

# CPSC 368 Research Project

Technology and Healthcare - by Aarav, Pushya, Suryansh

## 1. Introduction

### Research Question

We will explore the following Research Question through this project:

**How did socioeconomic and demographic factors influence telehealth adoption patterns during the COVID-19 pandemic in the United States (2020-2022)**

The impact on the telehealth adoption patterns will be explored via the following factors:

1. **Racial** and **ethnic disparities** in telehealth utilization rates
2. **Geographic differences** in adoption patterns (privileged and under-privileged)
3. The relationship between **educational attainment** and telehealth usage

Note: In our study, we will create consistent geographic units for rural/urban classification: Our assumption will classify “metropolitan” cities and “urban” areas as one classification which will represent more connected and privileged residential areas. Similarly, we will classify “non-metropolitan” and “rural” areas as one classification to represent less connected and generally under / less privileged residential areas.

### Motivation for the Research Question:

The COVID-19 pandemic catalyzed a dramatic shift toward telehealth services, potentially exacerbating existing healthcare disparities in the United States. While telehealth offered a crucial lifeline for continued healthcare access, preliminary studies suggest uneven adoption patterns across different demographic groups. We aim to use this research to shine a brighter light on the following areas:

- Quantify disparities in telehealth adoption across different populations
- Identify specific barriers to telehealth access for underserved communities
- Contribute to evidence-based policy recommendations for more equitable healthcare delivery systems
- Examine how existing healthcare inequities were either mitigated or amplified by telehealth expansion

## 2. Related work

The COVID-19 pandemic accelerated telehealth adoption in the U.S., but research highlights persistent disparities shaped by socioeconomic and demographic factors. **Racial and ethnic inequities** emerged as a critical concern, with studies showing Black patients were more likely to rely on audio-only telehealth compared to White patients, who disproportionately used video visits (*Luo, Jake, et al.*). Latino and Asian populations faced compounded barriers, including language access and lower broadband connectivity, leading to reduced overall telehealth engagement <sup>25</sup>. These disparities reflect systemic inequities in digital literacy and technology access, which risk exacerbating existing health gaps, particularly in chronic disease management and preventive care (*Chen, Evan M., et al.*). For example, non-English speakers were significantly less likely to complete video visits (OR 0.49), underscoring the need for culturally tailored telehealth solutions (*Chen, Evan M., et al.*).

**Geographic and socioeconomic disparities** further stratified telehealth access. ZIP code-level analyses revealed that areas with higher college education rates had stronger video visit adoption (coefficient 1.41), while low-income regions relied more on telephone-based care (*Luo, Jake, et al.*). Rural-urban divides were initially pronounced due to limited broadband infrastructure, though later studies (2022) suggested narrowing gaps as telehealth became normalized, with rural patients increasingly using messaging platforms for care (*Spaulding, Erin M., et al.*). However, regional inequities persisted; Midwestern states lagged in telehealth utilization (aPR 0.65), likely due to policy variability and provider density (*Spaulding, Erin M., et al.*). These patterns highlight how structural factors like education, income, and regional resource allocation create uneven access to high-quality telehealth, disadvantaging marginalized communities.

**Educational attainment and technology access** emerged as pivotal determinants of telehealth use. College-educated individuals were 24% more likely to engage in telehealth than those with lower education, reflecting advantages in navigating digital platforms and understanding insurance coverage (*Spaulding, Erin M., et al.*). Lower educational attainment correlated with perceived risks (e.g., privacy concerns, technical complexity), reducing adoption willingness (*Wu T., Ho C.*). Importantly, socioeconomic status intersected with race and geography; for instance, Black patients in low-income urban neighborhoods faced dual barriers of limited tech access and systemic underinsurance (*Chen, Evan M., et al.*). These findings underscore telehealth's dual role: it can democratize care but also reinforce inequities if not paired with investments in digital literacy and infrastructure (*Spaulding, Erin M., et al.*).

This research is vital because telehealth is now a permanent care modality. Without addressing these disparities, the shift to virtual care risks deepening health inequities, particularly for chronic disease management and preventive services. Studies emphasize the need for policies targeting

broadband expansion, multilingual platforms, and provider training to ensure equitable access (Chen, Evan M., et al).

