

Part II: Structured Prediction

Alex Schwing & Raquel Urtasun

University of Toronto

December 12, 2015

Roadmap Towards Learning Deep Structured Models

① Part I: Deep Models



② Part II: Structured Models



③ Part III: Deep Structured Models



Part II: Structured Prediction



- How to predict in structured spaces?
- How to learn in structured spaces?

How to predict in structured spaces?

Binary Classification

Prediction answers a Yes-No question

$$x \rightarrow y \in \{\text{Yes}, \text{No}\}$$

Binary Classification

Prediction answers a Yes-No question

- Spam filtering ([Is this eMail spam?](#))

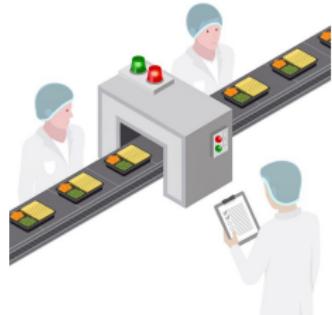
Subject: ICCV Tutorial
Text: Hey Raquel, ...

$$x \rightarrow y \in \{\text{Yes}, \text{No}\}$$

Binary Classification

Prediction answers a Yes-No question

- Spam filtering ([Is this eMail spam?](#))
- Quality control ([Can I sell this product?](#))



$$x \rightarrow y \in \{\text{Yes}, \text{No}\}$$

Binary Classification

Prediction answers a Yes-No question

- Spam filtering ([Is this eMail spam?](#))
- Quality control ([Can I sell this product?](#))
- Medical testing ([Does this lab analysis look normal?](#))

Dry Eye Progress Check	
Chart #	Date
Patient's Name:	(Last) _____ (First) _____ (Middle) _____
Current by Eye Treatment:	At present: _____ times per day. _____
Medical History (checkmark):	_____ pt notes to change: _____
Subjective Symptoms:	
None	Yes
Eye pain	Redness
Itching	Sensitivity
Watery eyes	Photophobia
Sticky eyes	Blurred vision
Other	Other
Visual Acuity:	
Left eye:	Right eye:
Distance	Intermediate
Close	Intermediate
See how clearly you can see the letters below:	
One line away	Two lines away
Three lines away	Four lines away
Five lines away	Six lines away
Other reading	Other
Total (in seconds) > 10 sec.	
Wet Eye Screening:	
Fluorescein	normal staining
Lurocortisone cream	1 2 3 4
Fluorometholone cream	1 2 3 4
Ranibizumab	1 2 3 4
Assessment:	
Plan:	

$$x \rightarrow y \in \{\text{Yes}, \text{No}\}$$

Binary Classification

Prediction answers a Yes-No question

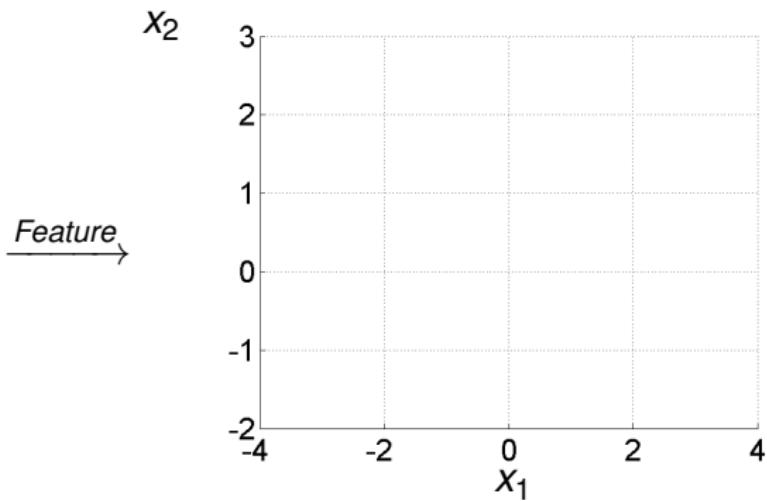
- Spam filtering ([Is this eMail spam?](#))
- Quality control ([Can I sell this product?](#))
- Medical testing ([Does this lab analysis look normal?](#))
- Driver assistance systems ([Is there an obstacle upfront?](#))



$$x \rightarrow y \in \{\text{Yes}, \text{No}\}$$

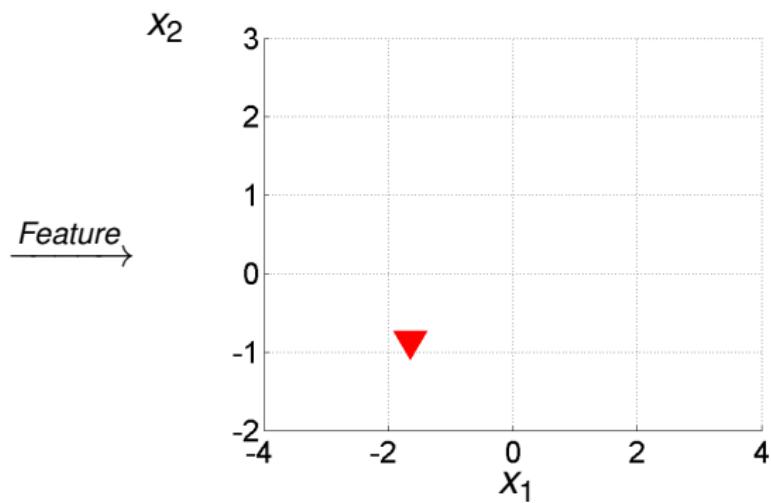
Binary Classification

How to find $y \in \{\text{Yes}, \text{No}\}$ given input data x ?



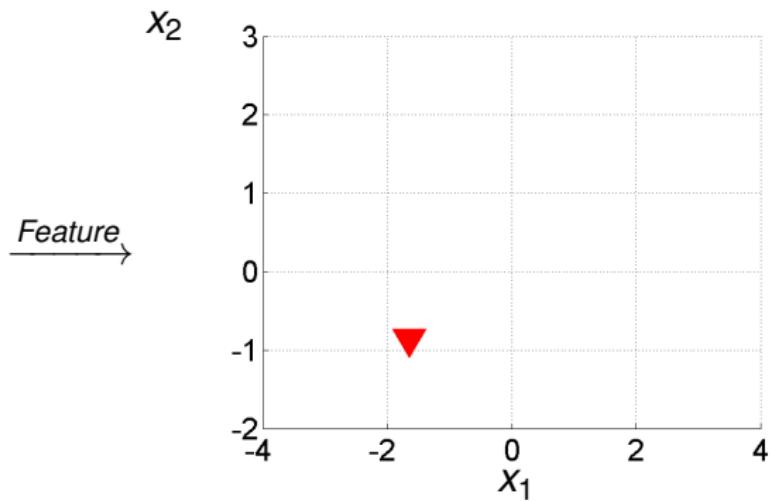
Binary Classification

How to find $y \in \{\text{Yes}, \text{No}\}$ given input data x ?



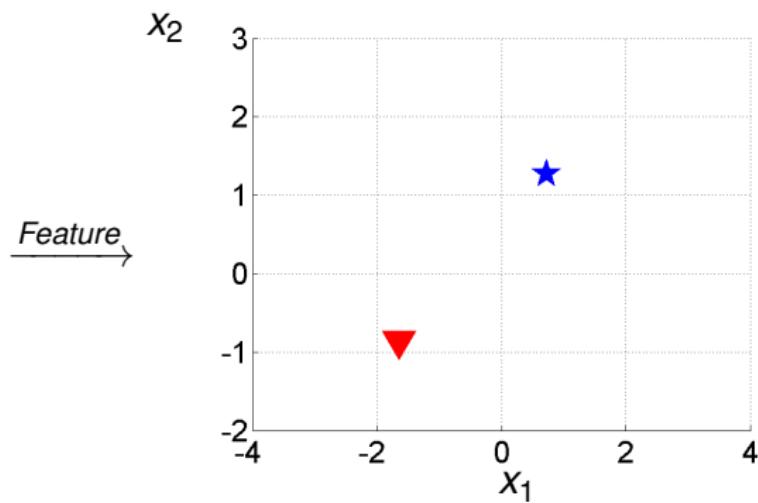
Binary Classification

How to find $y \in \{\text{Yes}, \text{No}\}$ given input data x ?



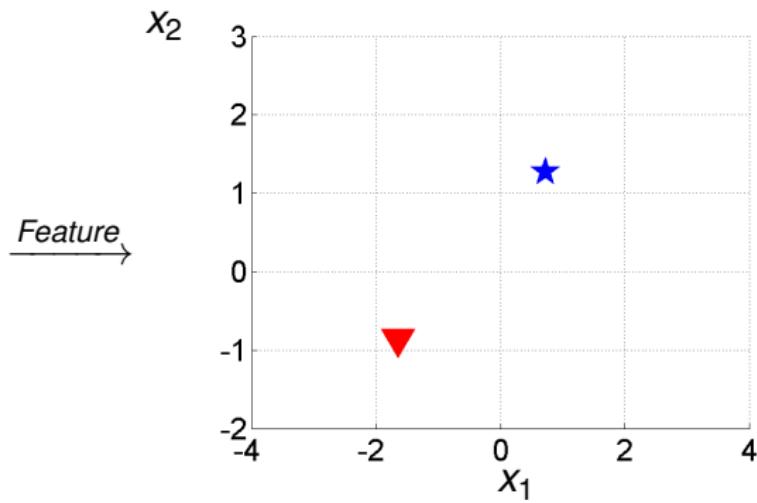
Binary Classification

How to find $y \in \{\text{Yes}, \text{No}\}$ given input data x ?



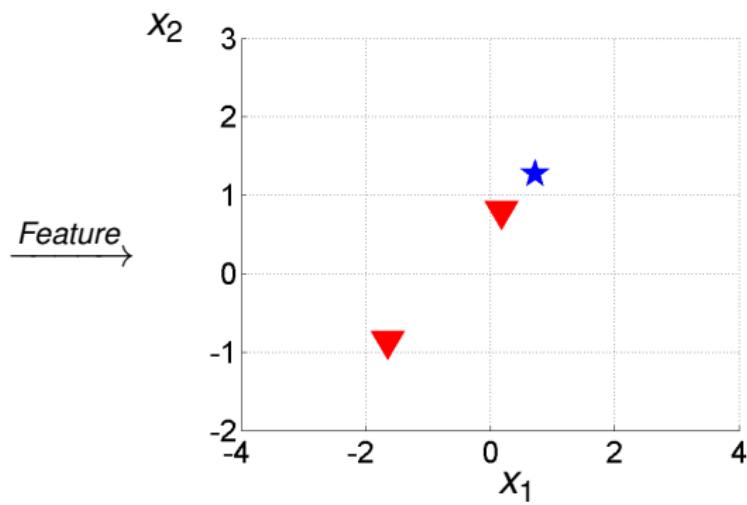
Binary Classification

How to find $y \in \{\text{Yes}, \text{No}\}$ given input data x ?



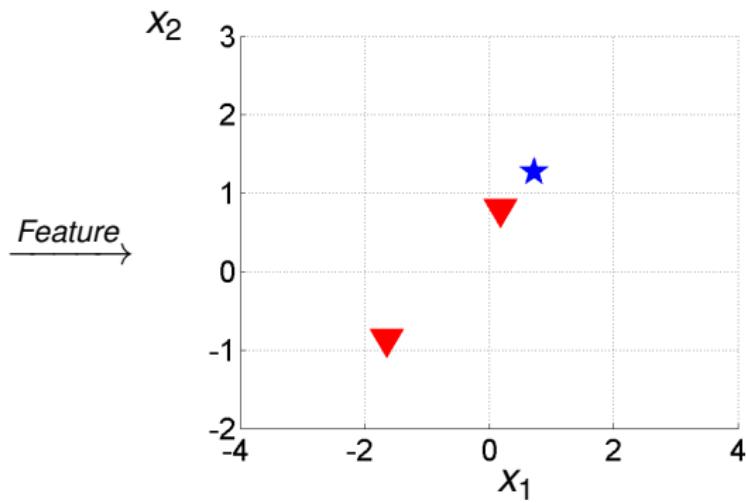
Binary Classification

How to find $y \in \{\text{Yes}, \text{No}\}$ given input data x ?



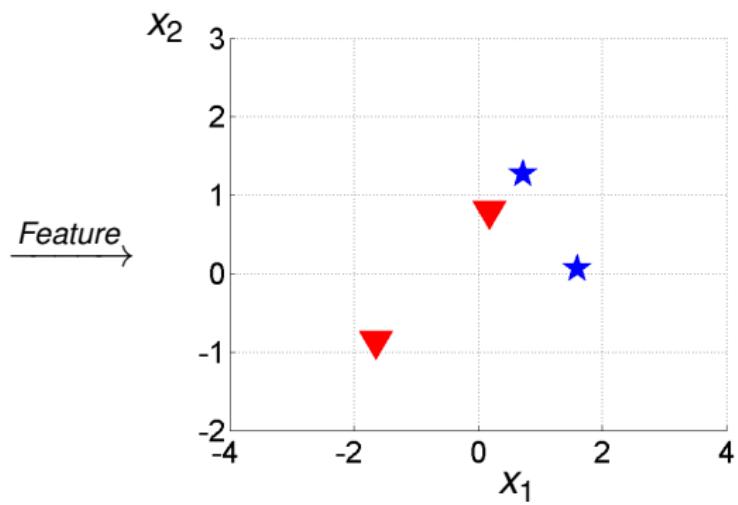
Binary Classification

How to find $y \in \{\text{Yes}, \text{No}\}$ given input data x ?



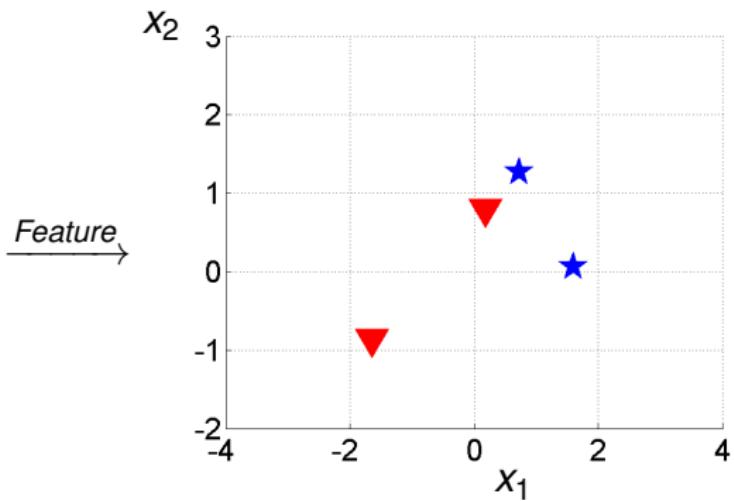
Binary Classification

How to find $y \in \{\text{Yes}, \text{No}\}$ given input data x ?



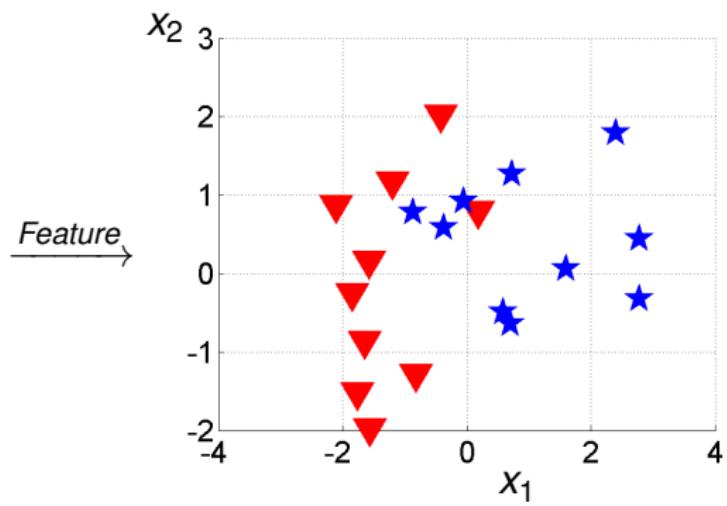
Binary Classification

How to find $y \in \{\text{Yes}, \text{No}\}$ given input data x ?



