

Part 4

Chapter 13

Linear Regression

Chapter Objectives

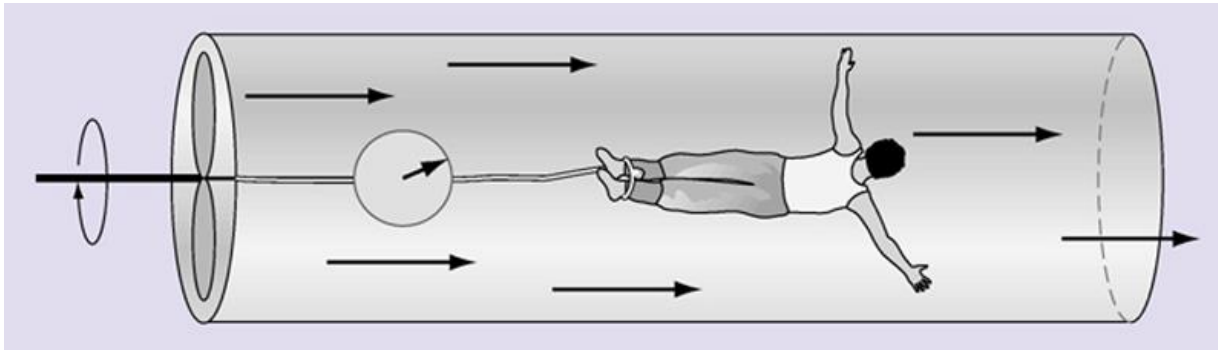
- Familiarizing yourself with some basic descriptive statistics and the normal distribution.
- Knowing how to compute the slope and intercept of a best fit straight line with linear regression.
- Knowing how to compute and understand the meaning of the coefficient of determination and the standard error of the estimate.
- Understanding how to use transformations to linearize nonlinear equations so that they can be fit with linear regression.
- Knowing how to implement linear regression with MATLAB.

Introduction (1/2)

- A free-falling bungee jumper is subjected to the air resistance force. This force was proportional to the square of velocity as in .

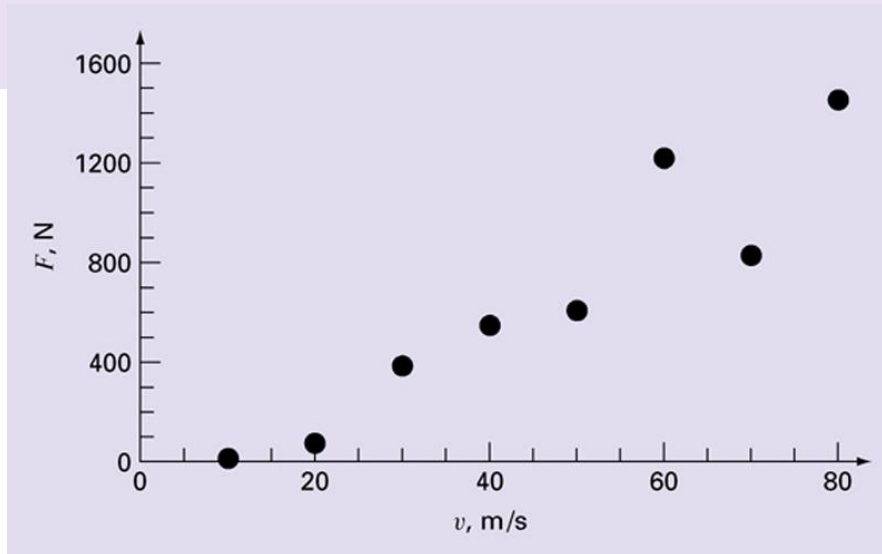
$$F_U = c_d v^2$$

- Experiments can formulation play a critical role in this formulation.



Wind tunnel experiment to measure how the force of air resistance depends on velocity.

Introduction (2/2)



Plot of forces vs. wind velocity for an object suspended in a wind tunnel

v (m/s)	10	20	30	40	50	60	70	80
F (N)	25	70	380	550	610	1220	830	1450

- The forces increase with increasing velocity.
- What kind of relationship between forces and velocities?
- Linear, square, or others?
- How to fit a “best” line or curve to these data?

Statistics Review

Measure of Location

- Arithmetic mean: the sum of the individual data points (y_i) divided by the number of points n :

$$\bar{y} = \frac{\sum y_i}{n}$$

- Median: the midpoint of a group of data.
- Mode: the value that occurs most frequently in a group of data.

Statistics Review

Measures of Spread

- Standard deviation:

$$s_y = \sqrt{\frac{S_t}{n-1}}$$

where S_t is the sum of the squares of the data residuals:

$$S_t = \sum (y_i - \bar{y})^2$$

and $n-1$ is referred to as the degrees of freedom.

