

Simulation to check the max purity

2019-03-05

We would like to firstly consider the scenario with two baseline covariates.

We have two treatment arms: placebo (pbo) and drug (drg). The outcomes of those two groups come from the formula:

$$\mathbf{y} = \mathbf{X}(\beta + \mathbf{b} + \mathbf{\Gamma}(\alpha' \mathbf{x})) + \epsilon.$$

We can define the covariate matrix of \mathbf{X} as \mathbf{z} . The \mathbf{z} contains both fixed effects and random effects.

$$\mathbf{z} = \beta + \mathbf{b} + \mathbf{\Gamma}x$$

Parameters:

Two groups

Set parameters:

- $\beta_{drg} = \begin{bmatrix} 0 \\ 25 \\ 1 \end{bmatrix}, \beta_{pbo} = \begin{bmatrix} 1 \\ -5 \\ -1 \end{bmatrix}$
- $\Gamma_{drg} = \begin{bmatrix} 0 \\ -2 \\ -1 \end{bmatrix}, \Gamma_{pbo} = \begin{bmatrix} 0 \\ 2 \\ 1 \end{bmatrix}$
- $\mathbf{b} \sim \begin{bmatrix} 1 & 0.1 & 0 \\ 1 & 0.3 & 0 \\ 2 & 0.2 & 0.1 \end{bmatrix} * N(3, 1)$
- $\epsilon_{drg} \sim N(3, 1); \epsilon_{pbo} \sim N(4, 1);$

Baselines

The baselines come from the same distributions

- Baseline covariate x_1, x_2 , iid $\sim N(0, 1)$
- A true coefficient vector $\alpha = (0, 1)$.
- A combination of baseline covariate w : $w = \alpha^T [x_1, x_2]$

Data generation

- n : number of subjects in each group (assume they have the same number of observations)
- x_1, x_2 : the baseline covariate x_1, x_2 , iid $\sim N(0, 1)$
- $w = \alpha_1 * x_1 + \alpha_2 * x_2$
- X : the independent covariates, $X = [1, t, t^2]$, where $t = 1, 2, \dots, 6$
- Y : the outcome is $\mathbf{y} = \mathbf{X}(\beta + \mathbf{b} + \mathbf{\Gamma}(\alpha' \mathbf{x})) + \epsilon$.

The code is

```
### data generation
true_generation = function(alpha){
  # alpha
  set.seed(123)
  alpha = as.matrix(alpha,p,1)
  dat_pbo = c()
  for(i in 1:n){
    pbo_temp = NULL
    pbo_temp$subj = rep(paste('pbo',i,sep=''),ni)
    pbo_temp$trt = rep('pbo',ni)
    baseline = as.matrix(rnorm(p),p,1)
    x1 = baseline[1]; x2 = baseline[2]
    w = rep(t(alpha) %*% baseline,ni)
    pbo_temp$x1 = rep(x1,ni); pbo_temp$x2 = rep(x2,ni)
    pbo_temp$w = w
    pbo_temp$tt = tt
    bi = randomeffcet %*%as.matrix(rnorm(3))
    yi = X%*%(beta_pbo+bi+gamma_pbo*w[1]) + sigma_pbo*rnorm(ni)
    pbo_temp$y = yi
    dat_pbo = rbind(dat_pbo, as.data.frame(pbo_temp))
  }

  dat_drg = c()
  for(i in 1:n){
    drg_temp = NULL
    drg_temp$subj = rep(paste('drg',i,sep=''),ni)
    drg_temp$trt = rep('drg',ni)
    baseline = as.matrix(rnorm(p),p,1)
    x1 = baseline[1]; x2 = baseline[2]
    w = rep(t(alpha) %*% baseline,ni)
    drg_temp$x1 = rep(x1,ni); drg_temp$x2 = rep(x2,ni)
    drg_temp$w = w
    drg_temp$tt = tt
    bi = randomeffcet %*%as.matrix(rnorm(3))
    yi = X%*%(beta_drg+bi+gamma_drg*w[1]) + sigma_pbo*rnorm(ni)
    drg_temp$y = yi
    dat_drg = rbind(dat_drg, as.data.frame(drg_temp))
  }
  return(list(dat_drg = dat_drg, dat_pbo = dat_pbo))
}
```

Purity calculation

$$p_w(x) = \frac{(f_1(x|w) - f_2(x|w))^2}{f_1(x|w) + f_2(x|w)}$$

where $f_1(x|w) \sim MVN(\beta_1 + \Gamma_1 * w, \mathbf{b}_1)$

- $f_2(x|w) \sim MVN(\beta_2 + \Gamma_2 * w, \mathbf{b}_2)$

1. Generate datasets based on the parameters and true α
2. Fit LME and estimate β , Γ and \mathbf{b}

3. Calculate the purity based on the above formula

With the true α , the purity should reach the max value.

