

DS 102 - HW 05

1. Simpson's Paradox

Adults		R	R ^c	R+R ^c	Recovery Rate
Drug (D)	0	1	1	0	
No Drug (D ^c)	14	24	39	0.359	
	14	25	40		

Kids		R	R ^c	R+R ^c	Recovery Rate
Drug (D)	20	19	39	0.513	
No Drug (D ^c)	2	0	1	2	
	22	19	40		

Simpson's paradox can occur when there is a large imbalance in the subgroups. By putting a large majority of kids as drug-users and adults as non-drug users, we can achieve tables s.t. equation 1 holds

b) If D and A are independent, prove D and K, D^c and A, D^c and K are inde

$$P(D) + P(D^c) = 1, P(A) + P(K) = 1 \text{ by definition}$$

$$P(D \cap A) = P(D) \cdot P(A); P(D) = P(D \cap A) + P(D \cap K)$$

$$P(D \cap K) = P(D) - P(D \cap A) = P(D) - P(D) \cdot P(A) = P(D)(1 - P(A)) = P(D)P(K)$$

So D is independent of A and K and by same argument, D^c is independent of A and K

c) Given D and A are independent, Simpson's Paradox is impossible

$$P(R|D, A) < P(R|D^c, A) \text{ and } P(R|D, K) < P(R|D^c, K) \text{ implies } P(R|D) < P(R|D^c)$$

$$\sum P(R|D, A_n)P(A_n|D) > \sum P(R|D^c, A_n)P(A_n|D^c)$$

$$\sum P(R|D, A_n)P(A_n) > \sum P(R|D^c, A_n)P(A_n)$$

$$\sum P(R|D, A_n) > \sum P(R|D^c, A_n)$$

Given $P(R|D, A) < P(R|D^c, A)$, the above statement cannot be true.

The same argument applies for Kids.

Therefore, if D and A, D and K, D^c and A, D^c and K are independent, then Simpson's paradox is not possible

$$(x_i - \bar{x})^2 = x_i^2 - 2x_i\bar{x} + \bar{x}^2$$



DS 102 - HW 5

11-4-19

2. Experiment Design for Linear Models

a) Show $\text{Var}(\hat{\alpha}_{OLS}) = \frac{\sigma^2}{n} \sigma^2$, $\text{Var}(\hat{\alpha}_{OLS}) = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$; $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

$$x \in 1, \dots, n \rightarrow x_i = 0 \text{ for } i = 1, \dots, \frac{n}{2}$$

$$x_i = 1 \text{ for } i = \frac{n}{2} + 1, \dots, n$$

$$\sum_{i=1}^n x_i = \sum_{i=1}^{n/2} x_i + \sum_{i=n/2+1}^n x_i = \frac{n}{2}$$

$$\bar{x} = \frac{1}{n} \left(\sum_{i=1}^{n/2} x_i + \sum_{i=n/2+1}^n x_i \right) = \frac{1}{n} \left(0 + \frac{n}{2} \right) = \frac{1}{2}$$

$$\sum_{i=1}^n (x_i - \frac{1}{2})^2 = \sum_{i=1}^n x_i^2 - x_i + \frac{1}{4}$$

$$\text{Var}(\hat{\alpha}_{OLS}) = \sigma^2 / \left(\sum_{i=1}^n x_i^2 - x_i + \frac{1}{4} \right) = \frac{\sigma^2}{\frac{n}{2} - \frac{n}{2} + \frac{n}{4}} = \frac{4}{n} \sigma^2$$

air) Show $\text{Var}(\hat{\alpha}_{quad}) = \frac{3}{2} \text{Var}(\hat{\alpha}_{OLS}) = \frac{3}{2} \cdot \frac{4}{n} \sigma^2 = \frac{6}{n} \sigma^2$

$$x_i = 0 \text{ for } i = 1, \dots, \frac{n}{3}$$

$$x_i = 0.5 \text{ for } i = \frac{n}{3} + 1, \dots, \frac{2n}{3}$$

$$x_i = 1 \text{ for } i = \frac{2n}{3} + 1, \dots, n$$

$$\bar{x} = \frac{1}{n} \left(\sum_{i=1}^{n/3} x_i + \sum_{i=n/3+1}^{2n/3} x_i + \sum_{i=2n/3+1}^n x_i \right) = \frac{n}{3} + \frac{n}{3} \cdot \frac{1}{2} + 0 = \frac{n}{3} + \frac{n}{6} = \frac{2n}{6} + \frac{n}{6} = \frac{3n}{6} = \frac{1}{2}$$

