

Final paper guide for PS 531

Jake Bowers

Spring 2017

Introduction

The paper in this class is a methods paper. And it is a term paper. As a term paper it should show that you understand the key concepts from the course. As a methods paper, it should explain the methods and evaluate them. Since it has the potential to something you publish as a part of another paper or even as a stand alone methods paper, it should look professional.

When I grade your paper, I ask myself: did the writer explain the methods clearly so that someone who is not a specialist can understand them? (this is the same question we ask of methods papers submitted for publication) did the writer show understanding of the key concepts in the course?

Here is a rough rubric. If you are doing a replication, then these are questions you may be asking about the paper that you are replicating and answering as you propose your own approaches (or as you agree with and justify the given approach). Each item receives equal weight in the final grade (except the code appendix, which receives half weight). When I came up with rubric, I was thinking especially about the reading from Berk (2004) and Achen (2004) Achen (2002), Achen (1986), Achen (1982):

- Writing

Can I read your writing? Is your paper written clearly? (I agree with Becker 1986's analysis of the evils of passive voice and the literature review.)

- Description

Why are you doing description? How useful is your description? If you are describing a kind of partial or adjusted relationship, how meaningful and believable is your adjustment?¹ Why did you summarize your outcome variable as you did? Recall that a linear model is a summary of an outcome variable as it relates to values of other variables. Subsidiary points here. Do

¹Recall that "adjustment" is "controlling for" in a linear model or "matching on" or "weighting by" or "stratifying on" etc..

you think about stuff like influential points, extrapolation, interpolation, meaningfulness of “controlling for”, overfitting.

- Statistical Inference I: Fundamentals

Why are you doing statistical inference? Recall “statistical inference” means hypothesis testing as single hypotheses using p -values or as collections of hypotheses not-rejected at some α -level of significance as confidence intervals. What are you inferring to (i.e. what is your target of statistical inference? is it potential outcomes? a population? a model (likelihood, fully Bayesian)?)?

- Statistical Inference II: Justification

How are you justifying the statistical inference that you are doing? What kind of repeated procedure (physical or theoretical (i.e. model based) or imaginary) makes it such that your study could have turned out otherwise?

- Statistical Inference III: Assessment

Can you provide evidence that your statistical inferential procedure fulfills its promises? Is the nominal size of Type I errors the same as the realized size of Type I errors? If you need a central limit theorem type of argument, can you justify it? If you are estimating something that you do not fully observe, is the estimation procedure/estimator unbiased and/or consistent and/or otherwise excellent in terms of mean squared error (i.e. trading bias for precision)?

- Final Paper Guide

Did you follow this guide?

- Code Appendix / Good Relationship with Future Self

Do you include an appendix that I could use to re-run all of the analyses presented in your paper at one go? That is, it must be a complete script, not a set of pieces that would require me to know how to cut and paste. Do you show that you have read Bowers and Voors (2017) on data analytic workflow?

What follows is more discussion about those general points.

Data summaries should reflect substantive questions.

First and foremost data summaries (like linear regression coefficients) ought to actually present the comparisons that we are using to make claims, that are implied by theory and explanation. Implications of this:

- If your claims are about democracy but your variable records amount of chocolate ice cream eaten per capita you have a lot to do to convince us that democracy is reasonably represented by chocolate ice cream eating. If

this statement seems overly self-evident or simple, read Adcock and Collier (2001) and see especially their Figure 1 on measurement.

- Is your concern about the total relationship between X and Y, or the partial relationship? Recall that a partial relationship is always defined relative to some specific other set of variables and functional form. If it is about the total relationship, then why are there control variables in your model? If it is about a partial relationship, then how can you guard against (a) interpolation and extrapolation of the kind discussed in Gelman and Hill (2007); (b) multi-dimensional influential points (i.e. points which would not be influential were it not for the addition of terms to your model);
(c) other micronumerosity and overfitting problems that occur as models become more complex [i.e. it does not take a lot of data to tell us something precise about a single difference of means, but it takes ever more data to tell us something precise about a difference of means conditioning on the kitchen sink.]²
- If your data summary mostly reflects the influence of just one observation, then you don't really have a very good summary, do you? That is, your summary does not reflect the substantive comparison you care about but rather it tells us a lot about Sweden. So here you must be alert to issues of influential points and overfitting, too.

