# MPI-accelerated Bootstrap Resampling for Hypothesis Testing in Small Sample Investigation
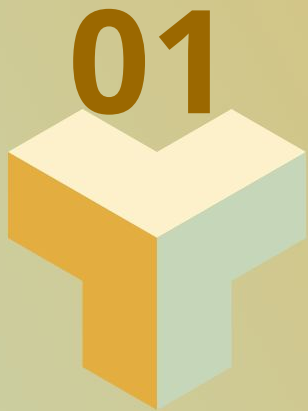
## Parallel Computing
## Final Report

M112040036 孫瑞鴻
M112040034 李祐瑄
2023-06-06

# CONTENTS

- **Data**
- **Bootstrap**
- **MPI Application**
- **Conclusion**
- **Extra try — Remote**

# 01

Data

# The data on throwing distances for dominant and non-dominant hand in pitching.

◆ 隨機抽樣七名學童，測量其以慣用手與非慣用手丟棒球，投出的距離 (公尺)如下：

| 學童 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|------|------|------|------|------|------|------|------|
| 慣用手 | 77.3 | 95.1 | 97.2 | 94.5 | 83.6 | 90.2 | 84.3 |
| 非慣用手 | 74.6 | 83.4 | 80.6 | 90.2 | 78.7 | 83.6 | 76.2 |

◆ 試找出學童兩手投球距離平均差異的 95% 信賴區間 在顯著水準 5% 之下，回答該差異是否明顯不同？

虛無假設:兩手投球距離無明顯差異

H0 = mean($X_1$) - mean($X_2$) = 0

02

**Bootstrap**

# parametric Bootstrap

- 觀測到的資料不知來自何分布
  - 計算經驗分布 $\hat{F}$，即每個資料點 $x_i$, $i=1,2,\ldots,n$ 的機率值都是 $1/n$
  - 將經驗分布 $\hat{F}$ 當作母體實際的分布
  - 從「母體」產生 $B$ 組樣本數為 $n$ 的 Bootstrap 樣本
  - 每個Bootstrap 樣本各自計算統計量 Bootstrap replications 的數值
  - 蒐集 $B$ 個Bootstrap replications，得到統計量的抽樣分布

# 03

## MPI Application

# MPI Application

```python
import numpy as np
from mpi4py import MPI
import time

# 初始時間
a = time.time()

# 定義Bootstrap
def bootstrap(data, n_bootstrap):
    n = len(data)
    bootstrap_samples = np.zeros((n_bootstrap, n))
    for i in range(n_bootstrap):
        bootstrap_sample = np.random.choice(data, size=n, replace=True)
        bootstrap_samples[i] = bootstrap_sample
    return bootstrap_samples

# 初始化MPI環境
comm = MPI.COMM_WORLD
rank = comm.Get_rank()
size = comm.Get_size()
```

