# Pattern Classification and Recognition
## ECE 681

### Spring 2019
### Homework #5: Nonlinear Dimensionality Reduction (t-SNE)
### (In-Class Assignment)

Due: 5:00 PM, Thursday, March 7, 2019
Grace Period Concludes: 11:30 PM, Tuesday, March 19, 2019

This homework assignment is worth **110 points**.
Each problem is worth some multiple of 10 points, and will be scored on the below letter scale.
The letter grades B through D may be modified by + (+3%) and A through D may be modified by a - (-3%).

    A+ = 100%: Exceeds expectations, and no issues identified
    A = 95%: Meets expectations, and (perhaps) minor/subtle issues
    B = 85%: Issues that need to be addressed
    C = 75%: Significant issues that must be addressed
    D = 65%: Major issues, but with noticeable perceived effort
    F = 50%: Major issues, and insufficient perceived effort
    Z = 30%: Minimal perceived effort
    N = 0%: Missing, or no (or virtually no) perceived effort

Your homework is not considered submitted until both components (**one self-contained pdf file** and your code) have been submitted. Please do not include a print-out of your code in the pdf file.

You should strive to submit this assignment by the due date/time. The grace period is intended to afford an opportunity to, for example, work through technical glitches, update your submission if you realize after the due date that you submitted the wrong file or would prefer to answer a question differently, and provide some flexibility to manage your workload as it ebbs and flows during the semester.

## Visualizing Flip Sequences from Pennies and Quarters

We are re-visiting the (impossible) machine learning problem we considered when exploring the curse of dimensionality: Distinguishing pennies from quarters based solely on the sequence of Heads/Tails flips they produce. The coin is in another room, where someone else will flip it some number of times and tell us the result. That is, the observations we get in order to classify the coin as either a penny or a quarter are a series of coin flip results, heads or tails.[1]

The data for this problem is the sequence of Heads/Tails flips from each coin. If $P$ pennies and $Q$ quarters ($N = P + Q$) were each flipped $D$ times the resulting dataset might look like this:

$$
X = \begin{bmatrix} T & H & \dots & T \\ H & T & \dots & T \\ \vdots & \vdots & \vdots & \vdots \\ H & H & \dots & H \end{bmatrix}
\qquad
Y = \begin{bmatrix} q \\ q \\ \vdots \\ p \end{bmatrix}
$$

where $X$ contains the results of the coin flips and is of size $N \times D$ (coins $\times$ flips), and $Y$ contains the truth as to which coin was flipped and is of size $N \times 1$ (coins $\times$ 1). The "data" ($X$) are the coin flip results: heads (H) or tails (T). The "truth" or "labels" ($Y$) are either quarter (q) or penny (p).

We can encode these data and associated labels by assigning 0 to represent tails and 1 to represent heads, and 0 to indicate the coin flipped was a penny and 1 to indicate the coin flipped was a quarter. With this encoding, the data and labels are:

$$
X = \begin{bmatrix} 0 & 1 & \dots & 0 \\ 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 1 & \dots & 1 \end{bmatrix}
\qquad
Y = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 0 \end{bmatrix}
$$

We are going to use t-distributed stochastic neighbor embedding (t-SNE) to visualize the data (D-dimensional sequences of coin flips), and investigate the effects of the number of observations and the number of dimensions on the resulting visualization. There are t-SNE implementations available for both Matlab and Python; I strongly encourage you to use one of these implementations over writing your own.

---

[1] We will assume all the coins are perfectly fair, that is there is no dynamical bias in these coin flips (*i.e.*, $p(H) = 0.5$ and $p(T) = 0.5$). See "Dynamical Bias in the Coin Toss" by Persi Diaconis, Susan Holmes, and Richard Montgomery for a discussion of dynamical bias in coin flips (posted to Sakai under Syllabus and Additional Resources $\rightarrow$ Interesting Articles).

(20)   1.  Consider first the case where each coin is flipped twice ($D = 2$).

(a)  Simulate sequences of coin flips for $P = 5$ pennies and $Q = 5$ quarters. Visualize the data using t-SNE, clearly denoting the data points that are associated with common coin flip sequences, *i.e.*, use color and or marker type to encode which data points are due to HH, which are due to HT, which are due to TH, and which are due to TT.



Figure: P and Q are not seperated



Figure: P and Q are seperated

(b) Simulate sequences of coin flips for $P = 10$ pennies and $Q = 10$ quarters. Visualize the data using t-SNE, clearly denoting the data points that are associated with common coin flip sequences.

Figure: P and Q  are not seperated



1(b) P=10, Q=10, D=2, t-SNE

Legend:
- TT: total points 5
- TH: total points 5
- HT: total points 5
- HH: total points 5

Figure: P and Q  are seperated



1(b) P=10, Q=10, D=2, t-SNE (P is square, 0 in Z)

Legend:
- TT: P total points 3
- TT: Q total points 2
- TH: P total points 2
- TH: Q total points 3
- HT: P total points 2
- HT: Q total points 3
- HH: P total points 3
- HH: Q total points 2

(c) Simulate sequences of coin flips for $P = 25$ pennies and $Q = 25$ quarters. Visualize the data using t-SNE, clearly denoting the data points that are associated with common coin flip sequences.

P and Q  are not seperated



1(c) P=25, Q=25, D=2, t-SNE

| | |
|---|---|
| ■ | TT: total points 13 |
| ▲ | TH: total points 11 |
| + | HT: total points 11 |
| ● | HH: total points 15 |

P and Q  are seperated



1(c) P=25, Q=25, D=2, t-SNE (P is square, 0 in Z)

| | |
|---|---|
| ■ | TT: P total points 8 |
| ▲ | TT: Q total points 5 |
| ■ | TH: P total points 5 |
| ▲ | TH: Q total points 6 |
| ■ | HT: P total points 8 |
| ▲ | HT: Q total points 3 |
| ■ | HH: P total points 4 |
| ▲ | HH: Q total points 11 |

(d) Simulate sequences of coin flips for $P = 50$ pennies and $Q = 50$ quarters. Visualize the data using t-SNE, clearly denoting the data points that are associated with common coin flip sequences.

P and Q are not seperated



1(d) P=50, Q=50, D=2, t-SNE

Legend:
- TT: total points 30
- TH: total points 15
- HT: total points 23
- HH: total points 32

P and Q are seperated



1(d) P=50, Q=50, D=2, t-SNE (P is square, 0 in Z)

Legend:
- TT: P total points 12
- TT: Q total points 18
- TH: P total points 6
- TH: Q total points 9
- HT: P total points 16
- HT: Q total points 7
- HH: P total points 16
- HH: Q total points 16

(20)  2. Now consider first the case where each coin is flipped three times ($D = 3$).

(a) Simulate sequences of coin flips for $P = 5$ pennies and $Q = 5$ quarters. Visualize the data using t-SNE, clearly denoting the data points that are associated with common coin flip sequences.

P and Q are not seperated



2(a) P=5, Q=5, D=3, t-SNE

| | |
|---|---|
| | TTT: total points 0 |
| ▲ | TTH: total points 4 |
| + | THT: total points 1 |
| ● | THH: total points 3 |
| | HTT: total points 0 |
| ▶ | HTH: total points 1 |
| ◀ | HHT: total points 1 |
| | HHH: total points 0 |

P and Q are seperated



2(a) P=5, Q=5, D=3, t-SNE (P is square,0 in Z)

| | |
|---|---|
| | TTT: P total points 0 |
| | TTT: Q total points 0 |
| ■ | TTH: P total points 3 |
| ▲ | TTH: Q total points 1 |
| | THT: P total points 0 |
| ▲ | THT: Q total points 1 |
| ■ | THH: P total points 1 |
| ▲ | THH: Q total points 2 |
| | HTT: P total points 0 |
| | HTT: Q total points 0 |
| ■ | HTH: P total points 1 |
| | HTH: Q total points 0 |
| | HHT: P total points 0 |
| ▲ | HHT: Q total points 1 |
| | HHH: P total points 0 |
| | HHH: Q total points 0 |

(b) Simulate sequences of coin flips for $P = 10$ pennies and $Q = 10$ quarters. Visualize the data using t-SNE, clearly denoting the data points that are associated with common coin flip sequences.

P and Q  are not seperated

2(b) P=10, Q=10, D=3, t-SNE



| | |
|---|---|
| ■ | TTT: total points 2 |
| ▲ | TTH: total points 5 |
| · | THT: total points 1 |
| ● | THH: total points 4 |
| ▼ | HTT: total points 2 |
| ▶ | HTH: total points 2 |
| ◀ | HHT: total points 2 |
| ▾ | HHH: total points 2 |

P and Q  are seperated

2(b) P=10, Q=10, D=3, t-SNE (P is square,0 in Z)



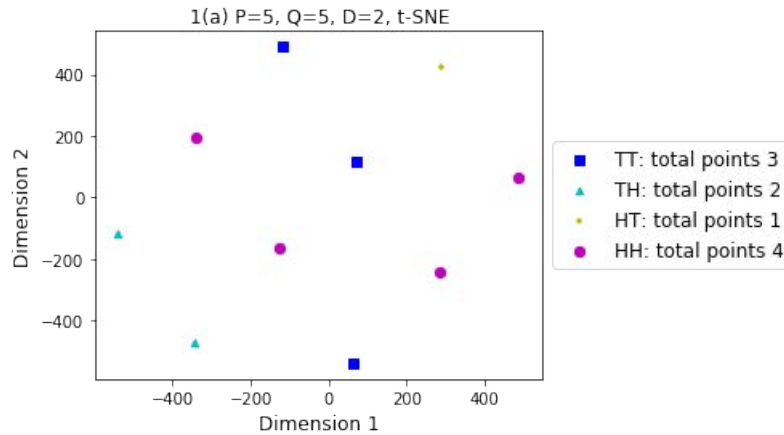| | |
|---|---|
| ■ | TTT: P total points 1 |
| ▲ | TTT: Q total points 1 |
| ■ | TTH: P total points 2 |
| ▲ | TTH: Q total points 3 |
| | THT: P total points 0 |
| ▴ | THT: Q total points 1 |
| ■ | THH: P total points 2 |
| ▲ | THH: Q total points 2 |
| ■ | HTT: P total points 1 |
| ▴ | HTT: Q total points 1 |
| ■ | HTH: P total points 2 |
| | HTH: Q total points 0 |
| ■ | HHT: P total points 1 |
| ▴ | HHT: Q total points 1 |
| ■ | HHH: P total points 1 |
| ▴ | HHH: Q total points 1 |

(c) Simulate sequences of coin flips for $P = 25$ pennies and $Q = 25$ quarters. Visualize the data using t-SNE, clearly denoting the data points that are associated with common coin flip sequences.

P and Q are not seperated



2(c) P=25, Q=25, D=3, t-SNE

| | |
|---|---|
| ■ | TTT: total points 3 |
| ▲ | TTH: total points 9 |
| + | THT: total points 8 |
| ● | THH: total points 4 |
| ▼ | HTT: total points 5 |
| ▶ | HTH: total points 6 |
| ◀ | HHT: total points 9 |
| ▼ | HHH: total points 6 |

P and Q are seperated



2(c) P=25, Q=25, D=3, t-SNE (P is square,0 in Z)

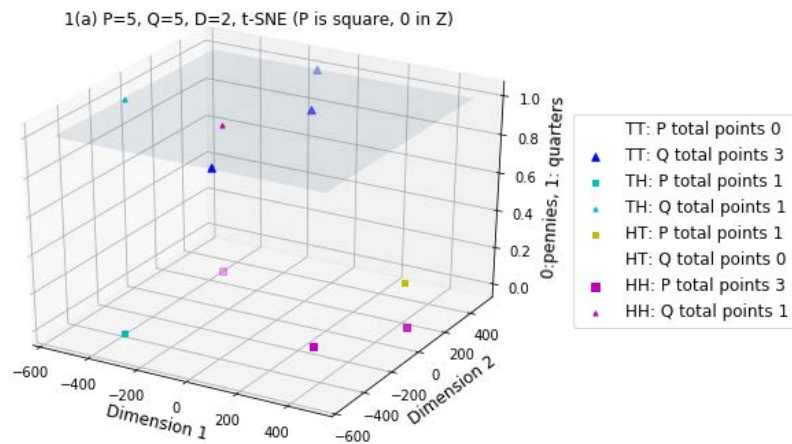| | |
|---|---|
| | TTT: P total points 0 |
| ▲ | TTT: Q total points 3 |
| ■ | TTH: P total points 6 |
| ▲ | TTH: Q total points 3 |
| ■ | THT: P total points 4 |
| ▲ | THT: Q total points 4 |
| ■ | THH: P total points 3 |
| ▴ | THH: Q total points 1 |
| ■ | HTT: P total points 3 |
| ▲ | HTT: Q total points 2 |
| ■ | HTH: P total points 4 |
| ▲ | HTH: Q total points 2 |
| ■ | HHT: P total points 2 |
| ▲ | HHT: Q total points 7 |
| ■ | HHH: P total points 3 |
| ▲ | HHH: Q total points 3 |

