

Chapter 5

Understanding TMLE

Sherri Rose, Mark J. van der Laan

This chapter focuses on understanding TMLE. We go into more detail than the previous chapter to demonstrate how this estimator is derived. Recall that TMLE is a two-step procedure where one first obtains an estimate of the data-generating distribution P_0 or the relevant portion Q_0 of P_0 . The second stage updates this initial fit in a step targeted toward making an optimal bias–variance tradeoff for the parameter of interest $\Psi(Q_0)$, instead of the overall density P_0 . The procedure is double robust and can incorporate data-adaptive-likelihood-based estimation procedures to estimate Q_0 and the treatment mechanism.

5.1 Conceptual Framework

We begin the discussion of TMLE at a conceptual level to give an overall picture of what the method achieves. In Fig. 5.1 we depict a flow chart for TMLE, and in this section, we walk the reader through the illustration and provide a conceptual foundation for TMLE. We start with our observed data and some (possibly) real valued function $\Psi()$, the target parameter mapping. These two objects are our inputs. We have an initial estimator of the probability distribution of the data (or something smaller than that – the relevant portion). This is P_n^0 and is estimated semiparametrically using super learning. This initial estimator is typically already somewhat informed about the target parameter of interest by, for example, only focusing on fitting the relevant part Q_0 of P_0 . P_n^0 falls within the statistical model, which is the set of all possible probability distributions of the data. P_0 , the true probability distribution, also falls within the statistical model, since it is assumed that the statistical model is selected to represent true knowledge. In many applications the statistical model is necessarily nonparametric. We update P_n^0 in a particular way, in a targeted way by incorporating the target parameter mapping Ψ , and now denote this targeted update as P_n^* . If we map P_n^* using our function $\Psi()$, we get our estimator $\Psi(P_n^*)$ and thereby a value on the real line. The updating step is tailored to result in values

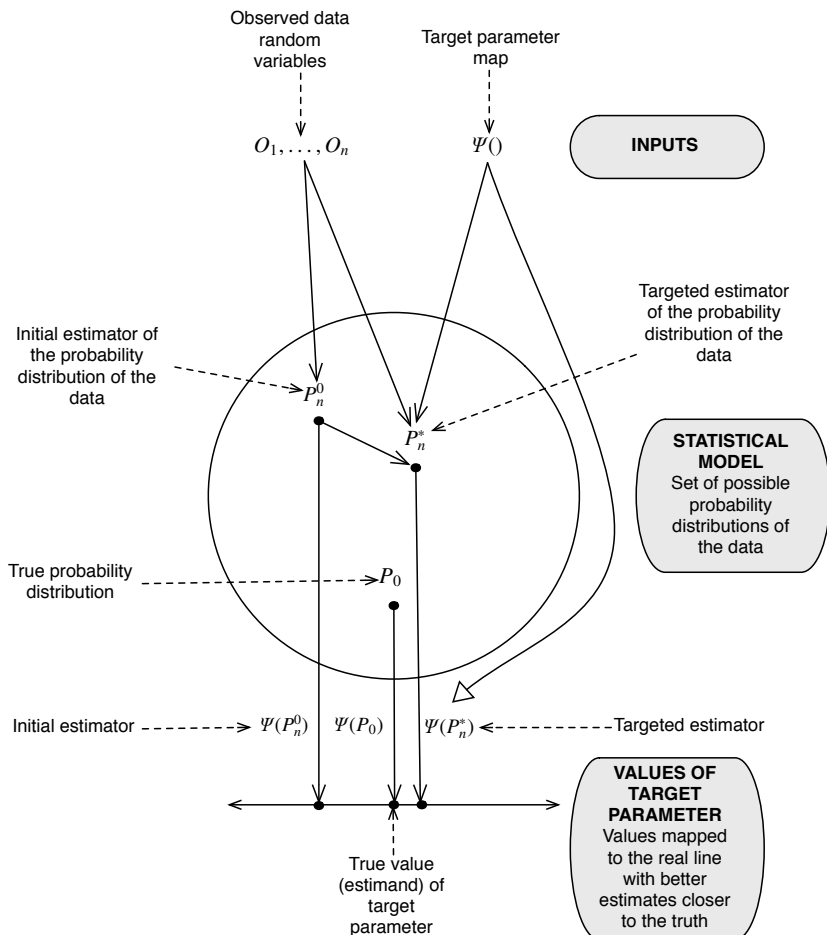


Fig. 5.1 TMLE flow chart.

$\Psi(P_n^*)$ that are closer to the truth than the value generated using the initial estimate P_n^0 : specifically, $\Psi(P_n^*)$ is less biased than $\Psi(P_n^0)$.

TMLE provides a concrete methodology for mapping the initial estimator P_n^0 into a targeted estimator P_n^* , which is described below in terms of an arbitrary statistical model \mathcal{M} and target parameter mapping $\Psi()$ defined on this statistical model. In order to make this more accessible to the reader, we then demonstrate this general template for TMLE with a nonparametric statistical model for a univariate random variable and a survival probability target parameter. Specifically, TMLE involves the following steps:

- Consider the target parameter $\Psi : \mathcal{M} \rightarrow \mathbb{R}$. Compute its pathwise derivative at P and its corresponding canonical gradient $D^*(P)$, which is also called the efficient

influence curve. This object $D^*(P)$, a function of O with mean zero under P , is now available for each possible probability distribution P .

