

# Today: 2-D Categorical Data Independence and Mosaic Plots

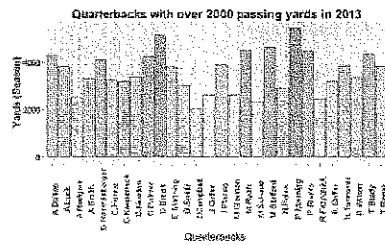
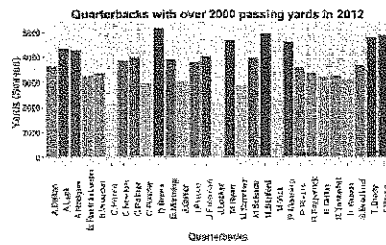
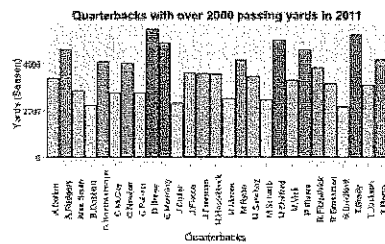
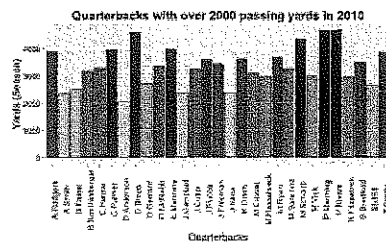
Sam Ventura  
36-315

Department of Statistics  
Carnegie Mellon University

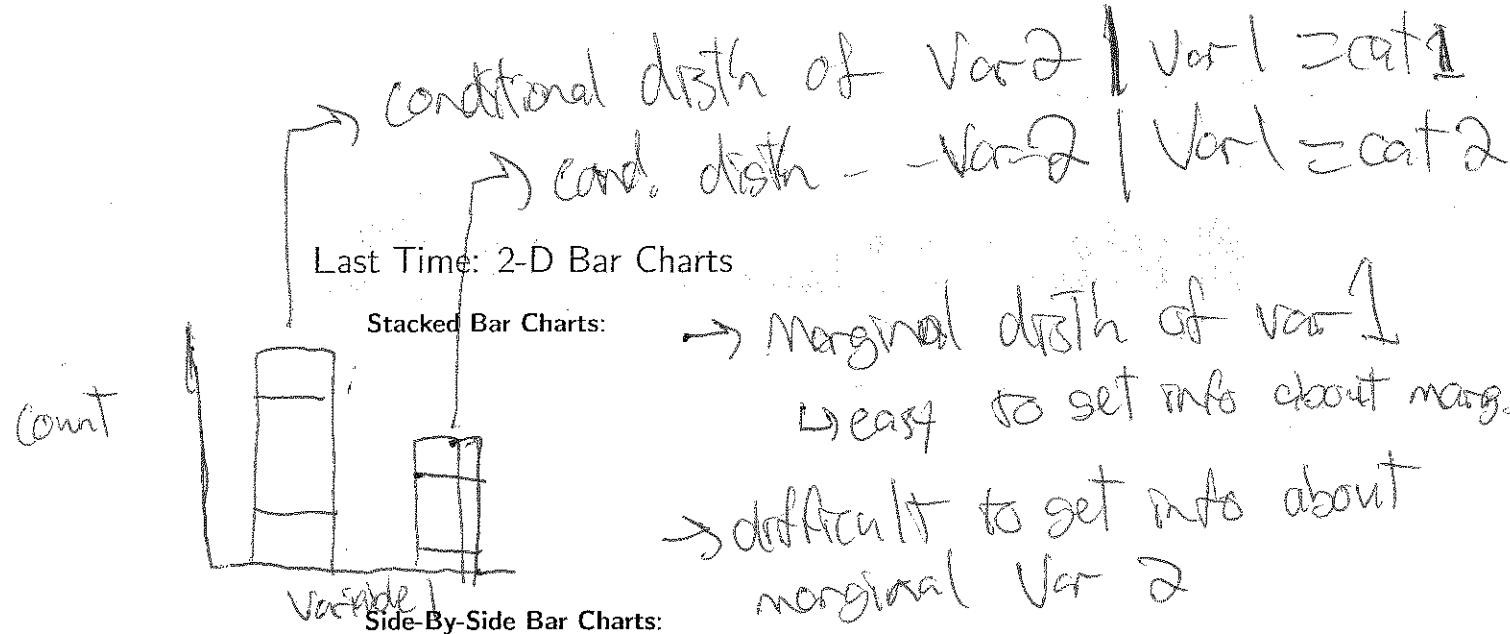
February 1, 2016

1/8

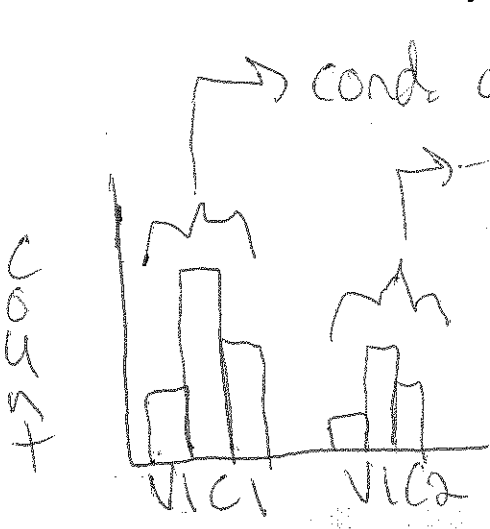
## Quarterbacks with Over 2000 Yards, 2010 – 2013



