



西安交通大学
XI'AN JIAOTONG UNIVERSITY

自然语言理解与机器翻译课程

词性标注与语言模型

李辰

cli@xjtu.edu.cn

2024年9月

2/3

概率基础



语言模型



NLU任务



课程提纲

一、词性标注

二、语言模型



词性标注

概念

- Part-of-speech (POS) tagging
- 为分词结果中的每个单词标注一个正确的词性。即确定每个词是名词、动词、形容词或其他词性的过程。
- 是高级自然语言处理任务的基础。

举例

人民网/nz 1月1日/t 讯/ng 据/p 《/w [纽约/nsf 时报/n]/nz 》 /w 报道/v , /w 美国 /nsf 华尔街/nsf 股市/n 在/p 2013年/t 的/ude1 最后/f 一天/mq 继续/v 上涨/vn , /w 和/cc [全球/n 股市/n]/nz 一样/uyy , /w 都/d 以/p [最高/a 纪录/n]/nz 或/c 接近/v [最高/a 纪录/n]/nz 结束/v 本/rz 年/qt 的/ude1 交易/vn 。 /w

Mr./NNP Vinken/NNP is/VBZ chairman/NN of/IN Elsevier/NNP N.V./NNP , / ,
the/DT Dutch/NNP publishing/VBG group/NN

汉语词性对照表

代码	名称	说明	举例
a	形容词	取英语形容词adjective的第1个字母 直接作状语的形容词.形容词代码a和副词代码d并在一起	最/d 大/a 的/u
ad	副形容词	形容词性语素.形容词代码为a, 语素代码g前面置以a	一定/d 能够/v 顺利/ad 实现/v 。/w
ag	形容词	具有名词功能的形容词.形容词代码a和名词代码n并在一起	喜/v 煞/ag 人/n 人民/n 的/u 根本/a 利益/n 和/c 国家/n 的/u 安稳/an 。/w
an	名词		
b	区别词	取汉字“别”的声母	副/b 书记/n 王/nr 思齐/nr
c	连词	取英语连词conjunction的第1个字母 取adverb的第2个字母, 因其第1个字母已用于形容词	全军/n 和/c 武警/n 先进/a 典型/n 代表/n
d	副词	副词性语素.副词代码为d, 语素代码g前面置以d	两侧/f 台柱/n 上/ 分别/d 雄踞/v 着/u 用/v 不/d 甚/dg 流利/a 的/u 中文/nz 主持/v 节目/n 。/w
dg	副语素		
e	叹词	取英语叹词exclamation的第1个字母	嗨/e ! /w
f	方位词	取汉字“方”的声母	从/p 一/m 大/a 堆/q 档案/n 中/f 发现/v 了/u
g	语素	绝大多数语素都能作为合成词的“词根”, 取汉字“根”的声母	例如dg 或ag
h	前接成分	取英语head的第1个字母	目前/t 各种/r 非/h 合作制/n 的/u 农产品/n
i	成语	取英语成语idiom的第1个字母	提高/v 农民/n 讨价还价/i 的/u 能力/n 。/w
j	简称略语	取汉字“简”的声母	民主/ad 选举/v 村委会/j 的/u 工作/vn
k	后接成分	权责/n 明确/a 的/u 逐级/d 授权/v 制/k 习用语尚未成为成语, 有点“临时性”, 取“临”的声母	是/v 建立/v 社会主义/n 市场经济/n 体制/n 的/u 重要/a 组成部分/l 。/w
l	习用语	取英语numeral的第3个字母, n, u 已有他用	科学技术/n 是/v 第一/m 生产力/n
m	数词		
n	名词	取英语名词noun的第1个字母 名词性语素.名词代码为n, 语素代码g前面置以n	希望/v 双方/n 在/p 市政/n 规划/vn
ng	名语素	名词代码n和“人(ren)”的声母并在一起	就此/d 分析/v 时/ng 认为/v
nr	人名		建设部/nt 部长/n 侯/nr 捷/nr
ns	地名	名词代码n和处所词代码s并在一起	北京/ns 经济/n 运行/vn 态势/n 喜人/a “团”的声母为t, 名词代码n和t并在一起
nt	机构团体		冶金/n 工业部/n 洛阳/ns 耐火材料/l 研究院/njnt
nx	字母专名		A T M/nx 交换机/n

nz	其他专名	“专”的声母的第1个字母为z, 名词代码n和z并在一起	德士古/nz 公司/n
o	拟声词	取英语拟声词onomatopoeia的第1个字母	汨汨/o 地/u 流/v 出来/v
p	介词	取英语介词prepositional的第1个字母	往/p 基层/n 跑/v 。/w 不止/v 一/m 次/q 地/u 听到/v , /w
q	量词	取英语quantity的第1个字母 取英语代词pronoun的第2个字母, 因p已用于介词	有些/r 部门/n
r	代词		
s	处所词	取英语space的第1个字母	移居/v 海外/s 。/w
t	时间词	取英语time的第1个字母 时间词性语素.时间词代码为t, 在语素的代码g前面置以t	当前/t 经济/n 社会/n 情况/n
tg	时语素	取英语助词auxiliary 的第2个字母, 因a已用于形容词	秋/Tg 冬/tg 连/d 早/a
u	助词		工作/vn 的/u 政策/n
ud	结构助词		有/v 心/n 裁/v 得/ud 梧桐树/n
ug	时态助词		你/r 想/v 过/ug 没有/v 迈向/v 充满/v 希望/n 的/uj 新/a 世纪/n
uj	结构助词的		
ul	时态助词了		完成/v 了/ ul 满怀信心/l 地/uv 开创/v 新/a 的/u 业绩/n
uv	结构助词地		
uz	时态助词着		眼看/v 着/uz 举行/v 老/a 干部/n 迎春/vn 团拜会/n
v	动词		
vd	副动词	动词性语素.动词代码为v. 在语素的代码g前面置以V	强调/vd 指出/v 做好/v 尊/vg 干/j 爱/v 兵/n 工作/vn
vg	动语素	指具有名词功能的动词.动词和名词的代码并在一起	股份制/n 这种/r 企业/n 组织/vn 形式/n , /w
vn	名动词		生产/v 的/u 5 G/nx 、/w 8 G /nx 型/k 燃气/n 热水器/n
w	标点符号	非语素字只是一个符号, 字母x通常用于代表未知数、符号	已经/d 3 0/m 多/m 年/q 了/y 。/w
x	非语素字		
y	语气词	取汉字“语”的声母	
z	状态词	取汉字“状”的声母的前一个字母	势头/n 依然/z 强劲/a ; /w

英语词性对照表

Tag	Description	Example
CC	conjunction, coordinating	and, or, but
CD	cardinal number	five, three, 13%
DT	determiner	the, a, these
EX	existential there	there were six boys
FW	foreign word	mais
IN	conjunction, subordinating or preposition	of, on, before, unless
JJ	adjective	nice, easy
JJR	adjective, comparative	nicer, easier
JJS	adjective, superlative	nicest, easiest
LS	list item marker	
MD	verb, modal auxillary	may, should
NN	noun, singular or mass	tiger, chair, laughter
NNS	noun, plural	tigers, chairs, insects
NNP	noun, proper singular	Germany, God, Alice
NNPS	noun, proper plural	we met two Christmases ago
PDT	predeterminer	both his children
POS	possessive ending	's
PRP	pronoun, personal	me, you, it
PRP\$	pronoun, possessive	my, your, our
RB	adverb	extremely, loudly

RBR	adverb, comparative	better
RBS	adverb, superlative	best
RP	adverb, particle	about, off, up
SYM	symbol	%
TO	infinitival to	what to do?
UH	interjection	oh, oops, gosh
VB	verb, base form	think
VBZ	verb, 3rd person singular present	she thinks
	verb, non-3rd person singular	
VBP	present	I think
VBD	verb, past tense	they thought
VCN	verb, past participle	a sunken ship
VBG	verb, gerund or present participle	thinking is fun
WDT	wh-determiner	which, whatever, whichever
WP	wh-pronoun, personal	what, who, whom
WP\$	wh-pronoun, possessive	whose, whosever
WRB	wh-adverb	where, when
.	punctuation mark, sentence closer	.;?*
,	punctuation mark, comma	,
:	punctuation mark, colon	:
(contextual separator, left paren	(
)	contextual separator, right paren)

词性标注

中科院计算所分词系统

NLPIR汉语分词系统

分词ictclas - Google 搜索 × ICTCLAS张华平博士的空间_百... × NLPIR汉语分词系统 × 中文分词_百度百科 × 什么是动态规划? - 已解决 - ... × +

ictclas.nlpir.org/onlineputong 百度 <K>

昔日NBA球星罗德曼访朝开启了两国间的篮球外交？至少目前看来，美国民众的回应应有褒有贬。据《纽约邮报》3月6日报道，“大虫”罗德曼回到美国后，在下榻的酒店聊天时力挺朝鲜领导人，结果遭到一片嘘声，在保镖的护送下才离开。

选择文件

普通分词

自适应分词

清除

昔日/t NBA/x 球星/n 罗德曼/nrf 访/v 朝/tg 开启/v 了/ule 两/m 国/n 间/f 的 /ude1 篮球/n 外交/n ? /ww 至少/d 目前/t 看来/v , /wd 美国/nsf 民众/n 的 /ude1 回应/vn 有/vyou 褒/vn 有/vyou 贬/v 。 /wj 据/p 《/wkz 纽约/nsf 邮报/n 》 /wky 3月/t 6日/t 报道/v , /wd “/wyz 大虫/n” /wyy 罗德曼/nrf 回到/v 美国 /nsf 后/f , /wd 在/p 下榻/v 的/ude1 酒店/n 聊天/vi 时/ng 力/n 挺/d 朝鲜/nsf 领 领导人/n , /wd 结果/d 遭到/v 一/m 片/q 嘘声/n , /wd 在/p 保镖/n 的/ude1 护送 /vn 下/f 才/d 离开/v 。 /wj

