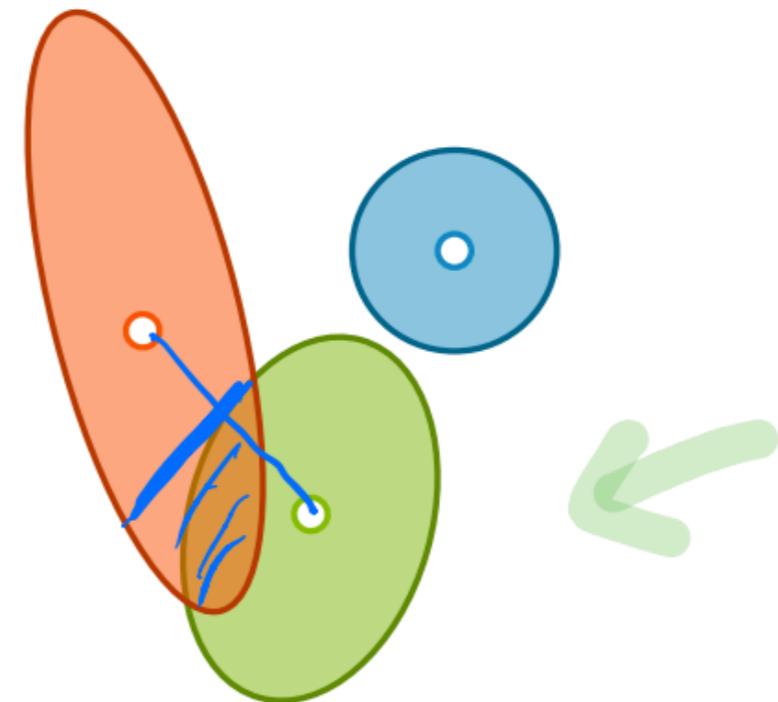
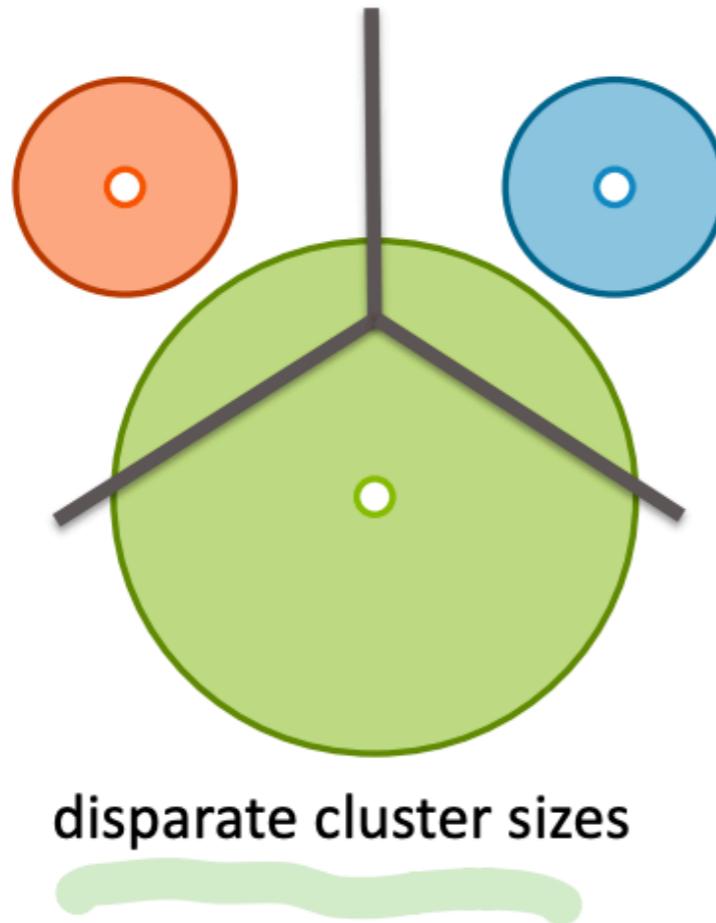


Expectation Maximization

Sewoong Oh

CSE446
University of Washington

- K-means algorithm fails, when



- one way to capture such clustering is by training the parameters of a **Gaussian Mixture Model (GMM)** that best captures the data

demo: <https://lukapopijac.github.io/gaussian-mixture-model/>

Gaussian Mixture Model.

input: $\{X_i\}_{i=1}^n$

, fix K : # of clusters

Parameters: $\pi = (\pi_1, \dots, \pi_K) \in R^K$: mixture weights

μ_j , $j \in \{1, \dots, K\} \in R^d$: mean

$C_j \in R^{d \times d}$,
: Covariance.

$d=1, K=2$.

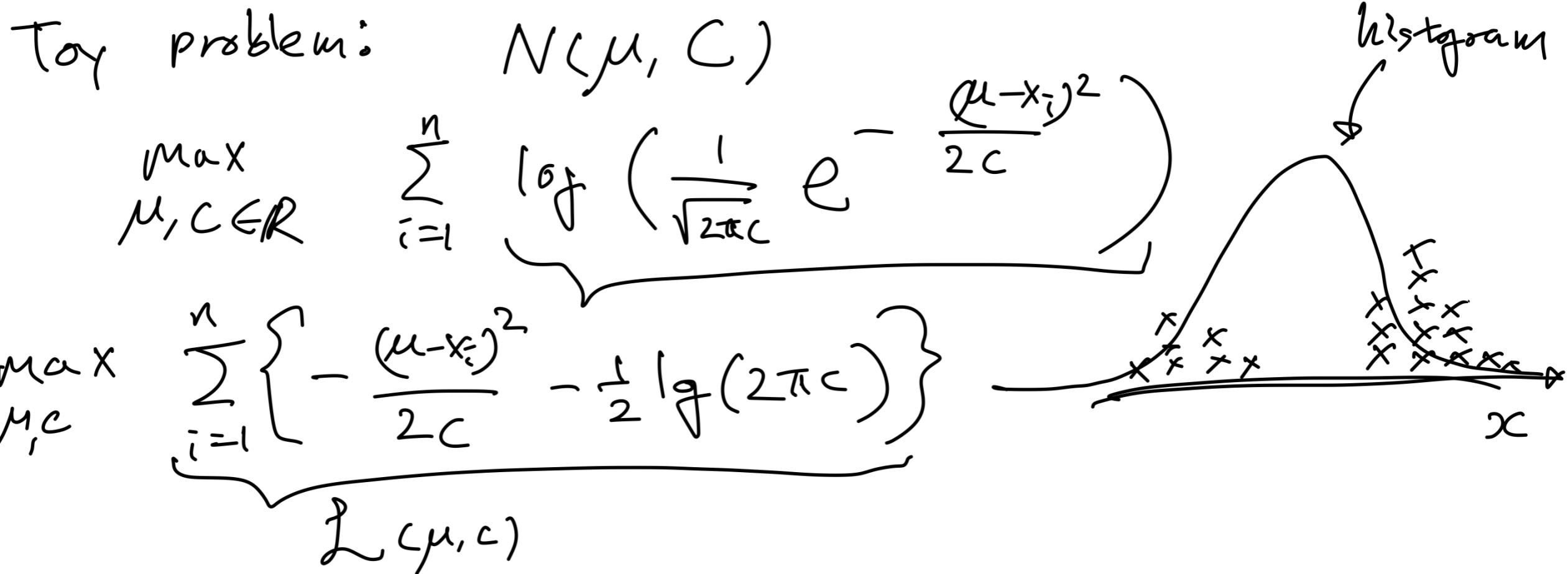
Parameters. $\underbrace{\pi_1, \pi_2, \mu_1, \mu_2, C_1, C_2 \in R}_1$

$$P(X_i | \text{Parameters}) = \pi_1 \frac{1}{\sqrt{2\pi C_1}} e^{-\frac{(X_i - \mu_1)^2}{2C_1}} + \pi_2 \frac{1}{\sqrt{2\pi C_2}} e^{-\frac{(X_i - \mu_2)^2}{2C_2}}$$