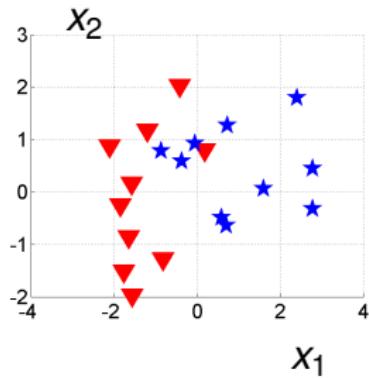
Binary Classification

How to find $y \in \{\text{Yes}, \text{No}\}$ given input data x ?



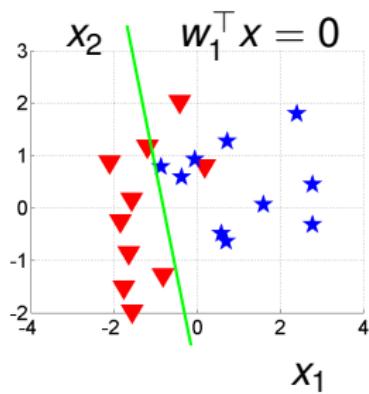
Linear Discriminant

$$y \in \{-1, 1\}$$



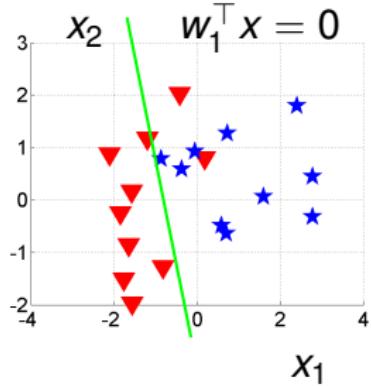
Linear Discriminant

$$y \in \{-1, 1\}$$



Linear Discriminant

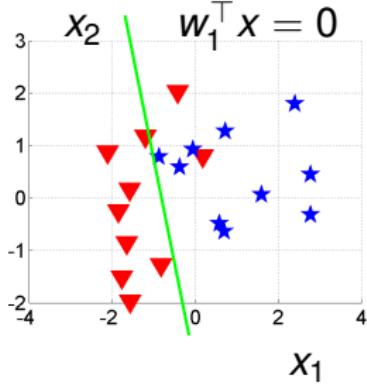
$$y \in \{-1, 1\}$$



$$y^* = \text{sign}(w_1^\top x)$$

Linear Discriminant

$$y \in \{-1, 1\}$$



$$y^* = \text{sign}(w_1^\top x)$$

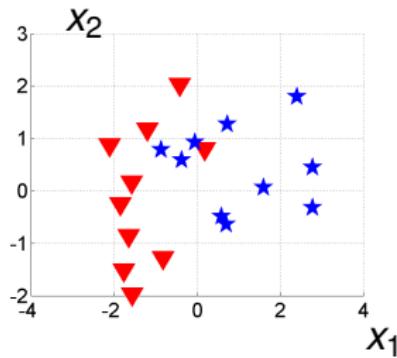
Equivalent:

$$y \in \{1, 2\}$$

$$y^* = \arg \max_{\hat{y}} \begin{bmatrix} w_1 \\ -w_1 \end{bmatrix}^\top \begin{bmatrix} x \delta(\hat{y} = 1) \\ x \delta(\hat{y} = 2) \end{bmatrix}$$

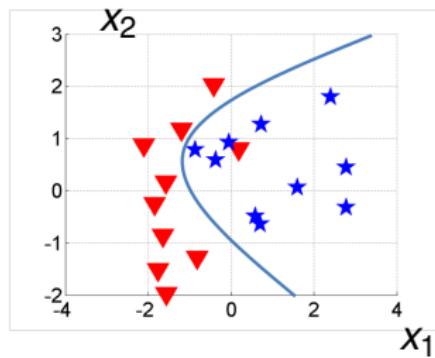
Non-linear Discriminant

$$y \in \{1, 2\}$$



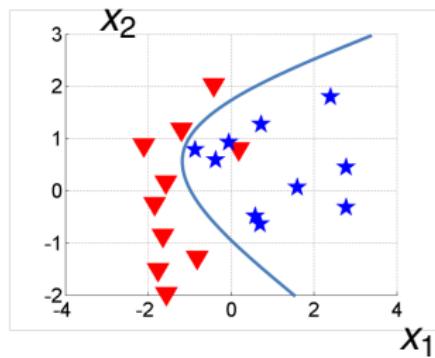
Non-linear Discriminant

$$y \in \{1, 2\}$$



Non-linear Discriminant

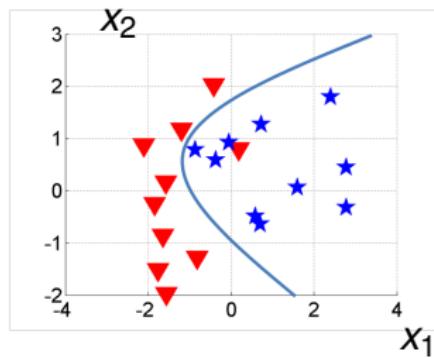
$$y \in \{1, 2\}$$



$$y^* = \arg \max_{\hat{y}} w^\top \phi(x, \hat{y})$$

Non-linear Discriminant

$$y \in \{1, 2\}$$



$$y^* = \arg \max_{\hat{y}} w^\top \phi(x, \hat{y})$$

More generally:

$$y^* = \arg \max_{\hat{y}} F(\hat{y}, x, w)$$

Multiclass Classification

Prediction answers a categorical question

$$x \rightarrow y \in \{1, \dots, K\}$$

Multiclass Classification

Prediction answers a categorical question

- Object classification ([Which object is in the image?](#))

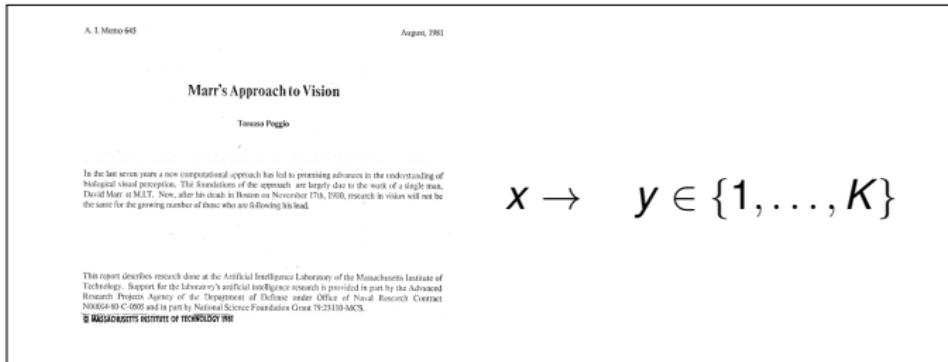


$$x \rightarrow y \in \{1, \dots, K\}$$

Multiclass Classification

Prediction answers a categorical question

- Object classification ([Which object is in the image?](#))
- Document retrieval ([What topic is this document about?](#))



Multiclass Classification

Prediction answers a categorical question

- Object classification ([Which object is in the image?](#))
- Document retrieval ([What topic is this document about?](#))
- Medical testing ([Which disease fits the symptoms?](#))


$$x \rightarrow y \in \{1, \dots, K\}$$

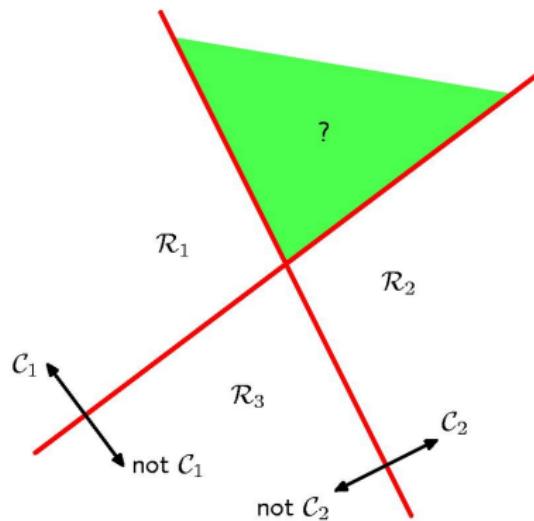
1 vs all

How to deal with more than two classes?

1 vs all

How to deal with more than two classes?

- Use $K - 1$ classifiers, each solving a two class problem of separating a point in class C_k from points not in the class.
- Known as **1 vs all** or **1 vs the rest** classifier



- Issue: more than one good answer

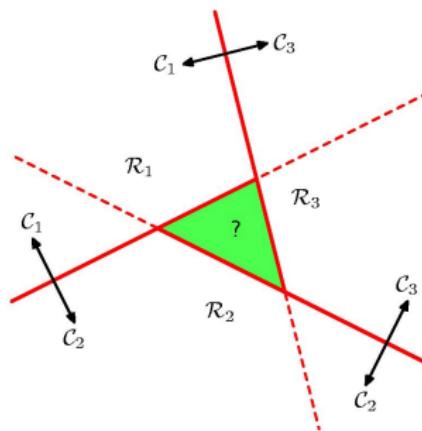
1 vs 1 classifier

How to deal with more than two classes?

1 vs 1 classifier

How to deal with more than two classes?

- Introduce $K(K - 1)/2$ two-way classifiers, one for each possible pair of classes
- Each point is classified according to majority vote amongst the discriminant function
- Known as the **1 vs 1 classifier**



- Issue: two-way preferences need not be transitive

Score Maximization

$$y \in \{1, 2, \dots, K\}$$

Prediction/Inference:

$$y^* = \arg \max_{\hat{y}} \begin{bmatrix} w_1 \\ \vdots \\ w_K \end{bmatrix}^\top \begin{bmatrix} x\delta(\hat{y} = 1) \\ \vdots \\ x\delta(\hat{y} = K) \end{bmatrix}$$

More Generally:

Score Maximization

$$y \in \{1, 2, \dots, K\}$$

Prediction/Inference:

$$y^* = \arg \max_{\hat{y}} \begin{bmatrix} w_1 \\ \vdots \\ w_K \end{bmatrix}^\top \begin{bmatrix} x\delta(\hat{y} = 1) \\ \vdots \\ x\delta(\hat{y} = K) \end{bmatrix}$$

More Generally:

$$y^* = \arg \max_{\hat{y}} w^\top \phi(x, \hat{y})$$

Even more generally:

Score Maximization

$$y \in \{1, 2, \dots, K\}$$

Prediction/Inference:

$$y^* = \arg \max_{\hat{y}} \begin{bmatrix} w_1 \\ \vdots \\ w_K \end{bmatrix}^\top \begin{bmatrix} x\delta(\hat{y} = 1) \\ \vdots \\ x\delta(\hat{y} = K) \end{bmatrix}$$

More Generally:

$$y^* = \arg \max_{\hat{y}} w^\top \phi(x, \hat{y})$$

Even more generally:

$$y^* = \arg \max_{\hat{y}} F(\hat{y}, x, w)$$

Score Maximization

$$y \in \{1, 2, \dots, K\}$$

Prediction/Inference:

$$y^* = \arg \max_{\hat{y}} \begin{bmatrix} w_1 \\ \vdots \\ w_K \end{bmatrix}^\top \begin{bmatrix} x\delta(\hat{y} = 1) \\ \vdots \\ x\delta(\hat{y} = K) \end{bmatrix}$$

More Generally:

$$y^* = \arg \max_{\hat{y}} w^\top \phi(x, \hat{y})$$

Even more generally:

$$y^* = \arg \max_{\hat{y}} F(\hat{y}, x, w)$$

Try all possible classes and return the highest scoring one.

Structured Prediction

Structured Prediction



Structured Prediction



V

Structured Prediction



Structured Prediction



Structured Prediction



Structured Prediction

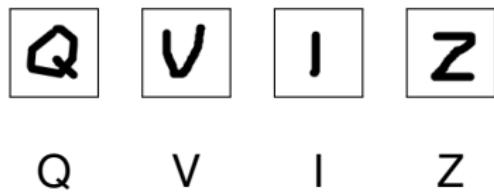


Z

Structured Prediction



Structured Prediction



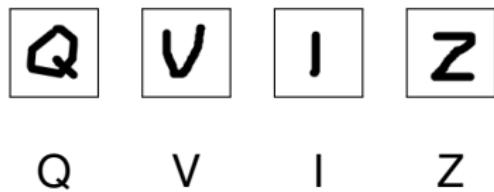
Q

V

I

Z

Structured Prediction



Q V I Z

Structured Prediction

   
Q V I Z

   
Q U I Z

Structured Prediction



Q V I Z



Q U I Z

Relationship not explicitly taken into account.

Structured Prediction

Example: Disparity map estimation

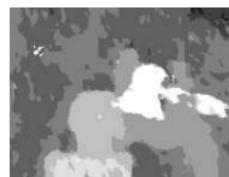
Why not to predict every variable separately:



Image



Independent Prediction



Structured Prediction

Structured Prediction

Prediction estimates a complex object

$$x \rightarrow \mathbf{y} = (y_1, \dots, y_D)$$

Structured Prediction

Prediction estimates a complex object

- Image segmentation ([estimate a labeling](#))

 x \rightarrow

$$\mathbf{y} = (y_1, \dots, y_D)$$



Structured Prediction

Prediction estimates a complex object

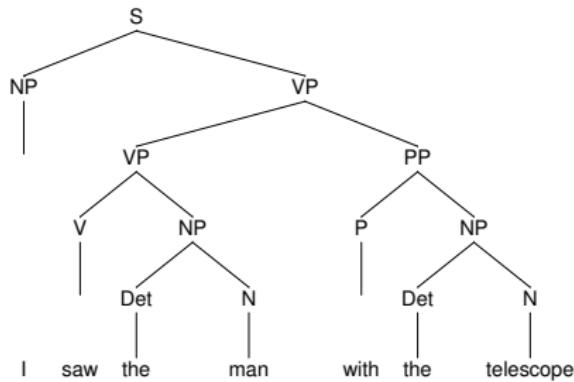
- Image segmentation (estimate a labeling)
- Sentence parsing (estimate a parse tree)

x

\rightarrow

$$\mathbf{y} = (y_1, \dots, y_D)$$

I saw the man with
the telescope.



Structured Prediction

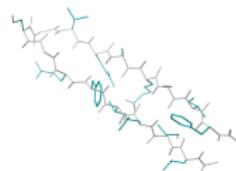
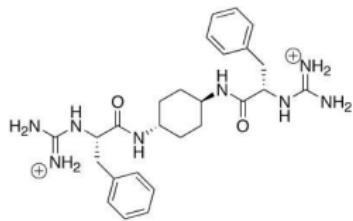
Prediction estimates a complex object

- Image segmentation ([estimate a labeling](#))
- Sentence parsing ([estimate a parse tree](#))
- Protein folding ([estimate a protein structure](#))

x

\rightarrow

$y = (y_1, \dots, y_D)$



Structured Prediction

Prediction estimates a complex object

- Image segmentation ([estimate a labeling](#))
- Sentence parsing ([estimate a parse tree](#))
- Protein folding ([estimate a protein structure](#))
- Stereo vision ([estimate a disparity map](#))

x

\rightarrow

$\mathbf{y} = (y_1, \dots, y_D)$



Structured Prediction

- “Standard” Prediction: output $y \in \mathcal{Y}$ is a scalar number

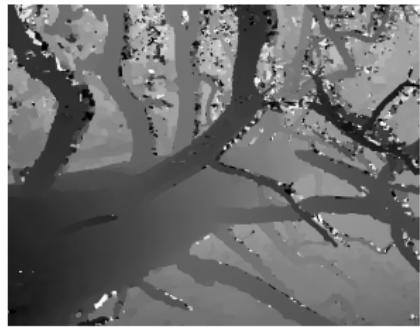
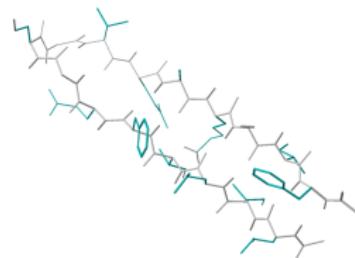
$$\mathcal{Y} = \{1, \dots, K\} \quad \text{or} \quad \mathcal{Y} = \mathbb{R}$$

Structured Prediction

- “Standard” Prediction: output $y \in \mathcal{Y}$ is a scalar number

$$\mathcal{Y} = \{1, \dots, K\} \quad \text{or} \quad \mathcal{Y} = \mathbb{R}$$

- “Structured” Prediction: output \mathbf{y} is a structured output:



Structured Prediction

Formally:

$$\mathbf{y} = (y_1, \dots, y_D) \quad y_i = \{1, \dots, K\}$$

Structured Prediction

Formally:

$$\mathbf{y} = (y_1, \dots, y_D) \quad y_i = \{1, \dots, K\}$$

Inference/Prediction:

$$\mathbf{y}^* = \arg \max_{\hat{\mathbf{y}}} F(\hat{y}_1, \dots, \hat{y}_D, x, w)$$

How many possibilities do we have to store and explore?

