

# Object-oriented Computation of Sandwich Estimators

Achim Zeileis

Wirtschaftsuniversität Wien

---

## Abstract

Some ideas about generalizing the tools available in **sandwich** to more models, in particular fully supporting `glm()` and maybe `survreg()`, `gam()`, `betareg()`, ...

*Keywords:* covariance matrix estimators, estimating functions, object orientation, R.

---

## 1. Theory

To fix notations, let us assume we have data in a regression setup, i.e.,  $(y_i, x_i)$  for  $i = 1, \dots, n$ , that follow some distribution that is controlled by a  $k$ -dimensional parameter vector  $\theta$ . In many situations, an estimating function  $\psi(\cdot)$  is available for this type of models such that  $E[\psi(y, x, \theta)] = 0$ . Then, under certain weak regularity conditions,  $\theta$  can be estimated using an M-estimator  $\hat{\theta}$  implicitly defined as

$$\sum_{i=1}^n \psi(y_i, x_i, \hat{\theta}) = 0. \quad (1)$$

This includes in particular maximum likelihood (ML) and ordinary least squares (OLS) estimation, where the estimating function  $\psi(\cdot)$  is the derivative of an objective function  $\Psi(\cdot)$ :

$$\psi(y, x, \theta) = \frac{\partial \Psi(y, x, \theta)}{\partial \theta}. \quad (2)$$

Inference about  $\theta$  is then typically performed relying on a central limit theorem of type

$$\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N(0, S(\theta)), \quad (3)$$

where  $\xrightarrow{d}$  denotes convergence in distribution. For the covariance matrix  $S(\theta)$ , a sandwich formula can be given

$$S(\theta) = B(\theta) M(\theta) B(\theta) \quad (4)$$

$$B(\theta) = (E[-\psi'(y, x, \theta)])^{-1} \quad (5)$$

$$M(\theta) = \text{VAR}[\psi(y, x, \theta)] \quad (6)$$

i.e., the “meat” of the sandwich  $M(\theta)$  is the variance of the estimating function and the “bread” is the inverse of the expectation of its first derivative (again with respect to  $\theta$ ).

In correctly specified models, estimated by ML, this sandwich expression for  $S(\theta)$  can be simplified because both  $B(\theta)^{-1}$  and  $M(\theta)$  correspond to the Fisher information matrix. However, typically covariance matrix estimates that are in some sense robust against misspecification can be computed by using estimators for  $B(\theta)$  and  $M(\theta)$  and plugging them into the sandwich formula.

Many of the models of interest to us, even have a bit more structure: the objective function  $\Psi(y, x, \theta)$  depends on  $x$  and  $\theta$  in a special way, namely it does only depend on the linear predictor  $\eta = x^\top \theta$ . The estimating function is then of type

$$\psi(y, x, \theta) = \frac{\partial \Psi}{\partial \eta} \cdot \frac{\partial \eta}{\partial \theta} = \frac{\partial \Psi}{\partial \eta} \cdot x. \quad (7)$$

The partial derivative  $r(y, \eta) = \partial \Psi(y, \eta) / \partial \eta$  is in some models also called “working residual” corresponding to the usual residuals in linear regression models. In these model, the meat of the sandwich can also be written as

$$M(\theta) = x \text{VAR}[r(y, x^\top \theta)] x^\top. \quad (8)$$

## 2. Implementation

There are already many model fitting functions which compute estimates  $\hat{\theta}$  for a multitude of regression models that can be seen as special cases of the framework outlined in the previous section. Many of these functions already have a `vcov()` method, which typically relies on the assumption of correctly specified models estimated via ML. Therefore, standard inference in these models (typically as reported by `summary()`) is essentially based on the estimated covariance matrix  $1/n B(\hat{\theta})$ .

To be able to compute (more robust) sandwich estimators in this general setup, we propose the following tools.

### *The meat*

Various estimators for  $M(\theta)$  have been suggested in the literature, most of which are based on the empirical values of estimating functions. Hence, a natural idea for object-oriented implementation of such estimators is the following: provide various functions that compute different estimators for the meat based only on an `estfun()` extractor function that extracts the empirical estimating function from a fitted model object. The `estfun()` method should return an  $n \times k$  matrix with

$$\begin{pmatrix} \psi(y_1, x_1, \hat{\theta}) \\ \vdots \\ \psi(y_n, x_n, \hat{\theta}) \end{pmatrix}.$$

Based on this the function `meat()` can simply compute crossproducts for deriving a naive estimator of  $M(\theta)$  which is called outer product of gradients in some communities. A simplified version is

```
meat <- function(x)
{
  psi <- estfun(x)
  crossprod(psi)/NROW(psi)
}
```

More elaborate estimators are also available: `meatHAC()` computes an estimate based on the weighted (empirical) autocorrelations of the (empirical) estimating function. An extremely simplified version of `meatHAC()` is

```
meatHAC <- function(x, weights)
{
  psi <- estfun(x)
  n <- NROW(psi)

  rval <- 0.5 * crossprod(psi) * weights[1]
  for(i in 2:length(weights))
    rval <- rval + weights[i] * crossprod(psi[1:(n-i+1),], psi[i:n,])

  (rval + t(rval))/n
}
```

