

# $n$ -Sample Test Classifier on Binary Outcomes of Stratified Randomized Experiments

MA 590 Special Topics: Causal Inference

Aukkawut Ammartayakun

Worcester Polytechnic Institute

27 February, 2023

# Two-Sample Case

## Two-Sample Case

Let say we have two group of data  $D$ :  $D_t$  and  $D_c$ . We want to test whether there is a difference between the two groups.

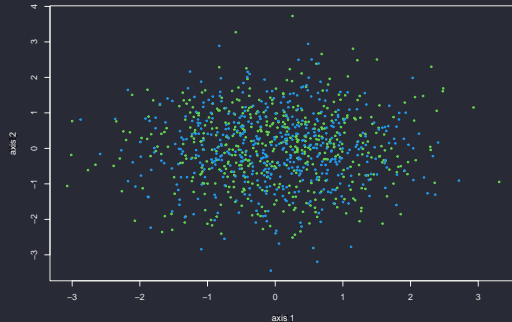
## Classifier Two-Sample Test (Lopez-Paz and Oquab 2017)

- Combined two dataset into one dataset  $D$ .
- Split the dataset into training and testing set.
- Fit the classifier (like logistic regression) to the training set and predict the testing set.
- Calculate the empirical loss  $l_e$  of the classifier. If  $|l_e - 0.5| < \epsilon$ , then  $\bar{\tau} = 0$ .

# Testing the Algorithm

```
set.seed(590)
# generate random multivariate gaussian data
n <- 1000
d <- 2
X <- matrix(rnorm(n*d), n, d)
y <- c(rep(0, n/2), rep(1, n/2))
c(c2st(X, y)$emp_loss, c2st(X,y)$pval)
```

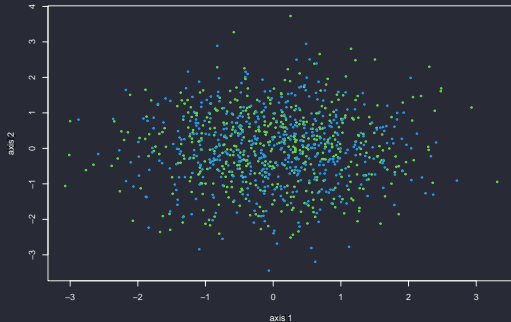
```
[1] 0.570 0.294
```



# Testing the Algorithm

```
set.seed(590)
# generate random multivariate gaussian data
n <- 1000
d <- 2
X <- matrix(rnorm(n*d), n, d)
y <- c(rep(0, n/2), rep(1, n/2))
c(c2st(X, y)$emp_loss, c2st(X,y)$pval)
```

```
[1] 0.570 0.294
```

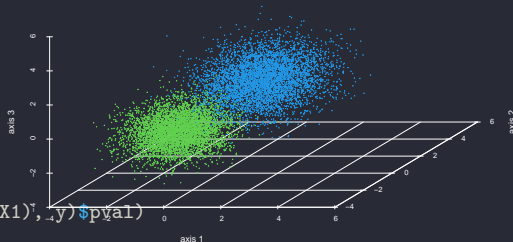


It is the *impossible* classification problem. Thus, the result should be close to near-chance level.

# Testing the Algorithm

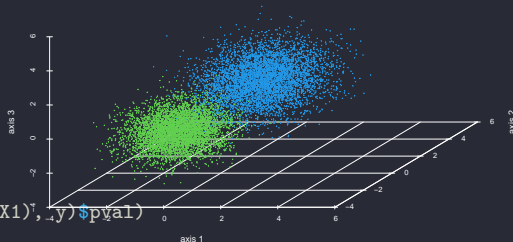
```
set.seed(590)
# generate two dataset: two gaussians
n <- 5000
d <- 3
X0 <- matrix(rnorm(n*d, -1,0.8), n, d)
X1 <- matrix(rnorm(n*d, 1,1), n, d)
y <- c(rep(0, n), rep(1, n))
c(c2st(rbind(X0, X1), y)$emp_loss, c2st(rbind(X0, X1), y)$pval)

[1] 0.0285 1.0000
```



# Testing the Algorithm

```
set.seed(590)
# generate two dataset: two gaussians
n <- 5000
d <- 3
X0 <- matrix(rnorm(n*d, -1, 0.8), n, d)
X1 <- matrix(rnorm(n*d, 1, 1), n, d)
y <- c(rep(0, n), rep(1, n))
c(c2st(rbind(X0, X1), y)$emp_loss, c2st(rbind(X0, X1), y)$pval)
[1] 0.0285 1.0000
```

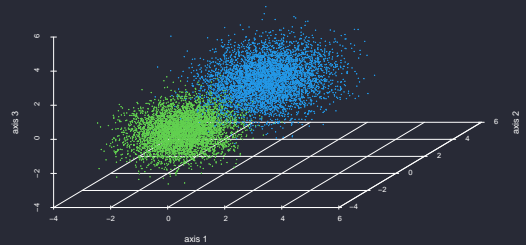


Test statistic diverges from near-chance level as there is an exist of “linear” decision boundary. The dimensionality problem in homogeneity test is solved.

**Question:** Does changing from GLM to other models increases (widen the range of) the testing power,  $\beta$ ?



# Testing the Algorithm



# Comparison to Maximum Mean Discrepancy

# Use C2ST on Causal Inference

- But, how can we use C2ST on causal inference?

## Example:

# Generalization of $n$ -Sample Case

## Possible Solution

- Combine both treatment and control group within each stratum into one dataset  $D_i^s$ .
- For each group, fit the classifier (like logistic regression) to the training set and predict the testing set.
- Calculate the empirical loss  $l_e$  of the classifier. If  $|l_e - 0.5| < \epsilon$ , then  $\bar{\tau}_{\text{within}} = 0$
- Find the way to infer  $\bar{\tau}_{\text{between}}$

# References

## References

Lopez-Paz, David, and Maxime Oquab. 2017. “Revisiting Classifier Two-Sample Tests.”  
In *International Conference on Learning Representations*.  
<https://openreview.net/forum?id=SJkXfE5xx>.