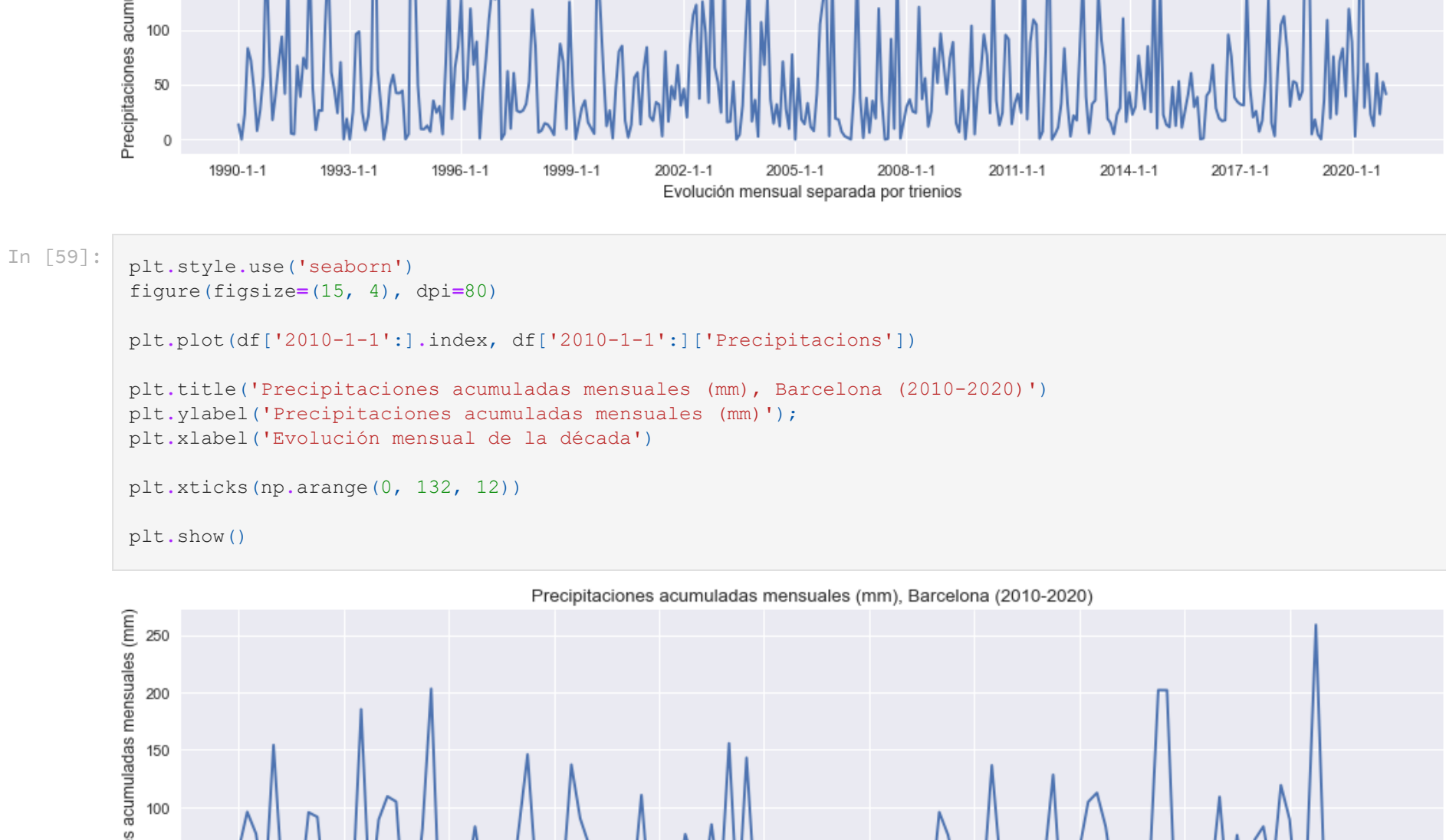



```
plt.style.use('seaborn')
figure(figsize=(15, 4), dpi=80)

plt.plot(df30.index[df30['Precipitaciones']])

plt.title('Precipitaciones acumuladas mensuales (mm), Barcelona (1990-2020)')
plt.ylabel('Precipitaciones acumuladas mensuales (mm)')
plt.xlabel('Evolución mensual separada por trienios')
plt.xticks(np.arange(0, 396, 36))
plt.show()
```



```
plt.style.use('seaborn')
figure(figsize=(15, 4), dpi=80)

plt.plot(df['2010-1-1':].index, df['2010-1-1':]['Precipitaciones'])

plt.title('Precipitaciones acumuladas mensuales (mm), Barcelona (1990-2020)')
plt.ylabel('Precipitaciones acumuladas mensuales (mm)')
plt.xlabel('Evolución mensual de la década')
plt.xticks(np.arange(0, 132, 12))
plt.show()
```



4. Clusters No-Supervisados

```
#librerías para clustering no-supervisado

from sklearn import preprocessing
from sklearn.model_selection import train_test_split
from sklearn.cluster import KMeans
from sklearn import metrics
from sklearn.preprocessing import StandardScaler
from sklearn.cluster import AgglomerativeClustering
import scipy.cluster.hierarchy as sch
from sklearn.preprocessing import normalize
from kneed import DataGenerator, KneedleLocator
from sklearn.cluster import DBSCAN
from sklearn.mixture import GaussianMixture
from sklearn.metrics import adjusted_rand_score
from sklearn.decomposition import PCA
from sklearn.preprocessing import LabelEncoder

In [60]:

df = pd.read_csv('precipitacionabarcelonadesde1786_format_long.csv')

In [62]:

a = [1 for x in range(len(df))]

df['Dia'] = a

df['Dia'] = df['Dia'].astype(str)

df['Mes'] = df.Any.apply(str)

df['Mes'] = df.Mes.apply(str)

df['Date'] = df['Any'] + '-' + df['Mes'] + '-' + df['Dia']

df['Date'] = df['Date'].astype(str)

pd.to_datetime(df['Date'])

df.drop(['Dia'], 1, inplace=True)

df = df[['Date', 'Any', 'Mes', 'Deca_Mes', 'Precipitations']]

df.set_index('Date', inplace=True)

In [63]:

df['Mes'] = df['Mes'].astype(int)
df['Any'] = df['Any'].astype(int)

In [64]:

df.reset_index(inplace=True)

In [65]:

df.drop(['Deca_Mes', 'Date'], 1, inplace=True)

In [66]:

df
```

