

Data Science and Big Data Analytics: Making Data-Driven Decisions

Module 2 Regression and Prediction

Lecture Notes

Overview

The tools and techniques developed in the previous module allowed us to take a peek at large and unstructured data. Once we get a better handle on what this data looks like and what information it contains, a natural question is “**Can we use it to make predictions?**”. The need to formulate predictions from data can arise in a variety of contexts, whether we are trying to predict future data or simply data that is hard or expensive to collect in the future but was available in the past. The general idea behind prediction is to analyze data to identify an input/output function. “Module 2: Regression and Prediction” covers many techniques for this task. Specifically, we will discuss:

1. Linear and nonlinear regression together with their extensions, including the important case of logistic regression for binary classification and causal inference where the goal is to understand the effects of actively manipulating a variable as opposed to passively measuring it.
2. Modern regression with high-dimensional data or finding a needle in a haystack: for large datasets, it becomes necessary to sort out which variables are relevant for prediction and which are not. Recent years have witnessed the development of new statistical techniques, such as Lasso or Random Forests, that are computationally scalable to large datasets and that automatically select relevant data.
3. Regression and causal inference to explain why “correlation does not imply causation” and how we can overcome this intrinsic limitation of regression by resorting to randomized control studies or controlling for confounding.

This module is taught by Victor Chernozhukov from the MIT Economics Department and MIT Institute for Data Systems and Society. Professor Chernozhukov’s expertise covers econometrics, machine learning, and statistics in which he has developed fundamental methodology, particularly to deal with high dimensional data. This module not only covers a description of the methodology involved but also a mathematical introduction to its underlying principles.

Goals

1. Understanding of basic regression for prediction and inferential purposes.
2. Understanding of modern linear and nonlinear regression for prediction and inferential purposes.

3. Getting practical experience of using classical and modern regression methods for prediction and inferential purposes.

Objectives

Students should be able to:

1. Recognize a regression problem
2. Propose an adequate model for inference purposes
3. Recognize the difference between correlation and causation
4. Recognize high-dimensional setups
5. Select appropriate method for high-dimensional learning
6. Understand the limitations of observational studies for causal inference with particular focus on confounding
7. Recognize and implement randomized controlled studies and perform causal inference in these

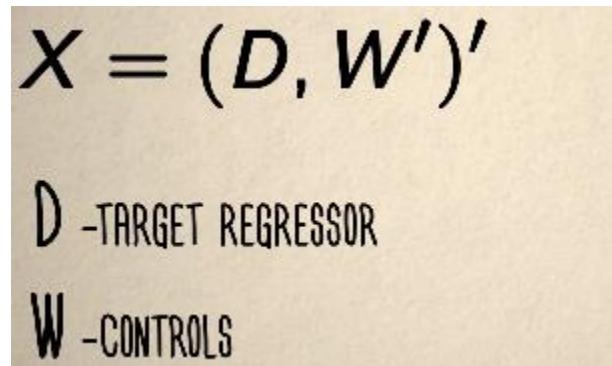
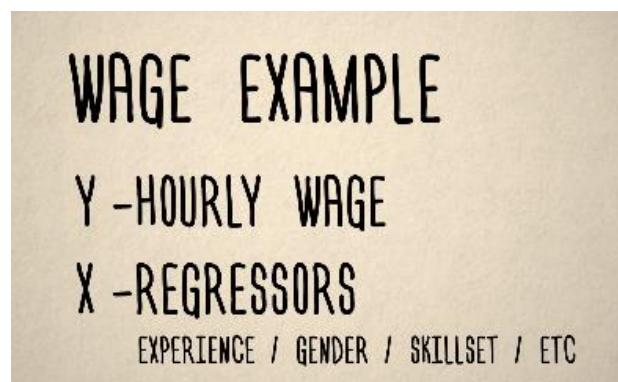
2.1.1 Introduction

Regression analysis is about discovering correlations between the outcome, y and the set of regressors x , which sometimes are also called features.

Outcome Y : real random variable

Regressors X or features: vector of random variables: $X = (X_1, X_2, \dots, X_p)$, assume $X_1 = 1$

Example:



The purpose of the regression analysis is to characterize the statistical relation y to x .

Two questions:

$Y \rightarrow X$

PREDICTION

-HOW CAN WE USE X TO PREDICT Y WELL?

INFERENCE

-HOW DOES THE PREDICTED VALUE OF Y CHANGE
IF WE CHANGE A COMPONENT OF X HOLDING
THE REST OF THE COMPONENTS OF X FIXED?

Example:

PREDICTION

CONSTRUCT A PREDICTION
RULE FOR HOURLY WAGE

INFERENCE

FIND GENDER PAY
DISCRIMINATION

2.1.2 Linear Regression for Prediction

LINEAR PREDICTION RULE FOR Y

$$X = (X_j)_{j=1}^p$$

$$\beta' X := \sum_{j=1}^p \beta_j X_j$$

$$\text{for } \beta = (\beta_j)_{j=1}^p$$

β : regression parameter

We minimize the expected or mean squared error that results from predicting Y , using linear rule given by $b'X$, Where b denotes a potential value for the parameter vector β .

BEST LINEAR PREDICTION BLP

$$\min_{b \in \mathbb{R}^p} E(Y - b'X)^2$$

$$b'X = \sum_{j=1}^p b_j X_j$$

$$b = (b_j)_{j=1}^p$$

The solution, $\beta'X$, is called the Best Linear Predictor of Y using X .

$\beta'X$

BEST LINEAR PREDICTOR
OF Y USING X

In words, this solution is the best linear prediction rule among all the possible linear prediction rules.

Next, we can compute an optimal beta by solving the first order conditions of the BLP problem, called the Normal Equations, where we have expected value of $(Y - \beta'X) = 0$

FIRST ORDER CONDITIONS OF THE BLP PROBLEM AKA NORMAL EQUATIONS

$$E(Y - \beta'X)X = 0$$

Minimize the objective function by setting the derivative to 0.

$$E(Y - \beta'X)X = 0$$

OBJECTIVE FUNCTION

$$b \mapsto E(Y - b'X)^2$$

Regression error: ϵ is uncorrelated to X .

REGRESSION ERROR

$$\epsilon := (Y - \beta'X)$$

$$EX\epsilon = 0 \quad Y = \beta'X + \epsilon$$

$$\beta'X$$

IS THE PREDICTED/EXPLAINED PART OF Y

$$\epsilon$$

IS THE RESIDUAL/UNEXPLAINED PART

In practice, we don't have access to the population data. Instead we observe a sample of observations of size n ,

FOR n OBSERVATIONS

$$(Y_i, X_i)_{i=1}^n = \{(Y_1, X_1), \dots, (Y_n, X_n)\}$$

DRAWN RANDOMLY FROM F

We assume that observations are generated as a random sample, from a distribution F , which is the population distribution of R , Y , and X . Formally this means that the observations are obtained as realizations of independently and identically distributed copies of the random vector Y, X .

We next construct the best linear prediction rule in sample for Y using X .

GIVEN X
PREDICTED Y

$$\sum_{j=1}^p \hat{\beta}_j X_j = \hat{\beta}' X$$
$$\hat{\beta} = (\hat{\beta}_j)_{j=1}^p$$

$\hat{\beta}$: sample regression coefficients

$$\sum_{j=1}^p \hat{\beta}_j X_j = \hat{\beta}' X$$
$$\hat{\beta} = (\hat{\beta}_j)_{j=1}^p$$

SAMPLE REGRESSION COEFFICIENTS

NEXT, DEFINE THE LINEAR REGRESSION IN THE SAMPLE
BY ANALOGY TO THE POPULATION PROBLEM

NEXT, DEFINE THE LINEAR REGRESSION IN THE SAMPLE
BY ANALOGY TO THE POPULATION PROBLEM

$\widehat{\beta}$ IS ANY SOLUTION TO THE BLP OR
ORDINARY LEAST SQUARES PROBLEM IN THE SAMPLE

$$\min_{b \in \mathbb{R}^p} \mathbb{E}_n (Y_i - b' X_i)^2$$

OPTIMAL $\widehat{\beta}$
SAMPLE NORMAL EQUATIONS

$$\mathbb{E}_n X_i (Y_i - \widehat{\beta}' X_i) = 0$$

FIRST ORDER CONDITIONS FOR THE BLP

IN-SAMPLE REGRESSION ERROR

$$\widehat{\epsilon}_i := (Y_i - \widehat{\beta}' X_i)$$

FIRST ORDER CONDITIONS FOR THE SAMPLE BLP

DECOMPOSITION

$$Y_i = \widehat{\beta}' X_i + \widehat{\epsilon}_i$$

RESIDUAL/UNEXPLAINED PART OF Y_i

THE BEST LINEAR PREDICTION RULE (OUT-OF-SAMPLE) IS $\beta'X$.
DOES $\hat{\beta}'X$ APPROXIMATE $\beta'X$?

WE ARE TRYING TO ESTIMATE P PARAMETERS β_1, \dots, β_p ,
WITHOUT IMPOSING ANY ASSUMPTIONS ON THESE PARAMETERS.

INTUITIVELY, TO ESTIMATE EACH PARAMETER WELL, WE NEED MANY
OBSERVATIONS PER PARAMETER.

THIS MEANS THAT n/p MUST BE LARGE, OR,
EQUIVALENTLY p/n MUST BE SMALL.

Theorem

Under regularity conditions

$$\sqrt{E_X(\beta'X - \hat{\beta}'X)^2} \leq const \cdot \sqrt{E\epsilon^2} \sqrt{\frac{p}{n}},$$

where E_X is expectation with respect to X and the inequality holds
with probability approaching 1 as $n \rightarrow \infty$.

If n is large and p is much smaller than n , for nearly all realizations
of the sample, the sample linear regression comes close to the
population linear regression.

WE DEFINE LINEAR REGRESSION IN THE POPULATION AND IN THE
SAMPLE THROUGH THE BEST LINEAR PREDICTION PROBLEMS SOLVED IN
THE POPULATION AND IN THE SAMPLE.

THE SAMPLE LINEAR REGRESSION (BEST LINEAR PREDICTOR)
APPROXIMATES THE POPULATION LINEAR REGRESSION (BEST LINEAR
PREDICTOR) VERY WELL WHEN THE RATIO p/n IS SMALL.

NEXT, WE WILL DISCUSS THE ASSESSMENT OF PREDICTION
PERFORMANCE IN PRACTICE.

2.1.3 Assessment of Prediction Quality

In this segment, we have two learning goals.

Goal 1 is to understand the analysis of variance, or ANOVA, in the population and in the sample.

Goal 2 is to assess the out of sample predictive performance of the sample linear regression.

We begin with the population case. Using the decomposition of Y into the explain and the residual part. And the normal equations that we derived in the previous segment, we can decompose variation of Y into the sum of explain variation and residual variation, as shown in the formula.

LINEAR REGRESSION IN POPULATION: ANALYSIS OF VARIANCE

USING THE DECOMPOSITION OF Y

$$Y = \beta'X + \epsilon$$

NORMAL EQUATIONS $E\mathbf{X}\epsilon = 0$

$$EY^2 = E(\beta'X)^2 + E\epsilon^2$$

We next define the population mean squared prediction error, or MSE, as the variance of the population residual. The expectation of epsilon squared. We also define the population R squared as the ratio of explained variation to the total variation. In other words, the population R squared is the variation of Y explained by the BLP. And as such, it is bounded below by 0 and above by 1.

DEFINE THE POPULATION MEAN SQUARED PREDICTION ERROR

$$\text{MSE}_{pop} = E\epsilon^2$$

$$R_{pop}^2 := \frac{E(\beta'X)^2}{EY^2} = 1 - \frac{E\epsilon^2}{EY^2} \in [0, 1]$$

The analysis of variants in the sample process analogous. We simply replace population expectations by the empirical expectations. Using the decomposition of

Y_i to explain and the residual part, and the normal equations for the sample. We can decompose the sample variation of Y_i into the sum of explained variation and residual variation. The former is given by the sample variation of this sample best in a predictor. And the latter is given by the sample variation of the residual.

$$Y = \beta' X + \epsilon$$

$$Y_i = \hat{\beta}' X_i + \hat{\epsilon}_i$$

$$Y_i = \hat{\beta}' X_i + \epsilon_i$$

$$\mathbb{E}_n X_i \hat{\epsilon}_i = 0$$

$$\mathbb{E}_n Y_i^2 = \mathbb{E}_n (\hat{\beta}' X_i)^2 + \mathbb{E}_n \hat{\epsilon}_i^2$$

We can define the sample MSE as the sample variance of the residuals and we can define the sample R squared as the ratio of the X-plane to the total variation in the sample.

$$\mathbb{E}_n Y_i^2 = \mathbb{E}_n (\hat{\beta}' X_i)^2 + \mathbb{E}_n \hat{\epsilon}_i^2$$

$$\text{MSE}_{\text{sample}} = \mathbb{E}_n \hat{\epsilon}_i^2$$

$$R_{\text{sample}}^2 := \frac{\mathbb{E}_n (\hat{\beta}' X_i)^2}{\mathbb{E}_n Y_i^2} = 1 - \frac{\mathbb{E}_n \hat{\epsilon}_i^2}{\mathbb{E}_n Y_i^2} \in [0, 1]$$

We know that when P/N is small, the sample linear predictor gets really close to the best linear predictor. Thus, when P/N is small, we expect that sample averages of y squared beta hat xi squared epsilon hat i squared to be close to the population averages of Y squared, beta X squared, and epsilon squared. So in this case, the

sample R squared and the sample MSE will be close to the true quantities, the population R squared and MSE.

