

Targeted Learning

Mark van der Laan

Human Art in Statistics

Role of Targeted Learning in Data Science

Roadmap for Targeted Learning

Theoretical Underpinnings

Adaptive Experimental Designs

Online Learning

Future of Targeted Learning

Targeted Learning

The bridge from machine learning to statistical and causal inference

Mark van der Laan

Jiann-Ping Hsu/Karl E. Peace Professor in Biostatistics & Statistics
University of California, Berkeley

September 23, 2020, Sentinel Seminar

Acknowledgements: Rachael Phillips, Ivana Malenica,
Chris Kennedy, Aurelien Bibaut, Nima Hejazi and Jonathan Levy

Traditional toolbox for statistics

Targeted Learning

Mark van der Laan

Human Art in Statistics

Role of Targeted Learning in Data Science

Roadmap for Targeted Learning

Theoretical Underpinnings

Adaptive Experimental Designs

Online Learning

Future of Targeted Learning

Goal	Type of Data			
	Measurement (from Gaussian Population)	Rank, Score, or Measurement (from Non-Gaussian Population)	Binomial (Two Possible Outcomes)	Survival Time
Describe one group	Mean, SD	Median, interquartile range	Proportion	Kaplan Meier survival curve
Compare one group to a hypothetical value	One-sample t test	Wilcoxon test	Chi-square or Binomial test**	
Compare two unpaired groups	Unpaired t test	Mann-Whitney test	Fisher's test (chi-square for large samples)	Log-rank test or Mantel-Haenszel*
Compare two paired groups	Paired t test	Wilcoxon test	McNemar's test	Conditional proportional hazards regression*
Compare three or more unmatched groups	One-way ANOVA	Kruskal-Wallis test	Chi-square test	Cox proportional hazard regression**
Compare three or more matched groups	Repeated-measures ANOVA	Friedman test	Cochrane Q**	Conditional proportional hazards regression**
Quantify association between two variables	Pearson correlation	Spearman correlation	Contingency coefficients**	
Predict value from another measured variable	Simple linear regression or Nonlinear regression	Nonparametric regression**	Simple logistic regression*	Cox proportional hazard regression*
Predict value from several measured or binomial variables	Multiple linear regression* or Multiple nonlinear regression**		Multiple logistic regression*	Cox proportional hazard regression*

