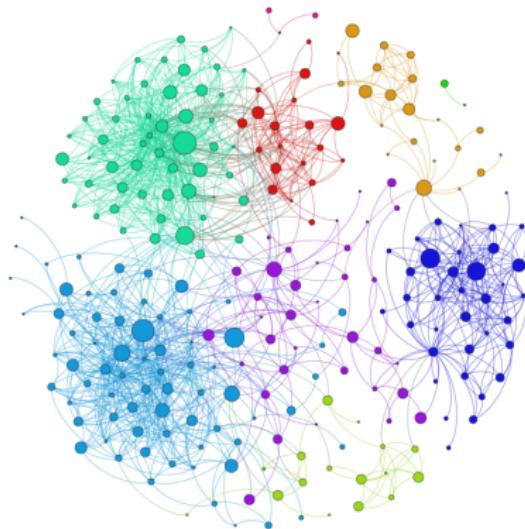


Regression

Jiaming Mao

Xiamen University



Copyright © 2017–2021, by Jiaming Mao

This version: Fall 2021

Contact: jmao@xmu.edu.cn

Course homepage: jiamingmao.github.io/data-analysis



All materials are licensed under the [Creative Commons Attribution-NonCommercial 4.0 International License](#).

Linear Regression

The linear regression model¹ is a discriminative model with $f(x) = \mathbb{E}[y|x]$ ² as the target function and $\mathcal{H} = \{h(x)\}$ consisting of linear functions³:

$$h(x) = x'\beta$$

, where $x = (1, x_1, \dots, x_p)'$ and $\beta = (\beta_0, \beta_1, \dots, \beta_p)'$.

The goal is to find $g \in \mathcal{H}$ that best approximates f .

¹Note on terminology: linear regression can refer broadly to the use of any linear models for regression purposes. Historically, however, it refers more narrowly to least squares linear regression, i.e., linear regression by minimizing in-sample MSE.

²The **conditional expectation function (CEF)**, $\mathbb{E}[y|x]$, is also known as the **regression function**.

³Since each $h(x)$ is associated with a unique β , $h(x)$ is said to be **parametrized** by β . In this case, choosing a hypothesis h is equivalent to choosing a parameter β .

Linear Regression

- Error measures:

$$E_{out}(h) = \mathbb{E} [(y - h(x))^2] \quad (1)$$

$$E_{in}(h) = \frac{1}{N} \sum_{i=1}^N (y_i - h(x_i))^2 \quad (2)$$

- The VC dimension of a linear model is $p + 1^4$. For $N \gg p$, the linear model generalizes well from E_{in} to E_{out} .

⁴ p is the dimension of the input space.

Linear Regression

Let

$$\begin{aligned}\beta^* &= \arg \min_{\beta} \mathbb{E} \left[(y - x'\beta)^2 \right] \\ &= \underbrace{\mathbb{E} [xx']^{-1}_{(p+1) \times (p+1)}}_{(p+1) \times (p+1)} \underbrace{\mathbb{E} [xy]}_{(p+1) \times 1} \quad (3)\end{aligned}$$

- β^* is the **population** regression coefficient.
- $x'\beta^*$ is the best⁵ **linear** predictor of y given x in the underlying population.

⁵in the sense of minimizing the L2 loss function.

Linear Regression

Recall that the CEF $f(x) = \mathbb{E}[y|x]$ is the best⁵ predictor of y given x in the class of all functions of x .

The function $x'\beta^*$ provides the best⁵ linear approximation to the CEF⁶:

$$\beta^* = \arg \min_{\beta} \mathbb{E} \left[(\mathbb{E}[y|x] - x'\beta)^2 \right]$$

⁶Generally,

$$\begin{aligned} \arg \min_h \mathbb{E} [(y - h(x))^2] &= \arg \min_h \mathbb{E} [(y - \mathbb{E}[y|x] + \mathbb{E}[y|x] - h(x))^2] \\ &= \arg \min_h \mathbb{E} [(y - \mathbb{E}[y|x])^2 + (\mathbb{E}[y|x] - h(x))^2 \\ &\quad + 2(y - \mathbb{E}[y|x])(\mathbb{E}[y|x] - h(x))] \\ &= \arg \min_h \mathbb{E} [(\mathbb{E}[y|x] - h(x))^2] \end{aligned}$$

Linear Regression

Let $e^* \equiv y - x'\beta^*$. By construction,

$$\underbrace{\mathbb{E}[xe^*]}_{(p+1) \times 1} = 0 \quad (4)$$

In particular, if x contains a constant term, then (4) $\Rightarrow \mathbb{E}[e^*] = 0$. In this case e^* and x are uncorrelated.

Linear Regression

We can separate the constant term and write the linear model as

$$y = \beta_0 + \tilde{x}'\tilde{\beta} + e$$

, where $\tilde{x} = (x_1, \dots, x_p)'$ and $\tilde{\beta} = (\beta_1, \dots, \beta_p)'$.

Then (3) \Rightarrow

$$\tilde{\beta}^* = \mathbb{V}(\tilde{x})^{-1} \text{Cov}(\tilde{x}, y) \quad (5)$$

$$\beta_0^* = \mathbb{E}[y] - \mathbb{E}[\tilde{x}]' \tilde{\beta}^*$$

Linear Regression

When $p = 1$,

$$y = \beta_0 + \beta_1 x + e$$

(5) \Rightarrow

$$\beta_1^* = \frac{\text{Cov}(x, y)}{\text{V}(x)}$$

$$\beta_0^* = \mathbb{E}[y] - \beta_1^* \mathbb{E}[x]$$

The OLS Estimator

Given observed data $\mathcal{D} = \{(x_1, y_1), \dots, (x_N, y_N)\} \sim^{i.i.d.} p(x, y)$, we have, for $i = 1, \dots, N$,

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + e_i \quad (6)$$

, which can be written as

$$Y = X\beta + e \quad (7)$$

, where $Y = [y_1, \dots, y_N]'$, $e = [e_1, \dots, e_N]'$, and

$$X = \begin{bmatrix} 1 & x_{11} & \cdots & x_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{N1} & \cdots & x_{Np} \end{bmatrix} = \begin{bmatrix} x'_1 \\ \vdots \\ x'_N \end{bmatrix}$$

, where $x_i = [1, x_{i1}, \dots, x_{ip}]'$.

The OLS Estimator

Minimizing the in-sample error (2) \Rightarrow

$$\begin{aligned}\hat{\beta} &= \left[\sum_{i=1}^N x_i x_i' \right]^{-1} \sum_{i=1}^N x_i y_i \\ &= (X' X)^{-1} X' Y\end{aligned}\tag{8}$$

$\hat{\beta}$ is the least squares regression coefficient – the sample estimate of β^* .

Geometric Interpretation

Consider two n -dimensional vectors: $a = (a_1, \dots, a_n)$ and $b = (b_1, \dots, b_n)$. The **Euclidean distance** between a and b is:

$$\|a - b\| = \sqrt{\sum_{i=1}^n (a_i - b_i)^2} = \sqrt{(a - b) \cdot (a - b)}$$

The cosine of the angle between a and b is:

$$\cos \theta = \frac{a \cdot b}{\|a\| \|b\|}$$

, where $\|a\| = \|a - 0\|$ is the length of a .

When $a \cdot b = 0$, a and b are **orthogonal**, denoted by $a \perp b$.

Geometric Interpretation

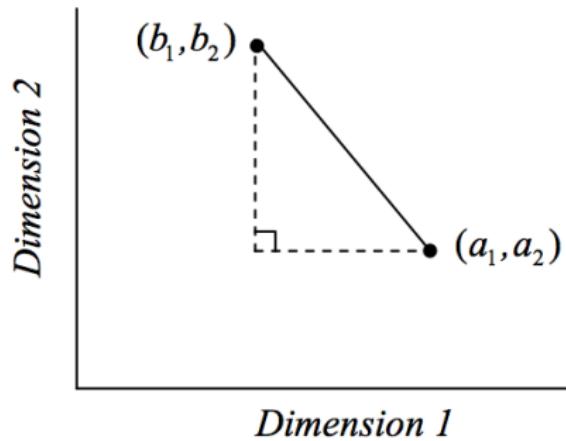
The linear space spanned by a , denoted by $\mathcal{R}(a)$, is the collection of points $\beta a = (\beta a_1, \dots, \beta a_n)$ for any real number β .

The **projection** of b onto $\mathcal{R}(a)$ is the point b^* in $\mathcal{R}(a)$ that is closest to b in terms of Euclidean distance:

$$b^* = \left(\frac{a \cdot b}{\|a\|^2} \right) a$$

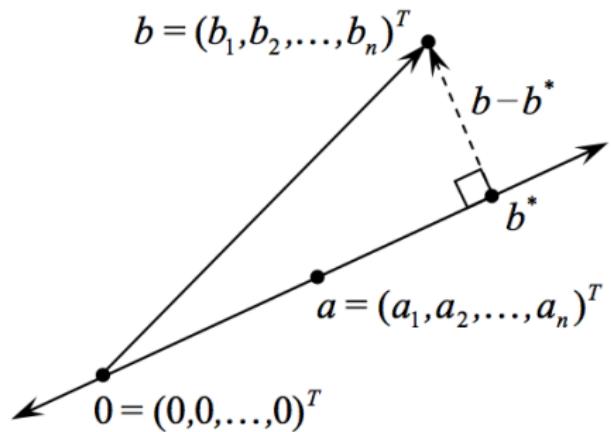
- $(b - b^*) \perp a$

Geometric Interpretation



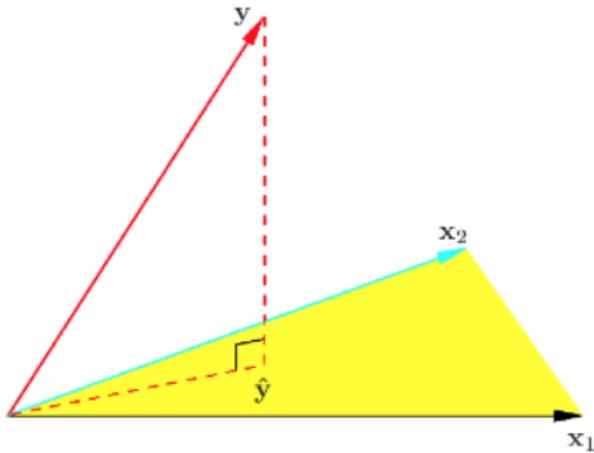
Euclidean Distance in two Dimensions

Geometric Interpretation



Geometric Interpretation

The least squares linear regression fit \hat{Y} is the projection of Y onto the linear space spanned by $\{\mathbf{1}, X_1, \dots, X_p\}$ ⁷.



⁷ $X_j = (x_{1j}, \dots, x_{Nj})'$ for $j = 1, \dots, p$.

Geometric Interpretation

- **Projection matrix** $\mathbb{H} = X(X'X)^{-1}X'$

$$\mathbb{H}Y = \hat{Y}$$

- ▶ \mathbb{H} is also called the **hat matrix**^{8,9}.
- $\hat{e} = Y - X\hat{\beta} = (\mathbb{I} - \mathbb{H})Y \perp \mathcal{R}(\mathbf{1}, X_1, \dots, X_p).$
 - ▶ $\hat{e} \perp X_j \forall j.$
 - ▶ $\hat{e} \perp \mathbf{1} \Rightarrow \sum_i \hat{e}_i = 0.$

⁸Since it “puts a hat” on Y .

⁹The hat matrix has many special properties such as: $\mathbb{H}^2 = \mathbb{H}$, $(\mathbb{I} - \mathbb{H})^2 = (\mathbb{I} - \mathbb{H})$, and $\text{trace}(\mathbb{H}) = 1 + p$.

Frisch-Waugh-Lovell Theorem

Frisch-Waugh-Lovell Theorem

Given linear model (6), the OLS solution is $\hat{\beta} = (\hat{\beta}_1, \dots, \hat{\beta}_p)$, where

$$\hat{\beta}_j = \frac{\sum_{i=1}^N \hat{u}_{ij} \hat{\varepsilon}_{ij}}{\sum_{i=1}^N \hat{u}_{ij}^2} = (\hat{u}'_j \hat{u}_j)^{-1} \hat{u}'_j \hat{\varepsilon}_j \quad (9)$$

, where $\hat{u}_j = (\hat{u}_{1j}, \dots, \hat{u}_{Nj})'$ is the estimated residual from a regression of x_j on x_{-j} ^a, and $\hat{\varepsilon}_j = (\hat{\varepsilon}_{1j}, \dots, \hat{\varepsilon}_{Nj})'$ is the estimated residual from a regression of y on x_{-j} .

Proof

$$^a x_{-j} = (x_1, \dots, x_{j-1}, x_{j+1}, \dots, x_p).$$

Frisch-Waugh-Lovell Theorem

- The Frisch-Waugh-Lovell theorem says that to obtain $\hat{\beta}_j$, we can¹⁰
 - ① regress y on all predictors other than x_j to obtain residual $\hat{\varepsilon}_j$
 - ② regress x_j on all the other predictors to obtain residual \hat{u}_j
 - ③ regress $\hat{\varepsilon}_j$ on \hat{u}_j
- $\hat{\beta}_j$ is the slope coefficient on a scatter plot with $\hat{\varepsilon}_j$ on the y -axis and \hat{u}_j on the x -axis.
- The geometric interpretation is that $(\hat{\varepsilon}_j, \hat{u}_j) \perp x_{-j}$ – they are the residuals in (x_j, y) after orthogonalizing away the influence of the other predictors.

¹⁰In fact, step 1 is not necessary. It is easy to prove an alternative version of the Frisch-Waugh-Lovell theorem:

$$\hat{\beta}_j = \frac{\sum_{i=1}^N \hat{u}_{ij} y_i}{\sum_{i=1}^N \hat{u}_{ij}^2} = (\hat{u}'_j \hat{u}_j)^{-1} \hat{u}'_j y$$

Frisch-Waugh-Lovell Theorem

Generate some data:

$$x_1 \sim U(0, 1)$$

$$x_2 = 0.5x_1 + 0.5r, r \sim U(0, 1)$$

$$y = 1 - 2.5x_1 + 5x_2 + e, e \sim \mathcal{N}(0, 1)$$

```
n <- 500
e <- rnorm(n)
x1 <- runif(n)
x2 <- 0.5*x1 + 0.5*runif(n)
y <- 1 - 2.5*x1 + 5*x2 + e
```

Frisch-Waugh-Lovell Theorem

```
require(AER)
reg <- lm(y ~ x1 + x2)
coeftest(reg)

##
## t test of coefficients:
##
##             Estimate Std. Error   t value Pr(>|t|)
## (Intercept)  1.01013   0.11884   8.4997 2.233e-16 ***
## x1          -2.59166   0.22529 -11.5039 < 2.2e-16 ***
## x2           5.06250   0.31213  16.2193 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Frisch-Waugh-Lovell Theorem

```
# to obtain beta1, we first regress y and x1 on x2
u1 <- residuals(lm(x1~x2))
e1 <- residuals(lm(y~x2))
b1 <- cov(u1,e1)/var(u1)

# to obtain beta2, we first regress y and x2 on x1
u2 <- residuals(lm(x2~x1))
e2 <- residuals(lm(y~x1))
b2 <- cov(u2,e2)/var(u2)

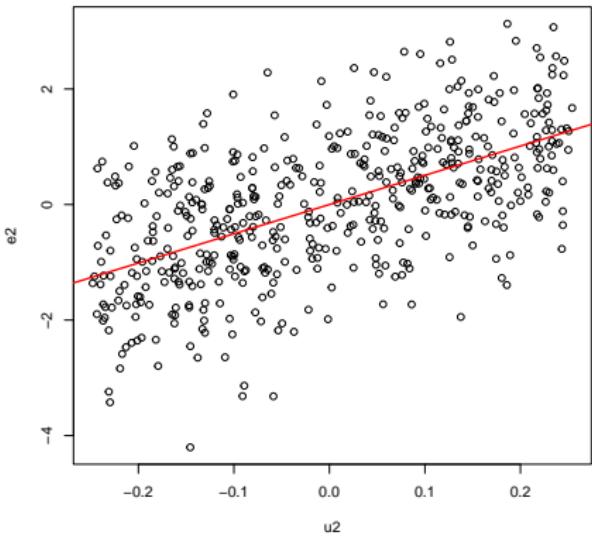
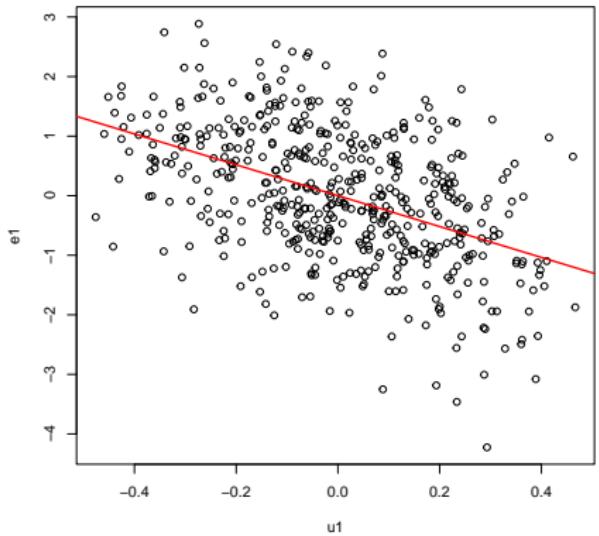
b0 <- mean(y) - b1*mean(x1) - b2*mean(x2)
cbind(b0,b1,b2)

##          b0          b1          b2
## [1,] 1.010133 -2.591657 5.062497
```

Frisch-Waugh-Lovell Theorem

```
# Alternatively,  
b1 <- cov(u1,y)/var(u1)  
b2 <- cov(u2,y)/var(u2)  
b0 <- mean(y) - b1*mean(x1) - b2*mean(x2)  
cbind(b0,b1,b2)  
  
