

LE MACHINE-LEARNING AVEC PYTHON

Etude de cas :
regroupement d'images
par similarités

Corentin Martel

Data Scientist Freelance

<https://www.corelyo.com>



CONTEXTE DU PROJET

Création d'une plateforme d'analyse de
positionnement marketing pour des marques de luxe

CONTEXTE DU PROJET

Objectif : mettre en évidence les différences de stratégies sur Instagram

CONTEXTE DU PROJET

Objectif : pouvoir afficher les images de plusieurs marques sur un même écran, en les regroupant par similarité

Problématique : comment comparer des images de manière non-supervisée?

LES DONNÉES

COLLECTER DES DONNÉES

Données disponible : liste d'url d'images Instagram

```
10 from PIL import Image
11 import requests
12 from io import BytesIO
13 import numpy as np
89 def download_from_url(url):
90     img_path = []
91     ret=[]
92     first = True
93     try :
94         if len(url.split('/')[1]) <=75 :
95             response = requests.get(url)
96             img = Image.open(BytesIO(response.content))
97             img = img.resize((224,224), Image.ANTIALIAS)
98             x = np.asarray(img)
99             ret = np.expand_dims(x, axis=0)
100             img_path = url.split('/')[1]
101             return img_path, ret
102     except Exception as e:
103         print(e)
```

```
1 import pandas as pd
2 from multiprocessing.pool import ThreadPool
```

```
1 # download pixels and shape it
2 %%time
3
4 df_url = pd.read_excel('./data/IG_Posts_Brands_prepared.xlsx')
5 url_list = [_ for _ in df_url.thumbnail[:500] if len(_)<=75]
6
7 if __name__ == '__main__':
8     num_threads = 7
9     p = ThreadPool(num_threads)
10    xs = p.map(download_from_url, url_list)
11    xs = [elt for elt in xs if elt is not None]
12    img = [_[0] for _ in xs]
13    X = np.concatenate([_[1] for _ in xs],axis=0)
```

```
1 print("Nombre d'images : {}\nX.shape : {}".format(len(img), X.shape))
```

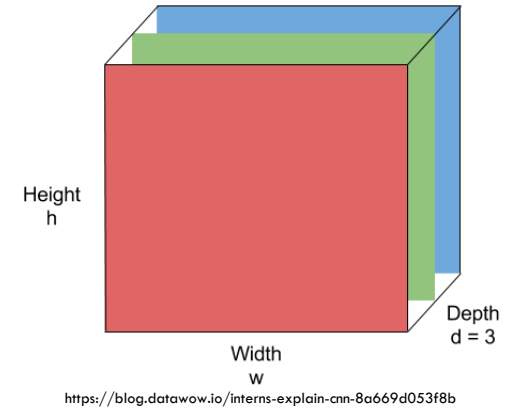
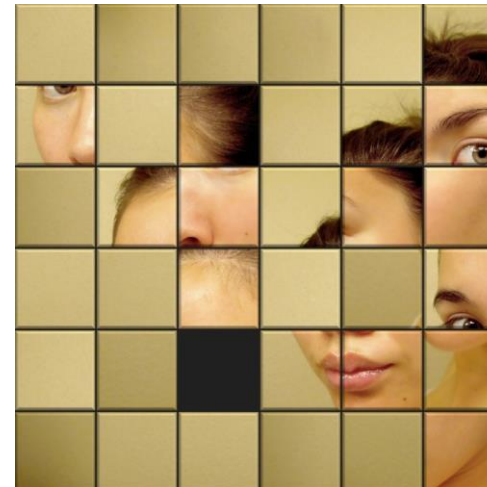
Nombre d'images : 10000
X.shape : (10000, 224, 224, 3)

COMMENT DÉCRIRE UNE IMAGE ?

Image = matrice de dimension $L * H * 3$

Utiliser les pixels ?

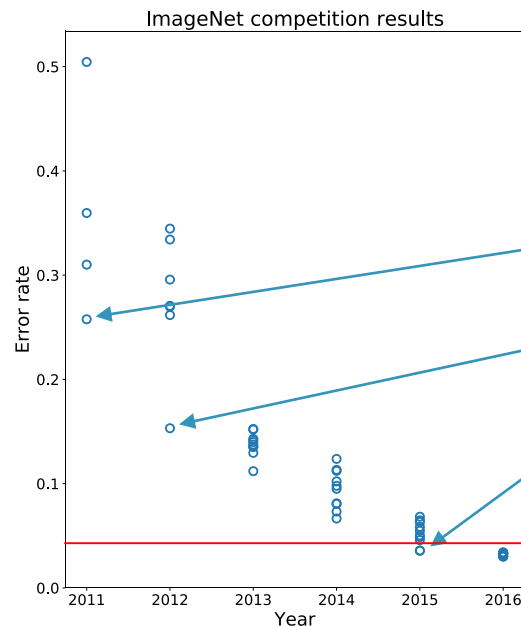
- Nombre de variables : 150 528 variables pour $224*224$ pixels
- Ces variables sont ordonnées : on ne peut pas les considérer comme un ensemble de variables indépendantes



UN PEU D'HISTOIRE : IMAGENET

Base de donnée de 14 millions d'images dans 21 000 catégories

Depuis 2010 : concours international : prédire la catégorie de l'image (1,2 millions d'images, 1000 catégories)



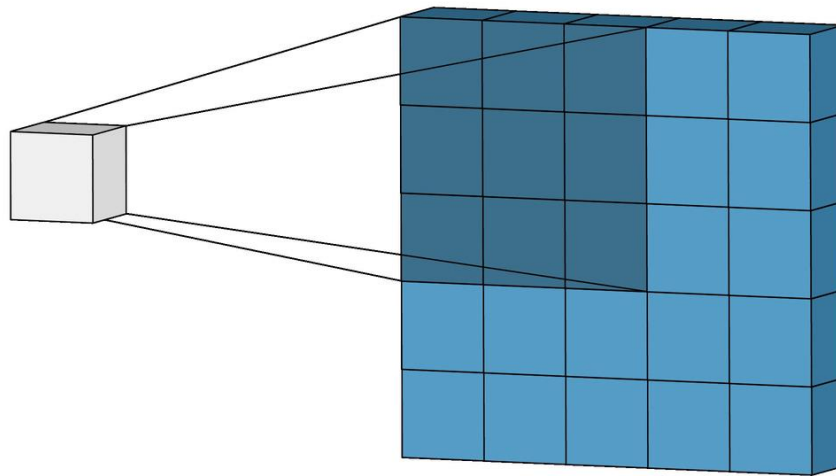
2011 : 25,8%, sans réseaux de neurones

2012 : AlexNet : 15,3 % (2eme place : 26,2%)

