

Expectation-Maximization (and related learning methods)

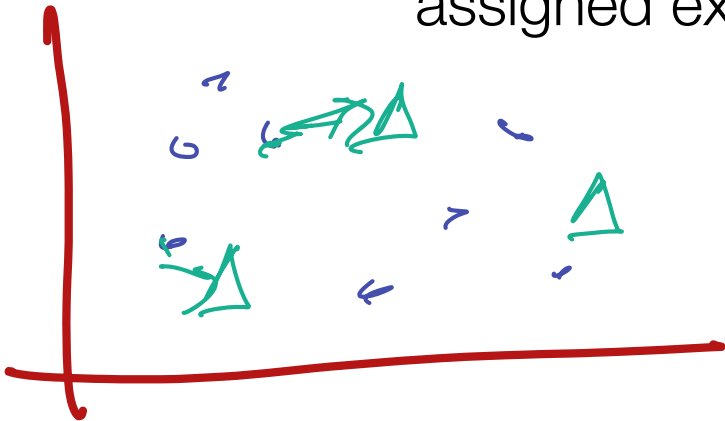
Lecture 18
CS 689, Spring 2023

Brendan O'Connor
College of Information and Computer Sciences
University of Massachusetts Amherst

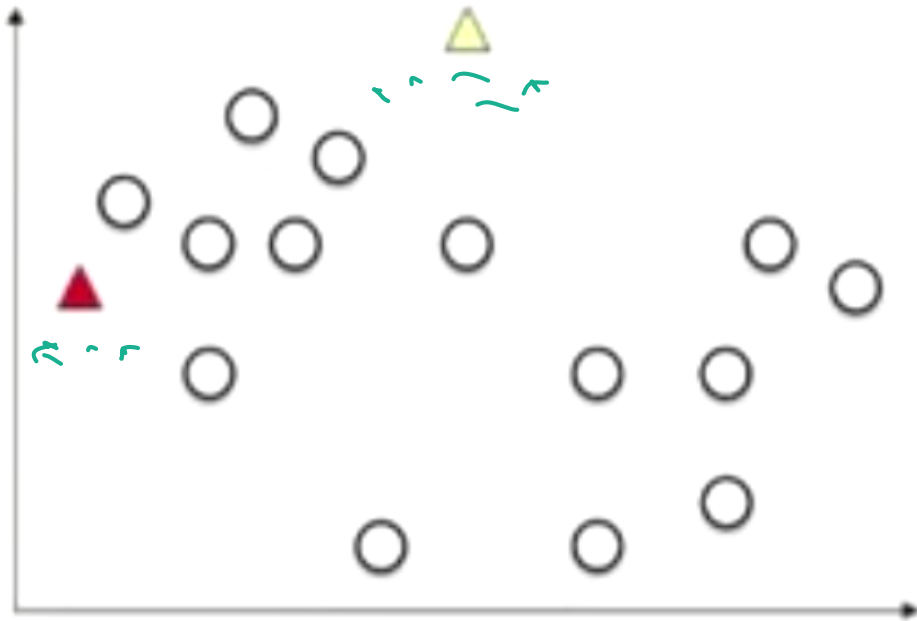
Clustering with (hard) EM

- ~~K~~-Means is the most basic example of a (basically) probabilistic unsupervised learning algorithm
 - 1. Randomly initialize cluster centroids
 - 2. Alternate until convergence:
 - (“E”): Assign each example to closest centroid
 - (“M”): Update centroids to means of these newly assigned examples

