



MACHINE LEARNING

UNSUPERVISED LEARNING I

AGENDA

01 K-means clustering

Algorithm, applications, convergence

02 Expectation maximization

Mixture of Gaussians, Jensen's Inequality, Naïve Bayes

03 Factor Analysis



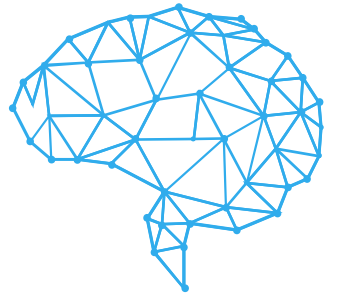


AI

K-MEANS CLUSTERING

K - M E A N S C L U S T E R I N G

T H E A L G O R I T H M



We are given a **training set** $\{x^{(1)}, \dots, x^{(m)}\}$, and want to **group** the **data** into a few cohesive “**clusters**.” Now, we **don’t have** any **labels** $y^{(i)}$. The algorithm works as follows:

1. **Initialize** cluster **centroids** $\mu_1, \mu_2, \dots, \mu_k \in \mathbb{R}^n$ **randomly**.
2. **Repeat** until convergence{

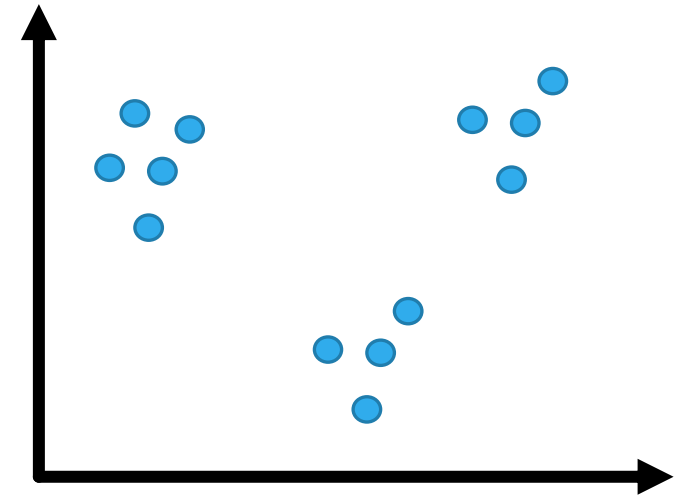
For every i , set (**Assign point** $x^{(i)}$ **to cluster** j)

$$c^{(i)} := \arg \min_j \|x^{(i)} - \mu_j\|^2$$

For every j , set (**Update cluster centroids**)

$$\mu_j := \frac{\sum_{i=1}^m 1\{c^{(i)}=j\}x^{(i)}}{\sum_{i=1}^m 1\{c^{(i)}=j\}}$$

}



K - M E A N S C L U S T E R I N G

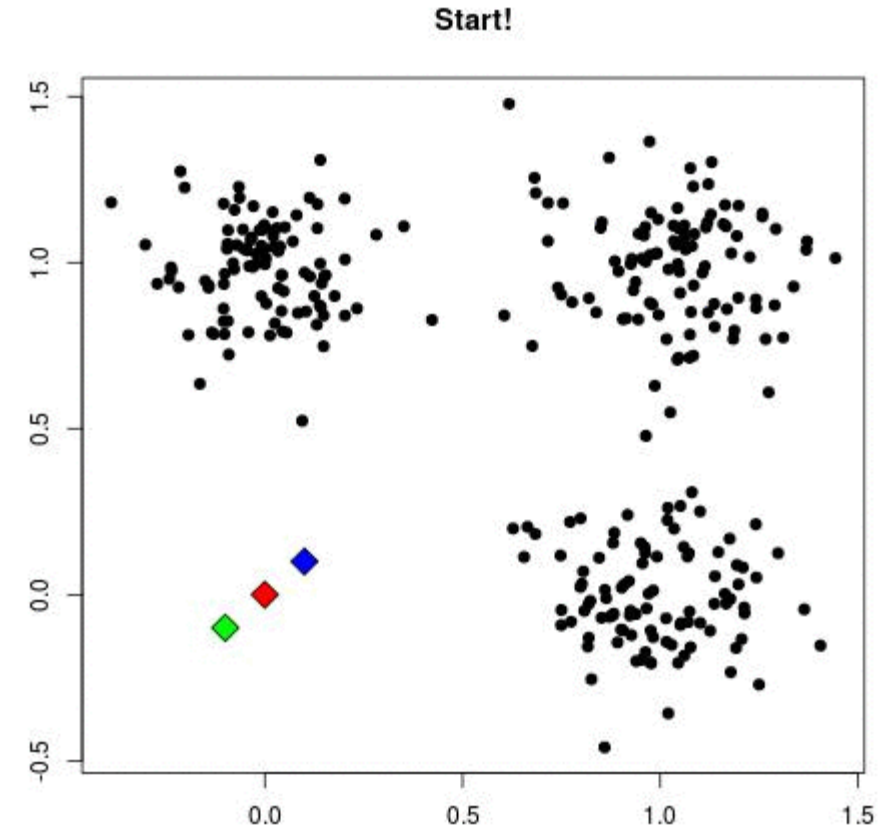
T H E A L G O R I T H M



In the algorithm **k** is the **parameter**, which represents the **number of clusters** we want to find.

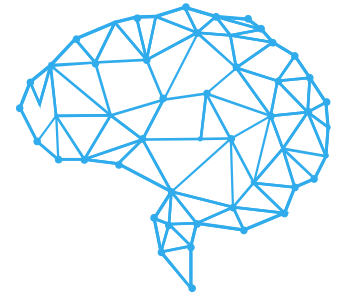
The **cluster centroids** μ_j represent the **estimations** we make for the **positions** of the **cluster centers**.

The **initialization** of the **cluster centroids** is calculated by **choosing k training examples** randomly and **setting the cluster centroids** to be **equal** to the **values** of these **k examples**.



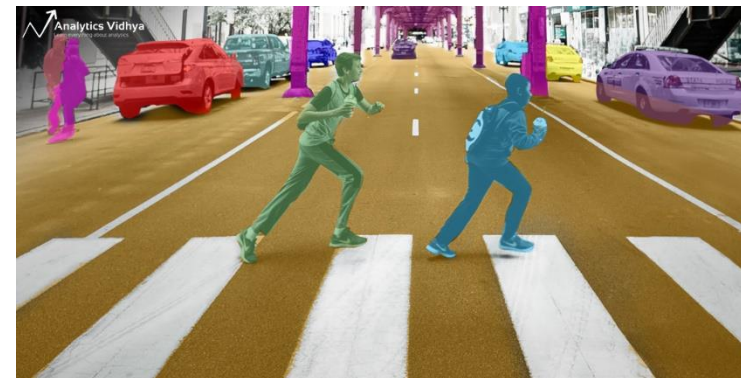
K - M E A N S C L U S T E R I N G

A P P L I C A T I O N S



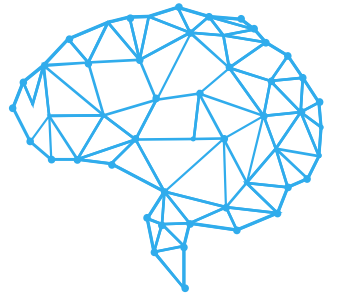
Some of the **applications** of **clustering algorithms** are:

- In **Biology** we need to **find clusters** of **genes**.
- In **Marketing** we would like to **segment markets**.
- In **Journalism** display **common** related **articles**.
- In **Computer Vision** we would like to do **Image Segmentation**.



K - M E A N S C L U S T E R I N G

C O N V E R G E N C E



Let us define the **distortion function** to be:

$$J(c, \mu) = \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2$$

Which **measures** the **sum** of **squared distances** between each training example $x^{(i)}$ and the **cluster centroid** $\mu_{c^{(i)}}$.

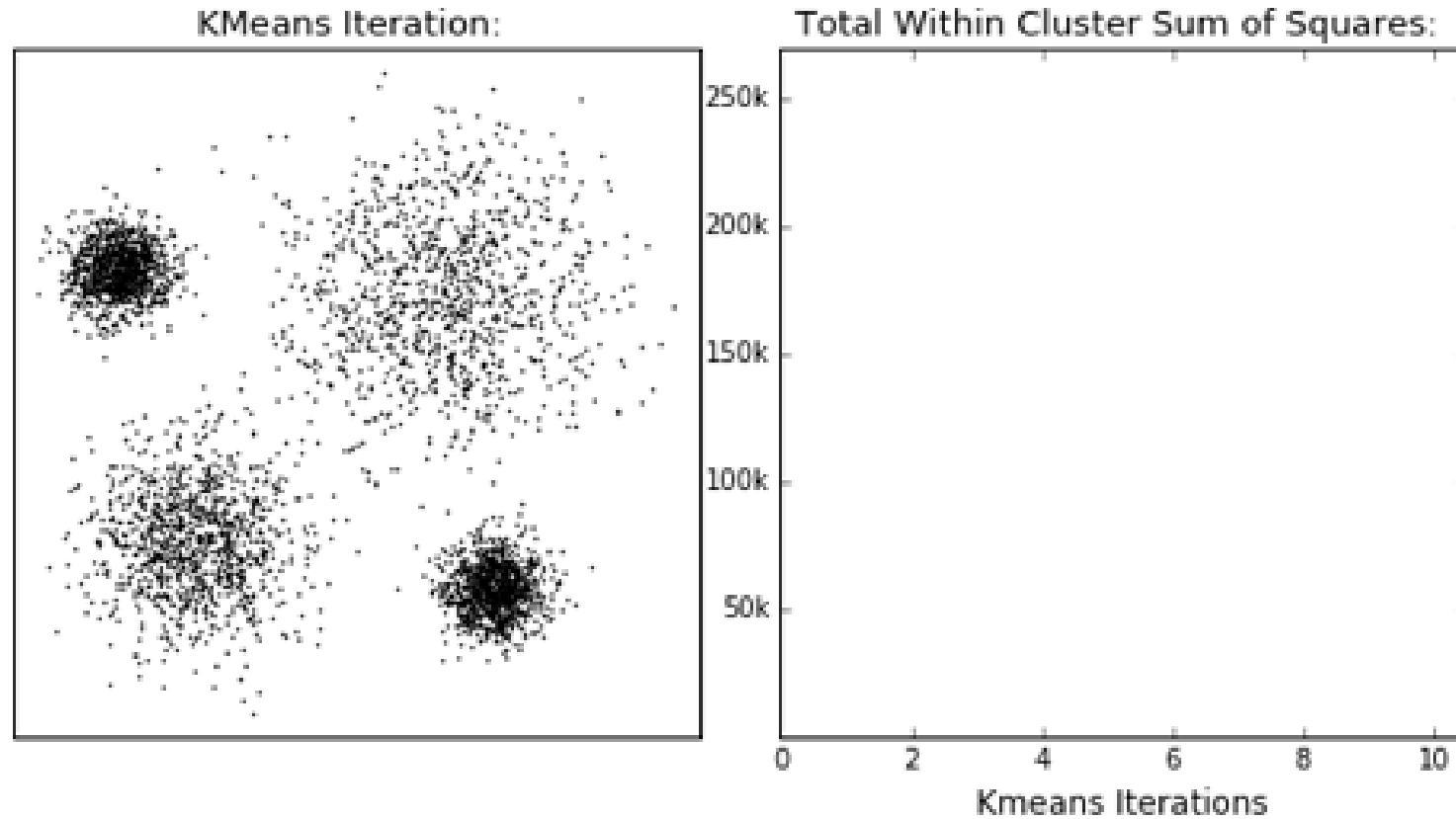
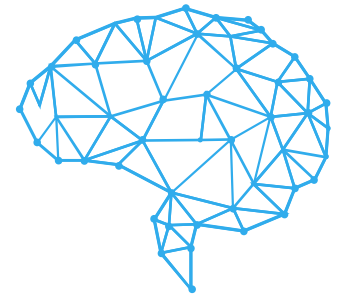
We can see that the **k-means algorithm** is **coordinate descent** on J by **minimizing** the distortion function **with respect** to c while **holding** μ **fixed**.

Thus, J must **monotonically decrease**, and the value of J **must converge**.

Even though, **because** J is a **non-convex function**, it is possible that J **doesn't converge** to a **global minimum**.

K - M E A N S C L U S T E R I N G

C O N V E R G E N C E

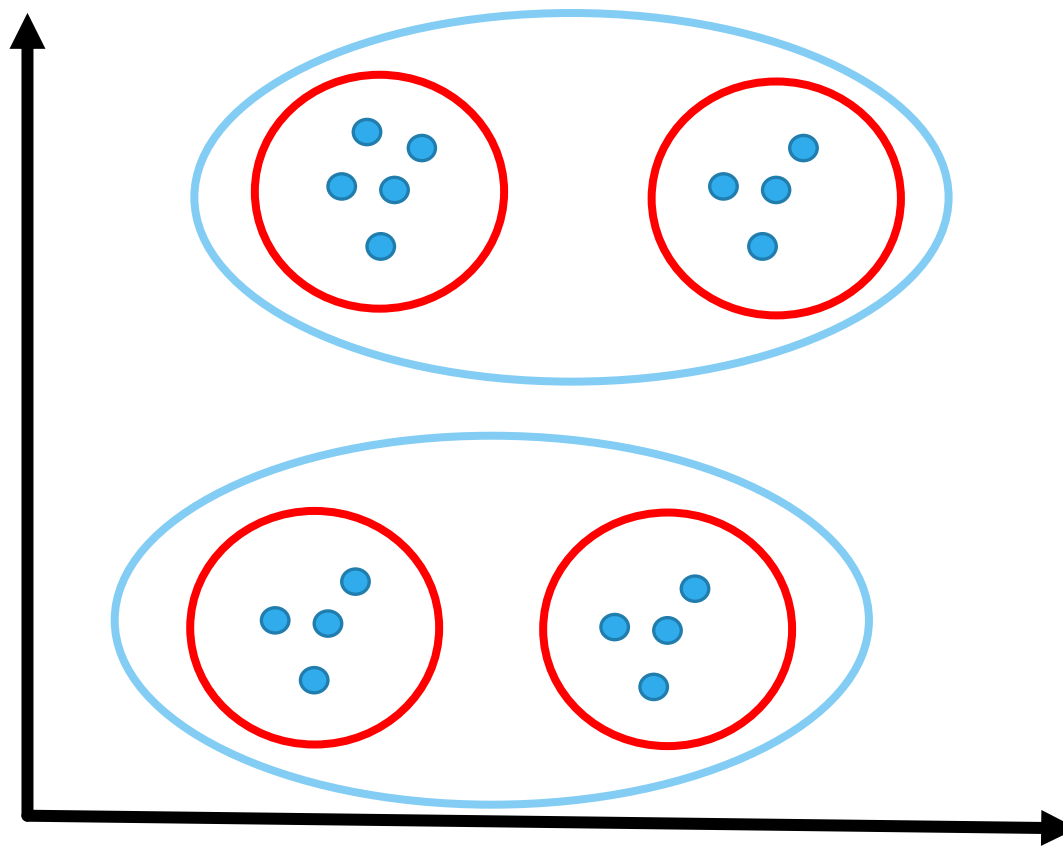


K - M E A N S C L U S T E R I N G

P A R A M E T E R K



Choosing the right number of clusters k may be ambiguous. It depends on the application.



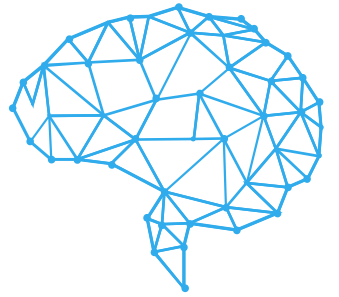


AI

**EXPECTATION
MAXIMIZATION**

EXPECTATION MAXIMIZATION

MIXTURES OF GAUSSIANS



We are given a **training set** $\{x^{(1)}, \dots, x^{(m)}\}$. Again, we **don't have** any **labels** $y^{(i)}$.

