# Generative models

   I  For classification

   II  For Clustering

   III  Approximate learning with ELBO

**Intro:**  $y = \{1, \dots, K\}$

$X = \mathbb{R}^d$    observed data

**Discriminative:**  $P(Y=k \mid X=x)$  $\overset{\mathbb{R}^{d\times k}}{\underset{\nearrow \mathbb{R}^k}{\uparrow}}$

$\rightarrow$ Logistic regression    $\Theta = \{W, b\}$

$$\mathbb{P}(k \mid x) = \frac{e^{W_{:k}^T x + b}}{\sum_{\ell=1}^{K} e^{W_{:\ell}^T x + b}}$$   **Direct**

**Generative:**  Model $x \neq$ for each $k$

$$\mathbb{P}_k(x) \rightarrow P(X=x, Y=k)$$

$$= \mathbb{P}(X=x \mid Y=k)$$ assumed to be normal

$\boxed{x \sim \mathcal{N}(\mu_k, \Sigma)}$

$$= \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu_k)^T \Sigma^{-1}(x-\mu_k)}$$

(LDA) Linear Discriminant Analysis

Posterior

**Classification:** $\mathbb{P}(Y=k \mid X=x) = \dfrac{\textcircled{1}\,(\mathbb{P}(X=x \mid Y=k)) * \textcircled{2}\,(\mathbb{P}(Y=k))}{\textcircled{3}\,\mathbb{P}(X=x)}$

$\textcircled{1} = \mathbb{P}_k(x)$ density function

$\textcircled{2}$ Prior $= \pi_k = \mathbb{P}(Y=k)$

$\textcircled{3}$ Evidence

$$P(X=x) = \sum_{\ell=1}^{K} P(Y=\ell) * P(X=x \mid Y=\ell)$$

LDA assumes $\quad x \sim \mathcal{N}(\mu_k, \Sigma)$

$$\Theta = \{ \mu_1, \dots, \mu_K, \underset{\substack{\downarrow \\ \in \mathbb{R}^d}}{\Sigma}, \underset{\substack{\downarrow \\ \in \mathbb{R}^{d\times d}}}{\Pi_1, \dots, \Pi_K} \}$$
$$\underset{\in \mathbb{R}}{\downarrow}$$

Largest $P_\Theta(Y=k \mid X=x)$

$$= \underset{k}{\arg\max}\left( \log \Pi_k * P_k(x) \right)$$

$$= \underline{\quad\quad} \underbrace{\left[ \log(\Pi_k) + x^T \Sigma^{-1}\mu_k - \frac{1}{2}\mu_k^T \Sigma^{-1}\mu_k \right]}_{\text{Discriminant score } \delta_k(x)}$$

$k_1, k_2 \quad \delta_{k_1}(x) = \delta_{k_2}(x) \quad (\_ - \_)x + (\_) = 0$

Learning $\Theta$ ? $\quad \frac{1}{n}\sum_{i=1}^{n} \log P_\Theta(x^{(i)}, y^{(i)}) = \mathcal{L}(\Theta)$

MLE $\qquad\qquad \underset{\Theta}{\text{Arg max }} \mathcal{L}(\Theta)$

$$P_\Theta(x, y) = \underset{\Pi}{P}(y) * \underset{M, \Sigma}{P}(x \mid y)$$

Systematic derivation $\quad \frac{d}{d\Pi_\ell}\left( P_\Theta(x, y) \right) = 0 \quad \Big| \quad P_\Theta(x^{(i)}, y^{(i)})$

$\frac{d}{d\mu_\ell}\left( P_\Theta(x, y) \right) = 0 \quad \Big|$

$\qquad\qquad y^{(i)} = k \qquad P_\Pi(y^{(i)} = k) = \Pi_k$

$$\left( \sum_{i=1}^{n} 1 \right) = nk$$

Frequency: $\Pi_k = \left( \underset{\substack{i=1 \\ y^{(i)}=k}}{\overset{n}{\sum}} 1 \right) / n$

$$\mu_k = \frac{1}{n_k} \sum_{\substack{i=1 \\ y^{(i)}=k}}^{n} x^{(i)} \qquad \sum = \frac{1}{n} \sum_{k=1}^{K} \sum_{\substack{i=1 \\ y^{(i)}=k}}^{m} (x^{(i)} - \mu_k)(x^{(i)} - \mu_k)^T$$

QDA : Quadratic DA

for $k$, $x \sim \mathcal{N}(\mu_k, \Sigma_{ik})$

$$\Sigma_{ik} = \frac{1}{n_k} \sum_{\substack{i=1 \\ y^{(i)}=k}}^{n} (x^{(i)} - \mu_k)(x^{(i)} - \mu_k)^T$$

$$\delta_k(x) = \log(\pi_k) + x^T \Sigma_{ik}^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma_{ik}^{-1} \mu_k$$

new term $\left( -\frac{1}{2} x^T \Sigma_{ik}^{-1} x \right.$

$$\delta_{k_1}(x) = \delta_{k_2}(x) \longrightarrow \underline{\phantom{a}} x^2 + \underline{\phantom{a}} x + \underline{\phantom{a}} = 0$$

LDA versus Logistic regression ?

