

In [11]:

```

1 import numpy as np
2 import pandas as pd
3 import seaborn as sns
4 from keras.models import Sequential
5 from keras.layers import Dense, Flatten, Conv2D, Conv2DTranspose
6 from keras.layers import BatchNormalization, Reshape, LeakyReLU, Dropout
7 from keras.optimizers import Adam, SGD, RMSprop
8 from sklearn.decomposition import PCA
9 from sklearn.utils import shuffle
10 from PIL import Image
11 from time import time
12 import matplotlib.pyplot as plt
13 %matplotlib inline

```

In [12]:

```

1 X_train = pd.read_csv("../data/ClustREFGenes-master/Data/Core_genome/Data_Core_Genome_Ecoli_log2.c
2                      index_col=0)
3 print("Dimensionalidade dos dados: ", X_train.shape)
4 X_train.head()

```

Dimensionalidade dos dados: (4051, 9)

Out[12]:

	BB9	BB10	BB17	BB19	BB20	BB21	BB11	BB12	BB18
<b>Genes</b>									
<b>accD</b>	6.875411	7.047582	7.431765	7.105877	6.516094	6.676126	6.304694	6.168221	6.245553
<b>aceF</b>	7.732412	7.674997	8.397717	7.455056	7.277269	6.525536	7.455730	6.403830	7.597941
<b>ackA</b>	7.231720	7.260976	8.033280	6.921924	6.920829	6.556644	6.358150	5.888768	6.359310
<b>agaV</b>	6.048825	6.250033	5.120269	5.559767	5.915593	6.279490	6.441998	6.553099	6.105364
<b>alaS</b>	7.811728	7.853890	8.622037	7.636451	7.641365	7.125920	7.164957	6.555678	7.098590

**PCA:**

- para ver a distribuição dos dados, estes serão reduzidos à só dois dimensões com PCA.

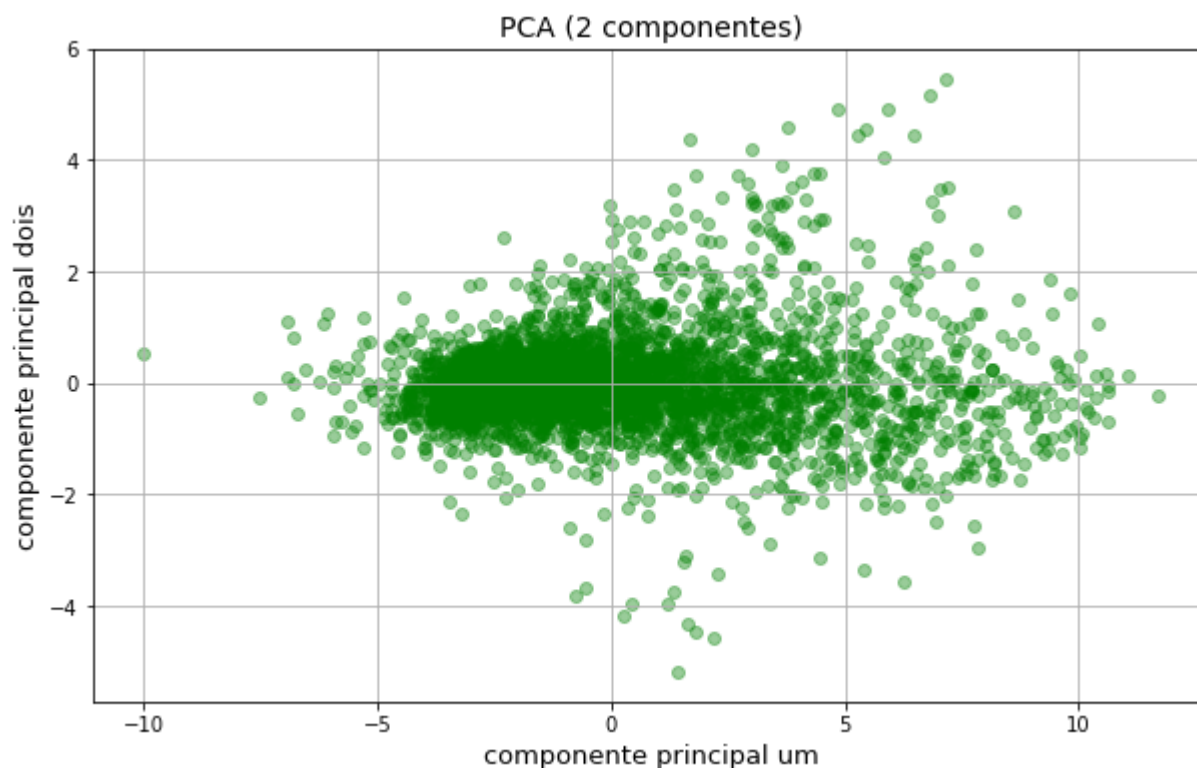
```
In [13]: 1 pca = PCA(n_components=2)
         2 pca.fit(X_train)
```

```
Out[13]: PCA(copy=True, iterated_power='auto', n_components=2, random_state=None,
          svd_solver='auto', tol=0.0, whiten=False)
```

```
In [14]: 1 X_pca = pca.transform(X_train)
         2 print("Dimensionalidade: ", X_pca.shape)
```

Dimensionalidade: (4051, 2)

```
In [15]: 1 plt.figure(figsize=(10,6))
         2 plt.title("PCA (2 componentes)", fontsize=14)
         3 plt.xlabel("componente principal um", fontsize=13)
         4 plt.ylabel("componente principal dois", fontsize=13)
         5 plt.grid()
         6 plt.scatter(X_pca[:,0], X_pca[:,1], color="green", alpha=.4);
```



## Normalização dos dados

- Para poder fazer uso das GAN's, a gente tem que normalizar os dados, para eso é usada a seguinte normalização:

$$X = \frac{x_i - \text{mean}(X)}{\text{mean}(X)}$$

```
In [16]: 1 X_train = (X_train - np.mean(X_train))/np.mean(X_train)
```

