













6383

















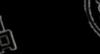




網易 NETEASE















[G] Fi



























网易云课堂 | 涨薪计划







HH





































500



DY DELEGRE ES-AND DA DELSE PROFILE BY DELSTRE







EDEE BA DELSTRE			
	功能	函数	DRE BA DELEUSE
PROFUEE BY DEFENSE	拆分或替换字符串	split、replace	
PROBUEE BY DEVE	17/12	cat	M. TOBE BY DEFEASE
BYCE WEE BA MR.	提取子字符串	extract	网络 是













MINISTER STATES



(INTERIOR

500

网易云课堂 | 涨薪计划

ES- AND BA DEL BA

















































500

TO DESTRUCTION OF THE PARTY OF

(S. 25.2)

replace

3 2.0

英国

```
import sys
v for i in range(0,10079):
      if str(dts.iloc[i,7]).find('(')>=0:
          dts.iloc[i,7]=str(re.findall('\((.*)\)',str(dts.iloc[i,7])))
      else:
         dts.iloc[i,7]=str(dts.iloc[i,7]).replace(',','')
  #替换掉其中的符号
 for i in range(0,10079):
     dts.iloc[i,7]=dts.iloc[i,7].replace(',','').replace('\'','').replace('[','').replace(']','')
  #转化为int数据类型
 dts['GDP(十亿美元)']=dts['GDP(十亿美元)'].astype(int)
 dts['GDP(十亿美元)']=round(dts['GDP(十亿美元)']/1000000000,1)
 dts
                                                        占世界%
                                                                    GDP(十亿美元)
      Unnamed: 0
                排名
                     国家/地区
                             所在洲
                                               GDP(美元)
                                                                                序列
                NaN
                        全世界
                               NaN
                                    .38万亿 (1,384,628,173,213)
                                                           NaN 1960
                                                                         1384.6
                                    5433.0(Z (543,300,000,000) 39.2380% 1960
                                                                          543.3
                         美国
                               美洲
                 1.0
                                   3644.48亿 (364,448,431,697) 26.3210% 1960
                                                                          364.4
                      欧盟地区
              2 NaN
                                                                                             6 36 E DEE BA DELEUSE
```

732.34(Z (73.233.967.692) 5.2891% 1960

PRESIBE BY DEFENSE

73.2



网易云课堂 | 涨薪计划

ES-AND BA BELFA

ALE SE











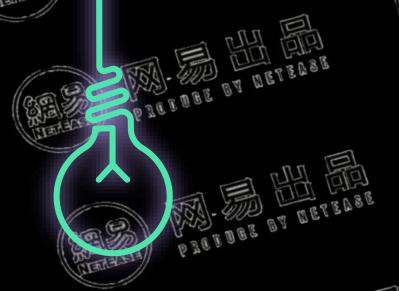


現 地 大 工 一









BACKURE BY DEFEARE





500