Approach 000 0000000 000000 Experiments 0000 000 Discussion 000 0000000000

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへで

# Learning (Convex Inference of Marginals)

Justin Domke

University of Maryland

Approach 000 0000000 000000 Experiments

Discussion 000 0000000000

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

# Outline

#### Introduction

Motivation Overview of the approach

### Approach

Inference Learning (Loss functions) Learning (Derivatives of beliefs)

## Experiments

Introduction Results

## Discussion

Summary

Approach 000 0000000 000000 Experiments

Discussion 000 0000000000

# Outline

#### Introduction

### Motivation

Overview of the approach

### Approach

Inference Learning (Loss functions) Learning (Derivatives of beliefs)

## Experiments

Introduction Results

## Discussion

Summary



<ロ> (四) (四) (三) (三) (三)

• We have only  $\{(\hat{\mathbf{x}}, \hat{\mathbf{y}})\} \sim p(\mathbf{x}, \mathbf{y})$ .



= nac

• We have only  $\{(\hat{\mathbf{x}}, \hat{\mathbf{y}})\} \sim p(\mathbf{x}, \mathbf{y}).$ 



Approach 000 0000000 000000 Experiments

Discussion 000 000000000

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

# Setting

• True (unknown) distribution  $p(\mathbf{x}, \mathbf{y})$ 

Maximum Likelihood Approach:

- 1. Fit a graphical model  $q(\mathbf{x}|\mathbf{y})$  by max (conditional) likelihood. (Learning)
- 2. Given y, compute  $q(x_i|y)$ . (Inference)

Justification: Given a correct model,  $q(\mathbf{x}|\mathbf{y}) \rightarrow p(\mathbf{x}|\mathbf{y})$ .



Approach 000 0000000 000000 Experiments 0000 000 Discussion 000 000000000

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

# Setting

• True (unknown) distribution  $p(\mathbf{x}, \mathbf{y})$ 

Maximum Likelihood Approach:

- 1. Fit a graphical model  $q(\mathbf{x}|\mathbf{y})$  by max (conditional) likelihood. (Learning)
- 2. Given y, compute  $q(x_i|y)$ . (Inference)

Justification: Given a correct model,  $q(\mathbf{x}|\mathbf{y}) \rightarrow p(\mathbf{x}|\mathbf{y})$ .

Approach 000 0000000 000000 Experiments 0000 000 Discussion 000 0000000000

◆ロト ◆得ト ◆ヨト ◆ヨト ヨー のくべ

# Problems with Max Likelihood

#### • In many cases, this works well.

• Problems arise, particularly when the model has high treewidth.

- Computational Intractability
- Model Defects

Approach 000 0000000 000000 Experiments 0000 000 Discussion 000 0000000000

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ のQ@

- In many cases, this works well.
- Problems arise, particularly when the model has high treewidth.
  - Computational Intractability
  - Model Defects

Introduction 0000●0 00 Approach 000 0000000 000000 Experiments 0000 000 Discussion 000 0000000000

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

- 1. Computational Intractability.
  - Often, max likelihood can't be done.
  - Even if it could, would the results be what we want under approximate inference?

Approach 000 0000000 000000 Experiments 0000 000 Discussion 000 0000000000

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

- 1. Computational Intractability.
  - Often, max likelihood can't be done.
  - Even if it could, would the results be what we want under <u>approximate</u> inference?

Introduction 00000● 00 Approach 000 0000000 000000 Experiments 0000 000 Discussion 000 0000000000

・ロト ・ 日 ・ ・ ヨ ・ ・ 日 ・ うらつ

- 2. Model Defects
  - Possible to set  $\theta$  such that  $q(\mathbf{x}|\mathbf{y}; \theta) = p(\mathbf{x}|\mathbf{y})$ ?
  - Max conditional likelihood  $\approx E[KL-divergence]$  $\arg\min_{\theta} \sum_{\mathbf{y}} p(\mathbf{y}) \sum_{\mathbf{x}} p(\mathbf{x}|\mathbf{y}) \log \frac{p(\mathbf{x}|\mathbf{y})}{q(\mathbf{x}|\mathbf{y};\theta)}$   $= \arg\min_{\theta} - \sum_{\mathbf{y}} p(\mathbf{y}) \sum_{\mathbf{x}} p(\mathbf{x}|\mathbf{y}) \log q(\mathbf{x}|\mathbf{y};\theta)$   $\approx \arg\max_{\theta} \sum_{\{(\hat{\mathbf{x}}, \hat{\mathbf{y}})\}} \log q(\hat{\mathbf{x}}|\hat{\mathbf{y}};\theta)$
  - min KL-divergence  $\neq$  best marginals.
    - (Even assuming exact inference.)

Approach 000 0000000 000000 Experiments 0000 000 Discussion 000 0000000000

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ のQ@

- 2. Model Defects
  - Possible to set  $\theta$  such that  $q(\mathbf{x}|\mathbf{y}; \theta) = p(\mathbf{x}|\mathbf{y})$ ?
  - Max conditional likelihood  $\approx E[KL-divergence]$  $\arg\min_{\theta} \sum_{\mathbf{y}} p(\mathbf{y}) \sum_{\mathbf{x}} p(\mathbf{x}|\mathbf{y}) \log \frac{p(\mathbf{x}|\mathbf{y})}{q(\mathbf{x}|\mathbf{y};\theta)}$   $= \arg\min_{\theta} - \sum_{\mathbf{y}} p(\mathbf{y}) \sum_{\mathbf{x}} p(\mathbf{x}|\mathbf{y}) \log q(\mathbf{x}|\mathbf{y};\theta)$   $\approx \arg\max_{\theta} \sum_{\{(\hat{\mathbf{x}}, \hat{\mathbf{y}})\}} \log q(\hat{\mathbf{x}}|\hat{\mathbf{y}};\theta)$
  - min KL-divergence  $\neq$  best marginals.
    - (Even assuming exact inference.)

