# Linear Regression
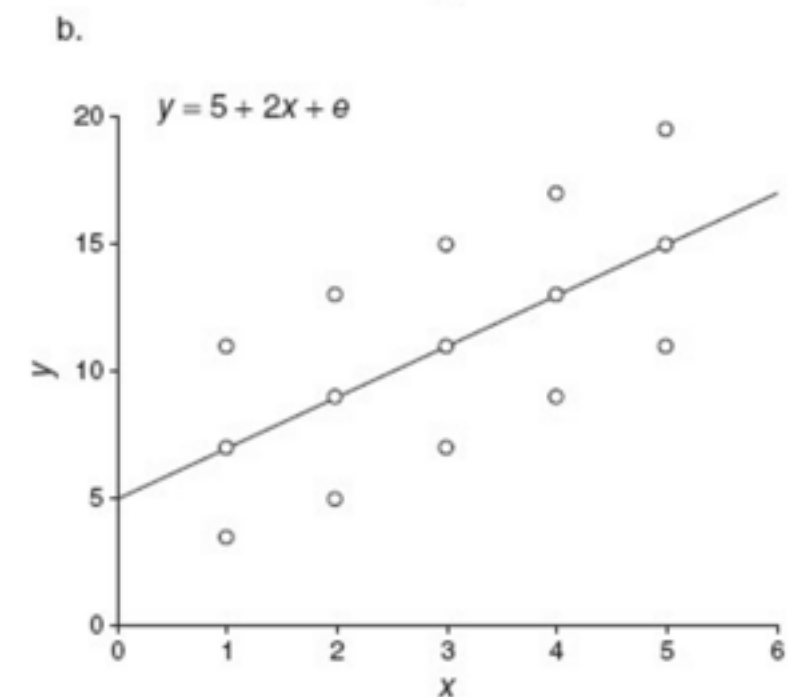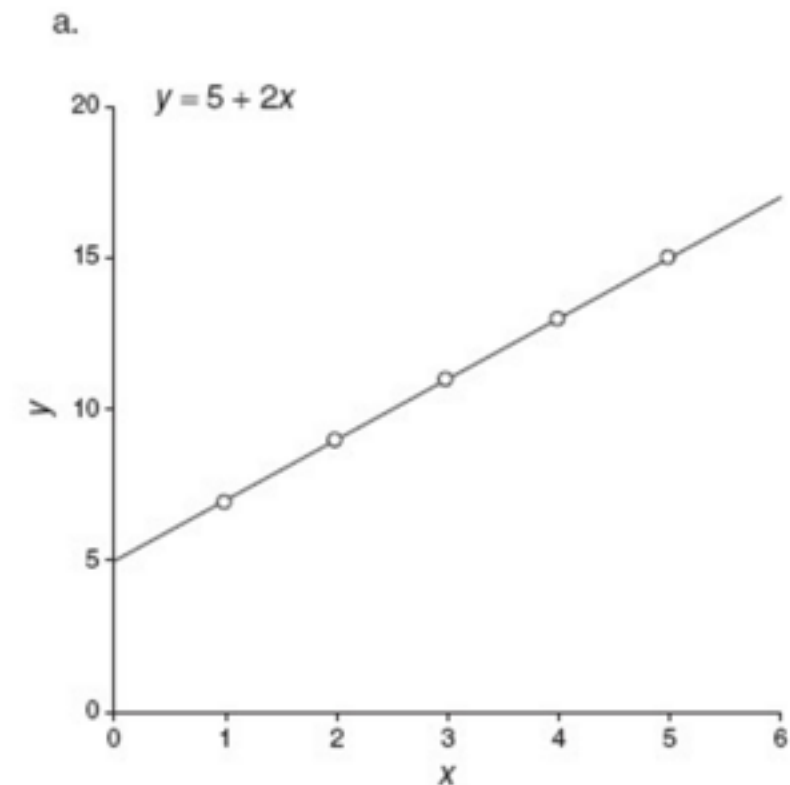
Fitting a line to data

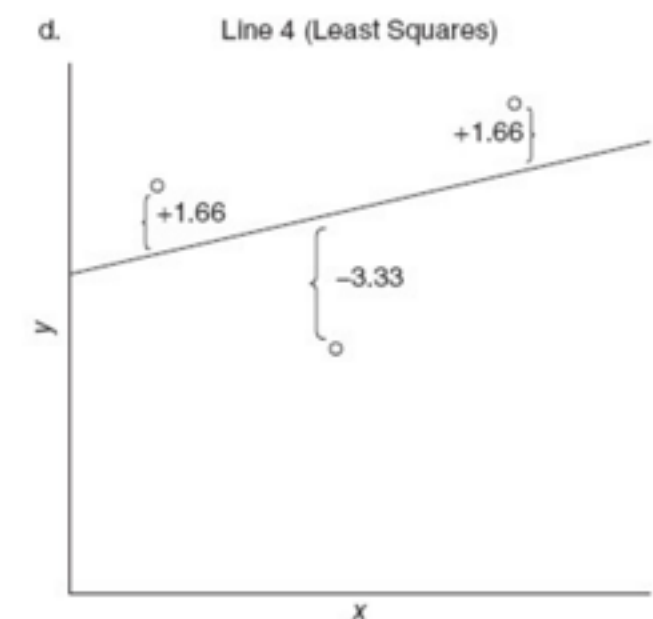Darren Reger Lecture for Galvanize DSI
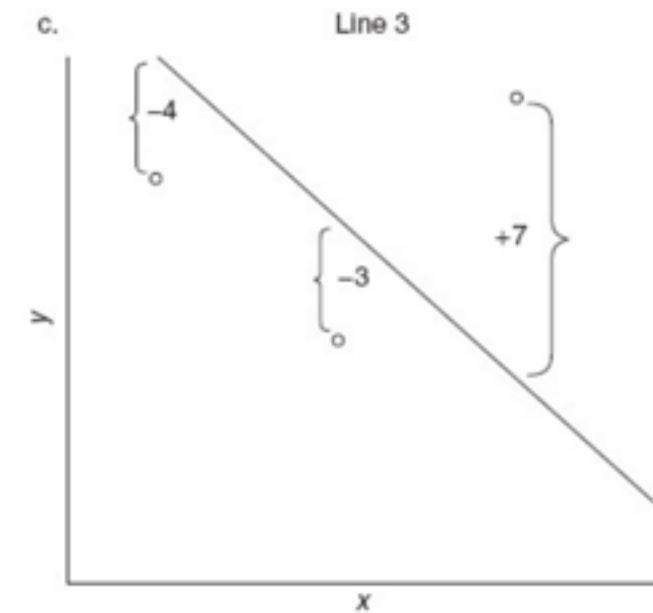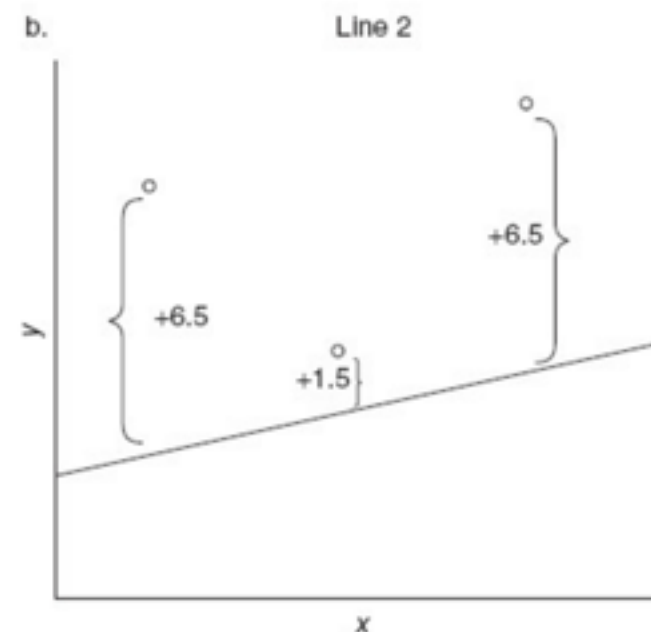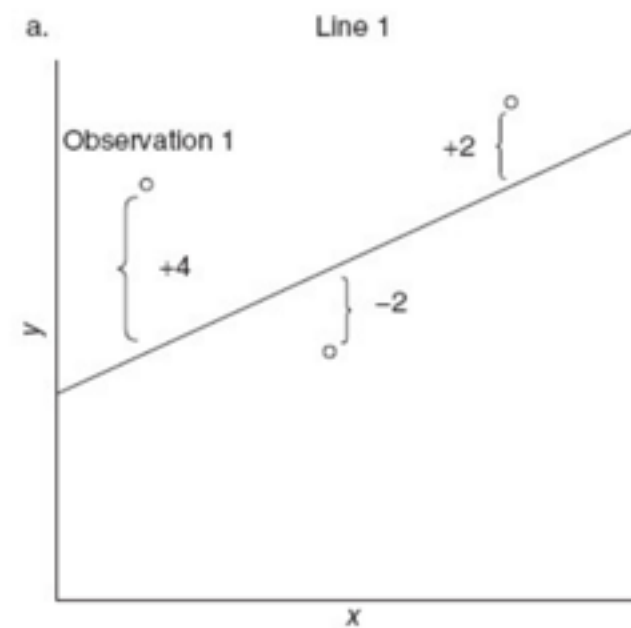
# Relationship Between X & Y

- Linear relationships

  - Exact vs. Inexact

- Why inexact?

a.

