

# TTIC 31230, Fundamentals of Deep Learning

David McAllester, Winter 2019

## Approximate Algorithms

for Exponential Softmax Backpropagation

# Distributions on Exponentially Large Sets

$$\Phi^* = \operatorname{argmin}_{\Phi} E_{(x,y) \sim \text{Pop}} - \ln P(y|x)$$

$$\Phi^* = \operatorname{argmin}_{\Phi} E_{y \sim \text{Pop}} - \ln P(y)$$

The structured case:  $y \in \mathcal{Y}$  where  $\mathcal{Y}$  is discrete but iteration over  $\hat{y} \in \mathcal{Y}$  is infeasible.

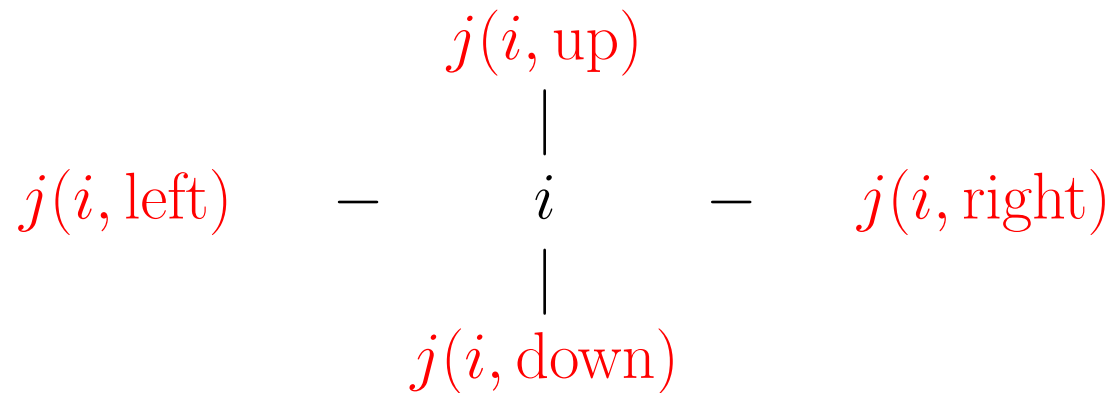
# Semantic Segmentation

We want to assign each pixel to one of  $C$  semantic classes.

For example “person”, “car”, “building”, “sky” or “other”.

## Constructing a Graph

We construct a graph whose nodes are the pixels and where there is an edges between each pixel and its four nearest neighboring pixels.



## Labeling the Nodes of the Graph

$\hat{y}$  assigns a semantic class  $\hat{y}[i]$  to each node (pixel)  $i$ .

We assign a score to  $\hat{y}$  by assigning a score to each node and each edge of the graph.

$$s(\hat{y}) = \sum_{i \in \text{Nodes}} s_n[i, \hat{y}[i]] + \sum_{\langle i, j \rangle \in \text{Edges}} s_e[\langle i, j \rangle, \hat{y}[i], \hat{y}[j]]$$

Node Scores                      Edge Scores

## Computing the Node and Edge Tensors

For input  $x$  we use a network to compute the score tensors.

$$s_n[I, C] = f_{\Phi}^n(x)$$

$$s_e[E, C, C] = f_{\Phi}^e(x)$$

## Exponential Softmax

$$\text{for } \hat{y} \quad s(\hat{y}) = \sum_i s_n[i, \hat{y}[i]] + \sum_{\langle i, j \rangle \in \text{Edges}} s_e[\langle i, j \rangle, \hat{y}[i], \hat{y}[j]]$$

$$\text{for } \hat{y} \quad P_s(\hat{y}) = \text{softmax}_{\hat{y}} s(\hat{y}) \quad \text{all possible } \hat{y}$$

$$\mathcal{L} = -\ln P_s(y) \quad \text{gold label (training label) } y$$

## Exponential Softmax is Typically Intractable

$\hat{y}$  assigns a label  $\hat{y}[i]$  to each node  $i$ .

