

# Kaggle setup

```
1 ! pip install -q kaggle
```

```
1 from google.colab import files  
2 files.upload()
```

Choose Files No file chosen Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.  
Saving kaggle.json to kaggle.json  
{'kaggle.json':  
h'{"username":"shreavanandigarim","key":"38cebed522d9a2c21bab1d8had8fa162"}'}

```
1 #getting api key  
2 ! mkdir ~/.kaggle  
3 ! cp kaggle.json ~/.kaggle/
```

```
1 #changing permission  
2 ! chmod 600 ~/.kaggle/kaggle.json
```

```
1 #loading ISL dataset  
2 !kaggle datasets download prathumarikeri/indian-sign-language-isl
```

Downloading indian-sign-language-isl.zip to /content  
98% 276M/281M [00:02<00:00, 76.8MB/s]  
100% 281M/281M [00:02<00:00, 98.6MB/s]

```
1 #unzip data  
2 !unzip indian-sign-language-isl.zip
```

**Streaming output truncated to the last 5000 lines.**

```
inflating: Indian/V/819.jpg  
inflating: Indian/V/82.jpg  
inflating: Indian/V/820.jpg  
inflating: Indian/V/821.jpg  
inflating: Indian/V/822.jpg  
inflating: Indian/V/823.jpg  
inflating: Indian/V/824.jpg  
inflating: Indian/V/825.jpg  
inflating: Indian/V/826.jpg  
inflating: Indian/V/827.jpg  
inflating: Indian/V/828.jpg  
inflating: Indian/V/829.jpg  
inflating: Indian/V/83.jpg  
inflating: Indian/V/830.jpg  
inflating: Indian/V/831.jpg
```

```
inflating: Indian/V/832.jpg
inflating: Indian/V/833.jpg
inflating: Indian/V/834.jpg
inflating: Indian/V/835.jpg
inflating: Indian/V/836.jpg
inflating: Indian/V/837.jpg
inflating: Indian/V/838.jpg
inflating: Indian/V/839.jpg
inflating: Indian/V/84.jpg
inflating: Indian/V/840.jpg
inflating: Indian/V/841.jpg
inflating: Indian/V/842.jpg
inflating: Indian/V/843.jpg
inflating: Indian/V/844.jpg
inflating: Indian/V/845.jpg
inflating: Indian/V/846.jpg
inflating: Indian/V/847.jpg
inflating: Indian/V/848.jpg
inflating: Indian/V/849.jpg
inflating: Indian/V/85.jpg
inflating: Indian/V/850.jpg
inflating: Indian/V/851.jpg
inflating: Indian/V/852.jpg
inflating: Indian/V/853.jpg
inflating: Indian/V/854.jpg
inflating: Indian/V/855.jpg
inflating: Indian/V/856.jpg
inflating: Indian/V/857.jpg
inflating: Indian/V/858.jpg
inflating: Indian/V/859.jpg
inflating: Indian/V/86.jpg
inflating: Indian/V/860.jpg
inflating: Indian/V/861.jpg
inflating: Indian/V/862.jpg
inflating: Indian/V/863.jpg
inflating: Indian/V/864.jpg
inflating: Indian/V/865.jpg
inflating: Indian/V/866.jpg
inflating: Indian/V/867.jpg
inflating: Indian/V/868.jpg
inflating: Indian/V/869.jpg
inflating: Indian/V/87.jpg
```

```
1 ls
```

```
Indian/ indian-sign-language-isl.zip kaggle.json sample_data/
```

## Import Modules

```
1 # import section
2
3 %matplotlib inline
```

```
4 from google.colab import files
5 import os
6
7 # TensorFlow and tf.keras
8 import tensorflow as tf
9 from tensorflow import keras
10
11 # Helper libraries
12 import numpy as np
13 from numpy import random
14 import matplotlib.pyplot as plt
15 import cv2
16 import pandas as pd
17 import scipy
18 from google.colab.patches import cv2_imshow
19
20 # Sklearn
21 from sklearn.model_selection import train_test_split
22 from sklearn.metrics import confusion_matrix
23 from sklearn.preprocessing import OneHotEncoder
24 from sklearn.utils import shuffle
25
26 print(tf.__version__)
```

2.7.0

```
1 import numpy as np
2 import cv2
3 import os
4 import csv
5 import sklearn.metrics as sm
6 from sklearn.cluster import MiniBatchKMeans
7 from sklearn.svm import SVC
8 import random
9 import warnings
10 import pickle
11 from sklearn.naive_bayes import GaussianNB as nb
12 from sklearn.neighbors import KNeighborsClassifier as knn
13 from sklearn.linear_model import LogisticRegression as lr
14 from sklearn.neural_network import MLPClassifier as mlp
15 import numpy as np
16 import sklearn.metrics as sm
```

## Loading Dataset

```
1 # walk through the directory to get label names
2 imagepaths = []
```

```

3
4 for root, dirs, files in os.walk(".", topdown=False):
5     for name in files:
6         path = os.path.join(root, name)
7         if path.endswith(".jpg"): # only pick the images with .jpg as per our c
8             imagepaths.append(path)
9
10 print(len(imagepaths))

```

42745

```

1 #debugging and check
2 def plot_image(path):
3     img = cv2.imread(path) # Reads the image into a numpy.array
4     img_cvt = cv2.cvtColor(img, cv2.COLOR_BGR2GRAY) # Converts into the corre
5     print(img_cvt.shape) # Prints the shape of the image just to check
6     plt.grid(False) # Without grid so we can see better
7     plt.imshow(img_cvt) # Shows the image
8     plt.xlabel("Width")
9     plt.ylabel("Height")
10    plt.title("Image " + path)

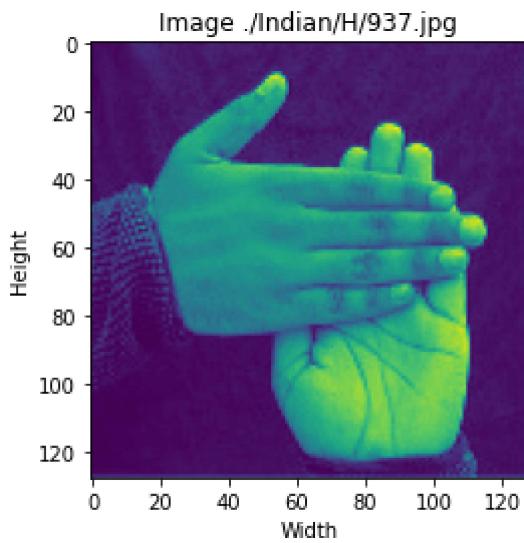
```

```

1 #check
2 plot_image(imagepaths[0])

```

(128, 128)



## Acquiring labels from folder names and preprocessing

```

1 images = [] # Image data
2 y = [] # Labels
3

```

```

4 # Loops through imagepaths to load images and labels into arrays
5 for path in imagepaths:
6     img = cv2.imread(path) # Reads image and returns np.array
7     img = cv2.cvtColor(img, cv2.COLOR_BGR2GRAY) # Converts into the correct c
8     img = cv2.resize(img, (200, 200)) # Reduce image size so training can be
9     images.append(img)
10
11 # Processing label in image path
12 category = path.split("/")[2]
13 #print(category)
14 label = category
15 y.append(label)
16
17 # Turn images and y into np.array to speed up train_test_split
18 images = np.array(images, dtype="uint8")
19 images = images.reshape(len(imagepaths), 200, 200, 1) # reshaping
20 y = np.array(y)
21
22 print("Images loaded: ", len(images))
23 print("Labels loaded: ", len(y))
24
25 print(y[42000], imagepaths[42000]) # Debugging

```

Images loaded: 42745

Labels loaded: 42745

3 ./Indian/3/197.jpg

```

1 #Test check
2 print(y[5500])
3 cv2_imshow(images[5500])

```

P



## Exploratory Data Analysis

```

1 #Getting mean image of each class
2 indices=[]

```

```
3 distinct_labels = set()
4 for i in range(len(y)):
5     if y[i] not in distinct_labels:
6         distinct_labels.add(y[i])
7         indices.append(i)
8 #print(indices)
9 print("      Sample images      Mean images of each class")
10 for i in range(len(indices)-1):
11     #cv2_imshow(X[indices[i]])
12     print(y[indices[i]])
13     mean_img = np.mean(images[indices[i]:indices[i+1]],axis=0)
14     mean_img = mean_img.reshape((200,200,1))
15     #cv2_imshow(mean_img)
16     img_concate_Hori=np.concatenate((images[indices[i]],mean_img),axis=1)
17     cv2_imshow(img_concate_Hori)
18
19
20
21 print(y[indices[len(indices)-1]])
22 mean_img = np.mean(images[indices[len(indices)-1]:len(images)-1],axis=0)
23 mean_img = mean_img.reshape((200,200,1))
24 #cv2_imshow(X[indices[len(indices)-1]])
25 #mean_img = np.mean(X[len(indices)-1:len(X)])
26 img_concate_Hori=np.concatenate((images[indices[len(indices)-1]],mean_img),
27 cv2_imshow(img_concate_Hori)
28
29
30
```



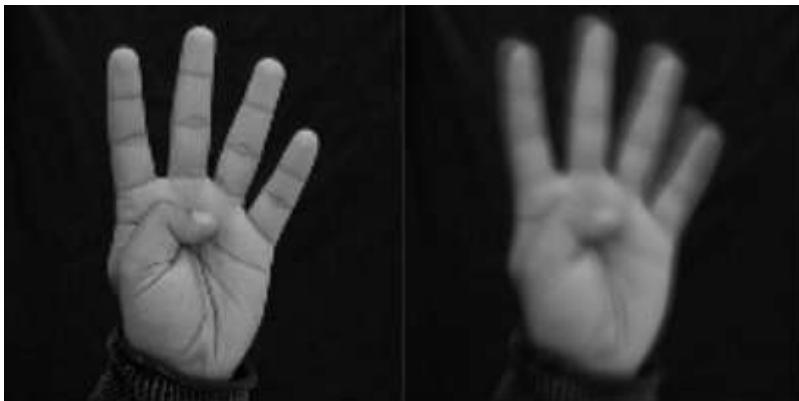
Sample images

Mean images of each class

H



4



D



T

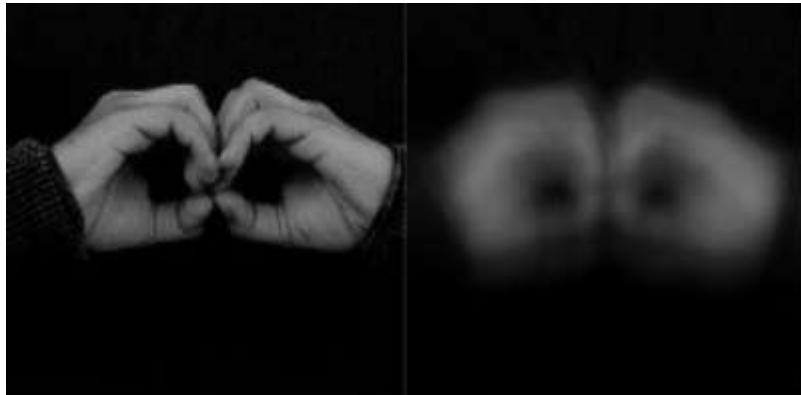


P





B



1

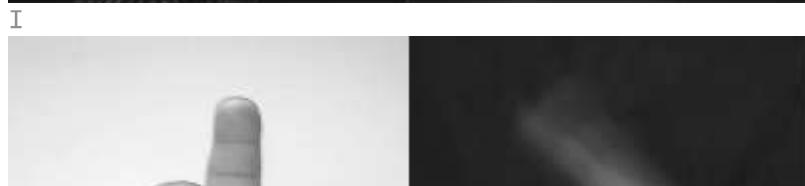
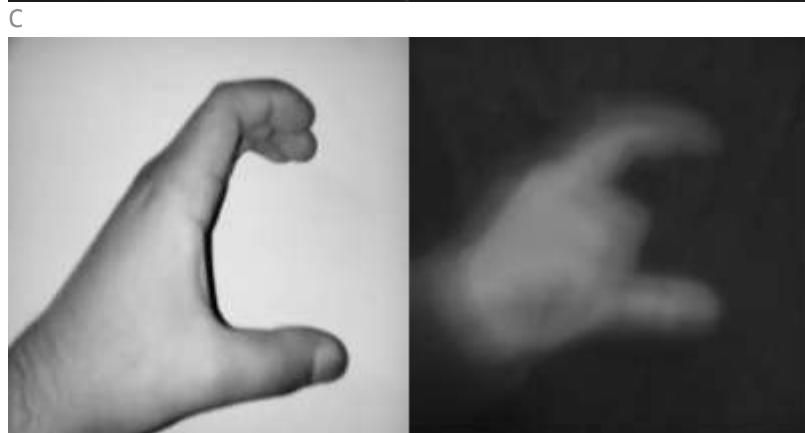


