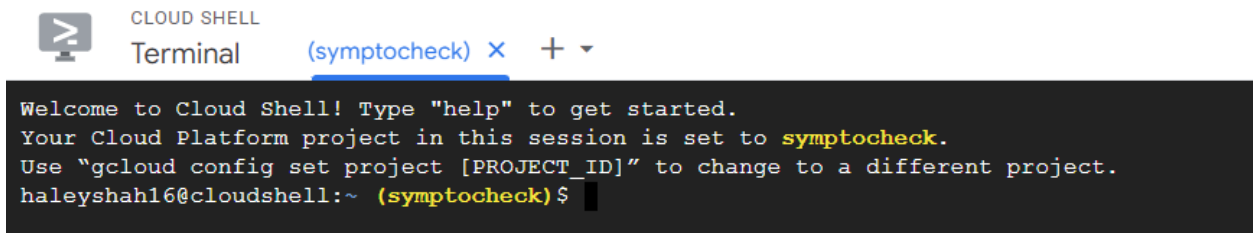


## GCP Connection:



```
Cloud Shell
Terminal (symptochek) × + ▾

Welcome to Cloud Shell! Type "help" to get started.
Your Cloud Platform project in this session is set to symptochek.
Use "gcloud config set project [PROJECT_ID]" to change to a different project.
haleyshah16@cloudshell:~ (symptochek)$
```

## DDL Commands:

```
CREATE TABLE Records(
    user_id VARCHAR(255),
    age INT,
    sex VARCHAR(10),
    country VARCHAR(10),
    checkin_date DATETIME,
    trackable_id VARCHAR(255),
    trackable_type VARCHAR(255),
    trackable_name VARCHAR(255),
    trackable_value INT
);

CREATE TABLE Symptoms (
    SymptomId INT AUTO_INCREMENT PRIMARY KEY,
    SymptomName VARCHAR(255) NOT NULL,
    Severity INT
);

CREATE TABLE Treatment (
    TreatmentId INT AUTO_INCREMENT PRIMARY KEY,
    Medicine VARCHAR(255) NOT NULL,
    Dosage VARCHAR(30)
);

CREATE TABLE Illness(
    IllnessId INT AUTO_INCREMENT PRIMARY KEY,
    IllnessName VARCHAR(255) NOT NULL
);

CREATE TABLE Has_Symptom(
    Illness_Symptom_Id INT AUTO_INCREMENT PRIMARY KEY,
    IllnessName VARCHAR(255),
    SymptomName VARCHAR(255),
    Severity INT
);
```

```
CREATE TABLE Has_Treatment(
    Illness_Treatment_Id INT AUTO_INCREMENT PRIMARY KEY,
    IllnessName VARCHAR(255),
    Medicine VARCHAR(255) NOT NULL,
    Dosage VARCHAR(30)
);
```

### Insert Into Tables:

```
mysql> SELECT COUNT(user_id) FROM Records;
+-----+
| COUNT(user_id) |
+-----+
|          7976223 |
+-----+
1 row in set (18.99 sec)
```

```
mysql> SELECT COUNT(SymptomId) FROM Symptoms;
+-----+
| COUNT(SymptomId) |
+-----+
|           63142 |
+-----+
1 row in set (0.33 sec)
```

```
mysql> SELECT COUNT(TreatmentId) FROM Treatment;
+-----+
| COUNT(TreatmentId) |
+-----+
|           19557 |
+-----+
1 row in set (0.24 sec)
```

```
mysql> SELECT COUNT(IllnessId) FROM Illness;
+-----+
| COUNT(IllnessId) |
+-----+
|           9203 |
+-----+
1 row in set (0.22 sec)
```

```
mysql> SELECT COUNT(Illness_Symptom_Id) FROM Has_Symptom;
+-----+
| COUNT(Illness_Symptom_Id) |
+-----+
|                14468863 |
+-----+
1 row in set (2.38 sec)
```

```
mysql> SELECT COUNT(Illness_Treatment_Id) FROM Has_Treatment;
+-----+
| COUNT(Illness_Treatment_Id) |
+-----+
|                297965 |
+-----+
1 row in set (0.44 sec)
```