Approach 000 0000000 000000 Experiments 0000 000 Discussion 000 0000000000

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ のQ@

- 2. Model Defects
  - Possible to set  $\theta$  such that  $q(\mathbf{x}|\mathbf{y}; \theta) = p(\mathbf{x}|\mathbf{y})$ ?
  - Max conditional likelihood  $\approx E[KL-divergence]$  $\arg\min_{\theta} \sum_{\mathbf{y}} p(\mathbf{y}) \sum_{\mathbf{x}} p(\mathbf{x}|\mathbf{y}) \log \frac{p(\mathbf{x}|\mathbf{y})}{q(\mathbf{x}|\mathbf{y};\theta)}$   $= \arg\min_{\theta} - \sum_{\mathbf{y}} p(\mathbf{y}) \sum_{\mathbf{x}} p(\mathbf{x}|\mathbf{y}) \log q(\mathbf{x}|\mathbf{y};\theta)$   $\approx \arg\max_{\theta} \sum_{\{(\hat{\mathbf{x}}, \hat{\mathbf{y}})\}} \log q(\hat{\mathbf{x}}|\hat{\mathbf{y}};\theta)$
  - min KL-divergence  $\neq$  best marginals.
    - (Even assuming exact inference.)

Approach 000 0000000 000000 Experiments

Discussion 000 0000000000

# Outline

#### Introduction

Motivation

#### Overview of the approach

#### Approach

Inference Learning (Loss functions) Learning (Derivatives of beliefs)

## Experiments

Introduction Results

## Discussion

Summary



Approach 000 0000000 000000 Experiments 0000 000

Discussion 000 0000000000

- Goal: Sidestep problems of intractability and model defects.
- In approximate inference,  $q(\mathbf{x}|\mathbf{y})$  is often used to create a free energy function.

$$F(\mathbf{y}, \{b_r(\mathbf{x}_r)\}) \quad \{b_r^*(\mathbf{x}_r)\} = \arg\min_{\{b_r\}} F(\mathbf{y}, \{b_r(\mathbf{x}_r)\})$$

- Idea: Think of F as mapping from y to {b<sup>\*</sup><sub>r</sub>(x<sub>r</sub>)}. Directly fit F to make the mapping as accurate as possible.
- Computational Tractability: Restrict F to be convex.
- Model Defects: Learn by minimizing <u>empirical risk</u>, where risk measures the accuracy of marginals.





Discussion 000 0000000000

# Overview

- Goal: Sidestep problems of intractability and model defects.
- In approximate inference,  $q(\mathbf{x}|\mathbf{y})$  is often used to create a free energy function.

 $F(\mathbf{y}, \{b_r(\mathbf{x}_r)\}) \quad \{b_r^*(\mathbf{x}_r)\} = \arg\min_{\{b_r\}} F(\mathbf{y}, \{b_r(\mathbf{x}_r)\})$ 

- Idea: Think of F as mapping from y to {b<sup>\*</sup><sub>r</sub>(x<sub>r</sub>)}. Directly fit F to make the mapping as accurate as possible.
- Computational Tractability: Restrict F to be convex.
- Model Defects: Learn by minimizing <u>empirical risk</u>, where risk measures the accuracy of marginals.





Discussion 000 0000000000

- Goal: Sidestep problems of intractability and model defects.
- In approximate inference,  $q(\mathbf{x}|\mathbf{y})$  is often used to create a free energy function.

$$F(\mathbf{y}, \{b_r(\mathbf{x}_r)\}) \quad \{b_r^*(\mathbf{x}_r)\} = \arg\min_{\{b_r\}} F(\mathbf{y}, \{b_r(\mathbf{x}_r)\})$$

- Idea: Think of F as mapping from y to  $\{b_r^*(\mathbf{x}_r)\}$ . Directly fit F to make the mapping as accurate as possible.
- Computational Tractability: Restrict F to be convex.
- Model Defects: Learn by minimizing <u>empirical risk</u>, where risk measures the accuracy of marginals.





Discussion 000 0000000000

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

- Goal: Sidestep problems of intractability and model defects.
- In approximate inference,  $q(\mathbf{x}|\mathbf{y})$  is often used to create a free energy function.

$$F(\mathbf{y}, \{b_r(\mathbf{x}_r)\}) \quad \{b_r^*(\mathbf{x}_r)\} = \arg\min_{\{b_r\}} F(\mathbf{y}, \{b_r(\mathbf{x}_r)\})$$

- Idea: Think of F as <u>mapping</u> from y to {b<sup>\*</sup><sub>r</sub>(x<sub>r</sub>)}. Directly fit F to make the mapping as accurate as possible.
- Computational Tractability: Restrict F to be convex.
- Model Defects: Learn by minimizing <u>empirical risk</u>, where risk measures the accuracy of marginals.





Discussion 000 0000000000

- Goal: Sidestep problems of intractability and model defects.
- In approximate inference,  $q(\mathbf{x}|\mathbf{y})$  is often used to create a free energy function.

$$F(\mathbf{y}, \{b_r(\mathbf{x}_r)\}) \quad \{b_r^*(\mathbf{x}_r)\} = \arg\min_{\{b_r\}} F(\mathbf{y}, \{b_r(\mathbf{x}_r)\})$$

- Idea: Think of F as <u>mapping</u> from y to {b<sup>\*</sup><sub>r</sub>(x<sub>r</sub>)}. Directly fit F to make the mapping as accurate as possible.
- Computational Tractability: Restrict F to be convex.
- Model Defects: Learn by minimizing <u>empirical risk</u>, where risk measures the accuracy of marginals.