Y



U







2



9

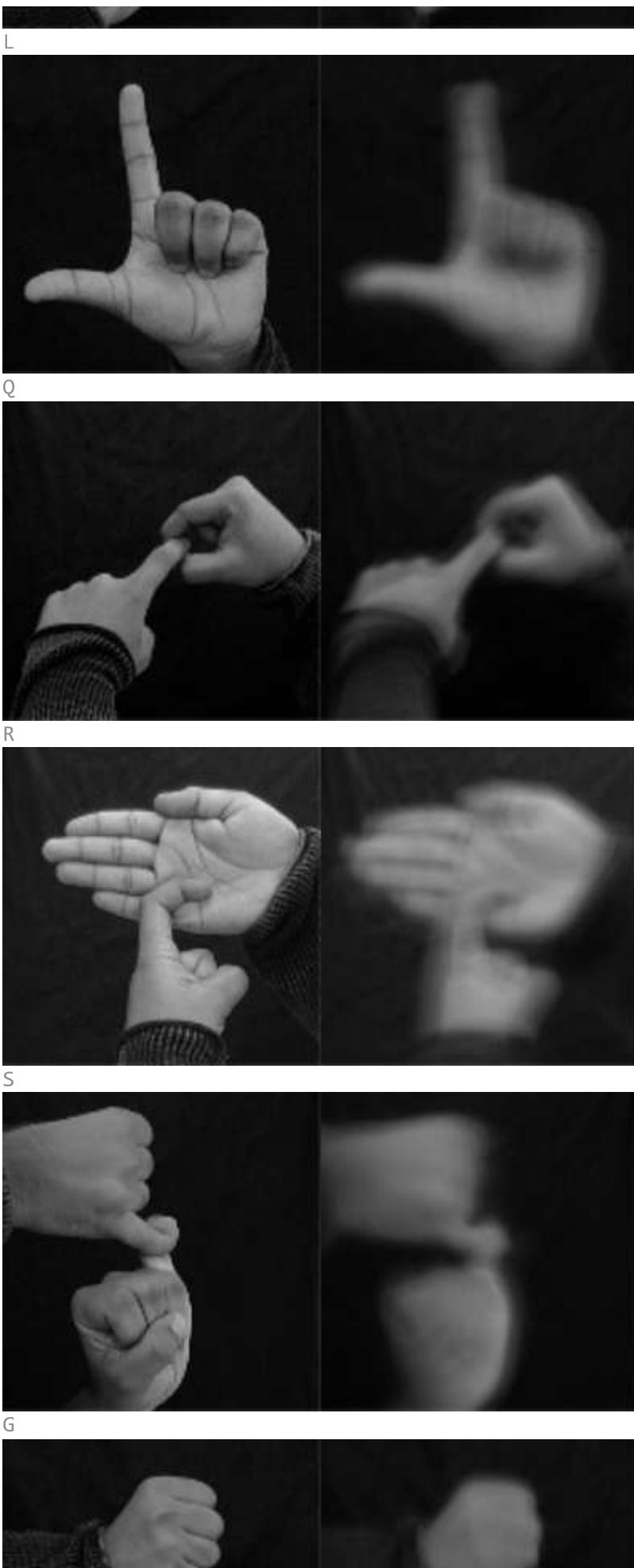


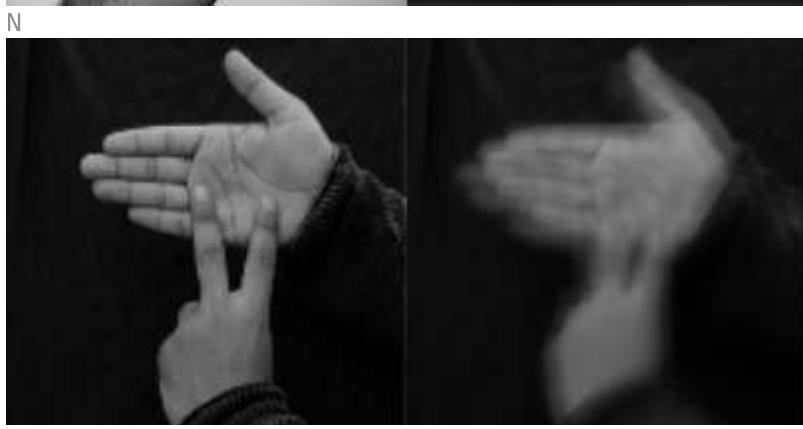
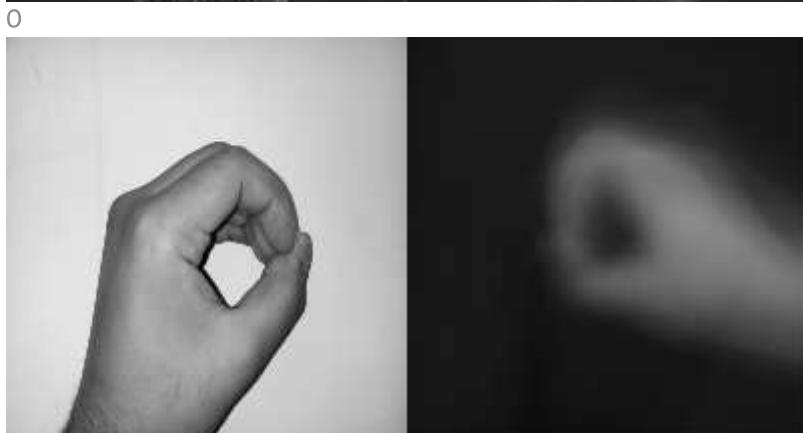
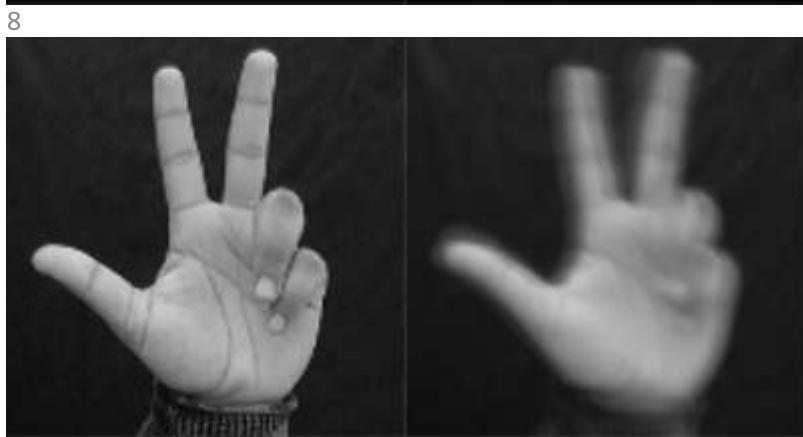
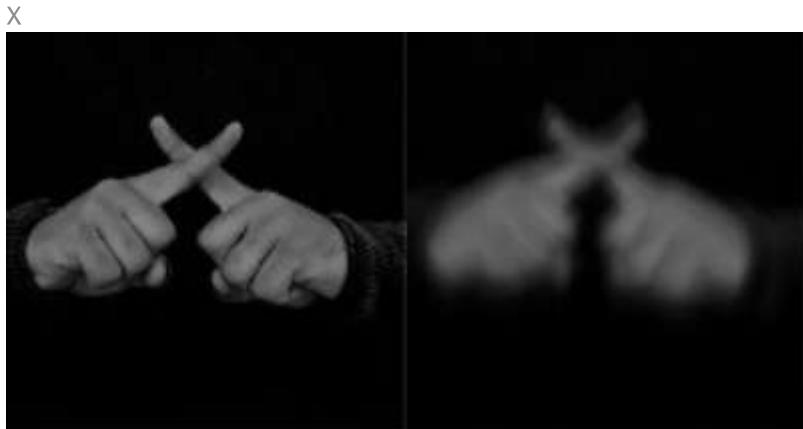
7



K







M



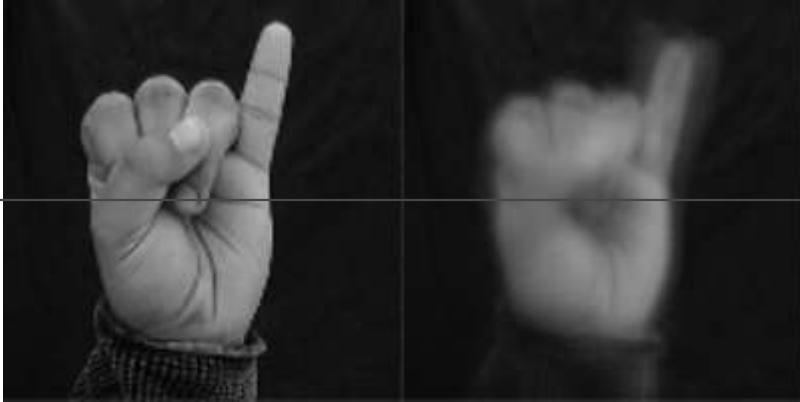
J



F



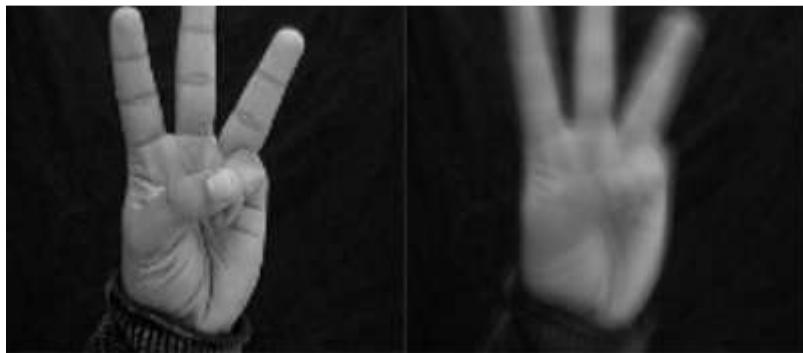
6



5



```
1 #Getting variance image of each class
2 indices=[]
3 distinct_labels = set()
4 for i in range(len(y)):
5     if y[i] not in distinct_labels:
6         distinct_labels.add(y[i])
7         indices.append(i)
8 #print(indices)
9 print("      Sample images      Variance images of each class")
10 for i in range(len(indices)-1):
11     #cv2_imshow(X[indices[i]])
12     print(y[indices[i]])
13     var_img = np.var(images[indices[i]:indices[i+1]],axis=0)
14     var_img = var_img.reshape((200,200,1))
15     #cv2_imshow(mean_img)
16     img_concatenate_Hori=np.concatenate((images[indices[i]],var_img),axis=1)
17     cv2_imshow(img_concatenate_Hori)
18
19
20
21 print(y[indices[len(indices)-1]])
22 var_img = np.var(images[indices[len(indices)-1]:len(images)-1],axis=0)
23 var_img = var_img.reshape((200,200,1))
24 #cv2_imshow(X[indices[len(indices)-1]])
25 #mean_img = np.mean(X[len(indices)-1:len(X)])
26 img_concatenate_Hori=np.concatenate((images[indices[len(indices)-1]],var_img),axis=1)
27 cv2_imshow(img_concatenate_Hori)
28
29
30
```

