

```
[51]: import matplotlib as plt
import seaborn as sns
import pandas as pd
import numpy as np
import json as js

In [52]: pwd

Out[52]: '/Users/mingyuanma/Desktop/HA1/analysis'

In [53]: before = pd.read_csv('.../../data/phase1/combined1.csv')
before = before.rename(columns = {"label":"label_dim1"})

In [54]: before_dim2 = pd.read_csv('.../../data/phase1/before_dim2.csv')
before_dim2 = before_dim2.rename(columns = {"Unnamed: 2":"label"}).loc[:,["ResponseId", "label"]]
before = before_dim2.merge(before, how="right", left_on="ResponseId", right_on="ResponseId")
before = before[before['round'] >= 3]
before

Out[54]:
```

	ResponseId	label	round	tick	orderid	taskid	taskTicks	workerid	workerTicks	label_dim1
135	R_0uOslE6BeLnJee5	3	3	1	2	1	2	0	0	3
136	R_0uOslE6BeLnJee5	3	3	1	1	1	2	1	0	3
137	R_0uOslE6BeLnJee5	3	3	1	0	0	0	0	0	3
138	R_0uOslE6BeLnJee5	3	3	2	0	0	0	0	0	3
139	R_0uOslE6BeLnJee5	3	3	2	0	0	0	0	0	3
...	...	...	...	...	...	...	...	...	...	...
100561	R_zTbYwRlBjcuhn	7	6	33	0	0	0	0	0	1
100562	R_zTbYwRlBjcuhn	7	6	33	0	0	0	0	0	1
100563	R_zTbYwRlBjcuhn	7	6	34	4	3	2	1	0	1
100564	R_zTbYwRlBjcuhn	7	6	34	0	0	0	0	0	1
100565	R_zTbYwRlBjcuhn	7	6	34	0	0	0	0	0	1

76941 rows x 10 columns

```
In [55]: after = pd.read_csv('.../../data/phase2/combined2.csv')
after = after.rename(columns = {"label":"label_dim1"})

In [56]: after_dim2 = pd.read_csv('.../../data/phase2/after_dim2.csv')
after_dim2 = after_dim2.rename(columns = {"Unnamed: 2":"label"}).loc[:,["ResponseId", "label"]]
after = after_dim2.merge(after, how="right", left_on="ResponseId", right_on="ResponseId")
after = after[after['round'] >= 3]
after

Out[56]:
```

	ResponseId	label	round	tick	orderid	taskid	taskTicks	workerid	workerTicks	label_dim1
117	R_DHBxkBVuS68D9GJ	4	3	1	1	1	2	0	0	2
118	R_DHBxkBVuS68D9GJ	4	3	1	2	1	2	1	0	2
119	R_DHBxkBVuS68D9GJ	4	3	1	0	0	0	0	0	2
120	R_DHBxkBVuS68D9GJ	4	3	2	3	1	2	1	0	2
121	R_DHBxkBVuS68D9GJ	4	3	2	0	0	0	0	0	2
...	...	...	...	...	...	...	...	...	...	...
141757	R_zZQbEOLF0I3yRX	5	6	37	0	0	0	0	0	3
141758	R_zZQbEOLF0I3yRX	5	6	37	0	0	0	0	0	3
141759	R_zZQbEOLF0I3yRX	5	6	38	4	3	2	1	0	3
141760	R_zZQbEOLF0I3yRX	5	6	38	0	0	0	0	0	3
141761	R_zZQbEOLF0I3yRX	5	6	38	0	0	0	0	0	3

108330 rows x 10 columns

## chi-square testing

```
In [57]: from scipy.stats import chi2_contingency
from collections import Counter

def chisquare(array1, array2, count=True):
    if not count:
        data = [array1, array2]
        stat, p, dof, expected = chi2_contingency(data)
    else:
        c1, c2 = Counter(array1), Counter(array2)
        before_dis, after_dis = [], []
        for i in set(c1).union(set(c2)):
            before_dis.append(c1[i])
            after_dis.append(c2[i])
        print(before_dis)
        data = [before_dis, after_dis]
        stat, p, dof, expected = chi2_contingency(data)

    alpha = 0.05
    print('p value is ' + str(p))
    if p <= alpha:
        print('difference between the two distributions (reject H0)')
    else:
        print('no difference between the two distributions (H0 holds true)')
```


## Distribution of label before and after

```
In [58]: d1 = before.groupby("ResponseId", as_index=False).agg(lambda x:x.lloc[0])
d2 = after.groupby("ResponseId", as_index=False).agg(lambda x:x.lloc[0])

In [59]: sns.histplot(data=d1, x="label", stat="probability", color="skyblue");
```



```
In [60]: sns.histplot(data=d2, x="label", stat="probability", color="skyblue");
```



