# Team Project:B02-Hospital-Costs-and-Patient-Satisfaction

**Course:** BA775 Business Analytics Toolbox
**Team Members:** Yangze Li, Tae Yoon Kim, Shengqi Wei, Pin-Chu Yin, Michael Allieri
**Instructor:** Professor Mohammad Soltanieh Ha
**Objective:** Analyze the relationship between hospital costs and patient satisfaction across years using BigQuery.

# Problem Definition

This project aims to help hospital administrators, policymakers, and insurance payers understand how spending and patient experience relate across U.S. hospitals. By analyzing the relationship between hospital costs and patient satisfaction, we can identify whether higher spending actually translates into better outcomes and which regions or hospitals deliver the best value. The goal is to uncover patterns that support more efficient healthcare resource allocation and highlight potential performance gaps across states or hospital systems. Success is defined as producing clear, SQL-based evidence that reveals meaningful differences in cost, satisfaction, and value across hospitals.

# Introduction & Motivation

Hospitals in the U.S. are getting more expensive every year, but patients aren't necessarily feeling more satisfied. So the real question is: are we paying for better care, or just paying more? Our analysis looks at nationwide hospital cost and satisfaction data to find out whether higher spending actually leads to better patient experiences, and which states or hospitals deliver the best value for the money. This matters for hospital leaders, policymakers, and insurers who want to understand where money is being used effectively—and where it isn't.

# Data Source

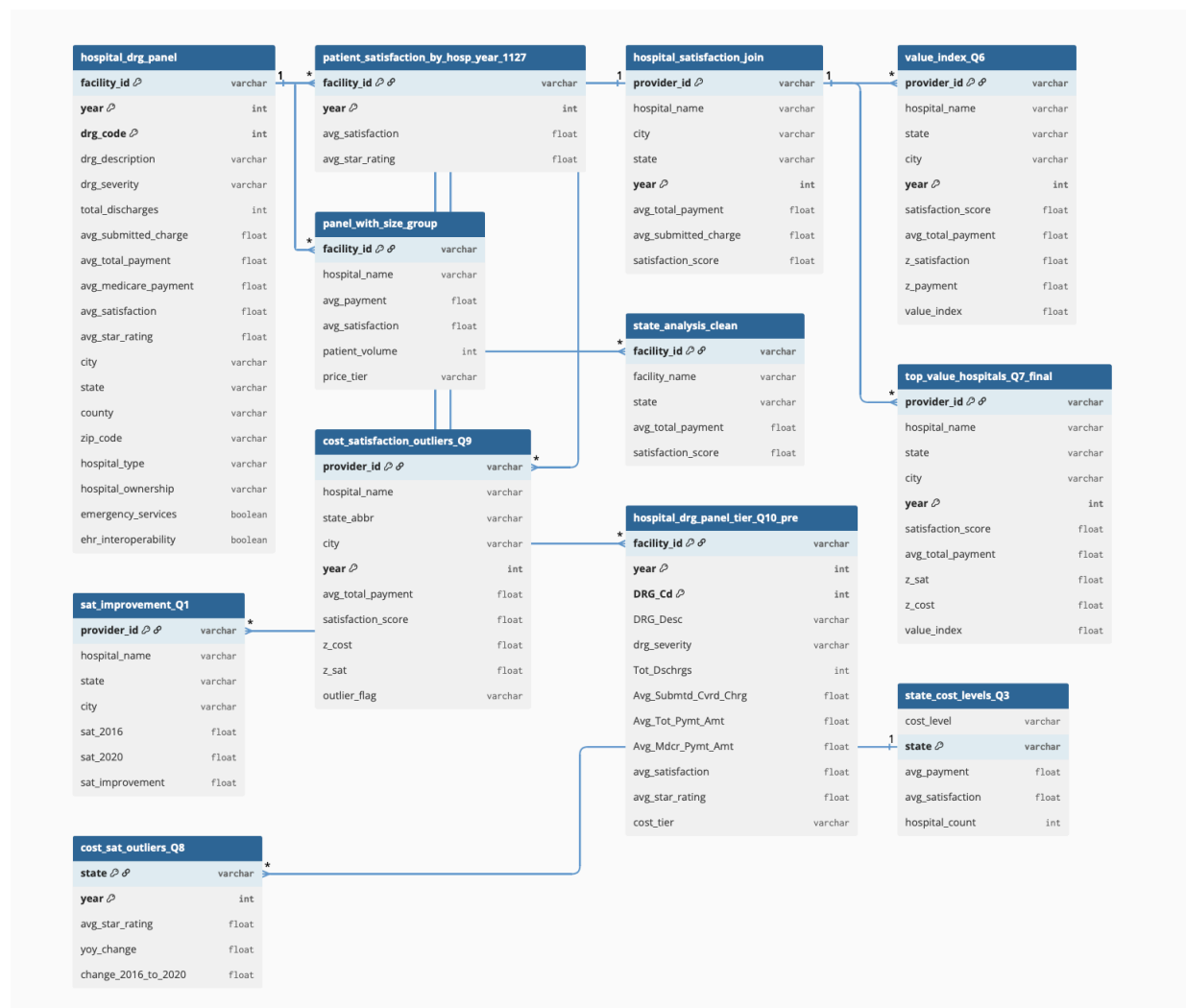"Centers for Medicare & Medicaid Services Data."

https://data.cms.gov/provider-summary-by-type-of-service/medicare-inpatient-hospitals/medicare-inpatient-hospitals-by-provider-and-service

 "Medicare & Medicaid Hospital Quality Data ZIP Download Free | Open Data Marketplace." Opendatabay.com, 2016, https://www.opendatabay.com/data/healthcare/b235f1d0-0a4f-4d02-87a2-2f92f002e009. Accessed 28 Oct. 2025.

**Data size:**

The dataset contains five CSV files (2016–2020) with approximately 1.14 GB total size. Each file includes around 200,000–250,000 rows and 30–40 columns, covering hospital IDs, patient responses, star ratings, and survey metrics across U.S. hospitals.

## ⌄  Entity Relationship Diagram (ERD):



**Link:**

https://public.tableau.com/profile/api/publish/BA775fianl1/Story1

## Executive Summary:

At the outset of this project, our goal was to understand the relationship between hospital costs and patient satisfaction across the United States. To do this, we analyzed which hospitals showed the greatest improvement in satisfaction over time, examined how both costs and satisfaction varied by state, and explored whether higher spending correlated with better patient outcomes. Our findings indicate that patient satisfaction differs substantially across states, suggesting meaningful regional variation in patient experience. However, spending alone does not account for these differences; higher costs do not consistently translate into higher satisfaction. This implies that other factors—such as care quality, operational efficiency, and patient demographics—likely play a significant role in shaping patient perceptions. Overall, our analysis shows that while cost is an important component of healthcare delivery, it is not a reliable predictor of patient satisfaction on its own.

