



UNIVERSITY OF  
COPENHAGEN



Targeted register  
analyses  
*PhD short course*

# A Gentle Introduction to Ensemble Machine Learning and Targeted Maximum Likelihood Estimation

**Andrew Mertens**

University of California, Berkeley, Division of Biostatistics

# Introduction to the Causal Roadmap

## 1. Causal question

- Translate scientific question into causal parameter (defined in terms of counterfactual outcomes)

## 2. Observed data & statistical model

- Model should reflect uncertainty

## 3. Identify

- Translate causal parameter to statistical parameter under explicit causal assumptions

## 4. Estimate

## 5. Interpret

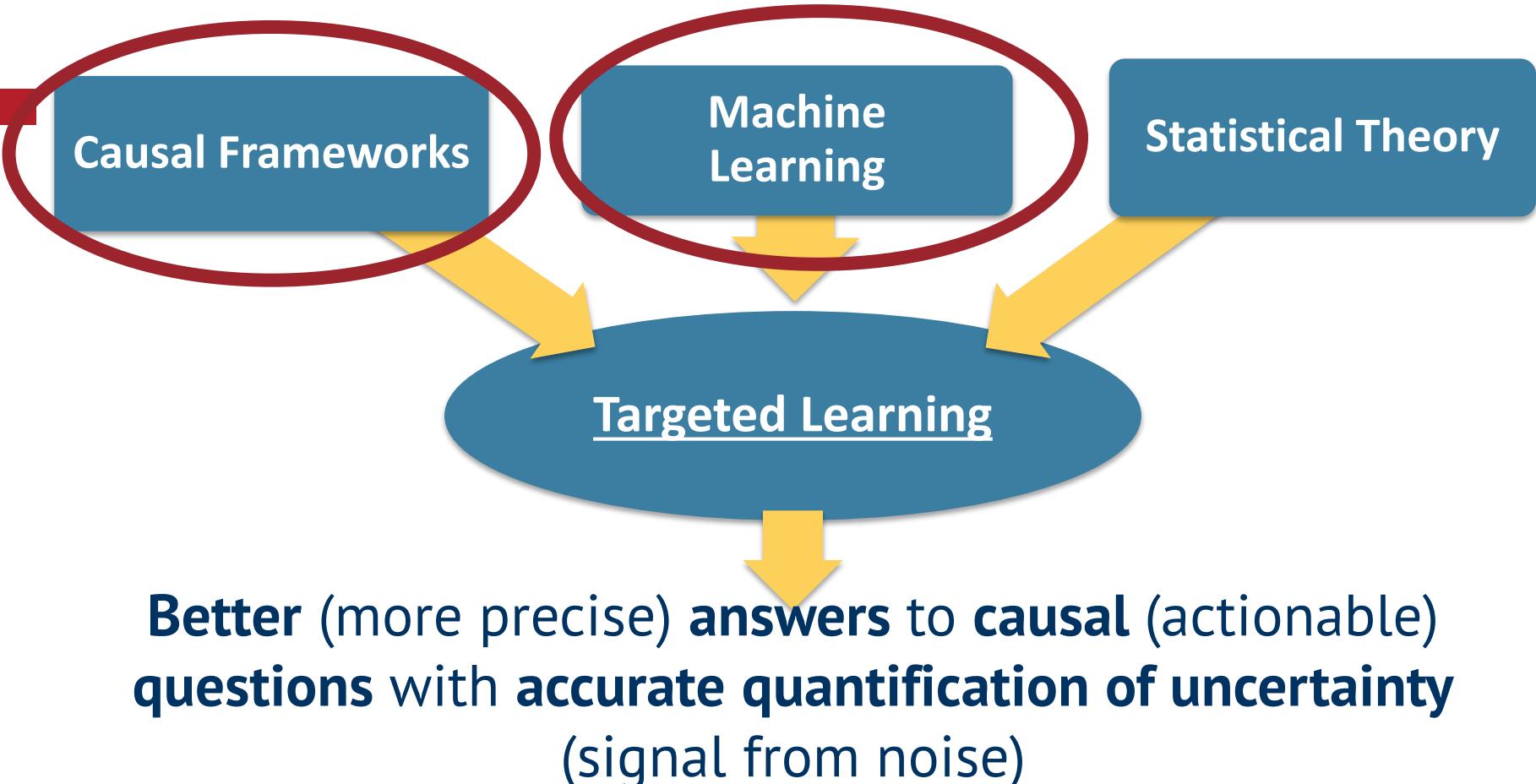
**Causal Question:** Ex. Difference in dementia risk (by, eg, 5 years) under different ideal longitudinal protocols



**Statistical Target Parameter**



- **Statistical Estimate**
- **Inference (95% CI, etc)**



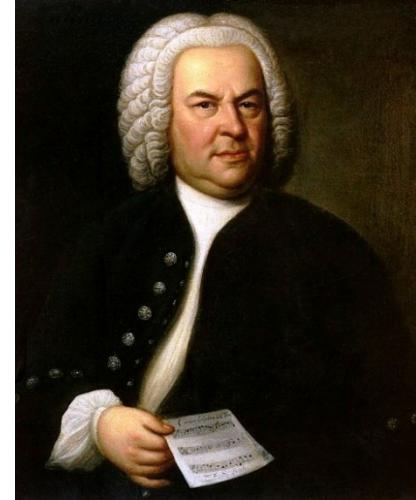
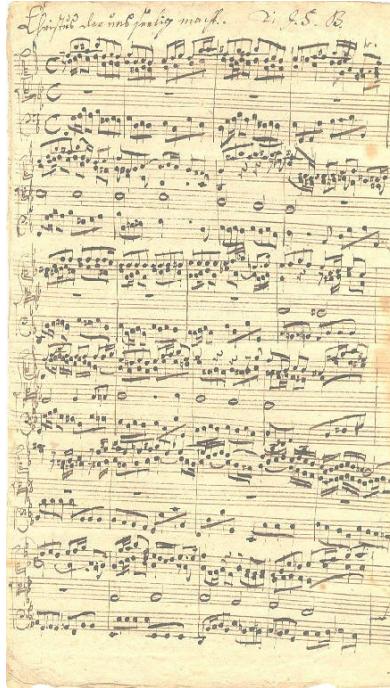
# Superlearner: the algorithm steps

1. Pick the learners
2. Split data into cross-validation folds
3. Fit learners on cross-validation folds
4. Obtain predictions from each fitted model within folds
5. Use a metalearner to combine predictions across learners
6. Repeat for the full dataset
7. Predict on new data

# Step 1: pick the learners

*Superlearner is a type of “ensemble machine learning”*

- The word “ensemble” from the arts. It means a collection of musicians or performers.



## Step 1: pick the learners

*Superlearner is a type of “ensemble machine learning”*

- Ensemble strengths:

# Step 1: pick the learners

*Superlearner is a type of “ensemble machine learning”*

- Ensemble strengths:
  - Diversity



# Step 1: pick the learners

*Superlearner is a type of “ensemble machine learning”*

- Ensemble strengths:
  - Diversity
  - Redundancy



© BayTaper.com

# Step 1: pick the learners

*Superlearner is a type of “ensemble machine learning”*

- Ensemble strengths:
  - Diversity
  - Redundancy
  - Synergy



# Step 1: pick the learners

*Superlearner is a type of “ensemble machine learning”*

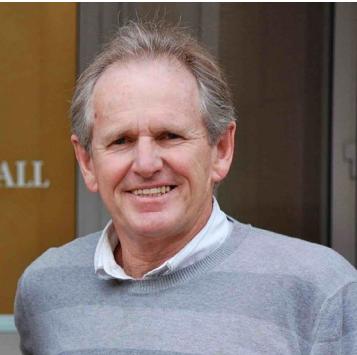
- Ensemble strengths:
  - Diversity
  - Redundancy
  - Synergy



