

SuperLearner and LTMLE

```
# Install packages
if (!require("pacman")) install.packages("pacman")

pacman::p_load(# Tidyverse packages including dplyr and ggplot2
  tidyverse,
  ggthemes,
  ltmle,
  tmle,
  SuperLearner,
  tidymodels,
  caret,
  e1071,
  rpart,
  frrrr,
  parallel,
  ranger)

set.seed(44)
```

Introduction

For our final lab, we will be looking at the [SuperLearner](#) library, as well as the [Targeted Maximum Likelihood Estimation \(TMLE\)](#) framework, with an extension to longitudinal data structures. This lab brings together a lot of what we learned about both machine learning and causal inference this year, and serves as an introduction to this intersection!

Data

For this lab, we will use the Boston dataset from the `MASS` library. This dataset is a frequently used toy dataset for machine learning. The main variables we are going to predict is `medv` which is the median home value for a house in Boston.

```
# Load Boston dataset from MASS package
data(Boston, package = "MASS")
glimpse(Boston)

## Rows: 506
## Columns: 14
## $ crim    <dbl> 0.00632, 0.02731, 0.02729, 0.03237, 0.06905, 0.02985, 0.08829, ~
## $ zn      <dbl> 18.0, 0.0, 0.0, 0.0, 0.0, 0.0, 12.5, 12.5, 12.5, 12.5, 12.5, 1~
## $ indus   <dbl> 2.31, 7.07, 7.07, 2.18, 2.18, 2.18, 7.87, 7.87, 7.87, 7.87, 7.~
## $ chas    <int> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ~
## $ nox     <dbl> 0.538, 0.469, 0.469, 0.458, 0.458, 0.458, 0.524, 0.524, 0.524, ~
## $ rm      <dbl> 6.575, 6.421, 7.185, 6.998, 7.147, 6.430, 6.012, 6.172, 5.631, ~
## $ age     <dbl> 65.2, 78.9, 61.1, 45.8, 54.2, 58.7, 66.6, 96.1, 100.0, 85.9, 9~
```

```
## $ dis      <dbl> 4.0900, 4.9671, 4.9671, 6.0622, 6.0622, 6.0622, 5.5605, 5.9505~
## $ rad      <int> 1, 2, 2, 3, 3, 3, 5, 5, 5, 5, 5, 5, 5, 4, 4, 4, 4, 4, 4, 4,~
## $ tax      <dbl> 296, 242, 242, 222, 222, 222, 311, 311, 311, 311, 311, 311, 31~
## $ ptratio  <dbl> 15.3, 17.8, 17.8, 18.7, 18.7, 18.7, 15.2, 15.2, 15.2, 15.2, 15~
## $ black    <dbl> 396.90, 396.90, 392.83, 394.63, 396.90, 394.12, 395.60, 396.90~
## $ lstat    <dbl> 4.98, 9.14, 4.03, 2.94, 5.33, 5.21, 12.43, 19.15, 29.93, 17.10~
## $ medv     <dbl> 24.0, 21.6, 34.7, 33.4, 36.2, 28.7, 22.9, 27.1, 16.5, 18.9, 15~
```

SuperLearner

First, we are going to introduce the SuperLearner package for machine learning in R. SuperLearner was developed here at Berkeley, and we are going to follow [the guide](#) put together by Chris Kennedy. There are some differences, mostly because of the introduction of new machine learning tools like [caret](#) and [tidymodels](#).

The basic idea underlying SuperLearner is that it combines cross-validation with ensemble learning to create a “meta-estimator” that is a weighted combination of constituent algorithms. This idea should sound familiar from when we explored ensemble methods in `sklearn` in Python.

Train/Test Split

Let’s start with our typical train/test split. There are several options for doing this, but we will use the `tidymodels` method. Let’s take a look:

```
# initial_split function from tidymodels/rsample
boston_split <- initial_split(Boston, prop = 3/4)

# Declare the training set with rsample::training()
train <- training(boston_split)
# y_train is medv where medv > 22 is a 1, 0 otherwise
y_train <- train %>%
  mutate(medv = ifelse(medv > 22,
                        1,
                        0)) %>%
  pull(medv)

# x_train is everything but the outcome
x_train <- train %>%
  select(-medv)

# Do the same procedure with the test set
test <- testing(boston_split)

y_test <- test %>%
  mutate(medv = ifelse(medv > 22,
                        1,
                        0)) %>%
  pull(medv)

x_test <- test %>%
  select(-medv)
```

