

Lesson 12: Model/Variable Selection

Nicky Wakim

2024-03-04

Learning Objectives

1. Understand the motivation for model selection, including bias-variance trade off and alignment of research goals (association vs. prediction)
2. Explain the general process or idea behind different model selection techniques
3. Recognize common model fit statistics and understand what they measure

Learning Objectives

1. Understand the motivation for model selection, including bias-variance trade off and alignment of research goals (association vs. prediction)
2. Explain the general process or idea behind different model selection techniques
3. Recognize common model fit statistics and understand what they measure

Why can't I just throw in all the variables into my model?

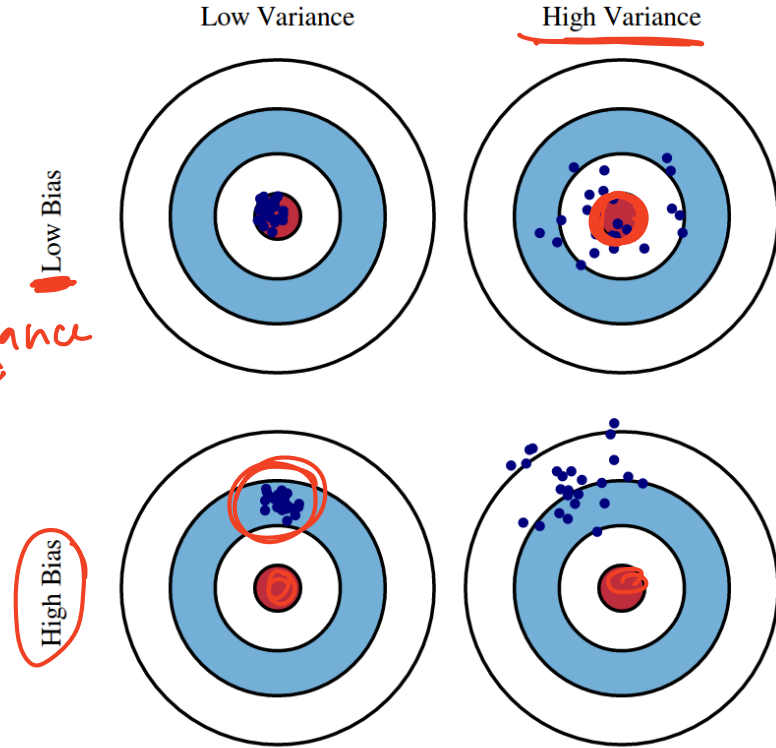
- First, let's think about the number of observations in our dataset
- For example: In the Gapminder dataset, I can use an indicator for each country
 - Remember that each country is an observation
 - So we have a perfectly fit model - a covariate for each observation
 - But we cannot generalize this to any other countries
 - And we haven't identified any meaningful relationships between life expectancy and other measured characteristics
- More covariates in the model is not always better
 - Overfitting the data limits our generalizability and prevents us from answering research questions

n countries
p = n covariates

Model Complexity vs. Parsimony

Suppose we have $p = 30$ covariates (in the true model) and $n = 50$ observations. We could consider the following two alternatives:

1. We could fit a model using all of the covariates. prediction 30
 - In this case, $\hat{\beta}$ is unbiased for β (in a linear model fit using OLS). But $\hat{\beta}$ has very high variance. \rightarrow high variance of \hat{y}
2. We could fit a model using only the five strongest covariates. association
 - In this case, $\hat{\beta}$ will be biased for β , but it will have lower variance (compared to the estimate including all covariates) $\hat{\beta}$ is interpretable



