

# Quantum Expectation Maximization and other quantum algorithms for learning representations

November 15, 2020



INSTITUT  
DE RECHERCHE  
EN INFORMATIQUE  
FONDAMENTALE

Atos

Cambridge University

Alessandro Luongo

## In short

We propose a quantum algorithm for  
**Expectation-Maximization**

that fits a

**Gaussian Mixture Model**

(using quantum linear algebra, QRAM, distance estimation, etc..)  
for a matrix  $X \in \mathbb{R}^{n \times d}$  in time:

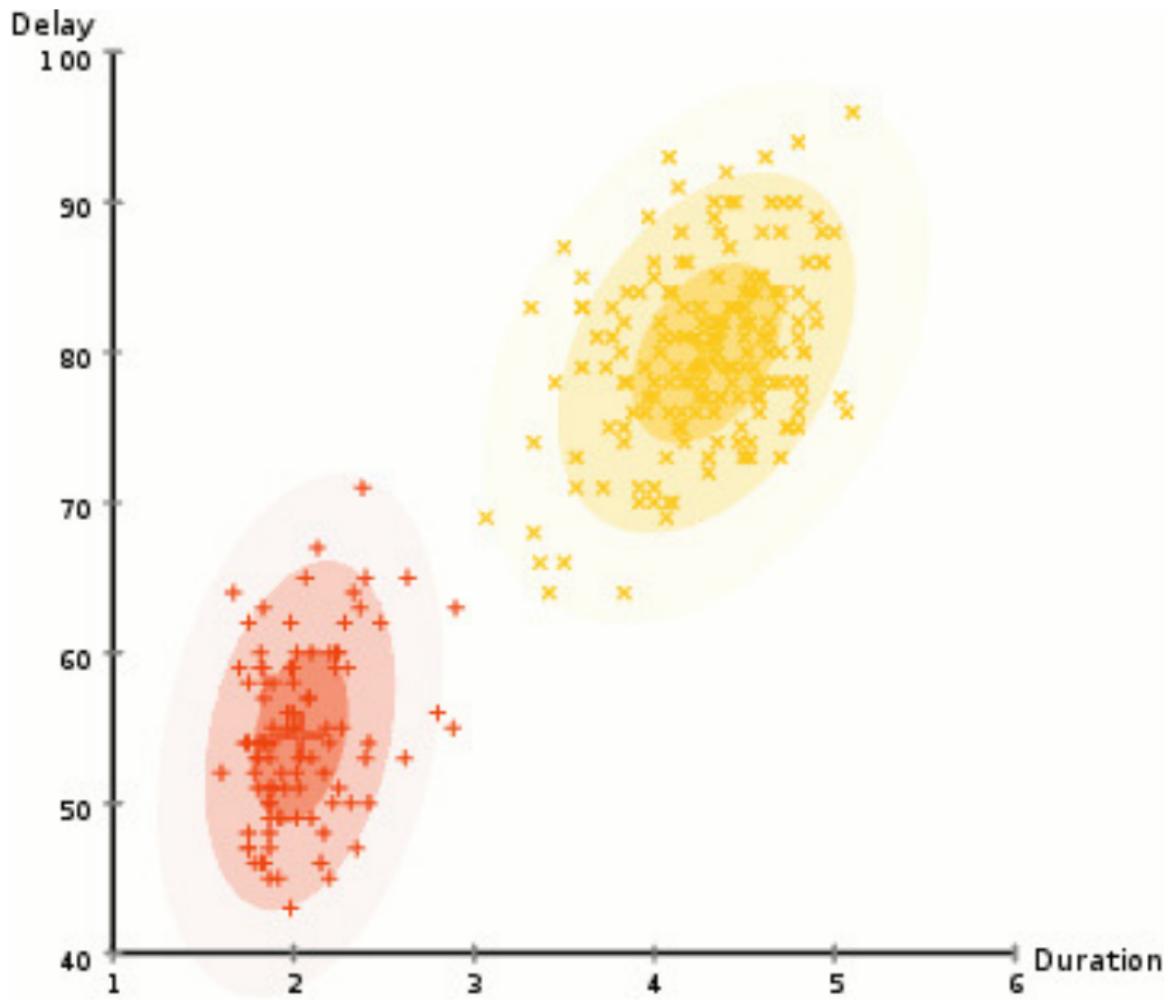
$$O\left(\frac{d^2 k^{4.5} \eta^3 \kappa^2(V) \kappa^2(\Sigma) \mu(\Sigma) \log n}{\delta_\mu^3}\right).$$

Particular case: **q-means**, a quantum algorithm for k-means:

$$\tilde{O}\left(k^2 d \frac{\eta^{2.5}}{\delta^3} + k^{2.5} \frac{\eta^2}{\delta^3}\right)$$

## Papers?

- ▶ **q-means**: <https://arxiv.org/abs/1812.03584>  
Iordanis Kerenidis, Jonas Landman, Anupam Prakash  
(Accepted at NeurIPS)
- ▶ **QEM**: <https://arxiv.org/pdf/1908.06657.pdf>  
Iordanis Kerenidis, Anupam Prakash
- ▶ ( Also **QGMM** <https://arxiv.org/abs/1908.06655> Hideyuki Miyahara, Kazuyuki Aihara, Wolfgang Lechner )



## Maximum Likelihood Estimation

Unsupervised data:  $V \in \mathbb{R}^{n \times d}$ .

Generative models: learn the **best** (i.e. with Maximum Likelihood )  $p(v_i)$  and **sample** from it.

