

TELECOM PARISTECH

APPRENTISSAGE POUR L'IMAGE ET LA RECONNAISSANCE
D'OBJETS

IMA 205

Rapport Challenge

Author :
de Sousa Silva, Raul Alfredo

Professor :
Gori, Pietro

Livré le : 21 avril 2019



Table des matières

1	Objectifs	2
2	Introduction	2
3	L'extraction des attributs	2
4	Mise en ouvre	4
4.1	Pré-traitement	4
4.2	KPPV	4
4.3	SVM	4
4.4	Arbres de décision	5
4.5	MLP	5
5	Résultats	6
5.1	KPPV	6
5.2	SVM	6
5.3	Arbres de décision	6
5.4	MLP	6
6	Conclusion	7
7	Références	8
	Annexes	8

1 Objectifs

Ce rapport a pour objectif de montrer les méthodes d'apprentissage qui ont été évalués pour l'exécution du challenge d'IMA 205, où l'objectif était de créer un algorithme capable de classer les images de lésions cutanées correctement entre cancérigènes et non-cancérigènes.

2 Introduction

Les méthodes pour apprentissage ont gagné beaucoup de place dans les dernières années surtout dans les domaines d'aide à un spécialiste, comme par exemple, le domaine médical.

Pour mettre en pratique certains de ces algorithmes, nous avons un challenge à exécuter à la fin de ce cours : Un ensemble de données avec 1000 images nous est mise à disposition, d'où 700 ont été déjà classées et 300 sont à classer à partir de l'apprentissage faite avec les premières.

On parle des images de lésions cutanées, entre elles nous avons des exemples d'images cancérigènes et non-cancérigènes. L'objectif c'est d'obtenir la meilleure classification possible de ces 300 dernières images en utilisant n'importe quelle des méthodes appris pendant le cours. La qualité de la classification est, alors, mesurée par la corrélation de Matthews, donné par :

$$MCC = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP + FP) \cdot (TP + FN) \cdot (TN + FP) \cdot (TN + FN)}} \quad (1)$$

Où,

- TP : proportion d'échantillons positives bien classés (comme positives) ;
- TN : proportion d'échantillons négatives bien classés (comme négatives) ;
- FP : proportion d'échantillons négatives mal classés (comme positives) ;
- FN : proportion d'échantillons positives mal classés (comme négatives) ;

Alors, pour essayer de faire le meilleur score possible seront testés 4 méthodes (les justificatives de pourquoi les essayer et que-est-ce que nous devons attendre sera explicité à chaque session) :

- K-plus proches voisins (KPPV) ;
- Support Vector Machines (SVM) ;
- Arbres de décision ;
- Multi-Layer Perceptron (MLP) ;

L'idée c'est d'augmenter la complexité de l'algorithme d'apprentissage peu à peu dans l'espoir d'avoir à chaque fois un algorithme plus performant et robuste, au fil des expérimentations. A la fin, nous allons discuter de la performance des différentes méthodes utilisées.

3 L'extraction des attributs

Pour que nous puissions faire la classification des images, il faut avoir des attributs que nous permettrons d'appliquer l'algorithme de classification. Pour l'instant, nous n'avons que des images des lésions et les cartes de segmentations. Alors, nous devons chercher certaines

caractéristiques qui puissent être intéressantes à détecter une lésion cancérigène (ou non-cancérigène).

Basée sur le travail de [1], nous avons obtenu 433 attributs liées à la forme, à la couleur et à la texture des images. En particulier, nous avons pris 13 attributs de forme, 348 attributs de couleur et 72 attributs de texture des 437 proposés par l'article original (certains attributs de forme ont été ajoutés à ceux qui ont été proposés par l'article dans l'espoir de créer des attributs non-linéaires qui puissent améliorer la classification). Alors, une liste un peu plus détaillée de ces attributs est présentée ci-dessous, la définition précise de chacune peut être facilement trouvée sur l'article original.

Attributs de forme

- Surface
- Périmètre
- Excentricité
- Asymétrie 1
- Longueur de l'axe majeur
- 3 moments d'ordre 2 ((0,2),(2,0),(1,1))
- Diamètre équivalente
- Ratio d'aspect
- Orientation de l'axe majeur
- Asymétrie 2
- Longueur de l'axe mineur

Attributs de couleur

Pour le calcul des attributs de couleur, l'image est projetée dans six espaces de couleurs : RGB, rgb (RGB normalisé), HSV, l1/2/3 (Ohta space), l1/2/3 and CIE L*u*v*. Dans chaque canal de chaque espace couleur ont été prises les attributs décrites ci-dessous. Alors, à l'aide de la carte de segmentation, les images sont découpées en trois parties d'intérêt : la lésion elle-même, la périphérie interne à la lésion et la périphérie externe à la lésion. Alors, nous avons $6 \text{ (espaces couleurs)} \times 3 \text{ (canaux)} \times 3 \text{ (régions d'intérêt)} = 54 \text{ régions d'où extraire des attributs}$.

- Moyenne et écart-type (D'ailleurs, les rapports et les différences des 2 statistiques sur les 3 régions ont également été calculés : (externe / interne), (externe / lésion), (interne / lésion), (externe - interne), (lésion - externe), et (lésion - interne)). (324)
- Asymétrie de couleur (Exclusivement en RGB) (6)
- Distance des centroïdes (1 pour chaque canal de chaque espace couleur) (18)

Attributs de texture

Afin de quantifier la texture présente dans une lésion, un ensemble de descripteurs statistiques de texture basées sur la co-occurrence de la matrice de niveau de gris (GLCM) ont été utilisés. Bien que de nombreuses statistiques puissent être dérivées du GLCM, huit statistiques invariantes par décalage de niveau de gris ont été utilisées. D'ailleurs, les rapports et les différences des statistiques sur les 3 régions ont également été calculés : (externe / interne), (externe / lésion), (interne / lésion), (externe - interne), (lésion - externe), et (lésion - interne).

- Probabilité maximale
- Énergie
- Entropie
- Dissimilabilité
- Contraste
- Différence inverse
- Moment de différence inverse
- Corrélation.

4 Mise en oeuvre

4.1 Pré-traitement

Pour permettre une classification adéquate sans être biaisé par l'extension où l'ordre de grandeur de chaque attribut, il faut avant tout, faire une normalisation des attributs. Cela est très facilement exécuté par les fonctions mises en place par la bibliothèque *sklearn*. Avec la fonction *StandardScaler* de *sklearn.preprocessing* la normalisation est faite en trois lignes de code.

Aussi pour accélérer l'entraînement et pour ne pas garder des attributs qui ne sont pas importants pour décrire les images, nous appliquons aussi une réduction de dimensionnalité, en particulier nous appliquons une Analyse en Composantes Principales (ACP) pour pouvoir obtenir les caractéristiques qui expliquent à 95% les échantillons de notre ensemble. Cette analyse est aussi facilement faite à partir de la fonction *PCA* de *sklearn.decomposition*.

4.2 KPPV

Pour la mise en oeuvre d'une classification avec KPPV nous n'avons pas grand-chose à faire. D'ailleurs, c'est une bonne référence pour évaluer la qualité des autres méthodes. La classification pour k plus proches voisins c'est aussi assez raisonnable, étant vu que nous attendions certainement que des images de lésion cancérigène (et non-cancérigènes) aient des caractéristiques en commun.

Alors, pour appliquer une classification pour K-plus proche voisins nous devons utiliser la bibliothèque *sklearn* qui contient entre autres méthodes de classification, le KPPV. Il nous reste donc d'appliquer la fonction *KNeighborsClassifier* aux échantillons avec le nombre de voisins désirés.

Pour améliorer la classification, nous pouvons faire une évaluation du score en fonction du hyperparamètre k (nombre de voisins) et ainsi, optimiser la classification. Cette amélioration est faite avec les fonctions *cross_val_score*, *GridSearchCV* and *KFold* de *sklearn.model_selection*.