$y = 5 + 2x$

b.

$y = 5 + 2x + e$

# Line Placement

- Why linear regression?

- Where to place the line?

- Why OLS?

# Simple Linear Regression

$$Y = \beta_0 + \beta_1 X + \epsilon$$



- The Model, what you're presuming the world looks like

- $\beta_0$ and $\beta_1$ are unknown constants that represent the intercept and slope.

- $\epsilon$ is the error term. $\epsilon \sim$ i.i.d. N(0, σ^2)

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

- $\beta_0$-hat and $\beta_1$-hat are model coefficient estimates for world presumed

- y-hat indicates the prediction of Y based on X=x

# Multiple Linear Regression

**Model in Matrix Form**

$$\mathbf{Y}_{n\times1} = \mathbf{X}_{n\times p}\beta_{p\times1} + \epsilon_{n\times1}$$
$$\epsilon \sim N(\mathbf{0}, \sigma^2\mathbf{I}_{n\times n})$$
$$\boxed{\mathbf{Y} \sim N(\mathbf{X}\beta, \sigma^2\mathbf{I})}$$

Design Matrix $\mathbf{X}$:

Target:

$$\mathbf{X} = \begin{bmatrix} 1 & X_{1,1} & X_{1,2} & \cdots & X_{1,p-1} \\ 1 & X_{2,1} & X_{2,2} & \cdots & X_{2,p-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n,1} & X_{n,2} & \cdots & X_{n,p-1} \end{bmatrix}$$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

Coefficient matrix $\beta$:

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{p-1} \end{bmatrix}$$

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^{\mathrm{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathrm{T}}\mathbf{y}$$

# Assessing Accuracy

## Residual Sum of Squares

$$RSS = \sum_{i=1}^{n}(y_i - \hat{y}_i)^2$$

Not great...

This is also what we use to estimate $\sigma = \sqrt{Var(\epsilon)}$

## Residual Standard Error

$$RSE = \sqrt{\frac{1}{n-p-1}RSS} = \sqrt{\frac{(y_i - \hat{y}_i)^2}{n-p-1}}$$

Better...can roughly think of as average amount that response will deviate from regression line



a. **Sample 1 (tight fit)**

$$\hat{y}_1 = b_{01} + b_{11}x_1$$

b. **Sample 2 (loose fit)**

$$\hat{y}_2 = b_{02} + b_{12}x_2$$

# R-squared

# Comparing Models

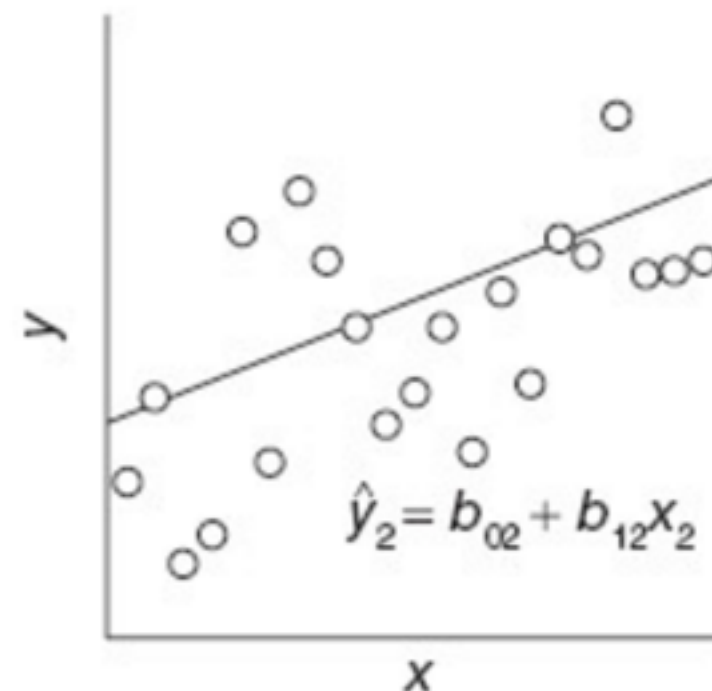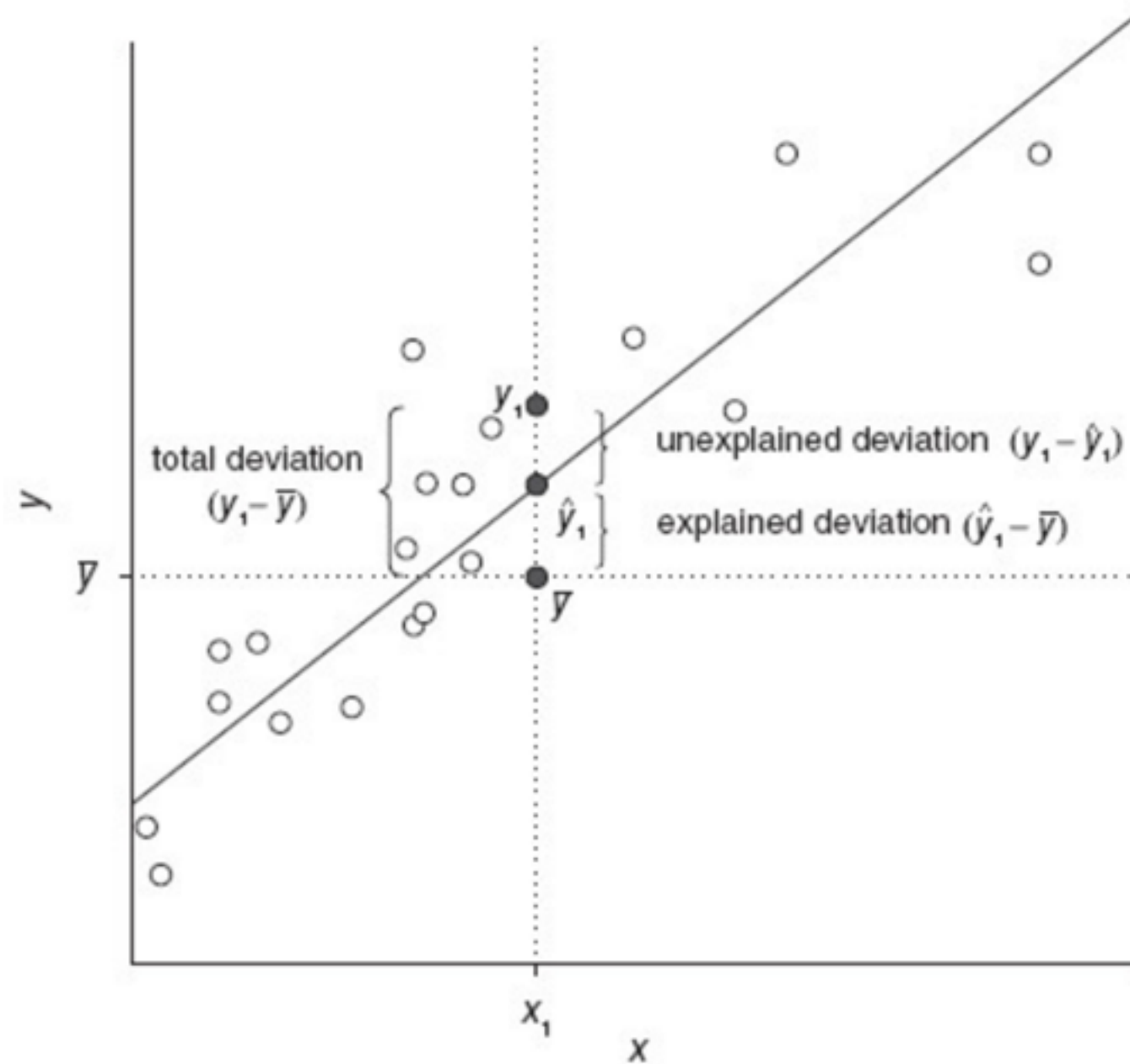(1) Set up comparison

**m_reduced:** $Y = \beta_0 + \beta_{weight} + \beta_{modelyear} + \beta_{cartype}$

**m_full:** $Y = \beta_0 + \beta_{weight} + \beta_{height} + \beta_{color} + \beta_{modelyear} + \beta_{cartype}$

(2) Compute F-statistic

$$F = \frac{(RSS_{reduced} - RSS_{full})/(p_{full} - p_{reduced})}{RSS_{full}/(n - p_{full} - 1)}$$

where F has degrees of freedom (p_full - p_reduced), (n - p_full – 1)

