

Learning Parameters of Undirected Models

Sargur Srihari

srihari@cedar.buffalo.edu

Topics

- Difficulties due to Global Normalization
- Likelihood Function
- Maximum Likelihood Parameter Estimation
 - Simple and Conjugate Gradient Ascent
- Conditionally Trained Models
- Parameter Learning with Missing Data
 - Gradient Ascent vs. EM
- Alternative Formulation of Max Likelihood
 - Maximum Entropy subject to Constraints
- Parameter Priors and Regularization

Local vs Global Normalization

- BN: Local normalization within each CPD

$$p(X) = \prod_{i=1}^n \ln p(X_i | pa_i)$$

- MN: Global normalization (partition function)

$$P(X_1, \dots, X_n) = \frac{1}{Z} \prod_{i=1}^K \phi_i(D_i)$$

$$Z = \sum_{X_1 \dots X_n} \phi_i(D_i)$$

- Global factor couples all parameters preventing decomposition
- Significant computational ramifications
 - M.L. parameter estimation has no closed-form soln.
 - Need iterative methods

Issues in Parameter Estimation

- Simple ML parameter estimation (even with complete data) cannot be solved in closed form
- Need iterative methods such as gradient ascent
- Good news:
 - Likelihood function is concave
 - Methods converge with global optimum
- Bad news:
 - Each step in iterative algorithm requires inference
 - Simple parameter estimation expensive/intractable
 - Bayesian estimation is practically infeasible

Discriminative Training

- Common use of MNs is in settings such as image segmentation where we have a particular inference task in mind
- Train the network discriminatively to get good performance for our particular inference task

Likelihood Function

- Basis for all discussion of learning
- How likelihood function can be optimized to find maximum likelihood parameter estimates
- Begin with form of likelihood function for Markov networks, its properties and their computational implications

Example of Coupled Likelihood Function

– Simple Markov network $A-B-C$

– Two potentials $\phi_1(A,B)$, $\phi_2(B,C)$:

- Gibbs: $P(A,B,C) = (1/Z) \phi_1(A,B) \phi_2(B,C)$ $Z = \sum_{a,b,c} \phi_1(a,b) \phi_2(b,c)$
- Log-likelihood of instance (a,b,c) is

$$\ln P(a,b,c) = \ln \phi_1(A,B) + \ln \phi_2(B,C) - \ln Z$$

- Log-likelihood of data set D with M instances:

$$\begin{aligned} \ell(\theta : D) &= \sum_{m=1}^M (\ln \phi_1(a[m], b[m]) + \ln \phi_2(b[m], c[m]) - \ln Z(\theta)) \\ &= \sum_{a,b} M[a,b] \ln \phi_1(a,b) + \sum_{b,c} M[b,c] \ln \phi_2(b,c) - M \ln Z(\theta) \end{aligned}$$

Parameter θ consists of all values of factors ϕ_1 and ϕ_2

Summing over different instances M

Summing over different values of A and B

- First term involves only ϕ_1 . Second only ϕ_2 . But third involves

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

- Which couples the two potentials in the likelihood function

- When we change one potential ϕ_1 , $Z(\theta)$ changes, possibly changing the value of ϕ_2 that maximizes $-\ln Z(\theta)$

Illustration of Coupled Likelihood

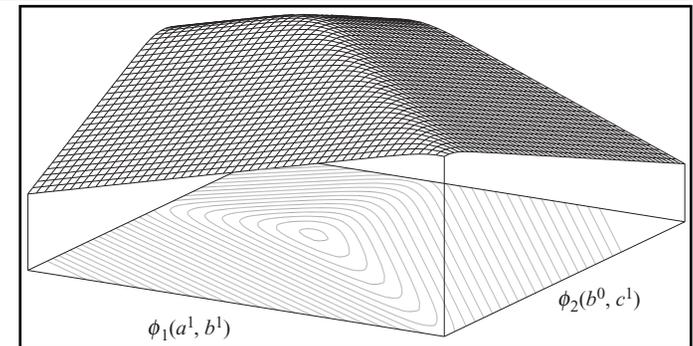
- Log-likelihood $A-B-C$ wrt two factors $\phi_1(A,B)$, $\phi_2(B,C)$
 - With binary variables we would have 8 parameters

- Log-likelihood surface $\ell(\theta : D) = \sum_{a,b} M[a,b] \ln \phi_1(a,b) + \sum_{b,c} M[b,c] \ln \phi_2(b,c) - M \ln Z(\theta)$

- when ϕ_1 changes, ϕ_2 that maximizes $-\ln Z(\theta)$ also changes

- In this case problem avoidable

- Equivalent to BN
- $A \rightarrow B \rightarrow C$ & estimate parameters of
- $\phi_1(A,B) = P(A)P(B|A)$, $\phi_2(B,C) = P(C|B)$



All other parameters set to 1
 Data Set has $M=100$
 $M[a^1, b^1]=40$, $M[b^0, c^1]=40$

- In general, cannot convert learned BN parameters into equivalent MN

- Optimal likelihood achievable by the two representations is not the same

Form of Likelihood Function

- Instead of Gibbs, use log-linear framework
 - Joint distribution of n variables X_1, \dots, X_n
 - k features $\mathcal{F} = \{f_i(D_i)\}_{i=1, \dots, k}$ k depends on how many values D_i take
 - where D_i is a sub-graph and f_i maps D_i to R

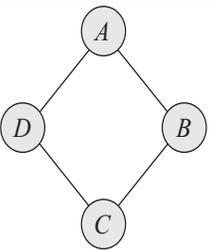
$$P(X_1, \dots, X_n; \theta) = \frac{1}{Z(\theta)} \exp \left\{ \sum_{i=1}^k \theta_i f_i(D_i) \right\}$$

