

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]


```
{ 'filename': '13', 'phone': ['145.973.4455']},
{ 'scan_results': {'email': [], 'phone': ['145.973.1386']}},
{ 'filename': '5', 'scan_results': { 'email': ['lair@enron.com', '28pgc.com', 'teresa.hess@enron.com', '28pgc.com', 'william.griffith@elpaso.com', 'adale.zuroff@elpaso.com', 'curate@elianteenergy.com', 'ernesto.ocha@elpaso.com', 'oyal@elpaso.com', 'ym.bla@enron.com', 'acey@enron.com', 'randy.young@usouthpl.com', 'michael.k.rasmussen@williams.com', 'sanie.k.nielsen@williams.com', 'jarrequ@mail.pnm.com', 'pdavidson@scocalgas.com', 'mark.gracey@elpaso.com', 'tom.gwilliams@iroquois.com', 'ay.story@pnp.org', 'ris.g.king@dom.com', '1111@elianteenergy.com', 'jwilliams@elianteenergy.com', 'curate@elianteenergy.com', 'hickman@nsource.com', 'tyoung@nsource.com', 'tbell@spc.com', 'gauga@tucacoraragas.com', 'acth@questar.com', 'kburc@duke-energy.com', 'milmccain@duke-energy.com', 'ken.schubert@transcanada.com', 'dale.m.davis@williams.com', 'enron.messaging.administration@enron.com', '415.973.1386', '713) 853-7637', '888-627-7024', '888-627-7024', '415.973.1386'}},
{ 'filename': '14', 'scan_results': {'email': [], 'phone': ['145.973.1386']}},
{ 'filename': '18', 'scan_results': { 'email': ['athaway@dom.com', 'athaway@dom.com', 'ercard@enron.com', 'enron.messaging.administration@enron.com'], 'phone': []}},
{ 'filename': '9', 'scan_results': {'email': [], 'phone': []}},
{ 'filename': '11', 'scan_results': {'email': [], 'phone': []}},
{ 'filename': '5', 'scan_results': {'email': [], 'phone': []}},
{ 'filename': '16', 'scan_results': { 'email': [], 'phone': ['816-756-1500', '816-756-1635', '713-853-1731', '713-853-5660']}},
{ 'filename': '6', 'scan_results': { 'email': ['steven.january@enron.com'], 'phone': ['713-853-5660']}},
{ 'filename': '17', 'scan_results': {'email': [], 'phone': []}},
{ 'filename': '1', 'scan_results': {'email': [], 'phone': []}},
{ 'filename': '10', 'scan_results': {'email': [], 'phone': []}},
{ 'scan_results': {'email': [], 'phone': ['713-853-5660']}},
{ 'filename': '8', 'scan_results': {'email': [], 'phone': []}},
{ 'filename': '2', 'scan_results': {'email': [], 'phone': []}},
{ 'filename': '13', 'scan_results': {'email': ['dave.tyler@enconenergy.com', 'phone': []]},
{ 'scan_results': { 'email': ['jmosher@oneok.com', 'tjones@oneok.com', 'enm18.lee@enron.com'], 'phone': ['713) 853-7121']}},
{ 'filename': '14', 'scan_results': {'email': [], 'phone': ['713-853-5660']}},
{ 'filename': '1', 'scan_results': {'email': ['mary.darveaux@enron.com', 'phone': []]},
{ 'filename': '8', 'scan_results': {'email': [], 'phone': []}},
{ 'filename': '13', 'scan_results': {'email': [], 'phone': []}},
{ 'filename': '5', 'scan_results': {'email': [], 'phone': []}},
{ 'filename': '1', 'scan_results': {'email': [], 'phone': []}},
{ 'scan_results': { 'email': ['ipayit@enron.com', 'ipayit@enron.com'], 'phone': ['713) 853-7121']}},
{ 'filename': '18', 'scan_results': {'email': [], 'phone': ['713-345-4727']}},
{ 'filename': '11', 'scan_results': {'email': [], 'phone': []}},
{ 'filename': '2', 'scan_results': {'email': [], 'phone': []}},
{ 'filename': '1', 'scan_results': {'email': [], 'phone': []}},
{ 'scan_results': { 'email': ['877490708@pagenetmessage.net'], 'phone': ['8774899']}},
{ 'filename': '6', 'scan_results': {'email': [], 'phone': []}},
{ 'filename': '2', 'scan_results': {'email': [], 'phone': []}},
{ 'filename': '10', 'scan_results': {'email': [], 'phone': []}},
{ 'scan_results': { 'email': [], 'phone': ['713-857-2103', '281-531-0385']}},
{ 'filename': '19', 'scan_results': { 'email': [], 'phone': ['877-550-7985', '402-896-2585']}},
{ 'filename': '8', 'scan_results': { 'email': [], 'phone': ['605-272-5680', '605-272-5296']}},
{ 'filename': '13', 'scan_results': {'email': [], 'phone': ['605-272-5680', '605-272-5296']}},
{ 'filename': '4', 'scan_results': {'email': [], 'phone': ['605-272-5680', '605-272-5296']}},
{ 'filename': '12', 'scan_results': { 'email': ['cynthia.biggs@trsvelpark.com'], 'phone': ['402 397 8000', '902 926 4036', '905 7033', '713) 853-5842', '800 523-6586', '215 245-4707']}},
{ 'filename': '2', 'scan_results': {'email': [], 'phone': []}},
{ 'filename': '13', 'scan_results': {'email': [], 'phone': []}},
{ 'filename': '5', 'scan_results': {'email': [], 'phone': ['605-272-5680']}},
{ 'scan_results': {'email': ['ulligan@enron.com'], 'phone': []}},
{ 'filename': '3', 'scan_results': {'email': ['stojic@gbmdc.com'], 'phone': []}}
```

