# User Manual
# For Graphical Interface of MALLET

# Contents

# 1. Introduction

## 1.1 Background

MALLET is a Java-based package for statistical natural language processing, including document classification and other machine learning applications to text. However, MALLET can only be run in cmd. That means users have to enter all complete and correct commands in cmd without any instructions, which is a really abstract way to use MALLET.

## 1.2 Purpose

For more convenience, the main purpose of this project is to design and build a graphical interface of MALLET, which is easy to input commands and select options. The most important thing is that the graphical interface is really practical for both common users and beginners.

# 2. Software Overview

## 2.1 Target

This document enables users to master the installation and use of this software.

## 2.2 Functions

### 2.2.1  Importing data

Before implementing some functions, the first step is to import data into unique MALLET format. There are two primary methods for importing data, the first is importing a directory, and the second is importing a single file.

Then the new generated MALLET file, such as web.mallet can be used in document classification.

### 2.2.2  Classification

Document classification means distinguishing between a fixed set of classes based on labeled training examples. MALLET includes several classification algorithms, such as Naïve Bayes, Maximum Entropy and Decision Tree, to train a classifier. Based on the classifier, users can apply a saved classifier to test new unlabeled data.

### 2.2.3  Sequence tagging

In natural language processing, sequence tagging means extracting information from the text, which aims to mark a sequence of words, such as a none or a non-none. Through analyzing the context, this function is used to judge the part of speech.

### 2.2.4  Topic modeling

Topic models provide a simple way to analyze large volumes of unlabeled text. That is, this

function is used to extract topics from a number of articles through the analysis of probability distribution of words. A 'topic' consists of a cluster of words that frequently occur together. Using contextual clues, topic models can connect words with similar meanings and distinguish between uses of words with multiple meanings. For example, in the article which is related to medical care, the presence of the probability of the words, such as medicine, cure, are very high. After importing text files into unique MALLET format, the topic models can be built. Besides, a topic inference tool can also be used based on the current, trained model.

### 2.2.5  Optimization

MALLET includes methods for numerical optimizing functions. The primary use for numerical optimization in machine learning is to find parameters that maximize a log-likelihood function given observed data. This software gives a demo of putting two parameters to achieve the function.

### 2.2.6  Graphical models

GRMM is short for Graphical Models in MALLET which is an add-on package to MALLET and contains support for performing inference and learning in graphical models of arbitrary structure. GRMM supports arbitrary factor graphs and it includes efficient implementations of several inference algorithms. In this software, GRMM can be achieved by selecting model path, train path and test path. All the paths should comply with specific format.
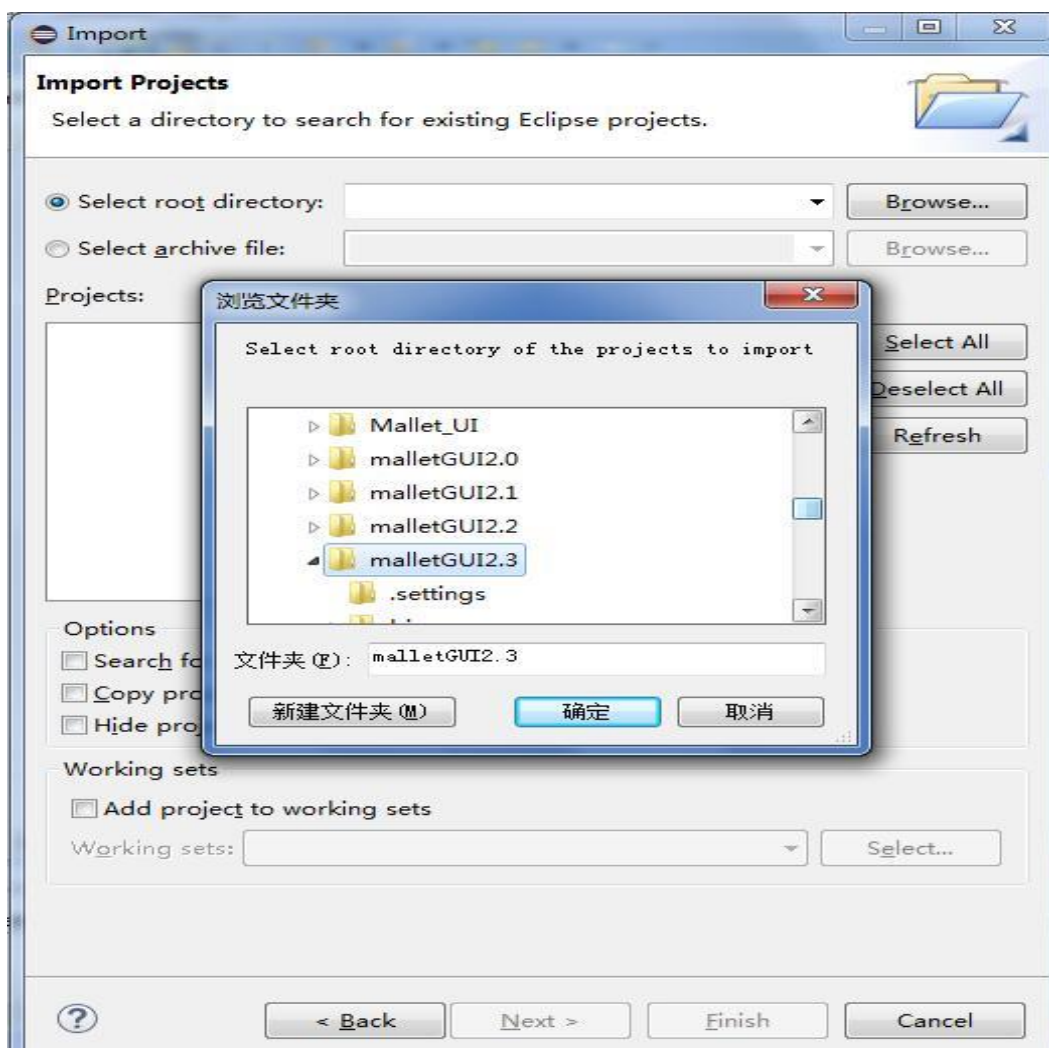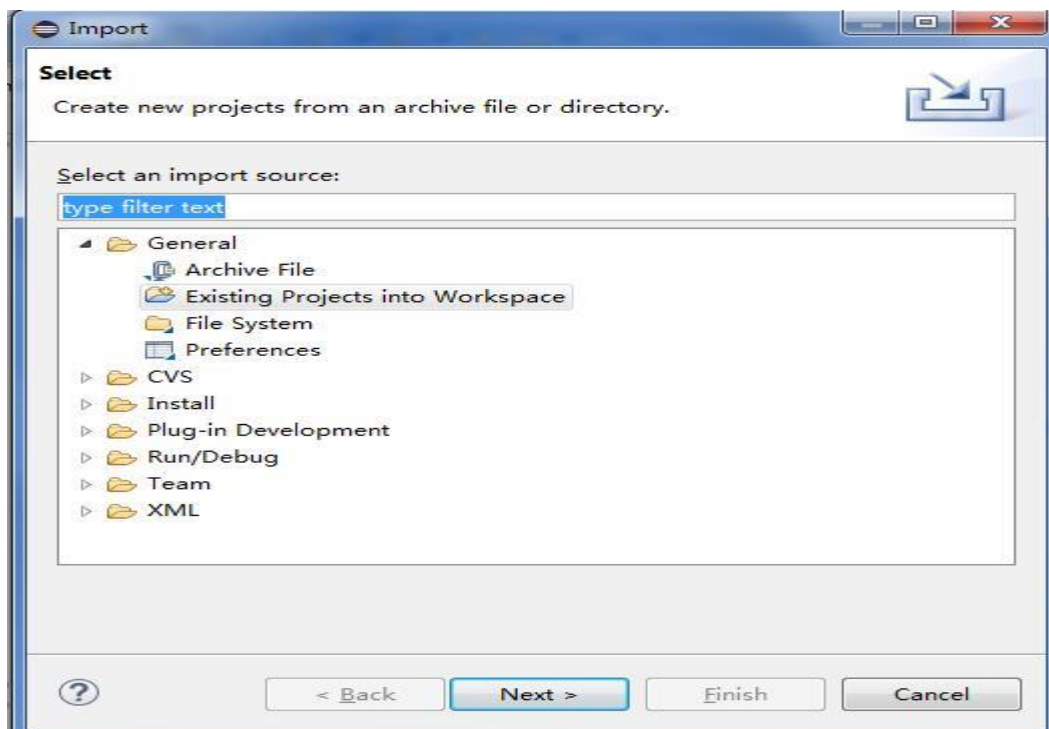
# 3. Operating Environment

