

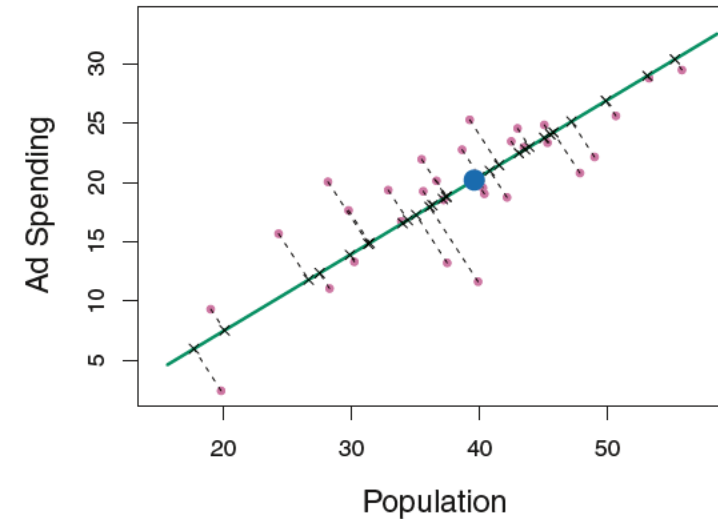
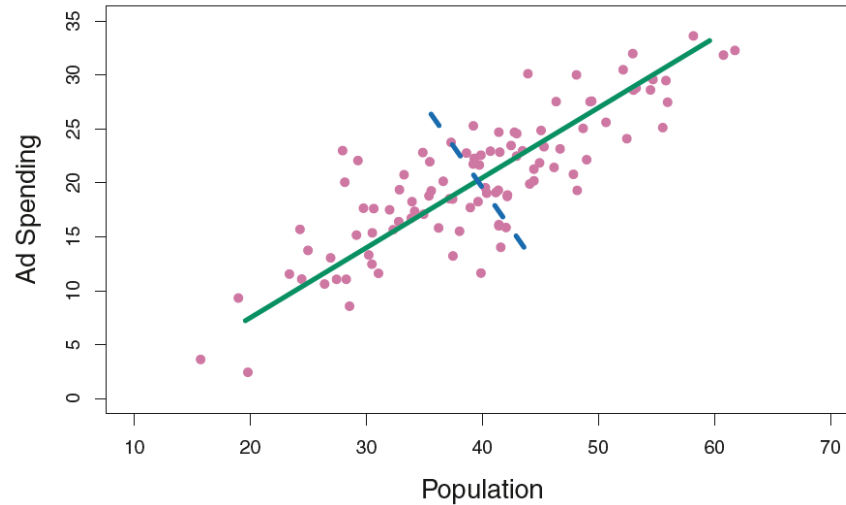
Reduce Dimension

$$\begin{array}{ccccc}
 n \times p & p \times m & n \times m & m \times 1 & n \times 1 \\
 \left(\begin{array}{c} \text{Data} \\ X \end{array} \right) & \left(\begin{array}{c} \Phi \end{array} \right) & = & \left(\begin{array}{c} Z \end{array} \right) & \left(\begin{array}{c} \theta \end{array} \right) = \left(\begin{array}{c} y \end{array} \right)
 \end{array}$$

Linear Regression
 $n \times p$ $p \times 1$ $n \times 1$
 $\left(\begin{array}{c} \text{Data} \\ X \end{array} \right) \quad \left(\begin{array}{c} \beta \end{array} \right) \cong \left(\begin{array}{c} y \end{array} \right)$
 Learn coefficient vector β such that $\|X\beta - y\|^2$ is minimal.

- Compress X into Z , then learn coefficient θ such that $\|Z\theta - y\|^2$ is minimal
- Since $m < p$, we reduced the number of predictors
- Require: Z contains most of the information in X
- *Principle Component Regression*: Z is set to be the first m principal components of X

Principal Component Analysis (PCA)



- Each principal component is a linear combination of the original variables.
- The first principal component direction (green line) is that along which the observations vary the most. The green line is the line *closest* (in squared perpendicular distance) to the data.
- The second principal component is orthogonal to the first (blue dashed line)

Principle Component Regression (PCR)

- In general, there are p distinct principal components for X , the first m will capture most of the information in X
- The best m can be determined by CV
- PCR is not a feature selection method
 - Each of the m principal components is a linear combination of all p of the original features
- When performing PCR, *standardizing* each predictor is recommended.

```
1 cars <- read.csv(file = "cars.csv",
2                   header = T, row.names = 1)
3
4 res.pca <- prcomp(cars[, 1:9], scale = T) # Perform PCA on numeric columns
5 sum.pca <- summary(res.pca)
6 barplot(sum.pca$importance[2,], xlab = "Principal Component", ylab = "Proportion of Variance Explained")
7 plot(sum.pca$importance[3,], type='b', xlab = "Dimension", ylab = "Cum. Prop. of Variance Explained")
8
9 biplot(res.pca)
10
11 library(ggfortify)
12 autoplot(res.pca, data = cars, colour = 'origin')
13 autoplot(res.pca, data = cars, colour = 'origin', label = T)
14 autoplot(res.pca, data = cars, colour = 'origin', loadings = T, loadings.label = T)
15
16
```