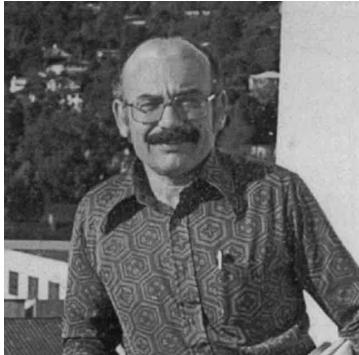
# ensemble of statisticians



Tibshirani



Hastie



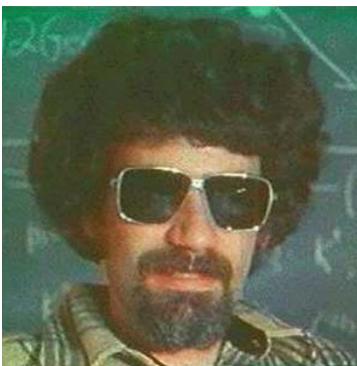
Breiman



Cleveland



Yman



Friedman



Gauss



Babylonians



van der Laan



Wolpert

*W. bancrofti* Bm14 Luminex Response

Model / algorithm

CV MSE

$10^5$   
 $10^4$   
 $10^3$   
 $10^2$   
 $10^1$   
 $10^0$

0 2 4 6 8 10 12

Age, years

## *W. bancrofti* Bm14 Luminex Response

Model / algorithm



Mean

CV MSE

1.67

$10^5$   
 $10^4$   
 $10^3$   
 $10^2$   
 $10^1$   
 $10^0$

0 2 4 6 8 10 12

Age, years

## *W. bancrofti* Bm14 Luminex Response

### Model / algorithm



Mean

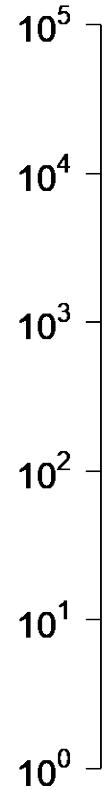


Linear regression

### CV MSE

1.67

1.13

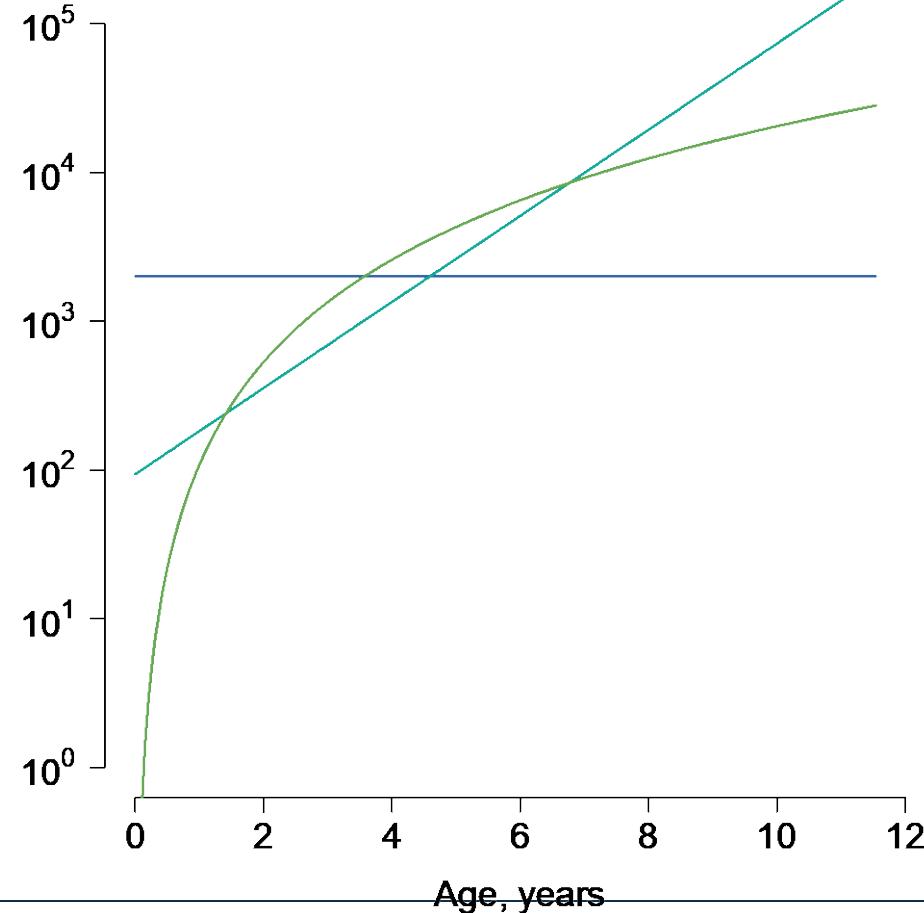


0 2 4 6 8 10 12

Age, years

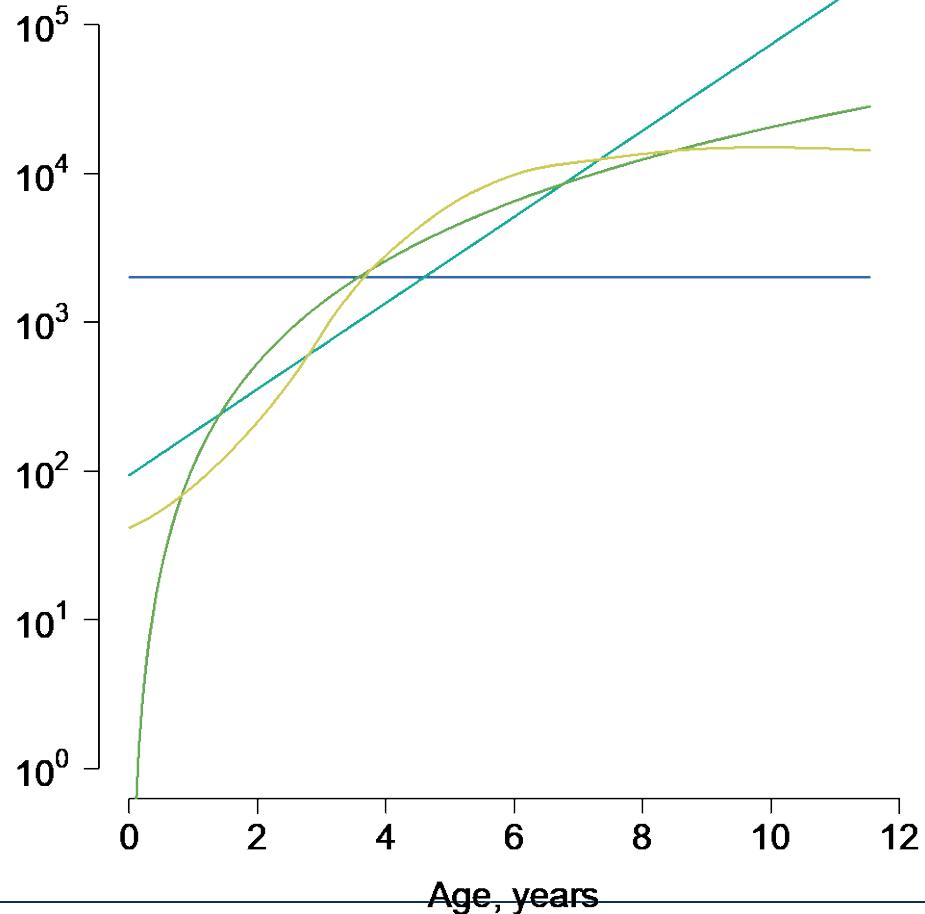
## *W. bancrofti* Bm14 Luminex Response

Model / algorithm	CV MSE
Mean	1.67
Linear regression	1.13
<b>Ab acquisition model</b>	<b>1.24</b>



## *W. bancrofti* Bm14 Luminex Response

Model / algorithm	CV MSE
Mean	1.67
Linear regression	1.13
Ab acquisition model	1.24
LOESS	1.03



## *W. bancrofti* Bm14 Luminex Response

### Model / algorithm



Mean



Linear regression



Ab acquisition model



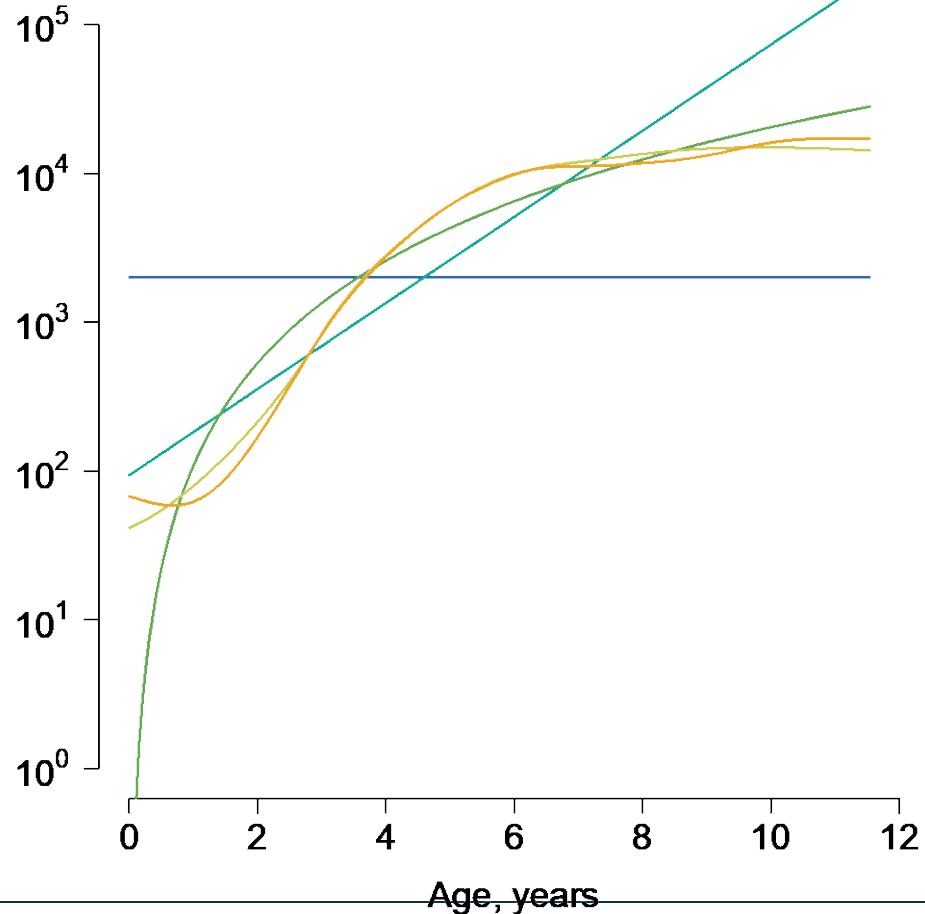
LOESS



Splines

### CV MSE

1.67  
1.13  
1.24  
1.03  
1.01



## *W. bancrofti* Bm14 Luminex Response

### Model / algorithm

### CV MSE



Mean

$10^5$

1.67



Linear regression

$10^4$

1.13



Ab acquisition model

$10^3$

1.24



LOESS

$10^2$

1.03



Splines

$10^1$

1.01



MARS

1.02

$10^0$

$10^5$

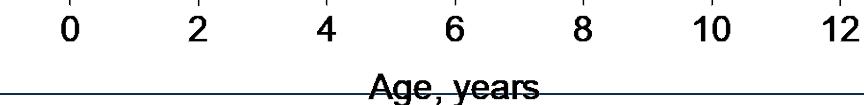
$10^4$

$10^3$

$10^2$

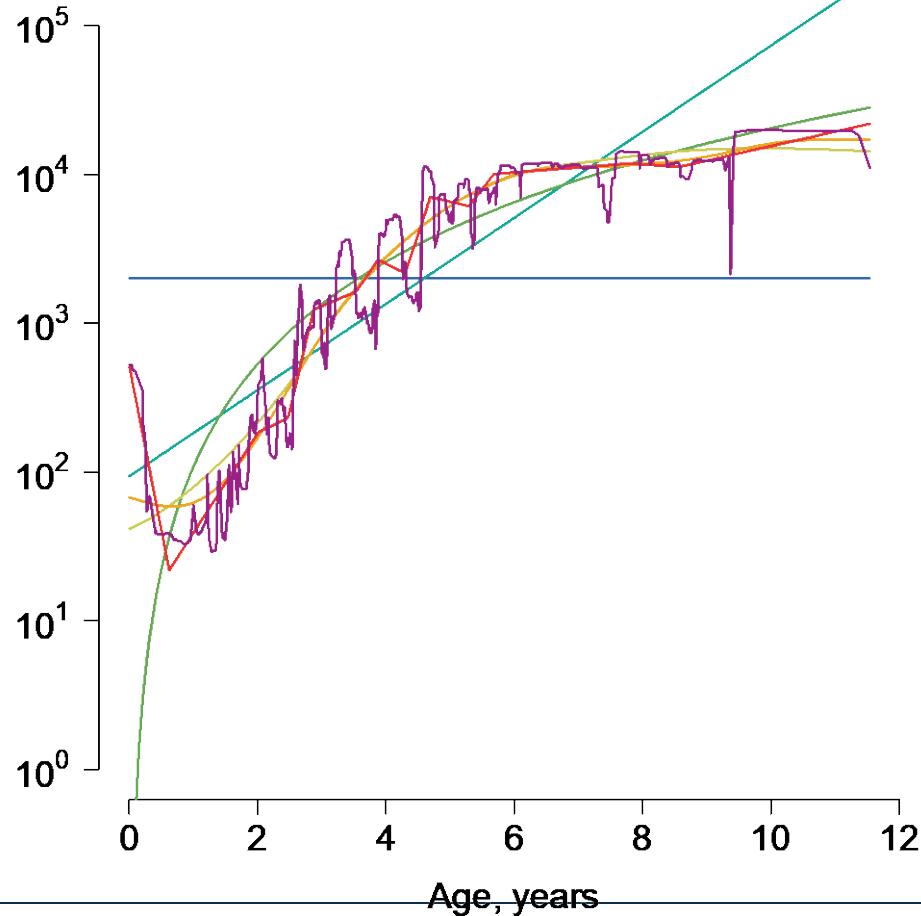
$10^1$

$10^0$



## *W. bancrofti* Bm14 Luminex Response

Model / algorithm	CV MS
Mean	1.67
Linear regression	1.13
Ab acquisition model	1.24
LOESS	1.03
Splines	1.01
MARS	1.02
<b>Random Forest</b>	<b>1.04</b>

