

- Hamming Distance

```
def hamming_distance(p, q):  
    mismatch = 0  
    for i in range(0, len(p)):  
        if p[i] != q[i]:  
            mismatch += 1
```

return mismatch

GTCG - 3

GC GCG - 5

```
def approx_find(pattern, text, d):  
    indices = []  
    for i in range(0, len(text) - len(pattern)  
                        + 1):  
        if hamming(text[i:i + len(pattern)],  
                    pattern) ≤ d:  
            append i to indices  
  
    return indices.
```

$$\begin{array}{r} \text{text} = \text{GACGACG} \\ \text{pattern} = \text{GACG} \\ \hline L10 \end{array}$$

$$\begin{array}{r} \text{GACG} \\ \text{C G C} \\ \hline L10 \end{array}$$

→ This approach will only work if pattern appears in the string.

```
def frequent_words_with_mismatches(text, k, d):
```

```
    k_mers = {}
```

```
    for i in range(len(text) - k + 1):
```

```
        k_mer = text[i:i+k]
```

```
        count = approxfind(text[i:i+k], text,
```

```
                             k_mers[k_mer] = count    d)
```

```
    maximum = max(k_mers.values())
```

```
    L = [key for key in k_mers if k_mers[key] == maximum]
```

```
    return L
```

AGCAT CCACTAAAA

CACTA

can be grouped using

group anagrams method.

⇒ Pattern might not appear in the string.
Therefore, will have to generate all k -mers of
length k as specified in the problem statement
and store their "occurrences" in a dictionary and
return the maximum or most frequent k -mer.

```
bases = ["A", "T", "G", "C"]  
def generate_kmers(k, bases):  
    if k == 1:  
        return bases
```

```
Small_kmers = generate_kmers(k-1, bases)
```

$k_mers = []$

for k_mer in $small_kmers$:

for b in $bases$

append ($k_mer + b$) to
 k_mers

return k_mers

4-length k_mer :

$_x_x_x_ = 4 \cdot 4 \cdot 4 \cdot 4$
 \uparrow
4 options
since bases
 $= ['A', 'G', 'T', 'C']$
 $= 4^4 = 256$
 k_mers

$= 4^k$ k_mers

$_x_x_x_x_ = 4^5 = 1024$
 k_mer

AAT
AAT
ATT

ATGCAT
TACGTA
ATGCCA

- Frequent Words with Mismatches and Reverse Complements

```
def foo(text, k, d):
```

```
    bases = ['A', 'T', 'G', 'C']
```

```
    s = ""
```

```
    visited = set()
```

```
    k_mers = generate_kmers(k, bases)
```

```
    k_mers_count = {}
```

```
    for kmer in k_mers:
```

```
        if kmer not in visited:
```

```
            r_c_kmer = reverse_complement(kmer)
```

```
            count1 = len(approx_find(kmer, text, d)[1])
```

```
            count2 = len(approx_find(r_c_kmer, text, d)[1])
```

```
            k_mers_count[kmer] = sum(count1, count2)
```

```
    # return kmer with most frequency.
```

Have to
make sure if
kmer and
its reverse
complement
same add once.