# MPI Application

```python
# 原始數據
strong_hand = np.array([77.3, 95.1, 97.2, 94.5, 83.6, 90.2, 84.3])
weak_hand = np.array([74.6, 83.4, 80.6, 90.2, 78.7, 83.6, 76.2])

#計算原始的平均數差異
u_diff = np.mean(strong_hand) - np.mean(weak_hand)

# 將10個bootstrap樣本平均分配給每個CPU
n_bootstrap = 1000000
n_bootstrap_per_cpu = n_bootstrap // size

# 每個CPU執行自己分配到的bootstrap樣本
bootstrap_samples_local1 = bootstrap(strong_hand, n_bootstrap_per_cpu)
bootstrap_samples_local2 = bootstrap(weak_hand, n_bootstrap_per_cpu)

# 將每個CPU的結果收集到根節點
bootstrap_samples_all1= comm.gather(bootstrap_samples_local1, root=0)
bootstrap_samples_all2 = comm.gather(bootstrap_samples_local2, root=0)
```

```python
41  if rank == 0:
42      '''
43      for i, samples in enumerate(bootstrap_samples_all):
44          for j, sample in enumerate(samples):
45              print(f"Bootstrap樣本{size*i + j + 1}：", sample)
46      '''
47      # 計算每個bootstrap樣本的平均數
48
49      bootstrap_means1 = np.mean(bootstrap_samples_all1, axis=(0, 2))
50      bootstrap_means2 = np.mean(bootstrap_samples_all2, axis=(0, 2))
51      mean_diff = bootstrap_means1 - bootstrap_means2
52
53      #計算std
54      std = np.std(mean_diff) / np.sqrt(n_bootstrap)
55
56      #計算信賴區間的上下界
57      lower = u_diff - (1.96 * std)
58      upper = u_diff + (1.96 * std)
```
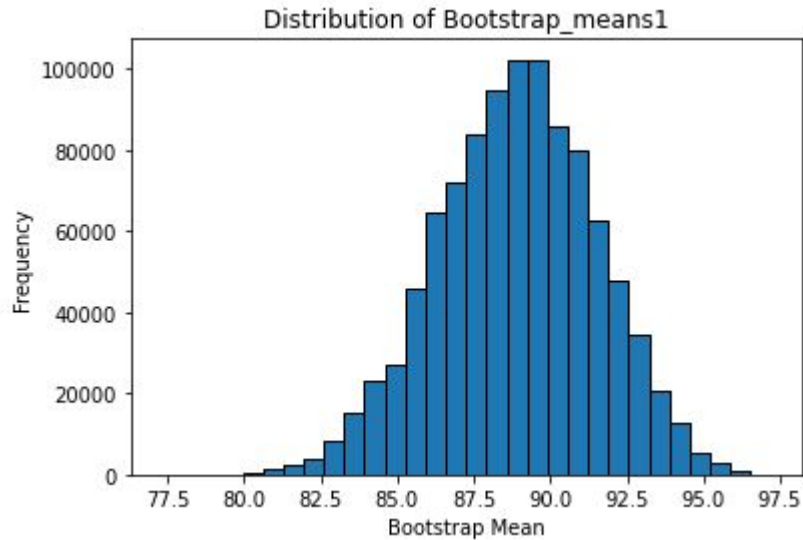
```
59
60        #利用信賴區間判斷平均數是否有顯著差異
61        if 0 < lower or 0 > upper:
62            decision = "reject the null hypothesis"
63        else:
64            decision = "don't reject the null hypothesis"
65
66    # 結束時間
67    b = time.time()
68    print(f"time {b-a:.3f}")
```
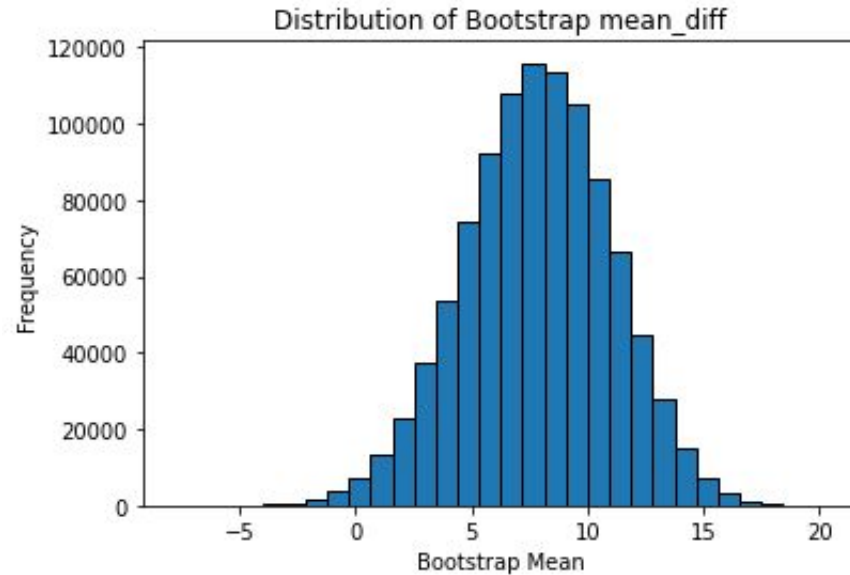
# 04

