

Quasi-likelihood and M-estimation

Patrick Breheny

December 9, 2024

Misspecified and partially specified models

- As you may recall, establishing the asymptotic normality of the MLE revolved around taking a Taylor series expansion of the score function
- This raises an interesting question: since our inferential methods (Score and Wald) rely entirely on the score and its properties. . . do we even need to specify the rest of the model?
- As we will see, answering this question (partial specification) also sheds light on what happens to the MLE when our model is wrong (misspecification)

Definition

- Let's formalize this idea: given data y_1, \dots, y_n , suppose we intend to estimate parameters θ by solving the equation

$$\sum_i \mathbf{u}(\theta|y_i) = \mathbf{0},$$

where $\mathbf{u} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is a known function

- Perhaps \mathbf{u} is the score function of some likelihood, but we are not bothering to specify that likelihood
- This idea goes by a few different names in the statistical literature:
 - Estimating equations*
 - Quasi-likelihood*
 - "M-estimation" (because it's kind of like an MLE)

Quasi-likelihood: Advantages

- Why might we choose to take this approach?
- One main reason is simplicity: in many applications such as longitudinal data, spatial statistics, and time series analysis, complex correlation structures are present and specifying a full likelihood is rather complex
- The other reason involves robustness: by focusing only on properties of the score, our results may hold for a wider class of models

Quasi-likelihood: Disadvantages

Obviously, there are also potential disadvantages:

- Our estimates may be less efficient (higher SE for a given sample size)
- Certain likelihood tools may be inaccessible, such as AIC and likelihood ratio tests
- Small-sample inference may be problematic; without an actual probability model, we have to rely on asymptotic approaches

Exponential dispersion families

- The term “quasi-likelihood” is typically used to refer to the application of this idea in the context of GLMs
- Recall that for an exponential dispersion family

$$\ell(\theta) \propto \frac{y\theta - \psi(\theta)}{\phi},$$

we have

$$\begin{aligned}\mathbb{E}(y) &= \nabla\psi(\theta) \equiv \mu \\ \mathbb{V}(y) &= \phi\nabla^2\psi(\theta) \equiv \phi v\end{aligned}$$

GLMs

- If we are in the modeling context where μ_i depends on a set of predictors \mathbf{x}_i through coefficients β , we have the score function

$$\sum_i \frac{\partial \theta_i}{\partial \beta} \frac{\partial \ell_i}{\partial \theta_i}$$

- Setting this equal to zero, we can rewrite the estimating equation so that it is solely a function of the mean and variance of y :

$$\phi^{-1} \sum_i \frac{\partial \mu_i}{\partial \beta} v_i^{-1} (y_i - \mu_i) = \mathbf{0}$$

Mean-variance modeling

- The appeal of this approach is that we can model

$$\mathbb{E}Y_i = \mu_i(\boldsymbol{\beta})$$

$$\mathbb{V}Y_i = \phi v(\mu_i)$$

without worrying about the full distribution of Y

- In other words, we can focus on modeling the mean and the only real distributional assumption we make is the mean-variance relationship $v(\mu_i)$

Generalized estimating equations

- These derivations are the same for multivariate outcomes, in which the estimating equations are

$$\sum_i \frac{\partial \boldsymbol{\mu}_i}{\partial \boldsymbol{\beta}} \mathbf{V}_i^{-1} (\mathbf{y}_i - \boldsymbol{\mu}_i) = \mathbf{0}$$

- In the multivariate context, this idea is known as *generalized estimating equations*, or GEE
- This is a popular approach for analyzing longitudinal data, and you will learn more about how it works in practice when you take Longitudinal Data Analysis

Properties of the “quasi-score”

- Does our usual likelihood theory hold for these quasi-likelihood models?
- Not by our previous arguments; recall that we needed a true likelihood (and some regularity conditions) to establish that $\mathbb{E}\mathbf{u}(\boldsymbol{\beta}^*) = \mathbf{0}$ and $\mathbb{V}\mathbf{u}(\boldsymbol{\beta}^*) = -\mathbb{E}\nabla\mathbf{u}(\boldsymbol{\beta}^*)$
- Let $\mathbf{u}_i(\boldsymbol{\beta}) = \phi^{-1}(\partial\mu_i/\partial\boldsymbol{\beta})v_i^{-1}(y_i - \mu_i)$, with $\mathbf{u}(\boldsymbol{\beta}) = \sum_i \mathbf{u}_i(\boldsymbol{\beta})$; what properties does this “quasi-score” statistic have?

