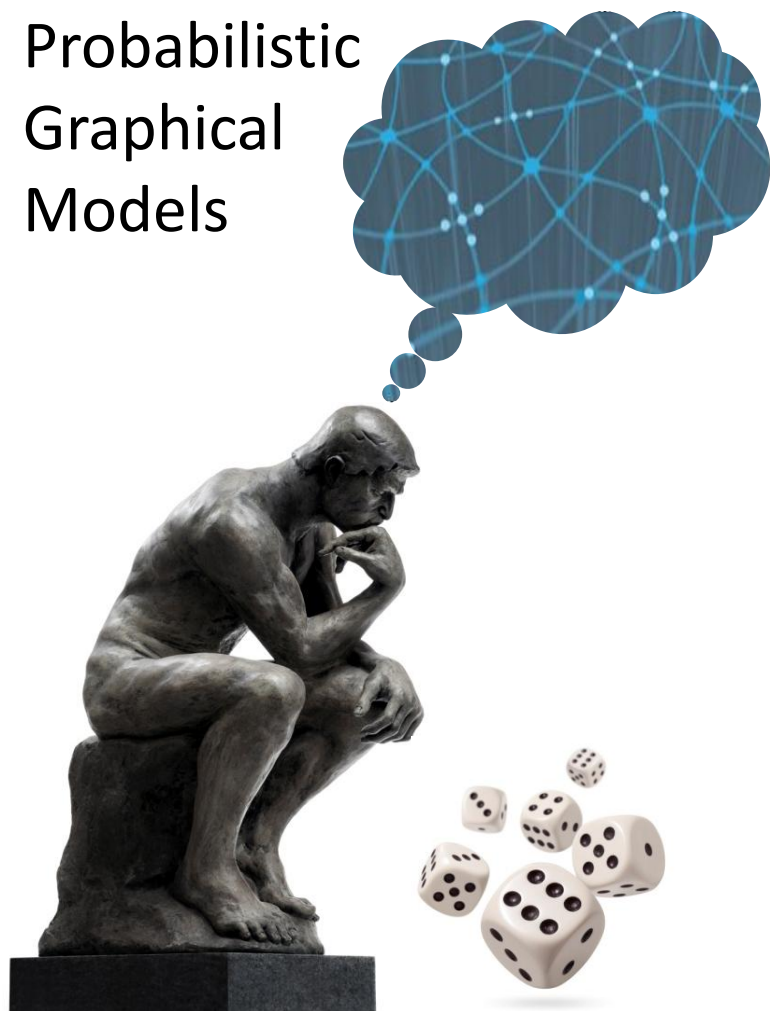


Probabilistic
Graphical
Models



Learning

Parameter Estimation

Max Likelihood for CRFs

Estimation for CRFs

$$P_{\theta}(\mathbf{Y} | \mathbf{x}) = \frac{1}{Z_{\mathbf{x}}(\theta)} \tilde{P}_{\theta}(\mathbf{x}, \mathbf{Y})$$

$$Z_{\mathbf{x}}(\theta) = \sum_{\mathbf{Y}} \tilde{P}_{\theta}(\mathbf{x}, \mathbf{Y})$$

log conditional likelihoods

$$\mathcal{D} = \{(\mathbf{x}[m], \mathbf{y}[m])\}_{m=1}^M$$

$$\ell_{\mathbf{Y}|\mathbf{X}}(\theta : \mathcal{D}) = \sum_{m=1}^M \ln P_{\theta}(\mathbf{y}[m] | \mathbf{x}[m], \theta)$$

$$\ell_{\mathbf{Y}|\mathbf{X}}(\theta : (\mathbf{x}[m], \mathbf{y}[m])) = \left(\sum_i \theta_i f_i(\mathbf{x}[m], \mathbf{y}[m]) \right) - \ln Z_{\mathbf{x}[m]}(\theta)$$