SuperLearner Models

Now that we have our train and test partitions set up, let's see what machine learning models we have available to us. We can use the `listWrappers()` function from SuperLearner:

```
listWrappers()

## [1] "SL.bartMachine"      "SL.bayesglm"      "SL.biglasso"
## [4] "SL.caret"           "SL.caret.rpart"   "SL.cforest"
## [7] "SL.earth"           "SL.extraTrees"    "SL.gam"
## [10] "SL.gbm"             "SL.glm"           "SL.glm.interaction"
## [13] "SL.glmnet"          "SL.ipredbagg"     "SL.kernelKnn"
## [16] "SL.knn"             "SL.ksvm"          "SL.lda"
## [19] "SL.leekasso"        "SL.lm"            "SL.loess"
## [22] "SL.logreg"          "SL.mean"          "SL.nnet"
## [25] "SL.nnls"            "SL.polymars"      "SL.qda"
## [28] "SL.randomForest"    "SL.ranger"        "SL.ridge"
## [31] "SL.rpart"           "SL.rpartPrune"    "SL.speedglm"
## [34] "SL.speedlm"         "SL.step"          "SL.step.forward"
## [37] "SL.step.interaction" "SL.stepAIC"       "SL.svm"
## [40] "SL.template"        "SL.xgboost"
## [1] "All"
## [1] "screen.corP"         "screen.corRank"   "screen.glmnet"
## [4] "screen.randomForest" "screen.SIS"        "screen.template"
## [7] "screen.ttest"        "write.screen.template"
```

Notice how we have both prediction algorithms for supervised learning, and screening algorithms for feature selection (some may be both).

Let's go ahead and fit a model. We'll start with a LASSO which we can call via `glmnet`:

```
sl_lasso <- SuperLearner(Y = y_train,
                        X = x_train,
                        family = binomial(),
                        SL.library = "SL.glmnet")

sl_lasso

##
## Call:
## SuperLearner(Y = y_train, X = x_train, family = binomial(), SL.library = "SL.glmnet")
##
##
##
##              Risk Coef
## SL.glmnet_All 0.09257247 1

Notice how it spits out a "Risk" and a "Coef". The "Coef" here is 1 because this is the only model in our ensemble right now. Risk is essentially a measure of accuracy, in this case something like mean squared error. We can see the model in our ensemble that had the lowest risk like this:

# Here is the risk of the best model (discrete SuperLearner winner).
# Use which.min boolean to find minimum cvRisk in list
sl_lasso$cvRisk[which.min(sl_lasso$cvRisk)]

## SL.glmnet_All
## 0.09257247
```

Multiple Models

Now let's extend this framework to multiple models. Within SuperLearner, all we need to do is add models to the `SL.library` argument:

```
sl = SuperLearner(Y = y_train,
                  X = x_train,
                  family = binomial(),
                  SL.library = c('SL.mean',
                                'SL.glmnet',
                                'SL.ranger'))

sl

##
## Call:
## SuperLearner(Y = y_train, X = x_train, family = binomial(), SL.library = c("SL.mean",
##   "SL.glmnet", "SL.ranger"))
##
##
##              Risk      Coef
## SL.mean_All    0.24840153 0.00000000
## SL.glmnet_All  0.09115271 0.05667652
## SL.ranger_All  0.07536380 0.94332348
```

Now let's move to our validation step. We'll use the `predict()` function to take our SuperLearner model (only keeping models that had weights) and generate predictions. We can then compare our predictions against our true observations.

```
preds <- predict(sl,
                  x_test,
                  onlySL = TRUE)

# start with y_test
validation <- y_test %>%
  # add our predictions
  bind_cols(preds$pred[,1]) %>%
  # rename columns
  rename(obs = `...1`,
         pred = `...2`) %>%
  mutate(pred = ifelse(pred >= .5,
                        1,
                        0))

head(validation)

## # A tibble: 6 x 2
##   obs pred
##   <dbl> <dbl>
## 1     1     1
## 2     1     1
## 3     1     1
## 4     0     0
## 5     0     0
## 6     0     0
```