##           b0          b1          b2  
## [1,] 1.010133 -2.591657 5.062497
```

Frisch-Waugh-Lovell Theorem



Numerical Solution

- (8) is the **analytical solution** to the problem:

$$\min_{\beta} \sum_{i=1}^N (y_i - x_i' \beta)^2 \quad (10)$$

- For many problems, however, such analytical solutions do not exist.
How do we solve (10) numerically?

Gradient Descent

- **Gradient descent** is a method that finds a minimum of a function by figuring out in which direction the function's slope is rising the most steeply, and moving in the opposite direction.
- If the function that we are trying to minimize is **differentiable** and **convex**, then it has a unique **global minimum**¹¹. In this case, gradient descent starting from any point is guaranteed to find the minimum.

¹¹ See [Appendix](#) for a comparison of convex and non-convex functions.

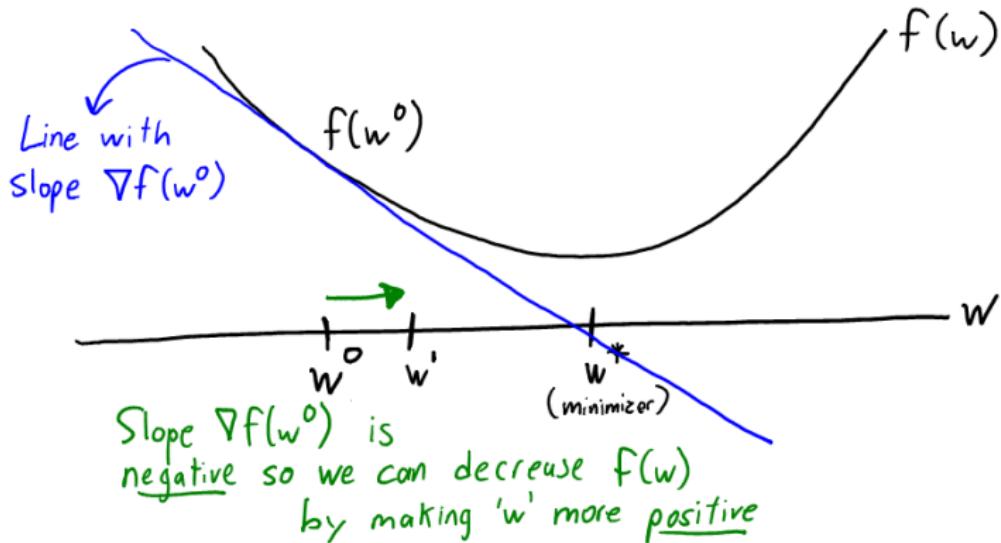
Gradient Descent

Let the function that we want to minimize be $f(\theta)$, so that we are looking for

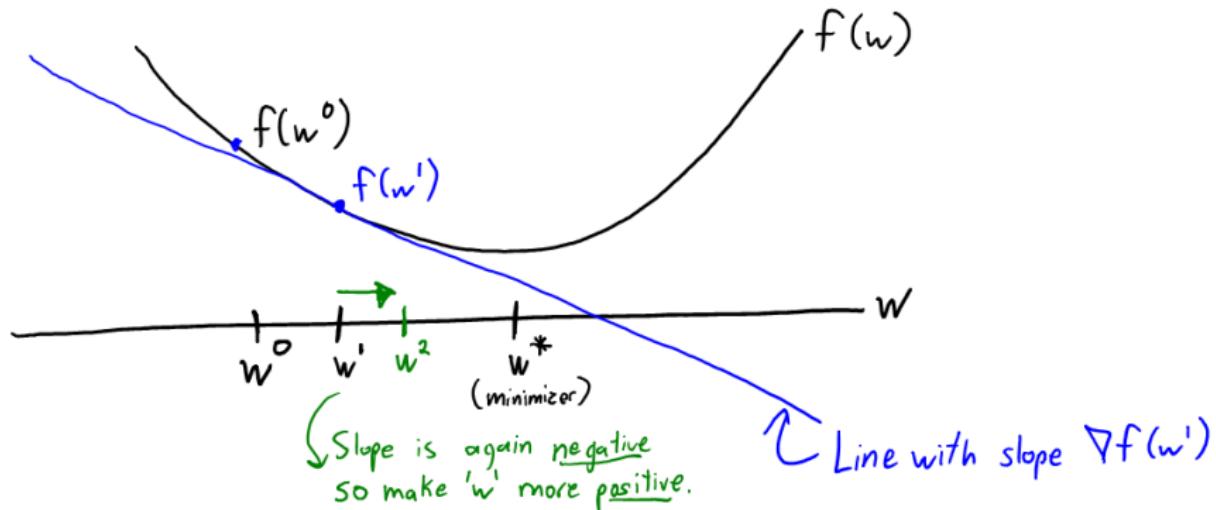
$$\theta^* = \arg \min_{\theta} f(\theta)$$

Gradient descent is based on a simple observation: at any point θ , $\nabla f(\theta)$ is a vector pointing in the direction of the greatest increase in f . Hence, to find θ^* , we need to move in the direction of $-\nabla f(\theta)$.

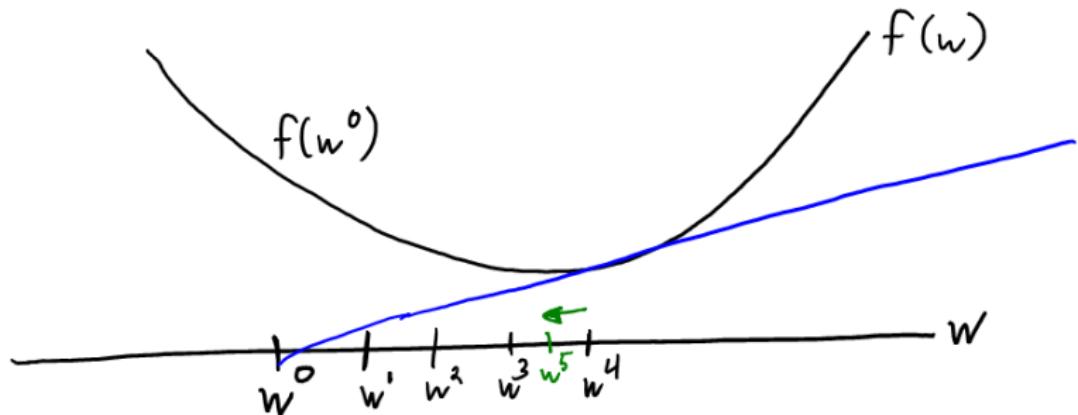
Gradient Descent



Gradient Descent



Gradient Descent



Now the slope $\nabla f(w^4)$ is positive
so we move in the negative direction.

Gradient Descent

Gradient Descent Algorithm

Start with $\theta^{(0)}$.

repeat

$$\theta^{(t)} := \theta^{(t-1)} - \alpha \cdot \nabla f \left(\theta^{(t-1)} \right)$$

until $\left\| \nabla f \left(\theta^{(t)} \right) \right\| \leq \text{tolerance}$

return $\theta^{(t)}$

Gradient Descent

- We can run gradient descent until the gradient is close to zero.
Alternatively, we can run the algorithm for a fixed number of steps.
- α is the step size, also called the **learning rate**. If it's too large, we may easily overshoot the minimum of f . If it's too small, it may take too long to get to the minimum. It is common to start with a higher learning rate and then slowly decrease it as we approach the minimum.

Gradient Descent for Linear Regression

For linear regression, the objective function we want to minimize is:

$$\begin{aligned}f(\beta) &= \sum_{i=1}^N (y_i - x_i' \beta)^2 \\&= (Y - X\beta)' (Y - X\beta)\end{aligned}$$

\Rightarrow

$$\nabla f(\beta) = 2(X\beta - Y)' X$$

Gradient Descent for Linear Regression

```
#####
# Objective Function #
#####
mse <- function(X,y,beta){
  N <- length(y)
  Z <- y - X%*%beta
  mse <- t(Z)%*%Z/N
}

#####
# Gradient Function #
#####
grad <- function(X,y,beta){
  N <- length(y)
  grad <- 2*t(X)%*%(X%*%beta-y)/N
}
```

Gradient Descent for Linear Regression

```
#####
# Gradient Descent Function #
#####
# alpha: learning rate
# niter: number of iterations
gradientDescent <- function(X,y,beta0,alpha,niter){
  beta <- beta0
  mse_hist <- rep(0,niter) # stores history of mse
  beta_hist <- list(niter) # stores history of beta
  for (i in 1:niter){
    beta_hist[[i]] <- beta
    mse_hist[i] <- mse(X,y,beta)
    beta <- beta - alpha*grad(X,y,beta) # update beta
  }
  result <- list("beta"=beta,"mse_hist"=mse_hist,"beta_hist"=beta_hist)
  return(result)
}
```

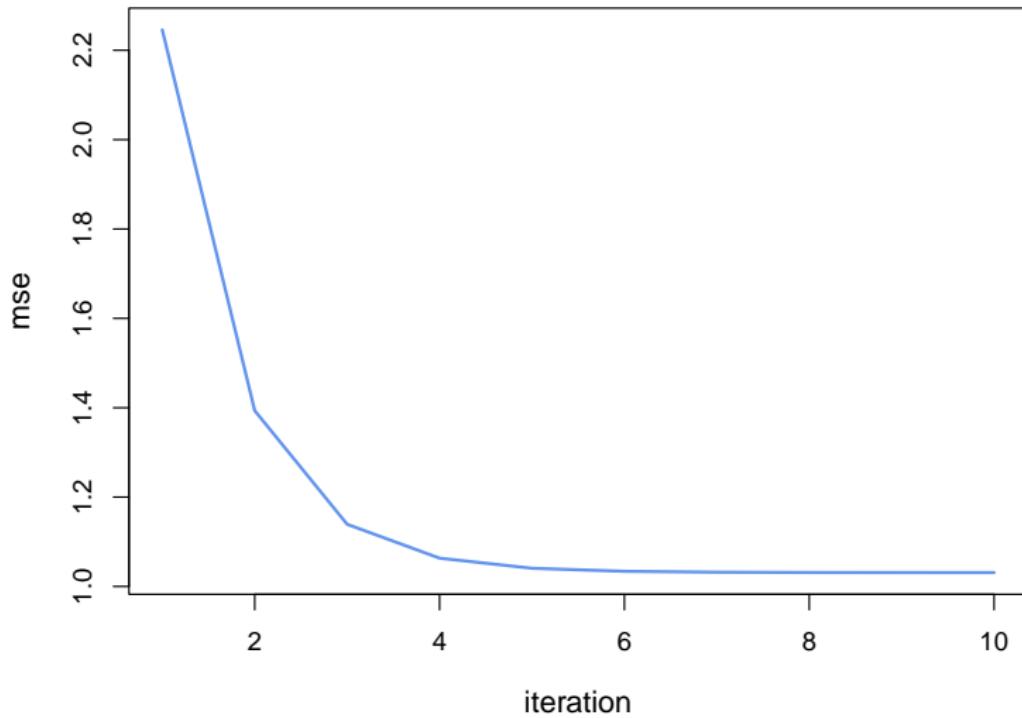
Gradient Descent for Linear Regression

```
#####
# Simulation 1 #
#####
# generate some data
## note: "true mse" = var(e) = 1
n <- 500
e <- rnorm(n)
x <- runif(n)
y <- 2*x + e

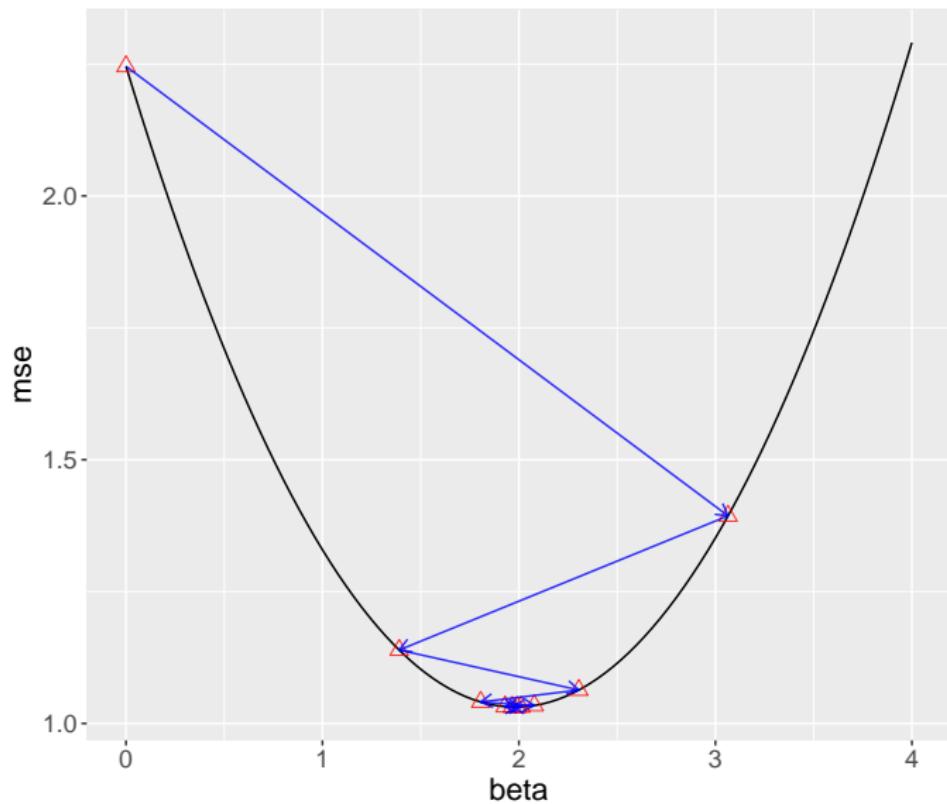
# Gradient Descent
## initial guess: 0; learning rate: 2.5; iteration: 10
X <- cbind(x) # make x column vector
result <- gradientDescent(X,y,0,2.5,10)
result$beta

##          [,1]
## x  1.977162
```

Gradient Descent for Linear Regression



Gradient Descent for Linear Regression



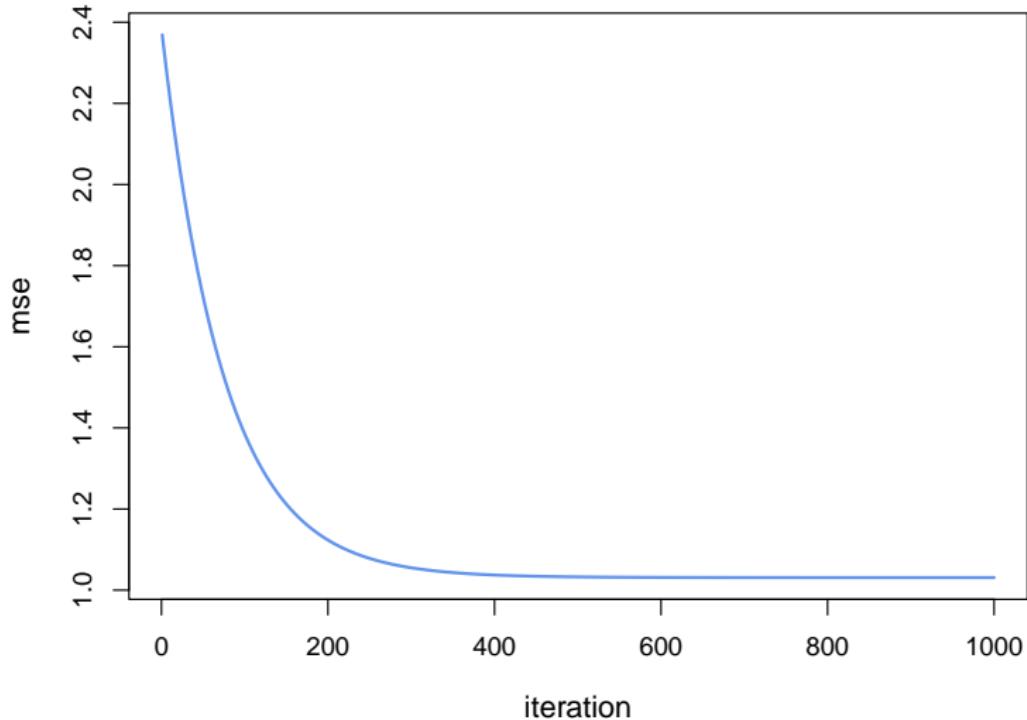
Gradient Descent for Linear Regression

```
#####
# Simulation 2 #
#####
n <- 500
e <- rnorm(n)
x <- runif(n)
y <- -2 + 4*x + e

# Gradient Descent
## learning rate: 0.05; iteration: 1000
X <- cbind(rep(1,n),x) # X matrix
result <- gradientDescent(X,y,c(0,0),0.05,1000)
result$beta