À la fin, nous pouvons gagner un peu plus de robustesse à l'aide d'une technique de boosting, qui permet d'améliorer sensiblement la qualité de la classification par l'utilisation de plusieurs classificateurs faibles, concept vu en IMA206, qui vont, chacun contribuer pour la classification avec un vote. Pour des classificateurs simples comme les KPPV et arbres de décision, qui ne demandent pas un grand effort de calcul, une telle stratégie pourrait être assez efficace.

4.3 SVM

La méthode du SVM est beaucoup moins évidente et peut être un peu naïve si certains changements ne sont pas faits. En particulier, la classification par SVM suppose la possibilité de diviser l'espace avec une droite, ce qui n'est pas tout à fait vrai. Mais quand même nous pouvons trouver une droite qui minimise l'erreur de la classification de telle sorte qu'on puisse avoir une classification acceptable.

Pour appliquer le SVM nous utilisons aussi la bibliothèque *sklearn* avec la fonction *LinearSVC*. Le hyperparamètre C de régularisation du SVM (pour qu'il puisse admettre une solution où la division de l'espace ne soit pas parfaite) est réglé avec une validation croisée de la classification avec différentes valeurs du hyperparamètre ainsi comme en KPPV.

4.4 Arbres de décision

Les arbres de décision sont, grosso modo, un raffinement des KPPV de manière un peu plus élégante. Tandis que KPPV va créer des régions où un éventuel nouveau échantillon serait classé comme positif (ou négatif) à cause de la proximité plutôt des exemples positifs (ou négatifs), les arbres de décision vont diviser l'espace en plusieurs régions à partir de seuils sur l'un des attributs de tel manière à minimiser l'erreur de classification.

Alors, les mêmes possibilités d'améliorations du KPPV s'ouvrent aux arbres de décision. Les arbres pourraient être réglés par certains hyperparamètres comme le nombre minimum d'échantillons dans chaque "brique" de l'espace ou le nombre maximum de "feuilles", c'est-à-dire divisions subséquentes d'un sous espace. Aussi, le boosting pourrait améliorer la qualité de la classification par le comptage de votes de classification de plusieurs arbres entraînés différemment.

4.5 MLP

Le Perceptron Multi-couches, est une technique de Deep Learning qui a aussi une bonne chance de faire une bonne classification sur les attributs obtenus. Le principe du MLP est l'application de plusieurs couches de perceptrons, chaque couche de taille variée ou tous les perceptrons de la couche antécédente participent de la fonction d'activation de chaque perceptron de la couche actuelle.

L'objectif est d'entraîner les poids des perceptrons à chaque couche pour améliorer la classification de l'ensemble d'entraînement pour permettre d'améliorer la classification dans l'ensemble de test.

Alors, le MLP n'a pas forcément de hyperparamètres qui peuvent être choisis pour optimiser la classification, mais par contre, l'architecture du réseau est complètement libre, c'est-à-dire que nous pouvons faire n'importe quel nombre de couches et n'importe quelle taille à chaque couche, et même il y a plusieurs possibilités de non-linéarités qui doivent être bien choisies pour garantir de succès de l'algorithme.

Tout ce que nous savons, c'est que le plus grand est le nombre de couches et le plus grand est le nombre de perceptrons à chaque couche, le plus de non-linéarités nous avons, cela complexifie le réseau de neurones, mais augmente aussi la précision de la classification.

Pour mettre en place le MLP nous avons utilisé une architecture en basée sur la bibliothèque *tensorflow* avec deux couches cachées de 256 neurones chacune.

5 Résultats

5.1 KPPV

Avec la mise en place de la classification par KPPV nous faisons un entraînement avec les 700 échantillons du training set, avec tous les nombres impaires de voisins entre 1 et 15 pour pouvoir choisir le meilleur nombre de voisins nécessaires pour classer un échantillon.

La valeur de K obtenu a été de 9 voisins, avec un score *balanced_accuracy* de 0.64.

En appliquant au test set, le résultat obtenu dans la plateforme kaggle a été d'un score en Corrélacion de Matthews de autour 0.14 ce qui signifie (considérant une certaine homogénéité entre les classes) une accuracy de autour 0.57.

En faisant une implémentation de adaboost pour améliorer le score du KPPV avec 21 classifieurs, nous obtenons un *balanced_accuracy* de 0.99 tandis que sur kaggle la Corrélacion de Matthews a été de 0.2.

5.2 SVM

Avec la mise en place de la classification par SVM nous faisons un entraînement avec les 700 échantillons du training set, avec un facteur de régularisation de 'C' : [0.001, 0.01, 0.1, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10].

La valeur de C obtenu a été de 1, avec un score *balanced_accuracy* de 0.67.

En appliquant au test set, le résultat obtenu dans la plateforme kaggle a été d'un score en Corrélacion de Matthews de autour 0.4 ce qui signifie (considérant une certaine homogénéité entre les classes) une *balanced_accuracy* de autour 0.8.

Pour améliorer la classification, nous mettons en place aussi un SVM non-linéaire à l'aide de la fonction *SVC*. Avec elle et un autre paramètre de courbure 'gamma' ajouté à l'analyse nous arrivons à une *balanced_accuracy* de 0.67 dans le training et une Corrélacion de Matthews de autour 0.40, ce qui signifie une très légère amélioration par rapport au SVM linéaire. Meilleurs paramètres : {'C' : 7, 'gamma' : 0.001}.

5.3 Arbres de décision

Avec le mise en place de la classification par arbres de décisions, sans influence d'aucun hyperparamètre nous obtenons une *balanced_accuracy* de 1 dans le training, mais une Corrélacion de Matthews de autour 0.17 dans le test set selon la plateforme kaggle.

5.4 MLP

En utilisant l'architecture décrite dans la section antérieure, nous avons une apprentissage presque complète dans le training, (accuracy 0.96) mais un score en Corrélacion de Matthews de 0.37.

6 Conclusion

L’objectif de ce challenge, outre que familiariser les élèves avec des défis d’un travail en science de données, était de détecter les lésions cancérigènes parmi un ensemble de 1000 images de lésions de peau par l’utilisation de quelques algorithmes d’apprentissage parmi ceux appris dans le cours.

L’extraction d’attributs s’est montré comme un facteur crucial pour permettre une détection efficace des lésions cancérigènes, ce n’est pas toujours fondamental extraire beaucoup d’attributs, mais c’est fondamental extraire de bons attributs. Alors, nous devons admettre que les attributs extraits n’ont pas forcément été super efficaces et peut-être que le calcul des attributs pourrait être beaucoup plus rapide et souffler si une autre référence était utilisée.

Quant aux algorithmes d’apprentissage, ils ont eu une performance assez moyenne, peut-être limité par les attributs disponibles. Le KPPV nous a donné une bonne idée de base, mais il a montré rapidement qu’il fallait penser à quelque chose de plus complexe pour permettre une prédiction plus qualifiée. Le résultat du boosting effectué a amélioré beaucoup le taux d’apprentissage, mais il a peu développé en termes de prédiction.

L’utilisation des arbres de décision ont amélioré un peu la qualité de prédiction par rapport aux KPPV, mais il n’a même pas dépassé le KPPV “boosté”. Alors, même s’il a eu l’apprentissage la plus efficace, il a créé un overfitting sur l’espace des attributs, ce qui a fini pour limiter la performance de cet algorithme d’apprentissage.

Le SVM a eu alors un très bon résultat de classification, par rapport aux deux premiers. Pourtant, les tentatives d’amélioration de la classification avec le paramètre de courbure ont peu contribué pour l’optimisation du classifieur. Alors, c’est lui que nous gardons comme le meilleur classifieur pour le contexte du challenge.