Notice that if *height* and *color* really don't matter much…
(RSS_reduced - RSS_full) will be small → F-statistic will be small

(3) Compute p-value



p-value = 0.1241

2.23

```
from scipy.stats import f
p_val = 1-f.cdf(calculated_F,
        p_full - p_reduced,
        n - p_full - 1)
```

Assuming α=0.05,
- if p < 0.05 reject null (that height and color don't matter)
- If p >= 0.05, fail to reject null (that height and color don't matter)

# Comparing Models

- F-test can be used super generally

- Two special use cases

① Is my model useful at all? i.e. Is at least one of my predictors X_1, X_2, ... X_p useful in predicting the response?

$$H_0 : \beta_1 = \beta_2 = ... = \beta_p = 0$$
$$H_A : at\ least\ one\ \beta_j\ is\ non-zero \longrightarrow F = \frac{(\text{TSS} - \text{RSS})/p}{\text{RSS}/(n - p - 1)} \sim F_{p,n-p-1}$$

① Equivalence to t-test in the Regression Output!

m_reduced: $Y = \beta_0 + \beta_{weight} + \beta_{height} + \beta_{color} + \beta_{cartype}$

m_full: $Y = \beta_0 + \beta_{weight} + \beta_{height} + \beta_{color} + \beta_{modelyear} + \beta_{cartype}$

$\longrightarrow$ Results in p-value associated with $\beta_{modelyear}$

# Interpreting Coefficients

| | Recall | Here |
|---|---|---|
| **Setup Hypothesis** | $H_0$: $\mu = 100$ | $H_0 : \beta_1 = 0$ |
| **Sample Statistic** | $\overline{x}$ | $\hat{\beta}_1$ |
| **Test Statistic** | $t = \dfrac{\overline{x} - \mu_0}{s/\sqrt{n}}$ | $t = \dfrac{\hat{\beta}_1 - 0}{\mathrm{SE}(\hat{\beta}_1)}$ |
| **Confidence Interval** | $(\overline{X} - t_{\alpha/2} \dfrac{S}{\sqrt{n}}, \overline{X} + t_{\alpha/2} \dfrac{S}{\sqrt{n}})$ | $\left[ \hat{\beta}_1 - 2 \cdot \mathrm{SE}(\hat{\beta}_1), \ \hat{\beta}_1 + 2 \cdot \mathrm{SE}(\hat{\beta}_1) \right]$ |

Test if X has effect on Y

# Assumptions

- Linearity
- Constant variance (homoscedasticity)
- Independence of errors
- Normality of errors
- Lack of multicollinearity

# Residual Plots

# Leverage

- Leverage point: an observation with an unusual X value

- Does not necessarily have a large effect on the regression model

- Most common measure, the hat value, $h_{ii} = (H)_{ii}$

- The $i$th diagonal of the hat matrix

$$H = X(X^T X)^{-1} X^T$$

# Studentized Residuals

$$H = X(X^T X)^{-1} X^T.$$

The **leverage** $h_{ii}$ is the $i$th diagonal entry in the hat matrix. The variance of the $i$th residual is

$$\text{var}(\hat{\varepsilon}_i) = \sigma^2 (1 - h_{ii}).$$

In case the design matrix $X$ has only two columns (as in the example above), this is equal to

$$\text{var}(\hat{\varepsilon}_i) = \sigma^2 \left( 1 - \frac{1}{n} - \frac{(x_i - \bar{x})^2}{\sum_{j=1}^{n}(x_j - \bar{x})^2} \right).$$
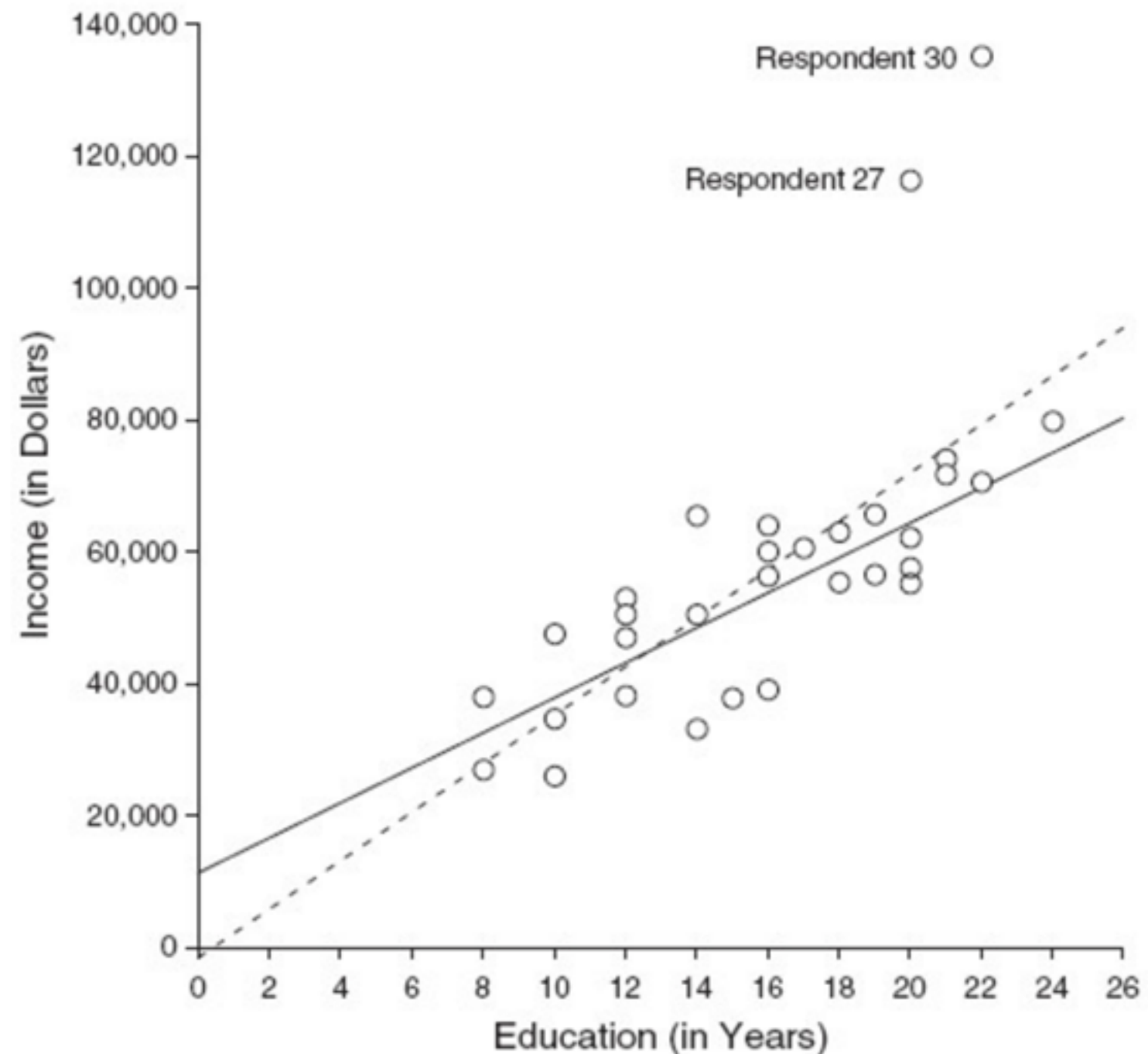
The corresponding **studentized residual** is then

$$t_i = \frac{\hat{\varepsilon}_i}{\hat{\sigma}\sqrt{1 - h_{ii}}}$$

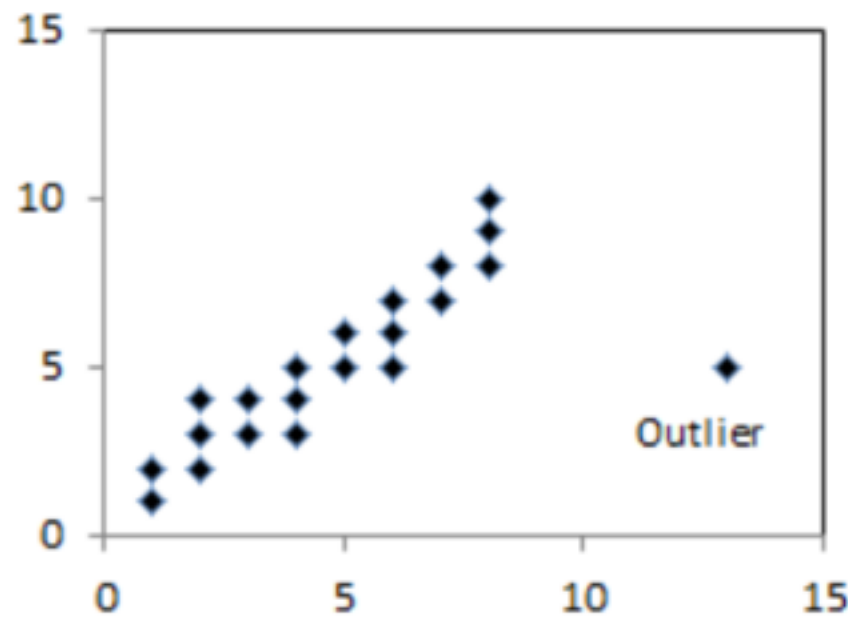where $\hat{\sigma}$ is an appropriate estimate of $\sigma$ (see below).

