Part 2: Introduction to Graphical Models

Sebastian Nowozin and Christoph H. Lampert

Providence, 21st June 2012



<ロト < 部 > < 注 > < 注 > こ つ < で</p>

Sebastian Nowozin and Christoph H. Lampert Part 2: Introduction to Graphical Models

Graphical Models		
000000		
Graphical Models		

- Model: relating observations x to quantities of interest y
- Example 1: given RGB image x, infer depth y for each pixel
- Example 2: given RGB image x, infer presence and positions y of all objects shown



Graphical Models		
000000		
Graphical Models		

- Model: relating observations x to quantities of interest y
- Example 1: given RGB image x, infer depth y for each pixel
- Example 2: given RGB image x, infer presence and positions y of all objects shown



・ロト ・伺 ト ・ヨト ・ヨト



 $\mathcal{X}:$ image, $\mathcal{Y}:$ object annotations

Graphical Models		
000000		
Graphical Models		

- General case: mapping $x \in \mathcal{X}$ to $y \in \mathcal{Y}$
- Graphical models are a concise language to define this mapping
- Mapping can be *ambiguous*: measurement noise, lack of well-posedness (e.g. occlusions)
- ▶ Probabilistic graphical models: define form p(y|x) or p(x, y) for all y ∈ Y



Graphical Models		
000000		
Graphical Models		

- General case: mapping $x \in \mathcal{X}$ to $y \in \mathcal{Y}$
- Graphical models are a concise language to define this mapping
- Mapping can be *ambiguous*: measurement noise, lack of well-posedness (e.g. occlusions)
- ▶ Probabilistic graphical models: define form p(y|x) or p(x, y) for all y ∈ Y



Graphical Models		
000000		
Graphical Models		

Graphical Models

A graphical model defines

- ▶ a family of probability distributions over a set of random variables,
- ▶ by means of a graph,
- so that the random variables satisfy *conditional independence assumptions* encoded in the graph.

Graphical Models		
000000		
Graphical Models		

Graphical Models

A graphical model defines

- ▶ a family of probability distributions over a set of random variables,
- by means of a graph,
- so that the random variables satisfy *conditional independence assumptions* encoded in the graph.

Popular classes of graphical models,

- Undirected graphical models (Markov random fields),
- Directed graphical models (Bayesian networks),
- Factor graphs,
- Others: chain graphs, influence diagrams, etc.



Graphical Models		
000000		
Graphical Models		

Bayesian Networks

• Graph:
$$G = (V, \mathcal{E}), \ \mathcal{E} \subset V \times V$$

- directed
- acyclic
- ► Variable domains \mathcal{Y}_i
- Factorization

$$p(Y = y) = \prod_{i \in V} p(y_i | y_{\mathrm{pa}_G(i)})$$

over distributions, by conditioning on parent nodes.



A simple Bayes net

3

Example

$$p(Y = y) = p(Y_l = y_l | Y_k = y_k) p(Y_k = y_k | Y_i = y_i, Y_j = y_j)$$

$$p(Y_i = y_i) p(Y_j = y_j).$$

Graphical Models		
000000		
Graphical Models		

Bayesian Networks

• Graph:
$$G = (V, \mathcal{E}), \ \mathcal{E} \subset V \times V$$

- directed
- acyclic
- Variable domains \mathcal{Y}_i
- Factorization

$$p(Y = y) = \prod_{i \in V} p(y_i | y_{\mathrm{pa}_G(i)})$$

over distributions, by conditioning on parent nodes.