(T7) Hashing with salt

- cryptographic package: <https://www.pycryptodome.org/en/latest/src/hash/hash.html>
- reverse lookup MD5 hash: <https://md5.gromweb.com/>

```
In [26]: from Crypto.Hash import MD5

In [27]: # MD5.new() only accepts binary values so by putting a 'b' in front of the string, it
         # easy_password_hash = MD5.new(b'password')

         easy_password_hash.hexdigest()

Out[27]: '5f4dccc3b5aa765d61d8327deb882cf99'
```

Avalanche effect

In cryptography, the avalanche effect is the desirable property of cryptographic algorithms wherein if an input is changed slightly (for example, flipping a single bit), the output changes significantly (e.g., half the output bits flip).

This property provides the following benefits:

- prevent people from making predictions about the input, being given only the output
- makes it obvious that the data's integrity is affected even if it is only a 1 bit change.

```
In [28]: print('The MD5 hash of \'' + '\') is {}'.format('password'), MD5.new(b'password').hexdigest()
         print('The MD5 hash of \'' + '\') is {}'.format('password', MD5.new(b'password').hexdigest())

         for i in range(10):
             plaintext_password = 'password{}'.format(i)
             print('The MD5 hash of \'' + '\') is {}'.format(plaintext_password, MD5.new(plaintext_password).hexdigest())

The MD5 hash of 'password' is 5f4dccc3b5aa765d61d8327deb882cf99
The MD5 hash of 'password' is a826176c6495c516189db91770e20ce
The MD5 hash of 'password0' is 305e4f53ce823e11a4e9d5f80bcb86c
The MD5 hash of 'password1' is 76da180b36699da0a8c2797eeab0e4c
The MD5 hash of 'password2' is 6cb75f652a9b52798eb6cf2201057c73
The MD5 hash of 'password3' is 819b06436b989dc9b579efdc9094f28e
The MD5 hash of 'password4' is 34cc93ace0ba9e3f62235d4a9f97b16c
The MD5 hash of 'password5' is db0e6d04aac5f0c7feda033ac5f5456
The MD5 hash of 'password6' is 218d427aebc6cecae69ad8a03d9a36bf
The MD5 hash of 'password7' is 00cd07b3942cf6b290ceb9706aca6a43
The MD5 hash of 'password8' is b37e00be3698480bc43d4e6c2c738
The MD5 hash of 'password9' is 5469dd95ac183c943780ed7027d128a
```

Reverse lookup hashes

Try searching for '5f4dccc3b5aa765d61d8327deb882cf99' in the internet. Are you able to find the original text from the hash?

Rainbow table attacks makes use of precomputed hashes tables to quickly identify your password from hashes.

