

Statistical Machine Learning Theory
Lecture 9
**From Multi-class Classification
to Structured Output Prediction**

Hisashi Kashima
kashima@i.Kyoto-u.ac.jp

Topics of the 2nd half of this course:

Advanced supervised learning and unsupervised learning

- Multi-class classification and structured output prediction
- Other variants of supervised learning problems
 - Semi-supervised learning
 - Active learning
 - Transfer learning
- On-line learning
 - Aggregating experts' predictions
 - Bandits and reinforcement learning
- Unsupervised learning
 - Gaussian graphical model

Homework

3

KYOTO UNIVERSITY

Homework: Supervised two-class classification

- Work on a supervised classification problem
 1. Implement at least one classification algorithm by yourself
 2. Use other publicly available implementations
- Participate into a competition
 - Dataset available at <http://universityofbigdata.net/?lang=en>
 - Register with the invitation code
 - Text classification competition lasting until Dec. 31th
 - Submit your predictions at least twice
(once with your implementation; once with another)

4

KYOTO UNIVERSITY

Report submission:

Submit a report summarizing your work

- Submission:
 - Due: Jan. 7th noon
 - Send to kashipong+report@gmail.com with title “SML2014 report 1” and confirm you receive an ack
- Report format:
 - Include:
 - Brief description of your implementation (not sourcecode)
 - Your approach, analysis pipeline, results, and discussions
 - Do not exceed 6 pages in LNCS format

Topics of the 2nd half of this course:

Advanced supervised learning and unsupervised learning

- ***Multi-class classification and structured output prediction***
- Other variants of supervised learning problems
 - Semi-supervised learning
 - Active learning
 - Transfer learning
- On-line learning
 - Aggregating experts’ predictions
 - Bandits and reinforcement learning
- Unsupervised learning
 - Gaussian graphical model

Multi-class Classification

7

KYOTO UNIVERSITY

Multi-class classification:

Generalization of supervised two-class classification

- Training dataset: $\{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(i)}, y^{(i)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})\}$
 - input $\mathbf{x}^{(i)} \in \mathcal{X} = \mathbb{R}^D$: D -dimensional real vector
 - output $y^{(i)} \in \mathcal{Y}$: one-dimensional scalar
- Estimate a *deterministic mapping* $f: \mathcal{X} \rightarrow \mathcal{Y}$ (often with a confidence value) or a *conditional probability* $P(y|\mathbf{x})$
- Classification
 - $\mathcal{Y} = \{+1, -1\}$: Two-class classification
 - $\mathcal{Y} = \{1, 2, \dots, K\}$: K -class classification
 - hand-written digit recognition, text classification, ...

8

KYOTO UNIVERSITY

Two-class classification model:

One model with one model parameter vector

- Two-class classification model
 - Linear classifier: $f(\mathbf{x}) = \text{sign}(\mathbf{w}^\top \mathbf{x})$
 - Logistic regression: $P(y|\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})}$
 - The model is specified by the parameter vector $\mathbf{w} = (w_1, w_2, \dots, w_D)^\top$
- Our goal is find the parameter $\hat{\mathbf{w}}$ by using the training dataset $\{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})\}$
 - Generalization: accurate prediction for future data sampled from some underlying distribution $\mathcal{D}_{\mathbf{x},y}$

Simple approaches to multi-class classification:

Reduction to two-class classification

- Reduction to a set of two-class classification problems
- Approach 1: One-versus-rest
 - Construct K two-class classifiers; each classifier $\text{sign}(\mathbf{w}^{(k)\top} \mathbf{x})$ discriminates class k from the others
 - Prediction: the most probable class with the highest $\mathbf{w}^{(k)\top} \mathbf{x}$
- Approach 2: One-versus-one
 - Construct $K(K - 1)/2$ two-class classifiers, each of which discriminates between a pair of two classes
 - Prediction by voting

Error Correcting Output Code (ECOC) :

An approach inspired by error correcting coding

- Approach 3: Error correcting output code (ECOC)
 - Construct a set of two-class classifiers, each of which discriminates between two groups of classes, e.g. AB vs. CD
 - Prediction by finding the nearest code in terms of Hamming distance

codes

class	two-class classification problems						
	1	2	3	4	5	6	
A	1	1	1	1	1	1	code for class A
B	1	-1	1	-1	-1	-1	
C	-1	-1	-1	1	-1	1	
D	-1	1	1	-1	-1	1	
prediction	1	-1	1	1	1	-1	

11

KYOTO UNIVERSITY

Design of ECOC :

Code design is the key for good classification

- Codes (row) should be apart from each other in terms of Hamming distance

codes

class	two-class classification problems					
	1	2	3	4	5	6
A	1	1	1	1	1	1
B	1	-1	1	-1	-1	-1
C	-1	-1	-1	1	-1	1
D	-1	1	1	-1	-1	1