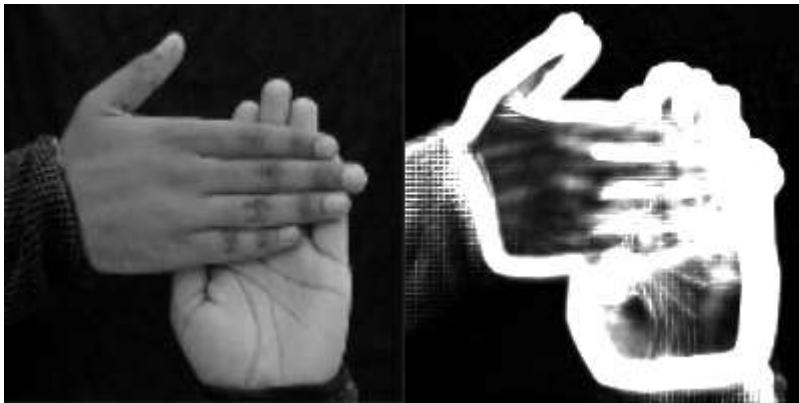




Sample images

Variance images of each class

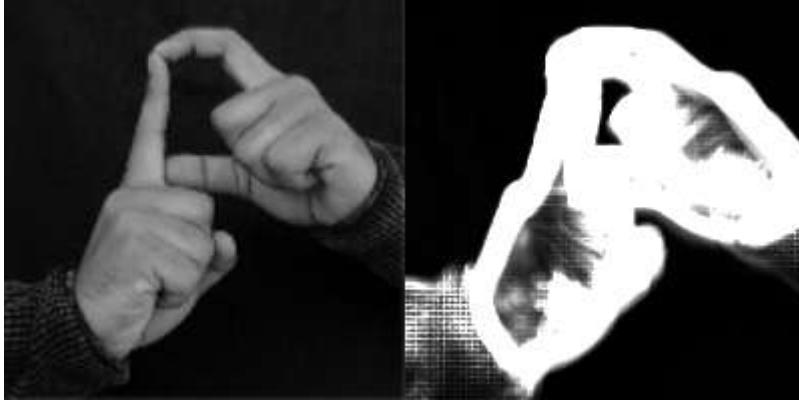
H



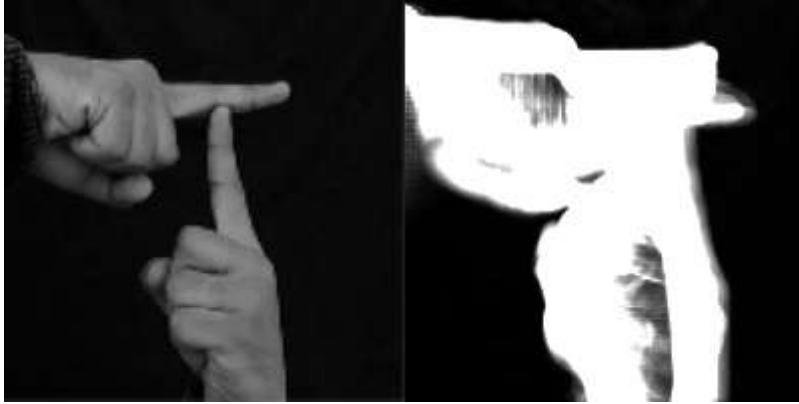
4



D



T

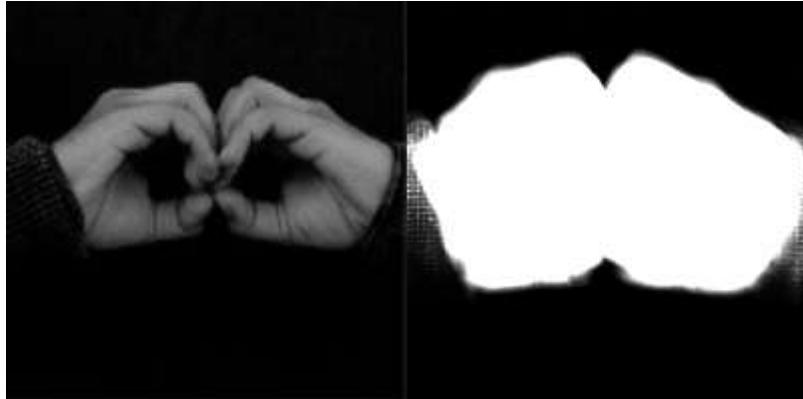


P

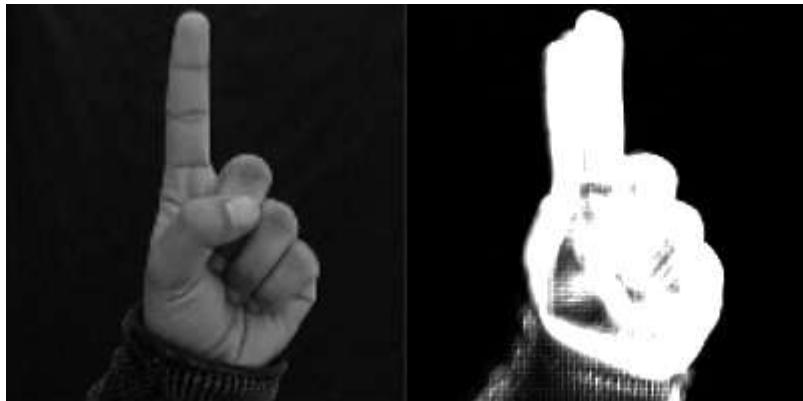




B



1

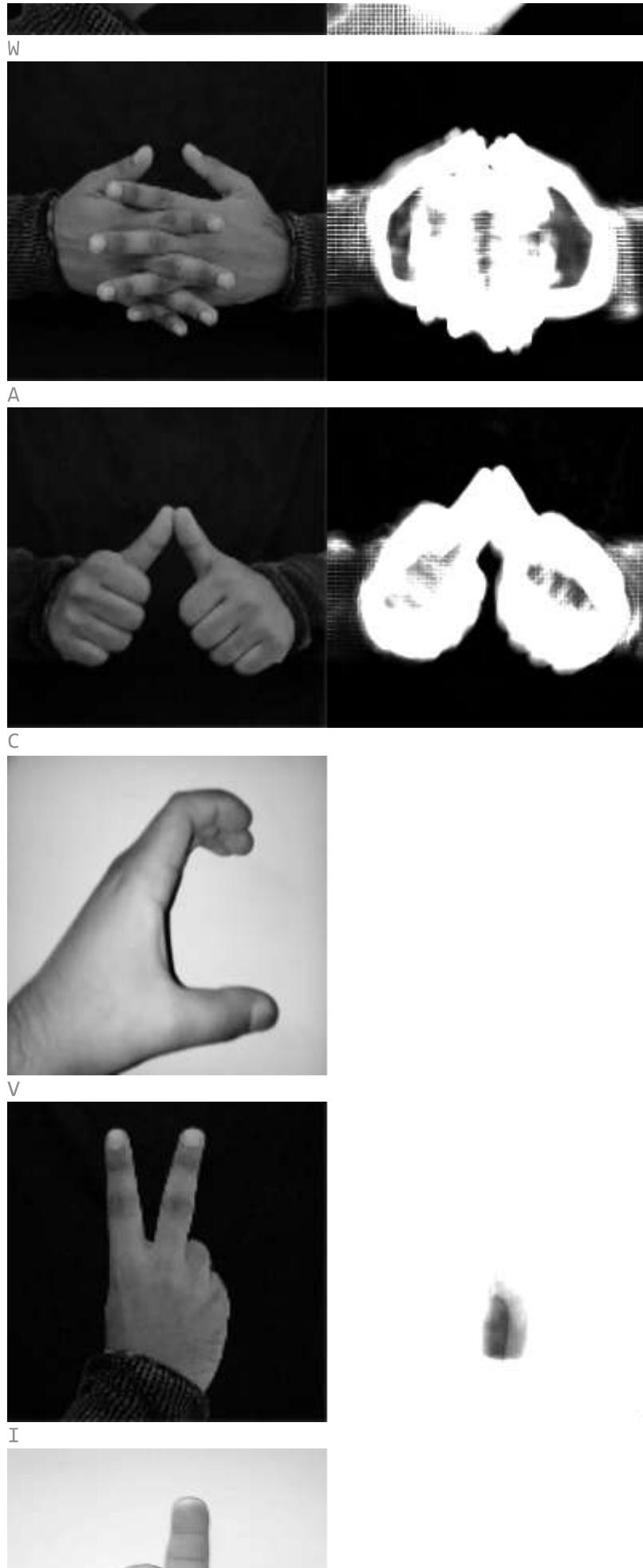


Y



U



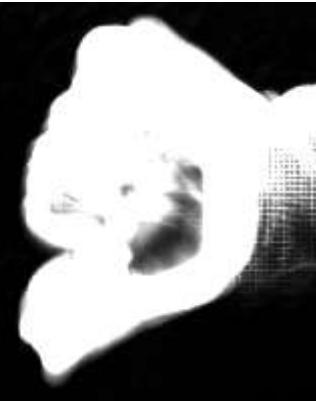




2



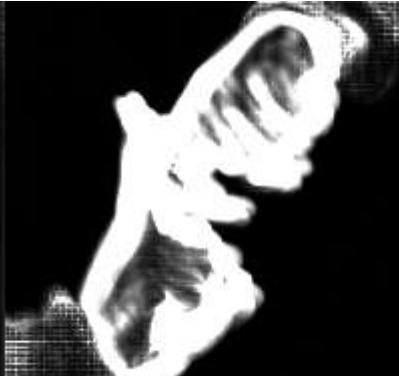
9

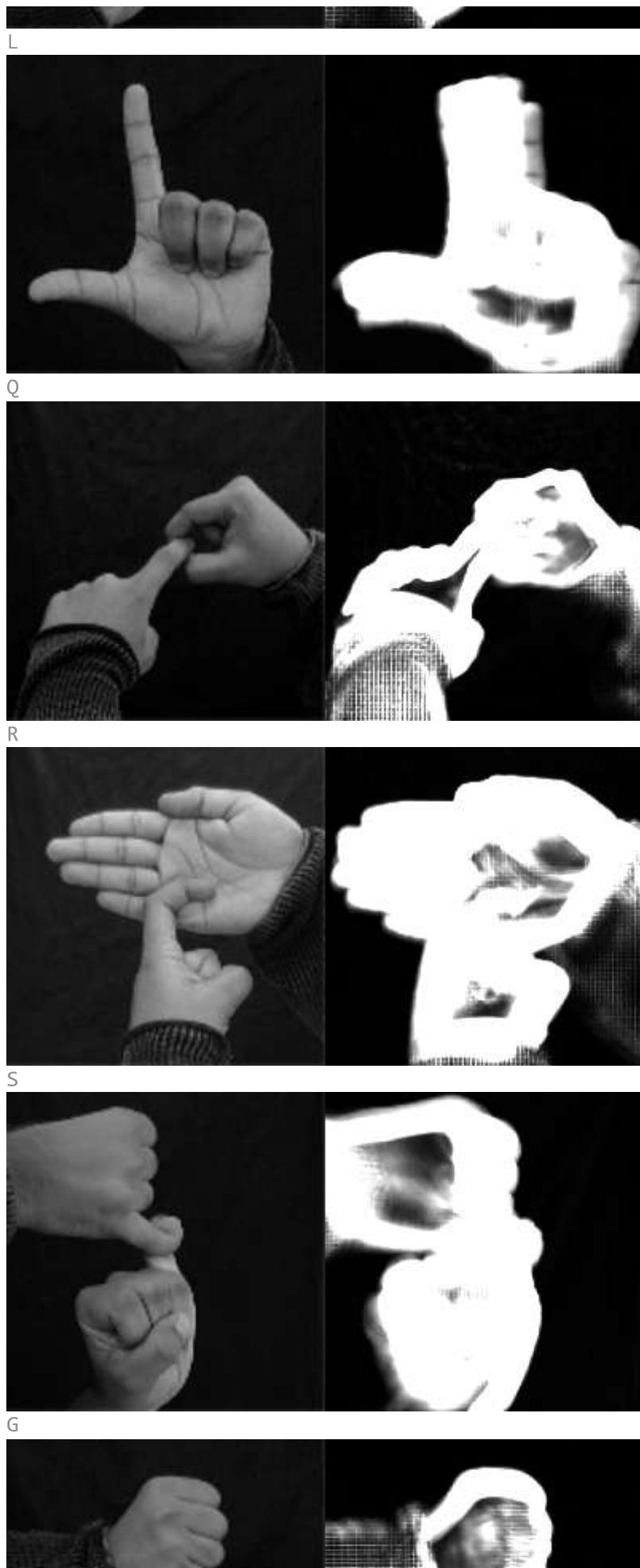


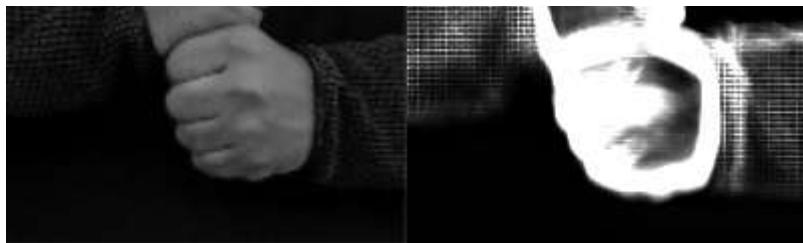
7



K



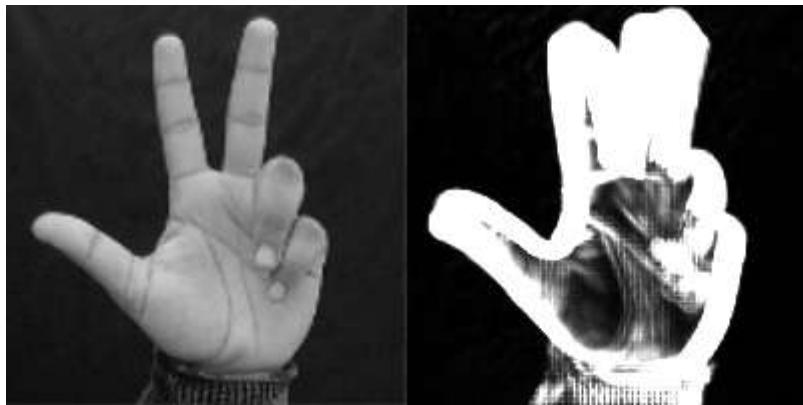




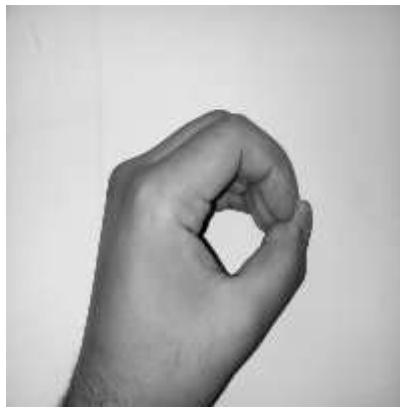
X



8

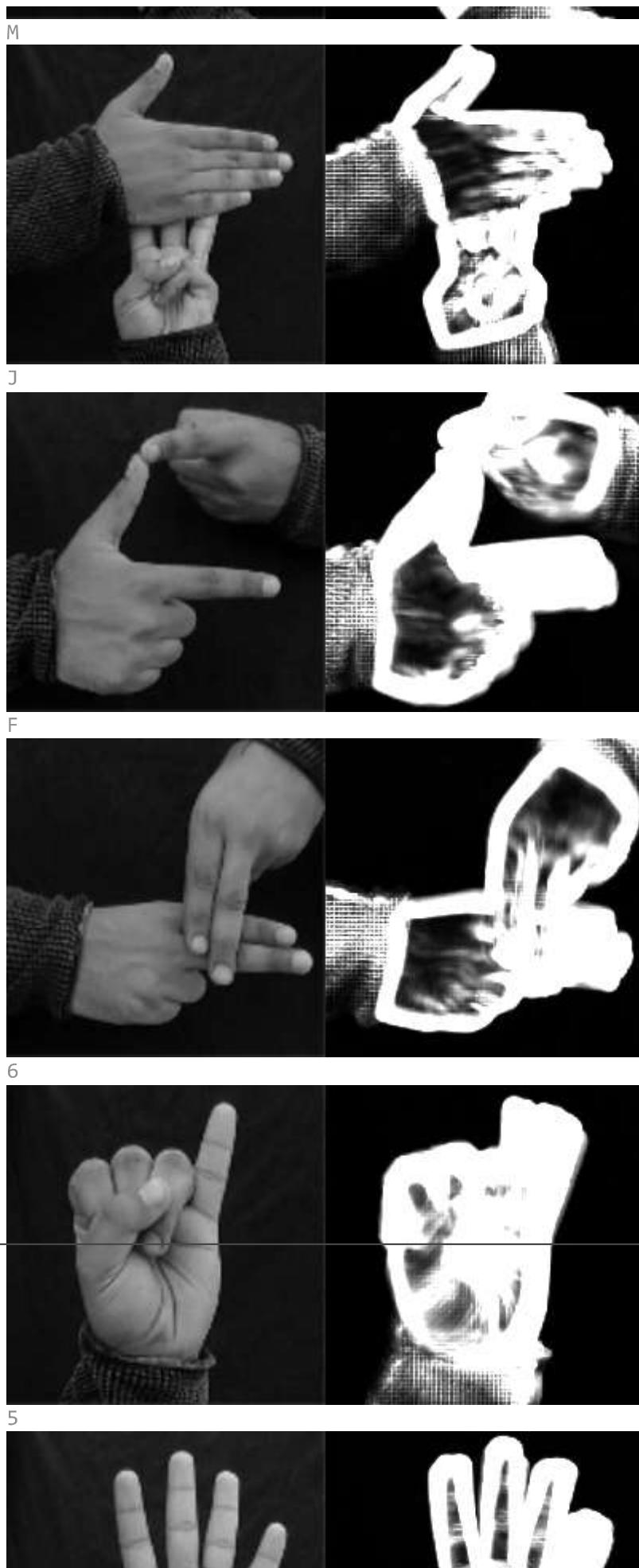


0

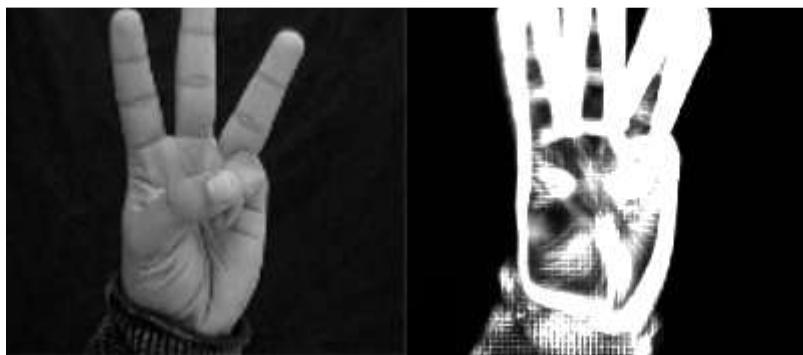


N





```
1 #Getting standard deviation image of each class
2 indices=[]
3 distinct_labels = set()
4 for i in range(len(y)):
5     if y[i] not in distinct_labels:
6         distinct_labels.add(y[i])
7         indices.append(i)
8 #print(indices)
9 print("      Sample images      SD images of each class")
10 for i in range(len(indices)-1):
11     #cv2_imshow(X[indices[i]])
12     print(y[indices[i]])
13     sd_img = np.std(images[indices[i]:indices[i+1]],axis=0)
14     sd_img = sd_img.reshape((200,200,1))
15     #cv2_imshow(mean_img)
16     img_concate_Hori=np.concatenate((images[indices[i]],sd_img),axis=1)
17     cv2_imshow(img_concate_Hori)
18
19
20
21 print(y[indices[len(indices)-1]])
22 sd_img = np.std(images[indices[len(indices)-1]:len(images)-1],axis=0)
23 sd_img = sd_img.reshape((200,200,1))
24 #cv2_imshow(X[indices[len(indices)-1]])
25 #mean_img = np.mean(X[len(indices)-1:len(X)])
26 img_concate_Hori=np.concatenate((images[indices[len(indices)-1]],sd_img),axis=1)
27 cv2_imshow(img_concate_Hori)
28
29
30
```