##          [,1]
## -2.005872
## x  3.989831
```

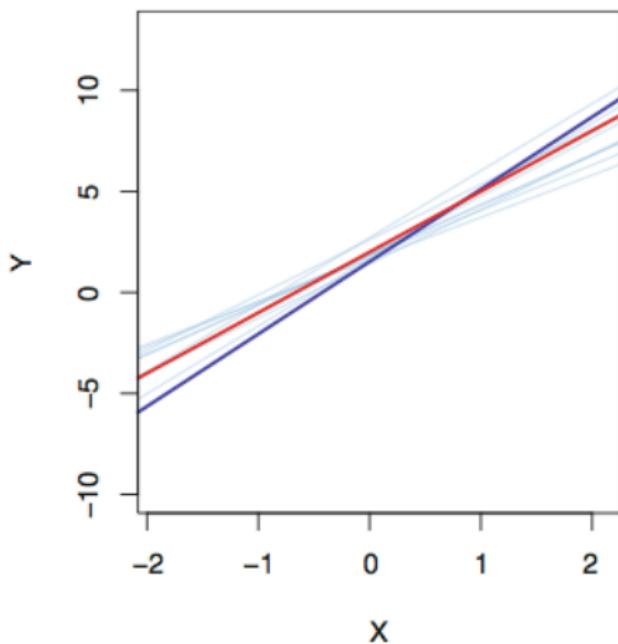
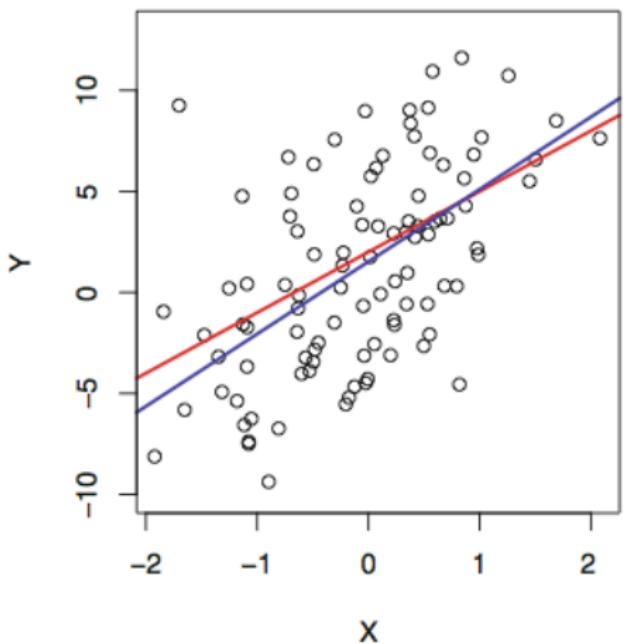
Gradient Descent for Linear Regression



Asymptotic Properties

- $\hat{\beta}$ is unbiased: $\mathbb{E}(\hat{\beta}) = \beta^*$.
- But how much does $\hat{\beta}$ vary around β^* ?

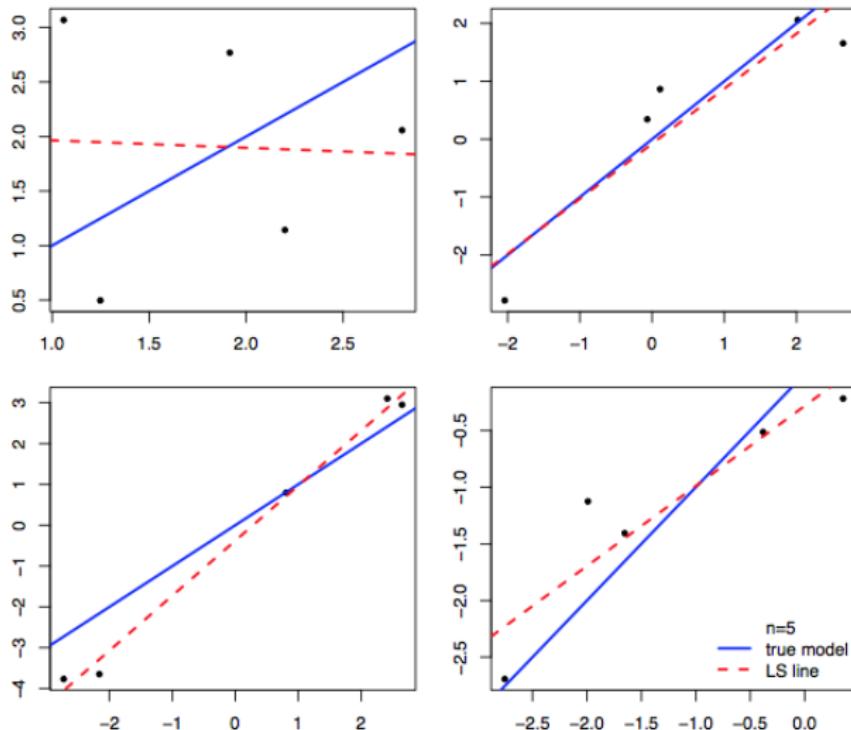
Asymptotic Properties



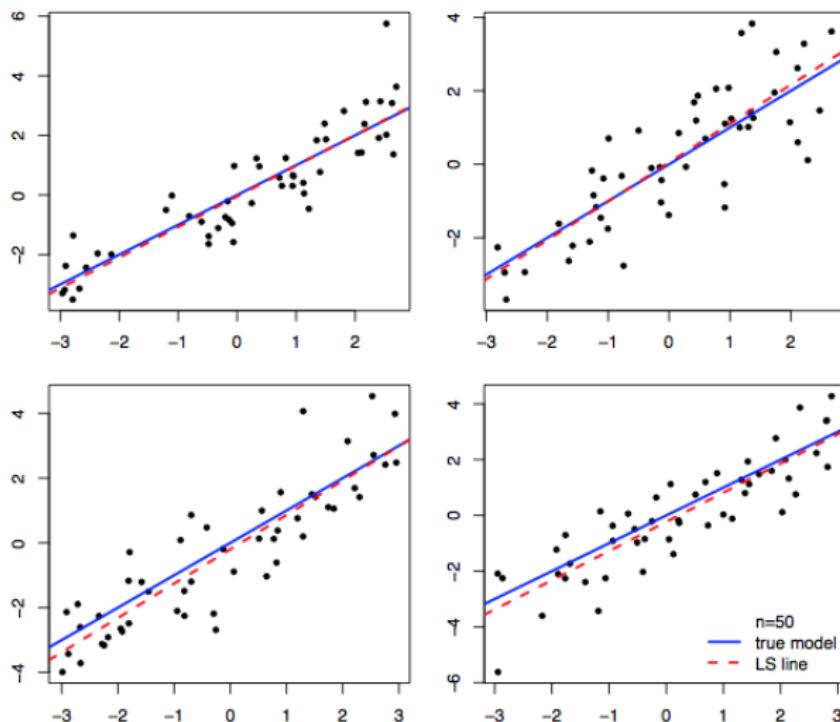
red: $x'\beta^*$. blue: $x'\hat{\beta}$

Right: $x'\hat{\beta}$ based on 10 random set of observations.

Asymptotic Properties



Asymptotic Properties



Asymptotic Properties

By the central limit theorem,

$$\sqrt{N} (\hat{\beta} - \beta^*) \xrightarrow{d} \mathcal{N} \left(0, E(xx')^{-1} E[xx' (e^*)^2] E(xx')^{-1} \right)$$

- $\text{V}(\hat{\beta}) \doteq \underbrace{N^{-1} E(xx')^{-1} E[xx' (e^*)^2] E(xx')^{-1}}_{(p+1) \times (p+1)}$ is the **asymptotic variance** of $\hat{\beta}$ conditional on x .
- $\text{V}(\hat{\beta})$ quantifies the uncertainty of $\hat{\beta}$ due to random sampling.

Asymptotic Properties

$$\begin{aligned}\hat{V}(\hat{\beta}) &= \left[\sum_{i=1}^N x_i x_i' \right]^{-1} \left(\sum_{i=1}^N x_i x_i' \hat{e}_i^2 \right) \left[\sum_{i=1}^N x_i x_i' \right]^{-1} \\ &= (X'X)^{-1} (X'\Omega X) (X'X)^{-1} \xrightarrow{p} V(\hat{\beta})\end{aligned}\tag{11}$$

, where $\Omega = \text{diag}(\hat{e}_1^2, \dots, \hat{e}_N^2) = \begin{bmatrix} \hat{e}_1^2 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \hat{e}_N^2 \end{bmatrix}$.

Asymptotic Properties

Homoskedasticity: $\mathbb{E} \left[(\epsilon^*)^2 \mid x \right] = \sigma^2$

Heteroskedasticity: $\mathbb{E} \left[(\epsilon^*)^2 \mid x \right] = \sigma^2(x)$

Under homoskedasticity,

$$\sqrt{N} \left(\hat{\beta} - \beta^* \right) \xrightarrow{d} \mathcal{N} \left(0, E(xx')^{-1} \sigma^2 \right)$$

$$\hat{V}(\hat{\beta}) = (X'X)^{-1} \hat{\sigma}^2 \quad (12)$$

Asymptotic Properties

From (9), we can also derive the homoskedastic asymptotic variance of $\hat{\beta}_j$ – the $(j + 1)$ -th diagonal element of $V(\hat{\beta})$ – as:

For $j = 1, \dots, p$,

$$\sqrt{N} (\hat{\beta}_j - \beta_j^*) \xrightarrow{d} \mathcal{N} \left(0, \frac{\sigma^2}{V(u_j)} \right)$$
$$\hat{V}(\hat{\beta}_j) = \frac{\hat{\sigma}^2}{\hat{u}'_j \hat{u}_j} \quad (13)$$

Asymptotic Properties

- **t-statistic**

$$t_j = \frac{\hat{\beta}_j - \beta_j^*}{\widehat{\text{se}}(\hat{\beta}_j)} \xrightarrow{d} \mathcal{N}(0, 1)$$

, where $\widehat{\text{se}}(\hat{\beta}_j) = \sqrt{\widehat{\text{V}}(\hat{\beta}_j)}$.

- 95% **confidence interval** for β_j^* :

$$[\hat{\beta}_j - 1.96 \times \widehat{\text{se}}(\hat{\beta}_j), \hat{\beta}_j + 1.96 \times \widehat{\text{se}}(\hat{\beta}_j)]$$

- ▶ The interval represents a **set estimate** of β_j^* .

Hypothesis Testing

$$\mathbb{H}_0 : \beta_j^* = 0 \text{ vs. } \mathbb{H}_1 : \beta_j^* \neq 0$$

Under \mathbb{H}_0 ,

$$t_j = \frac{\hat{\beta}_j}{\widehat{\text{se}}(\hat{\beta}_j)} \xrightarrow{d} \mathcal{N}(0, 1) \quad (14)$$

P-value: probability of observing any value more extreme than $|t_j|$ under \mathbb{H}_0 . (14) \Rightarrow in large sample,

$$p\text{-value} \approx 2(1 - \Phi(|t_j|)) \quad (15)$$

, where Φ is the CDF of $\mathcal{N}(0, 1)$.

Hypothesis Testing

For significance level α , reject \mathbb{H}_0 if $|t_j| > c_\alpha = \Phi^{-1}(1 - \alpha/2)$, or equivalently, if $p\text{-value} < \alpha^{12}$.

- c_α is called the **asymptotic critical value**.
- Common practice: $\alpha = 5\%$ ($c_{.05} \approx 1.96$), $\alpha = 10\%$ ($c_{.10} \approx 1.64$), $\alpha = 1\%$ ($c_{.01} \approx 2.58$).

¹²It is worth emphasizing that (15) is only valid *in large samples*, since it is based on the asymptotic distribution of t_j . Any p -values calculated using (15) on small samples should *not* be trusted. In general, hypothesis tests based on the asymptotic properties of test statistics are only valid for large samples.

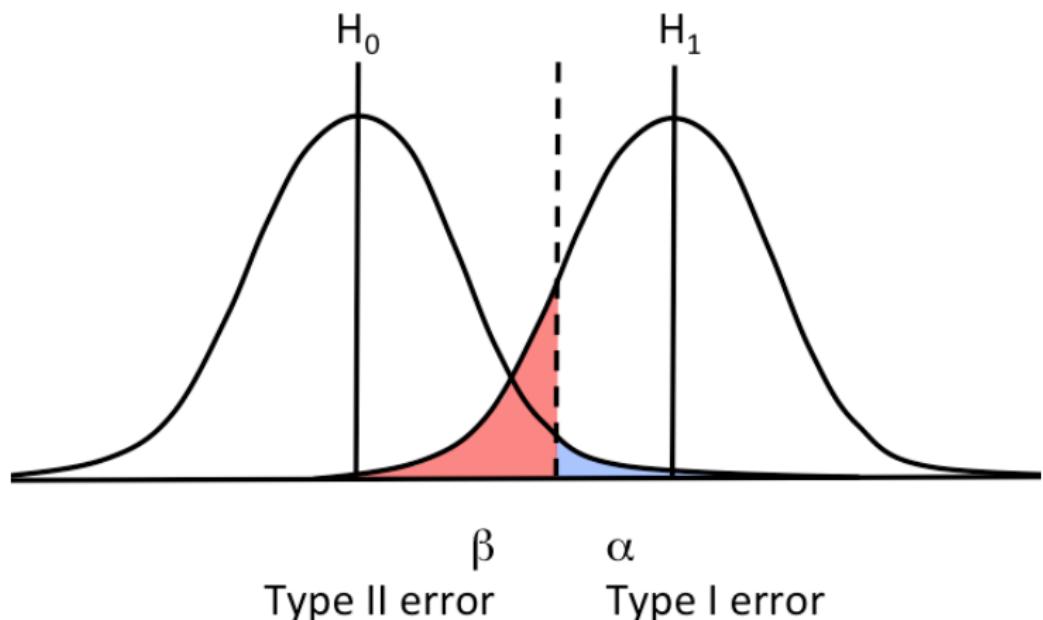
Hypothesis Testing

Hypothesis Testing Decisions

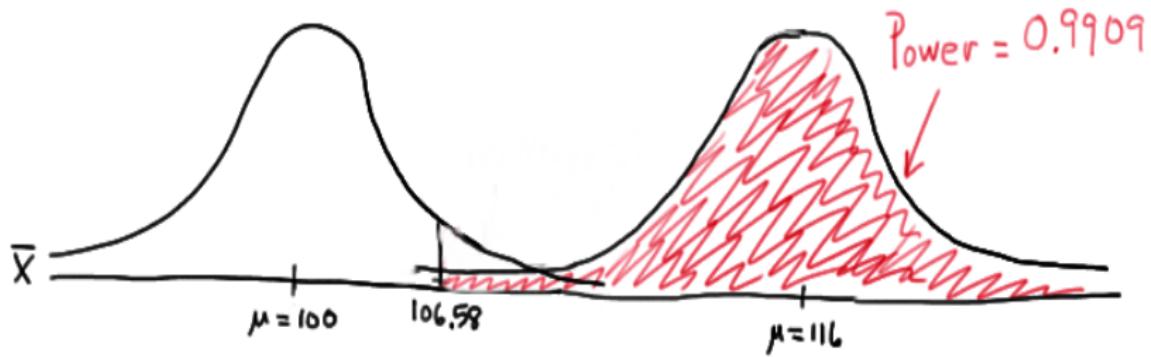
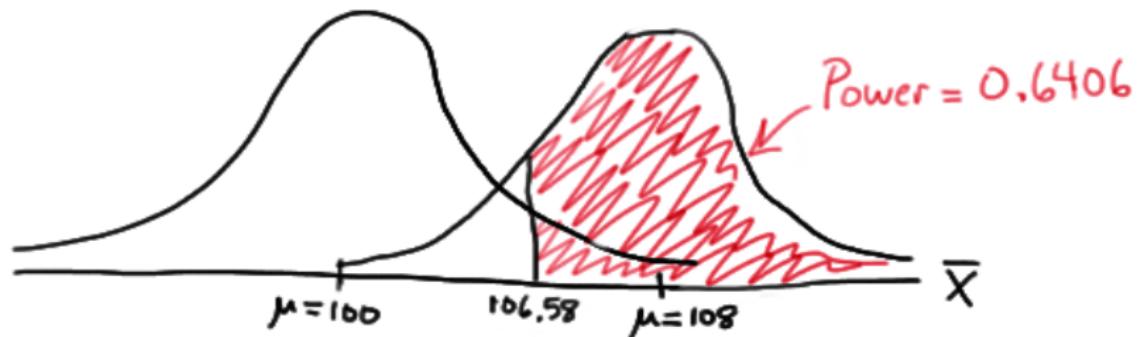
	Accept \mathbb{H}_0	Reject \mathbb{H}_0
\mathbb{H}_0 true	Correct Decision	Type I Error
\mathbb{H}_1 true	Type II Error	Correct Decision

- α is the **size** of the test – the probability of making a Type I error: $\Pr(\text{reject } \mathbb{H}_0 | \mathbb{H}_0 \text{ is true})$.
- The **power** or **sensitivity** of a test, is the probability of rejecting \mathbb{H}_0 when \mathbb{H}_1 is true. Thus $(1 - \text{power})$, denoted by β , is the probability of making a Type II error: $\Pr(\text{fail to reject } \mathbb{H}_0 | \mathbb{H}_1 \text{ is true})$.
 - ▶ Power \uparrow as $\alpha \uparrow$, or sample size $N \uparrow$, or the true (population) parameter value is further away from its hypothesized value under \mathbb{H}_0 .

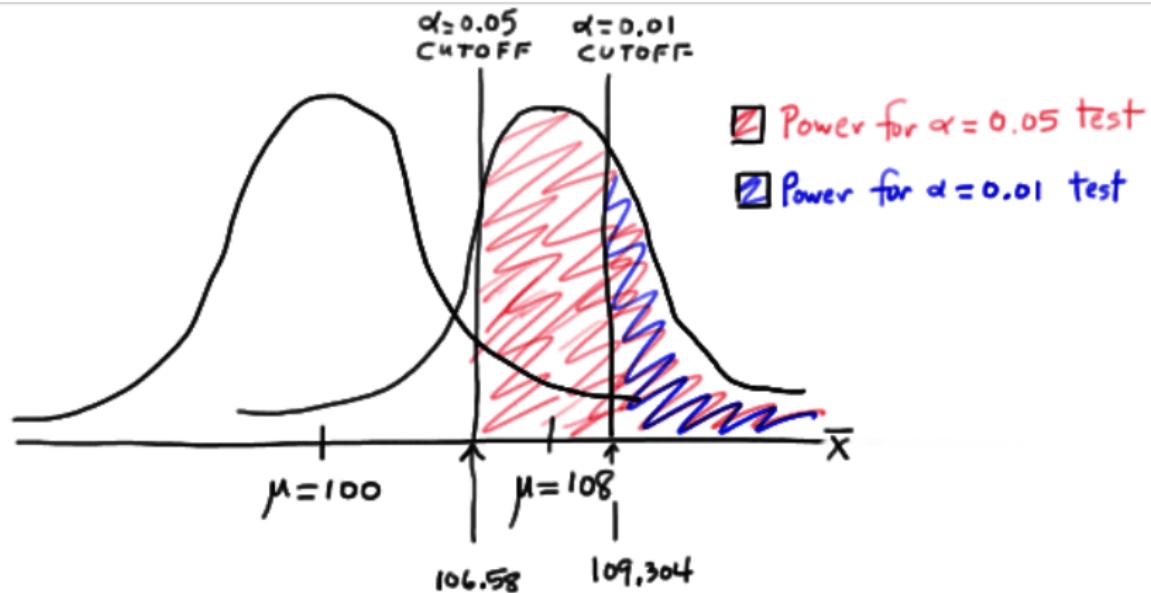
Hypothesis Testing



Hypothesis Testing



Hypothesis Testing



$$R^2 = \frac{\sum_{i=1}^N (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^N (y_i - \bar{y})^2} = 1 - \frac{\sum_{i=1}^N \hat{e}_i^2}{\sum_{i=1}^N (y_i - \bar{y})^2}$$

- measures the amount of variation in y_i accounted for by the model:
1 = perfect, 0 = perfect misfit.
- cannot go down when you add regressors.
 - ▶ Intuition: adding more regressors always allow us to fit the training data more accurately (i.e., reduce E_{in} , but not necessary E_{out})¹³.

¹³Technically, $\hat{\beta}$ is chosen to minimize $\sum_i \hat{e}^2$. if you add a regressor, you can always set the coefficient of that regressor equal to zero to get the same $\sum_i \hat{e}^2$. Therefore R^2 cannot go down.

Robust Standard Errors

(11) is known as **heteroskedasticity-consistent (HC) standard error**, **robust standard error**, or **White standard error**.

Let's generate some data:

$$x = U(0, 100)$$

$$y = 5x + e, \quad e \sim \mathcal{N}(0, \exp(x))$$

```
n <- 1e3
x <- 100*runif(n)
y <- rnorm(n, mean=5*x, sd=exp(x))
```

Robust Standard Errors

```
require(AER)
coeftest(lm(y~x)) # homoskedastic standard error

##
## t test of coefficients:
##
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept) 1.0116e+41 9.0736e+40 1.1148 0.26519
## x          -3.0822e+39 1.5634e+39 -1.9715 0.04895 *
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

coeftest(lm(y~x),vcov=vcovHC) # robust standard error

##
## t test of coefficients:
##
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept) 1.0116e+41 8.6253e+40 1.1728 0.2412
## x          -3.0822e+39 2.6314e+39 -1.1713 0.2417
```

The Bootstrap

- The bootstrap is a statistical tool that can be used to quantify the uncertainty associated with a given estimator or statistical method.
 - ▶ For example, it can provide an estimate of the standard error of a coefficient.
- The term is believed to derive from “The Surprising Adventures of Baron Munchausen” by Rudolph Erich Raspe¹⁴.

The Baron had fallen to the bottom of a deep lake. Just when it looked like all was lost, he thought to pick himself up by his own bootstraps.

¹⁴We also have the Munchausen number – a number that is equal to the sum of each digit raised to the power of itself. E.g., $3435 = 3^3 + 4^4 + 3^3 + 5^5$.

The Bootstrap



Münchhausen

O. Herrfurth pinx

Baron Munchausen
pulls himself out of
a mire by his own
hair (illustration by
Oskar Herrfurth)

The Bootstrap

- The idea is simple: suppose we want to estimate the standard error of $\hat{\beta}$. If we can generate many independent samples from the underlying population, then we can estimate our model on *each* sample, based on which we can compute an estimate of $V(\hat{\beta})$ ¹⁵.

¹⁵For an illustration, see page 42 - 44.

The Bootstrap

- In practice, we do not have access to the underlying population and do not know the true $p(x, y)$, but we can treat the **empirical distribution** as an estimate of the true distribution and draw new samples *out of the observed sample* itself.
- Each generated bootstrap sample contains the *same* number of observations as the original observed sample¹⁶. To accomplish this, we repeatedly draw observations from the observed sample **with replacement**.

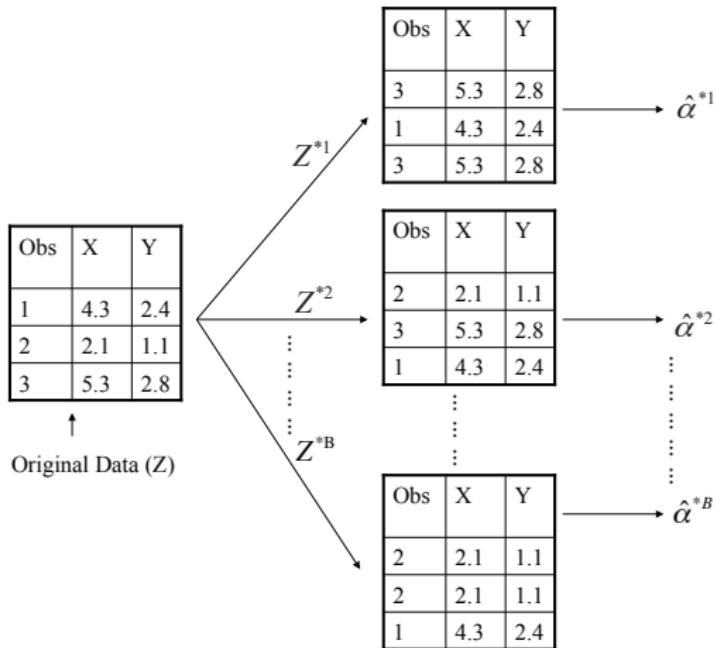
¹⁶This is because $\widehat{\beta}$ is estimated on a sample of size N . To quantify its uncertainty, we need to generate many samples of the same size.

The Bootstrap

Suppose we observe the following data set: $x = \{0, 0, 0, 1, 1\}$. Then the empirical distribution of x is $\Pr(x = 0) = 0.6$ and $\Pr(x = 1) = 0.4$.

If we believe this is a good approximation of the underlying true distribution of x , then we can generate new samples from this distribution. In practice, this can be accomplished by drawing from $\{0, 0, 0, 1, 1\}$ with replacement.

The Bootstrap



The bootstrap approach on a sample containing 3 observations.

Portfolio Choice

We wish to invest a fixed sum of money in two financial assets that yield returns of X and Y . Suppose our goal is to minimize the total risk, or variance, of our investment. Then the problem is to choose α such that

$$\alpha = \arg \min_{\gamma} \mathbb{V}(\gamma X + (1 - \gamma) Y) \quad (16)$$

(16) \Rightarrow

$$\alpha = \frac{\sigma_Y^2 - \sigma_{XY}}{\sigma_X^2 + \sigma_Y^2 - 2\sigma_{XY}} \quad (17)$$

Portfolio Choice

- Given data on (X, Y) , we can estimate $\hat{\sigma}_X^2, \hat{\sigma}_Y^2, \hat{\sigma}_{XY}$ and compute $\hat{\alpha}$.
- To estimate the standard error of $\hat{\alpha}$, we generate R bootstrap samples from the observed data and estimate α R times \Rightarrow

$$\hat{\alpha} = \frac{1}{R} \sum_{r=1}^R \hat{\alpha}_r$$
$$\hat{s.e.}(\hat{\alpha}) = \sqrt{\frac{1}{R-1} \sum_{r=1}^R (\hat{\alpha}_r - \hat{\alpha})^2}$$

Portfolio Choice

