mid \bullet] \\ = 1 - \mathbb{P}[\text{awesome} = 1 \mid \bullet]\end{aligned}$$

$$\mathbb{P}[\mathbf{x}^t, \bullet] = \mathbb{P}[\bullet] \cdot \prod_{j=1}^d \mathbb{P}[\mathbf{x}_j^t \mid \bullet]$$

# App: Email Categorizer using Naïve Bayes

$$\mathbb{P}[\bullet] = \frac{|\text{\#emails from supervisor}|}{|\text{\#total emails}|}$$

$$\mathbb{P}[\text{awesome} = 1 \mid \bullet] = \frac{|\text{\#emails from sup. with "awesome"}|}{|\text{\#total emails from supervisor}|}$$

At test time ...

$$\begin{aligned}\mathbb{P}[\text{awesome} = 0 \mid \bullet] \\ = 1 - \mathbb{P}[\text{awesome} = 1 \mid \bullet]\end{aligned}$$

$$\mathbb{P}[\mathbf{x}^t, \bullet] = \mathbb{P}[\bullet] \cdot \prod_{j=1}^d \mathbb{P}[\mathbf{x}_j^t \mid \bullet]$$

$$\hat{y}^t = \arg \max \{ \mathbb{P}[\mathbf{x}^t, \text{red}], \mathbb{P}[\mathbf{x}^t, \text{green}], \mathbb{P}[\mathbf{x}^t, \text{blue}] \}$$

# App: Email Categorizer using Naïve Bayes

$$\mathbb{P}[\bullet] = \frac{|\text{\#emails from supervisor}|}{|\text{\#total emails}|}$$

$$\mathbb{P}[\text{awesome} = 1 \mid \bullet] = \frac{|\text{\#emails from sup. with "awesome"}|}{|\text{\#total emails from supervisor}|}$$

At test time ...

$$\begin{aligned}\mathbb{P}[\text{awesome} = 0 \mid \bullet] \\ = 1 - \mathbb{P}[\text{awesome} = 1 \mid \bullet]\end{aligned}$$

$$\mathbb{P}[\mathbf{x}^t, \bullet] = \mathbb{P}[\bullet] \cdot \prod_{j=1}^d \mathbb{P}[\mathbf{x}_j^t \mid \bullet]$$

Will give the same result as  
 $\arg \max\{\mathbb{P}[\text{red} \mid \mathbf{x}^t], \mathbb{P}[\text{green} \mid \mathbf{x}^t], \mathbb{P}[\text{blue} \mid \mathbf{x}^t]\}$

$$\hat{y}^t = \arg \max\{\mathbb{P}[\mathbf{x}^t, \text{red}], \mathbb{P}[\mathbf{x}^t, \text{green}], \mathbb{P}[\mathbf{x}^t, \text{blue}]\}$$

# App: Automatic Email Generator!

Class proportions

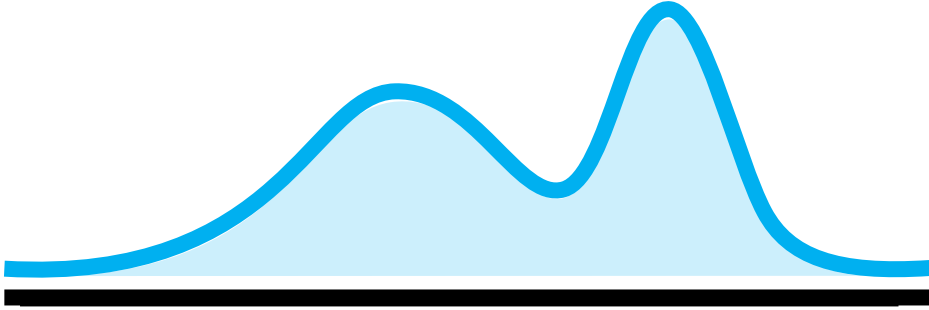
- Choose a category from {HEC, DoSA, Supervisor}
  - Toss a 3-sided "coin" aka categorical/multinoulli distribution using  $\mathbb{P}[\bullet]$
  - Say we chose DoSA
- For each word in your dictionary of  $d$  words, toss a Bernoulli coin to decide whether to include that word in the mail or not
  - For  $j \in [d]$ , toss a coin that lands heads with probability  $\mathbb{P}[x_j = 1 \mid \bullet]$
- Collect all words for which the toss landed heads
- Compose an email using only those words (and maybe a few articles, prepositions etc)
- Congratulations, you can now ask the dean to stop sending you emails – you will generate them yourself!

Already learnt from training data!



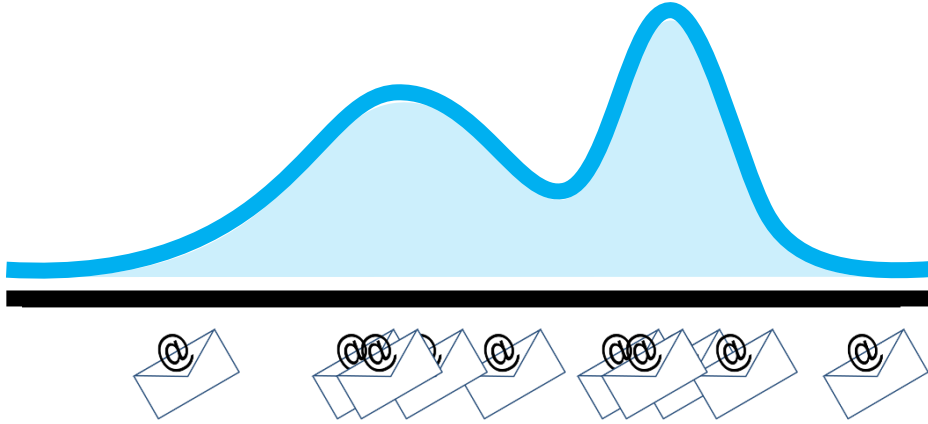
# The generative story for unlabelled data??

# The generative story for unlabelled data??

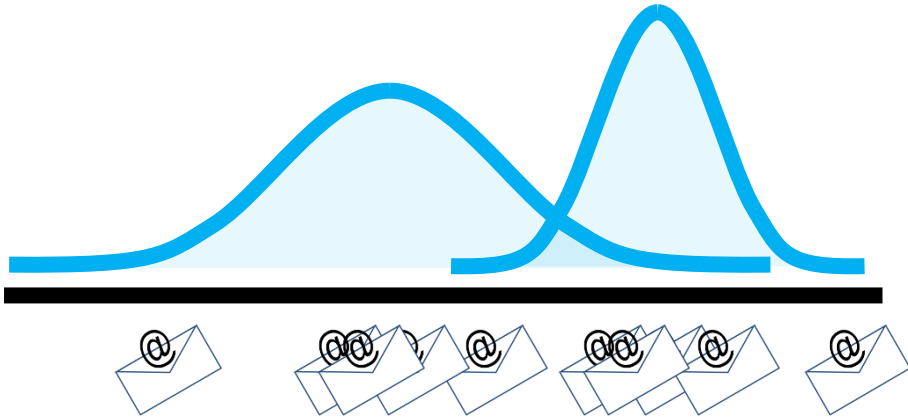




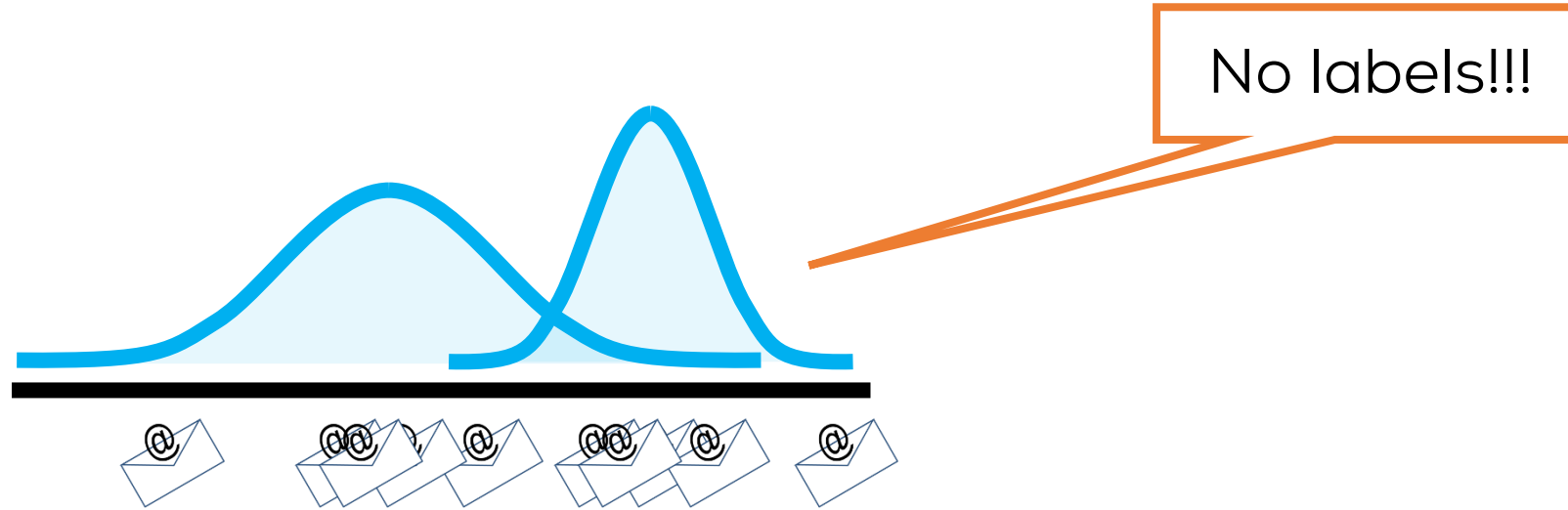
# The generative story for unlabelled data??



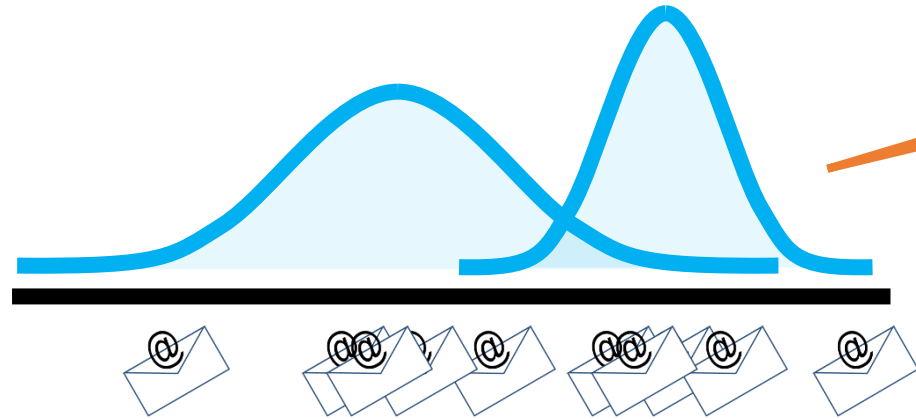
# The generative story for unlabelled data??



# The generative story for unlabelled data??



# The generative story for unlabelled data??



No labels!!!

Can we still recover the two Gaussian components in the mixture??

# How to grade a CS771 exam??

Sept 8, 2017



CS771: Intro to ML

# How to grade a CS771 exam??

 gradescope

# How to grade a CS771 exam??

 gradescope

Q1.  $\int x = ?$   
5 marks

# How to grade a CS771 exam??

 gradescope

Q1.  $\int x = ?$   
5 marks

Handwritten student answers for  $\int x = ?$ :

- $x^2 + d$
- $\frac{x^2}{2} + b$
- $\frac{x^2}{2}$
- $x^2/2$
- $x^2 + C$
- $\frac{x^2}{2} + C$
- $x^2/2 + a$
- $x^2/2$
- $x^2 + e$



# How to grade a CS771 exam??

 gradescope

Q1.  $\int x = ?$   
5 marks

$$x^2/2$$
$$x^2/2 \quad x^2/2$$

$$x^2/2 + a$$
$$\frac{x^2}{2} + b \frac{x^2}{2} + c$$

$$x^2 + C x^2 + d$$
$$x^2 + e$$

# How to grade a CS771 exam??

 gradescope

Q1.  $\int x = ?$   
5 marks

$\frac{x^2}{2} + a$   
 $\frac{x^2}{2} + b \frac{x^2}{2} + c$

$\frac{x^2}{2}$   
 $\frac{x^2}{2} \frac{x^2}{2}$

$x^2 + c x^2 + d$   
 $x^2 + e$

# How to grade a CS771 exam??

 gradescope

Q1.  $\int x = ?$   
5 marks

$x^2 / 2$   
4/5  
 $x^2 / 2$

$x^2 / 2 + a$   
5/5  
 $x^2 / 2 + b$   
 $x^2 / 2 + c$

$x^2 + c$   
2/5  
 $x^2 + d$   
 $x^2 + e$

# How to grade a CS771 exam??

 gradescope

Q1.  $\int x = ?$   
5 marks

Clustering!

Sept 8, 2017

$x^2/2$   
4/5  
 $x^2/2$

$x^2/2 + a$   
5/5  
 $x^2/2 + b$   
 $x^2/2 + c$

$x^2 + c$   
2/5  
 $x^2 + d$   
 $x^2 + e$



# How to grade a CS771 exam??

 gradescope

Q1.  $\int x = ?$   
5 marks

Clustering!

Like classification  
without labels ☺

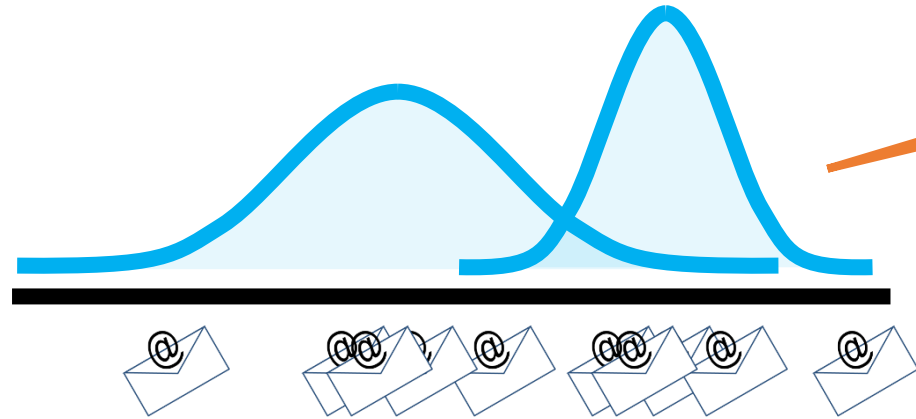
Sept 8, 2017

$x^2/2 + a$   
 $x^2/2 + b$   
 $x^2/2 + c$

$x^2/2 + a$   
 $x^2/2 + b$   
 $x^2/2 + c$

$x^2 + c$   
 $x^2 + d$   
 $x^2 + e$

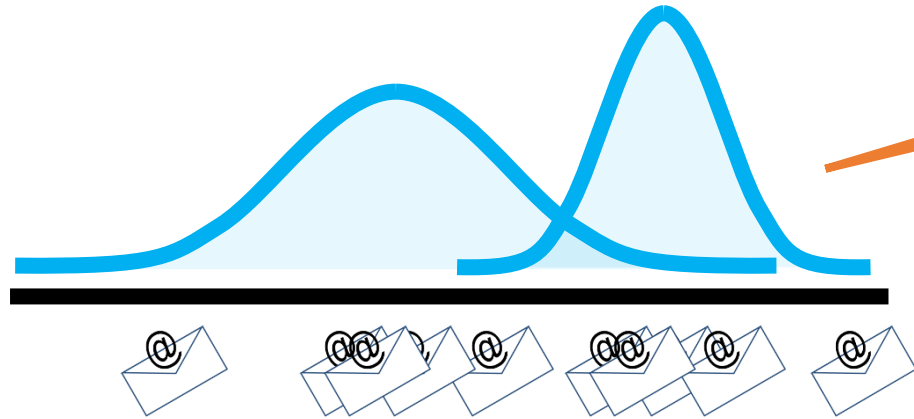
# The generative story for unlabelled data??



No labels!!!

Can we still recover the two Gaussian components in the mixture??

# The generative story for unlabelled data??

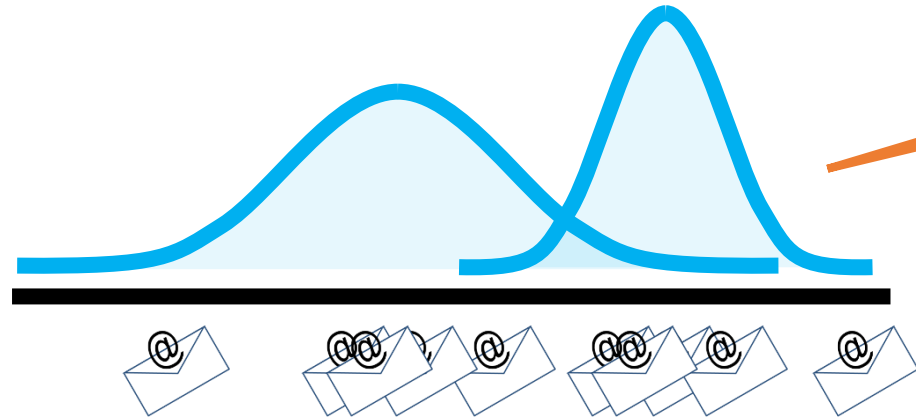


No labels!!!

Can we still recover the two Gaussian components in the mixture??

How do we know there are only two Gaussians?

# The generative story for unlabelled data??



No labels!!!

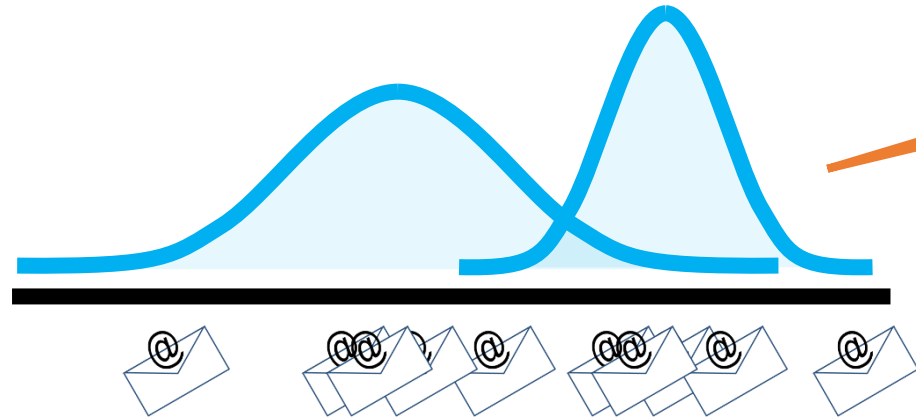
# Gaussians can be  
a hyper-parameter  
 $K$  you can tune

Can we still recover  
the two Gaussian  
components in the  
mixture??

How do we know there  
are only two Gaussians?



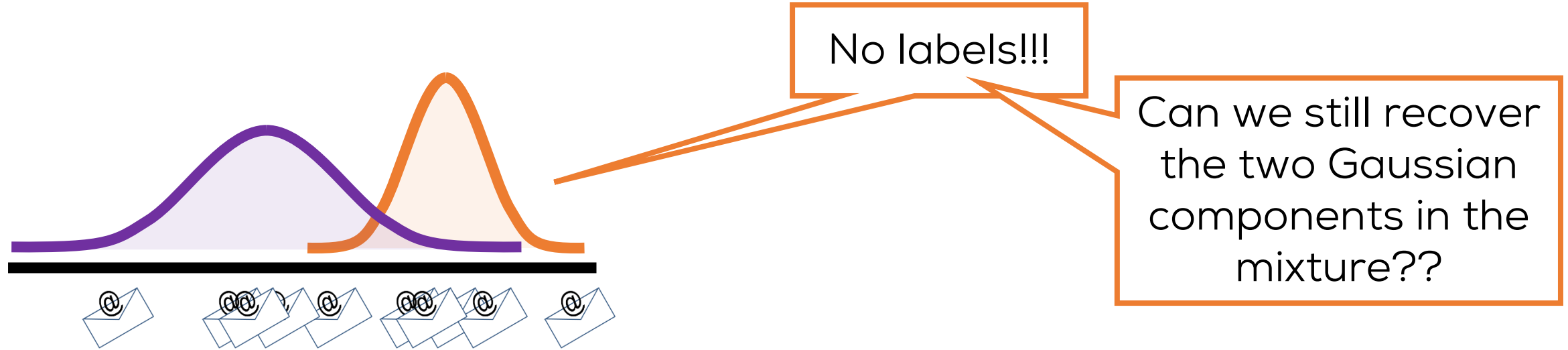
# The generative story for unlabelled data??



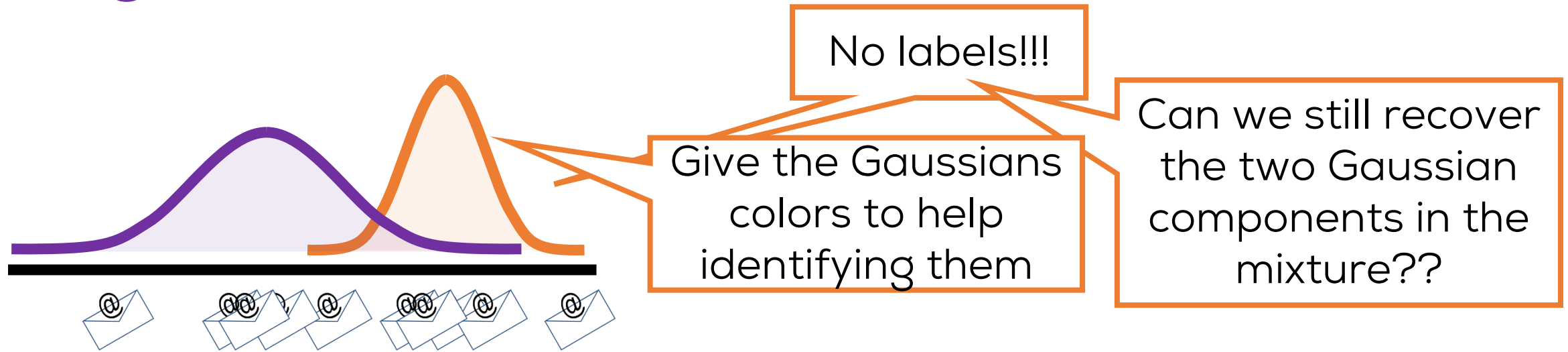
No labels!!!

Can we still recover the two Gaussian components in the mixture??

# The generative story for unlabelled data??



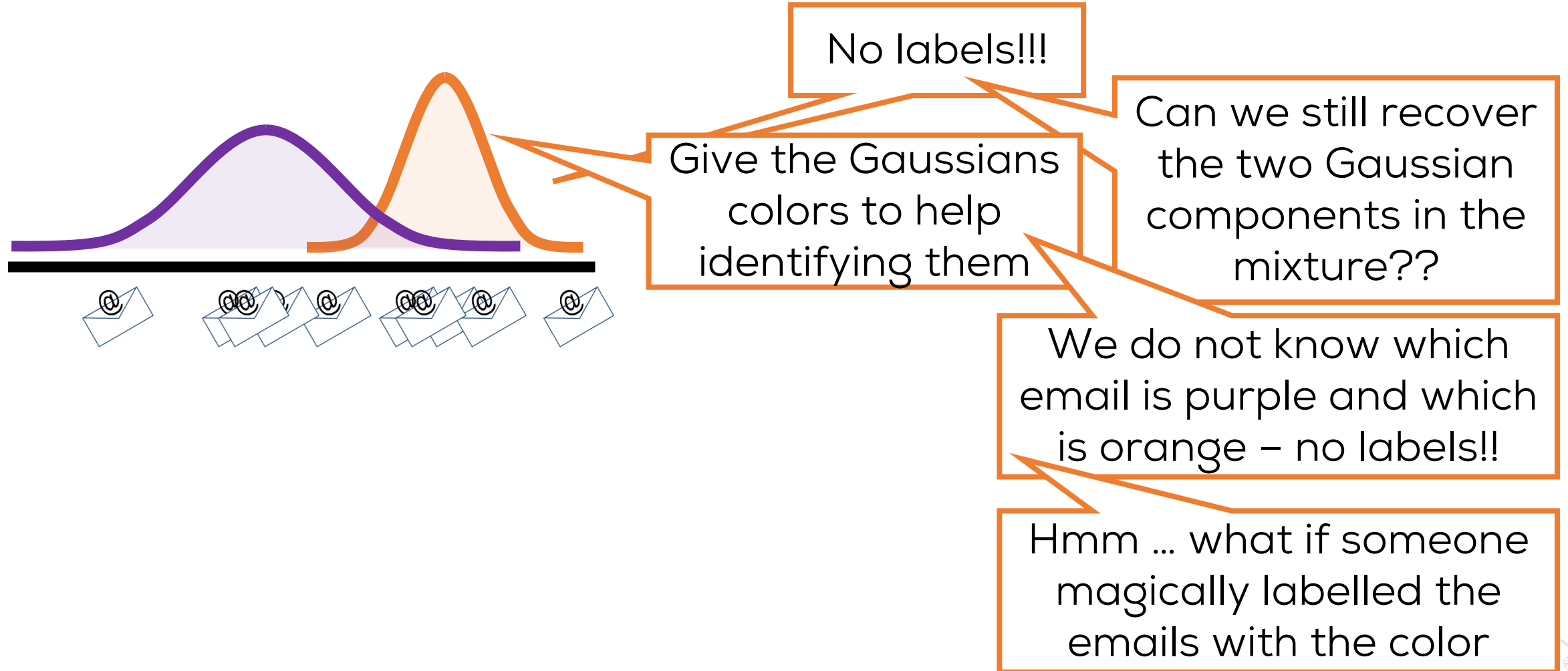
# The generative story for unlabelled data??



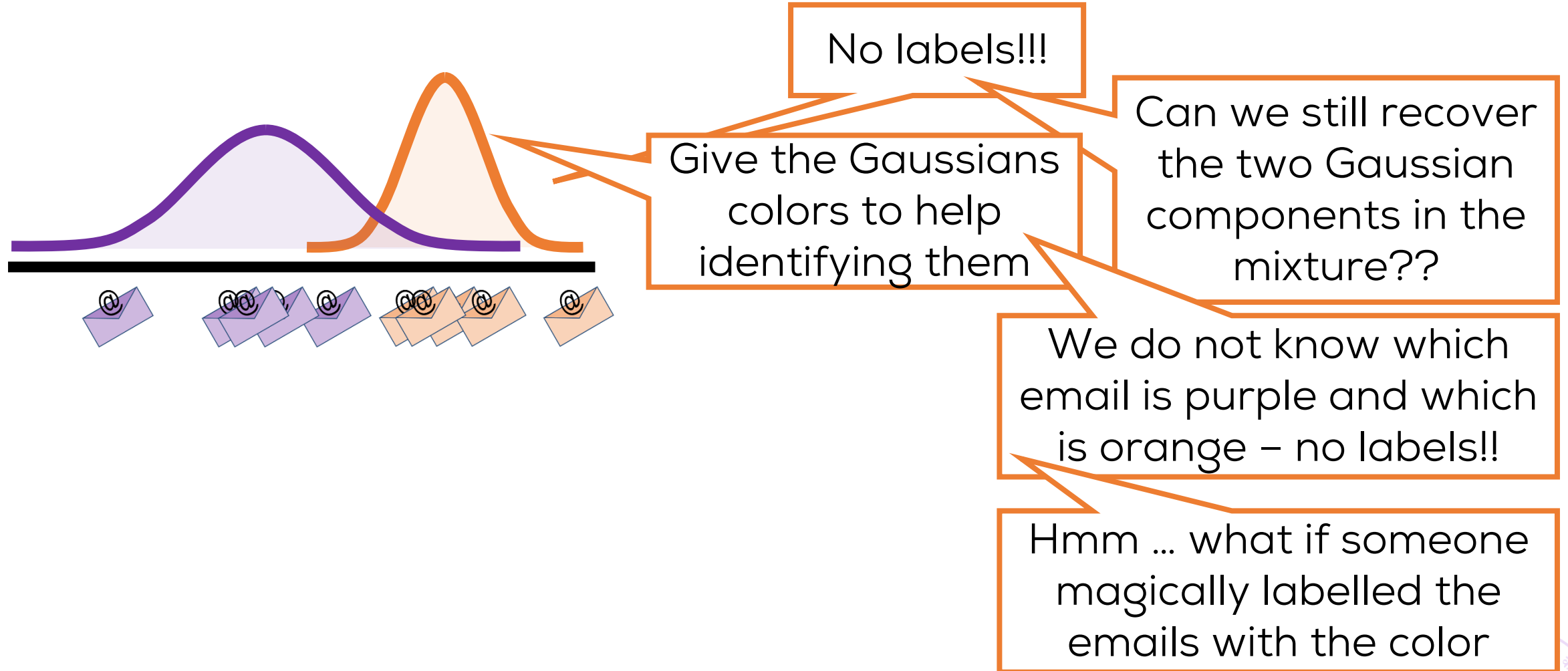
# The generative story for unlabelled data??



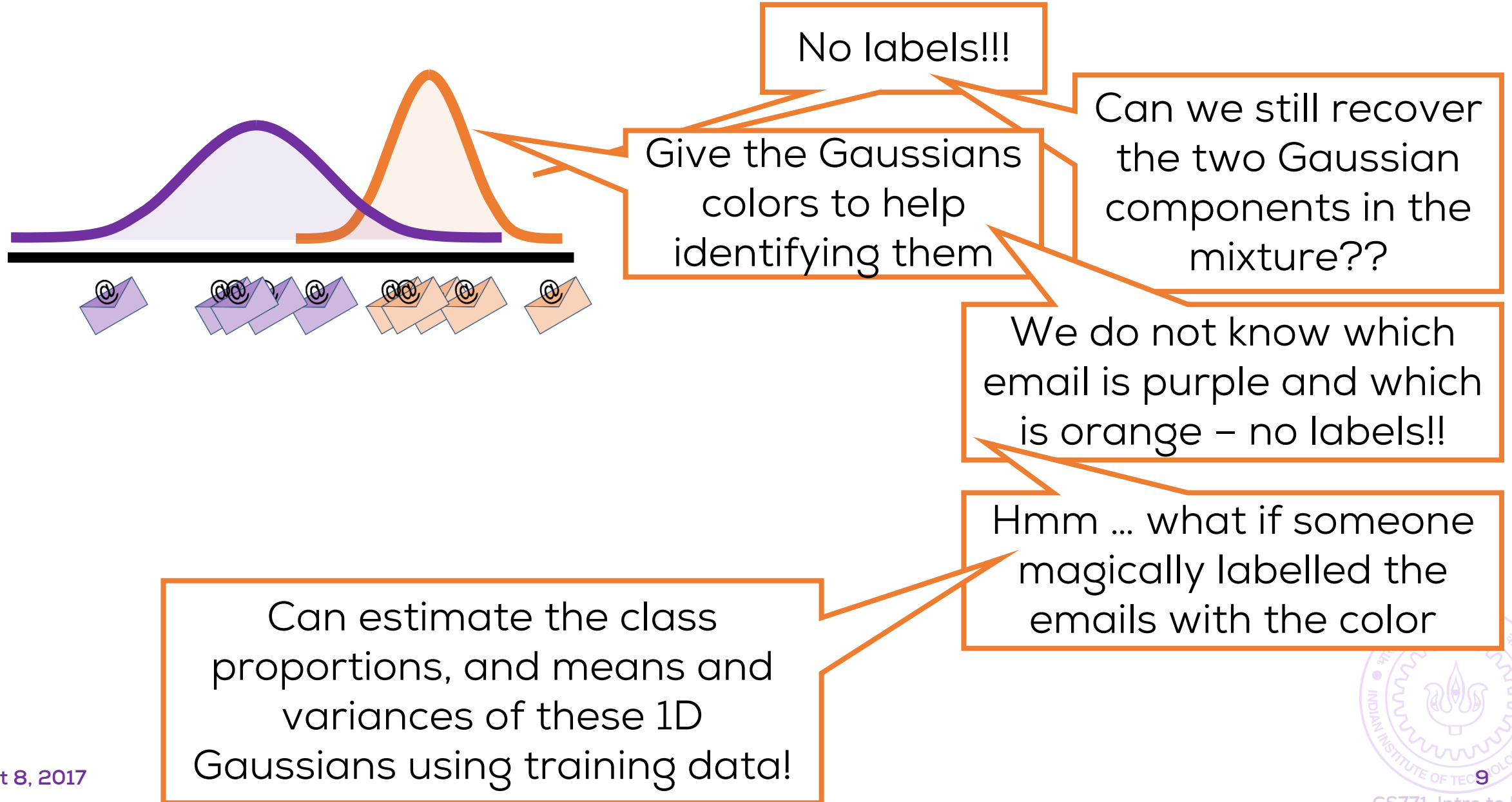
# The generative story for unlabelled data??



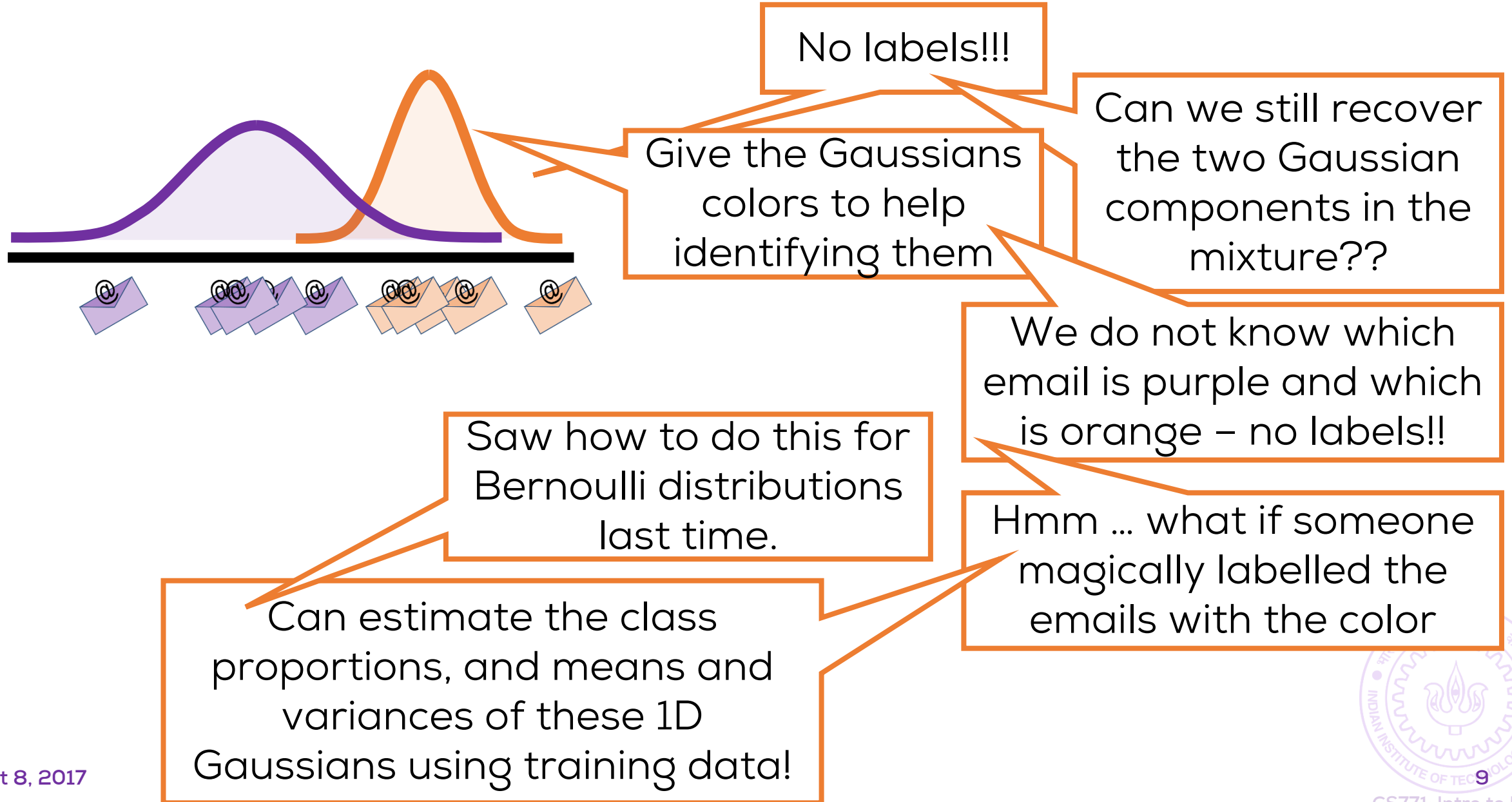
# The generative story for unlabelled data??