WHEN P/N IS SMALL

SAMPLE LINEAR PREDICTOR APPROACHES BLP

$$\mathbb{E}_n Y_i^2 \approx EY^2$$
$$\mathbb{E}_n (\hat{\beta}' X_i)^2 \approx E(\beta' X)^2$$
$$\mathbb{E}_n \hat{\epsilon}_i^2 \approx E\epsilon^2$$

WHEN P/N IS SMALL

SAMPLE LINEAR PREDICTOR APPROACHES BLP

$$R_{sample}^2 \approx R_{pop}^2$$
$$MSE_{sample} \approx MSE_{pop}$$

When P / N is not small, the discrepancy between the two sets of measures can be substantial. And the sample R squared and the sample MSE can be very poor measures of predictability.

The following simulation example will support our reason. In this example, Y and X are statistically independent and generated from the normal distributions with mean need 0 and variance 1. The means that the true linear predictor of Y given X is simply 0 and the true R squared is also 0. Suppose the number of observations is N, the number of regressors is P.

If p = n, then the typical sample R squared will be 1. Which is very far away from the true number of 0. This is an example of extreme over fitting. If p= n / 2 then the typical sample r squared is about half which is still far off from the truth. If p = n / 20, then the typical sample r squared is about 0.05 which is no longer far off from the truth.

$$Y \sim N(0, 1) \quad \text{BLP OF } Y \text{ IS}$$

$$X \sim N(0, I_p) \quad \beta' X = 0$$

$$\text{true } R_{pop}^2 = 0$$

n - OBSERVATIONS p - NUMBER OF REGRESSOR

n - OBSERVATIONS p - NUMBER OF REGRESSORS

IF $p = n$ THEN R_{sample}^2 IS $1 \gg 0$

IF $p = n/2$ THEN R_{sample}^2 IS $.5 \gg 0$

IF $p = n/20$ THEN R_{sample}^2 IS $.05$

We can now see that using sample R squared and MSE to judge predicted performance could be misleading when P over N is not small. So the question is, can we design better metrics for predicted performance, and the answer is yes. The first proposal is to use the adjusted R squared and MSE.

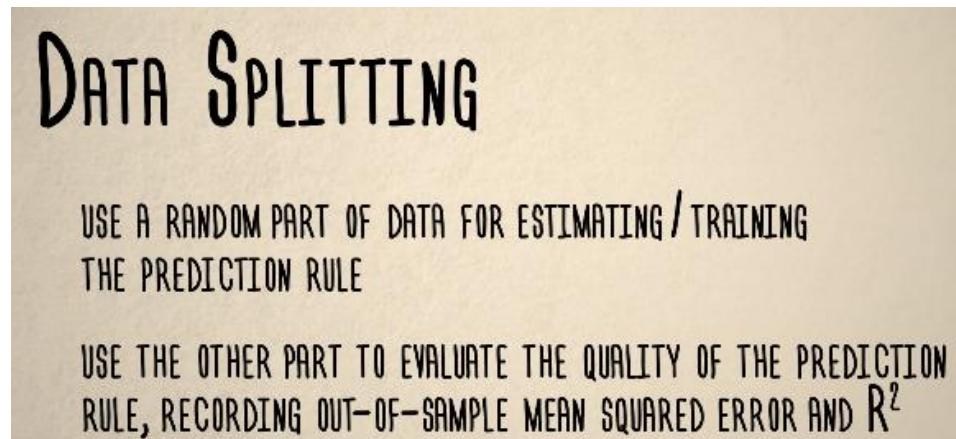
$$R_{adjusted}^2 := 1 - \frac{n}{n-p} \frac{\mathbb{E}_n \hat{\epsilon}_i^2}{\mathbb{E}_n Y_i^2}$$

$$\text{MSE}_{adjusted} = \frac{n}{n-p} \mathbb{E}_n \hat{\epsilon}_i^2$$

We first define the adjusted MSE by rescaling the sample MSE by a factor of n over $n-p$.

And then making the corresponding adjustment to the sample R squared. The re-scaling will correct for over-fitting but under rather strong assumptions.

A more universal way to measure predictive performance is to perform data splitting. First, we use a random part of data, say half of data, for estimating or training the prediction rules. Second, we use the other part of data to evaluate the predictive performance or the rule, recording the out-of-sample MSE or R squared. Accordingly, we call the first part of data, the training sample. And the second part, the testing or validation sample.



T TRAINING n OBSERVATIONS V INDEXES OBSERVATIONS IN THE TEST/VALIDATION SAMPLE	V VALIDATION m OBSERVATIONS
--	--

$$MSE_{test} = \frac{1}{m} \sum_{k \in V} (Y_k - \hat{\beta}' X_k)^2$$

$$R^2_{test} = 1 - \frac{MSE_{test}}{\frac{1}{m} \sum_{k \in V} Y_k^2}$$

2.1.4 Inference for Linear Regression

In this segment, we provide an answer to the inference question. To recall the question, we partition vector of regressors X into D and W, where D represents the target regressors of interest, and W represents other regressors, which we sometimes

call the controls. We write Y = the predicted value which is beta 1D + beta 2W plus noise.

$$X = (D, W')$$

"TARGET" REGRESSORS OF INTEREST
OTHER REGRESSORS OR CONTROLS

$$Y = \underbrace{\beta_1 D + \beta'_2 W}_{predicted\ value} + \underbrace{\epsilon}_{error}$$

And now we recall the inference question. Which is how does the predicted value of Y change if we increase D by a unit holding W fixed?

THE INFERENCE QUESTION HOW DOES THE PREDICTED VALUE OF Y CHANGE
IF WE INCREASE D BY A UNIT, HOLDING W FIXED?

IN THE WAGE EXAMPLE:

WHAT IS THE DIFFERENCE IN PREDICTED WAGES BETWEEN MEN AND WOMEN
WITH THE SAME JOB-RELEVANT CHARACTERISTICS?

POPULATION REGRESSION COEFFICIENT β_1
CORRESPONDING TO THE TARGET REGRESSOR D

D IS THE FEMALE INDICATOR AND β_1 IS THE GENDER WAGE GAP

Next, we attempt to understand better the meaning of beta 1 using a tool called 'partialling-out'.

"PARTIALLING-OUT"

IS AN IMPORTANT TOOL THAT PROVIDES CONCEPTUAL
UNDERSTANDING OF THE REGRESSION COEFFICIENT β_1

IN THE POPULATION, DEFINE THE PARTIALLING-OUT OPERATION AS A
PROCEDURE THAT TAKES A RANDOM VARIABLE V AND CREATES A
"RESIDUAL" \tilde{V} BY SUBTRACTING THE PART OF V THAT IS LINEARLY
PREDICTED BY W

$$\tilde{V} = V - \gamma'_{VW} W$$

$$\gamma_{VW} = \arg \min_{\gamma} E(V - \gamma' W)^2$$

It can be shown that the partialling-out operation is linear.

$$\tilde{V} = V - \gamma'_{VW} W \quad \gamma_{VW} = \arg \min_{\gamma} E(V - \gamma' W)^2$$

$$Y = V + U \implies \tilde{Y} = \tilde{V} + \tilde{U}$$

REGRESSION EQUATION $Y = \beta_1 D + \beta'_2 W + \epsilon$

PARTIALED-OUT $\tilde{Y} = \beta_1 \tilde{D} + \beta'_2 \tilde{W} + \tilde{\epsilon}$

We now apply partialling-out to both sides of our regression equation to get the partialled-out version which implies that we obtain the decomposition which reads tilde Y = beta 1 tilde D + epsilon, where epsilon is uncorrelated with tilde D.

$$Y = V + U \implies \tilde{Y} = \tilde{V} + \tilde{U}$$

REGRESSION EQUATION $Y = \beta_1 D + \beta'_2 W + \epsilon$

PARTIALED-OUT $\tilde{Y} = \beta_1 \tilde{D} + \beta'_2 \tilde{W} + \tilde{\epsilon}$

DECOMPOSITION $\tilde{Y} = \beta_1 \tilde{D} + \epsilon$

This decomposition follows because partialling out eliminates all the Ws, and also leaves epsilon untouched, since epsilon is linearly unpredictable by X, and therefore by W by definition. Moreover, since tilde D is a linear function of X, epsilon and tilde D are not correlated in our decomposition. Now recognize our decomposition as the normal equations for the population regression of tilde Y and tilde D.

$$\text{DECOMPOSITION} \quad \tilde{Y} = \beta_1 \tilde{D} + \epsilon \\ E\epsilon \tilde{D} = 0$$

Theorem (Frisch-Waugh-Lovell, FWL)

The population linear regression coefficient β_1 can be recovered from the population linear regression of \tilde{Y} on \tilde{D} :

$$\beta_1 = \arg \min_{b_1} E(\tilde{Y} - b_1 \tilde{D})^2 = (E\tilde{D}^2)^{-1} E\tilde{D}\tilde{Y},$$

where β_1 is uniquely defined if D cannot perfectly predict by W, i.e. $E\tilde{D}^2 > 0$.

That is, we've just obtained the Frisch-Waugh-Lovell theorem, which states that the population linear regression coefficient beta 1 can be recovered from the population linear regression of tilde Y on tilde D. Beta 1 is a solution to the best linear prediction problem, where we predict tilde Y by a linear function of tilde D. We can also give an explicit formula for beta 1, which is given by the ratio of two averages that you see in this formula. We also see that beta 1 is uniquely defined if D is not perfectly predicted by the Ws. So the tilde D has a non-zero variance.

The Frisch-Waugh-Lovell theorem is a remarkable fact which is useful for both interpretation and estimation.

IT ASSERTS THAT β_1 CAN BE INTERPRETED AS A (UNIVARIATE)
REGRESSION COEFFICIENT OF RESIDUALIZED Y ON RESIDUALIZED D ,
WHERE THE RESIDUALS ARE DEFINED BY PARTIALLING-OUT THE
LINEAR EFFECT OF W FROM Y AND D

It asserts that beta 1 can be interpreted as a regression coefficient of residualized Y on residualized D, where the residuals are defined by taking out or partialling-out the linear effect of W from Y and D.

We next proceed to discuss estimation of beta 1.

IN THE SAMPLE, WE WILL MIMIC
THE PARTIALLING-OUT IN THE POPULATION.

WHEN p/n IS SMALL, WE CAN DO THIS
BY SAMPLE LINEAR REGRESSION

WHEN p/n IS NOT SMALL, USING SAMPLE LINEAR REGRESSION FOR
PARTIALLING-OUT IS NOT A GOOD IDEA. CAN DO PENALIZED REGRESSION
AND DIMENSION REDUCTION INSTEAD. MORE ON THIS LATER.

When p over n is small we can rely on the sample linear regression, that is on ordinary least squares. When p over n is not small using sample linear regression for partialling-out is not a good idea. What we can do instead is do penalized regression or variables selection, which we will discuss later in this module.

ASSUME p/n IS SMALL, SO SAMPLE LINEAR REGRESSION
PROVIDES HIGH-QUALITY PARTIALLING-OUT

BY THE FWL THEOREM APPLIED TO THE SAMPLE INSTEAD OF
POPULATION, THE SAMPLE LINEAR REGRESSION OF Y ON D AND W
GIVES US ESTIMATOR $\hat{\beta}_1$ WHICH IS IDENTICAL TO THE ESTIMATOR
OBTAINED VIA SAMPLE PARTIALLING-OUT

It is still useful to give the formula for hat beta 1 in terms of sample partialling-out where we use checked quantities to denote the residuals that are left after predicting the variables in the sample with the controls.

IT IS STILL USEFUL TO GIVE THE FORMULA
FOR $\hat{\beta}_1$ IN TERMS OF SAMPLE PARTIALLING-OUT

$$\hat{\beta}_1 = \arg \min_b \mathbb{E}_n (\check{Y}_i - b\check{D}_i)^2 = (\mathbb{E}_n \check{D}_i^2)^{-1} \mathbb{E}_n \check{D}_i \check{Y}_i$$

WHERE \check{V}_i DENOTE THE RESIDUAL LEFT AFTER PREDICTING V_i WITH
CONTROLS W_i . HERE WE USE THE SAMPLE LINEAR REGRESSION

$$\check{V}_i = V_i - \hat{\gamma}'_{VW} W_i$$

$$\hat{\gamma}_{VW} = \arg \min_{\gamma} \mathbb{E}_n (V_i - \gamma' W_i)^2$$

The population partialling-out is replaced here by the sample partialling-out where V replace the population expectation by the empirical expectation.

Using the formula, it can be shown that the following results calls.

Theorem (Inference)

If p/n is small, then the estimation error in \check{D}_i and \check{Y}_i has no first order effect on $\hat{\beta}_1$, and

$$\hat{\beta}_1 \stackrel{d}{\sim} N(\beta_1, V/n)$$

where

$$V = (\mathbb{E} \check{D}^2)^{-1} \mathbb{E} (\check{D}^2 \epsilon^2) (\mathbb{E} \check{D}^2)^{-1}.$$

If p over n is small, then the estimation error in the estimated residualized quantities has a negligible effect on $\hat{\beta}_1$, and $\hat{\beta}_1$ is approximately distributed as normal variable with mean β_1 and variance V/n where the expression of the variance appears over here. In words, we can say that the estimator $\hat{\beta}_1$ concentrates in a square root of V/n neighborhood of β_1 with deviations controlled by the normal law.

We can now define the standard error for hat beta 1 as square root of hat V over n, where hat V is an estimator of V.