Hamming distances between codes

class	A	B	C	D
A	0	4	4	3
B		0	4	3
C			0	3
D				0

12

KYOTO UNIVERSITY

Multi-class classification model:

One model parameter vector for each class

- More direct modeling of multi-class classification
 - One parameter vector $\mathbf{w}^{(k)}$ for each class k
 - Multi-class linear classifier: $f(\mathbf{x}) = \operatorname{argmax}_{k \in \mathcal{Y}} \mathbf{w}^{(k)\top} \mathbf{x}$
 - Multi-class logistic regression: $P(k|\mathbf{x}) = \frac{\exp(\mathbf{w}^{(k)\top} \mathbf{x})}{\sum_{k' \in \mathcal{Y}} \exp(\mathbf{w}^{(k')\top} \mathbf{x})}$
 - converts real values into positive values, and then normalize them to obtain a probability $\in [0,1]$

Training multi-class classifier:

Constraints for correct classification

- Training multiclass linear classifier: $f(\mathbf{x}) = \operatorname{argmax}_{k \in \mathcal{Y}} \mathbf{w}^{(k)\top} \mathbf{x}$
 - One-versus-rest is O.K. but not perfect
- Constraints for correct classification of training data
$$\mathbf{w}^{(y^{(i)})\top} \mathbf{x}^{(i)} > \mathbf{w}^{(k)\top} \mathbf{x}^{(i)} \text{ for } \forall k \neq y^{(i)}$$

i.e. $\mathbf{w}^{(y^{(i)})\top} \mathbf{x}^{(i)} > \operatorname{argmax}_{k \in \mathcal{Y}, k \neq y^{(i)}} \mathbf{w}^{(k)\top} \mathbf{x}^{(i)}$

 - Learning algorithms find solutions satisfying (almost all) these constraints
 - Multi-class perceptron, multi-class SVM, ...

Multi-class perceptron:

Incremental learning algorithm of linear classifier

- Multi-class linear perceptron trains a classifier to meet the

$$\text{constraints } \mathbf{w}^{(y^{(i)})\top} \mathbf{x}^{(i)} > \max_{k \in \mathcal{Y}, y \neq y^{(i)}} \mathbf{w}^{(k)\top} \mathbf{x}^{(i)}$$

- Algorithm:

- Given $(\mathbf{x}^{(i)}, y^{(i)})$, make a prediction with :

$$f(\mathbf{x}^{(i)}) = \operatorname{argmax}_{k \in \mathcal{Y}} \mathbf{w}^{(k)\top} \mathbf{x}^{(i)}$$

- Update parameters only when the prediction is wrong:

- $\mathbf{w}^{(y^{(i)})} \leftarrow \mathbf{w}^{(y^{(i)})} + \mathbf{x}^{(i)}$: reinforce correct prediction

- $\mathbf{w}^{(f(\mathbf{x}^{(i)}))} \leftarrow \mathbf{w}^{(f(\mathbf{x}^{(i)}))} - \mathbf{x}^{(i)}$: discourage wrong prediction

15

KYOTO UNIVERSITY

Training multi-class logistic regression:

(Regularized) maximum likelihood estimation

- Find the parameters that minimizes the negative log-likelihood

$$J(\{\mathbf{w}^{(y)}\}_y) = - \sum_{i=1, \dots, N} \log p(y^{(i)} | \mathbf{x}^{(i)}) + \gamma \sum_{y \in \mathcal{Y}} \|\mathbf{w}^{(y)}\|_2^2$$

– $\|\mathbf{w}^{(y)}\|_2^2$: a regularizer to avoid overfitting

- For multi-class logistic regression $P(k|\mathbf{x}) = \frac{\exp(\mathbf{w}^{(k)\top} \mathbf{x})}{\sum_{k' \in \mathcal{Y}} \exp(\mathbf{w}^{(k')\top} \mathbf{x})}$

$$J = - \sum_i \mathbf{w}^{(k)\top} \mathbf{x}^{(i)} + \sum_i \log \sum_{k' \in \mathcal{Y}} \exp(\mathbf{w}^{(k)\top} \mathbf{x}^{(i)}) + \text{reg.}$$

– Minimization using gradient-based optimization methods

16

KYOTO UNIVERSITY

Difference of perceptron and ML estimation:
Perceptron needs only max operation; ML needs sum

- Perceptron
 - Training & prediction need only argmax operation $\underset{k \in \mathcal{Y}}{\operatorname{argmax}}$
 - SVM also does
- (Regularized) maximum likelihood estimation
 - Training: needs $\sum_{k' \in \mathcal{Y}}$ operation
 - Prediction: needs argmax operation $\underset{k \in \mathcal{Y}}{\operatorname{argmax}}$