- Variance:

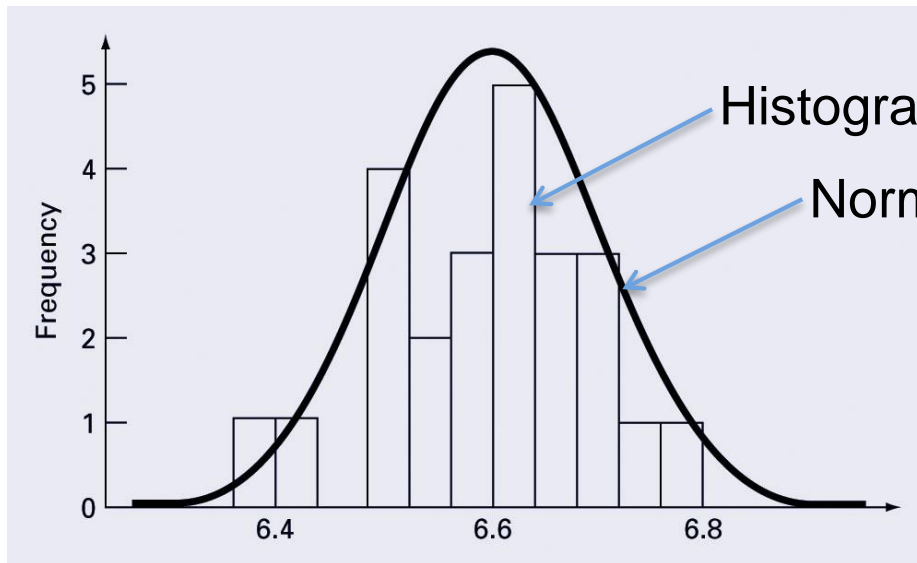
$$s_y^2 = \frac{\sum (y_i - \bar{y})^2}{n-1} = \frac{\sum y_i^2 - (\sum y_i)^2 / n}{n-1}$$

- Coefficient of variation: the ratio of standard deviation to the mean.

$$\text{c.v.} = \frac{s_y}{\bar{y}} \times 100\%$$

Normal Distribution

- Data distribution
 - Shape with which the data is spread around the mean.
- A histogram
 - Constructed by sorting the measurements into intervals, or bins.
- If we have a very large set of data, the histogram can be approximated by a smooth curve, which is symmetric, bell-shaped curve called the *normal distribution*.



The range between $\bar{y} - s_y$ and $\bar{y} + s_y$ will 68% of the total measurements.
95% for $\bar{y} - 2s_y$ and $\bar{y} + 2s_y$