Notice that the predictions are not binary 1/0s, but rather probabilities. So, we recode these so that if estimate $\geq .5$ it becomes a 1, and 0 otherwise. We can then use our standard classification metrics and

create a confusion matrix using `caret`:

TP: True Positives, predicted a 1 and observed a 1 *FP*: False Positives, predicted a 1 and observed a 0 *TN*: True Negatives, predicted a 0 and observed a 0 *FN*: False Negatives, predicted a 0 and observed a 1

```
caret::confusionMatrix(as.factor(validation$pred),
                        as.factor(validation$obs))
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction  0   1
##           0 64 11
##           1  7 45
##
##           Accuracy : 0.8583
##           95% CI : (0.7853, 0.9138)
##       No Information Rate : 0.5591
##       P-Value [Acc > NIR] : 4.827e-13
##
##           Kappa : 0.7103
##
##  Mcnemar's Test P-Value : 0.4795
##
##           Sensitivity : 0.9014
##           Specificity : 0.8036
##       Pos Pred Value : 0.8533
##       Neg Pred Value : 0.8654
##           Prevalence : 0.5591
##       Detection Rate : 0.5039
##   Detection Prevalence : 0.5906
##       Balanced Accuracy : 0.8525
##
##       'Positive' Class : 0
##
```

Ensemble Learning and Parallel Processing

SuperLearner can take a long time to run, but we can also speed this up with parallelization. Unfortunately this works slightly differently in Windows and Mac/Linux (R doesn't seem to recognize Windows Subsystem for Linux).

```
# Parallel backend Mac/Linux
n_cores <- availableCores() - 1

plan(multiprocess,
     workers = n_cores)
set.seed(44, "L'Ecuyer-CMRG")

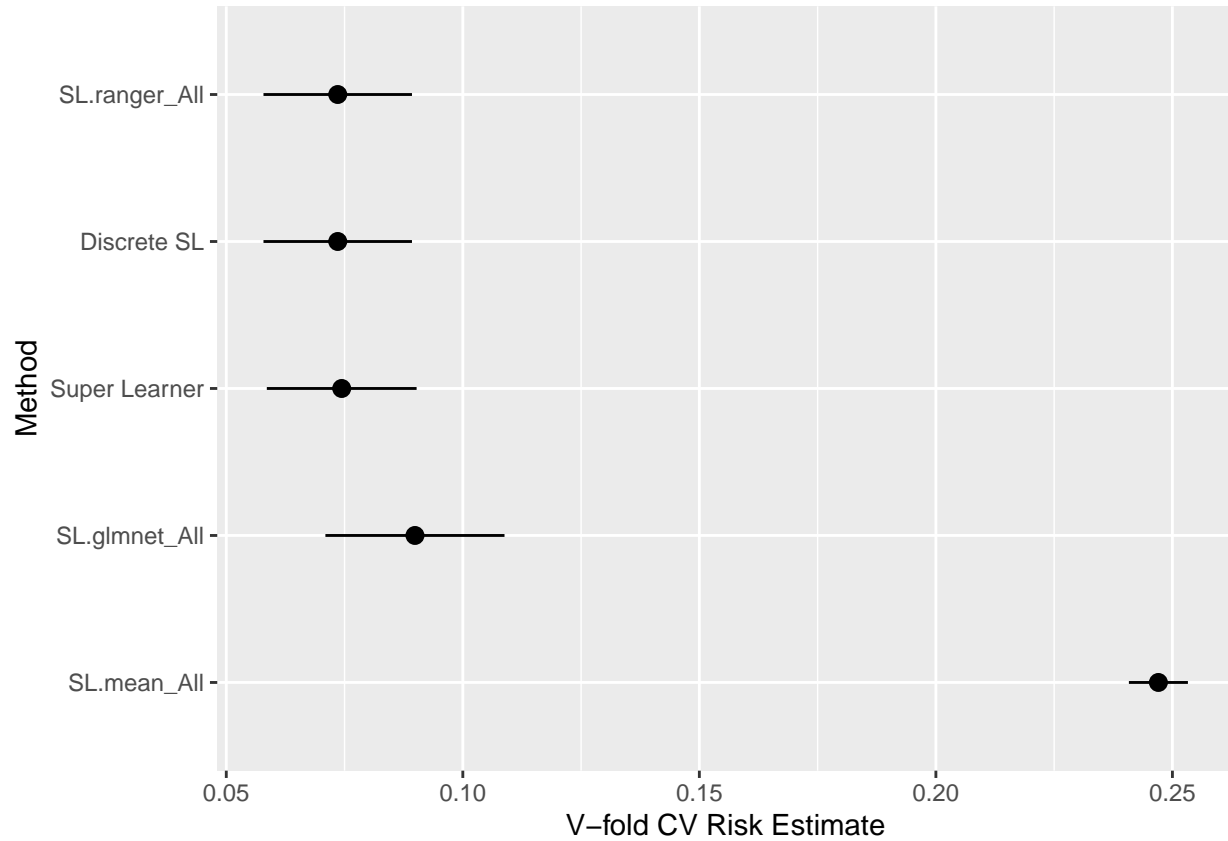
cv_sl = CV.SuperLearner(Y = y_train,
                       X = x_train,
                       family = binomial(),
                       V = 20,
                       parallel = 'multicore',
                       #parallel = cluster,
```

```

SL.library = c("SL.mean",
               "SL.glmnet",
               "SL.ranger"))

plot(cv_sl)

