



基于深度学习的微博情感分析

梁军 zhengdaxg@163.com

指导老师：柴玉梅 胥红英

郑州大学自然语言处理实验室

<http://nlp.zzu.edu.cn>

2014年10月18日



报告提纲

1. 背景介绍
2. 微博情感分析模型
 - 2.1 词语的表示
 - 2.2 递归自编码(Recursive AutoEncoder)
 - 2.3 基于RAE的情感极性转移模型
 - 2.4 情感分析模型算法
3. 实验结果分析
4. 总结



报告提纲

1. 背景介绍
2. 微博情感分析模型
 - 2.1 词语的表示
 - 2.2 递归自编码(Recursive AutoEncoder)
 - 2.3 基于RAE的情感极性转移模型
 - 2.4 情感分析模型算法
3. 实验结果分析
4. 总结

1. 背景介绍

* 情感分析

文本情感分析(sentiment analysis), 又称为意见挖掘, 是对带有情感色彩的主观性文本进行分析、处理、归纳和推理的过程。其中, 主观情感可以是他们的判断或者评价, 他们的情绪状态, 或者有意传递的情感信息。因此, 情感分析的一个主要任务就是情感倾向性的判断。

* 情感倾向性判断方法 or ?

1. 基于情感词典的方法
2. 转化为分类问题使用机器学习方法



报告提纲

1. 背景介绍

2. 微博情感分析模型

2.1 词语的表示

2.2 递归自编码(Recursive AutoEncoder)

2.3 基于RAE的情感极性转移模型

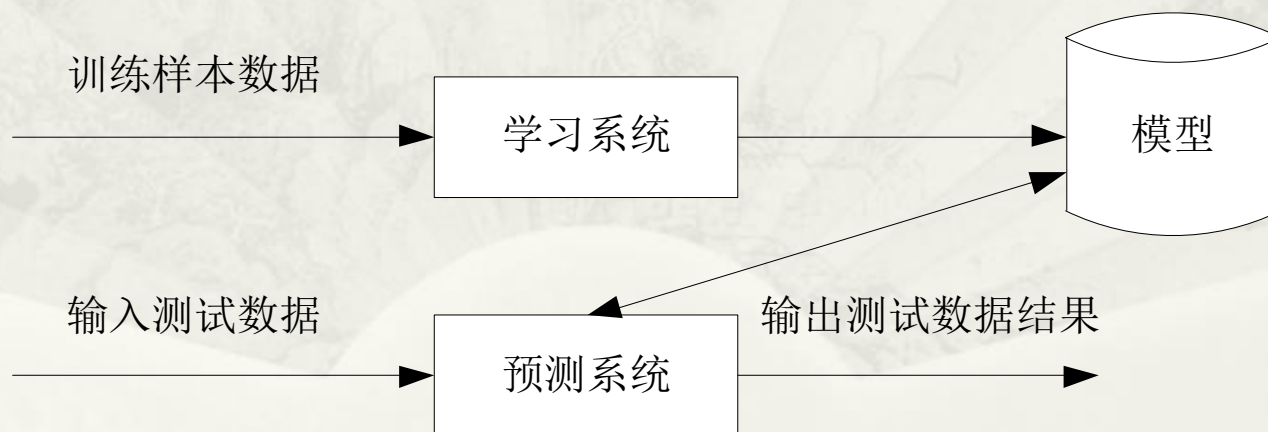
2.4 情感分析模型算法

3. 实验结果分析

4. 总结

2. 微博情感分析模型

本文提出一个基于RAE的情感极性转移模型，该模型首先将文本数据转化为低维实数向量表示，建立表示文本特征的矩阵，然后将其作为基于RAE的情感极性转移模型的输入，最后使用LBFGS算法迭代生成模型。





报告提纲

1. 背景介绍
2. 微博情感分析模型
 - 2.1 词语的表示
 - 2.2 递归自编码(Recursive AutoEncoder)
 - 2.3 基于RAE的情感极性转移模型
 - 2.4 情感分析模型算法
3. 实验结果分析
4. 总结

2.1 词语的表示(1/3)

* 传统的词语表示方法

- **向量空间表示方法：**将单词作为原子符号，其中词语通常被表示成由一串0、1代码组成的向量：[0 0 0 0 0 0 1 0 0 0]
- **缺点：**
 - 维度高：20K（语音）、50K（PTB）、13M（Google 1T）
 - 无法体现相关性：酒店[0 0 0 1 0 0 0 0 0] & 宾馆[0 0 0 0 0 1 0 0 0]=0

2.1 词语的表示(2/3)

* 词向量表示方法

- 在聚类模型中提出使用**分布式方法**表示词语 (distribution representation)
 - Latent Semantic Analysis (LSA/LSI), Random projections
 - Latent Dirichlet Analysis(LDA), HMM clustering
- Neural **word embeddings** as a distributed representation
 - Combine vector space semantics with the prediction of probabilistic models (Bengio et al. 2003, Collobert & Weston 2008, Turian et al. 2010)
 - 表示方法: 酒店[0.357 0.159 -0.526 0.610]

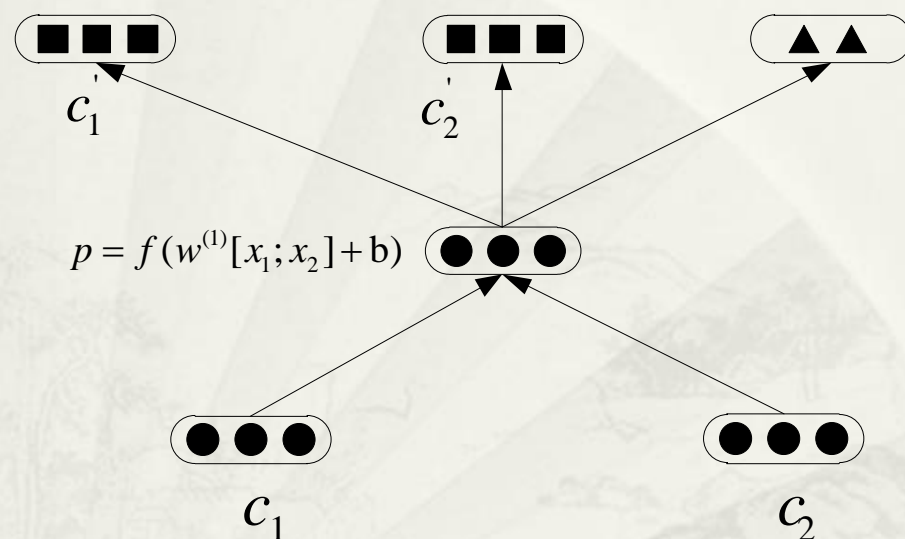
2.1 词语的表示(3/3)

* 自编码 (AutoEncoder)

