

# Multi-class Classification, Maximum Entropy and Structured Classification

# Outline

- Multi-class Classification
- Structured Learning



# Outline

- Multi-class Classification
- Structured Learning



# One-against-the rest

- Assume data in  $k$  classes:  $\{1, \dots, k\}$
- Train  $k$  binary SVMs:

1st class vs.  $(2, \dots, k)$ th class  
 2nd class vs.  $(1, 3, \dots, k)$ th class  
 $\vdots$

- $k$  decision functions

$$\begin{aligned} &(\mathbf{w}_1)^T \mathbf{x} \\ &\vdots \\ &(\mathbf{w}_k)^T \mathbf{x} \end{aligned}$$



# One-against-the rest (Cont'd)

- Prediction:

$$\arg \max_j (\mathbf{w}_j)^T \mathbf{x}$$

- Reason: If  $\mathbf{x} \in$  1st class, then we should have

$$(\mathbf{w}_1)^T \mathbf{x} \geq +1$$

$$(\mathbf{w}_2)^T \mathbf{x} \leq -1$$

$$\vdots$$

$$(\mathbf{w}_k)^T \mathbf{x} + b_k \leq -1$$



# Multi-class Classification (Cont'd)

- One-against-one: train  $k(k-1)/2$  binary SVMs  
 $(1, 2), (1, 3), \dots, (1, k), (2, 3), (2, 4), \dots, (k-1, k)$
- If 4 classes  $\Rightarrow$  6 binary SVMs

$y_i = 1$	$y_i = -1$	Decision functions
class 1	class 2	$f_{12}(\mathbf{x}) = (\mathbf{w}_{12})^T \mathbf{x} + b_{12}$
class 1	class 3	$f_{13}(\mathbf{x}) = (\mathbf{w}_{13})^T \mathbf{x} + b_{13}$
class 1	class 4	$f_{14}(\mathbf{x}) = (\mathbf{w}_{14})^T \mathbf{x} + b_{14}$
class 2	class 3	$f_{23}(\mathbf{x}) = (\mathbf{w}_{23})^T \mathbf{x} + b_{23}$
class 2	class 4	$f_{24}(\mathbf{x}) = (\mathbf{w}_{24})^T \mathbf{x} + b_{24}$
class 3	class 4	$f_{34}(\mathbf{x}) = (\mathbf{w}_{34})^T \mathbf{x} + b_{34}$



- For a testing data, predicting all binary SVMs

Classes		winner
1	2	1
1	3	1
1	4	1
2	3	2
2	4	4
3	4	3

- Select the one with **the largest vote**

class	1	2	3	4
# votes	3	1	1	1

- May use decision values as well



# Solving a Single Problem

- Example (Crammer and Singer, 2002)

$$\min_{\mathbf{w}_1, \dots, \mathbf{w}_k} \quad \frac{1}{2} \sum_{m=1}^k \|\mathbf{w}_m\|_2^2 + C \sum_{i=1}^l \xi(\{\mathbf{w}_m\}_{m=1}^k; \mathbf{x}_i, y_i),$$

where

$$\xi(\{\mathbf{w}_m\}_{m=1}^k; \mathbf{x}, y) \equiv \max_{m \neq y} \max(0, 1 - (\mathbf{w}_y - \mathbf{w}_m)^T \mathbf{x}).$$

- We hope the decision value of  $\mathbf{x}_i$  by the model  $\mathbf{w}_{y_i}$  is larger than others
- Prediction: same as one-against

$$\arg \max_j (\mathbf{w}_j)^T \mathbf{x}$$





# Maximum Entropy

- Maximum Entropy: a generalization of logistic regression for multi-class problems
- It is widely applied by NLP applications.
- Conditional probability of label  $y$  given data  $\mathbf{x}$ .

$$P(y|\mathbf{x}) \equiv \frac{\exp(\mathbf{w}_y^T \mathbf{x})}{\sum_{m=1}^k \exp(\mathbf{w}_m^T \mathbf{x})},$$



# Maximum Entropy (Cont'd)

- We then minimize regularized negative log-likelihood.

$$\min_{\mathbf{w}_1, \dots, \mathbf{w}_m} \frac{1}{2} \sum_{m=1}^k \|\mathbf{w}_m\|^2 + C \sum_{i=1}^l \xi(\{\mathbf{w}_m\}_{m=1}^k; \mathbf{x}_i, y_i),$$

where

$$\xi(\{\mathbf{w}_m\}_{m=1}^k; \mathbf{x}, y) \equiv -\log P(y|\mathbf{x}).$$



# Maximum Entropy (Cont'd)

- Is this loss function reasonable?
- If

$$\mathbf{w}_{y_i}^T \mathbf{x}_i \gg \mathbf{w}_m^T \mathbf{x}_i, \forall m \neq y_i,$$

then

$$\xi(\{\mathbf{w}_m\}_{m=1}^k; \mathbf{x}_i, y_i) \approx 0$$

That is, no loss

- In contrast, if

$$\mathbf{w}_{y_i}^T \mathbf{x}_i \ll \mathbf{w}_m^T \mathbf{x}_i, m \neq y_i,$$

then  $P(y_i|\mathbf{x}_i) \ll 1$  and the loss is large.