Performance of traditional tools

Targeted Learning

Mark van der Laan

Human Art in Statistics

Role of Targeted Learning in Data Science

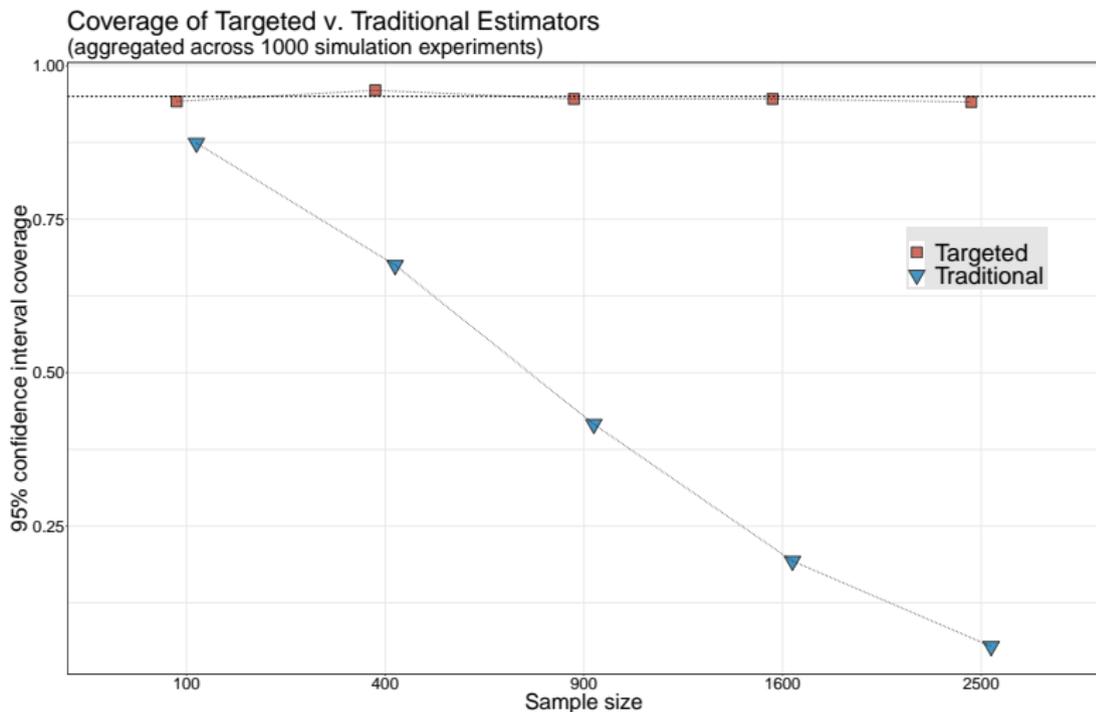
Roadmap for Targeted Learning

Theoretical Underpinnings

Adaptive Experimental Designs

Online Learning

Future of Targeted Learning



Performance of traditional tools

Targeted Learning

Mark van der Laan

Human Art in Statistics

Role of Targeted Learning in Data Science

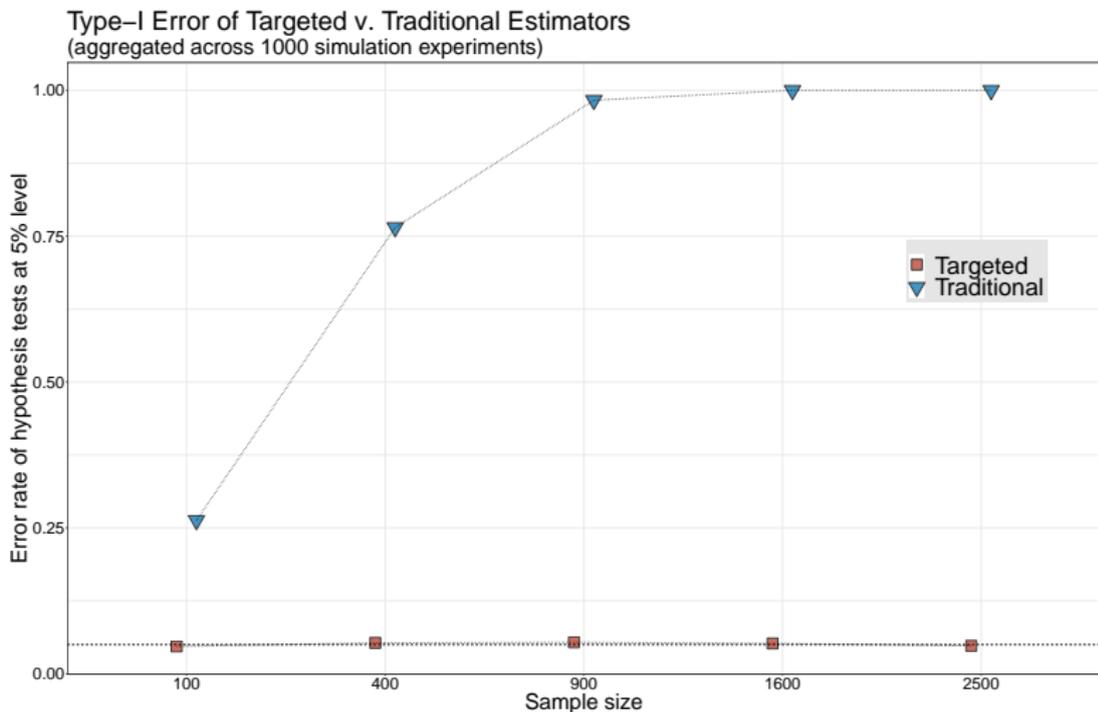
Roadmap for Targeted Learning

Theoretical Underpinnings

Adaptive Experimental Designs

Online Learning

Future of Targeted Learning



Post-hoc model manipulation

Targeted Learning

Mark van der Laan

Human Art in Statistics

Role of Targeted Learning in Data Science

Roadmap for Targeted Learning

Theoretical Underpinnings

Adaptive Experimental Designs

Online Learning

Future of Targeted Learning



Why care about statistical inference?

Targeted Learning

Mark van der Laan

Human Art in Statistics

Role of Targeted Learning in Data Science

Roadmap for Targeted Learning

Theoretical Underpinnings

Adaptive Experimental Designs

Online Learning

Future of Targeted Learning

Why Most Published Research Findings Are False

John P. A. Ioannidis

False-Positive Psychology: Undisclosed Flexibility in Data Collection and Analysis Allows Presenting Anything as Significant

Joseph P. Simmons¹, Leif D. Nelson², and Uri Simonsohn¹

¹The Wharton School, University of Pennsylvania, and ²Haas School of Business, University of California, Berkeley

The Statistical Crisis in Science

Data-dependent analysis—a “garden of forking paths”—explains why many statistically significant comparisons don’t hold up.

Andrew Gelman and Eric Loken

Targeted Learning for answering statistical and causal questions with confidence intervals

