

AI深度学习之自然语言处理顶级实战课程

四、依存句法和语义依存分析

讲师：aopu

自我介绍

- 天善商业智能和大数据社区[aopu](#) 讲师
- 天善社区ID- [aopu](#)主页
- <https://www.hellobi.com> – 学习过程中有任何相关的问题都可以提到技术社区 人工智能 版块。

4、4自定义语法与CFG

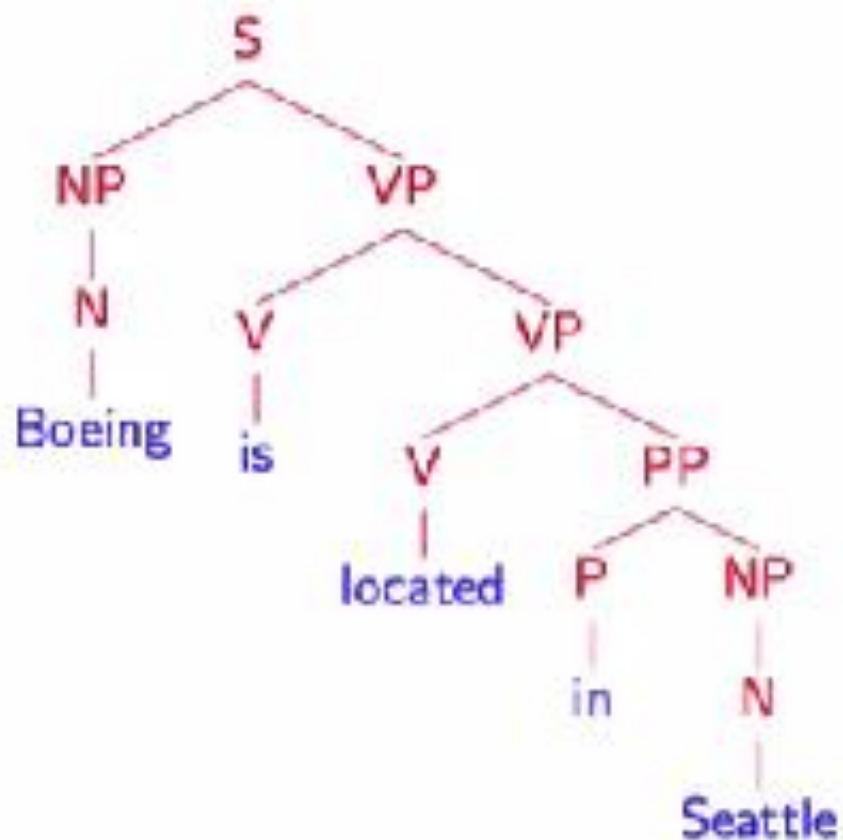
什么是语法解析？

- 在自然语言学习过程中，每个人一定都学过语法，例如句子可以用主语、谓语、宾语来表示。在自然语言的处理过程中，有许多应用场景都需要考虑句子的语法，因此研究语法解析变得非常重要。
- 语法解析有两个主要的问题，其一是句子语法在计算机中的表达与存储方法，以及语料数据集；其二是语法解析的算法。
- 对于第一个问题，我们可以用树状结构图来表示，如下图所示，S表示句子；NP、VP、PP是名词、动词、介词短语（短语级别）；N、V、P分别是名词、动词、介词。

4、4自定义语法与CFG

什么是语法解析？

Boeing is located in Seattle.



4、4自定义语法与CFG

- **上下文无关语法 (Context-Free Grammer)**
- 为了生成句子的语法树，我们可以定义如下的一套上下文无关语法。
- 1) N 表示一组非叶子节点的标注，例如{S、NP、VP、N...}
- 2) Σ 表示一组叶子结点的标注，例如{boeing、is...}
- 3) R 表示一组规则，每条规则可以表示为

$$X \rightarrow Y_1, Y_2, \dots, Y_n, X \in N, Y_i \in (N \cup \Sigma)$$

- 4) S 表示语法树开始的标注
- 举例来说，语法的一个语法子集可以表示为下图所示。当给定一个句子时，我们便可以按照从左到右的顺序来解析语法。例如，句子the man sleeps就可以表示为(S (NP (DT the) (NN man)) (VP sleeps))。

4、4自定义语法与CFG

• 上下文无关语法 (Context-Free Grammer)

$N = \{S, NP, VP, PP, DT, Vi, Vt, NN, IN\}$

$S = S$

$\Sigma = \{\text{sleeps, saw, man, woman, telescope, the, with, in}\}$

$R =$

S	→	NP	VP
VP	→	Vi	
VP	→	Vt	NP
VP	→	VP	PP
NP	→	DT	NN
NP	→	NP	PP
PP	→	IN	NP

Vi	→	sleeps
Vt	→	saw
NN	→	man
NN	→	woman
NN	→	telescope
DT	→	the
IN	→	with
IN	→	in

4、4自定义语法与CFG

- **概率分布的上下文无关语法 (Probabilistic Context-Free Grammar)**
- 上下文无关的语法可以很容易的推导出一个句子的语法结构，但是缺点是推导出的结构可能存在二义性。

由于语法的解析存在二义性，我们就需要找到一种方法从多种可能的语法树种找出最可能的一棵树。一种常见的方法既是PCFG (Probabilistic Context-Free Grammar)。如下图所示，除了常规的语法规则以外，我们还对每一条规则赋予了一个概率。对于每一棵生成的语法树，我们将其中所以规则的概率的乘积作为语法树的出现概率。

