

# Probabilistic Latent Semantic Analysis



Prithwijit Guha  
Dept. of EEE, IIT Guwahati

# Unsupervised Learning



Parametric  
Clustering  
Algorithms

Generic  
Clustering  
Algorithms

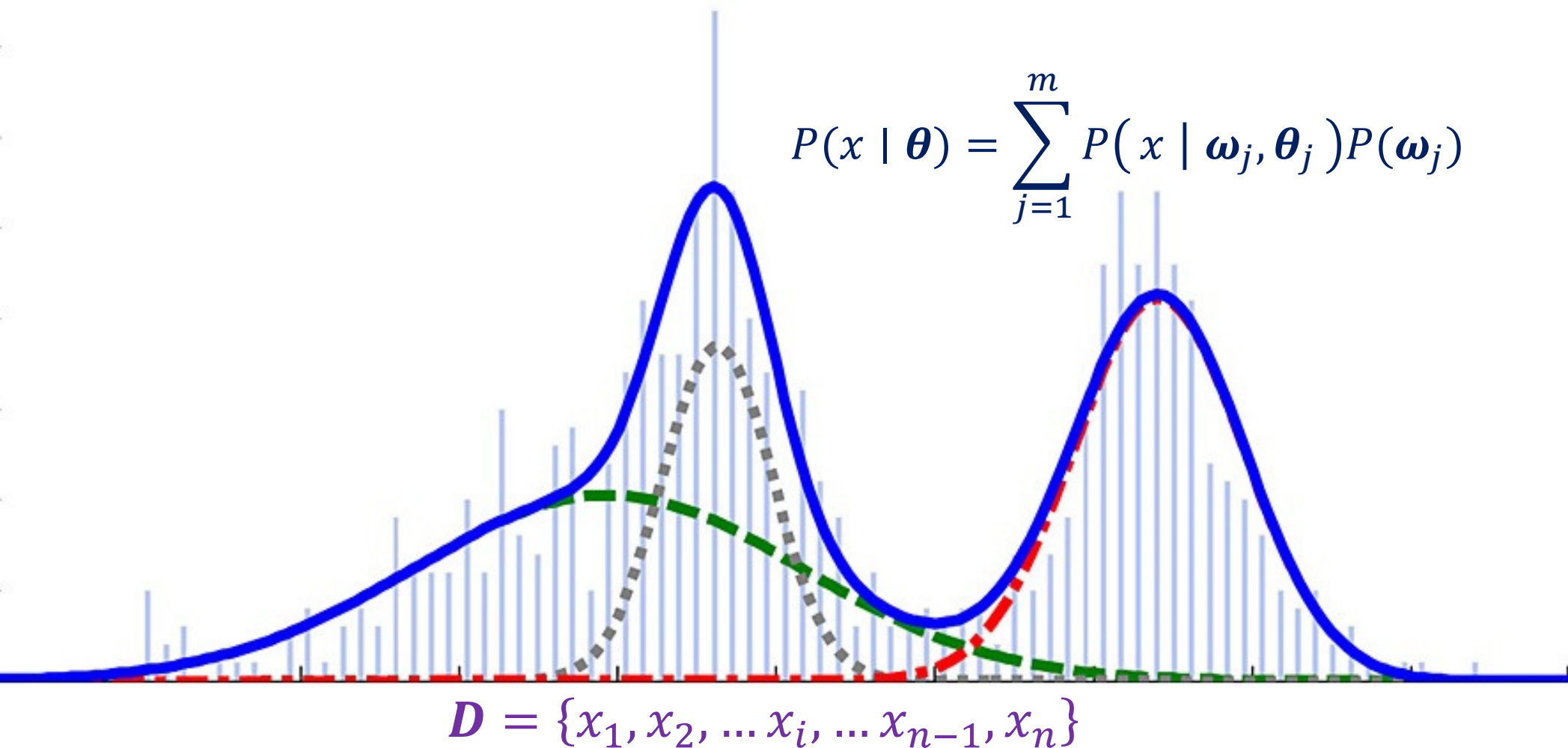
Estimation  
Theory

**Generative  
Models**

Pattern  
Mining

# Gaussian Mixture Models (GMM)

$$P(x | \boldsymbol{\theta}) = \sum_{j=1}^m P(x | \boldsymbol{\omega}_j, \boldsymbol{\theta}_j) P(\boldsymbol{\omega}_j)$$



