

Appendix: Proof of direction of Neyman's bias and counterexamples

D.M. Swanson^{1,3}, C.D. Anderson², and R.A. Betensky³

¹Oslo Centre for Biostatistics and Epidemiology, Oslo University Hospital, Oslo, NO 0407,

²Department of Neurology, Massachusetts General Hospital, 55 Fruit St, Boston, Massachusetts, 02114,

³Harvard School of Public Health, 655 Huntington Ave, Boston, Massachusetts, 02115.

In this supplementary material, we provide a theorem regarding the direction of Neyman's bias under certain modeling assumptions and examples of when Neyman's bias does or does not occur.

Theorem 1 *If G is associated with D such that $OR(t^*) \neq 1$, the distribution of $D \mid (G = 0)$ and $D \mid (G = 1)$ belong to the same location family, $pr(X > 0) = 1$, $pr(X < t^{**}) > 0$ (where t^{**} is defined as the time between t^* and the first possible presence of disease among the exposed or unexposed), and $X \perp\!\!\!\perp (D \mid G)^T$, then $OR_{ob}(t^*) \neq OR_{tr}(t^*)$. Specifically, if $D \mid (G = 0)$ is stochastically greater than $D \mid (G = 1)$ (alternatively, stochastically less than) so that exposure is a risk factor for disease (alternatively, protective against disease), then $OR_{ob}(t^*) < OR_{tr}(t^*)$ (alternatively, $OR_{ob}(t^*) > OR_{tr}(t^*)$).*

Proof Define $\partial F_{D \mid G=0}(x)/\partial x = f_0(x)$ and $\partial F_{D \mid G=1}(x)/\partial x = f_1(x)$, and suppose that $f_1(x) = f_0(x - k)$ for some k positive, without loss of generality. Such a scenario corresponds to exposure being protective against disease, though below we will also consider it a risk factor. $f_1(x)$ and $f_0(x)$ are in the same location family. Define $F(x)$ as the cumulative distribution function of X evaluated at x and remember $F(0) = 0$ and $F(t^*) > 0$. Consider the two quantities:

$$\frac{\int_0^{t^*} \{1 - F(t^* - x)\} f_0(x) \partial x}{\int_0^{t^*} f_0(x) \partial x} \quad \text{and} \quad \frac{\int_0^{t^*} \{1 - F(t^* - x)\} f_1(x) \partial x}{\int_0^{t^*} f_1(x) \partial x},$$

which we call the “percent erosion” of $\int_0^{t^*} f_0(x) \partial x$ and $\int_0^{t^*} f_1(x) \partial x$, respectively. Then

$$\begin{aligned} \frac{\int_0^{t^*} \{1 - F(t^* - x)\} f_1(x) \partial x}{\int_0^{t^*} f_1(x) \partial x} &= \frac{\int_0^{t^*} \{1 - F(t^* - x)\} f_0(x - k) \partial x}{\int_0^{t^*} f_0(x - k) \partial x} \\ &= \frac{\int_{-k}^{(t^*-k)} [1 - F\{t^* - (x + k)\}] f_0(x) \partial x}{\int_{-k}^{(t^*-k)} f_0(x) \partial x}. \end{aligned}$$

Since $F(\cdot)$ a cumulative distribution function and therefore increasing, we have

$$\begin{aligned} \frac{\int_0^{t^*} \{1 - F(t^* - x)\} f_1(x) \partial x}{\int_0^{t^*} f_1(x) \partial x} &= \frac{\int_{-k}^{(t^*-k)} [1 - F\{t^* - (x + k)\}] f_0(x) \partial x}{\int_{-k}^{(t^*-k)} f_0(x) \partial x} \\ &> \frac{\int_0^{t^*} \{1 - F(t^* - x)\} f_0(x) \partial x}{\int_0^{t^*} f_0(x) \partial x}, \quad (1) \end{aligned}$$

because at every “successive” ∂x in each integral, $1 - F\{t^* - (x + k)\} \geq 1 - F(t^* - x)$ and there is some $0 < x < t^*$ for which $1 - F\{t^* - (x + k)\} > 1 - F(t^* - x)$. Thus, the “percent erosion” of $f_0(x)$ will always be greater than that of $f_1(x) = f_0(x - k)$, which is intuitive since $f_1(\cdot)$ is located to the right of $f_0(\cdot)$ and thus subject to the corrosive