Out[66]:

	Any	Mes	Precipitations
0	1786	1	32.8
1	1786	2	28.4
2	1786	3	84.4
3	1786	4	42.3
4	1786	5	8.5
...
2815	2020	8	12.4
2816	2020	9	60.2
2817	2020	10	23.1
2818	2020	11	52.5
2819	2020	12	41.5

2820 rows x 3 columns

4.1 Hopkins Statistics

To understand if the dataset can be clustered, we will use the Hopkins statistic, which tests the spatial randomness of the data and indicates the cluster tendency or how well the data can be clustered. It calculates the probability that a given data is generated by a uniform distribution (Alkhoukadi Kassambara, n.d.). The inference is as follows for a data of dimensions 'd':

- If the value is around 0.5 or lesser, the data is uniformly distributed and hence it is unlikely to have statistically significant clusters.
- If the value is between (0.7, ..., 0.99), it has a high tendency to cluster and therefore likely to have statistically significant clusters.

```
In [67]:

from sklearn.neighbors import NearestNeighbors
from random import sample
from numpy.random import uniform
import numpy as np
from math import isnan

def hopkins(X):
    d = X.shape[1]
    #d = len(vars) # columns
    n = len(X) # rows
    m = int(0.1 * n) # heuristic from article [1]
    nbrs = NearestNeighbors(n_neighbors=m).fit(X.values)

    rand_X = sample(range(0, n, 1), m)

    ujd = []
    for j in range(0, m):
        u_dist, _ = nbrs.kneighbors(uniform(np.min(X,axis=0),np.max(X,axis=0),d).reshape(1, -1), 2, return_dist=True)
        ujd.append(u_dist[0][1])
        w_dist, _ = nbrs.kneighbors(X.iloc[rand_X[j]].values.reshape(1, -1), 2, return_distance=True)
        wjd.append(w_dist[0][1])

    H = sum(ujd) / (sum(ujd) + sum(wjd))
    if isnan(H):
        print(ujd, wjd)
        H = 0

    return H

In [68]:

round(hopkins(df), 5)

Out[68]:

0.84838
```

El cálculo de Hopkins Statistics nos devuelve 0.86129, un valor entre 0.7 y 0.99, por lo cual en principio nuestros datos son clusterizables.

4.2 Elbow Method for K-Means

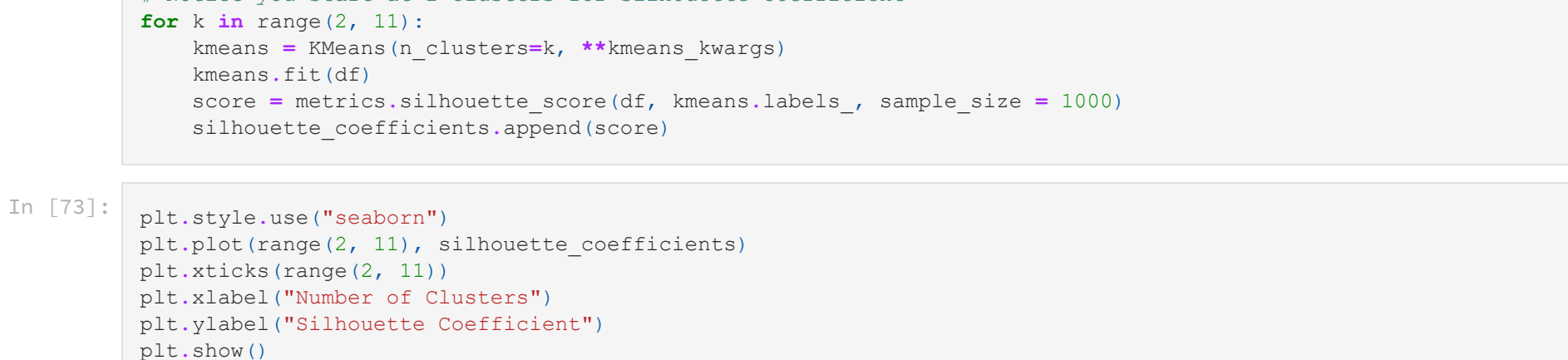
```
In [69]:

kmeans_kwards = {
    "init": "random",
    "n_init": 10,
    "max_iter": 300,
    "random_state": 42,
}

# A list holds the SSE values for each k
sse = []
for k in range(1, 11):
    kmeans = KMeans(n_clusters=k, **kmeans_kwards)
    kmeans.fit(df)
    sse.append(kmeans.inertia_)

In [70]:

plt.style.use('seaborn')
plt.plot(range(2, 11), sse)
plt.xticks(range(2, 11))
plt.xlabel('Number of Clusters')
plt.ylabel('SSE')
plt.show()
```



```
In [71]:

kl = KneedleLocator(
    range(1, 11), sse, curve="convex", direction="decreasing"
)

kl.elbow

Out[71]:

3
```

