

Causal Inference in Case-Control Studies: Vignette

Overview

This vignette describes how to use package “ciccr” that is based on the paper entitled “Causal Inference in Case-Control Studies” (Jun and Lee, 2020).

Causal Inference on Relative and Attributable Risk

Case-Control Sampling

We first load the ciccr and MASS packages.

```
library(ciccr)
library(MASS)
```

To illustrate the usefulness of the package, we use the dataset ACS_CC that is included the package. This dataset is an extract from American Community Survey (ACS) 2018, restricted to white males residing in California with at least a bachelor’s degree. The ACS is an ongoing annual survey by the US Census Bureau that provides key information about US population. We use the following variables:

```
y = ACS_CC$topincome
t = ACS_CC$baplus
x = ACS_CC$age
```

- The binary outcome ‘Top Income’ (Y) is defined to be one if a respondent’s annual total pre-tax wage and salary income is top-coded. In our sample extract, the top-coded income bracket has median income \$565,000 and the next highest income that is not top-coded is \$327,000.
- The binary treatment (T) is defined to be one if a respondent has a master’s degree, a professional degree, or a doctoral degree.
- The covariate (X) is age in years and is restricted to be between 25 and 70.

The original ACS sample is not a case-control sample but we construct one by the following procedure.

1. The case sample ($Y = 1$) is composed of 921 individuals whose income is top-coded.
2. The control sample ($Y = 0$) of equal size is randomly drawn without replacement from the pool of individuals whose income is not top-coded.

We now construct cubic b-spline terms with three inner knots using the age variable.

```
x = splines::bs(x, df = 6)
```

Causal Inference on Relative Risk Using Case-Control Samples

Define $\beta(y) = E[\log \text{OR}(X)|Y = y]$ for $y = 0, 1$, where $\text{OR}(x)$ is the odds ratio conditional on $X = x$:

$$\text{OR}(x) = \frac{P(T = 1|Y = 1, X = x) P(T = 0|Y = 0, X = x)}{P(T = 0|Y = 1, X = x) P(T = 1|Y = 0, X = x)}.$$

Using the retrospective sieve logistic regression model, we estimate $\beta(1)$ by

```

results_case = avg_RR_logit(y, t, x, 'case')
results_case$est
#>      y
#> 0.729
results_case$se
#>      y
#> 0.101

```

Here, option 'case' refers to conditioning on $Y = 1$.

Similarly, we estimate $\beta(0)$ by

```

results_control = avg_RR_logit(y, t, x, 'control')
results_control$est
#>      y
#> 0.547
results_control$se
#>      y
#> 0.152

```

Here, option 'control' refers to conditioning on $Y = 0$.

We carry out causal inference on relative risk by

```

results = cicc_RR(y, t, x, 'cc', 0.95)

```

Here, 'cc' refers to case-control sampling and 0.95 refers to the level of the uniform confidence band (0.95 is the default choice).

The S3 object **results** contains estimates **est**, standard errors **se**, and one-sided confidence bands **ci** at $p = 0$ and $p = 1$.

```

# point estimates
results$est
#>      y      y
#> 0.547 0.729
# standard errors
results$se
#>      y      y
#> 0.152 0.101
# confidence intervals
results$ci
#>      y      y
#> 0.845 1.026

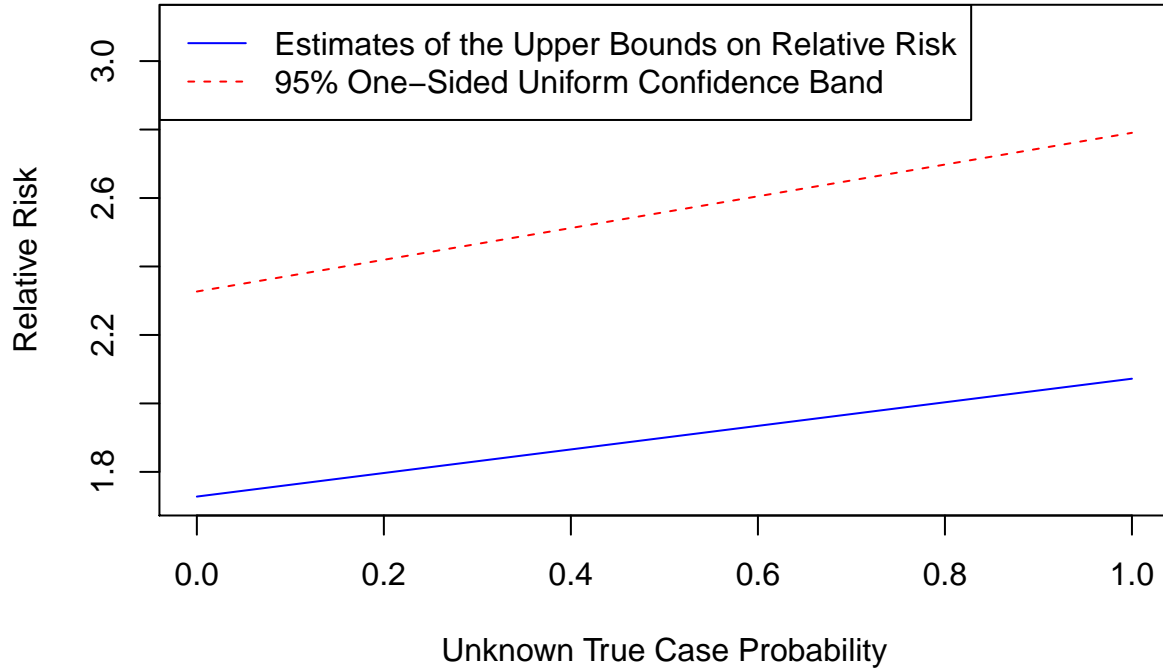
```

