

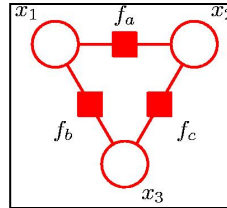
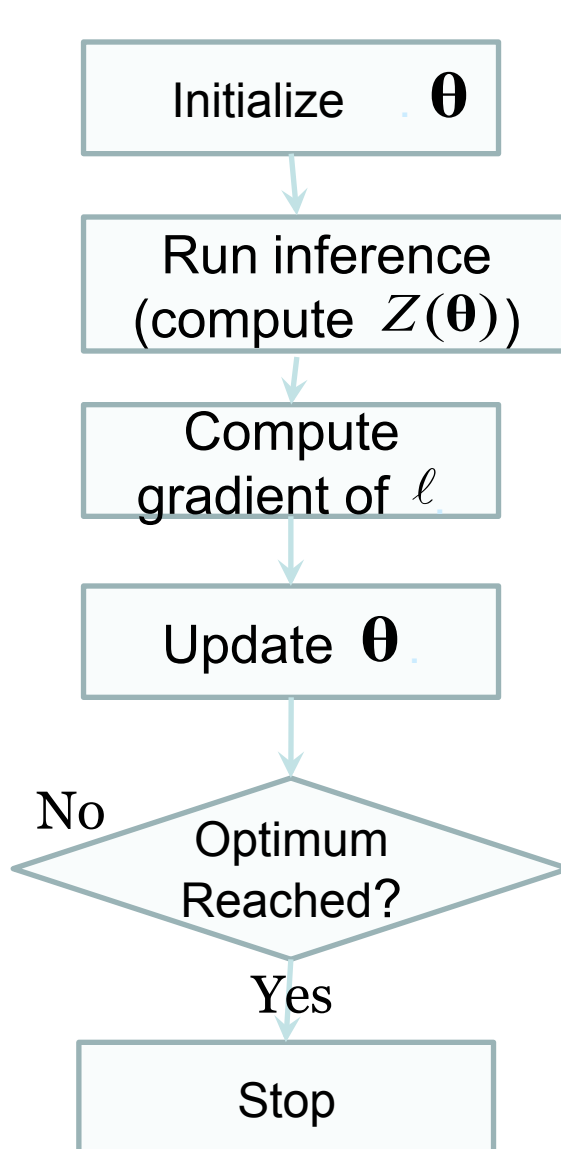
Learning MN Parameters with Approximation

Sargur Srihari
srihari@cedar.buffalo.edu

Topics

- Iterative exact learning of MN parameters
- Difficulty with exact methods
- Approximate methods
 - Approximate Inference
 - Belief Propagation
 - MAP-based Learning
 - Alternative Objective Functions
 - Pseudo-likelihood
 - Contrastive Objective
 - Divergence
 - Map-based

Exact Learning of MN parameters θ



$$P_{\theta}(\chi) = \frac{1}{Z(\theta)} \exp \left[- \sum_{i=1}^k \theta_i f_i(D_i) \right]$$

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

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

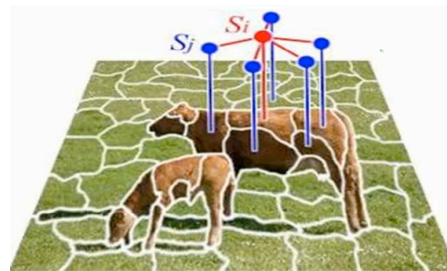
$$\frac{\partial}{\partial \theta_i} \ell(\theta; D) = E_D[f_i(\chi)] - E_{\theta}[f_i] \quad \nabla \ell(\theta) = \frac{d}{d\theta} \ell(\theta) = \left[\frac{\partial \ell(\theta)}{\partial \theta_1} \quad \frac{\partial \ell(\theta)}{\partial \theta_k} \right]$$

$$\theta^{t+1} \leftarrow \theta^t + \eta \nabla \ell(\theta^t; D)$$

$$\theta^{t+1} - \theta^t \leq \delta$$

Difficulty with Exact Methods

- Exact Parameter Estimation Methods assume ability to compute
 1. Partition function $Z(\theta)$ and
 2. Expectations $E_{p_\theta} [f_i]$
- In many applications structure of network does not allow exact computation of these terms
 - In image segmentation, grid networks lead to exponential size clusters for exact inference