To circumvent such attacks, add salt before hashing to significantly increase the size of the precomputed hash table needed to do reverse lookup of hashes. The longer the salt length, the larger the required precomputed hash table

Salt generation

Salt are additional characters that are randomly generated and added to the back of the original string before hashing to make dictionary attacks harder to execute (as the dictionary of known passwords would have to be multiple times bigger than dictionaries for unsalted password)

```
In [29]: import random

         def generate_salt(char_set: str = None, length: int = 16) -> str:
             """Generates a random string of alphanumeric characters as salt

             Args:
                 char_set: string containing the characters used for salt generation
                 length: the length of salt required

             Return:
                 randomly generated salt based on char_set and length

             """

             if char_set is None:
                 char_set = "0123456789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ"

             chars=[]
             for i in range(length):
                 chars.append(random.choice(char_set))

             return "".join(chars)

In [30]: salt = generate_salt()

         print(salt)

o3n3oyFtU0xwF6psr
```

```
In [31]: # append salt to the original easy password
         easy_password_hash.update(salt.encode('ascii'))

         easy_password_hash.hexdigest()

Out[31]: '562cc314626b51d93f15de1a40bfbb5e'
```

Reverse lookup hashes with salt

Try searching for the hash generated by adding salt to the value. Are you able to find the original text from the hash?

(T8) Field level encryption - Example using Symmetric Key Encryption

```
In [32]: import pandas as pd
         import requests
         import io

         # Get data from data.gov.sg

         api_response = requests.get('https://data.gov.sg/api/action/package_show?id=acra-info-
         resources = api_response['result']['resources']
         .json()

         data_set_url = ''

         for resource in resources:
             if resource['name'] == 'ACRA Information on Corporate Entities (\X')':
                 data_set_url = resource['url']
                 print('Data set URL:{}'.format(data_set_url))

         acra_csv_data = io.StringIO(requests.get(data_set_url).content.decode('utf-8'))
         acra_csv_df = pd.read_csv(acra_csv_data, dtype=str)

         acra_csv_df.head()
```

	Sole Proprietor	PUBS	na	CIRCULAR ROAD
1	Partnership	WHOLESALE TRADE OF A VARIETY OF GOODS WITHOUT	na	SENJA ROAD
2	na	CAFES AND COFFEE HOUSES	na	ROBINSON ROAD
3	Sole Proprietor	PASSENGER LAND TRANSPORT N.E.C (EG PRIVATE CA...	na	CHOA CHU KANG CRESCENT
4	na	OTHER HOLDING COMPANIES	INVESTMENT	CECIL STREET

5 rows x 94 columns

If internet access is not available. Use local copy of dataset.

```
In [34]: # read ACRA companies table starting with letter X (taken from data.gov.sg)
         # https://data.gov.sg/dataset/acra-information-on-corporate-entities
         # acra_csv_df = pd.read_csv('acra-information-on-corporate-entities-x.csv', dtype=str)
         # acra_csv_df.head()
```

```
In [35]: from Crypto.Cipher import AES
         from Crypto.Hash import MD5
         from typing import Tuple

         def aes_encrypt(value:bin, key:bin) -> Tuple(bin, bin):
             """Function to perform AES encryption using SIV mode (as SIV mode doesn't require
             Args:
                 value: plaintext value to be encrypted
                 key: encryption key used to perform encryption

             Return:
                 ciphertext and digest in binary data type

             """

             # convert key into MD5 as SIV mode requires minimum 32 byte key
             hash_key = MD5.new(key).hexdigest().encode('ascii')
             cipher = AES.new(hash_key, AES.MODE_SIV)

             return cipher.encrypt_and_digest(value)

         def aes_decrypt(ciphertext:bin, digest:bin, key:bin) -> bin:
             """Function to perform AES decryption using SIV mode (as SIV mode doesn't require
             Args:
                 ciphertext: ciphertext to be decrypted
                 digest: digest produced from encryption
                 key: encryption key used to perform decryption

             Return:
                 ciphertext and digest in binary data type

             """

             # convert key into MD5 as SIV mode requires minimum 32 byte key
             hash_key = MD5.new(key).hexdigest().encode('ascii')
             cipher = AES.new(hash_key, AES.MODE_SIV)

             return cipher.decrypt_and_verify(ciphertext, digest)
```

Perform encryption on simple message

```
In [36]: secret_key = b'secret'

         message = b'secret msg'

         ciphertext_l1, digest_l1 = aes_encrypt(message, secret_key)

         print('My encrypted message is:{}'.format(ciphertext_l1))
         print('My encrypted message digest is:{}'.format(digest_l1))
         print('My encrypted secret message is:{}'.format(aes_decrypt(ciphertext_l1, digest_l1, secret_key)))

         b'\xe9\xcf\xcd\x83\xad\x45\x4c\x4d'

         My encrypted message digest is:
         b'>2\x4f\xce\x85\x85\x04\x3b\x05\xbd\x4\xac\xca\x84\x8f'

         My decrypted secret message is:
         b'secret msg'
```