- vou trocar a dimensionalidade de cada gen, por uma dimesionalidade de 3x3, para assim ver o gen como uma matrix.

```
In [17]: 1 X_img = np.array(X_train).reshape((X_train.shape[0],3,3,1))  
2 X_lab = X_train.index
```

```
In [18]: 1 idx = np.random.randint(low=0, high=X_train.shape[0], size=24)
2 imgs = X_img[idx]
3 titles = X_lab[idx]
4 fig = plt.figure(figsize=(15,6))
5 p=0
6 #plt.title("Genomas representados na forma de uma imagem", fontsize=12)
7 plt.axis("off");
8 print("----- Genomas representados na forma de uma matriz -----")
9 for i in imgs:
10     ax=fig.add_subplot(3,8,p+1)
11     plt.title(titles[p])
12     plt.imshow(i.reshape(3,3))
13     plt.axis("off");
14     p += 1
```

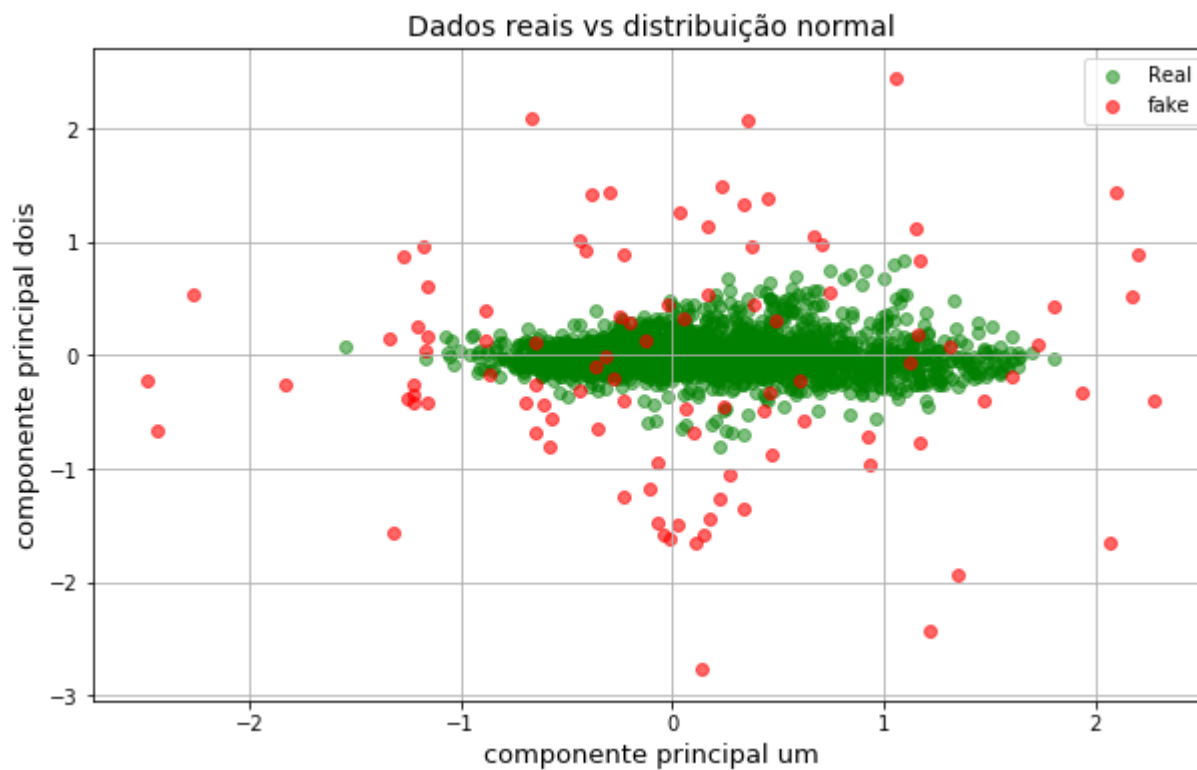
----- Genomas representados na forma de uma matriz -----  
 ----



## Dados fake

- Os dados que vão ser passados pro gerador, são dados de uma distribuição normal.
- O gráfico abaixo mostra os dados reais e os dados da distribuição normal os quais são para treinar a rede geradora.

```
In [19]: 1 pca2 = PCA(n_components=2)
2         pca2.fit(X_train)
3         X_real = pca2.transform(X_train)
4         X_fake = pca2.transform(np.random.normal(0,1,(100,9)))
5         plt.figure(figsize=(10,6))
6         plt.grid()
7         plt.title("Dados reais vs distribuição normal", fontsize=14)
8         plt.xlabel("componente principal um", fontsize=13)
9         plt.ylabel("componente principal dois", fontsize=13)
10        plt.scatter(X_real[:,0], X_real[:,1], label="Real", alpha=.5, color="green")
11        plt.scatter(X_fake[:,0], X_fake[:,1], label="fake", alpha=.6, color="red")
12        plt.legend();
```

