

# APRENDIZAJE NO SUPERVISADO

---

Transformaciones del dataset

# Aprendizaje no supervisado

## Problemas con el aprendizaje no supervisado

- Usamos datasets sin etiquetas (porque no contamos con ellas o porque no nos interesan)
- Hay dos tipos
  - Transformaciones
  - Clustering
- Es difícil saber cuando el algoritmo aprende cosas útiles

No tenemos manera de decirle al algoritmo que estamos buscando

# Transformaciones no supervisadas

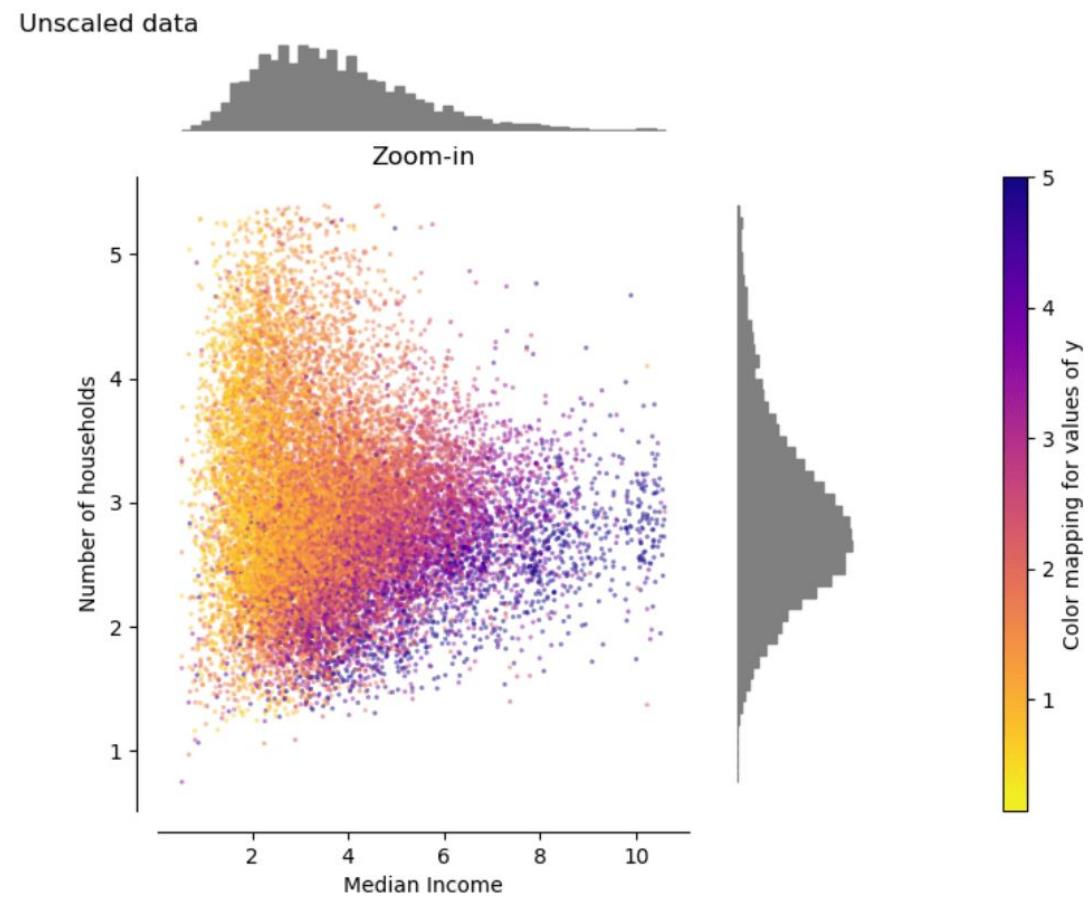
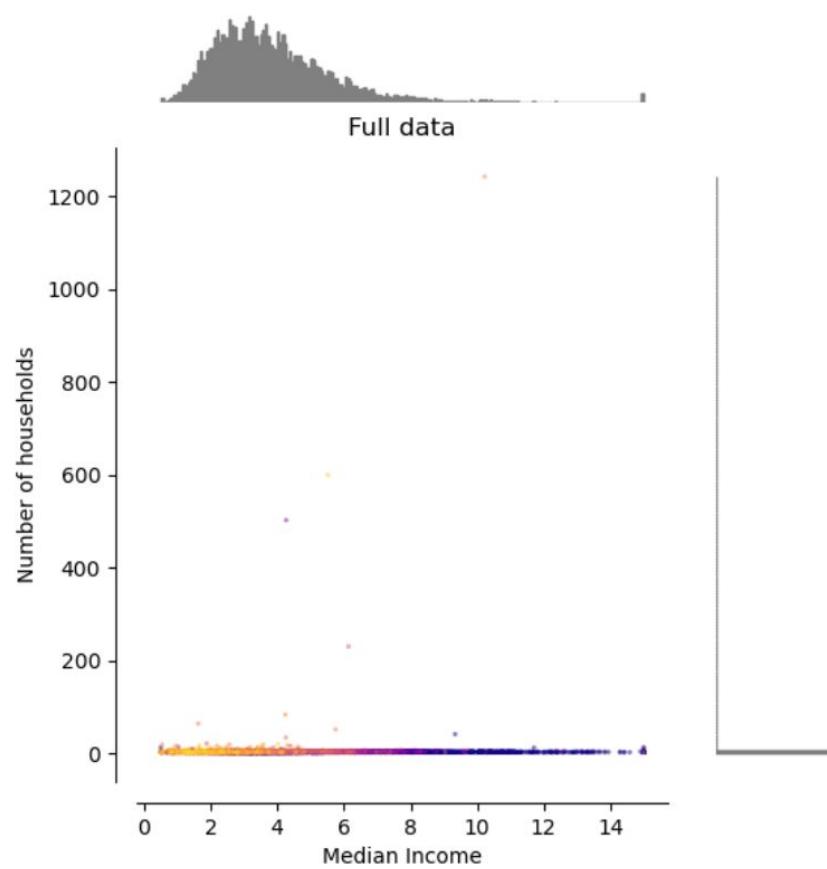
---

Crean nuevas representaciones de los datos para hacerlos más fáciles de trabajar, ya sea para las personas o para otros algoritmos.