17

KYOTO UNIVERSITY

Equivalent form of multi-class logistic regression:
Representation with one (huge) parameter vector

- Consider a joint feature space of \mathbf{x} and y :
 - $\boldsymbol{\varphi}(\mathbf{x}, y) = (\delta(y = 1)\mathbf{x}^\top, \delta(y = 2)\mathbf{x}^\top, \dots, \delta(y = K)\mathbf{x}^\top)^\top$
 - Corresponding parameter vector:
$$\mathbf{w} = (\mathbf{w}^{(1)\top}, \mathbf{w}^{(2)\top}, \dots, \mathbf{w}^{(K)\top})^\top$$
 - KD -dimensional feature space
- Multiclass LR model: $P(y|\mathbf{x}) = \frac{\exp(\mathbf{w}^\top \boldsymbol{\varphi}(\mathbf{x}, y))}{\sum_{k' \in \mathcal{Y}} \exp(\boldsymbol{\varphi}(\mathbf{x}, k'))}$
 - Equivalent to the previous model
 - Useful when we consider structured output prediction

18

KYOTO UNIVERSITY

Structured Output Prediction

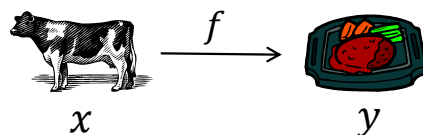
19

KYOTO UNIVERSITY

Ultimate predictive modeling:

Learn a mapping between general sets

- In supervised learning, what we want is a mapping $f: \mathcal{X} \rightarrow \mathcal{Y}$
 - $\mathcal{X} = \mathbb{R}^D$, $\mathcal{Y} = \mathbb{R}$ (regression) or a discrete set (classification)
- Ultimate predictor should take arbitrary \mathcal{X} and \mathcal{Y} sets
 - Indeed we restrict the classes of \mathcal{X} and \mathcal{Y}



- Especially cases with general output spaces are difficult to consider in the current framework
 - Classification with an infinite number of classes

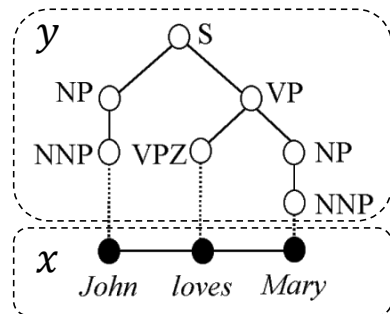
20

KYOTO UNIVERSITY

Structured output prediction: Outputs are sequences, trees, and graphs

- (Inputs and) outputs have complex structures such as sequences, trees, and graphs in many applications
 - Natural language processing: texts, parse trees, ...
 - Bioinformatics: sequences and structures of DNA/RNA/proteins

- Structured output prediction tasks:
 - Syntactic parsing: sequences to trees
 - $x = (\text{John}, \text{loves}, \text{Mary})$: sequence
 - $y = (\text{S}(\text{NP}(\text{NNP}))(\text{VP}(\text{VPZ})(\text{NP}(\text{NNP}))))$: tree



21

<http://en.wikipedia.org/wiki/Treebank#mediaviewer/File:Example-tree.png>

KYOTO UNIVERSITY

Sequence labeling: Structured prediction with sequential input & output

- Sequence labeling gives a label to each element of a sequence
 - $x = (x_1, x_2, \dots, x_T)$: input sequence of length T
 - $y = (y_1, y_2, \dots, y_T)$: output sequence with the same length
 - Simplest structured prediction problem

x_1	x_2	...	x_T
y_1	y_2	...	y_T

- Example. Part-of-speech tagging gives a part-of-speech tag to each word in a sentence
 - x : sentence (a sequence of words)
 - y : Part-of-speech tags (e.g. *noun*, *verb*,...)

22

KYOTO UNIVERSITY

Sequence labeling as multi-class classification: Impossible to work with exponentially many parameters

- Formulation as T independent classification problems
 - Predict y_t using surrounding words $(\dots, x_{t-1}, x_t, x_{t+1}, \dots)$
 - Sometimes quite works well and efficient
 - No guarantee of consistence among predicted labels
 - Might want to include dependencies among labels such as “a verb is likely to follow nouns”
- This problem can also be considered as one multi-class classification problem with K^T classes
 - $f(x) = \operatorname{argmax}_{k \in \mathcal{Y}} \mathbf{w}^{(k)\top} \mathbf{x}$ is almost impossible to work with exponentially many parameters

23

KYOTO UNIVERSITY

Key for solving structured output prediction: Formulation as a validation problem of in/output pairs