4.3 Silhouette coefficient en vez de SSE

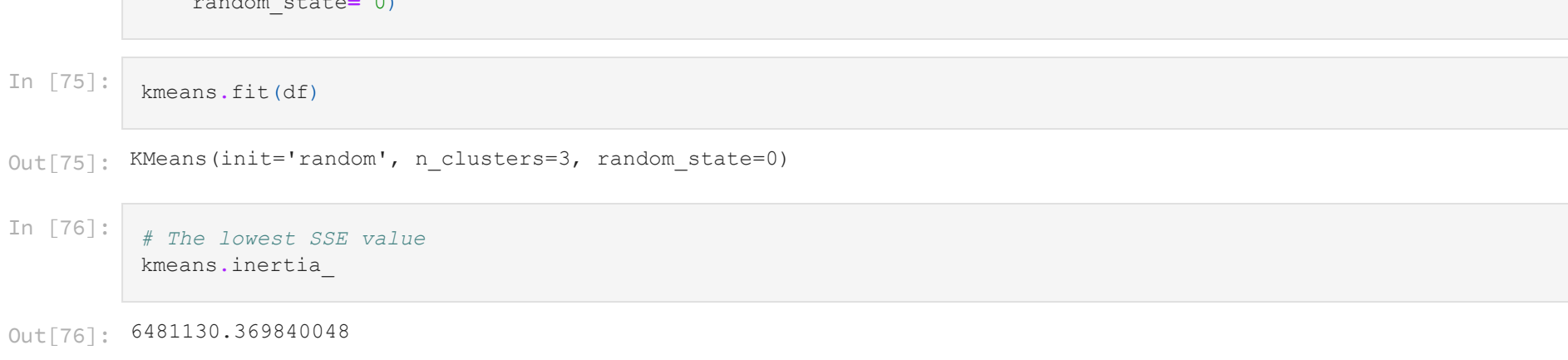
```
In [72]:

# A list holds the silhouette coefficients for each k
silhouette_coefficients = []

# Notice you start at 2 clusters for silhouette coefficient
for k in range(2, 11):
    kmeans = KMeans(n_clusters=k, **kmeans_kwards)
    score = metrics.silhouette_score(df, kmeans.labels_, sample_size = 1000)
    silhouette_coefficients.append(score)

In [73]:

plt.style.use('seaborn')
plt.plot(range(2, 11), silhouette_coefficients)
plt.xticks(range(2, 11))
plt.xlabel('Number of Clusters')
plt.ylabel('Silhouette Coefficient')
plt.show()
```



4.4 Agrupación

```
In [74]:

kmeans = KMeans(init="random",
    n_clusters=3,
    n_init=10,
    max_iter=300,
    random_state= 0)

In [75]:

kmeans.fit(df)

Out[75]:

KMeans(init='random', n_clusters=3, random_state=0)

In [76]:

# The lowest SSE value
kmeans.inertia_

Out[76]:

6481130.369840048

In [77]:

# Locations of the centroid
kmeans.cluster_centers_

Out[77]:

array([[1841.31965944, 6.40712074, 39.95162539],
       [1962.62340677, 6.35114504, 34.36044461],
       [1929.9179077, 7.34670487, 135.30143266]])

In [78]:

# The number of iterations required to converge
kmeans.n_iter_

Out[78]:

14

In [79]:

labels = kmeans.labels_

In [80]:

labels

Out[80]:

array([0, 0, 0, ..., 1, 1, 1])
```

Hacemos un scatter plot diferenciando por cluster de las precipitaciones por meses:

```
In [81]:

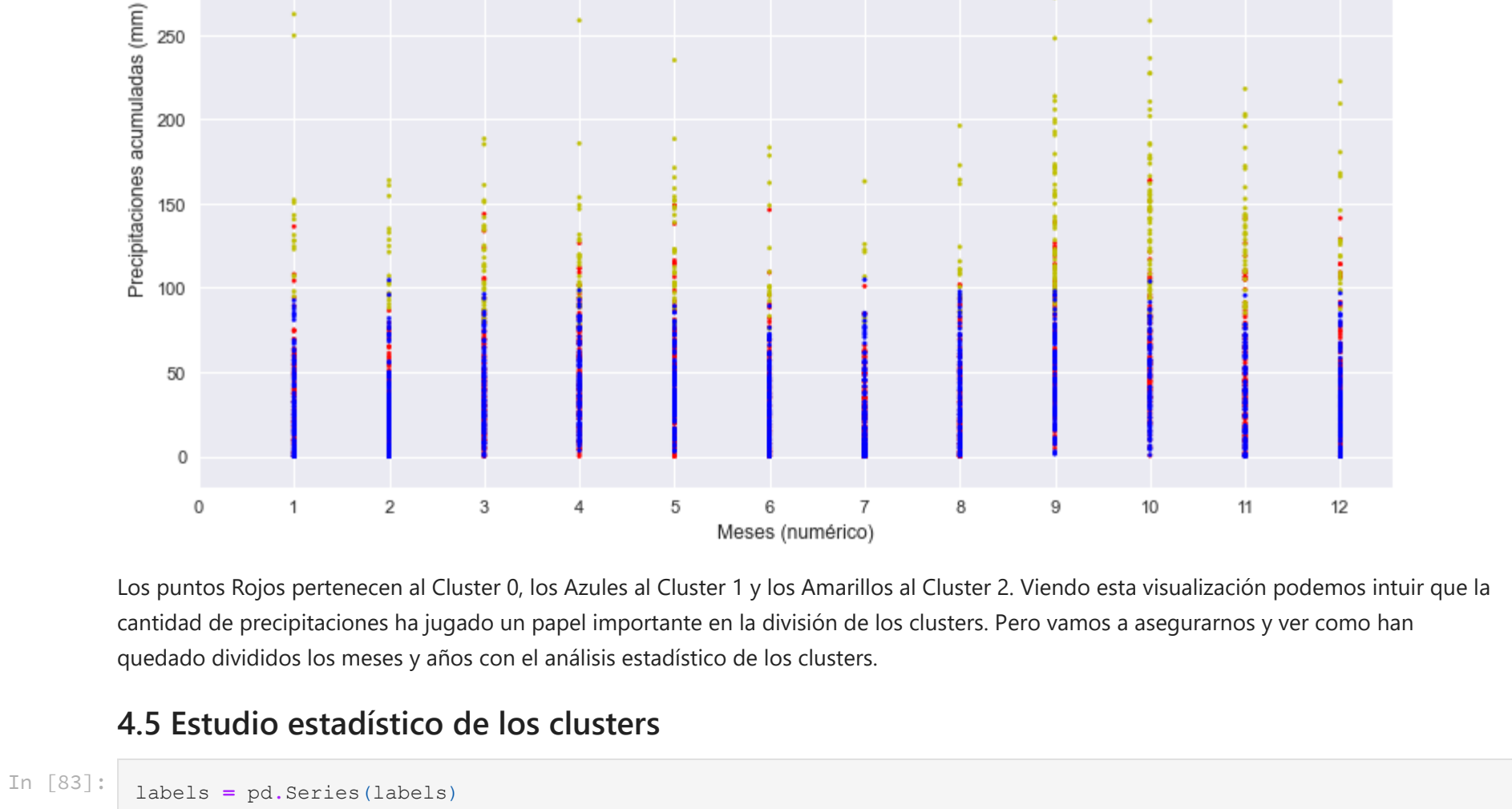
centroids = kmeans.cluster_centers_

In [82]:

plt.style.use('seaborn')
figure(figsize=(12, 7), dpi=80)
colors = ['r','b','y','g','c','m']
plt.scatter(df['Mes'], df['Precipitaciones'], color=colors[1] for i in labels, label=labels, s=5)

plt.title('Precipitación acumulada por Meses y Clusters', size=14)

plt.xticks(np.arange(0, 131))
plt.xlabel('Meses (numérico)')
plt.ylabel('Precipitaciones acumuladas (mm)')
plt.legend(['Cluster 0', 'Cluster 1', 'Cluster 2'])
plt.show()
```



