

**ECS7024 Statistics for Artificial Intelligence and Data
Science**

Topic 8: Contingency Tables

William Marsh

Outline

- Aim: Understand how to use a Contingency Table to look at the relationship between categorical variables
- Probability recap: joint and conditional
- Contingency tables (cross tabulation)
- Comparing categorical variables
- Odds and odds ratios
- Comparing distribution (of a continuous variable) for different categories

Context

- Correlation: do two variables 'vary together'?
- Strength of correlation
 - Approach: average of $(x_i - \mu_X)(y_i - \mu_Y)$
 - ... also written $(x_i - \bar{x})(y_i - \bar{y})$ *for sample values*
 - Only applies to continuous variables (with a mean)
- What about categorical variables?
- ANS: look at proportions of one variable, conditional on the values of others

Probability Recap

Recap on Joint and Conditional
Probability

What is $P(A)$?

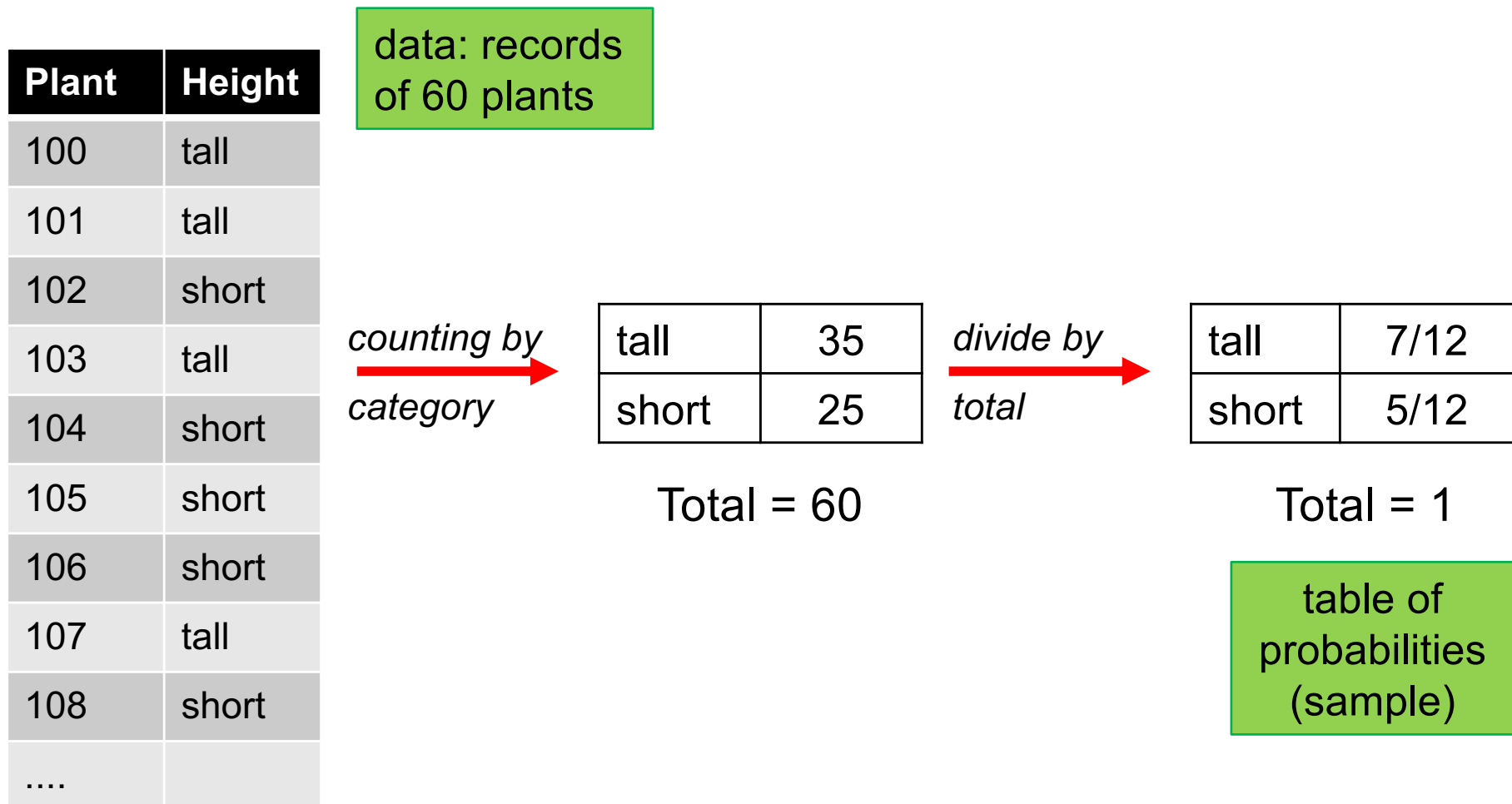
- Survey of plants
- Plant height (H)
 - Categorical variable
 - Values: tall (t), short (s)
- $P(H)$ is a table

Tall	7/12
Short	5/12

- How do we estimate probabilities from data?

Probabilities Come from Counts

- Counts in each category



Joint Probability $P(H, L)$

- Counts in two Categories

Leaves (L)	Stem Height (H)	
	Short (s)	Tall (t)
Broad (b)	9	21
Long (g)	16	14

Still 60 plants
in total

- Probabilities

$P(H, L)$ is also a table

Leaves (L)	Stem Height (H)	
	Short (s)	Tall (t)
Broad (b)	9/60	21/60
Long (g)	16/60	14/60

Total is 1

Conditional Probability

- Look at some entries in the table
 - The short ones

Leaves (L)	Stem Height (H)	
	Short (s)	Tall (t)
Broad (b)	9/60	21/60
Long (g)	16/60	14/60

25 short plants

- $P(L \mid H = \text{short})$

broad	9/25
long	16/25

*probability of leaf types
given the plant is short*

Conditional Probability

- Look at some entries in the table
 - The short ones or the tall ones

Leaves (L)	Stem Height (H)	
	Short (s)	Tall (t)
Broad (b)	9/60	21/60
Long (g)	16/60	14/60

25 short plants

35 tall plants

- $P(L \mid H = \text{short})$

broad	9/25
long	16/25

*probability of leaf types
given the plant is **short***

- $P(L \mid H = \text{tall})$

broad	21/35
long	14/35

*probability of leaf types
given the plant is **tall***

Conditional Probabilities Table

- $P(L|H)$ – a table of tables
 - Each column is a probability distribution

Leaves (L)	Stem Height (H)	
	Short (s)	Tall (t)
Broad (b)	9/25	21/35
Long (g)	16/25	14/35