Targeted Learning

Mark van der Laan

Human Art in Statistics

Role of Targeted Learning in Data Science

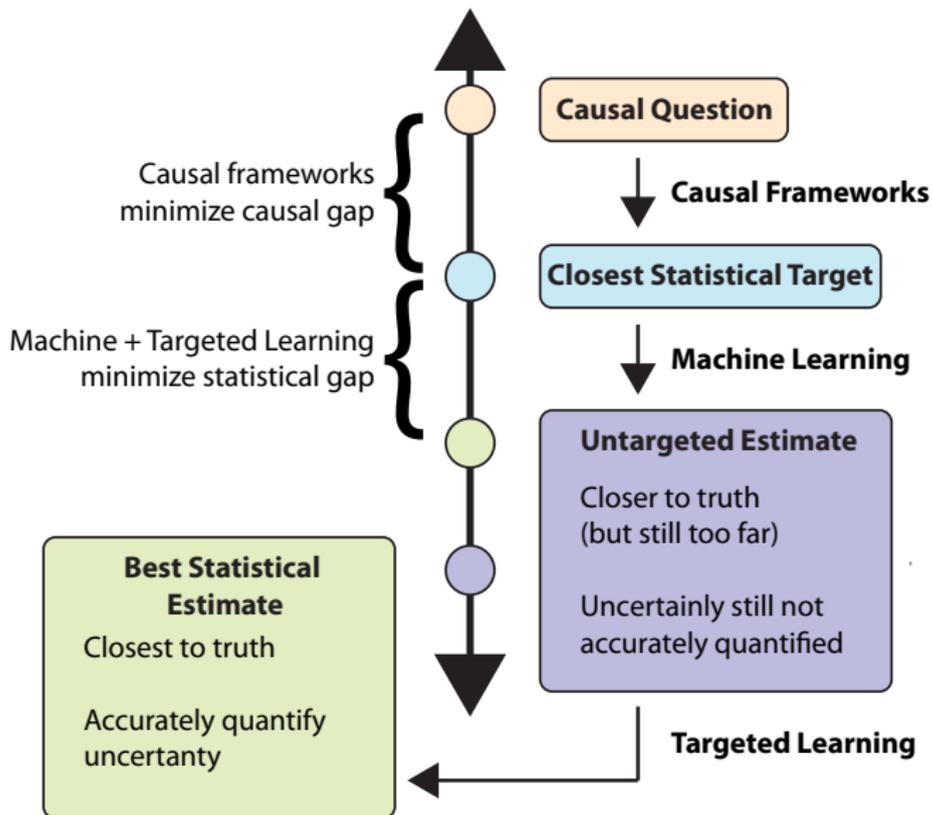
Roadmap for Targeted Learning

Theoretical Underpinnings

Adaptive Experimental Designs

Online Learning

Future of Targeted Learning



Targeted Learning is a subfield of statistics

Targeted Learning

Mark van der Laan

Human Art in Statistics

Role of Targeted Learning in Data Science

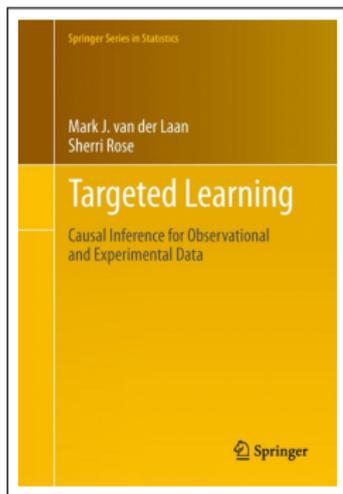
Roadmap for Targeted Learning

Theoretical Underpinnings

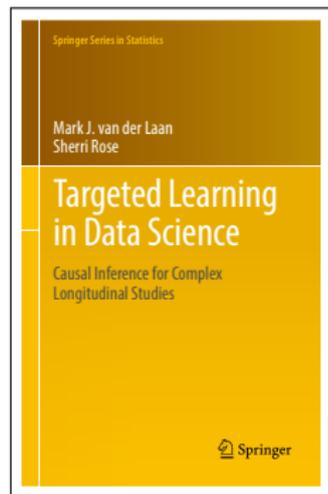
Adaptive Experimental Designs

Online Learning

Future of Targeted Learning



van der Laan & Rose, *Targeted Learning: Causal Inference for Observational and Experimental Data*. New York: Springer, 2011.



van der Laan & Rose, *Targeted Learning in Data Science: Causal Inference for Complex Longitudinal Studies*. New York: Springer, 2018.

The Hitchhiker's Guide to the tLverse

Better clinical decisions from observational data

Targeted Learning

Mark van der Laan

Human Art in Statistics

Role of Targeted Learning in Data Science

Roadmap for Targeted Learning

Theoretical Underpinnings

Adaptive Experimental Designs

Online Learning

Future of Targeted Learning

Statistics
in Medicine

Research Article

Received 24 May 2013,

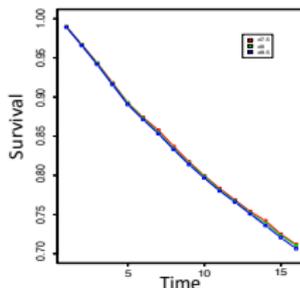
Accepted 5 January 2014

Published online 17 February 2014 in Wiley Online Library

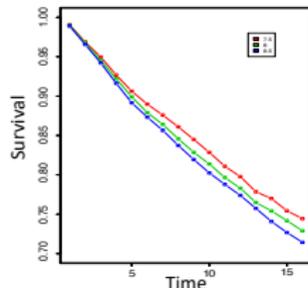
(wileyonlinelibrary.com) DOI: 10.1002/sim.6099

Targeted learning in real-world comparative effectiveness research with time-varying interventions

Romain Neugebauer,^{a,*†} Julie A. Schmittiel^a and Mark J. van der Laan^b



Standard methods: No benefit to more aggressive intensification strategy



Targeted Learning: More aggressive intensification protocols result in better outcomes

The roadmap for statistical learning

Targeted Learning

Mark van der Laan

Human Art in Statistics

Role of Targeted Learning in Data Science

Roadmap for Targeted Learning

Theoretical Underpinnings

Adaptive Experimental Designs

Online Learning

Future of Targeted Learning

STEP 1:
DESCRIBE
EXPERIMENT

STEP 2:
SPECIFY
STATISTICAL MODEL

STEP 3:
DEFINE
STATISTICAL QUERY

STEP 4:
CONSTRUCT
ESTIMATOR

STEP 5:
OBTAIN
INFERENCE

What is the experiment that generated the data?

**STEP 1:
DESCRIBE
EXPERIMENT**

STEP 2:
SPECIFY
STATISTICAL MODEL

STEP 3:
DEFINE
STATISTICAL QUERY

STEP 4:
CONSTRUCT
ESTIMATOR

STEP 5:
OBTAIN
INFERENCE

***Three multi-national RCTs assessing
impact of corticosteroids on mortality
among septic shock patients***

Targeted
Learning

Mark van der
Laan

Human Art in
Statistics

Role of
Targeted
Learning in
Data Science

Roadmap for
Targeted
Learning

Theoretical
Underpinnings

Adaptive
Experimental
Designs

Online
Learning

Future of
Targeted
Learning

What is the experiment that generated the data?

Targeted Learning

Mark van der Laan

Human Art in Statistics

Role of Targeted Learning in Data Science

Roadmap for Targeted Learning

Theoretical Underpinnings

Adaptive Experimental Designs

Online Learning

Future of Targeted Learning

**STEP 1:
DESCRIBE
EXPERIMENT**

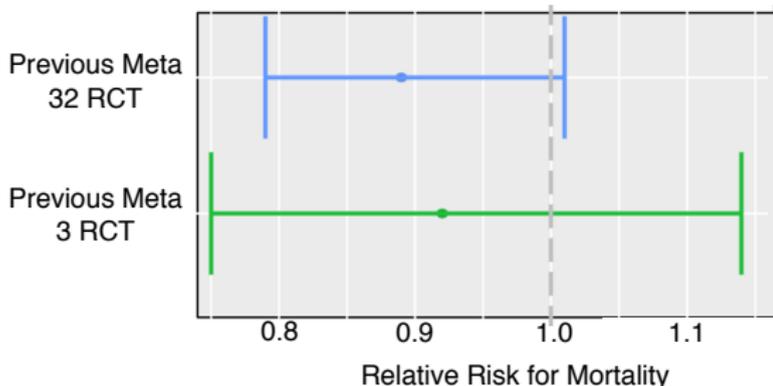
STEP 2:
SPECIFY
STATISTICAL MODEL

STEP 3:
DEFINE
STATISTICAL QUERY

STEP 4:
CONSTRUCT
ESTIMATOR

STEP 5:
OBTAIN
INFERENCE

Three multi-national RCTs assessing impact of corticosteroids on mortality among septic shock patients



What is the experiment that generated the data?