Statistical inferences about data summaries should have clear targets of inference

Why are you doing statistical inference at all? What is unobserved and stochastic that demands your guessing? Are you asking about an unobserved population? About an unobserved series of repeated experiments? About your own (or the field's own) prior beliefs? If you can answer this question without saying, "Journal editors make me do it." then you have a clear target of statistical inference. The point here is not for you to solve the deep problem of what statistical inference means, but to show thoughtful engagement with the issues.

Procedures for making statistical inferences ought to fulfill their promises.

Now that you know why you want to do statistical inference, the next question is whether or not your chosen statistical inference machine works well given your design and data. In this class we learned about three targets of inference and then we learned about four ways to approximate these targets. The main big question here is whether the actual performance of the testing procedure the **size**

²This is a partial unpacking of some of the reasons behind Achen's arguments against garbage can regressions Achen (2004; Achen 2002).

of the test accords with the promised false rejection rate the **level** of the test (recall that confidence intervals are collections of tests and that “coverage” is just another way to talk about these concepts). One way to engage that question is to ask whether the **sampling distribution, randomization distribution, or posterior distribution is the shape that our approximations assume**. Across repeated samples/repeated experiments are the t -statistics really t -distributed? Are the z -statistics really Normally distributed? Since we tend to use approximations which themselves require certain simplifications and assumptions, we ask whether such secondary and tertiary type of assumptions are credible *prima facie* and then we might assess such a hunch using simulations as we have done.

Target 1: a population of units (moments in time, people, country-years, ...)

In the ideal case, we can take many samples in the same way without replacement from the same population.

Since we never have this situation, instead we approximate this target of inference (the population) by repeatedly sampling from our sample in the same way that we sampled from the population [all we need is the rule by which we did the sampling, not the list of all of the units in the population]. (The bootstrap)

It turns out that if our rule was an iid kind of rule and our samples were large but very small relative to the population, then we can know that our sampling distribution is Normal (and if we have to estimate the spread of the Normal distribution then this sampling distribution is the t distribution). (Central Limit Theorems, Asymptotic theory)

And, even if didn't know how our units were sampled, but we did know how our outcome variable took on the values it does, then we can proceed as if we had a sampling process and central limit theorems can help us know the sampling distribution of our estimators. (Where the likelihood function takes the place of a population and sampling process).

Target 2: an experiment repeatedly administered to the same units (values of an experimental treatment variable, values of an explanatory variable)

We approximate this by repeatedly re-assigning treatment in the same way that we assigned treatment in the first place. So, again, we need the rule or algorithm by which we assigned treatment. (The permutation approach in our class proceeded by approximating the enumerated permutation distribution by sampling from it).

It turns out that if our rule was of an iid kind of rule and our experimental pool was fairly large and not too few people were assigned treatment, then we can

know that our randomization distribution is Normal. (Central Limit Theorems, Asymptotic theory)

Target 3: scientific uncertainty itself

We only really quickly and cursorily engaged with this target of inference — the posterior distribution of θ where θ (could be a vector) contains the quantities of scientific interest to us, and where we can express our prior uncertainty about θ in the form of a probability distribution.

My expectations

I assume that you will try to do statistical inference about one of the two frequentist targets of inference. I assume you'll use all of the different ways one may approximate such targets (for example, if you are interested in some population, you'll use the bootstrap, the likelihood, and the random-sampling iid+CLT approaches but you won't permute). In this case, you might prefer a Bayesian approach to likelihood or you might not.

That said, there are things we did in class which were meant to illustrate concepts and not to be used in your simulations (i.e. to make the assumptions of iid+CLT or other assumptions more vivid and meaningful). For example, I showed how taking larger and larger samples led to both tighter confidence intervals and I also did something similar when we talked about the central limit theorem (CLT) and law of large numbers. You should know what kinds of things we did in class to illustrate a theoretical point and which things we did to actually do statistical inference or assess the properties of one or another methods for doing statistical inference.

I assume you'll grapple with the question about which approximation does an adequate job fulfilling its promises using both what you know about the theory justifying the approximations (i.e. the stories you have to believe to believe the inferences) and also using simulations of some kind to convince yourself and others that the theory is useful or not given the particulars of your data and design and model.