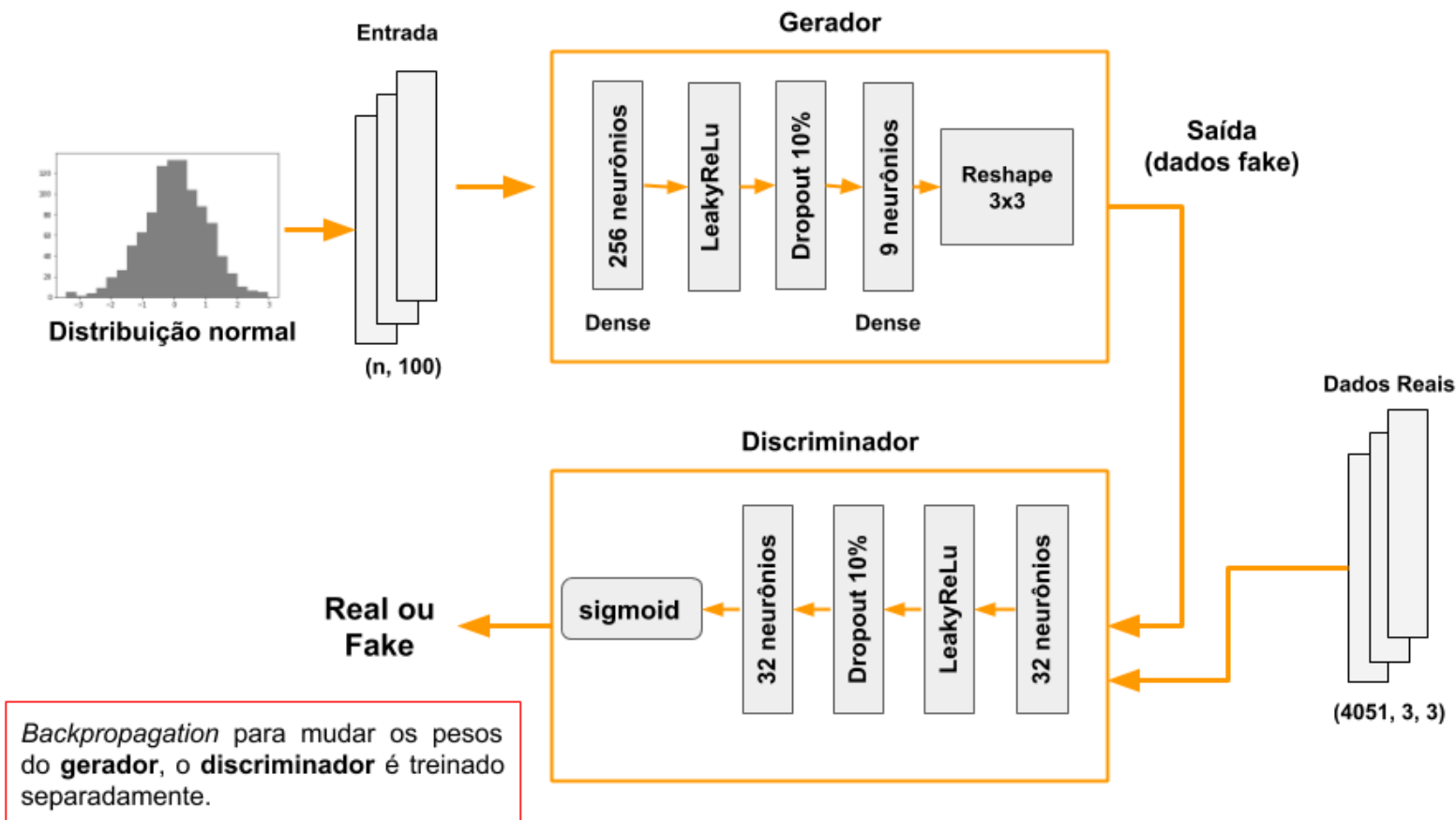


## Construção da GAN

- A GAN conta com duas redes, uma rede geradora e uma discriminadora.

```
In [20]: 1 Image.open("../images/image1.png")
```

```
Out[20]:
```



In [21]:

```

1  class GANs():
2      #inialização dos parâmetros
3      def __init__(self, width, height, channels, noise_input):
4          self.width = width
5          self.height = height
6          self.channels = channels
7          self.dim = (self.width, self.height, self.channels)
8          self.noise_input = noise_input
9          self.g_loss = []
10         self.d_loss = []
11         self.g_lpe = []
12         self.d_lpe = []
13         self.optimizerD = Adam(lr=0.0001, beta_1=0.5)
14         self.optimizerG = Adam(lr=0.0004, beta_1=0.5)
15         self.G = self.noise_generator()
16         print("Compilando o gerador...")
17         self.G.compile(loss='binary_crossentropy', optimizer=self.optimizerG)
18         self.D = self.discriminator()
19         print("Compilando o discriminador...")
20         self.D.compile(loss='binary_crossentropy', optimizer=self.optimizerD, metrics=['accuracy'])
21         self.stacked_generator_discriminator = Sequential()
22         self.stacked_generator_discriminator.add(self.G)
23         self.stacked_generator_discriminator.add(self.D)
24         self.D.trainable = False
25         self.stacked_generator_discriminator.compile(loss='binary_crossentropy', optimizer=self.
26
27     #criação do gerador de imagens fake
28     def noise_generator(self):
29         model = Sequential()
30         model.add(Dense(256, input_shape=(self.noise_input,)))
31         model.add(LeakyReLU(alpha=0.3))
32         model.add(Dropout(.1))
33         model.add(Dense(self.width*self.height*self.channels, activation="tanh"))
34         model.add(Reshape((self.width, self.height, self.channels)))
35         return model
36
37     #criação do discriminador
38     def discriminator(self):
39         model = Sequential()
40         model.add(Dense(32, input_shape=self.dim))
41         model.add(LeakyReLU(alpha=0.2)) #función rectificadora
42         model.add(Dropout(.1))

```

```
43     model.add(Dense(32))
44     model.add(Flatten())
45     model.add(Dense(1, activation='sigmoid'))
46
47     return model
48
49     #Para obter o sumary do gerador
50     def summary_gerador(self):
51         return self.G.summary()
52
53     #Para obter o sumary do gerador
54     def summary_discriminador(self):
55         return self.D.summary()
56
57     #pra obter os batches pra o treino
58     def get_batches(self, X_train, batch_size):
59         """
60         X_train: dataset para o treino
61         epochs: quantidade de epocas para o treino do gradiente
62         batch: tamanho to batch pra o treino de cada epochs
63         """
64         batches = []
65         num_bat = int(np.ceil(X_train.shape[0]/batch_size))
66         lim_i = 0
67         lim_s = batch_size
68         for i in range(num_bat):
69             if lim_s > X_train.shape[0]:
70                 lim_s = X_train.shape[0]
71             batches.append(X_train[lim_i:lim_s])
72             lim_i += batch_size
73             lim_s += batch_size
74
75         return batches
76
77     #devolve o loss do gerador e do discriminador
78     def get_loss(self):
79         return [self.g_loss, self.d_loss]
80
81     #treinamento da GAN
82     def train(self, X_train, epochs, batch_size):
83         self.d_loss = []
84         self.g_loss = []
85         self.g_lpe = []
```



```

86     self.d_lpe = []
87     for cnt in range(epochs):
88         batches = self.get_batches(X_train, batch_size)
89         count_b = 0
90         t_i = time()
91         for batch in batches:
92             gen_noise = np.random.normal(0, 1, (np.int64(batch.shape[0]), self.noise_input))
93             #gerando as imagens fake
94             syntetic_images = self.G.predict(gen_noise)
95             #criação do array de treinamento
96             x_combined_batch = np.concatenate((batch, syntetic_images))
97             y_combined_batch = np.concatenate((np.ones((batch.shape[0], 1)),
98                                                np.zeros((batch.shape[0], 1))))
99             #misturar os dados
100            #x_combined_batch, y_combined_batch = shuffle(x_combined_batch, y_combined_batch)
101            #treino do discriminador
102            d_l = self.D.train_on_batch(x_combined_batch, y_combined_batch)
103            self.d_loss.append(d_l[0])
104            # train generator
105            noise = np.random.normal(0, 1, (batch.shape[0], self.noise_input))
106            y_mislabeled = np.ones((batch.shape[0], 1))
107
108            g_l = self.stacked_generator_discriminator.train_on_batch(noise, y_mislabeled)
109            self.g_loss.append(g_l)
110            count_b += 1
111            if (count_b%len(batches))==0:
112                t_f = time()
113                t = t_f - t_i
114                t_i = time()
115                print ('epoch:[%d/%d] batch:[%d/%d], [Discriminator::d_loss: %f], [Generator
116                        % (cnt+1,epochs,count_b,len(batches),d_l[0],g_l,t))
117            self.g_lpe.append(g_l)
118            self.d_lpe.append(d_l[0])