Approach •00 •0000000 •000000 Experiments

(日) (四) (三) (三) (三)

Discussion 000 0000000000

# Outline

#### Introduction

Motivation Overview of the approach

### Approach

## Inference

Learning (Loss functions) Learning (Derivatives of beliefs)

## Experiments

Introduction Results

## Discussion

Summary



Experiments

Discussion 000 0000000000

Inference

$$F = \sum_{f \in \mathscr{F}} \sum_{r \in \mathscr{R}} \sum_{\mathbf{x}_r} w_f(\mathbf{x}_r, \mathbf{y}_r) f(b_r(\mathbf{x}_r))$$

• Typically,  $\mathscr{F} = \{b, b \log b\}$ . Anything convex over (0,1) is OK.

$$\{b_r^*(\mathbf{x}_r)\} = \arg\min_{\{b_r\}} F(\mathbf{y}, \{b_r(\mathbf{x}_r)\})$$

• Minimization is over some relaxation of the marginal polytope

ocal consistency: 
$$\sum_{\mathbf{x}_{c\setminus i}} b_c(\mathbf{x}_c) = b_i(x_i)$$
$$\sum_{\mathbf{x}_c} b_c(\mathbf{x}_c) = 1 \qquad \sum_{x_i} b_i(x_i) = 1$$
$$b_c(\mathbf{x}_c) \ge \mathbf{0} \qquad b_i(x_i) \ge 0$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ = 三 のへで



Experiments

Discussion 000 0000000000

Inference

$$F = \sum_{f \in \mathscr{F}} \sum_{r \in \mathscr{R}} \sum_{\mathbf{x}_r} w_f(\mathbf{x}_r, \mathbf{y}_r) f(b_r(\mathbf{x}_r))$$

• Typically,  $\mathscr{F} = \{b, b \log b\}$ . Anything convex over (0,1) is OK.

$$\{b_r^*(\mathbf{x}_r)\} = \arg\min_{\{b_r\}} F(\mathbf{y}, \{b_r(\mathbf{x}_r)\})$$

• Minimization is over some relaxation of the marginal polytope

ocal consistency: 
$$\sum_{\mathbf{x}_{c\setminus i}} b_c(\mathbf{x}_c) = b_i(x_i)$$
$$\sum_{\mathbf{x}_c} b_c(\mathbf{x}_c) = 1 \qquad \sum_{x_i} b_i(x_i) = 1$$
$$b_c(\mathbf{x}_c) \ge \mathbf{0} \qquad b_i(x_i) \ge 0$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへぐ



Experiments

Inference

$$F = \sum_{f \in \mathscr{F}} \sum_{r \in \mathscr{R}} \sum_{\mathbf{x}_r} w_f(\mathbf{x}_r, \mathbf{y}_r) f(b_r(\mathbf{x}_r))$$

• Typically,  $\mathscr{F} = \{b, b \log b\}$ . Anything convex over (0,1) is OK.

$$\{b_r^*(\mathbf{x}_r)\} = \arg\min_{\{b_r\}} F(\mathbf{y}, \{b_r(\mathbf{x}_r)\})$$

• Minimization is over some relaxation of the marginal polytope

$$\begin{array}{ll} \text{ocal consistency:} & \sum_{\mathbf{x}_c \setminus i} b_c(\mathbf{x}_c) = b_i(x_i) \\ & \sum_{\mathbf{x}_c} b_c(\mathbf{x}_c) = 1 & \sum_{x_i} b_i(x_i) = 1 \\ & b_c(\mathbf{x}_c) \ge \mathbf{0} & b_i(x_i) \ge 0 \end{array}$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへぐ



Experiments

Discussion 000 0000000000

Inference

$$F = \sum_{f \in \mathscr{F}} \sum_{r \in \mathscr{R}} \sum_{\mathbf{x}_r} w_f(\mathbf{x}_r, \mathbf{y}_r) f(b_r(\mathbf{x}_r))$$

• Typically,  $\mathscr{F} = \{b, b \log b\}$ . Anything convex over (0,1) is OK.

$$\{b_r^*(\mathbf{x}_r)\} = \arg\min_{\{b_r\}} F(\mathbf{y}, \{b_r(\mathbf{x}_r)\})$$

• Minimization is over some relaxation of the marginal polytope

ocal consistency: 
$$\sum_{\mathbf{x}_{c\setminus i}} b_c(\mathbf{x}_c) = b_i(x_i)$$
$$\sum_{\mathbf{x}_c} b_c(\mathbf{x}_c) = 1 \qquad \sum_{x_i} b_i(x_i) = 1$$
$$b_c(\mathbf{x}_c) \ge \mathbf{0} \qquad b_i(x_i) \ge 0$$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへで

Approach

Experiments

Discussion 000 0000000000

Inference

$$\{b_r^*(\mathbf{x}_r)\} = \arg\min_{\{b_r\}} \sum_{f \in \mathscr{F}} \sum_{r \in \mathscr{R}} \sum_{\mathbf{x}_r} w_f(\mathbf{x}_r, \mathbf{y}_r) f(b_r(\mathbf{x}_r))$$

such that (local consistency)

Equivalent, more convenient formulation:

$$\begin{split} \mathbf{b}^* &= \arg\min_{\mathbf{b}} \sum_{f \in \mathscr{F}} \mathbf{w}_f(\mathbf{y})^T f(\mathbf{b}) & \mathbf{b} \iff \{b_r(\mathbf{x}_r)\}, \forall r, \mathbf{x}_r \\ & \mathbf{w}_f(\mathbf{y}) \iff w_f(\mathbf{x}_r, \mathbf{y}_r), \forall r, \mathbf{x}_r \\ & \mathbf{b} \ge \mathbf{0}. \end{split}$$



Experiments 0000 000 Discussion 000 0000000000

Inference

$$\{b_r^*(\mathbf{x}_r)\} = \arg\min_{\{b_r\}} \sum_{f \in \mathscr{F}} \sum_{r \in \mathscr{R}} \sum_{\mathbf{x}_r} w_f(\mathbf{x}_r, \mathbf{y}_r) f(b_r(\mathbf{x}_r))$$

such that (local consistency)