THE ABOVE STATEMENT MEANS THAT $\hat{\beta}_1$ CONCENTRATES IN A $\sqrt{V/n}$ -NEIGHBORHOOD OF β_1 , WITH DEVIATIONS CONTROLLED BY THE NORMAL LAW.
THE STANDARD ERROR OF $\hat{\beta}_1$ IS $\sqrt{\hat{V}/n}$, WHERE \hat{V} IS AN ESTIMATOR OF V .

This result implies that the interval given by the estimate ± 2 standard errors covers the true value of beta 1, for most realizations of the data sample. Or more precisely approximately 95% of realizations of the data sample. If we replace 2 by other constants, we get other coverage probabilities. In other words, if our data sample is not extremely unusual, the interval covers the truth.

THE RESULT IMPLIES THAT THE INTERVAL
 $[\hat{\beta}_1 - 2\sqrt{\hat{V}/n}, \hat{\beta}_1 + 2\sqrt{\hat{V}/n}]$
COVERS β_1 FOR MOST REALIZATIONS OF THE SAMPLE, MORE PRECISELY,
APPROXIMATELY 95% OF THE REALIZATIONS OF THE SAMPLE.
(IF WE REPLACE 2 BY OTHER CONSTANTS, WE GET OTHER COVERAGE PROBABILITIES.)

For this reason, this interval is called the confidence interval.

IN OTHER WORDS, IF OUR SAMPLE IS NOT ATYPICAL, THE INTERVAL
COVERS THE TRUTH.

THIS INTERVAL IS CALLED THE CONFIDENCE INTERVAL.

FOR EXAMPLE, IN THE WAGE EXAMPLE,
OUR ESTIMATE OF GENDER HOURLY WAGE GAP IS ABOUT
-2\$ AND THE 95% CONFIDENCE INTERVAL
IS ABOUT [-1\$, -3\$].

Let us summarize, first we interpreted beta 1 as the regression coefficient in the univariate regression of the response variable and the target variable, after we have removed the linear effect of other variables. Second, we noted that this result is useful for interpretation and understanding of the regression coefficient. This result will also be super useful for setting up inference in modern high-dimensional settings which we will discussed later in this module. And next, we will carry out a case study for our wage example.

WE INTERPRETED β_1 AS THE REGRESSION COEFFICIENT IN THE UNIVARIATE
REGRESSION OF THE RESPONSE VARIABLE ON THE TARGET VARIABLE, AFTER
WE HAVE REMOVED THE LINEAR EFFECT OF THE OTHER VARIABLES

THIS RESULT IS USEFUL FOR INTERPRETATION AND UNDERSTANDING OF
THE REGRESSION COEFFICIENTS

IT WILL ALSO BE USEFUL FOR SETTING UP INFERENCE IN MODERN
HIGH-DIMENSIONAL SETTINGS

NEXT, WE WILL CARRY OUT A CASE STUDY FOR THE WAGE EXAMPLE

2.1.5 Other Types of Regression

REGRESSIONS INDUCED BY USING
NONLINEAR FUNCTIONAL FORMS,
FOR EXAMPLE, LOGISTIC REGRESSION

REGRESSIONS INDUCED BY USING DIFFERENT TYPES OF
LOSS FUNCTIONS INSTEAD OF THE SQUARED LOSS FUNCTION,
FOR EXAMPLE, MEDIAN AND QUANTILE REGRESSION

IN NONLINEAR REGRESSION, WE USE A NONLINEAR FUNCTION OF X ,
 $p(X, \beta)$, TO PREDICT Y . FOR INSTANCE, THE LOGISTIC REGRESSION EMPLOYS
$$p(X, \beta) = \exp(\beta' X) / (1 + \exp(\beta' X))$$

THIS FUNCTIONAL FORM IS PARTICULARLY ATTRACTIVE WHEN THE
RESPONSE VARIABLE Y IS BINARY, TAKING VALUES OF 0 AND 1, SO THAT
THE PREDICTED VALUES ARE NATURALLY CONSTRAINED BETWEEN 0 AND 1

THE PROBLEM OF PREDICTING BINARY Y IS A BASIC EXAMPLE OF A
CLASSIFICATION PROBLEM. WE CAN INTERPRET THE PREDICTED VALUES
IN THIS CASE AS APPROXIMATING PROBABILITY OF $Y=1$

THE PARAMETERS ARE ESTIMATED EITHER BY SOLVING THE SAMPLE
NONLINEAR LEAST SQUARES PROBLEM

$$\min_{b \in \mathbb{R}^p} \mathbb{E}_n(Y_i - p(X_i, b))^2$$

OR, IN THE CASE OF BINARY OUTCOMES, BY MAXIMIZING THE LOGISTIC
LOG-LIKELIHOOD FUNCTION

$$\max_{b \in \mathbb{R}^p} \mathbb{E}_n\{Y_i \ln p(X_i, b) + (1 - Y_i) \ln(1 - p(X_i, b))\}$$

WE PREVIOUSLY USED THE SQUARED ERROR LOSS TO SET UP THE BEST LINEAR PREDICTOR.
WE CAN CALL THIS APPROACH THE LINEAR MEAN REGRESSION,
BECAUSE THE BEST PREDICTOR OF Y UNDER SQUARED LOSS
IS THE CONDITIONAL EXPECTATION OF Y GIVEN X

WHAT IF WE USE THE ABSOLUTE DEVIATION LOSS INSTEAD? THEN WE
WILL OBTAIN THE MEDIAN REGRESSION, BECAUSE WE WILL IMPLICITLY
TRY TO ESTIMATE THE MEDIAN VALUE OF Y CONDITIONAL ON X

$$\min_{b \in \mathbb{R}^p} \mathbb{E}|Y - b'X|$$

IN THE POPULATION WE USE THE POPULATION EXPECTATION E AND IN
THE SAMPLE WE USE THE EMPIRICAL EXPECTATION \mathbb{E}_n INSTEAD

JUST LIKE MEAN REGRESSION, MEDIAN REGRESSION CAN ALSO BE MADE NONLINEAR, SO WE CAN REPLACE $b'X$ BY $p(X, b)$

MEDIAN REGRESSIONS ARE GREAT FOR CASES WHEN THE OUTCOMES HAVE HEAVY TAILS OR OUTLIERS. THEY ARE MUCH BETTER BEHAVED IN THIS CASE THAN THE MEAN REGRESSION

IF WE USE AN ASYMMETRIC ABSOLUTE DEVIATION LOSS, THEN WE END UP WITH THE ASYMMETRIC ABSOLUTE DEVIATION REGRESSION OR QUANTILE REGRESSION

$$\min_{b \in \mathbb{R}^p} E(\tau|Y - b'X|1(Y > b'X) + (1 - \tau)|Y - b'X|1(Y < b'X))$$

IN THE POPULATION WE USE THE POPULATION EXPECTATION AND IN THE SAMPLE WE USE THE EMPIRICAL EXPECTATION \mathbb{E}_n INSTEAD

THE WEIGHT τ , THAT APPEARS IN THE DEFINITION OF THE ABSOLUTE DEVIATION LOSS, SPECIFIES THE $\tau \times 100$ -TH PERCENTILE OF Y THAT WE ARE TRYING TO MODEL

QUANTILE REGRESSIONS ARE GREAT FOR DETERMINING THE FACTORS THAT AFFECT THE OUTCOME IN THE TAILS

FOR EXAMPLE

IN RISK MANAGEMENT, WE MIGHT FORECAST THE EXTREMAL CONDITIONAL PERCENTILES OF Y USING THE INFORMATION X . THIS IS CALLED CONDITIONAL VALUE-AT-RISK ANALYSIS

IN MEDICINE, WE COULD BE INTERESTED IN HOW SMOKING AND OTHER CONTROLLABLE FACTORS X AFFECT VERY LOW PERCENTILES OF INFANT BIRTH WEIGHTS Y

IN SUPPLY MANAGEMENT, WE COULD TRY TO PREDICT THE INVENTORY LEVEL OF A PRODUCT THAT IS ABLE TO MEET THE 90-TH PERCENTILE OF DEMAND Y GIVEN THE ECONOMIC CONDITIONS DESCRIBED BY X

WE HAVE BRIEFLY DISCUSSED NON-LINEAR REGRESSIONS AS WELL QUANTILE REGRESSIONS AND THEIR USES

IN THE NEXT BLOCK WE WILL CONSIDER MODERN LINEAR AND NONLINEAR REGRESSIONS WHICH ARE MOTIVATED BY HIGH-DIMENSIONAL DATA

2.2.1 Modern Linear Regression for High-Dimensional Data

We begin a new block of segments where we consider modern linear regression for high dimensional data. We consider the linear regression model Y equals Beta X plus Epsilon where Beta X is the population best linear predictor, or BLP, of Y using X . Or simply the population linear regression function, and epsilon is uncorrelated with X .

$$Y = \beta' X + \epsilon$$

$\beta' X$ IS THE POPULATION BEST LINEAR PREDICTOR OF Y
POPULATION LINEAR REGRESSION FUNCTION

Here, the regressor X is p-dimsensional, that is, there are P regressors and the sample size is N , and P is large, possibly larger than N .

$\beta' X$ IS THE POPULATION BEST LINEAR PREDICTOR OF Y
POPULATION LINEAR REGRESSION FUNCTION

$$X = (X_j)_{j=1}^p$$

p -DIMENSIONAL p REGRESSORS

p IS LARGE, POSSIBLY MUCH LARGER THAN n

One reason for having high dimensional regressors in our model is the increasing availability of modern rich data, or quote unquote, big data. Many modern data sets are rich in that they have many recorded features or regressors. For example, in house pricing and appraisal analysis, we can make use of numerous housing characteristics. In demand analysis, we can rely on price and product characteristics, and in the analysis of court decisions, we can employ many of the judges' characteristics.

BIG DATA

HOUSING CHARACTERISTICS IN HOUSE PRICING/APPRaisal ANALYSIS
PRICE AND PRODUCT CHARACTERISTICS AT THE POINT OF PURCHASE
JUDGE CHARACTERISTICS IN THE ANALYSIS OF JUDICIAL DECISIONS

Another reason for having high dimensional regressors in our model is the use of constructed regressors.

CONSTRUCTED REGRESSORS

IF Z ARE "RAW" REGRESSORS/FEATURES,
THE CONSTRUCTED REGRESSORS ARE

$$X = P(Z) = (P_1(Z), \dots, P_p(Z))'$$

WHERE THE SET OF TRANSFORMATIONS $P(Z)$
IS SOMETIMES CALLED THE "DICTIONARY"

IN THE WAGE EXAMPLE, WE USED QUADRATIC AND CUBIC
TRANSFORMATIONS OF EXPERIENCE, AS WELL AS
INTERACTIONS (PRODUCTS) OF THESE REGRESSORS
WITH EDUCATION & GEOGRAPHIC INDICATORS.

CONSTRUCTED REGRESSORS

THE USE OF CONSTRUCTED REGRESSORS ALLOWS US TO BUILD MORE
FLEXIBLE AND POTENTIALLY BETTER PREDICTION RULES

$$\beta' X = \beta' P(Z)$$

THE PREDICTION RULE $\beta' X$ IS STILL LINEAR WITH RESPECT TO THE
PARAMETER β AND WITH RESPECT TO
THE CONSTRUCTED REGRESSORS $X = P(Z)$

What is the best prediction rule for Y using Z ? It turns out that the best prediction rule is the conditional expectation of Y given Z . We denote this prediction rule by g of Z and we call it the regression function of Y on Z . This is the best prediction rule among all rules, and it is generally better than the best linear prediction rule.

THE BEST PREDICTION RULE

$$g(Z) = E(Y | Z)$$

THE CONDITIONAL EXPECTATION OF Y GIVEN Z

REGRESSION FUNCTION OF Y ON Z

IT CAN BE SHOWN THAT $g(Z)$
SOLVES THE BEST PREDICTION PROBLEM:

$$\min_{m(Z)} E(Y - m(Z))^2$$

WHERE WE MINIMIZE THE MEAN SQUARED PREDICTION ERROR (**MSE**)
AMONG ALL PREDICTION RULES $m(Z)$ (LINEAR OR NONLINEAR IN Z)

BY USING $\beta'P(Z)$ WE ARE IMPLICITLY APPROXIMATING
THE BEST PREDICTOR $g(Z)$

INDEED, IT CAN BE SHOWN THAT FOR ANY PARAMETER b

$$E(Y - b'P(Z))^2 = E(g(Z) - b'P(Z))^2 + \text{const.}$$

THAT IS, THE MEAN SQUARED PREDICTION ERROR IS EQUAL TO THE
MEAN SQUARED APPROXIMATION ERROR OF $b'P(Z)$ TO $g(Z)$, PLUS A
CONSTANT THAT DOES NOT DEPEND ON b , SO THAT MINIMIZING THE
FORMER IS THE SAME AS MINIMIZING THE LATTER

We now conclude that the BLP beta $P(Z)$ is the best linear approximation, or BLA, to the regression function $g(Z)$.

By using a richer and richer dictionary of transformations, $P(Z)$, we can make the BLP beta $P(Z)$ approximate $g(Z)$ better and better. We can illustrate this point by the following simulations.

Suppose $Z \sim \text{Uniform}(0, 1)$, and

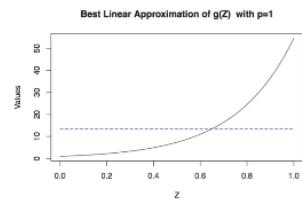
$$g(Z) = \exp(4 \cdot Z).$$

We use

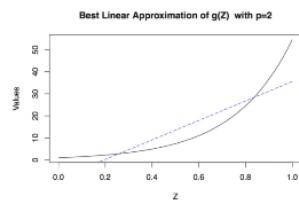
$$P(Z) = \underbrace{(1, Z, Z^2, \dots, Z^{p-1})'}_{p \text{ terms}}$$

to form the BLA/BLP, $\beta' P(Z)$.

