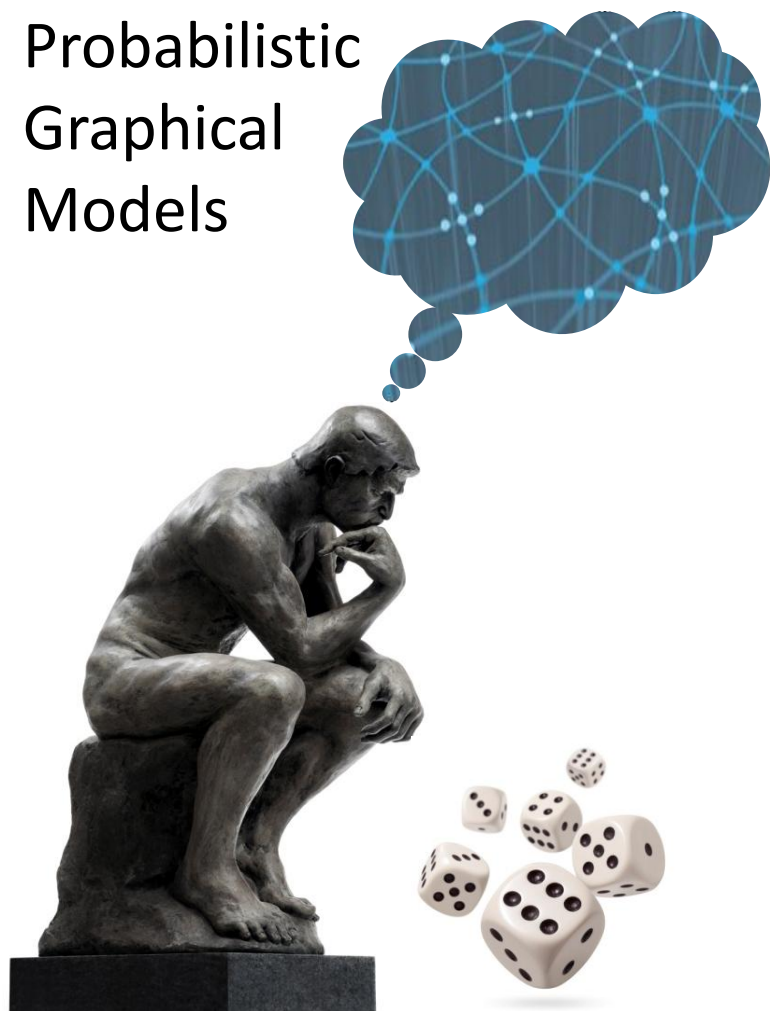


Probabilistic  
Graphical  
Models



Learning

---

Parameter Estimation

---

# Max Likelihood for Log-Linear Models

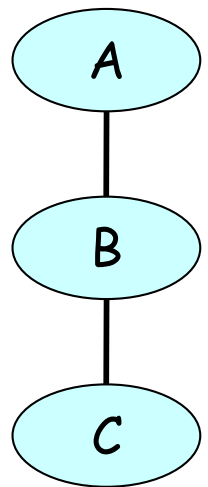
# Log-Likelihood for Markov Nets

$$P_{\theta}(a, b, c) = \frac{1}{Z} \phi_1(a, b) \cdot \phi_2(b, c)$$

$$\begin{aligned} \ell(\theta : \mathcal{D}) &= \sum_m (\ln \phi_1(a[m], b[m]) + \ln \phi_2(b[m], c[m]) - \ln Z(\theta)) \\ &= \sum_{a,b} M[a, b] \ln \phi_1(a, b) + \sum_{b,c} M[b, c] \ln \phi_2(b, c) - M \ln Z(\theta) \end{aligned}$$

