Data simulation experiments

BM

23/10/2020

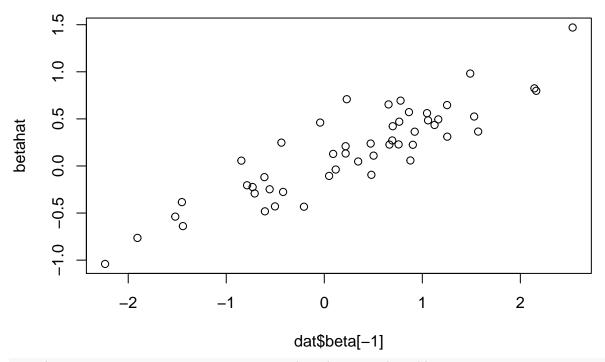
Case control simulation 1

- 1. Sample covariates on a population,
- 2. sample covariate effects from log-F prior,
- 3. calculate linear predictor of the logistic regression using all covariate effects,
- 4. simulate disease status, and
- 5. sample cases and controls from population.

```
ncase <- ncon <- 100
K <- 50 # num SNPs
m <- 4 # log-F parameter
N <- 20000 # popn size
beta0 <- (-4) # intercept
ccsim1 <- function(MAF) {</pre>
  X <- cbind(1,matrix(rbinom(N*K,size=2,p=MAF),ncol=K))</pre>
  beta <- c(beta0,log(rf(K,m,m)))</pre>
  linpred <- X %*% beta</pre>
  p <- exp(linpred)/(1+exp(linpred))</pre>
  cc <- rbinom(N,size=1,prob=p)</pre>
  caseind <- sample((1:N)[cc==1],size=ncase,replace=FALSE)</pre>
  conind <- sample((1:N)[cc==0],size=ncon,replace=FALSE)</pre>
  list(cc=c(rep(1,ncase),rep(0,ncase)),X = X[c(caseind,conind),],
        beta = beta)
}
```

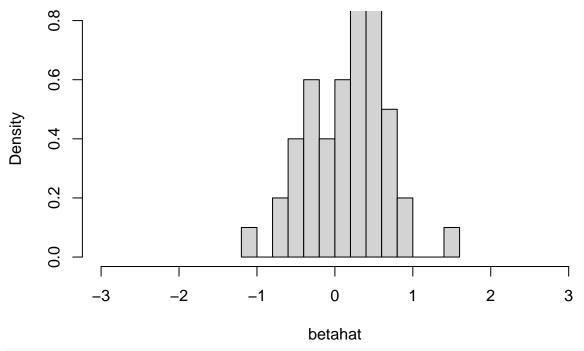
• Start with simulations for MAF = 0.5

```
set.seed(123)
dat <- ccsim1(MAF=.5)
betahat <- rep(NA,K)
for(i in 1:K) {
    XX <- dat$X[,1+i]
    betahat[i] <- coef(glm(dat$cc~XX,family=binomial()))[[2]]
}
plot(dat$beta[-1],betahat) # correlated, but betahats attenuated</pre>
```



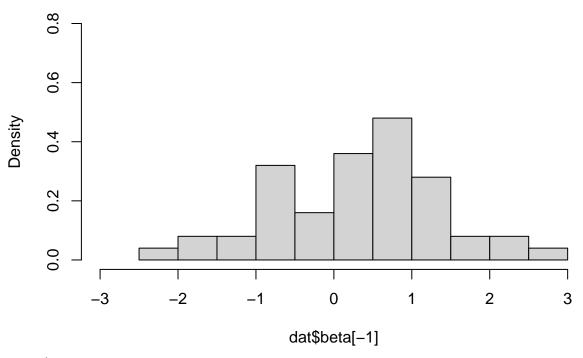
hist(betahat,freq=FALSE,nclass=10,ylim=c(0,.8),xlim=c(-3,3))