例如“深度”的向量表示为 $[0.3 \ 0.1 \ 0.6]^T$ ，“学习”的向量表示为 $[0.2 \ 0.5 \ 0.7]^T$ ，那么该如何表示“深度学习”呢？假定“深度学习”是“深度”、“学习”的父节点 p ，“深度”是第一个子节点 c_1 ，“学习”是第二个子节点 c_2 ，那么 p 可由函数 f 从 c_1 、 c_2 映射得到：

$$p = f(w^{(1)}[c_1; c_2] + b^{(1)})$$

$$s(p; \theta) = \mu(w^{(3)}p + b^{(3)})$$



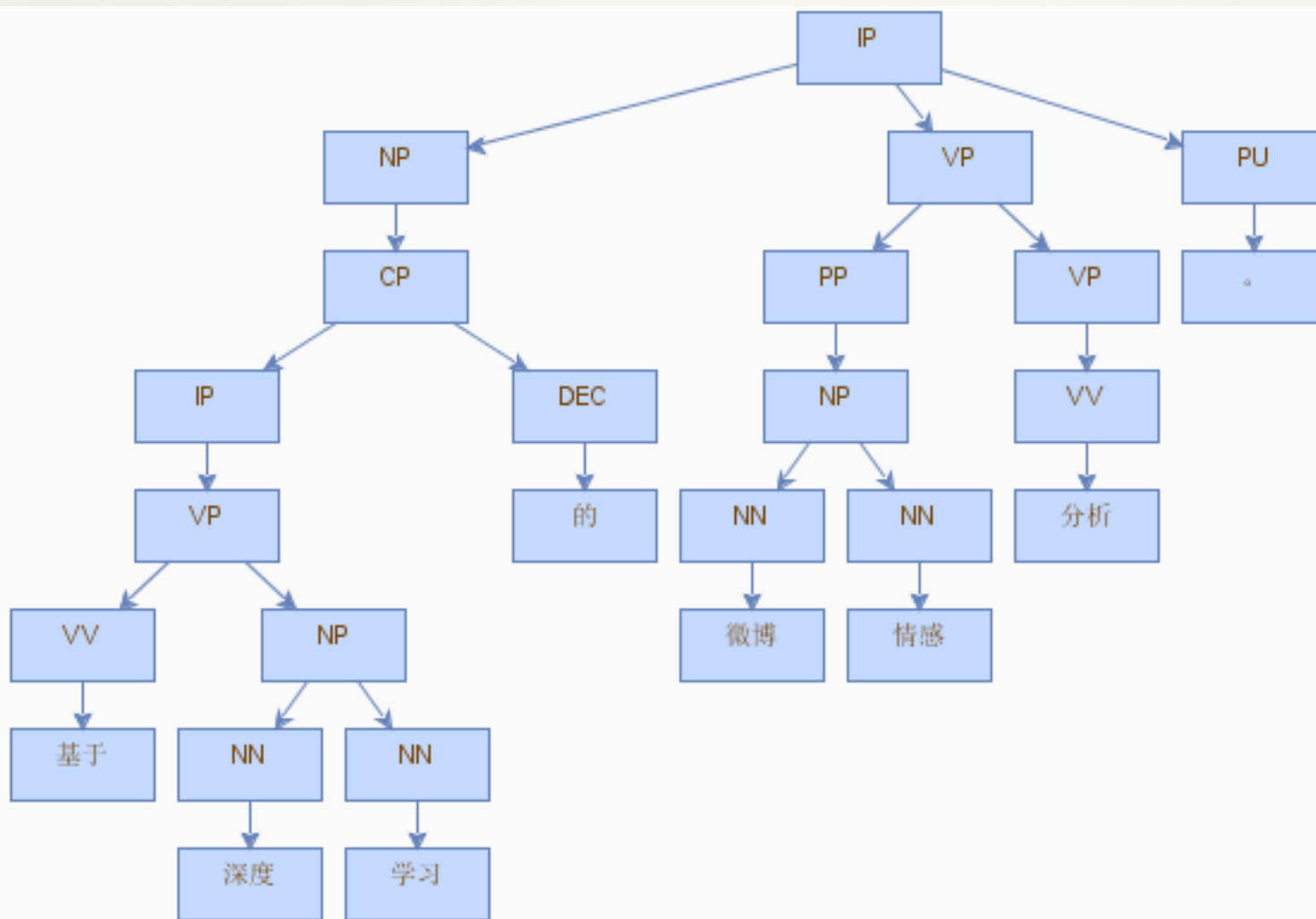
$$E_{rec} = \|c_1 - c_1'\|^2 + \|c_2 - c_2'\|^2$$



报告提纲

1. 背景介绍
2. 微博情感分析模型
 - 2.1 词语的表示
 - 2.2 递归自编码(Recursive AutoEncoder)
 - 2.3 基于RAE的情感极性转移模型
 - 2.4 情感分析模型算法
3. 实验结果分析
4. 总结

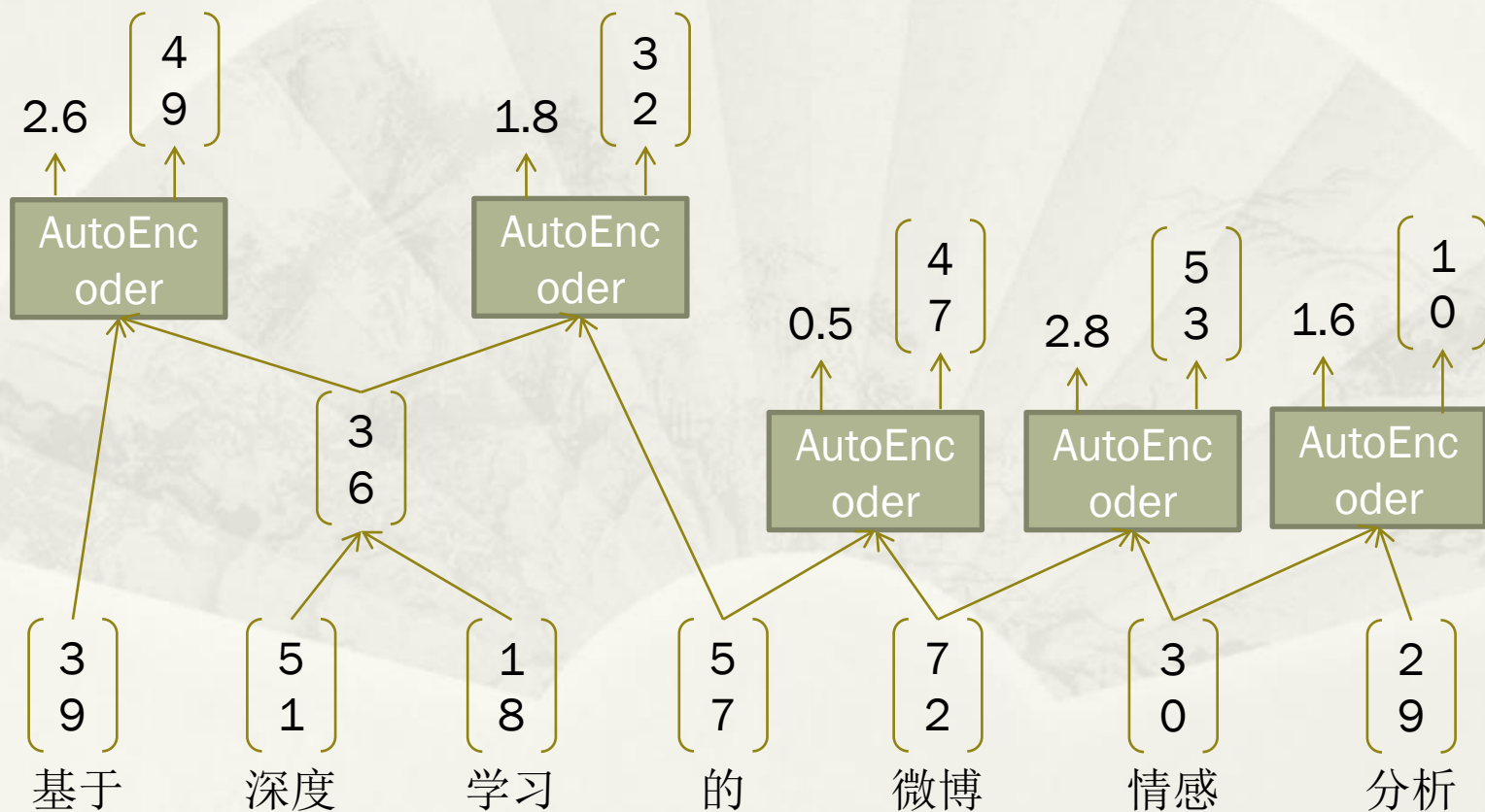
2.2 递归自编码(1/3)



