

Cluster Analysis of Student Performance

Objective:

The objective of this analysis is to group students in 3 sections based on their academic performance in five subjects using hierarchical clustering techniques.

Step 1:

Data Description:

The dataset includes information about students, including:

- Student ID: Unique identifier for each student.
- Read: Score in the reading subject.
- Write: Score in the writing subject.
- Math: Score in the mathematics subject.
- Science: Score in the science subject.
- Social Studies (Socst): Score in the social studies subject.

Step 2:

Data Preprocessing:

- Checked for missing values and confirmed data completeness.
- Extracted the last five columns (subject scores) from the dataset for analysis.
- Standardized the five columns (subject scores) from the dataset.

$$\text{std_value} = \frac{\text{value} - \text{mean}}{\text{standard deviation}}$$

Key Code Snippet (Data Preprocessing):

```
stdz_read    = (df1['read'] - df1['read'].mean() )/ df1['read'].std()
stdz_write   = (df1['write'] - df1['write'].mean() )/ df1['write'].std()
stdz_math    = (df1['math'] - df1['math'].mean() )/ df1['math'].std()
stdz_science = (df1['science'] - df1['science'].mean() )/ df1['science'].std()
stdz_socst   = (df1['socst'] - df1['socst'].mean() )/ df1['socst'].std()
```

Step 3:

Distance Calculation:

- Calculated pairwise distances between standardized values of subject marks using the Euclidean, Mankowski (P=3), Mankowski(P=4) and Manhattan distance metric.
- Constructed a distance matrix to represent the dissimilarity between students.

Key Code Snippet (Distance Calculation):

```
from scipy.spatial.distance import pdist, squareform

euclidean_dist = pd.DataFrame(squareform(pdist(d, metric='euclidean')))
minkowski_3_dist = pd.DataFrame(squareform(pdist(d, metric='minkowski', p = 3)))
minkowski_4_dist = pd.DataFrame(squareform(pdist(d, metric='minkowski', p = 4)))
cityblock_dist = pd.DataFrame(squareform(pdist(d, metric='cityblock')))
```

Step 4:

Clustering Technique:

- Applied hierarchical clustering using the Average method to minimize variance within clusters.
- Visualized clusters using dendrograms to determine optimal clusters.

Key Code Snippet (Clustering):

```
from scipy.cluster.hierarchy import linkage
from scipy.cluster.hierarchy import dendrogram

avc_link = linkage(student_sections[['Min_Distance']].fillna(0), method='average')
dend = dendrogram(avc_link)
```

Step 5:

Cluster Assignment:

- Divided data into clusters using the fcluster method based on the dendrogram analysis.

Key Code Snippet (F cluster):

```
from scipy.cluster.hierarchy import linkage, fcluster

# Perform hierarchical clustering
clusters = fcluster(avc_link, 3, criterion="maxclust")
student_sections["clusters"] = clusters
student_sections
```

	Student 1	Section 1	Min_Distance	Min_Distance_Type	clusters
0	1	B	0.489820	Minkowski (p=4)	1
1	2	C	1.319005	Minkowski (p=4)	2
2	3	C	0.547103	Minkowski (p=4)	1
3	4	A	0.546618	Minkowski (p=4)	1
4	5	A	0.792446	Minkowski (p=4)	1

Step 6:

Student Section Assignment:

- Listed each student along with their assigned section.
- Provided the count of students in each section.

Key Code Snippet (Section Assignment):

```
freq_table=pd.crosstab(student_sections['Section 1'],'count')
freq_table
```