```
# Function to calculate alpha
alpha <- function(data,index){
  X <- data$X[index]
  Y <- data$Y[index]
  return((var(Y)-cov(X,Y))/(var(X)+var(Y)-2*cov(X,Y)))
}

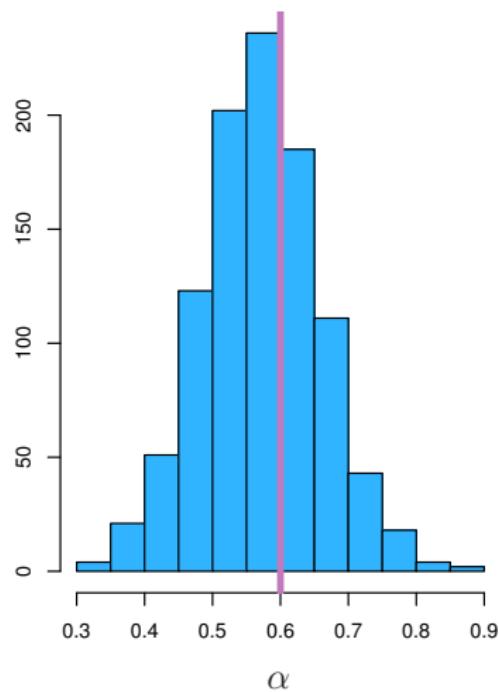
# 'Portfolio' is a simulated data set containing the returns of X and Y
require(ISLR) # contains 'Portfolio'
n <- nrow(Portfolio)
bootsample <- sample(n,n,replace=T) # generate one bootstrap sample
alpha(Portfolio,bootsample) # calculate alpha based on the bootstrap sample
## [1] 0.6618485
```

Portfolio Choice

```
# Calculate alpha based on 1000 bootstrap samples
require(boot)
boot(Portfolio,alpha,R=1000)

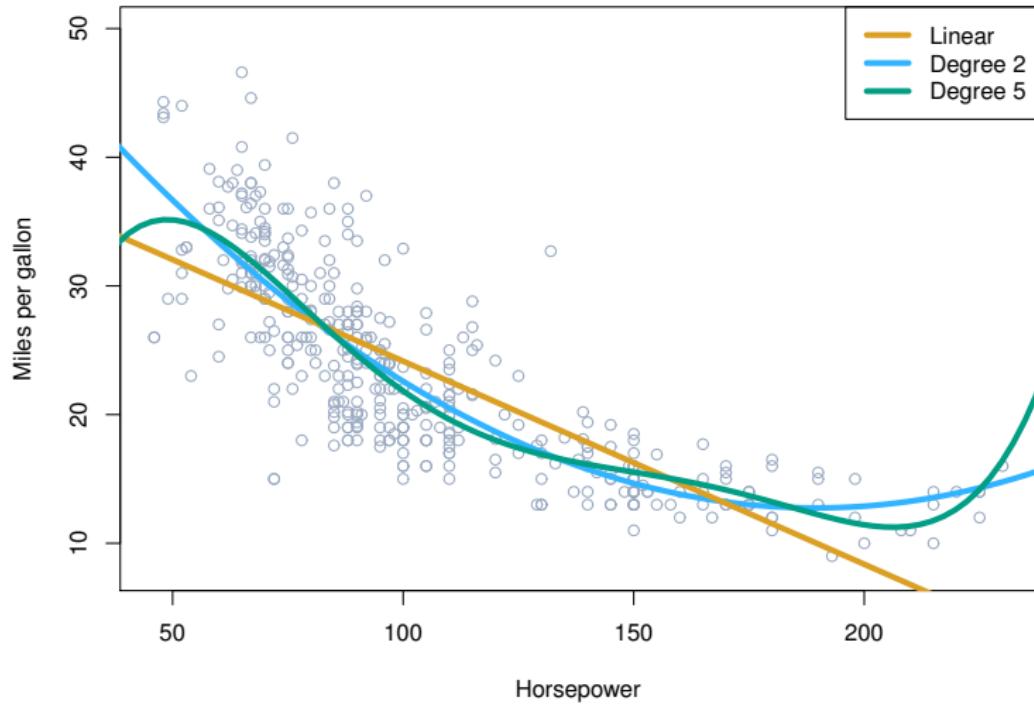
##
## ORDINARY NONPARAMETRIC BOOTSTRAP
##
##
## Call:
## boot(data = Portfolio, statistic = alpha, R = 1000)
##
##
## Bootstrap Statistics :
##      original     bias   std. error
## t1*  0.5758321 0.007540109  0.0895633
```

Portfolio Choice



Histogram of $\hat{\alpha}$ obtained from bootstrap samples

MPG and Horsepower



MPG and Horsepower

```
require(ISLR) # contains the data set 'Auto'  
require(boot)  
beta <- function(data,index){  
  coef(lm(mpg~horsepower,data=data,subset=index))  
}  
boot(Auto,beta,R=1000) # bootstrap std err  
  
##  
## ORDINARY NONPARAMETRIC BOOTSTRAP  
##  
##  
## Call:  
## boot(data = Auto, statistic = beta, R = 1000)  
##  
##  
## Bootstrap Statistics :  
##      original      bias    std. error  
## t1* 39.9358610  0.0378395594 0.877270914  
## t2* -0.1578447 -0.0004406186 0.007475398
```

MPG and Horsepower

```
require(AER)
coeftest(lm(mpg ~ horsepower, data=Auto)) # homoskedastic std err

##
## t test of coefficients:
##
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept) 39.9358610  0.7174987  55.660 < 2.2e-16 ***
## horsepower -0.1578447  0.0064455 -24.489 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

MPG and Horsepower

```
coeftest(lm(mpg ~ horsepower, data=Auto), vcov=vcovHC) # robust std err

##
## t test of coefficients:
##
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 39.9358610  0.8644903  46.196 < 2.2e-16 ***
## horsepower -0.1578447  0.0074943 -21.062 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Log-Linear Regression

When y changes on a multiplicative or percentage scale, it is often appropriate to use $\log(y)$ as the dependent variable¹⁷:

$$y = Ae^{\beta x + e} \Rightarrow \log(y) = \log(A) + \beta x + e$$

e.g.,

$$\log(\text{GDP}) = \alpha + g \times t + e$$

, where t = years, α = \log (base year GDP), and g = annual growth rate.

¹⁷ Suppose y grows at a rate i . If i is continuously compounded, then $y_t = y_0 \lim_{n \rightarrow \infty} \left(1 + \frac{i}{n}\right)^{nt} = y_0 e^{it} \Rightarrow \log(y_t) = \log(y_0) + i \times t$. If i is not continuously compounded, then $y_t = y_0 (1 + i)^t \Rightarrow \log(y_t) = \log(y_0) + t \log(1 + i) \approx \log(y_0) + i \times t$.

Log-Linear Regression

18

¹⁸ Note: in general, $E[f(y)] \neq f(E[y])$. In particular, by the Jensen's inequality, $E[\log(y)] < \log(E[y])$. Therefore, if $E[\log(y)|x] = \alpha + \beta x$, then $E[y|x] > \exp(\alpha + \beta x)$.

If we are willing to assume

$$\log(y) = \alpha + \beta x + e, \quad e \sim \mathcal{N}(0, \sigma^2)$$

, then we have: $E[y|x] = \exp(\alpha + \beta x + \frac{1}{2}\sigma^2) = \exp(E[\log(y)|x] + \frac{1}{2}\sigma^2)$.

Elasticity and Log-Log Regression

In a log-log model:

$$\log(y) = \beta_0 + \beta_1 \log(x) + e$$

β_1 can often be interpreted as an elasticity measure:

$$\beta_1 = \frac{\partial \log(y)}{\partial \log(x)} = \frac{\partial y / y}{\partial x / x} \approx \frac{\% \Delta y}{\% \Delta x}$$

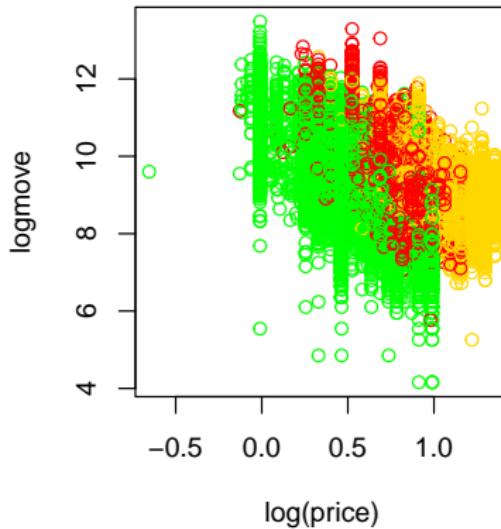
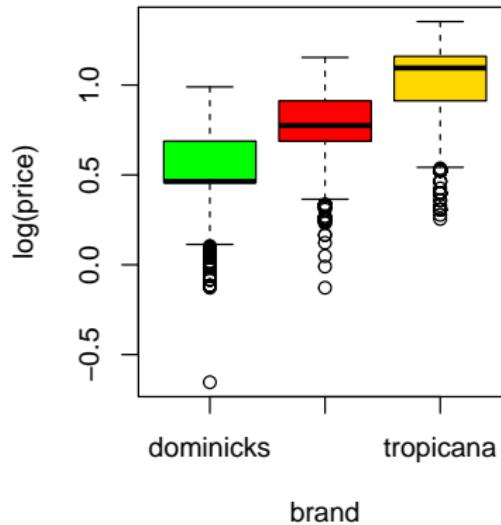
e.g.,

$$\log(\text{sales}) = \beta_0 + \beta_1 \log(\text{price}) + e$$

Orange Juice

Three brands: Tropicana, Minute Maid, Dominick's

Data from 83 stores on price, sales (units moved), and whether featured in the store



Orange Juice

$$\log(\text{sales}) = \alpha + \beta \log(\text{price}) + \epsilon$$

```
require(AER)
oj <- read.csv('oj.csv')
reg1 <- lm(logmove ~ log(price), data=oj)
coeftest(reg1)

##
## t test of coefficients:
##
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 10.4234    0.0154   679.0   <2e-16 ***
## log(price) -1.6013    0.0184   -87.2   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Orange Juice

$\log(\text{sales}) = \alpha_b + \beta_b \log(\text{price}) + e$, where b denotes brand

```
reg2 <- lm(logmove ~ log(price)*brand, data=obj)
coeftest(reg2)

## 
## t test of coefficients:
## 

##                               Estimate Std. Error t value Pr(>|t|)    
## (Intercept)             10.9547    0.0207 529.14   <2e-16 ***
## log(price)            -3.3775    0.0362 -93.32   <2e-16 ***
## brandminute.maid      0.8883    0.0416  21.38   <2e-16 ***
## brandtropicana        0.9624    0.0464  20.72   <2e-16 ***
## log(price):brandminute.maid 0.0568    0.0573   0.99     0.32  
## log(price):brandtropicana  0.6658    0.0535  12.44   <2e-16 ***
## --- 
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Orange Juice

$$\log(\text{sales}) = (\alpha_{0b} + \text{feature} \times a_{1b}) + (\beta_{0b} + \text{feature} \times \beta_{1b}) \times \log(\text{price}) + e$$

```
reg3 <- lm(logmove ~ log(price)*brand*feat, data=oj)
coeftest(reg3)
```

```
##  
## t test of coefficients:  
##  
##  
## (Intercept)           10.4066    0.0234   445.67 < 2e-16 ***  
## log(price)            -2.7742    0.0388  -71.45 < 2e-16 ***  
## brandminute.maid      0.0472    0.0466    1.01    0.31  
## brandtropicana         0.7079    0.0508   13.94 < 2e-16 ***  
## feat                  1.0944    0.0381   28.72 < 2e-16 ***  
## log(price):brandminute.maid 0.7829    0.0614   12.75 < 2e-16 ***  
## log(price):brandtropicana 0.7358    0.0568   12.95 < 2e-16 ***  
## log(price):feat        -0.4706    0.0741   -6.35  2.2e-10 ***  
## brandminute.maid:feat   1.1729    0.0820   14.31 < 2e-16 ***  
## brandtropicana:feat     0.7853    0.0987    7.95  1.9e-15 ***  
## log(price):brandminute.maid:feat -1.1092    0.1222  -9.07 < 2e-16 ***  
## log(price):brandtropicana:feat -0.9861    0.1241  -7.95  2.0e-15 ***  
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Orange Juice

- Elasticity¹⁹: -1.6
- Brand-specific elasticities:

Dominick's: -3.4 , Minute Maid: -3.4 , Tropicana: -2.7

- How does featuring a product affect its elasticity?

	Dominick's	Minute Maid	Tropicana
not featured	-2.8	-2.0	-2.0
featured	-3.2	-3.6	-3.5

¹⁹What economic assumptions need to be satisfied in order for the coefficients to be interpreted as price elasticities of demand?

CAPM

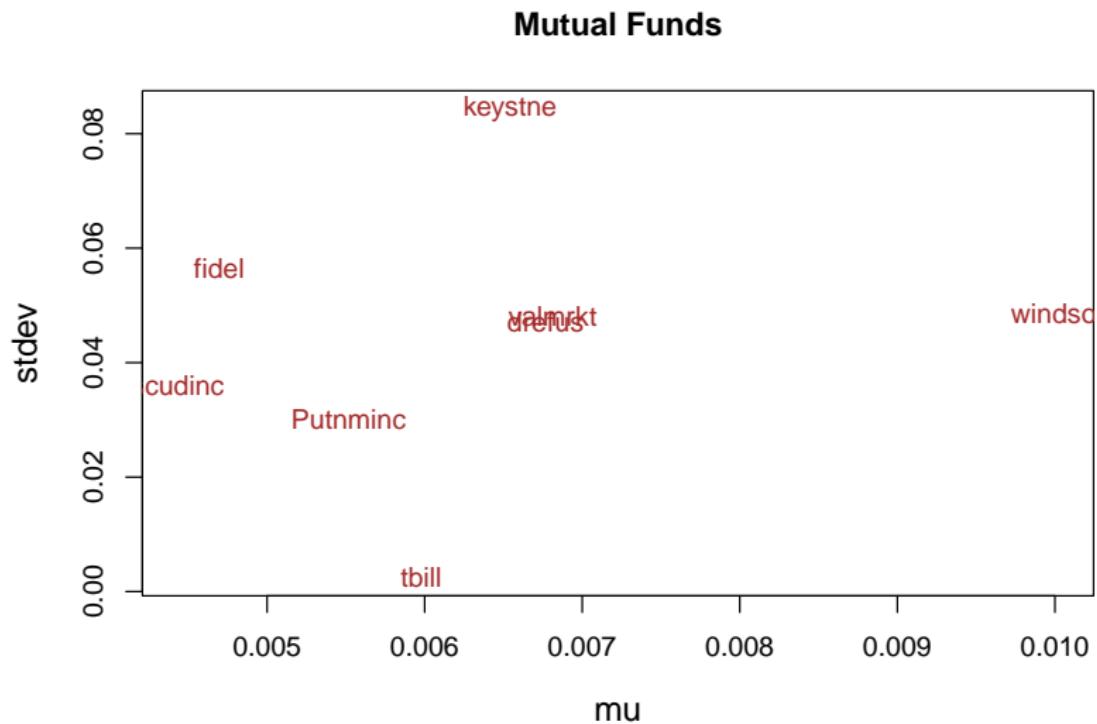
The Capital Asset Pricing Model (CAPM) for asset A relates return $R_{A,t}$ to the market return, $R_{M,t}$:

$$R_{A,t} = \alpha + \beta R_{M,t} + e$$

When asset A is a mutual fund, this CAPM regression can be used as a performance benchmark for fund managers.

```
# 'mfunds.csv' contains data on the historical returns of
#                   6 mutual funds as well as the market return
mfund <- read.csv('mfunds.csv')
mu <- apply(mfund, 2, mean)
stdev <- apply(mfund, 2, sd)
```