2015 : ResNet-152 : 3,6%

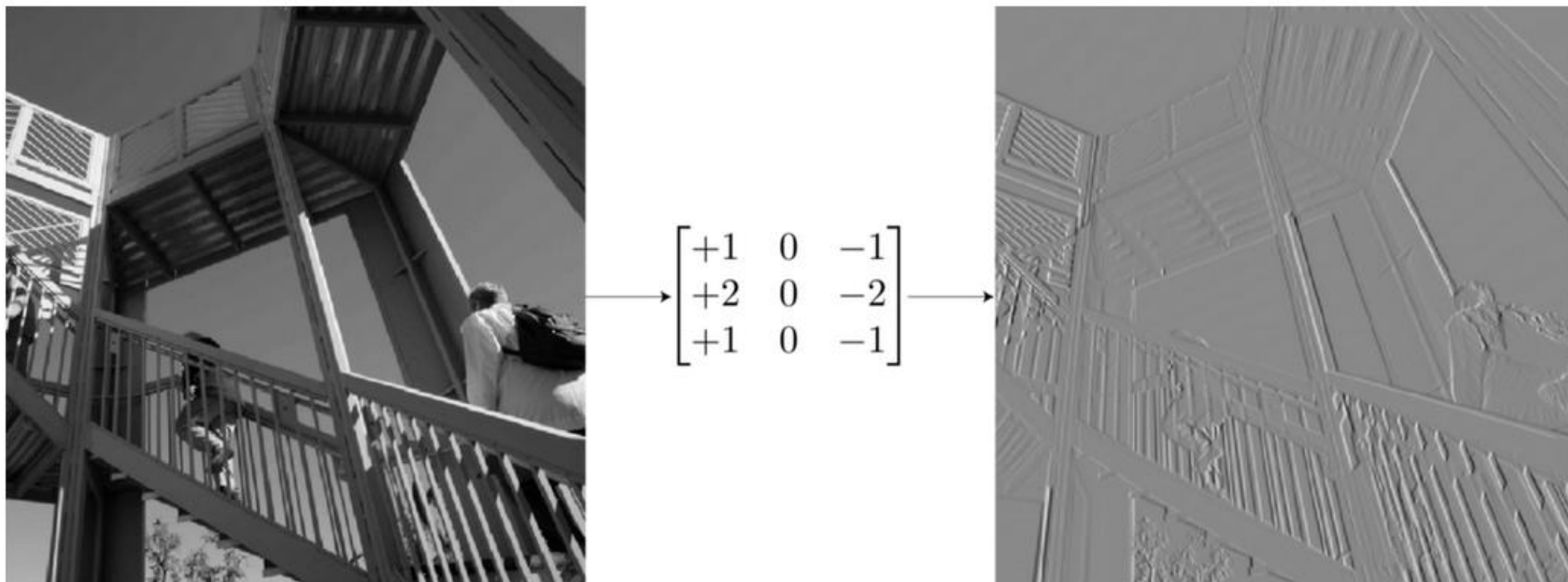
Humain (5,1%) : <https://cs.stanford.edu/people/karpathy/ilsvrc/>

LES CONVOLUTIONS



LES CONVOLUTIONS

Exemple de filtre de détection de lignes verticales :

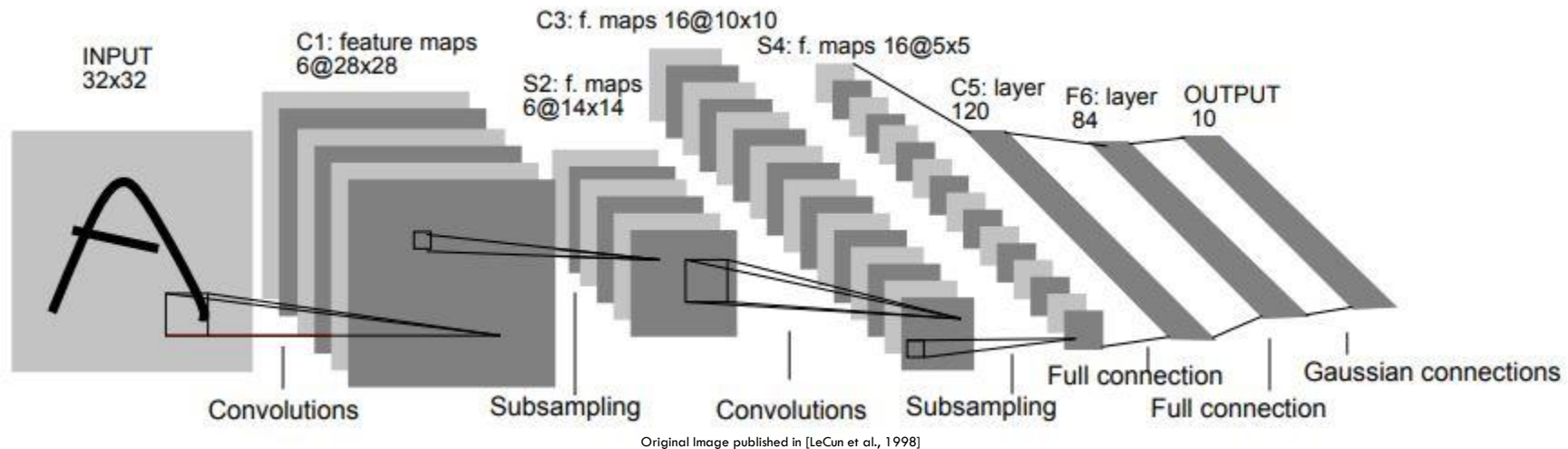


Applying a vertical edge detector kernel

<https://towardsdatascience.com/intuitively-understanding-convolutions-for-deep-learning-1f6f42faee1>

LES RÉSEAUX CONVOLUTIFS

1998 : LeNet-5 : reconnaissance de chiffres sur des chèques



DESCRIPTION DES IMAGES

CHOIX DU RÉSEAU CONVOLUTIF

ResNet50 : bons résultats

50 couches

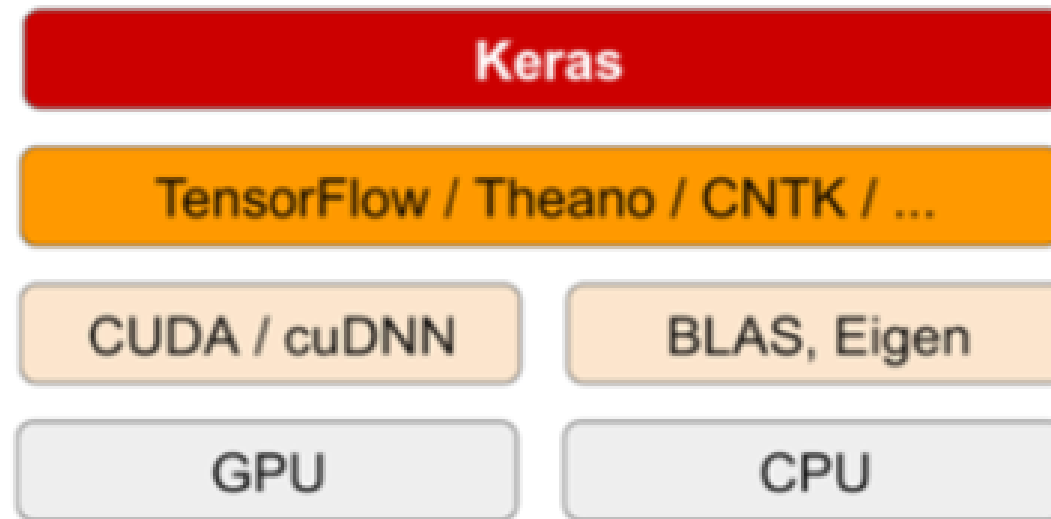
25,6 millions de paramètres (61k pour LeNet5)

Durée d'entraînement : 14 jours (Tesla P100 x8)



Utilisation d'un réseau pré-entraîné sur ImageNet pour extraire les variables

IMPLÉMENTATION : KERAS



<https://towardsdatascience.com/my-journey-into-deeplearning-using-keras-part-1-67cbb50f65e6>

KERAS

```
23 import keras
24 from keras.applications.imagenet_utils import preprocess_input
25 from keras.preprocessing import image
26 from keras.preprocessing.image import ImageDataGenerator
121 def extract_features(img,
122                     X,
123                     batch_size = 256,
124                     feature_save_file = './raw_features.h5'):
125     """
126     Extract features from an image using resnet and saves in HDF5 files.
127     img : list of img names
128     X : matrix of shape (Nb of images, 224, 224, 3)
129     batch size : for ResNet.
130     feature_save_file : HDF5 file where ResNet features of the images are stored
131     """
132
133     feat_extractor = keras.applications.resnet50.ResNet50(input_shape = (224,224,3),
134                                                         weights='imagenet',
135                                                         include_top=False,
136                                                         pooling='max')
137
138     datagen = ImageDataGenerator(preprocessing_function=preprocess_input)
139
140     generator = datagen.flow(
141         X,
142         batch_size=batch_size,
143         shuffle=False)
144
145     features = feat_extractor.predict_generator(generator)
146
147     print(features.shape)
148     df_features = pd.DataFrame(features, columns=["feat_"+str(i) for i in range((features.shape[1]))])
149     df_features['img'] = img
150     with pd.HDFStore(feature_save_file) as store:
151         store.append('df_features', df_features, data_columns = ['img'], min_itemsize=75)
```

```
1 %%time
2 extract_features(img,
3                 X,
4                 batch_size = 256,
5                 feature_save_file = feature_save_file)
(10000, 2048)
Wall time: 30.9 s
```


