

TP 1 - Introduction aux bibliothèques Python

Introduction

Le but de ce premier TP est de vous familiariser avec les bibliothèques Python de fouille données : `scikit-learn` (outils et algorithmes d'apprentissage automatique), `numpy` (pour les manipulations de nombres), `scipy` (pour les calculs scientifiques), `pandas` (pour manipuler les ensembles de données) et `matplotlib` (pour la visualisation).

1 Analyse descriptive des données

L'ensemble de données « Fruits » que vous allez étudier a été créé par Iain Murray de l'université d'Édimbourg. Il a été construit à partir de l'observation de quelques dizaines d'oranges, de mandarines, de citrons et de pommes de différentes variétés. Les mesures de chacun de ces fruits ont été enregistrées dans une table ; c'est un ensemble de taille réduite, établi à des fins pédagogiques.

Téléchargez le fichier `fruit_data_with_colors.txt` du site web du cours et mettez-le dans le répertoire de votre choix.

1. Vérifiez que

```
1 pip install -U scikit-learn sklearn
2 pip install -U pandas
3 pip install -U graphviz
4
```

2. Utilisez `pandas` pour lire le fichier et le convertir en *data frame*. Combien y a-t-il d'attributs ? Que représentent-ils ? Quel est leur type ? Quelle est la classe à prédire ? Combien y a-t-il d'instances dans le fichier ? Les classes sont-elles équilibrées quant au nombre d'instances ?

```
1 # read input text and put data inside a data frame
2 fruits = pd.read_table('fruit_data_with_colors.txt')
3 print(fruits.head())
4
5 # print nb of instances and features
6 print(fruits.shape)
7
8 # print feature types
9 print(fruits.dtypes)
10
11 # print balance between classes
12 print(fruits.groupby('fruit_name').size())
13
```

3. Examinez les données de manière graphique à l'aide d'un diagramme de dispersions (fonction `plotting.scatter_matrix()` de `pandas`) et d'histogrammes montrant la répartition selon un attribut et la classe.

```
1 pd.plotting.scatter_matrix(X, c = y, marker = 'o', s=40, hist_kwds={'bins':15}, figsize=(9,9), cmap
  = cmap)
2
3 for attr in feature_names:
4     pd.DataFrame({k: v for k, v in fruits.groupby('fruit_name')[attr]}).plot.hist(stacked=True)
5     plt.suptitle(attr)
6     plt.savefig('fruits_histogram_'+attr)
7
```

2 Prétraitement

Suivant les méthodes de fouille de données considérées, il peut être utile d'appliquer différents prétraitements. On s'intéresse ici à deux grandes classes de filtres : la discrétisation et la normalisation.

1. Appliquez à partir des données originelles deux variantes de discrétisation : une première reposant sur des intervalles de même taille (*equal-interval* avec la fonction `cut()` de `pandas`) ou de mêmes effectifs (*equal-sized bins* avec `qcut()`). En affichant les historiques correspondant à chacune des deux configurations à l'aide de `plot.bar()`, comparez les histogrammes obtenus pour l'attribut « mass ».

```
1 pd.cut(fruits[attr],10).value_counts(sort=False).plot.bar()
2 plt.suptitle('Histogram for '+attr+' discretized with equal-interval bins')
3
```

2. Reprenez les données originelles. Quel est l'effet du filtre de normalisation `preprocessing.MinMaxScaler` de `scikit-learn`? Dans quel cas d'utilisation est-il préférable d'appeler la fonction `transform()` de cet objet plutôt que `fit_transform()`?

3 Clustering¹

Le clustering peut être vu comme une méthode non supervisée qui cherche à regrouper les instances qui se ressemblent, le plus souvent après avoir fixé au préalable un certain nombre de groupes (ou *clusters*) à construire.

Plusieurs algorithmes de clustering sont disponibles dans `scikit-learn`. Dans cette séance, deux méthodes de clustering seront testées : l'une classique est basée sur l'algorithme de partitionnement des données *K-means* et est désignée par `KMeans`, l'autre est une méthode hiérarchique ascendante produisant un arbre donnant la proximité entre les instances et est appelée `AgglomerativeClustering` (pour ce dernier modèle, nous utiliserons en réalité des méthodes de bibliothèque `scipy` puisque nous nous contenterons d'un affichage).