More detailed points

On the assumptions behind canned lm/glm tables

If you want a list of the standard assumptions of iid+CLT least squares and/or mle based statistical inference consult Fox, Achen, Kennedy, or many other available sources. Yet, remember, Normality of outcomes or errors is **not** a

crucial assumption unless you are using some likelihood based approach (remember (Achen 1982) for how and why). For the MLE/Bayes based justification, Normality of outcomes does matter. But for the large-sample, iid+CLT based story of regression (as well articulated by Achen), Normality arises in the course of the iid+CLT.

If you are not using a simple linear regression but some glm model (logit, probit, poisson, negative binomial, ...)

Then it is not adequate justification to say, “I have counts...” or “My regression could, in theory, predict values lower than zero or higher than 1.” Nor is “heteroskedasticity” alone a reason to leave the simple linear regression model. The bootstrap, for example, makes no assumptions about homoskedasticity. And, in many cases, a linear fit to a binary outcome will not predict below 0 or 1 (or outside the range of the counts).

This is not to say that you can’t or shouldn’t do logit/Poisson/negative-binomial models, but that you should be know which justifications are strong and which are weak. For this paper, it is enough to say, “I found this negative binomial model” or “Others are using this.” — putting the burden of justification on someone else. It is fine to *start there* in order to deeply engage with / support / criticize such decisions. And if you want to justify said model, then that is fine. But I encourage you to spend a moment thinking about how such a move improves upon the basic smoothed difference of means of the simple regression. In the end, you may desire a more flexible functional form for \mathbf{Xb} rather than assuming that a particular link function is useful for you. For example, you might ask about the substantive reason for the logit versus the cloglog link versus probit link in the binomial glm.³

On assumptions:

It is not enough to claim that some model involves assumptions, and that said assumptions are not satisfied, in order to justify use of another model. The other model will have often as many if not more assumptions as the model left behind; it too would be easily discarded under the rhetoric of “makes assumptions.”

A very common move in this regard when using binary outcomes is to say, “Linear regression assumes linearity and allows predicted values beyond 0 or 1. Therefore I use the logistic regression/logit model.” Notice, that the logit model and the linear regression both assume the same \mathbf{Xb} , but the logit model *adds* the assumption that one should mash \mathbf{Xb} into the $[0,1]$ interval via the

³There are good answers to this question, by the way, especially in the economics literature on latent choices and in the epidemiology literature. And there is a large literature on nonparametric smoothing and other ways to fit functions to conditional moments of outcomes without making the strong assumptions of (generalized) linear models. See here for some nice alternative link functions for binary outcomes.

logit link function $g(\mathbf{X}\mathbf{b})$ where $g()$ is the logit function. So, (1) both require linearity in how \mathbf{X} relates to \mathbf{b} (ex. neither estimate b in x^b) and (2) logit adds another assumption on top of linearity and now requires a correct likelihood function whereas Achen shows that the OLS model does not require a correct likelihood function when n is large. In addition, if the linear regression does not predict beyond 0 or 1 in your particular data with your particular regression, then using the logit transform solely for the purpose of limiting $\hat{\mathbf{y}}$ to $[0,1]$ is no longer necessary. In this case, the homoskedasticity requirement for the standard errors of the linear model is assumed indirectly in the logit model in which all of the observations are assumed to come from the same binomial/Bernoulli data generating process with p or π (the probability of a 1) parameterized by $g(\mathbf{X}\mathbf{b})$. In fact concerns about the variance of the binomial outcomes are very common

Recall that we worked on different link functions in class. It will not be appropriate to just use a simple logit or probit model (or poisson or negative binomial, etc..) without discussing and thinking about the link function problem let alone directly confronting it.

Of course, we don't care about homoskedasticity if we are using design-respecting permutations/shuffling or bootstrap resampling. And, even if we are using iid+CLT justified statistical inference, as Fox noted in his textbook, heteroskedasticity is usually the very least of our worries.

There *are* good reasons to use logit (or negative binomial models) and this class has been meant, in part, to help you understand which assumptions matter for the linear regression model and for what reasons. The best reasons to do something are positive: logit may offer you something that would help you rather than merely be a way to not do simple linear regression.

On bad models

Some of you may be replicating what amount to be fairly shallow or basically bad linear models. That is fine and perhaps more fun for you. That said, you should feel free to improve upon those models. Remember, this paper is for you, not for the authors of the tables you are replicating (should you be scrutinizing other peoples' tables and not your own). So, feel free to change their models so that they are more meaningful.

R Code should be an appendix to the paper.

