

1. 实验内容与完成程度

自动词性标注就是用计算机来自动地给文本中的词标注词类。任务是根据一个词在某个特定句子中的上下文，为这个词标注正确的词性，要解决的主要问题是词性兼类歧义和未登录词词性的确定问题。词性标注的正确率直接影响到文本的后续工作，因为词义消歧和句法分析都以经过词性标注的句子为基础。

本实验中我选择了维特比算法实现中文分词。

本程序较好的完成了维特比算法，但是对未登录词的处理不是很好。

2. 实验原理

2.1 HMM 模型

从概率角度上看，词性标注问题可描述为：给定一个含有词语序列 w_1, w_2, \dots, w_m 的输入句子 W ，确定最有可能的词性标记序列 $T=t_1, t_2, \dots, t_m$ ，使得条件概率 $P(T|W)$ 最大。

根据 Bayes 公式， $P(T|W) = P(T) P(W|T) / P(W)$ ，由于词性标注中， W 是给定的， $P(W)$ 不依赖于 T ，所以 $P(T|W) \approx P(T) P(W|T)$ ，其中 $P(T) = P(t_1|t_0) P(t_2|t_1, t_0) \dots P(t_i|t_{i-1}, t_{i-2}, \dots)$ ，根据一阶 HMM 假设可得 $P(T) \approx P(t_1|t_0) P(t_2|t_1) \dots P(t_i|t_{i-1})$ ， $P(W|T)$ 是已知词性标记串产生词串的条件概率， $P(W|T) = P(w_1|t_1) P(w_2|t_2, t_1, w_1) \dots P(w_i|t_i, t_{i-1}, \dots, t_1, w_{i-1}, \dots, w_1)$ ，根据一阶 HMM 假设可得 $P(W|T) \approx P(w_1|t_1) P(w_2|t_2) \dots P(w_i|t_i)$ 。

2.2 维特比算法

维特比算法是一种动态规划算法用于寻找最有可能产生观测事件序列的-维特比路径-隐含状态序列，特别是在马尔可夫信息源上下文和隐马尔可夫模型中。维特比算法的基础可以概括成下面三点：

1. 如果概率最大的路径 p (或者说最短路径) 经过某个点，比如途中的 X_{22} ，那么这条路径上的起始点 S 到 X_{22} 的这段子路径 Q ，一定是 S 到 X_{22} 之间的最短路径。否则，用 S 到 X_{22} 的最短路径 R 替代 Q ，便构成一条比 P 更短的路径，这显然是矛盾的。证明了满足最优性原理。

2. 从 S 到 E 的路径必定经过第 i 个时刻的某个状态，假定第 i 个时刻有 k 个状态，那么如果记录了从 S 到第 i 个状态的所有 k 个节点的最短路径，最终的最短路径必经过其中一条，这样，在任意时刻，只要考虑非常有限的最短路径即可。

3. 结合以上两点，假定当我们从状态 i 进入状态 $i+1$ 时，从 S 到状态 i 上各个节点的最短路径已经找到，并且记录在这些节点上，那么在计算从起点 S 到第 $i+1$ 状态的某个节点 X_{i+1} 的最短路径时，只要考虑从 S 到前一个状态 i 所有的 k 个节点的最短路径，以及从这个节点到 X_{i+1} 的距离即可。

2.3 维特比算法处理词性标注

中文中，每个词会有多种词性（比如“希望”即是名字又是动词），给出一个句子后，我们需要给这个句子的每个词确定一个唯一的词性，实际上也就是在若干词性组合中选择一个合适的组合。动词、名词等词类的搭配是具有规律性的，比如动词+名词的形式是大量存在的，当我们看到句子“存在希望”，如果确定了“存在”是动词，那么由于动名词组合的概率较大，我们会认定“希望”是名词。viterbi 算法的原理就是基于此。我们需要计算所有名词+动词，名词+名词，动词+形容词。。。。等种种词性搭配的出现概率，然后从中选出最大概率的组合。设一句

