



華東師範大學

 软件学院
software engineering institute

个人邮件智能挖掘及可视化

Intelligent Mining and Visualization of Personal Email

黄一夫

软件学院

华东师范大学

上海, 中国

10092510437@ecnu.cn

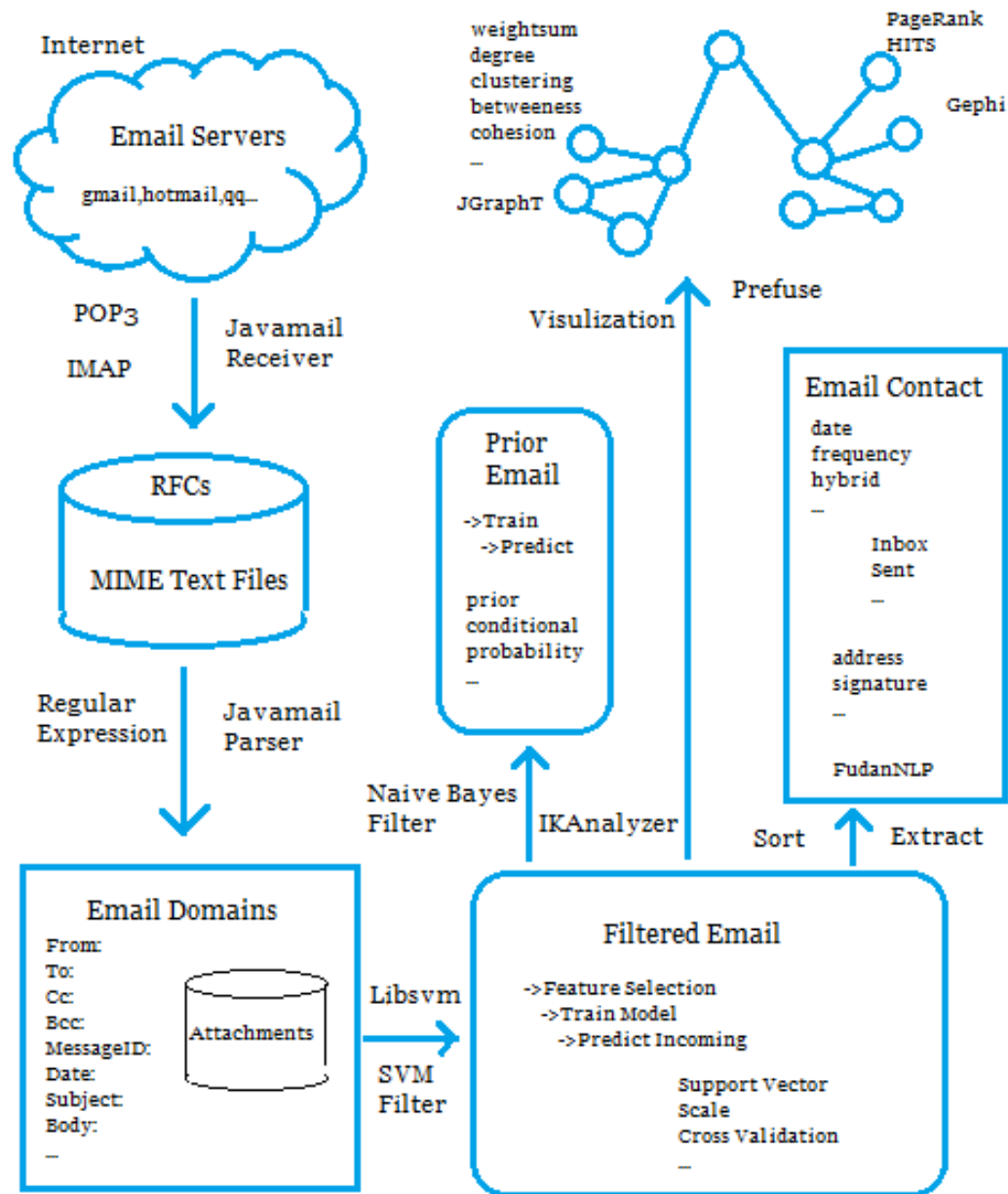
指导老师：罗远哉

目录

- **研究意义**
- 邮件预处理
- 邮件分类
- 邮件通讯录管理
- 邮件网络分析与可视化
- 代码重构
- 界面设计

研究意义

- 1. 每天我们都会收到大量的电子邮件，但查看它们却变成越来越头疼的问题。对此该项目不仅探讨垃圾邮件过滤，而且进一步向用户推荐重要邮件，从而提升邮件用户的收件效率。
- 2. 维护邮件通讯录是一个很重要的问题。该项目不仅从邮件签名，地址等中提取信息来提升邮件用户通讯录的完整性和可检索性，而且根据已给出的收件人预测推荐额外的收件人来提升邮件用户的发件效率。
- 3. 根据邮件收发关系，可以构建出个人社交网络。该项目分别从网络的结构特点和链接特点入手，计算多种系数值来衡量节点，并提供可视化展示，从而让邮件用户理清社交关系。

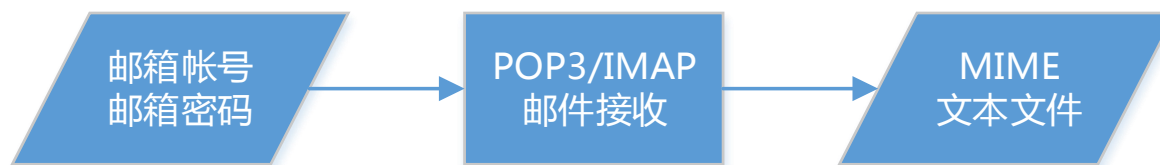


目录

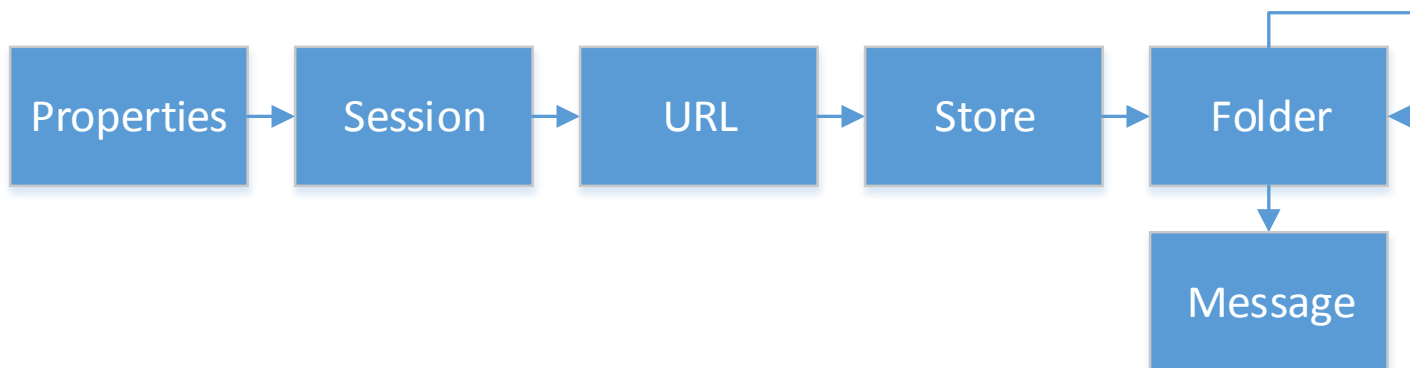
- 研究意义
- **邮件预处理**
- 邮件分类
- 邮件通讯录管理
- 邮件网络分析与可视化
- 代码重构
- 界面设计