col_0	count
Section 1	
A	71
B	63
C	65

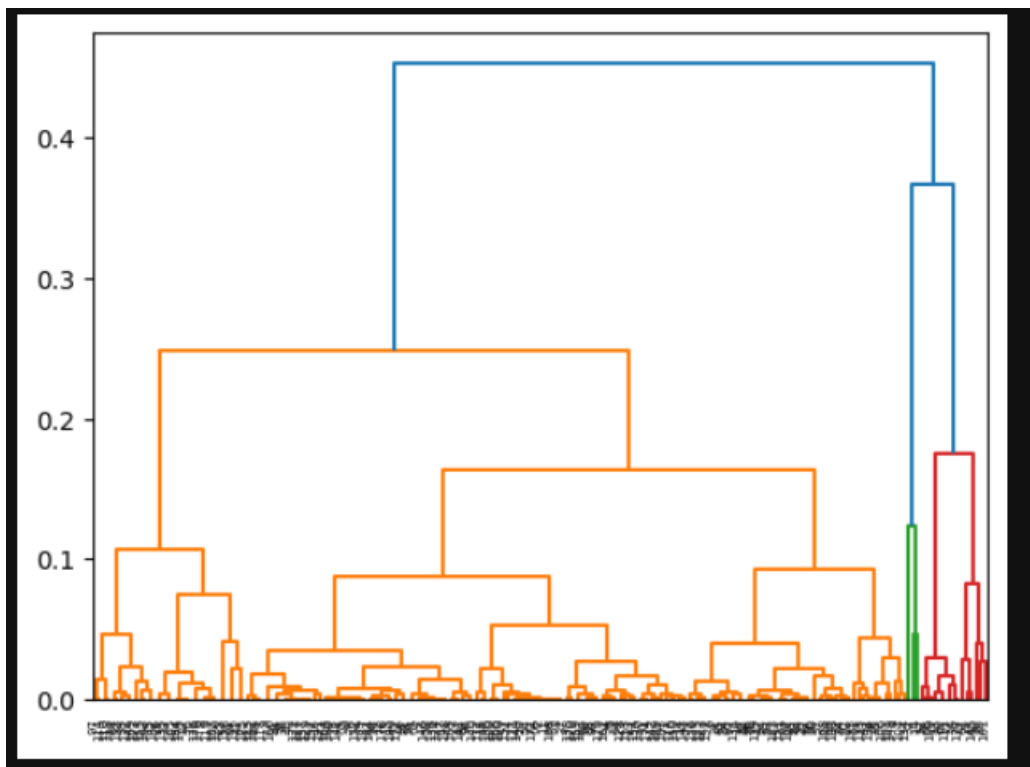
```
# Create a DataFrame for each student with their assign section
student_dataframe = student_sections[['Student 1', 'Section 1']]
student_dataframe
```

	Student 1	Section 1
0	1	B
1	2	C
2	3	C
3	4	A
4	5	A

Step 7:

Visualization:

- Generated dendrograms to visualize hierarchical clustering results.



Results:

1. Cluster Distribution:

- The clustering revealed that students were grouped into distinct clusters based on performance metrics.
- The uniform distribution of clusters indicated that students with similar academic profiles were grouped together effectively.

2. Section-Wise Distribution:

- The students were divided into three sections based on the cluster assignment:
 - **Section A: 71 students**
 - **Section B: 63 students**
 - **Section C: 65 students**
- The distribution shows a relatively balanced number of students across the three sections.

3. Dendrogram Interpretation:

- The dendrogram displays how students are progressively grouped based on academic performance, with vertical lines representing the dissimilarity between clusters.
- By cutting the dendrogram at an appropriate level, three distinct clusters (sections A, B, and C) were identified, effectively grouping students with similar performance levels.

Conclusion:

The cluster analysis successfully grouped students based on their academic performance, revealing patterns that can aid in personalized academic interventions and support. The balanced section distribution highlights the effectiveness of the clustering methodology.

```
[51]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import itertools
import seaborn as sns
from scipy.spatial.distance import pdist, squareform
from scipy.cluster.hierarchy import linkage, fcluster, dendrogram
```

```
[52]: raw_data = pd.read_csv("D:/hsb2.csv")
raw_data.head()
```

```
[52]:
```

	id	Gender	race	ses	schtyp	prog	read	write	math	science	socst
0	70	0	4	1	1	1	57	52	41	47	57
1	121	1	4	2	1	3	68	59	53	63	61
2	86	0	4	3	1	1	44	33	54	58	31
3	141	0	4	3	1	3	63	44	47	53	56
4	172	0	4	2	1	2	47	52	57	53	61

```
[53]: raw_data.shape # checking the number of rows and columns
```

```
[53]: (200, 11)
```

```
[54]: raw_data.isnull().sum() # checking for missing values
```

```
[54]: id          0
Gender         0
race          0
ses           0
schtyp        0
prog          0
read          0
write         0
math          0
science       0
socst         0
dtype: int64
```

```
[55]: df1 = raw_data.iloc[:, -5:]
df1.head(2)
```

```
[55]:
```