CAPM

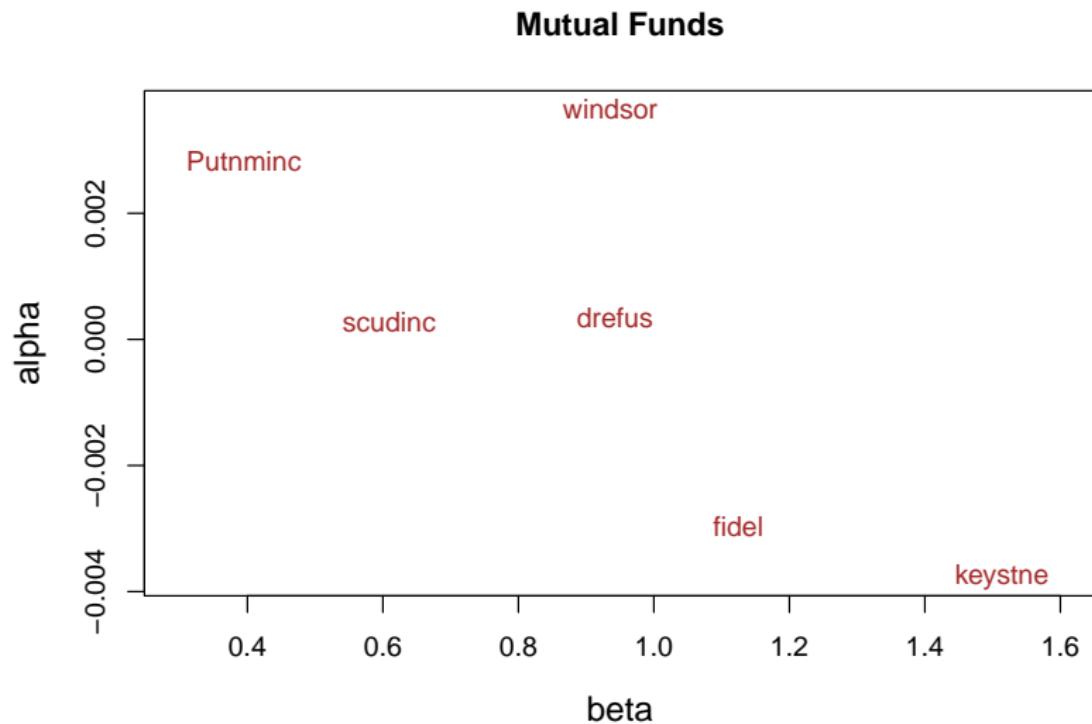


CAPM

```
CAPM <- lm(as.matrix(mfund[,1:6]) ~ mfund$valmrkt)
CAPM

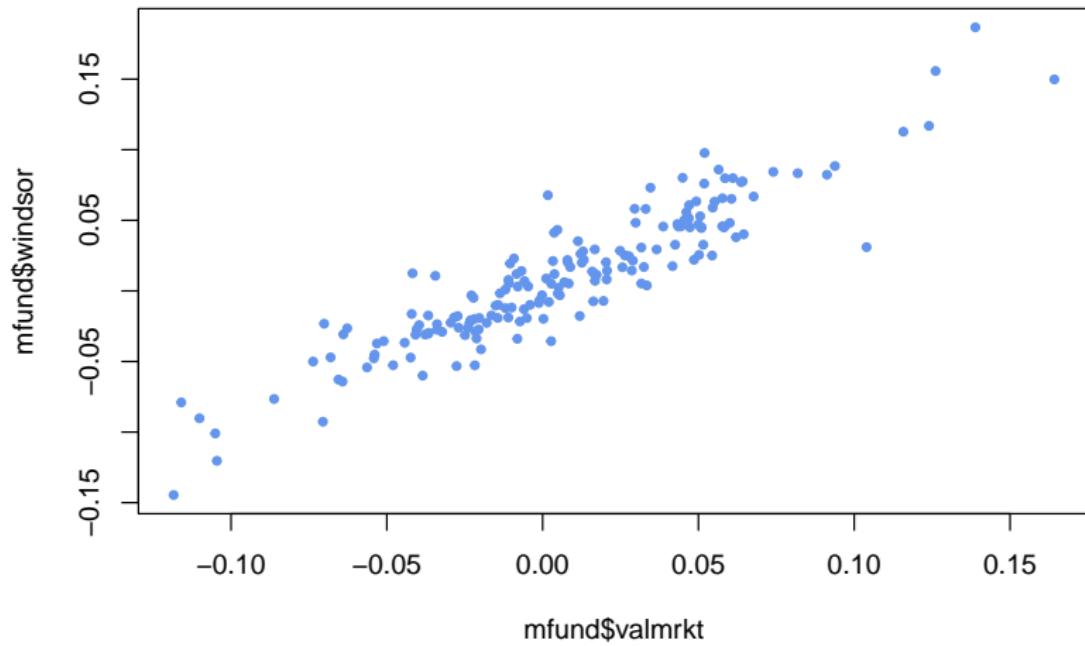
##
## Call:
## lm(formula = as.matrix(mfund[, 1:6]) ~ mfund$valmrkt)
##
## Coefficients:
##             drefus      fidel     keystnne     Putnminc     scudinc
## (Intercept) 0.0003462 -0.0029655 -0.0037704  0.0028271  0.000281
## mfund$valmrkt 0.9424286  1.1246549  1.5137186  0.3948280  0.609202
##             windsor
## (Intercept) 0.0036469
## mfund$valmrkt 0.9357170
```

CAPM



CAPM

Look at windsor (which dominates the market):



CAPM

Does Windsor have an “alpha” over the market?

$H_0 : \alpha = 0$ vs. $H_1 : \alpha \neq 0$

```
require(AER)
reg <- lm(mfund$windsor ~ mfund$valmrkt)
coeftest(reg)

##
## t test of coefficients:
##
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.0036469  0.0014094  2.5876  0.01046 *
## mfund$valmrkt 0.9357170  0.0291499 32.1002 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

CAPM

Now look at beta:

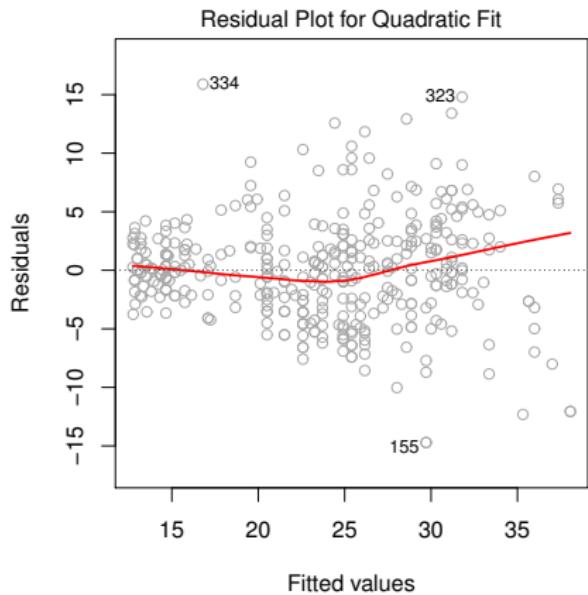
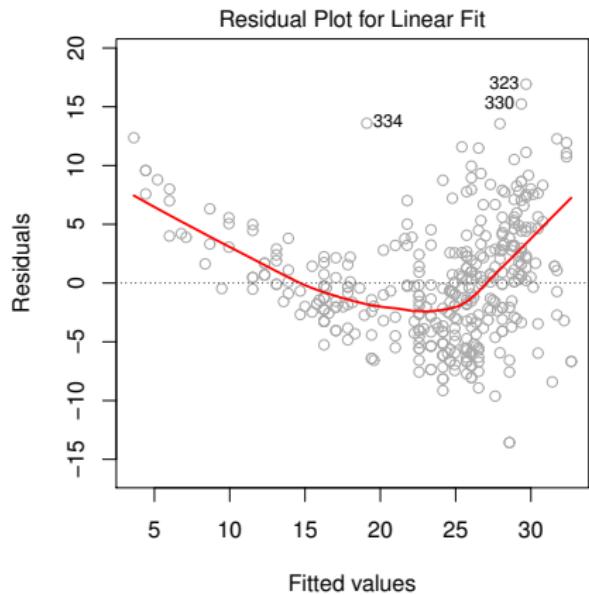
$H_0 : \beta = 1$, Windsor is just the market (+ alpha).

$H_1 : \beta \neq 1$, Windsor softens or exaggerates market moves.

```
linearHypothesis(reg, "mfund$valmrkt = 1")

## Linear hypothesis test
##
## Hypothesis:
## mfund$valmrkt = 1
##
## Model 1: restricted model
## Model 2: mfund$windsor ~ mfund$valmrkt
##
##          Res.Df      RSS Df Sum of Sq      F    Pr(>F)
## 1      179 0.064082
## 2      178 0.062378  1 0.0017042 4.8632 0.02872 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Regression Diagnostics: Is the CEF linear?



Regression Diagnostics: Is the CEF linear?

Anscombe's quartet comprises four datasets that have similar statistical properties ...

```
anscombe <- read.csv('anscombe.csv')
attach(anscombe)
c(x.m1=mean(x1),x.m2=mean(x2),x.m3=mean(x3),x.m4=mean(x4))

## x.m1 x.m2 x.m3 x.m4
##    9    9    9    9

c(y.m1=mean(y1),y.m2=mean(y2),y.m3=mean(y3),y.m4=mean(y4))

##      y.m1      y.m2      y.m3      y.m4
## 7.500909 7.500909 7.500000 7.500909
```

Regression Diagnostics: Is the CEF linear?

```
c(x.sd1=sd(x1),x.sd2=sd(x2),x.sd3=sd(x3),x.sd3=sd(x4))

##      x.sd1      x.sd2      x.sd3      x.sd3
## 3.316625 3.316625 3.316625 3.316625

c(y.sd1=sd(y1),y.sd2=sd(y2),y.sd4=sd(y3),y.sd4=sd(y4))

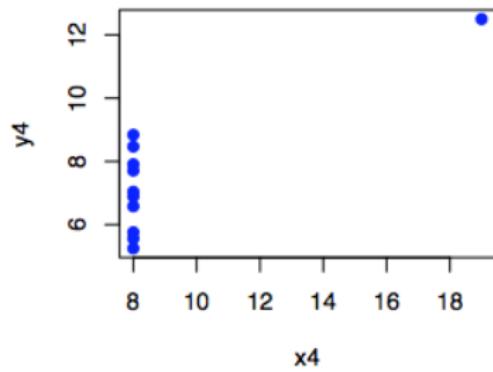
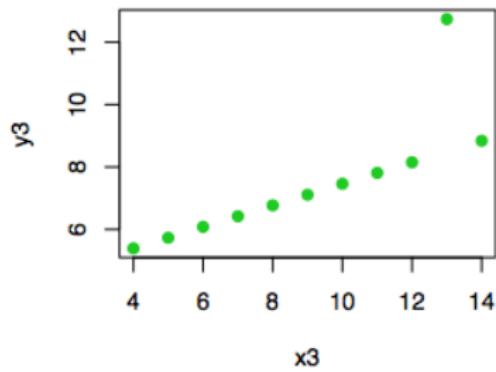
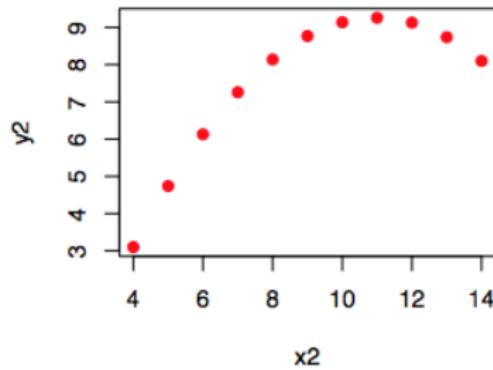
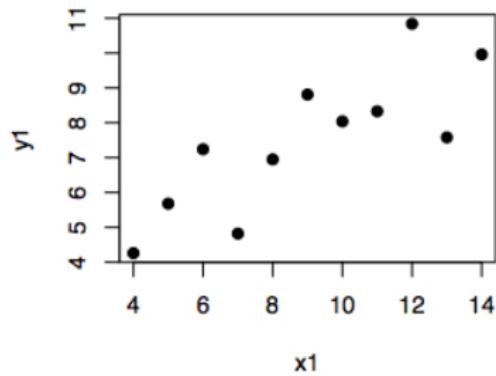
##      y.sd1      y.sd2      y.sd4      y.sd4
## 2.031568 2.031657 2.030424 2.030579

c(cor1=cor(x1,y1),cor2=cor(x2,y2),cor3=cor(x3,y3),cor4=cor(x4,y4))

##      cor1      cor2      cor3      cor4
## 0.8164205 0.8162365 0.8162867 0.8165214
```

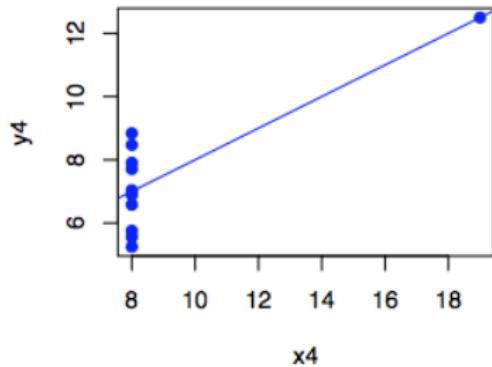
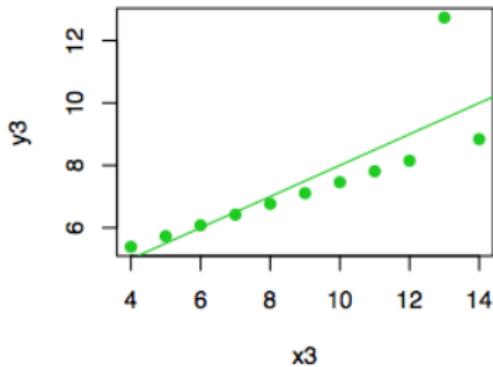
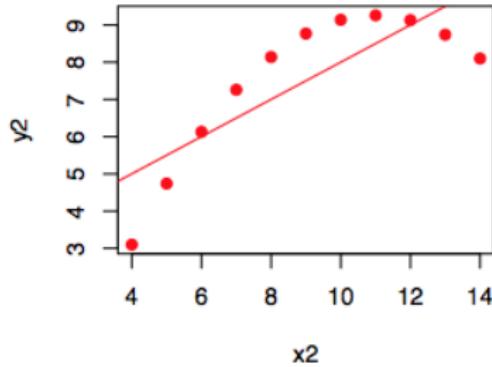
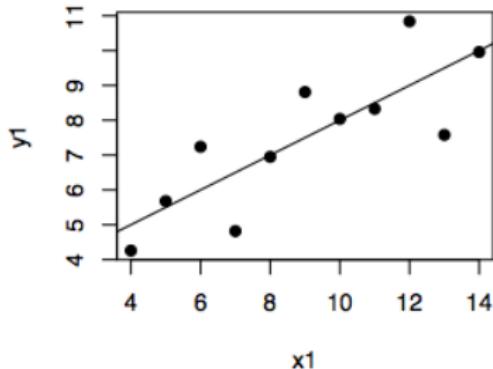
Regression Diagnostics: Is the CEF linear?

...but vary considerably when graphed:



Regression Diagnostics: Is the CEF linear?

Linear regression on each dataset:



Regression Diagnostics: Is the CEF linear?

The regression lines and R^2 values are the same...

```
areg <- list(areg1=lm(y1~x1),areg2=lm(y2~x2),areg3=lm(y3~x3),areg4=lm(y4~x4))
attach(areg)
cbind(areg1$coef,areg2$coef,areg3$coef,areg4$coef)

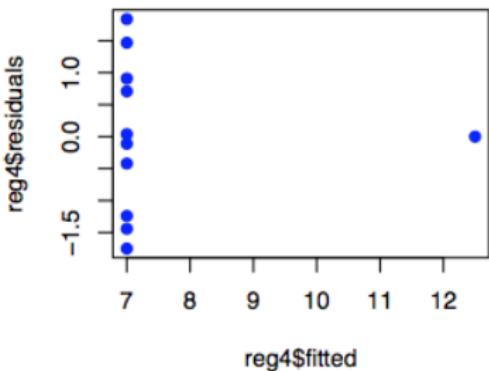
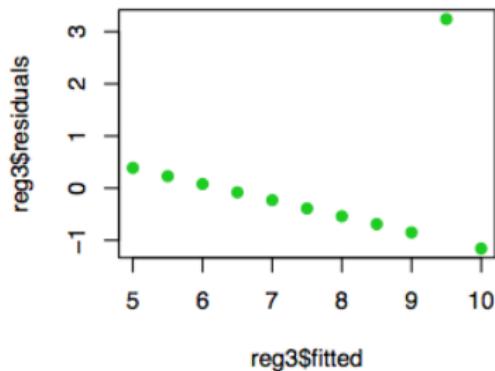
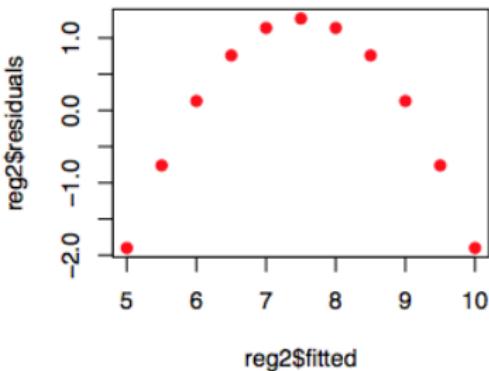
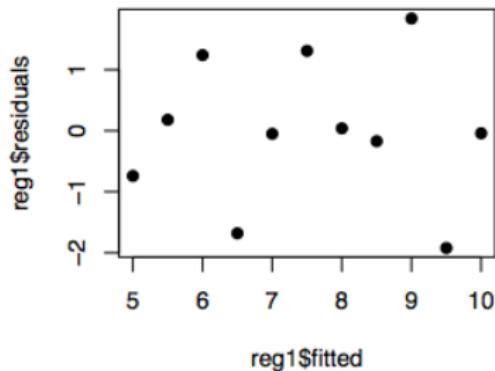
##           [,1]      [,2]      [,3]      [,4]
## (Intercept) 3.0000909 3.0000909 3.0024545 3.0017273
## x1          0.5000909 0.5000000 0.4997273 0.4999091

s <- lapply(areg,summary)
c(s$areg1$r.sq,s$areg2$r.sq,s$areg3$r.sq,s$areg4$r.sq)

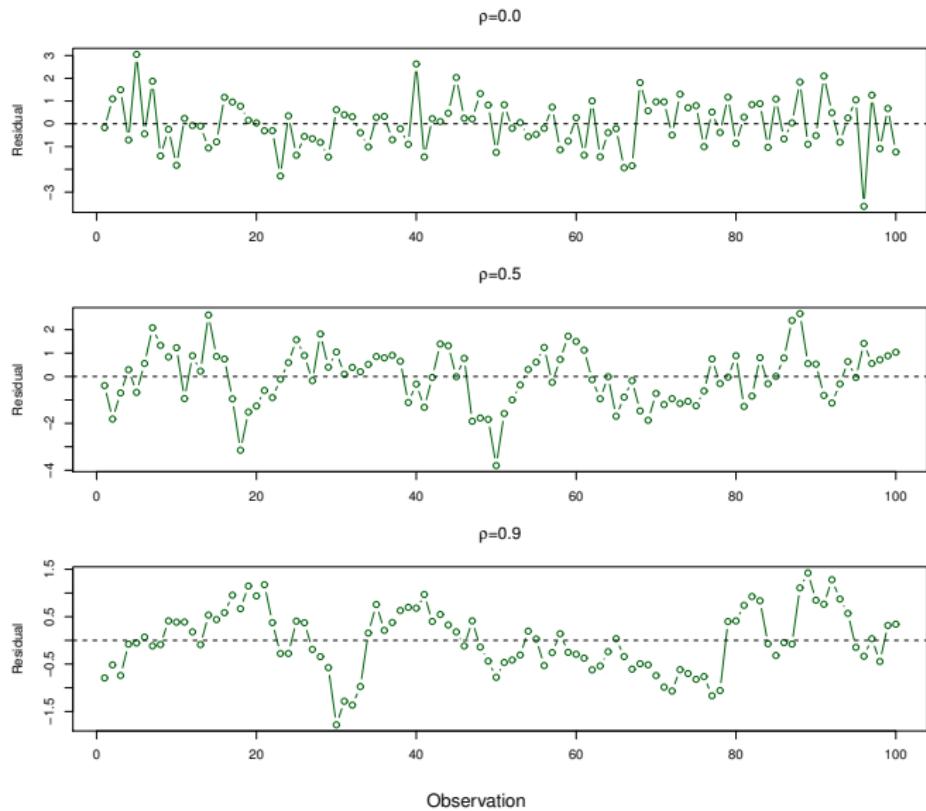
## [1] 0.6665425 0.6662420 0.6663240 0.6667073
```

Regression Diagnostics: Is the CEF linear?