- Parameters θ_i are weights we put on features
 - If we have a sample ξ then its features are $f_i(\xi(D_i))$ which has the shorthand $f_i(\xi)$.
- Representation is general
 - can capture Markov networks with global structure and local structure

Parameters θ , Factors ϕ , binary features f

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

A	B	$\phi_1(A, B)$
a^0	b^0	$\phi_1(a^0, b^0)$
a^0	b^1	$\phi_1(a^0, b^1)$
a^1	b^0	$\phi_1(a^1, b^0)$
a^1	b^1	$\phi_1(a^1, b^1)$



- Variables: $A - B - C - D -$

- A feature for every entry in every table

- $f_i(D_i)$ are sixteen indicator functions defined over clusters, AB, BC, CD, DA

$$f_{a^0 b^0}(A, B) = I\{A = a^0\} I\{B = b^0\}$$

$$\text{Val}(A) = \{a^0, a^1\} \quad \text{Val}(B) = \{b^0, b^1\}$$

$f_{a^0, b^0} = 1$ if $a = a^0, b = b^0$
0 otherwise, etc.

etc.

- With this representation

$$\theta_{a^0 b^0} = \ln \phi_1(a^0, b^0)$$

Parameters θ are potentials which are weights put on features

Log Likelihood and Sufficient Statistics

- Joint probability distribution:

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

$\theta = \{\theta_1, \dots, \theta_k\}$ are table entries
 f_i are features over instances of D_i

- Let \mathcal{D} be a data set of M samples $\xi [m]_{m=1, \dots, M}$

- Log-likelihood

– Log of product of probs of M indep. instances:

$$\ell(\theta : D) = \sum_i \theta_i \left(\sum_m f_i(\xi[m]) \right) - M \ln Z(\theta)$$

- Sufficient statistics (likelihood depends only on this)

Dividing
by no.
of samples

$$\frac{1}{M} \ell(\theta : D) = \sum_i \theta_i (E_D[f_i(d_i)]) - \ln Z(\theta)$$

$E_D[f_i(d_i)]$ is the average
in the data set

Properties of Log-Likelihood

- Log-likelihood is a sum of two functions

$$\ell(\theta : D) = \underbrace{\sum_i \theta_i \left(\sum_m f_i(\xi[m]) \right)}_{\text{First Term}} - \underbrace{M \ln Z(\theta)}_{\text{Second Term}}$$

- First term is linear in the parameters θ
 - Increasing parameters increases this term
 - But likelihood has upper-bound of probability 1

- Second term:

- balances first term

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

- Partition function $Z(\theta)$ is convex as seen next
 - » Its negative is concave

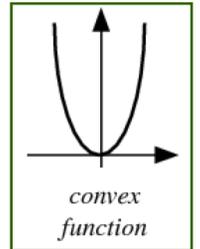
- Sum of linear fn. and concave is concave, So

- There are no local optima
 - Can use gradient ascent

Proof that $Z(\theta)$ is Convex

A function $f(\vec{x})$ is convex if for every $0 \leq \alpha \leq 1$

$$f(\alpha\vec{x} + (1-\alpha)\vec{y}) \leq \alpha f(\vec{x}) + (1-\alpha)f(\vec{y})$$



- The function is bowl-like
- Every interpolation between the images of two points is larger than the image of their interpolation
- One way to show that a Function is convex is to show that its Hessian (matrix of the function's second derivatives) is positive-semi-definite
 - Hessian (2nd der) of $\ln Z(\theta)$ computed as:

Hessian of $\ln Z(\theta)$

Given set of features \mathcal{F} with $\ln Z(\theta) = \ln \sum_{\xi} \exp \left\{ \sum_i \theta_i f_i(\xi) \right\}$

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

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

where $E_{\theta}[f_i]$ is shorthand for $E_{P(\chi, \theta)}[f_i]$

Proof:

Partial derivatives
wrt θ_i and θ_j

$$\frac{\partial}{\partial \theta_i} \ln Z(\theta) = \frac{1}{Z(\theta)} \sum_{\xi} \frac{\partial}{\partial \theta_i} \exp \left\{ \sum_j \theta_j f_j(\xi) \right\} = \frac{1}{Z(\theta)} \sum_{\xi} f_i(\xi) \exp \left\{ \sum_j \theta_j f_j(\xi) \right\} = E_{\theta}[f_i]$$

$$\frac{\partial^2}{\partial \theta_i \partial \theta_j} \ln Z(\theta) = \frac{\partial}{\partial \theta_j} \left[\frac{1}{Z(\theta)} \sum_{\xi} \frac{\partial}{\partial \theta_i} \exp \left\{ \sum_k \theta_k f_k(\xi) \right\} \right] = \text{Cov}_{\theta}[f_i; f_j]$$

– Since covariance matrix of features is pos. semi-definite, we have $-\ln Z(\theta)$ is a concave function of θ

- Corollary: log-likelihood function is concave

Non-unique Solution

- Since $\ln Z(\theta)$ is convex, $-\ln Z(\theta)$ is concave
- Implies that log-likelihood is unimodal
 - Has no local optima
- However does not imply uniqueness of global optimum
- Multiple parameterizations can result in same distribution
 - A feature for every entry in the table is always redundant, e.g.,
 - $f_{a0,b0} = 1 - f_{a0,b1} - f_{a1,b0} - f_{a1,b1}$
 - A continuum of parameterizations

Maximum (Conditional) Likelihood Parameter Estimation

- Task: Estimate parameters of a Markov network with a fixed structure given a fully observable data set \mathcal{D}
- Simplest variant of the problem is maximum likelihood parameter estimation
- Log-likelihood given features $\mathcal{F} = \{f_i, i=1, \dots, k\}$

is

$$\ell(\theta : D) = \sum_{i=1} \theta_i \left(\sum_m f_i(\xi[m]) \right) - M \ln Z(\theta)$$