## Advanced Queries:

1. Query to get the top conditions/illnesses in which anxiety is a symptom.

```
SELECT r2.trackable_name, COUNT(r2.user_id) AS
`count_user_id`
FROM Records r1 JOIN Records r2 ON r1.user_id = r2.user_id
WHERE r1.trackable_type = 'Symptom' AND r1.trackable_name =
'Anxiety' AND r2.trackable_type = 'Condition' AND
r1.checkin_date = r2.checkin_date
GROUP BY r2.trackable_name
ORDER BY `count_user_id` DESC
LIMIT 15;
```

trackable_name	count_user_id
Anxiety	12982
Fibromyalgia	12650
Depression	11673
Chronic fatigue syndrome	6651
Migraine	5921
IBS	3917
Generalized anxiety disorder	3828
Ehlers-Danlos syndrome	3170
Irritable bowel syndrome	3159
Lyme disease	3120
POTS	2835
Headaches	2663
Asthma	2509
Fatigue	2380
Allergies	2314

15 rows in set (56.83 sec)

EXPLAIN ANALYZE for Original Result:

```

-----+
| -> Limit: 15 row(s) (actual time=56052.569..56052.572 rows=15 loops=1)
-> Sort: count_user_id DESC, limit input to 15 row(s) per chunk (actual time=56052.567..56052.569 rows=15 loops=1)
    -> Table scan on <temporary> (actual time=0.003..0.669 rows=2252 loops=1)
        -> Aggregate using temporary table (actual time=56051.460..56052.282 rows=2252 loops=1)
            -> Inner hash join (r2.user_id = r1.user_id), (r2.checkin_date = r1.checkin_date) (cost=598843837.71 rows=619634)
(actual time=26884.533..55823.834 rows=22441 loops=1)
                -> Filter: (r2.trackable_type = 'Condition') (cost=6809.60 rows=7872) (actual time=0.082..26251.959 rows=1111517
loops=1)
                    -> Table scan on r2 (cost=6809.60 rows=7871685) (actual time=0.076..25205.263 rows=7976223 loops=1)
                        -> Hash
                            -> Filter: ((r1.trackable_type = 'Symptom') and (r1.trackable_name = 'Anxiety'))
(cost=851039.50 rows=78717) (actual time=0.115..26714.584 rows=61602 loops=1)
                                -> Table scan on r1 (cost=851039.50 rows=7871685) (actual time=0.084..25332.205
rows=7976223 loops=1)
|
+-----+

```

Indexing #1 with trackable\_id:

Trackable id is in the where clause and indexing this should give us a slight reduction on runtime as seen below.

```

mysql> CREATE INDEX trackable_id_idx on Records(trackable_id(10));
Query OK, 0 rows affected (49.95 sec)
Records: 0 Duplicates: 0 Warnings: 0

```

```

mysql> EXPLAIN ANALYZE SELECT r2.trackable_name,COUNT(r2.user_id) AS 'count user id' FROM Records r1 JOIN Records r2 ON r1.user_id = r
2.user_id WHERE r1.trackable_type = 'Symptom' AND r1.trackable_name = 'Anxiety' AND r2.trackable_type = 'Condition' AND r1.checkin_date
= r2.checkin_date GROUP BY r2.trackable_name ORDER BY 'count_user_id' DESC
LIMIT 15;

```

```

-----+
| EXPLAIN
+-----+
|
+-----+
| -> Limit: 15 row(s) (actual time=48282.251..48282.253 rows=15 loops=1)
    -> Sort: count_user_id DESC, limit input to 15 row(s) per chunk (actual time=48282.250..48282.251 rows=15 loops=1)
        -> Table scan on <temporary> (actual time=0.002..0.534 rows=2252 loops=1)
            -> Aggregate using temporary table (actual time=48281.346..48282.006 rows=2252 loops=1)
                -> Inner hash join (r2.user_id = r1.user_id), (r2.checkin_date = r1.checkin_date) (cost=598843837.71 rows=619634) (ac
tual time=22953.869..48063.539 rows=22441 loops=1)
                    -> Filter: (r2.trackable_type = 'Condition') (cost=6809.60 rows=7872) (actual time=0.075..22666.757 rows=1111517
loops=1)
                        -> Table scan on r2 (cost=6809.60 rows=7871685) (actual time=0.069..21691.393 rows=7976223 loops=1)
                            -> Hash
                                -> Filter: ((r1.trackable_type = 'Symptom') and (r1.trackable_name = 'Anxiety')) (cost=851039.50 rows=78717)
(actual time=0.117..22825.209 rows=61602 loops=1)
                                    -> Table scan on r1 (cost=851039.50 rows=7871685) (actual time=0.098..21531.255 rows=7976223 loops=1)
|
+-----+
1 row in set (48.31 sec)