This software is written in JAVA, which can be run in the eclipse.

# 4. Instructions

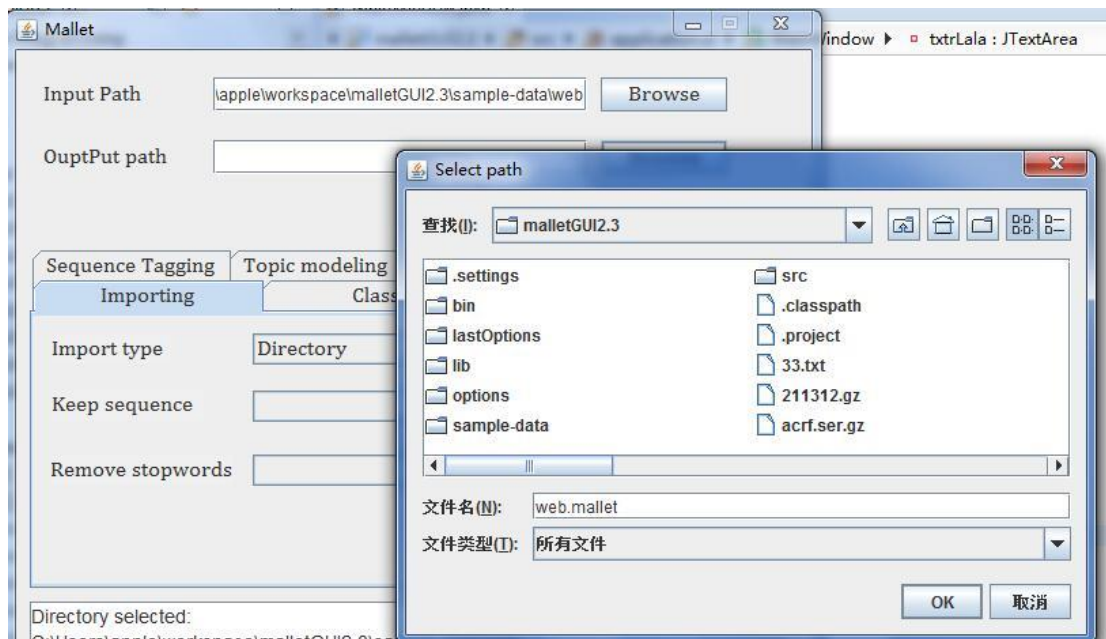## 4.1 Importing the folder into eclipse

First, put the 'malletGUI2.3' folder into the workspace, and import in the eclipse by clicking 'File', 'Import'. Then click 'Existing Projects into Workspace' ➔ 'Select root directory' can be chosen by the button 'Browse'. Select 'malletGUI2.3' in the folder 'workspace. Finally, click the button 'confirm', then the project is