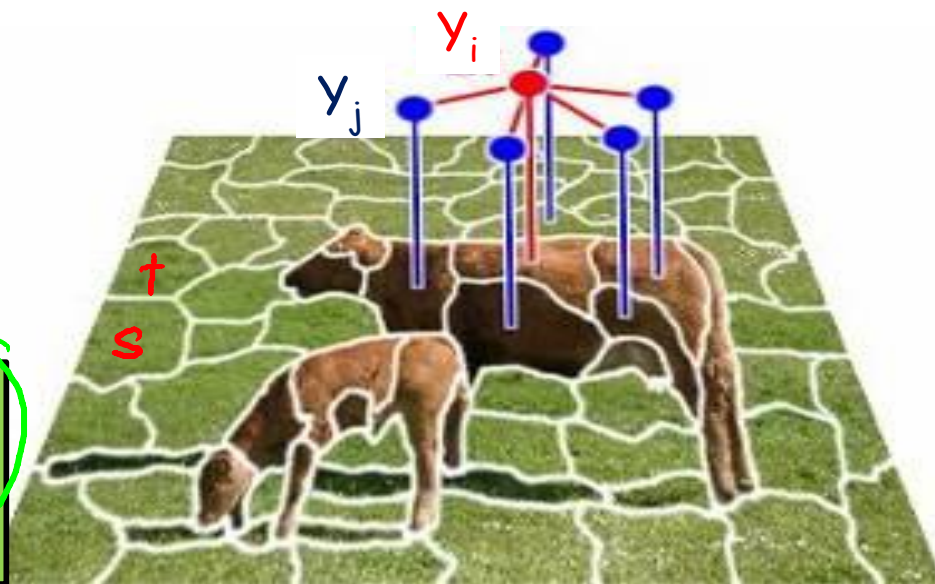
$$\frac{\partial}{\partial \theta_i} \frac{1}{M} \ell_{\mathbf{Y}|\mathbf{X}}(\theta : \mathcal{D}) = \frac{1}{M} \sum_{m=1}^M (f_i(\mathbf{x}[m], \mathbf{y}[m]) - E_{\theta}[f_i(\mathbf{x}[m], \mathbf{Y})])$$

Example

$$\underline{f_1(Y_s, X_s) = \underline{1(Y_s = g)} \times G_s}$$

$$\underline{f_2(Y_s, Y_t) = \underline{1(Y_s = Y_t)}}$$

average intensity of
green channel for
pixels in superpixel s



$$\frac{\partial}{\partial \theta_i} \ell_{\mathbf{Y}|\mathbf{X}}(\boldsymbol{\theta} : (\mathbf{x}[m], \mathbf{y}[m])) = (f_i(\mathbf{x}[m], \mathbf{y}[m]) - \mathbf{E}_{\boldsymbol{\theta}}[f_i(\mathbf{x}[m], \mathbf{Y})])$$

$$\frac{\partial}{\partial \theta_1} = \sum_s \underline{1\{y_s[m] = g\}} \underline{G_s[m]} - \sum_s \underline{P_{\boldsymbol{\theta}}(Y_s = g | \mathbf{x}[m])} \underline{G_s[m]}$$

$$\frac{\partial}{\partial \theta_2} = \sum_{(s,t) \in \mathcal{N}} \underline{1\{y_s[m] = y_t[m]\}} - \sum_{(s,t) \in \mathcal{N}} \underline{P_{\boldsymbol{\theta}}(Y_s = Y_t | \mathbf{x}[m])}$$

Computation

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

- Requires inference at each gradient step

CRF $\frac{\partial}{\partial \theta_i} \frac{1}{M} \ell_{\mathbf{Y}|\mathbf{X}}(\theta : \mathcal{D}) = \frac{1}{M} \sum_{m=1}^M (f_i(\underline{\mathbf{x}[m]}, \mathbf{y}[m]) - \mathbf{E}_{\theta}[f_i(\underline{\mathbf{x}[m]}, \mathbf{Y})])$

- Requires inference for each $\mathbf{x}[m]$ at each gradient step
 $M = \# \text{ training instances}$

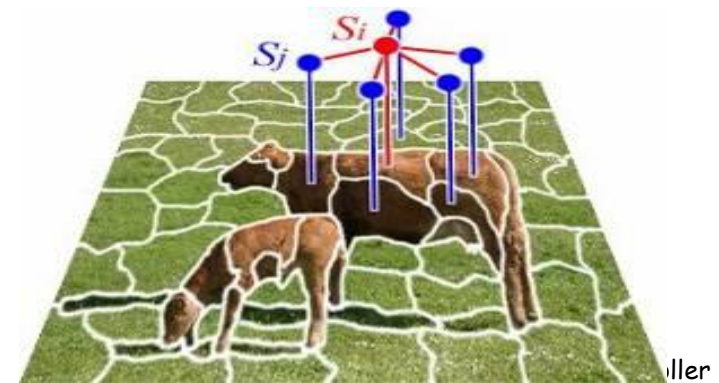
However...

- For inference of $P(Y | x)$, we need to compute distribution only over Y
- If we learn an MRF, need to compute $P(Y, X)$, which may be much more complex

$$f_1(Y_s, X_s) = \mathbf{1}(Y_s = g) \times G_s$$

$$f_2(Y_s, Y_t) = \mathbf{1}(Y_s = Y_t)$$

average intensity of
green channel for
pixels in superpixel i



Summary

- CRF learning very similar to MRF learning
 - Likelihood function is concave
 - Optimized using gradient ascent (usually L-BFGS)
- Gradient computation requires inference: one per gradient step, data instance
 - c.f., once per gradient step for MRFs
- But conditional model is often much simpler, so inference cost for CRF, MRF is not the same