Coord
interior
or
P_{xy}

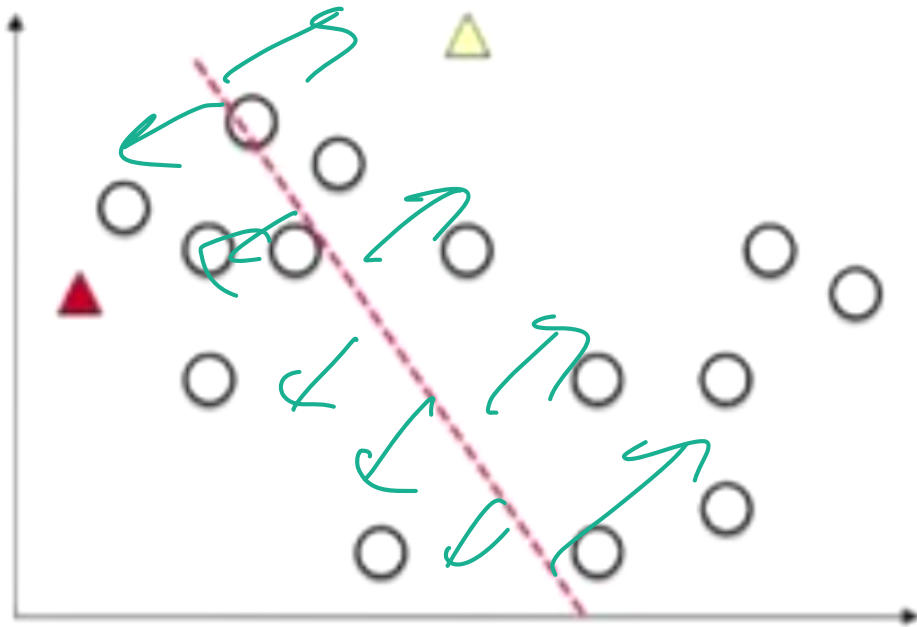


K-means clustering example



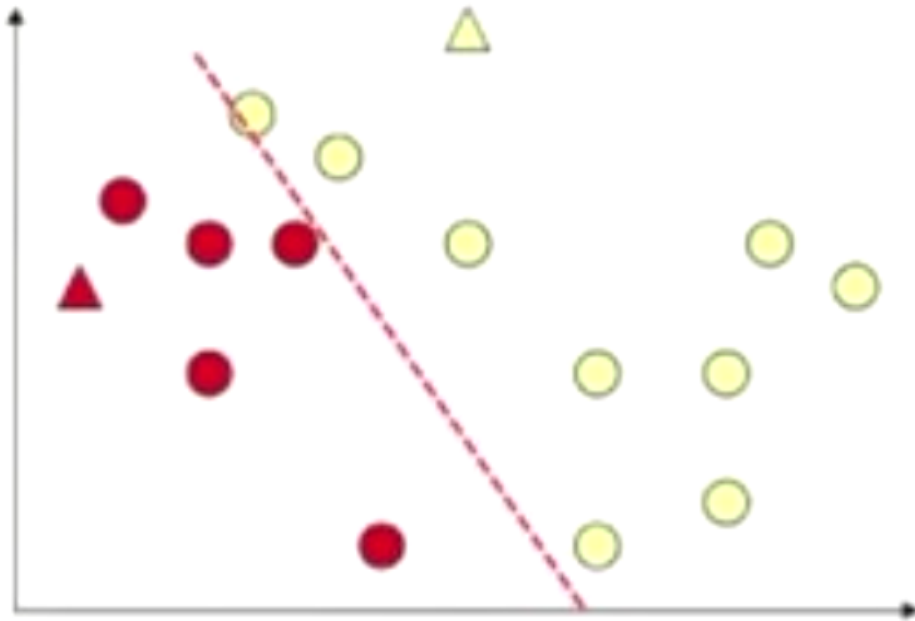
Slides from UMass alum Victor
Lavrenko, U. Edinburgh:
[https://www.youtube.com/watch?
v= aWzGGNrcic](https://www.youtube.com/watch?v=aWzGGNrcic)