Source: <http://scott.fortmann-roe.com/docs/BiasVariance.html>

Bias-variance trade off

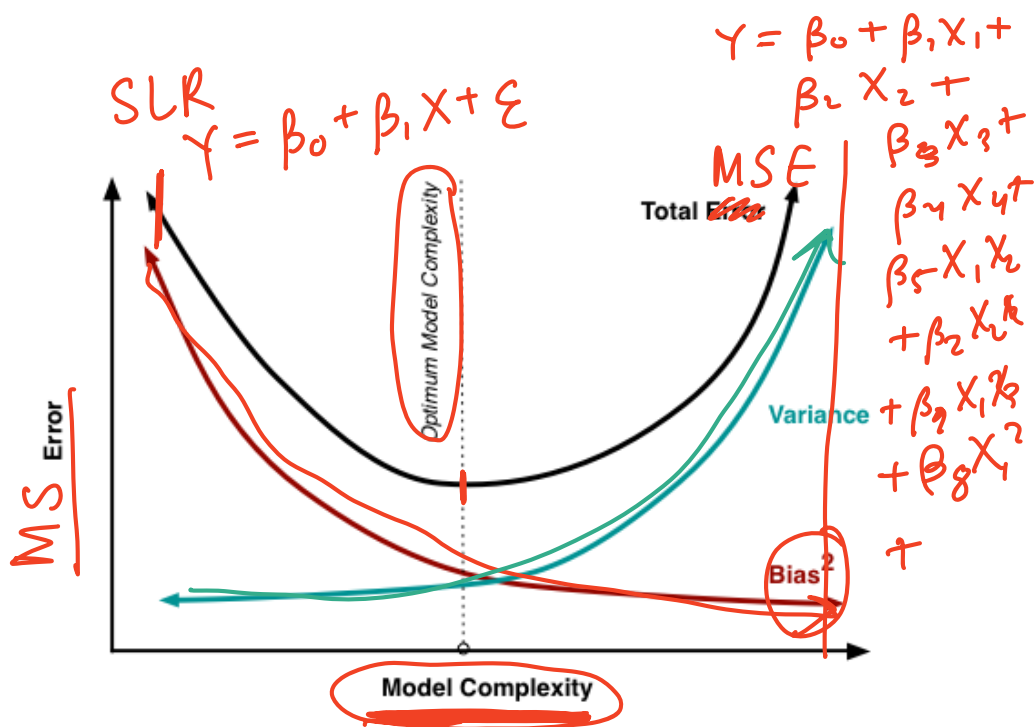
- Recall mean square error is a function of SSE (sum of squared residuals)

$$MSE = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

- MSE can also be written as a function of the bias and variance

$$MSE = \text{bias}(\hat{\beta})^2 + \text{variance}(\hat{\beta})$$

- For the same data:
 - More covariates in model: less bias, more variance
 - Less covariates in model: more bias, less variance
- Our goal: find a model with just the right amount of covariates to balance bias and variance



Source: <http://scott.fortmann-roe.com/docs/BiasVariance.html>

potentially NOT generalizable

inc # covariates
inc interactions

w/ new data (not fitted),
does our model still hold?

Model Selection basics (slide adjusted from Jodi Lapidus)

- “Because models always fall far short of the complex reality under study, there are no best or optimal strategies for modeling.”
 - From: **Statistical Foundations for Model-Based Adjustments**
- Not all statistical texts provide practical advice on model development
 - A lot of resources include methods/code to compare models, but does not include much advice re: selecting which model to ultimately use.
 - Other texts are sparse on details or incorporate simplistic approaches
- Model development strategy should **align with research goals**
 - Prediction vs. Estimating Association
 - Strategy may depend on study design and data set size

The goals of association vs. prediction

Association / Explanatory / One variable's effect

- **Goal:** Understand one variable's (or a group of variable's) effect on the response after adjusting for other factors
- Mainly interpret the coefficient of the variable that is the focus of the study
 - Interpreting the coefficients of the other variables is not important, but can help bring context
- Any variables not selected for the final model have still been adjusted for, since they had a chance to be in the model *Don't usually explicitly mention*
- Example: How is body mass of a penguin associated with flipper length?

Prediction

- **Goal:** to calculate the most precise prediction of the response variable
- Interpreting coefficients is not important
- Choose only the variables that are strong predictors of the response variable
 - Excluding irrelevant variables can help reduce widths of the prediction intervals
- Example: What is the flipper length of a penguin with body mass of 3000 g (and all its other characteristics)?