```



```

# Windows (for SL only, the above should work for tidymodels)
cluster = parallel::makeCluster(availableCores() - 1)
# Load SuperLearner onto all clusters
parallel::clusterEvalQ(cluster, library(SuperLearner))

```

```

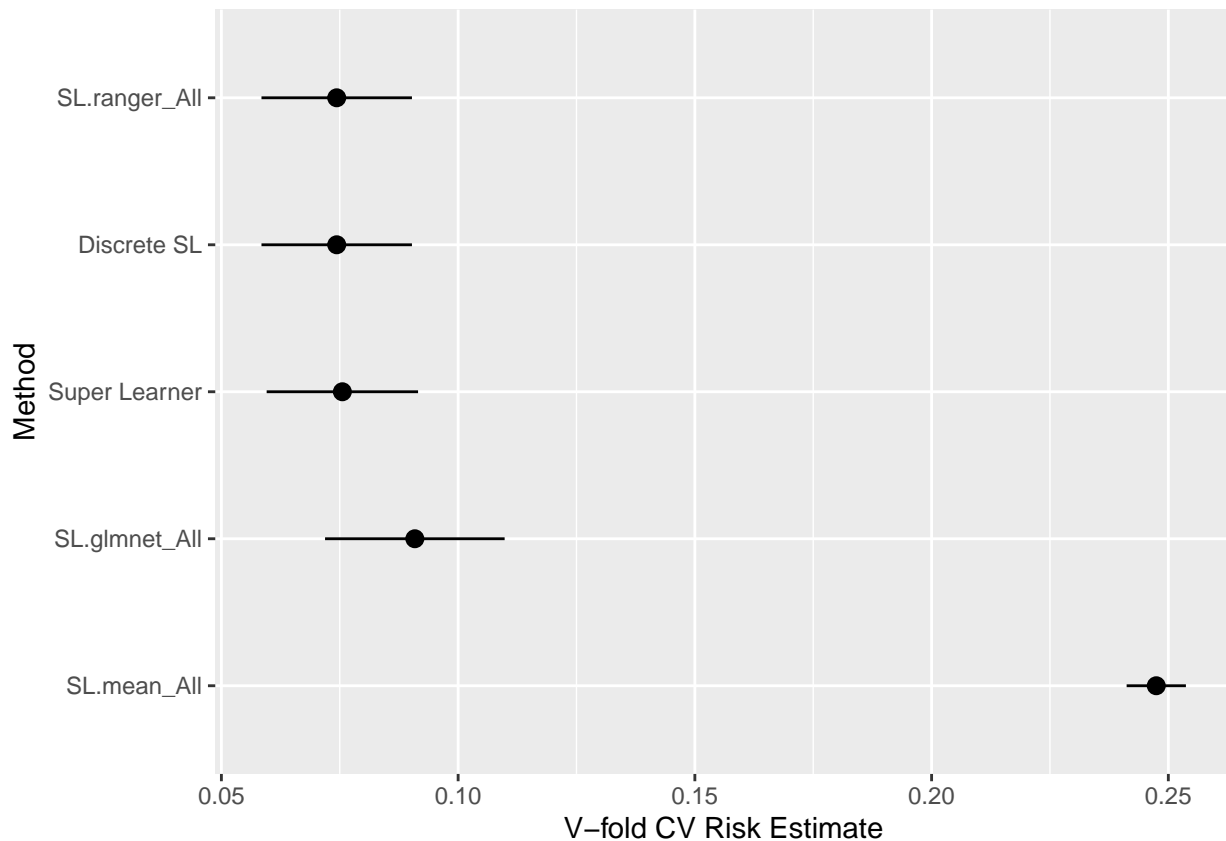
## [[1]]
## [1] "SuperLearner" "gam"          "foreach"      "splines"      "nnls"
## [6] "stats"        "graphics"     "grDevices"    "utils"        "datasets"
## [11] "methods"      "base"
##
## [[2]]
## [1] "SuperLearner" "gam"          "foreach"      "splines"      "nnls"
## [6] "stats"        "graphics"     "grDevices"    "utils"        "datasets"
## [11] "methods"      "base"
##
## [[3]]
## [1] "SuperLearner" "gam"          "foreach"      "splines"      "nnls"
## [6] "stats"        "graphics"     "grDevices"    "utils"        "datasets"
## [11] "methods"      "base"