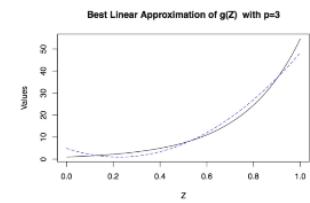
Simulation Example: $p = \dim(P(Z)) = 1$



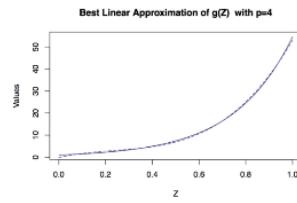
Simulation Example: $p = \dim(P(Z)) = 2$



Simulation Example: $p = \dim(P(Z)) = 3$



Simulation Example: $p = \dim(P(Z)) = 4$



Suppose that Z is uniformly distributed on the unit interval, and the true regression function $g(Z)$ is the exponential function of 4 times Z . Suppose we don't know this and we use the linear form beta $P(Z)$ to provide approximations for $g(Z)$? Suppose we use $P(Z)$ that consists of polynomial transformations of Z , consisting of the first P terms of the polynomial series 1, Z , Z squared, Z cubed, and so on. We use this dictionary to form the BLA or BLP beta $P(Z)$. We now provide a sequence of figures that illustrate the approximation properties of the BLA or BLP corresponding to $P(Z)$ equal to 1, 2, 3 and 4.

With P equal 1, we get a constant approximation to the regression function $g(Z)$, and as we can see, the quality of this approximation is very poor. With $P = 2$, we get a

linear inside approximation for $g(Z)$, and as the figure shows, the quality of this approximation is still very poor. With $P = 3$, we get a quadratic in Z approximation for $g(Z)$, and now, the quality is quite good all of a sudden. This explains why using the linear transformations over [INAUDIBLE] regressors is a good idea. Now with P equals 4 we get a cubic in Z approximation for $g(Z)$ and the quality of approximation becomes simply excellent.

WE PROVIDED TWO MOTIVATIONS FOR USING HIGH-DIMENSIONAL X IN PREDICTION:

1. MODERN DATA SETS HAVE HIGH-DIMENSIONAL FEATURES THAT CAN BE USED AS REGRESSORS.
2. NON-LINEAR TRANSFORMATIONS OF FEATURES/RAW REGRESSORS AND THEIR INTERACTIONS GENERATE CONSTRUCTED REGRESSORS.

USING MANY CONSTRUCTED REGRESSORS ALLOWS US TO BETTER APPROXIMATE THE BEST PREDICTION RULE – THE CONDITIONAL EXPECTATION OF THE OUTCOME Y GIVEN THE RAW REGRESSORS Z .

2.2.2 High-Dimensional Sparse Models and Lasso

In this segment, we will talk about high-dimensional sparse models and a penalized regression method called the Lasso.

POPULATION LINEAR REGRESSION FUNCTION

$$Y = \beta' X + \epsilon = \sum_{j=1}^p \beta_j X_j + \epsilon$$

$$X = (X_j)_{j=1}^p$$

Here, we consider the regression model $Y = \beta X + \epsilon$, where β is a population-best linear predictor of Y using X . Or simply the population linear regression function.

The regressor X is p -dimensional, with components denoted by X_j . That is, there are p regressors, and p is large, possibly much larger than N , where N is the sample size.

POPULATION LINEAR REGRESSION FUNCTION

$$X = (X_j)_{j=1}^p$$

p/N IS NOT SMALL

In order to simplify the discussion, we assume that each X_j has a unit variance in the sample.

Here, we are dealing with the high-dimensional setting, where p/N is not small. Classical linear regression, or ordinary least squares fails in these settings, because it overfits the data, as we have seen in part one of our module. We need to make some assumptions and modify the classical regression method to deal with the high-dimensional case.

APPROXIMATE SPARSITY

A SMALL GROUP OF REGRESSORS THAT
HAVE RELATIVELY LARGE COEFFICIENTS,

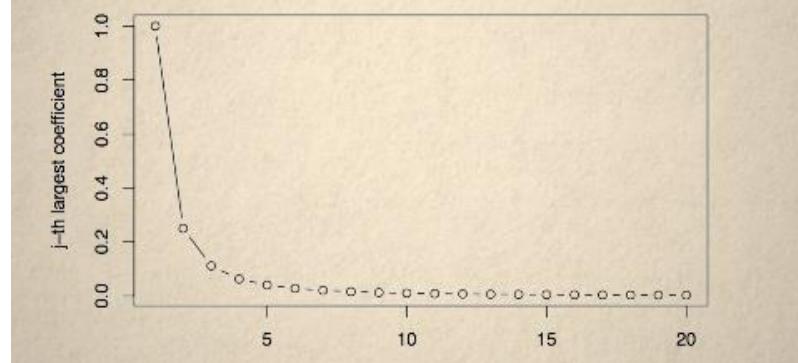
THAT CAN BE USED TO
APPROXIMATE THE $BLP \beta' X$

One intuitive assumption is approximate sparsity, which informally means that there is a small group of regressors that have relatively large coefficients, that can be used

to approximate the BLP beta 'X quite well. The rest of regressors have relatively small coefficients.

An example of an approximately sparse model is the linear regression model with the regression coefficients beta j given by $1/j^2$ as j ranges from 1 to p. The figure that you see here shows the graph of these coefficients.

$$\beta_j = 1/j^2, \quad j = 1, \dots, p.$$



As we can see, the coefficients has decreased quite fast, with only three or four regression coefficients that appear to be large.

$$|\beta|_{(j)} \leq A j^{-a}, \quad a > 1/2,$$

Formally, the approximate sparsity means that the sorted absolute values of the coefficients decrease to zero fast enough. Namely, the js largest in absolute value coefficient is at most of size j into the power of minus a times a constant, where a is greater than one half. Here, the constant a measures the speed of decay.

For estimation purposes, we have a random sample of Y_i and X_i , where i ranges from 1 to n. We seek to construct a good linear predictor, head-beta X, which works well when p over N is not small. We can construct head-beta as a solution of the following penalized regression problem called Lasso.

LASSO REGRESSION PROBLEM

$$\sum_i (Y_i - b'X_i)^2 + \lambda \cdot \sum_{j=1}^p |b_j|$$

Here, we are minimizing the sample mean-squared error that results from predicting Y_i with bX plus a penalty term, which penalizes the size of the coefficients, b_j s, by their absolute values. We control the degree of penalization by the penalty level λ . A theoretically justified penalty level for Lasso is given by the formula that you see here.

PENALTY LEVEL

$$\lambda = \sqrt{E\epsilon^2} 2 \sqrt{2n \log(pn)}.$$

THIS PENALTY LEVEL ENSURES THAT THE LASSO PREDICTOR $\hat{\beta}'X$ DOES NOT OVERFIT THE DATA AND DELIVERS GOOD PREDICTIVE PERFORMANCE UNDER APPROXIMATE SPARSITY

ANOTHER GOOD WAY TO PICK PENALTY LEVEL IS BY CROSS-VALIDATION
(WHICH IS A DATA-SPLITTING METHOD THAT WE WILL DISCUSS LATER IN THE MODULE)

Another good way to pick penalty level is by cross-validation, which uses repeated splitting of data into training and testing samples to measure predictive performance.

LASSO IMPOSES THE APPROXIMATE SPARSITY ON THE COEFFICIENTS
 $b = \hat{\beta}$, JUST LIKE IN THE ASSUMPTION

IT PASSES DOWN ALL OF THE COEFFICIENTS TO ZERO, AS MUCH AS POSSIBLE AND AT AN EQUAL RATE, WITHOUT SACRIFICING TOO MUCH FIT, AND IT ENDS UP SETTING MANY OF THESE COEFFICIENTS TO ZERO

Simulation Example:

$$Y = \beta' X + \epsilon, \quad X \sim N(0, I_p), \quad \epsilon \sim N(0, 1),$$

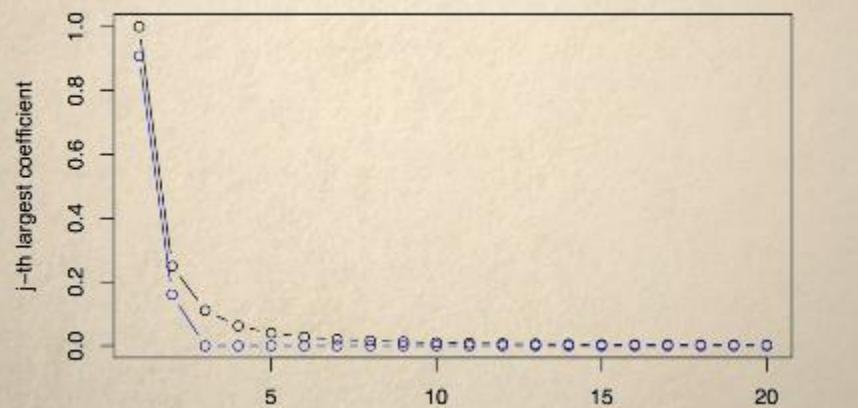
where

$$\beta_j = 1/j^2, \quad j = 1, \dots, p$$

and

$$n = 300, \quad p = 1000.$$

The figure shows that $\hat{\beta}$ is sparse and is close to β .



The figure shows that hat beta, in blue, is indeed sparse and is close to the true coefficient vector beta, in black, indeed. Most of hat-beta js are set to 0, except for several coefficients that align quite well with the largest coefficient beta j, or the true coefficient vector.

This really shows that Lasso is able to leverage approximate sparsity to deliver good approximation to the true coefficients.

From the figure, we see that Lasso sets most of the regression coefficients to zero. It basically figures out approximately, though not perfectly, the right set of regressors.

In practice, we often use the Lasso-selected set of regressors to refit the model by least squares. This method is called the least squares post Lasso, or simply post-Lasso.

WE SEE THAT LASSO SETS MOST OF REGRESSION COEFFICIENTS TO ZERO.
IT FIGURES OUT APPROXIMATELY THE RIGHT SET OF REGRESSORS

WE CAN THEN USE THE LASSO-SELECTED SET OF REGRESSORS
TO REFIT THE MODEL BY LEAST SQUARES.

THIS METHOD IS CALLED THE “LEAST SQUARES POST LASSO”
OR SIMPLY POST-LASSO

Post Lasso

Post-Lasso does not shrink large coefficients to zero as much as Lasso does, and often improves over Lasso in terms of prediction

Lasso

The best linear prediction rule (out-of-sample) is $\beta'X$.
Does $\hat{\beta}'X$ provide a good approximation to $\beta'X$?

We are trying to estimate p parameters β_1, \dots, β_p ,
imposing the approximate sparsity via penalization

LASSO

UNDER SPARSITY, ONLY A FEW, SAY s , PARAMETERS WILL BE "IMPORTANT".

WE CAN CALL s THE EFFECTIVE DIMENSION

LASSO WILL APPROXIMATELY FIGURE OUT WHICH PARAMETERS ARE IMPORTANT

INTUITIVELY, TO ESTIMATE EACH OF THE "IMPORTANT" s PARAMETERS WELL,
WE NEED MANY OBSERVATIONS PER SUCH PARAMETER

LASSO

THIS MEANS THAT n/s MUST BE LARGE,
OR, EQUIVALENTLY s/n MUST BE SMALL

Theorem

Under regularity conditions and under the approximate sparsity assumption, with probability approaching 1 as $n \rightarrow \infty$,

$$\sqrt{\text{E}_X(\beta'X - \hat{\beta}'X)^2} \leq \text{const} \cdot \sqrt{\text{E}\epsilon^2} \sqrt{\frac{s \log(pn)}{n}},$$

where E_X denotes expectation with respect to X , and the effective dimension is

$$s = \text{const} \cdot n^{\frac{1}{2\alpha}}.$$

If $s \log(pn)/n$ is small, Lasso and Post-Lasso regression come close to the population regression function/best linear predictor.

UNDER APPROXIMATE SPARSITY LASSO AND POST-LASSO WILL APPROXIMATE THE BEST LINEAR PREDICTOR WELL.

THIS MEANS THAT THEY WON'T OVERFIT THE DATA, AND WE CAN USE THE SAMPLE AND ADJUSTED R^2 AND MSE TO ASSESS OUT-OF-SAMPLE PREDICTIVE PERFORMANCE.

OF COURSE, IT IS ALWAYS A GOOD IDEA TO VERIFY THE OUT-OF-SAMPLE PREDICTIVE PERFORMANCE BY USING TEST/VALIDATION SAMPLES.

SUMMARY

WE HAVE DISCUSSED APPROXIMATE SPARSITY AS ONE ASSUMPTION THAT MAKES IT POSSIBLE TO PERFORM ESTIMATION/PREDICTION WITH HIGH-DIMENSIONAL DATA

LEAST SQUARED ERROR IS A REGRESSION METHOD THAT IMPOSES APPROXIMATE SPARSITY BY PENALIZATION

UNDER APPROXIMATE SPARSITY, LEAST SQUARED ERROR IS ABLE TO APPROXIMATE THE BEST LINEAR PREDICTOR AND THUS PRODUCES HIGH QUALITY PREDICTION.

2.2.3 Inference with Lasso

RECALL THE INFERENCE QUESTION:
HOW DOES THE PREDICTED VALUE OF Y CHANGE
IF WE INCREASE THE COMPONENT D OF X BY A UNIT,
HOLDING W , THE OTHER COMPONENTS X , FIXED?