It is handy to examine the results by plotting a graph.

```

cicc_plot(results)

```



To interpret the results, we assume both marginal treatment response (MTR) and marginal treatment selection (MTS). In this setting, MTR means that everyone will not earn less by obtaining a degree higher than bachelor's degree; MTS indicates that those who selected into higher education have higher potential to earn top incomes. Based on the MTR and MTS assumptions, we can conclude that the treatment effect lies in between 1 and the upper end point of the one-sided confidence interval with high probability. Thus, the estimates in the graph above suggest that the effect of obtaining a degree higher than bachelor's degree is anywhere between 1 and the upper end points of the uniform confidence bands. This roughly implies that the chance of earning top incomes may increase up to by a factor as large as the upper end points of the uniform confidence band, but allowing for possibility of no positive effect at all. The results are shown over the range of the unknown true case probability. See Jun and Lee, 2020 for more detailed explanations regarding how to interpret the estimation results.

Comparison with Logistic Regression

We can compare these results with estimates obtained from logistic regression.

```
logit = stats::glm(y~t+x, family=stats::binomial("logit"))
est_logit = stats::coef(logit)
ci_logit = stats::confint(logit, level = 0.9)
#> Waiting for profiling to be done...
# point estimate
exp(est_logit)
#> (Intercept)      t      x1      x2      x3      x4
#>    0.0546    2.0612    4.4218    12.9960    19.0396    26.8357
#>      x5      x6
#>    6.4238    26.1436
# confidence interval
exp(ci_logit)
#>      5 % 95 %
#> (Intercept) 0.0196 0.13
#> t          1.7517 2.43
#> x1          1.0568 21.66
#> x2          5.5058 33.89
```

```
#> x3      6.7946 61.33
#> x4     10.2294 78.74
#> x5      2.0054 22.85
#> x6      8.6698 87.63
```

Here, the relevant coefficient is 2.06 (t) and its two-sided 90% confidence interval is [1.75, 2.43]. If we assume strong ignorability, the treatment effect is about 2 and its two-sided confidence interval is between [1.75, 2.43]. However, it is unlikely that the higher BA treatment satisfies the strong ignorability condition.