Model building for association vs. prediction

More information on the two analysis goals:

5/2

5/3

Table 1. Summary of explanatory versus predictive models

| | Explanatory Models | Predictive Models |
|--|--|---|
| Goal | Establish <u>causal relationships</u> but mostly associations | Predict <u>current diagnoses</u> or future outcomes |
| <u>Threats to validity</u> | Chance <u>findings</u> (type I errors); <u>confounding</u> | <u>Overfitting</u> ; lack of generalizability to new populations |
| Candidate variables | A <u>limited set of prespecified</u> risk factors and confounders | A <u>larger set</u> of potential predictors; some predictors may not be causally related to the outcome |
| Variable selection | <u>Hypothesis driven</u> ; should not use automated selection procedures → <i>interpreting</i> | Exploratory; may use automated selection procedures, but validation is essential and newer automated procedures that incorporate shrinkage are preferred |
| Measures of model performance <i>model fit statistics</i> | Size of β coefficients for individual risk factors; level of significance for individual risk factors → | Discrimination (eg, <u>ROC analysis</u>); calibration (eg, <u>Hosmer-Lemeshow test</u>); goodness of fit (eg, <u>R^2</u> , AIC); reclassification (eg, net reclassification index); clinical utility |
| <u>Validation</u> | New studies are needed to confirm individual causal relationships | Internal validation: split-sample validation; cross validation; bootstrap validation; external validation |

ROC = receiver operating characteristic; AIC = Akaike information criterion.

If you ever get the chance, check out Dr. Kristin Sainani's series on Statistics

Poll Everywhere Question 1

Which of the following is the most likely consequence when selecting a model for association? 12

- Too many variables in the model, higher bias and lower variance
- Too many variables in the model, lower bias and higher variance
- Too few variables in the model, lower bias and higher variance
- Too few variables in the model, higher bias and lower variance

Model selection strategies for continuous outcomes

Association / Explanatory / One variable's effect

- Selection of potential models is tied more with the research context with some incorporation of prediction scores
- Pre-specification of multivariable model ✓
- ★ Purposeful model selection ✓ —
 - “Risk factor modeling”
- Change in Estimate (CIE) approaches →
 - Will learn in Survival Analysis (BSTA 514)

Prediction

- Selection of potential models is fully dependent on prediction scores
- Automated strategies
 - Stepwise selection (forward/backward) —
 - You'll see these a lot, but they're not really good methods —
 - Best subset
 - Regularization techniques (LASSO, Ridge, Elastic net)]

- For categorical outcomes, there are more prediction model selection strategies (will learn more in BSTA 513)
 - Examples: Decision trees, Random forest, Neural networks, K-means]

Y/N

Learning Objectives

1. Understand the motivation for model selection, including bias-variance trade off and alignment of research goals (association vs. prediction)
2. Explain the general process or idea behind different model selection techniques
3. Recognize common model fit statistics and understand what they measure

Pre-specification of multivariable model (slide adjusted from Jodi Lapidus)

- In a clinical trial, we often have to write and finalize a statistical analysis plan (SAP) before the trial starts
- If we wish to compare treatment effects adjusted for covariates, all covariates typically specified in advance
 - Example: Comparing effectiveness of 3-drug vs. 2-drug regimen for delaying AIDS onset or death.
Covariates such as severity of HIV infection at baseline would have been specified in advance.
 - Variables such as study site, as well as any randomization stratification variables are common covariates.
- In these cases, only a limited number of multivariable models are fit and reported
 - Do not perform all the model building steps outlined in Hosmer and Lemeshow texts

Purposeful model selection (slide adjusted from Jodi Lapidus)

- Can use this type of model selection for any type of regression
- Careful, well-thought out variable selection process
 - Considers both confounding and interaction, as well as checking model assumptions, fit, etc.
- Often a reasonable strategy, especially in epidemiology and more exploratory clinical studies
 - However, not always appropriate!
 - E.g. clinical trials with model specified in advance. (pre-specified model)
- **This is the selection process that we will focus on in this class!**

Change in estimate (CIE) approach (slide adjusted from Jodi Lapidus)