THE ANSWER IS THE POPULATION REGRESSION COEFFICIENT β_1
CORRESPONDING TO THE REGRESSOR D

AND HERE WE WOULD LIKE TO DISCUSS HOW TO USE LASSO TO ESTIMATE AND
CONSTRUCT CONFIDENCE INTERVALS FOR β_1 IN THE HIGH-DIMENSIONAL
SETTING.

We write the regression equation as Y equals beta 1 D plus beta 2 w plus epsilon, where d is the target regressor and w 's our d controls. We recall from part one of our module that after portioning out we ended up with a simplified regression equation. $\tilde{Y} = \tilde{D} \beta_1 + \epsilon$ where ϵ is uncorrelated with \tilde{D} . Here, \tilde{Y} and \tilde{D} , are the residuals that are left after parcelling out the linear effect of W .

As we argued in part one, this allows us to obtain β_1 as a coefficient in the linear regression of \tilde{Y} and \tilde{D} . β_1 in the population recovers the partialling-out procedure. For estimation purposes, we have a random sample of Y_i s and X_i s. Our main idea here is that, we will mimic in the sample the partialling-out procedure in the population.

$$Y = \beta_1 D + \beta_2' W + \epsilon$$

$$\tilde{Y} = \beta_1 \tilde{D} + \epsilon, \quad E\epsilon \tilde{D} = 0$$

$$\tilde{D} = D - \gamma_{DW}' W, \quad \gamma_{DW} = \arg \min_{\gamma \in \mathbb{R}^p} E(D - \gamma' W)^2$$

$$\tilde{Y} = Y - \gamma_{YW}' W, \quad \gamma_{YW} = \arg \min_{\gamma \in \mathbb{R}^p} E(Y - \gamma' W)^2$$

$$\beta_1 = \arg \min_{b_1 \in \mathbb{R}} E(\tilde{Y} - b_1 \tilde{D})^2 = (E\tilde{D}^2)^{-1} E\tilde{D} \tilde{Y}$$

$$(Y_i, X_i)_{i=1}^n$$

OUR MAIN IDEA HERE IS THAT, WE WILL MIMIC IN THE SAMPLE THE PARTIALLING-OUT PROCEDURE IN THE POPULATION.

PREVIOUSLY, WHEN ρ OVER n WAS SMALL, WE EMPLOYED LEAST SQUARES AS THE PREDICTION METHOD IN THE PARTIALLING-OUT STEPS.

HERE ρ OVER n IS NOT SMALL, AND WE EMPLOY INSTEAD THE LASSO METHOD IN THE PARTIALLING-OUT STEPS.

So let us explain this in more detail. First we run the lasso regression of Y_i on W_i and of D_i on W_i and keep the resulting residuals called check W_i and check D_i .

$$\hat{\gamma}_{YW} = \arg \min_{\gamma \in \mathbb{R}^p} \sum_i (Y_i - \gamma' W_i)^2 + \lambda_1 \sum_j |\gamma_j|$$

$$\hat{\gamma}_{DW} = \arg \min_{\gamma \in \mathbb{R}^p} \sum_i (D_i - \gamma' W_i)^2 + \lambda_2 \sum_j |\gamma_j|$$

$$\check{Y}_i = Y_i - \hat{\gamma}'_{YW} W_i$$

$$\check{D}_i = D_i - \hat{\gamma}'_{DW} W_i$$

Second we run the least squares of check W_i on check D_i . The resulting estimated regression coefficient is our estimator check beta one.

$$\check{\beta}_1 = \arg \min_{b_1 \in \mathbb{R}} \mathbb{E}_n (\check{Y}_i - b_1 \check{D}_i)^2 = (\mathbb{E}_n \check{D}_i \check{D}_i)^{-1} \mathbb{E}_n \check{D}_i \check{Y}_i$$

Our portioning out procedure was loss all. It relies on approximate sparsity in the two portioning out steps. Indeed, theoretically, the procedure will work well if the population regression coefficients and the two parceling out steps are approximately sparse, with a sufficiently high speed of degrees of assorted values, as shown in the conditions that you see here.

We need approximate sparsity of the population regression coefficients γ_{YW} and γ_{DW} , with a sufficiently high speed of decrease in the sorted coefficients, namely,

$$|\gamma_{YW}|_{(j)} \leq A j^{-a}, \quad |\gamma_{DW}|_{(j)} \leq A j^{-a} \quad a > 1, \quad j = 1, \dots, p$$

Remark. Approximate sparsity of γ_{YW} follows from approximate sparsity of β_2 , namely

$$|\beta_2|_{(j)} \leq A j^{-a}, \quad a > 1, \quad j = 1, \dots, p$$

Theorem (Inference in High-Dimensional Regression)

Under the stated approximate sparsity and other regularity conditions, the estimation error in \check{D}_i and \check{Y}_i has no first order effect on $\check{\beta}_1$, and

$$\check{\beta}_1 \stackrel{d}{\sim} N(\beta_1, V/n)$$

where

$$V = (E\check{D}^2)^{-1} E(\check{D}^2 \epsilon^2) (E\check{D}^2)^{-1}.$$

The above statement means that $\check{\beta}_1$ concentrates in a $\sqrt{V/n}$ -neighborhood of β_1 , with deviations controlled by the normal law.

We can use this result just like we used the analogous result for the low dimensional case in block one. We can define the standard error of check beta one as square root of hat v over n. Where hat v is an estimator of v. This is our quantification of uncertainty about that one. We then provide the approximate 95% confidence interval for beta 1, which is given by the estimate check beta 1 + or- two standard errors.

CONFIDENCE INTERVAL

95%

$$[\check{\beta}_1 - 2\sqrt{\hat{V}/n}, \check{\beta}_1 + 2\sqrt{\hat{V}/n}]$$

SUMMARY

WE HAVE ESTIMATED THE TARGET REGRESSION COEFFICIENT β_1 IN THE HIGH-DIMENSIONAL REGRESSION PROBLEM.

WE USE LASSO TO OBTAIN ESTIMATES OF THE OUTCOME Y AND TARGET REGRESSOR D NET OF THE EFFECT OF OTHER REGRESSORS W , AND THEN RUN LEAST SQUARES OF ONE ON THE OTHER.

THE RESULTING ESTIMATOR $\check{\beta}_1$ IS A HIGH-QUALITY ESTIMATOR OF β_1

2.2.5 Other Penalized Regression Methods. Cross-Validation

$$Y = \beta'X + \epsilon, \quad E(\epsilon|X) = 0, \quad \dim X = p$$

$$(Y_i, X_i)_{i=1}^n$$

$$\hat{f}(X) = \hat{\beta}'X$$

In this segment, we overview other penalized regression methods and also discuss cross-validation, which is a method to choose tuning parameters for the prediction rules. We are interested in prediction in the linear model $Y = \beta'X + \epsilon$, where ϵ isn't correlated with X . And we have random sample of (Y_i, X_i) . Our generic predictor will take the linear form $\hat{f}(x) = \hat{\beta}'x$. The idea of penalized regression is to choose the coefficients $\hat{\beta}$, to avoid overfitting in the sample. This is achieved by penalizing the size of the coefficients by various penalty functions.

Ideally, we should choose the level of penalization to minimize the out-of-sample mean squared prediction error. We first consider the Ridge method.

The Ridge method estimates coefficients by penalized least squares, where we minimize the sum of squared prediction error plus the penalty term given by the sum of the squared values of the coefficients times a penalty level, λ .

RIDGE METHOD

THE RIDGE METHOD ESTIMATES COEFFICIENTS BY PENALIZED LEAST SQUARES, WHERE WE MINIMIZE THE SUM OF SQUARED PREDICTION ERROR PLUS A PENALTY GIVEN BY THE SUM OF THE SQUARED VALUES OF THE COEFFICIENTS TIMES A PENALTY LEVEL

$$\hat{\beta}(\lambda) = \arg \min_{b \in \mathbb{R}^p} \sum_{i=1}^n (Y_i - b'X_i)^2 + \lambda \sum_j b_j^2$$

We can see this in the formula given here. If we look at the formula, we notice that analogous to the Lasso, Ridge penalty presses down or penalizes the coefficients without sacrificing too much fit. In contrast to Lasso, Ridge penalizes the large values of coefficients much more aggressively and small values much less aggressively.

RIDGE METHOD

UNLIKE LASSO, RIDGE DOES NOT DO VARIABLE SELECTION BECAUSE IT PENALIZES SMALL COEFFICIENTS VERY LIGHTLY

THE RIDGE PREDICTION $\hat{\beta}' X$ PERFORMS WELL FOR THE "DENSE" MODELS,
WHERE THE β_j ARE NOT ZERO BUT ARE ALL SMALL, WITHOUT
NECESSARILY BEING APPROXIMATELY SPARSE

IN PRACTICE, THE PENALTY LEVEL λ IS CHOSEN BY CROSS-VALIDATION

Because Ridge penalizes small coefficients very lightly, the Ridge fit is never sparse. And unlike Lasso, Ridge does not set estimated coefficients to zero, and so it does not do variable selection. Because of this property, Ridge predictor hat beta X is especially well suited for prediction in the dense models, where the beta js are all small without necessarily being approximately sparse. In this case, it can easily outperform the Lasso predictor. Finally, we noted in practice, we can choose the penalty level lambda in Ridge, by cross-validation which we will discuss later in this segment.

One popular modification is the method called the elastic net. Here we estimate the coefficients by penalized least squares with penalty given by the linear combination of the Lasso and the Ridge penalties as you see in this formula. In the formula, we see that the penalty function has two penalty levels, lambda 1 and lambda 2, which could be chosen by cross-validation in practice.

ELASTIC NET

THE ELASTIC NET METHOD ESTIMATES THE COEFFICIENTS BY PENALIZED LEAST SQUARES WITH PENALTY GIVEN BY A LINEAR COMBINATION OF THE LASSO AND RIDGE PENALTIES

$$\hat{\beta}(\lambda_1, \lambda_2) = \arg \min_{b \in \mathbb{R}^p} \sum_i (Y_i - b'X_i)^2 + \lambda_1 \sum_j b_j^2 + \lambda_2 \sum_j |b_j|$$

Now, let us look at the formula carefully. We see that the elastic net penalizes large coefficients as aggressively as Ridge. And we also see that it penalizes small coefficients as aggressively as Lasso. By selecting different values of penalty levels, lambda 1 and lambda 2, we could have more flexibility for building a good prediction rule than with just Ridge or with just Lasso. We also note that the elastic net doesn't perform variable selection unless we completely shut down the Lasso penalty by setting penalty level lambda 1 equals zero.

LAVA

$$\hat{\beta}(\lambda_1, \lambda_2) = \hat{\gamma}(\lambda_1, \lambda_2) + \hat{\delta}(\lambda_1, \lambda_2)$$

$$= \arg \min_{\gamma+\delta \in \mathbb{R}^p} \sum_i (Y_i - (\gamma + \delta)'X_i)^2 + \lambda_1 \sum_j \gamma_j^2 + \lambda_2 \sum_j |\delta_j|$$

Another useful modification of Lasso and Ridge is the lava method. In the lava method, we estimate the coefficients by the penalized least squares as shown in this form. If we look at the formula carefully, we see that, just like previously, we are minimizing the sum of squared prediction errors from predicting outcome observations Y_i with a linear and X_i prediction rule plus a penalty term. However, unlike previously, we are splitting the parameter components into γ_j and δ_j , and penalize γ_j like in Ridge and δ_j like in Lasso. There are two corresponding penalty levels, lambda 1 and lambda 2, which can be chosen by cross-validation in practice.

Now, because of the splitting, lava penalizes coefficients least aggressively compared to Ridge, Lasso, or elastic net, because it penalizes large coefficients like in Lasso and small coefficients like in Ridge. Lava never sets estimated coefficients to zero, and so it doesn't do variable selection. The lava method works really well in sparse + dense models, where there are several large coefficients and many small coefficients, which are not necessarily sparse. In these types of models, lava significantly outperforms Lasso, Ridge or elastic net.

COMPARED TO ELASTIC NET, LAVA PENALIZES LARGE AND SMALL COEFFICIENTS MUCH LESS AGGRESSIVELY, LARGE COEFFICIENTS LIKE IN LASO AND SMALL COEFFICIENTS LIKE IN RIDGE

LIKE RIDGE, LAVA DOES NOT DO VARIABLE SELECTION

LAVA WORKS REALLY WELL IN "SPARSE + DENSE" REGRESSION MODELS

CROSS-VALIDATION

CROSS-VALIDATION IS AN IMPORTANT AND COMMON PRACTICAL TOOL THAT PROVIDES A WAY TO CHOOSE TUNING PARAMETERS SUCH AS THE PENALTY LEVELS.

THE IDEA OF CROSS-VALIDATION IS TO RELY ON REPEATED SPLITTING OF THE TRAINING DATA TO ESTIMATE THE POTENTIAL OUT-OF-SAMPLE PREDICTIVE PERFORMANCE.

In step 1, we divide the data into K blocks called folds. For example was, K equal to five, we split the data into five parts.

In step 2, we begin by leaving one block out. We fit the prediction rule on all other blocks, we then predict outcome observations in the left out block and record the empirical mean squared prediction error. We repeat this procedure for each block.

In step 3, we average the empirical mean squared prediction errors over blocks. We do these steps for several or many values of the tuning parameters, and we choose the best tuning parameter to minimize the average mean squared prediction error.