## Chi-Square Testing

```
In [61]: chisquare(d1["label"], d2["label"])

[16, 0, 62, 26, 24, 14, 11, 4, 13]
[32, 2, 95, 17, 58, 6, 5, 1, 29]
p value is 0.000173093673541872
difference between the two distributions (reject H0)
```

## Analysis of Compliance

```
In [62]: def compliance(df):
    ids = []
    rounds = []
    server = []
    label = []
    for player in set(np.array(df["ResponseId"])):
        for i in np.arange(1,7):
            ten = df[(df["ResponseId"] == player) & (df["round"] == i)]
            if len(ten) != 0: # no response there
                ids.append(player)
                rounds.append(i)
                l = df[(df["ResponseId"] == player)][["label"]].lloc[0]
                label.append(l)
                if i <= 2:
                    num = sum((ten["workerId"] == 2) & (ten["taskId"] == 1))
                else:
                    num = sum((ten["workerId"] == 1) & (ten["taskId"] == 1))
                server.append(num)
            else:
                print(player, i)
    #
    d = {
        "ResponseId": ids,
        "round": rounds,
        "numServerCook": server,
        "label": label
    }
    return pd.DataFrame(d)
```

```
In [63]: before_compliance = compliance(before)
```

```
In [64]: after_compliance = compliance(after)
```

```
In [65]: before_compliance
```

	ResponseId	round	numServerCook	label
0	R_8nOxycIdgZm6kx	3	2	5
1	R_8nOxycIdgZm6kx	4	2	5
2	R_8nOxycIdgZm6kx	5	1	5
3	R_8nOxycIdgZm6kx	6	2	5
4	R_6JZ4qBJ3Xo9gt	3	1	5
...	...	...	...	...
675	R_BYHwHwPyAK6dP	6	3	4
676	R_1K1NAM6PzMKIU	3	2	3
677	R_1K1NAM6PzMKIU	4	2	3
678	R_1K1NAM6PzMKIU	5	2	3
679	R_1K1NAM6PzMKIU	6	2	3

680 rows x 4 columns

```
In [66]: after_compliance
```

	ResponseId	round	numServerCook	label
0	R_1MX1ThuDnBKGa8	3	2	5
1	R_1MX1ThuDnBKGa8	4	2	5
2	R_1MX1ThuDnBKGa8	5	2	5
3	R_1MX1ThuDnBKGa8	6	2	5
4	R_3jYXG134FVbha4S	3	2	3
...	...	...	...	...
975	R_ykHkGYUkpiMBh	6	4	3
976	R_3gMLEVUqubNFXW	3	10	4
977	R_3gMLEVUqubNFXW	4	3	4
978	R_3gMLEVUqubNFXW	5	1	4
979	R_3gMLEVUqubNFXW	6	3	4

980 rows x 4 columns

```
In [ ]:
```

## group by rounds

```
In [67]: round_before = before_compliance.groupby("round").agg(np.average).loc[:,["numServerCook"]]
round_before
```

	numServerCook
round	
3	2.552941
4	2.670588
5	2.617647
6	2.500000

```
In [68]: round_after = after_compliance.groupby("round").agg(np.average).loc[:,["numServerCook"]]
round_after
```

	numServerCook
round	
3	2.408163
4	2.232653
5	2.163265
6	2.044898

```
In [69]: a = np.array(round_before["numServerCook"])

In [70]: b = np.array(round_after["numServerCook"])
```

## chi-square testing

```
In [71]: chisquare(a, b, count=False)

p value is 0.9994561534639331
no difference between the two distributions (H0 holds true)
```

```
In [ ]:
```

