



Workshop Aims

Modelling
Strategies

Modelling to
Explain

Multicollinearity

Modelling to
Predict

Stepwise
Regression

Recap

Quantitative Social Research II

Workshop 2: Selecting Explanatory Variables

Jose Pina-Sánchez



Workshop Aims

Workshop Aims

Modelling
Strategies

Modelling to
Explain

Multicollinearity

Modelling to
Predict

Stepwise
Regression

Recap

- Discuss the difference between predicting and explaining
- Introduce stepwise regression methods
- Understand the implications of multicollinearity
 - Learn how to detect and tackle this problem



Workshop Aims: Recap

Workshop Aims

Modelling
StrategiesModelling to
Explain

Multicollinearity

Modelling to
*Predict*Stepwise
Regression

Recap

- Assumptions in the linear regression model ($Y = \alpha + \beta_k X_k + e$):
 - Normality: $N \sim (0, Var(e))$
 - Homoskedasticity: $Var(e_i) = Var(e)$
 - Independence: $Cov(e_i, e_j) = 0$
 - No endogeneity: $Cov(X_i, e_i) = 0$
 - Perfectly measured variables
 - No missing data (other than missing at random)
 - No omitted relevant variables
 - **No multicollinearity**
 - Linearity



Modelling Strategies

Workshop Aims

Modelling
Strategies

Modelling to
Explain

Multicollinearity

Modelling to
Predict

Stepwise
Regression

Recap

- Modelling strategies are first determined by the type of response variable (Y , aka dependent or outcome variable) to be explored
 - Last term: continuous (normally distributed), binary
 - Much more out there: duration data, count data, mixed data, etc.



Modelling Strategies

Workshop Aims

Modelling
Strategies

Modelling to
Explain

Multicollinearity

Modelling to
Predict

Stepwise
Regression

Recap

- Modelling strategies are first determined by the type of response variable (Y , aka dependent or outcome variable) to be explored
 - Last term: continuous (normally distributed), binary
 - Much more out there: duration data, count data, mixed data, etc.
- It is also crucial to think carefully about the right-hand side of the equation
 - Which set of explanatory variables (X_k , aka regressors, covariates, independent variables) to include
 - Question: What considerations have you been following so far?

Modelling Strategies

Workshop Aims

Modelling
Strategies

Modelling to
Explain

Multicollinearity

Modelling to
Predict

Stepwise
Regression

Recap

- Modelling strategies are first determined by the type of response variable (Y , aka dependent or outcome variable) to be explored
 - Last term: continuous (normally distributed), binary
 - Much more out there: duration data, count data, mixed data, etc.
- It is also crucial to think carefully about the right-hand side of the equation
 - Which set of explanatory variables (X_k , aka regressors, covariates, independent variables) to include
 - Question: What considerations have you been following so far?
- This is the focus of the next three workshops
 - Today: predictive/inductive vs explanatory/deductive strategies, multicollinearity
 - W3: confounders, mediators, moderators, colliders
 - W4: polynomial regression, LOWESS curves



What's the Research Aim

Workshop Aims

Modelling
Strategies

Modelling to
Explain

Multicollinearity

Modelling to
Predict

Stepwise
Regression

Recap

- Depending on whether we seek to *predict* or to *explain* we will adopt different strategies
- We can often figure out which one it is from the research question
- Question: Are the following research questions aiming at predicting or explaining?
 - Can the onset of riots be identified using real time Tweets?
 - Are riots caused by economic inequality?
 - Is sentencing an art or a science? (Can we forecast judicial decisions?)



Workshop Aims

Modelling
Strategies

Modelling to
Explain

Multicollinearity

Modelling to
Predict

Stepwise
Regression

Recap

Competing Strategies

Predicting

- Inductive / exploratory
- Data driven

Explaining

- Deductive
- Theory driven



Workshop Aims

Modelling
Strategies

Modelling to
Explain

Multicollinearity

Modelling to
Predict

Stepwise
Regression

Recap

Competing Strategies

Predicting

- Inductive / exploratory
- Data driven
- X_k chosen to maximise predictability
- Not interested in interpretations of β_k

Explaining

- Deductive
- Theory driven
- X_k choice theoretically determined
- Interested in interpretations of β_k (causal explanations)