```

```
parallel::clusterSetRNGStream(cluster, 1)
cv_sl = CV.SuperLearner(Y = y_train, X = x_train, family = binomial(),
  # For a real analysis we would use V = 10.
  V = 20,
  parallel = cluster,
  SL.library = c("SL.mean",
                 "SL.glmnet",
                 "SL.ranger"))

parallel::stopCluster(cluster)

plot(cv_sl)
```



Targeted Maximum Likelihood Estimation (TMLE)

We're now ready to move to the TMLE framework. TMLE combines machine learning and causal inference by using machine learning (i.e. data-adaptive models) to estimate some quantity of interest (so making inference). The key is that even though machine learning models oftentimes do not have outputs like coefficients, they can still be used to target the estimator (i.e. a regression) to that quantity of interest. This has a few different benefits, the primary one being that this framework creates a **doubly robust** estimator. Double robustness means that if we either:

1. Fit the right model to estimate the expected outcome correctly

OR

2. Fit the model to estimate the probability of treatment correctly

Then the final TMLE estimator will be **consistent**. Consistent means that as the sample size grows to infinity, the bias will drop to 0. We're going to explore this idea in depth. The example we will go through here is drawn from Katherine Hoffman's [blog](#). She also provides visual guides to [SuperLearner](#) and [TMLE](#). I recommend consulting the TMLE visual guide in particular as you work through these steps.

The basic procedure we will go through is:

1. Estimate the outcome model, $\bar{Q}_n^0(A, W)$
2. Estimate the probability of treatment model so we can “target” the initial estimate, $g(W) = P(A = 1|W)$
3. Extract the “clever covariate” or fluctuation parameter that tells us how to update our initial estimate
4. Update the initial estimate to get $\bar{Q}_n^1(A, W)$
5. Calculate ATE
6. Calculate confidence intervals

Step 1: Initial Estimate of the Outcome

The first step in TMLE is to estimate the **outcome model**, or the prediction of our target, Y , conditional on the treatment status, A and covariates/features, W . We could do this via a classic regression model, but we could also flexibly fit the model using machine learning. In this case, we'll create a SuperLearner ensemble. This Q step gives us our initial estimate. Do the following:

1. Create a SuperLearner library called `sl_libs`
2. Prepare an outcome vector, Y and treatment/covariate matrix, W_A
3. Fit the SuperLearner model for the Q step to obtain an initial estimate of the outcome

```
# SuperLearner libraries
sl_libs <- c('SL.glmnet', 'SL.ranger', 'SL.glm')

# Prepare data for SuperLearner/TMLE
# Mutate Y, A for outcome and treatment, use tax, age, and crim as covariates
data_obs <- Boston %>%
  mutate(Y = ifelse(medv > 22,
                    1,
                    0)) %>%
  rename(A = chas) %>%
  select(Y, A, tax, age, crim)

# Data Prep
# Outcome
Y <- data_obs %>% pull(Y)
# Covariates
W_A <- data_obs %>% select(-Y)

# Fit SL for Q step, initial estimate of the outcome
Q <- SuperLearner(Y = Y,
                  X = W_A,
                  family = binomial(),
                  SL.library = sl_libs)
```

Now that we have trained our model, we need to create predictions for three different scenarios:

1. Predictions assuming every unit had its observed treatment status.
2. Predictions assuming every unit was treated.
3. Predictions assuming every unit was control.