Los puntos Rojos pertenecen al Cluster 0, los Azules al Cluster 1 y los Amarillos al Cluster 2. Viendo esta visualización podemos intuir que la cantidad de precipitaciones ha jugado un papel importante en la división de los clusters. Pero vamos a asegurarnos y ver como han quedado divididos los meses y años con el análisis estadístico de los clusters.

4.5 Estudio estadístico de los clusters

```
In [83]:

labels = pd.Series(labels)

In [84]:

df['Labels'] = labels

In [85]:

df['Labels'].value_counts(normalize=True)*100

Out[85]:

0    45.815003
1    41.879433
2    12.304965
Name: Labels, dtype: float64
```

Número d'observaciones asignadas a cada cluster:

Cluster 0:	1292	45.82%
Cluster 1:	1181	42.88%
Cluster 2:	347	12.30%

Agrupamos por clusters

```
In [86]:

label_grp = df.groupby(['Labels'])
```

4.5.1 Precipitaciones

```
In [87]:

df['Precipitaciones'].mean()

Out[87]:

49.4150354609924

In [88]:

label_grp['Precipitaciones'].describe()

Out[88]:
```

	count	mean	std	min	25%	50%	75%	max
Labels								
0	12920	39.951625	31.275677	0.0	14.8	35.0	58.35	163.6
1	11810	34.452667	24.684668	0.0	14.2	30.4	51.40	104.8
2	3470	135.570029	49.610654	75.0	101.0	122.7	154.10	365.8

Podemos observar, a parte de la dispersión de registros en cada uno de los clusters, una desigualdad en las medias de precipitaciones. El primer (0) y el segundo (1) cluster tienen unas medias parecidas (39.95 y 34.45mm respectivamente), aunque inferiores a la media de la población (49.41). El tercer cluster (2) se destaca, con una media de 135.57mm y una desviación estándar de 49.61.

Es significativo observar que mientras el mínimo de los dos primeros clusters es 0, el del tercero es 75. Los máximos siguen un patrón parecido (1) 163.6 (2) 104.8 (3) 365.8.

Vamos a investigar como han quedado repartidos los meses y años para cada uno.

4.5.2 Meses

En el df original los valores de los meses son proporcionalmente iguales (=0.83333). A continuación podemos ver como han quedado repartidos por Clusters:

```
In [89]:

round(label_grp['Mes'].value_counts(normalize=True)*100,2)

Out[89]:
```

Labels	Mes	count	mean	std	min	25%	50%	75%	max
0	1	12	8.82						
	2	12	8.75						
	7	8.75							
	8	8.67							
	5	8.67							
	4	8.59							
	6	8.51							
	8	8.44							
	3	8.36							
	11	7.97							
	10	7.59							
	9	6.89							
1	7	12	8.74						
	6	9.31							
	8	9.14							
	1	9.06							
	2	8.98							
	12	8.95							
	4	8.30							
	9	7.96							
	5	7.71							
	11	7.62							
	9	7.28							
	10	6.35							
2	10	17.87							
	9	17.29							
	11	12.10							
	3	9.51							
	5	8.22							
	4	8.07							
	12	7.76							
	2	4.90							
	8	4.61							
	1	4.32							
	6	4.32							
	7	2.02							

Name: Mes, dtype: float64

Conociendo como se reparten las cantidades de precipitaciones anuales por meses podemos inferir qué lógica ha seguido el proceso de clustering no-supervisado.

Viendo la repartición proporcional de meses por cluster sabemos que si la aparición de ese mes en el cluster es superior a 8.33%, estará por encima de la media de la población, mientras que si es inferior estará por debajo.

En el primer cluster(0) los valores están repartidos casi equitativamente, destacan por debajo de la media:

Septiembre:	6.88%
Octubre:	7.58%
Noviembre:	7.97%

Y por encima, ningún mes aparece más de un 0.5% que la media de la población (están por debajo de 8.83%). Los más presentes son:

Diciembre:	8.82%
Enero:	8.75%
Julio:	8.75%

Conc. Vemos como los meses con menos presencia coinciden con los meses históricos de más precipitaciones.

En el segundo cluster(1) hay una mayor variabilidad en la presencia proporcional de cada mes, tanto por encima como por debajo de la media. Destacan por debajo:

Octubre:	6.25%
Septiembre:	7.28%
Noviembre:	7.62%

Y por encima de la media:

Julio:	9.74%
Junio:	9.31%
Agosto:	9.14%

Conc. En el caso del segundo cluster vemos que la distribución de los meses es parecida a la del primero, pero acentuando su caso. Los meses que menos aparecen corresponden también a los más lluviosos históricamente y los que más aparecen, a los menos lluviosos.

En el tercer cluster(2) las distribuciones son totalmente disparas, teniendo meses que aparecen hasta ocho veces más que otros (Octubre/Julio). Tenemos a octubre (17.87%) y septiembre (17.29%), que doblan la media, y noviembre significativamente por encima (12.10%).

El menos representado es julio (2.02%), seguido por junio, enero, agosto y febrero (todos alrededor del 4%).

Conc. En el tercer cluster se observa una tendencia clara a repetir aquellos registros de los meses que están asociados históricamente con la mayor cantidad de precipitaciones.