- Define a loss function $L()$ so that $P \rightarrow E_0 L(P)$ is minimized at the true probability distribution P_0 . One could select the log-likelihood loss function $L(P) = -\log P$. However, typically, this loss function is chosen so that it only depends on P through a relevant part $Q(P)$ and $Q \rightarrow L(Q)$ is minimized at $Q_0 = Q(P_0)$. This loss function could also be used to construct a super-learner-based initial estimator of Q_0 .
- For a P in our model \mathcal{M} , define a parametric working model $\{P(\epsilon) : \epsilon\}$ with finite-dimensional parameter ϵ so that $P(\epsilon = 0) = P$, and a “score” $\frac{d}{d\epsilon} L(P(\epsilon))$ at $\epsilon = 0$ for which a linear combination of the components of this “score” equals the efficient influence curve $D^*(P)$ at P . Typically, we simply choose the parametric working model so that this score equals the efficient influence curve $D^*(P)$. If the loss function $L()$ only depends on P through a relevant part $Q = Q(P)$, then this translates into a parametric working model $\{Q(\epsilon) : \epsilon\}$ chosen so that a linear combination of the components of the “score” $\frac{d}{d\epsilon} L(Q(\epsilon))$ at $\epsilon = 0$ equals the efficient influence curve $D^*(P)$ at P .
- Given an initial estimator P_n^0 of P_0 , we compute $\epsilon_n^0 = \arg \min_{\epsilon} \sum_{i=1}^n L(P_n^0(\epsilon))(O_i)$. This yields the first step TMLE $P_n^1 = P_n^0(\epsilon_n^0)$. This process is iterated: start with $k = 1$, compute $\epsilon_n^k = \arg \min_{\epsilon} \sum_{i=1}^n L(P_n^k(\epsilon))(O_i)$ and $P_n^{k+1} = P_n^k(\epsilon_n^k)$, increase k to $k + 1$, and repeat these updating steps until $\epsilon_n^k = 0$. The final update P_n^K at the final step K is denoted by P_n^* and is the TMLE of P_0 . The same algorithm can be directly applied to Q_n^0 of $Q_0 = Q(P_0)$ for the case that the loss function only depends on P through $Q(P)$.
- The TMLE of ψ_0 is now the substitution estimator obtained by plugging P_n^* into the target parameter mapping: $\psi_n^* = \Psi(P_n^*)$. Similarly, if $\psi_0 = \Psi(Q_0)$ and the above loss function $L()$ is a loss function for Q_0 , then we plug the TMLE Q_n^* into the target parameter mapping: $\psi_n^* = \Psi(Q_n^*)$.
- The TMLE P_n^* solves the efficient influence curve equation $0 = \sum_{i=1}^n D^*(P_n^*)(O_i)$, which provides a basis for establishing the asymptotic linearity and efficiency of the TMLE $\Psi(P_n^*)$.

For further presentation of TMLE at this general level we refer the interested reader to Appendix A.

Demonstration of TMLE template. In this section we demonstrate the TMLE template for estimation of survival probability. Suppose we observe n i.i.d. univariate random variables O_1, \dots, O_n with probability distribution P_0 , where O_i represents a time to failure such as death. Suppose that we have no knowledge about this probability distribution, so that we select as statistical model the nonparametric model \mathcal{M} . Let $\Psi(P) = P(O > 5)$ be the target parameter that maps any probability distribution in its survival probability at 5 years, and let $\psi_0 = P_0(O > 5)$ be our target parameter of the true data-generating distribution.

The pathwise derivative $\Psi(P(\epsilon))$ at $\epsilon = 0$ for a parametric submodel (i.e., path) $\{P_S(\epsilon) = (1 + \epsilon S(P))P : \epsilon\}$ with univariate parameter ϵ is given by

$$\left. \frac{d}{d\epsilon} \Psi(P_S(\epsilon)) \right|_{\epsilon=0} = E_P\{I(O > 5) - \Psi(P)\}S(P)(O).$$

Note that indeed, for any function S of O that has mean zero under P and is uniformly bounded, it follows that $P_S(\epsilon)$ is a probability distribution for a small enough choice of ϵ , so that the family of paths indexed by such functions S represents a valid family of submodels through P in the nonparametric model. By definition, it follows that the canonical gradient of this pathwise derivative at P (relative to this family of parametric submodels) is given by $D^*(P)(O) = I(O > 5) - \Psi(P)$. The canonical gradient is also called the efficient influence curve at P .

We could select the log-likelihood loss function $L(P) = -\log P(O)$ as loss function. A parametric working model through P is given by $P(\epsilon) = (1 + \epsilon D^*(P))P$, where ϵ is the univariate fluctuation parameter. Note that this parametric submodel includes P at $\epsilon = 0$ and has a score at $\epsilon = 0$ given by $D^*(P)$, as required for the TMLE algorithm. We are now ready to define the TMLE.

Let P_n^0 be an initial density estimator of the density P_0 . Let

$$\epsilon_n^0 = \arg \max_{\epsilon} \sum_{i=1}^n \log P_n^0(\epsilon)(O_i),$$

and let $P_n^1 = P_n^0(\epsilon_n^0)$ be the corresponding first-step TMLE of P_0 . It can be shown that the next iteration yields $\epsilon_n^1 = 0$, so that convergence of the iterative TMLE algorithm occurs in one step (van der Laan and Rubin 2006). The TMLE is thus given by $P_n^* = P_n^1$, and the TMLE of ψ_0 is given by the plug-in estimator $\psi_n^* = \Psi(P_n^*) = P_n^*(O > 5)$. Since P_n^* solves the efficient influence curve equation, it follows that $\psi_n^* = \frac{1}{n} \sum_{i=1}^n I(O_i > 5)$ is the empirical proportion of subjects that has a survival time larger than 5. This estimator is asymptotically linear with influence curve $D^*(P_0)$ since $\psi_n^* - \psi_0 = \frac{1}{n} \sum_{i=1}^n D^*(P_0)(O_i)$, which proves that the TMLE of ψ_0 is efficient for every choice of initial estimator: apparently, all bias of the initial estimator is removed by this TMLE update step.

Consider a kernel density estimator with an optimally selected bandwidth (e.g., based on likelihood-based cross-validation). Since this optimally selected bandwidth trades off bias and variance for the kernel density estimator as an estimate of the true density P_0 , it will, under some smoothness conditions, select a bandwidth that converges to zero in sample size at a rate $n^{-1/5}$. The bias of such a kernel density estimator converges to zero at the rate $n^{-2/5}$. As a consequence, the substitution estimator of the survival function at t for this kernel density estimator has a bias that converges to zero at a slower rate than $1/\sqrt{n}$ in the sample size n . We can conclude that the substitution estimator of a survival function at 5 years based on this optimal kernel density estimator will have an asymptotic relative efficiency of zero (!) relative to the empirical survival function at 5 years. This simple example demonstrates that a regularized maximum likelihood estimator of P_0 is not targeted toward the target parameter of interest and, by the same token, that current Bayesian inference is not targeted toward the target parameter. However, if we apply the TMLE step to the kernel density estimator, then the resulting TMLE of the survival function is

unbiased and asymptotically efficient, and it even remains unbiased and asymptotically efficient if the kernel density estimator is replaced by an incorrect guess of the true density.

The point is: the best estimator of a density is not a good enough estimator of a particular feature of the density, but the TMLE step takes care of this.