### 3. Datasets

We will get the publicly available datasets from the US government for the adoption of telehealth services which could be broken down by different attributes of the communities such as race, gender and educational background. Both of the data sources are accredited US government sites and focus on different issues and we will combine data from both of them to notice the impact of COVID-19 and accessibility to healthcare (via telehealth) for different groups.

1. RANDS (Research and Development Survey) data:  
*Access and use of telemedicine during COVID-19 | HealthData.gov.* (2021, February 25).  
[https://healthdata.gov/dataset/Access-and-Use-of-Telemedicine-During-COVID-19/c835-etjt/about\\_data](https://healthdata.gov/dataset/Access-and-Use-of-Telemedicine-During-COVID-19/c835-etjt/about_data)
2. Medicare Telehealth Trends dataset:  
Data.gov. (2025, February 3). *U.S. Department of Health & Human Services - Medicare Telehealth Trends.*  
[https://catalog.data.gov/dataset/medicare-telemedicine-snapshot?utm\\_source=chatgpt.com](https://catalog.data.gov/dataset/medicare-telemedicine-snapshot?utm_source=chatgpt.com)

Trustworthiness of the data sources:

- The RANDS (Research and Development Survey) data explicitly acknowledges its experimental nature and limitations, which demonstrates transparency. It uses probability-sampled commercial survey panels and has documented its methodology in technical notes. We notice such limitations and focus on healthcare datasets from several other publicly available health datasets for different countries.
- The Medicare Telehealth Trends dataset comes from actual Medicare service utilization records, providing concrete behavioural data rather than self-reported information.

Known limitations of the data sources:

- RANDS documentation openly acknowledges potential biases from different response patterns and sampling frames
- RANDS data has increased variability due to lower sample sizes compared to traditional NCHS surveys
- Medicare data primarily represents older Americans and those with disabilities, potentially underrepresenting other populations

## Data cleaning

We used a jupyter notebook to clean, analyse and visualise (EDA) the data. Present in the attached jupyter notebook file.

In both datasets, aggregated categories such as “Total” are being removed as it was checked in “Midway checkpoint” and preliminary research that the the subgroups forms a total set of the larger group.

Furthermore, in both datasets Unknowns or NULLS are replaced with “Unknown” and or “Other” category for the attribute as this will be used to classify missing or unknown data.

For cleaning the ***“Access\_and\_Use\_of\_Telemedicine\_During\_COVID-19.csv”***:

1. Remove all the irrelevant columns such as “Suppression, Significant 1 and Significant 2”, which were given in the Research proposal.
2. We will also convert 'High school graduate or less' to 'highschool or less', 'Some college' to 'college' and "Bachelor's degree or above" to 'bachelor or above' for convenience.
3. We will also ignore the “Urbanisation” (different from point b) and other subgroups besides the one mentioned above as it does not provide any useful information or any further subgrouping required for our research question.
4. Finally, add a column called “size” that shows not only the size instead of the percentage of the population which will make further calculation easier. This will be calculated per SubGroup of the relevant Groups per Indicator.
5. The dataset is **unpivoted** to ensure that it can be matched and joined with the second dataset which includes converting ‘Group’ turn into columns and their values be taken by ‘Subgroups’. This would also include splitting Response type: ['Yes', 'No', 'Do not know', 'No usual place of care', 'No telemedicine available']

For the *"TIMEDTREND\_PUBLIC\_241126.csv"*:

1. We first filter only the COVID years which is 2020-2022 (inclusive)
2. Look at all the Summary statistics for the numerical rows:
  - a. The summary statistics for all the numerical rows make sense as they are non-zero and meet expectations. This is done by comparing it to the values seen in the csv and don't indicate any major outliers.
3. Look at all the NULL values in the rows:
  - a. Almost all columns have complete data and are not NULL so no rows are removed.
4. We then aggregate the data for our relevant purposes which include 'Bene\_Race\_Desc', 'Bene\_RUCA\_Desc', 'Bene\_Mdcd\_Mdcr\_Enrl\_Stus' for each numerical column. While aggregating, every absolute value is taken as a sum except percentage which is taken as a mean as the summation would not make much sense for a relative metric.
5. We also renamed the column to 'Race', 'Urbanization', and 'Enrollement\_Status' for ease of understanding.
6. As during the Midway checkpoint, TA reminded us that one table needs to have a foreign key to the other table, we have decided to map NaN missing values to unknown as they fit there the best and during the analysis, these values are most likely to be ignored or treated as missing in both datasets.