(d) Simulate sequences of coin flips for $P = 50$ pennies and $Q = 50$ quarters. Visualize the data using t-SNE, clearly denoting the data points that are associated with common coin flip sequences.
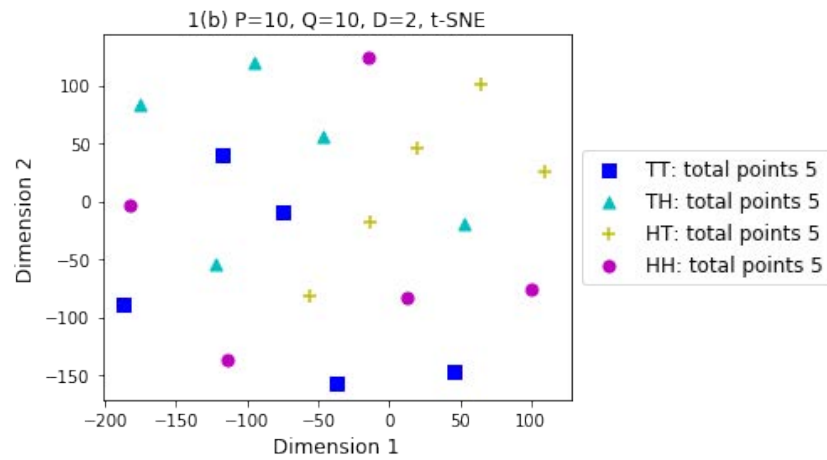
P and Q  are not seperated



2(d) P=50, Q=50, D=3, t-SNE

Legend:
- ■ TTT: total points 6
- ▲ TTH: total points 14
- + THT: total points 15
- ● THH: total points 11
- ▼ HTT: total points 13
- ▶ HTH: total points 14
- ◀ HHT: total points 11
- Y HHH: total points 16

P and Q  are seperated



2(d) P=50, Q=50, D=3, t-SNE (P is square,0 in Z)

Legend:
- ■ TTT: P total points 5
- ▴ TTT: Q total points 1
- ■ TTH: P total points 6
- ▲ TTH: Q total points 8
- ■ THT: P total points 7
- ▲ THT: Q total points 8
- ■ THH: P total points 7
- ▲ THH: Q total points 4
- ■ HTT: P total points 5
- ▲ HTT: Q total points 8
- ■ HTH: P total points 9
- ▲ HTH: Q total points 5
- ▪ HHT: P total points 4
- ▲ HHT: Q total points 7
- ■ HHH: P total points 7
- ▲ HHH: Q total points 9

(20)  3.  Now consider first the case where each coin is flipped four times ($D = 4$).

(a)  Simulate sequences of coin flips for $P = 5$ pennies and $Q = 5$ quarters. Visualize the data using t-SNE, clearly denoting the data points that are associated with common coin flip sequences.
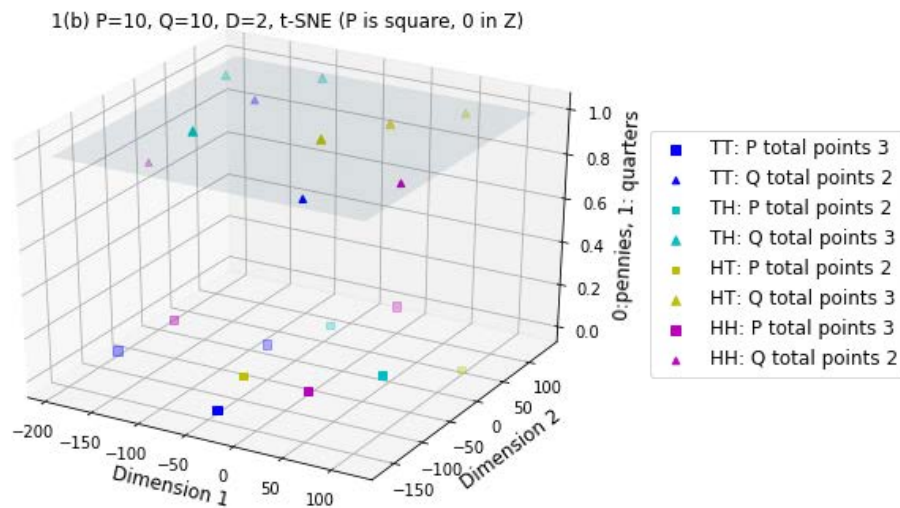
P and Q  are not seperated



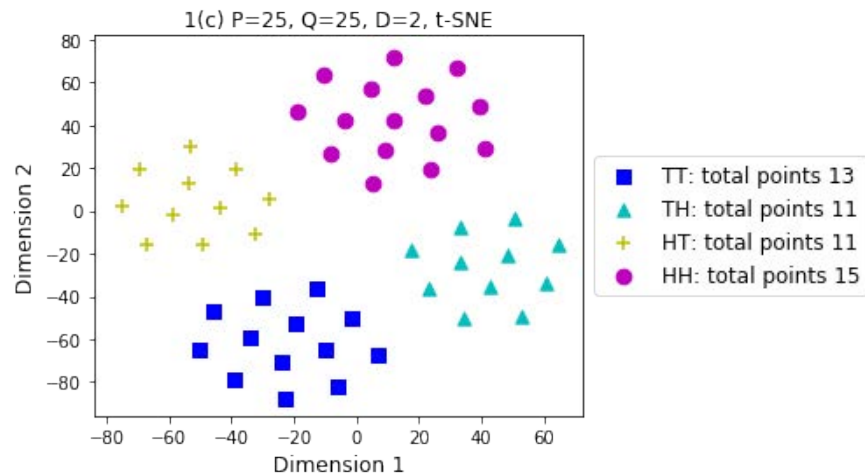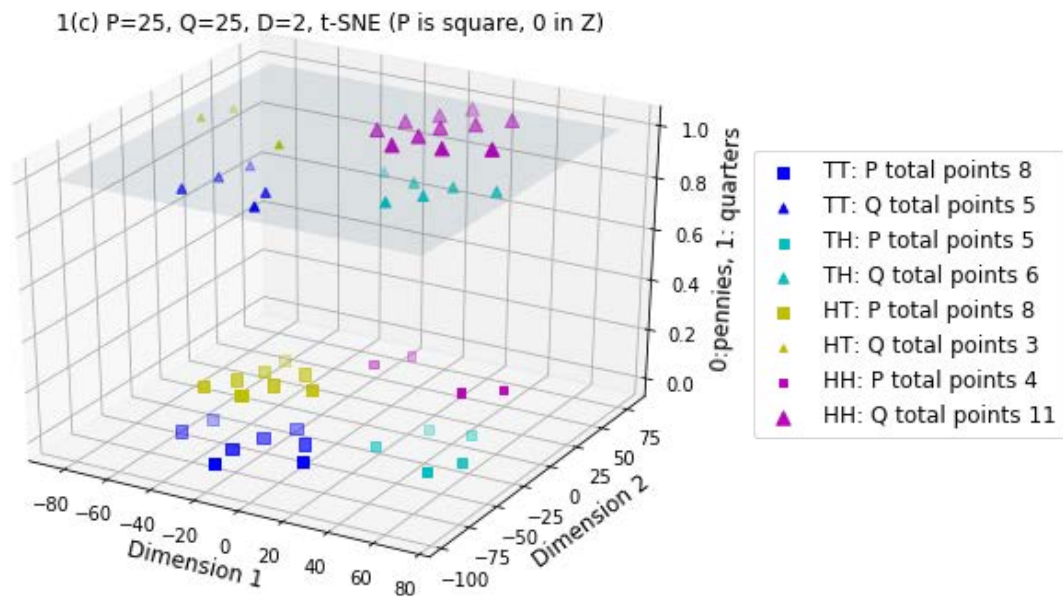3(a) P=5, Q=5, D=4, t-SNE

| | | | |
|---|---|---|---|
| ■ | TTTT: total points 1 | ⊥ | HTTT: total points 2 |
| | TTTH: total points 0 | | HTTH: total points 0 |
| + | TTHT: total points 3 | | HTHT: total points 0 |
| | TTHH: total points 0 | | HTHH: total points 0 |
| ▼ | THTT: total points 1 | ⬟ | HHTT: total points 2 |
| | THTH: total points 0 | | HHTH: total points 0 |
| ◄ | THHT: total points 1 | | HHHT: total points 0 |
| | THHH: total points 0 | | HHHH: total points 0 |

P and Q  are seperated



3(a) P=5, Q=5, D=4, t-SNE (P is square,0 in Z)

| | | | |
|---|---|---|---|
| | TTTT: P total points 0 | ■ | HTTT: P total points 1 |
| ▲ | TTTT: Q total points 1 | ▲ | HTTT: Q total points 1 |
| | TTTH: P total points 0 | | HTTH: P total points 0 |
| | TTTH: Q total points 0 | | HTTH: Q total points 0 |
| ■ | TTHT: P total points 3 | | HTHT: P total points 0 |
| | TTHT: Q total points 0 | | HTHT: Q total points 0 |
| | TTHH: P total points 0 | | HTHH: P total points 0 |
| | TTHH: Q total points 0 | | HTHH: Q total points 0 |
| ■ | THTT: P total points 1 | | HHTT: P total points 0 |
| | THTT: Q total points 0 | ▲ | HHTT: Q total points 2 |
| | THTH: P total points 0 | | HHTH: P total points 0 |
| | THTH: Q total points 0 | | HHTH: Q total points 0 |
| | THHT: P total points 0 | | HHHT: P total points 0 |
| ▲ | THHT: Q total points 1 | | HHHT: Q total points 0 |
| | THHH: P total points 0 | | HHHH: P total points 0 |
| | THHH: Q total points 0 | | HHHH: Q total points 0 |

(b) Simulate sequences of coin flips for $P = 10$ pennies and $Q = 10$ quarters. Visualize the data using t-SNE, clearly denoting the data points that are associated with common coin flip sequences.
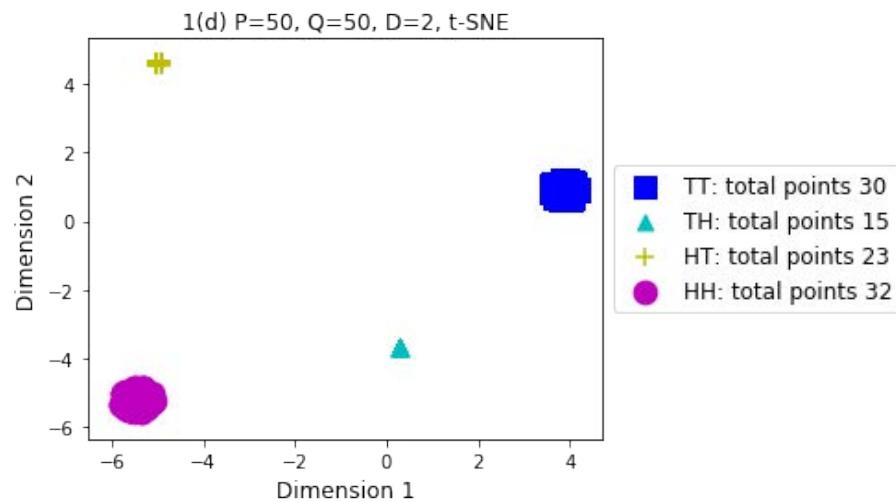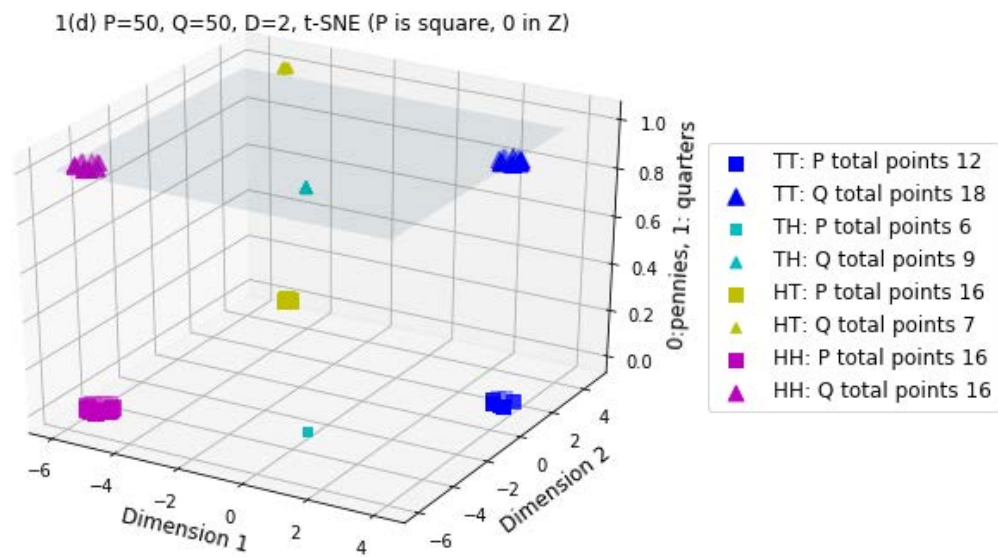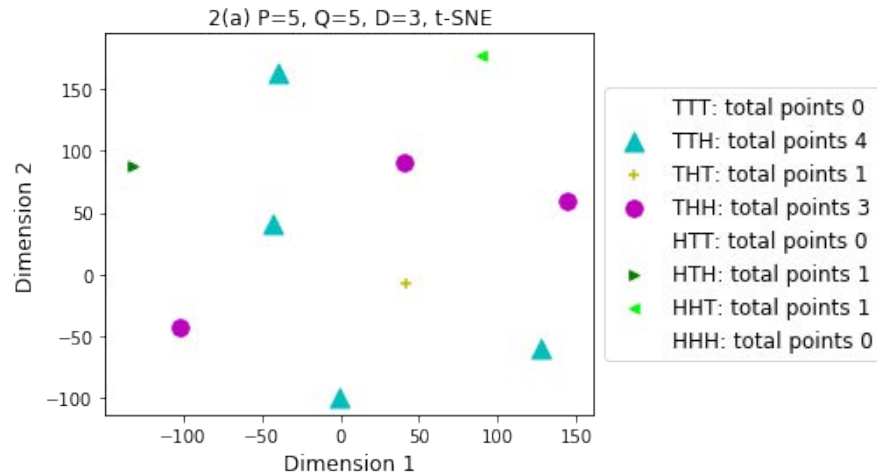
P and Q are not seperated