Targeted Learning

Mark van der Laan

Human Art in Statistics

Role of Targeted Learning in Data Science

Roadmap for Targeted Learning

Theoretical Underpinnings

Adaptive Experimental Designs

Online Learning

Future of Targeted Learning

**STEP 1:
DESCRIBE
EXPERIMENT**

**STEP 2:
SPECIFY
STATISTICAL MODEL**

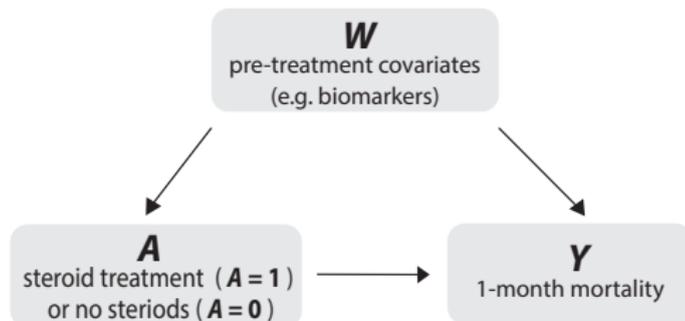
**STEP 3:
DEFINE
STATISTICAL QUERY**

**STEP 4:
CONSTRUCT
ESTIMATOR**

**STEP 5:
OBTAIN
INFERENCE**

Three multi-national RCTs assessing impact of corticosteroids on mortality among septic shock patients

Pooled sample of $n = 1,300$ adults in septic shock



What is known about stochastic relations for the observed variables?

Targeted Learning

Mark van der Laan

Human Art in Statistics

Role of Targeted Learning in Data Science

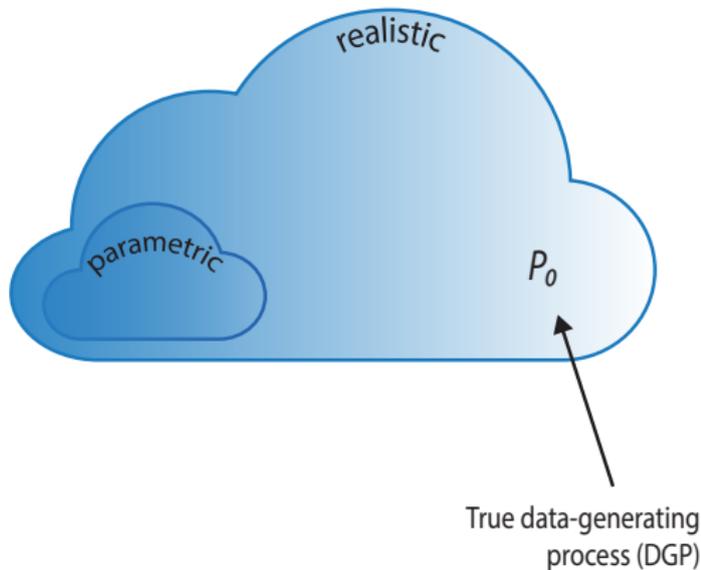
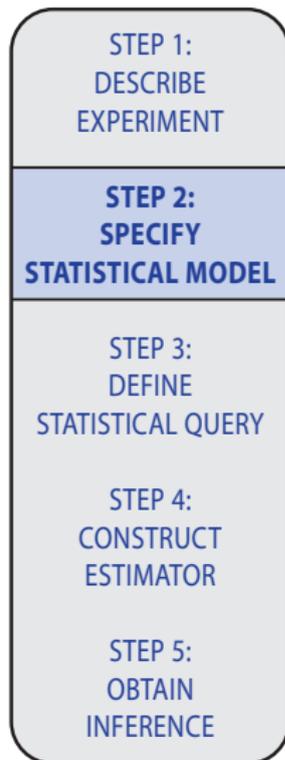
Roadmap for Targeted Learning

Theoretical Underpinnings

Adaptive Experimental Designs

Online Learning

Future of Targeted Learning



What is the target estimand that we want to learn from the data?

Targeted Learning

Mark van der Laan

Human Art in Statistics

Role of Targeted Learning in Data Science

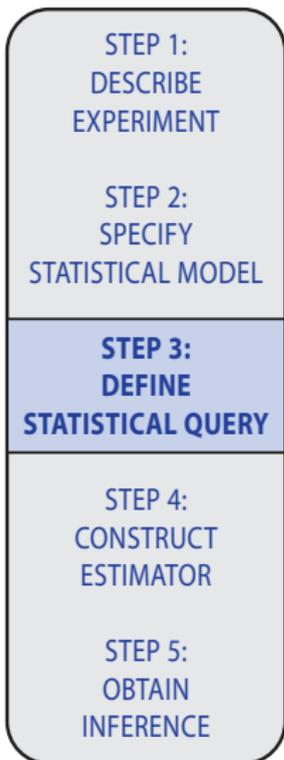
Roadmap for Targeted Learning

Theoretical Underpinnings

Adaptive Experimental Designs

Online Learning

Future of Targeted Learning



What is the average difference in mortality between treatment groups when adjusting for covariates?

$$\Psi(P_0) = E_0(E_0[Y|A = 1, W] - E_0[Y|A = 0, W])$$

How should we estimate the target estimand?

Targeted Learning

Mark van der Laan

Human Art in Statistics

Role of Targeted Learning in Data Science

Roadmap for Targeted Learning

Theoretical Underpinnings

Adaptive Experimental Designs

Online Learning

Future of Targeted Learning

STEP 1:
DESCRIBE
EXPERIMENT

STEP 2:
SPECIFY
STATISTICAL MODEL

STEP 3:
DEFINE
STATISTICAL QUERY

STEP 4:
**CONSTRUCT
ESTIMATOR**

STEP 5:
OBTAIN
INFERENCE

TARGETED MAXIMUM LIKELIHOOD ESTIMATION

- 1 Initial estimation of $E_0[Y|A, W]$ with super (machine) learning
- 2 Updating initial estimate to achieve optimal bias-variance trade-off for $\Psi(P_0)$

TMLE estimates are optimal:

plug-in, efficient, unbiased, finite sample robust

How should we approximate the sampling distribution of our estimator?