Cluster graph is clique tree with overlapping factors

Approximate Methods for Learning θ

1. Use approximate inference for queries of P_θ
 - Decouples inference from Learning Parameters
 - Inference is a black-box
 - But approximation may interfere with learning
 - Non-convergence of inference can lead to oscillating estimates of the gradient & no learning convergence
2. Use an approximate objective function
 - Whose optimization doesn't require much inference
 - *Approximately optimizing the likelihood function can be reformulated as exactly optimizing an approximate objective*

Approximate Methods (Taxonomy)

1. Approximate Inference methods

1. Belief Propagation
2. MAP-based Learning

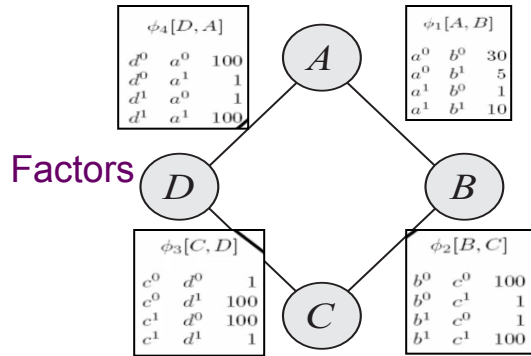
2. Alternative Objective Functions

1. Pseudo-likelihood and its generalizations
2. Contrastive Optimization Criteria
 1. Contrastive Divergence
 2. Margin-Based Training

Approximate Inference: Belief Propagation

- Popular Approach for Approximate Inference is Belief Propagation
 - Given a model resulting from a learning procedure, an algorithm from this family is used
- Model trained with same inference algorithm is better than model trained with exact inference!
- BP is run with each iteration of Gradient Ascent to compute expected feature count $E_{P_\theta} [f_i]$

Ex: Direct (Gibbs) vs BP (Clique tree)



1. Gibbs Distribution

$$P(A, B, C, D) = \frac{1}{Z} \phi_1(A, B) \cdot \phi_2(B, C) \cdot \phi_3(C, D) \cdot \phi_4(D, A)$$

where

$$Z = \sum_{A, B, C, D} \phi_1(A, B) \cdot \phi_2(B, C) \cdot \phi_3(C, D) \cdot \phi_4(D, A)$$

$$Z = 7,201,840$$

$$\tilde{P}_{\Phi}(A, B, C, D) = \phi_1(A, B) \phi_2(B, C) \phi_3(C, D) \phi_4(D, A)$$

Unnormalized Distribution

Assignment				Unnormalized
a^0	b^0	c^0	d^0	300000
a^0	b^0	c^0	d^1	300000
a^0	b^0	c^1	d^0	300000
a^0	b^0	c^1	d^1	30
a^0	b^1	c^0	d^0	500
a^0	b^1	c^0	d^1	500
a^0	b^1	c^1	d^0	5000000
a^0	b^1	c^1	d^1	500
a^1	b^0	c^0	d^0	100
a^1	b^0	c^0	d^1	1000000
a^1	b^0	c^1	d^0	100
a^1	b^0	c^1	d^1	100
a^1	b^1	c^0	d^0	10
a^1	b^1	c^0	d^1	100000
a^1	b^1	c^1	d^0	100000
a^1	b^1	c^1	d^1	100000

2. Clique Tree (triangulated):

Initial Potentials:

Each ψ has every factor involving its arguments

$$\psi_1(A, B, D) = \phi_1(A, B) \phi_2(B, C) \phi_3(C, D) \phi_4(D, A)$$

$$\psi_2(B, C, D) = \phi_1(A, B) \phi_2(B, C) \phi_3(C, D) \phi_4(D, A)$$

Clique/Cluster C_1

1. A, B, D

Sepset

$\{B, D\}$

Cluster C_2

2. B, C, D

Computing Clique Beliefs (β_i), Sepset Beliefs ($\mu_{i,j}$)

$$\beta_1(A, B, D) = \tilde{P}_{\Phi}(A, B, D) = \sum_C \psi_1(A, B, D) = \sum_C \phi_1(A, B) \phi_2(B, C) \phi_3(C, D) \phi_4(D, A)$$

e.g., $\beta_1(a^0, b^0, d^0) = 300,000 + 300,000 = 600,000$

$$\mu_{1,2}(B, D) = \sum_{C_1 - S_{1,2}} \beta_1(C_1) = \sum_A \beta_1(A, B, D)$$

e.g., $\mu_{1,2}(b^0, d^0) = 600,000 + 200 = 600,200$

$$\beta_2(B, C, D) = \tilde{P}_{\Phi}(B, C, D) = \sum_A \mu_{1,2}(B, D) \cdot \psi_2(B, C, D) = \sum_A \psi_2(B, C, D)$$

e.g., $\beta_2(b^0, c^0, d^0) = 300,000 + 100 = 300,100$

All Clique and Sepset Beliefs

Assignment			max _C	Assignment			max _{A,C}	Assignment			max _D
a ⁰	b ⁰	d ⁰	600,000	b ⁰	d ⁰	600,200	b ⁰	c ⁰	d ⁰	300,100	
a ⁰	b ⁰	d ¹	300,030	b ⁰	d ¹	1,300,130	b ⁰	c ¹	d ¹	1,300,100	
a ⁰	b ¹	d ⁰	5,000,500	b ¹	d ⁰	5,100,510	b ⁰	c ¹	d ⁰	300,100	
a ⁰	b ¹	d ¹	1,000	b ¹	d ¹	201,000	b ¹	c ⁰	d ⁰	100,100	
a ¹	b ⁰	d ⁰	200				b ¹	c ⁰	d ¹	100,100	
a ¹	b ⁰	d ¹	1,000,100				b ¹	c ¹	d ⁰	5,100,100	
a ¹	b ¹	d ⁰	100,010				b ¹	c ¹	d ¹	100,100	
a ¹	b ¹	d ¹	200,000								
β ₁ (A, B, D)				μ _{1,2} (B, D)				β ₂ (B, C, D)			

Verifying Inference

$$\tilde{P}_{\Phi}(a^1, b^0, c^1, d^0) = 100$$

$$\frac{\beta_1(a^1, b^0, d^0) \beta_2(b^0, c^1, d^0)}{\mu_{1,2}(b^0, d^0)} = \frac{200 \cdot 300 \cdot 100}{600 \cdot 200} = 100$$

Difficulty with Belief Propagation

- In principle, BP is run in each Gradient Ascent iteration
 - to compute $E_{P_\theta} [f_i]$ used in gradient computation
 - Due to family preservation property each feature must be a subset of a cluster C_i in the cluster graph
 - Hence to compute $E_{P_\theta} [f_i]$ we can compute BP marginals over C_i
- But BP often does not converge
 - Marginals derived often oscillate
 - Final results depend on where we stop
 - As a result gradient computed is unstable
 - Hurt convergence properties of gradient descent
 - Even more severe with line search

Convergent alternatives to Belief Propagation

- Three BP methods:

1. Pseudo- moment matching

- Reformulate task of learning with approximate inference as optimizing an alternative objective

2. Maximum Entropy approximation (CAMEL)

- General derivation that allows us to reformulate maximum likelihood with BP as a unified optimization problem with an approximate objective