- CIE strategies select covariates on the basis of how much their control changes exposure effect estimates
 - Observed change is presumed to measure confounding by the covariate. *→ explanatory*
- What estimate?
 - H/L text suggest using coefficients from the model
 - We typically use the coefficient estimate from the explanatory variable that we are most interested in
- What magnitude change is "important"?
 - H/L text suggest 10% *→ confounders*
- One must choose an effect measure to judge change importance, where "importance" needs to be evaluated along a contextually meaningful scale
- Accurate assessment of confounding may require examining changes from removing entire sets of covariates
 - Add in or eliminate candidate confounders one at time?
 - Add in or eliminate candidate confounders in sets?

Hosmer / Lemeshow 513

Stepwise selection (slide adjusted from Adrianna Westbrook)

- This is an incredibly common approach that statisticians use, often because it is an older and more recognized method
 - BUT IT IS ALSO ONE OF THE WORST MODEL SELECTION STRATEGIES!!
- Major disadvantages to stepwise selection:
 - Prone to overfitting ✓
 - Biased estimates ✓
 - Cements the wrong idea that we are looking for our "most significant" covariates
- Predictors/covariates are added or removed one at time if they are below a certain threshold (usually p-value below 0.10 to 0.20)

Stepwise selection: two common approaches

- I will introduce two of the approaches so that you understand the general process if a collaborator ever mentions stepwise selection
- Forward selection:
 - For p ⁰ covariates potential covariates, run all simple linear regressions:
 - $Y = \beta_0 + \beta_1 X_1 + \epsilon$ through $Y = \beta_0 + \beta_1 X_p + \epsilon$
 - Include the X_i with the lowest p-value (assuming it is below the threshold)
 - Now run $Y = \beta_0 + \beta_1 X_i + \beta_2 X_1 + \epsilon$ through $Y = \beta_0 + \beta_1 X_i + \beta_2 X_p + \epsilon$ and enter the next X_j with the lowest p-value
 - Continue process until no more predictors come back with a p-value below the threshold
- Backward selection:
 - Start with a full model ($Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \epsilon$) and remove predictor with the highest p-value (assuming it is above the threshold)
 - Repeatedly remove the variable with the highest p-value until all remaining variables meet the stopping criteria (are below the threshold)

p-value lowest < 0.2

p-val > 0.2 & highest

Best subset (slide adjusted from Adrianna Westbrook)

- I don't see this approach very often
- Quite literally making subsets of the data and using the “best” one
- General steps:
 - Run every possible model fitting 1 to all possible p predictors/covariates
 - You can limit number of potential predictors
 - 2^p = total number of models where p = number of predictors
 - You will get the best fitting model within each category (i.e., 1 predictor model, 2 predictor model,..., p predictor model)
 - Then have to find the best fitting model between the best models from each category
- Major disadvantages to best subset:
 - Does not account for interactions
 - Needs to run a lot of models (takes A LOT of time)

Regularization techniques

likelihood + penalty for more covariates

- Regularization techniques (LASSO, ridge, elastic net) adds a penalization that shrinks (or regularizes) coefficients down to reduce overfitting

| | <u>LASSO</u> (Least About Shrinkage and Selection Operator) | <u>Ridge</u> | <u>Elastic Net</u> |
|--------------|--|---|---|
| Penalization | <u>L-1 Norm</u> , uses absolute value | <u>L-2 Norm</u> , uses squared value | Best of both worlds, <u>L-1 and L-2 used</u> |
| Pro's | Reduces overfitting, will shrink coefficient to zero | Reduces overfitting, handles collinearity, can handle $k > n$ | Reduces overfitting, handles collinearity, handles $k > n$, shrinks coefficients to zero |
| Con's | Cannot <u>handle $k > n$</u> , doesn't handle multicollinearity well | Does not shrink coefficients to zero, difficult to interpret | More difficult for R to do than the other two (but not really that bad) |