Structured Prediction

Formally:

$$\mathbf{y} = (y_1, \dots, y_D) \quad y_i = \{1, \dots, K\}$$

Inference/Prediction:

$$\mathbf{y}^* = \arg \max_{\hat{\mathbf{y}}} F(\hat{y}_1, \dots, \hat{y}_D, x, w)$$

How many possibilities do we have to store and explore?

$$K^D$$

Structured Prediction

Formally:

$$\mathbf{y} = (y_1, \dots, y_D) \quad y_i = \{1, \dots, K\}$$

Inference/Prediction:

$$\mathbf{y}^* = \arg \max_{\hat{\mathbf{y}}} F(\hat{y}_1, \dots, \hat{y}_D, x, w)$$

How many possibilities do we have to store and explore?

$$K^D$$

That's a problem. What can we do?

Structured Prediction

Separate prediction:

$$\max_{\mathbf{y}} F(y_1, \dots, y_D, \mathbf{x}, \mathbf{w}) = \sum_{i=1}^D \max_{y_i} f_i(y_i, \mathbf{x}, \mathbf{w})$$

Why not predict every variable y_i from $\mathbf{y} = (y_1, \dots, y_D)$ separately?

Structured Prediction

Separate prediction:

$$\max_{\mathbf{y}} F(y_1, \dots, y_D, \mathbf{x}, \mathbf{w}) = \sum_{i=1}^D \max_{y_i} f_i(y_i, \mathbf{x}, \mathbf{w})$$

Why not predict every variable y_i from $\mathbf{y} = (y_1, \dots, y_D)$ separately?

Relationship not explicitly taken into account.

Structured Prediction

Discriminant function decomposes:

$$F(y_1, \dots, y_D, x, w) = \sum_r f_r(\mathbf{y}_r, x, w)$$

Restriction: every $r \subseteq \{1, \dots, D\}$

Structured Prediction

Discriminant function decomposes:

$$F(y_1, \dots, y_D, x, w) = \sum_r f_r(\mathbf{y}_r, x, w)$$

Restriction: every $r \subseteq \{1, \dots, D\}$

Discrete domain:

$$f_{\{1,2\}}(\mathbf{y}_{\{1,2\}}) = f_{\{1,2\}}(y_1, y_2) = [f_{\{1,2\}}(1, 1), f_{\{1,2\}}(1, 2), \dots]$$

Structured Prediction

Discriminant function decomposes:

$$F(y_1, \dots, y_D, x, w) = \sum_r f_r(\mathbf{y}_r, x, w)$$

Restriction: every $r \subseteq \{1, \dots, D\}$

Discrete domain:

$$f_{\{1,2\}}(\mathbf{y}_{\{1,2\}}) = f_{\{1,2\}}(y_1, y_2) = [f_{\{1,2\}}(1, 1), f_{\{1,2\}}(1, 2), \dots]$$

	Q	U	I	Z	V
Q	0	0.8	0.2	0.1	0.1
U			
I	:				
Z					
V					

Structured Prediction

Q

V

I

Z

Q

U

I

Z

Example:

$$\begin{aligned}F(y_1, \dots, y_4, x, w) = & f_1(y_1, x, w) + f_2(y_2, x, w) + f_3(y_3, x, w) + f_4(y_4, x, w) \\& + f_{1,2}(y_1, y_2, x, w) + f_{2,3}(y_2, y_3, x, w) + f_{3,4}(y_3, y_4, x, w)\end{aligned}$$

Structured Prediction

Q

V

I

Z

Q

U

I

Z

Example:

$$\begin{aligned}F(y_1, \dots, y_4, x, w) = & f_1(y_1, x, w) + f_2(y_2, x, w) + f_3(y_3, x, w) + f_4(y_4, x, w) \\& + f_{1,2}(y_1, y_2, x, w) + f_{2,3}(y_2, y_3, x, w) + f_{3,4}(y_3, y_4, x, w)\end{aligned}$$

How many function values need to be stored if $y_i \in \{1, \dots, 26\} \forall i$?

Structured Prediction

Q

V

I

Z

Q

U

I

Z

Example:

$$F(y_1, \dots, y_4, x, w) = f_1(y_1, x, w) + f_2(y_2, x, w) + f_3(y_3, x, w) + f_4(y_4, x, w) \\ + f_{1,2}(y_1, y_2, x, w) + f_{2,3}(y_2, y_3, x, w) + f_{3,4}(y_3, y_4, x, w)$$

How many function values need to be stored if $y_i \in \{1, \dots, 26\} \forall i$?

$$26^4 \quad \text{v.s.} \quad 3 \cdot 26^2 (+4 \cdot 26)$$

Structured Prediction

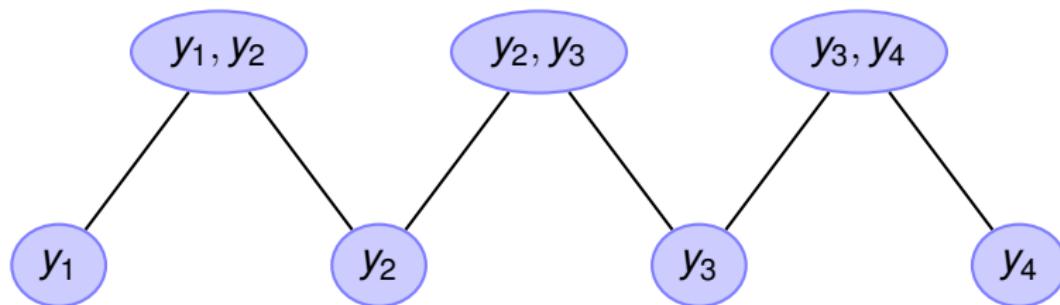
Visualization of the decomposition:

$$\begin{aligned} F(y_1, \dots, y_4, x, w) = & f_1(y_1, x, w) + f_2(y_2, x, w) + f_3(y_3, x, w) + f_4(y_4, x, w) \\ & + f_{1,2}(y_1, y_2, x, w) + f_{2,3}(y_2, y_3, x, w) + f_{3,4}(y_3, y_4, x, w) \end{aligned}$$

Structured Prediction

Visualization of the decomposition:

$$\begin{aligned} F(y_1, \dots, y_4, x, w) = & f_1(y_1, x, w) + f_2(y_2, x, w) + f_3(y_3, x, w) + f_4(y_4, x, w) \\ & + f_{1,2}(y_1, y_2, x, w) + f_{2,3}(y_2, y_3, x, w) + f_{3,4}(y_3, y_4, x, w) \end{aligned}$$



Edges denote subset relationship

Special cases

Predicting every variable separately:

$$F(y_1, \dots, y_D, x, w) = \sum_{i=1}^D f_i(y_i, x, w)$$

Structured Prediction

Special cases

Predicting every variable separately:

$$F(y_1, \dots, y_D, x, w) = \sum_{i=1}^D f_i(y_i, x, w)$$



Structured Prediction

Special cases

Predicting every variable separately:

$$F(y_1, \dots, y_D, x, w) = \sum_{i=1}^D f_i(y_i, x, w)$$



Markov random field with only unary variables

Structured Prediction

Special cases

Predicting every variable separately:

$$F(y_1, \dots, y_D, x, w) = \sum_{i=1}^D f_i(y_i, x, w)$$



Markov random field with only unary variables

Multi-variate prediction:

$$F(y_1, \dots, y_D, x, w) = f_{1,\dots,D}(\mathbf{y}_{1,\dots,D}, x, w)$$

Structured Prediction

Special cases

Predicting every variable separately:

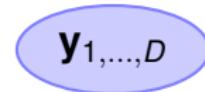
$$F(y_1, \dots, y_D, x, w) = \sum_{i=1}^D f_i(y_i, x, w)$$



Markov random field with only unary variables

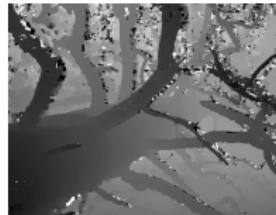
Multi-variate prediction:

$$F(y_1, \dots, y_D, x, w) = f_{1,\dots,D}(\mathbf{y}_{1,\dots,D}, x, w)$$

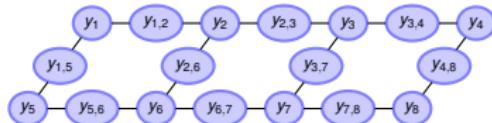


Structured Prediction

Example: stereo vision

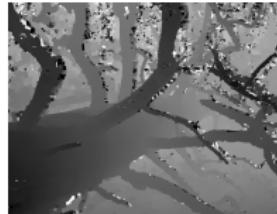


Markov/Conditional random field:

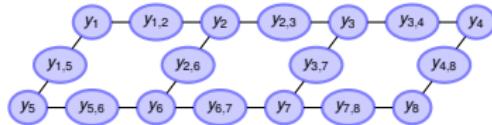


Structured Prediction

Example: stereo vision



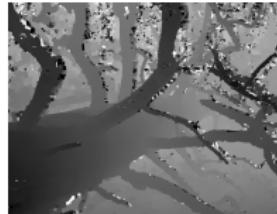
Markov/Conditional random field:



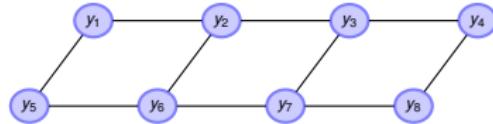
$$F(\mathbf{y}, \mathbf{x}, \mathbf{w}) = \sum_{i=1}^D f_i(y_i, \mathbf{x}, \mathbf{w}) + \sum_{i,j} f_{i,j}(y_i, y_j, \mathbf{x}, \mathbf{w})$$

Structured Prediction

Example: stereo vision



Markov/Conditional random field:

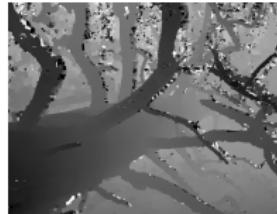


$$F(\mathbf{y}, \mathbf{x}, \mathbf{w}) = \sum_{i=1}^D f_i(y_i, \mathbf{x}, \mathbf{w}) + \sum_{i,j} f_{i,j}(y_i, y_j, \mathbf{x}, \mathbf{w})$$

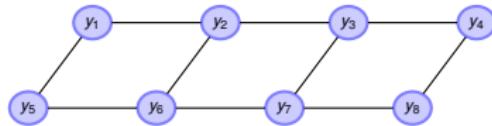
- Unary term:
- Pairwise term:

Structured Prediction

Example: stereo vision



Markov/Conditional random field:

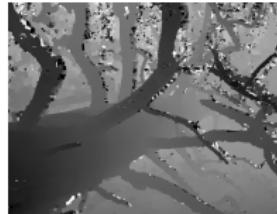


$$F(\mathbf{y}, \mathbf{x}, \mathbf{w}) = \sum_{i=1}^D f_i(y_i, \mathbf{x}, \mathbf{w}) + \sum_{i,j} f_{i,j}(y_i, y_j, \mathbf{x}, \mathbf{w})$$

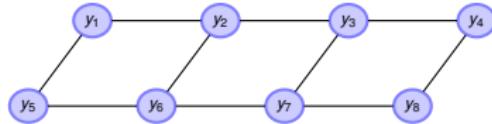
- Unary term: image evidence
- Pairwise term:

Structured Prediction

Example: stereo vision



Markov/Conditional random field:



$$F(\mathbf{y}, \mathbf{x}, \mathbf{w}) = \sum_{i=1}^D f_i(y_i, \mathbf{x}, \mathbf{w}) + \sum_{i,j} f_{i,j}(y_i, y_j, \mathbf{x}, \mathbf{w})$$

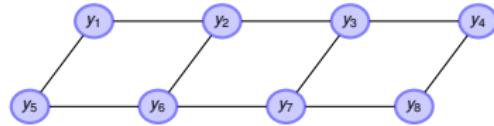
- Unary term: image evidence
- Pairwise term: smoothness prior

Structured Prediction

Example: semantic segmentation



Markov/Conditional random field:



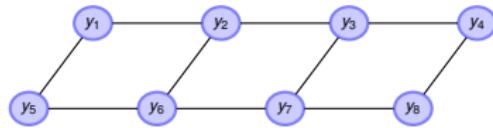
$$F(\mathbf{y}, \mathbf{x}, \mathbf{w}) = \sum_{i=1}^D f_i(y_i, \mathbf{x}, \mathbf{w}) + \sum_{i,j} f_{i,j}(y_i, y_j, \mathbf{x}, \mathbf{w})$$

Structured Prediction

Example: semantic segmentation



Markov/Conditional random field:



$$F(\mathbf{y}, \mathbf{x}, \mathbf{w}) = \sum_{i=1}^D f_i(y_i, \mathbf{x}, \mathbf{w}) + \sum_{i,j} f_{i,j}(y_i, y_j, \mathbf{x}, \mathbf{w})$$

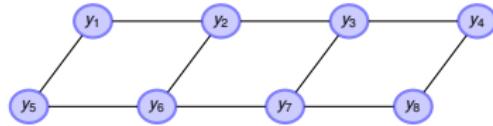
- Unary term:
- Pairwise term:

Structured Prediction

Example: semantic segmentation



Markov/Conditional random field:



$$F(\mathbf{y}, \mathbf{x}, \mathbf{w}) = \sum_{i=1}^D f_i(y_i, \mathbf{x}, \mathbf{w}) + \sum_{i,j} f_{i,j}(y_i, y_j, \mathbf{x}, \mathbf{w})$$

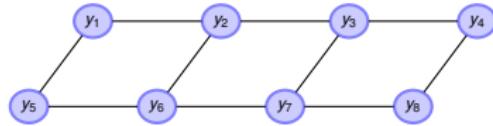
- Unary term: image evidence
- Pairwise term:

Structured Prediction

Example: semantic segmentation



Markov/Conditional random field:



$$F(\mathbf{y}, \mathbf{x}, \mathbf{w}) = \sum_{i=1}^D f_i(y_i, \mathbf{x}, \mathbf{w}) + \sum_{i,j} f_{i,j}(y_i, y_j, \mathbf{x}, \mathbf{w})$$