After cleaning the dataset, the following Database Schema was created:

For the cleaned version of the *"Access\_and\_Use\_of\_Telemedicine\_During\_COVID-19.csv"* dataset:

```
CREATE TABLE telemedicineprovider (  
  round INT,  
  indicator VARCHAR(255),  
  race VARCHAR(255),  
  urbanization VARCHAR(255),  
  education VARCHAR(255),  
  samplesize FLOAT,  
  sizevalue FLOAT,  
  
  -- Percent columns  
  percent_do_not_know FLOAT,  
  percent_no FLOAT,
```

```

percent_no_telemedicine_available FLOAT,
percent_no_usual_place_of_care FLOAT,
percent_yes FLOAT,

-- Standard Error columns
standarderror_do_not_know FLOAT,
standarderror_no FLOAT,
standarderror_no_telemedicine_available FLOAT,
standarderror_no_usual_place_of_care FLOAT,
standarderror_yes FLOAT,

-- Foreign key to beneficiarydata
FOREIGN KEY (race, urbanization)
  REFERENCES beneficiarydata(race, urbanization)
  ON DELETE SET NULL
  ON UPDATE CASCADE,

-- Primary Key
PRIMARY KEY (
  round, indicator, race, urbanization, education
)
);

```

For the cleaned version of the “***TIMEDTREND\_PUBLIC\_241126.csv***” dataset file:

```

CREATE TABLE beneficiarydata (
  race VARCHAR(255),
  urbanization VARCHAR(255),
  enrollment_status VARCHAR(255),
  total_bene_th_elig FLOAT,
  total_partb_enrl FLOAT,
  total_bene_telehealth FLOAT,
  pct_telehealth FLOAT,
  PRIMARY KEY (
    race, urbanization,
  )
);

```

## 4. Methodology

In this project, we aimed to explore telehealth usage patterns during the COVID-19 pandemic by analyzing two datasets: the RANDS survey, which included self-reported telemedicine usage, and Medicare telehealth utilization data. We specifically examined telehealth adoption rates across racial groups, urban versus rural regions, and different educational backgrounds. To facilitate analysis, we first cleaned and structured the data in Python. After cleaning, we loaded the datasets into an Oracle SQL database by dynamically generating INSERT statements.

Once the data was stored in the database, we used targeted SQL queries with specific filters and JOIN operations to combine relevant information from both datasets efficiently. For instance, one of our key queries combined telehealth usage by race and urbanization level with Medicare beneficiary data:

```
SELECT tp.race, tp.urbanization, AVG(tp.percent_yes), bd.total_bene_telehealth
FROM telemedicineprovider1 tp
JOIN beneficiarydata1 bd ON tp.urbanization = bd.urbanization
WHERE tp.indicator = 'Scheduled one or more telemedicine appointments'
GROUP BY tp.race, tp.urbanization, bd.total_bene_telehealth;
```

The research methodology employed hypothesis testing through carefully designed SQL queries executed in SQLPlus. Results were systematically captured using the “spool <filename>” and then “spool off” after the query, generating consistently formatted text files. These output files were subsequently parsed using a custom Python script that extracted the relevant data elements. The structured extraction approach facilitated comprehensive hypothesis testing on the accumulated data, enabling statistical analysis and evidence-based evaluation of the research questions under investigation.