Now, for `meatHC()` we need a bit more infrastructure. Relying on Equation 8 we want to estimate  $M(\theta)$  by a matrix of type  $1/n X^\top \Omega X$  where  $X$  is the regressor matrix and  $\Omega$  is a diagonal matrix estimating the variance of  $r(y, \eta)$ . Various forms for this diagonal matrix  $\Omega$  have been suggested, they all depend on the vector of the observed  $(r(y_1, x_1^\top \hat{\theta}), \dots, r(y_n, x_n^\top \hat{\theta}))^\top$ , the hat values and the degrees of freedom. Instead of basing this on a separate generic for this type of residuals, it is also possible to recover it from the estimating function. As  $\psi(y_i, x_i, \hat{\theta}) = r(y_i, x_i^\top \hat{\theta}) \cdot x_i$ , we can simply divide the empirical estimating function by  $x_i$  and obtain the residual. Even simpler, if the model includes an intercept, the residuals are in the first column of the estimating function matrix. Consequently, we only need a way to extract the regressor matrix from the fitted model (which is simple, there is standardized infrastructure available for most models) and get the hat values. For the latter, there is a generic function `hatvalues()` which has many methods and missing new methods are usually easy to write. A condensed version of `meatHC()` can then be given as

```
meatHC <- function(x, omega)
{
  X <- if(is.matrix(x$x)) x$x else model.matrix(x)
  n <- NROW(X)
  k <- NCOL(X)

  diaghat <- hatvalues(x)
  res <- if(attr(terms(x), "intercept") > 0) estfun(x)[,1]
    else rowMeans(estfun(x)/X, na.rm = TRUE)

  if(is.function(omega)) omega <- omega(res, diaghat, n-k)
  rval <- sqrt(omega) * X

  crossprod(rval)/n
}
```

### *The bread*

Estimating  $B(\theta)$  is typically easier and not dealt with in most of the publications (I know) about sandwich estimators. My feeling is that it suffices to simply provide a generic function `bread()` which has a method for each class of fitted models. For “lm” and “glm” all the necessary information can be conveniently obtained from the `summary()` method:

```
bread.lm <- function(x, ...)
{
  sx <- summary(x)
  sx$cov.unscaled * as.vector(sum(sx$df[1:2]))
}
```

### *The sandwich*

Computing the sandwich is easy given the previous building blocks. Currently, the function `sandwich()` computes an estimate for  $1/n S(\theta)$  via

```
sandwich <- function(x, bread. = bread, meat. = meat, ...)
{
  if(is.function(bread.)) bread. <- bread.(x)
  if(is.function(meat.)) meat. <- meat.(x, ...)
  1/NROW(estfun(x)) * (bread. %*% meat. %*% bread.)
}
```

and `meat.` could also be set to `meatHAC` or `meatHC`.

Therefore, all that a userR/developerR would have to do to make a new class of models, “foo” say, fit for this framework is: provide an `estfun()` method `estfun.foo()` and a `bread()` method `bread.foo()`. See also Figure 1.

Only for HC estimators, it has to be assured in addition that

- the model only depends on a linear predictor (this cannot be easily checked by the software, but has to be done by the user),
- a `hatvalues().foo()` method exists (for HC2–HC4), and
- the model matrix  $X$  is available (via a `model.matrix().foo()` method).

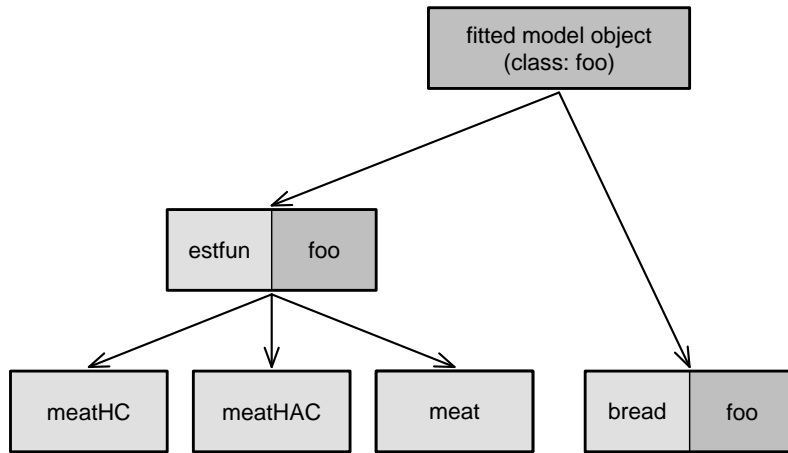


Figure 1: Structure of sandwich estimators

### 3. Open questions

#### *estfun()* extractors

To be able to plug in further models into the framework outlined above, we need extractor methods `estfun.foo()` for these model classes. In particular, it would be nice to provide such methods for “`survreg`” and “`coxph`” objects. As objects from both classes provide some information about their gradients, this seemed to be easy at a first glance. However, the results make me somewhat suspicious.

For “`survreg`” objects, there is a `type = "matrix"` argument to `residuals()` which returns a matrix with partial derivatives of the likelihood. Hence, this seemed to be a good first try

```

estfun.survreg <- function(x, ...)
{
  stopifnot(require(survival))
  X <- if (is.matrix(x$x)) x$x else model.matrix(terms(x), model.frame(x))
  attr(X, "assign") <- NULL
  wts <- if(!is.null(x$weights)) x$weights else 1
  res <- residuals(x, type = "matrix")
  rval <- as.vector(res[, "dg"]) * wts * X
}

```

```

if(NROW(x$var) > length(coefficients(x))) {
  rval <- cbind(rval, res[, "ds"])
  colnames(rval)[NCOL(rval)] <- "Log(scale)"
}
return(rval)
}

```

however, the "dg" component seems to be off by a factor of 2, whereas the "ds" component looks ok.

For “coxph” objects, computing the empirical estimating functions is taken on by the `type = "score"` residuals if I understand the documentation correctly. Hence:

```
estfun.coxph <- function(x, ...) residuals(x, type = "score", ...)
```

### *Appropriateness of HC estimators*

I am aware of only very few applications of HC estimators to regression models other than the linear regression model. Is there any theoretical or applied work that treats HC estimators in a general setup as done here? Are there arguments why HC estimators as used here might not be useful in other types of regression models?