

Two-phase designs in epidemiology

Thomas Lumley

March 28, 2025

This document explains how to analyse case-cohort and two-phase case-control studies with the “survey” package, using examples from <http://faculty.washington.edu/norm/software.html>. Some of the examples were published by Breslow & Chatterjee (1999).

The data are relapse rates from the National Wilm’s Tumor Study (NWTs). Wilm’s Tumour is a rare cancer of the kidney in children. Intensive treatment cures the majority of cases, but prognosis is poor when the disease is advanced at diagnosis and for some histological subtypes. The histological characterisation of the tumour is difficult, and histological group as determined by the NWTs central pathologist predicts much better than determinations by local institution pathologists. In fact, local institution histology can be regarded statistically as a pure surrogate for the central lab histology.

In these examples we will pretend that the (binary) local institution histology determination (`instit`) is available for all children in the study and that the central lab histology (`histol`) is obtained for a probability sample of specimens in a two-phase design. We treat the initial sampling of the study as simple random sampling from an infinite superpopulation. We also have data on disease stage, a four-level variable; on relapse; and on time to relapse.

Case-control designs

Breslow & Chatterjee (1999) use the NWTs data to illustrate two-phase case-control designs. The data are available at <http://faculty.washington.edu/norm/software.html> in compressed form; we first expand to one record per patient.

```
> library(survey)
> load(system.file("doc","nwts.rda",package="survey"))
> nwtsnb<-nwts
> nwtsnb$case<-nwts$case-nwtsb$case
> nwtsnb$control<-nwts$control-nwtsb$control
> a<-rbind(nwtsb,nwtsnb)
> a$in.ccs<-rep(c(TRUE,FALSE),each=16)
> b<-rbind(a,a)
> b$rel<-rep(c(1,0),each=32)
> b$n<-ifelse(b$rel,b$case,b$control)
> index<-rep(1:64,b$n)
> nwt.exp<-b[index,c(1:3,6,7)]
> nwt.exp$id<-1:4088
```

As we actually do know `histol` for all patients we can fit the logistic regression model with full sampling to compare with the two-phase analyses

```
> glm(rel~factor(stage)*factor(histol), family=binomial, data=nwt.exp)
```

```
Call: glm(formula = rel ~ factor(stage) * factor(histol), family = binomial,
data = nwt.exp)
```

Coefficients:

(Intercept)		factor(stage)2
-2.7066		0.7679
factor(stage)3		factor(stage)4
0.7747		1.0506
factor(histol)2	factor(stage)2:factor(histol)2	
1.3104		0.1477
factor(stage)3:factor(histol)2	factor(stage)4:factor(histol)2	
0.5942		1.2619

Degrees of Freedom: 4087 Total (i.e. Null); 4080 Residual

Null Deviance: 3306

Residual Deviance: 2943 AIC: 2959

The second phase sample consists of all patients with unfavorable histology as determined by local institution pathologists, all cases, and a 20% sample of the remainder. Phase two is thus a stratified random sample without replacement, with strata defined by the interaction of `instit` and `rel`.

```
> dccs2<-twophase(id=list(~id,~id),subset=~in.ccs,
+ strata=list(NULL,~interaction(instit,rel)),data=nwt.exp)
> summary(svyglm(rel~factor(stage)*factor(histol),family=binomial,design=dccs2))
```

Call:

```
svyglm(formula = rel ~ factor(stage) * factor(histol), design = dccs2,
family = binomial)
```

Survey design:

```
twophase(id = list(~id, ~id), subset = ~in.ccs, strata = list(NULL,
~interaction(instit, rel)), data = nwt.exp)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-2.5701	0.1288	-19.955	< 2e-16 ***
factor(stage)2	0.5482	0.1979	2.769	0.005708 **
factor(stage)3	0.4791	0.2032	2.359	0.018515 *
factor(stage)4	1.0037	0.2592	3.872	0.000114 ***
factor(histol)2	1.3505	0.3107	4.346	1.51e-05 ***
factor(stage)2:factor(histol)2	0.1152	0.4410	0.261	0.793876
factor(stage)3:factor(histol)2	0.5066	0.4241	1.194	0.232548
factor(stage)4:factor(histol)2	0.9785	0.6214	1.575	0.115615

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1.000876)

Number of Fisher Scoring iterations: 4

Disease stage at the time of surgery is also recorded. It could be used to further stratify the sampling, or, as in this example, to post-stratify. We can analyze the data either pretending that the sampling was stratified or using `calibrate` to post-stratify the design.