测试样例：

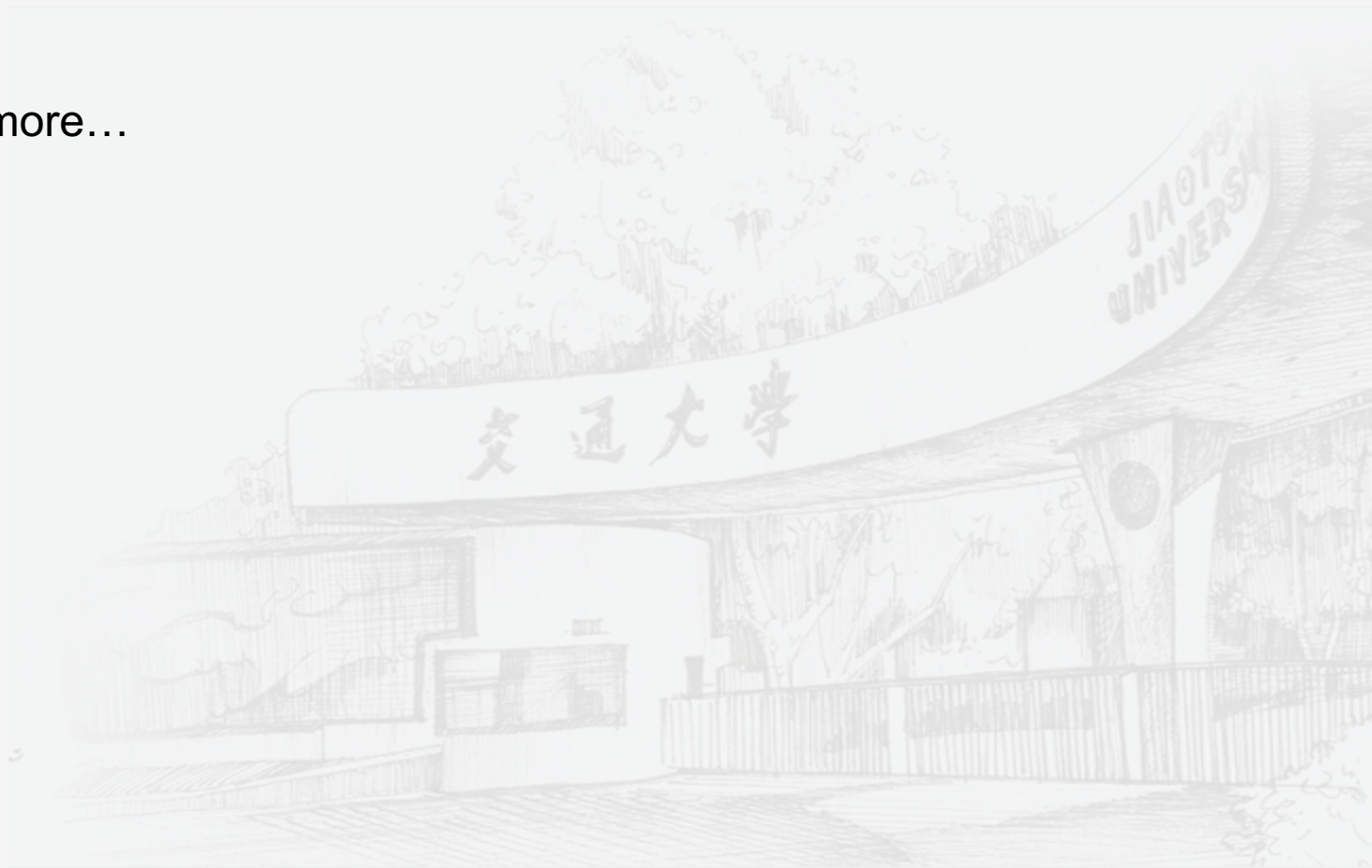
1. @ICTCLAS张华平博士 应各位ICTCLAS用户的要求，张华平博士提前发布ICTCLAS2013版本，为了与以前工作进行大的区隔，并推广NLPIR自然语言处理与信息检索共享平台，从本版本开始，系统名称调整为NLPIR汉语分词系统。

2. “屌丝”这个嘲讽意味的代词迅速爆红，迎合了大众的心理和趣味。因为你会发现从表面符合屌丝定义的人，到和屌丝属性八字打不着的人，都在争相认领这一名号。当人人都在忙着确认自己的屌丝身份，并乐此不疲时，屌丝一词一定与时代的什么特征实现了合拍。“屌丝”不是阿Q，他们公然比惨并乐在其中有评论认为，“屌丝”是新时代的阿Q，两者并不完全相同。首先，阿Q是文学巨匠鲁迅一己之力创造的，而“屌丝”则是网络群体狂欢的结果，它是真正由网民集体创作的形象；另外，阿Q最重要的特征是“精神胜利法”，梦想的是“银盔银甲”，意淫的是“我手持钢鞭将你打”。

词性标注

应用

- 统计自然语言处理基础
- 句法分析的基础
- 辅助词义消歧
- 机器翻译
- and many more...



词性标注歧义

- 一个词具有两个或者两个以上的词性
- 英文 Brown 语料库中，10.4% 的词是兼类词
 - the back door
 - on my back
 - promise to back the bill
- 汉语兼类词
 - 把门锁上 买了一把锁
 - 他研究与动物相关的研究工作
- 对兼类词消歧——词性标注的任务

词性标注

词性标注的信息来源是什么？

- 词汇信息
当前词本身提供了关于标注的信息
- 句法结构信息
考虑在当前词上下文中的词的词性

词性标注的性能指标

- 目前准确率大约在 97% 左右
- Baseline也可以达到 90%

Baseline算法：

对每一个词用它的最高频的词性进行标注

未登录词全部标为名词

马尔可夫模型标注器

原理

- W —— 标注序列, T —— 标注集, O —— 词集
- 马尔可夫模型标注器: 假定一个词的词性只依赖于前一个词的词性 (有限历史), 而且这个依赖不随时间变化 (时间不变)。

$$\operatorname{argmax}_{t_{1,\dots,n}} P(t_{1,\dots,n} \mid w_{1,\dots,n})$$

应用贝叶斯规则

$$\begin{aligned} &= \operatorname{argmax}_{t_{1,n}} \frac{P(w_{1,n} \mid t_{1,n}) P(t_{1,n})}{P(w_{1,n})} \\ &= \operatorname{argmax}_{t_{1,n}} P(w_{1,n} \mid t_{1,n}) P(t_{1,n}) \end{aligned}$$

马尔可夫模型标注器

原理

$$\begin{aligned} P(w_{1,n} | t_{1,n}) P(t_{1,n}) &= \prod_{i=1}^n P(w_i | t_{1,n}) \quad \text{词的相互独立性} \\ &\quad \times P(t_n | t_{1,n-1}) \times P(t_{n-1} | t_{1,n-2}) \times \dots \times P(t_2 | t_1) \\ &= \prod_{i=1}^n P(w_i | t_i) \quad \text{一个词的词形只依赖于它自身的词性} \\ &\quad \times P(t_n | t_{n-1}) \times P(t_{n-1} | t_{n-2}) \times \dots \times P(t_2 | t_1) \quad \text{有限历史} \\ &= \prod_{i=1}^n [P(w_i | t_i) \times P(t_i | t_{i-1})] \end{aligned}$$

最终一个句子的最优标注序列公式为：

$$\hat{t}_{1,n} = \operatorname{argmax}_{t_{1,n}} P(t_{1,n} | w_{1,n}) = \prod_{i=1}^n P(w_i | t_i) P(t_i | t_{i-1})$$