邮件预处理

- 邮件获取
 - 输入与输出



- 调用过程



邮件预处理

- 邮件获取
 - MIME文本文件示例

Subject: =?GB2312?Q?Project_update_=BB=E3=B1=A8?=

From: Lan Man <lanman.sg@gmail.com>

To: =?GB2312?B?wLzC/A==?= <mlan@cs.ecnu.edu.cn>

Content-Type: multipart/alternative; boundary=e89a8ff1c9a0e7631504bc1146ae

--e89a8ff1c9a0e7631504bc1146ae

Content-Type: text/plain; charset=GB2312

Content-Transfer-Encoding: base64

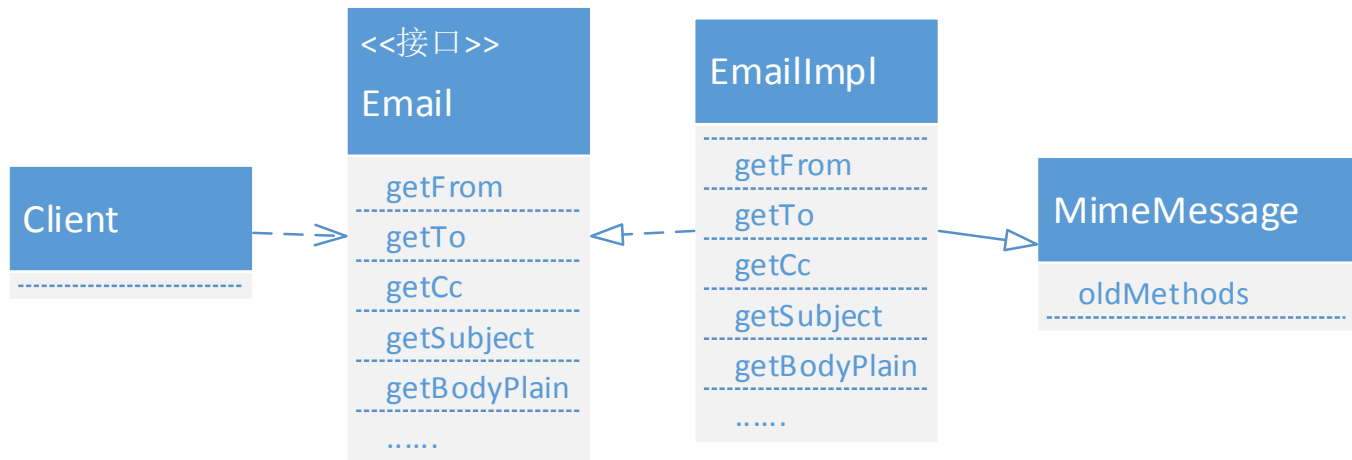
uPfOu9fps6SjrAoKvuDA68/uxL/X99K1tcTX7rrzzOG9u8jVxto01MI4yNXS0b6tuty9/MHLo6zO
qsHLyLexo7j3uPbQodfptcTP7sS/v6rM4rrNvfi2yMTcubvV/bOjvfjQ0KOszKz9bK9tqjT2s/C

邮件预处理

- 邮件解析
 - 输入与输出



- 适配器模式



邮件预处理

- 邮件解析

- 解析格式示例

- 垃圾邮件训练格式

```
0 1:22 2:0 3:1 4:6 5:0 6:0 7:  
0 1:35 2:0 3:1 4:48 5:0 6:0  
0 1:28 2:0 3:1 4:12 5:0 6:0  
0 1:23 2:0 3:1 4:14 5:0 6:0  
0 1:33 2:0 3:1 4:20 5:0 6:0  
0 1:30 2:0 3:1 4:31 5:0 6:0  
0 1:31 2:0 3:1 4:12 5:0 6:0
```

- 邮件网络格式

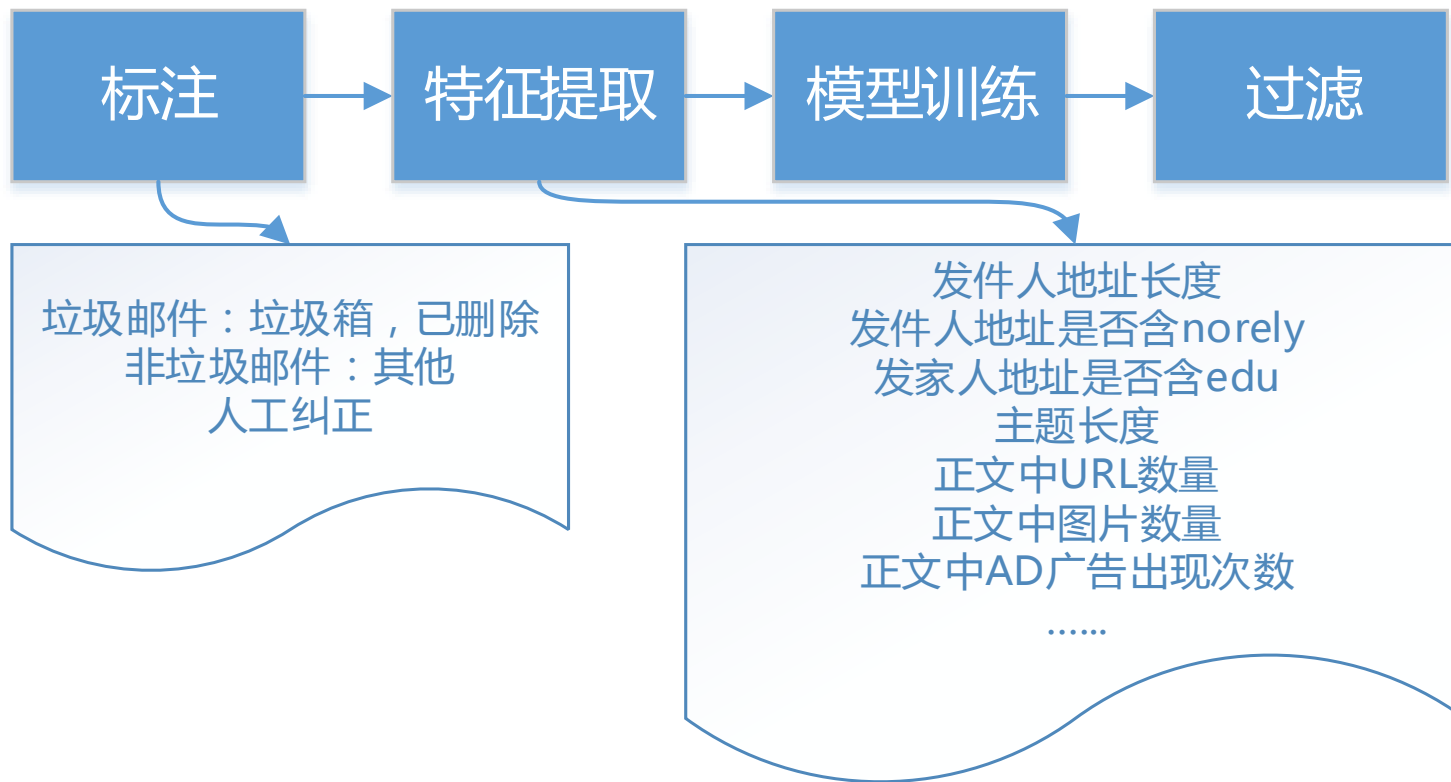
```
328679725@qq.com leegyao@gmail.com 1 2011-10-08 15  
10092510437@ecnu.cn sei_dzp@126.com 2 2012-03-01 16  
328679725@qq.com xiangnanhe@gmail.com 2 2012-01-2  
328679725@qq.com 10092510437@ecnu.cn 1 2011-08-30
```