Targeted Learning

Mark van der Laan

Human Art in Statistics

Role of Targeted Learning in Data Science

Roadmap for Targeted Learning

Theoretical Underpinnings

Adaptive Experimental Designs

Online Learning

Future of Targeted Learning

STEP 1:
DESCRIBE
EXPERIMENT

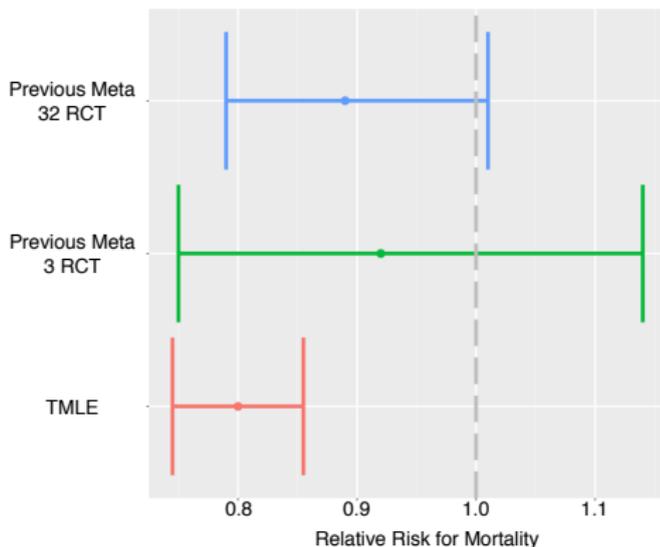
STEP 2:
SPECIFY
STATISTICAL MODEL

STEP 3:
DEFINE
STATISTICAL QUERY

STEP 4:
CONSTRUCT
ESTIMATOR

STEP 5:
OBTAIN
INFERENCE

Due to targeting (step ②), the TMLE behaves as the **sample mean** of efficient influence function



Effect of targeting on sampling distribution

Targeted Learning

Mark van der Laan

Human Art in Statistics

Role of Targeted Learning in Data Science

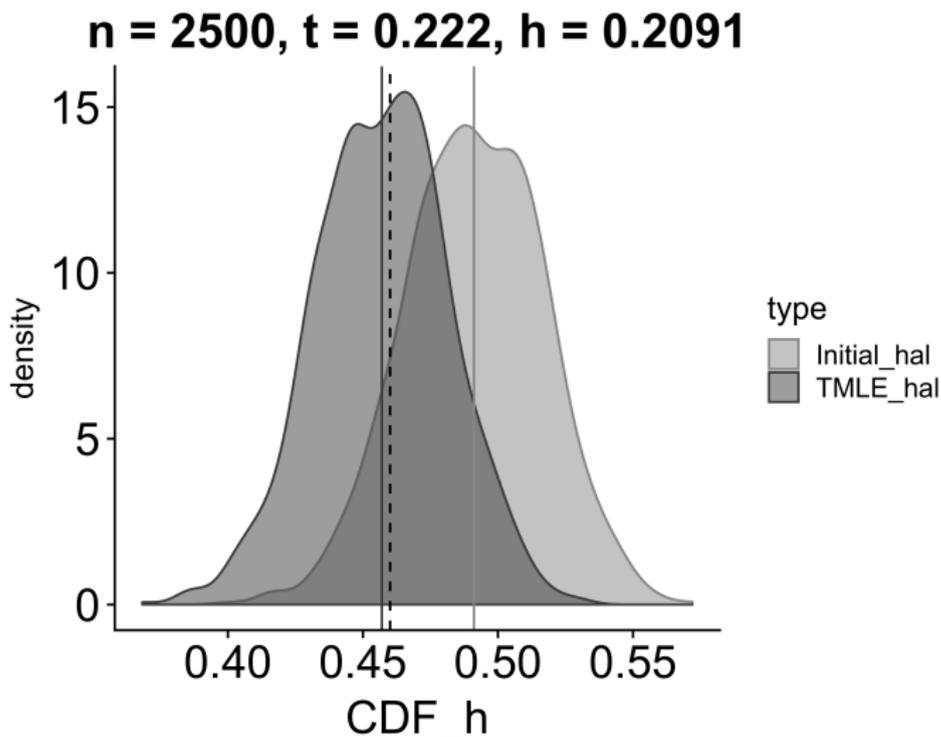
Roadmap for Targeted Learning

Theoretical Underpinnings

Adaptive Experimental Designs

Online Learning

Future of Targeted Learning



Truth = smoothed $\Pr(\text{TE}(W) \leq t)$ at dashed line

What is the optimal rule for assigning corticosteroids to patients in septic shock?