	read	write	math	science	socst
0	57	52	41	47	57
1	68	59	53	63	61

```
[56]: stdz_read  = (df1['read'] - df1['read'].mean() )/ df1['read'].std()
stdz_write   = (df1['write'] - df1['write'].mean() )/ df1['write'].std()
stdz_math    = (df1['math'] - df1['math'].mean() )/ df1['math'].std()
stdz_science = (df1['science'] - df1['science'].mean() )/ df1['science'].std()
stdz_socst   = (df1['socst'] - df1['socst'].mean() )/ df1['socst'].std()
```

```
[57]: dict = {
    'Student ID' : raw_data['id'],
    'std_read' : stdz_read,
    'std_write' : stdz_write,
    'std_math' : stdz_math,
    'std_science' : stdz_science,
    'std_socst' : stdz_socst
}
df2 = pd.DataFrame(dict)
df2.head(2)
```

```
[57]:
```

	Student ID	std_read	std_write	std_math	std_science	std_socst
0	70	0.465233	-0.081763	-1.243002	-0.489855	0.428007
1	121	1.538096	0.656744	0.037893	1.126161	0.800593

```
[58]: d = pd.DataFrame(df2.iloc[:, 1:].values)
d
```

```
[58]:
```

	0	1	2	3	4
0	0.465233	-0.081763	-1.243002	-0.489855	0.428007
1	1.538096	0.656744	0.037893	1.126161	0.800593
2	-0.802697	-2.086282	0.144634	0.621156	-1.993798

```
[59]: euclidean_dist = pd.DataFrame(squareform(pdist(d, metric='euclidean')))
      euclidean_dist
```

```
[59]:
```

	0	1	2	3	4	5	6	7	8	9	...	190	191	192	193	194	195
0	0.000000	2.467276	3.827631	1.356797	2.091440	2.344641	1.236516	3.177154	2.023103	1.372213	...	2.158941	2.607236	1.663841	2.242406	2.654417	1.491907
1	2.467276	0.000000	4.591137	2.095082	2.437808	2.463793	2.341126	4.990537	1.190121	1.982964	...	4.039649	4.044584	2.120681	4.214288	2.663643	2.516878
2	3.827631	4.591137	0.000000	3.319163	3.502842	3.490600	4.191776	2.766721	3.649071	3.341217	...	3.217336	3.190791	3.144266	4.312411	4.238974	4.093282
3	1.356797	2.095082	3.319163	0.000000	2.122233	2.359186	2.147963	3.682418	1.706246	1.510894	...	2.984094	3.338354	1.664624	3.360488	3.059596	2.152290
4	2.091440	2.437808	3.502842	2.122233	0.000000	1.231222	1.787341	3.327330	1.984618	1.515164	...	2.811400	2.748670	0.725000	2.828726	1.514037	1.700088
...
195	1.491907	2.516878	4.093282	2.152290	1.700088	2.534454	1.681796	3.193234	1.878117	1.042075	...	2.284009	2.687565	1.587012	1.968157	1.546474	0.000000
196	2.219857	4.501560	3.424240	2.996378	3.106757	3.428282	2.751811	1.458266	3.689574	2.737048	...	1.182811	2.021851	2.850149	1.526620	3.554929	2.637560
197	2.265691	2.509035	2.518049	1.327262	1.744028	2.172615	2.716533	3.432753	1.842521	1.652307	...	3.116411	3.303473	1.460831	3.621119	2.803356	2.406654
198	2.411854	1.499230	4.301472	2.462855	1.410064	1.754253	1.873448	4.251507	1.391760	1.585597	...	3.419763	3.279101	1.394425	3.409541	1.274498	1.827228
199	3.048334	1.816405	4.926860	2.969342	2.246217	2.927185	2.834715	4.898759	1.792663	2.083380	...	4.099938	4.151708	2.321012	3.947450	1.740347	2.094348

