

Computational Data Analytics for Economists

Lecture 1

Introduction to Machine Learning for Prediction

Helge Liebert Anthony Strittmatter

Outline

- 1 Prediction vs. Causal Inference
- 2 General Machine Learning Framework
- 3 Regularised Regressions
- 4 Trees and Random Forests
- 5 Other Machine Learning Methods

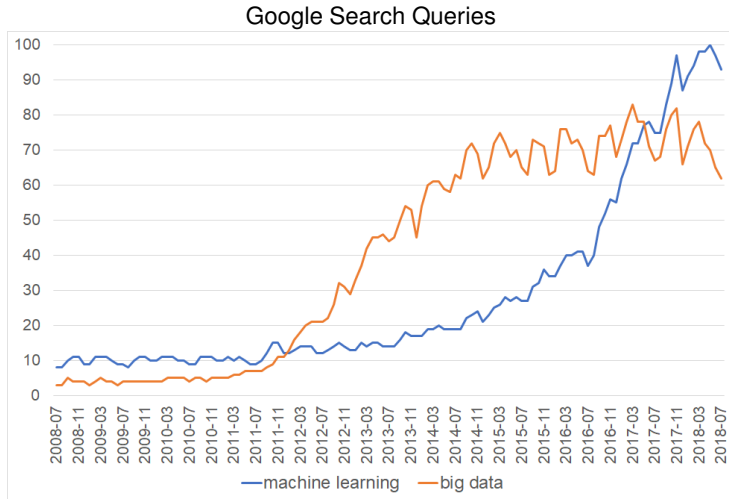
Literature

- Mullainathan and Spiess (2017): "Machine Learning: An Applied Econometric Approach", Journal of Economic Perspectives, 31 (2), pp. 87-106, [download](#).
- Hastie, Tibshirani, and Friedman (2009): "Elements of Statistical Learning", 2nd ed., Springer, [download](#).
- James, Witten, Hastie, and Tibshirani (2013): "An Introduction to Statistical Learning", Springer, [download](#).
- Kleinberg, Lakkaraju, Leskovec, Ludwig, and Mullainathan (2018): "Human Decisions and Machine Predictions", Quarterly Journal of Economics, 133 (1), pp. 237-293, [download](#).

Why Machine Learning (ML)?

- ML (or statistical learning) methods exist already since decades
- Currently "Machine Learning" is a buzz word with no clear definition
- Probably most people think of ML as some computational intensive methods that make data-driven modelling decisions
- Consider the structural model $Y = f(X, U)$
 - ML methods can be very powerful to predict $\hat{Y} = \hat{f}(X)$
 - ML methods are able to incorporate many (or even high-dimensional) covariates X in a convenient way
 - ML methods can make data-driven modelling decisions on $\hat{f}(\cdot)$
 - **Main disadvantage:** ML is a black-box approach and we loose the interpretability of $\hat{f}(X)$

Popularity of Machine Learning



Source: Google trends

Examples: Studies Using ML for Prediction

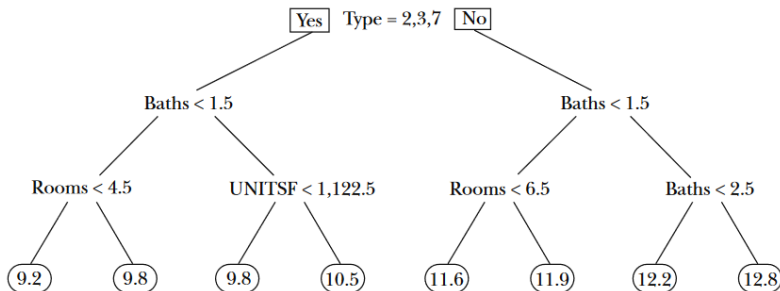
- [Glaeser, Kominers, Luca, and Naik \(2016\)](#) use images from Google Street View to measure block-level income in New York City and Boston
- [Jean et al. \(2016\)](#) train a neural net to predict local economic outcomes from satellite data in African countries
- [Chandler, Levitt, and List \(2011\)](#) predict shootings among high-risk youth so that mentoring interventions can be appropriately targeted
- [Kleinberg, Lakkaraju, Leskovec, Ludwig, and Mullainathan \(2018\)](#) predict the crime probability of defendants released from investigative custody to improve judge decisions
- [Kang, Kuznetsova, Luca, and Choi \(2013\)](#) use restaurant reviews on Yelp.com to predict the outcome of hygiene inspections
- [Huber and Imhof \(2018\)](#) use machine learning to detect bid-rigging cartels in Switzerland
- [Kogan, Levin, Routledge, Sagi, and Smith \(2009\)](#) predict volatility of firms from market-risk disclosure texts (annual 10-K forms)

LASSO Example: Predicting House Prices

	OLS	LASSO	Post-LASSO
Intercept	11.59	11.30	11.26
Baths	0.288	0.251	0.289
Bedrooms	0.011		
Kitchen	-0.270		
Living rooms	-0.126		
Lot size	-0.00000002		
Rooms	0.131	0.108	0.124
Garage	-0.540	-0.377	-0.549

Note: This table just serves for illustration. The sample contains 51,808 houses from American Housing Survey. House prices are measured in log dollars.