Equivalent, more convenient formulation:

$$\mathbf{b}^* = \arg\min_{\mathbf{b}} \sum_{f \in \mathscr{F}} \mathbf{w}_f(\mathbf{y})^T f(\mathbf{b}) \qquad \mathbf{b} \quad \leftrightarrow \quad \{b_r(\mathbf{x}_r)\}, \forall r, \mathbf{x}_r \\ \text{such that} \qquad A\mathbf{b} = \mathbf{d} \\ \mathbf{b} \ge \mathbf{0}. \end{cases}$$



Experiments 0000 000 Discussion 000 0000000000

Inference

$$\{b_r^*(\mathbf{x}_r)\} = \arg\min_{\{b_r\}} \sum_{f \in \mathscr{F}} \sum_{r \in \mathscr{R}} \sum_{\mathbf{x}_r} w_f(\mathbf{x}_r, \mathbf{y}_r) f(b_r(\mathbf{x}_r))$$

such that (local consistency)

Equivalent, more convenient formulation:

$$\begin{split} \mathbf{b}^* &= \arg\min_{\mathbf{b}} \sum_{f \in \mathscr{F}} \mathbf{w}_f(\mathbf{y})^T f(\mathbf{b}) & \mathbf{b} \iff \{b_r(\mathbf{x}_r)\}, \forall r, \mathbf{x}_r \\ & \mathbf{w}_f(\mathbf{y}) \iff w_f(\mathbf{x}_r, \mathbf{y}_r), \forall r, \mathbf{x}_r \\ & \mathbf{b} \ge \mathbf{0}. \end{split}$$

Approach 000 000000 000000 Experiments

Discussion 000 0000000000

# Outline

#### Introduction

Motivation Overview of the approach

## Approach

Inference Learning (Loss functions)

#### Experiments

Introduction Results

## Discussion

Summary



Learning

• Given  $\{(\hat{\mathbf{x}}, \hat{\mathbf{y}})\} \sim p(\mathbf{x}, \mathbf{y})$ , how to quantify the quality of predicted marginals?

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ のQ@

- Many possibilities. Two suggestions:
  - Log-loss
  - Quad-loss



Discussion 000 0000000000

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへで

# Log-loss

• Try to minimize "expected average univariate KL-divergence".

$$F^* = \arg\min_{F} \sum_{\mathbf{y}} p(\mathbf{y}) \sum_{i} \sum_{x_i} p(x_i | \mathbf{y}) \log \frac{p(x_i | \mathbf{y})}{b_i^*(x_i | \mathbf{y}, F)}.$$
  
=  $\arg\min_{F} - \sum_{\mathbf{y}} p(\mathbf{y}) \sum_{i} \sum_{x_i} p(x_i | \mathbf{y}) \log b_i^*(x_i | \mathbf{y}, F)$   
=  $\arg\min_{F} - \sum_{\mathbf{y}} \sum_{i} \sum_{x_i} p(x_i, \mathbf{y}) \log b_i^*(x_i | \mathbf{y}, F)$   
 $\approx \arg\min_{F} \sum_{\{(\hat{\mathbf{x}}, \hat{\mathbf{y}})\}} \underbrace{-\sum_{i} \log b_i^*(\hat{x}_i | \hat{\mathbf{y}}, F)}_{L_{\log}}$ 



Log-loss

• Try to minimize "expected average univariate KL-divergence".

$$L_{\log} = -\sum_{i} \log b_{i}^{*}(\hat{x}_{i}|\hat{\mathbf{y}}, F)$$
$$\frac{\partial L_{\log}}{\partial \theta_{j}} = -\sum_{i} \frac{\partial b_{i}^{*}(\hat{x}_{i}|\hat{\mathbf{y}}, F) / \partial \theta_{j}}{b_{i}^{*}(\hat{x}_{i}|\hat{\mathbf{y}}, F)}$$

(See also Kakade et al. "An Alternate Objective Function for Markovian Fields")

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ のQ@



Experiments

Discussion 000 0000000000

▲□▶ ▲録▶ ▲臣▶ ★臣▶ ―臣 …の�?

# Quad-loss

• Try to minimize "expected average univariate quadratic difference".

$$F^* = \arg\min_{F} \sum_{\mathbf{y}} p(\mathbf{y}) \sum_{i} \sum_{x_i} (p(x_i|\mathbf{y}) - b_i^*(x_i|\mathbf{y}, F))^2$$
  
= 
$$\arg\min_{F} \sum_{\mathbf{y}} \sum_{i} \sum_{x_i} (-2p(x_i, \mathbf{y})b_i^*(x_i|\mathbf{y}, F) + p(\mathbf{y})b_i^*(x_i|\mathbf{y}, F)^2)$$
  
$$\approx \arg\min_{F} \sum_{\{(\hat{\mathbf{x}}, \hat{\mathbf{y}})\}} \underbrace{\sum_{i} (-2b_i^*(\hat{x}_i|\hat{\mathbf{y}}, F) + \sum_{x_i} b_i^*(x_i|\hat{\mathbf{y}}, F)^2)}_{L_{guad}}$$

Approach 000 0000000 000000 Experiments

Discussion 000 000000000

# Quad-loss

• Try to minimize "expected average univariate quadratic difference".

$$L_{quad} = \sum_{i} \left( -2b_{i}^{*}(\hat{x}_{i}|\hat{y}, F) + \sum_{x_{i}} b_{i}^{*}(x_{i}|\hat{y}, F)^{2} \right)$$

$$\frac{\partial L_{quad}}{\partial \theta_j} = 2\sum_i \left( -\frac{\partial b_i^*(\hat{x}_i|\hat{\mathbf{y}},F)}{\partial \theta_j} + \sum_{x_i} b_i^*(x_i|\hat{\mathbf{y}},F) \frac{\partial b_i^*(x_i|\hat{\mathbf{y}},F)}{\partial \theta_j} \right)$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