话有三个词组成分，分别为 A、B、C，则该句的概率值可以使用如下公式表示：

$$P(A \text{ 为 } sxA1) * P(sxA1 \text{ 位于句首}) * P(B \text{ 为 } sxB1) * P(sxA1 \rightarrow sxB1) * P(C \text{ 为 } sxC1) * P(sxB1 \rightarrow sxC1)$$

A、B、C 有多个属性，对每个属性都用上式计算，去概率最大值为该句子的次序标注。

3. 训练语料库与开放环境

实验选用人民日报 98 语料库；

编程语言为 C++；

编程环境为 vs2017(windows)

评测数据 自己拟订

注：

读取训练文件的语句 `fin.open("F://train_half.txt")`;

读取测试文件的语句 `fin1.open("F://test_file.txt")`;

4. 实验设计与实现

4.1 语料库的准备

本实验采用人民日报 98 年语料库

中国/ns 人民/n 银行/n 一/m 项/q 调查/vn 表明/v 千/m 户/q 大型/b 企业/n 实力/n 增强/v 运作/v 呈现/v 五/m 个/q 特点/n : /w 产销/vn 利润/n 同步/vd 增长/v 盈利/n 向/p 少数/m 企业/n 集中/v 亏损/vn 状况/n 有所/v 改善/vn 兼并/vn 重组/vn 步伐/n 加快/v 在建/v 规模/n 增长/v 最近/t 对/p 1 2 5 4/m 户/q 大型/b 企业/n 的/u 监测/vn 分析/vn 表明/v 本报/t 北京/ns 1 月/t 4 日/t 讯/Ng 记者/n 施/nr 明珠/nr 报道/v : /w 中国/ns 人民/n 银行/n 最近/t 对/p 1 2 5 4/m 户/q 大型/b 企业/n 的/u 监测/vn 分析/vn 表明/v 据/p 中国/ns 人民/n 银行/n 调查/v , /w 1 9 9 7 年/t 千/m 户/q 大型/b 企业/n 的/u 运作/vn 轨道/n 呈现/v 出/v 如下/vn 特点/n : /w 产销/vn 利润/n 同步/vd 增长/v 截至/v 1 9 9 7 年/t 1 1 月/t 末/t , /w 千/m 户/q 企业/n 实现/v 工业/n 总产值/n 9 5 8 8 . 2 亿/m 元/q (/w 不/d 变价/v) , /w 盈利/vn 向/p 少数/m 特大型/b 企业/n 及/c 企业/n 集团/n 集中/v 截至/v 1 9 9 7 年/t 1 1 月/t 末/t , /w 千/m 户/q 企业/n 销售/vn 收入/n 1 0 0 亿/m 元/q 一些/m 企业/n 亏损/vn 状况/n 有所/v 改善/v 千/m 户/q 大型/b 企业/n 的/u 亏损面/n 为/v 3 0 . 7 %/m , /w 是/v 近年来/l 的/u 最低/a 水平/n 兼并/vn 重组/vn 步伐/n 加快/v 1 9 9 7 年/t 千/m 户/q 大型/b 企业/n 运用/v 资本/n 营运/vn 工具/n , /w 实施/v 重组/v 、/w 兼并/v , /w 加强/v 联合/v 在建/v 规模/n 增长/vn 较/d 多/a , /w 千/m 户/q 大型/b 企业/n 中/t 有/v 相当/vn 一/m 批/q 基础/n 产业/n 和/c 支柱/n 产业/n 企业/n 承担/v 了/u 国家/n 重点/n 目前/t 