# Learning MN Parameters with Alternative Objective Functions

Sargur Srihari

srihari@cedar.buffalo.edu

# Topics

- Max Likelihood & Contrastive Objectives
- Contrastive Objective Learning Methods
  - Pseudo-likelihood
    - Gradient descent
    - Ex: Use in Collaborative Filtering
  - Contrastive Optimization criteria
    - Contrastive Divergence
    - Margin-based Training

# Recapitulate ML learning

- ## Task: Estimate parameters of a MN

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

  - with a fixed structure given data set $\mathcal{D}$

- ## Simplest is to maximize log-likelihood objective given samples $\xi_1,\cdots\xi_M$, which is

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

  - Although concave, no analytical form for maximum

    - Can use iterative gradient ascent in param. space

$$\frac{\partial}{\partial\theta_i}\frac{1}{M}\ell(\theta:D) = E_D[f_i(\chi)] - E_\theta[f_i]$$   since   $$\frac{\partial}{\partial\theta_i}\ln Z(\theta) = \frac{1}{Z(\theta)}\sum_{\xi}f_i(\xi)\exp\left\{\sum_j\theta_j f_j(\xi)\right\} = E_\theta[f_i]$$

      – Good news: Likelihood is concave, Guaranteed to converge
      – Bad news:  Each step needs inference; estimation is intractable

3

# Replacing ML objective

- Previously seen approx. inference methods:
  - Belief propagation, MaxEnt, Sampling +SGD
- Now we look at replacing the objective function with one that is more tractable
- The ML objective is: $\ell(\theta:\mathcal{D})=\sum_i\theta_i\left(\sum_m f_i(\xi[m])\right)-M\ln Z(\theta)$  $Z(\theta)=\sum_\xi\exp\left\{\sum_i\theta_i f_i(\xi)\right\}$
- For simplicity focus on the case of a single data instance $\xi$ :

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

# Log-likelihood of a single sample

- In the case of a single instance $\xi$

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

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

- Expanding the partition function $Z(\theta)$

$$\ell\big(\theta : \xi\big) = \ln \tilde{P}\big(\xi \mid \theta\big) - \ln\left[\sum_{\xi'} \tilde{P}\big(\xi' \mid \theta\big)\right]$$

Summing over all possible values of dummy variable $\xi'$

- Maximizing $\ell$ is to increase distance (*contrast*) between the two terms

- Consider each of the two terms separately

$$\ln \tilde{P}(\xi \mid \theta)$$    and    $$\ln \tilde{P}(\xi \mid \theta) = -\sum_{i=1}^{k} \theta_i f_i(D_i)$$

# First Term of Objective

- First term is:

$$\boxed{\ln \tilde{P}(\xi \mid \theta)}$$

- Objective aims to increase the log measure
  - i.e., Logarithm of unnormalized probability of observed data instance $\xi$
  - log measure is a linear function of parameters in log-linear representation, $\boxed{\ln \tilde{P}(\xi \mid \theta) = -\sum_{i=1}^{k} \theta_i f_i(D_i)}$
  - Thus that goal can be accomplished by
    - Increasing all parameters associated with positive empirical expectations in $\xi$ and decreasing all parameters associated with negative empirical expectations
  - We can increase it unboundedly using this approach

6

# Second Term of Objective

- Second Term:   $\ln\left(\sum_{\xi'} \tilde{P}(\xi' \mid \theta)\right)$

- It is the logarithm of the sum of unnormalized measures of all possible instances in $Val(\chi)$

  – It is the aggregate of the measures of all instances

# The Contrastive Objective Approach

- Can view log-likelihood objective as
  - aiming to increase distance between log measure of $\xi$ and aggregate of the measures of all instances

- The key difficulty with this formulation
  - second term has exponential instances in $Val(\chi)$ and requires inference in the network

- It suggests an approach to approximation:
  - We can move our parameters in the right direction if we aim to increase the difference between
    1. The log-measure of data instances and
    2. 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**
   - Easiest method circumventing intractability of network inference
   - Simplifies likelihood by replacing exponential no. of summations with several summations, each more tractable

2. **Contrastive Optimization**
   - Drive probability of observed samples higher
   - Contrast data with a randomly perturbed set of neighbors

9

# Pseudolikelihood for Tractability

- Consider likelihood of single instance $\xi$ :

$$P\left(\xi\right) = \prod_{j=1}^{n} P\left(x_j \mid x_1,..,x_{j-1}\right)$$

From chain rule $P(x_1,x_2) = P(x_1) \ P(x_2|x_1)$ and
$P(x_1,x_2,x_3) = P(x_1) \ P(x_2|x_1) \ P(x_3|x_1.x_2)$