3(b) P=10, Q=10, D=4, t-SNE

| | | | |
|---|---|---|---|
| ■ | TTTT: total points 2 | | HTTT: total points 0 |
| ▲ | TTTH: total points 1 | ◄ | HTTH: total points 2 |
| | TTHT: total points 0 | ► | HTHT: total points 2 |
| ● | TTHH: total points 3 | ● | HTHH: total points 1 |
| ▼ | THTT: total points 6 | ● | HHTT: total points 1 |
| | THTH: total points 0 | | HHTH: total points 0 |
| ◄ | THHT: total points 1 | | HHHT: total points 0 |
| ▼ | THHH: total points 1 | | HHHH: total points 0 |

P and Q are seperated



3(b) P=10, Q=10, D=4, t-SNE (P is square,0 in Z)

| | | | |
|---|---|---|---|
| | TTTT: P total points 0 | | HTTT: P total points 0 |
| ▲ | TTTT: Q total points 2 | | HTTT: Q total points 0 |
| | TTTH: P total points 0 | ■ | HTTH: P total points 1 |
| ▲ | TTTH: Q total points 1 | ▲ | HTTH: Q total points 1 |
| | TTHT: P total points 0 | ■ | HTHT: P total points 1 |
| | TTHT: Q total points 0 | ▲ | HTHT: Q total points 1 |
| ■ | TTHH: P total points 3 | ■ | HTHH: P total points 1 |
| | TTHH: Q total points 0 | | HTHH: Q total points 0 |
| ■ | THTT: P total points 4 | | HHTT: P total points 0 |
| ▲ | THTT: Q total points 2 | ▲ | HHTT: Q total points 1 |
| | THTH: P total points 0 | | HHTH: P total points 0 |
| | THTH: Q total points 0 | | HHTH: Q total points 0 |
| | THHT: P total points 0 | | HHHT: P total points 0 |
| ▲ | THHT: Q total points 1 | | HHHT: Q total points 0 |
| | THHH: P total points 0 | | HHHH: P total points 0 |
| ▲ | THHH: Q total points 1 | | HHHH: Q total points 0 |

(c) Simulate sequences of coin flips for $P = 25$ pennies and $Q = 25$ quarters. Visualize the data using t-SNE, clearly denoting the data points that are associated with common coin flip sequences.
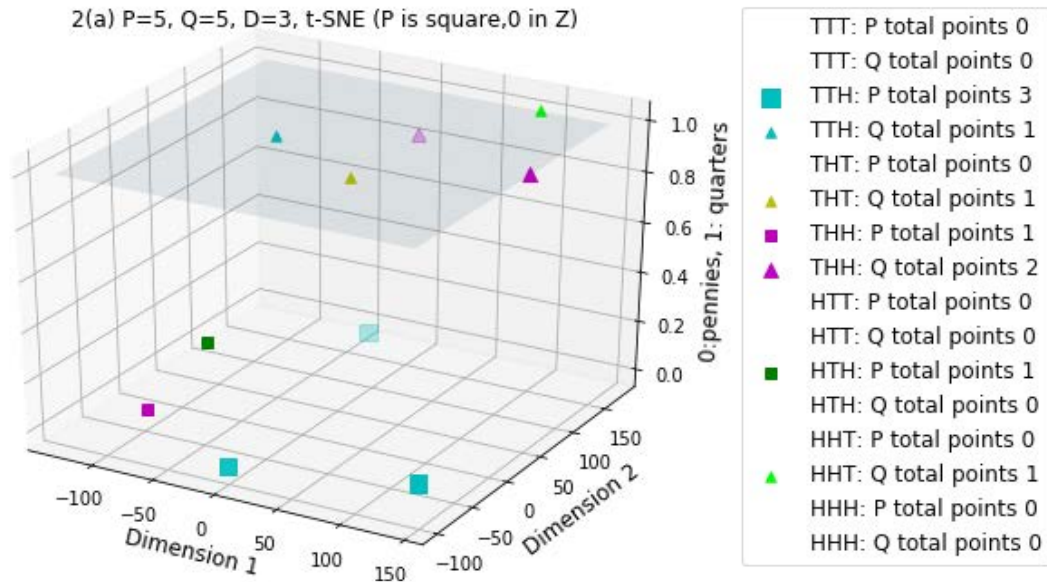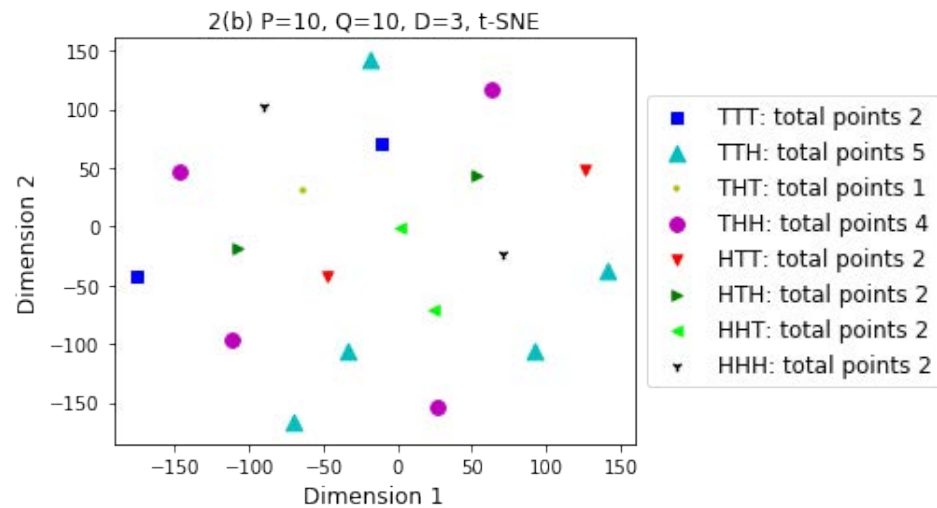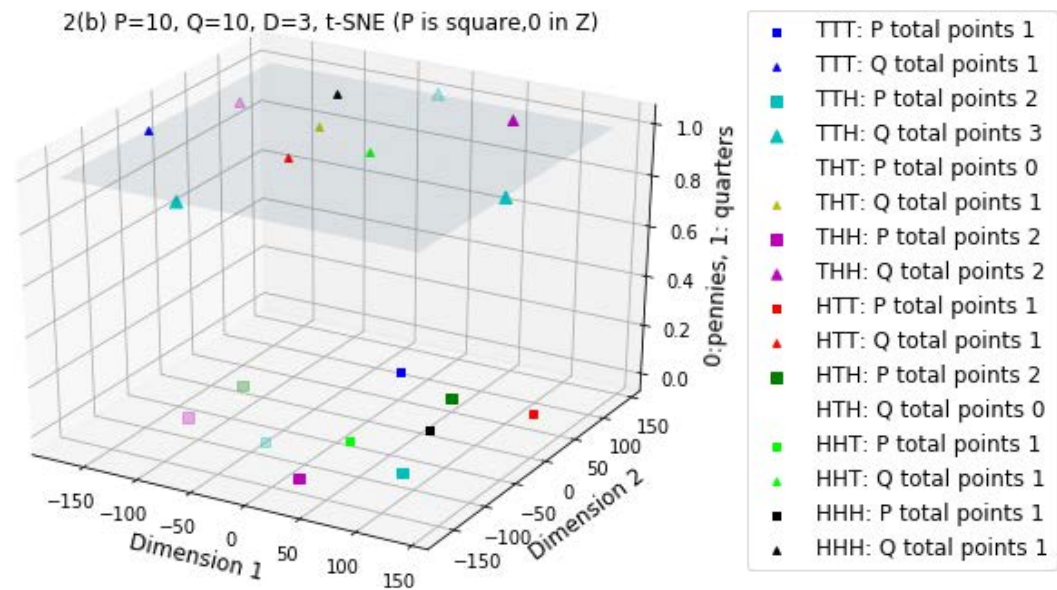
P and Q are not seperated



3(c) P=25, Q=25, D=4, t-SNE

| | | | | |
|---|---|---|---|---|
| ■ | TTTT: total points 2 | | · | HTTT: total points 1 |
| ▲ | TTTH: total points 4 | | · | HTTH: total points 1 |
| ✦ | TTHT: total points 3 | | ► | HTHT: total points 4 |
| ● | TTHH: total points 3 | | ● | HTHH: total points 3 |
| ▼ | THTT: total points 4 | | ● | HHTT: total points 2 |
| ► | THTH: total points 7 | | | HHTH: total points 3 |
| ◄ | THHT: total points 4 | | ★ | HHHT: total points 2 |
| ▼ | THHH: total points 3 | | ✕ | HHHH: total points 4 |

P and Q are seperated



3(c) P=25, Q=25, D=4, t-SNE (P is square,0 in Z)

| | | | | |
|---|---|---|---|---|
| | TTTT: P total points 0 | | | HTTT: P total points 0 |
| ▲ | TTTT: Q total points 2 | | ◦ | HTTT: Q total points 1 |
| ■ | TTTH: P total points 3 | | ■ | HTTH: P total points 1 |
| ▲ | TTTH: Q total points 1 | | | HTTH: Q total points 0 |
| ■ | TTHT: P total points 2 | | ■ | HTHT: P total points 4 |
| ▲ | TTHT: Q total points 1 | | | HTHT: Q total points 0 |
| ■ | TTHH: P total points 2 | | ■ | HTHH: P total points 1 |
| ▲ | TTHH: Q total points 1 | | ▲ | HTHH: Q total points 2 |
| ■ | THTT: P total points 1 | | | HHTT: P total points 0 |
| ▲ | THTT: Q total points 3 | | ▲ | HHTT: Q total points 2 |
| ■ | THTH: P total points 4 | | ■ | HHTH: P total points 2 |
| ▲ | THTH: Q total points 3 | | ◦ | HHTH: Q total points 1 |
| ■ | THHT: P total points 2 | | | HHHT: P total points 0 |
| ▲ | THHT: Q total points 2 | | ▲ | HHHT: Q total points 2 |
| ■ | THHH: P total points 1 | | ■ | HHHH: P total points 2 |
| ▲ | THHH: Q total points 2 | | ▲ | HHHH: Q total points 2 |

