

# Dimension Reduction

Vahid Partovi Nia

Lecture 06



Projection

Principal  
Components

Principal  
Component  
Regression

Partial Least  
Squares

Zip Code

- 1 Projection
- 2 Principal Components
- 3 Principal Component Regression
- 4 Partial Least Squares
- 5 Zip Code

Projection

Principal  
Components

Principal  
Component  
Regression

Partial Least  
Squares

Zip Code

Suppose

$$\mathbf{x} = (x_1, \dots, x_p)^\top.$$

Dimensions reduction is about projection of  $\mathbf{x}$  onto a new space

$$\mathbf{z} = (z_1, \dots, z_m)^\top$$

with a new dimension  $m < p$ .

Projection

Principal  
Components

Principal  
Component  
Regression

Partial Least  
Squares

Zip Code

If projection is linear

$$z_1 = \sum_{j=1}^p \phi_{j1} x_j$$

$$z_2 = \sum_{j=1}^p \phi_{j2} x_j$$

$$\vdots$$

$$z_m = \sum_{j=1}^p \phi_{jm} x_j$$

How to find good coefficients  $\phi_{jm}$

# What linear projection mean?

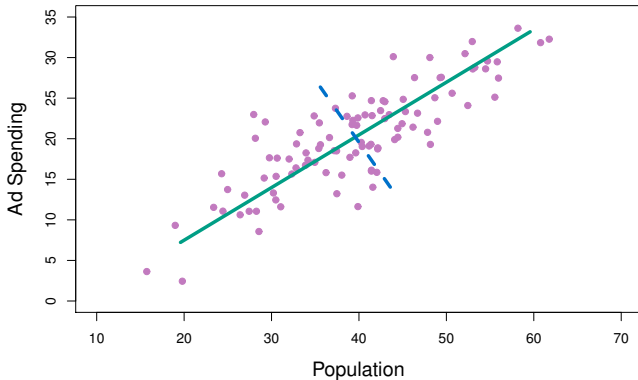
Projection

Principal  
Components

Principal  
Component  
Regression

Partial Least  
Squares

Zip Code



# Data projection

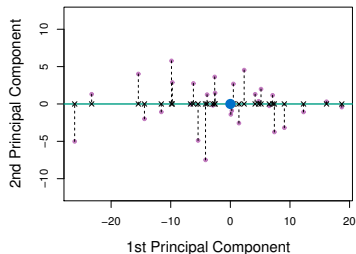
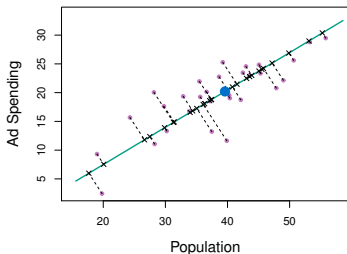
Projection

Principal Components

Principal Component Regression

Partial Least Squares

Zip Code



# How to find projection coefficients

Projection

Principal  
Components

Principal  
Component  
Regression

Partial Least  
Squares

Zip Code

Remember for a univariate random variable  $x$

$$E(\phi x) = \phi E(x)$$

$$V(\phi x) = \phi^2 V(x)$$

For a multivariate  $\mathbf{x}$ ,  $E(\mathbf{x}) = \boldsymbol{\mu}$ ,  $V(\mathbf{x}) = \boldsymbol{\Sigma}$

$$E(\boldsymbol{\phi}^\top \mathbf{x}) = \boldsymbol{\phi}^\top E(\mathbf{x}) = \boldsymbol{\phi}^\top \boldsymbol{\mu}$$

$$V(\boldsymbol{\phi}^\top \mathbf{x}) = \boldsymbol{\phi}^\top V(\mathbf{x}) \boldsymbol{\phi} = \boldsymbol{\phi}^\top \boldsymbol{\Sigma} \boldsymbol{\phi}$$

# Principal components

Projection

Principal  
Components

Principal  
Component  
Regression

Partial Least  
Squares

Zip Code

- Find  $\phi_1$  so that  $V(\mathbf{z}_1) = V(\phi_1^\top \mathbf{x})$  is maximized.
- Find  $\phi_2$  so that  $\phi_2 \perp \phi_1$  and  $V(\mathbf{z}_2) = V(\phi_2^\top \mathbf{x})$  is maximized.
- Find  $\phi_3$  so that  $\phi_3 \perp \phi_1, \phi_2$  and  $V(\mathbf{z}_3) = V(\phi_3^\top \mathbf{x})$  is maximized.
- $\vdots$
- Find  $\phi_m$  so that  $\phi_m \perp \phi_1, \dots, \phi_{m-1}$  and  $V(\mathbf{z}_m) = V(\phi_m^\top \mathbf{x})$  is maximized.