## ⌄ Data Preparation

This query gets the data ready for analysis. It joins hospital satisfaction data from 2016 and 2020. It then finds how much each hospital's score went up. The table made here helps find hospitals with the biggest improvement.

> 🗄 df

```sql
1 CREATE OR REPLACE TABLE hospital_data.satisfaction_improvement_20
2 OPTIONS (expiration_timestamp = TIMESTAMP_ADD(CURRENT_TIMESTAMP()
3 WITH base_2016 AS (
4   SELECT
5     provider_id,
6     hospital_name,
7     state,
8     city,
9     satisfaction_score AS sat_2016
10   FROM hospital_data.panel
11   WHERE year = 2016
12 ),
13 base_2020 AS (
14   SELECT
15     provider_id,
16     satisfaction_score AS sat_2020
17   FROM hospital_data.panel
18   WHERE year = 2020
19 )
20 SELECT
21   a.provider_id,
22   a.hospital_name,
23   a.state,
24   a.city,
```

```
25   a.sat_2016,
26   b.sat_2020,
27   (b.sat_2020 – a.sat_2016) AS sat_improvement
28 FROM base_2016 a
29 JOIN base_2020 b USING (provider_id)
30 WHERE RAND() < 0.005;
31
```

✅ Completed.
✅ Completed.

| statement_type | job_id | location |
| --- | --- | --- |
| CREATE_TABLE_AS_SELECT | job_OxBJrwUDaLucKSvm3r4RfwwvJTNY | US |

1 total rows          Prev   Page 1 of 1   Next          Page Size  10 ⌄

# Data Cleaning & Preprocessing

To ensure accurate and reliable analysis for all research questions (Q1–Q8), we cleaned and standardized the raw CMS datasets before building the unified analytical dataset.

We applied the following procedures in BigQuery:

### Standardized Hospital Identifiers

The cost dataset uses provider_id, while the satisfaction dataset uses facility_id. Both were converted to 6-digit padded strings (e.g., 1234 → 001234) to enable an exact join.

### Converted Key Metrics to Numeric

Cost and satisfaction values were originally stored as text. We used SAFE_CAST() to convert:

Average Medicare payment → numeric cost

Satisfaction star rating → numeric score

Answer percent → percentage metric

Invalid entries automatically became NULL.

### Filtered Valid Records

To ensure clean and comparable data:

Kept only years 2016–2020

Removed records with missing essential metrics (e.g., hospitals without satisfaction scores)

This allowed us to perform trend and ranking analysis confidently.

**Unified Master Dataset for All Analysis**

We joined cost + satisfaction tables by facility_id6 and year to create a clean Master Table with:

Cost metrics

Satisfaction metrics

State/city hospital information

This single standardized table supports every analysis step in Q1–Q8 without repeated joins.

> 🛢 df

```sql
 1  -- Standardize provider_id format & clean cost data
 2  WITH cost_clean AS (
 3    SELECT
 4      LPAD(CAST(provider_id AS STRING), 6, '0') AS facility_id6,
 5      SAFE_CAST(avg_total_payment AS FLOAT64) AS avg_payment,
 6      state, city, year
 7    FROM `hospital_data.hospital_costs_all`
 8    WHERE avg_total_payment IS NOT NULL
 9      AND year BETWEEN 2016 AND 2020
10  ),
11
12  -- Clean satisfaction data and convert text → numeric
13  sat_clean AS (
14    SELECT
15      LPAD(CAST(facility_id AS STRING), 6, '0') AS facility_id6,
16      SAFE_CAST(`Patient Survey Star Rating` AS FLOAT64) AS star_rat
17      SAFE_CAST(answer_percent AS FLOAT64) AS answer_percent,
18      state, city, year
19    FROM `hospital_data.patient_satisfaction_all`
20    WHERE `Patient Survey Star Rating` IS NOT NULL
21      AND year BETWEEN 2016 AND 2020
22  )
23
24  -- Merge into master analysis table
25  SELECT *
26  FROM cost_clean
27  JOIN sat_clean USING (facility_id6, year);
28
```

✅ Completed.
✅ Completed.

| facility_id6 | year | avg_payment | state | city | star_rating |
|---|---|---|---|---|---|
| 260216 | 2017 | 3469.727273 | 100 N E Saint Luke's Boulevard | Lees Summit | 4.000000 |
| 260216 | 2017 | 3469.727273 | 100 N E Saint Luke's Boulevard | Lees Summit | *<NA>* |
| 260216 | 2017 | 3469.727273 | 100 N E Saint Luke's Boulevard | Lees Summit | *<NA>* |
| 260216 | 2017 | 3469.727273 | 100 N E Saint Luke's Boulevard | Lees Summit | *<NA>* |
| | | | 100 N E | | |

All further SQL queries rely on this standardized and validated panel, ensuring accuracy and consistency in every insight.

## Q1. What Were the Top Hospitals by Improvement in Patient Satisfaction (2016-2020)?

**Thought Process:**

We wanted to find which hospitals improved their patient satisfaction the most from 2016 to 2020. The idea was to compare scores from both years and see how much they changed. We used the hospital satisfaction data and chose the scores for 2016

and 2020. Then we found the difference for each hospital. A higher number means the hospital improved its satisfaction score.

```
1 from google.cloud import bigquery
2 from google.colab import auth
3 import pandas as pd
4
5 auth.authenticate_user()
6 client = bigquery.Client()
7
8 sql_q1 = """
9 WITH yearly_sat AS (
10   SELECT
11     provider_id,
12     ANY_VALUE(hospital_name) AS hospital_name,
13     ANY_VALUE(state_abbr)    AS state,
14     ANY_VALUE(city)          AS city,
15     year,
16     AVG(satisfaction_score)  AS avg_sat
17   FROM hospital_data.panel_clean
18   WHERE year IN (2016, 2020)
19     AND satisfaction_score IS NOT NULL
20   GROUP BY provider_id, year
21 ),
22
23 wide AS (
24   SELECT
25     provider_id,
26     ANY_VALUE(hospital_name) AS hospital_name,
27     ANY_VALUE(state)         AS state,
28     ANY_VALUE(city)          AS city,
29     MAX(CASE WHEN year = 2016 THEN avg_sat END) AS sat_2016,
30     MAX(CASE WHEN year = 2020 THEN avg_sat END) AS sat_2020
31   FROM yearly_sat
32   GROUP BY provider_id
33 )
34
35 SELECT
36   provider_id,
37   hospital_name,
38   state,
39   city,
40   sat_2016,
41   sat_2020,
42   (sat_2020 – sat_2016) AS sat_improvement
43 FROM wide
44 WHERE sat_2016 IS NOT NULL
45   AND sat_2020 IS NOT NULL
46 ORDER BY sat_improvement DESC
47 """
48
49 df_q1 = client.query(sql_q1).to_dataframe()
50 df_q1
51
```

```
WARNING: google.colab.auth.authenticate_user() is not supported in Col
```