Then test whether it is correct or not.

1. Choose another α candidate: α' and calculate another baseline covariates combination w'
2. Fit the LME with w' and estimate β' , Γ' and \mathbf{b}'
3. Calculate the purity based on the above formula

With α' , the purity should be smaller then the purity calculated by the true α

Results

The purity calculated by true α : 0.5866095

Other α candidates:

- c(1.1,0): 0.5866102
- c(1,0.5): 0.4853827
- c(1,1): 0.3888147
- c(0,1): 0.3270925
- c(1,10): 0.3310953
- c(-1,1): 0.333868

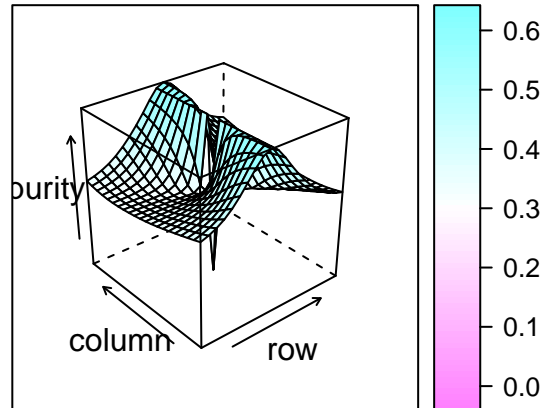
Find the max

I used the Newton Raphson method to find the max value. However, it still did not work well. We may try some other algorithm.

I just simply tried line search method, $\alpha = [\alpha_1, \alpha_2]$, vary α_1 for (-10,10,by =1); vary α_2 for (-10,10,by =1).

The purity looks like:

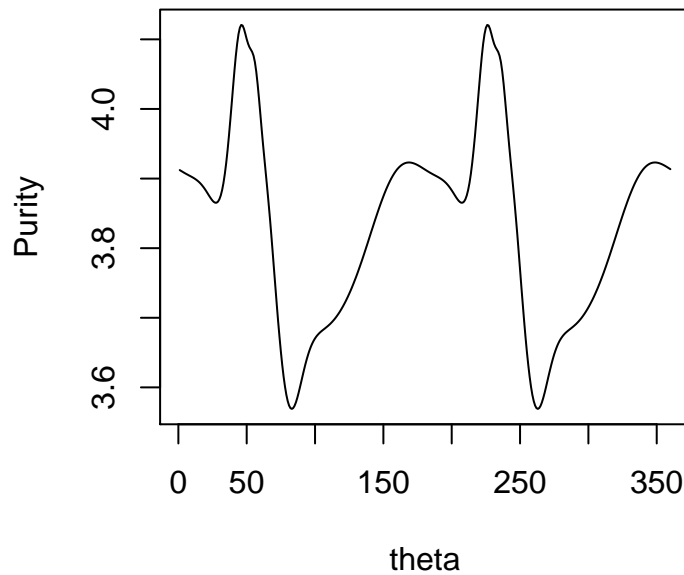
Table 1: 3D plot with different $\alpha = [\alpha_1, \alpha_2]$ values



Purity calculatin 2

- Make the $\alpha = [\sin(\theta), \cos(\theta)]$. The θ as the only input parameter.
- Set true purity as $\frac{3}{\pi}$

Table 2 purity v.s. theta, theta is from [0,360]



The max purity is:

```
data[data$purity == max(data$purity),]
```

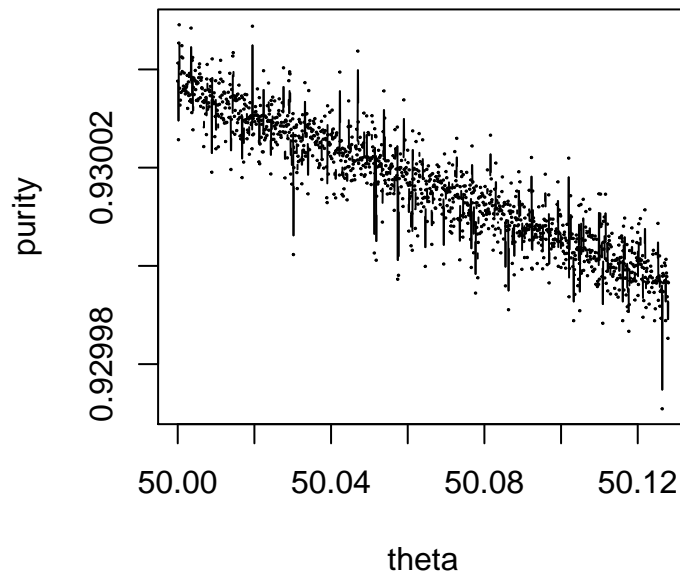
```
##      theta  purity
## 46      46 4.120739
## 226     226 4.120739
```

The plot looks smooth. However, it isn't. Since in the previous plot, the distance between two points is 1 degree. Let's make the distance smaller.