```

```

In [22]: 1 noise_input = 100
         2 gan = GANs(width=3, height=3, channels=1, noise_input=noise_input)

```

Compilando o gerador...

Compilando o discriminador...

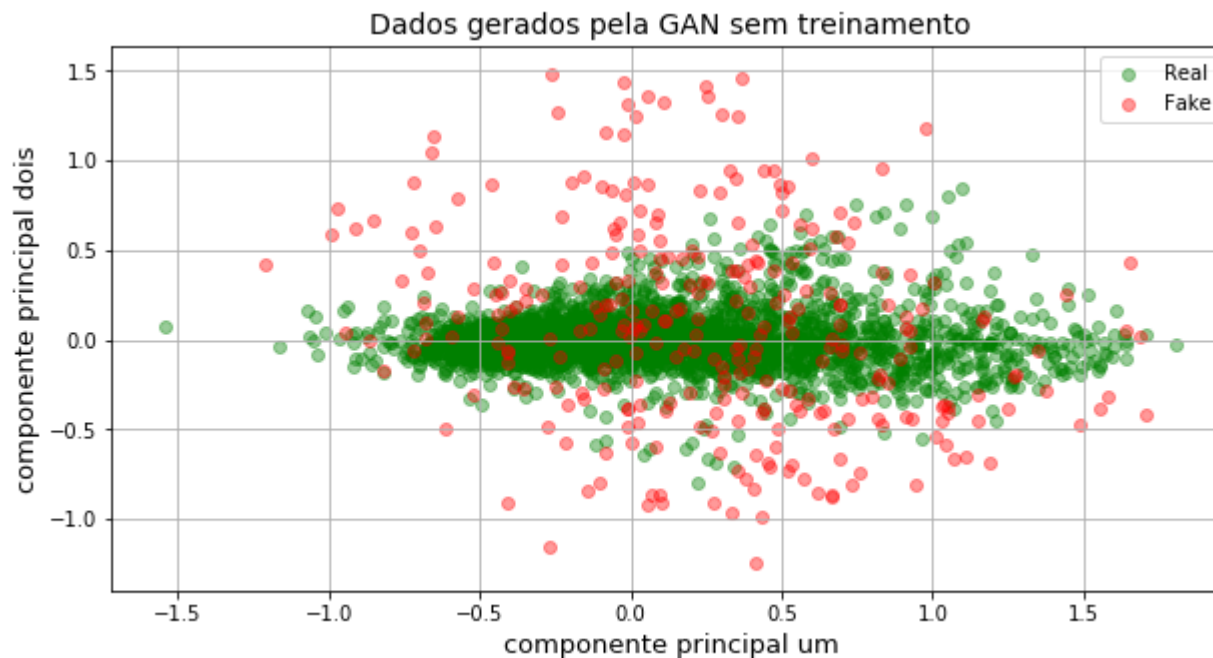
```

In [23]: 1 num_imgs = 300
2 fakes = gan.G.predict(np.random.normal(0,1,(num_imgs,noise_input)))
3 print("imagens fake: ", fakes.shape)
4 fakes = fakes.reshape(num_imgs,9)
5 print("re-dimesionalidade: ", fakes.shape)
6 pca2 = PCA(n_components=2)
7 pca2.fit(X_img.reshape(X_img.shape[0],9))
8
9 X_real = pca2.transform(X_img.reshape(X_img.shape[0],9))
10 X_fake = pca2.transform(fakes)
11 plt.figure(figsize=(10,5))
12 plt.title("Dados gerados pela GAN sem treinamento", fontsize=14)
13 plt.xlabel("componente principal um", fontsize=13)
14 plt.ylabel("componente principal dois", fontsize=13)
15 plt.grid()
16 plt.scatter(X_real[:,0], X_real[:,1], color="green", alpha=.4, label="Real")
17 plt.scatter(X_fake[:,0], X_fake[:,1], color="red", alpha=.4, label="Fake")
18 plt.legend();

```

imagens fake: (300, 3, 3, 1)

re-dimesionalidade: (300, 9)



```
In [24]: 1 print("-----Estrutura da rede generativa-----")
        2 gan.G.summary()
```

```
-----Estrutura da rede generativa-----
```

Layer (type)	Output Shape	Param #
=====		
dense_1 (Dense)	(None, 256)	25856
leaky_re_lu_1 (LeakyReLU)	(None, 256)	0
dropout_1 (Dropout)	(None, 256)	0
dense_2 (Dense)	(None, 9)	2313
reshape_1 (Reshape)	(None, 3, 3, 1)	0
=====		
Total params: 28,169		
Trainable params: 28,169		
Non-trainable params: 0		

```
In [25]: 1 print("-----Estrutura da rede Discriminadora-----")
        2 gan.D.summary()
```

-----Estrutura da rede Discriminadora-----

Layer (type)	Output Shape	Param #
=====		
dense_3 (Dense)	(None, 3, 3, 32)	64
leaky_re_lu_2 (LeakyReLU)	(None, 3, 3, 32)	0
dropout_2 (Dropout)	(None, 3, 3, 32)	0
dense_4 (Dense)	(None, 3, 3, 32)	1056
flatten_1 (Flatten)	(None, 288)	0
dense_5 (Dense)	(None, 1)	289
=====		
Total params: 2,818		
Trainable params: 1,409		
Non-trainable params: 1,409		

