

Graphical Models

Lecture 17:

Learning in Undirected Graphical Models

Andrew McCallum
mccallum@cs.umass.edu

Thanks to Noah Smith and Carlos Guestrin for some slide materials.

Learning

- learning input:
 - Graphical model with unknown parameters
 - Observations of variables (training data)
- learning output:
 - Parameters
- Maximum likelihood estimation (MLE):
 - Choose parameters that give highest probability to observed training data

Parameter Estimation in Bayesian Networks: Decomposability

$$\begin{aligned}\theta_{\text{MLE}} &= \arg \max_{\theta} \prod_t P(\mathbf{X} = \mathbf{x}^{(t)} \mid \theta) \\ &= \arg \max_{\theta} \prod_t \prod_i P(X_i = x_i^{(t)} \mid \text{Parents}(X_i) = \text{Parents}(x_i), \theta) \\ &= \arg \max_{\theta} \sum_t \sum_i \log P(X_i = x_i^{(t)} \mid \text{Parents}(X_i) = \text{Parents}(x_i), \theta)\end{aligned}$$

If the parameters θ are partitioned by CPT ...

$$= \arg \max_{\theta} \sum_i \sum_t \log P(X_i = x_i^{(t)} \mid \text{Parents}(X_i) = \text{Parents}(x_i), \theta_i)$$

Key Idea

- For known structure and fully observed data, MLE for a Bayesian network whose CPDs have disjoint parameters

equates to

MLE for each of its CPDs.

Bad News for Markov Networks

- The global normalization constant (Z) kills decomposability.

$$\begin{aligned}\theta_{\text{MLE}} &= \arg \max_{\theta} \prod_t P(\mathbf{X} = \mathbf{x}^{(t)} \mid \theta) \\ &= \arg \max_{\theta} \prod_t \frac{1}{Z} \prod_i \phi_i(\mathbf{x}_i^{(t)}) \\ &= \arg \max_{\theta} \left(\sum_t \sum_i \log \phi_i(\mathbf{x}_i^{(t)}) \right) - T \log Z\end{aligned}$$

- Solving for the parameters becomes more complicated.

Example Task: Entity Recognition

- “The outcome will help determine whether **Mr. Boehner** is leading his party or following the demands of the **Tea Party**.”

What are the Parameters?

- How do the factors ϕ get expressed as parameters θ ?
- Often, we adopt a log-linear parameterization.
- We covered this in lecture 9.

Log-Linear Markov Networks

- A **feature** is a function $f : \text{Val}(\mathbf{D}_i) \rightarrow \mathbb{R}$.
- Log-linear model:
$$\begin{aligned}P(\mathbf{X}) &= \frac{1}{Z} e^{\sum_i \log \phi_i(\mathbf{D}_i)} \\ &= \frac{1}{Z} e^{-\sum_i \psi_i(\mathbf{D}_i)} \\ &= \frac{1}{Z} e^{\sum_i \sum_j f_j(\mathbf{D}_i) w_j}\end{aligned}$$
- Features and weights can be *reused* for different factors.
 - Typical: features designed by expert, weights learned from data.
 - Note that this breaks parameter independence.

Log-Linear Markov Networks

- A **feature** is a function $f : \text{Val}(\mathbf{D}_i) \rightarrow \mathbb{R}$.
- Log-linear model:
$$\begin{aligned}P(\mathbf{X}) &= \frac{1}{Z} e^{\sum_i \log \phi_i(\mathbf{D}_i)} \\ &= \frac{1}{Z} e^{-\sum_i \psi_i(\mathbf{D}_i)} \\ &= \frac{1}{Z} e^{\sum_i \sum_j f_j(\mathbf{D}_i) w_j}\end{aligned}$$
- Log of the probability is *linear* in the weights **w**.
 - Ignoring Z , which is a constant for a given **w**.