- Remember another form of multi-class classifier using the joint feature space
 - $P(y|x) = \frac{\exp(\mathbf{w}^\top \boldsymbol{\varphi}(x, y))}{\sum_{k' \in \mathcal{Y}} \exp(\boldsymbol{\varphi}(x, k'))}$ or $f(x) = \operatorname{argmax}_{y \in \mathcal{Y}} \mathbf{w}^\top \boldsymbol{\varphi}(x, y)$
 - They evaluate the affinity of a input-output pair
- Still the problem is not solved....
but we can consider reducing the dimensionality of $\boldsymbol{\varphi}(x, y)$
 - Because the dimensionality of $\boldsymbol{\varphi}(x, y)$ is still huge

24

KYOTO UNIVERSITY

Features for sequence labeling:

First-order Markov assumption gives two feature types

- Two types of features for sequence labeling
 - Combination of one input label x_t and one output label y_t
 - Standard feature for multi-class classification
 - e.g. $x_t = \text{"loves"} \wedge y_t = \text{"verb"}$
 - Combination of two consecutive labels y_{t-1} and y_t
 - Markov assumption of output labels
 - e.g. $y_{t-1} = \text{"noun"} \wedge y_t = \text{"verb"}$

x_1	x_2	...	x_{t-1}	x_t	...	x_T
y_1	y_2	...	y_{t-1}	y_t	...	y_T

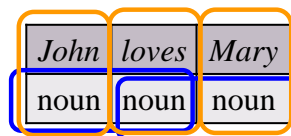
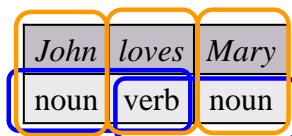
25

KYOTO UNIVERSITY

Feature vector definition:

The numbers of appearance of each pattern

- Each dimension of $\boldsymbol{\varphi}(x, y)$ is defined as the number of appearance of each pattern in the joint sequence (x, y) , e.g.
 - $-\varphi(x, y)_1 = \text{\#appearance of } [x_t = \text{"loves"} \wedge y_t = \text{"verb"}]$
 - $-\varphi(x, y)_2 = \text{\#appearance of } [y_{t-1} = \text{"noun"} \wedge y_t = \text{"verb"}]$
 - Features for all possible combination of POS tags and words



26

KYOTO UNIVERSITY

Impact of first-order Markov assumption: Reduced dimensionality of feature space

- Dimensionality of a feature vector was decreased from $O(K^T)$ to $O(K^2)$ (K is the number of labels for each position)
- Space problem was solved; we can calculate $\mathbf{w}^\top \boldsymbol{\varphi}(x, y)$
 - Prediction problem (i.e. $\operatorname{argmax}_{y \in \mathcal{Y}} \mathbf{w}^\top \boldsymbol{\varphi}(x, y)$) has not been solved
 - For sequential labeling, this can be done by using dynamic programming

27

KYOTO UNIVERSITY

Structured perceptron : Simple structured output learning algorithm

- Structured perceptron learns \mathbf{w} satisfying
$$\mathbf{w}^\top \boldsymbol{\varphi}(x^{(i)}, y^{(i)}) > \max_{y \in \mathcal{Y}, y \neq y^{(i)}} \mathbf{w}^\top \boldsymbol{\varphi}(x^{(i)}, y)$$

- Algorithm:

1. Given $(x^{(i)}, y^{(i)})$, make a prediction with :

$$f(x^{(i)}) = \operatorname{argmax}_{y \in \mathcal{Y}} \mathbf{w}^\top \boldsymbol{\varphi}(x^{(i)}, y)$$

2. Update parameters only when the prediction is wrong

$$\mathbf{w}^{\text{NEW}} \leftarrow \mathbf{w}^{\text{OLD}} + \boldsymbol{\varphi}(x^{(i)}, y^{(i)}) - \boldsymbol{\varphi}(x^{(i)}, f(x^{(i)}))$$

– Prediction can be done in polynomial time by using dynamic programming for sequence labeling

28

KYOTO UNIVERSITY

Conditional random field:

Conditional probabilistic model for structured prediction

- Conditional random field: conditional probabilistic model

$$P(y|x) = \frac{\exp(\mathbf{w}^\top \boldsymbol{\varphi}(x, y))}{\sum_{k' \in \mathcal{Y}} \exp(\boldsymbol{\varphi}(x, k'))}$$

- ML estimation needs the sum over all possible outputs

$$J = \sum_i \mathbf{w}^\top \boldsymbol{\varphi}(x^{(i)}, y^{(i)}) - \sum_i \log \sum_{y \in \mathcal{Y}} \exp(\mathbf{w}^\top \boldsymbol{\varphi}(x^{(i)}, y)) + \text{reg.}$$

- The sum can be taken with dynamic programming

Perceptron vs. CRF:

Perceptron needs only max operation; ML needs sum

- Just like in multi-class classification,
 - Structured perceptron can work only with argmax operation
 - Maximum likelihood estimation also needs sum operation
- There are some structured output problems where argmax operation is easy but sum operation is difficult
 - e.g. bipartite matching