## GMM: Mean & Variance Update

$$\hat{\mu}_r^{(t+1)} = \frac{\sum_{i=1}^n x_i P(\boldsymbol{\omega}_r \mid x_i, \boldsymbol{\theta}^{(t)})}{\sum_{i=1}^n P(\boldsymbol{\omega}_r \mid x_i, \boldsymbol{\theta}^{(t)})}$$

$$\hat{v}_r^{(t+1)} = \frac{\sum_{i=1}^n \left\{ x_i - \mu_r^{(t)} \right\}^2 P(\boldsymbol{\omega}_r \mid x_i, \boldsymbol{\theta}^{(t)})}{\sum_{i=1}^n P(\boldsymbol{\omega}_r \mid x_i, \boldsymbol{\theta}^{(t)})}$$

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

## GMM: Mean & Covariance Update

$$\hat{\boldsymbol{\mu}}_r^{(t+1)} = \frac{\sum_{i=1}^n \mathbf{x}_i P(\boldsymbol{\omega}_r \mid \mathbf{x}_i, \boldsymbol{\theta}^{(t)})}{\sum_{i=1}^n P(\boldsymbol{\omega}_r \mid \mathbf{x}_i, \boldsymbol{\theta}^{(t)})}$$

$$\hat{\mathbf{C}}_r^{(t+1)} = \frac{\sum_{i=1}^n \left\{ \mathbf{x}_i - \boldsymbol{\mu}_r^{(t)} \right\} \left\{ \mathbf{x}_i - \boldsymbol{\mu}_r^{(t)} \right\}^T P(\boldsymbol{\omega}_r \mid \mathbf{x}_i, \boldsymbol{\theta}^{(t)})}{\sum_{i=1}^n P(\boldsymbol{\omega}_r \mid \mathbf{x}_i, \boldsymbol{\theta}^{(t)})}$$