Example

$$p(Y = y) = p(Y_i = y_i | Y_k = y_k) p(Y_k = y_k | Y_i = y_i, Y_j = y_j)$$

$$p(Y_i = y_i) p(Y_j = y_j).$$



A simple Bayes net

<ロト <部ト < 注ト < 注ト 3

Graphical Models		
0000000		
Graphical Models		

Undirected Graphical Models

Markov random field (MRF) = Markov network

• Graph:
$$G = (V, \mathcal{E}), \ \mathcal{E} \subset V \times V$$

- undirected, no self-edges
- ► Variable domains \mathcal{Y}_i
- Factorization over potentials ψ at *cliques*,

$$p(y) = \frac{1}{Z} \prod_{C \in \mathcal{C}(G)} \psi_C(y_C)$$

• Constant
$$Z = \sum_{y \in \mathcal{Y}} \prod_{C \in \mathcal{C}(G)} \psi_C(y_C)$$

Example

$$p(y) = \frac{1}{Z}\psi_i(y_i)\psi_j(y_j)\psi_l(y_l)\psi_{i,j}(y_i, y_j)$$



A simple MRF

イロト イポト イヨト イヨト

3

Graphical Models		
0000000		
Graphical Models		

Undirected Graphical Models

 Markov random field (MRF) = Markov network

• Graph:
$$G = (V, \mathcal{E}), \ \mathcal{E} \subset V \times V$$

- undirected, no self-edges
- ► Variable domains \mathcal{Y}_i
- Factorization over potentials ψ at *cliques*,

$$p(y) = \frac{1}{Z} \prod_{C \in \mathcal{C}(G)} \psi_C(y_C)$$

• Constant $Z = \sum_{y \in \mathcal{Y}} \prod_{C \in \mathcal{C}(G)} \psi_C(y_C)$

Example

$$p(y) = \frac{1}{Z}\psi_i(y_i)\psi_j(y_j)\psi_l(y_l)\psi_{i,j}(y_i,y_j)$$



A simple MRF

イロト イポト イヨト

3

Graphical Models		
0000000		
Graphical Models		

Example 1

►



- Cliques C(G): set of vertex sets V' with $V' \subseteq V$, $\mathcal{E} \cap (V' \times V') = V' \times V'$
- Here $C(G) = \{\{i\}, \{i, j\}, \{j\}, \{j, k\}, \{k\}\}$

$$p(y) = \frac{1}{Z} \psi_i(y_i) \psi_j(y_j) \psi_l(y_l) \psi_{i,j}(y_i, y_j)$$

Graphical Models		
000000		
Graphical Models		

Example 2



Factor Graphs		
•000000000		

Factor Graphs

- Graph: $G = (V, \mathcal{F}, \mathcal{E})$, $\mathcal{E} \subseteq V \times \mathcal{F}$
 - variable nodes V,
 - ► factor nodes *F*,
 - \blacktriangleright edges ${\mathcal E}$ between variable and factor nodes.
 - ► scope of a factor, $N(F) = \{i \in V : (i, F) \in \mathcal{E}\}$
- Variable domains *Y_i*
- Factorization over potentials ψ at *factors*,

$$p(y) = \frac{1}{Z} \prod_{F \in \mathcal{F}} \psi_F(y_{N(F)})$$

• Constant
$$Z = \sum_{y \in \mathcal{Y}} \prod_{F \in \mathcal{F}} \psi_F(y_{N(F)})$$



Factor graph

3

イロト 不得下 イヨト イヨト

Factor Graphs		
•••••		

Factor Graphs

- Graph: $G = (V, \mathcal{F}, \mathcal{E})$, $\mathcal{E} \subseteq V \times \mathcal{F}$
 - ► variable nodes V,
 - factor nodes \mathcal{F} ,
 - \blacktriangleright edges ${\mathcal E}$ between variable and factor nodes.
 - ► scope of a factor, $N(F) = \{i \in V : (i, F) \in \mathcal{E}\}$
- Variable domains *Y_i*
- Factorization over potentials ψ at *factors*,

$$p(y) = \frac{1}{Z} \prod_{F \in \mathcal{F}} \psi_F(y_{N(F)})$$

• Constant
$$Z = \sum_{y \in \mathcal{Y}} \prod_{F \in \mathcal{F}} \psi_F(y_{N(F)})$$



Factor graph

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

	Factor Graphs		
	0000000000		
Factor Graphs			

Why factor graphs?