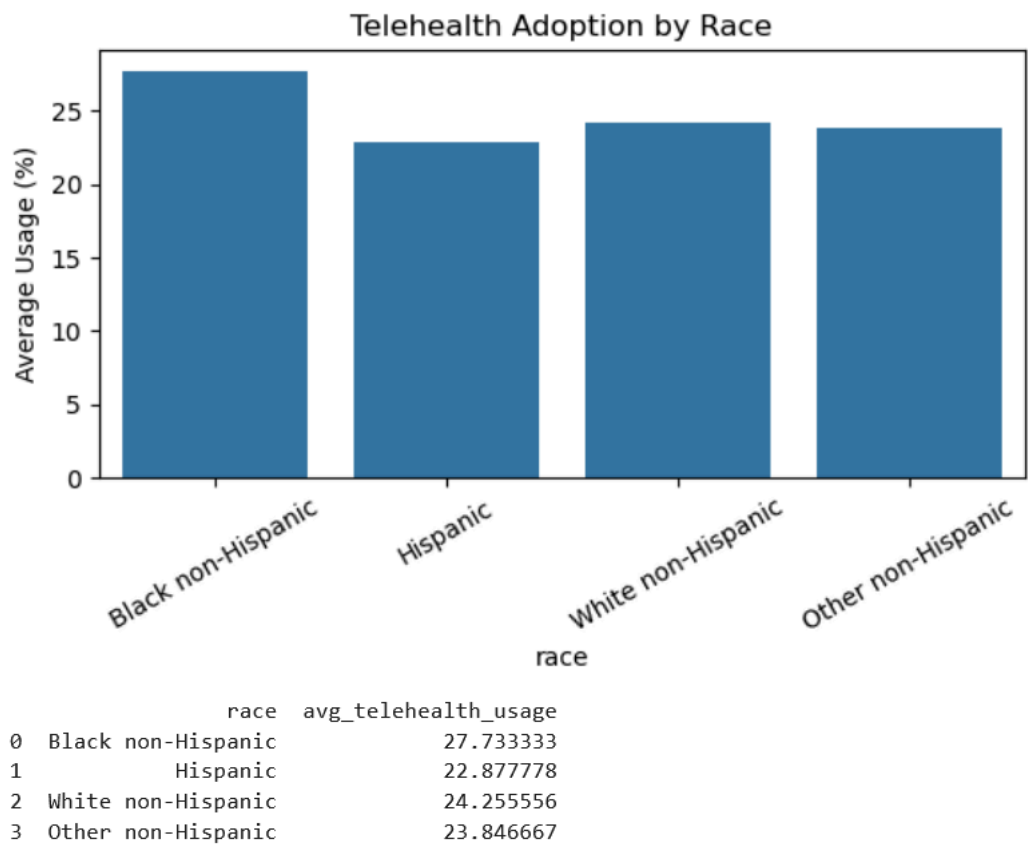
## Hypothesis testing

Our research was driven by three main hypotheses. Firstly, we hypothesized there would be significant disparities in telehealth adoption rates among racial groups compared to their respective shares of the overall population. Secondly, we sought to identify whether telehealth usage differed significantly between urban and rural populations. This relied more on percentage differences and EDA for inference, as we did not have access to more parameters such as the standard error of the population. Lastly, we hypothesized educational attainment would significantly influence telehealth adoption rates.

For the racial comparison, we conducted a series of z-tests. Specifically, we tested whether the proportion of telehealth users in each racial group significantly differed from their

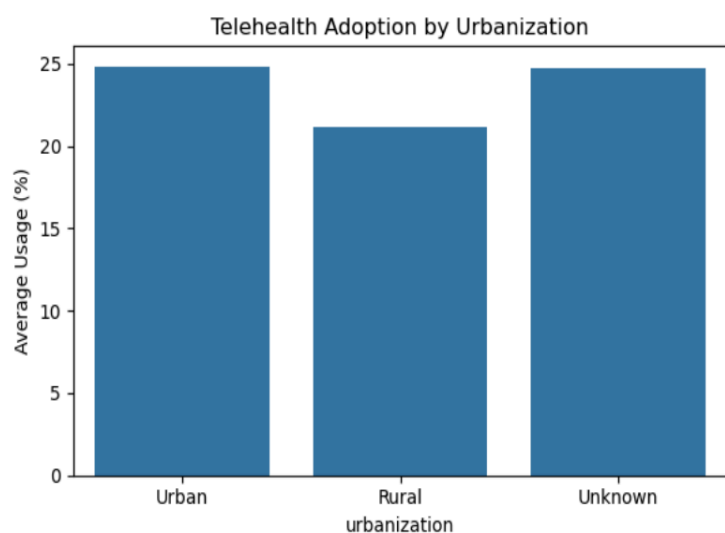
representation in the U.S. population according to the 2020 census. Our null hypothesis was that there would be no difference between telehealth usage proportions and population proportions. The z-tests were possible as we had access to the standard error for all the relevant racial demographics, as well as a large enough sample size for the population test to be relevant.