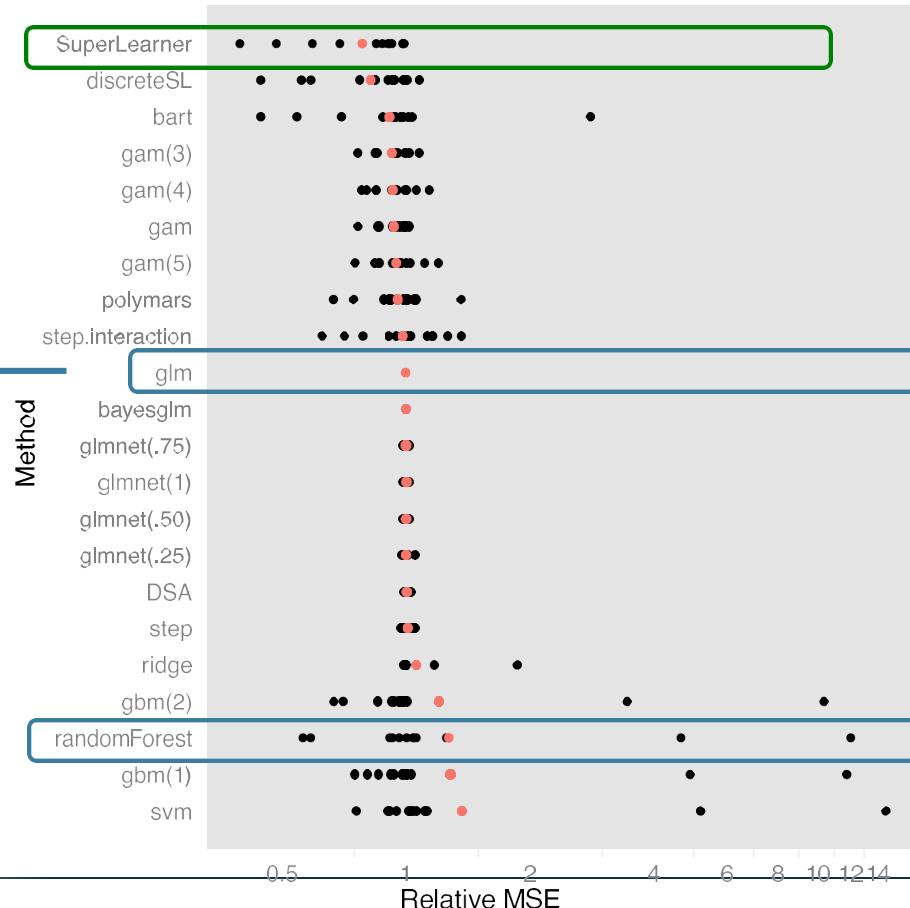


# Super Learning: Building a winning team!

Super Learner: ←  
Best weighted  
combination of  
algorithms for a  
given prediction  
problem

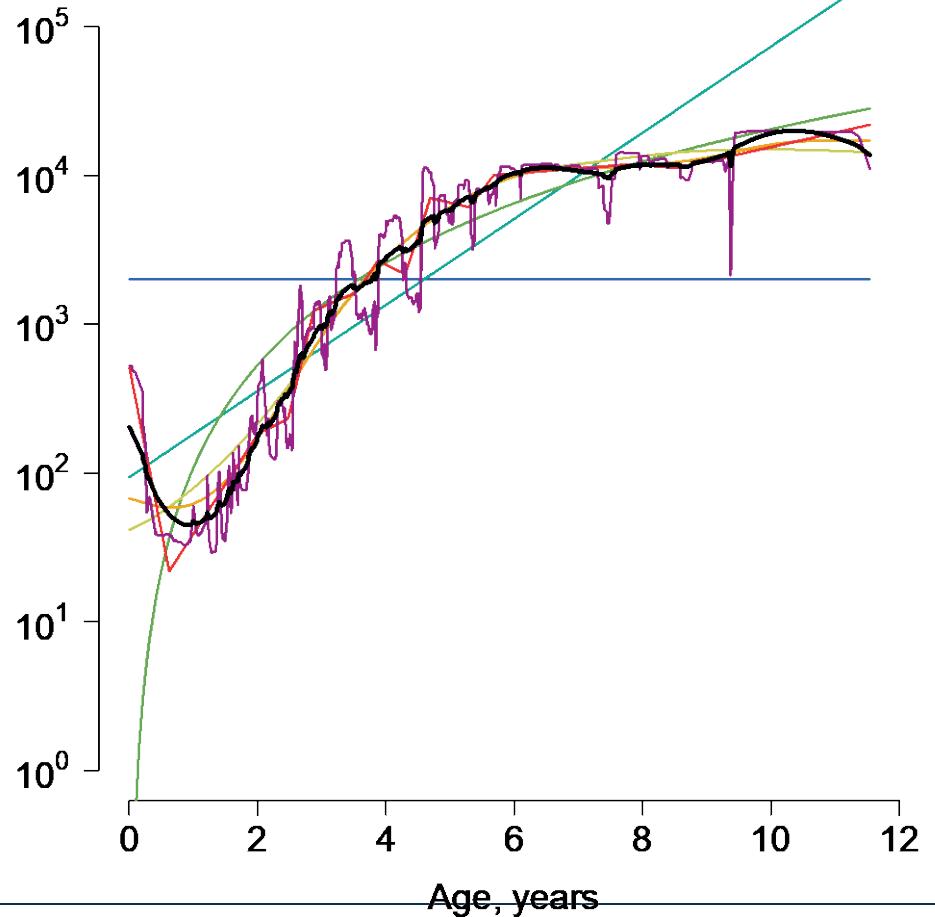
Example algorithm:  
Linear Main Term  
Regression

Example algorithm: ←  
Random Forest



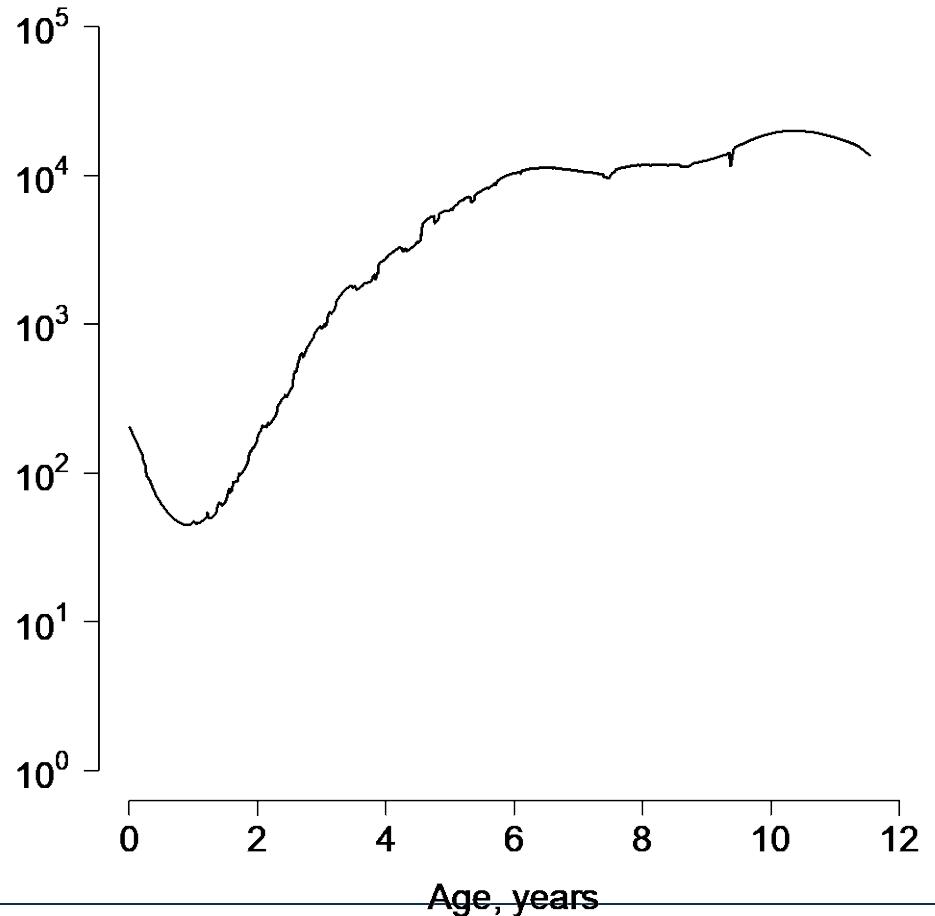
## *W. bancrofti* Bm14 Luminex Response

Model / algorithm	CV MSE
Mean	1.67
Linear regression	1.13
Ab acquisition model	1.24
LOESS	1.03
Splines	1.01
MARS	1.02
Random Forest	1.04
<b>Super Learner (stacked ensemble)</b>	<b>1.01</b>



## *W. bancrofti* Bm14 Luminex Response

Model / algorithm	CV MSE
Mean	1.67
Linear regression	1.13
Ab acquisition model	1.24
LOESS	1.03
Splines	1.01
MARS	1.02
Random Forest	1.04
<b>Super Learner (stacked ensemble)</b>	<b>1.01</b>



# Superlearner: the algorithm steps

1. Pick the learners
2. Split data into cross-validation folds
3. Fit learners on cross-validation folds
4. Obtain predictions from each fitted model within folds
5. Use a metalearner to combine predictions across learners
6. Fit base learners on the full dataset
7. Use metalearner fit to weight base learners and get SuperLearner prediction
8. Predict on new data