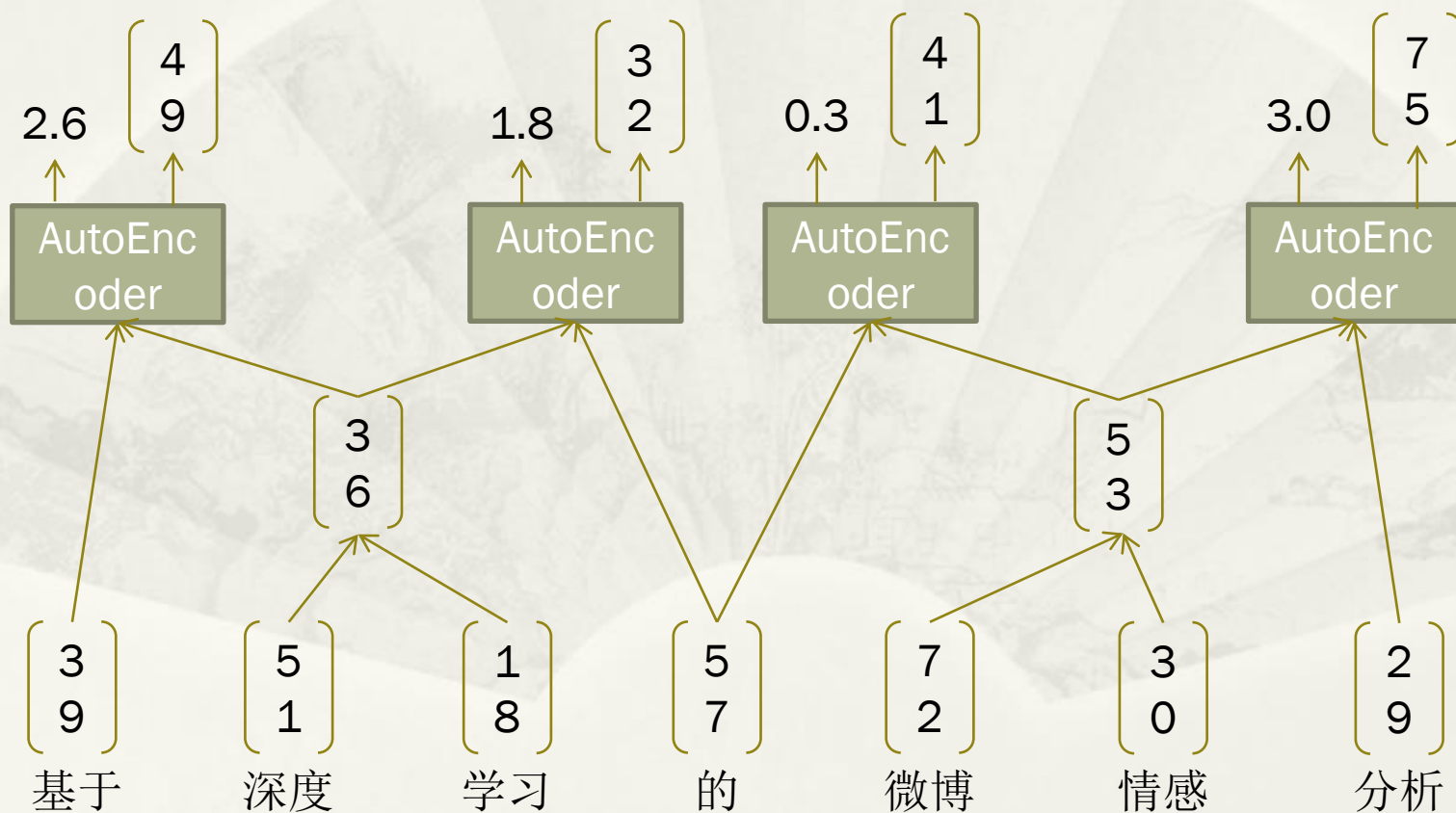
2.2 递归自编码(2/3)



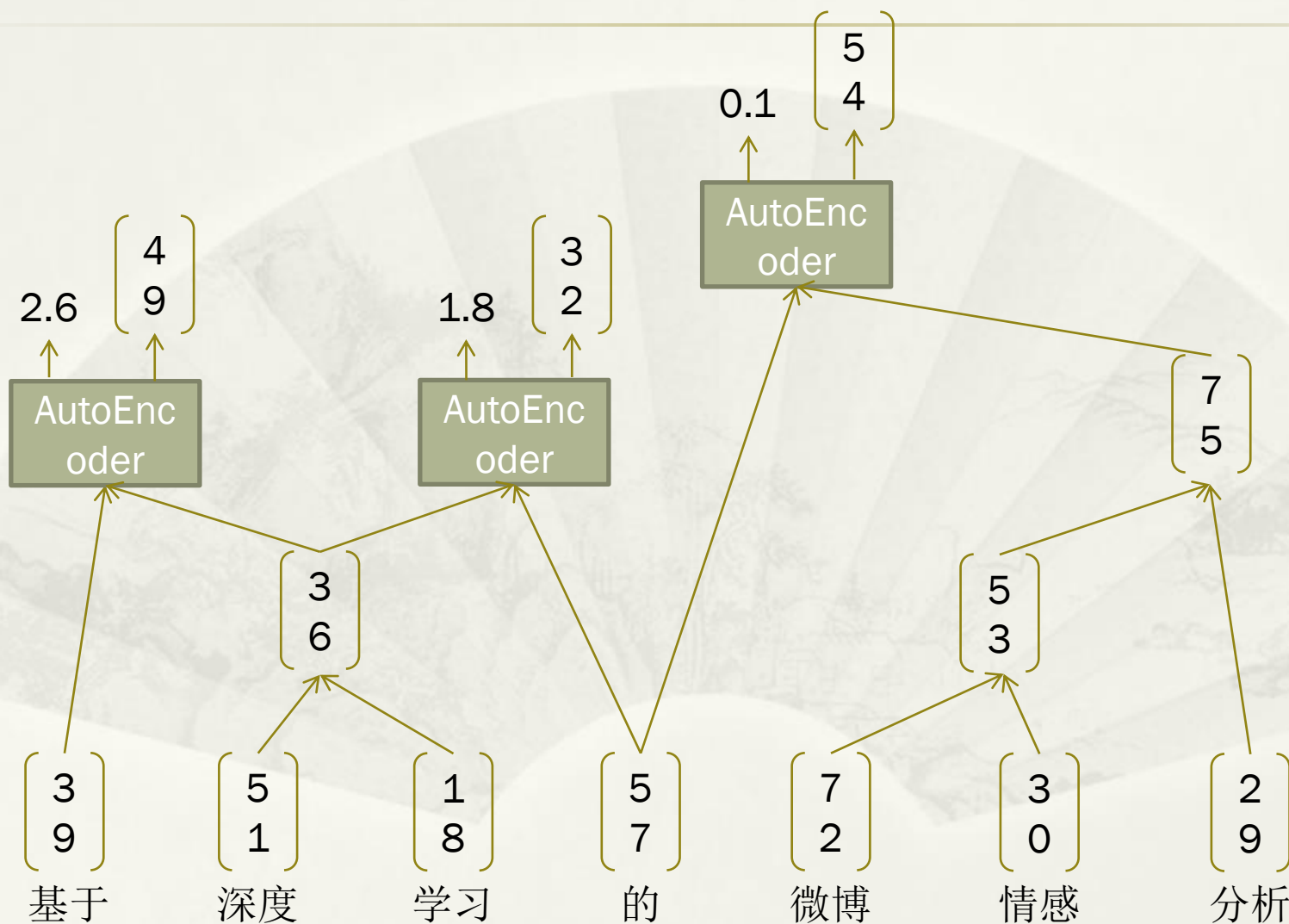
2.2 递归自编码(2/3)