1. Ouvrez le fichier de données en appliquant un filtre de normalisation. Appliquez `KMeans` sur ces données en fixant successivement le nombre de *clusters* à 3, 4 et 5. Examinez manuellement

1. Appelé aussi « classification » dans le cours.

les résultats du clustering en traçant les graphes comme suit. Quelle(s) classe(s) représente chacun des *clusters* ?

```

1 from mpl_toolkits.mplot3d import Axes3D
2 from sklearn.cluster import KMeans
3
4 # Plot clusters
5 lst_kmeans = [KMeans(n_clusters=n) for n in range(3,6)]
6 titles = [str(x)+' clusters' for x in range(3,6)]
7 fignum = 1
8 for kmeans in lst_kmeans:
9     fig = plt.figure(fignum, figsize=(8, 6))
10    ax = Axes3D(fig, rect=[0, 0, .95, 1], elev=48, azim=134)
11    kmeans.fit(X_norm)
12    labels = kmeans.labels_
13    ax.scatter(X['mass'], X['width'], X['color_score'],
14              c=labels.astype(np.float), edgecolor='k')
15
16    ax.w_xaxis.set_ticklabels([])
17    ax.w_yaxis.set_ticklabels([])
18    ax.w_zaxis.set_ticklabels([])
19    ax.set_xlabel('mass')
20    ax.set_ylabel('width')
21    ax.set_zlabel('color_score')
22    ax.set_title(titles[fignum - 1])
23    ax.dist = 12
24    plt.savefig('k-means_'+str(2+fignum)+'_clusters')
25    fignum = fignum + 1
26    plt.close(fig)
27
28 # Plot the ground truth
29 fig = plt.figure(fignum, figsize=(8, 6))
30 ax = Axes3D(fig, rect=[0, 0, .95, 1], elev=48, azim=134)
31 for label in fruits['fruit_name'].unique():
32     ax.text3D(fruits.loc[fruits['fruit_name']==label].mass.mean(),
33              fruits.loc[fruits['fruit_name']==label].width.mean(),
34              fruits.loc[fruits['fruit_name']==label].color_score.mean(),
35              label,
36              horizontalalignment='center',
37              bbox=dict(alpha=.2, edgecolor='w', facecolor='w'))
38 ax.scatter(X['mass'], X['width'], X['color_score'], c=y, edgecolor='k')
39 ax.w_xaxis.set_ticklabels([])
40 ax.w_yaxis.set_ticklabels([])
41 ax.w_zaxis.set_ticklabels([])
42 ax.set_xlabel('mass')
43 ax.set_ylabel('width')
44 ax.set_zlabel('color_score')
45 ax.set_title('Ground Truth')
46 ax.dist = 12
47 plt.savefig('k-means_ground_truth')
48 plt.close(fig)
49 plt.show()
50

```

2. Pour fixer le nombre de clusters, on propose d'utiliser la métrique R^2 calculée suivant les variances intra-classe et inter-classe. Cette métrique permet d'avoir une évaluation objective

de l'homogénéité des groupes construits artificiellement, sans avoir étiqueté manuellement chaque instance. Un module `R_square_clustering.py` effectuant le calcul de R^2 vous est fourni sur le site du cours. Dessinez le graphe montrant l'évolution de cette métrique en faisant varier le nombre de *clusters* entre 2 et 10.

```

1  lst_k=range(2,11)
2  lst_rsqu = []
3  for k in lst_k:
4      est=KMeans(n_clusters=k)
5      est.fit(X_norm)
6      lst_rsqu.append(r_square(X_norm, est.cluster_centers_,est.labels_,k))
7
8  fig = plt.figure()
9  plt.plot(lst_k, lst_rsqu, 'bx-')
10 plt.xlabel('k')
11 plt.ylabel('RSQ')
12 plt.title('The Elbow Method showing the optimal k')
13

```

Comment varie R^2 ? Quel est le nombre de clusters que vous retiendriez ?