Perform encryption on field 'entity_name'

```
In [37]: encrypted_df = acra_csv_df.copy()

         encrypted_df['entity_name'] = encrypted_df['entity_name'].apply(lambda x: aes_encrypt
         encrypted_df.head()
```

	Sole Proprietor	PUBS	na	CIRCULAR ROAD
1	Partnership	WHOLESALE TRADE OF A VARIETY OF GOODS WITHOUT	na	SENJA ROAD
2	na	CAFES AND COFFEE HOUSES	na	ROBINSON ROAD
3	Sole Proprietor	PASSENGER LAND TRANSPORT N.E.C (EG PRIVATE CA...	na	CHOA CHU KANG CRESCENT
4	na	OTHER HOLDING COMPANIES	INVESTMENT	CECIL STREET

5 rows x 94 columns

Perform decryption on field 'entity_name'

```
In [38]: encrypted_df['entity_name'].apply(lambda x: aes_decrypt(x[0], x[1], secret_key).decode
         .decode('utf-8'))

0 X
1 X & B TRADING
2 X & E PASSION PTE. LTD.
3 X & P SERVICES
4 X & H INTERNATIONAL PTE. LTD.
Name: entity_name, dtype: object
```

(T9) Tokenization

```
In [39]: # create functions for generating and storing tokens
         import random

         class TokenDict:
             """Class for storing randomly generated tokens

             Typical usage example:
                 token_dict = TokenDict()

             """

             def __init__(self):
                 # protected attributes
                 self.__token_dict = {}

             def generate_token(self) -> str:
                 NUMBERS = "0123456789"
                 chars=""
                 for i in range(16):
                     chars.append(random.choice(NUMBERS))

                 token = "".join(chars)
                 return token

             def tokenize_value(self, real_value:str) -> bin:
                 if self.retrieve_token(real_value) is not None:
                     raise ValueError('value is already in token dictionary')
                 token=self.generate_token()
                 while token in self.__token_dict.keys() or token==real_value:
                     token = self.generate_token()
                 self.__token_dict[token]=real_value
                 return token

             def retrieve_value(self, token:str) -> str:
                 return self.__token_dict[token]

             def retrieve_token(self, real_value: str) -> str:
                 for k,v in self.__token_dict.items():
                     if v == real_value:
                         return k
                 return None
```

Initialise token dictionary class and enter values into the dictionary

```
In [40]: token_dict_1 = TokenDict()

         token_l1 = token_dict_1.tokenize_value('1234567812345678')
         print('First Token: {}'.format(token_l1))
         token_2 = token_dict_1.tokenize_value('8888888888888888')
         print('Second Token: {}'.format(token_2))

         First Token: 1963812867436780
         Second token: 9256726201869807
```

Retrieve token from value and vice versa

```
In [41]: print('Retrieve second token from second token: {}'.format(token_dict_1.retrieve_valu
         print('Retrieve first token from first value: {}'.format(token_dict_1.retrieve_token
         Retrieve second token from second token: 8888888888888888
         Retrieve first token from first value: 1963812867436780
```

Try to extract dictionary from token dictionary class

```
In [42]: # prefix attribute with __ to protect it from access outside of the class

         try:
             new_dict = token_dict_1.__token_dict
         except AttributeError as e:
             print('Attribute Error: {}'.format(str(e)))

Attribute Error: 'TokenDict' object has no attribute '__token_dict'
```

(T13) Data File Integrity Verification - Example using Asymmetric Key Encryption

- cryptographic package: https://www.pycryptodome.org/en/latest/src/public_key/rsa.html
- Public-key cryptosystems are convenient in that we do not require the sender and receiver to share a common secret in order to communicate securely (among other useful properties).

However, they often rely on complicated mathematical computations and are thus generally much more inefficient than comparable symmetric-key cryptosystems. In many applications, the high cost of encrypting long messages in a public-key cryptosystem can be prohibitive.

The asymmetric encryption algorithm only allows encryption of data smaller than the key length. To break large files into chunks smaller than key length to perform asymmetric encryption would be computationally inefficient.

However, this can be addressed by hybrid systems by using a combination of both, where a symmetric session key is used to encrypt the file and the asymmetric encryption be done on the symmetric session key.