Gradient of Log-likelihood

- Log-likelihood:
$$\ell(\theta : D) = \sum_{i=1} \theta_i \left(\sum_m f_i(\xi[m]) \right) - M \ln Z(\theta)$$

- Average log-likelihood:
$$\frac{1}{M} \ell(\theta : D) = \sum_i \theta_i (E_D[f_i(d_i)]) - \ln Z(\theta)$$

Dividing by no. of samples

- Gradient of second term is

$$\frac{\partial}{\partial \theta_i} \ln Z(\theta) = \frac{1}{Z(\theta)} \sum_{\xi} \frac{\partial}{\partial \theta_i} \exp \left\{ \sum_j \theta_j f_j(\xi) \right\} = \frac{1}{Z(\theta)} \sum_{\xi} f_i(\xi) \exp \left\{ \sum_j \theta_j f_j(\xi) \right\} = E_{\theta}[f_i]$$

- Gradient of the average log-likelihood is

$$\frac{\partial}{\partial \theta_i} \frac{1}{M} \ell(\theta : D) = E_D[f_i(\chi)] - E_{\theta}[f_i]$$

First term is average value of f_i in data D . Second term is expected value from distribution

- Provides a precise characterization of m.l. parameters θ
- Theorem:** Let \mathcal{F} be a feature set. Then θ is a m.l. assignment if and only if $E_D[f_i(\chi)] = E_{\hat{\theta}}[f_i]$ for all i
 - i.e., expected value of each feature relative to P_{θ} matches its empirical expectation in D

Need for Iterative Method

- Although log-likelihood function

$$\ell(\theta : D) = \sum_i \theta_i \left(\sum_m f_i(\xi[m]) \right) - M \ln Z(\theta)$$

- is concave, there is no analytical form for the maximum
- Since no closed-form solution
 - Can use iterative methods, e.g. gradient ascent as shown next
- Gradient is

$$\frac{\partial}{\partial \theta_i} \frac{1}{M} \ell(\theta : D) = E_D[f_i(\chi)] - E_\theta[f_i]$$

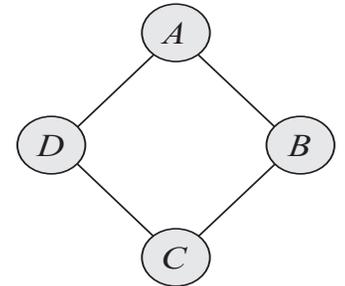
Computing the gradient

- Gradient wrt θ_i of log-likelihood $\ell(\theta; D)$ is

$$E_D[f_i(\chi)] - E_\theta[f_i]$$

– It is the difference between the

- feature's empirical count in data D and
- the expected count relative to our current parameterization



- Computing the empirical count $E_D[f_i(\chi)]$ is easy

– Ex: for feature

$$f_{a^0 b^0}(a, b) = I\{a = a^0\} I\{b = b^0\}$$

it is the empirical frequency in \mathcal{D} of the event a^0, b^0

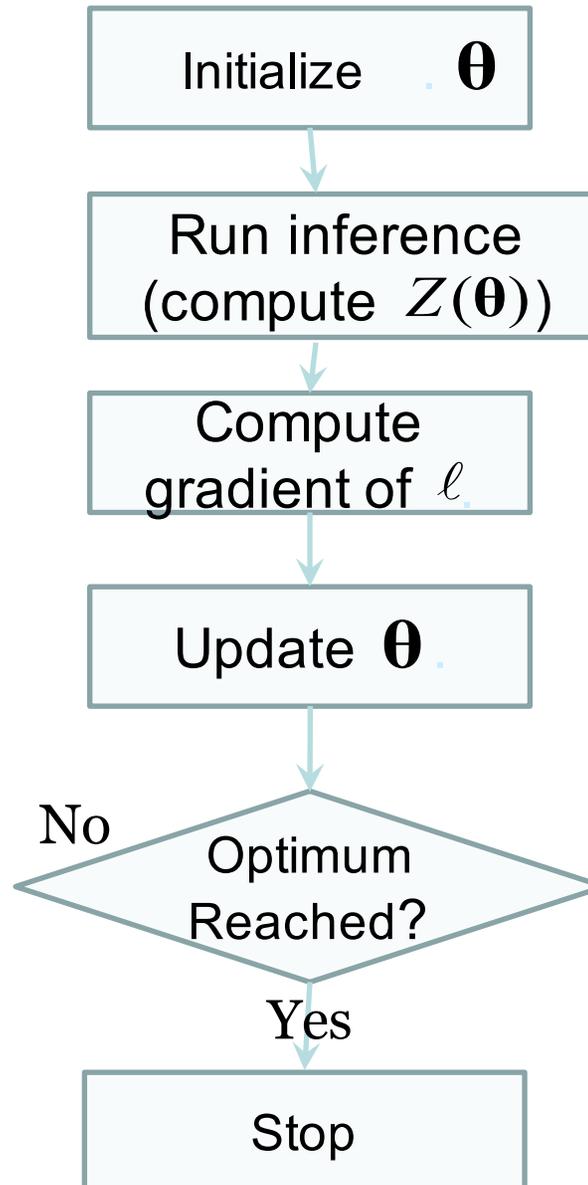
Computing the expected count

- Second term in the gradient is

$$E_{\theta}[f_i] = \frac{1}{Z(\theta)} \sum_{\xi} f_i(\xi) \exp \left\{ \sum_j \theta_j f_j(\xi) \right\}$$

- We need to compute the different probabilities of the form $P_{\theta^t}(a, b)$
 - Since expectation is a probability-weighted average
- Computing probability requires running inference over the network