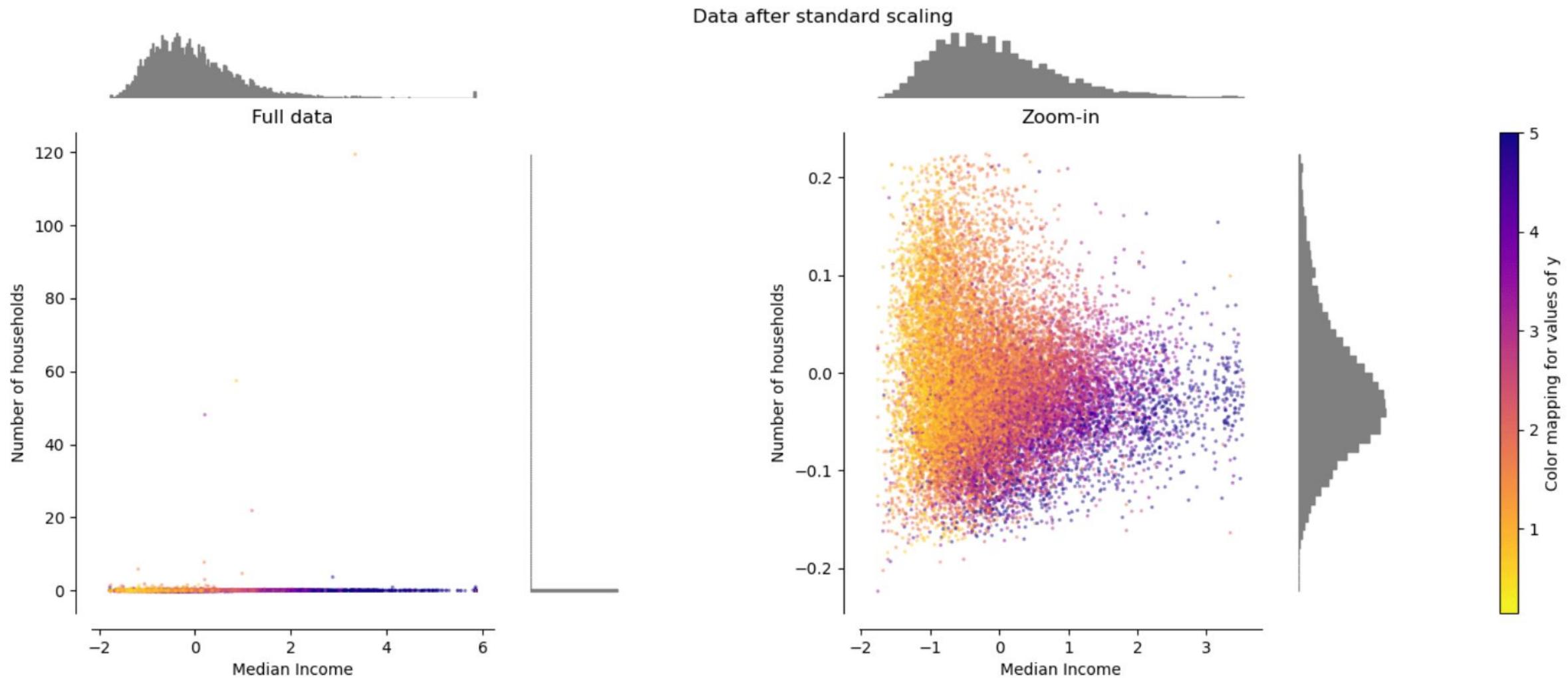
## Aplicaciones

- Reducir la dimensión
- Encontrar las partes que “componen” los datos
- Reducir memoria y tiempos

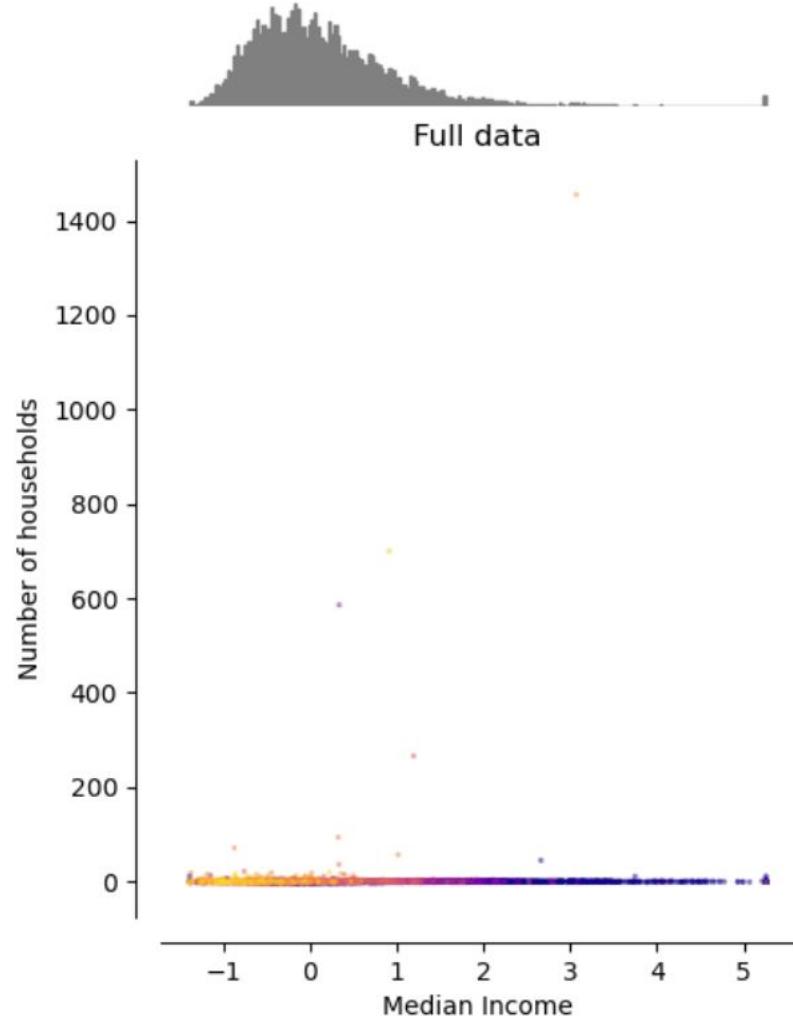
# Formas de procesar y escalar los datos