MLE:

Maximize
Parameters

$$\sum_{i=1}^n \log P(X_i | \text{Parameters})$$



$$\nabla_{\mu} L(\mu, \sigma) = \sum_{i=1}^n -\frac{2}{2\sigma^2} \cdot (\mu - x_i) = 0 \iff \frac{n\mu}{\sigma^2} = \sum_{i=1}^n x_i$$

$$\nabla_{\sigma} L(\mu, \sigma) = \sum_{i=1}^n \frac{(\mu - x_i)^2}{2\sigma^4} - \frac{n}{2\sigma^2} = 0 \iff \sigma^2 = \frac{1}{n} \sum_{i=1}^n (\mu - x_i)^2$$

MLE for GMM

Maximize

$\pi_1, \pi_2, \mu_1, \mu_2, c_1, c_2$

$$\sum_{i=1}^n \log \left(\pi_1 \frac{1}{\sqrt{2\pi c_1}} e^{-\frac{(x_i - \mu_1)^2}{2c_1}} + \pi_2 N(x_i | \mu_2, c_2) \right)$$

$$+ \sum_{k=1}^K \pi_k \frac{1}{\sqrt{2\pi c_k}} e^{-\frac{(x_i - \mu_k)^2}{2c_k}}$$

define $r_i \triangleq P(Z_i = 1 | X_i) = \frac{P(Z_i = 1, X_i)}{P(Z_i = 1, X_i) + P(Z_i = 2, X_i)}$ Bayes' Rule

$$= \frac{\pi_1 N(x_i | \mu_1, c_1)}{\pi_1 N(x_i | \mu_1, c_1) + \pi_2 N(x_i | \mu_2, c_2)}$$

$$1 - r_i = P(Z_i = 2 | X_i)$$

$$N_1 = \sum_{i=1}^n r_i, \quad N_2 = \sum_{i=1}^n (1 - r_i) \longrightarrow \pi_1 = \frac{N_1}{n}, \quad \pi_2 = \frac{N_2}{n}$$

$$\mu_2 = \frac{1}{N_2} \sum_i x_i (1 - r_i), \quad \mu_1 = \frac{1}{N_1} \sum_{i=1}^n x_i \cdot r_i$$

$$c_2 = \frac{1}{N_2} \sum_i (1 - r_i) (x_i - \mu_2)^2, \quad c_1 = \frac{1}{N_1} \sum_i r_i (x_i - \mu_1)^2$$

Gaussian Mixture Model

- input: data $\{x_i\}_{i=1}^n$ in \mathbb{R}^d
- parameters of a **Gaussian Mixture Model**
 - mixing weights:
 - $\pi_j = \mathbf{P}(\text{cluster membership} = j)$ for $j \in \{1, \dots, K\}$
 - means:
 - $\mu_j \in \mathbb{R}^d$ for $j \in \{1, \dots, K\}$
 - covariance matrices:
 - $\mathbf{C}_j \in \mathbb{R}^{d \times d}$ for $j \in \{1, \dots, K\}$
- we suppose that the given data has been generated from a GMM, and try to find the best GMM parameters (this naturally will define clustering of the training data)
- under the GMM, the i -th sample is drawn as follows
 - first sample a cluster $z_i \in \{1, \dots, K\}$, from $\pi = [\pi_1, \dots, \pi_K]$
 - conditioned on this cluster, x_i is sampled from
$$x_i \sim N(\mu_{z_i}, \mathbf{C}_{z_i})$$

Maximum likelihood estimation (MLE)

- we can find the best GMM, by MLE
- for simplicity, suppose $d = 1$ and $K = 2$
- Model parameters are $\pi_1, \pi_2, \mu_1, \mu_2, \mathbf{C}_1, \mathbf{C}_2 \in \mathbb{R}$
- the probability of observing a sample x_i can be written as

$$\mathbf{P}(x_i | \pi_1, \pi_2, \mu_1, \mu_2, \mathbf{C}_1, \mathbf{C}_2) = \underbrace{\pi_1 \frac{1}{\sqrt{2\pi\mathbf{C}_1}} e^{-\frac{(x_i - \mu_1)^2}{2\mathbf{C}_1}}}_{\triangleq N(x_i | \mu_1, \mathbf{C}_1)} + \underbrace{\pi_2 \frac{1}{\sqrt{2\pi\mathbf{C}_2}} e^{-\frac{(x_i - \mu_2)^2}{2\mathbf{C}_2}}}_{\triangleq N(x_i | \mu_2, \mathbf{C}_2)}$$

- MLE tries to find

$$\arg \max_{\pi_1, \pi_2, \mu_1, \mu_2, \mathbf{C}_1, \mathbf{C}_2} \sum_{i=1}^n \log \mathbf{P}(x_i | \pi_1, \pi_2, \mu_1, \mu_2, \mathbf{C}_1, \mathbf{C}_2)$$

- however, unlike least squared or logistic regression, this is not a concave function of the parameters (thus hard to find the optimal solution)
- in general, MLE of a mixture model is not convex/concave optimization

exercise: fitting a single Gaussian model

- given $\{x_i\}_{i=1}^n \in \mathbb{R}$, fit the best Gaussian model with mean $\mu \in \mathbb{R}$ and variance $C \in \mathbb{R}$
- using MLE we want to solve

$$\text{maximize}_{\mu, C} \mathcal{L}(\mu, C) = \sum_{i=1}^n \underbrace{\left(-\frac{(x_i - \mu)^2}{2C} - \log(\sqrt{2\pi C}) \right)}_{\log N(x_i | \mu, C)}$$

- we compute gradient and set it to zero:

$$\bullet \quad \nabla_{\mu} \mathcal{L}(\mu, C) = \frac{1}{C} \sum_{i=1}^n (\mu - x_i)$$

which is zero for $\mu = \boxed{\frac{1}{n} \sum_{i=1}^n x_i}$

(which makes sense as it is the empirical mean)

$$\bullet \quad \nabla_C \mathcal{L}(\mu, C) = \frac{\sum_{i=1}^n (x_i - \mu)^2}{2C^2} - \frac{n}{2C}$$

which is zero for $C = \boxed{\frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2}$

(which makes sense as it is the empirical variance)

MLE for GMM

- we want to fit a model by solving

$$\underset{\pi_1, \pi_2, \mu_1, \mu_2, \mathbf{C}_1, \mathbf{C}_2}{\text{maximize}} \sum_{i=1}^n \log \left(\underbrace{\pi_1 \frac{1}{\sqrt{2\pi\mathbf{C}_1}} e^{-\frac{(x_i - \mu_1)^2}{2\mathbf{C}_1}}}_{\triangleq N(x_i | \mu_1, \mathbf{C}_1)} + \underbrace{\pi_2 \frac{1}{\sqrt{2\pi\mathbf{C}_2}} e^{-\frac{(x_i - \mu_2)^2}{2\mathbf{C}_2}}}_{\triangleq N(x_i | \mu_2, \mathbf{C}_2)} \right)$$