imported. In 'src', choose the package which is called 'application.ui', and then run the 'MainWindow.java'. After that, a graphical interface will be displayed.
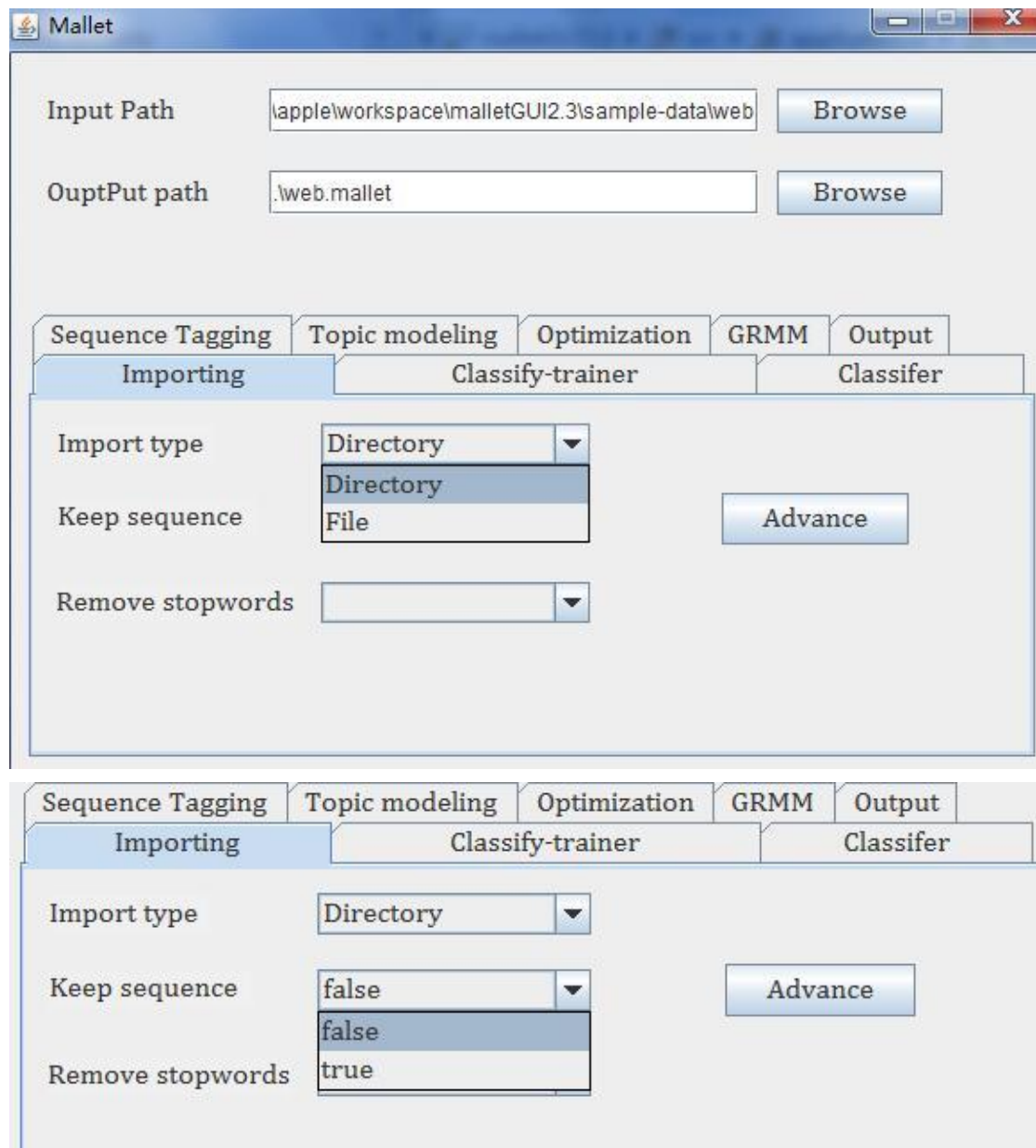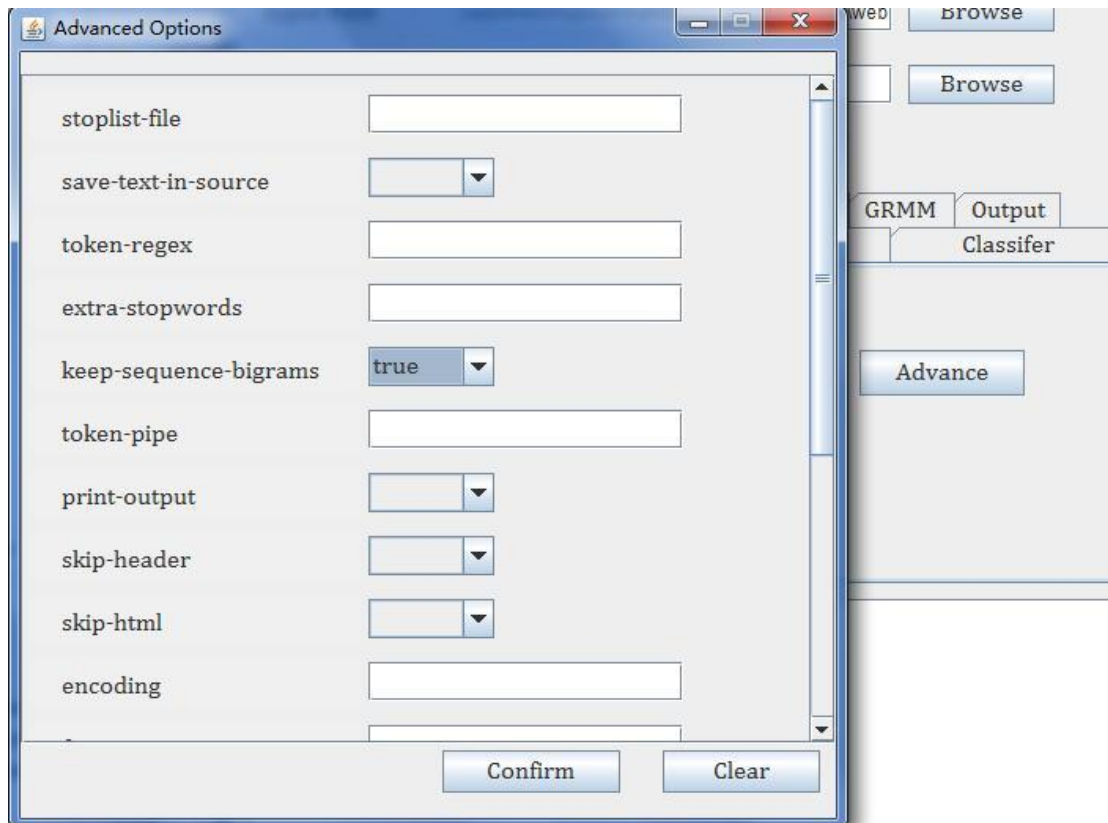
## 4.2 Importing data

At the top of the graphical interface, there are *input path* and *output path* that can be chosen by clicking the button '*Browse'* or entered by the users. *Output path* should be ended in .mallet.
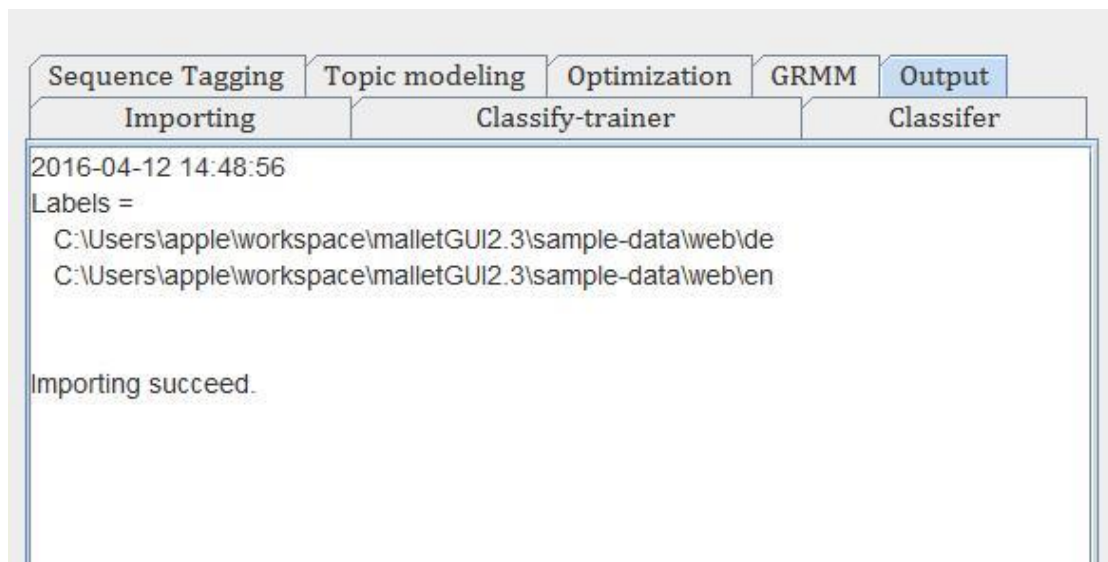