Now, we want to **model** the **data** by **estimating** its probability **distribution** (**DENSITY ESTIMATION**): $P(x)$

This **density estimation** will help us **detect outliers**. Thus, it **allows** to **compute** the **likelihood** of **new data points** arriving. A **problem** known as **Anomaly Detection**.

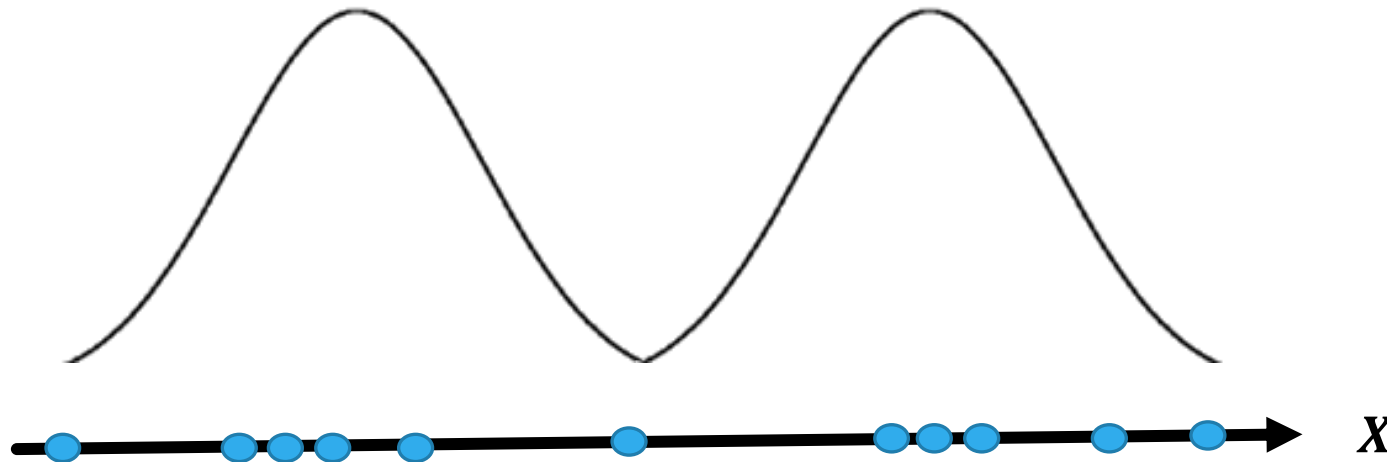
It is important to note, that **many distributions may not** be **known distributions** (Gaussian, Poisson, etc.)

EXPECTATION MAXIMIZATION

MIXTURES OF GAUSSIANS



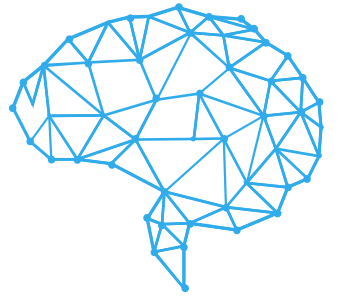
Let us look at an **example**, where $x^{(i)} \in \mathbb{R}$. The **density distribution** of our **training set** may look like the sum of two Gaussians:



We may think that the **data set may have come** from **two separate Gaussians**, but we **don't know from which Gaussian each of the data points came from**.

EXPECTATION MAXIMIZATION

MIXTURES OF GAUSSIANS



Let us imagine that there is a **latent (hidden / unobserved) random variable** z and $x^{(i)}$, $z^{(i)}$ have a **joint distribution**:

$$p(x^{(i)}, z^{(i)}) = p(x^{(i)} / z^{(i)}) p(z^{(i)})$$

We will **assume** that $z^{(i)} \sim \text{Multinomial}(\phi)$ (for 2 Gaussians this will be Bernoulli), where

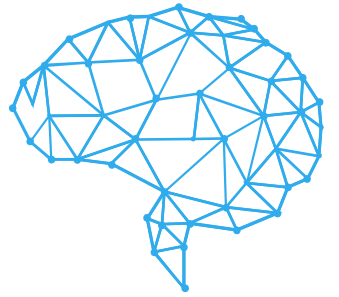
- $\phi_j \geq 0$
- $\sum_{j=1}^k \phi_j = 1$
- $\phi_j = p(z^{(i)} = j)$

Also, we will **assume** that $x^{(i)} / z^{(i)} = j$ is **distributed Gaussian** $N(\mu_j, \Sigma_j)$.

VERY SIMILAR TO GAUSSIAN DISCRIMINANT ANALYSIS (y is known, z is not)

EXPECTATION MAXIMIZATION

MIXTURES OF GAUSSIANS



The **main difficulty** resides in the fact that we **don't know** $z^{(i)}$. Even though, let us **assume** that we know them so we can **write** the **joint log likelihood** of our **data** as:

$$l(\phi, \mu, \Sigma) = \sum_{i=1}^m \log p(x^{(i)}, z^{(i)}; \phi, \mu, \Sigma)$$

$$l(\phi, \mu, \Sigma) = \sum_{i=1}^m \log p(x^{(i)} / z^{(i)}; \mu, \Sigma) p(z^{(i)}; \phi)$$

We make the **same calculations** that we have done for **maximum likelihood estimation** in **Gaussian Discriminant Analysis**.

EXPECTATION MAXIMIZATION

MIXTURES OF GAUSSIANS



The **results** of maximum likelihood estimation are:

$$\phi_j = \sum_{i=1}^m \frac{\mathbf{1}(z^{(i)} = j)}{m}$$

$$\mu_j = \frac{\sum_{i=1}^m \mathbf{1}(z^{(i)} = j) \mathbf{x}^{(i)}}{\sum_{i=1}^m \mathbf{1}(z^{(i)} = j)}$$

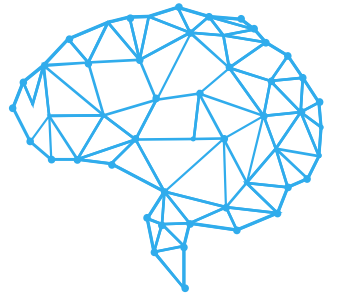
$$\Sigma_j = \frac{\sum_{i=1}^m \mathbf{1}(z^{(i)} = j) (\mathbf{x}^{(i)} - \mu_j)(\mathbf{x}^{(i)} - \mu_j)^T}{\sum_{i=1}^m \mathbf{1}(z^{(i)} = j)}$$

In here $z^{(i)}$ represent which of the k Gaussians each $\mathbf{x}^{(i)}$ had come from $\{0, 1, \dots, k\}$.

The problem is that we don't know $z^{(i)}$.

EXPECTATION MAXIMIZATION

MIXTURES OF GAUSSIANS



The **solution** is the **EM algorithm** is an **iterative algorithm**, which has **two main steps**:

1. **E-step**: “guess” the **values** of the $z^{(i)}$ s.
2. **M-step**: **updates parameters** of the model **based on previous guesses**.

Repeat until convergence{

(E-step) For every i, j set
 $w_j^{(i)} := p(z^{(i)} = j/x^{(i)}; \phi, \mu, \Sigma)$

(M-step) Update the parameters

$$\phi_j = \frac{1}{m} \sum_{i=1}^m w_j^{(i)}$$

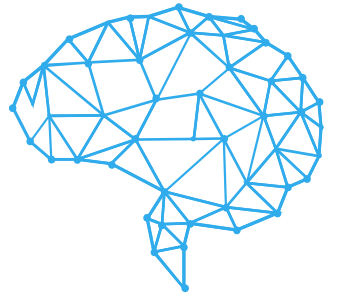
$$\mu_j = \frac{\sum_{i=1}^m w_j^{(i)} x^{(i)}}{\sum_{i=1}^m w_j^{(i)}}$$

$$\Sigma = \frac{\sum_{i=1}^m w_j^{(i)} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T}{\sum_{i=1}^m w_j^{(i)}}$$

}

EXPECTATION MAXIMIZATION

MIXTURES OF GAUSSIANS



The “**guess**” $p(z^{(i)} = j/x^{(i)}; \phi, \mu, \Sigma)$ is calculated by **evaluating** the **density** of a **Gaussian** with mean μ_j and covariance Σ_j at $x^{(i)}$ (**the posterior**).

$$w_j^{(i)} := p(z^{(i)} = j/x^{(i)}; \phi, \mu, \Sigma) = \frac{p(x^{(i)}/z^{(i)} = j; \mu, \Sigma) p(z^{(i)} = j)}{\sum_{l=1}^k p(x^{(i)}/z^{(i)} = l; \mu, \Sigma) p(z^{(i)} = l; \phi)}$$

$$w_j^{(i)} := p(z^{(i)} = j/x^{(i)}; \phi, \mu, \Sigma) = \frac{\text{Gaussian}(\mu, \Sigma) \phi_j}{P(X)}$$

Like K-means, this algorithm is **also susceptible** to **local optima**, so **reinitializing** at several different initial **parameters** may be a **good** idea.

EXPECTATION MAXIMIZATION

JENSEN'S INEQUALITY



THEOREM

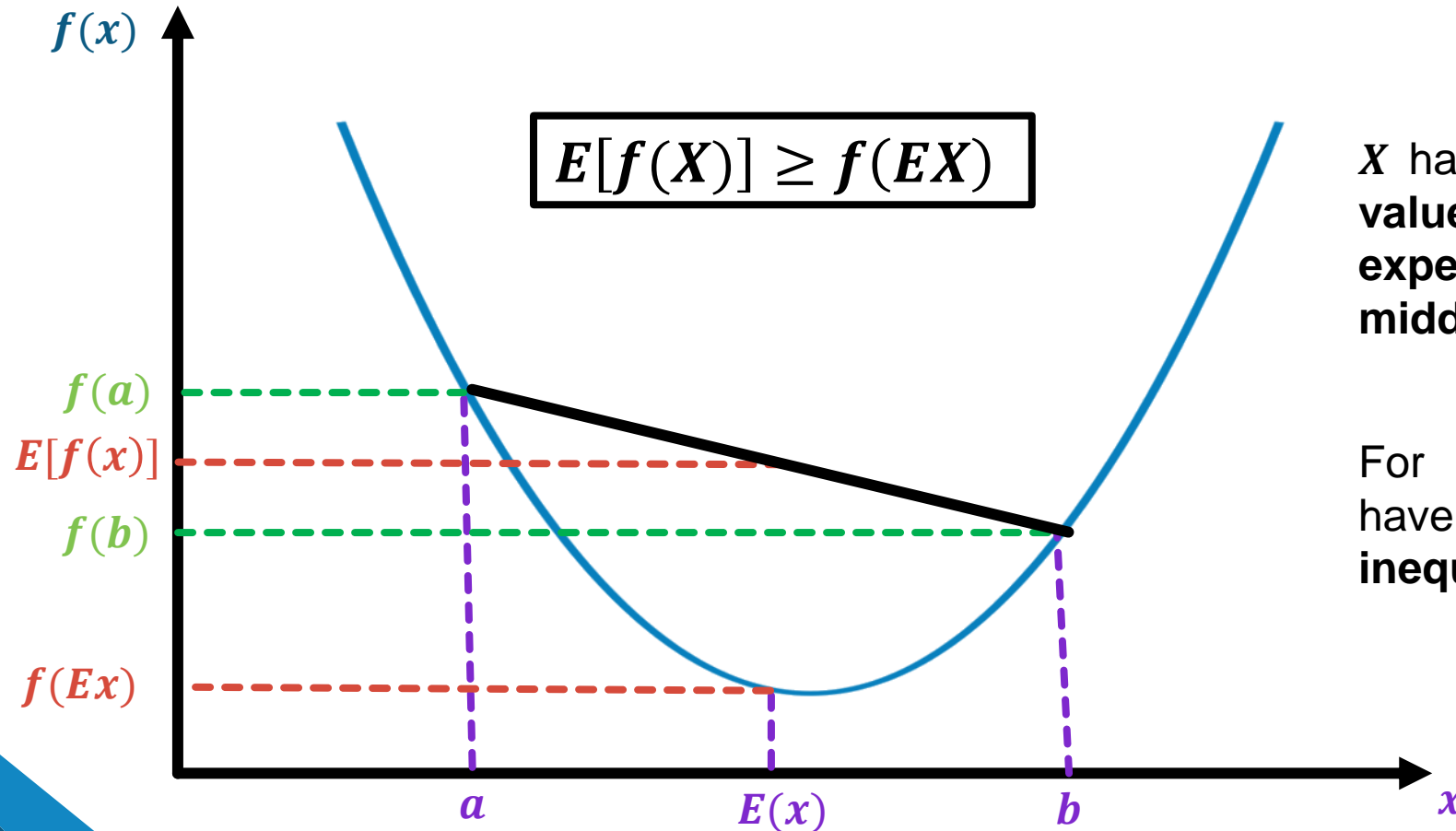
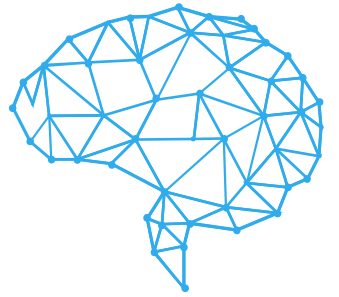
Let f be a **convex function** (if $f''(x) \geq 0; \forall x \in R$ or its **Hessian** is **positive semi-definite** for **vector inputs**) whose domain is the set of real numbers and let X be a **random variable**. Then:

$$E[f(X)] \geq f(EX)$$

If f is **strictly convex** ($f''(x) > 0$), then $E[f(X)] = f(EX)$ holds **true if and only if** $X = E[X]$ with **probability 1** (the **expected value** does **not change**).

EXPECTATION MAXIMIZATION

JENSEN'S INEQUALITY

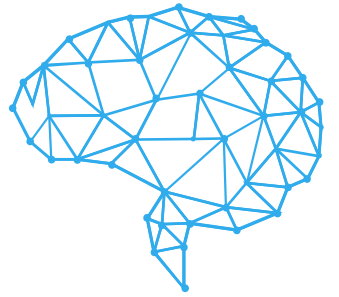


X has **50%** chance of **having value a** or **b** , thus its **expected value** is at the middle of both possibilities.

For **concave functions** we have the **direction** of the **inequality** will be **reversed**.

EXPECTATION MAXIMIZATION

E M A L G O R I T H M



Suppose we have an **estimation problem** in which we have a **training set** $\{x^{(1)}, \dots, x^{(m)}\}$ consisting of m **independent examples**.

The **objective** will be to **fit** the **parameters** of a **model** $p(x, z; w)$ to the **data**, where the **likelihood** is given by:

$$l(w) = \sum_{i=1}^m \log(p(x^{(i)}; w))$$

$$l(w) = \sum_{i=1}^m \log \sum_z p(x^{(i)}, z^{(i)}; w)$$

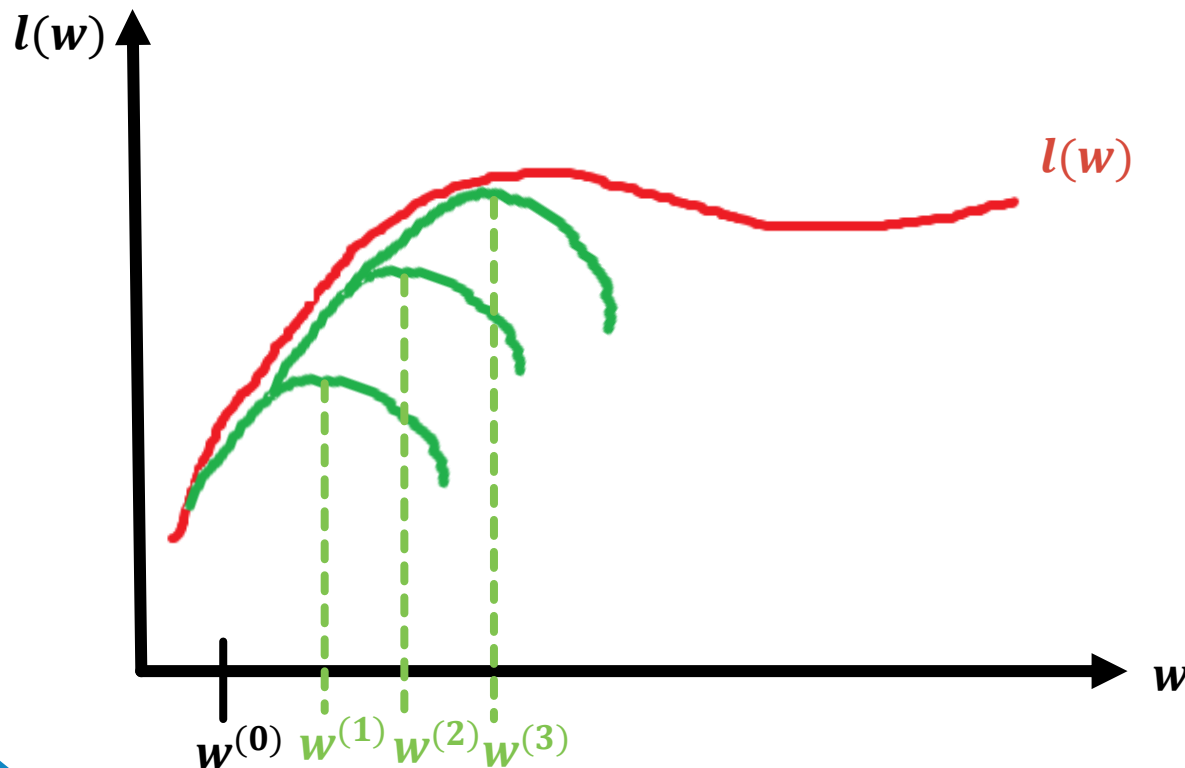
We want to **find** the **maximum likelihood estimates** of the **parameters** w .

