```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

    Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.m

```
1 !pip install textstat
2 !pip install PyPDF2
3 !pip install -U textblob
4 import os
5 import glob
6 import textstat
7 import csv
8 import PyPDF2
9 from textblob import TextBlob
```

    Requirement already satisfied: textstat in /usr/local/lib/python3.7/dist-packages (0
    Requirement already satisfied: pyphen in /usr/local/lib/python3.7/dist-packages (from
    Requirement already satisfied: PyPDF2 in /usr/local/lib/python3.7/dist-packages (1.26
    Requirement already up-to-date: textblob in /usr/local/lib/python3.7/dist-packages (6
    Requirement already satisfied, skipping upgrade: nltk>=3.1 in /usr/local/lib/python3
    Requirement already satisfied, skipping upgrade: six in /usr/local/lib/python3.7/dist

```
1 TARGET_FILES = r'/content/drive/My Drive/Control sample/*.pdf'
2 file_list = glob.glob(TARGET_FILES)
```

```
1 new_list=[]
2 name=[]
3 size=[]
4 pagen=[]
5 for file in file_list:
6   print(file)
7   n=(file.split('.pdf',0))
8   name.append(n)
9   size.append((os.stat(file).st_size)/1000000)
10
11
12   if file.endswith(".pdf"):
13     with open (file,'rb') as pdfobject:
14       pdfreader=PyPDF2.PdfFileReader(pdfobject)
15       page_numbers=pdfreader.numPages
16       pagen.append(page_numbers)
17       if pdfreader.isEncrypted == True:
18         pass
19       else:
20           currentpage = 0
21           text = ""
22           while (currentpage <page_numbers):
23             page=pdfreader.getPage(currentpage)
24             try:
25               text=text+page.extractText()
```

```
26            except:
27                print(currentpage) #to print the error page
28            currentpage +=1
29        new_list.append(text)
30
31
```

```
 1 gf=[]
 2 kg=[]
 3 ke=[]
 4 dc=[]
 5 sc=[]
 6 dif=[]
 7 st=[]
 8 ct=[]
 9 sep=[]
10 ses=[]
11 for test in new_list:
12   #print(textstat.gunning_fog(test))
13   #print(textstat.lexicon_count(test, removepunct=True))
14   gf.append(textstat.gunning_fog(test))
15   kg.append(textstat.flesch_kincaid_grade(test))
16   ke.append(textstat.flesch_reading_ease(test))
17   dc.append(textstat.dale_chall_readability_score(test))
18   sc.append(textstat.sentence_count(test))
19   dif.append(textstat.difficult_words(test))
20   st.append(textstat.text_standard(test))
21   ct.append(textstat.lexicon_count(test, removepunct=True))
22   sep.append(TextBlob(test).sentiment.polarity)
23   ses.append(TextBlob(test).sentiment.subjectivity)
```

```
1 ticker=[]
2 year=[]
3 for file in file_list:
4   ticker.append((file.rsplit('/content/drive/My Drive/Control sample/')[1][0:3]))
5   year.append((file.rsplit('/content/drive/My Drive/Control sample/')[1][4:8]))
```

```
1 from itertools import zip_longest
2 data = [ticker, year, gf, kg,ke,dc,sc,dif,st,name,ct, pagen, size,sep,ses]
3 export_data = zip_longest(*data, fillvalue = '')
4 with open('/content/drive/My Drive/MSFT/READABILITY_CONTROLS.csv', 'w', encoding="ISO-8
5        write = csv.writer(file)
6        write.writerow(("TICKER", "YEAR", "GUNNING FOG", "FLESCH-KINCAID GRADE","FLESCH-R
7        write.writerows(export_data)
```