Shallow Tree Example: Predicting House Prices



Note: Based on a sample from the 2011 American Housing Survey metropolitan survey. House-value predictions are in log dollars.

Source: Mullainathan and Spiess (2017)

Predictions vs. Causal Inference

- Outcome (e.g., house price): Y
 - Binary Treatment (e.g., incinerator in neighbourhood): $D \in \{0, 1\}$
 - Potential Outcome:
 - $Y(1)$ house price when an incinerator is in the neighbourhood
 - $Y(0)$ house price when no incinerator is in the neighbourhood
 - Only one potential outcome can be observed
 - Causal effect: $\delta = Y(1) - Y(0)$
- Predictions have the observable estimation target \hat{Y}
- Causal inference has the (partly) unobservable estimation target $\hat{\delta}$

Training of ML Algorithms

Out-of-Sample Mean-Squared-Error (MSE):

- Observed outcome Y :

$$MSE_{\hat{Y}} = E[(\hat{Y} - Y)^2] = \underbrace{E[(\hat{Y} - E[\hat{Y}])^2]}_{\text{Variance}} + \underbrace{(E[\hat{Y} - Y])^2}_{\text{Squared Bias}}$$

- Causal Effect δ :

$$MSE_{\hat{\delta}} = E[(\hat{\delta} - \delta)^2] = \underbrace{E[(\hat{\delta} - E[\hat{\delta}])^2]}_{\text{Variance}} + \underbrace{(E[\hat{\delta} - \delta])^2}_{\text{Squared Bias}}$$

→ δ is unobservable

Causal Machine Learning

Three possible Causal Machine Learning approaches:

- Some causal inference methods involve prediction problems (e.g., first stage IV, propensity score). Use ML only for the prediction steps.
- Transform the outcome (or regression weight) such that the conditional expectations roughly approximate the causal effect (e.g., modified outcome method).
- Adapt the optimization criteria of the ML algorithm (e.g., causal trees, causal forests).

When can Predictions be Useful?

- Tasks that have the purpose to predict an observable Y
 - stock or commodity prices
 - unemployment rate or GDP
 - pre-screening of job applications
 - consumer demand (shipping before the order occurs)
 - movie recommendations
 - handwriting, face, or voice recognition
- Tasks that have the purpose to estimate δ , but require an intermediate prediction step
 - first stage of an instrumental variable regression
 - effect heterogeneity of causal effects
 - modified outcome method
 - event study approach
 - sample selection model (e.g., Heckman equation)

What we NOT Learn from Machine Learning

Linear model: $Y = X\beta + U$

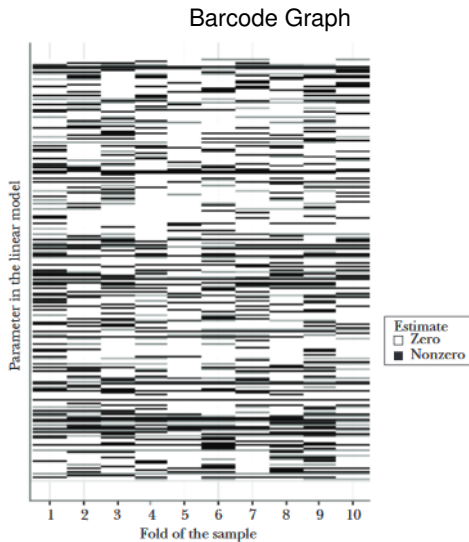
- $\hat{\beta}$ are inconsistently estimated in standard ML approaches (this could be the target of causal machine learning for low-dimensional β)

Non-parametric model: $Y = f(X, U)$

- We will not be able to identify the "true" structural function $f(X)$ without making very restrictive assumptions (see discussion of oracle properties below)

→ ML is a black-box approach

Selected Covariates Across 10 LASSO Regressions



Source: Mullainathan and Spiess (2017)

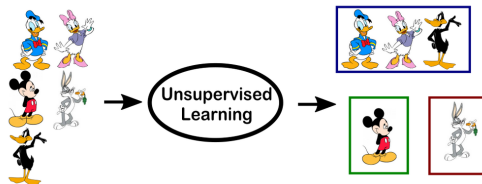
Supervised vs. Unsupervised Machine Learning

Supervised Machine Learning:

- We observe data on Y and X and want to learn the mapping $\hat{Y} = \hat{f}(X)$
- Classification when \hat{Y} is discrete, regression when \hat{Y} is continuous

Unsupervised Machine Learning:

- We observe only data on X and want to learn something about its structure
- Clustering: Partition data into homogeneous groups based on X



- Dimensionality reduction (e.g., principal component analysis)

Supervised Machine Learning in One Expression

Objective function:

$$\min \sum \underbrace{L(Y, f(X))}_{\text{loss function}} \quad \text{over} \quad \underbrace{f \in F}_{\text{function class}} \quad \text{s.t.} \quad \underbrace{R(f) \leq c}_{\text{complexity restriction}}$$

