

UCLA STAT 110A Applied Statistics

- **Instructor:** Ivo Dinov,
Asst. Prof. In Statistics and Neurology
- **Teaching Assistants:** Helen Hu, UCLA Statistics

University of California, Los Angeles, Spring 2002
<http://www.stat.ucla.edu/~dinov/>

STAT 110A, UCLA, Ivo Dinov Slide 1

Lines in 2D (Regression and Correlation)

- Vertical Lines
- Horizontal Lines
- Oblique lines
- Increasing/Decreasing
- Slope of a line
- Intercept
- $Y = \alpha X + \beta$, in general.

Math Equation for the Line?

STAT 110A, UCLA, Ivo Dinov Slide 2

Chapter 7: Lines in 2D (Regression and Correlation)

- Draw the following lines:
 - $Y = 2X + 1$
 - $Y = -3X - 5$
 - Line through (X_1, Y_1) and (X_2, Y_2) .
 - $(Y - Y_1) / (Y_2 - Y_1) = (X - X_1) / (X_2 - X_1)$.

Math Equation for the Line?

STAT 110A, UCLA, Ivo Dinov Slide 3

Approaches for modeling data relationships Regression and Correlation

- There are **random** and **nonrandom** variables
- **Correlation** applies if both variables (X/Y) are random (e.g., We saw a previous example, systolic vs. diastolic blood pressure SISVOL/DIAVOL) and are treated symmetrically.
- **Regression** applies in the case when you want to single out one of the variables (response variable, Y) and use the other variable as **predictor (explanatory variable, X)**, which explains the behavior of the response variable, Y.

STAT 110A, UCLA, Ivo Dinov Slide 4

Causal relationship? – infant death rate (per 1,000) in 14 countries

Infant death rate

% Breast feeding at 6 months

Strong evidence (linear pattern) of death rate increase with increasing level of breastfeeding (BF)? Naive conclusion breast feeding is bad? But high rates of BF is associated with lower access to H.O.

Predict behavior of Y (response) Based on the values of X (explanatory var.) Strategies for uncovering the reasons (causes) for an observed effect.

% Breast feeding at 6 mo.

% Access to safe water

STAT 110A, UCLA, Ivo Dinov Slide 5

Regression relationship = trend + residual scatter

Retail sales (\$)

Disposable income (\$)

(a) Sales/income

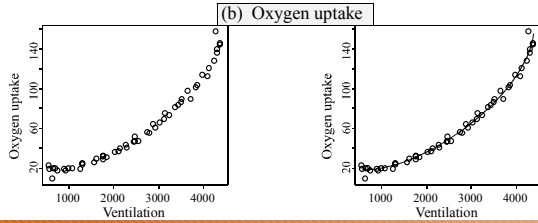
Retail sales (\$)

Disposable income (\$)

- **Regression** is a way of studying relationships between variables (random/nonrandom) for predicting or explaining behavior of 1 variable (**response**) in terms of others (**explanatory variables** or **predictors**).

STAT 110A, UCLA, Ivo Dinov Slide 6

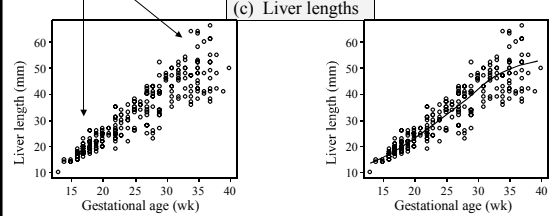
Trend (does not have to be linear) + scatter (could be of any type/distribution)



Slide 7 STAT 110A, UCLA, Jon Dinger

Trend + scatter (fetus liver length in mm)

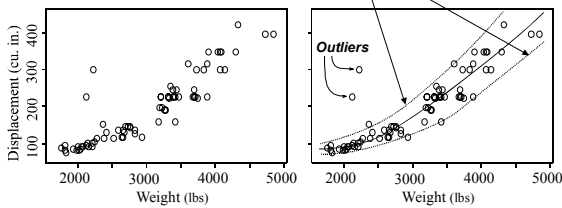
Change of scatter with age



Slide 8 STAT 110A, UCLA, Jon Dinger

Trend + scatter

Dotted curves (confidence intervals) represent the extend of the scatter.

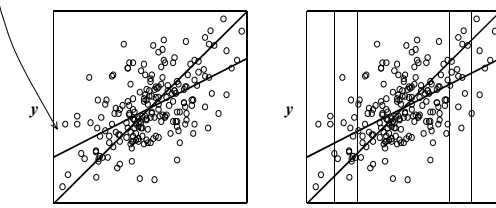


(a) Scatter plot (b) With trend plus scatter
Displacement versus weight for 74 models of automobile.

Slide 9 STAT 110A, UCLA, Jon Dinger

Looking vertically

Flatter line gives better prediction, since it approx. goes through the middle of the Y-range, for each fixed x-value (vertical line)

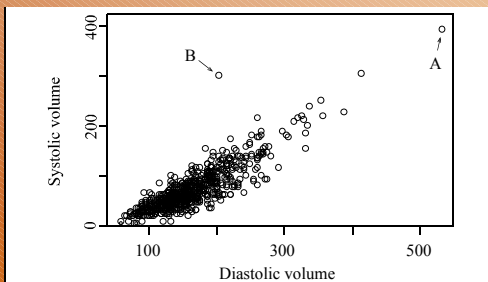


(a) Which line? (b) Flatter line gives better predictions.

----- Educating the eye to look vertically.

Slide 10 STAT 110A, UCLA, Jon Dinger

Outliers – odd, atypical, observations (errors, B, or real data, A)



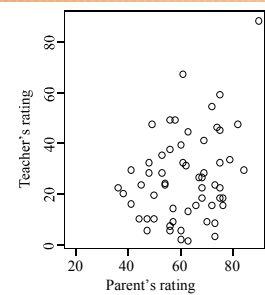
Scatter plot from the heart attack data.

Slide 11 STAT 110A, UCLA, Jon Dinger

A weak relationship

58 abused children are rated (by non-abusive parents and teachers) on a psychological disturbance measure.

How do we quantify weak vs. strong relationship?



Parent's rating versus teacher's rating for abused children.

Slide 12 STAT 110A, UCLA, Jon Dinger

A note of caution!

In **observational** data, **strong relationships** are **not** necessarily **causal**. It is virtually **impossible** to conclude a **cause-and-effect relationship** between variables using observational data!

Slide 13 STAT 1104, UCL, J. van Dine

Essential Points

1. What essential difference is there between the **correlation** and **regression** approaches to a relationship between two variables? (In **correlation** independent variables; **regression** response var depends on explanatory variable.)
2. What are the most common **reasons why people fit regression models** to data? (predict Y or unravel reasons/causes of behavior.)
3. Can you conclude that changes in X caused the changes in Y seen in a scatter plot if you have data from an observational study? (No, there could be **lurking variables**, hidden effects/predictors, also associated with the predictor X, itself, e.g., time is often a lurking variable, or may be that changes in Y cause changes in X, instead of the other way around.)

Slide 14 STAT 1104, UCL, J. van Dine

Essential Points

5. When can you reliably conclude that changes in X cause the changes in Y? (Only when **controlled randomized experiments** are used – levels of X are randomly distributed to available experimental units, or experimental conditions need to be identical for different levels of X, this includes **time**.)