- Factor graphs are *explicit* about the factorization
- Hence, easier to work with
- Universal (just like MRFs and Bayesian networks)

	Factor Graphs		
	0000000000		
Factor Graphs			

Capacity



イロト イヨト イヨト イヨト

э

- ► Factor graph defines family of distributions
- Some families are larger than others

	Factor Graphs		
	0000000000		
Factor Graphs			

Four remaining pieces

1. Conditional distributions (CRFs)

2. Parameterization

- 3. Test-time inference
- 4. Learning the model from training data

	Factor Graphs		
	0000000000		
Factor Graphs			

Four remaining pieces

- 1. Conditional distributions (CRFs)
- 2. Parameterization
- 3. Test-time inference
- 4. Learning the model from training data

	Factor Graphs		
	0000000000		
Factor Graphs			

Conditional Distributions

- We have discussed p(y),
- How do we define p(y|x)?
- Potentials become a function of x_{N(F}
- Partition function depends on x
- Conditional random fields (CRFs)
- x is not part of the probability model, i.e. not treated as random variable



conditional distribution

イロト イポト イヨト イヨト

	Factor Graphs		
	0000000000		
Factor Graphs			

Conditional Distributions

- We have discussed p(y),
- How do we define p(y|x)?
- Potentials become a function of x_{N(F)}
- Partition function depends on x
- Conditional random fields (CRFs)
- x is not part of the probability model, i.e. not treated as random variable



conditional distribution

$$p(y) = \frac{1}{Z} \prod_{F \in \mathcal{F}} \psi_F(y_{N(F)})$$

$$p(y|x) = \frac{1}{Z(x)} \prod_{F \in \mathcal{F}} \psi_F(y_{N(F)}; x_{N(F)})$$

	Factor Graphs		
	0000000000		
Factor Graphs			

Conditional Distributions

- We have discussed p(y),
- How do we define p(y|x)?
- Potentials become a function of $x_{N(F)}$
- Partition function depends on x
- Conditional random fields (CRFs)
- x is not part of the probability model, i.e. not treated as random variable

$$p(y) = \frac{1}{Z} \prod_{F \in \mathcal{F}} \psi_F(y_{N(F)})$$

$$p(y|x) = \frac{1}{Z(x)} \prod_{F \in \mathcal{F}} \psi_F(y_{N(F)}; x_{N(F)})$$



conditional distribution

イロト イポト イヨト イヨト

	Factor Graphs		
	0000000000		
Eactor Graphs			

Potentials and Energy Functions

► For each factor
$$F \in \mathcal{F}$$
, $\mathcal{Y}_F = \underset{i \in N(F)}{X} \mathcal{Y}_i$,

$$E_F: \mathcal{Y}_{N(F)} \to \mathbb{R},$$

• Potentials and energies (assume $\psi_F(y_F) > 0$)

 $\psi_F(y_F) = \exp(-E_F(y_F)), \text{ and } E_F(y_F) = -\log(\psi_F(y_F)).$

• Then p(y) can be written as

$$p(Y = y) = \frac{1}{Z} \prod_{F \in \mathcal{F}} \psi_F(y_F)$$
$$= \frac{1}{Z} \exp(-\sum_{F \in \mathcal{F}} E_F(y_F)),$$

► Hence, p(y) is completely determined by $E(y) = \sum_{F \in \mathcal{F}} E_F(y_F)$

Sebastian Nowozin and Christoph H. Lampert Part 2: Introduction to Graphical Models

	Factor Graphs		
	0000000000		
Eactor Graphs			

Potentials and Energy Functions

► For each factor
$$F \in \mathcal{F}$$
, $\mathcal{Y}_F = \underset{i \in N(F)}{X} \mathcal{Y}_i$,

$$E_F: \mathcal{Y}_{N(F)} \to \mathbb{R},$$

• Potentials and energies (assume $\psi_F(y_F) > 0$)

