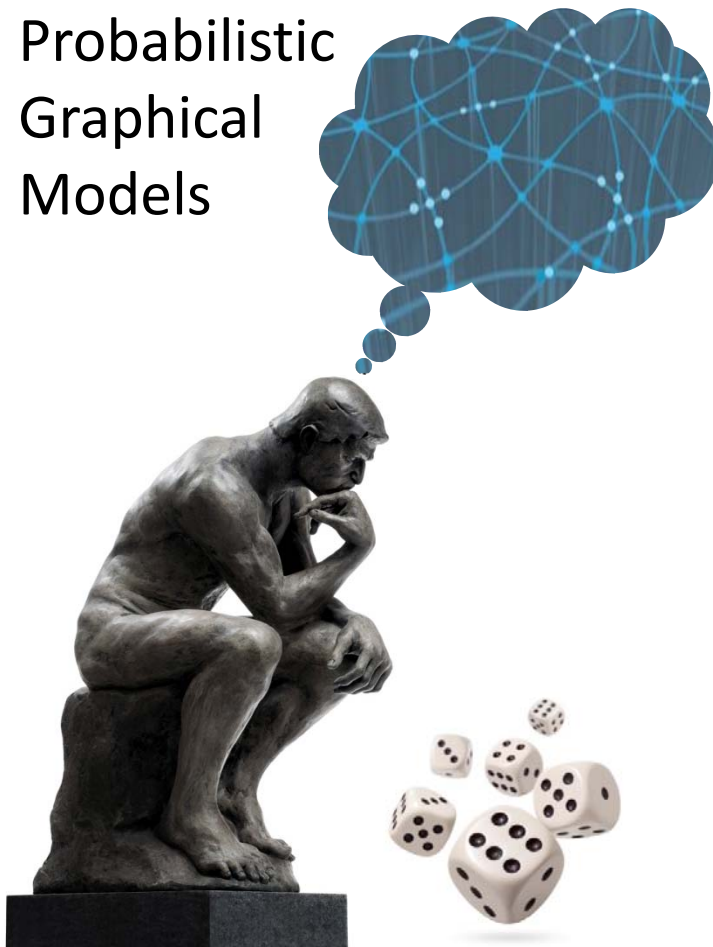


Probabilistic
Graphical
Models



Learning

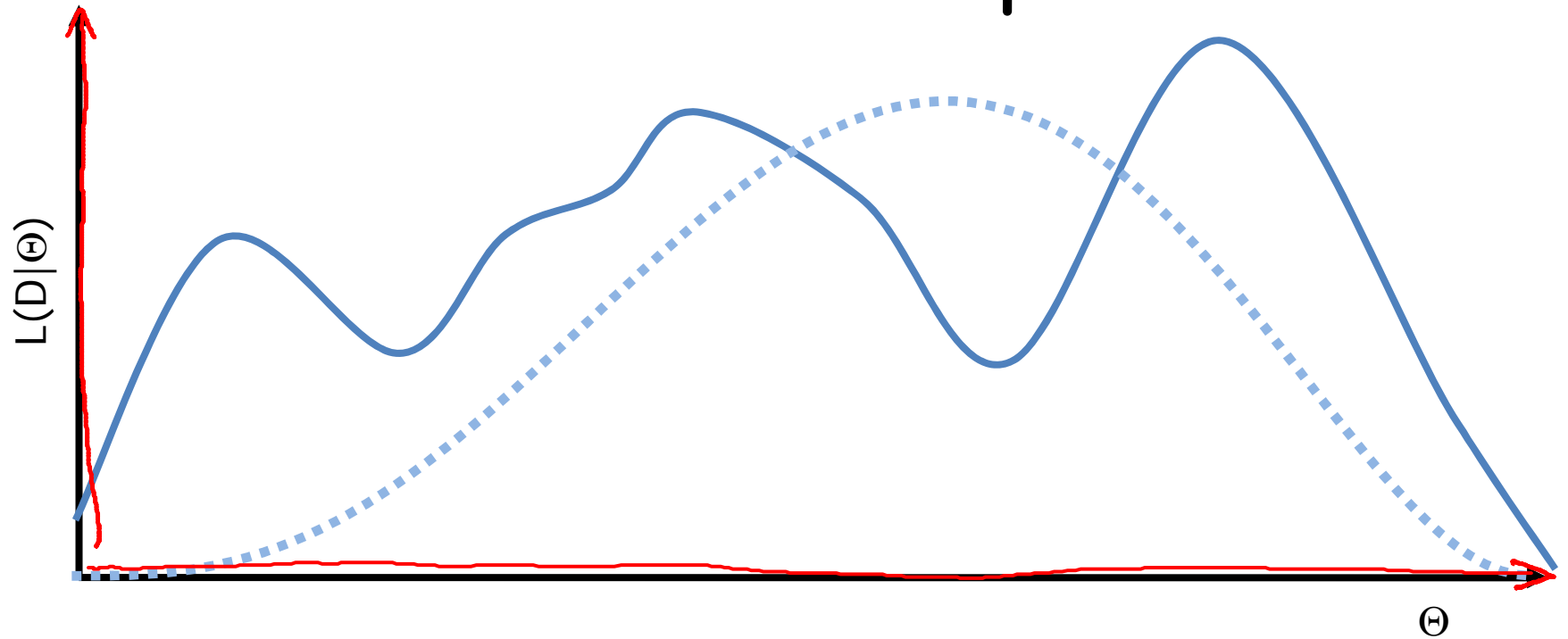
Incomplete Data

Likelihood

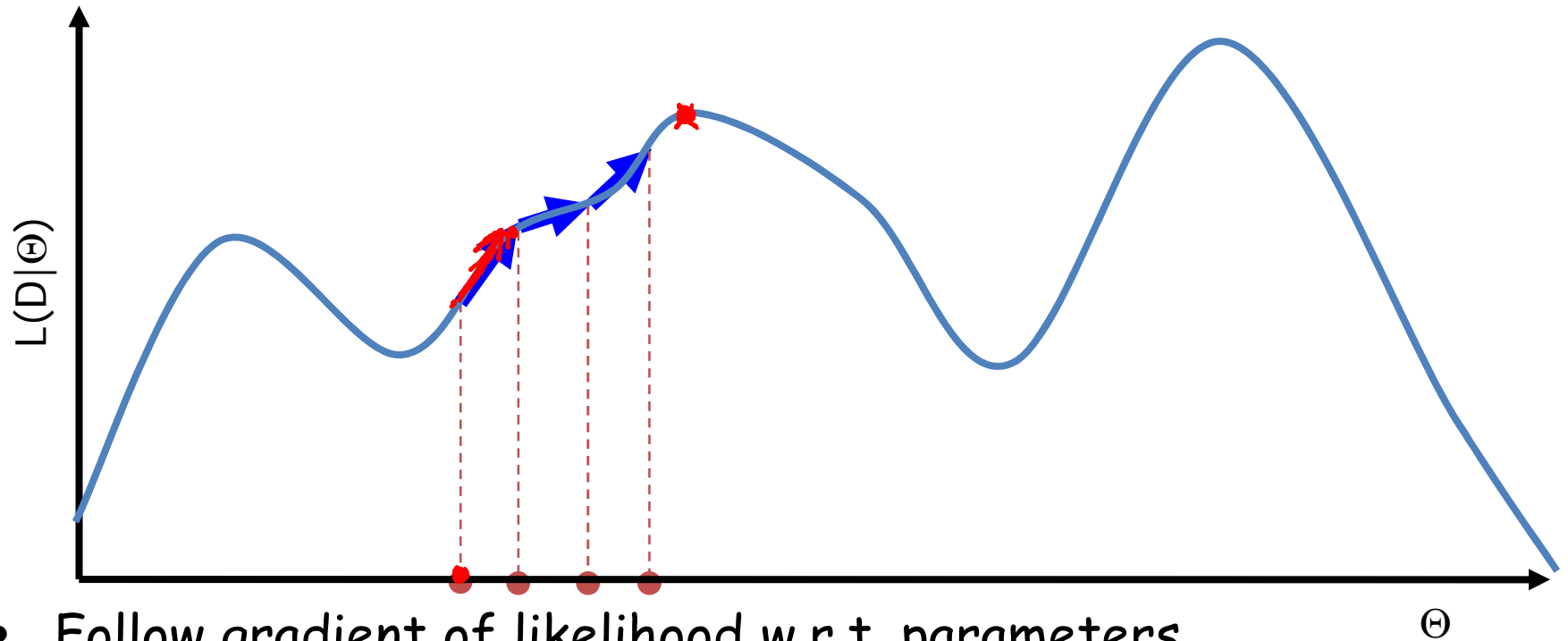
Optimization

Methods

Likelihood with Incomplete Data



Gradient Ascent



- Follow gradient of likelihood w.r.t. parameters
- Line search & conjugate gradient methods for fast convergence

Gradient Ascent

- Theorem:

$$\frac{\partial \log P(D | \Theta)}{\partial \theta_{x_i | u_i}} = \frac{1}{\theta_{x_i | u_i}} \sum_m P(\underbrace{x_i, u_i}_{\text{data instances } m} | \underbrace{d[m]}_{\text{evidence in } m^{\text{th}} \text{ instance}}, \underbrace{\Theta}_{\text{current param value}})$$