- Loss function $L(\cdot)$: quadratic, absolute, etc.
- Function class F : linear (e.g., LASSO), non-linear (e.g., Logit-LASSO), non-parametric (e.g., tree, neural net)
- Complexity restriction $R(f) \leq c$: depth of tree, number of hidden layers, LASSO penalty, etc.
- Complexity level c : tuning parameter for max. complexity

Source: Francis Diebold's "No Hesitations" [blog](#)

General ML Procedure

- 1 Select a ML method (e.g., LASSO)
- 2 Draw randomly a hold-out-sample from the data
- 3 Estimate the machine learning model using different complexity levels c
- 4 Select the optimal complexity level c^*
- 5 Predict \hat{Y} using c^* and extrapolate the fitted values to the retarded hold-out-sample
- 6 Evaluate the prediction power of the ML method in the hold-out-sample

Selection of the Optimal Complexity

- [Van der Vaart, Dudoit, and van der Laan \(2006\)](#) show that cross-validation (CV) tuning approximates the optimal complexity
- Concrete implementation of CV:
 - Number of CV-folds: 5-20?
 - Performance measure: MSE, Variance, Gini-index?
 - Ad-hoc extensions: "one standard-error rule" for tuning the LASSO
- Alternatives to CV:
 - Some alternatives (e.g., AIC, BIC) rely on assumptions about the sampling process
 - Rare theoretical guidance applies only to specific methods (see, e.g., [Belloni, Chen, Chernozhukov, and Hansen, 2012](#), for the LASSO)

The Firewall Principle

Why do we use the hold-out-sample to evaluate the prediction power?

- If we try many tuning parameter values, we may end up overfitting even in cross-validation samples
- The cross-validation samples are smaller than the hold-out-sample, which may lead to underestimation of the prediction power
- The cross-validation performance is an aggregation over multiple different prediction functions, which differs from the single prediction function we finally estimate

→ Often we use MSE, RMSE, or R^2 to assess the out-of-sample prediction power

Prediction Performance in Housing Price Example

	In-Sample		Hold-Out-Sample	
	MSE	R^2	MSE	R^2
OLS	0.589	47.30%	0.674	41.70%
Shallow Tree	0.675	39.60%	0.758	34.50%
LASSO	0.603	46.00%	0.656	43.30%
Forest	0.167	85.10%	0.631	45.50%
Ensemble	0.219	80.40%	0.626	45.90%

Source: Mullainathan and Spiess (2017)

Honest Inference

- Using the same data to train and estimate the ML model might lead to overfitting
- For example, there is a risk that extreme outliers are grouped in the same leave of a tree which may increase the variance
- To avoid this, we split the data in training and estimation samples
- The training sample is used to select the ML model, e.g., with a CV procedure
- The selected ML model is used for predicting Y in the estimation sample
- The fitted values \hat{Y} are extrapolated to the hold-out-sample for an assessment of the prediction power
- **Potential gain:** Better coverage of confidence intervals
- **Price to pay:** Smaller samples, less precise predictions

Simulation from Athey and Imbens (2016)

Sample Size	Design 1		Design 2		Design 3	
	500	1,000	500	1,000	500	1,000
Coverage of 90% confidence interval						
Standard	0.82	0.85	0.78	0.81	0.69	0.74
Honest	0.90	0.90	0.90	0.89	0.89	0.90
MSE Ratio						
Standard/Honest	1.021		0.754		0.717	

- Results are for tree estimator
- Only designs 2 and 3 have extreme outliers

Cross-Fitting

- To improve the finite sample performance [Chernozhukov et al. \(2017\)](#) propose different cross-fitting procedures
- For example, we could partition the data in two samples A and B
- Train the ML model in each sample separately
- Fit \hat{Y}_A in sample A using the model trained in sample B
- Fit \hat{Y}_B in sample B using the model trained in sample A
- Finally, $\hat{Y} = \frac{1}{2}(\hat{Y}_A + \hat{Y}_B)$

Data transformations

- The way how we encode and transform our data can affect the performance of the ML method
- Transformations of outcome variable Y :
 - ML methods do not guide us regarding the appropriate transformation (e.g., levels, logarithm)
- Transformations of covariates X :
 - Different measures for same variable: We can include many different measures of the same variable (e.g., level, logarithm, squared, cubic) and the ML algorithm will disregard them if not needed
 - Same variable on different aggregation levels: ML methods select an appropriate aggregation level
 - Standardisation: Usually, we standardise all X variables by the means and standard deviations in the training sample

Regularised Regressions

$$\min_{\beta} \left\{ \sum_{i=1}^N \left(Y_i - \beta_0 - \sum_{j=1}^p X_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p p(\beta_j) \right\}$$

where $\lambda \geq 0$ is the tuning parameter and the number of covariates p can be high-dimensional ($p \gg N$)

Choice of penalty function $p(\cdot)$:

- Ridge: $p(\beta_j) = \beta_j^2$
 - LASSO: $p(\beta_j) = |\beta_j|$ (Least Absolute Shrinkage and Selection Operator)
 - Elastic Net: $p(\beta_j) = \alpha |\beta_j| + (1 - \alpha) \beta_j^2$
 - Best Subset Selection: $p(\beta_j) = 1 \{\beta_j \neq 0\}$
- Note that coefficient size depends on the scaling of X_j . It is best practice to standardise X_j .