目录

- 研究意义
- 邮件预处理
- **邮件分类**
- 邮件通讯录管理
- 邮件网络分析与可视化
- 代码重构
- 界面设计

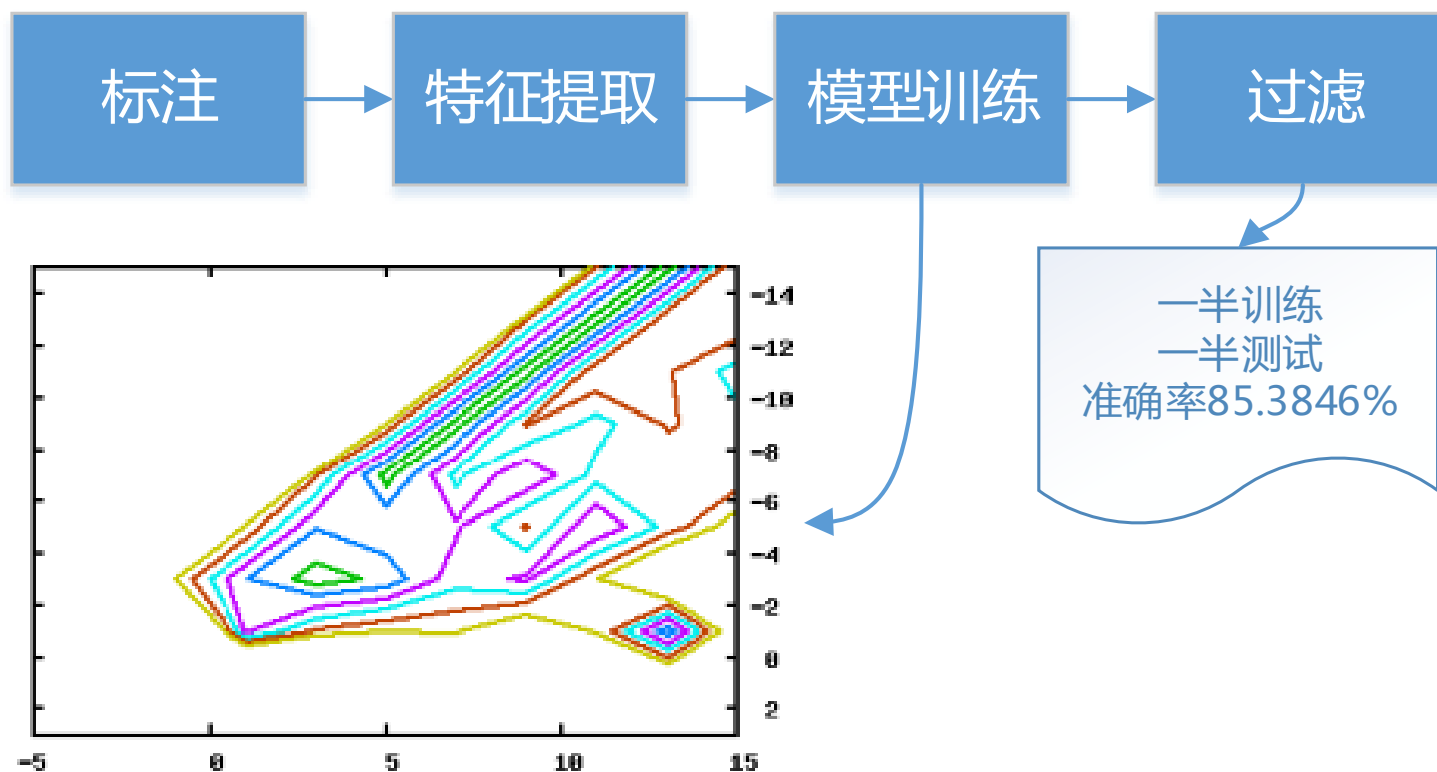
邮件分类

- 垃圾邮件过滤



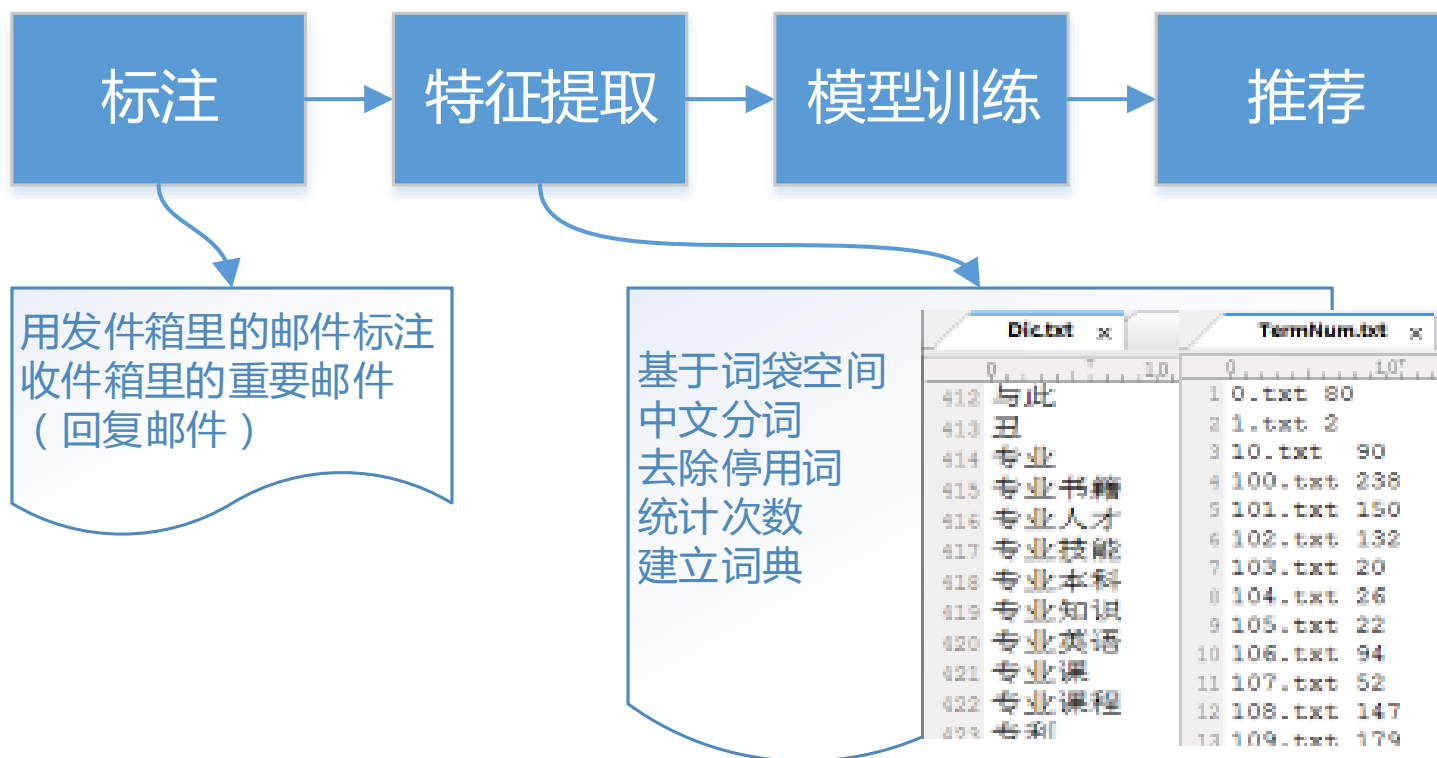
邮件分类

- 垃圾邮件过滤



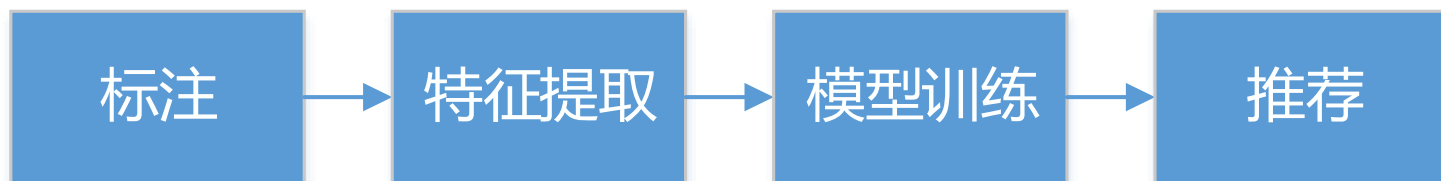
邮件分类

- 重要邮件推荐



邮件分类

- 重要邮件推荐



计算类别的先验概率和
词项在类别下的条件概率

```
model.txt x
1 prior:
2 0 0.07917888563049853
3 1 0.9208211143695014
4 condprob:
5 0 7.980209081477934E-5 6.416220204677425E-5
6 1 7.980209081477934E-5 3.2081101023387125E-5
7 2 7.980209081477934E-5 3.2081101023387125E-5
8 - - - - -
```