5.2 Definition of TMLE in Context of the Mortality Example

This section presents the definition of TMLE in the context of our mortality example, thereby allowing the reader to derive the TMLE presented in the previous chapter. The reader may recognize the general recipe for TMLE as presented in Sect. 5.1 that can be applied in any semiparametric model with any target parameter. After having read this section, the reader might consider revisiting this general TMLE presentation. Our causal effect of interest is the causal risk difference, and the estimand is the corresponding statistical W -adjusted risk difference, which can be interpreted as the causal risk difference under causal assumptions. The data structure in the illustrative example is $O = (W, A, Y) \sim P_0$. TMLE follows the basic steps enumerated below, which we then illustrate in more detail.

TMLE for the Risk Difference

1. Estimate \bar{Q}_0 using super learner to generate our prediction function \bar{Q}_n^0 . Let $Q_n^0 = (\bar{Q}_n^0, Q_{W,n}^0)$ be the estimate of $Q_0 = (\bar{Q}_0, Q_{W,0})$, where $Q_{W,n}$ is the empirical probability distribution of W_1, \dots, W_n .
2. Estimate the treatment mechanism using super learning. The estimate of g_0 is g_n .
3. Determine a parametric family of fluctuations $\{Q_n^0(\epsilon) : \epsilon\}$ of the initial estimator Q_n^0 with fluctuation parameter ϵ , and a loss function $L(Q)$ so that a linear combination of the components of the derivative of $L(Q_n^0(\epsilon))$ at $\epsilon = 0$ equals the efficient influence curve $D^*(Q_n^0, g_n)$ at any initial estimator $Q_n^0 = (\bar{Q}_n^0, Q_{W,n}^0)$ and g_n . Since the initial estimate $Q_{W,n}^0$ of the marginal distribution of W is the empirical distribution (i.e., nonparametric maximum likelihood estimator), the TMLE using a separate ϵ for fluctuating $Q_{W,n}^0$ and \bar{Q}_n^0 will only fluctuate \bar{Q}_n^0 . The parametric family of fluctuations of \bar{Q}_n^0 is defined by parametric regression including a clever covariate chosen so that the above derivative condition holds with ϵ playing the role of the coefficient in front of the clever covariate. This “clever covariate” $H_n^*(A, W)$ depends on (Q_n^0, g_n) only through g_n , and in the TMLE procedure it needs to be evaluated for each observation (A_i, W_i) , and at $(0, W_i)$, $(1, W_i)$.

4. Update the initial fit $\bar{Q}_n^0(A, W)$ from step 1. This is achieved by holding $\bar{Q}_n^0(A, W)$ fixed (i.e., as intercept) while estimating the coefficient ϵ for $H_n^*(A, W)$ in the parametric working model using maximum likelihood estimation. Let ϵ_n be this parametric maximum likelihood estimator. The updated regression is given by $\bar{Q}_n^1 = \bar{Q}_n^0(\epsilon_n)$. For the risk difference, no iteration is necessary, since the next iteration will not result in any change: that is, the next ϵ_n will be equal to zero. The TMLE of Q_0 is now $Q_n^* = (\bar{Q}_n^1, Q_{W,n}^0)$, where only the conditional mean estimator \bar{Q}_n^0 was updated.
5. Obtain the substitution estimator of the causal risk difference by application of the target parameter mapping to Q_n^* :

$$\psi_n = \Psi(Q_n^*) = \frac{1}{n} \sum_{i=1}^n \{\bar{Q}_n^1(1, W_i) - \bar{Q}_n^1(0, W_i)\}.$$

6. Calculate standard errors based on the influence curve of the TMLE ψ_n , and then calculate p -values and confidence intervals.

There are several concepts in this enumerated step-by-step list that may be somewhat opaque for the reader: the parametric working model coding the fluctuations of the initial estimator, the corresponding clever covariate, the efficient influence curve, and the influence curve. We expand upon the list, including these topics, below. For the nontechnical reader, we provide gray boxes so that you can read these to understand the essential topics relevant to each step. The white boxes outlined in black contain additional technical information for the more theoretical reader.

5.2.1 Estimating \bar{Q}_0

The first step in TMLE is obtaining an estimate \bar{Q}_n^0 for \bar{Q}_0 . This initial fit is achieved using super learning, avoiding assuming a misspecified parametric statistical model.

5.2.2 Estimating g_0

The TMLE procedure uses the estimate of \bar{Q}_0 obtained above in conjunction with an estimate of g_0 . We estimate g_0 with g_n , again using super learning.

5.2.3 Determining the Efficient Influence Curve $D^*(P)$

To obtain such a parametric working model to fluctuate the initial estimator Q_n^0 we need to know the efficient influence curve of the target parameter mapping at a particular P in the statistical model. This is a mathematical exercise that takes as input the definition of the statistical model \mathcal{M} (i.e., the nonparametric model) and the target parameter mapping from this statistical model to the real line (i.e., $\Psi : \mathcal{M} \rightarrow \mathbb{R}$). We refer to Appendix A for required background material. It follows that the efficient influence curve at P_0 only depends on (Q_0, g_0) and is given by

$$D^*(Q_0, g_0)(W, A, Y) = \left(\frac{I(A = 1)}{g_0(1 | W)} - \frac{I(A = 0)}{g_0(0 | W)} \right) (Y - \bar{Q}_0(A, W)) \\ + \bar{Q}_0(1, W) - \bar{Q}_0(0, W) - \Psi(Q_0).$$

More on the efficient influence curve. Calculation of the efficient influence curve, and of components of the efficient influence curve, requires calculations of projections of an element onto a subspace within a Hilbert space. These projections are defined in the Hilbert space $L_0^2(P)$ of functions of O that have mean zero under P endowed with an inner product $\langle S_1, S_2 \rangle_P = E_P S_1(O) S_2(O)$, being the covariance of two functions of O . Two elements in an Hilbert space are orthogonal if the inner product equals zero: so two functions of O are defined as orthogonal if their correlation or covariance equals zero. Recall that a projection of a function S onto a subspace of $L_0^2(P)$ is defined as follows: (1) the projection is an element of the subspace and (2) the difference of S minus the projection is orthogonal to the subspace. The subspaces on which one projects are so-called tangent spaces and subtangent spaces. The tangent space at P is defined as the closure of the linear span of all scores of submodels through P . The tangent space is a subspace of $L_0^2(P)$. The tangent space of a particular variation-independent parameter of P is defined as the closure of the linear span of all scores of submodels through P that only vary this particular factor. We can denote the tangent spaces by $T(P)$ and a projection of a function S onto a $T(P)$ by $\Pi(S | T(P))$.