200 rows x 200 columns

```
[60]: # Extract all unique combinations of indices (i, j) where i < j
      index_pairs = list(itertools.combinations(range(len(euclidean_dist)), 2))

      # Convert combinations to DataFrame
      distances = [{"ID Pair": (i, j), "Distance": euclidean_dist.iloc[i, j]} for i, j in index_pairs]
      distance_df = pd.DataFrame(distances)

      # Find the pair with the smallest distance
      min_distance_pair = distance_df.loc[distance_df["Distance"].idxmin()]

      # Display all combinations and the smallest one
      print(distance_df) # Displays pairs
      print("\nSmallest Distance Pair:\n", min_distance_pair)
```

```
[60]: # Extract all unique combinations of indices (i, j) where i < j
      index_pairs = list(itertools.combinations(range(len(euclidean_dist)), 2))

      # Convert combinations to DataFrame
      distances = [{"ID Pair": (i, j), "Distance": euclidean_dist.iloc[i, j]} for i, j in index_pairs]
      distance_df = pd.DataFrame(distances)

      # Find the pair with the smallest distance
      min_distance_pair = distance_df.loc[distance_df["Distance"].idxmin()]

      # Display all combinations and the smallest one
      print(distance_df) # Displays pairs
      print("\nSmallest Distance Pair:\n", min_distance_pair)
```

	ID Pair	Distance
0	(0, 1)	2.467276
1	(0, 2)	3.827631
2	(0, 3)	1.356797
3	(0, 4)	2.091440
4	(0, 5)	2.344641
...
19895	(196, 198)	3.987989
19896	(196, 199)	4.628121
19897	(197, 198)	2.398450
19898	(197, 199)	2.868166
19899	(198, 199)	1.233818

[19900 rows x 2 columns]

Smallest Distance Pair:
ID Pair (26, 159)
Distance 0.255789
Name: 4981, dtype: object

```
[13]: minkowski_3_dist = pd.DataFrame(squareform(pdist(d, metric='minkowski', p = 3)))
minkowski_3_dist.head(2)
```

```
[13]:
```

	0	1	2	3	4	5	6	7	8	9	...	190	191	192	193	194	195	
0	0.000000	2.000933	3.04884	1.087936	1.835463	1.959738	0.998833	2.697663	1.658883	1.228224	...	1.873856	2.294405	1.433889	1.879792	2.281648	1.292368	1.8
1	2.000933	0.000000	3.81250	1.762302	2.161829	2.365623	2.004988	4.052225	1.020746	1.636539	...	3.257382	3.441455	1.860714	3.578305	2.265511	2.269058	3.5

2 rows × 200 columns

```
[33]: # Extract all unique combinations of indices (i, j) where i < j
index_pairs_3 = list(itertools.combinations(range(len(minkowski_3_dist)), 2))

# Convert combinations to DataFrame
distances_3 = [{"ID Pair": (i, j), "Distance": minkowski_3_dist.iloc[i, j]} for i, j in index_pairs_3]
distance_df_3 = pd.DataFrame(distances_3)

# Find the pair with the smallest distance
min_distance_pair_3 = distance_df_3.loc[distance_df_3["Distance"].idxmin()]

# Display all combinations and the smallest one
print(distance_df_3) # Displays pairs
print("\nSmallest Distance Pair:\n", min_distance_pair_3)
```

```

      ID Pair  Distance
0      (0, 1)  2.000933
1      (0, 2)  3.048840
2      (0, 3)  1.087936
3      (0, 4)  1.835463
4      (0, 5)  1.959738
...      ...      ...
19895  (196, 198)  3.116737
19896  (196, 199)  3.643663
19897  (197, 198)  2.257215
19898  (197, 199)  2.603823
19899  (198, 199)  1.017260

[19900 rows x 2 columns]
```

```
[34]: minkowski_4_dist = pd.DataFrame(squareform(pdist(d, metric='minkowski', p = 4)))
minkowski_4_dist.head(2)
```

```
[34]:
```

	0	1	2	3	4	5	6	7	8	9	...	190	191	192	193	194	195
0	0.000000	0.924452	0.441724	0.650402	0.834382	0.768606	1.003978	0.784186	0.735969	0.647368	...	0.588528	0.797580	0.532084	0.470814	0.412349	0.719878
1	0.924452	0.000000	0.943114	0.468729	0.631342	0.776780	0.829162	0.856780	0.866970	0.631731	...	0.804613	0.854799	0.720708	0.740210	0.972702	0.850895

