



Data Engineering

Project (Student Health & Data Processing)

Students Name:

Layan Al-Joudi 444000810 (2)

Table of Contents

1. Phase 1: Relational Database (SQL)	3
2. Phase 2: NoSQL Database (TinyDB)	4
3. Phase 3: Stream Processing (PySpark)	4
4. Diagrams and Screenshots.....	5
1. Phase 1	5
2. Phase 2	6
3. Phase 3	7
4. Dashboard	8
5. Data From Excel.....	10
6. ER Model: Chen's Notation.....	11
7. Normalisation	11
5. Lessons Learned	12
6. Project Demo	12

Table of Figures

<i>Figure 1 Table creation using SQLite</i>	<i>5</i>
<i>Figure 2 CRUD operation Update and Delete example.....</i>	<i>5</i>
<i>Figure 3 Student with poor sleep quality.....</i>	<i>5</i>
<i>Figure 4 Sample of JSON document after transformation</i>	<i>6</i>
<i>Figure 5 TinyDB query result Count students with risk = Moderate.....</i>	<i>6</i>
<i>Figure 6 Sample student records</i>	<i>7</i>
<i>Figure 7 Filtered students with high stress.....</i>	<i>7</i>
<i>Figure 8 Phase1 Dashboard.....</i>	<i>8</i>
<i>Figure 9 Phase2 Dashboard.....</i>	<i>9</i>
<i>Figure 10 Phase3 Dashboard.....</i>	<i>9</i>
<i>Figure 11 Excel data before normalisation</i>	<i>10</i>
<i>Figure 12 Chen's Notation.....</i>	<i>11</i>
<i>Figure 13 Table after normalisation with primary and foreign key.....</i>	<i>11</i>

1. Phase 1: Relational Database (SQL)

Design Overview:

We designed a normalized relational schema consisting of five connected tables:

- **Students:** ID, Age, Gender
- **Health_Metrics:** Heart rate, blood pressure
- **Stress_Levels:** Biosensor and self-reported stress
- **Lifestyle:** Physical activity, sleep quality, mood
- **Study_And_Project:** Study hours, project hours, risk level

We used Python and SQLite to:

- Create the schema and insert data from a CSV file
- Perform CRUD operations (Insert, Read, Update, Partial Delete)
- Apply indexes on common fields (e.g., Mood, Study_Hours)
- Execute optimized SQL queries (e.g., top 5 heart rates, students with poor sleep)

2. Phase 2: NoSQL Database (TinyDB)

Design Overview:

To demonstrate flexibility in data modeling, we converted part of our dataset (Student info + Study_And_Project) into a JSON format. We used **TinyDB**, a lightweight NoSQL document database.

What we did:

- Converted selected fields into a list of JSON documents
- Stored the data in `tinydb_students.json`
- Ran NoSQL queries to extract meaningful information

3. Phase 3: Stream Processing (PySpark)

Design Overview:

We simulated real-time data stream processing using PySpark, loading the same dataset (`student_health_data.csv`) as a simulated stream.

Processing Logic:

- **Filter 1:** Students with Systolic BP ≥ 130
- **Filter 2:** Students with Biosensor Stress > 7
- **Filter 3:** Students studying more than 40 hours

4. Diagrams and Screenshots

1. Phase 1

```
# Create tables
cursor.execute('''
CREATE TABLE IF NOT EXISTS Students (
    Student_ID INTEGER PRIMARY KEY AUTOINCREMENT,
    Age INTEGER,
    Gender TEXT
)
''')

cursor.execute('''
CREATE TABLE IF NOT EXISTS Health_Metrics (
    Health_Metrics_ID INTEGER PRIMARY KEY AUTOINCREMENT,
    Student_ID INTEGER,
    Heart_Rate REAL,
    Blood_Pressure_Systolic REAL,
    Blood_Pressure_Diastolic REAL,
    FOREIGN KEY (Student_ID) REFERENCES Students(Student_ID)
)
''')

cursor.execute('''
CREATE TABLE IF NOT EXISTS Stress_Levels (
    Stress_Levels_ID INTEGER PRIMARY KEY AUTOINCREMENT,
    Student_ID INTEGER,
    Stress_Level_Biosensor REAL,
    Stress_Level_Self_Report REAL,
    FOREIGN KEY (Student_ID) REFERENCES Students(Student_ID)
)
''')
```