Histogram of betahat



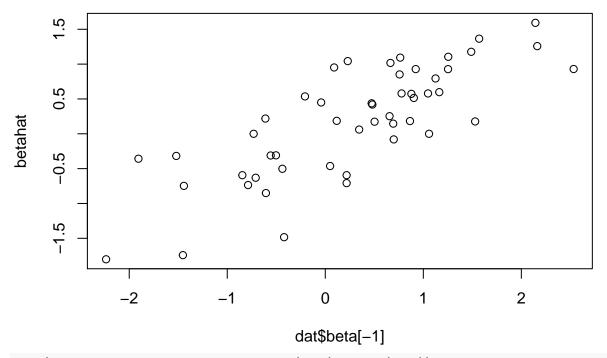
hist(dat\$beta[-1],freq=FALSE,nclass=10,ylim=c(0,.8),xlim=c(-3,3))

Histogram of dat\$beta[-1]



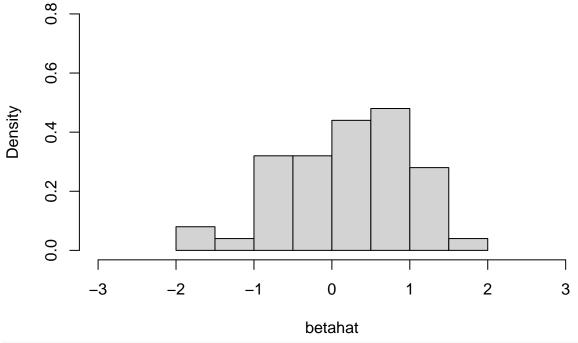
- $\hat{\beta}$'s are correlated with the β 's but are attenuated. I guess this is why DY and SC's methods are suggesting that m is very large (variance of log-F small).
- Next try MAF = 0.5

```
set.seed(123)
dat <- ccsim1(MAF=.05)
betahat <- rep(NA,K)
for(i in 1:K) {
    XX <- dat$X[,1+i]
    betahat[i] <- coef(glm(dat$cc~XX,family=binomial()))[[2]]
}
plot(dat$beta[-1],betahat) # betahats still attenuated, but less so</pre>
```



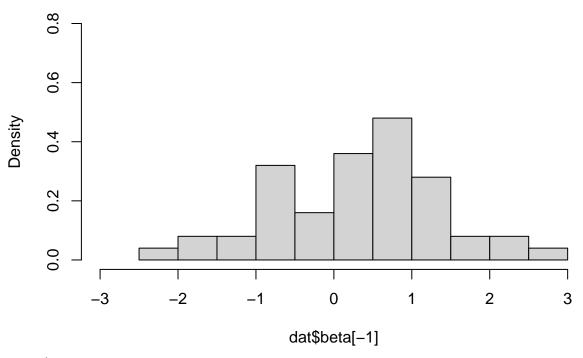
hist(betahat,freq=FALSE,nclass=10,ylim=c(0,.8),xlim=c(-3,3))

Histogram of betahat



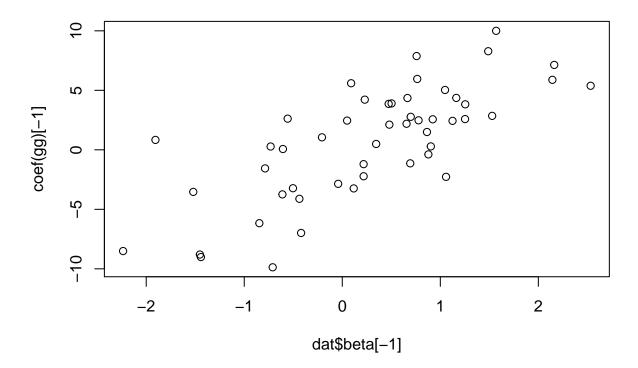
hist(dat\$beta[-1],freq=FALSE,nclass=10,ylim=c(0,.8),xlim=c(-3,3))