(d) Simulate sequences of coin flips for $P = 50$ pennies and $Q = 50$ quarters. Visualize the data using t-SNE, clearly denoting the data points that are associated with common coin flip sequences.
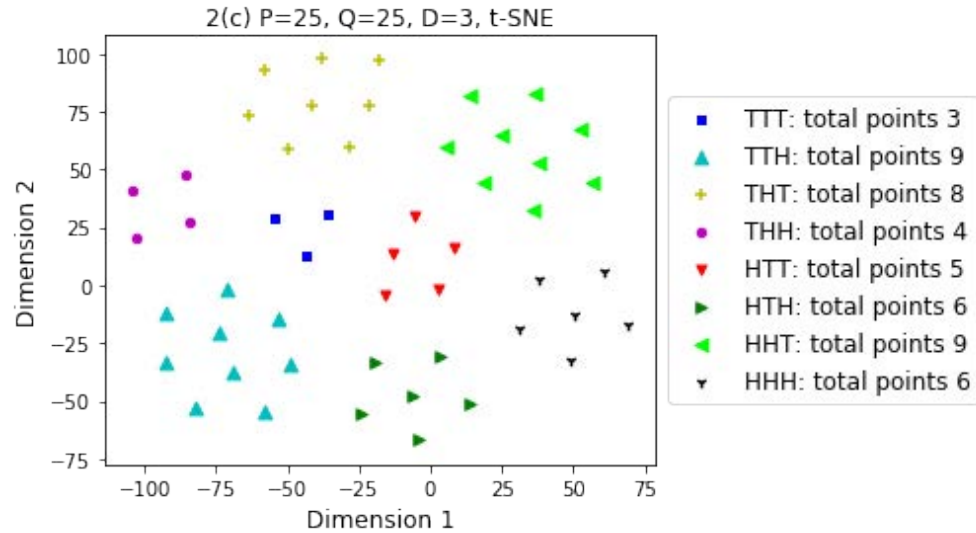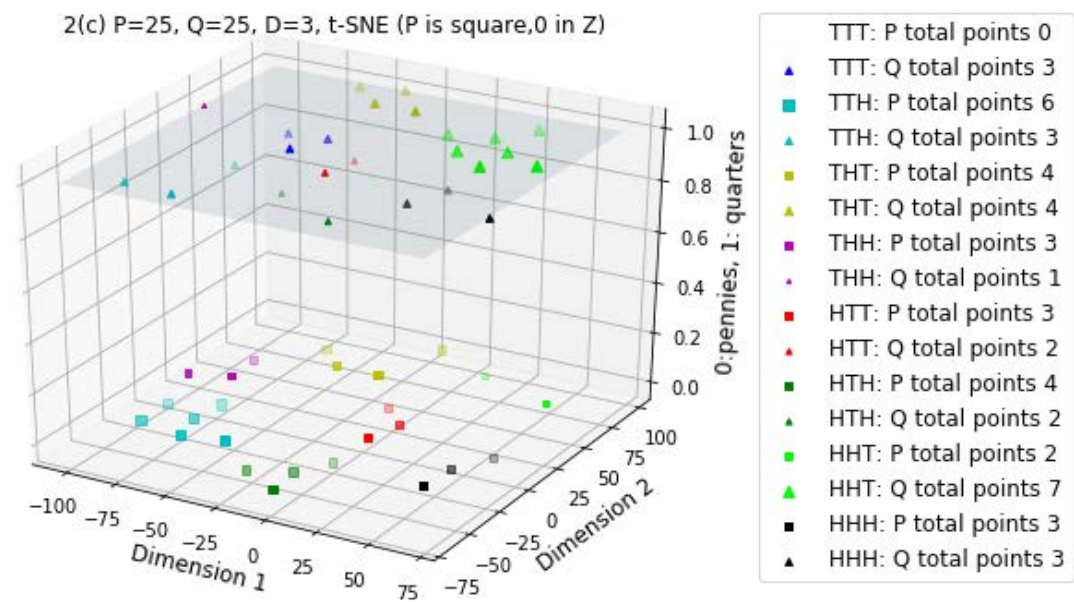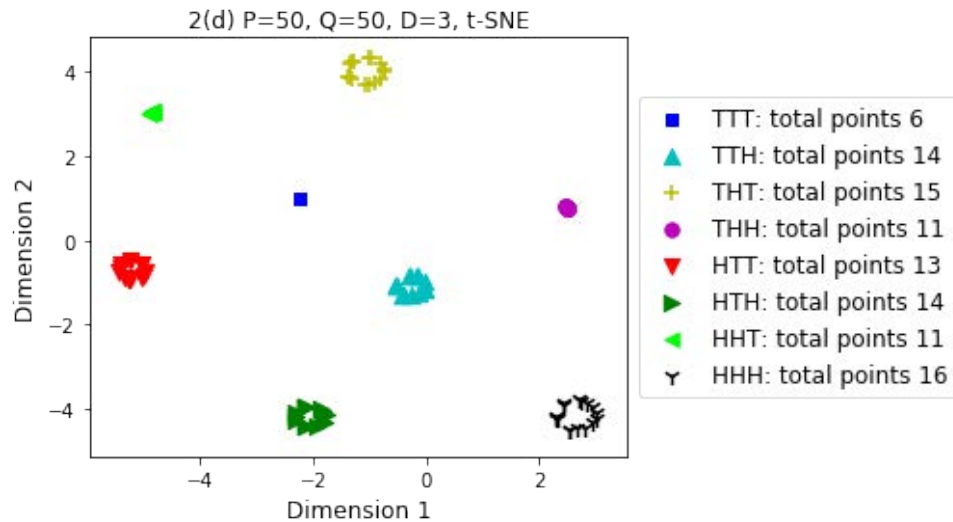
P and Q are not seperated



3(d) P=50, Q=50, D=4, t-SNE

| | | | |
|---|---|---|---|
| ■ | TTTT: total points 9 | ⅄ | HTTT: total points 8 |
| ▲ | TTTH: total points 4 | ⊣ | HTTH: total points 9 |
| + | TTHT: total points 5 | ⊢ | HTHT: total points 10 |
| • | TTHH: total points 2 | ● | HTHH: total points 7 |
| ▼ | THTT: total points 4 | ⬟ | HHTT: total points 7 |
| ▶ | THTH: total points 6 | | HHTH: total points 4 |
| ◀ | THHT: total points 7 | ★ | HHHT: total points 4 |
| ⅄ | THHH: total points 4 | ✕ | HHHH: total points 10 |

P and Q are seperated



3(d) P=50, Q=50, D=4, t-SNE (P is square,0 in Z)

| | | | |
|---|---|---|---|
| ■ | TTTT: P total points 4 | ■ | HTTT: P total points 5 |
| ▲ | TTTT: Q total points 5 | ▲ | HTTT: Q total points 3 |
| ■ | TTTH: P total points 2 | ■ | HTTH: P total points 5 |
| ▲ | TTTH: Q total points 2 | ▲ | HTTH: Q total points 4 |
| ■ | TTHT: P total points 3 | ■ | HTHT: P total points 7 |
| ▲ | TTHT: Q total points 2 | ▲ | HTHT: Q total points 3 |
| ■ | TTHH: P total points 1 | ■ | HTHH: P total points 3 |
| ▲ | TTHH: Q total points 1 | ▲ | HTHH: Q total points 4 |
| ■ | THTT: P total points 2 | ■ | HHTT: P total points 1 |
| ▲ | THTT: Q total points 2 | ▲ | HHTT: Q total points 6 |
| ■ | THTH: P total points 3 | | HHTH: P total points 3 |
| ▲ | THTH: Q total points 3 | | HHTH: Q total points 1 |
| ■ | THHT: P total points 6 | | HHHT: P total points 0 |
| ▲ | THHT: Q total points 1 | ▲ | HHHT: Q total points 4 |
| ■ | THHH: P total points 1 | ■ | HHHH: P total points 4 |
| ▲ | THHH: Q total points 3 | ▲ | HHHH: Q total points 6 |

Homework #5: Nonlinear Dimensionality Reduction (t-SNE)

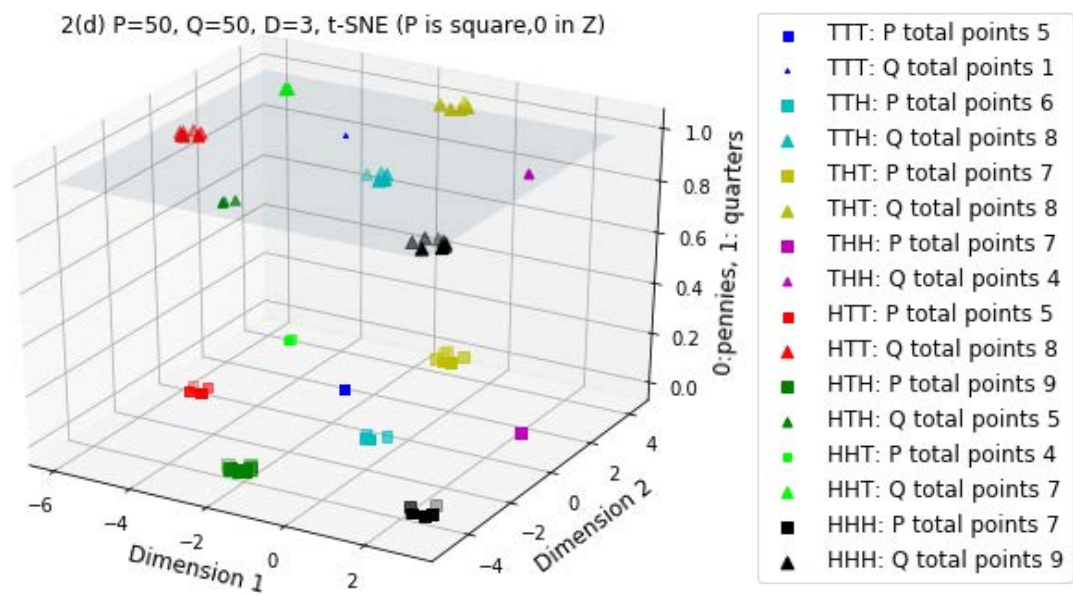(20)  4.  Now consider first the case where each coin is flipped four times ($D = 5$).

(a)  Simulate sequences of coin flips for $P = 5$ pennies and $Q = 5$ quarters.[2] Visualize the data using t-SNE, clearly denoting the data points that are associated with common coin flip sequences.
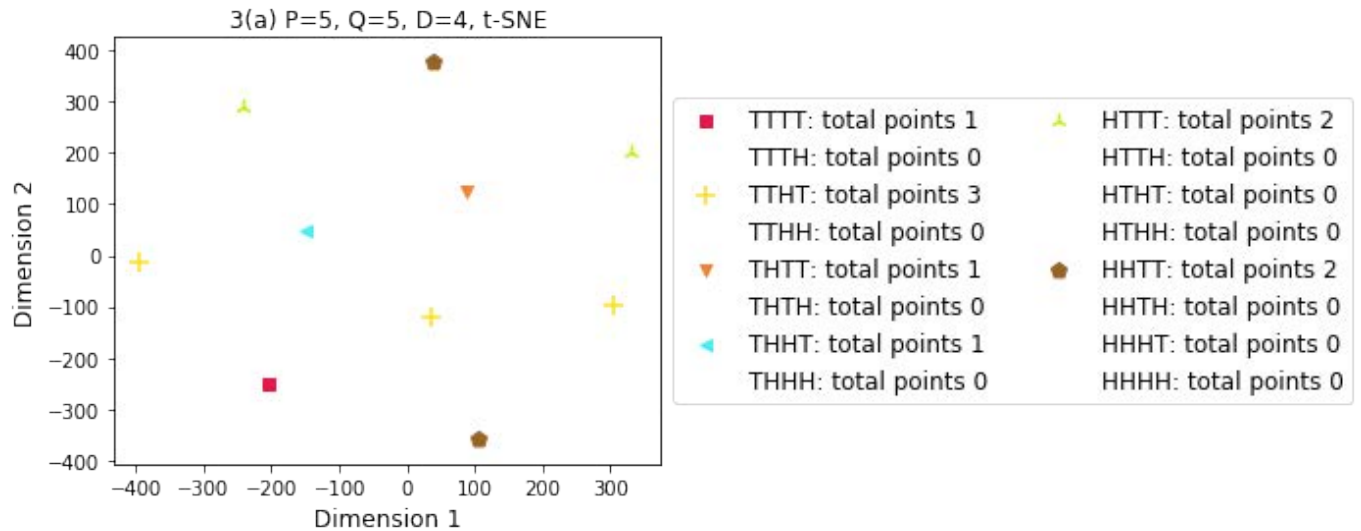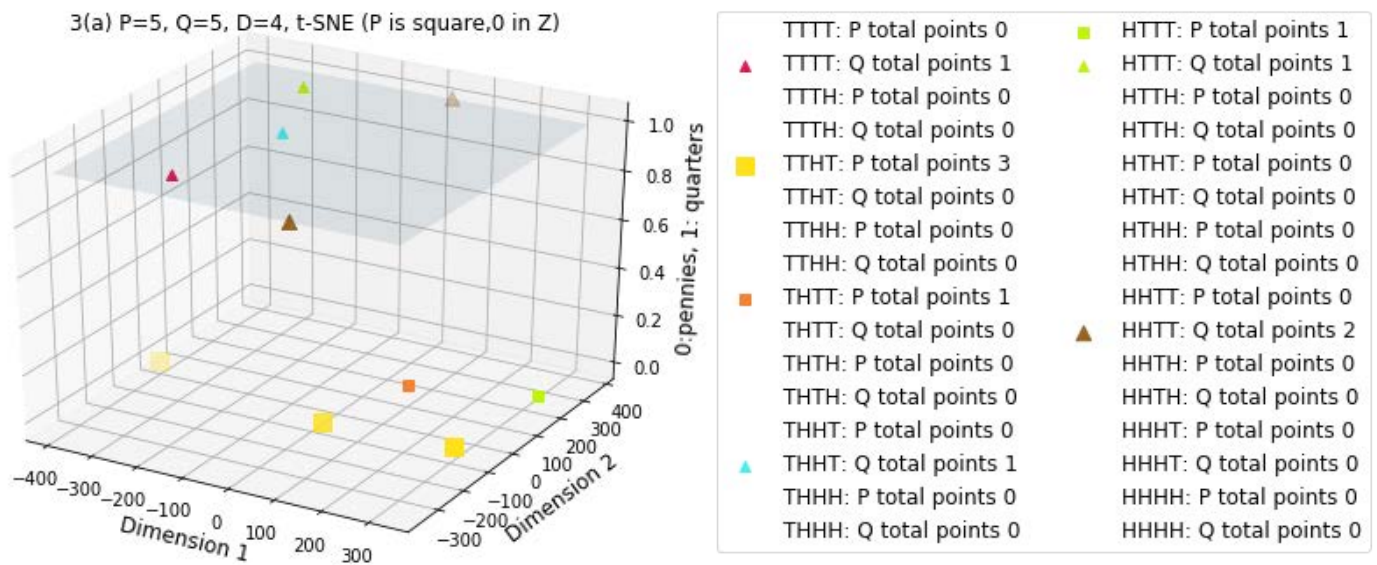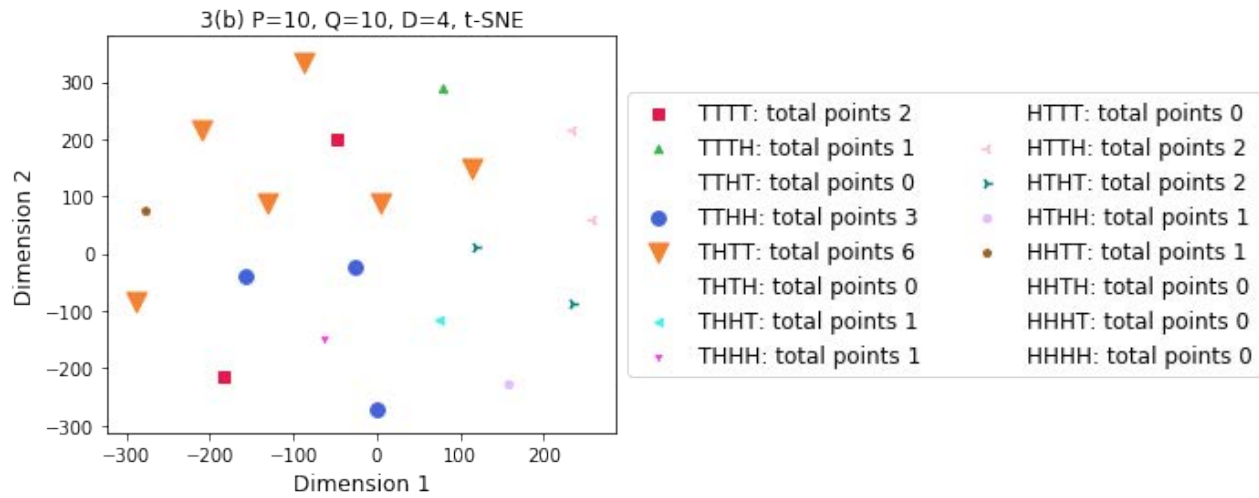
P and Q  are not seperated