Slide 15 STAT 1104, UCL, J. van Dine

Correlation Coefficient

Correlation coefficient ($-1 \leq R \leq 1$): a measure of linear association, or clustering around a line of multivariate data.

Relationship between two variables (X, Y) can be summarized by: (μ_X, σ_X) , (μ_Y, σ_Y) and the correlation coefficient, R . $R=1$, **perfect positive correlation** (straight line relationship), $R=0$, **no correlation** (random cloud scatter), $R=-1$, **perfect negative correlation**.

Computing $R(X,Y)$: (standardize, multiply, average)

$$R(X, Y) = \frac{1}{N-1} \sum_{k=1}^N \left(\frac{x_k - \mu}{\sigma} \right) \left(\frac{y_k - \mu}{\sigma} \right)$$

$X = \{x_1, x_2, \dots, x_N\}$
 $Y = \{y_1, y_2, \dots, y_N\}$
 $(\mu_X, \sigma_X), (\mu_Y, \sigma_Y)$
 sample mean / SD.

Slide 16 STAT 1104, UCL, J. van Dine

Correlation Coefficient

Example:

$$R(X, Y) = \frac{1}{N-1} \sum_{k=1}^N \left(\frac{x_k - \mu}{\sigma} \right) \left(\frac{y_k - \mu}{\sigma} \right)$$

Student	Height	Weight	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})^2$	$(y_i - \bar{y})^2$	$(x_i - \bar{x})(y_i - \bar{y})$
1	167	60	6	4.67	36	21.8089	28.02
2	170	64	9	6.67	81	44.4889	59.03
3	160	57	-1	1.67	1	2.7889	-1.67
4	152	46	-8	-6.33	64	40.0689	40.97
5	157	55	-4	-3.33	16	11.0889	13.2
6	160	50	-1	-5.33	1	28.4089	5.33
Total	966	332	0	0	216	215.3394	195.0

Slide 17 STAT 1104, UCL, J. van Dine

Correlation Coefficient

Example:

$$R(X, Y) = \frac{1}{N-1} \sum_{k=1}^N \left(\frac{x_k - \mu}{\sigma} \right) \left(\frac{y_k - \mu}{\sigma} \right)$$

$$\mu_x = \frac{966}{6} = 161 \text{ cm}, \quad \mu_y = \frac{332}{6} = 55 \text{ kg},$$

$$\sigma_x = \sqrt{\frac{216}{5}} = 6.573, \quad \sigma_y = \sqrt{\frac{215.3}{5}} = 6.563,$$

$$\text{Corr}(X, Y) = R(X, Y) = 0.904$$

Slide 18 STAT 1104, UCL, J. van Dine

Correlation Coefficient - Properties

Correlation is invariant w.r.t. linear transformations of X or Y

$$R(X, Y) = \frac{1}{N-1} \sum_{k=1}^N \left(\frac{x_k - \mu_x}{\sigma_x} \right) \left(\frac{y_k - \mu_y}{\sigma_y} \right) = R(aX + b, cY + d), \text{ since}$$

$$\left(\frac{ax_k + b - \mu_{ax+b}}{\sigma_{ax+b}} \right) = \left(\frac{ax_k + b - (a\mu_x + b)}{a \times \sigma_x} \right) = \left(\frac{a(x_k - \mu_x) + b - b}{a \times \sigma_x} \right) = \left(\frac{x_k - \mu_x}{\sigma_x} \right)$$

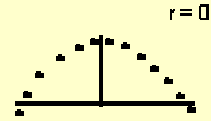
Slide 19 STAT 1104, UCL, Lee Dinger

Correlation Coefficient - Properties

Correlation is Associative

$$R(X, Y) = \frac{1}{N} \sum_{k=1}^N \left(\frac{x_k - \mu_x}{\sigma_x} \right) \left(\frac{y_k - \mu_y}{\sigma_y} \right) = R(Y, X)$$

Correlation measures linear association, NOT an association in general!!! So, Corr(X,Y) could be misleading for X & Y related in a non-linear fashion.

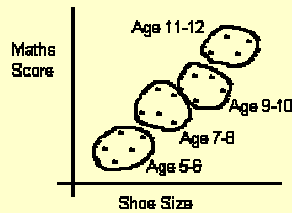


Slide 20 STAT 1104, UCL, Lee Dinger

Correlation Coefficient - Properties

$$R(X, Y) = \frac{1}{N} \sum_{k=1}^N \left(\frac{x_k - \mu_x}{\sigma_x} \right) \left(\frac{y_k - \mu_y}{\sigma_y} \right) = R(Y, X)$$

1. R measures the extent of linear association between two continuous variables.
2. Association does not imply causation - both variables may be affected by a third variable - age was a confounding variable.



Slide 21 STAT 1104, UCL, Lee Dinger

Essential Points

6. If the experimenter has control of the levels of X used, how should these levels be allocated to the available experimental units?

At random! Example, testing **hardness of concrete**, Y, based on **levels of cement**, X, incorporated. Factors effecting Y: amount of H₂O, ratio stone-chips to sand, drying conditions, etc. To prevent uncontrolled differences in batches of concrete in confounding our impression of cement effects, we should choose which batch (H₂O levels, sand, dry-conditions) gets what amount of cement at random! Then investigate for X-effects in Y observations. If some significance test indicates observed trend is significantly different from a random pattern → we have evidence of causal relationship, which may strengthen even further if the results are replicable.

Slide 22 STAT 1104, UCL, Lee Dinger

Essential Points

7. What theories can you explore using regression methods?

Prediction, explanation/causation, testing a scientific hypothesis/mathematical model:

- a. **Hooke's spring law**: amount of stretch in a spring, Y, is related to the applied weight X by $Y = \alpha + \beta X$, a, b are spring constants.
- b. **Theory of gravity**: force of gravity F between 2 objects is given by $F = \alpha/D^\beta$, where D=distance between objects, a is a constant related to the masses of the objects and $\beta = 2$, according to the **inverse square law**.
- c. **Economic production function**: $Q = \alpha L^\beta K^\gamma$, Q=production, L=quantity of labor, K=capital, α, β, γ are constants specific to the market studied.

Slide 23 STAT 1104, UCL, Lee Dinger

Essential Points

8. People fit theoretical models to data for three main purposes.

- a. To test the model, itself, by checking if the data is reasonably close agreement with the relationship predicted by the model.
- b. Assuming the model is correct, to test if theoretically specified values of a parameter are consistent with the data ($y=2x+1$ vs. $y=2.1x-0.9$).
- c. Assuming the model is correct, to estimate unknown constants in the model so that the relationship is completely specified ($y=ax+5$, $a=?$)

Slide 24 STAT 1104, UCL, Lee Dinger

Trend and Scatter - Computer timing data

- The major components of a regression relationship are **trend** and **scatter** around the trend.
- To investigate a trend – fit a math function to data, or smooth the data.
- Computer timing data: a mainframe computer has X users, each running jobs taking Y min time. The main CPU swaps between all tasks. Y* is the total time to finish all tasks. **Both Y and Y* increase with increase of tasks/users, but how?**

X = Number of terminals:	40	50	60	45	40	10	30	20
Y* = Total Time (mins):	6.6	14.9	18.4	12.4	7.9	0.9	5.5	2.7
Y = Time Per Task (secs):	9.9	17.8	18.4	16.5	11.9	5.5	11	8.1
X = Number of terminals:	50	30	65	40	65	65		
Y* = Total Time (mins):	12.6	6.7	23.6	9.2	20.2	21.4		
Y = Time Per Task (secs):	15.1	13.3	21.8	13.8	18.6	19.8		

Slide 25 STAT 110A, UCLA, Jon Dinger

Trend and Scatter - Computer timing data

Y* = Total time (min)

X = Number of terminals

Y = Time per task (s)

X = Number of terminals

We want to find reasonable models (descriptions) for these data!