EXPECTATION MAXIMIZATION

EM ALGORITHM



Again, **finding** the **parameters** w with MLE is **not** an **easy task** because we **don't know** the latent $z^{(i)}$'s. The **EM algorithm** will help us **overcome** the **problem**.



- Repeatedly **construct** a **lower-bound** on l (**E-step**).
- **Optimize** that lower-bound (**M-step**).

EXPECTATION MAXIMIZATION

E M A L G O R I T H M



Thus **we have** the following:

$$l(\mathbf{w}) = \sum_{i=1}^m \log \sum_{\mathbf{z}} p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)}; \mathbf{w})$$

We will **build** a **probability distribution** Q_i over the latent variables $\mathbf{z}^{(i)}$, where $Q_i(\mathbf{z}^{(i)})$ and $\sum Q_i(\mathbf{z}^{(i)}) = 1$.

$$l(\mathbf{w}) = \sum_{i=1}^m \log \sum_{\mathbf{z}} \frac{Q_i(\mathbf{z}^{(i)}) p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)}; \mathbf{w})}{Q_i(\mathbf{z}^{(i)})}$$

$$l(\mathbf{w}) = \sum_{i=1}^m \log \sum_{\mathbf{z}} E_{\mathbf{z}^{(i)} \sim Q_i} \left[\frac{p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)}; \mathbf{w})}{Q_i(\mathbf{z}^{(i)})} \right]$$

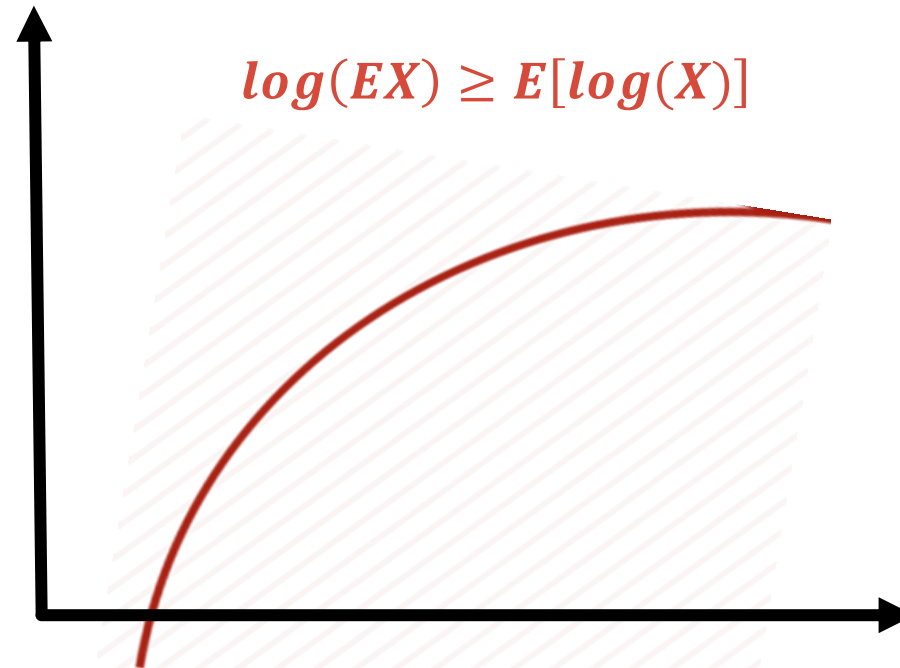
EXPECTATION MAXIMIZATION

EM ALGORITHM



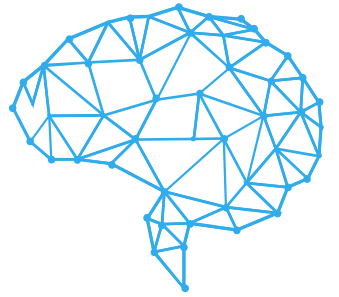
Now, we can see that the **log function** is **concave**, which allow us to **define** the **following**:

$$\log \sum_z E_{z^{(i)} \sim Q_i} \left[\frac{p(x^{(i)}, z^{(i)}; w)}{Q_i(z^{(i)})} \right] \geq \sum_z E_{z^{(i)} \sim Q_i} \left[\log \frac{p(x^{(i)}, z^{(i)}; w)}{Q_i(z^{(i)})} \right]$$



EXPECTATION MAXIMIZATION

E M A L G O R I T H M



Expanding out we have:

$$\sum_{i=1}^m \log \sum_z E_{z^{(i)} \sim Q_i} \left[\frac{p(x^{(i)}, z^{(i)}; w)}{Q_i(z^{(i)})} \right] \geq \sum_{i=1}^m \sum_z E_{z^{(i)} \sim Q_i} \left[\log \frac{p(x^{(i)}, z^{(i)}; w)}{Q_i(z^{(i)})} \right]$$

$$\sum_{i=1}^m \log \sum_z Q_i(z^{(i)}) \left[\frac{p(x^{(i)}, z^{(i)}; w)}{Q_i(z^{(i)})} \right] \geq \sum_{i=1}^m \sum_z Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; w)}{Q_i(z^{(i)})}$$

$$l(w) \geq \sum_{i=1}^m \sum_z Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; w)}{Q_i(z^{(i)})}$$

We observe that we **have a lower bound** over the **likelihood** $l(w)$.

EXPECTATION MAXIMIZATION

EM ALGORITHM

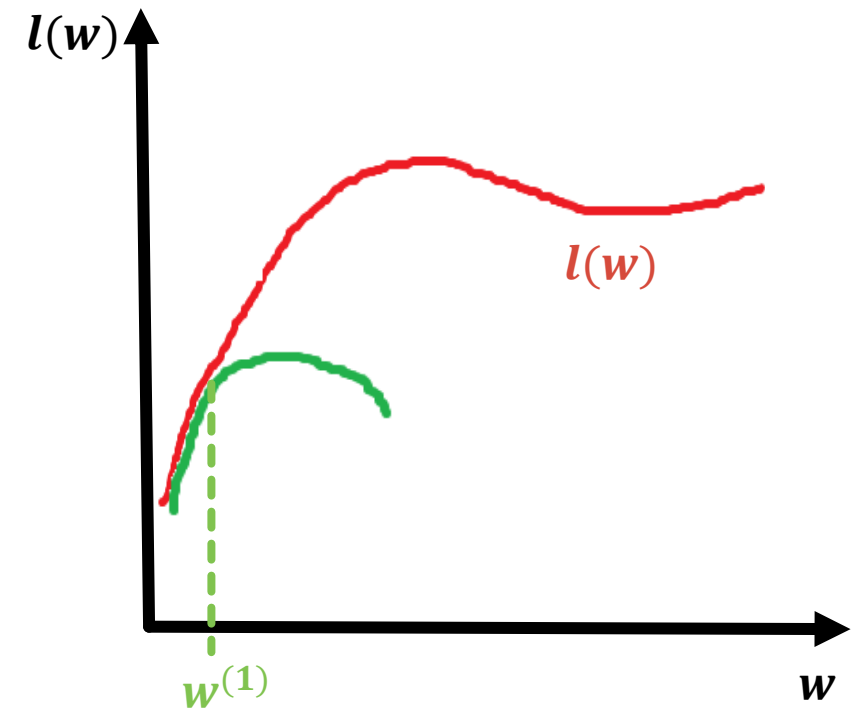


What we **want** is the **inequality** to **turn into** an **equality** for a current value of w . Thus, when we **optimize** the **lower bound**, we are **also optimizing** the **true** $l(w)$.

The **objective** is to find probability distribution Q_i that will **transform** the **inequality** to **equality**.

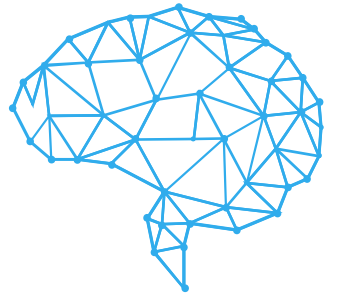
The only way of doing this is to take the **expectation** of a **constant value** (remembering Jensen's inequality).

$$\frac{p(x^{(i)}, z^{(i)}; w)}{Q_i(z^{(i)})} = c$$



EXPECTATION MAXIMIZATION

E M A L G O R I T H M



Elaborating further in the **equation**, we have:

$$Q_i(z^{(i)}) = \frac{p(x^{(i)}, z^{(i)}; w)}{c}$$

$$Q_i(z^{(i)}) \propto p(x^{(i)}, z^{(i)}; w)$$

Since we **know** that $\sum Q_i(z^{(i)}) = 1$

$$1 = \frac{p(x^{(i)}, z^{(i)}; w)}{Q_i(z^{(i)}) c}$$

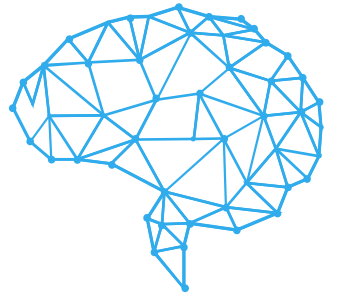
$$\frac{1}{\sum Q_i(z^{(i)})} = \frac{p(x^{(i)}, z^{(i)}; w)}{Q_i(z^{(i)}) c}$$

$$\frac{Q_i(z^{(i)})}{\sum Q_i(z^{(i)})} = \frac{p(x^{(i)}, z^{(i)}; w)}{c}$$

$$c = \sum_z p(x^{(i)}, z^{(i)}; w)$$

EXPECTATION MAXIMIZATION

E M A L G O R I T H M



In **conclusion** we have that:

$$Q_i(z^{(i)}) = \frac{p(x^{(i)}, z^{(i)}; w)}{\sum_z p(x^{(i)}, z^{(i)}; w)}$$

$$Q_i(z^{(i)}) = \frac{p(x^{(i)}, z^{(i)}; w)}{p(x^{(i)}; w)}$$

$$Q_i(z^{(i)}) = p(z^{(i)} / x^{(i)}; w)$$

The distribution $Q_i(z^{(i)})$ is just the **posterior distribution** of the **latent random variables** $z^{(i)}$'s **given** we have **observed** the **data**.

EXPECTATION MAXIMIZATION

E M A L G O R I T H M



The **EM algorithm** therefore is as follows:

Repeat until **convergence** {

1. **(E-step)**: for each i , set the lower bound

$$Q_i(\mathbf{z}^{(i)}) := p(\mathbf{z}^{(i)} / \mathbf{x}^{(i)}; \mathbf{w})$$

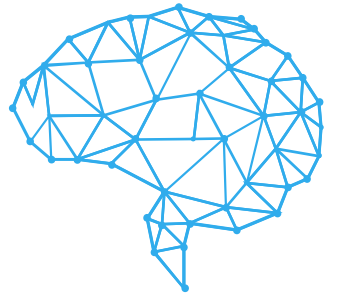
2. **(M-step)**: optimize the lower bound

$$\mathbf{w} := \arg \max_{\mathbf{w}} \sum_{i=1}^m \sum_{\mathbf{z}} Q_i(\mathbf{z}^{(i)}) \log \frac{p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)}; \mathbf{w})}{Q_i(\mathbf{z}^{(i)})}$$

}

EXPECTATION MAXIMIZATION

E M A L G O R I T H M



We can **define**:

$$J(\mathbf{w}, Q) = Q_i(\mathbf{z}^{(i)}) \log \frac{p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)}; \mathbf{w})}{Q_i(\mathbf{z}^{(i)})}$$

$$l(\mathbf{w}) \leq J(\mathbf{w}, Q)$$

Thus, the **EM algorithm** can be viewed as **coordinate ascent** on J , in which the **E-step maximizes** with respect to Q and the **M-step maximizes** it with respect to \mathbf{w} .

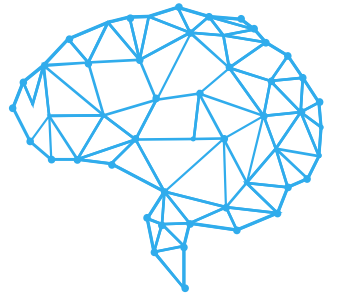


AI

**EM and MIXTURE OF
GAUSSIANS**

EXPECTATION MAXIMIZATION

EM AND MIXTURE OF GAUSSIANS



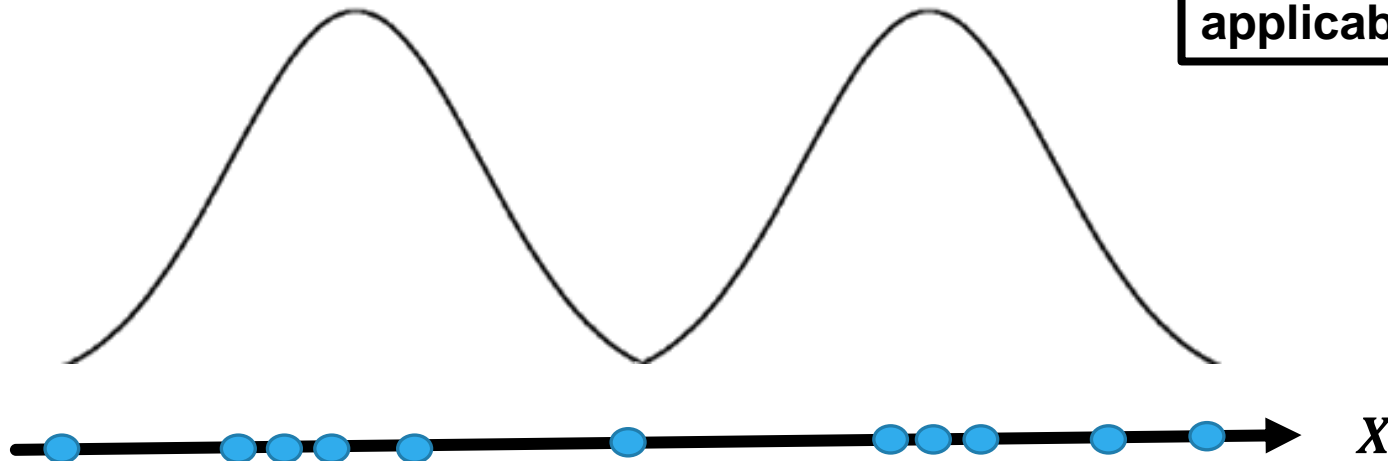
As we have seen, we can **estimate** the **probability density distribution** of a set of **data points** $\{x^{(1)}, \dots, x^{(m)}\}$ using a **mixture of Gaussians**.

$$p(x^{(i)}, z^{(i)}) = p(x^{(i)} / z^{(i)}) p(z^{(i)})$$

$$z^{(i)} \sim \text{Multinomial}(\phi)$$

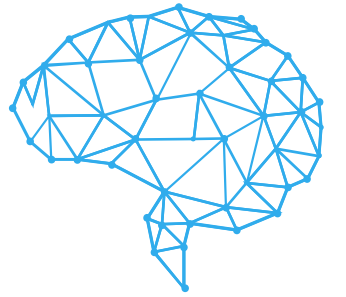
$$x^{(i)} / z^{(i)} = j \sim N(\mu_j, \Sigma_j)$$

NOTE: the mixture of Gaussians model is applicable when $m \gg n$.



