

<http://goo.gl/Jv7Vj9>

Course website

KYOTO UNIVERSITY

Statistical Machine Learning Theory

From Multi-class Classification to Structured Output Prediction

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

DEPARTMENT OF INTELLIGENCE SCIENCE
AND TECHNOLOGY

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:
 - Follow the leader, on-line gradient descent, perceptron
 - Regret analysis
- Sparse modeling:
 - L_1 regularization, Lasso, & reduced rank regression
- Model evaluation

Homework

Homework:

Supervised regression

- Work on a supervised regression problem:
 1. Implement at least one linear regression by yourself
 2. Use publicly available nonlinear regression implementations
- Participate into a competition at <http://universityofbigdata.net>
 - Register with your Google account (if you have not yet)



- The competition will last until Dec. 31th
- Submit your predictions at least twice
(once with your implementation; once with another)

Submitting your prediction: <http://goo.gl/3BMpf4>

- See the instructions at <http://universityofbigdata.net/competition/5757948332605440?lang=en>

The screenshot shows the submission interface for a competition. On the left, there is a sidebar with the competition title and details. The main content area is divided into three sections: Submission, Note (optional), and Intermediate ranking.

Submission

管理者アカウントには提出回数制限はありません。
You can upload a file of up to 20MB. You can compress your submission using the .zip compression format.

Note (optional)

You can add a note to your submission. Notes are shown in the bottom of this page and only you can see your note.

Intermediate ranking

Intermediate rank	Nickname	Intermediate score
1	University of Big Data	0.0240

This leaderboard is calculated on the latest submissions.
The intermediate scores are calculated using 50% of the test dataset, and the final scores are calculated using the other 50%.
Final ranks are determined according to the final scores.

Report submission:

Submit a report summarizing your work

- Submission:

- Due: Jan. 7th noon, 2016

- Send your report to kashipong+report@gmail.com with subject “SML2015 competition report” and confirm you receive an ack before 8th

- Report format:

- Must include:

- Brief description of your implementation (not sourcecode)
 - Your approach, analysis pipeline, results, and discussions

- At least 3 pages, but 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:
 - Follow the leader, on-line gradient descent, perceptron
 - Regret analysis
- Sparse modeling:
 - L_1 regularization, Lasso, & reduced rank regression
- Model evaluation

Multi-class Classification

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, ...

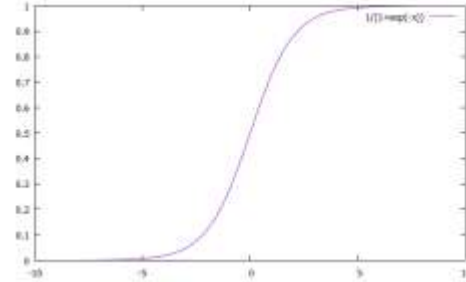
Two-class classification model:

One model with one model parameter vector

- Two-class classification model

- Linear classifier: $f(\mathbf{x}) = \text{sign}(\mathbf{w}^\top \mathbf{x}) \in \{+1, -1\}$

- 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



confidence

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
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

code for class A

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

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 normalizes them to obtain a probability value $\in [0,1]$

Training multi-class classifier:

Constraints for correct classification

- Training multiclass linear classifier: $f(\mathbf{x}) = \operatorname{argmax}_{k \in \mathcal{Y}} \mathbf{w}^{(k)\top} \mathbf{x}$
 - can use the one-versus-rest method, 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

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:

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

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

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

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

- Perceptron

- Training & prediction need only $\operatorname{argmax}_{k \in \mathcal{Y}}$ operation
- SVM also does

- (Regularized) maximum likelihood estimation

- Training: needs $\sum_{k' \in \mathcal{Y}}$ operation
- Prediction: needs $\operatorname{argmax}_{k \in \mathcal{Y}}$ operation

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 $P(k|\mathbf{x}) = \frac{\exp(\mathbf{w}^{(k)\top} \mathbf{x})}{\sum_{k' \in \mathcal{Y}} \exp(\mathbf{w}^{(k')\top} \mathbf{x})}$
 - Useful when we consider structured output prediction

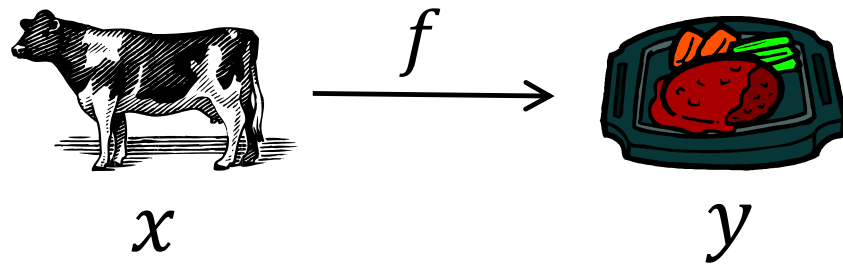
Structured Output Prediction

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

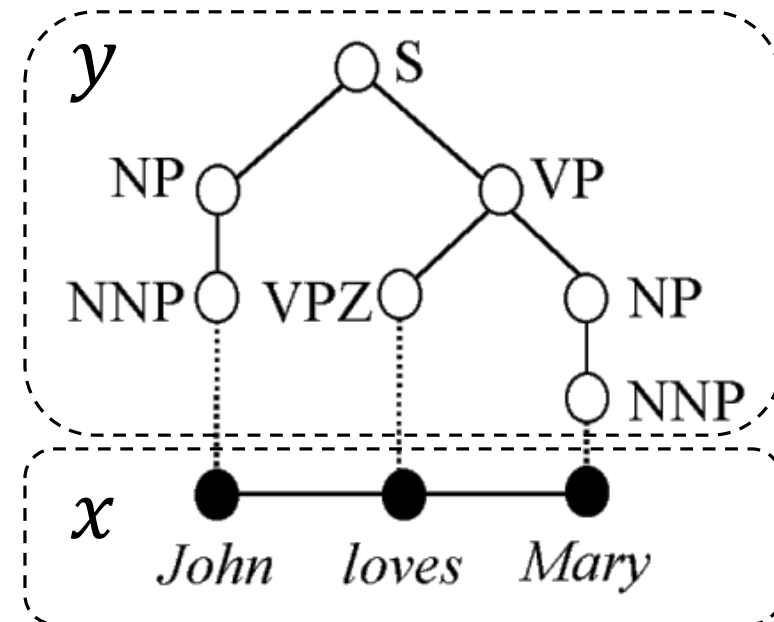


- We have to restrict the classes of \mathcal{X} and \mathcal{Y} in practice
 - Especially, cases with general output spaces are difficult to consider in the current framework
 - Classification with an infinite number of classes

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
 - $x = (\textit{John}, \textit{loves}, \textit{Mary})$: sequence
 - $y = (\text{S}(\text{NP}(\text{NNP}))(\text{VP}(\text{VPZ})(\text{NP}(\text{NNP}))))$: tree



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*,...)

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

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

$$- P(y|x) = \frac{\exp(\mathbf{w}^\top \boldsymbol{\varphi}(x,y))}{\sum_{k' \in \mathcal{Y}} \exp(\boldsymbol{\varphi}(x,k'))} \text{ or } f(x) = \operatorname{argmax}_{y \in \mathcal{Y}} \mathbf{w}^\top \boldsymbol{\varphi}(x,y)$$

– They evaluate the affinity of an 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

Features for sequence labeling:

First-order Markov assumption gives two feature types

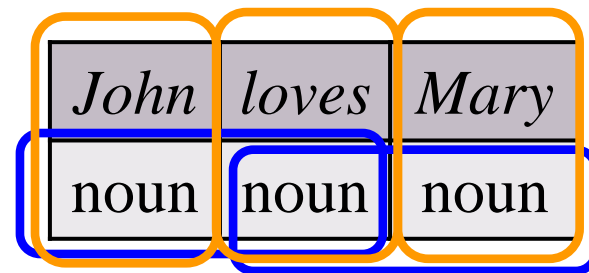
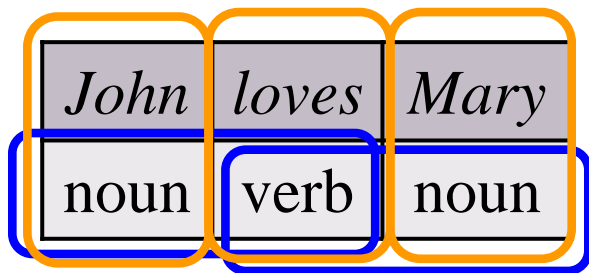
- Two types of features for sequence labeling
 1. Combination of one input label x_t and one output label y_t
 - Standard feature for multi-class classification
 - e.g. $x_t = \text{“loves”} \wedge y_t = \text{“verb”}$
 2. Combination of two consecutive labels y_{t-1} and y_t
 - Markov assumption of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

Feature vector definition:

The numbers of appearance of each pattern

- Each dimension of $\varphi(x, y)$ 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



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

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:

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

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

- 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

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