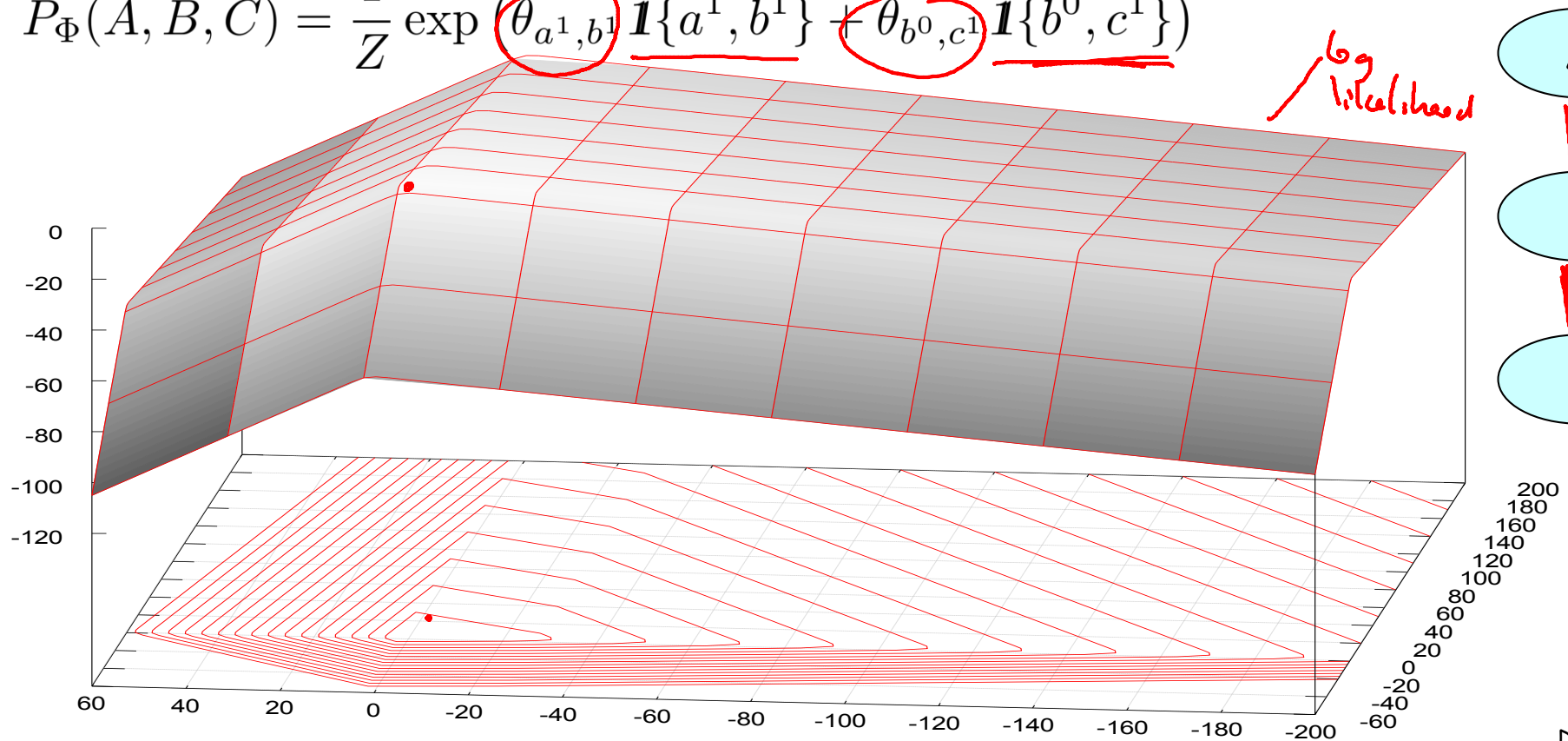
$$Z(\theta) = \sum_{a,b,c} \phi_1(a, b) \phi_2(b, c)$$

- Partition function couples the parameters
  - No decomposition of likelihood
  - No closed form solution



# Example: Log-Likelihood Function

$$P_{\Phi}(A, B, C) = \frac{1}{Z} \exp(\underbrace{\theta_{a^1, b^1}}_{\text{log likelihood}} \mathbf{1}\{a^1, b^1\} + \underbrace{\theta_{b^0, c^1}}_{\text{log likelihood}} \mathbf{1}\{b^0, c^1\})$$



# Log-Likelihood for Log-Linear Model

$$\underline{P(X_1, \dots, X_n : \theta)} = \frac{1}{Z(\theta)} \exp \left\{ \sum_{i=1}^k \theta_i f_i(\underline{D_i}) \right\}$$