4(a) P=5, Q=5, D=5, t-SNE

| | | | | | | |
|---|---|---|---|---|---|---|
| | TTTTT: total points 0 | | THTHH: total points 0 | | HTHHT: total points 0 |
| ▲ | TTTTH: total points 1 | | THHTT: total points 0 | | HTHHH: total points 0 |
| | TTTHT: total points 0 | | THHTH: total points 0 | | HHTTT: total points 0 |
| | TTTHH: total points 0 | ★ | THHHT: total points 2 | | HHTTH: total points 0 |
| | TTHTT: total points 0 | | THHHH: total points 0 | | HHTHT: total points 0 |
| ▶ | TTHTH: total points 1 | | HTTTT: total points 0 | | HHTHH: total points 0 |
| | TTHHT: total points 0 | | HTTTH: total points 0 | ▲ | HHHTT: total points 2 |
| | TTHHH: total points 0 | | HTTHT: total points 0 | | HHHTH: total points 0 |
| | THTTT: total points 0 | | HTTHH: total points 0 | ◀ | HHHHT: total points 1 |
| | THTTH: total points 2 | ● | HTHTT: total points 1 | | HHHHH: total points 0 |
| | THTHT: total points 0 | | HTHTH: total points 0 | | |

P and Q  are seperated



4(a) P=5, Q=5, D=5, t-SNE (P is square,0 in Z)

| | | | | | | |
|---|---|---|---|---|---|---|
| | TTTTT: P total points 0 | | THTHH: P total points 0 | | HTHTH: Q total points 0 |
| | TTTTT: Q total points 0 | | THTHH: Q total points 0 | | HTHHT: P total points 0 |
| ■ | TTTTH: P total points 1 | | THHTT: P total points 0 | | HTHHT: Q total points 0 |
| | TTTTH: Q total points 0 | | THHTT: Q total points 0 | | HTHHH: P total points 0 |
| | TTTHT: P total points 0 | | THHTH: P total points 0 | | HTHHH: Q total points 0 |
| | TTTHT: Q total points 0 | | THHTH: Q total points 0 | | HHTTT: P total points 0 |
| | TTTHH: P total points 0 | | THHHT: P total points 0 | | HHTTT: Q total points 0 |
| | TTTHH: Q total points 0 | ▲ | THHHT: Q total points 2 | | HHTTH: P total points 0 |
| | TTHTT: P total points 0 | | THHHH: P total points 0 | | HHTTH: Q total points 0 |
| | TTHTT: Q total points 0 | | THHHH: Q total points 0 | | HHTHT: P total points 0 |
| | TTHTH: P total points 0 | | HTTTT: P total points 0 | | HHTHT: Q total points 0 |
| ▲ | TTHTH: Q total points 1 | | HTTTT: Q total points 0 | | HHTHH: P total points 0 |
| | TTHHT: P total points 0 | | HTTTH: P total points 0 | | HHTHH: Q total points 0 |
| | TTHHT: Q total points 0 | | HTTTH: Q total points 0 | ■ | HHHTT: P total points 1 |
| | TTHHH: P total points 0 | | HTTHT: P total points 0 | ▲ | HHHTT: Q total points 1 |
| | TTHHH: Q total points 0 | | HTTHT: Q total points 0 | | HHHTH: P total points 0 |
| | THTTT: P total points 0 | | HTTHH: P total points 0 | | HHHTH: Q total points 0 |
| | THTTT: Q total points 0 | | HTTHH: Q total points 0 | | HHHHT: P total points 0 |
| ■ | THTTH: P total points 2 | ■ | HTHTT: P total points 1 | | HHHHT: Q total points 0 |
| | THTTH: Q total points 0 | | HTHTT: Q total points 0 | ▲ | HHHHT: Q total points 1 |
| | THTHT: P total points 0 | | HTHTH: P total points 0 | | HHHHH: P total points 0 |
| | THTHT: Q total points 0 | | | | | HHHHH: Q total points 0 |

_____

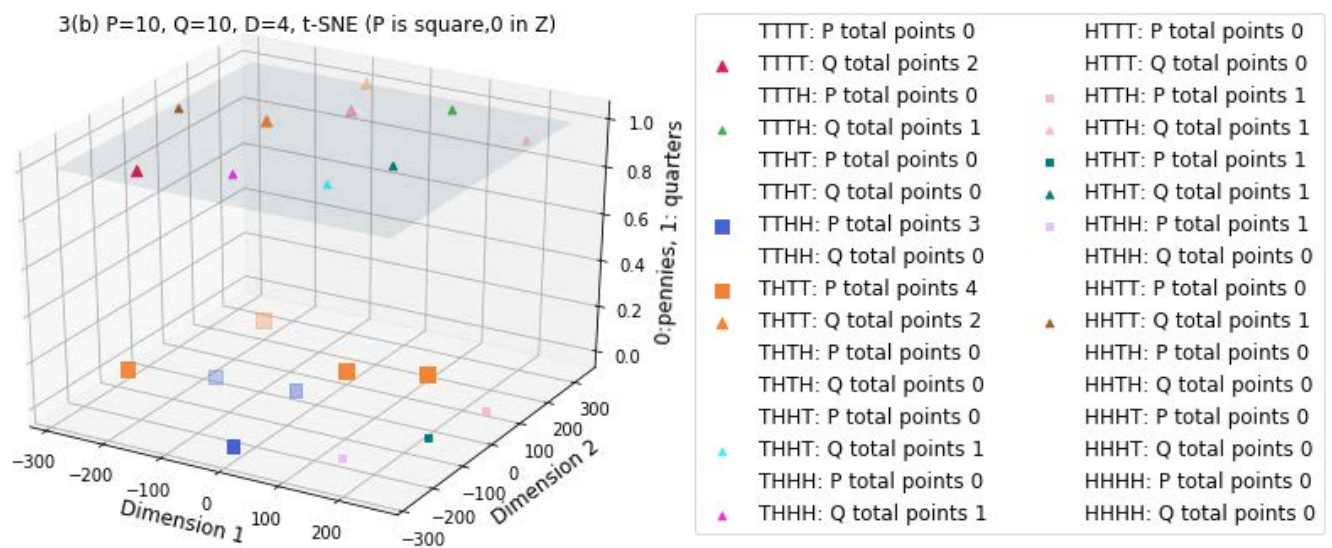[2]This is the scenario which you explored via Monte Carlo simulation.