$$L(V; \gamma) = \prod_{i=1}^n p(v_i; \gamma)$$

$$\gamma_{\textcolor{red}{MLE}}^* = \arg \max_{\gamma} L(V; \gamma)$$

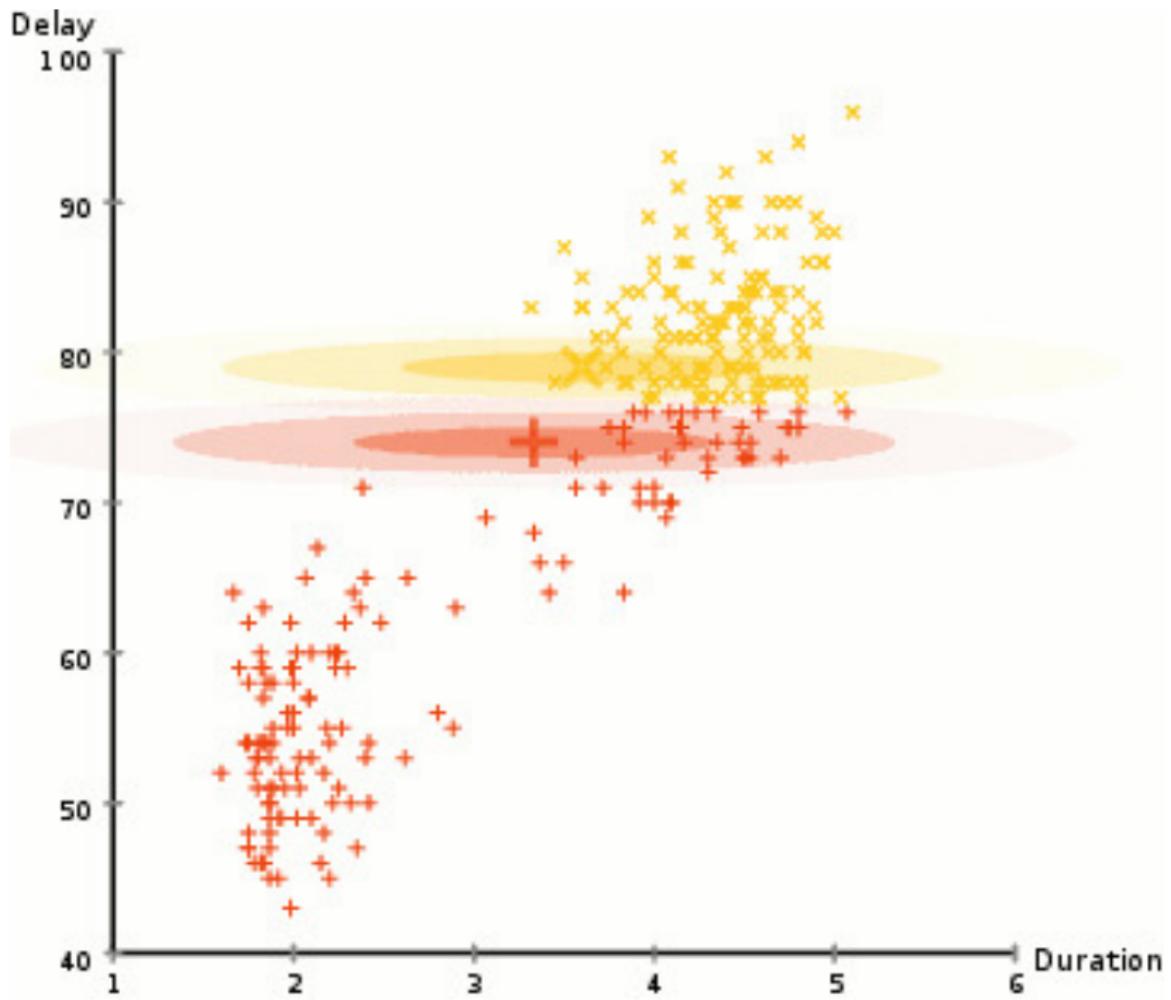
## Maximum A Posteriori estimate

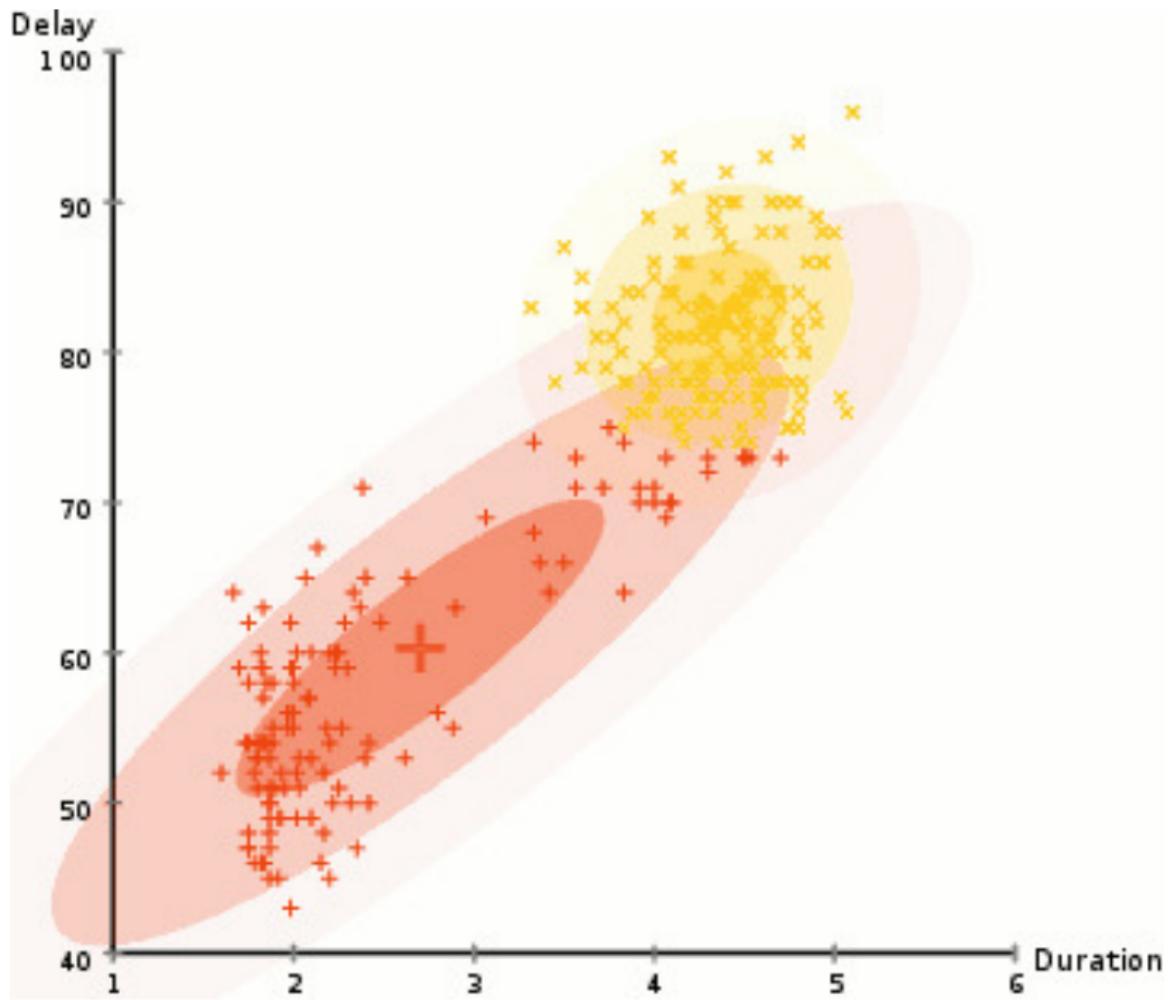
MAP uses a prior  $p(\gamma)$  on the model:

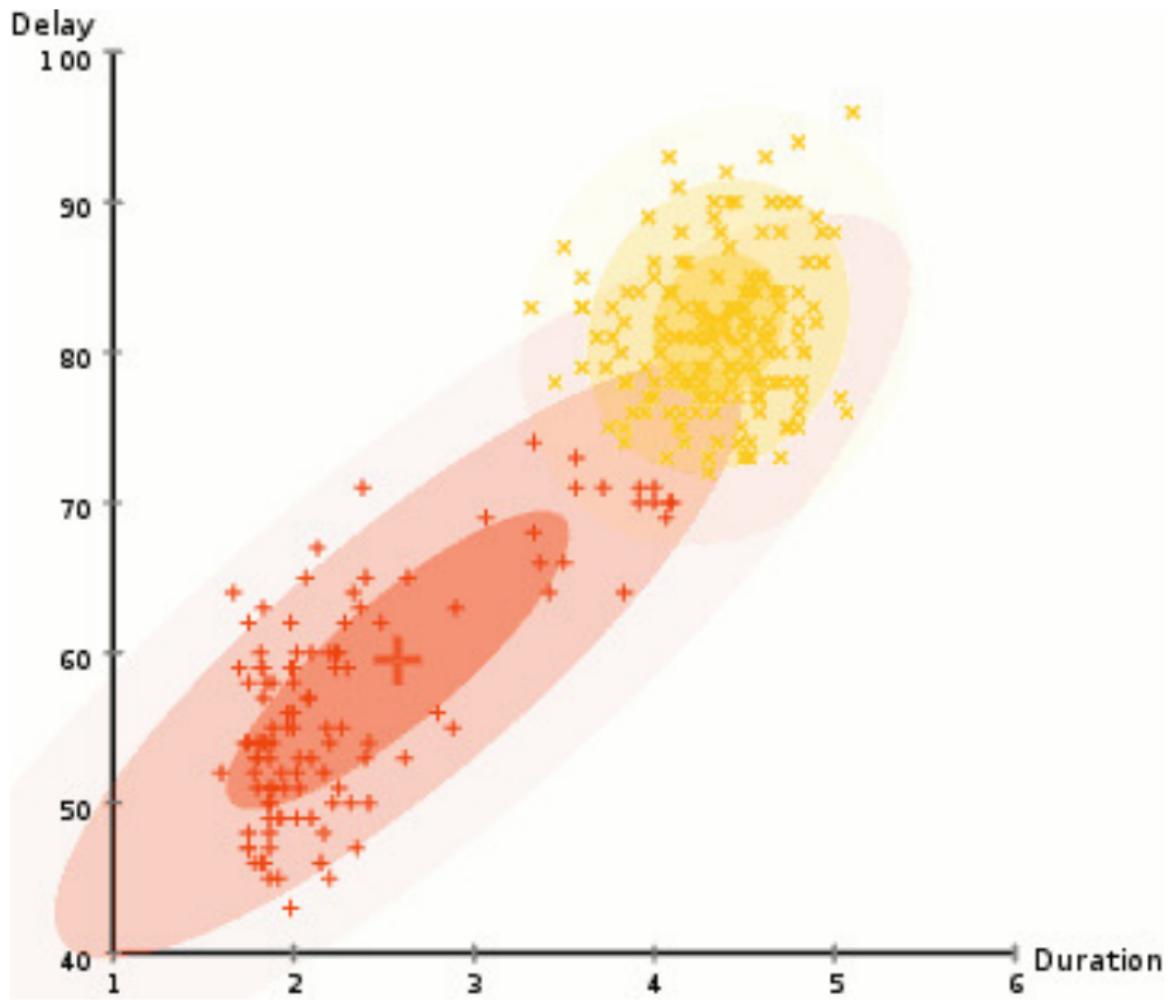
$$L(\gamma; V) = \prod_{i=1}^n p(v_i|\gamma)p(\gamma)$$