Properties of the “quasi-score” (cont'd)

- As it turns out, $\mathbf{u}(\boldsymbol{\beta})$ has the same theoretical properties as the usual score:

$$\begin{aligned}\mathbb{E}\mathbf{u}(\boldsymbol{\beta}^*) &= \mathbf{0} \\ \mathbb{V}\mathbf{u}_i(\boldsymbol{\beta}^*) &= -\mathbb{E}\nabla\mathbf{u}_i(\boldsymbol{\beta}^*)\end{aligned}$$

- Thus, we can apply our previous theoretical arguments (again, assuming Lindeberg condition, an interior neighborhood, and a suitably smooth \mathbf{u}) to obtain the asymptotic distribution

$$(\mathbf{X}^\top \mathbf{W} \mathbf{X})^{1/2}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*) \xrightarrow{d} \mathbf{N}(\mathbf{0}, \mathbf{I}),$$

where \mathbf{W} is a diagonal matrix with entries $(\partial\mu_i/\partial\eta_i)^2/(\phi v_i)$

- One can also use a robust/sandwich estimator for the variance (we'll talk more about this shortly)

Poisson and quasi-Poisson

- To see an example of how this works, let's consider the Poisson distribution
- As you may have seen in other courses, the Poisson distribution is a convenient distribution for modeling counts, but in practice there are usually extra sources of variability such that the relationship $\text{Var} Y_i = \mathbb{E} Y_i$ often does not hold in practice
- A simple remedy is a quasi-Poisson model in which $\text{Var} Y_i = \phi \mu_i$

Quasi-Poisson: Estimates and standard errors

- Note that ϕ cancels out of the estimating equation – Poisson and quasi-Poisson models give the exact same estimates $\hat{\beta}$
- The standard errors, however, are different
- The variance-covariance matrix is $(\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1}$ in both cases, although
 - Poisson: $w_i = \mu_i$
 - Quasi-Poisson: $w_i = \mu_i / \phi$
- The dispersion parameter ϕ can be estimated with

$$\hat{\phi} = \frac{\sum_i (y_i - \mu_i)^2 / \mu_i}{n},$$

although typically $n - d$ is used to account for degrees of freedom

Simulation: Setup

- To see how this works, let's simulate some data in which the mean model is correct, but the variance is incorrect
- Specifically, let

$$\begin{aligned}g_i &\sim \text{Exp}(1) \\ \log(\mu_i) &= x_i\beta \\ Y_i|g_i &\sim \text{Pois}(\mu_i g_i)\end{aligned}$$

- Note that the quasi-Poisson model is also wrong here, but at least it has a dispersion parameter ϕ that allows for extra variability beyond what the model can account for

Simulation: Results

Over 1,000 independent replications, for 95% confidence intervals:

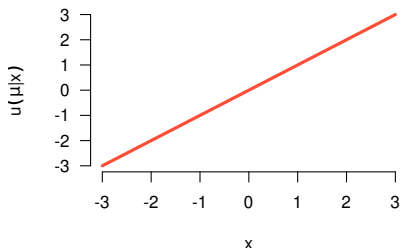
	Coverage	Average SE
Poisson	0.749	0.275
Quasi	0.939	0.445

General quasi-likelihood

- Thus far, we have considered quasi-likelihood exclusively as it pertains to regression models of the mean
- In the time we have left, let's look at this idea more broadly, without assuming that $\mathbf{u}(\boldsymbol{\theta})$ can be written in a form involving $y_i - \mu_i$ (in this context, the idea is often called M-estimation instead of quasi-likelihood)
- To make the discussion a bit more specific, we'll focus on the use of quasi-likelihood as it pertains to robust estimation

That non-robust mean

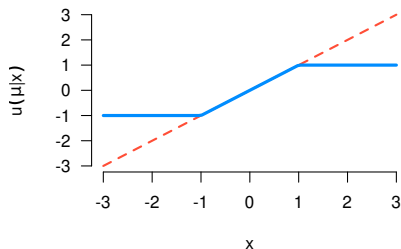
- As you know, the mean is not robust to outlying observations
- One way of visualizing this is to look at it as an M-estimate, with $u(\mu|x) = x - \mu$:



the influence that x has over the solution grows without bound as x becomes far from μ

The Huber function

- Consider instead the idea of “capping” the influence of x :



