# Objective: Learn to do clustering and noise reduction in data using PCA

In [2]:
```python
import matplotlib.pyplot as plt
import numpy as np
from numpy.linalg import svd
from sklearn.datasets import load_digits

digits = load_digits()
digits.data.shape
```
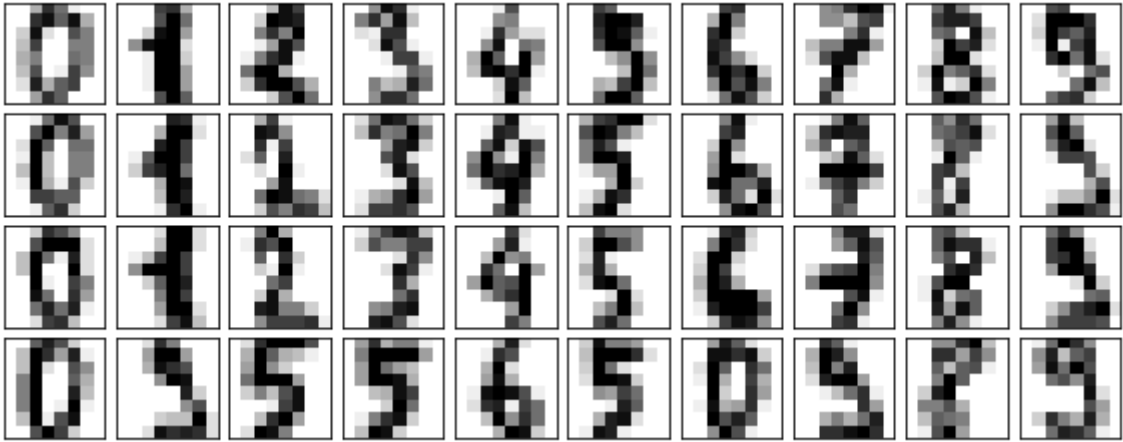
Out[2]: (1797, 64)

## PCA using SVD

In [3]:
```python
def pca(X):
    U, S, PT = svd(X, full_matrices=False)
    Sigma = np.diag(S)
    T = np.dot(U, Sigma)
    return T, Sigma, PT.T # Score, Variance, Loadings
```

In [4]:
```python
def plot_digits(data):
    fig, axes = plt.subplots(
        4, 10, figsize=(10, 4),
        subplot_kw={'xticks':[], 'yticks':[]},
        gridspec_kw=dict(hspace=0.1, wspace=0.1)
    )
    for i, ax in enumerate(axes.flat):
        ax.imshow(
            data[i].reshape(8, 8),
            cmap='binary', interpolation='nearest',
            clim=(0, 16)
        )
```

In [5]:
```python
# Find out the original dimension of the data
X = digits.data
y = digits.target
print("Shape of X", X.shape)
print("Shape of y", y.shape)
```

Shape of X (1797, 64)
Shape of y (1797,)

In [6]:
```python
#Visualize the original data
plot_digits(X)
```

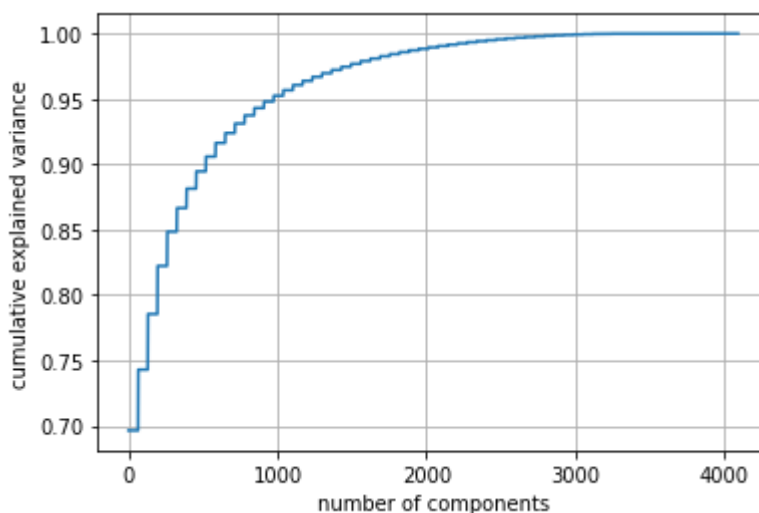## Task 1: Dimensionality reduction: Conduct PCA on the the matrix $X$ to find out the dimension required to capture 80% of the variance