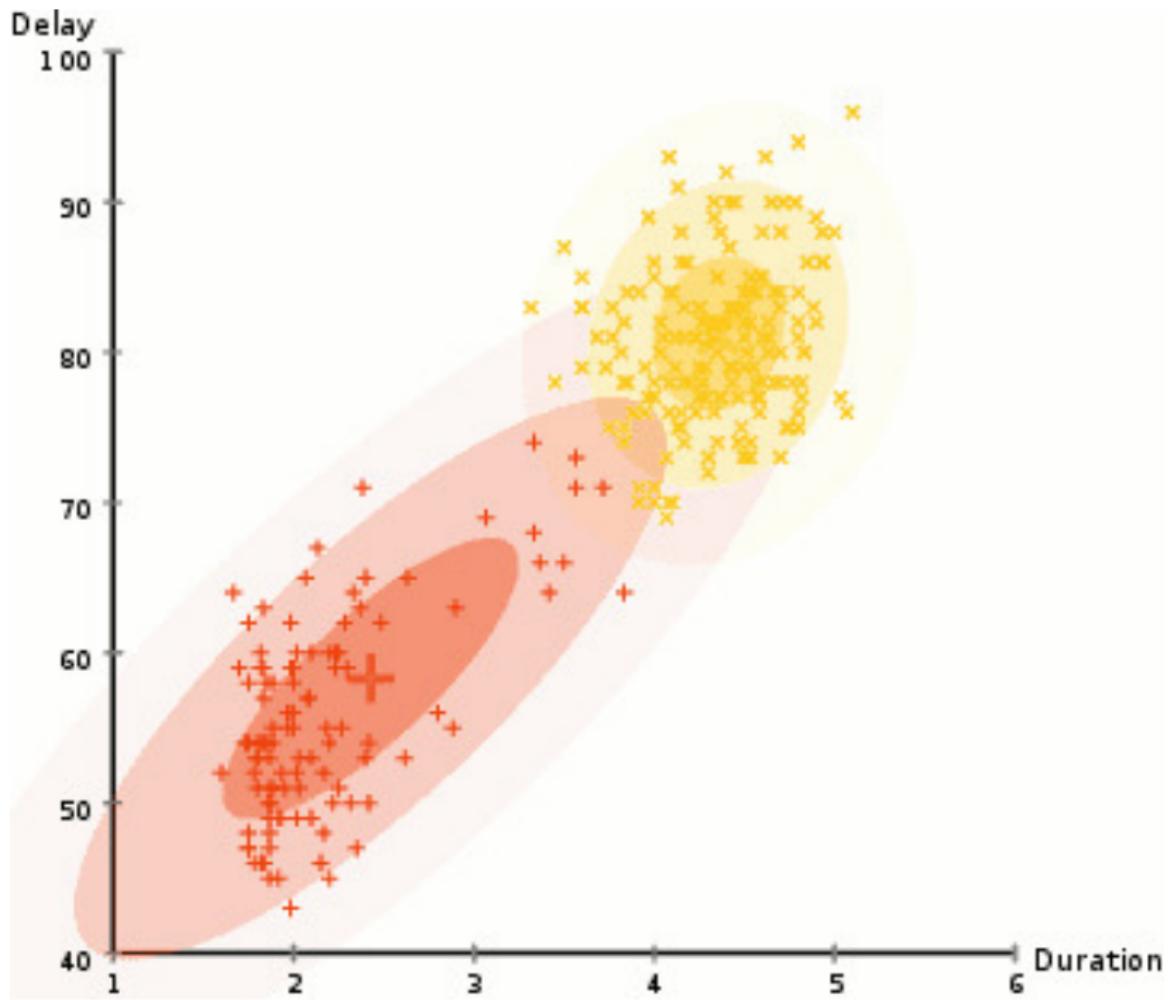
$$\gamma_{MAP} := \arg \max_{\gamma} L(\gamma; V)$$

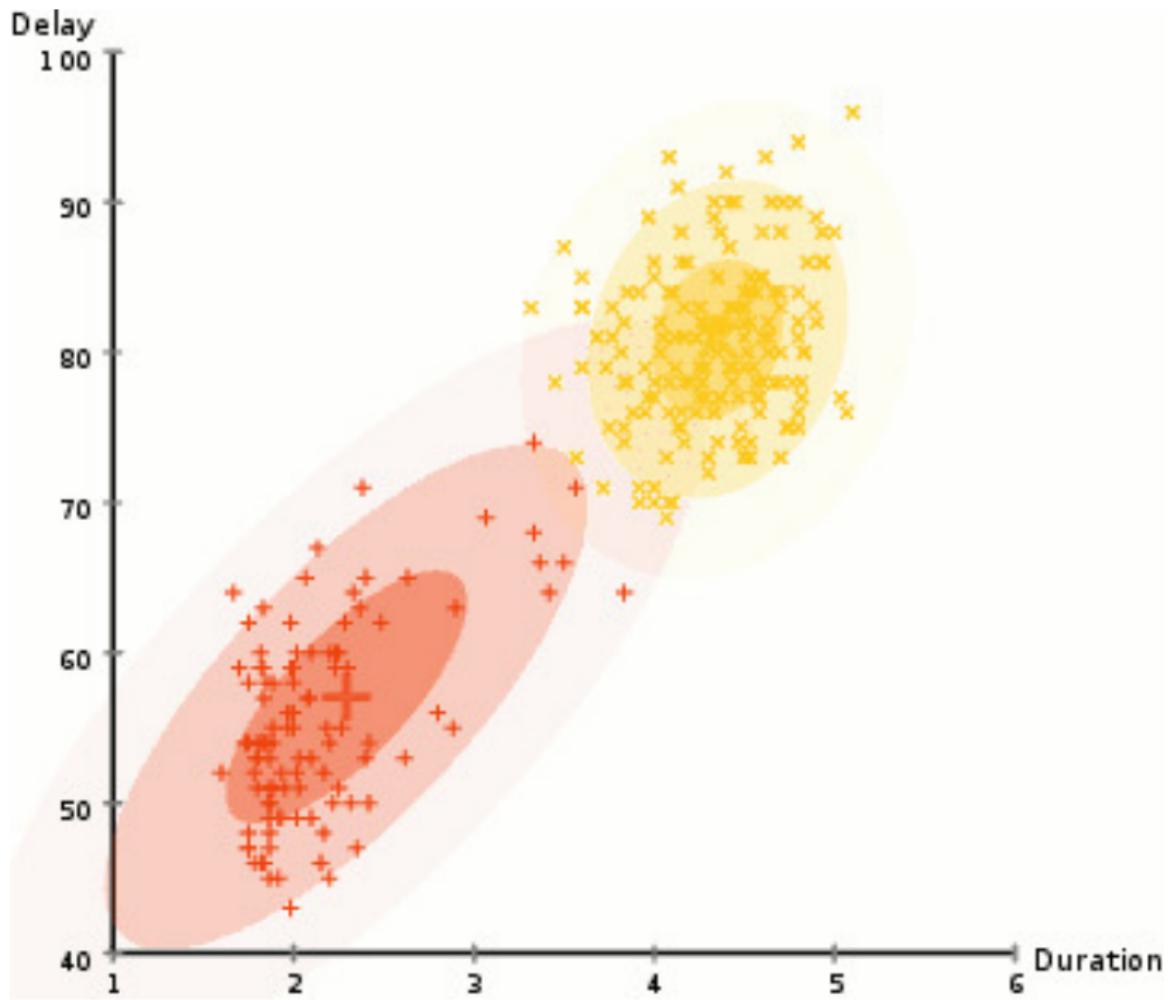
**Intuition:** Putting a prior (from domain experts) often avoid overfitting and wrong local minima (decrease number of iterations);