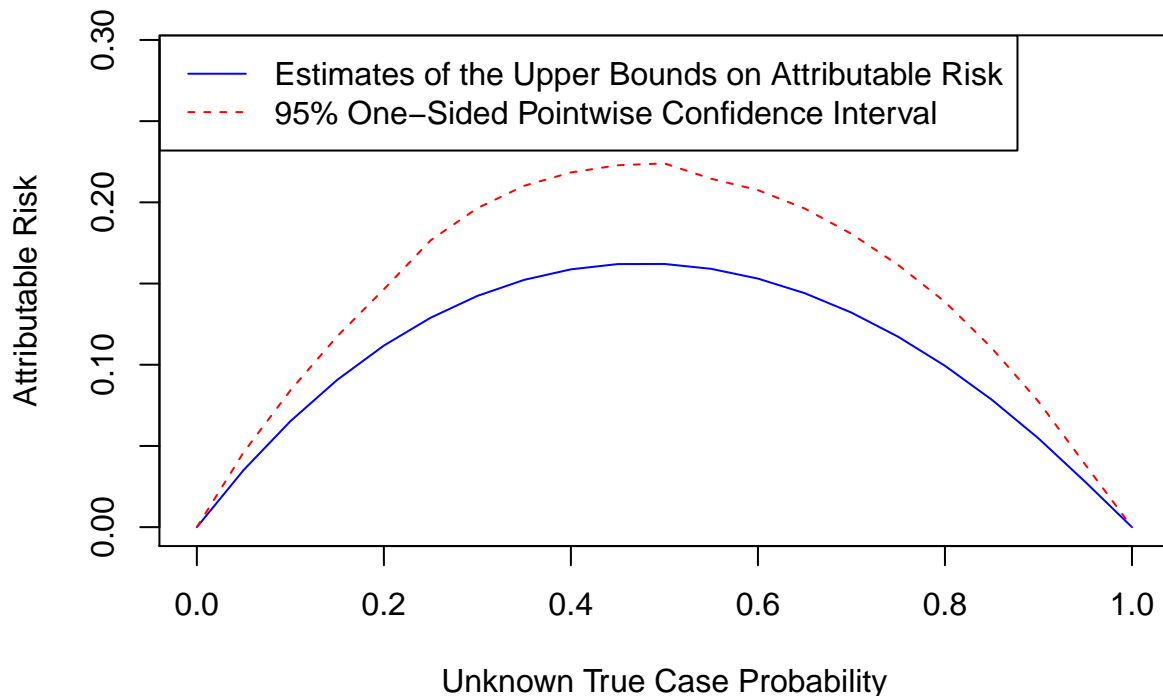
Causal Inference on Attributable Risk Using Case-Control Samples

We now consider attributable risk, that is the absolute difference in probabilities. We carry out causal inference on attributable risk by

```
results_AR = cicc_AR(y, t, x, sampling = 'cc', no_boot = 100)
```

The results can be plotted as before.

```
cicc_plot(results_AR, parameter = 'AR')
```



The upper bounds are approximately inverted-U shaped. When $p = 0$ or $p = 1$, there could no causal effect; the upper bound is maximized around $p = 0.5$.

Causal Inference on Relative Risk Using Case-Population Samples

We now consider an example of case-population samples. For this purpose, we use the dataset ACS_CP that is included in the package. This dataset is again an extract from American Community Survey (ACS) 2018. The original ACS sample is not a case-population sample but we construct one by the following procedure.

1. The case sample ($Y = 1$) is composed of 921 individuals whose income is top-coded.
2. The control sample ($Y = 0$) of equal size is randomly drawn with replacement from all observations and its top-coded status is coded missing.

We use the following variables:

```

y = ACS_CP$topincome
t = ACS_CP$baplust
x = ACS_CP$age

```

We print y to see how the outcome variable is coded.

```

print(head(y))
#> [1] NA  1  1 NA NA NA

```

We now code missing Y by 0 in the population sample.

```

y = as.integer(is.na(y)==FALSE)

```

We estimate $\beta(0)$ by

```

results_control = avg_RR_logit(y, t, x, 'control')
results_control$est
#>      y
#> 0.609
results_control$se
#>      y
#> 0.0987

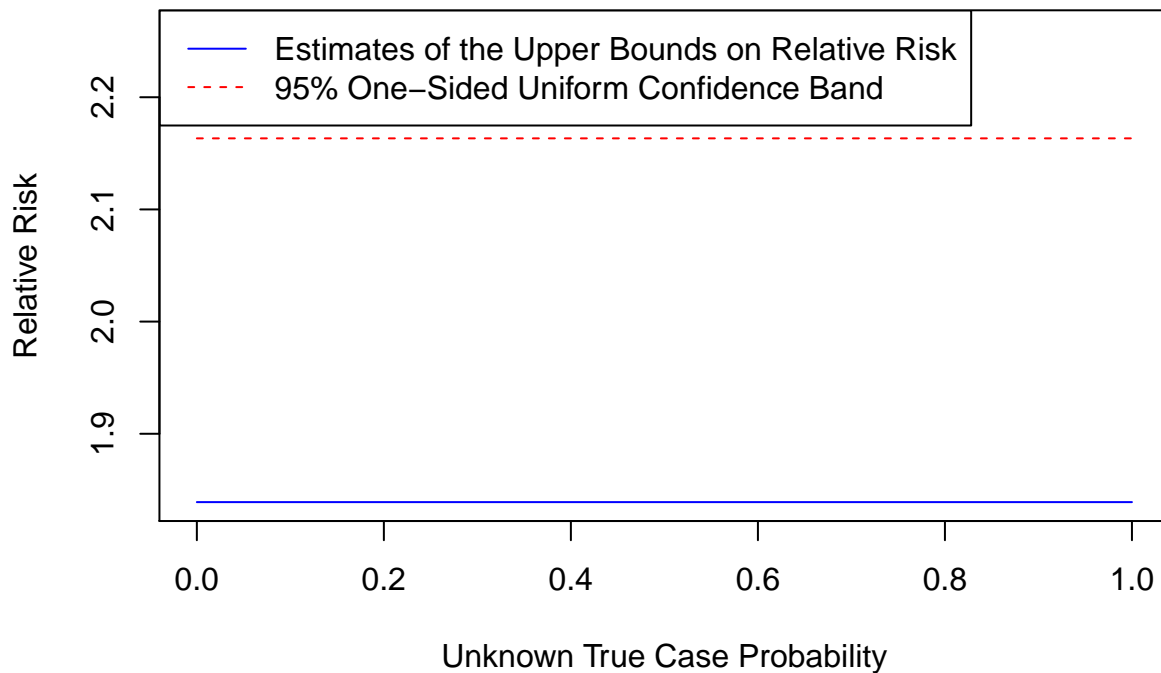
```

