

```
1 import pandas as pd
2 import math as m
3 import numpy as np
4 import scipy.stats as stats
5 import time
6 import csv
7
8 from rdkit import Chem
9 from rdkit.Chem import Descriptors
10
11
12
13 def main():
14     start_time = time.clock()
15
16     # Read in train and test as Pandas DataFrames
17     df_train_all = pd.read_csv('train.csv',usecols=["smiles"])
18     # df_test_all = pd.read_csv('test.csv',usecols=["smiles"])
19
20     # df_all = pd.concat([df_train_all,df_test_all])
21     functions = [i for i in dir(Descriptors) if not (i.startswith("__") or i
22
23     with open("train_rdkit_new.csv",'wb') as w:
24         writer = csv.writer(w)
25         writer.writerow(list(['smiles'])+list(functions))
26         for i, row in df_train_all.iterrows():
27             holder = list([row.smiles])
28             mol = Chem.MolFromSmiles(row.smiles)
29             print "%sth row is processing"%i
30             # print mol.GetNumAtoms()
31             for each in functions:
32                 item = getattr(Descriptors,each)
33                 if callable(item):
34                     try:
35                         holder.extend([item(mol)])
36                     except AttributeError:
37                         print "AttributeError"
38                     except ValueError:
39                         print "ValueError"
40             writer.writerow(holder)
41
42     print "Runtime is ", time.clock() - start_time, "seconds"
43
44
45
46
```

```
47 if __name__ == "__main__":  
48     # execute only if run as a script  
49     main()  
50
```