Workshop Aims

Modelling
StrategiesModelling to
Explain

Multicollinearity

Modelling to
*Predict*Stepwise
Regression

Recap

Competing Strategies

Predicting

- Inductive / exploratory
- Data driven
- X_k chosen to maximise predictability
- Not interested in interpretations of β_k
- Not so worried about violating assumptions
- Can employ unsupervised model selection

Explaining

- Deductive
- Theory driven
- X_k choice theoretically determined
- Interested in interpretations of β_k (causal explanations)
- Very worried about violating assumptions
- Model selection should be supervised



Good Practices in Variable Selection

- Last term you reviewed good practices for selecting explanatory variables
 - Let theory dictate the selection process
 - Aim to include only relevant variables (variables of interest but also potential confounders)
 - The principle of parsimony (simple is better)

Workshop Aims

Modelling
Strategies

Modelling to
Explain

Multicollinearity

Modelling to
Predict

Stepwise
Regression

Recap



Good Practices in Variable Selection

Workshop Aims

Modelling
Strategies

Modelling to
Explain

Multicollinearity

Modelling to
Predict

Stepwise
Regression

Recap

- Last term you reviewed good practices for selecting explanatory variables
 - Let theory dictate the selection process
 - Aim to include only relevant variables (variables of interest but also potential confounders)
 - The principle of parsimony (simple is better)
- This is to facilitate the interpretation of our results and to avoid some ‘unpleasant’ outcomes
 - Prevent P-hacking (1 in 20 coefficient estimates will be significant by chance even if they are just noise)
 - Prevent HARKing (hypothesising after results are known)

Good Practices in Variable Selection

Workshop Aims

Modelling
Strategies

Modelling to
Explain

Multicollinearity

Modelling to
Predict

Stepwise
Regression

Recap

- Last term you reviewed good practices for selecting explanatory variables
 - Let theory dictate the selection process
 - Aim to include only relevant variables (variables of interest but also potential confounders)
 - The principle of parsimony (simple is better)
- This is to facilitate the interpretation of our results and to avoid some ‘unpleasant’ outcomes
 - Prevent P-hacking (1 in 20 coefficient estimates will be significant by chance even if they are just noise)
 - Prevent HARKing (hypothesising after results are known)
 - Facilitate model estimation and minimise computational time
 - Models difficult to interpret

Good Practices in Variable Selection

Workshop Aims

Modelling
Strategies

Modelling to
Explain

Multicollinearity

Modelling to
Predict

Stepwise
Regression

Recap

- Last term you reviewed good practices for selecting explanatory variables
 - Let theory dictate the selection process
 - Aim to include only relevant variables (variables of interest but also potential confounders)
 - The principle of parsimony (simple is better)
- This is to facilitate the interpretation of our results and to avoid some ‘unpleasant’ outcomes
 - Prevent P-hacking (1 in 20 coefficient estimates will be significant by chance even if they are just noise)
 - Prevent HARKing (hypothesising after results are known)
 - Facilitate model estimation and minimise computational time
 - Models difficult to interpret
 - Overfitting (loss of degrees of freedom)
 - Multicollinearity



Multicollinearity

Workshop Aims

Modelling
Strategies

Modelling to
Explain

Multicollinearity

Modelling to
Predict

Stepwise
Regression

Recap

- Absence of severe multicollinearity is one of the assumptions we invoke when specifying regression models
- Can arise as result of using too many and/or too highly correlated explanatory variables
- The model cannot identify the variability on Y associated to each X_k
 - Regression coefficient estimates (β_k) are unstable (likely biased)
 - Their measures of uncertainty (anything derived from their SE_k) are overestimated \rightarrow false negatives (type-II errors) more likely

Which Model Is Affected by Multicollinearity?

Workshop Aims

Modelling
Strategies

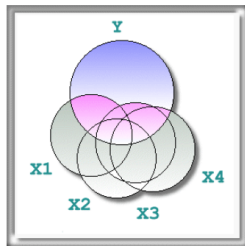
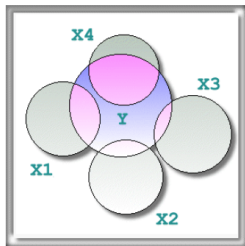
Modelling to
Explain

Multicollinearity

Modelling to
Predict

Stepwise
Regression

Recap



Source: Quantitative Methods for Linguistic Data

Which Model Is Affected by Multicollinearity?