4.5.1 Años

```
In [90]:

len(df['Any'].unique())

Out[90]:

235

La media de la población en este caso es igual a 1/número de años distintos = 1/235 = 0.00425 = 0.42%

In [91]:

label_grp['Any'].describe()

Out[91]:
```

	count	mean	std	min	25%	50%	75%	max
Labels								
0	12920	1841.319659	32.945499	1786.0	1813.0	1840.0	1869.0	1902.0
1	11810	1962.616427	34.023512	1901.0	1940.0	1963.0	1992.0	2020.0
2	3470	1929.755043	49.542362	1786.0	1897.0	1926.0	1969.0	2020.0

```
In [92]:

round(label_grp['Any'].value_counts(normalize=True)*100,2).unique()

Out[92]:

array([11.44, 1.15, 0.86, 0.58, 0.29])

In [93]:

round(label_grp['Any'].value_counts(normalize=True)*100,2)[2].value_counts()

Out[93]:
```

	count	mean	std	min	25%	50%	75%	max
Labels								
0	2.29	56						
1	0.58	52						
2	0.86	33						
3	1.15	12						
4	1.44	8						
Name: Any, dtype: int64								

```
In [94]:

grp_144 = [1901, 1906, 1907, 1921, 1923, 1943, 1971, 1996]
```

```
In [95]:

df[df['Any'] == 1901].describe()

Out[95]:
```

	Any	Mes	Precipitaciones	Labels
count	120	12.000000	12.000000	12.000000
mean	1901.0	6.500000	85.883333	1.000000
std	0.0	3.605551	69.310642	0.953463
min	1901.0	1.000000	9.200000	0.000000
25%	1901.0	3.750000	31.000000	0.000000
50%	1901.0	6.500000	65.850000	1.000000
75%	1901.0	9.250000	130.750000	2.000000
max	1901.0	12.000000	209.300000	2.000000

```
In [96]:

df.isin({'Any': [grp_144] }).value_counts()

Out[96]:
```

	Any	Mes	Precipitaciones	Labels
False	False	False	False	2820
dtype:	int64			

```
In [99]:

round(label_grp['Any'].value_counts(normalize=True)*100,2)[2].head()

Out[99]:
```

	Any	Mes	Precipitaciones	Labels
1901	1.44			
1906	1.44			
1907	1.44			
1921	1.44			
1923	1.44			
Name: Any, dtype: float64				

```
In [100]:

round(label_grp['Any'].value_counts(normalize=True)*100,2)[2].keys().sort_values()

Out[100]:
```

	Any	Mes	Precipitaciones	Labels																	
int64[1]	1786	1787	1790	1793	1810	1811	1820	1840	1843	1844	...	2010	2011	2012	2013	2014	2017	2018	2019	2020	
dtype:	int64											int64									
name:	Any											Any									
length:	161																				

```
In [101]:

len((label_grp['Any'].value_counts(normalize=True)*100)[0].keys())

Out[101]:

17
```

```
In [102]:

len((label_grp['Any'].value_counts(normalize=True)*100)[1].keys())

Out[102]:

120
```

```
In [103]:

len((label_grp['Any'].value_counts(normalize=True)*100)[2].keys())

Out[103]:

161
```

Esta variable es más difícil de describir porque tiene 235 valores distintos. Del mismo modo que antes, compararemos su proporción en cada cluster con la media de la población. Podemos ver que en los 3 casos ocurre un fenómeno similar. La mayoría de los años aparecidos en cada cluster tiene una proporción superior o muy superior a la de la media. Esta hecho nos indica que es muy probable que no aparezcan registros de todos los años en cada cluster es decir, que los años han quedado repartidos por clusters.

Años distintos por cluster (sobre 235):

Cluster 0:	117
Cluster 1:	120
Cluster 2:	161

Ordenando por fecha ascendente vemos en qué rango están los años de cada cluster:

Cluster 0:	1786 – 1902
Cluster 1:	1901 – 2020
Cluster 2:	1786 – 2020

El tercer cluster incluye todos los años pero tiene una distribución particular de estos. Los años quedan agrupados en 5 tipos de distribuciones: 1.44, 1.15, 0.86, 0.58, 0.29. Recordamos que la media era 0.42%. Tenemos cuatro grupos que aparecen por encima la media, destacando por encima el de 1.44. Hay un grupo por debajo la media, 0.29.

Tipo	años distintos
0.29	56
0.58	52
0.86	33
1.15	12
1.44	8

4.6 Conclusiones estudio de clusters

Después de analizar los valores de cada una de las variables en cada cluster podemos inferir que el algoritmo no-supervisado ha dividido:

1. Por año: en el caso del primer cluster(0) y segundo cluster(1) hay una clara división en los datos. Los años incluidos en cada uno según un rango temporal. Siendo el del primero de 1786 a 1902 y el del segundo de 1901 a 2020. Esta división no se ha aplicado al tercer cluster, que incluía valores de todos los años. Ahora bien,
2. en el caso del tercer cluster, teniendo en cuenta la variable 'Año', hemos encontrado otro patrón. Aún conteniendo registros de todos los años, sus distribuciones estaban repartidas en 5 tipos. Podríamos decir que hemos encontrado clusters dentro de un cluster. Vemos una tendencia, hay más años con menos presencia (56 con 0.29, poco debajo la media) y viceversa (8


```
[117..
models = [{"DecTree": RandomForestRegressor(),
            "DecTree": DecisionTreeRegressor(),
            "XGBBoost": xgb.XGBRegressor(),
            "LinearReg": LinearRegression()]}

results = {}
names = []

for name, model in models:
    kfold = model_selection.KFold(n_splits=10, random_state=2, shuffle = True)
    cv_results = model_selection.cross_val_score(model, X, y, cv=kfold)
    results.append(cv_results)
    names.append(name)
    msg = "%s: %f (%f) %s" % (name, cv_results.mean(), cv_results.std())
    print(msg)

fig = plt.figure()
fig.suptitle("Algorithm Comparison")
ax = fig.add_subplot(111)
plt.boxplot(results)
ax.set_xticklabels(names)
plt.show()
```