EXPECTATION MAXIMIZATION

EM AND MIXTURE OF GAUSSIANS



E-STEP:

Applying the **EM algorithm** to the problem of the **mixture of Gaussians** we can perform the **E step** by getting the **posterior** of the **latent variables**:

$$Q_i(\mathbf{z}^{(i)}) := p(\mathbf{z}^{(i)} / \mathbf{x}^{(i)}; \mathbf{w})$$

$$w^{(i)}_j = Q_i(\mathbf{z}^{(i)} = j) = p(\mathbf{z}^{(i)} = j / \mathbf{x}^{(i)}; \phi, \mu, \Sigma)$$

Where $Q_i(\mathbf{z}^{(i)} = j)$ denotes the **probability** of $\mathbf{z}^{(i)}$ **taking the value j** under the **distribution Q_i** .

EXPECTATION MAXIMIZATION

EM AND MIXTURE OF GAUSSIANS



E-STEP:

Expanding the formula of the **E** step by using Bayes' rule:

$$w^{(i)}_j = Q_i(z^{(i)}) := p(z^{(i)} / x^{(i)}; w)$$

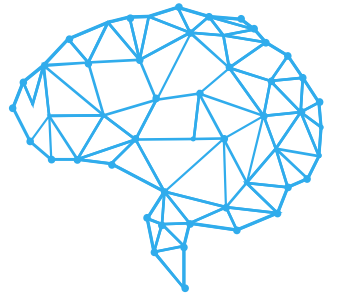
$$Q_i(z^{(i)} = j) = \frac{p(x^{(i)} / z^{(i)} = j) P(z^{(i)} = j)}{\sum_k p(x^{(i)} / z^{(i)} = k) P(z^{(i)} = k)}$$

We know that $p(x^{(i)} / z^{(i)} = j) \sim \text{Gaussian}$ and $P(z^{(i)} = j) \sim \text{Multinomial}$

$$Q_i(z^{(i)} = j) = \frac{\frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)} \phi_j}{\sum_k \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(x^{(i)} - \mu_k)^T \Sigma_k^{-1} (x^{(i)} - \mu_k)} \phi_k}$$

EXPECTATION MAXIMIZATION

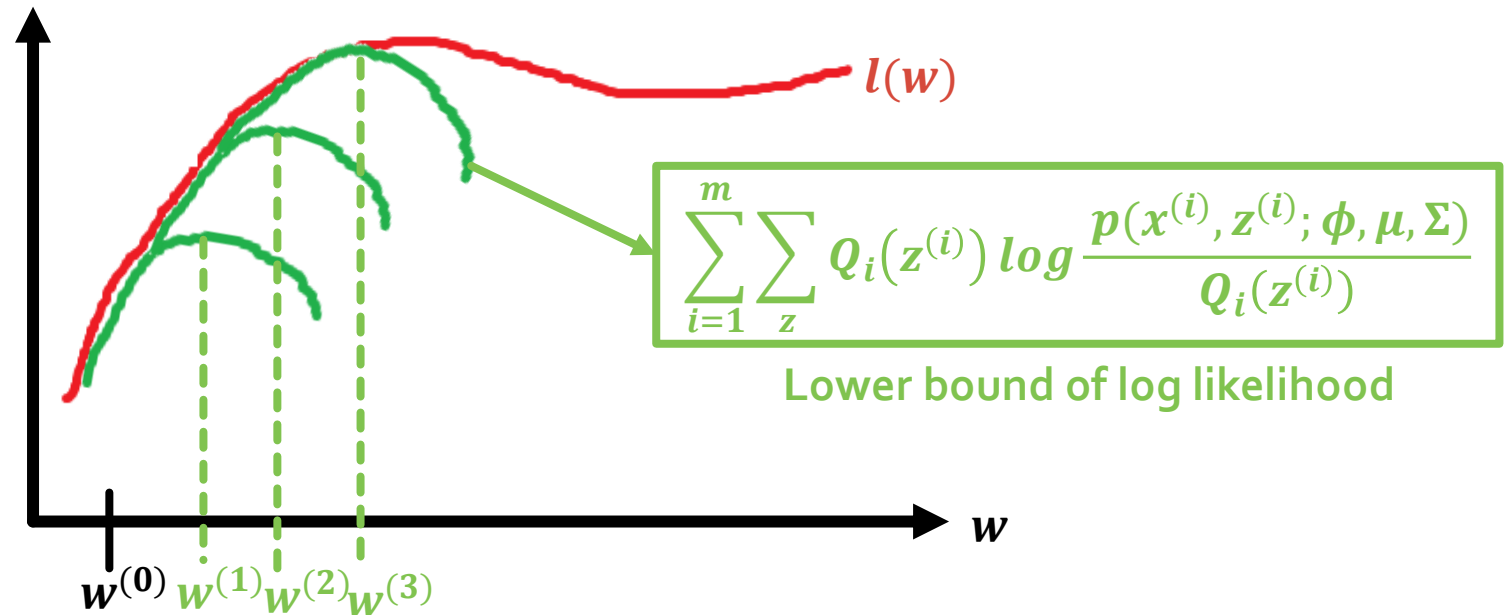
EM AND MIXTURE OF GAUSSIANS



M-STEP:

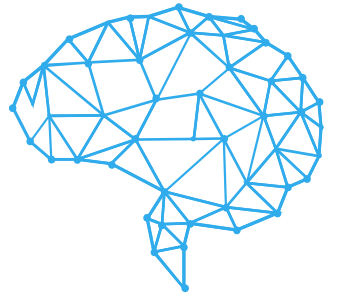
For the **M step** we need to **maximize** with respect to the parameters ϕ, μ, Σ the following:

$$w := \arg \max_w \sum_{i=1}^m \sum_z Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \phi, \mu, \Sigma)}{Q_i(z^{(i)})}$$



EXPECTATION MAXIMIZATION

EM AND MIXTURE OF GAUSSIANS



M-STEP: maximize for μ

Expanding terms out we have:

$$\sum_{i=1}^m \sum_z Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \phi, \mu, \Sigma)}{Q_i(z^{(i)})} = \sum_{i=1}^m \sum_{j=1}^k w^{(i)}_j \log \frac{\frac{1}{(2\pi)^{n/2} |\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2}(x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)} \phi_j}{w^{(i)}_j}$$

Simplifying:

$$= \sum_{i=1}^m \sum_{j=1}^k w^{(i)}_j \left[\log \left(\frac{1}{w^{(i)}_j (2\pi)^{n/2} |\Sigma|^{\frac{1}{2}}} \right) + \log \left(e^{-\frac{1}{2}(x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)} \phi_j \right) \right]$$

EXPECTATION MAXIMIZATION

EM AND MIXTURE OF GAUSSIANS



M-STEP: maximize for μ

Taking the derivative with respect to μ_l we have:

$$\begin{aligned} \nabla_{\mu_l} \sum_{i=1}^m \sum_{j=1}^k w^{(i)}_j & \left[\log \left(\frac{1}{w^{(i)}_j (2\pi)^{n/2} |\Sigma|^{1/2}} \right) + \log \left(e^{-\frac{1}{2}(x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)} \phi_j \right) \right] \\ &= \nabla_{\mu_l} \sum_{i=1}^m \sum_{j=1}^k w^{(i)}_j \left[-\frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j) \right] \\ &= -\frac{1}{2} \sum_{i=1}^m w^{(i)}_l \nabla_{\mu_l} \left[-2\mu_l^T \Sigma_l^{-1} x^{(i)} + \mu_l^T \Sigma_l^{-1} \mu_l \right] \end{aligned}$$

EXPECTATION MAXIMIZATION

EM AND MIXTURE OF GAUSSIANS



M-STEP: maximize for μ

$$= -\frac{1}{2} \sum_{i=1}^m w^{(i)}_l \nabla_{\mu_l} \left[-2\mu_l^T \Sigma_l^{-1} x^{(i)} + \mu_l^T \Sigma_l^{-1} \mu_l \right]$$

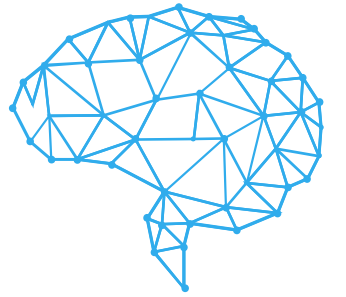
$$= -\frac{1}{2} \sum_{i=1}^m w^{(i)}_l \left[-2\Sigma_l^{-1} x^{(i)} + 2\mu_l^T \Sigma_l^{-1} \right]$$

$$\sum_{i=1}^m w^{(i)}_l \left[\Sigma_l^{-1} x^{(i)} - \mu_l^T \Sigma_l^{-1} \right] = 0$$

$$\sum_{i=1}^m w^{(i)}_l \Sigma_l^{-1} x^{(i)} - \sum_{i=1}^m w^{(i)}_l \mu_l^T \Sigma_l^{-1} = 0$$

EXPECTATION MAXIMIZATION

EM AND MIXTURE OF GAUSSIANS



M-STEP: maximize for μ

$$\sum_{i=1}^m w^{(i)}_l \mu_l^T \Sigma_l^{-1} = \sum_{i=1}^m w^{(i)}_l \Sigma_l^{-1} x^{(i)}$$

$$\mu_l^T \sum_{i=1}^m w^{(i)}_l = \Sigma_l^{-1} \sum_{i=1}^m w^{(i)}_l x^{(i)}$$

$$\mu_l^T = \Sigma_l \Sigma_l^{-1} \frac{\sum_{i=1}^m w^{(i)}_l x^{(i)}}{\sum_{i=1}^m w^{(i)}_l}$$

$$\mu_l^T = \frac{\sum_{i=1}^m w^{(i)}_l x^{(i)}}{\sum_{i=1}^m w^{(i)}_l}$$

EXPECTATION MAXIMIZATION

EM AND MIXTURE OF GAUSSIANS



M-STEP: maximize for ϕ

$$\begin{aligned} & \sum_{i=1}^m \sum_{j=1}^k w^{(i)}_j \log \frac{\frac{1}{(2\pi)^{n/2} |\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2}(x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)} \phi_j}{w^{(i)}_j} \\ &= \sum_{i=1}^m \sum_{j=1}^k w^{(i)}_j \left[\log \left(\frac{1}{w^{(i)}_j (2\pi)^{n/2} |\Sigma|^{\frac{1}{2}}} \right) + \log \left(e^{-\frac{1}{2}(x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)} \phi_j \right) \right] \\ &= \sum_{i=1}^m \sum_{j=1}^k w^{(i)}_j \left[\frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j) + \log(\phi_j) \right] \\ &= \sum_{i=1}^m \sum_{j=1}^k w^{(i)}_j \log(\phi_j) \end{aligned}$$

EXPECTATION MAXIMIZATION

EM AND MIXTURE OF GAUSSIANS



M-STEP: maximize for ϕ

Because $\phi \sim \text{Multinomial}$ we have an additional constraint that $\sum_{j=1}^k \phi_j = 1$.
Thus we need to construct our Lagrangian:

$$L(\phi) = \sum_{i=1}^m \sum_{j=1}^k w^{(i)}_j \log(\phi_j) + \beta \left(\sum_{j=1}^k \phi_j - 1 \right)$$

$$\frac{\partial L(\phi)}{\partial \phi} = \sum_{i=1}^m \frac{w^{(i)}_j}{\phi_j} + \beta = 0$$

$$\phi_j = \sum_{i=1}^m \frac{w^{(i)}_j}{-\beta}$$

EXPECTATION MAXIMIZATION

EM AND MIXTURE OF GAUSSIANS



M-STEP: maximize for ϕ

Using the constraint $\sum_{j=1}^k \phi_j = 1$. we have:

$$\phi_j = \sum_{i=1}^m \frac{w^{(i)}_j}{-\beta}$$

$$\sum_{j=1}^k \phi_j = \sum_{j=1}^k \sum_{i=1}^m \frac{w^{(i)}_j}{-\beta}$$

$$1 = \sum_{j=1}^k \sum_{i=1}^m \frac{w^{(i)}_j}{-\beta}$$

$$1 = \sum_{i=1}^m \frac{1}{-\beta} \rightarrow -\beta = m$$

EXPECTATION MAXIMIZATION

EM AND MIXTURE OF GAUSSIANS



M-STEP: maximize for ϕ

Finally we have that:

$$\phi_j = \sum_{i=1}^m \frac{w^{(i)}_j}{-\beta}$$

$$\phi_j = \frac{1}{m} \sum_{i=1}^m w^{(i)}_j$$

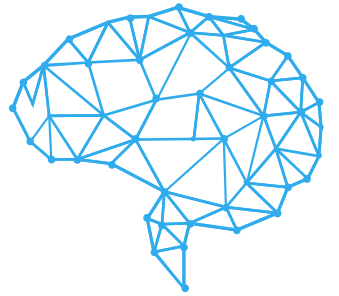


AI

EM and NAÏVE BAYES

EXPECTATION MAXIMIZATION

EM AND NAÏVE BAYES



As we have seen, **Naïve Bayes classifier** runs for **input** training **examples** that take **discrete values**.

Therefore, given a training set $\{x^{(1)}, \dots, x^{(m)}\}$ where $x^{(i)} \in \{0, 1\}^n$ and $x^{(i)}_j = 1\{\text{word } j \text{ appears in document } i\}$.

Suppose we want to find **two clusters** $z^{(i)} = \{0, 1\}$ (“spam” or “not spam”).

The **assumptions** are the following:

- $z^{(i)} \sim \text{Bernoulli}(\phi) \rightarrow$ probability that document i comes from cluster 1 or 2.
- $x^{(i)} \sim \text{Multinomial}$.
- $P(x^{(i)} / z^{(i)}) = \prod_{i=1}^n P(x_j^{(i)} / z^{(i)})$ the appearance of words is independent from each other.
- $P(x_j^{(i)} = 1 / z^{(i)} = 0) = \phi_{j/z=0}$

EXPECTATION MAXIMIZATION

EM AND NAÏVE BAYES



Computing the **EM algorithm** we find that:

- **E-Step:** find the posterior distribution (estimate where the document comes from).

$$w^{(i)}_j = Q_i(z^{(i)}):= p(z^{(i)} = 1/x^{(i)}; \phi_{j/z}, \phi)$$

- **M-Step:** maximize the lower bound .

$w^{(i)}$ captures uncertainty of cluster membership

$$\phi_{j/z=1} = \frac{\sum_{i=1}^m w^{(i)} \mathbf{1}\{x_j^{(i)} = 1\}}{\sum_{i=1}^m w^{(i)}} = \frac{\# \text{ times word } j \text{ is in documents that we think are in cluster 1}}{\# \text{ documents (estimated) we think are in cluster 1}}$$

$$\phi_{j/z=0} = \frac{\sum_{i=1}^m (1 - w^{(i)}) \mathbf{1}\{x_j^{(i)} = 1\}}{\sum_{i=1}^m (1 - w^{(i)})}$$

$$\phi_{j/z=0} = \frac{\sum_{i=1}^m w^{(i)}}{m}$$

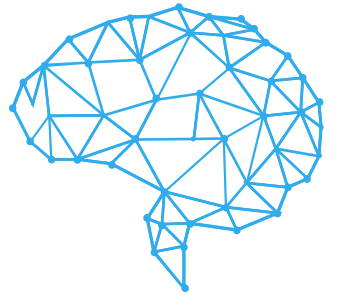


AI

FACTOR ANALYSIS

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



Given a training set $\{x^{(1)}, \dots, x^{(m)}\}$ where $x^{(i)} \in \mathbb{R}^n$ it would be **very difficult** to **model data** with a **mixture of Gaussians** when $n \gg m$.