Figure 1 Table creation using SQLite

```
Lifestyle before update: (40, 40, 'Moderate', 'Moderate', 'Stressed')
Lifestyle after update (changed Physical_Activity, Sleep_Quality, Mood): (40, 40, 'Moderate', 'Good', 'Calm')
Health_Metrics before update: (25, 25, 62.947892636702186, 160.1716481045759, 78.58198375003913)
Health_Metrics after update (changed Heart_Rate, Blood_Pressure_Systolic): (25, 25, 82.0, 125.0, 78.58198375003913)
Before partial delete: (205, 205, 64.97744283805876, 132.65058933235588, 84.03300390245282)
After partial delete (cleared Blood_Pressure_Diastolic): (205, 205, 64.97744283805876, 132.65058933235588, None)
Study_And_Project before delete: (700, 700, 30.04910558563553, 15.720909957544643, 'Moderate')
Study_And_Project after delete (cleared Project_Hours, Health_Risk_Level): (700, 700, 30.04910558563553, None, None)
✅ CRUD operations completed.
```

Figure 2 CRUD operation Update and Delete example

```
Query 2: Students with 'Poor' sleep quality
(4, 'Poor')
(5, 'Poor')
(16, 'Poor')
(25, 'Poor')
(29, 'Poor')
(30, 'Poor')
(33, 'Poor')
(43, 'Poor')
(49, 'Poor')
(53, 'Poor')
(59, 'Poor')
(61, 'Poor')
(71, 'Poor')
```

Figure 3 Student with poor sleep quality

2. Phase 2

```
{
  "Student_ID": 1,
  "Age": 24,
  "Gender": "M",
  "Study_And_Project": {
    "Study_Hours": 34.520972884506875,
    "Project_Hours": 16.80095639050803,
    "Health_Risk_Level": "Moderate"
  }
},
{
  "Student_ID": 2,
  "Age": 21,
  "Gender": "F",
  "Study_And_Project": {
    "Study_Hours": 16.763846015109607,
    "Project_Hours": 15.79115434826643,
    "Health_Risk_Level": "Moderate"
  }
},
{
```

Figure 4 Sample of JSON document after transformation

Number of students with moderate health risk: 672

Figure 5 TinyDB query result Count students with risk = Moderate

3. Phase 3

Student_ID	Age	Gender	Blood_Pressure_Systolic	Stress_Level_Biosensor	Study_Hours	Health_Risk_Level
1	24	M	135.0	3.13	34.52	Moderate
2	21	F	110.77	8.2	16.76	High
3	22	M	109.37	6.78	44.2	Moderate
4	24	M	125.14	9.0	21.77	High
5	20	M	140.51	7.26	41.96	High

Figure 6 Sample student records

Students with high stress (biosensor > 7):						
Student_ID	Stress_Level_Biosensor					
2	8.2					
4	9.0					
5	7.26					