QDA = LDA + interaction between coordinates of $x$

LDA : Assumption of normality

- $\checkmark$ more efficient
- $\times$ LR does MLE $\pi(y=k \mid x)$
  $\hookrightarrow$ it's more efficient

Generative versus Discriminative
- — Include prior knowledge
- — Density function $\searrow \rightarrow$ Sample
  $\rightarrow$ Detect anomalies $\qquad \rightarrow$ missing values

# II Clustering    K clusters

GMM : Gaussian Mixture Model

$$P_\theta(x,z) = P_\theta(x|z) P_\theta(z)$$

Assume generative process $(i)$

$$z^{(i)} \sim \text{Multi}(1, \dots K)$$

$$x^{(i)} | z^{(i)} = k \sim \mathcal{N}(\mu_k, \Sigma_k)$$

- You know $z^{(i)} \rightarrow$ QDA !

- You don't know $z^{(i)} \rightarrow$ Clustering

  Observations $x^{(i)}$

  Latent variables $z^{(i)}$

$$\theta = \{ \pi_1, -, \pi_k, \mu_1, -, \mu_K, \Sigma_1, -, \Sigma_k \}$$

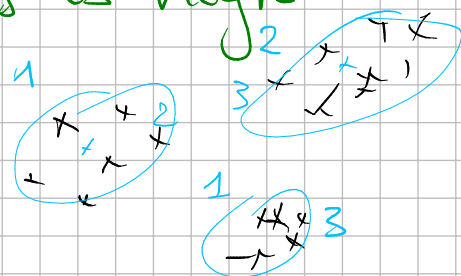MLE   $\underset{\theta}{\text{Argmax}} \left( \dfrac{1}{n} \sum_{i=1}^{n} \log P_\theta(x^{(i)}|) \right) \rightarrow$ Unsupervised learning

K !

$$\log P_\theta(x^{(i)}) = \log \left( \sum_{k=1}^{K} P_\theta(x^{(i)}, z^{(i)} = k) \right)$$

at least one $k$ for which this is high

Difficulty :

- no closed formed solution
- not identifiable

Reasoning for smarter learning:

 1 - $z^{(i)}$ are latent;

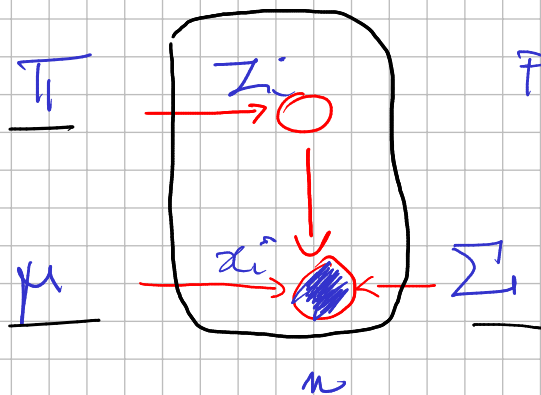  But if we know them, easy to find $\Theta$

 2 - Estimating the $z^{(i)}$ ?

Expectation - Maximization ( EM )

 1. Find the expected clusters ( E-step) from $\Theta$

 2. Maximize the expected likelihood (M-step)
  for $\Theta$

EM: general estimation for incomplete data.



$\Pi$   $Z_i$    PGM: probabilistic
         Graphical Model

$\mu$   $x_i$   $\Sigma$

   $n$

- K-means $\left\{ \begin{array}{l} \text{finding clusters} \to \text{Hard assignments} \\ \text{finding new centers} \to \text{Isotropic gaussian} \end{array} \right.$

              $\Sigma = \sigma^2 Id$

- EM, for $t = 0$ ... until convergence
 1. E-step; for each $x^{(i)}$
   $\mathbb{E}[z^{(i)} | x^{(i)}] \to \mathbb{P}_{\Theta^t}(z | x^{(i)})$