2 rows × 200 columns

```
[16]: # Extract all unique combinations of indices (i, j) where i < j
index_pairs_4 = list(itertools.combinations(range(len(minkowski_4_dist)), 2))

# Convert combinations to DataFrame
distances_4 = [{"ID Pair": (i, j), "Distance": minkowski_4_dist.iloc[i, j]} for i, j in index_pairs_4]
distance_df_4 = pd.DataFrame(distances_4)

# Find the pair with the smallest distance
min_distance_pair_4 = distance_df_4.loc[distance_df_4["Distance"].idxmin()]

# Display all combinations and the smallest one
print(distance_df_4) # Displays pairs
print("\nSmallest Distance Pair:\n", min_distance_pair_4)
```

```

      ID Pair  Distance
0      (0, 1)  1.827479
1      (0, 2)  2.763950
2      (0, 3)  0.981467
3      (0, 4)  1.758695
4      (0, 5)  1.809540
...      ...      ...
19895  (196, 198)  2.779240
19896  (196, 199)  3.276402
19897  (197, 198)  2.227015
19898  (197, 199)  2.549469
19899  (198, 199)  0.932897

[19900 rows x 2 columns]
```



```
[35]: cityblock_dist = pd.DataFrame(squareform(pdist(d, metric='cityblock')))
cityblock_dist.head(2)
```

	0	1	2	3	4	5	6	7	8	9	...	190	191	192	193	194	195
0	0.000000	1.799869	0.885726	1.346769	2.020730	1.495248	2.839978	2.083988	1.819068	1.219695	...	1.646782	2.139813	1.467986	1.037694	0.812147	1.691597
1	1.799869	0.000000	1.995993	1.201002	1.598056	1.793036	1.970445	2.018475	2.060335	1.683886	...	2.344879	1.971981	1.870290	1.696427	2.583091	1.673993

2 rows x 200 columns

```
[46]: # Extract all unique combinations of indices (i, j) where i < j
index_pairs_CB = list(itertools.combinations(range(len(cityblock_dist)), 2))

# Convert combinations to DataFrame
distances_CB = [{"ID Pair": (i, j), "Distance": cityblock_dist.iloc[i, j]} for i, j in index_pairs_CB]
distance_df_CB = pd.DataFrame(distances_CB)

# Find the pair with the smallest distance
min_distance_pair_CB = distance_df_CB.loc[distance_df_CB["Distance"].idxmin()]

# Display all combinations and the smallest one
print(distance_df_CB) # Displays pairs
print("\nSmallest Distance Pair:\n", min_distance_pair_CB)
```

	ID Pair	Distance
0	(0, 1)	1.799869
1	(0, 2)	0.885726
2	(0, 3)	1.346769
3	(0, 4)	2.020730
4	(0, 5)	1.495248
...
19895	(196, 198)	1.589242
19896	(196, 199)	1.576566
19897	(197, 198)	1.274858
19898	(197, 199)	2.007761
19899	(198, 199)	1.510816

[19900 rows x 2 columns]

```
[19]: np.random.seed(42) # For reproducibility
d = np.random.rand(200, 5) # Simulating 200 students with 5 feature attributes

# Sample section mapping (Assigning random sections: A, B, C)
sections = {i: np.random.choice(['A', 'B', 'C']) for i in range(len(d))}
```

```
[36]: # Creating a distance Matrix DataFrame from a Dictionary
dist_matrix = {
    "ID Pair": distance_df["ID Pair"],
    "Student 1": [i for i, j in distance_df["ID Pair"]],
    "Student 2": [j for i, j in distance_df["ID Pair"]],
    "Section 1": [sections[i] for i, j in distance_df["ID Pair"]],
    "Section 2": [sections[j] for i, j in distance_df["ID Pair"]],
    "Euclidean": distance_df["Distance"],
    "Minkowski (p=3)": distance_df_3["Distance"],
    "Minkowski (p=4)": distance_df_4["Distance"],
    "Manhattan/Cityblock": distance_df_CB["Distance"]
}
data = pd.DataFrame(dist_matrix)
data.head(2)
```

	ID Pair	Student 1	Student 2	Section 1	Section 2	Euclidean	Minkowski (p=3)	Minkowski (p=4)	Manhattan/Cityblock
0	(0, 1)	0	1	C	B	2.467276	2.000933	1.827479	5.080867
1	(0, 2)	0	2	C	C	3.827631	3.048840	2.763950	8.192901

```
[55]: # Define Distance Columns
distance_columns = ["Euclidean", "Minkowski (p=3)", "Minkowski (p=4)", "Manhattan/Cityblock"]

# Find Minimum Distance for Each Row
data["Min_Distance"] = data[distance_columns].min(axis=1)

# Find Which Distance Metric is the Minimum
data["Min_Distance_Type"] = data[distance_columns].idxmin(axis=1)

# Display Results
print(data[["ID Pair", "Student 1", "Student 2", "Min_Distance", "Min_Distance_Type"]].head())
```