(b) Simulate sequences of coin flips for $P = 10$ pennies and $Q = 10$ quarters. Visualize the data using t-SNE, clearly denoting the data points that are associated with common coin flip sequences.
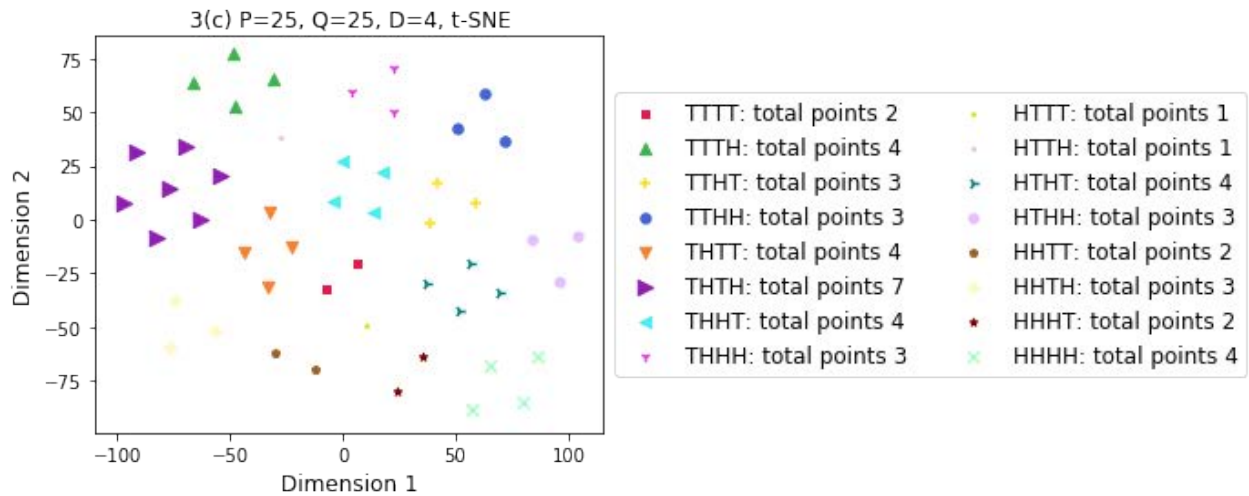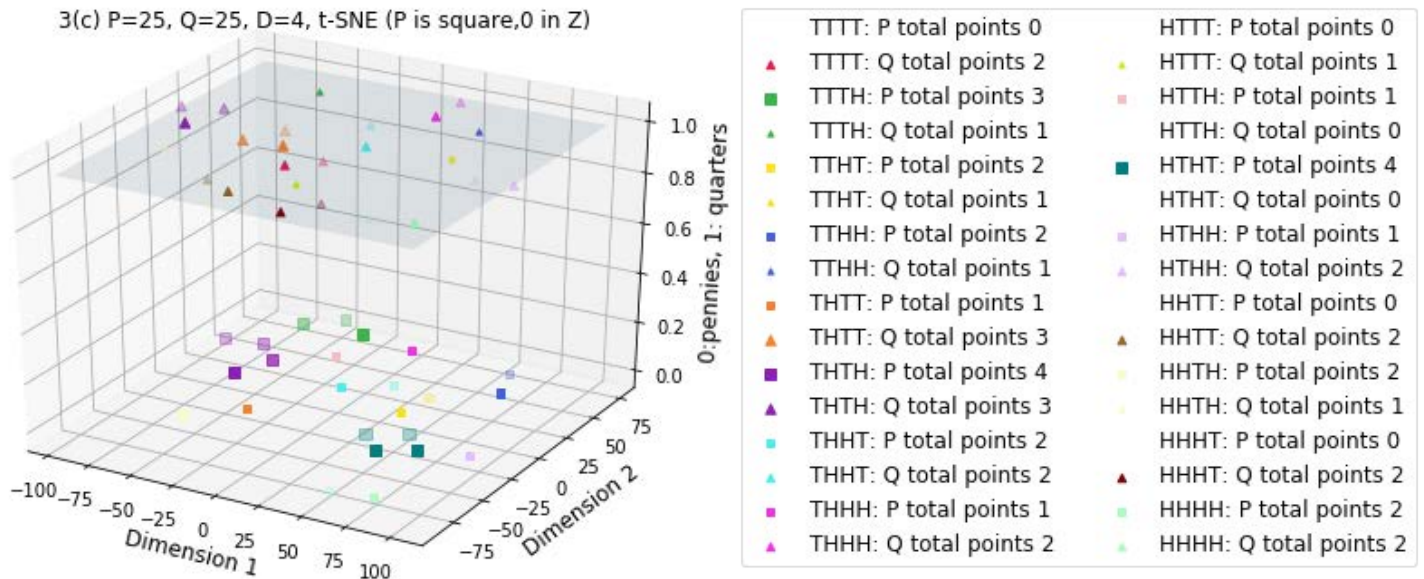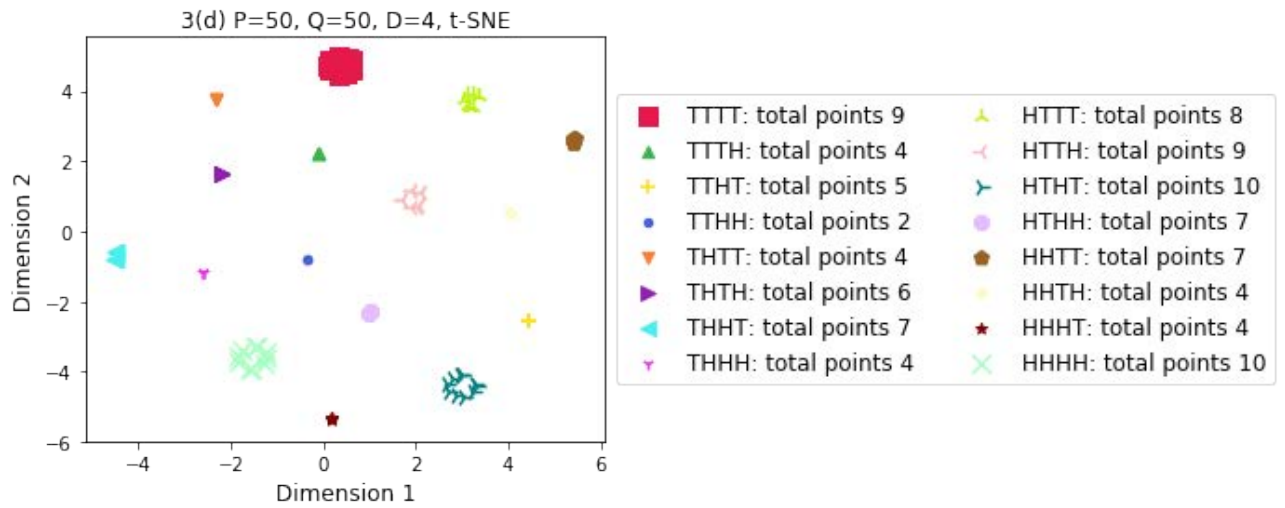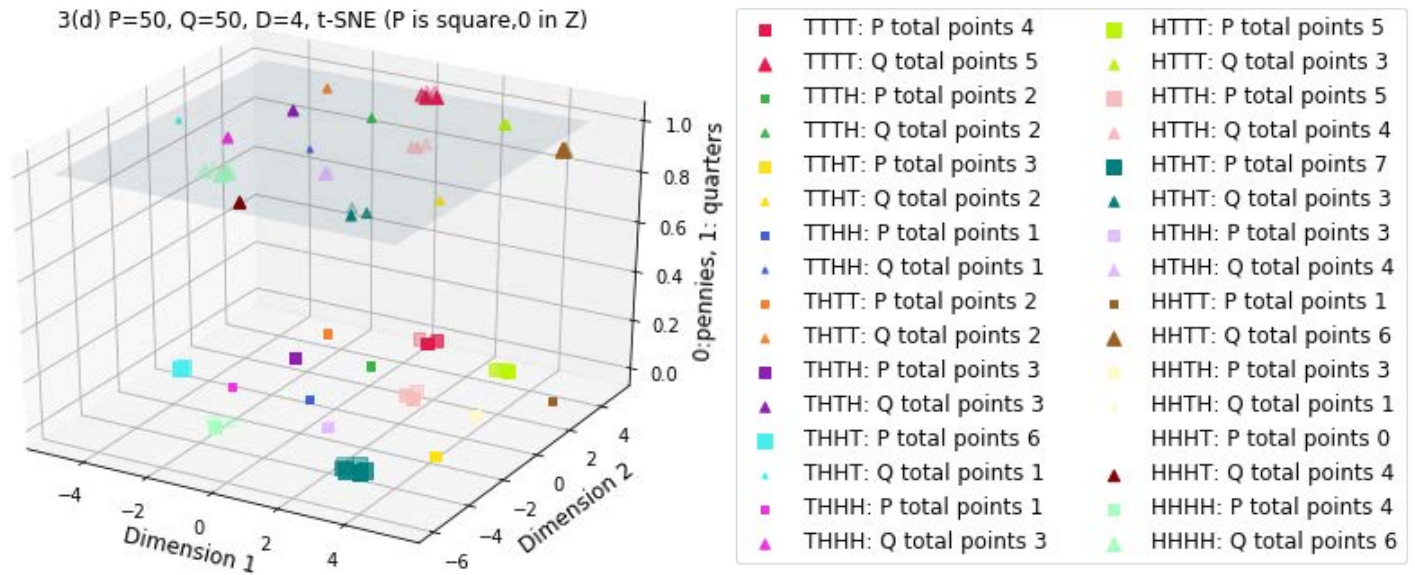
P and Q are not seperated



4(b) P=10, Q=10, D=5, t-SNE

| | | | | | |
|---|---|---|---|---|---|
| TTTTT: total points 0 | | THTHH: total points 1 | – | HTHHT: total points 1 | |
| TTTTH: total points 0 | ● | THHTT: total points 1 | | HTHHH: total points 0 | |
| TTटHT: total points 0 | | THHTH: total points 2 | | HHTTT: total points 0 | |
| TTTHH: total points 1 | ★ | THHHT: total points 2 | ǀ | HHTTH: total points 2 | |
| TTHTT: total points 0 | | THHHH: total points 0 | ◄ | HHTHT: total points 1 | |
| TTHTH: total points 0 | ✳ | HTTTT: total points 1 | ► | HHTHH: total points 1 | |
| TTHHT: total points 0 | ◆ | HTTTH: total points 1 | ▲ | HHHTT: total points 1 | |
| TTHHH: total points 0 | | HTTHT: total points 0 | | HHHTH: total points 0 | |
| THTTT: total points 0 | | HTTHH: total points 0 | | HHHHT: total points 0 | |
| THTTH: total points 1 | | HTHTT: total points 0 | | HHHHH: total points 0 | |
| THTHT: total points 0 | ● | HTHTT: total points 1 | | | |
| | ǀ | HTHTH: total points 3 | | | |

P and Q are seperated



4(b) P=10, Q=10, D=5, t-SNE (P is square,0 in Z)

| | | | | | |
|---|---|---|---|---|---|
| TTTTT: P total points 0 | ■ | THTHH: P total points 1 | ▲ | HTHTH: Q total points 2 | |
| TTTTT: Q total points 0 | | THTHH: Q total points 0 | ■ | HTHHT: P total points 1 | |
| TTTTH: P total points 0 | ■ | THHTT: P total points 1 | | HTHHT: Q total points 0 | |
| TTTTH: Q total points 0 | | THHTT: Q total points 0 | | HTHHH: P total points 0 | |
| TTTHT: P total points 0 | | THHTH: P total points 1 | | HTHHH: Q total points 0 | |
| TTTHT: Q total points 0 | | THHTH: Q total points 1 | | HHTTT: P total points 0 | |
| TTTHH: P total points 1 | ■ | THHHT: P total points 2 | | HHTTT: Q total points 0 | |
| TTTHH: Q total points 0 | | THHHT: Q total points 0 | | HHTTH: P total points 0 | |
| TTHTT: P total points 0 | | THHHH: P total points 0 | ▲ | HHTTH: Q total points 2 | |
| TTHTT: Q total points 0 | | THHHH: Q total points 0 | | HHTHT: P total points 0 | |
| TTHTH: P total points 0 | | HTTTT: P total points 0 | ▲ | HHTHT: Q total points 1 | |
| TTHTH: Q total points 0 | ▲ | HTTTT: Q total points 1 | | HHTHH: P total points 0 | |
| TTHHT: P total points 0 | | HTTTH: P total points 0 | ▲ | HHTHH: Q total points 1 | |
| TTHHT: Q total points 0 | ▲ | HTTTH: Q total points 1 | | HHHTT: P total points 0 | |
| TTHHH: P total points 0 | | HTTHT: P total points 0 | ▲ | HHHTT: Q total points 1 | |
| TTHHH: Q total points 0 | | HTTHT: Q total points 0 | | HHHTH: P total points 0 | |
| THTTT: P total points 0 | | HTTHH: P total points 0 | | HHHTH: Q total points 0 | |
| THTTT: Q total points 0 | | HTTHH: Q total points 0 | | HHHHT: P total points 0 | |
| THTTH: P total points 1 | ■ | HTHTT: P total points 1 | | HHHHT: Q total points 0 | |
| THTTH: Q total points 0 | | HTHTT: Q total points 0 | | HHHHH: P total points 0 | |
| THTHT: P total points 0 | ■ | HTHTH: P total points 1 | | HHHHH: Q total points 0 | |
| THTHT: Q total points 0 | | | | | |