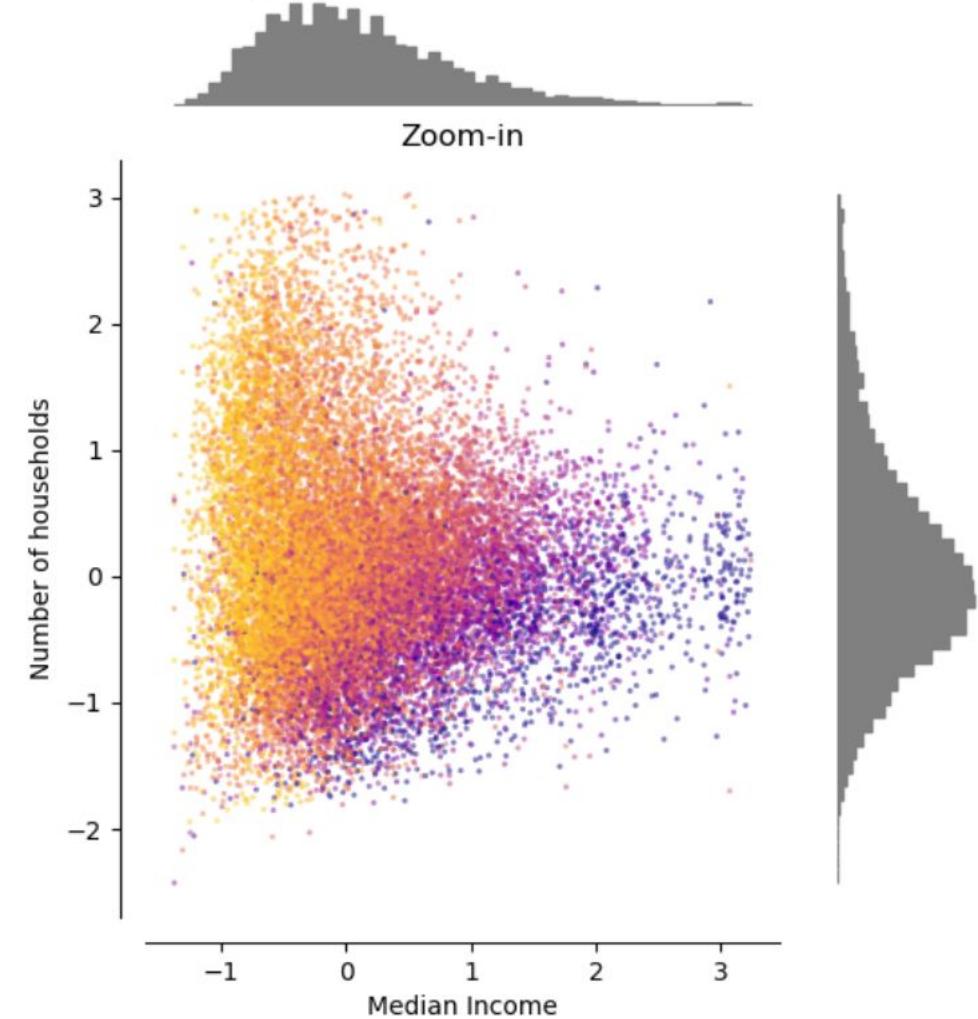
# StandardScaling



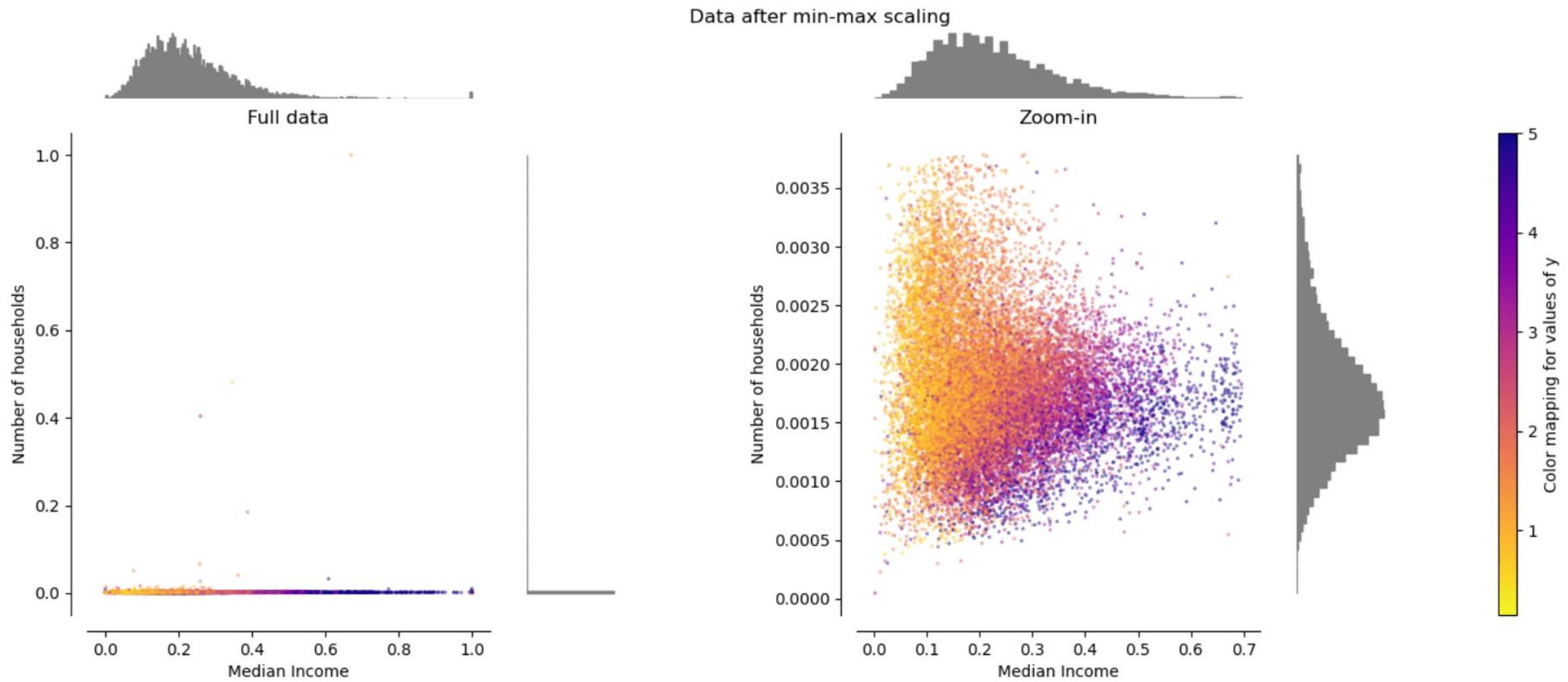
# RobustScaling



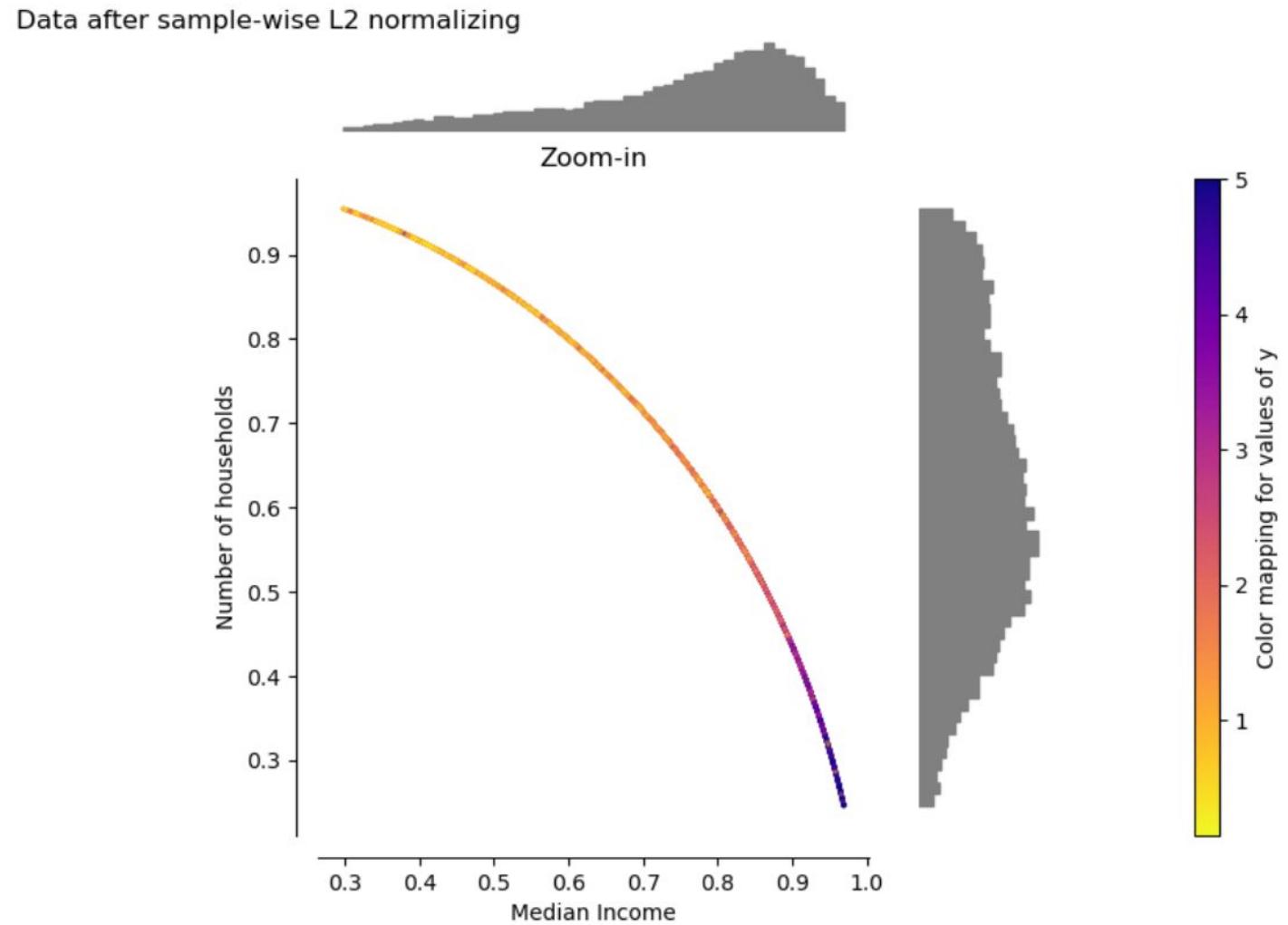
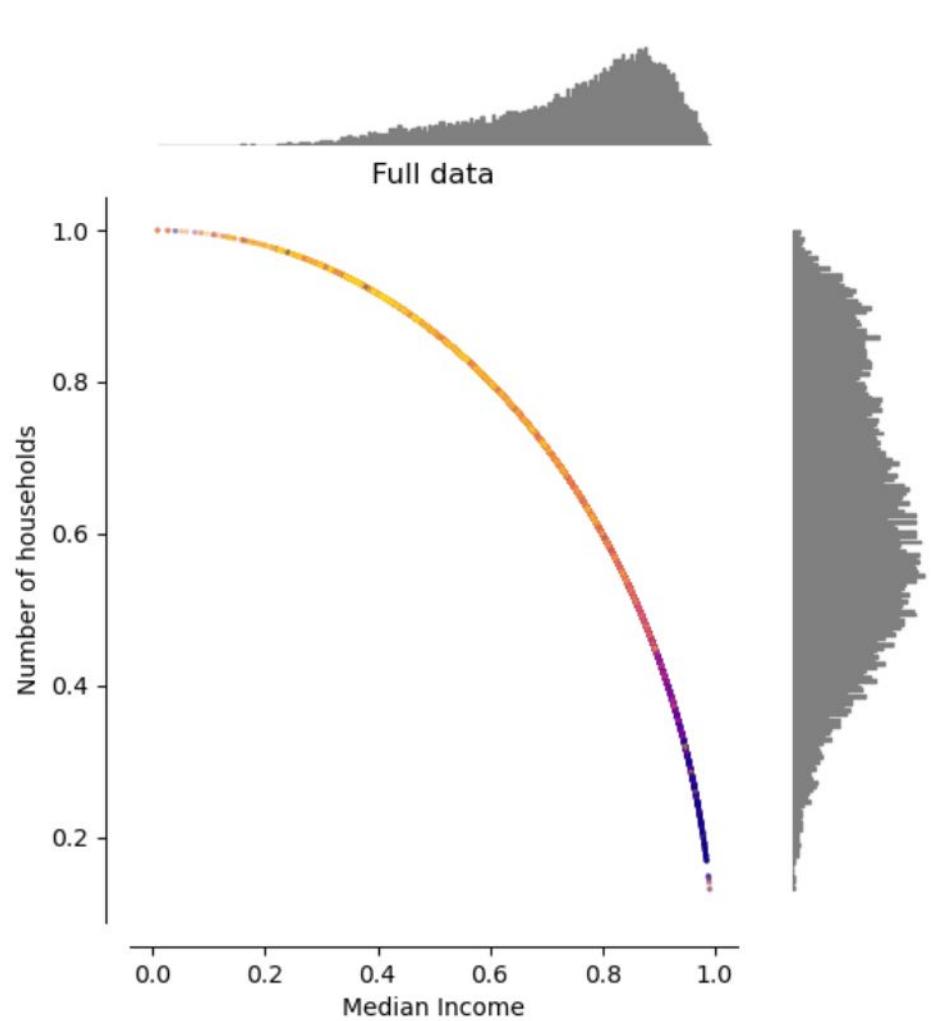
Data after robust scaling