```
> dccs8<-twophase(id=list(~id,~id),subset=~in.ccs,
+               strata=list(NULL,~interaction(instit,stage,rel)),data=nwt.exp)
> gccs8<-calibrate(dccs2,phase=2,formula=~interaction(instit,stage,rel))
> summary(svyglm(rel~factor(stage)*factor(histol),family=binomial,design=dccs8))
```

Call:

```
svyglm(formula = rel ~ factor(stage) * factor(histol), design = dccs8,
       family = binomial)
```

Survey design:

```
twophase(id = list(~id, ~id), subset = ~in.ccs, strata = list(NULL,
~interaction(instit, stage, rel)), data = nwt.exp)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-2.71604	0.10827	-25.085	< 2e-16 ***
factor(stage)2	0.78141	0.14726	5.306	1.35e-07 ***
factor(stage)3	0.80093	0.15250	5.252	1.80e-07 ***
factor(stage)4	1.07293	0.17817	6.022	2.33e-09 ***
factor(histol)2	1.45836	0.31780	4.589	4.96e-06 ***
factor(stage)2:factor(histol)2	-0.04743	0.43495	-0.109	0.913
factor(stage)3:factor(histol)2	0.28064	0.41298	0.680	0.497
factor(stage)4:factor(histol)2	0.90983	0.63774	1.427	0.154

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1.000876)

Number of Fisher Scoring iterations: 4

```
> summary(svyglm(rel~factor(stage)*factor(histol),family=binomial,design=gccs8))
```

Call:

```
svyglm(formula = rel ~ factor(stage) * factor(histol), design = gccs8,
       family = binomial)
```

Survey design:

```
calibrate(dccs2, phase = 2, formula = ~interaction(instit, stage,
rel))
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-2.71604	0.10878	-24.968	< 2e-16 ***
factor(stage)2	0.78141	0.14729	5.305	1.35e-07 ***
factor(stage)3	0.80093	0.15212	5.265	1.68e-07 ***
factor(stage)4	1.07293	0.17905	5.993	2.77e-09 ***
factor(histol)2	1.45836	0.31757	4.592	4.88e-06 ***

```
factor(stage)2:factor(histol)2 -0.04743    0.43432   -0.109    0.913
factor(stage)3:factor(histol)2  0.28064    0.41231    0.681    0.496
factor(stage)4:factor(histol)2  0.90983    0.63187    1.440    0.150
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

(Dispersion parameter for binomial family taken to be 1.000876)

Number of Fisher Scoring iterations: 4

Case-cohort designs

In the case-cohort design for survival analysis, a $P\%$ sample of a cohort is taken at recruitment for the second phase, and all participants who experience the event (cases) are later added to the phase-two sample.

Viewing the sampling design as progressing through time in this way, as originally proposed, gives a double sampling design at phase two. It is simpler to view the process *sub specie aeternitatis*, and to note that cases are sampled with probability 1, and controls with probability $P/100$. The subcohort will often be determined retrospectively rather than at recruitment, giving stratified random sampling without replacement, stratified on case status. If the subcohort is determined prospectively we can use the same analysis, post-stratifying rather than stratifying.

There have been many analyses proposed for the case-cohort design (Therneau & Li, 1999). We consider only those that can be expressed as a Horvitz-Thompson estimator for the Cox model.