Summation Notation

- OLS residual sum of squares (RSS):

$$RSS = \sum_{i=1}^N (Y_i - \beta_0 - \sum_{j=1}^p X_{ij}\beta_j)^2$$

- Penalized regression:

- Lagrangian operator

$$\min_{\beta} \{RSS + \lambda \sum_{j=1}^p p(\beta_j)\}$$

- Constrained regression

$$\min_{\beta} \{RSS\} \text{ s.t. } \sum_{j=1}^p p(\beta_j) \leq c$$

Matrix Notation

- **Ridge:**

$$\min_{\beta} \left\{ (Y - X\beta)'(Y - X\beta) + \lambda \|\beta\|_2^2 \right\}$$

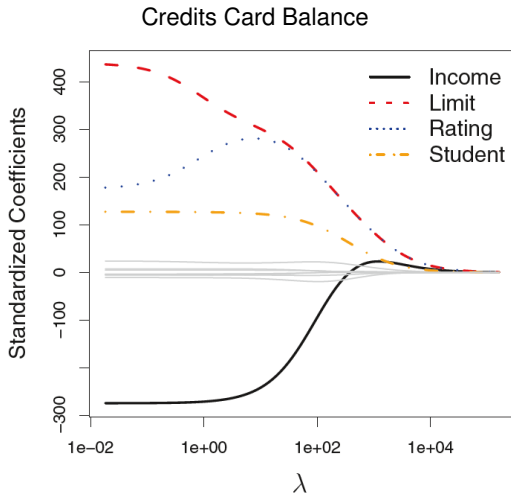
with $\|\beta\|_2^2 = \beta'\beta = \sum_{j=1}^p \beta_j^2$ (squared l_2 -norm), $\hat{\beta} = (X'X + \lambda I)^{-1}X'Y$, and I being the identity matrix

- **LASSO:**

$$\min_{\beta} \left\{ (Y - X\beta)'(Y - X\beta) + \lambda \|\beta\|_1 \right\}$$

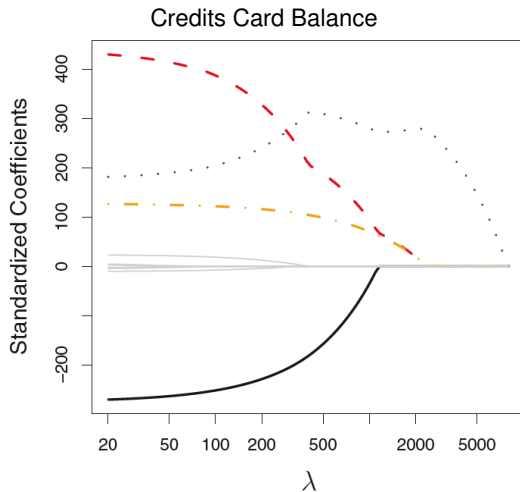
with $\|\beta\|_1 = \sum_{j=1}^p |\beta_j|$ (l_1 -norm)

Ridge Coefficients



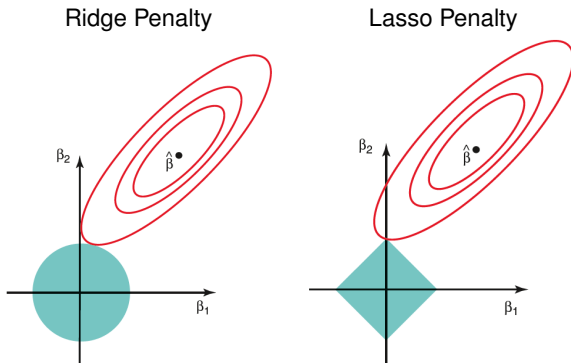
Source: James, Witten, Hastie, Tibshirani (2013)

LASSO Coefficients



Source: James, Witten, Hastie, Tibshirani (2013)

Constraint Regions



Source: James, Witten, Hastie, Tibshirani (2013)

Simple Example

- Consider $X = I$ with dimension $p = N$
- OLS model

$$\min \sum_{j=1}^p (Y_j - \beta_j)^2,$$

such that the estimated OLS coefficients are $\hat{\beta}_j = Y_j$

- Ridge model

$$\min \sum_{j=1}^p (Y_j - \beta_j)^2 + \lambda \sum_{j=1}^p \beta_j^2,$$

such that the estimated Ridge coefficients are $\hat{\beta}_j^R = \hat{\beta}_j / (1 + \lambda)$

Simple Example (cont.)

