## Outline

### Stat 849: ggplot2 graphics

Douglas Bates

University of Wisconsin - Madison and R Development Core Team <Douglas.Bates@R-project.org>

Sept 08, 2010

#### ggplot2

The pima data set from the faraway package

Univariate summary plots

Bivariate plots

Simple regression or ancova lines

Ancova

#### The ggplot2 graphics package

- Another advanced graphics package for R is ggplot2 by Hadley Wickham (a recent Iowa State Stats Ph.D., now at Rice).
- His book is listed as one of the references on the course web site.
- The core chapter introducing the basic function called qplot can be obtained from the URL in the links section on the course web site.
- I will use data from the faraway package to accompany Julian Faraway's freely available book "Practical Regression and Anova using R" to illustrate the use of qplot.

#### Examining the pima data

#### > head(pima)

	pregnant	glucose	diastolic	triceps	insulin	bmi	diabetes	age	test
1	6	148	72	35	0	33.6	0.627	50	1
2	1	85	66	29	0	26.6	0.351	31	0
3	8	183	64	0	0	23.3	0.672	32	1
4	1	89	66	23	94	28.1	0.167	21	0
5	0	137	40	35	168	43.1	2.288	33	1
6	5	116	74	0	0	25.6	0.201	30	0

## Recoding the missing data

- As Faraway indicates, several of the values of variables that cannot reasonably be zero are recorded as zero.
- A bit of research shows that these are missing data values.
   Also the test variable is a factor, not numeric.

```
> pima <- within(pima, {</pre>
      diastolic[diastolic == 0] <- glucose[glucose ==</pre>
+
           0] <- triceps[triceps == 0] <- insulin[insulin ==
+
           0] <- bmi[bmi == 0] <- NA
+
      test <- factor(test, labels = c("negative", "positive"))</pre>
+
+ })
> head(pima, 3)
  pregnant glucose diastolic triceps insulin bmi diabetes age
               148
                          72
                                  35
                                          NA 33.6
                                                      0.627 50
1
         6
2
                85
                          66
                                  29
                                          NA 26.6
                                                     0.351 31
         1
                                  NA
                                          NA 23.3
3
         8
               183
                          64
                                                     0.672 32
      test
1 positive
2 negative
3 positive
```

## Histogram of diastolic bp by test



# Histogram of diastolic blood pressure

> qplot(diastolic, data = pima, geom = "histogram")



## Empirical density plot



### Empirical density of diastolic by test

#### > qplot(diastolic, data = pima, geom = "density", linetype = tes



## Simple scatterplot, c.f. Fig. 1.2a, p. 13

#### > qplot(diastolic, diabetes, data = pima, xlab = ...)



### Adding a scatterplot smoother

> qplot(diastolic, diabetes, data = pima, geom = c("point", + "smooth"))



### Multiple smoothers by group

- > qplot(diastolic, diabetes, data = pima, geom = c("point",
- + "smooth"), shape = test)



## Comparative boxplots - apparently only vertical





## Adding a reference line - c.f. Fig. 1.3, p. 14

final -1 -2 midterm

> p + geom\_abline(intercept = 0, slope = 1, color = "red")

Adding a simple linear regression line - c.f. Fig. 1.3, p. 14

> (p <- qplot(midterm, final, data = stat500, geom = c("point",</pre> "smooth"), method = "lm")) +



## Suppressing the confidence band

- It happens that the defaults are intercept=0 and slope=1
  > (p <- qplot(midterm, final, data = stat500, geom = c("point",</pre>
- "smooth"), method = "lm", se = FALSE) + geom\_abline(color





Nave Height (ft)

Plotting multiple groups in separate panels