Figure 7 Filtered students with high stress

4. Dashboard

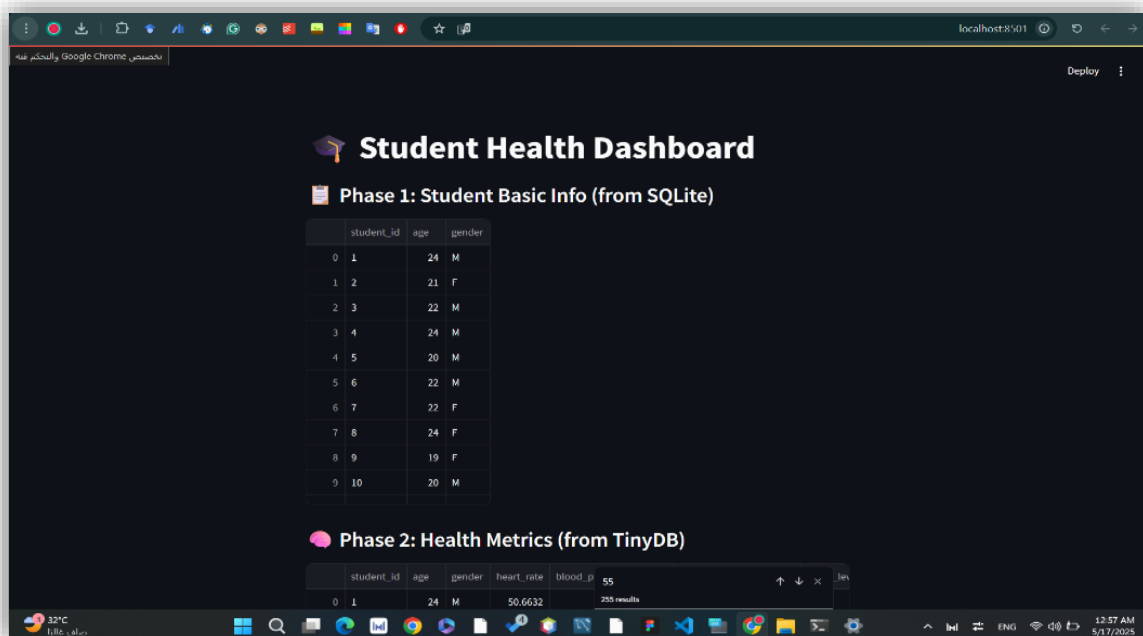
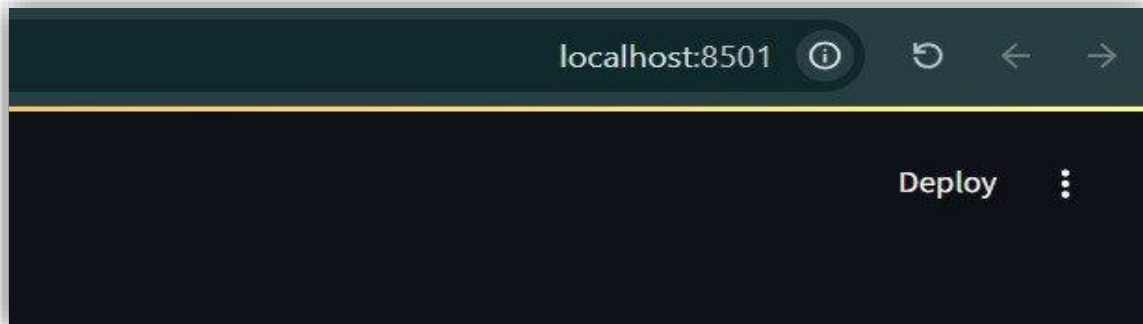
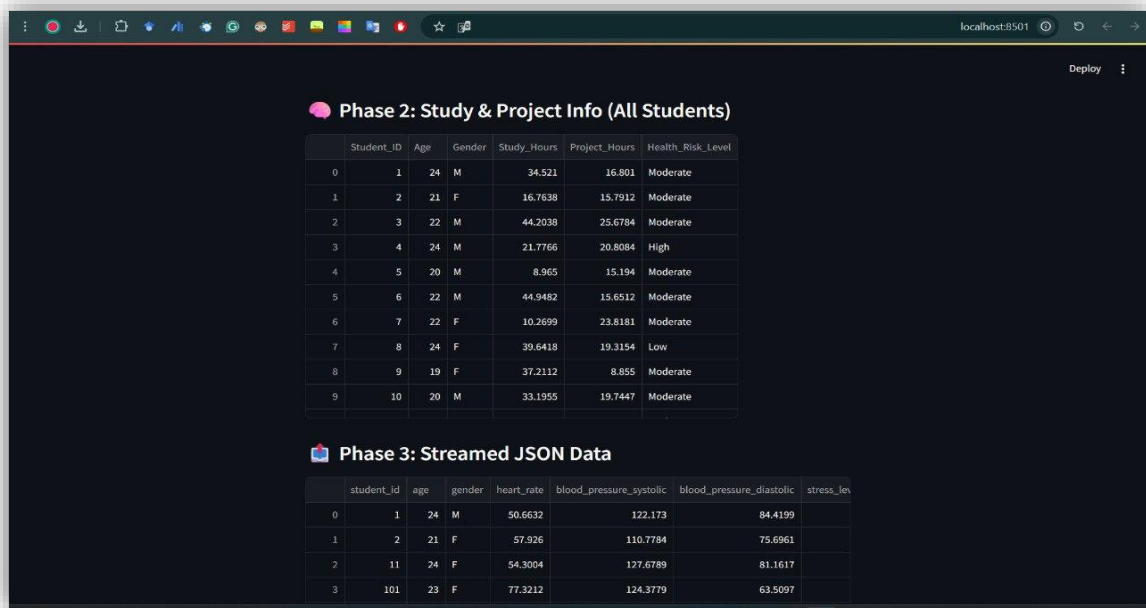