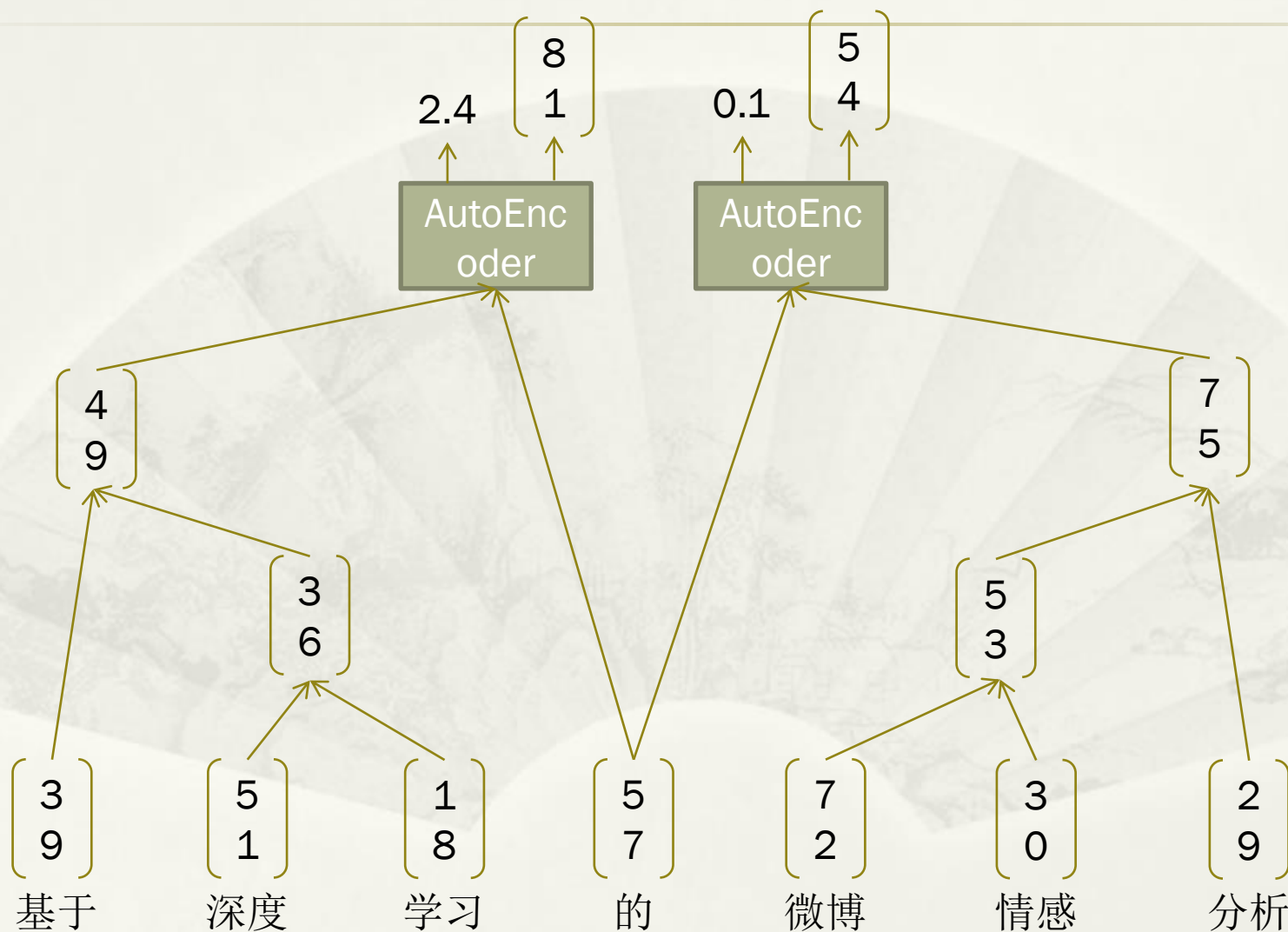
2.2 递归自编码(2/3)



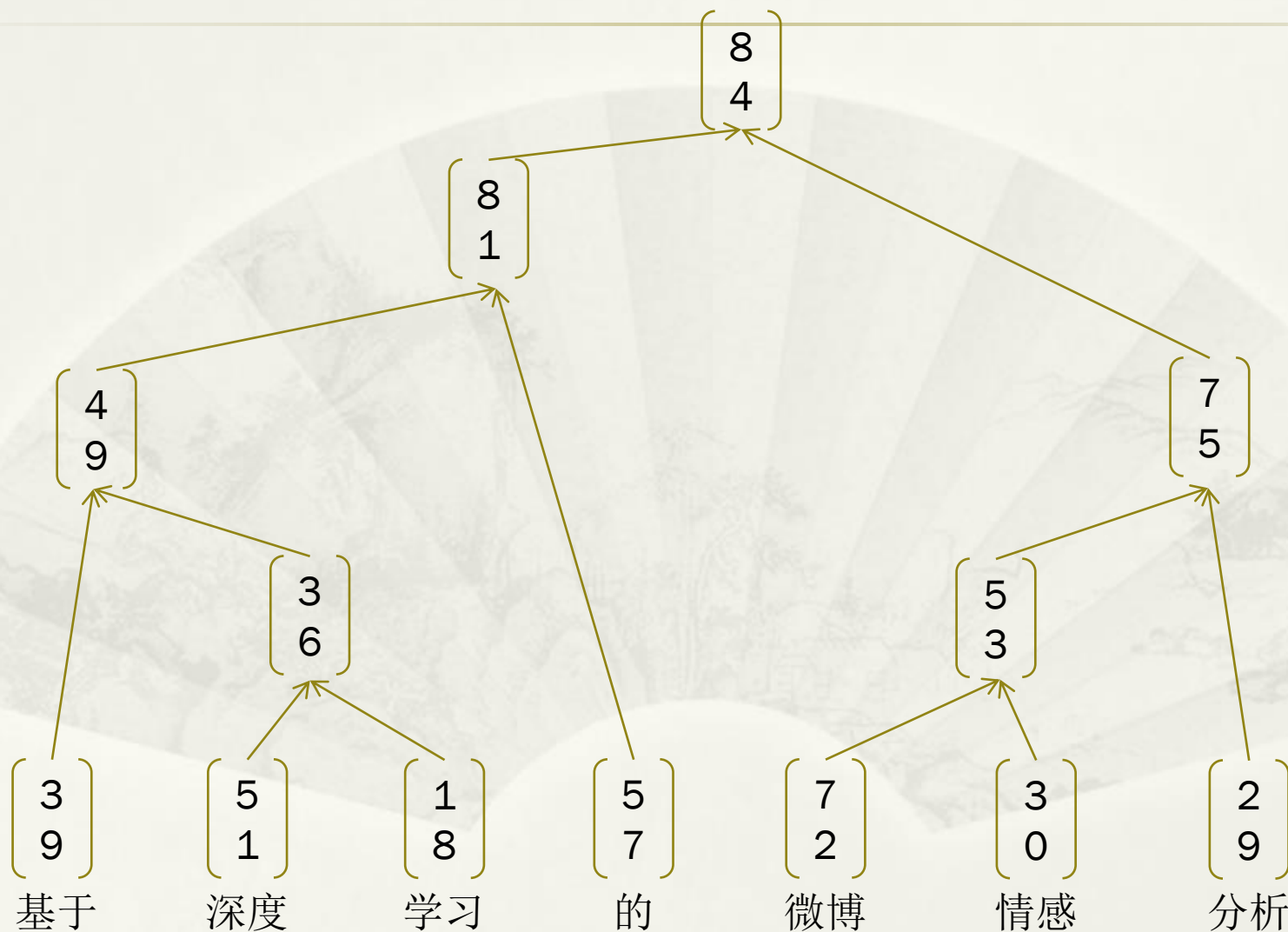
2.2 递归自编码(2/3)