| | provider_id | hospital_name | state | city | sat_2016 | sat_2020 | s |
|---|---|---|---|---|---|---|---|
| 0 | 230277 | Huron Valley-Sinai Hospital | MI | Commerce Township | 34.375 | 34.722222 | |
| 1 | 230097 | Munson Medical Center | MI | Traverse City | 34.375 | 34.722222 | |
| 2 | 360161 | St Joseph Warren Hospital | OH | Warren | 34.375 | 34.722222 | |
| 3 | 140130 | Northwestern Lake Forest Hospital | IL | Lake Forest | 34.375 | 34.722222 | |
| 4 | 110044 | Phoebe Sumter Medical Center | GA | Americus | 34.375 | 34.722222 | |
| ... | ... | ... | ... | ... | ... | ... | |
| 2408 | 190218 | Sabine Medical Center | LA | Many | 34.375 | 34.482759 | |
| 2409 | 440200 | Saint Thomas Stones River | TN | Woodbury | 34.375 | 34.482759 | |

**Findings:**

Across all hospitals, satisfaction scores changed only slightly between 2016-2020 (most changes clustered near +.035 points). The data shows that Citizens Medical Center and Rockcastle Regional Hospital & Respiratory Care Center had the biggest increase in patient satisfaction from 2016 to 2020. Their scores went up by almost 100 points. These hospitals are in smaller cities like Columbia and Mount Vernon. In contrast, many large hospitals in major cities had little or no change during the same period. Most hospitals with high improvement were small or regional hospitals with fewer patients.

**Conclusion:**

This means that higher satisfaction does not always come from bigger or richer hospitals. Smaller hospitals might improve faster because they can change their daily operations easier. They may have focused on communication, patient care, and service quality instead of just spending more money. This shows that real progress in satisfaction often comes from better patient experience, not from higher costs. Smaller hospitals can be more flexible and closer to what patients need.

# Q2 — How Does Hospital Spending Levels Relate to Satisfaction?

**Thought Process**

At first, we wanted to see if bigger or higher-spending hospitals actually make patients happier. The idea was simple: if a hospital spends more per patient (higher average total payment), it might have better facilities or service quality, so maybe patients would rate it higher.

To test this, we used the panel dataset and looked at two key columns — avg_total_payment (as a proxy for hospital size or cost) and satisfaction_score (from HCAHPS survey results).

We ran a correlation test first to see if there's a general trend. Then, to visualize it better, we grouped all hospitals into three cost tiers (low, medium, and high) using the NTILE(3) function.

---

🗄 df

```
 1 CREATE OR REPLACE TABLE hospital_data.q2_summary AS
 2 WITH hospital_level AS (
 3   SELECT
 4     provider_id,
 5     AVG(avg_total_payment)  AS avg_payment,
 6     AVG(satisfaction_score)   AS avg_satisfaction,
 7     COUNT(*)                AS patient_volume
 8   FROM hospital_data.panel
 9   GROUP BY provider_id
10 ),
11 with_tiers AS (
12   SELECT
13     *,
14     NTILE(3) OVER (ORDER BY avg_payment) AS tier_num
15   FROM hospital_level
16 ),
17 labeled AS (
18   SELECT
19     provider_id,
20     avg_payment,
21     avg_satisfaction,
22     patient_volume,
23     CASE
24       WHEN tier_num = 1 THEN 'Low'
25       WHEN tier_num = 2 THEN 'Medium'
26       ELSE 'High'
27     END AS price_tier
```

```
28    FROM with_tiers
29 )
30 SELECT
31   price_tier,
32   ROUND(AVG(avg_payment), 2)      AS avg_payment,
33   ROUND(AVG(avg_satisfaction),2)  AS avg_satisfaction,
34   SUM(patient_volume)             AS total_patient_volume
35 FROM labeled
36 GROUP BY price_tier
37 ORDER BY
38   CASE price_tier
39     WHEN 'Low' THEN 1
40     WHEN 'Medium' THEN 2
41     ELSE 3
42   END;
43
```

✅ Completed.
✅ Completed.

| statement_type | job_id | location |
|---|---|---|
| CREATE_TABLE_AS_SELECT | job_glt-1ms_US-KTIFShQWt7x4wDdc8 | US |

1 total rows     Prev  Page 1 of 1  Next     Page Size  10 ⌄

**Findings**: The overall correlation between cost and satisfaction was basically zero (r≈ 0.0004). Across the three groups, the results were almost identical:

**Conclusion**: Therefore, higher spending doesn't lead to more satisfied patients. Even though hospitals spend more on average, the satisfaction score stays around 34.6 — almost flat.

That means；1.The extra money goes to things patients don't directly feel (like administrative or tech costs), or 2.Patient experience depends more on service quality than financial scale.

# Q3.How do States Compare in Satisfaction at Similar Cost Levels?

**Thought Process**

After finding that higher hospital spending doesn't necessarily lead to higher satisfaction, we wanted to see if location makes a difference — specifically, whether hospitals in some states consistently achieve higher patient satisfaction even when costs are similar.

To control for cost, we divided hospitals into three spending tiers using NTILE(3) — low, medium, and high — based on avg_total_payment. Within each cost level, we then compared the average satisfaction score across states. This allows us to ask: Do some states perform better at the same spending level?

We used the state_analysis_clean table, which contains state, avg_total_payment, and satisfaction_score, ensuring each hospital's cost and satisfaction were paired correctly.