K-means clustering example



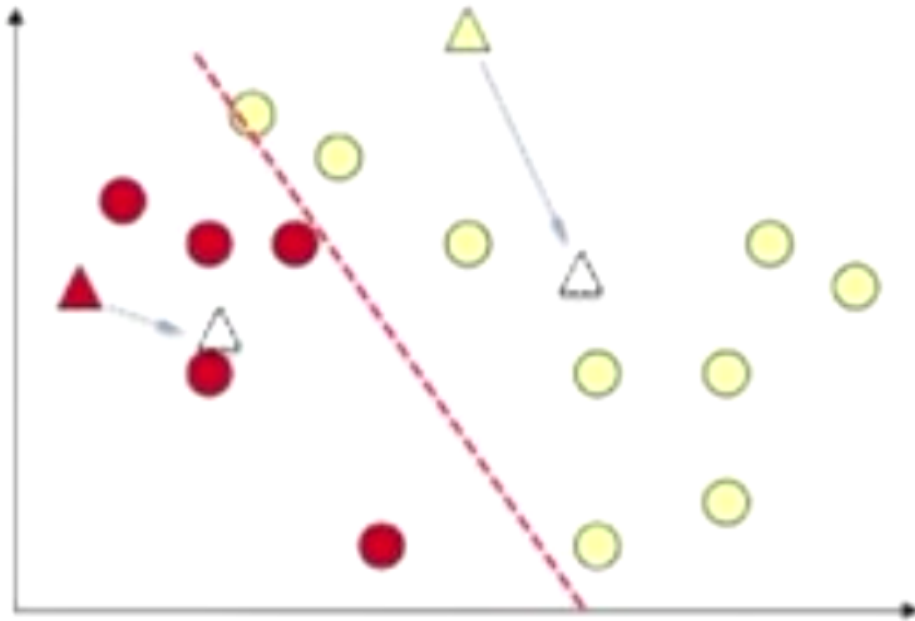
Slides from UMass alum Victor
Lavrenko, U. Edinburgh:
[https://www.youtube.com/watch?
v= aWzGGNrcic](https://www.youtube.com/watch?v=aWzGGNrcic)

K-means clustering example



E-step: infer labels
M-step: update centroids

K-means clustering example

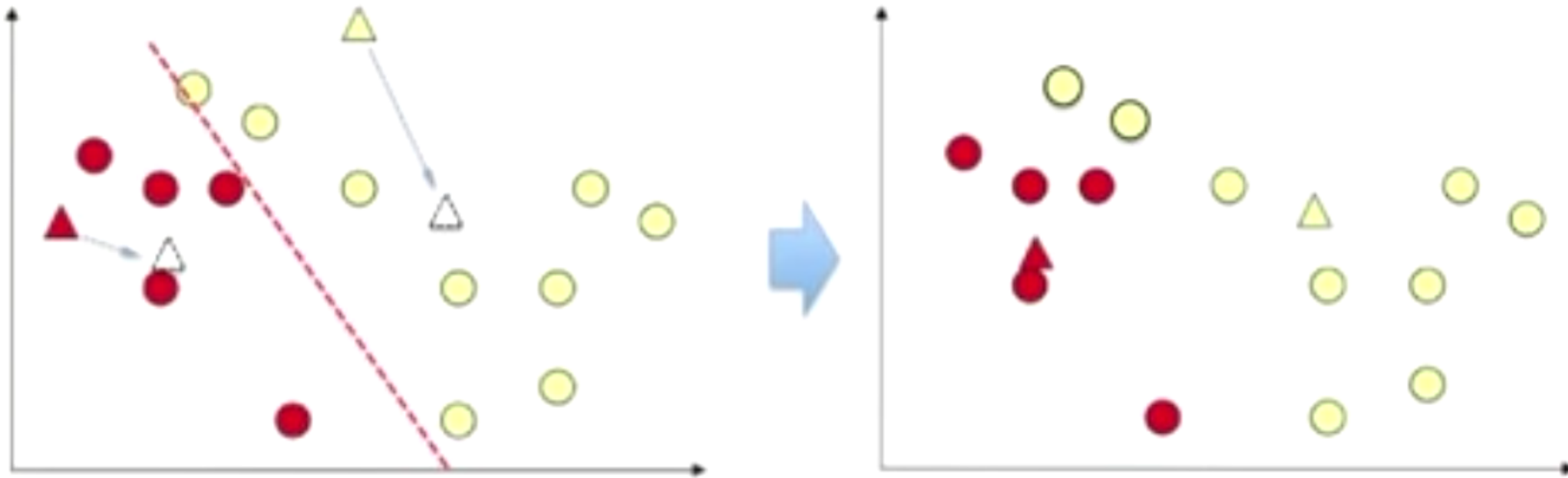


Slides from UMass alum Victor
Lavrenko, U. Edinburgh:

[https://www.youtube.com/watch?](https://www.youtube.com/watch?v=aWzGGNrcic)

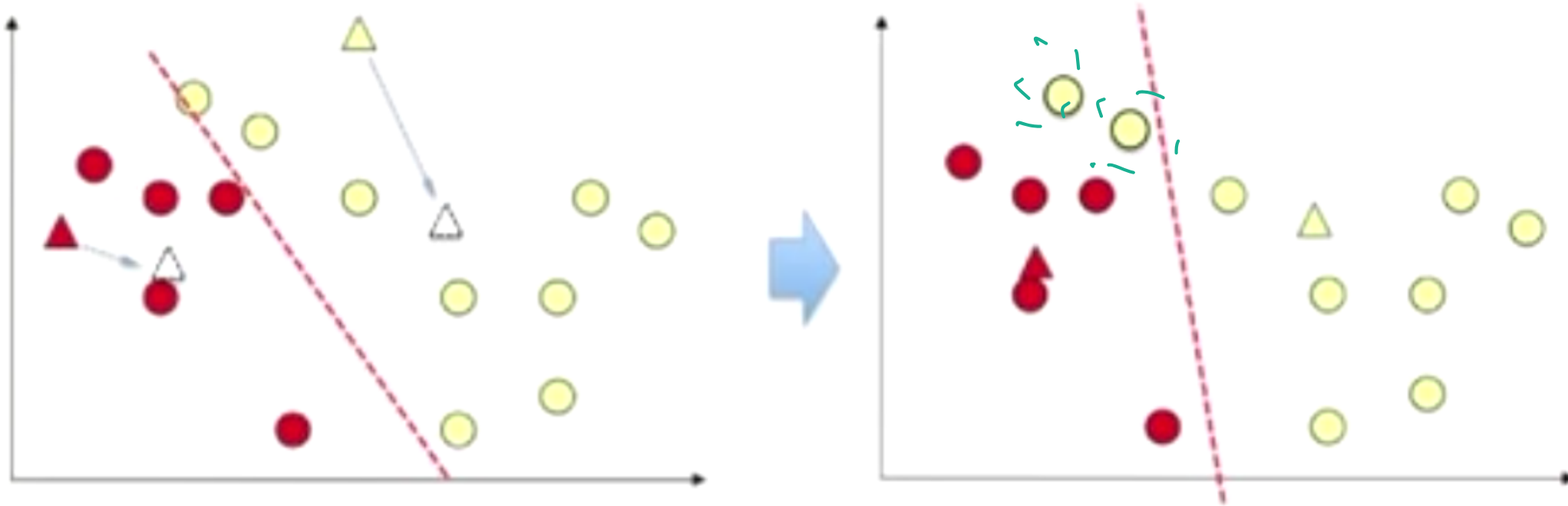
[v= aWzGGNrcic](https://www.youtube.com/watch?v=aWzGGNrcic)