Figure 8 Phase1 Dashboard



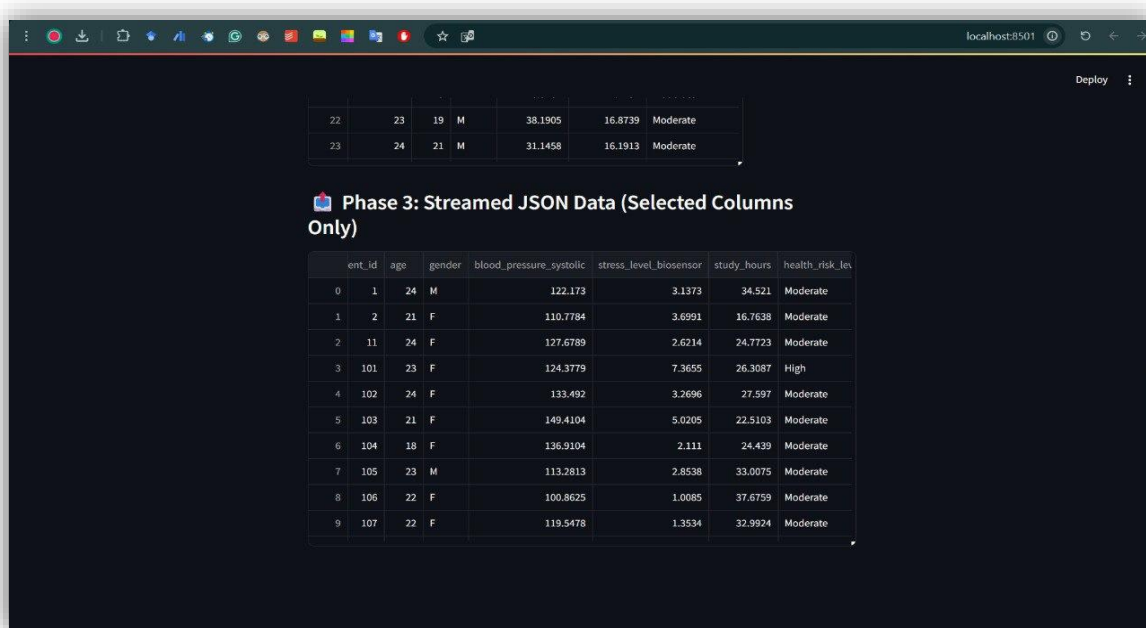
Phase 2: Study & Project Info (All Students)

	Student_ID	Age	Gender	Study_Hours	Project_Hours	Health_Risk_Level
0	1	24	M	34.521	16.801	Moderate
1	2	21	F	16.7638	15.7912	Moderate
2	3	22	M	44.2038	25.6784	Moderate
3	4	24	M	21.7766	20.8084	High
4	5	20	M	8.965	15.194	Moderate
5	6	22	M	44.9482	15.6512	Moderate
6	7	22	F	10.2699	23.8181	Moderate
7	8	24	F	39.6418	19.3154	Low
8	9	19	F	37.2112	8.855	Moderate
9	10	20	M	33.1955	19.7447	Moderate

Phase 3: Streamed JSON Data

	student_id	age	gender	heart_rate	blood_pressure_systolic	blood_pressure_diastolic	stress_level
0	1	24	M	50.6632	122.173	84.4199	
1	2	21	F	57.926	110.7784	75.6961	
2	11	24	F	54.3004	127.6789	81.1617	
3	101	23	F	77.3212	124.3779	63.5097	