$s(\hat{y})$  is defined by a sum over node and edge tensor scores.

$P_s(\hat{y})$  is defined by an exponential softmax over  $s(\hat{y})$ .

Computing  $Z$  in general is  $\#P$  hard (there is an easy direct reduction from SAT).



## Compactly Representing Scores on Exponentially Many Labels

The tensor  $s_n[I, C]$  holds  $IC$  scores.

The tensor  $s_e[E, C, C]$  holds  $EC^2$  scores where  $e$  ranges over edges  $\langle i, j \rangle \in \text{Edges}$ .

## Back-Propagation Through Exponential Softmax

$$\begin{aligned}s_n[I, C] &= f_{\Phi}^n(x) \\ s_e[E, C, C] &= f_{\Phi}^e(x)\end{aligned}$$

$$s(\hat{y}) = \sum_i s_n[i, \hat{y}[i]] + \sum_{\langle i, j \rangle \in \text{Edges}} s_e[\langle i, j \rangle, \hat{y}[i], \hat{y}[j]]$$

$$P_s(\hat{y}) = \text{softmax}_{\hat{y}} s(\hat{y}) \quad \text{all possible } \hat{y}$$

$$\mathcal{L} = -\ln P_s(y) \quad \text{gold label } y$$

We want the gradients  $s_n.\text{grad}[I, C]$  and  $s_e.\text{grad}[E, C, C]$ .

## Model Marginals Theorem

Theorem:

$$s_n.\text{grad}[i, c] = P_{\hat{y} \sim P_s}(\hat{y}[i] = c) - \mathbf{1}[y[i] = c]$$

$$s_e.\text{grad}[\langle i, j \rangle, c, c'] = P_{\hat{y} \sim P_s}(\hat{y}[i] = c \wedge \hat{y}[j] = c') - \mathbf{1}[y[i] = c \wedge y[j] = c']$$

We need to compute (or approximate) the model marginals.

## Proof of Model Marginals Theorem

We consider the case of node marginals, The case of edge marginals is similar.

$$\begin{aligned} s_n.\text{grad}[i, c] &= \partial \mathcal{L}(\Phi, x, y) / \partial s_n[i, c] \\ &= \partial \left( -\ln \frac{1}{Z} \exp(s(y)) \right) / \partial s_n[i, c] \\ &= \partial (\ln Z - s(y)) / \partial s_n[i, c] \\ &= \left( \frac{1}{Z} \sum_{\hat{y}} e^{s(\hat{y})} (\partial s(\hat{y}) / \partial s_n[i, c]) \right) - (\partial s(y) / \partial s_b[i, c]) \end{aligned}$$

## Proof of Model Marginals Theorem

$$\begin{aligned}
 s_n.\text{grad}[i, c] &= \left( \frac{1}{Z} \sum_{\hat{y}} e^{s(\hat{y})} (\partial s(\hat{y}) / \partial s_n[i, c]) \right) - (\partial s(y) / \partial s_b[i, c]) \\
 &= \left( \sum_{\hat{y}} P_s(\hat{y}) (\partial s(\hat{y}) / \partial s_n[i, c]) \right) - (\partial s(y) / \partial s_n[i, c]) \\
 s(\hat{y}) &= \sum_i s_n[i, \hat{y}[i]] + \sum_{\langle i, j \rangle \in \text{Edges}} s_e[\langle i, j \rangle, \hat{y}[i], \hat{y}[j]] \\
 \frac{\partial s(\hat{y})}{\partial s_n[i, c]} &= \mathbf{1}[\hat{y}[i] = c]
 \end{aligned}$$