## group by tip

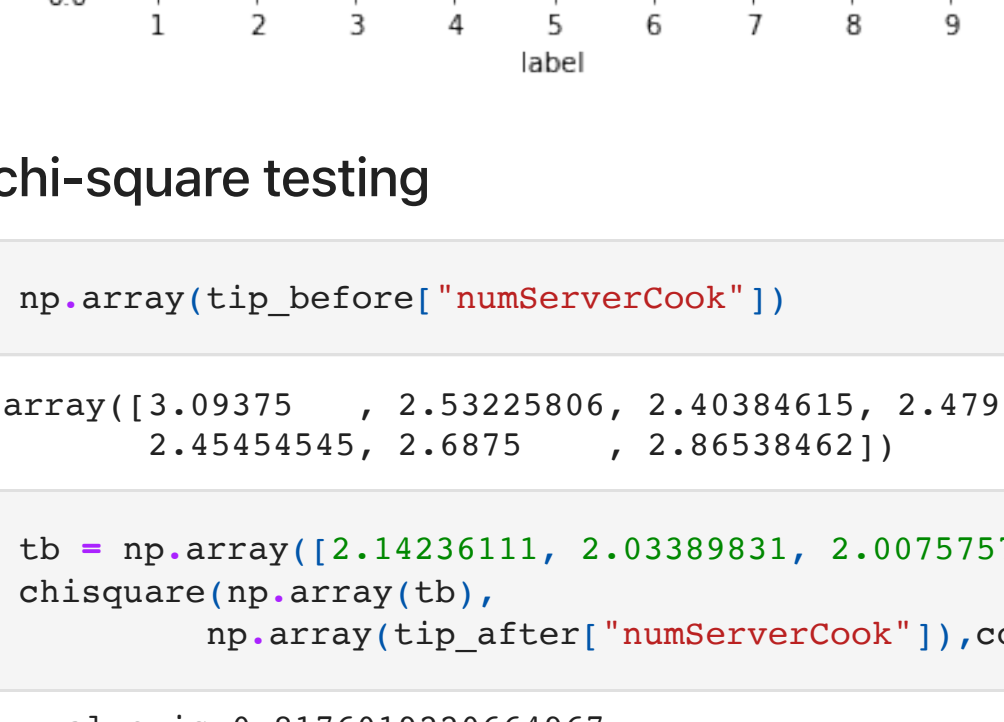
```
In [72]: tip_before = before_compliance.groupby("label", as_index=False).agg(np.average).loc[:,["label", "numServerCook"]]
tip_before
```

	label	numServerCook
0	1	3.093750
1	3	2.532258
2	4	2.403846
3	5	2.479167
4	6	2.571429
5	7	2.454545
6	8	2.687500
7	9	2.865385

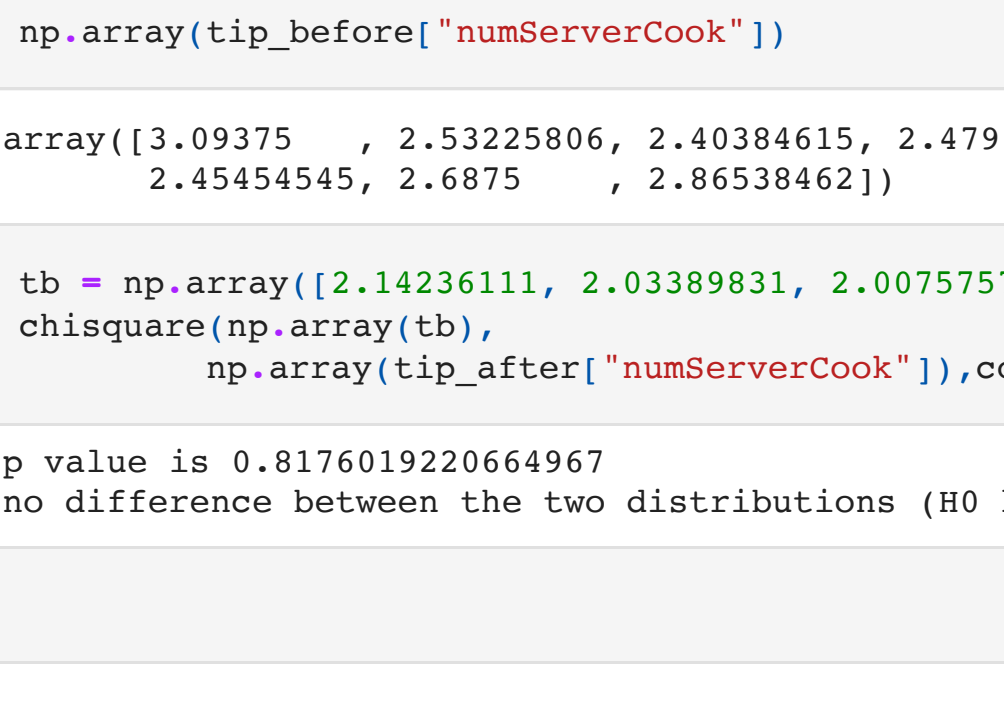
```
In [73]: tip_after = after_compliance.groupby("label", as_index=False).agg(np.average).loc[:,["label", "numServerCook"]]
tip_after
```

	label	numServerCook
0	1	1.914062
1	2	2.375000
2	3	2.215789
3	4	2.383265
4	5	2.219828
5	6	2.166667
6	7	1.850000
7	8	2.000000
8	9	2.517241

```
In [74]: sns.barplot(x="label", y="numServerCook", data=tip_before, alpha=0.8);
```



```
In [75]: sns.barplot(x="label", y="numServerCook", data=tip_after, alpha=0.8);
```