CROSS-VALIDATION

1. WE PARTITION THE DATA INTO K BLOCKS CALLED "FOLDS", FOR EXAMPLE, WITH $K = 5$, WE SPLIT THE DATA INTO 5 NON-OVERLAPPING BLOCKS
2. LEAVE ONE BLOCK OUT. FIT A PREDICTION RULE ON ALL THE OTHER BLOCKS. PREDICT THE OUTCOME OBSERVATIONS IN THE LEFT OUT BLOCK, AND RECORD THE EMPIRICAL MEAN SQUARED PREDICTION ERROR. REPEAT THIS FOR EACH BLOCK

CROSS-VALIDATION

3. AVERAGE THE EMPIRICAL MEAN SQUARED PREDICTION ERRORS OVER BLOCKS

Let us now consider more formal description of cross-validation. We randomly select equal sized blocks B_1 through B_k to randomly partition the observation indices one through M .

$$B_1, \dots, B_K$$

$$\hat{f}_{-k}(X; \theta)$$

We then fit a prediction rule denoted by \hat{f}_{-k} or X and θ . Where θ denotes tuning parameters such as penalty levels, and \hat{f}_{-k} depends only on observations that are not in the block B_k .

The empirical mean squared error for the block of observations B_k is given by the average squared prediction error for this block, as shown in this formula. In this formula M is the size of the block.

$$\text{MSE}_k(\theta) = \frac{1}{m} \sum_{i \in B_k} (Y_i - \hat{f}_{-k}(X_i; \theta))^2$$

The cross validated MSE is the average of MSEs for each block as shown again, in this formula.

$$\text{CV-MSE}(\theta) = \frac{1}{K} \sum_{k=1}^K \text{MSE}_k(\theta)$$

Finally, the best tuning parameter θ is chosen by minimizing the cross validated MSE.

UNLIKE LASSO, THE THEORETICAL PROPERTIES OF RIDGE AND THE OTHER PENALIZED PROCEDURES ARE LESS WELL UNDERSTOOD IN THE HIGH-DIMENSIONAL CASE (YET). SO IT IS A GOOD IDEA TO RELY ON DATA SPLITTING TO ASSESS THEIR PREDICTIVE PERFORMANCE

CROSS-VALIDATION IS A GOOD WAY TO CHOOSE PENALTY LEVELS, BUT ITS THEORETICAL PROPERTIES ARE ALSO NOT COMPLETELY UNDERSTOOD IN THE HIGH-DIMENSIONAL CASE (YET)

HOW DO THESE METHODS PERFORM IN PRACTICE?

We now provide some concluding remarks.

First, we note that in contrast to Lasso, the theoretical properties of Ridge and other penalized procedures are less well understood in the high-dimensional case, yet. So it is a good idea to rely on test data to assess their predictive performance.

Second, we note that cross validation is a good way to choose penalty levels, but its theoretical properties are not completely understood in high-dimensional case yet. So again, it is a good idea to rely on task data to assess the predictive performance of cross-validated rules. Finally, we may want to ask a question here.

How do the penalize regression methods work in practice?

In the next part of our module we will assess the predictive performance of these methods in a real example, where we will also compare these methods to modern nonlinear regression methods.

2.3.1 Modern Nonlinear Regression. Trees, Random Forests, and Boosted Trees

We begin a new block of our module devoted modern nonlinear regression. Here we are interested in predicting the outcome y using the raw regressor Z which are k -dimensional.

Recall that the best prediction rule $g(Z)$ is the conditional expectation of y given Z . In this module so far, we have used best linear prediction rules to approximate $g(Z)$ and linear regression or lost regression for estimation. Now we consider nonlinear prediction rules to approximate $g(Z)$, including tree-based methods and neural networks. In this segment we discussed tree-based methods.

MODERN NONLINEAR REGRESSION

WE ARE INTERESTED IN PREDICTING THE OUTCOME Y USING THE RAW REGRESSORS Z , WHICH ARE K -DIMENSIONAL

Recall that the best prediction rule $g(Z)$ is the conditional expectation of Y given Z :

$$g(Z) = E(Y|Z)$$

NOW WE CONSIDER NONLINEAR PREDICTION RULES TO APPROXIMATE $g(Z)$, INCLUDING TREE-BASED METHODS AND NEURAL NETWORKS

REGRESSION TREES

THE IDEA IS TO PARTITION THE REGRESSOR SPACE WHERE Z TAKES VALUES INTO A SET OF REGIONS, AND THEN PROVIDE A PREDICTED VALUE FOR EACH REGION

SUPPOSE WE HAVE A SAMPLE OF N OBSERVATIONS $(Z_i, Y_i)_{i=1}^n$

REGRESSION TREES

SUPPOSE WE HAVE A SAMPLE OF N OBSERVATIONS $(Z_i, Y_i)_{i=1}^n$

Given a partition of the regressor space into M regions R_1, \dots, R_M , which will be determined by data, the regression tree is a prediction rule of the form:

$$\hat{g}(Z) = \sum_{m=1}^M \hat{\beta}_m 1(Z \in R_m)$$

where $1(Z \in R_m)$ is an indicator for the region R_m , and $\hat{\beta}_m$ is the predicted value for the region R_m .

In this segment we discussed tree-based methods. The idea of regression trees is to partition the regressor space where Z takes values into a set of rectangles, and then for each rectangle provide a predicted value. Suppose we have n observations, Z_i , Y_i for i ranging from 1 to n.

Given the partition of the space into M rectangles, R_1 through R_M , which will be determined by data. The regression rule is a prediction of the form. Hat g(Z) is equal to the sum over m from 1 to M of hat beta m times the indicator that Z belongs to the rectangle R_m .

The estimated coefficients are obtained by minimizing the sample MSE, as shown in this formula. From the formula we can conclude that hat Beta m is set equal to average of Y_i with Z_i falling in the rectangle R_m , as shown in the formula that you see.

The rectangles or regions R_m are called nodes and each node has predicted value associated with it.

The predicted values $\hat{\beta} = (\hat{\beta}_1, \dots, \hat{\beta}_M)$ are obtained by minimizing the sample MSE:

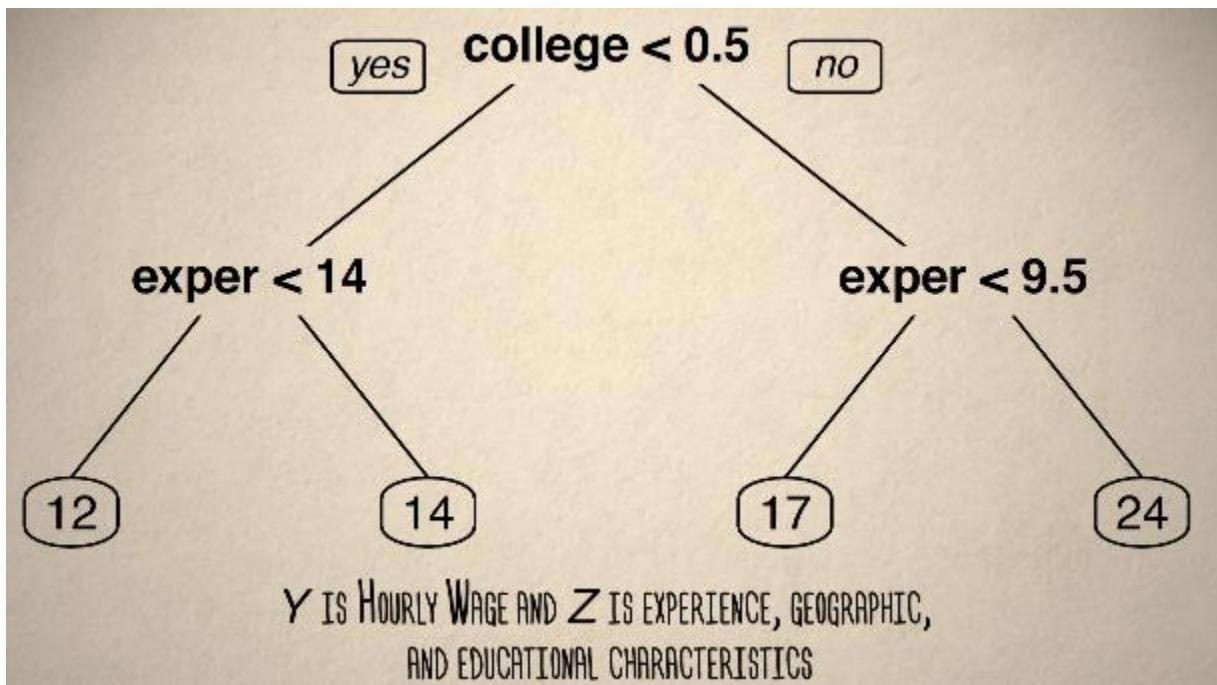
$$\hat{\beta} = \arg \min_{b_1, \dots, b_M} \sum_{i=1}^n \left(Y_i - \sum_{m=1}^M b_m 1(Z_i \in R_m) \right)^2$$

so that

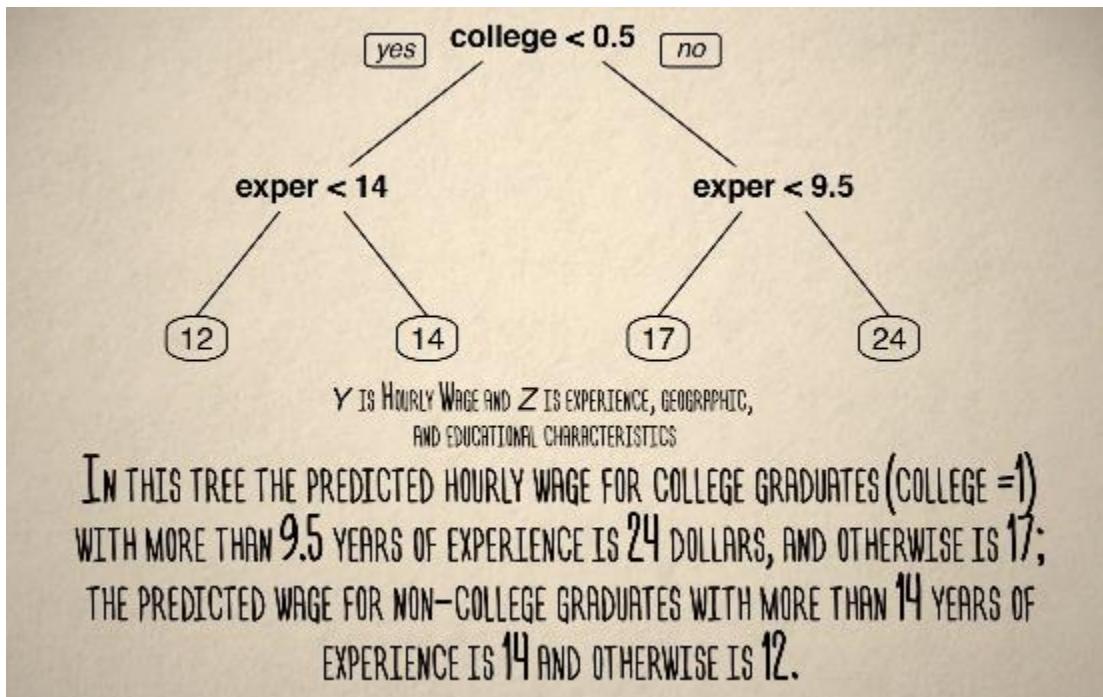
$$\hat{\beta}_m = \text{average of } Y_i \text{ such that } Z_i \in R_m.$$

The regions R_1, \dots, R_M are called nodes, and each node R_m has a predicted value $\hat{\beta}_m$ associated with it.

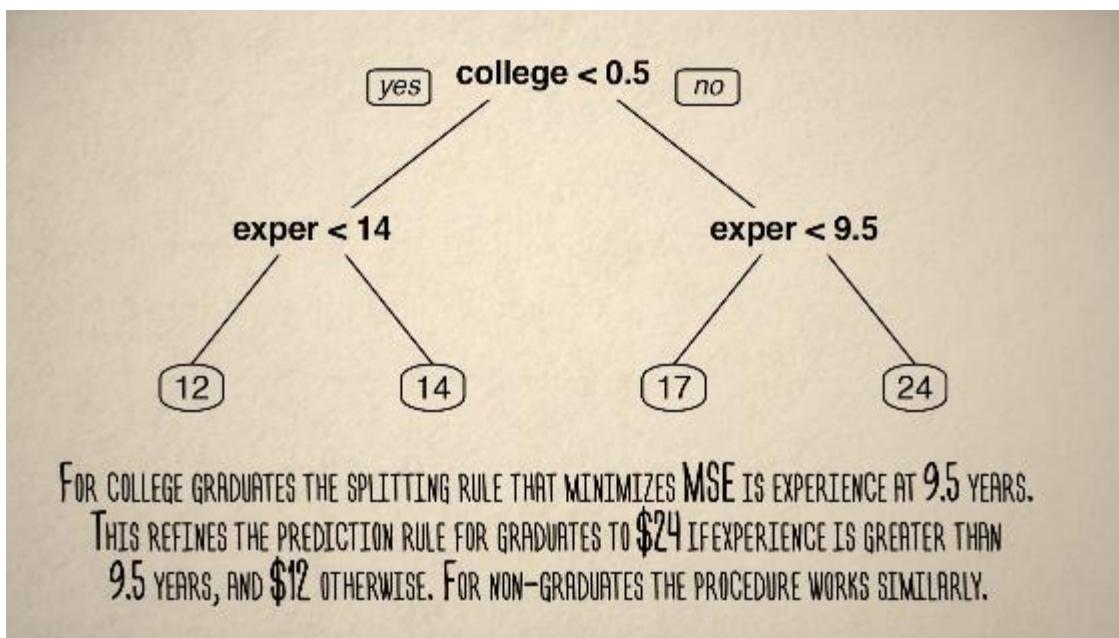
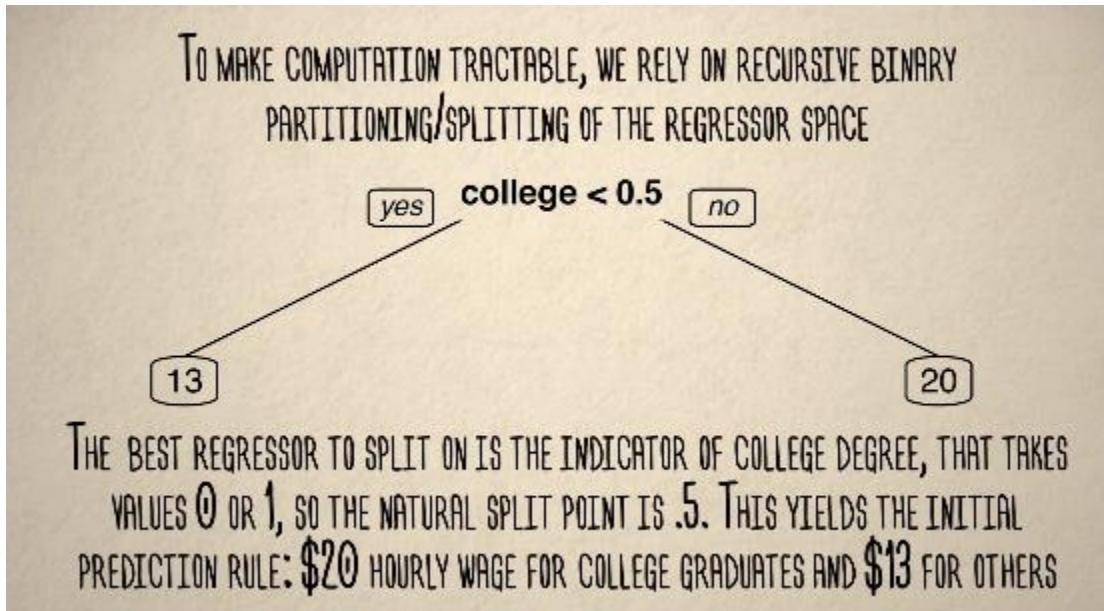
And nice thing about regression trees is that you get to draw cool pictures like this one. Here we show a tree based prediction rule for our wage example where Y is wage and Z are experience, geographic, and educational characteristics.