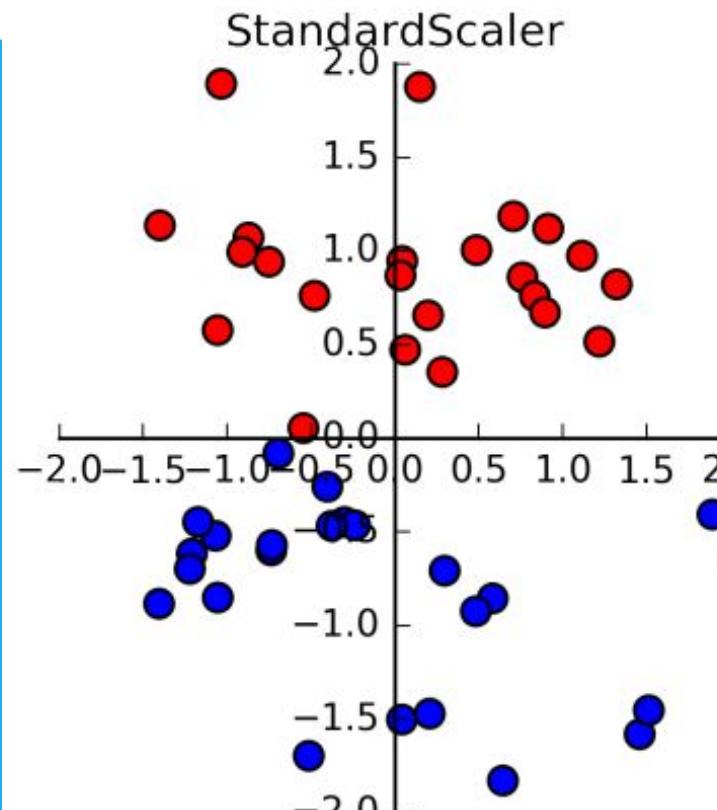
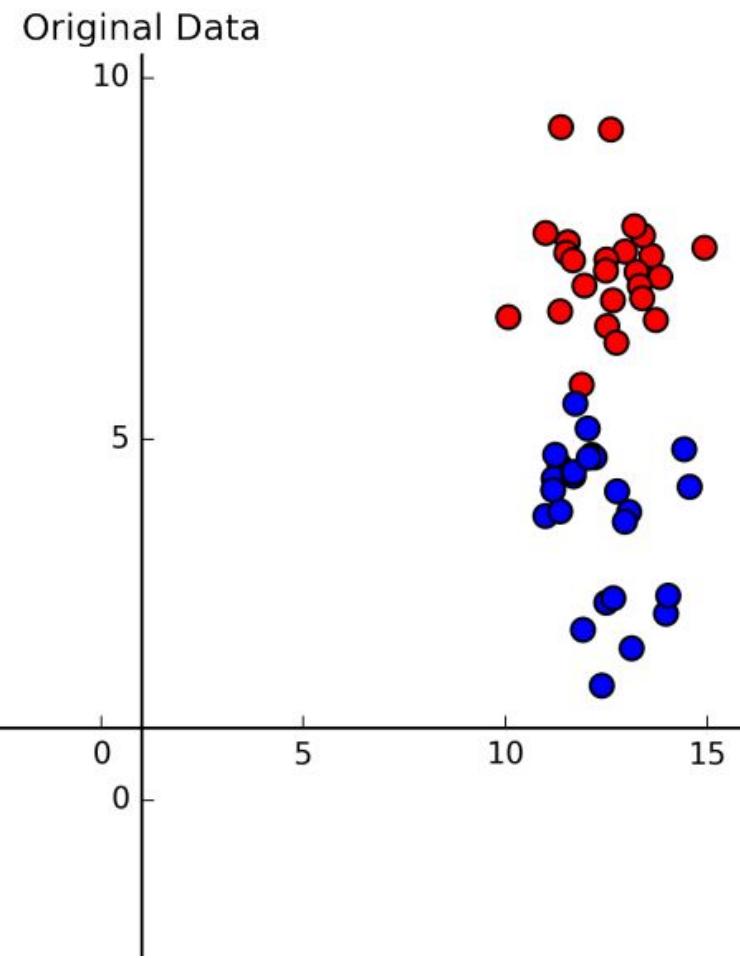
# MinMaxScaling



