# 🩺 Web Scraping Project: U.S. Infusion Centers Locator

## 📌 Overview

This project focuses on web scraping healthcare facility data from the [Infusion Center Locator](), a public-facing directory that lists infusion centers across the United States. Infusion centers provide essential outpatient services for patients with autoimmune diseases, cancer, and chronic illnesses.

The goal of this project is to automate the extraction of all center names and their physical addresses into a structured format (Excel/CSV), enabling future use in analytics, mapping, and healthcare accessibility research.

---

## 🎯 Objectives

- Scrape the full list of infusion centers from the locator website.
- Extract the **center name** and **address** for each location.
- Save the collected data into an Excel file for easy access and analysis.
- Build a clean, modular scraping solution suitable for a data science portfolio.

---

## 🛠️ Tools and Technologies

- **Python** – For scripting and automation
- **Selenium** – To automate interaction with the website and extract data
- **Pandas** – To store, manipulate, and export the collected data
- **Jupyter Notebook** – For documenting the entire process interactively

---

## 📁 Output

A downloadable Excel file named `infusion_centers_all.xlsx` containing:

- `Center Name`
- `Address`

---

## ✅ Why This Project?

- Demonstrates real-world web scraping skills
- Showcases automation of dynamic web content using Selenium
- Builds a portfolio-worthy dataset in the healthcare domain
- Provides potential use cases in public health analysis and business intelligence

# 💼 Section 2: Tools and Library Imports

In this project, we will use the following Python libraries:

- **Selenium**: Automates browser interactions (used to navigate the website, click buttons, and scrape dynamic content).
- **Pandas**: Used for organizing the scraped data into a structured DataFrame and exporting it to Excel.
- **Time**: Adds delays to allow dynamic content to load before scraping.

  > ⚠️ Make sure these libraries are installed before running the code.

### ✅ Installation (if not already installed)

You can install these libraries using pip:

```
pip install selenium pandas openpyxl
```

```python
# 📦 Import necessary libraries
from selenium import webdriver
from selenium.webdriver.chrome.options import Options
from selenium.webdriver.common.by import By

import pandas as pd
import time
```

## 🗣️ Section 3: Setting Up Selenium WebDriver (with webdriver-manager)

Instead of manually downloading and managing the ChromeDriver, we'll use `webdriver-manager` to automatically fetch the correct version. This makes the setup smoother and platform-independent.

```
In [37]:
# 🛠️ Set up headless Chrome WebDriver
chrome_options = Options()
chrome_options.add_argument("--headless")          # Run browser in headless mode (no GUI)
chrome_options.add_argument("--disable-gpu")        # Disable GPU rendering
chrome_options.add_argument("--no-sandbox")         # Required for some Linux environments

# Launch the browser
driver = webdriver.Chrome(options=chrome_options)

# Confirm it's working
print("✅ Chrome WebDriver launched successfully.")
```

✅ Chrome WebDriver launched successfully.

```
In [38]:
# 📦 Install webdriver-manager if not already installed
!pip install selenium webdriver-manager --quiet
```

```
In [39]:
# 🛠️ Set up WebDriver using webdriver-manager
from selenium import webdriver
from selenium.webdriver.chrome.service import Service
from selenium.webdriver.chrome.options import Options
from webdriver_manager.chrome import ChromeDriverManager
from selenium.webdriver.common.by import By
import pandas as pd
import time

# Configure Chrome to run headless
chrome_options = Options()
#chrome_options.add_argument("--headless")
chrome_options.add_argument("--disable-gpu")
chrome_options.add_argument("--no-sandbox")

# Initialize driver using ChromeDriverManager
driver = webdriver.Chrome(service=Service(ChromeDriverManager().install()), options=chrome_options)

print("✅ Chrome WebDriver launched successfully (via webdriver-manager).")
```

✅ Chrome WebDriver launched successfully (via webdriver-manager).

## 🌎 Section 4: Load the Website and Trigger the Search

Now that the browser is running, we'll:

1. Navigate to the Infusion Center Locator website.
2. Wait for the page to fully load.
3. Click the **Search** button to display all available centers across the U.S.

> We use `time.sleep()` to ensure the page loads properly before we interact with elements.

```
In [41]:
from selenium.webdriver.support.ui import WebDriverWait
from selenium.webdriver.support import expected_conditions as EC

# 🌐 Reload the site
driver.get("https://locator.infusioncenter.org/")
time.sleep(3)

try:
    # 🕵️ Wait until the backdrop disappears
    WebDriverWait(driver, 20).until(
        EC.invisibility_of_element_located((By.CLASS_NAME, "MuiBackdrop-root"))
    )

    # 🟢 Wait until the 'Search' button is clickable and enabled
    search_button = WebDriverWait(driver, 20).until(
        EC.element_to_be_clickable((By.XPATH, "//button[contains(text(), 'Search')]"))
    )

    # ✅ Click the button
    search_button.click()
    print("✅ 'Search' button clicked successfully.")
    time.sleep(10)  # Wait for results to load

except Exception as e:
    print("❌ Error:", e)
```

✅ 'Search' button clicked successfully.