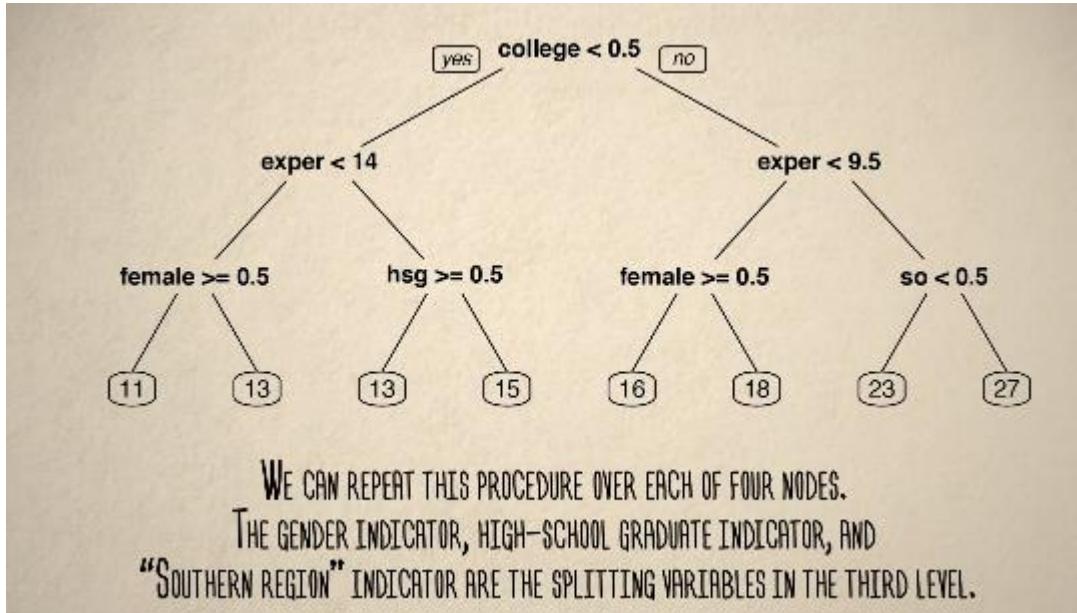
As you can see from looking at the terminal nodes of the tree, in this tree, the predicted hourly wage for college graduates with more than 9.5 years of experience is \$24, and otherwise it is 17. The predicted wage for non-college graduates with more than 14 years of experience is 14, and otherwise it is 12. Now how do we grow this tree?



First, we cut that regressor space into two regions by choosing a regressor and a split point that achieve the best improvement in the sample MSE. This gives us the tree of depth one, that you see in the figure. The best variable to split on here is the indicator of college degree and it takes values of 0 or 1, so the natural split point is 0.5. This gives us a starting tree-based prediction rule, which predicts a \$20 hourly wage for college graduates and 13 for all others.



Second, we repeat this procedure over two resulting regions or nodes, college graduates and non-college graduates, obtaining in this step four new nodes that we see in this figure. For college graduates the splitting rule that minimizes MSE is the experience regressor at 9.5 years. This refines the prediction rule for graduates to \$24 if experience is greater than 9.5 years and \$12 otherwise. For non-graduates the procedure works similarly.



Third, we repeat this procedure over each of the four nodes. The tree of that three now has eight nodes. We see that in the final level we are now splitting our gender indicator, high school graduate indicator, and the south indicator. Finally, we stop growing the tree when the desired depths of the tree is reached. Or when the minimal number of observations per node, called minimal node, size is reached.

The deeper we grow the tree, the better is our approximation to the prediction function $g(Z)$.

On the other hand, the deeper the tree, the more noisy our estimate $\hat{g}(Z)$ becomes, since there are fewer observations per terminal node to estimate the predicted values.

From a prediction point of view we can try to find the right depth or structure of the tree by cross-validation.

For example, in the wage example the tree of depth 2 performs better in terms of cross-validated MSE than the trees of depth 3 or 1.

PRUNING THE TREE

The process of cutting down the branches of the tree to improve predictive performance is called Pruning The Tree.

However in practice pruning the tree often doesn't give satisfactory performance because a single prune tree provides a very crude approximation to the regression function $g(Z)$.

PRUNING THE TREE RARELY GIVES SATISFACTORY PERFORMANCE IN PRACTICE,
BECAUSE A SINGLE PRUNED TREE PROVIDES A VERY CRUDE APPROXIMATION TO
THE REGRESSION FUNCTION $g(Z)$

A much more powerful and one use approach to improve simple regression trees is to build the Random Forest.

Random Forest

The idea of the Random Forest is to grow many different trees and average their prediction rules.

The trees are grown over bootstrap samples: artificial samples randomly drawn from the original data with replacement.

The trees are grown deep to keep the approximation error low, and averaging reduces the noise of the individual trees.

The idea of Random Forest is to grow many different deep trees and then average prediction rules based on them. The trees are grown over different artificial data samples generated by sampling randomly with replacement from the original data. This way of creating artificial samples is called the bootstrap statistics.

The trees are growing deep to keep the approximation error low and averaging is meant to reduce the noisiness of the individual trees.

Each bootstrap sample is created by sampling from our data on pairs (Y_i, Z_i) randomly with replacement, so that some pairs get drawn multiple times and some don't get redrawn.

Given a bootstrap sample, indexed by b , we build a tree-based prediction rule $\hat{g}_b(Z)$.

We repeat the procedure B times, and then average the prediction rules that result from each of the bootstrap samples:

$$\hat{g}_{\text{random forest}}(Z) = \frac{1}{B} \sum_{b=1}^B \hat{g}_b(Z)$$

Each bootstrap sample is created by sampling from our data on pairs (Y_i, Z_i) randomly with replacement, so some observations get drawn multiple times and some don't get redrawn at all. Given a bootstrap sample numbered by numbered by numeral b , we build a tree-based prediction rule $\hat{g}_b(Z)$. We repeat the procedure

capital B times in total and then average the prediction rules, that result from each of the bootstrap samples.

Using bootstrap here is an unusual, yet an intuitive idea: if we would have many independent copies of the data, we could average the prediction rules obtained over these copies to reduce the noise. Since we don't have such copies, we rely on bootstrap to create many quasi-copies of the data.

THE KEY IDEA IS THAT THE TREES ARE GROWN DEEP TO KEEP THE APPROXIMATION ERROR LOW, AND AVERAGING REDUCES THE NOISE OF THE INDIVIDUAL TREES.

THE PROCEDURE OF AVERAGING NOISY PREDICTION RULES OVER THE BOOTSTRAP SAMPLES IS CALLED BOOTSTRAP AGGREGATION OR BAGGING

Using bootstrap here is an intuitive idea. If we could have many independent copies of the data, we could average the prediction rules obtained over these copies to reduce the noisiness. Since we don't have such copies, we rely on bootstrap copies of the data. The key underlying idea here is that the trees are grown deep to keep the approximation error low and averaging is meant to reduce the noisiness of the individual trees. The procedure of averaging noisy prediction rules over the bootstrap samples is called Bootstrap Aggregation or Bagging.

BOOSTING

WE ESTIMATE A SHALLOW TREE-BASED PREDICTION RULE, THEN TAKE THE RESIDUALS AND ESTIMATE ANOTHER SHALLOW TREE-BASED PREDICTION RULE FOR THESE RESIDUALS

THE SUM OF THE PREDICTION RULES FOR THE RESIDUALS IS THE PREDICTION RULE FOR OUTCOME

WE USE SHALLOW TREES TO KEEP NOISE LOW, AND EACH STEP REDUCES THE APPROXIMATION ERROR.

Another approach to improve simple regression tree is by boosting. The idea of boosting is that of request of fatigue where estimated tree-based prediction rule, then we take the residuals and estimate another shallow tree-based prediction rule for these residuals and so on. Summing up the prediction rules for the residuals gives us the prediction rule for the outcome.

Unlike in a random forest, we use shallow trees, rather than deep trees, to keep the noise low and each step of boosting aims to reduce the approximation error. In order to avoid overfitting and boosting, we can stop the procedure once we don't improve the cross-validated mean square error.

BOOSTING

IN ORDER TO AVOID OVERFITTING, WE CAN STOP THE PROCEDURE ONCE WE DON'T IMPROVE THE CROSS-VALIDATED MSE

In order to avoid overfitting and boosting, we can stop the procedure, once we don't improve the cross-validated mean square error.

Step 1. Initialize the residuals $R_i = Y_i, i = 1, \dots, n$

Step 2. For $j = 1, \dots, J$

1. Fit a tree-based prediction rule $\hat{g}_j(Z)$ to the data $(Z_i, R_i)_{i=1}^n$;
2. Update the residuals $R_i \leftarrow R_i - \lambda \hat{g}_j(Z_i)$

Step 3. Construct the boosted prediction rule:

$$\hat{g}(Z) = \sum_{j=1}^J \lambda \hat{g}_j(Z).$$

The tuning parameters J and λ can be chosen by cross-validation.

Formerly, the boosting algorithm looks as follows.

In step 1, we initialize the residuals $R_i = Y_i$, for i that runs from 1 to n .

In step 2, we perform the following operation over index j running from 1 to capital J .

We fit a tree-based prediction rule, $g_j(Z)$ to the data (Z_i, R_i) with i from 1 to n and we have data residuals as R'_i equals previous R_i minus lambda times $\hat{g}_j(Z_i)$. Finally, in step 3, we output the boosted prediction rule. $\hat{g}(Z)$ equals sum over J of lambda $\hat{g}_j(Z)$. The capital J and lambda that you see here are the tuning parameters. Representing the number of boosting steps, and the degree of updating the residuals. In particular, lambda equals 1, gives us the full update. We can choose j and lambda by cross validation.

SUMMARY

WE HAVE DISCUSSED THE TREE-BASED PREDICTION RULES AND WAYS TO IMPROVE THEM BY BAGGING AND BOOSTING

BAGGING IS BOOTSTRAP AGGREGATION OF THE PREDICTION RULES. BOOTSTRAP AGGREGATION OF DEEP REGRESSION TREES GIVES RANDOM FORESTS

BOOSTING USES RECURSIVE FITTING OF RESIDUALS BY A SEQUENCE OF TREE-BASED PREDICTION MODELS. THE SUM OF THESE PREDICTION RULES GIVES THE PREDICTION FOR OUTCOME

So, let us summarize. We discussed the tree-based prediction rules, and ways to improve them by bagging or boosting. Bagging is bootstrap aggregation of the prediction rules. Bootstrap aggregation of deep regression trees gives us random forests. Boosting uses recursive fitting of residuals by a sequence of tree-based prediction models. The sum of these prediction rules gives us the prediction for outcome.

2.3.2 Modern Nonlinear Regression. Neural Networks

The idea of the neural network is to use parameterized nonlinear transformations of linear combinations of the raw regressors as constructed regressors, called neurons. And produce the predicted value as a linear function of these regressors.

WE ARE INTERESTED IN PREDICTING THE OUTCOME Y USING THE RAW REGRESSORS Z , WHICH ARE K -DIMENSIONAL. RECALL THAT THE BEST PREDICTION RULE IS THE CONDITIONAL EXPECTATION OF Y GIVEN Z

$$g(Z) = E(Y|Z)$$

IDEA: USE PARAMETERIZED NONLINEAR TRANSFORMATIONS OF LINEAR COMBINATIONS OF THE RAW REGRESSORS AS CONSTRUCTED REGRESSORS (CALLED NEURONS), AND PRODUCE THE PREDICTED VALUE AS A LINEAR FUNCTION OF THESE REGRESSORS

The prediction rule $\hat{g}(Z)$ is nonlinear in some parameters and with respect to the raw regressors. With sufficiently many neurons, $\hat{g}(Z)$ can approximate the best prediction rule $g(Z)$.

