

UCLA STAT 110 A 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

Chapter 2: Data Summaries, Plots

- Types of variables
- Presentation of data
- Simple plots
- Numerical summaries
- Repeated and grouped data
- Qualitative variables

STAT 110A, UCLA, Ivo Dinov Slide 2

TABLE 2.1.1 Data on Male Heart Attack Patients

A subset of the data collected at a Hospital is summarized in this table. Each patient has measurements recorded for a number of variables – ID, Ejection factor (ventricular output), blood systolic/diastolic pressure, etc.

- Reading the table
- Which of the measured variables (age, ejection etc.) are useful in predicting how long the patient may live.
- Are there relationships between these predictors?
- variability & noise in the observations hide the message of the data.

Slide 3 STAT 110A, UCLA, Ivo Dinov

Data on Male Heart Attack Patients												
ID	EJEC	SYS- VOL	DIA- VOL	OCCLU	STEN	TIME	COME	AGE	SMOKE	BETA	CHOL	SURV
390	72	36	131	0	0	143	0	49	2	2	59	0
279	52	74	155	37	63	140	0	54	2	2	66	1
391	62	52	137	33	47							
201	50	165	329	33	30							
202	50	47	95	0	100							
69	27	124	170	77	23							
310	60	86	215	7	50							
392	72	37	132	40	10							
311	60	65	163	0	40							
288	59	39	94	0	0							
407	67	39	117	0	73							

NA = Not Available (missing data code)

Ivo Dinov

Data & Variables

- **Variable** is the name (label) given to the object being measured, counted, observed or recorded in any way. E.g., ID, EjectionVolume, Sys/Dia pressure, etc.
- **Data** are the actual recording values. E.g., 120/80 (for the arterial pressure).

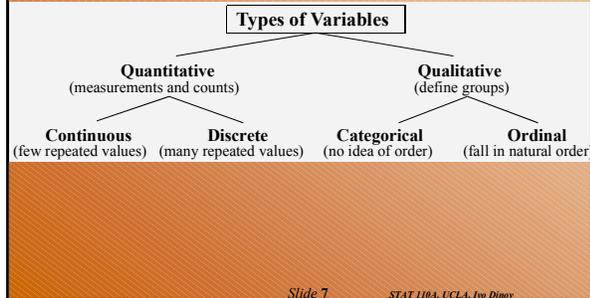
Slide 5 STAT 110A, UCLA, Ivo Dinov

Types of variable

- **Quantitative** variables are *measurements* and counts
 - Variables with *few repeated values* are treated as *continuous*.
 - Variables with *many repeated values* are treated as *discrete*
- **Qualitative** variables (a.k.a. *factors* or *class-variables*) describe *group membership*

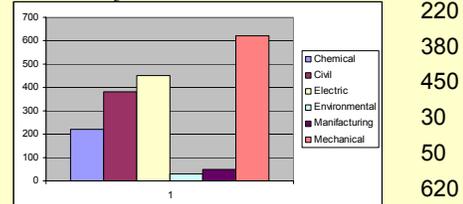
Slide 6 STAT 110A, UCLA, Ivo Dinov

Distinguishing between types of variable



Bar Chart

- List all possible categories the data is classified in!
 - Represents the frequency of occurrence of the data in each category
- Example: Number of engineering students enrolled in different majors:



Frequency Histograms

- Frequency Histograms - protocol:
 - Determine the **RANGE** of values $[a : b]$
 - Determine the **numbers of bars (bins)** to plot; d
 - Width** of each bar: $(b-a)/d$
 - Count the **frequency** of your data in each bin, subinterval
 - Draw** the histogram
- Example:

Slide 9 STAT 110A, UCLA, Jon Dinger

Frequency Histograms - Heights

- Example: Gender, Age, Heights (in inches):