  – Approximate by replacing each product term by conditional probability of $x_j$ given all other variables

$$P\left(\xi\right) = \prod_{j=1}^{n} P\left(x_j \mid x_1,..,x_{j-1}, \ x_{j+1},..,x_n\right)$$

  – Gives the pseudolikelihood (PL) objective

$$\ell_{PL}\left(\boldsymbol{\theta} : \mathcal{D}\right) = \frac{1}{M}\sum_{m}\sum_{j}\ln P\left(x_j\left[m\right] \mid \boldsymbol{x}_{-j}\left[m\right],\boldsymbol{\theta}\right)$$

   where $\boldsymbol{x}_{-j}$ stands for $x_1,.., \ x_{j-1}, \ x_{j+1},.., \ x_n$

  – This objective measures ability to predict each variable given full observation of all other variables

10

# Pseudolikelihood and Multinomial

- The predictive model takes a form that generalizes the multinomial logistic CPD

$P(Y|X_1,\ldots,X_k)$ such that
$$P(y^j \mid X_1,\ldots,X_k) = \frac{\exp(\ell_j(X_1,\ldots,X_k))}{\sum_{j'=1}^{m}\exp(\ell_{j'}(X_1,\ldots,X_k))}$$
$$\ell_j(X_1,\ldots,X_k) = w_{j,0} + \sum_{i=1}^{k} w_{j,i}X_i$$

  – Identical to it when the network consists of only pairwise features

    • Factors over edges in the network

  – We can use conditional independence properties in the network to simplify the expression,

$$\ell_{PL}\big(\boldsymbol{\theta}:\mathcal{D}\big) = \frac{1}{M}\sum_{m}\sum_{j}\ln P\big(x_j[m] \mid \boldsymbol{x}_{-j}[m],\boldsymbol{\theta}\big)$$

  – removing from the rhs of $P(X_j|\boldsymbol{X}_{-j})$ any variable that is not a neighbor of $X_j$

# Simplification in Pseudolikelihood

- The pseudolikelihood ($\mathrm{PL}$) objective is

$$\ell_{PL}\left(\boldsymbol{\theta}:\mathcal{D}\right)=\frac{1}{M}\sum_{m}\sum_{j}\ln P\left(x_{j}\left[m\right]\mid \boldsymbol{x}_{-j}\left[m\right],\boldsymbol{\theta}\right)$$

  – Whereas likelihood is $\ell(\boldsymbol{\theta}:\mathcal{D})=\sum_{i}\theta_{i}\left(\sum_{m}f_{i}(\xi[m])\right)-M\ln Z(\boldsymbol{\theta})$      $\ln Z(\theta)=\ln\sum_{\xi}\exp\left\{\sum_{i}\theta_{i}f_{i}(\xi)\right\}$

- Pseudolikelihood eliminates exponential summation over instances with several summations, each of which is more tractable

  – In particular $P(x_{j}\mid \boldsymbol{x}_{-j})=\dfrac{P(x_{j},\boldsymbol{x}_{-j})}{P(\boldsymbol{x}_{-j})}=\dfrac{\tilde{P}(x_{j},\boldsymbol{x}_{-j})}{\sum_{x_{j}'}\tilde{P}(x_{j}',\boldsymbol{x}_{-j})}$

    - Global partition function has disappeared
    - Requires only summation over $X_{j}$

- But there is a contrastive perspective

  – Described next

12

# Contrastive perspective of PL

- Pseudolikelihood objective of single data $\xi$:

$$\sum_j \ln P\left(x_j \mid \boldsymbol{x}_{-j}\right) = \sum_j \left( \ln \tilde{P}\left(x_j, \boldsymbol{x}_{-j}\right) - \ln \sum_{x'_j} \tilde{P}\left(x'_j, \boldsymbol{x}_{-j}\right) \right)$$

$$= \sum_j \left( \ln \tilde{P}\left(\xi\right) - \ln \sum_{x'_j} \tilde{P}\left(x'_j, \boldsymbol{x}_{-j}\right) \right)$$

- Each term of final sum is a contrastive term
  - Where we aim to increase difference between log-measure of training instance $\xi$ and an aggregate of log-measures of instances that differ from $\xi$ in the assignment to precisely one variable
    - In other words we are increasing the contrast between our training instance $\xi$ and the instances in a local neighborhood around it