PROJECTION EN 2D

PROJECTION

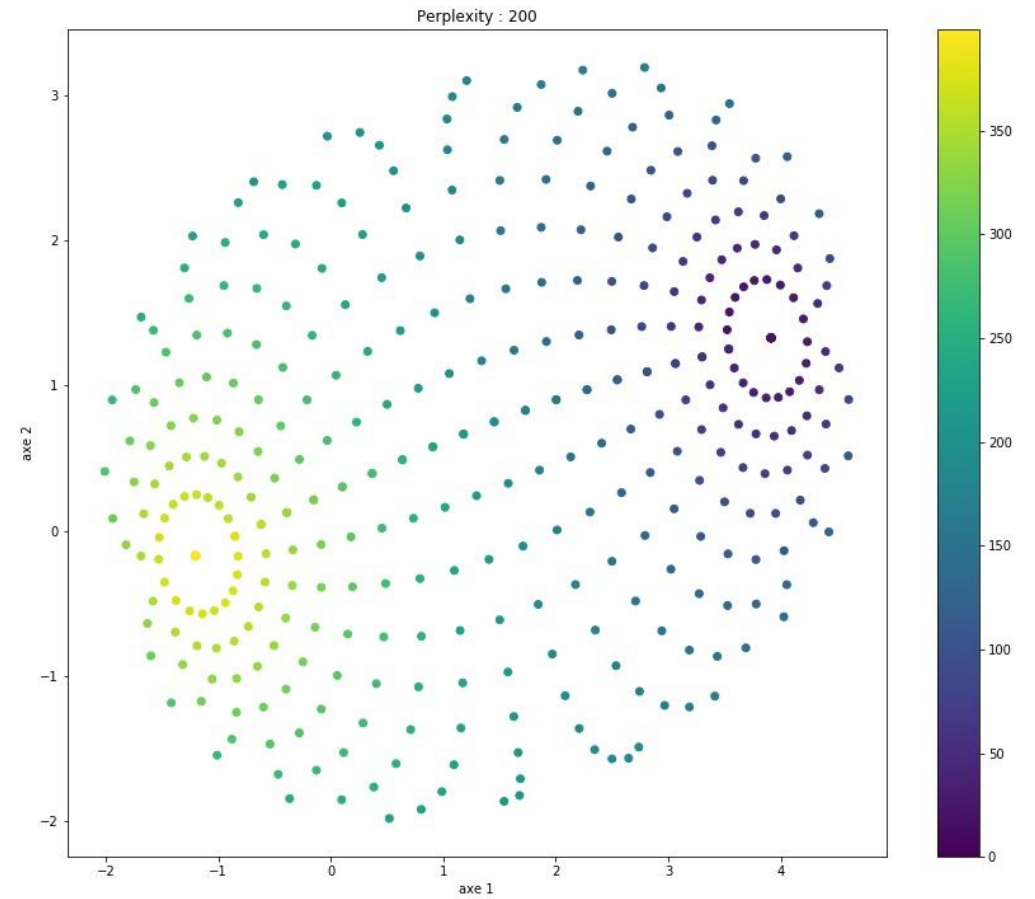
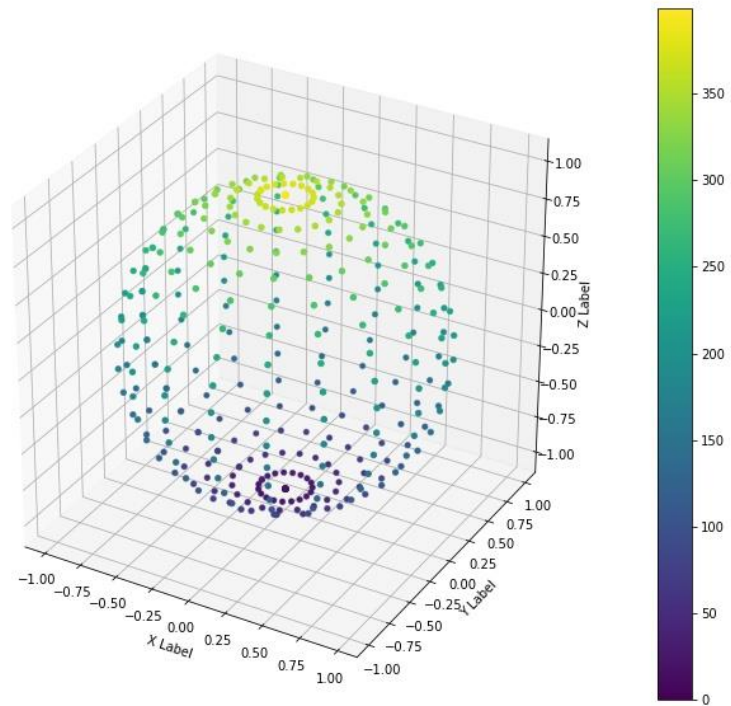
Objectif : créer des coordonnées (x,y) pour chaque image

Contrainte : des images proches doivent être similaires



T-SNE : algorithme de réduction de dimensions conservant les relations de proximité

PROJECTION : T-SNE



PROJECTION : T-SNE

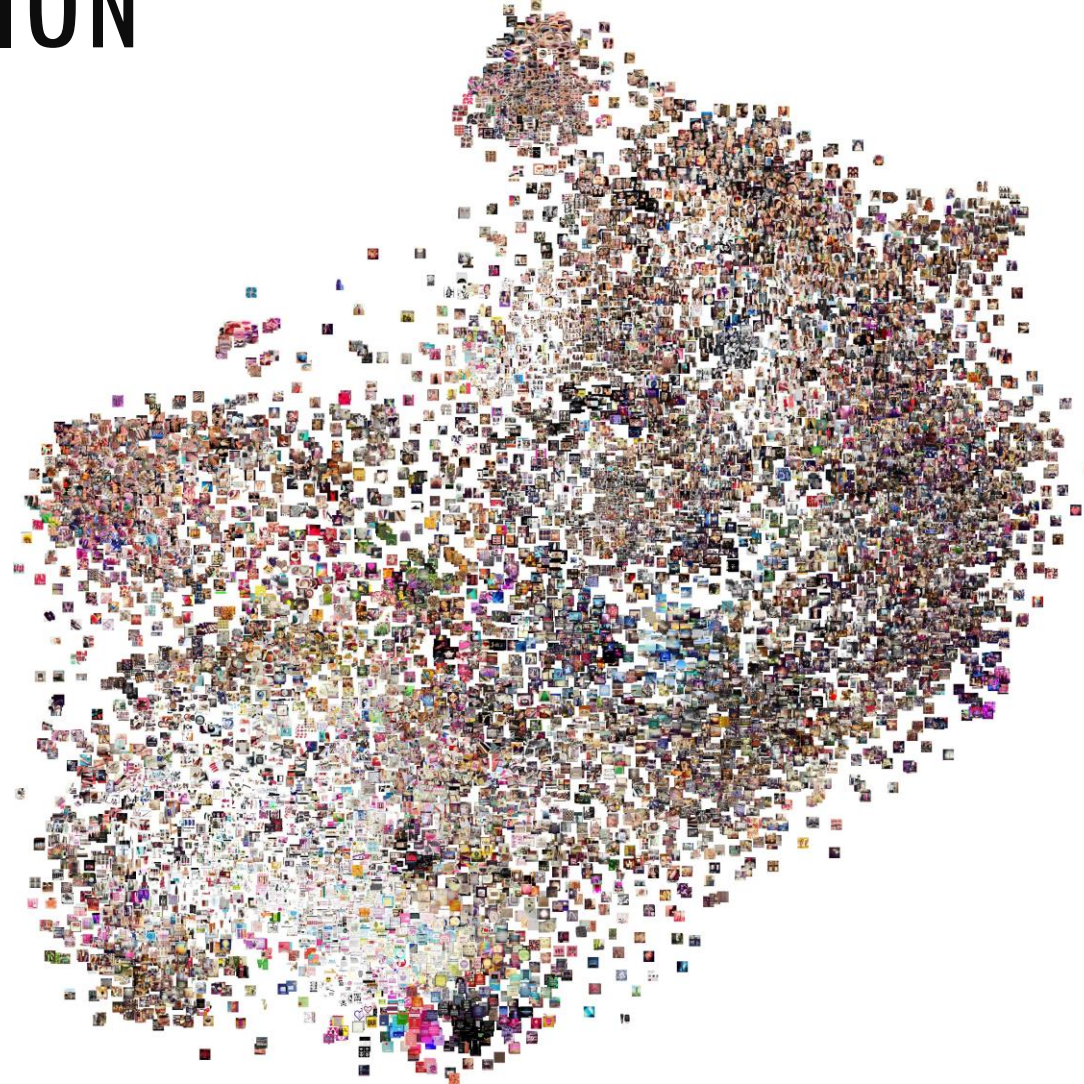
```
1 def do_tsne(img,
2             pca_feature_save_file = './pca_features.h5',
3             tsne_feature_save_file = './tsne_features.h5',
4             perplexity = 10
5             ):
6     """
7     Load PCA features of an image, transform it using parametric t-SNE and save it as HDF5.
8     img : list of img names, should be in the pca_feature_save_file.
9     pca_feature_save_file : HDF5 file where PCA features of the images are stored
10    ptsne_model_path : path to the ptsne model file (HDF5)
11    tsne_feature_save_file : HDF5 file to save transformed images
12    training : Boolean : if True, trains the ptsne model and saves it.
13    """
14
15    with pd.HDFStore(pca_feature_save_file) as store:
16        # for df in store.select('df_pca_features', where = "img in img", chunksize=chunksize):
17        df = store.select('df_pca_features', where = "img in img")
18
19        # load images
20        img = df['img'].as_matrix()
21        df.drop(['img'], inplace=True, axis=1)
22        pca_features = df.as_matrix()
23
24        print(pca_features.shape, img.shape)
25        high_dims = pca_features.shape[1]
26
27        output_res = TSNE(n_components=2, perplexity=perplexity).fit_transform(pca_features)
28
29        df_tsne_features = pd.DataFrame(output_res, columns=["tx", "ty"])
30        df_tsne_features['img'] = img
31
32    with pd.HDFStore(tsne_feature_save_file) as store:
33        store.append('df_tsne_features', df_tsne_features, data_columns = ['img'], min_itemsize=75)
```