# Outliers

- Y values very from from our predictions
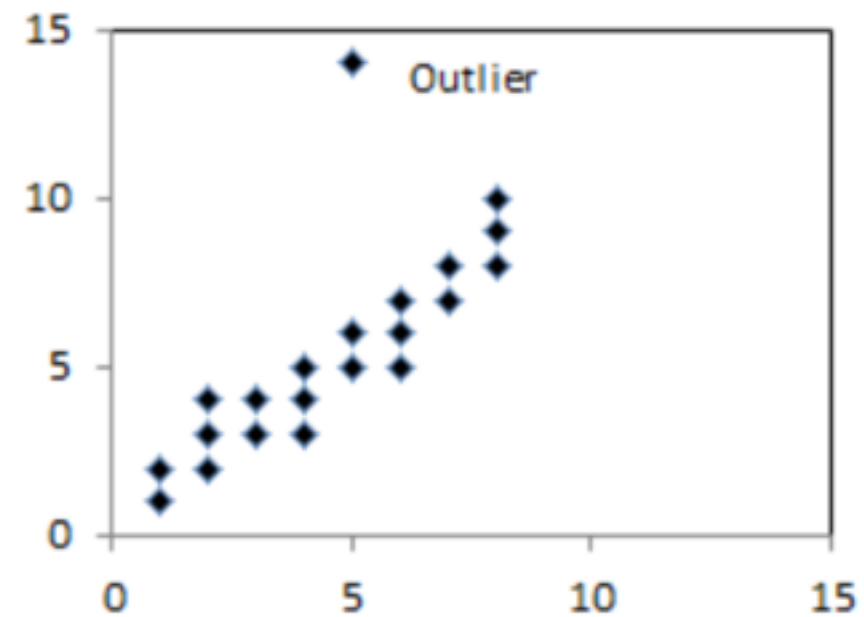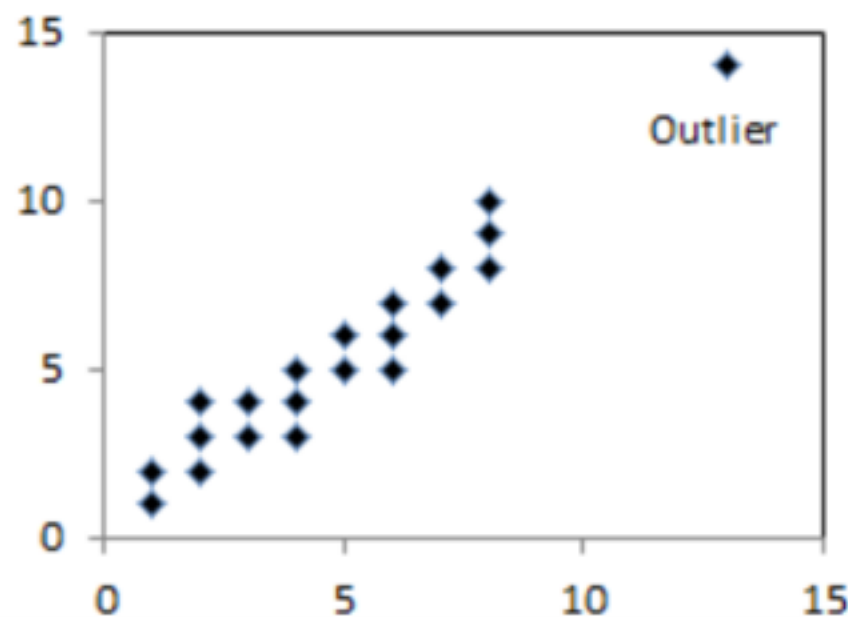
- Reasons they occur

- OLS sensitivity

# Types of Outliers

# Detecting Outliers

- Residual plots can help identify outliers
  - Recall that residuals are $e_i = y_i - \hat{y}_i$

  - and that $\varepsilon \sim$ i.i.d. $N(0, \sigma^2)$

    → "Studentized" residuals: Dividing each residual by its standard error, should result in a "studentized residual" between -3 and 3. Studentized residuals outside this range indicate outliers.

# Multicollinearity

- Perfect multicollinearity

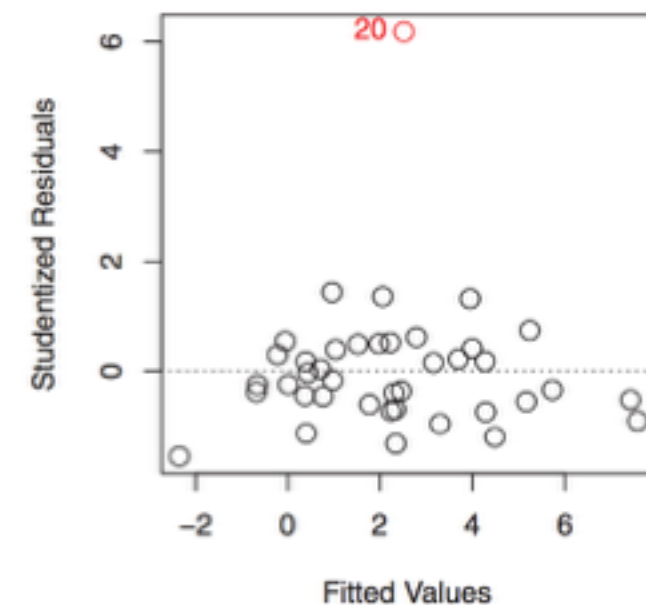    - Easily detectable because your model will fail to run

    - Unlikely to occur in practice, unless you goof

- Partial Multicollinearity

    - Uncertainty in the model becomes large

    - Does not affect model accuracy or bias coefficients

# Multicollinearity

- ## Correlation Matrix / Scatterplot Matrix



Downside is can only pick up pairwise effects ☹

- ## Variance Inflation Factors (VIF)

  – Run ordinary least squares for each predictor as function of all the other predictors.  k times for k predictors

$$X_1 = \alpha_2 X_2 + \alpha_3 X_3 + \cdots + \alpha_k X_k + c_0 + e$$

$$VIF = \frac{1}{1 - R_i^2}$$

Looks at all predictors together! ☺

Rule of Thumb, > 10 is problematic

# QQ Plots

# Normal QQ Plot

- Check out this explanation

- http://emp.byui.edu/BrownD/Stats-intro/dscrptv/graphs/qq-plot_egs.htm

# Break for Morning Sprint

# Categorical Variables

- Interested in **Credit Card Balances** (y)
- Suspect it may be related to *Gender* or *Ethnicity*

Modeling with just *Gender*

$$x_i = \begin{cases} 1 & \text{if } i\text{th person is female} \\ 0 & \text{if } i\text{th person is male} \end{cases}$$

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i = \begin{cases} \beta_0 + \beta_1 + \epsilon_i & \text{if } i\text{th person is female} \\ \beta_0 + \epsilon_i & \text{if } i\text{th person is male.} \end{cases}$$

# Categorical Variables

## Modeling with *Ethnicity* (more than 2 Levels)

$$x_{i1} = \begin{cases} 1 & \text{if } i\text{th person is } \underline{\text{Asian}} \\ 0 & \text{if } i\text{th person is not Asian} \end{cases}$$

$$x_{i2} = \begin{cases} 1 & \text{if } i\text{th person is } \underline{\text{Caucasian}} \\ 0 & \text{if } i\text{th person is not Caucasian} \end{cases}$$

$$y_i = \beta_0 + \beta_1 \underline{x_{i1}} + \beta_2 \underline{x_{i2}} + \epsilon_i = \begin{cases} \beta_0 + \beta_1 + \epsilon_i & \text{if } i\text{th person is Asian} \\ \beta_0 + \beta_2 + \epsilon_i & \text{if } i\text{th person is Caucasian} \\ \beta_0 + \epsilon_i & \text{if } i\text{th person is AA.} \end{cases}$$

**Data**

| Ones | Ethnicity |
|------|-----------|
| 1 | AA |
| 1 | Asian |
| 1 | Asian |
| 1 | Caucasian |
| 1 | AA |
| 1 | AA |
| 1 | Asian |
| 1 | Caucasian |
| 1 | AA |
| ... | ... |

**Recode Design Matrix**

| Ones | Asian | Caucasian |
|------|-------|-----------|
| 1 | 0 | 0 |
| 1 | 1 | 0 |
| 1 | 1 | 0 |
| 1 | 0 | 1 |
| 1 | 0 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 0 |
| 1 | 0 | 1 |
| 1 | 0 | 0 |
| ... | ... | ... |

- $\beta 0$ as average credit card balance for AA

- $\beta 1$ as <u>difference</u> in average balance between Asian and AA

- $\beta 2$ as <u>difference</u> in average balance between Caucasian and AA

So what if $\beta 1 = -23.1$?