/home/ejrueada/anaconda3/lib/python3.6/site-packages/keras/engine/training.py:490: UserWarning: Discrepancy between trainable weights and collected trainable weights, did you set `model.trainable` with out calling `model.compile` after ?  
'Discrepancy between trainable weights and collected trainable'

```
In [26]: 1 import warnings
        2 warnings.filterwarnings('ignore')
```

```
In [ ]: 1 t_i = time()
        2 gan.train(X_img, epochs=400, batch_size=32)
        3 t_f = time()
```

```
In [28]: 1 print("tempo de execução: ", (t_f-t_i)/60, "[min]")
```

tempo de execução: 2.7460547208786013 [min]

```

In [29]: 1 g_loss, d_loss = gan.get_loss()
          2
          3 plt.figure(figsize=(15,6))
          4 #plt.plot(range(len(g_loss)), g_loss)
          5 plt.title("Loss GAN", fontsize=14)
          6 plt.ylabel("loss", fontsize=13.5)
          7 plt.xlabel("batch/epoch", fontsize=13.5)
          8 plt.plot(range(np.array(d_loss).shape[0]), np.array(d_loss), label="loss discriminador", alpha=.8)
          9 plt.plot(range(np.array(g_loss).shape[0]), np.array(g_loss), label="loss gerador", alpha=.8)
         10 #plt.plot(range(np.array(g_loss).shape[0]), 0.5*np.ones(np.array(g_loss).shape[0]),
         11 #         color="black", label="objetivo", linestyle='--')
         12 plt.grid()
         13 #plt.yticks([0, 0.5,1,1.5,2,2.5,3,3.5,4,4.5,5])
         14 plt.legend();

```

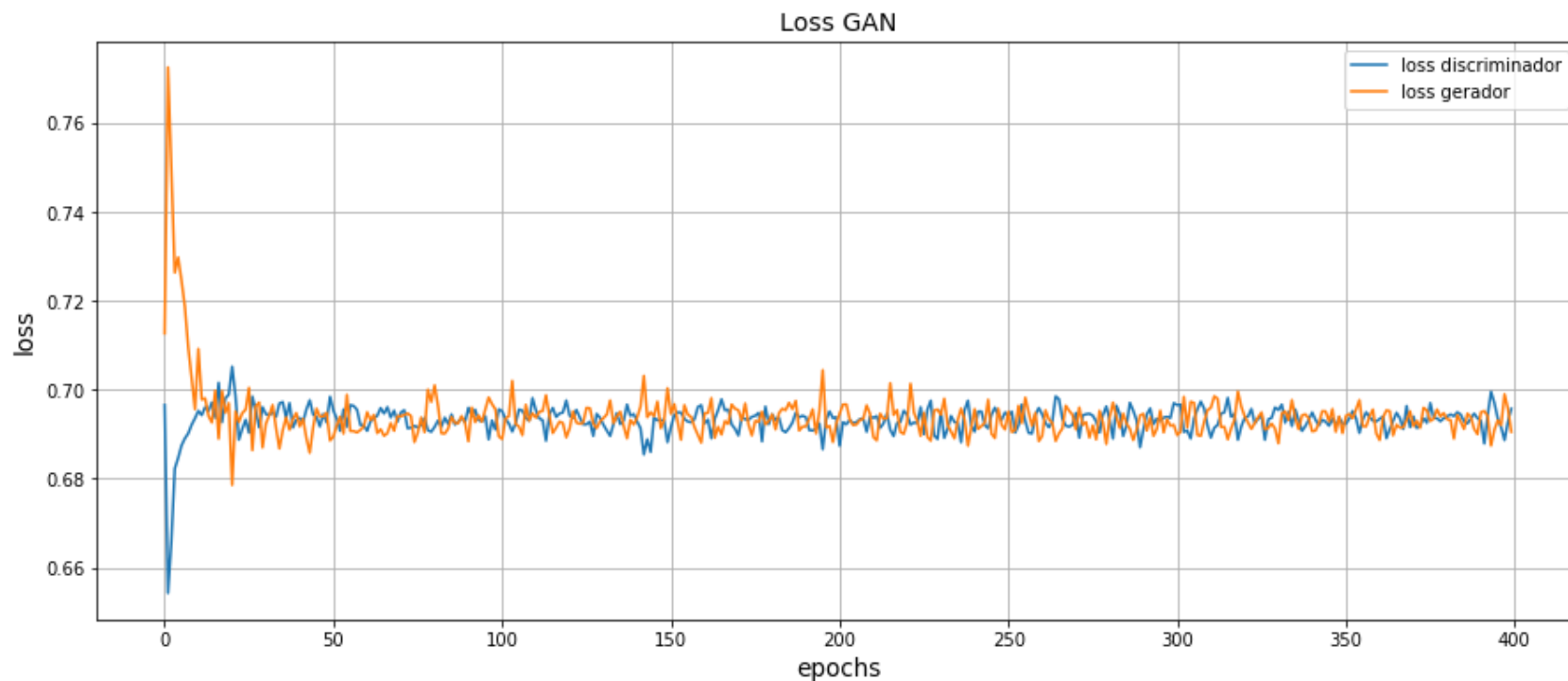


```

In [ ]: 1

```

```
In [30]: 1 plt.figure(figsize=(15,6))
2 plt.title("Loss GAN", fontsize=14)
3 plt.ylabel("loss", fontsize=13.5)
4 plt.xlabel("epochs", fontsize=13.5)
5 plt.plot(range(len(gan.d_lpe)), np.array(gan.d_lpe), label="loss discriminador")
6 plt.plot(range(len(gan.g_lpe)), np.array(gan.g_lpe), label="loss gerador")
7 #plt.plot(range(len(gan.g_lpe)), 0.5*np.ones(len(gan.g_lpe)),
8 #         color="black", label="objetivo", linestyle='--')
9 plt.grid()
10 plt.legend();
```



```
In [ ]: 1
```

```
In [31]: 1 num_imgs = 24 #número de imágenes a mostrar aleatoriamente
2 img_pre = gan.G.predict(np.random.normal(0,1,(num_imgs, noise_input)))
3 fig = plt.figure(figsize=(15,6))
4 for i in range(num_imgs):
5     ax=fig.add_subplot(3,8,i+1)
6     img = img_pre[i]
7     plt.imshow(img.reshape((3,3)))
8     ax.axis("off")
9 plt.show()
```



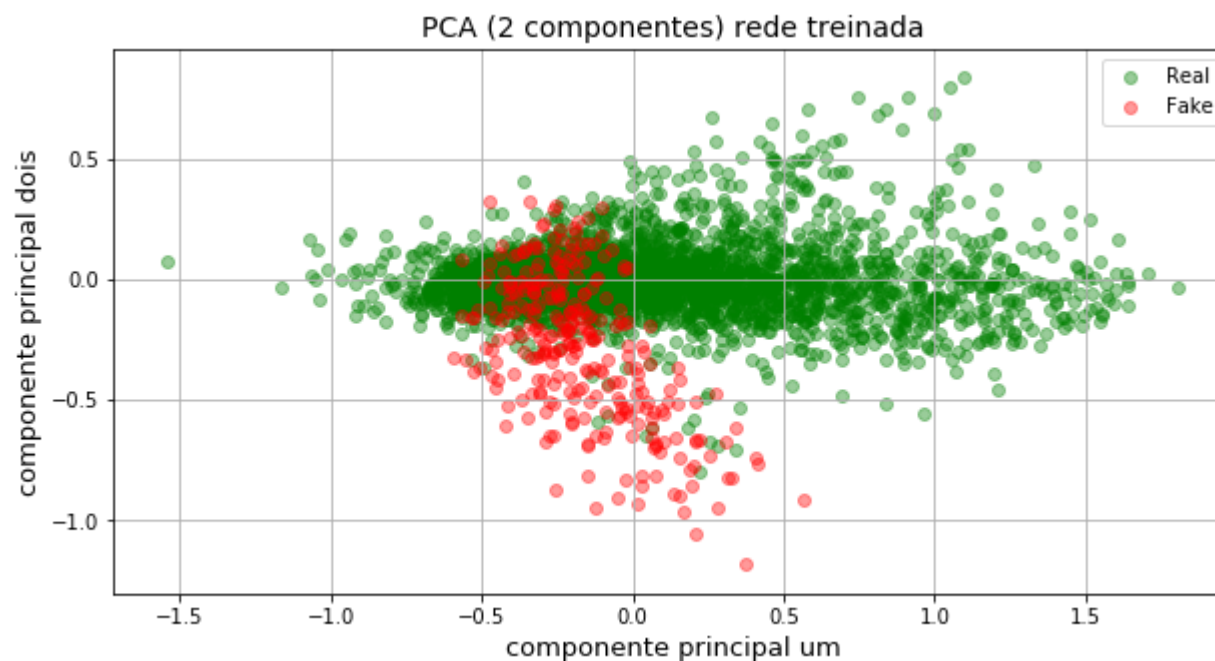
```
In [32]: 1 num_imgs = 300
2 fakes = gan.G.predict(np.random.normal(0,1,(num_imgs,noise_input)))
3 print("imagens fake: ", fakes.shape)
4 fakes = fakes.reshape(num_imgs,9)
5 print("re-dimesionalidade: ", fakes.shape)
6 pca2 = PCA(n_components=2)
7 pca2.fit(X_img.reshape(X_img.shape[0],9))
8
9 X_real = pca2.transform(X_img.reshape(X_img.shape[0],9))
10 X_fake = pca2.transform(fakes)
```

imagens fake: (300, 3, 3, 1)

re-dimesionalidade: (300, 9)



```
In [33]: 1 plt.figure(figsize=(10,5))
2 plt.title("PCA (2 componentes) rede treinada", fontsize=14)
3 plt.xlabel("componente principal um", fontsize=13)
4 plt.ylabel("componente principal dois", fontsize=13)
5 plt.grid()
6 plt.scatter(X_real[:,0], X_real[:,1], color="green", alpha=.4, label="Real")
7 plt.scatter(X_fake[:,0], X_fake[:,1], color="red", alpha=.4, label="Fake")
8 plt.legend();
```



```
In [ ]: 1
```

```
In [34]: 1 num_imgs = 300
2 fakes = gan.G.predict(np.random.normal(0,1,(num_imgs,noise_input)))
3 print("Dados fake: ", fakes.shape)
4 y_predict = gan.D.predict_classes(fakes)
5 print("----- Discriminador -----")
6 print("porcentagem de dados fake que o discriminador acredita reais: ", np.mean(y_predict==1))
7 print("porcentagem de dados fake que o discriminador acredita fakes: ", np.mean(y_predict==0))
8 print()
9 y_predict2 = gan.D.predict_classes(X_img)
10 print("----- Discriminador com dados reais -----")
11 print("porcentagem de dados reais que o discriminador acredita reais: ", np.mean(y_predict2==1))
12 print("porcentagem de dados reais que o discriminador acredita fakes: ", np.mean(y_predict2==0))
```

Dados fake: (300, 3, 3, 1)

----- Discriminador -----

porcentagem de dados fake que o discriminador acredita reais: 0.74

porcentagem de dados fake que o discriminador acredita fakes: 0.26

----- Discriminador com dados reais -----

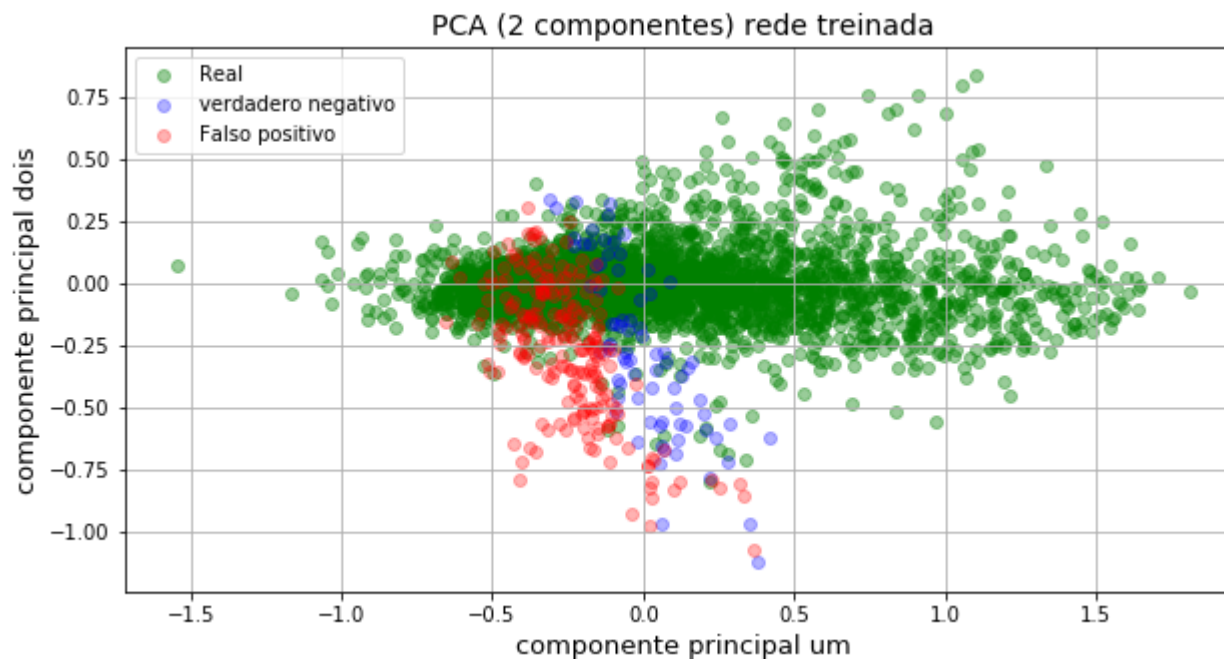
porcentagem de dados reais que o discriminador acredita reais: 0.44235991113305356

porcentagem de dados reais que o discriminador acredita fakes: 0.5576400888669464

```

In [35]: 1 fakes = fakes.reshape(num_imgs,9)
          2
          3 pca2 = PCA(n_components=2)
          4 pca2.fit(X_img.reshape(X_img.shape[0],9))
          5
          6 X_real = pca2.transform(X_img.reshape(X_img.shape[0],9))
          7 X_fake = pca2.transform(fakes)
          8 X_fp = X_fake[np.where(y_predict==1)[0]]
          9 X_vn = X_fake[np.where(y_predict==0)[0]]
         10
         11 plt.figure(figsize=(10,5))
         12 plt.title("PCA (2 componentes) rede treinada", fontsize=14)
         13 plt.xlabel("componente principal um", fontsize=13)
         14 plt.ylabel("componente principal dois", fontsize=13)
         15 plt.grid()
         16 plt.scatter(X_real[:,0], X_real[:,1], color="green", alpha=.4, label="Real")
         17 plt.scatter(X_vn[:,0], X_vn[:,1], color="blue", alpha=.3, label="verdadero negativo")
         18 plt.scatter(X_fp[:,0], X_fp[:,1], color="red", alpha=.3, label="Falso positivo")
         19 plt.legend();

```



```
In [36]: 1 bp_fp = []
          2 bp_vn = []
          3 for i in range(300):
          4     fakes = gan.G.predict(np.random.normal(0,1,(num_imgs,noise_input)))
          5     y_predict = gan.D.predict_classes(fakes)
          6     bp_fp.append(np.mean(y_predict==1))
          7     bp_vn.append(np.mean(y_predict==0))
```