If you are using R markdown or knitr to produce your paper you can ensure that each code chunk has a label, like the following:

The Knitr Way

The regression of y on x is simple and shows no relationship, as you can see in Table~\ref{tab:lm1}.

```
<<mymainregression,eval=TRUE,echo=FALSE,results='hide'>>=
## Fit the model y~x on the data.frame mydata
lm1<-lm(y~x,data=mydata)
@
```

```
<<tableformainregression,echo=FALSE,results='asis'>>=
## Make a table of the regression results
xtable(summary(lm1)$coef,label="tab:lm1")
@
```

Then you can repeat `<<mymainregression>>=` by just referring to it in your appendix.

```
\appendix
\section{The Appendix}
Here is the code appendix.
```

```
<<appendix,eval=FALSE,echo=TRUE,results=verbatim>>=
<<mymainregression>>
<<tableformainregression>>
@
```

The R+Markdown way

The regression of y on x is simple and shows no relationship, as you can see in Table~\ref{tab:lm1}.

```
```{r mymainregression,eval=TRUE,echo=FALSE,results='hide'}
Fit the model y~x on the data.frame mydata
lm1<-lm(y~x,data=mydata)
```
```

```
```{r tableformainregression,echo=FALSE,results='asis'}
Make a table of the regression results
xtable(summary(lm1)$coef,label="tab:lm1")
```
```

Then you can repeat `<<mymainregression>>=` by just referring to it in your appendix.

```
# The Appendix
Here is the code appendix.
```



```

```{r appendix,eval=FALSE,echo=TRUE,results=verbatim}
<<mymainregression>>
<<tableformainregression>>
```

```

Alternatively, if you are using knitr you can collect the code using `purl("mypaper.Rnw")` from within R. Edit the resulting `mypaper.R` file adding `\begin{Verbatim}` at the top and `\end{Verbatim}` at the bottom. Then you can just do something like:

```

\section{R Code Appendix}
\input{mypaper.R}

```

You can do something similar with `.Rmd` files.

If you are not using LaTeX, then you can just add the R command file to your Word/OpenOffice/Pages document. I am assuming that you will be able to rerun all of the analyses in your papers without cutting and pasting by just doing “`source(mycode.R)`” in R or R CMD BATCH `mycode.R` from the unix command line (I’m not sure what the Windows command line is like these days).

More on literate programming.

- OpenOffice has it’s own kind of literate programming possible called `odfWeave`.
- You can use LyX too LyX with knitr
- Of course you can use LyX with Sweave.
- IPython looks very promising, too.

In case you are wondering about editors, I just switched back to MacVim and also just Vim in the terminal with this plugin. I’ve also happily used Emacs and to Emacs from Aquamacs as my editor and others really like Textmate (the `emacs`en required `auctex` and `ESS` (`emacs` speaks statistics) plugins). I don’t know the state of editors for Windows (although Emacs and Vi will run there as they do on Linux too). Once social science computing environments left the VAX, I used Emacs on Windows and Emacs and `vi/vim` on different flavors of Linux and Solaris.⁴

⁴Jeff Gill and I periodically maintain a list of stuff to help folks who use Macs make them useful scientific workstations: <https://github.com/jwbowers/SocialScienceMacConfig>. Feel free to fork and contribute.

About writing:

If you don't already have it, get and read Becker's Writing for Social Scientists (Becker 1986). Do whatever he says to do.

Here is a list of some of what I think are best practices for writing scholarly papers. Not everything on this list applies directly to the final paper for this class. So, in no particular order, here are my thoughts:

- Dump all sections entitled "Literature Review"

Read Becker (1986) for his chapter on using literature. The point of using past literature is to *motivate* your paper. Demonstrating knowledge of the literature occurs as you *situate* and *frame* your paper — as you show how your ideas are great and important and novel and different from and building upon but in conversation with past ideas. Talking about a literature without linking each and every sentence to your own particular project is like putting an undergraduate mid-term exam into the middle of a scholarly paper. It is boring, distracting, and inappropriate. Don't do it. If you find yourself labeling a section, "Literature Review" rather than "Why my ideas are so great and interesting" stop and ask yourself why you are doing this.

- Motivate, motivate, motivate

Readers should not have to ask "So?" or "So what?" or "Who cares?" I call paper that generates a lot of "So what?" questions *unmotivated*.