累加类别先验概率和词项
条件概率的对数项

```
数据挖掘分组.txt x
1 score:
2 0 -282.69634001041527
3 1 -301.99844208827534

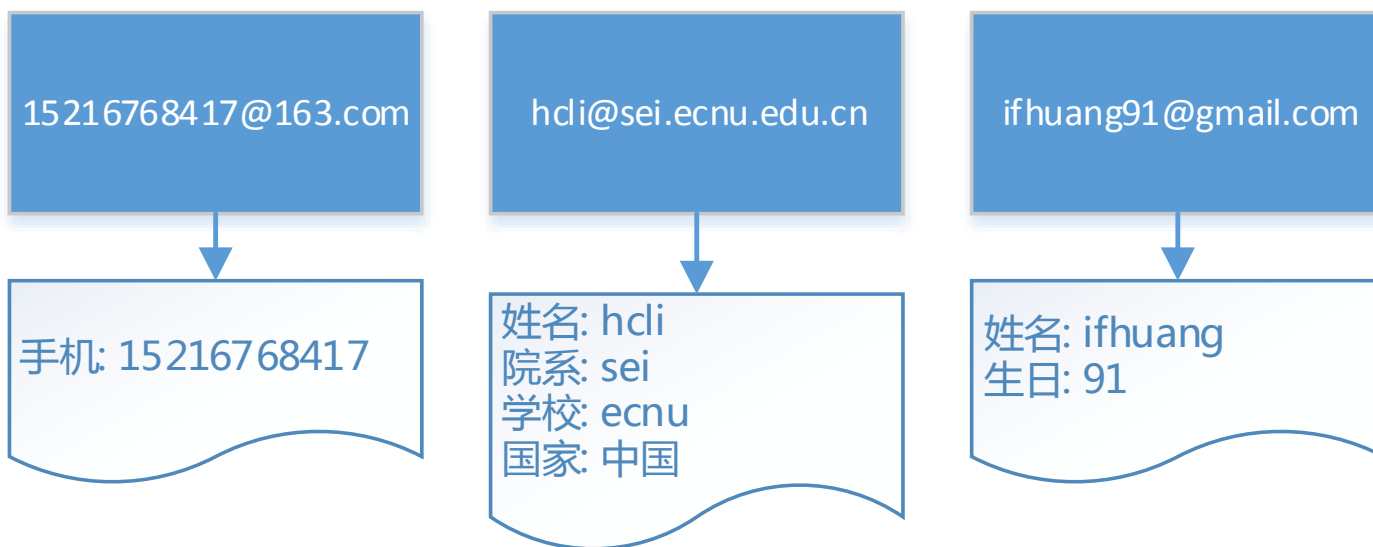
百度招聘实习生.txt x
1 score:
2 0 -269.5113104795328
3 1 -246.99610634796377
```