# 📜 Section 5: Scroll to Load All Results

Since the website loads more centers as you scroll, we automate scrolling to the bottom repeatedly until no new centers load.

This ensures we collect the full list before extraction.

```python
In [44]:  # Scroll to load all centers dynamically
          prev_count = -1
          while True:
              driver.execute_script("window.scrollTo(0, document.body.scrollHeight);")
              time.sleep(3)

              containers = driver.find_elements(By.CSS_SELECTOR, "div.locator-result-container")
              current_count = len(containers)
              print(f"Centers loaded: {current_count}")

              if current_count == prev_count:
                  print("✅ All centers loaded.")
                  break
              prev_count = current_count

          # Extract name and address
          data = []
          for idx, container in enumerate(containers):
              try:
                  # Use WebDriverWait to ensure elements are present inside this container
                  WebDriverWait(container, 2).until(
                      EC.presence_of_element_located((By.CLASS_NAME, "locator-result-name"))
                  )

                  # Try to get the name from either desktop or mobile view
                  name = ""
                  try:
                      name = container.find_element(
                          By.XPATH, ".//div[contains(@class, 'hidden') and contains(@class, 'sm:block')]//div[contains(@class, 'locator-
                      ).text.strip()
                  except:
                      try:
                          name = container.find_element(
                              By.XPATH, ".//div[contains(@class, 'sm:hidden')]//div[contains(@class, 'locator-result-name')]"
                          ).text.strip()
                      except:
                          name = ""

                  # Extract address
                  address = ""
                  try:
                      address_div = container.find_element(By.CLASS_NAME, "locator-result-address")
                      address_lines = address_div.find_elements(By.TAG_NAME, "div")
                      street = address_lines[0].text.strip() if len(address_lines) > 0 else ""
                      city_state = address_lines[1].text.strip() if len(address_lines) > 1 else ""
                      address = f"{street}, {city_state}" if street and city_state else ""
                  except:
                      address = ""

                  if name and address:
                      data.append({"Center Name": name, "Address": address})

                  else:
                      print(f"⚠️ Skipped at index {idx} | Name: {name} | Address: {address}")

              except Exception as e:
                  print(f"❌ Error at index {idx}: {e}")

          # Convert and preview
          import pandas as pd
          df = pd.DataFrame(data)
          print("\n📊 Preview of Top 5:")
          print(df.head())
          print(f"\n🧊 Total Valid Centers Extracted: {len(df)}")
```

```
Centers loaded: 960
Centers loaded: 960
✅ All centers loaded.
⚠️ Skipped at index 240 | Name:  | Address:
⚠️ Skipped at index 241 | Name:  | Address:
⚠️ Skipped at index 242 | Name:  | Address:
⚠️ Skipped at index 243 | Name:  | Address:
⚠️ Skipped at index 244 | Name:  | Address:
⚠️ Skipped at index 245 | Name:  | Address:
⚠️ Skipped at index 246 | Name:  | Address:
⚠️ Skipped at index 247 | Name:  | Address:
⚠️ Skipped at index 248 | Name:  | Address:
⚠️ Skipped at index 249 | Name:  | Address:
⚠️ Skipped at index 250 | Name:  | Address:
⚠️ Skipped at index 251 | Name:  | Address:
⚠️ Skipped at index 252 | Name:  | Address:
⚠️ Skipped at index 253 | Name:  | Address:
⚠️ Skipped at index 254 | Name:  | Address:
⚠️ Skipped at index 255 | Name:  | Address:
⚠️ Skipped at index 256 | Name:  | Address:
⚠️ Skipped at index 257 | Name:  | Address:
⚠️ Skipped at index 258 | Name:  | Address:
⚠️ Skipped at index 259 | Name:  | Address:
⚠️ Skipped at index 260 | Name:  | Address:
⚠️ Skipped at index 261 | Name:  | Address:
⚠️ Skipped at index 262 | Name:  | Address:
⚠️ Skipped at index 263 | Name:  | Address:
⚠️ Skipped at index 264 | Name:  | Address:
⚠️ Skipped at index 265 | Name:  | Address:
⚠️ Skipped at index 266 | Name:  | Address:
⚠️ Skipped at index 267 | Name:  | Address:
⚠️ Skipped at index 268 | Name:  | Address:
⚠️ Skipped at index 269 | Name:  | Address:
⚠️ Skipped at index 270 | Name:  | Address:
⚠️ Skipped at index 271 | Name:  | Address:
⚠️ Skipped at index 272 | Name:  | Address:
⚠️ Skipped at index 273 | Name:  | Address:
⚠️ Skipped at index 274 | Name:  | Address:
⚠️ Skipped at index 275 | Name:  | Address:
⚠️ Skipped at index 276 | Name:  | Address:
⚠️ Skipped at index 277 | Name:  | Address:
⚠️ Skipped at index 278 | Name:  | Address:
⚠️ Skipped at index 279 | Name:  | Address:
⚠️ Skipped at index 280 | Name:  | Address:
⚠️ Skipped at index 281 | Name:  | Address:
⚠️ Skipped at index 282 | Name:  | Address:
⚠️ Skipped at index 283 | Name:  | Address:
⚠️ Skipped at index 284 | Name:  | Address:
⚠️ Skipped at index 285 | Name:  | Address:
⚠️ Skipped at index 286 | Name:  | Address:
⚠️ Skipped at index 287 | Name:  | Address:
⚠️ Skipped at index 288 | Name:  | Address:
⚠️ Skipped at index 289 | Name:  | Address:
⚠️ Skipped at index 290 | Name:  | Address:
⚠️ Skipped at index 291 | Name:  | Address:
⚠️ Skipped at index 292 | Name:  | Address:
⚠️ Skipped at index 293 | Name:  | Address:
⚠️ Skipped at index 294 | Name:  | Address:
⚠️ Skipped at index 295 | Name:  | Address:
⚠️ Skipped at index 296 | Name:  | Address:
⚠️ Skipped at index 297 | Name:  | Address:
⚠️ Skipped at index 298 | Name:  | Address:
⚠️ Skipped at index 299 | Name:  | Address:
⚠️ Skipped at index 300 | Name:  | Address:
⚠️ Skipped at index 301 | Name:  | Address:
⚠️ Skipped at index 302 | Name:  | Address:
⚠️ Skipped at index 303 | Name:  | Address:
⚠️ Skipped at index 304 | Name:  | Address:
⚠️ Skipped at index 305 | Name:  | Address:
⚠️ Skipped at index 306 | Name:  | Address:
⚠️ Skipped at index 307 | Name:  | Address:
⚠️ Skipped at index 308 | Name:  | Address:
⚠️ Skipped at index 309 | Name:  | Address:
⚠️ Skipped at index 310 | Name:  | Address:
⚠️ Skipped at index 311 | Name:  | Address:
⚠️ Skipped at index 312 | Name:  | Address:
⚠️ Skipped at index 313 | Name:  | Address:
⚠️ Skipped at index 314 | Name:  | Address:
⚠️ Skipped at index 315 | Name:  | Address:
⚠️ Skipped at index 316 | Name:  | Address:
⚠️ Skipped at index 317 | Name:  | Address:
⚠️ Skipped at index 318 | Name:  | Address:
⚠️ Skipped at index 319 | Name:  | Address:
⚠️ Skipped at index 320 | Name:  | Address:
⚠️ Skipped at index 321 | Name:  | Address:
⚠️ Skipped at index 322 | Name:  | Address:
⚠️ Skipped at index 323 | Name:  | Address:
```