# Features as Functions

- NLP applications often use a function  $\mathbf{f}(\mathbf{x}, y)$  to generate the feature vector

$$P(y|\mathbf{x}) \equiv \frac{\exp(\mathbf{w}^T \mathbf{f}(\mathbf{x}, y))}{\sum_{y'} \exp(\mathbf{w}^T \mathbf{f}(\mathbf{x}, y'))}. \quad (1)$$

- The earlier probability model is a special case by

$$\mathbf{f}(\mathbf{x}_i, y) = \left[ \begin{array}{c} 0 \\ \vdots \\ 0 \\ \mathbf{x}_i \\ 0 \\ \vdots \\ 0 \end{array} \right] \left. \vphantom{\begin{bmatrix} 0 \\ \vdots \\ 0 \\ \mathbf{x}_i \\ 0 \\ \vdots \\ 0 \end{bmatrix}} \right\} y - 1 \in \mathbf{R}^{n_k} \text{ and } \mathbf{w} = \left[ \begin{array}{c} \mathbf{w}_1 \\ \vdots \\ \mathbf{w}_k \end{array} \right].$$



# Outline

- Multi-class Classification
- Structured Learning



# Structured Data

- So far we assume that the label  $y_i$  is a **single value**
- In some applications, the label may be a more sophisticated object.
- For example, in part-of-speech (POS) tagging, a training instance is a sentence and a label is a sequence of POS tags of words.
- For  $l$  sentences, training instances are

$$(\mathbf{y}_i, \mathbf{x}_i) \in Y^{n_i} \times X^{n_i}, \forall i = 1, \dots, l,$$

where  $\mathbf{x}_i$  is the  $i$ th sentence,  $\mathbf{y}_i$  is a sequence of tags,



# Structured Data (Cont'd)

- $X$ : set of unique words in the context
- $Y$ : set of candidate tags for each word
- $n_i$ : number of words in the  $i$ th sentence.



# Conditional Random Fields

- CRF solves

$$\min_{\mathbf{w}} \quad \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^I \xi(\mathbf{w}; \mathbf{x}_i, \mathbf{y}_i),$$

where

$$\xi(\mathbf{w}; \mathbf{x}_i, \mathbf{y}_i) \equiv -\log P(\mathbf{y}_i | \mathbf{x}_i), \text{ and}$$

$$P(\mathbf{y} | \mathbf{x}) \equiv \frac{\exp(\mathbf{w}^T \mathbf{f}(\mathbf{x}, \mathbf{y}))}{\sum_{\mathbf{y}'} \exp(\mathbf{w}^T \mathbf{f}(\mathbf{x}, \mathbf{y}'))}. \quad (2)$$

- Dynamic programming used to handle **exponentially** many  $\mathbf{y}$





# Structured SVM

- Structured SVM solves

$$\min_{\mathbf{w}} \quad \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^I \xi(\mathbf{w}; \mathbf{x}_i, \mathbf{y}_i),$$

where

$$\begin{aligned} & \xi(\mathbf{w}; \mathbf{x}_i, \mathbf{y}_i) \\ & \equiv \max_{\mathbf{y} \neq \mathbf{y}_i} \left( \max \left( 0, \Delta(\mathbf{y}_i, \mathbf{y}) - \mathbf{w}^T (f(\mathbf{x}_i, \mathbf{y}_i) - f(\mathbf{x}_i, \mathbf{y})) \right) \right) \end{aligned}$$

- $\Delta(\cdot)$  is a distance function



# Structured SVM (Cont'd)

- $\Delta(\cdot)$  should satisfy

$$\Delta(\mathbf{y}_i, \mathbf{y}_i) = 0 \text{ and } \Delta(\mathbf{y}_i, \mathbf{y}_j) = \Delta(\mathbf{y}_j, \mathbf{y}_i)$$

- If

$$\Delta(\mathbf{y}_i, \mathbf{y}_j) = \begin{cases} 0 & \text{if } \mathbf{y}_i = \mathbf{y}_j \\ 1 & \text{otherwise,} \end{cases}$$

and

$$\mathbf{y}_i \in \{1, \dots, k\}, \forall i,$$

then structured SVM becomes the multi-class formulation discussed earlier



# More Information

- See Sections V and VIII in Yuan et al. (2012) and references therein



# References I

- K. Crammer and Y. Singer. On the learnability and design of output codes for multiclass problems. *Machine Learning*, (2–3):201–233, 2002.
- G.-X. Yuan, C.-H. Ho, and C.-J. Lin. Recent advances of large-scale linear classification. *Proceedings of the IEEE*, 100(9):2584–2603, 2012. URL <http://www.csie.ntu.edu.tw/~cjlin/papers/survey-linear.pdf>.