马尔可夫模型标注器

训练

- 有一个大的带标训练集
- 最大似然估计

$$\forall t^j \in tag, t^k \in tag, P(t^k | t^j) = \frac{C(t^j, t^k)}{C(t^j)}$$

$$\forall t^j \in tag, w^l \in word, P(w^l | t^j) = \frac{C(w^l, t^j)}{C(t^j)}$$

C(...) 是出现次数

课程提纲

一、词性标注（平滑技术）

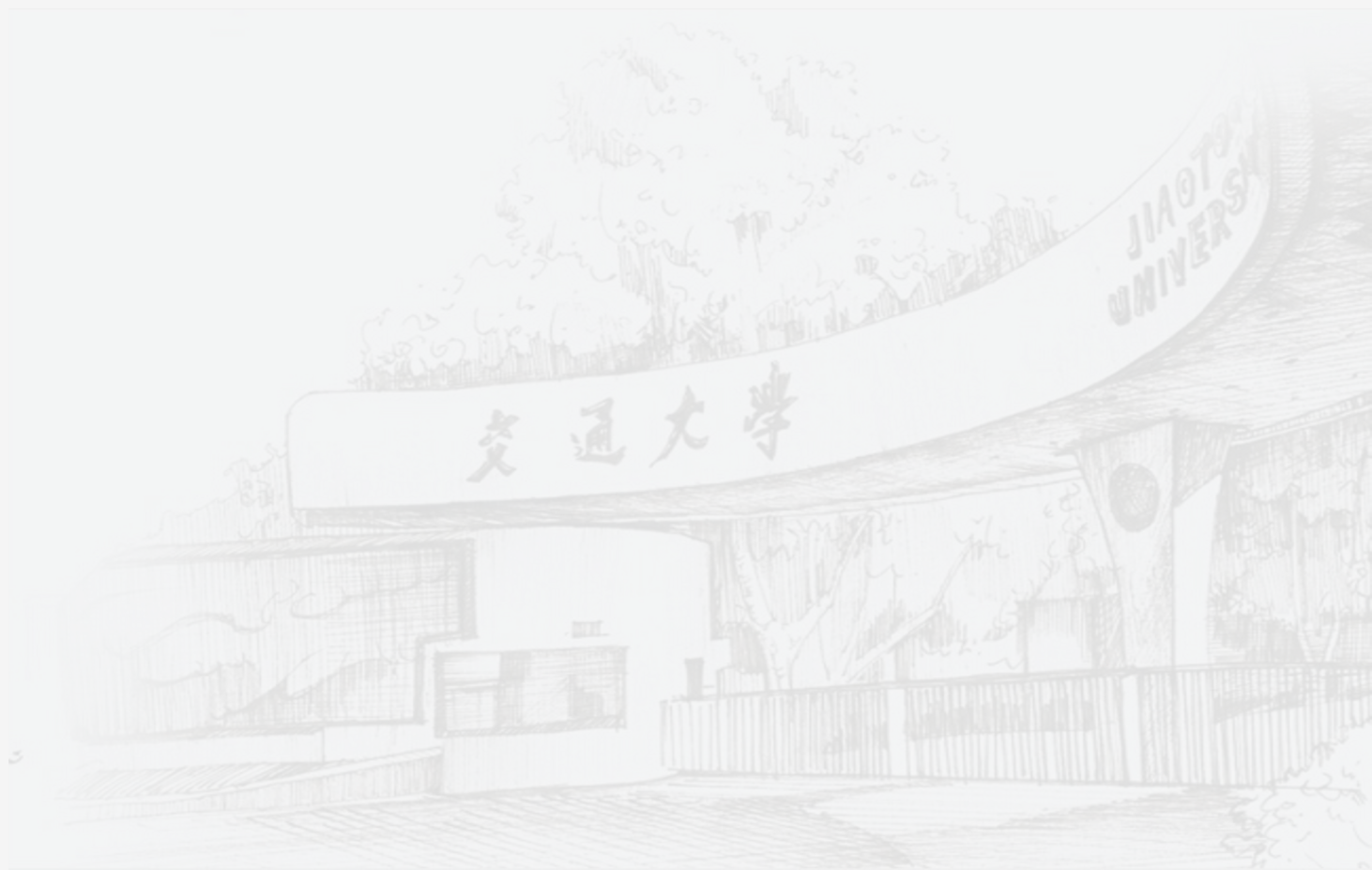
二、语言模型



平滑处理

概念

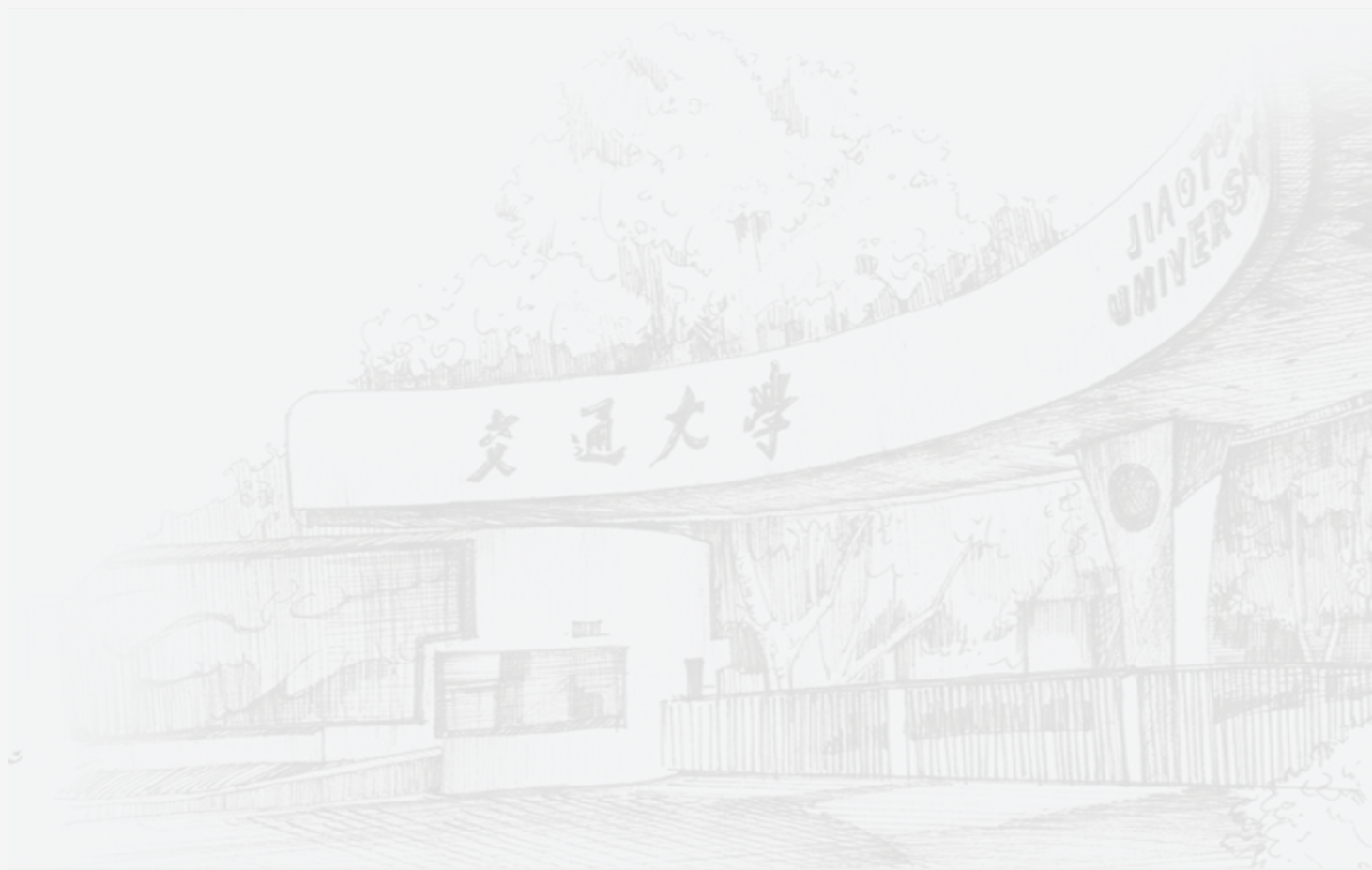
- 为什么要做平滑处理？



平滑处理

概念

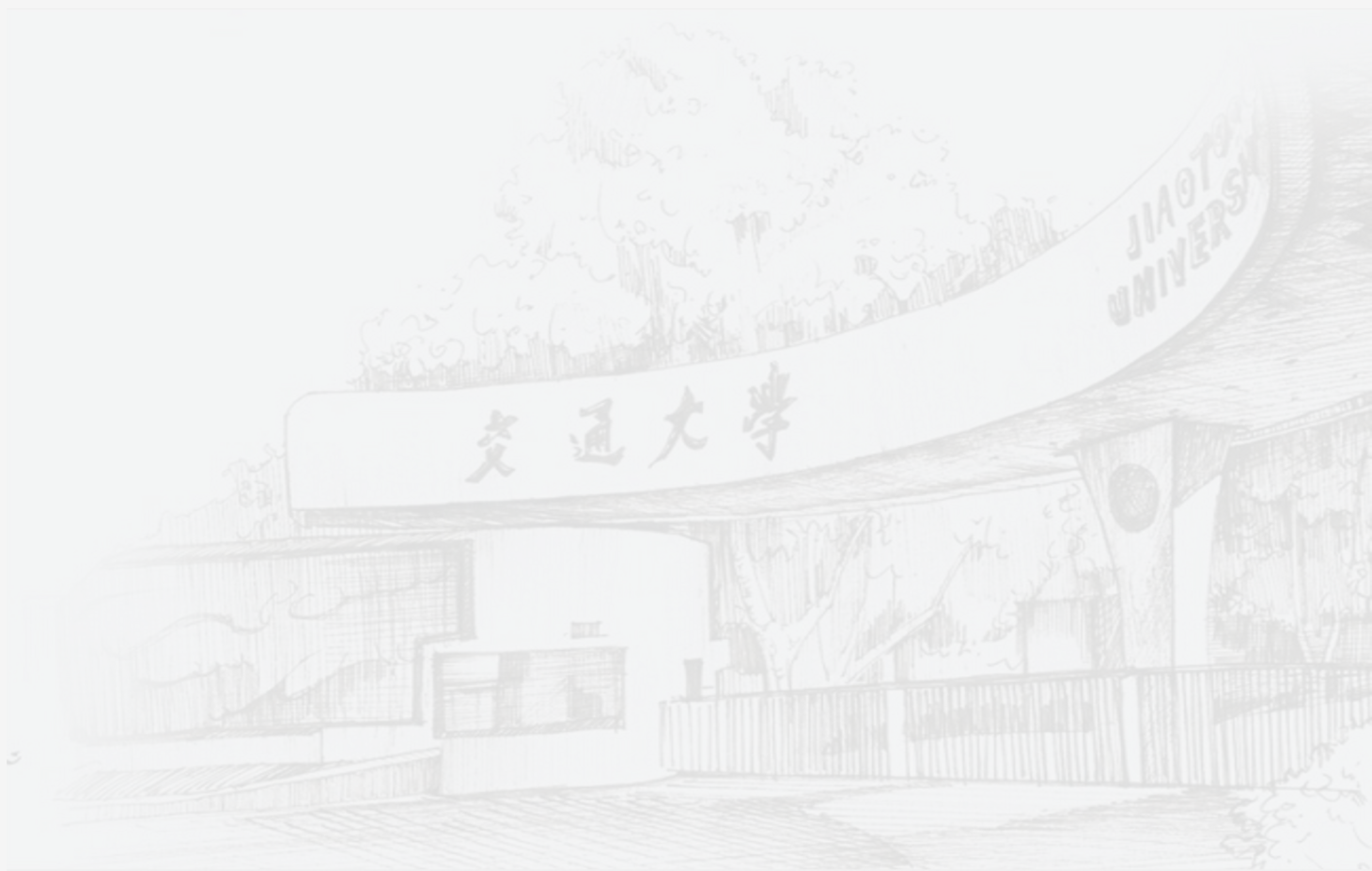
- 为什么要做平滑处理？
处理数据矩阵的稀疏问题；



平滑处理

概念

- 为什么要做平滑处理？
处理数据矩阵的稀疏问题；
- 为什么会稀疏？



平滑处理

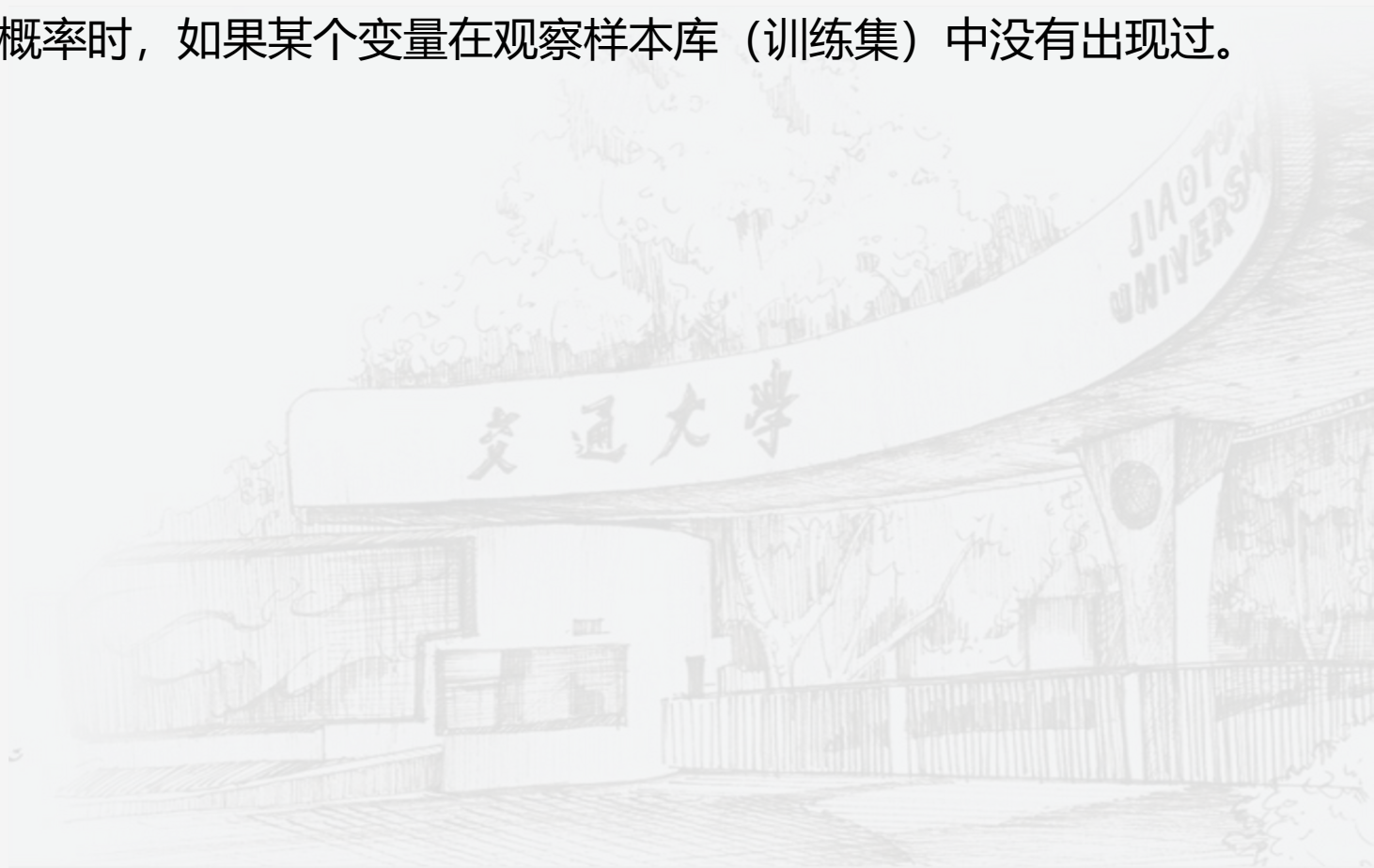
概念

- 为什么要做平滑处理？

处理数据矩阵的稀疏问题；

- 为什么会稀疏？

计算实例的概率时，如果某个变量在观察样本库（训练集）中没有出现过。



拉普拉斯平滑

概念

- Laplace smoothing 或者 Additive smoothing
- N 个样本中, 样本 i 的频率是 x_i , 得到频率观察 $X = \langle x_1, x_2, \dots, x_i, \dots, x_d \rangle$

$$p_{i, \text{ empirical}} = \frac{x_i}{N}$$



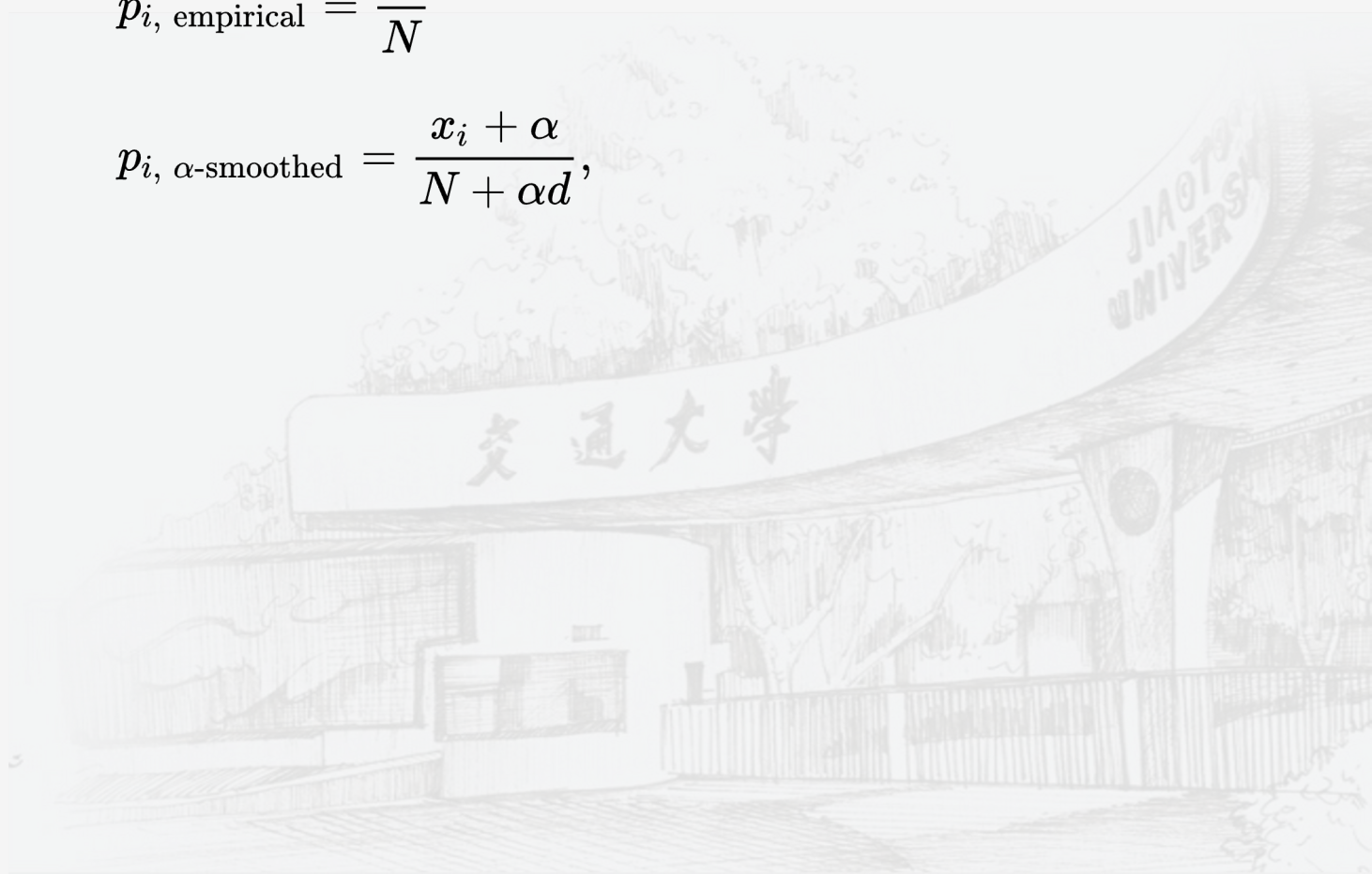