Log-Likelihood Function for Log-Linear Models

$$\begin{aligned}\boldsymbol{\theta}_{\text{MLE}} &= \arg \max_{\boldsymbol{\theta}} \left(\sum_t \sum_i \log \phi_i(\mathbf{x}_i^{(t)}) \right) - T \log Z \\ &= \arg \max_{\mathbf{w}} \left(\sum_t \sum_i \sum_j w_j f_j(\mathbf{x}_i^{(t)}) \right) - T \log Z\end{aligned}$$

- The first term is linear in \mathbf{w} .
- The second term is also a function of \mathbf{w} :

$$\log Z = \log \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right)$$

Log-Likelihood Function for Log-Linear Models

- $\log Z$ does not decompose.
 - No closed form solution.
 - Even *computing* the likelihood requires inference!
- It is, however, *concave*.
- The weights \mathbf{w} are *unconstrained*.

$$\log Z = \log \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right)$$

Optimization Returns

- We talked about two abstract optimization problems last time:
 - Integer linear programming (NP hard)
 - Linear programming (solvable in poly time)
- **Convex** optimization: globally concave or globally convex function
- **Unconstrained** optimization: $\mathbf{w} \in \mathbb{R}^d$

Solving Unconstrained Convex Optimization Problems

- Gradient descent and variations
 - Stochastic gradient descent
 - Coordinate descent
 - Conjugate gradient descent
 - Newton, Quasi-Newton methods
- Specialized algorithms
 - For Markov networks, iterative proportional fitting, a.k.a. iterative scaling.

The Gradient of Log-Likelihood

$$\begin{aligned} & \frac{\partial}{\partial w_k} \left[\left(\sum_t \sum_i \sum_j w_j f_j(\mathbf{x}_i^{(t)}) \right) - T \log Z \right] \\ = & \frac{\partial}{\partial w_k} \left[\left(\sum_t \sum_i \sum_j w_j f_j(\mathbf{x}_i^{(t)}) \right) - T \log \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \right] \end{aligned}$$

Flesh out Z.

$$\begin{aligned} & = \sum_t \sum_i f_k(\mathbf{x}_i^{(t)}) - T \frac{\partial}{\partial w_k} \left[\log \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \right] \\ & = \sum_t \sum_i f_k(\mathbf{x}_i^{(t)}) - T \frac{\frac{\partial}{\partial w_k} \left[\sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \right]}{Z} \\ & = \sum_t \sum_i f_k(\mathbf{x}_i^{(t)}) - \frac{T}{Z} \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \frac{\partial}{\partial w_k} \left[\exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \right] \\ & = \sum_t \sum_i f_k(\mathbf{x}_i^{(t)}) - \frac{T}{Z} \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \left(\sum_i f_k(\mathbf{x}_i) \right) \\ & = T \left(\frac{\sum_t \sum_i f_k(\mathbf{x}_i^{(t)})}{T} - \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \frac{\exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right)}{Z} \left(\sum_i f_k(\mathbf{x}_i) \right) \right) \\ & = T \left(\mathbb{E}_{\tilde{P}} \left[\sum_i f_k(\mathbf{x}_i) \right] - \mathbb{E}_{P_w} \left[\sum_i f_k(\mathbf{x}_i) \right] \right) \end{aligned}$$

The Gradient of Log-Likelihood

$$\begin{aligned} & \frac{\partial}{\partial w_k} \left[\left(\sum_t \sum_i \sum_j w_j f_j(\mathbf{x}_i^{(t)}) \right) - T \log Z \right] \\ = & \frac{\partial}{\partial w_k} \left[\left(\sum_t \sum_i \sum_j w_j f_j(\mathbf{x}_i^{(t)}) \right) - T \log \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \right] \\ = & \sum_t \sum_i f_k(\mathbf{x}_i^{(t)}) - T \frac{\partial}{\partial w_k} \left[\log \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \right] \end{aligned}$$

Linear function is easy.

T is constant with respect to \mathbf{w} .

$$\begin{aligned} = & \sum_t \sum_i f_k(\mathbf{x}_i^{(t)}) - T \frac{\frac{\partial}{\partial w_k} \left[\sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \right]}{Z} \\ = & \sum_t \sum_i f_k(\mathbf{x}_i^{(t)}) - \frac{T}{Z} \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \frac{\partial}{\partial w_k} \left[\exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \right] \\ = & \sum_t \sum_i f_k(\mathbf{x}_i^{(t)}) - \frac{T}{Z} \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \left(\sum_i f_k(\mathbf{x}_i) \right) \\ = & T \left(\frac{\sum_t \sum_i f_k(\mathbf{x}_i^{(t)})}{T} - \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \frac{\exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right)}{Z} \left(\sum_i f_k(\mathbf{x}_i) \right) \right) \\ = & T \left(\mathbb{E}_{\tilde{P}} \left[\sum_i f_k(\mathbf{x}_i) \right] - \mathbb{E}_{P_{\mathbf{w}}} \left[\sum_i f_k(\mathbf{x}_i) \right] \right) \end{aligned}$$