$$\sum_{i=1}^n x_i^2 = \sum_{i=1}^{n/3} 0^2 + \sum_{i=n/3+1}^{2n/3} \left(\frac{1}{2}\right)^2 + \sum_{i=2n/3+1}^n 1^2 = \frac{n}{3} + \frac{1}{4} \cdot \frac{n}{3} = \frac{n}{3} + \frac{n}{12} = \frac{4n}{12} + \frac{n}{12} = \frac{5n}{12}$$

$$\sum_{i=1}^n x_i = \frac{n}{3} + \frac{n}{3} \cdot \frac{1}{2} = \frac{n}{3} + \frac{n}{6} = \frac{2n}{6} + \frac{n}{6} = \frac{3n}{6} = \frac{n}{2}$$

$$\text{Var}(\hat{\alpha}_{quad}) = \sigma^2 / \left(\frac{5n}{12} - \frac{n}{2} + \frac{n}{4} \right) = \sigma^2 / \frac{n}{6} = \frac{6}{n} \sigma^2$$

$$\frac{5n}{12} - \frac{n}{2} + \frac{n}{4} = \frac{5n}{12} - \frac{6n}{12} + \frac{3n}{12} = \frac{2n}{12} = \frac{n}{6}$$

aiii) Show $\text{Var}(\hat{\alpha}_{even}) = \frac{3(n-1)}{n+1} \text{Var}(\hat{\alpha}_{OLS}) = \frac{3(n-1)}{n+1} \cdot \frac{4\sigma^2}{n} = \frac{12\sigma^2(n-1)}{n(n+1)}$

evenly spaced $\Rightarrow \sum_{i=1}^n i/n = x_i$

$$\bar{x} = \frac{1}{n} \left(\sum_{i=1}^n x_i \right) = \frac{1}{n} \left(\sum_{i=1}^n i/n \right) = \frac{n+1}{2} \cdot \frac{1}{n} = \frac{n+1}{2n}$$

$$\text{Var}(\hat{\alpha}_{even}) = \sigma^2 / \left(\sum_{i=1}^n x_i^2 - 2x_i\bar{x} + \bar{x}^2 \right) = \sigma^2 / \left(\sum_{i=1}^n x_i^2 - \frac{x_i(n+1)}{n} + \left(\frac{n+1}{2n} \right)^2 \right)$$

$$\sum_{i=1}^n x_i^2 = \sum_{i=1}^n \left(\frac{i}{n} \right)^2 = \frac{1}{n^2} \sum_{i=1}^n i^2 = \frac{(n+1)(2n+1)}{6n}$$

$$\text{Var}(\hat{\alpha}_{even}) = \sigma^2 / \left(\frac{(n+1)(2n+1)}{6n} - \frac{(n+1)}{2} \cdot \frac{(n+1)}{n} + \frac{(n+1)^2}{4n^2} \right) = \sigma^2 / \left(\frac{(n+1)(2n+1)}{6n} - \frac{(n+1)^2}{4n} + \frac{(n+1)^2}{4n^2} \right)$$

$$= \frac{\sigma^2}{\frac{(n-1)(n+1)}{12n}} = \frac{12n\sigma^2}{(n-1)(n+1)}$$

11-6-19

PS102 - HW 5 X_i

3. K arms, P_i for $i=1, \dots, K$ reward distribution
 mean $\mu_i = E_{P_i}[X_i]$, $P(a \leq X_i \leq b) = 1$ for all $a \leq b$

pulls each K arms, c timesAfter cK rounds $\rightarrow \hat{\mu}_i = \frac{1}{c} \sum_{s=1}^c X_{i,s}$
 for $t=1, 2, \dots$ A_t - denote choice of arm to pullmax. # explore pulls, c $t \leq cK$ $t > cK$

$$A_t = \begin{cases} (t \bmod K) + 1 \\ \arg \max_{i \in \{1, \dots, K\}} \hat{\mu}_i \end{cases}$$