## 5. Results and Discussion

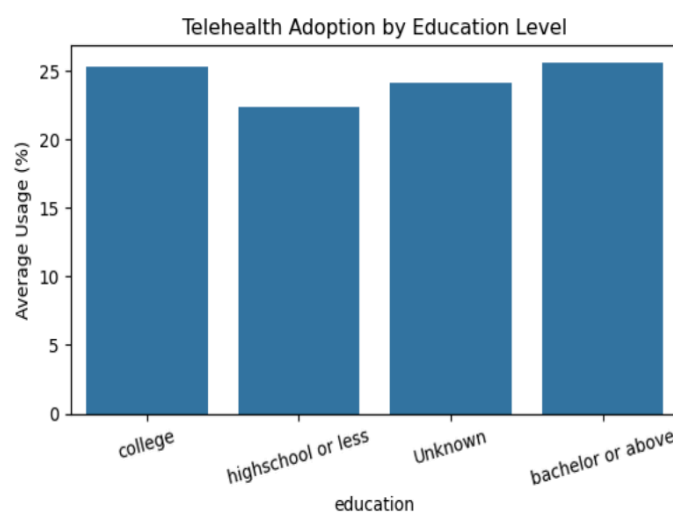


Our analysis revealed significant racial disparities in telehealth usage compared to population proportions. Notably, Black non-Hispanic respondents had significantly higher telehealth usage (26.00%) than their population share (12.4%), with a z-score of 5.30 and a highly significant p-value ( $<0.0001$ ). Similarly, the Other non-Hispanic group also showed significantly higher usage (18.63%) compared to their smaller population proportion (7.1%), with a z-score of 4.68. In contrast, White non-Hispanic respondents showed a significantly lower telehealth usage (21.97%) relative to their substantial population share (61.6%), reflected in a large negative z-score of -33.03. Hispanic individuals' telehealth usage (20.27%) was roughly proportional to their population share (18.9%), with no statistically significant difference (p-value = 0.5345).





urbanization	avg_telehealth_usage
0 Urban	24.811111
1 Rural	21.166667
2 Unknown	24.687037



education	avg_telehealth_usage
0 college	25.311111
1 highschool or less	22.388889
2 Unknown	24.168889
3 bachelor or above	25.555556

When analyzing telehealth adoption by urbanization, we observed slightly higher usage rates in urban areas (~24.8%) compared to rural areas (~21.1%). However, this difference was not statistically significant due to limitations in the data. Additionally, education level analyses indicated higher telehealth adoption among those with higher educational attainment, though statistical significance could not be conclusively determined from the available data due to small sample sizes in some education subgroups.

These results highlight critical disparities in telehealth adoption during the pandemic, reflecting broader patterns of healthcare access and equity. These results may be limited in their scope due to the limited data and as it is based on secondary evidence (lack of knowledge on the details of the surveys and tests conducted), it still agrees with the consensus of the other results measured in relevant studies (Chen, Evan M., et al).

## 6. Discussion of Results

Our analysis indicates that telehealth adoption during the COVID-19 pandemic varied considerably among different demographic groups, reflecting longstanding systemic inequalities. For instance, usage rates were significantly higher among Black patients and individuals categorized as Other non-Hispanic. At first glance, this pattern might seem counterintuitive. However, these higher rates likely reflect greater healthcare needs and fewer alternative options in these communities during the pandemic (Chen et al., 2021). Limited access to in-person services, higher disease burden, and existing disparities in healthcare infrastructure may have driven heavier reliance on telehealth, especially audio-only services (Luo et al., 2021).