Check the points between 50, by 0.0001.

The plots:

Table 3 purity v.s. theta, theta is from [50,50]



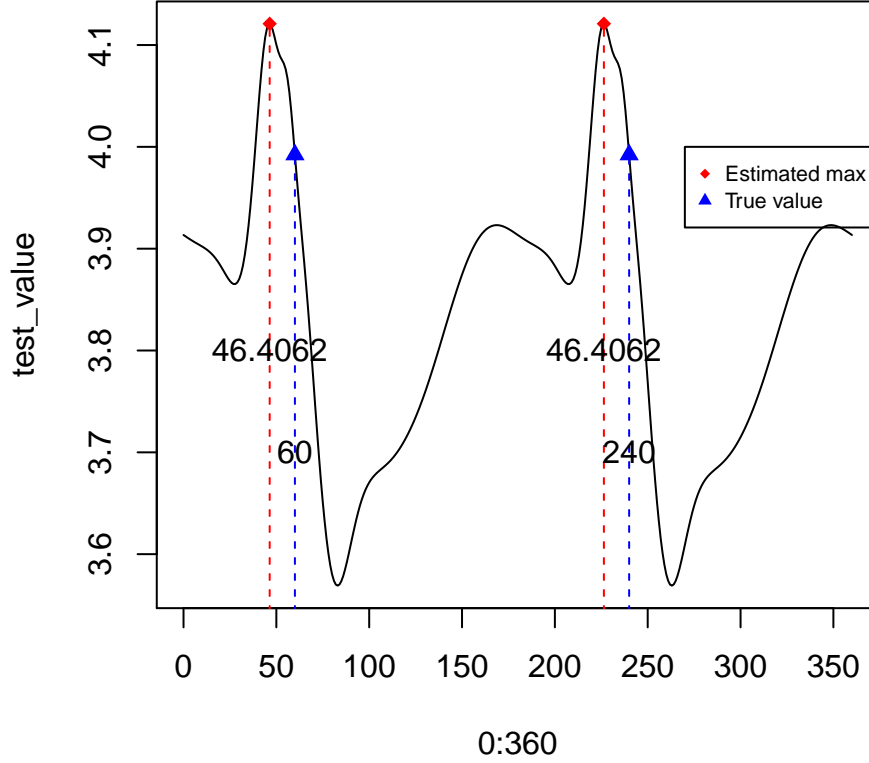
We can see that the plot has a trend, but actually very rough. It is hard to calculate the max value through Newton method.

Global optimization algorithms

There are some algorithms to find the global extreme values instead Newton Raphson method or gradient descent.

Two algorithms: genetic algorithm (GA) and simulated annealing (SA) are commonly used. The genetic algorithm seems to be the most accurate method of the two to find both the maximum and minimum of any function. Therefore, I tried to use GA to find the max value in our purity function. The results is showing below.

Table 4: purity vs theta with marked max value



The red points is the max value find by the genetic algorithm. We can see this algorithm works well. However, the max value did not match the true value.

Monte Carlo simulation to calculate the integral

step 1: The simulated values:

$$x = [x_1, x_2], \quad x_1 \sim UNI(-1, 1), \quad x_2 \sim UNI(-1, 1)$$

step 2: With the simulated x value, to calculate the purity function

$$p_w(x) = \frac{(f_1(x|w) - f_2(x|w))^2}{f_1(x|w) + f_2(x|w)}$$

step 3: Calculate the mean value of the purities

To make it consistent, the integral of $p_w(x)$ with different w values were calculated with the same input x , that is, the integrals were calculated in the same range.

```
Xstart = cbind(runif(1000,-1,1),runif(1000,-1,1))
monta_carlo_pdf = function(Xstart, mu1, D1, mu2, D2){

  mu1 = matrix(rep(mu1, 1000),1000,2,byrow = TRUE)
  mu2 = matrix(rep(mu2, 1000),1000,2,byrow = TRUE)

  Q1 = diag((-1/2)*(Xstart-mu1)%*%solve(D1)%*%t(Xstart-mu1))
  Q2 = diag((-1/2)*(Xstart-mu2)%*%solve(D2)%*%t(Xstart-mu2))
```

```

f1 = (1/(2*pi))*(1/sqrt(det(D1)))*exp(Q1)
f2 = (1/(2*pi))*(1/sqrt(det(D2)))*exp(Q2)

f3 = (f1 + f2)
f3 = ifelse(f3 == 0, 1e-4, f3)
purity = c((f1 - f2)^2 / f3) * 4
return(purity)
}

```

The other procedures are the same as the previous method.

purity vs theta after the integral is calculated by monte c