- Unary term: image evidence
- Pairwise term: smoothness prior

Structured Prediction

Inference:

$$\mathbf{y}^* = \arg \max_{\hat{\mathbf{y}}} F(\hat{\mathbf{y}}, x, w)$$

Structured Prediction

Inference:

$$\mathbf{y}^* = \arg \max_{\hat{\mathbf{y}}} F(\hat{\mathbf{y}}, x, w)$$

Probability of a configuration \mathbf{y} :

$$p(\mathbf{y} | x, w) = \frac{1}{Z(x, w)} \exp F(\mathbf{y}, x, w)$$

Structured Prediction

Inference:

$$\mathbf{y}^* = \arg \max_{\hat{\mathbf{y}}} F(\hat{\mathbf{y}}, x, w)$$

Probability of a configuration \mathbf{y} :

$$p(\mathbf{y} | x, w) = \frac{1}{Z(x, w)} \exp F(\mathbf{y}, x, w)$$

Normalization constant/[partition function](#):

$$Z(x, w) = \sum_{\mathbf{y}} \exp F(\mathbf{y}, x, w)$$

Structured Prediction

Inference:

$$\mathbf{y}^* = \arg \max_{\hat{\mathbf{y}}} F(\hat{\mathbf{y}}, x, w)$$

Probability of a configuration \mathbf{y} :

$$p(\mathbf{y} | x, w) = \frac{1}{Z(x, w)} \exp F(\mathbf{y}, x, w)$$

Normalization constant/[partition function](#):

$$Z(x, w) = \sum_{\mathbf{y}} \exp F(\mathbf{y}, x, w)$$

Inference as probability maximization:

$$\mathbf{y}^* = \arg \max_{\hat{\mathbf{y}}} p(\hat{\mathbf{y}} | x, w)$$

Structured Prediction

Inference:

$$\mathbf{y}^* = \arg \max_{\hat{\mathbf{y}}} \sum_r f_r(\hat{\mathbf{y}}_r, x, w)$$

Some inference algorithms:

Structured Prediction

Inference:

$$\mathbf{y}^* = \arg \max_{\hat{\mathbf{y}}} \sum_r f_r(\hat{\mathbf{y}}_r, x, w)$$

Some inference algorithms:

- Exhaustive search
- Dynamic programming
- Integer linear program
- Linear programming relaxation
- Message passing
- Graph-cut

Structured Prediction - Exhaustive Search

$$\mathbf{y}^* = \arg \max_{\hat{\mathbf{y}}} \sum_r f_r(\hat{\mathbf{y}}_r, x, w)$$

Structured Prediction - Exhaustive Search

$$\mathbf{y}^* = \arg \max_{\hat{\mathbf{y}}} \sum_r f_r(\hat{\mathbf{y}}_r, x, w)$$

Algorithm:

- try all possible configurations $\hat{\mathbf{y}} \in \mathcal{Y}$
- keep highest scoring element

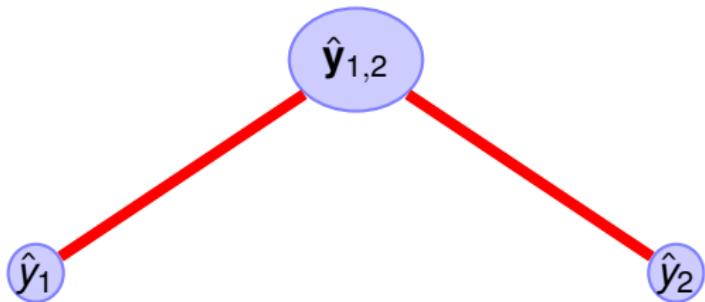
Structured Prediction - Exhaustive Search

$$\mathbf{y}^* = \arg \max_{\hat{\mathbf{y}}} \sum_r f_r(\hat{\mathbf{y}}_r, x, w)$$

Algorithm:

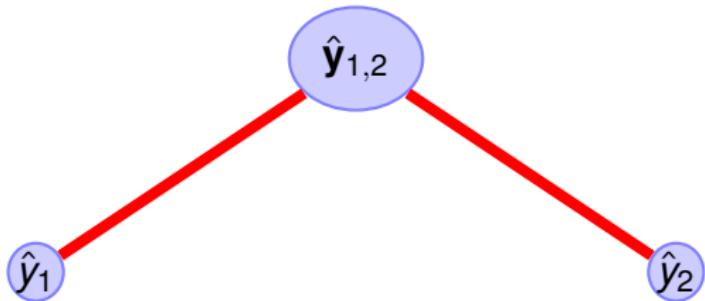
- try all possible configurations $\hat{\mathbf{y}} \in \mathcal{Y}$
 - keep highest scoring element
-
- **Advantage:** very simple to implement
 - **Disadvantage:** very slow for reasonably sized problems

Structured Prediction - Dynamic Programming



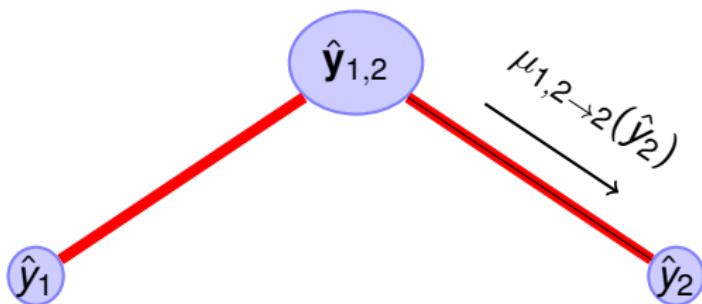
$$\max_{\hat{\mathbf{y}}} F(\hat{\mathbf{y}}, x, w) = \max_{\hat{y}_1, \hat{y}_2} f_2(\hat{y}_2) + f_1(\hat{y}_1) + f_{1,2}(\hat{y}_1, \hat{y}_2)$$

Structured Prediction - Dynamic Programming



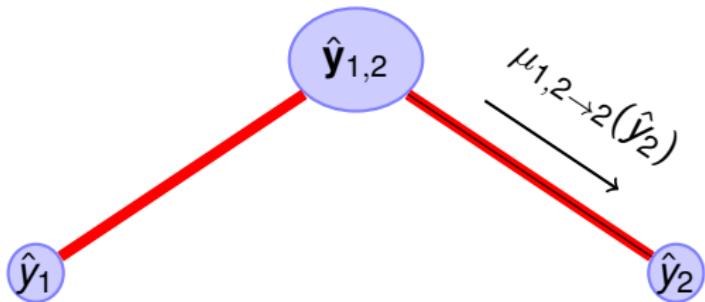
$$\begin{aligned}\max_{\hat{\mathbf{y}}} F(\hat{\mathbf{y}}, \mathbf{x}, \mathbf{w}) &= \max_{\hat{y}_1, \hat{y}_2} f_2(\hat{y}_2) + f_1(\hat{y}_1) + f_{1,2}(\hat{y}_1, \hat{y}_2) \\ &= \max_{\hat{y}_2} f_2(\hat{y}_2) + \max_{\hat{y}_1} \{f_1(\hat{y}_1) + f_{1,2}(\hat{y}_1, \hat{y}_2)\}\end{aligned}$$

Structured Prediction - Dynamic Programming



$$\begin{aligned}\max_{\hat{\mathbf{y}}} F(\hat{\mathbf{y}}, x, w) &= \max_{\hat{y}_1, \hat{y}_2} f_2(\hat{y}_2) + f_1(\hat{y}_1) + f_{1,2}(\hat{y}_1, \hat{y}_2) \\ &= \max_{\hat{y}_2} f_2(\hat{y}_2) + \underbrace{\max_{\hat{y}_1} \{ f_1(\hat{y}_1) + f_{1,2}(\hat{y}_1, \hat{y}_2) \}}_{\mu_{1,2 \rightarrow 2}(\hat{y}_2)}\end{aligned}$$

Structured Prediction - Dynamic Programming



$$\begin{aligned}\max_{\hat{\mathbf{y}}} F(\hat{\mathbf{y}}, x, w) &= \max_{\hat{y}_1, \hat{y}_2} f_2(\hat{y}_2) + f_1(\hat{y}_1) + f_{1,2}(\hat{y}_1, \hat{y}_2) \\ &= \max_{\hat{y}_2} f_2(\hat{y}_2) + \underbrace{\max_{\hat{y}_1} \{ f_1(\hat{y}_1) + f_{1,2}(\hat{y}_1, \hat{y}_2) \}}_{\mu_{1,2 \rightarrow 2}(\hat{y}_2)} \\ &= \max_{\hat{y}_2} f_2(\hat{y}_2) + \mu_{1,2 \rightarrow 2}(\hat{y}_2)\end{aligned}$$

Structured Prediction - Dynamic Programming

We can reorganize terms whenever the graph is a **tree**.

What to do for general loopy graphs?

Structured Prediction - Dynamic Programming

We can reorganize terms whenever the graph is a **tree**.

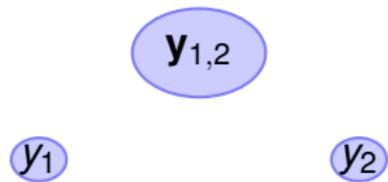
What to do for general loopy graphs?

- Dynamic programming extensions (message passing)
- Graph cut algorithms

Structured Prediction - Integer Linear Program

Example:

$$\mathbf{y}^* = \arg \max_{\hat{\mathbf{y}}} \sum_r f_r(\hat{\mathbf{y}}_r)$$



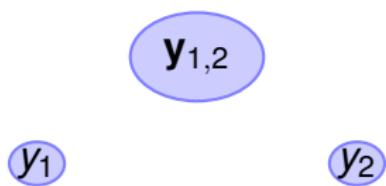
Integer Linear Program (LP) equivalence: variables $b_r(\mathbf{y}_r)$

$$\max_{b_1, b_2, b_{12}} \begin{bmatrix} b_1(1) \\ b_1(2) \\ b_2(1) \\ b_2(2) \\ b_{12}(1, 1) \\ b_{12}(2, 1) \\ b_{12}(1, 2) \\ b_{12}(2, 2) \end{bmatrix}^\top \begin{bmatrix} f_1(1) \\ f_1(2) \\ f_2(1) \\ f_2(2) \\ f_{12}(1, 1) \\ f_{12}(2, 1) \\ f_{12}(1, 2) \\ f_{12}(2, 2) \end{bmatrix}$$

Structured Prediction - Integer Linear Program

Example:

$$\mathbf{y}^* = \arg \max_{\hat{\mathbf{y}}} \sum_r f_r(\hat{\mathbf{y}}_r)$$



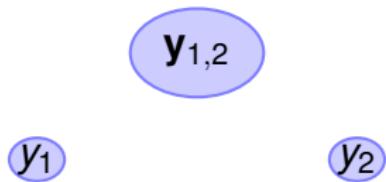
Integer Linear Program (LP) equivalence: variables $b_r(\mathbf{y}_r)$

$$\begin{aligned} & \max_{b_1, b_2, b_{12}} && \left[\begin{array}{c} b_1(1) \\ b_1(2) \\ b_2(1) \\ b_2(2) \\ b_{12}(1, 1) \\ b_{12}(2, 1) \\ b_{12}(1, 2) \\ b_{12}(2, 2) \end{array} \right]^\top \left[\begin{array}{c} f_1(1) \\ f_1(2) \\ f_2(1) \\ f_2(2) \\ f_{12}(1, 1) \\ f_{12}(2, 1) \\ f_{12}(1, 2) \\ f_{12}(2, 2) \end{array} \right] \\ & \text{s.t.} && b_r(\mathbf{y}_r) \in \{0, 1\} \end{aligned}$$

Structured Prediction - Integer Linear Program

Example:

$$\mathbf{y}^* = \arg \max_{\hat{\mathbf{y}}} \sum_r f_r(\hat{\mathbf{y}}_r)$$



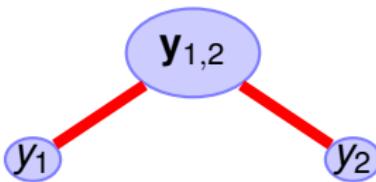
Integer Linear Program (LP) equivalence: variables $b_r(\mathbf{y}_r)$

$$\max_{b_1, b_2, b_{12}} \begin{bmatrix} b_1(1) \\ b_1(2) \\ b_2(1) \\ b_2(2) \\ b_{12}(1, 1) \\ b_{12}(2, 1) \\ b_{12}(1, 2) \\ b_{12}(2, 2) \end{bmatrix}^\top \begin{bmatrix} f_1(1) \\ f_1(2) \\ f_2(1) \\ f_2(2) \\ f_{12}(1, 1) \\ f_{12}(2, 1) \\ f_{12}(1, 2) \\ f_{12}(2, 2) \end{bmatrix} \quad \text{s.t.} \quad \begin{aligned} b_r(\mathbf{y}_r) &\in \{0, 1\} \\ \sum_{\mathbf{y}_r} b_r(\mathbf{y}_r) &= 1 \end{aligned}$$

Structured Prediction - Integer Linear Program

Example:

$$\mathbf{y}^* = \arg \max_{\hat{\mathbf{y}}} \sum_r f_r(\hat{\mathbf{y}}_r)$$



Integer Linear Program (LP) equivalence: variables $b_r(\mathbf{y}_r)$

$$\begin{aligned} \max_{b_1, b_2, b_{12}} & \left[\begin{array}{c} b_1(1) \\ b_1(2) \\ b_2(1) \\ b_2(2) \\ b_{12}(1, 1) \\ b_{12}(2, 1) \\ b_{12}(1, 2) \\ b_{12}(2, 2) \end{array} \right]^\top \left[\begin{array}{c} f_1(1) \\ f_1(2) \\ f_2(1) \\ f_2(2) \\ f_{12}(1, 1) \\ f_{12}(2, 1) \\ f_{12}(1, 2) \\ f_{12}(2, 2) \end{array} \right] \\ & \text{s.t. } \begin{aligned} b_r(\mathbf{y}_r) &\in \{0, 1\} \\ \sum_{\mathbf{y}_r} b_r(\mathbf{y}_r) &= 1 \\ \sum_{\mathbf{y}_p \setminus \mathbf{y}_r} b_p(\mathbf{y}_p) &= b_r(\mathbf{y}_r) \end{aligned} \end{aligned}$$

Structured Prediction - Integer Linear Program

$$\mathbf{y}^* = \arg \max_{\hat{\mathbf{y}}} \sum_r f_r(\hat{\mathbf{y}}_r)$$

Integer linear program:

$$\begin{aligned} \max_{b_1, b_2, b_{12}} & \left[\begin{array}{c} b_1(1) \\ b_1(2) \\ b_2(1) \\ b_2(2) \\ b_{12}(1, 1) \\ b_{12}(2, 1) \\ b_{12}(1, 2) \\ b_{12}(2, 2) \end{array} \right]^\top \left[\begin{array}{c} f_1(1) \\ f_1(2) \\ f_2(1) \\ f_2(2) \\ f_{12}(1, 1) \\ f_{12}(2, 1) \\ f_{12}(1, 2) \\ f_{12}(2, 2) \end{array} \right] \\ & \text{s.t. } \begin{aligned} b_r(\mathbf{y}_r) &\in \{0, 1\} \\ b_r(\mathbf{y}_r) &\geq 0 \\ \sum_{\mathbf{y}_r} b_r(\mathbf{y}_r) &= 1 \\ \sum_{\mathbf{y}_p \setminus \mathbf{y}_r} b_p(\mathbf{y}_p) &= b_r(\mathbf{y}_r) \end{aligned} \end{aligned}$$

Structured Prediction - Integer Linear Program

$$\mathbf{y}^* = \arg \max_{\hat{\mathbf{y}}} \sum_r f_r(\hat{\mathbf{y}}_r)$$

Integer linear program:

$$b_r(\mathbf{y}_r) \in \{0, 1\}$$

$$b_r(\mathbf{y}_r) \geq 0$$

$$\max_{b_r} \quad \sum_{r, \mathbf{y}_r} b_r(\mathbf{y}_r) f_r(\mathbf{y}_r) \quad \text{s.t.} \quad \sum_{\mathbf{y}_r} b_r(\mathbf{y}_r) = 1$$

$$\sum_{\mathbf{y}_p \setminus \mathbf{y}_r} b_p(\mathbf{y}_p) = b_r(\mathbf{y}_r)$$

Structured Prediction - Integer Linear Program

$$\mathbf{y}^* = \arg \max_{\hat{\mathbf{y}}} \sum_r f_r(\hat{\mathbf{y}}_r)$$

Integer linear program:

$$b_r(\mathbf{y}_r) \in \{0, 1\}$$

$$b_r(\mathbf{y}_r) \geq 0$$

$$\max_{b_r} \quad \sum_{r, \mathbf{y}_r} b_r(\mathbf{y}_r) f_r(\mathbf{y}_r) \quad \text{s.t.} \quad \sum_{\mathbf{y}_r} b_r(\mathbf{y}_r) = 1$$