## Step 1: Pick the learners

- Better to pick a variety of different learners
- Can be guided by what's normally used in your field + flexible ML algorithms
- You we learn more about individual algorithms tomorrow
- The **only downside** to using more learners is computation time

## Step 2: Split data into cross-validation folds

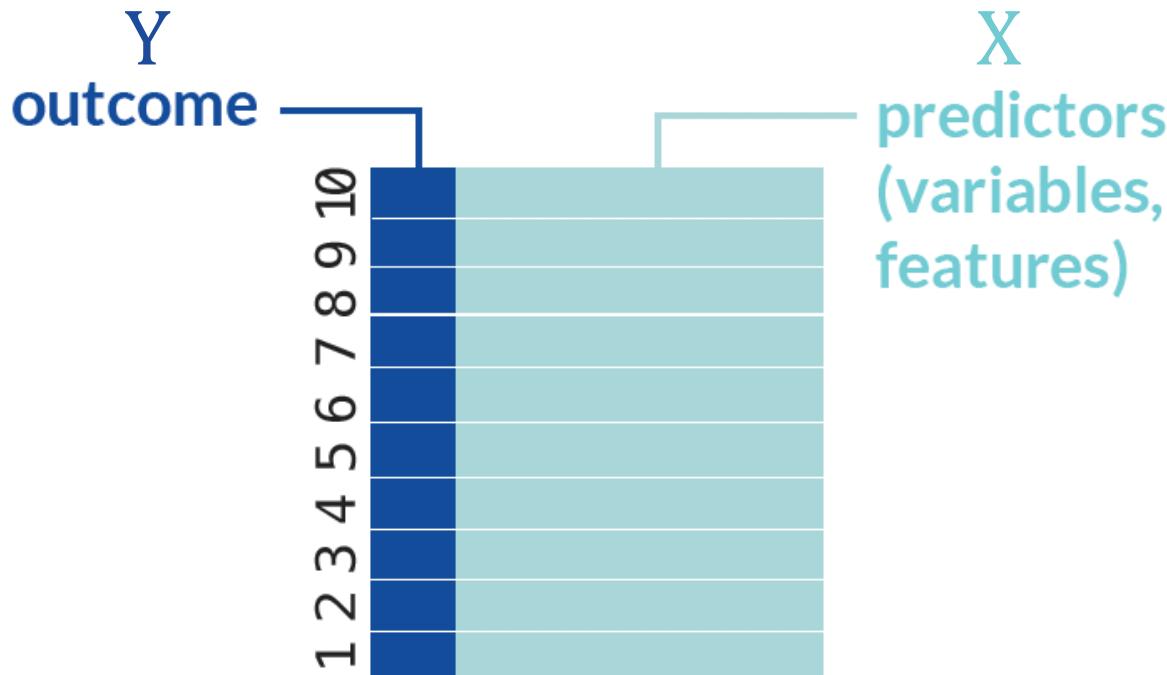
Example data

- ID variable identifying observation, individual or cluster
- Set of X predictors
- Y outcome

Simulated data set					
id	x1	x2	x3	x4	Y
1.000	2.287	1.000	1.000	1.385	5.270
2.000	-1.197	0.000	0.000	0.000	-1.197
3.000	-0.694	0.000	0.000	0.000	-0.694
4.000	-0.412	0.000	1.000	-0.541	-0.928
5.000	-0.971	0.000	0.000	0.000	-0.971
6.000	-0.947	0.000	1.000	-0.160	-1.107

## Step 2: Split data into cross-validation folds

Create indices for each of the K folds of size  $N/K$



# V-Fold Cross Validation

1	1	1	1	1	1	1	1	1	1
2	2	2	2	2	2	2	2	2	2
3	3	3	3	3	3	3	3	3	3
4	4	4	4	4	4	4	4	4	4
5	5	5	5	5	5	5	5	5	5
6	6	6	6	6	6	6	6	6	6
7	7	7	7	7	7	7	7	7	7
8	8	8	8	8	8	8	8	8	8
9	9	9	9	9	9	9	9	9	9
10	10	10	10	10	10	10	10	10	10
Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Fold 6	Fold 7	Fold 8	Fold 9	Fold 10

## Step 3: Fit learners on first cross-validation folds

Fit model for each learner on the training data.

Here, in the first CV fold, data fold 10 is held out and 3 learners are fit to folds 1-9.

$$Y = f(X)$$

```
fit_1a <- lrnr_a(
```



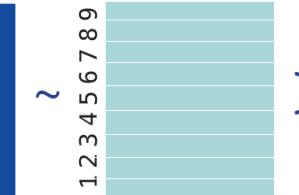
```
)
```

```
fit_1b <- lrnr_b(
```



```
)
```

```
fit_1c <- lrnr_c(
```



```
)
```

## Step 4: Obtain predictions from each fitted model within folds

Get the predictions for each held out fold.

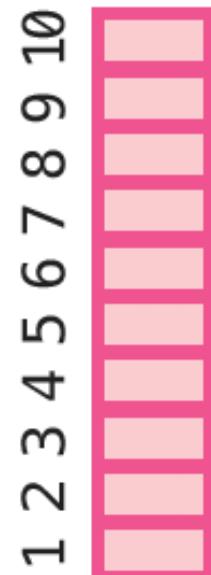
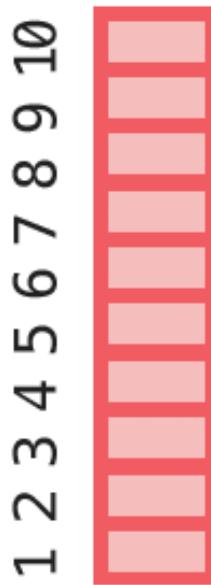
```
 $\hat{Y}_1$  10  <- predict(fit_1a,  
           newdata =  )
```

```
 $\hat{Y}_2$  10  <- predict(fit_1b,  
           newdata =  )
```

```
 $\hat{Y}_3$  10  <- predict(fit_1c,  
           newdata =  )
```

## Step 4: Obtain predictions from each fitted model within folds

Repeat for the entire dataset, so now there are outcome predictions for every row of data



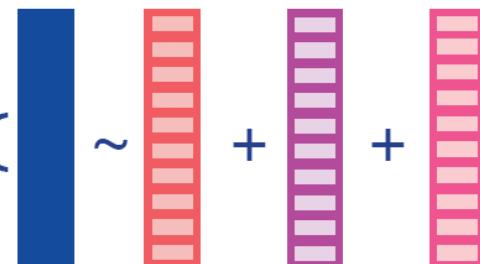
## Step 5: Use a metalearner to combine predictions across learners

- This is a function combining predictions across learners to minimize the loss function.
- Example: CV-MSE
- Could use something as simple as a linear regression, predicting the true outcome from the set of learner-specific prediction.
  - Coefficients become the weights of
- Non-negative least squares is the R package default, but variety of options.

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

$$Y = \beta_1 \hat{Y}_1 + \beta_2 \hat{Y}_2 + \beta_3 \hat{Y}_3$$

```
SL_fit <- meta_lrnrr( [ ] ~ [ ] + [ ] + [ ] )
```



## Step 5: Use a metalearner to combine predictions across learners

- You can compare the risk of different learners and the weights chosen by the metalearners using the SuperLearner R package

Call:

```
SuperLearner(Y = obs$y, X = x_df, family = gaussian(),  
SL.library = c("SL.ranger",  
"SL.glmnet", "SL.earth"))
```

	Risk (MSE)	Coef (weight)
SL.ranger_All	0.013672503	0.1606329
SL.glmnet_All	0.097257031	0
SL.earth_All	0.003181357	0.8393671

## Step 6: Fit base learners on the full dataset and get predictions

- A. Get a single model fit (not cross-validated) for each learner on the full data
- B. Get full-data predictions

A.

```
fit_a <- lrnr_a( [ ] ~ [ ] )
```

B.

```
[ ] <- predict(fit_a)
```

```
fit_b <- lrnr_b( [ ] ~ [ ] )
```

```
[ ] <- predict(fit_b)
```

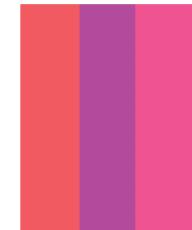
```
fit_c <- lrnr_c( [ ] ~ [ ] )
```

```
[ ] <- predict(fit_c)
```

## Step 7: Use metalearner fit to weight base learners and get SuperLearner prediction

Combine the learner-specific predictions with the estimated weights to get a single predicted outcome per observation

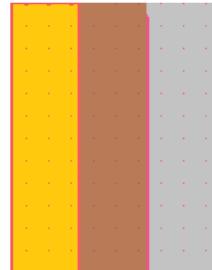
$$\widehat{Y}_{SL} = \beta_1 \widehat{Y}_1 + \beta_2 \widehat{Y}_2 + \beta_3 \widehat{Y}_3$$

```
 <- predict(SL_fit,  
           newdata =  )
```

## Step 7: (Optional) Predict on new data

- Repeat steps 5 and 6 on any new data
- This is how to use SuperLearner as a clinical screening/prediction tool



```
<- predict(SL_fit,  
          newdata =  )
```

## Extensions of SuperLearner

### Discrete SuperLearner:

- Simpler form: chose the single learner with the lowest CV risk
- Can be useful for hyperparameter tuning: put models with all