Iterative solution for MRF parameters



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

$$E_D[f_i(\chi)] - E_{\theta}[f_i]$$

$$\theta^{t+1} \leftarrow \theta^t + \eta \nabla f_{obj}(\theta^t)$$

Difficulty with Iterative method

- Gradient ascent over parameter space
- Good news:
 - likelihood function is concave
 - Guaranteed to converge to global optimum
- Bad news:
 - each step needs inference
 - Simple parameter estimation is intractable
 - Bayesian parameter estimation even harder
 - Integration done using MCMC

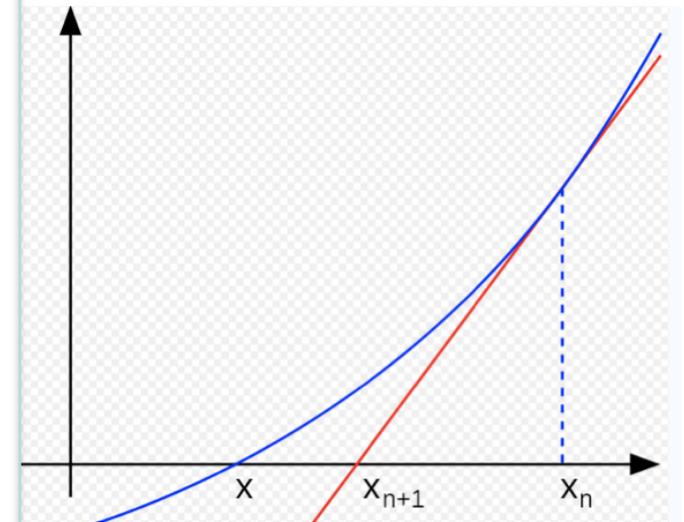
Use of Second Order Solution

- Computational cost of parameter estimation is very high
 - Gradient ascent is not efficient
- Much faster convergence using second order methods based on Hessian

Intuition for Second-Order Solution

- Newton's method
 - finds zeroes of a function using derivatives
- More efficient than simple gradient descent
- Quasi Newton's method
 - uses an approximation to the gradient
- Since we are solving for derivative of $l(\theta, D)$
 - need second derivative (Hessian)

Newton in one-dim.



$$f'(x_n) = \frac{\text{rise}}{\text{run}} = \frac{\Delta y}{\Delta x} = \frac{f(x_n) - 0}{x_n - x_{n+1}}$$

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$$

Derivation of Hessian of the Log-Likelihood

$$\begin{aligned}
 \frac{\partial}{\partial \theta_i} \ln Z(\theta) &= \frac{1}{Z(\theta)} \sum_{\xi} \frac{\partial}{\partial \theta_i} \exp \left\{ \sum_j \theta_j f_j(\xi) \right\} = \frac{1}{Z(\theta)} \sum_{\xi} f_i(\xi) \exp \left\{ \sum_j \theta_j f_j(\xi) \right\} = E_{\theta}[f_i] \\
 \frac{\partial^2}{\partial \theta_i \partial \theta_j} \ln Z(\theta) &= \frac{\partial}{\partial \theta_j} \left[\frac{1}{Z(\theta)} \sum_{\xi} \frac{\partial}{\partial \theta_i} \exp \left\{ \sum_k \theta_k f_k(\xi) \right\} \right] \\
 &= -\frac{1}{Z(\theta)^2} \left(\frac{\partial}{\partial \theta_j} Z(\theta) \right) \sum_{\xi} f_i(\xi) \exp \left\{ \sum_k \theta_k f_k(\xi) \right\} + \frac{1}{Z(\theta)} \sum_{\xi} f_i(\xi) f_j(\xi) \exp \left\{ \sum_k \theta_k f_k(\xi) \right\} \\
 &= -\frac{1}{Z(\theta)^2} Z(\theta) E_{\theta}[f_j] \sum_{\xi} f_i(\xi) \tilde{P}(\xi : \theta) + \frac{1}{Z(\theta)} \sum_{\xi} f_i(\xi) f_j(\xi) \tilde{P}(\xi : \theta) \\
 &= -E_{\theta}[f_j] \sum_{\xi} f_i(\xi) P(\xi : \theta) + \sum_{\xi} f_i(\xi) f_j(\xi) P(\xi : \theta) \\
 &= E_{\theta}[f_i f_j] - E_{\theta}[f_i] E_{\theta}[f_j] \\
 &= \text{Cov}_{\theta}[f_i; f_j]
 \end{aligned}$$

$$\frac{\partial}{\partial \theta_i \partial \theta_j} \ell(\theta, D) = -M \text{Cov}_{\theta}(f_i, f_j)$$

Computation of Hessian

- Log-likelihood has the form

$$\ell(\theta : D) = \sum_{i=1} \theta_i \left(\sum_m f_i(\xi[m]) \right) - M \ln Z(\theta)$$

- Solution 1: Hessian

$$\frac{\partial}{\partial \theta_i \partial \theta_j} \ell(\theta, D) = -M \text{Cov}_\theta(f_i, f_j)$$

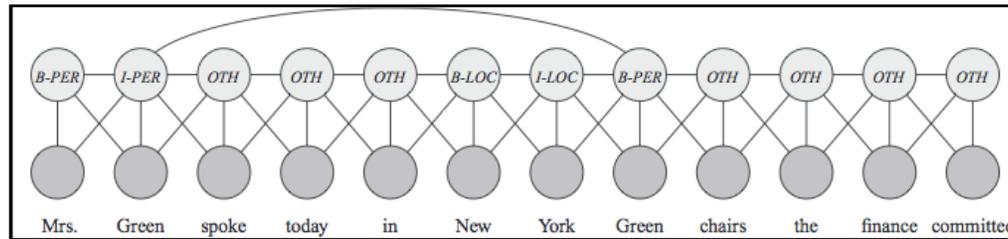
- Requires joint expectation of two features, often computationally infeasible
- Solution 2: Commonly used
 - *L-BFGS* (a quasi-Newton algorithm)
 - uses gradient ascent line search to avoid computing the Hessian