In [7]:
```python
import seaborn as sns

T, S, P = pca(X)
#plt.plot(sorted(S, reverse=True));
```

In [8]:
```python
#SS, _ = np.linalg.eig(S)
SS = S
explained_variance = (SS ** 2) / 4
explained_variance_ratio = (explained_variance / explained_variance.sum())
plt.plot(np.cumsum(explained_variance_ratio))
plt.grid()
plt.xlabel('number of components')
plt.ylabel('cumulative explained variance');
```



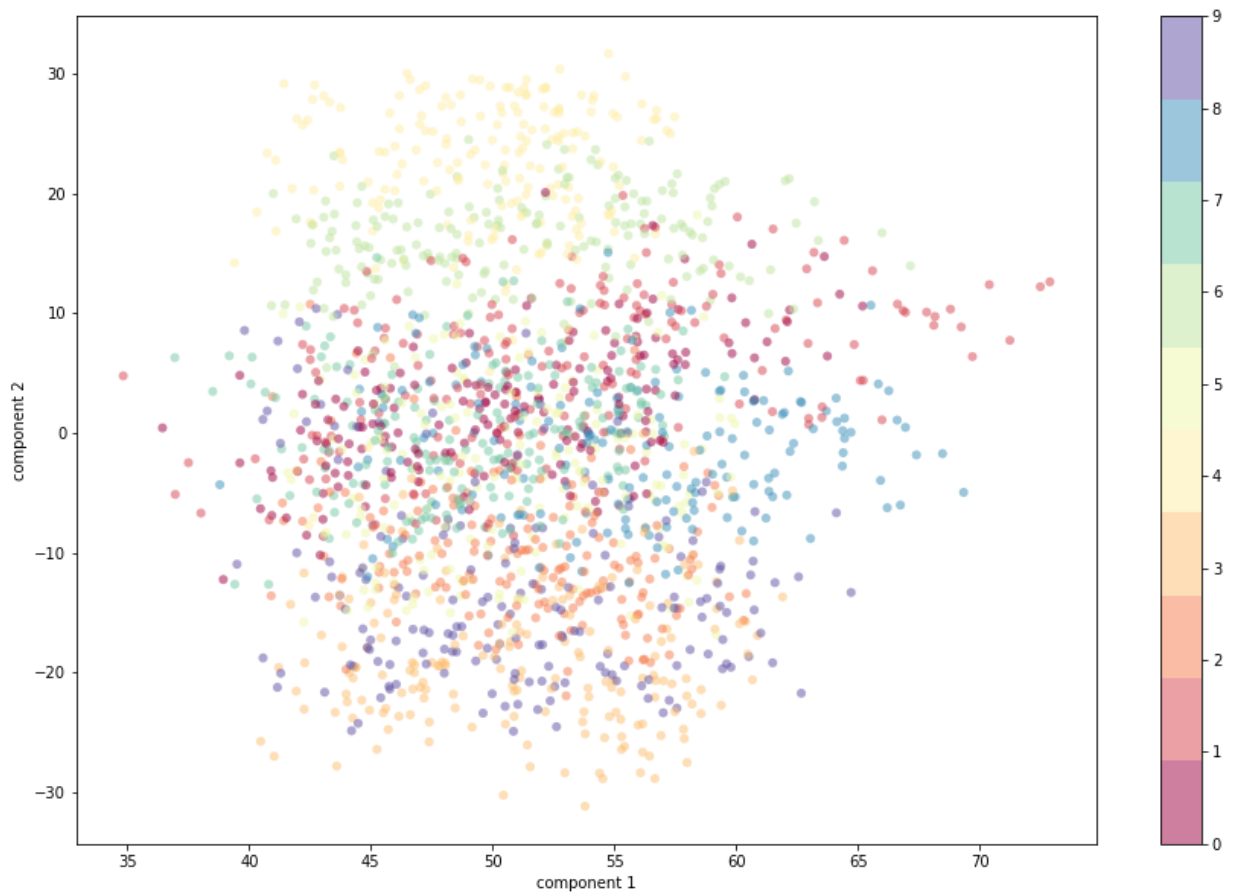## Task 2: Clustering: Project the original data matrix X on the first two PCs and draw the scalar plot