effects of $F(\cdot)$ for less “time.” Then using the inequality in (1),

$$\begin{aligned} 1 &> \left[\frac{\int_0^{t^*} (1 - F(t^* - x)) f_0(x) \partial x}{\int_0^{t^*} f_0(x) \partial x} \right] / \left[\frac{\int_0^{t^*} (1 - F(t^* - x)) f_1(x) \partial x}{\int_0^{t^*} f_1(x) \partial x} \right] \\ &= \frac{\int_0^{t^*} f_1(x) \partial x p}{\int_0^{t^*} f_0(x) \partial x (1 - p)} \times \frac{\int_0^{t^*} \{1 - F(t^* - x)\} f_0(x) \partial x (1 - p)}{\int_0^{t^*} \{1 - F(t^* - x)\} f_1(x) \partial x p} \\ &= \frac{pr(\text{Case, Exposed})}{pr(\text{Case, Unexposed})} \times \frac{pr(\text{Case, Unexposed, Observed})}{pr(\text{Case, Exposed, Observed})}, \end{aligned}$$

which implies that

$$\frac{pr(\text{Case, Exposed, Observed})}{pr(\text{Case, Unexposed, Observed})} > \frac{pr(\text{Case, Exposed})}{pr(\text{Case, Unexposed})} \quad \text{and} \quad OR_{ob}(t^*) > OR_{tr}(t^*)$$

since $pr(X > 0)$ implies $pr(\text{Control, Exposed, Observed}) = pr(\text{Control, Exposed})$ and $pr(\text{Control, Unexposed, Observed}) = pr(\text{Control, Unexposed})$. Again, these inequalities only hold when exposure is protective against disease. When exposure is a risk factor for disease and therefore shifts the mean age of disease onset to the left under the above assumptions,

$$\frac{pr(\text{Case, Exposed, Observed})}{pr(\text{Case, Unexposed, Observed})} < \frac{pr(\text{Case, Exposed})}{pr(\text{Case, Unexposed})} \quad \text{and} \quad OR_{ob}(t^*) < OR_{tr}(t^*)$$

using analogous results. So we see that the bias is not toward the null, but in a definite direction depending on model assumptions.

Example 1 Consider $D \mid (G = 1)$ uniform on $(0, 2)$, $D \mid (G = 0)$ uniform on $(0, 1)$, and X uniform on $(0, 3)$, independent of G . Clearly the distributions of disease for exposed and unexposed are not in the same location family in this case, and the model for X corresponds to disease-induced mortality necessarily occurring within 3 times units after disease, D . We need only consider cases when investigating the odds ratio since we assume $pr(X > 0) = 1$, implying $pr(D < M_d) = 1$. Taking $t^* = 1$,

$$\begin{aligned} \frac{pr(\text{Case, Exposed, Observed})}{pr(\text{Case, Unexposed, Observed})} &= \frac{\int_0^1 (2/3 + x/3) (1/2) p \partial x}{\int_0^1 (2/3 + x/3) 1 (1 - p) \partial x} \\ &= \frac{1/2 \int_0^1 (2/3 + x/3) p \partial x}{1 \int_0^1 (2/3 + x/3) (1 - p) \partial x} = \frac{1 p}{2 (1 - p)} = \frac{pr(\text{Case, Exposed})}{pr(\text{Case, Unexposed})}. \end{aligned}$$

So we have X independent of exposure status and time of disease-onset, as was the case above, but here $OR_{ob} = OR_{tr}$.

Example 2 Consider again $D \mid (G = 1)$ uniform on $(0, 2)$, and $D \mid (G = 0)$ uniform on $(0, 1)$. However, consider $X \mid (G = 1)$ uniform on $(0, 3)$ and $X \mid (G = 0)$ with density $f_{X|G=0}(x) = 2/3 (1 - x)^2$ on $[0, 1 + (9/2)^{1/3}]$. Again, we need only consider cases when investigating potential bias of the odds ratio since we assume $pr(D < M_d) = 1$ so that controls are not subject to the bias-inducing mortality event. Taking $t^* = 1$,

$$\begin{aligned} \frac{pr(\text{Case, Exposed, Observed})}{pr(\text{Case, Unexposed, Observed})} &= \frac{\int_0^1 (2/3 + x/3) (1/2) p \partial x}{\int_0^1 (7/9 + 2x^3/9) 1 (1 - p) \partial x} \\ &= \frac{1/2 \cdot \int_0^1 (2/3 + x/3) p \partial x}{1 \int_0^1 (7/9 + 2x^3/9) (1 - p) \partial x} = \frac{1/2 (5/6) p}{1 (5/6) (1 - p)} = \frac{1 p}{2 (1 - p)} = \frac{pr(\text{Case, Exposed})}{pr(\text{Case, Unexposed})}, \end{aligned}$$

and so here we have no bias again.