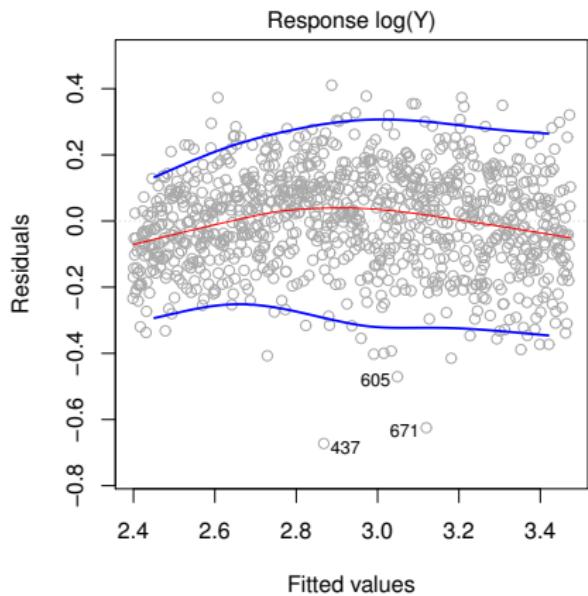
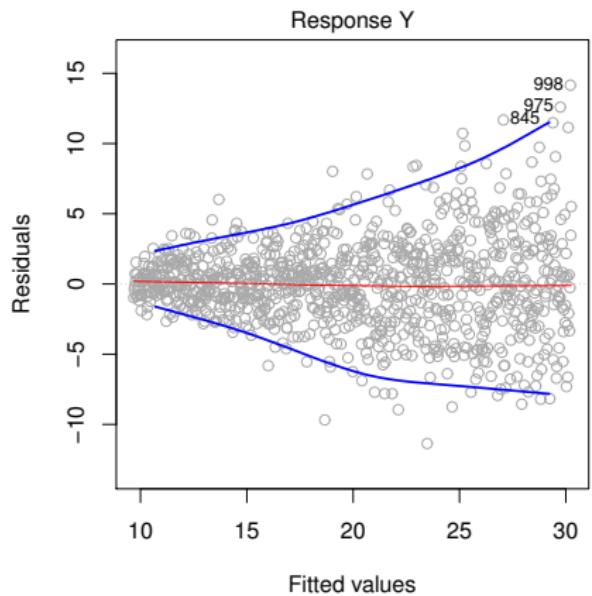
...but residual plots show the differences:



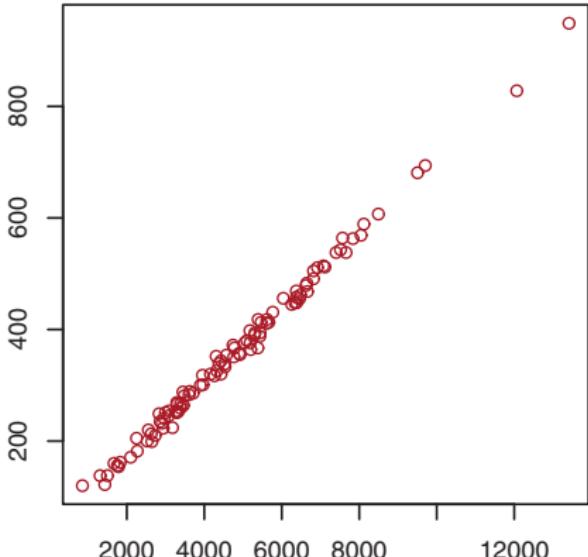
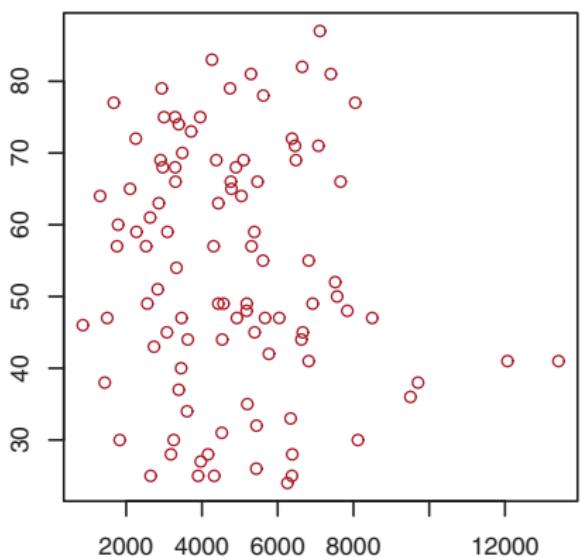
Regression Diagnostics: Nonrandom Sampling



Regression Diagnostics: Heteroskedasticity

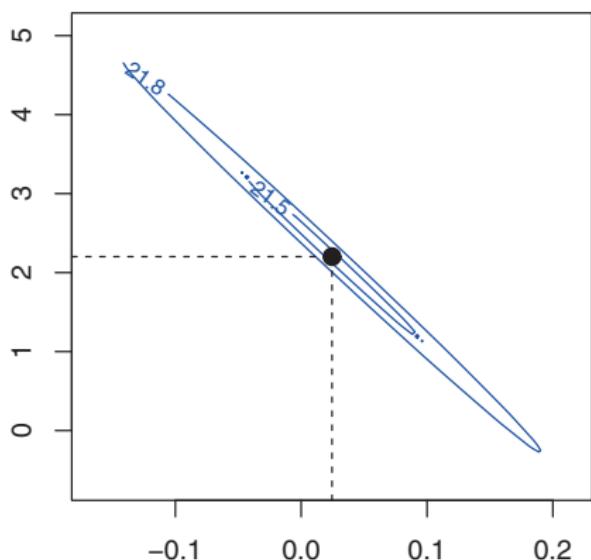
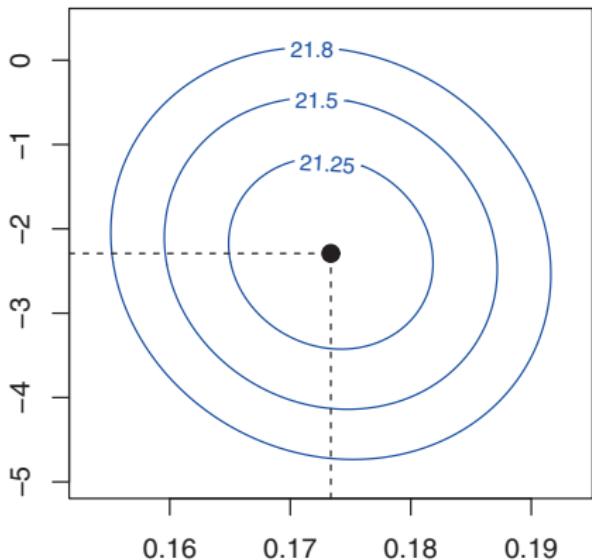


Regression Diagnostics: Collinearity



Left: two variables are not collinear; Right: there is high collinearity.

Regression Diagnostics: Collinearity



Contour plots for the RSS as a function of parameters β .

Left: regression on two variables that are not collinear; Right: regression on two variables that are highly collinear. Because of the collinearity, there are many pairs of β with a similar value for the RSS.

Regression Diagnostics: Collinearity

From (9), we can see that:

- If X_1, \dots, X_p are orthogonal, then $\hat{\beta}_j$ is equal to the simple linear regression coefficient of y on X_j .
- If X_1, \dots, X_p are correlated – in particular – if X_j is highly correlated with the other predictors, then \hat{u}_j will be close to 0. This makes $\hat{\beta}_j$ **unstable**, as both the denominator and the numerator are small.

From (13), we can see that:

- If X_j is highly correlated with the other predictors, the variance of $\hat{\beta}_j$ is **inflated**, making it less likely to be significant.

Regression Diagnostics: Collinearity

- A simple way to detect collinearity is to look at the correlation matrix of the predictors.
- However, it is possible for collinearity to exist between three or more variables even if no pair of variables has a particularly high correlation. This is called **multicollinearity**.
- **Variance inflation factor (VIF)**:

$$\text{VIF}(\hat{\beta}_j) = \frac{1}{1 - R_{X_j|X_{-j}}^2}$$

, where $R_{X_j|X_{-j}}^2$ is the R^2 from a regression of X_j onto all of the other predictors.

- ▶ $\text{VIF} \geq 1$. Large VIF indicates a problematic amount of collinearity.

Regression Diagnostics: Collinearity

- When faced with the problem of collinearity, a simple solution is to drop one of the problematic variables.
- Suppose two variables both contribute in explaining y , but are highly correlated with each other.
 - ▶ Both will be insignificant if both are included in the regression model.
 - ▶ Dropping one will likely make the other significant.
- This is why we can't remove two (or more) supposedly insignificant predictors *at a time*: significance depends on what other predictors are in the model!

Maximum Likelihood Estimation

- While least squares regression learns a deterministic function $f(x)$ that directly maps each x into a prediction of y , an alternative approach is to learn the conditional distribution $p(y|x)$ and use the estimated $p(y|x)$ to form a prediction of y .
- To do so, let $\mathcal{H} = \{q_\theta(y|x) : \theta \in \Theta\}$, where the hypotheses $q_\theta(y|x)$ are conditional distributions parametrized by $\theta \in \Theta$.
- We select a $q_\theta(y|x) \in \mathcal{H}$, or equivalently, a $\theta \in \Theta$, by minimizing the empirical KL divergence, or equivalently, by maximizing the (log) likelihood function.

Maximum Likelihood Estimation

The log likelihood function²⁰:

$$\log \mathcal{L}(\theta) = \sum_{i=1}^N \log q_\theta(y_i | x_i)$$

The maximum likelihood estimator chooses

$$\hat{\theta} = \arg \max_{\theta \in \Theta} \log \mathcal{L}(\theta)$$

²⁰Also written as $\log \mathcal{L}(\theta | \mathcal{D})$ to emphasize its dependence on sample \mathcal{D} .

Normal Linear Model

The normal linear regression model is $\mathcal{H} = \{q_\theta(y|x)\}$, where

$$q_\theta(y|x) = \mathcal{N}(x'\beta, \sigma^2) \quad (18)$$

, where $\theta = (\beta, \sigma)$.

This is equivalent to assuming²¹:

$$y = x'\beta + e, \quad e \sim \mathcal{N}(0, \sigma^2) \quad (19)$$

²¹ Notice the strong assumptions imposed by (18) and (19). In addition to assuming a linear regression function, we are now assuming that (1) at each x , the scatter of y around the regression function is normally distributed (**Gaussianity**); (2) the variance of this scatter is constant (**homoskedasticity**); and (3) there is no dependence between this scatter and anything else (**error independence**).

Normal Linear Model

Given sample \mathcal{D} and model (18),

$$\begin{aligned}\log \mathcal{L} &= \sum_{i=1}^N \log \left\{ \frac{1}{\sqrt{2\pi}\sigma} \exp \left(-\frac{1}{2\sigma^2} (y_i - x_i' \beta)^2 \right) \right\} \\ &= -\frac{N}{2} \log (2\pi) - N \log \sigma - \underbrace{\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - x_i' \beta)^2}_{\text{RSS}}\end{aligned}\tag{20}$$

Normal Linear Model

Maximizing (20) with respect to β and σ \Rightarrow

$$\frac{\partial \log \mathcal{L}}{\partial \beta} = 0 \Rightarrow \hat{\beta} = \left[\sum_{i=1}^N x_i x_i' \right]^{-1} \sum_{i=1}^N x_i y_i = (X' X)^{-1} X' Y$$

$$\frac{\partial \log \mathcal{L}}{\partial \sigma} = 0 \Rightarrow \hat{\sigma} = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - x_i' \hat{\beta})^2}$$

Thus, maximum likelihood estimation of the normal linear model produces the same estimate of β as least squares regression.

Normal Linear Model

Let's fit the normal linear model (18) on the data we generated on [page 20](#):

```
# Define the negative log likelihood function
nll <- function(theta){
  beta0 <- theta[1]
  beta1 <- theta[2]
  beta2 <- theta[3]
  sigma <- theta[4]
  N <- length(y)
  z <- (y - beta0 - beta1*x1 - beta2*x2)/sigma
  logL <- -1*N*log(sigma) - 0.5*sum(z^2)
  return(-logL)}

# Minimize the negative likelihood function
mlefit <- optim(c(0,0,0,1),nll) # initial value for theta: (0,0,0,1)
mlefit$par # parameter estimate

## [1] 1.010153 -2.591790 5.062709 1.004935
```

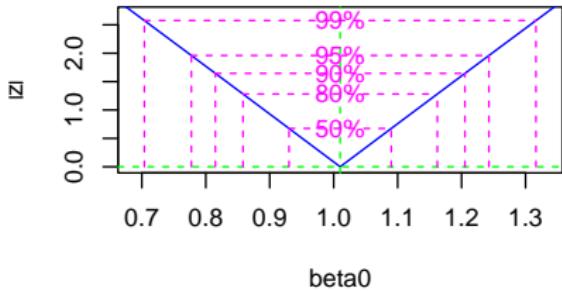
Normal Linear Model

```
# Alternatively,
require(bbmle)
parnames(nll) <- c("beta0", "beta1", "beta2", "sigma")
result <- mle2(nll, start=c(beta0=0, beta1=0, beta2=0, sigma=1))
coeftest(result)

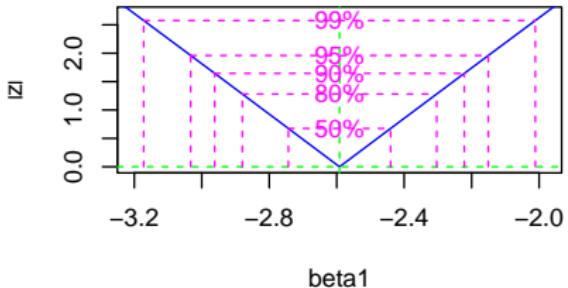
##
## z test of coefficients:
##
##          Estimate Std. Error   z value Pr(>|z|)
## beta0    1.010134   0.118487   8.5253 < 2.2e-16 ***
## beta1   -2.591654   0.224609 -11.5385 < 2.2e-16 ***
## beta2    5.062493   0.311189  16.2682 < 2.2e-16 ***
## sigma   1.004913   0.031778  31.6227 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Normal Linear Model

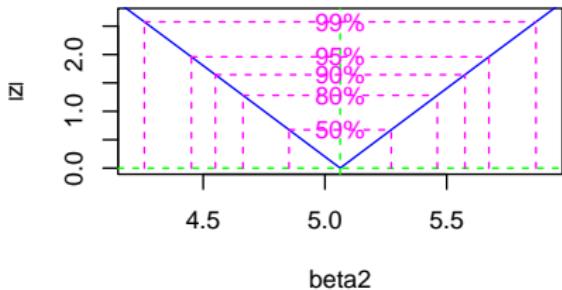
Likelihood profile: β_0



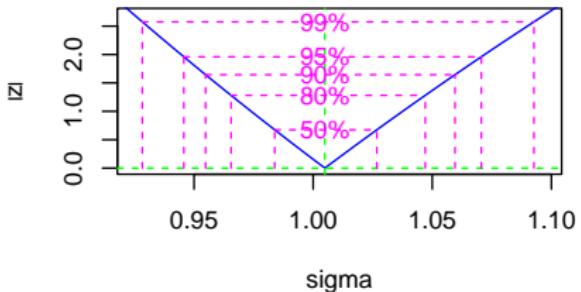
Likelihood profile: β_1



Likelihood profile: β_2



Likelihood profile: σ

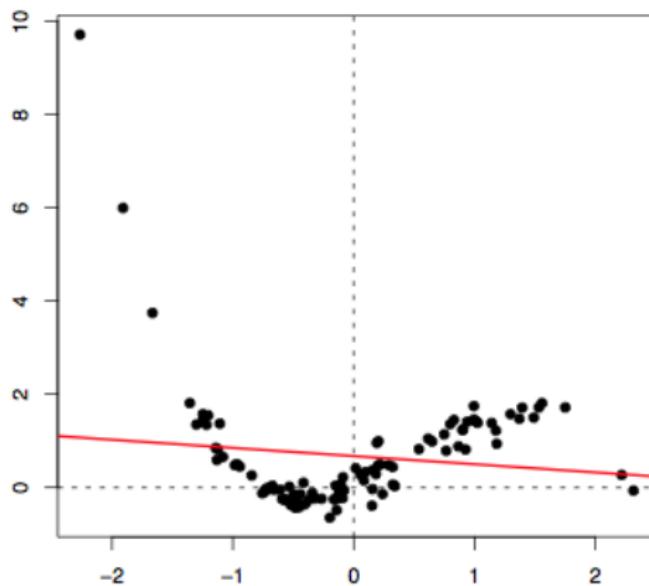


Moving Beyond Linearity

- The CEF $f(x) = \mathbb{E}(y|x)$ is seldom linear. The least squares linear regression model, however, doesn't have to be linear in x either. We can move beyond linearity in inputs x as long as we retain linearity in parameters β^{22} .
- Polynomial regression is a standard way to extend linear regression to settings in which the relationship between x and y is nonlinear.

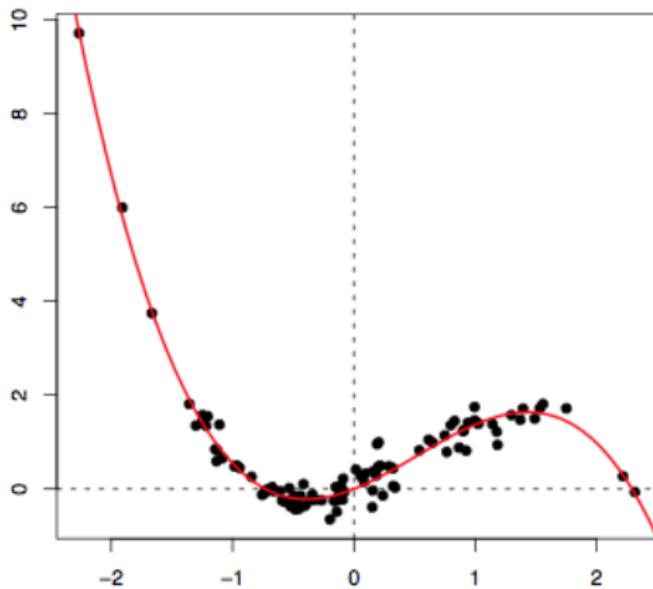
²²We have already seen examples of including nonlinear terms in x such as $\log(x)$ and interaction effects (x_1x_2) in the regression model.

Polynomial Regression



$$h(x) = \beta_0 + \beta_1 x$$

Polynomial Regression



$$h(x) = \beta_0 + \beta_1x + \beta_2x^2 + \beta_3x^3$$

Wage Profile

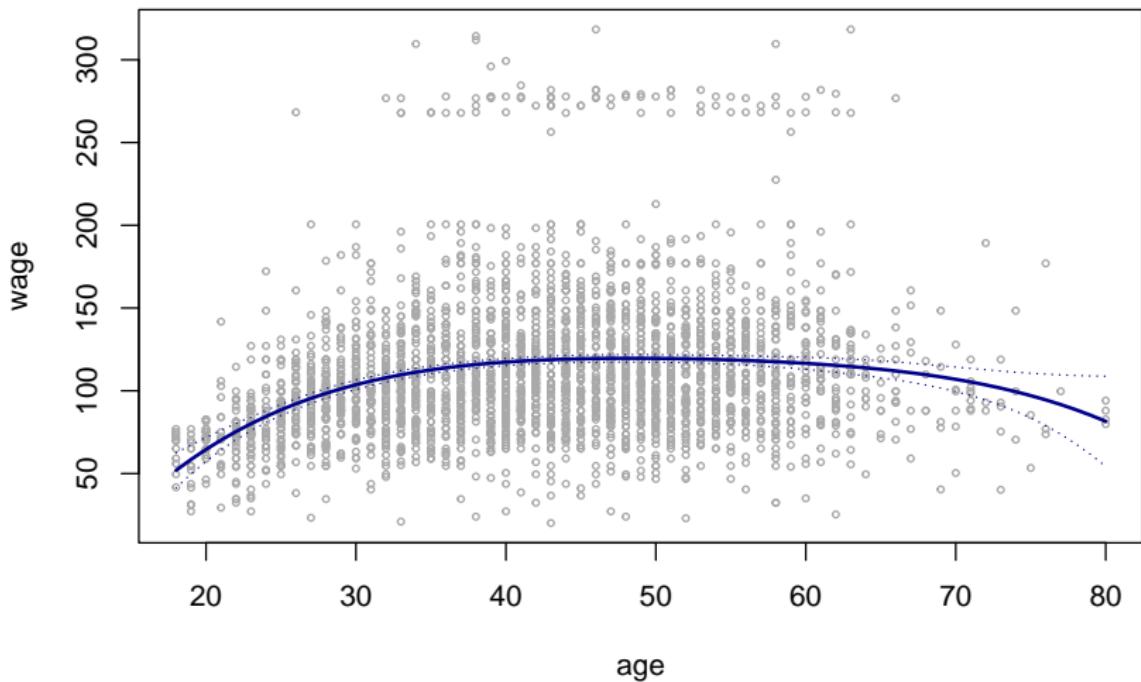
Data: income survey for men in central Atlantic region of USA