> 🗄 df

```
 1 CREATE OR REPLACE TABLE hospital_data.state_cost_levels_Q3 AS
 2 WITH cost_groups AS (
 3     SELECT
 4         state,
 5         NTILE(3) OVER (ORDER BY avg_total_payment) AS cost_level_
 6         avg_total_payment,
 7         satisfaction_score
 8     FROM hospital_data.state_analysis_clean
 9 ),
10
11 with_labels AS (
12     SELECT
13         state,
14         CASE
15             WHEN cost_level_num = 1 THEN 'Low Cost'
16             WHEN cost_level_num = 2 THEN 'Mid Cost'
17             ELSE 'High Cost'
18         END AS cost_level,
19         avg_total_payment,
20         satisfaction_score
21     FROM cost_groups
22 )
23
24 SELECT
25     cost_level,
26     state,
27     ROUND(AVG(avg_total_payment), 2) AS avg_payment,
28     ROUND(AVG(satisfaction_score), 2) AS avg_satisfaction,
29     COUNT(*) AS hospital_count
30 FROM with_labels
31 GROUP BY cost_level, state
32 HAVING hospital_count >= 10
33 ORDER BY
34     CASE
35         WHEN cost_level = 'Low Cost' THEN 1
36         WHEN cost_level = 'Mid Cost' THEN 2
37         ELSE 3
38     END,
39     avg_satisfaction DESC;
40
```

✅ Completed.
✅ Completed.

| statement_type | job_id | location |
|---|---|---|
| CREATE_TABLE_AS_SELECT | job_miguULNq3XS8S5W4_2UcIgysI-fx | US |

1 total rows          Prev  Page 1 of 1  Next          Page Size  10 ∨

**Findings:** Even within the same cost level, patient satisfaction still varies slightly across states. For example, in the low-cost group, Hawaii (HI) and Wisconsin (WI) show higher satisfaction (around 88) compared to states like Louisiana (LA), which scores below 87. However, the differences are relatively small — for instance, Wisconsin's satisfaction rises from 87.8 in low-cost hospitals to 88.1 in high-cost hospitals, showing that more spending within a state does not make a major difference. This suggests that while state-level patterns exist, the overall satisfaction trend remains quite stable across cost levels.

**Conclusion:** Patient satisfaction shows small but consistent differences across states, even after controlling for cost levels. That means spending alone can't explain why some states achieve slightly better results. It's possible that factors like hospital culture, communication quality, and patient expectations play a larger role than financial input.

# Q4. What is the Cross-Sectional Relationship Between Cost and Satisfaction from 2016-2020?

🗄 df

```
 1 SELECT
 2   h.provider_id,
 3   h.hospital_name,
 4   ROUND(AVG(SAFE_CAST(p.`Patient Survey Star Rating` AS FLOAT64))
 5   ROUND(AVG(h.avg_total_payment), 2) AS avg_total_payment
 6 FROM hospital_data.hospital_satisfaction_join AS h
 7 LEFT JOIN hospital_data.patient_satisfaction_all AS p
 8   ON LPAD(CAST(h.provider_id AS STRING), 6, '0')
 9     = LPAD(CAST(p.facility_id AS STRING), 6, '0')
10   AND p.`Patient Survey Star Rating` IS NOT NULL
11 GROUP BY h.provider_id, h.hospital_name
12 ORDER BY avg_total_payment DESC;
```

✅ Completed.
✅ Completed.

| provider_id | hospital_name | avg_star_rating | avg_total_payment |
|---|---|---|---|
| 670116 | Wise Health System | 4.640000 | 76230.470000 |
| 260004 | Cooper County Community Hospital | <NA> | 66570.850000 |
| 450797 | Us Pain & Spine Hospital | <NA> | 65318.810000 |
| 450289 | Harris Health System | 2.890000 | 63362.970000 |
| 450674 | Womans Hospital Of Texas,the | 2.700000 | 63103.700000 |
| 050076 | Kaiser Foundation Hospital - San Francisco | 2.950000 | 62063.890000 |
| 450803 | United Memorial Medical Center | 1.670000 | 57857.660000 |
|  | Ronald Reagan U C |  |  |

**Finding:**

When we sort the table by Average Total Payment from high to low, the hospitals at the top (the most expensive) usually have low patient satisfaction scores. This shows up again and again, not just in one row, so it looks like paying more doesn't go hand in hand with happier patients in this view of the data.

**Conclusion:**

Therefore, patient satisfaction isn't closely linked to Average Total Payment. If anything, the highest costs often come with lower satisfaction.

# Q5. When Comparing Hospital Cost and Satisfaction, What are High-Value Providers?

🗄 df

```
1 CREATE OR REPLACE TABLE `ba775-fall25-b02.hospital_data.value_ho
2
3 WITH cost AS (
4   SELECT
5     LPAD(CAST(provider_id AS STRING), 6, '0') AS facility_id6,
6     hospital_name, city, state,
7     SAFE_CAST(avg_total_payment AS FLOAT64)   AS avg_total_payme
```

```
 8       year
 9     FROM `ba775-fall25-b02.hospital_data.hospital_costs_all`
10   ),
11   sat AS (
12     SELECT
13       LPAD(CAST(facility_id AS STRING), 6, '0')  AS facility_id6,
14       facility_name, city, state,
15       SAFE_CAST(`Patient Survey Star Rating` AS FLOAT64) AS star_r
16       SAFE_CAST(linear_mean_value AS FLOAT64)         AS linear_
17       SAFE_CAST(answer_percent AS FLOAT64)            AS answer_
18       CAST(answer_percent_footnote AS STRING)         AS answer_
19       measure_id, year
20     FROM `ba775-fall25-b02.hospital_data.patient_satisfaction_all`
21   ),
22
23   sat_16_20 AS (
24     SELECT
25       facility_id6, facility_name, state, city,
26       star_rating, linear_mean_value, answer_percent, answer_perce
27     FROM sat
28     WHERE year BETWEEN 2016 AND 2020
29   ),
30
31   sat_latest AS (
32     SELECT facility_id6, facility_name, state, city, answer_percen
33     FROM (
34       SELECT
35         facility_id6, facility_name, state, city, answer_percent_f
36         ROW_NUMBER() OVER (PARTITION BY facility_id6 ORDER BY year
37       FROM sat_16_20
38       WHERE answer_percent_footnote IS NOT NULL
39     )
40     WHERE rn = 1
41   ),
42
43   sat_avg AS (
44     SELECT
45       facility_id6,
46       AVG(star_rating)       AS avg_star,
47       AVG(linear_mean_value) AS avg_linear_mean,
48       AVG(answer_percent)    AS avg_answer_percent
49     FROM sat_16_20
50     GROUP BY facility_id6
51   ),
52
53   sat_fac AS (
54     SELECT
55       a.facility_id6,
56       l.facility_name,
57       l.state,
58       l.city,
59       a.avg_star,
60       a.avg_linear_mean,
61       a.avg_answer_percent,
62       l.answer_percent_footnote
```

```
63    FROM sat_avg AS a
64    LEFT JOIN sat_latest AS l USING (facility_id6)
65  ),
66
67
68  cost_16_20 AS (
69    SELECT
70      facility_id6, hospital_name, state, city, avg_total_payment,
71    FROM cost
72    WHERE year BETWEEN 2016 AND 2020
73  ),
74
75  cost_latest AS (
76    SELECT facility_id6, hospital_name, state, city
77    FROM (
78      SELECT
79        facility_id6, hospital_name, state, city, year,
80        ROW_NUMBER() OVER (PARTITION BY facility_id6 ORDER BY year
81      FROM cost_16_20
82    )
83    WHERE rn = 1
84  ),
85
86  cost_avg AS (
87    SELECT
88      facility_id6,
89      AVG(avg_total_payment) AS avg_payment
90    FROM cost_16_20
91    GROUP BY facility_id6
92  ),
93
94  cost_fac AS (
95    SELECT
96      a.facility_id6,
97      l.hospital_name,
98      l.state,
99      l.city,
100     a.avg_payment
101   FROM cost_avg AS a
102   LEFT JOIN cost_latest AS l USING (facility_id6)
103 ),
104
105
106 joined AS (
107   SELECT
108     s.facility_id6,
109     COALESCE(s.facility_name, c.hospital_name) AS hospital_name,
110     COALESCE(s.state, c.state)                 AS state,
111     COALESCE(s.city, c.city)                   AS city,
112     s.avg_star, s.avg_linear_mean,
113     s.avg_answer_percent, s.answer_percent_footnote,
114     c.avg_payment
115   FROM sat_fac s
116   JOIN cost_fac c USING (facility_id6)
117 ),
```

```
118
119
120 ranked AS (
121   SELECT
122     *,
123     PERCENT_RANK() OVER (ORDER BY avg_star)   AS star_pr,
124     PERCENT_RANK() OVER (ORDER BY avg_payment) AS cost_pr
125   FROM joined
126 )
127
128
129 SELECT
130   facility_id6            AS facility_id,
131   hospital_name, state, city,
132   ROUND(avg_star, 2)        AS avg_star_rating,
133   ROUND(avg_linear_mean, 2) AS avg_linear_mean,
134   ROUND(avg_payment, 2)     AS avg_total_payment,
135   ROUND(star_pr, 3)         AS star_percentile,
136   ROUND(cost_pr, 3)         AS cost_percentile,
137   ROUND(avg_answer_percent, 2) AS avg_answer_percent,
138   answer_percent_footnote
139 FROM ranked
140 ORDER BY avg_star_rating DESC, avg_total_payment ASC;
141
```