Slide 26 STAT 110A, UCLA, Jon Dinger

Equation for the straight line – linear/affine function

β_0 =Intercept (the y-value at x=0)
 β_1 =Slope of the line (rise/run), change of y for every unit of increase for x.

Slide 27 STAT 110A, UCLA, Jon Dinger

The quadratic curve

Quadratic Curve

β_2 positive

β_2 negative

$Y = \beta_0 + \beta_1 x + \beta_2 x^2$

Slide 28 STAT 110A, UCLA, Jon Dinger

The quadratic curve

Segments of the curve

$Y = \beta_0 + \beta_1 x + \beta_2 x^2$

Slide 29 STAT 110A, UCLA, Jon Dinger

The exponential curve, $y = a e^{bx}$

b positive

b negative

Used in population growth/decay models.

Slide 30 STAT 110A, UCLA, Jon Dinger

Effects of changing x for different functions/curves

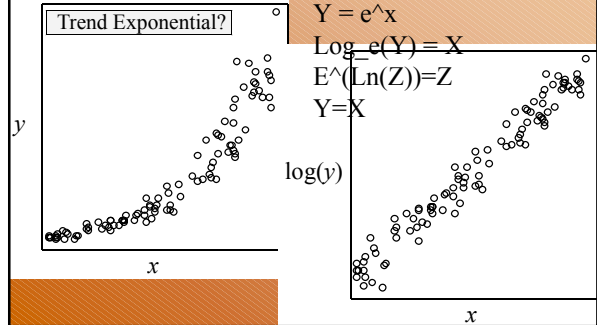
A straight **line** changes by a fixed **amount** with each unit change in x .

An **exponential** changes by a fixed **percentage** with each unit change in x .

Slide 31 STAT 110A, UCLA, Jon Dinger

To tell whether a trend is exponential

check whether a plot of **log(y) versus x** has a linear trend.

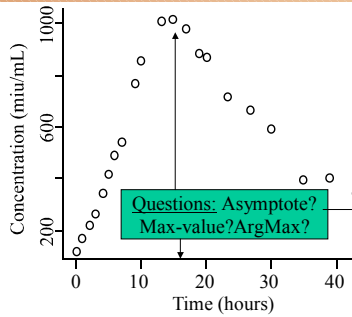


Creatine kinase concentration in patient's blood

You should not let the questions you want to ask be dictated by the tools you know how to use.

Here Y = creatine kinase concentration in blood for a set of heart attack patients vs. the time, X .

No symmetry so X^2 models won't work!



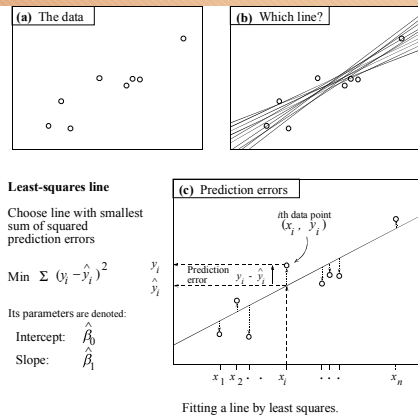
Slide 33 STAT 110A, UCLA, Jon Dinger

Comments

1. In statistics what are the **two main approaches** to summarizing **trends** in data? (model fitting; smoothing – done by the eye!)
2. In $y = 5x + 2$, what information do the 5 and the 2 convey? (slope, y-intercept)
3. In $y = 7 + 5x$, what change in y is associated with a 1-unit increase in x ? with a 10-unit increase? (5; 50)
How about for $y = 7 - 5x$. (-5; -50)
5. How can we tell whether a trend in a scatter plot is exponential? (plot $\log(Y)$ vs. X , should be linear)

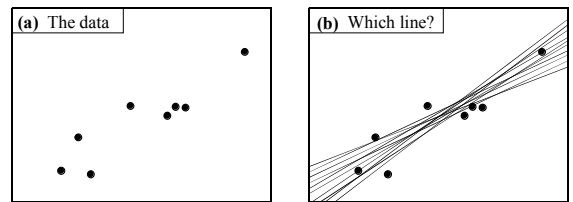
Slide 34 STAT 110A, UCLA, Jon Dinger

Choosing the "best-fitting" line



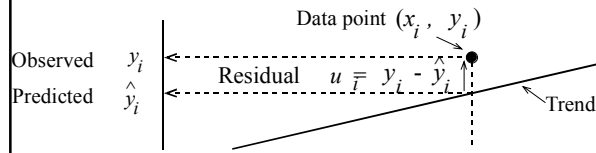
Slide 35 STAT 110A, UCLA, Jon Dinger

Fitting a line through the data



Slide 36 STAT 110A, UCLA, Jon Dinger

The idea of a residual or prediction error



Slide 37 STAT 110A, UCLA, Jon Dinger

Least squares criterion

Least squares criterion: Choose the values of the parameters to *minimize the sum of squared prediction errors* (or sum of squared residuals),

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

Slide 38 STAT 110A, UCLA, Jon Dinger

The least squares line

Least-squares line

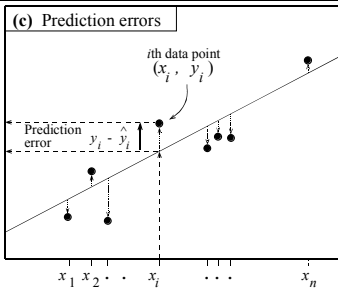
Choose line with smallest sum of squared prediction errors

$$\text{Min } \sum (y_i - \hat{y}_i)^2$$

Its parameters are denoted:

Intercept: $\hat{\beta}_0$

Slope: $\hat{\beta}_1$



Least-squares line: $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

Slide 39 STAT 110A, UCLA, Jon Dinger

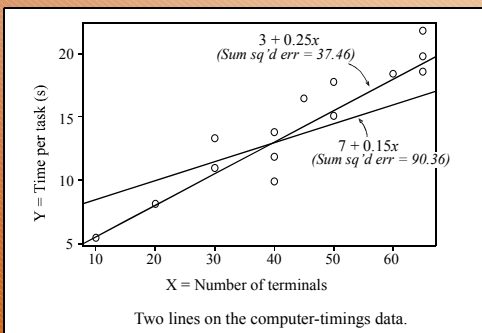
The least squares line

Least-squares line: $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n [(x_i - \bar{x})(y_i - \bar{y})]}{\sum_{i=1}^n (x_i - \bar{x})^2}; \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