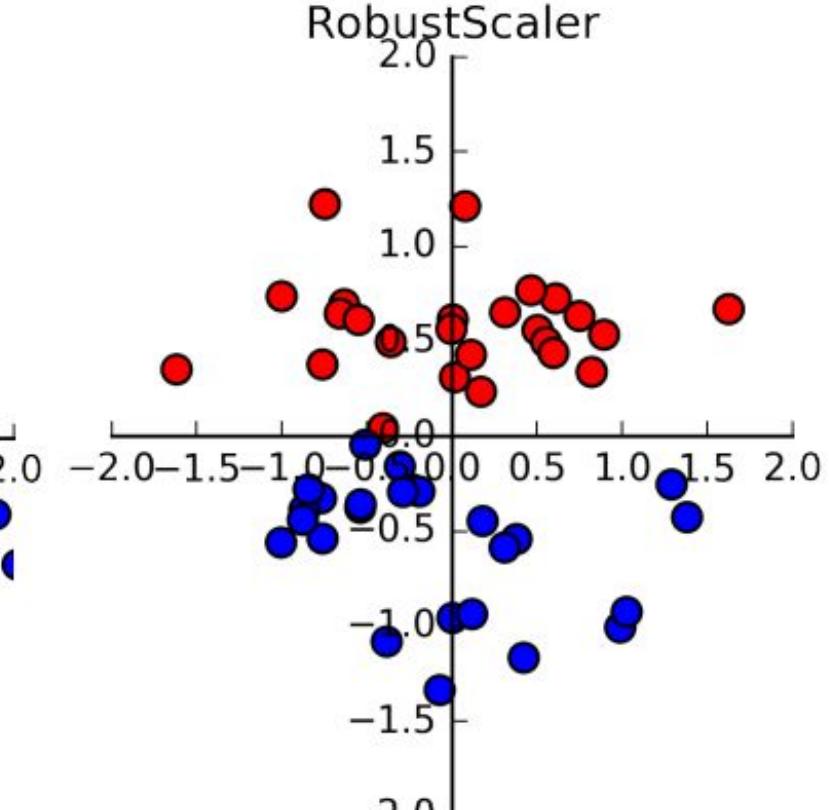
# Normilizer



# Formas de procesar y escalar los datos

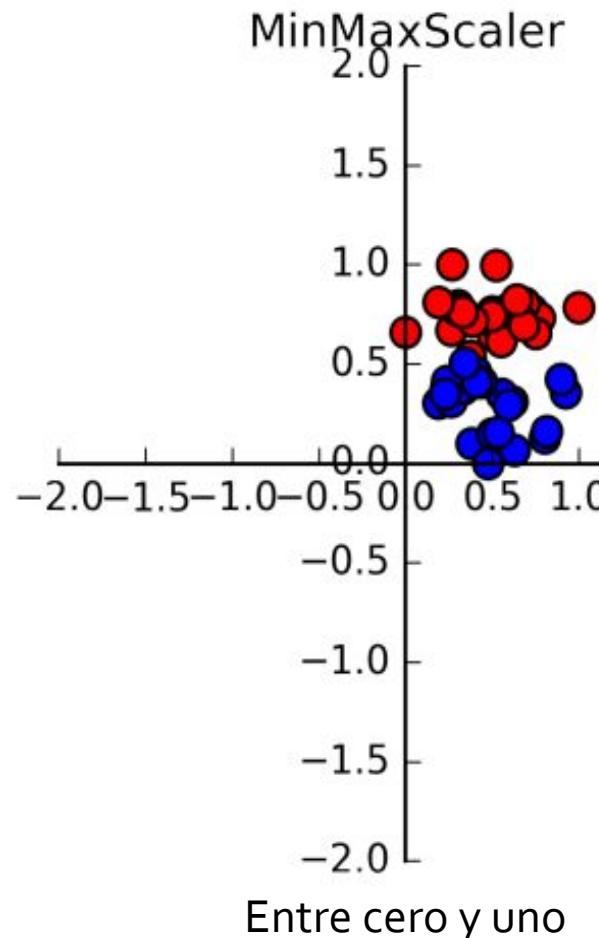
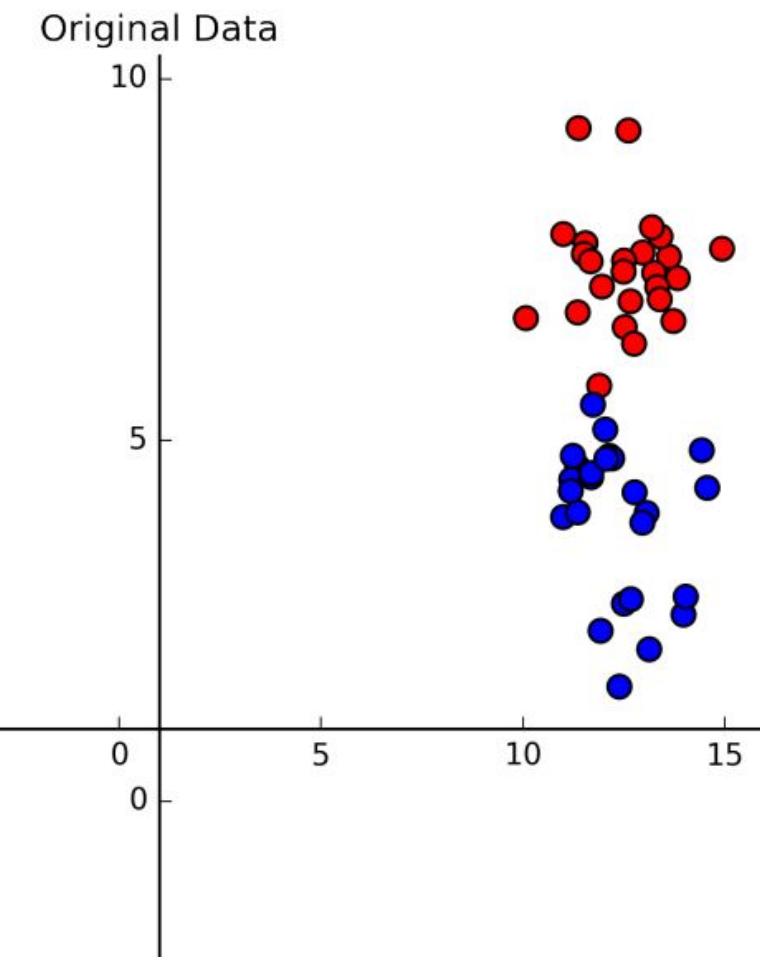


Todo alrededor del  
cero

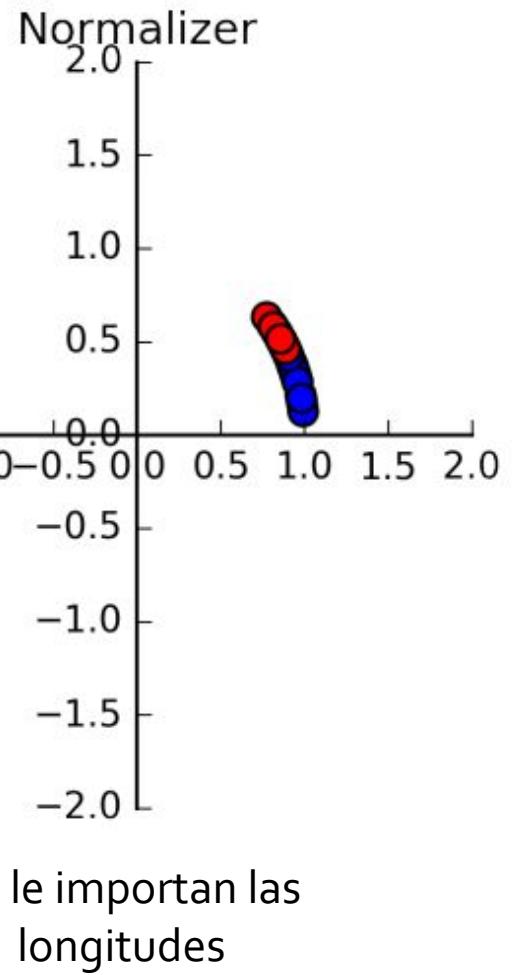


Es bueno si existen puntos  
atípicos

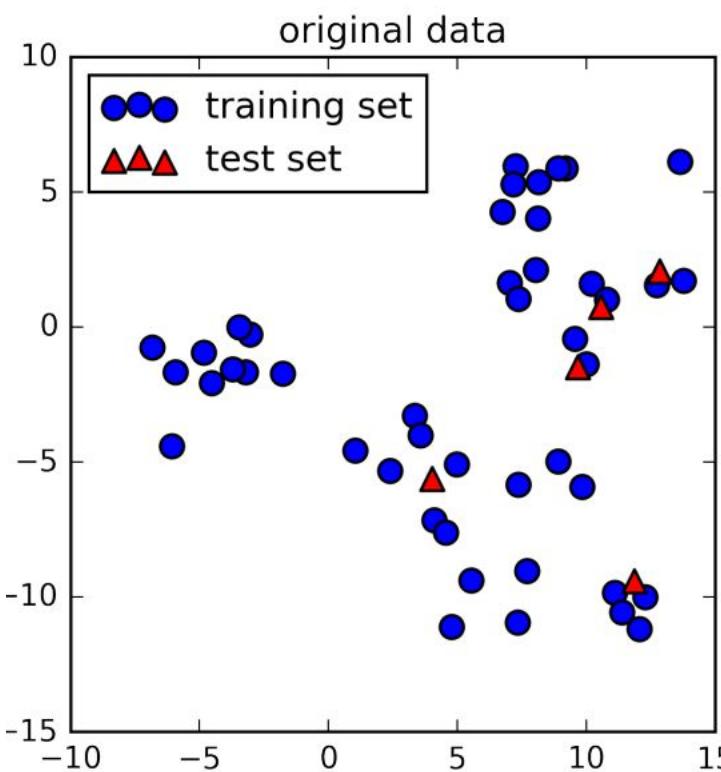
# Formas de procesar y escalar los datos