L-BFGS Algorithm

- Limited-memory BFGS
- Approximates the Broyden-Fletcher-Goldfarb-Shanno (BFGS)
- Popular algorithm for ML parameter estimation
 - Algorithm of choice for log-linear models and CRFs
- Uses an estimation of the inverse Hessian to steer through variable space

Conditionally Trained Models

- Often we want to perform an inference task where we have a known set of variables, or features, X



- We want to query a pre-determined set of variables Y
- We prefer to use *discriminative training*
- Train the network as a *Conditional Random Field (CRF)* that encodes a conditional distribution $P(Y|X)$

CRF Training

- Train the network as a CRF that encodes a conditional distribution $P(Y|X)$

- Training set consists of M pairs

$$D = \{\mathbf{y}[m], \mathbf{x}[m]\}, m = 1, \dots, M$$

Example:

$y[m]$ =word category, *B-PER*
 $x[m]$ =word, *Mrs.*

- Objective Function: Log-Conditional likelihood

$$E_{(\mathbf{x}, \mathbf{y}) \sim P^*} \left[\log \tilde{P}(\mathbf{y} | \mathbf{x}) \right]$$

We are not interested in the distribution of \mathbf{x} variables; only in predicting \mathbf{y} given \mathbf{x}

- Log-Conditional likelihood is

$$\begin{aligned} \ell_{Y|X}(\boldsymbol{\theta} : D) &= \ln P(\mathbf{y}[1, \dots, M] | \mathbf{x}[1, \dots, M], \boldsymbol{\theta}) \\ &= \sum_{m=1}^M \ln P(y[m] | \mathbf{x}[m], \boldsymbol{\theta}) \end{aligned}$$

Optimizing likelihood of each observed assignment $\mathbf{y}[m]$ given observed assignment $\mathbf{x}[m]$

–It is concave since each term is concave

Gradient of Conditional Likelihood

- A reduced MN is itself an MN.
- We use log-linear representation with features f_i and parameters θ
 - Analogous to gradient for full MN

$$\frac{\partial}{\partial \theta_i} \frac{1}{M} \ell(\theta : D) = E_D[f_i(\chi)] - E_\theta[f_i]$$

we can write gradient for reduced MN

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

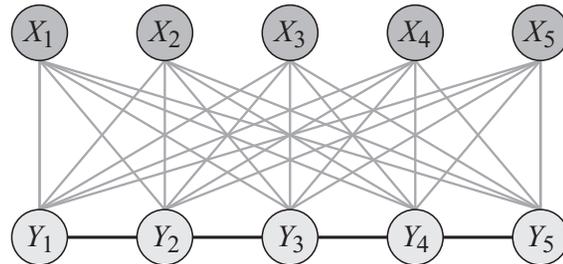
First term is empirical count conditioned on $x[m]$

Second is based on running inference on each data case

Conditional Training Complexity

- In the unconditional case, each gradient step required only a single execution of inference
- When training CRF we must execute inference for every single data case, conditioning on $\mathbf{x}[m]$
- On the other hand inference is executed on a simpler model
 - Since conditioning on evidence can only reduce computational cost

Simplification due to Conditioning



- Full MN encodes

$$\tilde{P}(\mathbf{X}, \mathbf{Y}) = \prod_{i=1}^4 \phi_i(Y_i, Y_{i+1}) \prod_{i=1}^5 \phi_i(Y_i, X_1, X_2, X_3, X_4, X_5)$$

- Edges disappear in a reduced Markov network

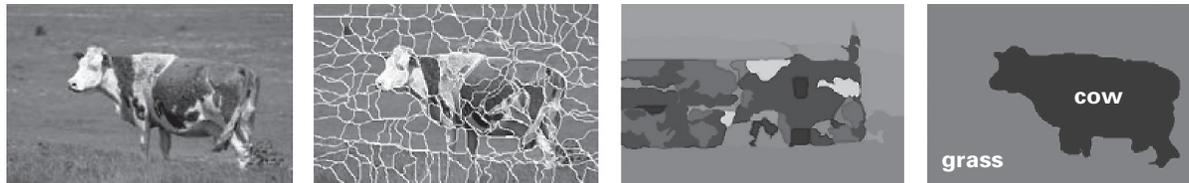
– After conditioning on \mathbf{X} , remaining edges form a simple chain

$$\tilde{P}(\mathbf{Y} | \mathbf{X}) = \prod_{i=1}^5 \phi_i(Y_i, Y_{i+1} | X_1, X_2, X_3, X_4, X_5)$$

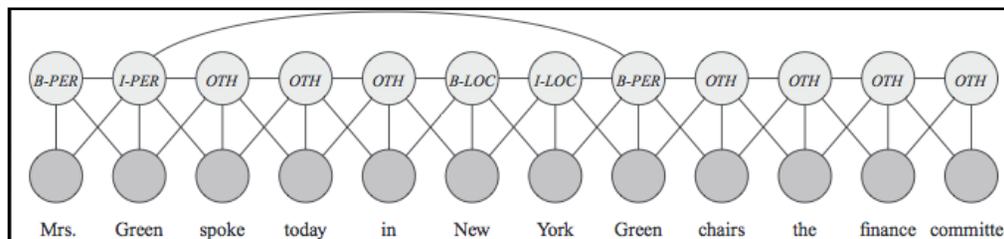
– Next slide compares different sequential models, CRF, HMM and MEMM wrt parameter learning complexity

Generative and Discriminative Models for Sequence Labeling

- A main task of PGMs: taking a set of inter-related instances and jointly labeling them
- Also called collective classification



Non-sequential



Sequential

- We look at trade-offs in using different models for instances organized sequentially

Models for Sequence Labeling

Given: sequence of observations $\mathbf{X} = \{X_1, \dots, X_k\}$.