A JTabbedPane can be used to select functions, such as '*Importing'*. *Import type* must be selected depending on what users have input in the above, a *file* or a *directory*. *Keep sequence* and *Remove stopwords* are optional which can be selected through JComboBox or just leave them empty.

There are also other options can be seen by clicking the button '*Advance*'.

At the end, clicking the button 'Run', then the data is imported and the result can be displayed in the 'output'.



## 4.3 Classification

### 4.3.1 Train a classifier

As for the tag 'Classify-trainer', users can input data which is ended in .mallet and output data which is ended in .classifier.

The classification algorithm (such as Naïve Bayes, Maximum Entropy, C45 and Decision Tree), *training-portion*, *trials number* and *cross-validation* are optional. Users can also leave them empty. The button '*Run*' is used to start the function and the output can be observed in '*output*' panel.

**Mallet**

| Input Path | \Users\apple\workspace\malletGUI2.3\web.mallet | Browse |
| OuptPut path | .\web.classifier | Browse |

Sequence Tagging | Topic modeling | Optimization | GRMM | Output
Importing | Classify-trainer | Classifer

| Classifier-trainer | NaiveBayes ▼ |
| training-portion | 0.8 |
| trials-number | 2 |
| cross-validation | 0 |

```
.web.mallet
File selected:
.\web.mallet
File selected:
C:\Users\apple\workspace\malletGUI2.3\web.mallet
File selected:
.\web.classifier
```

Run



Sequence Tagging | Topic modeling | Optimization | GRMM | Output
Importing | Classify-trainer | Classifer

```
Trial 1 Trainer NaiveBayesTrainer test data accuracy= 1.0

NaiveBayesTrainer
Summary. train accuracy mean = 1.0 stddev = 0.0 stderr = 0.0
Summary. test accuracy mean = 1.0 stddev = 0.0 stderr = 0.0

NaiveBayesTrainer
Summary. train accuracy mean = 1.0 stddev = 0.0 stderr = 0.0
Summary. test accuracy mean = 1.0 stddev = 0.0 stderr = 0.0

NaiveBayesTrainer
Summary. train accuracy mean = 1.0 stddev = 0.0 stderr = 0.0
Summary. test accuracy mean = 1.0 stddev = 0.0 stderr = 0.0

Train Classifier succeed.
```

Run

### 4.3.2  Forecast new data

The tag '*classifier*' is related to apply the saved classifier to test new data. Users only need to choose '*Sample Path*' and '*Classifier Path*'. Choose *file* or *directory* according to the type of new data, then clicking '*Run*'.
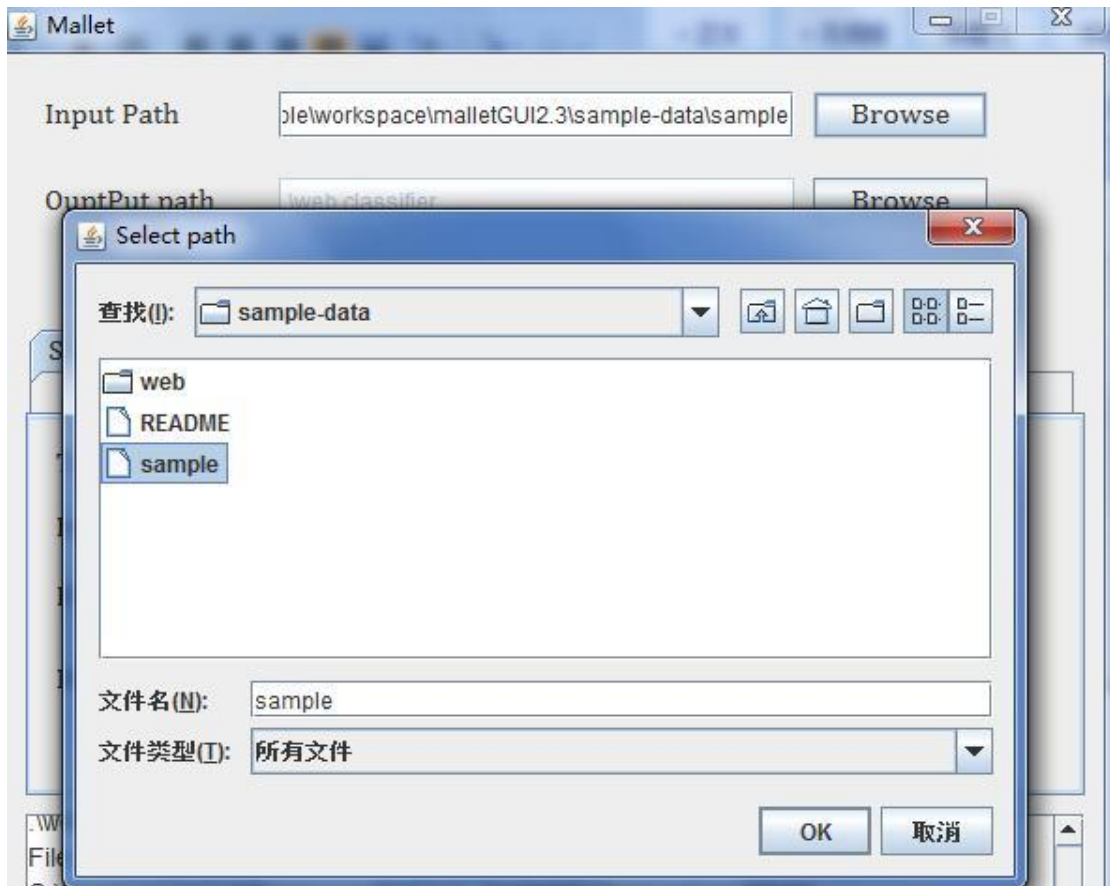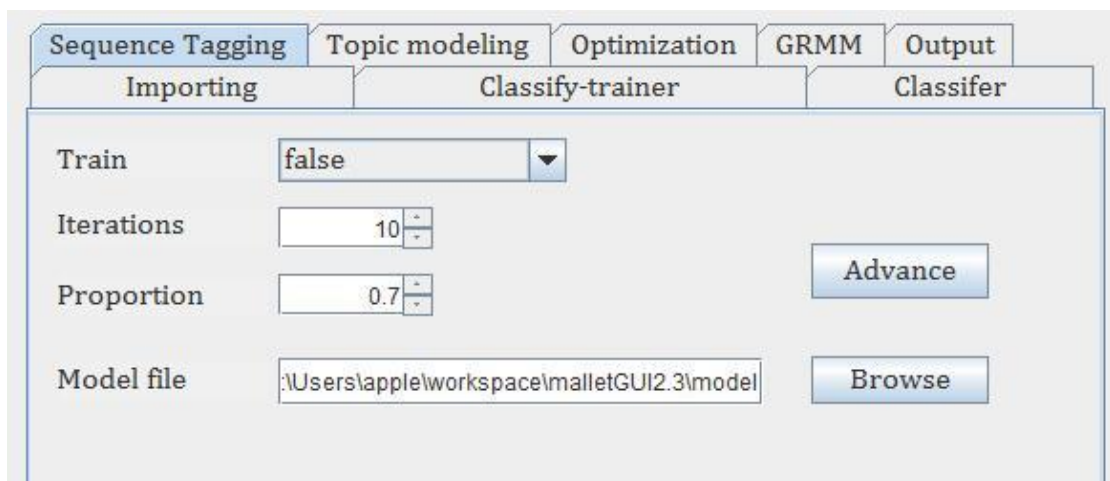