Workshop Aims

Modelling Strategies

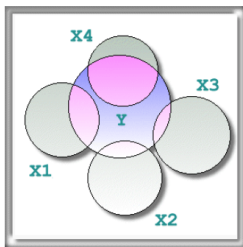
Modelling to Explain

Multicollinearity

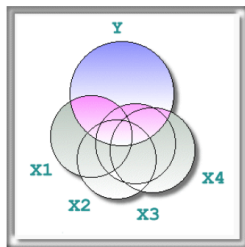
Modelling to Predict

Stepwise Regression

Recap



No collinearity



Substantial collinearity

Source: Quantitative Methods for Linguistic Data



Detecting Multicollinearity

- Most commonly detected by looking at a correlation matrix with your potential explanatory variables
 - The rule of thumb is to look out for correlations $> .8$
 - Yet, this diagnostic is based just on pairwise comparisons
 - Multicollinearity can also take place when variables are moderately correlated, but there are plenty of them

Workshop Aims

Modelling
Strategies

Modelling to
Explain

Multicollinearity

Modelling to
Predict

Stepwise
Regression

Recap



Detecting Multicollinearity

Workshop Aims

Modelling
Strategies

Modelling to
Explain

Multicollinearity

Modelling to
Predict

Stepwise
Regression

Recap

- Most commonly detected by looking at a correlation matrix with your potential explanatory variables
 - The rule of thumb is to look out for correlations $> .8$
 - Yet, this diagnostic is based just on pairwise comparisons
 - Multicollinearity can also take place when variables are moderately correlated, but there are plenty of them
- Better to rely on the Variance Inflation Factor (VIF)
 - The VIF estimates how much the variance of a regression coefficient is inflated due to multicollinearity in the model
 - $VIF_k = \frac{1}{1 - R_k^2}$, where the R_k^2 is obtained by taking a predictor (k) and regressing it against every other predictor in the model
 - Interpretation: the factor by which the variance of a coefficient (SE_K^2) is inflated compared to what it would be if there was no correlation with other predictors
 - Rule of thumb: if $VIF_k > 5$ then k is considered problematic



Tackling Multicollinearity

Workshop Aims

Modelling
Strategies

Modelling to
Explain

Multicollinearity

Modelling to
Predict

Stepwise
Regression

Recap

- Question: How do you deal with problems of multicollinearity?



Tackling Multicollinearity

Workshop Aims

Modelling
Strategies

Modelling to
Explain

Multicollinearity

Modelling to
Predict

Stepwise
Regression

Recap

- Question: How do you deal with problems of multicollinearity?
- Drop variables with a $VIF > 5$
 - This can lead to arbitrary choices
 - Difficult to justify when the correlated variables are theoretically important
- Aggregate variables into an index/scale
 - A possibility if various variables are tapping on the same latent concept
 - We can include a new single variable (the index) in the model, and remove all other variables used to create it (the items)
 - E.g. in exploring the presence of labour discrimination we can simply use a scale of social class, rather than employment status, level of education, salary, etc.
 - You saw how to create indexes using averages last term, in W6 we will learn how to use latent variable estimation



Modelling to *Predict*

Workshop Aims

Modelling
Strategies

Modelling to
Explain

Multicollinearity

Modelling to
Predict

Stepwise
Regression

Recap

- We do not care about the quality of the regression coefficients estimates since we do not need to interpret them
- All we care about is the accuracy with which the model predicts the outcome variable
- We should use as many useful predictors as possible
- Still, we need to be careful not to include noise



Model Selection

Workshop Aims

Modelling
Strategies

Modelling to
Explain

Multicollinearity

Modelling to
Predict

Stepwise
Regression

Recap

- Question: How can we determine whether the predictors we introduce are useful?