(c) Simulate sequences of coin flips for $P = 25$ pennies and $Q = 25$ quarters. Visualize the data using t-SNE, clearly denoting the data points that are associated with common coin flip sequences.
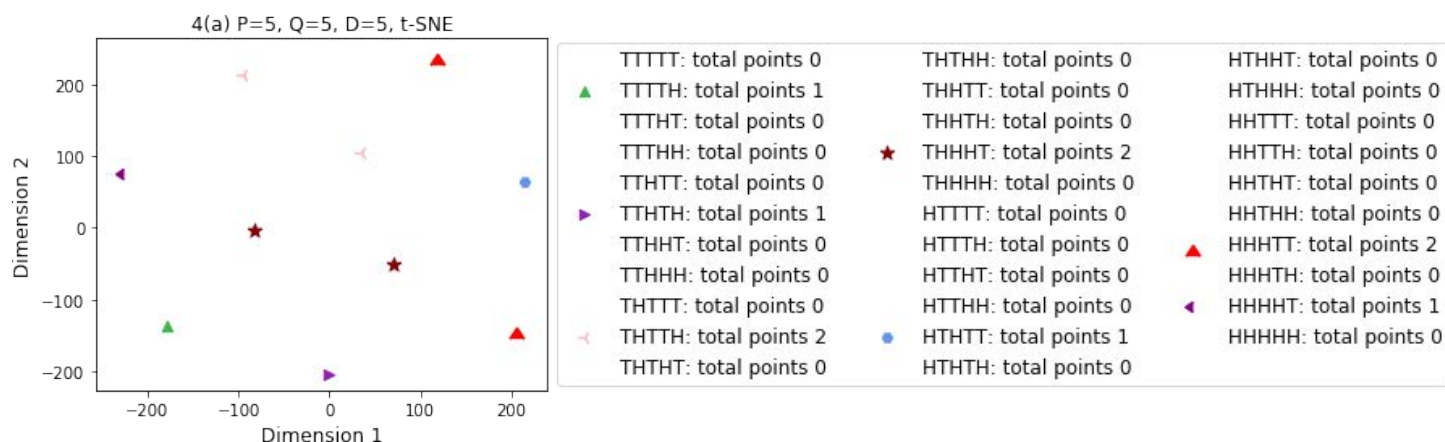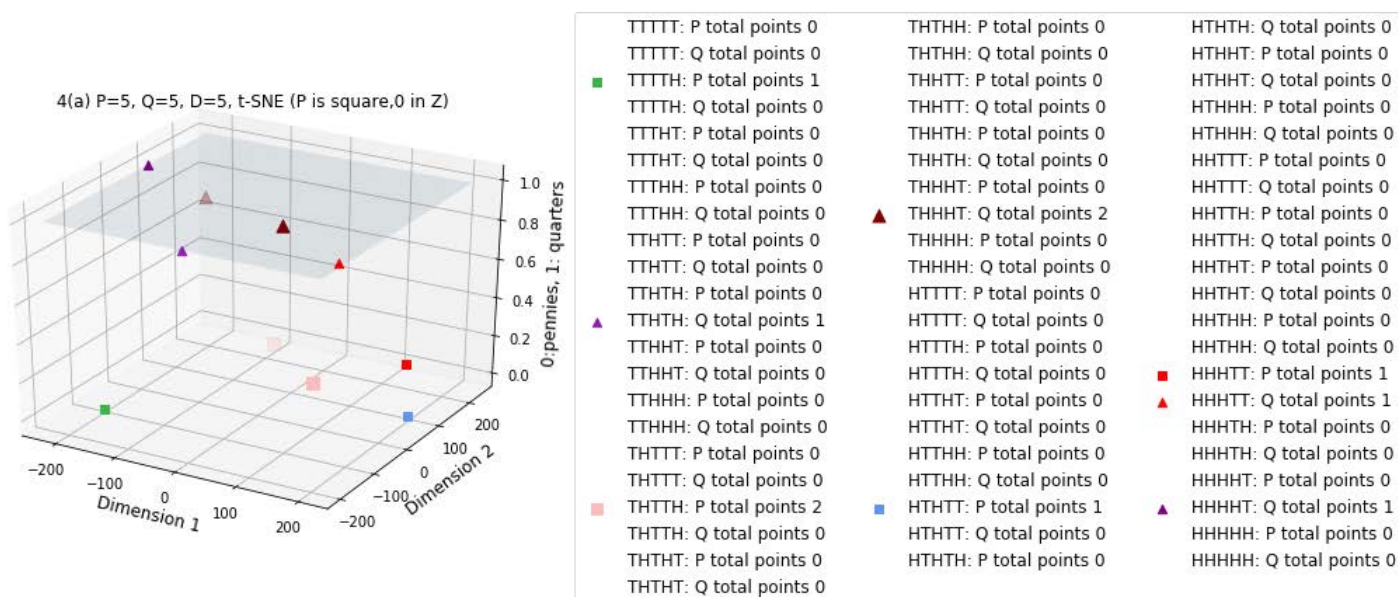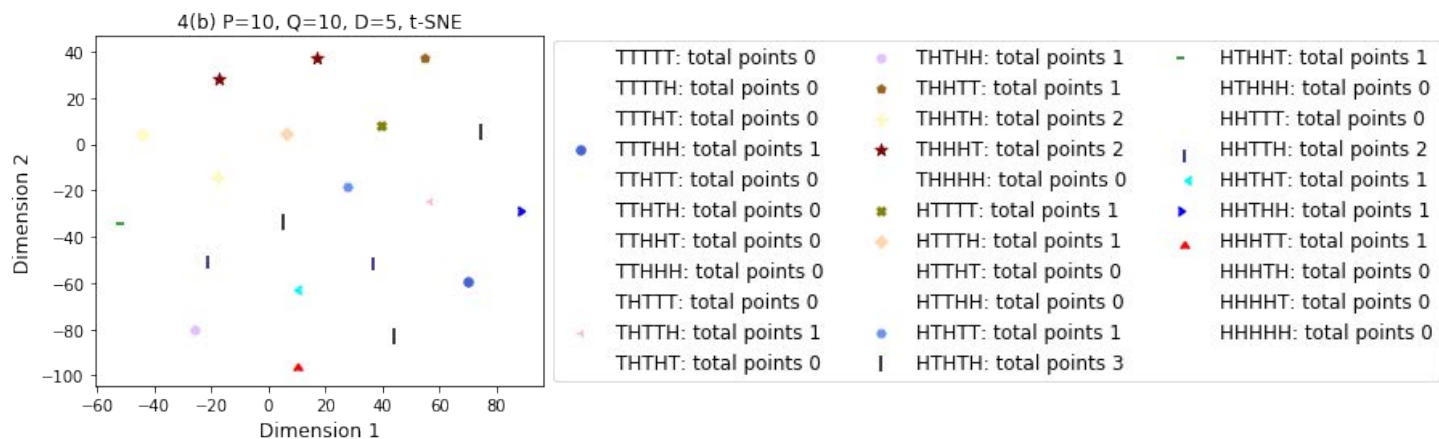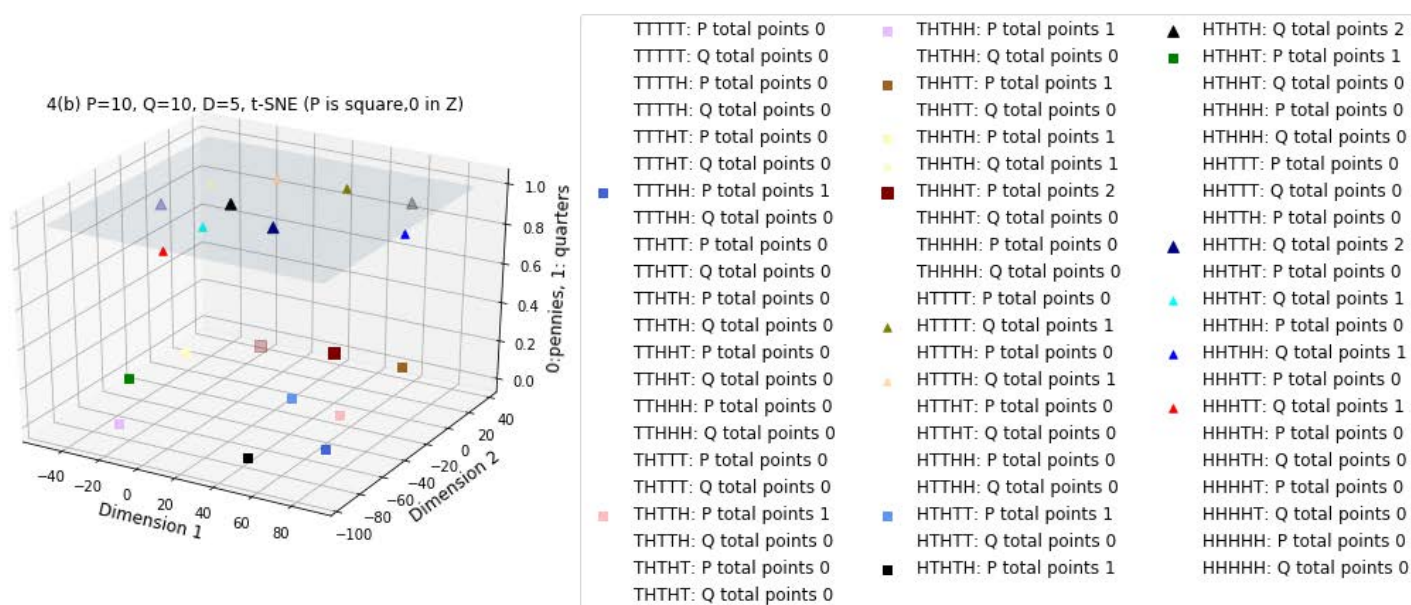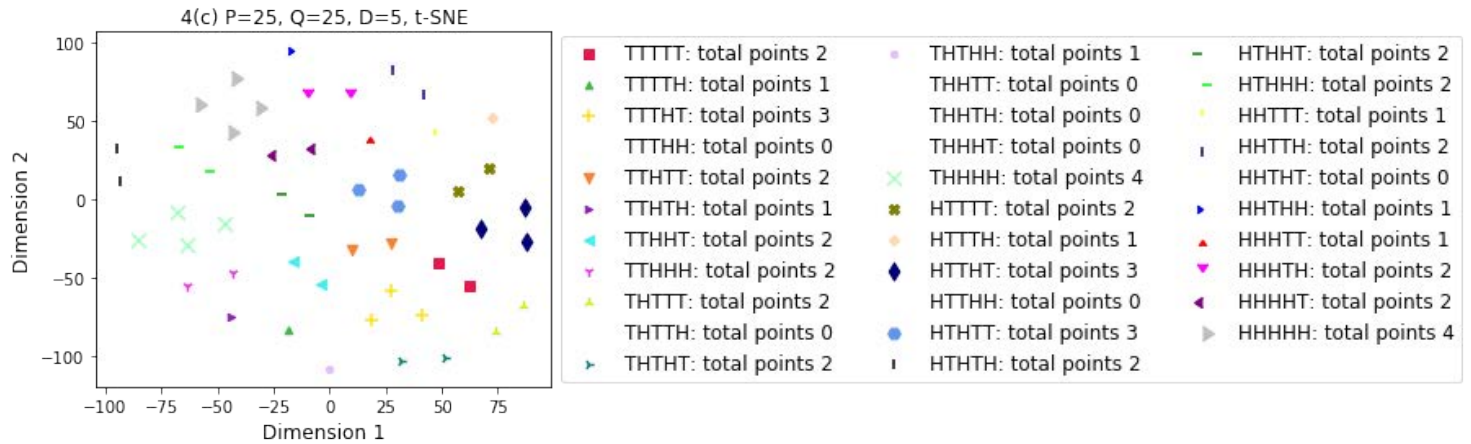
P and Q are not seperated



4(c) P=25, Q=25, D=5, t-SNE

| | |
|---|---|
| ■ TTTTT: total points 2 | • THTHH: total points 1 |
| ▲ TTTTH: total points 1 | THHTT: total points 0 |
| + TTTHT: total points 3 | THHTH: total points 0 |
| TTTHH: total points 0 | THHHT: total points 0 |
| ▼ TTHTT: total points 2 | × THHHH: total points 4 |
| ▶ TTHTH: total points 1 | ✱ HTTTT: total points 2 |
| ◀ TTHHT: total points 2 | • HTTTH: total points 1 |
| ▼ TTHHH: total points 2 | ◆ HTTHT: total points 3 |
| ▲ THTTT: total points 2 | HTTHH: total points 0 |
| THTTH: total points 0 | ● HTHTT: total points 3 |
| ▶ THTHT: total points 2 | ∎ HTHTH: total points 2 |
| − HTHHT: total points 2 |
| − HTHHH: total points 2 |
| • HHTTT: total points 1 |
| ∎ HHTTH: total points 2 |
| HHTHT: total points 0 |
| ▶ HHTHH: total points 1 |
| ▲ HHHTT: total points 1 |
| ▼ HHHTH: total points 2 |
| ◀ HHHHT: total points 2 |
| ▶ HHHHH: total points 4 |

P and Q are seperated



4(c) P=25, Q=25, D=5, t-SNE (P is square,0 in Z)

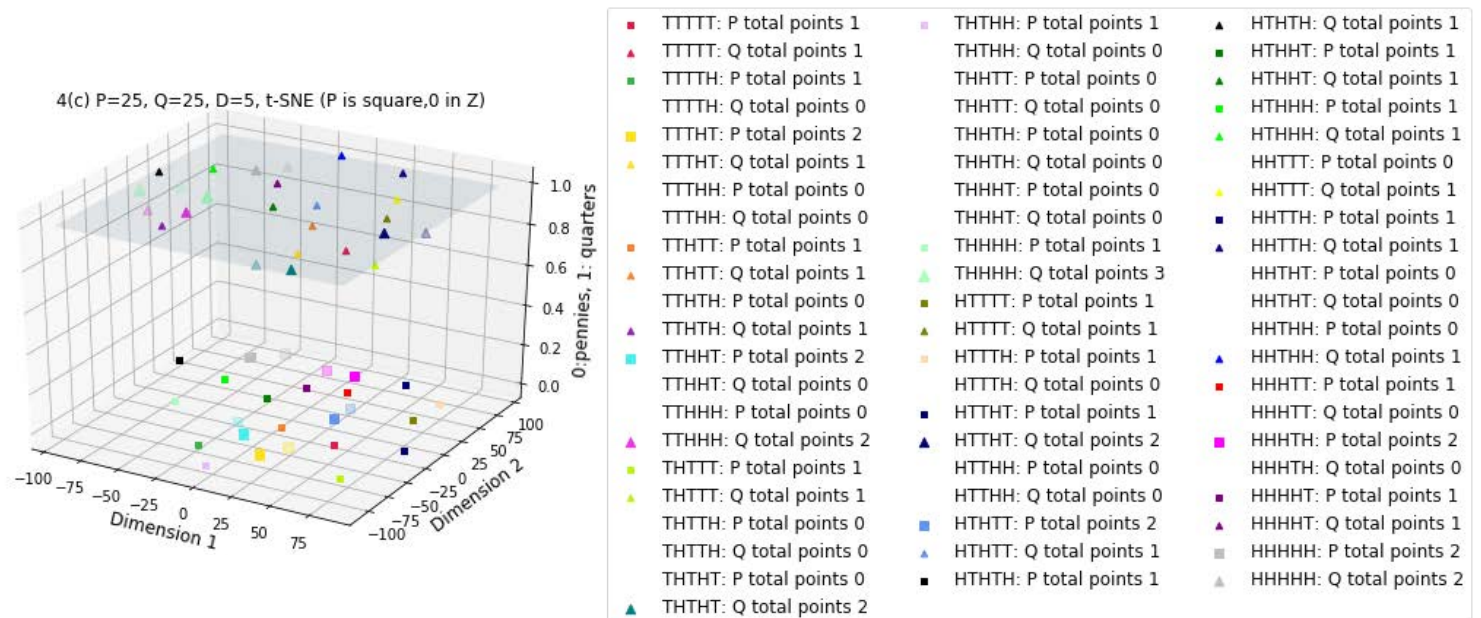| | | |
|---|---|---|
| ■ TTTTT: P total points 1 | ∎ THTHH: P total points 1 | ▲ HTHTH: Q total points 1 |
| ▲ TTTTT: Q total points 1 | THTHH: Q total points 0 | ∎ HTHHT: P total points 1 |
| ∎ TTTTH: P total points 1 | THHTT: P total points 0 | ▲ HTHHT: Q total points 1 |
| TTTTH: Q total points 0 | THHTT: Q total points 0 | ∎ HTHHH: P total points 1 |
| ■ TTTHT: P total points 2 | THHTH: P total points 0 | ▲ HTHHH: Q total points 1 |
| ▲ TTTHT: Q total points 1 | THHTH: Q total points 0 | HHTTT: P total points 0 |
| TTTHH: P total points 0 | THHHT: P total points 0 | • HHTTT: Q total points 1 |
| TTTHH: Q total points 0 | THHHT: Q total points 0 | ∎ HHTTH: P total points 1 |
| ∎ TTHTT: P total points 1 | ∎ THHHH: P total points 1 | ▲ HHTTH: Q total points 1 |
| ▲ TTHTT: Q total points 1 | ▲ THHHH: Q total points 3 | HHTHT: P total points 0 |
| TTHTH: P total points 0 | ∎ HTTTT: P total points 1 | HHTHT: Q total points 0 |
| ▲ TTHTH: Q total points 1 | ▲ HTTTT: Q total points 1 | HHTHH: P total points 0 |
| ∎ TTHHT: P total points 2 | ∎ HTTTH: P total points 1 | ▲ HHTHH: Q total points 1 |
| TTHHT: Q total points 0 | HTTTH: Q total points 0 | ∎ HHHTT: P total points 1 |
| TTHHH: P total points 0 | ∎ HTTHT: P total points 1 | HHHTT: Q total points 0 |
| ▲ TTHHH: Q total points 2 | ▲ HTTHT: Q total points 2 | ∎ HHHTH: P total points 2 |
| • THTTT: P total points 1 | HTTHH: P total points 0 | HHHTH: Q total points 0 |
| ▲ THTTT: Q total points 1 | HTTHH: Q total points 0 | ∎ HHHHT: P total points 1 |
| THTTH: P total points 0 | ∎ HTHTT: P total points 2 | ▲ HHHHT: Q total points 1 |
| THTTH: Q total points 0 | ▲ HTHTT: Q total points 1 | ∎ HHHHH: P total points 2 |
| THTHT: P total points 0 | ∎ HTHTH: P total points 1 | ▲ HHHHH: Q total points 2 |
| ▲ THTHT: Q total points 2 | | |