Figure 9 Phase2 Dashboard



Phase 3: Streamed JSON Data (Selected Columns Only)

	ent_id	age	gender	blood_pressure_systolic	stress_level_biosensor	study_hours	health_risk_level
0	1	24	M	122.173	3.1373	34.521	Moderate
1	2	21	F	110.7784	3.6991	16.7638	Moderate
2	11	24	F	127.6789	2.6214	24.7723	Moderate
3	101	23	F	124.3779	7.3655	26.3087	High
4	102	24	F	133.492	3.2696	27.597	Moderate
5	103	21	F	149.4104	5.0205	22.5103	Moderate
6	104	18	F	136.9104	2.111	24.439	Moderate
7	105	23	M	113.2813	2.8538	33.0075	Moderate
8	106	22	F	100.8625	1.0085	37.6759	Moderate
9	107	22	F	119.5478	1.3534	32.9924	Moderate

Figure 10 Phase3 Dashboard

5. Data From Excel

Student_ID	Age	Gender	Heart_Rate	Blood_Pressure_Systolic	Blood_Pressure_Diastolic	Stress_Level_Biosensor	Stress_Level_Self_Report	Physical_Activity	Sleep_Quality	Mood	Study_Hours	Project_Hours	Health_Risk_Level
1	24	M	50.663216988244	122.1730149588300	84.41986006395030	3.137349746524460	9.028669159842590	High	Moderate	Happy	34.520972884506900	16.80055639050800	Moderate
2	21	F	57.926041683536300	110.77840702264600	75.69614523602980	3.690978336950830	5.819697244773670	Moderate	Good	Stressed	16.763846015109600	15.79115438286400	Moderate
3	22	M	59.29421920322830	109.37567306957500	83.80381448216310	6.785155634556780	5.892360174201630	Low	Moderate	Happy	44.20379848028350	25.67843704706690	Moderate
4	24	M	78.82623236457150	125.14222743603900	78.09158681356650	6.408509360839470	6.884001306905410	High	Poor	Happy	21.77864528647280	29.80839115661620	High
5	20	M	68.34276946853390	107.51559175322700	80.67495681247280	7.264718754103350	4.483450107413980	Moderate	Poor	Happy	8.964999082744990	15.19404500227400	Moderate
6	22	M	61.74415156882540	90.0	84.45086471353400	4.2625178607838200	6.825000550580700	Moderate	Good	Happy	44.94822898892330	15.651195419247700	Moderate
7	22	F	93.09722755406190	106.67778156409800	76.49981545408860	8.415978866919800	1.4029154981963800	Moderate	Moderate	Happy	10.26895000015160	23.81809603875400	Moderate
8	24	F	63.38103086453330	115.39678276968800	70.03950213998230	2.8367889504018400	1.4864286461475400	Low	Good	Neutral	39.64178773383430	19.315438930026200	Low
9	19	F	81.67102914383640	142.82097746747900	74.67687740740210	5.221366727863390	5.115073778575900	Moderate	Moderate	Neutral	37.21119508758500	8.854985697527700	Moderate
10	20	M	86.21108557787720	112.63593825193000	91.95139179833440	8.208448493872870	5.946600709954310	Moderate	Moderate	Happy	33.195541012198800	19.744741140335500	Moderate
11	24	F	54.30039682342690	127.67886138543900	81.16173084076550	2.621427829593570	8.990358763022370	High	Moderate	Happy	24.772250517834400	21.214005410299700	Moderate
12	20	M	70.73962640735700	117.47537649058900	71.551015468589530	1.9887680236670500	1.9419457051710700	Moderate	Good	Stressed	35.02797294525580	23.088585378878900	Low
13	20	F	64.82468490917880	138.39607715805600	84.01560608387720	8.953155709191910	8.4038313447450000	Moderate	Good	Neutral	40.48484873704500	12.992541591962600	High