Fill in the code to obtain these predictions.


```

# observed treatment
Q_A <- as.vector(predict(Q)$pred)

# if every unit was treated
W_A1 <- W_A %>% mutate(A = 1)
Q_1 <- as.vector(predict(Q, newdata = W_A1)$pred)

# if everyone was control
W_A0 <- W_A %>% mutate(A = 0)
Q_0 <- as.vector(predict(Q, newdata = W_A0)$pred)

```

We can take our predictions and put them all into one dataframe:

```

dat_tmle <- tibble(Y = Y, A = W_A$A, Q_A, Q_0, Q_1)
head(dat_tmle)

```

```

## # A tibble: 6 x 5
##       Y     A   Q_A   Q_0   Q_1
##   <dbl> <int> <dbl> <dbl> <dbl>
## 1     1     0 0.802 0.802 0.849
## 2     0     0 0.236 0.236 0.315
## 3     1     0 0.676 0.676 0.727
## 4     1     0 0.841 0.841 0.878
## 5     1     0 0.889 0.889 0.916
## 6     1     0 0.847 0.847 0.876

```

We *could* go ahead and grab our ATE now using G-computation, which is the difference in expected outcomes under treatment and control conditional on covariates. However, the estimate we get below is targeted (optimized the bias-variance tradeoff) at the predictions of the outcome, not the ATE.

```

ate_gcomp <- mean(dat_tmle$Q_1 - dat_tmle$Q_0)
ate_gcomp

```

```
## [1] 0.1282819
```

Step 2: Probability of Treatment

Now that we have an initial estimate, we want to “target” it. To do this, we first need to estimate the probability of every unit receiving treatment. This step should look similar to how we estimated the propensity score during matching. Fit a g model using SuperLearner. **Hint:** Think carefully about what should be supplied to the Y and X arguments here.

```

A <- W_A$A
W <- Boston %>% select(tax, age, crim)

g <- SuperLearner(Y = A,
                  X = W,
                  family=binomial(),
                  SL.library=sl_libs)

```

Now that we have a model for our propensity score, we can go ahead and calculate both the inverse probability of receiving treatment and the negative inverse probability of not receiving treatment (basically the probability of not receiving treatment). Fill in the code below to obtain both of these quantities. **Hint:** The name “negative inverse probability of not receiving treatment” is a mouthful but should give you a hint about how to set up the calculation.

```

# Prediction for probability of treatment
g_w <- as.vector(predict(g)$pred) # Pr(A=1|W)

# probability of treatment
H_1 <- 1/g_w

# probability of control
H_0 <- -1/(1-g_w)

```

We can then create the “clever covariate” which assigns the probability of treatment to treated variables, and the probability of control to control variables.

```

#
dat_tmle <- # add clever covariate data to dat_tmle
  dat_tmle %>%
  bind_cols(
    H_1 = H_1,
    H_0 = H_0) %>%
  mutate(H_A = case_when(A == 1 ~ H_1,
    A == 0 ~ H_0))

```

We now have the initial estimate of the outcome, Q_A , and the estimates for the probability of treatment, H_A .

Step 3: Fluctuation Parameter

We are now ready to update our initial estimate, Q_A with the information from our propensity score estimates, H_A . We perform this update by fitting a logistic regression where we predict our outcome after passing our initial estimate, Q_A through a logistic transformation and then fitting a logistic regression (the `-1 + offset()` here is necessary because our intercept term is not constant).

```
glm_fit <- glm(Y ~ -1 + offset(qlogis(Q_A)) + H_A, data=dat_tmle, family=binomial)
```

We can then grab the coefficient from our model:

```
eps <- coef(glm_fit)
```

eps (or $\text{epsilon}/\hat{\epsilon}$) is the **fluctuation parameter**. This is the coefficient on our “clever covariate”, and basically tells us how much to update our initial model, Q by.