# The generative story for unlabelled data??

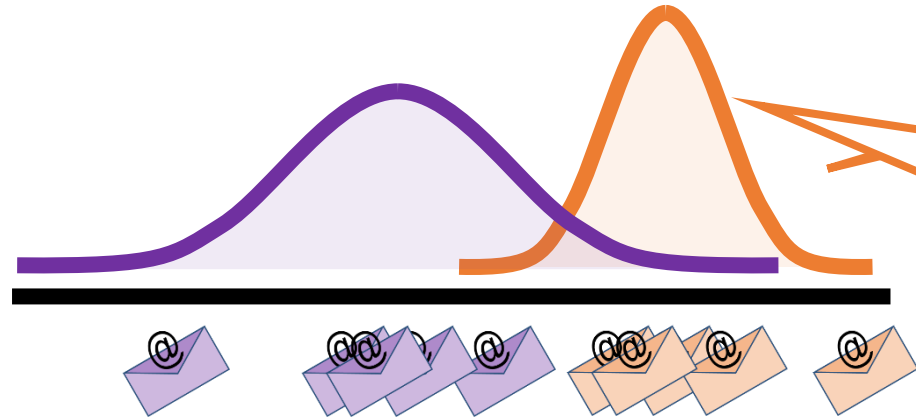


# The generative story for unlabelled data??





# The generative story for unlabelled data??



No labels!!!

Give the Gaussians colors to help identifying them

Can we still recover the two Gaussian components in the mixture??

We do not know which email is purple and which is orange – no labels!!

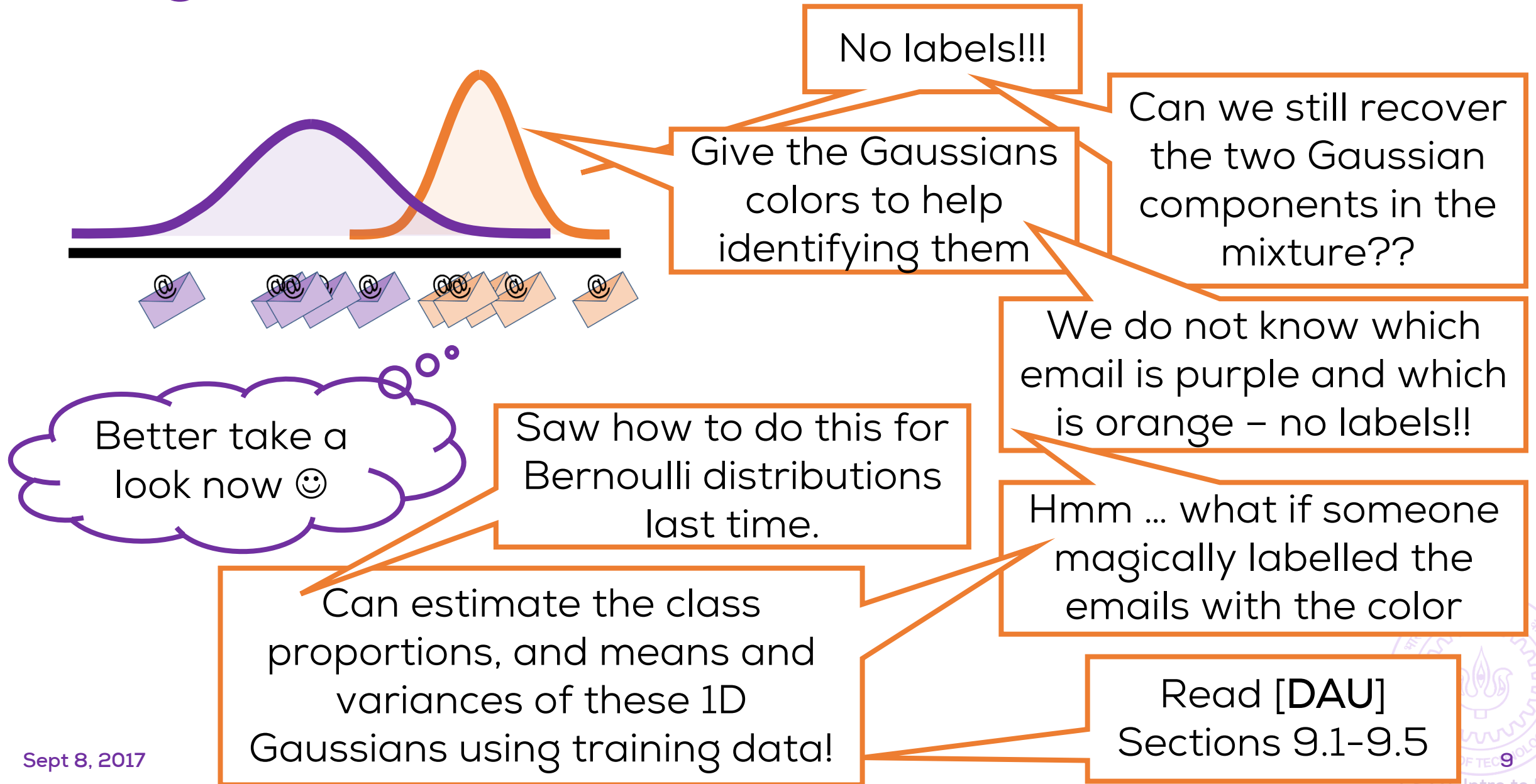
Saw how to do this for Bernoulli distributions last time.

Hmm ... what if someone magically labelled the emails with the color

Can estimate the class proportions, and means and variances of these 1D Gaussians using training data!

Read [DAU]  
Sections 9.1-9.5

# The generative story for unlabelled data??

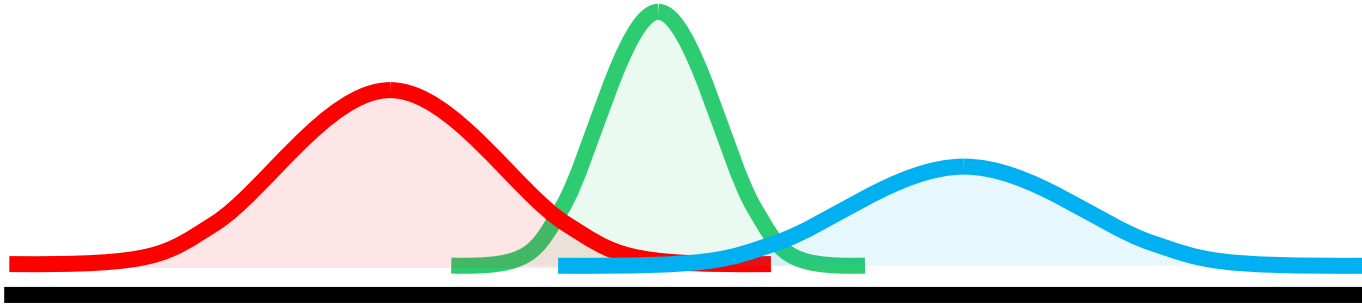


# <detour>

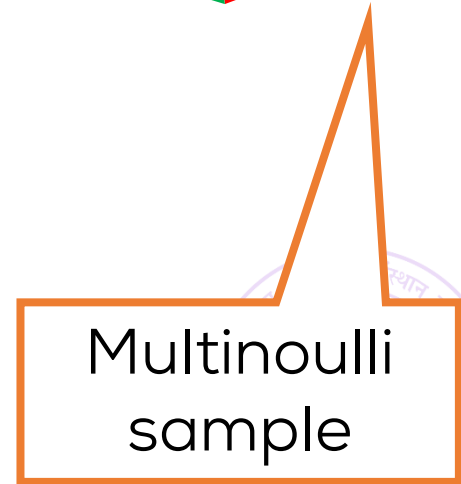
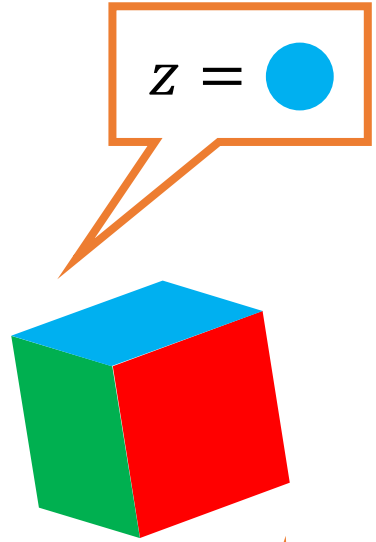
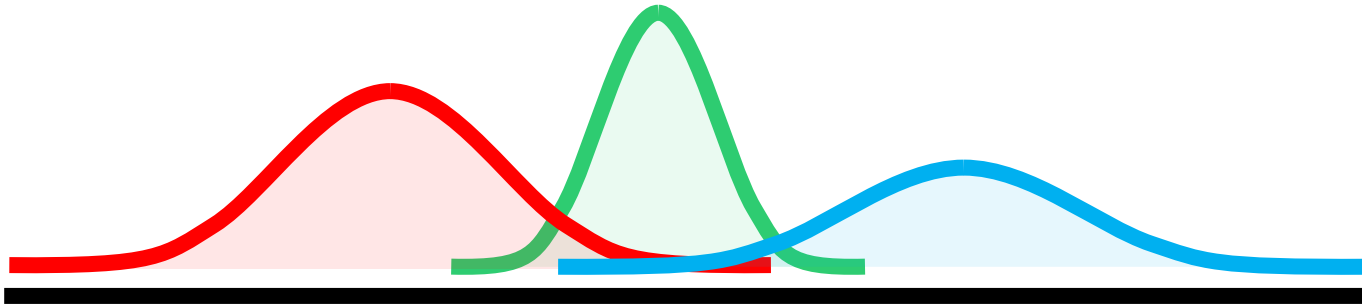
Learning a mixture of Gaussians in presence of labels

# The generative story for labelled data

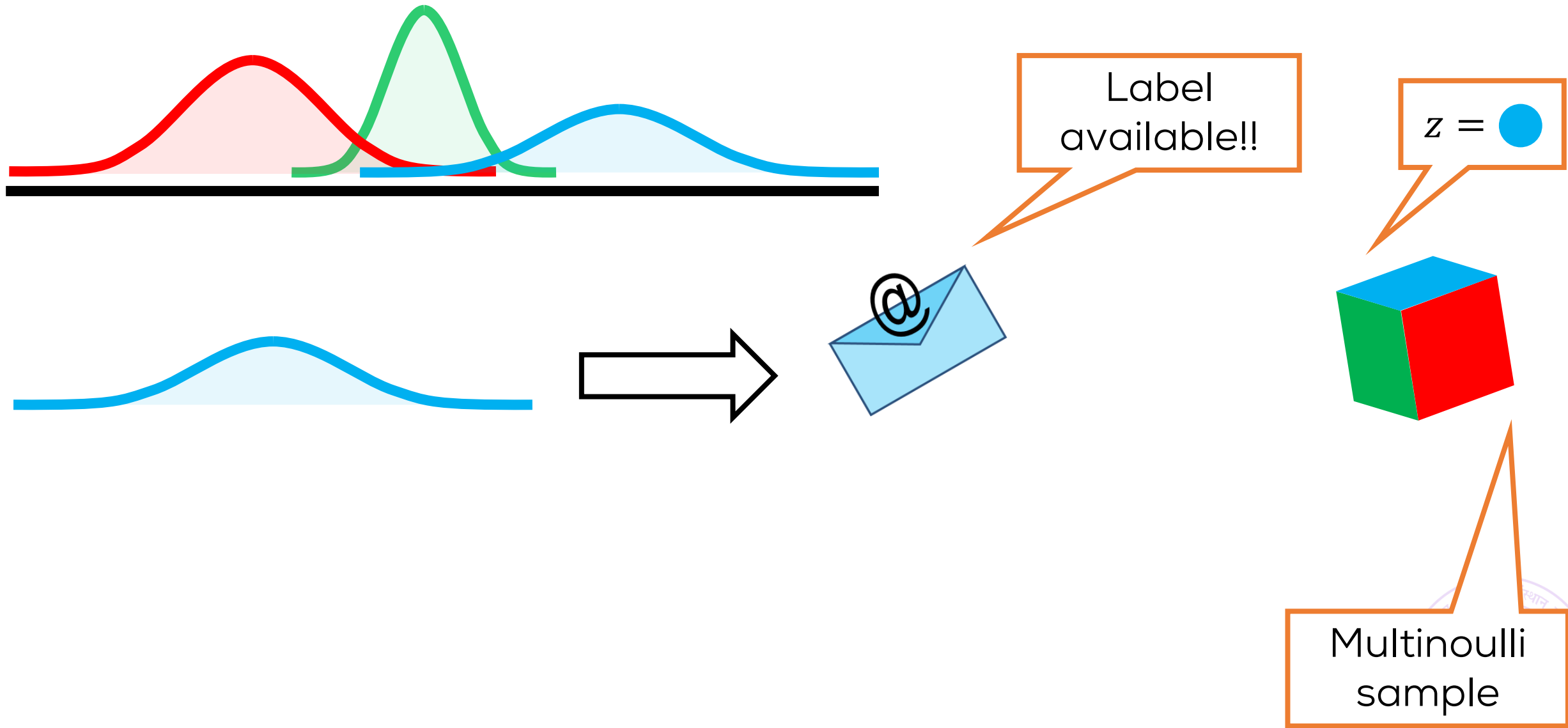
# The generative story for labelled data



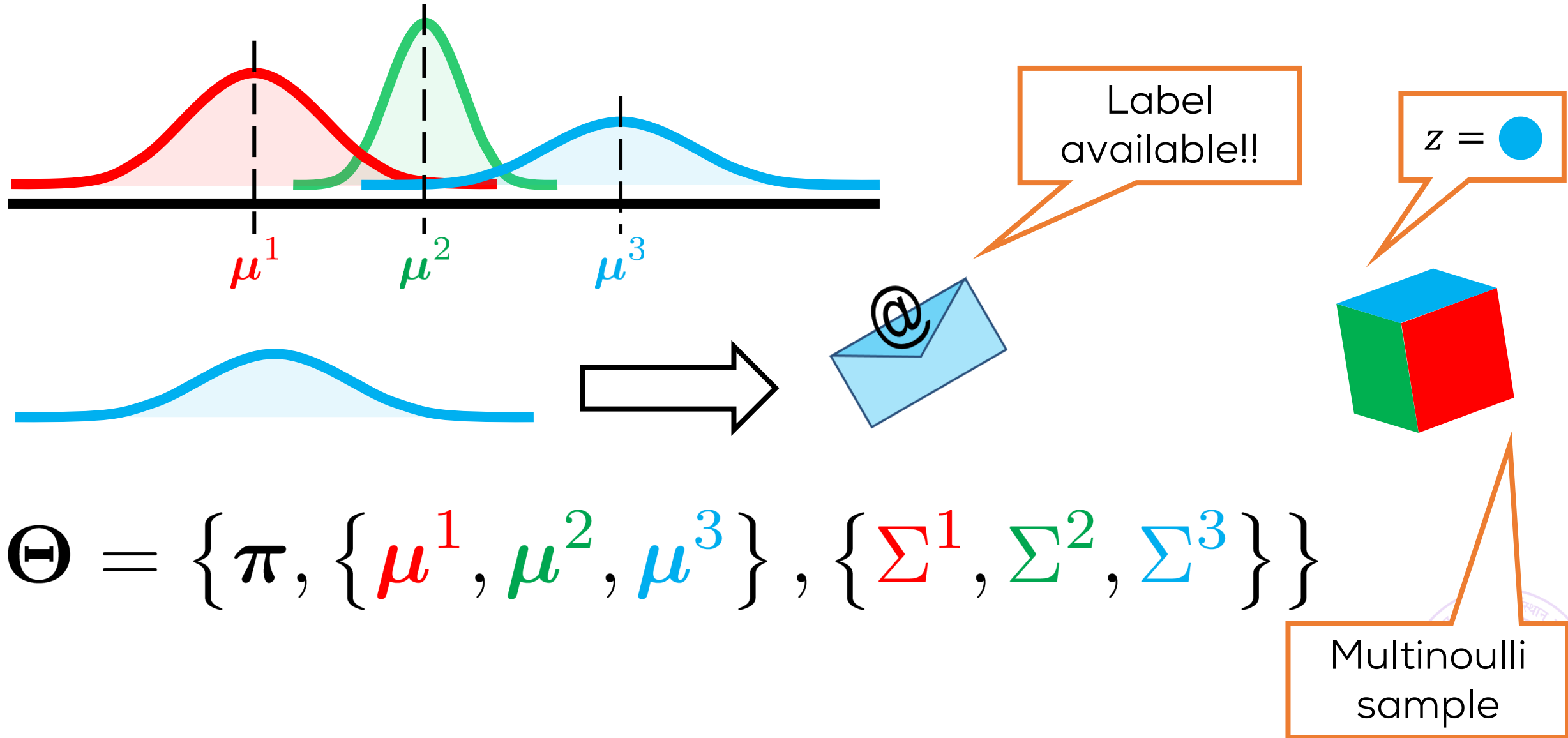
# The generative story for labelled data



# The generative story for labelled data

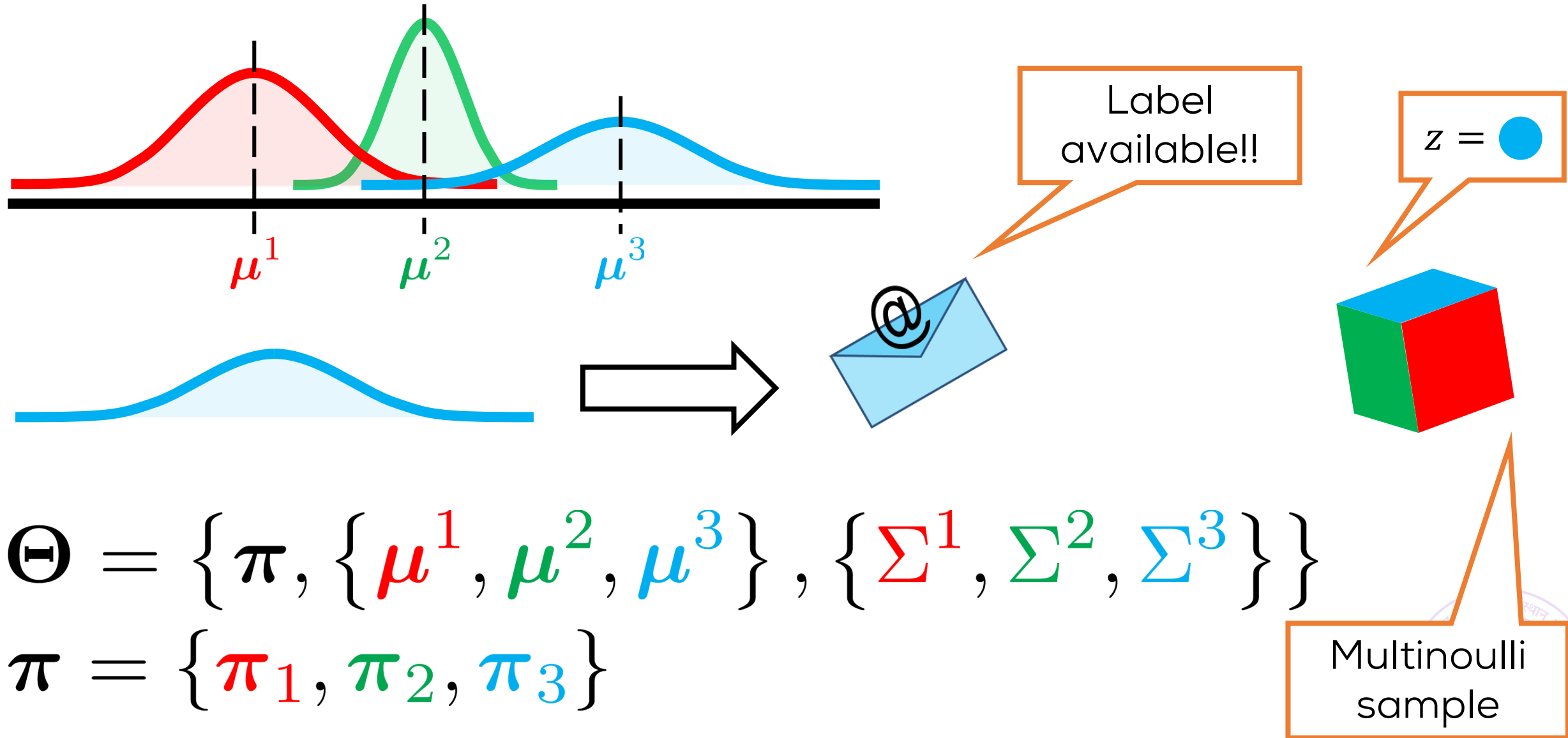


# The generative story for labelled data

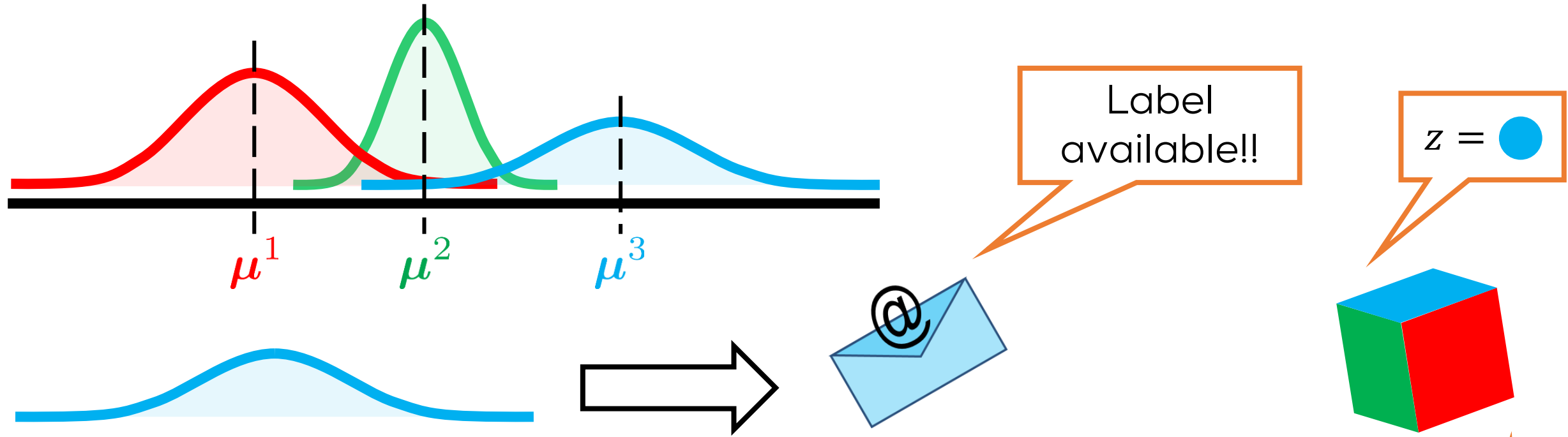




# The generative story for labelled data



# The generative story for labelled data

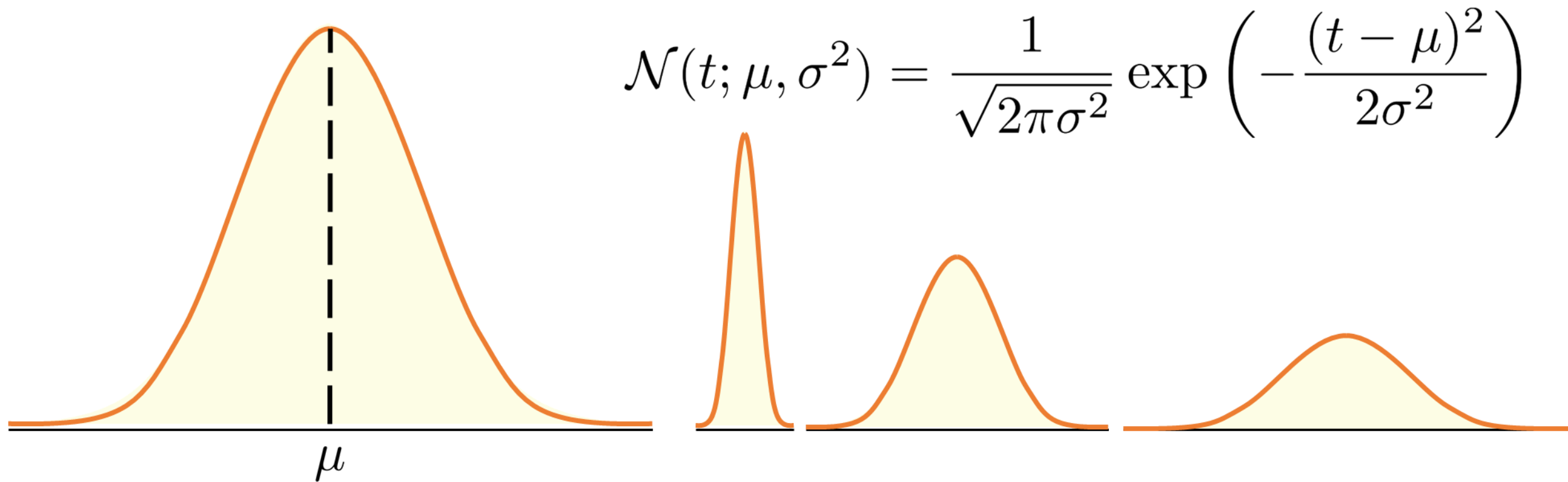


$$\Theta = \left\{ \pi, \left\{ \mu^1, \mu^2, \mu^3 \right\}, \left\{ \Sigma^1, \Sigma^2, \Sigma^3 \right\} \right\}$$

$$\pi = \left\{ \pi_1, \pi_2, \pi_3 \right\} \quad \pi_1 = \mathbb{P} \left[ z = \bullet \right]$$

Multinoulli  
sample

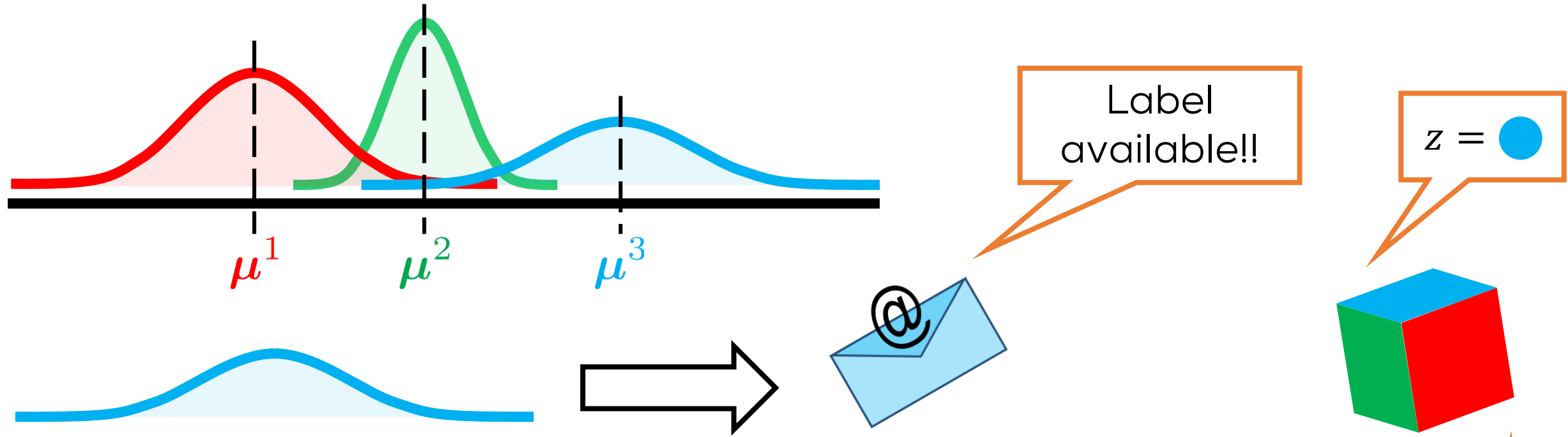
# The Gaussian Distribution



**Multivariate  
Gaussian**

$$\mathcal{N}(\mathbf{v}; \boldsymbol{\mu}, \Sigma) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2}(\mathbf{v} - \boldsymbol{\mu})^\top \Sigma^{-1}(\mathbf{v} - \boldsymbol{\mu})\right)$$

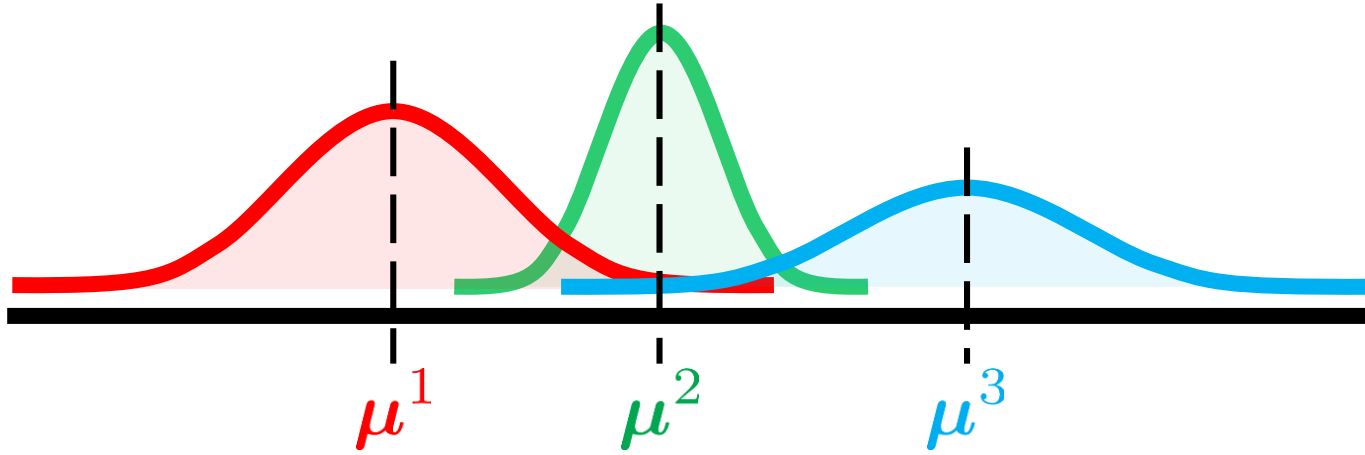
# The generative story for labelled data



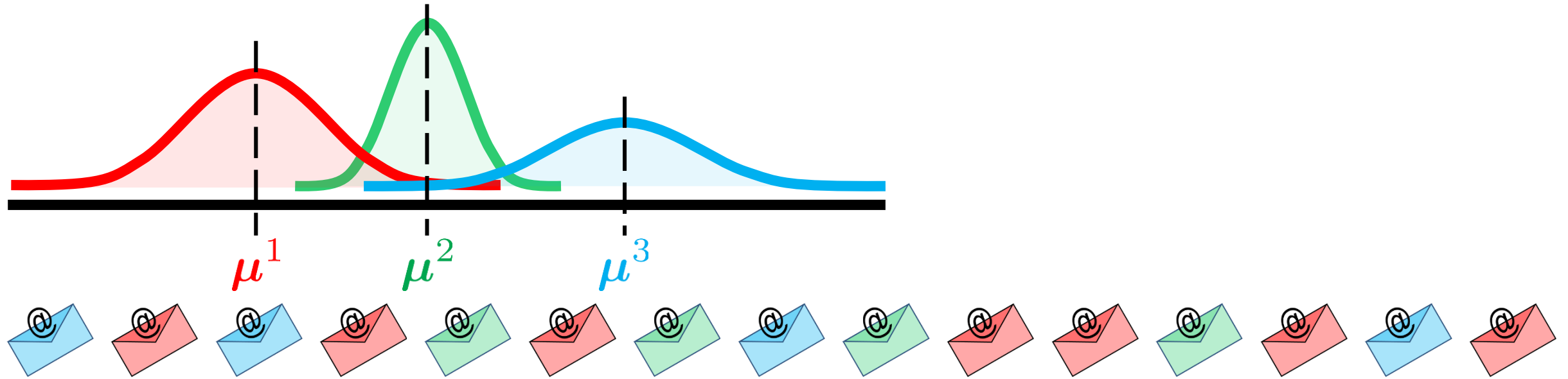
$$\Theta = \left\{ \pi, \left\{ \mu^1, \mu^2, \mu^3 \right\}, \left\{ \Sigma^1, \Sigma^2, \Sigma^3 \right\} \right\}$$

$$\pi = \left\{ \pi_1, \pi_2, \pi_3 \right\} \quad \pi_1 = \mathbb{P} \left[ z = \bullet \right]$$

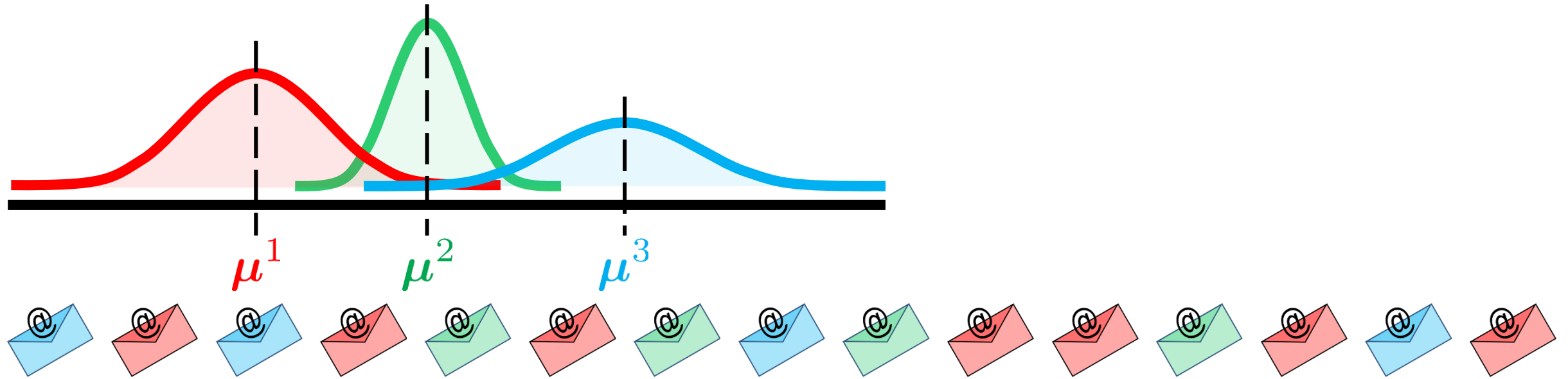
# The generative story for labelled data