14	22	F	64.01075603440660	137.3564457819590	84.68315584939140	4.086366648327800	2.8904260356286300	Moderate	Moderate	Neutral	20.14189503524530	12.627747322655200	Low
15	21	F	75.55378816808910	137.6426265785990	85.19999075797810	7.154831124181290	4.5467328365902900	High	Moderate	Neutral	39.7637927949800	10.290924892783300	Moderate
16	20	M	63.77129503576720	93.99024705873970	70.22162799628710	1.2281374383754500	3.849421913446720	High	Poor	Happy	39.30874511937240	16.5403732168070	High
17	23	F	65.04030959673530	123.10086070867500	73.17090005714110	7.2104219958685050	4.601397455788590	Moderate	Good	Happy	32.930510909215600	18.664337912278500	Moderate
18	22	F	68.31885104374710	131.3840790018290	73.23107836129990	3.844911803922740	2.35445482653270	Moderate	Moderate	Happy	24.620251219188900	12.433354066612500	Low
19	20	M	66.26607318234300	125.72307843884000	76.39352087933810	1.5016676822669800	5.532130669632000	Moderate	Good	Stressed	39.35457364848470	14.29023511654200	Moderate
20	21	F	58.26429782987880	129.54894488732700	85.389919185143740	6.488314853933050	8.214783098482310	Moderate	Good	Stressed	35.083129561625200	4.7314110219046800	Moderate
21	23	M	77.9506054847470	120.87526506428200	82.98184804731010	3.167234346858560	4.973809093410070	Low	Moderate	Neutral	19.840430388054300	22.5472452501330	Low
22	23	F	80.48530938507600	108.8282368991020	83.75216957157160	4.48126517846692	5.415222183083790	Moderate	Moderate	Happy	26.32844191231750	7.175118142499520	Moderate
23	19	M	73.43818589353220	94.54353922520620	93.54019142371520	1.7252273105788200	8.63890826751830	Low	Good	Happy	38.190549384207000	17.867389876896400	Moderate
24	21	M	66.23447348977730	110.31645857027300	86.88890932491630	2.324469259854510	7.354523031742280	High	Good	Happy	31.145806171008400	16.19127416686170	Moderate
25	22	F	62.947892636702200	160.1716481045760	78.58198375003910	3.6817951194625900	5.414089676678060	Moderate	Poor	Neutral	40.859828944453300	15.229039842658000	Moderate
26	18	M	62.963799370735500	112.71523235664000	68.16114823947500	1.8078355050414100	3.319260833702620	High	Moderate	Happy	26.583533043516000	15.343984391696400	Low
27	21	M	93.57373635032710	90.6091466063880	72.85452991529130	9.197436382176840	8.96434072160470	Moderate	Moderate	Stressed	20.001047077519000	9.941094121385630	High
28	19	M	79.37838301566640	130.98974041093600	77.38294932321770	1.0544233886581010	6.193183615215030	Moderate	Moderate	Happy	51.93602615610070	13.359673979027500	Moderate
29	23	M	61.598899214166600	132.83190013024000	88.29129827336110	5.354741969778800	2.9903284909962400	Moderate	Poor	Stressed	13.191049172535100	19.844753261260900	Moderate
30	22	M	71.84697827437140	139.309886585228200	79.72668065893910	2.109227123192960	3.4220777737420600	Low	Poor	Neutral	25.001604025983400	21.3257750218320	Low
31	21	M	68.37464974900640	127.0646775556020	73.43804045111560	3.7691118053014300	8.228197800426210	Moderate	Good	Neutral	32.79748000699980	12.895495307125300	Moderate