目录

- 研究意义
- 邮件预处理
- 邮件分类
- **邮件通讯录管理**
- 邮件网络分析与可视化
- 代码重构
- 界面设计

邮件通讯录管理

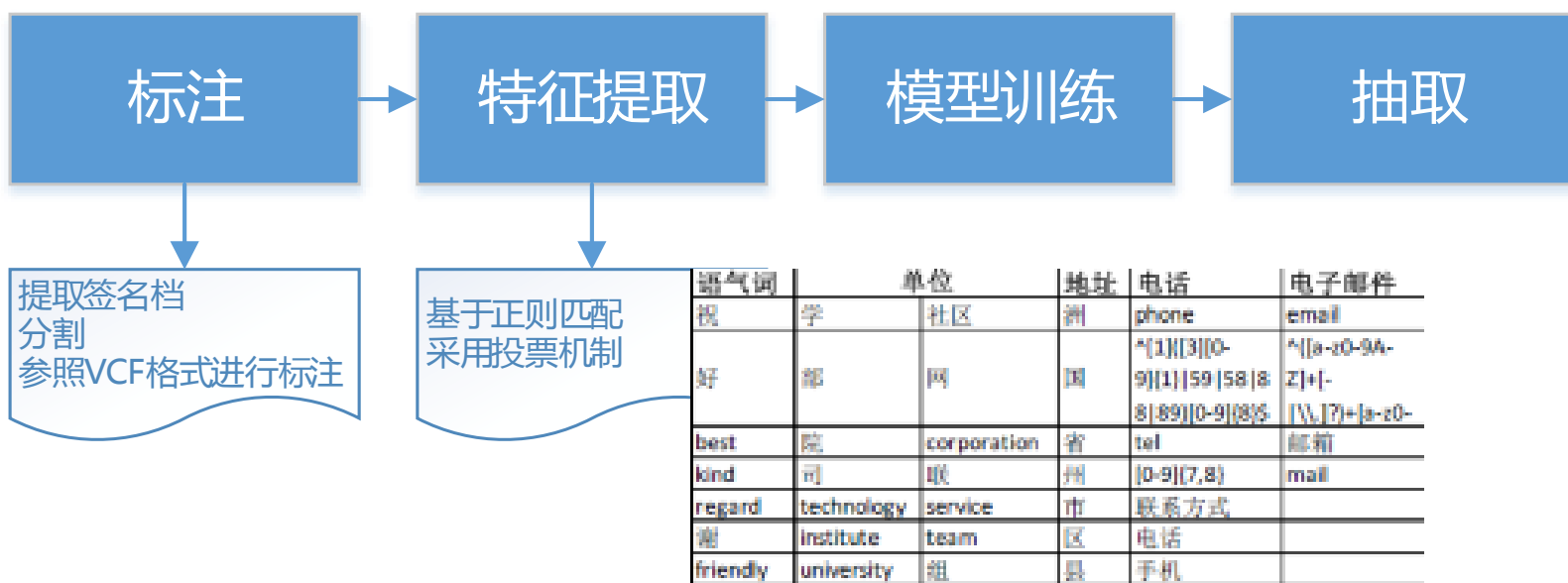
- 联系人信息抽取
 - 邮件地址信息抽取
 - 基于规则，正则匹配



邮件通讯录管理

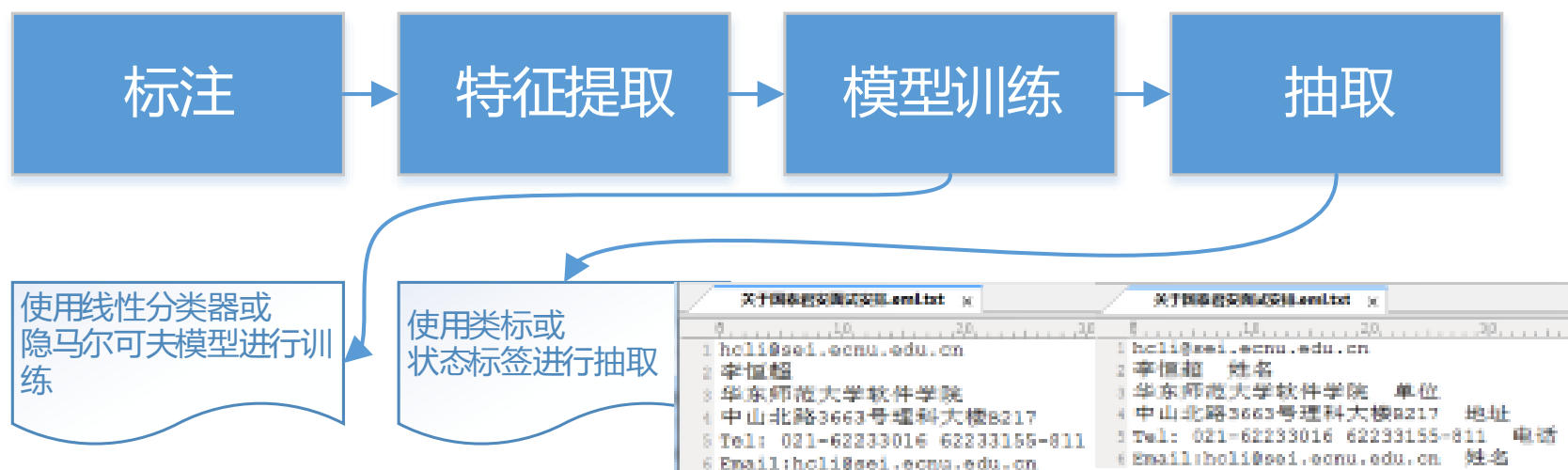
• 联系人信息抽取

一 邮件签名抽取



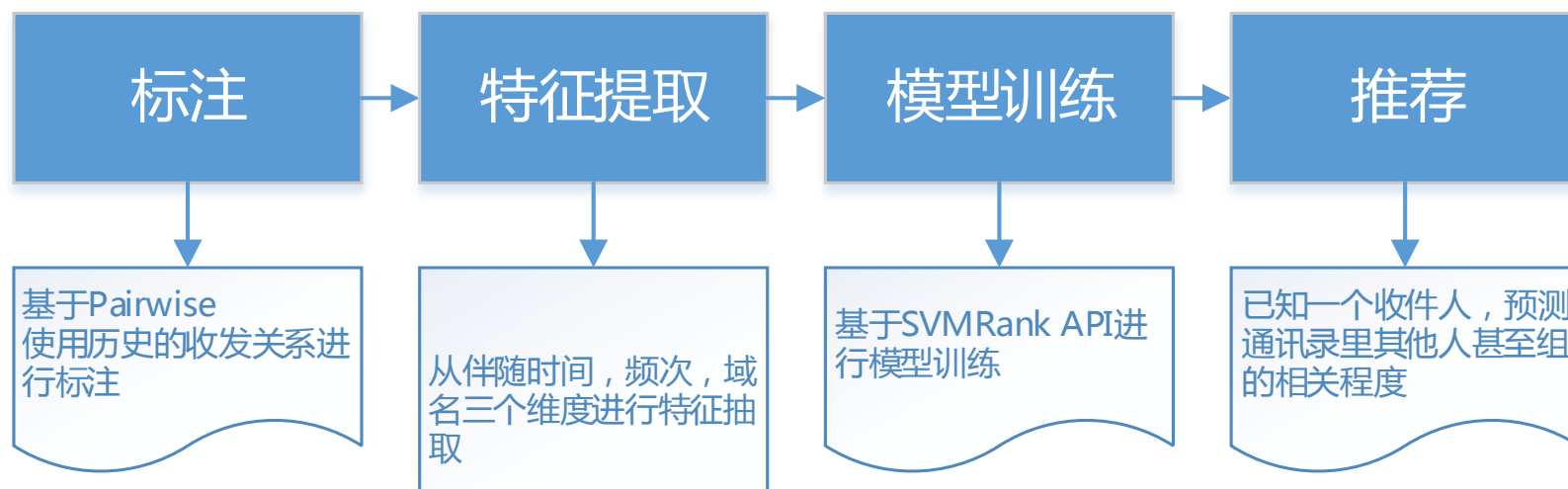
邮件通讯录管理

- 联系人信息抽取
 - 邮件签名抽取



邮件通讯录管理

• 收件人推荐



对addressbook里的每一个地址, 提取其与第一个地址的特征如下

1. 发件箱里, 该邮件之前, 出现的总次数
2. 发件箱里, 该邮件之前, 一天内出现的总次数
3. 发件箱里, 该邮件之前, 七天内出现的总次数
4. 发件箱里, 该邮件之前, 一个月出现的总次数
5. 发件箱里, 该邮件之前, 最近一次到此间间隔的天数, 若无则用无穷大表示

目录

- 研究意义
- 邮件预处理
- 邮件分类
- 邮件通讯录管理
- **邮件网络分析与可视化**
- 代码重构
- 界面设计

邮件网络分析与可视化

- 邮件网络分析

EmailNetworkParser

EmailNetworkComputer

- 结构

- 度数：衡量节点的活跃程度
 - 聚集系数：衡量节点的聚团程度
 - 介数：衡量节点的中介程度
 - 凝聚度：衡量节点的中心程度

- 链接

- PageRank：衡量节点的重要程度
 - HITS：衡量节点的枢纽程度和权威程度

邮件网络分析与可视化

- 邮件网络分析

- 重要节点

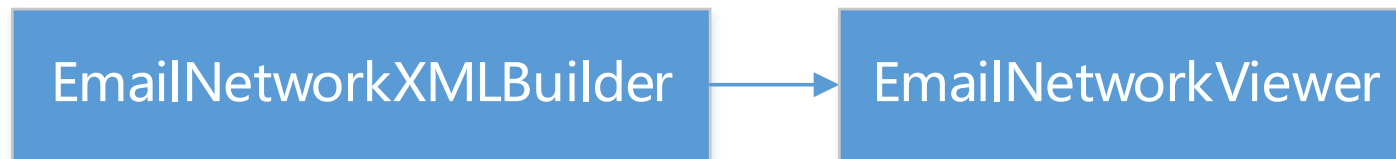
- 活跃节点：与很多节点通信的节点
 - 领导节点：仅与一些活跃节点通信的节点
 - 意外节点：外围节点和桥接节点

- 社团发现

- 话题发现：不同社团的话题不同
 - 专家发现：汇总并分类历史知识

邮件网络分析与可视化

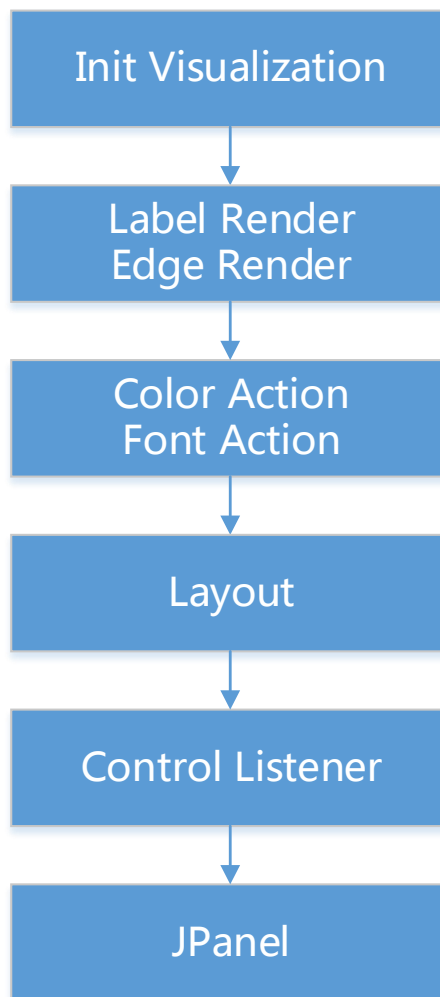
- 邮件网络可视化



- 根据From, To, Cc建立邮件网络
- 采用物理拉力引擎进行绘制
- 使用不同颜色代表不同类型的节点
- 边长与通信次数成反比
- 从边可以查看相关通信邮件

