

# Introduction; problems with classical methods

Patrick Breheny

January 21, 2025

# Introduction

- This course concerns the analysis of data in which we are attempting to predict an outcome  $Y$  using a number of explanatory factors  $X_1, X_2, X_3, \dots$ , some of which may not be particularly useful
- Although the methods we will discuss can be used solely for prediction (i.e., as a “black box”), I will adopt the perspective that we would like the statistical methods to be interpretable and to explain something about the relationship between the  $X$  and  $Y$
- Regression models are an attractive framework for approaching problems of this type, and the majority of the course will focus on extending classical regression modeling to deal with high-dimensional data

# High-dimensional data

Modern computation has changed the way science is conducted, and enabled researchers to easily collect, store, and access data for large numbers of features (ballpark number of features in parentheses):

- Advances in information technology such as REDCap ( $\sim 100$ )
- Adoption of electronic medical records ( $> 100$ )
- Molecular biology technologies such as microarrays and RNA-Seq ( $> 10,000$ )
- Advances in genotyping and genetic sequencing ( $> 100,000$ )

## High-dimensional data (cont'd)

- This type of data is known as *high dimensional data*
- Throughout the course, we will let
  - $n$  denote the number of independent sampling units (e.g., number of patients)
  - $p$  denote the number of features recorded for each unit
- In high-dimensional data,  $p$  is large with respect to  $n$ 
  - This certainly includes the case where  $p > n$
  - However, the ideas we discuss in this course are also relevant to many situations in which  $p < n$ ; for example, if  $n = 100$  and  $p = 80$ , we probably don't want to use ordinary least squares
- Note that “high dimensional” is not the same as “big data”:
  - $n = 50, p = 60$ : High dimensional, but not big
  - $n = 50,000,000, p = 60$ : Big, but not high dimensional

## More notation

- We will use  $\mathbf{X}$  to denote the  $n \times p$  matrix containing the predictor variables, with element  $x_{ij}$  recording the value of the  $j$ th feature for the  $i$ th independent unit
- We will let  $\mathbf{y}$  denote the length- $n$  vector of response values
- For the sake of simplicity, for most of the course we will assume that  $Y$  is normally distributed, but we will consider other types of responses in the “Other likelihoods” topic

# Univariate analysis

- A simple, widely used approach to analyzing high-dimensional data is to split the problem up into a large number of low-dimensional problems
- For example, rather than trying to regress  $y$  simultaneously on all the features, we can carry out  $p$  separate single-variable regressions, one for each feature:

$$y_i = \alpha_j + \beta_j x_{ij} + \epsilon_{ij}$$
$$\epsilon_{ij} \stackrel{\perp\!\!\!\perp}{\sim} N(0, \sigma_j^2);$$

this approach is also known as *marginal regression*

# Univariate analysis: Challenges

- The appeal of this approach is that classical regression can be easily applied to the separate analyses to yield estimates  $\{\hat{\beta}_j\}$ , confidence intervals, and test hypotheses to produce  $p$ -values  $\{p_j\}$
- The major complication, however, is that this approach involves a large number of separate analyses that must somehow be combined into a single set of results
- Thus, while standard methods can be used for the initial analyses, there has been a great deal of innovation over the past 30 years in terms of how to combine these results; we will discuss these innovations during the “Large scale testing” topic

## Limitations of univariate models

Marginal regression is straightforward, but has several drawbacks:

- Fails to account for correlation among the features
- Provides no way to estimate the independent effect of a feature while other features remain unchanged
- Diminished power
- No good way to combine the predictions of separate regressions into a single overall prediction
- No way of assessing the overall proportion of the variability in the outcome that may be explained by the features



## Joint modeling

- These issues can only be resolved by considering a joint model of the relationship between  $\mathbf{y}$  and the full set of features:

$$y_i = \beta_0 + \sum_{j=1}^p \beta_j x_{ij} + \epsilon_i$$
$$\epsilon_i \stackrel{\text{i.i.d.}}{\sim} \text{N}(0, \sigma^2)$$

- The maximum likelihood approach involves solving for the value of  $\beta$ , known as the maximum likelihood estimator (MLE), that minimizes the residual sum of squares  $\|\mathbf{y} - \mathbf{X}\beta\|_2^2$
- Here,  $\|\mathbf{v}\|_2 = \sqrt{\sum_i v_i^2}$  denotes the Euclidean norm; we will use this notation frequently throughout the course

# OLS

- The solution is determined by the linear system of equations

$$\mathbf{X}^T \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{X}^T \mathbf{y}$$

- Provided that  $\mathbf{X}^T \mathbf{X}$  is invertible, the system has the unique solution