In [37]:

```

1 from bokeh.plotting import figure, show, output_file, output_notebook
2 from bokeh.layouts import row
3 from bokeh.models import ColumnDataSource
4
5 output_notebook()
6 x_ticks = ["falsos_positivos"]
7 p1 = figure(tools="", background_fill_color="#efefef", toolbar_location=None, x_range=x_ticks,
8             title="Falsos positivos do discriminador", width=450, height=500)
9 q1 = np.quantile(bp_fp, q=0.25)
10 q2 = np.quantile(bp_fp, q=0.5)
11 q3 = np.quantile(bp_fp, q=0.75)
12 iqr = q3 - q1
13 upper = q3 + 1.5*iqr
14 lower = q1 - 1.5*iqr
15 outliers = np.array(bp_fp)[(bp_fp>upper) + (bp_fp<lower)]
16 source1 = ColumnDataSource(dict(x=x_ticks, upper=[upper], lower=[lower], q1=[q1], q2=[q2], q3=[q3]
17 source2 = ColumnDataSource(dict(x=x_ticks*len(outliers), y=outliers))
18 #Para graficar las lineas superiores del boxplot
19 p1.rect("x", "upper", 0.2, 0.0003, line_color="black", fill_color="black", source=source1)
20 #Para graficar las lineas inferiores del boxplot
21 p1.rect("x", "lower", 0.2, 0.0003, line_color="black", fill_color="black", source=source1)
22 #Para graficar los segmentos del boxplot
23 p1.segment("x", "lower", "x", "q1", line_color="black", source=source1)
24 p1.segment("x", "upper", "x", "q3", line_color="black", source=source1)
25 #Para graficar las barras
26 p1.vbar("x", 0.3, "q1", "q2", fill_color="#3B8686", line_color="black", source=source1, legend="c
27 p1.vbar("x", 0.3, "q2", "q3", fill_color="#E08E79", line_color="black", source=source1, legend="c
28 p1.circle("x", "y", size=8, color="#F38630", fill_alpha=0.6, source=source2)
29
30 x_ticks2 = ["verdaderos negativos"]
31 p2 = figure(tools="", background_fill_color="#efefef", toolbar_location=None, x_range=x_ticks2,
32             title="Verdaderos negativos do discriminador", width=450, height=500)
33 p2_q1 = np.quantile(bp_vn, q=0.25)
34 p2_q2 = np.quantile(bp_vn, q=0.5)
35 p2_q3 = np.quantile(bp_vn, q=0.75)
36 p2_iqr = p2_q3 - p2_q1
37 p2_upper = p2_q3 + 1.5*iqr
38 p2_lower = p2_q1 - 1.5*iqr
39 outliers2 = np.array(bp_vn)[(bp_vn>p2_upper) + (bp_vn<p2_lower)]
40 source3 = ColumnDataSource(dict(x=x_ticks2, upper=[p2_upper], lower=[p2_lower],
41                                 q1=[p2_q1], q2=[p2_q2], q3=[p2_q3]))
42 source4 = ColumnDataSource(dict(x=x_ticks2*len(outliers2), y=outliers2))

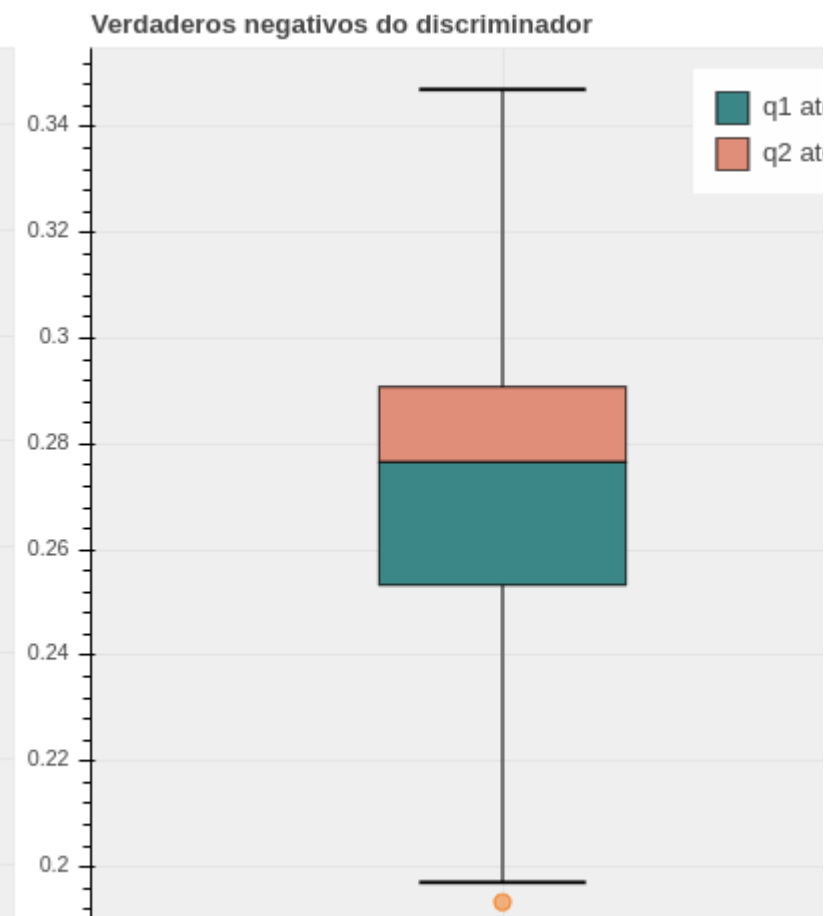
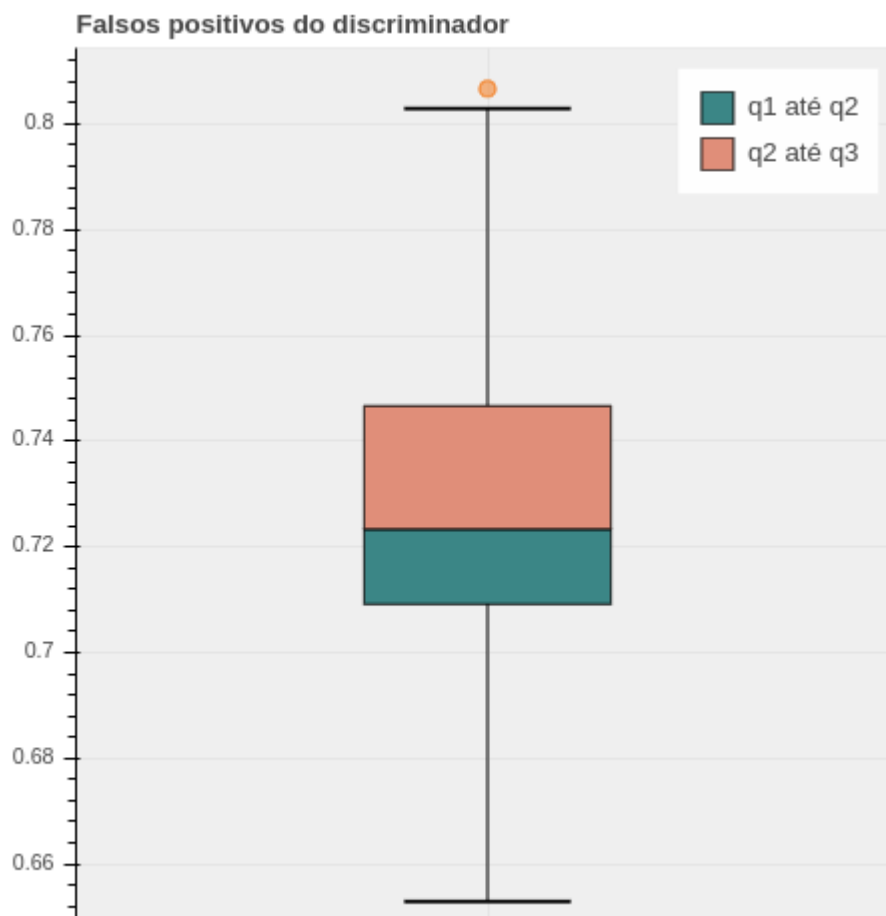
```

```

43 #Para graficar las lineas superiores del boxplot
44 p2.rect("x", "upper", 0.2, 0.0003, line_color="black", fill_color="black", source=source3)
45 #Para graficar las lineas inferiores del boxplot
46 p2.rect("x", "lower", 0.2, 0.0003, line_color="black", fill_color="black", source=source3)
47 #Para graficar los segmentos del boxplot
48 p2.segment("x", "lower", "x", "q1", line_color="black", source=source3)
49 p2.segment("x", "upper", "x", "q3", line_color="black", source=source3)
50 #Para graficar las barras
51 p2.vbar("x", 0.3, "q1", "q2", fill_color="#3B8686", line_color="black", source=source3, legend="c
52 p2.vbar("x", 0.3, "q2", "q3", fill_color="#E08E79", line_color="black", source=source3, legend="c
53 p2.circle("x", "y", size=8, color="#F38630", fill_alpha=0.6, source=source4)
54
55 show(row([p1,p2]))

```

(<http://book5.pydata.org/>) successfully loaded.





In [ ]: 1

In [ ]: 1

In [ ]: 1

In [ ]: 1

In [ ]: 1

In [ ]: 1

In [187]: 1 `#gan.D.save_weights("D_weights.h5")`  
2 `#gan.G.save_weights("G_weights.h5")`  
3 `#gan.stacked_generator_discriminator.save_weights("stacked_weights.h5")`

In [ ]: 1