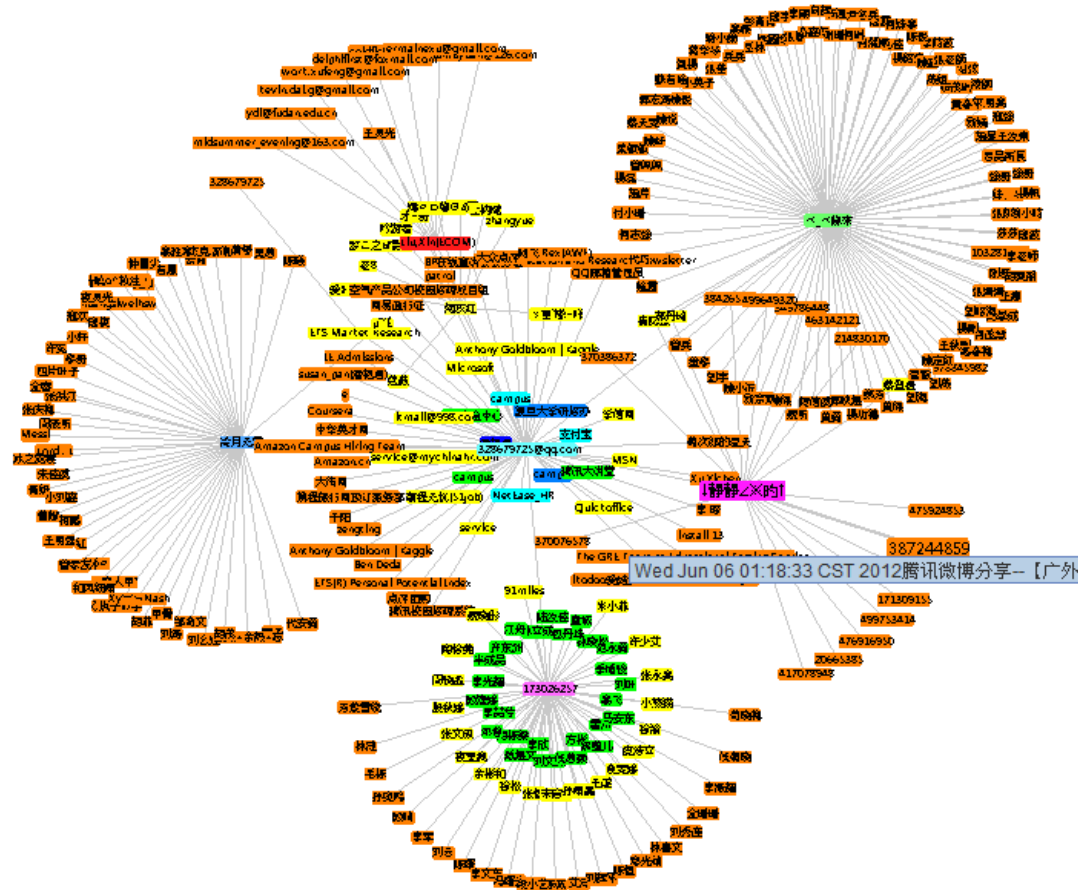
邮件网络分析与可视化

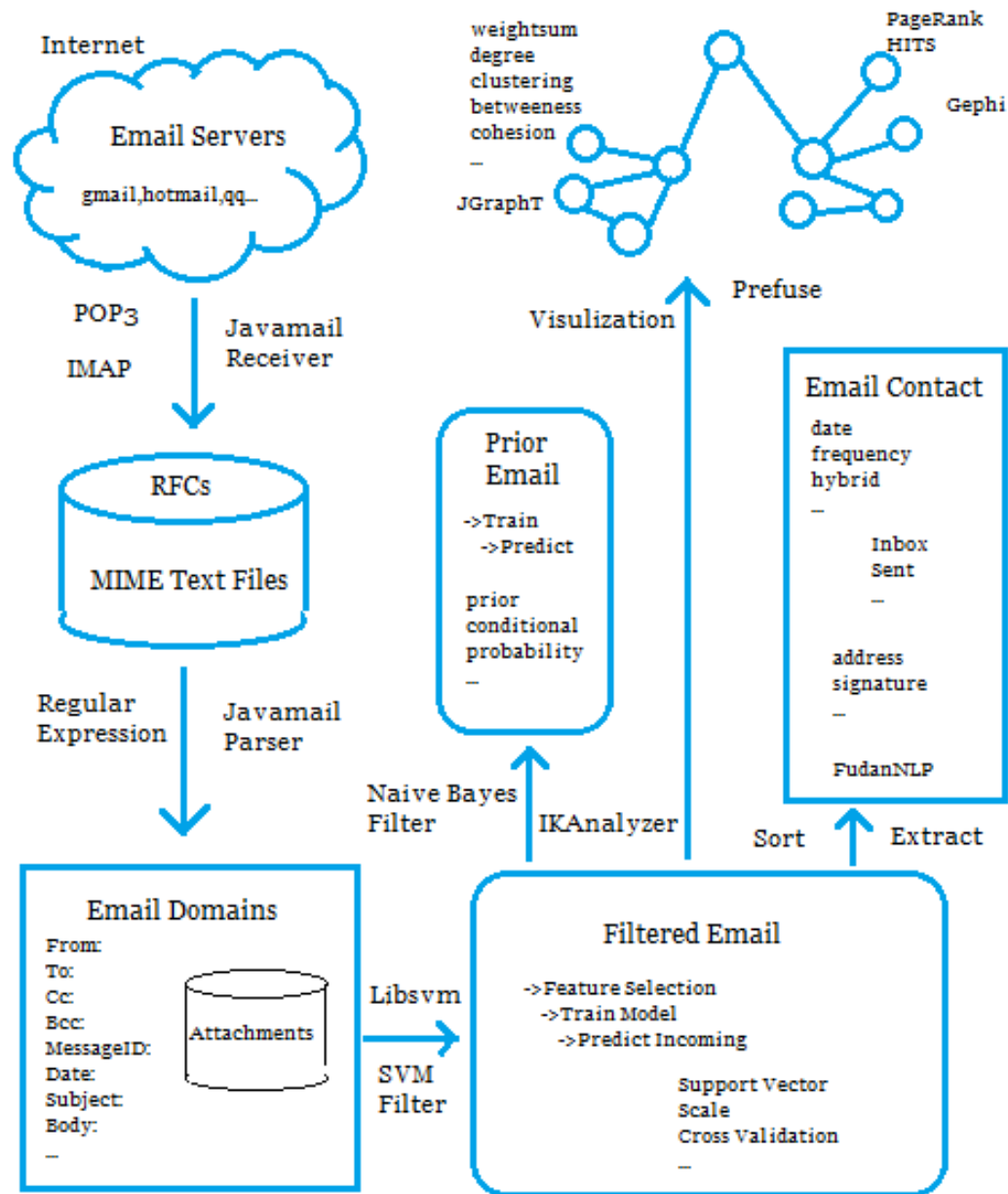
- 邮件网络可视化



邮件网络分析与可视化

- 邮件网络可视化



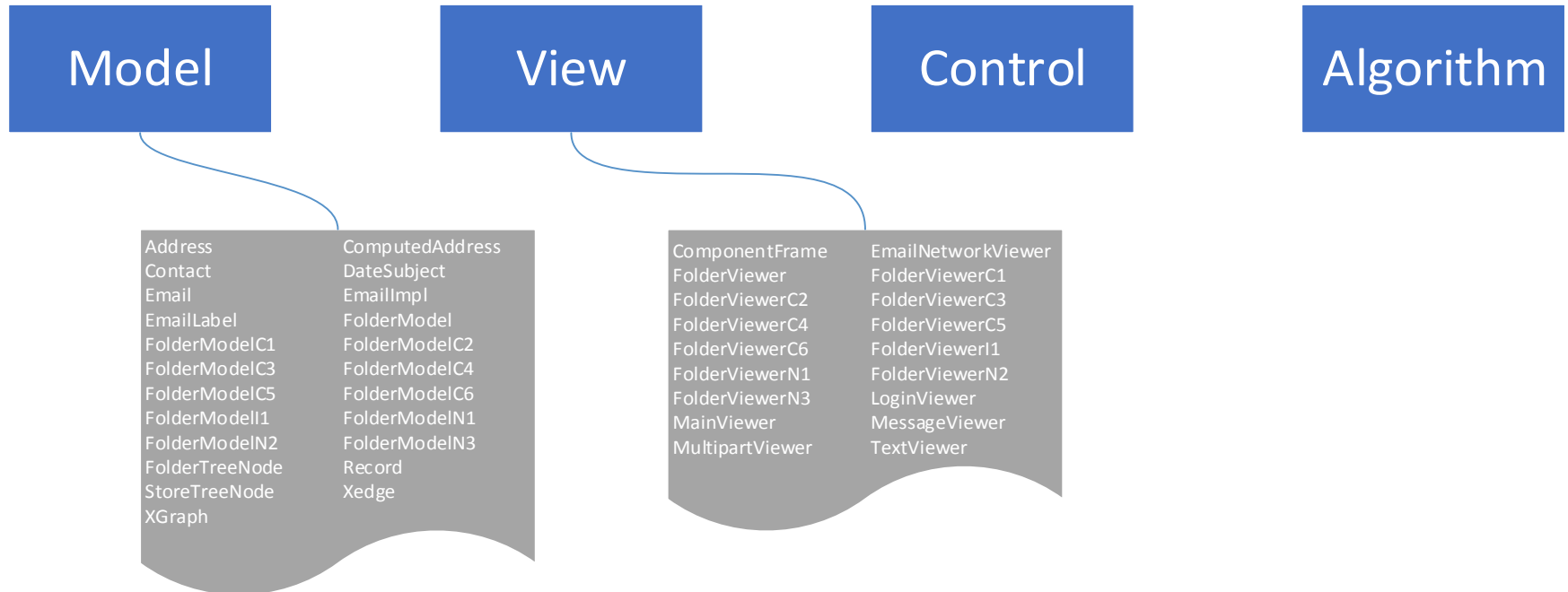


目录

- 研究意义
- 邮件预处理
- 邮件分类
- 邮件通讯录管理
- 邮件网络分析与可视化