The Gradient of Log-Likelihood

$$\begin{aligned}
 & \frac{\partial}{\partial w_k} \left[\left(\sum_t \sum_i \sum_j w_j f_j(\mathbf{x}_i^{(t)}) \right) - T \log Z \right] \\
 = & \frac{\partial}{\partial w_k} \left[\left(\sum_t \sum_i \sum_j w_j f_j(\mathbf{x}_i^{(t)}) \right) - T \log \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \right] \\
 = & \sum_t \sum_i f_k(\mathbf{x}_i^{(t)}) - T \frac{\partial}{\partial w_k} \left[\log \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \right] \\
 = & \sum_t \sum_i f_k(\mathbf{x}_i^{(t)}) - T \frac{\frac{\partial}{\partial w_k} \left[\sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \right]}{Z}
 \end{aligned}$$

Logarithm rule for derivatives.

Use “Z” shorthand.

$$\begin{aligned}
 & = \sum_t \sum_i f_k(\mathbf{x}_i^{(t)}) - \frac{T}{Z} \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \frac{\partial}{\partial w_k} \left[\exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \right] \\
 & = \sum_t \sum_i f_k(\mathbf{x}_i^{(t)}) - \frac{T}{Z} \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \left(\sum_i f_k(\mathbf{x}_i) \right) \\
 & = T \left(\frac{\sum_t \sum_i f_k(\mathbf{x}_i^{(t)})}{T} - \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \frac{\exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right)}{Z} \left(\sum_i f_k(\mathbf{x}_i) \right) \right) \\
 & = T \left(\mathbb{E}_{\tilde{P}} \left[\sum_i f_k(\mathbf{x}_i) \right] - \mathbb{E}_{P_w} \left[\sum_i f_k(\mathbf{x}_i) \right] \right)
 \end{aligned}$$

The Gradient of Log-Likelihood

$$\begin{aligned}
 & \frac{\partial}{\partial w_k} \left[\left(\sum_t \sum_i \sum_j w_j f_j(\mathbf{x}_i^{(t)}) \right) - T \log Z \right] \\
 = & \frac{\partial}{\partial w_k} \left[\left(\sum_t \sum_i \sum_j w_j f_j(\mathbf{x}_i^{(t)}) \right) - T \log \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \right] \\
 = & \sum_t \sum_i f_k(\mathbf{x}_i^{(t)}) - T \frac{\partial}{\partial w_k} \left[\log \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \right] \\
 = & \sum_t \sum_i f_k(\mathbf{x}_i^{(t)}) - T \frac{\frac{\partial}{\partial w_k} \left[\sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \right]}{Z} \\
 = & \sum_t \sum_i f_k(\mathbf{x}_i^{(t)}) - \frac{T}{Z} \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \frac{\partial}{\partial w_k} \left[\exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \right]
 \end{aligned}$$

Use the sum rule.

$$\begin{aligned}
 = & \sum_t \sum_i f_k(\mathbf{x}_i^{(t)}) - \frac{T}{Z} \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \left(\sum_i f_k(\mathbf{x}_i) \right) \\
 = & T \left(\frac{\sum_t \sum_i f_k(\mathbf{x}_i^{(t)})}{T} - \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \frac{\exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right)}{Z} \left(\sum_i f_k(\mathbf{x}_i) \right) \right) \\
 = & T \left(\mathbb{E}_{\tilde{P}} \left[\sum_i f_k(\mathbf{x}_i) \right] - \mathbb{E}_{P_w} \left[\sum_i f_k(\mathbf{x}_i) \right] \right)
 \end{aligned}$$