# Categorical Variables

Card_Balance ~ Age + Years_of_Education + Gender + Ethnicity + ....

- Intercept β0 loses nice interpretation

- Now what's it mean if β1 = -23.1?

- What if you wanted to compare groups to Caucasians as a baseline?

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \epsilon_i = \begin{cases} \beta_0 + \beta_1 + \epsilon_i & \text{if } i\text{th person is Asian} \\ \beta_0 + \beta_2 + \epsilon_i & \text{if } i\text{th person is Caucasian} \\ \beta_0 + \epsilon_i & \text{if } i\text{th person is AA.} \end{cases}$$

# Categorical Variables

Card_Balance ~ Age + Years_of_Education + Gender + Ethnicity + ….

- Intercept $\beta 0$ loses nice interpretation

- Now what's it mean if $\beta 1 = -23.1$?

  ✓ Still interpret as difference between Asian and AA…*holding all other predictors constant.* Again, beware of interpretation.

- What if you wanted to compare groups to Caucasians as a baseline?

  ✓

| **Data** | | **Recode Design Matrix** | | |
|---|---|---|---|---|
| Ones | Ethnicity | Ones | **AA** | **Asian** |
| 1 | AA | 1 | 1 | 0 |
| 1 | Asian | 1 | 0 | 1 |
| 1 | Asian | 1 | 0 | 1 |
| 1 | Caucasian | 1 | 0 | 0 |
| 1 | AA | 1 | 1 | 0 |
| 1 | AA | 1 | 1 | 0 |
| 1 | Asian | 1 | 0 | 1 |
| 1 | Caucasian | 1 | 0 | 0 |
| 1 | AA | 1 | 1 | 0 |
| … | … | … | 0 | 0 |

# Varying Intercepts

- 2 Formulations

  - Baseline and alternative

  - Individual fit



$x_2 = 1$ ( **Galv** ): $\hat{y} = (b_0 + b_2) + b_1 x_1$

$b_2$

$x_2 = 0$ ( **GA** ): $\hat{y} = b_0 + b_1 x_1$
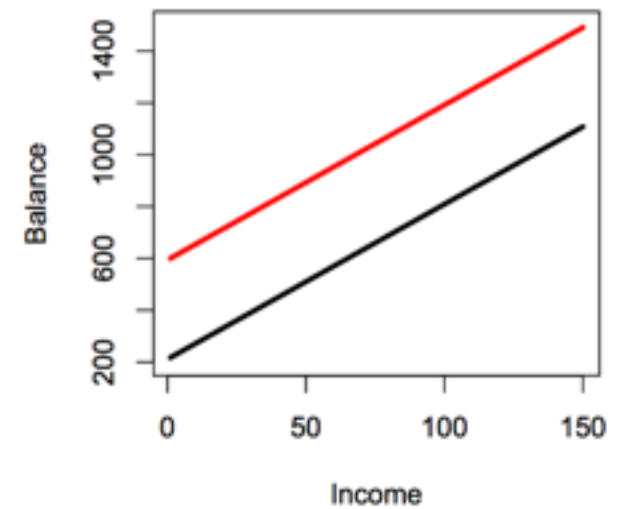
Income ($y$)

Education ($x_1$)

# Interactions

Interacting **student** (qualitative) and **income** (quantitative)

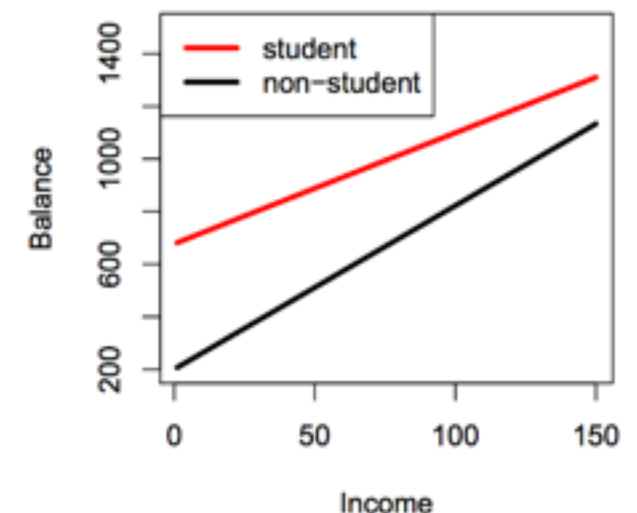## No Interaction

$$balance_i = \beta_0 + \beta_1 * income_i + \beta_2 * student_i$$

$$balance_i \approx \beta_0 + \beta_1 \times income_i + \begin{cases} \beta_2 & \text{if } i\text{th person is a student} \\ 0 & \text{if } i\text{th person is not a student} \end{cases}$$

$$= \beta_1 \times income_i + \begin{cases} \underline{\beta_0 + \beta_2} & \text{if } i\text{th person is a student} \\ \beta_0 & \text{if } i\text{th person is not a student} \end{cases}$$
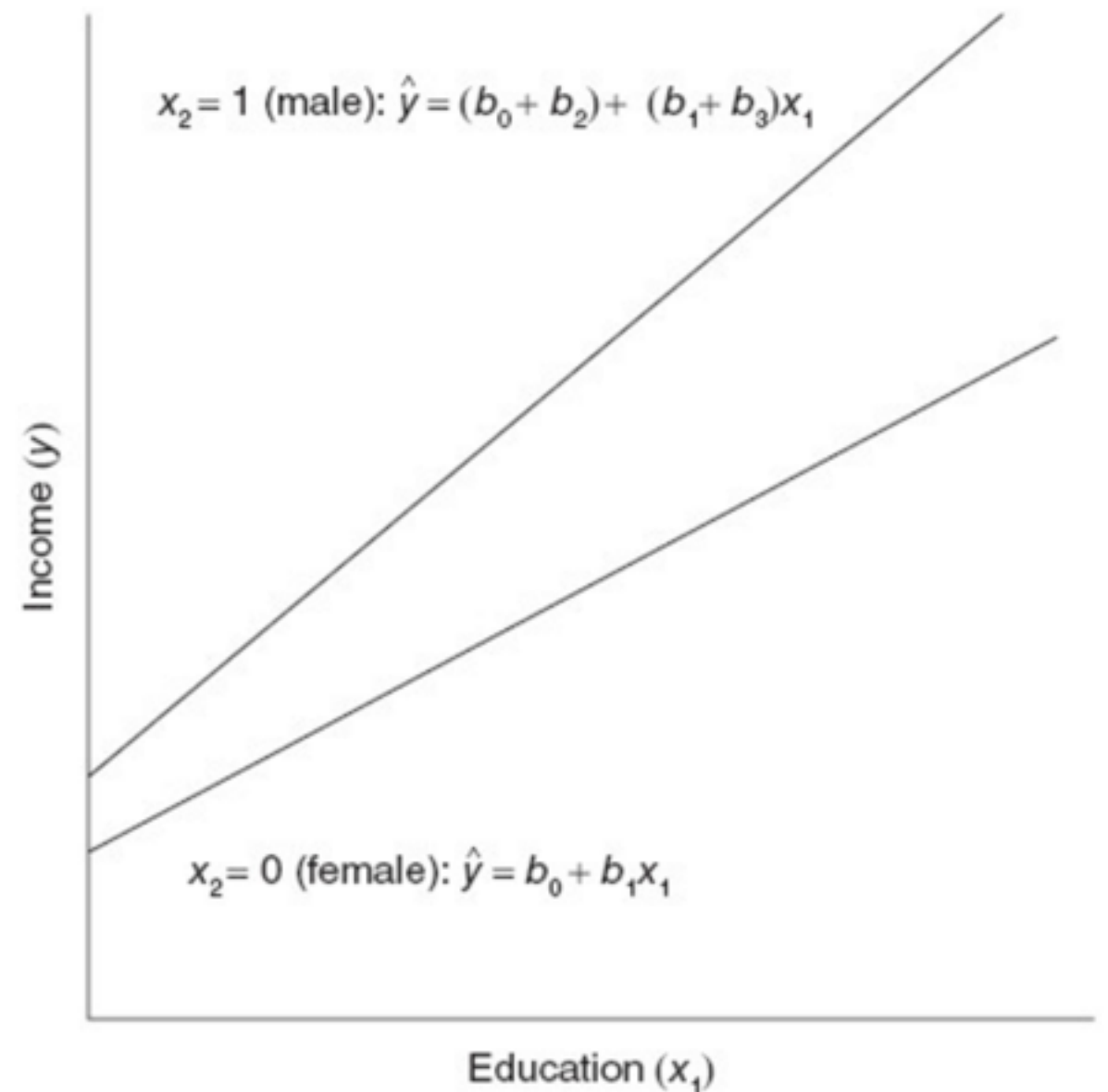