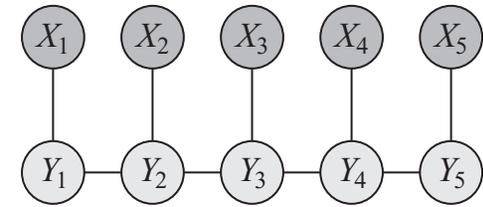
Need: a joint label $\mathbf{Y} = \{Y_1, \dots, Y_k\}$

Model Trade-offs: expressive power, learnability

- HMM is a generative model
That needs joint probability $P(\mathbf{X}, \mathbf{Y})$
- MEMM is a discriminative model
conditional directed model
 $Y_1 \perp X_2$ if not given Y_2 , by *D-separation*
More generally, $Y_i \perp X_j \mid \mathbf{X}_{-j} \ j > i$
- CRF is a discriminative model

CRF

Directly obtains $P(\mathbf{Y}|\mathbf{X})$



Note: $Z(\mathbf{X})$ is marginal of un-normalized measure

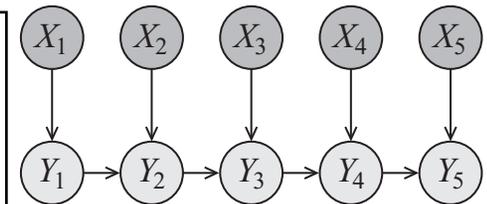
$$P(\mathbf{Y}|\mathbf{X}) = \frac{1}{Z(\mathbf{X})} \tilde{P}(\mathbf{Y}, \mathbf{X})$$

$$\tilde{P}(\mathbf{Y}, \mathbf{X}) = \prod_{i=1}^{k-1} \phi_i(Y_i, Y_{i+1}) \prod_{i=1}^k \phi_i(Y_i, X_i)$$

$$Z(\mathbf{X}) = \sum_{\mathbf{Y}} \tilde{P}(\mathbf{Y}, \mathbf{X})$$

MEMM

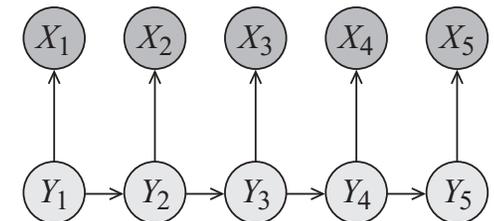
Does not model distribution over \mathbf{X}



$$P(\mathbf{Y}|\mathbf{X}) = \prod_{i=1}^k P(Y_i | X_i, Y_{i-1})$$

HMM

Needs joint distribution



$$P(\mathbf{X}, \mathbf{Y}) = \prod_{i=1}^k P(X_i | Y_i) P(Y_i | Y_{i-1})$$

$$P(\mathbf{Y} | \mathbf{X}) = \frac{P(\mathbf{X}, \mathbf{Y})}{P(\mathbf{X})}$$

Comparison of Sequential Models

- Trade-offs: Training, Expression, Independence
- Training Effort
 - HMM, MEMM easily learned (they are BNs)
 - CRF: gradient-based inference for every sequence
difficult with large data sets
- Expressibility
 - HMM is generative, training over features is hard
 - CRF can use a rich feature set, which is important
- Independence Assumptions

MEMM Has label bias problem: $Y_i \perp X_j \mid X_j \ j > I$

Later observation has no effect on posterior probability of current state. In activity recognition in video sequence: frames are labeled as running/walking. Earlier frames may be blurry but later ones clearer.

Markov Network Parameter Estimation with Missing Data

Parameter Estimation with Missing Data

- We look at estimating parameters θ
- When some data fields are missing
- Missing data examples
 - Some data fields omitted or not collected
 - Some hidden variables

Methods for Parameter Estimation with Missing Data

- Learning problem has difficulties
 - Parameters θ may not be identifiable
 - Coupling between different parameters
 - Likelihood is not concave (has local maxima)
- Two alternative methods
 1. Gradient Ascent (assume missing data is random)
 2. Expectation-Maximization

Parameter Method 1 for Missing Data: Gradient Ascent

- Assume data is missing at random in data \mathcal{D}
- In the m^{th} instance, let $o[m]$ be observed entries and $h[m]$ random variables that are missing

- Log-likelihood has the form

$$\frac{1}{M} \ln P(D|\theta) = \frac{1}{M} \sum_{m=1}^M \ln \left(\sum_{h[m]} \tilde{P}(o[m], h[m]|\theta) \right) - \ln Z$$

- Gradient wrt θ_i for feature f_i has the form

$$\frac{\partial}{\partial \theta_i} \frac{1}{M} \ell(\theta; D) = \frac{1}{M} \left[\sum_{m=1}^M E_{h[m] \sim P(h[m]|o[m], \theta)} [f_i] \right] - E_{\theta} [f_i]$$

- Which is a difference between two expectations of f_i
 - Expectation over data and hidden variables
 - Expectation in current distribution $P(\chi|\theta)$

Gradient Ascent Complexity: Full vs Missing

1. With full data, gradient of log-likelihood is

$$\frac{\partial}{\partial \theta_i} \frac{1}{M} \ell(\boldsymbol{\theta} : D) = E_D[f_i(\boldsymbol{\chi})] - E_{\boldsymbol{\theta}}[f_i]$$

- For second term we need inference over current distribution $P(\boldsymbol{\chi}|\boldsymbol{\theta})$
- First term is aggregate over data \mathcal{D} .