# The generative story for labelled data

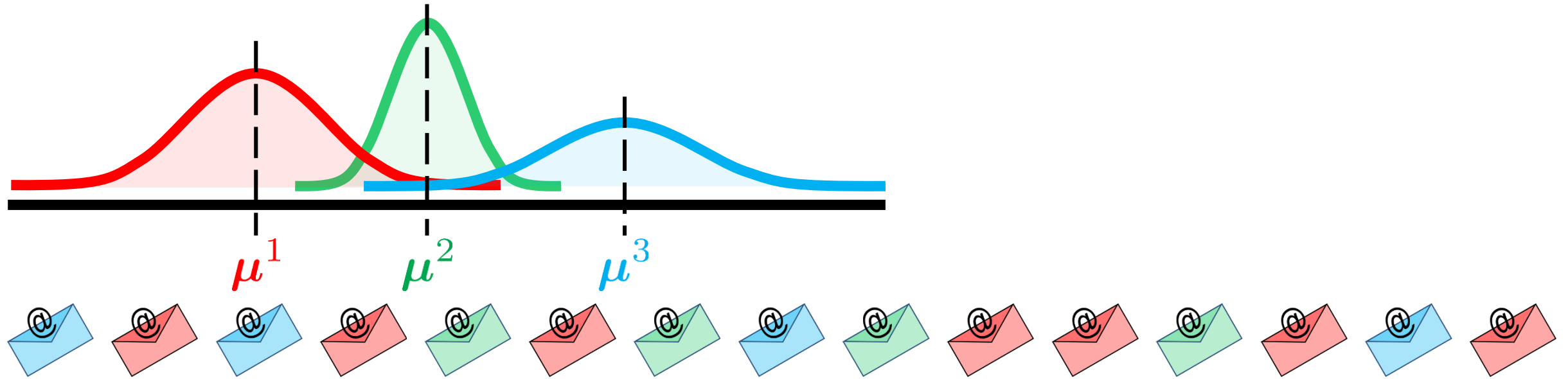


# The generative story for labelled data



$$\mathbb{P}[\mathbf{x}^i, z^i | \Theta] = \mathbb{P}[z^i | \Theta] \cdot \mathbb{P}[\mathbf{x}^i | z^i, \Theta] = \pi_{z^i} \cdot \mathcal{N}(\mathbf{x}^i | \mu^{z^i}, \Sigma^{z^i})$$

# The generative story for labelled data

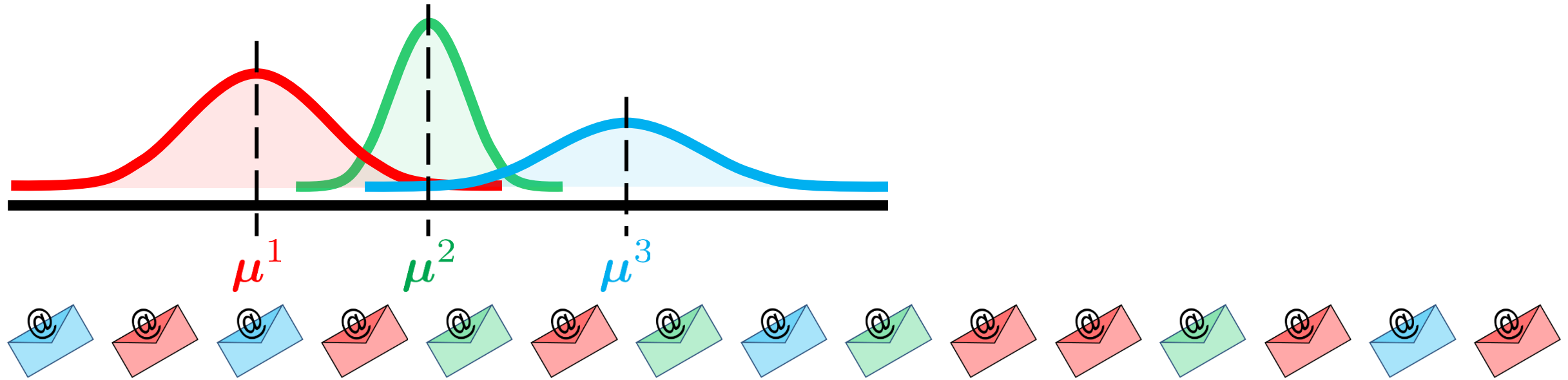


$$\mathbb{P}[\mathbf{x}^i, z^i \mid \Theta] = \mathbb{P}[z^i \mid \Theta] \cdot \mathbb{P}[\mathbf{x}^i \mid z^i, \Theta] = \pi_{z^i} \cdot \mathcal{N}(\mathbf{x}^i \mid \mu^{z^i}, \Sigma^{z^i})$$

$$\mathbb{P}[X, \{z^i\} \mid \Theta] = \prod_{i=1}^n \mathbb{P}[\mathbf{x}^i, z^i \mid \Theta]$$



# The generative story for labelled data

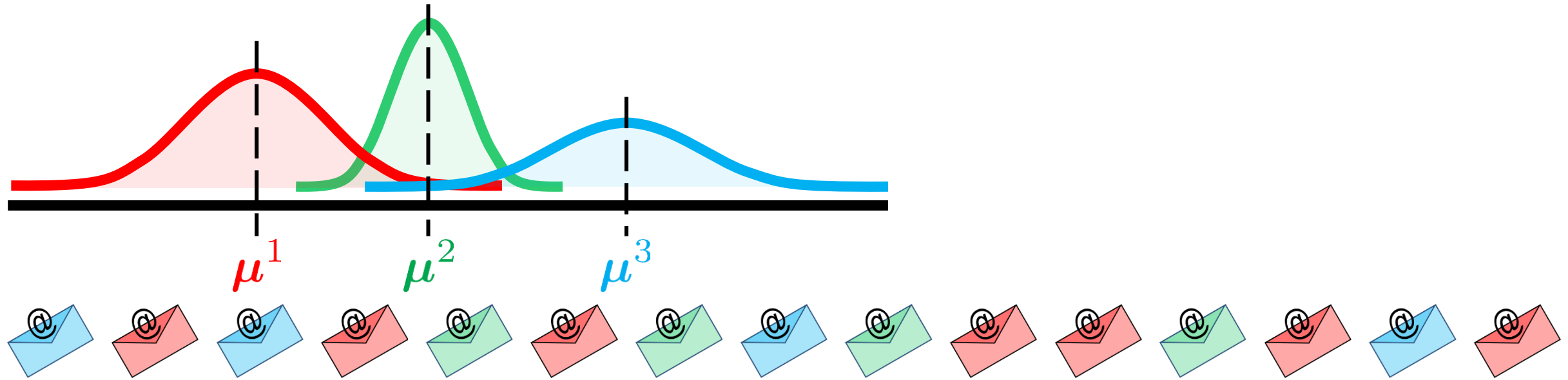


$$\mathbb{P}[\mathbf{x}^i, z^i \mid \Theta] = \mathbb{P}[z^i \mid \Theta] \cdot \mathbb{P}[\mathbf{x}^i \mid z^i, \Theta] = \pi_{z^i} \cdot \mathcal{N}(\mathbf{x}^i \mid \mu^{z^i}, \Sigma^{z^i})$$

$$\mathbb{P}[X, \{z^i\} \mid \Theta] = \prod_{i=1}^n \mathbb{P}[\mathbf{x}^i, z^i \mid \Theta]$$

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X, \{z^i\} \mid \Theta]$$

# The generative story for labelled data



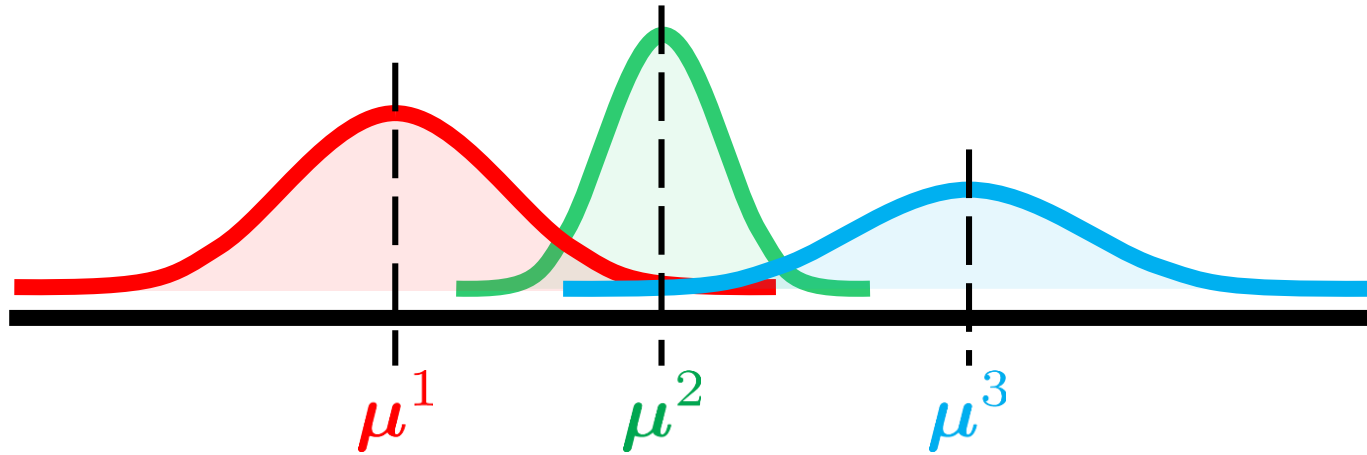
$$\mathbb{P}[\mathbf{x}^i, z^i | \Theta] = \mathbb{P}[z^i | \Theta] \cdot \mathbb{P}[\mathbf{x}^i | z^i, \Theta] = \pi_{z^i} \cdot \mathcal{N}(\mathbf{x}^i | \mu^{z^i}, \Sigma^{z^i})$$

$$\mathbb{P}[X, \{z^i\} | \Theta] = \prod_{i=1}^n \mathbb{P}[\mathbf{x}^i, z^i | \Theta]$$

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X, \{z^i\} | \Theta]$$

Take log and apply 1<sup>st</sup> order optimality

# The generative story for labelled data



$$\mathbb{P}[\mathbf{x}^i, z^i | \Theta] = \mathbb{P}[z^i | \Theta] \cdot \mathbb{P}[\mathbf{x}^i | z^i, \Theta] = \pi_{z^i} \cdot \mathcal{N}(\mathbf{x}^i | \boldsymbol{\mu}^{z^i}, \Sigma^{z^i})$$

$$\mathbb{P}[X, \{z^i\} | \Theta] = \prod_{i=1}^n \mathbb{P}[\mathbf{x}^i, z^i | \Theta]$$

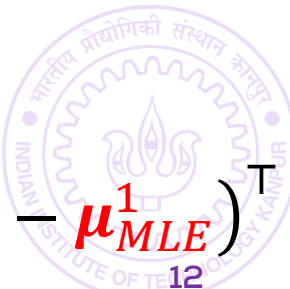
$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X, \{z^i\} | \Theta]$$

Take log and apply 1<sup>st</sup> order optimality

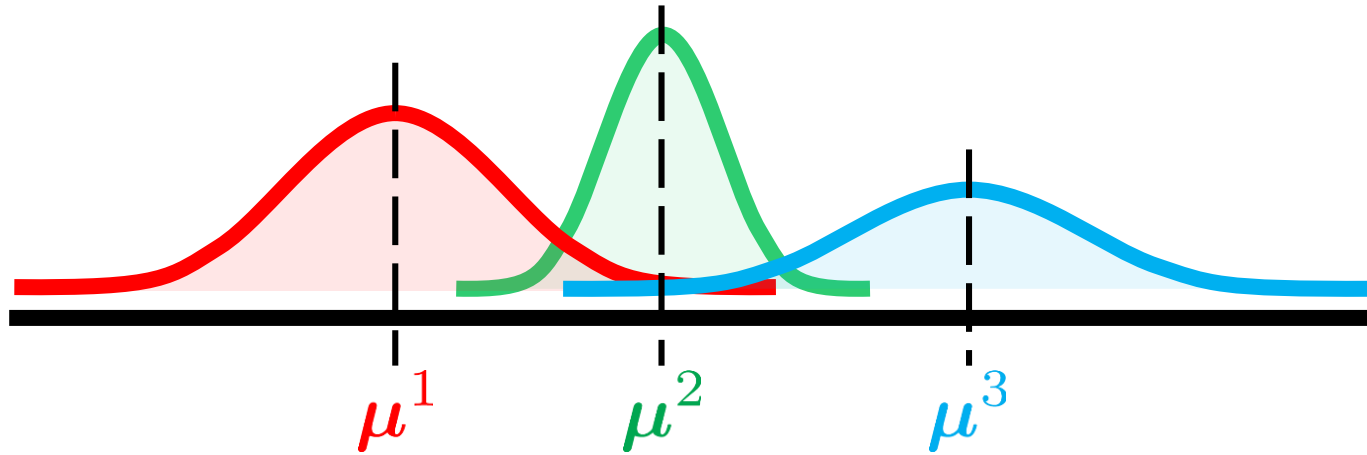
$$\pi_1^{\text{MLE}} = \frac{\# \text{ red emails}}{\# \text{ total emails}} = \frac{n_r}{n}$$

$$\boldsymbol{\mu}_{\text{MLE}}^1 = \frac{1}{n_r} \sum_{i: z^i = \bullet} \mathbf{x}^i$$

$$\Sigma_{\text{MLE}}^1 = \frac{1}{n_r} \sum_{i: z^i = \bullet} (\mathbf{x}^i - \boldsymbol{\mu}_{\text{MLE}}^1)(\mathbf{x}^i - \boldsymbol{\mu}_{\text{MLE}}^1)^\top$$



# The generative story for labelled data



$$\mathbb{P}[\mathbf{x}^i, z^i | \Theta] = \mathbb{P}[z^i | \Theta] \cdot \mathbb{P}[\mathbf{x}^i | z^i, \Theta] = \pi_{z^i} \cdot \mathcal{N}(\mathbf{x}^i | \boldsymbol{\mu}^{z^i}, \Sigma^{z^i})$$

$$\mathbb{P}[X, \{z^i\} | \Theta] = \prod_{i=1}^n \mathbb{P}[\mathbf{x}^i, z^i | \Theta]$$

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X, \{z^i\} | \Theta]$$

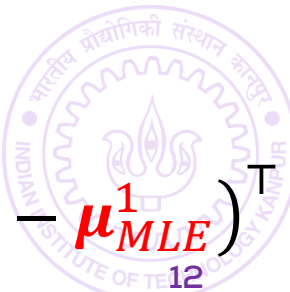
Take log and apply 1<sup>st</sup> order optimality

Read [DAU]  
Sections 9.1-9.5

$$\pi_1^{\text{MLE}} = \frac{\# \text{ red emails}}{\# \text{ total emails}} = \frac{n_r}{n}$$

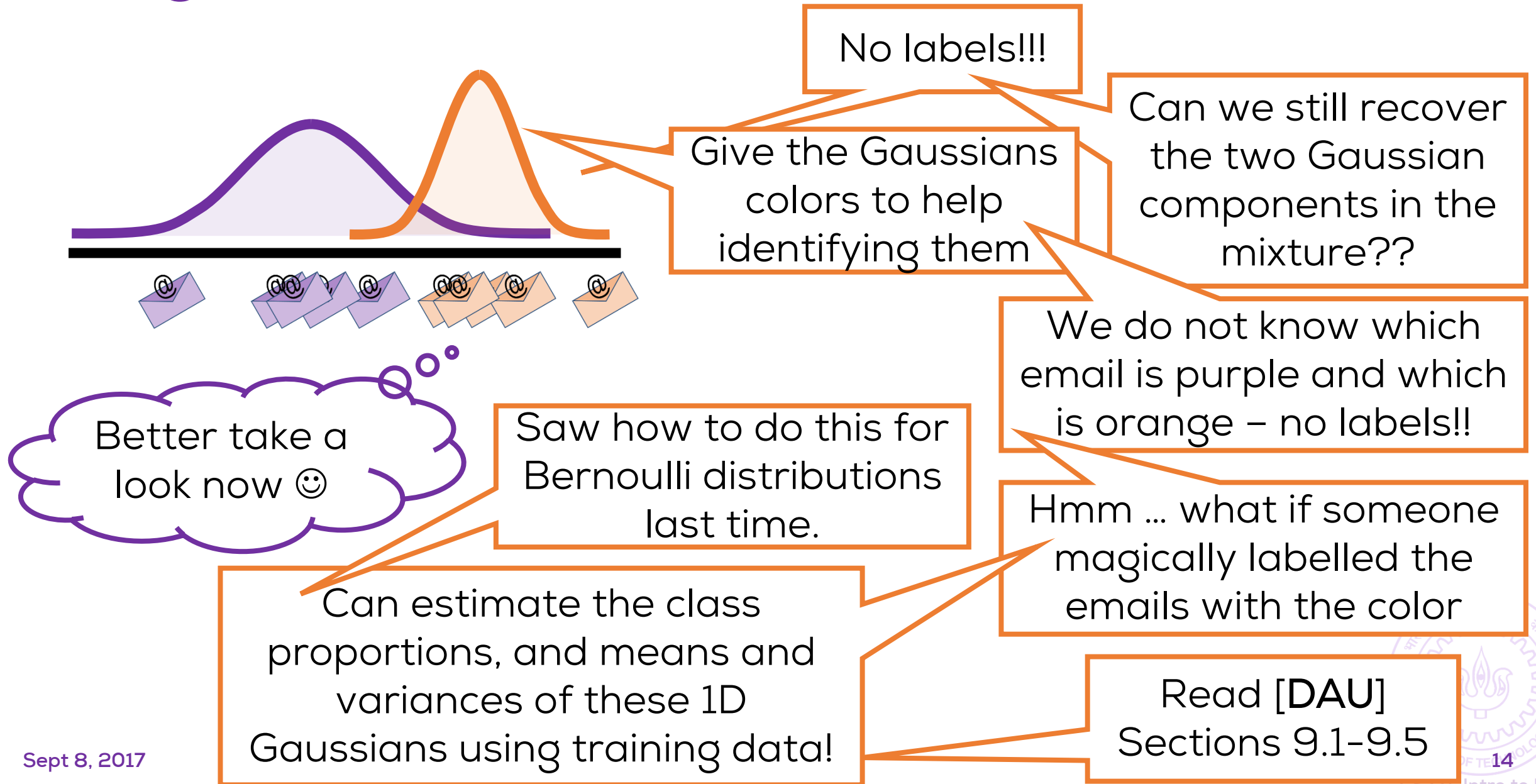
$$\boldsymbol{\mu}_{\text{MLE}}^1 = \frac{1}{n_r} \sum_{i: z^i = \bullet} \mathbf{x}^i$$

$$\Sigma_{\text{MLE}}^1 = \frac{1}{n_r} \sum_{i: z^i = \bullet} (\mathbf{x}^i - \boldsymbol{\mu}_{\text{MLE}}^1)(\mathbf{x}^i - \boldsymbol{\mu}_{\text{MLE}}^1)^\top$$



</detour>

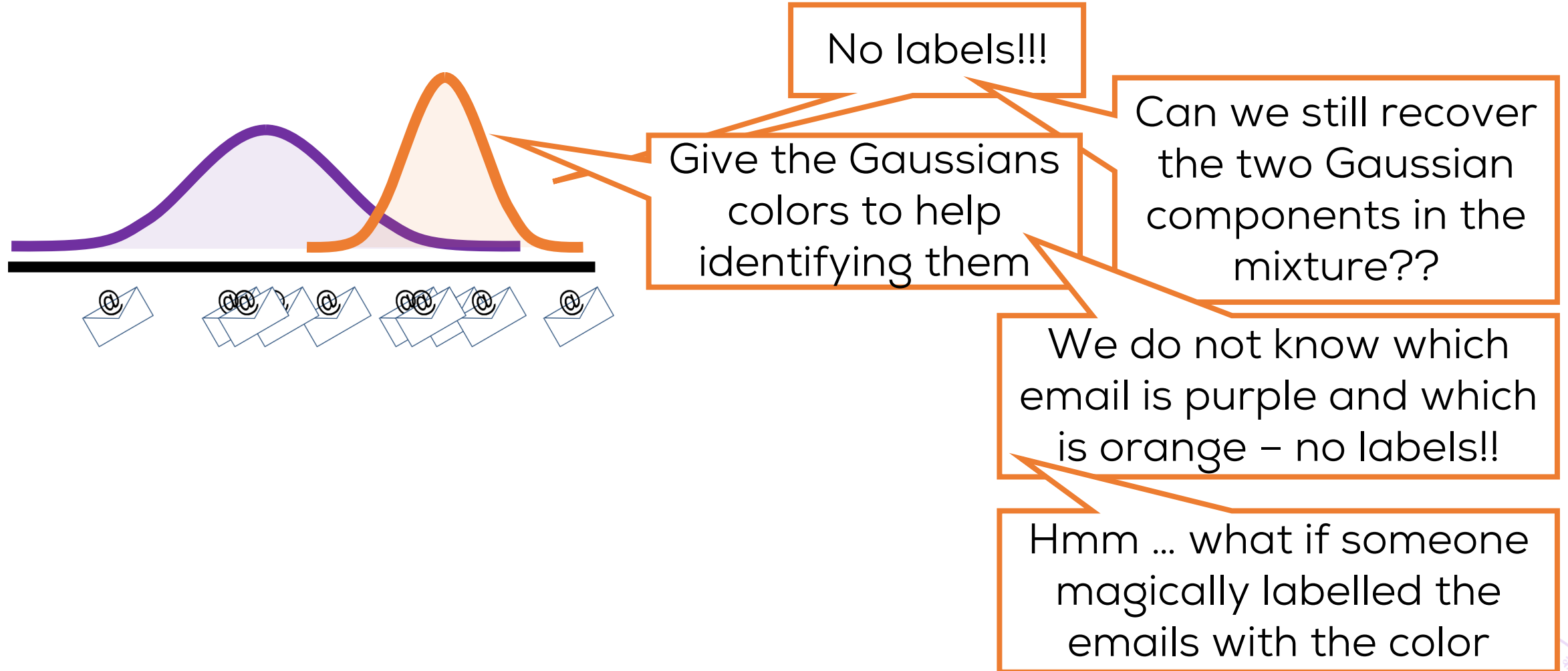
# The generative story for unlabelled data??



# The generative story for unlabelled data??

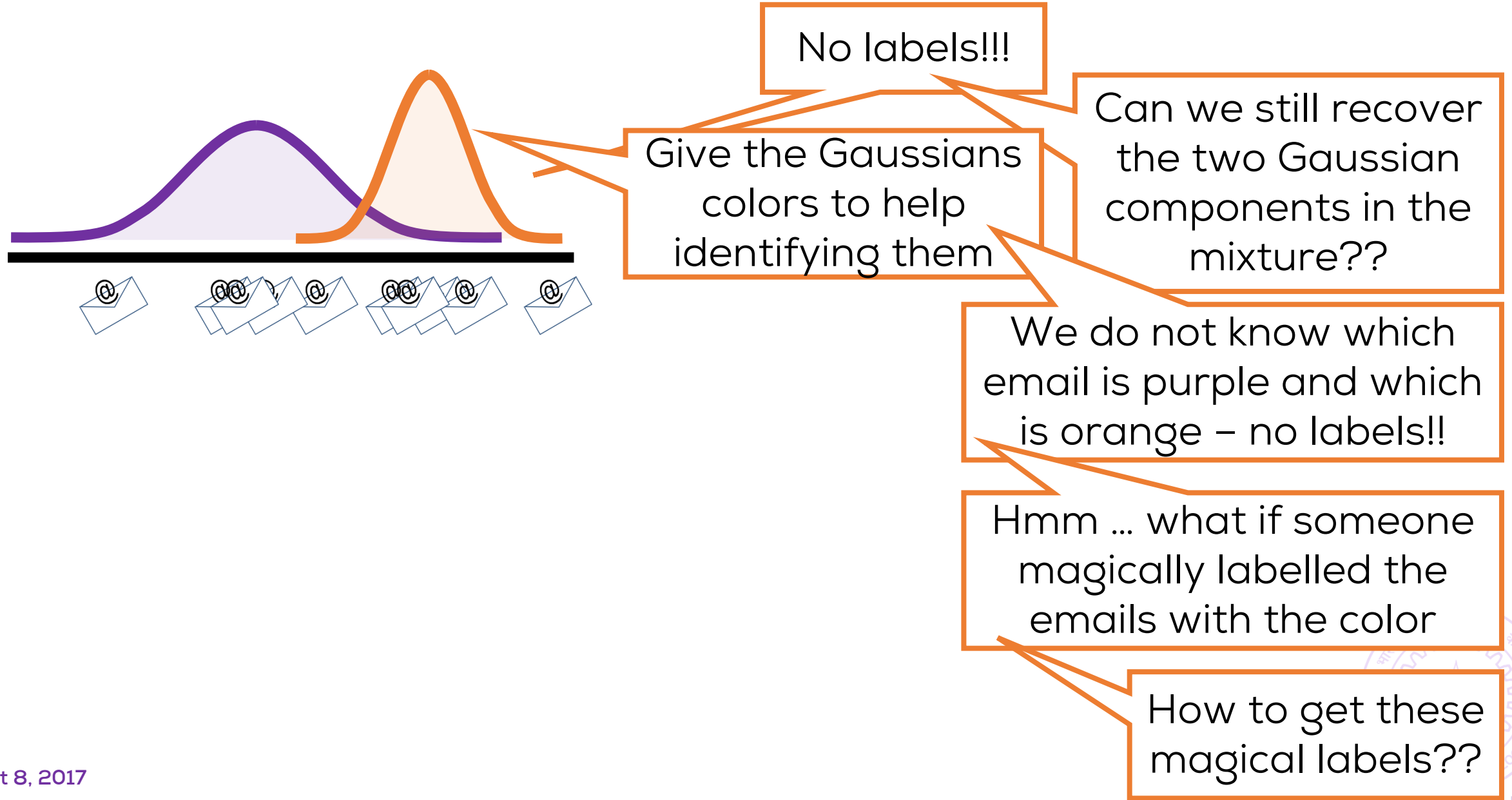


# The generative story for unlabelled data??

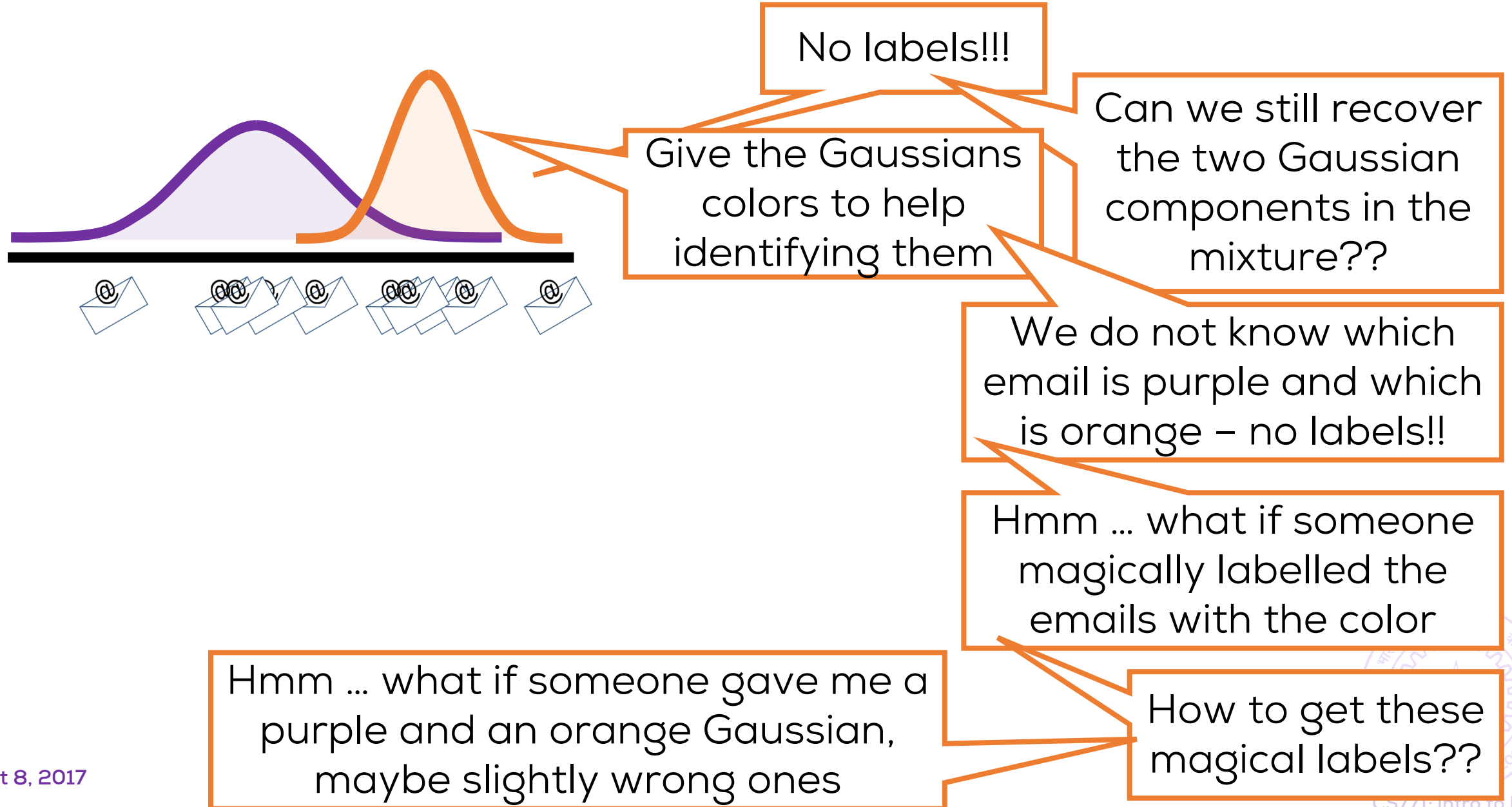




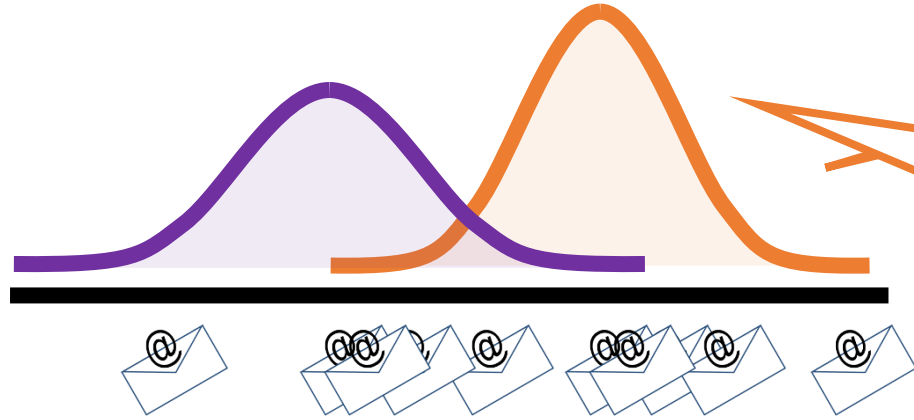
# The generative story for unlabelled data??



# The generative story for unlabelled data??



# The generative story for unlabelled data??



No labels!!!

Give the Gaussians colors to help identifying them

Can we still recover the two Gaussian components in the mixture??

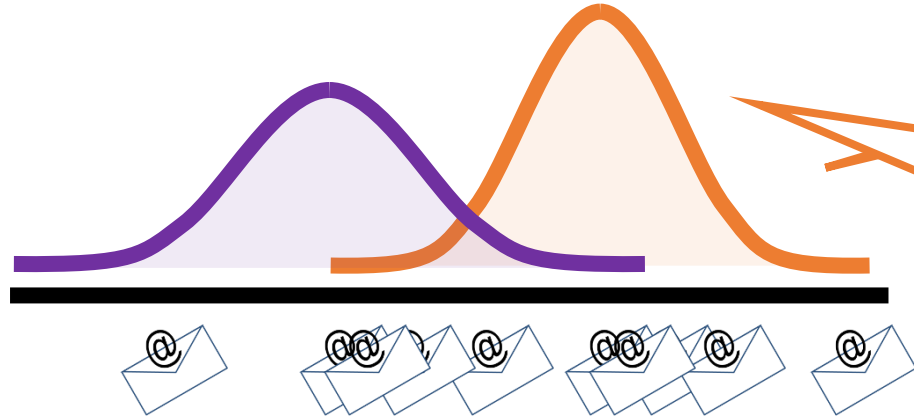
We do not know which email is purple and which is orange – no labels!!

Hmm ... what if someone magically labelled the emails with the color

Hmm ... what if someone gave me a purple and an orange Gaussian, maybe slightly wrong ones

How to get these magical labels??

# The generative story for unlabelled data??



No labels!!!

Give the Gaussians colors to help identifying them

Can we still recover the two Gaussian components in the mixture??

We do not know which email is purple and which is orange – no labels!!

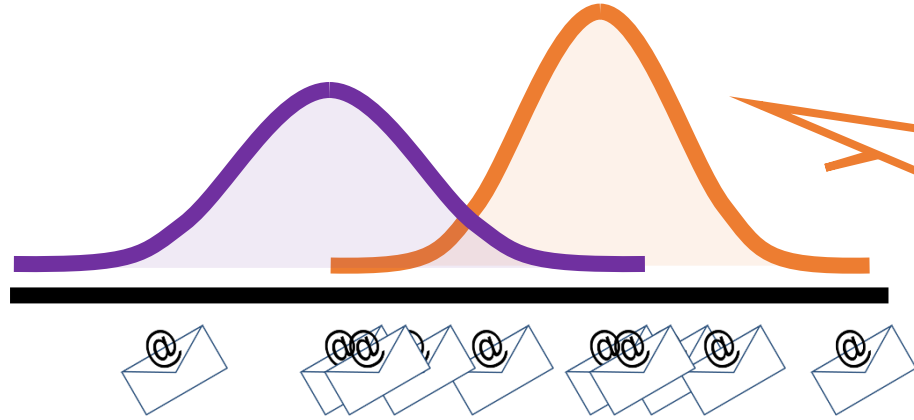
Can use these to label emails!

Hmm ... what if someone magically labelled the emails with the color

Hmm ... what if someone gave me a purple and an orange Gaussian, maybe slightly wrong ones

How to get these magical labels??

# The generative story for unlabelled data??



No labels!!!

Give the Gaussians colors to help identifying them

Can we still recover the two Gaussian components in the mixture??

We do not know which email is purple and which is orange – no labels!!

$$\mathbb{P}[\text{purple} \mid \text{envelope}] \geq \mathbb{P}[\text{orange} \mid \text{envelope}] \Rightarrow \text{purple}$$

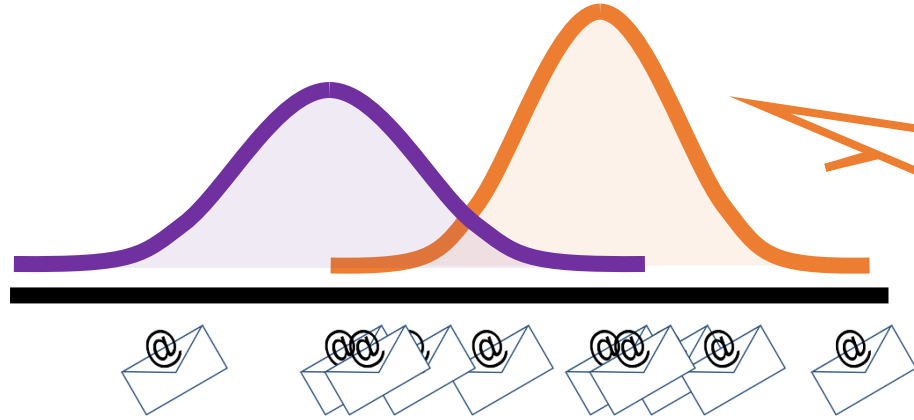
Can use these to label emails!

Hmm ... what if someone magically labelled the emails with the color

Hmm ... what if someone gave me a purple and an orange Gaussian, maybe slightly wrong ones

How to get these magical labels??

# The generative story for unlabelled data??



No labels!!!

Give the Gaussians colors to help identifying them

Can we still recover the two Gaussian components in the mixture??