$P(L | H = \text{short})$

total is 1

$P(L | H = \text{tall})$

total is 1

Probabilities and Tables: Summary

- Assume variable with two categories (binary)
 - Generalises to more categories

Probability	Table Entries	Total
$P(H)$	2: short & tall	1
$P(L)$	2: broad & long	1
$P(H, L)$	4: (s, b) & (s, g) & (t, b) & (t, g)	1
$P(L, H)$	Same table as $P(H, L)$	
$P(H L)$	Broad column: s & t	1
	Long column: s & t	1
$P(L H)$	Short column: b & g	1
	Tall column: b & g	1

Relationship between Joint and Conditional Probabilities

Contingency Tables

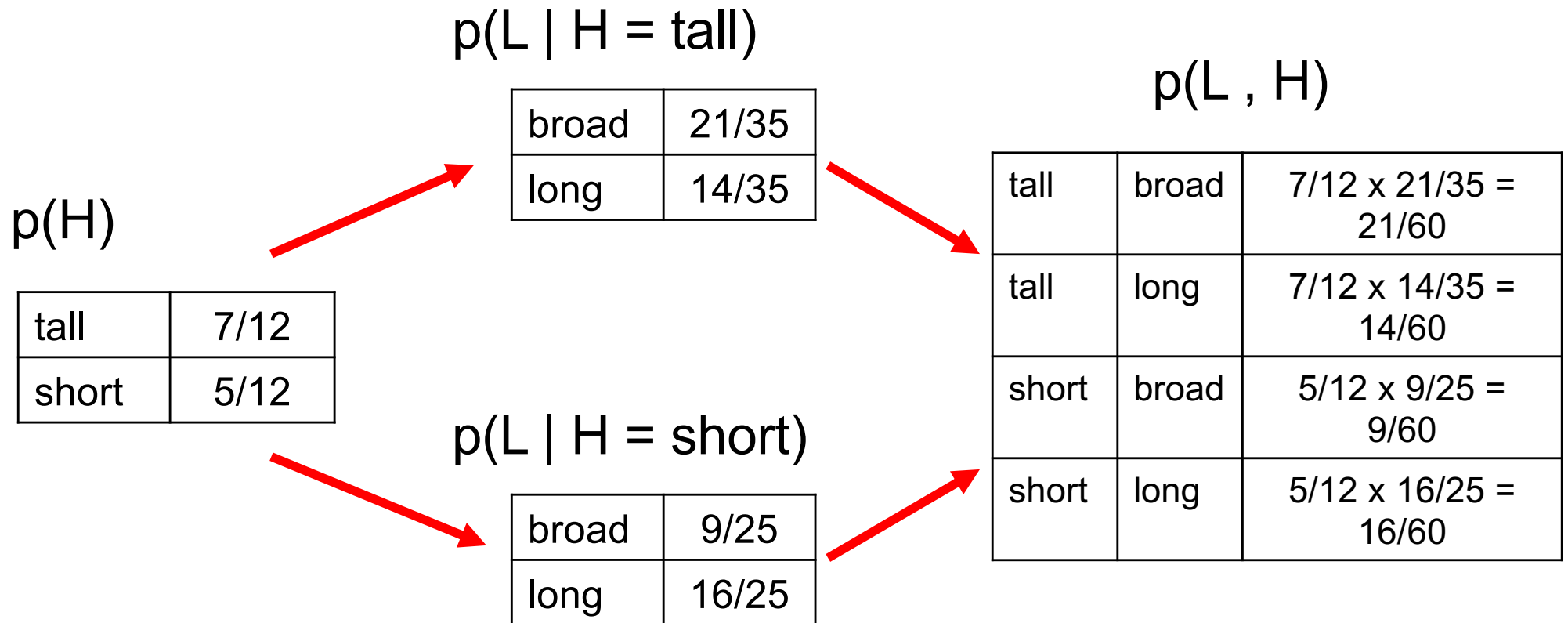
Recall Rule for P(A,B)

$$P(A, B) = P(A) \cdot P(B|A)$$

$$P(A, B) = P(B) \cdot P(A|B)$$

- *The joint probability of A & B is the probability of A multiplied by the probability of B given A*
- *(also the other way around)*

$P(H).P(L|H)$ Using Tables



Marginalisation

- Going from $p(H, L)$ to $p(H)$ or $p(L)$

Leaves (L)	Stem Height (H)	
	Short (s)	Tall (t)
Broad (b)	9/60	21/60
Long (g)	16/60	14/60



Add up value with same height

	25/60	35/60
--	-------	-------

$$p(H) = \sum_{l \in \{b, g\}} p(H, L = l)$$

$$p(L) = \sum_{h \in \{s, t\}} p(H = h, L)$$

Independence

If two variables A , B are independent then:

- $p(A \mid B) = p(A)$
 - *knowing B makes not difference to what A to expect*
- $p(B \mid A) = p(B)$
 - *Same the other way around*
- $p(A, B) = p(A).p(B)$
 - *The joint probability is given by the product*
- Dependent?
 - If $p(A \mid B=b_1)$ differs from $p(A \mid B=b_2)$

Quiz 1



Gold (not chicken)

Python Nugget

Keyword parameters

Keyword Parameters

- Function parameters can be given by
 - Position
 - Keyword
- Keywords have a default value
- Also allows 'extra' parameters

```
def addxy(x, y=1):  
    return x + y  
  
addxy(2, 3) # result 5  
addxy(2, y=3) # result 5  
addxy(2)      # result 3 - default  
addxy(y = 3) # error
```

Examples of Keyword Arguments

pandas.DataFrame

`class pandas.DataFrame(data=None, index=None, columns=None, dtype=None, copy=False)`

[\[source\]](#)

Two-dimensional, size-mutable, potentially heterogeneous tabular data.

Data structure also contains labeled rows and columns, thought of as a dict-like container for ordered dicts.

Parameters: **data** : ndarray (structured or unstructured), dict-like object, or pandas DataFrame

Dict can contain arrays, scalars, or lists. Note that None values are allowed for the column values. Default is None, which means the column will not be included in the DataFrame. Changed in version 0.17.0: Data can be a Series.

index : Index or array-like

Index to use for the new DataFrame. If None, the index will be the same as the data's index. If no index is provided, the index will be a RangeIndex from 0 to n - 1 where n is the number of rows.

columns : Index or array-like

Column labels to use for the new DataFrame. If None, the columns will be the same as the data's columns. If no columns are provided, the columns will be the same as the data's columns.

dtype : dtype, default None

Data type to force. If None, the data type will be inferred from the data.

copy : bool, default False

Copy data from objects. If True, a copy of the data is made. If False, the data is referenced.

pandas.DataFrame.plot

`DataFrame.plot(*args, **kwargs)`

Make plots of Series or DataFrame.

Uses the backend specified by the option `plotting.backend`. By default, matplotlib is used.

Parameters: **data** : Series or DataFrame

The object for which the method is called.

x : label or position, default None

Only used if data is a DataFrame.

y : label, position or list of label, positions, default None

Allows plotting of one column versus another. Only used if data is a DataFrame.

kind : str

The kind of plot to produce:

- 'line': line plot
- 'bar': vertical bar chart

****kwargs**

Options to pass to matplotlib plotting method.

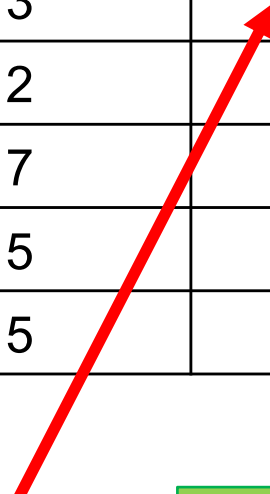
Cross Tabulation

Contingency Tables

Generalising to More Variables

- Contingency tables can show counts over many categories

Leaves (L)	(Seed) Size (Z)	Stem Height (H)	
		Short (s)	Tall
Broad (b)	Large (l)	4	7
	Medium (m)	3	10
	Small (s)	2	4
Long (l)	Large (l)	7	3
	Medium (m)	5	4
	Small (s)	5	6


$$p(L=b, H=t, Z=m) = 10/60$$

Probability tables given
by 'normalising'

Contingency Table Summary

- Also called 'cross tabulation'
 - Pandas has 'crosstab'
- Table shape not fixed
 - Can have multiple variables on each axis
- Table contains
 - Counts

Comparing Categories

Heart Data Again

Heart Disease: Categorical

Variable	Meaning	Type
Sex	1 = male, 0 = female	Categorical
ChestPain	The chest pain experienced (1: typical angina, 2: atypical angina, 3: non-anginal pain, 4: asymptomatic)	Categorical
Bsugar	The person's fasting blood sugar (> 120 mg/dl, 1 = true; 0 = false)	Binary
Angina	Exercise induced angina (1 = yes; 0 = no)	Binary
RestECG	Resting electrocardiographic measurement (0 = normal, 1 = having ST-T wave abnormality, 2 = showing probable e left ventricular hypertrophy)	Ordinal (?)
ECG_ST_slope	The slope of the peak exercise ST segment (1: upsloping, 2: flat, 3: downsloping)	Categorical
Vessels	The number of major vessels (0-3) coloured by fluoroscopy	Ordinal
Thallium	Thallium update test (0 = normal; 1 = fixed defect; 2 = reversible defect)	Categorical
Disease	Heart disease (0 = no, 1 = yes)	Binary

Questions: Do They Vary Together?

- Two categorical variables
 - Does knowing one change value the distribution of the other
- ‘Dependent’ or ‘response’ variable: Disease
 - Often one of the variables

Cross Tabulation Example

- 'Crosstab' counts occurrences in categories
- Can 'normalise'
- Results are joint or conditional probability distributions

$p(\text{Sex}, \text{Disease})$

Sex	M	F
Disease		
False	0.30	0.24
True	0.38	0.08

'Normalise all'

$p(\text{Disease} \mid \text{Sex})$

Sex	M	F
Disease		
False	0.44	0.74
True	0.56	0.26

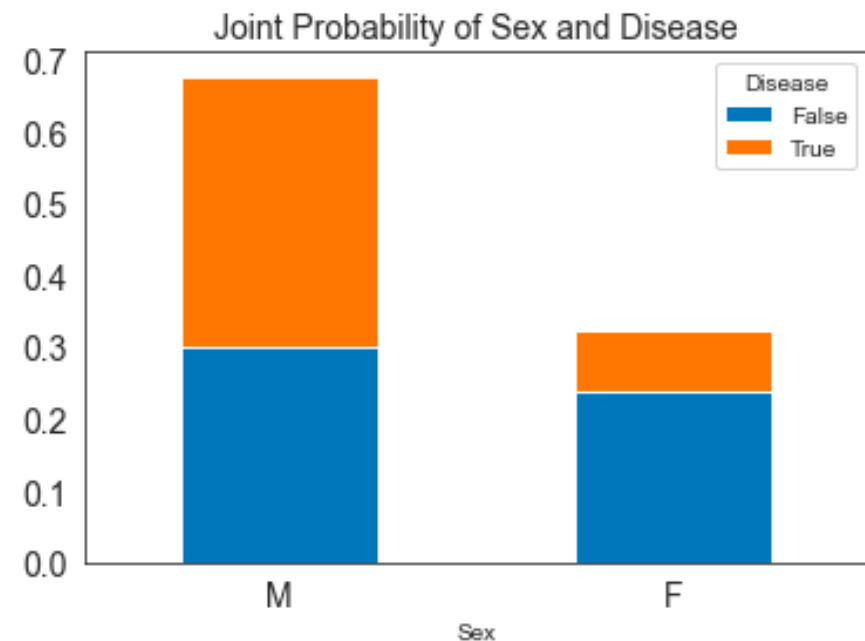
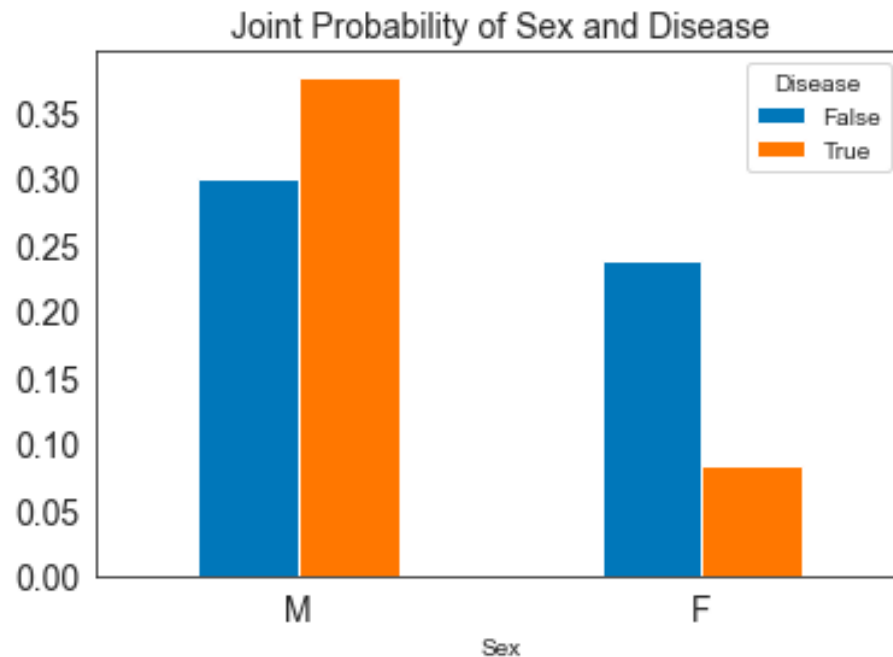
'Normalise columns'

Cross Tabulation Example

- Joint probability

$p(\text{Sex}, \text{Disease})$

Sex	M	F
Disease		
False	0.30	0.24
True	0.38	0.08

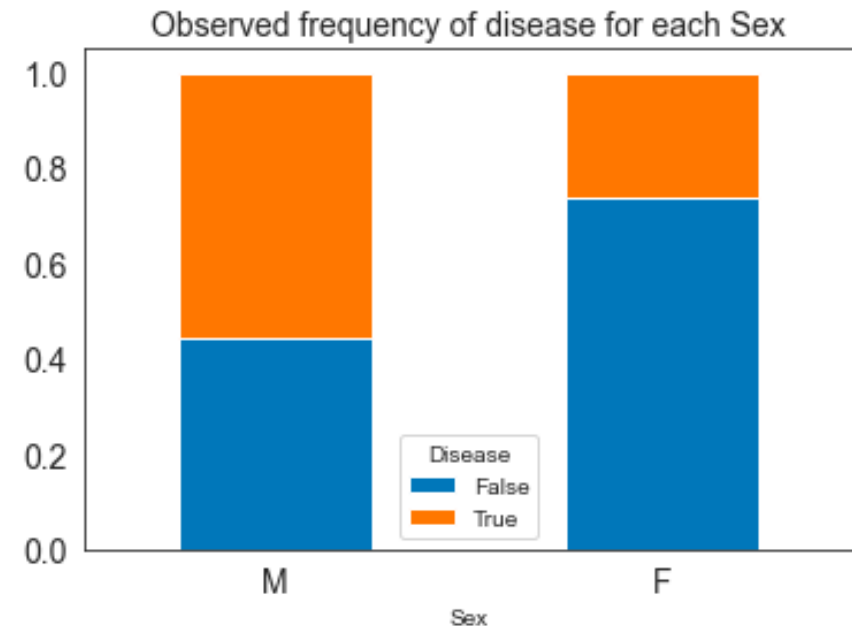
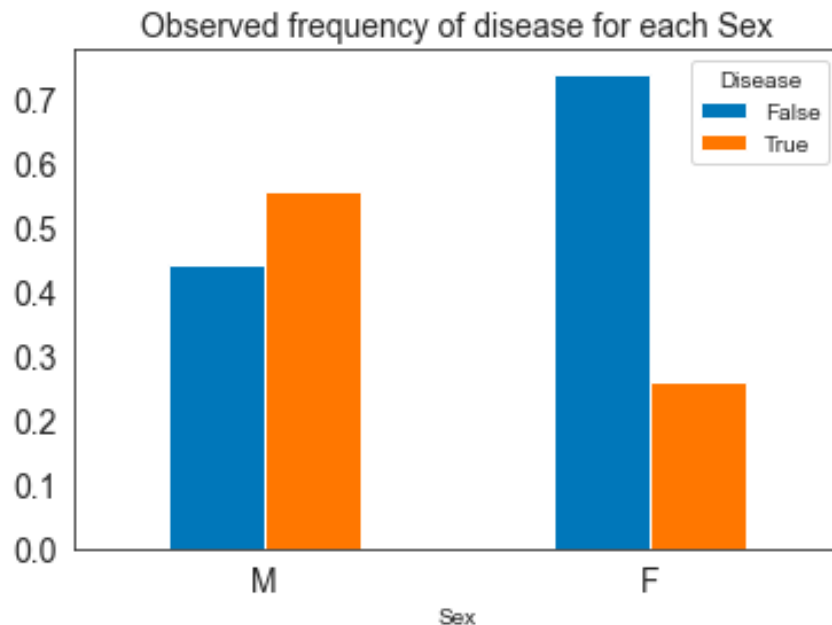


Cross Tabulation Example

- Conditional probability

$$p(\text{Disease} \mid \text{Sex})$$

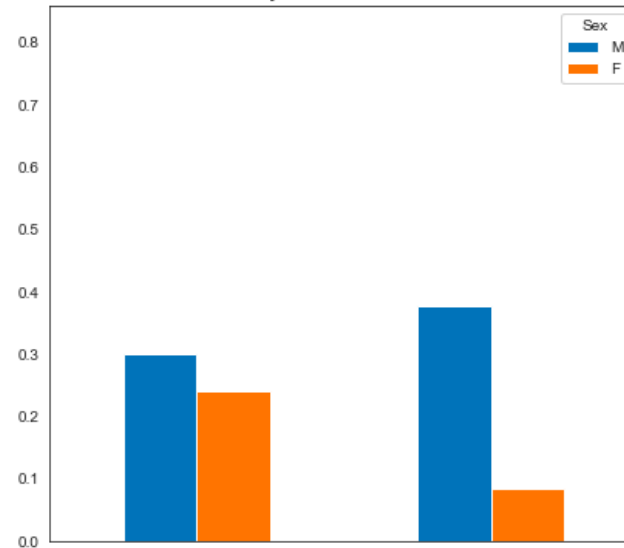
Sex	M	F
Disease		
False	0.44	0.74
True	0.56	0.26



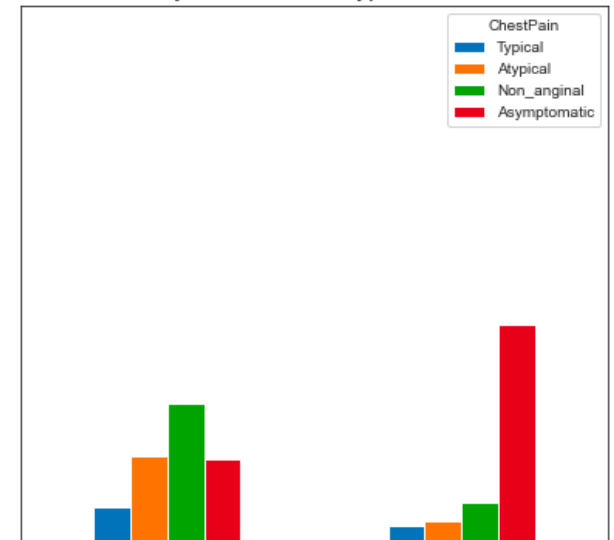
Which Conditional Probability?

Joint distribution:
easy to understand,
but harder to see
'association'

Joint Probability of Sex and Disease Status



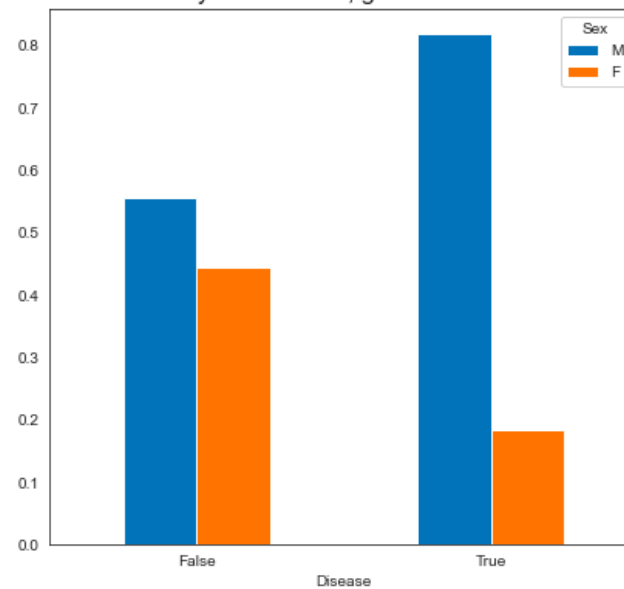
Joint Probability of Chest Pain type and Disease Status



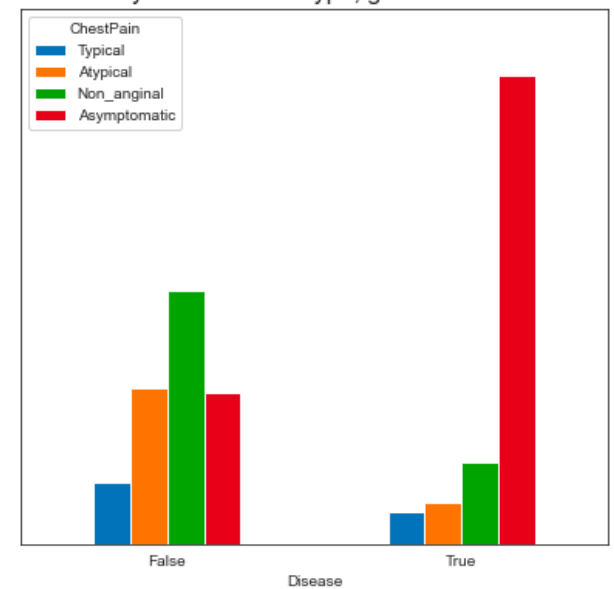
$p(? | \text{Disease})$

Is this the best
way around?

Probability of each Sex, given Disease Status

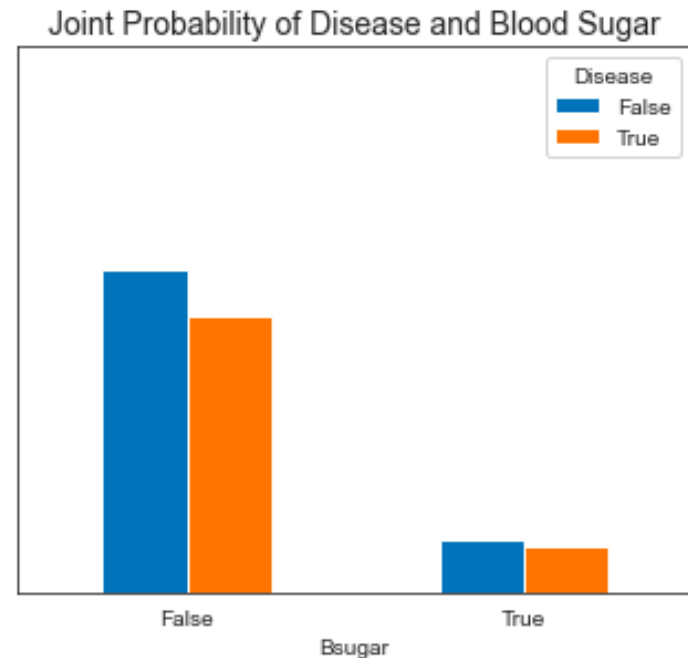
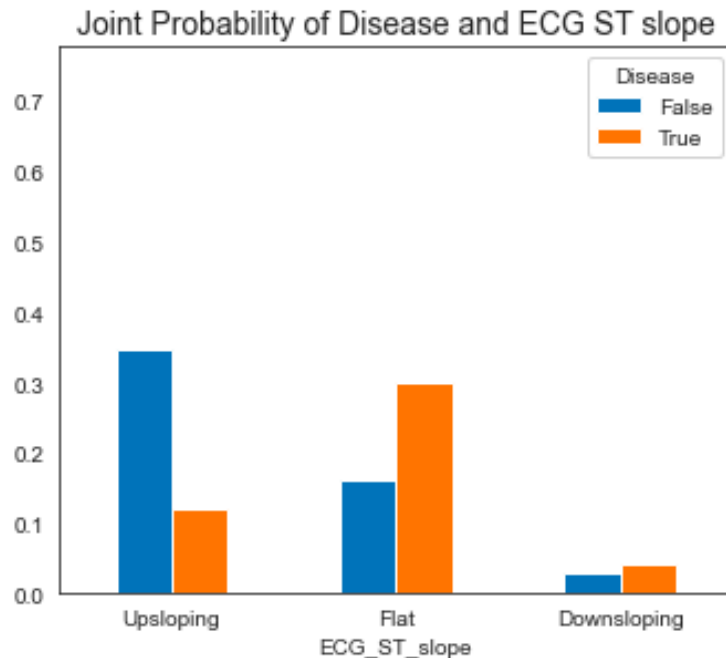


Probability of Chest Pain type, given Disease Status

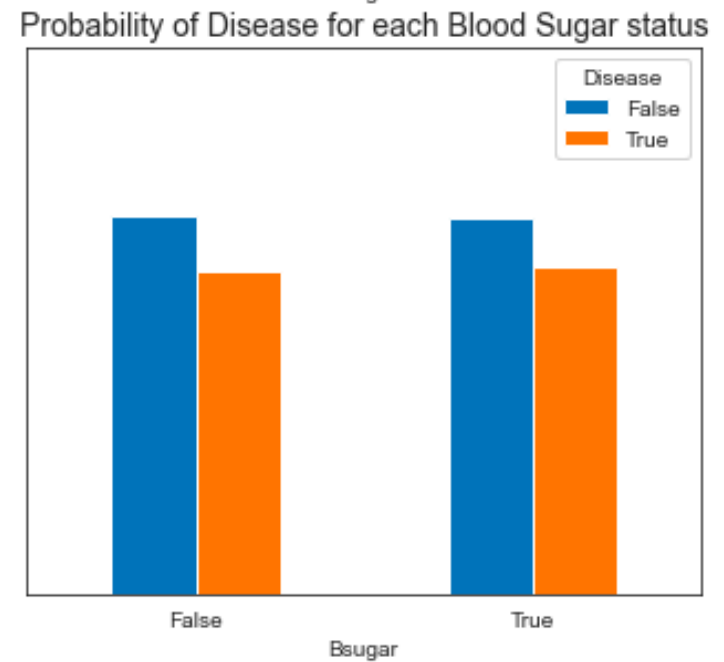
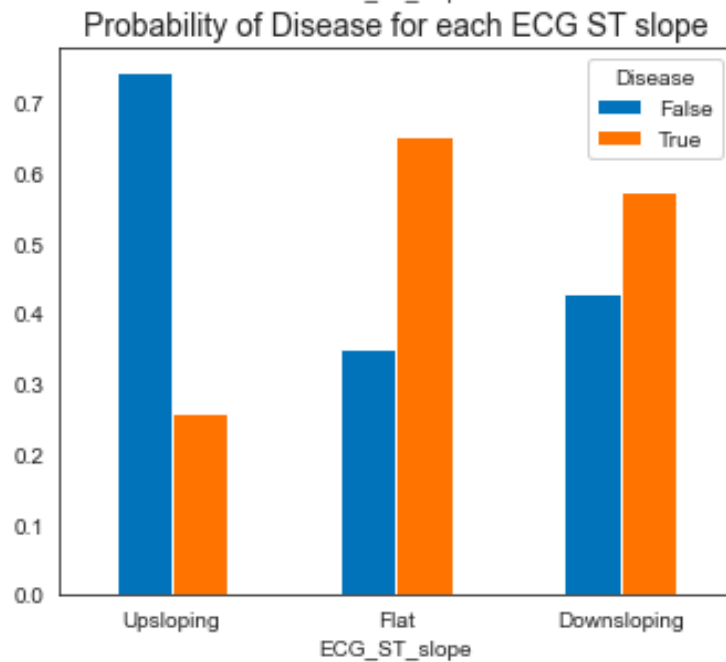


Prefer: $p(\text{Disease} \mid ?)$

Joint



Conditional



Probability and Odds

Odds is Another Way to Write a Probability

- Two rules of probability
 - $0 \leq p(A) \leq 1$
 - $p(A) + p(\text{not } A) = 1$ (we write 'not A' as \bar{A})
- Definition of odds: $o_A = \frac{p(A)}{p(\bar{A})}$
 - Odds ranges from zero upwards
 - $o_{\bar{A}} = 1/o_A$ so that $o_A \cdot o_{\bar{A}} = 1$
- Example: $p(A) = 75\%$ then $\text{odds}_A = 75/25 = 3$
 - Odds > 1 implies probability $> 50\%$
 - Odds < 1 implies probability $< 50\%$

Odds / Odds Ratio of Heart Disease

- From sample probabilities
 - $\text{odds}_{\text{Disease} | \text{M}} = 56/44 = 1.27$
 - $\text{odds}_{\text{Disease} | \text{F}} = 26/74 = 0.35$

Sex	M	F
Disease		
False	0.44	0.74
True	0.56	0.26

- Observation:
 - Men have increased chance of heart disease
 - Odds ratio: $\text{odds}_{\text{Disease} | \text{M}} / \text{odds}_{\text{Disease} | \text{F}} = 3.6$
- ‘Odds ratio’ can measure strength of one variable on another (for binary variables)

Quiz 2

Comparing Continuous Distributions for Categories

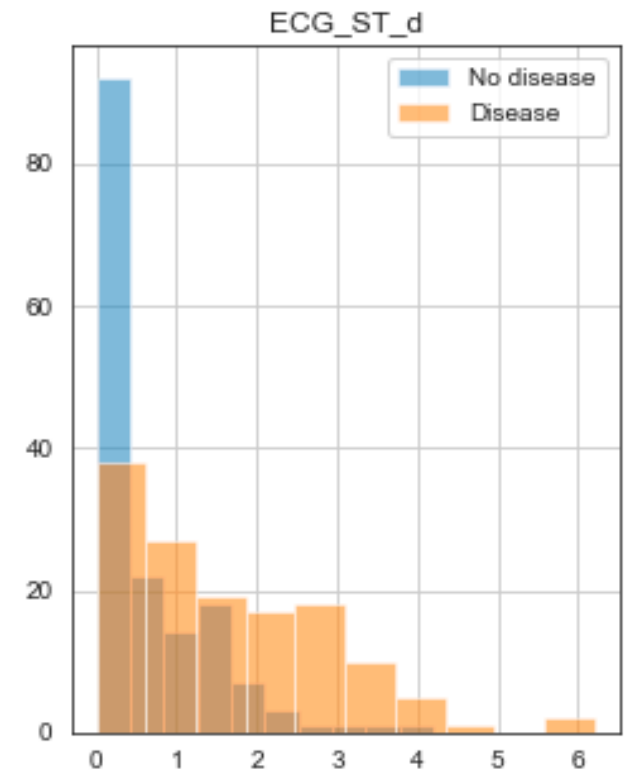
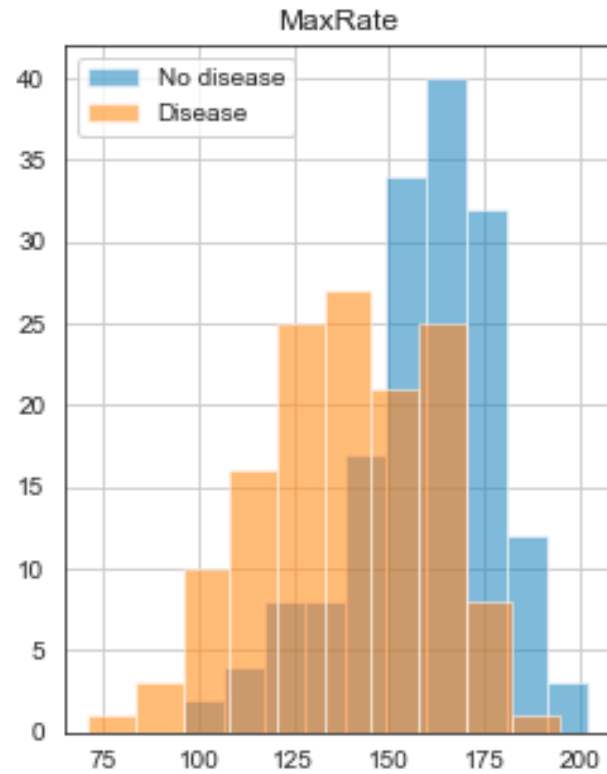
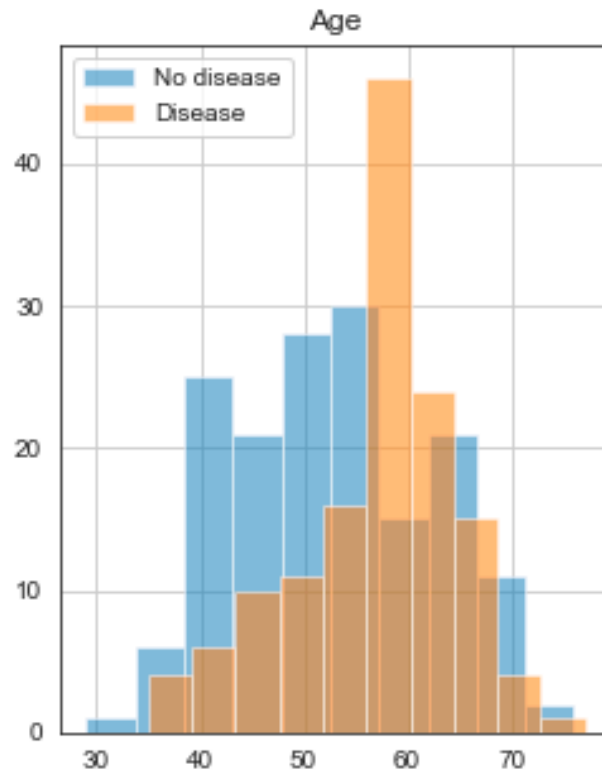
Heart Disease: Continuous

Variable	Meaning	Type
Age	The person's age in years	Continuous
RestBP	The person's resting blood pressure (mm Hg on admission to the hospital)	Continuous
Chol	The person's cholesterol measurement in mg/dl	Continuous
MaxRate	The person's maximum heart rate achieved	Continuous
ECG_ST_d	ST depression induced by exercise relative to rest ('ST' relates to positions on the ECG plot)	Continuous

Question: Do They Vary Together?

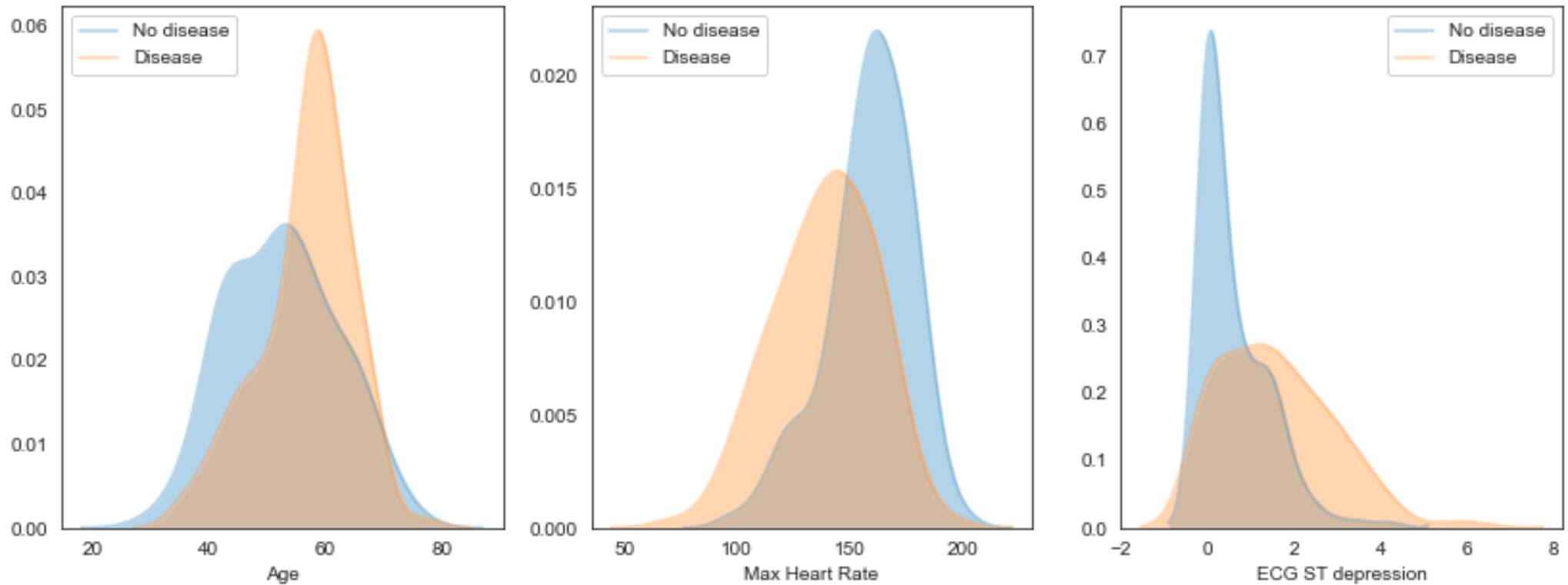
- A categorical variable and a continuous one
 - Is the continuous variable distribution different for the different categorical values
- Example: look at the different distribution of 3 continuous variables by disease status
 - Histograms
 - Kernel density
 - Boxplot

Histograms



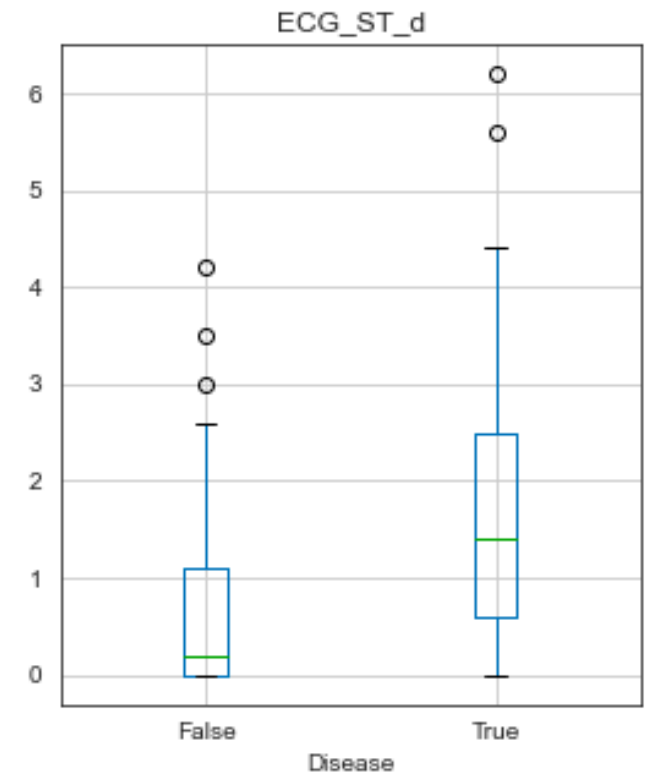
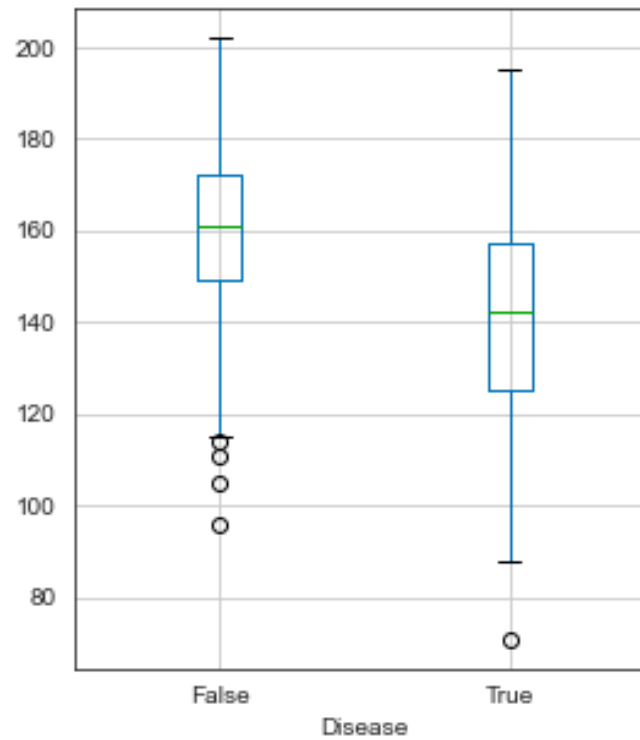
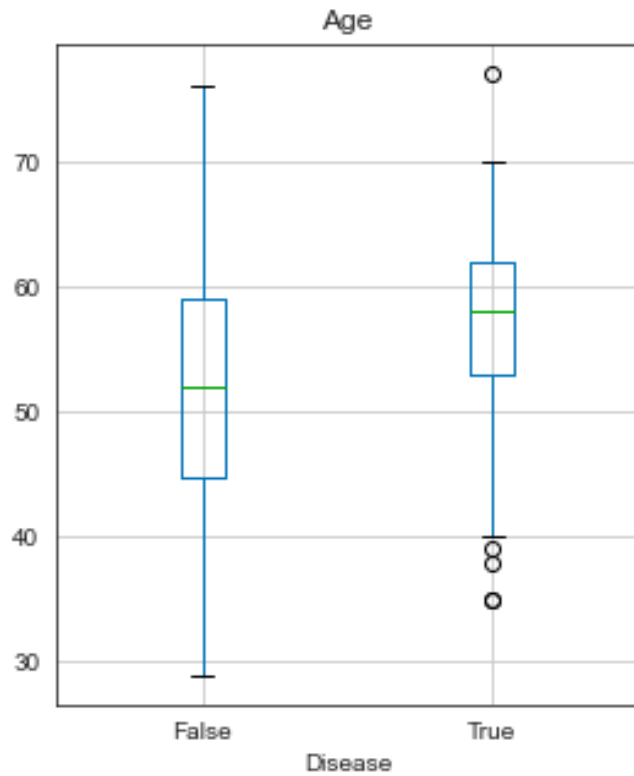
Using Kernel Density

- Kernel Density Estimator (KDE): smoothed histogram



BoxPlots

Boxplot grouped by Disease



Quiz 3

Summary

- Comparing categorical variables
 - Look at conditional probability or odds ratio
- Contingency tables
 - Counting, by category
 - Normalised to give probabilities (joint or conditional)
- Which conditional?
 - Very easy to get confused – label clearly
- Continuous by category
 - use boxplot, KDE or histogram
- *Future: are the differences we see 'reliable'?*