Le réseau de neurone, le MLP a été peut-être un peu frustrante. Malgré son efficacité proche de celle du MLP (mais même pas supérieur), nous attendions quelque chose de plus robuste et efficace, étant vu qu’il a une tendance à contourner les non-linéarités éventuelles du système. Il est possible que l’architecture a été peu développée et que l’imposition d’autres couches pourrait beaucoup améliorer ce classifieur.

Pour conclure, nous avons eu une bonne occasion de mettre en oeuvre tout ce que nous avons appris dans de cours dans un contexte pratique. Nous sommes confrontés avec les défis et difficultés que nous aurons très probablement dans l’univers professionnel, alors, cela nous a permis de les avoir déjà en tête et concevoir certaines idées pour les résoudre partiel ou complètement.

7 Références

- [1] M. E. Celebi, H. A. Kingravi, B. Uddin, H. Iyatomi, Y. A. Aslandogan, W. V. Stoecker, R. H. Moss, "A methodological approach to the classification of dermoscopy images", *Comput. Med. Imag. Grap.*, vol. 31, no. 6, pp. 362-373, 2007.
- [2] Gómez, W., Pereira, W. & Infantosi, A. F. C. Analysis of co-occurrence texture statistics as a function of gray-level quantization for classifying breast ultrasound. *IEEE Trans. Med. Imag.* 31, 1889–1899 (2012).
- [3] R. Oliveira, J. Papa, A. Pereira, J. Tavares Computational methods for pigmented skin lesion classification in images : review and future trends *Neural Compu. Appl.* (2016), pp. 1-24
- [4] R. Garnavi, M. Aldeen, J. Bailey, "Computer-aided diagnosis of melanoma using border and wavelet-based texture analysis", *IEEE Trans. Inf. Technol. Biomed.*, vol. 16, no. 6, pp. 1239-1252, Nov. 2012.

Annexes

```

1  # -*- coding: utf-8 -*-
2  """
3  Created on Sun Mar  3 15:38:17 2019
4
5  @author: Raul Alfredo de Sousa Silva
6  Features extraction
7
8  """
9  # Imports (used libraries)
10
11  import cv2
12  import numpy as np
13  import matplotlib.pyplot as plt
14  import skimage.morphology as morpho
15  from skimage.transform import resize
16  from mpl_toolkits.axes_grid1 import AxesGrid
17
18
19  #####
20  # Basic function from ancient TPs of IMA coures
21  #####
22  def Get_values_without_error(im,XX,YY):
23      """ retouren une image de la taille de XX et YY
24          qui vaut im[XX,YY] mais en faisant attention a ce que XX et YY ne debordent
25          pas """
26      sh=XX.shape
27      defaultval=0;
28      if len(im.shape)>2: #color image !
29          defaultval=np.asarray([0,0,0])
30          sh=[*sh,im.shape[2]]
31      imout=np.zeros(sh)
32      (ty,tx)=XX.shape[0:2]
33      for k in range(ty):
34          for l in range(tx):
35              posx=int(XX[k,l]-0.5)
36              posy=int(YY[k,l]-0.5)
37              if posx<0 or posx>=im.shape[1] or posy<0 or posy>=im.shape[0]:
38                  valtmp=defaultval
39              else:
40                  valtmp=im[posy,posx]
41              imout[k,l]=valtmp
42
43      return imout
44
45  def rotation(im,theta,alpha=1.0,x0=None,y0=None,ech=0,clip=True):
46      """
47          %
48          %Effectue la transformation geometrique d'une image par
49          %une rotation + homothetie
50          %
51          % x' = alpha*cos(theta)*(x-x0) - alpha*sin(theta)*(y-y0) + x0
52          % y' = alpha*sin(theta)*(x-x0) + alpha*cos(theta)*(y-y0) + y0
53          %
54          % theta : angle de rotation en degres
55          % alpha : facteur d'homothetie (defaut=1)
56          % x0, y0 : centre de la rotation (defaut=centre de l'image)
57          % ech : plus proche voisin (defaut=0) ou bilineaire (1)
58          % clip : format de l'image originale (defaut=True), image complete (False)
59          %
60
61          """
62      dy=im.shape[0]
63      dx=im.shape[1]
64
65      if x0 is None:
66          x0=dx/2.0
67      if y0 is None:

```

```

68     y0=dy/2.0
69     v0=np.asarray([x0,y0]).reshape((2,1))
70     theta=theta/180*np.pi
71     ct=alpha*np.cos(theta)
72     st=alpha*np.sin(theta)
73     matdirect=np.asarray([[ct,-st],[st,ct]])
74     if clip==False:
75         #ON CALCULE exactement la transformee des positions de l'image
76         # on cree un tableau des quatre points extremes
77         tabextreme=np.asarray([[0,0,dx,dx],[0,dy,0,dy]])
78         tabextreme_trans= matdirect@(tabextreme-v0)+v0
79         xmin=np.floor(tabextreme_trans[0].min())
80         xmax=np.ceil(tabextreme_trans[0].max())
81         ymin=np.floor(tabextreme_trans[1].min())
82         ymax=np.ceil(tabextreme_trans[1].max())
83
84     else:
85         xmin=0
86         xmax=dx
87         ymin=0
88         ymax=dy
89     if len(im.shape)>2:
90         shout=(int(ymax-ymin),int(xmax-xmin),im.shape[2]) # image couleur
91     else:
92         shout=(int(ymax-ymin),int(xmax-xmin))
93     dyout=shout[0]
94     dxout=shout[1]
95     eps=0.0001
96     Xout=np.arange(xmin+0.5,xmax-0.5+eps)
97     Xout=np.ones((dyout,1))@Xout.reshape((1,-1))
98
99     Yout=np.arange(ymin+0.5,ymax-0.5+eps)
100    Yout=Yout.reshape((-1,1))@np.ones((1,dxout))
101
102    XY=np.concatenate((Xout.reshape((1,-1)),Yout.reshape((1,-1))),axis=0)
103    XY=np.linalg.inv(matdirect)@(XY-v0)+v0
104    Xout=XY[0,:].reshape(shout)
105    Yout=XY[1,:].reshape(shout)
106    if ech==0: # plus proche voisin
107        out=Get_values_without_error(im,Xout,Yout)
108    else: #bilineaire
109        assert ech == 1 , "Vous avez choisi un echantillonnage inconnu"
110        Y0=np.floor(Yout-0.5)+0.5 # on va au entier+0.5 inferieur
111        X0=np.floor(Xout-0.5)+0.5
112        Y1=np.ceil(Yout-0.5)+0.5
113        X1=np.ceil(Xout-0.5)+0.5
114        PoidsX=Xout-X0
115        PoidsY=Yout-Y0
116        PoidsX[X0==X1]=1 #points entiers
117        PoidsY[Y0==Y1]=1 #points entiers
118        I00=Get_values_without_error(im,X0,Y0)
119        I01=Get_values_without_error(im,X0,Y1)
120        I10=Get_values_without_error(im,X1,Y0)
121        I11=Get_values_without_error(im,X1,Y1)
122
123        out=I00*(1.0-PoidsX)*(1.0-PoidsY)+I01*(1-PoidsX)*PoidsY+I10*PoidsX*(1-PoidsY)+I11
124        *PoidsX*PoidsY
125    return out
126
127 def strel(forme,taille,angle=45):
128     """renvoie un element structurant de forme
129     'diamond' boule de la norme 1 fermee de rayon taille
130     'disk'     boule de la norme 2 fermee de rayon taille
131     'square'   carre de cote taille (il vaut mieux utiliser taille=impair)
132     'line'     segment de longueur taille et d'orientation angle (entre 0 et 180 en
133     degres)
134     (Cette fonction n'est pas standard dans python)