- define $r_i = \mathbf{P}(z_i = 1 | x_i) = \frac{\mathbf{P}(z_i = 1, x_i)}{\mathbf{P}(z_i = 1, x_i) + \mathbf{P}(z_i = 2, x_i)}$

$$= \frac{\pi_1 N(x_i | \mu_1, \mathbf{C}_1)}{\pi_1 N(x_i | \mu_1, \mathbf{C}_1) + \pi_2 N(x_i | \mu_2, \mathbf{C}_2)}$$

- setting the gradient to zero, we get

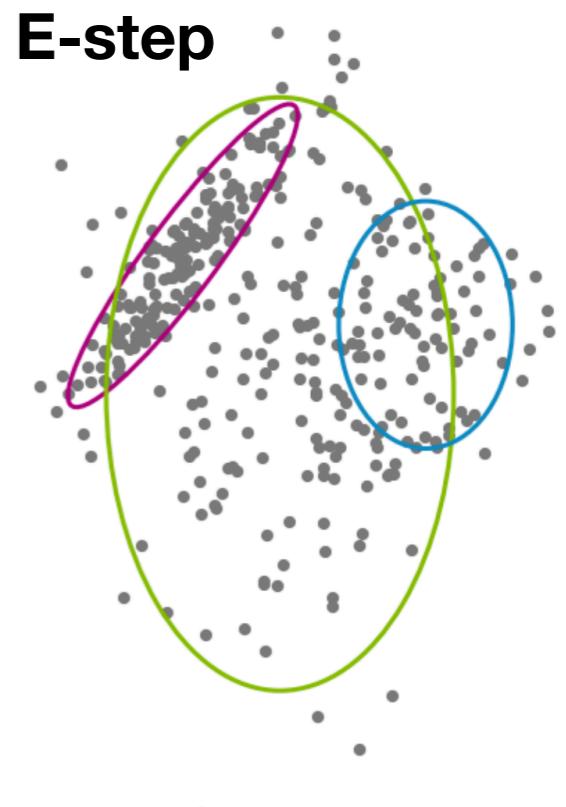
- $\pi_1 = \frac{N_1}{n}$ where $N_1 = \sum_{i=1}^n r_i$, and $\pi_2 = \frac{N_2}{n}$ where $N_2 = \sum_{i=1}^n (1 - r_i)$
- $\mu_1 = \frac{1}{N_1} \sum_{i=1}^n r_i x_i$ and $\mu_2 = \frac{1}{N_2} \sum_{i=1}^n (1 - r_i) x_i$
- $\mathbf{C}_1 = \frac{1}{N_1} \sum_{i=1}^n r_i (x_i - \mu_1)^2$ and $\mathbf{C}_2 = \frac{1}{N_2} \sum_{i=1}^n (1 - r_i) (x_i - \mu_2)^2$

- both LHS and RHS depend on the parameters, and no closed form solution exists
- **note that if we know r_i 's it is trivial to compute parameters, and vice versa**

Expectation Maximization (EM) algorithm

- EM is a popular method to solve MLE for mixture models
- input: training data $\{x_i\}_{i=1}^n$
- output: $\pi_1, \pi_2, \mu_1, \mu_2, \mathbf{C}_1, \mathbf{C}_2 \in \mathbb{R}$
- initialization: randomly initialize the parameters
- repeat
 - **E-step (Expectation):** parameters \rightarrow soft membership

$$r_i = \frac{\pi_1 N(x_i | \mu_1, \mathbf{C}_1)}{\pi_1 N(x_i | \mu_1, \mathbf{C}_1) + \pi_2 N(x_i | \mu_2, \mathbf{C}_2)}$$