Model Selection

Workshop Aims

Modelling
Strategies

Modelling to
Explain

Multicollinearity

Modelling to
Predict

Stepwise
Regression

Recap

- Question: How can we determine whether the predictors we introduce are useful?
 - We have considered p-values and the R^2 (or the predictive accuracy of a logistic model)
 - By definition, the more variables included in the model, the higher its R^2 (or the predictive accuracy of a logistic model)...
 - but that is only when we make predictions based on the sample used to build the model



Model Selection

Workshop Aims

Modelling
Strategies

Modelling to
Explain

Multicollinearity

Modelling to
Predict

Stepwise
Regression

Recap

- Question: How can we determine whether the predictors we introduce are useful?
 - We have considered p-values and the R^2 (or the predictive accuracy of a logistic model)
 - By definition, the more variables included in the model, the higher its R^2 (or the predictive accuracy of a logistic model)...
 - but that is only when we make predictions based on the sample used to build the model
 - Including noisy predictors can reduce predictive accuracy
 - We can see that by using two samples of the same population, or by splitting our sample into a *train* and a *test* sample



Model Selection

Workshop Aims

Modelling
Strategies

Modelling to
Explain

Multicollinearity

Modelling to
Predict

Stepwise
Regression

Recap

- Question: How can we determine whether the predictors we introduce are useful?
 - We have considered p-values and the R^2 (or the predictive accuracy of a logistic model)
 - By definition, the more variables included in the model, the higher its R^2 (or the predictive accuracy of a logistic model)...
 - but that is only when we make predictions based on the sample used to build the model
 - Including noisy predictors can reduce predictive accuracy
 - We can see that by using two samples of the same population, or by splitting our sample into a *train* and a *test* sample
 - We should undertake variable selection based on criteria that penalises adding variables, such as AIC,...
 - rather than on the statistical significance of individual predictors



Stepwise Regression

Workshop Aims

Modelling
Strategies

Modelling to
Explain

Multicollinearity

Modelling to
Predict

**Stepwise
Regression**

Recap

- ok, so we use a *train* and a *test* sample, AIC to select useful predictors...
- but how do we undertake the model comparison process that will give us the best model?



Stepwise Regression

Workshop Aims

Modelling
Strategies

Modelling to
Explain

Multicollinearity

Modelling to
Predict

Stepwise
Regression

Recap

- ok, so we use a *train* and a *test* sample, AIC to select useful predictors...
- but how do we undertake the model comparison process that will give us the best model?
 - Do we go one by one until adding new predictors does not improve the AIC? (forward selection)
 - Do we throw them all in the model and proceed to remove them sequentially until the AIC stops improving? (backward selection)
 - If the list of predictors is long this could take us a long time, that's why the above strategies are normally unsupervised



Stepwise Regression

Workshop Aims

Modelling
Strategies

Modelling to
Explain

Multicollinearity

Modelling to
Predict

Stepwise
Regression

Recap

- ok, so we use a *train* and a *test* sample, AIC to select useful predictors...
- but how do we undertake the model comparison process that will give us the best model?
 - Do we go one by one until adding new predictors does not improve the AIC? (forward selection)
 - Do we throw them all in the model and proceed to remove them sequentially until the AIC stops improving? (backward selection)
 - If the list of predictors is long this could take us a long time, that's why the above strategies are normally unsupervised
 - Also, how do we choose which variables go in/out first?
 - Predictors can become much more or less important depending on what other variables are already in the model
 - E.g. years of experience might appear less important in predicting salary if workers' age is already in the model, and vice versa



Stepwise Selection

Workshop Aims

Modelling
Strategies

Modelling to
Explain

Multicollinearity

Modelling to
Predict

Stepwise
Regression

Recap

- Stepwise selection consists on iteratively adding and removing predictors, a combination of forward and backward selection
- There are three procedures involved in the algorithm
 - Starts with no predictors, then sequentially add the most consequential variables (forward selection)
 - After adding each new variable, remove any variables that no longer provide an improvement in the model fit (backward selection)
 - Until the model cannot be improved by adding or removing variables
- Model selection is a key area in machine learning, with new methods being developed at speed, e.g.
 - Random forests
 - Bayesian model averaging

Workshop Aims

Modelling
StrategiesModelling to
Explain

Multicollinearity

Modelling to
*Predict*Stepwise
Regression

Recap

- Think about what are you trying to accomplish through your research: *explain* or *predict*
 - The first step in designing your variable selection strategy
- If we seek to explain then parsimony is key
 - Avoiding problems of multicollinearity and overfitted models
 - Next week we will learn the importance of distinguishing between confounder, mediator and collider effects
- If we seek to predict we will include as many useful predictors as we can gather
 - The model selection can be unsupervised
 - Using methods such as stepwise regression
- To learn more about stepwise regression you can read:
Ruczinsky *Variable selection*