Targeted Learning

Mark van der Laan

Human Art in Statistics

Role of Targeted Learning in Data Science

Roadmap for Targeted Learning

Theoretical Underpinnings

Adaptive Experimental Designs

Online Learning

Future of Targeted Learning

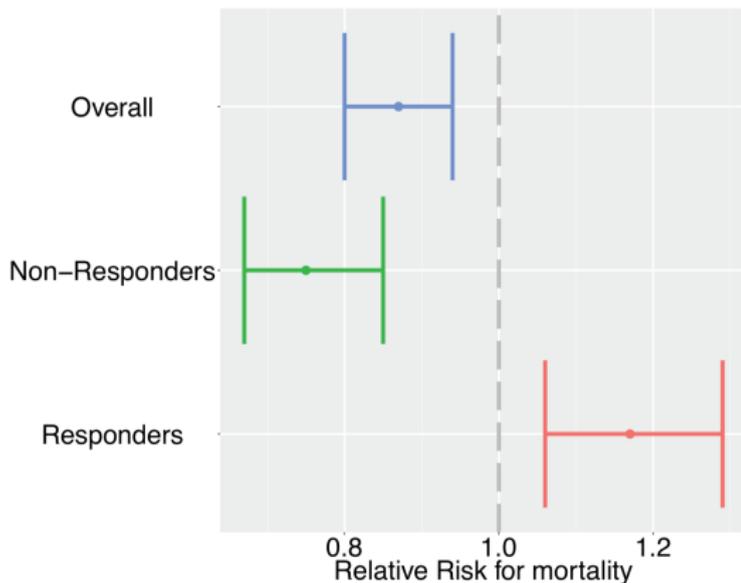
STEP 1:
DESCRIBE
EXPERIMENT

STEP 2:
SPECIFY
STATISTICAL MODEL

STEP 3:
DEFINE
STATISTICAL QUERY

STEP 4:
CONSTRUCT
ESTIMATOR

STEP 5:
OBTAIN
INFERENCE



Super learner

Targeted Learning

Mark van der Laan

Human Art in Statistics

Role of Targeted Learning in Data Science

Roadmap for Targeted Learning

Theoretical Underpinnings

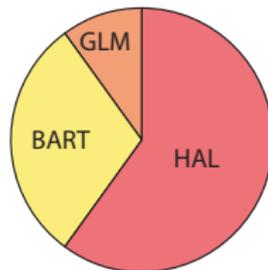
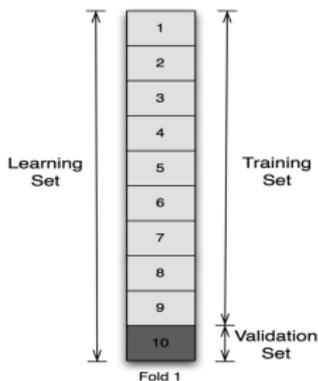
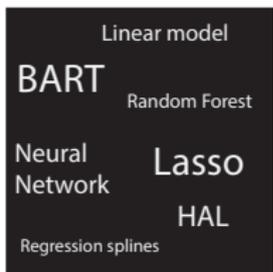
Adaptive Experimental Designs

Online Learning

Future of Targeted Learning



Cross-validated performance of learners + ensembles



Oracle inequality tells us cross-validation is optimal for selection among estimators

Super learner performance in practice

Targeted Learning

Mark van der Laan

Human Art in Statistics

Role of Targeted Learning in Data Science

Roadmap for Targeted Learning

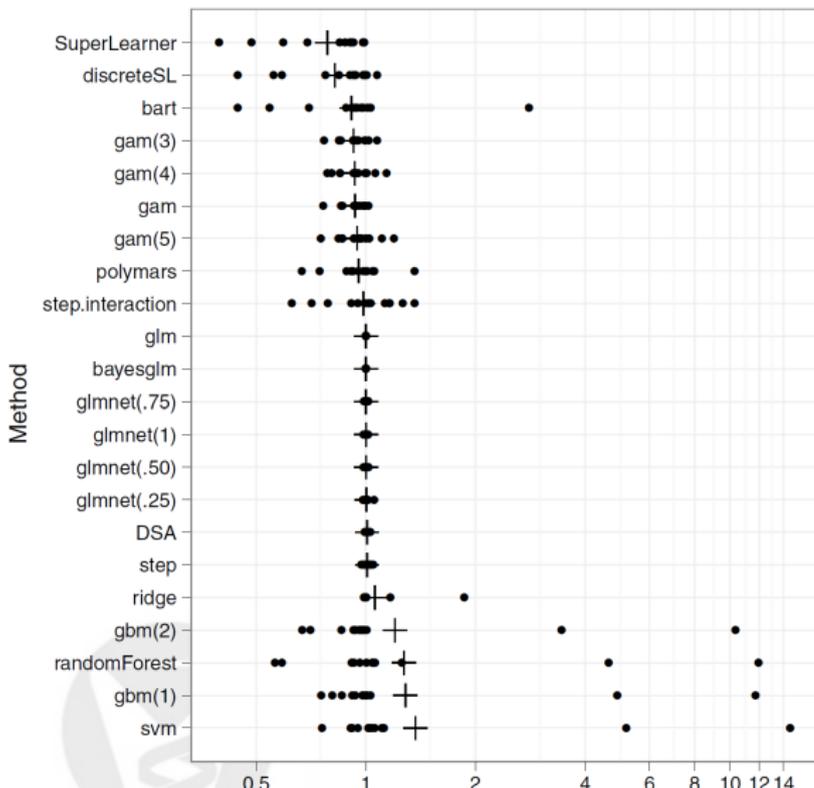
Theoretical Underpinnings

Adaptive Experimental Designs

Online Learning

Future of Targeted Learning

SuperLearner comparison.png



Highly Adaptive Lasso (HAL)

Targeted Learning

Mark van der Laan

Human Art in Statistics

Role of Targeted Learning in Data Science

Roadmap for Targeted Learning

Theoretical Underpinnings

Adaptive Experimental Designs

Online Learning

Future of Targeted Learning

Key Idea

- Any d -dimensional cadlag function (i.e. right-continuous) can be represented as a possibly infinite linear combination of spline basis functions.
- The variation norm / complexity of a function is the L_1 -norm of the vector of coefficients.