企业/n 集团/n 发展/vn 值得/v 注意/v 的/u 几/m 个/q 问题/n 国家/n 计委/j 宏观/n 经济/n 研究院/n 实施/v 大/a 公司/n 、/w 大/a 集团/n 战略/n 是/v 国有/vn 大中型/b 企业/n 改革/v 和/c 发展/v 的/u 一/m 项/q 重大/a 选择/vn , /w 对于/p 国民经济/n 实现/v 两/m 个/c 一/m 、/w 企业/n 集团/n 发展/vn 食/v 大/a 团/v 快/a 自前/t , /w 一些/m 企业/n 集团/n 和/c 政府部门/n 明确/ad 提出/v 要/v 在/p 本世纪/t 末/t 下/t 一/m 世纪/n 初/t 使/v 一些/m 企业/n (/w 集团/n) /w 进入/v 企业/n 发展/v 、/w 扩张/v 可以/v 采取/v 两/m 种/q 途径/n : /w 一/m 是/v 内部/t 扩张/vn , /w 通过/p 资本/n 积累/vn , /w 凭借/p 自己/t 的/u 技术/n 优势/n 、/w 沿江/nr 泽民/nr 总书记/n 在/p 十五大/j 报告/n 中/t 提出/v 发展/v “/w 四跨/j ”/w 企业/n 集团/n 时/Ng , /w 同时/c 提到/p 了/u 一些/m 前提/n 条件/n , /w 就是/v 要/v 二/m 、/w 把/p 规模/n 经济/n 等同/v 于/p 经济/n 规模/n 一些/m 企业/n 集团/n 和/c 有关/vn 部门/n 的/u 同志/n 看到/v 我国/n 企业/n 集团/n 与/p 国外/s 大/a 企业/n 相比/v , /w 从/p 资产/n 、/w 销售额/n 等/u 方面/n 看/v 所谓/v 规模/n 经济/n , /w 是/v 指/v 在/p 技术/n 水平/n 不/d 变/v 时/Ng , /w n/tx 倍/q 的/u 投入/vn 产生/v 了/u 大于/v n/tx 倍/q 的/u 产出/vn , /w 这/t 也/d 但是/c , /w 规模/n 大/a 不/d 等于/v 规模/n 经济/n , /w 这/t 可以/v 从/p “三/m 个/q 方面/n 考察/v : /w (/w 1 m) /w 生产能力/l 的/u 限度/n , /w 投入/vn 增加/vn (/w A/nx 、/w B/nx) /w 过去/t 我国/n 曾/d 不止/v 一/m 次/q 出于/v 实现/v 规模/n 经济/n 和/c 专业化/vn 分工/vn 等/u 方面/n 的/u 考虑/vn , /w 组建/v 大/a 企业/n 、/w 大/a 集团/n , /w 三/m 、/w 过分/ad 追求/v 多元化/vn 经营/vn 目前/t 不少/m 企业/n 集团/n 为了/p 迅速/ad 扩张/v , /w 不仅/c 在/p 本/t 行业/n 大量/m 购并/v , /w 而且/c 进入/v 别的/t 行业/n 不少/m 企业/n 集团/n 提出/v : 应该/v 说/v , /w 多元化/vn 经营/vn 战略/n 是/v 大型/b 企业/n 集团/n 发展/v 的/u 重要/a 战略/n 选择/vn , /w 在/p 美国/ns , /w 特别/d 是/v 进入/v 本世纪/t 6 0 m 但是/c , /w 多元化/v 不/d 一定/d 会/v 减少/v 企业/n 的/u 经营/vn 风险/n , /w 表面/n 上/t 看/v , /w 多元化/vn 使/v 