## With Interaction

$$balance_i = \beta_0 + \beta_1 * income_i + \beta_2 * student_i + \beta_3 * income_i * student_i$$

$$balance_i \approx \beta_0 + \beta_1 \times income_i + \begin{cases} \beta_2 + \beta_3 \times income_i & \text{if student} \\ 0 & \text{if not student} \end{cases}$$

$$= \begin{cases} (\underline{\beta_0 + \beta_2}) + (\beta_1 + \boxed{\beta_3}) \times income_i & \text{if student} \\ \beta_0 + \beta_1 \times income_i & \text{if not student} \end{cases}$$

# Varying Slopes

- 2 Formulations

  - Baseline and alternative

  - Individual fit



$x_2 = 1$ (male): $\hat{y} = (b_0 + b_2) + (b_1 + b_3)x_1$

$x_2 = 0$ (female): $\hat{y} = b_0 + b_1 x_1$

Income ($y$)

Education ($x_1$)

# Interactions

$$\begin{aligned}
\texttt{sales} &= \beta_0 + \beta_1 \times \texttt{TV} + \beta_2 \times \texttt{radio} + \underline{\beta_3 \times (\texttt{radio} \times \texttt{TV})} + \epsilon \\
&= \beta_0 + \underline{(\beta_1 + \beta_3 \times \texttt{radio})} \times \texttt{TV} + \beta_2 \times \texttt{radio} + \epsilon.
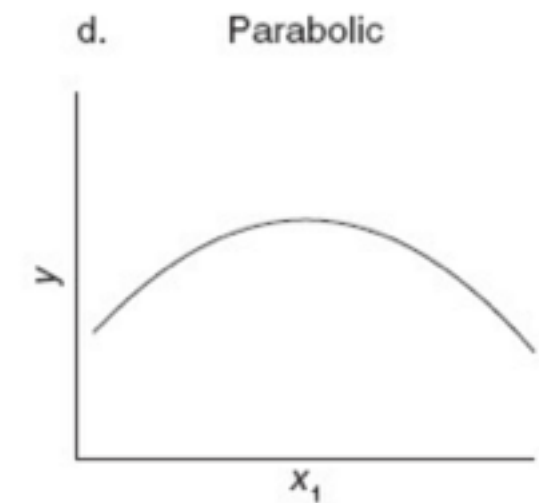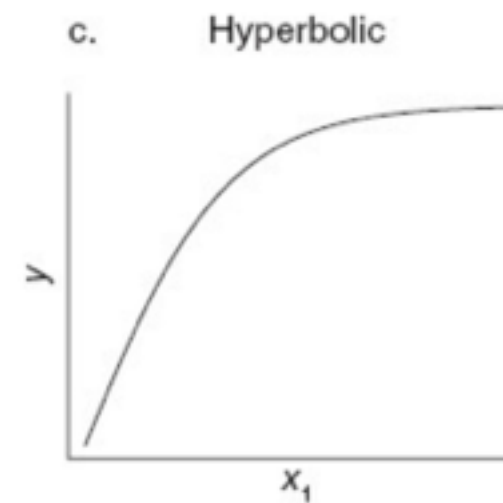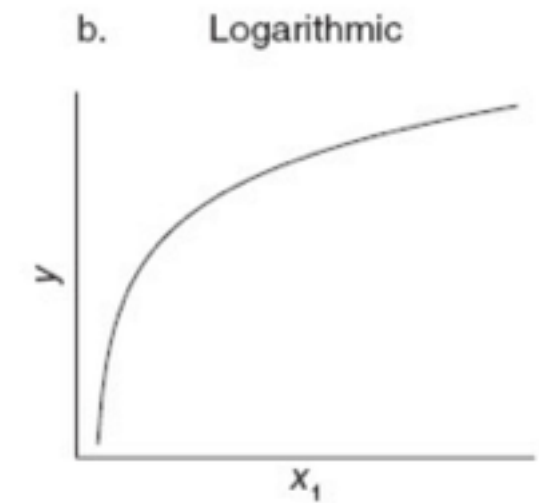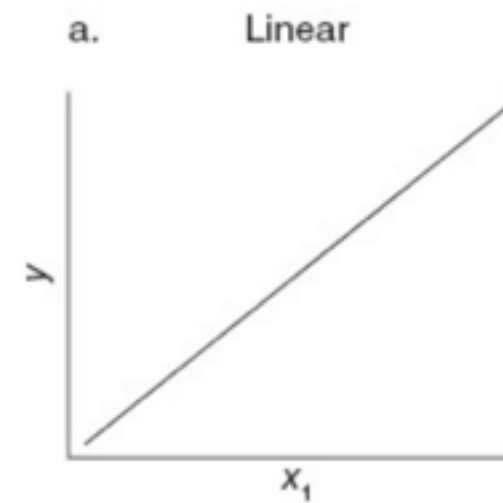\end{aligned}$$

Results:

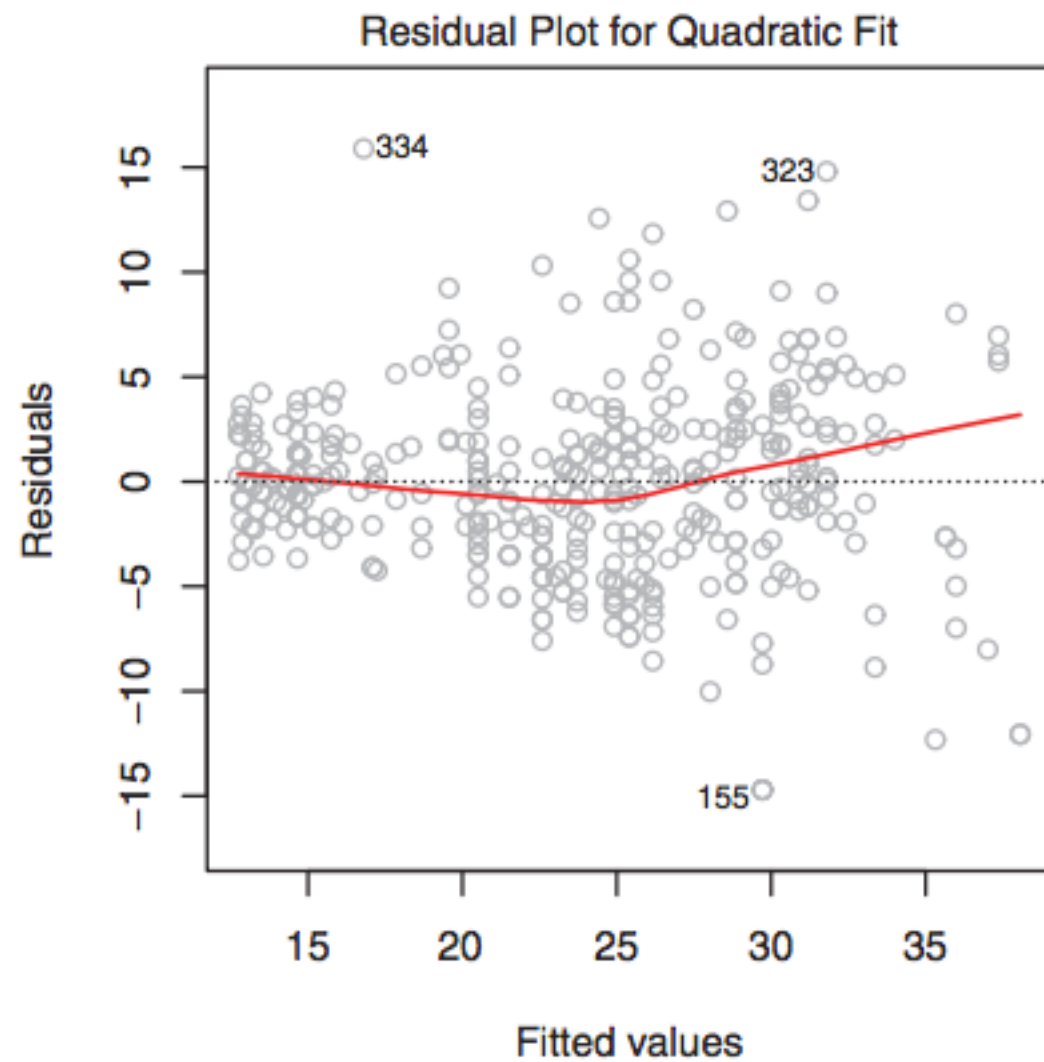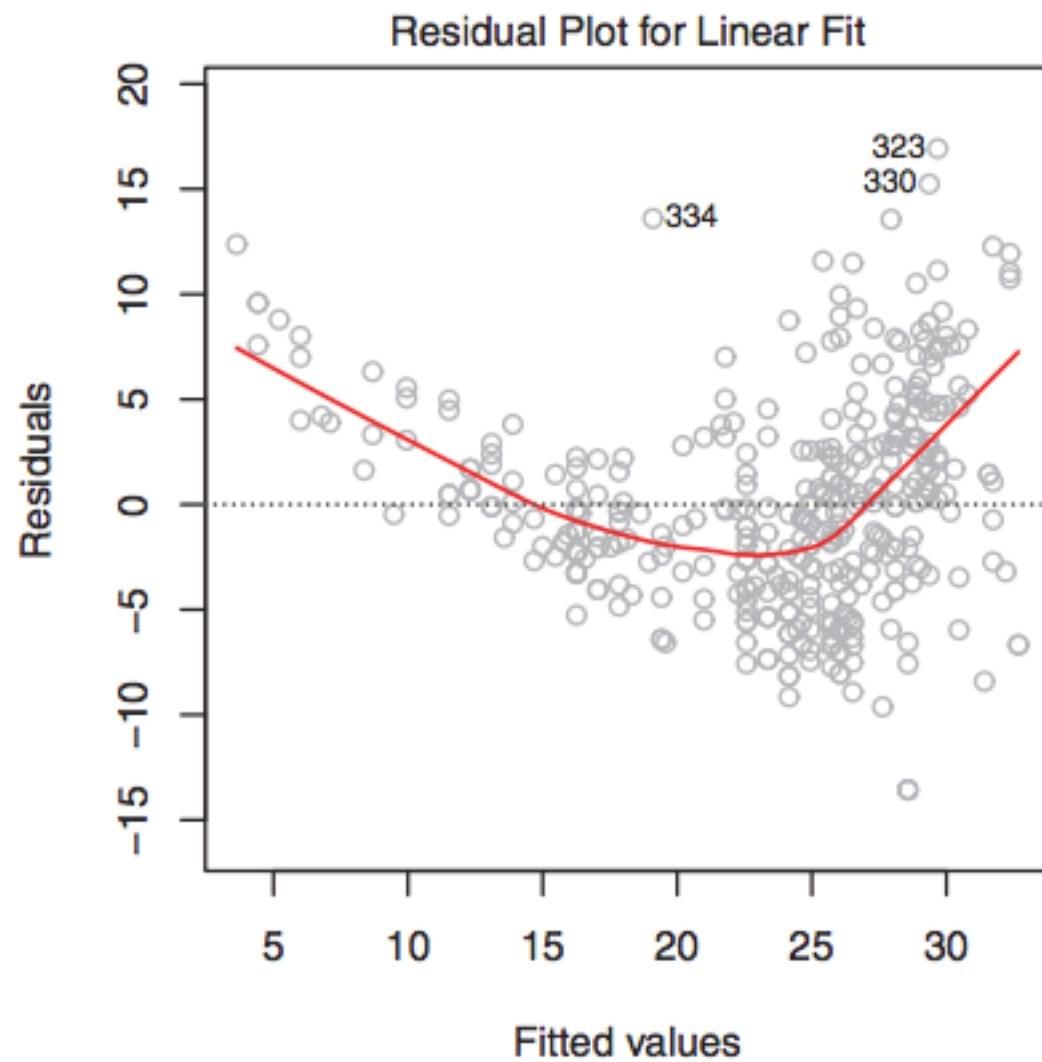| | Coefficient | Std. Error | t-statistic | p-value |
|---|---|---|---|---|
| Intercept | 6.7502 | 0.248 | 27.23 | $< 0.0001$ |
| TV | 0.0191 | 0.002 | 12.70 | $< 0.0001$ |
| radio | 0.0289 | 0.009 | 3.24 | 0.0014 |
| TV×radio | 0.0011 | 0.000 | 20.73 | $< 0.0001$ |

← Improvement!

The coefficient estimates in the table suggest that an increase in TV advertising of $\$1,000$ is associated with increased sales of
$$(\hat{\beta}_1 + \hat{\beta}_3 \times \texttt{radio}) \times 1000 = 19 + 1.1 \times \texttt{radio} \text{ units}.$$
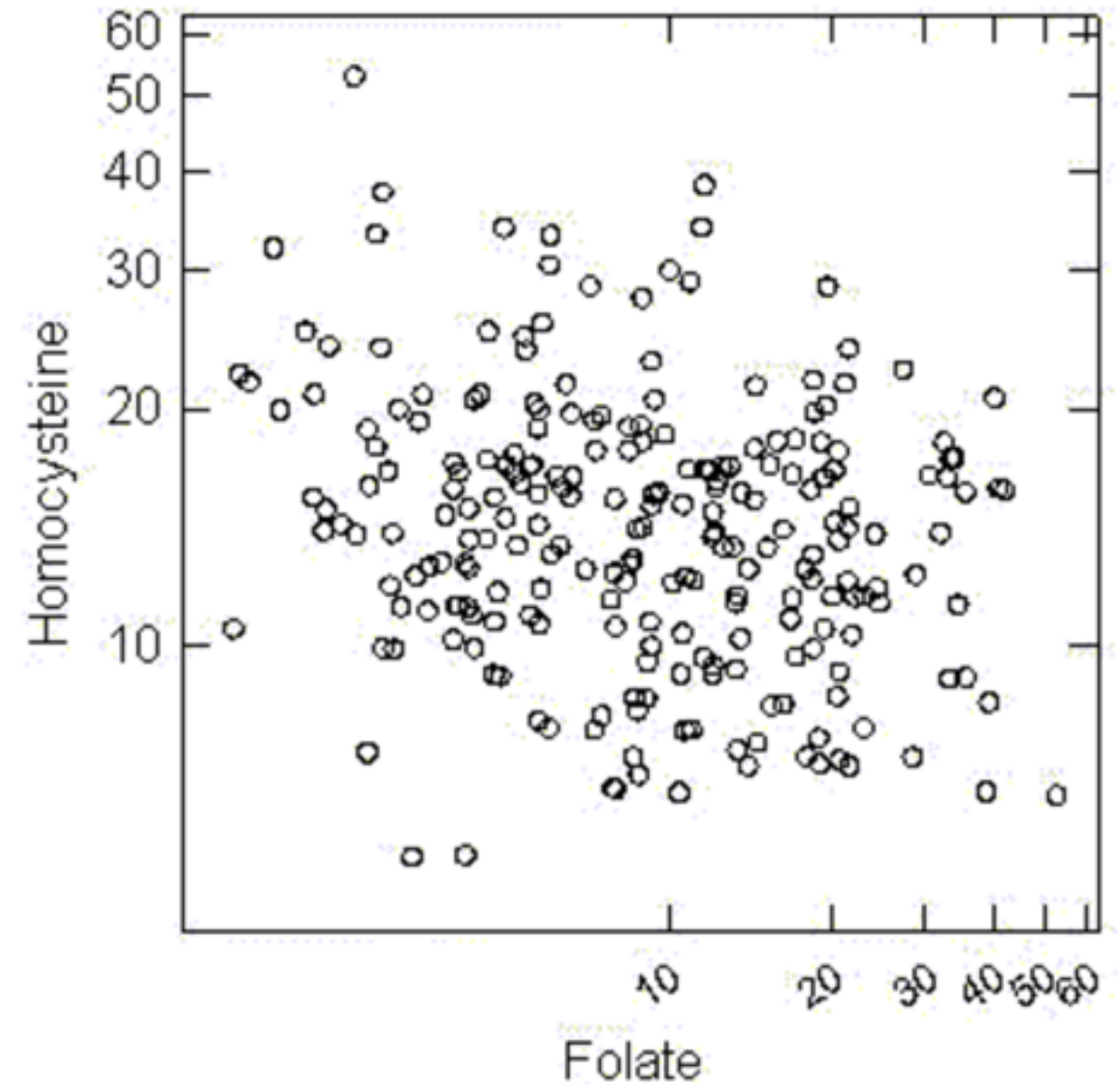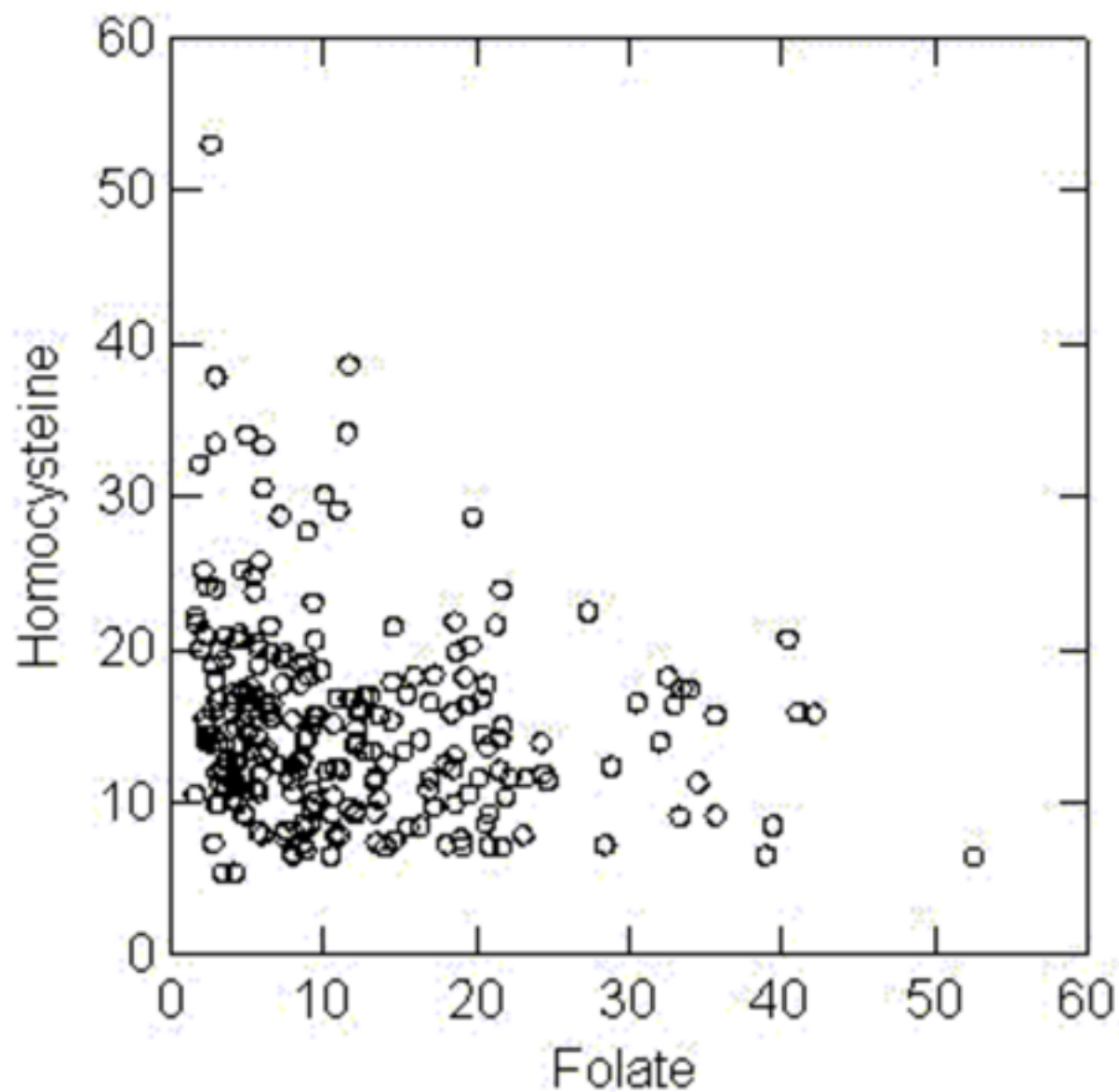
# Non-linear Features