13

# Pseudolikelihood is concave

- Further simplification of summands in the expression $\sum_j \ln P\left(x_j \mid \boldsymbol{x}_{-j}\right) = \sum_j \left[ \ln \tilde{P}(\xi) - \ln \sum_{x'_j} \tilde{P}\left(x'_j, \boldsymbol{x}_{-j}\right) \right]$ obtains

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

  - Each term is a log-conditional-likelihood term for a MN for a single variable $X_j$ conditioned on rest
  - Thus it follow that the function is concave in the parameters $\theta$

- Since a sum of concave functions is concave the pseudolikelihood objective

  - $\ell_{PL}\left(\boldsymbol{\theta} : \mathcal{D}\right) = \dfrac{1}{M} \sum_m \sum_j \ln P\left(x_j[m] \mid \boldsymbol{x}_{-j}[m], \boldsymbol{\theta}\right)$ is also concave

# Gradient of pseudolikelihood

- ## To compute gradient we use $\ln P(x_j \mid x_{-j}) = \left( \sum\limits_{i:Scope[f_i] \ni X_j} \theta_i f_i\left(x_j, x_{-j}\right) \right) - \ln\left( \sum\limits_{x_j'} \exp\left\{ \sum\limits_{i:Scope[f_i] \ni X_j} \theta_i f_i\left(x'_j, x_{-j}\right) \right\} \right)$

  - to obtain

    $$\frac{\partial}{\partial \theta_i} \ln\left(x_j \mid \boldsymbol{x}_{-j}\right) = f_i(x_j, \boldsymbol{x}_{-j}) - E_{x_j' \sim P_\theta\left(X_j \mid x_{-j}\right)}\left[ f_i(x_j', \boldsymbol{x}_{-j}) \right]$$

    - If $X_j$ is not in the scope of $f_i$ then $f_i(x_j, \boldsymbol{x}_{-j}) = f_i(x_j', \boldsymbol{x}_{-j})$ for any $x_j'$ and the two terms are identical making the derivative $0$. Inserting this into $\ell_{\text{PL}}$ we get

    $$\frac{\partial}{\partial \theta_i} \ell_{\text{PL}}\left(\boldsymbol{\theta} : \mathcal{D}\right) = \sum_{j:X_j \in Scope[f_i]} \left( \frac{1}{M} \sum_m f_i(\xi[m]) - E_{x_j' \sim P_\theta(X_j \mid \boldsymbol{x}_{-j}[m])}\left[ f_i(x_j', \boldsymbol{x}_{-j}[m]) \right] \right)$$

  - It is much easier to compute this than $\frac{\partial}{\partial \theta_i} \frac{1}{M} \ell(\boldsymbol{\theta} : D) = E_D[f_i(\chi)] - E_\theta[f_i]$

  - Each expectation term requires a summation over a single random variable $X_j$, conditioned on all of its neighbors

    - A computation that can be performed very efficiently

15

# Relationship between likelihood and pseudolikelihood

- *Theorem*: Assume that our data are generated by a log-linear model $P_{\theta*}$, i.e.,

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

  – Then as $M$ goes to infinity, with probability approaching $1$, $\theta*$ is a global optimum of the pseudolikelihood objective

$$\ell_{PL}(\boldsymbol{\theta}:\mathcal{D}) = \frac{1}{M}\sum_m\sum_j \ln P\left(x_j[m] \mid \boldsymbol{x}_{-j}[m],\boldsymbol{\theta}\right)$$

- Result is an important, but there are limitations

  – Model being learnt must be expressive

    • But model never perfectly represents true distribution

  – Data distribution must be near generating distribution

    • Often not enough data to approach large sample limit [16]

# Type of queries determine method

- How good is the pseudolikelihood objective depends on the type of queries
  - If we condition on most of variables and condition on few, pseudolikelihood is a very close match to the type of predictions we would like to make
    - So pseudolikelihood may well be better than likelihood
    - Ex: in learning a MN for collaborative filtering, we take user's preference for all items as observed  except the query item
  - Conversely, if a typical query involves most or all variables, the likelihood objective is better
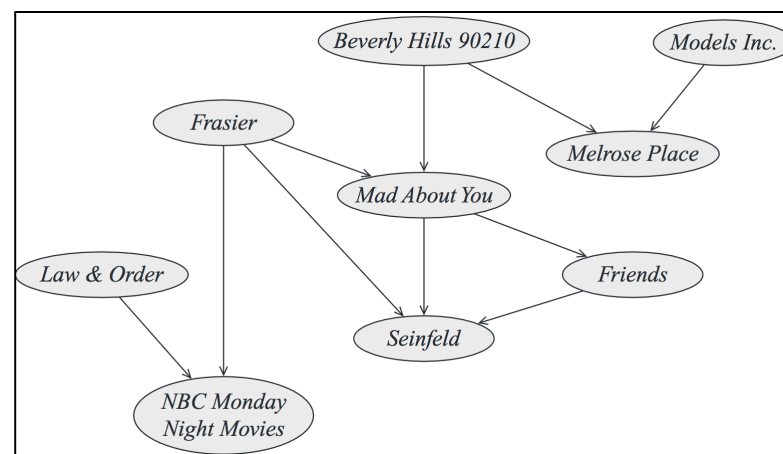    - In learning a CRF model for image segmentation