Step 4: Update Initial Estimates

We can now update the expected outcome for all of observations, conditional on their observed treatment status and their covariates:

```

H_A <- dat_tmle$H_A
Q_A_update <- plogis(qlogis(Q_A) + eps*H_A)

```

Do the same for outcome under treatment:

```
Q_1_update <- plogis(qlogis(Q_1) + eps*H_1)
```

Do the same for outcome under control:

```
Q_0_update <- plogis(qlogis(Q_0) + eps*H_0)
```

We now have updated expected outcomes for each unit for their actual treatment status, as well as simulated if they were all in treatment and all in control.

Step 5: Compute the Statistical Estimand of Interest

Fill in the code to calculate the ATE using our updated information for Q:

```
tmle_ate <- mean(Q_1_update - Q_0_update)
tmle_ate
```

```
## [1] 0.2051534
```

Step 6: Calculate Standard Errors for CIs and p-values

And finally calculate the standard errors and p-values:

```
infl_fn <- (Y - Q_A_update) * H_A + Q_1_update - Q_0_update - tmle_ate

tmle_se <- sqrt(var(infl_fn)/nrow(data_obs))

conf_low <- tmle_ate - 1.96*tmle_se
conf_high <- tmle_ate + 1.96*tmle_se
pval <- 2 * (1 - pnorm(abs(tmle_ate / tmle_se)))
```

Via TMLE Package

Thankfully, we do not need to take all of these steps every time we use the TMLE framework. The `tmle()` function handles all of this for us! Notice how we get a similar ATE from running this function:

```
tmle_fit <-
  tmle::tmle(Y = Y,
            A = A,
            W = W,
            Q.SL.library = sl_libs,
            g.SL.library = sl_libs)

tmle_fit

## Additive Effect
## Parameter Estimate: 0.20382
## Estimated Variance: 0.00061811
## p-value: 2.441e-16
## 95% Conf Interval: (0.15509, 0.25255)
##
## Additive Effect among the Treated
## Parameter Estimate: 0.12966
## Estimated Variance: 0.0071194
## p-value: 0.12436
## 95% Conf Interval: (-0.035716, 0.29504)
##
## Additive Effect among the Controls
## Parameter Estimate: 0.43978
## Estimated Variance: 0.00059728
## p-value: <2e-16
## 95% Conf Interval: (0.39188, 0.48768)
```

LTMLE

One of the main shortcomings of the TMLE framework is that it can only handle baseline covariates (covariates measured before treatment). It also requires that intervention happens at a single time point. The `ltmle` library allows us to relax some of these restrictions so that we can handle time-dependent confounders. The basic setup for the LTMLE function is that we need to set up our W , A , and Y , and then pass it to the `ltmle()` function:

```
data_obs_ltmle <- data_obs %>%
  rename(W1 = tax, W2 = age, W3 = crim) %>%
  select(W1, W2, W3, A, Y)

result <- ltmle(data_obs_ltmle, Anodes = "A", Ynodes = "Y", abar = 1)
```

Single Time Point

Imagine if there was a time ordering to our data, in particular a model like:

$$W1- > W2- > W3- > A- > Y$$

Let's simulate some data that captures this relationship. Notice how we added two arguments here, `Lnodes` and `SL.library`. `SL.library` should look familiar now, and `Lnodes` refers to a "time varying covariate". For now we're going to leave this as a `NULL`, but we'll see how to use it soon.

```
rexpfit <- function(x) rbinom(n=length(x), size=1, prob=plogis(x))

n <- 1000
W1 <- rnorm(n)
W2 <- rbinom(n, size=1, prob=0.3)
W3 <- rnorm(n)
A <- rexpfit(-1 + 2 * W1 + W3)
Y <- rexpfit(-0.5 + 2 * W1^2 + 0.5 * W2 - 0.5 * A + 0.2 * W3 * A - 1.1 * W3)
data <- data.frame(W1, W2, W3, A, Y)

result <- ltmle(data, Anodes="A", Lnodes=NULL, Ynodes="Y", abar=1, SL.library=sl_libs)
```