Even if you think that you are answering the "So what?" question, often you are not: people often misunderstand what it means to adequately answer the "So what?" question.

Two examples of common misguided attempts at motivation illustrate:

- (1) Motivation attempt # 1 is to say "The literature has a hole." and even to point to the hole. Of course, the existence of something unstudied does not, in fact, mean it ought to be studied — or if it is studied that it will yield great and deep insights into the causes of war or voting or violence or democratic transitions. To stop at noticing the hole in the literature is to leave the paper unmotivated. To explain why the hole ought to be filled is to motivate the paper.
- (2) Motivation attempt #2 consists of claiming that some variation in (presumably causal) effects across some subgroup has not been studied before. This is a form of the first failed motivation attempt. The hole consists of a lack of knowledge about subgroup specific effects or about extent of variation across subgroups. Again, however, merely not having looked at such variation doesn't mean that so looking ought to teach us a lot about life, the universe, and everything (let alone about some important bit of scientific understanding).

I do *not* mean that you have to talk about Rousseau and Hobbes (or Douglas Adams) necessarily, but that you have to talk about the theoretical and practical/policy reasons why anyone should be reading your paper. If the real reason you wrote the paper is “Because (for some reason) I want a Ph.D. and this paper was required for a class.” then you should consider writing a different paper. There is very little room in graduate school (let alone life) for you to write papers that you don’t care about and that you can’t convince other people to care about, too.

- Answer “So what?” constantly

Even if you have a page or two at the beginning of your paper explaining how it contributes important knowledge about important topics (raises important questions, generates important new concepts, etc.), readers should not have to ask *why* they are reading a giving sentence, paragraph, or phrase. That is, the need for motivation is not over once you’ve written a couple of paragraphs. Every piece of the paper should be directly and clearly intelligible as contributing to the “So what?” question answer.

Even if I am not a specialist in a particular area, I should still be able to say why this paper matters, why this section contributes to the paper (in a way that is linked to it being important), why this paragraph is here, why this sentence, and why this word. That is, the motivation for the paper in general ought to permeate the entire paper. Each sentence should impel the narrative forward. People should *want* to read the next sentence. They should be grateful for the next section since they are so desperate to find out the answer to the important mystery you posed at the beginning. And they should be convinced that this next section is precisely the perfect section to be reading at exactly that moment in the paper.

- Think about clarity rather than sophistication

Again, I refer you to Becker (1986). Let us say that you write something like, “X can explain Y.” What does this mean? Are you referring to a particular kind of statistical model? If so, are you talking descriptively? Predictively? Causally? What about X “explains” what about Y? What does it mean for some thing to “explain” something else? Is X like a professor trying to give a lecture about something complicated called Y? This is *jargon*. And worse, it is *unclear*.

What you probably mean is that X and Y are things called variables — columns of numbers, vectors — and that when X has higher values, Y also tends to have higher values. Another way to say this is that they co-vary. Or, using scientific language in a clear way, you can say that, the R^2 from a linear regression of Y on X is substantial (XX%). Perhaps you mean that if we could change the values of X for some unit of observation then we would almost surely see, as a result, a change in the values of Y for that unit. Notice that *scientific language* is not a problem. But when does scientific language becomes jargon? It is fine to use technical language

and theoretical language, but it is not fine to use such language either incorrectly or obscurely.

Hannah Pitkin has said that she writes for college freshmen. This means that she aims to talk to smart and motivated people who do not know a lot of particular technical language or particular literature. I like this rule of thumb. If I can say “A bright and interested but inexperienced person can understand this sentence”, then I know that my writing has a fighting chance at clarity.

- Interpret, Interpret, Interpret and Explain, Explain, Explain

A regression table or a plot doesn’t talk about itself. That is your job. And your job is to *interpret* them — that is, *make them meaningful* to the reader and the field at large. And, you might need to *explain* them a bit as you interpret.

This exhortation also holds for claims and arguments that you state. You should be citing appropriate literature, explicating and explaining and defending and justifying the things that you say. (Notice that if I write “means what?” or “So?” next to some text then it is both unmotivated as well as unexplained, undefended, unjustified, and thus uninterpreted.)