Sample images

SD images of each class

H



4



D



T



P





B



1



Y



U



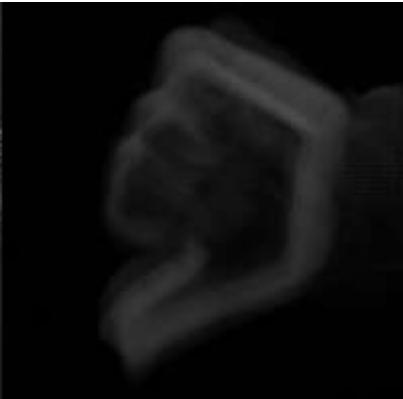




2



9

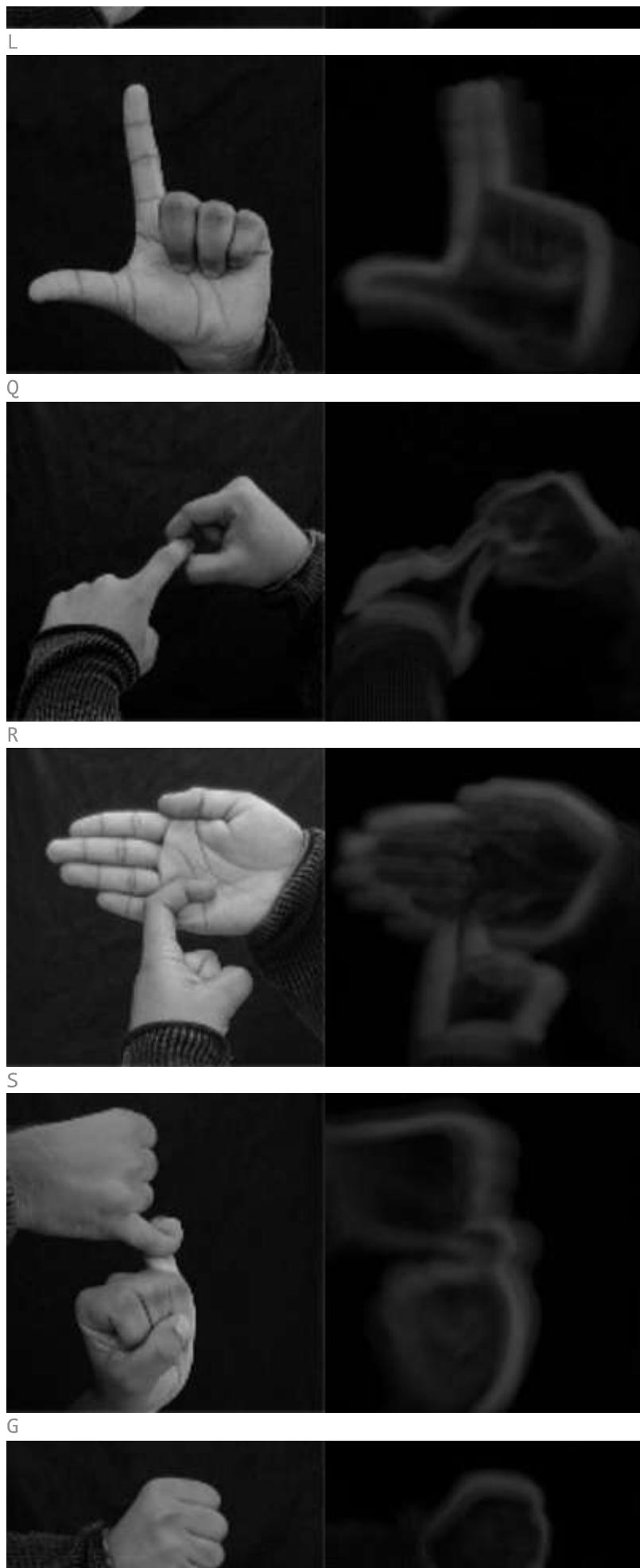


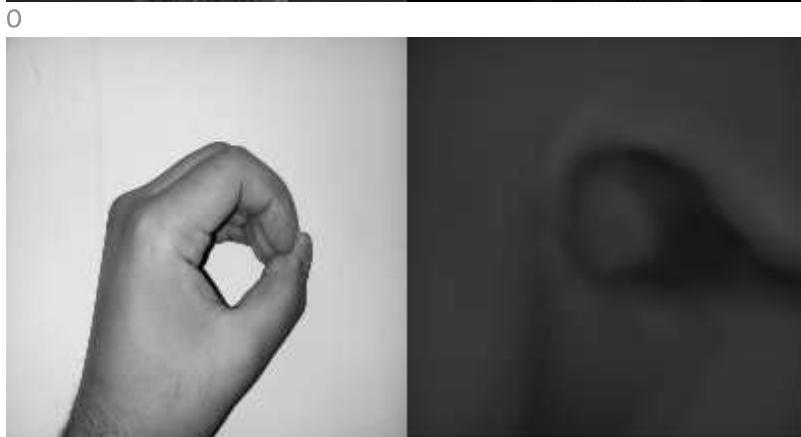
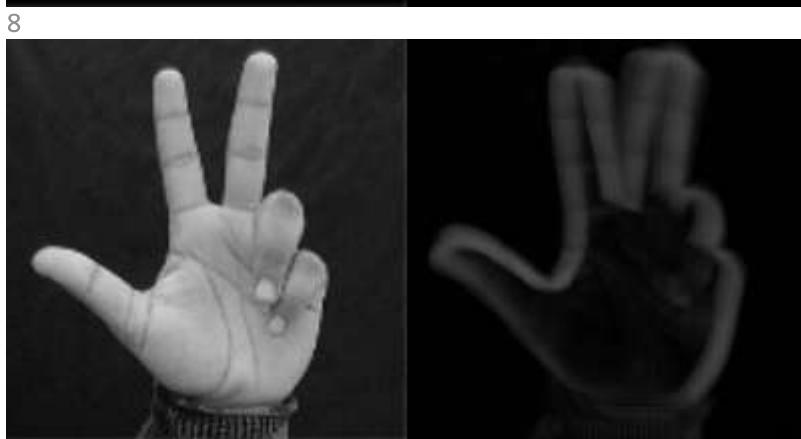
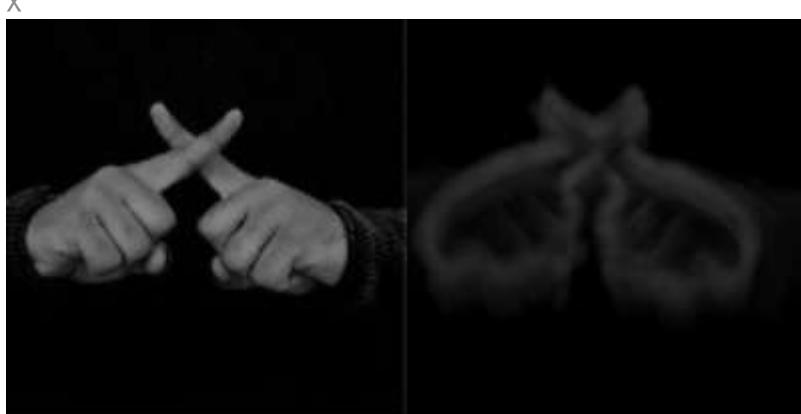
7



K







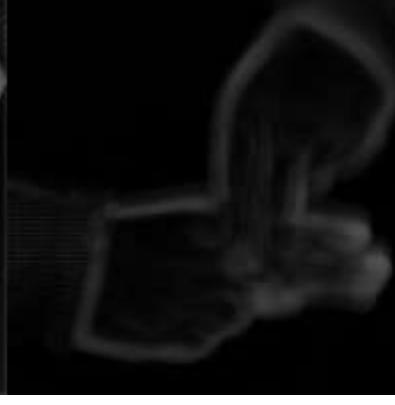
M



J



F



6



5



# Feature Extraction using ORB

```
1 #Extracting the keypoints and their descriptors
2 keypoints=[]
3 descriptors=[]
4 for i in range(len(images)):
5     img = images[i]
6     orb = cv2.ORB_create(nfeatures=2000)
7     kp = orb.detect(img,None)
8     kp,des = orb.compute(img,kp)
9     keypoints.append(kp)
10    descriptors.append(des)
11 print(len(keypoints))
12 #print(keypoints[0])
13 print(len(descriptors))
14 #print(descriptors[0])
```

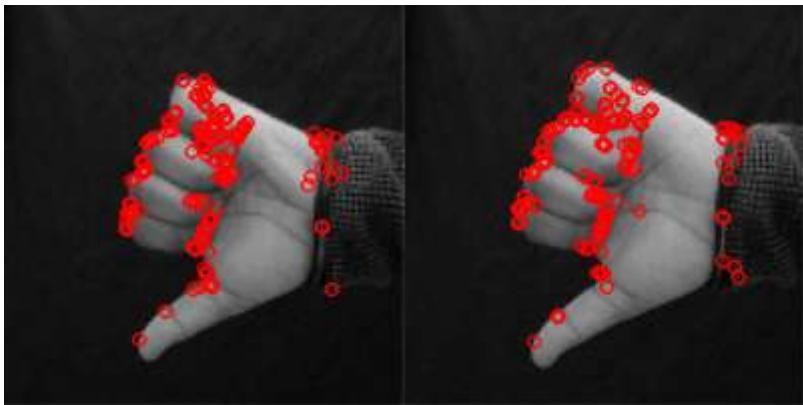
42745  
42745

```
1 #displaying the keypoints for a few random images and another sample from t
2 for num in range(10):
3     rand = random.randint(0,len(images))
4     print(y[rand])
5     op_img1 = cv2.drawKeypoints(images[rand],keypoints[rand],None,color=(0,0,
6     op_img2 = cv2.drawKeypoints(images[rand+1],keypoints[rand+1],None,color=(
7     #cv2_imshow(op_img1)
8     #op_img = op_img.reshape((200,200,1))
9     img_concat = np.concatenate((op_img1,op_img2),axis=1)
10    cv2_imshow(img_concat)
11
```

