

# Pseudonymization

Requires: numpy

Tested on Python 2.7.

## Introduction

The main purpose of this library is to provide a secure means of pseudonymizing token sets (sets of hashable objects) so that Jaccard similarity scores can be estimated from token set pairs. Typically the token sets will be derived from strings. The motivating application is privacy-preserving record linkage.

Instances of the *C\_hash* class are pseudo-random integers, where the expected proportion of collisions on the bits of a pair of hashes is a function of the underlying Jaccard score. Estimation is (practically) unbiased and larger token sets do not result in a larger number of set bits. This offers two distinct advantages over Bloom filters.

They are more expensive to generate than Bloom filters, but more resistant to attacks on the token set sizes and frequency attacks on the tokens.

## Concatenated class

Create a concatenated 1-bit minwise hasher of length  $m$  using an underlying  $q$ -bit hash function. The *seed* argument seeds the RNG so that multiple instances can be created with identical behaviour.

```
>>> import pseudo
>>> h = pseudo.Concatenated(1024, q=64, seed=6)
>>>
```

The returned object can be used to hash token sets via its *hash* method. We can use padded bigrams for tokenization.

```
>>> import tokenization
>>> tokens1 = tokenization.n_grams('thompson', 2, pad=True)
>>> tokens2 = tokenization.n_grams('timpson', 2, pad=True)
>>> tokens1
['_t', 'th', 'ho', 'om', 'mp', 'ps', 'so', 'on', 'n_']
>>> tokens2
['_t', 'ti', 'im', 'mp', 'ps', 'so', 'on', 'n_']
>>>
```

We can generate the true Jaccard score for these tokens sets.

```
>>> pseudo.Jaccard(tokens1, tokens2)
0.5454545454545454
>>>
```

We can securely tokenize the token sets using the *Concatenated* instance created earlier. In fact we can create a second instance using the same seed to simulate pseudonymization in distinct locations.

```
>>> hash1 = h.hash(tokens1)
>>> hash2 = pseudo.Concatenated(1024, q=64, seed=6).hash(tokens2)
```

Hashes print as long ints.

```
>>> hash1
1718952352000389468904029704285699447798704608855737624884462738203935658658
4148534364561616172834226109268120168857065074418720433775920765909048566531
0779799536599898998509149858262826292348966905322667620008405806604077273651
3831543477516390027969214799552854907435709347042965904195673546553281199007
12436L
>>>
```

In fact, these *C\_hash* instances are derived from Python's long integer type and are associated with two extra attributes - the number of significant bits and the XOR compression factor (which are both required for Jaccard score estimation).

Given two hashes we can estimated the Jaccard score.

```
>>> pseudo.J_hat_from_conc(hash1, hash2, 1024)
0.572265625
>>>
```

The above exploits the fact that concatenated 1-bit hashes are subclasses of *long*. The above call will also work with any Python *int* or *long* instances. It is the user's responsibility to ensure that hashes are comparable (same length, generated by the same 1-bit hashes etc.).

As we know the true score we can calculate the variance of estimation.

```
>>> pseudo.var_J_hat(0.5454545454545454, 1024, N=1)
0.0006860149793388431
>>>
```

## C\_hash class

This class contains a couple of methods for compressing hashes, and a property that returns the bitstring representation of the hash.

The *compressed* method returns a new hash with fewer significant bits. The returned value is simply the value we would have had if we had used a smaller number of minwise hash functions.

```
>>> comp1 = hash1.compressed(256)
>>> comp2 = hash2.compressed(256)
>>> comp1.m, comp1.N
(256, 1)
>>> pseudo.J_hat_from_conc(comp1, comp2, 256, N=1)
0.65625
>>> pseudo.var_J_hat(0.5454545454545454, 256, N=1)
0.0027440599173553723
>>>
```

The *XOR* method returns XOR compressed hashes with compression factor N (which must be a positive power of 2).

```
>>> comp1 = hash1.XOR(N=4)
>>> comp2 = hash2.XOR(N=4)
>>> comp1.m, comp1.N
(256, 4)
>>> pseudo.J_hat_from_conc(comp1, comp2, 256, N=4)
0.5149417859767794
>>> pseudo.var_J_hat(0.5454545454545454, 256, N=4)
0.009197553658364058
>>>
```

Here the variance is greater than for the 'bitstring slice' form of compression for this Jaccard score, but would be lower for sufficiently high Jaccard scores.

We can print the binary representation of a hash.

```
>>> comp1.digits
'011011111011011111100110001010010011100010101001010111110110010000010000011
1000110001011001101010000110010001110001001110111100110000101101111101100101
001111010011011111110101110110010010101100010111011010010010010100111110111
01100000111001000011000001000'
```

This shows that compression still produces what is (for all practical purposes) a random integer.

## BloomFilter class

We need to generate a suitable list of hash functions, supplying the number of hash functions, the bit

length of the Bloom filter and an optional seed.

```
>>> funcs = set_like.k_hashes(4, 1024, seed=6)
>>>
```

For the Bloom filter we need to supply the length and the hash functions.

```
>>> bf1 = set_like.BloomFilter(1024, funcs)
>>>
```

We can generate a second Bloom filter with distinct, but identically behaving, hash functions. This simulates the generation of comparable Bloom filters in distinct locations.

```
>>> bf2 = set_like.BloomFilter(1024, set_like.k_hashes(4, 1024, seed=6))
>>>
```

We can add our tokens (or we could have supplied them as an optional keyword argument 'items' when creating the Bloom filters).

```
>>> for token in tokens1:
    bf1.add(token)
```

```
>>> for token in tokens2:
    bf2.add(token)
```

We can use the standard Jaccard score estimator or the bias-corrected estimator due to Swamidass and Baldi (2007). Note the extra argument required for the corrected estimator.

```
>>> pseudo.J_hat_from_bf(bf1.bits, bf2.bits)
0.5348837209302325
>>> pseudo.J_hat_from_bf_corrected(bf1.bits, bf2.bits, 1024)
0.5272637368168046
>>>
```

Of course, we can simply use the Bloom filter as a probabilistic set data structure.

```
>>> 'th' in bf1
True
>>>
```

Note: we would generally use a larger number of hash functions for this purpose (to minimize the

false positive rate). For a 1024 bit Bloom filter and 9 items we might use 79 hash functions.

```
>>> set_like.opt_k(1024, 9)
78.86474587704267
```

Both concatenated hashes and Bloom filters can be (partially) serialized using their *hex* or *digits* methods. Only the *digits* method automatically pads the returned string to convey the number of significant bits. The builtins *str* and *repr* can also be used.

## References

Swamidass, S. J. and Baldi, P. (2007), *Mathematical correction for fingerprint similarity measures to improve chemical retrieval*, Journal of chemical information and modeling (ACS Publications) **47** (3): pp.952-964