RForest: 0.543114 (0.057645)
DecTree: 0.199874 (0.181715)
XGBost: 0.468808 (0.082404)
LinearReg: 0.587405 (0.063472)

Algorithm Comparison

Model	Score (Mean)	Score (Std)
RForest	0.543114	0.057645
DecTree	0.199874	0.181715
XGBost	0.468808	0.082404
LinearReg	0.587405	0.063472



	Actual	Predicted
--	--------	-----------

```

1990    88.5    32.573583
2041    17.4    139.723966
2261    39.3    43.443573
2683    0.4    147.712848
2201    26.6    39.537782
2701    24.4    34.730412
2282    10.2    45.124738
2063    35.3    42.293652
2266    4.3    36.972625
2519    84.4    27.963137

```

```

In [128]: print('Mean Absolute Error:', metrics.mean_absolute_error(y_test, lr_pred).round(2))
          print('Mean Squared Error:', metrics.mean_squared_error(y_test, lr_pred).round(2))
          print('Root Mean Squared Error:', np.sqrt(metrics.mean_squared_error(y_test, lr_pred)).round(2))

Mean Absolute Error: 23.66
Mean Squared Error: 855.84
Root Mean Squared Error: 29.25

```

5.4 Random Forest Regressor

```

In [129]: forecast_col = 'Precipitaciones'

```

```

In [130]: X = df2.drop([forecast_col], 1)
          y = df2[forecast_col]

```

```

In [131]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2)

```

```

In [132]: steps = [('scaler', StandardScaler()), ('rf', RandomForestRegressor())]
          pipeline = Pipeline(steps)

```

```

In [133]: parameters = {'rf_n_estimators': [100, 500, 700], 'rf_max_features': ['auto',
          'sqrt', 'log2']}

```

```

In [134]: grid = GridSearchCV(pipeline, param_grid=parameters, cv=5)

```

```

In [135]: grid.fit(X_train, y_train)
          print("score = %.4f" % (grid.score(X_test, y_test)))
          print(grid.best_params_)

score = 0.5571
('rf_max_features': 'log2', 'rf_n_estimators': 100)

```

```

In [136]: rf_pred = grid.predict(X_test)

```

```

In [137]: df_rf = pd.DataFrame({'Actual': y_test, 'Predicted': rf_pred})
          df_rf[:10]

```

```

Out[138]:
          Actual  Predicted
1990    88.5    146.037
2041    17.4    58.685
2261    39.3    21.786
2683    0.4    151.571
2201    26.6    32.710
2701    24.4    30.853
2282    10.2    29.998
2063    35.3    121.351
2266    4.3    43.467
2519    84.4    20.891

```

```

In [138]: print('Mean Absolute Error:', metrics.mean_absolute_error(y_test, rf_pred).round(2))
          print('Mean Squared Error:', metrics.mean_squared_error(y_test, rf_pred).round(2))
          print('Root Mean Squared Error:', np.sqrt(metrics.mean_squared_error(y_test, rf_pred)).round(2))

Mean Absolute Error: 23.82

```

5.4 Neural network

```
In [139]: forecast_col = 'Precipitaciones'
```

```
In [140]: x = df2.drop([forecast_col],1)
y = df2[forecast_col]
```

```
In [141]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2)
```

```
In [142]: scaler = StandardScaler()
scaler.fit(X_train)
StandardScaler(copy=True, with_mean=True, with_std=True)
X_train = scaler.transform(X_train)
X_test = scaler.transform(X_test)
```

```
In [143]: mlp = MLPRegressor(hidden_layer_sizes=(6,6,6),max_iter=500)
```

```
In [144]: mlp.fit(X_train,y_train)
```

```
Out[144]: MLPRegressor(hidden_layer_sizes=(6, 6, 6), max_iter=500)
```

```
In [145]: mlp.score(X_train,y_train)
```

```
Out[145]: 0.6130842436614741
```

```
In [146]: mlp_pred = mlp.predict(X_test)
```

```
In [150]: df_mlp=pd.DataFrame({'Actual':y_test, 'Predicted':mlp_pred})
df_mlp[10]
```

```
Out[150]:
```