Conclusion

# Conclusion



Distribution of Bootstrap_means1

Distribution of Bootstrap_means2

# Conclusion



Distribution of Bootstrap mean_diff

**According to the Central Limit Theorem, it can be observed that the distribution approaches a normal distribution.**

# Conclusion

```
difference of observed mean = 7.843
95% confidence interval : [7.837,7.849]
decision: reject the null hypothesis
```

**The confidence interval does not include zero, therefore the null hypothesis is rejected.**

# Conclusion

```python
import os

use_one_pros = []
for i in range(5):
    cmd = os.popen(r'mpiexec -n 1 python D:\code\workplace_py\bootstrap.py "%s"' % str("5003"))
    output = cmd.read().strip().splitlines()  # 將輸出按行分割為list
    last_line = output[-1]  # 獲取最後一行
    time = float(last_line)
    use_one_pros.append(time)
avg_time_1 = sum(use_one_pros)/len(use_one_pros)
print("The time taken by 1 process" , round(avg_time_1, 4), "s")

use_four_pros = []
for i in range(5):
    cmd = os.popen(r'mpiexec -n 4 python D:\code\workplace_py\bootstrap.py "%s"' % str("5003"))
    output = cmd.read().strip().splitlines()  # 將輸出按行分割為list
    last_line = output[-1]  # 獲取最後一行
    time = float(last_line)
    use_four_pros.append(time)
avg_time_4 = sum(use_four_pros)/len(use_four_pros)
print("The time taken by 4 process" , round(avg_time_4, 4), "s")

speedup = avg_time_1 / avg_time_4
print("speedup:", round(speedup,2))
```

# Conclusion

**Speed up : The time taken by 1 process / The time taken by 4 process**

**The time taken by 1 process 25.8166 s**
**The time taken by 4 process 9.9266 s**
**speedup: 2.6**

```
The time taken by 1 process 25.8166 s
The time taken by 4 process 9.9266 s
speedup: 2.6
```

**05**

Extra try — Remote

# Extra try — Remote

# Extra try — Remote

# Extra try — Remote

```python
import ipyparallel as ipp
c = ipp.Client(profile='remote', cluster_id='remote')
dview = c[:]
print(c.ids)
```

```
[0, 1, 2, 3, 4, 5]
```

```python
dview.run('D://code/workplace_py/Parallel/PythonCode/rank_mpi.py')
%px totalhost = hostname()
dview['totalhost']
```

```
['LENOVO', 'LENOVO', 'LENOVO', 'LENOVO', 'LAPTOP-NBMNPREL', 'LAPTOP-NBMNPREL']
```

```
%%px

a = time.time()
# 假設有一組原始數據
strong_hand = np.array([77.3, 95.1, 97.2, 94.5, 83.6, 90.2, 84.3])
weak_hand = np.array([74.6, 83.4, 80.6, 90.2, 78.7, 83.6, 76.2])

#計算原始的平均數差異
u_diff = np.mean(strong_hand) - np.mean(weak_hand)

# 將10個bootstrap樣本平均分配給每個CPU
n_bootstrap = int(1000000/2)
n_bootstrap_per_cpu = n_bootstrap // size

# 每個CPU執行自己分配到的bootstrap樣本
bootstrap_samples_local1 = bootstrap(strong_hand, n_bootstrap_per_cpu)
bootstrap_samples_local2 = bootstrap(weak_hand, n_bootstrap_per_cpu)

# 將每個CPU的結果收集到根節點
bootstrap_samples_all1= comm.gather(bootstrap_samples_local1, root=0)
bootstrap_samples_all2 = comm.gather(bootstrap_samples_local2, root=0)

if rank == 0:

    # 計算每個bootstrap樣本的平均數
    bootstrap_means1 = np.mean(bootstrap_samples_all1, axis=(0, 2))
    bootstrap_means2 = np.mean(bootstrap_samples_all2, axis=(0, 2))
    mean_diff = bootstrap_means1 - bootstrap_means2

# 結束時間
b = time.time()
print(f"time {b-a:.3f}")
```

# Extra try — Remote

```
[stdout:0] time 4.666
[stdout:1] time 4.634
[stdout:2] time 4.636
[stdout:3] time 4.638
[stdout:4] time 15.813
[stdout:5] time 15.868
```

THE END

**THANKS**