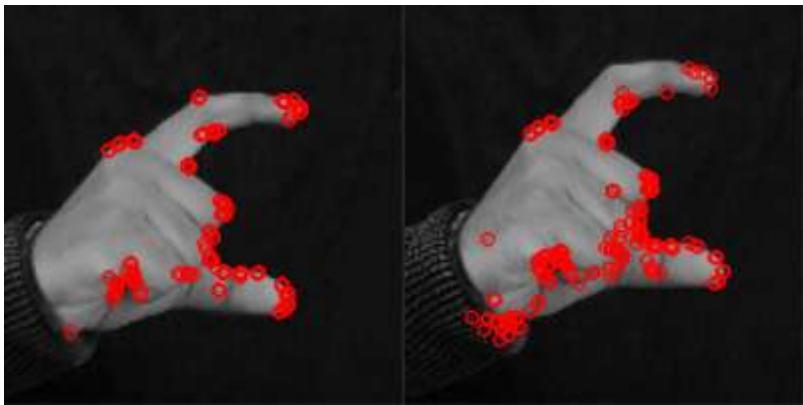




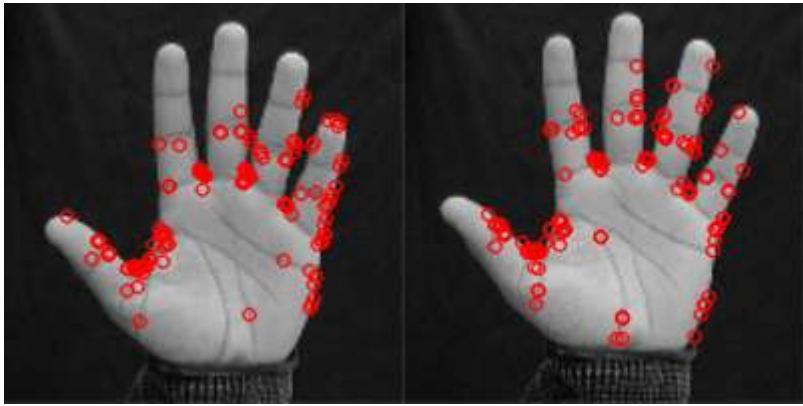
9



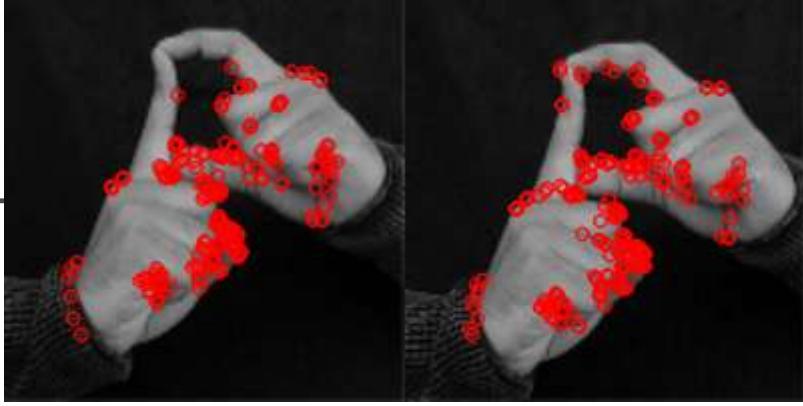
C



5



D

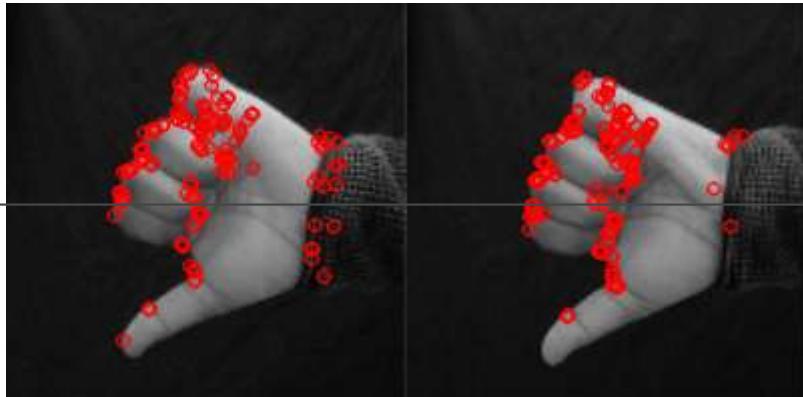


F

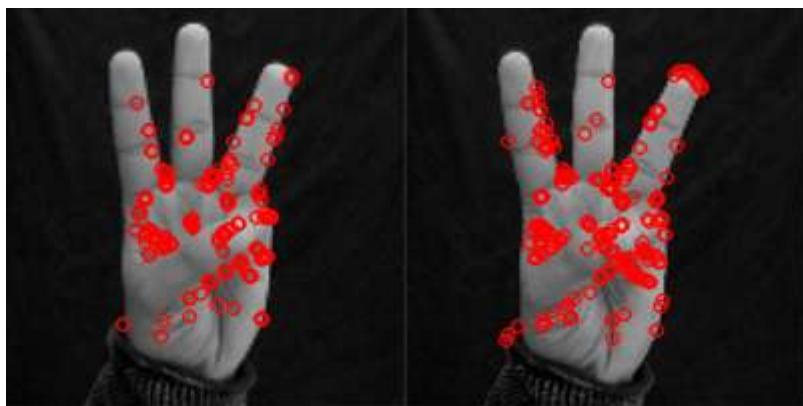




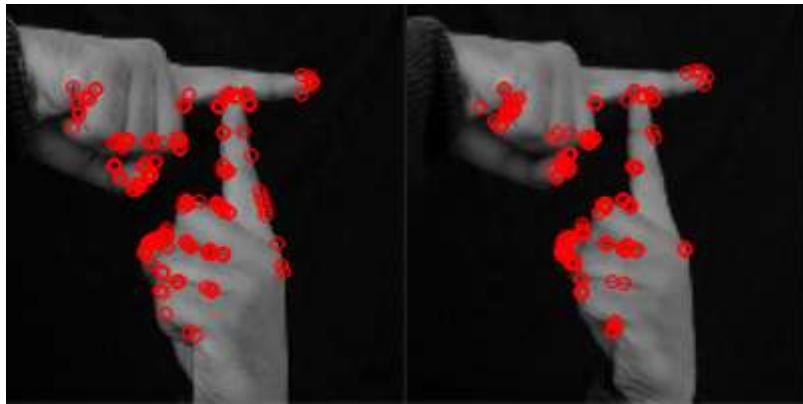
9



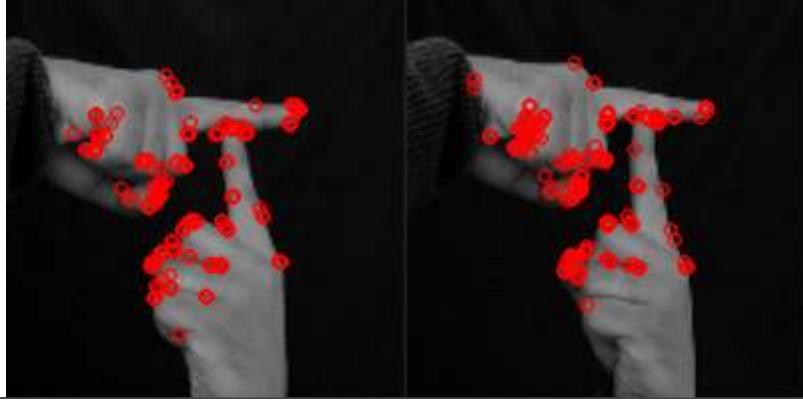
3



T



T



5

## Train Test Split Indexes



```

1 def train_test_val_split_idxs(total_rows, percent_test, percent_val):
2     if percent_test + percent_val >= 1.0:
3         raise ValueError('percent_test and percent_val must sum to less than 1.0')
4     row_range = range(total_rows)
5     no_test_rows = int(total_rows*(percent_test))
6     test_idxs = np.random.choice(row_range, size=no_test_rows, replace=False)
7     row_range = [idx for idx in row_range if idx not in test_idxs]
8     no_val_rows = int(total_rows*(percent_val))
9     val_idxs = np.random.choice(row_range, size=no_val_rows, replace=False)
10    training_idxs = [idx for idx in row_range if idx not in val_idxs]
11
12    print('Train-test-val split: %i training rows, %i test rows, %i validation rows' % (len(training_idxs), len(test_idxs), len(val_idxs)))
13
14    return training_idxs, test_idxs, val_idxs
15
16 training_idxs, test_idxs, val_idxs = train_test_val_split_idxs(len(descriptors), 0.2, 0.1)

```

Train-test-val split: 25647 training rows, 17098 test rows, 0 validation rows

## Bag of Visual Words

```

1 def cluster_features(img_descs, training_idxs, cluster_model):
2
3     n_clusters = cluster_model.n_clusters
4     training_descs = [img_descs[i] for i in training_idxs]
5     all_train_descriptors = []
6     for desc_list in training_descs:
7         if desc_list is not None:
8             for desc in desc_list:
9                 all_train_descriptors.append(desc)
10    all_train_descriptors = [desc for desc_list in training_descs for desc in desc_list]
11    all_train_descriptors = np.array(all_train_descriptors)
12
13    print ('%i descriptors before clustering' % all_train_descriptors.shape)
14
15
16    print ('Using clustering model %s...' % repr(cluster_model))
17    print ('Clustering on training set to get codebook of %i words' % n_clusters)
18    cluster_model.fit(all_train_descriptors)
19    print ('done clustering. Using clustering model to generate Bow histograms')
20    new = []
21    for i in img_descs:

```

```

22     if i is not None:
23         new.append(i)
24
25     img_clustered_words = [cluster_model.predict(raw_words) for raw_words in
26
27     img_bow_hist = np.array(
28         [np.bincount(clustered_words, minlength=n_clusters) for clustered_words in
29
30     X = img_bow_hist
31     print ('done generating BoW histograms.')
32
33     return X, cluster_model
34 X, cluster_model = cluster_features(descriptors, training_idxs, MiniBatchKMeans())

```

4399616 descriptors before clustering  
Using clustering model MiniBatchKMeans(n\_clusters=150)...  
Clustering on training set to get codebook of 150 words  
done clustering. Using clustering model to generate BoW histograms for each image  
done generating BoW histograms.

## Train Test Split

```

1 def perform_data_split(X, y, training_idxs, test_idxs):
2     temp1=[]
3     temp2 =[]
4     for i in training_idxs:
5         if(i<len(X)):
6             temp1.append(i)
7     for i in test_idxs:
8         if(i<len(X)):
9             temp2.append(i)
10
11     X_train = X[temp1]
12     y_train = y[temp1]
13
14     X_test = X[temp2]
15     y_test = y[temp2]
16
17
18     return X_train, X_test, y_train, y_test
19 X_train, X_test, y_train, y_test = perform_data_split(X, y, training_idxs,

```

```

1 #shuffling the data to ensure better training
2 X_train,y_train = shuffle(X_train,y_train)
3 X_test,y_test = shuffle(X_test,y_test)

```

# Evaluation Metrics

```
1 beforepca_accuracy = []
2 beforepca_precision = []
3 beforepca_f1 = []
4 beforepca_recall = []
```

```
1 afterpca_accuracy = []
2 afterpca_precision = []
3 afterpca_f1 = []
4 afterpca_recall = []
```

```
1 def calc_accuracy(method,label_test,pred):
2
3     acc = sm.accuracy_score(label_test,pred)
4     print("accuracy score for ",method,acc)
5     beforepca_accuracy.append(acc)
6
7     pre = sm.precision_score(label_test,pred,average='macro')
8     print("precision_score for ",method,pre)
9     beforepca_precision.append(pre)
10
11    fone = sm.f1_score(label_test,pred,average='macro')
12    print("f1 score for ",method,fone)
13    beforepca_f1.append(fone)
14
15    rec = sm.recall_score(label_test,pred,average='macro')
16    print("recall score for ",method,rec)
17    beforepca_recall.append(rec)
18
```

```
1 def calc_accuracy_pca(method,label_test,pred):
2
3     acc = sm.accuracy_score(label_test,pred)
4     print("accuracy score for ",method,acc)
5     afterpca_accuracy.append(acc)
6
7     pre = sm.precision_score(label_test,pred,average='macro')
8     print("precision_score for ",method,pre)
9     afterpca_precision.append(pre)
10
11    fone = sm.f1_score(label_test,pred,average='macro')
12    print("f1 score for ",method,fone)
13    afterpca_f1.append(fone)
14
15    rec = sm.recall_score(label_test,pred,average='macro')
```

```

16     print("recall score for ",method,rec)
17     afterpca_recall.append(rec)

```

## Models before PCA

```

1 def predict_svm(X_train, X_test, y_train, y_test):
2     svc=SVC(kernel='linear')
3     print("svm started")
4     svc.fit(X_train,y_train)
5     y_pred=svc.predict(X_test)
6     calc_accuracy("SVM",y_test,y_pred)
7 predict_svm(X_train, X_test,y_train, y_test)

```

```

svm started
accuracy score for SVM 0.9989471837164415
precision_score for SVM 0.9989730849493961
f1 score for SVM 0.9989719815393505
recall score for SVM 0.9989725108346736

```

```

1 def predict_lr(X_train, X_test, y_train, y_test):
2     clf = lr()
3     print("lr started")
4     clf.fit(X_train,y_train)
5     y_pred=clf.predict(X_test)
6     calc_accuracy("Logistic regression",y_test,y_pred)
7 predict_lr(X_train, X_test,y_train, y_test)

```

```

lr started
accuracy score for Logistic regression 0.9988886939229105
precision_score for Logistic regression 0.998920243912007
f1 score for Logistic regression 0.9989216491596813
recall score for Logistic regression 0.9989252652609963
/usr/local/lib/python3.7/dist-packages/scikit-learn/_linear_model/logistic.py:818: ConvergenceWarning: STOP: TOTAL NO. OF ITERATIONS REACHED LIMIT.