企业/n “/w 不/d 把/p 所有/b 的/u 鸡蛋/n 放在, 深圳/ns 赛格/nz 集团/n 近年/t 的/u 发展/vn 提供/v 了/u 一个/m 很/d 好/a 的/u 案例/n , /w 前/t 几/m 年/q , /w 赛格/nz 集团/n 大量/m 铺摊/v 设点/v 、/w 收购/v 兼 总之/c , /w 多元化/vn 经营/vn 不/d 一定/d 会/v 减弱/v 风险/n , /w 全面/ad 出击/v 可能/v 不/d 如/v 重点/d 出击/v , /w “/w 伤其十指不如断其一指/l ”/w “/w 把/p 主, 四/m 、/w 过分/ad 强调/v 低/a 成本/n 扩张/v 从/p 整体/n 上/t 搞活/v 国有/vn 经济/n , /w 实行/v 国有/vn 企业/n 的/u 战略性/n 改组/vn , /w 是/v 国有/vn 企业/n 改革/vn 的/u 重要/a 举措/n , /w 很多/m 企业/n : 应当/v 看到/v , /w 收购/vn 兼并/vn 是/v 一/m 种/q 风险/n 很/d 大/a 的/u 经营/vn 活动/vn , /w 国内/s 不乏/v 因/c 兼并/v 收购/v 后/t 救/v 不/d 活/a 别人/t , /w 最/d 重要/a 的/u 是/v , /w 企业/n 扩张/vn 成本/n 的/u 高低/n 必须/d 同/p 能否/v 实施/v 正确/a 的/u 发展/vn 战略/n 相/d 比较/v , /w 在/p 涉及/v 企业/n 长远/a 发 五/m 、/w 认为/v 资产/n 经营/vn 基于/v 生产/v 经营/vn 目前/t , /w 资产/n 经营/vn 概念/n 很/d 时髦/a , /w 相当/d 一部分/m 人士/n 认为/v 资产/n 经营/vn 是/v 一/m 种/q 高级/a 经营/vn 形式/n , /w 企业/n 要/v 从/p 生产, 资产/n 经营/vn 属于/v 投资/vn 银行/n 等/u 中介/n 机构/n 的/u 业务/n , /w 旨在/v 进行/v 有效/a 的/u 证券/n 组合/n (/w p o r t f o l i o / n x) /w , /w 规避/v 风险, 美国/ns 6 0 m 年代/n 资产/n 经营/vn 热/n 的/u 兴起/v 和/c 消退/v , /w 很/d 值得/v 我们/t 借鉴/v , /w 当时/t , /w 美国/ns 购并/v 浪潮/n 兴起/v , /w 企业/n 向/p 六/m 、/w 大/a 企业/n 应当/v 强化/v 自身/t 建设/v 企业/n 的/u 十五大/j 为/p 企业/n 集团/n 的/u 进一步/p 发展/vn 提供/v 了/u 新/a 的/u 机遇/n , /w 要/v 抓住/v 大/d 发展/v 的/u 机遇/n , /w 必须/d 全面/ad 把握/v 十 企业/n 集团/n 总部/n 还要/v 强化/v 功能/n 建设/vn , /w 建立/v 完整/a 的/u 功能/n 体系/n , /w 传统/n 体制/n 下/t 的/u 企业/n 基本上/d 功能/n 单一/a , /w 或者/c : 总之/c , /w 党/n 的/u 十五大/j 进一步/d 指明/v 了/u 企业/n 集团/n 改革/v 和/c 发展/v 的/u 方向/n , /w 提供/v 了/u 新/a 的/u 机遇/n , /w 企业/n 集团/n 目前/t 应