GENDER	AGE	HEIGHT	GENDER	AGE	HEIGHT	GENDER	AGE	HEIGHT
M	35	95.6	F	16	77.0	F	98	128.3
F	78	112.7	F	5	64.5	M	6	63.0
M	85	122.8	F	38	91.7	F	59	103.5
M	30	91.1	M	70	107.9	M	16	77.5
M	5	68.6	F	68	112.7	F	91	127.3
M	57	93.1	F	90	127.2	M	14	69.7
F	7	63.8	F	81	121.4	M	35	91.0
F	100	144.0	M	95	130.9	M	51	102.2
M	45	97.5	M	15	73.0	F	83	126.7
M	77	120.1	M	72	112.5	M	30	86.0
F	16	70.6	M	30	89.2	F	91	120.8
F	37	93.1	F	60	108.2	M	93	125.5
F	41	93.8	F	17	79.1	M	53	102.8
M	6	67.0	F	37	91.0	F	15	72.5
F	6	65.0	M	92	140.9	M	31	91.2
M	51	104.9	F	5	58.3	F	88	114.0
M	25	86.3	M	2	43.5	M	32	97.1
M	11	73.5	M	41	90.4	F	15	87.0
M	76	108.0	M	19	76.5	M	17	67.5
						M	8	68.0

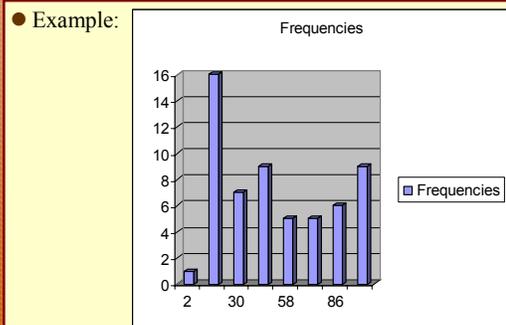
Slide 10 STAT 110A, UCLA, Jon Dinger

Frequency Histograms - Heights

- Frequency Histograms - protocol:
 - Determine the **RANGE** of values $[0 : 130]$
 - Determine the **numbers of bars (bins)** to plot; 8
 - Width** of each bar: $(130 - 0)/8 = 16.3$
 - Count the **frequency** of your data in each bin, subinterval
 - Draw** the histogram

Slide 11 STAT 110A, UCLA, Jon Dinger

Frequency Histograms - Heights



Density Histograms

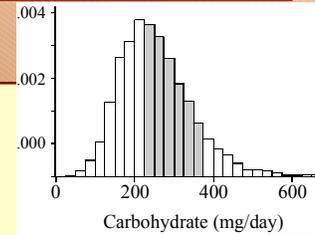
- The **height** of the histogram bar (bin):
 - f_i = frequency of the data
 - on the i -th interval,
 - w_i = width of the i -th interval
 - n – total number of data.
$$h_i = \frac{f_i}{n \times w_i}$$
- Then the **Area** of the histogram bar is: $A_i = \frac{f_i}{n}$
- And the Total area of all histogram bins is: **1 (100%)**

Slide 13 STAT 110A, UCLA, Ivo Dinov

Density (standardized) histograms

For a **standardized histogram**:

- The vertical scale is **Relative frequency / Interval width**
- Total area** under histogram = 1
- Proportion** of the data between a and b is the **area** under histogram between a and b



Slide 14 STAT 110A, UCLA, Ivo Dinov

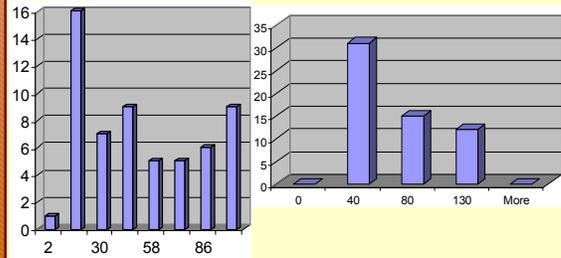
Uni- vs. Multi-modal histograms

- Number of clear humps on the frequency histogram plot determines the modality of a histogram plot.
- Note:** Modality of the histogram is histogram parameter specific! Changing the width of the bins changes its appearance!

Slide 15 STAT 110A, UCLA, Ivo Dinov

Uni- vs. Multi-modal histograms

- Number of clear humps on the frequency histogram plot determines the modality of a histogram plot.



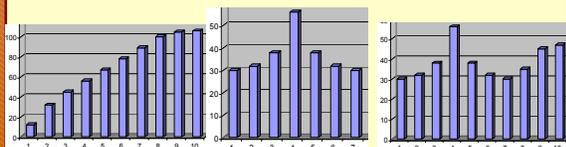
Slide 16 STAT 110A, UCLA, Ivo Dinov

Skewness & Symmetry of histograms

- A histogram is **symmetric** if the bars (bins) to the left of some point (mean) are approximately mirror images of those to the right of the mean.

file:///C:/Ivo.dir/UCLA_Classes/Applets.dir/HistogramApplet.html

- Histogram is **skewed** if it is not symmetric, the histogram is heavy to the left or right, or non-identical on both sides of the mean.



Slide 17 STAT 110A, UCLA, Ivo Dinov

Analyzing Histogram Plots

- Modality** – uni- vs. multi-modal (Why do we care?)
- Symmetry** – how skewed is the histogram?
- Center of gravity** for the Histogram plot – does it make sense?
- If center-of-gravity exists quantify the **spread of the frequencies** around this point.
- Strange patterns** – gaps, atypical frequencies lying away from the center.

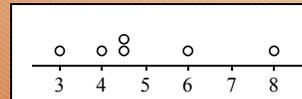
Slide 18 STAT 110A, UCLA, Ivo Dinov

Caution: Storing and Reporting data

- Round numbers for presentation
- Maintain complete accuracy in numbers to be used in calculations. If you need to round-off, this should be the very last operation ...

Slide 20 STAT 110A, UCLA, Jon Dinger

The dot plot



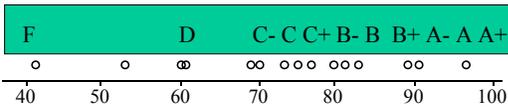
Dot plot.



Dot plot showing special features.

Slide 24 STAT 110A, UCLA, Jon Dinger

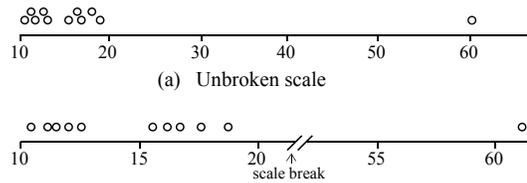
Example of exploiting gaps and clusters



Grading of a university course.

Slide 25 STAT 110A, UCLA, Jon Dinger

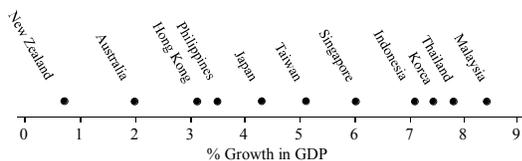
Scale breaks



Dot plot with and without a scale break.

Slide 26 STAT 110A, UCLA, Jon Dinger

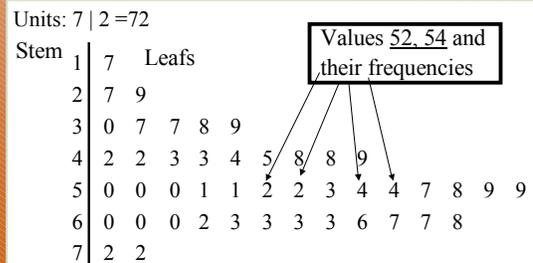
A labeled dot plot



Forecast of percent growth in GDP for 1990 for some South-East Asian and Pacific countries.

Slide 27 STAT 110A, UCLA, Jon Dinger

Example of a stem-and-leaf plot



Stem-plot of the 45 obs's of the Ejection variable in the Heart Attack data table.

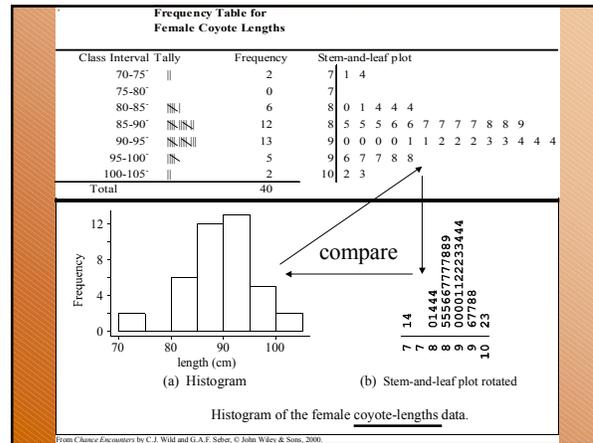
Slide 28 STAT 110A, UCLA, Jon Dinger

Coyote Lengths Data (cm)											
Females											
93.0	97.0	92.0	101.6	93.0	84.5	102.5	97.8	91.0	98.0	93.5	91.7
90.2	91.5	80.0	86.4	91.4	83.5	88.0	71.0	81.3	88.5	86.5	90.0
84.0	89.5	84.0	85.0	87.0	88.0	86.5	96.0	87.0	93.5	93.5	90.0
85.0	97.0	86.0	73.7								
Males											
97.0	95.0	96.0	91.0	95.0	84.5	88.0	96.0	96.0	87.0	95.0	100.0
101.0	96.0	93.0	92.5	95.0	98.5	88.0	81.3	91.4	88.9	86.4	101.6
83.8	104.1	88.9	92.0	91.0	90.0	85.0	93.5	78.0	100.5	103.0	91.0
105.0	86.0	95.5	86.5	90.5	80.0	80.0					

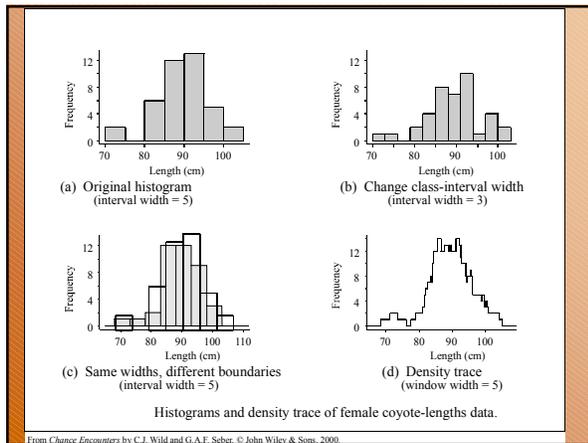
Coyotes captured in Nova Scotia, Canada. Data courtesy of Dr Vera Eastwood.

Frequency Table for Female Coyote Lengths

Class Interval	Tally	Frequency	Stem-and-leaf plot
70-75		2	7 1 4
75-80		0	7
80-85		6	8 0 1 4 4 4
85-90		12	8 5 5 5 6 6 7 7 7 7 8 8 9
90-95		13	9 0 0 0 0 1 1 2 2 2 2 3 3 4 4 4
95-100		5	9 6 7 7 8 8
100-105		2	10 2 3
Total		40	



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



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

Questions ...

- What advantages does a stem-and-leaf plot have over a histogram? (S&L Plots return info on individual values, quick to produce by hand, provide data sorting mechanisms. But, Hist's are more attractive and more understandable).
- The shape of a histogram can be quite drastically altered by choosing different class-interval boundaries. What type of plot does not have this problem? (density trace) What other factor affects the shape of a histogram? (bin-size)
- What was another reason given for plotting data on a variable, apart from interest in how the data on that variable behaves? (shows features, cluster/gaps, outliers; as well as trends)

Slide 34 STAT 110A, UCLA, Jon Dineen

Interpreting Stem-plots and Histograms

(a) Unimodal (b) Bimodal (c) Trimodal (d) Symmetric (e) Positively skewed (long upper tail) (f) Negatively skewed (long lower tail) (g) Symmetric (h) Bimodal with gap (i) Exponential shape

Slide 35 STAT 110A, UCLA, Jon Dineen

Interpreting Stem-plots and Histograms

(j) Spike in pattern (k) Outliers (l) Truncation plus outlier

Features to look for in histograms and stem-and-leaf plots.

Slide 36 STAT 110A, UCLA, Jon Dineen

Descriptive statistics ...

- The sample **mean** is denoted by \bar{x} .

The **sample mean** = $\frac{\text{Sum of the observations}}{\text{Number of observations}}$

Slide 37 STAT 110A, UCLA, Jon Dinger

The **sample mean** is where the dot plot balances

(a) (b) (c)

Mechanical construction representing a dot plot:
(a) shows a balanced rod while (b) and (c) show unbalanced rods.

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

Slide 38 STAT 110A, UCLA, Jon Dinger

The **sample median**

For n observations, $\{x_1, x_2, x_3, \dots, x_n\}$. Suppose we order the observations min-to-max to get $\{x(1), x(2), x(3), \dots, x(n)\}$.

Then the **sample median** is the $[(n+1)/2]$ -st largest Observation $X^{(\frac{n+1}{2})}$.

If $\frac{n+1}{2}$ is not a whole number, the median is the average of the two observations on either side.

Slide 39 STAT 110A, UCLA, Jon Dinger

Effect of **outliers** on the mean and median

(a) Data symmetric about P

(b) Two largest points moved to the right

The mean and the median.

[Grey disks in (b) are the "ghosts" of the points that were moved.]

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

Slide 40 STAT 110A, UCLA, Jon Dinger

Beware of inappropriate averaging

Welcome to MEANSTOWN	
Founded	1867
Area	20
Altitude	584
Population	372
Average	711

Suggested by a 1977 cartoon in *The New Yorker* magazine by Dana Fradon.
From Chance Encounters by C.J. Wild and G.A.F. Seber. © John Wiley & Sons, 1999.

Quantiles (vs. quartiles)

- The **q^{th} quantile** ($100 \times q^{\text{th}}$ percentile) is a value, in the range of our data, so that proportion of at least q of the data lies at or below it and a proportion of at least $(1-q)$ lies at or above it.
- E.g., $X = \{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\}$. The **20th percentile** (**0.2 quantile**) is the value **2**, since 20% of the data is below it and 80% above it. The **70th percentile** is the value 7, etc.
- We could have also selected **2.5** and **7.5** for the **20th** and **70th** percentile, above. There is no agreement on the exact definitions of quantiles.

Slide 46 STAT 110A, UCLA, Jon Dinger

Outliers and Atypical observations

- **Outliers** – an extremely unrepresentative (of the process) data point
- **Example:** A paper mill Co. calibrates a device for measuring water depth, using ultrasound, by timing the echos. In calibrating the equipment they ran a simulation of a water tank with known dept. results:

ACTUAL	ULTRASOUND	DIFFERENCE
■ 0.250	0.2505	0.0005
■ 0.265	0.3011	0.0361
■ 0.280	0.2944	0.0144
■		

Slide 54 STAT 110A, UCLA, Jon Dinger

Outliers and Atypical observations

- **Measuring water depth, using ultrasound.**

Slide 55 STAT 110A, UCLA, Jon Dinger

Outliers and Atypical observations

- Measuring water depth, using ultrasound. There are 4 clear outliers, see histogram (of differences).

Slide 56 STAT 110A, UCLA, Jon Dinger

Outliers and Atypical observations

- Measuring water depth, using ultrasound. There are 4 clear outliers, see histogram (of differences).
- One engineer noticed that at depths below 0.3 ft, the radiation pattern of the ultrasound device intersected the wall of the tank, which appeared to have disturbed the measurements.
- They repositioned the ultrasound device to that the path of the sound was completely within the tank. Repeat of the measurements produced a better histogram without the initial 4 outliers.

Slide 57 STAT 110A, UCLA, Jon Dinger

Outliers and Atypical observations

- Measuring water depth, using ultrasound. After repositioning.

Slide 58 STAT 110A, UCLA, Jon Dinger

Trimmed, Winsorized means and Resistency

- A data-driven **parameter estimate** is said to be **resistant** if it does not greatly change in the presence of outliers.

Order statistic

$$\bar{y}_{tk} = \frac{1}{n - 2k} \sum_{i=k+1}^{n-k} y_{(i)}$$

- **K-times trimmed mean**
- **Winsorized k-times mean:**

$$\bar{y}_{wk} = \frac{1}{n} \left[(k+1)y_{(k+1)} + \sum_{i=k+2}^{n-k-1} y_{(i)} + (k+1)y_{(n-k)} \right]$$

Slide 59 STAT 110A, UCLA, Jon Dinger