### Cross-validated SuperLearner :

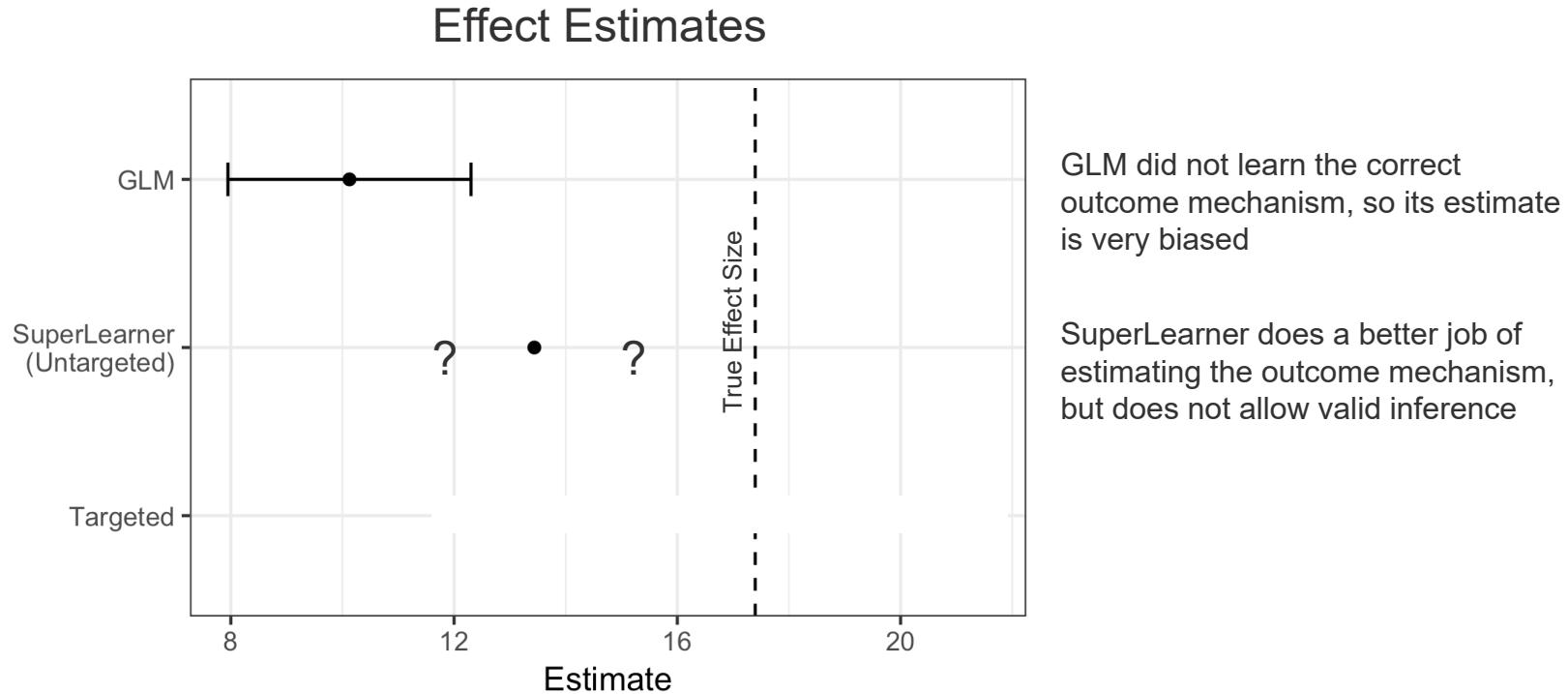
- Cross-validate the entire model fitting process
- Gets us the valid risk (ex. CV-MSE) of the SuperLearner ensemble in addition to the individual learners

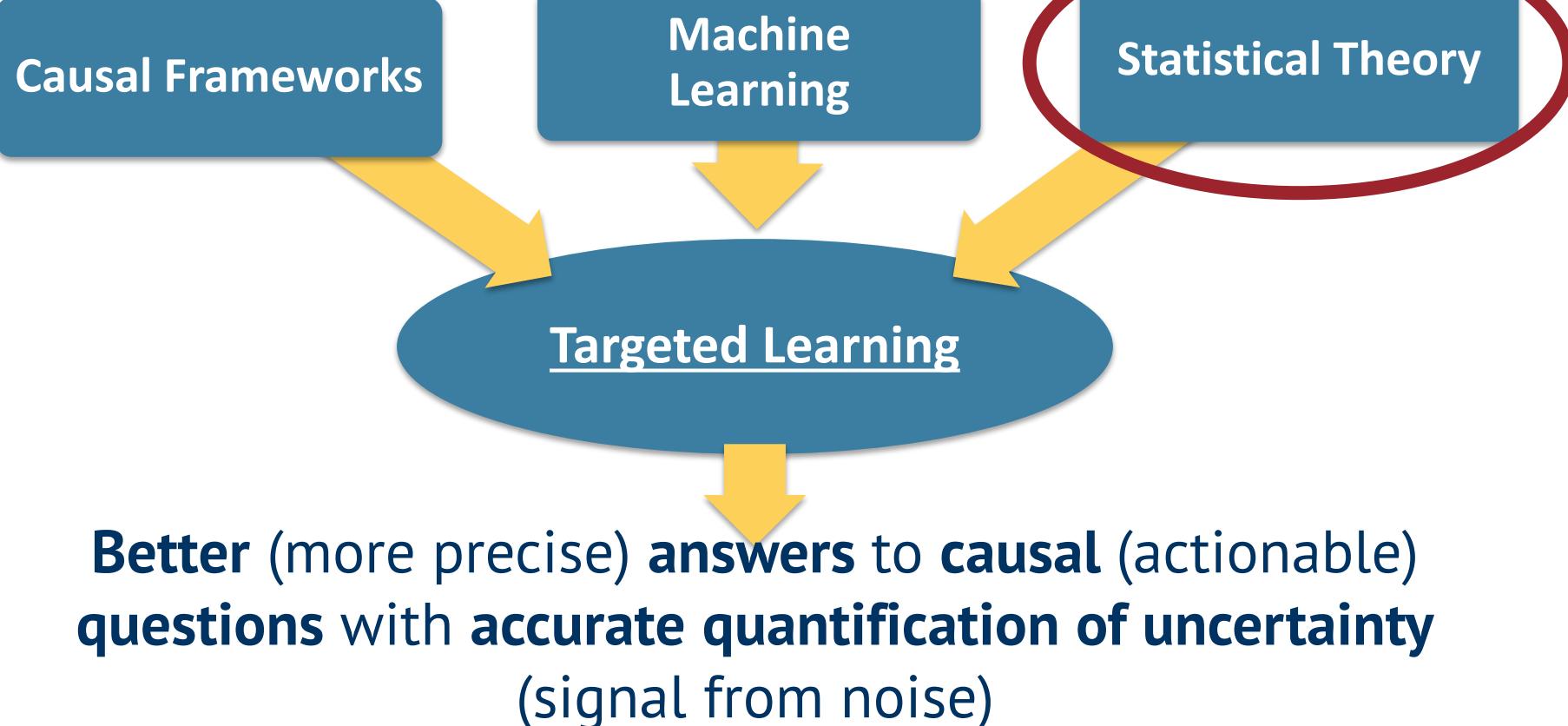
### Online SuperLearner:

- Continuously update the model and predictions as new data comes in

# Results: removing bias compared to individual prediction

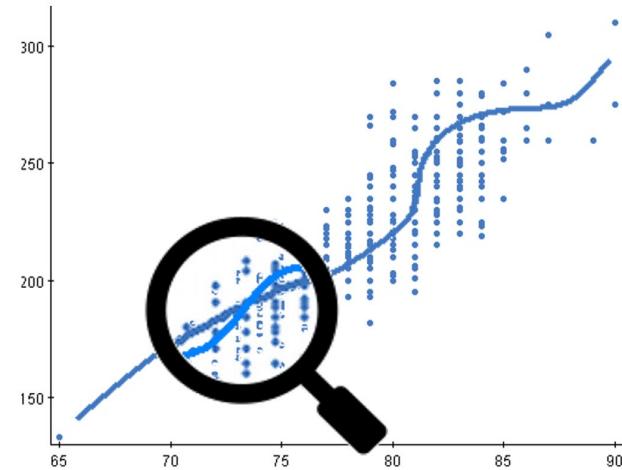
But what about inference?





# Beyond Machine Learning...

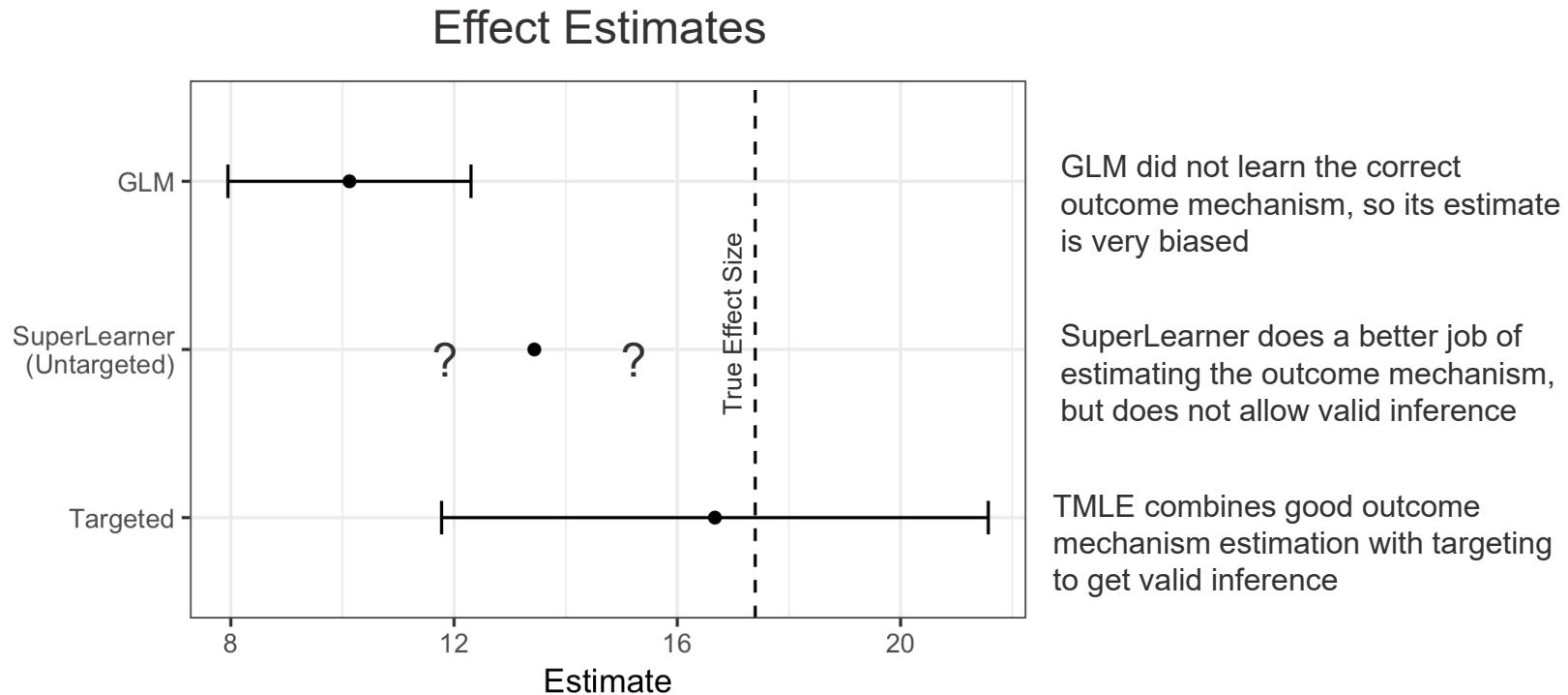
- Could use SuperLearner for more than prediction
  - Fit SuperLearner
  - Evaluate using a “plug in” estimator for average treatment effect
- Paying a price for too much ambition!
  - Super Learner- trying to do a good job on a harder question
- The costs: Increased bias and variance, no inference