Converges to true function at rate $n^{-1/3}(\log n)^{d/2}$

HAL performance for d=3

Targeted Learning

Mark van der Laan

Human Art in Statistics

Role of Targeted Learning in Data Science

Roadmap for Targeted Learning

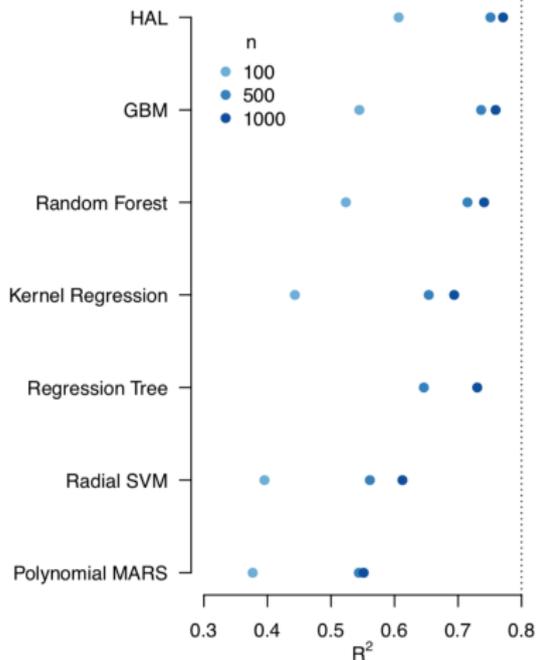
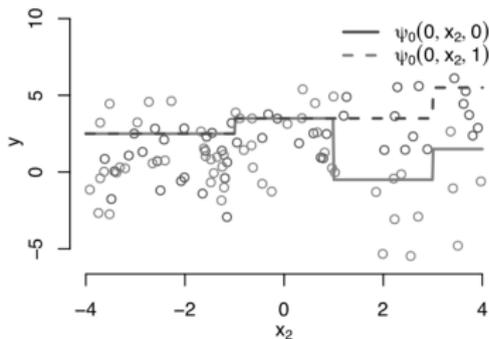
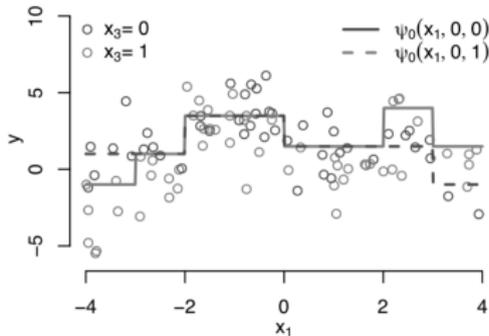
Theoretical Underpinnings

Adaptive Experimental Designs

Online Learning

Future of Targeted Learning

$$\psi_0(x) = -2x_3 I(x_1 < -3) + 2.5 I(x_1 > -2) - 2 I(x_1 > 0) + 2.5x_3 I(x_1 > 2) - 2.5 I(x_1 > 3) + I(x_2 > -1) - 4x_3 I(x_2 > 1) + 2 I(x_2 > 3)$$



HAL metalearner

Targeted Learning

Mark van der Laan

Human Art in Statistics

Role of Targeted Learning in Data Science

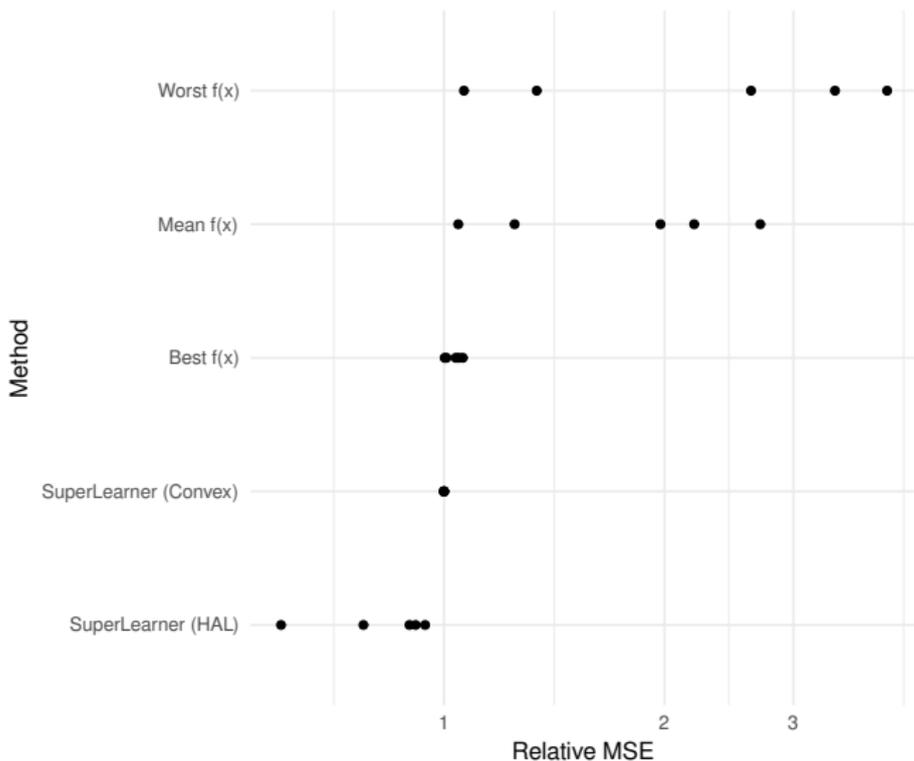
Roadmap for Targeted Learning

Theoretical Underpinnings

Adaptive Experimental Designs

Online Learning

Future of Targeted Learning



TMLE follows a path of maximal change in target estimand per unit of information/likelihood

Targeted Learning

Mark van der Laan

Human Art in Statistics

Role of Targeted Learning in Data Science

Roadmap for Targeted Learning

Theoretical Underpinnings

Adaptive Experimental Designs

Online Learning

Future of Targeted Learning

Can we break HAL-TMLE?