2.2 递归自编码(2/3)

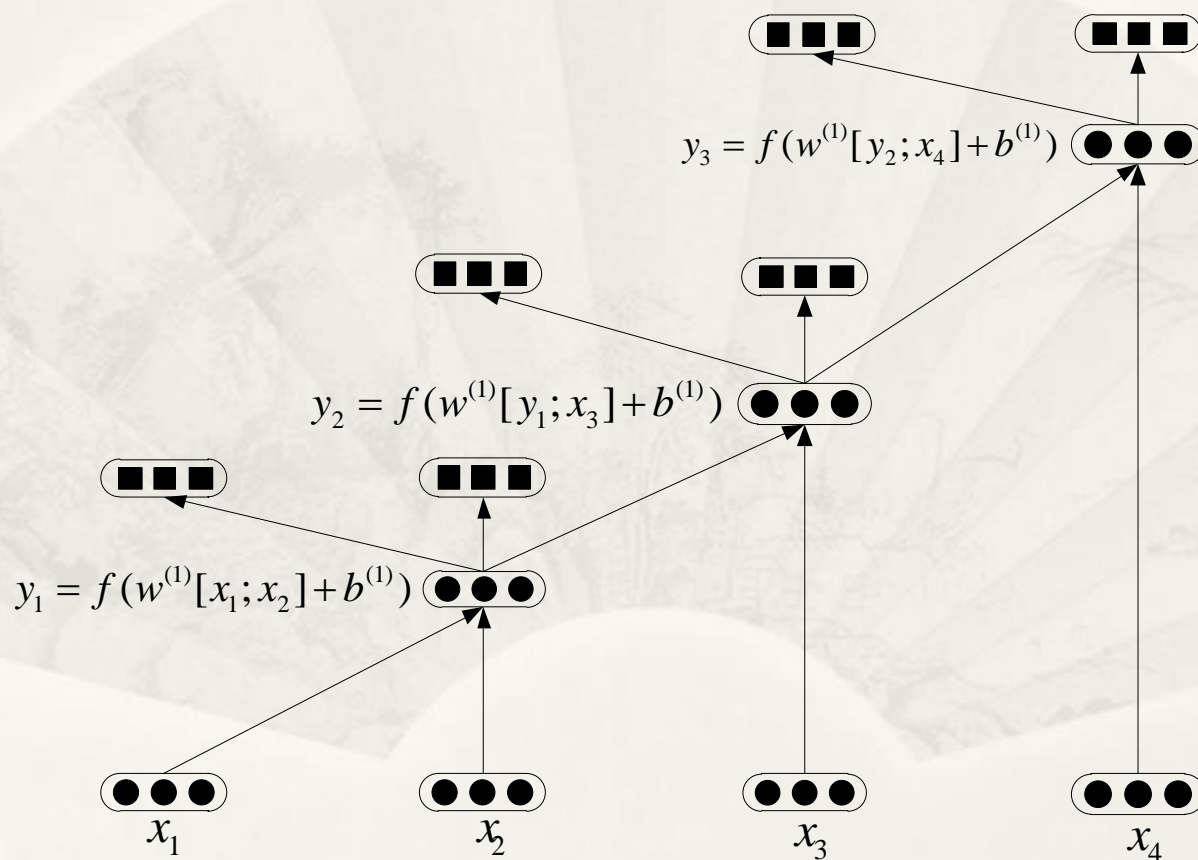


2.2 递归自编码(2/3)



2.2 递归自编码(3/3)

* 递归自编码 (Recursive AutoEncoder)





报告提纲

1. 背景介绍
2. 微博情感分析模型
 - 2.1 词语的表示
 - 2.2 递归自编码(Recursive AutoEncoder)
 - 2.3 基于RAE的情感极性转移模型
 - 2.4 情感分析模型算法
3. 实验结果分析
4. 总结

2.3 基于RAE的情感极性转移模型 (1/3)

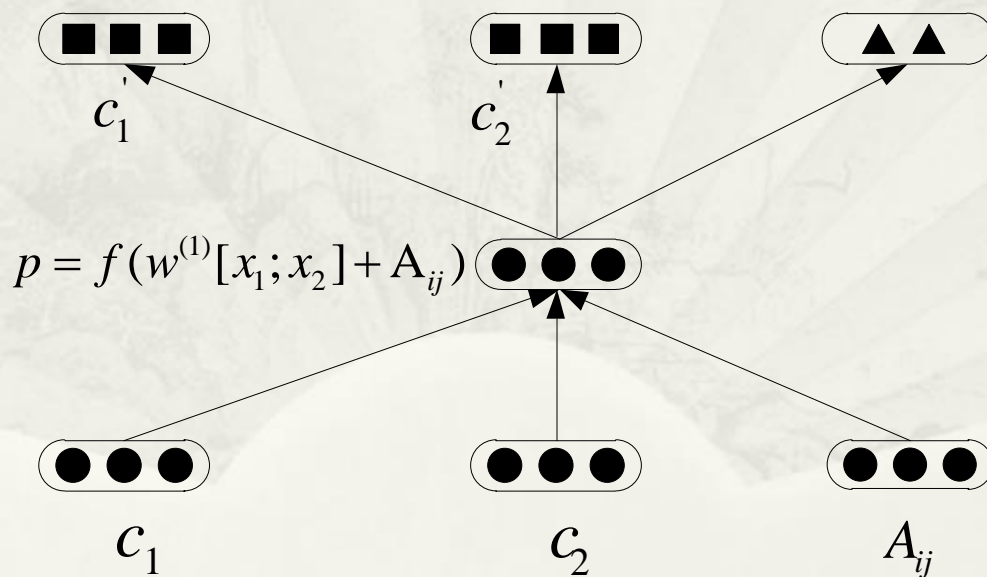
* 情感极性转移现象

极性转移类型实例	
正向+正向	她很漂亮！
正向+负向	他很自私。
负向+正向	他学习不好。
负向+负向	他学习不坏。

2.3 基于RAE的情感极性转移模型 (2/3)

* 情感极性转移模型

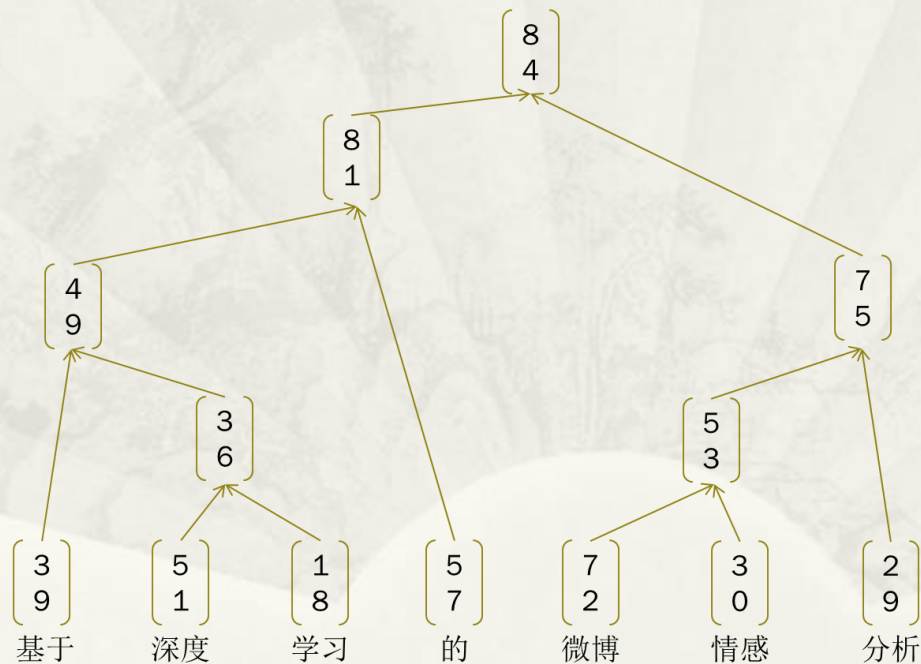
假定情感标注集 $T=\{\text{正}, \text{负}\}$ ，当连续两个词语的情感标签由 T 中的第 i 个 t_i 变为第 j 个 t_j 时，引入转换分数 A_{ij} ($A_{ij} \in \{A_{01}, A_{00}, A_{11}, A_{10},\}$)。