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y},$$

known as the *ordinary least squares* (OLS) estimate

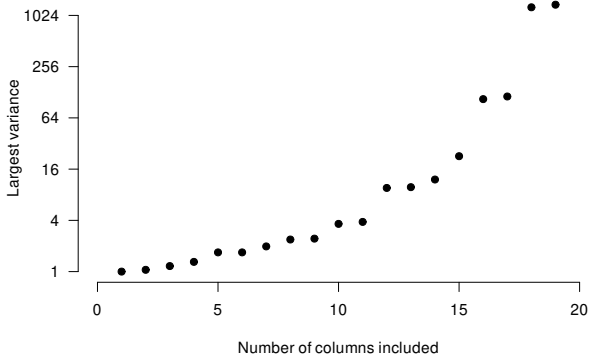
- The OLS estimate resolves all of the issues on slide 8 and has many well-recognized benefits such as yielding best linear unbiased estimates of  $\boldsymbol{\beta}$

## MLE problems

- However, there are many drawbacks to the use of maximum likelihood for estimating  $\beta$  when  $p$  is large
- Most dramatically, when  $p \geq n$  the matrix  $\mathbf{X}^T \mathbf{X}$  is not invertible and the MLE is not unique
- However, even if  $\mathbf{X}^T \mathbf{X}$  can be inverted and a unique maximum identified, as  $p$  increases and  $\mathbf{X}^T \mathbf{X}$  approaches singularity, the likelihood surface becomes very flat
- This means that a wide range of values of  $\beta$  are consistent with the data and wide confidence intervals required to achieve, say, 95% coverage

# An example

Consider a matrix  $\mathbf{X}$  with  $n = 20$  and whose elements consist of independent, normally distributed random numbers; the figure below plots the largest variance of the  $\hat{\beta}_j$  estimates as we increase the number of columns in  $\mathbf{X}$ :



## Remarks

- As  $p \rightarrow n$ ,  $\mathbb{V}(\hat{\beta})$  increases without bound; the increase is substantial as  $p$  approaches  $n$ , and infinite when  $p \geq n$
- Clearly, maximum likelihood cannot accommodate high-dimensional data without running into serious problems of identifiability and inefficiency

## The oracle model

- Suppose, however, that many features are unrelated to the outcome (in the sense that  $\beta_j = 0$ ), and only a few features are important
- If we knew in advance which elements of  $\beta$  are zero and which are not, then we could modify maximum likelihood without abandoning it completely, and avoid all of the earlier problems
- Specifically, we could apply maximum likelihood only to the variables for which  $\beta_j \neq 0$ ; this is known as the *oracle* model

# Model selection

- Obviously, the oracle model is a theoretical gold standard, not a realistic approach to data analysis, as it would require access to an oracle that could tell you which features are related to the outcome and which are not
- In the real world, we have to use the data in order to make empirical decisions about which features are related to the outcome and which are not; this is known as *model selection*

# Prediction

- Most approaches to model selection are based on prediction: if model A predicts future observations better than model B, then we should prefer model A to model B
- However, evaluating prediction error is not as straightforward as it may seem, and there are many competing approaches
- In particular, it is dangerous to use the same data for two purposes – to select the model and also to carry out inference with respect to the model's parameters – as we will see, this can introduce substantial bias



# Residual sum of squares

- By fitting the general regression model

$$y_i = f(\mathbf{x}_i) + \varepsilon_i, \quad i = 1, \dots, n,$$

we obtain  $\hat{y}_i = \hat{f}(\mathbf{x}_i)$ , the model's prediction for  $y_i$

- For linear regression,  $\hat{f}(\mathbf{x}_i) = \mathbf{x}_i^\top \hat{\boldsymbol{\beta}}$
- It is misleading to evaluate predictive accuracy by comparing  $\hat{y}_i$  to  $y_i$ , however, as  $y_i$  has already been used to calculate  $\hat{y}_i$  (i.e., it is not a genuine prediction):

$$\text{RSS} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

## Prediction error

- Because this underestimates the prediction error (PE, also known as the *test error*), we must examine how well  $\hat{f}(\mathbf{x})$  predicts *new* observations
- However, there are two ways of defining what exactly we mean by “new” observations:
  - $PE_X$ : Fit using  $(\mathbf{X}, \mathbf{y})$ , predict  $(\mathbf{X}, \mathbf{y}^{\text{new}})$
  - PE: Fit using  $(\mathbf{X}, \mathbf{y})$ , predict  $(\mathbf{X}^{\text{new}}, \mathbf{y}^{\text{new}})$
- Information criteria approaches tend to focus on  $PE_X$ , whereas cross-validation approaches focus on PE
- Both measures are reasonable; in this course, however, we'll use the second one, as it makes more sense in the high-dimensional setting