✅ Completed.
✅ Completed.

| statement_type | job_id | location |
|---|---|---|
| CREATE_TABLE_AS_SELECT | job_CpxOZDAA8zeFnjThFIbDD-BBaMK0 | US |

1 total rows          Prev  Page 1 of 1  Next          Page Size  10 ⌄

🗄 df

```
1 WITH cost AS (
2   SELECT
3     LPAD(CAST(provider_id AS STRING), 6, '0') AS facility_id6,
4     AVG(SAFE_CAST(avg_total_payment AS FLOAT64)) AS avg_payment_1(
5   FROM `ba775-fall25-b02.hospital_data.hospital_costs_all`
6   WHERE year BETWEEN 2016 AND 2020
7   GROUP BY facility_id6
8 ),
9 sat AS (
10   SELECT
11     LPAD(CAST(facility_id AS STRING), 6, '0') AS facility_id6,
12     AVG(SAFE_CAST(`Patient Survey Star Rating` AS FLOAT64)) AS av
13     AVG(SAFE_CAST(linear_mean_value AS FLOAT64))          AS av
14     AVG(SAFE_CAST(answer_percent AS FLOAT64))             AS av
15   FROM `ba775-fall25-b02.hospital_data.patient_satisfaction_all`
16   WHERE year BETWEEN 2016 AND 2020
17   GROUP BY facility_id6
18 ),
```

```
19 features AS (
20   SELECT s.facility_id6, s.avg_star_16_20, s.avg_linear_16_20, s.
21   FROM sat s JOIN cost c USING (facility_id6)
22 )
23 SELECT
24   CORR(avg_star_16_20, avg_payment_16_20)    AS r_star_vs_paymen
25   CORR(avg_star_16_20, avg_answer_pct_16_20) AS r_star_vs_answer_
26   CORR(avg_star_16_20, avg_linear_16_20)     AS r_star_vs_linear
27 FROM features;
```

✅ Completed.
✅ Completed.

| r_star_vs_payment | r_star_vs_answer_pct | r_star_vs_linear |
|---|---|---|
| -0.050830 | 0.035801 | 0.985655 |

1 total rows          Prev  Page 1 of 1  Next          Page Size  10  ∨

**Finding:**

We checked how different factors relate to patient satisfaction using a simple correlation test. The results showed that Average Total Payment and the Star Answer Percent both have very weak relationships with satisfaction—close to no link. In contrast, the HCAHPS Linear Mean Value moves strongly in the same direction as satisfaction, which means higher linear mean scores usually come with higher star ratings.

**Conclusion:**

Spending more or having a higher share of certain answers doesn't appear to make patients happier, based on this snapshot. The continuous HCAHPS Linear Mean Value seems to be a much better signal of patient experience than cost or answer percent.

# Q5-2. What is the Correlation between different features?

🗄 df

```
1 WITH corr_calc AS (
2   SELECT
3     corr(avg_satisfaction, avg_total_payment)  AS corr_cost,
4     corr(avg_satisfaction, bedsize_proxy)      AS corr_bedsize
5
6     corr(avg_satisfaction, own_government)     AS corr_own_gov
7     corr(avg_satisfaction, own_private)        AS corr_own_pri
8     corr(avg_satisfaction, own_nonprofit)      AS corr_own_nor
9
```

```
10      corr(avg_satisfaction, type_acute)           AS corr_type_ac
11      corr(avg_satisfaction, type_critical_access) AS corr_type_cr
12      corr(avg_satisfaction, type_childrens)       AS corr_type_ch
13
14      corr(avg_satisfaction, reg_NE)               AS corr_region_I
15      corr(avg_satisfaction, reg_MW)                AS corr_region_
16      corr(avg_satisfaction, reg_S)                AS corr_region_!
17      corr(avg_satisfaction, reg_W)                AS corr_region_\
18
19      corr(avg_satisfaction, sev_low)              AS corr_severit;
20      corr(avg_satisfaction, sev_medium)           AS corr_severit;
21      corr(avg_satisfaction, sev_high)             AS corr_severit;
22
23      corr(avg_satisfaction, emergency_num)        AS corr_emergen
24      corr(avg_satisfaction, ehr_num)              AS corr_ehr
25    FROM hospital_data.satisfaction_feature_matrix
26    WHERE avg_satisfaction IS NOT NULL
27 )
28
29 SELECT 'avg_total_payment' AS feature, corr_cost AS correlation FI
30 SELECT 'bedsize_proxy', corr_bedsize FROM corr_calc UNION ALL
31
32 SELECT 'ownership_government', corr_own_government FROM corr_calc
33 SELECT 'ownership_private', corr_own_private FROM corr_calc UNION
34 SELECT 'ownership_nonprofit', corr_own_nonprofit FROM corr_calc UI
35
36 SELECT 'type_acute', corr_type_acute FROM corr_calc UNION ALL
37 SELECT 'type_critical_access', corr_type_critical_access FROM cor
38 SELECT 'type_childrens', corr_type_childrens FROM corr_calc UNION
39
40 SELECT 'region_NE', corr_region_NE FROM corr_calc UNION ALL
41 SELECT 'region_MW', corr_region_MW FROM corr_calc UNION ALL
42 SELECT 'region_S', corr_region_S FROM corr_calc UNION ALL
43 SELECT 'region_W', corr_region_W FROM corr_calc UNION ALL
44
45 SELECT 'severity_low', corr_severity_low FROM corr_calc UNION ALL
46 SELECT 'severity_medium', corr_severity_medium FROM corr_calc UNI(
47 SELECT 'severity_high', corr_severity_high FROM corr_calc UNION A
48
49 SELECT 'emergency_services', corr_emergency FROM corr_calc UNION /
50 SELECT 'ehr_interoperability', corr_ehr FROM corr_calc
51
52 ORDER BY ABS(correlation) DESC;
```