5.2.4 Determining the Fluctuation Working Model

Now, can we slightly modify the initial estimator \bar{Q}_n^0 to reduce bias for the additive causal effect? Let $Q_{W,n}^0$ be the empirical probability distribution of W_1, \dots, W_n . We refer to the combined conditional probability distribution of Y and the marginal probability distribution of W as Q_0 . $Q_n^0 = (\bar{Q}_n^0, Q_{W,n}^0)$ denotes the initial estimator of this Q_0 . We also remind the reader that the target parameter ψ_0 only depends on

P_0 through \bar{Q}_0 and $Q_{W,0}$. Since the empirical distribution $Q_{W,n}^0$ is already a nonparametric maximum likelihood estimator of the true marginal probability distribution of W , for the sake of bias reduction for the target parameter, we can focus on only updating \bar{Q}_n^0 , as explained below.

We want to reduce the bias of our initial estimator, where the initial estimator is a random variable that has bias and variance. We only need to update \bar{Q}_n^0 since the empirical distribution $Q_{W,n}^0$ is a nonparametric maximum likelihood estimator (and can thus not generate bias for our target parameter).

Our parametric working model is denoted as $\{\bar{Q}_n^0(\epsilon) : \epsilon\}$, which is a small parametric statistical model, a one-dimensional submodel that goes through the initial estimate $\bar{Q}_n^0(A, W)$ at $\epsilon = 0$. If we use the log-likelihood loss function

$$L(\bar{Q})(O) = -\log \bar{Q}(A, W)^Y (1 - \bar{Q}(A, W))^{1-Y},$$

then the parametric working model for fluctuating the conditional probability distribution of Y , given (A, W) , needs to have the property

$$\frac{d}{d\epsilon} \log \bar{Q}_n^0(\epsilon)(A, W)^Y (1 - \bar{Q}_n^0(A, W))^{1-Y} \Big|_{\epsilon=0} = D_Y^*(Q_n^0, g_n)(W, A, Y), \quad (5.1)$$

where $D_Y^*(Q_n^0, g_n)$ is the appropriate component of the efficient influence curve $D^*(Q_n^0, g_n)$ of the target parameter mapping at (Q_n^0, g_n) . Formally, the appropriate component D_Y^* is the component of the efficient influence curve that equals a score of a fluctuation of a conditional distribution of Y , given (A, W) . These components of the efficient influence curve that correspond with scores of fluctuations that only vary certain parts of factors of the probability distribution can be computed with Hilbert space projections. We provide the required background and tools in Appendix A and various subsequent chapters.

More on fluctuating the initial estimator. If the target parameter ψ_0 depends on different variation-independent parts $(Q_{W,0}, \bar{Q}_0)$ of the probability distribution P_0 , then one can decide to fluctuate the initial estimators $(Q_{W,n}^0, \bar{Q}_n)$ with separate submodels and separate loss functions $L(Q_W) = -\log Q_W$ and $L(\bar{Q})$, respectively. The submodels $\{Q_{W,n}^0(\epsilon) : \epsilon\}$, $\{\bar{Q}_n(\epsilon) : \epsilon\}$ and their corresponding loss functions $L(Q_W)$ and $L(\bar{Q})$ need to be chosen such that a linear combination of the components of the derivative $\frac{d}{d\epsilon} L(Q_n^0(\epsilon)) \Big|_{\epsilon=0}$ equals $D^*(Q_n^0, g_n)$ for the sum-loss function $L(Q) = L(Q_W) + L(\bar{Q})$. This corresponds with requiring that each of the two loss functions generates a “score” so that the sum of these two “scores” equals the efficient influence curve. If the initial estimator $Q_{W,n}^0$ is a nonparametric maximum likelihood estimator, the TMLE using a separate ϵ_1 and ϵ_2 for the two submodels will not update $Q_{W,n}^0$.

Following the protocol of TMLE, we also need to fluctuate the marginal distribution of W . For that purpose we select as loss function of $Q_{W,0}$ the log-likelihood loss function $-\log Q_W$. Then we would select a parametric working model coding fluctuations $Q_{W,n}^0(\epsilon)$ of $Q_{W,n}^0$ so that

$$\left. \frac{d}{d\epsilon} \log Q_{W,n}^0(\epsilon) \right|_{\epsilon=0} = D_W^*(Q_n^0, g_n),$$

where D_W^* is the component of the efficient influence curve that is a score of a fluctuation of the marginal distribution of W .

Tangent spaces. Since Q_W and \bar{Q} represent parameters of different factors P_W and $P_{Y|A,W}$ in a factorization of $P = P_W P_{A|W} P_{Y|A,W}$, these components $D_W^*(P)$ and $D_Y^*(P)$ can be defined as the projection of the efficient influence curve $D^*(P)$ onto the tangent space of P_W at P and $P_{Y|A,W}$ at P , respectively. The tangent space T_W of P_W is given by all functions of W with mean zero. The tangent space T_Y of $P_{Y|A,W}$ is given by all functions of W, A, Y for which the conditional mean, given A, W , equals zero. The tangent space T_A of $P_{A|W}$ is given by all functions of A, W , with conditional mean zero, given W . These three tangent spaces are orthogonal, as a general consequence of the factorization of P into the three factors. The projection of a function S onto these three tangent spaces is given by $\Pi(S | T_W) = E_P(S(O) | W)$, $\Pi(S | T_Y) = S(O) - E_P(S(O) | A, W)$, and $\Pi(S | T_A) = E_P(S(O) | A, W) - E_P(S | W)$, respectively. From these projection formulas and setting $S = D^*(P)$, the explicit forms of $D_W^*(P) = \Pi(D^*(P) | T_W)$ and $D_Y^*(P) = \Pi(D^*(P) | T_Y)$ can be calculated as provided below, and for each choice of P . It also follows that the projection of $D^*(P)$ onto the tangent space of $P_{A|W}$ equals zero: $\Pi(D^*(P) | T_A) = 0$. The latter formally explains that the TMLE does not require fluctuating the initial estimator of g_0 . It follows that the efficient influence curve $D^*(P)$ at P can be decomposed as:

$$D^*(P) = D_Y^*(P) + D_W^*(P).$$

Our loss function for Q is now $L(Q) = L(\bar{Q}) + L(Q_W)$, and with this parametric working model coding fluctuations $Q_n^0(\epsilon) = (Q_{W,n}^0(\epsilon), \bar{Q}_n^0(\epsilon))$ of Q_n^0 , we have that the derivative of $\epsilon \rightarrow L(Q_n^0(\epsilon))$ at $\epsilon = 0$ equals the efficient influence curve at (Q_n^0, g_n) . If we use different ϵ for each component of Q_n^0 , then the two derivatives span the efficient influence curve, since the efficient influence curve equals the sum of the two scores D_Y^* and D_W^* . Either way, the derivative condition is satisfied:

$$\left\langle \left. \frac{d}{d\epsilon} L(Q_n^0(\epsilon)) \right|_{\epsilon=0} \right\rangle \supset D^*(Q_n^0, g_n), \quad (5.2)$$

where $D^*(Q_n^0, g_n) = D_Y^*(Q_n^0, g_n) + D_W^*(Q_n^0, g_n)$. Here we used the notation $\langle (h_1, \dots, h_k) \rangle$ for the linear space consisting of all linear combinations of the functions h_1, \dots, h_k . That is, the task of obtaining a loss function and parametric working model for fluctuating Q_n^0 so that the derivative condition holds has been completed.