VISUALISATION

VISUALISATION

```
1 def visualize(img,  
2             tsne_feature_save_file = './tsne_features.h5',  
3             save_img_path = ".ptSNE_img.png",  
4             img_path = './data/img/small_sample/file/',  
5             width = 5000,  
6             height = 5000,  
7             max_dim = 25):  
8  
9     import matplotlib.pyplot as plt  
10    with pd.HDFStore(tsne_feature_save_file) as store:  
11        df = store.select('df_tsne_features', where = "img in img")  
12  
13    tx = df['tx']  
14    ty = df['ty']  
15    img = df['img']  
16  
17    min_x = np.min(tx)  
18    max_x = np.max(tx)  
19    min_y = np.min(ty)  
20    max_y = np.max(ty)  
21    tx = (tx-min_x) / (max_x - min_x)  
22    ty = (ty-min_y) / (max_y - min_y)  
23  
24  
25    # t0 = time.time()  
26    full_image = Image.new('RGB', (width, height), color=(255,255,255)) # RGBA for png  
27    i = 0  
28    for img_, x, y in zip(img[:], tx[:], ty[:]):  
29        if i%int(len(tx[:])/10.)==0:  
30            print('img {} / {}'.format(i, len(tx[:])))  
31            i+=1  
32            tile = Image.open(img_path + img_)  
33            rs = max(1, tile.width/max_dim, tile.height/max_dim)  
34            tile = tile.resize((int(tile.width/rs), int(tile.height/rs)), Image.ANTIALIAS)  
35            full_image.paste(tile, (int((width-max_dim)*1*x), int((height-max_dim)*1*y)), mask=tile.convert('RGBA'))  
36  
37    # print("total time : {}".format(time.time()-t0))  
38    plt.figure(figsize = (16,12))  
39    # plt.imshow(full_image)  
40    full_image.save(save_img_path)  
  
1 %%time  
2 save_img_path = ".tests_perplexity/ptSNE_img_train_p{}_V2.jpg".format(perplexity)  
3 ##### below is visualisation#####  
4 visualize(img[:max_img_number],  
5         tsne_feature_save_file = tsne_feature_save_file,  
6         save_img_path = save_img_path,  
7         img_path = img_path,  
8         width = 5000,  
9         height = 5000,  
10        max_dim = 50)
```