In contrast, White non-Hispanic individuals—despite their large population share—used telehealth less than expected. This underrepresentation may be explained by better access to in-person care or greater hesitancy to transition to remote healthcare. Meanwhile, Hispanic individuals used telehealth at rates roughly in line with their demographic share. This balance could indicate the presence of subtle cultural or technological barriers, such as language access or broadband availability, which have been identified in prior research (Chen et al., 2021).

Though we did not find statistically significant differences between urban and rural populations, rural communities still appear to face distinct challenges. These often include limited broadband coverage and inadequate digital infrastructure, which are well-documented barriers to telehealth (Spaulding et al., 2024). Our findings on education echoed trends reported in the literature—individuals with higher levels of educational attainment were more likely to engage in telehealth, likely due to stronger digital literacy and greater confidence navigating healthcare technologies (Wu & Ho, 2021).

While our conclusions align with these established patterns, we recognize that our analysis was limited by the use of secondary data and uneven subgroup representation. Nonetheless, this study reinforces the broader understanding that digital health tools, if not implemented equitably, risk amplifying rather than alleviating healthcare disparities. Addressing these disparities requires thoughtful intervention tailored to the needs of underserved communities (Spaulding et al., 2024).

## 7. Future work

With additional time and resources, this project could expand in multiple directions. One important step would be to disaggregate telehealth usage by modality—examining differences between video visits, audio-only calls, and asynchronous communication. Prior studies have shown that access to video visits, in particular, is unequally distributed and may reflect broader structural barriers (Luo et al., 2021). Understanding which modes are most used—and by whom—could inform efforts to improve quality and accessibility of care.

Another key extension would be to incorporate more demographic variables into the analysis. Information such as household income, insurance coverage, language proficiency, and employment status could provide a clearer picture of how socioeconomic factors intersect to shape telehealth usage (Chen et al., 2021). Additionally, working with post-2022 data would allow us to track whether the disparities observed during the pandemic have persisted or changed as the healthcare system has evolved.

We also see potential for predictive modeling. Using machine learning algorithms trained on the cleaned datasets, we could estimate the likelihood of telehealth adoption across different demographic profiles. These models could help healthcare providers and policymakers identify communities with low expected uptake and tailor interventions accordingly. This approach could support more targeted outreach and better allocation of resources.

Lastly, we recommend expanding the scope of analysis to include international comparisons. Countries such as Canada, the UK, and Australia have also experienced rapid telehealth expansion during the pandemic. Studying how these nations addressed issues of access and equity could offer valuable lessons. Cross-national comparisons could help U.S. policymakers adopt more effective, evidence-based approaches to ensure long-term, equitable access to telehealth services (Spaulding et al., 2024).

## 8. References

Luo, Jake, et al. “Telemedicine Adoption During the COVID-19 Pandemic: Gaps and Inequalities.” *Applied Clinical Informatics*, vol. 12, no. 04, Aug. 2021, pp. 836–44. <https://doi.org/10.1055/s-0041-1733848>.

Chen, Evan M., et al. “Socioeconomic and Demographic Disparities in the Use of Telemedicine for Ophthalmic Care During the COVID-19 Pandemic.” *Ophthalmology*, vol. 129, no. 1, July 2021, pp. 15–25. <https://doi.org/10.1016/j.ophtha.2021.07.003>.

Spaulding, Erin M., et al. “Prevalence and Disparities in Telehealth Use Among US Adults Following the COVID-19 Pandemic: National Cross-Sectional Survey.” *Journal of Medical Internet Research*, vol. 26, May 2024, p. e52124. <https://doi.org/10.2196/52124>.

Wu T., Ho C. “Telemedicine Adoption During the COVID-19 Pandemic: Gaps and Inequalities.” *Applied Clinical Informatics*, vol. 12, no. 04, Aug. 2021, pp. 836–44. <https://doi.org/10.1055/s-0041-1733848>.

## 9. AI Declaration

We have not used any AI or GenAI tools in this research project so far, All the steps have been completed with the materials learned in class or researched online (Links provided above).

- Aarav, Pushya, Suryansh