## Proof of Model Marginals Theorem

$$\begin{aligned} s_n.\text{grad}[i, c] &= \left( \frac{1}{Z} \sum_{\hat{y}} e^{s(\hat{y})} (\partial s(\hat{y}) / \partial s_n[i, c]) \right) - (\partial s(y) / \partial s_b[i, c]) \\ &= \left( \sum_{\hat{y}} P_s(\hat{y}) (\partial s(\hat{y}) / \partial s_n[i, c]) \right) - (\partial s(y) / \partial s_n[i, c]) \\ &= E_{\hat{y} \sim P_s} \mathbf{1}[\hat{y}[i] = c] - \mathbf{1}[y[i] = c] \\ &= P_{\hat{y} \sim P_s}(\hat{y}[i] = c) - \mathbf{1}[y[i] = c] \end{aligned}$$

# Model Marginals Theorem

Theorem:

$$s_n.\text{grad}[i, c] = P_{\hat{y} \sim P_s} ( \textcolor{red}{\hat{y}}[i] = c ) \\ - \mathbf{1}[ \textcolor{red}{y}[i] = c ]$$

$$s_e.\text{grad}[\langle i, j \rangle, c, c'] = P_{\hat{y} \sim P_s} ( \textcolor{red}{\hat{y}}[i] = c \wedge \textcolor{red}{\hat{y}}[j] = c' ) \\ - \mathbf{1}[ \textcolor{red}{y}[i] = c \wedge \textcolor{red}{y}[j] = c' ]$$

# Methods of Approximating Model Marginals

Monte Carlo Markov Chain (MCMC) Sampling

Pseudolikelihood

Contrastive Divergence

Loopy Belief Propagation (loopy BP)



## MCMC Sampling

The model marginals, such as the node marginals  $P_s(\hat{y}[i] = c)$ , can be estimated by sampling  $\hat{y}$  from  $P_s(\hat{y})$ .

There are various ways to design a Markov process whose states are node labelings  $\hat{y}$  and whose stationary distribution is  $P_s$ .

Given such a process we can sample  $\hat{y}$  from  $P_s$  by running the process past its mixing time.

We will consider Metropolis MCMC and the Gibbs MCMC. But there are more (like Hamiltonian MCMC).

## Metropolis MCMC

We assume a neighbor relation on node assignments and let  $N(\hat{y})$  be the set of neighbors of assignment  $\hat{y}$ .

For example,  $N(\hat{y})$  can be taken to be the set of assignments  $\hat{y}'$  that differ from  $\hat{y}$  on exactly one node.

For the correctness of Metropolis MCMC we need that all states have the same number of neighbors and that the neighbor relation is symmetric —  $\hat{y}' \in N(\hat{y})$  if and only if  $\hat{y} \in N(\hat{y}')$ .

## Metropolis MCMC

Pick an initial state  $\hat{y}_0$  and for  $t \geq 0$  do

1. Pick a neighbor  $\hat{y}' \in N(\hat{y}_t)$  uniformly at random.

2. If  $P_s(\hat{y}') > P_s(\hat{y}_t)$  then  $\hat{y}_{t+1} = \hat{y}'$

3. If  $P_s(\hat{y}') \leq P_s(\hat{y}_t)$  then with probability

$$e^{-\Delta s} = e^{-(s(\hat{y}) - s(\hat{y}'))} = \frac{e^{s(\hat{y}')}}{e^{s(\hat{y})}} = \frac{P_s(\hat{y}')}{P_s(\hat{y})}$$

do  $\hat{y}_{t+1} = \hat{y}'$  and otherwise  $\hat{y}_{t+1} = \hat{y}_t$

## The Metropolis Markov Chain

We need to show that  $P_s$  is a stationary distribution of this process.

We must show that if we select  $\hat{y}_t$  from  $P_s$ , and then select  $\hat{y}_{t+1}$  using the transition probabilities, then the distribution on  $\hat{y}_{t+1}$  is also  $P_s$ .