2. With missing data we have

$$\frac{\partial}{\partial \theta_i} \frac{1}{M} \ell(\boldsymbol{\theta} : D) = \frac{1}{M} \left[\sum_{m=1}^M E_{h[m] \sim P(h[m]|o[m], \boldsymbol{\theta})} [f_i] \right] - E_{\boldsymbol{\theta}} [f_i]$$

- We have to run inference separately for every instance m conditioning on $o[m]$
- Cost is much higher than with full data

Parameter Method 2 for Missing Data: Expectation Maximization

- As for any probabilistic model an alternative method of estimating parameters θ with missing data is EM
 - E step is to use current parameters to estimate missing values
 - M step is to re-estimate the parameters
- For BN it has significant advantages
 - M-step has closed-form solution
- For MN, the answer is not clear-cut
 - M-step requires running inference multiple times
 - As seen next

EM for MN parameter learning

- E-step

- using current parameters θ^t to compute expected sufficient statistics, i.e., expected feature counts

- At iteration t expected sufficient statistic for feature f_i is

$$M_{\theta^{(t)}}[f_i] = \frac{1}{M} \left[E_{h[m] \sim P(h[m]|o[m], \theta)} [f_i] \right]$$

- M-step

- done using maximum likelihood parameter estimation (using gradient ascent!)

- Requires running inference multiple times, once for each iteration of gradient ascent procedure

- At step k of this “inner loop” optimization, we have a gradient of the form

$$M_{\theta^{(t)}}[f_i] - E_{\theta^{(t,k)}}[f_i]$$

MN Parameter Learning: Maximum Entropy Formulation of Maximum Likelihood Objective

Maximum Entropy Formulation of Markov Network Parameter Estimation

- Basic formulation is: maximum likelihood parameter estimation
- Alternative formulation provides insight
- Insight is that

“Distribution that maximizes entropy subject to constraints (of expectation) is the Gibbs distribution”

Maximizing entropy subject to constraints is dual of maximizing likelihood given structural constraints

Entropy Definition

- Entropy is a measure of information
 - It is the *Expected value* (average value) of information generated by the distribution
 - Where Information is (log) reciprocal of probability
 - For a random variable X with distribution $P(X)$

$$H_P(X) = E_P \left[\log \frac{1}{P(X)} \right] = \sum_x P(x) \log \frac{1}{P(X)}$$

- where we treat $0 \log 1/0 = 0$
- Entropy is the Lower bound on the average number of bits needed to encode values of X
 - Central Theorem of Information Theory by Shannon
 - i.e., Expected code length cannot be less than $H_P(X)$

Motivation of Max Entropy with Constraints

- Suppose we have some summary statistics of an empirical distribution
 - Ex 1: published in a census report
 - Marginal distributions of single variables or certain pairs, or other events of interest
 - Ex 2: average final grade of students in class (70), correlation of final grades with homework scores
 - (correlation coefficient, $r=0.6$, on scale of -1 to +1)
- But no access to full data set
- Find a typical distribution that satisfies these constraints and use it to answer other queries

Maximum Entropy Formulation

- Goal is to select a distribution that:
 1. Satisfies given constraints
 2. Has no additional structure or information
- We formulate the following problem:
 - Find $Q(\chi)$
 - Maximizing $H_Q(\chi)$
 - Subject to $E_Q[f_i] = E_{\mathcal{D}}[f_i], i=1, \dots, k$
- These are called *expectation constraints*

Solution to Max Entropy Estimation

- Somewhat surprisingly
- The distribution Q^* , the maximum entropy distribution satisfying the expectation constraints, $E_Q[f_i] = E_{\mathcal{D}}[f_i]$, is a Gibbs distribution $Q^* = P_{\hat{\theta}}$ where

$$P_{\hat{\theta}}(\chi) = \frac{1}{Z(\hat{\theta})} \exp \left\{ \sum_i \hat{\theta}_i f_i(\chi) \right\}$$

where $\hat{\theta}$ is the m.l.e. parameterization relative to \mathcal{D}

Duality of Max Ent and Max Likelihood

- The two problems
 1. *Maximizing entropy* subject to expectation constraints
 2. *Maximizing likelihood* given structural constraints on the distribution
 - are *convex duals of each other*
- For the maximum likelihood parameters $\hat{\theta}$

$$H_{P_{\hat{\theta}}}(\mathcal{X}) = -\frac{1}{M} \ell(\hat{\theta} : D)$$

Convex Duality

- Optimization problems are viewed as
- Primal Problem
 - E.g., Lagrangian form with specified constraints
- Dual Problem
 - Provides lower bounds to the primal problem
- In general optimal values are not the same
- For convex optimization the duality gap is zero

MN Parameter Estimation: Parameter Priors and Regularization

Parameter Priors and Regularization

- Maximum likelihood estimation of θ is prone to over-fitting
- This is true of maximum likelihood of parameters of Bayesian networks as well
- But it is not as apparent with MNs due to lack of direct correspondence between empirical counts and parameters
- Overfitting is as much of a problem with MNs

Parameter Priors

- Can introduce prior distribution $P(\theta)$
 - Defined over model parameters θ
- Due to non-decomposable likelihood function, a fully Bayesian approach is infeasible

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

- In a fully Bayesian approach we integrate out the parameters to get the next prediction
- Instead, we perform MAP estimation
 - Find parameters that maximize $P(\theta)P(\mathcal{D} | \theta)$

MAP estimation of θ

- Instead, we perform MAP estimation
 - Find parameters that maximize $P(\theta)P(\mathcal{D} | \theta)$
 - Where we have $\ln P(\mathcal{D} | \theta)$ expressed in log-space as

$$\ell(\theta : D) = \sum_{i=1} \theta_i \left(\sum_m f_i(\xi[m]) \right) - M \ln Z(\theta)$$

- Which in turn is derived from the joint distribution

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

Choice of Prior $P(\theta)$

- No constraints on conjugacy of prior and likelihood. Only a few priors used in practice