- This quasi-score function was proposed by Peter Huber in 1964

Remarks

- As you would imagine, the resulting M-estimate is much more accurate than the mean when outliers or contamination is present
- So is the median, of course, but one big advantage of the Huber estimate is that unlike the median, it is continuous in the sense that small changes to the data produce small changes in the estimate (unless we're in the capped region)
- The u function on the preceding slide would be the score function of a distribution that was normal near the mean, but at some point the tails of the distribution became exponential

Theoretical setup

- With this in mind as a potentially motivating example (there are many, many other examples of robust location estimators and u functions for a wide variety of problems), let's consider the theoretical properties of these estimators
- First, some notation:

$$\lambda_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_i u(\boldsymbol{\theta} | \mathbf{x}_i)$$

$$\lambda(\boldsymbol{\theta}) = \mathbb{E}u(\boldsymbol{\theta} | X)$$

- Note, of course, that $\lambda_n(\boldsymbol{\theta}) \xrightarrow{P} \lambda(\boldsymbol{\theta})$ for all $\boldsymbol{\theta}$

But what's θ ?

- Second, let's think about θ^* . . . what is it?
- We can't really think about it as the “true” value in the probability model, since we're not even specifying a full model anymore
- For the theory to work, we have to define θ^* as the (unique) solution to $\lambda(\theta) = \mathbf{0}$
- In the case of a misspecified likelihood, it can be shown that θ^* is the value of θ that minimizes the Kullback-Liebler distance to the true data generating process

Main result

Theorem: Let $\{x_i\}_{i=1}^n$ be an iid sample, with $\hat{\boldsymbol{\theta}}$ satisfying $\sum_i \mathbf{u}(\boldsymbol{\theta}|x_i) = \mathbf{0}$. Suppose

- (i) $\mathbf{u}(\boldsymbol{\theta}|x_i)$ is monotone
- (ii) $\lambda(\boldsymbol{\theta})$ is differentiable at $\boldsymbol{\theta}^*$ and $-\nabla\lambda(\boldsymbol{\theta}^*)$ is positive definite
- (iii) $\mathbb{E}\mathbf{u}(\boldsymbol{\theta})\mathbf{u}(\boldsymbol{\theta})^\top$ is finite and continuous in a neighborhood of $\boldsymbol{\theta}^*$
- (iv) $\nabla^2\mathbf{u}(\boldsymbol{\theta})$ are bounded in a neighborhood of $\boldsymbol{\theta}^*$

Then

$$\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) \xrightarrow{d} N(\mathbf{0}, \mathbf{A}^{-1}\mathbf{B}\mathbf{A}^{-1}),$$

where

$$\mathbf{A} = -\nabla\lambda(\boldsymbol{\theta}^*)$$

$$\mathbf{B} = \mathbb{E}\mathbf{u}(\boldsymbol{\theta}^*|X)\mathbf{u}(\boldsymbol{\theta}^*|X)^\top.$$

Remarks

- Condition (i) ensures that solutions to $\lambda(\boldsymbol{\theta}) = \mathbf{0}$ and $\lambda_n(\boldsymbol{\theta}) = \mathbf{0}$ are unique; as with the corresponding MLE theorem, this can be relaxed to a conclusion about the existence of an asymptotically normal solution
- For a correctly and fully specified model, $\mathbf{A} = \mathbf{B} = \mathcal{J}$ and we simply have the usual asymptotic normality result
- However, the “sandwich estimator” $\mathbf{A}^{-1}\mathbf{B}\mathbf{A}^{-1}$ is valid in a wider range of models, and is therefore often recommended as a better way of estimating variance in order to obtain inferential results that are (more) robust against model misspecification

Estimating \mathbf{A} and \mathbf{B}

- The empirical estimator of $\mathbf{A} = -\nabla\lambda(\boldsymbol{\theta}^*)$ is

$$\frac{1}{n} \sum_{i=1}^n \nabla u(\hat{\boldsymbol{\theta}}|x_i);$$

i.e., the average of the observed information

- The empirical estimator of $\mathbf{B} = \mathbb{E}\mathbf{u}(\boldsymbol{\theta}^*|X)\mathbf{u}(\boldsymbol{\theta}^*|X)^\top$ is

$$\frac{1}{n} \sum_{i=1}^n u(\hat{\boldsymbol{\theta}}|x_i)u(\hat{\boldsymbol{\theta}}|x_i)^\top;$$

for a correctly specified model, this would also be a consistent estimator for \mathcal{F} although it is less efficient than the standard estimators $\frac{1}{n}\mathcal{I}_n(\boldsymbol{\theta})$ and $\mathcal{F}(\boldsymbol{\theta})$