Approach 000 000000 000000 Experiments 0000 000 Discussion 000 0000000000

Recap

$$\frac{\partial L_{\mathsf{log}}}{\partial \theta_j} = -\sum_i \frac{\partial b_i^*(\hat{x}_i | \hat{\mathbf{y}}, F) / \partial \theta_j}{b_i^*(\hat{x}_i | \hat{\mathbf{y}}, F)}$$

$$\frac{\partial L_{quad}}{\partial \theta_j} = 2\sum_i \left(-\frac{\partial b_i^*(\hat{x}_i|\hat{\mathbf{y}}, F)}{\partial \theta_j} + \sum_{x_i} b_i^*(x_i|\hat{\mathbf{y}}, F) \frac{\partial b_i^*(x_i|\hat{\mathbf{y}}, F)}{\partial \theta_j}\right)$$

If we could calculate  $\frac{\partial b_i^*(\hat{\mathbf{x}}_i|\hat{\mathbf{y}},F)}{\partial \theta_j}$ , we could optimize  $L_{\{\log, quad\}}$ . However, recall that  $b_i^*(\hat{\mathbf{x}}_i|\hat{\mathbf{y}},F)$  is <u>implicit</u>.

$$\mathbf{b}^* = \arg\min_{\mathbf{b}} \sum_{f \in \mathscr{F}} \mathbf{w}_f(\mathbf{y})^T f(\mathbf{b})$$
  
such that  $A\mathbf{b} = \mathbf{d}$   
 $\mathbf{b} \ge \mathbf{0}.$ 

▲□▶ ▲□▶ ▲□▶ ▲□▶ = 三 のへで

Approach ○○○ ○○○○○○● Experiments 0000 000 Discussion 000 0000000000

Recap

$$\frac{\partial L_{\log}}{\partial \theta_j} = -\sum_i \frac{\partial b_i^*(\hat{x}_i | \hat{\mathbf{y}}, F) / \partial \theta_j}{b_i^*(\hat{x}_i | \hat{\mathbf{y}}, F)}$$

$$\frac{\partial L_{\mathsf{quad}}}{\partial \theta_j} = 2\sum_i \big( -\frac{\partial b_i^*(\hat{x}_i|\hat{\mathbf{y}}, F)}{\partial \theta_j} + \sum_{x_i} b_i^*(x_i|\hat{\mathbf{y}}, F) \frac{\partial b_i^*(x_i|\hat{\mathbf{y}}, F)}{\partial \theta_j} \big)$$

If we could calculate  $\frac{\partial b_i^*(\hat{x}_i|\hat{y},F)}{\partial \theta_j}$ , we could optimize  $L_{\{\log, quad\}}$ . However, recall that  $b_i^*(\hat{x}_i|\hat{y},F)$  is <u>implicit</u>.

$$\mathbf{b}^* = \arg\min_{\mathbf{b}} \sum_{f \in \mathscr{F}} \mathbf{w}_f(\mathbf{y})^T f(\mathbf{b})$$
  
such that  $A\mathbf{b} = \mathbf{d}$   
 $\mathbf{b} \ge \mathbf{0}.$ 

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへで
Approach 000 000000 000000 Experiments 0000 000 Discussion 000 0000000000

Recap

$$\frac{\partial L_{\log}}{\partial \theta_j} = -\sum_i \frac{\partial b_i^*(\hat{x}_i | \hat{\mathbf{y}}, F) / \partial \theta_j}{b_i^*(\hat{x}_i | \hat{\mathbf{y}}, F)}$$

$$\frac{\partial L_{\mathsf{quad}}}{\partial \theta_j} = 2\sum_i \big( -\frac{\partial b_i^*(\hat{x}_i|\hat{\mathbf{y}}, F)}{\partial \theta_j} + \sum_{x_i} b_i^*(x_i|\hat{\mathbf{y}}, F) \frac{\partial b_i^*(x_i|\hat{\mathbf{y}}, F)}{\partial \theta_j} \big)$$

If we could calculate  $\frac{\partial b_i^*(\hat{x}_i|\hat{\mathbf{y}},F)}{\partial \theta_j}$ , we could optimize  $L_{\{\log,quad\}}$ . However, recall that  $b_i^*(\hat{x}_i|\hat{\mathbf{y}},F)$  is <u>implicit</u>.

$$egin{array}{lll} \mathbf{b}^* = rg\min_{\mathbf{b}} \sum_{f \in \mathscr{F}} \mathbf{w}_f(\mathbf{y})^{\mathcal{T}} f(\mathbf{b}) \ ext{such that} & A\mathbf{b} = \mathbf{d} \ extbf{b} \geq \mathbf{0}. \end{array}$$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへで

Approach 000 000000 000000 Experiments 0000 000

Recap

$$\frac{\partial L_{\log}}{\partial \theta_j} = -\sum_i \frac{\partial b_i^*(\hat{x}_i | \hat{\mathbf{y}}, F) / \partial \theta_j}{b_i^*(\hat{x}_i | \hat{\mathbf{y}}, F)}$$

$$\frac{\partial L_{quad}}{\partial \theta_j} = 2\sum_i \big( -\frac{\partial b_i^*(\hat{x}_i|\hat{\mathbf{y}}, F)}{\partial \theta_j} + \sum_{x_i} b_i^*(x_i|\hat{\mathbf{y}}, F) \frac{\partial b_i^*(x_i|\hat{\mathbf{y}}, F)}{\partial \theta_j} \big)$$

If we could calculate  $\frac{\partial b_i^*(\hat{x}_i|\hat{\mathbf{y}},F)}{\partial \theta_j}$ , we could optimize  $L_{\{\log, quad\}}$ . However, recall that  $b_i^*(\hat{x}_i|\hat{\mathbf{y}},F)$  is <u>implicit</u>.

