

```
In [19]: import pandas as pd
from wordcloud import WordCloud
import matplotlib.pyplot as plt
import nltk
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize

df1 = pd.read_csv('考研weibo.csv', encoding='utf-8')
# df1 = pd.read_csv('计算机专业考研weibo.csv', encoding='utf-8')
df1 = df1.fillna('') # 将NaN值替换成空字符串

df2 = pd.read_csv('考研zhihu.csv', encoding='utf-8')
# df2 = pd.read_csv('计算机专业考研zhihu.csv', encoding='utf-8')
df2 = df2.fillna('') # 将NaN值替换成空字符串
#
text = ' '.join(df1['content'].tolist() + df1['topic'].tolist() + df2['title'].tolist())
# wc = WordCloud(font_path = r'./MSYH.TTC')
# wc.generate(text)

# plt.imshow(wc)
# plt.axis("off") # 不显示坐标轴
# plt.show()

# 分词并去除停用词
nltk.download('stopwords')
nltk.download('punkt')
stop_words = set(stopwords.words('chinese'))
tokens = word_tokenize(text.lower())
keywords = [word for word in tokens if word.isalpha() and word not in stop_words]

# 生成词云
wordcloud = WordCloud(font_path = r'MSYH.TTC', width=800, height=800, background_color='

# 显示词云
plt.figure(figsize=(8, 8), facecolor=None)
plt.imshow(wordcloud)
plt.axis("off")
plt.tight_layout(pad=0)
plt.show()
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...
```

```
[nltk_data] Package stopwords is already up-to-date!
```

```
[nltk_data] Downloading package punkt to /root/nltk_data...
```

```
[nltk_data] Package punkt is already up-to-date!
```



```
In [20]: import pandas as pd
from wordcloud import WordCloud
import matplotlib.pyplot as plt
import nltk
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize

# df1 = pd.read_csv('考研weibo.csv', encoding='utf-8')
df1 = pd.read_csv('计算机专业考研weibo.csv', encoding='utf-8')
df1 = df1.fillna('') # 将NaN值替换成空字符串

# df2 = pd.read_csv('考研zhihu.csv', encoding='utf-8')
df2 = pd.read_csv('计算机专业考研zhihu.csv', encoding='utf-8')
df2 = df2.fillna('') # 将NaN值替换成空字符串
#
text = ' '.join(df1['content'].tolist() + df1['topic'].tolist() + df2['title'].tolist() + df2['content'].tolist())
# wc = WordCloud(font_path = r'./MSYH.TTC')
# wc.generate(text)

# plt.imshow(wc)
# plt.axis("off") # 不显示坐标轴
# plt.show()

# 分词并去除停用词
nltk.download('stopwords')
nltk.download('punkt')
stop_words = set(stopwords.words('chinese'))
tokens = word_tokenize(text.lower())
keywords = [word for word in tokens if word.isalpha() and word not in stop_words]

# 生成词云
wordcloud = WordCloud(font_path = r'MSYH.TTC', width=800, height=800, background_color='white')
wordcloud.generate_from_keywords(keywords, font_size=24, min_font_size=10)

# 显示词云
plt.figure(figsize=(8, 8), facecolor=None)
plt.imshow(wordcloud)
plt.axis("off")
plt.tight_layout(pad=0)
plt.show()
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
```

