# Probabilistic Graphical Models
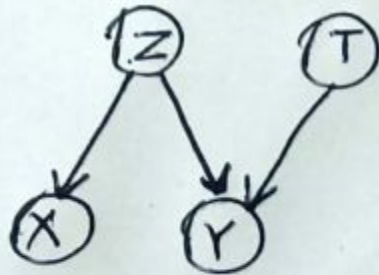
# Lecture 25,26

Learning with incomplete data
Latent Variables Models
Expectation Maximization

# Learning with incomplete data

# Learning with incomplete data

(I) complete: $L(\theta) = \prod_{i=1}^{4} P_\theta(X^i, Y^i, Z^i, T^i) = \prod_{i=1}^{4} \underbrace{P(Z^i)}_{\theta_1} \underbrace{P(T^i)}_{\theta_2} \underbrace{P(Y^i | Z^i, T^i)}_{\theta_3} \underbrace{P(X^i | Z^i)}_{\theta_4}$

(II) incomplete: $Pr(data | \theta) = P_\theta(X^1, Y^1, Z^1, T^1) \quad P_\theta(X^2, Y^2, T^2)$

$P_\theta(X^3, Y^3, Z^3) \quad P_\theta(X^4, Y^4)$

$L(\theta) = P_\theta(X^1, Y^1, Z^1, T^1) \left( \sum_Z P_\theta(X^2, Y^2, Z, T^2) \right) \left( \sum_T P_\theta(X^3, Y^3, Z^3, T) \right)$

$\left( \sum_Z \sum_T P(X^4, Y^4, Z, T) \right)$

**Incomplete** (II)

$X^1, Y^1, Z^1, T^1$

$X^2, Y^2, ?, T^2$

$X^3, Y^3, Z^3, ?$

$X^4, Y^4, ?, ?$

# Learning with incomplete data



K. N. Toosi

$$
\boxed{\text{II}}
$$

incomplete: $Pr(\text{data}|\theta) = \blacksquare\; P_\theta(x^1, y^1, z^1, T^1) \quad P_\theta(x^2, y^2, T^2) \;\text{\textbullet}$

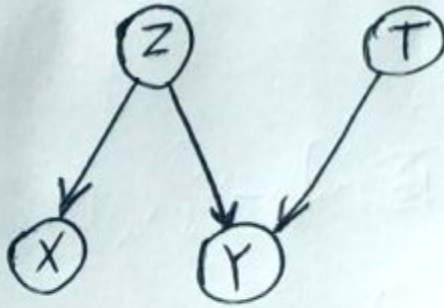$\qquad\qquad\qquad\qquad\qquad\qquad\qquad P_\theta(x^3, y^3, z^3) \quad P_\theta(x^4, y^4)$

$L(\theta) = P_\theta(x^1, y^1, z^1, T^1) \left( \sum_Z P_\theta(x^2, y^2, Z, T^2) \right) \left( \sum_T P_\theta(x^3, y^3, z^3, T) \right)$

$\qquad\qquad\qquad\qquad \left( \sum_Z \sum_T P(x^4, y^4, Z, T) \right)$

$= \;\;\xcancel{P(x^4)} \sum_Z \sum_T \underset{\theta_1}{P(Z)} \; \underset{\theta_2}{P(T)} \; \underset{\theta_3}{P(y^4|Z,T)} \; \underset{\theta_4}{P(x^4|Z)}$

Incomplete $\boxed{\text{I}}$

$x^1, y^1, z^1, T^1$

$x^2, y^2, ?, T^2$

$x^3, y^3, z^3, ?$

$x^4, y^4, ?, ?$

4

# Latent Variables



Latent variables
{ present in model
  absent in data

| | Data | | |
|---|---|---|---|
| X | Y | Z | T |
| $X^1$ | $Y^1$ | ? | ? |
| $X^2$ | $Y^2$ | ? | ? |
| $\vdots$ | | | |
| $X^m$ | $Y^m$ | ? | ? |

X Y    Z T

observed variables — latent (hidden) variables

# Latent Variables: Example



$$P(X_1 - X_n, Z_1 - Z_n) = \prod_{i=2}^{n} P(Z_i | Z_{i-1}) \prod_{i=1}^{n} P(X_i | Z_i)$$

$$P(X_1 - X_n) = \sum_{Z_n} \cdots \sum_{Z_2} \sum_{Z_1} \prod_{i=2}^{n} P(Z_i | Z_{i-1}) \prod_{i=1}^{n} P(X_i | Z_i)$$

$Z_i$'s not observed $\Rightarrow$ all variables $X_1 - X_n$ are dependent

without latent variables

$\Rightarrow$ form a clique of size $\underline{n}$

pgm 25

$P(\text{speed})$

$p(d | \text{speed})$

$Z = (P_n, P_y, \theta)$

6

# Latent Variables



latent variables

$z^i \in \mathbb{R}^p$
$x^i \in \mathbb{R}^n$

$z$ → $x$

directed (most common)

$P(z,x) = \frac{1}{Z(\theta)} \Phi(z,x)$

$z$ — $x$

undirected (less common)

# Latent Variables



latent variables

$z^i \in \mathbb{R}^p$
$x^i \in \mathbb{R}^n$

directed (most common)

$$P_\theta(X, Z) = P_{\theta_1}(Z) \, P_{\theta_2}(X \mid Z)$$

$P_{\theta_1}(Z)$, $P_{\theta_2}(X \mid Z)$ are given (easy)

prior

$$P_\theta(X, Z) = P(Z \mid X) \, P(X)$$

posterior → needed for leany

→ data likehood

hard

$P(X, Z)$ $P(Z)$ $P(X \mid Z)$
easy!