mean of optimal arm: $\mu^* = \max_{i \in \{1, \dots, K\}} \mu_i$ Sub-optimality gap: $\Delta_i = \mu^* - \mu_i$ Pseudo-Regret $\Rightarrow R(n) = n\mu^* - E\left[\sum_{t=1}^n X_{A_t}\right]$ a) $T_i(t) \Rightarrow$ # times arm i , pulled by time t Show regret $\Rightarrow R(n) = \sum_{i=1}^K \Delta_i E(T_i(n))$

$$R(n) = n\mu^* - \sum_{t=1}^n E[X_{A_t}] = \sum_{t=1}^n [\mu^* - E[X_{A_t}]]$$

$$\sum_{t=1}^n E[X_{A_t}] = \sum_{i=1}^K \mu_i E(T_i(n))$$

 $\sum_{t=1}^n X_{A_t} \Rightarrow$ Sum of reward distributions for arm A , up to time step n

$$\mu_i = E[X_i]$$

 \hookrightarrow Avg. reward of Arm i $E[X_{A_t}] =$ Avg. reward of Arm at time step t

$$E[\mu_i T_i(n)] = \mu_i E(T_i(n))$$

$$R(n) = \sum_{i=1}^K (\mu^* - \mu_i) E(T_i(n)) = \sum_{i=1}^K \Delta_i E(T_i(n))$$

"Expectation of arm i at time step n "b) Show if $n > cK$, $E[T_i(n)] = c + (n - cK) \Pr(\hat{\mu}_i > \max_{j \neq i} \hat{\mu}_j)$

$$T_i(n) = \sum_{s=1}^n \mathbb{1}\{A_s = i\} \Rightarrow E\left[\sum_{s=1}^n \mathbb{1}\{A_s = i\}\right] = \sum_{s=1}^n E[\mathbb{1}\{A_s = i\}]$$

$$= \sum_{s=1}^{cK} E[\mathbb{1}\{A_s = i\}] + \sum_{s=cK+1}^n E[\mathbb{1}\{A_s = i\}]$$

$$= \sum_{s=1}^{cK} \Pr((s \bmod K) + 1 = i) + \sum_{s=cK+1}^n \Pr(\hat{\mu}_i > \max_{j \neq i} \hat{\mu}_j)$$

$$= c + (n - cK) \Pr(\hat{\mu}_i > \max_{j \neq i} \hat{\mu}_j)$$

c) Given optimal arm $i^* = 1$, show for any sub-optimal arm i :

$$\Pr(\hat{\mu}_i > \max_{j \neq i} \hat{\mu}_j) \leq \Pr(\hat{\mu}_i > \hat{\mu}_1) \quad \hat{\mu}_1 = \mu^*$$

Given optimal arm 1, and any sub-optimal arm i :

$$\max_{j=1, \dots, K, j \neq i} \hat{\mu}_j = \hat{\mu}_1 = \mu^* \Rightarrow \Pr(\hat{\mu}_i > \hat{\mu}_1) \leq \Pr(\hat{\mu}_i > \mu^*)$$

$$\Pr(\hat{\mu}_i > \hat{\mu}_1) > \Pr(\hat{\mu}_i > \hat{\mu}_j) \text{ for all } j=1, \dots, K$$

DS 102 - HW 5

11-9-19

3. d) Use Hoeffding to show $P(\hat{\mu}_1 > \hat{\mu}_1) \leq \exp\left\{-\frac{C \Delta_1^2}{(b-a)^2}\right\}$ $\Delta_1 = \mu^* - \mu_1$
 Given: $E[T_i(n)] \leq c + (n - K_c) P(\hat{\mu}_1 > \hat{\mu}_1)$. All X_i bounded on $[a, b]$