```
require(AER)
require(ISLR) # contains the data set 'Wage'
fit <- lm(wage ~ poly(age,4,raw=T), data=Wage) # degree-4 polynomial
coeftest(fit)

##
## t test of coefficients:
##
##                               Estimate Std. Error t value Pr(>|t|)
## (Intercept)           -1.8415e+02 6.0040e+01 -3.0672 0.0021803 **
## poly(age, 4, raw = T)1 2.1246e+01 5.8867e+00  3.6090 0.0003124 ***
## poly(age, 4, raw = T)2 -5.6386e-01 2.0611e-01 -2.7357 0.0062606 **
## poly(age, 4, raw = T)3  6.8107e-03 3.0659e-03  2.2214 0.0263978 *
## poly(age, 4, raw = T)4 -3.2038e-05 1.6414e-05 -1.9519 0.0510386 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Wage Profile

Degree-4 Polynomial



Piecewise Constant Regression

- For the following analysis, consider modeling the relationship between y and a single input variable x .
- So far we have imposed a global structure on the relationship between x and y .
- Piecewise regression breaks the input space into distinct regions and fit a different relationship in each region.

Piecewise Constant Regression

How it works:

1. Divide the range of x into M regions by creating $M - 1$ cutpoints, or **knots**, ξ_1, \dots, ξ_{M-1} . Then construct the following dummy variables:

Region	$\phi(x)$
R_1	$\phi_1(x) = \mathcal{I}(x < \xi_1)$
R_2	$\phi_2(x) = \mathcal{I}(\xi_1 \leq x < \xi_2)$
\vdots	\vdots
R_M	$\phi_M(x) = \mathcal{I}(\xi_{M-1} \leq x)$

Piecewise Constant Regression

How it works:

2. Fit the following model:

$$y = \beta_1\phi_1(x) + \beta_2\phi_2(x) + \cdots + \beta_M\phi_M(x) + e \quad (21)$$

$\sum_{m=1}^M \beta_m\phi_m(x)$ is a **step function** or **piecewise constant function**, and (21) is called a **piecewise constant regression** model.

Piecewise Constant Regression

Solving (21) by least squares \Rightarrow

$$\hat{\beta}_m = \bar{y}_m$$

, where $\bar{y}_m \equiv \frac{1}{n_m} \sum_{x_i \in R_m} y_i^{23}$.

i.e., for every $x \in R_m$, we make the same prediction, which is simply the mean of the response values for the training observations in R_m .

²³ n_m is the number of observations in R_m .

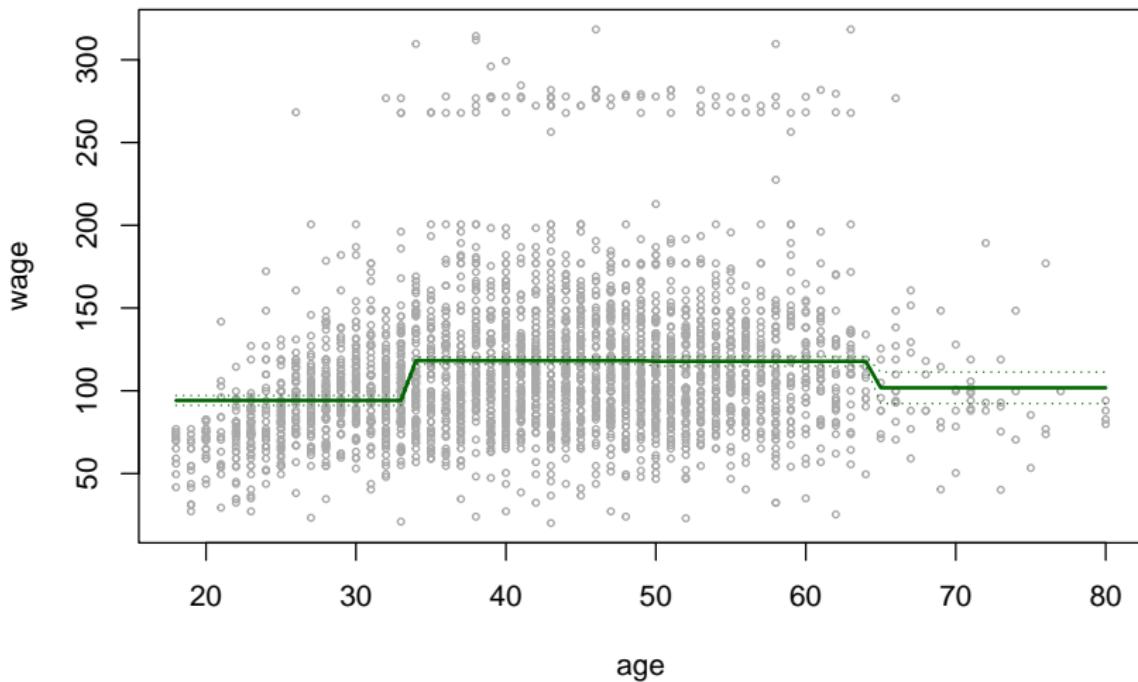
Wage Profile

```
# cut(x,M) divides x into M pieces of equal length
#           and generates the corresponding dummy variables
fit <- lm(wage ~ 0 + cut(age,4), data=Wage) # no intercept
coeftest(fit)

##
## t test of coefficients:
##
##                               Estimate Std. Error t value Pr(>|t|)
## cut(age, 4)(17.9,33.5]    94.1584     1.4761  63.790 < 2.2e-16 ***
## cut(age, 4)(33.5,49]     118.2119     1.0808 109.379 < 2.2e-16 ***
## cut(age, 4)(49,64.5]     117.8230     1.4483  81.351 < 2.2e-16 ***
## cut(age, 4)(64.5,80.1]   101.7990     4.7640  21.368 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Wage Profile

Piecewise Constant



Basis Functions

In general, $\phi(x)$ are called **basis functions** and do not have to be dummy variables. They can be any functions of x .

A **linear basis function model** is defined as²⁴:

$$y = \beta_1\phi_1(x) + \beta_2\phi_2(x) + \cdots + \beta_M\phi_M(x) + e = \beta'\Phi(x) + e \quad (22)$$

, where $\beta = (\beta_1, \dots, \beta_M)'$ and $\Phi = (\phi_1, \dots, \phi_M)'$.

Solving (22) by least squares \Rightarrow

$$\hat{\beta} = (\Phi'\Phi)^{-1}\Phi'Y$$

, where $\Phi = \Phi(X)$.

²⁴Notice that (22) is the same as (21), except now $\phi(x)$ can be any function of x .

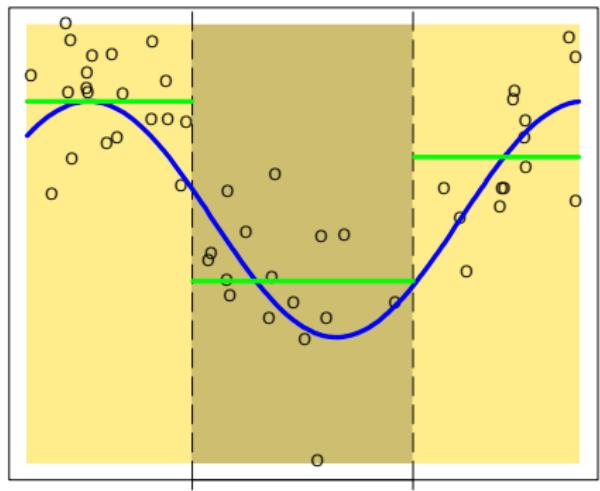
Regression Splines

- Polynomial and piecewise constant regression models are special cases of linear basis function models²⁵.
- We can also do **piecewise polynomial regression**, which involves fitting different polynomials over different regions of x .

²⁵For example, for K -degree polynomial regressions,
 $\phi_1(x) = 1, \phi_2(x) = x, \phi_3(x) = x^2, \dots, \phi_K(x) = x^K$.

Regression Splines

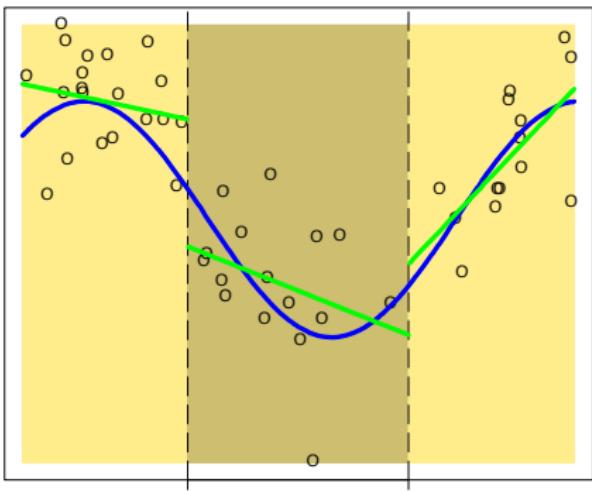
Piecewise Constant



ξ_1

ξ_2

Piecewise Linear



ξ_1

ξ_2

Regression Splines

Oftentimes it is desired that the fitted curve is **continuous** over the range of x , i.e., there should be no jump at the knots.

For piecewise linear regression with one knot (ξ), this means:

$$y = \begin{cases} \alpha_{10} + \alpha_{11}x + e & x < \xi \\ \alpha_{20} + \alpha_{21}(x - \xi) + e & x \geq \xi \end{cases} \quad (23)$$

under the constraint that

$$\alpha_{10} + \alpha_{11}\xi = \alpha_{20} \quad (24)$$

Regression Splines

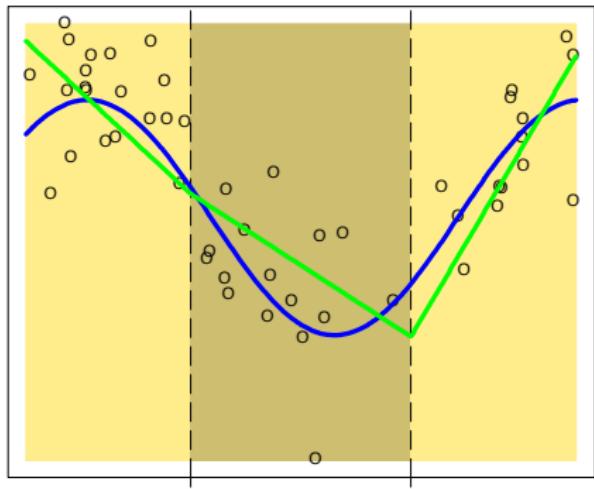
(23) and (24) \Rightarrow the continuous piecewise linear model can be parametrized as

$$y = \beta_0 + \beta_1 x + \beta_2 (x - \xi)_+ + e \quad (25)$$

, where $\beta_0 = \alpha_{10}$, $\beta_1 = \alpha_{11}$, $\beta_2 = \alpha_{21} - \alpha_{11}$, and
 $(x - \xi)_+ \equiv (x - \xi) \mathcal{I}(x \geq \xi)$.

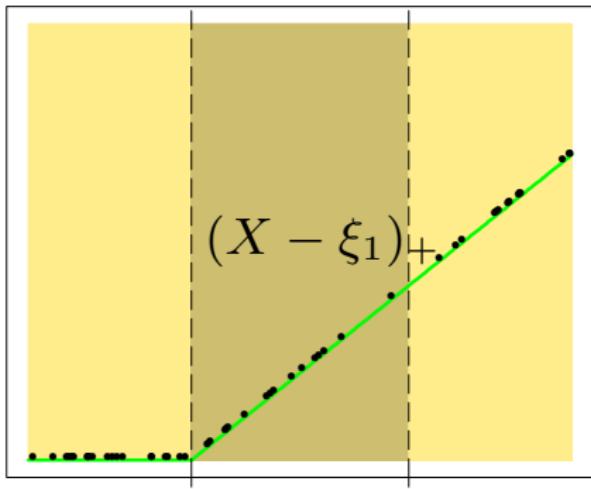
Regression Splines

Continuous Piecewise Linear



$$\xi_1 \quad \xi_2$$

Piecewise-linear Basis Function



$$\xi_1 \quad \xi_2$$

Regression Splines

For higher-order piecewise polynomial regression, in addition to the fitted curve being continuous, we may also want it to be **smooth** by requiring the derivatives of the piecewise polynomials to be also continuous at the knots.

For piecewise cubic polynomial regression with one knot (ξ), this means:

$$y = \begin{cases} \alpha_{10} + \alpha_{11}x + \alpha_{12}x^2 + \alpha_{13}x^3 + e & x < \xi \\ \alpha_{20} + \alpha_{21}(x - \xi) + \alpha_{22}(x - \xi)^2 + \alpha_{23}(x - \xi)^3 + e & x \geq \xi \end{cases} \quad (26)$$

subject to the constraints that the piecewise polynomials as well as their 1^{st} and 2^{nd} derivatives are continuous at ξ :

$$\alpha_{10} + \alpha_{11}\xi + \alpha_{12}\xi^2 + \alpha_{13}\xi^3 = \alpha_{20} \quad (27)$$

$$\alpha_{11} + 2\alpha_{12}\xi + 3\alpha_{13}\xi^2 = \alpha_{21}$$

$$\alpha_{12} + 3\alpha_{13}\xi = \alpha_{22}$$

Regression Splines

(26) and (27) \Rightarrow

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \beta_4 (x - \xi)_+^3 + e \quad (28)$$

, where $\beta_0 = \alpha_{10}$, $\beta_1 = \alpha_{11}$, $\beta_2 = \alpha_{12}$, $\beta_3 = \alpha_{13}$, and $\beta_4 = \alpha_{23} - \alpha_{13}$.

Regression Splines

(25) and (28) are examples of **regression splines**. (25) is called a *linear spline* and (28) is called a *cubic spline*.

Regression Spline

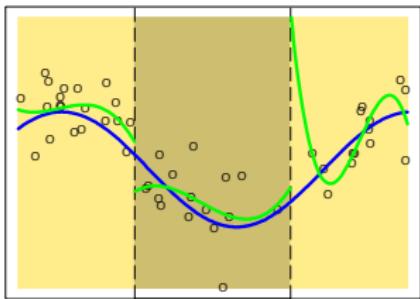
A degree- d spline is a piecewise degree- d polynomial, with continuity in derivatives up to degree $d - 1$ at each knot.

- In general, a degree- d spline with $M - 1$ knots has $d + M$ degrees of freedom²⁶.

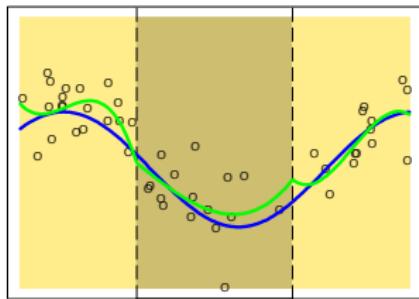
²⁶For example, a linear spline has $1 + M$ degrees of freedom (see (25)). A cubic spline has $3 + M$ degrees of freedom (see (28)). In comparison, a degree- d polynomial has $d + 1$ degrees of freedom.

Piecewise Cubic Polynomials

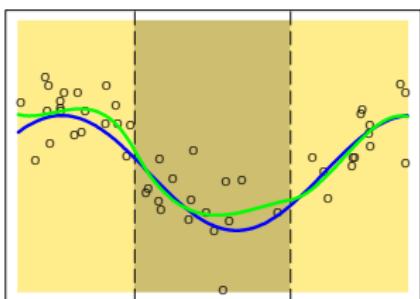
Discontinuous



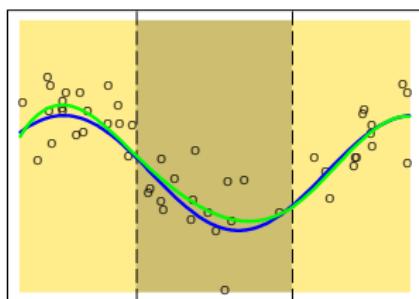
Continuous



Continuous First Derivative



Continuous Second Derivative



ξ_1

ξ_2

ξ_1

ξ_2

Natural Splines

- Splines tend to have high variance at the boundary ($x < \xi_1$ or $x \geq \xi_{M-1}$, where $M - 1$ is the total number of knots).
- A **natural spline** is a regression spline with additional boundary constraints: the function is required to be *linear* beyond the boundary knots, in order to produce more stable estimates.

Wage Profile

```
require(splines)

# Cubic Spline
# -----
# bs() generates B-spline basis functions with specified degrees
#      of polynomials and knots. Here, knots at age 25,40,60
fit.cs <- lm(wage ~ bs(age,knots=c(25,40,60),degree=3), data=Wage)

# Natural Cubic Spline
# -----
# ns() fits a natural cubic spline
fit.ncs <- lm(wage ~ ns(age,knots=c(25,40,60)))
```

Wage Profile

```
coeftest(fit.cs)