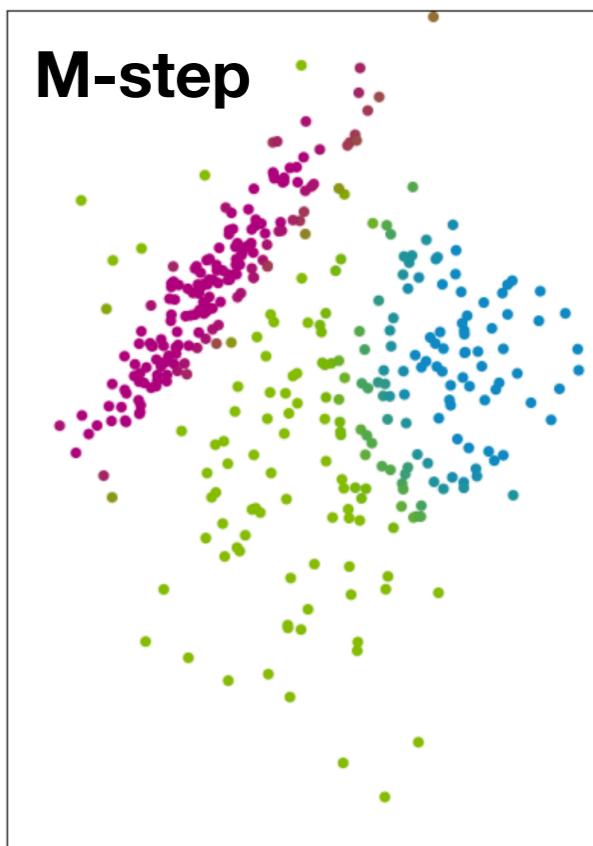


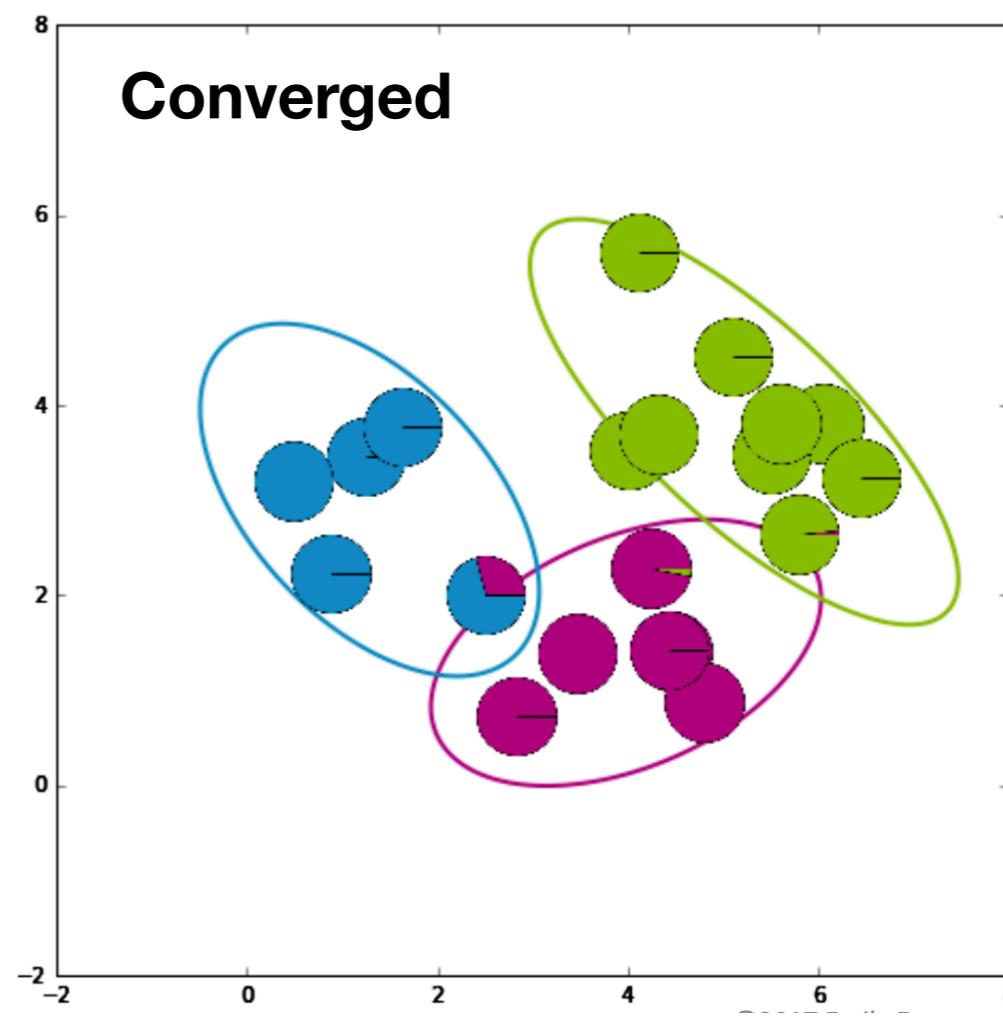
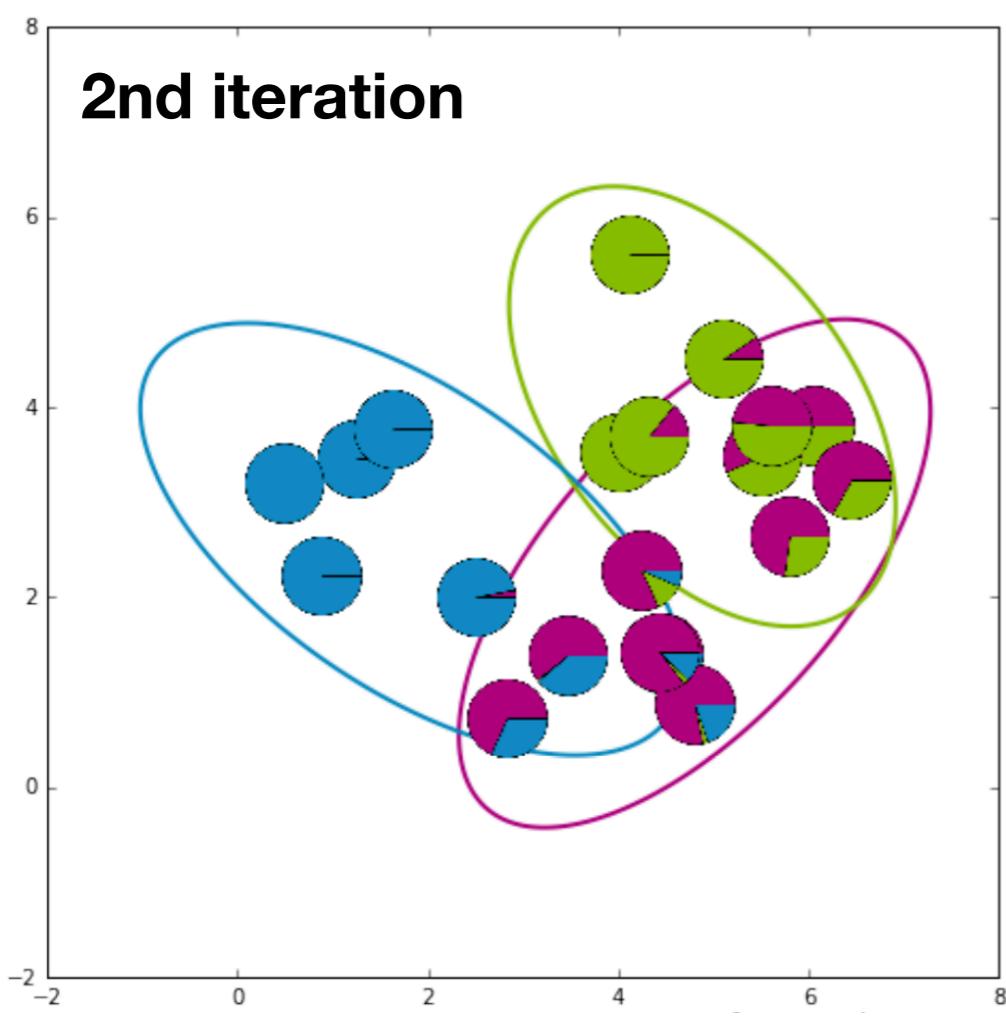
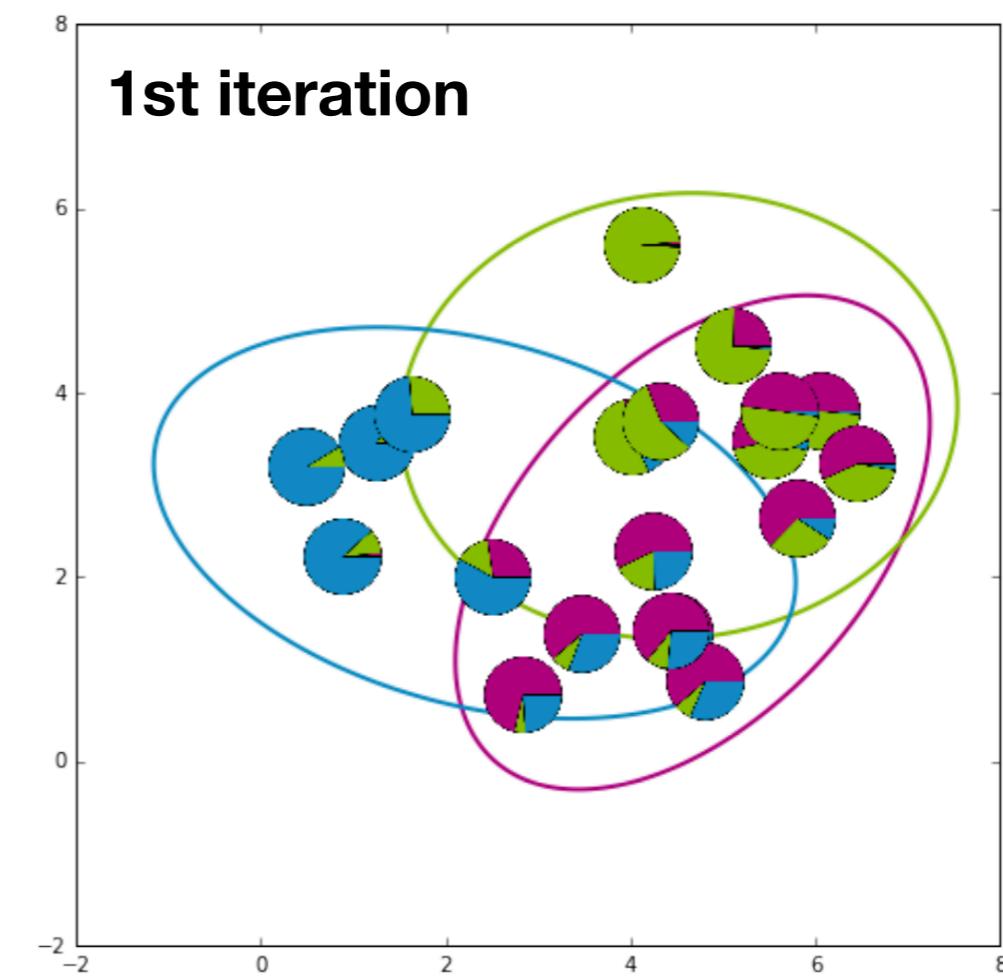
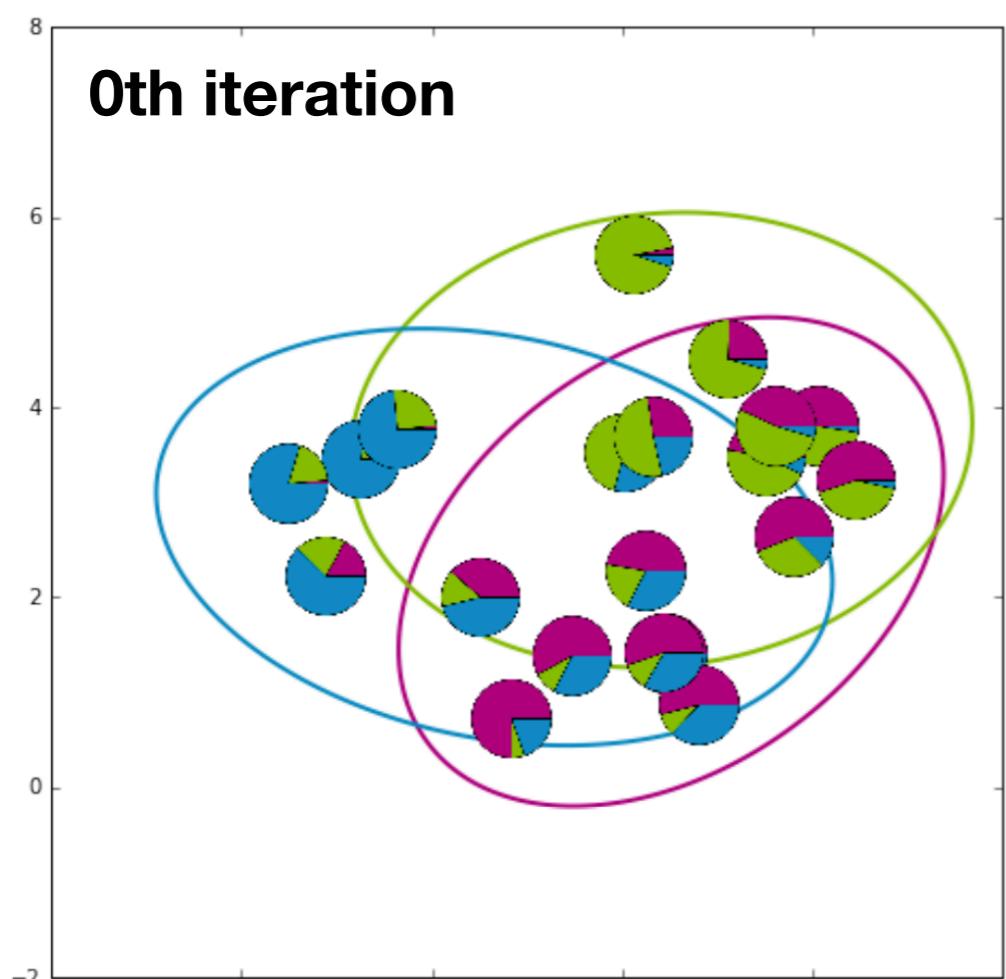
- **M-step (Maximization):** soft membership \rightarrow parameters