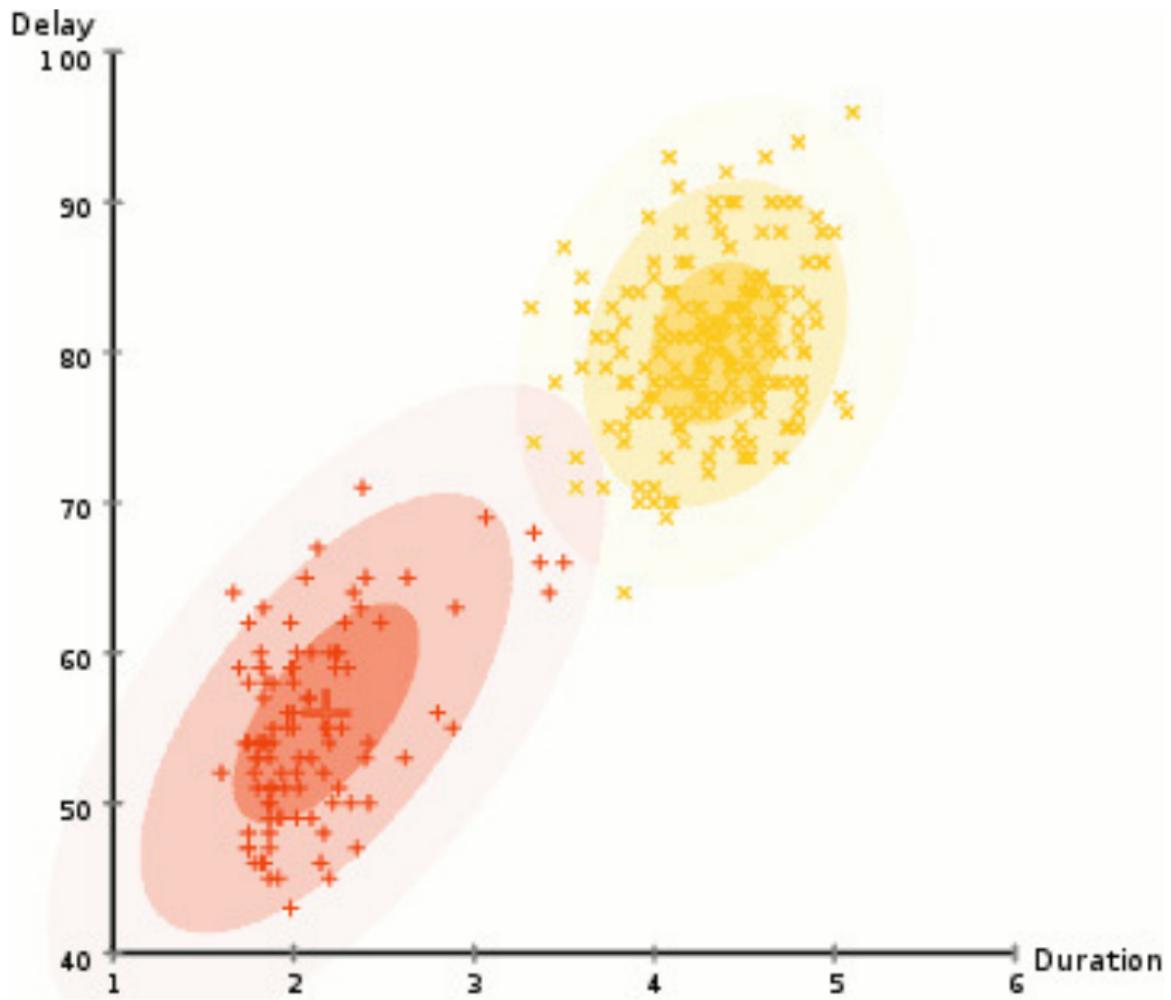


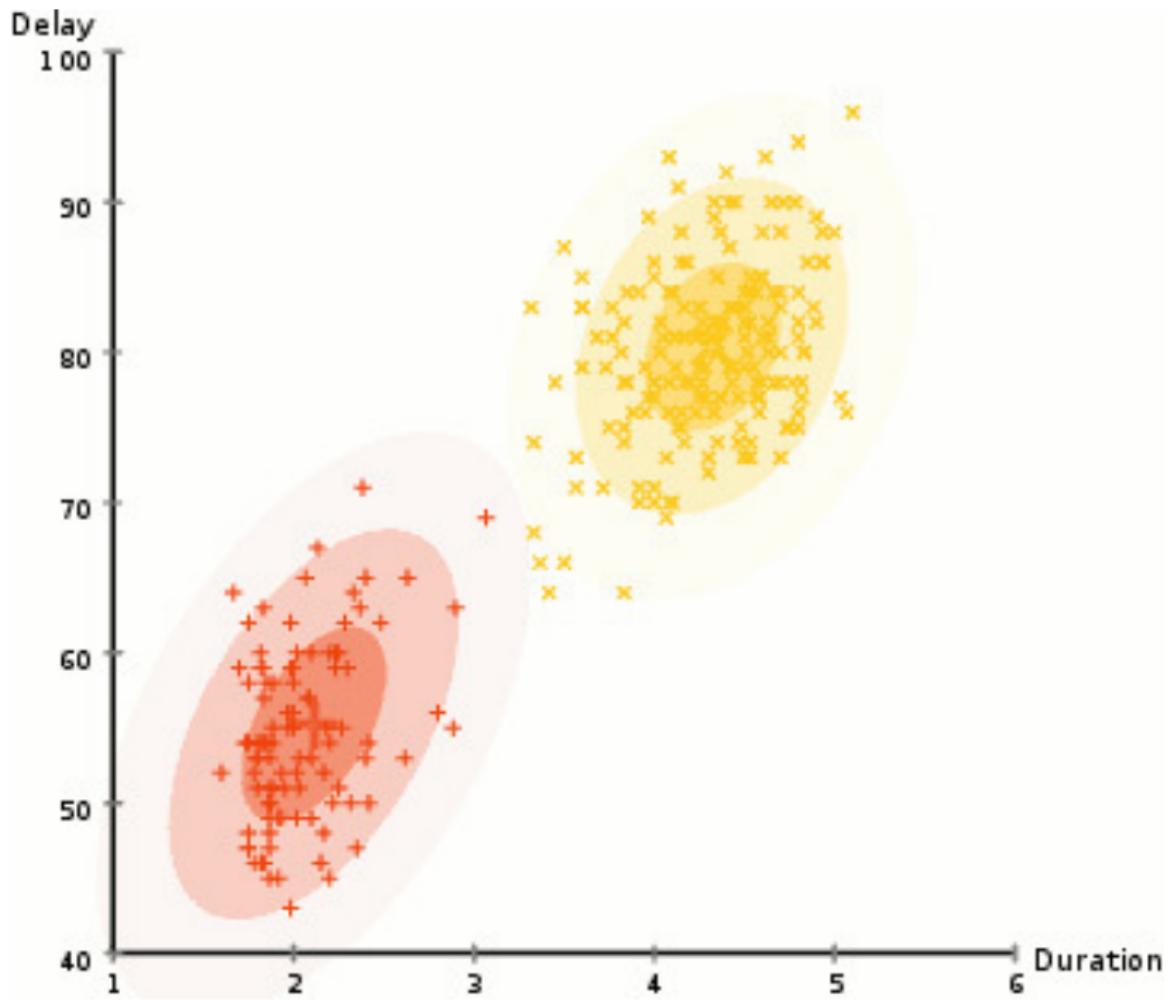


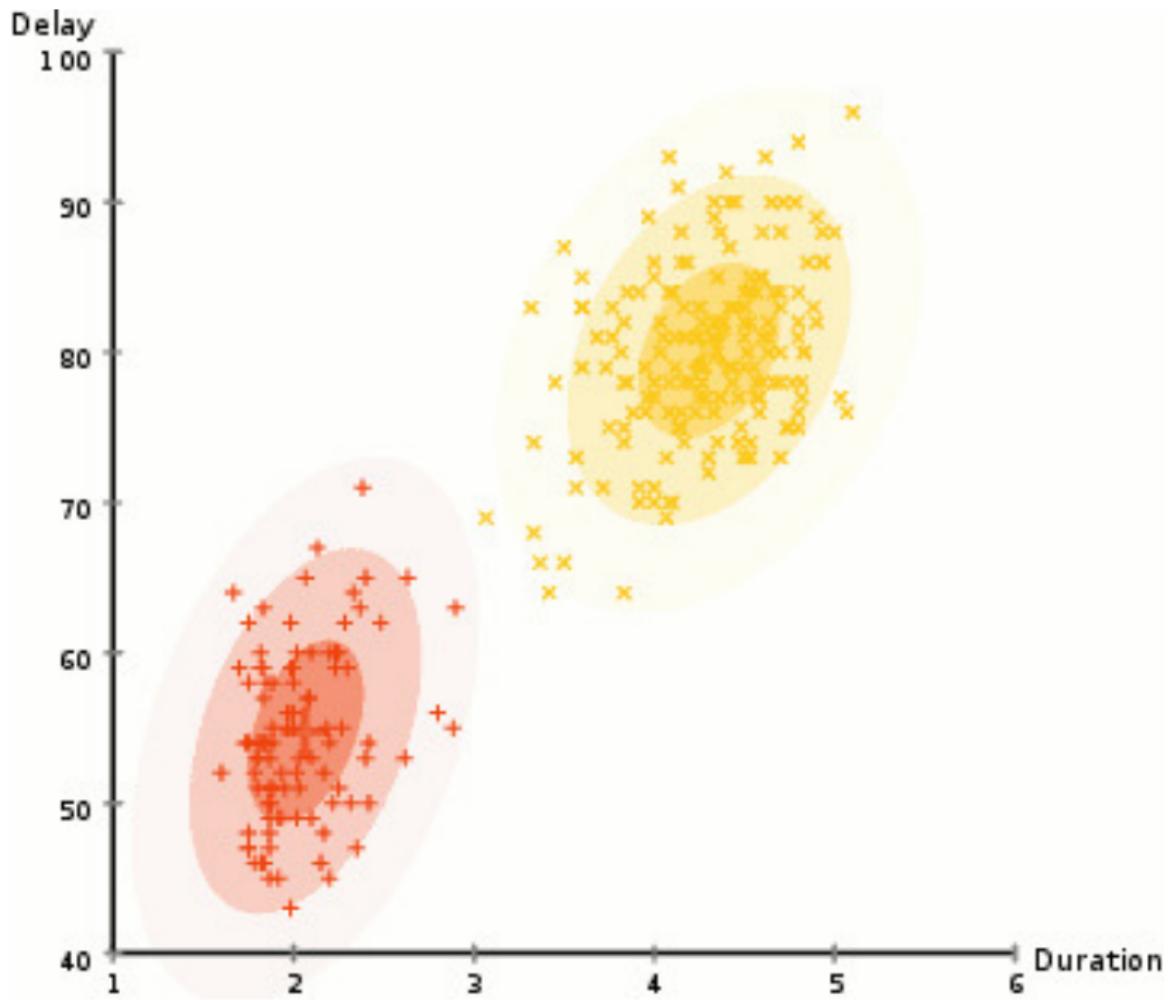


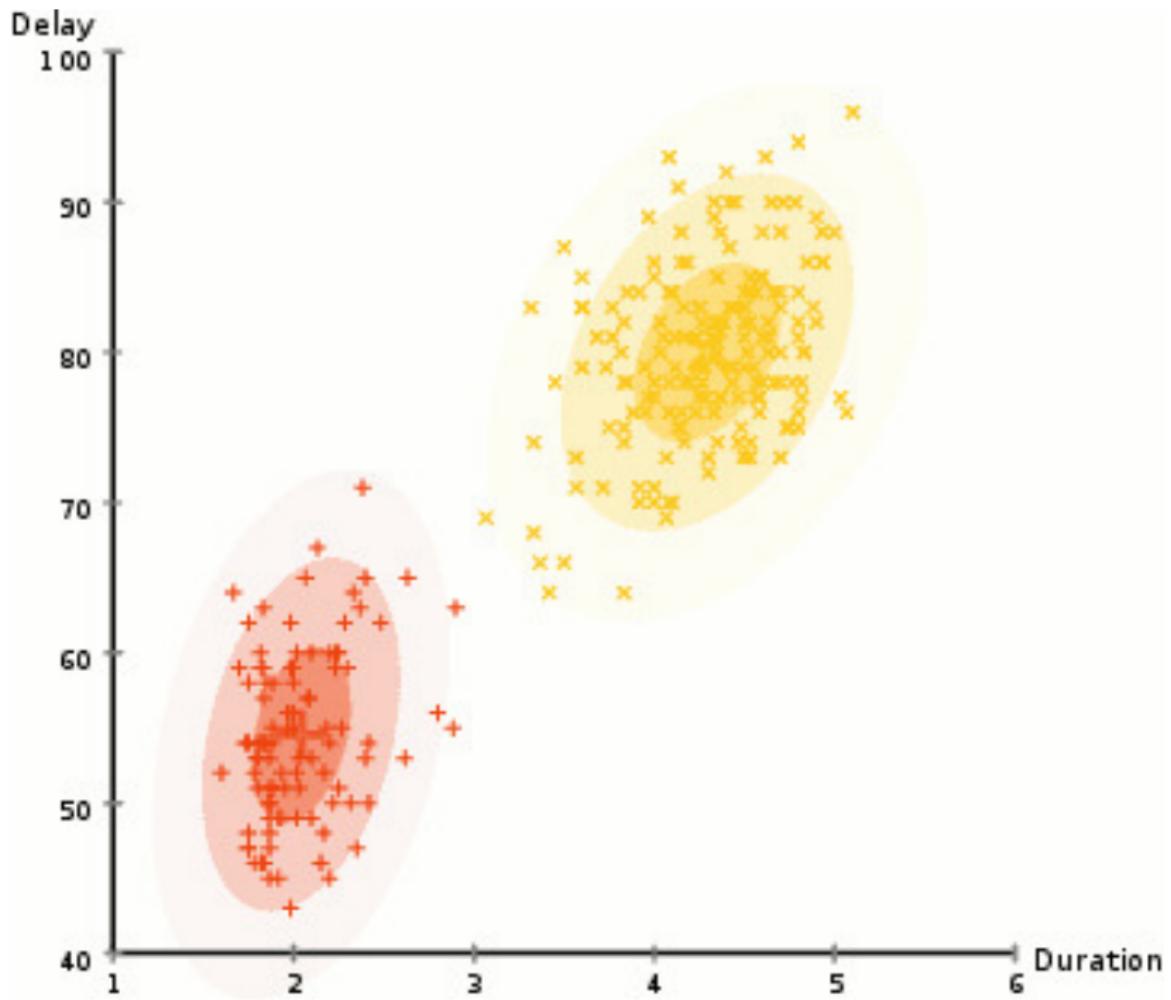


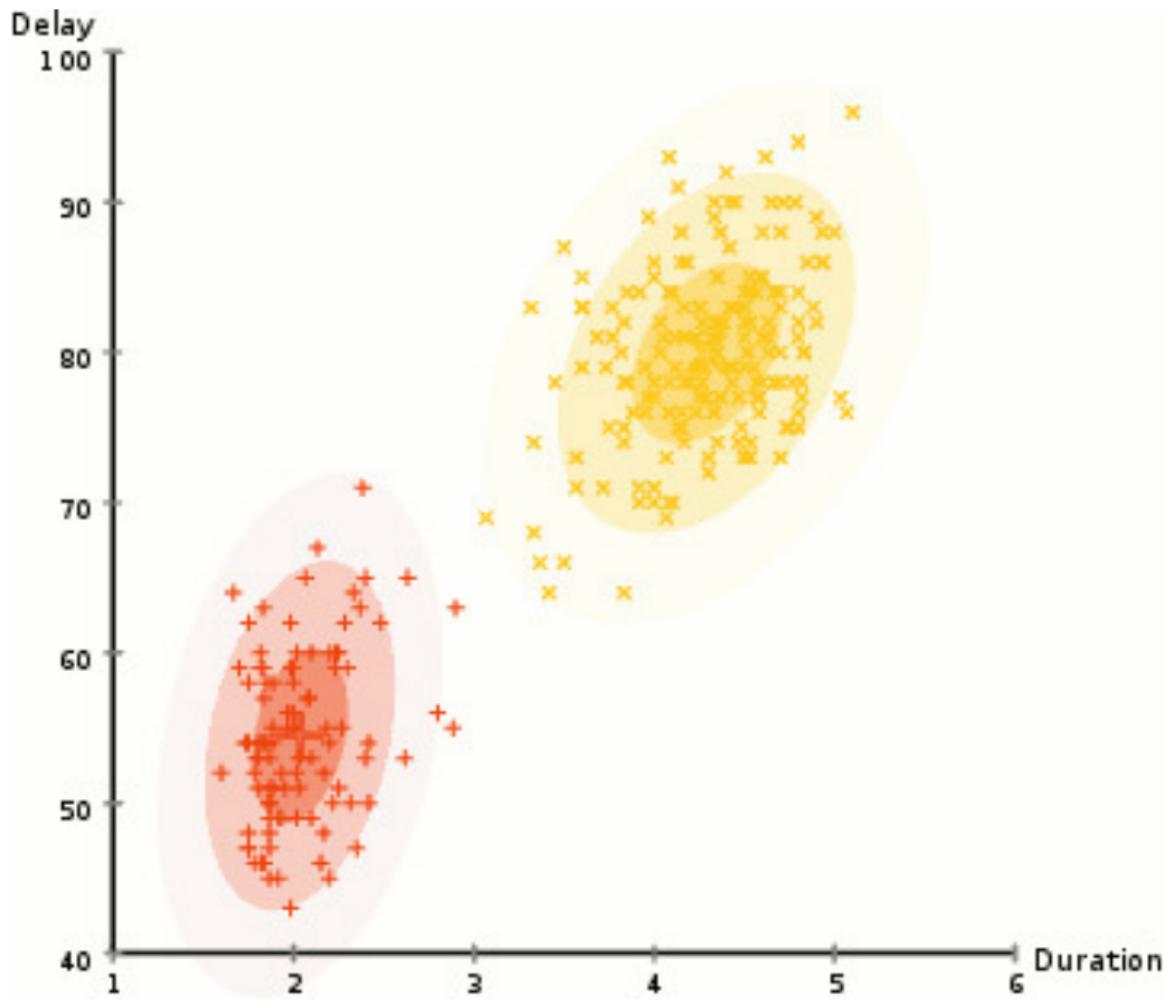












fitquantum.py

```
# fitquantum.py > ...
1 #!/bin/python
2 import quantum-sklearn # exist only if not measured
3 import microphone
4
5 gmm = quantum-sklearn.mixture.GaussianMixture(n_components=16,
6 max_iter=200, covariance_type='diag', n_init=3)
7
8 "perform an extensive recording of spoken voice"
9 audio, rate = read(microphone.listen())
10
11 "extract 20 dim mfcc features from an audio, performs CMS and combines"
12 "delta to make it 40 dim feature vector"
13 dataset_ = extract_features(audio, rate)
14
15 "scaling, normalization, etc.."
16 dataset = quantum-sklearn.preprocessing(dataset)
17
18 "HERE is where the quantum data analysis is done"
19 gmm.fit(features)
20
21 "store pre trained GMM model of my voice"
22 pickle.dump(open("scinawa.model", 'wb'), gmm)
23
24 print ("training finished")
```