$$\begin{aligned} P(\hat{\mu}_1 > \hat{\mu}_1) &= P(\hat{\mu}_1 - E(\hat{\mu}_1) > \hat{\mu}_1 - E(\hat{\mu}_1)) \\ &= P(\hat{\mu}_1 - E[\frac{1}{c} \sum_{s=1}^c X_s] > \hat{\mu}_1 - \hat{\mu}_1) \\ &= P(\frac{1}{c} \sum_{s=1}^c X_s - E[\frac{1}{c} \sum_{s=1}^c X_s] > \Delta_1) = P(\frac{1}{c} \sum_{s=1}^c (X_s - E[X_s]) > \Delta_1) \\ &= P(\sum_{s=1}^c (X_s - E[X_s]) > c \Delta_1) \\ &\leq \exp\left\{-\frac{c^2 \Delta_1^2}{\sum_{s=1}^c (b-a)^2}\right\} = \exp\left\{-\frac{c^2 \Delta_1^2}{c(b-a)^2}\right\} = \exp\left\{-\frac{C \Delta_1^2}{(b-a)^2}\right\} \end{aligned}$$

$$\hat{\mu}_1 = \frac{1}{c} \sum_{s=1}^c X_s \Rightarrow \text{Constant}$$

$$\mu^* = \max_{i \in \{1, \dots, K\}} \mu_i = \hat{\mu}_1$$

e) $E[T_i(n)] \leq c + (n - K_c) \exp\left\{-\frac{C \Delta_i^2}{(b-a)^2}\right\}$

Min. Sub-optimality: $\Delta = \min_{i \geq 2} \Delta_i$. Then for each sub-optimal arm $i = 2, \dots, K$:

$$E[T_i(n)] \leq c + n \cdot \exp\left\{-\frac{C \Delta^2}{(b-a)^2}\right\} \text{ w/ } (n - K_c) \text{ upper bounded by } n$$

Find value of C st. $\exp\left(-\frac{C \Delta^2}{(b-a)^2}\right) \leq \frac{1}{n}$

$$-\frac{C \Delta^2}{(b-a)^2} \leq \ln \frac{1}{n} = -\ln(n)$$

$$-C \Delta^2 \leq -\ln(n) \cdot (b-a)^2; \quad C \Delta^2 \geq \ln(n) \cdot (b-a)^2$$

$$C \geq \frac{\ln(n) \cdot (b-a)^2}{\Delta^2}$$

$$\exp\left(-\frac{2 \ln(n) (b-a)^2 \Delta^2}{\Delta^2}\right) \leq \frac{1}{n}$$

$$\exp(-2 \ln(n)) \leq \frac{1}{n}$$

$$\exp\left(\ln\left(\frac{1}{n^2}\right)\right) \leq \frac{1}{n}$$

$$\frac{1}{n^2} \leq \frac{1}{n} \quad \checkmark$$

Pseudo-Regret: $R(n) = \sum_{i=1}^K \Delta_i E(T_i(n))$

$$\leq \sum_{i=1}^K \Delta_i \left(c + n \cdot \exp\left(-\frac{C \Delta_i^2}{(b-a)^2}\right)\right)$$

```
In [1]: import matplotlib.pyplot as plt
import numpy as np
import scipy.stats
import seaborn as sns
import math
```

Part B

```
In [2]: n = 100
even_spacing = [i/n for i in range(n)]
dumbbell = [0] * int(n/2) + [1] * int(n/2)
quad = [0] * int(n/3) + [1] * int(n/3) + [0.5] * int(n/3) + [0.5]

list_of_lists1 = [even_spacing, dumbbell, quad]
```

```
In [26]: a_const = 5
b_const = 8
# eps = np.random.normal(loc=0, scale=0.5, size=1)
```

```
In [27]: y1_even = []
y2_dumb = []
y3_quad = []
list_of_lists2 = [y1_even, y2_dumb, y3_quad]
for i in range(3):
    for x in list_of_lists1[i]:
        y = a_const * x + np.random.normal(loc=0, scale=0.5, size=1)
        list_of_lists2[i].append(y[0])
```

```
In [28]: a_ols = []

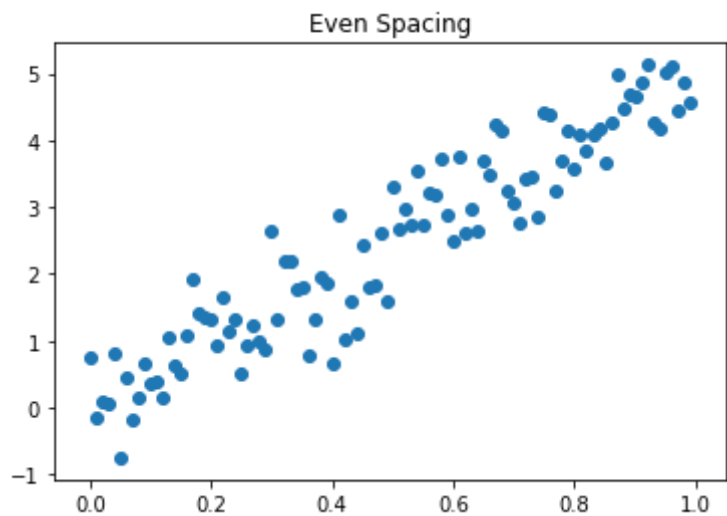
for i in range(3):
    xs = list_of_lists1[i]
    ys = list_of_lists2[i]
    x_bar = np.mean(xs)
    y_bar = np.mean(ys)
    alpha = 0
    for j in range(n):
        alpha += (xs[j] - x_bar) * (ys[j] - y_bar) / (xs[j] - x_bar)**2
#     a = np.sum([x - x_bar for x in xs])*np.sum([y - y_bar for y in y
s])/np.sum([(x - x_bar)**2 for x in xs])
    a_ols.append(alpha)
```

```
/usr/local/anaconda3/lib/python3.7/site-packages/ipykernel_launcher.py:
10: RuntimeWarning: invalid value encountered in double_scalars
# Remove the CWD from sys.path while we load stuff.
```