$$\pi_1 = \frac{N_1}{n} \text{ where } N_1 = \sum_{i=1}^n r_i, \text{ and } \pi_2 = \frac{N_2}{n} \text{ where } N_2 = \sum_{i=1}^n (1 - r_i)$$

$$\mu_1 = \frac{1}{N_1} \sum_{i=1}^n r_i x_i \quad \text{and} \quad \mu_2 = \frac{1}{N_2} \sum_{i=1}^n (1 - r_i) x_i$$

$$\mathbf{C}_1 = \frac{1}{N_1} \sum_{i=1}^n r_i (x_i - \mu_1)^2 \text{ and } \mathbf{C}_2 = \frac{1}{N_2} \sum_{i=1}^n (1 - r_i) (x_i - \mu_2)^2$$





For general number of clusters K and dimension d

- we can derive EM for general case, in an analogous way
- Initialize parameters: $\pi_1, \dots, \pi_K, \mu_1, \dots, \mu_K, \mathbf{C}_1, \dots, \mathbf{C}_K$

- **E-step:**

- For $k=1, \dots, K$

$$r_{i,k} = \frac{\pi_k N(x_i | \mu_k, \mathbf{C}_k)}{\sum_{j=1}^K \pi_j N(x_i | \mu_j, \mathbf{C}_j)}$$

- **M-step:**

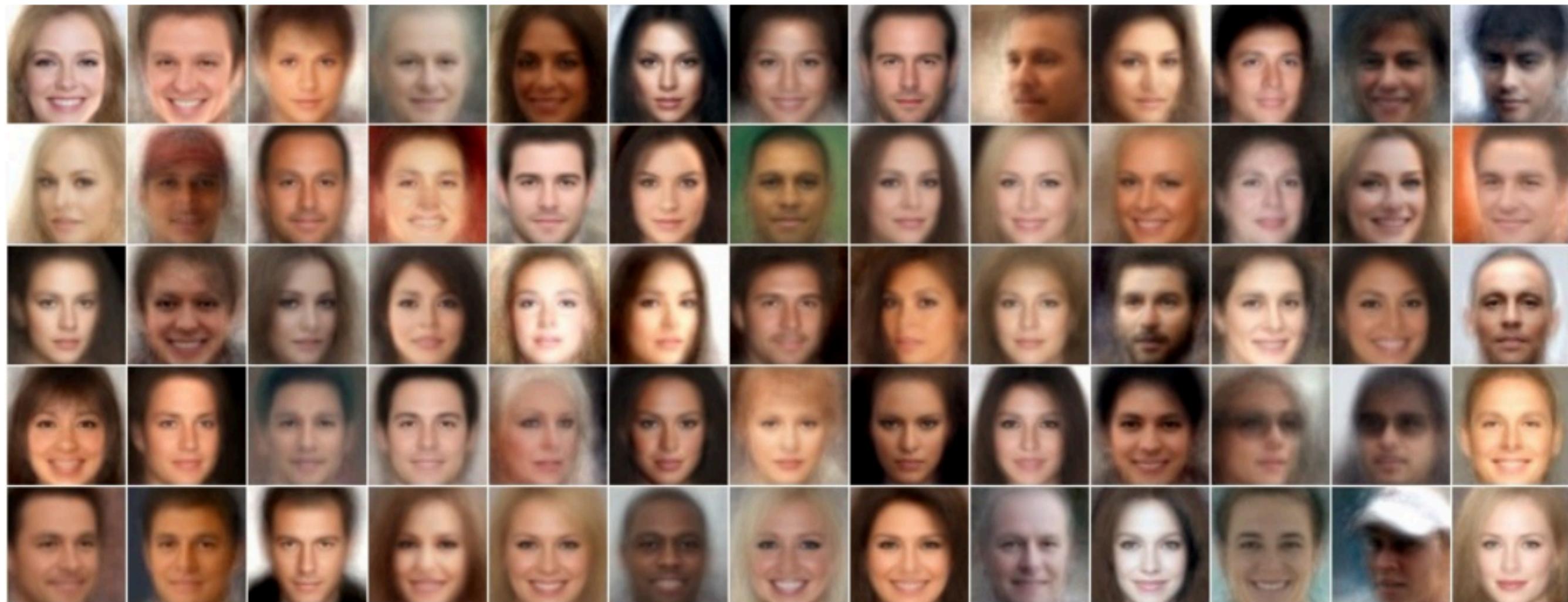
- For $k=1, \dots, K$

$$\pi_k = \frac{N_k}{n} \quad \text{where} \quad N_k = \frac{\sum_{i=1}^n r_{i,k}}{n}$$