The Gradient of Log-Likelihood

$$\begin{aligned}
 & \frac{\partial}{\partial w_k} \left[\left(\sum_t \sum_i \sum_j w_j f_j(\mathbf{x}_i^{(t)}) \right) - T \log Z \right] \\
 = & \frac{\partial}{\partial w_k} \left[\left(\sum_t \sum_i \sum_j w_j f_j(\mathbf{x}_i^{(t)}) \right) - T \log \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \right] \\
 = & \sum_t \sum_i f_k(\mathbf{x}_i^{(t)}) - T \frac{\partial}{\partial w_k} \left[\log \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \right] \\
 = & \sum_t \sum_i f_k(\mathbf{x}_i^{(t)}) - T \frac{\frac{\partial}{\partial w_k} \left[\sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \right]}{Z} \\
 = & \sum_t \sum_i f_k(\mathbf{x}_i^{(t)}) - \frac{T}{Z} \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \frac{\partial}{\partial w_k} \left[\exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \right] \\
 = & \sum_t \sum_i f_k(\mathbf{x}_i^{(t)}) - \frac{T}{Z} \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \left(\sum_i f_k(\mathbf{x}_i) \right)
 \end{aligned}$$

Exponential rule.

$$\begin{aligned}
 = & T \left(\frac{\sum_t \sum_i f_k(\mathbf{x}_i^{(t)})}{T} - \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \frac{\exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right)}{Z} \left(\sum_i f_k(\mathbf{x}_i) \right) \right) \\
 = & T \left(\mathbb{E}_{\tilde{P}} \left[\sum_i f_k(\mathbf{x}_i) \right] - \mathbb{E}_{P_w} \left[\sum_i f_k(\mathbf{x}_i) \right] \right)
 \end{aligned}$$

The Gradient of Log-Likelihood

$$\begin{aligned}
 & \frac{\partial}{\partial w_k} \left[\left(\sum_t \sum_i \sum_j w_j f_j(\mathbf{x}_i^{(t)}) \right) - T \log Z \right] \\
 = & \frac{\partial}{\partial w_k} \left[\left(\sum_t \sum_i \sum_j w_j f_j(\mathbf{x}_i^{(t)}) \right) - T \log \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \right] \\
 = & \sum_t \sum_i f_k(\mathbf{x}_i^{(t)}) - T \frac{\partial}{\partial w_k} \left[\log \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \right] \\
 = & \sum_t \sum_i f_k(\mathbf{x}_i^{(t)}) - T \frac{\frac{\partial}{\partial w_k} \left[\sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \right]}{Z} \\
 = & \sum_t \sum_i f_k(\mathbf{x}_i^{(t)}) - \frac{T}{Z} \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \frac{\partial}{\partial w_k} \left[\exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \right] \\
 = & \sum_t \sum_i f_k(\mathbf{x}_i^{(t)}) - \frac{T}{Z} \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \left(\sum_i f_k(\mathbf{x}_i) \right) \\
 = & T \left(\frac{\sum_t \sum_i f_k(\mathbf{x}_i^{(t)})}{T} - \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \frac{\exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right)}{Z} \left(\sum_i f_k(\mathbf{x}_i) \right) \right)
 \end{aligned}$$

Rearrange terms.

$$= T \left(\mathbb{E}_{\tilde{P}} \left[\sum_i f_k(\mathbf{x}_i) \right] - \mathbb{E}_{P_w} \left[\sum_i f_k(\mathbf{x}_i) \right] \right)$$

The Gradient of Log-Likelihood

$$\begin{aligned}
 & \frac{\partial}{\partial w_k} \left[\left(\sum_t \sum_i \sum_j w_j f_j(\mathbf{x}_i^{(t)}) \right) - T \log Z \right] \\
 = & \frac{\partial}{\partial w_k} \left[\left(\sum_t \sum_i \sum_j w_j f_j(\mathbf{x}_i^{(t)}) \right) - T \log \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \right] \\
 = & \sum_t \sum_i f_k(\mathbf{x}_i^{(t)}) - T \frac{\partial}{\partial w_k} \left[\log \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \right] \\
 = & \sum_t \sum_i f_k(\mathbf{x}_i^{(t)}) - T \frac{\frac{\partial}{\partial w_k} \left[\sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \right]}{Z} \\
 = & \sum_t \sum_i f_k(\mathbf{x}_i^{(t)}) - \frac{T}{Z} \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \frac{\partial}{\partial w_k} \left[\exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \right] \\
 = & \sum_t \sum_i f_k(\mathbf{x}_i^{(t)}) - \frac{T}{Z} \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right) \left(\sum_i f_k(\mathbf{x}_i) \right) \\
 = & T \left(\frac{\sum_t \sum_i f_k(\mathbf{x}_i^{(t)})}{T} - \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \frac{\exp \left(\sum_i \sum_j w_j f_j(\mathbf{x}_i) \right)}{Z} \left(\sum_i f_k(\mathbf{x}_i) \right) \right) \\
 = & T \left(\mathbb{E}_{\tilde{P}} \left[\sum_i f_k(\mathbf{x}_i) \right] - \mathbb{E}_{P_w} \left[\sum_i f_k(\mathbf{x}_i) \right] \right)
 \end{aligned}$$