- **代码重构**
- 界面设计

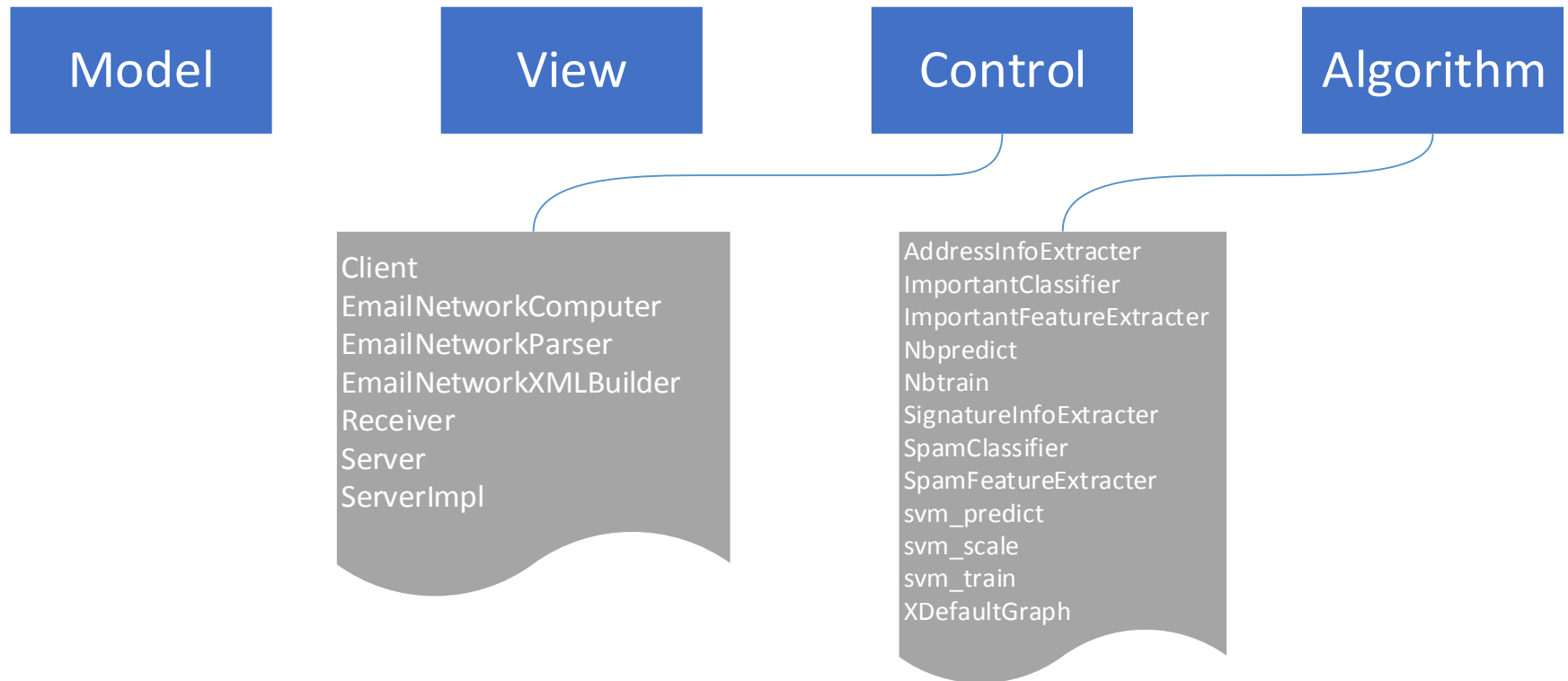
代码重构

- MVC架构



代码重构

- MVC架构

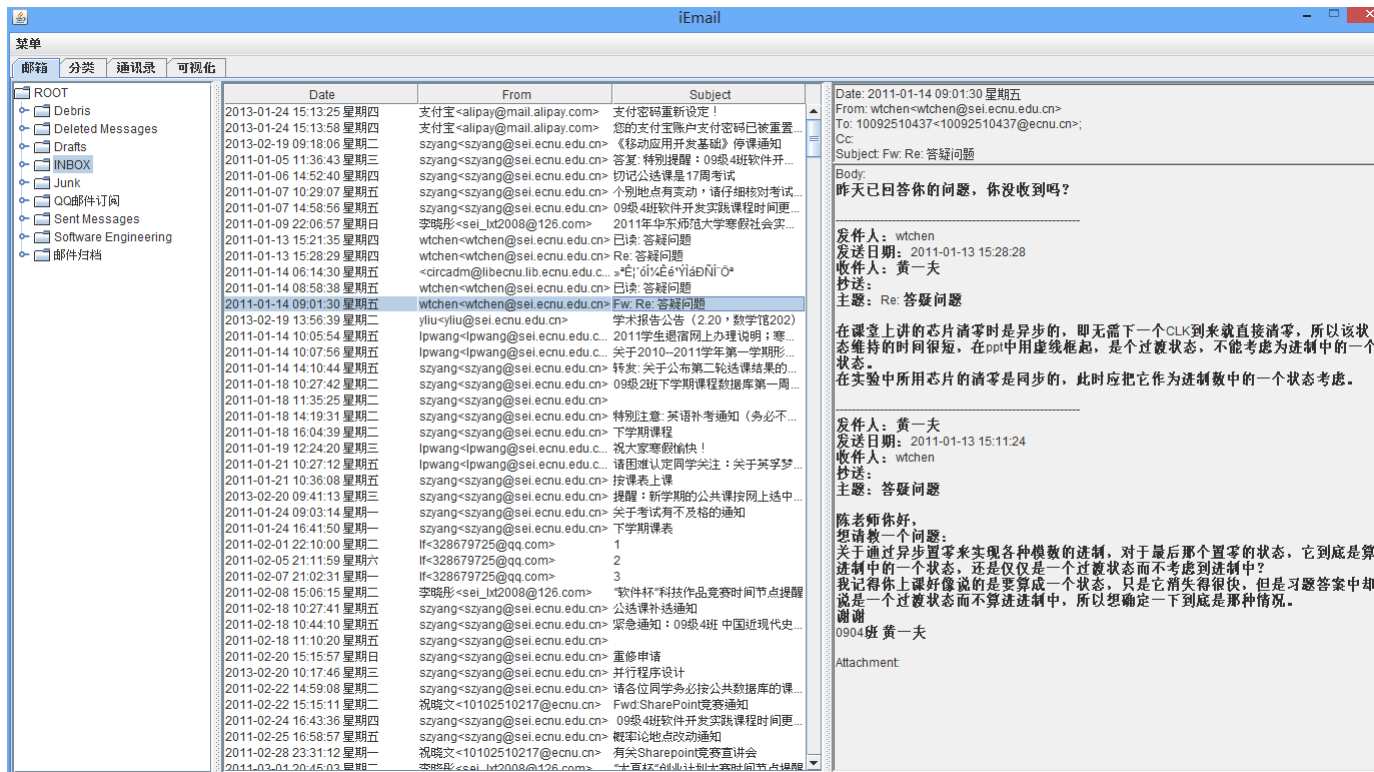


目录

- 研究意义
- 邮件预处理
- 邮件分类
- 邮件通讯录管理
- 邮件网络分析与可视化
- 代码重构
- **界面设计**

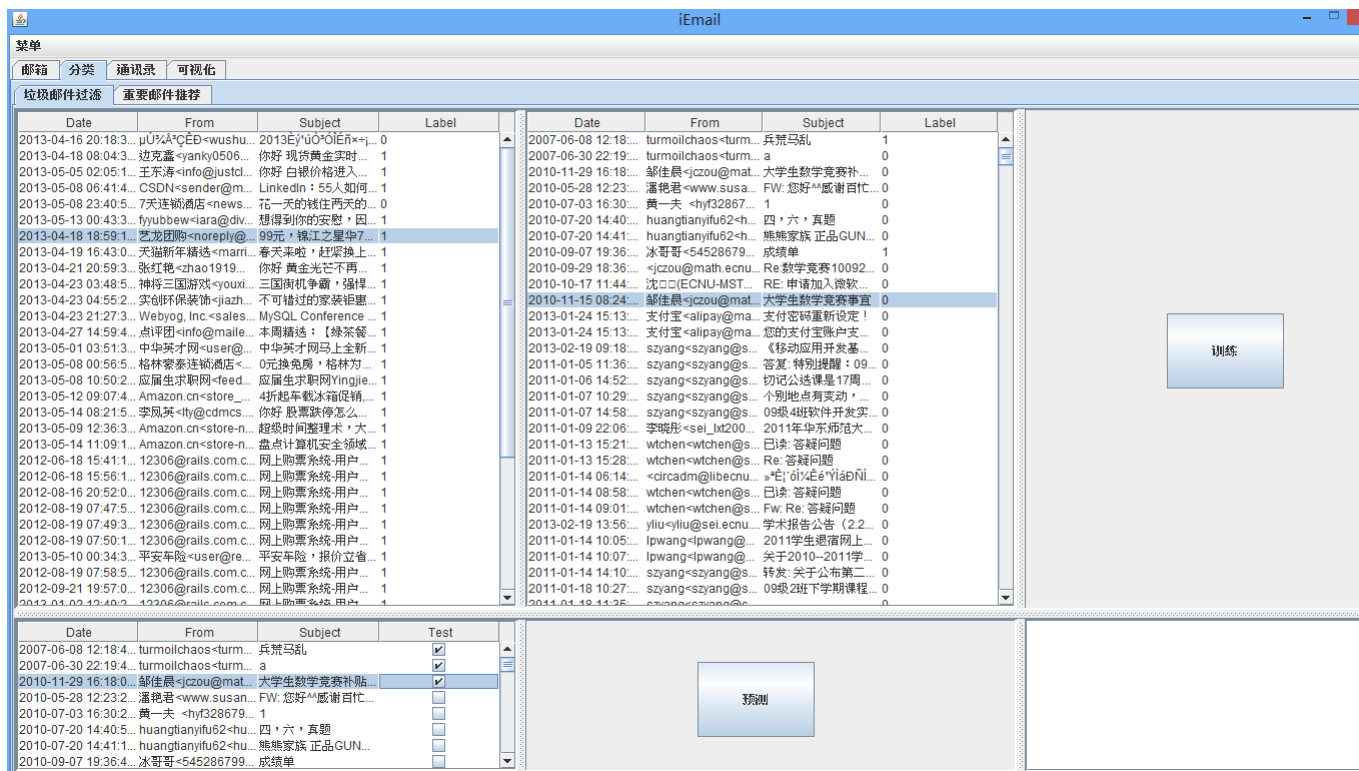