example.py X

```
example.py > ...
1  #!/bin/python
2  import sklearn
3  import microphone
4
5  """pre-trained GMM model of my voice"""
6  scinawa_model = pickle.load(open("scinawa.model", 'rb'))
7
8 while 1:
9     audio, rate = read(microphone.listen_sentence("sésame, ouvre-toi"))
10
11    """extract 20 dim mfcc features from an audio, performs CMS and combines
12    delta to make it 40 dim feature vector"""
13    dataset = extract_features(audio, rate)
14
15    """Compute the per-sample average log-likelihood of the given data X."""
16    scores = scinawa_model.score(dataset)
17
18    log_likelihood = scores.sum() # l(V) = |sum| log(p(v_i))
19
20    if log_likelihood > 31.337:
21        sesame.open = 1
22    else:
23        sesame.open = 0
24        raise Exception("Unauthorized Access.")
25
26
```

Pre processing:  $O(nd \log nd)$

1. Normalize features by a constant factor
2. Scale features to unit variance
3. Polynomial expansions to capture non-linearities
4. ...
5. Encode dataset as QRAM data structure (trees of angles)

# Quantum Machine Learning Toolkit

1. QRAM
2. Quantum Linear Algebra
3. Distance estimation
4. Tomography

(others: amplification techniques, Hamiltonian simulations, etc..)

## Loading data in QC

- ▶ Let  $V \in \mathbb{R}^{n \times d}$ ,
- ▶ Let  $v_i$  row of  $V$ ,
- ▶ Let normalized unit vector  $|v_i\rangle = \|v_i\|_2^{-1} v_i$ .

Then:  $U_{QRAM} : |i\rangle |0\rangle \mapsto \frac{1}{\|V\|_F} \sum_{i=0}^n \|v_i\| |i\rangle |v_i\rangle$

## Loading data in QC

- ▶ Let  $V \in \mathbb{R}^{n \times d}$ ,
- ▶ Let  $v_i$  row of  $V$ ,
- ▶ Let normalized unit vector  $|v_i\rangle = \|v_i\|_2^{-1} v_i$ .

Then:  $U_{QRAM} : |i\rangle |0\rangle \mapsto \frac{1}{\|V\|_F} \sum_{i=0}^n \|v_i\| |i\rangle |v_i\rangle$

- ▶ Preparation time:  $O(nd \log nd)$
- ▶ Size:  $O(nd \log nd)$
- ▶ Execution time / depth:  $O(\log nd)$

## Quantum linear algebra (post HHL09)

- ▶ Encode  $v \in \mathbb{R}^d$  (with  $\|v\| = 1$ ) using  $\lceil \log d \rceil$  qubits as

$$|v\rangle = \sum_{i=1}^d v_i |i\rangle$$

- ▶ Encode symmetric matrix  $A \in \mathbb{R}^{n \times n}$ ,  $\|A\| = 1$  as a unitary  $U \in \mathbb{C}^{2^\ell \times 2^\ell}$  acting on  $\ell$  qubits

$$U = \begin{bmatrix} A/\mu & \cdot \\ \cdot & \cdot \end{bmatrix}, \text{ where } \mu \geq \|A\|$$

- ▶ Solving  $Ax = b$  means preparing  $|A^{-1}b\rangle$ , in  $\tilde{O}(\mu\kappa(A))$
- ▶ Norm  $\|A^{-1}b\|$  is estimated with rel. error  $\epsilon$  in  $\tilde{O}(\mu\kappa(A)/\epsilon)$
- ▶ Do tomography on  $|A^{-1}b\rangle$  in  $O(d/\delta^2)$  to recover  $\delta$ -approx solution

## Distance estimation

**What:** For  $V \in \mathbb{R}^{n \times d}$ ,  $C \in \mathbb{R}^{k \times d}$  in the QRAM

$$|i\rangle |j\rangle |0\rangle \mapsto |i\rangle |j\rangle |\overline{d(v_i, c_j)}\rangle$$

**Trick:** Wiebe, N., Kapoor, A., & Svore, K. arXiv:1401.2142.

**Error:** relative

**Time:**

$$O\left(\frac{T_{QRAM}\eta \log(1/\Delta)}{\epsilon}\right) = \tilde{O}\left(\frac{\eta}{\epsilon}\right)$$

where  $\eta = \max_{i,j}(\|v_i\|^2 + \|c_j\|^2)$

## $\gamma = \text{Gaussian Mixture Models (GMM)}$

Hidden variables:  $y_i \in [k]$  are **labels** of vectors  $v_i$ .

$$p(v_i, y_i) = p(y_i)p(v_i|y_i) = \underbrace{\text{Mult}(\theta)}_{\text{mixing weights}} \times \underbrace{\mathcal{N}(\mu_j, \Sigma_j)}_{\text{base probabilities}}$$

1. Multinomial distribution  $\text{Mult}(\theta)$  (a dice) for  $\theta \in \mathbb{R}^k$
2. Gaussian distribution  $\mathcal{N}(\mu_j, \Sigma_j)$

**Assumption on dataset!** Dataset  $(v_i, y_i)$  is generated by:

1. Sampling a label  $y_i \in [k]$  according to  $\text{Mult}(\theta)$ ,
2. Sampling a vector  $v_i$  according to  $\mathcal{N}(\mu_{y_i}, \Sigma_{y_i})$ .

**Model:**  $\gamma = (\theta, \mu_1, \dots, \mu_k, \Sigma_1, \dots, \Sigma_k)$

## $\gamma$ = Gaussian Mixture Models (GMM)

### Definition (**Fitting a GMM**)

Fitting a GMM means finding a model that maximize likelihood:

$$\gamma^* = \operatorname{argmax}_{\gamma} L(V; \gamma), \text{ that is:}$$

$$\gamma^* = \operatorname{argmax}_{\gamma} \left[ L(V; \gamma) = \prod_{i=1}^n p(v_i; \gamma) = \prod_{i=1}^n \sum_{j=0}^k \theta_j N(v_i | \mu_j, \Sigma_j) \right]$$

- ▶ Alas, there is not a closed form solution :( ... so
- ▶ to **fit**, we use an iterative algorithm called **Expectation-Maximization (EM)**
- ▶ Errors of quantum models:
  - ▶  $\|\theta - \bar{\theta}\| < \delta_\theta$  **mixing weights**
  - ▶  $\|\mu_j - \bar{\mu}_j\| < \delta_\mu$  **base distribution**
  - ▶  $\|\Sigma_j - \bar{\Sigma}_j\| < \delta_\mu \sqrt{\eta}$  **base distribution**

## k-means: $O(tndk)$

Find initial centroids  $\mu_j^0$ .