- LASSO model

$$\min \sum_{j=1}^p (Y_j - \beta_j)^2 + \lambda \sum_{j=1}^p |\beta_j|,$$

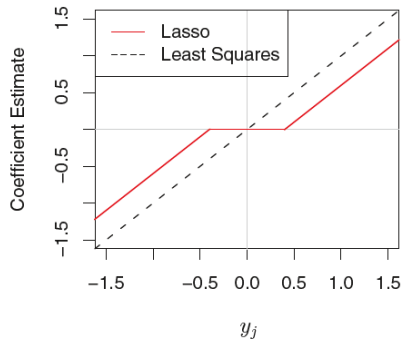
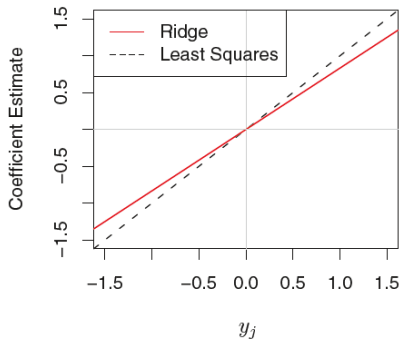
such that the estimated LASSO coefficients are

$$\hat{\beta}_j^L = \begin{cases} \hat{\beta}_j - \lambda/2 & \text{if } \hat{\beta}_j > \lambda/2, \\ \hat{\beta}_j + \lambda/2 & \text{if } \hat{\beta}_j < -\lambda/2, \\ 0 & \text{if } |\hat{\beta}_j| < \lambda/2, \end{cases}$$

which corresponds to the soft-thresholding operator

$$\hat{\beta}_j^L = \text{sign}(\hat{\beta}_j)(|\hat{\beta}_j| - \lambda/2)_+$$

Simple Example (cont.)



Source: James, Witten, Hastie, Tibshirani (2013)

Coordinate Descent Algorithm for LASSO

$$\min_{\beta} \left\{ \frac{1}{2N} \sum_{i=1}^N (Y_i - \beta_0 - \sum_{j=1}^p X_{ij} \beta_j)^2 + \lambda_s \sum_{j=1}^p |\beta_j| \right\}$$

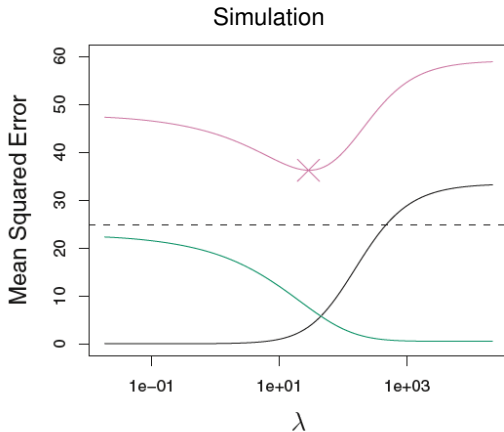
- (1) Specify a grid of $s = 1, \dots, S$ tuning parameters $\lambda_s \in \{\lambda_1, \lambda_2, \dots, \lambda_S\}$
- (2) Take residuals $Y_i^* = Y_i - \frac{1}{N} \sum_{i=1}^N Y_i$ and initialise $\beta_j = 0$
- (3) Circulate repeatedly over all $j = 1, \dots, p$ until convergence:
 - (a) Compute the partial residuals by $r_{ij} = Y_i^* - \sum_{k \neq j} X_{ik} \beta_k$
 - (b) Calculate the simple univariate OLS coefficient $\tilde{\beta}_j = \frac{1}{N} \sum_{i=1}^N X_{ij} r_{ij}$
 - (c) Update β_j with the soft-thresholding operator:

$$\beta_j = \text{sign}(\tilde{\beta}_j) (|\tilde{\beta}_j| - \lambda_s)_+$$

- (4) Repeat (3) for $s = 1, \dots, S$

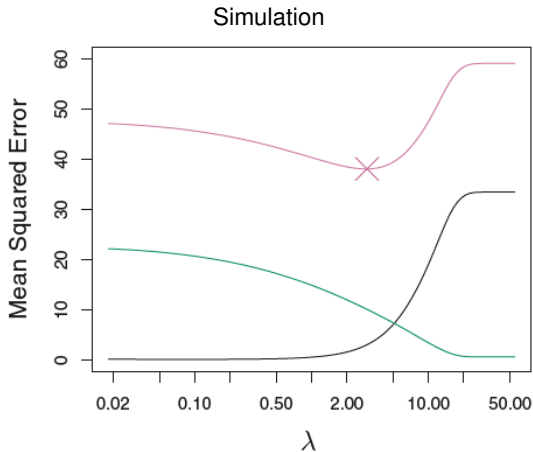
Note: Standardisation of X is required

Ridge: Variance-Bias Trade-Off



Source: James, Witten, Hastie, Tibshirani (2013)

LASSO: Variance-Bias Trade-Off



Source: James, Witten, Hastie, Tibshirani (2013)

Oracle Property

Setup:

- Structural model: $Y_i = \beta_0 + \sum_{j=1}^K X_{ij}\beta_j + U_i$
- The vector X_i has $p > K$ dimensions
- Estimated LASSO model: $\hat{Y}_i = \hat{\beta}_0 + \sum_{j=1}^p X_{ij}\hat{\beta}_j + \lambda \sum_{j=1}^p |\hat{\beta}_j|$

Oracle property:

- Model selection consistency: When $N \rightarrow \infty$,

$$Pr(\hat{\beta}_j = 0) \xrightarrow{p} 1$$

for all $\hat{\beta}_j \in \{\hat{\beta}_{K+1}, \dots, \hat{\beta}_p\}$ (irrelevant covariates)

- Coefficient estimation consistency: When $N \rightarrow \infty$,

$$Pr(|\hat{\beta}_j - \beta_j| > \varepsilon) \xrightarrow{p} 0$$

for any $\varepsilon > 0$ and all $\hat{\beta}_j \in \{\hat{\beta}_0, \dots, \hat{\beta}_K\}$ (relevant covariates)

Minimum assumptions required:

- Sparsity ($K \ll N$) and relevant covariates have to be roughly orthogonal to the irrelevant ones ("irrepresentability condition")

Post-LASSO

- Coefficients of LASSO $\hat{\beta}_j$ are inconsistent when $\lambda_N > 0$ (asymptotically)
- Post-LASSO:

$$\min_{\alpha} \sum_{i=1}^N \left(Y_i - \alpha_0 - \sum_{j=1}^p 1\{\hat{\beta}_j \neq 0\} X_{ij} \alpha_j \right)^2$$

- Coefficients of Post-LASSO are consistent
- Model selection consistency depends on the first-step LASSO estimates
- Alternatives:
 - Adaptive LASSO: $\lambda_j^* = \lambda / |\hat{\beta}_j|^\gamma$ ([Zou, 2006](#))
 - Conservative LASSO: $\lambda_j^* = \lambda / \max(|\hat{\beta}_j|, \lambda)$ ([Caner and Kock, 2018](#))

Logit-Lasso

- Logistic log-likelihood function:

$$L = \sum_{i=1}^N (Y_i X_i \beta - \log(1 + \exp(X_i \beta)))$$

- Logit-LASSO:

$$\min_{\beta} \left\{ -L + \lambda \sum_{j=1}^p |\beta_j| \right\}$$

- Log likelihood function can be approximated by repeated application of weighted LASSO (proximal Newton iterations, [Lee, Sun, and Saunders, 2014](#))
- Can be estimated with a coordinate descend algorithm with weighted soft-thresholding
- If $p \gg N$ we cannot let λ go down to zero
- See [Hastie, Tibshirani and Wainwright \(2016\)](#) for a comprehensive summary of logit-LASSO

Computation Time

Computation time (in seconds) for 100 λ values ($N = 5,000$ and $p = 100$)

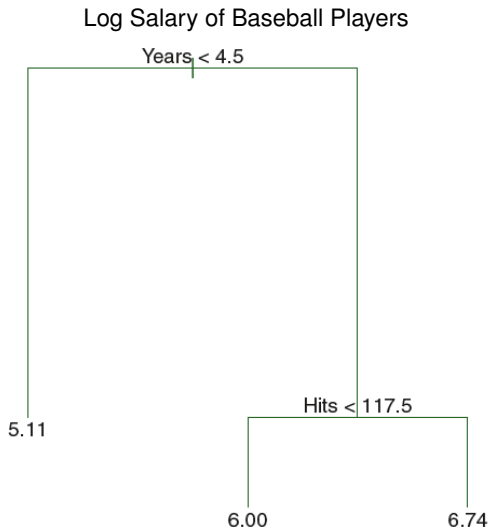
		Average Correlation between Covariates					
		0	0.1	0.2	0.5	0.9	0.95
Coordinate Descent	Linear	0.05	0.05	0.05	0.05	0.05	0.05
	Logistic	7.89	8.48	9.01	13.39	26.68	26.36
Least Angle	Linear	0.29	0.29	0.29	0.3	0.29	0.29

Source: [Hastie \(2009\)](#)

Tree

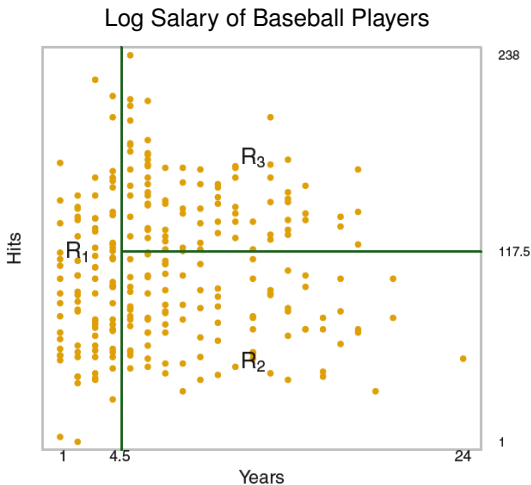
- Trees partition the sample into mutually exclusive groups l_j , which are called leaves
- Let $\pi = \{l_1, \dots, l_{\#(\pi)}\}$ be the terminal leaves of a specific tree or sample partition
- Let $l_j \equiv l_j(x, \pi)$ be the respective leaf (for $j = 1, \dots, \#(\pi)$)
- The leaf $l_j(x, \pi)$ of tree π is a function of the covariates X such that $x \in l_j$
- Let $\#(\pi)$ be the number of terminal leaves in tree π

Example: Shallow Tree



Source: James, Witten, Hastie, Tibshirani (2013)

Example: Shallow Tree (cont.)



