

# 统计软件 HW9

邵智轩

1400012141

物理学院

```
HLM.url <- "http://www.xiexingcun.net/honglouloumeng/"
index.url <- paste(HLM.url, "index.html", sep = "")
chapters.url <- paste(HLM.url, c(formatC(1:9, width = 2, flag = 0), 10:120),
  ".htm", sep = "")

outfcon <- file("honglouloumeng.txt", open = "wt") # 输出的 txt 文件

chchar.substr <- function(c) {
  # 找出汉字子串
  find.Chinese <- gregexpr("[^\\x00-\\xff]|[,!?:;]", c, perl = TRUE)[[1]]
  substr(c, start = find.Chinese[1], stop = find.Chinese[length(find.Chinese)])
}

subpunc <- function(c) {
  # 替换标点符号
  c <- gsub("&quot;", "\"", c, fixed = T)
  c <- gsub("[ ;]?&lt;?*[ ]?", "<", c, perl = T)
  c <- gsub("[ ;]?&gt;?*[ ]?", ">", c, perl = T)
  c
}

# 利用 index 页制作 txt 的 Contents 页
makeContents <- function(index.url = index.url) {
```

```

index.text <- readLines(url(index.url))
Contents <- grep("\\d+[.]htm", index.text, perl = TRUE, value = TRUE)
Contents <- vapply(Contents, chchar.substr, "s")
Contents <- c(" 红楼梦", " 清 曹雪芹 著", Contents)
Contents
}

# 读取单章
addChapter <- function(chapter.url) {
  chapter <- readLines(url(chapter.url))
  # 找出含汉字的行
  chchar <- grep("[^\\x00-\\xff]", chapter, perl = TRUE, value = TRUE)
  # 去掉汉字行中不需要的行
  chchar <- chchar[-grep("<title>|html| 作者: 曹雪芹", chchar, perl = TRUE)]
  # chchar<-vapply(chchar,chchar.substr,'s') 先处理标题第几回和题目
  chapnum <- chchar.substr(chchar[1])
  chapname <- chchar.substr(chchar[2])

  # 删掉回目，下面处理正文
  chchar <- chchar[-(1:2)]

  # 没有<td> 或<p> 或<br> 的段落连接到上一行
  newlines.index <- grep("      ", chchar, perl = TRUE) # 重启一行的
  chchar <- vapply(chchar, chchar.substr, "s")

  chnew <- c()
  for (i in (1:length(chchar))) {
    if (i %in% newlines.index) {
      chnew <- c(chnew, chchar[i])
    } else {
      chnew[length(chnew)] <- paste(chnew[length(chnew)], chchar[i], sep = "")
    }
  }
}

```

```

# 替换掉标点符号
chnew <- subpunc(chnew)

# 加上回目，以及回目和正文间的空行
c("", chapnum, chapname, "", chnew)
}

time.start <- proc.time()

writeLines(makeContents(index.url), con = outfcon)
for (chapter.url in chapters.url) {
  writeLines(addChapter(chapter.url), con = outfcon)
}

proc.time() - time.start

##      user  system elapsed
##    9.53    2.78    18.07

close(outfcon)

```

效果图：