*parameters* (pointing to  $\theta_i$ )  
*features* (pointing to  $f_i$ )  
*partition function* (pointing to  $Z(\theta)$ )

$$\ell(\theta : \mathcal{D}) = \sum_i \theta_i \left( \sum_m f_i(\underline{x[m]}) \right) - M \ln Z(\theta)$$

*feature  $f_i$  applied to the  $m$ th instance* (pointing to  $f_i(\underline{x[m]})$ )

$$\ln Z(\theta) = \ln \sum_{\underline{x}} \exp \left\{ \sum_i \theta_i f_i(\underline{x}) \right\}$$

*exponentially large space* (pointing to  $\underline{x}$ )  
*log-sum-exp* (pointing to the outer expression)

# The Log-Partition Function

Theorem:  $\frac{\partial}{\partial \theta_i} \ln Z(\theta) = E_{\theta}[f_i]$

*vector at derivative* *expectation of  $f_i$  relative to  $P_{\theta}$*

$\sum_x P_{\theta}(x) f_i(x)$

*matrix (Hessian)*  $\frac{\partial^2}{\partial \theta_i \partial \theta_j} \ln Z(\theta) = \text{Cov}_{\theta}[f_i; f_j]$

Proof:  $\frac{\partial}{\partial \theta_i} \ln Z(\theta) = \frac{1}{Z(\theta)} \sum_x \frac{\partial}{\partial \theta_i} \exp \left\{ \sum_j \theta_j f_j(x) \right\}$

$= \frac{1}{Z(\theta)} \sum_x f_i(x) \exp \left\{ \sum_j \theta_j f_j(x) \right\}$

$= \sum_x \left[ \frac{1}{Z(\theta)} \exp \left\{ \sum_j \theta_j f_j(x) \right\} \right] f_i(x) = \sum_x P_{\theta}(x) f_i(x)$

$\frac{\partial}{\partial \theta_i} \theta_j f_j = \begin{cases} 0 & i \neq j \\ f_i & i = j \end{cases}$

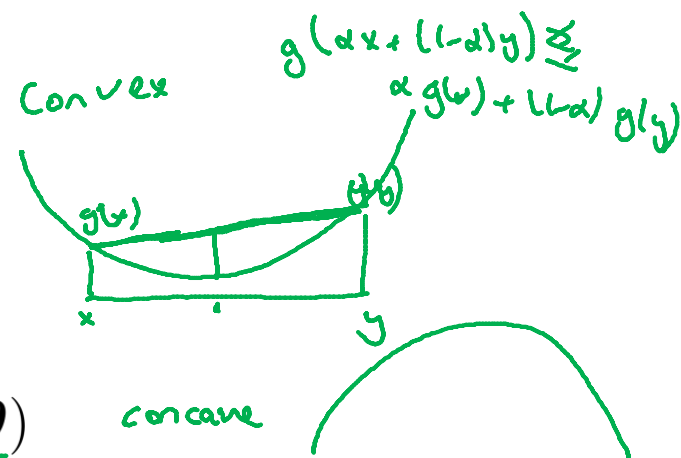
# The Log-Partition Function

Theorem:  $\frac{\partial}{\partial \theta_i} \ln Z(\theta) = E_{\theta}[f_i]$

Hessian:  $\frac{\partial^2}{\partial \theta_i \partial \theta_j} \ln Z(\theta) = \text{Cov}_{\theta}[f_i; f_j]$

$$\ell(\theta : \mathcal{D}) = \sum_i \theta_i \left( \sum_m f_i(\mathbf{x}[m]) \right) - \underline{M \ln Z(\theta)}$$