K-means clustering example



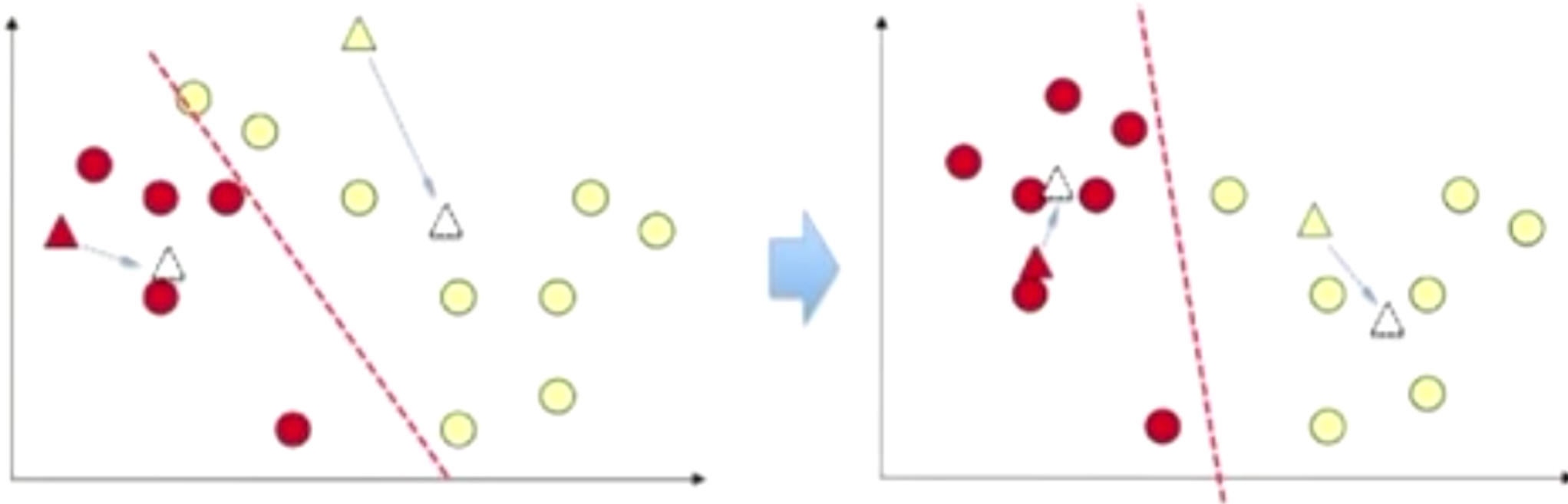
Slides from UMass alum Victor
Lavrenko, U. Edinburgh:
[https://www.youtube.com/watch?
v= aWzGGNrcic](https://www.youtube.com/watch?v=aWzGGNrcic)