2.3 基于RAE的情感极性转移模型 (3/3)

* 模型输出

$$J = \text{softmax}(\lambda \cdot x + b), \quad \text{softmax} = \frac{e^{\theta_j^T x^{(i)}}}{\sum_{l=1}^k e^{\theta_l^T x^{(i)}}}$$





报告提纲

1. 背景介绍
2. 微博情感分析模型
 - 2.1 词语的表示
 - 2.2 递归自编码(Recursive AutoEncoder)
 - 2.3 基于RAE的情感极性转移模型
 - 2.4 情感分析模型算法
3. 实验结果分析
4. 总结

2.4 情感分析模型算法

情感分析模型

输入：训练语料 S 及其对应标签 t

输出： $\theta = \{w^{(1)}, w^{(2)}, w^{(3)}, A, b^{(2)}, b^{(3)}\}$

1) 使用高斯分布初始化训练语料的词向量表示；

2) **while** 不收敛 **do**

$\nabla J = 0$

for all $\langle s, t \rangle \in S$ **do**

 使用贪心算法生成句子二叉树结构

 计算 $\nabla J_i = \partial J(s, t) / \partial \theta$

 更新 $\nabla J = \nabla J + \nabla J_i$

end for

 更新 $\theta = \frac{1}{N} \nabla J + \lambda \theta$

end while

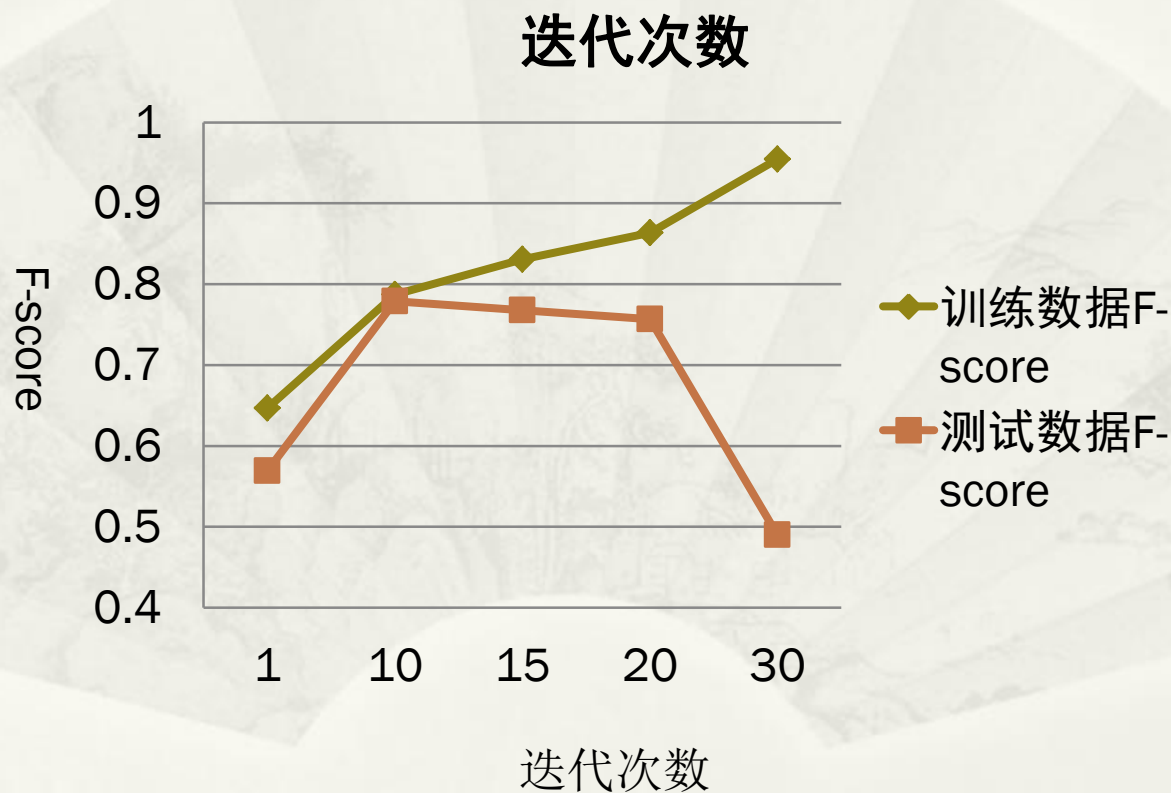


报告提纲

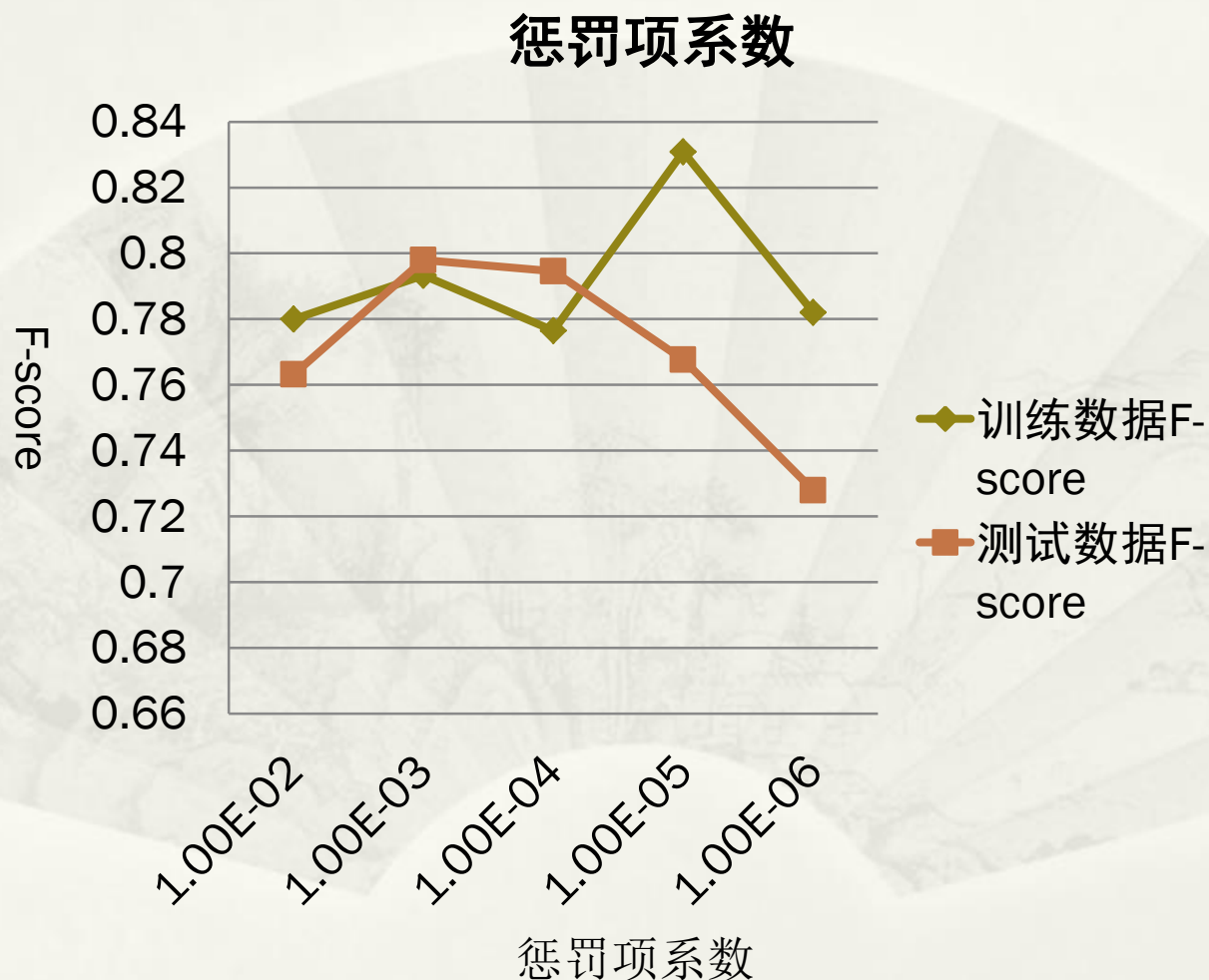
1. 背景介绍
2. 微博情感分析模型
 - 2.1 词语的表示
 - 2.2 递归自编码(Recursive AutoEncoder)
 - 2.3 基于RAE的情感极性转移模型
 - 2.4 情感分析模型算法
3. 实验结果分析
4. 总结

3. 实验结果分析(1/3)

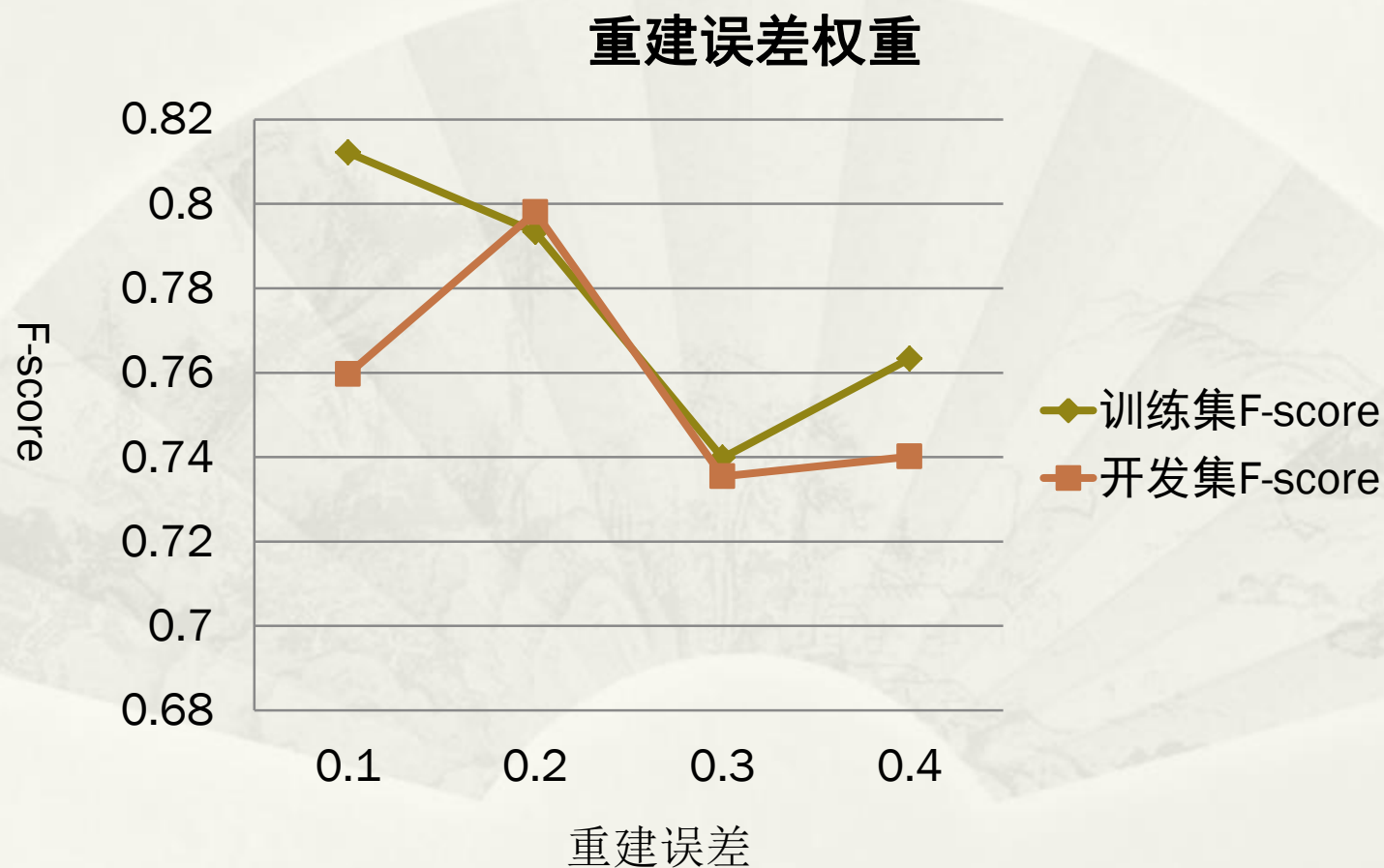
* 实验一：可调参数选择



3. 实验结果分析(1/3)

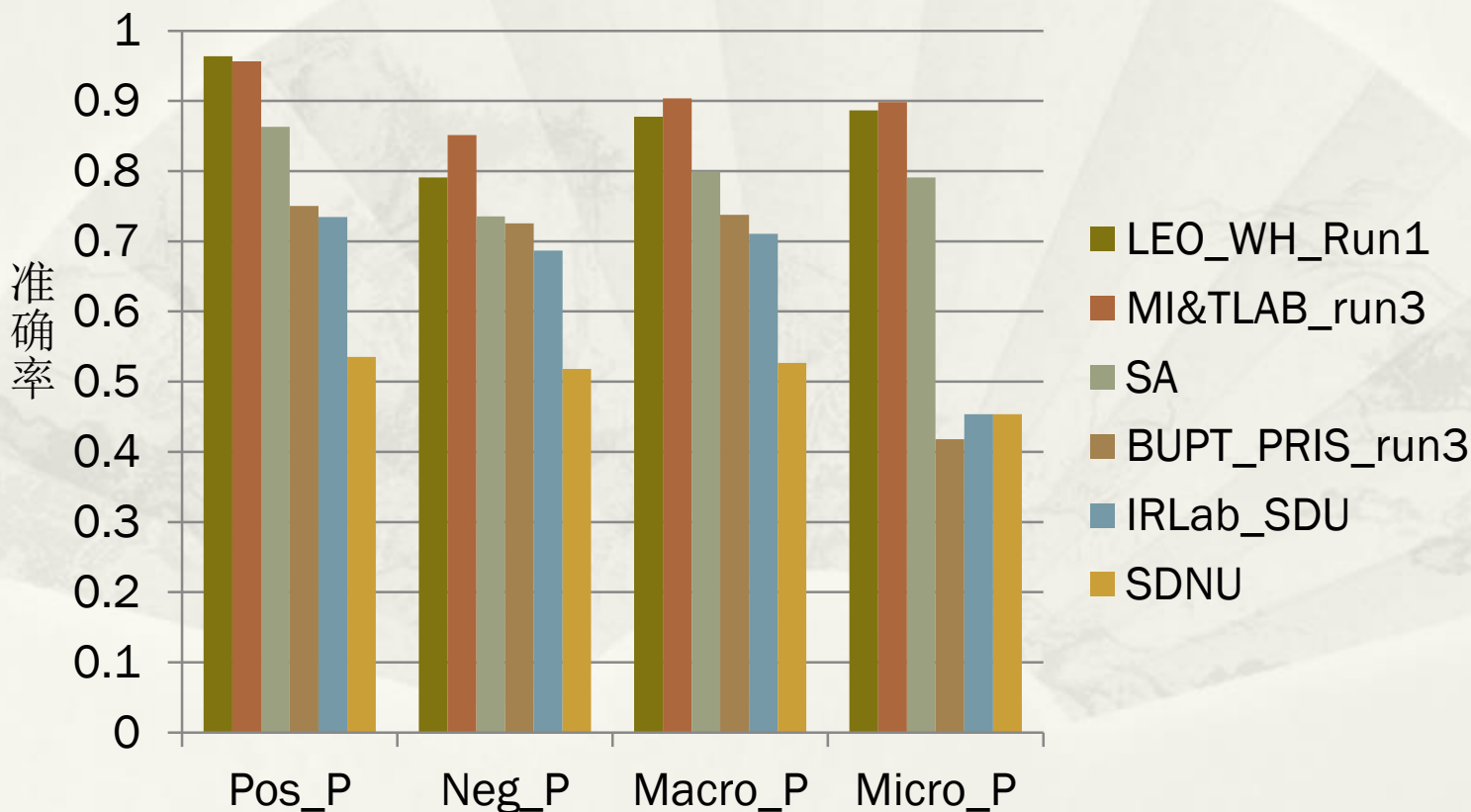


3. 实验结果分析(1/3)



3. 实验结果分析(2/3)

* 实验二：COAE2014数据集封闭测试



3. 实验结果分析(3/3)

* 实验三：对比试验

参数设定

可调参数	值
词向量维数	50
迭代次数	10
正则惩罚项系数	1e-03
重建误差与交叉熵误差权重比	0.2

3. 实验结果分析(3/3)

模型结果对比

系统	准确率(%)	召回率(%)	F-Score(%)
RAE	78.75	78.44	77.88
情感极性转移模型	80.59	80.59	79.80



报告提纲

1. 背景介绍
2. 微博情感分析模型
 - 2.1 词语的表示
 - 2.2 递归自编码(Recursive AutoEncoder)
 - 2.3 基于RAE的情感极性转移模型
 - 2.4 情感分析模型算法
3. 实验结果分析
4. 总结

4. 总结

- * 本文将深度学习方法应用到中文微博情感分析达到与使用手工标注特征相当的水平。
- * 根据语言语法现象改进现有模型提出情感极性转移模型，较原模型效果有一定提升。

4. 进一步工作(1/2)

* 中文中的语法现象

否定：机器的智能**不会**超过人脑。

加强：深度学习的应用会**更加**广泛。

比较：小明**比**小红的个子高。

* 中文中的多义现象

例如：先帝不以臣**卑鄙**。

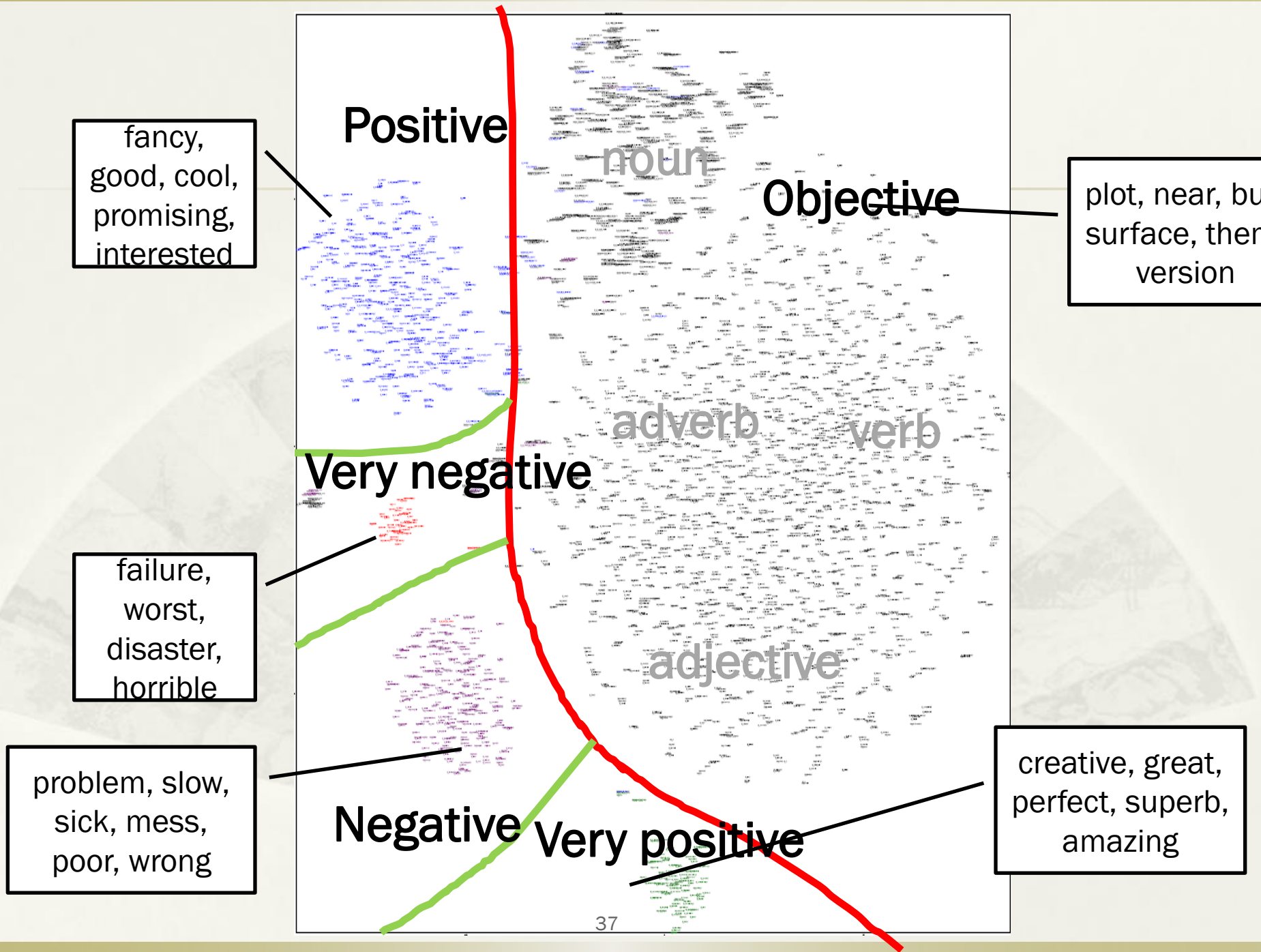
——摘自 诸葛亮《前出师表》

一些势利小人**卑鄙**，惯会兴风作浪，挑拨离间。

——摘自 茅盾《霜叶红似二月花》

4. 进一步工作(2/2)

* 深度学习在聚类上的应用



敬请各位老师同学指正！
谢谢！

梁军

邮箱：*zhengdaxg@163.com*

郑州大学自然语言处理实验室

实验室网站：*http://nlp.zzu.edu.cn*