Due to this property (5.2) of the parametric working model, the TMLE has the important feature that it solves the efficient influence curve equation $0 = \sum_i D^*(Q_n^*, g_n)(O_i)$ (also called the efficient score equation). Why is this true? Because at the next iteration of TMLE, the parametric maximum likelihood estimator $\epsilon_n = 0$, and a parametric maximum likelihood estimator solves its score equation, which exactly yields this efficient score equation. This is a strong feature of the procedure as it implies that TMLE is double robust and (locally) efficient under regularity conditions. In other words, TMLE is consistent and asymptotically linear if either Q_n or g_n is a consistent estimator, and if both estimators are asymptotically consistent, then TMLE is asymptotically efficient.

However, if one uses a separate ϵ_W and ϵ for the two parametric working models through $Q_{W,n}^0$ and \bar{Q}_n^0 , respectively, then the maximum likelihood estimator of ϵ_W equals zero, showing that TMLE will only update \bar{Q}_n^0 . Therefore, it was never necessary to update the part of Q_n^0 that was already nonparametrically estimated.

If the initial estimator of $Q_{W,0}$ is a nonparametric maximum likelihood estimator, then the TMLE does not update this part of the initial estimator Q_n^0 .

Of course, we have not been explicit yet about how to construct this submodel $\bar{Q}_n^0(\epsilon)$ through \bar{Q}_n^0 . For that purpose, we now note that $D_Y^*(Q_n^0, g_n)$ equals a function $H_n^*(A, W)$ times the residual $(Y - \bar{Q}_n^0(A, W))$, where

$$H_n^*(A, W) \equiv \left(\frac{I(A = 1)}{g_n(A = 1 | W)} - \frac{I(A = 0)}{g_n(A = 0 | W)} \right).$$

Here $I(A = 1)$ is an indicator variable that takes the value 1 when $A = 1$. One can see that for $A = 1$ the second term disappears, and for $A = 0$ the first term disappears.

It can be shown (and it is a classical result for parametric logistic main term regression in a parametric statistical model) that the score of a coefficient in front of a covariate in a logistic linear regression in a parametric statistical model for a conditional distribution of a binary Y equals the covariate times the residual. Therefore, we can select the following parametric working model for fluctuating the initial estimate of the conditional probability distribution of Y , given (A, W) , or, equivalently, for the estimate of the probability of $Y = 1$, given (A, W) :

$$\bar{Q}_n^0(\epsilon)(Y = 1 | A, W) = \frac{1}{1 + \exp\left(-\log \frac{\bar{Q}_n^0}{(1 - \bar{Q}_n^0)}(A, W) - \epsilon H_n^*(A, W)\right)}.$$

By this classical result, it follows that indeed the score of ϵ of this univariate logistic regression submodel at $\epsilon = 0$ equals $D_Y^*(Q_n^0, g_n)$. That is, we now have really fully succeeded in finding a parametric submodel through the initial estimator Q_n^0 that satisfies the required derivative condition. Since $H_n^*(A, W)$ now just plays the role of a covariate in a logistic regression, using an offset, this explains why we call the covariate $H_n^*(A, W)$ a clever covariate.

More on constructing the submodel. If one needs a submodel through an initial estimator of a conditional distribution of a binary variable Y , given a set of parent variables $Pa(Y)$, and it needs to have a particular score D_Y^* , then one can define this submodel as a univariate logistic regression model, using the initial estimator as offset, with univariate clever covariate defined as $H^*(Pa(Y)) = E(D_Y^* | Y = 1, Pa(Y)) - E(D_Y^* | Y = 0, Pa(Y))$. Application of this general result to the above setting yields the clever covariate $H^*(A, W)$ presented above.

If our goal was to target $P_0(Y_1 = 1)$ or $P_0(Y_0 = 1)$, then going through the same protocol for the TMLE shows that one would use as clever covariate

$$H_{0,n}^*(A, W) \equiv \left(\frac{I(A = 0)}{g_n(A = 0 | W)} \right) \text{ or } H_{1,n}^*(A, W) \equiv \left(\frac{I(A = 1)}{g_n(A = 1 | W)} \right).$$

By targeting these two parameters simultaneously, using a two-dimensional clever covariate with coefficients ϵ_1, ϵ_2 , one automatically obtains a valid TMLE for parameters that are functions of these two marginal counterfactual probabilities, such as a causal relative risk and causal odds ratio.

By computing the TMLE that targets a multidimensional target parameter, one also obtains a valid TMLE for any (say) univariate summary measure of the multidimensional target parameter. By valid we mean that this TMLE will still satisfy the same asymptotic properties, such as efficiency and double robustness, as the TMLE that directly targets the particular summary measure. The TMLE that targets the univariate summary measure of the multidimensional parameter may have a better finite sample performance than the TMLE that targets the whole multidimensional target parameter, in particular, if the dimension of the multidimensional parameter is large.

5.2.5 Updating \bar{Q}_n^0

We first perform a logistic linear regression of Y on $H_n^*(A, W)$ where $\bar{Q}_n^0(A, W)$ is held fixed (i.e., used as an offset), and an additional intercept is suppressed in order to estimate the coefficient in front of $H_n^*(A, W)$, denoted ϵ . The TMLE procedure

is then able to incorporate information from g_n , through $H_n^*(A, W)$, into an updated regression. It does this by extracting ϵ_n , the maximum likelihood estimator of ϵ , from the fit described above, and updating the estimate \bar{Q}_n^0 according to the logistic regression working model. This updated regression is then given by \bar{Q}_n^1 :

$$\text{logit } \bar{Q}_n^1(A, W) = \text{logit } \bar{Q}_n^0(A, W) + \epsilon_n H_n^*(A, W).$$

One iterates this updating process until the next $\epsilon_n = 0$ or has converged to zero, but, in this example, convergence is achieved in one step. The TMLE of Q_0 is now $Q_n^* = (Q_{W,n}^0, \bar{Q}_n^1)$. Note that this step is equivalent to $(\epsilon_{1n}, \epsilon_{2n}) = \arg \min_{\epsilon_1, \epsilon_2} \sum_i L(Q_n^0(\epsilon_1, \epsilon_2))(O_i)$, and setting $Q_n^1 = Q_n^0(\epsilon_{1n}, \epsilon_{2n})$, where, as noted above, $\epsilon_{1n} = 0$, so that only \bar{Q}_n^0 is updated.