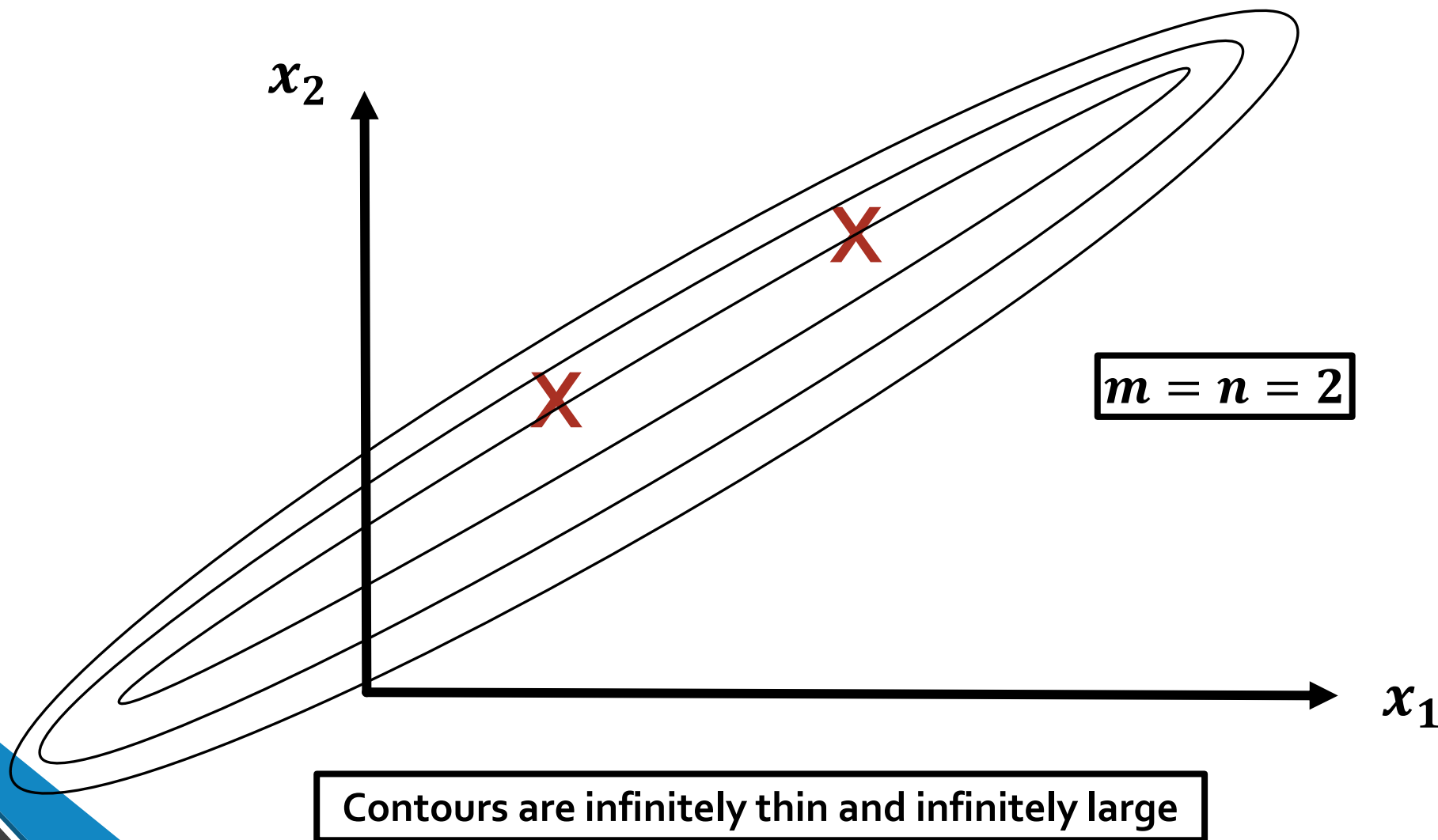
If we **model** the **data** as a **Gaussian**, and **estimate** the **mean** and **covariance** using the usual **maximum likelihood** estimators we would find that the **matrix Σ** is **singular** $\left(\frac{1}{|\Sigma|^{1/2}} = \frac{1}{0}\right)$.

Thus, the **maximum likelihood estimates** of the parameters **result** in a **Gaussian** that places **all** its **probability** in the affine **space spanned** by the **data**, resulting in a **singular covariance matrix**.

$$f_{\mathbf{X}}(x_1, \dots, x_k) = \frac{\exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)}{\sqrt{(2\pi)^k |\boldsymbol{\Sigma}|}}$$

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



Unless $m \gg n$, by a **reasonable amount**, the **maximum likelihood estimates** of the **mean** and **covariance** may be quite **poor**.

Factor analysis will be used when $m \approx n$ or when $n \gg m$. We are going to look at **two restrictions** on Σ that will **allow us** to **fit** Σ with **small amounts of data**.

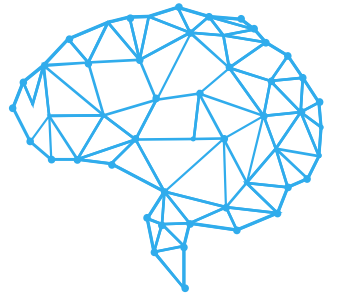
1. **Constrain Σ to be diagonal.**

$$\Sigma = \begin{pmatrix} \sigma_1^2 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_n^2 \end{pmatrix}$$

2. **The diagonal entries must be equal: $\Sigma = \sigma^2 I$.**

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



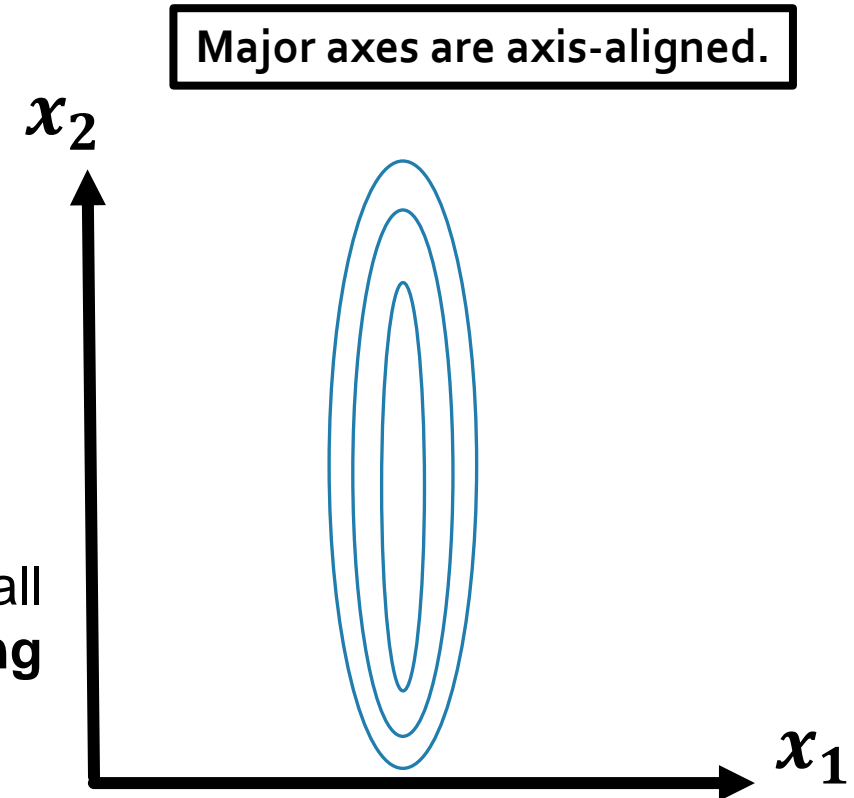
1. CONSTRAIN Σ TO BE DIAGONAL

The **maximum likelihood** estimate **would be**:

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j) (x_j^{(i)} - \mu_j)^T$$

$$\Sigma_{jj} = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$

The **main problem** is that you are **removing** all **correlations** between **features**. → We are **assuming** that the **features** are **independent** between them.



EXPECTATION MAXIMIZATION

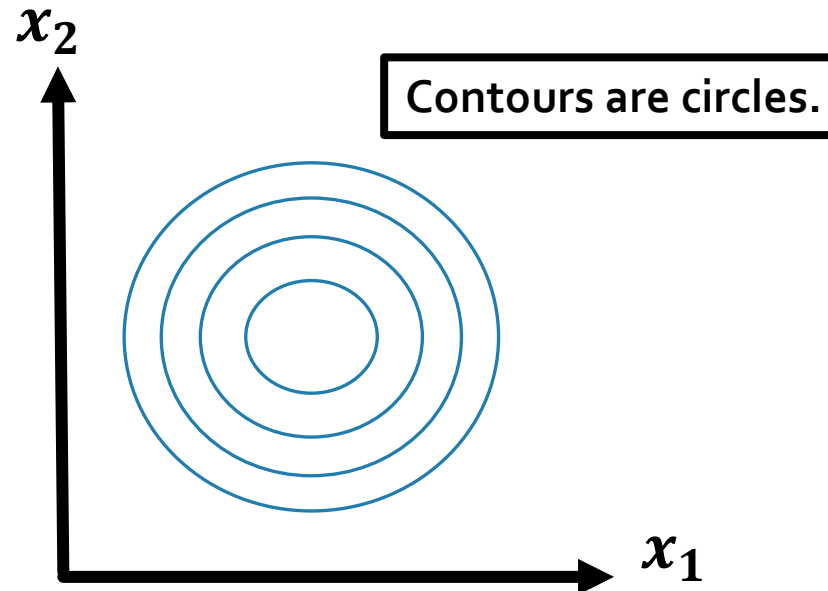
EM AND FACTOR ANALYSIS



2. THE DIAGONAL ENTRIES MUST BE EQUAL: $\Sigma = \sigma^2 I$.

The **maximum likelihood** estimate **would be**:

$$\sigma^2 = \frac{1}{mn} \sum_{j=1}^n \sum_{i=1}^m \left(x_j^{(i)} - \mu_j \right)^2$$



EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



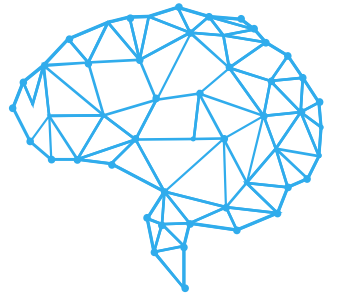
If we are **fitting** a **full, unconstrained**, covariance matrix Σ to data, it is necessary that $m \geq n + 1$ for the **maximum likelihood** estimate of Σ **not** to be **singular**.

Under either of the **two restrictions presented**, we may **obtain** a **non-singular** Σ when $n \geq 2$.

The **problem** is that in **many occasions** we **want** to be able to **capture** some **interesting correlation structure** in the data.

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



MARGINALS AND CONDITIONALS OF GAUSSIANS

Before talking about the **factor analysis model**, we will discuss how to **find conditional and marginal distributions** of random variables with a **joint multivariate Gaussian distribution**.

Suppose we have a **vector-valued random variable**:

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

Where $x_1 \in \mathbb{R}^r$, $x_2 \in \mathbb{R}^s$, and $x \in \mathbb{R}^{r+s}$. **Suppose** $x \sim N(\mu, \Sigma)$, where:

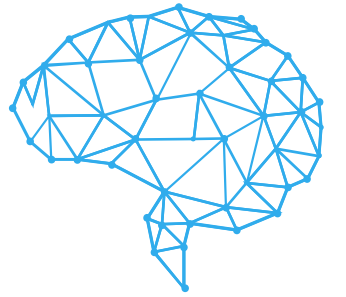
$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} \quad \Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$$

x_1 and x_2 are jointly distributed multivariate Gaussian

$\mu_1 \in \mathbb{R}^r$, $\mu_2 \in \mathbb{R}^s$, $\Sigma_{11} \in \mathbb{R}^{r \times r}$, $\Sigma_{12} \in \mathbb{R}^{r \times s}$, $\Sigma_{21} \in \mathbb{R}^{s \times r}$ and $\Sigma_{22} \in \mathbb{R}^{s \times s}$.

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



MARGINALS AND CONDITIONALS OF GAUSSIANS

To **obtain** the **marginal distribution** of x_1 , we **can see** that:

$$E[x_1] = \mu_1$$

$$Cov[x_1] = E[(x_1 - \mu_1)(x_1 - \mu_1)] = \Sigma_{11}$$

To **demonstrate** the **previous statement**, we can see the **joint covariance** of x_1 and x_2 :

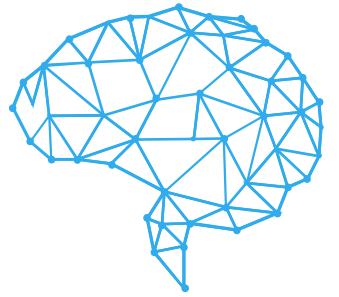
$$Cov[x] = \Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} = E[(x - \mu)(x - \mu)^T]$$

$$Cov[x] = E \left[\begin{pmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{pmatrix} \begin{pmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{pmatrix}^T \right] = E \left[\begin{array}{cc} (x_1 - \mu_1)(x_1 - \mu_1)^T & (x_1 - \mu_1)(x_2 - \mu_2)^T \\ (x_2 - \mu_2)(x_1 - \mu_1)^T & (x_2 - \mu_2)(x_2 - \mu_2)^T \end{array} \right]$$

Therefore the **marginal distribution** of x_1 is $N(\mu_1, \Sigma_{11})$.

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



MARGINALS AND CONDITIONALS OF GAUSSIANS

We can also **obtain** the **marginal conditional distribution** of x_1/x_2 :

$$P(x_1/x_2) = \frac{P(x_1, x_2)}{P(x_2)} = \frac{N(\mu, \Sigma)}{N(\mu_2, \Sigma_{22})}$$

Substituting the **formulas** for **both gaussians**, the **joint** and the **marginal** of x_2 , you would **obtain** the **following** (these **computations** are **non-trivial**):

$$x_1/x_2 \sim N(\mu_{1/2}, \Sigma_{1/2})$$

$$\mu_{1/2} = \mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(x_2 - \mu_2)$$

$$\Sigma_{1/2} = \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}$$

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



FACTOR ANALYSIS MODEL

We will **create** a **joint distribution** (x, z) by **assuming** that:

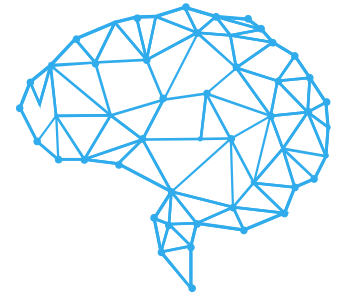
- There is a **latent random variable** $z \sim N(\mathbf{0}, I)$, where $z \in \mathbb{R}^k$ such that $k < m$.
- $x/z \sim N(\mu + \Lambda z, \Psi)$
- $x = \mu + \Lambda z + \varepsilon$, where $\varepsilon \sim N(\mathbf{0}, \Psi)$

Therefore, the **parameters** of the **model** are:

- $\mu \in \mathbb{R}^n$, $\Lambda \in \mathbb{R}^{n \times k}$, $\Psi \in \mathbb{R}^{n \times n}$ and Ψ is diagonal (**usually** $k < n$).