Marginalization

Structured Prediction - Integer Linear Program

$$\mathbf{y}^* = \arg \max_{\hat{\mathbf{y}}} \sum_r f_r(\hat{\mathbf{y}}_r)$$

Integer linear program:

$$b_r(\mathbf{y}_r) \in \{0, 1\}$$

$$\max_{b_r} \quad \sum_{r, \mathbf{y}_r} b_r(\mathbf{y}_r) f_r(\mathbf{y}_r) \quad \text{s.t.} \quad \begin{array}{l} \text{Local probability } b_r \\ \text{Marginalization} \end{array}$$

Structured Prediction - Integer Linear Program

$$\mathbf{y}^* = \arg \max_{\hat{\mathbf{y}}} \sum_r f_r(\hat{\mathbf{y}}_r)$$

Integer linear program:

$$b_r(\mathbf{y}_r) \in \{0, 1\}$$

$$\max_{b_r} \quad \sum_{r, \mathbf{y}_r} b_r(\mathbf{y}_r) f_r(\mathbf{y}_r) \quad \text{s.t.} \quad \begin{array}{l} \text{Local probability } b_r \\ \text{Marginalization} \end{array}$$

- **Advantage:** very good solvers available
- **Disadvantage:** very slow for larger problems

Structured Prediction - Linear Programming Relaxation

$$\mathbf{y}^* = \arg \max_{\hat{\mathbf{y}}} \sum_r f_r(\hat{\mathbf{y}}_r)$$

LP relaxation:

$$\cancel{b_r(\mathbf{y}_r) \in \{0, 1\}}$$

$$\max_{b_r} \quad \sum_{r, \mathbf{y}_r} b_r(\mathbf{y}_r) f_r(\mathbf{y}_r) \quad \text{s.t.} \quad \begin{array}{l} \text{Local probability } b_r \\ \text{Marginalization} \end{array}$$

$$\underbrace{\quad}_{\text{s.t.}} \quad \underbrace{b \in \mathcal{C}}_{\text{b.c.}}$$

Structured Prediction - Linear Programming Relaxation

$$\mathbf{y}^* = \arg \max_{\hat{\mathbf{y}}} \sum_r f_r(\hat{\mathbf{y}}_r)$$

LP relaxation:

$$\underline{b_r(\mathbf{y}_r) \in \{0, 1\}}$$

$$\max_{b_r} \quad \sum_{r, \mathbf{y}_r} b_r(\mathbf{y}_r) f_r(\mathbf{y}_r) \quad \text{s.t. Local probability } b_r$$

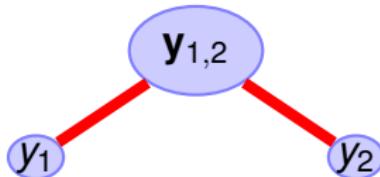
Marginalization

$$\underbrace{\quad}_{\text{s.t.}} \quad \underbrace{b \in \mathcal{C}}$$

- **Advantage:** very good solvers available
- **Disadvantage:** slow for larger problems

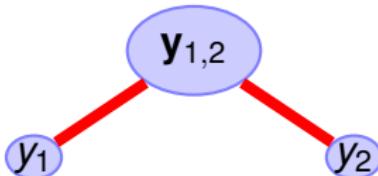
Structured Prediction - Message Passing

Graph structure defined via marginalization constraints



Structured Prediction - Message Passing

Graph structure defined via marginalization constraints



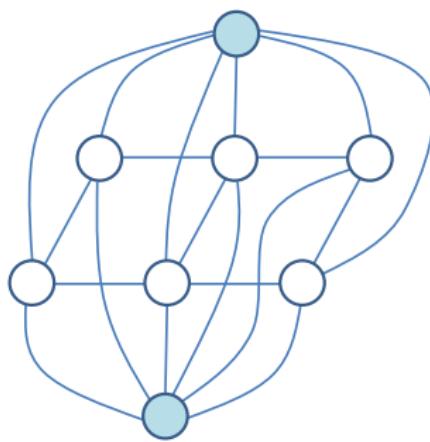
Message passing solvers:

- Advantage: Efficient due to analytically computable sub-problems
- Problem: Special care required to find global optimum

Structured Prediction - Graph-cut Solvers

For $y_i \in \{1, 2\}$:

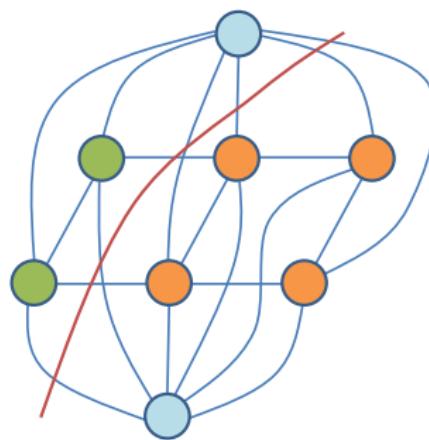
- Convert scoring function F into auxiliary graph



Structured Prediction - Graph-cut Solvers

For $y_i \in \{1, 2\}$:

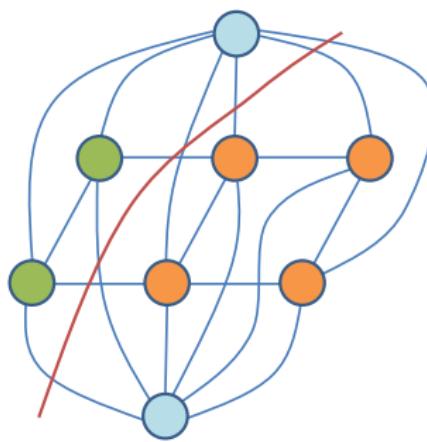
- Convert scoring function F into auxiliary graph
- Compute a weighted cut cost corresponding to the labeling score



Structured Prediction - Graph-cut Solvers

For $y_i \in \{1, 2\}$:

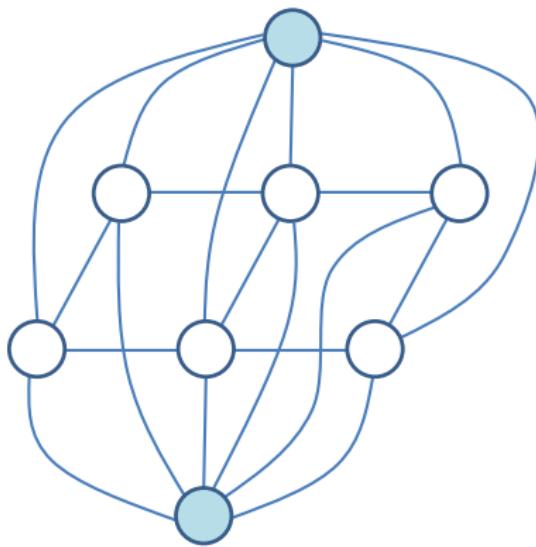
- Convert scoring function F into auxiliary graph
- Compute a weighted cut cost corresponding to the labeling score



What are the nodes and what are the weights on the edges?

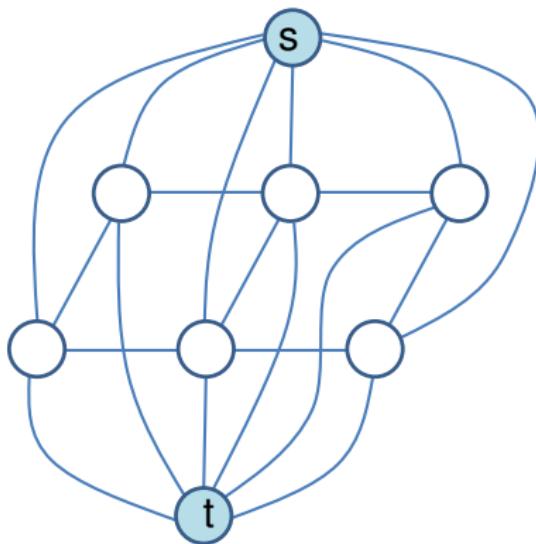
Structured Prediction - Graph-cut Solvers

What are the nodes?



Structured Prediction - Graph-cut Solvers

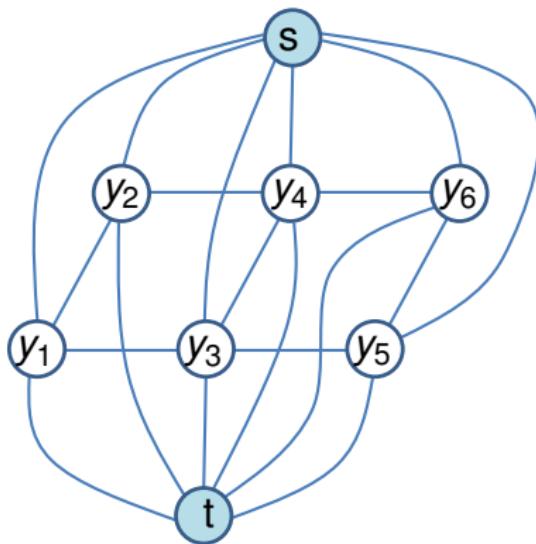
What are the nodes?



- Two special nodes called ‘source’ and ‘terminal’

Structured Prediction - Graph-cut Solvers

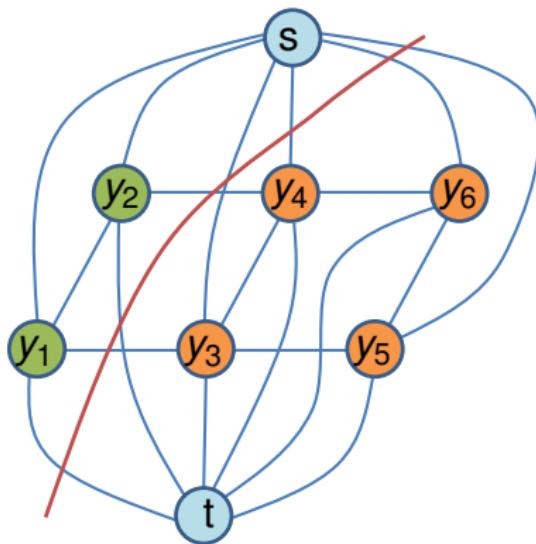
What are the nodes?



- Two special nodes called ‘source’ and ‘terminal’
- Variables y_i as nodes

Structured Prediction - Graph-cut Solvers

What are the nodes?



- Two special nodes called ‘source’ and ‘terminal’
- Variables y_i as nodes

Structured Prediction - Graph-cut Solvers

What weights do we assign to edges?

Structured Prediction - Graph-cut Solvers

What weights do we assign to edges?

Recall that local scoring functions are arrays:

$$[\ f_1(y_1 = 1) \quad f_1(y_1 = 2) \]$$

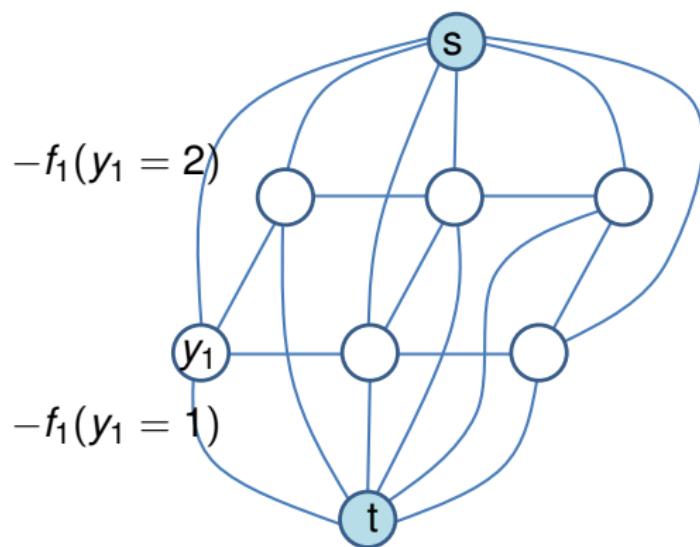
Structured Prediction - Graph-cut Solvers

What weights do we assign to edges?

Recall that local scoring functions are arrays:

$$[\ f_1(y_1 = 1) \quad f_1(y_1 = 2) \]$$

Graph-cut solvers compute a min-cut:



Structured Prediction - Graph-cut Solvers

What weights do we assign to edges?

Recall that local scoring functions are arrays:

$$\begin{bmatrix} f_{12}(1,1) & f_{12}(1,2) \\ f_{12}(2,1) & f_{12}(2,2) \end{bmatrix} =$$

Structured Prediction - Graph-cut Solvers

What weights do we assign to edges?

Recall that local scoring functions are arrays:

$$\begin{bmatrix} f_{12}(1,1) & f_{12}(1,2) \\ f_{12}(2,1) & f_{12}(2,2) \end{bmatrix} = f(1,1) - f(2,1) + f(2,2)$$
$$+ \begin{bmatrix} 0 & 0 \\ f(2,1) - f(1,1) & f(2,1) - f(1,1) \end{bmatrix}$$
$$+ \begin{bmatrix} f(2,1) - f(2,2) & 0 \\ f(2,1) - f(2,2) & 0 \end{bmatrix}$$
$$+ \begin{bmatrix} 0 & f(1,2) + f(2,1) - f(1,1) - f(2,2) \\ 0 & 0 \end{bmatrix}$$

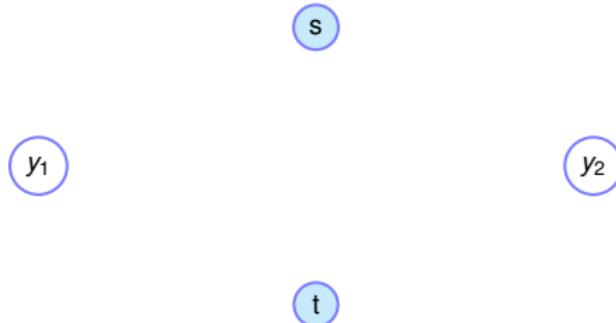
Structured Prediction - Graph-cut Solvers

What weights do we assign to edges?

Recall that local scoring functions are arrays:

$$\begin{bmatrix} f_{12}(1,1) & f_{12}(1,2) \\ f_{12}(2,1) & f_{12}(2,2) \end{bmatrix} = f(1,1) - f(2,1) + f(2,2)$$
$$+ \begin{bmatrix} 0 & 0 \\ f(2,1) - f(1,1) & f(2,1) - f(1,1) \end{bmatrix}$$
$$+ \begin{bmatrix} f(2,1) - f(2,2) & 0 \\ f(2,1) - f(2,2) & 0 \end{bmatrix}$$
$$+ \begin{bmatrix} 0 & f(1,2) + f(2,1) - f(1,1) - f(2,2) \\ 0 & 0 \end{bmatrix}$$

Graph-cut solvers compute a min-cut:



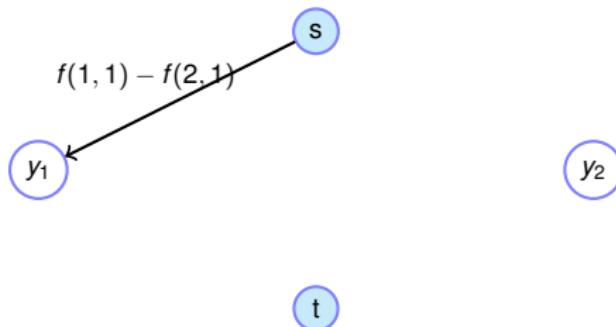
Structured Prediction - Graph-cut Solvers

What weights do we assign to edges?