- Il est possible d'obtenir une évaluation quantitative plus complète du clustering lorsque l'on dispose des classes à obtenir dans les données (c'est un cas peu réaliste de l'apprentissage non supervisé mais ici nous sommes dans la situation luxueuse où chaque instance est étiquetée). Un module `purity.py` disponible sur le site du cours détermine un score de pureté en attribuant chaque cluster à la classe majoritaire. Utilisez-le pour tracer un graphe montrant l'évolution de ce score en fonction du nombre de *clusters* variant entre 2 et 10. Quel est le nombre de groupes que vous retiendriez suivant ce critère ?
- Remplacez l'utilisation du calcul des k-moyennes par une classification ascendante hiérarchique en établissant le dendrogramme des données normalisées comme suit.

```

1  from scipy.cluster.hierarchy import dendrogram, linkage
2
3  lst_labels = map(lambda pair: pair[0]+str(pair[1]), zip(fruits['fruit_name'].values, fruits.index))
4  linkage_matrix = linkage(X_norm, 'ward')
5  fig = plt.figure()
6  dendrogram(
7      linkage_matrix,
8      color_threshold=0,
9      labels=lst_labels,
10 )
11 plt.title('Hierarchical Clustering Dendrogram (Ward)')
12 plt.xlabel('sample index')
13 plt.ylabel('distance')
14 plt.tight_layout()
15

```

Commentez l'arbre obtenu et proposez une partition en sélectionnant une hauteur appropriée de coupure.

4 Classement

Un des intérêts de `scikit-learn` est d'incorporer un grand nombre d'algorithmes d'apprentissage automatique permettant de prédire la classe des individus. On se propose ici de tester quelques méthodes simples de classement selon deux protocoles d'évaluation.

1. Utilisez dans premier temps un découpage des données en $\frac{3}{4}$ entraînement et $\frac{1}{4}$ test à l'aide de la fonction `train_test_split` de `scikit-learn`. Comparez les performances des méthodes suivantes de classement :
 - classifieur élémentaire associant chaque instance à la classe majoritaire (`DummyClassifier`),
 - méthode bayésienne naïve (`GaussianNB`),
 - arbres de décision (`DecisionTreeClassifier`,
 - méthode de régression logistique (`LogisticRegression`).

Quel est le taux de classification obtenu pour chaque méthode ? Quelle est la classe pour laquelle la prédiction fait le plus d'erreurs ? Dans le cas de l'arbre de décision, visualisez et commentez l'arbre construit par le modèle.

```

1  for clf, name_clf in zip(lst_classif, lst_classif_names):
2      clf.fit(X_train, y_train)
3      y_pred = clf.predict(X_test)
4      print('Accuracy of '+name_clf+' classifier on training set: {:.2f}'
5            .format(clf.score(X_train, y_train)))
6      print('Accuracy of '+name_clf+' classifier on test set: {:.2f}'
7            .format(clf.score(X_test, y_test)))
8      print(confusion_matrix(y_test, y_pred))
9
10 # print decision tree
11 import graphviz
12 dot_data = tree.export_graphviz(dectree, out_file=None,
13                                feature_names=feature_names,
14                                class_names=fruits['fruit_name'].unique(),
15                                filled=True, rounded=True,
16                                special_characters=True)
17 graph = graphviz.Source(dot_data)
18 graph.render("fruits")
19

```

2. Utilisez maintenant une validation croisée en 5 strates à l'aide de la fonction `cross_val_score`. Parmi les quatre méthodes précédemment testées, quelle est la plus performante ? Le classement des méthodes est-il le même que lors du protocole d'évaluation précédent ?

5 Classement et discrétisation

On examine maintenant l'effet de la discrétisation sur la construction de modèles de classement.

1. Discrétisez les attributs numériques à l'aide d'intervalles de mêmes tailles. Appliquez les mêmes méthodes de classement que précédemment. Obtenez-vous les mêmes performances qu'avant la discrétisation ?

```

1  prefix = 'eqintevalued_bins_'
2  for attr in feature_names:
3      fruits[prefix+attr]=pd.cut(fruits[attr],10)

```

```

4      # use pd.concat to join the new columns with your original dataframe
5      fruits = pd.concat([fruits, pd.get_dummies(fruits[prefix+attr], prefix=prefix+attr)], axis=1)
6      # now drop the original column (you don't need it anymore)
7      fruits.drop(prefix+attr, axis=1, inplace=True)
8
9      feature_names_bins = filter(lambda x: x.startswith(prefix) and x.endswith(']'), list(fruits))
10     X_discret = fruits[feature_names_bins]

```

2. Mêmes questions que précédemment en utilisant cette fois-ci une discrétisation à l'aide d'intervalles de mêmes effectifs.