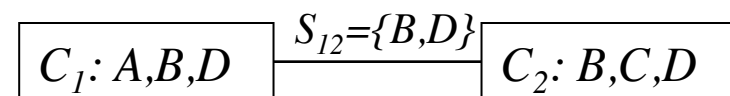
3. MAP-based Learning

- Approximate expected feature counts with their counts in the single MAP assignment in current MN

Pseudo-moment Matching

- Begin with analysis of fixed points in learning
 - Converged BP beliefs must satisfy $E_{\beta_i(C_i)}[f_{C_i}] = E_D[f_i(C_i)]$
 - Or $\beta_i(c_i^j) = \hat{P}(c_i^j)$ Convergent point is a set of beliefs that match the data
- We define for each sepset $S_{i,j}$ between C_i and C_j

$$\phi_i \leftarrow \frac{\beta_i}{\mu_{i,j}}$$



- We use the final set of potentials as the parameterization of the Markov Network
- Provides a closed-form solution for both inference and learning
 - Cannot be used with parameter regularization, non-table factors or CRFs

BP and Maximum Entropy

- Maximum entropy dual of max likelihood

Find $Q(\chi)$

Maximizing $H_Q(\chi)$

Subject to $E_Q[f_i] = E_D[f_i], i=1, \dots, k$

- Made tractable by approximating H_Q and E_Q :
 - H_Q : Given a cluster graph U with clusters C_i and sepset $S_{i,j}$

$$H_Q(\chi) \approx \sum_{C_i \in U} H_{\beta_i}(C_i) - \sum_{C_i - C_j \in S} H_{\mu_{i,j}}(S_{i,j})$$
 - $E_Q[f_i]$: replace by $E_{\beta_i}[f_i]$
 - Approximation is exact when cluster graph is a tree
 - Method is known as CAMEL (*Constrained Approximate Maximum Entropy Learning*)

Example of Max Ent Learning

- Pairwise Markov Network $\overbrace{A-B-C}$
 - Variables are binary
 - Three clusters: $C_1 = \{A, B\}$, $C_2 = \{B, C\}$, $C_3 = \{C, A\}$
 - Log-linear model with features
 - $f_{00}(x, y) = 1$ if $x=0, y=0$ and 0 otherwise for x, y , instance of C_i
 - $f_{11}(x, y) = 1$ if $x=1, y=1$ and 0 otherwise
 - Three data instances (A, B, C) : $(0, 0, 0), (0, 1, 0), (1, 0, 0)$
 - Unnormalized Feature counts, pooled over all clusters, are

$$E_{\hat{P}}[f_{00}] = (3 + 1 + 1) / 3 = 5 / 3$$

$$E_{\hat{P}}[f_{11}] = (0 + 0 + 0) / 3 = 0$$

CAMEL Optimization Problem

- Optimization problem takes the form
 - with two types of constraints:

Find $Q = \{\beta_1, \beta_2, \beta_3, \mu_{1,2}, \mu_{2,3}, \mu_{1,3}\}$
 maximizing $H_{\beta_1}(A, B) + H_{\beta_2}(B, C) + H_{\beta_3}(A, C)$
 subject to $-H_{\mu_{1,2}}(B) - H_{\mu_{2,3}}(C) - H_{\mu_{1,3}}(A)$

$$\sum_i E_{\beta_i}[f_{00}] = 5/3$$

$$\sum_i E_{\beta_i}[f_{11}] = 0$$

$$\sum_a \beta_1(a, b) - \sum_c \beta_2(b, c) = 0 \quad \forall b$$

$$\sum_b \beta_2(b, c) - \sum_a \beta_3(a, c) = 0 \quad \forall c$$

$$\sum_c \beta_3(a, c) - \sum_b \beta_1(a, b) = 0 \quad \forall a$$

$$\sum_{\mathbf{c}_i} \beta_i(\mathbf{c}_i) = 1 \quad i = 1, 2, 3$$

$$\beta_i \geq 0 \quad i = 1, 2, 3.$$

Type 1 Constraints:

$$E_Q[f_i] = E_D[f_i], i=1, \dots, k$$

Type 2 Constraints:
 Marginals from
 Cluster-graph
 approximation

CAMEL Solutions

- CAMEL optimization is a constrained maximization problem with
 - linear constraints and
 - a nonconcave objective
- Several solution algorithms, one of which is
 - Lagrange multipliers for all constraints and optimize over resulting new variables