$$\mathbf{x} \in \mathbb{R}^d$$

# Expectation-Maximization (EM) Algorithm

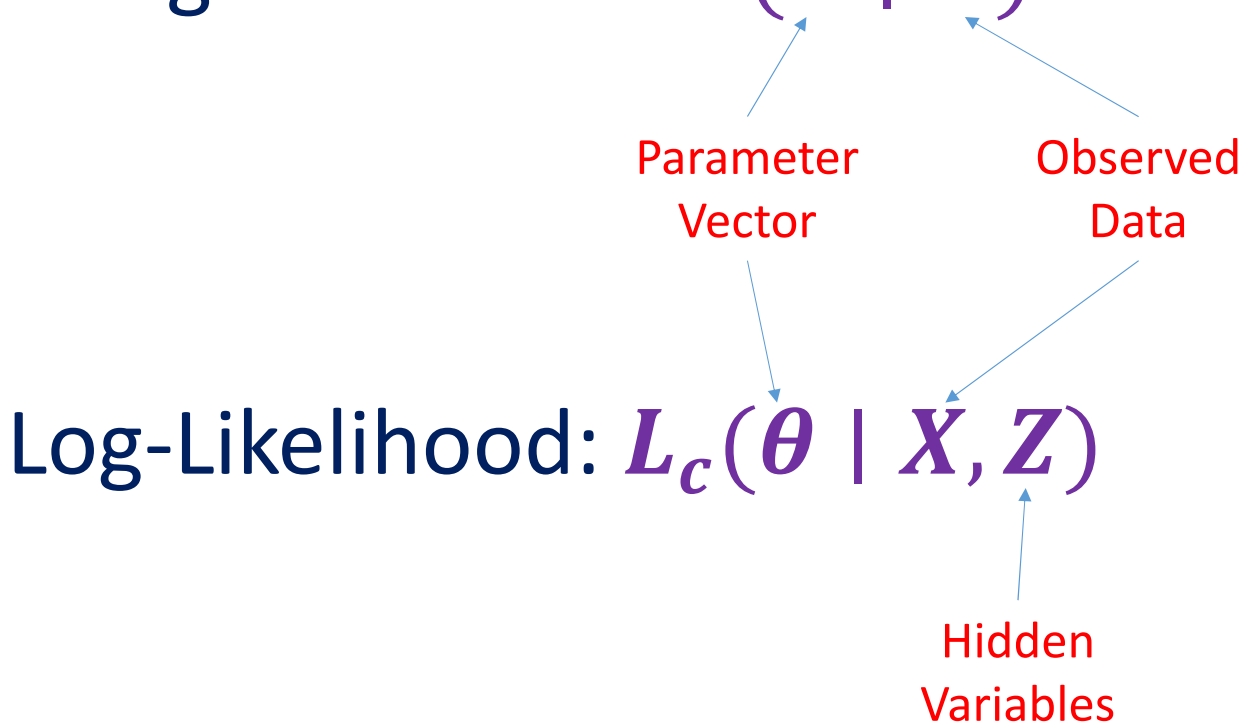
Log-Likelihood:  $L(\theta | X)$

Parameter  
Vector

Observed  
Data

Log-Likelihood:  $L_c(\theta | X, Z)$

Hidden  
Variables



# Expectation-Maximization (EM) Algorithm

$$\text{E-Step: } Q(\boldsymbol{\theta} \mid \boldsymbol{\theta}^{(t)}) = \mathbb{E}[L_c(\boldsymbol{\theta} \mid \boldsymbol{X}, \boldsymbol{Z}) \mid \boldsymbol{X}, \boldsymbol{\theta}^{(t)}]$$

$$\text{M-Step: } \boldsymbol{\theta}^{(t+1)} = \operatorname{argmax}_{\boldsymbol{\theta}} Q(\boldsymbol{\theta} \mid \boldsymbol{\theta}^{(t)})$$

**Chennai:** A video of government workers in Puducherry hurriedly throwing the body of a COVID-19 positive man into a pit has caused massive outrage, prompting the administration to order a probe into the incident.

The COVID pandemic has hit the world and the Vedanta Group business. It has incurred losses in oil, gas and mining sectors.

Yuvraj Singh was diagnosed with a cancerous tumor in his left lung following India's World Cup triumph in 2011. He had scored 362 runs and claimed 15 wickets in the tournament and was bestowed with the Player of the Tournament award in the end

Sonu Sood, known as a Bollywood actor, has emerged as the superhero during the COVID-19 pandemic. For thousands of migrants, he is the man who helped them at a time when they were gripped with fear and were walking an uncertain path with no support or even an assurance that all would be well.

Toyota Kirloskar Motor (TKM) has announced one or two percent price increase in India for Toyota Glanza, Yaris, Innova Crysta, and the Fortuner.

Facebook will invest Rs 43,574 crore in Jio Platforms, a unit of Reliance Industries Ltd (RIL), for a 9.99% stake, an allcash deal that will help the oil-to-retail conglomerate reduce debt and strengthen the social media company's presence in its largest market, especially for its WhatsApp unit.