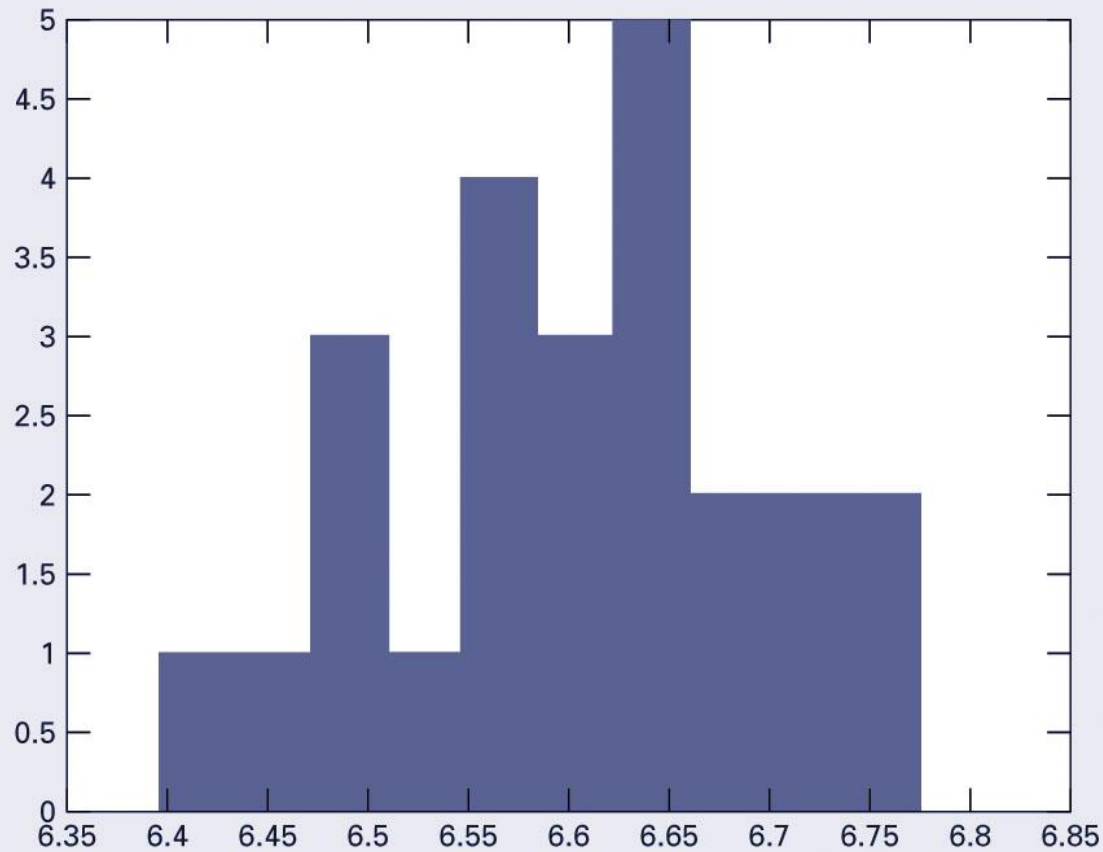
Descriptive Statistics in MATLAB

- MATLAB has several built-in commands to compute and display descriptive statistics.
- Assuming some column vector s :
 - `mean(s)`, `median(s)`, `mode(s)`
 - Calculate the mean, median, and mode of s .
mode is a part of the statistics toolbox.
 - `min(s)`, `max(s)`
 - Calculate the minimum and maximum value in s .
 - `var(s)`, `std(s)`
 - Calculate the variance and standard deviation of s
- Note - if a matrix is given, the statistics will be returned for each column.

Histograms in MATLAB

- $[n, X] = \text{hist}(s, x)$
 - Determine the number of elements in each bin of data in s .
 x is a vector containing the center values of the bins.
- $[n, x] = \text{hist}(s, m)$
 - Determine the number of elements in each bin of data in s using m bins. x will contain the centers of the bins.
 - The default case is $m=10$
- $\text{hist}(s, x)$ or $\text{hist}(s, m)$ or $\text{hist}(s)$
 - With no output arguments, hist will actually produce a histogram.

Histogram Example



Linear Least-Squares Regression

- Linear least-squares regression is a method to determine the “best” coefficients in a linear model for given data set.
- “Best” for least-squares regression means minimizing the sum of the squares of the estimate residuals.
- For a straight line model, this gives:

$$S_r = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - a_0 - a_1 x_i)^2$$

$$y = a_0 + a_1 x$$

- This method will yield a unique line for a given set of data.

Least-Squares Fit of a Straight Line

- For a minimum to occur, it is necessary that

$$\frac{\partial S_r}{\partial a_0} = -2 \sum (y_i - a_0 - a_1 x_i) = 0 \quad \frac{\partial S_r}{\partial a_1} = -2 \sum [(y_i - a_0 - a_1 x_i) x_i] = 0$$

$$0 = \sum y_i - \sum a_0 - \sum a_1 x_i$$

$$0 = \sum x_i y_i - \sum a_0 x_i - \sum a_1 x_i^2$$

$$n a_0 + \left(\sum x_i \right) a_1 = \sum y_i$$

$$\left(\sum x_i \right) a_0 + \left(\sum x_i^2 \right) a_1 = \sum x_i y_i$$

- These two equations can be solved simultaneously for

$$a_1 = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - \left(\sum x_i \right)^2}$$

$$a_0 = \bar{y} - a_1 \bar{x}$$

Where \bar{x} and \bar{y} are the means of x and y , respectively.