K-means clustering example



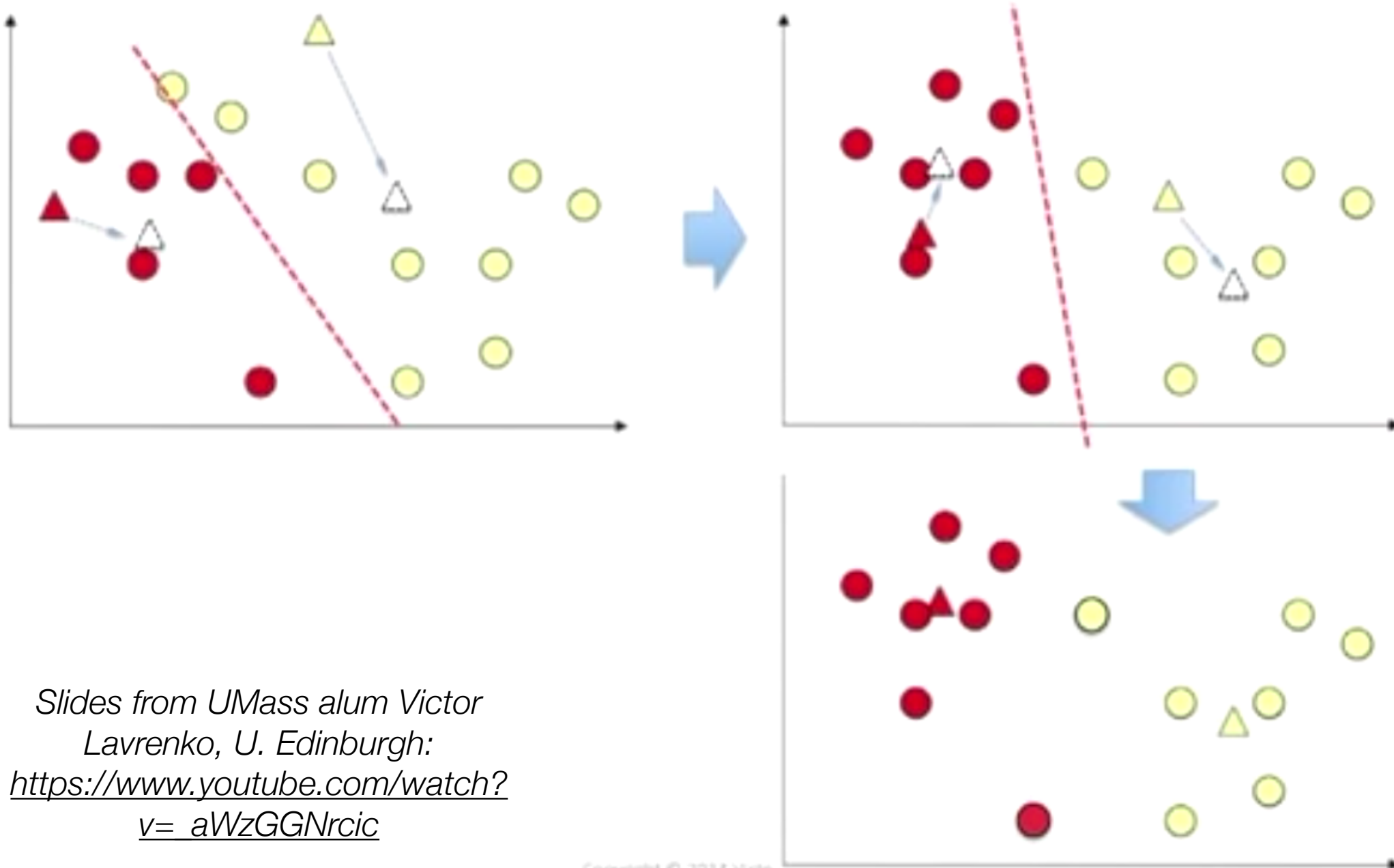
Slides from UMass alum Victor
Lavrenko, U. Edinburgh:
[https://www.youtube.com/watch?
v= aWzGGNrcic](https://www.youtube.com/watch?v=aWzGGNrcic)

K-means clustering example



Slides from UMass alum Victor
Lavrenko, U. Edinburgh:
[https://www.youtube.com/watch?
v= aWzGGNrcic](https://www.youtube.com/watch?v=aWzGGNrcic)

K-means clustering example



Slides from UMass alum Victor
Lavrenko, U. Edinburgh:
[https://www.youtube.com/watch?
v= aWzGGNrcic](https://www.youtube.com/watch?v=aWzGGNrcic)

K-Means as Gaussian Mixture

- Observed data $\mathbf{x}_1 \dots \mathbf{x}_n$
- Latent variables: cluster labels $\mathbf{z}_1 \dots \mathbf{z}_n$
- Parameters: Gaussian centroids $\mu_1 \dots \mu_k$
- Assume Gaussian mixture model
 $p(x_i | z_i) \sim N(\mu_{z_i}, \text{var})$ *constant* μ_{z_i}
- Learning algorithm: alternate until convergence,

- ("E"): Assign each example to closest centroid
 $\Rightarrow z_i := \arg\max_k P(z_i=k | x, \mu_k)$

- ("M"): Update centroids to averages
 $\Rightarrow \mu_k := \arg\max_m P(x | z, \mu=m)$
 $= (1/n_k) \sum_{i: z_i=k} x_i$