Define functions:

```
In [43]: from Crypto.PublicKey import RSA
         from Crypto.Cipher import PKCS1_OAEP
         from Crypto.Signature import pkcs1_15
         from Crypto.Hash import MD5
         import base64
         from typing import Tuple

         def generate_keys() -> Tuple(RSA.RsaKey, RSA.RsaKey):
             """Generate asymmetric key pairs

             Args: None

             Return:
                 private key, public key

             """

             modulus_length = 1024

             key = RSA.generate(modulus_length)

             pub_key = key.publickey()

             return key, pub_key

         def aes_encrypt(plaintext_message: bin, public_key: RSA.RsaKey) -> bin:
             """Asymmetric encryption function

             Note: encryption can only be done using public key for security reasons.
             For signing, please refer to Crypto.Signature

             Args:
                 plaintext_message: plaintext message that you want to encrypt
                 public_key: recipient's public key to be used for encrypting plaintext message

             Return:
                 ciphertext_message

             """

             encryptor = PKCS1_OAEP.new(public_key)
             encrypted_msg = encryptor.encrypt(plaintext_message)
             encoded_encrypted_msg = base64.b64encode(encrypted_msg)
             return encoded_encrypted_msg

         def aes_decrypt(ciphertext_message: bin, private_key: RSA.RsaKey) -> bin:
             """Asymmetric decryption function

             Note: decryption can only be done using private key for security reasons.
             For signing, please refer to Crypto.Signature

             Args:
                 ciphertext_message: encrypted message sent from sender encrypted with recipie
                 private_key: recipient's private key to be used for decrypting ciphertext message

             """

             decryptor = PKCS1_OAEP.new(private_key)
             encrypted_decrypted_msg = base64.b64decode(ciphertext_message)
             decoded_decrypted_msg = decryptor.decrypt(decoded_decrypted_msg)
             return decoded_decrypted_msg

         def generate_session_key(char_set: str = None, length: int = 16) -> str:
             """Generates a random string of alphanumeric characters as session key for AES encryption

             Args:
                 char_set: string containing the characters used for session key generation
                 length: the length of session key required

             Return:
                 randomly generated salt based on char_set and length

             """

             if char_set is None:
                 char_set = "0123456789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ"

             chars=[]
             for i in range(length):
                 chars.append(random.choice(char_set))

             return "".join(chars)
```

Perform encryption and decryption on simple message

```
In [44]: private, public = generate_keys()
         message = b'Hello world'
         encoded = aes_encrypt(message, public)
         sym_decrypt(encoded, private)

         b'Hello world'
```

Perform public key cryptography example on files

```
In [45]: # prepare key pairs

         # sender key pairs
         s_private, s_public = generate_keys()

         # receiver key pairs
         r_private, r_public = generate_keys()
```

Sender's perspective

Has knowledge of:

- sender's private key
- sender's public key
- receiver's public key

Perform encryption on file

As asymmetric encryption can only work for messages smaller than the asymmetric key length, we need to use symmetric encryption to encrypt the file and use the asymmetric encryption to encrypt the symmetric encryption key (also known as the **session key**).

```
In [46]: # prepare files and filenames
         filename = 'acra-information-on-corporate-entities-x.csv'
         signed_md5_filename = 'signed_md5'
         encrypted_data_filename = 'encrypted_file'
         encrypted_key_filename = 'encrypted_key'

         # generate asymmetric key for file encryption. This is also known as a session key as
         symm_session_key = generate_session_key().encode('ascii')

         Start symmetric encryption using session key and write to file
```

```
In [47]: # read file to be encrypted. We will use the same data.gov.sg acra companies dataset.
         with open(filename, 'rb') as f:
             data = f.read()

         # write encrypted data into file
         with open(encrypted_data_filename, 'wb') as f:
             cipher, digest = aes_encrypt(data, symm_session_key)
             # print(len(encrypted_data))
             # print(len(cipher))
             f.write(digest)
             f.write(cipher)

         print('File () is encrypted using AES and saved as {}'.format(filename, encrypted_data))

         File acra-information-on-corporate-entities-x.csv is encrypted using AES and saved as encrypted_file
```

Start asymmetric encryption of session key using recipient's public key and write into file

=> only recipient's private key can be used to decrypt the encrypted session key

```
In [48]: # write encrypted session key into file
         with open(encrypted_key_filename, 'wb') as f:
             f.write(aes_encrypt(symm_session_key, r_public))

         print('Session is encrypted using recipient's public key and saved as {}'.format(enco
         Session is encrypted using recipient's public key and saved as encrypted_key
```