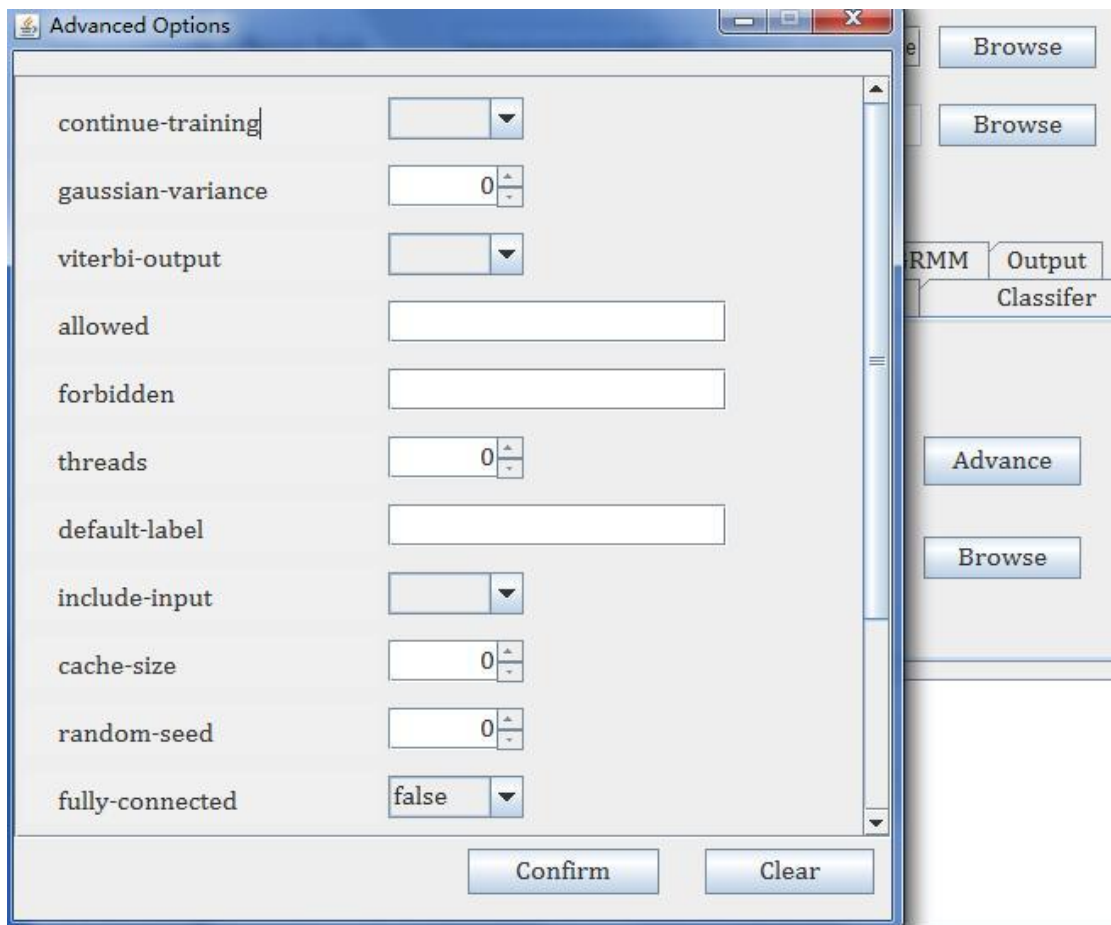
The output can be seen.



## 4.4 Sequence tagging

### 4.4.1  Train a model

The file '*sample*' which is in a unique train format is selected by clicking the button '*Browse*' to fill in the '*Input Pat*h', while the 'Output Path' does not need to be filled.

To train a sequence tagging model, the '*Train*' must be *true*. An existing model can be selected or create a new model in '*Model file*'. *Iterations* and *Proportion* are optional and there are also other options can be seen by clicking the button '*Advance*'.

After clicking the button '*Run*', the function can be implemented and the output can be seen.

## 4.4.2 Forecast text

Users can use the trained model to forecast the new text. The new text, such as 'hill.txt', should be selected and put in the '*Input Path*'.