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへで

Approach

Experiments

Discussion 000 0000000000

# Outline

#### Introduction

Motivation Overview of the approach

### Approach

Inference Learning (Loss functions) Learning (Derivatives of beliefs)

#### Experiments

Introduction Results

### Discussion

Approach

Experiments

Discussion 000 0000000000

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

## Claim 1

Claim 1: Let  $F(\mathbf{b}, \theta)$  be a continuous function such that for all  $\theta$ , F that has a unique stationary point in **b**.Define  $\mathbf{b}^*(\theta)$  such that  $\frac{\partial F(\mathbf{b}^*(\theta), \theta)}{\partial \mathbf{b}} = \mathbf{0}$ .Then,

$$\frac{\partial \mathbf{b}^*(\theta)}{\partial \theta_j} = -\left(\frac{\partial^2 F(\mathbf{b}^*(\theta), \theta)}{\partial \mathbf{b} \partial \mathbf{b}^T}\right)^{-1} \frac{\partial^2 F(\mathbf{b}^*(\theta), \theta)}{\partial \mathbf{b} \partial \theta_j}.$$

**Proof**: (Implicit Function Theorem.)

Approach

Experiments

Discussion 000 0000000000

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

# Claim 1

**Claim 1**: Let  $F(\mathbf{b}, \theta)$  be a continuous function such that for all  $\theta$ , F that has a unique stationary point in **b**.Define  $\mathbf{b}^*(\theta)$  such that  $\frac{\partial F(\mathbf{b}^*(\theta), \theta)}{\partial \mathbf{b}} = \mathbf{0}$ .Then,

$$\frac{\partial \mathbf{b}^*(\theta)}{\partial \theta_j} = -\left(\frac{\partial^2 F(\mathbf{b}^*(\theta), \theta)}{\partial \mathbf{b} \partial \mathbf{b}^T}\right)^{-1} \frac{\partial^2 F(\mathbf{b}^*(\theta), \theta)}{\partial \mathbf{b} \partial \theta_j}$$

• Not good enough, since F is minimized under constraints.

Approach

Experiments

Discussion 000 000000000

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ の��

### Claim 2

**Claim 2**: Define  $\mathbf{b}^*(\theta) \doteq \operatorname{argmin}_{\mathbf{b}} F(\mathbf{b}, \theta)$ , such that  $A\mathbf{b} = d$  for some convex function F. Then,

$$\frac{\partial \mathbf{b}^*(\theta)}{\partial \theta_j} = (D^{-1}A^T (AD^{-1}A^T)^{-1}AD^{-1} - D^{-1})\frac{\partial^2 F}{\partial \mathbf{b} \partial \theta_j}, \ D = (\frac{\partial^2 F}{\partial \mathbf{b} \partial \mathbf{b}^T}).$$

**Proof**: (Make a Lagrangian, apply claim 1, do algebra.)



Experiments 0000 000 Discussion 000 0000000000

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへで

# Claim 2

**Claim 2**: Define  $\mathbf{b}^*(\theta) \doteq \operatorname{argmin}_{\mathbf{b}} F(\mathbf{b}, \theta)$ , such that  $A\mathbf{b} = d$  for some convex function F. Then,

$$\frac{\partial \mathbf{b}^*(\theta)}{\partial \theta_j} = (D^{-1}A^T (AD^{-1}A^T)^{-1}AD^{-1} - D^{-1})\frac{\partial^2 F}{\partial \mathbf{b} \partial \theta_j}, \ D = (\frac{\partial^2 F}{\partial \mathbf{b} \partial \mathbf{b}^T}).$$

In our case,  

$$D = \operatorname{diag}(\sum_{f} \mathbf{w}_{f}(\mathbf{y}) \odot f''(\mathbf{b}))$$

$$\frac{\partial^{2} F}{\partial \mathbf{b} \partial \theta_{j}} = \sum_{f} \frac{\partial \mathbf{w}_{f}(\mathbf{y})}{\partial \theta_{j}} \odot f'$$

We still need  $\frac{\partial \mathbf{w}_f(\mathbf{y})}{\partial \theta_i}$ .



Experiments

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

### Claim 2

**Claim 2**: Define  $\mathbf{b}^*(\theta) \doteq \operatorname{argmin}_{\mathbf{b}} F(\mathbf{b}, \theta)$ , such that  $A\mathbf{b} = d$  for some convex function F. Then,

$$\frac{\partial \mathbf{b}^*(\theta)}{\partial \theta_j} = (D^{-1}A^T (AD^{-1}A^T)^{-1}AD^{-1} - D^{-1})\frac{\partial^2 F}{\partial \mathbf{b} \partial \theta_j}, \ D = (\frac{\partial^2 F}{\partial \mathbf{b} \partial \mathbf{b}^T}).$$

In our case,  $D = \operatorname{diag}(\sum_{f} \mathbf{w}_{f}(\mathbf{y}) \odot f''(\mathbf{b}))$   $\frac{\partial^{2} F}{\partial \mathbf{b} \partial \theta_{j}} = \sum_{f} \frac{\partial \mathbf{w}_{f}(\mathbf{y})}{\partial \theta_{j}} \odot f'$ 

We still need  $\frac{\partial \mathbf{w}_f(\mathbf{y})}{\partial \theta_j}$ .