Recall that local scoring functions are arrays:

$$\begin{bmatrix} f_{12}(1,1) & f_{12}(1,2) \\ f_{12}(2,1) & f_{12}(2,2) \end{bmatrix} = f(1,1) - f(2,1) + f(2,2)$$
$$+ \begin{bmatrix} 0 & 0 \\ f(2,1) - f(1,1) & f(2,1) - f(1,1) \end{bmatrix}$$
$$+ \begin{bmatrix} f(2,1) - f(2,2) & 0 \\ f(2,1) - f(2,2) & 0 \end{bmatrix}$$
$$+ \begin{bmatrix} 0 & f(1,2) + f(2,1) - f(1,1) - f(2,2) \\ 0 & 0 \end{bmatrix}$$

Graph-cut solvers compute a min-cut:



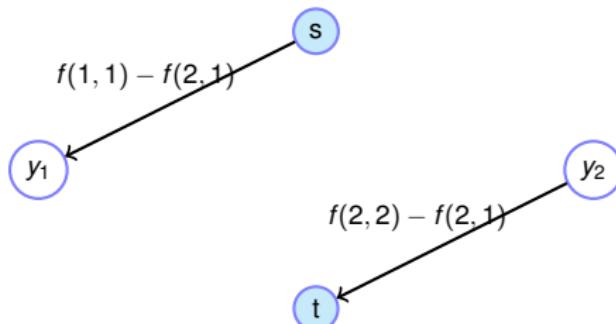
Structured Prediction - Graph-cut Solvers

What weights do we assign to edges?

Recall that local scoring functions are arrays:

$$\begin{bmatrix} f_{12}(1,1) & f_{12}(1,2) \\ f_{12}(2,1) & f_{12}(2,2) \end{bmatrix} = f(1,1) - f(2,1) + f(2,2) \\ + \begin{bmatrix} 0 & 0 \\ f(2,1) - f(1,1) & f(2,1) - f(1,1) \end{bmatrix} \\ + \begin{bmatrix} f(2,1) - f(2,2) & 0 \\ f(2,1) - f(2,2) & 0 \end{bmatrix} \\ + \begin{bmatrix} 0 & f(1,2) + f(2,1) - f(1,1) - f(2,2) \\ 0 & 0 \end{bmatrix}$$

Graph-cut solvers compute a min-cut:



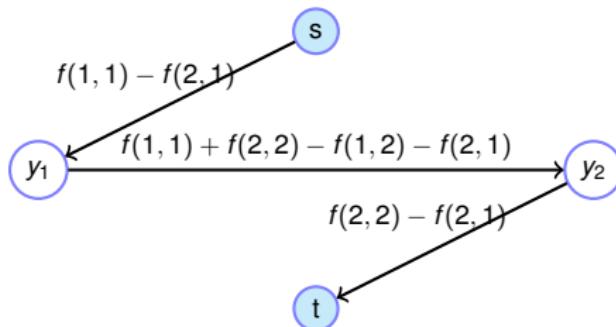
Structured Prediction - Graph-cut Solvers

What weights do we assign to edges?

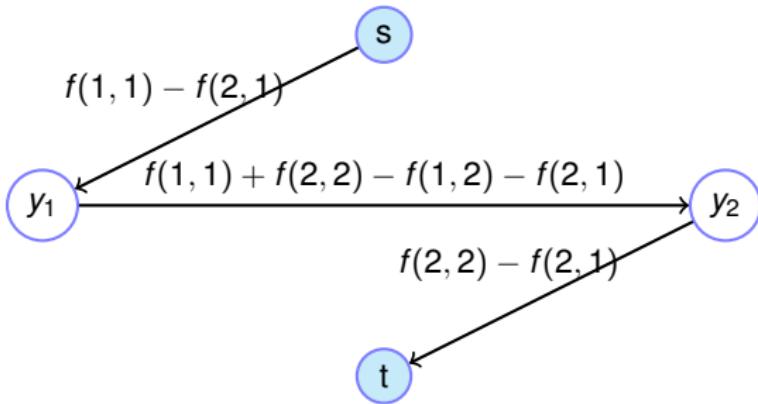
Recall that local scoring functions are arrays:

$$\begin{bmatrix} f_{12}(1,1) & f_{12}(1,2) \\ f_{12}(2,1) & f_{12}(2,2) \end{bmatrix} = f(1,1) - f(2,1) + f(2,2)$$
$$+ \begin{bmatrix} 0 & 0 \\ f(2,1) - f(1,1) & f(2,1) - f(1,1) \end{bmatrix}$$
$$+ \begin{bmatrix} f(2,1) - f(2,2) & 0 \\ f(2,1) - f(2,2) & 0 \end{bmatrix}$$
$$+ \begin{bmatrix} 0 & f(1,2) + f(2,1) - f(1,1) - f(2,2) \\ 0 & 0 \end{bmatrix}$$

Graph-cut solvers compute a min-cut:

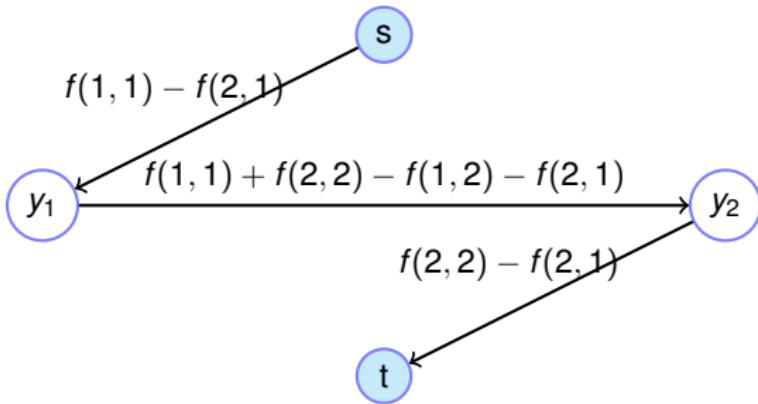


Structured Prediction - Graph-cut Solvers



Requirement for optimality:

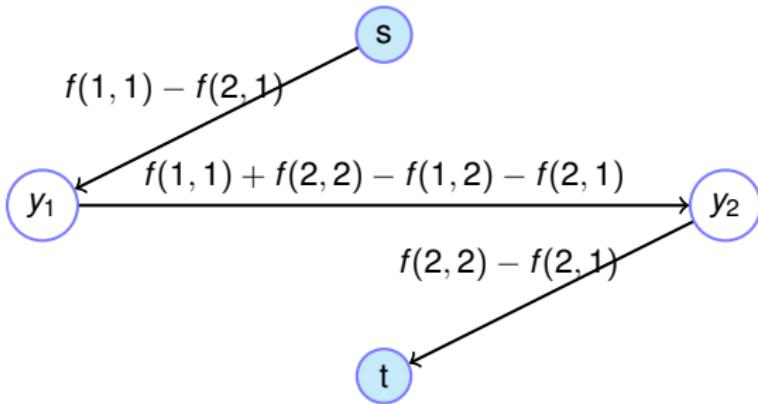
Structured Prediction - Graph-cut Solvers



Requirement for optimality: Pairwise edge weights are positive

$$f(1, 1) + f(2, 2) - f(1, 2) - f(2, 1) \geq 0 \quad \text{sub-modularity}$$

Structured Prediction - Graph-cut Solvers

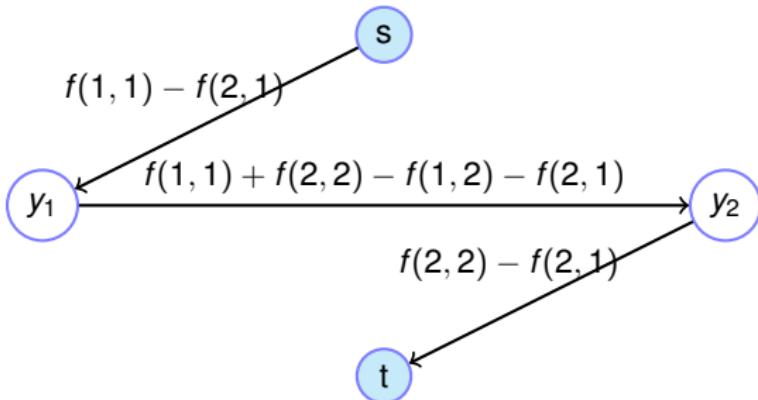


Requirement for optimality: Pairwise edge weights are positive

$$f(1, 1) + f(2, 2) - f(1, 2) - f(2, 1) \geq 0 \quad \text{sub-modularity}$$

For higher order functions?

Structured Prediction - Graph-cut Solvers

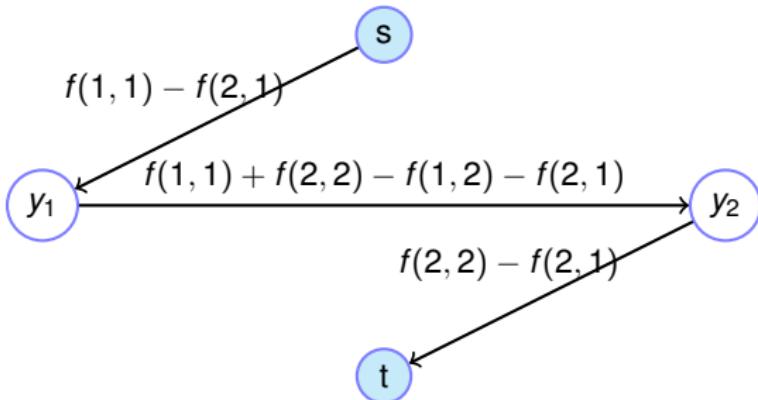


Requirement for optimality: Pairwise edge weights are positive

$$f(1,1) + f(2,2) - f(1,2) - f(2,1) \geq 0 \quad \text{sub-modularity}$$

For higher order functions? More complicated graph constructions

Structured Prediction - Graph-cut Solvers

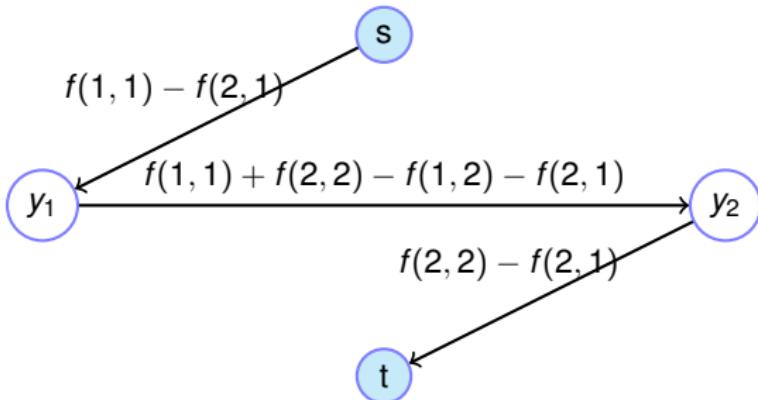


Requirement for optimality: Pairwise edge weights are positive

$$f(1,1) + f(2,2) - f(1,2) - f(2,1) \geq 0 \quad \text{sub-modularity}$$

For higher order functions? More complicated graph constructions
For more than two labels?

Structured Prediction - Graph-cut Solvers



Requirement for optimality: Pairwise edge weights are positive

$$f(1,1) + f(2,2) - f(1,2) - f(2,1) \geq 0 \quad \text{sub-modularity}$$

For higher order functions? More complicated graph constructions
For more than two labels? Move making algorithms

Structured Prediction

Inference:

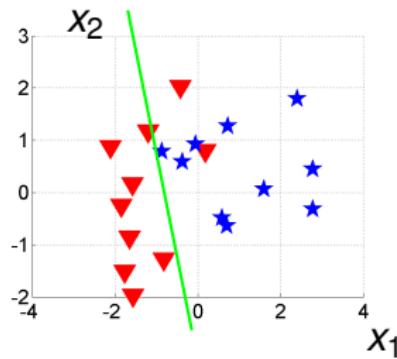
$$\mathbf{y}^* = \arg \max_{\hat{\mathbf{y}}} \sum_r f_r(\hat{\mathbf{y}}_r, x, w)$$

Efficiency and accuracy of inference algorithms is problem dependent:

- Exhaustive search
- Dynamic programming
- Integer linear program
- Linear programming relaxation
- Message passing
- Graph-cut

How to learn in structured spaces?

Binary SVM



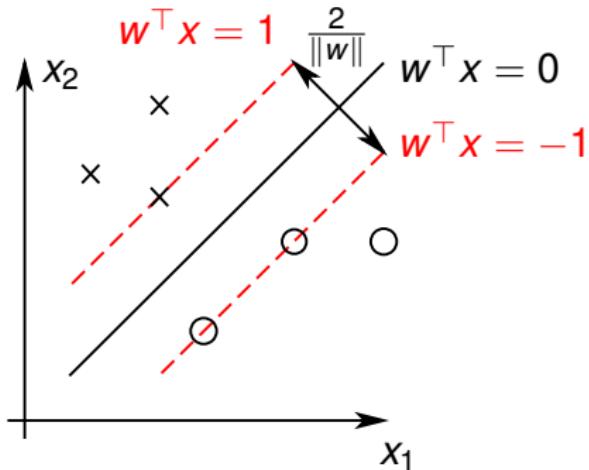
Given: dataset

$$\mathcal{D} = \{(x, y)\}, \quad N = |\mathcal{D}|$$

containing independent and identically drawn pairs

Binary SVM

Intuitively:

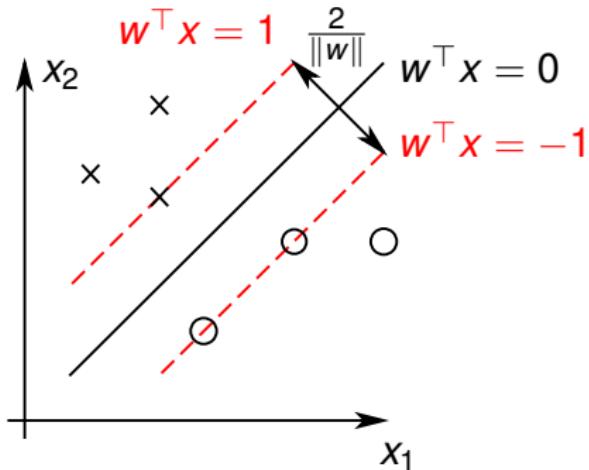


Maximize margin $\frac{2}{\|w\|}$:

$$\min_w \frac{1}{2} \|w\|_2^2 \quad \text{s.t.} \quad yw^\top x \geq 1 \quad \forall (x, y) \in \mathcal{D}$$

Binary SVM

Intuitively:



Maximize margin $\frac{2}{\|w\|}$:

$$\min_w \frac{1}{2} \|w\|_2^2 \quad \text{s.t.} \quad yw^\top x \geq 1 \quad \forall (x, y) \in \mathcal{D}$$

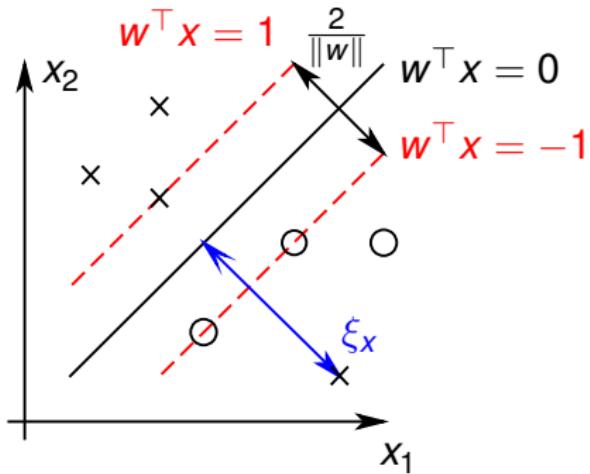
Issue: what if data not linearly separable?