First we load the data and the necessary packages. The version of the NWTs data that includes survival times is not identical to the data set used for case-control analyses above.

```
> library(survey)
> library(survival)
> data(nwtco)
> ntwco<-subset(nwtco, !is.na(edrel))
```

Again, we fit a model that uses `histol` for all patients, to compare with the two-phase design

```
> coxph(Surv(edrel, rel)~factor(stage)+factor(histol)+I(age/12),data=ntwco)
```

Call:

```
coxph(formula = Surv(edrel, rel) ~ factor(stage) + factor(histol) +
      I(age/12), data = ntwco)
```

	coef	exp(coef)	se(coef)	z	p
factor(stage)2	0.66730	1.94898	0.12156	5.490	4.03e-08
factor(stage)3	0.81737	2.26455	0.12077	6.768	1.31e-11
factor(stage)4	1.15373	3.16999	0.13490	8.553	< 2e-16
factor(histol)2	1.58389	4.87387	0.08869	17.859	< 2e-16
I(age/12)	0.06789	1.07025	0.01492	4.549	5.39e-06

```
Likelihood ratio test=395.4  on 5 df, p=< 2.2e-16
n= 4028, number of events= 571
```

We define a two-phase survey design using simple random superpopulation sampling for the first phase, and sampling without replacement stratified on `rel` for the second phase. The `subset`

argument specifies that observations are in the phase-two sample if they are in the subcohort or are cases. As before, the data structure is rectangular, but variables measured at phase two may be NA for participants not included at phase two.

We compare the result to that given by `survival::cch` for Lin & Ying's (1993) approach to the case-cohort design.

```
> (dcch<-twophase(id=list(~seqno,~seqno), strata=list(NULL,~rel),
+               subset=~I(in.subcohort | rel), data=nwtco))
```

Two-phase sparse-matrix design:

```
twophase(id = list(~seqno, ~seqno), strata = list(NULL, ~rel),
        subset = ~I(in.subcohort | rel), data = nwtco)
```

Phase 1:

Independent Sampling design (with replacement)

```
svydesign(ids = ~seqno)
```

Phase 2:

Stratified Independent Sampling design

```
svydesign(ids = ~seqno, strata = ~rel, fpc = `*phase1*`)
```

```
> svycoxph(Surv(edrel,rel)~factor(stage)+factor(histol)+I(age/12),
+          design=dcch)
```

Call:

```
svycoxph(formula = Surv(edrel, rel) ~ factor(stage) + factor(histol) +
        I(age/12), design = dcch)
```

	coef	exp(coef)	se(coef)	robust se	z	p
factor(stage)2	0.69266	1.99902	0.22688	0.16279	4.255	2.09e-05
factor(stage)3	0.62685	1.87171	0.22873	0.16823	3.726	0.000194
factor(stage)4	1.29951	3.66751	0.25017	0.18898	6.877	6.13e-12
factor(histol)2	1.45829	4.29861	0.16844	0.14548	10.024	< 2e-16
I(age/12)	0.04609	1.04717	0.02732	0.02302	2.002	0.045233

Likelihood ratio test= on 5 df, p=

n= 1154, number of events= 571

```
> subcoh <- nwtco$in.subcohort
> selccoh <- with(nwtco, rel==1/subcoh==1)
> ccoh.data <- nwtco[selccoh,]
> ccoh.data$subcohort <- subcoh[selccoh]
> cch(Surv(edrel, rel) ~ factor(stage) + factor(histol) + I(age/12),
+     data=ccoh.data, subcoh = ~subcohort, id=~seqno,
+     cohort.size=4028, method="LinYing")
```

Case-cohort analysis,x\$method, LinYing

with subcohort of 668 from cohort of 4028

```
Call: cch(formula = Surv(edrel, rel) ~ factor(stage) + factor(histol) +
        I(age/12), data = ccoh.data, subcoh = ~subcohort, id = ~seqno,
        cohort.size = 4028, method = "LinYing")
```

Coefficients:

	Value	SE	Z	p
factor(stage)2	0.69265646	0.16287906	4.252581	2.113204e-05
factor(stage)3	0.62685179	0.16746144	3.743260	1.816478e-04
factor(stage)4	1.29951229	0.18973707	6.849016	7.436052e-12
factor(histol)2	1.45829267	0.14429553	10.106291	0.000000e+00
I(age/12)	0.04608972	0.02230861	2.066006	3.882790e-02

Barlow (1994) proposes an analysis that ignores the finite population correction at the second phase. This simplifies the standard error estimation, as the design can be expressed as one-phase stratified superpopulation sampling. The standard errors will be somewhat conservative. More data preparation is needed for this analysis as the weights change over time.

```
> nwtco$eventrec<-rep(0,nrow(nwtco))
> nwtco.extra<-subset(nwtco, rel==1)
> nwtco.extra$eventrec<-1
> nwtco.expd<-rbind(subset(nwtco,in.subcohort==1),nwtco.extra)
> nwtco.expd$stop<-with(nwtco.expd,
+                         ifelse(rel & !eventrec, edrel-0.001,edrel))
> nwtco.expd$start<-with(nwtco.expd,
+                         ifelse(rel & eventrec, edrel-0.001, 0))
> nwtco.expd$event<-with(nwtco.expd,
+                         ifelse(rel & eventrec, 1, 0))
> nwtco.expd$ppts<-ifelse(nwtco.expd$event, 1, 1/with(nwtco,mean(in.subcohort | rel)))
```

The analysis corresponds to a cluster-sampled design in which individuals are sampled stratified by subcohort membership and then time periods are sampled stratified by event status. Having individual as the primary sampling unit is necessary for correct standard error calculation.

```
> (dBarlow<-svydesign(id=~seqno+eventrec, strata=~in.subcohort+rel,
+                   data=nwtco.expd, weight=~ppts))
```

Stratified 2 - level Cluster Sampling design (with replacement)

With (1154, 1239) clusters.

```
svydesign(id = ~seqno + eventrec, strata = ~in.subcohort + rel,
         data = nwtco.expd, weight = ~ppts)
```

```
> svycoxph(Surv(start,stop,event)~factor(stage)+factor(histol)+I(age/12),
+          design=dBarlow)
```

Call:

```
svycoxph(formula = Surv(start, stop, event) ~ factor(stage) +
         factor(histol) + I(age/12), design = dBarlow)
```

	coef	exp(coef)	se(coef)	robust se	z	p
factor(stage)2	0.73589	2.08734	0.18571	0.16985	4.333	1.47e-05
factor(stage)3	0.59763	1.81780	0.18876	0.17529	3.409	0.000651
factor(stage)4	1.39068	4.01757	0.20500	0.20777	6.693	2.18e-11
factor(histol)2	1.50450	4.50191	0.13945	0.16407	9.170	< 2e-16
I(age/12)	0.04315	1.04410	0.02228	0.02425	1.779	0.075191

Likelihood ratio test= on 5 df, p=
n= 1239, number of events= 571

In fact, as the finite population correction is not being used the second stage of the cluster sampling could be ignored. We can also produce the stratified bootstrap standard errors of Wacholder et al (1989), using a replicate weights analysis

```
> (dWacholder <- as.svrepdesign(dBarlow,type="bootstrap",replicates=500))

Call: as.svrepdesign.default(dBarlow, type = "bootstrap", replicates = 500)
Survey bootstrap with 500 replicates.

> svycoxph(Surv(start,stop,event)~factor(stage)+factor(histol)+I(age/12),
+          design=dWacholder)

Call:
svycoxph.svyrep.design(formula = Surv(start, stop, event) ~ factor(stage) +
  factor(histol) + I(age/12), design = dWacholder)
```

	coef	exp(coef)	se(coef)	z	p
factor(stage)2	0.73589	2.08734	0.17278	4.259	2.05e-05
factor(stage)3	0.59763	1.81780	0.17845	3.349	0.000811
factor(stage)4	1.39068	4.01757	0.20783	6.691	2.21e-11
factor(histol)2	1.50450	4.50191	0.16131	9.327	< 2e-16
I(age/12)	0.04315	1.04410	0.02522	1.711	0.087090