## Targeting: Update initial Super learner fit

- Dont try to do a good job for all (causal) questions at once
- Focus estimation where it matters most for the question at hand
- Update depends on other “nuisance parameters” like propensity score

# Results: removing bias AND robust inference



# Why does TMLE work?

*Based on semi-parametric theory*

1. Some estimands allow for asymptotically linear estimation.
  1. AKA estimators can be represented as sample averages + term that converges to 0
2. The quantities being averaged are called influence functions
  1. quantifies how much influence each observation has on the estimator
  2. useful to characterize the variance of the estimator.
3. The efficient influence function (EIF) is the influence function that achieves the efficiency bound
  1. Cramer Rao Lower Bound from parametric maximum likelihood estimation
  2. Can be used to create efficient estimators.
4. Can use EIF to construct an estimator that is efficient,
5. Our estimand of interest admits asymptotically linear estimation
  1. We are using properties of the EIF to construct an estimator with optimal statistical properties (e.g. double robustness).
  2. Allows the use of machine learning models while still obtaining asymptotic properties for inference

## Example data for a walkthrough of TMLE steps

### Parameter of interest

- Parameter of interest: Average Treatment Effect  $Y_1 - Y_0$  for a point-treatment (not longitudinal) effect

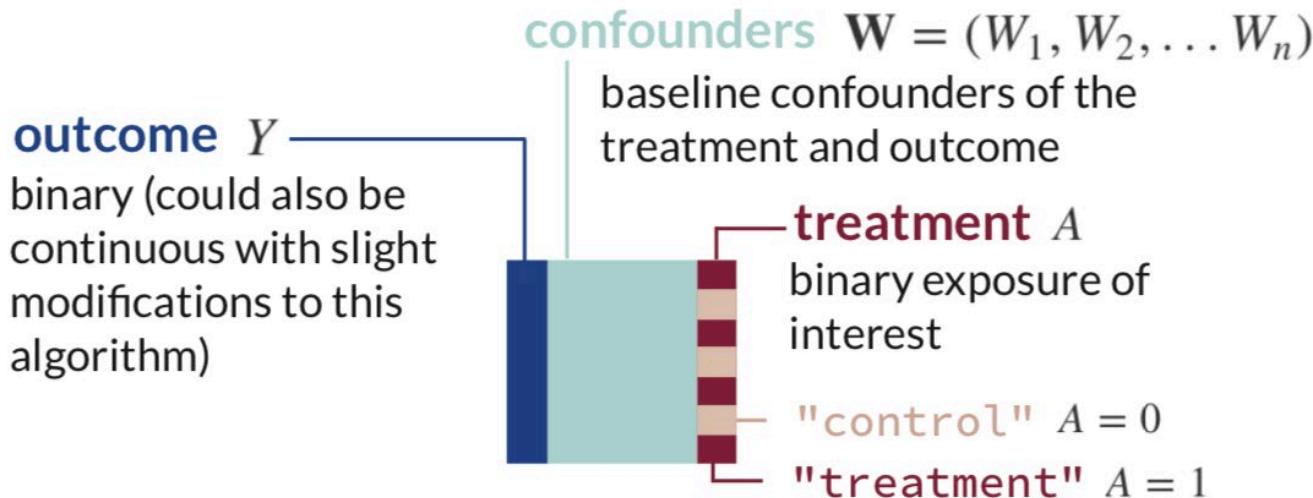
$$ATE = \Psi = E_W[E[Y|A = 1, \mathbf{W}] - E[Y|A = 0, \mathbf{W}]]$$

- Many other parameters possible (causal relative risk, treatment specific mean, more advanced discussed at the end)
- Longitudinal TMLE coming Thursday (Thanks Zeyi!)
- Remember, must meet identifiability assumptions outlined in the causal roadmap for causal inference, otherwise this is an associational analysis

## Example data for a walkthrough of TMLE steps

### Structure of data

$$ATE = \Psi = E_W[\mathbb{E}[Y|A = 1, \mathbf{W}] - \mathbb{E}[Y|A = 0, \mathbf{W}]]$$



## Example data for a walkthrough of TMLE steps

Example data

Simulated data set.					
Y	W1	W2	W3	W4	A
1	1	0	6	1	1
1	0	1	6	3	0
0	0	0	3	2	0
1	0	1	5	1	1
1	0	0	5	2	0
1	0	1	6	1	1

# TMLE: the algorithm steps

1. Estimate the Outcome
2. Estimate the Probability of Treatment
3. Estimate the Fluctuation Parameter
4. Update the Initial Estimates of the Expected Outcome
5. Compute the Statistical Estimand of Interest
6. Calculate the Standard Errors for Confidence Intervals and P-values

# TMLE: the algorithm steps

1. Estimate the Outcome
2. Estimate the Probability of Treatment
3. Estimate the Fluctuation Parameter
4. Update the Initial Estimates of the Expected Outcome
5. Compute the Statistical Estimand of Interest
6. Calculate the Standard Errors for Confidence Intervals and P-values

## Step 1: Estimate the Outcome

Predict the outcome

$$Q(A, W) = E[Y|A, W]$$

- This could be done with a parametric regression, or any other algorithm of interest
- Recommend using SuperLearner

```
outcome_fit <- fit(  ~  )
```

  $Q(A, W)$   `<- predict(outcome_fit)`

## Step 1: Estimate the Outcome

Predict the counterfactual

If you could roll back the clock and **treat everyone** in your sample, what would their outcome be?

$$Q(1, W) = E[Y|A = 1, W]$$

```
  <- predict(outcome_fit, newdata= )
```

What if everyone was **not treated**?

$$Q(0, W) = E[Y|A = 0, W]$$

```
  <- predict(outcome_fit, newdata= )
```

## Step 1: Estimate the Outcome

Dataset after Step 1

Y	A	Q_A	Q_0	Q_1
1	1	0.8461853	0.6770917	0.8461853
1	0	0.6986440	0.6986440	0.8589257
0	0	0.4932538	0.4932538	0.7188934
1	1	0.8213403	0.6363132	0.8213403
1	0	0.62666258	0.62666258	0.8151742
1	1	0.8578239	0.69666588	0.8578239

- $Q_A$ = predicted probability of outcome under observed value of A
- $Q_0$ = predicted probability of outcome under  $A=0$
- $Q_1$ = predicted probability of outcome under  $A=1$

## Step 1: Estimate the Outcome

We could estimate the ATE using G-computation at this point

- Introduced by Christian yesterday
- Can get 95% CI's via bootstrap
  - Computationally expensive with SuperLearner

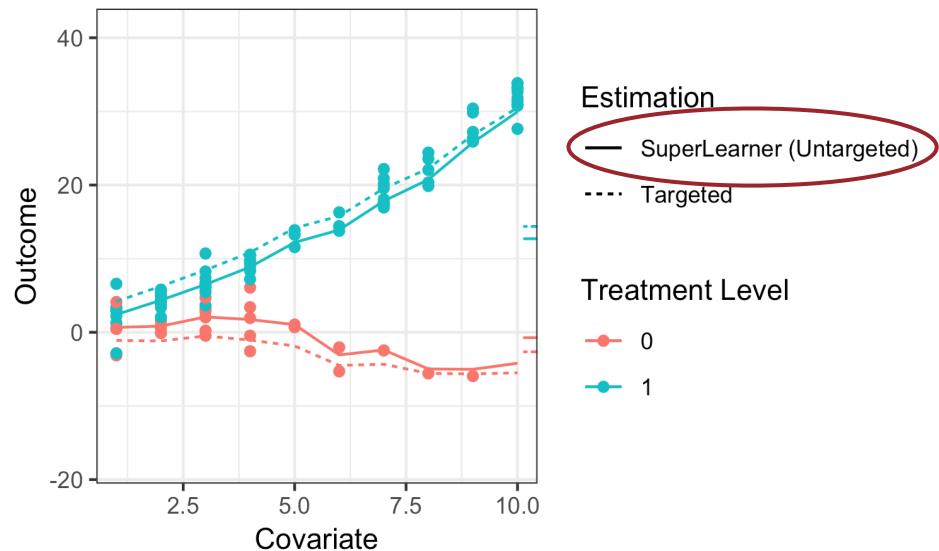
```
ATE  
G-comp <- mean(  -  )
```

$$\text{ATE} = \mathbb{E}[Y_1] - \mathbb{E}[Y_0] = \frac{\sum_{i=1}^N (\mathbb{E}[Y|A = 1, W] - \mathbb{E}[Y|A = 0, W])}{N}$$

## So why not stop here?

- Difference in solid lines on right margin is the G-computation based ATE
- This does not have the more optimal bias/variance tradeoff for the ATE that TMLE
  - Optimized for the outcome prediction
- Does not have the double-robust properties of TMLE

Difference in solid lines on right margin  
is the G-computation based ATE



# TMLE: the algorithm steps

1. Estimate the Outcome
2. Estimate the Probability of Treatment
3. Estimate the Fluctuation Parameter
4. Update the Initial Estimates of the Expected Outcome
5. Compute the Statistical Estimand of Interest
6. Calculate the Standard Errors for Confidence Intervals and P-values

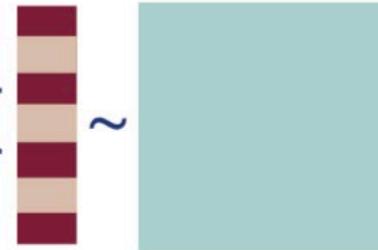
## Step 2: Estimate the Probability of Treatment

Predict the probability of treatment (the propensity score)

$$g(A) = \Pr(A|W)$$

Estimated with new SuperLearner model

```
treatment_fit <- fit(
```



```
)
```

## Step 2: Estimate the Probability of Treatment

Now we need to compute three different quantities from this model fit:

1. The inverse **probability of receiving treatment**

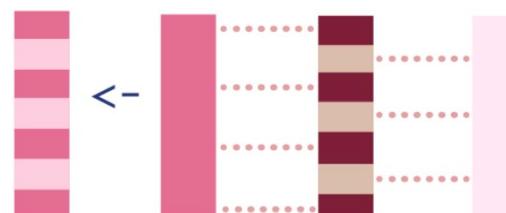
```
<- 1/predict(treatment_fit)
```

2. The negative inverse **probability of not receiving treatment**

```
<- -1/(1-predict(treatment_fit))
```

3. Combine these:

- If the observation was **treated**, the inverse **probability of receiving treatment**
- If they were **not treated**, the negative **inverse probability of not receiving treatment**
- In the TMLE literature this is called the **clever covariate**



## Step 2: Estimate the Probability of Treatment

Dataset after Step 2

Y	A	Q_A	Q_0	Q_1	H_1	H_0	H_A
1.00	1.00	0.85	0.68	0.85	2.18	-1.85	2.18
1.00	0.00	0.70	0.70	0.86	1.60	-2.67	-2.67
0.00	0.00	0.49	0.49	0.72	3.39	-1.42	-1.42
1.00	1.00	0.82	0.64	0.82	2.39	-1.72	2.39
1.00	0.00	0.63	0.63	0.82	2.31	-1.76	-1.76
1.00	1.00	0.86	0.70	0.86	2.14	-1.88	2.14

# TMLE: the algorithm steps

1. Estimate the Outcome
2. Estimate the Probability of Treatment
3. Estimate the Fluctuation Parameter
4. Update the Initial Estimates of the Expected Outcome
5. Compute the Statistical Estimand of Interest
6. Calculate the Standard Errors for Confidence Intervals and P-values

## Step 3: Estimate the Fluctuation Parameter

We will now use information about the treatment mechanism (from Step 2) to **optimize the bias-variance trade-off for the ATE** (rather than the outcome) so we can obtain more valid inference.