```

```

132     """
133
134     if forme == 'diamond':
135         return morpho.selem.diamond(taille)
136     if forme == 'disk':
137         return morpho.selem.disk(taille)
138     if forme == 'square':
139         return morpho.selem.square(taille)
140     if forme == 'line':
141         angle=int(-np.round(angle))
142         angle=angle%180
143         angle=np.float32(angle)/180.0*np.pi
144         x=int(np.round(np.cos(angle)*taille))
145         y=int(np.round(np.sin(angle)*taille))
146         if x**2+y**2 == 0:
147             if abs(np.cos(angle))>abs(np.sin(angle)):
148                 x=int(np.sign(np.cos(angle)))
149                 y=0
150             else:
151                 y=int(np.sign(np.sin(angle)))
152                 x=0
153             rr,cc=morpho.selem.draw.line(0,0,y,x)
154             rr=rr-rr.min()
155             cc=cc-cc.min()
156             img=np.zeros((rr.max()+1,cc.max()+1) )
157             img[rr,cc]=1
158             return img
159     raise RuntimeError('Erreur dans fonction strel: forme incomprie')
160
161 def extract(image,seg,img_id):
162     #####
163     # Color channels - 6 space colors
164     #####
165     # Creating channels in the 6 space colors
166
167     #RGB
168     R = image[:, :, 0]
169     G = image[:, :, 1]
170     B = image[:, :, 2]
171
172     #Normalized RGB
173     image_n = image.astype(np.uint16)
174     r = image_n[:, :, 0]/(image_n[:, :, 0]+image_n[:, :, 1]+image_n[:, :, 2]+0.00001)
175     g = image_n[:, :, 1]/(image_n[:, :, 0]+image_n[:, :, 1]+image_n[:, :, 2]+0.00001)
176     b = image_n[:, :, 2]/(image_n[:, :, 0]+image_n[:, :, 1]+image_n[:, :, 2]+0.00001)
177
178     #HSV
179     image_hsv = cv2.cvtColor(image, cv2.COLOR_RGB2HSV)
180     h = image_hsv[:, :, 0]
181     s = image_hsv[:, :, 1]
182     v = image_hsv[:, :, 2]
183
184     #I1/2/3
185     I1 = (R+G+B)/3
186     I2 = (R+B)/2
187     I3 = (2*G-R-B)/4
188     #image_I = np.transpose(np.array([I1,I2,I3]))
189
190     #l1/2/3
191     l1 = (R-G)**2/((R-G)**2+(R-B)**2+(G-B)**2+.000001)
192     l2 = (R-B)**2/((R-G)**2+(R-B)**2+(G-B)**2+.000001)
193     l3 = (G-B)**2/((R-G)**2+(R-B)**2+(G-B)**2+.000001)
194     #image_l = np.transpose(np.array([l1,l2,l3]))
195
196     # CIE l*a*b*
197     image_luv = cv2.cvtColor(image, cv2.COLOR_RGB2Luv)
198     l = image_luv[:, :, 0]

```

```

199     u = image_luv[:, :, 1]
200     v2 = image_luv[:, :, 2]
201
202     channels = [G, B, r, g, b, h, s, v, I1, I2, I3, l1, l2, l3, l, u, v2]
203     colors = np.expand_dims(R, axis=2)
204     for element in channels:
205         colors = np.append(colors, np.expand_dims(element, axis=2), axis=2)
206     del channels, R, G, B, r, g, b, h, s, v, I1, I2, I3, l1, l2, l3, l, u, v2
207     del image_hsv, image_luv, image_n
208     #####
209     # Shape Features - 13
210     #####
211     [m, n] = np.shape(seg)
212     ind = np.where(seg==1)
213     row = np.expand_dims(ind[0], axis=1)
214     col = np.expand_dims(ind[1], axis=1)
215
216     # First order geometric moments
217     m10 = sum(row)
218     m01 = sum(col)
219     m00 = sum(sum(seg))
220     r0 = m10/m00
221     c0 = m01/m00
222
223     # Second order central moments
224     mu02 = sum((col-c0)**2)
225     mu20 = sum((row-r0)**2)
226     mu11 = sum((col-c0)*(row-r0))
227     m10c = sum(row*colors[ind])
228     m01c = sum(col*colors[ind])
229     m00c = sum(colors[ind])
230     mu02c = sum((col-c0)**2*image[ind])
231     mu20c = sum((row-r0)**2*image[ind])
232     mu11c = sum((col-c0)*(row-r0)*image[ind])
233
234     # Perimeter
235     perimeter = 0
236     for k in range(0, len(ind[0])):
237         i = ind[0][k]
238         j = ind[1][k]
239         zeroc = max(0, j-1)
240         infic = min(n-1, j+1)
241         zeror = max(0, i-1)
242         infir = min(m-1, i+1)
243         if (seg[i, infic]+seg[i, zeroc]+seg[infir, j]+seg[zeror, j])<4:
244             perimeter+=1
245
246
247     # Aspect ratio
248     L1 = (8*(mu02+mu20 + ((mu02-mu20)**2 + 4*mu11)**(1/2) ))**(1/2)
249     L2 = (8*(mu02+mu20 - ((mu02-mu20)**2 + 4*mu11)**(1/2) ))**(1/2)
250     Ar = L1/L2
251     epsilon = ((mu02-mu20)**2 + 4*mu11)/((mu02-mu20)**2)
252     theta = 0.5*np.tan(2*mu11/(mu20-mu02))
253
254     mu20 = mu20
255     mu02 = mu02
256     mu11 = mu11
257     # Assimetry 1 and 2
258     # Rotation around the centroid of the lesion
259     boolean_x = rotation(seg, 180*theta[0]/np.pi, x0=round(c0[0]), y0=round(r0[0]))
260     boolean_x = boolean_x.astype(np.uint8) # Cast to integer
261     # number of pixels that belongs to both images (original an rotated)
262     diffx = sum(sum(boolean_x*seg))
263
264     boolean_y = rotation(seg, 180*theta[0]/np.pi-90, x0=round(c0[0]), y0=round(r0[0]))
265     boolean_y = boolean_y.astype(np.uint8) # Cast to integer

```

```

266         # number of pixels that belongs to both images (original an rotated)
267     diffy = sum(sum(boolean_y*seg))
268     Ax = m00-diffx
269     Ay = m00-diffy
270     A1 = min(Ax,Ay)/m00
271     A2 = (Ax+Ay)/m00
272
273     #Equivalent diameter parametrised by the axis L1 ond (or) L2
274     eqd = round(2*np.sqrt(m00/np.pi))
275     L1 = L1
276     L2 = L2
277     # Shape features
278     shapes =
[m00,eqd,perimeter,Ar[0],epsilon[0],theta[0],A1,A2,L1[0],L2[0],mu20[0],mu02[0],mu11[0]
]
279
280     # Grouping all features
281     metrics = [img_id]
282     for element in shapes:
283         metrics.append(element)
284
285     del shapes
286     #####
287     # Color features - 348
288     #####
289     se=strel('disk',round(0.05*eqd/2))
290     se2=strel('disk',round(0.1*eqd/2))
291     # dilatation
292     ignore=morpho.dilation(seg,se)
293     consider = morpho.dilation(seg,se2)
294     inner = consider-ignore
295     # erosion
296     ignore=morpho.erosion(seg,se)
297     consider = morpho.erosion(seg,se2)
298     outer = ignore-consider
299
300     seg_ex = np.expand_dims(seg, axis=2)
301     inner_ex = np.expand_dims(inner, axis=2)
302     outer_ex = np.expand_dims(outer, axis=2)
303
304     col_seg = (seg_ex*colors)
305     col_inner = (inner_ex*colors)
306     col_outer = (outer_ex*colors)
307
308     # 324 features of mean and standard deviation
309     for i in range(0,colors.shape[2]):
310         metrics.append(np.mean(col_seg[:, :, i]))
311         metrics.append(np.std(col_seg[:, :, i]))
312         metrics.append(np.mean(col_inner[:, :, i]))
313         metrics.append(np.std(col_inner[:, :, i]))
314         metrics.append(np.mean(col_outer[:, :, i]))
315         metrics.append(np.std(col_outer[:, :, i]))
316         metrics.append(np.mean(col_outer[:, :, i])/np.mean(col_inner[:, :, i]))
317         metrics.append(np.std(col_outer[:, :, i])/np.std(col_inner[:, :, i]))
318         metrics.append(np.mean(col_outer[:, :, i])/np.mean(col_seg[:, :, i]))
319         metrics.append(np.std(col_outer[:, :, i])/np.std(col_seg[:, :, i]))
320         metrics.append(np.mean(col_inner[:, :, i])/np.mean(col_seg[:, :, i]))
321         metrics.append(np.std(col_inner[:, :, i])/np.std(col_seg[:, :, i]))
322         metrics.append(np.mean(col_outer[:, :, i])-np.mean(col_inner[:, :, i]))
323         metrics.append(np.std(col_outer[:, :, i])-np.std(col_inner[:, :, i]))
324         metrics.append(np.mean(col_outer[:, :, i])-np.mean(col_seg[:, :, i]))
325         metrics.append(np.std(col_outer[:, :, i])-np.std(col_seg[:, :, i]))
326         metrics.append(np.mean(col_inner[:, :, i])-np.mean(col_seg[:, :, i]))
327         metrics.append(np.std(col_inner[:, :, i])-np.std(col_seg[:, :, i]))
328
329     # 6 features of color asymmetry
330     thetacol = 0.5*np.tan(2*mu11c/(mu20c-mu02c))