2/8



→ Marginal distn of  $Var1$   
 ↳ easy to get info about marg.  
 → difficult to get info about marginal  $Var2$



conditional distn of  $V2 | V1 = 1$   
 $V2 | V1 = 2$   
 → hard to see marginal distn of  $Var1$  in side-by-side.  
 → easy to see cond. distn of  $V2 | V1$

### Contingency Tables

#### Two Categorical Variables:

- Give us
- Marginal distns (of one variable)
  - Joint distns (of two variables)
  - cond. distns ( $Var2 | Var1, Var1/Var2$ )

#### Asymptotics

### Two Cat. Variables:

$Var1$  has  $k_1$  categories  
 $Var2$  has  $k_2$  categories

#### Interested in:

- observed counts in each cell  
 -  $C_{ij} = \# \text{ obs w/ } Var1 = i, Var2 = j$
- Expected counts in each cell  
 -  $\hat{C}_{ij} = \frac{n_{i.} \times n_{.j}}{n}$

whole thing = joint dist'n

# Contingency Tables and Marginal/Conditional Distributions

	Nas2 Cat1	Cat 2	...	Cat K	
Cat1	$C_{11}$	$C_{12}$	...	$C_{1K}$	$n_{1\cdot}$
cat 2	$C_{21}$	$C_{22}$	...	$C_{2K}$	$n_{2\cdot}$
⋮	...	...	...	...	...
Cat K	$C_{K1}$	$C_{K2}$	...	$C_{KK}$	$n_{K\cdot}$
al drsh	$n_{\cdot 1}$	$n_{\cdot 2}$	...	$n_{\cdot K}$	$N$

→ marginal dist'n of Var 1

one row:  $\rightarrow$  row  $i$   
cond. dist'n of  $V_2 | V_1 = i$

one column (column  $j$ ):  
cond. dist'n of  $V_1 | V_2 = j$

marginal dist'n of  $V_2$

Recall: Independence Rules from Probability

$$P(A|B) = P(A) \rightarrow \text{conditional of } V_1 | V_2 = \text{marginal of } V_1$$

$$P(B|A) = P(B) \rightarrow \text{--- } V_2 | V_1 = \text{--- } V_2$$

$$P(A \cap B) = P(A)P(B) \rightarrow \text{joint} = (\text{marginal of } V_1) \times (\text{marginal of } V_2)$$

Can input contingency tables into chi-square tests for independence

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

More on this in Lab 04

→ contingency table in R  
→ lab 03

Pearson Residuals

→ deviations in what we observe vs. what we expected.

Pearson Residuals: Scaled difference between observed/expected

$$r_{ij} = \frac{C_{ij} - \hat{C}_{ij}}{\sqrt{\hat{C}_{ij}}} = \frac{\text{Obs.} - \text{Exp.}}{\sqrt{\text{Exp.}}}$$

$r_{ij} > 0$  : "too many" observed

$r_{ij} < 0$  : "too few" observed

$r_{ij}$ 's are asymptotically Normal!

Asymptotics

$|r_{ij}| > 2 \Rightarrow$  signif. at  $\alpha = 0.05$  level

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

Mosaic Plots → visualizes contingency tables

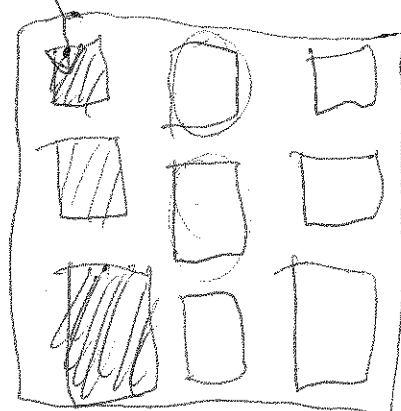
Mosaic Plots: Area plot for two categorical variables

each cell in contingency table gets a box in the plot  
area  $\propto$  % of obs. in the corresponding cell of the contingency table

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

height  $\propto$  % in Var 2 cat.  $j \mid \text{Var 1} = \text{cat } i \rightarrow$  conditional distns

Can color the boxes by their differences from what was expected:



Wednesday: Mosaic Plot Demo