

基于查找树的 IP 地址分类算法研究

王晓勇¹ 邱玉辉²

(西南大学信息中心 重庆 400715)¹(西南大学计算机与信息科学学院 重庆 400715)²

摘要 随着 Internet 的大规模发展,越来越多的网络业务需要对 IP 地址进行适时、快速分类。在分析二叉 Trie 树的基础上,改进了其结构,提出了基于 256-叉查找树的 IP 地址分类算法,并详细介绍了其实现过程,比较了它们的优缺点。该算法在满足空间要求的情况下,提高了查找分类时间,具有通用性和实用价值。

关键词 Trie 树, IP 地址, 分类, 查找树

IP Classification Algorithm Based on Search-tree

WANG Xiao-Yong¹ QIU Yu-Hui²

(Information Center, Southwest University, Chongqing 400715)¹

(Faculty of Computer & Information Science, Southwest University, Chongqing 400715)²

Abstract With the development of Internet, real-time and fast IP address classifications have been applied to more and more services. After analyzing the binary-Trie tree, an algorithm of IP address classification bases on 256-branch search tree is introduced. The analysis shows that this algorithm is prior in time and is acceptable in space.

Keywords Trie tree, IP address, Classification, Search tree

1 引言

随着互联网的大规模发展,越来越多的业务需要对 IP 地址进行适时、快速的分类,比如网络计费系统、网络流量监控、多出口的策略路由等等^[1]。IP 地址分类问题描述为:当 Internet 的 IP 数据包到达计费网关时,计费网关要决定是否允许该 IP 包通过,并完成计费。计费网关主要根据 IP 数据包头的目的地址和源地址,在内存的所有类别的 IP 地址列表中查找,根据查找结果决定丢弃或通过 IP 包,如允许通过,则根据地址的类别按不同的费率完成计费。内存中的 IP 地址列表由各种类别的 IP 地址块组成,如局域网内部地址、免费数据库地址、国内地址等。互联网一般采用 CIDR^[2](Classless Inter-Domain Routing)表示一个 IP 地址块,CIDR 可把多个 IP 地址块聚合成一个更大的地址块,减少地址列表的长度。地址块 192.168.8.0/24 的网络前缀长度为 24 位,表示 IP 地址范围:192.168.8.0—192.168.8.255。

因此内存中的 IP 地址列表是一个个的长度不等的网络地址前缀,而不是单个的 IP 地址。IP 地址分类的过程就是查找数据包的 IP 地址究竟属于哪块 IP 地址,查找采用最长网络前缀匹配。比如 192.168.8.1 即在地址块 192.168.0.0/16 范围内,也在地址块 192.168.8.0/24 范围内,但后者网络前缀 24 更长,则 192.168.8.1 只属于后者,而不属于前者。

2 Trie 树

Trie 树是一种用于快速检索的多叉树结构。Trie 树把要查找的关键词看作一个字符序列,根据这一序列构造用于检索的树结构。在 Trie 树上进行检索类似于查阅英语词典。

定义 1 $S = \{s_1, s_2, \dots, s_n\}$ 是定义在字符集 Σ 上的字符

串集合。S 的一个 Trie 树是一棵 k 分支树 T (其中 $k = |\Sigma|$),除根以外的所有非叶子节点(或每条边)表示 Σ 的一个字符,每个叶子节点 l_i 对应一个字符串 s_i ,而且从根到叶子节点 l_i 的路径上的所有节点(不包括叶子节点)表示的字符连接起来就是字符串 s_i 。Trie 树 T 具有以下特点:

(1) T 最多有 n 个叶子节点;(2) T 的深度为 S 中最长字符串的长度;(3) T 的节点个数为 $O(n \times m)$, n 为 S 中字符串的个数, m 为 S 中最长字符串的长度。

根据 Trie 树的以上特点同样可得其优缺点。优点:在 Trie 树中查找关键字的时间与树中包含的节点树无关,而只与组成关键字的字符数有关。如果所有关键字的字符数较少而关键字的数量较大时,采用 Trie 树具有明显速度上的优势。缺点:存储空间浪费较大。

3 基于 Trie 树的 IP 地址分类方法

3.1 二叉 Trie 树算法

有许多研究人员提出采用二叉 Trie 数来加速 IP 查找^[3,4]。一个 IP 地址可看作是 32 个 0 和 1 组成的字符串,即 $\Sigma = \{0, 1\}$,IP 地址块的网络前缀是由不大于 32 个 0 和 1 字符组成的字符串。表 1 是一个 IP 地址列表,其对应的二叉 Trie 树如图 1。

在分类一个 IP 地址时,从二叉 Trie 树的根开始查找,从 IP 地址的第一位开始依次往后,IP 地址的该位是 0 则指针移到左分支节点,是 1 则指针移到右分支节点。当下一指针是空指针时,则根据当前节点的标志即可判断该 IP 地址的类别。

二叉 Trie 树查找结构简单,但在最坏的情况下,对 IPv4 来说,查找比较次数多达 32 次。二叉 Trie 树的 IP 地址分类

*)本文受到教育部重点课题资助(No. 104262)。王晓勇 工程师,主要研究方向:计算机应用技术、网络流量监控。邱玉辉 博士研究生导师,主要研究方向:模糊逻辑、多 Agent 系统、网格计算等。