Binary SVM

Introduce slack variables ξ :

$$\min_{w, \xi_x \geq 0} \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y) \in \mathcal{D}} \xi_x \quad \text{s.t.} \quad yw^\top x \geq 1 - \xi_x \quad \forall (x, y) \in \mathcal{D}$$

Intuitively:



Binary SVM

$$\min_{w, \xi_x \geq 0} \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y) \in \mathcal{D}} \xi_x \quad \text{s.t.} \quad yw^\top x \geq 1 - \xi_x \quad \forall (x,y) \in \mathcal{D}$$

Binary SVM

$$\min_{w, \xi_x \geq 0} \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y) \in \mathcal{D}} \xi_x \quad \text{s.t.} \quad yw^\top x \geq 1 - \xi_x \quad \forall (x,y) \in \mathcal{D}$$

Equivalent:

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y) \in \mathcal{D}} \max\{0, 1 - yw^\top x\}$$

Binary SVM

$$\min_{w, \xi_x \geq 0} \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y) \in \mathcal{D}} \xi_x \quad \text{s.t.} \quad yw^\top x \geq 1 - \xi_x \quad \forall (x,y) \in \mathcal{D}$$

Equivalent:

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y) \in \mathcal{D}} \max\{0, 1 - yw^\top x\}$$

Generally:

$$\min_w R(w) + \frac{C}{N} \sum_{(x,y) \in \mathcal{D}} \bar{\ell}(x, y, w)$$

Binary SVM

$$\min_{w, \xi_x \geq 0} \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y) \in \mathcal{D}} \xi_x \quad \text{s.t.} \quad yw^\top x \geq 1 - \xi_x \quad \forall (x,y) \in \mathcal{D}$$

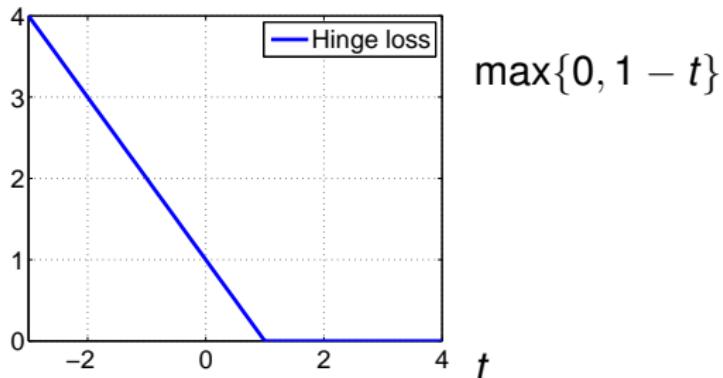
Equivalent:

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y) \in \mathcal{D}} \max\{0, 1 - yw^\top x\}$$

Generally:

$$\min_w R(w) + \frac{C}{N} \sum_{(x,y) \in \mathcal{D}} \bar{\ell}(x, y, w)$$

Hinge-Loss $\bar{\ell}$:



Multiclass SVM

Multiclass SVM

Recap of Inference:

$$y \in \{1, 2, \dots, K\}$$

$$y^* = \arg \max_{\hat{y}} \begin{bmatrix} w_1 \\ \vdots \\ w_K \end{bmatrix}^\top \begin{bmatrix} x\delta(\hat{y}=1) \\ \vdots \\ x\delta(\hat{y}=K) \end{bmatrix}$$

More generally:

Multiclass SVM

Recap of Inference:

$$y \in \{1, 2, \dots, K\}$$

$$y^* = \arg \max_{\hat{y}} \begin{bmatrix} w_1 \\ \vdots \\ w_K \end{bmatrix}^\top \begin{bmatrix} x\delta(\hat{y}=1) \\ \vdots \\ x\delta(\hat{y}=K) \end{bmatrix}$$

More generally:

$$y^* = \arg \max_{\hat{y}} w^\top \phi(x, \hat{y})$$

Even more generally:

Multiclass SVM

Recap of Inference:

$$y \in \{1, 2, \dots, K\}$$

$$y^* = \arg \max_{\hat{y}} \begin{bmatrix} w_1 \\ \vdots \\ w_K \end{bmatrix}^\top \begin{bmatrix} x\delta(\hat{y}=1) \\ \vdots \\ x\delta(\hat{y}=K) \end{bmatrix}$$

More generally:

$$y^* = \arg \max_{\hat{y}} w^\top \phi(x, \hat{y})$$

Even more generally:

$$y^* = \arg \max_{\hat{y}} F(\hat{y}, x, w)$$

Multiclass SVM

What do we want?

Multiclass SVM

What do we want? Groundtruth y scores higher than any other \hat{y}

Multiclass SVM

What do we want? Groundtruth y scores higher than any other \hat{y}

$$\mathbf{w}^\top \phi(x, y) \geq \mathbf{w}^\top \phi(x, \hat{y}) \quad \forall (x, y) \in \mathcal{D}, \hat{y} \in \{1, \dots, K\}$$

Multiclass SVM

What do we want? Groundtruth y scores higher than any other \hat{y}

$$\mathbf{w}^\top \phi(x, y) \geq \mathbf{w}^\top \phi(x, \hat{y}) \quad \forall (x, y) \in \mathcal{D}, \hat{y} \in \{1, \dots, K\}$$

Let's employ margin and slack:

$$\begin{aligned} \min_{\mathbf{w}, \xi_x \geq 0} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + \frac{C}{N} \sum_{(x,y)} \xi_x \\ \text{s.t.} \quad & \mathbf{w}^\top (\phi(x, y) - \phi(x, \hat{y})) \geq 1 - \xi_x \quad \forall (x, y), \hat{y} \end{aligned}$$

Multiclass SVM

$$\begin{aligned} \min_{w, \xi_x \geq 0} \quad & \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y)} \xi_x \\ \text{s.t.} \quad & w^\top (\phi(x, y) - \phi(x, \hat{y})) \geq 1 - \xi_x \quad \forall (x, y), \hat{y} \end{aligned}$$

Properties:

Multiclass SVM

$$\begin{aligned} \min_{w, \xi_x \geq 0} \quad & \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y)} \xi_x \\ \text{s.t.} \quad & w^\top (\phi(x, y) - \phi(x, \hat{y})) \geq 1 - \xi_x \quad \forall (x, y), \hat{y} \end{aligned}$$

Properties:

- Quadratic cost function
- Linear constraints
- Convex program

Algorithms:

Multiclass SVM

$$\begin{aligned} \min_{w, \xi_x \geq 0} \quad & \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y)} \xi_x \\ \text{s.t.} \quad & w^\top (\phi(x, y) - \phi(x, \hat{y})) \geq 1 - \xi_x \quad \forall (x, y), \hat{y} \end{aligned}$$

Properties:

- Quadratic cost function
- Linear constraints
- Convex program

Algorithms:

- Primal algorithms (Sub-gradient, cutting plane)
- Dual algorithms

Multiclass SVM

$$\begin{aligned} \min_{w, \xi_x \geq 0} \quad & \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y)} \xi_x \\ \text{s.t.} \quad & w^\top (\phi(x, y) - \phi(x, \hat{y})) \geq 1 - \xi_x \quad \forall (x, y), \hat{y} \end{aligned}$$

Properties:

- Quadratic cost function
- Linear constraints
- Convex program

Algorithms:

- Primal algorithms (Sub-gradient, cutting plane)
- Dual algorithms

How does this fit into regularized surrogate loss minimization?

$$\min_w R(w) + \frac{C}{N} \sum_{(x,y)} \bar{\ell}(x, y, w)$$

Multiclass SVM

$$\begin{aligned} \min_{w, \xi_x \geq 0} \quad & \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y)} \xi_x \\ \text{s.t.} \quad & w^\top (\phi(x, y) - \phi(x, \hat{y})) \geq 1 - \xi_x \quad \forall (x, y), \hat{y} \end{aligned}$$

Let's get rid of the slack variable:

Multiclass SVM

$$\begin{aligned} \min_{w, \xi_x \geq 0} \quad & \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y)} \xi_x \\ \text{s.t.} \quad & w^\top (\phi(x, y) - \phi(x, \hat{y})) \geq 1 - \xi_x \quad \forall (x, y), \hat{y} \end{aligned}$$

Let's get rid of the slack variable:

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y)} \max \left\{ 0, \max_{\hat{y}} \left(1 - w^\top (\phi(x, y) - \phi(x, \hat{y})) \right) \right\}$$

Multiclass SVM

$$\begin{aligned} \min_{w, \xi_x \geq 0} \quad & \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y)} \xi_x \\ \text{s.t.} \quad & w^\top (\phi(x, y) - \phi(x, \hat{y})) \geq 1 - \xi_x \quad \forall (x, y), \hat{y} \end{aligned}$$

Let's get rid of the slack variable:

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y)} \max \left\{ 0, \underbrace{\max_{\hat{y}} \left(1 - w^\top (\phi(x, y) - \phi(x, \hat{y})) \right)}_{\geq 0} \right\}$$

Multiclass SVM

$$\begin{aligned} \min_{w, \xi_x \geq 0} \quad & \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y)} \xi_x \\ \text{s.t.} \quad & w^\top (\phi(x, y) - \phi(x, \hat{y})) \geq 1 - \xi_x \quad \forall (x, y), \hat{y} \end{aligned}$$

Let's get rid of the slack variable:

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y)} \max \left\{ 0, \underbrace{\max_{\hat{y}} \left(1 - w^\top (\phi(x, y) - \phi(x, \hat{y})) \right)}_{\geq 0} \right\}$$

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y)} \max_{\hat{y}} (1 + w^\top \phi(x, \hat{y})) - w^\top \phi(x, y)$$

Multiclass SVM

$$\begin{aligned} \min_{w, \xi_x \geq 0} \quad & \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y)} \xi_x \\ \text{s.t.} \quad & w^\top (\phi(x, y) - \phi(x, \hat{y})) \geq 1 - \xi_x \quad \forall (x, y), \hat{y} \end{aligned}$$

Let's get rid of the slack variable:

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y)} \max \left\{ 0, \underbrace{\max_{\hat{y}} \left(1 - w^\top (\phi(x, y) - \phi(x, \hat{y})) \right)}_{\geq 0} \right\}$$

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y)} \underbrace{\max_{\hat{y}} (1 + w^\top \phi(x, \hat{y})) - w^\top \phi(x, y)}_{\text{Loss-augmented inference}}$$

Multiclass SVM

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y)} \underbrace{\max_{\hat{y}} (1 + w^\top \phi(x, \hat{y})) - w^\top \phi(x, y)}_{\text{Loss-augmented inference}}$$

Fits into the general formulation:

Multiclass SVM

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y)} \underbrace{\max_{\hat{y}} (1 + w^\top \phi(x, \hat{y})) - w^\top \phi(x, y)}_{\text{Loss-augmented inference}}$$

Fits into the general formulation:

$$\min_w R(w) + \frac{C}{N} \sum_{(x,y)} \bar{\ell}(x, y, w)$$

Multiclass SVM

How to optimize this?

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y)} \underbrace{\max_{\hat{y}} (1 + w^\top \phi(x, \hat{y})) - w^\top \phi(x, y)}_{\text{Loss-augmented inference}}$$

E.g., with gradient descent:

Iterate:

Multiclass SVM

How to optimize this?

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y)} \underbrace{\max_{\hat{y}} (1 + w^\top \phi(x, \hat{y})) - w^\top \phi(x, y)}_{\text{Loss-augmented inference}}$$

E.g., with gradient descent:

Iterate:

① Loss-augmented inference:

$$\arg \max_{\hat{y}_i} (1 + w^\top \phi(x, \hat{y}_i))$$

Multiclass SVM

How to optimize this?

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y)} \underbrace{\max_{\hat{y}} (1 + w^\top \phi(x, \hat{y})) - w^\top \phi(x, y)}_{\text{Loss-augmented inference}}$$

E.g., with gradient descent:

Iterate:

- ① Loss-augmented inference:

$$\arg \max_{\hat{y}_i} (1 + w^\top \phi(x, \hat{y}_i))$$

- ② Perform gradient step:

$$w \leftarrow w - \alpha \nabla_w L$$

Multiclass SVM

How to optimize this?

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y)} \underbrace{\max_{\hat{y}} (1 + w^\top \phi(x, \hat{y})) - w^\top \phi(x, y)}_{\text{Loss-augmented inference}}$$

E.g., with gradient descent:

Iterate:

- ① Loss-augmented inference:

$$\arg \max_{\hat{y}_i} (1 + w^\top \phi(x, \hat{y}_i))$$

- ② Perform gradient step:

$$w \leftarrow w - \alpha \nabla_w L$$

How complicated is this?

Solve one multi-class inference task per sample per iteration.

Structured SVM

Structured SVM

From Multiclass SVM:

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y)} \max_{\hat{y}} (1 + w^\top \phi(x, \hat{y})) - w^\top \phi(x, y)$$

Structured SVM

From Multiclass SVM:

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y)} \max_{\hat{y}} (1 + w^\top \phi(x, \hat{y})) - w^\top \phi(x, y)$$

via general losses:

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y)} \max_{\hat{y}} (\Delta(y, \hat{y}) + w^\top \phi(x, \hat{y})) - w^\top \phi(x, y)$$

Structured SVM

From Multiclass SVM:

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y)} \max_{\hat{y}} (1 + w^\top \phi(x, \hat{y})) - w^\top \phi(x, y)$$

via general losses:

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y)} \max_{\hat{y}} (\Delta(y, \hat{y}) + w^\top \phi(x, \hat{y})) - w^\top \phi(x, y)$$

to general functions and structures:

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y)} \underbrace{\max_{\hat{y}} (\Delta(\mathbf{y}, \hat{\mathbf{y}}) + F(\hat{\mathbf{y}}, x, w))}_{\text{Loss-augmented inference}} - F(\mathbf{y}, x, w)$$

Structured SVM

How to optimize this?

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y) \in \mathcal{D}} \underbrace{\max_{\hat{y}} (\Delta(y, \hat{y}) + F(\hat{y}, x, w)) - F(y, x, w)}_{\text{Loss-augmented inference}}$$

E.g., with gradient descent:

Iterate:

Structured SVM

How to optimize this?

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x, \mathbf{y}) \in \mathcal{D}} \underbrace{\max_{\hat{\mathbf{y}}} (\Delta(\mathbf{y}, \hat{\mathbf{y}}) + F(\hat{\mathbf{y}}, x, w)) - F(\mathbf{y}, x, w)}_{\text{Loss-augmented inference}}$$

E.g., with gradient descent:

Iterate:

① Loss-augmented inference:

$$\arg \max_{\hat{\mathbf{y}}} (\Delta(\mathbf{y}, \hat{\mathbf{y}}) + F(\hat{\mathbf{y}}, x, w))$$

Structured SVM

How to optimize this?

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y) \in \mathcal{D}} \underbrace{\max_{\hat{y}} (\Delta(y, \hat{y}) + F(\hat{y}, x, w)) - F(y, x, w)}_{\text{Loss-augmented inference}}$$

E.g., with gradient descent:

Iterate:

- ① Loss-augmented inference:

$$\arg \max_{\hat{y}} (\Delta(y, \hat{y}) + F(\hat{y}, x, w))$$

- ② Perform gradient step:

$$w \leftarrow w - \alpha \nabla_w L$$

Structured SVM

How to optimize this?

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y) \in \mathcal{D}} \underbrace{\max_{\hat{y}} (\Delta(y, \hat{y}) + F(\hat{y}, x, w)) - F(y, x, w)}_{\text{Loss-augmented inference}}$$

E.g., with gradient descent:

Iterate:

- ① Loss-augmented inference:

$$\arg \max_{\hat{y}} (\Delta(y, \hat{y}) + F(\hat{y}, x, w))$$

- ② Perform gradient step:

$$w \leftarrow w - \alpha \nabla_w L$$

How complicated is this?

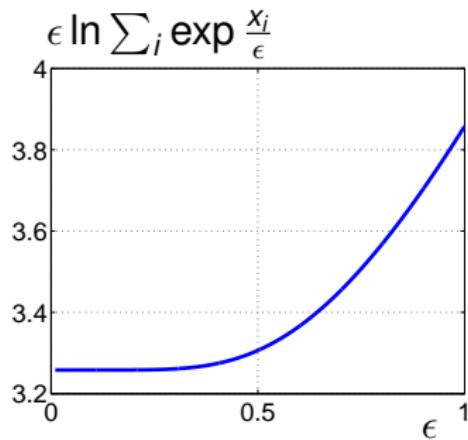
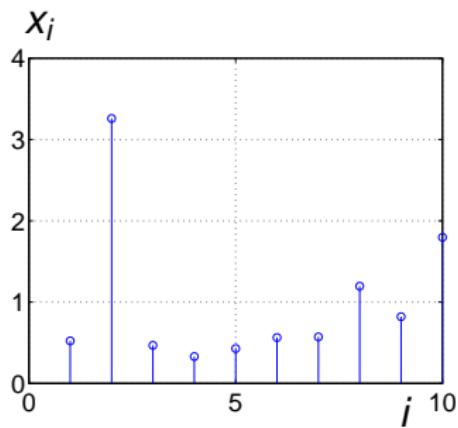
Solve one structured prediction task per sample per iteration.

