#### **Chapter I**

#### Introduction to data



Monday, January 7, 13

#### **Boston housing data**



MTH 183 -Prof. Bradic Winter 2013

#### Consider the following dataset:

The Boston house-price data of Harrison, D. and Rubinfeld, D.L

	crim	chas	•	age	town	rad	medv
	0.00632	0		65.2	Nahant	1	24.0
2	0.02731	0		78.9	Swampscott	2	21.6
3	0.02729	1		61.1	Swampscott	2	34.7
4	0.03237	0		45.8	Marblehead	6	33.4
•••							
505	0.10959	1		89.3	Winthrop	1	22.0
506	0.04741	0		80.8	Winthrop	1	11.9



#### **Boston housing data**

MTH 183 -Prof. Bradic Winter 2013



### 506 observations

506 census tracts



#### **Boston housing data**

	crim	chas	•••	age	town	rad	medv
	0.00632	0		65.2	Nahant	1	24.0
2	0.02731	0		78.9	Swampscott	2	21.6
3	0.02729	1		61.1	Swampscott	2	34.7
4	0.03237	0		45.8	Marblehead	6	33.4
•••							
505	0.10959	1		89.3	Winthrop	1	22.0
506	0.04741	0		80.8	Winthrop	1	11.9







#### We see 6 variables



	crim	chas	•••	age	town	rad	medv
I	0.00632	0		65.2	Nahant	1	24.0
2	0.02731	0		78.9	Swampscott	2	21.6
3	0.02729	1		61.1	Swampscott	2	34.7
4	0.03237	0		45.8	Marblehead	6	33.4
•••							
505	0.10959	1		89.3	Winthrop	1	22.0
506	0.04741	0		80.8	Winthrop	1	11.9





MTH 183 -

Prof. Bradic

Winter 2013









MTH 183 -

**Prof. Bradic** 

Winter 2013

	crim	chas	•••	age	town	rad	medv
I	0.00632	0		65.2	Nahant	1	24.0
2	0.02731	0		78.9	Swampscott	2	21.6
3	0.02729	1		61.1	Swampscott	2	34.7
4	0.03237	0		45.8	Marblehead	6	33.4
505	0.10959	1		89.3	Winthrop	1	22.0
506	0.04741	0		80.8	Winthrop	1	11.9









#### variable rad ——[index of accessibility to radial highways

Numerical (it's a number)
Discrete (can take only a discrete set of values. Here: 6 or more)



	crim	chas	•••	age	town	rad	medv
I	0.00632	0		65.2	Nahant	1	24.0
2	0.02731	0		78.9	Swampscott	2	21.6
3	0.02729	1		61.1	Swampscott	2	34.7
4	0.03237	0		45.8	Marblehead	6	33.4
•••							
505	0.10959	1		89.3	Winthrop	1	22.0
506	0.04741	0		80.8	Winthrop	1	11.9



11

#### **Types of variables**

MTH 183 -Prof. Bradic Winter 2013



MTH 183 -Prof. Bradic Winter 2013

# location of house •Categorical (it's a label)

11





town

Swampscott

Swampscott

Marblehead

. . .

Winthrop

Winthrop

Nahant

MTH 183 -Prof. Bradic Winter 2013

# 

variable town

#### •Categorical (it's a label)

Categorical variables are ALWAYS discrete. The discrete/continuous distinction is made only for numerical variables





	crim	chas	•	age	town	rad	medv
I	0.00632	0		65.2	Nahant	1	24.0
2	0.02731	0		78.9	Swampscott	2	21.6
3	0.02729	1		61.1	Swampscott	2	34.7
4	0.03237	0		45.8	Marblehead	6	33.4
•••							
505	0.10959	1		89.3	Winthrop	1	22.0
506	0.04741	0		80.8	Winthrop	1	11.9









#### 

Numerical (it's a number)
Discrete (can take only a discrete set of values. Here: 2)





#### 

Numerical (it's a number)
Discrete (can take only a discrete set of values. Here: 2)
Could also be considered categorical

•Categories:YES or NO



	crim	chas	•••	age	town	rad	medv
I	0.00632	0		65.2	Nahant	1	24.0
2	0.02731	0		78.9	Swampscott	2	21.6
3	0.02729	1		61.1	Swampscott	2	34.7
4	0.03237	0		45.8	Marblehead	6	33.4
•••							
505	0.10959	1		89.3	Winthrop	1	22.0
506	0.04741	0		80.8	Winthrop	1	11.9







What is the variable of interest here?