Slide 40 STAT 110A, UCLA, Jon Dinger

Computer timings data – linear fit

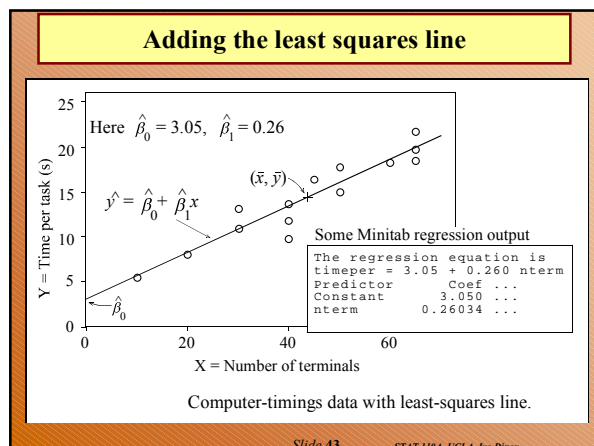


Slide 41 STAT 110A, UCLA, Jon Dinger

Computer timings data

Prediction Errors		Computer timings data			
		$3 + 0.25x$		$7 + 0.15x$	
x	y	\hat{y}	$y - \hat{y}$	\hat{y}	$y - \hat{y}$
40	9.90	13.00	-3.10	13.00	-3.10
50	17.80	15.50	2.30	14.50	3.30
60	18.40	18.00	0.40	16.00	2.40
45	16.50	14.25	2.25	13.75	2.75
40	11.90	13.00	-1.10	13.00	-1.10
10	5.50	5.50	0.00	8.50	-3.00
30	11.00	10.50	0.50	11.50	-0.50
20	8.10	8.00	0.10	10.00	-1.90
50	15.10	15.50	-0.40	14.50	0.60
30	13.30	10.50	2.80	11.50	1.80
65	21.80	19.25	2.55	16.75	5.05
40	13.80	13.00	0.80	13.00	0.80
65	18.60	19.25	-0.65	16.75	1.85
65	19.80	19.25	0.55	16.75	3.05
Sum of squared errors			37.46		90.36

Slide 42 STAT 110A, UCLA, Jon Dinger



Hands – on worksheet !

1. $X = \{-1, 2, 3, 4\}$, $Y = \{0, -1, 1, 2\}$,

X	Y	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})^2$	$(y - \bar{y})^2$	$\frac{(x - \bar{x})x}{(y - \bar{y})}$
-1	0					
2	-1					
3	1					
4	2					

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n [(x_i - \bar{x})(y_i - \bar{y})]}{\sum_{i=1}^n (x_i - \bar{x})^2}; \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

Slide 44 STAT 1104, UCLA, Jon Dinger

Hands – on worksheet !

1. $X = \{-1, 2, 3, 4\}$, $Y = \{0, -1, 1, 2\}$, $\bar{x} = 2$, $\bar{y} = 0.5$

X	Y	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})^2$	$(y - \bar{y})^2$	$\frac{(x - \bar{x})x}{(y - \bar{y})}$
-1	0	-3	-0.5	9	0.25	1.5
2	-1	0	-1.5	0	2.25	0
3	1	1	0.5	1	0.25	0.5
4	2	2	1.5	4	2.25	3
2	0.5			14	5	5

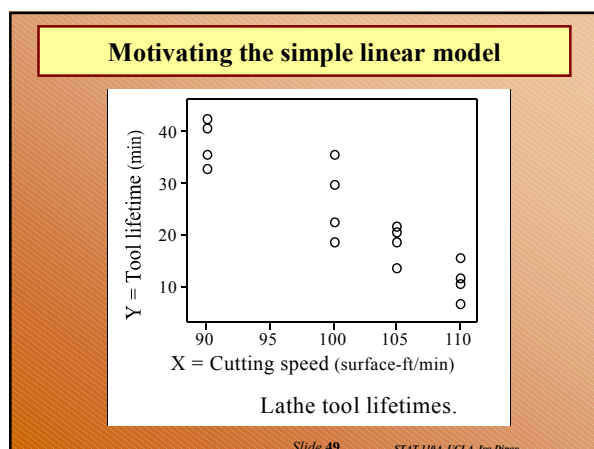
$\hat{\beta}_1 = 5/14$
 $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = 0.5 - 10/14$

Slide 45 STAT 1104, UCLA, Jon Dinger

Review

1. What are the quantities that specify a particular line?
2. Explain the idea of a prediction error in the context of fitting a line to a scatter plot. To what visual feature on the plot does a prediction error correspond? (scatter-size)
3. What property is satisfied by the line that fits the data best in the least-squares sense?
4. The **least-squares line** $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$ passes through the points $(x = 0, \hat{y} = ?)$ and $(x = \bar{x}, \hat{y} = ?)$. Supply the missing values.

Slide 48 STAT 1104, UCLA, Jon Dinger



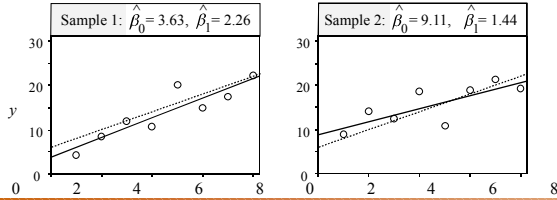
The simple linear model

(a) The simple linear model (b) Data sampled from the model

When $X = x$, $Y \sim \text{Normal}(\mu_y, \sigma)$ where $\mu_y = \beta_0 + \beta_1 x$, OR
 when $X = x$, $Y = \beta_0 + \beta_1 x + U_x$ where $U \sim \text{Normal}(0, \sigma)$

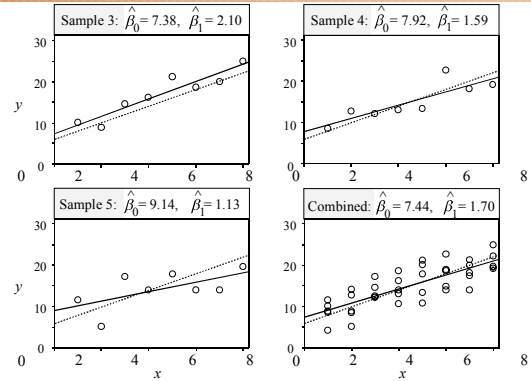
Slide 50 STAT 1104, UCLA, Jon Dinger

Data generated from $Y = 6 + 2x + \text{error}(U)$
 Dotted line is true line and
 solid line — is the data-estimated LS line.
Note differences between true $\beta_0=6, \beta_1=2$ and
 their estimates $\hat{\beta}_0^\wedge$ & $\hat{\beta}_1^\wedge$.



Slide 51 STAT 110A, UCLA, Jon Dimez

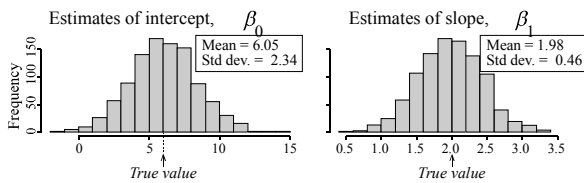
Data generated from $Y = 6 + 2x + \text{error}(U)$



Slide 52 STAT 110A, UCLA, Jon Dimez

Data generated from $Y = 6 + 2x + \text{error}(U)$

Histograms of least-squares estimates from 1,000 data sets



Data generated from the model $Y = 6 + 2x + U$
 where $U \sim \text{Normal}(\mu = 0, \sigma = 3)$.

Slide 53 STAT 110A, UCLA, Jon Dimez

Summary

For the simple linear model, *least-squares estimates are unbiased* [$E(\hat{\beta}) = \beta$] and *Normally distributed*.

Noisier data produce *more-variable* least-squares estimates.