Histogram of dat\$beta[-1]



- $\hat{\beta}$'s are less attenuated, but still too small in magnitude.
- Consistent with the fact that DY and SC's likelihoods for m sometimes have a max.
- Can we fit a joint model in X's, like the one used to simulate the data? No, it fails with current sample size.

```
XX <- dat$X[,-1]</pre>
colnames(XX)<- paste0("X",1:K)</pre>
dd <- data.frame(cc=dat$cc,XX)</pre>
gg <- glm(cc ~ ., data=dd, family=binomial())</pre>
head(round(cbind(dat$beta,coef(gg),c(1,betahat)),3)) # too unstable
                          [,2]
##
                  [,1]
                                  [,3]
## (Intercept) -4.000 -7.541
                                1.000
## X1
                 2.143 5.882
                                1.592
## X2
                 1.163
                         4.362
                                0.598
## X3
                -0.557
                         2.620 -0.311
## X4
                 1.529 2.853 0.177
                 0.693 -1.136
                               0.146
plot(dat$beta[-1],coef(gg)[-1])
```



Case control simulation 2

- Use the Qin and Zhang method.
- If g(x) is the covariate distribution in controls, the distribution in cases is proportional to $g(x) \exp(x\beta)$.
- If g(x) is binomial(2, p), where p is the MAF (under a rare disease, the distribution in controls is about that in the population), then

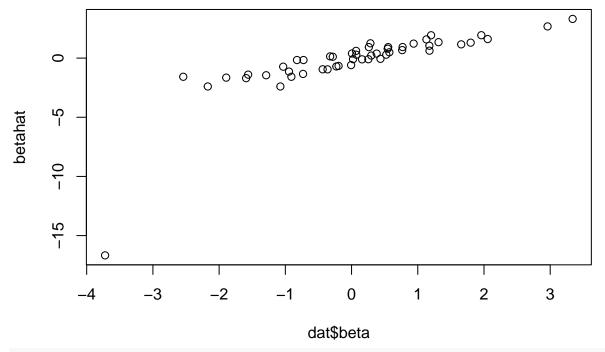
$$g(x)\exp(x\beta) = \begin{cases} (1-p)^2 & x = 0\\ 2p(1-p)\exp(\beta) & x = 1\\ p^2\exp(2\beta) & x = 2 \end{cases},$$

which has normalizing constant $(1-p)^2 + 2p(1-p)\exp(\beta) + p^2\exp(2\beta)$.

```
rcase <- function(n,beta,MAF) {</pre>
  p0 <- (1-MAF)^2
  p1 <- 2*MAF*(1-MAF)*exp(beta)
  p2 <- MAF^2*exp(2*beta)
  pp <- c(p0,p1,p2)
  sample(0:2,size=n,replace=TRUE,prob=pp)
}
ccsim2 <- function(MAF) {</pre>
  Xcon <- matrix(rbinom(ncon*K,size=2,prob=MAF),ncol=K)</pre>
  beta <- log(rf(K,m,m))</pre>
  Xcase <- matrix(NA,ncol=K,nrow=ncase)</pre>
  for(i in 1:K) {
    Xcase[,i] <- rcase(ncase,beta[i],MAF)</pre>
  X <- rbind(Xcase, Xcon)</pre>
  list(cc=c(rep(1,ncase),rep(0,ncase)),X = X,beta = beta)
}
```