## Estimation error

- Throughout the course, I will often show simulation results to illustrate how well various methods perform
- In simulations, we will often calculate estimation error (doing so with real data is impossible, of course, since the true parameter values  $\beta^*$  are unknown)
- Estimation error is usually summarized using the mean squared error (MSE)

$$\text{MSE} = \mathbb{E} \|\hat{\beta} - \beta^*\|_2^2$$

or root mean squared error  $\text{RMSE} = \sqrt{\text{MSE}}$

- “Mean” here refers to averaging over *data sets*, not averaging over features
- The name MSE is unfortunate and leads to confusion, as it doesn't explicitly state that we are looking at estimation error

## Model error

- Lastly, we will sometimes refer to the model error (ME), which is the error at level of the linear predictors:

$$\text{ME} = \mathbb{E}\{f(\mathbf{x}) - \hat{f}(\mathbf{x})\}^2$$

- For (homoskedastic) linear regression models,  $\text{PE} = \text{ME} + \sigma^2$
- However, the distinction is useful when considering consistency, since ideally ME will go to zero for a model as we collect more data, but PE will never go to zero due to inherent variability in the outcome

# The model selection problem

- Let's return now to what I call the “model selection problem”: what's wrong with using the data to both select a model and then fit that model?
- To illustrate, consider the following simulation:

$$x_{ij} \stackrel{iid}{\sim} \text{Unif}(0, 1) \quad \text{for } j \text{ in } 1, 2, \dots, 100$$
$$y_i \stackrel{iid}{\sim} \text{N}(0, 1)$$

for  $i$  in  $1, 2, \dots, 25$

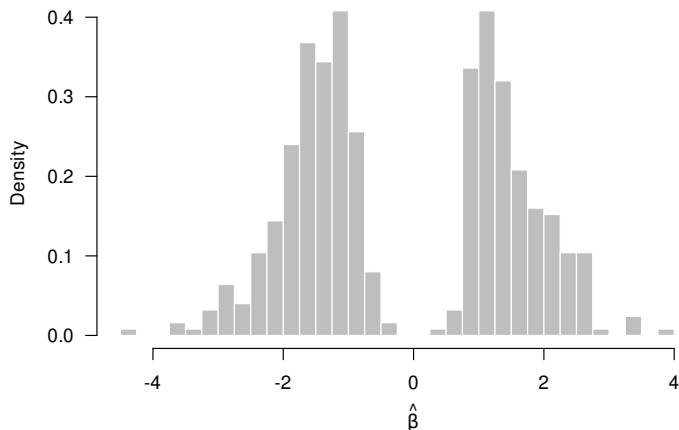
- We will use BIC to select the 5 most important variables, then use OLS with only those variables, and repeat this 100 times

## Remarks

- As we will see, this approach performs terribly
- By using the data set for model selection as well as estimation and inference, we have grossly distorted the sampling distribution of  $\hat{\beta}$
- This has dramatic consequences in terms of estimation, prediction, variable selection, and the validity of inference

# Results

A histogram of the 500  $\hat{\beta}_j$  estimates we obtain:



# Estimation

- The model selection process heavily biases the estimates of the regression coefficients away from zero
- In our simulation, most estimates were approximately  $\pm 1.5$  instead of being close to 0, the true value
- In particular, the average MSE is 2.7, compared to 0.52 for marginal regression, roughly a 5-fold increase
- This phenomenon is sometimes referred to as the “winner’s curse”



## Variable selection

- Here, we imposed an upper bound of 5 on the number of variables we allowed to be selected by the BIC-guided forward selection process; in all 100 replications, this upper bound was reached
- Obviously, since the true model in this case is the null (intercept-only) model, the model selection process we have employed here results in systematic overfitting
- While it is true that asymptotically, BIC will select the true model with probability tending to 1, that asymptotic argument relies on  $p$  remaining fixed while  $n \rightarrow \infty$ , or in other words, on  $n \gg p$
- Clearly, BIC cannot be relied on for accurate variable selection in high-dimensional problems

# Prediction

- On average, the selected models achieved a mean squared prediction error of 2.15, compared to a prediction error of  $\sigma^2 = 1$  for the null model
- Thus, by carrying out model selection, we have reduced the predictive accuracy of the model by half (doubled its error)

# Inference

- Finally, let us consider the validity of the inferences that we obtain from the post-selection OLS model:
  - The median  $p$ -value for testing  $H_0 : \beta_j = 0$  was  $p = 0.0013$
  - The actual coverage achieved by constructing 95% confidence intervals was under 5%
- Ignoring selection effects when carrying out post-selection inference produces conclusions that are far too liberal, with actual errors accumulating at a much higher rate than the statistical inferential approaches would indicate
- In summary, this approach is wildly optimistic and overconfident

## Final remarks

- These problems are widely recognized; unfortunately, they are also widely ignored
- The problem of developing statistical methods capable of simultaneous variable selection and inference has challenged statisticians for decades, from Scheffé (1953) to the present
- One of the primary goals of this course is to demonstrate the extent to which recent developments in penalized regression address and alleviate the concerns about simultaneous selection and inference we have raised today