

# **TTIC 31230, Fundamentals of Deep Learning**

David McAllester, Winter 2019

## **Deep Graphical Models**

# 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



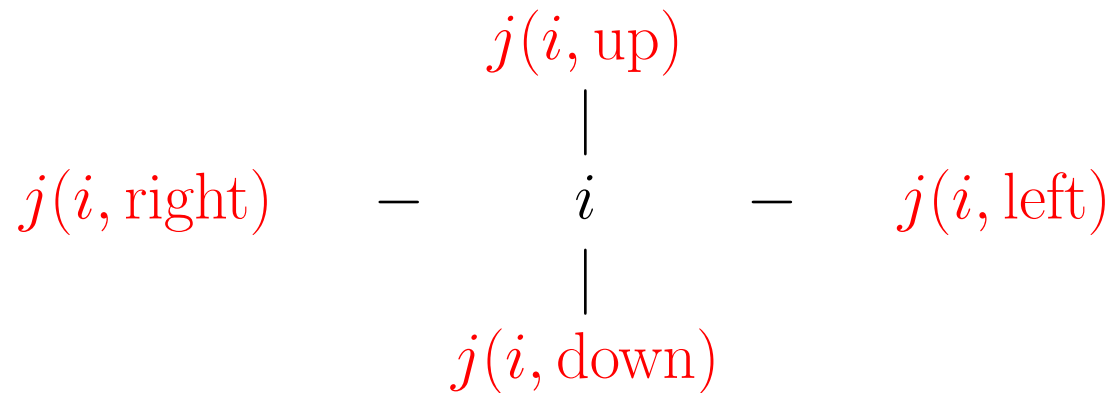
SLIC superpixels, Achanta et al.

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.

$$s(\hat{y}) = \sum_{i \in \text{Nodes}} s_n[i, \hat{y}[i]] + \sum_{i \in \text{Nodes}, d \in \{L, R, U, D\}} s_e[i, d, \hat{y}[i], \hat{y}[j(i, d)]]$$

Node Scores                      Edge Scores

## Computing the Node and Edge Scores

We assume a CNN computing node and edge score tensors

input image  
:  
CNN  
:  
for  $i, c$   $s_n[i, c] = \dots$   
for  $i, d, c, c'$   $s_e[i, d, c, c'] = \dots$

The tensor  $s_n[i, c]$  holds  $PC$  scores.

The tensor  $s_e[i, d, c, c']$  holds  $4PC^2$  scores.

# Computationally Infeasible Exponential Softmax

input image

⋮

CNN

⋮

for  $i, c$   $s_n[i, c] = \dots$

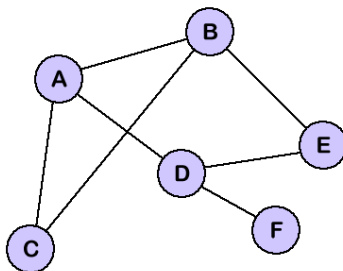
for  $i, d, c, c'$   $s_e[i, d, c, c'] = \dots$

for  $\hat{y}$   $s(\hat{y}) = \sum_i s_n[i, \hat{y}[i]] + \sum_{i,d} s_e[i, d, \hat{y}[i], \hat{y}[j[i, d)]]$

for  $\hat{y}$   $P(\hat{y}) = \text{expsoftmax}_{\hat{y}} s(\hat{y})$  all possible  $\hat{y}$

$\mathcal{L} = -\ln P(y)$  gold 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.

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

Computing  $Z$  in general is #P hard.



## Back-Propagation Through Intractable Softmax

```

                                ⋮
for  $i, c$        $s_n[i, c] = \dots$ 
for  $i, d, c, c'$   $s_e[i, d, c, c'] = \dots$ 

for  $\hat{y}$   $s(\hat{y}) = \sum_i s_n[i, \hat{y}[i]] + \sum_{i,d} s_e[i, d, \hat{y}[i], \hat{y}[j[i, d]]]$ 
for  $\hat{y}$   $P(\hat{y}) = \text{expsoftmax}_{\hat{y}} s(\hat{y})$  all possible  $\hat{y}$ 
 $\mathcal{L} = -\ln P(y)$  gold label  $y$ 

```

We need to compute  $s_n.\text{grad}[i, c]$  and  $s_e.\text{grad}[i, d, c, c']$ .