Experiments 0000 000

# Derivatives of beliefs- the bottom line

- Want to calculate  $\frac{\partial L}{\partial \theta_i}$ . Procedure:
  - 1. Run some optimization

$$\begin{split} \mathbf{b}^* &= \arg\min_{\mathbf{b}} \sum_{f \in \mathscr{F}} \mathbf{w}_f(\mathbf{y})^T f(\mathbf{b}) \\ \text{such that} \qquad A\mathbf{b} &= \mathbf{d} \\ \mathbf{b} &\geq \mathbf{0}. \end{split}$$

2. Solve the linear system given by Claim 2 to get  $\frac{\partial L}{\partial \theta_i}$ .

- When learning, <u>two</u> optimizations:
  - 1. "inner" optimization (over  $\{b_r\}$ )
  - 2. "outer" optimization (over  $\theta$ )



Experiments 0000 000

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

### Derivatives of beliefs- the bottom line

- Want to calculate  $\frac{\partial L}{\partial \theta_i}$ . Procedure:
  - 1. Run some optimization

$$\begin{split} \mathbf{b}^* &= \arg\min_{\mathbf{b}} \sum_{f \in \mathscr{F}} \mathbf{w}_f(\mathbf{y})^T f(\mathbf{b}) \\ \text{such that} \qquad & A\mathbf{b} = \mathbf{d} \\ & \mathbf{b} \geq \mathbf{0}. \end{split}$$

2. Solve the linear system given by Claim 2 to get  $\frac{\partial L}{\partial \theta_i}$ .

- When learning, <u>two</u> optimizations:
  - 1. "inner" optimization (over  $\{b_r\}$ )
  - 2. "outer" optimization (over  $\theta$ )



Experiments 0000 000

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

### Derivatives of beliefs- the bottom line

- Want to calculate  $\frac{\partial L}{\partial \theta_i}$ . Procedure:
  - 1. Run some optimization

$$\begin{split} \mathbf{b}^* &= \arg\min_{\mathbf{b}} \sum_{f \in \mathscr{F}} \mathbf{w}_f(\mathbf{y})^T f(\mathbf{b}) \\ \text{such that} \qquad & A\mathbf{b} = \mathbf{d} \\ & \mathbf{b} \geq \mathbf{0}. \end{split}$$

2. Solve the linear system given by Claim 2 to get  $\frac{\partial L}{\partial \theta_i}$ .

- When learning, <u>two</u> optimizations:
  - 1. "inner" optimization (over  $\{b_r\}$ )
  - 2. "outer" optimization (over  $\theta$ )

Approach 000 0000000 000000 Experiments

(日) (四) (三) (三) (三)

Discussion 000 0000000000

# Outline

#### Introduction

Motivation Overview of the approach

#### Approach

Inference Learning (Loss functions) Learning (Derivatives of beliefs)

# Experiments

### Introduction

Results

### Discussion



Experiments

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへで

# Problem

- "Denoising" of 10 binary images of each class (1-9) from the MNIST database.
  - **y** is the observed, noisy image



• x is the unobserved, clean image



Approach 000 0000000 000000 Experiments

- Use regions consisting of individual variables, and neighboring pairs.
- Total of 40 parameters. (Somewhat redundant)
  - $w_b(\mathbf{x}_c, \mathbf{y}_c)$  (initialize to 0)
  - $w_b(\mathbf{x}_i, \mathbf{y}_i)$  (initialize to 0)
  - $w_{b \log b}(\mathbf{x}_c, \mathbf{y}_c)$  (initialize to 1)
  - $w_{b \log b}(\mathbf{x}_i, \mathbf{y}_i)$  (initialize to 1)
- "Inner loop" optimization uses PDCO interior method.
  - Tolerances very strict.
- "Outer loop" optimization uses Matlab's BFGS.



Experiments

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

- Use regions consisting of individual variables, and neighboring pairs.
- Total of 40 parameters. (Somewhat redundant)
  - $w_b(\mathbf{x}_c, \mathbf{y}_c)$  (initialize to 0)
  - $w_b(\mathbf{x}_i, \mathbf{y}_i)$  (initialize to 0)
  - $w_{b \log b}(\mathbf{x}_c, \mathbf{y}_c)$  (initialize to 1)
  - $w_{b \log b}(\mathbf{x}_i, \mathbf{y}_i)$  (initialize to 1)
- "Inner loop" optimization uses PDCO interior method.
  - Tolerances very strict.
- "Outer loop" optimization uses Matlab's BFGS.



Experiments

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

- Use regions consisting of individual variables, and neighboring pairs.
- Total of 40 parameters. (Somewhat redundant)
  - $w_b(\mathbf{x}_c, \mathbf{y}_c)$  (initialize to 0)
  - $w_b(\mathbf{x}_i, \mathbf{y}_i)$  (initialize to 0)
  - $w_{b \log b}(\mathbf{x}_c, \mathbf{y}_c)$  (initialize to 1)
  - $w_{b \log b}(\mathbf{x}_i, \mathbf{y}_i)$  (initialize to 1)
- "Inner loop" optimization uses PDCO interior method.
  - Tolerances very strict.
- "Outer loop" optimization uses Matlab's BFGS.



Experiments

Discussion 000 0000000000

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ のQ@

- Use regions consisting of individual variables, and neighboring pairs.
- Total of 40 parameters. (Somewhat redundant)
  - $w_b(\mathbf{x}_c, \mathbf{y}_c)$  (initialize to 0)
  - $w_b(\mathbf{x}_i, \mathbf{y}_i)$  (initialize to 0)
  - $w_{b \log b}(\mathbf{x}_c, \mathbf{y}_c)$  (initialize to 1)
  - $w_{b \log b}(\mathbf{x}_i, \mathbf{y}_i)$  (initialize to 1)
- "Inner loop" optimization uses PDCO interior method.
  - Tolerances very strict.
- "Outer loop" optimization uses Matlab's BFGS.

Approach 000 0000000 000000 Experiments

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ のQ@

## Comparison

- Compare to CRF toolbox.
  - Vishwanathan et al., ICML 2006
  - www.cs.ubc.ca/~murphyk/Software/CRF/crf.html
- Inference: mean-field or (loopy) belief propagation.
- Learning: Pseudolikelihood, or surrogate to likelihood where inference algorithm is used to approximate marginals defining the gradient.
- Features are constant + indicator for each possible configuration of each variable/pair.