## Entre cero y uno



# Hay que transformar el conjunto prueba con cuidado



# Ejemplo de MinMaxScaler

```
from sklearn.datasets import load_breast_cancer
from sklearn.model_selection import train_test_split
cancer = load_breast_cancer()
X_train, X_test, y_train, y_test = train_test_split(cancer.data, cancer.target,
    random_state=1)
print(X_train.shape)
print(X_test.shape)
```

```
(426, 30)
(143, 30)
```

```
from sklearn.preprocessing import MinMaxScaler
scaler = MinMaxScaler()
```

```
scaler.fit(X_train)
MinMaxScaler(copy=True, feature_range=(0, 1))
```

```
MinMaxScaler()
```



# Diferencias al escalar los datos

```
from sklearn.svm import SVC
X_train, X_test, y_train, y_test = train_test_split(cancer.data, cancer.target,
    random_state=4)
svm = SVC(C=100)
svm.fit(X_train, y_train)
print(svm.score(X_test, y_test))

0.8881118881118881
```

Con los datos originales

```
# preprocessing using 0-1 scaling
scaler = MinMaxScaler()
scaler.fit(X_train)
X_train_scaled = scaler.transform(X_train)
X_test_scaled = scaler.transform(X_test)
# Learning an SVM on the scaled training data
svm.fit(X_train_scaled, y_train)
# scoring on the scaled test set
svm.score(X_test_scaled, y_test)

0.951048951048951
```

Al escalar los datos

# Escalar los datos de otra manera

Con MinMaxScaler

```
# preprocessing using 0-1 scaling
scaler = MinMaxScaler()
scaler.fit(X_train)
X_train_scaled = scaler.transform(X_train)
X_test_scaled = scaler.transform(X_test)
# Learning an SVM on the scaled training data
svm.fit(X_train_scaled, y_train)
# scoring on the scaled test set
svm.score(X_test_scaled, y_test)
```

0.951048951048951

Con StandardScaler

```
# preprocessing using zero mean and unit variance scaling
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
scaler.fit(X_train)
X_train_scaled = scaler.transform(X_train)
X_test_scaled = scaler.transform(X_test)
# Learning an SVM on the scaled training data
svm.fit(X_train_scaled, y_train)
# scoring on the scaled test set
svm.score(X_test_scaled, y_test)
```

0.958041958041958

# Escalar con RobustScaler y Normalizer

Con RobustScaler

```
# preprocessing using robustScaler
from sklearn.preprocessing import RobustScaler
scaler = RobustScaler()
scaler.fit(X_train)
X_train_scaled = scaler.transform(X_train)
X_test_scaled = scaler.transform(X_test)
# Learning an SVM on the scaled training data
svm.fit(X_train_scaled, y_train)
# scoring on the scaled test set
svm.score(X_test_scaled, y_test)
```

0.972027972027972

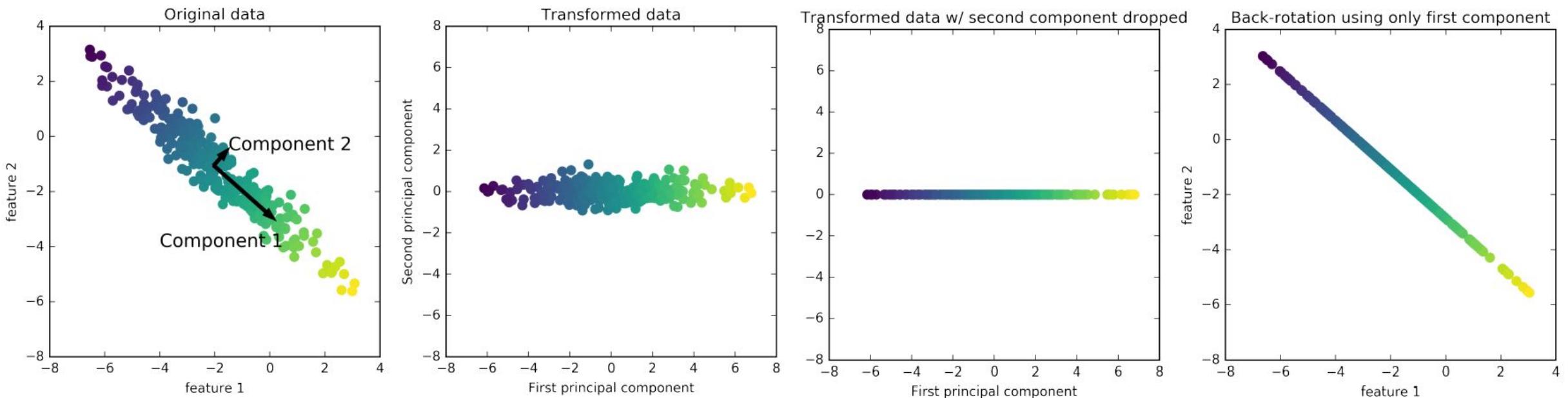
Con Normalizer

```
# preprocessing using normalizer
from sklearn.preprocessing import Normalizer
scaler = Normalizer()
scaler.fit(X_train)
X_train_scaled = scaler.transform(X_train)
X_test_scaled = scaler.transform(X_test)
# Learning an SVM on the scaled training data
svm.fit(X_train_scaled, y_train)
# scoring on the scaled test set
svm.score(X_test_scaled, y_test)
```

0.8671328671328671

# Principal Component Analysis (PCA)

## Análisis de componentes principales



# PCA para visualizar datos

```
from sklearn.datasets import load_breast_cancer
cancer = load_breast_cancer()
scaler = StandardScaler()
scaler.fit(cancer.data)
X_scaled = scaler.transform(cancer.data)
```