17

# Ex: Collaborative Filtering

- We wan't to recommend to a user an item based on previous items he bought/liked

- Because we don't have enough data for any single user to determine his/her preference, we use *collaborative filtering*

  – i.e., use observed preferences of others to determine preferences for any other user

  – One approach: learn dependency structure between different purchases as observed in the population

    - Item $i$ is a variable $X_i$ in a joint distribution and each user is an instance

    - View purchase/non-purchase as values of a variable
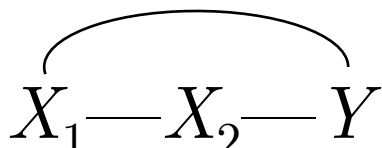
18

# A Collaborative Filtering Task

- Given a set of purchases $S$ we can compute the probability that user would like new item $i$

- This is a probabilistic inference task where all purchases other than $i$ are set to $False$; thus all variables other than $X_i$ are observed

If we condition on most of variables and condition on few,
Pseudolikelihood ($\text{PL}$) is a
very close match to the type of predictions we would like to make



Network learned from Nielsen TV rating data capturing viewing record of sample viewers. Variables denote whether a TV program was watched

19

# Limitation of pseudolikelihood (PL)

- ## Ex: MN with 3 variables: $\overparen{X_1 \!-\! X_2 \!-\! Y}$

  - $X_1$ and $X_2$ highly correlated, Both not as correlated to $Y$
  - Best predictor for $X_1$ is $X_2$ and vice versa
    - Pseudolikelihood likely to overestimate $X_1 \!-\! X_2$ parameters and almost entirely dismiss $X_1 \!-\! Y$ and $X_2 \!-\! Y$
    - Excellent predictor for $X_2$ when $X_1$ is observed, but useless when only Y, not $X_1$, is observed

- ## PL assumes local neighborhood is observed

  - cannot exploit weaker or long-range dependencies

- ## Solution: generalized pseudolikelihood ($\mathrm{GPL}$)

We define a set of subsets of variables $\{X_s : s \in S\}$ and define an objective:

$$\ell_{GPL}(\boldsymbol{\theta} : \mathcal{D}) = \frac{1}{M} \sum_m \sum_s \ln P(\boldsymbol{x}_s[m] \mid \boldsymbol{x}_{-s}[m], \boldsymbol{\theta})$$

where $\boldsymbol{X}_{-s} = \chi - \boldsymbol{X}_s$

20

# Contrastive Optimization Criteria

- Both Likelihood and Pseudolikelihood

  – attempt to increase log-probability gap between probability of observed set of $m$ instances $\mathcal{D}$ and logarithm of the aggregate probability of a set of instances:

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

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

$$\ell_{PL}(\boldsymbol{\theta}:\mathcal{D}) = \frac{1}{M}\sum_m \sum_j \ln P\left( x_j[m]\mid \boldsymbol{x}_{-j}[m], \boldsymbol{\theta} \right)$$

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

- Based on this intuition, a range of methods developed to increase log-probability gap

  – By driving probability of observed data higher relative to other instances,

    - we are tuning the parameters to predict the data better

# Contrastive Objective Definition

- Consider a single training instance $\xi$

  – We aim to maximize the *log-probability gap*

  $$\ln \tilde{P}\left(\xi \mid \theta\right) - \ln \tilde{P}\left(\xi' \mid \theta\right)$$

    - where $\xi'$ is some other instance, whose selection we discuss shortly

- Importantly, the expression takes a simple form

  $$\ln \tilde{P}\left(\xi \mid \boldsymbol{\theta}\right) - \ln \tilde{P}\left(\xi' \mid \boldsymbol{\theta}\right) = \boldsymbol{\theta}^{T}[\boldsymbol{f}(\xi) - \boldsymbol{f}(\xi')]$$

  – For a fixed instantiation of $\xi'$, this expression is a linear function of $\theta$ and hence is unbounded