In part 2 of our module, we already saw that many constructed regressors are useful in the high-dimensional linear setting to approximate $g(Z)$.

Neural networks also rely on many constructed regressors to approximate $g(Z)$.

The method and the name "neural networks" were loosely inspired by the mode of operation of the human brain, and developed by scientists working on Artificial Intelligence.

They can be represented by cool graphs and diagrams that we will discuss shortly, so please stay tuned.

Here we focus on single layer neural network to discuss the idea.

The estimated prediction rule will take the form:

$$\hat{g}(Z) = \sum_{m=1}^M \hat{\beta}_m X_m(\hat{\alpha}_m)$$

where the $X_m(\hat{\alpha}_m)$'s are constructed regressors called neurons

The M neurons are generated by

$$X_m(\alpha_m) = \sigma(\alpha'_m Z), \quad m = 1, \dots, M,$$

where α_m 's are neuron-specific vectors of parameters called weights, and σ is the activation function, for example:

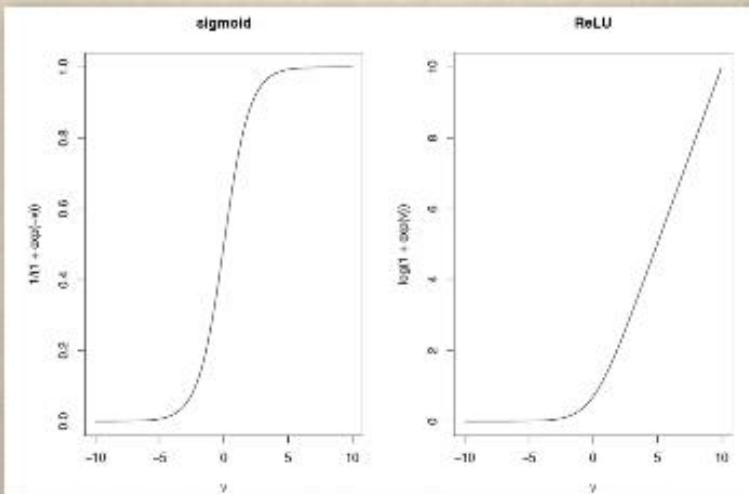
the sigmoid function:

$$\sigma(v) = \frac{1}{1 + e^{-v}}$$

or the rectified linear unit function (ReLU):

$$\sigma(v) = \log(1 + \exp(v))$$

Activation Functions



The figure shows the two graphs of the two Activation Functions. The sigmoid function and the rectified linear unit function. The horizontal axis shows the value of the argument. And the vertical axis the value of the function.

THE ESTIMATORS $\hat{\alpha}_m$ AND $\hat{\beta}_m$ FOR $M = 1, \dots, M$, ARE OBTAINED AS THE SOLUTION TO THE PENALIZED NONLINEAR LEAST SQUARES PROBLEM

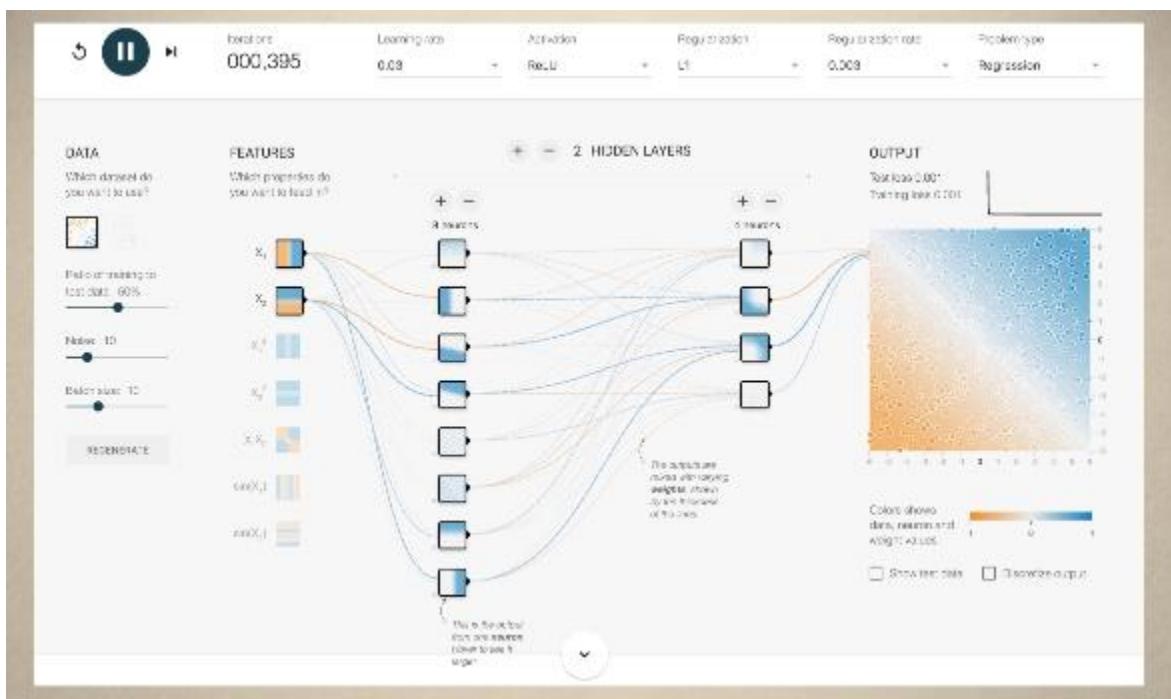
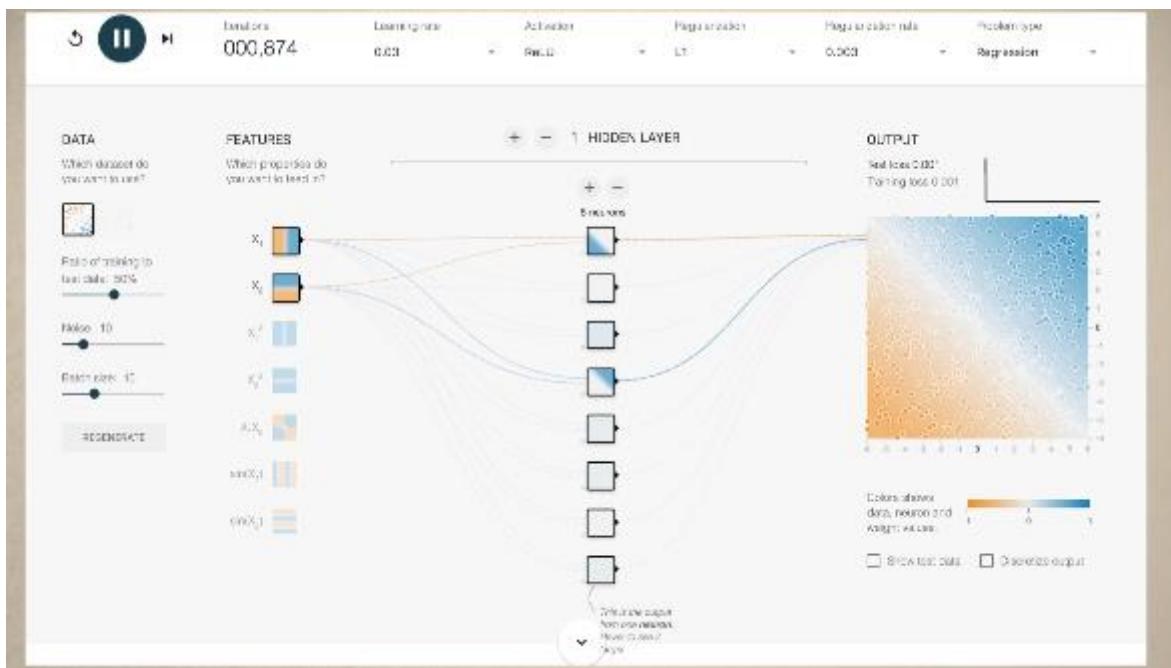
$$\min_{\{\alpha_m\}, \{\beta_m\}} \sum_i \left(Y_i - \sum_{m=1}^M \beta'_m X_{im}(\alpha_m) \right)^2 + \lambda \left(\sum_m \sum_j |\alpha_{mj}| + \sum_m |\beta_m| \right)$$

The estimators had alpha m, and beta m for each M are obtained as the solution to the penalized nonlinear least squares problem shown by this formula. Here, we are minimizing the sum of squared prediction errors in the sample, plus a penalty term given by the sum of the absolute values of components of alpha m and beta m, multiplied by the penalty level, lambda. In this formula, we use the lasso type penalty but we can also use the ridge, another type of penalties.

THE ESTIMATES ARE COMPUTED USING SOPHISTICATED GRADIENT DESCENT ALGORITHMS, WHERE SOPHISTICATION IS NEEDED BECAUSE THIS OPTIMIZATION IS GENERALLY NOT A CONVEX PROBLEM, MAKING THE COMPUTATION A DIFFICULT TASK.

THE PROCEDURE HAS MANY TUNING PARAMETERS, AND IN PRACTICE WE CAN CHOOSE THEM BY CROSS-VALIDATION. THE MOST IMPORTANT CHOICES CONCERN THE NUMBER OF NEURONS AND THE NUMBER OF NEURON LAYERS.

HAVING MORE NEURONS GIVES US FLEXIBILITY, JUST LIKE HAVING MORE CONSTRUCTED REGRESSORS WAS GIVING US MORE FLEXIBILITY WITH HIGH-DIMENSIONAL LINEAR MODELS. TO PREVENT OVERFITTING WE CAN RELY ON PENALIZATION AS IN THE CASE OF LINEAR REGRESSION.



In order to visualize working of the neural network, we rely on the resource called [Playground.TensorFlow.org](https://playground.tensorflow.org) using which we produce a prediction model given by simple single layer neural network model. We now see the graphical representation of this network.

Here we have a regression problem and the network depicts the process of taking row regressors and transforming them into predicted values. In the second column on the left, we see the inputs are features. These features are our two row regressors. The third column shows eight neurons. The neurons are constructed as linear combinations of the row regressors transformed by an activation function. That is, the neurons are along linear transformations of the row regressors. Here, we set the activation function, to be the rectified linear unit function, RELU.

The neurons are connected to the inputs and the connections represent the coefficients $\hat{\alpha} M$, which are coefficients of the neuron specific linear transformations of raw regressors. The coloring represents the sign or the coefficients, orange negative and the blue positive.

And with all the connections represents the size of the coefficients. The neurons are then linearly combined to produce the output, the prediction rule. In the diagram we see the connections going outwards from the neurons to the output. These connections represent the coefficients $\hat{\beta} M$, or the linear combination of the neurons that produce the final output. The coloring and the widths represent the sign and the size of these coefficients.

In the diagram, the prediction rule is shown by the heatmap in the box on the right. On the horizontal and vertical axis, we see the values of the two inputs. The color and its intensity in the heatmap represent the predicted value.

We also see on the top that the penalty function is L1, which stands for the lasso type penalty. Another option is to use L2, which stands for the rich type penalty. The penalty level is called here the regularization rate. In this example we are using a single layer neural network.

If we add one or two additional layers of neurons constructed from the previous layer of neurons, we get a different network, which we show in the following diagram.

Prediction methods based on neural networks with several layers of neurons called the deep learning methods.

This diagrams showcases that one of the major benefits of doing prediction with neural networks is that you can adapt with pretty cool artwork.

Let us summarize, in this segment we have discussed neural networks that have recently gathered much attention under the umbrella of deep learning. Neural

networks represent a powerful and all-purpose method for prediction and regression analysis. Using many neurons and multiple layers gives rise to networks that are very flexible and can approximate the best prediction rules quite well.

IN THIS SEGMENT, WE HAVE DISCUSSED NEURAL NETWORKS, THAT HAVE RECENTLY GATHERED MUCH ATTENTION UNDER THE UMBRELLA OF DEEP LEARNING.

NEURAL NETWORKS RECENTLY EMERGED AS A POWERFUL AND ALL-PURPOSE METHOD FOR PREDICTION AND REGRESSION ANALYSIS.

USING MANY NEURONS AND MULTIPLE LAYERS GIVES RISE TO DEEPER NETWORKS THAT ARE VERY FLEXIBLE AND CAN APPROXIMATE THE BEST PREDICTION RULES QUITE WELL.

Aggregation of Predictors

THE MEAN SQUARED APPROXIMATION ERROR CAN BE SMALL ONCE THE SAMPLE SIZE n IS SUFFICIENTLY LARGE, NAMELY

$$E_Z(\hat{g}(Z) - g(Z))^2 \rightarrow 0, \quad \text{as } n \rightarrow \infty$$

where E_Z denotes the expectation taken over Z , holding everything else fixed.

UNDER THESE CONDITIONS WE EXPECT THAT THE SAMPLE MSE AND R^2 WOULD AGREE WITH THE OUT-OF-SAMPLE MSE AND R^2

DATA SPLITTING

1. We use a random part of data for estimating/training the prediction rule,
2. We use the other part to evaluate the quality of the prediction rule, recording out-of-sample mean squared error (can also look at R^2).

n OBSERVATIONS FOR TRAINING

m FOR TESTING / VALIDATION

V DENOTE THE INDICES OF THE OBSERVATIONS IN THE TEST SAMPLE