We do not know which email is purple and which is orange – no labels!!

$$\begin{aligned} \mathbb{P}[\text{purple} \mid \text{envelope}] &\geq \mathbb{P}[\text{orange} \mid \text{envelope}] \Rightarrow \text{purple} \\ \mathbb{P}[\text{orange} \mid \text{envelope}] &\geq \mathbb{P}[\text{purple} \mid \text{envelope}] \Rightarrow \text{orange} \end{aligned}$$

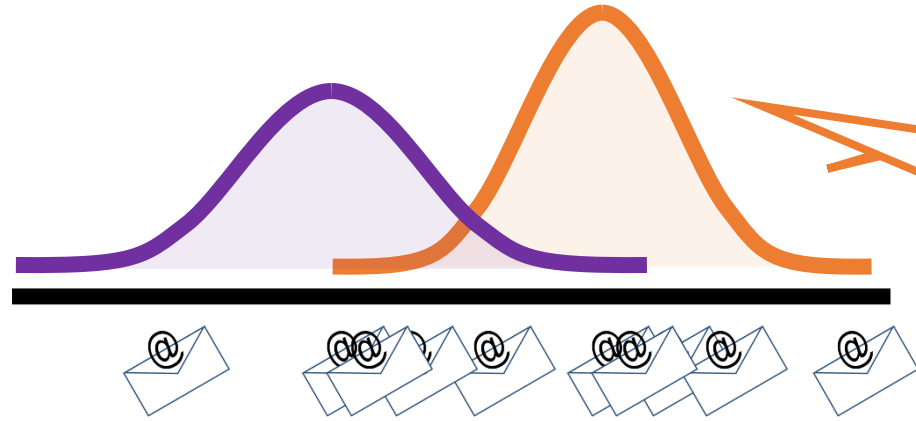
Can use these to label emails!

Hmm ... what if someone magically labelled the emails with the color

Hmm ... what if someone gave me a purple and an orange Gaussian, maybe slightly wrong ones

How to get these magical labels??

# The generative story for unlabelled data??



No labels!!!

Give the Gaussians colors to help identifying them

Can we still recover the two Gaussian components in the mixture??

The dist. are wrong but trust them for now

We do not know which email is purple and which is orange – no labels!!

$$\begin{aligned} \mathbb{P}[\text{purple} \mid \text{envelope}] &\geq \mathbb{P}[\text{orange} \mid \text{envelope}] \Rightarrow \text{purple} \\ \mathbb{P}[\text{orange} \mid \text{envelope}] &\geq \mathbb{P}[\text{purple} \mid \text{envelope}] \Rightarrow \text{orange} \end{aligned}$$

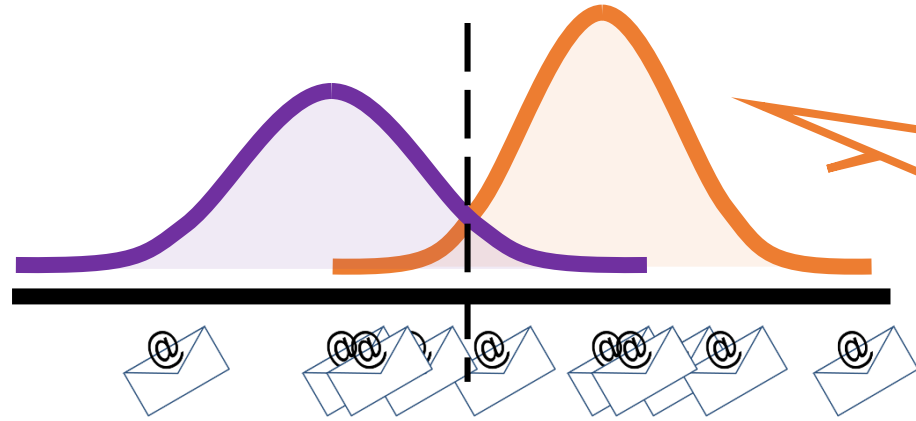
Can use these to label emails!

Hmm ... what if someone magically labelled the emails with the color

Hmm ... what if someone gave me a purple and an orange Gaussian, maybe slightly wrong ones

How to get these magical labels??

# The generative story for unlabelled data??



No labels!!!

Give the Gaussians colors to help identifying them

Can we still recover the two Gaussian components in the mixture??

The dist. are wrong but trust them for now

We do not know which email is purple and which is orange – no labels!!

$$\begin{aligned} \mathbb{P}[\text{purple} \mid \text{envelope}] &\geq \mathbb{P}[\text{orange} \mid \text{envelope}] \Rightarrow \text{purple} \\ \mathbb{P}[\text{orange} \mid \text{envelope}] &\geq \mathbb{P}[\text{purple} \mid \text{envelope}] \Rightarrow \text{orange} \end{aligned}$$

Can use these to label emails!

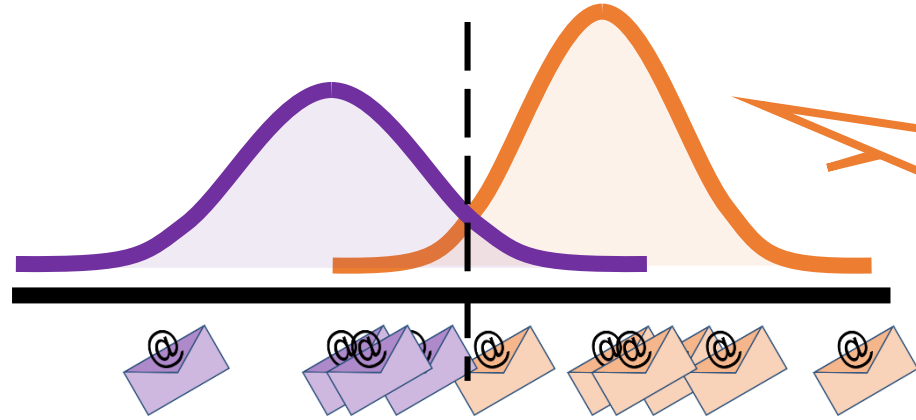
Hmm ... what if someone magically labelled the emails with the color

Hmm ... what if someone gave me a purple and an orange Gaussian, maybe slightly wrong ones

How to get these magical labels??



# The generative story for unlabelled data??



No labels!!!

Give the Gaussians colors to help identifying them

Can we still recover the two Gaussian components in the mixture??

The dist. are wrong but trust them for now

We do not know which email is purple and which is orange – no labels!!

$$\begin{aligned} \mathbb{P}[\text{purple} \mid \text{envelope}] &\geq \mathbb{P}[\text{orange} \mid \text{envelope}] \Rightarrow \text{purple} \\ \mathbb{P}[\text{orange} \mid \text{envelope}] &\geq \mathbb{P}[\text{purple} \mid \text{envelope}] \Rightarrow \text{orange} \end{aligned}$$

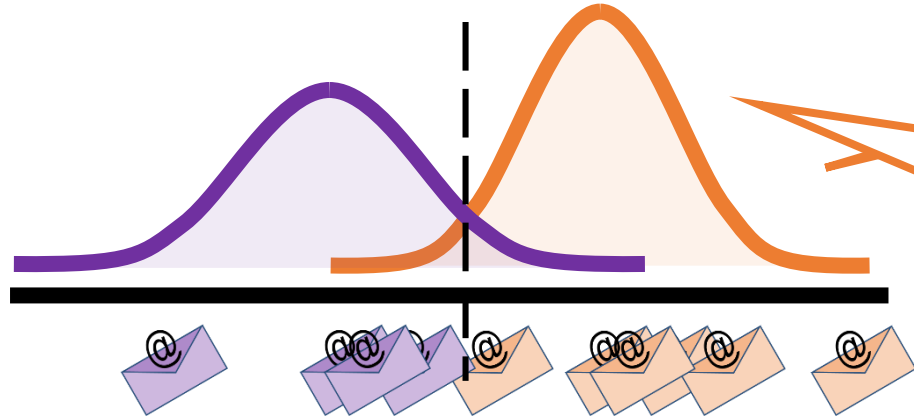
Can use these to label emails!

Hmm ... what if someone magically labelled the emails with the color

Hmm ... what if someone gave me a purple and an orange Gaussian, maybe slightly wrong ones

How to get these magical labels??

# The generative story for unlabelled data??



No labels!!!

Give the Gaussians colors to help identifying them

Can we still recover the two Gaussian components in the mixture??

We do not know which email is purple and which is orange – no labels!!

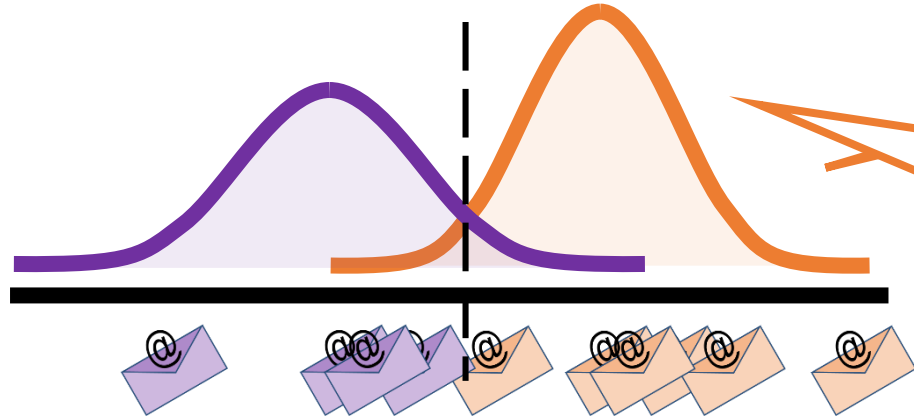
Can use these to label emails!

Hmm ... what if someone magically labelled the emails with the color

Hmm ... what if someone gave me a purple and an orange Gaussian, maybe slightly wrong ones

How to get these magical labels??

# The generative story for unlabelled data??



No labels!!!

Give the Gaussians colors to help identifying them

Can we still recover the two Gaussian components in the mixture??

We do not know which email is purple and which is orange – no labels!!

What if I use these “magical” labels to learn two Gaussians?

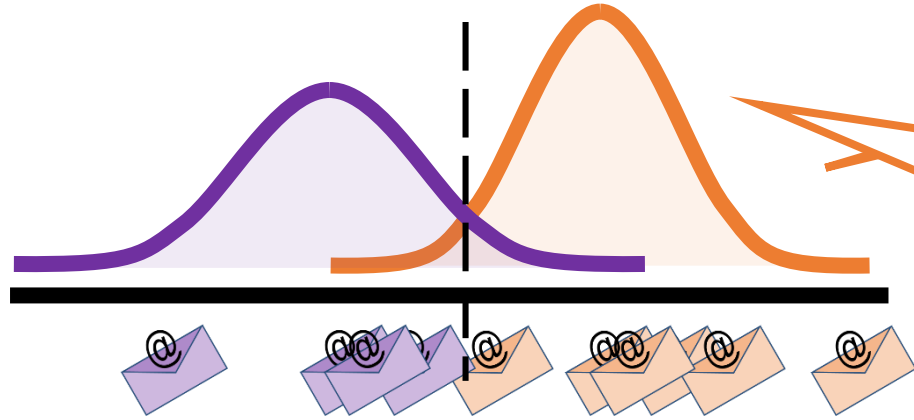
Can use these to label emails!

Hmm ... what if someone magically labelled the emails with the color

Hmm ... what if someone gave me a purple and an orange Gaussian, maybe slightly wrong ones

How to get these magical labels??

# The generative story for unlabelled data??



No labels!!!

Give the Gaussians colors to help identifying them

Can we still recover the two Gaussian components in the mixture??

We do not know which email is purple and which is orange – no labels!!

Won't I just learn the wrong ones that generated the labels?

What if I use these "magical" labels to learn two Gaussians?

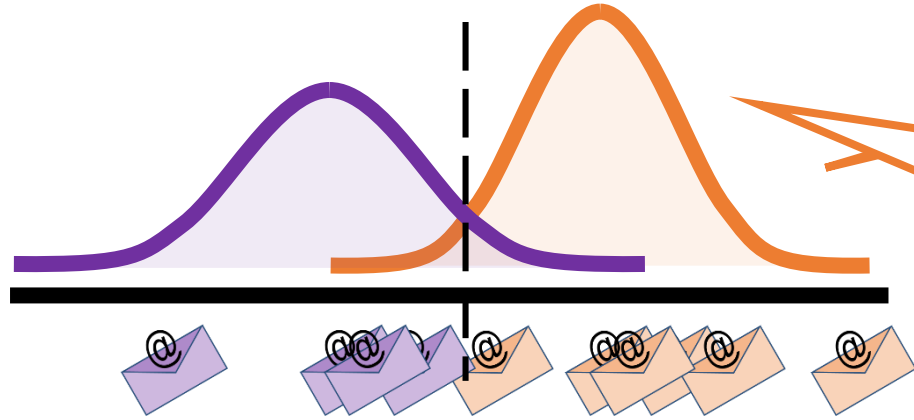
Can use these to label emails!

Hmm ... what if someone magically labelled the emails with the color

Hmm ... what if someone gave me a purple and an orange Gaussian, maybe slightly wrong ones

How to get these magical labels??

# The generative story for unlabelled data??



No labels!!!

Give the Gaussians colors to help identifying them

Can we still recover the two Gaussian components in the mixture??

In practice: you will learn slightly "less wrong" ones ☺

Won't I just learn the wrong ones that generated the labels?

We do not know which email is purple and which is orange – no labels!!

What if I use these "magical" labels to learn two Gaussians?

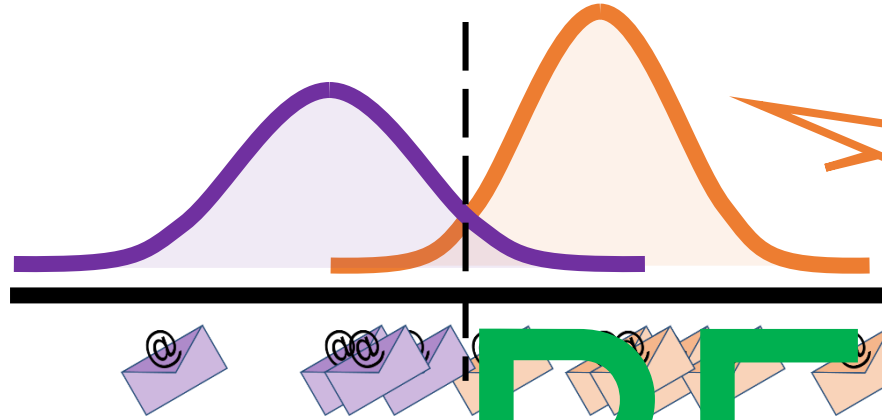
Can use these to label emails!

Hmm ... what if someone magically labelled the emails with the color

Hmm ... what if someone gave me a purple and an orange Gaussian, maybe slightly wrong ones

How to get these magical labels??

# The generative story for unlabelled data??



No labels!!!

Give the Gaussians colors to help identifying them

Can we still recover the two Gaussian components in the mixture??

In practice: you will learn slightly "less wrong" ones ☺

Would it just learn the wrong ones that generated the labels?

We do not know which email is purple and which is orange – no labels!!

What if I use these "magical" labels to learn two Gaussians?

Can use these to label emails!

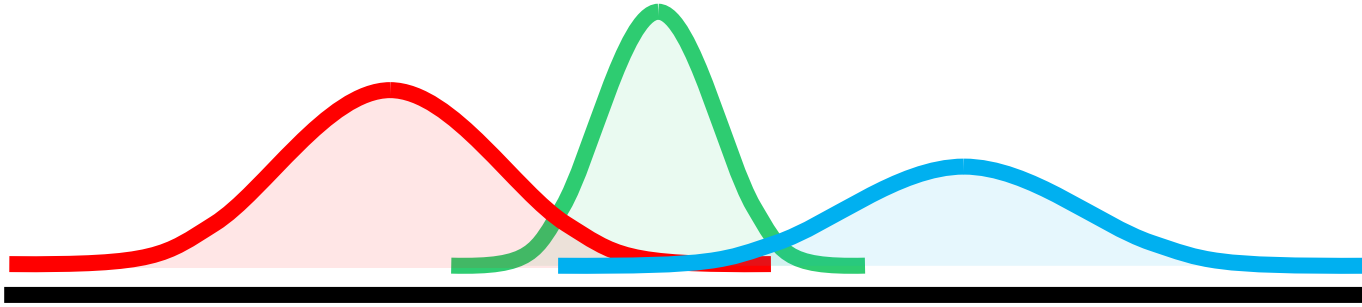
Hmm ... what if someone magically labelled the emails with the color

Hmm ... what if someone gave me a purple and an orange Gaussian, maybe slightly wrong ones

How to get these magical labels??

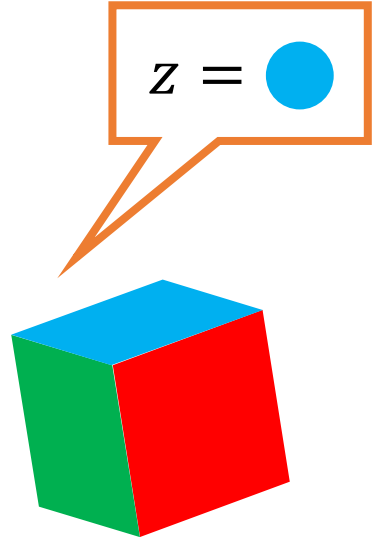
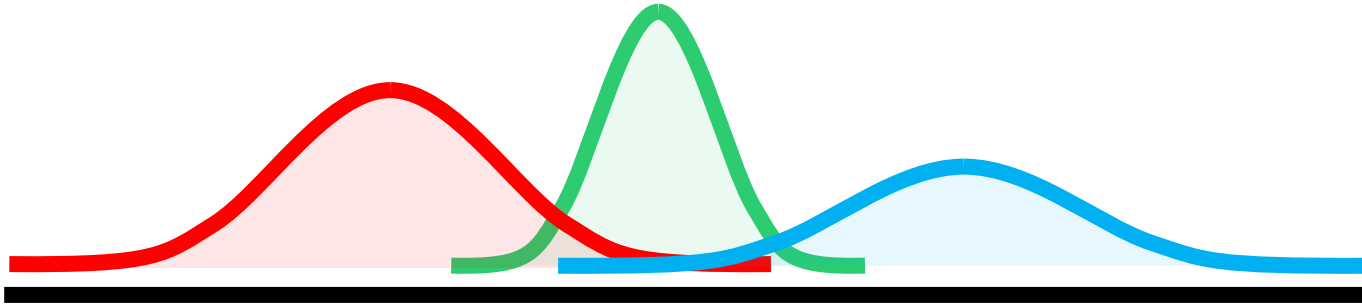
# The generative story for unlabelled data

# The generative story for unlabelled data

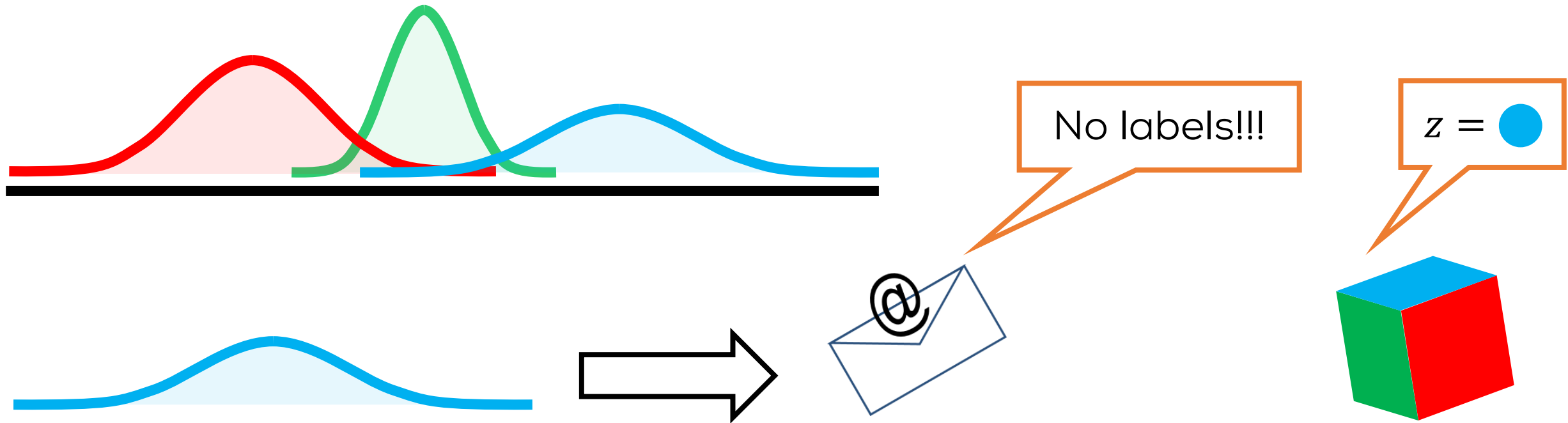




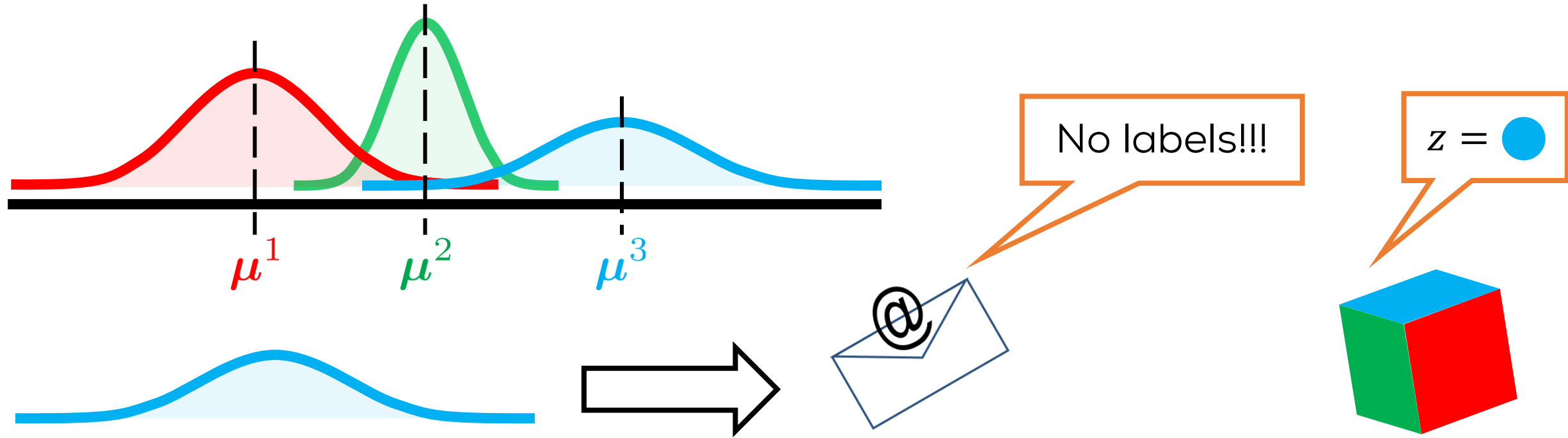
# The generative story for unlabelled data



# The generative story for unlabelled data

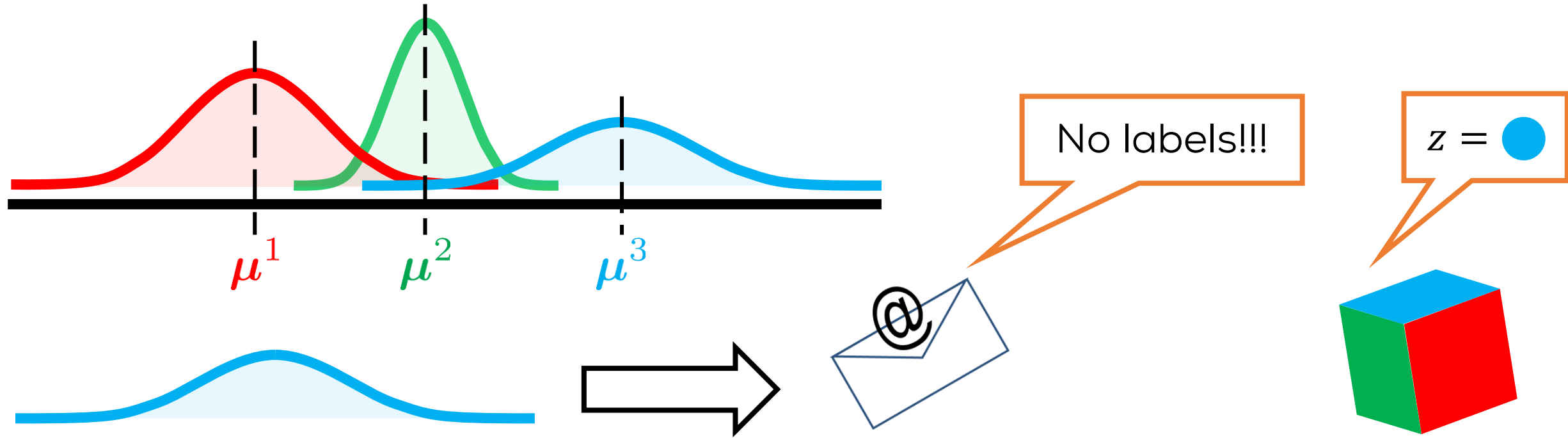


# The generative story for unlabelled data



$$\Theta = \left\{ \pi, \left\{ \mu^1, \mu^2, \mu^3 \right\}, \left\{ \Sigma^1, \Sigma^2, \Sigma^3 \right\} \right\}$$

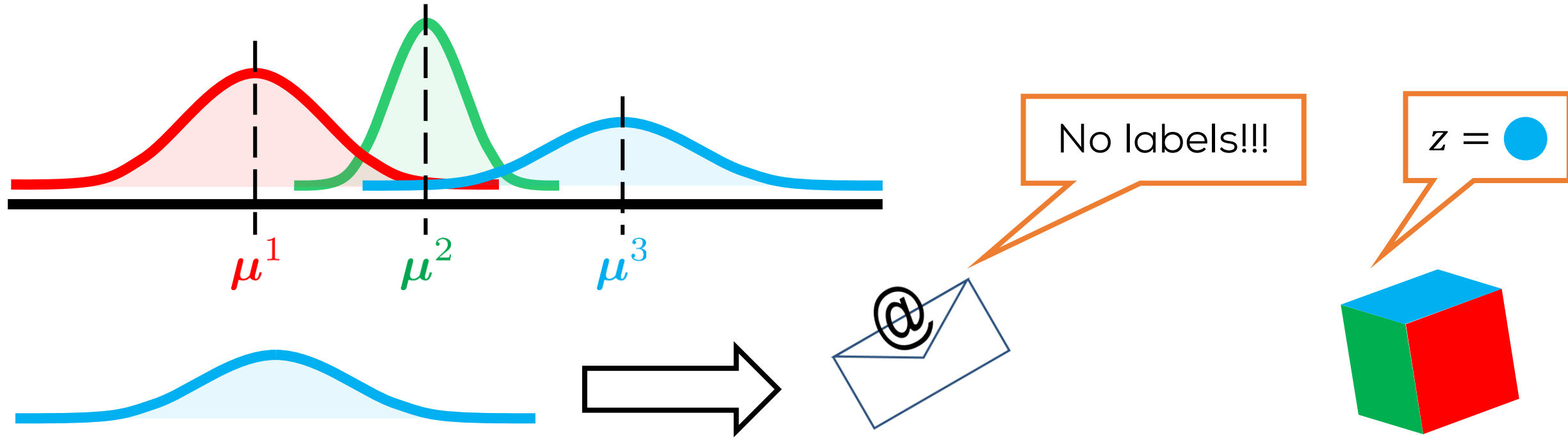
# The generative story for unlabelled data



$$\Theta = \left\{ \pi, \left\{ \mu^1, \mu^2, \mu^3 \right\}, \left\{ \Sigma^1, \Sigma^2, \Sigma^3 \right\} \right\}$$

$$\pi = \left\{ \pi_1, \pi_2, \pi_3 \right\}$$

# The generative story for unlabelled data

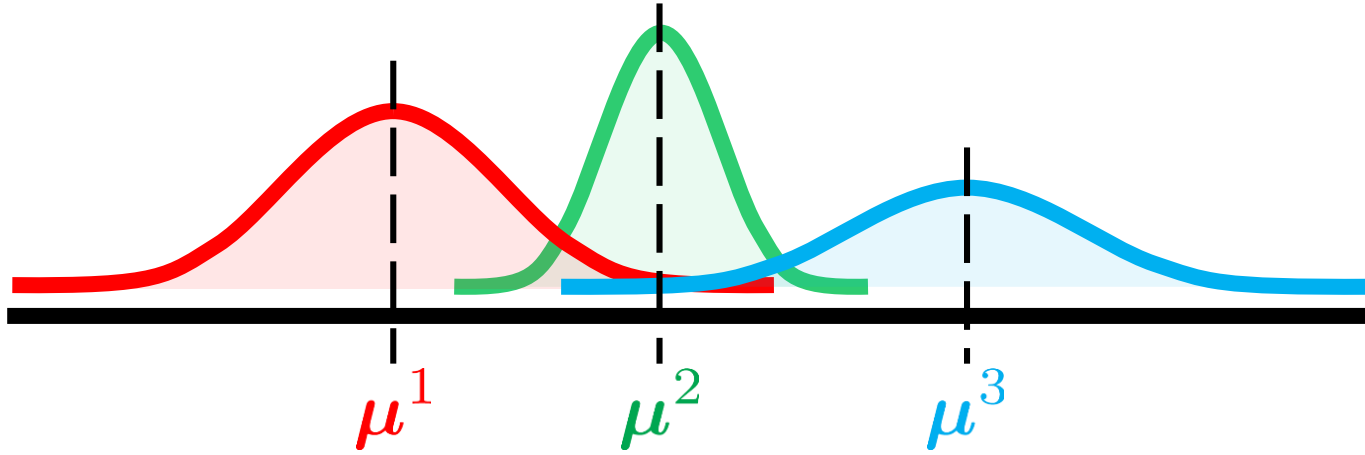


$$\Theta = \left\{ \pi, \left\{ \mu^1, \mu^2, \mu^3 \right\}, \left\{ \Sigma^1, \Sigma^2, \Sigma^3 \right\} \right\}$$

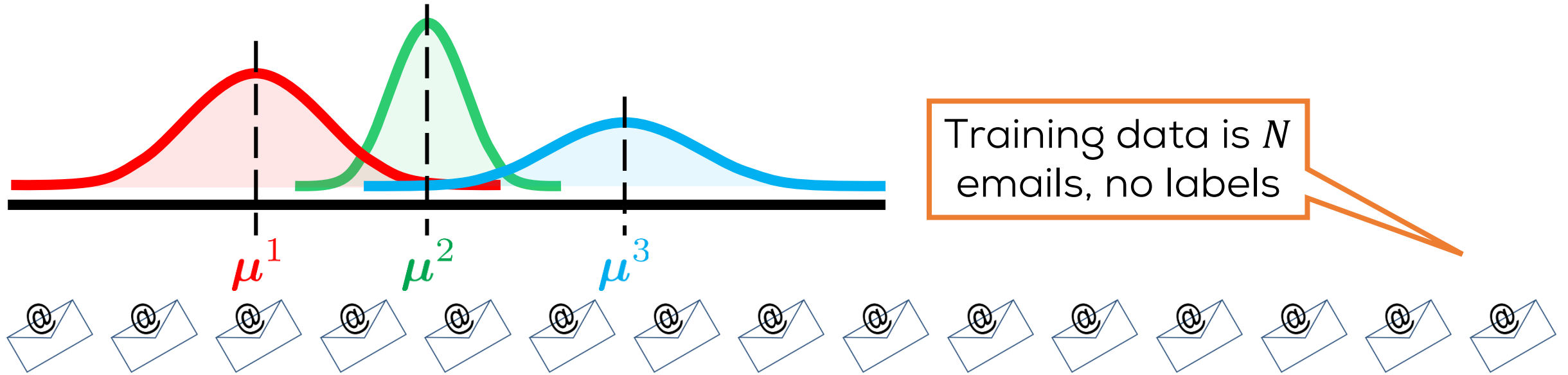
$$\pi = \left\{ \pi_1, \pi_2, \pi_3 \right\} \quad \pi_1 = \mathbb{P} \left[ z = \bullet \right]$$

Just as before!