EXPECTATION MAXIMIZATION

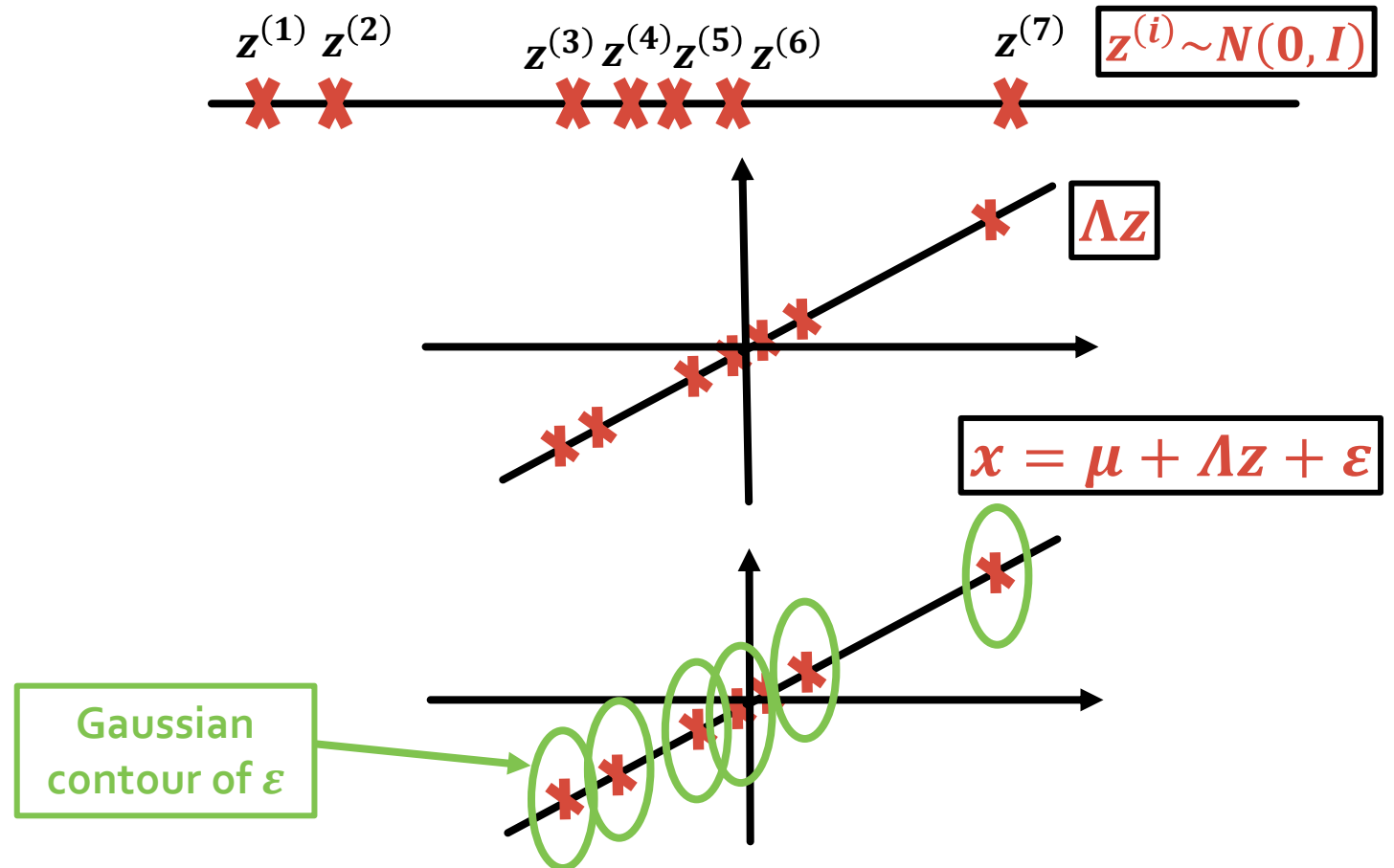
EM AND FACTOR ANALYSIS



FACTOR ANALYSIS MODEL

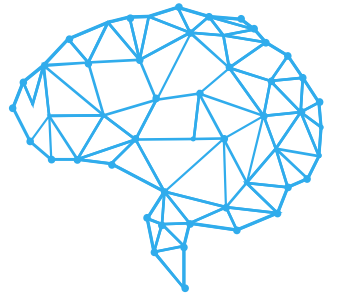
Let us give an **example**.

- $z \in \mathbb{R}^1$
- $x \in \mathbb{R}^2$
- $\Lambda = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$
- $\Psi = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$
- $\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$



EXPECTATION MAXIMIZATION

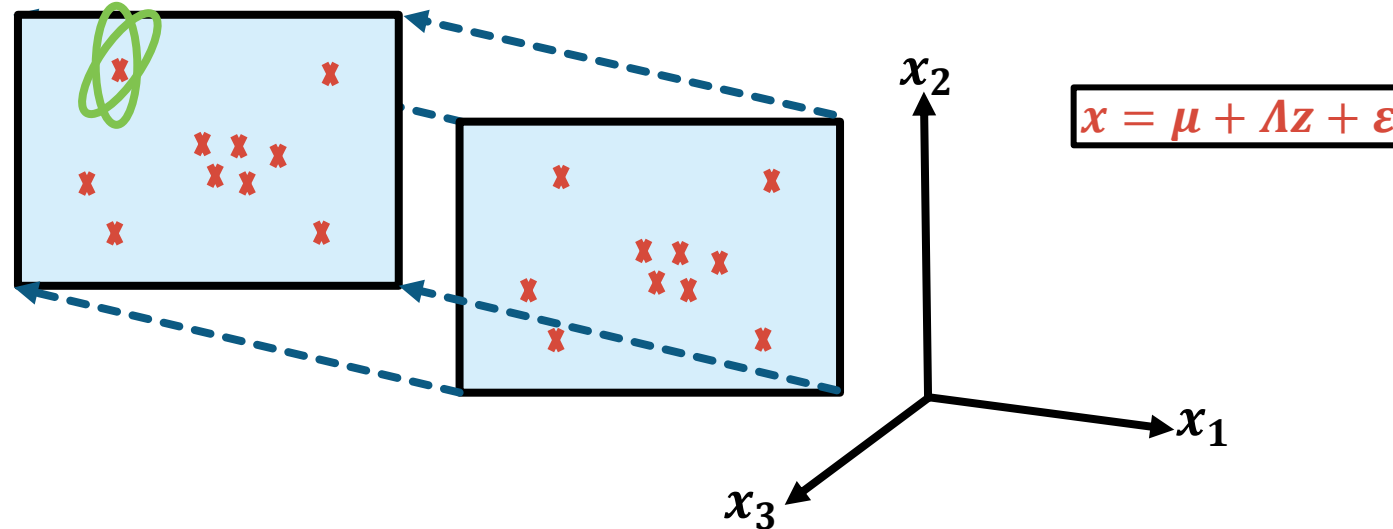
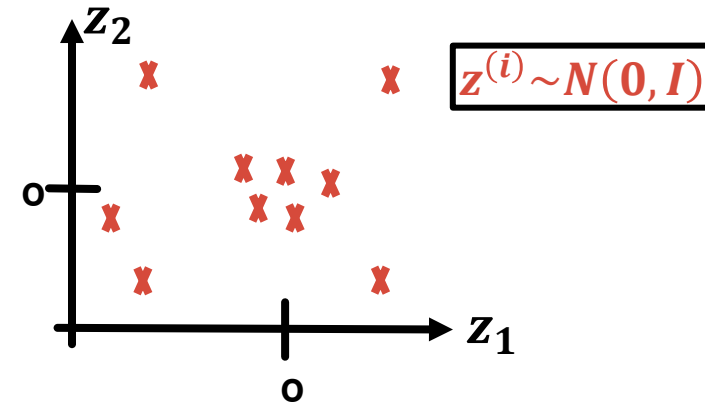
EM AND FACTOR ANALYSIS



FACTOR ANALYSIS MODEL

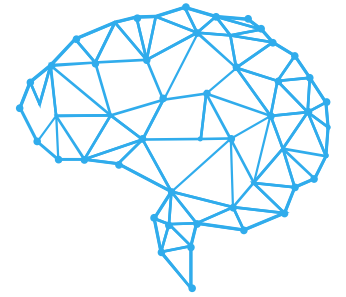
Let us give another **example**.

- $z \in \mathbb{R}^2$
- $x \in \mathbb{R}^3$



EXPECTATION MAXIMIZATION

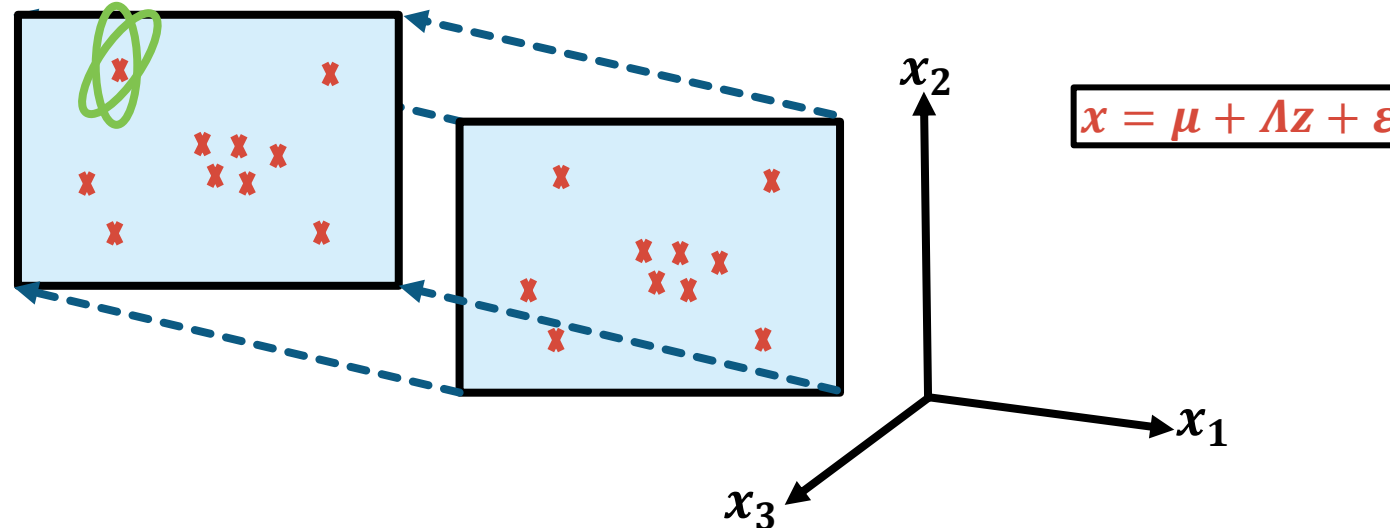
EM AND FACTOR ANALYSIS



FACTOR ANALYSIS MODEL

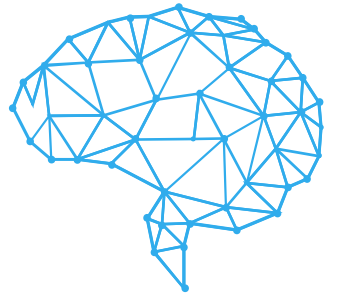
In as **summary**:

- Each **datapoint** $x^{(i)}$ is **generated** by **sampling** a k dimension **multivariate Gaussian** $z^{(i)}$.
- Then, it is **mapped** to a n -dimensional **affine space** of \mathbb{R}^n by **computing** $\mu + \Lambda z^{(i)}$.
- Lastly, $x^{(i)}$ is **generated** by **adding covariance** Ψ **noise** to $\mu + \Lambda z^{(i)}$.



EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



FACTOR ANALYSIS MODEL

We **will work out** exactly what **distribution** our model defines using **marginal** and **conditional distributions**.

Our **random variables** z and x have a **joint Gaussian distribution**:

$$\begin{bmatrix} z \\ x \end{bmatrix} \sim N(\mu_{zx}, \Sigma)$$

We **want to find** μ_{zx} and Σ .

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



FACTOR ANALYSIS MODEL

To find μ_{zx} , let us remember that $E[z] = 0$ because $z \sim N(0, I)$.

Also we can compute the expected value of x :

$$E[x] = E[\mu + \Lambda z + \varepsilon] = E[\mu] + \Lambda E[z] + E[\varepsilon]$$

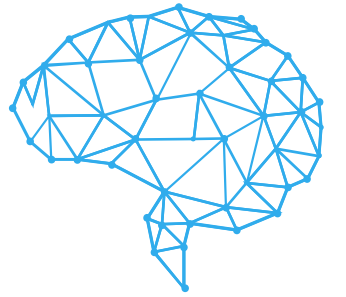
$$E[x] = \mu$$

Therefore we have :

$$\mu_{zx} = \begin{bmatrix} E[z] \\ E[x] \end{bmatrix} = \begin{bmatrix} \vec{0} \\ \mu \end{bmatrix}$$

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



FACTOR ANALYSIS MODEL

To **find** Σ , we need to compute:

- $\Sigma_{zz} = E[(z - E[z])(z - E[z])^T]$
- $\Sigma_{zx} = E[(z - E[z])(x - E[x])^T]$
- $\Sigma_{xx} = E[(x - E[x])(x - E[x])^T]$
- Σ_{xz} (it is **not needed** because Σ_{zx} and Σ_{zx} are **symmetric** $\Sigma_{zx} = \Sigma_{xz}^T$).

Since $z \sim N(0, I)$ we have $\Sigma_{zz} = I$.

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



FACTOR ANALYSIS MODEL

Now we find Σ_{zx} :

$$\Sigma_{zx} = E[(z - E[z])(x - E[x])^T]$$

$$\Sigma_{zx} = E[(z - \vec{0})(x - \mu)^T]$$

$$\Sigma_{zx} = E[(z)(x - \mu)^T]$$

$$\Sigma_{zx} = E[(z)(\mu + \Lambda z + \varepsilon - \mu)^T]$$

$$\Sigma_{zx} = \Lambda^T E[zz^T] + E[z\varepsilon^T]$$

Because z and ε are **independent** $E[z\varepsilon^T] = E[z]E[\varepsilon^T] = \mathbf{0}$. Also we have $E[zz^T] = \text{Cov}(z) = I$:

$$\Sigma_{zx} = \Lambda^T$$

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



FACTOR ANALYSIS MODEL

Now we find Σ_{xx} :

$$\Sigma_{xx} = E[(\mu + \Lambda z + \varepsilon - \mu)(\mu + \Lambda z + \varepsilon - \mu)^T]$$

$$\Sigma_{xx} = E[(\Lambda z + \varepsilon)(\Lambda z + \varepsilon)^T]$$

$$\Sigma_{xx} = E[\Lambda z z^T \Lambda^T + \varepsilon z^T \Lambda^T + \Lambda z \varepsilon^T + \varepsilon \varepsilon^T]$$

$$\Sigma_{xx} = E[\Lambda z z^T \Lambda^T] + E[\varepsilon z^T \Lambda^T] + E[\Lambda z \varepsilon^T] + E[\varepsilon \varepsilon^T]$$

$$\Sigma_{xx} = E[\Lambda z z^T \Lambda^T] + \Lambda^T E[\varepsilon] E[z^T] + \Lambda E[z] E[\varepsilon^T] + E[\varepsilon \varepsilon^T]$$

$$\Sigma_{xx} = \Lambda E[zz^T] \Lambda^T + \mathbf{0} + \mathbf{0} + \text{Cov}(\varepsilon)$$

$$\Sigma_{xx} = \Lambda \Lambda^T + \Psi$$

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



FACTOR ANALYSIS MODEL

With Σ_{xx} , Σ_{zz} , Σ_{zx} , and Σ_{xz} we can now build Σ :

$$\Sigma = \begin{bmatrix} \Sigma_{zz} & \Sigma_{zx} \\ \Sigma_{xz} & \Sigma_{xx} \end{bmatrix} = \begin{bmatrix} I & \Lambda^T \\ \Lambda & \Lambda\Lambda^T + \Psi \end{bmatrix}$$

Therefore, the **joint distribution** of (z, x) is **defined** as:

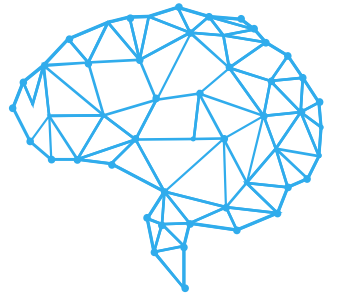
$$\begin{bmatrix} z \\ x \end{bmatrix} \sim N \left(\begin{bmatrix} \vec{0} \\ \mu \end{bmatrix}, \begin{bmatrix} I & \Lambda^T \\ \Lambda & \Lambda\Lambda^T + \Psi \end{bmatrix} \right)$$

And the **marginal distribution** of x would be:

$$x \sim N(\mu, \Lambda\Lambda^T + \Psi)$$

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



FACTOR ANALYSIS MODEL

The **parameters** of our **model** would be μ , Λ , and Ψ .

