Generalizing diagnostics from single-variance linear models

Many papers have proposed individual methods to look for outliers, check normality of residuals, detect influential observations, etc.

I know of two attempts to produce a <u>system</u> of diagnostics for MLMs. One is Yuan & Johnson (2012 *Biometrics*). I'll present the other one – Hodges (1998 *JRSSB*).

RWC and Lee et al (LNP) use residuals to assess non-linearity and non-constant error variance.

RWC seem unaware of problems with residuals (discussed below); LNP acknowledge them but blow them off.

I'll focus on Hodges (1998) with brief comments on RWC and LNP.

Some preliminary items

The main analysis is Bayesian; assume we have MCMC draws.

I follow Weisberg (1983) in seeking these qualities in diagnostics:

- A diagnostic should aim to detect a specific problem.
- Diagnostics should compute quickly.
- A diagnostic should have a corresponding plot, so you can assess the effect of individual observations.
- Graphics should let users to look at the data as directly as possible.

The last item argues against using a Bayesian approach to *diagnostics*, even when using a Bayesian approach to *analysis*.

This lecture uses the constraint-case formulation.

We'll use "the HMO dataset"

Assembled in the mid-1990s, it describes 341 HMOs serving U.S. Government employees.

It has plans in 42 states, DC, Guam, and Puerto Rico, all called "states."

• Between 1 and 31 plans per state, median 5.

Of particular interest: each HMO's monthly premium (\$) for individuals.

• We have some plan-level and state-level covariates.

The next slide summarizes the data.



For each state: the numeral is its number of plans, \bullet is its average premium, \triangle is its posterior mean in the simple analysis (next).

Begin with an unbalanced one-way RE model:

$$\begin{array}{rcl} y_{ij} &=& \theta_i + \epsilon_{ij}, & \text{for state } i = 1, \dots, 45 \\ \theta_i &=& \mu + \delta_i, & \text{for plan within state } j = 1, \dots, n_i \\ & \epsilon_{ij} \sim \mathcal{N}(0, \sigma_e^2), \, \delta_i \sim \mathcal{N}(0, \sigma_s^2); \, \text{flat prior on } \mu. \end{array}$$

Rewrite (1) as $0 = -\theta_i + \mu + \delta_i$ to give the constraint-case form:



y and ϵ are 341 \times 1, ordered with *j* changing fastest; δ is 45 \times 1.

Here's the simple model for the HMO data again:

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{0}_{45\times 1} \end{bmatrix} = \begin{bmatrix} \mathbf{1}_{n_1} & \cdots & \mathbf{0}_{n_1} \\ \vdots & \ddots & \vdots \\ \mathbf{0}_{n_{45}} & \cdots & \mathbf{1}_{n_{45}} \\ \hline \mathbf{-I}_{45} & \mathbf{1}_{45} \end{bmatrix} \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_{45} \\ \hline \mu \end{bmatrix} + \begin{bmatrix} \epsilon \\ \overline{\delta} \end{bmatrix}$$

To derive and present results in more generality, I'll use the form

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{0} \\ \mathbf{M} \end{bmatrix} = \begin{bmatrix} \mathbf{K} & \mathbf{0} \\ \mathbf{S}_1 & \mathbf{S}_2 \\ \mathbf{W}_1 & \mathbf{W}_2 \end{bmatrix} \begin{bmatrix} \mathbf{\Theta}_1 \\ \mathbf{\Theta}_2 \end{bmatrix} + \begin{bmatrix} \frac{\epsilon}{\delta} \\ \frac{\delta}{\xi} \end{bmatrix},$$

I'll summarize this as $\mathbf{Y} = \mathbf{H}\Theta + \mathbf{E}$, where $cov(\mathbf{E})$ has a simple form.

Added variable plots (AVPs) in the usual linear model

An AVP displays the evidence for adding a predictor to a linear model. Suppose we've included predictors in **A** and are considering **B**.

The usual linear model is $\mathbf{y} = \mathbf{A}\beta + \mathbf{B}\omega + \epsilon$.

Premultiply both sides of this by $\mathbf{I} - \mathbf{A}(\mathbf{A}'\mathbf{A})^{-1}\mathbf{A}'$:

$$\hat{\mathbf{e}} = (\mathbf{I} - \mathbf{A}(\mathbf{A}'\mathbf{A})^{-1}\mathbf{A}')\mathbf{B}\omega + (\mathbf{I} - \mathbf{A}(\mathbf{A}'\mathbf{A})^{-1}\mathbf{A}')\epsilon$$

If $E(\epsilon) = 0$ as usual, then

$$E(\hat{\mathbf{e}}) = (\mathbf{I} - \mathbf{A}(\mathbf{A}'\mathbf{A})^{-1}\mathbf{A}')\mathbf{B}\omega.$$

The AVP plots $\hat{\mathbf{e}}$ vs $(\mathbf{I} - \mathbf{A}(\mathbf{A}'\mathbf{A})^{-1}\mathbf{A}')\mathbf{B}$.