32	18	F	75.83800376156290	132.68402271700400	72.67107114272770	9.334536038364390	1.3973100567871000	Low	Moderate	Neutral	39.860838891777710	7.1979023476154300	Moderate
33	18	M	75.43499102789340	147.60325702263000	79.30421791078770	2.4552509794658800	9.928829683256380	Moderate	Poor	Neutral	36.46449368542260	11.106630734351480	Moderate
34	20	M	80.83096053905470	108.81789761906800	88.8998756221710	9.692490877464990	2.173165117433130	Moderate	Moderate	Neutral	53.11259004942880	24.28348547426440	Moderate
35	20	F	61.49443541337800	105.84223941864400	95.91019305775120	2.8521493929048000	8.164616728572440	Moderate	Good	Neutral	28.824015473418700	11.404277550060390	Moderate
36	24	F	71.04732989638130	103.58120190670600	77.94231401970330	7.6740949762014700	4.552977850945150	Low	Good	Neutral	21.270167388494700	14.513981648701100	Moderate
37	19	F	52.51453015042470	116.01345620259000	79.55256819237400	5.920095384075810	5.429945649587200	Low	Moderate	Stressed	27.03039757835500	22.079561810099100	Moderate
38	21	F	57.42656996711600	141.50133629737800	67.80058293842620	8.9603026548370830	8.92585519718120	Moderate	Moderate	Neutral	26.80023470649800	12.365519973314500	High
39	21	F	78.84100836559180	127.92102551619400	99.39192569099470	5.985939299528180	4.716200706302920	High	Good	Neutral	29.317048527098600	14.172323051761600	Moderate
40	24	F	60.49316208894670	127.33037929063800	65.9897161843920	5.811032972406780	5.837443783059110	Moderate	Moderate	Stressed	20.175118771007200	18.97249145378160	Moderate
41	23	F	61.590912019521500	92.71038107002910	91.71662311780550	2.759122003315600	1.4436969616825300	Moderate	Moderate	Happy	37.84056038250360	23.559034337647500	Moderate
42	23	M	69.91721925449950	109.60219759145400	76.34664046987530	2.4960642584649300	1.0211272648070800	Moderate	Moderate	Happy	33.70830209338210	11.772240676525100	Low
43	24	M	76.94761030846540	92.46513550706950	74.00933505159680	9.2484022360039630	3.203198857015530	Low	Poor	Neutral	17.36119170925890	21.960476220056900	Moderate
44	23	F	59.39173204165900	137.56868683982300	68.89996428114070	3.9505437896503100	7.125247836268910	High	Moderate	Neutral	39.68025075870580	13.592356450678400	Moderate
45	20	M	50.0	143.784697393286	104.10048649407100	4.3487402946660600	2.0913840402322300	Moderate	Good	Stressed	21.865184630155000	10.570359472675300	Low
46	21	M	61.29776757856490	131.69350090860600	82.21491623278400	7.8384129258655500	2.0238576292409400	Low	Good	Neutral	18.17184497665710	16.247782151537800	Moderate
47	24	M	82.05628077952840	105.19820931390700	80.21280437104060	4.049411341805270	4.698645472323130	Moderate	Moderate	Neutral	30.87328703921650	9.212159711570200	Low
48	21	F	73.43135051016870	142.23033740423700	88.86732359554540	9.588562019751020	6.830533347474040	Low	Good	Neutral	12.421565506009090	17.360057157314400	Moderate
49	18	F	83.23674803787210	116.20117634962200	73.49167274370040	7.641349920989200	7.408647930231380	High	Poor	Happy	35.89503049877250	30.7482405567270	High