Therefore, **given** a **training set** $\{x^{(1)}, \dots, x^{(m)}\}$ we would like to **make maximum likelihood estimation** for the **parameters**.

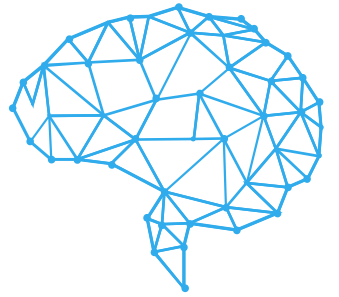
$$l(\mu, \Lambda, \Psi) = \log \prod_{i=1}^m P(x^{(i)}) = \log \prod_{i=1}^m \text{Gaussian}(\mu, \Lambda \Lambda^T + \Psi)$$

The **procedure** would be the **same**, **take derivatives** with **respect** to **each parameter**, **equal to 0** and **solve** for **each parameter**.

But **maximizing** this **formula explicitly** is **hard**, we **have not found** an **algorithm** that **does it** in **closed-form**.

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



EM FOR FACTOR ANALYSIS: E-STEP

We will use the **EM algorithm** instead of **maximum likelihood estimation**.

We need to compute $Q_i(z^{(i)}) = p(z^{(i)}/x^{(i)}; \mu, \Lambda, \Psi)$. We already know how to compute a **conditional probability** for a **Gaussian**.

$$z^{(i)}/x^{(i)} \sim N\left(\mu_{z^{(i)}/x^{(i)}}, \Sigma_{z^{(i)}/x^{(i)}}\right)$$

$$\mu_{1/2} = \mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(x_2 - \mu_2) \rightarrow \mu_{z^{(i)}/x^{(i)}} = \Lambda^T(\Lambda\Lambda^T + \Psi)^{-1}(x^{(i)} - \mu)$$

$$\Sigma_{1/2} = \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21} \rightarrow \Sigma_{z^{(i)}/x^{(i)}} = I - \Lambda^T(\Lambda\Lambda^T + \Psi)^{-1}\Lambda$$

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



EM FOR FACTOR ANALYSIS: E-STEP

Therefore, we have that distribution $Q_i(\mathbf{z}^{(i)})$ is defined as:

$$\mathbf{z}^{(i)} / \mathbf{x}^{(i)} \sim N\left(\boldsymbol{\mu}_{\mathbf{z}^{(i)} / \mathbf{x}^{(i)}}, \boldsymbol{\Sigma}_{\mathbf{z}^{(i)} / \mathbf{x}^{(i)}}\right)$$

$$Q_i(z^{(i)}) = \frac{1}{(2\pi)^{k/2} |\boldsymbol{\Sigma}_{\mathbf{z}^{(i)} | \mathbf{x}^{(i)}}|^{1/2}} \exp\left(-\frac{1}{2} (z^{(i)} - \mu_{\mathbf{z}^{(i)} | \mathbf{x}^{(i)}})^T \boldsymbol{\Sigma}_{\mathbf{z}^{(i)} | \mathbf{x}^{(i)}}^{-1} (z^{(i)} - \mu_{\mathbf{z}^{(i)} | \mathbf{x}^{(i)}})\right)$$

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



EM FOR FACTOR ANALYSIS: M-STEP

Remember that the **joint log likelihood** was **expressed** as **follows**.

$$l(w) = \sum_{i=1}^m \sum_z Q_i(z^{(i)}) \log \left(\frac{p(x^{(i)}, z^{(i)}; \mu, \Lambda, \Psi)}{Q_i(z^{(i)})} \right)$$

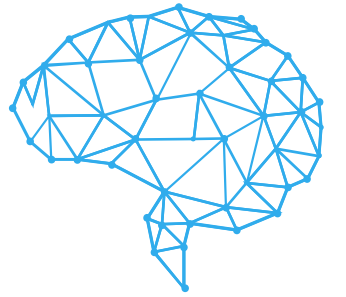
Because now, z is a **continuous variable**, we have that the **joint log likelihood** will be **defined** as:

$$l(w) = \sum_{i=1}^m \log \int_{z^{(i)}} Q_i(z^{(i)}) \log \left(\frac{p(x^{(i)}, z^{(i)}; \mu, \Lambda, \Psi)}{Q_i(z^{(i)})} \right) dz^{(i)}$$

We **want to maximize** this **expression**.

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



EM FOR FACTOR ANALYSIS: M-STEP

We will **expand** the **logarithm**.

$$l(\mathbf{w}) = \sum_{i=1}^m \int_{\mathbf{z}^{(i)}} Q_i(\mathbf{z}^{(i)}) \log \left(\frac{p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)}; \boldsymbol{\mu}, \boldsymbol{\Lambda}, \boldsymbol{\Psi})}{Q_i(\mathbf{z}^{(i)})} \right) d\mathbf{z}^{(i)}$$

$$l(\mathbf{w}) = \sum_{i=1}^m \int_{\mathbf{z}^{(i)}} Q_i(\mathbf{z}^{(i)}) \left[\log(p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)}; \boldsymbol{\mu}, \boldsymbol{\Lambda}, \boldsymbol{\Psi})) - \log(Q_i(\mathbf{z}^{(i)})) \right] d\mathbf{z}^{(i)}$$

We **know** that $p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)}) = p(\mathbf{x}^{(i)} / \mathbf{z}^{(i)}) p(\mathbf{z}^{(i)})$:

$$l(\mathbf{w}) = \sum_{i=1}^m \int_{\mathbf{z}^{(i)}} Q_i(\mathbf{z}^{(i)}) \left[\log(p(\mathbf{x}^{(i)} / \mathbf{z}^{(i)}; \boldsymbol{\mu}, \boldsymbol{\Lambda}, \boldsymbol{\Psi}) p(\mathbf{z}^{(i)})) - \log(Q_i(\mathbf{z}^{(i)})) \right] d\mathbf{z}^{(i)}$$

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



EM FOR FACTOR ANALYSIS: M-STEP

Expanding the equation we have:

$$l(\mathbf{w}) = \sum_{i=1}^m \int_{\mathbf{z}^{(i)}} Q_i(\mathbf{z}^{(i)}) \left[\log \left(p(\mathbf{x}^{(i)} / \mathbf{z}^{(i)}; \boldsymbol{\mu}, \boldsymbol{\Lambda}, \boldsymbol{\Psi}) p(\mathbf{z}^{(i)}) \right) - \log \left(Q_i(\mathbf{z}^{(i)}) \right) \right] d\mathbf{z}^{(i)}$$

$$l(\mathbf{w}) = \sum_{i=1}^m \int_{\mathbf{z}^{(i)}} Q_i(\mathbf{z}^{(i)}) \left[\log \left(p(\mathbf{x}^{(i)} / \mathbf{z}^{(i)}; \boldsymbol{\mu}, \boldsymbol{\Lambda}, \boldsymbol{\Psi}) \right) + \log \left(p(\mathbf{z}^{(i)}) \right) - \log \left(Q_i(\mathbf{z}^{(i)}) \right) \right] d\mathbf{z}^{(i)}$$

Applying the definition of the **expected** value of $\mathbf{z}^{(i)}$ under the distribution Q_i :

$$E_{\mathbf{z}^{(i)} \sim Q_i}[\mathbf{z}^{(i)}] = \int_{\mathbf{z}^{(i)}} Q_i(\mathbf{z}^{(i)}) \mathbf{z}^{(i)} d\mathbf{z}^{(i)}$$

$$l(\mathbf{w}) = \sum_{i=1}^m E_{\mathbf{z}^{(i)} \sim Q_i} \left[\log \left(p(\mathbf{x}^{(i)} / \mathbf{z}^{(i)}; \boldsymbol{\mu}, \boldsymbol{\Lambda}, \boldsymbol{\Psi}) \right) + \log \left(p(\mathbf{z}^{(i)}) \right) - \log \left(Q_i(\mathbf{z}^{(i)}) \right) \right]$$

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



EM FOR FACTOR ANALYSIS: M-STEP

$$l(w) = \sum_{i=1}^m E_{z^{(i)} \sim Q_i} \left[\log \left(p(x^{(i)} / z^{(i)}; \mu, \Lambda, \Psi) \right) + \log \left(p(z^{(i)}) \right) - \log \left(Q_i(z^{(i)}) \right) \right]$$

Substituting the distributions we notice that the **only term depending on the parameters** is:

$$\log \left(p(x^{(i)} / z^{(i)}; \mu, \Lambda, \Psi) \right).$$

Notice that $Q_i(z^{(i)})$ is a **fixed Gaussian** (known parameters that **resulted from the previous maximization or initialization**).

$$l(w) = \sum_{i=1}^m E_{z^{(i)} \sim Q_i} \left[\log \left(p(x^{(i)} / z^{(i)}; \mu, \Lambda, \Psi) \right) + \log \left(\text{Gaussian}(\vec{0}, I) \right) - \log \left(\text{Gaussian}(\mu_{z^{(i)} / x^{(i)}}, \Sigma_{z^{(i)} / x^{(i)}}) \right) \right]$$

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



EM FOR FACTOR ANALYSIS: M-STEP

We **drop** the **terms not depending** on the **parameters**:

$$l(\mathbf{w}) = \sum_{i=1}^m E_{\mathbf{z}^{(i)} \sim q_i} \left[\log \left(p(\mathbf{x}^{(i)} / \mathbf{z}^{(i)}; \boldsymbol{\mu}, \boldsymbol{\Lambda}, \boldsymbol{\Psi}) \right) + \log(\text{Gaussian}(\vec{\mathbf{0}}, I)) - \log \left(\text{Gaussian}(\boldsymbol{\mu}_{\mathbf{z}^{(i)} / \mathbf{x}^{(i)}}, \boldsymbol{\Sigma}_{\mathbf{z}^{(i)} / \mathbf{x}^{(i)}}) \right) \right]$$

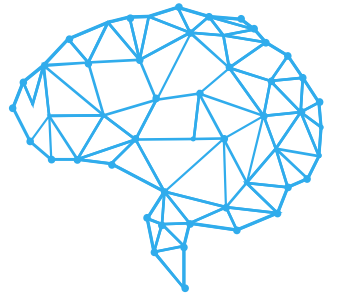
$$l(\mathbf{w}) = \sum_{i=1}^m E_{\mathbf{z}^{(i)} \sim q_i} \left[\log \left(p(\mathbf{x}^{(i)} / \mathbf{z}^{(i)}; \boldsymbol{\mu}, \boldsymbol{\Lambda}, \boldsymbol{\Psi}) \right) \right]$$

Where **we know** that:

$$\mathbf{x}^{(i)} / \mathbf{z}^{(i)} \sim \text{Gaussian}(\boldsymbol{\mu}_{\mathbf{x}^{(i)} / \mathbf{z}^{(i)}}, \boldsymbol{\Sigma}_{\mathbf{x}^{(i)} / \mathbf{z}^{(i)}})$$

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



EM FOR FACTOR ANALYSIS: M-STEP

We **obtain** the **mean** and **covariance matrix** of the **distribution** of $x^{(i)}/z^{(i)}$:

$$x^{(i)}/z^{(i)} \sim \text{Gaussian}(\mu_{x^{(i)}/z^{(i)}}, \Sigma_{x^{(i)}/z^{(i)}})$$

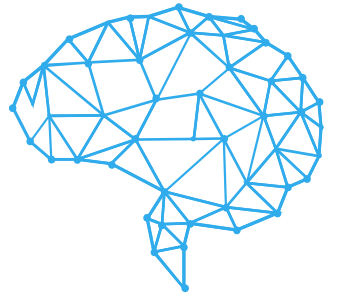
$$\begin{bmatrix} z \\ x \end{bmatrix} \sim N \left(\begin{bmatrix} \vec{0} \\ \mu \end{bmatrix}, \begin{bmatrix} I & \Lambda^T \\ \Lambda & \Lambda \Lambda^T + \Psi \end{bmatrix} \right)$$

$$\mu_{2/1} = \mu_2 + \Sigma_{21}\Sigma_{11}^{-1}(x_1 - \mu_1) \rightarrow \mu_{x^{(i)}/z^{(i)}} = \mu + \Lambda I(z^{(i)} - \vec{0}) = \mu + \Lambda z^{(i)}$$

$$\Sigma_{2/1} = \Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12} \rightarrow \Sigma_{x^{(i)}/z^{(i)}} = \Lambda \Lambda^T + \Psi - (\Lambda I^{-1} \Lambda^T) = \Psi$$

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



EM FOR FACTOR ANALYSIS: M-STEP

Substituting the mean and covariance matrix in the function that we want to maximize:

$$l(\mathbf{w}) = \sum_{i=1}^m E_{\mathbf{z}^{(i)} \sim q_i} \left[\log \left(p(\mathbf{x}^{(i)} / \mathbf{z}^{(i)}; \boldsymbol{\mu}, \boldsymbol{\Lambda}, \boldsymbol{\Psi}) \right) \right]$$

$$l(\mathbf{w}) = \sum_{i=1}^m E_{\mathbf{z}^{(i)} \sim q_i} \left[\log \left(\text{Gaussian}(\boldsymbol{\mu}_{\mathbf{x}^{(i)} / \mathbf{z}^{(i)}}, \boldsymbol{\Sigma}_{\mathbf{x}^{(i)} / \mathbf{z}^{(i)}}) \right) \right]$$

$$l(\mathbf{w}) = \sum_{i=1}^m E_{\mathbf{z}^{(i)} \sim q_i} \left[\log(\text{Gaussian}(\boldsymbol{\mu} + \boldsymbol{\Lambda} \mathbf{z}^{(i)}, \boldsymbol{\Psi})) \right]$$

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



EM FOR FACTOR ANALYSIS: M-STEP

Representing the **explicit form** of the **Gaussian**:

$$l(\mathbf{w}) = \sum_{i=1}^m E_{\mathbf{z}^{(i)} \sim Q_i} [\log(\text{Gaussian}(\boldsymbol{\mu} + \boldsymbol{\Lambda} \mathbf{z}^{(i)}, \boldsymbol{\Psi}))]$$

$$l(\mathbf{w}) = \sum_{i=1}^m E_{\mathbf{z}^{(i)} \sim Q_i} \left[\log \left(\frac{1}{\sqrt{2\pi} |\boldsymbol{\Psi}|^{1/2}} \exp \left(-\frac{1}{2} (\mathbf{x}^{(i)} - \boldsymbol{\mu} - \boldsymbol{\Lambda} \mathbf{z}^{(i)})^T \boldsymbol{\Psi}^{-1} (\mathbf{x}^{(i)} - \boldsymbol{\mu} - \boldsymbol{\Lambda} \mathbf{z}^{(i)}) \right) \right) \right]$$

$$l(\mathbf{w}) = \sum_{i=1}^m E_{\mathbf{z}^{(i)} \sim Q_i} \left[-\frac{1}{2} \log(|\boldsymbol{\Psi}|) - \frac{1}{2} \log(2\pi) - \frac{1}{2} \left[(\mathbf{x}^{(i)} - \boldsymbol{\mu} - \boldsymbol{\Lambda} \mathbf{z}^{(i)})^T \boldsymbol{\Psi}^{-1} (\mathbf{x}^{(i)} - \boldsymbol{\mu} - \boldsymbol{\Lambda} \mathbf{z}^{(i)}) \right] \right]$$

$$l(\mathbf{w}) = \sum_{i=1}^m -E_{\mathbf{z}^{(i)} \sim Q_i} \left[\frac{1}{2} \log(|\boldsymbol{\Psi}|) + \frac{1}{2} \left[(\mathbf{x}^{(i)} - \boldsymbol{\mu} - \boldsymbol{\Lambda} \mathbf{z}^{(i)})^T \boldsymbol{\Psi}^{-1} (\mathbf{x}^{(i)} - \boldsymbol{\mu} - \boldsymbol{\Lambda} \mathbf{z}^{(i)}) \right] \right]$$