6 ACP et sélection de variables

L'analyse en composantes principales (ACP) peut être utilisée en méthode de prétraitements pour réduire le nombre de variables, ce qui est souvent très utile quand le nombre de dimensions des données est très grand (par exemple plusieurs dizaines). Dans cette section, on se propose d'étudier l'ACP sur les données et de comparer son influence sur les performances de classement par rapport à une sélection arbitraire ou par élimination récursive des variables.

1. Ouvrez le fichier de données en appliquant un filtre de normalisation puis une ACP. Quelles sont les valeurs propres des facteurs de l'ACP ? Dessinez le graphe montrant l'évolution des valeurs propres en fonction de leur rang. Quelles variables originelles sont corrélées avec chacun des nouveaux facteurs ? Tracez les instances dans le plan principal constitué des deux premiers facteurs. Les classes des fruits sont-ils bien séparés dans ce nouvel espace ?

```

1  from sklearn.decomposition import PCA
2
3  acp = PCA(svd_solver='full')
4  coord = acp.fit_transform(X_norm)
5  # explained variance scores
6  print(acp.explained_variance_ratio_)
7  # plot eigen values
8  n = np.size(X_norm, 0)
9  p = np.size(X_norm, 1)
10 eigval = float(n-1)/n*acp.explained_variance_
11 fig = plt.figure()
12 plt.plot(np.arange(1,p+1), eigval)
13 plt.title("Scree plot")
14 plt.ylabel("Eigen values")
15 plt.xlabel("Factor number")
16 plt.savefig('acp_eigen_values')
17 plt.close(fig)
18
19 # print eigen vectors
20 print(acp.components_)
21 # lines: factors
22 # columns: variables
23
24 # print correlations between factors and original variables
25 sqrt_eigval = np.sqrt(eigval)
26 corvar = np.zeros((p,p))
27 for k in range(p):
28     corvar[:, k] = acp.components_[k,:] * sqrt_eigval[k]

```

```
29 print(corvar)
30 # lines: variables
31 # columns: factors
32
33 # plot instances on the first plan ( first 2 factors )
34 fig, axes = plt.subplots( figsize =(12,12))
35 axes.set_xlim(-1,1)
36 axes.set_ylim(-1,1)
37 for i in range(n):
38     plt.annotate(y.values[i ],( coord[i ,0], coord[i ,1]) )
39 plt.plot ([ - 1,1],[0,0], color='silver ', linestyle ='-',linewidth=1)
40 plt.plot ([0,0],[ - 1,1], color='silver ', linestyle ='-',linewidth=1)
41 plt.savefig('acp_instances_1st_plan')
42 plt.close(fig)
```

2. Calculez les performances en validation croisée des quatres méthodes de classement utilisées dans la section 4 sur les données, en ne considérant non plus les quatre attributs initiaux mais les deux facteurs principaux de l'ACP. Obtenez-vous les mêmes performances qu'en utilisant deux attributs du fichier original normalisé (par exemple les deux premiers attributs) ?
3. Dans la question précédente, les deux attributs ont été choisis au hasard. `scikit-learn` met à disposition plusieurs méthodes de sélection de variables. Utilisez comme suit `RFECV` pour réaliser une élimination récursive des variables à l'aide d'un estimateur évaluant l'importance de chaque attribut. Quels sont les deux attributs jugés les plus importants par cette méthode ?

```
1 from sklearn.feature_selection import RFECV
2 from sklearn.svm import SVR
3 estimator = SVR(kernel="linear")
4 selector = RFECV(estimator=estimator, cv=5)
5 selector.fit(X_norm, y)
6 print("Optimal number of features: %d" % selector.n_features_)
7 print(selector.ranking_)
```

Calculez les performances en validation croisée des quatre classifieurs utilisés sur les deux meilleures variables retenues et comparez les résultats précédents calculés sur les deux premiers facteurs de l'ACP.