Create file hash for file integrity check by recipient and sign it with sender's private key

```
In [49]: # create file digest
         md5_digest = MD5.new(data).digest()

         # write md5 hash of data to file
         with open(signed_md5_filename, 'wb') as f:
             # print(len(md5_digest))
             # print(len(pkcs1_15.new(s_private).sign(MD5.new(md5_digest))))
             f.write(pkcs1_15.new(s_private).sign(MD5.new(md5_digest)))

         print('Hash digest of file is signed using sender's private key and saved as {}'.form
         Hash digest of file is signed using sender's private key and saved as signed_md5
```

Receiver's perspective

Has knowledge of:

- receiver's private key
- receiver's public key
- sender's public key

Decrypt session key using receiver's private key

```
In [50]: # read encrypted session key from file and decrypt it
         with open(encrypted_key_filename, 'rb') as f:
             decrypted_file_symm_key=aes_decrypt(f.read(), r_private)

         print('Successfully decrypted session key: \nSession Key:{}'.format(decrypted_file_sym
         Successfully decrypted session key:
         Session Key:b'\x12\x4frrvU5'
```

Use session key to decrypt file

```
In [51]: decrypted_filename = 'decrypted.csv'

         # read encrypted data with decrypted symmetric key
         with open(encrypted_data_filename, 'rb') as f:
             sym_digest = f.read(16)
             sym_data = f.read()
             decrypted_data = aes_decrypt(sym_data, sym_digest, decrypted_file_symm_key)

         # write decrypted data into csv file
         with open(decrypted_filename, 'wb') as f:
             f.write(decrypted_data)
```

	Sole Proprietor	PUBS	na	CIRCULAR ROAD
1	Partnership	WHOLESALE TRADE OF A VARIETY OF GOODS WITHOUT	na	SENJA ROAD
2	na	CAFES AND COFFEE HOUSES	na	ROBINSON ROAD
3	Sole Proprietor	PASSENGER LAND TRANSPORT N.E.C (EG PRIVATE CA...	na	CHOA CHU KANG CRESCENT
4	na	OTHER HOLDING COMPANIES	INVESTMENT	CECIL STREET

5 rows x 94 columns

Validate sender identity

```
In [52]: # validate signed msg digest from file
         with open(signed_md5_filename, 'rb') as f:
             received_md5_digest = f.read(16)
             received_digest_signature = f.read(16)

         # verify signature of md5 digest to validate sender identity
         try:
             pkcs1_15.new(s_public).verify(MD5.new(received_md5_digest), received_digest_signature)
             print('Signature is valid. Sender identity confirmed')
         except (ValueError, TypeError):
             print('Signature is not valid. Sender identity invalid')

         Signature is valid. Sender identity confirmed
```

Validate file integrity

```
In [53]: # perform md5 hash on data and compare with received md5 hash, to verify file integrity
         decrypted_data_digest = MD5.new(decrypted_data).digest()

         if received_md5_digest == decrypted_data_digest:
             print('File integrity verified')
         else:
             print('File integrity failed')
```

File integrity verified

Format preserving encryption

```
In [54]: # library for format preserving encryption using ffx method
         import pyffx

         # create fpe object to perform encryption and decryption
         fpe = pyffx.Integer(b'secret-key', length=16)
```

```
In [55]: fpe.ciphertext = fpe.encrypt('1234567812345678')

         fpe.ciphertext

Out[55]: 964510889077199
```

```
In [56]: e.decrypt(fpe.ciphertext)

Out[56]: 1234567812345678
```

Synthetic Data Generation

Generate data points for linear regression

```
In [57]: from sklearn import datasets
         import matplotlib.pyplot as plt
         import numpy as np

         # generate regression dataset
         inputs, outputs = datasets.make_regression(n_features=1, noise=0.0)
```

```
# reshape the x values
input_x = np.concatenate(inputs, axis = 0)

# get gradient and coefficients
m, b = np.polyfit(input_x, outputs, 1)

# scatter plot
plt.plot(input_x, outputs, 'o')
```

```
# plot regression line
plt.plot(input_x, np.array(input_x) * m + b)

Out[57]: [Cmplotlib.lines.Line2D at 0x7f80f99ee80c]
```


Add some noise to the regression

```
In [58]: # generate regression dataset
         inputs, outputs = datasets.make_regression(n_features=1, noise=20.0)

         # reshape the x values
         input_x = np.concatenate(inputs, axis = 0)

         # get gradient and coefficients
         m, b = np.polyfit(input_x, outputs, 1)

         # scatter plot
         plt.plot(input_x, outputs, 'o')
```

```
# plot regression line
plt.plot(input_x, np.array(input_x) * m + b)

Out[58]: [Cmplotlib.lines.Line2D at 0x7f80f0c5d070c]
```