拉普拉斯平滑

概念

- Laplace smoothing 或者 Additive smoothing
- N 个样本中, 样本 i 的频率是 x_i , 得到频率观察 $X = \langle x_1, x_2, \dots, x_i, \dots, x_d \rangle$

$$p_{i, \text{empirical}} = \frac{x_i}{N}$$

$$p_{i, \alpha\text{-smoothed}} = \frac{x_i + \alpha}{N + \alpha d},$$



图灵估计

概念

- Good-Turing frequency estimation、Good-Turing 平滑法
- I. J. Good 1953 年提出。
- 基本思想：用观察计数较高的变量重新估计概率量的大小，并把它指派给那些具有零计数或者较低计数的变量。



图灵估计

例子

- 假设钓鱼过程中抓到了 18 条鱼，种类如下：
10 carp, 3 perch, 2 whitefish, 1 trout, 1 salmon, 1 eel

- 那么，下一条鱼是 trout 的概率是多少？应该是 $\frac{1}{18}$

- 那么，下一条鱼是新品种的概率是多少？

不考虑其他 0

根据只出现一次的情况来估计 $\frac{3}{18}$

- 在此基础上，下一条鱼是 trout 的概率是多少？肯定小于 $\frac{1}{18}$
- 根据 Good Turing，对每一个计数 r ，做一个调整，变为 r^*

$$r^* = (r + 1) \frac{n_{r+1}}{n_r}$$

- 所以， $c = 1$ 时，

$$C^*(trout) = 2 \times \frac{1}{3}; \quad P^*(trout) = \frac{\frac{2}{3}}{18} = \frac{1}{27}$$