   $\mathbb{P}_{\Theta^t}(z = k | x^{(i)}) = \dfrac{\Pi_k * \mathbb{P}_{\Theta_k}(x^{(i)})}{\sum\limits_{\ell=1}^{K} \Pi_\ell \, \mathbb{P}_{\Theta_\ell}(x^{(i)})}$

2 - M-step: Compute $\theta^{t+1}$

$$\theta^{t+1} = \underset{\theta}{\arg\max} \sum_{i=1}^{n} \mathbb{E}_{z^{(i)} \sim P_{\theta^t}(z|x^{(i)})} \left[ \log P_\theta(x^{(i)}, z^{(i)}) \right]$$

→ solve this, for example with GD

− Why EM ?  → Very simple / easy to implement
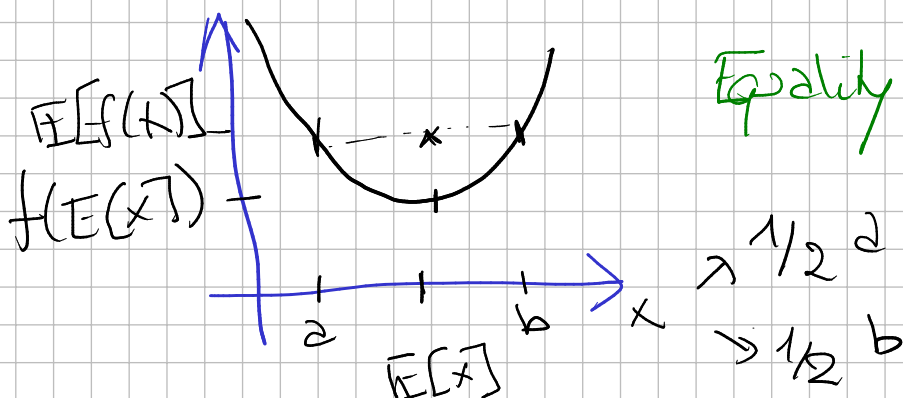
→ Guaranteed to converge

Issue: local optima

**III** Why does EM work ?

→ More general formulation: ELBO

(Evidence Lower bound).

Jensen's inequality : $f$ convex

$$\mathbb{E}[f(x)] \geq f(\mathbb{E}[x])$$

Equality $\Longleftrightarrow$ $x$ constant

$\mathbb{E}[f(x)]$

$f(\mathbb{E}[x])$



$x$

$\mathbb{E}[x]$

$\geq \frac{1}{2} a$

$\geq \frac{1}{2} b$

$x, z \rightarrow \log P_\theta(x) = \log \sum_{k=1}^{K} P_\theta(x, z=k)$

In GMM: $P_\theta(x, z=k) = P_\theta(z=k) \times P_\theta(x|z=k)$

Easy !

$\rightarrow$ We don't know how to write $P_\theta(x, z=k)$

$\qquad \rightarrow$ We parametrize the value of $z$

$$\theta = \{\pi, \mu, \Sigma\}$$

$$\phi = \{z\}$$

Distribution $\qquad Q_\phi(z=k) = \phi_k$

(proposal) $\qquad \sum_{k=1}^{K} \phi_k = 1$

$$\log \sum_{k=1}^{K} P_\theta(x, z=k) = \log \sum_{k=1}^{K} Q_\phi(z=k) * \frac{P_\theta(x, z=k)}{Q_\phi(z=k)}$$

$$\overset{||}{\log\left( \mathop{\mathbb{E}}_{z \sim Q_\phi} \left[ \frac{P_\theta(x, z)}{Q_\phi(z)} \right] \right)}$$

Apply Jensen's:

$$\log P_\theta(x) \underset{\text{Jensen}}{\geq} \mathop{\mathbb{E}}_{z \sim Q_\phi} \left[ \log\left( \frac{P_\theta(x, z)}{Q_\phi(z)} \right) \right]$$

$$\forall \phi, \theta, x \quad , \rightarrow \text{ Lower bound on}$$
$$\text{my evidence}$$

$\rightarrow$ Choose $Q$ : sample $z$ from $Q$

$\rightarrow$ What $Q_\phi$ should we choose?