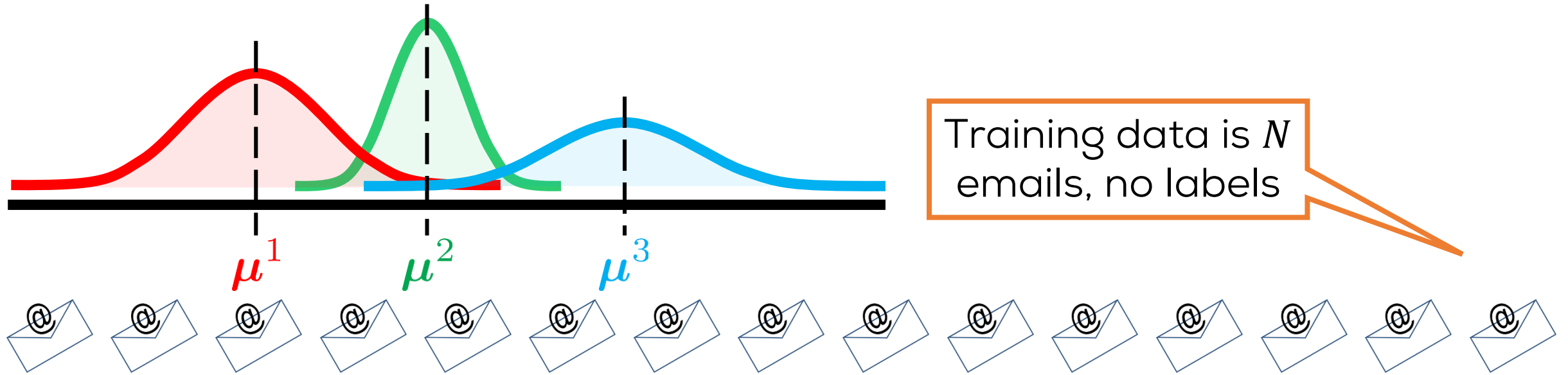
# The generative story for unlabelled data



# The generative story for unlabelled data



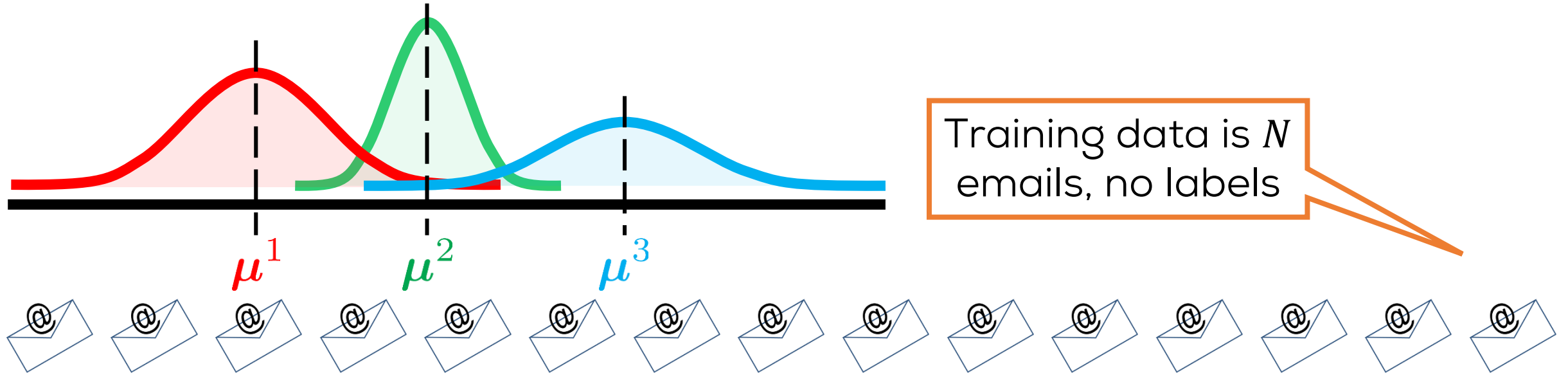
# The generative story for unlabelled data



$$\mathbb{P} [\mathbf{x}^i \mid \Theta]$$

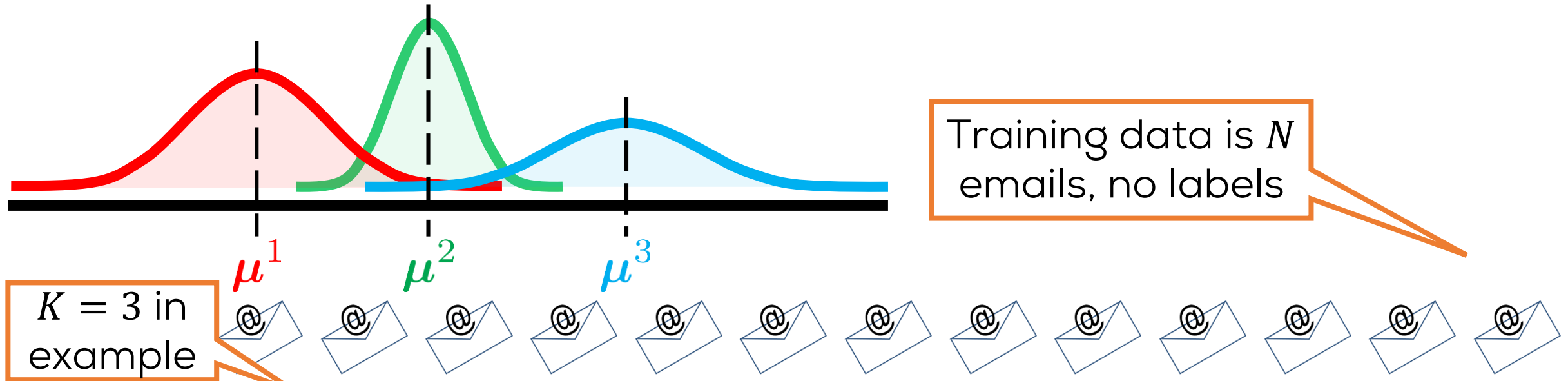


# The generative story for unlabelled data



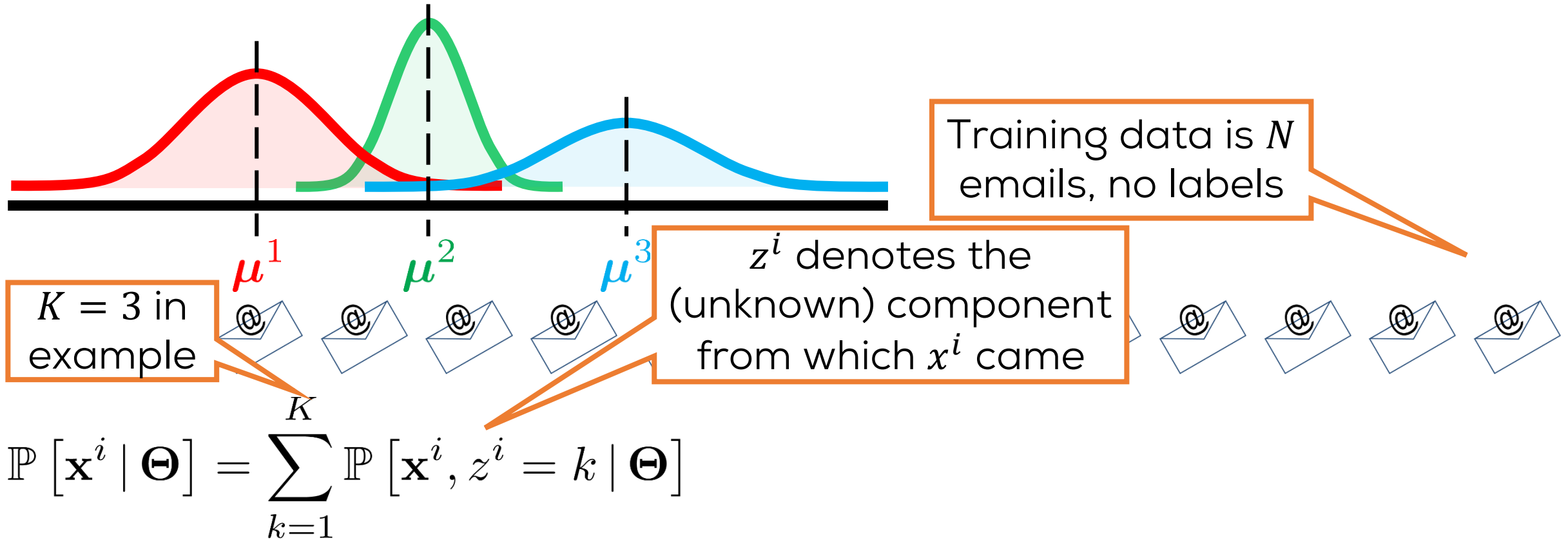
$$\mathbb{P} [\mathbf{x}^i \mid \Theta] = \sum_{k=1}^K \mathbb{P} [\mathbf{x}^i, z^i = k \mid \Theta]$$

# The generative story for unlabelled data

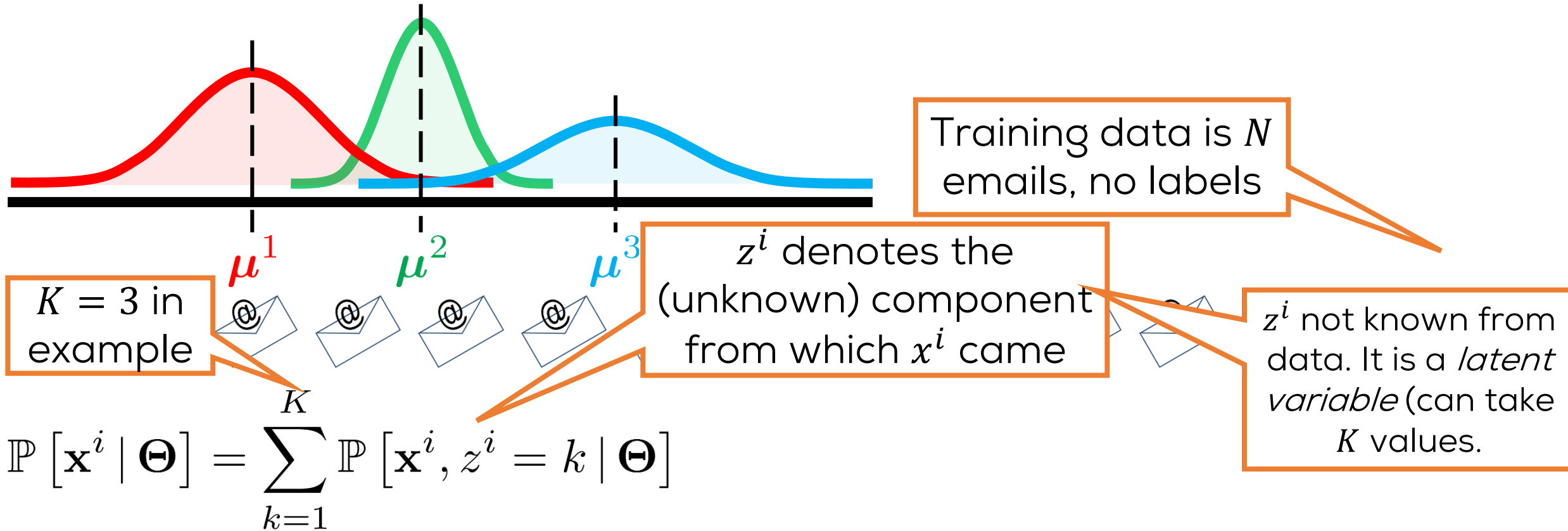


$$\mathbb{P} [\mathbf{x}^i \mid \Theta] = \sum_{k=1}^K \mathbb{P} [\mathbf{x}^i, z^i = k \mid \Theta]$$

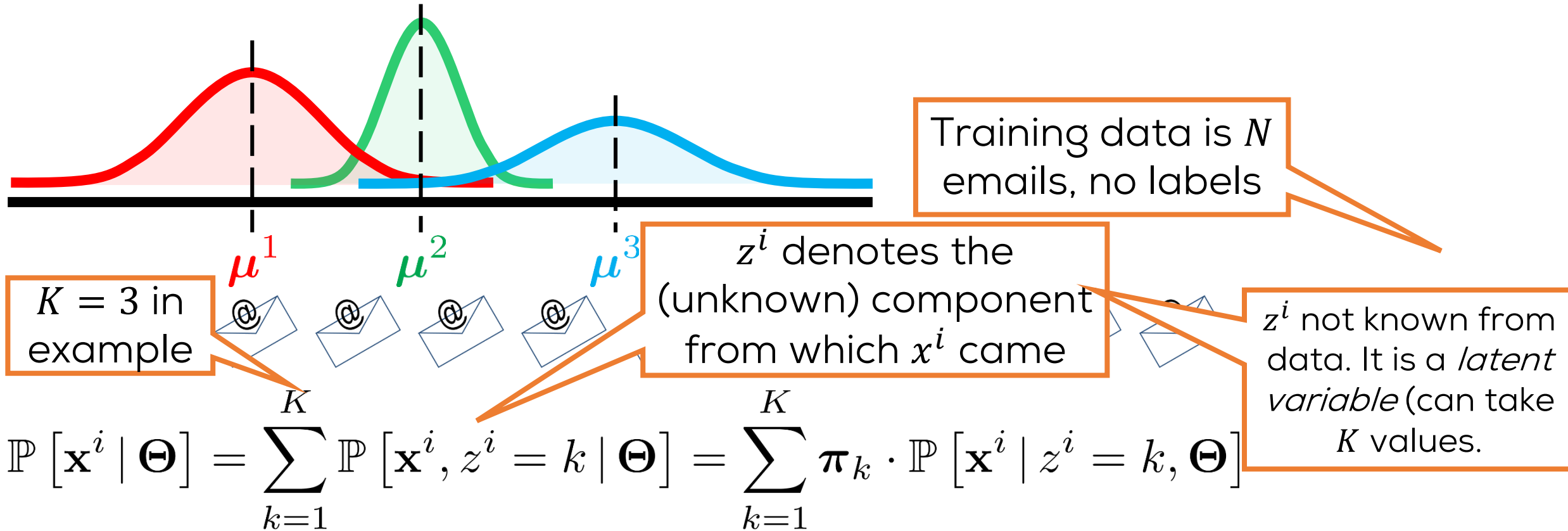
# The generative story for unlabelled data



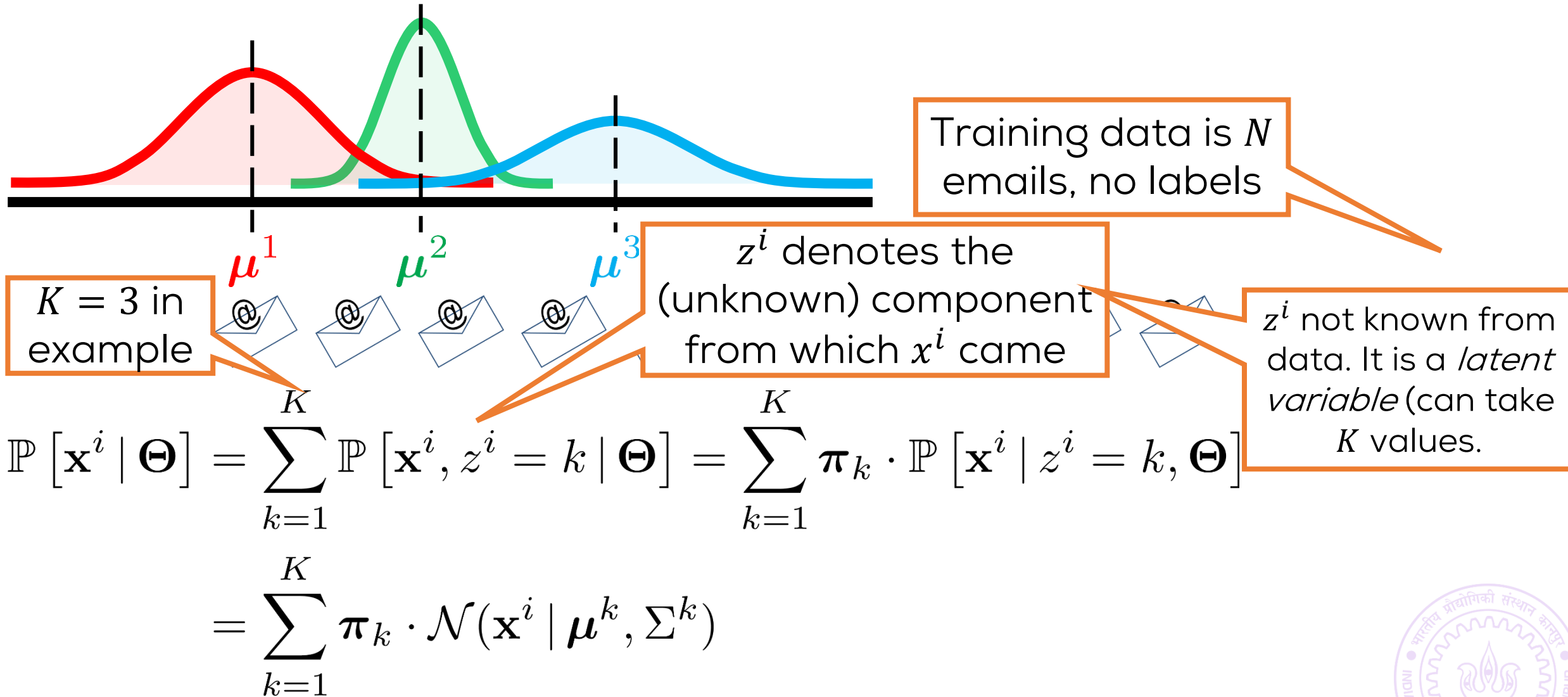
# The generative story for unlabelled data



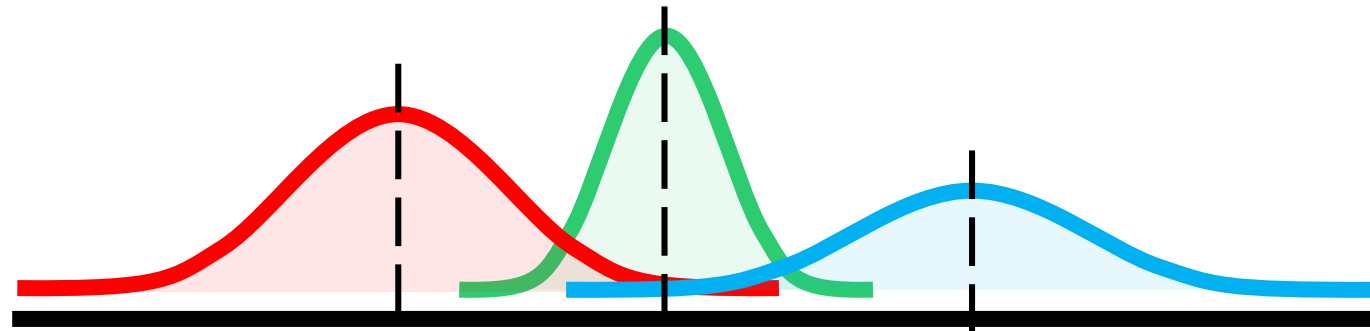
# The generative story for unlabelled data



# The generative story for unlabelled data

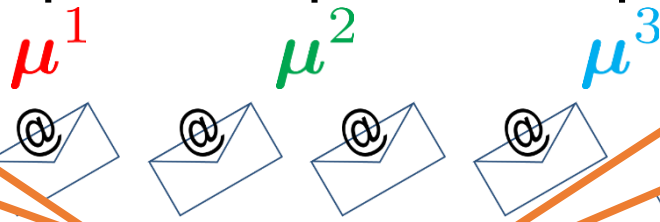


# The generative story for unlabelled data



Training data is  $N$  emails, no labels

$K = 3$  in example



$z^i$  denotes the (unknown) component from which  $x^i$  came

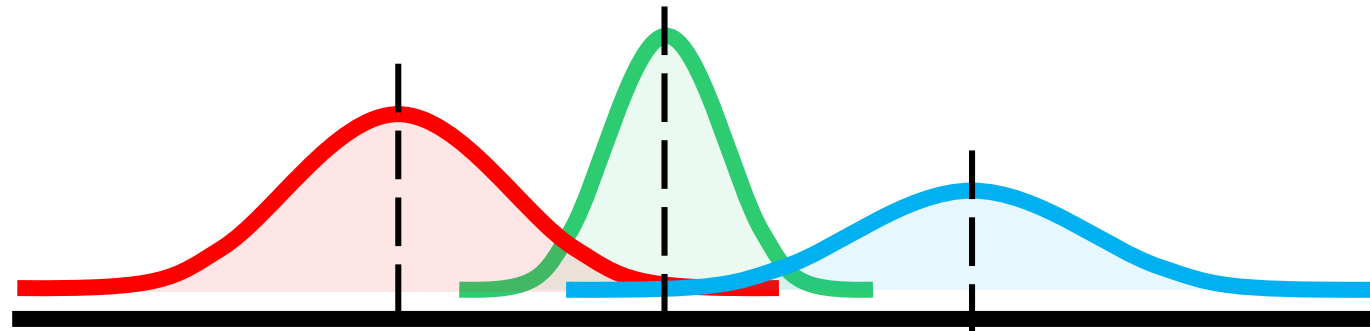
$z^i$  not known from data. It is a *latent variable* (can take  $K$  values).

$$\mathbb{P}[\mathbf{x}^i | \Theta] = \sum_{k=1}^K \mathbb{P}[\mathbf{x}^i, z^i = k | \Theta] = \sum_{k=1}^K \pi_k \cdot \mathbb{P}[\mathbf{x}^i | z^i = k, \Theta]$$

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

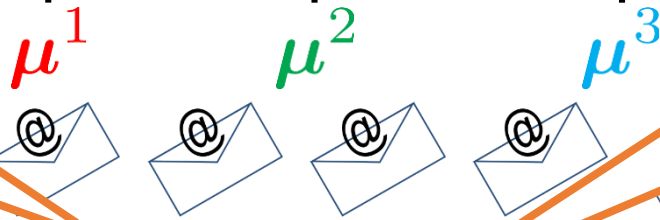
We assumed each component is a multidim. Gaussian

# The generative story for unlabelled data



Training data is  $N$  emails, no labels

$K = 3$  in example



$z^i$  denotes the (unknown) component from which  $x^i$  came

$z^i$  not known from data. It is a *latent variable* (can take  $K$  values).

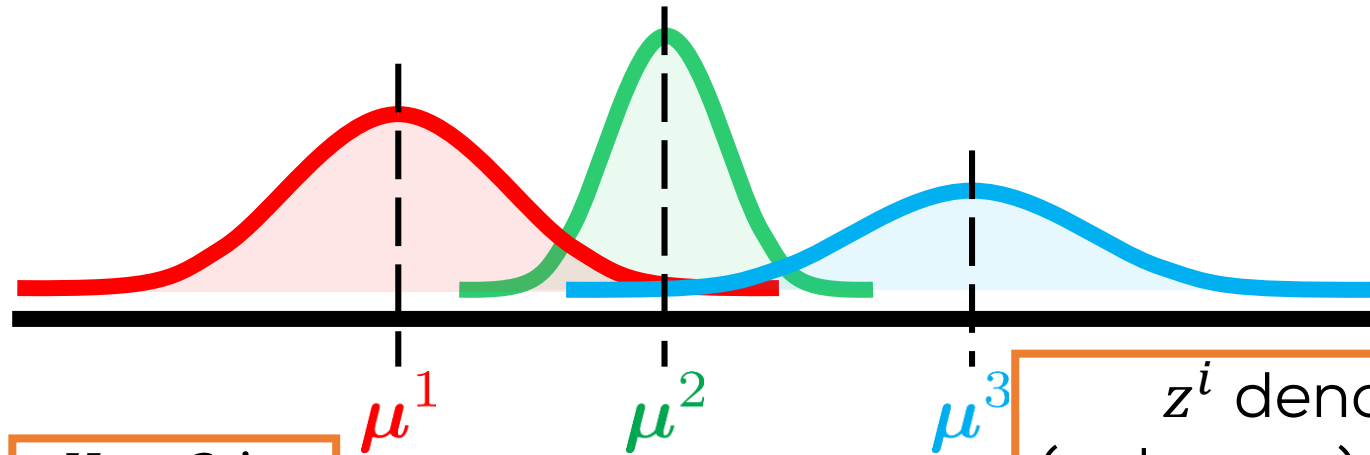
$$\begin{aligned}\mathbb{P}[\mathbf{x}^i | \Theta] &= \sum_{k=1}^K \mathbb{P}[\mathbf{x}^i, z^i = k | \Theta] = \sum_{k=1}^K \pi_k \cdot \mathbb{P}[\mathbf{x}^i | z^i = k, \Theta] \\ &= \sum_{k=1}^K \pi_k \cdot \mathcal{N}(\mathbf{x}^i | \boldsymbol{\mu}^k, \Sigma^k)\end{aligned}$$

Gaussian Mixture Model (GMM) with  $K$  components

We assumed each component is a multidim. Gaussian



# The generative story for unlabelled data



Training data is  $N$  emails, no labels

$K = 3$  in example

$z^i$  denotes the (unknown) component from which  $x^i$  came

$z^i$  not known from data. It is a *latent variable* (can take  $K$  values).

$$\mathbb{P}[\mathbf{x}^i | \Theta] = \sum_{k=1}^K \mathbb{P}[\mathbf{x}^i, z^i = k | \Theta] = \sum_{k=1}^K \pi_k \cdot \mathbb{P}[\mathbf{x}^i | z^i = k, \Theta]$$

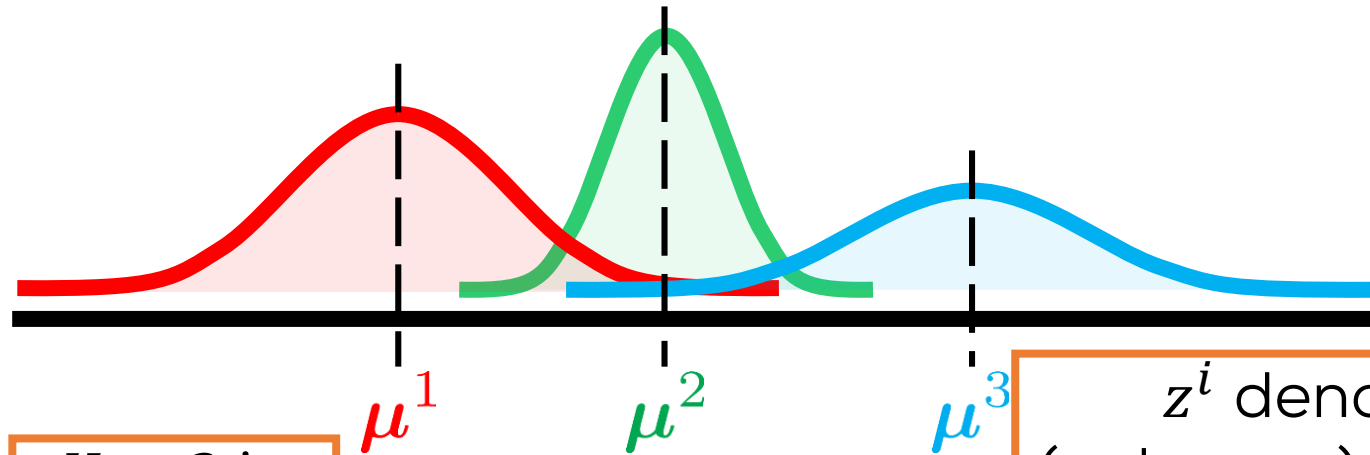
$$= \sum_{k=1}^K \pi_k \cdot \mathcal{N}(\mathbf{x}^i | \boldsymbol{\mu}^k, \Sigma^k)$$

Gaussian Mixture Model (GMM) with  $K$  components

$\mathbb{P}[z^i = k]$  prior prob. of  $\mathbf{x}^i$  coming from  $k$ -th component

We assumed each component is a multidim. Gaussian

# The generative story for unlabelled data



Training data is  $N$  emails, no labels

$K = 3$  in example

$z^i$  denotes the (unknown) component from which  $x^i$  came

$z^i$  not known from data. It is a *latent variable* (can take  $K$  values).

$$\mathbb{P}[\mathbf{x}^i | \Theta] = \sum_{k=1}^K \mathbb{P}[\mathbf{x}^i, z^i = k | \Theta] = \sum_{k=1}^K \pi_k \cdot \mathbb{P}[\mathbf{x}^i | z^i = k, \Theta]$$

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

Gaussian Mixture Model (GMM) with  $K$  components

Goal: incomplete data, learn  $\boldsymbol{\mu}^k, \Sigma^k, \mathbb{P}[z = k]$

$\mathbb{P}[z^i = k]$  prior prob. of  $\mathbf{x}^i$  coming from  $k$ -th component

We assumed each component is a multidim. Gaussian

# The Likelihood Expression

# The Likelihood Expression

$$\mathbb{P} [\mathbf{x}^i \mid \Theta] = \sum_{k=1}^K \pi_k \cdot \mathcal{N}(\mathbf{x}^i \mid \boldsymbol{\mu}^k, \Sigma^k)$$

# The Likelihood Expression

$$\mathbb{P} [\mathbf{x}^i \mid \Theta] = \sum_{k=1}^K \pi_k \cdot \mathcal{N}(\mathbf{x}^i \mid \boldsymbol{\mu}^k, \Sigma^k)$$

$$\mathbb{P} [X \mid \Theta] = \prod_{i=1}^n \mathbb{P} [\mathbf{x}^i \mid \Theta]$$

# The Likelihood Expression

$$\mathbb{P} [\mathbf{x}^i \mid \Theta] = \sum_{k=1}^K \pi_k \cdot \mathcal{N}(\mathbf{x}^i \mid \boldsymbol{\mu}^k, \Sigma^k)$$

$$\mathbb{P} [X \mid \Theta] = \prod_{i=1}^n \mathbb{P} [\mathbf{x}^i \mid \Theta]$$

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P} [X \mid \Theta]$$

$$= \arg \max_{\Theta} \prod_{i=1}^n \sum_{k=1}^K \pi_k \cdot \mathcal{N}(\mathbf{x}^i \mid \boldsymbol{\mu}^k, \Sigma^k)$$

# The Likelihood Expression

$$\mathbb{P}[\mathbf{x}^i | \Theta] = \sum_{k=1}^K \pi_k \cdot \mathcal{N}(\mathbf{x}^i | \boldsymbol{\mu}^k, \Sigma^k)$$

$$\mathbb{P}[X | \Theta] = \prod_{i=1}^n \mathbb{P}[\mathbf{x}^i | \Theta]$$

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

$$= \arg \max_{\Theta} \prod_{i=1}^n \sum_{k=1}^K \pi_k \cdot \mathcal{N}(\mathbf{x}^i | \boldsymbol{\mu}^k, \Sigma^k)$$

Cannot apply first order  
Optimality to get solution

# The Likelihood Expression

$$\mathbb{P}[\mathbf{x}^i | \Theta] = \sum_{k=1}^K \pi_k \cdot \mathcal{N}(\mathbf{x}^i | \boldsymbol{\mu}^k, \Sigma^k)$$

$$\mathbb{P}[X | \Theta] = \prod_{i=1}^n \mathbb{P}[\mathbf{x}^i | \Theta]$$

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

$$= \arg \max_{\Theta} \prod_{i=1}^n \sum_{k=1}^K \pi_k \cdot \mathcal{N}(\mathbf{x}^i | \boldsymbol{\mu}^k, \Sigma^k)$$

Cannot apply first order Optimality to get solution

Horribly non-convex problem. Initialization matters!



# The Likelihood Expression

$$\mathbb{P}[\mathbf{x}^i | \Theta] = \sum_{k=1}^K \pi_k \cdot \mathcal{N}(\mathbf{x}^i | \boldsymbol{\mu}^k, \Sigma^k)$$

$$\mathbb{P}[X | \Theta] = \prod_{i=1}^n \mathbb{P}[\mathbf{x}^i | \Theta]$$

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

$$= \arg \max_{\Theta} \prod_{i=1}^n \sum_{k=1}^K \pi_k \cdot \mathcal{N}(\mathbf{x}^i | \boldsymbol{\mu}^k, \Sigma^k)$$

Cannot apply first order Optimality to get solution

Horribly non-convex problem. Initialization matters!

Getting to global optimum may be NP-hard ☹

# The Likelihood Expression

$$\mathbb{P}[\mathbf{x}^i | \Theta] = \sum_{k=1}^K \pi_k \cdot \mathcal{N}(\mathbf{x}^i | \boldsymbol{\mu}^k, \Sigma^k)$$

$$\mathbb{P}[X | \Theta] = \prod_{i=1}^n \mathbb{P}[\mathbf{x}^i | \Theta]$$

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

$$= \arg \max_{\Theta} \prod_{i=1}^n \sum_{k=1}^K \pi_k \cdot \mathcal{N}(\mathbf{x}^i | \boldsymbol{\mu}^k, \Sigma^k)$$

Things were so nice  
when we had the  
labels ☺

Cannot apply first order  
Optimality to get solution

Horribly non-convex problem.  
Initialization matters!

Getting to global optimum  
may be NP-hard ☹

# The Likelihood Expression

$$\mathbb{P}[\mathbf{x}^i | \Theta] = \sum_{k=1}^K \pi_k \cdot \mathcal{N}(\mathbf{x}^i | \boldsymbol{\mu}^k, \Sigma^k)$$

$$\mathbb{P}[X | \Theta] = \prod_{i=1}^n \mathbb{P}[\mathbf{x}^i | \Theta]$$

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

$$= \arg \max_{\Theta} \prod_{i=1}^n \sum_{k=1}^K \pi_k \cdot \mathcal{N}(\mathbf{x}^i | \boldsymbol{\mu}^k, \Sigma^k)$$

Things were so nice  
when we had the  
labels ☺

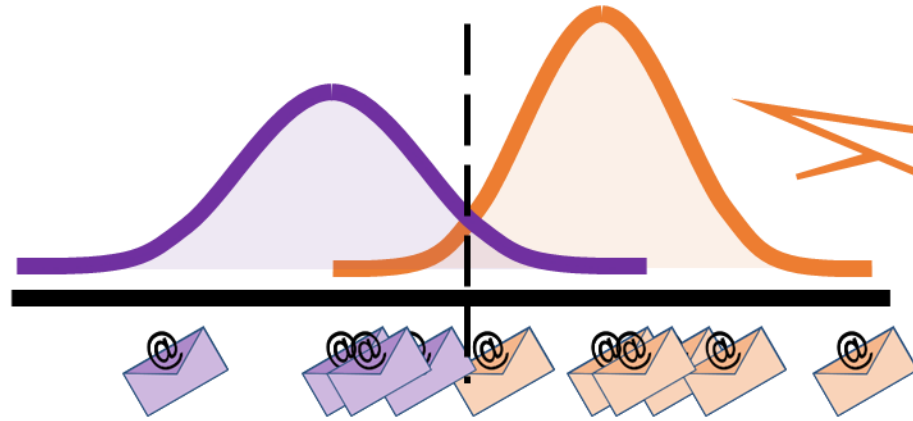
wait ...

Cannot apply first order  
Optimality to get solution

Horribly non-convex problem.  
Initialization matters!

Getting to global optimum  
may be NP-hard ☹

# The generative story for unlabelled data??



No labels!!!

Give the Gaussians colors to help identifying them

Can we still recover the two Gaussian components in the mixture??

In practice: you will learn slightly "less wrong" ones 😊

Won't I just learn the wrong ones that generated the labels?

We do not know which email is purple and which is orange – no labels!!

What if I use these "magical" labels to learn two Gaussians?

Can use these to label emails!

Hmm ... what if someone magically labelled the emails with the color

Hmm ... what if someone gave me a purple and an orange Gaussian, maybe slightly wrong ones

How to get these magical labels??

# The Likelihood Expression

$$\mathbb{P}[\mathbf{x}^i | \Theta] = \sum_{k=1}^K \pi_k \cdot \mathcal{N}(\mathbf{x}^i | \boldsymbol{\mu}^k, \Sigma^k)$$

$$\mathbb{P}[X | \Theta] = \prod_{i=1}^n \mathbb{P}[\mathbf{x}^i | \Theta]$$

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

$$= \arg \max_{\Theta} \prod_{i=1}^n \sum_{k=1}^K \pi_k \cdot \mathcal{N}(\mathbf{x}^i | \boldsymbol{\mu}^k, \Sigma^k)$$

Things were so nice  
when we had the  
labels ☺

wait ...

Cannot apply first order  
Optimality to get solution

Horribly non-convex problem.  
Initialization matters!

Getting to global optimum  
may be NP-hard ☹

# A Ray of Hope

Sept 8, 2017



# A Ray of Hope

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

# A Ray of Hope

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

## ALTERNATING OPTIMIZATION

1. Initialize  $\Theta^0$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\Theta^t$
3. Update  $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} | \Theta]$
4. Repeat until convergence



# A Ray of Hope

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

## ALTERNATING OPTIMIZATION

1. Initialize  $\Theta^0$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\Theta^t$
3. Update  $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} | \Theta]$
4. Repeat until convergence

Can use a method like the one discussed in the “detour”!

# A Ray of Hope

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

## ALTERNATING OPTIMIZATION

1. Initialize  $\Theta^0$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\Theta^t$
3. Update  $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} | \Theta]$
4. Repeat until convergence

Various ways of  
updating  $z^i$

Can use a method like the one  
discussed in the “detour”!

# A Ray of Hope

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

Looks like block coordinate descent with  $\Theta, \{z^i\}$  being two blocks of “coordinates”

## ALTERNATING OPTIMIZATION

1. Initialize  $\Theta^0$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\Theta^t$
3. Update  $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} | \Theta]$
4. Repeat until convergence

Various ways of updating  $z^i$

Can use a method like the one discussed in the “detour”!

# Hard Assignment

The K-means algorithm

# A Ray of Hope

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

Looks like block coordinate descent with  $\Theta, \{z^i\}$  being two blocks of “coordinates”

## ALTERNATING OPTIMIZATION

1. Initialize  $\Theta^0$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\Theta^t$
3. Update  $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} | \Theta]$
4. Repeat until convergence

Various ways of updating  $z^i$

Can use a method like the one discussed in the “detour”!