Questions we may ask:

- •Are houses on the river more expensive?
- •Does the crime rate impact the value?
- •Is the crime rate affected by the index rad?
- Is the town of winthrop more expensive than nahant?

attach(BostonHousing2)
mean(medv[town=="Winthrop"])
mean(medv[town=="Nahant"])



#### What is the variable of interest here? medy

Questions we may ask:

- •Are houses on the river more expensive?
- •Does the crime rate impact the value?
- •Is the crime rate affected by the index rad?
- Is the town of winthrop more expensive than nahant?

attach(BostonHousing2)
mean(medv[town=="Winthrop"])
mean(medv[town=="Nahant"])



#### Variables relationships

MTH 183 -Prof. Bradic Winter 2013









MTH 183 -Prof. Bradic Winter 2013























### Summarizing numerical data

MTH 183 -Prof. Bradic Winter 2013

The nice things about numerical data: I.We can add (average) them 2.There is a natural order

We will consider two types of summaries

- •Numerical: mean, median, standard deviation, quantiles
- •Graphical: scatter plots, boxplots.

They give the "big picture" about the dataset. Hereafter, we work on the variable medv



MTH 183 -Prof. Bradic Winter 2013

# •A histogram is probably the most informative graphical summary about a variable.



Histogram of medv

medv



# •A histogram is probably the most informative graphical summary about a variable.



Comments:

- Peak in 20-25
- 25 seems to be a cut-off
- Few houses below 10



# •A histogram is probably the most informative graphical summary about a variable.



Comments:

- Peak in 20-25
- 25 seems to be a cut-off
- Few houses below 10

To read it better, we need to know how it's made.



#### Histogram of medv




### Histograms





### Histograms





### Histograms with frequencies



#### Histogram of medv

We can read on the histogram the proportion of census tracts in a bin/class





hist(medv, freq=FALSE, breaks=c(0, 10, 15, 17, 20, 22, 25, 30, 40, 50))



MTH 183 -Prof. Bradic Winter 2013







![](_page_41_Picture_3.jpeg)

![](_page_42_Figure_2.jpeg)

### Questions:

I.where is the largest proportion? 20-22 or 22-25 2.What is the proportion of medv>30?

![](_page_42_Picture_5.jpeg)

### Shape

MTH 183 -Prof. Bradic Winter 2013

0.07 0.06 0.05 0.04 Density 0.03 0.02 0.01 0.00 10 20 0 30 40 50 medv

![](_page_43_Picture_6.jpeg)

### Shape

MTH 183 -Prof. Bradic Winter 2013

![](_page_44_Figure_2.jpeg)

![](_page_44_Picture_3.jpeg)

### Modes

A mode is a peak in the distribution. It indicates a highly populated bin.

We can have **one**, **two**, or **more than two** modes indicating that there may be one, two or more sub-populations in our data.

![](_page_45_Figure_4.jpeg)

### MTH 183 -Prof. Bradic Winter 2013

### **Summary statistics**

	crim	chas	•••	age	town	rad	medv
Ι	0.00632	0		65.2	Nahant	1	24.0
2	0.02731	0		78.9	Swampscott	2	21.6
3	0.02729	1		61.1	Swampscott	2	34.7
4	0.03237	0		45.8	Marblehead	6	33.4
•••							
505	0.10959	1		89.3	Winthrop	1	22.0
506	0.04741	0		80.8	Winthrop	1	11.9

![](_page_46_Picture_3.jpeg)

### **Summary statistics**

crim	Numbers th	at summa	rize crim:	
0.00632	Mean	3.61		
0.02731	Median	0.26	<pre>summary(crim)</pre>	
0.02729	Min	0.01	<pre>mean(crim) modian(crim)</pre>	
0.03237	Max	88.98	min(crim) max(crim)	
	l st quartile	0.08	<pre>quantile(crim,.25) quantile(crim,.75)</pre>	
0.40050	3rd quartile	3.68		
0.10959	Standard deviation	8.60	sd(crim)	
0.04741	Variance	73.99	var(crim)	

![](_page_47_Picture_3.jpeg)

### A little reminder

MTH 183 -Prof. Bradic Winter 2013

Consider a list of numbers:  $x_1, x_2, \ldots, x_n$ 

$$\sum_{i=1}^{n} x_i = x_1 + x_2 + \dots + x_n$$

$$\sum_{i=1}^{4} i = 1 + 2 + 3 + 4 = 10$$
$$\left(\sum_{i=1}^{4} i\right)^2 = (1 + 2 + 3 + 4)^2 = 100$$

$$\sum_{i=1}^{4} i^2 = 1^2 + 2^2 + 3^2 + 4^4 = 30$$

### A little reminder

MTH 183 -Prof. Bradic Winter 2013

![](_page_49_Figure_2.jpeg)

![](_page_49_Picture_3.jpeg)

MTH 183 -Prof. Bradic Winter 2013

### **Measures of location**

Mean:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i = \frac{x_1 + x_2 + \dots + x_n}{n}$$

### Median: the number in the middle

 $\mathbf{n}$ 

Order the list  $x_1, x_2, ..., x_n$  into  $x_{(1)} \le x_{(2)} \le ... \le x_{(n)}$ 

Then median: 
$$x_{med} = \begin{cases} x_{(\frac{n+1}{2})} & \text{if n odd} \\ \frac{x_{(\frac{n}{2})} + x_{(\frac{n}{2}+1)}}{2} & \text{if n even} \end{cases}$$

Example: 65.2 78.9 61.1 45.8 89.3 80.3 variable age

![](_page_50_Picture_8.jpeg)

### Median: n even

MTH 183 -Prof. Bradic Winter 2013

### 45.8 61.1 65.2 78.9 80.3 89.3

![](_page_51_Picture_5.jpeg)

![](_page_52_Picture_0.jpeg)

### $45.8 \leq 61.1 \leq 65.2 \leq 78.9 \leq 80.3 \leq 89.3$

![](_page_52_Picture_3.jpeg)

### Median: n even

MTH 183 -Prof. Bradic Winter 2013

### $x_{(1)}$ $x_{(2)}$ $x_{(3)}$ $x_{(4)}$ $x_{(5)}$ $x_{(6)}$ 45.8 $\leq$ 61.1 $\leq$ 65.2 $\leq$ 78.9 $\leq$ 80.3 $\leq$ 89.3

![](_page_53_Picture_4.jpeg)

37

### Median: n even

## 

![](_page_54_Picture_5.jpeg)

MTH 183 -Prof. Bradic Winter 2013

### Median: n even

![](_page_55_Figure_3.jpeg)

### Median: n odd

MTH 183 -Prof. Bradic Winter 2013

# 

**n=5** odd:  $x_{med} = x_{(3)} = 62.5$ 

![](_page_56_Picture_4.jpeg)

![](_page_57_Picture_0.jpeg)

## Median splits the list in half. Quartiles split it into I/4 - 3/4 and 3/4 - I/4

![](_page_57_Figure_3.jpeg)

First quartile

Third quartile

What is the second quartile?

A fourth of the data is smaller than  $Q_1$ A fourth of the data is larger than  $Q_3$ 

![](_page_57_Picture_8.jpeg)

MTH 183 -Prof. Bradic Winter 2013

### On a dot plot

![](_page_58_Figure_2.jpeg)

![](_page_58_Picture_3.jpeg)

![](_page_58_Picture_4.jpeg)

### Shape, mean and median

MTH 183 -Prof. Bradic Winter 2013

![](_page_59_Figure_2.jpeg)

![](_page_59_Picture_3.jpeg)

Variance: 
$$s^2 = \frac{1}{n-1} \sum_{i=1}^{n} (x_i - \bar{x})^2$$

**Standard deviation:** 
$$s = \sqrt{s^2}$$
 same units as  $x_i$ 

![](_page_60_Picture_4.jpeg)

MTH 183 -Prof. Bradic Winter 2013

Variance: 
$$s^2 = \frac{1}{n-1} \sum_{i=1}^{n} (x_i - \bar{x})^2$$

**Standard deviation:**  $s = \sqrt{s^2}$  same units as  $x_i$ 

If all the numbers in the list are close to each other, they are close to their mean  $\bar{x}$  and  $(x_i - \bar{x})^2$  is close to zero. The variance and (standard deviation) measure how the data is "clustered" around its mean  $\bar{x}$ 

![](_page_61_Picture_5.jpeg)

MTH 183 -Prof. Bradic Winter 2013

Variance: 
$$s^2 = \frac{1}{n-1} \sum_{i=1}^{n} (x_i - \bar{x})^2$$

**Standard deviation:**  $s = \sqrt{s^2}$  same units as  $x_i$ 

If all the numbers in the list are close to each other, they are close to their mean  $\bar{x}$  and  $(x_i - \bar{x})^2$  is close to zero. The variance and (standard deviation) measure how the data is "clustered" around its mean  $\bar{x}$ 

![](_page_62_Figure_5.jpeg)

MTH 183 -Prof. Bradic Winter 2013

### Interquartile range (IQR) = $Q_3 - Q_1$

Is also a measure a dispersion.

## Measure of location and dispersion say nothing about shape. Data x, y and z all have mean 0 and SD 1.

![](_page_63_Figure_5.jpeg)

![](_page_63_Picture_6.jpeg)

### **Boxplots and outliers**

![](_page_64_Figure_2.jpeg)

### **Boxplots and outliers**

![](_page_65_Figure_2.jpeg)

![](_page_65_Figure_3.jpeg)

### Outliers Observations outside of the whiskers

Note: The whiskers draw at the last observation that is not an outlier. This may differ from the exact value  $Q_3 + 1.5 \cdot IQR$  or  $Q_1 - 1.5 \cdot IQR$  significantly

![](_page_65_Picture_6.jpeg)

### Outliers

MTH 183 -Prof. Bradic Winter 2013

It is important to understand the data to identify why outliers are outliers. Possible sources: •Number incorrectly reported •Heterogeneous data (several populations)

The median and IQR are robust to outliers. The mean and SD are not.

![](_page_66_Picture_4.jpeg)

Boxplots become very useful when we want to compare one variable across two categories. For example, we want to know if census tracts along the Charles river have more expansive houses.

![](_page_67_Figure_3.jpeg)

![](_page_67_Picture_4.jpeg)

Boxplots become very useful when we want to compare one variable across two categories. For example, we want to know if census tracts along the Charles river have more expansive houses.

![](_page_68_Figure_3.jpeg)

![](_page_68_Picture_4.jpeg)

MTH 183 -Prof. Bradic Winter 2013

### **Contingency tables**

•For categorical data we cannot compute the mean, variance, median, ...

![](_page_69_Figure_3.jpeg)

### **Two-way contingency tables**

For two variables, we can cross them. Consider this dataset about 2,201 passengers of the RMS Titanic.

### Survived

		No	Yes	Total
Class	l st class	122	203	325
	2nd class	167	118	285
	3rd class	528	178	706
	crew	673	212	885
	Total	I,490	711	2,201

ftable(Titanic, row.vars="Class", col.vars="Survived")

![](_page_70_Picture_6.jpeg)

MTH 183 -Prof. Bradic Winter 2013

### **Mosaic plots**

### Useful represent categorical data. Displays the content of a contingency table.

![](_page_71_Figure_3.jpeg)

Titanic

Class

table(rad, chas)

![](_page_71_Figure_6.jpeg)

mosaicplot(~ Class + Survived, data = Titanic, color = TRUE)

mosaicplot(table(rad, chas), color = TRUE)

![](_page_71_Picture_9.jpeg)
MTH 183 -Prof. Bradic Winter 2013

#### **Mosaic plots**



# We can see that Class and Survived are associated variables.



chas and rad are associated variables



MTH 183 -Prof. Bradic Winter 2013

### **Mosaic plots**



## These variables seem to be independent



 $\times$ 

### Why mosaic plots?

Arguably, mosaic plots are not the most accurate plot to assess independence. But compare this to the following scatter plot of chas vs rad





•Our Boston Housing dataset is a sample of census tracts in the Boston area.

- •All census tracts in the greater Boston area represent the population.
- •The census tracts were selected randomly from the population.



•Our Boston Housing dataset is a sample of census tracts in the Boston area.

•All census tracts in the greater Boston area represent the population.

•The census tracts were selected randomly from the population.

It is important to trust the data collection process: no census tract should be preferred over another in the selection.

This way we can extrapolate to other census

tracts.



Without this uniform sampling at random things could go wrong. For instance, if we selected US census tracts that vote republican, we could end up with mostly houses with high median value.





#### Republican



Without this uniform sampling at random things could go wrong. For instance, if we selected US census tracts that vote republican, we could end up with mostly houses with high median value.



#### Democrat

#### **Representative sample**

- Thus, the sample should be representative of the population.
- In this course, we will always assume that it is but if you design your own statistical experiment our if you use data that has been collected by someone else, you should make sure that you have a sample representative of the population that you are interested in.
- Note: you can always consider a smaller population. For example, the population of republican voting census tracts.



MTH 183 -Prof. Bradic Winter 2013

# Random sapling and probability

Assuming uniform sampling at random allows us to use probability.

For example, if we know that 10% of the census tracts are on the Charles river, probability tells us how many census tracts among the 506 sampled are on the Charles river (in average): 50.6 The idea of statistics is to replace this 10% by a variable p to be estimated. In our sample, there are 35 census tracts on the Charles river so we estimate p to be 35/506=6.9%. Probability will allow us to do much more.