```

```

331     # Rotation around the centroid of the lesion
332     for i in range(0,3):
333         thetac = thetacol[i]
334         boolean_x =
335             rotation(seg,180*thetac/np.pi,x0=round(m01c[i]/m00),y0=round(m10c[i]/m00))
336         boolean_x = boolean_x.astype(np.uint8) # Cast to integer
337         diffx = sum(sum(boolean_x*seg))
338         boolean_y =
339             rotation(seg,180*thetac/np.pi-90,x0=round(m01c[i]/m00),y0=round(m10c[i]/m00))
340         boolean_y = boolean_y.astype(np.uint8) # Cast to integer
341         diffy = sum(sum(boolean_y*seg))
342         Ax = m00-diffx
343         Ay = m00-diffy
344         A1 = min(Ax,Ay)/m00
345         A2 = (Ax+Ay)/m00
346         metrics.append(A1)
347         metrics.append(A2)
348
349     # 18 features of centroidal distances
350     r0c = m10c/m00c
351     c0c = m01c/m00c
352     dist = (r0-r0c)**2+(c0-c0c)**2
353     # Normalizing
354     dist = dist/max(dist)
355     for i in range(0,len(dist)):
356         metrics.append(dist[i])
357     del colors
358     #####
359     # Texture features - 72
360     #####
361     gray = cv2.cvtColor(image, cv2.COLOR_BGR2GRAY)
362     #imshow(gray)
363     comp = np.array([seg,inner,outer])
364     ind = np.where(comp == 1)
365     # Requantized image
366     ncolor = 64
367     img_c = np.round(np.array([ncolor-1])*gray/255).astype(np.uint8)
368     # Co-occurrence matrix for 0°(0), 45°(1), 90°(2), 135°(3)
369     C = np.zeros([64,64,3,4])
370     # Size of the masks (to normalization)
371     size = np.array([[[sum(sum(seg)),sum(sum(inner)),sum(sum(outer))]]])
372     for l in range(0,len(ind[0])):
373         k = ind[0][l]
374         i = ind[1][l]
375         j = ind[2][l]
376         maxi = min(n-1,j+1)
377         minij = max(0,j-1)
378         minii = max(0,i-1)
379         C[img_c[i,j],img_c[i,maxi],k,0] +=1
380         C[img_c[i,j],img_c[minii,maxi],k,1] +=1
381         C[img_c[i,j],img_c[minii,j],k,2] +=1
382         C[img_c[i,j],img_c[minii,minij],k,3] +=1
383
384     Cmean = np.mean(C,axis=3)/size
385
386     # Maximum probability
387     mp = np.max(Cmean,axis=(0,1))
388     # Energy
389     E = sum(sum(Cmean**2))
390     # Entropy
391     S = -sum(sum(Cmean*np.log(Cmean+0.00001)))
392     # Auxiliary matrix with |i-j| to each (i,j)
393     pos = np.linspace(0,ncolor-1,ncolor)
394     m1 = np.matlib.repmat(pos,ncolor,1)
395     o = np.ones([ncolor,ncolor])
396     mpos = abs(m1 - np.transpose(pos*o))
397     mpos = np.expand_dims(mpos, axis=2)

```

```

396 m = np.expand_dims(m1, axis=2)
397 # Dissimilarity
398 D = sum(sum(mpos*Cmean))
399 # Contrast
400 C = sum(sum(mpos**2*Cmean))
401 # Inverse difference
402 ID = sum(sum(Cmean/(1+mpos)))
403 # Inverse difference moment
404 IDM = sum(sum(Cmean/(1+mpos**2)))
405 # Auxiliary variables
406 muj = sum(sum(m*Cmean))
407 mui = sum(sum(np.transpose(m, (1,0,2))*Cmean))
408 sigmaj = sum(sum((m-muj)**2*Cmean))
409 sigmai = sum(sum((np.transpose(m, (1,0,2))-mui)**2*Cmean))
410 # Correlation
411 COR = sum(sum((np.transpose(m, (1,0,2))-mui)*(m-muj)*Cmean/(sigmai*sigmaj)))
412 # Texture features for seg, inner and outer
413 t = np.array([mp,E,S,D,C,ID,IDM,COR])
414 del mp,E,S,D,C,ID,IDM,COR,mpos,o,m,pos,m1,Cmean
415 # Expanding texture features: same operations made for color (- and /)
416 new1 = np.transpose(np.array([t[:,2]-t[:,1],t[:,2]-t[:,0],t[:,1]-t[:,0]]))
417 new2 = np.transpose(np.array([t[:,2]/t[:,1],t[:,2]/t[:,0],t[:,1]/t[:,0]]))
418
419 for i in range(0,8):
420     for j in range(3):
421         metrics.append(t[i,j])
422         metrics.append(new1[i,j])
423         metrics.append(new2[i,j])
424 return metrics
425

```



```

1  # -*- coding: utf-8 -*-
2  """
3  Created on Mon Mar  4 23:23:08 2019
4
5  @author: Raul Alfredo de Sousa Silva
6  """
7  import csv
8  import pandas as pd
9  import numpy as np
10 from skimage.io import imread
11 from skimage.transform import resize
12 import features_extractor as fs
13
14 def
15 features(filename1,filename2,address1,address2,featurefile1,featurefile2,lim=3000,a=0,b=0
16 ):
17     '''
18     Given the images located in the 'filenamex', which names are included in
19     'addressx', this fuction returns all the
20     433 features of your dataset into a file with the name given in 'featurefilex'.
21     The features are:
22         13 features of shape
23         348 features of color
24         72 features of texture
25     To see more details about this features please refer to the report
26     Variables:
27     # Change this to access training set images
28         filename1 = 'XXX'
29     # Change this to access test set images
30         filename2 = 'XXX'
31     # Change this to access training set names
32         address1 = 'XXX'
33     # Change this to access test set names
34         address2 = 'XXX'
35     # Change this to rename the training set features file
36         featurefile1 = 'XXX'
37     # Change this to rename the test set features file
38         featurefile2 = 'XXX'
39     # A reduction is applied to images in which one of the dimension is greater
40     then 3000 pixels by default (to reduce the computational time). You are able
41     to change it.
42     # Threshold for reduction
43     lim = 3000
44     # New dimensions in case of reduction
45     a=0 size of rows
46     b=0 size of columns
47     If they were let at 0 the default reduction will be applied, which means
48     each dimension divided by 2.
49     Syntax:
50         features(filename1,filename2,address1,address2,featurefile1,
51                 featurefile2,lim=3000,a=0,b=0)
52     Observation1: To reduce computational time images with more than 3000 pixels
53     are reduced by 2 in the two dimensions
54     Observation3: We suppose to have .jpg images
55     Observation2: We suppose to have a segmentation map of the image with
56     filename    <image>_segmentation.jpg
57     '''
58     # Change this to access training set images
59     #filename1 = 'C:\\Users\\raul-\\Documents\\DataSet-Challenge-IMA205-2019\\data\\im\\'
60     # Change this to access test set images
61     #filename2 = 'C:\\Users\\raul-\\Documents\\DataSet-Challenge-IMA205-2019\\data\\im\\'
62     # Change this to access training set names
63     #address1 =
64     'C:\\Users\\raul-\\Documents\\DataSet-Challenge-IMA205-2019\\data\\train.csv'
65     # Change this to access test set names
66     #address2 =
67     'C:\\Users\\raul-\\Documents\\DataSet-Challenge-IMA205-2019\\data\\test.csv'