Given a parametric working model $Q_n^0(\epsilon)$ with fluctuation parameter ϵ , and a loss function $L(Q)$ satisfying (5.2), the first-step TMLE is defined by determining the minimum ϵ_n^0 of $\sum_{i=1}^n L(Q_n^0(\epsilon))(O_i)$ and setting $Q_n^1 = Q_n^0(\epsilon_n^0)$. This updating process is iterated until convergence of $\epsilon_n^k = \arg \min_{\epsilon} \sum_{i=1}^n L(Q_n^k(\epsilon))$ to zero, and the final update Q_n^* is referred to as the TMLE of Q_0 . In this case, the next $\epsilon_n^1 = 0$, so that convergence is achieved in one step and $Q_n^* = Q_n^1$.

5.2.6 Estimating the Target Parameter

The estimate $\bar{Q}_n^* = \bar{Q}_n^1$ obtained in the previous step is now plugged into our target parameter mapping, together with the empirical distribution of W , resulting in the targeted substitution estimator given by

$$\psi_n = \Psi(Q_n^*) = \frac{1}{n} \sum_{i=1}^n \{\bar{Q}_n^1(1, W_i) - \bar{Q}_n^1(0, W_i)\}.$$

This mapping is accomplished by evaluating $\bar{Q}_n^1(1, W_i)$ and $\bar{Q}_n^1(0, W_i)$ for each observation i and plugging these values into the above equation.

5.2.7 Calculating Standard Errors

The calculation of standard errors for TMLE can be based on the central limit theorem, relying on δ -method conditions. (See Appendix A for an advanced introduction to these topics.) Under such regularity conditions, the asymptotic behavior of the estimator, that is, its asymptotic normal limit distribution, is completely characterized

by the so-called influence curve of the estimator in question. In our example, we need to know the influence curve of the TMLE of its estimand.

Note that, in order to recognize that an estimator is a random variable, an estimator should be represented as a mapping from the data into the parameter space, where the data O_1, \dots, O_n can be represented by the empirical probability distribution function P_n . Therefore, let $\hat{\Psi}(P_n)$ be the TMLE described above. Since the TMLE is a substitution estimator, we have $\hat{\Psi}(P_n) = \Psi(P_n^*)$ for a targeted estimator P_n^* of P_0 . An estimator $\hat{\Psi}(P_n)$ of ψ_0 is asymptotically linear with influence curve $IC(O)$ if it satisfies:

$$\sqrt{n}(\hat{\Psi}(P_n) - \psi_0) = \frac{1}{\sqrt{n}} \sum_{i=1}^n IC(O_i) + o_{P_0}(1).$$

Here the remainder term, denoted by $o_{P_0}(1)$, is a random variable that converges to zero in probability when the sample size converges to infinity. The influence curve $IC(O)$ is a random variable with mean zero under P_0 .

More on estimators and the influence curve. An estimator $\hat{\Psi}(P_n)$ is a function $\hat{\Psi}$ of the empirical probability distribution function P_n . Specifically, one can express the estimator as a function $\hat{\Psi}$ of a large family of empirical means $1/n \sum_{i=1}^n f(O_i)$ of functions f of O varying over a class of functions \mathcal{F} . We say the estimator is a function of $P_n = (P_n f : f \in \mathcal{F})$, where we use the notation $P_n f \equiv 1/n \sum_{i=1}^n f(O_i)$. By proving that the estimator is a differentiable function $\hat{\Psi}$ of $P_n = (P_n f : f \in \mathcal{F})$ at $P_0 = (P_0 f : f \in \mathcal{F})$, and that a uniform central limit theorem applies to P_n based on empirical process theory, it follows that the estimator minus its estimand $\psi_0 = \hat{\Psi}(P_0)$ behaves in first order as an empirical mean of $IC(O_i)$: we write $\hat{\Psi}(P_n) - \psi_0 = (P_n - P_0)IC + o_P(1/\sqrt{n})$. This function $IC(O)$ is called the influence curve of the estimator, and it is uniquely determined by the derivative of $\hat{\Psi}$. Specifically, $IC(O) = \sum_{f \in \mathcal{F}} \frac{d}{dP_0 f} \hat{\Psi}((P_0 f : f)(f(O) - P_0 f))$, where, formally, the \sum becomes an integral when \mathcal{F} is not finite.

Asymptotic linearity is a desirable property as it indicates that the estimator behaves like an empirical mean, and, as a consequence, its bias converges to zero in sample size at a rate faster than $1/\sqrt{n}$, and, for n large enough, it is approximately normally distributed. The influence curve of an estimator evaluated as a function in O measures how robust the estimator is toward extreme values. The influence curve $IC(O)$ has mean zero under sampling from the true probability distribution P_0 , and its (finite) variance is the asymptotic variance of the standardized estimator $\sqrt{n}(\hat{\Psi}(P_n) - \psi_0)$.

In other words, the variance of $\hat{\Psi}(P_n)$ is well approximated by the variance of the influence curve, divided by sample size n . If ψ_0 is multivariate, then the

covariance matrix of $\hat{\Psi}(P_n)$ is well approximated by the covariance matrix of the multivariate influence curve divided by sample size n . More importantly, the probability distribution of $\hat{\Psi}(P_n)$ is well approximated by a normal distribution with mean ψ_0 and the covariance matrix of the influence curve, divided by sample size.

An estimator is asymptotically efficient if its influence curve is equal to the efficient influence curve, $IC(O) = D^*(O)$. The influence curve of the TMLE indeed equals D^* if Q_n^* is a consistent estimator of Q_0 , and g_n is a consistent estimator of g_0 . A complete technical understanding of influence curve derivation is not necessary to implement the TMLE procedure. However, we provide Appendix A for a detailed methodology for deriving the influence curve of an estimator.