We want the bound to be tight

$$\rightarrow \text{Equality} \iff \frac{\mathbb{P}_\theta(x,z)}{Q_\phi(z)} = C \iff Q_\phi(z) \propto \mathbb{P}_\theta(x,z)$$

$$\text{But } \sum_{k=1}^{K} Q_\phi(z=k) = 1 \rightarrow Q_\phi(z=k) = \frac{\mathbb{P}_\theta(x,z=k)}{\underbrace{\sum_{\ell=1}^{K} \mathbb{P}_\theta(x,z=\ell)}_{\mathbb{P}_\theta(x)}}$$

$$Q_\phi(z=k) = \mathbb{P}_\theta(z=k \mid x)$$
$$\text{posterior !}$$

$$ELBO(x,\phi,\theta) = \mathbb{E}_{z \sim Q_\phi} \left[ \log \left( \frac{\mathbb{P}_\theta(x,z)}{Q_\phi(z)} \right) \right]$$

- EM?  E-step:  setting $Q_\phi(z) = \mathbb{P}_\theta(z \mid x)$

$$\rightarrow \log \mathbb{P}_\theta(x) = ELBO(x,\phi,\theta)$$

M-step:  $\underset{\theta}{\text{Argmax}} \; ELBO(x,\phi,\theta)$

- Many examples $x^{(1)} \ldots, x^{(n)}$

$$Q_{\phi^{(i)}}(z=k) = \phi_k^{(i)} \rightarrow \text{We have}$$
$$n \times K$$
$$\text{parameters.}$$

$$\sum_{i=1}^{n} \log \mathbb{P}_\theta(x^{(i)}) \geq \sum_{i=1}^{M} ELBO(x^{(i)}, \phi^{(i)}, \theta)$$

- Generalized EM:

Block
Coordinate
Ascent

- $\forall i, \; \phi^{(i)^{t+1}} = \underset{\phi \in \Delta(K)}{\text{argmax}} \; ELBO\left(x^{(i)}, \phi, \theta^t\right)$

- $\theta^{t+1} = \underset{\theta}{\text{argmax}} \sum_{i=1}^{n} ELBO\left(x^{(i)}, \phi^{(i)^{t+1}}, \theta\right)$

- Show convergence: $\ell(\theta^t) \leq \ell(\theta^{t+1})$ ?

$$\ell(\theta^t) \doteq \sum_{i=1}^{n} \log P_{\theta^t}(x^{(i)})$$

(E-step + Jensen) $\qquad = \sum_{i=1}^{n} ELBO(x^{(i)}, \phi^{(i)^t}, \theta^t)$

$\downarrow$ max w.r.t $\theta$

$$\ell(\theta^{t+1})$$

- For GMM :

M-step : $\forall p \in \theta$

Optim : $\qquad \dfrac{\partial}{\partial p} \sum_{i=1}^{n} ELBO(x^{(i)}, \phi^{(i)}, \theta^t) = 0$

$\rightarrow p^{t+1}$

- Interpret

$$\log P_\theta(x) - ELBO(x, \phi, \theta) = D_{KL}(Q_\phi \| P_{\theta}z|x)$$

- What's next ?

- E-step $\rightarrow$ $Q_\phi(z) = \dfrac{P_\theta(z) \, P_\theta(x|z)}{\left( \sum_{k=1}^{K} P_\theta(z=k) \, P_\theta(x|z=k) \right)}$

1 What when it's intractable

Mean field assumption :

$$Q_\phi(z) = \prod_{k=1}^{K} Q_\phi^k(z_k) \quad \left( \begin{array}{l} \text{Independent} \\ \text{latent} \\ \text{coordinates} \end{array} \right)$$

$\rightarrow$ Sigmoid Belief Network $\rightarrow$ $2^K$

$Z$ discrete $\rightarrow$ $Q$ multinomial $(\pi)$

$Z$ continuous $\rightarrow$

$$\phi = \{\xi, \psi\} \qquad Q_\phi = \mathcal{N}\left(q_\xi(x), \text{diag}\left(v_\psi(x)\right)^2\right)$$

$Z \sim Q$

$$\begin{pmatrix} q_\xi \\ v_\psi \end{pmatrix} \rightarrow \text{Neural Networks}$$

$\rightarrow$ Variational Auto-encoders.

$\hookrightarrow$ Encoder

$$x \mid z \sim \mathcal{N}\left(g_\theta(z), \sigma^2 I_d\right)$$

$\downarrow$

decoder

Learning:

E step: Evaluate $Q^{(i)}(z^{(i)}) \rightarrow$ Compute the density (encoder)

M step: $\mathbb{E}_{Z^{(i)} \sim Q^{(i)}} \rightarrow$ sample from $Q^{(i)}$

$\downarrow$

ELBO

$$\theta = \theta + \eta \nabla_\theta \text{ELBO} \qquad \rightarrow \text{easy}$$

$$\xi = \text{-----} \left.\begin{matrix} \\ \\ \end{matrix}\right\} \rightarrow \text{Difficult}$$

$$\psi = \text{-----}$$

$\rightarrow$ reparametrization trick.