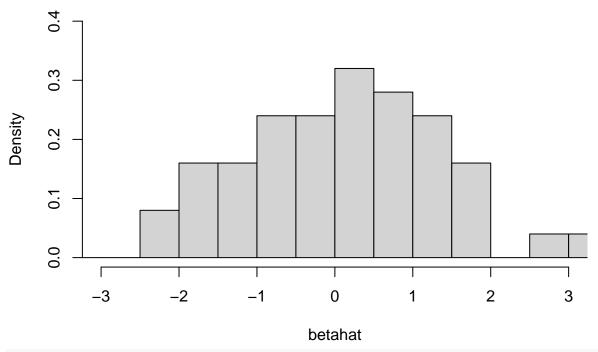
Simulate covariates with MAF = 0.05 and fit single-SNP models to see what the coefficients look like.
 We see one estimated coefficient of about -15, which I guess means non-convergence or monotone likelihood for that covariate.

```
set.seed(123)
dat <- ccsim2(MAF=.05)
betahat <- rep(NA,K)
for(i in 1:K) {
   XX <- dat$X[,i]
   betahat[i] <- coef(glm(dat$cc~XX,family=binomial()))[[2]]
}
plot(dat$beta,betahat)</pre>
```



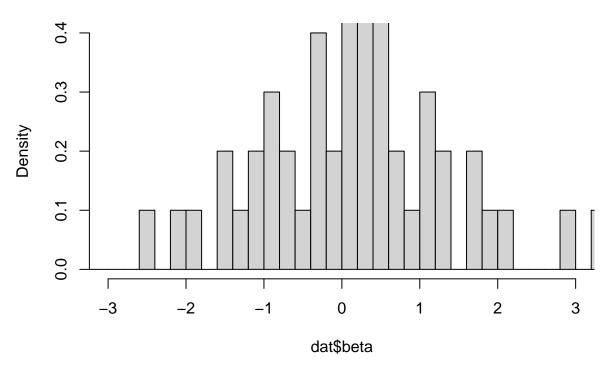
hist(betahat,freq=FALSE,nclass=30,ylim=c(0,.4),xlim=c(-3,3))

Histogram of betahat



hist(dat\$beta,freq=FALSE,nclass=30,ylim=c(0,.4),xlim=c(-3,3))

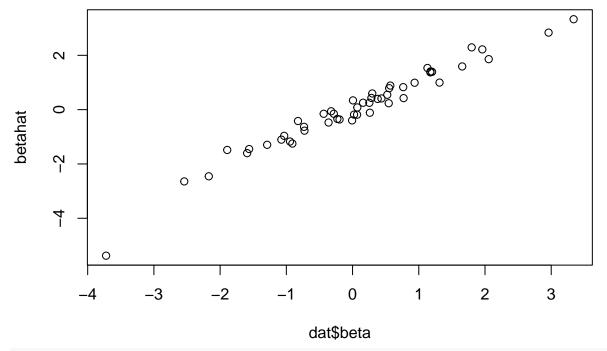
Histogram of dat\$beta



• Repeat with MAF = 0.5

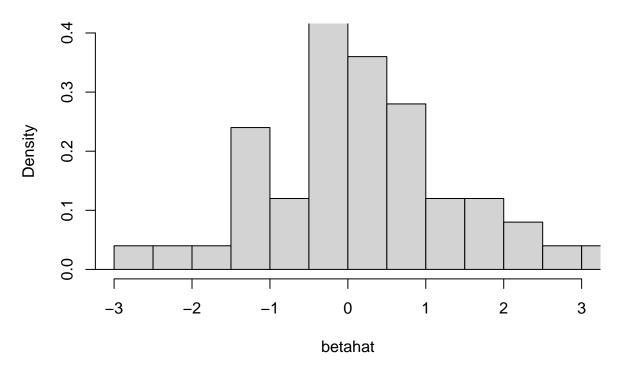
```
set.seed(123)
dat <- ccsim2(MAF=.5)</pre>
```

```
betahat <- rep(NA,K)
for(i in 1:K) {
    XX <- dat$X[,i]
    betahat[i] <- coef(glm(dat$cc~XX,family=binomial()))[[2]]
}
plot(dat$beta,betahat)</pre>
```



hist(betahat,freq=FALSE,nclass=30,ylim=c(0,.4),xlim=c(-3,3))

Histogram of betahat



Histogram of dat\$beta