课程提纲

一、词性标注（继续）

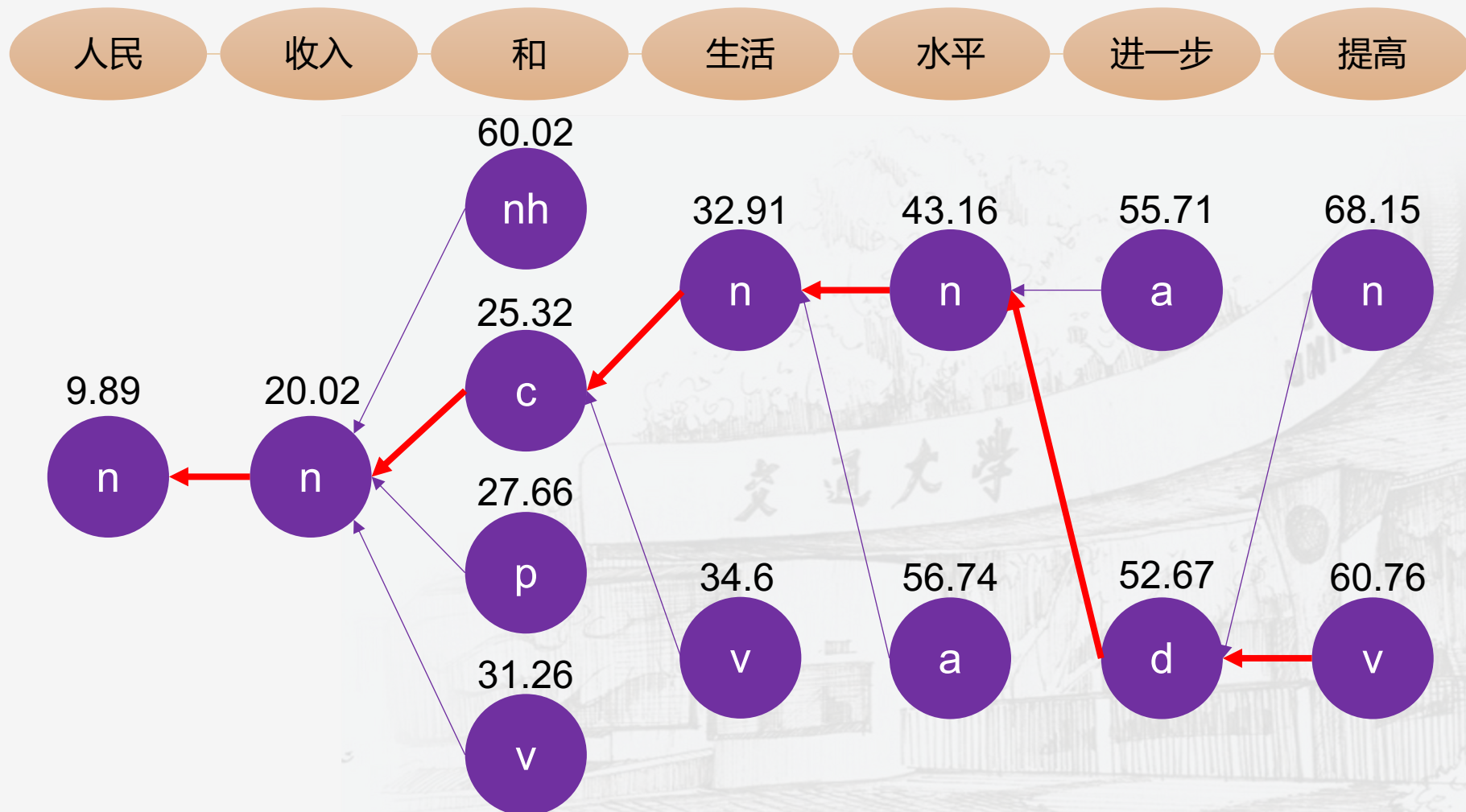
二、语言模型



马尔可夫模型标注器

预测

- 如何求全局最优解



马尔可夫模型标注器

算法变形

- 未登录词
 - 考虑所有词性
 - 只考虑开放类词性
 - Uniform (平均分配概率)
 - Unigram (考虑每个词性独立出现的概率)
 - 根据未登录词的前缀和后缀猜测其词性
- 平滑
 - 平滑的原因：数据稀疏！
 - 收集更多的数据
 - 平滑
 - 估计在训练样本中没有出现情况的概率，降低已出现情况的概率来“分给”未出现的情况

马尔可夫模型标注器

注意

- 在训练时，我们能够观察到马尔可夫模型的状态，但是在标注时我们只能观察到词。所以我们说在马尔可夫模型标注时，实际上使用的是一个混合的方法。
- 在训练时构造VMMs，但是在标注时把他们当作HMMs。

为什么不直接称其为隐马尔科夫标注器？

- 隐马尔可夫标注器标注过程和马尔可夫标注器相同。
- 两者的区别在于怎样训练模型。
- 如果没有足够的训练语料，可以使用HMM来学习标注序列（使用前向后向算法）

基于转换的词性标注

基本思想

- 正确结果是通过不断修正错误得到的。
- 修正错误的过程是有迹可寻的。
- 让计算机学习修正错误的过程，这个过程可以用转换规则形式记录下来，然后用学习得到的转换规则进行词性标注。



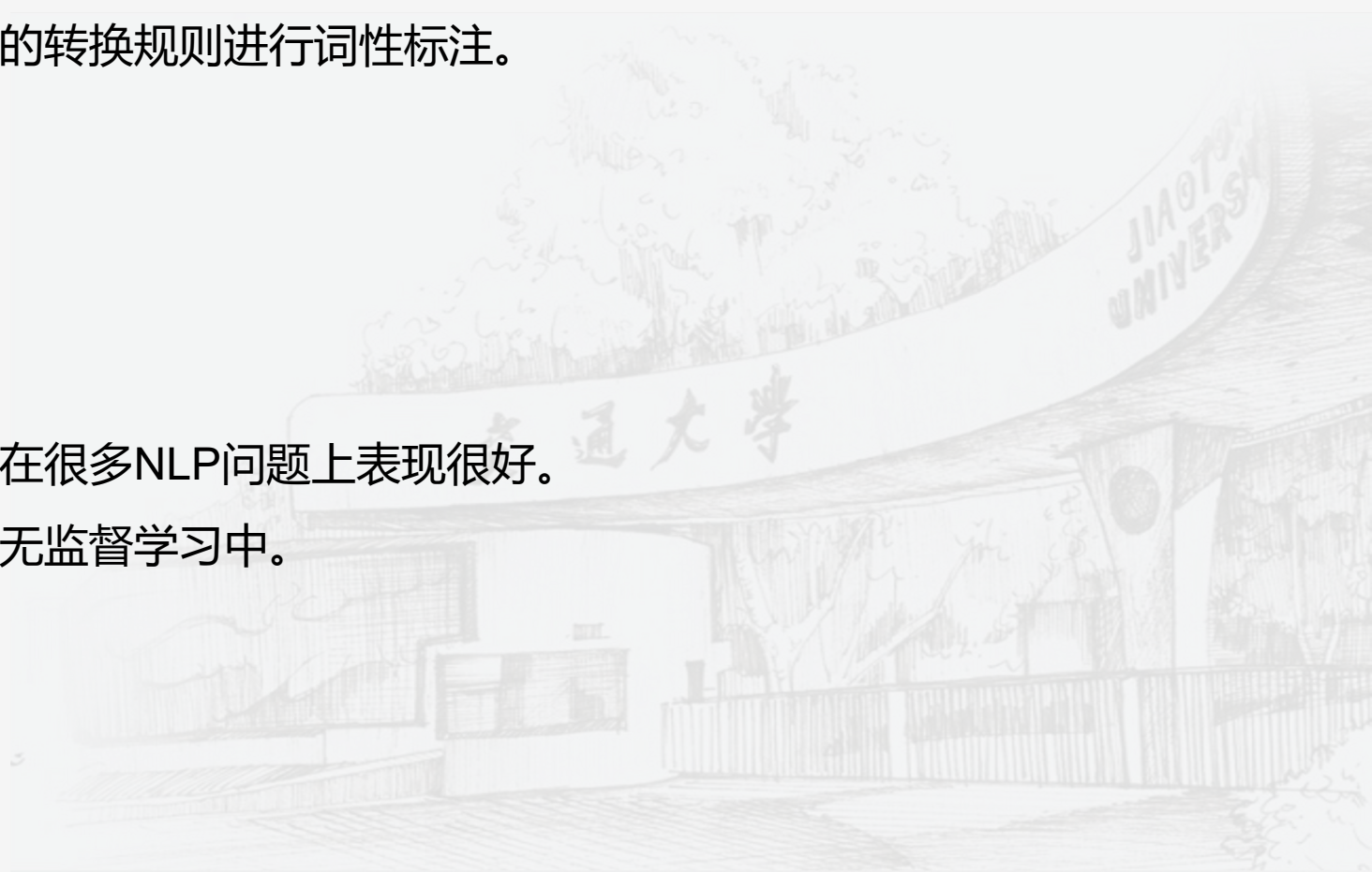
基于转换的词性标注

基本思想

- 正确结果是通过不断修正错误得到的。
- 修正错误的过程是有迹可寻的。
- 让计算机学习修正错误的过程，这个过程可以用转换规则形式记录下来，然后用学习得到的转换规则进行词性标注。

方法评价

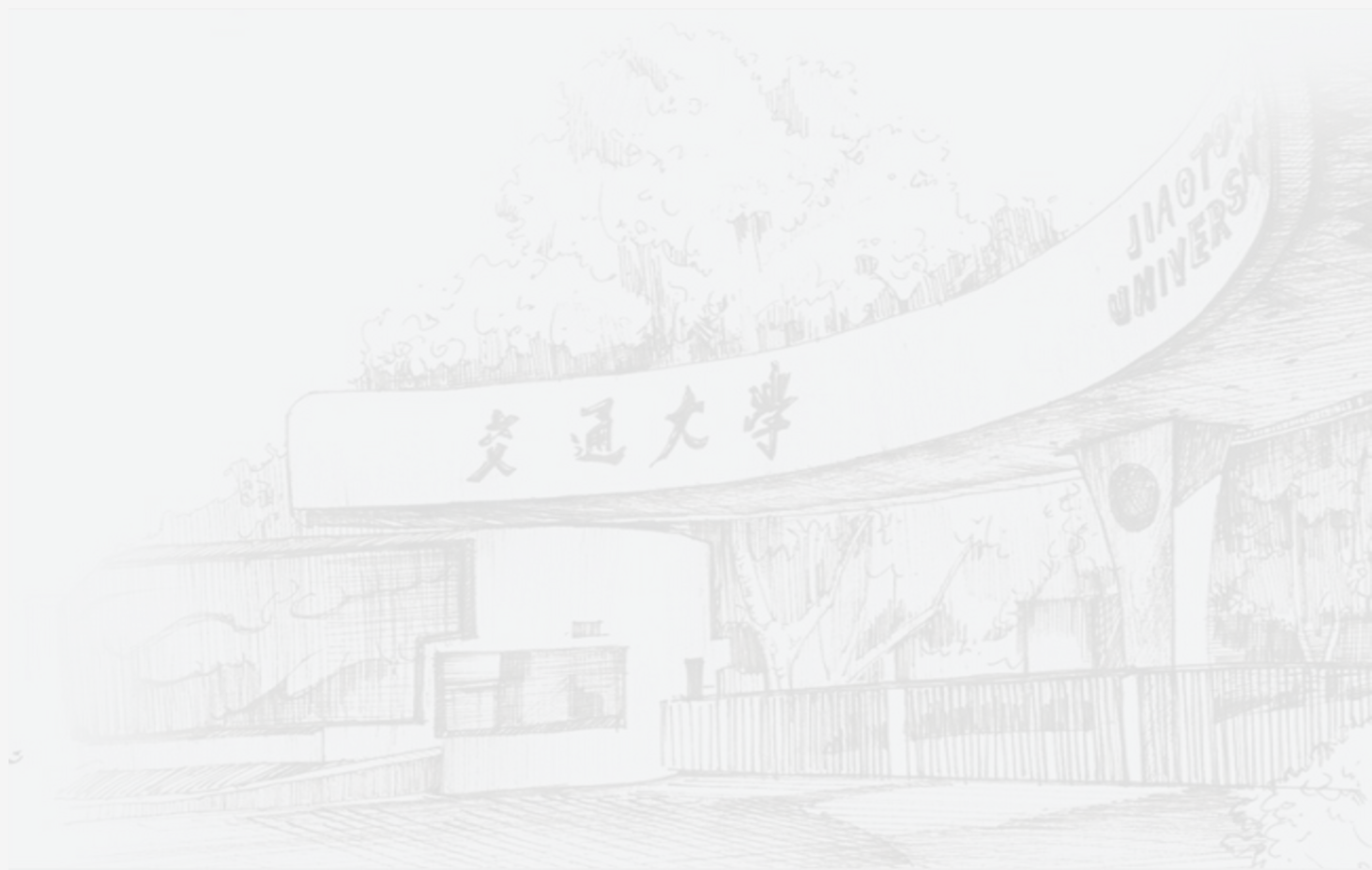
- 效果很好，在很多NLP问题上表现很好。
- 可以扩展到无监督学习中。



转换规则的形式

转换规则的组成

- 改写规则：将一个词性转换成另一个词性
- 激活环境：激发改写规则的条件



转换规则的形式

转换规则的组成

- 改写规则：将一个词性转换成另一个词性
- 激活环境：激发改写规则的条件

转换规则(T1)