**Warning:** this step is easy to code, but difficult to understand.

- Requires knowledge of semi-parametric estimation theory.
- We will first work through **how** to do this
- Later we will seek to understand **why** this helps

```
update <- glm( [redacted] ~ -1+offset(qlogis([redacted]))+[redacted], family=binomial)
```

## Step 3: Estimate the Fluctuation Parameter

The point of this step is to solve an estimating equation for the efficient influence function (EIF) of our estimand of interest.

Skipping details of EIF or estimating equations, know that the help us:

- Update our initial outcome estimates so that our estimate of the ATE is asymptotically unbiased (under certain conditions)
- Calculate the variance (allowing for non-bootstrapped 95% CI's)

```
update <- glm( [redacted] ~ -1+offset(qlogis([redacted]))+[redacted], family=binomial)
```

## Step 3: Estimate the Fluctuation Parameter

- The model that will help us solve an EIF estimating equation and then update our estimates:

$$\text{logit}(E[Y|A, W]) = \text{logit}(\hat{E}[Y|A, W]) + \epsilon H(A, W)$$

Our estimating equation looks *a lot* like a simple logistic regression:

$$\text{logit}(E[Y|X]) = \beta_0 + \beta_1 X.$$

- With the clever covariate  $H(A, W)$  as an added offset (or fixed intercept) with coefficient  $\epsilon$

```
update <- glm( [redacted] ~ -1+offset(qlogis([redacted]))+[redacted], family=binomial)
```

## Step 3: Estimate the Fluctuation Parameter

```
update <- glm(  ~ -1+offset(qlogis())+, family=binomial)
```

- Save the coefficient  $\epsilon$  from this model (called the fluctuation parameter) for the next step:

```
 <- coef(update)
```

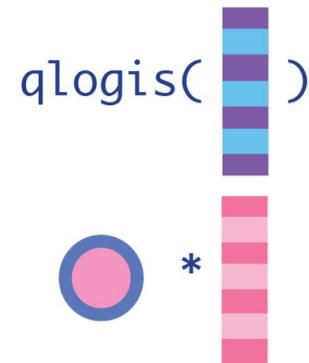
# TMLE: the algorithm steps

1. Estimate the Outcome
2. Estimate the Probability of Treatment
3. Estimate the Fluctuation Parameter
4. Update the Initial Estimates of the Expected Outcome
5. Compute the Statistical Estimand of Interest
6. Calculate the Standard Errors for Confidence Intervals and P-values

## Step 4: Update the Initial Estimates of the Expected Outcome

Calculate how much we need to fluctuate our initial predictions of  $E[Y|A,W]$  to optimize the estimation of the ATE rather than the outcome

1. Put the initial outcome predictions on the logit scale because that was the scale we solved the EIF on (`qlogis()` in R).
2. Multiple the outputs of steps 2 and 3 (the **clever covariate** and **fluctuation parameter**)
3. Sum these 2 quantiles and then transform back to the true outcome scale with the `expit` function (`plogis()` in R):



```
<- plogis(qlogis( ) +  *  )
```

## Step 4: Update the Initial Estimates of the Expected Outcome

Update all 3 predictions

1. Update the expected outcomes of all observations, given the treatment they actually received and their baseline confounders.


$$\text{observed_outcomes} \leftarrow \text{plogis}(\text{qlogis}(\text{observed_outcomes}) + \text{treatment} * \text{baseline_cofounders})$$

2. Update the expected outcomes, conditional on baseline confounders and everyone receiving the treatment.


$$\text{expected_outcomes} \leftarrow \text{plogis}(\text{qlogis}(\text{observed_outcomes}) + \text{treatment} * \text{baseline_cofounders})$$

3. Update the expected outcomes, conditional on baseline confounders and no one receiving the treatment.


$$\text{expected_outcomes} \leftarrow \text{plogis}(\text{qlogis}(\text{observed_outcomes}) + \text{treatment} * \text{baseline_cofounders})$$

# TMLE: the algorithm steps

1. Estimate the Outcome
2. Estimate the Probability of Treatment
3. Estimate the Fluctuation Parameter
4. Update the Initial Estimates of the Expected Outcome
5. Compute the Statistical Estimand of Interest
6. Calculate the Standard Errors for Confidence Intervals and P-values

## Step 5: Compute the Statistical Estimand of Interest

Calculate the more optimized ATE

```
ATE  
TMLE <- mean(  -  )
```

# TMLE: the algorithm steps

1. Estimate the Outcome
2. Estimate the Probability of Treatment
3. Estimate the Fluctuation Parameter
4. Update the Initial Estimates of the Expected Outcome
5. Compute the Statistical Estimand of Interest
6. Calculate the Standard Errors for Confidence Intervals and P-values

# Step 6: Calculate the Standard Errors for Confidence Intervals and P-values

## Compute the Influence Function (IF)

- This is the empirical version of the EIF we used our estimating equation to figure out in Step 3
- The IF shows how much each observation influences the final estimate
- The 95% confidence intervals can be then calculated as normal from the SE
  - These are called influence curve based confidence intervals

$$\hat{IF} = (Y - \hat{E}^*[Y|A, \mathbf{W}])H(A, \mathbf{W}) + \hat{E}^*[Y|A = 1, \mathbf{W}] - \hat{E}^*[Y|A = 0, \mathbf{W}] - \hat{ATE}$$

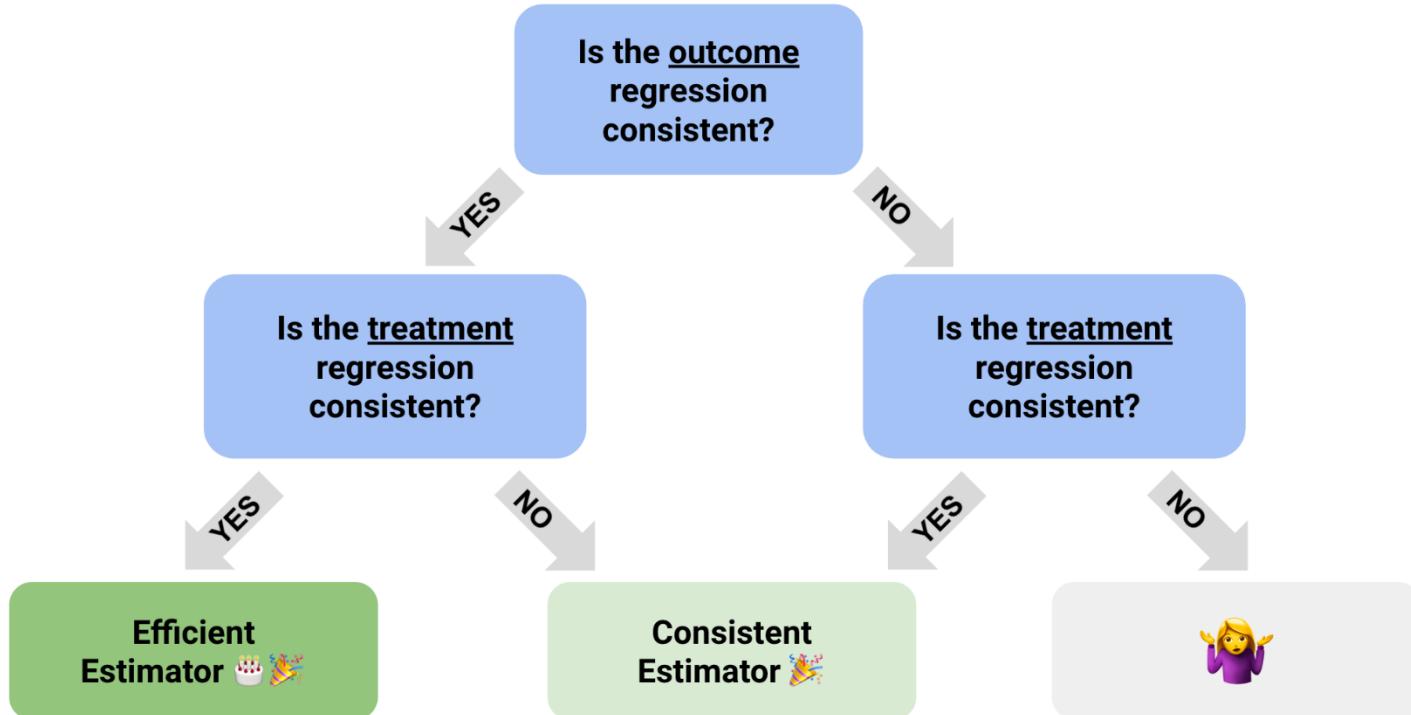
```
st_error <- sqrt(var(( [REDACTED] - [REDACTED] ) * [REDACTED] + [REDACTED] - [REDACTED] - [REDACTED] ATE TMLE ) / N))
```

# Properties of TMLE

*So why did we go through all that effort?*

- Can incorporate Machine Learning
  - G-computation and IPTW approaches also allow for machine learning for inference
    - Methods exist for inference, including bootstrapping
    - But both are susceptible to model misspecification, especially if using parametric methods
- Efficient
  - Lowest (asymptotic) variance among reasonable estimators if both  $\mathbf{Q}(Y|A, W)$  AND  $g(A|W)$  estimated consistently at reasonable rates
- Substitution (aka "plug in") Estimator
  - Its estimates will always stay within the bounds of the original outcome
  - Improved robustness to sparse data compared to estimating equation alternatives
- Double Robust
  - Consistent if either  $\mathbf{Q}(Y|A, W)$  or  $g(A|W)$  estimated consistently