If **B** should be added, this will fit a line through the origin with slope ω .

AVPs for mixed linear models in constraint-case form

Predictors can be added for data or constraint cases.

The AVP for linear models can be applied directly to either kind of predictor using four steps (explained on following slides)

- $1. \ \mbox{Reformulate the candidate predictor and call it } {\bf B}.$
- 2. Write the model with the candidate predictor as $\mathbf{Y} = \mathbf{H}\Theta + \mathbf{B}\phi + \mathbf{E}$.
- 3. Pre-multiply that equation by $\Gamma^{-1/2}$ for a suitable Γ , to make it a homoscedastic-errors problem.
- 4. Draw the usual AVP.

 $\label{eq:step 1. Reformulate the candidate predictor and call it$ **B**. If**B** $_1 is the 341-vector of plan enrollments,$ **B** $is <math display="inline">[\mathbf{B}_1', \mathbf{0}_{1\times 45}]'$. If **B**_2 is the 45-vector of state populations, **B** is $[\mathbf{0}_{1\times 341}, \mathbf{B}_2']'$.

<u>Step 3.</u> Make it homoskedastic: The AVP is not too sensitive to Γ . Start with $\mathbf{Y} = \mathbf{H}\Theta + \mathbf{B}\omega + \mathbf{E}$; pre-multiply by $\Gamma^{-1/2}$ to give

 $\mathbf{Y} = \mathbf{H}\boldsymbol{\Theta} + \mathbf{B}\boldsymbol{\omega} + \mathbf{E},$

where sans-serif font indicates pre-multiplication by $\Gamma^{-1/2}$.

Pre-multiply by $I - H(H'H)^{-1}H'$, to give

$$\hat{\mathsf{E}} = (\mathbf{I} - \mathbf{V})\mathsf{Y} = (\mathbf{I} - \mathbf{V})\mathsf{B}\omega + (\mathbf{I} - \mathbf{V})\mathsf{E}.$$

For any $\mathbf{\Gamma}$, the AVP will show the correct slope ω .

Step 4. Draw the usual added-variable plot.

This plot includes the data- and constraint-case residuals (and prior-case residuals, if present).

The AVP's vertical axis is

$$\hat{\mathsf{E}} = \mathbf{\Gamma}^{-1/2} \hat{\mathbf{E}} = \mathbf{\Gamma}^{-1/2} (\mathbf{Y} - \mathbf{H} \hat{\Theta}).$$

This scaling puts the different kinds of residuals on the same scale.

It may seem odd to include the different residuals on the same plot

- but the data and constraint cases convey distinct information,
- though the two kinds of residuals are linearly related.

AVP for state average expenses per hospital admission



Slope 0.007 (SE 0.002), strongly influenced by Guam (two points at left).

Transforming the outcome or predictors

Transforming y or a predictor can remove curvature or interactions and can make residuals look normally distributed.

The Box-Cox transformation family for y:

$$y^{(\lambda)} = \left\{ egin{array}{cc} (y^\lambda - 1)/\lambda & \lambda
eq 0, \ \log(y) & \lambda = 0. \end{array}
ight.$$

Here's a graphical method to examine the evidence about λ :

- Expand $y^{(\lambda)}$ as a linear Taylor series around $\lambda = 1$.
- ▶ \Rightarrow a new predictor **B** with coefficient that's linear in λ .
 - And rews: $\mathbf{B}_l = \hat{y}_l \log(\hat{y}_l) \hat{y}_l + 1$; \hat{y}_l is the fit for untransformed y.
- Draw the added variable plot for B.

Andrews AVP for the HMO data



OLS fit: slope 0.33 (SE 0.04) $\Rightarrow \lambda = 0.67$ (SE 0.04).

Case influence: Delete a case; how do estimates change?

For the linear model $\mathbf{y} = \mathbf{A}\beta + \epsilon$, there's a handy updating formula.

If (I) indicates deleting cases indexed by I, then

$$(A'_{(I)}A_{(I)})^{-1} = (A'A)^{-1} + (A'A)^{-1}A'_{I}(I - H_{I})^{-1}A_{I}(A'A)^{-1}$$

where A_I is the deleted rows of A and $H_I = A_I (A'A)^{-1} A'_I$.

