

```

1 !pip install bitarray
2 !pip install mmh3
3
4 !wget https://bit.ly/3umrYIV -O positive-words.txt
5 !wget https://bit.ly/39HgJmk -O negative-words.txt

```

```

1 f = open("positive-words.txt")
2 S = f.readlines()
3 S = [x.strip() for x in S if not x.startswith(';')]
4 print(S)
5
6 f = open("negative-words.txt")
7 T = f.readlines()
8 T = [x.strip() for x in T if not x.startswith(';')]
9 print(T)

```

```

['', 'a+', 'abound', 'abounds', 'abundance', 'abundant', 'accessable', 'accessible', 'ac
['', '2-faced', '2-faces', 'abnormal', 'abolish', 'abominable', 'abominably', 'abominate

```

```

1 #insert
2
3 import math
4 import mmh3
5 from bitarray import bitarray
6
7 m = len(S)
8
9 n = 8*m
10
11 k = 5
12
13 B = bitarray(n)
14 B.setall(False)
15
16 for s in S:
17     for i in range(k):
18         # with differet seed, hash function is different
19         h_i_s = mmh3.hash(s, i) % n
20         B[h_i_s] = 1
21
22 print("Fraction of 1's in B", B.count()/n)

```

```

Fraction of 1's in B 0.4624439461883408

```

```

1 # lookup (filtering the stream)
2
3 import random

```

```
4
5 S_sample = random.sample(S, 100)
6 T_sample = random.sample(T, 100) # all these words aren't in S
7
8 print(S_sample)
9 print(T_sample)
10
11 # Let's check these words against the Bloom filter
12 cnt_not_in = 0
13 for x in T_sample:
14     not_in = False
15
16     for i in range(k):
17         h_i_x = mmh3.hash(x, i) % n
18         if B[h_i_x] == 0:
19             not_in = True
20             cnt_not_in += 1
21             break
22
23 # how many elements did not pass the filter
24 print('We discarded', cnt_not_in)

['rejoicingly', 'dreamland', 'refunded', 'envy', 'enthuse', 'precise', 'sumptuously', 'c
['maladjustment', 'infamously', 'shun', 'gossip', 'hopelessness', 'cramping', 'dastard',
We discarded 95
```