# Non-linear Features

# Y-variable Transform

# Potential Transformations

| Method | Transformation(s) | Regression equation | Predicted value ($\hat{y}$) |
|---|---|---|---|
| Standard linear regression | None | $y = b_0 + b_1 x$ | $\hat{y} = b_0 + b_1 x$ |
| Exponential model | Dependent variable = log(y) | $\log(y) = b_0 + b_1 x$ | $\hat{y} = 10^{b_0 + b_1 x}$ |
| Quadratic model | Dependent variable = sqrt(y) | $\mathrm{sqrt}(y) = b_0 + b_1 x$ | $\hat{y} = (b_0 + b_1 x)^2$ |
| Reciprocal model | Dependent variable = 1/y | $1/y = b_0 + b_1 x$ | $\hat{y} = 1 / (b_0 + b_1 x)$ |
| Logarithmic model | Independent variable = log(x) | $y = b_0 + b_1 \log(x)$ | $\hat{y} = b_0 + b_1 \log(x)$ |
| Power model | Dependent variable = log(y) Independent variable = log(x) | $\log(y) = b_0 + b_1 \log(x)$ | $\hat{y} = 10^{b_0 + b_1 \log(x)}$ |

# Standard Errors

$$\widehat{se}(\hat{b}) = \sqrt{\frac{n\hat{\sigma}^2}{n\sum x_i^2 - (\sum x_i)^2}}.$$

The denominator can be written as

$$n \sum_i (x_i - \bar{x})^2$$

Thus,

$$\widehat{se}(\hat{b}) = \sqrt{\frac{\hat{\sigma}^2}{\sum_i (x_i - \bar{x})^2}}$$

With

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_i \hat{\epsilon}_i^2$$
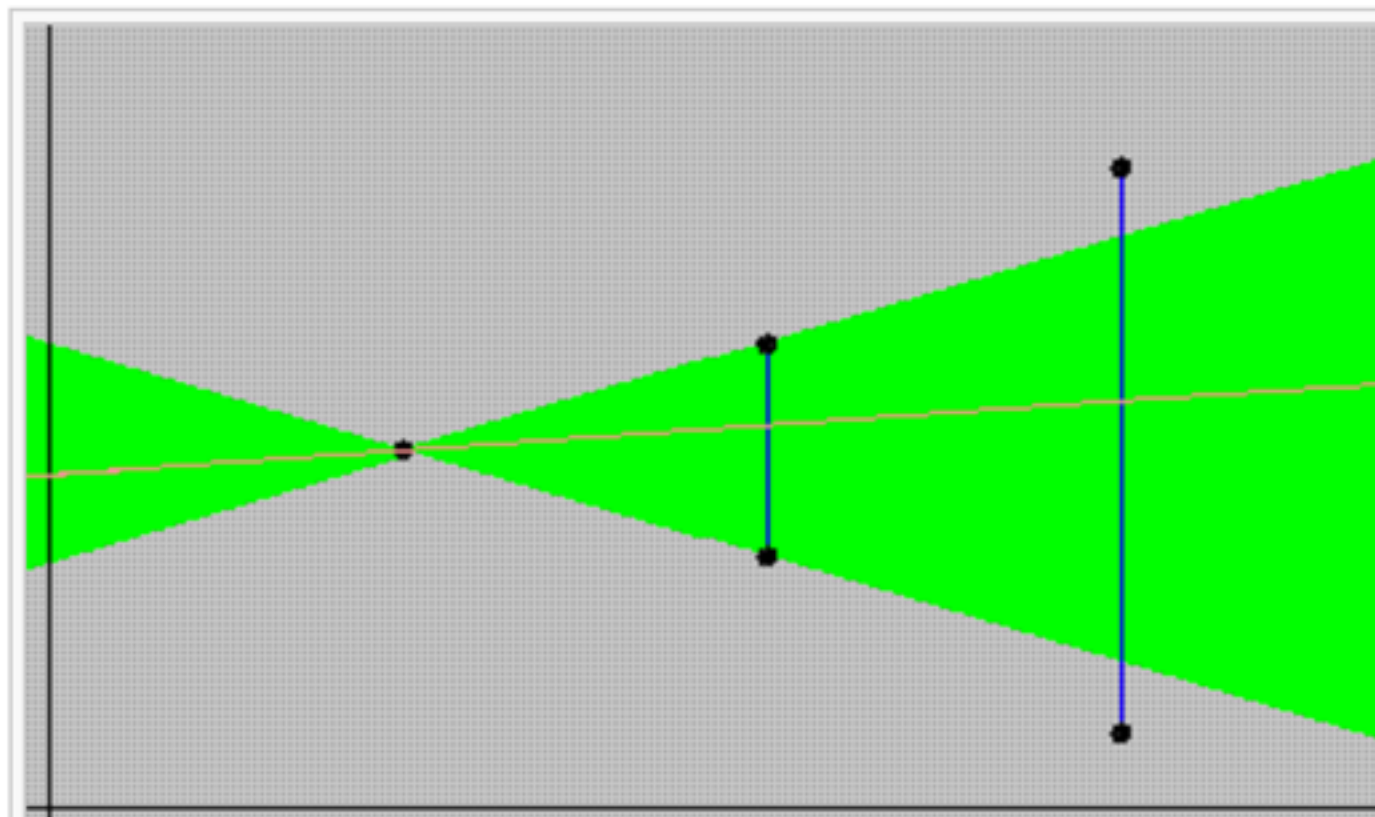
# Why LAD gives multiple solutions



Figure A: A set of data points with reflection symmetry and multiple least absolute deviations solutions. The "solution area" is shown in green. The vertical blue lines represent the absolute errors from the pink line to each data point. The pink line is one of infinitely many solutions within the green area.