✅ Completed.
✅ Completed.

| feature | correlation |
|---|---|
| region_MW | 0.192348 |
| ownership_private | -0.146957 |
| region_NE | -0.142459 |
| ownership_nonprofit | 0.116035 |
| region_W | -0.093123 |
| emergency_services | -0.031789 |
| severity_low | 0.030729 |
| severity_high | -0.024137 |
| avg_total_payment | -0.024117 |
| severity_medium | 0.019032 |

17 total rows     Prev Page 1 of 2 Next     Page Size 10 ⌄

🗄 df

```sql
 1 CREATE OR REPLACE TABLE hospital_data.feature_correlations AS
 2 WITH corr_calc AS (
 3 SELECT
 4   corr(avg_satisfaction, avg_total_payment) AS corr_cost,
 5   corr(avg_satisfaction, bedsize_proxy) AS corr_bedsize,
 6   corr(avg_satisfaction, own_government) AS corr_own_government,
 7   corr(avg_satisfaction, own_private) AS corr_own_private,
 8   corr(avg_satisfaction, own_nonprofit) AS corr_own_nonprofit,
 9   corr(avg_satisfaction, type_acute) AS corr_type_acute,
10   corr(avg_satisfaction, type_critical_access) AS corr_type_criti
11   corr(avg_satisfaction, type_childrens) AS corr_type_childrens,
12   corr(avg_satisfaction, reg_NE) AS corr_region_NE,
13   corr(avg_satisfaction, reg_MW) AS corr_region_MW,
14   corr(avg_satisfaction, reg_S) AS corr_region_S,
15   corr(avg_satisfaction, reg_W) AS corr_region_W,
16   corr(avg_satisfaction, sev_low) AS corr_severity_low,
17   corr(avg_satisfaction, sev_medium) AS corr_severity_medium,
18   corr(avg_satisfaction, sev_high) AS corr_severity_high,
19   corr(avg_satisfaction, emergency_num) AS corr_emergency,
20   corr(avg_satisfaction, ehr_num) AS corr_ehr
21 FROM hospital_data.satisfaction_feature_matrix
22 WHERE avg_satisfaction IS NOT NULL
23 )
24 SELECT 'avg_total_payment' AS feature, corr_cost AS correlation FI
25 SELECT 'bedsize_proxy', corr_bedsize FROM corr_calc UNION ALL
26 SELECT 'ownership_government', corr_own_government FROM corr_calc
27 SELECT 'ownership_private', corr_own_private FROM corr_calc UNION
```

```
28 SELECT 'ownership_nonprofit', corr_own_nonprofit FROM corr_calc UI
29 SELECT 'type_acute', corr_type_acute FROM corr_calc UNION ALL
30 SELECT 'type_critical_access', corr_type_critical_access FROM cor
31 SELECT 'type_childrens', corr_type_childrens FROM corr_calc UNION
32 SELECT 'region_NE', corr_region_NE FROM corr_calc UNION ALL
33 SELECT 'region_MW', corr_region_MW FROM corr_calc UNION ALL
34 SELECT 'region_S', corr_region_S FROM corr_calc UNION ALL
35 SELECT 'region_W', corr_region_W FROM corr_calc UNION ALL
36 SELECT 'severity_low', corr_severity_low FROM corr_calc UNION ALL
37 SELECT 'severity_medium', corr_severity_medium FROM corr_calc UNI(
38 SELECT 'severity_high', corr_severity_high FROM corr_calc UNION A
39 SELECT 'emergency_services', corr_emergency FROM corr_calc UNION /
40 SELECT 'ehr_interoperability', corr_ehr FROM corr_calc
41 ORDER BY ABS(correlation) DESC;
```

✅ Completed.
✅ Completed.

| statement_type | job_id | location |
|---|---|---|
| CREATE_TABLE_AS_SELECT | job_NS9-siSl3HL4r8K-05qGBOiNKuQf | US |

1 total rows     Prev  Page 1 of 1  Next          Page Size  10 ⌄

**Findings:**

None of the features show a strong linear correlation with total_charges. Most relationships are weak or moderate.

**Conclusion:**

This indicates that the charges are influenced by multiple independent hospital and patient factors.

# Q6. Which hospitals achieve high satisfaction at relatively low cost? (Constructing the Value Index (z-scores per year))

🗄 df

```
1 CREATE OR REPLACE TABLE hospital_data.value_index_Q6 AS
2 WITH yearly_stats AS (
3     SELECT
4         year,
5         AVG(satisfaction_score) AS mean_satisfaction,
6         STDDEV_POP(satisfaction_score) AS std_satisfaction,
7         AVG(avg_total_payment) AS mean_payment,
8         STDDEV_POP(avg_total_payment) AS std_payment
```

```
 9      FROM hospital_data.panel_clean
10      GROUP BY year
11  )
12
13  SELECT
14      p.provider_id,
15      p.hospital_name,
16      p.state_abbr AS state,
17      p.city,
18      p.year,
19      p.satisfaction_score,
20      p.avg_total_payment,
21
22      -- Z Scores
23      (p.satisfaction_score - y.mean_satisfaction) / y.std_satisfac
24      (p.avg_total_payment   - y.mean_payment)     / y.std_payment
25
26      -- Value index = z_sat - z_cost
27      ((p.satisfaction_score - y.mean_satisfaction) / y.std_satisfa
28      - (p.avg_total_payment   - y.mean_payment)     / y.std_paymen
29  FROM hospital_data.panel_clean p
30  FROM hospital_data.panel_clean p
31  JOIN yearly_stats y USING (year);
32
```

✅ Completed.
✅ Completed.

| statement_type | job_id | location |
|---|---|---|
| CREATE_TABLE_AS_SELECT | job_RLfqBtvjylj3PdY4yaDjiqS9FgKM | US |

1 total rows          Prev  Page 1 of 1  Next               Page Size  10  ∨

**Findings:** The hospitals with the highest value-index scores all operate at low costs, yet achieve high satisfaction. The lowest scoring hospitals, such as Westchester Medical Center had high costs and lower satisfaction rates.