Slide 54 STAT 110A, UCLA, Jon Dimez

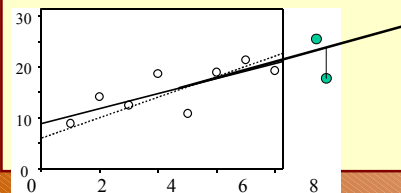
Summary

1. Before considering using the simple linear model, what sort of pattern would you be looking for in the scatter plot? (linear trend with constant scatter spread across the range of X)
2. What assumptions are made by the simple linear model, **SLM**? (X is linearly related to the mean value of the Y obs's at each X, $\mu_y = \beta_0 + \beta_1 x$; where β_0 & β_1 are the true values of the intercept and slope of the SLM; The LS estimates $\hat{\beta}_0^\wedge$ & $\hat{\beta}_1^\wedge$ estimate the true values of β_0 & β_1 ; and the random errors $U = Y - \mu_y \sim N(\mu, \sigma)$.)
3. If the simple linear model holds, what do you know about the sampling distributions of the least-squares estimates? (Unbiased and Normally distributed)

Slide 55 STAT 110A, UCLA, Jon Dimez

Summary

4. In the simple linear model, what behavior is governed by σ ? (the spread of scatter of the data around trend)
5. Our estimate of σ can be thought of as a sample standard deviation for the set of **prediction errors** from the **least-squares line**.



Slide 56 STAT 110A, UCLA, Jon Dimez

RMS Error for regression

- Error = Actual value - Predicted value

$Y = \beta_0 + \beta_1 X$

- The RMS Error for the regression line $Y = \beta_0 + \beta_1 X$ is

$$\sqrt{\frac{(y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2 + (y_3 - \hat{y}_3)^2 + (y_4 - \hat{y}_4)^2 + (y_5 - \hat{y}_5)^2}{5-1}}$$

where $\hat{y}_k = \beta_0 + \beta_1 x_k, \quad 1 \leq k \leq 5$

Slide 57 STAT 1104, UCL4, Jon Dinger

Compute the RMS Error for this regression line

- Error = Actual value - Predicted value

X	Y
1	9
2	15
3	12
4	19
5	11
6	20
7	22
8	18

- The RMS Error for the regression line $Y = \beta_0 + \beta_1 X$ is

$$\sqrt{\frac{(y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2 + (y_3 - \hat{y}_3)^2 + (y_4 - \hat{y}_4)^2 + (y_5 - \hat{y}_5)^2}{5-1}}$$

where $\hat{y}_k = \beta_0 + \beta_1 x_k, \quad 1 \leq k \leq 5$

Slide 58 STAT 1104, UCL4, Jon Dinger

Compute the RMS Error for this regression line

- Error = Actual value - Predicted value

X	Y
1	9
2	15
3	12
4	19
5	11
6	20
7	22
8	18

- The RMS Error for the regression line $Y = \beta_0 + \beta_1 X$ is

$$\sqrt{\frac{(y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2 + (y_3 - \hat{y}_3)^2 + (y_4 - \hat{y}_4)^2 + (y_5 - \hat{y}_5)^2}{5-1}}$$

where $\hat{y}_k = \beta_0 + \beta_1 x_k, \quad 1 \leq k \leq 5$

- First compute the LS linear fit (estimate $\beta_0^{\wedge} + \beta_1^{\wedge}$)
- Then Compute the individual errors
- Finally compute the cumulative RMS measure.

Slide 59 STAT 1104, UCL4, Jon Dinger

Compute the RMS Error for this regression line

- First compute the LS linear fit (estimate $\beta_0^{\wedge} + \beta_1^{\wedge}$), $\mu_x = 4.5, \mu_y = 15.75$

X	Y	X - μ_x	Y - μ_y	(X - μ_x) ²	(Y - μ_y) ²	(X - μ_x) ² * (Y - μ_y) ²
1	9					
2	15					
3	12					
4	19					
5	11					
6	20					
7	22					
8	18					

Total:

- Compute

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n [(x_i - \bar{x})(y_i - \bar{y})]}{\sum_{i=1}^n (x_i - \bar{x})^2}; \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

Slide 60 STAT 1104, UCL4, Jon Dinger

Compute the RMS Error for this regression line

- Then Compute the individual errors

X	Y
1	9
2	15
3	12
4	19
5	11
6	20
7	22
8	18

$(y_k - \hat{y}_k)^2$, where $\hat{y}_k = \beta_0 + \beta_1 x_k, \quad 1 \leq k \leq 8$

- Finally compute the cumulative RMS measure.

$$\sqrt{\frac{(y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2 + (y_3 - \hat{y}_3)^2 + (y_4 - \hat{y}_4)^2 + (y_5 - \hat{y}_5)^2}{5-1}}$$

where $\hat{y}_k = \beta_0 + \beta_1 x_k, \quad 1 \leq k \leq 5$

- Note on the Correlation coefficient formula,**

$$R(X, Y) = \frac{1}{N-1} \sum_{k=1}^N \left(\frac{x_k - \mu}{\sigma_x} \right) \left(\frac{y_k - \mu}{\sigma_y} \right)$$

Slide 61 STAT 1104, UCL4, Jon Dinger

Compute the RMS Error for this regression line

- The RMS Error for the regression line $Y = \beta_0 + \beta_1 X$ says how far away from the (model/predicting) regression line is each observation.
- Observe that the SD(Y) is also a RMS Error measure of another specific line - horizontal line through the average of the Y values. This line may also be taken for a regression line, but often it's not the best linear fit.

$SD(Y) = \sqrt{\frac{1}{N-1} \sum_{i=1}^N (Y_i - \bar{Y})^2}$ vs.

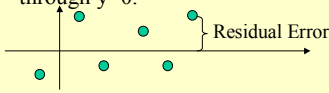
$RMS E(Y, \hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X) = \sqrt{\frac{1}{N-1} \sum_{i=1}^N (Y_i - \hat{Y}_i)^2}$

Predicted vs. Observed

Slide 62 STAT 1104, UCL4, Jon Dinger

Plotting the Residuals

- The Residuals=Observed –Predicted for the regression line $Y = \beta_0 + \beta_1 X$ (just like the error).
- Residuals average to zero, mathematically, and the regression line for the residuals is a horizontal line through $y=0$.



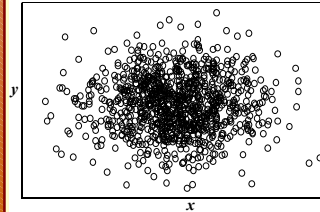
When $X = x$, $Y \sim \text{Normal}(\mu_y, \sigma)$ where $\mu_y = \beta_0 + \beta_1 x$, **OR**
 when $X = x$, $Y = \beta_0 + \beta_1 x + U$, where $U \sim \text{Normal}(0, \sigma)$

Random error

Slide 63 STAT 110A, UCLA, Jon Dineen

Plotting the Residuals – patterns?

- The Residuals=Observed –Predicted for the regression line $Y = \beta_0 + \beta_1 X + U$ should show no clear trend or pattern, for our linear model to be a good and useful approximation to the unknown process.