```

```

63     # Change this to rename the training set file
64     #featurefile1 = 'features_train.csv'
65     # Change this to rename the test set file
66     #featurefile2 = 'features_test.csv'
67     #a=0
68     #b=0
69     #lim = 3000
70
71     # Creating labels
72     nfeat = 433
73     labels = ["ImageID"]
74     for i in range(0,nfeat):
75         labels.append('f{}'.format(i))
76     # Loading features of the training-set into a .csv file
77     df = pd.read_csv(address1)
78     X_df = df['ImageId']
79     X_train = X_df.values
80     # Creating archive and writing labels
81     with open(featurefile1, 'w',newline='') as csvfile:
82         filewriter = csv.writer(csvfile)
83         # Writing labels
84         filewriter.writerow(labels)
85     # Writing training-set features in the file
86     i=0
87     for name_im in X_train:
88         filename = filename1+'{}.jpg'.format(name_im)
89         image = imread(filename)
90         filename_Segmentation = filename1+'{}_segmentation.jpg'.format(name_im)
91         image_Segmentation = imread(filename_Segmentation) # Value 0 or 255
92         # Use if necessary
93         (h,w,c) = image.shape
94         if (h > lim or w > lim):
95             if (a == 0 or b==0):
96                 a = int(h/2)
97                 b = int(w/2)
98                 h_n = a
99                 w_n = b
100                image = resize(image, (h_n,w_n), mode='reflect')
101                image_Segmentation = resize(image_Segmentation, (h_n,w_n), mode='reflect')
102                # To get uint8
103                seg = (np.round(image_Segmentation)).astype(np.uint8)
104                image = (np.round(255*image)).astype(np.uint8)
105            else:
106                seg = (image_Segmentation/255).astype(np.uint8)
107            print()
108            features = fs.extract(image,seg,name_im)
109            with open(featurefile1, 'a',newline='') as csvfile:
110                filewriter = csv.writer(csvfile)
111                filewriter.writerow(features)
112            i+=1
113            print(i,"out of 700")
114
115     # Loading features of the trainingset into a .csv file
116     df = pd.read_csv(address2)
117     X_df = df['ImageId']
118     y_df = df['Malignant']
119     X_test = X_df.values
120     y_test = y_df.values
121
122     # Creating labels
123
124     with open(featurefile2, 'w',newline='') as csvfile:
125         filewriter = csv.writer(csvfile)
126         # Writing labels
127         filewriter.writerow(labels)
128     # Writing features of the test-set into a .csv file
129     i=0

```

```

130
131     for name_im in X_test:
132         filename = filename2+'{}.jpg'.format(name_im)
133         image = imread(filename)
134         filename_Segmentation = filename2+'{}_segmentation.jpg'.format(name_im)
135         image_Segmentation = imread(filename_Segmentation) # Value 0 or 255
136         # Use if necessary
137         (h,w,c) = image.shape
138         if (h > lim or w > lim):
139             if (a == 0 or b==0):
140                 a = int(h/2)
141                 b = int(w/2)
142             h_n = a
143             w_n = b
144             image = resize(image, (h_n,w_n), mode='reflect')
145             image_Segmentation = resize(image_Segmentation, (h_n,w_n), mode='reflect')
146             # To get uint8
147             seg = (np.round(image_Segmentation)).astype(np.uint8)
148             image = (np.round(255*image)).astype(np.uint8)
149         else:
150             seg = (image_Segmentation/255).astype(np.uint8)
151         features = fs.extract(image,seg,name_im)
152         with open(featurefile2, 'a',newline='') as csvfile:
153             filewriter = csv.writer(csvfile)
154             filewriter.writerow(features)
155         i+=1
156     print(i,"out of 300")
157
158

```

```

1  # -*- coding: utf-8 -*-
2  """
3  Created on Mon Apr 15 22:20:31 2019
4
5  @author: Raul Alfredo de Sousa Silva
6  """
7
8  ### Printing versions
9  print("Library versions used in this work:")
10 print("csv version: 1.0")
11 print("numpy version: 1.15.4")
12 print("Tensorflow version:1.13.1")
13 print("matplotlib version: 3.0.2")
14 print("sklearn version: 0.20.1")
15 ###
16 import csv
17 import itertools
18 import numpy as np
19 import tensorflow as tf
20 import matplotlib.pyplot as plt
21 from sklearn import decomposition
22 from sklearn.svm import SVC, LinearSVC
23 from sklearn.tree import DecisionTreeClassifier
24 from sklearn.preprocessing import StandardScaler
25 from sklearn.neighbors import KNeighborsClassifier
26 from sklearn.model_selection import cross_val_score, GridSearchCV, KFold
27
28 import warnings
29 warnings.filterwarnings("ignore", category=DeprecationWarning)
30 warnings.simplefilter(action='ignore', category=FutureWarning)
31
32 from sklearn.exceptions import ConvergenceWarning
33 warnings.filterwarnings(action='ignore', category=ConvergenceWarning)
34
35
36 ###
37 # Code from scikit-learn
38 def plot_confusion_matrix(cm, classes,
39                           normalize=False,
40                           title='Confusion matrix',
41                           cmap=plt.cm.Blues):
42     """
43     This function prints and plots the confusion matrix.
44     Normalization can be applied by setting `normalize=True`.
45     """
46     if normalize:
47         cm = cm.astype('float') / cm.sum(axis=1)[:, np.newaxis]
48         print("Normalized confusion matrix")
49     else:
50         print('Confusion matrix, without normalization')
51
52     print(cm)
53
54     plt.imshow(cm, interpolation='nearest', cmap=cmap)
55     plt.title(title)
56     plt.colorbar()
57     tick_marks = np.arange(len(classes))
58     plt.xticks(tick_marks, classes, rotation=45)
59     plt.yticks(tick_marks, classes)
60
61     fmt = '.2f' if normalize else 'd'
62     for i, j in itertools.product(range(cm.shape[0]), range(cm.shape[1])):
63         plt.text(j, i, format(cm[i, j], fmt),
64                  horizontalalignment="center",
65                  color="red")
66
67     plt.ylabel('True label')