Structured Learning

Structured Learning

Soft-max extension:

$$\epsilon \ln \sum_i \exp \frac{x_i}{\epsilon} \xrightarrow{\epsilon \rightarrow 0} \max_i x_i$$



Structured Learning

Soft-max extension:

$$\epsilon \ln \sum_i \exp \frac{x_i}{\epsilon} \xrightarrow{\epsilon \rightarrow 0} \max_i x_i$$

Let's generalize

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x, \mathbf{y}) \in \mathcal{D}} \max_{\hat{\mathbf{y}}} (\Delta(\mathbf{y}, \hat{\mathbf{y}}) + F(\hat{\mathbf{y}}, x, w)) - F(\mathbf{y}, x, w)$$

even further

Structured Learning

Soft-max extension:

$$\epsilon \ln \sum_i \exp \frac{x_i}{\epsilon} \xrightarrow{\epsilon \rightarrow 0} \max_i x_i$$

Let's generalize

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x, \mathbf{y}) \in \mathcal{D}} \max_{\hat{\mathbf{y}}} (\Delta(\mathbf{y}, \hat{\mathbf{y}}) + F(\hat{\mathbf{y}}, x, w)) - F(\mathbf{y}, x, w)$$

even further

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x, \mathbf{y}) \in \mathcal{D}} \epsilon \ln \sum_{\hat{\mathbf{y}}} \exp \frac{\Delta(\mathbf{y}, \hat{\mathbf{y}}) + F(\hat{\mathbf{y}}, x, w)}{\epsilon} - F(\mathbf{y}, x, w)$$

Structured Learning

A pretty general formulation:

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y) \in \mathcal{D}} \epsilon \ln \sum_{\hat{y}} \exp \frac{\Delta(y, \hat{y}) + F(\hat{y}, x, w)}{\epsilon} - F(y, x, w)$$

Structured Learning

A pretty general formulation:

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y) \in \mathcal{D}} \epsilon \ln \sum_{\hat{y}} \exp \frac{\Delta(y, \hat{y}) + F(\hat{y}, x, w)}{\epsilon} - F(y, x, w)$$

Generalizes:

- any form of SVM
- conditional random fields
- maximum-likelihood
- logistic regression
- convolutional neural networks

Structured Learning

A pretty general formulation:

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y) \in \mathcal{D}} \epsilon \ln \sum_{\hat{y}} \exp \frac{\Delta(y, \hat{y}) + F(\hat{y}, x, w)}{\epsilon} - F(y, x, w)$$

Generalizes:

- any form of SVM
- conditional random fields
- maximum-likelihood
- logistic regression
- convolutional neural networks

How does this correspond to a binary SVM?

Structured Learning

A pretty general formulation:

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y) \in \mathcal{D}} \epsilon \ln \sum_{\hat{y}} \exp \frac{\Delta(y, \hat{y}) + F(\hat{y}, x, w)}{\epsilon} - F(y, x, w)$$

Generalizes:

- any form of SVM
- conditional random fields
- maximum-likelihood
- logistic regression
- convolutional neural networks

How does this correspond to a binary SVM? Reverse slides.

Structured Learning

A pretty general formulation:

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y) \in \mathcal{D}} \epsilon \ln \sum_{\hat{y}} \exp \frac{\Delta(y, \hat{y}) + F(\hat{y}, x, w)}{\epsilon} - F(y, x, w)$$

Generalizes:

- any form of SVM
- conditional random fields
- maximum-likelihood
- logistic regression
- convolutional neural networks

How does this correspond to a binary SVM? Reverse slides.

How does this correspond to logistic regression?

Structured Learning

A pretty general formulation:

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y) \in \mathcal{D}} \epsilon \ln \sum_{\hat{y}} \exp \frac{\Delta(y, \hat{y}) + F(\hat{y}, x, w)}{\epsilon} - F(y, x, w)$$

Generalizes:

- any form of SVM
- conditional random fields
- maximum-likelihood
- logistic regression
- convolutional neural networks

How does this correspond to a binary SVM? Reverse slides.

How does this correspond to logistic regression? Up next...

Structured Learning

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x, \mathbf{y}) \in \mathcal{D}} \epsilon \ln \sum_{\hat{\mathbf{y}}} \exp \frac{\Delta(\mathbf{y}, \hat{\mathbf{y}}) + F(\hat{\mathbf{y}}, x, w)}{\epsilon} - F(\mathbf{y}, x, w)$$

Structured Learning

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x, \mathbf{y}) \in \mathcal{D}} \epsilon \ln \sum_{\hat{\mathbf{y}}} \exp \frac{\Delta(\mathbf{y}, \hat{\mathbf{y}}) + F(\hat{\mathbf{y}}, x, w)}{\epsilon} - F(\mathbf{y}, x, w)$$

Equivalent formulation:

$$\min_w \frac{1}{2} \|w\|_2^2 - \frac{C}{N} \sum_{(x, \mathbf{y}) \in \mathcal{D}} \epsilon \ln \frac{\exp \frac{\Delta(\mathbf{y}, \mathbf{y}) + F(\mathbf{y}, x, w)}{\epsilon}}{\sum_{\hat{\mathbf{y}}} \exp \frac{\Delta(\mathbf{y}, \hat{\mathbf{y}}) + F(\hat{\mathbf{y}}, x, w)}{\epsilon}}$$

Structured Learning

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x, \mathbf{y}) \in \mathcal{D}} \epsilon \ln \sum_{\hat{\mathbf{y}}} \exp \frac{\Delta(\mathbf{y}, \hat{\mathbf{y}}) + F(\hat{\mathbf{y}}, x, w)}{\epsilon} - F(\mathbf{y}, x, w)$$

Equivalent formulation:

$$\min_w \frac{1}{2} \|w\|_2^2 - \frac{C}{N} \sum_{(x, \mathbf{y}) \in \mathcal{D}} \epsilon \ln \frac{\exp \frac{\Delta(\mathbf{y}, \mathbf{y}) + F(\mathbf{y}, x, w)}{\epsilon}}{\sum_{\hat{\mathbf{y}}} \exp \frac{\Delta(\mathbf{y}, \hat{\mathbf{y}}) + F(\hat{\mathbf{y}}, x, w)}{\epsilon}}$$

Annealed loss-augmented probability:

$$p_\epsilon^\Delta(\mathbf{y} \mid x, w) = \left(\frac{\exp \frac{\Delta(\mathbf{y}, \mathbf{y}) + F(\mathbf{y}, x, w)}{\epsilon}}{\sum_{\hat{\mathbf{y}}} \exp \frac{\Delta(\mathbf{y}, \hat{\mathbf{y}}) + F(\hat{\mathbf{y}}, x, w)}{\epsilon}} \right)^\epsilon$$

Structured Learning

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x, \mathbf{y}) \in \mathcal{D}} \epsilon \ln \sum_{\hat{\mathbf{y}}} \exp \frac{\Delta(\mathbf{y}, \hat{\mathbf{y}}) + F(\hat{\mathbf{y}}, x, w)}{\epsilon} - F(\mathbf{y}, x, w)$$

Equivalent formulation:

$$\min_w \frac{1}{2} \|w\|_2^2 - \frac{C}{N} \sum_{(x, \mathbf{y}) \in \mathcal{D}} \epsilon \ln \frac{\exp \frac{\Delta(\mathbf{y}, \mathbf{y}) + F(\mathbf{y}, x, w)}{\epsilon}}{\sum_{\hat{\mathbf{y}}} \exp \frac{\Delta(\mathbf{y}, \hat{\mathbf{y}}) + F(\hat{\mathbf{y}}, x, w)}{\epsilon}}$$

Annealed loss-augmented probability:

$$p_\epsilon^\Delta(\mathbf{y} | x, w) = \left(\frac{\exp \frac{\Delta(\mathbf{y}, \mathbf{y}) + F(\mathbf{y}, x, w)}{\epsilon}}{\sum_{\hat{\mathbf{y}}} \exp \frac{\Delta(\mathbf{y}, \hat{\mathbf{y}}) + F(\hat{\mathbf{y}}, x, w)}{\epsilon}} \right)^\epsilon$$

Equivalent maximum-likelihood formulation:

$$\min_w \frac{1}{2} \|w\|_2^2 - \frac{C}{N} \ln \prod_{(x, \mathbf{y}) \in \mathcal{D}} p_\epsilon^\Delta(\mathbf{y} | x, w)$$

$$\min_w \frac{1}{2} \|w\|_2^2 - \frac{C}{N} \sum_{(x, \mathbf{y}) \in \mathcal{D}} \epsilon \ln \frac{\exp \frac{\Delta(\mathbf{y}, \mathbf{y}) + F(\mathbf{y}, x, w)}{\epsilon}}{\sum_{\hat{\mathbf{y}}} \exp \frac{\Delta(\mathbf{y}, \hat{\mathbf{y}}) + F(\hat{\mathbf{y}}, x, w)}{\epsilon}}$$

Simplifications:

$$\min_w \frac{1}{2} \|w\|_2^2 - \frac{C}{N} \sum_{(x, \mathbf{y}) \in \mathcal{D}} \epsilon \ln \frac{\exp \frac{\Delta(\mathbf{y}, \mathbf{y}) + F(\mathbf{y}, x, w)}{\epsilon}}{\sum_{\hat{\mathbf{y}}} \exp \frac{\Delta(\mathbf{y}, \hat{\mathbf{y}}) + F(\hat{\mathbf{y}}, x, w)}{\epsilon}}$$

Simplifications:

- Binary classification: $y \in \{-1, 1\}$
- Simple feature-function: $F(\mathbf{y}, x, w) = yw^\top x$
- Epsilon: $\epsilon = 1$

$$\min_w \frac{1}{2} \|w\|_2^2 - \frac{C}{N} \sum_{(x, \mathbf{y}) \in \mathcal{D}} \epsilon \ln \frac{\exp \frac{\Delta(\mathbf{y}, \mathbf{y}) + F(\mathbf{y}, x, w)}{\epsilon}}{\sum_{\hat{\mathbf{y}}} \exp \frac{\Delta(\mathbf{y}, \hat{\mathbf{y}}) + F(\hat{\mathbf{y}}, x, w)}{\epsilon}}$$

Simplifications:

- Binary classification: $y \in \{-1, 1\}$
- Simple feature-function: $F(\mathbf{y}, x, w) = yw^\top x$
- Epsilon: $\epsilon = 1$

$$\min_w \frac{1}{2} \|w\|_2^2 - \frac{C}{N} \ln \prod_{(x, y)} \frac{\exp yw^\top x}{\exp(w^\top x) + \exp(-w^\top x)}$$

$$\min_w \frac{1}{2} \|w\|_2^2 - \frac{C}{N} \sum_{(x,y) \in \mathcal{D}} \epsilon \ln \frac{\exp \frac{\Delta(y, y) + F(y, x, w)}{\epsilon}}{\sum_{\hat{y}} \exp \frac{\Delta(y, \hat{y}) + F(\hat{y}, x, w)}{\epsilon}}$$

Simplifications:

- Binary classification: $y \in \{-1, 1\}$
- Simple feature-function: $F(y, x, w) = yw^\top x$
- Epsilon: $\epsilon = 1$

$$\min_w \frac{1}{2} \|w\|_2^2 - \frac{C}{N} \ln \prod_{(x,y)} \frac{\exp yw^\top x}{\exp(w^\top x) + \exp(-w^\top x)}$$

$$\min_w \frac{1}{2} \|w\|_2^2 - \frac{C}{N} \ln \prod_{(x,y)} \frac{1}{1 + \exp(-2yw^\top x)}$$

$$\min_w \frac{1}{2} \|w\|_2^2 - \frac{C}{N} \sum_{(x,y) \in \mathcal{D}} \epsilon \ln \frac{\exp \frac{\Delta(y, y) + F(y, x, w)}{\epsilon}}{\sum_{\hat{y}} \exp \frac{\Delta(y, \hat{y}) + F(\hat{y}, x, w)}{\epsilon}}$$

Simplifications:

- Binary classification: $y \in \{-1, 1\}$
- Simple feature-function: $F(y, x, w) = yw^\top x$
- Epsilon: $\epsilon = 1$

$$\min_w \frac{1}{2} \|w\|_2^2 - \frac{C}{N} \ln \prod_{(x,y)} \frac{\exp yw^\top x}{\exp(w^\top x) + \exp(-w^\top x)}$$

$$\min_w \frac{1}{2} \|w\|_2^2 - \frac{C}{N} \ln \prod_{(x,y)} \frac{1}{1 + \exp(-2yw^\top x)}$$

$$\min_w \frac{1}{2} \|w\|_2^2 - \frac{C}{N} \sum_{(x,y)} \ln \sigma(2yw^\top x) \quad \text{with} \quad \sigma(z) = \frac{1}{1 + \exp(-z)}$$

We started with hinge loss:

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y)}^n \max\{0, 1 - yw^\top x\}$$

We started with hinge loss:

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y)}^n \max\{0, 1 - yw^\top x\}$$

Derived the general structured prediction framework:

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,\mathbf{y}) \in \mathcal{D}} \epsilon \ln \sum_{\hat{\mathbf{y}}} \exp \frac{\Delta(\mathbf{y}, \hat{\mathbf{y}}) + F(\hat{\mathbf{y}}, x, w)}{\epsilon} - F(\mathbf{y}, x, w)$$

We started with hinge loss:

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y)}^n \max\{0, 1 - yw^\top x\}$$

Derived the general structured prediction framework:

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,\mathbf{y}) \in \mathcal{D}} \epsilon \ln \sum_{\hat{\mathbf{y}}} \exp \frac{\Delta(\mathbf{y}, \hat{\mathbf{y}}) + F(\hat{\mathbf{y}}, x, w)}{\epsilon} - F(\mathbf{y}, x, w)$$

And showed how to simplify it to logistic regression:

$$\min_w \frac{1}{2} \|w\|_2^2 - \frac{C}{N} \sum_{(x,y)} \ln \sigma(2yw^\top x)$$

We started with hinge loss:

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y)}^n \max\{0, 1 - yw^\top x\}$$

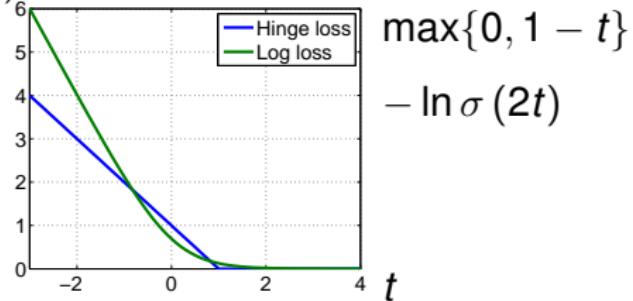
Derived the general structured prediction framework:

$$\min_w \frac{1}{2} \|w\|_2^2 + \frac{C}{N} \sum_{(x,y) \in \mathcal{D}} \epsilon \ln \sum_{\hat{y}} \exp \frac{\Delta(y, \hat{y}) + F(\hat{y}, x, w)}{\epsilon} - F(y, x, w)$$

And showed how to simplify it to logistic regression:

$$\min_w \frac{1}{2} \|w\|_2^2 - \frac{C}{N} \sum_{(x,y)} \ln \sigma(2yw^\top x)$$

$$\min_w R(w) + \frac{C}{N} \sum_{(x,y)} \bar{\ell}(x, y, w)$$



Structured Learning

How does this framework simplify to neural networks?

- Multiclass setting: $y \in \{1, \dots, K\}$
- Composite feature function:

$$F(y, x, w) = f_1\left(y, w_1, f_2(w_2, f_3(\dots))\right)$$

- Epsilon: $\epsilon = 1$

Stay tuned for part III