	ID Pair	Student 1	Student 2	Min_Distance	Min_Distance_Type
0	(0, 1)	0	1	1.827479	Minkowski (p=4)
1	(0, 2)	0	2	2.763950	Minkowski (p=4)
2	(0, 3)	0	3	0.981467	Minkowski (p=4)
3	(0, 4)	0	4	1.758695	Minkowski (p=4)
4	(0, 5)	0	5	1.809540	Minkowski (p=4)

```
[56]: data.head(1)
```

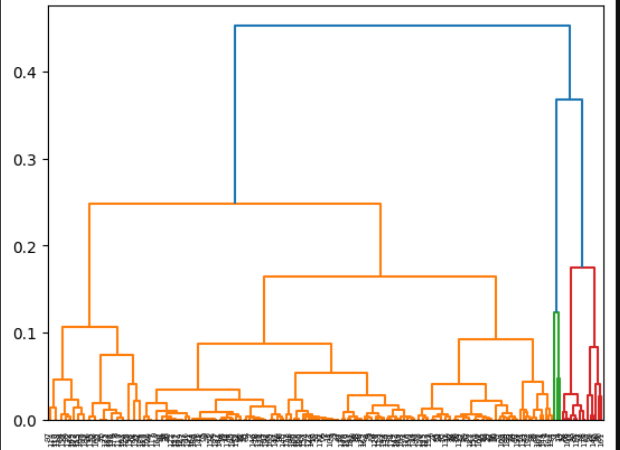
	ID Pair	Student 1	Student 2	Section 1	Section 2	Euclidean	Minkowski (p=3)	Minkowski (p=4)	Manhattan/Cityblock	Min_Distance	Min_Distance_Type
0	(0, 1)	0	1	C	B	2.467276	2.000933	1.827479	5.080867	1.827479	Minkowski (p=4)

```
[57]: # Extract student IDs and their sections
student_sections = pd.concat([
    data[['Student 1', 'Section 1', 'Min_Distance', 'Min_Distance_Type']],
    data[['Student 2', 'Section 2', 'Min_Distance', 'Min_Distance_Type']].rename(columns={'Student 2': 'Student 1',
                                                                                          'Section 2': 'Section 1'})
], ignore_index=True)

[58]: # Remove duplicates and keep only the minimum distance per student
student_sections = student_sections.groupby('Student 1', as_index=False).agg({'Section 1': 'first', 'Min_Distance': 'min',
                                                                              'Min_Distance_Type': 'first'})

# Ensure student IDs remain 1-200
student_sections = student_sections[student_sections['Student 1'].between(1, 200)].sort_values(by='Student 1').reset_index(drop=True)
student_sections.head(2)
```

```
[26]: avc_link = linkage(student_sections[['Min_Distance']].fillna(0), method='average')
dend = dendrogram(avc_link)
```



```
[41]: # Perform hierarchical clustering
clusters = fcluster(avc_link, 3, criterion="maxclust")
student_sections["clusters"] = clusters
student_sections.head(2)
```

	Student 1	Section 1	Min_Distance	Min_Distance_Type	clusters
0	1	B	0.489820	Minkowski (p=4)	1
1	2	C	1.319005	Minkowski (p=4)	2

```
[42]: student_sections.head(2)
```

	Student 1	Section 1	Min_Distance	Min_Distance_Type	clusters
0	1	B	0.489820	Minkowski (p=4)	1
1	2	C	1.319005	Minkowski (p=4)	2

```
[43]: # Create a DataFrame for each student with their assign section
student_dataframe = student_sections[['Student 1', 'Section 1']]
student_dataframe.head(2)
```

```
[43]:
```

	Student 1	Section 1
0	1	B
1	2	C

```
[30]: freq_table=pd.crosstab(student_sections['Section 1'],'count')
freq_table
```

```
[30]:
```

	col_0	count
Section 1		
A		71
B		63
C		65