- 改写规则：将一个词性从动词(v)改为名词(n)
- 激活环境：该词左边第一个词的词性是量词(q)

第二个词的词性是数词(m);

举例：

他/r 做/v 了/u 一/m 个/q 报告/v

他/r 做/v 了/u 一/m 个/q 报告/n

转换规则的形式

转换规则的组成

- 改写规则：将一个词性转换成另一个词性
- 激活环境：激发改写规则的条件

转换规则(T1)

自然科学和社会科学的一个链接

- 改写规则：将一个词性从动词(V)改为名词(n)
- 激活环境：该词左边第一个词的词性是量词(q)

第二个词的词性是数词(m);

举例：

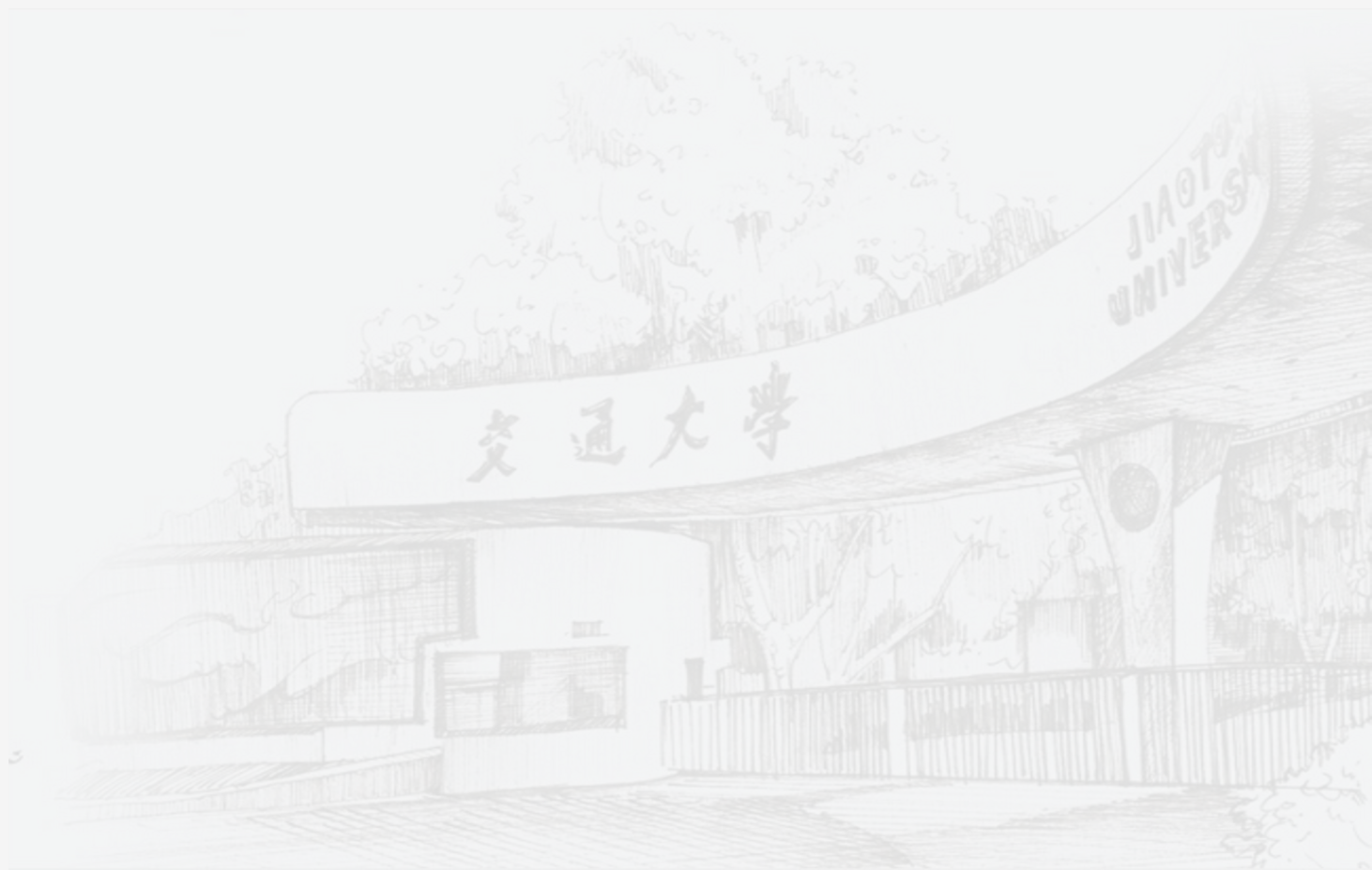
他/r 做/v 了/u 一/m 个/q 报告/v

他/r 做/v 了/u 一/m 个/q 报告/n

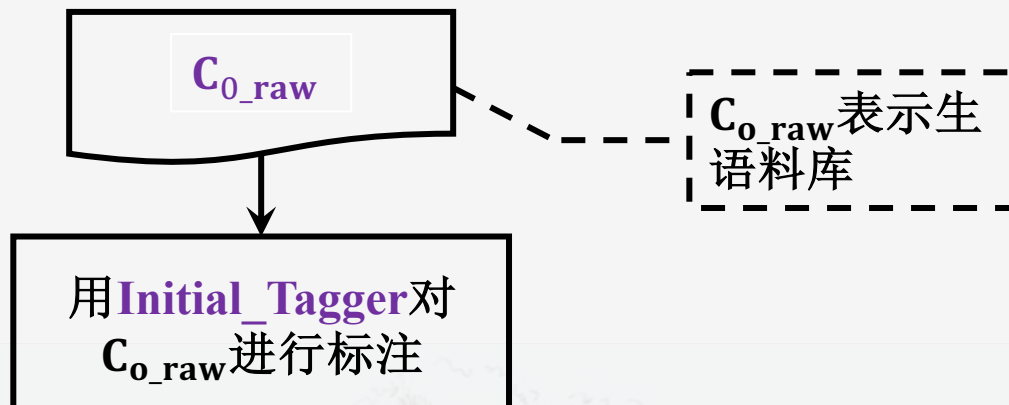
转换规则的学习流程

C_{0_raw}

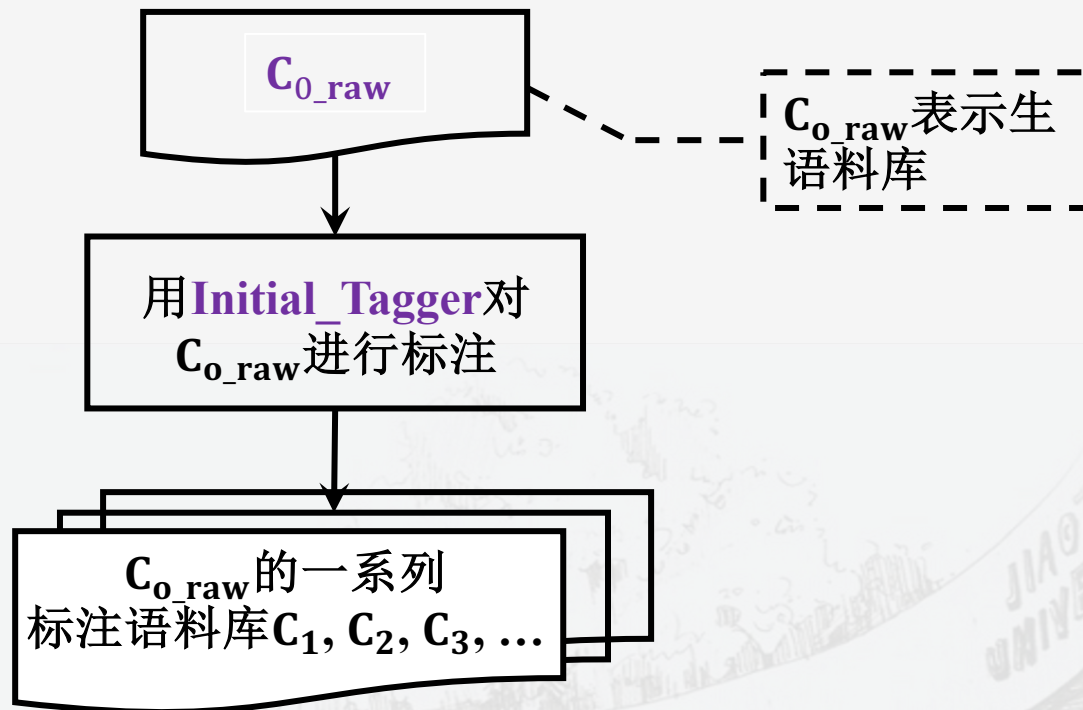
C_{0_raw} 表示生
语料库



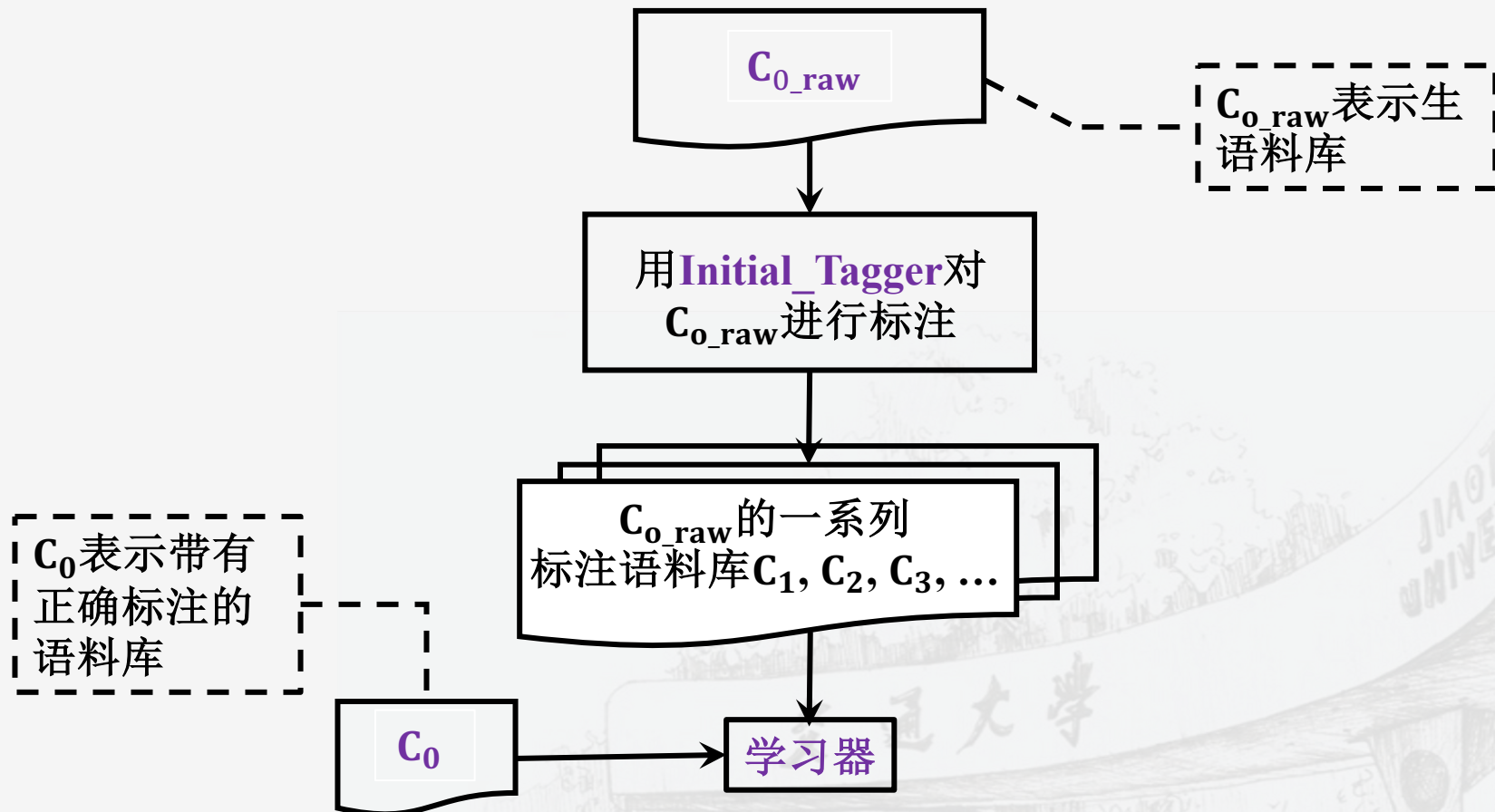
转换规则的学习流程



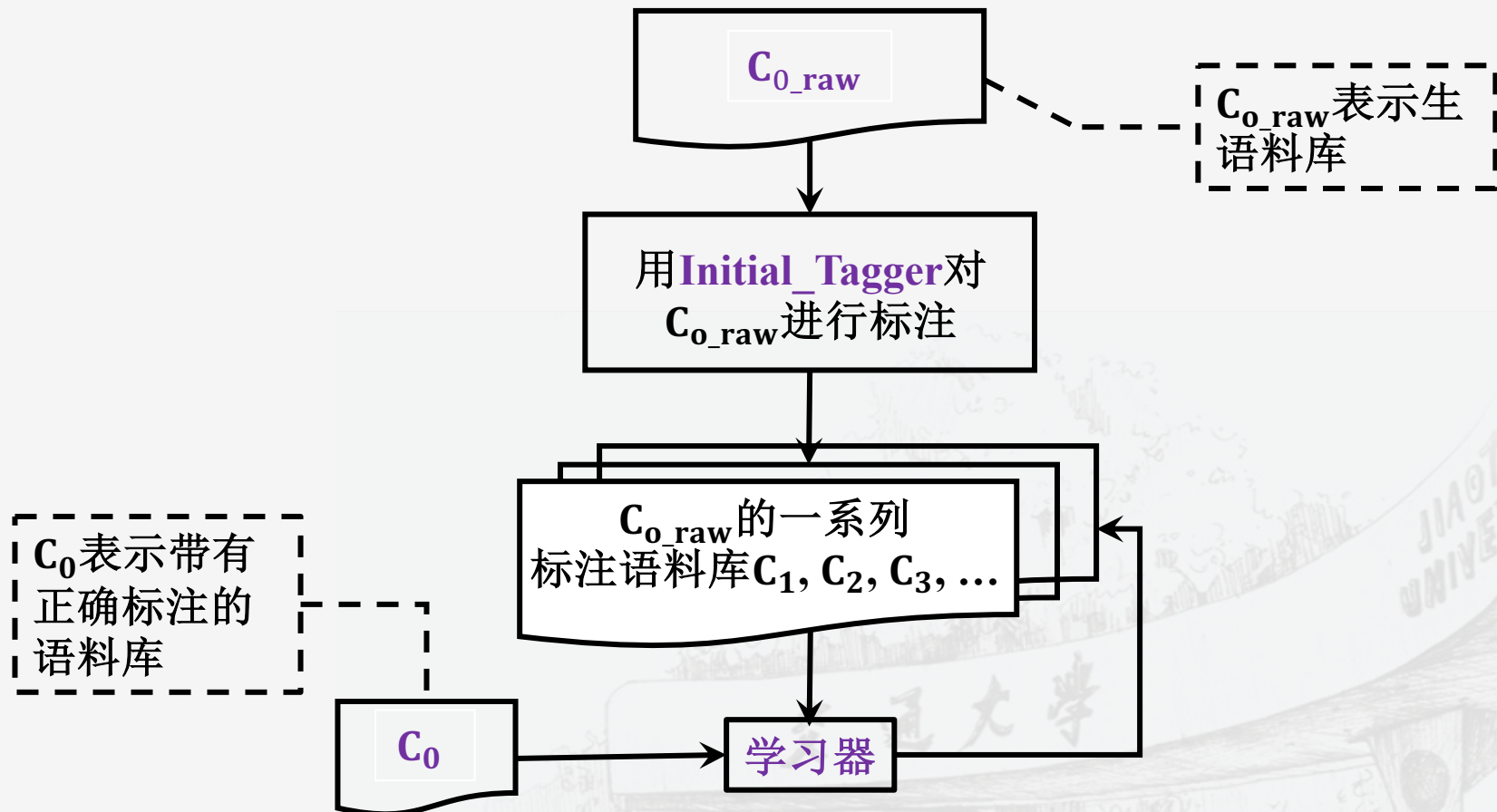
转换规则的学习流程



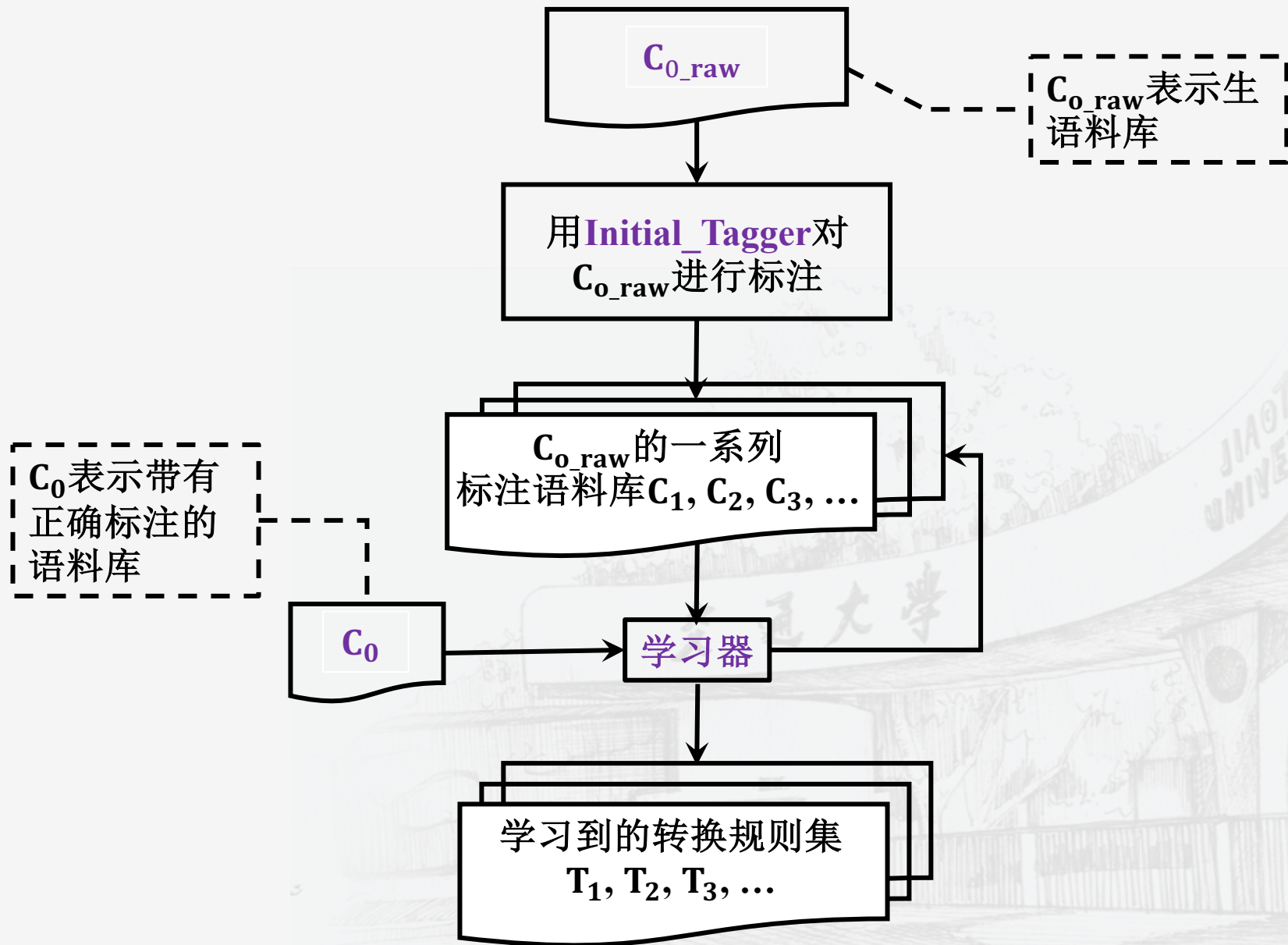