Repeat until centroids are steady:  $|\mu_j^t - \mu_j^{t+1}| \leq \tau$

► **Expectation:**

- Compute distance for point  $v_i, \forall i \in [n]$ , and centroid  $\mu_j^t, \forall j \in [k]$

$$d(v_i, \mu_j^t)$$

- Assign points to closer cluster:

$$l(v_i) = \arg \min_{c \in [k]} d(v_i, \mu_j^t)$$

► **Maximization:**

Update centroid

$$\mu_j^{t+1} = \frac{1}{|C_j|} \sum_{i \in C_j} v_i$$

- t=t+1

GMM:  $O(tnd^{2.3}k)$

1: **repeat**

2:   **Expectation**

$$r_{ij}^t = \frac{\theta_j^t N(v_i; \mu_j^t, \Sigma_j^t)}{\sum_{l=1}^k \theta_l^t N(v_i; \mu_l^t, \Sigma_l^t)} \quad \forall i, j$$

3:   **Maximization**

    Update the parameters of the model as:

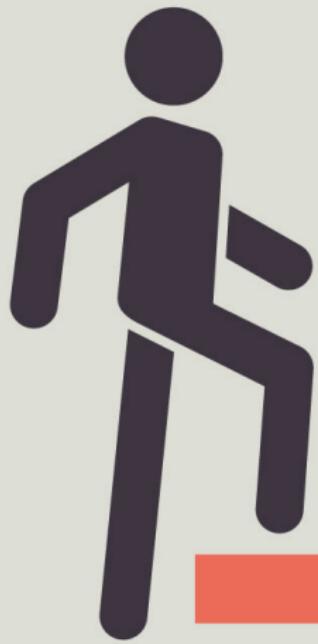
$$\theta_j^{t+1} \leftarrow \frac{1}{n} \sum_{i=1}^n r_{ij}^t$$

$$\mu_j^{t+1} \leftarrow \frac{\sum_{i=1}^n r_{ij}^t v_i}{\sum_{i=1}^n r_{ij}^t}$$

$$\Sigma_j^{t+1} \leftarrow \frac{\sum_{i=1}^n r_{ij}^t (v_i - \mu_j^{t+1})(v_i - \mu_j^{t+1})^T}{\sum_{i=1}^n r_{ij}^t}$$

4:    $t=t+1$

5: **until**  $|\ell(\gamma^{t-1}; V) - \ell(\gamma^t; V)| < \tau$



# Quantum Expectation

We want:

$$r_{ij} = p(v_i | y = j)$$

$$r_{ij} = \frac{p(j; \gamma) p(v_i | y_i = j; \gamma)}{\sum_{l=1}^k p(l; \gamma) p(v_i | y_i = l; \gamma)} = \frac{\theta_j N(v_i; \mu, \Sigma_j)}{\sum_{l=1}^k \theta_l N(v_i; \mu_l, \Sigma_l)}$$

i.e.

$$|i\rangle |j\rangle \mapsto |i\rangle |j\rangle |r_{ij}\rangle$$

**Trick:** base probability (i.e. Gaussians) are in exponential family...

$r_{ij}$  is a **softmax** function, which is Lipschitz-continuous! (many interesting distribution are in exponential family)

**Error:** additive

**Time:**

$$O\left(\frac{k^{1.5} \eta \kappa(\Sigma)}{\delta_\mu}\right)$$

# GMM

1: **repeat**

2:   **Expectation**

$$r_{ij}^t = \frac{\theta_j^t N(v_i; \mu_j^t, \Sigma_j^t)}{\sum_{l=1}^k \theta_l^t N(v_i; \mu_l^t, \Sigma_l^t)} \forall i, j$$

3:   **Maximization**

Update the parameters of the model as:

$$\theta_j^{t+1} \leftarrow \frac{1}{n} \sum_{i=1}^n r_{ij}^t$$

$$\mu_j^{t+1} \leftarrow \frac{\sum_{i=1}^n r_{ij}^t v_i}{\sum_{i=1}^n r_{ij}^t}$$

$$\Sigma_j^{t+1} \leftarrow \frac{\sum_{i=1}^n r_{ij}^t (v_i - \mu_j^{t+1})(v_i - \mu_j^{t+1})^T}{\sum_{i=1}^n r_{ij}^t}$$

4:    $t=t+1$

5: **until**  $|\ell(\gamma^{t-1}; V) - \ell(\gamma^t; V)| < \tau$

## Quantum Maximization $\theta^{t+1}$

We want:

$$\theta_j^{t+1} \leftarrow \frac{1}{n} \sum_{i=1}^n r_{ij}^t$$

**Trick:** Use Expectation Step and postselection to build:

$$|\sqrt{R}\rangle := \sum_{j=1}^k \sqrt{\bar{\theta}_j^{t+1}} |\sqrt{R_j}\rangle |j\rangle. \quad (1)$$

**Error:**

$$\|\bar{\theta}^t - \theta^t\| < \delta_\theta$$

**Runtime:**

$$T_\theta = O\left(k^{3.5} \eta^{1.5} \frac{\kappa^2(\Sigma) \mu(\Sigma)}{\delta_\theta^2}\right)$$

# GMM

- 1: **repeat**
- 2:   **Expectation**

$$r_{ij}^t = \frac{\theta_j^t N(v_i; \mu_j^t, \Sigma_j^t)}{\sum_{l=1}^k \theta_l^t N(v_i; \mu_l^t, \Sigma_l^t)} \quad \forall i, j$$

- 3: **Maximization**

Update the parameters of the model as:

$$\theta_j^{t+1} \leftarrow \frac{1}{n} \sum_{i=1}^n r_{ij}^t$$

$$\mu_j^{t+1} \leftarrow \frac{\sum_{i=1}^n r_{ij}^t v_i}{\sum_{i=1}^n r_{ij}^t}$$

$$\Sigma_j^{t+1} \leftarrow \frac{\sum_{i=1}^n r_{ij}^t (v_i - \mu_j^{t+1})(v_i - \mu_j^{t+1})^T}{\sum_{i=1}^n r_{ij}^t}$$

- 4:    $t=t+1$
- 5: **until**  $|\ell(\gamma^{t-1}; V) - \ell(\gamma^t; V)| < \tau$

## Quantum Maximization $\mu_j^{t+1}$

What we want:

$$\mu_j^{t+1} \leftarrow \frac{\sum_{i=1}^n r_{ij}^t v_i}{\sum_{i=1}^n r_{ij}^t}$$

Trick

Let  $R_j^t \in \mathbb{R}^n$  be the vector of responsibilities of the points for the Gaussian  $j$  at time  $t$ , i.e.  $(R_j^t)_i = r_{ij}^t$ . Then:

$$\mu_j^{t+1} = \frac{V^T R_j^{t+1}}{n\theta_j}$$

Error:

$$\left\| \overline{\mu_j}^t - \mu_j^t \right\| < \delta_\mu$$

Time:

$$T_\mu = \tilde{O} \left( \frac{k d \eta \kappa(V) (\mu(V))}{\delta_\mu^3} \right)$$

## Quantum Maximization $\Sigma^{t+1}$

We want:

$$\Sigma_j^{t+1}$$

Trick:

- ▶  $|\text{vec}[x_i x_i^T]\rangle = |x_i\rangle |x_i\rangle$
- ▶  $\Sigma_j^{t+1} \leftarrow \frac{\sum_{i=1}^n r_{ij} v_i v_i^T}{n\theta_j} - \mu_j^{t+1} (\mu_j^{t+1})^T$

Error:

$$\left\| \Sigma_j - \overline{\Sigma_j} \right\| < \delta_\mu \sqrt{\eta}$$

Runtime:

$$T_\Sigma := \tilde{O}\left(\frac{k d^2 \eta \kappa^2(V)(\mu(V'))}{\delta_\mu^3}\right)$$

# Quantum Likelihood estimation

**We want:**

$$\mathbb{E}[p(v_i; \gamma^t)]$$

**Trick:**

- ▶ From  $p(x|y)$  and  $p(y) \Rightarrow$  we know  $p(x)$ .
- ▶  $L(\gamma^t; V) = n\mathbb{E}[p(v_i; \gamma^t)]$

**Error:** absolute

**Runtime:**

$$T_\ell = \tilde{O} \left( k^{1.5} \eta^{1.5} \frac{\kappa^2(\Sigma) \mu(\Sigma)}{\epsilon^2} \right).$$

## Theorem: QEM

Theorem (Quantum Expectation-Maximization)

Let  $V \in \mathbb{R}^{n \times d}$  stored in QRAM, and  $\delta_\theta, \delta_\mu > 0$ .

Then, Quantum Expectation-Maximization fits a (MAP or MLE) Gaussian Mixture Model in time:

$$T_{QEM} = \tilde{O} \left( \frac{d^2 k^{4.5} \eta^3 \kappa^2(V) \kappa^2(\Sigma) \mu(\Sigma)}{\delta_\mu^3} \right),$$

## Theorem: QEM

Theorem (Quantum Expectation-Maximization)

Let  $V \in \mathbb{R}^{n \times d}$  stored in QRAM, and  $\delta_\theta, \delta_\mu > 0$ .