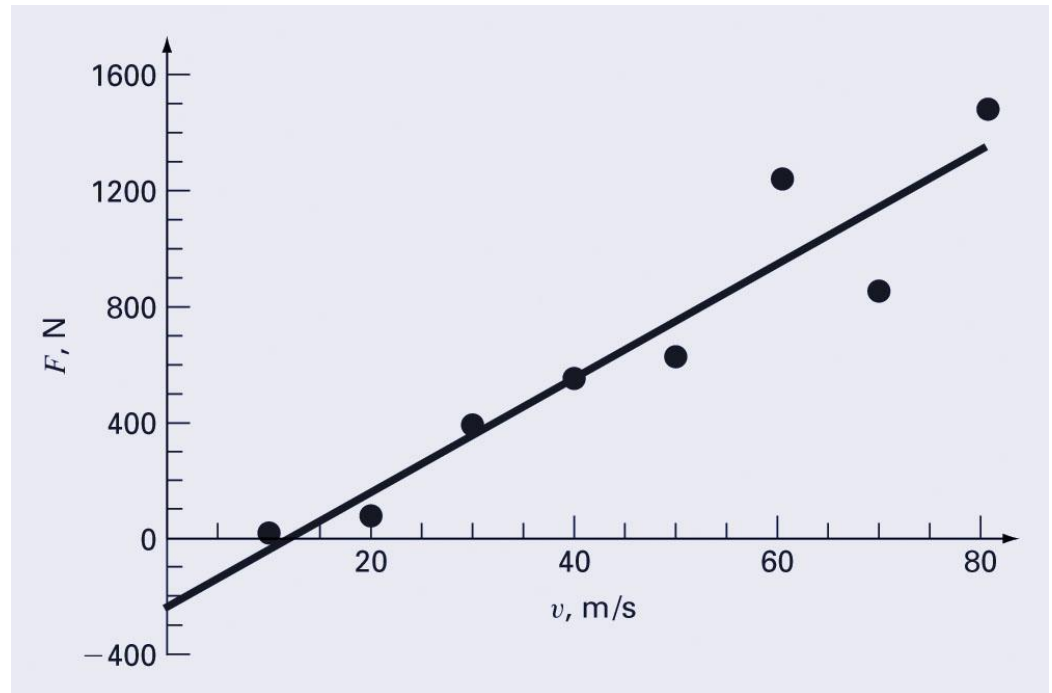
Example 13.2

	V (m/s)	F (N)		
i	x_i	y_i	$(x_i)^2$	$x_i y_i$
1	10	25	100	250
2	20	70	400	1400
3	30	380	900	11400
4	40	550	1600	22000
5	50	610	2500	30500
6	60	1220	3600	73200
7	70	830	4900	58100
8	80	1450	6400	116000
Σ	360	5135	20400	312850

$$a_1 = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2} = \frac{8(312850) - (360)(5135)}{8(20400) - (360)^2} = 19.47024$$

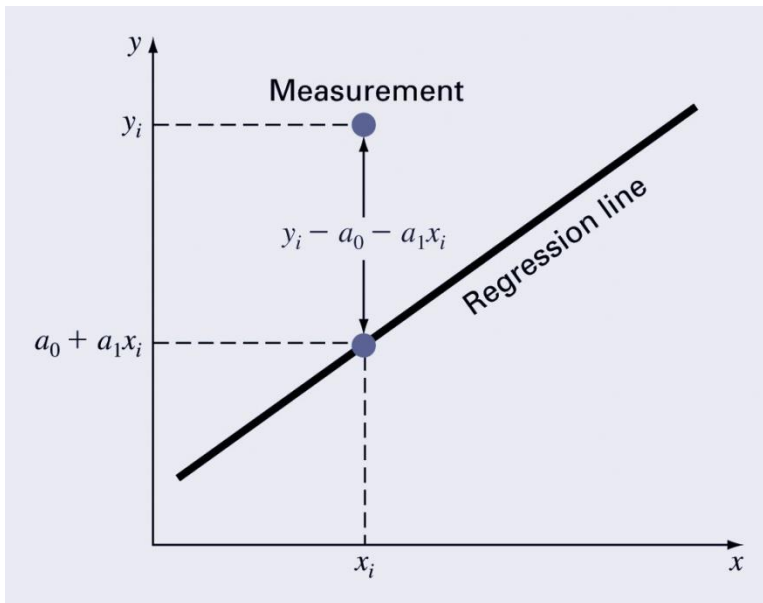
$$a_0 = \bar{y} - a_1 \bar{x} = 641.875 - 19.47024(45) = -234.2857$$

$$F_{est} = -234.2857 + 19.47024v$$



Quantification of Error

- Recall for a straight line, the sum of the squares of the estimate residuals:



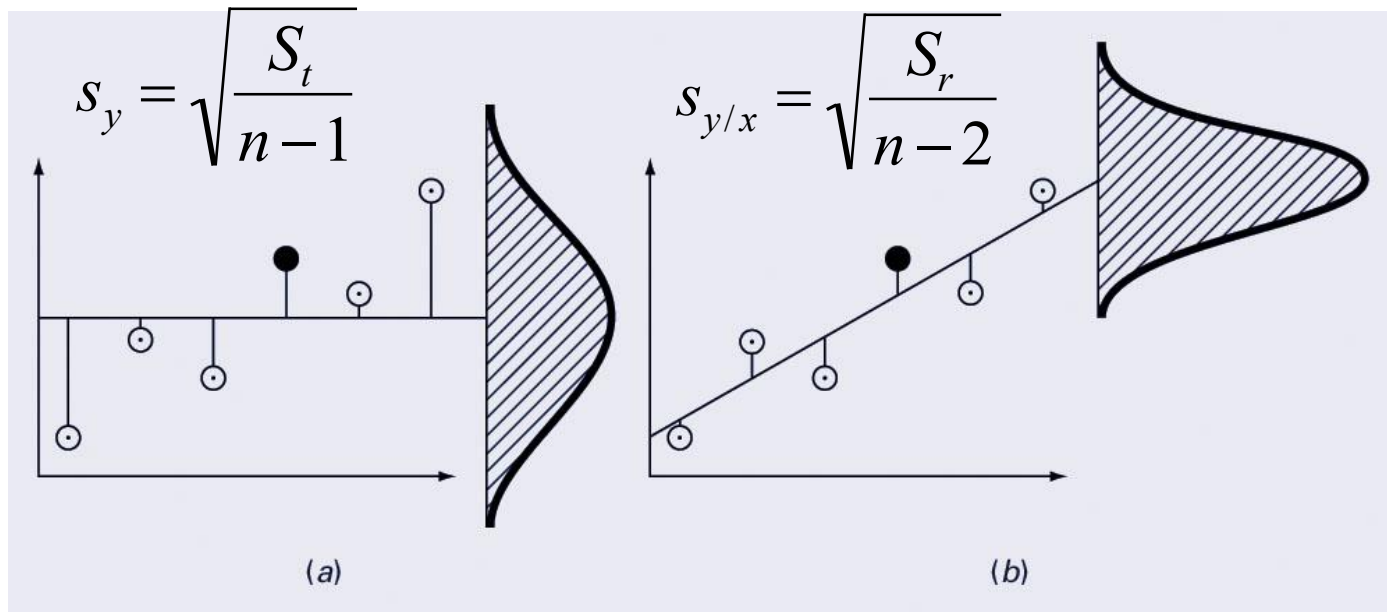
$$S_r = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - a_0 - a_1 x_i)^2$$