界面设计

• 邮箱界面



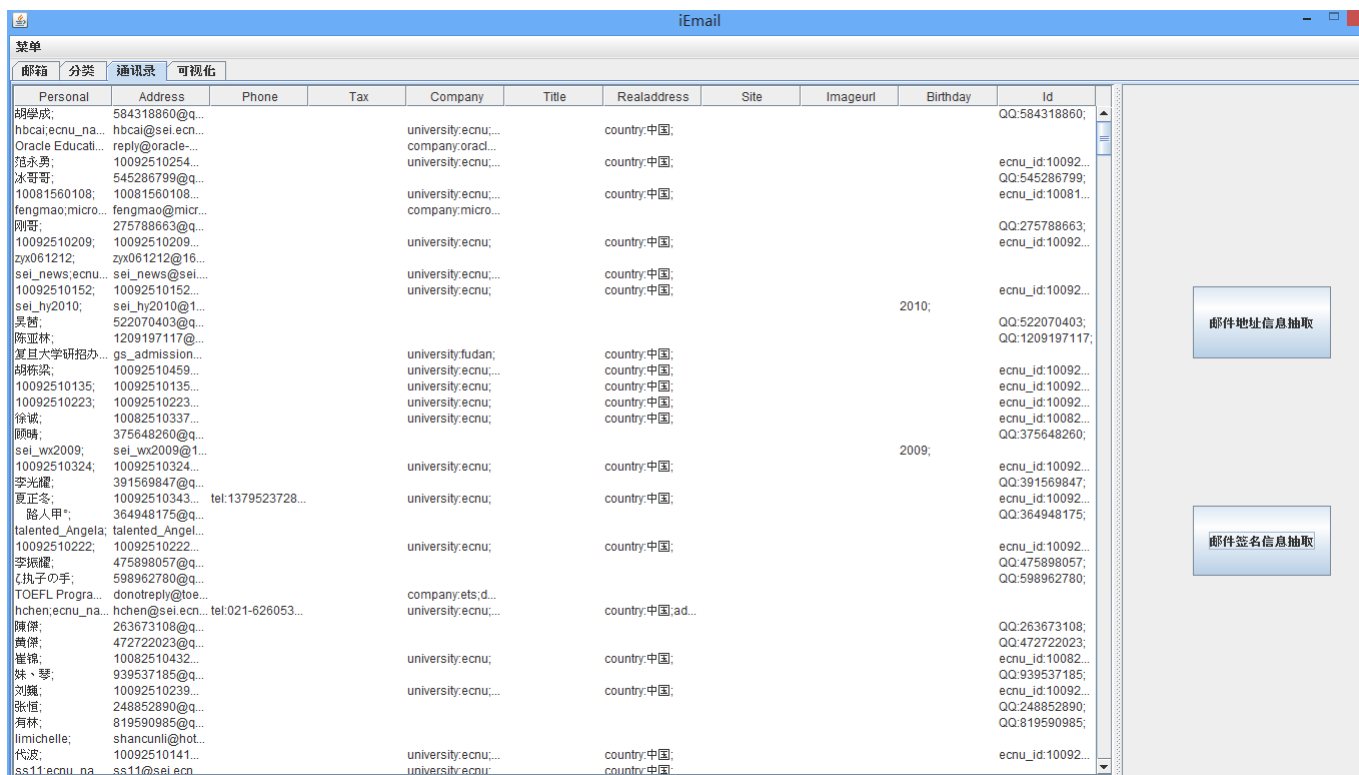
界面设计

- 分类界面



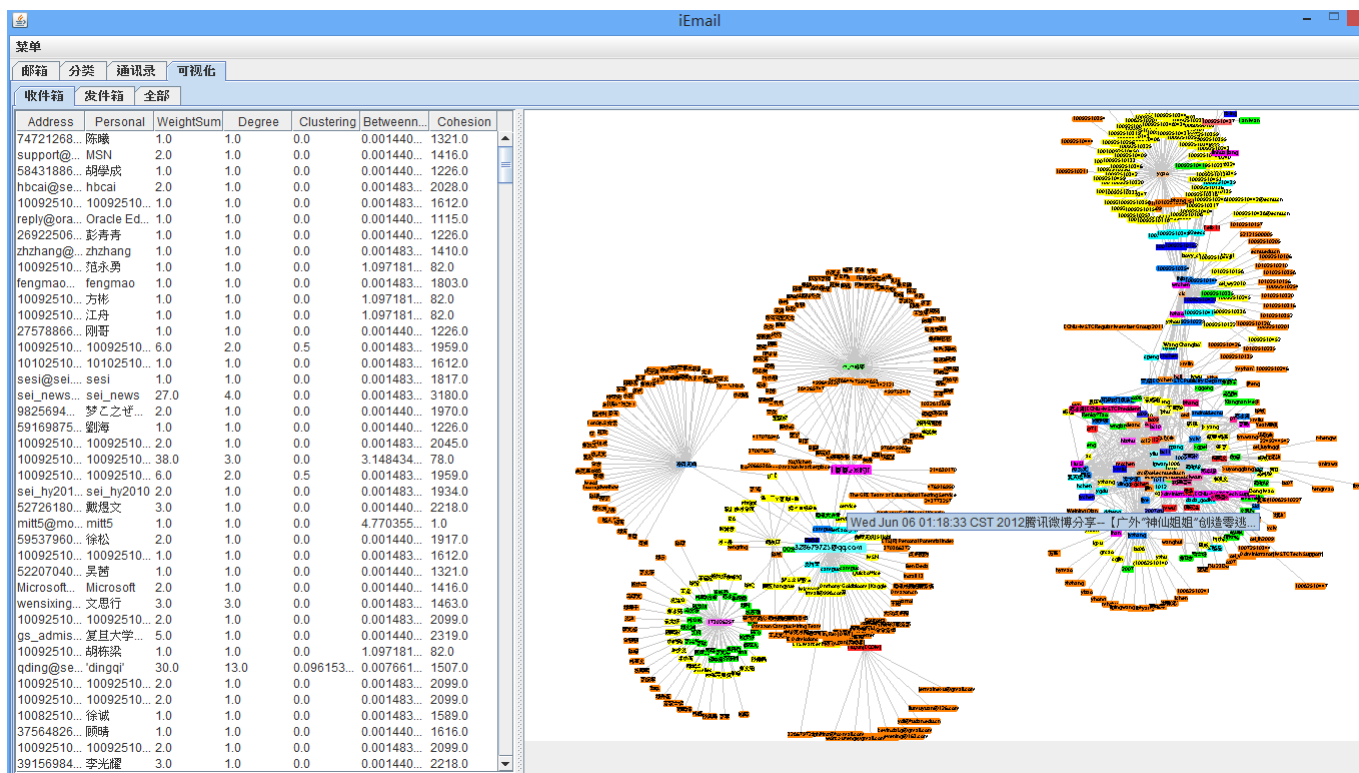
界面设计

- 通讯录界面



界面设计

- 可视化界面





華東師範大學

软件学院
software engineering institute

Thank You!