4.2 读取并记录数据

对语料库的内容进行统计。需要得到以下数据。

(1) 所有可能的词性。

(2) 所有出现的词语。

- (3) 每个词语以不同词性出现的次数。
- (4) 记录句首词为不同词性的次数。
- (5) 记录不同词性间转换的次数。(如遇到：“看电影”这个句子，则有[动词][名词]的值加一。)

4.3 获取概率结果

- (1) 计算每类词性作为句首出现的比例(比如:动词为句首, 占有所有不同词性为句首中的比例), 记录到 `double head_sx[sx_num];`。
- (2) 计算前词固定为词性[n]时, 后词为词性[x] 占总情况的比例(如: 前词固定为[动词]时, 后词[名词]出现的次数占有所有[x][动词]的比例), 记录到 `double bet[sx_num][sx_num];`
- (3) 计算每一个词作为不同类词性出现的次数, 占有该类词出现总数的比例(如: “震惊”作为动词出现的次数占有所有名词的比例), 记录到 `double sx[sx_num];`

4.4 标注

由以上得到了 `head_sx[],bet[],sx[]` 数组。先将测试的一条语句全部读入, 求得其词组数目 `len_test`。然后使用函数 `head_handle()`, 处理第一个词组 A。对 A 的所有属性 `sx`, 计算

`pro = P(sx)*P(sx 在句首)`, 调用函数 `body_handle(int last,double probability, int level)`。解释一下该函数的几个参数的含义: `last`->上一个词组的属性, `probability`->到上一个词组的概率积, `level`->目前的深度(`head_handle()`中 `level = 0`)。 `body_handle()`处理词组 B, 对 B 的所有属性 `sx`, 计算 `new_pro = pro* P(sx)*bet[A][B]`, 递归调用函数 `body_handle(int last,double probability, int level)`。在该过程中, 记录每一个 `level` 所取得的属性值。

当 `level>=len_test`, 即该句子处理完毕。比较此时的概率与最大概率的值, 若此时的概率大于最大概率, 则记录此时每一个 `level` 对应的属性值。

最后 `head_handle()`返回, 即得到了当前语句的最优词性标注。

继续处理下一句, 直至文本结束。

4.5 部分代码展示

读取训练文本函数 `read_head()`

```

bool vtbl::read_head()
{
    //c = ' ';
    while (c == ' ' || c == '\n' || c == '\r')
        c = fin.get();

    if (fin.eof())
        return false;

    char temp_value[50];
    int len_temp_value = 0;
    while (c != '/') //读取词组
    {
        if (c != ' ' && c != '\n' && c != '\r')
            temp_value[len_temp_value++] = c;

        c = fin.get();
    }
    temp_value[len_temp_value] = '\0';

    cout << temp_value << " ";

    /*****
    *处理head*
    *****/
    count_head++; //首词组总数加1
    int pos = word_pos(temp_value); //if(pos == len_w)->新词组, 将改词添加后还需要对w[len_w]做初始化

    /*新词组处理*/
    if (pos == len_w)
    {
        strcpy_s(w[len_w].value, temp_value);
        w[len_w].total_value = 1; //出现总数为1
        for (int i = 0; i < sx_num; i++)
        {
            w[len_w].count_sx[i] = 0;
            w[len_w].sx[i] = 0;
        }
        len_w++;
    }

    /*已有词处理*/
    else

```

读取训练文本函数 read_czsx()

```

bool vtbl::read_czsx()
{
    int sentence_end = 0;
    while (c == ' ' || c == '\n')
        c = fin.get();
    if (fin.eof())
        return false;
    char temp_value[50];
    int len_temp_value = 0;
    while (c != '/')
    {
        if (c == '\n')
            return false;

        temp_value[len_temp_value++] = c;
        c = fin.get();
    }
    temp_value[len_temp_value] = '\0';
    cout << temp_value << " ";

    /*****
    *句子终止处理*
    *****/
    if (strcmp(temp_value, ".") == 0 || strcmp(temp_value, "?") == 0 || strcmp(temp_value, "!") == 0 || strcmp(temp_value, ":") == 0 || strcmp(temp_value, ";") == 0)
        sentence_end = 1;

    /*****
    *处理w*
    *****/
    //count_head++; //首词组总数加1
    int pos = word_pos(temp_value); //if(pos == len_w)->新词组, 将改词添加后还需要对w[len_w]做初始化

    /*新词组处理*/
    if (pos == len_w)
    {
        strcpy_s(w[len_w].value, temp_value);
        w[pos].total_value = 1; //出现总数为1
        for (int i = 0; i < sx_num; i++)
        {
            w[pos].count_sx[i] = 0;
            w[pos].sx[i] = 0;
        }
    }