```

```

68     plt.xlabel('Predicted label')
69     plt.tight_layout()
70
71 class_names = ["healthy", "cancer"]
72 ### Writing predictions
73
74 def generate_file(y, filename):
75     with open('features_test.csv', 'r') as csvfile:
76         csv_reader = csv.DictReader(csvfile)
77         A = []
78         for row in csv_reader:
79             A.append(row["ImageID"])
80
81     labels = [['ImageId', 'Malignant']]
82     for i in range(0, 300):
83         labels.append([A[i], int(y[i])])
84
85     with open(filename, 'w', newline='') as csvfile:
86         filewriter = csv.writer(csvfile)
87         filewriter.writerows(labels)
88
89 ### Reading features
90
91 #####
92 '''
93 Observation: Features were obtained based on the two other .py files delivered
94 with this one. It's a very long procedure even if I did all I could to
95 accelerate it in terms of vectorization of the code. It took me around 17 hours
96 to finish so I do not recommend you to repeat, but if you really would like to
97 evaluate it you can uncomment the next lines below to turn it.
98 If not, just continue, the file with the featured should be attached with this
99 code too.
100 '''
101 #import main_features as mf
102 ## Change this to access training set images
103 #filename1 = 'C:\\Users\\raul-\\Documents\\DataSet-Challenge-IMA205-2019\\data\\im\\'
104 ## Change this to access test set images (they can eventually be in different folders)
105 #filename2 = 'C:\\Users\\raul-\\Documents\\DataSet-Challenge-IMA205-2019\\data\\im\\'
106 ## Change this to access training set names
107 #address1 = 'C:\\Users\\raul-\\Documents\\DataSet-Challenge-IMA205-2019\\data\\train.csv'
108 ## Change this to access test set names
109 #address2 = 'C:\\Users\\raul-\\Documents\\DataSet-Challenge-IMA205-2019\\data\\test.csv'
110 ## Change this to rename the training set features file
111 #featurefile1 = 'features_train_t.csv'
112 ## Change this to rename the test set features file
113 #featurefile2 = 'features_test_t.csv'
114 #
115 ## A reduction is applied to images in which one of the dimension is greater
116 ## then 3000 pixels by default (to reduce the computational time). You are able
117 ## to change it.
118 ## Threshold for reduction
119 #lim = 3000
120 ## New dimensions in case of reduction
121 #a=0 # size of rows
122 #b=0 # size of columns
123 ## If they were let at 0 the default reduction will be applied, which means
124 ## each dimension divided by 2.
125 #mf.features(filename1, filename2, address1, address2, featurefile1, featurefile2,
126 #             lim=3000, a=0, b=0)
127
128 #####
129
130 with open('features_train.csv', mode='r') as csv_file:
131     csv_reader = csv.DictReader(csv_file)
132     line_count = 0
133     row_count = 0
134     A = []

```

```

135     for row in csv_reader:
136         for element in row:
137             if line_count == 0:
138                 row_count += 1
139                 if row[element].find('IM') < 0:
140                     A.append(float(row[element]))
141             line_count += 1
142
143 with open('features_test.csv', mode='r') as csv_file:
144     csv_reader = csv.DictReader(csv_file)
145     line_count = 0
146     row_count = 0
147     for row in csv_reader:
148         for element in row:
149             if line_count == 0:
150                 row_count += 1
151                 if row[element].find('IM') < 0:
152                     A.append(float(row[element]))
153             line_count += 1
154
155
156 X = np.reshape(A, (1000, (row_count-1)))
157 %% Reading classes
158
159 with open('train.csv', mode='r') as csv_file:
160     csv_reader = csv.DictReader(csv_file)
161     line_count = 0
162     row_count = 0
163     B = []
164     for row in csv_reader:
165         for element in row:
166             if line_count == 0:
167                 row_count += 1
168                 if row[element].find('IM') < 0:
169                     B.append(int(row[element]))
170             line_count += 1
171 y = np.reshape(B, (700))
172 %%
173 #####
174 #                               Preprocessing                               #
175 #####
176
177 # Scale data (each feature will have average equal to 0 and unit variance)
178 scaler = StandardScaler()
179 scaler.fit(X)
180 X = scaler.transform(X)
181
182 # Use number of components take explain 95% of variability
183 pca = decomposition.PCA(n_components=0.95)
184 pca.fit(X)
185 X_pca = pca.transform(X)
186
187
188 X_train = X_pca[:700,:]
189 X_test = X_pca[700:,:]
190
191 %%
192 #####
193 #                               First test: KNN                               #
194 #####
195
196 # Direct training
197
198 print("Fitting K-nearest neighbour to the training set")
199 p_grid_KNN = {'n_neighbors': [1,3,5,7,9,11,13,15]}
200 KNN = KNeighborsClassifier()
201 inner_cv = KFold(n_splits=5, shuffle=True, random_state=42)

```

```

202 outer_cv = KFold(n_splits=5, shuffle=True, random_state=42)
203
204 grid_KNN = GridSearchCV(estimator=KNN, param_grid=p_grid_KNN,
205                          scoring='balanced_accuracy', cv=5)
206 grid_KNN.fit(X_train, y)
207 print("Best training Score: {}".format(grid_KNN.best_score_))
208 print("Best training params: {}".format(grid_KNN.best_params_))
209
210 y_pred = grid_KNN.predict(X_test)
211 name = 'test_KNN.csv'
212 generate_file(y_pred, name)
213
214 ###
215 #####
216 #                               Applying boosting to knn                               #
217 #####
218 # Creating multiple predictors (21)
219 p_grid_KNN = {'n_neighbors': [1,3,5,7,9,11,13,15]}
220 KNN = KNeighborsClassifier()
221
222 T = 21
223 m = len(X_train)
224 ys = y - (y==0)
225 p = 1/m*np.ones(m)
226 index = np.zeros(400)
227 alpha = np.zeros(T)
228 grid_KNN = []
229 S = X_train
230 for t in range(T):
231     Z = sum(p)
232     p /= Z
233     pcum = np.cumsum(p)
234     for i in range(400):
235         index[i] = np.where(pcum >= np.random.rand())[0][0]
236     index = index.astype(np.int)
237     St = X_train[index,:]
238     yt = ys[index]
239     grid_KNN.append(GridSearchCV(estimator=KNN, param_grid=p_grid_KNN,
240                                  scoring='balanced_accuracy', cv=5))
241     grid_KNN[-1].fit(St, yt)
242     et = sum((grid_KNN[-1].predict(St) != yt))/len(St)
243     alpha[t] = 0.5*np.log((1-et)/(et+1e-10))
244     correct1 = (grid_KNN[-1].predict(S) != ys).astype(np.int)
245     correct2 = (grid_KNN[-1].predict(S) == ys).astype(np.int)
246     correct = correct1 - correct2
247     p = p/Z*np.e**(correct*alpha[t])
248
249 y_pred = np.zeros(300)
250 for t in range(T):
251     y_pred += alpha[t]*grid_KNN[t].predict(X_test)
252 y_pred_fin = (y_pred > 0)
253 y_pred_fin = y_pred_fin.astype(np.int)
254
255 name = 'test_KNNboost.csv'
256 generate_file(y_pred_fin, name)
257
258
259 ###
260
261 #####
262 #                               Second test: SVM                               #
263 #####
264
265 # Fitting Linear SVM on original data
266 print("Fitting Linear SVM to the training set")
267
268 p_grid_lsvm = {'C': [1e-3,1e-2,1e-1,1,2,3,4,5,6,7,8,9,1e1]}