 $\psi_F(y_F) = \exp(-E_F(y_F)), \quad \text{and} \quad E_F(y_F) = -\log(\psi_F(y_F)).$

• Then p(y) can be written as

$$p(Y = y) = \frac{1}{Z} \prod_{F \in \mathcal{F}} \psi_F(y_F)$$
$$= \frac{1}{Z} \exp(-\sum_{F \in \mathcal{F}} E_F(y_F)),$$

► Hence, p(y) is completely determined by $E(y) = \sum_{F \in \mathcal{F}} E_F(y_F)$

	Factor Graphs		
	0000000000		
Eactor Graphs			

Potentials and Energy Functions

► For each factor
$$F \in \mathcal{F}$$
, $\mathcal{Y}_F = \underset{i \in N(F)}{X} \mathcal{Y}_i$,

$$E_F: \mathcal{Y}_{N(F)} \to \mathbb{R},$$

• Potentials and energies (assume $\psi_F(y_F) > 0$)

 $\psi_F(y_F) = \exp(-E_F(y_F)), \quad \text{and} \quad E_F(y_F) = -\log(\psi_F(y_F)).$

• Then p(y) can be written as

$$p(Y = y) = \frac{1}{Z} \prod_{F \in \mathcal{F}} \psi_F(y_F)$$
$$= \frac{1}{Z} \exp(-\sum_{F \in \mathcal{F}} E_F(y_F)),$$

• Hence, p(y) is completely determined by $E(y) = \sum_{F \in \mathcal{F}} E_F(y_F)$

	Factor Graphs		
	0000000000		
Factor Graphs			

Energy Minimization

$$\begin{aligned} \underset{y \in \mathcal{Y}}{\operatorname{argmax}} p(Y = y) &= \operatorname{argmax}_{y \in \mathcal{Y}} \frac{1}{Z} \exp(-\sum_{F \in \mathcal{F}} E_F(y_F)) \\ &= \operatorname{argmax}_{y \in \mathcal{Y}} \exp(-\sum_{F \in \mathcal{F}} E_F(y_F)) \\ &= \operatorname{argmax}_{y \in \mathcal{Y}} - \sum_{F \in \mathcal{F}} E_F(y_F) \\ &= \operatorname{argmin}_{y \in \mathcal{Y}} \sum_{F \in \mathcal{F}} E_F(y_F) \\ &= \operatorname{argmin}_{y \in \mathcal{Y}} E(y). \end{aligned}$$

Energy minimization can be interpreted as solving for the most likely state of some factor graph model

	Factor Graphs 00000000000		
Factor Graphs			

Parameterization

- ► Factor graphs define a family of distributions
- ▶ Parameterization: identifying individual members by parameters w

	Factor Graphs 00000000000		
Factor Graphs			

Parameterization

- Factor graphs define a family of distributions
- ▶ Parameterization: identifying individual members by parameters w



Sebastian Nowozin and Christoph H. Lampert

	Factor Graphs		
	00000000000		
Factor Graphs			

Example: Parameterization

- Image segmentation model
- ► Pairwise "Potts" energy function E_F(y_i, y_j; w₁),

$$E_F: \{0,1\} \times \{0,1\} \times \mathbb{R} \to \mathbb{R},$$

•
$$E_F(0,0; w_1) = E_F(1,1; w_1) = 0$$

•
$$E_F(0,1;w_1) = E_F(1,0;w_1) = w_1$$



image segmentation model

	Factor Graphs		
	00000000000		
Factor Graphs			

Example: Parameterization (cont)

- Image segmentation model
- Unary energy function $E_F(y_i; x, w)$,

$$E_F: \{0,1\} \times \mathcal{X} \times \mathbb{R}^{\{0,1\} \times D} \to \mathbb{R},$$