Likelihood ratio test=NA on 5 df, p=NA
n= 1239, number of events= 571

Exposure-stratified designs

Borgan et al (2000) propose designs stratified or post-stratified on phase-one variables. The examples at <http://faculty.washington.edu/norm/software.html> use a different subcohort sample for this stratified design, so we load the new subcohort variable

```
> load(system.file("doc","nwtco-subcohort.rda",package="survey"))
> nwtco$subcohort<-subcohort
> d_BorganII <- twophase(id=list(~seqno,~seqno),
+                        strata=list(NULL,~interaction(instit,rel)),
+                        data=nwtco, subset=~I(rel |subcohort))
> (b2<-svycoxph(Surv(edrel,rel)~factor(stage)+factor(histol)+I(age/12),
+              design=d_BorganII))
```

```
Call:
svycoxph(formula = Surv(edrel, rel) ~ factor(stage) + factor(histol) +
  I(age/12), design = d_BorganII)
```

	coef	exp(coef)	se(coef)	robust se	z	p
factor(stage)2	0.46286	1.58861	0.23762	0.18087	2.559	0.01049
factor(stage)3	0.58309	1.79156	0.23965	0.17848	3.267	0.00109
factor(stage)4	1.05967	2.88541	0.26182	0.20524	5.163	2.43e-07
factor(histol)2	1.59744	4.94035	0.17688	0.13342	11.973	< 2e-16
I(age/12)	0.02994	1.03039	0.02942	0.03337	0.897	0.36972

Likelihood ratio test= on 5 df, p=
n= 1062, number of events= 571

We can further post-stratify the design on disease stage and age with `calibrate`

```
> d_BorganIIps <- calibrate(d_BorganII, phase=2, formula=~age+interaction(instit,rel,stage))
> svycoxph(Surv(edrel,rel)~factor(stage)+factor(histol)+I(age/12),
+          design=d_BorganIIps)
```

Call:

```
svycoxph(formula = Surv(edrel, rel) ~ factor(stage) + factor(histol) +
I(age/12), design = d_BorganIIps)
```

	coef	exp(coef)	se(coef)	robust se	z	p
factor(stage)2	0.67006	1.95436	0.23776	0.14263	4.698	2.63e-06
factor(stage)3	0.75935	2.13689	0.23952	0.14228	5.337	9.45e-08
factor(stage)4	1.27046	3.56249	0.26150	0.15176	8.371	< 2e-16
factor(histol)2	1.57302	4.82121	0.17627	0.12999	12.101	< 2e-16
I(age/12)	0.03135	1.03185	0.02984	0.03346	0.937	0.349

```
Likelihood ratio test= on 5 df, p=
n= 1062, number of events= 571
```

References

- Barlow WE (1994). Robust variance estimation for the case-cohort design. *Biometrics* 50: 1064-1072
- Borgan Ø, Langholz B, Samuelson SO, Goldstein L and Pogoda J (2000). Exposure stratified case-cohort designs, *Lifetime Data Analysis* 6:39-58
- Breslow NW and Chatterjee N. (1999) Design and analysis of two-phase studies with binary outcome applied to Wilms tumour prognosis. *Applied Statistics* 48:457-68.
- Lin DY, and Ying Z (1993). Cox regression with incomplete covariate measurements. *Journal of the American Statistical Association* 88: 1341-1349.
- Therneau TM and Li H., Computing the Cox model for case-cohort designs. *Lifetime Data Analysis* 5:99-112, 1999
- Wacholder S, Gail MH, Pee D, and Brookmeyer R (1989) Alternate variance and efficiency calculations for the case-cohort design *Biometrika*, 76, 117-123