	Actual	Predicted
1990	88.5	142.947632
2041	17.4	53.033183

2683	0.4	28.965057
2201	26.6	35.724590

2701	24.4	37.210239
2282	10.2	28.454192
2063	35.3	47.709205
2266	4.3	38.692868
2519	84.4	32.948265

```
In [154]: print('Mean Absolute Error:', metrics.mean_absolute_error(y_test, mlp_pred).round(2))
print('Mean Squared Error:', metrics.mean_squared_error(y_test, mlp_pred).round(2))
print('Root Mean Squared Error:', np.sqrt(metrics.mean_squared_error(y_test, mlp_pred)).round(2))

Mean Absolute Error: 24.57
Mean Squared Error: 922.53
Root Mean Squared Error: 30.37

5.4.1 Hyperparameter tuning
```

```
In [155]: mlp = MLPRegressor(max_iter=500)
parameter_space = {
    'hidden_layer_sizes': [(50,50,50), (50,100,50), (100,)],
    'activation': ['tanh', 'relu'],
    'solver': ['sgd', 'adam'],
    'alpha': [0.0001, 0.05],
    'learning_rate': ['constant', 'adaptive'],
}

In [156]: clf = GridSearchCV(mlp, parameter_space, n_jobs=-1, cv=3)
clf.fit(X_train, y_train)

Out[156]: GridSearchCV(cv=3, estimator=MLPRegressor(max_iter=500, n_jobs=-1,
            param_grid={'activation': ['tanh', 'relu'],
                        'alpha': [0.0001, 0.05],
                        'hidden_layer_sizes': [(50, 50, 50), (50, 100, 50),
                                                (100,)],
                        'learning_rate': ['constant', 'adaptive'],
                        'solver': ['sgd', 'adam']}))

In [157]: # Best parameter set
print('Best parameters found:\n', clf.best_params_)

Best parameters found:
{'activation': 'relu', 'alpha': 0.0001, 'hidden_layer_sizes': (100,), 'learning_rate': 'adaptive', 'solver': 'sgd'}

In [158]: # All results
means = clf.cv_results_['mean_test_score']
stds = clf.cv_results_['std_test_score']
for mean, std, params in zip(means, stds, clf.cv_results_['params']):
    print("%0.15f (+/-0.03f) for %r" % (mean, std, 2. * params))

0.595 (+/-0.111) for {'activation': 'tanh', 'alpha': 0.0001, 'hidden_layer_sizes': (50, 50, 50), 'learning_rate': 'constant', 'solver': 'sgd'}
0.158 (+/-0.066) for {'activation': 'tanh', 'alpha': 0.0001, 'hidden_layer_sizes': (50, 50, 50), 'learning_rate': 'constant', 'solver': 'adam'}
0.373 (+/-0.155) for {'activation': 'tanh', 'alpha': 0.0001, 'hidden_layer_sizes': (50, 50, 50), 'learning_rate': 'adaptive', 'solver': 'sgd'}
0.178 (+/-0.049) for {'activation': 'tanh', 'alpha': 0.0001, 'hidden_layer_sizes': (50, 50, 50), 'learning_rate': 'adaptive', 'solver': 'sgd'}
0.595 (+/-0.107) for {'activation': 'tanh', 'alpha': 0.0001, 'hidden_layer_sizes': (50, 100, 50), 'learning_rate': 'constant', 'solver': 'sgd'}
0.185 (+/-0.030) for {'activation': 'tanh', 'alpha': 0.0001, 'hidden_layer_sizes': (50, 100, 50), 'learning_rate': 'constant', 'solver': 'adam'}
0.583 (+/-0.077) for {'activation': 'tanh', 'alpha': 0.0001, 'hidden_layer_sizes': (50, 100, 50), 'learning_rate': 'adaptive', 'solver': 'sgd'}
0.160 (+/-0.057) for {'activation': 'tanh', 'alpha': 0.0001, 'hidden_layer_sizes': (50, 100, 50), 'learning_rate': 'constant', 'solver': 'adam'}
0.578 (+/-0.116) for {'activation': 'tanh', 'alpha': 0.0001, 'hidden_layer_sizes': (100,), 'learning_rate': 'constant', 'solver': 'sgd'}
0.562 (+/-0.100) for {'activation': 'tanh', 'alpha': 0.0001, 'hidden_layer_sizes': (100,), 'learning_rate': 'constant', 'solver': 'adam'}
```

```

aptive', 'solver': 'sgd
0.563 (+/-0.105) for {'
aptive', 'solver': 'ada
0.592 (+/-0.134) for {'

```

```

'constant', 'solver': 'sgd'
'tanh', 'alpha': 0.05, 'hidden_layer_sizes': (50, 50, 50), 'learning_rate':
'adaptive', 'solver': 'adam'
0.601 (-0.0795) for 'Activation': 'tanh', 'alpha': 0.05, 'hidden_layer_sizes': (50, 50, 50), 'learning_rate':
'adaptive', 'solver': 'sgd'
0.159 (-0.0335) for 'Activation': 'tanh', 'alpha': 0.05, 'hidden_layer_sizes': (50, 50, 50), 'learning_rate':
'adaptive', 'solver': 'adam'
0.601 (-0.0988) for 'Activation': 'tanh', 'alpha': 0.05, 'hidden_layer_sizes': (50, 50, 50), 'learning_rate':
'constant', 'solver': 'sgd'
0.420 (-0.0484) for 'Activation': 'tanh', 'alpha': 0.05, 'hidden_layer_sizes': (50, 100, 50), 'learning_rate':
'adaptive', 'solver': 'adam'
0.178 (-0.0474) for 'Activation': 'tanh', 'alpha': 0.05, 'hidden_layer_sizes': (50, 100, 50), 'learning_rate':
'adaptive', 'solver': 'sgd'
0.578 (-0.0270) for 'Activation': 'tanh', 'alpha': 0.05, 'hidden_layer_sizes': (100, ), 'learning_rate': 'constant', 'solver': 'sgd'
'tanh', 'alpha': 0.05, 'hidden_layer_sizes': (100, ), 'learning_rate': 'constant', 'solver': 'adam'
0.573 (-0.0229) for 'Activation': 'tanh', 'alpha': 0.05, 'hidden_layer_sizes': (100, ), 'learning_rate': 'adaptive', 'solver': 'sgd'
0.563 (-0.0105) for 'Activation': 'tanh', 'alpha': 0.05, 'hidden_layer_sizes': (100, ), 'learning_rate': 'adaptive', 'solver': 'adam'
0.209 (-0.0139) for 'Activation': 'relu', 'alpha': 0.0001, 'hidden_layer_sizes': (50, 50, 50), 'learning_rate': 'constant', 'solver': 'sgd'
0.584 (-0.0707) for 'Activation': 'relu', 'alpha': 0.0001, 'hidden_layer_sizes': (50, 50, 50), 'learning_rate': 'constant', 'solver': 'adam'
0.383 (-0.0106) for 'Activation': 'relu', 'alpha': 0.0001, 'hidden_layer_sizes': (50, 50, 50), 'learning_rate': 'adaptive', 'solver': 'sgd'
0.583 (-0.0119) for 'Activation': 'relu', 'alpha': 0.0001, 'hidden_layer_sizes': (50, 50, 50), 'learning_rate': 'adaptive', 'solver': 'adam'
0.571 (-0.0133) for 'Activation': 'relu', 'alpha': 0.0001, 'hidden_layer_sizes': (50, 100, 50), 'learning_rate': 'constant', 'solver': 'sgd'
0.582 (-0.0103) for 'Activation': 'relu', 'alpha': 0.0001, 'hidden_layer_sizes': (50, 100, 50), 'learning_rate': 'constant', 'solver': 'adam'
0.592 (-0.0101) for 'Activation': 'relu', 'alpha': 0.0001, 'hidden_layer_sizes': (50, 100, 50), 'learning_rate': 'adaptive', 'solver': 'sgd'
0.592 (-0.0122) for 'Activation': 'relu', 'alpha': 0.0001, 'hidden_layer_sizes': (50, 100, 50), 'learning_rate': 'adaptive', 'solver': 'adam'
0.598 (-0.0108) for 'Activation': 'relu', 'alpha': 0.0001, 'hidden_layer_sizes': (100, ), 'learning_rate': 'constant', 'solver': 'sgd'
0.592 (-0.0105) for 'Activation': 'relu', 'alpha': 0.0001, 'hidden_layer_sizes': (100, ), 'learning_rate': 'constant', 'solver': 'adam'
0.611 (-0.0102) for 'Activation': 'relu', 'alpha': 0.0001, 'hidden_layer_sizes': (100, ), 'learning_rate': 'adaptive', 'solver': 'sgd'
0.592 (-0.0107) for 'Activation': 'relu', 'alpha': 0.0001, 'hidden_layer_sizes': (100, ), 'learning_rate': 'adaptive', 'solver': 'adam'
0.553 (-0.0162) for 'Activation': 'relu', 'alpha': 0.05, 'hidden_layer_sizes': (50, 50, 50), 'learning_rate': 'constant', 'solver': 'sgd'
0.585 (-0.0094) for 'Activation': 'relu', 'alpha': 0.05, 'hidden_layer_sizes': (50, 50, 50), 'learning_rate': 'constant', 'solver': 'adam'
0.575 (-0.0134) for 'Activation': 'relu', 'alpha': 0.05, 'hidden_layer_sizes': (50, 50, 50), 'learning_rate': 'adaptive', 'solver': 'sgd'
0.582 (-0.0098) for 'Activation': 'relu', 'alpha': 0.05, 'hidden_layer_sizes': (50, 50, 50), 'learning_rate': 'adaptive', 'solver': 'adam'
0.585 (-0.0122) for 'Activation': 'relu', 'alpha': 0.05, 'hidden_layer_sizes': (50, 100, 50), 'learning_rate': 'constant', 'solver': 'sgd'
0.582 (-0.0133) for 'Activation': 'relu', 'alpha': 0.05, 'hidden_layer_sizes': (50, 100, 50), 'learning_rate': 'constant', 'solver': 'adam'
0.559 (-0.0103) for 'Activation': 'relu', 'alpha': 0.05, 'hidden_layer_sizes': (50, 100, 50), 'learning_rate': 'adaptive', 'solver': 'sgd'
0.586 (-0.0096) for 'Activation': 'relu', 'alpha': 0.05, 'hidden_layer_sizes': (50, 100, 50), 'learning_rate': 'adaptive', 'solver': 'adam'
0.585 (-0.0104) for 'Activation': 'relu', 'alpha': 0.05, 'hidden_layer_sizes': (100, ), 'learning_rate': 'constant', 'solver': 'sgd'
0.590 (-0.0105) for 'Activation': 'relu', 'alpha': 0.05, 'hidden_layer_sizes': (100, ), 'learning_rate': 'constant', 'solver': 'adam'
0.600 (-0.0136) for 'Activation': 'relu', 'alpha': 0.05, 'hidden_layer_sizes': (100, ), 'learning_rate': 'adaptive', 'solver': 'sgd'
0.590 (-0.0102) for 'Activation': 'relu', 'alpha': 0.05, 'hidden_layer_sizes': (100, ), 'learning_rate': 'adaptive', 'solver': 'adam'