```

Increase the number of iterations (max\_iter) or scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html>

Please also refer to the documentation for alternative solver options:

[https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression)  
extra\_warning\_msg=\_LOGISTIC\_SOLVER\_CONVERGENCE\_MSG,

```

1 def predict_nb(X_train, X_test, y_train, y_test):
2     clf = nb()
3     print("nb started")
4     clf.fit(X_train,y_train)
5     y_pred=clf.predict(X_test)
6     calc_accuracy("Naive Bayes",y_test,y_pred)
7 predict_nb(X_train, X_test,y_train, y_test)

```

nb started  
accuracy score for Naive Bayes 0.9916359595250629  
precision\_score for Naive Bayes 0.9922005444950185  
f1 score for Naive Bayes 0.9921034124468571  
recall score for Naive Bayes 0.9922439188894564

```

1 def predict_knn(X_train, X_test, y_train, y_test):
2     clf=knn(n_neighbors=8)
3     print("knn started")
4     clf.fit(X_train,y_train)
5     y_pred=clf.predict(X_test)
6     calc_accuracy("K nearest neighbours",y_test,y_pred)
7 predict_knn(X_train, X_test,y_train, y_test)

```

knn started  
accuracy score for K nearest neighbours 0.9989471837164415  
precision\_score for K nearest neighbours 0.9989603780712668  
f1 score for K nearest neighbours 0.9989671554355628  
recall score for K nearest neighbours 0.9989763643927331

```

1 def predict_mlp(X_train, X_test, y_train, y_test):
2     clf=mlp(random_state=0)
3     print("mlp started")
4     clf.fit(X_train,y_train)
5     y_pred=clf.predict(X_test)
6     calc_accuracy("MLP classifier",y_test,y_pred)
7 predict_mlp(X_train, X_test,y_train, y_test)

```

mlp started  
accuracy score for MLP classifier 0.9989471837164415  
precision\_score for MLP classifier 0.9989603742903763  
f1 score for MLP classifier 0.9989616189942739  
recall score for MLP classifier 0.9989644531108172

```

1 for i in range(5):
2     beforepca_accuracy[i]*=10000
3     beforepca_precision[i]*=10000
4     beforepca_f1[i]*=10000
5     beforepca_recall[i]*=10000

```

# PCA

```
1 from sklearn.decomposition import PCA
2 pca = PCA(n_components = 25)
3 X_train_pca = pca.fit_transform(X_train)
```

## Models after PCA

```
1 def predict_svm(X_train_pca, X_test, y_train, y_test):
2     svc=SVC(kernel='linear')
3     print("svm started")
4     X_test_pca = pca.transform(X_test)
5     svc.fit(X_train_pca,y_train)
6     y_pred=svc.predict(X_test_pca)
7     calc_accuracy_pca("SVM",y_test,y_pred)
8 predict_svm(X_train_pca, X_test,y_train, y_test)
```

svm started  
accuracy score for SVM 0.995964204246359  
precision\_score for SVM 0.9960886410377979  
f1 score for SVM 0.9960959556484338  
recall score for SVM 0.9961145972035224

```
1 def predict_lr(X_train, X_test, y_train, y_test):
2     clf_pca = lr()
3     print("lr started")
4     X_test_pca = pca.transform(X_test)
5     clf_pca.fit(X_train_pca,y_train)
6     y_pred_pca=clf_pca.predict(X_test_pca)
7     calc_accuracy_pca("Logistic regression",y_test,y_pred_pca)
8 predict_lr(X_train_pca, X_test,y_train, y_test)
```

lr started  
accuracy score for Logistic regression 0.9953793063110488  
precision\_score for Logistic regression 0.9955200711205675  
f1 score for Logistic regression 0.9955134454976512  
recall score for Logistic regression 0.9955150238138464  
/usr/local/lib/python3.7/dist-packages/sklearn/linear\_model/\_logistic.py:818: Co  
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max\_iter) or scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html>

Please also refer to the documentation for alternative solver options:

[https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression)  
extra\_warning\_msg=\_LOGISTIC\_SOLVER\_CONVERGENCE\_MSG,

```

1 def predict_nb(X_train, X_test, y_train, y_test):
2     clf_pca = nb()
3     print("nb started")
4     X_test_pca = pca.transform(X_test)
5     clf_pca.fit(X_train_pca,y_train)
6     y_pred_pca=clf_pca.predict(X_test_pca)
7     calc_accuracy_pca("Naive Bayes",y_test,y_pred_pca)
8 predict_nb(X_train_pca, X_test,y_train, y_test)

```

nb started  
accuracy score for Naive Bayes 0.9832719190501258  
precision\_score for Naive Bayes 0.9848901501897884  
f1 score for Naive Bayes 0.9843945340935885  
recall score for Naive Bayes 0.9841127561810826

```

1 def predict_knn(X_train, X_test, y_train, y_test):
2     clf_pca=knn(n_neighbors=8)
3     print("knn started")
4     X_test_pca = pca.transform(X_test)
5     clf_pca.fit(X_train_pca,y_train)
6     y_pred_pca=clf_pca.predict(X_test_pca)
7     calc_accuracy_pca("K nearest neighbours",y_test,y_pred_pca)
8 predict_knn(X_train_pca, X_test,y_train, y_test)

```

knn started  
accuracy score for K nearest neighbours 0.9987717143358484  
precision\_score for K nearest neighbours 0.9987868400669715  
f1 score for K nearest neighbours 0.9987888147576622  
recall score for K nearest neighbours 0.9987936237521077

```

1 def predict_mlp(X_train, X_test, y_train, y_test):
2     clf_pca=mlp(random_state=0)
3     print("mlp started")
4     X_test_pca = pca.transform(X_test)
5     clf_pca.fit(X_train_pca,y_train)
6     y_pred_pca=clf_pca.predict(X_test_pca)
7     calc_accuracy_pca("MLP classifier",y_test,y_pred_pca)
8 predict_mlp(X_train_pca, X_test,y_train, y_test)

```

mlp started  
accuracy score for MLP classifier 0.9977773878458209  
precision\_score for MLP classifier 0.9978160503621177  
f1 score for MLP classifier 0.9978157801995978  
recall score for MLP classifier 0.9978209727848361

```

1 for i in range(5):
2     afterpca_accuracy[i]*=10000
3     afterpca_precision[i]*=10000

```

```

4     afterpca_f1[i]*=10000
5     afterpca_recall[i]*=10000

```

## ▼ Visualization

### ▼ Before PCA

```

1 import matplotlib.cm as cm
2 import matplotlib.pyplot as plt
3 import math

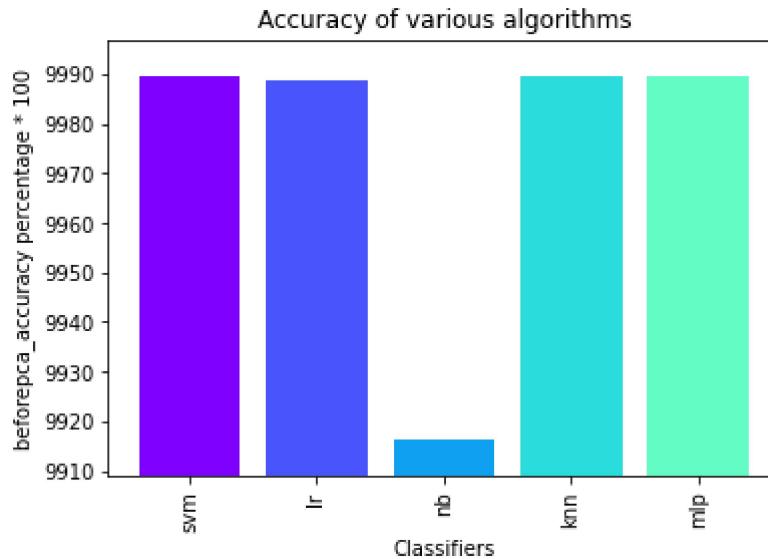
```

```

1 colors = cm.rainbow(np.linspace(0, 1, 10))
2 labels = ['svm', 'lr', 'nb', 'knn', 'mlp']
3 low = min(beforepca_accuracy)
4 high = max(beforepca_accuracy)
5 plt.ylim([math.floor(low-0.1*(high-low)), math.ceil(high+0.1*(high-low))])
6 plt.bar(labels,
7          beforepca_accuracy,
8          color = colors)
9 plt.xlabel('Classifiers')
10 plt.ylabel('beforepca_accuracy percentage * 100')
11 plt.title('Accuracy of various algorithms')
12 plt.xticks(rotation=90)

```

([0, 1, 2, 3, 4], <a list of 5 Text major ticklabel objects>)



```

1 svm_plot = []
2 lr_plot = []
3 nb_plot = []

```

```
4 knn_plot = []
5 mlp_plot = []
```

```
1 svm_plot.append(beforepca_accuracy[0])
2 svm_plot.append(beforepca_precision[0])
3 svm_plot.append(beforepca_f1[0])
4 svm_plot.append(beforepca_recall[0])
5
6 lr_plot.append(beforepca_accuracy[1])
7 lr_plot.append(beforepca_precision[1])
8 lr_plot.append(beforepca_f1[1])
9 lr_plot.append(beforepca_recall[1])
10
11 nb_plot.append(beforepca_accuracy[2])
12 nb_plot.append(beforepca_precision[2])
13 nb_plot.append(beforepca_f1[2])
14 nb_plot.append(beforepca_recall[2])
15
16 knn_plot.append(beforepca_accuracy[3])
17 knn_plot.append(beforepca_precision[3])
18 knn_plot.append(beforepca_f1[3])
19 knn_plot.append(beforepca_recall[3])
20
21 mlp_plot.append(beforepca_accuracy[4])
22 mlp_plot.append(beforepca_precision[4])
23 mlp_plot.append(beforepca_f1[4])
24 mlp_plot.append(beforepca_recall[4])
```