4、4自定义语法与CFG

- 概率分布的上下文无关语法 (Probabilistic Context-Free Grammar)

S	\Rightarrow	NP	VP	1.0
VP	\Rightarrow	Vi		0.4
VP	\Rightarrow	Vt	NP	0.4
VP	\Rightarrow	VP	PP	0.2
NP	\Rightarrow	DT	NN	0.3
NP	\Rightarrow	NP	PP	0.7
PP	\Rightarrow	P	NP	1.0

Vi	\Rightarrow	sleeps	1.0
Vt	\Rightarrow	saw	1.0
NN	\Rightarrow	man	0.7
NN	\Rightarrow	woman	0.2
NN	\Rightarrow	telescope	0.1
DT	\Rightarrow	the	1.0
IN	\Rightarrow	with	0.5
IN	\Rightarrow	in	0.5

- 当我们或得多颗语法树时，我们可以分别计算每颗语法树的概率 $p(t)$ ，出现概率最大的那颗语法树就是我们希望得到的结果，即 $\arg \max p(t)$ 。

4、4自定义语法与CFG

• 训练算法

- 我们已经定义了语法解析的算法，而这个算法依赖于CFG中对于 N 、 Σ 、 R 、 S 的定义以及PCFG中的 $p(x)$ 。上文中我们提到了Penn Treebank通过手工的方法已经提供了一个非常大的语料数据集，我们的任务就是从语料库中训练出PCFG所需要的参数。
- 1) 统计出语料库中所有的 N 与 Σ ；
- 2) 利用语料库中的所有规则作为 R ；
- 3) 针对每个规则 $A \rightarrow B$ ，从语料库中估算 $p(x) = p(A \rightarrow B) / p(A)$ ；
- 在CFG的定义的基础上，我们重新定义一种叫Chomsky的语法格式。这种格式要求每条规则只能是 $X \rightarrow Y_1 Y_2$ 或者 $X \rightarrow Y$ 的格式。实际上Chomsky语法格式保证生产的语法树总是二叉树的格式，同时任意一棵语法树总是能够转化成Chomsky语法格式。

4、4自定义语法与CFG

- 语法树预测算法

- 假设我们已经有一个PCFG的模型，包含 N 、 Σ 、 R 、 S 、 $p(x)$ 等参数，并且语法树总数Chomsky语法格式。当输入一个句子 x_1, x_2, \dots, x_n 时，我们要如何计算句子对应的语法树呢？
- 第一种方法是暴力遍历的方法，每个单词 x 可能有 $m = \text{len}(N)$ 种取值，句子长度是 n ，每种情况至少存在 n 个规则，所以在时间复杂度 $O(m^n n)$ 的情况下，我们可以判断出所有可能的语法树并计算出最佳的那个。

4、4自定义语法与CFG

- 语法树预测算法
- 第二种方法当然是动态规划，我们定义 $w[i, j, X]$ 是第 i 个单词至第 j 个单词由标注 X 来表示的最大概率。直观来讲，例如 x_i, x_{i+1}, \dots, x_j ，当 $X=PP$ 时，子树可能是多种解释方式，如 $(P\ NP)$ 或者 $(PP\ PP)$ ，但是 $w[i, j, PP]$ 代表的是继续往上一层递归时，我们只选择当前概率最大的组合方式。特殊情况下。因此，动态规划的方程可以表示为
- 关于动态规划方法，leetcode里有不少案例可以说明。
- 语法解析按照上述的算法过程便完成了。虽说PCFG也有一些缺点，例如：1）缺乏词法信息；2）连续短语（如名词、介词）的处理等。但总体来讲它给语法解析提供了一种非常有效的实现方法。

秦路主讲

七周成为数据分析师

七周为期，Get一条数据分析师职业黄金通道！



Python

数据分析与挖掘

集Python爬虫、数据采集、数据处理、数据分析与数据挖掘于一体，打造Python全栈工程师

主讲老师: 韦玮

VIP会员群+在线答疑+录播复习+1年反复观看

案例为师,实战为王

开启Python机器学习之路

科学规划全套课程体系，从入门到进阶，从理论到技巧，嵌入丰富课程案例讲解，逐步推进

讲师: 唐宇迪 深度学习领域多年一线实践研究专家

独一无二的数据库建模指南系列教程升级版

- 从企业视角进行数据规划以及数据库模型的搭建
- 高质量的数据库模型和技巧，以及丰富的例子
- 数据库架构理论和实践要领

资深讲师: BAO胖子 15年+BI从业经验
涉足电力、快消品、医药、信息服务行业的BI老兵

业务知识一站通

技术+业务，挣钱有门路！

讲师: 陈文



自己动手 丰衣足食

Python3网络爬虫实战案例

一循序渐进，案例为王，诠释全面，思路制胜一

讲师: 崔庆才 北航硕士，百万级热度爬文博主



讲师 丘祐玮

人人都爱数据科学家

Python数据科学精华实战课程

数据分析报告制作

秘籍升级版

讲师: 陈丹奕 知乎大神，前百度资深数据分析师

先机致胜 破冰AI

深度学习模型/框架与实战

讲师: 唐宇迪 同济大学硕士
深度学习领域多年一线实践研究专家



BI、商业智能
数据挖掘 大数据
数据分析师
R语言 Python
机器学习
深度学习
人工智能
Hive Hadoop
Tableau
BIEE ETL
数据科学家
PowerBI