Example 3 Assume the same models of D conditional on G , and suppose $X \mid (G = 1)$ is uniform on $(0, 3)$ and $X \mid (G = 0)$ has density $f_{X|G=0}(x) = 5/2 (1 - x)^4$ on $[0, 1 + 2^{1/5}]$. For the reasons given above, we again only consider cases for investigating the bias of the odds ratio. Taking $t^* = 1$,

$$\begin{aligned} \frac{pr(\text{Case, Exposed, Observed})}{pr(\text{Case, Unexposed, Observed})} &= \frac{\int_0^1 (2/3 + x/3) (1/2) p \partial x}{\int_0^1 (1/2 + x^5/2) 1 (1 - p) \partial x} \\ &= \frac{1/2 \int_0^1 (2/3 + x/3) p \partial x}{1 \int_0^1 (1/2 + x^5/2) (1 - p) \partial x} = \frac{1/2 (5/6) p}{1 (7/12) (1 - p)} \neq \frac{1 p}{2 (1 - p)} = \frac{pr(\text{Case, Exposed})}{pr(\text{Case, Unexposed})}, \end{aligned}$$

and so here we have bias.

Example 4 Take $D \mid (G = 1)$ with density $f_{D|G=1}(x) = x^2/4$ on $[0, 12^{1/3}]$, $D \mid (G = 0)$ with density $f_{D|G=0}(x) = x/3$ $[0, 6^{1/2}]$. Then let $X \mid (G = 1)$ have density $f_{X|G=1}(x) = (2-x)^2/4$ on $[0, 2+4^{1/3}]$ and $X \mid (G = 0)$ be uniform on $[0, 2]$. As before, we need only consider cases when investigating the odds ratio since we assume $pr(D < M_d) = 1$ so that controls are not subject to the bias-inducing mortality event. Taking $t^* = 2$,

$$\begin{aligned} \frac{pr(\text{Case, Exposed, Observed})}{pr(\text{Case, Unexposed, Observed})} &= \frac{\int_0^2 (1/3 + 1/12 x^3) (x^2/4) p \partial x}{\int_0^2 (x/2) x/3 (1-p) \partial x} \\ &= \frac{(4/9) p}{4/9 (1-p)} = \frac{p \int_0^2 (x^2/4) \partial x}{(1-p) \int_0^2 x/3 \partial x} = \frac{p}{1-p} = \frac{pr(\text{Case, Exposed})}{pr(\text{Case, Unexposed})}. \end{aligned}$$

Remember that $pr(\text{Case, Exposed})/pr(\text{Case, Unexposed}) = p/(1-p)$ implies $OR_{tr}(t^*) = 1$ when $pr(D < M_d) = 1$, which is assumed from condition 3.

Example 5 On the other hand, we can obtain a biased odds ratio using the same conditional disease models as in the previous example and having $X \mid (G = 1)$ with density $f_{X|G=1}(x) = (2-x)^2/4$ on $[0, 2+4^{1/3}]$ and $X \mid (G = 0)$ uniform on $[0, 2]$. We again assume $pr(D < M_d) = 1$ from condition 3. Taking $t^* = 2$,

$$\begin{aligned} \frac{pr(\text{Case, Exposed, Observed})}{pr(\text{Case, Unexposed, Observed})} &= \frac{\int_0^2 (1/2 + 1/16 x^3) (x^2/4) p \partial x}{\int_0^2 (x/2) x/3 (1-p) \partial x} = \frac{p(1/2)}{(1-p)4/9} \\ &\neq \frac{(4/9) p}{4/9 (1-p)} = \frac{p \int_0^2 (x^2/4) \partial x}{(1-p) \int_0^2 x/3 \partial x} = \frac{p}{1-p} = \frac{pr(\text{Case, Exposed})}{pr(\text{Case, Unexposed})}. \end{aligned}$$