- Standard error of the estimate:

$$s_{y/x} = \sqrt{\frac{S_r}{n-2}}$$

Standard Error of the Estimate

- Regression data showing (a) the spread of data around the mean of the dependent data and (b) the spread of the data around the best fit line:



- The reduction in spread represents the improvement due to linear regression.

Coefficient of Determination

- The coefficient of determination r^2 is the difference between the sum of the squares of the data residuals and the sum of the squares of the estimate residuals, normalized by the sum of the squares of the data residuals:

$$r^2 = \frac{S_t - S_r}{S_t}$$

r^2 represents the percentage of the original uncertainty explained by the model.

- For a perfect fit, $S_r=0$ and $r^2 =1$.
- If $r^2 =0$, there is no improvement over simply picking the mean.
- If $r^2 <0$, the model is worse than simply picking the mean!

Example 13.3

	V (m/s)	F (N)			
i	x_i	y_i	$a_0+a_1x_i$	$(y_i - \bar{y})^2$	$(y_i - a_0 - a_1x_i)^2$
1	10	25	-39.58	380535	4171
2	20	70	155.12	327041	7245
3	30	380	349.82	68579	911
4	40	550	544.52	8441	30
5	50	610	739.23	1016	16699
6	60	1220	933.93	334229	81837
7	70	830	1128.63	35391	89180
8	80	1450	1323.33	653066	16044
Σ	360	5135		1808297	216118

$$F_{est} = -234.2857 + 19.47024v$$

$$S_t = \sum (y_i - \bar{y})^2 = 1808297$$

$$S_r = \sum (y_i - a_0 - a_1x_i)^2 = 216118$$

$$s_y = \sqrt{\frac{1808297}{8-1}} = 508.26$$

$$s_{y/x} = \sqrt{\frac{216118}{8-2}} = 189.79$$

$$r^2 = \frac{1808297 - 216118}{1808297} = 0.8805$$

88.05% of the original uncertainty has been explained by the linear model

Nonlinear Relationships

- Linear regression is predicated on the fact that the relationship between the dependent and independent variables is linear - this is not always the case.
- Three common examples are:

exponential: $y = \alpha_1 e^{\beta_1 x}$

power: $y = \alpha_2 x^{\beta_2}$

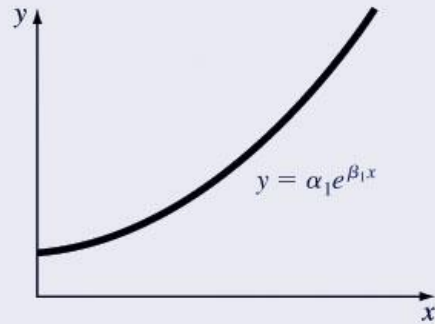
saturation-growth-rate: $y = \alpha_3 \frac{x}{\beta_3 + x}$

Linearization of Nonlinear Relationships

- One option for finding the coefficients for a nonlinear fit is to linearize it. For the three common models, this may involve taking logarithms or inversion:

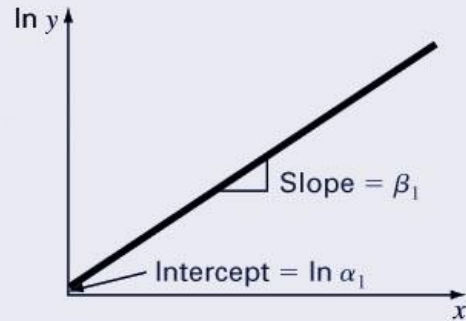
Model	Nonlinear	Linearized
exponential :	$y = \alpha_1 e^{\beta_1 x}$	$\ln y = \ln \alpha_1 + \beta_1 x$
power:	$y = \alpha_2 x^{\beta_2}$	$\log y = \log \alpha_2 + \beta_2 \log x$
saturation - growth - rate :	$y = \alpha_3 \frac{x}{\beta_3 + x}$	$\frac{1}{y} = \frac{1}{\alpha_3} + \frac{\beta_3}{\alpha_3} \frac{1}{x}$

