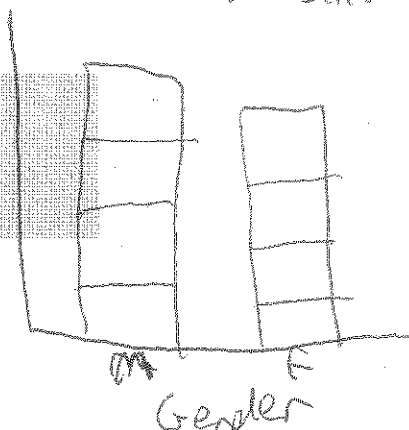


$\text{aes}(x = \text{gender}, \text{fill} = \text{class\_year})$   
 ↳ marginal      ↳ conditional

Today: 2-D Categorical Data  
 Independence and Mosaic Plots  
 1-D Continuous Data



Sam Ventura  
 36-315

Department of Statistics  
 Carnegie Mellon University

September 19, 2016

= cond. dist'n of Var 2 given Var 1 = cat 1

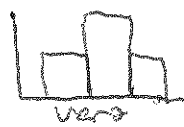
2 categorical variables: Var 1 + Var 2  
 $O_{ij}$  = # of obs. with category  $i$  in Var 1 & cat  $j$  in Var 2  
 $K_1$        $K_2$   
 categories

Contingency Tables and Marginal/Conditional Distributions

		Var 2			
		Cat 1	Cat 2	...	Cat $K_2$
Var 1	Cat 1	$O_{11}$	$O_{12}$	...	$O_{1K_2}$
	Cat 2	$O_{21}$	$O_{22}$	...	$O_{2K_2}$
	⋮	⋮	⋮	⋮	⋮
	Cat $K_1$	$O_{K_1 1}$	$O_{K_1 2}$	...	$O_{K_1 K_2}$
		$N_{.1}$	$N_{.2}$	...	$N_{.K_2}$

Marginal dist'n of Var 1 (points to the rightmost column)  
 Marginal dist'n of Var 2 (points to the bottom row)

Suppose we have two ~~events~~ events A, B



Recall: Independence Rules from Probability

$A \perp\!\!\!\perp B \Leftrightarrow$

$P(A) = P(A|B) \rightarrow$  marginal dist<sup>n</sup> of ~~Var 1~~ = conditional dist<sup>n</sup> of Var 1 | Var 2

$P(B) = P(B|A) \rightarrow$   $\text{Var 2} = \text{Var 2} | \text{Var 1}$

$P(A \cap B) = P(A) \cdot P(B) \rightarrow$  Joint dist<sup>n</sup> of Var 1 and Var 2 = (marginal of Var 1)  $\times$  (marginal of Var 2)

Can input contingency tables into chi-square tests for independence

E.g. `chisq.test(table(var1, var2))`

More on this in Lab 04

$\rightarrow$  contingency table in R

Pearson Residuals

$\rightarrow$  assumes  $\text{Var 1} \perp\!\!\!\perp \text{Var 2}$

Pearson Residuals: Scaled difference between observed/expected

$\rightarrow$  independent

$$r_{ij} = \frac{O_{ij} - E_{ij}}{\sqrt{E_{ij}}} = \frac{\text{observed} - \text{Expected}}{\sqrt{\text{expected}}}$$

$r_{ij} > 0$ : "too many" observed  $\} \text{ assuming Var 1 and Var 2 are independent}$   
 $r_{ij} < 0$ : "too few" observed

$r_{ij}$ 's are Asymptotically Normally distributed!

$|r_{ij}| > 2 \Rightarrow$  significant at the  $\alpha = 0.05$  level

$|r_{ij}| > 4 \Rightarrow \dots \alpha = 0.0001$  level

color Mosaic plot according to the Pearson residuals in each cell

Mosaic Plots  $\rightarrow$  visualizes contingency tables

Mosaic Plots: Area plot for two categorical variables

each cell in the contingency table gets a box in the plot

area of each box  $\propto$  % of obs. in the corresponding cell of the conf. table

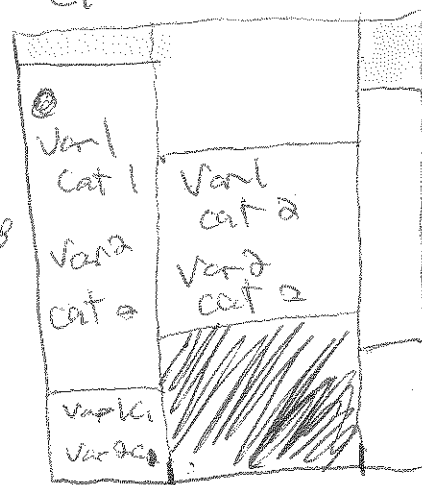
width of box  $\propto$  % of obs. in Var 1 cat  $i$   
 $\rightarrow$  marginal distn of Var 1

height of box  $\propto$  % of obs. in Var 2 cat  $j$   
 given Var 1 is in cat  $i \rightarrow$  conditional distribution

Can color the boxes by their differences from what was expected

Friday: Mosaic Plots in R

5/8



## 1-D Continuous Data

Ex: time,  
age, temperature  
distance, height  
weight, rate  
probabilities

Structure:

$X = \{X_1, X_2, \dots, X_n\}, X_i \in \mathbb{R}$   
 $n \times 1$  vector (column of our data), each obs. is a real #

Summary:

mean / average, median (or other measures of center)

spread: standard deviation, variance, IQR

range: min, max, quantile, percentiles

In R:

range(), min(), max()

mean(), sd(), median()

summary(), var(), quantile()

modality  $\rightarrow$  unimodal, bimodal, multimodal

## 1-D Continuous Distributions

How do we describe continuous distributions?

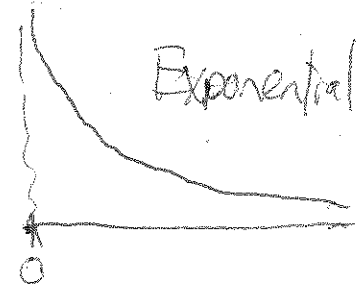
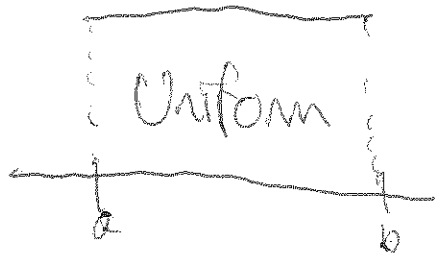
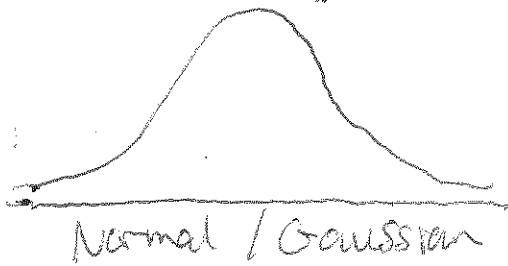
outliers

Shape (skew)  
Center (mean, median)  
Spread (sd, var)

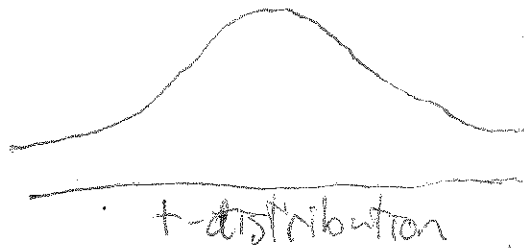
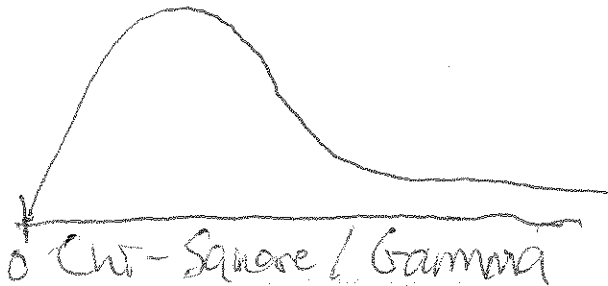
symmetry

ranges of values that are most common  $\rightarrow$  "high density"  
or least common  $\rightarrow$  "low density"

Do the data appear to fit some common/known distribution:



7/8



## Visualizing 1-D Continuous Distributions

How do we visualize continuous distributions?

$\rightarrow$  divide range into bins

Histogram  $\sim$  Bar Chart

For continuous data  
 $\rightarrow$  count obs in each bin