$P(X)$ $P(Z \mid X)$
Difficult!

# Incomplete data introduces complexities

$$data \ (X^1, X^2, \ldots, X^m)$$

$$ll(\theta) = \log P \sum_{i=1}^{m} P_\theta(X^i)$$

$$ll(\theta) = \sum_{i=1}^{m} \log P_\theta(X^i)$$

$$= \sum_{i=1}^{m} \log \sum_{z} P_\theta(X^i, z)$$

$$= \sum_{i=1}^{m} \log \sum_{z} P_{\theta_1}(z) \, P_{\theta_2}(X|z)$$

$\theta_1, \theta_2$ are entangled

# Learning with Incomplete data



Solution 1

$$\frac{\partial \ell\ell(\theta)}{\partial \theta_1}, \quad \frac{\partial \ell\ell(\theta)}{\partial \theta_2} \implies \text{gradient} \quad \text{ascent}$$

$$\frac{\partial}{\partial \theta_1} = \sum_{i=1}^{m} \frac{\sum_{z} \frac{\partial}{\partial \theta_1} P_{\theta_1}(z) \; P_{\theta_2}(X|z)}{\sum_{z} P_{\theta_1}(z) \, P_{\theta_2}(X|z)}$$
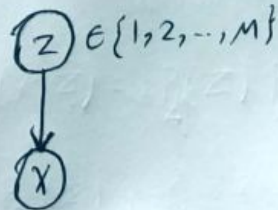
# Example Gaussian Mixture Models



$$P(X|Z) \quad \text{simple}$$

$$P(X) = \int_{\Sigma} P(X,Z) \, d\theta = \int_{\Sigma} P(X|Z) \, P(Z) \, d\theta$$

Gaussian Mixture $\quad Z \in \{1, 2, \cdots, M\}$

$Z \in \{1, 2, \cdots, M\}$

$P(Z=k) = \pi_k$

$P(X|Z) = N(X; \mu_z, \sigma_z^2)$

$P(X|Z=k) = N(X, \mu_k, \sigma_k)$

parameters $\theta = \{(\pi_1, \mu_1, \sigma_1), (\pi_2, \mu_2, \sigma_2), \cdots, (\pi_M, \mu_M, \sigma_M)\}$

$$P(X) = \sum_{k=1}^{M} Pr(Z=k) \, P(X|Z=k) = \sum_{k=1}^{M} \pi_k \, N(X|\mu_k, \sigma_k)$$

# How to learn latent variable models?

$$P_\theta(x)$$

$$x^1, x^2, \dots, x^m$$

$$P_\theta(X, Z)$$

$$P_\theta(X, Z) = P_{\theta_1}(X|Z)\, P_{\theta_2}(Z)$$

$$\text{log-likelihood} \quad \ell\ell(\theta) = \sum_{i=1}^{m} \log P_\theta(x^i) = \sum_{i=1}^{m} \log \sum_{Z} P_\theta(x^i, Z)$$

12

# How to learn latent variable models?

$$ll(\theta_0) \leqslant ll(\theta_1) \leqslant ll(\theta_2) \leqslant ll(\theta_3) \leqslant \cdots \leqslant ll(\theta_t) \leqslant \cdots$$

$\boxed{\theta_t} \to$ current $\theta$

we dont have access to $z$!

heuristic:

$z^i = \underset{z}{\operatorname{argmax}}\ P_{\theta_t}(z \mid x^i)$ $\xrightarrow{\text{Maximize}}$ posterior distribution

OR

$z^i \sim P_{\theta_t}(z \mid x^i)$ take a sample

$$\theta_{t+1} = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^{m} \log P_\theta(x^i; z^i)$$

data

$\begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$

# How to learn latent variable models?

How to compute the posterior?

$$P_\theta(X,z) = P_\alpha(X|z) \, P_\beta(z) \qquad \theta = (\alpha, \beta)$$

$$\theta_t = (\alpha_t, \beta_t)$$

$$P_{\theta_t}(z|x^i) = \frac{P_{\theta_t}(x^i, z)}{\sum_z P_{\theta_t}(x^i, z)} = \frac{P_{\alpha_t}(x^i|z) \, P_{\beta_t}(z)}{\sum_z P_{\alpha_t}(x^i|z) \, P_{\beta_t}(z)}$$

# Expectation Maximization (EM)

Expectation - Maximization

Compute $P_{\theta_t}(Z|x^i)$ for all $Z$.     Expectation step (E-step)

$$\theta_{t+1} = \underset{\theta}{\text{argmax}} \sum_{i=1}^{m} E_{P_{\theta_t}(z|x^i)} \left\{ \log P_{\theta}(x^i, z) \right\} \rightarrow \text{expected log-likelihood}$$

$$\theta_{t+1} = \underset{\theta}{\text{argmax}} \sum_{i=1}^{m} \sum_{Z} P_{\theta_t}(z|x^i) \log P_{\theta}(x^u, z)$$

Maximization-Step
~~M-S~~ (M-step)

EM (Expectation-Maximization): Alternate betweem E-step and M-step.