Although exact calculation is intractable, in practice the gradients can be approximated.

## Model Marginals Theorem

Theorem:

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

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

To approximately back-propagate log loss of an intractable graphical model it suffices to approximate **the model marginals** in red above.

## Proof of Model Marginals Theorem

We consider the case of node marginals.

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

# Methods of Approximating Model Marginals

MCMC Sampling

Pseudolikelihood

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  $P_s(\hat{y})$ .

We will design a Markov process whose states are segmentations  $\hat{y}$  and whose stationary distribution is  $P_s$ .

We will run the process past its mixing time to get a sample  $\hat{y}$  from  $P_s$ .

## A Neighbor Relation on States

We will say that segmentations (states)  $\hat{y}$  and  $\hat{y}'$  are neighbors if they differ in exactly one pixel.

Note that the number of neighbors of  $\hat{y}$  is  $P(C - 1)$ .

We will write  $N(\hat{y})$  for the set of neighbors of  $\hat{y}$ .

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

# The Metropolis Markov Chain

We need to define state transition probabilities.

In the **Metropolis algorithm** we do the following.

Pick an initial state  $\hat{y}$  and then repeat:

1. Pick a neighbor  $\hat{y}' \in N(\hat{y})$  uniformly at random.
2. If  $P_s(\hat{y}') \geq P_s(\hat{y})$  then do  $\hat{y} = \hat{y}'$
3. If  $P_s(\hat{y}') < P_s(\hat{y})$  then do  $\hat{y} = \hat{y}'$  with probability  $\frac{P_s(\hat{y}')}{P_s(\hat{y})}$

## The Metropolis Markov Chain

Pick an initial state  $\hat{y}$  and then repeat:

1. Pick a neighbor  $\hat{y}' \in N(\hat{y})$  uniformly at random.
2. If  $P_s(\hat{y}') \geq P_s(\hat{y})$  then do  $\hat{y} = \hat{y}'$
3. If  $P_s(\hat{y}') < P_s(\hat{y})$  then do  $\hat{y} = \hat{y}'$  with probability  $\frac{P_s(\hat{y}')}{P_s(\hat{y})}$

Note that we can determine which probability is larger just by comparing scores — we do not need to know  $Z$ .

Also note that the ratio  $P_s(\hat{y}')/P_s(\hat{y})$  can be computed from the scores without knowledge of  $Z$ .



## 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  $P_s$ .

# The Metropolis Markov Chain

$$\begin{aligned} P'(\hat{y}) &= \sum_{\hat{y}'} P_s(\hat{y}') P_{\text{Trans}}(\hat{y} \mid \hat{y}') \\ &= P_s(\hat{y}) P_{\text{Trans}}(\hat{y} \mid \hat{y}) + \sum_{\hat{y}' \in N(\hat{y})} P_s(\hat{y}') P_{\text{Trans}}(\hat{y} \mid \hat{y}') \\ &= P_s(\hat{y}) \left( 1 - \sum_{\hat{y}' \in N(\hat{y})} P_{\text{Trans}}(\hat{y}' \mid \hat{y}) \right) + \sum_{\hat{y}' \in N(\hat{y})} P_s(\hat{y}') P_{\text{Trans}}(\hat{y} \mid \hat{y}') \\ &= P_s(\hat{y}) + \sum_{\hat{y}' \in N(\hat{y})} P_s(\hat{y}') P_{\text{Trans}}(\hat{y} \mid \hat{y}') - P_s(\hat{y}) P_{\text{Trans}}(\hat{y}' \mid \hat{y}) \end{aligned}$$

## Detailed Balance

$$P'(\hat{y}) = P_s(\hat{y}) + \sum_{\hat{y}' \in N(\hat{y})} P_s(\hat{y}')P_{\text{Trans}}(\hat{y} \mid \hat{y}') - P_s(\hat{y})P_{\text{Trans}}(\hat{y}' \mid \hat{y})$$

$P_s(\hat{y}')P_{\text{Trans}}(\hat{y} \mid \hat{y}')$  is an amount of probability mass that moves from  $\hat{y}'$  to  $\hat{y}$ .