and carry out causal inference on relative risk by

```

results = cicc_RR(y, t, x, 'cp', 0.95)
cicc_plot(results)

```



Note that the estimates and upper bounds are constant across the unknown true case probability. This is because they do not depend on the value of the case probability in the case-population sample.

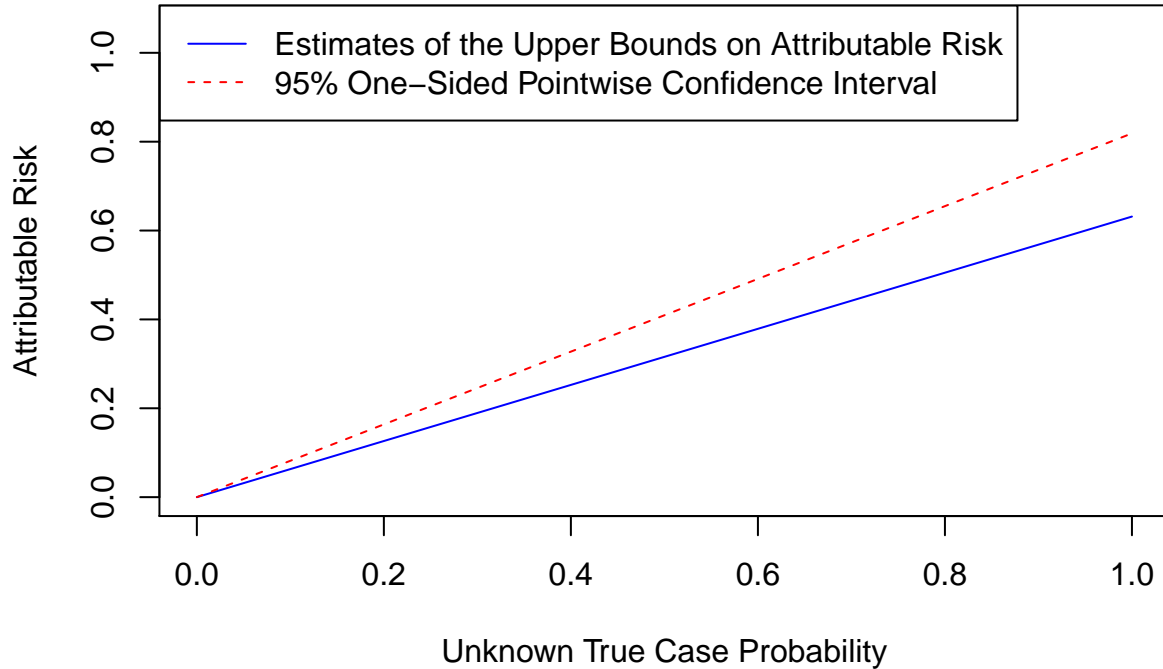
Causal Inference on Attributable Risk Using Case-Population Samples

We now consider causal inference on attributable risk by

```
results_AR = cicc_AR(y, t, x, sampling = 'cp', no_boot = 100)
```

The results can be plotted as before.

```
cicc_plot(results_AR, parameter = 'AR')
```



For case-population sampling, the upper bound on attributable risk is a linear function of the unknown true case probability; as a result, it increases linearly, as shown in the figure.

References

- Sung Jae Jun and Sokbae Lee. Causal Inference in Case-Control Studies. <https://arxiv.org/abs/2004.08318>.
- Manski, C.F. (1997). Monotone Treatment Response. *Econometrica*, 65(6), 1311-1334.
- Manski, C.F. and Pepper, J.V. (2000). Monotone Instrumental Variables: With an Application to the Returns to Schooling. *Econometrica*, 68(4), 997-1010.