Transfer Learning for Text Categorization

Chuong B. Do and Andrew Y. Ng Presenter: Lei Tang

Department of Computer Science & Engineering Arizona State University

6th October 2006

() < </p>

- also Meta Learning, Multi-task Learning, Structural Learning
- For linear classifiers, past work(naive Bayes, TFIDF) focus on learning the weight (actually try to identify better functions mapping from statistics to weight)
- This work focus on learning the MAPPING FUNCTION from related classification problems

・ 同 ・ ・ ヨ ・ ・ ヨ ・

Linear Classifier

• Linear Classifier: $f_k(x) = \sum \theta_{ki} x_i$, and

$$y = \arg \max_k f_k(x)$$

where x_i is usually a frequency corresponding to term i.Naïve Bayes:

$$f_k^{NB}(x) = \log \hat{p}(y = k) + \sum_{i=1}^n x_i \log \hat{p}(w_i | y = k)$$

Rocchio Alorithm

$$f_k^{Rocchio} = \sum_{i=1}^n (\bar{x}_i | y = k \cdot \log idf) (x_i \cdot \log idf)$$

臣

Mapping Function(1)

• Linear Classifier: $f_k(x) = \sum \theta_{ki} x_i$, and where x_i is usually a frequency corresponding to term *i*. We can rewrite $f_k(x) = \sum \theta_{ki} x_i$ as

$$f_k(x) = \sum g(\mathbf{u}_{ki}) x_i$$

where $\mathbf{u}_k i$ is a vector of some statistics computed from the training set. (Similar to the concept of *sufficient statistics*)

۲

$$\mathbf{u}_{ki} = \begin{bmatrix} u_{ki1} \\ u_{ki2} \\ u_{ki3} \\ u_{ki4} \\ u_{ki5} \end{bmatrix} = \begin{bmatrix} \#w_i \text{ appear in documents of class } k \\ \#documents \text{ of class } k \text{ containing } w_i \\ \#total \text{ words in documents of class } k \\ \#documents \text{ of class } k \\ \#documents \text{ of class } k \\ \#total \text{ documents} \end{bmatrix}$$

A (1) > A (2) > A (2) >

Mapping Function(2)

$$f_k^{NB}(x) = \log \hat{p}(y=k) + \sum_{i=1}^n x_i \log \hat{p}(w_i|y=k)$$

 $\mathbf{u}_{ki} = \begin{bmatrix} u_{ki1} \\ u_{ki2} \\ u_{ki3} \\ u_{ki4} \\ u_{ki5} \end{bmatrix} = \begin{bmatrix} \#w_i \text{ appear in documents of class } k \\ \#documents \text{ of class } k \text{ containing } w_i \\ \#total \text{ words in documents of class } k \\ \#documents \text{ of class } k \\ \#documents \text{ of class } k \\ \#total \text{ documents} \end{bmatrix}$

$$g_{NB}(\mathbf{u}_{ki}) = \log \frac{u_{ki1} + \varepsilon}{u_{ki3} + n\varepsilon}$$

(1日) (1日) (1日)

Similarly, TFIDF can also be written as linear combination of $g(u_{ki}) \cdot x_i$

- Why not learn the mapping function g automatically?
- The authors assume the mapping function to be linear (Can be extended later to nonlinear case by kernel trick)

$$g(\mathbf{u}_{ki}) = \beta^T \mathbf{u}_{ki}$$

• Reformulate the logistic regression:

$$p(y = k | \mathbf{x}; \{\theta_{ki}\}) := \frac{\exp(\sum_{i} \theta_{ki} x_{i})}{\sum_{k'} \exp(\sum_{i} \theta_{k'i} x_{i})}$$
$$p(y = k | \mathbf{x}; \beta) = \frac{\exp(\sum_{i} \beta^{T} \mathbf{u}_{ki} x_{i})}{\sum_{k'} \exp(\sum_{i} \beta^{T} \mathbf{u}_{k'i} x_{i})}$$

To maximize the log likelihood

$$\ell(\beta:\Omega) = \sum_{i=1}^{m} \log p(y^{(j)}|x^{(j)};\beta) - C \|\beta\|^2$$

Chuong B. Do and Andrew Y. NgPresenter: Lei Tang

- Why not learn the mapping function g automatically?
- The authors assume the mapping function to be linear (Can be extended later to nonlinear case by kernel trick)

$$g(\mathbf{u}_{ki}) = \beta^T \mathbf{u}_{ki}$$

• Reformulate the logistic regression:

$$p(y = k | \mathbf{x}; \{\theta_{ki}\}) := \frac{\exp(\sum_{i} \theta_{ki} x_{i})}{\sum_{k'} \exp(\sum_{i} \theta_{k'i} x_{i})}$$
$$p(y = k | \mathbf{x}; \beta) = \frac{\exp(\sum_{i} \beta^{T} \mathbf{u}_{ki} x_{i})}{\sum_{k'} \exp(\sum_{i} \beta^{T} \mathbf{u}_{k'i} x_{i})}$$

To maximize the log likelihood

$$\ell(\beta:\Omega) = \sum_{i=1}^{m} \log p(y^{(j)}|x^{(j)};\beta) - C \|\beta\|^2$$

Chuong B. Do and Andrew Y. NgPresenter: Lei Tang

- Why not learn the mapping function g automatically?
- The authors assume the mapping function to be linear (Can be extended later to nonlinear case by kernel trick)

$$g(\mathbf{u}_{ki}) = \beta^T \mathbf{u}_{ki}$$

• Reformulate the logistic regression:

$$p(y = k | \mathbf{x}; \{\theta_{ki}\}) := \frac{\exp(\sum_{i} \theta_{ki} x_{i})}{\sum_{k'} \exp(\sum_{i} \theta_{k'i} x_{i})}$$
$$p(y = k | \mathbf{x}; \beta) = \frac{\exp(\sum_{i} \beta^{T} \mathbf{u}_{ki} x_{i})}{\sum_{k'} \exp(\sum_{i} \beta^{T} \mathbf{u}_{k'i} x_{i})}$$

To maximize the log likelihood

$$\ell(\beta:\Omega) = \sum_{i=1}^{m} \log p(y^{(j)}|x^{(j)};\beta) - C \|\beta\|^2$$

Chuong B. Do and Andrew Y. NgPresenter: Lei Tang

- Why not learn the mapping function g automatically?
- The authors assume the mapping function to be linear (Can be extended later to nonlinear case by kernel trick)

$$g(\mathbf{u}_{ki}) = \beta^T \mathbf{u}_{ki}$$

• Reformulate the logistic regression:

$$p(y = k | \mathbf{x}; \{\theta_{ki}\}) := \frac{\exp(\sum_{i} \theta_{ki} x_{i})}{\sum_{k'} \exp(\sum_{i} \theta_{k'i} x_{i})}$$
$$p(y = k | \mathbf{x}; \beta) = \frac{\exp(\sum_{i} \beta^{T} \mathbf{u}_{ki} x_{i})}{\sum_{k'} \exp(\sum_{i} \beta^{T} \mathbf{u}_{k'i} x_{i})}$$

To maximize the log likelihood

$$\ell(\beta:\Omega) = \sum_{i=1}^{m} \log p(y^{(j)}|x^{(j)};\beta) - C \|\beta\|^2$$

Chuong B. Do and Andrew Y. NgPresenter: Lei Tang

- Why not learn the mapping function g automatically?
- The authors assume the mapping function to be linear (Can be extended later to nonlinear case by kernel trick)

$$g(\mathbf{u}_{ki}) = \beta^T \mathbf{u}_{ki}$$

• Reformulate the logistic regression:

$$p(y = k | \mathbf{x}; \{\theta_{ki}\}) := \frac{\exp(\sum_{i} \theta_{ki} x_{i})}{\sum_{k'} \exp(\sum_{i} \theta_{k'i} x_{i})}$$
$$p(y = k | \mathbf{x}; \beta) = \frac{\exp(\sum_{i} \beta^{T} \mathbf{u}_{ki} x_{i})}{\sum_{k'} \exp(\sum_{i} \beta^{T} \mathbf{u}_{k'i} x_{i})}$$

To maximize the log likelihood

$$\ell(\beta:\Omega) = \sum_{i=1}^m \log p(y^{(j)}|x^{(j)};\beta) - C \|\beta\|^2$$

Chuong B. Do and Andrew Y. NgPresenter: Lei Tang

$$\beta^* = \sum_{j=1}^m \sum_k \alpha_{jk}^* \sum_i \mathbf{u}_{ki}^{(j)} x_i^{(j)}$$

User kernel trick, we can extend the mapping function g to nonlinear cases(Details skipped here)

(1日) (日) (日)

How to evaluate?

- Typical cross validation
- dmoz with 16 top-level categories: Test on one category based on built 450 classification problems from the other 15 categories
- Four corpora: learned from one corpora and test on the other three corpora
- 2 Results:

Better than NBC, logistic regression, 1-vs-all SVM, MC-SVM

(4月) (4日) (4日)

- Each learning task is 10-class with 2 instances in each category for training, and 1 instance in each category for testing. Not very surprising that it outperforms SVM.
- O Why regularization?

Why use logistic regression formulation?

Because logistic regression can model any decision boundaries?

Transfer Learning: Learn how to learn!

A new machine learning problem?

- Each learning task is 10-class with 2 instances in each category for training, and 1 instance in each category for testing. Not very surprising that it outperforms SVM.
- Why regularization?

Why use logistic regression formulation?

Because logistic regression can model any decision boundaries?

Transfer Learning: Learn how to learn!

A new machine learning problem?

- Each learning task is 10-class with 2 instances in each category for training, and 1 instance in each category for testing. Not very surprising that it outperforms SVM.
- Why regularization?

Why use logistic regression formulation?

Because logistic regression can model any decision boundaries?

Transfer Learning: Learn how to learn!

A new machine learning problem?

- Each learning task is 10-class with 2 instances in each category for training, and 1 instance in each category for testing. Not very surprising that it outperforms SVM.
- Why regularization?

Why use logistic regression formulation?

Because logistic regression can model any decision boundaries?

Transfer Learning: Learn how to learn!

A new machine learning problem?

A (1) > A (2) > A (2) >