Slide 64 STAT 110A, UCLA, Jon Dineen

Inference – just a glance at statistical inference

- The regression **intercept** β_0 and **slope** β_1 are usually called **regression coefficients**
 - The least squares estimates of their values are found in the coefficients column of program printouts
- **Confidence intervals** for a true regression coefficient (whether intercept or slope) is given by
estimated coefficient \pm *t std errors*
- **t-test statistic** $df = n - 2$

$$t_0 = \frac{\text{estimated coefficient} - \text{hypothesized value}}{\text{standard error}}$$

Slide 65 STAT 110A, UCLA, Jon Dineen

Inferences

- **Confidence intervals** for a true regression coefficient (whether intercept or slope) is given by
estimated coefficient \pm *t std errors*

$$\hat{\beta}_1 \pm t \text{SE}(\hat{\beta}_1)$$

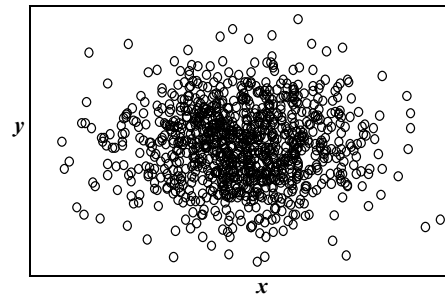
- **t-test statistic** $df = n - 2$ **Ho:** $\beta_1 = c$

$$t_0 = \frac{\hat{\beta}_1 - c}{\text{SE}(\hat{\beta}_1)}$$

Slide 66 STAT 110A, UCLA, Jon Dineen

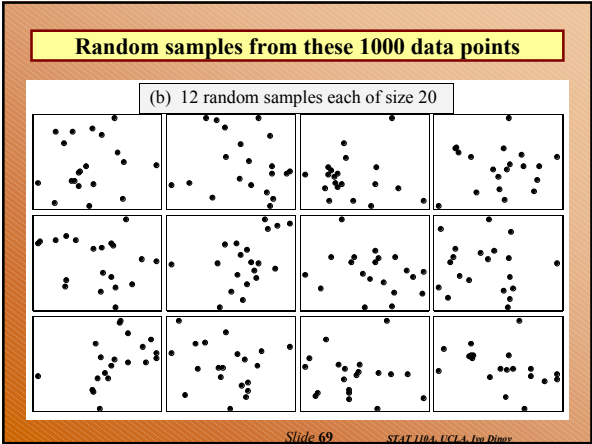
Is there always an X Y relationship? Linear Relationship ?

(a) 1000 data points with no relationship between X and Y



Slide 68 STAT 110A, UCLA, Jon Dineen

Slide 67 STAT 110A, UCLA, Jon Dineen



Testing for **no linear relationship – trend of Y w.r.t. X is trivial!**

$H_0: \text{true slope} = 0$

OR

$H_0: \beta_1 = 0$

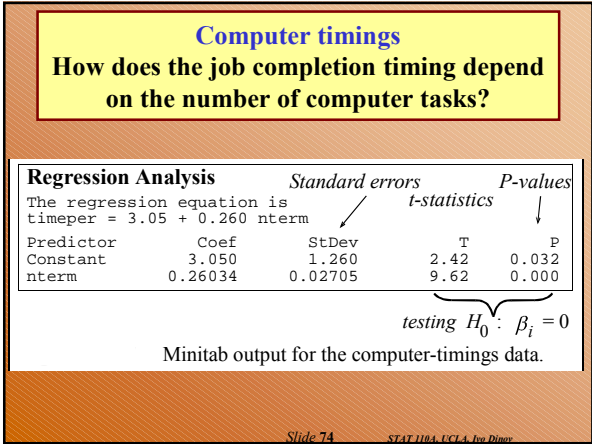
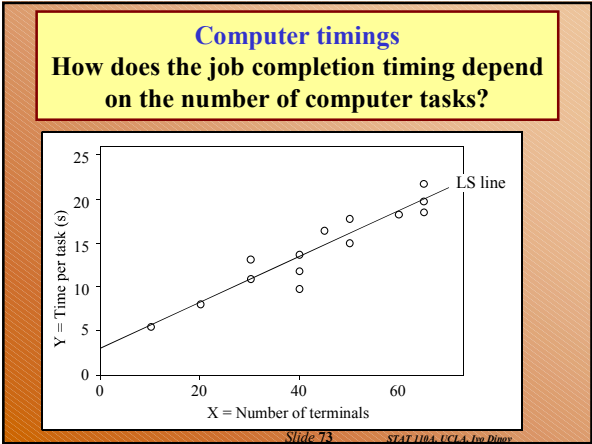
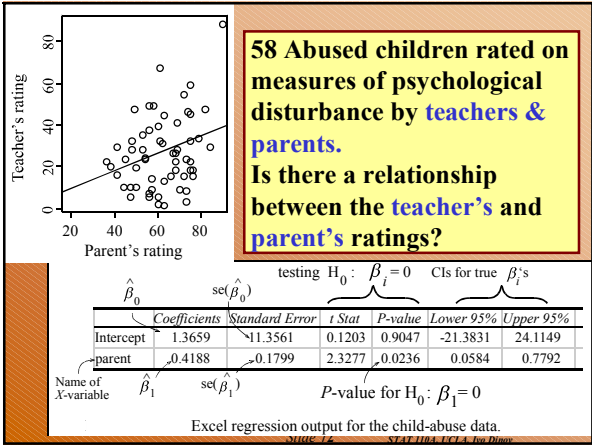
Slide 70 STAT 110A, UCLA, Jon Dinger

58 Abused children rated on measures of psychological disturbance by teachers & parents. Is there a relationship between the teacher's and parent's ratings?

H_0 : parent's and teacher's ratings are identical
 $H_0: \beta_1 = 1$, $df = 58 - 2 = 56$,

H_0 : No relation between parent's and teacher's ratings. $H_0: \beta_1 = 0$, $df = 58 - 2 = 56$,

Slide 71 STAT 110A, UCLA, Jon Dinger



CI for true slope

Predictor	Coef	Standard errors StDev	t-statistics T	P-values P
Constant	3.050	1.260	2.42	0.032
nterm	0.26034	0.02705	9.62	0.000

The regression equation is
timeper = 3.05 + 0.260 nterm

testing $H_0: \beta_i = 0$

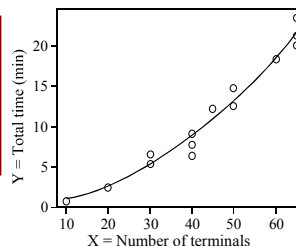
Minitab output for the computer-timings data.

For a 95% CI with $df = n - 2 = 12$, $t = 2.179$

CI: estimate $\pm t$ std errors
= $0.26034 \pm 2.179 \times 0.02705 = [0.20, 0.32]$

Slide 75 STAT 110A, UCLA, Jon Dinger

Computer timings: Is the trend for $Y = \text{Total time}$ curved?



	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.215067	1.941166	0.111	0.91378
nterm	0.036714	0.100780	0.364	0.72254
ntermsq	0.004526	0.001209	3.745	0.00324 **

Quadratic model for $Y^* = \text{Total Time}$.

Slide 76 STAT 110A, UCLA, Jon Dinger

Remarks

1. What value of df is used for inference for $\hat{\beta}_0$ and $\hat{\beta}_1$?
2. Within the context of the simple linear model, what formal hypothesis is tested when you want to test for no linear relationship between X and Y ?
3. What hypotheses do the t -test statistics and associated P -values on regression output test?
4. What is the form of a confidence interval for the true slope?
5. What is the form of the test statistic for testing $H_0: \beta_1 = c$?