To forecast text, the '*Train*' must be *false*. '*Model file'* should be a trained sequence tagging model which has already existed.
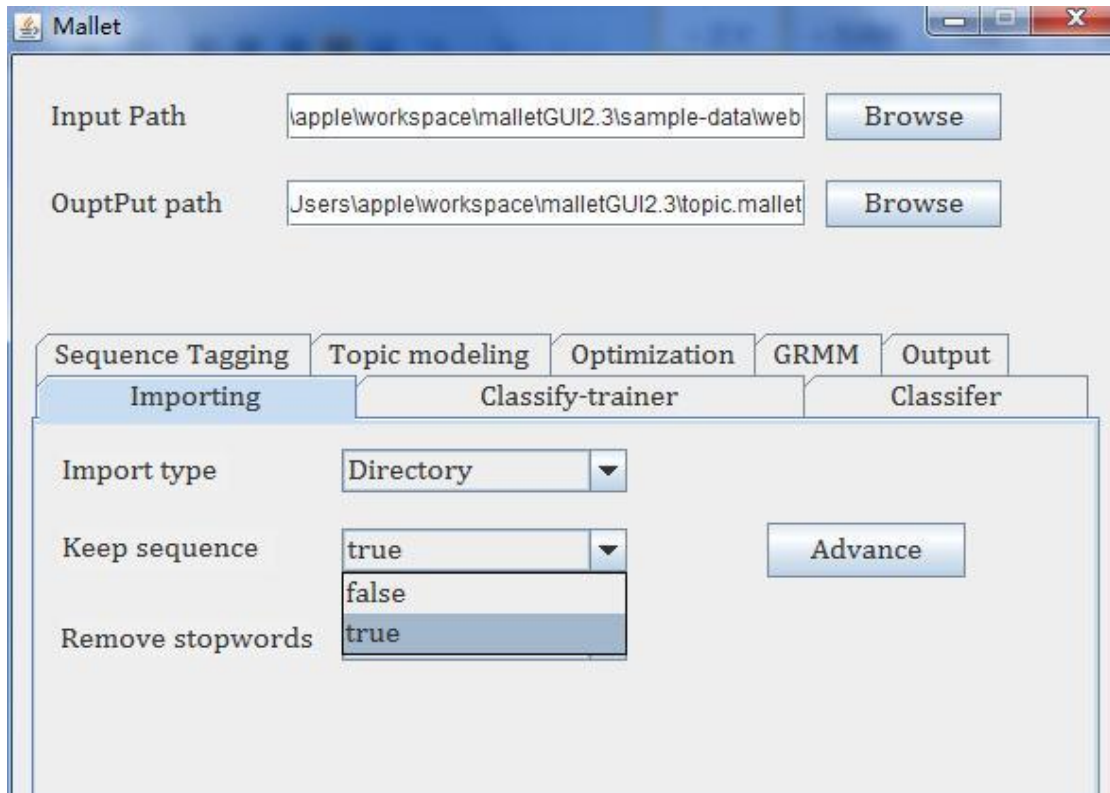
The output can be seen.
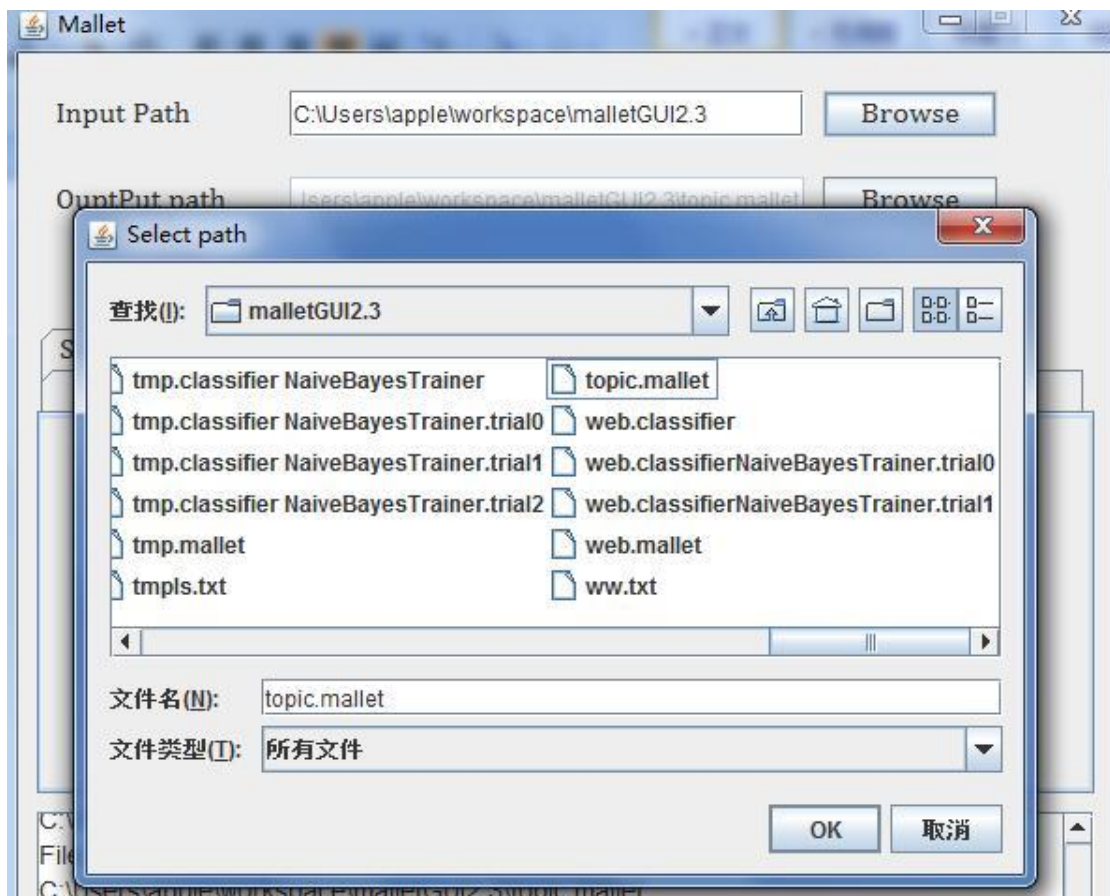
```
2016-04-12 22:01:26
Number of predicates: 140
```
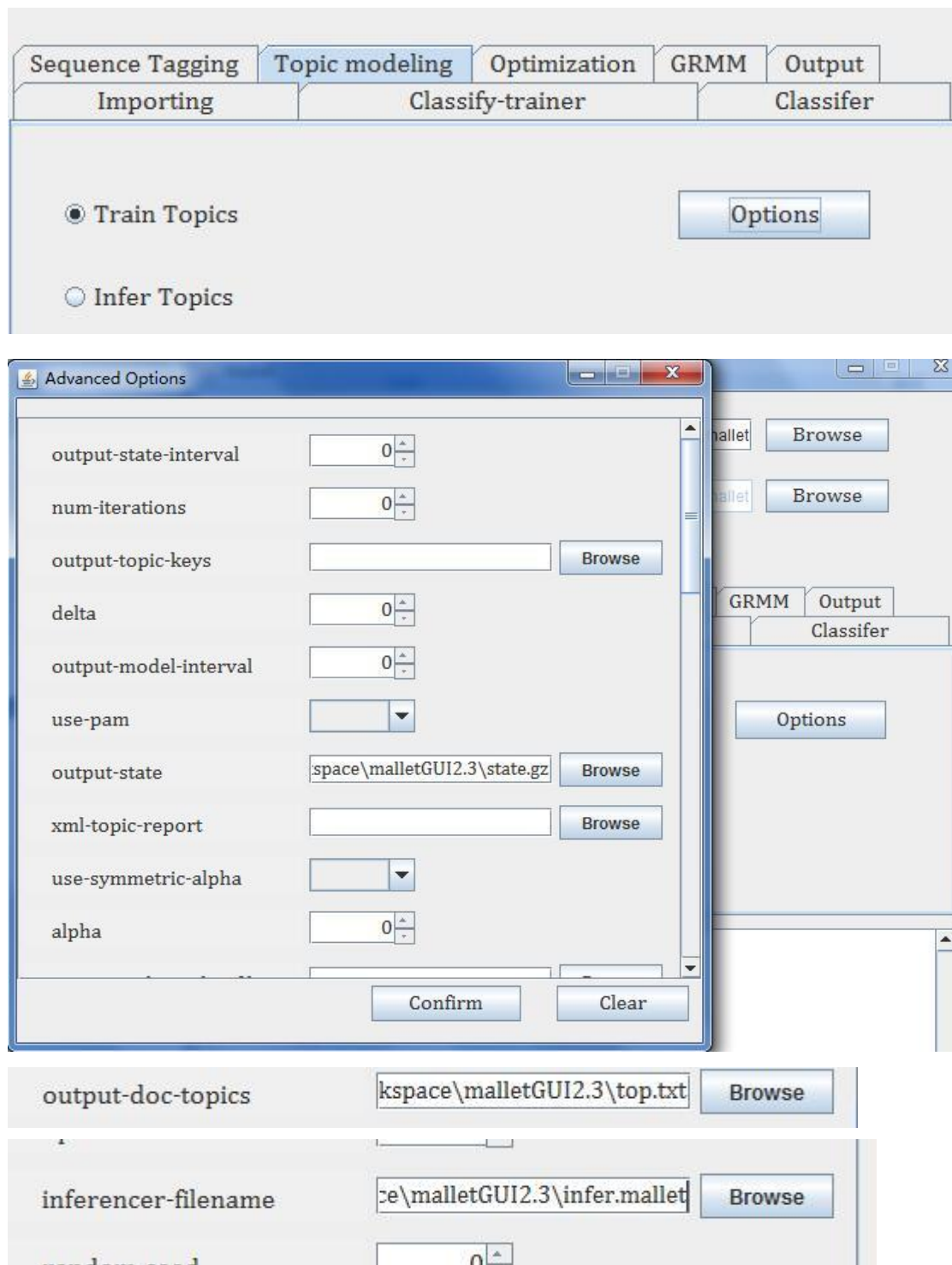
## 4.5 Topic modeling