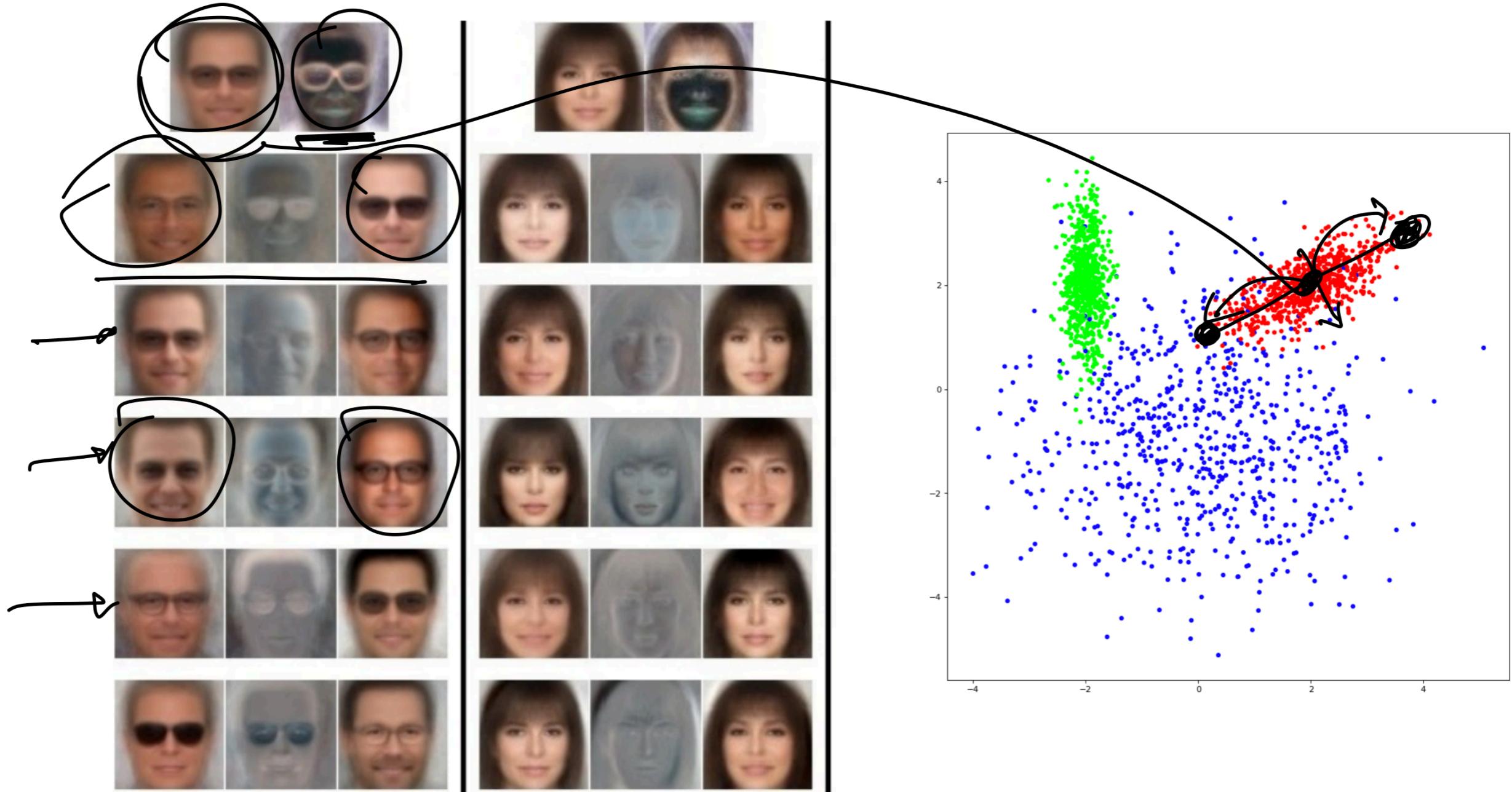
$$\mu_k = \frac{1}{N_k} \sum_{i=1}^n r_{i,k} x_i \quad \text{and} \quad \mathbf{C}_k = \frac{1}{N_k} \sum_{i=1}^n r_{i,k} (x_i - \mu_k)(x_i - \mu_k)^T$$

- once GMM is learned, clustering is straight forward: cluster according to the $r_{i,k}$'s