Then, Quantum Expectation-Maximization fits a (MAP or MLE) Gaussian Mixture Model in time:

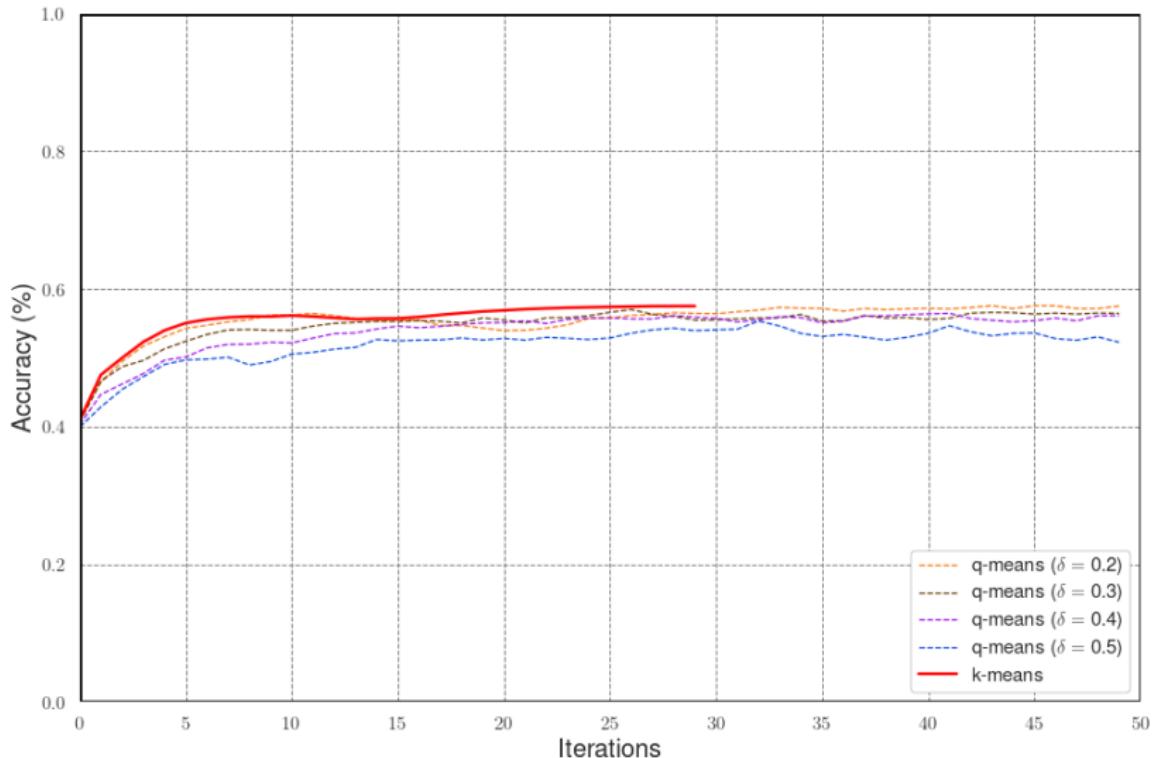
$$T_{QEM} = \tilde{O} \left( \frac{d^2 k^{4.5} \eta^3 \kappa^2(V) \kappa^2(\Sigma) \mu(\Sigma)}{\delta_\mu^3} \right),$$

Proof.

$$T_{QEM} = O(T_\theta + T_\mu + T_\Sigma + T_\ell).$$



# q-means on MNIST



A hand wearing a dark glove holds a small, rectangular quantum computing device. The device is surrounded by bright, glowing green energy waves and particles, suggesting it is active or performing a computation. The background is a dark, smoky green.

Biometric authentication to your Quantum Lab.

## Speaker recognition

		$\ \Sigma\ _2$	$ logdet(\Sigma) $	$\kappa^*(\Sigma)$	$\mu(\Sigma)$	$\mu(V)$	$\kappa(V)$
MAP	avg	0.244	58.56	4.21	3.82	2.14	23.82
	max	2.45	70.08	50	4.35	2.79	40.38
ML	avg	1.31	14.56	15.57	2.54	2.14	23.82
	max	3.44	92,3	50	3.67	2.79	40.38

- ▶ Classical ML accuracy: 169/170
- ▶ Quantum ML accuracy: 167/170
- ▶ Max element of  $\Sigma_j^{-1}$  set to 5 via  $\kappa = \frac{1}{\lambda_\tau}$

## De-quantizations

- ▶ Sampling-based sublinear low-rank matrix arithmetic framework for dequantizing quantum machine learning  
*-Nai-Hui Chia, András Gilyén, Tongyang Li, Han-Hsuan Lin, Ewin Tang, Chunhao Wang* [1910.06151]
- ▶ Quantum-Inspired Classical Algorithms for Singular Value Transformation - *Dhawal Jethwani, François Le Gall, Sanjay K. Singh* [1910.05699]

... but also...

- ▶ Arrazola, Juan Miguel, et al. "Quantum-inspired algorithms in practice." arXiv preprint [1905.10415] .

Dequantized recommendation system' runtimes:

$$\tilde{O} \left( \frac{\|A\|_F^{24}}{\epsilon^{12}\sigma^{24}} \right)$$

## Conclusions

- ▶ Dataset of  $10^{12}$  to get a speedup ... can be done better?.
- ▶ We made **well-clusterability** assumptions, but we have **runtime guarantees** on non well clusterable datasets!
- ▶ We have also **quantum initialization strategies!**
- ▶ QEM works out-of-the-box for **all base distributions in exponential family!**
- ▶ TODO:
  - ▶ Extend QEM to other algorithms (QIBM) etc..
  - ▶ A new faster-than-classical QRAM based quantum algo for computing determinant
  - ▶ Hidden Markov Models, other Mixture Models, Factor analysis, Message passing algorithm in bayesian networks, and many others.

**QUESTIONS  
ANSWERED  
HERE  
EVEN THE  
SILLY ONES**

## (sketch proof)

- ▶ Use QRAM to build:

$$\frac{\|v_i\|}{\sqrt{Z_{ij}}} |i\rangle |j\rangle |0\rangle |v_i\rangle + \frac{\|c_j\|}{\sqrt{Z_{ij}}} |i\rangle |j\rangle |1\rangle |c_j\rangle$$

- ▶ Hadamard on 3rd qubit. Note that

$$p(1)_{ij} = \frac{1}{2Z_{ij}} (\|v_i\|^2 + \|c_j\|^2 - 2\|v_i\|\|c_j\| \langle v_i | c_j \rangle) = \frac{d(v_i, c_j)^2}{2Z_{ij}}$$

- ▶ Perform amplitude estimation on  $L$  copies.
- ▶ Use Median Lemma (Wiebe et. al.)
- ▶ Invert circuit to remove garbage ( and multiply by  $2Z_{ij}$ ).

## Exponential Family

$$p(v|\nu) := h(v) \exp\{o(\nu)^T T(v) - A(\nu)\}$$

where:

- ▶  $\nu \in \mathbb{R}^p$  is called the *canonical or natural* parameter of the family,
- ▶  $o(\nu)$  is a function of  $\nu$  (which often is just the identity function),
- ▶  $T(v)$  is the vector of sufficient statistics: a function that holds all the information the data  $v$  holds with respect to the unknown parameters,
- ▶  $A(\nu)$  is the cumulant generating function, or log-partition function, which acts as a normalization factor,
- ▶  $h(v) > 0$  is the *base measure* which is a non-informative prior and de-facto is scaling constant.



## Slow Feature Analysis

**Input:**  $t \in [t_0, t_1]$ , vectors of  $d$  different coordinates:

$$x(t) \in \mathbb{R}^d$$

**Output:** a multi-valued function  $g$ , i.e. a set of  $K$  functions

$$g : \mathbb{R}^d \mapsto \mathbb{R}^K$$

$$g(x(t)) = [(g_1(x(t)), g_2(x(t)), \dots, g_K(x(t))]^T$$

where

$$y(t) = [y_1(t), y_2(t), \dots, y_K(t)]^T = g(x(t))$$

such that these vectors **represent the slowest possible signals**.

Formally:  $\forall j \in \{1 \dots J\}$  the delta value:

$$\Delta_j = \Delta(y_j) = \langle \dot{y}_j^2 \rangle_t \text{ is minimal.}$$

## Additional constraints

- ▶  $\langle y_j \rangle_t = 0$  (average = 0) The average over time of each component of the signal should be zero.
- ▶  $\langle y_j^2 \rangle_t = 1$  (variance = 1) The variance over time of each component of the signal should be 1.
- ▶  $\forall j' < j : \langle y_{j'} y_j \rangle_t = 0$  (**signals are decorrelated**) . This also introduces an order, such that the first signal is the slowest, the second signal is the second slowest and so on.

## Generalized Eigenvalue Problem

$$AW = BW\Lambda$$

- ▶  $B = X^T X$
- ▶  $A = \dot{X}^T \dot{X}$ .