Transformation Examples

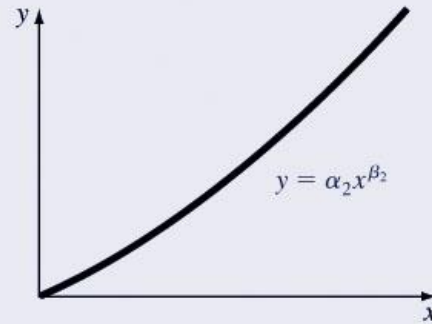


(a)

Linearization

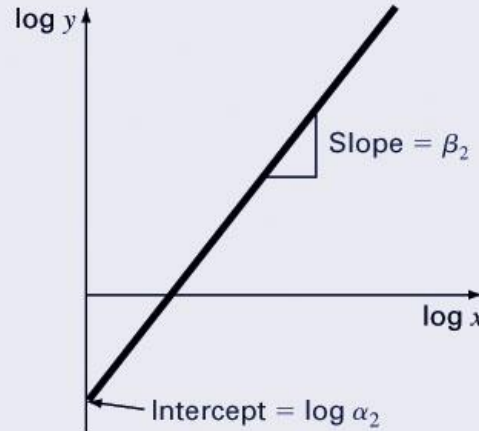


(d)

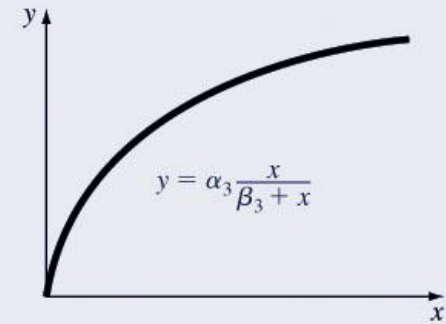


(b)

Linearization

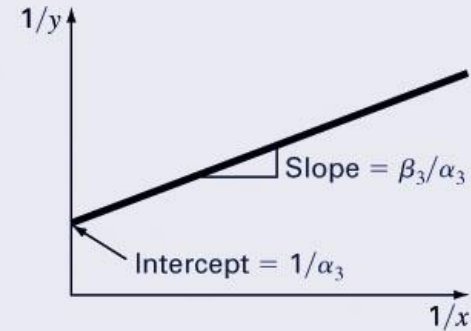


(e)



(c)

Linearization



(f)

Example 13.4 (1/2)

- Q. Fit Eq. (13.23) to the data below using log transformation.

i	x_i	y_i	$\log x_i$	$\log y_i$	$(\log x_i)^2$	$\log x_i \log y_i$
1	10	25	1.000	1.398	1.000	1.398
2	20	70	1.301	1.845	1.693	2.401
3	30	380	1.477	2.580	2.182	3.811
4	40	550	1.602	2.740	2.567	4.390
5	50	610	1.699	2.785	2.886	4.732
6	60	1,220	1.778	3.086	3.162	5.488
7	70	830	1.845	2.919	3.404	5.386
8	80	1,450	1.903	3.161	3.622	6.016
Σ			12.606	20.515	20.516	33.622

Example 13.4 (2/2)

$$a_1 = \frac{n \sum \log x_i \log y_i - \sum \log x_i \sum \log y_i}{n \sum (\log x_i)^2 - (\sum \log x_i)^2} \quad a_0 = \bar{y} - a_1 \bar{x}$$

$$\bar{x} = \frac{12.606}{8} = 1.5757 \quad \bar{y} = \frac{20.515}{8} = 2.5644$$

$$a_1 = \frac{8(33.622) - 12.606(20.515)}{8(20.516) - (12.606)^2} = 1.9842$$

$$a_0 = 2.5644 - 1.9842(1.5757) = -0.5620$$

$$\log y = -0.5620 + 1.9842 \log x \quad F = 0.2741v^{1.9842}$$

$$\log y = \log \alpha_2 + \beta_2 \log x \quad y = \alpha_2 x^{\beta_2}$$

Linear Regression Program

```
function [a, r2] = linregr(x,y)
% linregr: linear regression curve fitting
% [a, r2] = linregr(x,y):Least squares fit of straight
% line to data by solving the normal equations

% input:
% x = independent variable
% y = dependent variable
% output:
% a = vector of slope, a(1), and intercept, a(2)
% r2 = coefficient of determination

n = length(x);
if length(y)~=n, error('x and y must be same length'); end
x = x(:); y = y(:); % convert to column vectors
sx = sum(x); sy = sum(y);
sx2 = sum(x.*x); sxy = sum(x.*y); sy2 = sum(y.*y);
a(1) = (n*sxy-sx*sy)/(n*sx2-sx^2);
a(2) = sy/n-a(1)*sx/n;
r2 = ((n*sxy-sx*sy)/sqrt(n*sx2-sx^2)/sqrt(n*sy2-sy^2))^2;
% create plot of data and best fit line
xp = linspace(min(x),max(x),2);
yp = a(1)*xp+a(2);
plot(x,y,'o',xp,yp)
grid on
```