转换规则的学习流程



转换规则的学习流程

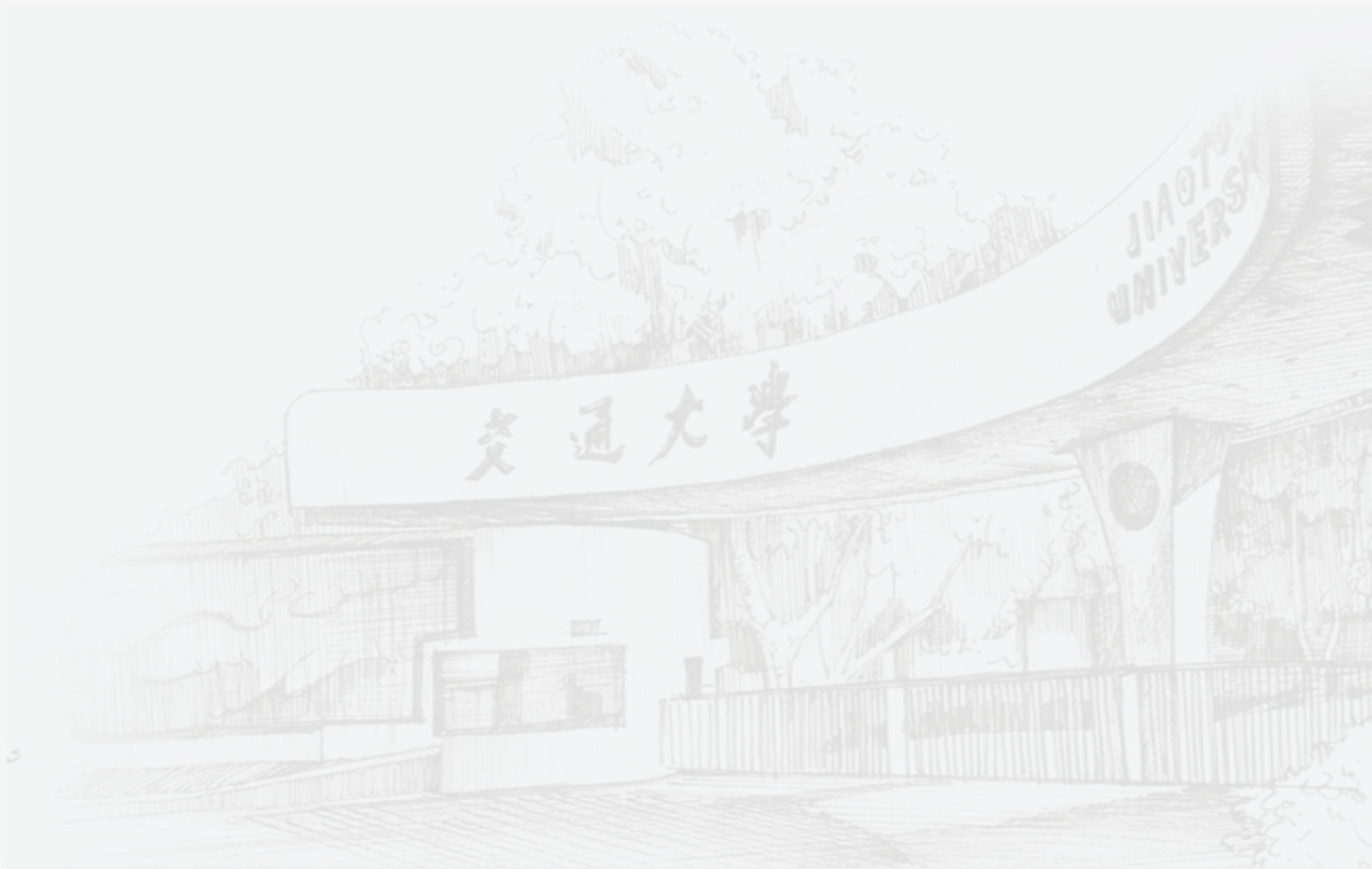


转换规则的学习流程



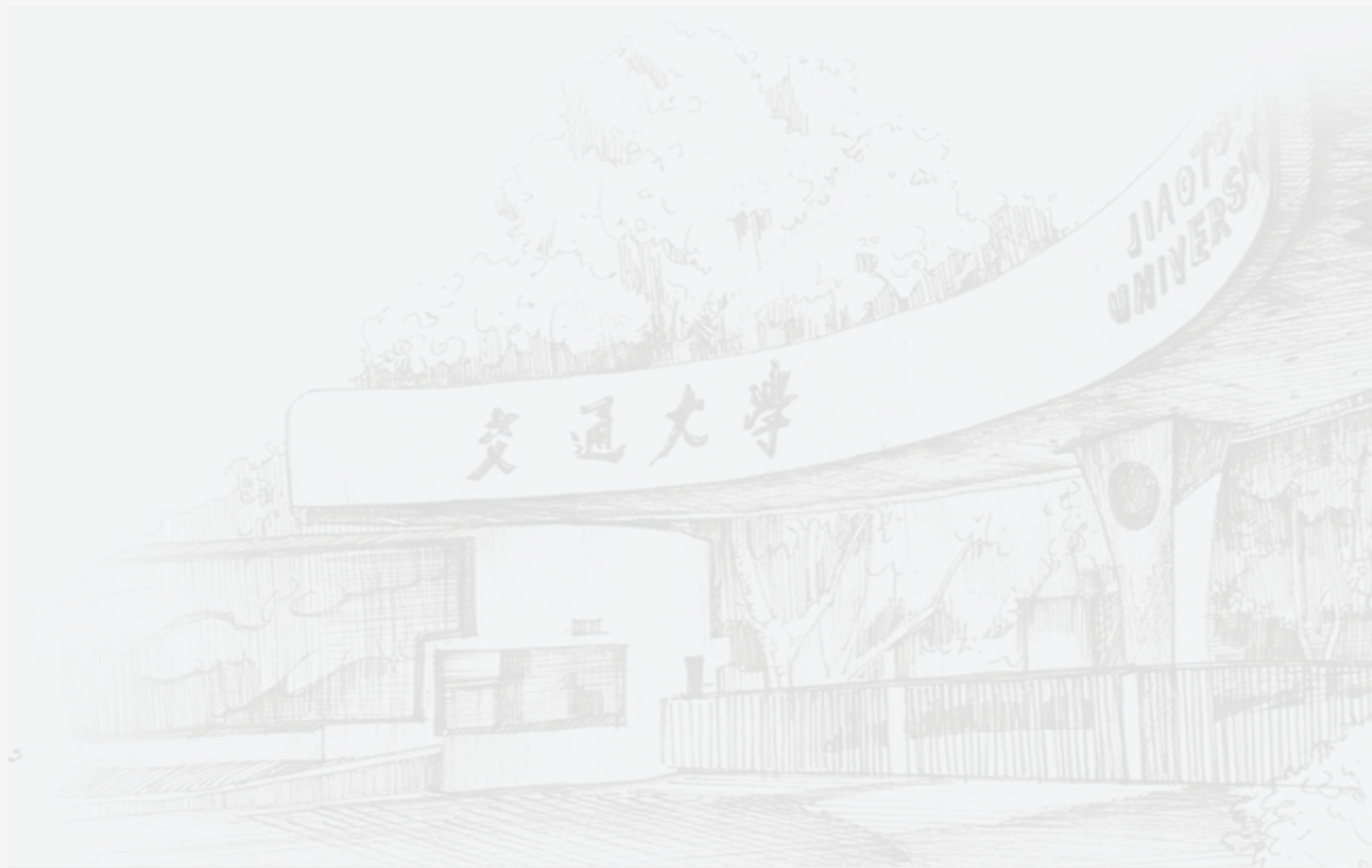
转换规则学习器算法

1. 首先用初始标注器对 C_{0_raw} 进行标注，得到带有词性标记的语料 C_i ($i=1$);



转换规则学习器算法

1. 首先用初始标注器对 C_{0_raw} 进行标注，得到带有词性标记的语料 C_i ($i=1$);
2. 将 C_i 跟正确的语料标注结果 C_0 比较，可以得到 C_i 中总的词性标注错误数;



转换规则学习器算法

1. 首先用初始标注器对 C_{0_raw} 进行标注，得到带有词性标记的语料 C_i ($i=1$);
2. 将 C_i 跟正确的语料标注结果 C_0 比较，可以得到 C_i 中总的词性标注错误数;
3. 依次从候选规则中取出一条规则 T_m ($m=1,2,\dots$)，每用一条规则对 C_i 中的词性标注结果进行一次修改，就会得到一个新版本的语料库，不妨记做 C_i^m ($m=1,2,3,\dots$)，将每个 C_i^m 跟 C_0 比较，可计算出每个 C_i^m 中的词性标注错误数。假定其中错误数最少的那个是 C_i^j (可预期 C_i^j 中的错误数一定少于 C_i 中的错数)，产生它的规则 T_j 就是这次学习得到的转换规则；此时 C_i^j 成为新的待修改语料库，即 $C_i^j = C_i^j$