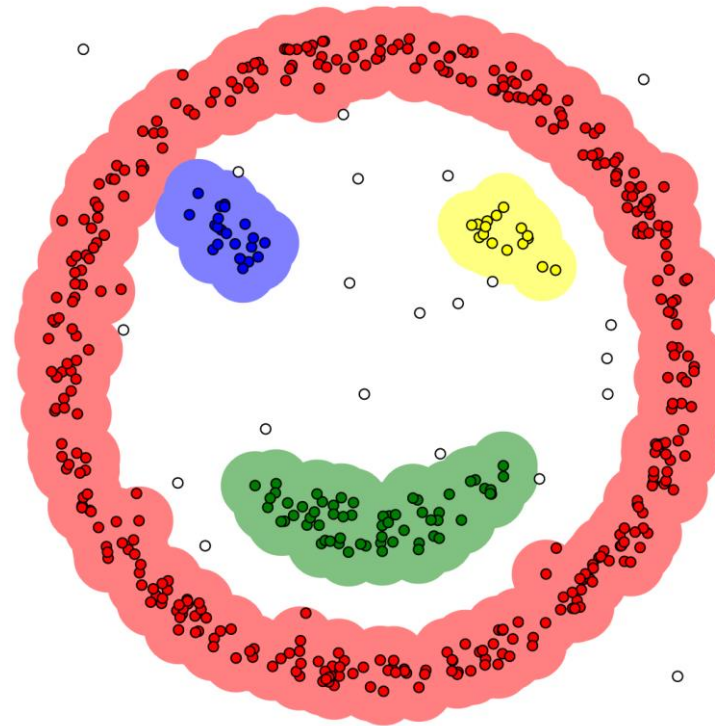
VISUALISATION



CLUSTERING & NETTOYAGE

CLUSTERING : DBSCAN

<https://www.naftaliharris.com/blog/visualizing-dbscan-clustering/>



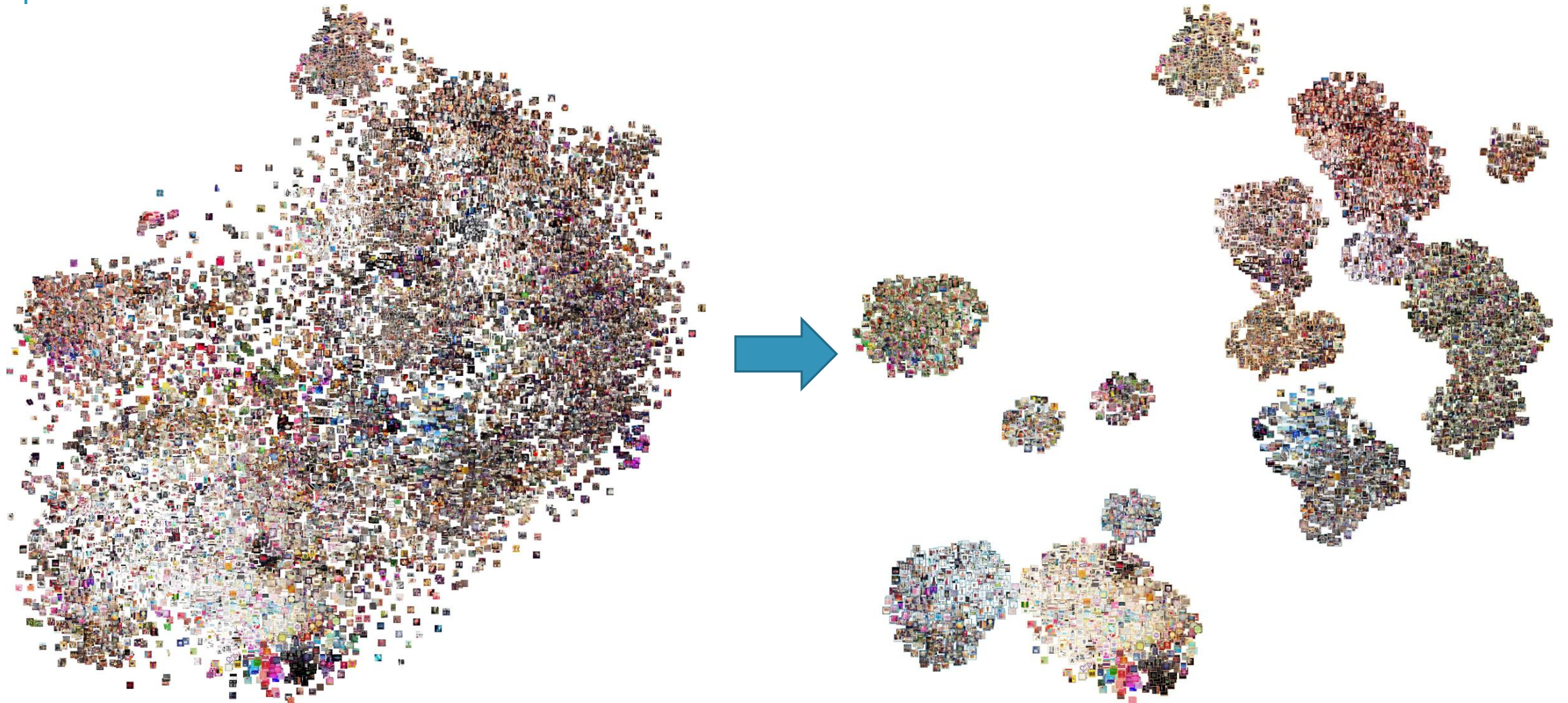
CLUSTERING : DBSCAN

```
32 from sklearn.cluster import DBSCAN
1  def visualize_DB(df,
2      save_img_path = "./ptSNE_img_DB.png",
3      img_path = './data/img/small_sample/file/',
4      width = 5000,
5      height = 5000,
6      max_dim = 25):
7
8      import matplotlib.pyplot as plt
9
10     df = df[df['cluster']!=-1] # remove unclustered img
11     tx = df['tx']
12     ty = df['ty']
13     img = df['img']
14     cluster = df['cluster']
15
16     min_x = np.min(tx)
17     max_x = np.max(tx)
18     min_y = np.min(ty)
19     max_y = np.max(ty)
20     tx = (tx-min_x) / (max_x - min_x) # normalize
21     ty = (ty-min_y) / (max_y - min_y)
22
23     color = ["#a6cee3",
24             "#1f78b4",
25             "#b2df8a",
26             "#33a02c",
27             "#fb9a99",
28             "#e31a1c",
29             "#fdbf6f",
30             "#ff7f00",
31             "#cab2d6",
32             "#6a3d9a",
33             "#ffff99",
34             "#b15928"]
35     color.extend(color)
36
37     full_image = Image.new('RGB', (width, height), color=(255,255,255))
38     i = 0
39     for img_, x, y, cluster in zip(img[:], tx[:], ty[:], cluster[:]):
40         if i%int(len(tx[:])/10.)==0:
41             print('img {} / {}'.format(i, len(tx[:])))
42             i+=1
43             tile = Image.open(img_path + img_)
44             rs = max(1, tile.width/max_dim, tile.height/max_dim)
45             tile = tile.resize((int(tile.width/rs), int(tile.height/rs)), Image.ANTIALIAS)
46             clr = color[cluster]
47             tile = ImageOps.expand(tile, border=2, fill=clr)
48             full_image.paste(tile, (int((width-max_dim)*1*x), int((height-max_dim)*1*y)), mask=tile.convert('RGBA'))
49
50     plt.figure(figsize = (16,12))
51     full_image.save(save_img_path)
```

```
1  eps = 3
2  min_sample = 90
3  clustering = DBSCAN(eps=eps, min_samples=min_sample).fit(df[['tx', 'ty']])
4  cnt = Counter(clustering.labels_)
5  print("eps = {}, n_clusters = {}, unclassified = {}".format(eps, len(cnt), cnt[-1]))
6  df['cluster'] = clustering.labels_
7
8  visualize_DB(df,
9      save_img_path = "./tests_DB/ptSNE_img_DB_eps{}_min_sample{}_v3.1.jpg".format(eps, min_sample),
10     img_path = img_path,
11     width = 5000,
12     height = 5000,
13     max_dim = 50)
```