```

Record id isn't really seen in this query, so we should not expect a significant reduction in runtime, as seen below.

```
_id> EXPLAIN ANALYZE SELECT r2.trackable_name,COUNT(r2.user_id) AS 'count_user'
-> FROM Records r1 JOIN Records r2 ON r1.user_id = r2.user_id
-> WHERE r1.trackable_type = 'Symptom' AND r1.trackable_name = 'Anxiety' AND r2.trackable_type = 'Condition' AND r1.checkin_date = r2.checkin_date
-> GROUP BY r2.trackable_name
-> ORDER BY 'count_user_id' DESC
-> LIMIT 15;
```

```
| EXPLAIN                                     |
+-----+-----+
| -> Limit: 15 row(s)  (actual time=52148.384..52148.388 rows=15 loops=1)
|   -> Sort: count_user_id DESC, limit input to 15 row(s) per chunk  (actual time=52148.382..52148.385 rows=15 loops=1)
|     -> Table scan on <temporary>  (actual time=0.002..0.683 rows=2252 loops=1)
|       -> Aggregate using temporary table  (actual time=52147.081..52147.938 rows=2252 loops=1)
|         -> Inner hash join (r2.user_id = r1.user_id), (r2.checkin_date = r1.checkin_date)  (cost=598843837.71 rows=619634) (actual time=25153.772..51908.223 rows=224441 loop
ps=1)
|           -> Filter: (r2.trackable_type = 'Condition')  (cost=6809.60rows=7872) (actual time=0.981..24165.746 rows=1111517 loops=1)
|             -> Table scan on r2  (cost=6809.60 rows=7871685) (actual time=0.974..23149.448 rows=7976223 loops=1)
|               Hash
|                 -> Filter: ((r1.trackable_type = 'Symptom') and (r1.trackable_name = 'Anxiety'))  (cost=851039.50 rows=78717) (actual time=0.095..24993.066 rows=61602 loops
=1)
|                   -> Table scan on r1  (cost=851039.50 rows=7871685) (actual time=0.072..23647.995 rows=7976223 loops=1)
|
+-----+-----+
1 row in set (52.18 sec)
```

The where clause uses `trackable_name` from `Records` and we know that what's in the where clause affects the runtime. By adding an index, we significantly reduce the time as seen in the `EXPLAIN ANALYZE`:

```
mysql> EXPLAIN ANALYZE SELECT r2.trackable_name,COUNT(r2.user_id) AS `count_user_id` FROM Records r1 JOIN Records r2 ON r1.user_id = r2.user_id WHERE r1.trackable_type = 'Symptom' AND r1.trackable_name = 'Anxiety' AND r2.trackable_type = 'Condition' AND r1.checkin_date = r2.checkin_date GROUP BY r2.trackable_name ORDER BY `count_user_id` DESC LIMIT 15;
```

```
1 row in set (35.37 sec)
```

2. Query to show the most frequent treatment method for some symptoms

```
SELECT t.Medicine FROM Has_Treatment t JOIN (SELECT IllnessName FROM Has_Symptom
WHERE SymptomName LIKE 'Headache' AND Severity = 4 GROUP BY IllnessName ORDER BY
COUNT(IllnessName) DESC LIMIT 15) AS s ON t.IllnessName = s.IllnessName GROUP BY
t.Medicine ORDER BY COUNT(t.Medicine) DESC LIMIT 15;
```

When we run `explain analyze` with no indices, we get a time of 7.077 seconds:

### Indexing analysis 1: Add index on Medicine (improved time: 5.91 sec)

We will be showing users potential medicines more frequently, and relative to the number of entries (which is in the millions) medicine will be much smaller. Adding an index on Medicine should speed this up.

```
mysql> EXPLAIN ANALYZE SELECT t.Medicine FROM HasTreatment t JOIN (SELECT IllnessName FROM HasSymptom WHERE SymptomName LIKE 'Headache' AND Severity = 4 GROUP BY IllnessName ORDER BY COUNT(IllnessName) DESC LIMIT 15) AS s ON t.IllnessName = s.IllnessName GROUP BY t.Medicine ORDER BY COUNT(t.Medicine) DESC LIMIT 15;
```



### Indexing analysis 2: Index on SymptomName (time: 3.017s)

The Where clause uses Symptom name in this advanced query and we can add an index here to reduce the time searching. This is proven true when we run EXPLAIN ANALYZE:

[illegible]

### Indexing analysis 3: Index on IllnessName

IllnessName is used in the subquery, so indexing this column should help reduce our time. Looking at the EXPLAIN ANALYZE, we can see that there is a slight time reduction, although not as significant as adding an index on SymptomName.

[illegible]