# Find the first projection

Projection

Principal  
Components

Principal  
Component  
Regression

Partial Least  
Squares

Zip Code

Suppose  $V(\mathbf{x}) = \Sigma$  has an eigen value decomposition  
 $\Sigma = \mathbf{P}\mathbf{\Lambda}\mathbf{P}^\top$

- $\mathbf{\Lambda} = \text{diag}(\lambda_i)$
- $\mathbf{P}^\top\mathbf{P} = \mathbf{P}\mathbf{P}^\top = \mathbf{I}$

It is easy to show that  $\lambda_{\max} = \max \frac{V(\phi^\top \mathbf{x})}{\phi^\top \phi}$  and the maximizer is  $\hat{\phi} = \mathbf{e}_{\max}$

Projection

Principal  
ComponentsPrincipal  
Component  
RegressionPartial Least  
Squares

Zip Code

```
import pandas as pd
path='data/'
filename = path+'Auto.csv'
auto = pd.read_csv(filename,
                    na_values=['?'], na_filter=True)
auto = auto.dropna()
```

Projection

Principal  
ComponentsPrincipal  
Component  
RegressionPartial Least  
Squares

Zip Code

```
import pandas as pd
path='data/'
filename = path+'Auto.csv'
auto = pd.read_csv(filename,
                    na_values=['?'], na_filter=True)
auto = auto.dropna()

X = auto[['cylinders', 'displacement',
          'horsepower', 'weight',
          'acceleration']]

y = auto['mpg']
```

## Projection

Principal  
ComponentsPrincipal  
Component  
RegressionPartial Least  
Squares

## Zip Code

```
import pandas as pd
path='data/'
filename = path+'Auto.csv'
auto = pd.read_csv(filename,
                    na_values=['?'], na_filter=True)
auto = auto.dropna()
```

```
X = auto[['cylinders', 'displacement',
          'horsepower', 'weight',
          'acceleration']]
```

```
y = auto['mpg']
```

```
from sklearn.decomposition import PCA
pca = PCA(n_components=2)
pca.fit(X.values)
Z = pca.transform(X)
```

## Projection

Principal  
ComponentsPrincipal  
Component  
RegressionPartial Least  
Squares

## Zip Code

```
import pandas as pd
path='data/'
filename = path+'Auto.csv'
auto = pd.read_csv(filename,
                    na_values=['?'], na_filter=True)
auto = auto.dropna()
```

```
X = auto[['cylinders', 'displacement',
          'horsepower', 'weight',
          'acceleration']]
```

```
y = auto['mpg']
```

```
from sklearn.decomposition import PCA
pca = PCA(n_components=2)
pca.fit(X.values)
Z = pca.transform(X)
```

```
import matplotlib.pyplot as plt
%matplotlib inline
plt.scatter(Z[:,0], Z[:,1]);
```

Projection

Principal  
Components

Principal  
Component  
Regression

Partial Least  
Squares

Zip Code

```
from sklearn.preprocessing import scale
X_std = scale(X.values)
pca_std = PCA(n_components=2)
pca_std.fit(X_std)
Z_std = pca_std.transform(X_std)
plt.scatter(Z_std[:,0], Z_std[:,1]);
```

# Principal component regression

Projection

Principal  
Components

Principal  
Component  
Regression

Partial Least  
Squares

Zip Code

After finding the coefficients, projecting  $x$  to the new dimension provides  $\mathbf{Z}_{n \times m}$ .

Now one can use the projected dimensions to predict a response variable  $y$ .

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

$$\mathbf{Z}_{n \times m} = \mathbf{X}_{n \times p} \boldsymbol{\Phi}_{p \times m}$$

$$\mathbf{y} = \mathbf{Z}\boldsymbol{\theta} + \boldsymbol{\varepsilon}$$

This is called principal component regression,  $\mathbf{Z}$  involves no collinearity!