More on asymptotic linearity and efficiency. The TMLE is a consistent estimator of ψ_0 if either \bar{Q}_n is consistent for \bar{Q}_0 or g_n is consistent for g_0 . The TMLE is asymptotically linear under additional conditions. For a detailed theorem establishing asymptotic linearity and efficiency of the TMLE, we refer the reader to Chap. 27. In particular, if for some $\delta > 0$, $\delta < g_0(1 | W) < 1 - \delta$, and the product of the L^2 -norm of $\bar{Q}_n - \bar{Q}_0$ and the L^2 -norm of $g_n - g_0$ converges to zero at faster rate than $1/\sqrt{n}$, then the TMLE is asymptotically efficient. If g_n is a consistent estimator of g_0 , then the influence curve of the TMLE $\hat{\Psi}(P_n)$ equals $IC = D^*(Q^*, g_0) - \Pi(D^*(Q^*, g_0) | T_g)$, the efficient influence curve at the possibly misspecified limit of Q_n^* minus its projection on the tangent space of the model for the treatment mechanism g_0 . The projection term makes $D^*(Q^*, g_0)$ a conservative working influence curve, and the projection term equals zero if either $Q^* = Q_0$ or g_0 was known and $g_n = g_0$.

From these formal asymptotic linearity results for the TMLE it follows that if g_n is a consistent estimator of g_0 , then the TMLE $\hat{\Psi}(P_n)$ is asymptotically linear with an influence curve that can be conservatively approximated by $D^*(Q^*, g_0)$, where Q^* denotes the possibly misspecified estimand of Q_n^* . If g_0 was known, as in a randomized controlled trial, and g_n was not estimated, then the influence curve of the TMLE equals $D^*(Q^*, g_0)$. If, on the other hand, g_n was estimated under a correctly specified model for g_0 , then the influence curve of the TMLE has a smaller variance than the variance of $D^*(Q^*, g_0)$, except if $Q^* = Q_0$, in which case the influence curve of the TMLE equals the efficient influence curve $D^*(Q_0, g_0)$. As a consequence, we can use as a working estimated influence curve for the TMLE

$$IC_n(O) = \left(\frac{I(A=1)}{g_n(1|W)} - \frac{I(A=0)}{g_n(0|W)} \right) (Y - \bar{Q}_n^1(A, W)) + \bar{Q}_n^1(1, W) - \bar{Q}_n^1(0, W) - \psi_n.$$

Even if \bar{Q}_n^1 is inconsistent, but g_n is consistent, this influence curve can be used to obtain an asymptotically *conservative* estimator of the variance of the TMLE

$\hat{\Psi}(P_n)$. This is very convenient since the TMLE requires calculation of $D^*(Q_n^*, g_n)$, and apparently we can use the latter as influence curve to estimate the normal limit distribution of the TMLE.

If one assumes that g_n is a consistent maximum-likelihood-based estimator of g_0 , then one can (asymptotically) conservatively estimate the variance of the TMLE with the sample variance of the estimated efficient influence curve $D^*(Q_n^*, g_n)$.

An estimate of the asymptotic variance of the standardized TMLE, $\sqrt{n}(\hat{\Psi}(P_n) - \psi_0)$, viewed as a random variable, using the estimate of the influence curve $IC_n(O)$ is thereby given by

$$\sigma_n^2 = \frac{1}{n} \sum_{i=1}^n IC_n^2(o_i).$$

5.3 Foundation and Philosophy of TMLE

TMLE in semiparametric statistical models for P_0 is the extension of maximum likelihood estimation in parametric statistical models. Three key ingredients are needed for this extension. Firstly, one needs to define the parameter of interest semiparametrically as a function of the data-generating distribution varying over the (large) semiparametric statistical model. Many practitioners are used to thinking of their parameter in terms of a regression coefficient, but that luxury is not available in semi- or nonparametric statistical models. Instead, one has to carefully think of what feature of the distribution of the data one wishes to target.

Secondly, one needs to estimate the true distribution P_0 , or at least its relevant factor or portion as needed to evaluate the target parameter, and this estimate should respect the actual semiparametric statistical model. As a consequence, nonparametric maximum likelihood estimation is often ill defined or results in a complete overfit, and thereby results in estimators of the target parameter that are too variable. We discussed this issue in Chap. 3. The theoretical results obtained for the cross-validation selector (discrete super learner) inspired the general super learning methodology for estimation of probability distributions of the data, or factors of other high-dimensional parameters of the probability distributions of the data. In the sequel, a reference to a true probability distribution of the data is meant to refer to this relevant part of the true probability distribution of the data. This super learning methodology takes as input a collection of candidate estimators of the distribution of the data and then uses cross-validation to determine the best weighted combination of these estimators. It is assumed or arranged that the loss function is uniformly bounded so that oracle results for the cross-validation selector apply. The super learning methodology results now in an estimator of the distribution of the

data that will be used as an initial estimator in the TMLE procedure. The oracle results for this super learner teach us that the initial estimator is optimized with respect to a global loss function such as the log-likelihood loss function and is thereby not targeted toward the target parameter, $\Psi(P_0)$. That is, it will be too biased for $\Psi(P_0)$ due to a bias–variance tradeoff with respect to the more ambitious full P_0 (or relevant portion thereof) instead of having used a bias–variance tradeoff with respect to $\Psi(P_0)$. The targeted maximum likelihood step is tailored to remove bias due to the nontargeting of the initial estimator.

The targeted maximum likelihood step involves now updating this initial (super-learning-based) estimator P_n^0 of P_0 to tailor its fit to estimation of the target ψ_0 , the value of the parameter $\Psi(P_0)$. This is carried out by determining a cleverly chosen parametric working model modeling fluctuations $P_n^0(\epsilon)$ of the initial estimator P_n^0 with a (say) univariate fluctuation parameter ϵ . The value $\epsilon = 0$ corresponds with no fluctuation so that $P_n^0(0) = P_n^0$. One now estimates ϵ with maximum likelihood estimation, treating the initial estimator as a fixed offset, and updates the initial estimator accordingly. If needed, this updating step is iterated to convergence, and the final update P_n^* is called the TMLE of P_0 , while the resulting substitution estimator $\hat{\Psi}(P_n^*)$ of $\Psi(P_0)$ is the TMLE of ψ_0 . This targeted maximum likelihood step thus uses a parametric maximum likelihood estimator, accordingly to a cleverly chosen parametric working model that includes the initial estimator, to obtain a bias reduction for the target $\Psi(P_0)$.

