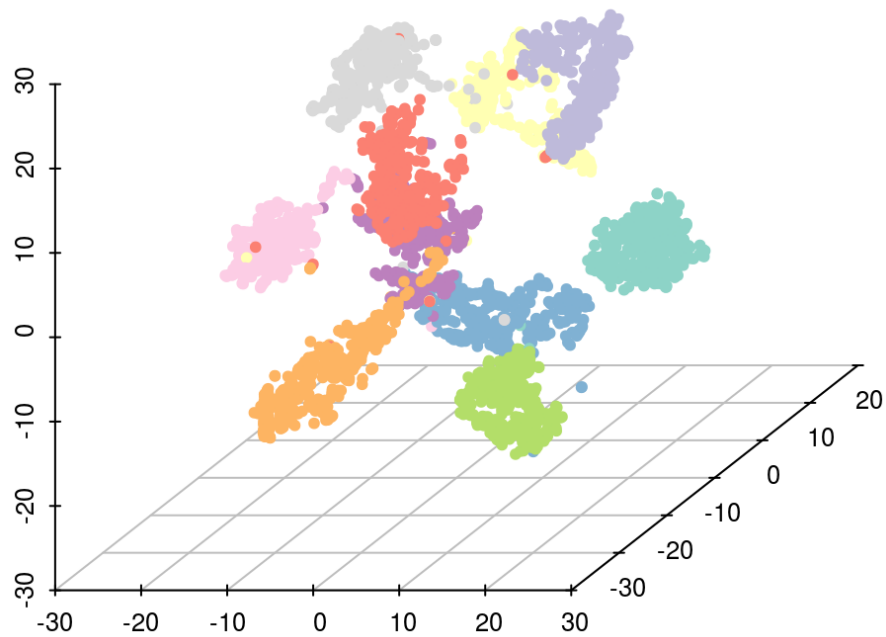


GEOMETRÍA COMPUTACIONAL

PCA Y ANALOGÍA



LUCAS DE TORRE

Índice

1. Introducción	2
2. Material	2
3. Resultados	3
4. Conclusiones	3
5. Anexo A: Código	4

1. Introducción

Consideramos la atmósfera como un sistema de 6 variables: tiempo (t), longitud (x , discretizada en 144 valores), latitud (y , discretizada en 73 valores), presión (p , discretizada en 17 valores), temperatura (T) y altura geopotencial (Z), siendo las dos últimas las únicas variables dinámicas respecto a la variable tiempo. A partir de ahí, considerando el sistema $S = \{a_d, X_d\}_{d=1}^{365}$, de los días a_d del año y las variables de estado $X_d = \{Z_{i,j,k}\}_{i=1, j=1, k=1}^{i=144, j=73, k=17}$, estimaremos las 4 componentes principales, las representaremos espacialmente en (x, y) y veremos el porcentaje de varianza explicado.

Además, consideraremos el subsistema $\sigma \in S$ tal que $x \in (340^\circ, 20^\circ)$ e $y \in (30^\circ, 50^\circ)$. En este subsistema σ , teniendo en cuenta solo la altura geopotencial, hallaremos los 4 días más análogos al día $a_0 = 2020/01/20$ y calcularemos el error absoluto de la temperatura prevista para a_0 por la media de esos 4 análogos.

2. Material

Para la estimación de las 4 componentes principales, primero cargamos los datos de altura geopotencial de 2019. De esos datos, nos quedamos con aquellos cuyo nivel de presión sea de 500hPa.

Después, haciendo uso de la biblioteca *sklearn*, utilizamos el método PCA (con 4 componentes principales como argumento) para estimar las 4 componentes principales. Por último, utilizamos otro método de la librería *sklearn* para mostrar el ratio de varianza explicado y dibujamos las componentes principales en el plano.

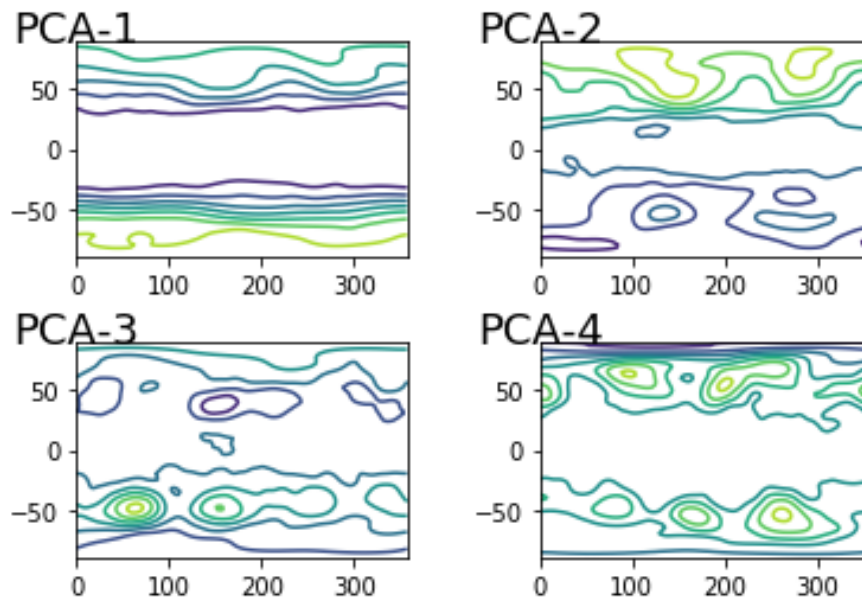
Procedemos ahora a buscar los 4 días más análogos al día a_0 . Para ello, cargamos los datos de altura geopotencial de dicho día (para lo que hallamos el índice que corresponde a dicho día en la matriz). Después, cargamos los datos de altura geopotencial de 2019 y calculamos la distancia euclídea de cada uno de los días de 2019 al día a_0 . Nos quedamos con los 4 días más análogos al día a_0 (aquellos cuya distancia a a_0 fuera menor), digamos $d1$, $d2$, $d3$, $d4$.

Después, cargamos los datos de temperatura del día a_0 y los datos de temperatura de 2019 (estos datos y los del día a_0 los transformamos a grados kelvin). De estos últimos, hacemos la media (\bar{d}) únicamente de los 4 días ($d1$, $d2$, $d3$, $d4$) calculados antes (eran los más parecidos respecto a altura geopotencial, y ahora calculamos su media respecto a temperatura).

Por último, calculamos la media del valor absoluto de las diferencias entre \bar{d} y a_0 . De esta manera, hemos calculado el error absoluto de la temperatura prevista para a_0 por la media \bar{d} de los 4 días más análogos a a_0 según la altura geopotencial.

3. Resultados

Los valores explicados en función del número de componentes principales son: 0.86579346, 0.93789939, 0.94300602 y 0.94717256, en función de si el número de componentes principales es 1, 2, 3 ó 4, respectivamente. Por tanto, los porcentajes son 86.579346, 93.789939, 94.300602 y 94.717256, respectivamente. Su representación espacial en (x,y) es:



Por otro lado, los 4 días de 2019 más análogos al día a_0 según la altura geopotencial son: 78, 79, 109 y 335. El error absoluto de la temperatura prevista para a_0 por la media de esos 4 análogos es: 4.2725863.

4. Conclusiones

Llama la atención que los días más análogos al día a_0 se reparten bastante a lo largo del año, habiendo 2 en invierno, 1 en primavera y otro en otoño. Sin embargo, los 4 días son muy próximos al comienzo o al final del invierno (lo que es lógico, porque tendrán unas condiciones parecidas).

5. Anexo A: Código

```
1 # -*- coding: utf-8 -*-
2 """
3 Referencias:
4
5     Fuente primaria del rean lisis
6     https://www.esrl.noaa.gov/psd/data/gridded/data.ncep.
       reanalysis2.pressure.html
7
8     Altura geopotencial en niveles de presi n
9     https://www.esrl.noaa.gov/psd/cgi-bin/db_search/DBListFiles.pl?
       did=59&tid=81620&vid=1498
10
11     Temperatura en niveles de presi n:
12     https://www.esrl.noaa.gov/psd/cgi-bin/db_search/DBListFiles.pl?
       did=59&tid=81620&vid=1497
13 """
14 #import os
15 import datetime as dt # Python standard library datetime module
16 import numpy as np
17 import matplotlib.pyplot as plt
18 from scipy.io import netcdf as nc
19 from sklearn.decomposition import PCA
20 import math
21
22 """
23 Apartado 1
24 """
25
26 def apartado1():
27     print("Apartado 1:\n")
28
29
30     #Cargamos los datos de altura geopotencial del 2019
31     f = nc.netcdf_file("hgt.2019.nc", 'r')
32     time = f.variables['time'][:].copy()
33     level = f.variables['level'][:].copy() #los valores de p
34     lats = f.variables['lat'][:].copy() #los valores de y
35     lons = f.variables['lon'][:].copy() #los valores de x
36     hgt = f.variables['hgt'][:].copy() #los valores de cada d a
37     f.close()
38
39     """
40     Ejemplo de evoluci n temporal de un elemento de hgte
41
42     plt.plot(time, hgt[:, 1, 1, 1], c='r')
43     plt.show()
44     """
```

```

45
46     PRESION = 500
47     for i in range(len(level)):
48         if level[i]==PRESION:
49             p500=i
50             break
51
52
53     plt.title("Distribuci n espacial de la altura geopotencial en
54 el nivel de 500hPa, para el primer d a")
55     plt.contour(lons, lats, hgt[0,p500,:,:])
56     plt.show()
57
58     #Seleccionamos los datos con p = 500hPa y redimensionamos la
59 matriz
60     hgt2 = hgt[:,p500,:,:].reshape(len(time),len(lats)*len(lons))
61
62     n_components=4
63
64     Y = hgt2.transpose()
65     pca = PCA(n_components=n_components)
66
67     #Aplicamos el algortimo PCA
68     pca.fit(Y)
69     print("Valores explicados en funci n del n mero de
70 componentes:",pca.explained_variance_ratio_.cumsum())
71
72     #Ejercicio de la pr ctica
73     #Representaci n espacial de las componentes principales
74     Element_pca = pca.fit_transform(Y)
75     Element_pca = Element_pca.transpose(1,0).reshape(n_components,
76 len(lats),len(lons))
77     fig = plt.figure()
78     fig.subplots_adjust(hspace=0.4, wspace=0.4)
79     for i in range(1, n_components+1):
80         ax = fig.add_subplot(2, 2, i)
81         ax.text(0.5, 90, 'PCA-' +str(i),
82             fontsize=18, ha='center')
83         plt.contour(lons, lats, Element_pca[i-1,:,:])
84     plt.show()
85     print("\n\n\n\n")
86
87
88 """
89 apartado 2

```

```

90 """
91 IND_500=5
92 IND_1000=0
93
94 #Calcula la distancia eucl dea entre dos d as en un subconjunto
  del sistema
95 def dist_euc(d0,d,min_lat,max_lat,min_lon,max_lon,lons):
96     dist=0
97     for i in range(min_lat+1,max_lat):
98         for j in range(min_lon+1,max_lon):
99             dist+=0.5*(d0[IND_500,i,j%lons]-d[IND_500,i,j%lons])
100             **2+0.5*(d0[IND_1000,i,j%lons]-d[IND_1000,i,j%lons])**2
101         return math.sqrt(dist)
102
103 def apartado2():
104     print("Apartado 2:\n")
105
106     #Cargamos los datos de altura geopotencial del 2020
107     f = nc.netcdf_file("hgt.2020.nc", 'r')
108     time = f.variables['time'][:].copy()
109     lats = f.variables['lat'][:].copy() #los valores de y
110     lons = f.variables['lon'][:].copy() #los valores de x
111     hgt = f.variables['hgt'][:].copy() #los valores de cada d a
112     f.close()
113
114     #Calculamos el ndice correspondiente al d a 2020/01/20
115     for t in range(len(time)):
116         if dt.date(1800, 1, 1) + dt.timedelta(hours=time[t]) == dt.
date(2020,1,20):
117             break
118
119     double_lons=np.concatenate((lons, lons), axis=None)
120     min_lon=-1
121     for max_lon in range(len(double_lons)):
122         if double_lons[max_lon]==340:
123             min_lon=max_lon
124         if double_lons[max_lon]==20 and not min_lon==-1:
125             break
126
127     for max_lat in range(len(lats)):
128         if lats[max_lat]==50:
129             min_lat=max_lat
130         if lats[max_lat]==30:
131             break
132
133
134     #cargamos los datos de altura geopotencial del d a 2020/01/20
135     d0h = hgt[t,:,:,:]

```

```
136
137
138     #cargamos los datos de temperatura del 2020
139     f = nc.netcdf_file("air.2020.nc", 'r')
140     air = f.variables['air'][:].copy() #los valores de cada d a
141     tem_scale_factor=f.variables['air'].scale_factor.copy()
142     tem_add_offset=f.variables['air'].add_offset.copy()
143     f.close()
144
145     #cargamos los datos de temperatura del d a 2020/01/20
146     d0a = air[t,0,:,:]*tem_scale_factor+tem_add_offset
147
148
149     #cargamos los datos de altura geopotencial del 2019
150     f = nc.netcdf_file("hgt.2019.nc", 'r')
151     lons = f.variables['lon'][:].copy() #los valores de x
152     hgt = f.variables['hgt'][:].copy() #los valores de cada d a
153     f.close()
154
155
156
157     #Calculamos la distancia de cada d a de nuestro subsistema al
158     d a 2020/01/20
159     distancias=[[dist_euc(d0h,hgt[i,:,:,:],min_lat,max_lat,min_lon,
160     max_lon,len(lons)),i]for i in range(hgt.shape[0])]
161     distancias.sort()
162
163
164     take=[j for i,j in distancias[0:4]]
165     print("Los 4 d a s m s a n l o g o s s o n:",take,"\n")
166
167
168     #cargamos los datos de temperatura de 2019
169     f = nc.netcdf_file("air.2019.nc", 'r')
170     air = f.variables['air'][:].copy() #los valores de cada d a
171     tem_scale_factor=f.variables['air'].scale_factor.copy()
172     tem_add_offset=f.variables['air'].add_offset.copy()
173     f.close()
174
175     #Transformamos los datos para que est n en grados Kelvin
176     air=air*tem_scale_factor+tem_add_offset
177
178
179     #Calculamos la media de los 4 d a s m s a n l o g o s a 2020/01/20
180     mediaDias=np.mean(air[take,0,:,:),axis=0)
181
182     print("El error absoluto medio de la temperatura prevista para
183     el d a a0 es",np.mean(abs(mediaDias-d0a)))
184
185
186 apartado1()
187 apartado2()
```