Generate data points for clustering

```
In [59]: # generate circles dataset
         labels = datasets.make_circles()

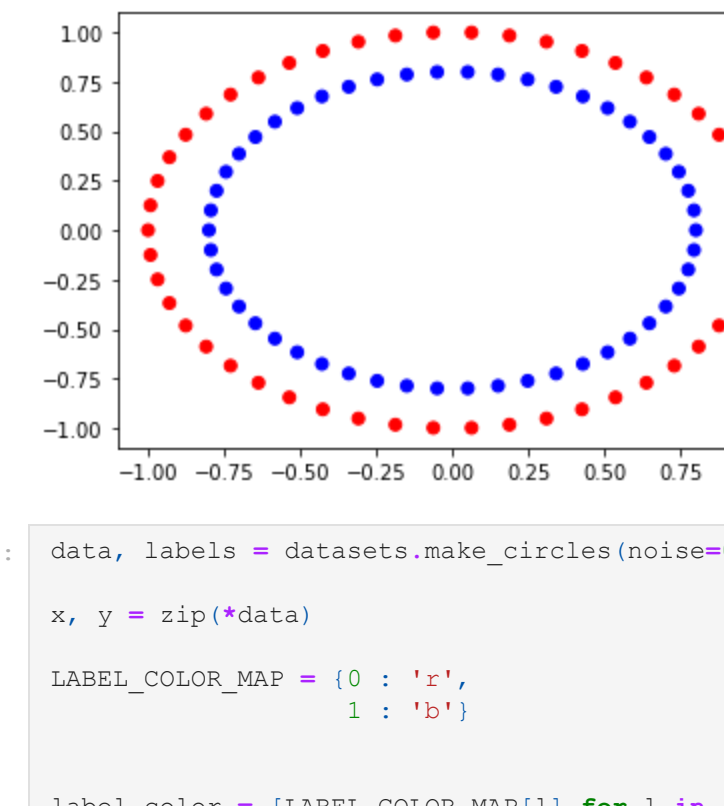
         # split axis
         x, y = zip(*data)

         # define color map for labels
         LABEL_COLOR_MAP = {0: 'r',
                             1: 'b'}
```

```
# label_color = [LABEL_COLOR_MAP[i] for i in list(labels)]

plt.scatter(x, y, c=label_color)

Out[59]: <matplotlib.collections.PathCollection at 0x7f80f0ee280>
```

```
In [60]: data, labels = datasets.make_circles(noise=0.09)

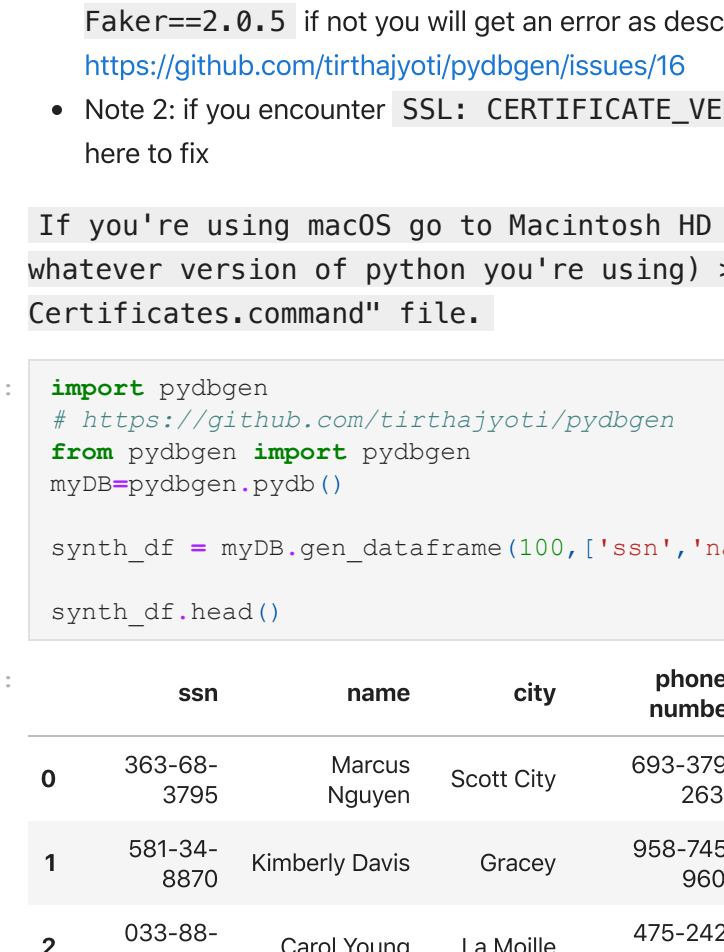
x, y = zip(*data)

LABEL_COLOR_MAP = {0 : 'r',
                    1 : 'b'}

label_color = [LABEL_COLOR_MAP[l] for l in list(labels)]

plt.scatter(x, y, c=label_color)

Out[60]: <matplotlib.collections.PathCollection at 0x7f80f0a3dac0>
```