- Do not focus on “statistical significance”. Focus on the Real Story

- (1) When it comes to OLS, the assumptions required to justify hypothesis testing (and thus the production of p-values and confidence intervals) are much more stringent than those required to justify the estimation of coefficients. Thus, focusing on statistical significance rests your argument on the weakest parts of your statistical model.
- (2) In addition, statistical significance is not the real story. The real story is that if people differ on X (say $X=0$ vs $X=1$) then, on average, they also tend to say $Y=1$ rather than $Y=0$ — and this difference is the difference between elections won and lost, between people fed and starving, between views of the world, etc... You should never present a regression table without talking about the substantive implications of the predictions and coefficients of the model. Never. *Never*. I would much much prefer that you not mention the results of hypothesis tests rather than you not mention the substantively important predictions, effects, differences implied by your model.
- (3) Talking about statistical significance allows you to avoid talking about your research design at precisely the moment when your design ought to be helping you make your case. Talking about the predicted changes and effects and differences that come out of your model as they relate to your research design allows you to use multiple strong arguments together — rather than just one weak argument.
- (4) In general, focusing on statistical significance in current Political Science is a sign of lack of substantive knowledge and lack of theoretical

engagement let alone lack of methodological depth.⁵

- Excessive passive voice allows sloppiness and cowardice

Again, read Becker (1986) on this. By using the passive voice you allow someone other than you to act: you abdicate responsibility for your words without attribution. Of course, sometimes passive voice is fine. The problem occurs with excessive or ill-placed passivity (For example, when the crucial sentences justifying your decisions suddenly become passive, then the reader knows that you don't feel very confident about your assertions).

Logical Paper Structure should be Evocatively Denoted

- This means that section headings should describe and suggest. "Literature Review" does not impart forward motion to your story. "A common misunderstanding about war and democracy, and how it can be clarified" is much more exciting and foreshadows your argument.

Tables and figures should stand alone

- A lazy reader should be able to look at the table or figure, and from the notes, caption, and label be able to make sense of it. You shouldn't tend to force readers to comb through the text to understand a figure or table.

All specifications of all models should be written out in equation format

- There is just too much possibility for slippage between a table and what is actually produced at the end of some computational algorithm. And talking at length in a paragraph can also be misinterpreted. If you are using a linear model (be it logistic, Poisson, or OLS), a simple sketch of what goes into your linear specification (\mathbf{Xb}) really clarifies exactly what you are estimating — especially since not all coefficients from all models must be presented in all tables.

Use full and informative labels in all figures, tables and in text

- Acronyms like EducCat3 or CONFDUR obfuscate. Those kinds of labels work fine for memos between collaborators but not for papers for general reading and publication. The general rule, as I see it, is to make life easy for readers. Be grateful that anyone is reading your paper at all. Make life easy and exciting and pleasant for them — or else they will stop reading, or if forced to continue reading (because they are reviewing your article for publication, tenure, or funding) they will skim, or worse, get annoyed.

Assumptions of the statistical model should be directly engaged

⁵Note: For PS531 we are engaging directly with questions about statistical significance and the meaningfulness of such statements. So, you **must** talk about statistical significance in our class, but you must do it thoughtfully, backed up by statistical theory and your own simulations. Even so, the real story is the substantive one, not the statistical one.

- I mean this for every paper — not just papers for this class. Great writers and scholars can do this in one paragraph and a couple of footnotes. Every statistical procedure requires a set of logical, mathematical, and algorithmic commitments in addition to commitments about how the data at hand arrived into your hands (the research design): you should be up-front about such commitments and what they mean for the substantive interpretation of your results and for the status of your contribution to our shared understanding about the world.

Formatting

- (1) Papers should have abstracts. (2) Fonts should not be less than 11 pt font, with no less than 1 inch margins (I like 1.5 inch margins since I write in the margins rather than in between lines). (3) Tables within papers should have as few vertical lines as possible and never use double rules (some might say, never use vertical lines) Chicago (2003, sec. 13.51–§13.53). (4) If at all possible, ensure that you don't have annoying justification or footnoting problems in your manuscripts (i.e. single lines with excessive space between words because the justification algorithm has failed; footnotes ill-placed; etc. . .)

Consider your Ink to Information Ratio

- See Tufte (1983) for more on this. Big areas of your figures ought not consist of color or other ink that does not carry information. For example, the default Stata graphics with a big colorful border has a very low ink to information ratio: What does that border mean? If it is not meaningful, why is it there? Why present it to me? Default Excel graphics with the body of the graph in gray are even worse on this metric since the ink is all over the plotting area itself. And you already know about 3D bar charts: Ugh!