## Stationarity Condition

$$\begin{aligned} P'(\hat{y}) &= \sum_{\hat{y}'} P_s(\hat{y}') P_{\text{Trans}}(\hat{y} \mid \hat{y}') \\ &= P_s(\hat{y}) + \text{flow-in} - \text{flow-out} \\ &= P_s(\hat{y}) + \sum_{\hat{y}' \in N(\hat{y})} P_s(\hat{y}') P_{\text{Trans}}(\hat{y} \mid \hat{y}') - P_s(\hat{y}) P_{\text{Trans}}(\hat{y}' \mid \hat{y}) \end{aligned}$$

## Detailed Balance

Detailed balance means that for each pair of neighboring assignments  $\hat{y}, \hat{y}'$  we have equal flows in both directions.

$$P_s(\hat{y}')P_{\text{Trans}}(\hat{y} \mid \hat{y}') = P_s(\hat{y})P_{\text{Trans}}(\hat{y}' \mid \hat{y})$$

Without loss generality assume  $P_s(\hat{y}') \geq P_s(\hat{y})$ .

Metropolis is defined by

$$P_{\text{Trans}}(\hat{y} \mid \hat{y}') = e^{-\Delta s} P_{\text{Trans}}(\hat{y}' \mid \hat{y}) = \frac{P_s(\hat{y})}{P_s(\hat{y}')} P_{\text{Trans}}(\hat{y}' \mid \hat{y})$$

## Gibbs Sampling

The Metropolis algorithm wastes time by rejecting proposed moves.

Gibbs sampling avoids this move rejection.

In Gibbs sampling we select a node  $i$  at random and change that node by drawing a new node value conditioned on the current values of the other nodes.

We let  $\hat{y} \setminus i$  be the assignment of labels given by  $\hat{y}$  except that no label is assigned to node  $i$ .

We let  $\hat{y}[N(i)]$  be the assignment that  $\hat{y}$  gives to the nodes (pixels) that are the neighbors of node  $i$  (connected to  $i$  by an edge.)

# Gibbs Sampling

Markov Blanket Property:

$$P_s(\hat{y}[i] \mid \hat{y} \setminus i) = P_s(\hat{y}[i] \mid \hat{y}[N(i)])$$

Gibbs Sampling, Repeat:

- Select  $i$  at random
- draw  $c$  from  $P_s(\hat{y}[i] \mid y \setminus i) = P_s(\hat{y}[i] \mid \hat{y}[N(i)])$
- $\hat{y}[i] = c$

This algorithm does not require knowledge of  $Z$ .

The stationary distribution is  $P_s$ .



## Pseudolikelihood

For any distribution  $Q$  on assignments of labels to nodes (segmentations), and any assignment  $\hat{y}$ , we define  $\tilde{Q}(\hat{y})$  as follows.

$$\tilde{Q}(\hat{y}) = \prod_i Q(\hat{y}[i] \mid \hat{y}/i) = \prod_i Q(\hat{y}[i] \mid \hat{y}[N(i)])$$

We then train a graphical model with pseudolikelihood loss.

$$\Phi^* = \operatorname{argmin}_{\Phi} E_{y \sim \text{Pop}} - \ln \tilde{P}_{\Phi}(y)$$

## Pseudolikelihood

$$\mathcal{L}_{\text{PL}} = -\ln \tilde{P}_s(y)$$

We note that by the Markov blanket property for Markov random fields we have

$$\tilde{P}_s(\hat{y}) = \prod_i P_s(\hat{y}[i] \mid \hat{y}[N(i)])$$

Since the loss is directly computed we can directly back-propagate on the loss.