Sampling-based Learning

- Partition function $Z(\theta)$ is a summation over an exponentially large space
 - Reformulate as expectation wrt a distribution $Q(\chi)$

$$\begin{aligned} Z(\theta) &= \sum_{\xi} \exp \left\{ \sum_i \theta_i f_i(\xi) \right\} \\ &= \sum_{\xi} \frac{Q(\xi)}{Q(\xi)} \exp \left\{ \sum_i \theta_i f_i(\xi) \right\} \\ &= E_Q \left[\frac{1}{Q(\chi)} \exp \left\{ \sum_i \theta_i f_i(\chi) \right\} \right] \end{aligned}$$

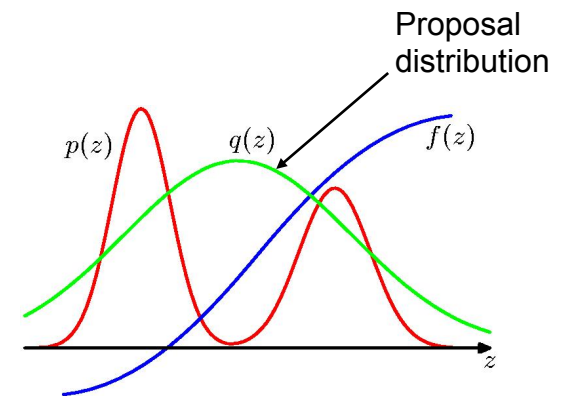
- Precisely the form of importance sampling estimator
 - Can generate samples from Q and correct by weights
 - Simplify by choosing Q to be P_{θ^0} for some θ^0

$$\begin{aligned} Z(\theta) &= E_{P_{\theta^0}} \left[\frac{Z(\theta^0) \exp \left\{ \sum_i \theta_i f_i(\chi) \right\}}{\exp \left\{ \sum_i \theta_i^0 f_i(\chi) \right\}} \right] \\ &= Z(\theta^0) E_{P_{\theta^0}} \left[\exp \left\{ \sum_i (\theta_i - \theta_i^0) f_i(\xi_k) \right\} \right] \end{aligned}$$

Importance sampling

- Samples $\{z^{(l)}\}$ are drawn from simpler dist. $q(z)$

$$\begin{aligned} E[f] &= \int f(z) p(z) dz \\ &= \int f(z) \frac{p(z)}{q(z)} q(z) dz \\ &= \frac{1}{L} \sum_{l=1}^L \frac{p(z^{(l)})}{q(z^{(l)})} f(z^{(l)}) \end{aligned}$$



Unlike rejection sampling
All of the samples are retained

- Samples are weighted by ratios $r_l = p(z^{(l)}) / q(z^{(l)})$
 - Known as importance weights
 - Which corrects the bias introduced by wrong distribution

Samples to approximate $\ln Z(\theta)$

- Choose Q to be P_{θ^0} for some parameters θ^0 .
- Sample instances ξ^1, \dots, ξ^K from P_{θ^0} , to approximate the log-partition function as

$$\ln Z(\theta) \approx \ln \left(\frac{1}{K} \sum_{k=1}^K \exp \left\{ \sum_i (\theta_i - \theta_i^0) f_i(\xi^k) \right\} \right) + \ln Z(\theta^0)$$