```

处理统计数据，获取相应概率

```
void vtbl1::train()
{
    int end_flag = read_head();
    //cout << "over0";
    while (end_flag)
    {
        read_rest();
        end_flag = read_head();
    }
    //cout << "over1";
    /*训练结果汇总*/

    //head_sx[]
    for (int i = 0; i < sx_num; i++)
    {
        head_sx[i] = count_head_sx[i] / count_head;
    }

    //w[]
    for (int i = 0; i < len_w; i++)
    {
        for (int j = 0; j < sx_num; j++)
        {
            w[i].sx[j] = w[i].count_sx[j] / w[i].total_value;
        }
    }

    //bet[][]
    for (int i = 0; i < sx_num; i++)
    {
        for (int j = 0; j < sx_num; j++)
        {
            bet[i][j] = count_bet[i][j] / total_bet;
        }
    }
}
```

head_handle()

```
void vtbl::head_handle()
{
    int pos = word_pos(test[0]);           //句首词组在w[]中的标号
    int level = 0;                         //当前属性数组长度(层数)记为0
    max_pro = 0;                           //最大属性置0

    int i;
    int exist = 0;
    int exist1 = 0;
    for (i = 0; i < sx_num; i++)
    {
        if (w[pos].sx[i] != 0 && head_sx[i] != 0)
        {
            exist = 1;
            pre_sx[level] = i;
            double pro = w[pos].sx[i] * head_sx[i];
            body_handle(i, pro, level+1);
        }
    }
    if (exist == 0)                       //不存在两个都满足
    {
        for (i = 0; i < sx_num; i++)
        {
            if (w[pos].sx[i] != 0)
            {
                exist1 = 1;
                pre_sx[level] = i;
                //cout << "i = " << i << endl;
                double pro = w[pos].sx[i];
                body_handle(i, pro, level+1);
            }
        }
        if (exist1 == 0)                  //未登陆词
        {
            pre_sx[level] = 16;           //名词
            double pro = 1.0;
            body_handle(16, pro, level + 1);
        }
    }
}
```

body_handle()

```

void vtbl::body_handle(int last,double pro, int level)
{
    //level++;

    if (level >= len_test)
    {
        //比较当前的pro与最大pro, 判断是否需要记录此时的属性数组
        if (pro > max_pro)
        {
            for (len_best_sx = 0; len_best_sx < level; len_best_sx++)
            {
                best_sx[len_best_sx] = pre_sx[len_best_sx];
            }
            max_pro = pro;
            //level--;
            return;
        }
        int pos = word_pos(test[level]);

        //cout << "go into 1.5" << endl;

        int i;
        int exist = 0; //不存在满足if (w[pos].sx[i] != 0 && bet[last][i] != 0)的i
        int exist1 = 0;
        for (i = 0; i < sx_num; i++)
        {
            if (w[pos].sx[i] != 0 && bet[last][i] != 0)
            {
                //cout << "go into 1.75" << endl;
                exist = 1;
                //last = i;
                pre_sx[level] = i;
                double new_pro = pro * w[pos].sx[i] * bet[last][i];
                //cout << "go into 2" << endl;
                body_handle(i,new_pro,level+1);
                //cout << "go into 3" << endl;
                //level--;
            }
        }
    }
}

```

输出标注结果

```

void vtbl::analyze()
{
    cout << endl << endl << endl << "-----测试结果-----" << endl << endl;
    int flag = read_a_sentence();

    //cout <<endl<< "analyze start" << endl;
    //cout << "flag = " << flag;
    while (flag == true)
    {
        head_handle();
        int i;
        for (i = 0; i < len_test; i++)
        {
            fout << test[i] << "/"<<int_to_char(best_sx[i]) << " ";
            cout << test[i] << "/"<<int_to_char(best_sx[i]) << " ";
        }
        cout << endl;
        flag = read_a_sentence();
    }
}