Generate pandas dataframe of customer (categorical/attribute) data

- Note 1: when doing package installation, do downgrade faker package `pip install Faker==2.0.5` if not you will get an error as described in <https://github.com/tirthajyoti/pydbgen/issues/16>
- Note 2: if you encounter `SSL: CERTIFICATE_VERIFY_FAILED error` do follow instructions here to fix


If you're using macOS go to Macintosh HD > Applications > Python3.6 folder (or whatever version of python you're using) > double click on "Install Certificates.command" file.

```
In [61]: import pydbgen
import pprint
# https://github.com/tirthajyoti/pydbgen
from pydbgen import pydbgen
myDB=pydbgen.pydb()

synth_df = myDB.gen_dataframe(100,['ssn','name','city','phone','date','zipcode','job_title'])
synth_df.head()
```

	ssn	name	city	phone-number	date	zipcode	job_title
0	363-68-3795	Marcus Nguyen	Scott City	693-379-2636	1999-05-23	93035	Environmental consultant
1	581-34-8870	Kimberly Davis	Gracey	958-745-9607	1997-11-30	81349	Civil Service fast streamer
2	033-88-0392	Carol Young	La Moille	475-242-0524	1973-03-11	79846	Armed forces technical officer
3	874-59-7154	Thomas Johnson	Miami Lakes	297-233-7340	1976-11-26	86730	Environmental education officer
4	454-93-1995	Jessica Hess	Fern Creek	524-985-2616	1994-10-17	02340	Scientist, marine

Homomorphic Encryption

 Drawing Fully homomorphic encryption is where any function applied on the ciphertext when decrypted would give the same results as the function applied on the plaintext value.

Where only a few types of functions are supported, the type of encryption is called partially homomorphic encryption.

In this example we explore the partially homomorphic encryption cryptosystem called paillier.

```
In [62]: from phe import paillier
import pprint
import numpy as np

pp = pprint.PrettyPrinter(indent=4)

# generate key pair for encryption and decryption
public_key, private_key = paillier.generate_paillier_keypair()

# perform encryption
secret_number_list = [3.141592653, 300, -4.6e-12]
encrypted_number_list = [public_key.encrypt(x) for x in secret_number_list]

# perform functions
encrypted_number_list_multiply = [x * 5 for x in encrypted_number_list]
encrypted_number_list_addition = [x + 5 for x in encrypted_number_list]
encrypted_number_list_mean = np.mean(encrypted_number_list)

# show results
print('Original data:\n {}'.format(secret_number_list))

print('\nMultiplication example (x * 5):')
pp.pprint(encrypted_number_list_multiply)
print('\nDecrypts to ')
pp.pprint([private_key.decrypt(x) for x in encrypted_number_list_multiply])

print('\nAddition example (x + 5):')
pp.pprint(encrypted_number_list_addition)
print('\nDecrypts to ')
pp.pprint([private_key.decrypt(x) for x in encrypted_number_list_addition])

print('\nMean example:\n{}\n\ndecrypts to\n{}'.format(
    encrypted_number_list_mean, private_key.decrypt(encrypted_number_list_mean)))

Original data:
[3.141592653, 300, -4.6e-12]

Multiplication example (x * 5):
[ <phe.paillier.EncryptedNumber object at 0x7f80f293a670>,
  <phe.paillier.EncryptedNumber object at 0x7f80f293a460>,
  <phe.paillier.EncryptedNumber object at 0x7f80f293a6a0>]

decrypts to
[15.707963265, 1500, -2.2999999999999998e-11]

Addition example (x + 5):
[ <phe.paillier.EncryptedNumber object at 0x7f80f293a790>,
  <phe.paillier.EncryptedNumber object at 0x7f80f293a8b0>,
  <phe.paillier.EncryptedNumber object at 0x7f80f293a850>]

decrypts to
[8.141592653, 305, 4.99999999999954]

Mean example:
<phe.paillier.EncryptedNumber object at 0x7f80f293a970>

decrypts to
101.04719755099846
```

In [] :