# Linear Regression

18.05 Spring 2014

# Agenda

- Fitting curves to bivariate data

- Measuring the goodness of fit

- The fit vs. complexity tradeoff

- Regression to the mean

- Multiple linear regression

# Modeling bivariate data as a function + noise

**Ingredients**

- Bivariate data $(x_1, y_1)$, $(x_2, y_2)$, ..., $(x_n, y_n)$.

- Model: $\quad y_i = f(x_i) + E_i$

  where $f(x)$ is some function, $E_i$ random error.

- Total squared error: $\displaystyle \sum_{i=1}^{n} E_i^2 = \sum_{i=1}^{n} (y_i - f(x_i))^2$

Model allows us to predict the value of $y$ for any given value of $x$.

- $x$ is called the independent or predictor variable.

- $y$ is the dependent or response variable.

# Examples of $f(x)$

- lines: $\quad\quad\quad\quad y = ax + b + E$

- polynomials: $\quad y = ax^2 + bx + c + E$

- other: $\quad\quad\quad y = a/x + b + E$

- other: $\quad\quad\quad y = a\sin(x) + b + E$

# Simple linear regression: finding the best fitting line

- Bivariate data $(x_1, y_1), \ldots, (x_n, y_n)$.
- Simple linear regression: fit a line to the data

$$y_i = ax_i + b + E_i, \quad \text{where} \quad E_i \sim \mathsf{N}(0, \sigma^2)$$

  and where $\sigma$ is a fixed value, the same for all data points.

- Total squared error: $\displaystyle\sum_{i=1}^{n} E_i^2 \;=\; \sum_{i=1}^{n}(y_i - ax_i - b)^2$

- Goal: Find the values of $a$ and $b$ that give the 'best fitting line'.

- Best fit: (least squares)
  The values of $a$ and $b$ that minimize the total squared error.

# Linear Regression: finding the best fitting polynomial

- Bivariate data: $(x_1, y_1), \ldots, (x_n, y_n)$.

- Linear regression: fit a parabola to the data

$$y_i = ax_i^2 + bx_i + c + E_i, \quad \text{where} \quad E_i \sim N(0, \sigma^2)$$

  and where $\sigma$ is a fixed value, the same for all data points.

- Total squared error: $\displaystyle\sum_{i=1}^{n} E_i^2 = \sum_{i=1}^{n}(y_i - ax_i^2 - bx_i - c)^2$.
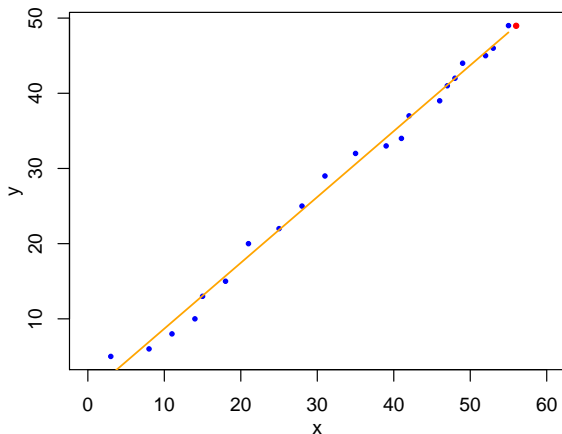
- Goal:
  Find the values of $a$, $b$, $c$ that give the 'best fitting parabola'.

- Best fit: (least squares)
  The values of $a$, $b$, $c$ that minimize the total squared error.
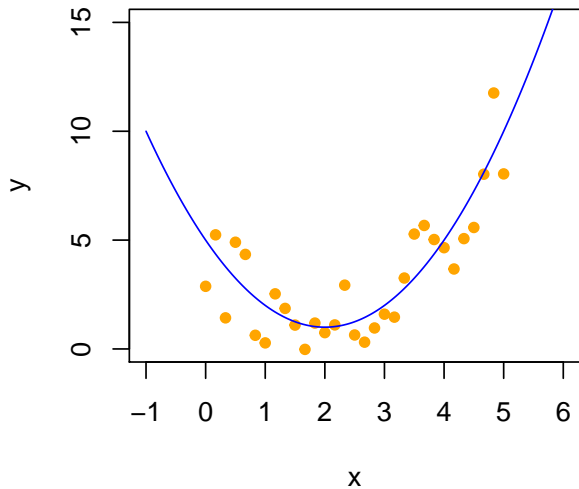
  Can also fit higher order polynomials.