```


5. 实验结果展示

5.1 词性分类

```
1 a; 形容词
2 Ag; 形容词语素
3 an; 名形词 //具有名词功能的形容词
4 b; 区别词
5 c; 连词
6 d; 副词
7 Dg; 副词语素
8 e; 叹词
9 Eg; 叹词语素
10 f; 方位词
11 g; 语素
12 h; 前接成分
13 i; 成语
14 j; 缩略语
15 k; 后接成分
16 l; 习用语
17 m; 数词
18 n; 名词
19 Ng; 名词语素
20 ns; 名词-表处所
21 nt; 名词-表时间
22 nx; 名词-英文字母
23 nz; 名词-其他专有名词
24 o; 像声词
25 p; 介词
26 q; 量词
27 Qg; 量词语素
28 r; 代词
29 Rg; 代词语素
30 s; 处所词
31 t; 时间词
32 Tg; 时间词语素
33 u; 助词
34 Ug; 助词语素
35 v; 动词
36 Vg; 动词语素
37 vd; 副动词
38 w; 标点
39 x; 字
40 y; 语气词
41 Yg; 语气词语素
42 z; 状态词
```


5.2 训练过程（读取数据并统计）

产业异军突起，呈现出快速化、国际化的发展趋势。我国的体育产业自80年代初开始萌芽，因产业发展前景看好，市场潜力巨大，使其在起步发展阶段便焕发出勃勃生机。目前，我国的体育主体产业、相关产业及多种经营业均取得长足进展，体育产业有望成为新的消费热点和新的经济增长点。体育产业也因此成为社会各界关注的热点。今天的研讨会开得热烈、活跃，来自社会各界的20位演讲者从各自不同的角度，就如何开发体育市场的巨大潜力、如何发挥新闻媒介的优势推进体育产业发展、体育中介机构如何发挥作用、企业如何运作体育赞助活动和体育产业发展前景等问题，进行了深入、细致的探讨。大家在今天的研讨活动中还达成这样一个共识：体育主管部门、新闻传媒、商家企业和体育中介机构，是推动体育产业发展的四个轮子，四者缺一不可。同时，只有动员社会力量合力推动体育这一新兴产业，才能确保其快速、持续、稳定的发展。此次研讨活动为期两天，与会代表明天将分组进行座谈、研讨，1998年国内体育竞赛项目的发布、洽谈活动也将同时进行。八运失利痛定思痛安徽体育卧薪尝胆再图强本报合肥1月5日电记者刘杰报道：“体育就是形象，就是士气。八运失利，与人口大省不相符，与各项工作争先进的态势不相符。要卧薪尝胆，奋发图强”。4日，新年上班第一天，在安徽省政府第一会议室内，省长、副省长、省长助理、秘书长们，集体听取了《省体委关于八运会工作总结和迎接九运工作的报告》，大家众口一词，誓图体育强省。安徽省代表团在八运会上获5.5枚金牌，4.5枚银牌，10枚铜牌，总分369分，分别在金牌榜上和总分榜上列全国第二十位和第二十一位。尽管一些项目有了突破，一些新秀崭露头角，但位次较上届全运会有所后移，尤其是金牌位次有较大后退。“这么大一个省，体育上不去，成绩不突出，脸上无光啊！”回忆全运会上的情景，代表团的头头脑脑们仍觉心寒。“气可鼓而不可泄”。曾任省体委主任的省委常委、常务副省长汪洋提高了嗓门：“要以市场的眼光办体育，抓选才，抓教练，领导要与运动员、教练员滚在一块，艰苦奋斗，进军九运，再创辉煌。”省长回良玉说，立足省情，突出重点，形成争分夺牌的尖子群，下功夫抓好教练员队伍建设，抓基层，打基础，扩大业余训练规模，提高业余训练质量，这些都需要横下一条心，把标杆定在体育强省上，奋力拼搏才行。变与不变间——男篮联赛第一循环印象（附图片1张）本报记者薛原要说变化，这个赛季的希尔顿全国男篮甲A联赛不可谓不多。有8支球队替换了主教练，4支球队更改了主场，在内外援的引进上，各队也做足了功夫。但联赛打过1997年，打完了第一循环，从积分榜上看，同上个赛季这一阶段相比，变化甚少。只有广东队掉出前六名，北京队跻身前六名，两队互换位置。除此之外，上个赛季忙着保级的球队看来还得忙这事，而八一队再度卫冕的前景也是一片光明。第一循环中最令人兴奋的一场球要属江苏南钢队主场扳倒八一队，江苏队内有查理·曼德、詹姆斯大力扣篮，拼抢篮板，外有胡卫东、苑志南、李青山等人三分远投，个人突破，水平发挥淋漓尽致。不过，江苏队除了遇强不弱的优点外，还有遇弱不强的痼疾。这个赛季能赢八一队、辽宁队，却输给济南军区队和四川蓝剑队，表现忽起忽落，捉摸不定，在积分榜上也总是徘徊于四名至六名间。遇强不弱的并不止江苏队，浙江中欣队、空军联航队、四川蓝剑队都有这个劲头，联赛也因此爆出不少冷门。但碰上实力相当或稍差于自己的球队时，又往往“不会打球”，这一现象不是出现在同强队的比赛中，确实耐人寻味。