•
$$E_F(0; x, w) = \langle w(0), \psi_F(x) \rangle$$

•
$$E_F(1; x, w) = \langle w(1), \psi_F(x) \rangle$$

▶ Features $\psi_F : \mathcal{X} \to \mathbb{R}^D$, e.g. image filters



image segmentation model

	Factor Graphs		
	0000000000		
Factor Graphs			

Example: Parameterization (cont)



Sebastian Nowozin and Christoph H. Lampert Part 2: Introduction to Graphical Models
 Graphical Models
 Factor Graphs
 Test-time Inference
 Training
 Software

 0000000
 000000000
 000000000
 0000
 000
 000

 Factor Graphs
 000000000
 000000000
 0000
 000
 000

Example: Parameterization (cont)



- Total number of parameters: D + D + 1
- ▶ Parameters are *shared*, but energies differ because of different $\psi_F(x)$
- ▶ General form, linear in w,

$$E_F(y_F; x_F, w) = \langle w(y_F), \psi_F(x_F) \rangle$$

	Test-time Inference	
	●0000000	
Test-time Inference		

Making Predictions

- Making predictions: given $x \in \mathcal{X}$, predict $y \in \mathcal{Y}$
- How to measure quality of prediction? (or function $f : \mathcal{X} \to \mathcal{Y}$)

	Test-time Inference	
	0000000	

Loss function

► Define a loss function

$$\Delta: \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}^+,$$

so that $\Delta(y, y^*)$ measures the loss incurred by predicting y when y^* is true.

▶ The *loss function* is application dependent



	Test-time Inference 00●000000	

Test-time Inference

• Loss function $\Delta(y, f(x))$: correct label y, predict f(x)

 $\Delta:\mathcal{Y}\times\mathcal{Y}\to\mathbb{R}$

- ▶ True joint distribution d(X, Y) and true conditional d(y|x)
- Model distribution p(y|x)

Expected loss: quality of prediction

$$\mathcal{R}_{f}^{\Delta}(x) = \mathbb{E}_{y \sim d(y|x)} \Delta(y, f(x))$$
$$= \sum_{y \in \mathcal{Y}} d(y|x) \Delta(y, f(x)).$$
$$\approx \mathbb{E}_{y \sim p(y|x|y)} \Delta(y, f(x))$$

イロト イポト イヨト イヨト

• Assuming that $p(y|x;w) \approx d(y|x)$

	Test-time Inference	
Test-time Inference		

Test-time Inference

• Loss function $\Delta(y, f(x))$: correct label y, predict f(x)

$$\Delta:\mathcal{Y}\times\mathcal{Y}\to\mathbb{R}$$

- True joint distribution d(X, Y) and true conditional d(y|x)
- Model distribution p(y|x)
- Expected loss: quality of prediction

$$\mathcal{R}_{f}^{\Delta}(x) = \mathbb{E}_{y \sim d(y|x)} \Delta(y, f(x))$$
$$= \sum_{y \in \mathcal{Y}} d(y|x) \Delta(y, f(x)).$$
$$\approx \mathbb{E}_{y \sim p(y|x;w)} \Delta(y, f(x))$$

イロト 不得下 イヨト イヨト 二日

• Assuming that $p(y|x;w) \approx d(y|x)$

	Test-time Inference	
Test-time Inference		

Test-time Inference

• Loss function $\Delta(y, f(x))$: correct label y, predict f(x)

$$\Delta:\mathcal{Y}\times\mathcal{Y}\to\mathbb{R}$$

- True joint distribution d(X, Y) and true conditional d(y|x)
- Model distribution p(y|x)
- Expected loss: quality of prediction

$$\begin{aligned} \mathcal{R}_{f}^{\Delta}(x) &= \mathbb{E}_{y \sim d(y|x)} \Delta(y, f(x)) \\ &= \sum_{y \in \mathcal{Y}} d(y|x) \Delta(y, f(x)). \\ &\approx \mathbb{E}_{y \sim p(y|x;w)} \Delta(y, f(x)) \end{aligned}$$