```

```

269 Lsvm = LinearSVC(class_weight='balanced')
270 inner_cv = KFold(n_splits=5, shuffle=True, random_state=42)
271 outer_cv = KFold(n_splits=5, shuffle=True, random_state=42)
272 # Nested CV with parameter optimization
273 clf = GridSearchCV(estimator=Lsvm, param_grid=p_grid_lsvm, cv=inner_cv)
274 nested_score = cross_val_score(clf, X_train, y, cv=outer_cv,
275                                 scoring='balanced_accuracy')
276
277 print("Average and std Nested Cv score : {0} +- {1}".format(nested_score.mean(),
278                                                             nested_score.std() ))
279
280 # Looking for the best hyperparameters
281 grid_lsvm = GridSearchCV(estimator=Lsvm, param_grid=p_grid_lsvm,
282                         scoring='balanced_accuracy', cv=5)
283 grid_lsvm.fit(X_train, y)
284 print("Best Score: {}".format(grid_lsvm.best_score_))
285 print("Best params: {}".format(grid_lsvm.best_params_))
286
287 y_pred = grid_lsvm.predict(X_test)
288 name = 'test_SVM.csv'
289 generate_file(y_pred, name)
290
291 ###
292 #####
293 #                               Applying non linearity to SVM                               #
294 #####
295 # Fitting Non-linear SVM
296 print("Fitting Non-linear SVM to the training set")
297 p_grid_nlsvm = {'C': [1e-3, 1e-2, 1e-1, 1, 2, 3, 4, 5, 6, 7, 8, 9, 1e1],
298                'gamma': [0.0001, 0.0005, 0.001, 0.005, 0.01, 0.1], }
299 NLsvm = SVC(kernel='rbf', class_weight='balanced')
300 inner_cv = KFold(n_splits=5, shuffle=True, random_state=42)
301 outer_cv = KFold(n_splits=5, shuffle=True, random_state=42)
302
303 # Looking for the best hyperparameters
304 grid_nlsvm = GridSearchCV(estimator=NLsvm, param_grid=p_grid_nlsvm,
305                           scoring="balanced_accuracy", cv=5)
306 grid_nlsvm.fit(X_train, y)
307 print("Best Score: {}".format(grid_nlsvm.best_score_))
308 print("Best params: {}".format(grid_nlsvm.best_params_))
309
310
311 y_pred = grid_nlsvm.predict(X_test)
312 name = 'test_SVM_nl.csv'
313 generate_file(y_pred, name)
314 ###
315 #####
316 #                               Third test: Decision Trees                               #
317 #####
318 # Fitting Decision Trees
319 Tree = DecisionTreeClassifier()
320 Tree.fit(X_train, y)
321 # Score in the training set
322 print('Score in the training set is {0}'.format(Tree.score(X_train, y)) )
323
324 y_pred = Tree.predict(X_test)
325 name = 'test_tree.csv'
326 generate_file(y_pred, name)
327
328 ###
329 #####
330 #                               Fourth test: MLP                               #
331 #####
332 y_train = tf.keras.utils.to_categorical(y)
333
334 def init_weights_and_biases(shape, stddev=0.1, seed_in=None):
335     """

```



```

336 This function should return Tensorflow Variables containing the initialised
337 weights and biases of the network,
338 using a normal distribution for the initialisation, with stddev of the
339 normal as an input argument
340
341 Parameters
342 -----
343 shape : tuple, (n_input_features,n_hidden_1, n_hidden_2,n_classes)
344         sizes necessary for defining the weights and biases
345
346 Returns
347 -----
348 w1, b1, w2, b2, w3, b3: initialised weights and biases, with correct shapes
349 """
350
351 # BEGIN STUDENT CODE
352 w1 = tf.Variable(tf.random_normal([shape[0],shape[1]], stddev=stddev,
353                                   seed=seed_in))
354 w2 = tf.Variable(tf.random_normal([shape[1],shape[2]], stddev=stddev,
355                                   seed=seed_in))
356 w3 = tf.Variable(tf.random_normal([shape[2],shape[3]], stddev=stddev,
357                                   seed=seed_in))
358 b1 = tf.Variable(tf.zeros([shape[1]]))
359 b2 = tf.Variable(tf.zeros([shape[2]]))
360 b3 = tf.Variable(tf.zeros([shape[3]]))
361 # END STUDENT CODE
362 return w1, b1, w2, b2, w3, b3
363
364 def forward_prop_multi_layer(X, w1, b1, w2, b2, w3, b3):
365     """
366     This function should define the network architecture, explained above
367
368     Parameters
369     -----
370     X : input to network
371     w1, w2, w3 : Tensorflow Variables
372                 network weights
373     b1, b2, b3 : Tensorflow Variables
374                 network biases
375
376     Returns
377     -----
378     Y_pred :
379     the output layer of the network, the classification prediction
380     """
381
382     # BEGIN STUDENT CODE
383     A1 = tf.nn.sigmoid(tf.matmul(X, w1)+b1)
384     A2 = tf.nn.softmax(tf.matmul(A1,w2)+b2)
385     Y_pred = tf.nn.softmax(tf.matmul(A2,w3)+b3)
386     # END STUDENT CODE
387     return Y_pred
388
389 def accuracy(Y_pred, Y_true):
390     """
391     This function calculates the network's accuracy, ie the percentage of
392     correct classifications.
393     Here, we consider the class with the highest score in the network's output
394
395     Parameters
396     -----
397     Y_pred : Tensorflow variable
398              predicted classification of network
399     Y_true : Tensorflow variable
400              true classes of input
401
402     """

```

```

403
404     acc = tf.equal(tf.argmax(Y_pred, 1), tf.argmax(Y_true, 1))
405     acc = tf.reduce_mean(tf.cast(acc, tf.float32))
406
407     return acc
408
409 RANDOM_SEED = 52#42#
410 tf.set_random_seed(RANDOM_SEED)
411
412 # Network Parameters
413 n_hidden_1 = 256 # 1st layer number of neurons
414 n_hidden_2 = 256 # 2nd layer number of neurons
415 n_input = X_train.shape[1]
416 n_classes = y_train.shape[1] # MNIST total classes (0-9 digits)
417
418 # tf Graph input
419 X_input = tf.placeholder("float", [None, n_input])
420 Y_true = tf.placeholder("float", [None, n_classes])
421
422 # Store layers weight & bias
423 stddev = 0.1
424
425 w1,b1,w2,b2,w3,b3 = init_weights_and_biases([n_input, n_hidden_1, n_hidden_2,
426                                             n_classes], stddev=0.1, seed_in=RANDOM_SEED)
427
428 # Construct model
429 Y_pred = forward_prop_multi_layer(X_input,w1,b1,w2,b2,w3,b3)
430
431 # Define loss and optimizer
432 cross_entropy = -tf.reduce_sum(Y_true * tf.log(Y_pred),axis=1)
433
434 loss = tf.reduce_mean(cross_entropy)
435 acc = accuracy(Y_pred, Y_true)
436
437 learning_rate = 0.001
438 optimizer = tf.train.AdamOptimizer(learning_rate=learning_rate)
439 training_variables = optimizer.minimize(loss)
440
441 # Parameters
442 n_epochs = 500
443 train_accuracy = []
444 test_accuracy = []
445 batch_size = 700
446 display_step = 1
447 n_batches = int(np.ceil(X_train.shape[0]/batch_size))
448 print(n_batches)
449
450
451 with tf.Session() as sess:
452     # Initializing the variables
453     init = tf.global_variables_initializer()
454     sess.run(init)
455
456     for epoch in range(n_epochs):
457         # Loop over all batches
458         for batch_idx in range(n_batches):
459             #get the next batch in the MNIST dataset and carry out training
460             poss = batch_size*batch_idx
461             posf = batch_size*(batch_idx+1)
462
463             #BEGIN STUDENT CODE
464             sess.run([training_variables,loss], feed_dict={X_input:
465                 X_train[poss:posf,:], Y_true: y_train[poss:posf]})
466             #END STUDENT CODE
467
468             # calculate accuracy for this epoch
469             train_accuracy.append(sess.run(acc, feed_dict={X_input: X_train,

```

```
470                                     Y_true:y_train)))
471     u = sess.run([Y_pred], feed_dict={X_input: X_test})
472     print(".", end='')
473
474     print("Training finished")
475
476
477     ypred = np.array(u)
478     ypred = ypred.reshape(300,2)
479     ypred = ypred>0.5
480     ypred = ypred[:,1]
481     print(sum(ypred.astype(np.int)))
482     name = 'test_MLP.csv'
483     generate_file(ypred,name)
484
485
```