Difference of
expectations!



The Gradient of Log-Likelihood

- Difference of expectations!

$$T \left(\mathbb{E}_{\tilde{P}} \left[\sum_i f_k(\mathbf{x}_i) \right] - \mathbb{E}_{P_{\mathbf{w}}} \left[\sum_i f_k(\mathbf{x}_i) \right] \right)$$

- At a maximum of the likelihood function
- This form helps us prove the global concavity of the log-likelihood function.
 - Second derivative matrix (Hessian) is a correlation matrix of the features; it is positive semidefinite.
- The first term is simple; what about the second?

Feature Expectations

$$\mathbb{E}_{P_{\mathbf{w}}} \left[\sum_i f_k(\mathbf{x}_i) \right] = \sum_i \mathbb{E}_{P_{\mathbf{w}}} [f_k(\mathbf{x}_i)] = \sum_i \sum_{\mathbf{c} \in \text{Val}(\mathbf{X}_i)} P_{\mathbf{w}}(\mathbf{X}_i = \mathbf{c}) f_k(\mathbf{c})$$

- Linearity of expectation.
- Feature expectations are easily obtained from *marginals*.
 - We spent seven lectures on that problem!

Bayesian Learning

Maximum likelihood estimation: $\max_{\theta} P(\mathbf{X} | \theta)$

Bayes' rule:

$$\begin{aligned}\max_{\theta} P(\theta | \mathbf{X}) &= \max_{\theta} \frac{P(\mathbf{X} | \theta)P(\theta)}{P(\mathbf{X})} \\ &= \max_{\theta} \frac{P(\mathbf{X} | \theta)P(\theta)}{\int P(\mathbf{X} | \theta)P(\theta) d\theta} \\ &= \max_{\theta} P(\mathbf{X} | \theta)P(\theta)\end{aligned}$$

In Log Space

MLE

- $\max_{\theta} \sum_t \log P(\mathbf{X}^{(t)} | \theta)$
- “Make the data likely.”
- Closed form or convex in many cases.

Bayesian

- $\max_{\theta} \sum_t \log P(\mathbf{X}^{(t)} | \theta) + \log P(\theta)$
- “Make the data likely ... and make the model likely, too.”
- Closed form or convex in many cases.

Priors for Log-Linear Parameters

- For Bayesian networks, we fixated on conjugacy.
 - Conjugate priors for log-linear parameters don't decompose as nicely as the Dirichlet. See K&F 20.4.2.
- Not here; we already have to solve an optimization problem, so we're pretty open to any prior where $\log P(\mathbf{w})$ is concave.

Gaussian Prior on \mathbf{w}

- Let each w_j have a prior that says its mean is 0 and its variance is σ^2 .

$$\begin{aligned}\log P(\mathbf{w}) &= \sum_j \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{w_j^2}{2\sigma^2} \right) \right) \\ &= -\sum_j \frac{w_j^2}{2\sigma^2} + \text{constant} \\ &= -\frac{1}{2\sigma^2} \|\mathbf{w}\|_2^2\end{aligned}$$

- Result: quadratic/Euclidean/ L_2 penalty on likelihood.
- Generalizes **ridge regression**.