Longitudinal Data Structure

Now imagine we instead of a time ordering like this:

$$W- > A1- > L- > A2- > Y$$

We can simulate some more data to reflect this structure, and then fit a `ltmle()` function. Notice how now we have "L" nodes. In the data structure specified above, we have some baseline covariates that occur before the first treatment, some time dependent covariate that comes after the first treatment but before the second, then a second treatment and then outcome. Notice how we can now add L nodes that occur in between the two treatments so that we can deal with this "time-dependent confounding".

```
n <- 1000
W <- rnorm(n)
A1 <- rexpfit(W)
L <- 0.3 * W + 0.2 * A1 + rnorm(n)
A2 <- rexpfit(W + A1 + L)
```



```
## one multinomial or binomial class has 1 or 0 observations; not allowed
## Error in pred[, "1"] : subscript out of bounds
## Error in lognet(xd, is.sparse, ix, jx, y, weights, offset, alpha, nobs, :
## one multinomial or binomial class has 1 or 0 observations; not allowed
## Error in pred[, "1"] : subscript out of bounds
## Error in lognet(xd, is.sparse, ix, jx, y, weights, offset, alpha, nobs, :
## one multinomial or binomial class has 1 or 0 observations; not allowed
## Error in pred[, "1"] : subscript out of bounds
## Error in lognet(xd, is.sparse, ix, jx, y, weights, offset, alpha, nobs, :
## one multinomial or binomial class has 1 or 0 observations; not allowed
## Error in pred[, "1"] : subscript out of bounds
## Error in lognet(xd, is.sparse, ix, jx, y, weights, offset, alpha, nobs, :
## one multinomial or binomial class has 1 or 0 observations; not allowed
## Error in pred[, "1"] : subscript out of bounds
## Error in lognet(xd, is.sparse, ix, jx, y, weights, offset, alpha, nobs, :
## one multinomial or binomial class has 1 or 0 observations; not allowed
## Error in pred[, "1"] : subscript out of bounds
## Call:
## ltmle(data = data, Anodes = c("A1", "A2"), Lnodes = "L", Ynodes = "Y",
## abar = c(1, 1), SL.library = sl_libs)
##
## TMLE Estimate: 0.297777
```

The [ltmle vignette](#) provides even more examples, including for censored data. The main concept here is that ltmle allows for causal estimates in observational data in complicated treatment regimes such as staggered adoption, subject dropout, multiple treatments, etc.

Concluding Thoughts

We have covered a lot of ground this year! To recap:

- Reproducible Data Science
 - GitHub/version control
 - Collaborative data science
- Machine Learning
 - Supervised Learning
 - * Regression
 - * Classification
 - Unsupervised Learning
 - * Dimensionality Reduction
 - * Clustering/grouping
- Text Analysis
 - Text preprocessing
 - Dictionary methods, word2vec, sentiment analysis
 - Prediction
- Causal Inference
 - Matching
 - Regression Discontinuity
 - Diff-in-Diffs/Synthetic Control
 - Sensitivity Analysis
 - SuperLearner/Targeted Maximum Likelihood Estimation

This covers a lot of computational social science! But there's still more. As you're thinking about where to go next, I recommend exploring these areas:

- Computational Tools/Data Acquisition
 - Web scraping (selenium)
 - APIs
 - Bash/CLI
 - XML/HTML parsing (BeautifulSoup)
- High Performance Computing
 - Parallel processing. We mostly covered “embarrassingly parallel” techniques, but there are other more advanced ones as well
 - Amazon Web Services (AWS)/Microsoft Azure cloud computing
 - Secure File Transfer Protocol (SFTP)
- Deep Learning
 - We covered some of the basics but more applications would be good
 - GPU acceleration
 - Deep Learning for text analysis (Bidirectional Encoder Representations from Transformers (BERT))
- Machine Learning and Causal Inference
 - Meta-estimators (like TMLE)
 - Heterogeneous Treatment Effects

Appendix for SL & LTMLE

a) We can now plot our [AUC-ROC curve](#):

$$Sensitivity = TruePositiveRate = Recall = \frac{TP}{TP + FN}$$

$$Specificity = TrueNegativeRate = \frac{TN}{TN + FP}$$

```
roc_df <- roc_curve(validation,
  pred,
  truth = as.factor(obs))

roc_df %>%
  ggplot() +
  geom_point(aes(x = specificity, y = 1 - sensitivity)) +
  geom_line(aes(x = specificity, y = 1 - sensitivity)) +
  theme_fivethirtyeight() +
  theme(axis.title = element_text()) +
  ggtitle('AUC-ROC Curve for SuperLearner') +
  xlab('Specificity') +
  ylab('Sensitivity')
```