## Pseudolikelihood Theorem

$$\operatorname{argmin}_Q E_{y \sim \text{Pop}} - \ln \tilde{Q}(y) = \text{Pop}$$

or equivalently

$$\min_Q E_{y \sim \text{Pop}} - \ln \tilde{Q}(y) = E_{y \sim \text{Pop}} - \ln \widetilde{\text{Pop}}(y)$$

## Proof I

We have

$$\min_Q E_{y \sim \text{Pop}} - \ln \tilde{Q}(y) \leq E_{y \sim \text{Pop}} - \ln \widetilde{\text{Pop}}(y)$$

So it suffices to show

$$\min_Q E_{y \sim \text{Pop}} - \ln \tilde{Q}(y) \geq E_{y \sim \text{Pop}} - \ln \widetilde{\text{Pop}}(y)$$

## Proof II

We will prove the case of two nodes.

$$\begin{aligned} & \min_Q E_{y \sim \text{Pop}} - \ln Q(y[1]|y[2]) Q(y[2]|y[1]) \\ & \geq \min_{P_1, P_2} E_{y \sim \text{Pop}} - \ln P_1(y[1]|y[2]) P_2(y[2]|y[1]) \\ & = \min_{P_1} E_{y \sim \text{Pop}} - \ln P_1(y[1]|y[2]) + \min_{P_2} E_{y \sim \text{Pop}} - \ln P_2(y[2]|y[1]) \\ & = E_{y \sim \text{Pop}} - \ln \text{Pop}(y[1]|y[2]) + E_{y \sim \text{Pop}} - \ln \text{Pop}(y[2]|y[1]) \\ & = E_{y \sim \text{Pop}} - \ln \widetilde{\text{Pop}}(y) \end{aligned}$$

## Contrastive Divergence (CDk)

In contrastive divergence we first construct an MCMC process whose stationary distribution is  $P_s$ . This could be Metropolis or Gibbs or something else.

**Algorithm CDk:** Given a gold segmentation  $y$ , start the MCMC process from initial state  $y$  and run the process for  $k$  steps to get  $\hat{y}$ . Then take the loss to be

$$\mathcal{L}_{\text{CD}} = s(\hat{y}) - s(y)$$

If  $P_s = \text{Pop}$  then the the distribution on  $\hat{y}$  is the same as the distribution on  $y$  and the expected loss gradient is zero.

## Gibbs CD1

CD1 for the Gibbs MCMC process is a particularly interesting special case.

**Algorithm (Gibbs CD1):** Given  $y$ , select a node  $i$  at random and draw  $c \sim P(y[i] \mid y[N(i)])$ . Define  $y[i = c]$  to be the assignment (segmentation) which is the same as  $y$  except that node  $i$  is assigned label  $c$ . Take the loss to be

$$\mathcal{L}_{\text{CD}} = s(y[i = c]) - s(y)$$

## Gibbs CD1 Theorem

Gibbs CD1 is equivalent in expectation to pseudolikelihood.

$$\begin{aligned}\mathcal{L}_{\text{PL}} &= E_{y \sim \text{Pop}} \sum_i -\ln P_s(y[i] = c \mid y \setminus i) \\ &= E_{y \sim \text{Pop}} \sum_i -\ln \frac{e^{s(y)}}{Z_i} \quad Z_i = \sum_{c'} e^{s(y[i=c'])} \\ &= E_{y \sim \text{Pop}} \sum_i (\ln Z_i - s(y)) \\ \nabla_{\Phi} \mathcal{L}_{\text{PL}} &= E_{y \sim \text{Pop}} \sum_i \left( \frac{1}{Z_i} \sum_{c'} e^{s(y[i=c'])} \nabla_{\Phi} s(y[i=c']) \right) - \nabla_{\Phi} s(y) \\ &= E_{y \sim \text{Pop}} \sum_i \left( \sum_{c'} P(y[i=c'] \mid y \setminus i) \nabla_{\Phi} s(y[i=c']) \right) - \nabla_{\Phi} s(y)\end{aligned}$$



