

1. 对内容安全的需求主要体现在哪几个方面?

- ① 数字版权侵权及其控制: 数字内容产业主要指影视和音乐的数字化制作和发行行业。人们逐渐发现, 对数字版权的侵权仅仅依靠非技术手段是不够的。数字内容制作企业、内容创作者及管理部门也迫切需要有遏制版权侵权的技术手段。
- ② 不良内容传播及其控制: 不良内容的肆意传播是另外一个与内容相关的安全问题。
- ③ 敏感内容泄露及其控制: 敏感信息主要包括保密文件和与知识产权相关的资料等。
- ④ 内容伪造及其控制: 人们需要能够核实数字内容的真伪。

2. 不良文本过滤主要有哪些方法? 各自的工作原理是什么?

- ① 基于关键字的过滤方法: 首先由专业人员编制一个不良文本关键词库, 关键字词库中出现的字词都是经常不良文本中的敏感词汇, 能够很大程度地代表不良文本。当有文本到来之后, 对文本全文进行检索, 通过比较简单的布尔逻辑运算进行匹配。当匹配超过一定阈值之后, 系统就认为这篇文本是不良文本, 给予过滤。
- ② 基于分级标签的过滤方法: 分级标签过滤方法通过对不同的网页根据内容赋予不同级别, 以实现过滤。根据网页内容的不同, 分为普通级、一般限制级、严格限制级。
- ③ 基于地址库的过滤方法: 基于地址库过滤方法可分为以下三个类别: IP 地址过滤、URL 过滤以及 IP 和 URL 相结合的过滤方法。IP 过滤是指通过封锁指定网站的 IP 地址, 以实现包含有大量不良文本网站的过滤目的。URL 过滤方法直接定位不良文本在互联网上的具体位置, 直接对该网站进行屏蔽。IP 和 URL 的结合对于一个网站下大部分的网页都是不良信息的情况, 则采用 IP 过滤。如果是一个网站下只有极少一部分是不良文本, 则采用 URL 过滤。
- ④ 基于内容的动态过滤法: 在使用关键字过滤方法进行过滤时, 会有很高的误判率, 就需要进一步根据文本的具体内容来判断指定文本的具体类别属性。

3. 话题发现主要包括哪几个步骤?

- ① 信息采集: 利用网络爬虫工具, 从指定的网站把 Web 网页等互联网信息资源抓取到计算机本机存储。
- ② 网络文本预处理: 即把互联网网页源码信息进行处理, 包含剔除无关字符清洗源码, 提取正文和必要的附带信息。
- ③ 文本内容分词: 是在汉语文本处理, 且选择词语作为文档特征表达的特定情况下必要的步骤。
- ④ 文本向量化: 是汇总分词后文本中的词语, 将这些词语作为空间向量的维度构建文本表示的多维向量, 然后将各词的 TFIDF 公式转换为一个权重值, 用以表示文本在这个词语代表的维度的值, 进而将文本表示为一组关键词及其词频为权重的空间向量。
- ⑤ 网络文本聚类: 采取一定的组织策略调度文本向量参与相似度计算, 并建立话题的向量表示方式。
- ⑥ 话题热度评估: 综合考虑话题中所有报道的点击数、回帖数、报道频率等参数, 来评估该话题受关注的程度。

4. 多媒体内容安全面临的主要问题有哪些?

- ① 如何鉴别一个数字媒体作品的创建者?
- ② 如何确定数字媒体作品创建者的版权声明?
- ③ 如何公证一个数字作品的签名与版权声明?
- ④ 如何控制用户访问数字媒体作品的权限?