イロト イポト イヨト イヨト 二日

• Assuming that $p(y|x;w) \approx d(y|x)$

	Test-time Inference 000●00000	

Example 1: 0/1 loss

Loss 0 iff perfectly predicted, 1 otherwise:

$$\Delta_{0/1}(y,y^*) = I(y
eq y^*) = \left\{egin{array}{cc} 0 & ext{if } y = y^* \ 1 & ext{otherwise} \end{array}
ight.$$

Plugging it in,

$$\begin{array}{rcl} y^{*} & := & \operatorname*{argmin}_{y' \in \mathcal{Y}} \mathbb{E}_{y \sim p(y|x)} \left[\Delta_{0/1}(y, y') \right] \\ & = & \operatorname*{argmax}_{y' \in \mathcal{Y}} p(y'|x) \\ & = & \operatorname*{argmin}_{y' \in \mathcal{Y}} E(y', x). \end{array}$$

Minimizing the expected 0/1-loss → MAP prediction (energy minimization)

イロト イポト イヨト イヨト

3

	Test-time Inference 000●00000	

Example 1: 0/1 loss

Loss 0 iff perfectly predicted, 1 otherwise:

$$\Delta_{0/1}(y,y^*) = I(y
eq y^*) = \left\{egin{array}{cc} 0 & ext{if } y = y^* \ 1 & ext{otherwise} \end{array}
ight.$$

Plugging it in,

$$\begin{array}{lll} y^* & := & \operatorname*{argmin}_{y' \in \mathcal{Y}} \mathbb{E}_{y \sim p(y|x)} \left[\Delta_{0/1}(y, y') \right] \\ & = & \operatorname*{argmax}_{y' \in \mathcal{Y}} p(y'|x) \\ & = & \operatorname*{argmin}_{y' \in \mathcal{Y}} E(y', x). \end{array}$$

Minimizing the expected 0/1-loss → MAP prediction (energy minimization)

イロト イポト イヨト イヨト

	Test-time Inference	
	00000000	

Example 2: Hamming loss

Count the number of mislabeled variables:

$$\Delta_H(y,y^*) = \frac{1}{|V|} \sum_{i \in V} I(y_i \neq y_i^*)$$



Plugging it in,

$$y^* := \operatorname{argmin}_{y' \in \mathcal{Y}} \mathbb{E}_{y \sim p(y|x)} \left[\Delta_H(y, y') \right]$$
$$= \left(\operatorname{argmax}_{y'_i \in \mathcal{Y}_i} p(y'_i|x) \right)_{i \in V}$$

► Minimizing the expected Hamming loss → maximum posterior marginal (MPM, Max-Marg) prediction

	Test-time Inference	
	00000000	

Example 2: Hamming loss

Count the number of mislabeled variables:

$$\Delta_H(y,y^*) = \frac{1}{|V|} \sum_{i \in V} I(y_i \neq y_i^*)$$



Plugging it in,

$$y^* := \operatorname{argmin}_{y' \in \mathcal{Y}} \mathbb{E}_{y \sim p(y|x)} \left[\Delta_H(y, y') \right]$$
$$= \left(\operatorname{argmax}_{y'_i \in \mathcal{Y}_i} \frac{p(y'_i|x)}{y'_i \in \mathcal{Y}_i} \right)_{i \in V}$$

► Minimizing the expected Hamming loss → maximum posterior marginal (MPM, Max-Marg) prediction