COVID

Business

Bollywood

Health

Migrant Crisis

Cricket



# Word-Document Co-occurrence Matrix

$\mathbb{N} =$

	$d_1$	$d_2$	...	$d_m$	...	$d_{M-1}$	$d_M$
$w_1$							
$w_2$							
$\vdots$							
$w_n$				$\eta(d_m, w_n)$			
$\vdots$							
$w_{N-1}$							
$w_N$							

# Articles, Topics & Words

Documents:  $d_m; m = 1, \dots M$

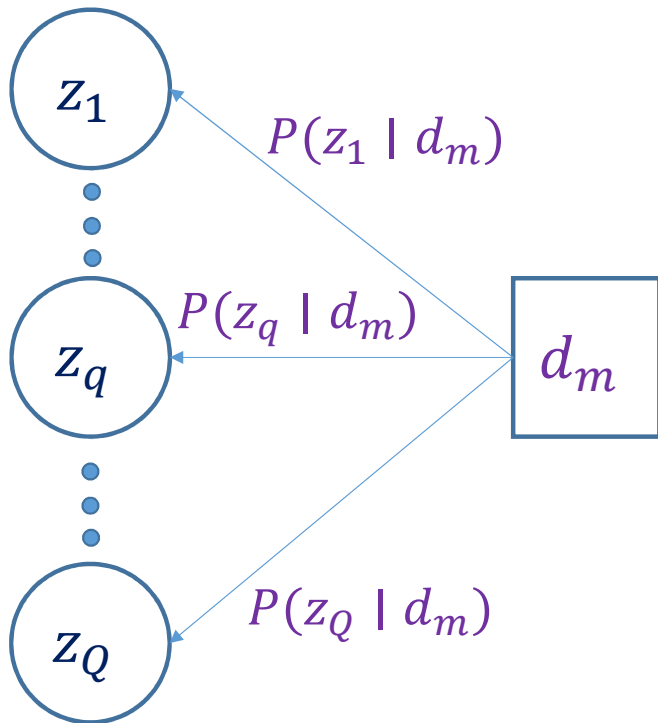
Words:  $w_n; n = 1, \dots N$

Topics:  $z_q; q = 1, \dots Q$

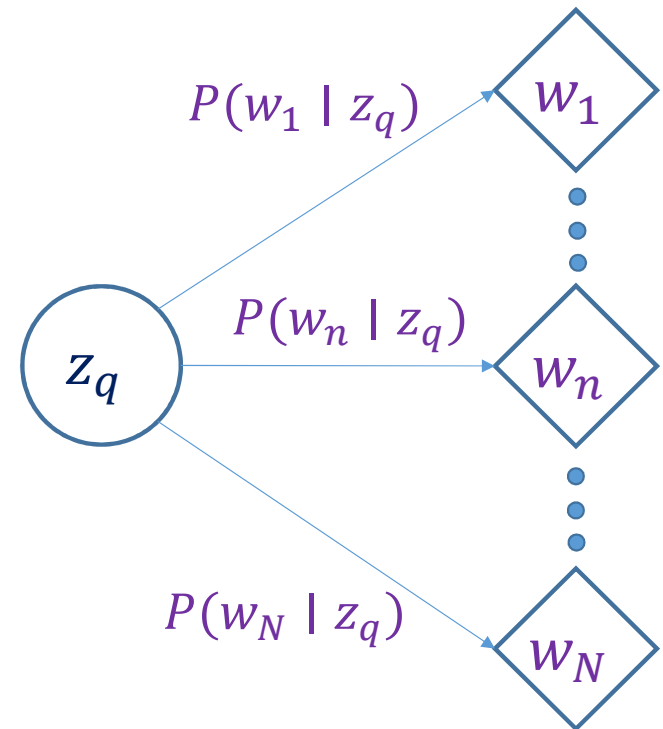
$$Q \ll M$$

$$Q \ll N$$

# Articles, Topics & Words

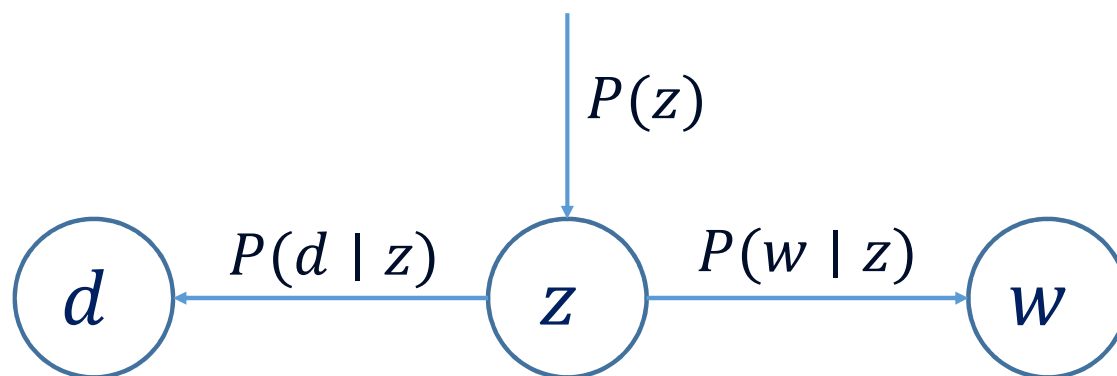


Document as a  
Mixture of Topics



Topics Identified  
by Words

# Aspect Model



$$P(d, w) = \sum_z P(z) P(d | z) P(w | z)$$

# Log-Likelihood Function

$$\mathbf{L}(\mathbf{z}) = \sum_{m=1}^M \sum_{n=1}^N \eta(d_m, w_n) \log P(d_m, w_n)$$

# Estimating $P(\mathbf{z} \mid \mathbf{d}, \mathbf{w})$

$$P(\underset{\text{Topic}}{z_q} \mid \underset{\text{Document}}{d_m}, \underset{\text{Word}}{w_n}) = \frac{P(z_q)P(d_m \mid z_q)P(w_n \mid z_q)}{\sum_{r=1}^Q P(z_r)P(d_m \mid z_r)P(w_n \mid z_r)}$$

Probability of Occurrence of Topic  $z_q$  for Given Word  $w_n$  in Document  $d_m$

# Estimating $P(\mathbf{w} \mid \mathbf{z})$

$$P(\underset{\text{Word}}{w_n} \mid \underset{\text{Topic}}{z_q}) = \frac{\sum_{m=1}^M \eta(d_m, w_n) P(z_q \mid d_m, w_n)}{\sum_{m=1}^M \sum_{i=1}^N \eta(d_m, w_i) P(z_q \mid d_m, w_i)}$$