Source: James, Witten, Hastie, Tibshirani (2013)

Recursive Partitioning

- Trees select the leaves with a top-down, greedy algorithm, which is called *recursive partitioning*
- *Top-down* because we start with a root (tree without leaves) and successively add splits
- *Greedy* because at each step of the tree building we add the split that improves the prediction power best (instead of looking ahead)

Tree Building Algorithm

- (1) Start with the entire sample (root)
- (2) For each predictor X_j and cut-point s define the pair of half-planes

$$l_1^{(j,s)} = \{X|X_j < s\} \text{ and } l_2^{(j,s)} = \{X|X_j \geq s\}$$

- Calculate the mean outcomes \bar{Y}_1 and \bar{Y}_2 in each half-plane, respectively
- Seek the covariate X_{j1}^* and the cut-point s_1^* that minimise

$$\sum_{i: X_i \in l_1^{(j,s)}} (Y_i - \bar{Y}_1)^2 + \sum_{i: X_i \in l_2^{(j,s)}} (Y_i - \bar{Y}_2)^2$$

Tree Building Algorithm (cont.)

(2) For each predictor X_j and cut-point s define the triple of half-planes

$$l_1^{(j,s)} = \{X|X_{j1}^* < s_1^*, X_j < s\}, l_2^{(j,s)} = \{X|X_{j1}^* < s_1^*, X_j \geq s\}, \text{ and} \\ l_3^{(j,s)} = \{X|X_{j1}^* \geq s_1^*\}$$

and

$$l_1^{(j,s)} = \{X|X_{j1}^* \geq s_1^*, X_j < s\}, l_2^{(j,s)} = \{X|X_{j1}^* \geq s_1^*, X_j \geq s\}, \text{ and} \\ l_3^{(j,s)} = \{X|X_{j1}^* < s_1^*\}$$

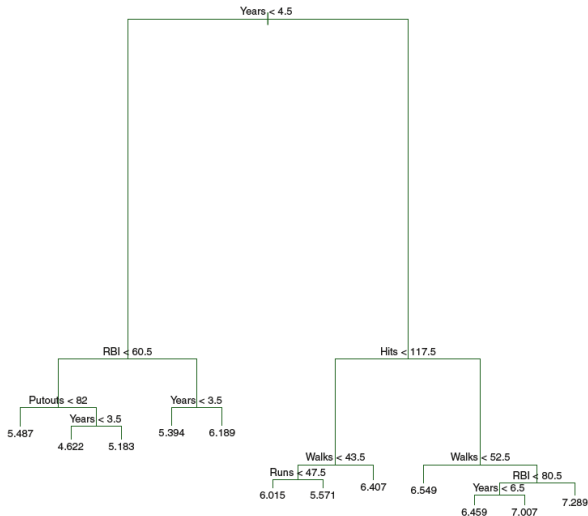
- Calculate the mean outcomes \bar{Y}_1 , \bar{Y}_2 , and \bar{Y}_3 in each half-plane, respectively
- Seek the covariate X_{j2}^* and the cut-point s_2^* that minimise

$$\sum_{i: X_i \in l_1^{(j,s)}} (Y_i - \bar{Y}_1)^2 + \sum_{i: X_i \in l_2^{(j,s)}} (Y_i - \bar{Y}_2)^2 + \sum_{i: X_i \in l_3^{(j,s)}} (Y_i - \bar{Y}_3)^2$$

(3) Continue until some stopping rule is reached (e.g., max. tree size, min. terminal leave size, min. MSE gain)

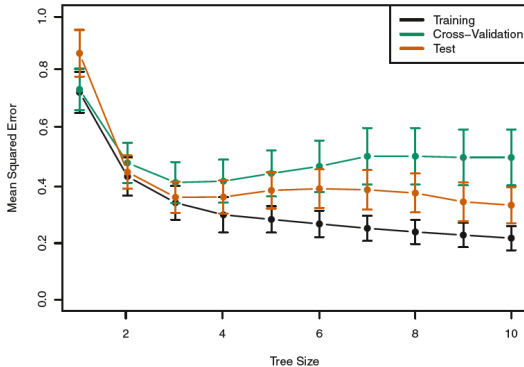
"Deep" Tree

Log Salary of Baseball Players



Source: James, Witten, Hastie, Tibshirani (2013)

Selecting Optimal Tree Size



- Pruning the complexity of trees can improve out-of-sample prediction power
- Select optimal tree size π with cross-validation