## 4.5.1 Build the topic model

The first step is to import data and turn it into unique mallet format (.mallet), which is the same as the above. However, in order to build a topic model, one option '*Keep sequence*' must be *true*.

Then in '*Input Path*', the unique mallet format file should be input.

Select '*Train Topics*' to build the topic model, and then click the button '*Options*' to complete the output, such as *output-state*, *output-doc-topics*, *inference-filename* and etc.
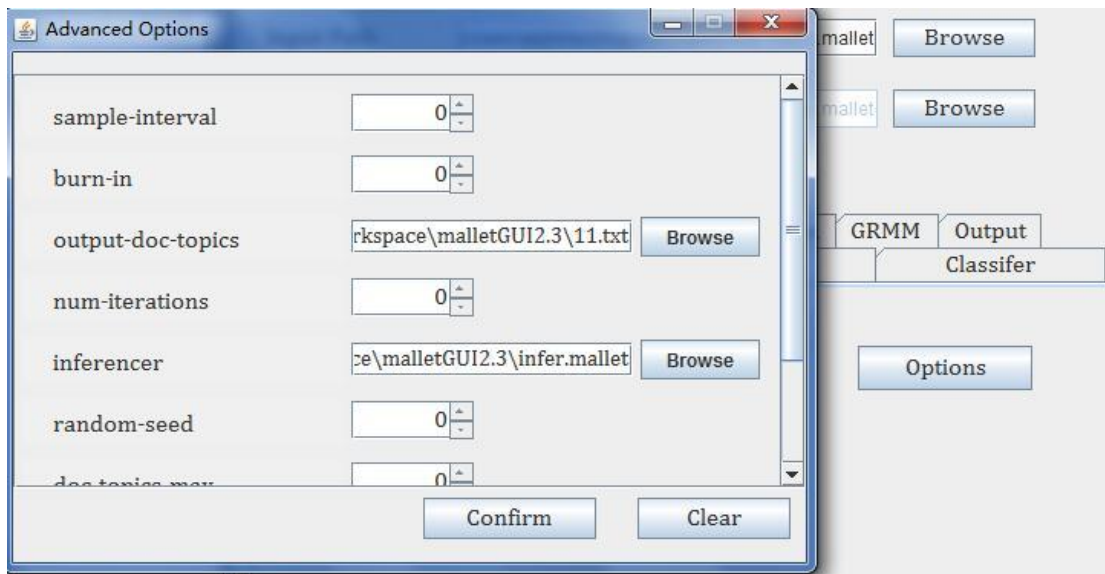


The output can be seen.

```
gamma
C:\Users\apple\workspace\malletGUI2.3\state.gz
C:\Users\apple\workspace\malletGUI2.3\top.txt
C:\Users\apple\workspace\malletGUI2.3\infer.mallet
{output-state=C:\Users\apple\workspace\malletGUI2.3\state.gz, output-doc-topics=C:\U
last options:train_topics

Topic modeling succeed.
```
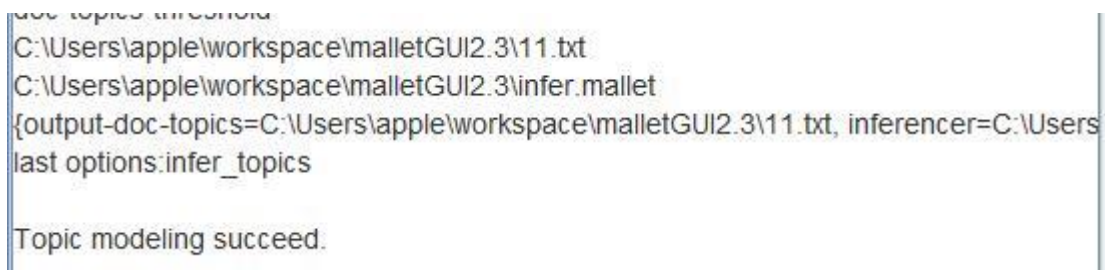
## 4.5.2 Model inference