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



EM FOR FACTOR ANALYSIS: M-STEP

Maximizing with respect to Λ :

$$\nabla_{\Lambda} \sum_{i=1}^m -E_{z^{(i)} \sim q_i} \left[\frac{1}{2} \log(|\Psi|) + \frac{1}{2} \left[(x^{(i)} - \mu - \Lambda z^{(i)})^T \Psi^{-1} (x^{(i)} - \mu - \Lambda z^{(i)}) \right] \right] = 0$$

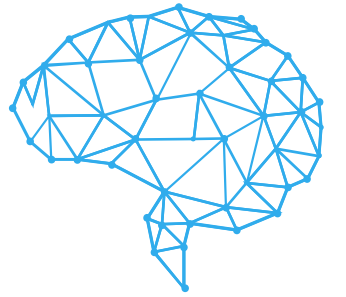
$$\nabla_{\Lambda} \sum_{i=1}^m -E_{z^{(i)} \sim q_i} \left[\frac{1}{2} \left[(x^{(i)} - \mu - \Lambda z^{(i)})^T \Psi^{-1} (x^{(i)} - \mu - \Lambda z^{(i)}) \right] \right] = 0$$

Distributing operations and getting only **results depending only** on Λ we have:

$$\nabla_{\Lambda} \sum_{i=1}^m -E_{z^{(i)} \sim q_i} \left[\frac{1}{2} \left[-x^{(i)T} \Psi^{-1} \Lambda z^{(i)} + \mu^T \Psi^{-1} \Lambda z^{(i)} - z^{(i)T} \Lambda^T \Psi^{-1} x^{(i)} + z^{(i)T} \Lambda^T \Psi^{-1} \mu + z^{(i)T} \Lambda^T \Psi^{-1} \Lambda z^{(i)} \right] \right] = 0$$

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



EM FOR FACTOR ANALYSIS: M-STEP

Reducing terms:

$$\nabla_{\Lambda} \sum_{i=1}^m -E_{z^{(i)} \sim Q_i} \left[\frac{1}{2} \left[-x^{(i)T} \Psi^{-1} \Lambda z^{(i)} + \mu^T \Psi^{-1} \Lambda z^{(i)} - z^{(i)T} \Lambda^T \Psi^{-1} x^{(i)} + z^{(i)T} \Lambda^T \Psi^{-1} \mu + z^{(i)T} \Lambda^T \Psi^{-1} \Lambda z^{(i)} \right] \right] = 0$$

$$\nabla_{\Lambda} \sum_{i=1}^m -E_{z^{(i)} \sim Q_i} \left[\frac{1}{2} \left[z^{(i)T} \Lambda^T \Psi^{-1} \Lambda z^{(i)} + 2z^{(i)T} \Lambda^T \Psi^{-1} \mu - 2z^{(i)T} \Lambda^T \Psi^{-1} x^{(i)} \right] \right] = 0$$

$$\nabla_{\Lambda} \sum_{i=1}^m -E_{z^{(i)} \sim Q_i} \left[\frac{1}{2} \left[z^{(i)T} \Lambda^T \Psi^{-1} \Lambda z^{(i)} + 2z^{(i)T} \Lambda^T \Psi^{-1} (\mu - x^{(i)}) \right] \right] = 0$$

$$\nabla_{\Lambda} \sum_{i=1}^m E_{z^{(i)} \sim Q_i} \left[-\frac{1}{2} z^{(i)T} \Lambda^T \Psi^{-1} \Lambda z^{(i)} + z^{(i)T} \Lambda^T \Psi^{-1} (x^{(i)} - \mu) \right] = 0$$

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



EM FOR FACTOR ANALYSIS: M-STEP

We apply the following properties $\nabla_x b^T x = b$ and $\nabla_x x^T A x = 2Ax$:

$$\nabla_{\Lambda} \sum_{i=1}^m E_{z^{(i)} \sim Q_i} \left[-\frac{1}{2} z^{(i)T} \Lambda^T \Psi^{-1} \Lambda z^{(i)} + z^{(i)T} \Lambda^T \Psi^{-1} (x^{(i)} - \mu) \right] = 0$$

$$\nabla_{\Lambda} \sum_{i=1}^m E_{z^{(i)} \sim Q_i} \left[-\frac{1}{2} \Lambda^T \Psi^{-1} z^{(i)} z^{(i)T} \Lambda + \Psi^{-1} (x^{(i)} - \mu) z^{(i)T} \Lambda \right] = 0$$

$$\sum_{i=1}^m E_{z^{(i)} \sim Q_i} \left[-\frac{1}{2} 2 \left(\Psi^{-1} z^{(i)} z^{(i)T} \Lambda \right) + \Psi^{-1} (x^{(i)} - \mu) z^{(i)T} \right] = 0$$

$$\sum_{i=1}^m E_{z^{(i)} \sim Q_i} \left[-\Psi^{-1} z^{(i)} z^{(i)T} \Lambda + \Psi^{-1} (x^{(i)} - \mu) z^{(i)T} \right] = 0$$

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



EM FOR FACTOR ANALYSIS: M-STEP

We **distribute** the **expected value** and the **sum**:

$$\sum_{i=1}^m E_{z^{(i)} \sim Q_i} \left[-\Psi^{-1} z^{(i)} z^{(i)T} \Lambda + \Psi^{-1} (x^{(i)} - \mu) z^{(i)T} \right] = 0$$

$$-\sum_{i=1}^m E_{z^{(i)} \sim Q_i} \left[\Psi^{-1} z^{(i)} z^{(i)T} \Lambda \right] + \sum_{i=1}^m E_{z^{(i)} \sim Q_i} \left[\Psi^{-1} (x^{(i)} - \mu) z^{(i)T} \right] = 0$$

$$\Psi^{-1} \Lambda \sum_{i=1}^m E_{z^{(i)} \sim Q_i} \left[z^{(i)} z^{(i)T} \right] = \Psi^{-1} \sum_{i=1}^m (x^{(i)} - \mu) E_{z^{(i)} \sim Q_i} \left[z^{(i)T} \right]$$

$$\Lambda \sum_{i=1}^m E_{z^{(i)} \sim Q_i} \left[z^{(i)} z^{(i)T} \right] = \sum_{i=1}^m (x^{(i)} - \mu) E_{z^{(i)} \sim Q_i} \left[z^{(i)T} \right]$$

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



EM FOR FACTOR ANALYSIS: M-STEP

We **solve** for Λ :

$$\Lambda \sum_{i=1}^m E_{z^{(i)} \sim Q_i} [z^{(i)} z^{(i)T}] = \sum_{i=1}^m (x^{(i)} - \mu) E_{z^{(i)} \sim Q_i} [z^{(i)T}]$$

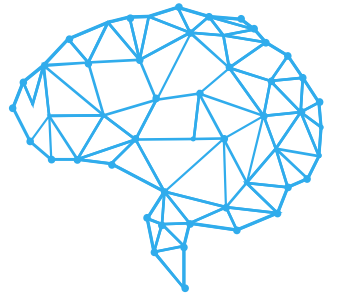
$$\Lambda = \left(\sum_{i=1}^m (x^{(i)} - \mu) E_{z^{(i)} \sim Q_i} [z^{(i)T}] \right) \left(\sum_{i=1}^m E_{z^{(i)} \sim Q_i} [z^{(i)} z^{(i)T}] \right)^{-1}$$

It is **interesting** to **note** the **close relationship** between this equation and the **normal equation** that we'd **derived** for **least squares regression**:

$$w^T = (y^T X)(X^T X)^{-1}$$

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



EM FOR FACTOR ANALYSIS: M-STEP

The **analogy** is that **here**, the **x**'s are a **linear function of** the **z**'s (plus noise).

Given the “**guesses**” for **z** that the **E-step** has found, we will now try to **estimate** the **unknown linearity** Λ relating the **x**'s and **z**'s.

$$\mathbf{w}^T = (\mathbf{y}^T \mathbf{X})(\mathbf{X}^T \mathbf{X})^{-1}$$

$$\Lambda = \left(\sum_{i=1}^m (\mathbf{x}^{(i)} - \mu) E_{\mathbf{z}^{(i)} \sim Q_i} [\mathbf{z}^{(i)T}] \right) \left(\sum_{i=1}^m E_{\mathbf{z}^{(i)} \sim Q_i} [\mathbf{z}^{(i)} \mathbf{z}^{(i)T}] \right)^{-1}$$

The **final step** is to **work out** the **expectations** $E_{\mathbf{z}^{(i)} \sim Q_i} [\mathbf{z}^{(i)T}]$ and $E_{\mathbf{z}^{(i)} \sim Q_i} [\mathbf{z}^{(i)} \mathbf{z}^{(i)T}]$.

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



EM FOR FACTOR ANALYSIS: M-STEP

The **expectation** $E_{\mathbf{z}^{(i)} \sim Q_i} [\mathbf{z}^{(i)T}]$ we **already have it** from our **previous definition**:

$$E_{\mathbf{z}^{(i)} \sim Q_i} [\mathbf{z}^{(i)T}] = \boldsymbol{\mu}_{\mathbf{z}^{(i)} / \mathbf{x}^{(i)}}^T$$

$$E_{\mathbf{z}^{(i)} \sim Q_i} [\mathbf{z}^{(i)T}] = \left(\boldsymbol{\Lambda}^T (\boldsymbol{\Lambda} \boldsymbol{\Lambda}^T + \boldsymbol{\Psi})^{-1} (\mathbf{x}^{(i)} - \boldsymbol{\mu}) \right)^T$$

For the **expectation** $E_{\mathbf{z}^{(i)} \sim Q_i} [\mathbf{z}^{(i)} \mathbf{z}^{(i)T}]$, we **have the following definition**:

$$\text{Cov}(\mathbf{Y}) = E[\mathbf{Y}\mathbf{Y}^T] - E[\mathbf{Y}]E[\mathbf{Y}]^T$$

$$E_{\mathbf{z}^{(i)} \sim Q_i} [\mathbf{z}^{(i)} \mathbf{z}^{(i)T}] = \boldsymbol{\mu}_{\mathbf{z}^{(i)} / \mathbf{x}^{(i)}} \boldsymbol{\mu}_{\mathbf{z}^{(i)} / \mathbf{x}^{(i)}}^T + \boldsymbol{\Sigma}_{\mathbf{z}^{(i)} / \mathbf{x}^{(i)}}$$

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



EM FOR FACTOR ANALYSIS: M-STEP

Substituting both expectations back in the **expression** of Λ we have:

$$\Lambda = \left(\sum_{i=1}^m (x^{(i)} - \mu) E_{z^{(i)} \sim q_i} [z^{(i)T}] \right) \left(\sum_{i=1}^m E_{z^{(i)} \sim q_i} [z^{(i)} z^{(i)T}] \right)^{-1}$$

$$\Lambda = \left(\sum_{i=1}^m (x^{(i)} - \mu) \mu_{z^{(i)}/x^{(i)}}^T \right) \left(\sum_{i=1}^m \mu_{z^{(i)}/x^{(i)}} \mu_{z^{(i)}/x^{(i)}}^T + \Sigma_{z^{(i)}/x^{(i)}} \right)^{-1}$$

In this **last equation** $\Sigma_{z^{(i)}/x^{(i)}}$ is the **covariance matrix** of the **posterior**, which **represents** the **uncertainty** of our **estimations**.

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



EM FOR FACTOR ANALYSIS: M-STEP

Now, that we have the parameter Λ , we want to maximize with respect to μ :

$$\nabla_{\mu} \sum_{i=1}^m -E_{z^{(i)} \sim q_i} \left[\frac{1}{2} \left[(x^{(i)} - \mu - \Lambda z^{(i)})^T (x^{(i)} - \mu - \Lambda z^{(i)}) \right] \right] = 0$$

To save time in the computations, the result of the maximization for is μ :

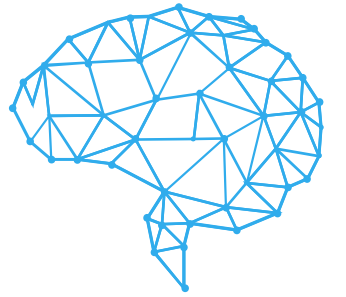
$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

An important thing is that this parameter does not change as the parameters vary, which is different for the parameter Λ which we have just calculated.

Therefore, μ can be calculated just once and needs not be further updated as the algorithm is run.

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



EM FOR FACTOR ANALYSIS: M-STEP

Finally, if you **maximize** for Ψ you **obtain** the **following matrix**:

$$\Phi = \frac{1}{m} \sum_{i=1}^m \mathbf{x}^{(i)} \mathbf{x}^{(i)T} - \mathbf{x}^{(i)} \boldsymbol{\mu}_{\mathbf{z}^{(i)} / \mathbf{x}^{(i)}}^T \boldsymbol{\Lambda}^T - \boldsymbol{\Lambda} \boldsymbol{\mu}_{\mathbf{z}^{(i)} / \mathbf{x}^{(i)}} \mathbf{x}^{(i)T} + \boldsymbol{\Lambda} \left(\boldsymbol{\mu}_{\mathbf{z}^{(i)} / \mathbf{x}^{(i)}} \boldsymbol{\mu}_{\mathbf{z}^{(i)} / \mathbf{x}^{(i)}}^T + \boldsymbol{\Sigma}_{\mathbf{z}^{(i)} / \mathbf{x}^{(i)}} \right) \boldsymbol{\Lambda}^T$$

The **diagonal** of Ψ is **obtained** by **setting** $\Psi_{ii} = \Phi_{ii}$

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



EM FOR FACTOR ANALYSIS: SUMMARY

1. Initialize parameters μ , Λ , and Ψ .
2. Estimate the posterior distribution $Q_i(z^{(i)})$:

$$\begin{aligned}\mu_{z^{(i)}/x^{(i)}} &= \Lambda^T (\Lambda \Lambda^T + \Psi)^{-1} (x^{(i)} - \mu) \\ \Sigma_{z^{(i)}/x^{(i)}} &= I - \Lambda^T (\Lambda \Lambda^T + \Psi)^{-1} \Lambda\end{aligned}$$

$$z^{(i)}/x^{(i)} \sim N(\mu_{z^{(i)}/x^{(i)}}, \Sigma_{z^{(i)}/x^{(i)}})$$

$$\begin{aligned}E_{z^{(i)} \sim Q_i}[z^{(i)T}] &= \mu_{z^{(i)}/x^{(i)}}^T \\ E_{z^{(i)} \sim Q_i}[z^{(i)T}] &= \left(\Lambda^T (\Lambda \Lambda^T + \Psi)^{-1} (x^{(i)} - \mu) \right)^T\end{aligned}$$

$$Q_i(z^{(i)}) = \frac{1}{(2\pi)^{k/2} |\Sigma_{z^{(i)}|x^{(i)}}|^{1/2}} \exp \left(-\frac{1}{2} (z^{(i)} - \mu_{z^{(i)}|x^{(i)}})^T \Sigma_{z^{(i)}|x^{(i)}}^{-1} (z^{(i)} - \mu_{z^{(i)}|x^{(i)}}) \right)$$

3. Compute parameters and go back to step 2 (repeat until convergence):

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)} \quad \Lambda = \left(\sum_{i=1}^m (x^{(i)} - \mu) E_{z^{(i)} \sim Q_i}[z^{(i)T}] \right) \left(\sum_{i=1}^m E_{z^{(i)} \sim Q_i}[z^{(i)} z^{(i)T}] \right)^{-1}$$

$$\Psi_{ii} = \Phi_{ii}$$

Computed only one time

Check previous slide

EXPECTATION MAXIMIZATION

EM AND FACTOR ANALYSIS



EM FOR FACTOR ANALYSIS: SUMMARY

To **compute** the **probability** of a **new sample** we would have $x^{(i)}$:

$$x/z \sim N(\mu + \Lambda z, \Psi)$$

$$p(x) = \frac{1}{\sqrt{2\pi}|\Psi|^{1/2}} \exp\left(-\frac{1}{2}(x^{(i)} - \mu - \Lambda z^{(i)})^T \Psi^{-1}(x^{(i)} - \mu - \Lambda z^{(i)})\right)$$