(d) Simulate sequences of coin flips for $P = 75$ pennies and $Q = 75$ quarters. Visualize the data using t-SNE, clearly denoting the data points that are associated with common coin flip sequences.
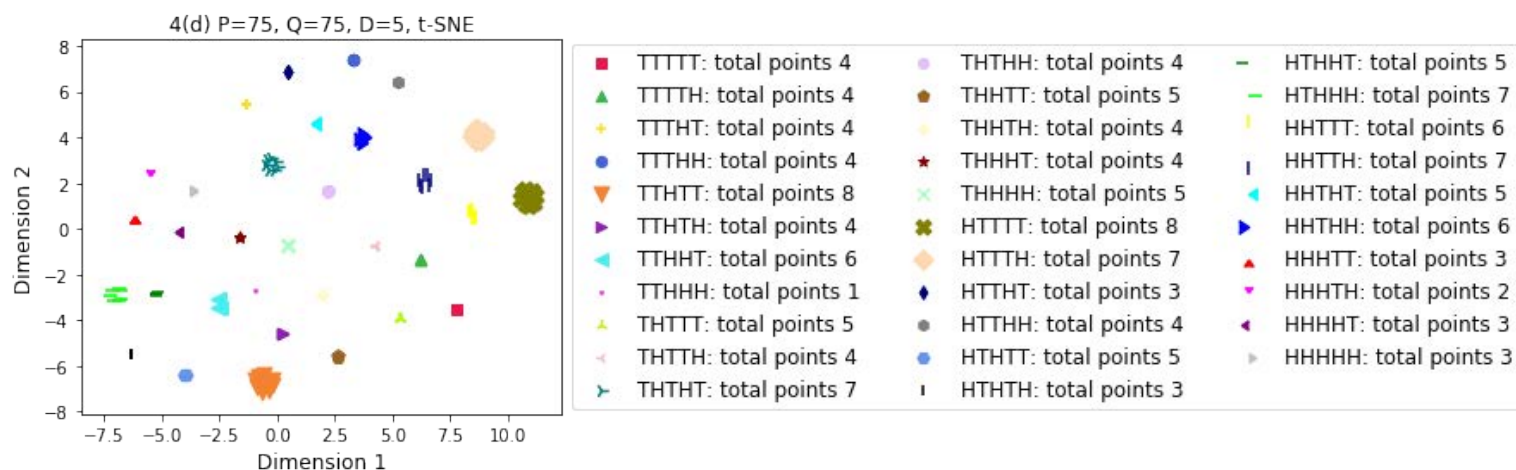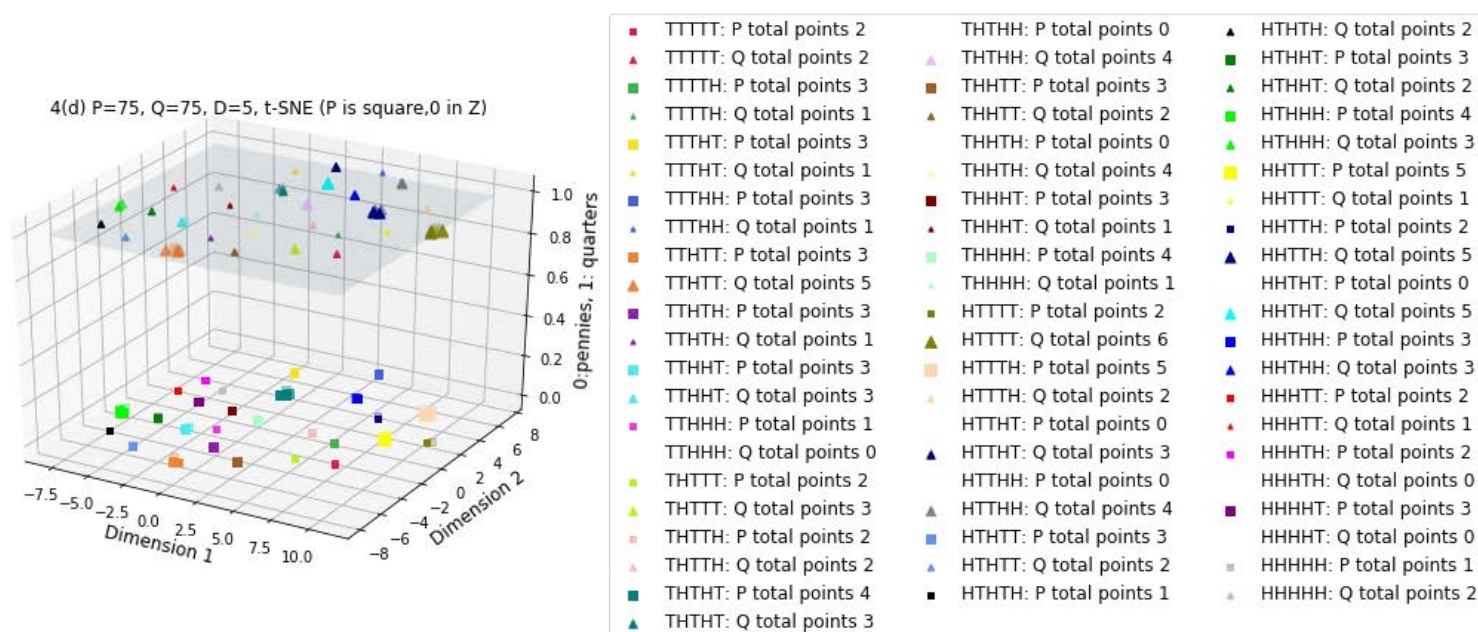
P and Q are not seperated



4(d) P=75, Q=75, D=5, t-SNE

| | | | |
|---|---|---|---|
| ■ TTTTT: total points 4 | ● THTHH: total points 4 | — HTHHT: total points 5 | |
| ▲ TTTTH: total points 4 | ⬟ THHTT: total points 5 | — HTHHH: total points 7 | |
| ✦ TTTHT: total points 4 | THHTH: total points 4 | HHTTT: total points 6 | |
| ● TTTHH: total points 4 | ★ THHHT: total points 4 | ┃ HHTTH: total points 7 | |
| ▼ TTHTT: total points 8 | ✕ THHHH: total points 5 | ◄ HHTHT: total points 5 | |
| ► TTHTH: total points 4 | ✖ HTTTT: total points 8 | ► HHTHH: total points 6 | |
| ◄ TTHHT: total points 6 | ◇ HTTTH: total points 7 | ▲ HHHTT: total points 3 | |
| · TTHHH: total points 1 | ◆ HTTHT: total points 3 | ▼ HHHTH: total points 2 | |
| ⅃ THTTT: total points 5 | ● HTTHH: total points 4 | ◄ HHHHT: total points 3 | |
| ◁ THTTH: total points 4 | ● HTHTT: total points 5 | ▷ HHHHH: total points 3 | |
| ⊢ THTHT: total points 7 | ┃ HTHTH: total points 3 | | |

P and Q are seperated



4(d) P=75, Q=75, D=5, t-SNE (P is square,0 in Z)

| | | |
|---|---|---|
| ■ TTTTT: P total points 2 | THTHH: P total points 0 | ▲ HTHTH: Q total points 2 |
| ▲ TTTTT: Q total points 2 | ▲ THTHH: Q total points 4 | ■ HTHHT: P total points 3 |
| ■ TTTTH: P total points 3 | ■ THHTT: P total points 3 | ▲ HTHHT: Q total points 2 |
| ▲ TTTTH: Q total points 1 | ▲ THHTT: Q total points 2 | ■ HTHHH: P total points 4 |
| ■ TTTHT: P total points 3 | THHTH: P total points 0 | ▲ HTHHH: Q total points 3 |
| ▲ TTTHT: Q total points 1 | THHTH: Q total points 4 | ■ HHTTT: P total points 5 |
| ■ TTTHH: P total points 3 | ■ THHHT: P total points 3 | HHTTT: Q total points 1 |
| ▲ TTTHH: Q total points 1 | ▲ THHHT: Q total points 1 | ■ HHTTH: P total points 2 |
| ■ TTHTT: P total points 3 | ■ THHHH: P total points 4 | ▲ HHTTH: Q total points 5 |
| ▲ TTHTT: Q total points 5 | THHHH: Q total points 1 | HHTHT: P total points 0 |
| ■ TTHTH: P total points 3 | ■ HTTTT: P total points 2 | ▲ HHTHT: Q total points 5 |
| ▲ TTHTH: Q total points 1 | ▲ HTTTT: Q total points 6 | ■ HHTHH: P total points 3 |
| ■ TTHHT: P total points 3 | ■ HTTTH: P total points 5 | ▲ HHTHH: Q total points 3 |
| ▲ TTHHT: Q total points 3 | HTTTH: Q total points 2 | ■ HHHTT: P total points 2 |
| ■ TTHHH: P total points 1 | HTTHT: P total points 0 | ▲ HHHTT: Q total points 1 |
| TTHHH: Q total points 0 | ▲ HTTHT: Q total points 3 | ■ HHHTH: P total points 2 |
| ■ THTTT: P total points 2 | HTTHH: P total points 0 | HHHTH: Q total points 0 |
| ▲ THTTT: Q total points 3 | ▲ HTTHH: Q total points 4 | ■ HHHHT: P total points 3 |
| THTTH: P total points 2 | ■ HTHTT: P total points 3 | HHHHT: Q total points 0 |
| ▲ THTTH: Q total points 2 | ▲ HTHTT: Q total points 2 | HHHHH: P total points 1 |
| ■ THTHT: P total points 4 | ■ HTHTH: P total points 1 | ▲ HHHHH: Q total points 2 |
| ▲ THTHT: Q total points 3 | | |

(30)  5.  (a)  What trend do you see in the clusters produced by t-SNE as the number of observations[3] increases?

Solution:
Under the same number of features, as the number of observations increases, the data points that are associated with common coin flip sequences are more likely to be clustered. Different classes are more likely to be separable and clustered.

(b)  What trend do you see in the clusters produced by t-SNE as the number of features[4] increases?

Solution:
Under the same number of observations, as the number of features increases, the data points that are associated with common coin flip sequences are less likely to be clustered. It becomes difficult to separate and cluster different classes.

(c)  How might changes in the number of observations and features impact the conclusions drawn from the visualizations produced by t-SNE, particularly in the absence of ground truth to help guide interpreting the visualization?[5]

Solution:
As the number of observations increases and the number of features decreases, the data points that are associated with common coin flip sequences are more likely to be clustered. It becomes easier to separate and cluster different classes, vice versa.
This is because when the number of observations increases, the probability of seeing identical(or substantially similar, in the case of continuous variables) data under both classes increases. And when the number of features increases, the probability of seeing identical (or substantially similar, in the case of continuous variables) data under both classes decreases.
In the absence of ground truth, it will make it even harder to distinguish different classes when the number of observations decreases and the number of features increases.

---

[3]In the distinguishing coins problem, the number of coins ($P$ and $Q$) is the number of observations.
[4]In the distinguishing coins problem, the number of flips of each coin ($D$) is the number of features.
[5]Feel free to run Monte Carlo simulations before answering this question, if you think that would be helpful.