In [9]:
```python
t1 = T[:, 0]
t2 = T[:, 1]

plt.figure(figsize=(15,10))
plt.scatter(t1, t2,
            c=digits.target, edgecolor='none', alpha=0.5,
```

```
                    cmap=plt.cm.get_cmap('Spectral', 10))
plt.xlabel('component 1')
plt.ylabel('component 2')
plt.colorbar();
```
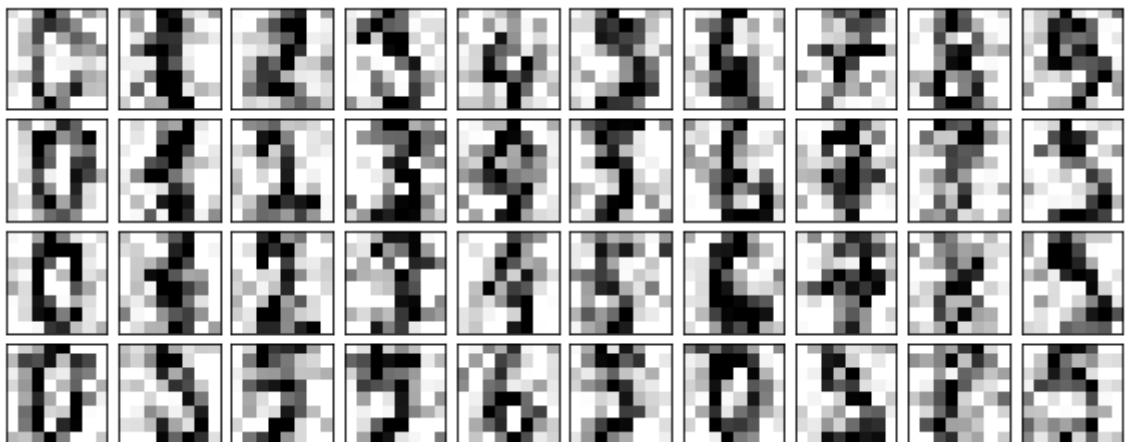


## Task 3: Denoising: Remove noise from the noisy data

In [10]:
```
# Adding noise to the original data
X = digits.data
y = digits.target
np.random.seed(42)
noisy = np.random.normal(X, 4)
plot_digits(noisy)
```
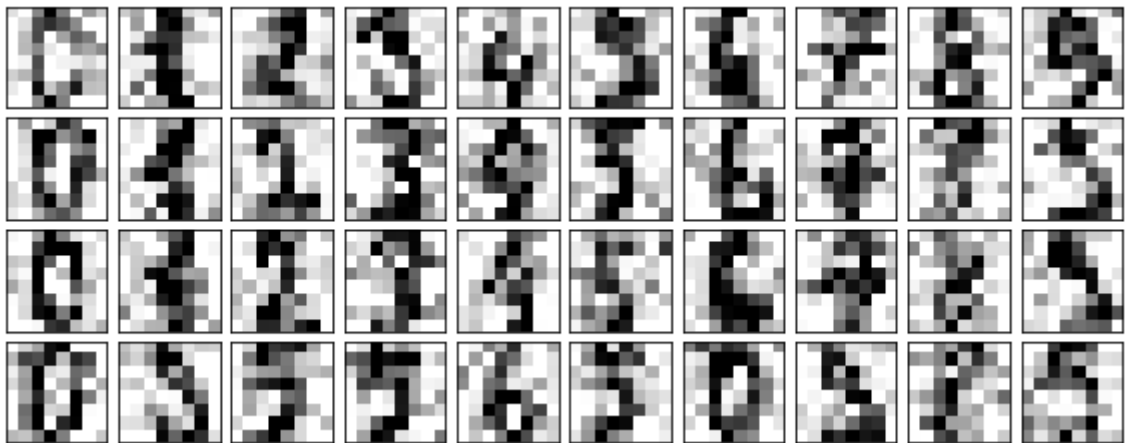


Tips:

- Decompose the noisy data using PCA

- Reconstruct the data using just a few dominant components.For eg. check the variance plot

Since the nature of the noise is more or less similar across all the digits, they are not the fearues with enough variance to discriminate between the digits.