- Plug it into $\frac{1}{M} \ell(\theta : D) = \sum_i \theta_i (E_D[f_i(d_i)]) - \ln Z(\theta)$ and optimize
 - Use a sampling procedure to generate samples from current parameter set θ^t
 - Then use gradient descent to find θ^{t+1} that improves log-likelihood based on the samples

MAP-based Learning

- Another approach to inference in learning
 - Approximating expected feature counts with the counts in the single MAP assignment to current MN
- Approximate gradient at assignment θ is

$$E_D[f_i(\chi)] - f_i(\xi^{MAP}(\theta))$$

- where $\xi^{MAP}(\theta) = \arg \max_{\xi} P(\xi|\theta)$ is the MAP assignment given the current set of parameters θ
- Approach also called as *Viterbi training*
- Equivalent to exact optimization of approximate objective

$$\frac{1}{M} \ell(\theta : D) - \ln P(\xi^{MAP}(\theta) | \theta)$$

Alternative Objectives

- A class of approximations obtained by replacing objective with one more tractable

- In the case of a single instance ξ

$$\ell(\theta : \xi) = \ln \tilde{P}(\xi | \theta) - \ln Z(\theta)$$

- Expanding the partition function

$$\ell(\theta : \xi) = \ln \tilde{P}(\xi | \theta) - \ln \left(\sum_{\xi'} \tilde{P}(\xi' | \theta) \right)$$

Summing over all possible values of dummy variable

- Maximizing ℓ is to increase distance (*contrast*) between two terms
 - First term: Log measure
 - Unnormalized measure (log-probability) of ξ
 - Second Term: Aggregate measure of all instances

Contrastive Objective

- Since log-measure increases with parameters
 - increase parameters with positive empirical expectations in ξ and decrease parameters with negative empirical expectations
- However second term balances the first
 - Key difficulty
 - second term has exponential instances in $Val(\chi)$
 - and requires inference in the network
- Approach
- Increase log-measure of data instances and a more tractable set of other instances, one not requiring summation over an exponential space

Two Approaches to increase probability gap

1. Pseudolikelihood and its generalizations

- Simplifies likelihood by considering only pairwise dependencies

2. Contrastive Optimization

- Contrast data with a randomly perturbed set of neighbors

Pseudolikelihood for Tractability

- Consider likelihood of single instance ξ :

$$P(\xi) = \prod_{j=1}^n P(x_j | x_1, \dots, x_{j-1})$$

Since from product rule $P(x_1, x_2) = P(x_1) P(x_2 | x_1)$ and $P(x_1, x_2, x_3) = P(x_1) P(x_2, x_3 | x_1) = P(x_1) P(x_2 | x_1) P(x_3 | x_1, x_2)$

- Approximate by

$$P(\xi) = \prod_{j=1}^n P(x_j | x_1, \dots, x_{j-1}, x_{j+1}, \dots, x_n)$$

Replace each product term above by conditional probability of x_j given all the other variables

- Gives the pseudolikelihood objective

$$\ell_{PL}(\theta : D) = \frac{1}{M} \sum_m \sum_j \ln P(x_j[m] | x_{-j}[m], \theta)$$

where x_{-j} stands for $x_1, \dots, x_{j-1}, x_{j+1}, \dots, x_n$

– which eliminates exponential summation

$$P(x_j | x_{-j}) = \frac{P(x_j, x_{-j})}{P(x_{-j})} = \frac{\tilde{P}(x_j, x_{-j})}{\sum_{x'_j} \tilde{P}(x'_j, x_{-j})}$$

Requires only summation over X_j

Pseudolikelihood is concave

- Pseudolikelihood objective of a single data ξ :

$$\begin{aligned}\sum_j \ln P(x_j | x_{-j}) &= \sum_j \left(\ln \tilde{P}(x_j, x_{-j}) - \ln \sum_{x'_j} \tilde{P}(x'_j, x_{-j}) \right) \text{ second term is partition fn.} \\ &= \sum_j \left(\ln \tilde{P}(\xi) - \ln \sum_{x'_j} \tilde{P}(x'_j, x_{-j}) \right)\end{aligned}$$