- Is there a "Soft" EM (close variant) iteratively optimizes $P(x, z | \mu)$
- This is "Hard EM". A very close variant

$$x_i \in \mathbb{R}^d$$

$$z_i \in \{1, \dots, k\}$$

$$\mu_k \in \mathbb{R}^d$$

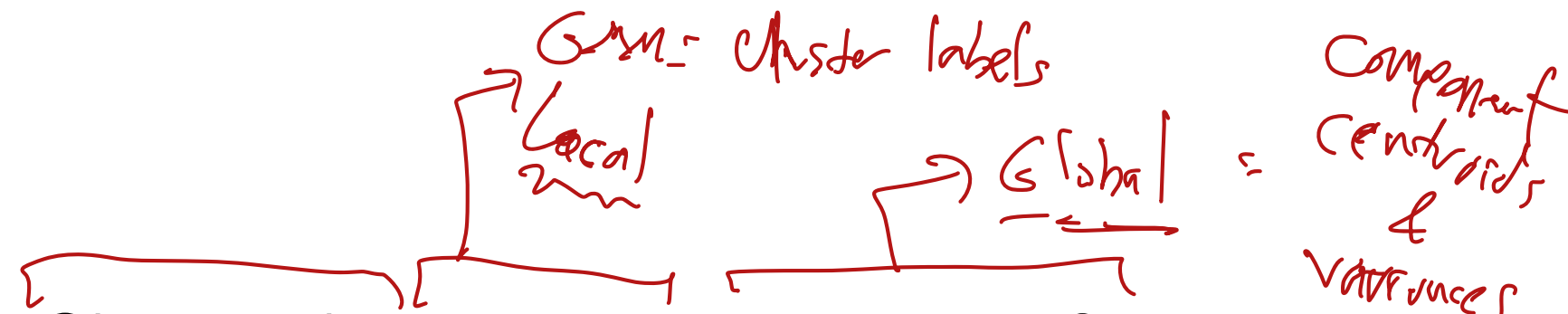


$$P(z) = \text{uniform}$$

$$P(z=1|x, \mu) \propto p(x|z=1)p(z=1)$$

$$P(z=2|x, \mu) \propto p(x|z=2)p(z=2)$$

Expectation-Maximization

- 
- Observed \mathbf{x} , latent \mathbf{z} , parameters $\boldsymbol{\theta}$
 - EM is a meta-algorithm for settings where MLE for $\boldsymbol{\theta}$ is easy, if only you knew \mathbf{z}

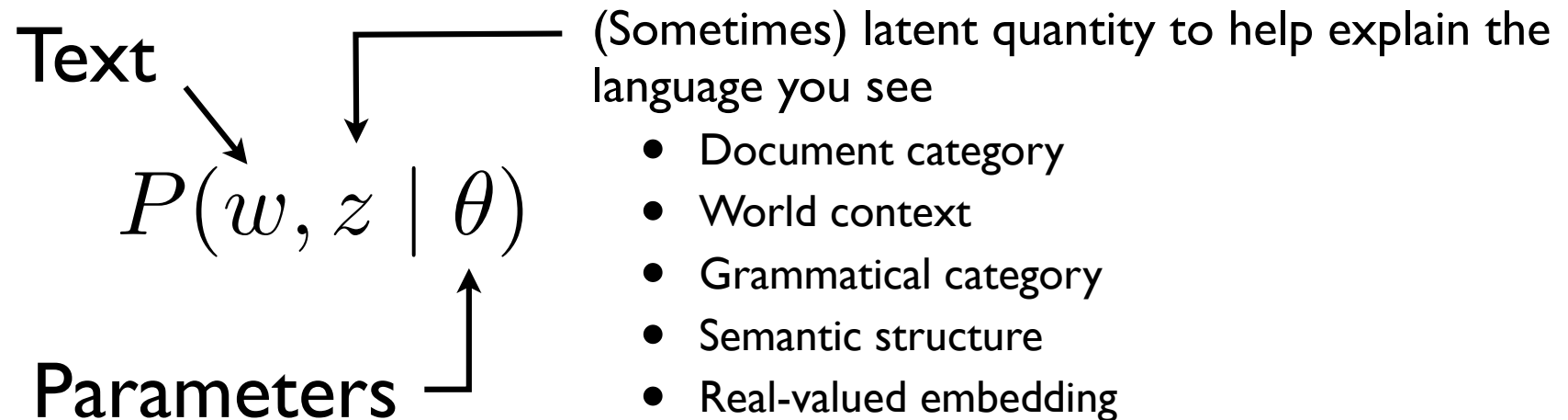
$$\max_{\boldsymbol{\theta}} \phi(\mathbf{x}|\boldsymbol{\theta}) \quad \text{is}$$