  – For a coherent optimization objective, choice of $\xi'$ has to change throughout the optimization

    - Even then, take care to prevent unbounded parameters

22

# Two Contrastive Optimization Methods

- One can construct many variants of this method
  - i,e., methods to increase log-probability gap
- Two methods for choosing $\xi$' that have been useful in practice:
  1. Contrastive divergence
     - Popularity of method has grown
     - Used in *deep learning*– for training layers of RBMs
  2. Margin-based training

# Contrastive Divergence ($CD$)

- In this method we contrast our data instances $\mathcal{D}$ with set of randomly perturbed *neighbors* $\mathcal{D}^{-}$

  – We aim to maximize

  $$\ell_{CD}\left(\theta : \mathcal{D} \,\|\, \mathcal{D}^{-}\right) = E_{\xi \sim \tilde{P}_{\mathcal{D}}}\left[\ln \tilde{P}_{\theta}\left(\xi\right)\right] - E_{\xi \sim \tilde{P}_{\mathcal{D}^{-}}}\left[\ln \tilde{P}_{\theta}\left(\xi\right)\right]$$

  - where $P_{\mathcal{D}}$ and $P_{\mathcal{D}^{-}}$ are the empirical distributions relative to $\mathcal{D}$ and $\mathcal{D}^{-}$

- The set of contrasted instances $\mathcal{D}^{-}$ will necessarily differ at different stages in the search

  – Choosing instances to contrast is next

# Choice of Data instances

- Given current $\theta$ , what instances to which we want to contrast our data instances $\mathcal{D}$?

- One intuition: move $\theta$ in a direction that increases probability of instances in $\mathcal{D}$ relative to "typical" instances in current distribution

  – i.e., increase probability gap between instances $\xi \varepsilon$ $\mathcal{D}$ and instances $\xi$ sampled randomly from $P_\theta$

- Thus, we can generate a contrastive set $\mathcal{D}^-$ by sampling from $P_\theta$ and then maximizing the objective in

$$\ell_{CD}\left(\theta : \mathcal{D} \parallel \mathcal{D}^-\right) = E_{\xi \sim \tilde{P}_{\mathcal{D}}}\left[\ln \tilde{P}_\theta\left(\xi\right)\right] - E_{\xi \sim \tilde{P}_{\mathcal{D}^-}}\left[\ln \tilde{P}_\theta\left(\xi\right)\right]$$

25

# How to sample from $P_\theta$ ?

- We can run a Markov chain defined by the MN $P_\theta$ using Gibbs sampling and initializing from the instances in $\mathcal{D}$

  – Once the chain mixes we can collect samples from the distribution

  – Unfortunately, sampling from the chain for long enough to achieve mixing takes far too long for the inner loop of a learning algorithm

  – So we initialize from instances in $\mathcal{D}$ and run the chain only for a few steps

    • Instances from these short sampling runs define $\mathcal{D}^-$

# Updating parameters in CD

- We want model to give high probability to instances in $\mathcal{D}$ relative to the perturbed instances in $\mathcal{D}$-

  - Thus we want to move our parameters in a direction that increases the probability of instances in $\mathcal{D}$ relative to the perturbed instances in $\mathcal{D}$-

- Gradient of objective is easy to compute

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

  - In practice approximation we get by taking only a few steps in the Markov chain provides a good direction for the search

# Margin-Based Training

- Very different intuition in settings where our goal is to use the network for predicting a MAP assignment

- Ex: in image segmentation we want the lerned network to predict a single high probability assignment to the pixels that will encode pur final segmentation output

- Occurs only when queries are conditional

- So we describe objective for a CRF

# Margin-Based Training

- Training set consists of pairs $D = \left\{ \left( \boldsymbol{y}[m], \boldsymbol{x}[m] \right) \right\}_{m=1}^{M}$

- Given observation $\boldsymbol{x}[m]$ we would like learned model to give highest probability to $\boldsymbol{y}[m]$

  - i.e., we would like $P_\theta \left( \boldsymbol{y}[m] | \boldsymbol{x}[m] \right)$ to be higher than any other probability $P_\theta \left( \boldsymbol{y} | \boldsymbol{x}[m] \right)$ for $\boldsymbol{y} \neq \boldsymbol{y}[\mathrm{m}]$

  - i.e., Maximize the margin

    $$\ln P_\theta \left( y[m] \,|\, x[m] \right) - \left[ \max_{y \neq y[m]} \ln P_\theta \left( y[m] \,|\, x[m] \right) \right]$$

    - the difference between the log-probability of the target assignment $y[m]$ and "next best" assignment