- Simplify each term in summation as

$$\ln P(x_j | x_{-j}) = \left(\sum_{i: \text{Scope}[f_i] \ni X_j} \theta_i f_i(x_j, x_{-j}) \right) - \ln \left(\sum_{x'_j} \exp \left\{ \sum_{i: \text{Scope}[f_i] \ni X_j} \theta_i f_i(x'_j, x_{-j}) \right\} \right)$$

- Each term is a log-conditional-likelihood term for a MN for a single variable X_j conditioned on rest

- Function is concave in the parameters θ
- Gradient is

$$\frac{\partial}{\partial \theta_i} \ln(x_j | x_{-j}) = f_i(x_j, x_{-j}) - E_{x'_j \sim P_\theta(X_j | x_{-j})} [f_i(x'_j, x_{-j})]$$

Contrastive Optimization

- Likelihood and Pseudolikelihood

$$\ell(\theta : \xi) = \ln \tilde{P}(\xi | \theta) - \ln \left(\sum_{\xi'} \tilde{P}(\xi' | \theta) \right)$$

- Both attempt to increase log-probability gap between D and probability of a set of instances
- Base on this intuition, a range of methods developed:
 - By driving the probability of the observed data higher relative to other instances, we are tuning the parameters to predict the data better

Contrastive Optimization Criteria

- Aim to maximize the log-probability gap
- Consider a single training instance ξ

– Maximize the *contrastive objective*

$$\ln \tilde{P}(\xi | \theta) - \ln \tilde{P}(\xi' | \theta)$$

- where ξ' is some other instance
- This takes a simple form

$$\theta^T [f(\xi) - f(\xi')]$$

–which is a linear function of θ and hence unbounded

- Choice of ξ' has to change throughout the optimization
- Two methods for choosing ξ' :
 1. Contrastive divergence
 2. Margin-based training

Contrastive Divergence

- Popularity of method has grown
 - Used in *deep learning*– for training layers of RBMs
- Contrast data instances \mathcal{D} with set of randomly perturbed neighbors \mathcal{D}^-
 - Maximize
$$\ell_{CD}(\theta : \mathcal{D} || \mathcal{D}^-) = E_{\xi \sim \tilde{P}_D} [\ln \tilde{P}_\theta(\xi)] - E_{\xi \sim \tilde{P}_{D^-}} [\ln \tilde{P}_\theta(\xi)]$$
 - Where P_D and P_{D^-} are distributions relative to D and D^-
 - We want model to give high probability to instances in D relative to the perturbed instances in D^-

Generating \mathcal{D}^- for Contrastive Divergence

- Contrasted instances \mathcal{D}^- will differ at different stages in the search
- Given current parameters θ generate samples \mathcal{D}^- from P_θ using Gibbs sampling
 - Initialize from the instances in \mathcal{D} , run the Markov chain only a few steps to define \mathcal{D}^-
- Gradient of objective is

$$\frac{\partial}{\partial \theta_i} \ell_{CD}(\theta : \mathcal{D} || \mathcal{D}^-) = E_{\hat{P}_D} [f_i(\chi)] - E_{\hat{P}_{D^-}} [f_i(\chi)]$$

Margin-Based Training

- When the goal is MAP assignment, an SVM-based method can be used
- Training set consists of pairs $D = \left\{ \left(y[m], x[m] \right) \right\}_{m=1}^M$
- Given observation $x[m]$ we would like learned model to give highest probability to $y[m]$
- Maximize the margin
$$\ln P_{\theta}(y[m] | x[m]) - \left[\max_{y \neq y[m]} \ln P_{\theta}(y[m] | x[m]) \right]$$
 - the difference between the log-probability of the target assignment $y[m]$ and “next best” assignment