	Test-time Inference	
	000000000	

Example 3: Squared error

Assume a vector space on \mathcal{Y}_i (pixel intensities, optical flow vectors, etc.). Sum of squared errors

$$\Delta_Q(y, y^*) = \frac{1}{|V|} \sum_{i \in V} ||y_i - y_i^*||^2.$$



Plugging it in,

$$\begin{array}{lll} y^{*} & := & \operatorname*{argmin}_{y' \in \mathcal{Y}} \mathbb{E}_{y \sim p(y|x)} \left[\Delta_{Q}(y, y') \right] \\ & = & \left(\sum_{y'_{i} \in \mathcal{Y}_{i}} p(y'_{i}|x) y'_{i} \right)_{i \in V} \end{array}$$

► Minimizing the expected squared error → minimum mean squared error (MMSE) prediction

Sebastian Nowozin and Christoph H. Lampert Part 2: Introduction to Graphical Models

	Test-time Inference	
	000000000	

Example 3: Squared error

Assume a vector space on \mathcal{Y}_i (pixel intensities, optical flow vectors, etc.). Sum of squared errors

$$\Delta_Q(y, y^*) = \frac{1}{|V|} \sum_{i \in V} ||y_i - y_i^*||^2.$$



Plugging it in,

$$\begin{array}{lll} y^* & := & \operatorname*{argmin}_{y' \in \mathcal{Y}} \mathbb{E}_{y \sim p(y|x)} \left[\Delta_Q(y, y') \right] \\ & = & \left(\sum_{y'_i \in \mathcal{Y}_i} p(y'_i|x) y'_i \right)_{i \in V} \end{array}$$

► Minimizing the expected squared error → minimum mean squared error (MMSE) prediction

Sebastian Nowozin and Christoph H. Lampert Part 2: Introduction to Graphical Models

	Test-time Inference	
	000000000	

Inference Task: Maximum A Posteriori (MAP) Inference

Definition (Maximum A Posteriori (MAP) Inference)

Given a factor graph, parameterization, and weight vector w, and given the observation x, find

$$y^* = \underset{y \in \mathcal{Y}}{\operatorname{argmax}} p(Y = y | x, w) = \underset{y \in \mathcal{Y}}{\operatorname{argmin}} E(y; x, w).$$

	Test-time Inference	
	000000000	

Inference Task: Probabilistic Inference

Definition (Probabilistic Inference)

Given a factor graph, parameterization, and weight vector w, and given the observation x, find

$$\begin{split} \log Z(x,w) &= & \log \sum_{y \in \mathcal{Y}} \exp(-E(y;x,w)), \\ \mu_F(y_F) &= & p(Y_F = y_f | x, w), \quad \forall F \in \mathcal{F}, \forall y_F \in \mathcal{Y}_F. \end{split}$$

This typically includes variable marginals

$$\mu_i(y_i) = p(y_i|x,w)$$

イロト イポト イヨト イヨト

Sebastian Nowozin and Christoph H. Lampert Part 2: Introduction to Graphical Models

	Test-time Inference	
	00000000	
Test-time Inference		

Example: Man-made structure detection



- ► Left: input image x,
- Middle: ground truth labeling on 16-by-16 pixel blocks,
- Right: factor graph model
- ► Features: gradient and color histograms
- Estimate model parameters from \approx 60 training images

	Test-time Inference	
	00000000	
Test-time Inference		

Example: Man-made structure detection



- ► Left: input image x,
- ► Middle (probabilistic inference): visualization of the variable marginals p(y_i = "manmade" | x, w),

< ロ > < 同 > < 三 > < 三 > <

► Right (MAP inference): joint MAP labeling y* = argmax_{y∈y} p(y|x, w).

	Training	
	0000	

Training the Model

What can be learned?

- Model structure: factors
- Model variables: observed variables fixed, but we can add unobserved variables
- ► Factor energies: parameters

	Training	
	0000	

Training the Model

What can be learned?