If only case *l* is deleted, and subscript (*l*) indicates deleting case *l*:

$$\hat{\beta}_{(l)} - \hat{\beta} = -(\mathbf{A}'\mathbf{A})^{-1}a_lr_l/(1-h_l)$$

where a_l is the l^{th} row of **A**, r_l is the l^{th} residual in the full-data fit, $h_l \in (0,1)$ is the l^{th} diagonal element of $\mathbf{A}(\mathbf{A}'\mathbf{A})^{-1}\mathbf{A}'$.

Case influence measures for simple linear models

Cook's distance measures how $\hat{\beta}$ changes from deleting case I,

$$d_{l} = (\hat{\beta}_{(l)} - \hat{\beta})' \mathbf{A}' \mathbf{A} (\hat{\beta}_{(l)} - \hat{\beta}) / p \hat{\sigma}^{2}$$
$$= r_{l}^{2} \frac{h_{l}}{(1 - h_{l})^{2}} \frac{1}{p \hat{\sigma}^{2}}$$

where **A** has *p* columns and $\hat{\sigma}^2$ is the full-data error-variance estimate. $\rightarrow d_l$ is large if both r_l^2 and the leverage h_l are large.

The change in $\hat{\sigma}^2$ from deleting case *I*:

$$\hat{\sigma}_{(l)}^2 - \hat{\sigma}^2 = \frac{1}{n-p-1} (\hat{\sigma}^2 - \frac{r_l^2}{1-h_l}).$$

Under the assumed model, $E(r_l^2) = (1 - h_l)\sigma^2$ and $E(r_l) = 0$

 \Rightarrow deleting case *I* reduces $\hat{\sigma}^2$ if r_I^2 is large relative to its variance.

These measures apply to MLM data and constraint cases

<u>Objection</u>: Constraint and prior cases are different from data cases. But any modeling choice inserts information into the analysis. And it's handy to use one diagnostic for all kinds of cases.

Deleting a prior case = setting the prior variance to infinity.

Deleting a constraint case:

Simple RE model for the HMO dataset: Deleting MN's constraint case implies MN's posterior mean doesn't shrink toward μ . (Boring.)

It may also affect the posteriors of θ_i from other states. (Boring.)

When the state-level model is more complicated, the influence of MN's constraint case on MN's θ_i reflects on the state-level model.

Single-parameter influence measures are better for MLMs

Cook's distance describes case influence on the mean-parameter vector.

For an MLM, the mean parameter vector has many elements; a case's influence is often felt strongly by one parameter and weakly by others.

- \Rightarrow Cook's distance can miss important effects.
- \Rightarrow Use a single-parameter analog to Cook's distance, "relative change":

RC(i, I) = Change in parameter *i*'s estimate from deleting case I \div full-data posterior SD of parameter *i* Case deletion has non-linear effects, thru the variances

Deleting a case changes the information about $\pmb{\Gamma}$ and thus affects the posterior mean of Θ non-linearly.

If Cook's distance is used with fixed $\pmb{\Gamma},$ it's linear in the omitted case.

This suggests re-running the MCMC for each deleted case, but doing so would be onerous computationally.

 \Rightarrow I considered two alternatives:

- ► Fix **Γ** anyway (i.e., ignore the non-linearity).
- ► Re-use the MCMC draws using importance-sampling weighting.

Importance sampling the MCMC draws

Label the MCMC sequence by $k = 1, \ldots, m$.

Let $\mathbf{Y}_{(I)}$ be \mathbf{Y} with case I deleted, $\mathbf{\Psi} = (\Theta, \mathbf{\Gamma})$, and g be a function.

Then
$$E(g(\Psi)|\mathbf{Y}_{(l)}) = \int g(\Psi)f(\Psi|\mathbf{Y}_{(l)})d\Psi$$

 $= \int g(\Psi)f(\Psi|\mathbf{Y})\frac{f(\Psi|\mathbf{Y}_{(l)})}{f(\Psi|\mathbf{Y})}d\Psi$
 $\approx \frac{1}{m}\sum_{k=1}^{m}g(\Psi_k)\frac{f(\Psi_k|\mathbf{Y}_{(l)})}{f(\Psi_k|\mathbf{Y})},$

where $\Psi_k, k = 1, \ldots, m$, is the MCMC sequence.

 $f(\Psi_k|\mathbf{Y}_{(l)})/f(\Psi_k|\mathbf{Y})$ is simple, involving one row of $\mathbf{Y} = \mathbf{H}\Theta + \mathbf{E}$.

How exactly case deletion is non-linear

For a simplified version of the HMO dataset ($n_i = 8$), make one of ME's observations an outlier by the amount on the horizontal axis:



Solid line: ME's average; dotted line: ME's posterior mean; bunch of lines: 44 other posterior means.

Fixing **Г** (linear approx'n) vs. importance sampling

Both work for cases with small influence, but we care about detecting cases with large influence.

