→ Sign-up for a presentation slot
( random assignment after
tonight 11/16 midnight)

→ HW4 Dec 1st , Project

—

Last time : Generative assumptions
on the data

→ ML estimation , ELBO objective

→ EM algorithm (today)

→ Variational Inference

Summary : Learn data distribution

Variational

Sampling

**Parametric assumption on the distribution and learn the unknown parameters that best fit the data**

E.g.
ML * (today)
* VI (Bayesian)
(today)

Input: data
Incomplete description of distribution

↳ Generate more samples from the distribution

E.g.
Ⓐ→ * MCMC
Markov Chain
Monte Carlo
variants

Ⓑ→ * GAN , VAE
(CNN representations
of distributions)

Ⓒ→ Stein Variational Gradient
Descent (SVGD) and particle-based
deterministic/ stochastic algorithms *

Ⓓ→ Variational perspective "Transport" methods
Optimal transport, Normalizing flows etc

Ⓔ→ Score- generative models

## $\underline{LDA}$    Linear Discriminant Analysis

Introduces idea of making generative
    assumptions on the data

Setting: binary classification
    $h(x) = 1$ or $-1$

Bayes optimal $\underline{classifier}$

$\underline{Assumption}$ : • $P(Y) = \begin{cases} \frac{1}{2} & Y = 1 \\ \frac{1}{2} & Y = -1 \end{cases}$ (Uniform)

$\underline{Conditional\ Gaussian}$

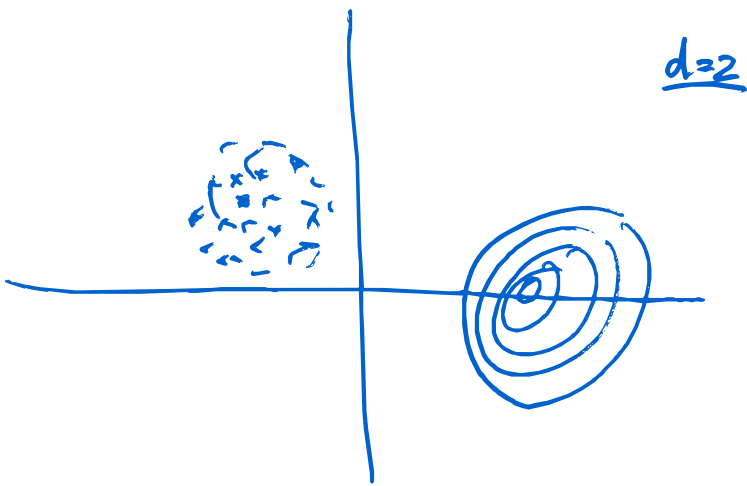• $P(X = x \mid Y = y) = \dfrac{e^{-(x-\mu_y)^T \Sigma^{-1}(x-\mu_y)}}{(2\pi)^{d/2}\,|\Sigma|^{1/2}}$

$x \in \mathbb{R}^d$     $\pm 1$

**Box B**

$$h_{Bayes}(x) = sgn(w \cdot x + b)$$

$$w = (\mu_1 - \mu_{-1})^T \Sigma^{-1}$$

$$b = \frac{1}{2}\left(\mu_{-1}^T \Sigma^{-1} \mu_{-1} - \mu_1^T \Sigma^{-1} \mu_1\right)$$

$\underline{d=2}$



$$h_{Bayes}(x) = \underset{y}{arg\,max}\left\{ P(Y = y \mid X = x)\right\}$$

$$P(Y = y \mid X = x) = \frac{P(Y = y)\, P(X = x \mid Y = y)}{P(X = x)}$$

$$h_{Bayes}(x) = \underset{y}{arg\,max}\left\{ P(Y = y) P(X = x \mid Y = y)\right\}$$

$$(\because P(Y = y) = \tfrac{1}{2} \quad y = \pm 1)$$

$$= \underset{y}{arg\,max}\left\{ P(X = x \mid Y = y)\right\}$$

$$h_{Bayes}(x) = sgn\left(\log \frac{P(X = x \mid Y = 1)}{P(X = x \mid Y = -1)}\right)$$

$$= sgn\left(\log \frac{e^{-(x-\mu_1)^T \Sigma^{-1}(x-\mu_1)}}{e^{-(x-\mu_{-1})^T \Sigma^{-1}(x-\mu_{-1})}}\right)$$

$$= sgn\left(-(x-\mu_1)^T \Sigma^{-1}(x-\mu_1) + (x-\mu_{-1})^T \Sigma^{-1}(x-\mu_{-1})\right)$$

$$(\text{to get box B})$$

Takeaway: Can make parametric generative
assumptions on data distribution
   and solve for $Y$ "easily"

# Gaussian Mixture model

## Generative assumption / Probabilistic model

$$P_\Theta(X, Z) = \sum_{j=1}^{k} \frac{\pi_j \cdot e^{-(x-\mu_j)^T \Sigma_j^{-1}(x-\mu_j)}}{3}$$

$\Theta \uparrow$ set of parameters

$X \uparrow$ observed

$Z \uparrow$ Latent variable

$3 \downarrow$ normalization constant

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

$Z \in [k]$      $k$ : no of components

$$\Theta = \left( \underset{\underset{\mathbb{R}^d}{m}}{\mu_1}, \ldots, \mu_k , \underset{\underset{\mathbb{R}^{d\times d}}{n}}{\Sigma_1}, \ldots, \Sigma_k , \pi_1 \ldots \pi_k \right)$$