- Model structure: factors
- Model variables: observed variables fixed, but we can add unobserved variables
- ► Factor energies: parameters

		Training	
		0000	
Training			

Training: Overview

 Assume a fully observed, independent and identically distributed (iid) sample set

$$\{(x^n, y^n)\}_{n=1,...,N}, \qquad (x^n, y^n) \sim d(X, Y)$$

- ► Goal: predict well,
- ► Alternative goal: first model d(y|x) well by p(y|x, w), then predict by minimizing the expected loss

	Training	
	0000	

Probabilistic Learning

Problem (Probabilistic Parameter Learning)

Let d(y|x) be the (unknown) conditional distribution of labels for a problem to be solved. For a parameterized conditional distribution p(y|x, w) with parameters $w \in \mathbb{R}^{D}$, probabilistic parameter learning is the task of finding a point estimate of the parameter w^* that makes $p(y|x, w^*)$ closest to d(y|x).

▶ We will discuss probabilistic parameter learning in detail.

	Training	
	0000	

Probabilistic Learning

Problem (Probabilistic Parameter Learning)

Let d(y|x) be the (unknown) conditional distribution of labels for a problem to be solved. For a parameterized conditional distribution p(y|x, w) with parameters $w \in \mathbb{R}^{D}$, probabilistic parameter learning is the task of finding a point estimate of the parameter w^* that makes $p(y|x, w^*)$ closest to d(y|x).

• We will discuss probabilistic parameter learning in detail.

		Training	
		0000	
Training			

Loss-Minimizing Parameter Learning

Problem (Loss-Minimizing Parameter Learning)

Let d(x, y) be the unknown distribution of data in labels, and let $\Delta : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$ be a loss function. Loss minimizing parameter learning is the task of finding a parameter value w^* such that the expected prediction risk

$$\mathbb{E}_{(x,y)\sim d(x,y)}[\Delta(y,f_p(x))]$$

イロト イポト イヨト

is as small as possible, where $f_p(x) = \operatorname{argmax}_{y \in \mathcal{Y}} p(y|x, w^*)$.

- Requires loss function at training time
- Directly learns a prediction function $f_p(x)$

		Training	
		0000	
Training			

Loss-Minimizing Parameter Learning

Problem (Loss-Minimizing Parameter Learning)

Let d(x, y) be the unknown distribution of data in labels, and let $\Delta : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$ be a loss function. Loss minimizing parameter learning is the task of finding a parameter value w^* such that the expected prediction risk

$$\mathbb{E}_{(x,y)\sim d(x,y)}[\Delta(y,f_p(x))]$$

イロト 不得下 イヨト イヨト 二日

is as small as possible, where $f_p(x) = \operatorname{argmax}_{y \in \mathcal{Y}} p(y|x, w^*)$.

- Requires loss function at training time
- Directly learns a prediction function $f_p(x)$

		Software	
		0	
Software			

Vision Software: Graphical Models

Inference-only

- OpenGM, University of Heidelberg
 C++, discrete factor graphs, irregular, higher-order, probabilistic inference and energy minimization, MIT license
- libDAI, Joris Mooij
 C++, discrete factor graphs, irregular, higher-order, mainly probabilistic inference, BSD license

► ALE, Lubor Ladicky

 $C{++},$ discrete factor graphs, regular/irregular, higher-order, energy minimization, proprietary license

		Software
		00

Vision Software: Graphical Models (cont)

Inference and Estimation

 JGMT, Justin Domke C++/Matlab, discrete factor graphs, regular/irregular, pairwise only, probabilistic inference, loss-based learning, license?

grante, Microsoft Research UK
 C++ with Matlab wrappers, discrete factor graphs,
 regular/irregular, higher-order, prob. inference and energy
 minimization, likelihood- and loss-based estimation, MSR-LA license

Factorie, UMass

Scala (Java), imperative discrete factor graphs, continuous/discrete/any-order, likelihood-based, Apache license

- Infer.Net, Microsoft Research UK C#, discrete/continuous, any-order (probabilistic programming), full Bayesian inference, MSR-LA license
- svm-struct-matlab, Andrea Vedaldi
 Matlab wrapper for SVMstruct (Thorsten Joachims)