Laplacian Prior on \mathbf{w}

- Let each w_j have a Laplacian prior with parameter β .

$$\begin{aligned}\log P(\mathbf{w}) &= \sum_j \log \left(\frac{1}{2\beta} \exp \left(-\frac{|w_j|}{\beta} \right) \right) \\ &= -\frac{|w_j|}{\beta} + \text{constant} \\ &= -\frac{1}{\beta} \|\mathbf{w}\|_1\end{aligned}$$

- Result: absolute value/ L_1 penalty on likelihood.
 - Still concave, but not everywhere differentiable.
- Generalizes **lasso regression**.

Priors for Log-Linear Parameters

- Both the Gaussian and the Laplacian priors push the weights toward zero.
 - Laplacian pushes “harder” near zero, resulting in many weights being zero at the optimum (“**sparsity**”).
- Weights near zero make the distribution over **X** flatter (smoother).
- Choice of the hyperparameters (σ^2 or β) can have a big effect.
 - Cross-validation.

Discriminative Learning

Goal of Learning?

- **Density estimation:** return a model M that precisely captures P^*
- **Knowledge discovery:** reveal facts about the domain.
- **Prediction:** optimize quality of answers to specific queries

Generative vs. Discriminative

- Every model we've looked at so far this semester is generative, defining a distribution $P(\mathbf{X})$.
- Often we are less interested in density estimation than accurate performance on some query: $P(\mathbf{Y} \mid \mathbf{X})$.
- Discriminative models seek to perform well on a particular query.

Understanding Conditional Random Fields

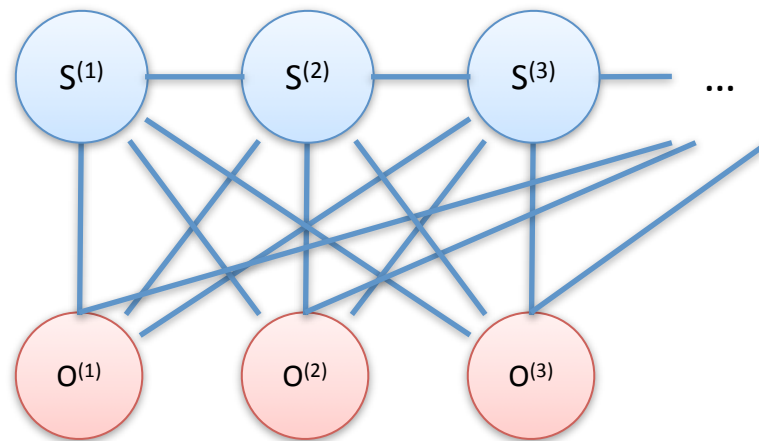
- We know that graphical models can be used to define conditional distributions rather than joint ones.
 - This is not quite the same as having a joint distribution and conditioning on evidence.
 - The model really has nothing to say about $P(\mathbf{X})$, only $P(\mathbf{Y} \mid \mathbf{X})$.
- Intuitive motivation: don't waste your time learning the density of something that will always be in evidence.

Conditional Random Fields

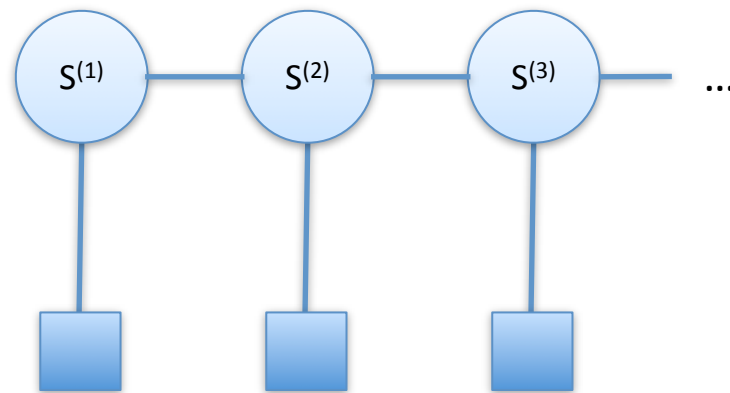
- Normalization now depends on \mathbf{X} . Because \mathbf{X} is always observed, every factor can depend on any part of \mathbf{X} .

$$P(\mathbf{Y} | \mathbf{X}) = \frac{1}{Z(\mathbf{X})} \prod_i \phi_i(\mathbf{Y}_i, \mathbf{X})$$
$$Z(\mathbf{X}) = \sum_{\mathbf{y} \in \text{Val}(\mathbf{Y})} \prod_i \phi_i(\mathbf{y}_i, \mathbf{X})$$