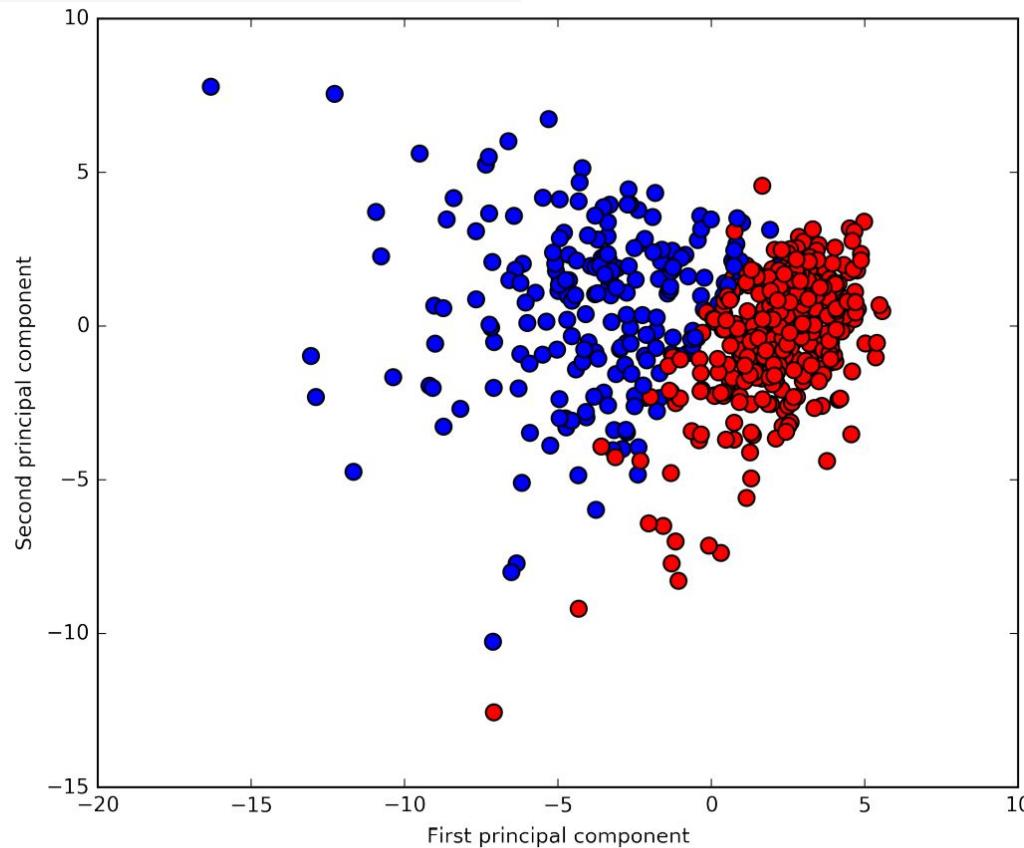
```
from sklearn.decomposition import PCA
# keep the first two principal components of the data
pca = PCA(n_components=2)
# fit PCA model to breast cancer data
pca.fit(X_scaled)
# transform data onto the first two principal components
X_pca = pca.transform(X_scaled)
print("Original shape: %s" % str(X_scaled.shape))
print("Reduced shape: %s" % str(X_pca.shape))
```

Original shape: (569, 30)

Reduced shape: (569, 2)

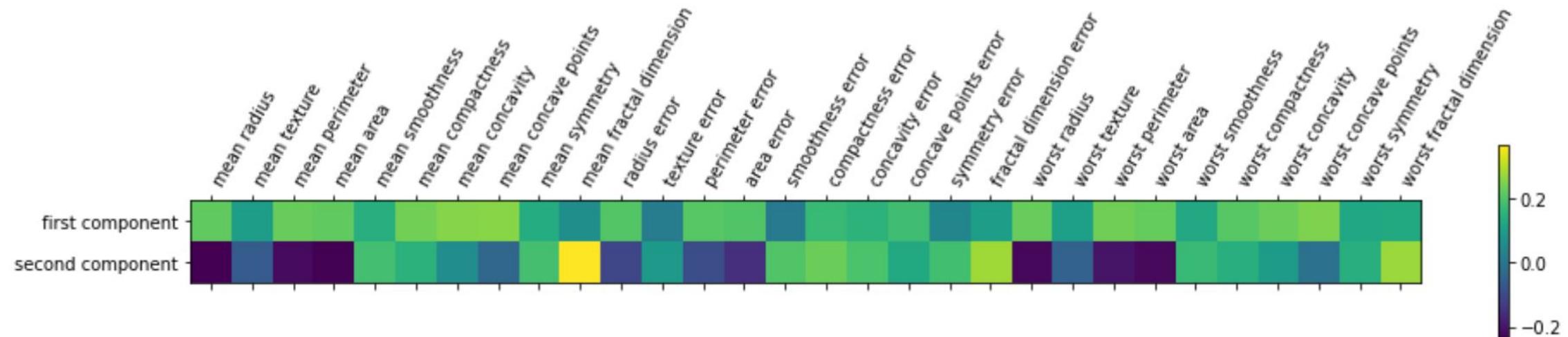
# PCA para visualizar datos

```
# plot first vs second principal component, color by class
plt.figure(figsize=(8, 8))
plt.scatter(X_pca[:, 0], X_pca[:, 1], c=cancer.target, cmap=mglearn.tools.cm, s=60)
plt.gca().set_aspect("equal")
plt.xlabel("First principal component")
plt.ylabel("Second principal component")
```



# Interpretando las componentes

```
plt.matshow(pca.components_, cmap='viridis')
plt.yticks([0, 1], ["first component", "second component"])
plt.colorbar()
plt.xticks(range(len(cancer.feature_names)),
           cancer.feature_names, rotation=60, ha='left');
# plt.suptitle("pca_components_cancer")
```

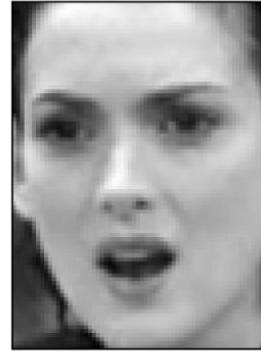


# Feature extraction

## Reducción de dimensionalidad

Es muy útiles con imágenes, donde cada pixel es una característica

Winona Ryder



Jean Chretien



Carlos Menem



Ariel Sharon



Alvaro Uribe



Colin Powell



Recep Tayyip Erdogan



Gray Davis



George Robertson



Silvio Berlusconi

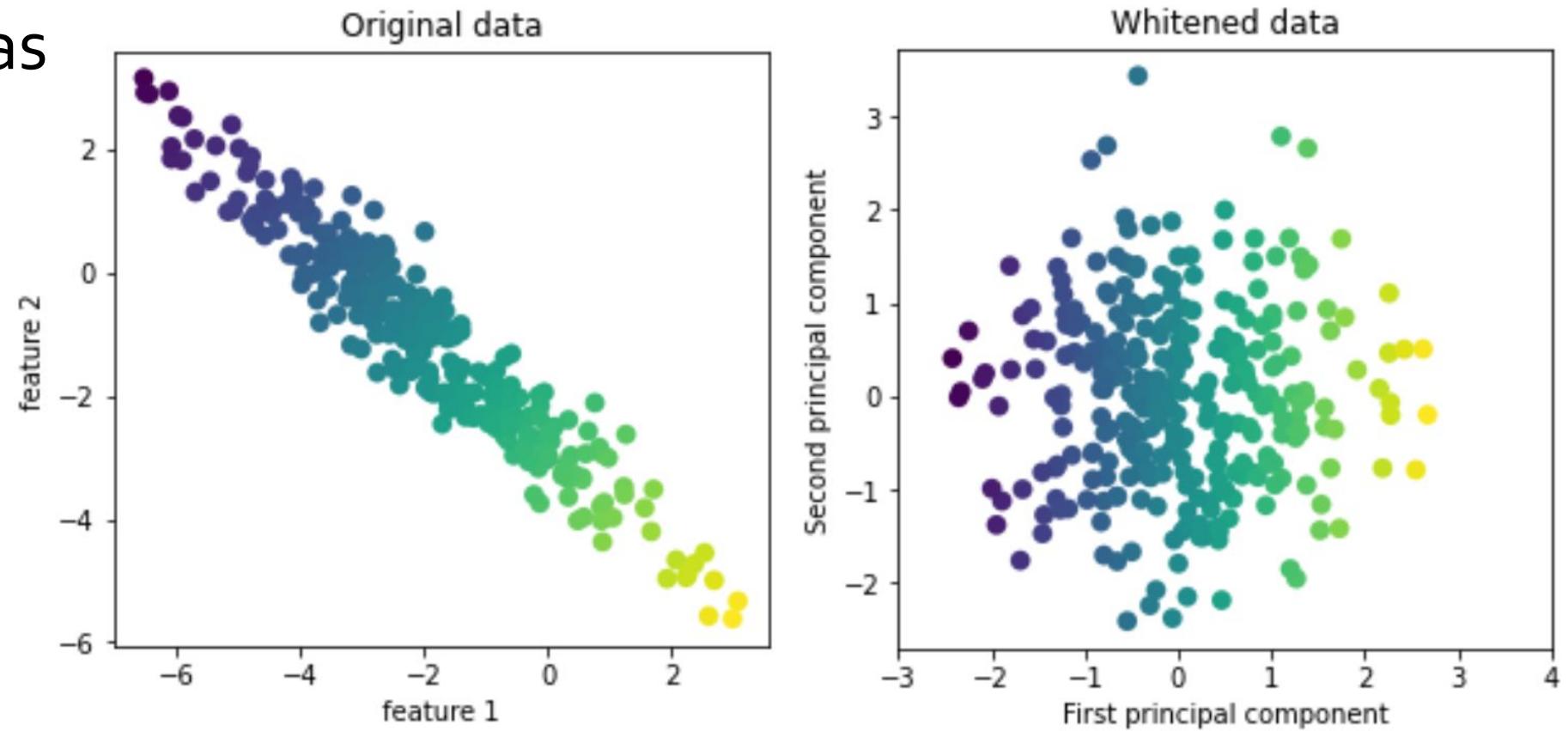
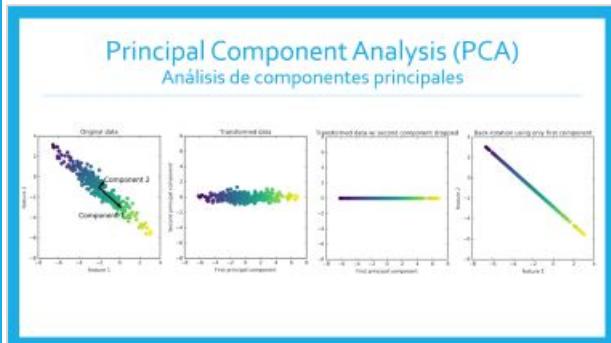