- Requires computing $P(\underline{X_i}, \underline{U_i} | \underline{d[m]}, \Theta)$ for all $\underline{i}, \underline{m}$
- Can be done with clique-tree algorithm, since $\underline{X_i}, \underline{U_i}$ are in the same clique

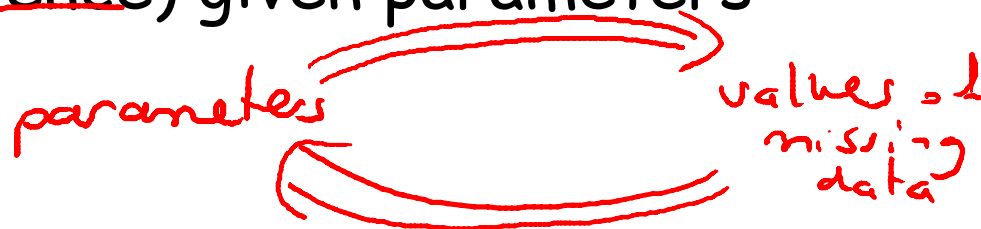
Gradient Ascent Summary

- Need to run inference over each data instance at every iteration
- Pros
 - Flexible, can be extended to non table CPDs
- Cons
 - Constrained optimization: need to ensure that parameters define legal CPDs
 - For reasonable convergence, need to combine with advanced methods (conjugate gradient, line search)

chain rule for derivatives

Expectation Maximization (EM)

- Special-purpose algorithm designed for optimizing likelihood functions
- **Intuition**
 - Parameter estimation is easy given complete data
 - Computing probability of missing data is “easy” (=inference) given parameters



EM Overview

- Pick a starting point for parameters
- Iterate:
 - E-step (Expectation): "Complete" the data using current parameters
 - M-step (Maximization): Estimate parameters relative to data completion
- Guaranteed to improve $L(\theta : D)$ at each iteration

Expectation Maximization (EM)

- Expectation (E-step):

- For each data case $\underline{d[m]}$ and each family $\underline{X, U}$ compute

- Compute the expected sufficient statistics for each x, u

soft completion

$$\frac{P(X, U | d[m], \theta^t)}{m(x, u)}$$

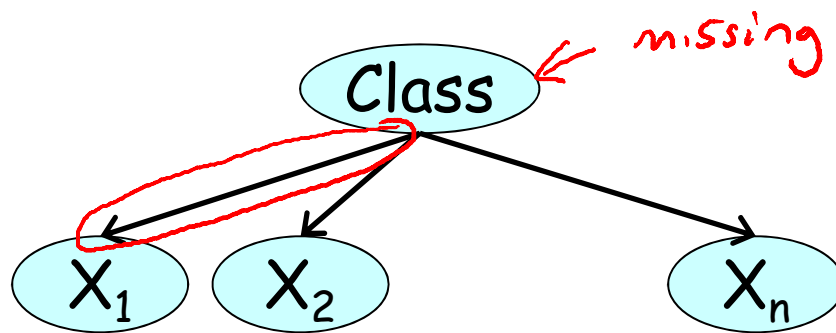
$$\bar{M}_{\theta^t}[x, u] = \sum_{m=1}^M P(x, u | d[m], \theta^t)$$

- Maximization (M-step):

- Treat the expected sufficient statistics (ESS) as if real
- Use MLE with respect to the ESS

$$\theta_{x|u}^{t+1} = \frac{\bar{M}_{\theta^t}[x, u]}{\bar{M}_{\theta^t}[u]}$$

Example: Bayesian Clustering



$$\begin{aligned}
 \bar{M}_{\theta}[c] &:= \sum_m P(c \mid \underline{x_1[m], \dots, x_n[m]}, \theta^t) & \theta_c^{t+1} &= \frac{\bar{M}_{\theta}[c]}{M} \\
 \bar{M}_{\theta}[x_i, c] &:= \sum_m P(c, x_i \mid \underline{x_1[m], \dots, x_n[m]}, \theta^t) & \theta_{x_i|c}^{t+1} &:= \frac{\bar{M}_{\theta}[x_i, c]}{\bar{M}_{\theta}[c]}
 \end{aligned}$$

EM Summary

- Need to run inference over each data instance at every iteration
- Pros
 - Easy to implement on top of MLE for complete data
 - Makes rapid progress, especially in early iterations
- Cons
 - Convergence slows down at later iterations