Slide 77 STAT 110A, UCLA, Jon Dinger

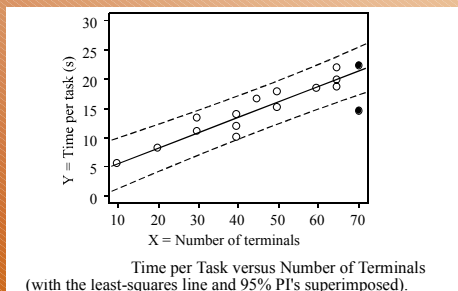
Prediction

Predicting at $X = x_p$

- The **confidence interval for the mean** estimates the **average** Y -value at $X = x_p$.
 - (averaged over many repetitions of the experiment.)
- The **prediction interval (PI)** tries to predict the next **actual** Y -value at x_p , in the future.

Slide 78 STAT 110A, UCLA, Jon Dinger

Predicting time-per-task for 70 terminals



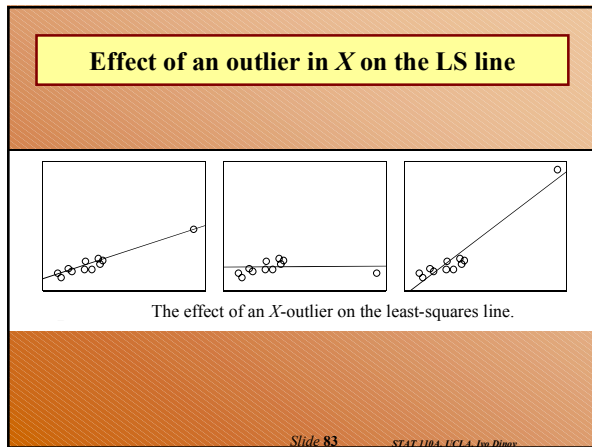
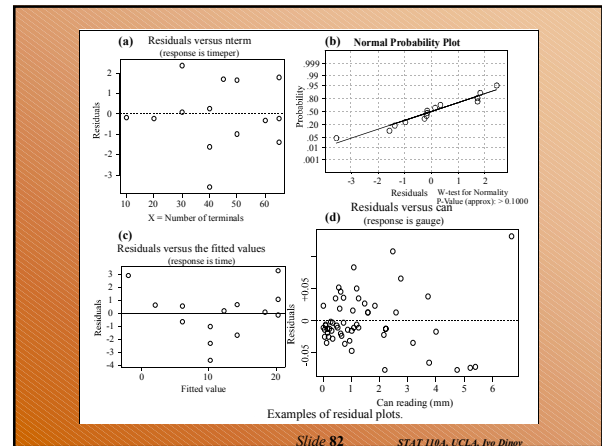
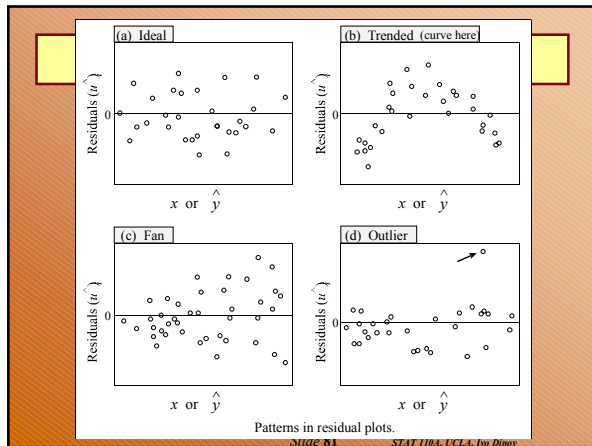
Time per Task versus Number of Terminals
(with the least-squares line and 95% PI's superimposed).

Slide 79 STAT 110A, UCLA, Jon Dinger

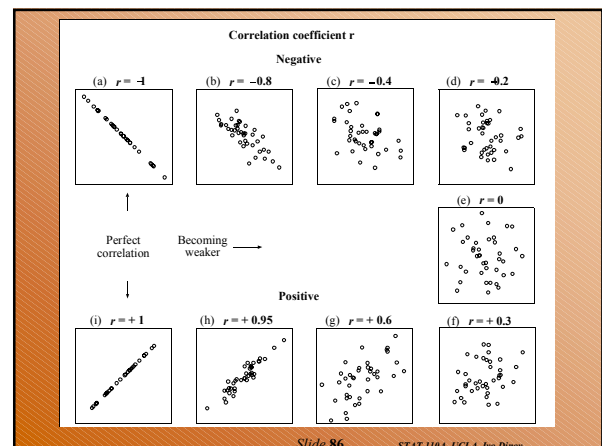
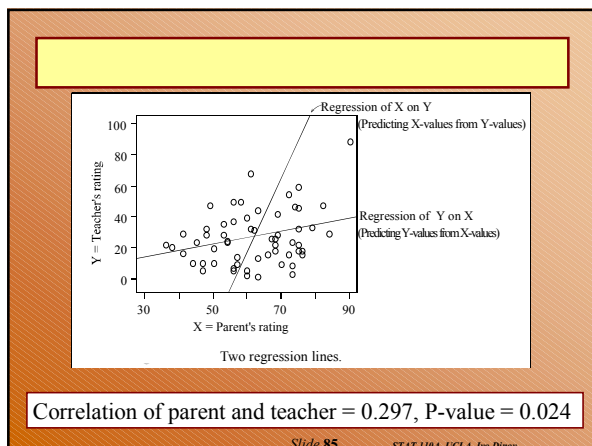
Review

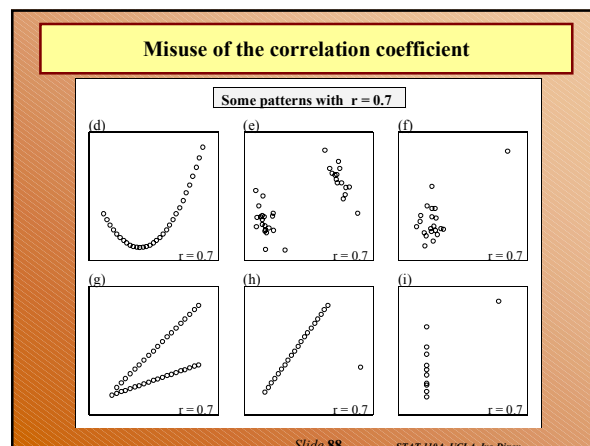
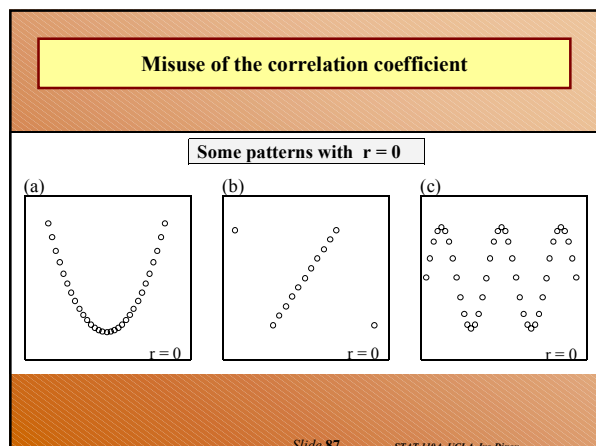
1. What is the difference between a confidence interval for the mean and a prediction interval?
2. Prediction intervals make allowances for two sources of uncertainty. What are they? How does a confidence interval for the mean differ in this regard?
3. At what point along the X -axis are these intervals narrowest?
4. We gave some general warnings about prediction earlier. They are relevant here as well. What were those warnings?

Slide 80 STAT 110A, UCLA, Jon Dinger



- ### Review
1. What assumptions are made by the simple linear model?
 2. Which assumptions are critical for all types of inference?
 3. What types of inference are relatively robust against departures from the Normality assumption?
 4. Four types of residual plot were described. What were they, and what can we learn from each?
 5. What is an outlier in X, and why do we have to be on the lookout for such observations?
- Slide 84 STAT 110A, UCLA, Jon Dinger





Correlation does not necessarily imply causation.

Slide 89 STAT 110A, UCLA, Jon Dineen

- ### Review
1. Describe a fundamental difference between the way regression treats data and the way correlation treats data.
 2. What is the correlation coefficient intended to measure?
 3. For what shape(s) of trend in a scatter plot does it make sense to calculate a correlation coefficient?
 4. What is the meaning of a correlation coefficient of $r = +1$? $r = -1$? $r = 0$?
- Slide 90* STAT 110A, UCLA, Jon Dineen

Summary

Slide 91 STAT 110A, UCLA, Jon Dineen

- ### Concepts
- Relationships between quantitative variables should be explored using **scatter plots**.
 - Usually the Y variable is continuous
(or behaves like one in that there are few repeated values)
 - and the X variable is discrete or continuous.
 - **Regression** singles out one variable (Y) as the response and uses the explanatory variable (X) to explain or predict its behavior.
 - **Correlation** treats both variables symmetrically.
- Slide 92* STAT 110A, UCLA, Jon Dineen

Concepts cont'd

In practical problems, regression models may be fitted for any of the following reasons:

- To understand a **causal relationship** better.
- To find relationships which may be **causal**.
- To make **predictions**.
 - But be cautious about predicting outside the range of the data
- To **test theories**.
- To **estimate parameters** in a theoretical model.

Slide 93 STAT 110A, UCLA, Jon Dineen

Concepts cont'd

- In observational data, strong relationships are not necessarily causal.
- We can only have reliable evidence of causation from controlled experiments.
- Be aware of the possibility of **lurking** variables which may effect both X and Y .

Slide 94 STAT 110A, UCLA, Jon Dineen

Concepts cont'd

- Two important trend curves are the **straight line** and the **exponential curve**.
 - A straight line changes by a *fixed amount* with each unit change in x .
 - An exponential curve changes by a *fixed percentage* with each unit change in x .
- You should not let the questions you want to ask of your data be dictated by the tools you know how to use. You can always ask for help.

Slide 95 STAT 110A, UCLA, Jon Dineen

Concepts cont'd

- The two main approaches to summarizing trends in data are using *smoothers* and *fitting mathematical curves*.
- The **least-squares criterion** for fitting a mathematical curve is to choose the values of the parameters (e.g. β_0 and β_1) to minimize the sum of squared prediction errors, $\sum (y_i - \hat{y}_i)^2$.

Slide 96 STAT 110A, UCLA, Jon Dineen

Linear Relationship

- We fit the linear relationship $\hat{y} = \beta_0 + \beta_1 x$.
- The slope β_1 is the change in \hat{y} associated with a one-unit increase in x .

Least-squares estimates

- The least-squares estimates, $\hat{\beta}_0$ and $\hat{\beta}_1$ are chosen to minimize $\sum (y_i - \hat{y}_i)^2$.
- The **least-squares regression line** is $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$.

Slide 97 STAT 110A, UCLA, Jon Dineen

Model for statistical inference

- Our theory assumes the model $Y_i = \beta_0 + \beta_1 x_i + U_i$,
- where the random errors, U_1, U_2, \dots, U_n , are a random sample from a Normal(0, σ) distribution.
- This means that the random errors
 - are Normally distributed (each with mean 0),
 - all have the same standard deviation σ regardless of the value of x , and
 - are all independent.

Slide 98 STAT 110A, UCLA, Jon Dineen

Residuals and outliers

- These assumptions should be checked using residual plots (Section 12.4.4). The i th *residual* (or *prediction error*) is

$$y_i - \hat{y}_i = \text{observed} - \text{predicted.}$$

- An **outlier** is a data point with an unexpectedly large residual (positive or negative).

Slide 99 STAT 1104, UCL4, Jon Dinger

Inference

- Inferences for the intercept and slope** are just as in Chapters 8 and 9, with confidence intervals being of the form $\text{estimate} \pm t \text{ std errors}$ and test statistics of the form

$$t_0 = (\text{estimate} - \text{hypothesized value}) / \text{standard error.}$$

- We use $df = n - 2$.
- To test for **no linear association**, we test $H_0: \beta_1 = 0$.

Slide 100 STAT 1104, UCL4, Jon Dinger

*Prediction

- The predicted value for a new Y at $X = x_p$ is

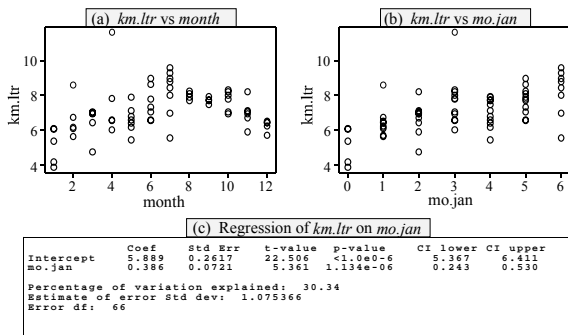
$$\hat{y}_p = \hat{\beta}_0 + \hat{\beta}_1 x_p$$
- The *confidence interval for the mean* estimates the average Y -value at $X = x_p$.
 - averaged over many repetitions of the experiment.
- The *prediction interval* tries to predict the next *actual* Y -value at $X = x_p$.
- The prediction interval is wider than the corresponding confidence interval for the mean.

Slide 101 STAT 1104, UCL4, Jon Dinger

Correlation coefficient

- The **correlation coefficient** r is a measure of linear association with $-1 \leq r \leq 1$.
- If $r = 1$, then X and Y have a perfect positive linear relationship.
- If $r = -1$, then X and Y have a perfect negative linear relationship.
- If $r = 0$, then there is no linear relationship between X and Y .
- Correlation does not necessarily imply causation.

Slide 102 STAT 1104, UCL4, Jon Dinger



Fuel consumption data.

Slide 103 STAT 1104, UCL4, Jon Dinger

Regression of Log(price) on Age

	Coef	Std Err	t-value	p-value	CI lower	CI upper
Intercept	3.8511	0.0494	78.02	0	---	---
age	-0.2164	0.0095	-22.67	0	-0.24	-0.20

Percent of variation explained: 90.02
 Estimate of error Std dev: 0.2433205
 Error df: 57

Age	0	1	2	3	4	5	6	7	8	9	10
Predicted	3.85	3.63	3.42	3.20	2.99	2.77	---	2.34	2.12	1.90	1.69
Pred lower	3.35	3.14	2.93	2.71	2.49	2.28	---	1.84	1.62	1.40	1.18
Pred upper	4.35	4.13	3.91	3.69	3.48	3.26	---	2.83	2.62	2.40	2.19

From *Chance Encounters* by C.J. Wild and G.A.F. Seber. © John Wiley & Sons, 2000.

Slide 104 STAT 1104, UCL4, Jon Dinger