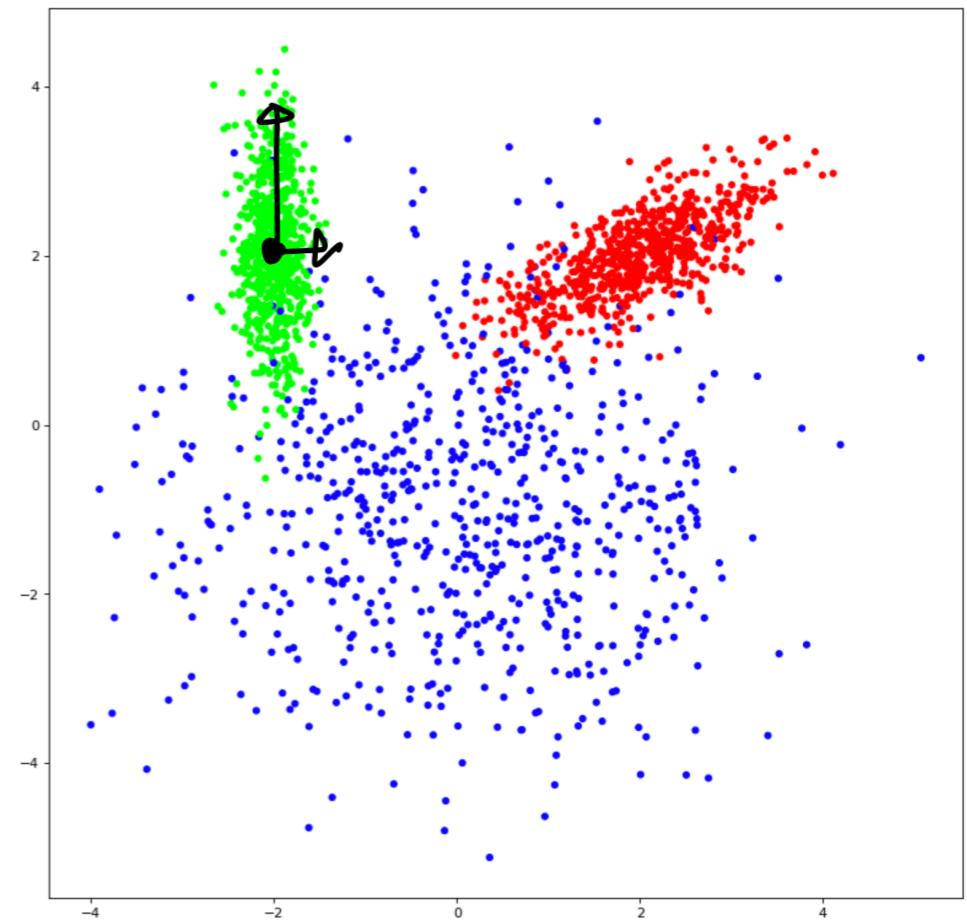
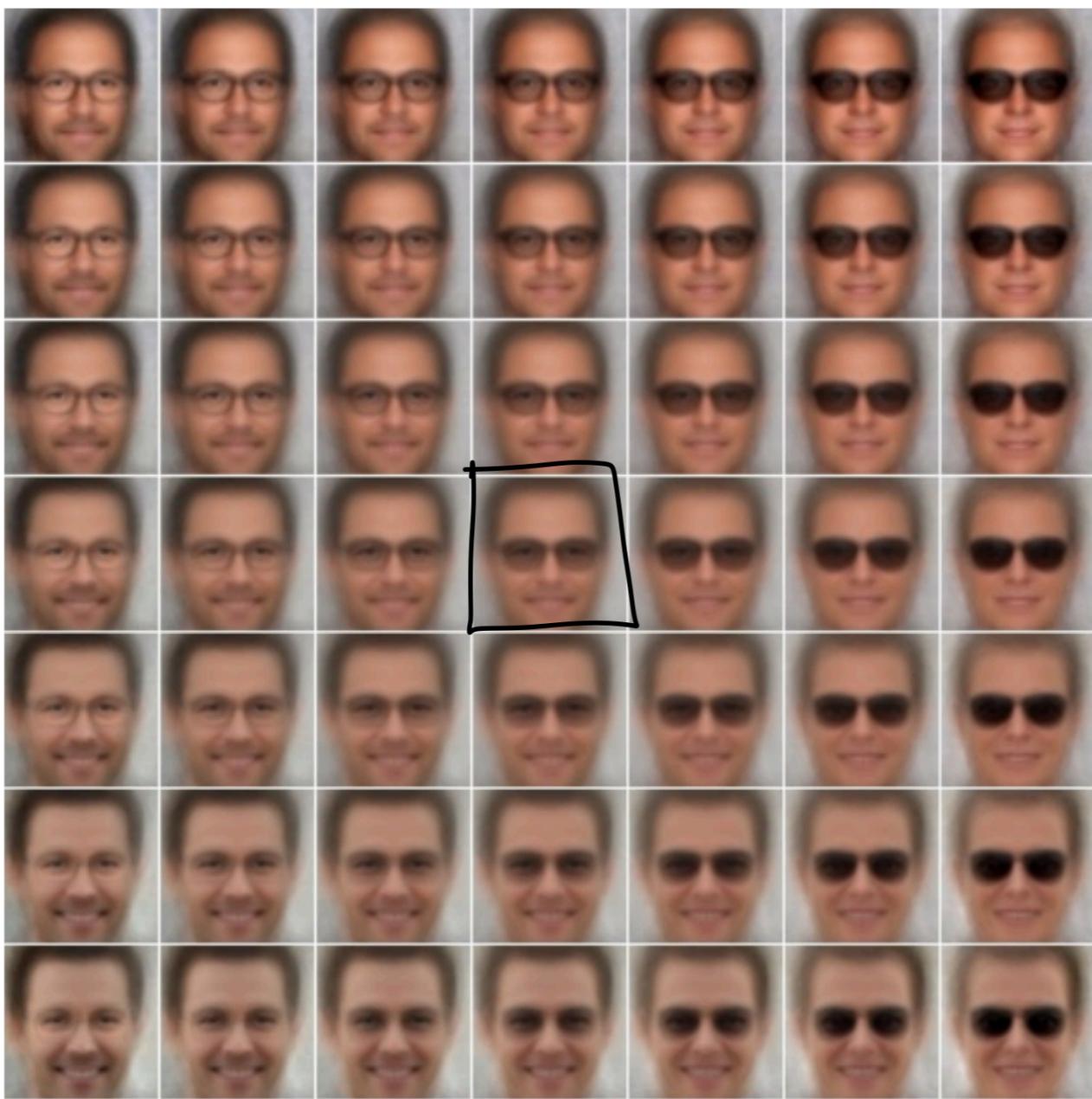
GMM for real data



- these are generated samples, from GMM trained on CelebA dataset
- image: $64 \times 64 \times 3 = 288$ dimension
- covariance: restricted to rank-10 matrices
- mixture: $K=1,000$



- top: center of a cluster μ_k and the diagonal entries of the covariance matrix \mathbf{C}_k
- note that we have trained 10-dimensional covariance matrix $\mathbf{C}_k = \mathbf{A}\mathbf{A}^T$, with $\mathbf{A} \in \mathbb{R}^{288 \times 10}$, and let $\mathbf{A}^{(j)}$ be the j-th column
- bottom: each row corresponds to different j , and we show $\mu_k + \mathbf{A}^{(j)}, 0.5 + \mathbf{A}^{(j)}, \mu_k - \mathbf{A}^{(j)}$



- middle: μ_k
- Each row: middle + $c \times A^{(1)}$
- Each column: middle + $c \times A^{(2)}$

Mixture model for documents

- Input: n documents $\{x_i\}_{i=1}^n$
- Each document is a sequence of words of length T
 $x_i = (w_1, w_2, \dots, w_T)$
- Bag-of-words model:
 - parameters:
 - mixing weights: $\pi_k = \mathbf{P}(\text{topic} = k)$ for $k \in \{1, \dots, K\}$
 - word probability: $b_{wk} = \mathbf{P}(\text{word} = w \mid \text{topic} = k)$
 - the generative model
 - first sample topic from $\pi = (\pi_1, \dots, \pi_K)$
 - next sample T words i.i.d. from $b_k = (b_{w_1k}, \dots, b_{w_{200,000}k})$
 - to make the problem tractable, this completely ignores the order of the words in the document (but still very successful in document clustering)

$$\mathbf{P}(\text{topic } z_i = k, x_i = (w_1, \dots, w_T)) = \pi_k b_{w_1k} \cdots b_{w_Tk}$$

Topic modeling

- to fit a topic model, we solve the following

$$\underset{b \in \mathbb{R}^{K \times T}, \pi \in \mathbb{R}^K}{\text{maximize}} \sum_{i=1}^n \log \mathbf{P}(x_i | b, \pi)$$

- we can apply EM algorithm
- initialize b, π
- **E-step:** parameters \rightarrow soft assignments

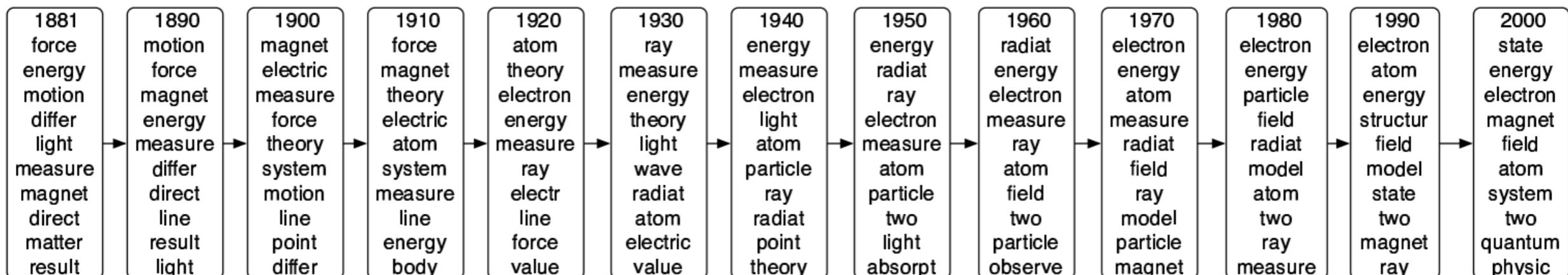
$$r_{ik} = \mathbf{P}(\text{topic } z_i = k | x_i) = \frac{\pi_k b_{w_1 k} \cdots b_{w_T k}}{\sum_{k'=1}^K \pi_{k'} b_{w_1 k'} \cdots b_{w_T k'}}$$