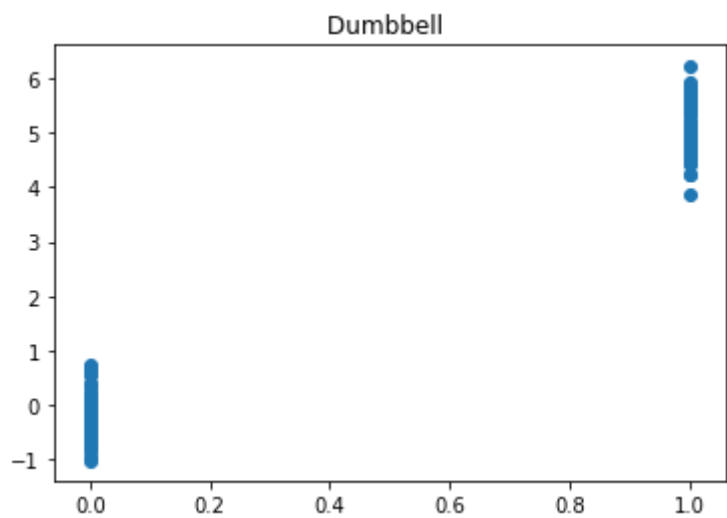
```
In [29]: a_ols
```

```
Out[29]: [955.5048739574764, 524.3193985660329, nan]
```