## Testing Aversion Distribution

```
In [76]: ave_before = before_compliance[before_compliance["numServerCook"] >= 3]
        (before_compliance["numServerCook"] == 0).groupby(
            "label").count().lloc[:,("round")]
ave_before = ave_before.rename(columns = {"round":"numAversion"})
```

```
In [79]: ave_after = after_compliance[after_compliance["numServerCook"] >= 3]
        (after_compliance["numServerCook"] == 0).groupby(
            "label").count().lloc[:,("round")]
ave_after = ave_after.rename(columns = {"round":"numAversion"})
```

```
In [80]: def counting(df):
    count = []
    for i in np.arange(1,10):
        cou = len(df[df["label"] == i])
        count.append(cou)
    return count
```

```
In [81]: counting(before_compliance), counting(after_compliance)
[_, for _ in counting(before_compliance) if _]

Out[81]: [64, 248, 104, 96, 56, 44, 16, 52]
```

```
In [82]: ave_before["total_num_label"] = [_, for _ in counting(before_compliance) if _]
ave_before["proportion_ersion"] = ave_before["numAversion"] / ave_before["total_num_label"]
ave_before["proportion_label"] = ave_before["total_num_label"] / before_compliance.shape[0]
ave_before
```

	numAversion	total_num_label	proportion_ersion	proportion_label
label				
1	35	64	0.546875	0.094118
3	96	248	0.387097	0.364706
4	37	104	0.355769	0.152941
5	43	96	0.447917	0.141176
6	23	56	0.410714	0.082353
7	20	44	0.454545	0.064706
8	7	16	0.437500	0.023529
9	25	52	0.480769	0.076471

```
In [84]: ave_after["total_num_label"] = [_, for _ in counting(after_compliance) if _]
ave_after["proportion_ersion"] = ave_after["numAversion"] / ave_after["total_num_label"]
ave_after["proportion_label"] = ave_after["total_num_label"] / after_compliance.shape[0]
ave_after
```

	numAversion	total_num_label	proportion_ersion	proportion_label
label				
1	26	64	0.406250	0.065306
2	3	248	0.012097	0.030061
3	108	104	1.038462	0.106122
4	21	96	0.218750	0.097959
5	57	56	1.017857	0.057143
6	5	44	0.113636	0.044898
7	2	16	0.125000	0.016327
8	4	52	0.076923	0.023529
9	51	52	0.980769	0.053061

```
In [ ]:
```

## Testing Compliance Distribution

```
In [85]: com_before = before_compliance[before_compliance["numServerCook"] ==2].groupby(
        "label").count().lloc[:,("round")]
com_before = com_before.rename(columns = {"round":"numCompliance"})
```

```
In [86]: com_after = after_compliance[after_compliance["numServerCook"] ==2].groupby(
        "label").count().lloc[:,("round")]
com_after = com_after.rename(columns = {"round":"numCompliance"})
com_after
```

	numCompliance
label	
1	61
2	4
3	193
4	39
5	130
6	18
7	12
8	4
9	51

```
In [87]: com_before["total_num_label"] = [_, for _ in counting(before_compliance) if _]
com_before["proportion_compliance"] = com_before["numCompliance"] / com_before["total_num_label"]
com_before["proportion_label"] = com_before["total_num_label"] / before_compliance.shape[0]
com_before
```

	numCompliance	total_num_label	proportion_compliance	proportion_label
label				
1	27	64	0.421875	0.094118
2	145	248	0.584677	0.364706
3	64	104	0.615385	0.152941
4	51	96	0.531250	0.141176
5	31	56	0.553571	0.082353
6	23	44	0.522727	0.064706
7	8	16	0.500000	0.023529
8	23	52	0.442308	0.076471

```
In [88]: com_after["total_num_label"] = counting(after_compliance)
com_after["proportion_compliance"] = com_after["numCompliance"] / com_after["total_num_label"]
com_after["proportion_label"] = com_after["total_num_label"] / after_compliance.shape[0]
com_after
```

	numCompliance	total_num_label	proportion_compliance	proportion_label
label				
1	61	128	0.476562	0.130612
2	4	8	0.500000	0.008163
3	193	380	0.507895	0.387755
4	39	68	0.573529	0.069388
5	130	232	0.560345	0.236735
6	18	24	0.750000	0.024490
7	12	20	0.600000	0.020408
8	4	4	1.000000	0.004082
9	51	116	0.439655	0.118367

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```