Example



Example



Conditional Random Fields

- Normalization now depends on \mathbf{X} . Because \mathbf{X} is always observed, every factor can depend on any part of \mathbf{X} .

$$P(\mathbf{Y} | \mathbf{X}) = \frac{1}{Z(\mathbf{X})} \prod_i \phi_i(\mathbf{Y}_i, \mathbf{X})$$

$$Z(\mathbf{X}) = \sum_{\mathbf{y} \in \text{Val}(\mathbf{Y})} \prod_i \phi_i(\mathbf{y}_i, \mathbf{X})$$

- Log-linear form:

$$\phi_i(\mathbf{Y}_i, \mathbf{X}) = \exp \mathbf{w}^\top \mathbf{f}(\mathbf{Y}_i, \mathbf{X})$$

Maximizing Conditional Likelihood

$$\begin{aligned}\boldsymbol{\theta}_{\text{MLE}} &= \arg \max_{\boldsymbol{\theta}} \prod_t P(\mathbf{y}^{(t)} \mid \mathbf{x}^{(t)}, \boldsymbol{\theta}) \\ &= \arg \max_{\boldsymbol{\theta}} \sum_t \log P(\mathbf{y}^{(t)} \mid \mathbf{x}^{(t)}, \boldsymbol{\theta}) \\ \mathbf{w}_{\text{MLE}} &= \arg \max_{\mathbf{w}} \sum_t \left(\sum_i \mathbf{w}^\top \mathbf{f}(\mathbf{y}_i^{(t)}, \mathbf{x}_i^{(t)}) - \log Z_{\mathbf{w}}(\mathbf{x}_i^{(t)}) \right)\end{aligned}$$

Compare to MLE for the classic Markov network:

$$\mathbf{w}_{\text{MLE}} = \arg \max_{\mathbf{w}} \sum_t \left(\sum_i \mathbf{w}^\top \mathbf{f}(\mathbf{y}_i^{(t)}, \mathbf{x}_i^{(t)}) - \log Z_{\mathbf{w}} \right)$$

Training the CRF

- Everything is the same as doing MLE in classic Markov networks, except now we have T different log partition functions.
 - Each requires marginalizing over \mathbf{Y} for a single \mathbf{x} , rather than over all random variables.
- CRF likelihood first derivatives:

$$\sum_t \left(\sum_i \mathbf{f}(\mathbf{y}_i^{(t)}, \mathbf{x}^{(t)}) - \mathbb{E}_{P_{\mathbf{w}}}[\mathbf{f}(\mathbf{Y}_i, \mathbf{x}^{(t)})] \right)$$

Additional Notes: CRF

- CRFs appear to be *way* more widely used than classic Markov networks.
 - The \mathbf{x} -specific partition functions are much less painful.
- Same training methods apply as before.
 - Same advice: L-BFGS, stochastic gradient descent.
 - Same priors: Gaussian, Laplacian

My Advice on Maximizing Likelihood

- If inference is relatively fast (on the whole dataset): use a **quasi-Newton** method like limited memory BFGS (L-BFGS).
 - “Batch” algorithm.
- If inference is relatively slow or you have massive amounts of training data, use **stochastic gradient descent**.
 - “Online” algorithm.
 - Lets you avoid Z, because you only need to calculate the additive updates.

CRF Pseudocode: Value and Gradient

- for $t = \{1, \dots, T\}$ // for each training example
- likelihood += log_score($y^{(t)}$)
- (marginals, logz) = inference ($x^{(t)}$)
- likelihood -= logz
- for i, y_i, j // for each clique, assignment, feature
- gradient(j) -= marginals(i, y_i) * $f_j(x^{(t)}, y_i)$
- gradient += constraints // constraints cached
- likelihood += prior(w)
- gradient += prior_gradient(w)

More implementation advice

- Cod structure: Optimizer asks for the likelihood and gradient at with parameters
- Efficiency:
 - Cache “constraints” (data expectations)
 - Cache features
 - Parallelize inference step
 - Online learning with stochastic gradient
- Gaussian prior typical values: 1, 10