Can prove: $ll(\theta_0) \le ll(\theta_1) \le ll(\theta_2) \le \cdots$
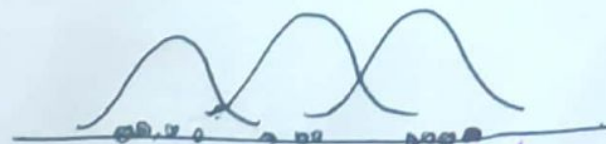
15

# Example: Mixture of Gaussians

pgm 26

Mixture of Gaussians

training data $x^1, x^2, \dots, x^m$

$$P(x) = \sum_{k=1}^{K} w_k \, N(x \mid \mu_k, \sigma_k^2)$$

$$\sum_{k=1}^{K} w_k = 1$$

$P(z) = Pr(Z=z) = w_z$

$$P(x,z) = \underline{P(x \mid z)} \; P(z)$$

$w_1 \, w_2 \; - \; w_k$

$N(x; \mu_k, \sigma_k^2)$

16

# Example: Mixture of Gaussians

$$ll(\theta) = \log \sum_{i=1}^{m} \log P(x^i) = \sum_{i=1}^{m} \log \sum_{k=1}^{K} Pr(Z=k) \underbrace{N(x | \mu_k, \sigma_k^2)}_{w_k}$$

$$\theta = \left( \{\mu_k\}, \{\sigma_k^2\}, w_1, \ldots, w_k \right)$$

$$\theta^t$$

$$P(Z | x^i) = \frac{P(x^i | z) \, P(z)}{\sum_{z} P(x^i | z) \, P(z)}$$

**E-step**

**find**

$$Pr(Z=k | x^i) = \frac{N(x^i; \mu_k^t, \sigma_k^{2\,t}) \, w_k^t}{\sum_{k'=1}^{K} N(x^i; \mu_{k'}^t, \sigma_{k'}^{t\,2}) \, w_{k'}^t} = \alpha_{ik}^t$$

for $x^1, x^2, \ldots, x^m$
for $k = 1, 2, \ldots, K$

# Example: Mixture of Gaussians



**M-step**

$$E_{\underset{\theta_t}{P(z|x^i)}}\left\{ \sum_{i=1}^{m} \log P_\theta(x^i, z) \right\} = \sum_{i=1}^{m} \sum_{z} P_{\theta^t}(z|x^i) \log P_\theta(x^i, z)$$

$$= \sum_{i=1}^{m} \sum_{z=1}^{K} P_{\theta^t}(z|x^i) \left[ \log P_\theta(x^i|z) + \log P_\theta(z) \right]$$

$$\theta^{t+1} = \arg\max_{\theta} \sum_{i=1}^{m} \sum_{k=1}^{K} \alpha_{ik}^t \left[ \log \mathcal{N}(x^i | \mu_k, \sigma_k) + \log w_k \right]$$

$$\{\mu_k\}, \{\sigma_k\}, \{w_k\}$$

$$\frac{\partial}{\partial w_j} \sum_{k=1}^{K} \left( \underbrace{\sum_{i=1}^{m} \alpha_{ik}^t}_{\beta_k} \right) \log w_k \quad \underline{+\lambda\left( \sum w_k - 1 \right)} \qquad \boxed{w_j = \frac{\beta_j}{\sum_k \beta_k}}$$

$$\frac{\beta_j}{w_j} = -\lambda \Rightarrow \frac{w_j}{\beta_j} \text{ cost.} \qquad \sum w_j = 1$$

$$\sum \frac{\beta_j}{-\lambda} = 1$$

18

# Limitations of EM

K. N. Toosi
University of Technology

take a sample from $P_\theta(X,Z)$

$z^i$

1- $z^i \sim P(Z) = N(Z; \mu = 0, \Sigma = I)$    Image of noise

2- Feed $z^i$ to neural nets to get $\mu_\theta(z^i)$

$\Sigma_\theta(z^i)$

3- sample $x^i \sim N(x; \mu_\theta(z), \Sigma_\theta(z))$

$P(X) = \sum_Z P_\theta(X,Z) = \sum_Z P(X|Z)\, P(Z)$

$= \sum_Z N(x|\underbrace{\mu(z)}_{\theta}, \underline{\Sigma_\theta(z)})\, N(z; 0, I)$

Very hard to compute the posterior!

$P(Z|\hat{X})$

Data $x^1, x^2, \underline{\quad\quad}, x^m$

20