##
## t test of coefficients:
##
##                               Estimate Std. Error t value
## (Intercept)                  60.4937   9.4604  6.3944
## bs(age, knots = c(25, 40, 60), degree = 3)1    3.9805  12.5376  0.3175
## bs(age, knots = c(25, 40, 60), degree = 3)2   44.6310   9.6263  4.6364
## bs(age, knots = c(25, 40, 60), degree = 3)3   62.8388  10.7552  5.8426
## bs(age, knots = c(25, 40, 60), degree = 3)4   55.9908  10.7063  5.2297
## bs(age, knots = c(25, 40, 60), degree = 3)5   50.6881  14.4018  3.5196
## bs(age, knots = c(25, 40, 60), degree = 3)6   16.6061  19.1264  0.8682
##                               Pr(>|t|)
## (Intercept)                1.863e-10 ***
## bs(age, knots = c(25, 40, 60), degree = 3)1  0.7508987
## bs(age, knots = c(25, 40, 60), degree = 3)2  3.698e-06 ***
## bs(age, knots = c(25, 40, 60), degree = 3)3  5.691e-09 ***
## bs(age, knots = c(25, 40, 60), degree = 3)4  1.815e-07 ***
## bs(age, knots = c(25, 40, 60), degree = 3)5  0.0004387 ***
## bs(age, knots = c(25, 40, 60), degree = 3)6  0.3853380
```

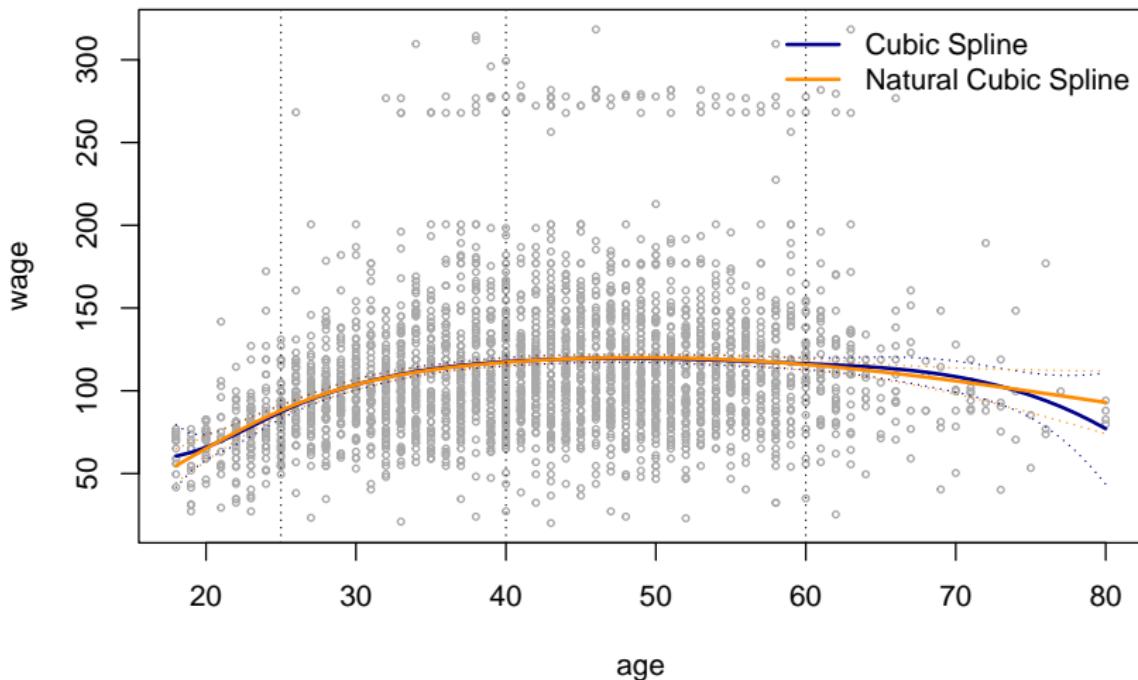
Wage Profile

```
coeftest(fit.ncs)

##
## t test of coefficients:
##
##                               Estimate Std. Error t value Pr(>|t|) 
## (Intercept)                  54.7595   5.1378 10.6581 < 2.2e-16 ***
## ns(age, knots = c(25, 40, 60))1 67.4019   5.0134 13.4442 < 2.2e-16 ***
## ns(age, knots = c(25, 40, 60))2 51.3828   5.7115  8.9964 < 2.2e-16 ***
## ns(age, knots = c(25, 40, 60))3 88.5661  12.0156  7.3709 2.181e-13 ***
## ns(age, knots = c(25, 40, 60))4 10.6369   9.8332  1.0817   0.2795
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Wage Profile

Cubic and Natural Cubic Spline



Generalized Additive Models

So far we have been dealing with a single input x in our discussion of polynomial regression and regression splines. A natural way to extend this discussion to multiple inputs is to assume the following model:

$$y = \omega_0 + \omega_1(x_1) + \omega_2(x_2) + \cdots + \omega_p(x_p) + e \quad (29)$$

, where

$$\omega_j(x_j) = \sum_{m=1}^{M_j} \beta_{jm} \phi_{jm}(x_j)$$

(29) is called a **generalized additive model (GAM)**.

Generalized Additive Models

The GAM allows for flexible nonlinear relationships in each dimension of the input space while maintaining the additive structure of linear models.

- For example, we can fit a linear relationship in x_1 , a polynomial in x_2 , a cubic spline in x_3 , etc.
- The GAM remains a linear basis function model and therefore can be fit by least squares²⁷.

²⁷(29) is equivalent to

$$y = \beta' \Phi(x) + e$$

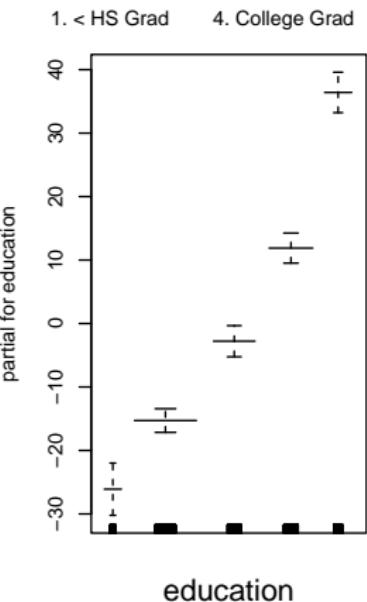
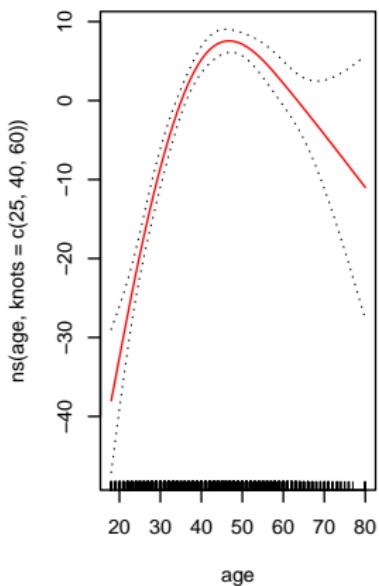
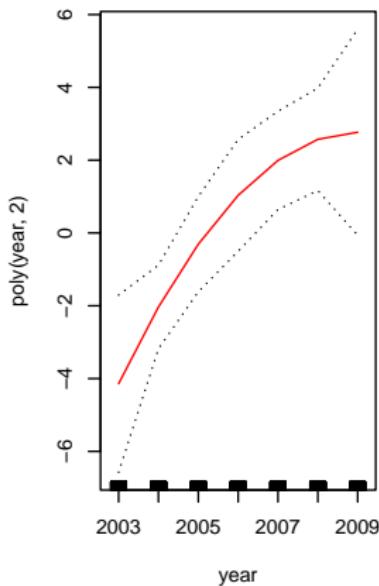
, where $\Phi = (\mathbf{1}, \phi_{11}, \dots, \phi_{1M_1}, \dots, \phi_{p1}, \dots, \phi_{pM_p})'$.

Wage Profile

```
fit <- lm(wage ~ poly(year,2) + ns(age,knots=c(25,40,60)) + education)
coeftest(fit)

##
## t test of coefficients:
##
##                               Estimate Std. Error t value Pr(>|t|) 
## (Intercept)                  47.5751   4.8992  9.7108 < 2.2e-16 ***
## poly(year, 2)1                130.4942  35.2930  3.6974 0.0002217 ***
## poly(year, 2)2                 -36.3005  35.2579 -1.0296 0.3032959  
## ns(age, knots = c(25, 40, 60))1  51.1072   4.4572 11.4662 < 2.2e-16 ***
## ns(age, knots = c(25, 40, 60))2  33.1989   5.0767  6.5394 7.237e-11 ***
## ns(age, knots = c(25, 40, 60))3  53.5004  10.6621  5.0178 5.532e-07 ***
## ns(age, knots = c(25, 40, 60))4  12.3733   8.6866  1.4244 0.1544320  
## education2. HS Grad           10.8174   2.4305  4.4507 8.871e-06 ***
## education3. Some College      23.3191   2.5626  9.0997 < 2.2e-16 ***
## education4. College Grad     37.9867   2.5464 14.9176 < 2.2e-16 ***
## education5. Advanced Degree  62.5184   2.7629 22.6275 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Wage Profile



A GAM model of wage with a quadratic polynomial in year, a natural cubic spline in age, and a step function in education

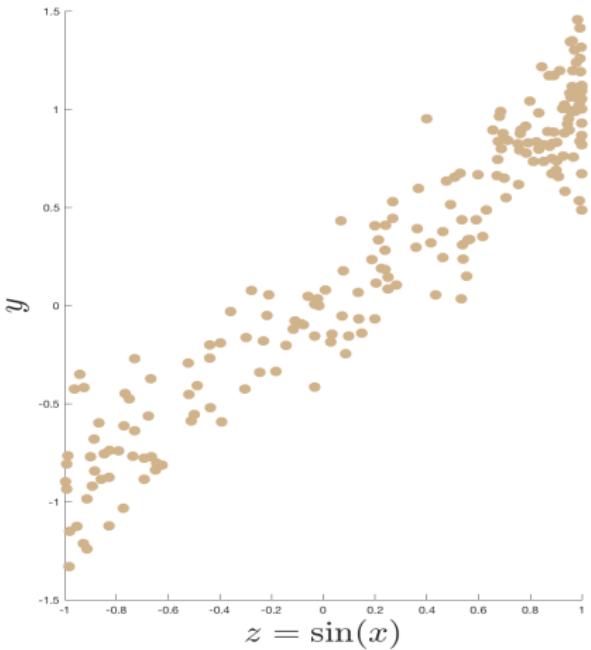
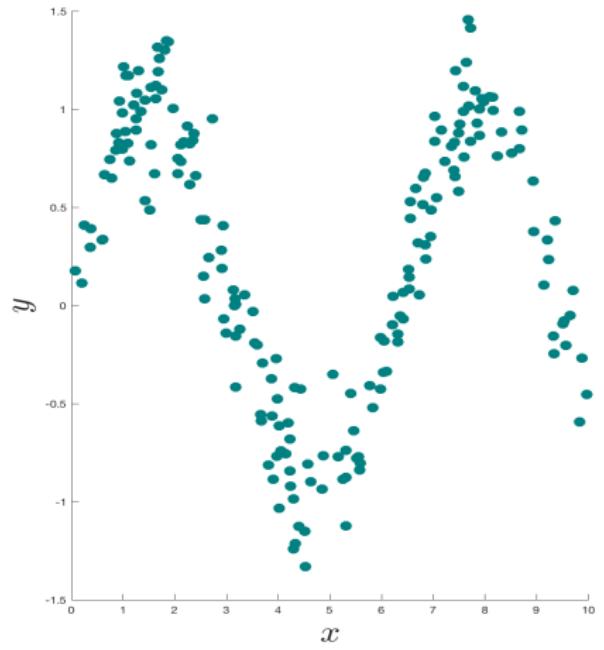
Generalization Issues

Fitting a linear basis function model (22) can be thought of as a two-step process:

- ① Transform x into $\Phi(x)$ ²⁸.
 - ▶ Let $z = \Phi(x) \in \mathcal{Z}$. $\Phi : \mathcal{X} \rightarrow \mathcal{Z}$ is called a **feature transform**.
- ② Fit the linear model: $\mathcal{H}_\Phi = \{h : h(z) = \beta' z\}$, where \mathcal{H}_Φ denotes the hypothesis set corresponding to the feature transform Φ .

²⁸ x can be multi-dimensional: $x = (x_1, \dots, x_p)$

Feature Transform



Left: data in \mathcal{X} -space; Right: data in \mathcal{Z} -space

Generalization Issues

If we decide on the feature transform Φ *before* seeing the data, then the VC generalization bound holds with $d_{VC}(\mathcal{H}_\Phi)$ as the VC dimension.

I.e., for any $g \in \mathcal{H}_\Phi$, with probability at least $1 - \delta$,

$$\begin{aligned} E_{out}(g) &\leq E_{in}(g) + \sqrt{\frac{8}{N} \ln \frac{4((2N)^{d_{VC}} + 1)}{\delta}} \\ &= E_{in}(g) + \mathcal{O}\left(\sqrt{\frac{d_{VC}}{N} \ln N}\right) \end{aligned} \tag{30}$$

, where $d_{VC} = d_{VC}(\mathcal{H}_\Phi)$.

Generalization Issues

Therefore, when choosing a high-order polynomial, or a spline with many degrees of freedom, or a GAM with complex nonlinearities in many dimensions, we cannot avoid the approximation-generalization tradeoff:

- More complex \mathcal{H}_Φ ($d_{VC}(\mathcal{H}_\Phi) \uparrow$) $\Rightarrow E_{in} \downarrow$
- Less complex \mathcal{H}_Φ ($d_{VC}(\mathcal{H}_\Phi) \downarrow$) $\Rightarrow |E_{out} - E_{in}| \downarrow$

Generalization Issues

What if we try a transformation Φ_1 first, and then, finding the results unsatisfactory, decide to use Φ_2 ? Then we are effectively using a model that contains both $\{\beta' \Phi_1(x)\}$ and $\{\beta' \Phi_2(x)\}$.

- For example, if we try a linear model first, then a quadratic polynomial, then a piecewise constant model, before settling on a cubic spline, then d_{VC} in (30) should be the VC dimension of a hypothesis set that contains not only the cubic spline model, but all of the aforementioned models.
- The process of trying a series of models until we get a satisfactory result is called **specification search** or **data snooping**. In general, the more models you try, the poorer your final result will generalize out of sample.

Appendix: Frisch-Waugh-Lovell Theorem

Proof.

Define $\mathbb{M} \doteq \mathbb{I} - \mathbb{H} = \mathbb{I} - X(X'X)^{-1}X'$. Then for any Y , we can write

$$Y = X\hat{\beta} + \hat{e} = \mathbb{H}Y + \mathbb{M}Y$$

, where $\hat{e} = Y - X(X'X)^{-1}X'Y = (\mathbb{I} - \mathbb{H})Y = \mathbb{M}Y$ is orthogonal to X .

Now consider the model

$$Y = X_1\beta_1 + X_2\beta_2 + e$$

OLS estimation \Rightarrow

$$Y = X_1\hat{\beta} + X_2\hat{\beta} + \hat{e}$$



Appendix: Frisch-Waugh-Lovell Theorem

Proof. (cont.)

Define $\mathbb{M}_{X_2} \doteq \mathbb{I} - X_2 (X_2' X_2)^{-1} X_2'$. Then

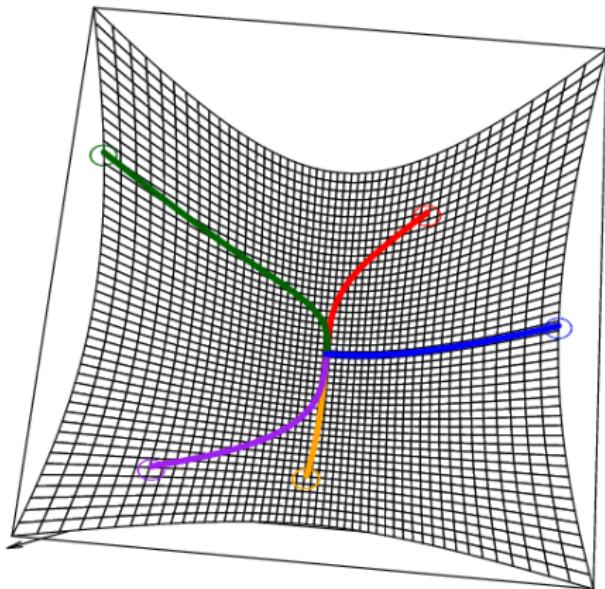
$$\begin{aligned}\mathbb{M}_{X_2} Y &= \mathbb{M}_{X_2} X_1 \hat{\beta}_1 + \mathbb{M}_{X_2} X_2 \hat{\beta}_2 + \mathbb{M}_{X_2} \hat{e} \\ &\stackrel{[1]}{=} \mathbb{M}_{X_2} X_1 \hat{\beta}_1 + \hat{e}\end{aligned}\tag{31}$$

, where [1] follows since $\mathbb{M}_{X_2} X_2 = 0$ and $\mathbb{M}_{X_2} \hat{e} = \hat{e}$ because $\hat{e} \perp X_2$.

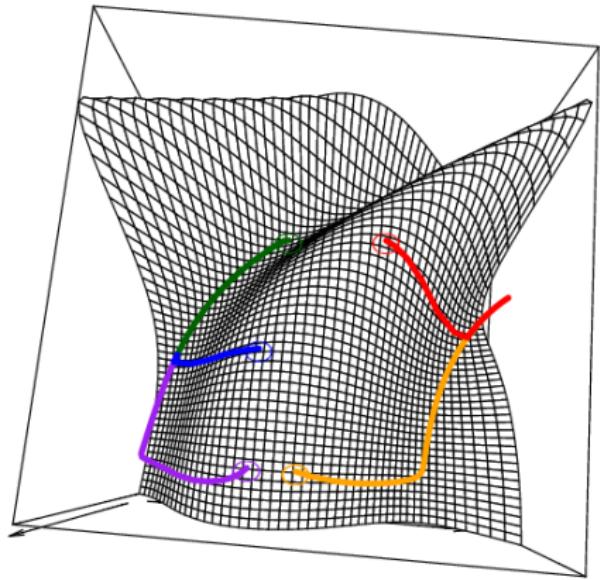
Finally, notice that because $\hat{e} \perp X_1$, we also have $\hat{e} \perp \mathbb{M}_{X_2} X_1$. Hence (31) represents an orthogonal projection of $\mathbb{M}_{X_2} Y$ on $\mathcal{R}(\mathbb{M}_{X_2} X_1)$. Therefore, OLS regression of $\mathbb{M}_{X_2} Y$ on $\mathbb{M}_{X_2} X_1$ will produce the same $\hat{\beta}_1$ as in (31).



Appendix: Gradient Descent



(a) Convex function



(b) Non-convex function

Gradient descent paths with different starting points

Acknowledgement I

Part of this lecture is based on the following sources:

- Gramacy, R. B. *Applied Regression Analysis*. Lecture at the University of Chicago Booth School of Business, retrieved on 2017.01.01. [[link](#)]
- Hastie, T., R. Tibshirani, and J. Friedman. 2008. *The Elements of Statistical Learning* (2nd ed.). Springer.
- James, G., D. Witten, T. Hastie, and R. Tibshirani. 2013. *An Introduction to Statistical Learning: with Applications in R*. Springer.
- Penn State University. *Probability Theory and Mathematical Statistics*. Online course, retrieved on 2017.01.01. [[link](#)]
- Schmidt, M. *Machine Learning*. Lecture at the University of British Columbia, retrieved on 2021.01.01. [[link](#)]
- Shalizi, C. R. 2019. *Advanced Data Analysis from an Elementary Point of View*. Manuscript.

Acknowledgement II

- Taddy, M. *Big Data*. Lecture at the University of Chicago Booth School of Business, retrieved on 2017.01.01. [[link](#)]
- Tibshirani, R. *Convex Optimization*. Lecture at Carnegie Mellon University, retrieved on 2021.01.01. [[link](#)]