Queremos saber si una cara pertenece a  
una persona en la base de datos

¿Cómo hacemos tal cosa?

# Usando KNeighborsClassifier

Nos conviene tener buenas distancias



# Diferencias al no usar PCA y al si usarlo

```
from sklearn.neighbors import KNeighborsClassifier
# split the data in training and test set
X_train, X_test, y_train, y_test = train_test_split(
    X_people, y_people, stratify=y_people, random_state=0)
# build a KNeighborsClassifier with using one neighbor:
knn = KNeighborsClassifier(n_neighbors=1)
knn.fit(X_train, y_train)
knn.score(X_test, y_test)
```

0.23255813953488372

```
pca = PCA(n_components=100, whiten=True).fit(X_train)
X_train_pca = pca.transform(X_train)
X_test_pca = pca.transform(X_test)
print(X_train_pca.shape)
```

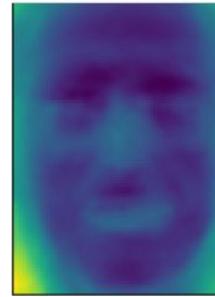
(1547, 100)

```
knn = KNeighborsClassifier(n_neighbors=1)
knn.fit(X_train_pca, y_train)
knn.score(X_test_pca, y_test)
```

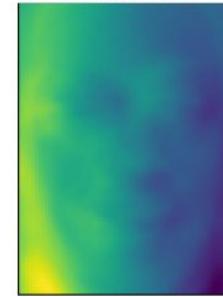
0.31007751937984496

# Interpretamos las componentes

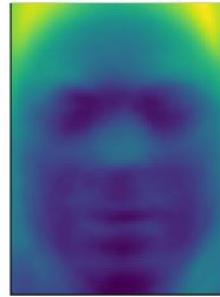
1. component



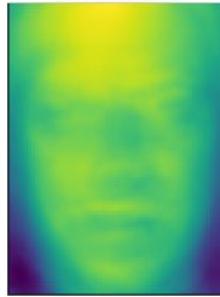
2. component



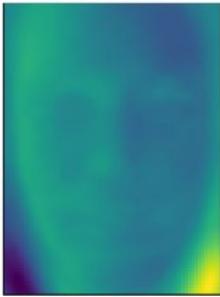
3. component



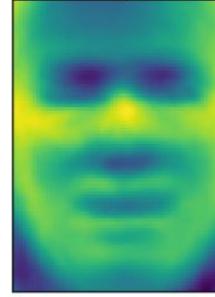
4. component



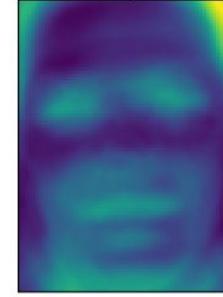
5. component



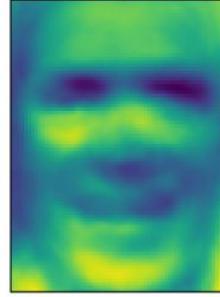
6. component



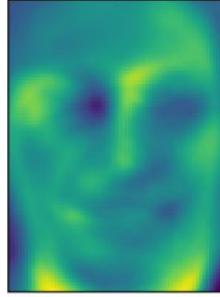
7. component



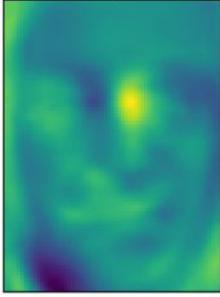
8. component



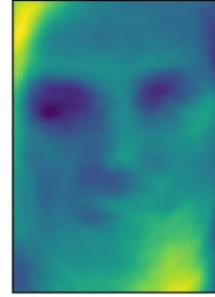
9. component



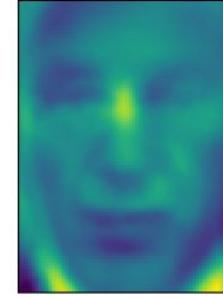
10. component



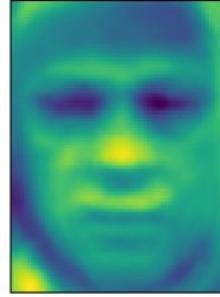
11. component



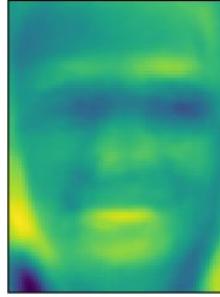
12. component



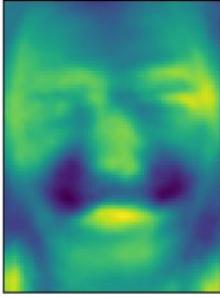
13. component



14. component

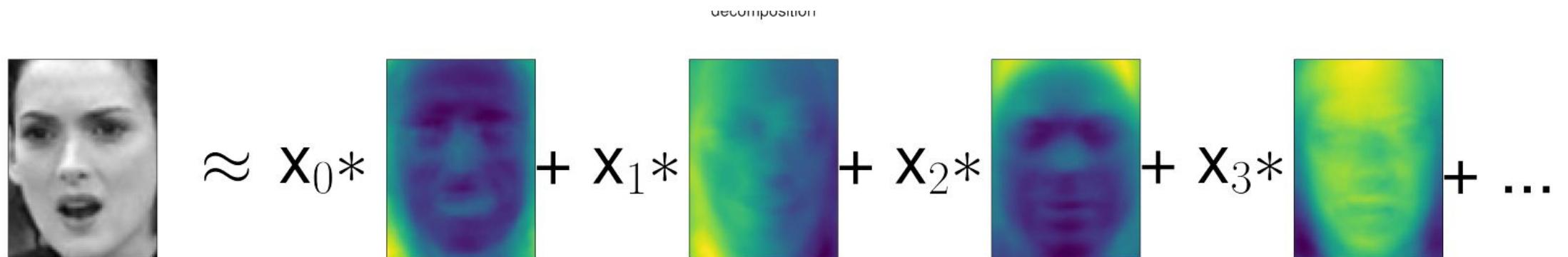


15. component



# Entendiendo el modelo de otro modo

decomposition

$$\text{Image} \approx x_0 * \text{Image} + x_1 * \text{Image} + x_2 * \text{Image} + x_3 * \text{Image} + \dots$$


# La influencia de las componentes

```
mglearn.plots.plot_pca_faces(X_train, X_test, image_shape)  
plt.suptitle("pca_reconstructions")
```



¿Podremos visualizar los datos en dos dimensiones?

# ¿Podremos visualizar los datos en dos dimensiones?