## Gibbs CD1 Theorem

$$\begin{aligned}
\nabla_{\Phi} \mathcal{L}_{\text{PL}} &= E_{y \sim \text{Pop}} \sum_i \left( \sum_{c'} P(y[i = c' \mid y \setminus i]) \nabla_{\Phi} s(y[i] = c') \right) - \nabla_{\Phi} s(y) \\
&= E_{y \sim \text{Pop}} \sum_i \left( E_{c' \sim P(y[i=c' \mid y \setminus i])} \nabla_{\Phi} s(y[i] = c') \right) - \nabla_{\Phi} s(y) \\
&\propto E_{y \sim \text{Pop}} E_i E_{c' \sim P(y[i=c' \mid y \setminus i])} \left( \nabla_{\Phi} s(y[i] = c') - \nabla_{\Phi} s(y) \right) \quad \text{Gibbs CD(1)}
\end{aligned}$$

## Loopy Belief Propagation (Loopy BP)

We design an algorithm that is correct for tree graphs and use it on non-tree (loopy) graphs.

## Belief Propagation on Trees

Belief Propagation is a message passing procedure (actually dynamic programming).

For each edge  $\{i, j\}$  and possible value  $\tilde{y}$  for node  $i$  we define  $Z_{j \rightarrow i}[c]$  to be the partition function for the subtree attached to  $i$  through  $j$  and with  $\hat{y}[i]$  restricted to  $c$ .

The function  $Z_{j \rightarrow i}$  on the possible values of node  $i$  is called the **message** from  $j$  to  $i$ .

The reverse direction message  $Z_{i \rightarrow j}$  is defined similarly.

## Dynamic Programming Computes the Messages

$$Z_{j \rightarrow i}[c] = \sum_{c'} e^{s_n[j, c'] + s_e[j, i, c', c]} \left( \prod_{k \in N(j), k \neq i} Z_{k \rightarrow j}[c'] \right)$$

## Loopy BP

In a Loopy Graph we can initialize all messages  $Z_{i \rightarrow j}[c] = 1$  and then repeating (until convergence) the updates

$$\tilde{Z}_{j \rightarrow i}[c] = \frac{1}{Z_{j \rightarrow i}} Z_{j \rightarrow i}[c] \quad Z_{j \rightarrow i} = \sum_c Z_{j \rightarrow i}[c]$$

$$Z_{j \rightarrow i}[c] = \sum_{c'} e^{s_n[j, c'] + s_e[j, i, c', c]} \left( \prod_{k \in N(j), k \neq i} \tilde{Z}_{k \rightarrow j}[c'] \right)$$

## Computing Node Marginals from Messages

$$\begin{aligned} Z_i(c) &\doteq \sum_{\hat{y}: \hat{y}[i]=c} e^{s(\hat{y})} \\ &= e^{s_i[c]} \left( \prod_{j \in N(i)} Z_{j \rightarrow i}[c] \right) \\ \textcolor{red}{P_i(c)} &= Z_i(c)/Z, \quad Z = \sum_c Z_i(c) \end{aligned}$$

## Computing Edge Marginals from Messages

$$\begin{aligned} Z_{i,j}(c, c') &\doteq \sum_{\hat{y}: \hat{y}[i]=c, \hat{y}[j]=c'} e^{s(\hat{y})} \\ &= e^{s_n[i,c]+s_n[j,c'] + s_e[i,j,c,c']} \\ &\quad \prod_{k \in N(i), k \neq j} Z_{k \rightarrow i}[c] \\ &\quad \prod_{k \in N(j), k \neq i} Z_{k \rightarrow j}[c'] \\ \textcolor{red}{P}_{i,j}(c, c') &= Z_{i,j}(c, c') / Z \quad Z = \sum_{c, c'} Z_{i,j}(c, c') \end{aligned}$$

## Summary

We are often interested in probability distributions on structured objects such as sentence or images.

Graphical models define softmax distributions on structured values.

It is infeasible to enumerate all sentences or all images.

However, pseudolikelihood provides a reasonable training algorithm and loopy BP can be used for both training time and test time inference.



**END**