# A Ray of Hope

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

## ALTERNATING OPTIMIZATION

1. Initialize  $\Theta^0$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\Theta^t$
3. Update  $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} | \Theta]$
4. Repeat until convergence

# A Ray of Hope

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

## ALTERNATING OPTIMIZATION

1. Initialize  $\Theta^0$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\Theta^t$
3. Update  $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} | \Theta]$
4. Repeat until convergence

$$z^{i,t} = \arg \max_{k \in [K]} \mathbb{P}[k | \mathbf{x}^i, \Theta^t]$$

# A Ray of Hope

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

## ALTERNATING OPTIMIZATION

1. Initialize  $\Theta^0$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\Theta^t$
3. Update  $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} | \Theta]$
4. Repeat until convergence

$$z^{i,t} = \arg \max_{k \in [K]} \mathbb{P}[k | \mathbf{x}^i, \Theta^t]$$

$$\Theta^t = \left\{ \pi^t, \left\{ \mu^{1,t}, \mu^{2,t}, \mu^{3,t} \right\}, \left\{ \Sigma^{1,t}, \Sigma^{2,t}, \Sigma^{3,t} \right\} \right\}$$



# A Ray of Hope

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

## ALTERNATING OPTIMIZATION

1. Initialize  $\Theta^0$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\Theta^t$
3. Update  $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} | \Theta]$
4. Repeat until convergence

Bayes  
Rule!

$$z^{i,t} = \arg \max_{k \in [K]} \mathbb{P}[k | \mathbf{x}^i, \Theta^t] = \arg \max_{k \in [K]} \mathbb{P}[k | \Theta^t] \cdot \mathbb{P}[\mathbf{x}^i | k, \Theta^t]$$

$$\Theta^t = \left\{ \pi^t, \left\{ \mu^{1,t}, \mu^{2,t}, \mu^{3,t} \right\}, \left\{ \Sigma^{1,t}, \Sigma^{2,t}, \Sigma^{3,t} \right\} \right\}$$

# A Ray of Hope

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

## ALTERNATING OPTIMIZATION

1. Initialize  $\Theta^0$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\Theta^t$
3. Update  $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} | \Theta]$
4. Repeat until convergence

Bayes Rule!

$$z^{i,t} = \arg \max_{k \in [K]} \mathbb{P}[k | \mathbf{x}^i, \Theta^t] = \arg \max_{k \in [K]} \pi_k^t \cdot \mathbb{P}[\mathbf{x}^i | k, \Theta^t]$$

$$\Theta^t = \left\{ \pi^t, \left\{ \mu^{1,t}, \mu^{2,t}, \mu^{3,t} \right\}, \left\{ \Sigma^{1,t}, \Sigma^{2,t}, \Sigma^{3,t} \right\} \right\}$$

# A Ray of Hope

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

## ALTERNATING OPTIMIZATION

1. Initialize  $\Theta^0$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\Theta^t$
3. Update  $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} | \Theta]$
4. Repeat until convergence

Bayes Rule!

$$z^{i,t} = \arg \max_{k \in [K]} \mathbb{P}[k | \mathbf{x}^i, \Theta^t] = \arg \max_{k \in [K]} \pi_k^t \cdot \mathcal{N}(\mathbf{x}^i | \mu^{k,t}, \Sigma^{k,t})$$

$$\Theta^t = \left\{ \pi^t, \left\{ \mu^{1,t}, \mu^{2,t}, \mu^{3,t} \right\}, \left\{ \Sigma^{1,t}, \Sigma^{2,t}, \Sigma^{3,t} \right\} \right\}$$

# Towards the K-Means Algorithm

## ALTERNATING OPTIMIZATION

1. Initialize  $\Theta^0$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\Theta^t$ 
  1. Let  $z^{i,t} = \arg \max_k \pi_k^t \cdot \mathcal{N}(\mathbf{x}^i \mid \boldsymbol{\mu}^{k,t}, \Sigma^{k,t})$
3. Update  $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} \mid \Theta]$ 
  1. Let  $\pi_k^{t+1} = \frac{n_k^t}{n}$ , where  $n_k^t = |\{i: z^{i,t} = k\}|$
  2. Let  $\boldsymbol{\mu}^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
  3. Let  $\Sigma_k^{t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} (\mathbf{x}^i - \boldsymbol{\mu}^{k,t+1})(\mathbf{x}^i - \boldsymbol{\mu}^{k,t+1})^\top$
4. Repeat until convergence



# A few simplifications

- Fix  $\boldsymbol{\pi}_k^t = \frac{1}{K}$  for all iterations. Don't update it.
- Fix  $\boldsymbol{\Sigma}^{k,t} = I$  for all iterations. Don't update it.

## K-MEANS/LLOYD'S ALGORITHM

1. Initialize means  $\{\boldsymbol{\mu}^{k,0}\}_{k=1,\dots,K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\boldsymbol{\mu}^{k,t}$ 
  1. Let  $z^{i,t} = \arg \max_k \mathcal{N}(\mathbf{x}^i \mid \boldsymbol{\mu}^{k,t}, I)$
3. Update  $\boldsymbol{\mu}^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence

# A few simplifications

- Fix  $\boldsymbol{\pi}_k^t = \frac{1}{K}$  for all iterations. Don't update it.
- Fix  $\boldsymbol{\Sigma}^{k,t} = I$  for all iterations. Don't update it.

## K-MEANS/LLOYD'S ALGORITHM

1. Initialize means  $\{\boldsymbol{\mu}^{k,0}\}_{k=1,\dots,K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\boldsymbol{\mu}^{k,t}$ 
  1. Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \boldsymbol{\mu}^{k,t}\|_2$
3. Update  $\boldsymbol{\mu}^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence

# A few simplifications

- Fix  $\boldsymbol{\pi}_k^t = \frac{1}{K}$  for all iterations. Don't update it.
- Fix  $\boldsymbol{\Sigma}^{k,t} = I$  for all iterations. Don't update it.

## K-MEANS/LLOYD'S ALGORITHM

1. Initialize means  $\{\boldsymbol{\mu}^{k,0}\}_{k=1,\dots,K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\boldsymbol{\mu}^{k,t}$ 
  1. Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \boldsymbol{\mu}^{k,t}\|_2^2$
3. Update  $\boldsymbol{\mu}^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence

# The K-Means Objective



# The K-Means Objective

$$\hat{\Theta}_{\text{km}} = \arg \min_{\substack{\{\boldsymbol{\mu}^k\}_{k=1,\dots,K} \\ \{z^i\}_{i=1,\dots,n}}} \sum_{k=1}^K \sum_{i: z^i=k} \|\mathbf{x}^i - \boldsymbol{\mu}^k\|_2^2$$

# The K-Means Objective

$$\hat{\Theta}_{\text{km}} = \arg \min_{\substack{\{\boldsymbol{\mu}^k\}_{k=1,\dots,K} \\ \{z^i\}_{i=1,\dots,n}}} \sum_{k=1}^K \sum_{i: z^i=k} \|\mathbf{x}^i - \boldsymbol{\mu}^k\|_2^2$$

For cluster  $k$ , we have  
cluster center  $\boldsymbol{\mu}^k$

# The K-Means Objective

$$\hat{\Theta}_{\text{km}} = \arg \min_{\substack{\{\boldsymbol{\mu}^k\}_{k=1,\dots,K} \\ \{z^i\}_{i=1,\dots,n}}} \sum_{k=1}^K \sum_{i:z^i=k} \|\mathbf{x}^i - \boldsymbol{\mu}^k\|_2^2$$

For cluster  $k$ ,

$$\sum_{i:z^i=k} \|\mathbf{x}^i - \boldsymbol{\mu}^k\|_2^2$$

is a measure of how much variance is in the cluster

For cluster  $k$ , we have cluster center  $\boldsymbol{\mu}^k$

# The K-Means Objective

$$\hat{\Theta}_{\text{km}} = \arg \min_{\substack{\{\boldsymbol{\mu}^k\}_{k=1,\dots,K} \\ \{z^i\}_{i=1,\dots,n}}} \sum_{k=1}^K \sum_{i:z^i=k} \|\mathbf{x}^i - \boldsymbol{\mu}^k\|_2^2$$

For cluster  $k$ ,

$$\sum_{i:z^i=k} \|\mathbf{x}^i - \boldsymbol{\mu}^k\|_2^2$$

is a measure of how much variance is in the cluster

For cluster  $k$ , we have cluster center  $\boldsymbol{\mu}^k$

Want to minimize sum of variance across all clusters.  
Want tight clusters

# The K-Means Objective

$$\hat{\Theta}_{\text{km}} = \arg \min_{\substack{\{\boldsymbol{\mu}^k\}_{k=1,\dots,K} \\ \{z^i\}_{i=1,\dots,n}}} \sum_{k=1}^K \sum_{i: z^i=k} \|\mathbf{x}^i - \boldsymbol{\mu}^k\|_2^2$$

# The K-Means Objective

$$\hat{\Theta}_{\text{km}} = \arg \min_{\substack{\{\boldsymbol{\mu}^k\}_{k=1,\dots,K} \\ \{z^i\}_{i=1,\dots,n}}} \sum_{k=1}^K \sum_{i: z^i = k} \|\mathbf{x}^i - \boldsymbol{\mu}^k\|_2^2$$

## K-MEANS/LLOYD'S ALGORITHM

1. Initialize means  $\{\boldsymbol{\mu}^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\boldsymbol{\mu}^{k,t}$ 
  1. Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \boldsymbol{\mu}^{k,t}\|_2^2$
3. Update  $\boldsymbol{\mu}^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence

# The K-Means Objective

$$\hat{\Theta}_{\text{km}} = \arg \min_{\substack{\{\boldsymbol{\mu}^k\}_{k=1,\dots,K} \\ \{z^i\}_{i=1,\dots,n}}} \sum_{k=1}^K \sum_{i: z^i = k} \|\mathbf{x}^i - \boldsymbol{\mu}^k\|_2^2$$

Alternates between updating  $\{z^i\}$  and  $\{\boldsymbol{\mu}^k\}$

## K-MEANS/LLOYD'S ALGORITHM

1. Initialize means  $\{\boldsymbol{\mu}^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\boldsymbol{\mu}^{k,t}$ 
  1. Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \boldsymbol{\mu}^{k,t}\|_2^2$
3. Update  $\boldsymbol{\mu}^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence

# The K-Means Objective

$$\hat{\Theta}_{\text{km}} = \arg \min_{\substack{\{\mu^k\}_{k=1,\dots,K} \\ \{z^i\}_{i=1,\dots,n}}} \sum_{k=1}^K \sum_{i: z^i=k} \|\mathbf{x}^i - \mu^k\|_2^2$$

An FA approach to solving a data modelling task!

Alternates between updating  $\{z^i\}$  and  $\{\mu^k\}$

## K-MEANS/LLOYD'S ALGORITHM

1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$ 
  1. Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



# The K-Means Objective

NP-hard problem!

An FA approach to solving a data modelling task!

$$\hat{\Theta}_{\text{km}} = \arg \min_{\substack{\{\mu^k\}_{k=1,\dots,K} \\ \{z^i\}_{i=1,\dots,n}}} \sum_{k=1}^K \sum_{i: z^i = k} \|\mathbf{x}^i - \mu^k\|_2^2$$

Alternates between updating  $\{z^i\}$  and  $\{\mu^k\}$

## K-MEANS/LLOYD'S ALGORITHM

1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$ 
  1. Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence

# The K-Means Objective

NP-hard problem!

An FA approach to solving a data modelling task!

$$\hat{\Theta}_{\text{km}} = \arg \min_{\substack{\{\mu^k\}_{k=1,\dots,K} \\ \{z^i\}_{i=1,\dots,n}}} \sum_{k=1}^K \sum_{i: z^i=k} \|\mathbf{x}^i - \mu^k\|_2^2$$

Alternates between updating  $\{z^i\}$  and  $\{\mu^k\}$

Very scalable but sensitive to initialization!

## K-MEANS/LLOYD'S ALGORITHM

1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$ 
  1. Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence

# The K-Means Objective

NP-hard problem!

An FA approach to solving a data modelling task!

$$\hat{\Theta}_{\text{km}} = \arg \min_{\substack{\{\mu^k\}_{k=1,\dots,K} \\ \{z^i\}_{i=1,\dots,n}}} \sum_{k=1}^K \sum_{i: z^i=k} \|\mathbf{x}^i - \mu^k\|_2^2$$

Alternates between updating  $\{z^i\}$  and  $\{\mu^k\}$

Very scalable but sensitive to initialization!

k-means++ initialization

1. Sample  $i_1 \sim [n]$ , let  $\mu^{1,0} = \mathbf{x}^{i_1}$
2. For  $k = 2, \dots, K$ 
  - Sample  $i_k \propto \text{min distance from } \{\mu^{1,0}, \dots, \mu^{k-1,0}\}$
  - Let  $\mu^{k,0} = \mathbf{x}^{i_k}$

## K-MEANS/LLOYD'S ALGORITHM

1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$ 
  1. Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence

# K-Means Algorithm in action!

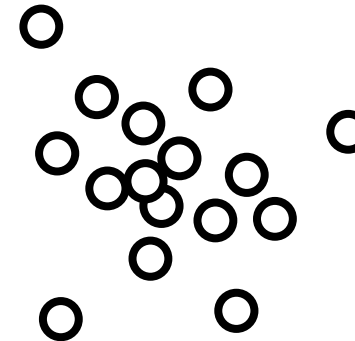
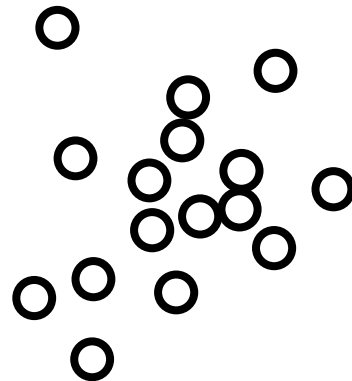
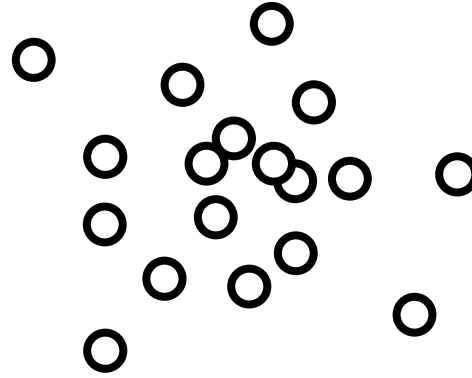
## K-MEANS/LLOYD'S ALGORITHM

1. Initialize means  $\{\boldsymbol{\mu}^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\boldsymbol{\mu}^{k,t}$   
Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \boldsymbol{\mu}^{k,t}\|_2^2$
3. Update  $\boldsymbol{\mu}^{k,t+1} = \frac{1}{n_k^t} \sum_{i:z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence

# K-Means Algorithm in action!

## K-MEANS/LLOYD'S ALGORITHM

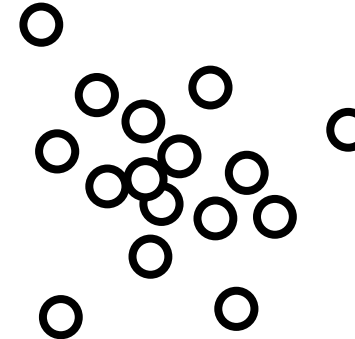
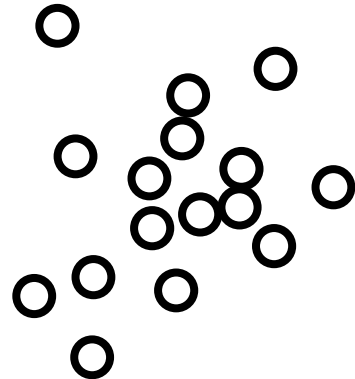
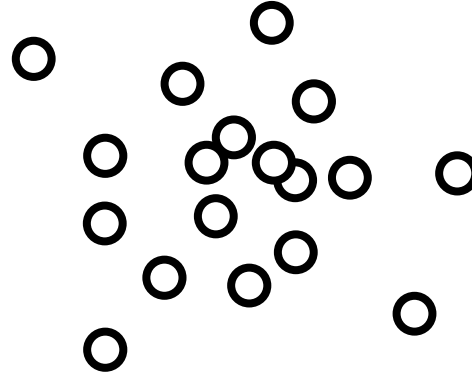
1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$   
Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i:z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



# K-Means Algorithm in action!

## K-MEANS/LLOYD'S ALGORITHM

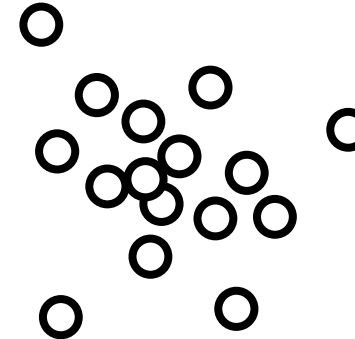
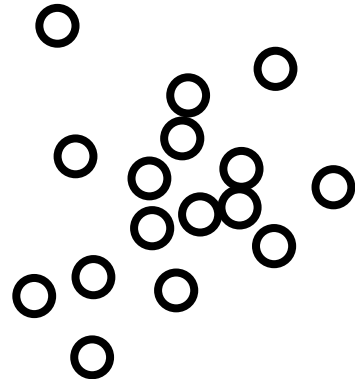
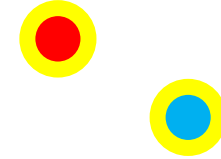
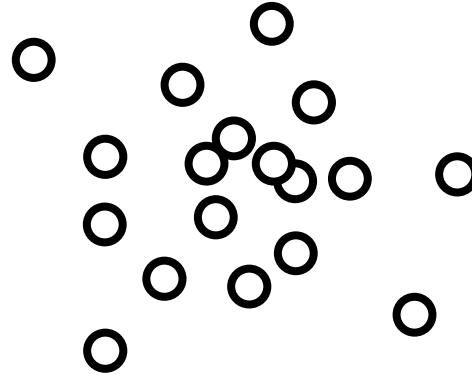
1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$   
Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i:z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



# K-Means Algorithm in action!

## K-MEANS/LLOYD'S ALGORITHM

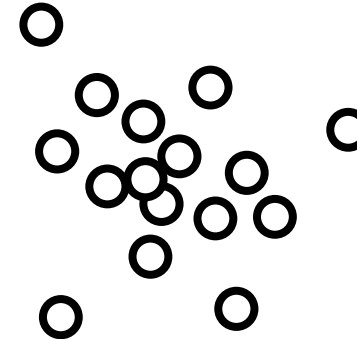
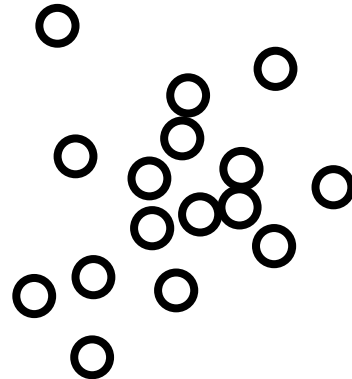
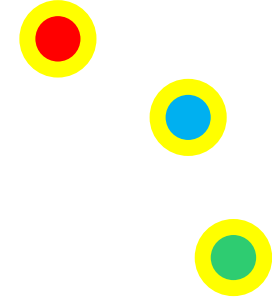
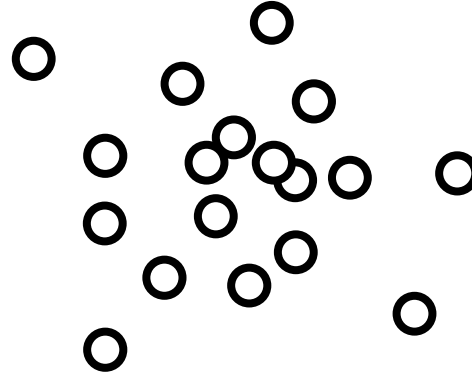
1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$   
Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i:z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



# K-Means Algorithm in action!

## K-MEANS/LLOYD'S ALGORITHM

1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$   
Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i:z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence

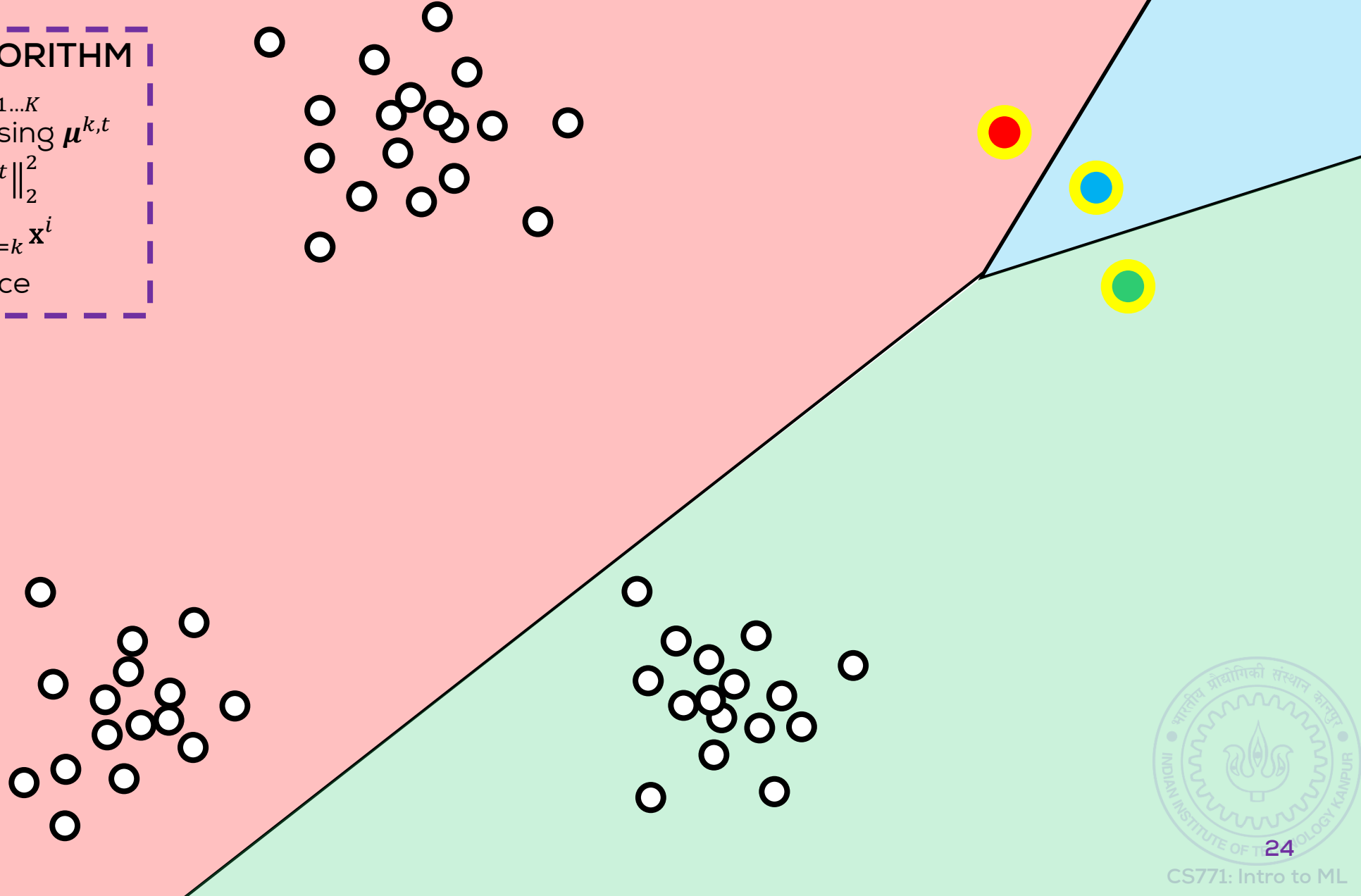




# K-Means Algorithm in action!

## K-MEANS/LLOYD'S ALGORITHM

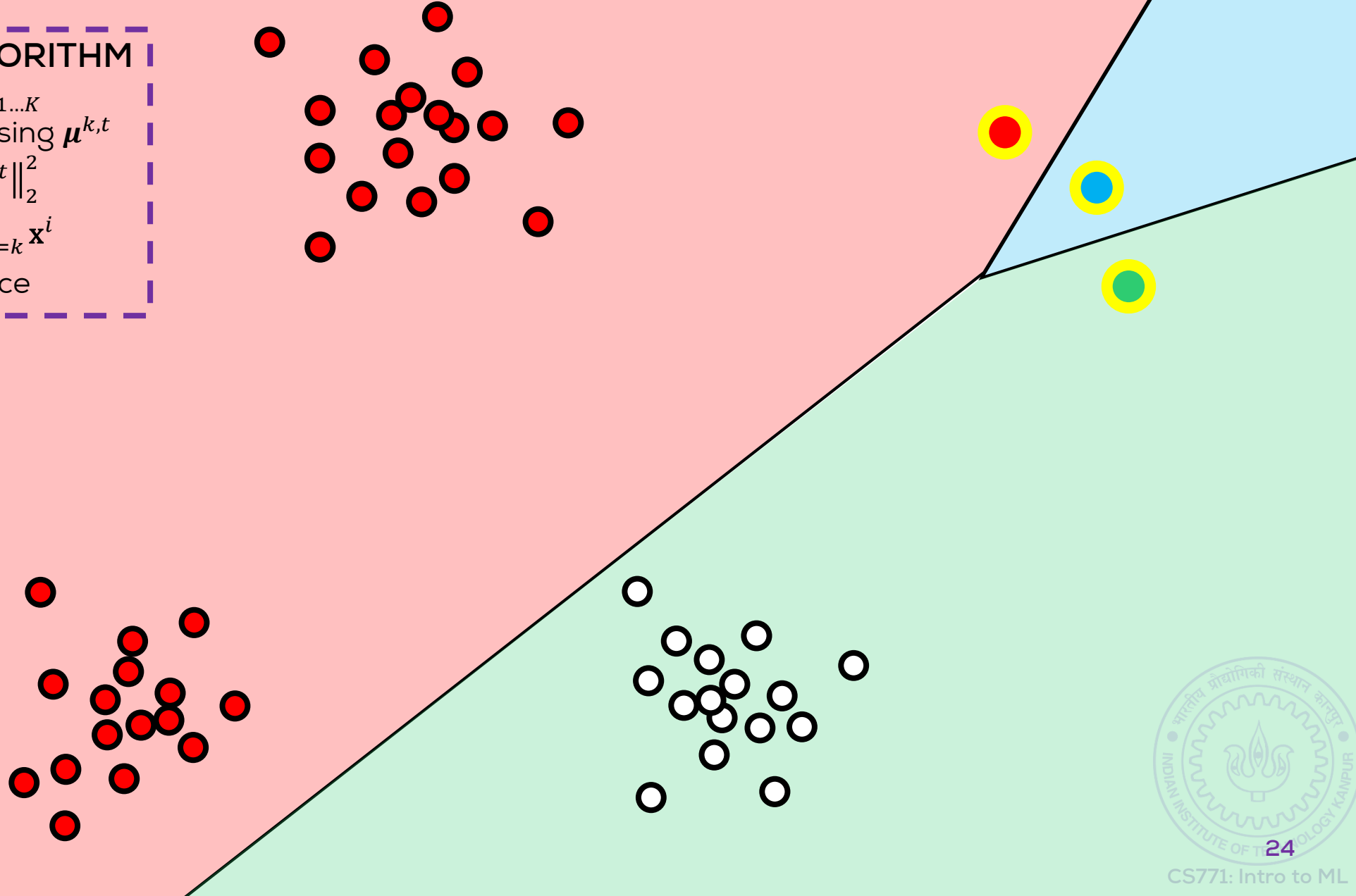
1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$   
Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



# K-Means Algorithm in action!

## K-MEANS/LLOYD'S ALGORITHM

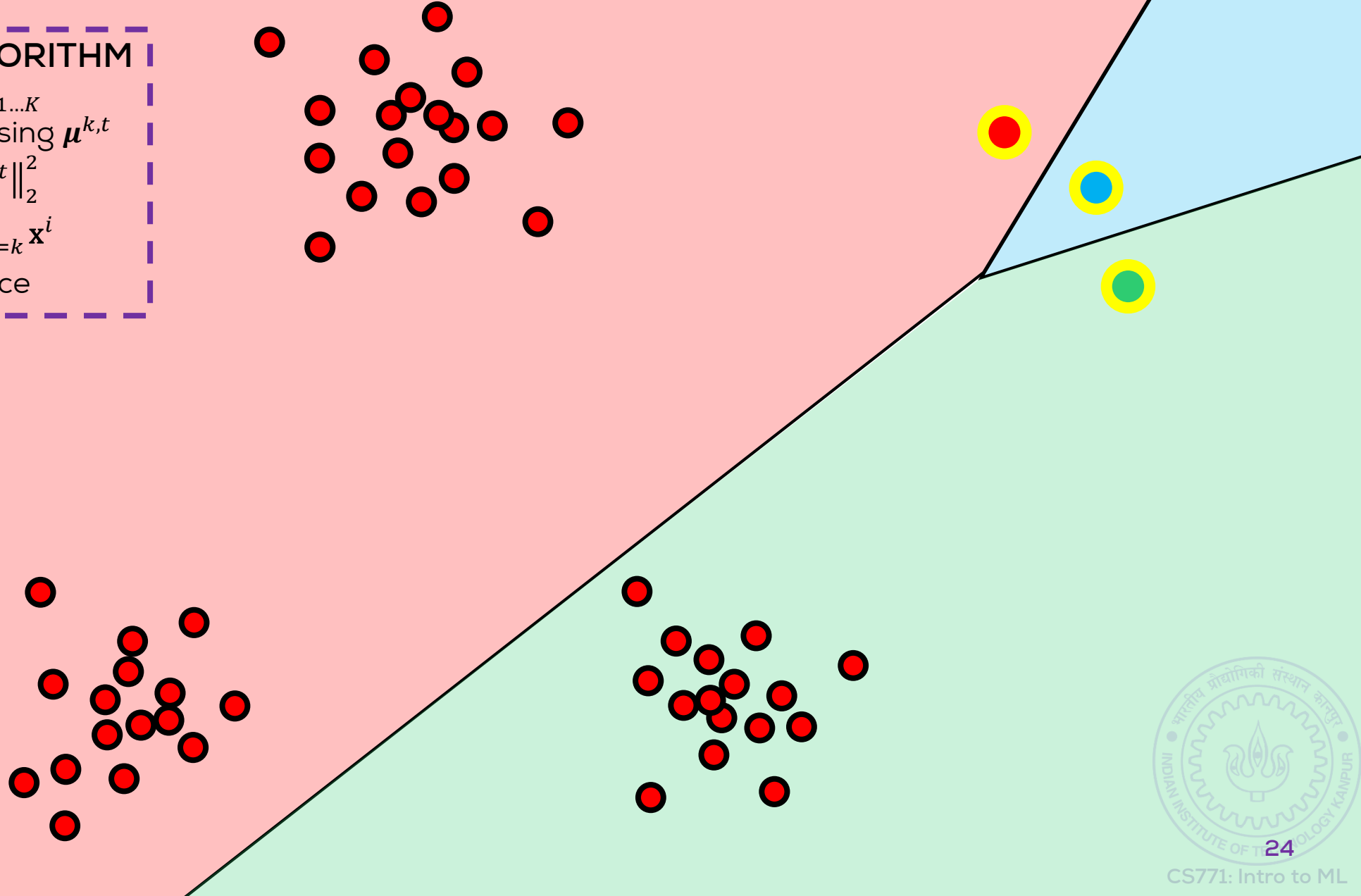
1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$   
Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



# K-Means Algorithm in action!

## K-MEANS/LLOYD'S ALGORITHM

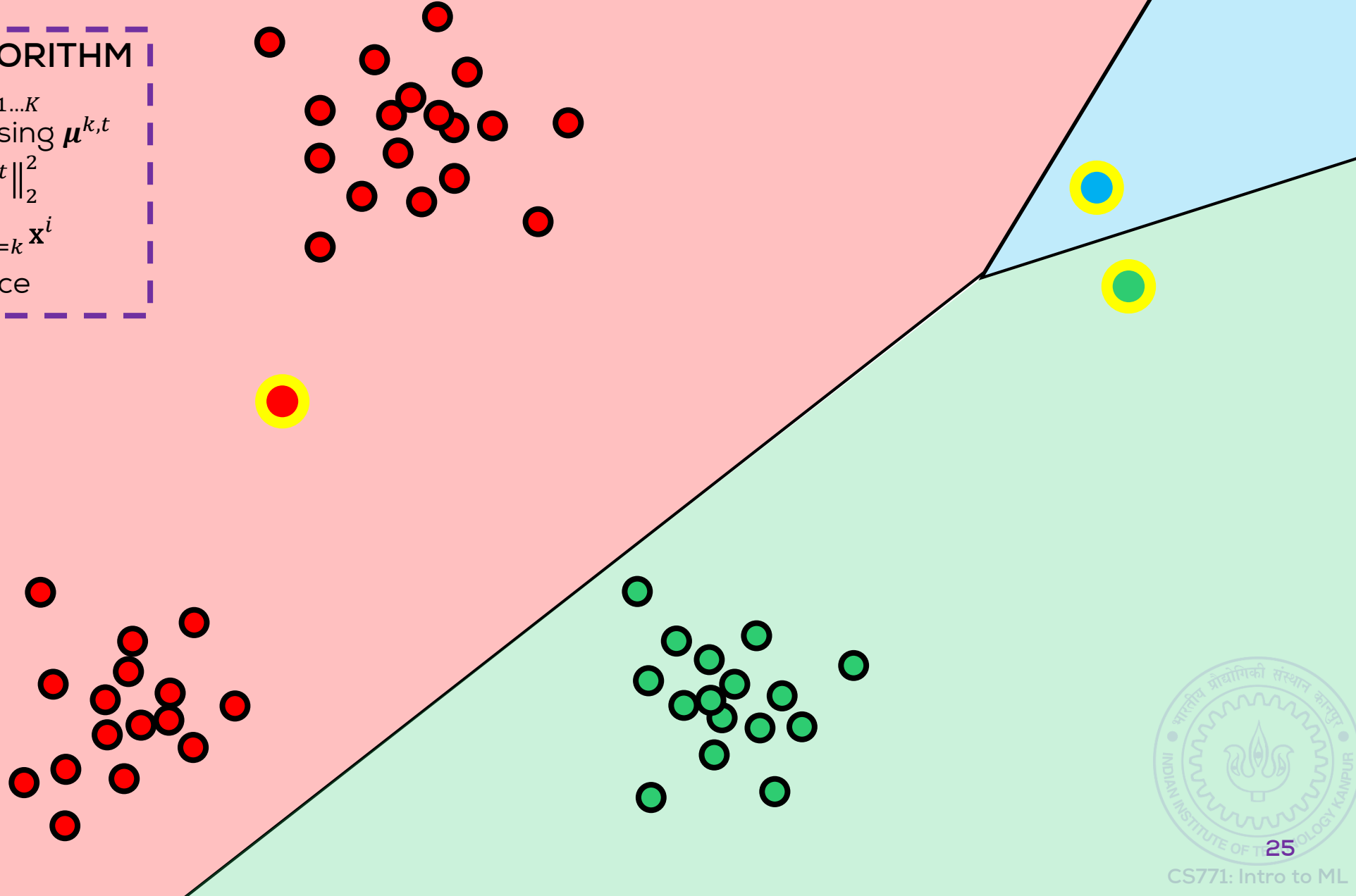
1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$   
Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



# K-Means Algorithm in action!

## K-MEANS/LLOYD'S ALGORITHM

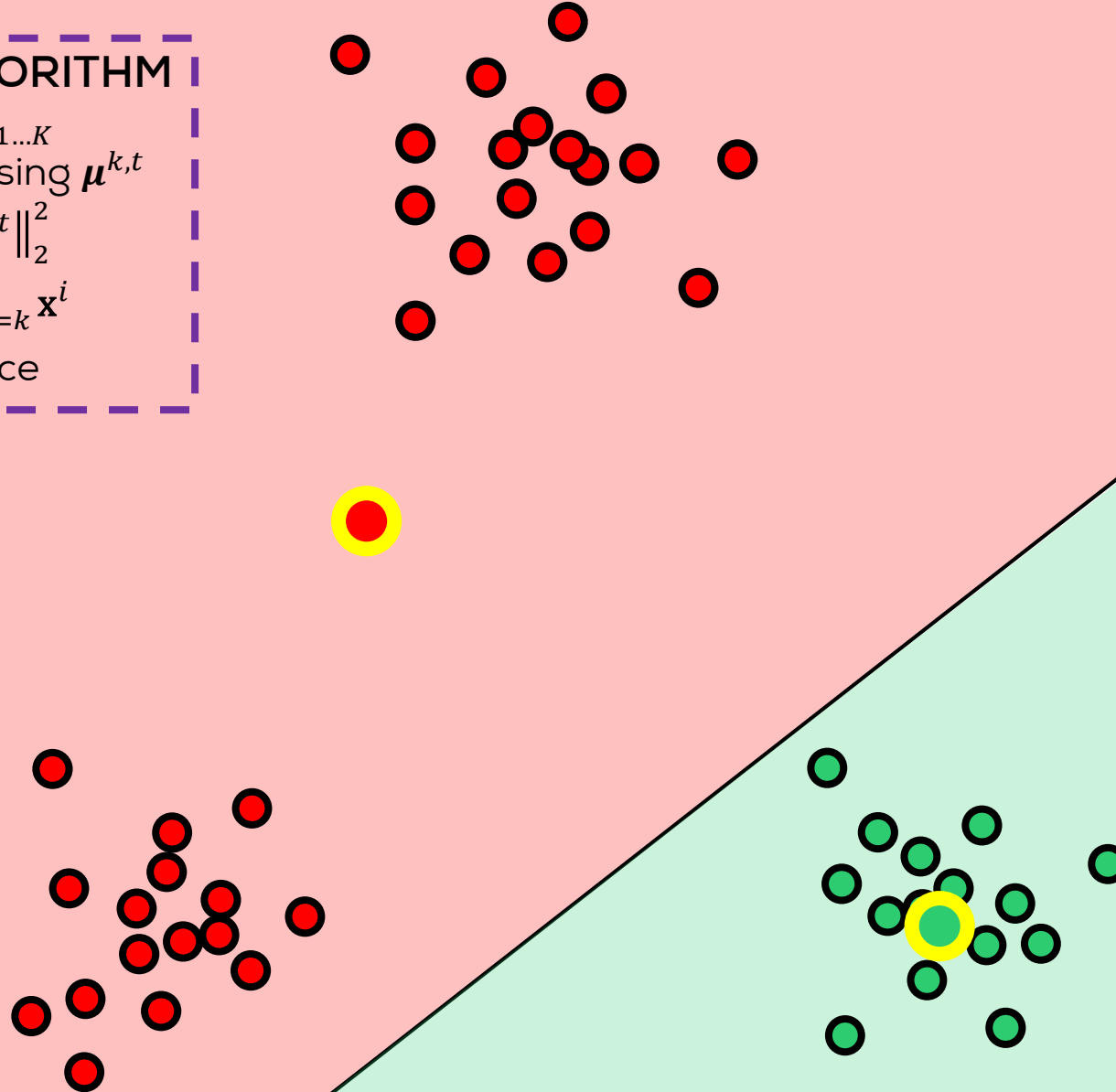
1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$   
Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i:z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