Figure 11 Excel data before normalisation

6. ER Model: Chen's Notation

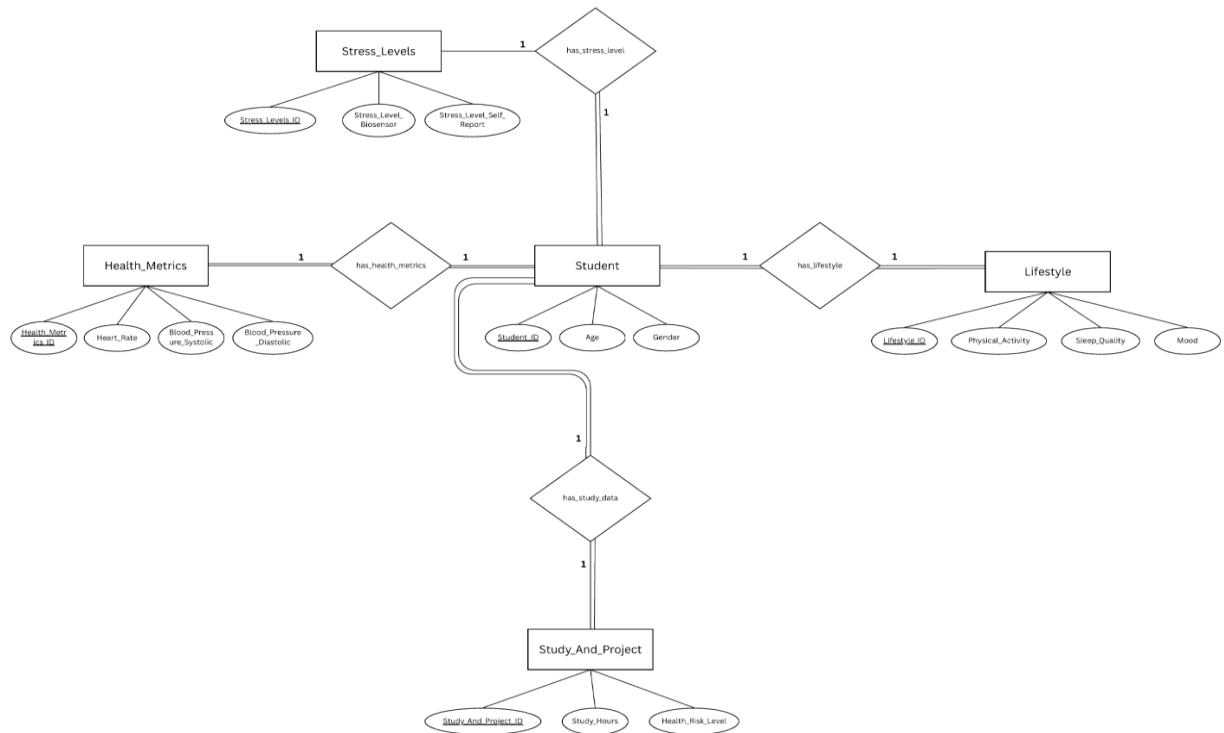


Figure 12 Chen's Notation

7. Normalisation

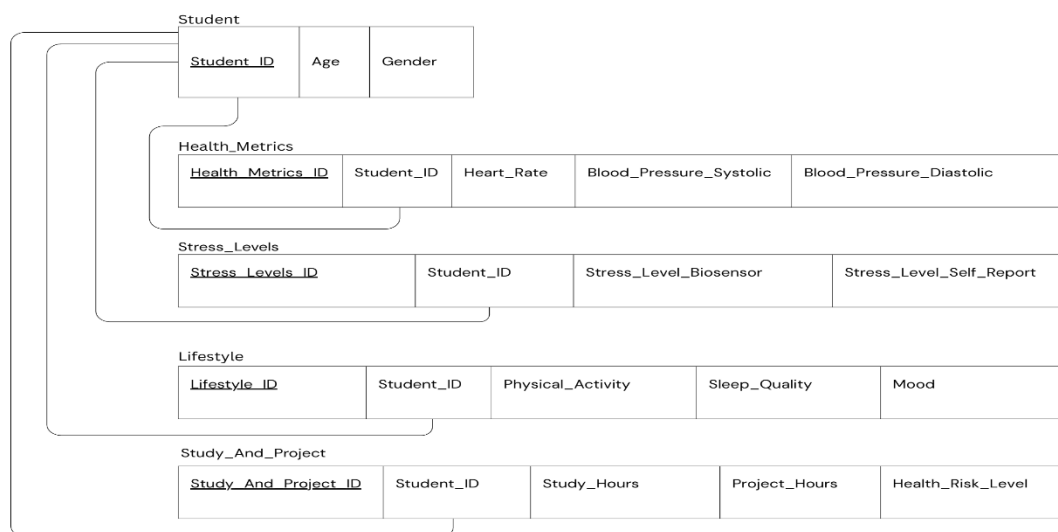


Figure 13 Table after normalisation with primary and foreign key

5. Lessons Learned

- How to build a full data pipeline from ingestion to analysis
- Difference between structured (SQL) and semi-structured (NoSQL) models
- Importance of indexing and normalized design
- Hands-on experience with PySpark stream filtering
- Real teamwork in dividing tasks across phases

6. Project Demo

We have prepared a live demo using Google Colab, which includes:

- **Table creation and SQL queries** using SQLite (**Phase 1**)
Relational schema, data insertion, indexing, and analytical queries.
- **TinyDB JSON integration and search** (**Phase 2**)
JSON conversion, document insertion, and NoSQL-style filtering.
- **Stream simulation and filters** using PySpark (**Phase 3**)
Real-time-like processing to detect high blood pressure, stress, and heavy study loads.
- **Interactive dashboard** using Streamlit (**Phase 4**)
Combines all data sources into one visual, user-friendly interface to explore and filter student health insights.