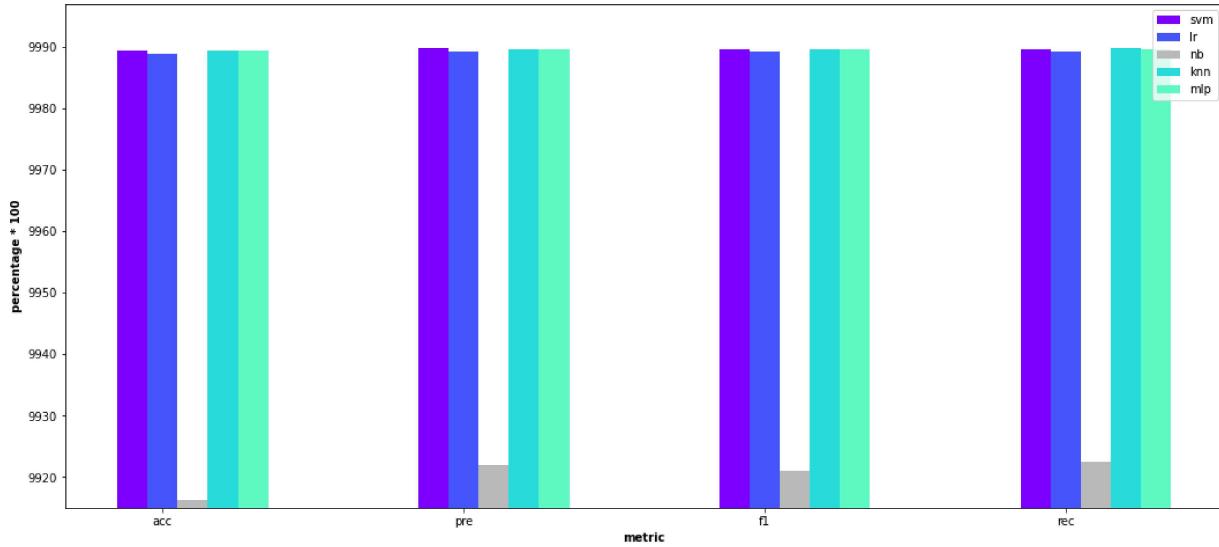
```
1 barWidth = 0.1
2 f, ax = plt.subplots(figsize=(18,8)) # set the size that you'd like (width,
3 # Set position of bar on X axis
4 pos1 = np.arange(4)
5 pos2 = [x + barWidth for x in pos1]
6 pos3 = [x + barWidth for x in pos2]
7 pos4 = [x + barWidth for x in pos3]
8 pos5 = [x + barWidth for x in pos4]
9
10 low = min(beforepca_recall)
11 high = max(beforepca_recall)
12 plt.ylim([math.floor(low-0.1*(high-low)), math.ceil(high+0.1*(high-low))])
13
14 # Make the plot
15 plt.bar(pos1, svm_plot, color='#8000FF', width=barWidth, label='svm')
16 plt.bar(pos2, lr_plot, color='#4856FB', width=barWidth, label='lr')
17 plt.bar(pos3, nb_plot, color='#BABABA', width=barWidth, label='nb')
18 plt.bar(pos4, knn_plot, color='#2ADDAA', width=barWidth, label='knn')
19 plt.bar(pos5, mlp_plot, color='#62FBC4', width=barWidth, label='mlp')
20
21 # Add xticks on the middle of the group bars
22 plt.xlabel('metric', fontweight='bold')
```

```

23 plt.ylabel('percentage * 100', fontweight='bold')
24 plt.xticks([i + barWidth for i in range(4)], ['acc', 'pre', 'f1', 'rec'])
25
26 ax.legend(fontsize = 14)
27 plt.legend(loc=1)

```

<matplotlib.legend.Legend at 0x7f7e1bab3f10>



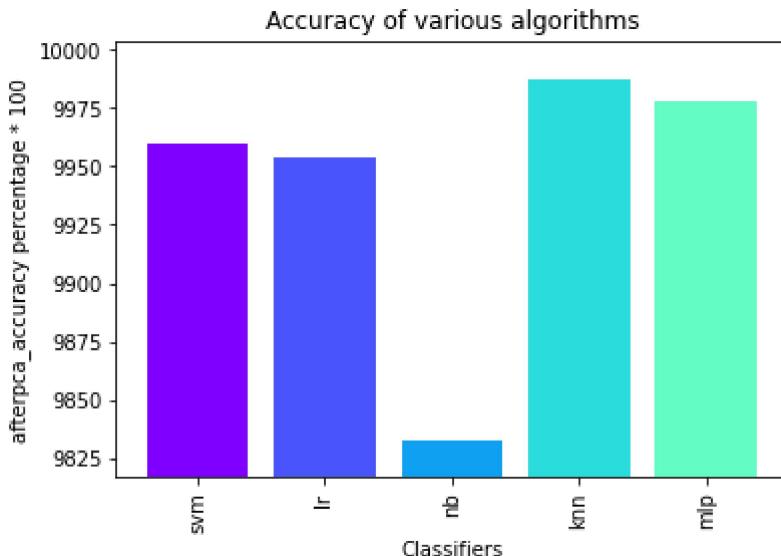
## ▼ After PCA

```

1 colors = cm.rainbow(np.linspace(0, 1, 10))
2 labels = ['svm', 'lr', 'nb', 'knn', 'mlp']
3 low = min(afterpca_accuracy)
4 high = max(afterpca_accuracy)
5 plt.ylim([math.floor(low-0.1*(high-low)), math.ceil(high+0.1*(high-low))])
6 plt.bar(labels,
7          afterpca_accuracy,
8          color = colors)
9 plt.xlabel('Classifiers')
10 plt.ylabel('afterpca_accuracy percentage * 100')
11 plt.title('Accuracy of various algorithms')
12 plt.xticks(rotation=90)

```

```
([0, 1, 2, 3, 4], <a list of 5 Text major ticklabel objects>)
```



```

1 svm_plot_after = []
2 lr_plot_after = []
3 nb_plot_after = []
4 knn_plot_after = []
5 mlp_plot_after = []

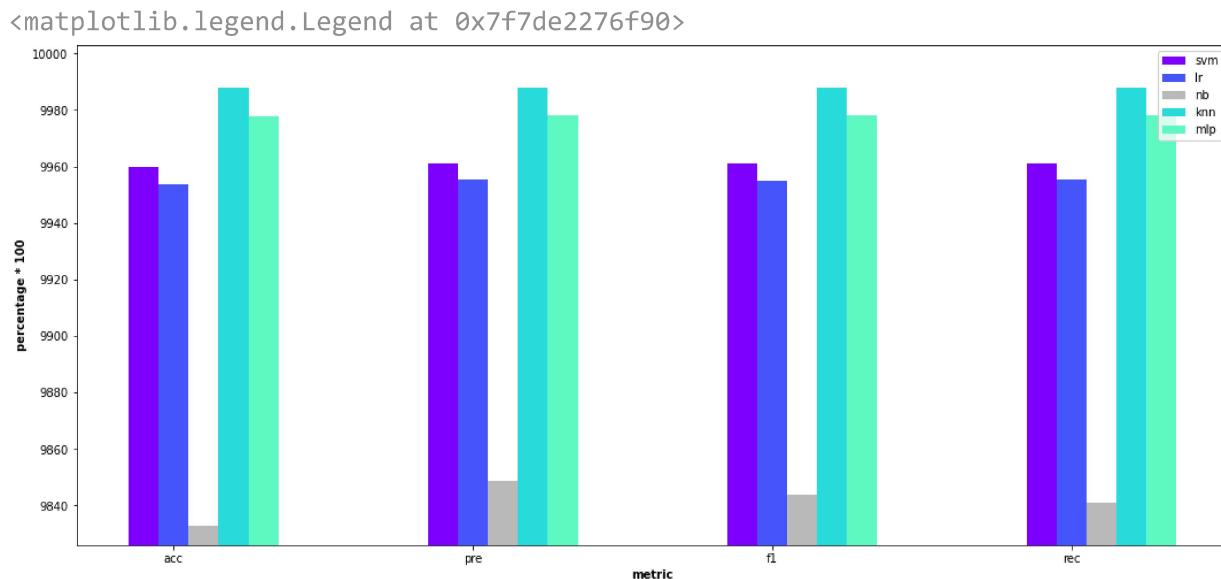
```

```

1 svm_plot_after.append(afterpca_accuracy[0])
2 svm_plot_after.append(afterpca_precision[0])
3 svm_plot_after.append(afterpca_f1[0])
4 svm_plot_after.append(afterpca_recall[0])
5
6 lr_plot_after.append(afterpca_accuracy[1])
7 lr_plot_after.append(afterpca_precision[1])
8 lr_plot_after.append(afterpca_f1[1])
9 lr_plot_after.append(afterpca_recall[1])
10
11 nb_plot_after.append(afterpca_accuracy[2])
12 nb_plot_after.append(afterpca_precision[2])
13 nb_plot_after.append(afterpca_f1[2])
14 nb_plot_after.append(afterpca_recall[2])
15
16 knn_plot_after.append(afterpca_accuracy[3])
17 knn_plot_after.append(afterpca_precision[3])
18 knn_plot_after.append(afterpca_f1[3])
19 knn_plot_after.append(afterpca_recall[3])
20
21 mlp_plot_after.append(afterpca_accuracy[4])
22 mlp_plot_after.append(afterpca_precision[4])
23 mlp_plot_after.append(afterpca_f1[4])
24 mlp_plot_after.append(afterpca_recall[4])

```

```
1 barWidth = 0.1
2 f, ax = plt.subplots(figsize=(18,8)) # set the size that you'd like (width,
3 # Set position of bar on X axis
4 pos1 = np.arange(4)
5 pos2 = [x + barWidth for x in pos1]
6 pos3 = [x + barWidth for x in pos2]
7 pos4 = [x + barWidth for x in pos3]
8 pos5 = [x + barWidth for x in pos4]
9
10 low = min(afterpca_recall)
11 high = max(afterpca_recall)
12 plt.ylim([math.floor(low-0.1*(high-low)), math.ceil(high+0.1*(high-low))])
13
14 # Make the plot
15 plt.bar(pos1, svm_plot_after, color='#8000FF', width=barWidth, label='svm')
16 plt.bar(pos2, lr_plot_after, color='#4856FB', width=barWidth, label='lr')
17 plt.bar(pos3, nb_plot_after, color='#BABABA', width=barWidth, label='nb')
18 plt.bar(pos4, knn_plot_after, color='#2ADDAA', width=barWidth, label='knn')
19 plt.bar(pos5, mlp_plot_after, color='#62FBC4', width=barWidth, label='mlp')
20
21 # Add xticks on the middle of the group bars
22 plt.xlabel('metric', fontweight='bold')
23 plt.ylabel('percentage * 100', fontweight='bold')
24 plt.xticks([i + barWidth for i in range(4)], ['acc', 'pre', 'f1', 'rec'])
25
26 ax.legend(fontsize = 14)
27 plt.legend(loc=1)
```



## ▼ comparison

## ▼ Accuracy

```
1 ## Accuracy comparison
2
3 barWidth = 0.2
4 f, ax = plt.subplots(figsize=(18,8)) # set the size that you'd like (width, height)
5 # Set position of bar on X axis
6 pos1 = np.arange(5)
7 pos2 = [x + barWidth for x in pos1]
8
9 low = min(afterpca_accuracy)
10 high = max(beforepca_accuracy)
11 plt.ylim([math.floor(low-0.1*(high-low)), math.ceil(high+0.1*(high-low))])
12
13 # Make the plot
14 plt.bar(pos1, beforepca_accuracy, color='#8000FF', width=barWidth, label='beforePCA')
15 plt.bar(pos2, afterpca_accuracy, color='#2ADDAA', width=barWidth, label='afterPCA')
```

```
17 # Add xticks on the middle of the group bars
18 plt.xlabel('metric', fontweight='bold')
19 plt.ylabel('percentage * 100', fontweight='bold')
20 plt.xticks([i + barWidth for i in range(5)], ['svm', 'lr', 'nb', 'knn', 'mip'])
21
22 ax.legend(fontsize = 14)
23 plt.legend(loc=1)
```

```
<matplotlib.legend.Legend at 0x7f7de95cf10>
```