The importance sampler understates case-deletion effects when it fails.

- It reweights the MCMC draws it has; if few are near the case-deleted estimate, importance sampling gives a too-small change.
- It's easy to make this happen in real data.

The linear approximation appears much less sensitive to the non-linearity.

Referee: "It is difficult to imagine an outlier big enough to produce this effect that you couldn't detect by casual examination of the data."

- ► The previous slide: The weird stuff happens when the outlier is $\approx 45\sigma_e$ from the average of ME's other plans.
- ► The approximation's fragility depends some on the sample sizes.

The linear approximation applied to the HMO data

For research purposes, I computed the "exact" relative changes (RCs).

For the simple one-way RE model, four constraint cases had |RC| > 2 on their own θ_i :

ME (RC 3.4), CT (RC 2.1), PR (RC -2.1), ND (RC -2.4)

The linear approximation was effectively exact.

The importance sampler was not so good:

• For |RC| < 1, it was accurate.

► The four biggest |*RC*| were substantially underestimated; the importance weights were dominated by one MCMC draw.

Residuals, the diagnostic everybody seems to like

For the ordinary linear model $\mathbf{y} = \mathbf{A}\beta + \epsilon$,

- the residuals are $\hat{\epsilon} = (\mathbf{I} \mathbf{A}(\mathbf{A}'\mathbf{A})^{-1}\mathbf{A}')\mathbf{y};$
- $E(\hat{\epsilon}) = 0$ and $\hat{\epsilon}'\hat{y} = 0$ by construction.

Residuals are plotted

- vs. fitted values to show outliers or non-constant variance;
- vs. candidate predictors: in AVPs; to check for non-linearity; or to check for variance that's a function of the predictor; and
- ▶ in quantile plots to check the assumed normal error distribution.

All these uses ignore correlations among the $\hat{\epsilon},$ which are usually small.

Some of these tools don't work in the constraint-case form

For the general model in constraint case form,

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{0} \\ \mathbf{M} \end{bmatrix} = \begin{bmatrix} \mathbf{K} & \mathbf{0} \\ \mathbf{S}_1 & \mathbf{S}_2 \\ \mathbf{W}_1 & \mathbf{W}_2 \end{bmatrix} \begin{bmatrix} \mathbf{\Theta}_1 \\ \mathbf{\Theta}_2 \end{bmatrix} + \begin{bmatrix} \frac{\epsilon}{\delta} \\ \frac{\delta}{\xi} \end{bmatrix},$$

the biggest problem is that for the constraint cases:

 $\begin{array}{ll} \mbox{residuals:} & 0 - \left(\bm{S}_1 \hat{\Theta}_1 + \bm{S}_2 \hat{\Theta}_2 \right) \\ \mbox{fitted values:} & \bm{S}_1 \hat{\Theta}_1 + \bm{S}_2 \hat{\Theta}_2. \end{array}$

These are the correct residuals for this plot; the correct fitted values are $S_2\hat{\Theta}_2$.

 \Rightarrow Unlike the previous three diagnostics, data and constraint cases need separate residual plots.

A different unhappy thing happens for data-case residuals



For this model & others, residuals & fitted values are positively correlated.

Here's how this looks in the HMO dataset

The solid line is the LS fit of studentized residuals on fitted values.



Residuals are biased and $\hat{\epsilon}'\hat{y} \neq 0$; what to do?

In MLMs generally, the fit is biased \Rightarrow residuals are biased.

- ► There's a tidy formula for the bias (RWC; Hodges 2013, Sec. 8.5).
- ▶ E.g., penalized splines: the fit is low at peaks, high in valleys.

Hilden-Minton (1995):

- ▶ Fit an unshrunk model and use the residuals from that.
- Construct linear transforms of the data-case residuals, l'Ê_d that are independent of each other and unbiased.

LNP (2006):

- ▶ Plot the residuals from the full fit, $\hat{e}_l = y_l X_l \hat{\beta} Z_l \hat{\mathbf{u}}$, vs. the fitted *fixed effects* $X_l \hat{\beta}$.
- Their result isn't true in the generality they claim.
- For many RPMs, $X_i \hat{\beta}$ is almost completely divorced from the fitted values of interest, e.g., penalized spline.

Hodges view: There's not much you can do

Probably the best solution is to try to ignore this structure in residuals.

 For ordinary linear models, we ignore correlations among the residuals.

It's *possible* to estimate the bias and remove it.

 Zhao, Hodges, Carlin (2017), in network meta-analysis, estimating inconsistency of direct and indirect evidence.

It's not clear that this is a good idea in general

Cunanan, Dai (2014, 2016 projects): Estimating bias of a penalized spline fit works poorly for constructing point-wise Cls.