**Conclusion**:

The value index analysis supports the idea that increases in price does not equate to increases in care level or patient satisfaction. For example, Westchester Medical Center had transactions over $600 000 with a range of satisfaction scores as low as 7. This may suggest that operational differences between the hospitals simultaneously create differences in care level and operation costs. It may also suggest that increased prices lower patient satisfaction as expectations are increased as well.

## Q7. What are the Hospitals in the Top 10% of the Value Index?

🗄 df

```
 1 CREATE OR REPLACE TABLE hospital_data.top_value_hospitals_Q7_fina
 2 WITH hospital_year AS (
 3   -- 1) Collapse to one row per hospital × year
 4   SELECT
 5     provider_id,
 6     ANY_VALUE(hospital_name) AS hospital_name,
 7     ANY_VALUE(state_abbr)    AS state,
 8     ANY_VALUE(city)          AS city,
 9     year,
10     AVG(satisfaction_score) AS satisfaction_score,   -- average sa
11     AVG(avg_total_payment)  AS avg_total_payment     -- average co
12   FROM hospital_data.panel_clean
13   GROUP BY provider_id, year
14 ),
15
16 stats AS (
17   -- 2) Compute yearly mean & sd at the hospital-year level
18   SELECT
19     year,
20     AVG(satisfaction_score)        AS mean_satisfaction,
21     STDDEV_POP(satisfaction_score) AS std_satisfaction,
22     AVG(avg_total_payment)         AS mean_payment,
23     STDDEV_POP(avg_total_payment)  AS std_payment
24   FROM hospital_year
25   GROUP BY year
26 ),
27
28 z_scores AS (
29   -- 3) Calculate z-scores and value_index (still one row per hosp
30   SELECT
31     h.*,
32     (h.satisfaction_score - s.mean_satisfaction) / s.std_satisfac
33     (h.avg_total_payment  - s.mean_payment)     / s.std_payment
34     (
35       (h.satisfaction_score - s.mean_satisfaction) / s.std_satisfa
36       - (h.avg_total_payment  - s.mean_payment)   / s.std_payment
37     ) AS value_index
38   FROM hospital_year h
39   JOIN stats s USING (year)
40 ),
41
42 ranked AS (
43   -- 4) Split each year into 10 deciles based on value_index → top
44   SELECT
45     *,
46     NTILE(10) OVER (PARTITION BY year ORDER BY value_index DESC)
```

```
47    FROM z_scores
48 ),
49
50 dedup AS (
51    -- 5) Safety: ensure only one row per hospital-year (in case of
52    SELECT
53      *,
54      ROW_NUMBER() OVER (
55        PARTITION BY provider_id, year
56        ORDER BY value_index DESC
57      ) AS rn
58    FROM ranked
59    WHERE decile = 1          -- only keep the top 10% hospitals in
60 )
61
62 SELECT
63    provider_id,
64    hospital_name,
65    state,
66    city,
67    year,
68    satisfaction_score,
69    avg_total_payment,
70    z_sat,
71    z_cost,
72    value_index
73 FROM dedup
74 WHERE rn = 1
75 ORDER BY year, value_index DESC;
76
```

✅ Completed.
✅ Completed.

| statement_type | job_id | location |
|---|---|---|
| CREATE_TABLE_AS_SELECT | job_3jIVDKWt_nBLz-ILrnhcYhtbG0LG | US |

1 total rows        Prev  Page 1 of 1  Next        Page Size  10  ⌄

**Findings:** The top hospitals all had high satisfaction and low costs. Hospitals such as Oklahoma Heart Hospital, Llc, Caldwell Memorial Hospital, Inc, and Citizens Medical Center were top scorers that appeared multiple times. The value-indexes range from 2.90 to 1.81, suggesting that there are still significant differences within the top decile. Hospitals with higher costs also tend to be lower value-index scorers, even though satisfaction scores are relatively high within the top decile.

**Conclusion:** Because an increase in price leads to an decrease in value-index score, it means that there is either a redundantly circular relation between price and satisfaction, or there is a different variable that was not considered such as hospital

size. As previously mentioned, hospital size may be one such candidate as smaller hospitals may be able to lower operating costs and provide personalized care in ways that are not scaleable to larger hospitals. Any further research should aim to clarify this relationship as well as explore other seemingly unrelated variables.

## Q8. How did each state and territories change from 2016-2020?

⛁ df

```
 1 CREATE OR REPLACE TABLE hospital_data.cost_sat_outliers_Q8 AS
 2 WITH base AS (
 3   SELECT
 4     state,
 5     year,
 6     ROUND(AVG(SAFE_CAST(`Patient Survey Star Rating` AS FLOAT64))
 7       AS avg_star_rating,
 8     ROUND(AVG(SAFE_CAST(`Patient Survey Star Rating` AS FLOAT64))
 9     ─ LAG(ROUND(AVG(SAFE_CAST(`Patient Survey Star Rating` AS FLO
10       OVER (PARTITION BY state ORDER BY year)
11       AS yoy_change
12   FROM hospital_data.patient_satisfaction_all
13   WHERE SAFE_CAST(`Patient Survey Star Rating` AS FLOAT64) IS NOT
14   GROUP BY state, year
15 )
16
17 SELECT
18   *,
19   MAX(CASE WHEN year = 2020 THEN avg_star_rating END)
20     OVER (PARTITION BY state)
21   ─
22   MAX(CASE WHEN year = 2016 THEN avg_star_rating END)
23     OVER (PARTITION BY state)
24       AS change_2016_to_2020
25 FROM base
26 ORDER BY state, year;
27
```

顯示隱藏的輸出內容

**Findings:**

Most states only saw modest improvements (.1–.3) from 2016–2020. It was also found that yoy_change fluctuated — some years showed positive movement, while others showed negative. There was also no significant jump in rating. No state saw an increase of +1.0 star from 2016–2020. There are meaningful differences in star ratings between states. For example, Nevada had some of the lowest ratings among

states and territories with almost 2 stars, while Wisconsin had some of the highest ratings with close to 4 stars.

**Conclusion:**

From 2016–2020, satisfaction across states saw gradual improvement; however, the improvement in each state was not large. The change was not linear, and annual fluctuation was common, with the direction of change shifting year to year.

It is more important to note that differences between states are much larger than changes within states over time. This suggests that patient satisfaction is shaped more by regional factors than by short-term operational changes.