MATLAB Functions

- MATLAB has a built-in function **polyfit** that fits a least-squares n^{th} order polynomial to data:

```
>> p = polyfit(x, y, n)
```

```
% x: independent data
```

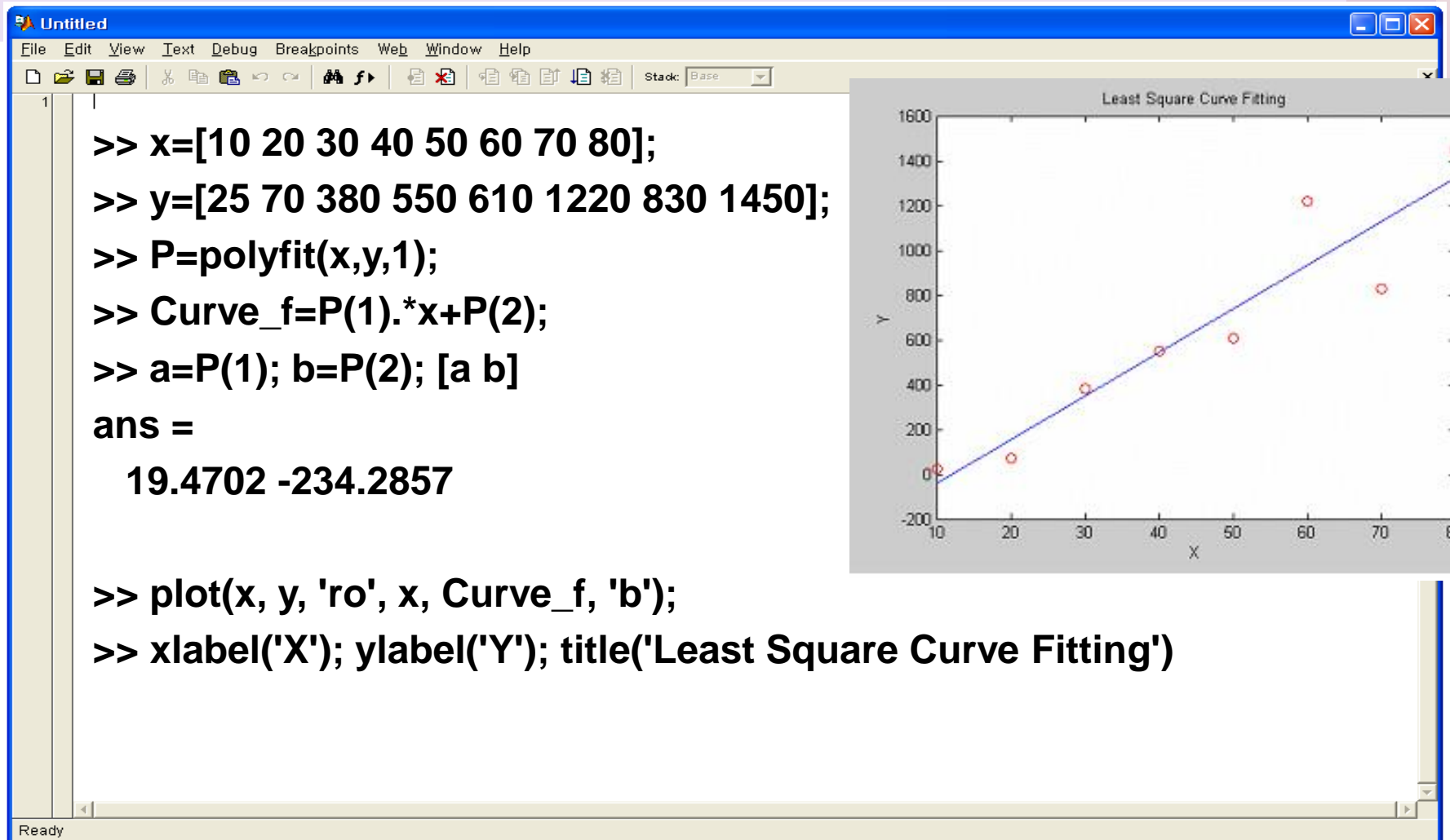
```
% y: dependent data
```

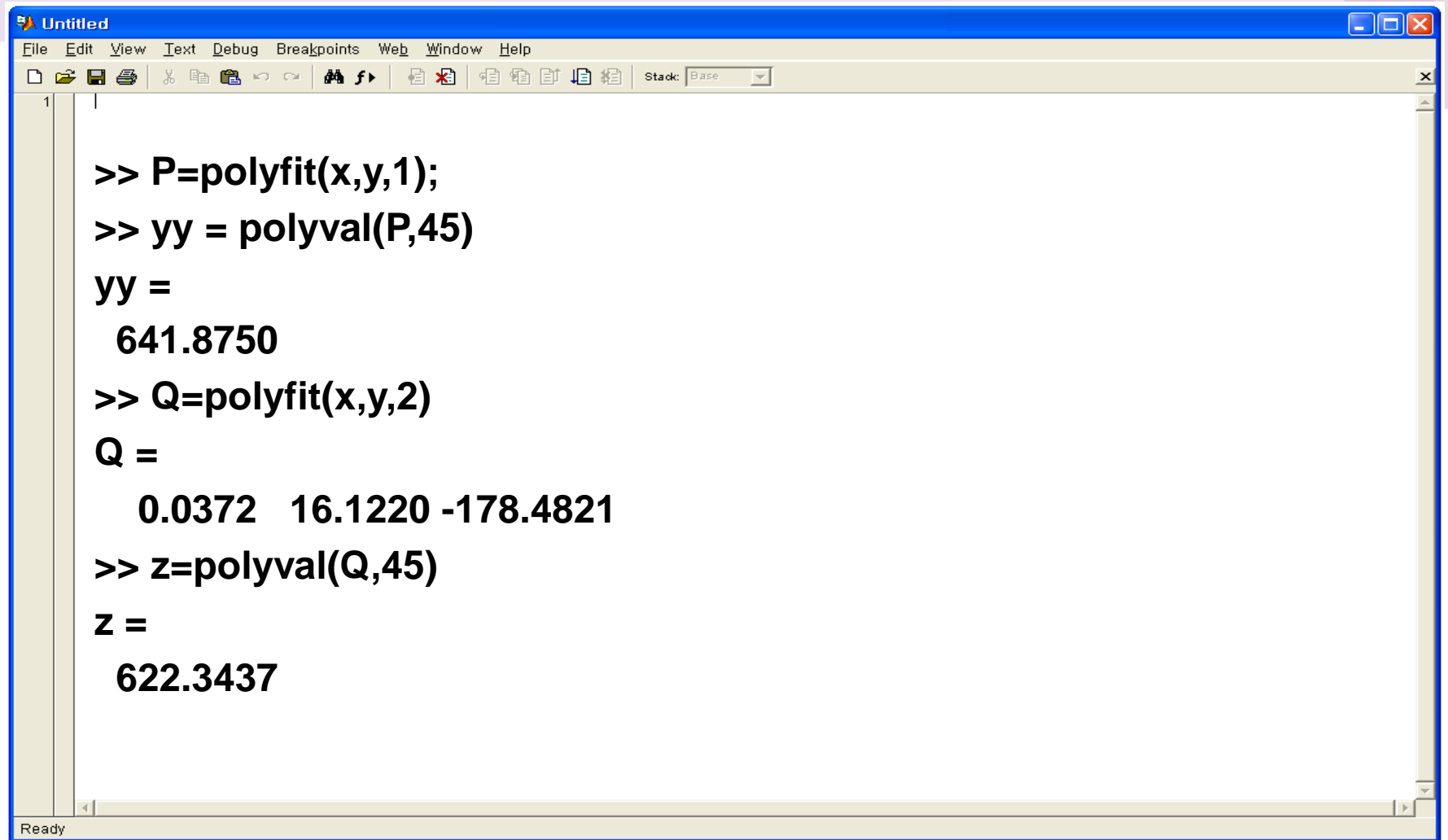
```
% n: order of polynomial to fit
```

```
% p: coefficients of polynomial
```

```
%  $f(x) = p_1x^n + p_2x^{n-1} + \dots + p_nx + p_{n+1}$ 
```

- MATLAB's **polyval** command can be used to compute a value using the coefficients.
- $y = \text{polyval}(p, x)$



A screenshot of a MATLAB command window titled "Untitled". The window has a blue title bar and a menu bar with options: File, Edit, View, Text, Debug, Breakpoints, Web, Window, Help. Below the menu bar is a toolbar with various icons for file operations and execution. The main area of the window contains the following MATLAB commands and their outputs:

```
1 |  
  
>> P=polyfit(x,y,1);  
>> yy = polyval(P,45)  
yy =  
    641.8750  
  
>> Q=polyfit(x,y,2)  
Q =  
    0.0372    16.1220   -178.4821  
  
>> z=polyval(Q,45)  
z =  
    622.3437
```

The status bar at the bottom left of the window shows "Ready".