Predict the unknown text topics through the established model. The unknown text should be in the unique mallet format. Select '*Infer Topics*' to implement model inference and then click the button '*Options*' to complete the output including selecting which model to be used to predict the text.
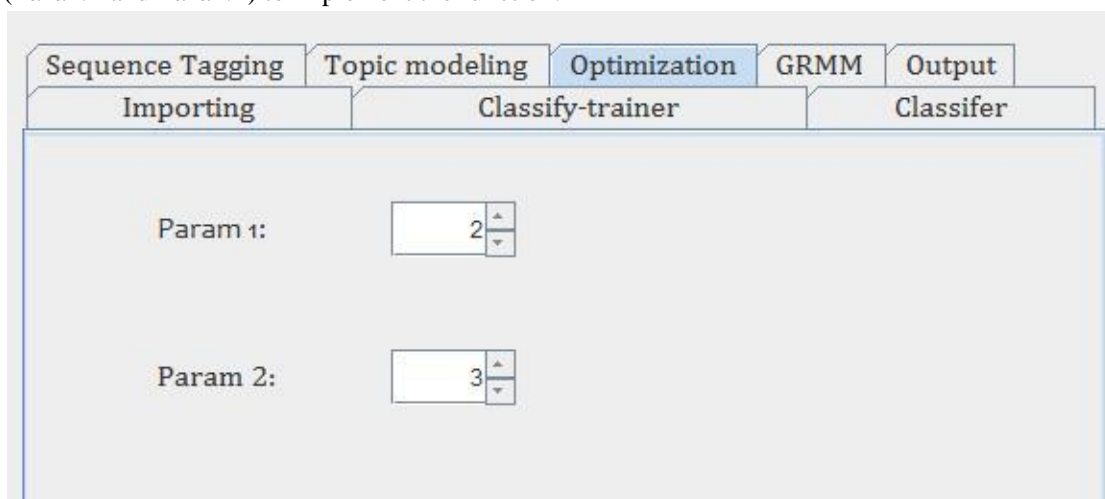
The output can be seen.



## 4.6 Optimization

Optimization needs a specific environment which is used in the above functions. In this section, optimization is a demo rather than a solution. Users only need to input the two parameters (*Param1* and *Param2*) to implement the function.
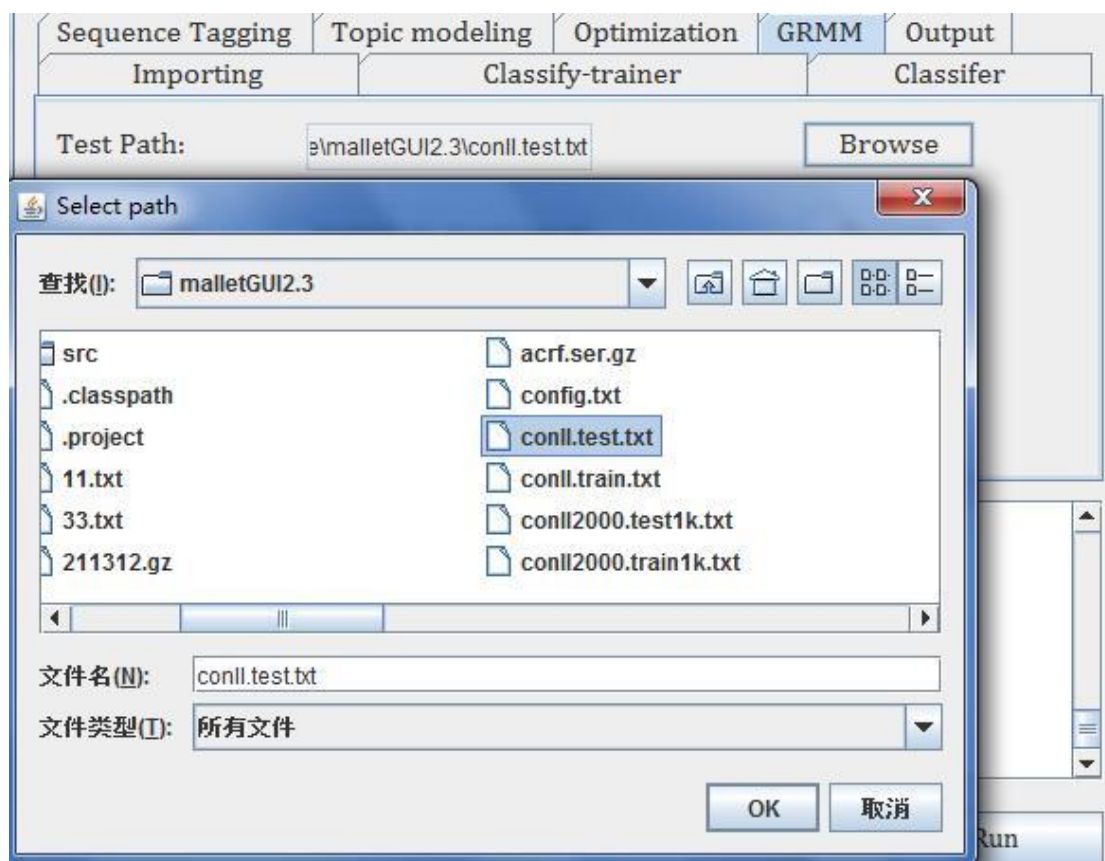


The output can be seen.

```
Exiting L-BFGS on termination #1:
value difference below tolerance (oldValue: 19.332526081620024 newValue: 19.3333

Parameter:
(0.333333365875235,-0.49999990064442473)
1.0
2.0
1.0
2.0
3.0

Optimization succeed.
```

## 4.7 Graphical Models in MALLET

*Test Path*, *Train Path* and *Model Path* should be selected by clicking the button '*Browse*'. The format for all files must conform to provisions. And the official website has the examples of the files.

The output can be seen.