Overall, the data shows that hospitals across the country are making incremental improvements, but long-standing regional disparities remain, and there needs to be broader, systemic state-level change.

# Q9. What are Cost-Satisfaction Outliers Using Quadrants?

**Thought Process**

In the previous questions, we mainly focused on the overall trend. But, the TA gave us some inspiration: Are there some hospitals that are particularly "strange"? For instance, they spend a lot of money but have very low satisfaction rates? Or conversely, are there hospitals that spend very little but provide excellent patient services? To ensure a fair comparison, we did not directly use the original "payment/satisfaction" figures. Instead, we adopted what TA suggested: converting both into z-scores. This approach enables us to place different hospitals on the same standard and assess how much they deviate from the average. We are specifically looking for two types of special hospitals:

• High-Cost Low-Satisfaction (spending a lot but with a poor experience)—— Hospitals that spend much more than average ($z\_cost > 1.5$) but have below-average satisfaction ($z\_sat < -1.0$)

• Low-Cost High-Satisfaction (spending less but with a great experience)—— Hospitals that spend much less ($z\_cost < -1.0$) but provide better patient experience ($z\_sat > 1.5$)

```
  🗄 df

    1 CREATE OR REPLACE TABLE hospital_data.cost_satisfaction_outliers_(
    2 WITH stats AS (
```

```
 3     SELECT
 4         AVG(avg_total_payment) AS mean_cost,
 5         STDDEV_POP(avg_total_payment) AS sd_cost,
 6         AVG(satisfaction_score) AS mean_sat,
 7         STDDEV_POP(satisfaction_score) AS sd_sat
 8     FROM hospital_data.panel_clean
 9 ),
10
11 zscores AS (
12     SELECT
13         p.*,
14         (p.avg_total_payment - s.mean_cost) / s.sd_cost AS z_cost
15         (p.satisfaction_score - s.mean_sat) / s.sd_sat AS z_sat
16     FROM hospital_data.panel_clean p
17     CROSS JOIN stats s
18 ),
19
20 final AS (
21     SELECT
22         *,
23         CASE
24             WHEN z_cost > 1.5 AND z_sat < -1.0 THEN 'High-Cost Lo
25             WHEN z_cost < -1.0 AND z_sat > 1.5 THEN 'Low-Cost Hig
26             ELSE 'Normal'
27         END AS outlier_flag
28     FROM zscores
29 )
30
31 SELECT *
32 FROM final
33 WHERE outlier_flag != 'Normal'
34 ORDER BY year, outlier_flag, z_cost DESC;
```

✅ Completed.
✅ Completed.

| statement_type | job_id | location |
|---|---|---|
| CREATE_TABLE_AS_SELECT | job_LV3VTrf69whbtKN91hCVM1GdVgic | US |

1 total rows      Prev   Page 1 of 1   Next      Page Size 10 ∨

**Findings:**

After tagging the hospitals, the two outlier groups actually showed some clear patterns. The High-Cost Low-Satisfaction hospitals spend above the national average, but their satisfaction scores are still below 0 (after standardization).

The Low-Cost High-Satisfaction hospitals were the opposite. They spent noticeably less, but their satisfaction was far above average. This part stood out after using z-scores, because the raw numbers didn't look that obvious before.

The outlier detection helped uncover which hospitals are "paying more but not performing" and which ones are "doing more with less."

**Conclusion:**

Facilities with higher charges tend to treat a larger share of severe or complex cases. The pattern differs across states and hospital groups, but cost levels generally rise with case difficulty.

# Q10 - How do Hospitals Rank in Cost and Satisfaction by Hospital Case-Mix Severity Tier?

**Thought Process**

In earlier questions, we compared hospitals without considering how difficult the cases were. But this could be misleading, because hospitals treating more severe patients naturally spend more. So for Q10, we grouped hospitals by DRG severity to make the comparison fair. The goal was to see whether higher case complexity actually leads to higher costs and higher patient satisfaction.

We calculated the average payment and satisfaction for each severity tier—High, Medium, and Low. If higher spending reflected better care, the High-severity group should have higher satisfaction or a better value index. But the results showed the opposite: the High-severity tier had the highest costs but similar or even lower satisfaction. Meanwhile, the Low-severity tier provided the best experience at the lowest cost.

---

🗄 df

```sql
1 CREATE OR REPLACE TABLE hospital_data.hospital_drg_panel_tier_Q10
2 SELECT
3   *,
4   CASE
5     WHEN Avg_Tot_Pymt_Amt < 5000 THEN 'Low cost'
6     WHEN Avg_Tot_Pymt_Amt BETWEEN 5000 AND 15000 THEN 'Mid cost'
7     ELSE 'High cost'
8   END AS cost_tier
9 FROM hospital_data.hospital_drg_panel;
10
```

✅ Completed.
✅ Completed.

| statement_type | job_id | location |
|---|---|---|
| CREATE_TABLE_AS_SELECT | job_Uezaj1qzFjPENSRYi64Ztg34ERF1 | US |

1 total rows      Prev  Page 1 of 1  Next      Page Size 10 ∨

**Findings:**

As case complexity increases, hospital costs rise sharply, yet patient satisfaction remains largely unchanged, with the highest-severity group even showing slightly lower star ratings. In contrast, the lowest-severity group has the lowest costs but the highest satisfaction, suggesting that tougher cases do not lead to better perceived care. The value delivered per $1,000 spent drops substantially in high-severity hospitals—almost half of that in the lowest tier. Overall, higher costs reflect treating more difficult cases rather than delivering higher-quality patient experience.

**Conclusion:**

Higher case complexity is associated with higher payments, while patient satisfaction does not increase in the same way. Hospitals managing more complex cases face higher cost intensity without a clear pattern in satisfaction scores.

# Final Conclusion

Several conclusions in EDA questions suggests that higher costs does not equate to higher patient satisfaction or higher value in care. In fact, higher costs often led to lower patient perceptions of care value. Whether examined across time or location, it held true that cost and satisfaction moved independently when examining all of the hospital within the dataset.

Smaller regional hospital such as Citizens Medical Center, Oklahoma Heart Hospital, Caldwell Memorial Hospital, and Rockcastle consistently demonstrated high levels of care for lower costs; suggesting that high value-index care is a repeatable process. Whether or not it can be scaled for larger hospitals remains unclear.

Analysis on the state level showed that there were little differences in care quality within the different cost tiers. Meanwhile, the overall correlation between cost and satisfaction ($r \approx 0.0004$) confirms the idea that care can be provided at afforable costs.

The next steps of this project should aim to cover the project shortcomings. For example, treatment for different illnesses were compared on equal footing. A

treatment method that has a low success rate and high price could be compared to