# K-Means Algorithm in action!

## K-MEANS/LLOYD'S ALGORITHM

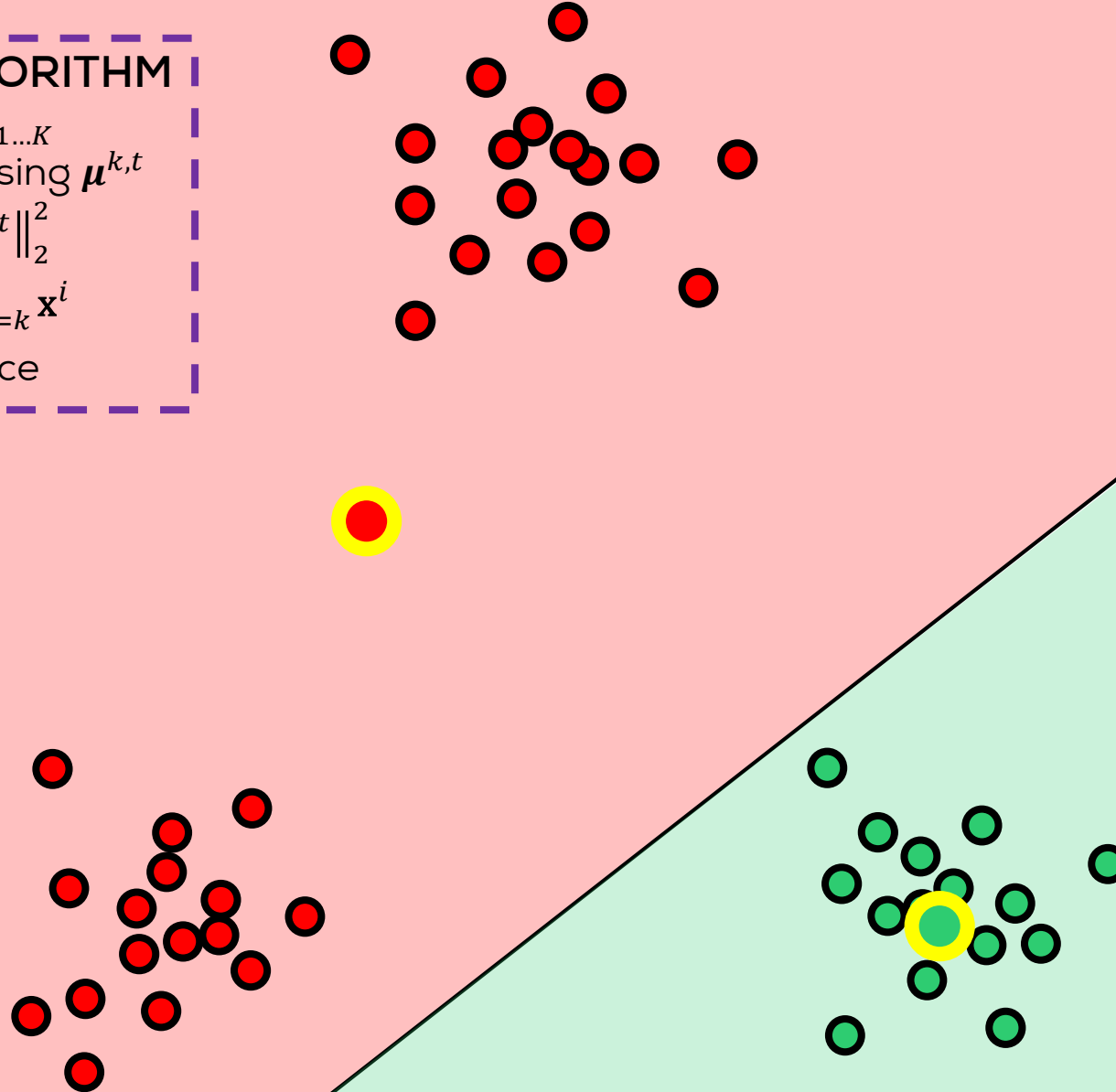
1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$   
Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



# K-Means Algorithm in action!

## K-MEANS/LLOYD'S ALGORITHM

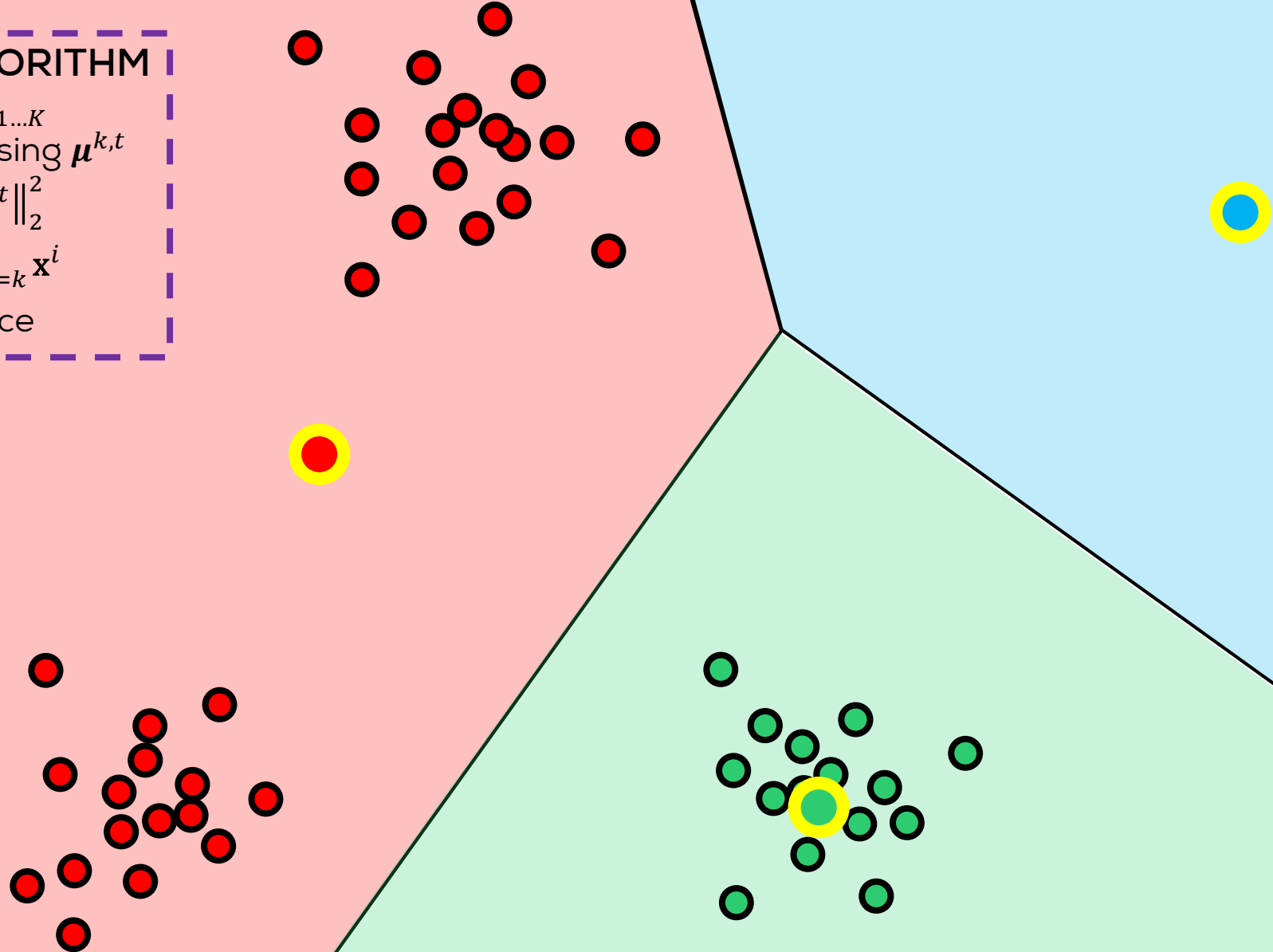
1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$   
Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



# K-Means Algorithm in action!

## K-MEANS/LLOYD'S ALGORITHM

1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$   
Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence

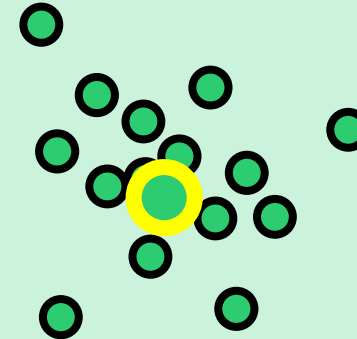
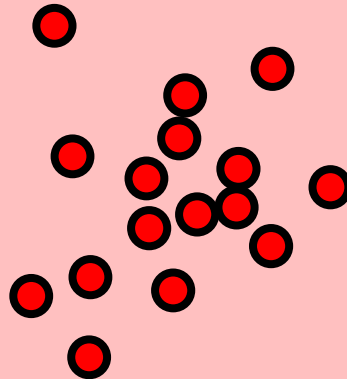
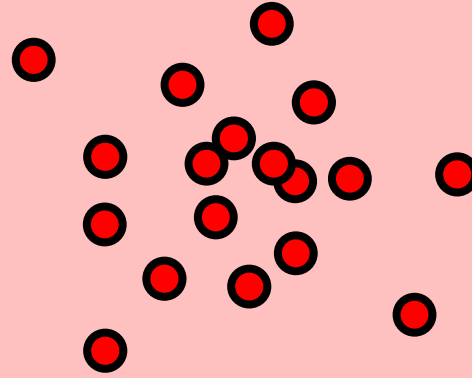


# K-Means Algorithm in action!

## K-MEANS/LLOYD'S ALGORITHM

1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$   
Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence

Stuck!!!





# K-Means Algorithm in action!

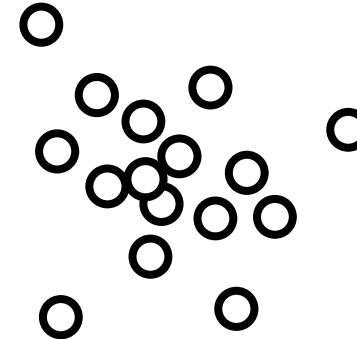
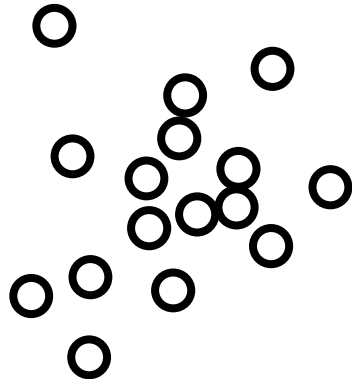
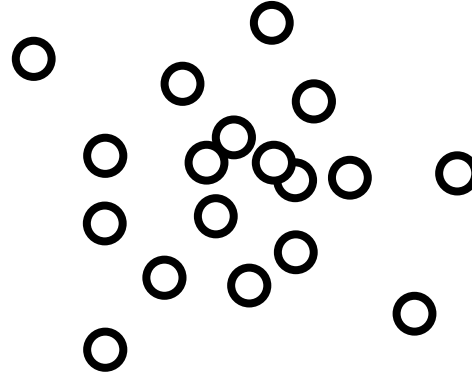
## K-MEANS/LLOYD'S ALGORITHM

1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$   
Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence

# K-Means Algorithm in action!

## K-MEANS/LLOYD'S ALGORITHM

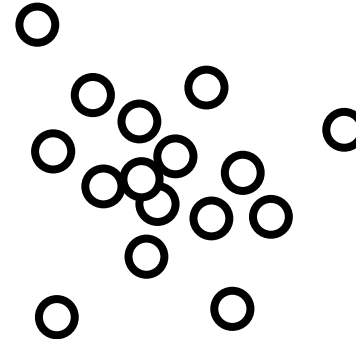
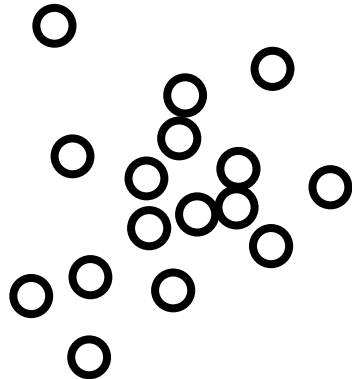
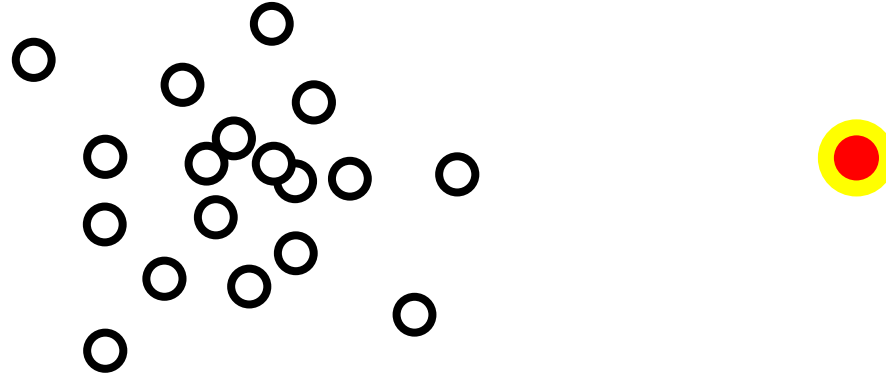
1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$   
Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i:z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



# K-Means Algorithm in action!

## K-MEANS/LLOYD'S ALGORITHM

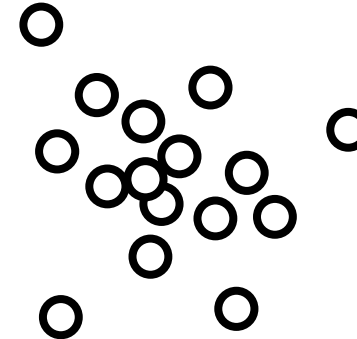
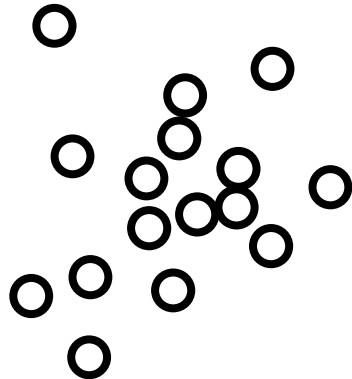
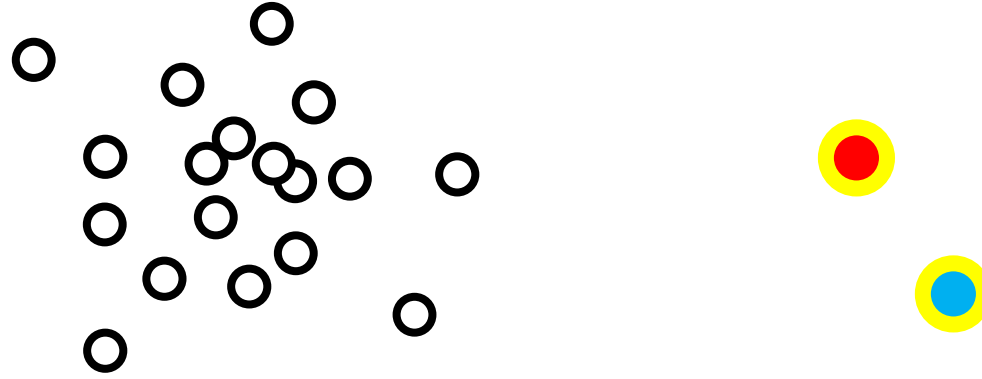
1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$   
Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i:z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



# K-Means Algorithm in action!

## K-MEANS/LLOYD'S ALGORITHM

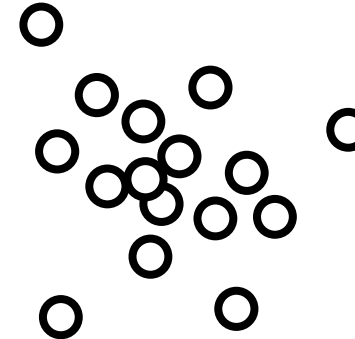
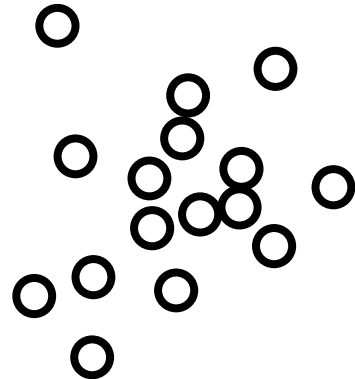
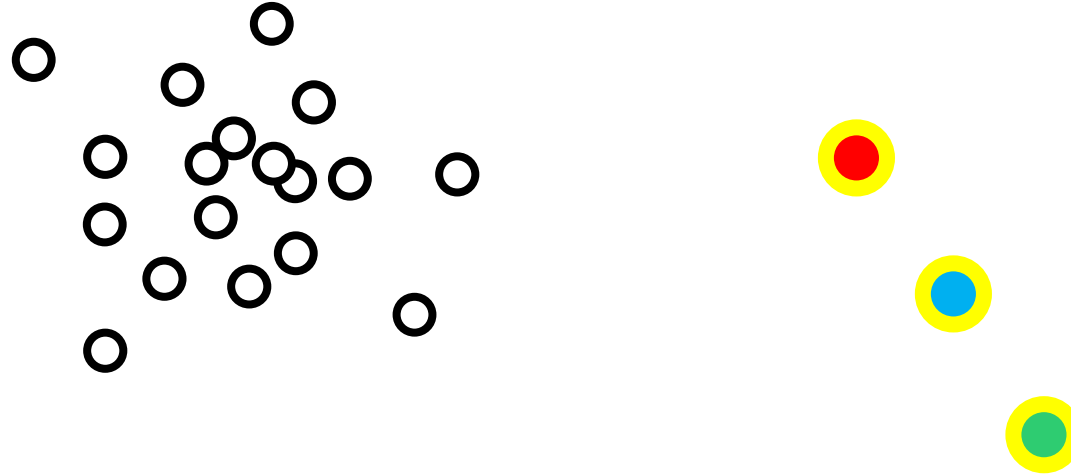
1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$   
Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i:z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



# K-Means Algorithm in action!

## K-MEANS/LLOYD'S ALGORITHM

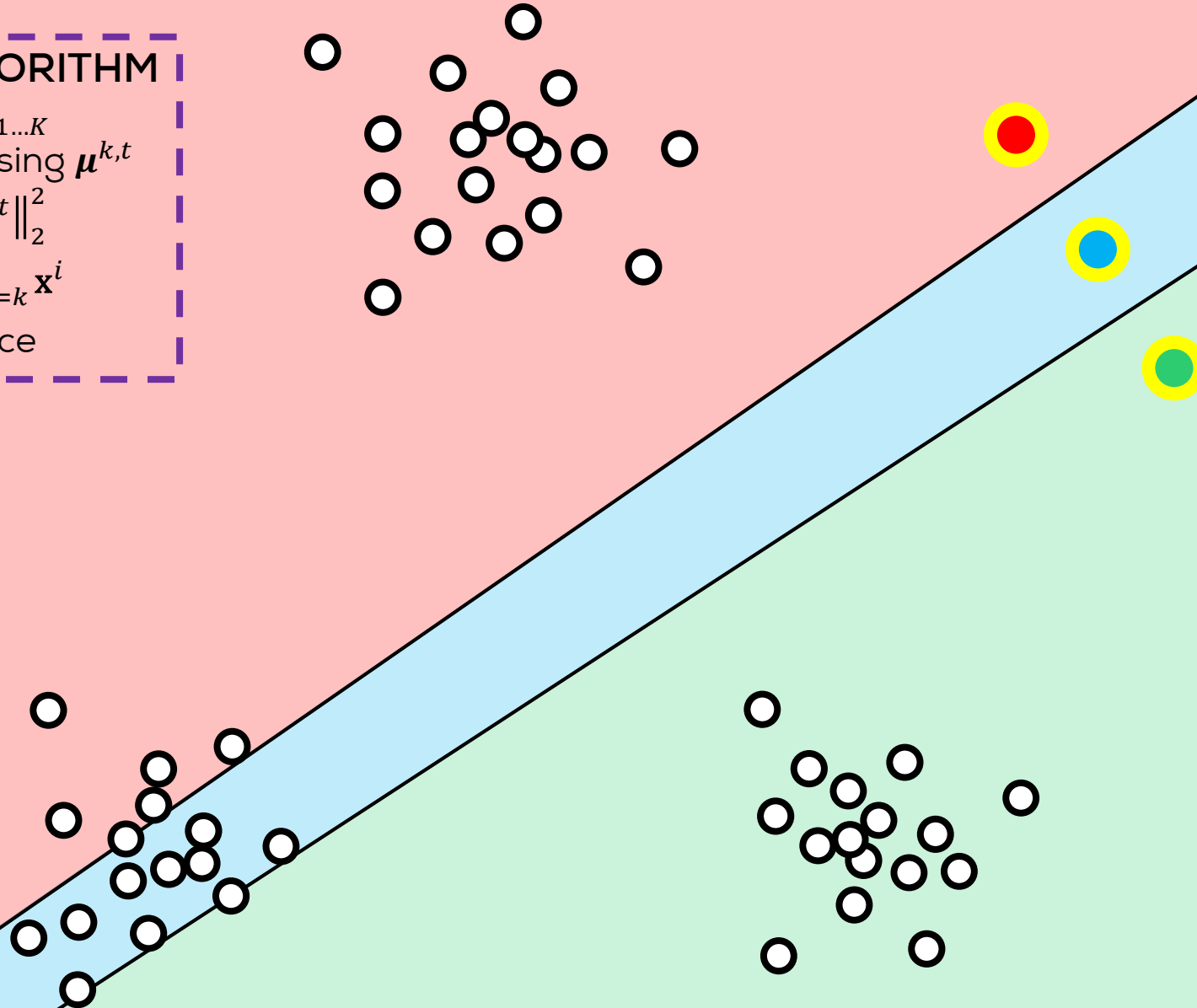
1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$   
Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i:z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



# K-Means Algorithm in action!

## K-MEANS/LLOYD'S ALGORITHM

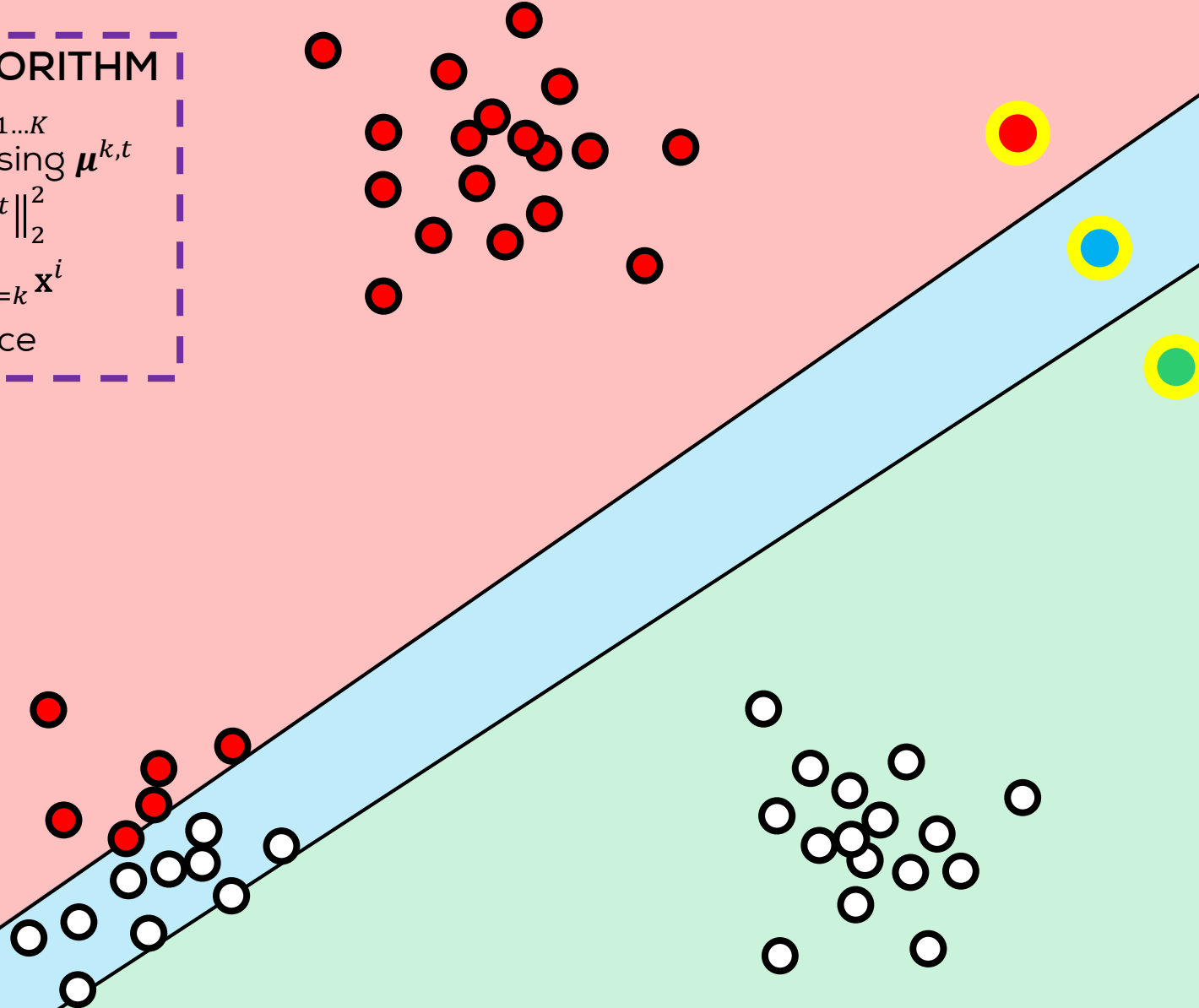
1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$   
Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i:z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



# K-Means Algorithm in action!

## K-MEANS/LLOYD'S ALGORITHM

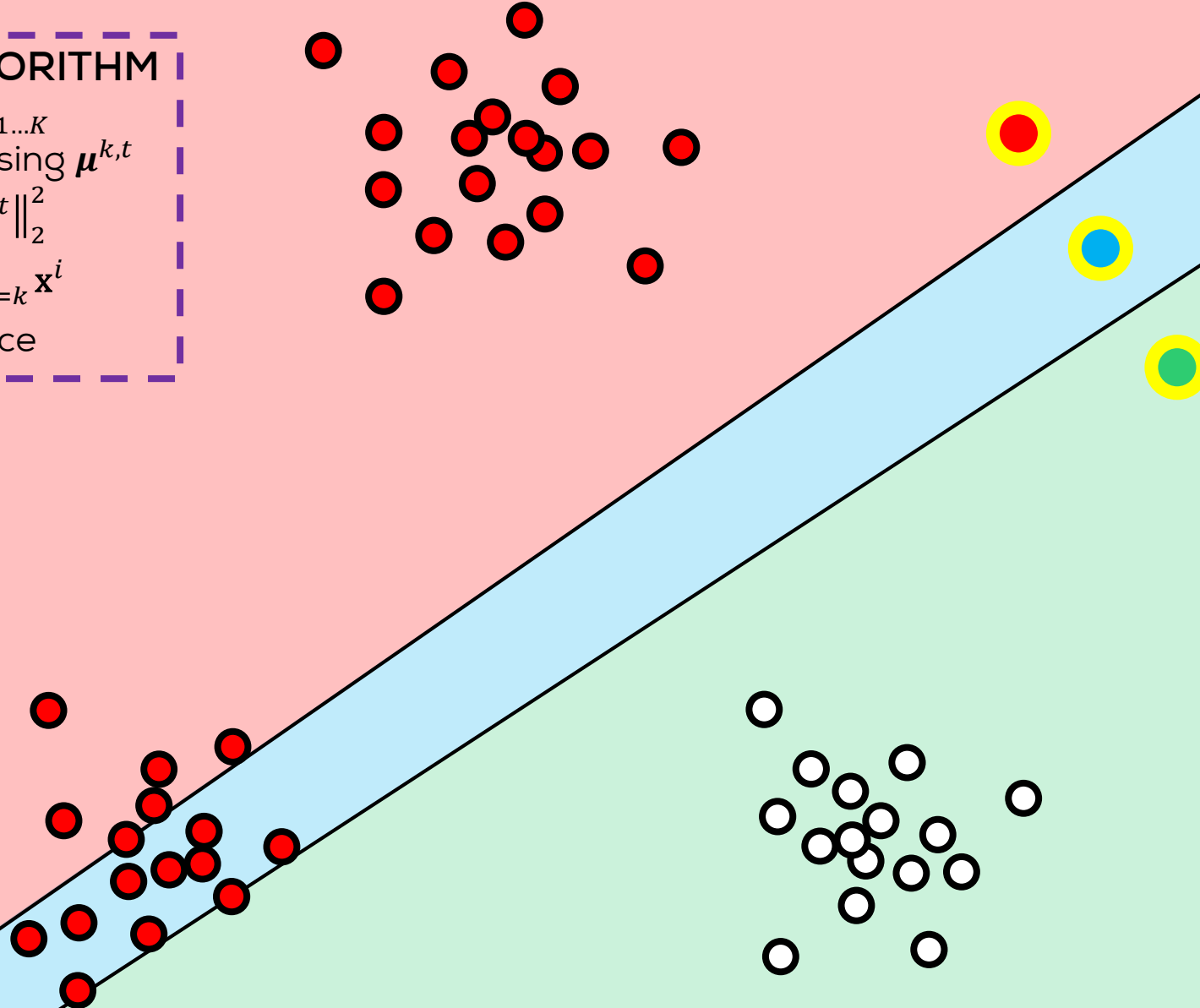
1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$   
Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i:z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



# K-Means Algorithm in action!

## K-MEANS/LLOYD'S ALGORITHM

1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$   
Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence

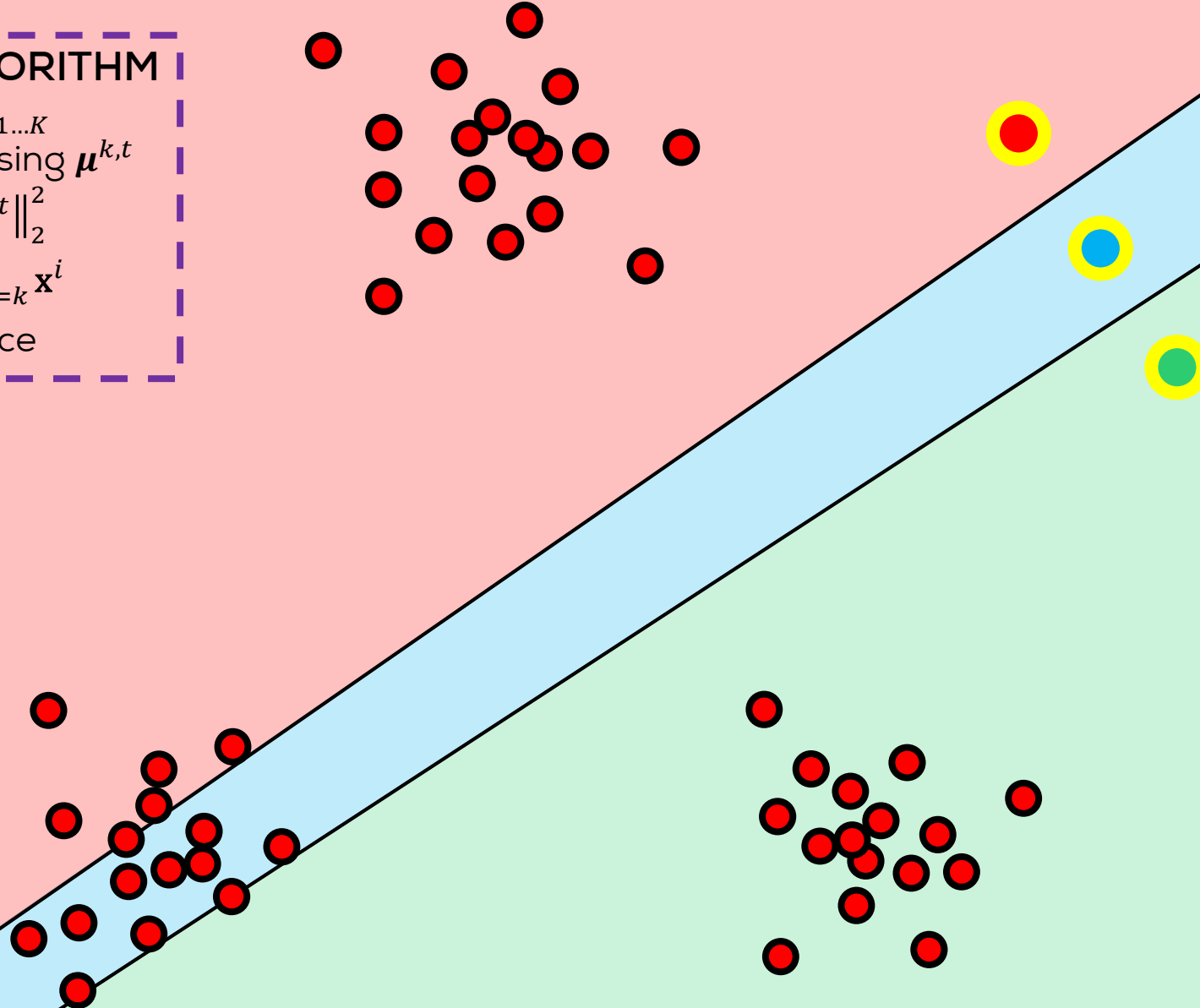




# K-Means Algorithm in action!

## K-MEANS/LLOYD'S ALGORITHM

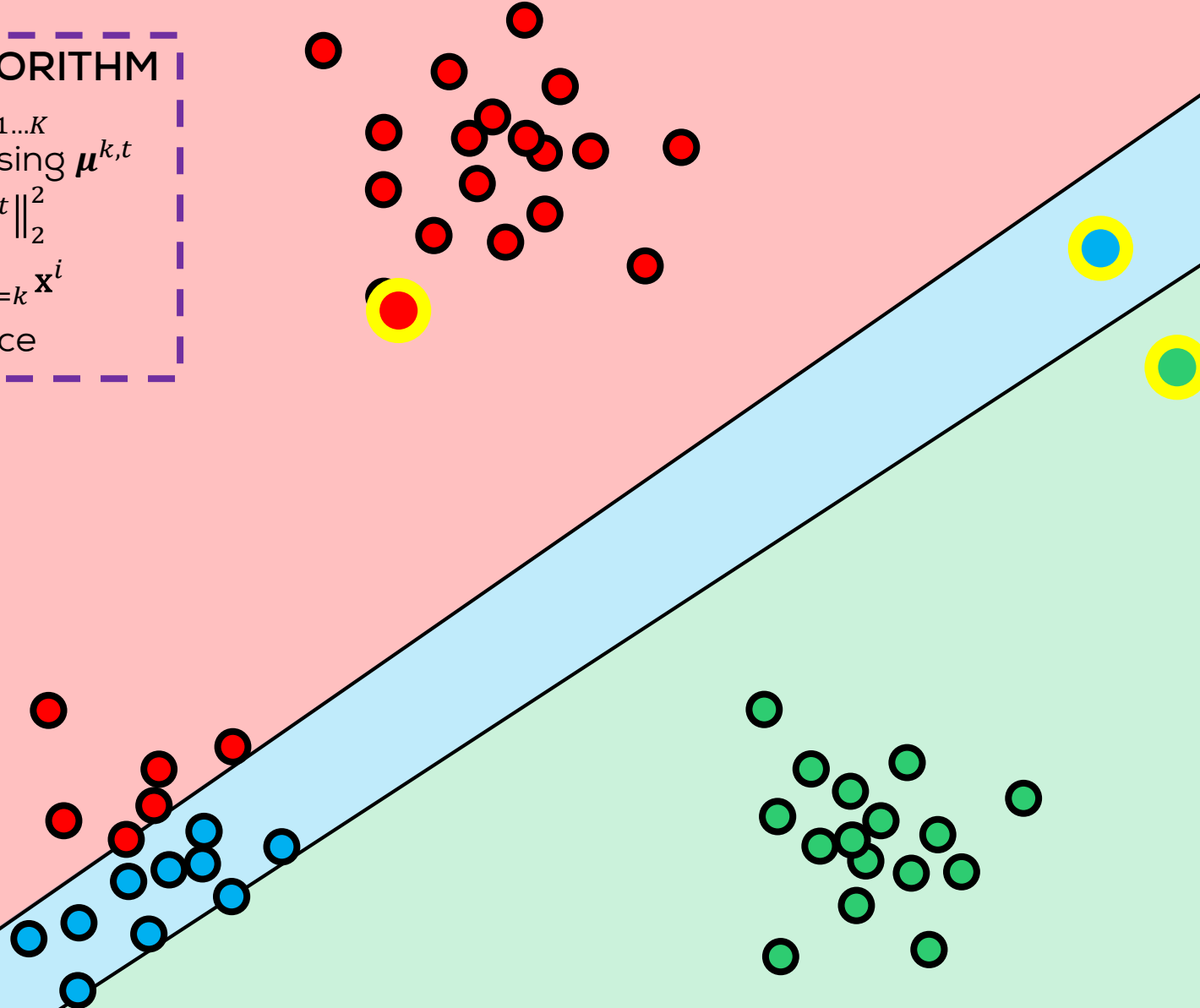
1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$   
Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