Targeted Learning

Mark van der Laan

Human Art in Statistics

Role of Targeted Learning in Data Science

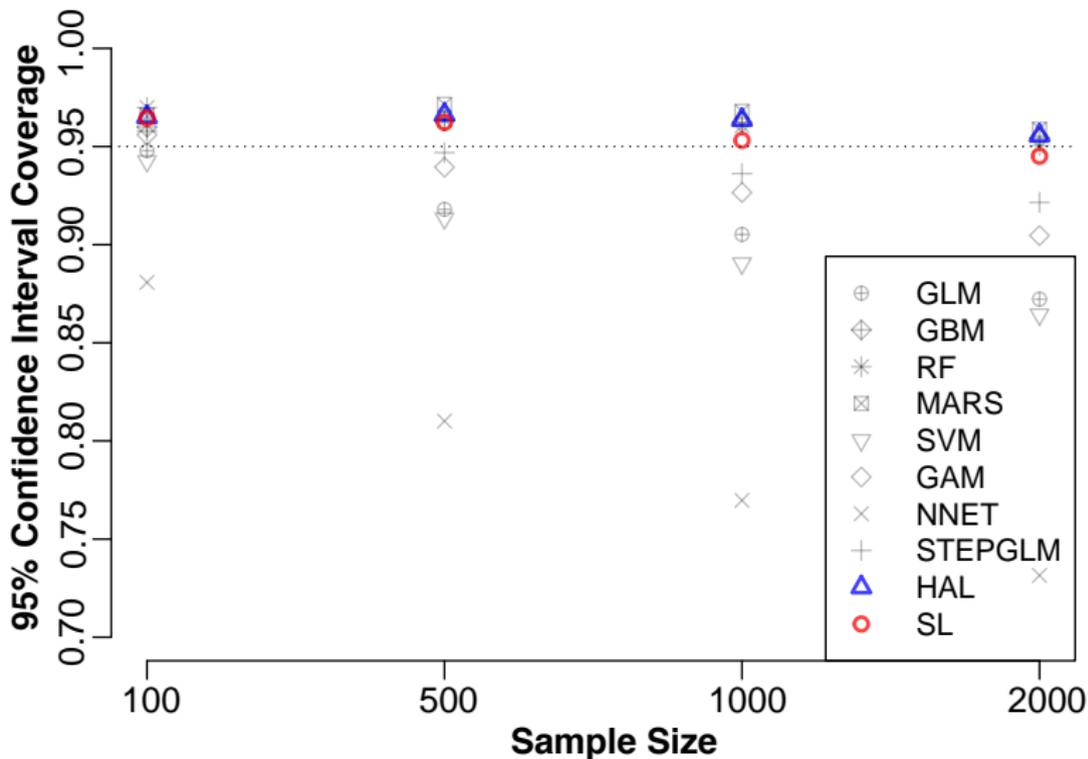
Roadmap for Targeted Learning

Theoretical Underpinnings

Adaptive Experimental Designs

Online Learning

Future of Targeted Learning



Robust inference for adaptive sequential RCTs

Targeted Learning

Mark van der Laan

Human Art in Statistics

Role of Targeted Learning in Data Science

Roadmap for Targeted Learning

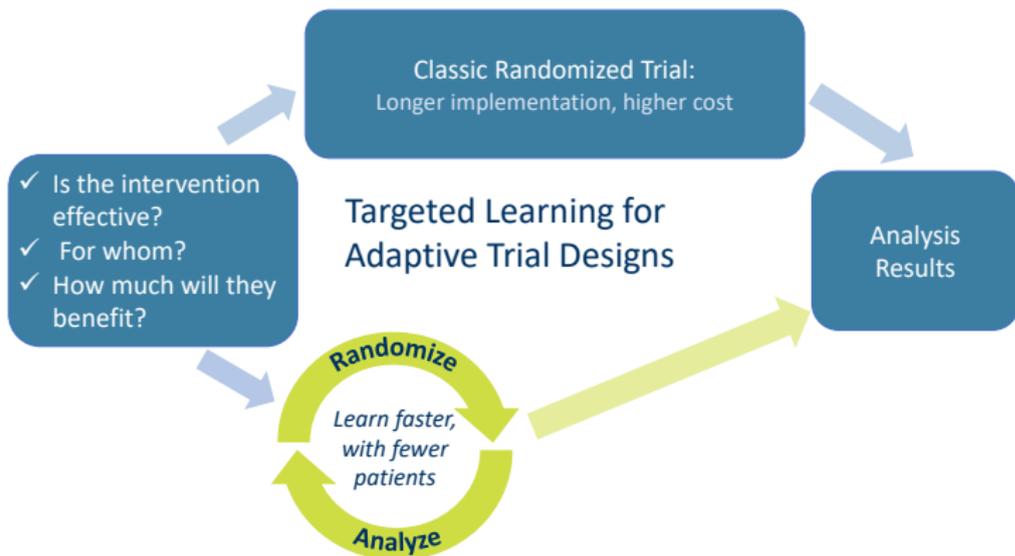
Theoretical Underpinnings

Adaptive Experimental Designs

Online Learning

Future of Targeted Learning

Optimal intervention allocation: “Learn as you go”



Balanced vs. adaptive sequential design

Targeted Learning

Mark van der Laan

Human Art in Statistics

Role of Targeted Learning in Data Science

Roadmap for Targeted Learning

Theoretical Underpinnings

Adaptive Experimental Designs

Online Learning

Future of Targeted Learning

Balanced vs. adaptive sequential design

Targeted Learning

Mark van der Laan

Human Art in Statistics

Role of Targeted Learning in Data Science

Roadmap for Targeted Learning

Theoretical Underpinnings

Adaptive Experimental Designs

Online Learning

Future of Targeted Learning

Online super learning in the ICU

Targeted Learning

Mark van der Laan

Human Art in Statistics

Role of Targeted Learning in Data Science

Roadmap for Targeted Learning

Theoretical Underpinnings

Adaptive Experimental Designs

Online Learning

Future of Targeted Learning

Adaptive algorithm

- Regularly updated with batches of new data
- Learns from both
 - ① within individual time series, and
 - ② across patients
- Uncertainty of forecasts assessed with prediction intervals

15-minute ahead forecasts with prediction intervals for patient with hypotensive episodes

Targeted Learning

Mark van der Laan

Human Art in Statistics

Role of Targeted Learning in Data Science

Roadmap for Targeted Learning

Theoretical Underpinnings

Adaptive Experimental Designs

Online Learning

Future of Targeted Learning

Future of Targeted Learning

Targeted Learning

Mark van der Laan

Human Art in Statistics

Role of Targeted Learning in Data Science

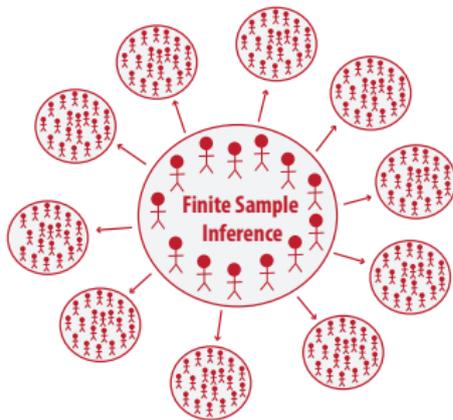
Roadmap for Targeted Learning

Theoretical Underpinnings

Adaptive Experimental Designs

Online Learning

Future of Targeted Learning



 **accenture**