$Z$ : Latent variables

$$P_\Theta(X|Z=j) = \frac{1}{(2\pi)^{d/2}|\Sigma|^{1/2}} e^{-(x-\mu_j)^T \Sigma(x-\mu_j)}$$

$Z$ : clustering

# Maximum Likelihood estimation of parameters

$$\ell(\theta) \quad = \quad \max_{\theta} \ \log P_\theta(X_1, \ldots, X_m)$$

iid data $(X_1, \ldots, X_m) = (x_1, x_2, \ldots, x_m)$

$$\ell(\theta) = \max_{\theta} \ \log \prod_{i=1}^{m} P_\theta(X_i)$$

$$= \max_{\theta} \ \sum_{i=1}^{m} \log P_\theta(X_i = x_i)$$

Sources of error

$\longrightarrow$ non-convex, only solved approximately

$\longrightarrow$ $P_\theta$ model could be wrong

## Variational perspective

Latent variable $z$

Data $\quad X$

Parameters $\theta$ $\quad$ that describe $P_\theta(X, z)$

---

Want $\quad$ Variational Inference

$$\underset{q \in Q}{\arg\max} \quad -D_{KL}\left(q \,\|\, P_\theta(\cdot | x)\right)$$

$\downarrow$

set of
probability
distributions
over $Z$

$$D_{KL}\left(q \,\|\, P_\theta(\cdot | x)\right) =$$

$$\underset{q}{E} \log q \quad - \quad \underset{q}{E} \log P_\theta(z | x) \cdot$$

$$= \quad \underset{q}{E} \log q \quad - \quad \underset{q}{E} \log \frac{P_\theta(z, x)}{P_\theta(x)}$$

$$= \quad \underset{q}{E} \log q \quad - \quad \underset{q}{E} \log P_\theta(z, x) \quad +$$

$$\underset{q}{E} \log P_\theta(x)$$

$$\left(\because \underset{q}{E} \log P_\theta(x)\right.$$

$$= \sum_z q(z) \log P_\theta(x)$$

$$\left. = \log P_\theta(x)\right)$$

$$= \quad \underset{q}{E} \log q \quad - \quad \underset{q}{E} \log P_\theta(z, x) \quad + \quad \log P_\theta(x)$$

$$ELBO(q, \theta, x) = \quad \underset{q}{E} \log P_\theta(z, x) \quad -$$

$$\underset{q}{E} \log q$$

Variational objective function

$$D_{KL}\left(q \,\|\, P_\theta(\cdot | x)\right) = -ELBO(q, \theta, x) + \log P_\theta(x)$$

$$\log P_\theta(x) = ELBO(q, \theta, x) + \underbrace{D_{KL}(\quad)}_{\geq 0}$$

For any $q$,

$$\log P_\theta(x) \geq ELBO(q, \theta, x).$$

# EM algorithm

An algorithm for ML estimation in
the presence of Latent variables
and a probabilistic model

$$\ell(\theta) = \sum_{i=1}^{m} \log P_\theta(x_i)$$

## Fix $x$

$$\ell(\theta, x) = \log P_\theta(x) \quad \hookrightarrow \text{Marginal of } x$$

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

$$= \log \sum_{j=1}^{k} P_\theta(z = z_j) P_\theta(x \mid z_j)$$

Iterative algorithm:

E - step : $\quad q_{t+1, i}(z) = P_{\theta_t}(z \mid x_i)$

M - step : $\quad \theta_{t+1} = \underset{\theta}{\text{argmax}} \sum_{i=1}^{m} ELBO(q_{t+1,i}, \theta, x_i)$

$$ELBO(q, \theta, x) = -\underset{q}{E} \log q + \underset{q}{E} \log P_\theta(z, x)$$

Why does EM algorithm
amount to $\underset{\theta}{\max} \sum_{i=1}^{m} \ell(\theta, x_i)$

- For $q(z) = P_\theta(z \mid x)$ :

$$\underset{\theta}{\max} \ell(\theta, x) = \underset{\theta}{\max} ELBO(q, \theta, x)$$

$ELBO(q, \theta, x) =$
$-\underset{q}{E} \log q + \underset{q}{E} \log P_\theta(x) P_\theta(z \mid x)$

$$= -\underset{q}{E} \log q + \underset{q}{E} \log P_\theta(z \mid x) + \log P_\theta(x)$$

$$= \log P_\theta(x)$$

- For a fixed $\theta$,
$$q(z) = P_\theta(z \mid x) \text{ is}$$

$$\underset{q \in Q}{\text{argmax}} \; ELBO(q, \theta, x)$$

EM algorithm always increasing the
$\log P_\theta(x)$

$$\ell(\theta_{t+1}) \gg \ell(\theta_t)$$

$$\downarrow$$

$$\sum_{i=1}^{m} \ell(\theta_{t+1}, x_i)$$

$\ell(\theta_{t+1}) \underset{(E\text{-step})}{=} ELBO(q_{t+2}, \theta_{t+1})$

$\underset{(E\text{-step})}{\geqslant} ELBO(q_{t+1}, \theta_{t+1})$

$\underset{(M\text{-step})}{\geqslant} ELBO(q_{t+1}, \theta_t)$

$$= \ell(\theta_t)$$