5.3 测试结果展示

状态转移矩阵（#表示该状态不存在）

```
29 a->s 1.59406e-05
30 a->t 3.18811e-05
31 a->Tg 1.59406e-05
32 a->u 0.00787464
33 a->Ug 0
34 a->v 0.000812969
35 a->Vg 0
36 a->w 0.00589801
37 a->x 0
38 a->y 0.000207227
39 a->Yg 0
40 a->z 0
41 a->ad 0
42 a->vn 0.00216792
43 a->vd 0
44 a->nx 0
45 a->an 0.000127525
46 a->nz 0
47 a-># 0
48 a-># 0
49 a-># 0
50 a-># 0
51 a-># 0
52 Ag->a 1.59406e-05
53 Ag->Ag 0
54 Ag->b 0
55 Ag->c 1.59406e-05
56 Ag->d 0
57 Ag->Dg 0
58 Ag->e 0
59 Ag->Eg 0
60 Ag->f 1.59406e-05
61 Ag->g 0
62 Ag->h 0
63 Ag->i 0
64 Ag->j 0
65 Ag->k 0
66 Ag->l 0
67 Ag->m 0
68 Ag->n 7.97029e-05
69 Ag->Ng 3.18811e-05
70 Ag->nr 0
71 Ag->ns 0
```

测试文件

6. 遇到的问题以及解决方案

1、在处理语料库文件，要注意编码方式，避免中文信息容易出现的乱码情况。Windows 环境下的换行是/n/r，因此在判断换行情况的时候条件应该是 if(==n/r)。在获取训练集的

时候要在每个句子的结尾加上/n。

2、if (w[pos].sx[i] != 0 && head_sx[i] != 0)，A 为句首词，但 A 的所有属性都不合放句首的情况。导致无法继续调用 body_handle()。if (w[pos].sx[i] != 0 && bet[last][i] != 0)，没有考虑不存在这种条件的情况，导致递归无法继续，最后标注结果全部为/a。

解决方法是增加了对这两种情况的处理：

```
if (exist == 0) //
{
    //cout << "go into 3.75" << endl;
    for (i = 0; i < sx_num; i++)
    {
        if (w[pos].sx[i] != 0)
        {
            //last = i;
            //cout << "i ==" << i<<endl;
            exist1 = 1;
            pre_sx[level] = i;
            double new_pro = pro * w[pos].sx[i];
            //cout << "go into 4" << endl;
            body_handle(i, new_pro, level + 1);
            //cout << "go into 5" << endl;
            //level--;
        }
    }
}
```

3、未登陆词词性为 a（较少出现）导致后序分析出错。解决方法：
未登陆词的属性随机选取出现频率最高的三种属性之一。

7. 心得体会

通过本次词性标注实验，对模型训练有了初步的了解。也知道了怎么处理得到训练数据，体会到了中文信息处理的乐趣。