1. Zero-mean diagonal Gaussian with equal variances for each of the weights
2. Zero-mean Laplacian distribution

- They are referred to as

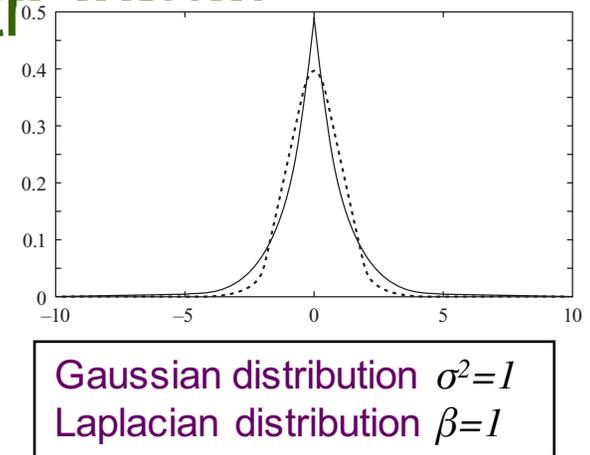
- L_2 regularization,

- Weight penalty is quadratic in θ (Euclidean norm)

- L_1 regularization

- Weight penalty is linear in $|\theta|$

- Both priors penalize parameters whose magnitude (positive or negative) is large



Gaussian Prior and L_2 -Regularization

- Most common is Gaussian prior on log-linear parameters θ

$$P(\theta | \sigma^2) = \prod_{i=1}^k \frac{1}{\sqrt{2\pi\sigma}} \exp\left\{-\frac{\theta_i^2}{2\sigma^2}\right\}$$

– For some choice of hyper-parameter (variance) σ^2

- Analogous to α_i in Dirichlet prior for multinomial

- Converting MAP objective $P(\theta)P(\mathcal{D} | \theta)$ to log-space, gives $\ln P(\theta) + \ell(\theta : D)$ whose first term

$$\ln P(\theta) = -\frac{1}{2\sigma^2} \sum_{i=1}^k \theta_i^2$$

– is called an L_2 -regularization term

Laplacian Prior and L_1 -Regularization

- Zero-mean Laplacian distribution

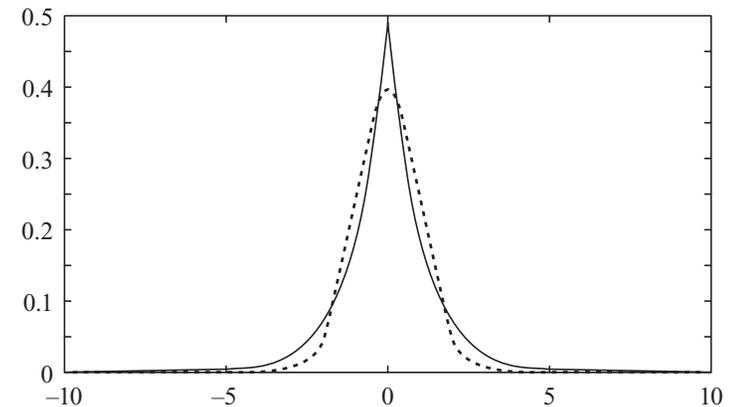
$$P_{Laplacian}(\boldsymbol{\theta}) = \frac{1}{2\beta} \exp\left\{-\frac{|\boldsymbol{\theta}|}{\beta}\right\}$$

- Taking log we obtain a term

$$\ln P(\boldsymbol{\theta}) = -\frac{1}{\beta} \sum_{i=1}^k |\theta_i|$$

– which we wish to minimize

- Generally called L_1 -regularization
- Both forms of regularization penalize parameters whose magnitude is large



Laplacian distribution $\beta=1$
Gaussian distribution $\sigma^2=1$

Why prefer low magnitude parameters ?

- Properties of prior
 - To pull distribution towards an uninformed one
 - To smooth fluctuations in the data
- A distribution is *smooth* if
 - Probabilities assigned to different assignments are not radically different
- Consider two assignments ξ and ξ'
 - An assignment is an instance of variables X_1, \dots, X_n
- We consider ratio of their probabilities next

Smoothness resulting from small θ

- Given two assignments ξ and ξ'

– Their relative probability is

$$\frac{P(\xi)}{P(\xi')} = \frac{\tilde{P}(\xi) / Z(\theta)}{\tilde{P}(\xi') / Z(\theta)} = \frac{\tilde{P}(\xi)}{\tilde{P}(\xi')}$$

- where the un-normalized probabilities are

$$\tilde{P}(\xi) = \exp\left\{\sum_{i=1}^k \theta_i f_i(\xi)\right\}$$

- In log-space, log-probability ratio is

$$\ln \frac{P(\xi)}{P(\xi')} = \sum_{i=1}^k \theta_i f_i(\xi) - \sum_{i=1}^k \theta_i f_i(\xi') = \sum_{i=1}^k \theta_i (f_i(\xi) - f_i(\xi'))$$

- When θ_i 's have small magnitude, this log-ratio is also bounded, i.e., probabilities are similar
- This results in a smooth distribution

Comparison of L_1 and L_2 Regularization

- Both L_1 and L_2 penalize parameter magnitude
 - Encode belief that model weights should be small (Close to zero)
- In Gaussian case (L_2), penalty grows quadratically with parameters
 - An increase in θ_i from 3 to 3.1 is penalized more than θ_i from 0 to 0.1
 - Leads to many small parameters
- In Laplacian case (L_1), penalty grows linearly
 - Results in fewer edges and is more tractable

Efficiency of Optimization

- Both L_1 and L_2 Regularization terms are Concave
- Because Log-likelihood is also Concave, resulting posterior is also concave
- Can be optimized using gradient descent methods
- Introduction of penalty terms eliminates multiple equivalent minima

Choice of Hyper-parameters

- Regularization hyper-parameters σ^2 (L_1) and β (L_2)
 - Larger hyper-parameters mean broader prior
 - Choice of prior has effect on learned model
- Standard method of selecting this parameter is via cross-validation
 - Repeatedly partition training set
 - Learn model over one part with some choice of parameters
 - ensure the performance on the held-out fragment