Complexity Pruning

- (A) Use recursive partitioning to grow the deep tree π_0 in the training data
- (B) For each value of α a subtree $\pi \subseteq \pi_0$ minimises

$$\sum_{j=1}^{\#(\pi)} \sum_{i: X_i \in I_j} (Y_i - \bar{Y}_j)^2 + \alpha \#(\pi) \quad (1)$$

Obtain a sequence of best subtrees

- (C) Use cross-validation to choose α . Partition the sample in k folds. For each fold:
 - (a) Repeat steps (A) and (B) in the k th-fold
 - (b) Evaluate the MSE using equation (1)
 - (c) Average the MSE across the k folds for each value of α and select the α that minimises the average MSE
- (D) Return to the subtree from (B) with the selected value of α

Prediction

- For the selected tree π^* use the estimation sample to predict

$$\hat{Y}_i = \sum_{j=1}^{\#(\pi^*)} \frac{1}{\sum_{i=1}^N 1\{X_i \in l_j(x, \pi^*)\}} \sum_{i=1}^N 1\{X_i \in l_j(x, \pi^*)\} \cdot Y_i$$

- This is equivalent to the linear regression

$$\min_{\beta} \sum_{i=1}^N \left(Y_i - \sum_{j=1}^{\#(\pi^*)} 1\{X_i \in l_j(x, \pi^*)\} \beta_j \right)^2$$

Advantages and Disadvantages of Trees

Advantages

- Shallow trees are very easy to understand
- Shallow trees can be displayed graphically in a nice way
- Trees automatically handle interactions between covariates
- It is not necessary to transform covariates as long as they have an order

Disadvantages

- Often trees are unstable and have a high variance

Random Forests

- Build G deep trees π_g on different subsets of the data (subsampling or bootstrapping) and covariates
→ decorrelated trees
- These trees are overfitted in the sample and will have a high out-of-sample variance
- To overcome this problem, we aggregate the trees

$$\hat{Y}_i^{RF} = \frac{1}{G} \sum_{g=1}^G \hat{Y}_i^{\pi_g}$$

- This procedure is often called bootstrap-aggregation ("bagging")
- We loose interoperability but gain prediction power compared to (shallow) trees
- Tuning parameters: Forest size, subsample selection, covariate selection, tree size, honest inference, etc

Random Forests: Weighted Representation

- Tree weights:

$$\alpha_{ig}(x) = \frac{1\{X_i \in l_j(x, \pi_g)\}}{\sum_{i=1}^N 1\{X_i \in l_j(x, \pi_g)\}}$$

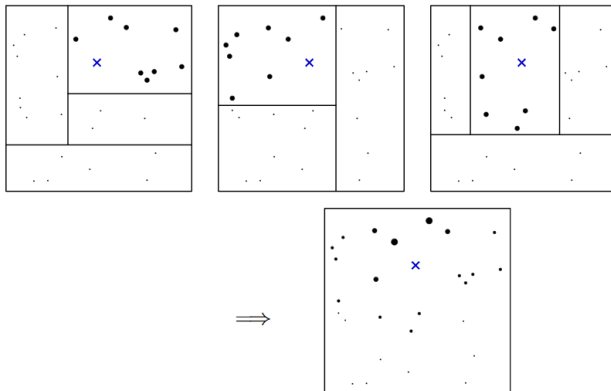
- Forest weights:

$$\alpha_i(x) = \frac{1}{G} \sum_{g=1}^G \alpha_{ig}(x)$$

- Predicted outcome:

$$\hat{\mu}(x) = \sum_{i=1}^N \alpha_i(x) Y_i \text{ and } \hat{Y}_i = \hat{\mu}(X_i)$$

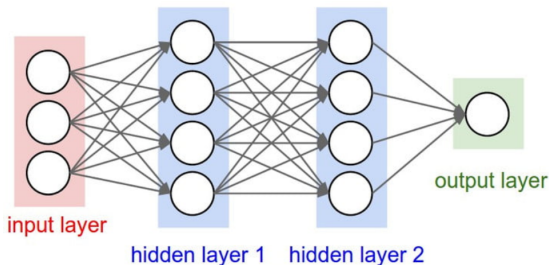
Random Forest: Weighted Representation (cont.)



Source: [Athey, Tibshirani, Wager \(2018\)](#)

Neural Networks

- Input layer, output layer, hidden layer, neurons



- Neurons network function: $f(X) = F(\sum_j w_j g(x_j))$
- Types of neural nets: Feedforward, Recurrent, Autoencoder, etc.
- Econometric references: [Farell, Liang, and Misra \(2018\)](#),
[Hartford, Lewis, Leyton-Brown, and Taddy \(2016\)](#)

Other (Supervised) Machine Learning Methods

- Best Subset Selection, Forward Selection
- Boosting ([Luo and Spindler, 2017](#))
- k-Nearest Neighbours, Kernel, Splines
- Ensemble Methods ([Künzel, Sekhon, Bickel, and Yu, 2018](#))
- etc.