Projection

Principal  
Components

Principal  
Component  
Regression

Partial Least  
Squares

Zip Code

```
from sklearn.linear_model import LinearRegression
lr = LinearRegression()
X_simple = X[['horsepower ']]
lr.fit(X_simple, y)
lr.score(X_simple, y)
```



Projection

Principal  
Components

Principal  
Component  
Regression

Partial Least  
Squares

Zip Code

```
from sklearn.linear_model import LinearRegression
lr = LinearRegression()
X_simple = X[['horsepower']]
lr.fit(X_simple, y)
lr.score(X_simple, y)
```

```
pcr = LinearRegression()
Z_simple = Z[:,0].reshape(-1,1)
pcr.fit(Z_simple, y)
pcr.score(Z_simple, y)
```

Projection

Principal  
Components

Principal  
Component  
Regression

Partial Least  
Squares

Zip Code

If regression with  $y$  is the reason of projection, it makes sense to project  $x$  to orthogonal axes, while correlation to  $y$  is accounted for.

- PCA:  $\max V(\phi^\top \mathbf{x})$
- PLS:  $\max V(\phi^\top \mathbf{x}) \text{cor}^2(\phi \mathbf{x}, y)$

Projection

Principal  
Components

Principal  
Component  
Regression

Partial Least  
Squares

Zip Code

```
from sklearn.cross_decomposition import PLSRegression
pls = PLSRegression(n_components=2)
pls.fit(X, y)
W = pls.transform(X)
plt.scatter(W[:,0], W[:,1]);
```

Projection

Principal  
Components

Principal  
Component  
Regression

Partial Least  
Squares

Zip Code

```
from sklearn.cross_decomposition import PLSRegression
pls = PLSRegression(n_components=2)
pls.fit(X, y)
W = pls.transform(X)
plt.scatter(W[:,0], W[:,1]);

pls.score(X, y)
```

Projection

Principal  
Components

Principal  
Component  
Regression

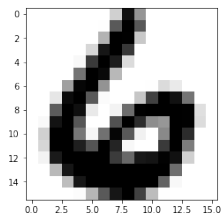
Partial Least  
Squares

Zip Code

```
path='data/'  
filename = path+'ziptrain.csv'  
import numpy as np  
zipdata = np.loadtxt(filename)
```

```
zipdata.shape
```

```
plt.imshow(-zipdata[0, 1:].reshape(16,16), "gray");
```



Projection

Principal  
Components

Principal  
Component  
Regression

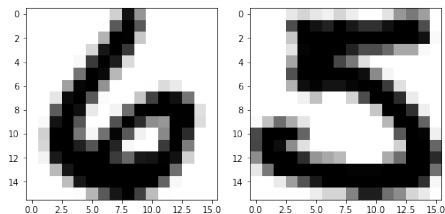
Partial Least  
Squares

Zip Code

```
path='data/'  
filename = path+'ziptrain.csv'  
import numpy as np  
zipdata = np.loadtxt(filename)
```

```
zipdata.shape
```

```
plt.imshow(-zipdata[0, 1:].reshape(16,16), "gray");
```



Projection

Principal  
Components

Principal  
Component  
Regression

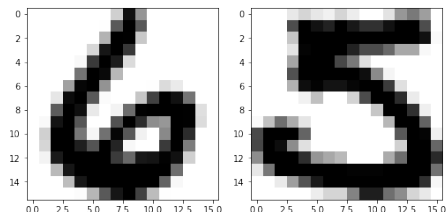
Partial Least  
Squares

Zip Code

```
path='data/'  
filename = path+'ziptrain.csv'  
import numpy as np  
zipdata = np.loadtxt(filename)
```

```
zipdata.shape
```

```
plt.imshow(-zipdata[0, 1:].reshape(16,16), "gray");
```



```
zipdata3=zipdata[zipdata[:, 0] == 3]  
zipdata8=zipdata[zipdata[:, 0] == 8]  
zipdata38 = np.vstack([zipdata3, zipdata8])
```

```
pca = PCA(n_components=2)  
pca.fit(zipdata38[:, 1:])  
Z = pca.transform(zipdata38[:, 1:])  
plt.scatter(Z[:, 0], Z[:, 1], c= zipdata38[:, 0], alpha=0.3);
```

Projection  
Principal  
Components  
Principal  
Component  
Regression  
Partial Least  
Squares  
Zip Code