eps = 3, n_clusters = 15, unclassified = 4381

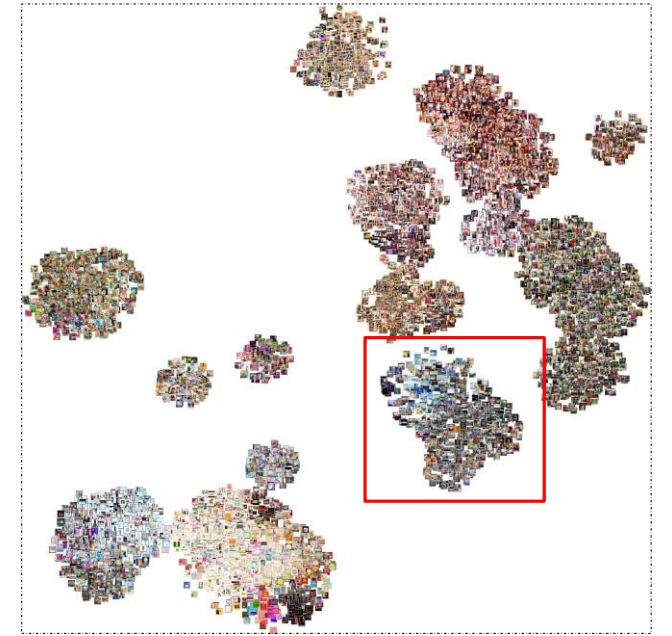
VISUALISATION



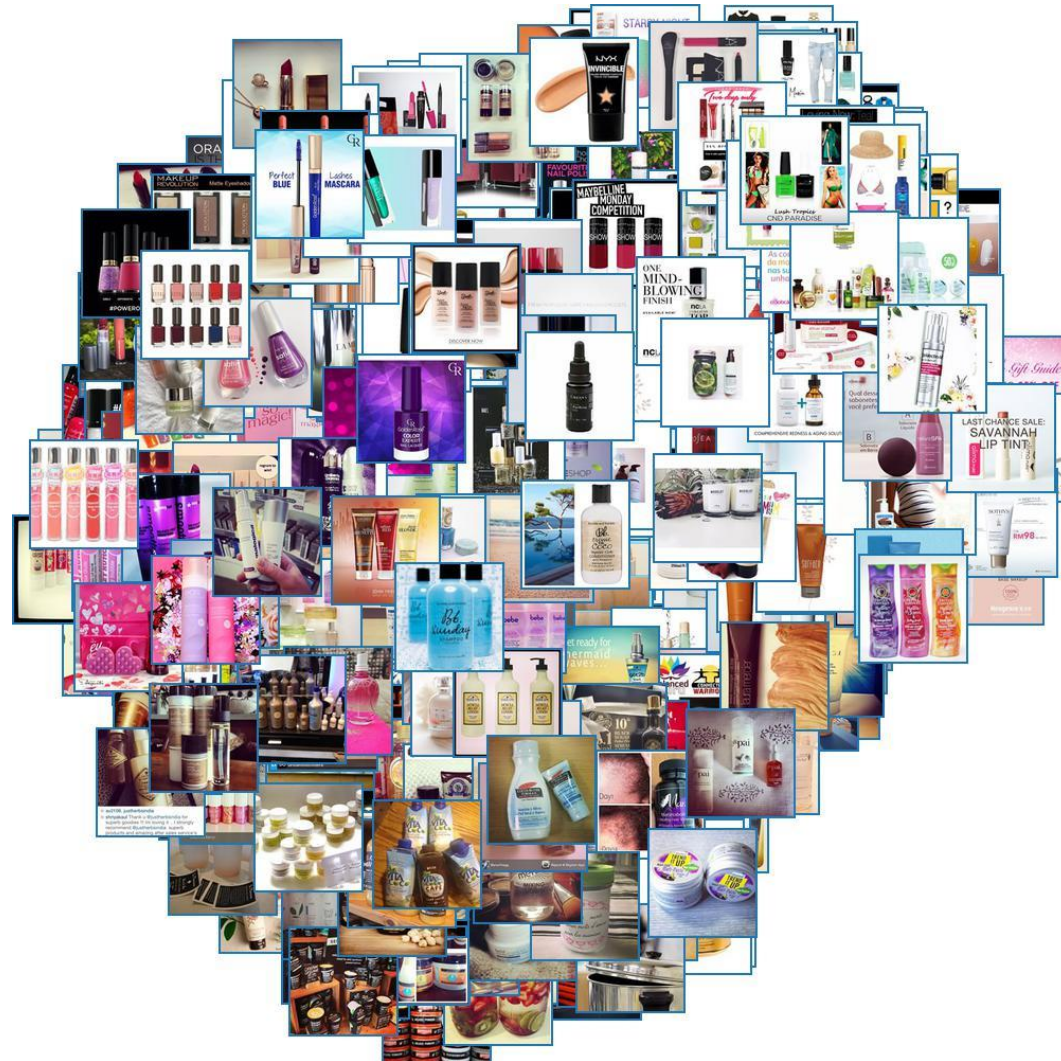
VISUALISATION



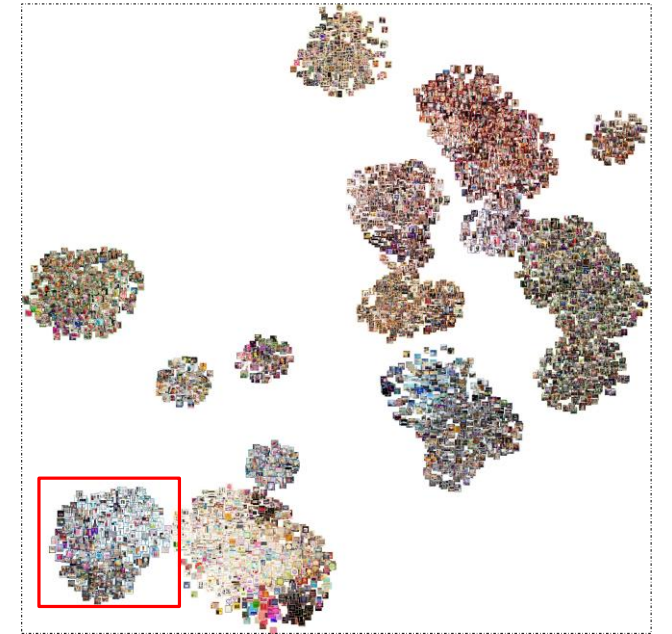
Cluster 0



VISUALISATION



Cluster 1



VISUALISATION



Cluster 3

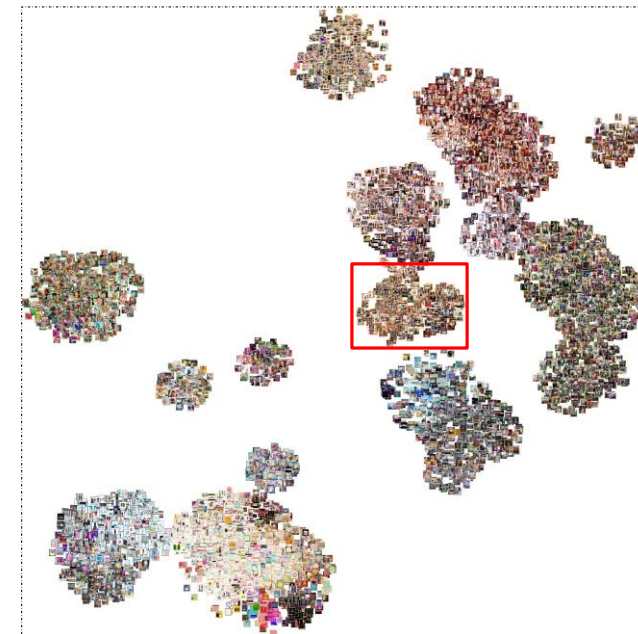




VISUALISATION



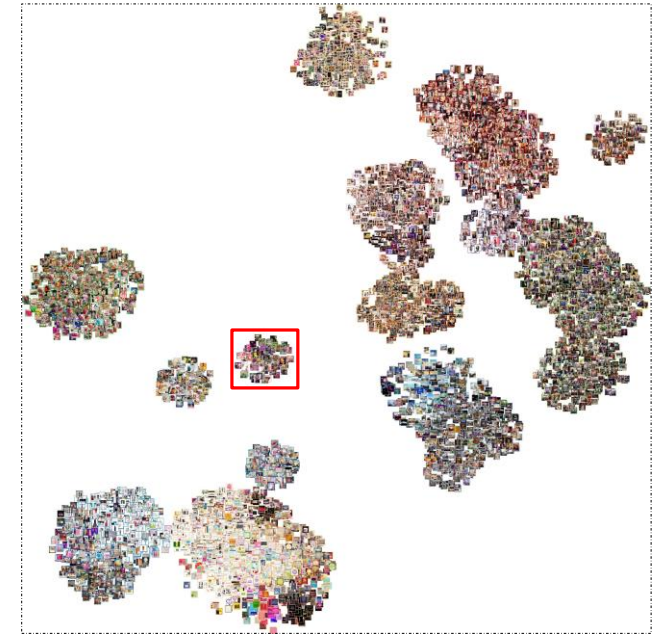
Cluster 7



VISUALISATION

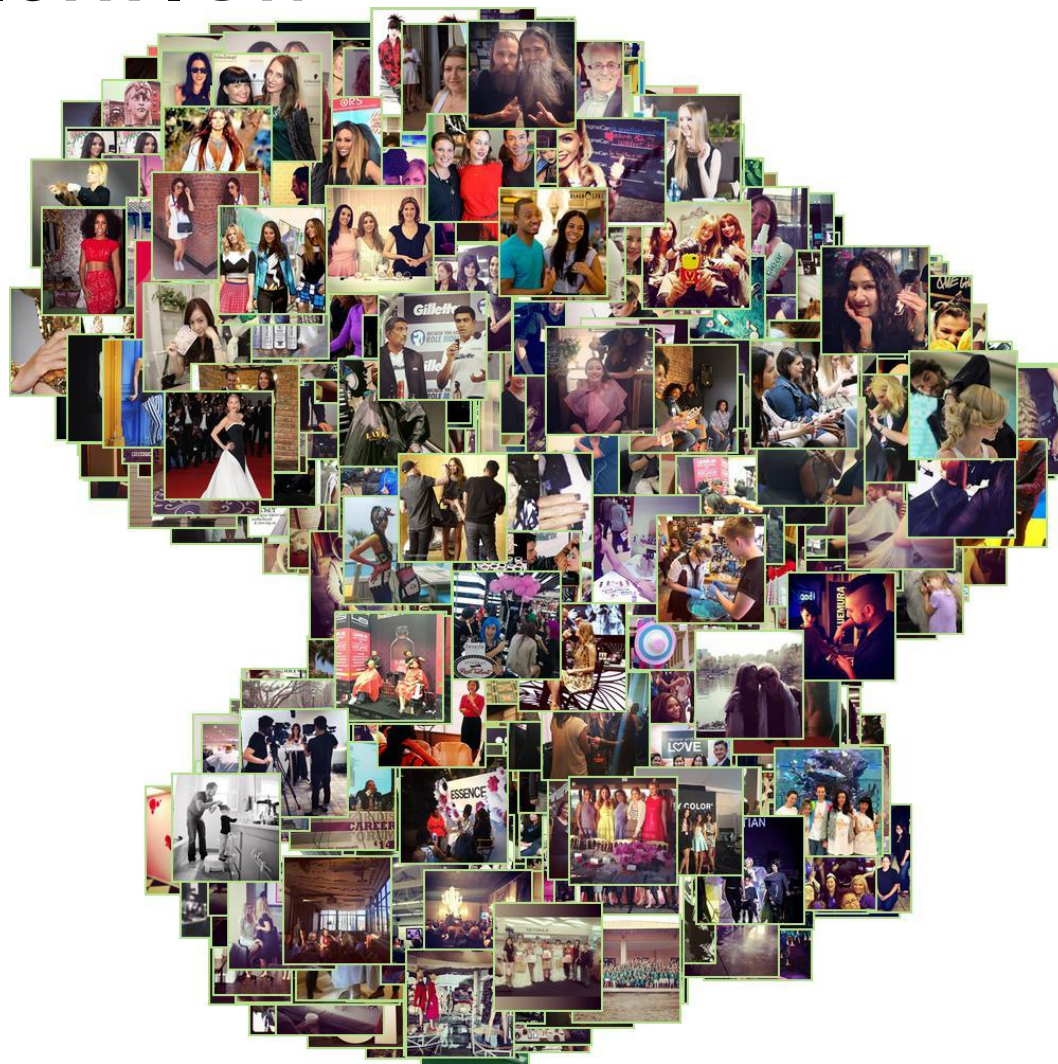


Cluster 8

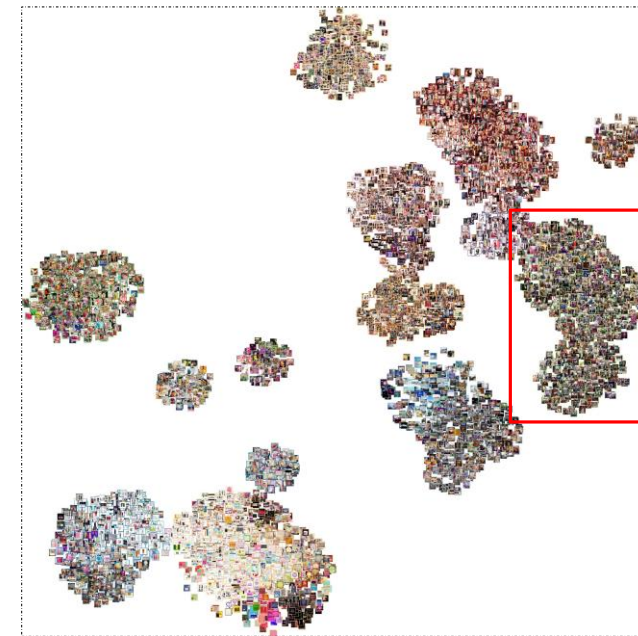




VISUALISATION



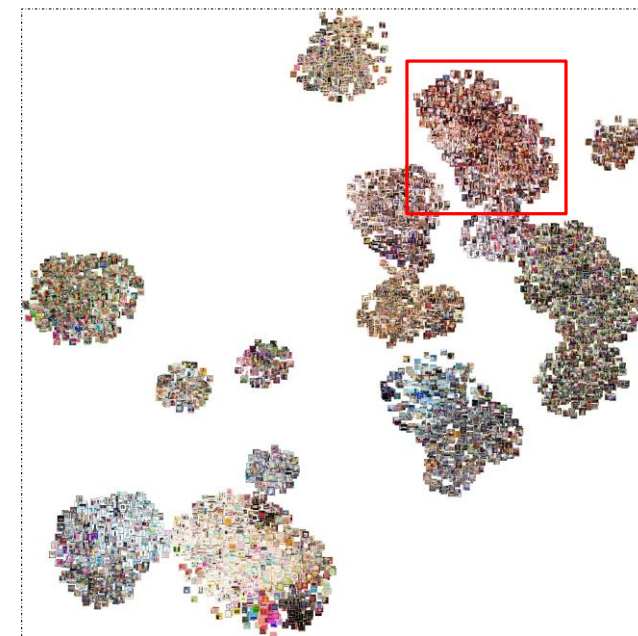
Cluster 2



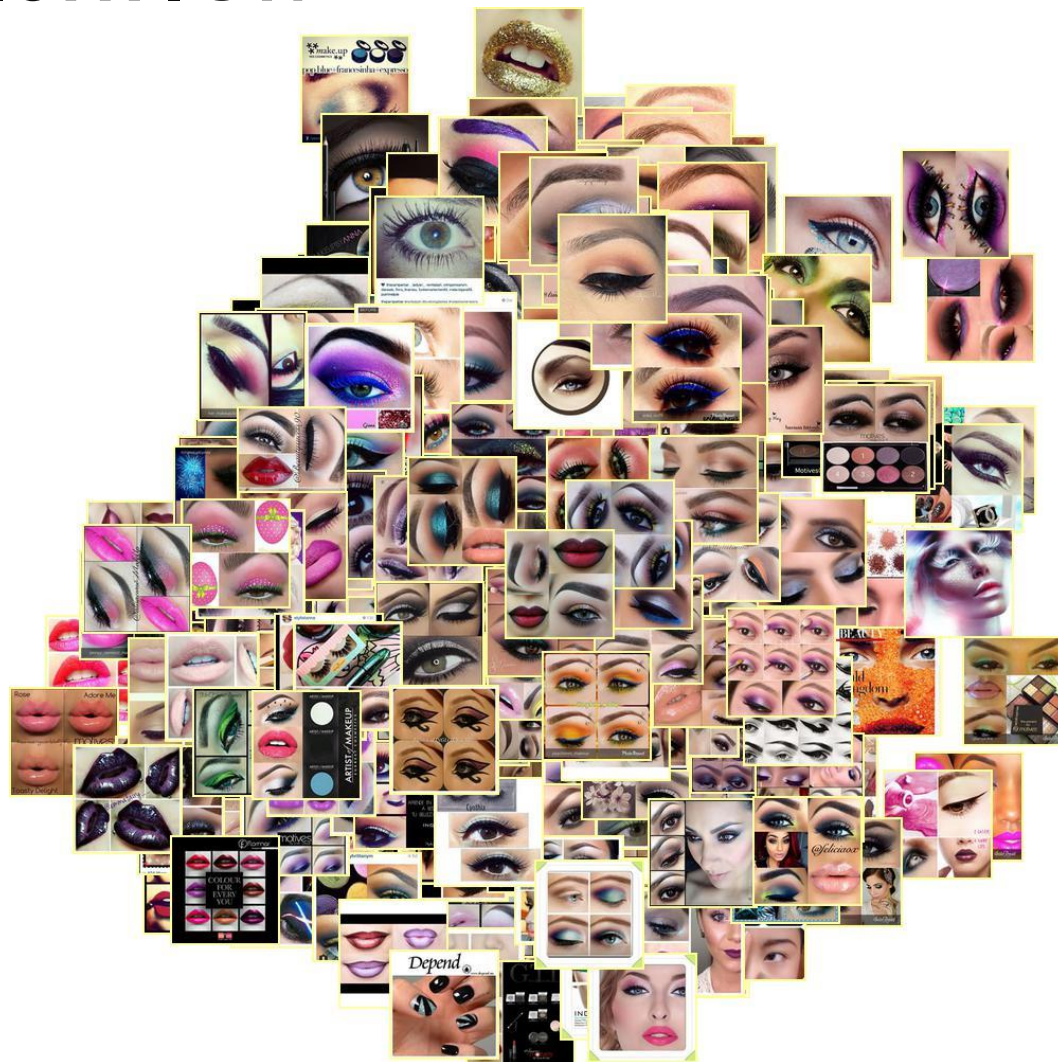
VISUALISATION



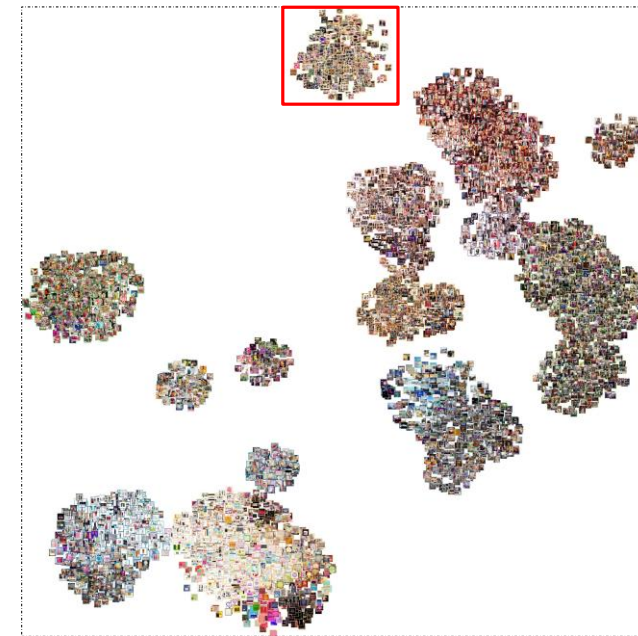