Figures and Tables and Notes should be in the text itself

- I know that some journals like tables and figures at the end of the paper. Some even like endnotes. I hate this. It is terrible for the reader to have bounce back and forth through your document.. When you give stuff to me to read, please put the tables and figures in the text itself and avoid endnotes in favor of footnotes.

Be polite

- Notice something about these points about ink, formatting, captions, tables and figures: They amount to a requirement that you be **polite** to a reader, that you create a paper that you would find easy to read (cf. The Golden Rule). I think that readers become annoyed sometimes because they feel that writers are not respecting their time, attention, and energy. So, be a polite writer. Show that you've spent time and energy making the reading experience as nice as possible.

Decide on number of digits based on substance

- If your data allows you to distinguish between .001 and .002 in substantive terms, then, by all means use 3 digits, if not, then use the number of significant digits that is substantively meaningful. I have never seen a political science problem where one could detect substantively important differences at 5 (let alone 8) decimal places. But, if you are really interpreting your modeling results and your summary statistics keeping the theory and substance of your problem in the forefront of your reader's attention, then, even 8 decimal places will make sense if they are in fact sensible. (If you are cursory and hasty and technical in your interpretations, of course, you will just annoy and distract your reader with your misplaced precision.)

Understand Dashes and Hyphens

- A hyphen (-) combines two words into one word or breaks words at end of lines: short- or long-term; pick-me-up. The em dash (—) and the en dash (–) are forms of punctuation. Use the en dash to indicate ranges like 1990–2000 or relationships like Student–teacher. And use the em dash to indicate a parenthetical thought — as suggested in (Chicago 2003)[§ 6.80–§ 6.96].

The hyphen, the en dash, and the em dash are the most commonly used [of the hyphens and dashes] and must be typeset correctly; an en dash appearing where a hyphen is called for bespeaks editorial or typographic confusion. Chicago (2003, sec. 6.80).

“Impact” does not produce much impact

- This from *New Oxford American Dictionary* on my computer:

USAGE The phrasal verb **impact on**, as in: *when produce is lost, it always impacts on the bottom line*, has been in the language since the 1960s. Many people disapprove of it despite its relative frequency, saying that **make an impact on** or other equivalent wordings should be used instead. This may be partly because, in general, new formations of verbs from nouns (as in the case of impact) are regarded as somehow inferior. As a verb, impact remains rather vague and rarely carries the noun's original sense of forceful collision. Careful writers are advised to use more exact verbs that will leave their readers in no doubt about the intended meaning. In addition, since the use of impact is associated with business and commercial writing, it has a peripheral status of ‘jargon,’ which makes it doubly disliked.

References

- Achen, Christopher H. 1982. *Interpreting and Using Regression*. Newbury Park, CA: Sage.
- . 1986. *The Statistical Analysis of Quasi-Experiments*. Berkeley, CA: University of California Press.
- . 2002. “Toward A New Political Methodology: Microfoundations and ART.” *Annual Review of Political Science* 5 (1). [http://arjournals.annualreviews.org.proxy.lib.umich.edu/doi/abs/10.1146/annurev.polisci.5.112801.080943?prevSearch=allfield:\(achen\):423-50](http://arjournals.annualreviews.org.proxy.lib.umich.edu/doi/abs/10.1146/annurev.polisci.5.112801.080943?prevSearch=allfield:(achen):423-50).
- . 2004. “Let’s Put Garbage—Can Regressions and Garbage—Can Probits Where They Belong.”
- Adcock, Robert, and David Collier. 2001. “Measurement Validity: A Shared Standard for Qualitative and Measurement Validity: A Shared Standard for Qualitative and Quantitative Research.” *American Political Science Review* 95 (3): 529–46.
- Becker, Howard S. 1986. *Writing for Social Scientists : How to Start and Finish Your Thesis, Book, or Article (Chicago Guides to Writing, Editin*. University Of Chicago Press.
- Berk, Richard. 2004. *Regression Analysis: A Constructive Critique*. Sage.
- Chicago, The University of. 2003. *The Chicago Manual of Style*. 15th ed. The University of Chicago press.
- Gelman, A., and J. Hill. 2007. *Data Analysis Using Regression and Multi-level/Hierarchical Models*. Cambridge University Press.
- Tufte, Edward. 1983. *The Visual Display of Quantative Information*. Cheshire, CT: Graphics Press.