$$\max_{\boldsymbol{\theta}} \underbrace{p(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta})}_{\text{is}} \underbrace{p(\mathbf{z}|\boldsymbol{\theta})}_{\text{is}} \quad [\text{for GMM}]$$

Derivation of EM

- (new page)

Latent-variable generative models



Easy stuff

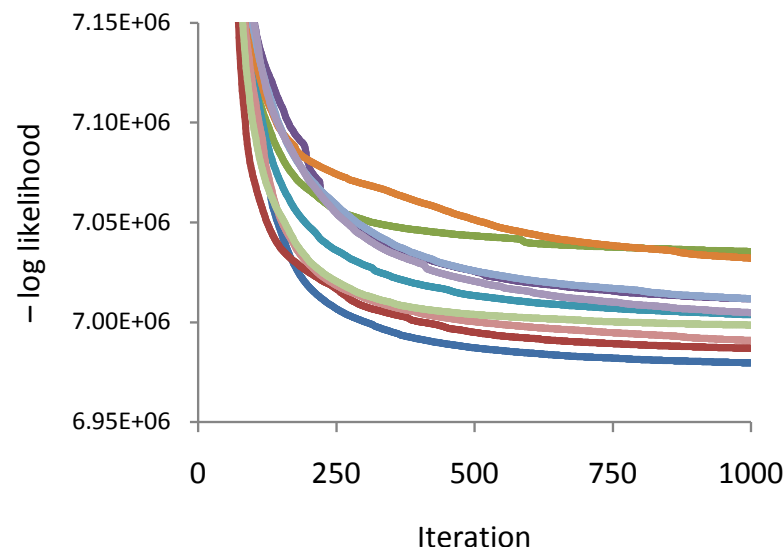
- Supervised learning: **$\text{argmax}_{\theta} \mathbf{P}(\mathbf{w}^{\text{train}}, \mathbf{z}^{\text{train}} | \theta)$**
- Prediction (via posterior inference): **$\mathbf{P}(\mathbf{z} | \mathbf{w}^{\text{input}}, \theta)$**

Unsupervised stuff with *marginal inference*

- Latent (unsupervised) learning: **$\text{argmax}_{\theta} \mathbf{P}(\mathbf{w}^{\text{train}} | \theta)$**
- Language modeling (via marginal inference): **$\mathbf{P}(\mathbf{w}^{\text{input}} | \theta)$**

EM performance

- Guaranteed to find a locally maximum likelihood solution. Guaranteed to converge.
- But can take a while
- Dependent on initialization



Johnson 2007, “Why doesn’t EM find good HMM POS-taggers?”

Figure 1: Variation in negative log likelihood with increasing iterations for 10 EM runs from different random starting points.

EM pros/cons

- Works best for a simple model with rapid E/M-step inference
- Requires probabilistic modeling assumptions
- Dependent on initialization
 - Many alternative methods (e.g. MCMC), but can have similar issues with local optima
- EM originally invented for Hidden Markov Models in speech recognition
- E-step infers structured posteriors
- General issue: Closed form M-steps only available for pretty simple models (Gaussians, count-based multinomials...)

EM versus direct gradients

- What if the M-step requires gradient ascent?
 - Running LBFGS or many iterations of GD inside the M-step can be slow
- Partial/incremental EM variants (Neal and Hinton, 1998): Why not just 1 gradient step? Gradient step on only a few examples?
- Or... consider the direct gradient. We can interpret it as an EM-like method itself.