Cluster 5



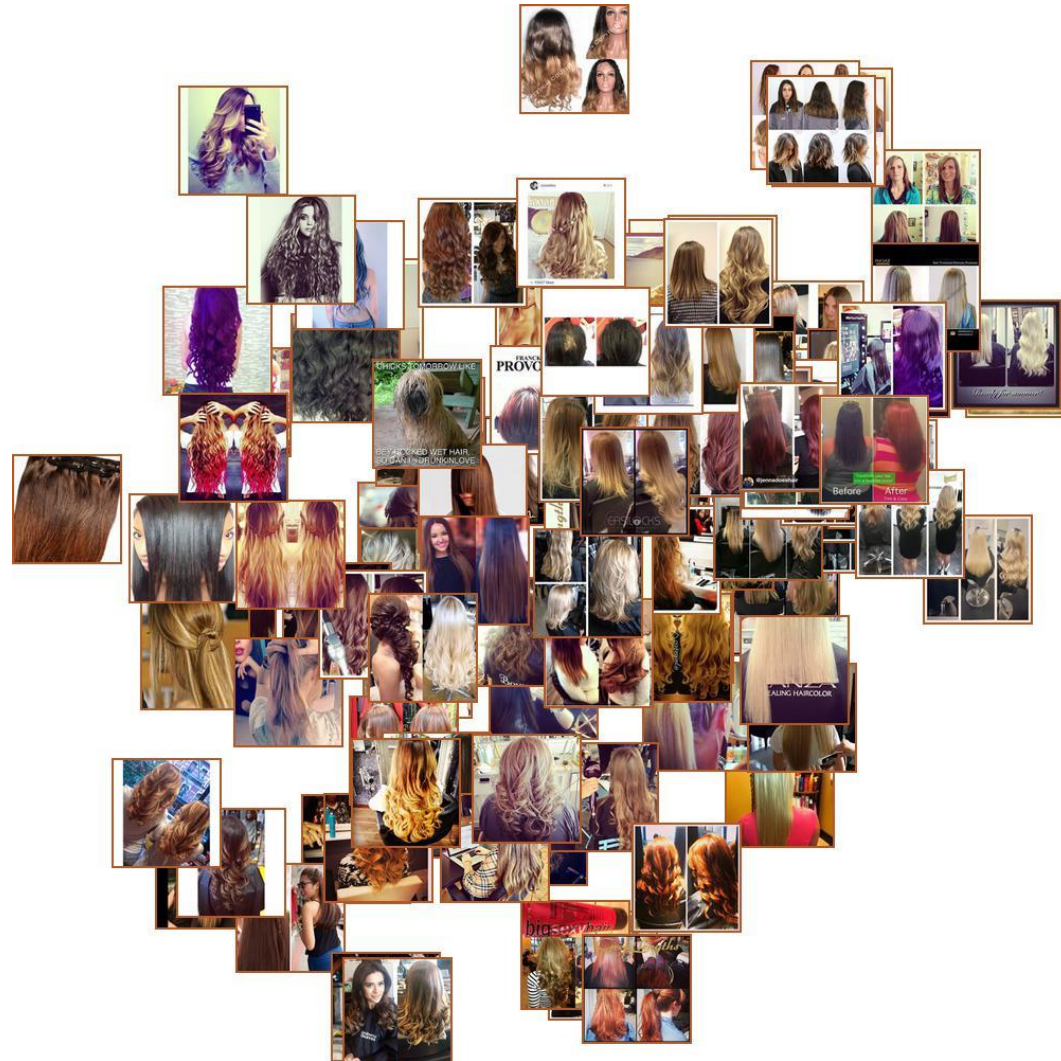
VISUALISATION



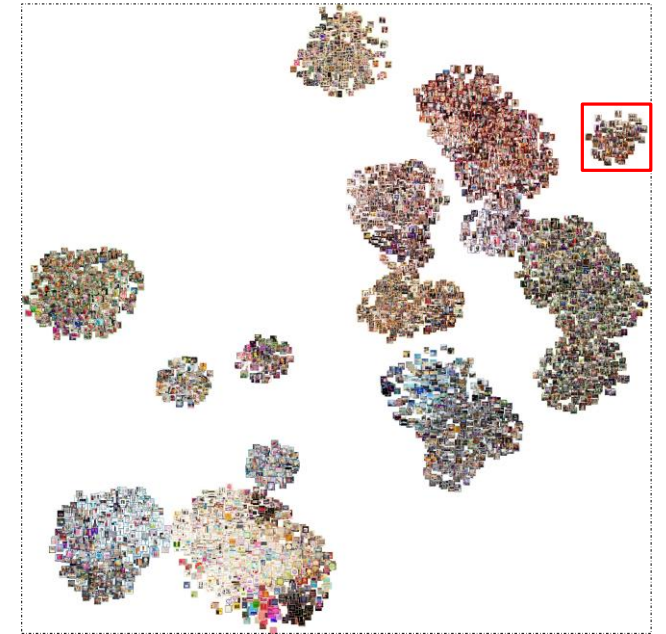
Cluster 10



VISUALISATION



Cluster 11

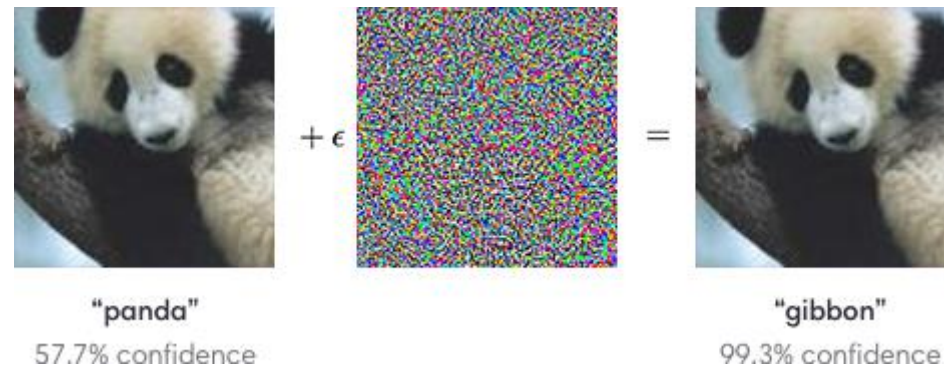


LIMITES

Images similaires au dataset d'entraînement ImageNet

Pas de contrôle de la définition de similarité : il peut y avoir des « erreurs » par rapport à notre perception d'images similaires

Sensible au bruit :



<https://arxiv.org/abs/1412.6572>

MERCI !

Corentin Martel

Data Scientist Freelance

<https://www.corelyo.com>



RESSOURCES

<https://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf>

<http://huboqiang.cn/2017/02/20/DeepLearningSkinCancer>

<https://cv-tricks.com/cnn/understand-resnet-alexnet-vgg-inception/>

<https://arxiv.org/pdf/1512.03385.pdf>

O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. arXiv:1409.0575, 2014.

https://nnabla.org/paper/imagenet_in_224sec.pdf

<https://towardsdatascience.com/intuitively-understanding-convolutions-for-deep-learning-1f6f42faee1>

http://cs231n.stanford.edu/slides/2018/cs231n_2018_lecture05.pdf

<https://towardsdatascience.com/how-to-visualize-convolutional-features-in-40-lines-of-code-70b7d87b0030>

<https://distill.pub/2016/misread-tsne/>

<http://alexhwilliams.info/itsneuronalblog/2016/03/27/pca/>

<https://www.qwertee.io/blog/deep-learning-with-point-clouds/>