## Determining your Doubly Robust Estimator's Properties



[khstats.com/blog/tmle/tutorial](http://khstats.com/blog/tmle/tutorial)

**Consistent:** bias decreases to zero as sample size grows large

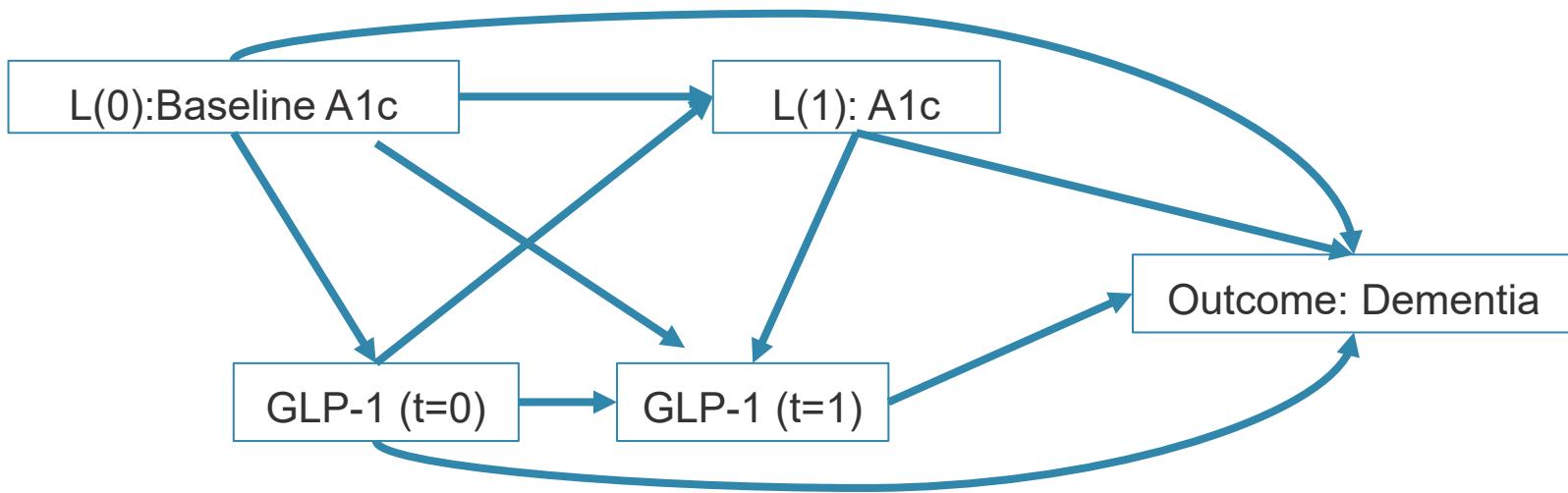
**Efficient:** bias decreases to zero *and* variance decreases to smallest possible value for a fixed number of observations

# TMLE Extensions

1. Longitudinal TMLE
2. Population and Optimal intervention effects
3. Stochastic interventions

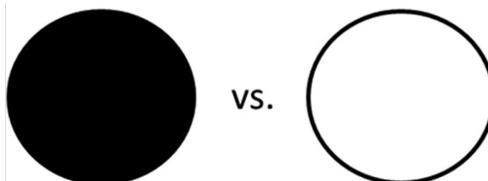
# Longitudinal TMLE

- Coming Thursday!
- Allows for untangling time-dependent confounding through iterated Q and g models, and retains all the TMLE properties



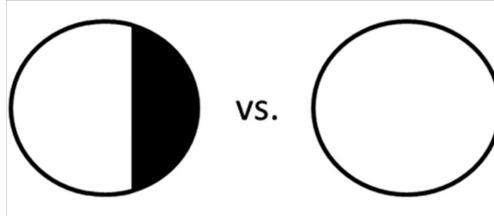
# Targeted learning to directly estimate the parameters we care about

Average Treatment Effect



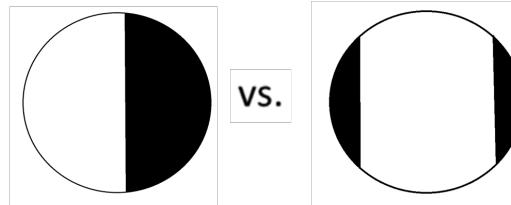
$$\varphi = E(Y_1 - Y_0) = E[Y|A = 1, W] - E[Y|A = 0, W]$$

Population Intervention Impact



$$\varphi_{PI} = E(Y - Y_0) = E[Y] - E[Y|A = 0, W]$$

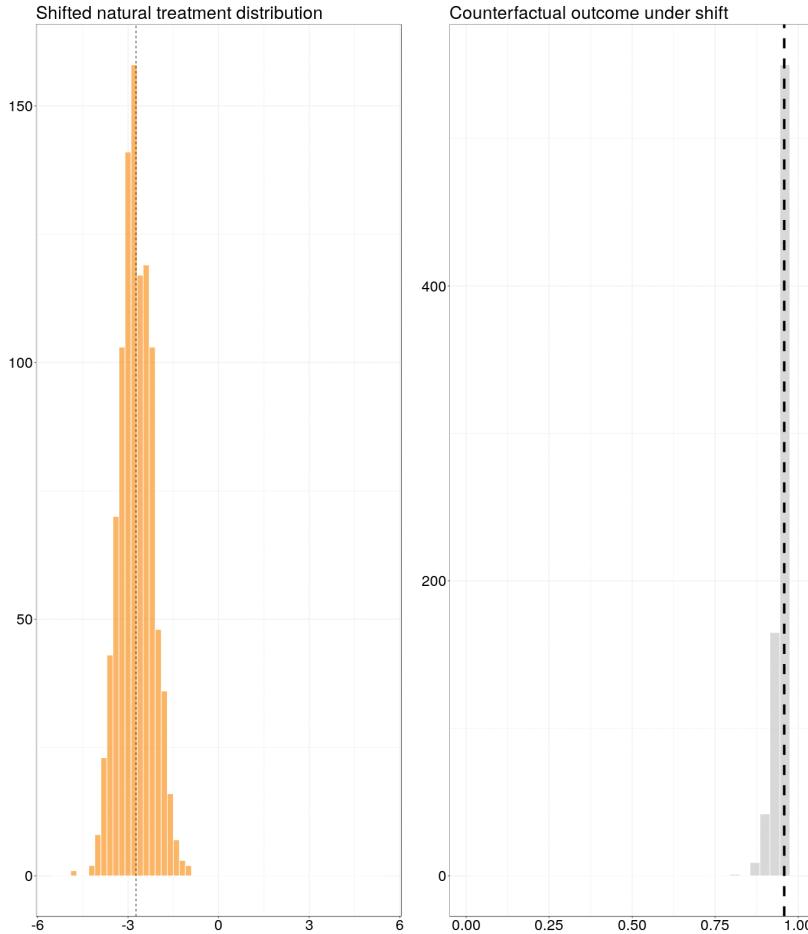
Optimal Dynamic Treatment Rule



$$\varphi_{ODTR} = E(Y - Y_{d_{opt}(W)}) = E[Y] - E[Y|A = d_{opt}(W), W]$$

# Stochastic interventions

- Present a relatively simple yet extremely flexible manner by which *realistic* causal effects (and contrasts thereof) may be defined.
- Allows for estimating the effect of shifting a distribution of a variable
  - Rather than contrasting two specific levels



# Resources

- Some lecture material adapted from
  - Tlverse handbook (<https://tlverse.org/tlverse-handbook/>)
  - Illustrated guide to TMLE (<https://www.khstats.com/blog/tmle/tutorial-pt1>)
  - Slides shared by Maya Petersen and Ben Arnold
- Other good learning references
  - <https://achambaz.github.io/tlide/>
  - <https://multithreaded.stitchfix.com/blog/2021/07/23/double-robust-estimator/>

# References

- H. Bang and J.M. Robins. Doubly-robust estimation in missing data and causal inference models. *Biometrics*, 61:962 – 972, 2005.
- L. Breiman. Bagging predictors. *Machine Learning*, 24:123 – 140, 1996.
- M A Hernán, E Lanoy, D Costagliola, and J M Robins. Comparison of dynamic treatment regimes via inverse probability weighting. *Basic & Clinical Pharmacology & Toxicology*, 98:237 – 242, 2006.
- Judea Pearl. Causal diagrams for empirical research. *Biometrika*, 82(4): 669–688, 1995.
- Judea Pearl and James M Robins. Probabilistic evaluation of sequential plans from causal models with hidden variables. In UAI, volume 95, pages 444–453. Citeseer, 1995.
- Maya L Petersen and Mark J van der Laan. Causal models and learning from data: integrating causal modeling and statistical estimation. *Epidemiology (Cambridge, Mass.)*, 25(3):418, 2014.
- J. Robins and A. Rotnitzky. Recovery of information and adjustment for dependent censoring using surrogate markers. In AIDS Epidemiology, pages 297–331. Springer, 1992.

## References

- James Robins. The control of confounding by intermediate variables.  
Statistics in medicine, 8(6):679–701, 1989.
- J.M. Robins. A new approach to causal inference in mortality studies with sustained exposure periods – application to control of the healthy worker survivor effect. Mathematical Modelling, 7:1393 – 1512, 1986.
- J.M. Robins. Robust estimation in sequentially ignorable missing data and causal inference models. In Proceedings of the American Statistical Association on Bayesian Statistical Science, 1999, pages 6–10, 2000.
- J. Schwab, S. Lendle, M. Petersen, and M. van der Laan. Itmle: Longitudinal Targeted Maximum Likelihood Estimation, 2013. R package version 0.9.3, <http://cran.r-project.org/web/packages/itmle/>.
- L Tran, Constantin T Yiannoutsos, B Musick, K Wools-Kaloustian, A Siika, S Kimaiyo, M van der Laan, and M Petersen. Evaluating the impact of a hiv low-risk express care task-shifting program: a case study of the targeted learning roadmap. Epidemiologic Methods, 5(1):69–91, 2016.

## References

- M.J. van der Laan and S. Gruber. Targeted minimum loss based estimation of causal effects of multiple time point interventions. *The International Journal of Biostatistics*, 8(1):Article 8, 2012.
- MJ van der Laan and S Rose. Targeted learning: causal inference for observational and experimental data. Springer Science & Business Media, 2011.
- M.J. van der Laan, E. Polley, and A. Hubbard. Super learner. *Statistical Applications in Genetics and Molecular Biology*, 6(25), 2007.