In [11]:
```python
T, S, P = pca(noisy)

n_comps = None
X_denoised = np.dot(noisy, P)
plot_digits(X_denoised @ P.T)
```



## Task 4: Study the impact of normalization of the dataset before conducting PCA. Discuss if it is critical to normalize this particular data compared to the dataset in other notebooks
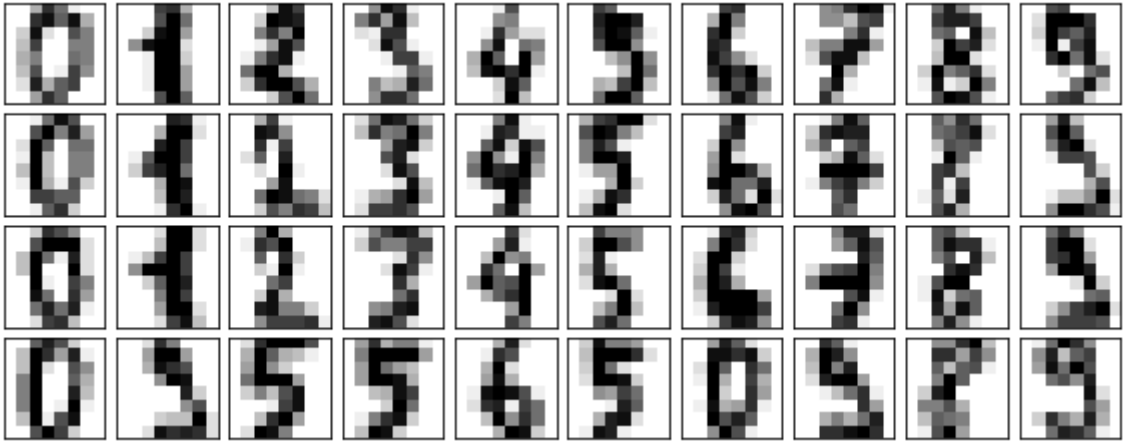
# All the above excercise can be done using the SKLEAR library as follows