转换规则学习器算法

1. 首先用初始标注器对 C_{0_raw} 进行标注，得到带有词性标记的语料 C_i ($i=1$);
2. 将 C_i 跟正确的语料标注结果 C_0 比较，可以得到 C_i 中总的词性标注错误数;
3. 依次从候选规则中取出一条规则 T_m ($m=1,2,\dots$)，每用一条规则对 C_i 中的词性标注结果进行一次修改，就会得到一个新版本的语料库，不妨记做 C_i^m ($m=1,2,3,\dots$)，将每个 C_i^m 跟 C_0 比较，可计算出每个 C_i^m 中的词性标注错误数。假定其中错误数最少的那个是 C_i^j (可预期 C_i^j 中的错误数一定少于 C_i 中的错数)，产生它的规则 T_j 就是这次学习得到的转换规则；此时 C_i^j 成为新的待修改语料库，即 $C_i^j = C_i^j$
4. 重复第3步的操作，得到一系列的标注语料库 $C_2^k, C_3^l, C_4^m, \dots$ 后一个语料库中的标注错误数都少于前一个中的错误数，每一次都学习到一条令错误数降低最多的转换规则。直至运用所有规则后，都不能降低错误数，学习过程结束。这时得到一个有序的转换规则集合 $\{T_a, T_b, T_c, \dots\}$



西安交通大学
XI'AN JIAOTONG UNIVERSITY

Q & A