$P_s(\hat{y})P_{\text{Trans}}(\hat{y}' \mid \hat{y})$  is an amount of probability mass that moves from  $\hat{y}$  to  $\hat{y}'$ .

If these two flows are equal we call it **detailed balance**.

Detailed balance implies  $P'(\hat{y}) = P_s(\hat{y})$  and hence  $P_s$  is the stationary distribution.

## Detailed Balance

For  $P_s(\hat{y}) \geq P_s(\hat{y}')$  we have

$$\begin{aligned} P_s(\hat{y}) P_{\text{Trans}}(\hat{y}' \mid \hat{y}) &= P_s(\hat{y}) \frac{1}{|N(\hat{y})|} \frac{P_s(\hat{y}')}{P_s(\hat{y})} \\ &= \frac{1}{N} P_s(\hat{y}') \\ &= P_s(\hat{y}') P_{\text{Trans}}(\hat{y} \mid \hat{y}') \end{aligned}$$

## 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  $P_s(\hat{y}[i] \mid \hat{y}[N(i)])$
- $\hat{y}[i] = c$

This algorithm does not require knowledge of  $Z$ .

The stationary distribution is  $P_s$ .

## Pseudolikelihood

Pseudolikelihood is related to Gibbs Sampling. Rather than sample, we replace the objective  $-\ln P_s(\hat{y})$  with the objective  $-\ln \tilde{P}_s(\hat{y})$ .

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

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

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

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



## 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.

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

## Contrastive Divergence

Contrastive Divergence is closely related to pseudolikelihood.

**Algorithm (Gibbs CD1):** Select a node  $i$  at random. 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  $c$ .

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

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

## Contrastive Divergence Theorem

Gibbs CD1 is equivalent in expectation to pseudolikelihood.

$$\mathcal{L}_{\text{PL}} = \sum_i -\ln \frac{e^{s(y)}}{\sum_c e^{s(y[i=c])}}$$

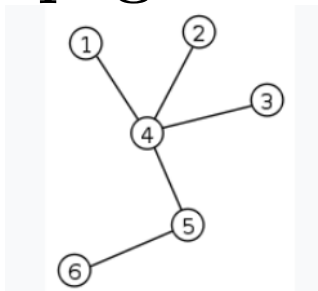
$$= \sum_i \left( \ln \left( \sum_c e^{s(y[i=c])} \right) - s(y) \right)$$

$$\begin{aligned} \nabla_{\Phi} \mathcal{L}_{\text{PL}} &= \sum_i \left( E_{c|i} \nabla_{\Phi} s(y[i=c]) - \nabla_{\Phi} s(y) \right) \\ &= N E_{i,c} \nabla_{\Phi} \mathcal{L}_{\text{CD}} \end{aligned}$$

## 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



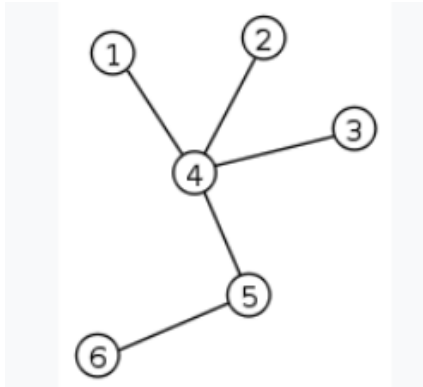
Belief Propagation is a message passing procedure (actually dynamic programming).

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

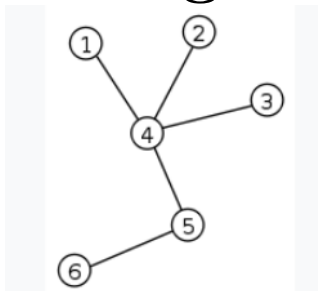
In a Loopy Graph we can initialize all message  $Z_{i \rightarrow j}[c] = 1$  and then repeating (until convergence) the updates

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

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



## 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**