```

Los mejores resultados del hyperparameter tuning no son diferentes de los que ya teníamos.

```
In [162]: results = pd.DataFrame(data)

In [162]: results.set_index('model', inplace=True)

In [163]: results

Out[163]:
```

	R2 score	MAE	MSE	RMSE
model				
LinearRegression	0.6187	26.05	1262.32	35.53
RandomForestRegressor	0.6115	24.69	1153.99	33.97
MLPRegressor	0.6286	23.08	908.52	30.14

6. Conclusiones

El primer hecho para evaluar las estimaciones de precipitaciones en un futuro es el estudio de los ciclos pasados. Para ello hemos investigado sobre el clima de Barcelona, sus variaciones históricas y las influencias - más recientes - del cambio climático. Para aproximarnos al futuro hemos estudiado los datos dividiéndolos en dos grupos: de 1786-1990 y 1990-2020.

Tanto en el primer periodo largo de tiempo como en el segundo, y considerando las características del clima de Barcelona, los meses más lluviosos han sido y serán octubre, septiembre, noviembre y abril

Creemos que el ejercicio requiere una hipótesis más definida y un periodo de predicción determinado. Aún para el caso y sin hipótesis definida, trataramos de determinar cuándo se prevén más lluvias durante los próximos 3 años. Dado que el ejercicio de regresión lineal no ha sido un gran éxito no podemos apoyarnos en sus datos para las predicciones. No obstante, gracias a otros trabajos parecidos, a la información de expertos y nuestra propia exploración trataramos de llegar a algunas conclusiones:

1. Teniendo como referencia de algún actual y próximo los datos de los últimos treinta años, podemos

esperar una
en orden descendiente
irregularidad de a
podemos esperar m

podemos esperar mayores descargas puntuales, tormentas y lluvias intensas.

2. El análisis de los clusters del modelo no-supervisado nos ha acercado a los datos y nos ha dado otro punto de referencia para comprender su distribución. Como extra, hemos podido usar las agrupaciones para mejorar las scores de los modelos supervisados.
3. El modelo supervisado deja mucho que desear con el MAE más bajo de 23.08, un margen de error inaceptable ya que hay meses los cuales, su media histórica de precipitaciones es inferior a esta métrica. Aún así, estamos satisfechos por el hecho de haber incrementado en gran medida el R2 score, que inicialmente era de -0.02, hasta 0.628.

Para incrementar las predicciones de precipitaciones recomendamos seguir el artículo *"Rainfall forecasting model using machine learning methods: Case study Terengganu, Malaysia"*, donde llegan a conseguir un R2 de 0.999 en la predicción de precipitaciones con modelos supervisados de regresión lineal.

Nos habría gustado profundizar en el estudio de los outliers, así como en el de períodos más reducidos de tiempo, cada 3 o cinco años y el estudio de meses por separado. Hubiéramos probado otros algoritmos de clusterización y habríamos trabajado más el feature engineering: obtención de nuevas variables y datos para aumentar el R2 score.

Fuentes: 1. <https://www.sciencedirect.com/science/article/pii/S20908447920302069>