- Log likelihood function
  - No local optima
  - Easy to optimize



# Maximum Likelihood Estimation

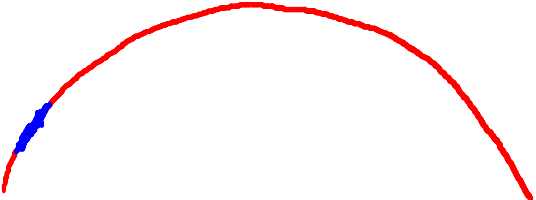
$$\frac{1}{M} \ell(\boldsymbol{\theta} : \mathcal{D}) = \sum_i \theta_i \underbrace{\left( \frac{1}{M} \sum_m f_i(x[m]) \right)}_{\text{empirical expectation of } f_i \text{ in } \mathcal{D}} - \underbrace{\ln Z(\boldsymbol{\theta})}_{\text{expectation of } f_i \text{ in } p_\theta}$$

$$\frac{\partial}{\partial \theta_i} \frac{1}{M} \ell(\boldsymbol{\theta} : \mathcal{D}) = \mathbf{E}_{\mathcal{D}}[f_i(\mathbf{X})] - \mathbf{E}_{\boldsymbol{\theta}}[f_i]$$

Theorem:  $\hat{\boldsymbol{\theta}}$  is the MLE if and only if

$$\forall i \quad \underbrace{\mathbf{E}_{\mathcal{D}}[f_i(\mathbf{X})]}_{\text{expectation in } \mathcal{D}} = \underbrace{\mathbf{E}_{\hat{\boldsymbol{\theta}}}[f_i]}_{\text{expectation relative to } \hat{\boldsymbol{\theta}}}$$

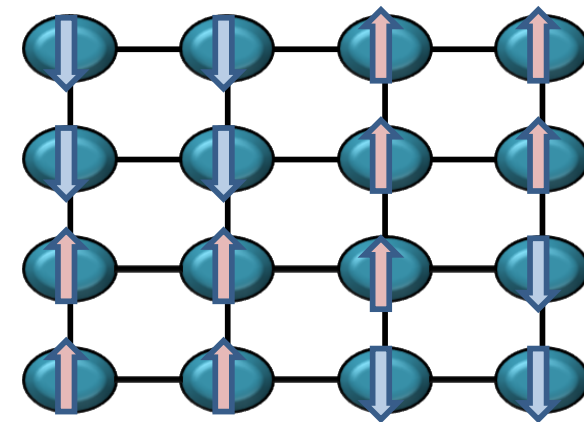
# Computation: Gradient Ascent

$$\frac{\partial}{\partial \theta_i} \frac{1}{M} \ell(\theta : \mathcal{D}) = \underbrace{E_{\mathcal{D}}[f_i(\mathbf{X})]}_{\text{in data}} - \underbrace{E_{\theta}[f_i]}_{\text{relative to current model}}$$


- Use gradient ascent:
  - typically L-BFGS - a quasi-Newton method
- For gradient, need expected feature counts:
  - in data
  - relative to current model
- Requires inference at each gradient step



# Example: Ising Model



$$x_i \in \{-1, +1\}$$

$$E(x_1, \dots, x_n) = - \sum_{i < j} w_{i,j} x_i x_j - \sum_i u_i x_i$$

$$\frac{\partial}{\partial \theta_i} \frac{1}{M} \ell(\theta : \mathcal{D}) = \mathbf{E}_{\mathcal{D}}[f_i(\mathbf{X})] - \mathbf{E}_{\theta}[f_i]$$

$$\frac{\partial}{\partial u_i} = \frac{1}{M} \sum_m \underbrace{x_i[m]}_{+1} - (P_{\theta}(X_i = 1) - P_{\theta}(X_i = -1))$$

$$\frac{\partial}{\partial w_{ij}} = \frac{1}{M} \sum_m \underbrace{x_i[m] x_j[m]}_{+1} - \left( \frac{P_{\theta}(X_i = 1, X_j = 1) + P_{\theta}(X_i = -1, X_j = -1)}{-P_{\theta}(X_i = 1, X_j = -1) - P_{\theta}(X_i = -1, X_j = -1)} \right)$$

# Summary

- Partition function couples parameters in likelihood
- No closed form solution, but convex optimization
  - Solved using gradient ascent (usually L-BFGS) *global opt.*
- Gradient computation requires inference at each gradient step to compute expected feature counts
- Features are always within clusters in cluster-graph or clique tree due to family preservation
  - One calibration suffices for all feature expectations