- **M-step:** soft assignments \rightarrow parameters

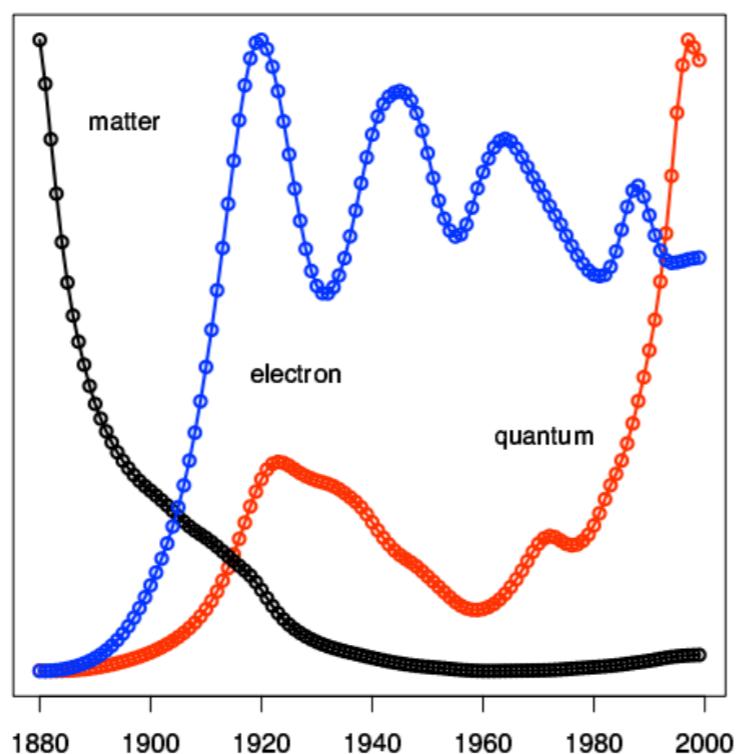
$$\pi_k = \frac{N_k}{n} \quad \text{where} \quad N_k = \sum_{i=1}^n r_{ik}$$

$$b_{w k} = \frac{1}{N_k} \sum_{i=1}^n r_{ik} \frac{\text{Count}(w \text{ in } x_i)}{T}$$

Dynamic topic modeling (over time)

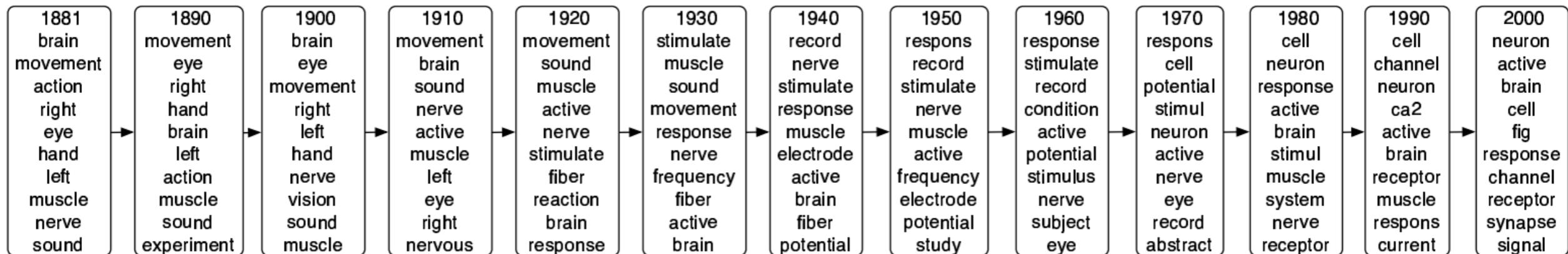


"Atomic Physics"

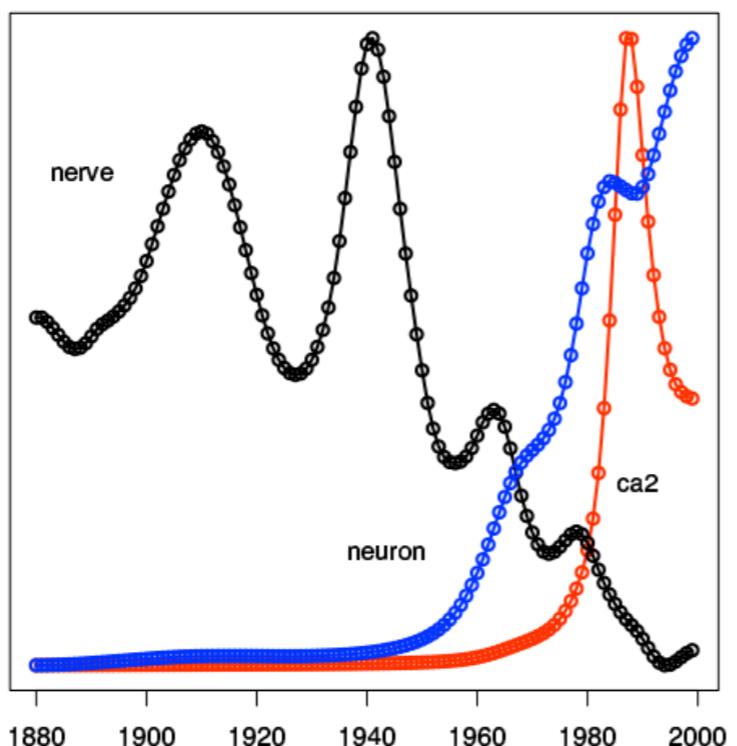


- 1881 On Matter as a form of Energy
- 1892 Non-Euclidean Geometry
- 1900 On Kathode Rays and Some Related Phenomena
- 1917 ``Keep Your Eye on the Ball''
- 1920 The Arrangement of Atoms in Some Common Metals
- 1933 Studies in Nuclear Physics
- 1943 Aristotle, Newton, Einstein. II
- 1950 Instrumentation for Radioactivity
- 1965 Lasers
- 1975 Particle Physics: Evidence for Magnetic Monopole Obtained
- 1985 Fermilab Tests its Antiproton Factory
- 1999 Quantum Computing with Electrons Floating on Liquid Helium

Dynamic topic modeling (over time)



"Neuroscience"



- 1887 Mental Science
1900 Hemianopsia in Migraine
1912 A Defence of the ``New Phrenology''
1921 The Synchronal Flashing of Fireflies
1932 Myoesthesia and Imageless Thought
1943 Acetylcholine and the Physiology of the Nervous System
1952 Brain Waves and Unit Discharge in Cerebral Cortex
1963 Errorless Discrimination Learning in the Pigeon
1974 Temporal Summation of Light by a Vertebrate Visual Receptor
1983 Hysteresis in the Force-Calcium Relation in Muscle
1993 GABA-Activated Chloride Channels in Secretory Nerve Endings