In [12]:
```python
from sklearn.decomposition import PCA

X = digits.data
y = digits.target
```
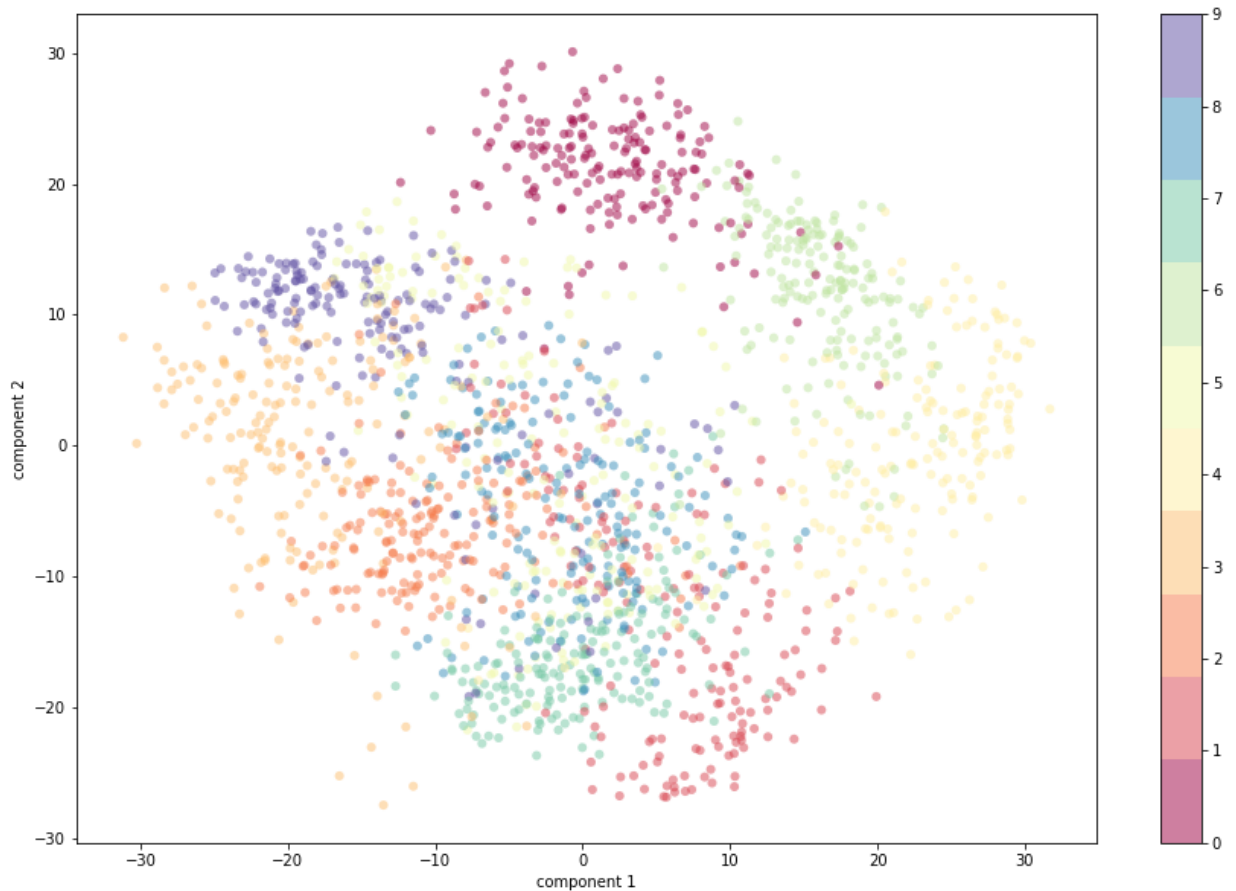
In [13]:
```python
pca = PCA(n_components=2)   # project from 64 to 2 dimensions
projected = pca.fit_transform(digits.data)
print(digits.data.shape)
print(projected.shape)
plot_digits(digits.data)
```
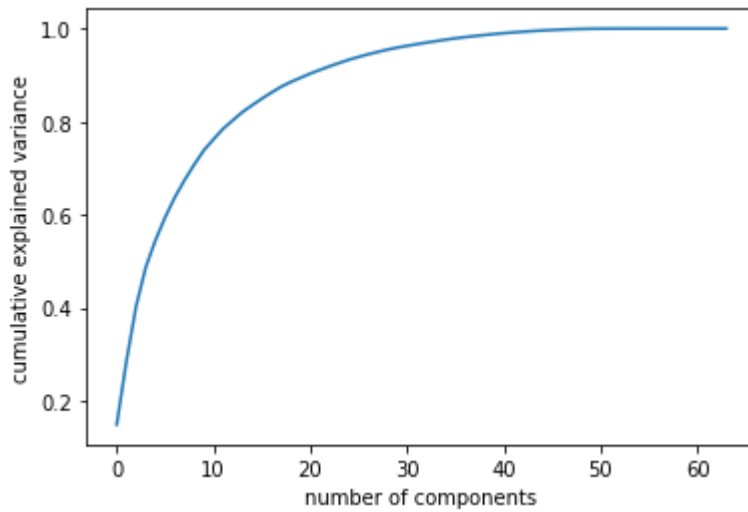
```
(1797, 64)
(1797, 2)
```

In [14]:
```python
plt.figure(figsize=(15,10))
plt.scatter(projected[:, 0], projected[:, 1],
            c=digits.target, edgecolor='none', alpha=0.5,
            cmap=plt.cm.get_cmap('Spectral', 10))
plt.xlabel('component 1')
plt.ylabel('component 2')
plt.colorbar();
```
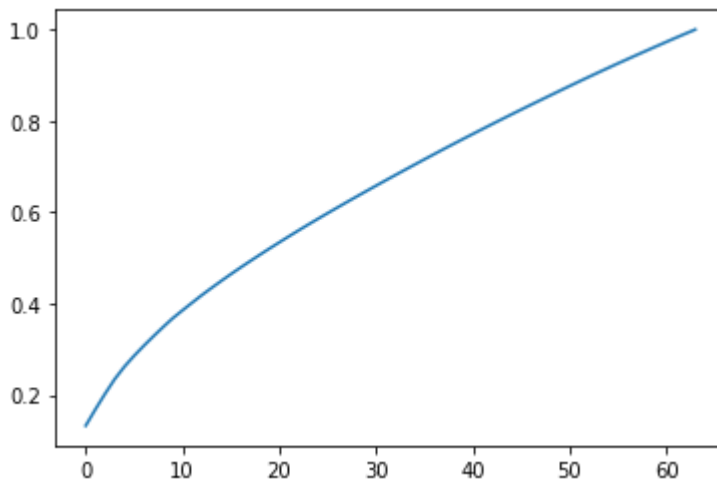