Approach 000 0000000 000000 Experiments

Discussion 000 0000000000

# Outline

#### Introduction

Motivation Overview of the approac

#### Approach

Inference Learning (Loss functions) Learning (Derivatives of beliefs)

### Experiments

Introduction

### Results

### Discussion

Approach 000 0000000 000000 Experiments

Discussion 000 000000000

# 50% noise



▲ロ▶ ▲圖▶ ▲画▶ ▲画▶ 三回 - のQで

Approach 000 0000000 0000000 Experiments

Discussion 000 0000000000

## 50% noise





◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへで

Approach 000 0000000 000000 Experiments

Discussion

# Outline

#### Introduction

Motivation Overview of the approach

#### Approach

Inference Learning (Loss functions) Learning (Derivatives of beliefs)

#### Experiments

Introduction Results

### Discussion Summary

▲□▶ ▲圖▶ ▲≣▶ ▲≣▶ = 三 のQQ

Approach 000 0000000 0000000 Experiments

<ロ> (四) (四) (三) (三) (三)

Discussion 000

- Common to fit classifiers by minimizing a loss close to classification error.
- Graphical models, however, usually use a score (the likelihood) that is both
  - Difficult to optimize.
  - Remote from empirical risk.
- This paper presents an approach for "fitting a free energy" to directly give good marginals.

Approach 000 0000000 0000000 Experiments

Discussion

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

- Common to fit classifiers by minimizing a loss close to classification error.
- Graphical models, however, usually use a score (the likelihood) that is both
  - Difficult to optimize.
  - Remote from empirical risk.
- This paper presents an approach for "fitting a free energy" to directly give good marginals.

Approach 000 0000000 000000 Experiments

Discussion 000

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

- Common to fit classifiers by minimizing a loss close to classification error.
- Graphical models, however, usually use a score (the likelihood) that is both
  - Difficult to optimize.
  - Remote from empirical risk.
- This paper presents an approach for "fitting a free energy" to directly give good marginals.

Approach 000 0000000 000000 Experiments

Discussion

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ □ の Q @

# The End

### Thank you.

- 1. Different loss functions.
- 2. Better entropy term.
- 3. Exploit hidden variables.
- 4. Better relaxation of marginal polytope.
- 5. Better "inner loop" optimization.
- 6. Better F.

Approach 000 0000000 000000 Experiments

Discussion

# Outline

#### Introduction

Motivation Overview of the approach

#### Approach

Inference Learning (Loss functions) Learning (Derivatives of beliefs)

### Experiments

Introduction Results

### Discussion

Summary

▲□▶ ▲□▶ ▲目▶ ▲目▶ 三目 - のへで

Approach 000 0000000 000000 Experiments

Discussion

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

### Future Work

- 1. Different loss functions.
  - For example, common to find  $\operatorname{argmax}_{x_i} b^*(x_i | \mathbf{y})$ .
  - Could use a loss function (approximating) the univariate classification error.

(Gross et al. "Training CRFs for maximum labelwise accuracy", NIPS 2006)

Approach 000 0000000 000000 Exp eriments

Discussion

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

# Future Work

- 2. More flexible entropy term.
  - As described, this approach requires fixing  $w_{b\log b} > 0$  to guarantee convexity.
  - However, negative entropy terms can be allowed while preserving convexity over the locally consistent marginals.
  - Easiest approach- just pick a better (fixed) entropy approximation.

(Heskes, "Convexity Arguments for Efficient Minimization of the Bethe and Kikuchi Free Energies", JAIR 2006)

Approach 000 0000000 0000000 Experiments 0000 000 Discussion

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへで

# Entropy

 $w_{b\log b}(\mathbf{x}_c, \mathbf{y}_c)$  and  $w_{b\log b}(x_i, y_i)$  for  $L_{quad}$  with 50% noise

$\mathbf{x}_c \setminus \mathbf{y}_c$	(0,0)	(0,1)	(1,0)	(1,1)			
(0,0)	4.86	0.04	0.05	0.02	$x_i \setminus y_i$	0	1
(0,1)	4.22	3.86	4.54	5.00	0	4.47	0.02
(1,0)	4.13	4.49	2.14	5.13	1	0.03	0.03
(1,1)	0.06	0.02	0.03	0.02			

Approach 000 0000000 000000 Exp eriments

Discussion

◆ロト ◆得ト ◆ヨト ◆ヨト ヨー のくべ

- 3. Hidden variables.
  - The same algorithm can be used with hidden variables, by taking the sum over the variables in L<sub>{log,quad}</sub> over the observed variables.



Experiments 0000 000 Discussion

◆ロト ◆得ト ◆ヨト ◆ヨト ヨー のくべ

- 4. Better relaxation of the marginal polytope.
  - Currently, the model must try to compensate during the learning stage for defects in the marginal polytope.

Approach 000 0000000 000000 Experiments 0000 000 Discussion

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへで

- 5. Better "inner loop" optimization.
  - Generic optimization is OK for 28x28 images.
  - Derive a message passing algorithm?

Approach 000 0000000 000000 Exp eriments

Discussion

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

- 6. More general function F.
  - Any function  $F(\mathbf{y}, \{b_r(\mathbf{x}_r)\})$  can give an implicit mapping.
    - As long as it is convex and continuous, learning should be doable by implicit differentiation.
  - A larger set  $\mathscr{F}$  than  $\{b, b \log b\}$ ?
  - Some F based on different principles?

Approach 000 0000000 000000 Experiments 0000 000 Discussion

# 30% noise



▲ロ▶ ▲圖▶ ▲画▶ ▲画▶ 三回 - のQで

Approach 000 0000000 0000000 Experiments 0000 000 Discussion

## 30% noise





◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへで