

431 Class 26

thomaseLove.github.io/431

2021-12-02

Today's Agenda

- 1 Graphical and Numerical Summaries of Data
- 2 It's Just a Linear Model
- 3 432 preview: A `tidymodels` example with interaction
- 4 Takeaways from 431

Today's R Setup

```
library(knitr)
library(janitor)
library(magrittr)
library(mosaic)
library(equationmatic)
library(patchwork)
library(broom)
library(tidyverse)

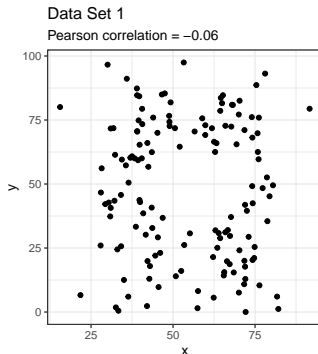
opts_chunk$set(comment=NA)
options(dplyr.summarise.inform = FALSE)
theme_set(theme_bw())
```

and a couple of secrets, hidden for now.

Visualizing Data

New Data Set 1

response	n	missing	mean	sd
y	142	0	47.83	26.94
x	142	0	54.27	16.77



13 Data Sets (summarized) in the df tibble:

	set	n	mean_x	sd_x	mean_y	sd_y	corr_xy
1	1	142	54.27	16.77	47.83	26.94	-0.06
2	2	142	54.27	16.77	47.83	26.94	-0.07
3	3	142	54.27	16.76	47.84	26.93	-0.07
4	4	142	54.26	16.77	47.83	26.94	-0.06
5	5	142	54.26	16.77	47.84	26.93	-0.06
6	6	142	54.26	16.77	47.83	26.94	-0.06
7	7	142	54.27	16.77	47.84	26.94	-0.07
8	8	142	54.27	16.77	47.84	26.94	-0.07
9	9	142	54.27	16.77	47.83	26.94	-0.07
10	10	142	54.27	16.77	47.84	26.93	-0.06
11	11	142	54.27	16.77	47.84	26.94	-0.07
12	12	142	54.27	16.77	47.83	26.94	-0.07
13	13	142	54.26	16.77	47.84	26.93	-0.07

New Data: Model for Set 1

```
set_1 <- lm(y ~ x, data = df %>% filter(set == 1))
```

```
tidy(set_1, conf.int = T, conf.level = 0.9) %>%  
  select(-statistic, -p.value) %>% kable(digits = 2)
```

term	estimate	std.error	conf.low	conf.high
(Intercept)	53.43	7.69	40.69	66.16
x	-0.10	0.14	-0.33	0.12

```
glance(set_1) %>%  
  select(r.squared, adj.r.squared, sigma, BIC, p.value) %>%  
  kable(digits = 3)
```

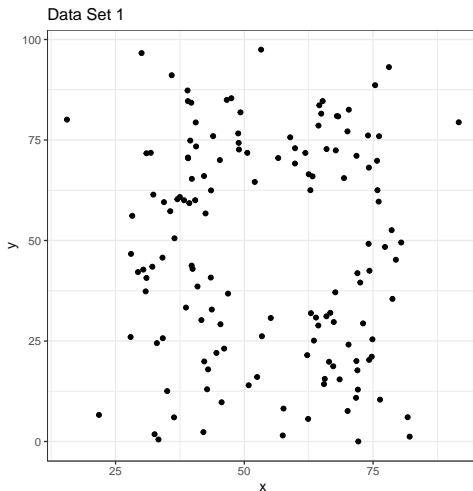
r.squared	adj.r.squared	sigma	BIC	p.value
0.004	-0.003	26.98	1351.64	0.448

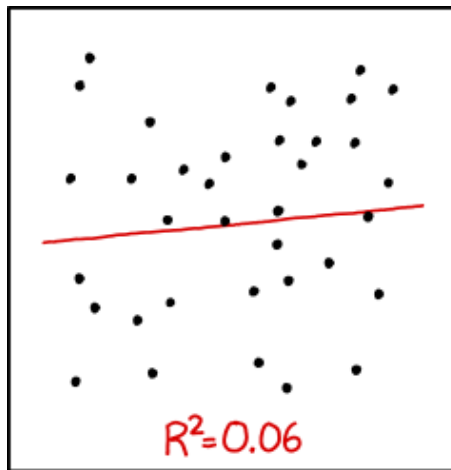
All 13 Models, at a glance

dataset	r.squared	adj.r.squared	sigma	AIC	BIC
1	0.004	-0.003	27	1343	1352
2	0.005	-0.002	27	1343	1352
3	0.005	-0.002	27	1343	1351
4	0.004	-0.003	27	1343	1352
5	0.004	-0.003	27	1343	1352
6	0.004	-0.003	27	1343	1352
7	0.005	-0.002	27	1343	1352
8	0.005	-0.002	27	1343	1352
9	0.005	-0.002	27	1343	1352
10	0.004	-0.003	27	1343	1352
11	0.005	-0.002	27	1343	1352
12	0.004	-0.003	27	1343	1352
13	0.004	-0.003	27	1343	1352

Plot for Data Set 1

Does a linear model for y using x seem appropriate?

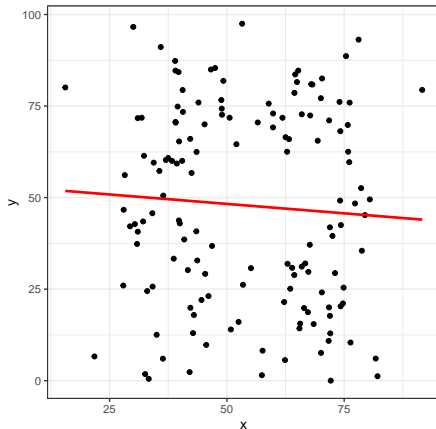




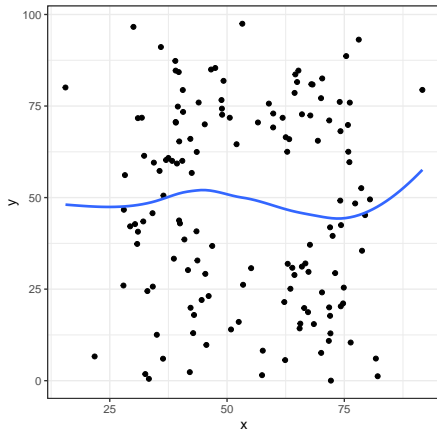
I DON'T TRUST LINEAR REGRESSIONS WHEN IT'S HARDER
TO GUESS THE DIRECTION OF THE CORRELATION FROM THE
SCATTER PLOT THAN TO FIND NEW CONSTELLATIONS ON IT.

Plot for Set 1 (with linear model and loess smooth)

Data Set 1 with lm fit



Data Set 1 with default loess smooth

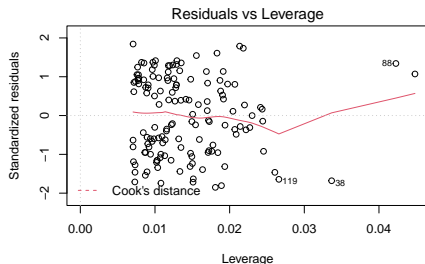
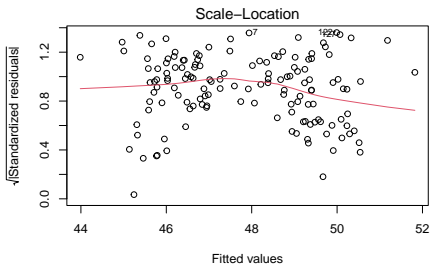
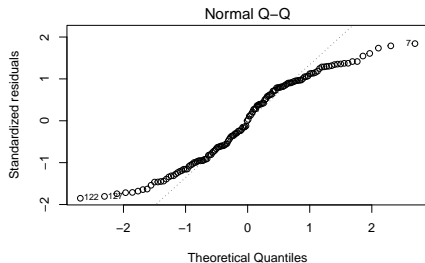
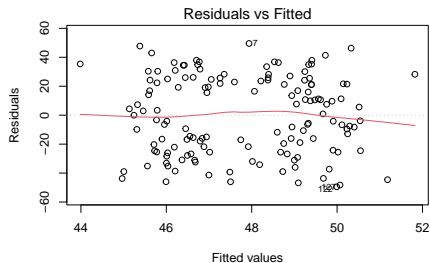


Model 1 (linear model for Set 1)

$$\hat{y} = 53.43 - 0.1(x)$$

(1)

Residual Plots for Set 1 Model

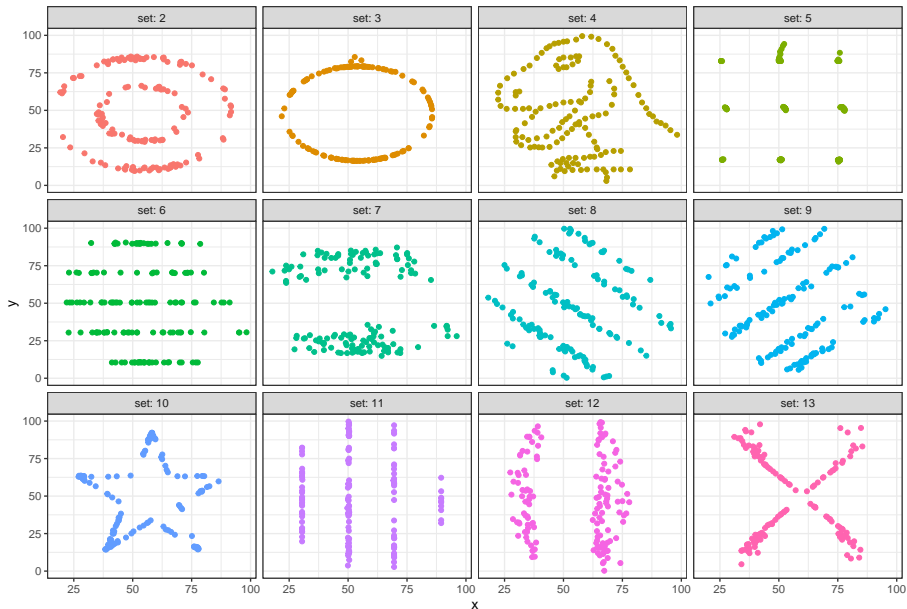


Models 2-13

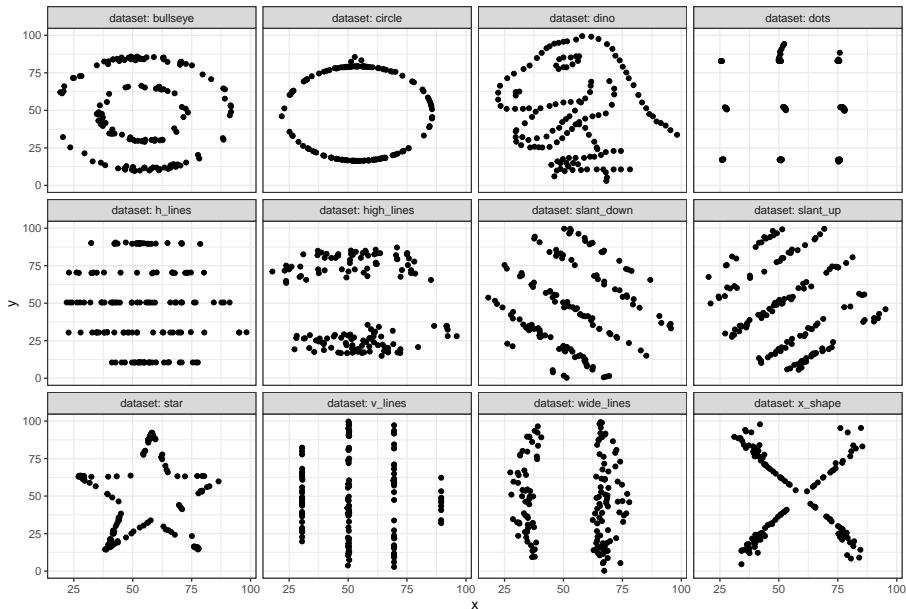
Models 2-13 all look about the same in terms of means, medians, correlations, regression models, but what happens if we plot the data?

```
df %>%  
  filter(set != 1) %>%  
  ggplot(., aes(x = x, y = y, color = dataset)) +  
  geom_point(show.legend = FALSE) +  
  facet_wrap(~ set, labeller = "label_both")
```

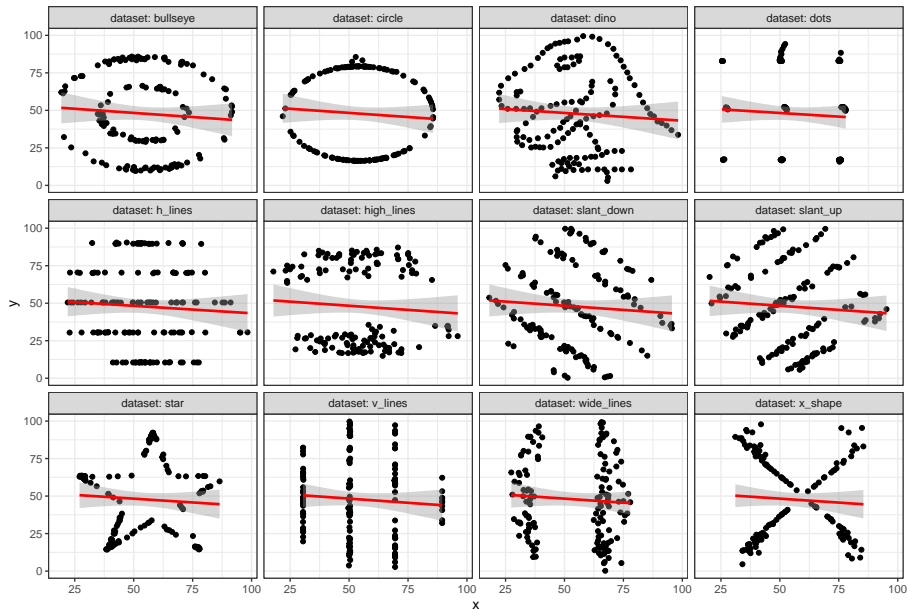
The Other 12 Data Sets



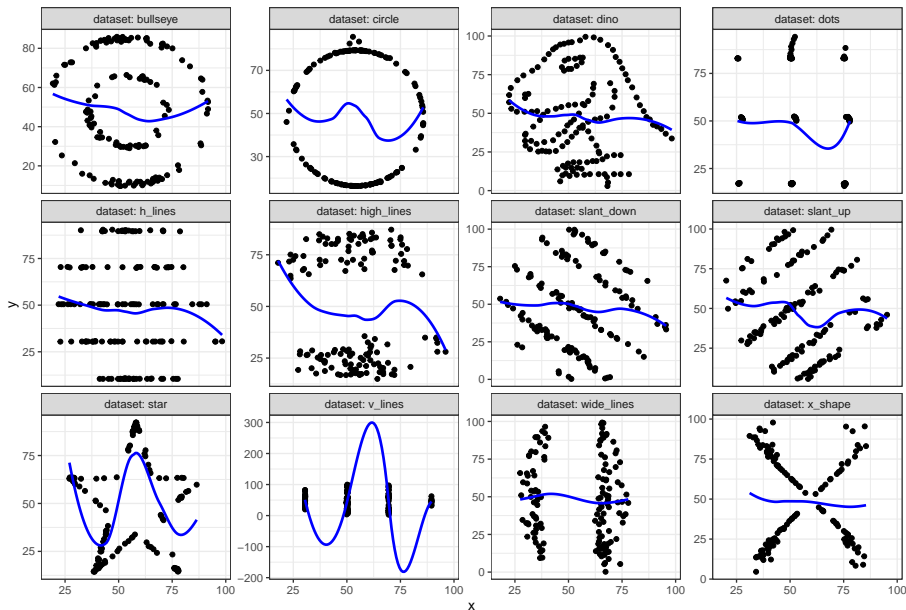
Actually, each of these sets has a name



And a linear model yields the same fit for each



How about a loess smooth with default span?



And the data come from

These are the datasauRus dozen data sets, available in the datasauRus package, which you can install from CRAN, thanks to the work of Steph Locke.

```
library(datasauRus)
df <- datasaurus_dozen
```

- You may recognize these data from their brief discussion in Spiegelhalter.
- These were created by Alberto Cairo, who has some great books like *How Charts Lie*

The moral of the story: **Never trust summary statistics alone, always visualize your data**

Two Cool Things, available online...

- 1 We'll visit Tomas Westlake's work at <https://r-mageddon.netlify.app/post/reanimating-the-datasaurus/>

```
library(datasauRus)
library(ggplot2)
library(gganimate)

ggplot(datasaurus_dozen, aes(x=x, y=y))+
  geom_point()+
  theme_minimal() +
  transition_states(dataset, 3, 1)
```

Two Cool Things, available online...

- 2 Next, we'll visit <https://www.autodesk.com/research/publications/same-stats-different-graphs>

This is Same Stats, Different Graphs: Generating Datasets with Varied Appearance and Identical Statistics through Simulated Annealing by Justin Matejka and George Fitzmaurice.

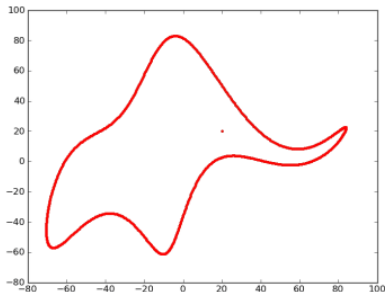
We'll look at a couple of the animated plots they generate there.

John von Neumann famously said

With four parameters I can fit an elephant, and with five I can make him wiggle his trunk.

By this he meant that one should not be impressed when a complex model fits a data set well. With enough parameters, you can fit any data set.

It turns out you can literally fit an elephant with four parameters if you allow the parameters to be complex numbers.



It's Just a Linear Model

Common Statistical Tests are Linear Models

Jonas Kristoffer Lindelov has built up a terrific resources to explain this at

<https://lindelov.github.io/tests-as-linear/>

What's the point?

Common statistical tests are linear models

Last updated: 02 April, 2019

See worked examples and more details at the accompanying notebook: <https://lindelev.github.io/tests-as-linear>

	Common name	Built-in function in R	Equivalent linear model in R	Exact?	The linear model in words	Icon
Simple regression: $\text{lm}(y \sim 1 + x)$	y is independent of x P: One-sample t-test N: Wilcoxon signed-rank	t.test(y) wilcox.test(y)	$\text{lm}(y \sim 1)$ $\text{lm}(\text{signed_rank}(y) \sim 1)$	✓ for $N \geq 14$	One number (intercept, i.e., the mean) predicts y. - (Same, but it predicts the <i>signed rank</i> of y.)	
	P: Paired-sample t-test N: Wilcoxon matched pairs	t.test(y1, y2, paired=TRUE) wilcox.test(y1, y2, paired=TRUE)	$\text{lm}(y_2 - y_1 \sim 1)$ $\text{lm}(\text{signed_rank}(y_2 - y_1) \sim 1)$	✓ for $N \geq 14$	One intercept predicts the pairwise $y_2 - y_1$ differences. - (Same, but it predicts the <i>signed rank</i> of $y_2 - y_1$.)	
	y ~ continuous x P: Pearson correlation N: Spearman correlation	cor.test(x, y, method='Pearson') cor.test(x, y, method='Spearman')	$\text{lm}(y \sim 1 + x)$ $\text{lm}(\text{rank}(y) \sim 1 + \text{rank}(x))$	✓ for $N \geq 10$	One intercept plus x multiplied by a number (slope) predicts y. - (Same, but with <i>ranked x</i> and y)	
	y ~ discrete x P: Two-sample t-test P: Welch's t-test N: Mann-Whitney U	t.test(y1, y2, var.equal=TRUE) t.test(y1, y2, var.equal=FALSE) wilcox.test(y1, y2)	$\text{lm}(y \sim 1 + G_1)^k$ $\text{glm}(y \sim 1 + G_2, \text{weights}=\dots)^k$ $\text{lm}(\text{signed_rank}(y) \sim 1 + G_2)^k$	✓ ✓ for $N \geq 11$	An intercept for group 1 (plus a difference if group 2) predicts y. - (Same, but with one variance <i>per group</i> instead of one common.) - (Same, but it predicts the <i>signed rank</i> of y.)	
	P: One-way ANOVA N: Kruskal-Wallis	aov(y ~ group) kruskal.test(y ~ group)	$\text{lm}(y \sim 1 + G_2 + G_3 + \dots + G_k)^k$ $\text{lm}(\text{rank}(y) \sim 1 + G_2 + G_3 + \dots + G_k)^k$	✓ for $N \geq 11$	An intercept for group 1 (plus a difference if $\text{group} \neq 1$) predicts y. - (Same, but it predicts the <i>rank</i> of y.)	
Multiple regression: $\text{lm}(y \sim 1 + x_1 + x_2 + \dots)$	P: One-way ANCOVA	aov(y ~ group + x)	$\text{lm}(y \sim 1 + G_2 + G_3 + \dots + G_k + x)^k$	✓	- (Same, but plus a slope on x.) <i>Note: this is discrete AND continuous. ANCOVAs are ANOVAs with a continuous x.</i>	
	P: Two-way ANOVA	aov(y ~ group * sex)	$\text{lm}(y \sim 1 + G_2 + G_3 + \dots + G_k + S_2 + S_3 + \dots + S_k + G_2*S_2 + G_3*S_3 + \dots + G_k*S_k)$	✓	Interaction term: changing sex changes the y - group parameters. <i>Note: $G_{2:k}$ is an indicator (0 or 1) for each non-intercept levels of the group variable. Similarly for $S_{2:k}$ for sex. The first line (with G_i) is main effect of group, the second (with S_i) for sex and the third is the group * sex interaction. For two levels (e.g. male/female), line 2 would just be "S_2" and line 3 would be S_2 multiplied with each G_i.</i>	[Coming]
	Counts ~ discrete x N: Chi-square test	chisq.test(group*xsex_table)	Equivalent log-linear model $\text{glm}(y \sim 1 + G_2 + G_3 + \dots + G_k + S_2 + S_3 + \dots + S_k + G_2*S_2 + G_3*S_3 + \dots + G_k*S_k, \text{family}=\dots)^k$	✓	Interaction term: (Same as Two-way ANOVA.) <i>Note: Run glm using the following arguments: glm(model, family=poisson())</i> As linear-model, the Chi-square test is $\log(y) = \log(N) + \log(\alpha) + \log(\beta) + \log(\alpha\beta)$ where α and β are proportions. See more info in the accompanying notebook .	Same as Two-way ANOVA
	N: Goodness of fit	chisq.test(y)	$\text{glm}(y \sim 1 + G_2 + G_3 + \dots + G_k, \text{family}=\dots)^k$	✓	(Same as One-way ANOVA and see Chi-Square note.)	1W-ANOVA

List of common parametric (P) non-parametric (N) tests and equivalent linear models. The notation $y \sim 1 + x$ is R shorthand for $y = 1 + b + a \cdot x$ which most of us learned in school. Models in similar colors are highly similar, but really, notice how similar they *all* are across colors! For non-parametric models, the linear models are reasonable approximations for non-small sample sizes (see "Exact" column and click links to see simulations). Other less accurate approximations exist, e.g., Wilcoxon for the sign test and Goodness-of-fit for the binomial test. The signed rank function is `signed_rank = function(x) sign(x) * rank(abs(x))`. The variables G_i and S_i are "dummy coded" indicator variables (either 0 or 1) exploiting the fact that when $\Delta x = 1$ between categories the difference equals the slope. Subscripts (e.g., G_2 or y_1) indicate different columns in data. `lm` requires long-format data for all non-continuous models. All of this is exposed in greater detail and worked examples at <https://lindelev.github.io/tests-as-linear>.

* See the note to the two-way ANOVA for explanation of the notation.

* Same model, but with one variance per group: `glm(value ~ 1 + G2, weights = varIdent(form = ~1|group), method="ML")`.



Jonas Kristoffer Lindelev
<https://lindelev.net>

Consider Study 1 from Project B.

- Analysis A. Compare two means/medians using paired samples
 - This is a linear model. See Section 4.2 of Lindelov

4.2.2 R code: Paired sample t-test

```
a = t.test(y, y2, paired = TRUE) # Built-in paired t-test  
b = lm(y ~ y2 ~ 1) # Equivalent linear model
```

Results:

model	mean	p.value	df	t	conf.low	conf.high
t.test	-0.5952	0.0934	49	-1.7108	-1.2944	0.104
lm	-0.5952	0.0934	49	-1.7108	-1.2944	0.104

Project B Study 1?

- Analysis B. Compare two means/medians using independent samples
 - This is a linear model, and not just for the t test. See Section 5 of Lindelov
- Analysis C. Compare 3-6 means/medians using independent samples
 - ANOVA is obviously a linear model, but actually we can generate (essentially) the Kruskal-Wallis this way, too. See Section 6.1 of Lindelov
- Analysis D. Create and analyze a 2x2 table
 - Yes, the chi-square test of independence can emerge from a linear model. See Section 7.2 of Lindelov
- Analysis E. Create and analyze a JxK table, where $2 \leq J \leq 5$ and $3 \leq K \leq 5$
 - Linear model, as in the 2x2 case. See Section 7.2 of Lindelov

Analyses D-E are more commonly thought about in the context of generalized linear models, as we'll see in 432.



📷 A recent count found 350 million purple sea urchins on one Oregon reef alone. (Shutterstock)

Sea Urchins and Tidy Modeling (a taste of 432)

Constable (1993) compared the inter-radial suture widths of urchins maintained on one of three food regimes

- Initial: no additional food supplied above what was in the initial sample
- Low: food supplied periodically
- High: food supplied ad libitum (as often as desired)

In an attempt to control for substantial variability in urchin sizes, the initial body volume of each urchin was measured as a covariate.

- This example comes from <https://www.tidymodels.org/start/models/>
- Another key source is <https://www.flutterbys.com.au/stats/tut/tut7.5a.html>
- Data from Constable, A.J. The role of sutures in shrinking of the test in *Helicidaris erythrogramma* (Echinoidea: Echinometridae). *Marine Biology* 117, 423-430 (1993). <https://doi.org/10.1007/BF00349318>

Package Load / Data ingest (Sea Urchins)

```
library(tidymodels)
library(readr)
library(broom.mixed)

urchins <-
  # Data were assembled for a tutorial
  # at https://www.flutterbys.com.au/stats/tut/tut7.5a.html
  read_csv("https://tidymodels.org/start/models/urchins.csv")
  # Change the names to be a little more verbose
  setNames(c("food_regime", "initial_volume", "width")) %>%
  mutate(food_regime =
    factor(food_regime,
           levels = c("Initial", "Low", "High")))
```

The urchins data

For each of 72 sea urchins, we know their

- experimental feeding regime group (`food_regime`: either Initial, Low, or High),
- size in milliliters at the start of the experiment (`initial_volume`), and
- suture width at the end of the experiment (`width`).

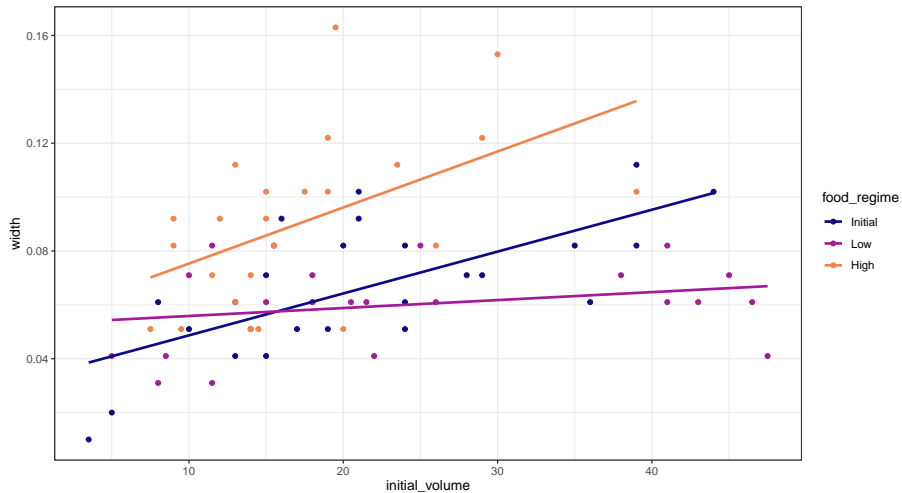
```
glimpse(urchins)
```

```
Rows: 72
```

```
Columns: 3
```

```
$ food_regime    <fct> Initial, Initial, Initial, Ini~  
$ initial_volume <dbl> 3.5, 5.0, 8.0, 10.0, 13.0, 13.~  
$ width          <dbl> 0.010, 0.020, 0.061, 0.051, 0.~
```

Plot the Data



How should we model the data?

Since the slopes appear to be different for at least two of the feeding regimes, let's build a model that allows for two-way interactions. We'll use a linear model for width which allows each food regime to generate a different slope and intercept for the effect of initial volume.

```
lm(width ~ initial_volume * food_regime,  
  data = urchins) %>% tidy() %>%  
  select(term, estimate) %>% kable(dig = 4)
```

term	estimate
(Intercept)	0.0331
initial_volume	0.0016
food_regimeLow	0.0198
food_regimeHigh	0.0214
initial_volume:food_regimeLow	-0.0013
initial_volume:food_regimeHigh	0.0005

Setting up a linear regression with tidymodels

```
lm_mod <-  
  linear_reg() %>%  
  set_engine("lm")  
  
lm_mod
```

Linear Regression Model Specification (regression)

Computational engine: lm

It turns out that we'll have several options for engines here.

We can estimate or train the model with `fit()`

```
lm_fit <-  
  lm_mod %>%  
  fit(width ~ initial_volume * food_regime, data = urchins)
```

We'll look at the results on the next slide.

What's in lm_fit?

```
tidy(lm_fit, conf.int = TRUE) %>%  
  select(term, estimate, conf.low, conf.high) %>%  
  kable(dig = 4)
```

term	estimate	conf.low	conf.high
(Intercept)	0.0331	0.0139	0.0523
initial_volume	0.0016	0.0008	0.0023
food_regimeLow	0.0198	-0.0061	0.0457
food_regimeHigh	0.0214	-0.0076	0.0504
initial_volume:food_regimeLow	-0.0013	-0.0023	-0.0002
initial_volume:food_regimeHigh	0.0005	-0.0009	0.0019

Make Predictions

Suppose that, for a publication, it would be particularly interesting to make a plot of the mean body size for urchins that started the experiment with an initial volume of 20ml. To create such a graph, we start with some new example data that we will make predictions for.

```
new_points <- expand.grid(  
  initial_volume = 20,  
  food_regime = c("Initial", "Low", "High"))
```

```
new_points
```

	initial_volume	food_regime
1	20	Initial
2	20	Low
3	20	High

Obtain Predicted Results for these new_points

We'll develop mean predictions and uncertainty intervals.

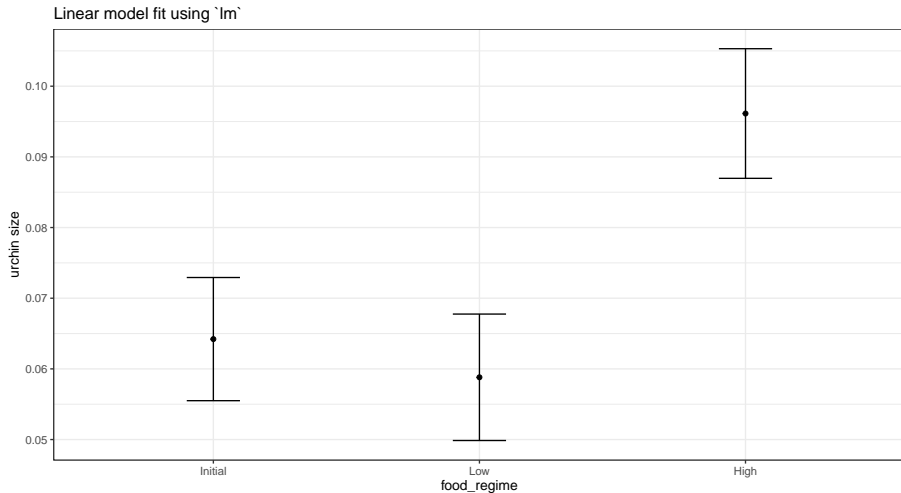
```
mean_pred <- predict(lm_fit, new_data = new_points)
conf_int_pred <- predict(lm_fit,
                          new_data = new_points,
                          type = "conf_int")

plot_data <-
  new_points %>%
  bind_cols(mean_pred) %>%
  bind_cols(conf_int_pred)
```

Plot the plot_data results (code)

```
ggplot(plot_data, aes(x = food_regime, y = .pred)) +  
  geom_point() +  
  geom_errorbar(aes(ymin = .pred_lower,  
                    ymax = .pred_upper),  
                width = .2) +  
  labs(y = "urchin size",  
        title = "Linear model fit using `lm`")
```

Plot the plot_data results (result)



But, I've just read something about Bayes?

Would the results be different if we used a Bayesian approach?

- Need to select a prior.
- Let's use bell-shaped priors on the intercepts and slopes, using a Cauchy distribution (works out to be the same as a t distribution with one degree of freedom)
- The `stan_glm()` function can be used, and this is available as an engine in `tidymodels`, where we need to specify `prior` and `prior_intercept` to fit a linear model.

Setting up a Bayesian Model

```
prior_dist <- rstanarm::student_t(df = 1)

set.seed(123)

bayes_mod <-
  linear_reg() %>%
  set_engine("stan",
             prior_intercept = prior_dist,
             prior = prior_dist)
```

Training the Bayes Model

```
bayes_fit <-  
  bayes_mod %>%  
  fit(width ~ initial_volume * food_regime, data = urchins)  
  
tidy(bayes_fit, conf.int = TRUE) %>%  
  select(term, estimate, conf.low, conf.high) %>%  
  kable(dig = 4)
```

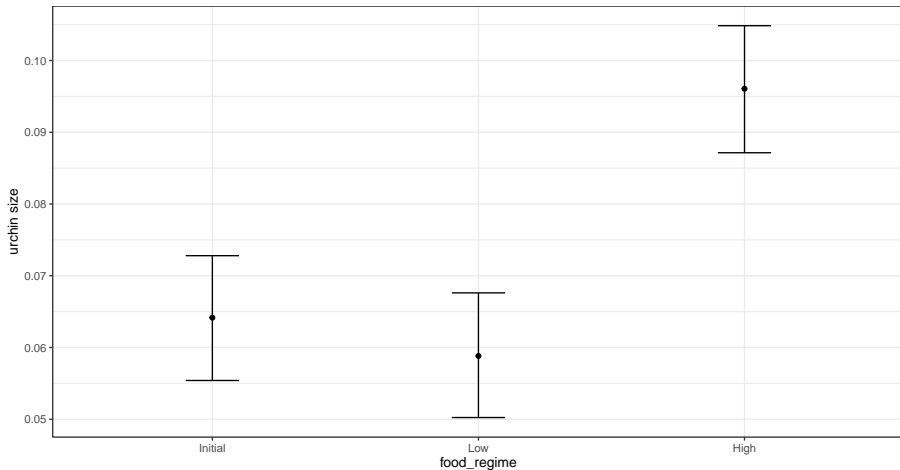
term	estimate	conf.low	conf.high
(Intercept)	0.0331	0.0172	0.0491
initial_volume	0.0016	0.0009	0.0022
food_regimeLow	0.0200	-0.0021	0.0415
food_regimeHigh	0.0216	-0.0046	0.0457
initial_volume:food_regimeLow	-0.0013	-0.0021	-0.0004
initial_volume:food_regimeHigh	0.0005	-0.0007	0.0017

Building the plot for the Bayes model (code)

```
bayes_plot_data <-  
  new_points %>%  
  bind_cols(predict(bayes_fit, new_data = new_points)) %>%  
  bind_cols(predict(bayes_fit, new_data = new_points,  
                    type = "conf_int"))  
  
ggplot(bayes_plot_data, aes(x = food_regime, y = .pred)) +  
  geom_point() +  
  geom_errorbar(aes(ymin = .pred_lower, ymax = .pred_upper),  
               width = .2) +  
  labs(y = "urchin size",  
       title = "Bayesian model with t(1) prior distribution")
```

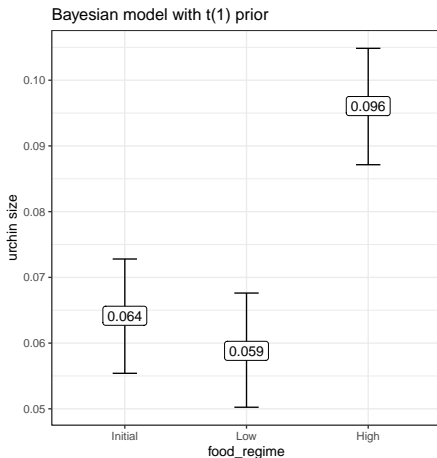
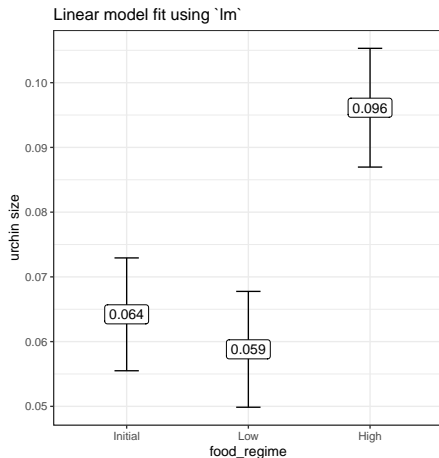
Plot for the Bayes model (result)

Bayesian model with $t(1)$ prior distribution



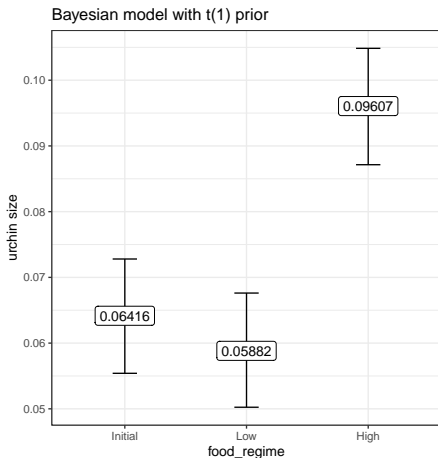
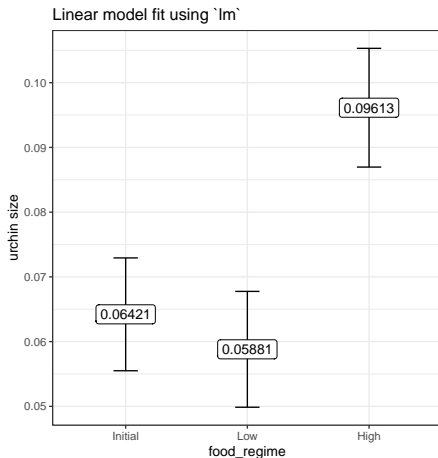
Comparing the Models

Comparing linear models for urchins data



Note that the models aren't actually the same

Comparing linear models for urchins data



What are we plotting, actually?

```
plot_data %>% kable(dig = 4)
```

initial_volume	food_regime	.pred	.pred_lower	.pred_upper
20	Initial	0.0642	0.0555	0.0729
20	Low	0.0588	0.0499	0.0678
20	High	0.0961	0.0870	0.1053

```
bayes_plot_data %>% kable(dig = 4)
```

initial_volume	food_regime	.pred	.pred_lower	.pred_upper
20	Initial	0.0642	0.0554	0.0728
20	Low	0.0588	0.0502	0.0676
20	High	0.0961	0.0871	0.1048

What do we take away from 431 at the end of the day?

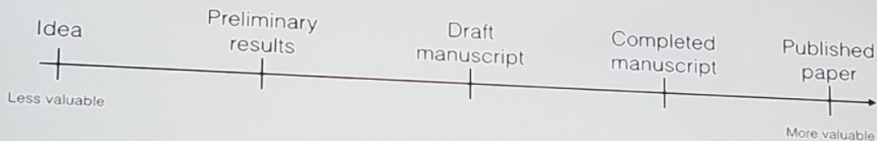
Ten Simple Rules for Effective Statistical Practice

From *PLoS Comput Biol* [link](#)

- 1 Statistical Methods Should Enable Data to Answer Scientific Questions
- 2 Signals Always Come with Noise
- 3 Plan Ahead, Really Ahead
- 4 Worry About Data Quality
- 5 Statistical Analysis Is More Than a Set of Computations
- 6 Keep it Simple
- 7 Provide Assessments of Variability
- 8 Check Your Assumptions
- 9 When Possible, Replicate!
- 10 Make Your Analysis Reproducible

A Tip from David Robinson

How I thought of my goals in grad school:



How I should have been thinking of them:

Anything still
on your computer

(Data, code, results,
draft, finished paper)

Anything out
in the world

(Paper, preprint, product,
blog post, open source,
tweet)



Build Tidy Data Sets

- Each variable you measure should be in one column.
- Each different observation of that variable should be in a different row.
- There should be one table for each “kind” of variable.
- If you have multiple tables, they should include a column in the table that allows them to be linked.
- Include a row at the top of each data table that contains real row names. `Age_at_Diagnosis` is a much much better name than `ADx`.
- Build useful codebooks.

Jeff Leek: “How to share data with a statistician” [link](#)

The Impact of Study Design (Gelman)

Applied statistics is hard.

- Doing a statistical analysis is like playing basketball, or knitting a sweater. You can get better with practice.
- Incompetent statistics does not necessarily doom a research paper: some findings are solid enough that they show up even when there are mistakes in the data collection and data analyses. But we've also seen many examples where incompetent statistics led to conclusions that made no sense but still received publication and publicity.
- We should be thinking not just about data analysis, but also data quality.

To consult the statistician after an experiment is finished is often merely to ask him to conduct a post mortem examination. He can perhaps say what the experiment died of. (R. A. Fisher)

$P > 0.05$



GAME OVER, TRY AGAIN

imgflip.com

What does sad p-value bear lead to?

So you collected data and analyzed the results. Now you want to do an after data gathering (post hoc) power analysis.

- ❶ What will you use as your “true” effect size?
 - Often, point estimate from data - yuck - results very misleading - power is generally seriously overestimated when computed on the basis of statistically significant results.
 - Much better (but rarer) to identify plausible effect sizes based on external information rather than on your sparkling new result.
- ❷ What are you trying to do? (too often)
 - get researcher off the hook (I didn't get $p < 0.05$ because I had low power - an alibi to explain away non-significant findings) or
 - encourage overconfidence in the finding.

None of this is particularly smart.

The Course So Far

- 1 Statistics is too important to be left to statisticians.
- 2 Models and visualization are the big takeaways, but don't forget about methods for making statistical inferences.
- 3 Reproducible research is the current wave.
- 4 Things are changing quickly. We live in interesting times.



That's all Folks!