In [15]:
```python
pca = PCA().fit(digits.data)
plt.plot(np.cumsum(pca.explained_variance_ratio_))
plt.xlabel('number of components')
plt.ylabel('cumulative explained variance');
```
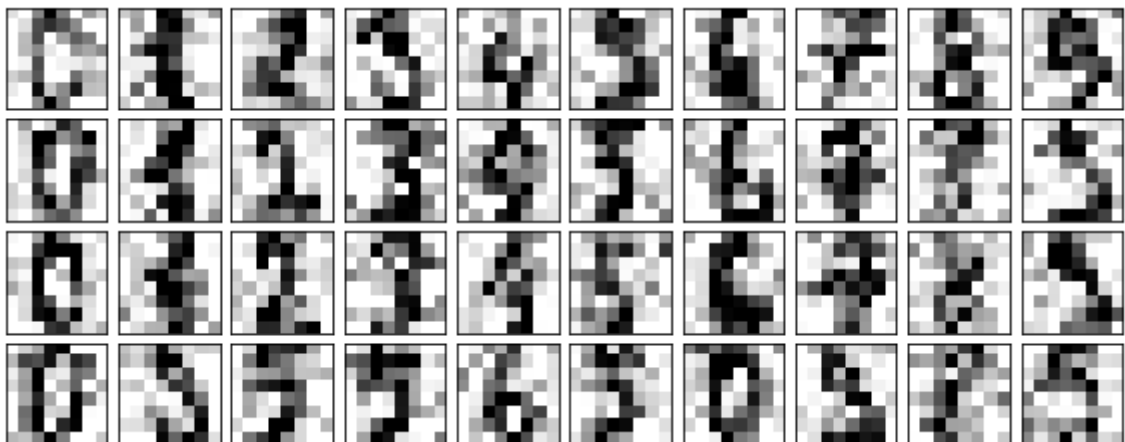
```
In [16]:   S_diag = np.diag(S)#**2 / 4
           SS = S_diag / np.sum(S_diag)
           plt.plot(np.cumsum(SS))
```