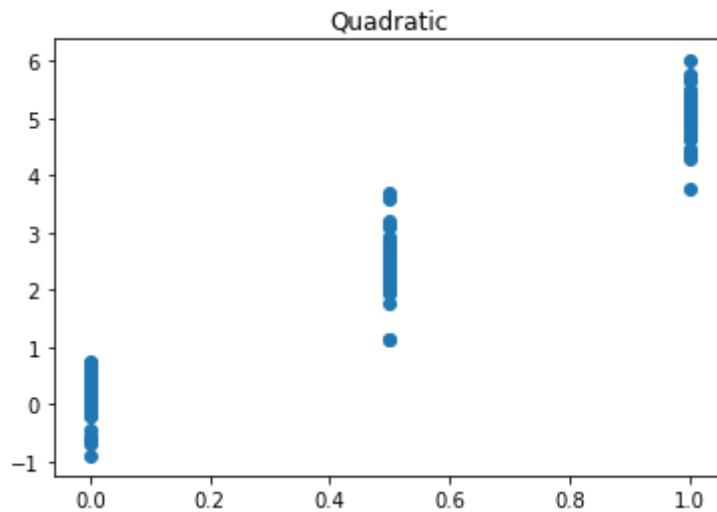
```
In [34]: plt.scatter(even_spacing, y1_even)
# plt.plot(even_spacing, y1_even)
plt.title('Even Spacing')
plt.show()
```



```
In [31]: plt.scatter(dumbbell, y2_dumb)
# plt.plot(even_spacing, y1_even)
plt.title('Dumbbell ')
plt.show()
```



```
In [32]: plt.scatter(quad, y3_quad)
plt.title('Quadratic')
plt.show()
```



```
In [35]: y1_even = []
y2_dumb = []
y3_quad = []
list_of_lists2 = [y1_even, y2_dumb, y3_quad]
for i in range(3):
    for x in list_of_lists1[i]:
        y = b_const * (x - 0.5)**2 + np.random.normal(loc=0, scale=0.5,
size=1)
        list_of_lists2[i].append(y[0])
```

```
In [36]: a_ols = []

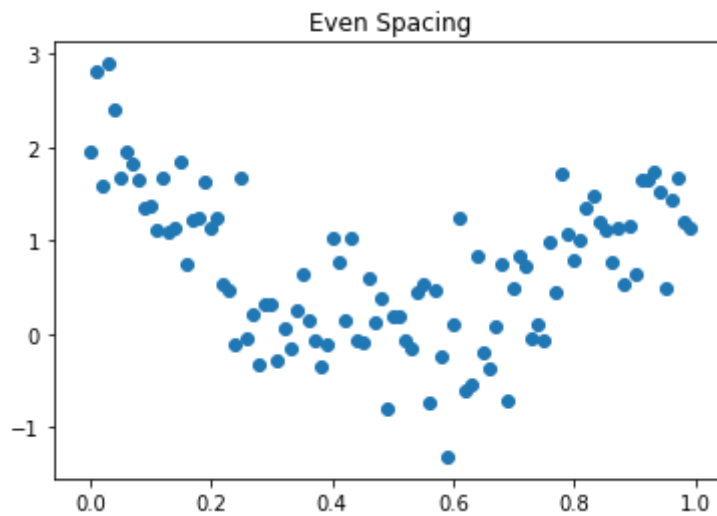
for i in range(3):
    xs = list_of_lists1[i]
    ys = list_of_lists2[i]
    x_bar = np.mean(xs)
    y_bar = np.mean(ys)
    alpha = 0
    for j in range(n):
        alpha += (xs[j] - x_bar) * (ys[j] - y_bar) / (xs[j] - x_bar)**2
    a_ols.append(alpha)
```

```
/usr/local/anaconda3/lib/python3.7/site-packages/ipykernel_launcher.py:
10: RuntimeWarning: invalid value encountered in double_scalars
# Remove the CWD from sys.path while we load stuff.
```

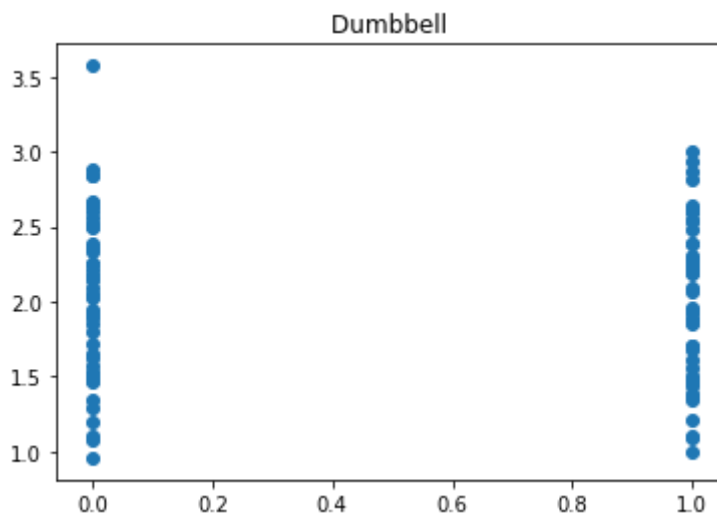
```
In [37]: a_ols
```

```
Out[37]: [105.8799360342669, -2.7380566902176735, nan]
```

```
In [38]: plt.scatter(even_spacing, y1_even)
# plt.plot(even_spacing, y1_even)
plt.title('Even Spacing')
plt.show()
```



```
In [39]: plt.scatter(dumbbell, y2_dumb)
# plt.plot(even_spacing, y1_even)
plt.title('Dumbbell ')
plt.show()
```




```
In [40]: plt.scatter(quad, y3_quad)
plt.title('Quadratic')
plt.show()
```