Probability of Occurrence of Word  $w_n$  for Given Topic  $z_q$

Can be Used to Understand the Words that Describe a Topic

# Estimating $P(d \mid z)$

$$P(\underset{\text{Document}}{d_m} \mid \underset{\text{Topic}}{z_q}) = \frac{\sum_{n=1}^N \eta(d_m, w_n) P(z_q \mid d_m, w_n)}{\sum_{j=1}^M \sum_{n=1}^N \eta(d_j, w_n) P(z_q \mid d_j, w_n)}$$

Probability of Occurrence of Word  $w_n$  for Given Topic  $z_q$

Can be Used to Identify the Documents that Subscribe to a Topic



# Estimating $P(\mathbf{z})$

$$P(z_q) = \frac{\sum_{m=1}^M \sum_{n=1}^N \eta(d_m, w_n) P(z_q | d_m, w_n)}{\sum_{m=1}^M \sum_{n=1}^N \eta(d_m, w_n)}$$

The Distribution of Topics

# Estimating $P(\mathbf{z}|\mathbf{d})$

$$\overset{\text{Topic}}{P(z_q \mid d_m)} = \sum_{n=1}^N \underset{\text{Document}}{P(z_q \mid d_m, w_n)}$$

The Distribution of Topics in a Document

# Estimating $P(\mathbf{z}|\mathbf{w})$

$$P(\overset{\text{Topic}}{z_q} \mid \underset{\text{word}}{w_n}) = \sum_{m=1}^M P(z_q \mid d_m, w_n)$$

The Distribution of **Topics** with respect to a **Word**

# Summary

- EM Algorithms
- Aspect Model of PLSA
- Topic Distribution as Dimensionality Reduction
- E-M Equations of PLSA
- Interpreting Different Distributions



# Thank You