Out[16]: [<matplotlib.lines.Line2D at 0x7fc119270ac0>]
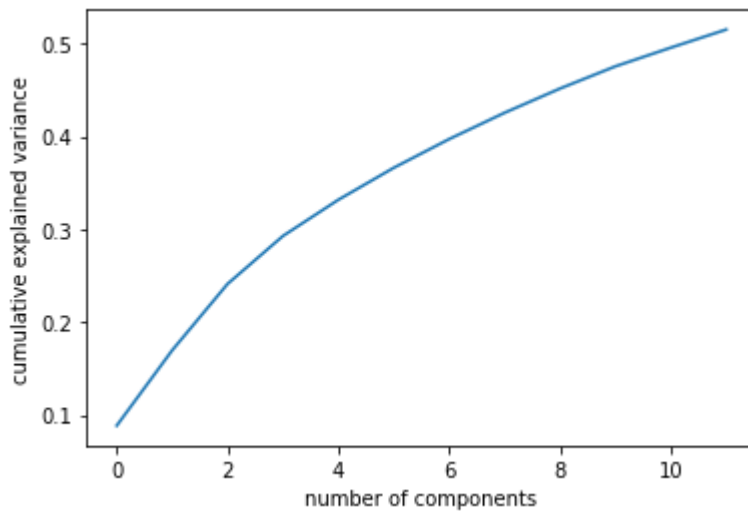


```
In [17]:   np.random.seed(42)
           noisy = np.random.normal(digits.data, 4)
           plot_digits(noisy)
```
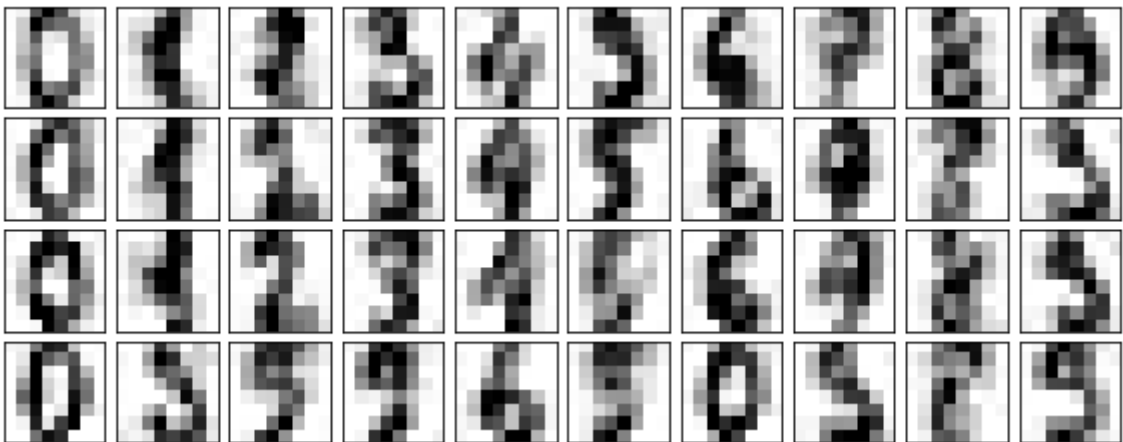


```
In [18]:   pca = PCA(0.50).fit(noisy) # 50% of the variance amounts to 12 principal comp
           pca.n_components_
           plt.plot(np.cumsum(pca.explained_variance_ratio_))
```

```python
plt.xlabel('number of components')
plt.ylabel('cumulative explained variance');
```
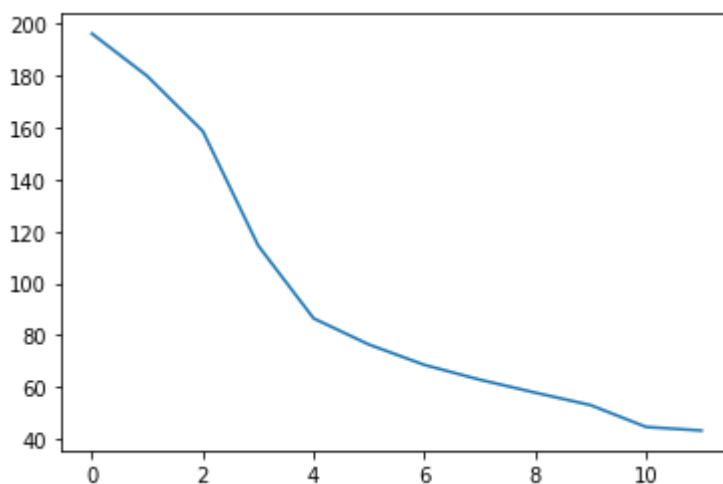


In [19]:
```python
components = pca.transform(noisy)
filtered = pca.inverse_transform(components)
plot_digits(filtered)
```
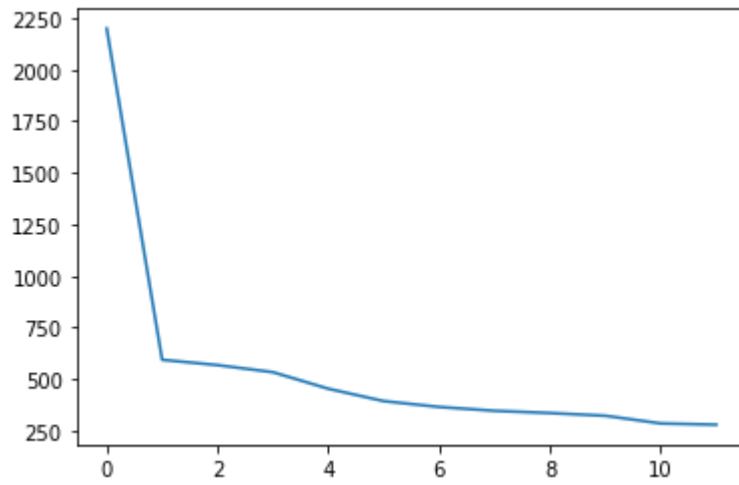


In [20]:
```python
plt.plot(pca.explained_variance_)
```

Out[20]: [<matplotlib.lines.Line2D at 0x7fc117b2bac0>]



In [21]:
```python
plt.plot(np.diag(S)[:12])
```

Out[21]: [<matplotlib.lines.Line2D at 0x7fc118a17be0>]



In [ ]:

In [ ]: