

# Active Learning for Cross-domain Sentiment Classification



Shoushan Li, Yunxia Xue, Zhongqing Wang, and Guodong Zhou  
 Natural Language Processing Lab, School of Computer Science and Technology  
 Soochow University, Suzhou, China



## Introduction

● **Sentiment Classification** is a task of determining the sentimental orientation of a given textual document.

● **Domain Adaptation** use some labeled data from the source domain and a large amount of unlabeled data from the target domain to classify the target domain data.

□ **Negative transfer of cross-domain sentiment classification**, the distributions of the source and target domains become too different to make the adaptation algorithm useful.

● **Active Learning**, One possible solution to such dilemma is to annotate a small amount of good labeled data in the target domain to quickly reduce the huge difference between the two domains.

□ **Major challenges of active learning for cross-domain sentiment classification**

✓ The newly-added labeled data from the target domain may become too weak to affect the **selection decision**.

✓ The newly-added labeled data from the target domain are normally too few to quickly affect the **classification decision**.

● **We address these challenges:**

1. We use the newly-added labeled data from the target domain to **train a separate classifier** and apply it in both the sample **selection** strategy and the **classification** algorithm

2. We propose a **label propagation (LP)** - based classification algorithm, which leverages both the labeled and unlabeled data, and apply it to both the source and target classifiers.

## Symbol definition

Symbol	Definition
$L_S$	Labeled source-domain data
$L_T$	Labeled target-domain data
$f_S$	The source Classifier
$f_T$	The target Classifier
$U_T$	Unlabeled target-domain data
$\Delta L_T$	Newly-added data at each iteration
$f_{LP-S}$	The LP-based source Classifier
$f_{LP-T}$	The LP-based target Classifier
$T_T$	Test data in the target domain

## QBC-based Selection Strategy

● **Query by Committee (QBC)** is a group of active learning approaches that employ many copies of “**hypotheses**” to select an unlabeled example at which their classification predictions are maximally spread.

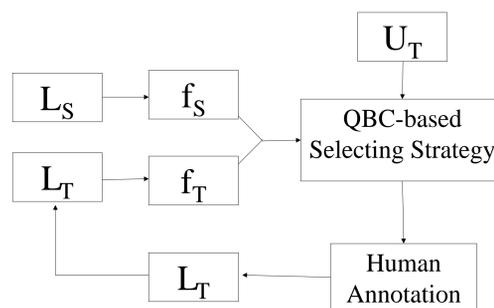
● **In our approach,**

1. We first use the source classifier and the target classifier to collaboratively **select label-disagreed samples** as the selection candidates.

2. We then **rank the label-disagreed samples** according to their uncertainty values by the source classifier.

3. We finally **select the top-N uncertainty samples** as the newly-added data for human annotation.

Figure: Sample selection in our approach



## LP-based Classification

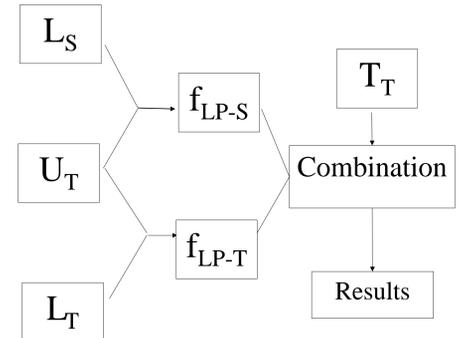
1. Construct the **document-word bipartite graph** with both the labeled data and the unlabeled data (in the target domain) and get the **transition probability matrix**.

2. Run the **LP algorithm** as shown in following figure to obtain the **labels of unlabeled data**.

3. Consider the unlabeled data with the **predicted labels** as pseudo-labeled data.

4. **Emerge** the labeled data and the pseudo-labeled data to train a classifier.

Figure: Sample classification in our approach



## Experimentation

● **Dataset:** Book, DVD, Electronics and Kitchen appliances.

✓ We randomly select 1600 instances from the source domain as **labeled data**, 1600 instances from the target domain as **unlabeled data**, and the remaining 400 instances from the target domain are reserved as **test data**.

● **Classification algorithm:** the **maximum entropy classifier** implemented with the public tool, Mallet Toolkits.

Figure: Experiment Results