# K-Means Algorithm in action!

## K-MEANS/LLOYD'S ALGORITHM

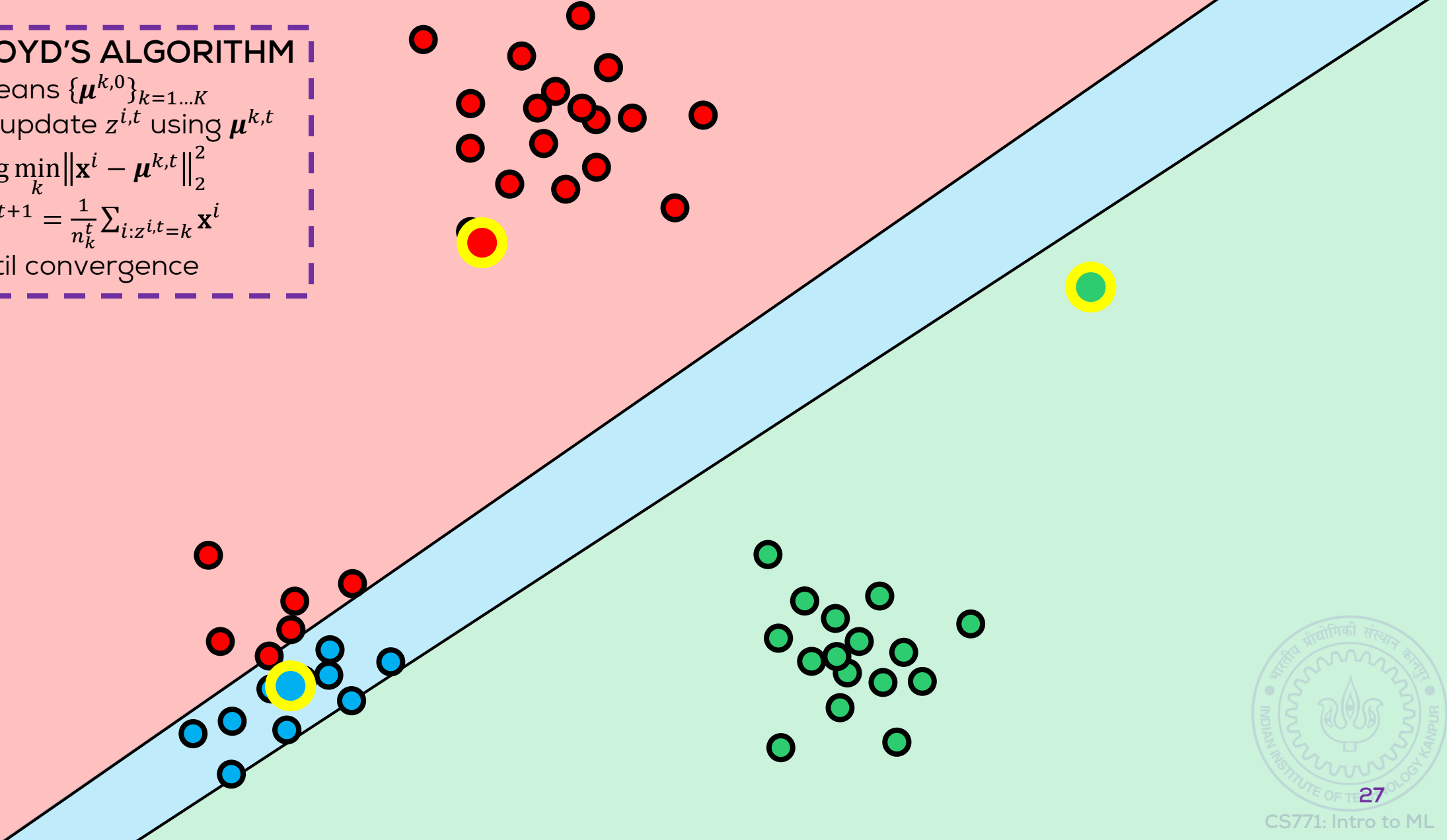
1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$   
Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



# K-Means Algorithm in action!

## K-MEANS/LLOYD'S ALGORITHM

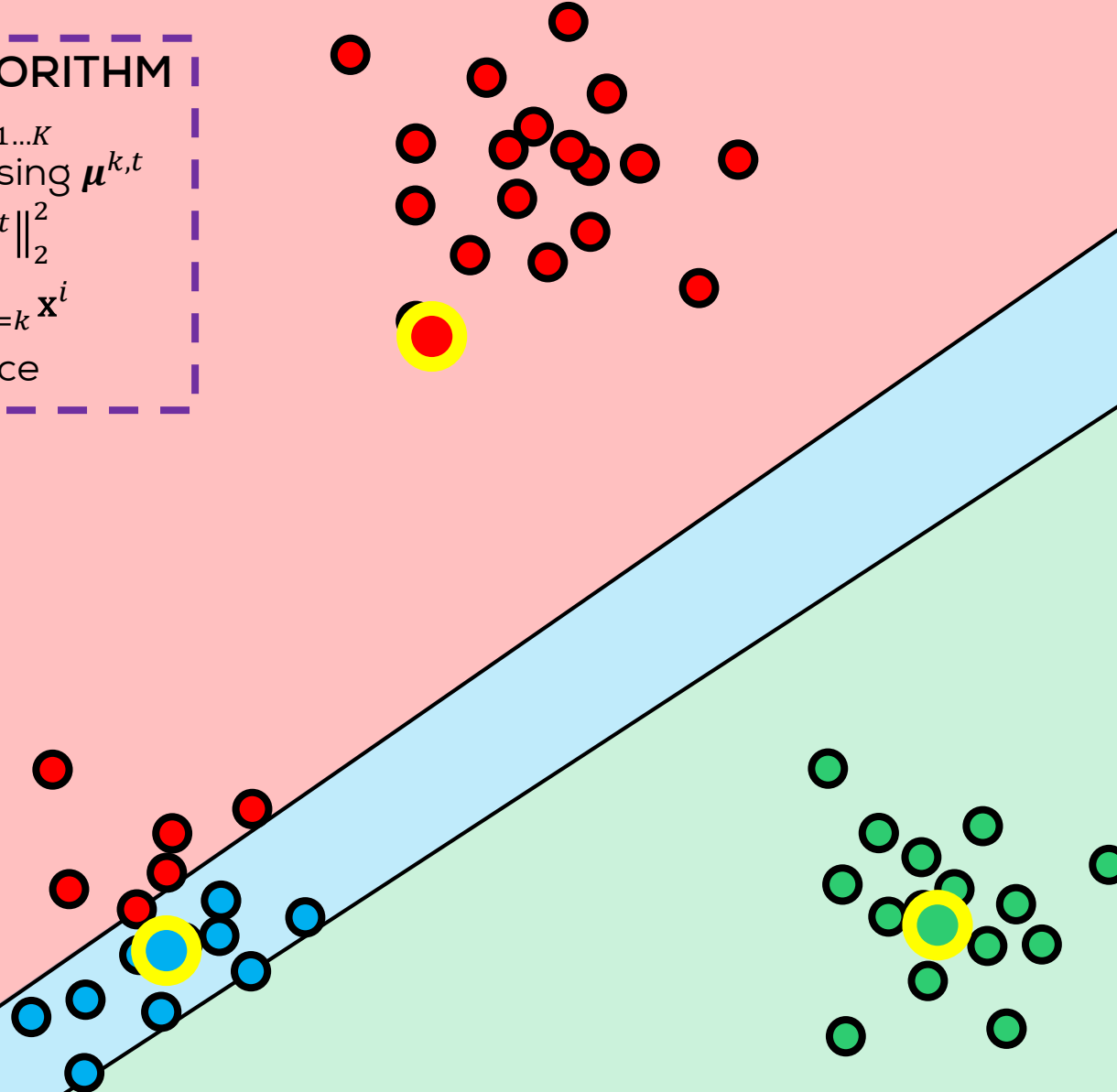
1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$   
Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



# K-Means Algorithm in action!

## K-MEANS/LLOYD'S ALGORITHM

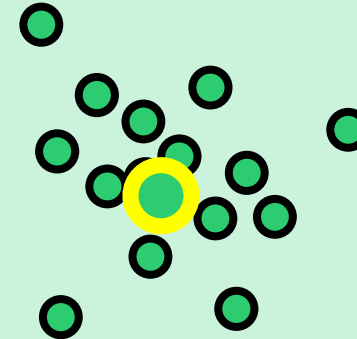
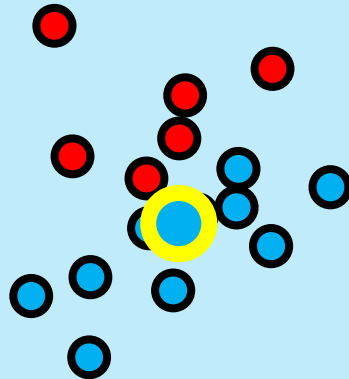
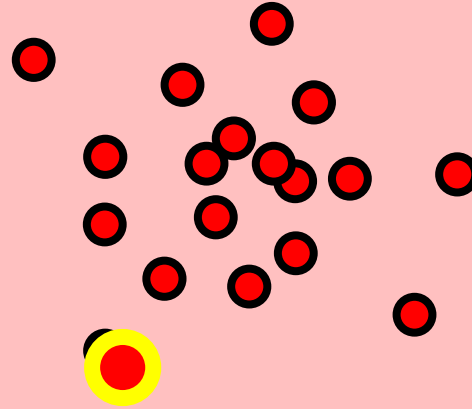
1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$   
Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



# K-Means Algorithm in action!

## K-MEANS/LLOYD'S ALGORITHM

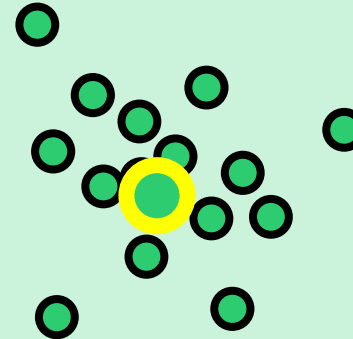
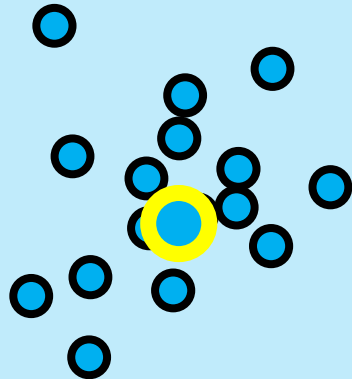
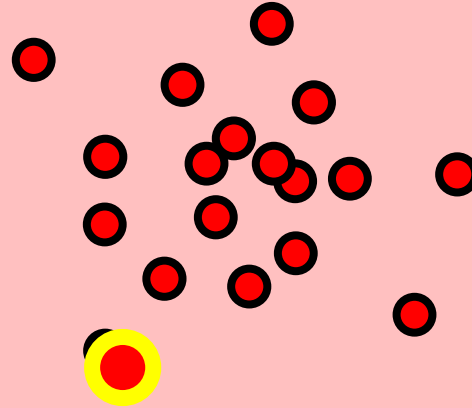
1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$   
Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



# K-Means Algorithm in action!

## K-MEANS/LLOYD'S ALGORITHM

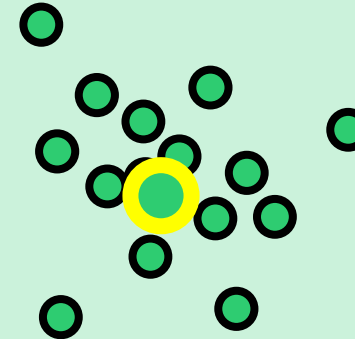
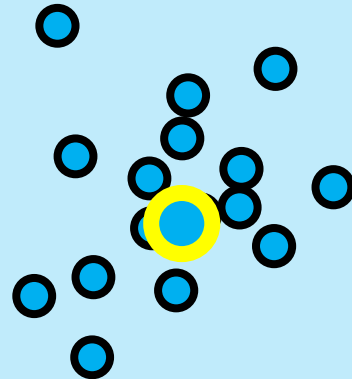
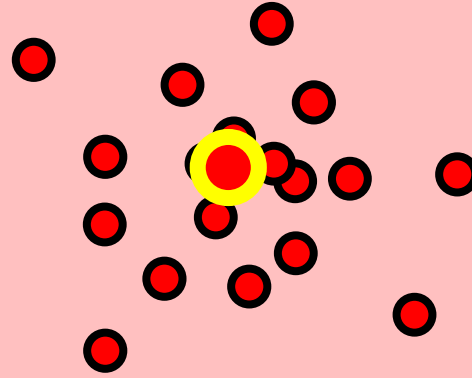
1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$   
Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



# K-Means Algorithm in action!

## K-MEANS/LLOYD'S ALGORITHM

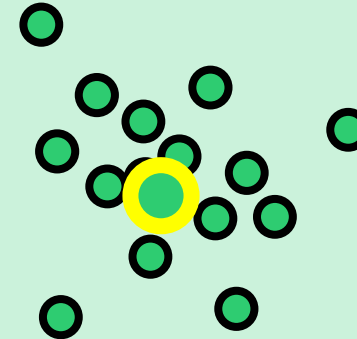
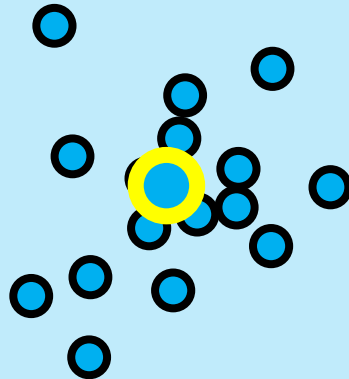
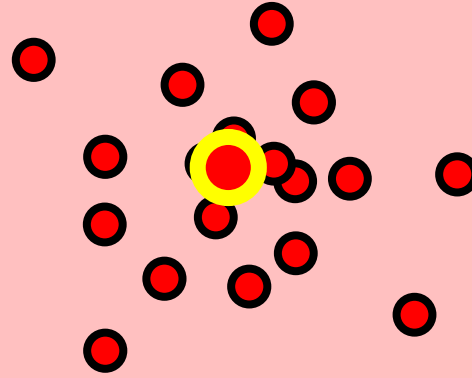
1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$   
Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i:z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



# K-Means Algorithm in action!

## K-MEANS/LLOYD'S ALGORITHM

1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$   
Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i:z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence

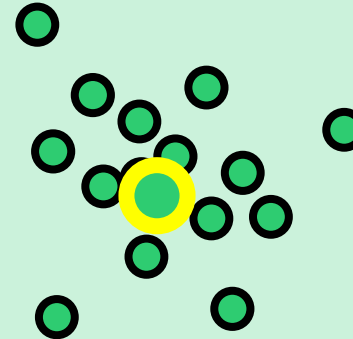
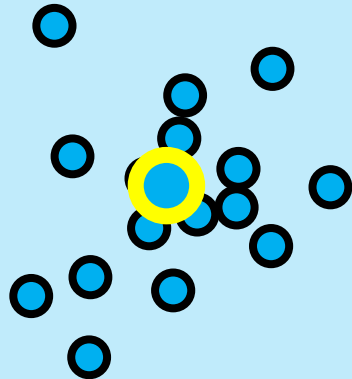
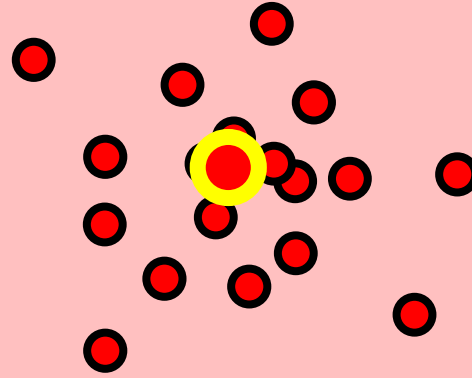




# K-Means Algorithm in action!

## K-MEANS/LLOYD'S ALGORITHM

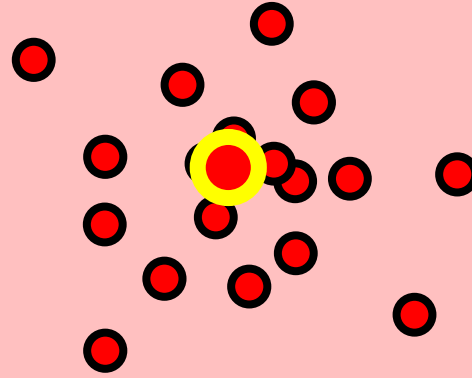
1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$   
Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



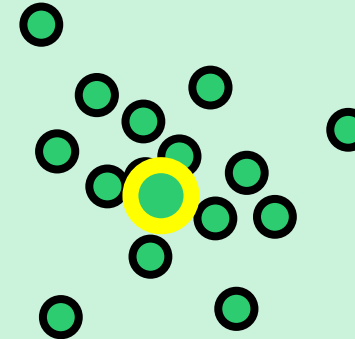
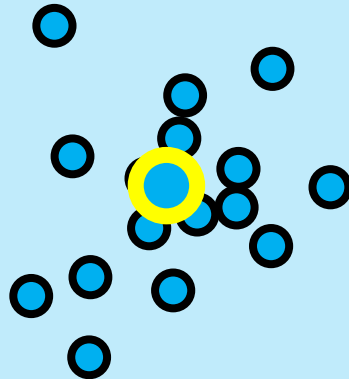
# K-Means Algorithm in action!

## K-MEANS/LLOYD'S ALGORITHM

1. Initialize means  $\{\mu^{k,0}\}_{k=1\dots K}$
2. For  $i \in [n]$ , update  $z^{i,t}$  using  $\mu^{k,t}$   
Let  $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update  $\mu^{k,t+1} = \frac{1}{n_k} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



Stuck!!! ... but  
at the global  
optimum 😊



# Generating data from GMM learnt using Lloyd

- $K$  clusters with means  $\mu^1, \mu^2, \dots, \mu^K$  learnt using k-means
- To generate a data point from this GMM
  1. Select a cluster  $k \sim [K]$  uniformly at random
  2. Select a point from the Gaussian  $\mathcal{N}(\mu^k, I)$

K-means  
uses  $\pi_k = \frac{1}{K}$

K-means  
use  $\Sigma^k = I$

# The K-means clustering algorithm

# The K-means clustering algorithm

- Extremely popular

# The K-means clustering algorithm

- Extremely popular
- Helps make sense of data

# The K-means clustering algorithm

- Extremely popular
- Helps make sense of data
- Can speed up k-NN computation quite a bit (assignment 😊)

# The K-means clustering algorithm

- Extremely popular
- Helps make sense of data
- Can speed up k-NN computation quite a bit (assignment 😊)
- Submission deadline: September 10, 2017, 2359 hrs, IST



# The K-means clustering algorithm

- Extremely popular
- Helps make sense of data
- Can speed up k-NN computation quite a bit (assignment 😊)
- Submission deadline: September 10, 2017, 2359 hrs, IST
- Works well if all clusters are balanced (comparable)

# The K-means clustering algorithm

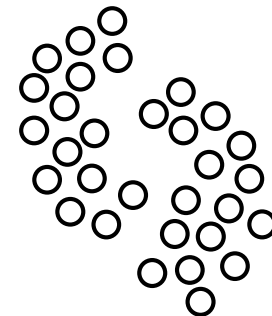
- Extremely popular
- Helps make sense of data
- Can speed up k-NN computation quite a bit (assignment 😊)
- Submission deadline: September 10, 2017, 2359 hrs, IST
- Works well if all clusters are balanced (comparable)
- Brittle in high dimensions and if many many small clusters

# The K-means clustering algorithm

- Extremely popular
- Helps make sense of data
- Can speed up k-NN computation quite a bit (assignment 😊)
- Submission deadline: September 10, 2017, 2359 hrs, IST
- Works well if all clusters are balanced (comparable)
- Brittle in high dimensions and if many many small clusters
- Learns clusters that are ball shaped, does not do well on non-convex shaped clusters

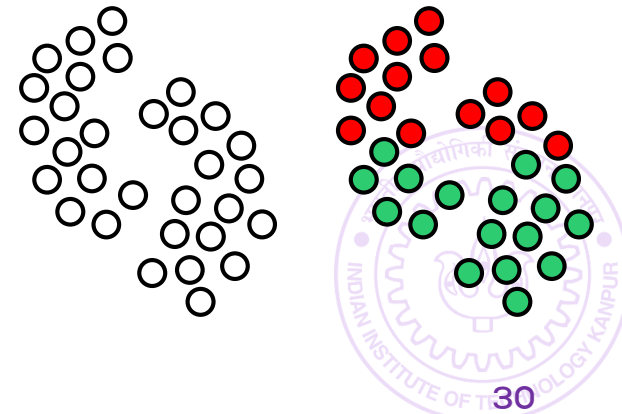
# The K-means clustering algorithm

- Extremely popular
- Helps make sense of data
- Can speed up k-NN computation quite a bit (assignment 😊)
- Submission deadline: September 10, 2017, 2359 hrs, IST
- Works well if all clusters are balanced (comparable)
- Brittle in high dimensions and if many many small clusters
- Learns clusters that are ball shaped, does not do well on non-convex shaped clusters



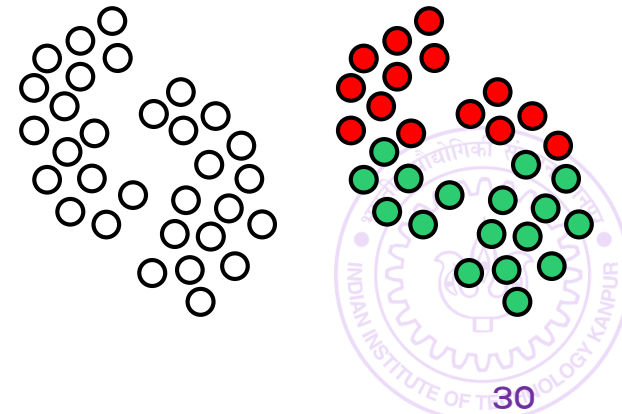
# The K-means clustering algorithm

- Extremely popular
- Helps make sense of data
- Can speed up k-NN computation quite a bit (assignment 😊)
- Submission deadline: September 10, 2017, 2359 hrs, IST
- Works well if all clusters are balanced (comparable)
- Brittle in high dimensions and if many many small clusters
- Learns clusters that are ball shaped, does not do well on non-convex shaped clusters



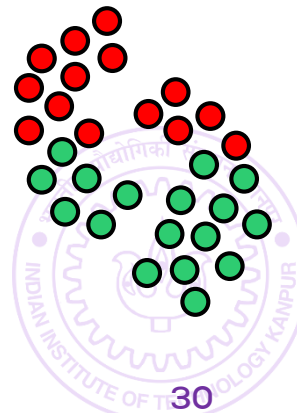
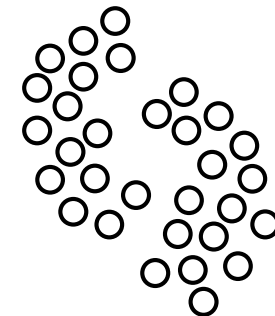
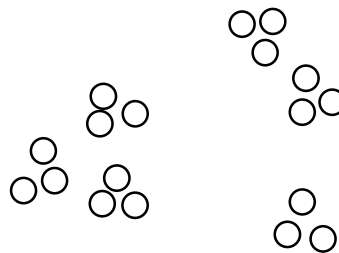
# The K-means clustering algorithm

- Extremely popular
- Helps make sense of data
- Can speed up k-NN computation quite a bit (assignment 😊)
- Submission deadline: September 10, 2017, 2359 hrs, IST
- Works well if all clusters are balanced (comparable)
- Brittle in high dimensions and if many many small clusters
- Learns clusters that are ball shaped, does not do well on non-convex shaped clusters
- Hierarchical clustering



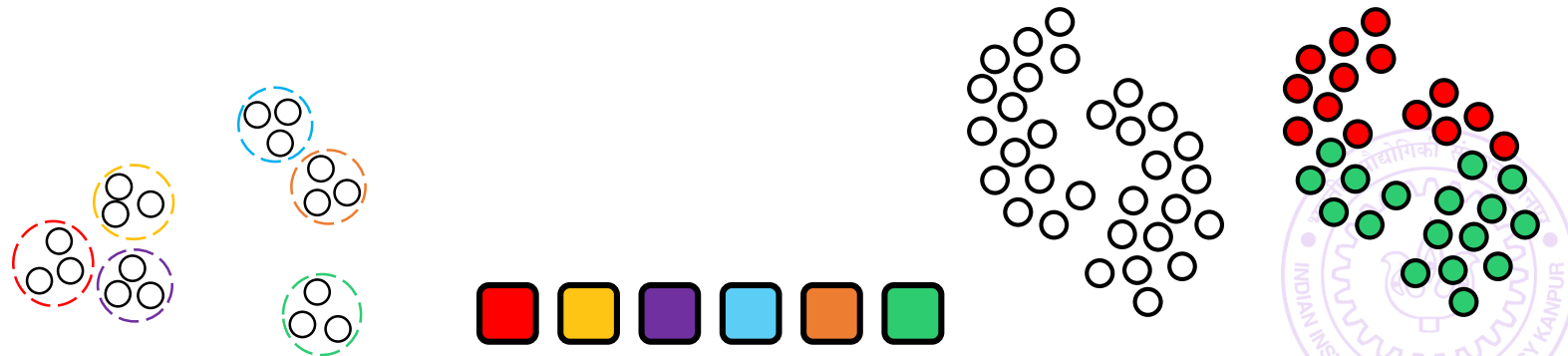
# The K-means clustering algorithm

- Extremely popular
- Helps make sense of data
- Can speed up k-NN computation quite a bit (assignment 😊)
- Submission deadline: September 10, 2017, 2359 hrs, IST
- Works well if all clusters are balanced (comparable)
- Brittle in high dimensions and if many many small clusters
- Learns clusters that are ball shaped, does not do well on non-convex shaped clusters
- Hierarchical clustering



# The K-means clustering algorithm

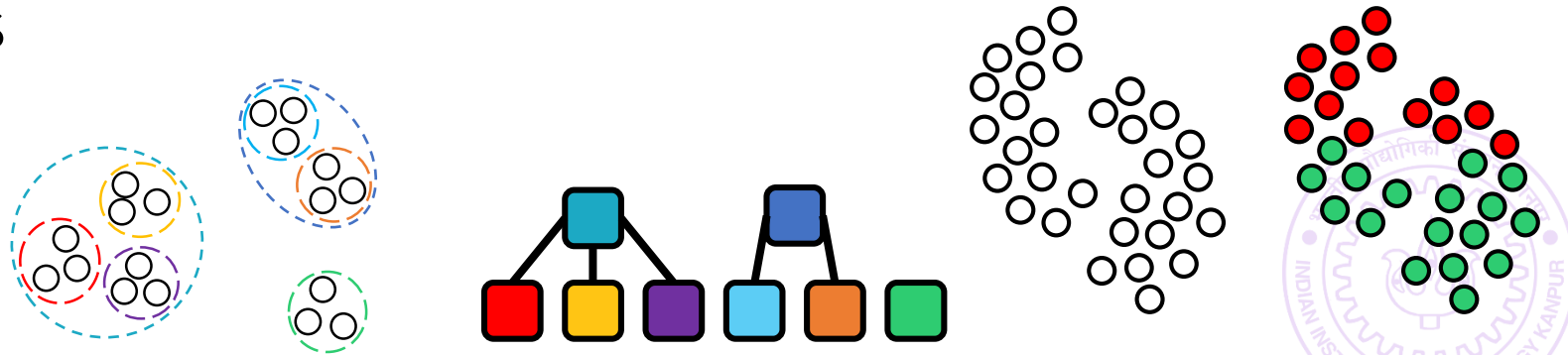
- Extremely popular
- Helps make sense of data
- Can speed up k-NN computation quite a bit (assignment 😊)
- Submission deadline: September 10, 2017, 2359 hrs, IST
- Works well if all clusters are balanced (comparable)
- Brittle in high dimensions and if many many small clusters
- Learns clusters that are ball shaped, does not do well on non-convex shaped clusters
- Hierarchical clustering





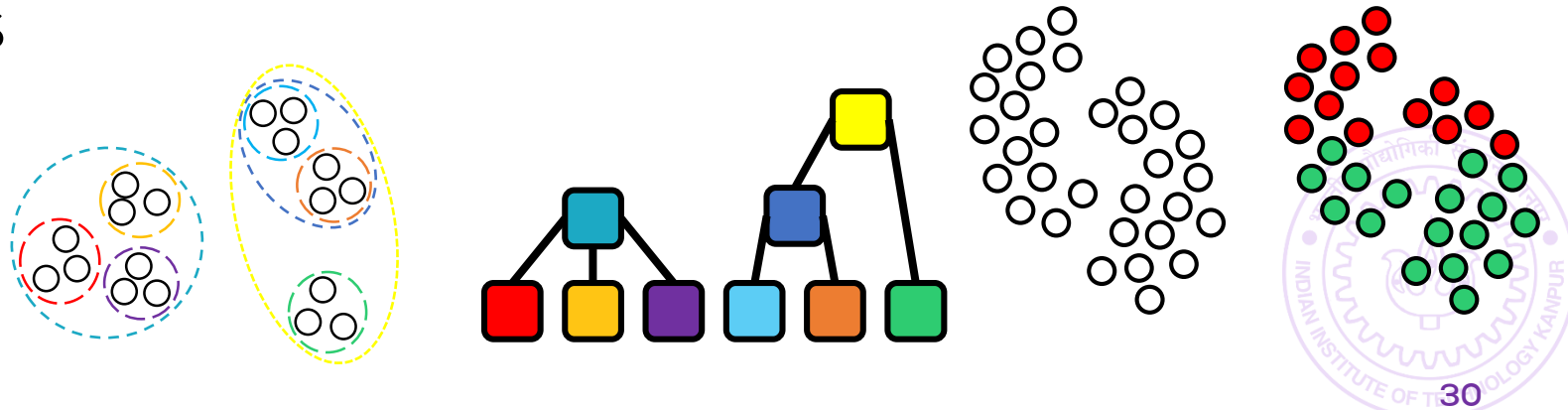
# The K-means clustering algorithm

- Extremely popular
- Helps make sense of data
- Can speed up k-NN computation quite a bit (assignment 😊)
- Submission deadline: September 10, 2017, 2359 hrs, IST
- Works well if all clusters are balanced (comparable)
- Brittle in high dimensions and if many many small clusters
- Learns clusters that are ball shaped, does not do well on non-convex shaped clusters
- Hierarchical clustering



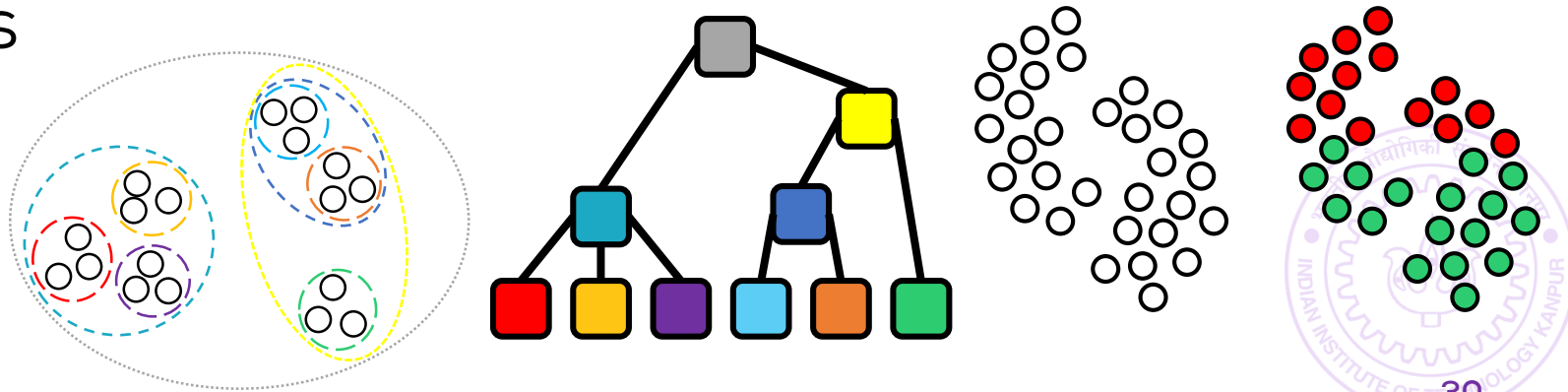
# The K-means clustering algorithm

- Extremely popular
- Helps make sense of data
- Can speed up k-NN computation quite a bit (assignment 😊)
- Submission deadline: September 10, 2017, 2359 hrs, IST
- Works well if all clusters are balanced (comparable)
- Brittle in high dimensions and if many many small clusters
- Learns clusters that are ball shaped, does not do well on non-convex shaped clusters
- Hierarchical clustering



# The K-means clustering algorithm

- Extremely popular
- Helps make sense of data
- Can speed up k-NN computation quite a bit (assignment 😊)
- Submission deadline: September 10, 2017, 2359 hrs, IST
- Works well if all clusters are balanced (comparable)
- Brittle in high dimensions and if many many small clusters
- Learns clusters that are ball shaped, does not do well on non-convex shaped clusters
- Hierarchical clustering



# Please give your Feedback

<http://tinyurl.com/ml17-18afb>