# Stamps



Stamp cost (cents) vs. time (years since 1960)
(Red dot = 49 cents is predicted cost in 2016.)

(Actual cost of a stamp dropped from 49 to 47 cents on 4/8/16.)

# Parabolic fit

## Board question: make it fit

Bivariate data:

$$(1, 3), \ (2, 1), \ (4, 4)$$

**1.** Do (simple) linear regression to find the best fitting line.
Hint: minimize the total squared error by taking partial derivatives with respect to $a$ and $b$.

**2.** Do linear regression to find the best fitting parabola.

**3.** Set up the linear regression to find the best fitting cubic. but don't take derivatives.

**4.** Find the best fitting exponential $y = e^{ax+b}$.
Hint: take $\ln(y)$ and do simple linear regression.

# What is linear about linear regression?

Linear in the parameters $a$, $b$, . . ..
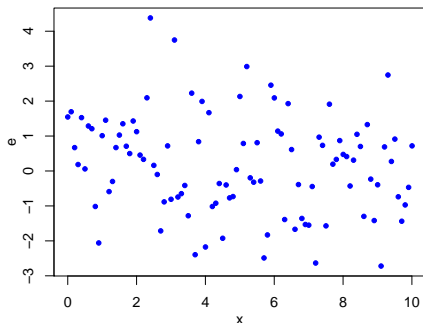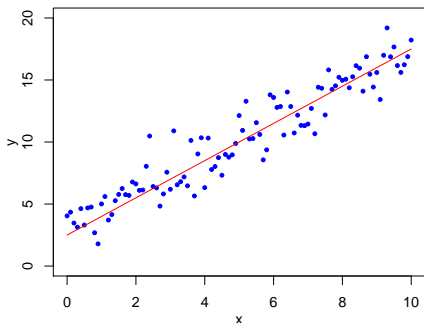
$$y = ax + b.$$
$$y = ax^2 + bx + c.$$

It is **not** because the curve being fit has to be a straight line –although this is the simplest and most common case.

Notice: in the board question you had to solve a system of simultaneous linear equations.

Fitting a line is called simple linear regression.

# Homoscedastic
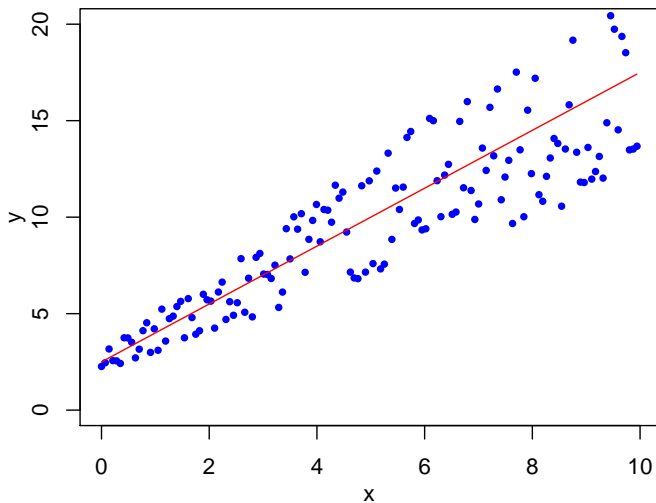
**BIG ASSUMPTIONS**: the $E_i$ are independent with the same variance $\sigma^2$.



Regression line (left) and residuals (right).
Homoscedasticity = uniform spread of errors around regression line.