Poll Everywhere Question 2

Learning Objectives

1. Understand the motivation for model selection, including bias-variance trade off and alignment of research goals (association vs. prediction)
2. Explain the general process or idea behind different model selection techniques
3. Recognize common model fit statistics and understand what they measure

Introduction to model fit statistics

- So far we have compared models using the F-test
- The F-test is a great way to compare models that are **nested**
 - Basically, this means that the “full” model contains all the covariates that the “reduced” model contains
 - The full model will have additional covariates, but the covariates in the reduced is a subset of the covariates in the full
- What if we want to compare models that are not nested?
 - There is a special group of fit statistics that can help us compare models
 - Note: these are sometimes used in the model building process (within one strategy)
 - Helpful if we want to compare selected models across strategies
 - Helpful if we have a few “final” models with different covariates that we want to compare

Common model fit statistics

- The following model fit statistics combine information about the SSE, the number of parameters in the model, and the sample size
- For these fit statistics, smaller values indicate better model fit!

| Fit statistic | Equation | R code |
|--------------------------------------|--|---|
| R-squared / Adjusted R-squared | $Adj. R^2 = 1 - \frac{SSE/(n-p-1)}{SSY/(n-1)}$ | <code>Within summary(model_name)</code> |
| Mallow's C_p | $C_p = \left[\frac{\hat{\sigma}_p^2}{\hat{\sigma}_{max}^2} - 1 \right] (n - p) + p$ | <code>ols_mallows_cp()</code> |
| Akaike information criterion (AIC) | $AIC = n \log(SSE) - n \log(n) + 2(p + 1)$ | <code>AIC(model_name)</code> |
| Bayesian information criterion (BIC) | $BIC = n \log(SSE) - n \log(n) + \log(n) \cdot (p + 1)$ | <code>BIC(model_name)</code> |

- We don't need to know the exact formulas for them!

Common model fit statistics

- There is no hypothesis testing for these fit statistics
 - Only helpful if you are comparing models
 - Works for nested and non-nested models
- Common to report all or some of them
- All of the fit statistics will not necessarily reach a consensus about the best fitting model
 - Each weigh SSE, number of parameters, and number of observations differently

| Time point(s) | Model | χ^2 (df) | AIC | Sample size adjusted BIC | CFI | TLI | RMSEA | | | |
|---------------|-----------------------------|-----------------------------|-----------|--------------------------|-----|-----|----------------------|-------------------------|------------------|------|
| | | | | | | | RMSEA [95% CI] | Prob. Close Fit (< .05) | Null Model RMSEA | SRMR |
| T1 | 1 factor | 304.56 (82), $p < .001$ | 33,700.01 | 33,782.35 | .94 | .92 | .069 [.061, .077] | .000 | .217 | .066 |
| | 2 correlated factors | 258.91 (80), $p < .001$ | 33,658.36 | 33,743.05 | .95 | .93 | .062 [.054, .071] | .008 | .217 | .080 |
| | Bifactor | 201.99 (76), $p < .001$ | 33,609.44 | 33,698.84 | .97 | .95 | .054 [.045, .063] | .234 | .238 | .044 |
| T2 | 1 factor | 201.66 (78), $p < .001$ | 29,622.57 | 29,702.88 | .96 | .94 | .055 [.046, .065] | .179 | .197 | .074 |
| | 2 correlated factors | 201.17 (80), $p < .001$ | 29,618.07 | 29,696.22 | .96 | .94 | .054 [.045, .063] | .239 | .197 | .054 |
| | Bifactor | 177.93 (74), $p < .001$ | 29,606.83 | 29,691.49 | .96 | .94 | .052 [.042, .062] | .365 | .216 | .049 |
| T1-T2 | Regression structural model | 746.23 (370), $p < .001$ | 60,432.23 | 60,655.73 | .96 | .95 | .042 [.038, .046] | .999 | .186 | .054 |
| T1-T2 | Trait structural model | 817.17 (378), $p < .001$ | 60,487.16 | 60,701.25 | .96 | .94 | .045 [.041, .049] | .974 | .186 | .061 |

https://www.researchgate.net/figure/Model-Fit-Statistics_tbl1_308844501