This is not just any parametric working model. That is, we wish to select a parametric working model such that the parametric maximum likelihood estimator is maximally effective in removing bias for the target parameter, at minimal increase in variance. So if ϵ_n is the parametric maximum likelihood estimator of ϵ , then we want the mean squared error of $\Psi(P_n^0(\epsilon_n)) - \psi_0$ to be as small as possible. We want this parametric working model to really listen to the information in the data that is relevant for the target parameter. In fact, we would like the parametric maximum likelihood estimator to be as responsive to the information in the data that is relevant for the target parameter as an estimator that is asymptotically efficient in the semiparametric model.

To get insight into what kind of choice of parametric working model may be as adaptive to such target-parameter-specific features in the data as a semiparametric efficient estimator, we make the following observations. Suppose one is interested in determining the parametric working model coding fluctuations $P_0(\epsilon)$ of P_0 so that the maximum likelihood estimator of $\psi_0 = \Psi(P_0(\epsilon = 0))$ according to this parametric working model is asymptotically equivalent to an efficient estimator in the large semiparametric model. Note that this parametric working model is not told that the true value of ϵ equals zero. It happens to be the case that from an asymptotic efficiency perspective this can be achieved as follows. Among all possible parametric working models that code fluctuations $P_0(\epsilon)$ of the true P_0 we chose the one for which the Cramer–Rao lower bound for the target parameter $\Psi(P_0(\epsilon))$ at $\epsilon = 0$ is equivalent to the semiparametric information bound for the target parameter at P_0 . The Cramer–Rao lower bound for a parametric working model $P_0(\epsilon)$ is given by

$$\frac{\left\{ \frac{d}{d\epsilon} \Psi(P_0(\epsilon)) \Big|_{\epsilon=0} \right\}^2}{I(0)},$$

where $I(0)$ denotes the variance of the score of the parametric working model at $\epsilon = 0$. In parametric model theory $I(0)$ is called the information at parameter value 0. The semiparametric information bound for the target parameter at P_0 is defined as the supremum over all these possible Cramer–Rao lower bounds for the parametric working models. That is, the semiparametric information bound is defined as the Cramer–Rao lower bound for the hardest parametric working model. Thus, the parametric working model for which the parametric maximum likelihood estimator is as responsive to the data with respect to the target parameter as a semiparametric efficient estimator is actually given by this hardest parametric working model. Indeed, the TMLE selects this hardest parametric working model, but through P_n^0 .

Note also that this hardest working parametric model can also be interpreted as the one that maximizes the change of the target parameter relative to a change $P_0(\epsilon) - P_0$ under small amounts of fluctuations. Thus this hardest working parametric model through an initial estimator P_n^0 will maximize the change of the target parameter relative to the initial value $\Psi(P_n^0)$ for small values of ϵ .

Beyond the practical appeal of this TMLE update that uses the parametric likelihood to fit the target parameter of interest, an important feature of the TMLE is that it solves the efficient influence curve equation, also called the efficient score equation, of the target parameter. We refer the reader to Sect. 5.2 and Appendix A for relevant material on the efficient influence curve. For now, it suffices to know that an estimator is semiparametric efficient if the estimator minus the true target parameter behaves as an empirical mean of $D^*(P_0)(O_i)$, $i = 1, \dots, n$, showing the incredible importance of this transformation $D^*(P_0)$ of O , which somehow captures all the relevant information of O for the sake of learning the statistical parameter $\Psi(P_0)$. If $D^*(P)(O)$ is the efficient influence curve at P , a possible probability distribution for O in the statistical model, and P_n^* is the TMLE of P_0 , then, $0 = \sum_{i=1}^n D^*(P_n^*)(O_i)$.

Just as a parametric maximum likelihood estimator solves a score equation by virtue of its maximizing the likelihood over the unknown parameters, a TMLE solves the target-parameter-specific score equation for the target parameter by virtue of maximizing the likelihood in a targeted direction. This can then be used to establish that the TMLE is asymptotically efficient if the initial estimator is consistent and remarkably robust in the sense that for many data structures and semiparametric statistical models, the TMLE of ψ_0 remains consistent even if the initial estimator is inconsistent. By using submodels that have a multivariate fluctuation parameter ϵ , the TMLE will solve the score equation implied by each component of ϵ . In this manner, one can obtain TMLEs that solve not only the efficient influence curve/efficient score equation for the target parameter, but also an equation that characterizes other interesting properties, such as being an imputation estimator (Gruber and van der Laan 2010a).

In particular, in semiparametric models used to define causal effect parameters, the TMLE is a double robust estimator. In such semiparametric models the probability distribution function P_0 can be factorized as $P_0(O) = Q_0(O)g_0(O)$, where g_0

is the treatment mechanism and Q_0 is the relevant factor that defines the g-formula for the counterfactual distributions. The TMLE $\Psi(Q_n^*)$ of $\psi_0 = \Psi(Q_0)$ is consistent if either Q_n^* or g_n is consistent. In our example, g_n is the estimator of the treatment mechanism $g_0(A | W) = P_0(A | W)$, and Q_n^* is the TMLE of Q_0 .

5.4 Summary

TMLE of a parameter $\Psi(Q_0)$ distinguishes from nonparametric or regularized maximum likelihood estimation by fully utilizing the power of cross-validation (super learning) to fine-tune the bias–variance tradeoff with respect to the part Q_0 of the data-generating distribution, thereby increasing adaptivity to the true Q_0 , and by targeting the fit to remove bias with respect to ψ_0 . The loss-based super learner of Q_0 already outperforms with respect to bias and variance a regularized maximum likelihood estimator for the semiparametric statistical model with respect to estimation of Q_0 itself by its asymptotic equivalence to the oracle selector: one could include the regularized maximum likelihood estimator in the collection of algorithms for the super learner. Just due to using the loss-based super learner it already achieves higher rates of convergence for Q_0 itself, thereby improving both in bias and variance for Q_0 as well as $\Psi(Q_0)$. In addition, due to the targeting step, which again utilizes super learning for estimation of the required g_0 in the fluctuation function, it is less biased for ψ_0 than the initial loss-function-based super learner estimator, and, as a bonus, the statistical inference based on the central limit theorem is also heavily improved relative to just using a nontargeted regularized maximum likelihood estimator.

Overall it comes down to the following: the TMLE is a semiparametric efficient substitution estimator. This means it fully utilizes all the information in the data (super learning and asymptotic efficiency), in addition to fully using knowledge about global constraints implied by the statistical semiparametric statistical model \mathcal{M} and the target parameter mapping (by being a substitution estimator), thereby making it robust under sparsity with respect to the target parameter. It fully incorporates the power of super learning for the benefit of getting closer to the truth in finite samples.