# Heteroscedastic



Heteroscedastic Data

# Formulas for simple linear regression

Model:
$$y_i = ax_i + b + E_i \quad \text{where} \quad E_i \sim \mathsf{N}(0, \sigma^2).$$

Using calculus or algebra:
$$\hat{a} = \frac{s_{xy}}{s_{xx}} \quad \text{and} \quad \hat{b} = \bar{y} - \hat{a}\,\bar{x},$$

where
$$\bar{x} = \frac{1}{n} \sum x_i \quad s_{xx} = \frac{1}{n-1} \sum (x_i - \bar{x})^2$$
$$\bar{y} = \frac{1}{n} \sum y_i \quad s_{xy} = \frac{1}{n-1} \sum (x_i - \bar{x})(y_i - \bar{y}).$$

**WARNING:** This is just for simple linear regression. For polynomials and other functions you need other formulas.

# Board Question: using the formulas plus some theory

Bivariate data: $(1, 3)$, $(2, 1)$, $(4, 4)$

**1.(a)** Calculate the sample means for $x$ and $y$.

**1.(b)** Use the formulas to find a best-fit line in the $xy$-plane.

$$\hat{a} = \frac{s_{xy}}{s_{xx}} \qquad\qquad \hat{b} = \overline{y} - \hat{a}\overline{x}$$

$$s_{xy} = \frac{1}{n-1} \sum (x_i - \overline{x})(y_i - \overline{y}) \quad s_{xx} = \frac{1}{n-1} \sum (x_i - \overline{x})^2.$$

**2.** Show the point $(\overline{x}, \overline{y})$ is always on the fitted line.

**3.** Under the assumption $E_i \sim N(0, \sigma^2)$ show that the least squares method is equivalent to finding the MLE for the parameters $(a, b)$.

Hint: $f(y_i \mid x_i, a, b) \sim N(ax_i + b, \sigma^2)$.

# Measuring the fit

$y = (y_1, \cdots, y_n) =$ data values of the response variable.

$\hat{y} = (\hat{y}_1, \cdots, \hat{y}_n) =$ 'fitted values' of the response variable.

- TSS $= \sum (y_i - \overline{y})^2 =$ total sum of squares $=$ total variation.

- RSS $= \sum (y_i - \hat{y}_i)^2 =$ residual sum of squares.
  RSS $=$ unexplained by model squared error (due to random fluctuation)

- $RSS/TSS =$ unexplained fraction of the total error.

- $R^2 = 1 - RSS/TSS$ is measure of goodness-of-fit

- $R^2$ is the fraction of the variance of $y$ explained by the model.

# Overfitting a polynomial

- Increasing the degree of the polynomial increases $R^2$

- Increasing the degree of the polynomial increases the complexity of the model.

- The optimal degree is a tradeoff between goodness of fit and complexity.

- If all data points lie on the fitted curve, then $y = \hat{y}$ and $R^2 = 1$.

R demonstration!

# Outliers and other troubles

Question: Can one point change the regression line significantly?

Use mathlet
http://mathlets.org/mathlets/linear-regression/

## Regression to the mean

- Suppose a group of children is given an IQ test at age 4. One year later the same children are given another IQ test.

- Children's IQ scores at age 4 and age 5 should be positively correlated.

- Those who did poorly on the first test (e.g., bottom 10%) will tend to show improvement (i.e. regress to the mean) on the second test.

- A completely useless intervention with the poor-performing children might be misinterpreted as causing an increase in their scores.

- Conversely, a reward for the top-performing children might be misinterpreted as causing a decrease in their scores.

This example is from Rice *Mathematical Statistics and Data Analysis*

# A brief discussion of multiple linear regression

Multivariate data: $(x_{i,1}, x_{i,2}, \ldots, x_{i,m}, y_i)$ ($n$ data points: $i = 1, \ldots, n$)

Model $\hat{y}_i = a_1 x_{i,1} + a_2 x_{i,2} + \ldots + a_m x_{i,m}$

$x_{i,j}$ are the explanatory (or predictor) variables.

$y_i$ is the response variable.

The total squared error is

$$\sum_{i=1}^{n} (y_i - \hat{y}_i)^2 = \sum_{i=1}^{n} (y_i - a_1 x_{i,1} - a_2 x_{i,2} - \ldots - a_m x_{i,m})^2$$

18.05 Introduction to Probability and Statistics
Spring 2014