⚠️ Skipped at index 324 | Name: | Address:
⚠️ Skipped at index 325 | Name: | Address:
⚠️ Skipped at index 326 | Name: | Address:
⚠️ Skipped at index 327 | Name: | Address:
⚠️ Skipped at index 328 | Name: | Address:
⚠️ Skipped at index 329 | Name: | Address:
⚠️ Skipped at index 330 | Name: | Address:
⚠️ Skipped at index 331 | Name: | Address:
⚠️ Skipped at index 332 | Name: | Address:
⚠️ Skipped at index 333 | Name: | Address:
⚠️ Skipped at index 334 | Name: | Address:
⚠️ Skipped at index 335 | Name: | Address:
⚠️ Skipped at index 336 | Name: | Address:
⚠️ Skipped at index 337 | Name: | Address:
⚠️ Skipped at index 338 | Name: | Address:
⚠️ Skipped at index 339 | Name: | Address:
⚠️ Skipped at index 340 | Name: | Address:
⚠️ Skipped at index 341 | Name: | Address:
⚠️ Skipped at index 342 | Name: | Address:
⚠️ Skipped at index 343 | Name: | Address:
⚠️ Skipped at index 344 | Name: | Address:
⚠️ Skipped at index 345 | Name: | Address:
⚠️ Skipped at index 346 | Name: | Address:
⚠️ Skipped at index 347 | Name: | Address:
⚠️ Skipped at index 348 | Name: | Address:
⚠️ Skipped at index 349 | Name: | Address:
⚠️ Skipped at index 350 | Name: | Address:
⚠️ Skipped at index 351 | Name: | Address:
⚠️ Skipped at index 352 | Name: | Address:
⚠️ Skipped at index 353 | Name: | Address:
⚠️ Skipped at index 354 | Name: | Address:
⚠️ Skipped at index 355 | Name: | Address:
⚠️ Skipped at index 356 | Name: | Address:
⚠️ Skipped at index 357 | Name: | Address:
⚠️ Skipped at index 358 | Name: | Address:
⚠️ Skipped at index 359 | Name: | Address:
⚠️ Skipped at index 360 | Name: | Address:
⚠️ Skipped at index 361 | Name: | Address:
⚠️ Skipped at index 362 | Name: | Address:
⚠️ Skipped at index 363 | Name: | Address:
⚠️ Skipped at index 364 | Name: | Address:
⚠️ Skipped at index 365 | Name: | Address:
⚠️ Skipped at index 366 | Name: | Address:
⚠️ Skipped at index 367 | Name: | Address:
⚠️ Skipped at index 368 | Name: | Address:
⚠️ Skipped at index 369 | Name: | Address:
⚠️ Skipped at index 370 | Name: | Address:
⚠️ Skipped at index 371 | Name: | Address:
⚠️ Skipped at index 372 | Name: | Address:
⚠️ Skipped at index 373 | Name: | Address:
⚠️ Skipped at index 374 | Name: | Address:
⚠️ Skipped at index 375 | Name: | Address:
⚠️ Skipped at index 376 | Name: | Address:
⚠️ Skipped at index 377 | Name: | Address:
⚠️ Skipped at index 378 | Name: | Address:
⚠️ Skipped at index 379 | Name: | Address:
⚠️ Skipped at index 380 | Name: | Address:
⚠️ Skipped at index 381 | Name: | Address:
⚠️ Skipped at index 382 | Name: | Address:
⚠️ Skipped at index 383 | Name: | Address:
⚠️ Skipped at index 384 | Name: | Address:
⚠️ Skipped at index 385 | Name: | Address:
⚠️ Skipped at index 386 | Name: | Address:
⚠️ Skipped at index 387 | Name: | Address:
⚠️ Skipped at index 388 | Name: | Address:
⚠️ Skipped at index 389 | Name: | Address:
⚠️ Skipped at index 390 | Name: | Address:
⚠️ Skipped at index 391 | Name: | Address:
⚠️ Skipped at index 392 | Name: | Address:
⚠️ Skipped at index 393 | Name: | Address:
⚠️ Skipped at index 394 | Name: | Address:
⚠️ Skipped at index 395 | Name: | Address:
⚠️ Skipped at index 396 | Name: | Address:
⚠️ Skipped at index 397 | Name: | Address:
⚠️ Skipped at index 398 | Name: | Address:
⚠️ Skipped at index 399 | Name: | Address:
⚠️ Skipped at index 400 | Name: | Address:
⚠️ Skipped at index 401 | Name: | Address:
⚠️ Skipped at index 402 | Name: | Address:
⚠️ Skipped at index 403 | Name: | Address:
⚠️ Skipped at index 404 | Name: | Address:
⚠️ Skipped at index 405 | Name: | Address:
⚠️ Skipped at index 406 | Name: | Address:
⚠️ Skipped at index 407 | Name: | Address:
⚠️ Skipped at index 408 | Name: | Address:
⚠️ Skipped at index 409 | Name: | Address:
⚠️ Skipped at index 410 | Name: | Address:

⚠️ Skipped at index 411 | Name: | Address:
⚠️ Skipped at index 412 | Name: | Address:
⚠️ Skipped at index 413 | Name: | Address:
⚠️ Skipped at index 414 | Name: | Address:
⚠️ Skipped at index 415 | Name: | Address:
⚠️ Skipped at index 416 | Name: | Address:
⚠️ Skipped at index 417 | Name: | Address:
⚠️ Skipped at index 418 | Name: | Address:
⚠️ Skipped at index 419 | Name: | Address:
⚠️ Skipped at index 420 | Name: | Address:
⚠️ Skipped at index 421 | Name: | Address:
⚠️ Skipped at index 422 | Name: | Address:
⚠️ Skipped at index 423 | Name: | Address:
⚠️ Skipped at index 424 | Name: | Address:
⚠️ Skipped at index 425 | Name: | Address:
⚠️ Skipped at index 426 | Name: | Address:
⚠️ Skipped at index 427 | Name: | Address:
⚠️ Skipped at index 428 | Name: | Address:
⚠️ Skipped at index 429 | Name: | Address:
⚠️ Skipped at index 430 | Name: | Address:
⚠️ Skipped at index 431 | Name: | Address:
⚠️ Skipped at index 432 | Name: | Address:
⚠️ Skipped at index 433 | Name: | Address:
⚠️ Skipped at index 434 | Name: | Address:
⚠️ Skipped at index 435 | Name: | Address:
⚠️ Skipped at index 436 | Name: | Address:
⚠️ Skipped at index 437 | Name: | Address:
⚠️ Skipped at index 438 | Name: | Address:
⚠️ Skipped at index 439 | Name: | Address:
⚠️ Skipped at index 440 | Name: | Address:
⚠️ Skipped at index 441 | Name: | Address:
⚠️ Skipped at index 442 | Name: | Address:
⚠️ Skipped at index 443 | Name: | Address:
⚠️ Skipped at index 444 | Name: | Address:
⚠️ Skipped at index 445 | Name: | Address:
⚠️ Skipped at index 446 | Name: | Address:
⚠️ Skipped at index 447 | Name: | Address:
⚠️ Skipped at index 448 | Name: | Address:
⚠️ Skipped at index 449 | Name: | Address:
⚠️ Skipped at index 450 | Name: | Address:
⚠️ Skipped at index 451 | Name: | Address:
⚠️ Skipped at index 452 | Name: | Address:
⚠️ Skipped at index 453 | Name: | Address:
⚠️ Skipped at index 454 | Name: | Address:
⚠️ Skipped at index 455 | Name: | Address:
⚠️ Skipped at index 456 | Name: | Address:
⚠️ Skipped at index 457 | Name: | Address:
⚠️ Skipped at index 458 | Name: | Address:
⚠️ Skipped at index 459 | Name: | Address:
⚠️ Skipped at index 460 | Name: | Address:
⚠️ Skipped at index 461 | Name: | Address:
⚠️ Skipped at index 462 | Name: | Address:
⚠️ Skipped at index 463 | Name: | Address:
⚠️ Skipped at index 464 | Name: | Address:
⚠️ Skipped at index 465 | Name: | Address:
⚠️ Skipped at index 466 | Name: | Address:
⚠️ Skipped at index 467 | Name: | Address:
⚠️ Skipped at index 468 | Name: | Address:
⚠️ Skipped at index 469 | Name: | Address:
⚠️ Skipped at index 470 | Name: | Address:
⚠️ Skipped at index 471 | Name: | Address:
⚠️ Skipped at index 472 | Name: | Address:
⚠️ Skipped at index 473 | Name: | Address:
⚠️ Skipped at index 474 | Name: | Address:
⚠️ Skipped at index 475 | Name: | Address:
⚠️ Skipped at index 476 | Name: | Address:
⚠️ Skipped at index 477 | Name: | Address:
⚠️ Skipped at index 478 | Name: | Address:
⚠️ Skipped at index 479 | Name: | Address:
⚠️ Skipped at index 480 | Name: | Address:
⚠️ Skipped at index 481 | Name: | Address:
⚠️ Skipped at index 482 | Name: | Address:
⚠️ Skipped at index 483 | Name: | Address:
⚠️ Skipped at index 484 | Name: | Address:
⚠️ Skipped at index 485 | Name: | Address:
⚠️ Skipped at index 486 | Name: | Address:
⚠️ Skipped at index 487 | Name: | Address:
⚠️ Skipped at index 488 | Name: | Address:
⚠️ Skipped at index 489 | Name: | Address:
⚠️ Skipped at index 490 | Name: | Address:
⚠️ Skipped at index 491 | Name: | Address:
⚠️ Skipped at index 492 | Name: | Address:
⚠️ Skipped at index 493 | Name: | Address:
⚠️ Skipped at index 494 | Name: | Address:
⚠️ Skipped at index 495 | Name: | Address:
⚠️ Skipped at index 496 | Name: | Address:
⚠️ Skipped at index 497 | Name: | Address:

⚠️ Skipped at index 498 | Name: | Address:
⚠️ Skipped at index 499 | Name: | Address:
⚠️ Skipped at index 500 | Name: | Address:
⚠️ Skipped at index 501 | Name: | Address:
⚠️ Skipped at index 502 | Name: | Address:
⚠️ Skipped at index 503 | Name: | Address:
⚠️ Skipped at index 504 | Name: | Address:
⚠️ Skipped at index 505 | Name: | Address:
⚠️ Skipped at index 506 | Name: | Address:
⚠️ Skipped at index 507 | Name: | Address:
⚠️ Skipped at index 508 | Name: | Address:
⚠️ Skipped at index 509 | Name: | Address:
⚠️ Skipped at index 510 | Name: | Address:
⚠️ Skipped at index 511 | Name: | Address:
⚠️ Skipped at index 512 | Name: | Address:
⚠️ Skipped at index 513 | Name: | Address:
⚠️ Skipped at index 514 | Name: | Address:
⚠️ Skipped at index 515 | Name: | Address:
⚠️ Skipped at index 516 | Name: | Address:
⚠️ Skipped at index 517 | Name: | Address:
⚠️ Skipped at index 518 | Name: | Address:
⚠️ Skipped at index 519 | Name: | Address:
⚠️ Skipped at index 520 | Name: | Address:
⚠️ Skipped at index 521 | Name: | Address:
⚠️ Skipped at index 522 | Name: | Address:
⚠️ Skipped at index 523 | Name: | Address:
⚠️ Skipped at index 524 | Name: | Address:
⚠️ Skipped at index 525 | Name: | Address:
⚠️ Skipped at index 526 | Name: | Address:
⚠️ Skipped at index 527 | Name: | Address:
⚠️ Skipped at index 528 | Name: | Address:
⚠️ Skipped at index 529 | Name: | Address:
⚠️ Skipped at index 530 | Name: | Address:
⚠️ Skipped at index 531 | Name: | Address:
⚠️ Skipped at index 532 | Name: | Address:
⚠️ Skipped at index 533 | Name: | Address:
⚠️ Skipped at index 534 | Name: | Address:
⚠️ Skipped at index 535 | Name: | Address:
⚠️ Skipped at index 536 | Name: | Address:
⚠️ Skipped at index 537 | Name: | Address:
⚠️ Skipped at index 538 | Name: | Address:
⚠️ Skipped at index 539 | Name: | Address:
⚠️ Skipped at index 540 | Name: | Address:
⚠️ Skipped at index 541 | Name: | Address:
⚠️ Skipped at index 542 | Name: | Address:
⚠️ Skipped at index 543 | Name: | Address:
⚠️ Skipped at index 544 | Name: | Address:
⚠️ Skipped at index 545 | Name: | Address:
⚠️ Skipped at index 546 | Name: | Address:
⚠️ Skipped at index 547 | Name: | Address:
⚠️ Skipped at index 548 | Name: | Address:
⚠️ Skipped at index 549 | Name: | Address:
⚠️ Skipped at index 550 | Name: | Address:
⚠️ Skipped at index 551 | Name: | Address:
⚠️ Skipped at index 552 | Name: | Address:
⚠️ Skipped at index 553 | Name: | Address:
⚠️ Skipped at index 554 | Name: | Address:
⚠️ Skipped at index 555 | Name: | Address:
⚠️ Skipped at index 556 | Name: | Address:
⚠️ Skipped at index 557 | Name: | Address:
⚠️ Skipped at index 558 | Name: | Address:
⚠️ Skipped at index 559 | Name: | Address:
⚠️ Skipped at index 560 | Name: | Address:
⚠️ Skipped at index 561 | Name: | Address:
⚠️ Skipped at index 562 | Name: | Address:
⚠️ Skipped at index 563 | Name: | Address:
⚠️ Skipped at index 564 | Name: | Address:
⚠️ Skipped at index 565 | Name: | Address:
⚠️ Skipped at index 566 | Name: | Address:
⚠️ Skipped at index 567 | Name: | Address:
⚠️ Skipped at index 568 | Name: | Address:
⚠️ Skipped at index 569 | Name: | Address:
⚠️ Skipped at index 570 | Name: | Address:
⚠️ Skipped at index 571 | Name: | Address:
⚠️ Skipped at index 572 | Name: | Address:
⚠️ Skipped at index 573 | Name: | Address:
⚠️ Skipped at index 574 | Name: | Address:
⚠️ Skipped at index 575 | Name: | Address:
⚠️ Skipped at index 576 | Name: | Address:
⚠️ Skipped at index 577 | Name: | Address:
⚠️ Skipped at index 578 | Name: | Address:
⚠️ Skipped at index 579 | Name: | Address:
⚠️ Skipped at index 580 | Name: | Address:
⚠️ Skipped at index 581 | Name: | Address:
⚠️ Skipped at index 582 | Name: | Address:
⚠️ Skipped at index 583 | Name: | Address:
⚠️ Skipped at index 584 | Name: | Address:

⚠️ Skipped at index 585 | Name: | Address:
⚠️ Skipped at index 586 | Name: | Address:
⚠️ Skipped at index 587 | Name: | Address:
⚠️ Skipped at index 588 | Name: | Address:
⚠️ Skipped at index 589 | Name: | Address:
⚠️ Skipped at index 590 | Name: | Address:
⚠️ Skipped at index 591 | Name: | Address:
⚠️ Skipped at index 592 | Name: | Address:
⚠️ Skipped at index 593 | Name: | Address:
⚠️ Skipped at index 594 | Name: | Address:
⚠️ Skipped at index 595 | Name: | Address:
⚠️ Skipped at index 596 | Name: | Address:
⚠️ Skipped at index 597 | Name: | Address:
⚠️ Skipped at index 598 | Name: | Address:
⚠️ Skipped at index 599 | Name: | Address:
⚠️ Skipped at index 600 | Name: | Address:
⚠️ Skipped at index 601 | Name: | Address:
⚠️ Skipped at index 602 | Name: | Address:
⚠️ Skipped at index 603 | Name: | Address:
⚠️ Skipped at index 604 | Name: | Address:
⚠️ Skipped at index 605 | Name: | Address:
⚠️ Skipped at index 606 | Name: | Address:
⚠️ Skipped at index 607 | Name: | Address:
⚠️ Skipped at index 608 | Name: | Address:
⚠️ Skipped at index 609 | Name: | Address:
⚠️ Skipped at index 610 | Name: | Address:
⚠️ Skipped at index 611 | Name: | Address:
⚠️ Skipped at index 612 | Name: | Address:
⚠️ Skipped at index 613 | Name: | Address:
⚠️ Skipped at index 614 | Name: | Address:
⚠️ Skipped at index 615 | Name: | Address:
⚠️ Skipped at index 616 | Name: | Address:
⚠️ Skipped at index 617 | Name: | Address:
⚠️ Skipped at index 618 | Name: | Address:
⚠️ Skipped at index 619 | Name: | Address:
⚠️ Skipped at index 620 | Name: | Address:
⚠️ Skipped at index 621 | Name: | Address:
⚠️ Skipped at index 622 | Name: | Address:
⚠️ Skipped at index 623 | Name: | Address:
⚠️ Skipped at index 624 | Name: | Address:
⚠️ Skipped at index 625 | Name: | Address:
⚠️ Skipped at index 626 | Name: | Address:
⚠️ Skipped at index 627 | Name: | Address:
⚠️ Skipped at index 628 | Name: | Address:
⚠️ Skipped at index 629 | Name: | Address:
⚠️ Skipped at index 630 | Name: | Address:
⚠️ Skipped at index 631 | Name: | Address:
⚠️ Skipped at index 632 | Name: | Address:
⚠️ Skipped at index 633 | Name: | Address:
⚠️ Skipped at index 634 | Name: | Address:
⚠️ Skipped at index 635 | Name: | Address:
⚠️ Skipped at index 636 | Name: | Address:
⚠️ Skipped at index 637 | Name: | Address:
⚠️ Skipped at index 638 | Name: | Address:
⚠️ Skipped at index 639 | Name: | Address:
⚠️ Skipped at index 640 | Name: | Address:
⚠️ Skipped at index 641 | Name: | Address:
⚠️ Skipped at index 642 | Name: | Address:
⚠️ Skipped at index 643 | Name: | Address:
⚠️ Skipped at index 644 | Name: | Address:
⚠️ Skipped at index 645 | Name: | Address:
⚠️ Skipped at index 646 | Name: | Address:
⚠️ Skipped at index 647 | Name: | Address:
⚠️ Skipped at index 648 | Name: | Address:
⚠️ Skipped at index 649 | Name: | Address:
⚠️ Skipped at index 650 | Name: | Address:
⚠️ Skipped at index 651 | Name: | Address:
⚠️ Skipped at index 652 | Name: | Address:
⚠️ Skipped at index 653 | Name: | Address:
⚠️ Skipped at index 654 | Name: | Address:
⚠️ Skipped at index 655 | Name: | Address:
⚠️ Skipped at index 656 | Name: | Address:
⚠️ Skipped at index 657 | Name: | Address:
⚠️ Skipped at index 658 | Name: | Address:
⚠️ Skipped at index 659 | Name: | Address:
⚠️ Skipped at index 660 | Name: | Address:
⚠️ Skipped at index 661 | Name: | Address:
⚠️ Skipped at index 662 | Name: | Address:
⚠️ Skipped at index 663 | Name: | Address:
⚠️ Skipped at index 664 | Name: | Address:
⚠️ Skipped at index 665 | Name: | Address:
⚠️ Skipped at index 666 | Name: | Address:
⚠️ Skipped at index 667 | Name: | Address:
⚠️ Skipped at index 668 | Name: | Address:
⚠️ Skipped at index 669 | Name: | Address:
⚠️ Skipped at index 670 | Name: | Address:
⚠️ Skipped at index 671 | Name: | Address:

⚠️ Skipped at index 672 | Name: | Address:
⚠️ Skipped at index 673 | Name: | Address:
⚠️ Skipped at index 674 | Name: | Address:
⚠️ Skipped at index 675 | Name: | Address:
⚠️ Skipped at index 676 | Name: | Address:
⚠️ Skipped at index 677 | Name: | Address:
⚠️ Skipped at index 678 | Name: | Address:
⚠️ Skipped at index 679 | Name: | Address:
⚠️ Skipped at index 680 | Name: | Address:
⚠️ Skipped at index 681 | Name: | Address:
⚠️ Skipped at index 682 | Name: | Address:
⚠️ Skipped at index 683 | Name: | Address:
⚠️ Skipped at index 684 | Name: | Address:
⚠️ Skipped at index 685 | Name: | Address:
⚠️ Skipped at index 686 | Name: | Address:
⚠️ Skipped at index 687 | Name: | Address:
⚠️ Skipped at index 688 | Name: | Address:
⚠️ Skipped at index 689 | Name: | Address:
⚠️ Skipped at index 690 | Name: | Address:
⚠️ Skipped at index 691 | Name: | Address:
⚠️ Skipped at index 692 | Name: | Address:
⚠️ Skipped at index 693 | Name: | Address:
⚠️ Skipped at index 694 | Name: | Address:
⚠️ Skipped at index 695 | Name: | Address:
⚠️ Skipped at index 696 | Name: | Address:
⚠️ Skipped at index 697 | Name: | Address:
⚠️ Skipped at index 698 | Name: | Address:
⚠️ Skipped at index 699 | Name: | Address:
⚠️ Skipped at index 700 | Name: | Address:
⚠️ Skipped at index 701 | Name: | Address:
⚠️ Skipped at index 702 | Name: | Address:
⚠️ Skipped at index 703 | Name: | Address:
⚠️ Skipped at index 704 | Name: | Address:
⚠️ Skipped at index 705 | Name: | Address:
⚠️ Skipped at index 706 | Name: | Address:
⚠️ Skipped at index 707 | Name: | Address:
⚠️ Skipped at index 708 | Name: | Address:
⚠️ Skipped at index 709 | Name: | Address:
⚠️ Skipped at index 710 | Name: | Address:
⚠️ Skipped at index 711 | Name: | Address:
⚠️ Skipped at index 712 | Name: | Address:
⚠️ Skipped at index 713 | Name: | Address:
⚠️ Skipped at index 714 | Name: | Address:
⚠️ Skipped at index 715 | Name: | Address:
⚠️ Skipped at index 716 | Name: | Address:
⚠️ Skipped at index 717 | Name: | Address:
⚠️ Skipped at index 718 | Name: | Address:
⚠️ Skipped at index 719 | Name: | Address:
⚠️ Skipped at index 720 | Name: | Address:
⚠️ Skipped at index 721 | Name: | Address:
⚠️ Skipped at index 722 | Name: | Address:
⚠️ Skipped at index 723 | Name: | Address:
⚠️ Skipped at index 724 | Name: | Address:
⚠️ Skipped at index 725 | Name: | Address:
⚠️ Skipped at index 726 | Name: | Address:
⚠️ Skipped at index 727 | Name: | Address:
⚠️ Skipped at index 728 | Name: | Address:
⚠️ Skipped at index 729 | Name: | Address:
⚠️ Skipped at index 730 | Name: | Address:
⚠️ Skipped at index 731 | Name: | Address:
⚠️ Skipped at index 732 | Name: | Address:
⚠️ Skipped at index 733 | Name: | Address:
⚠️ Skipped at index 734 | Name: | Address:
⚠️ Skipped at index 735 | Name: | Address:
⚠️ Skipped at index 736 | Name: | Address:
⚠️ Skipped at index 737 | Name: | Address:
⚠️ Skipped at index 738 | Name: | Address:
⚠️ Skipped at index 739 | Name: | Address:
⚠️ Skipped at index 740 | Name: | Address:
⚠️ Skipped at index 741 | Name: | Address:
⚠️ Skipped at index 742 | Name: | Address:
⚠️ Skipped at index 743 | Name: | Address:
⚠️ Skipped at index 744 | Name: | Address:
⚠️ Skipped at index 745 | Name: | Address:
⚠️ Skipped at index 746 | Name: | Address:
⚠️ Skipped at index 747 | Name: | Address:
⚠️ Skipped at index 748 | Name: | Address:
⚠️ Skipped at index 749 | Name: | Address:
⚠️ Skipped at index 750 | Name: | Address:
⚠️ Skipped at index 751 | Name: | Address:
⚠️ Skipped at index 752 | Name: | Address:
⚠️ Skipped at index 753 | Name: | Address:
⚠️ Skipped at index 754 | Name: | Address:
⚠️ Skipped at index 755 | Name: | Address:
⚠️ Skipped at index 756 | Name: | Address:
⚠️ Skipped at index 757 | Name: | Address:
⚠️ Skipped at index 758 | Name: | Address:

⚠️ Skipped at index 759 | Name: | Address:
⚠️ Skipped at index 760 | Name: | Address:
⚠️ Skipped at index 761 | Name: | Address:
⚠️ Skipped at index 762 | Name: | Address:
⚠️ Skipped at index 763 | Name: | Address:
⚠️ Skipped at index 764 | Name: | Address:
⚠️ Skipped at index 765 | Name: | Address:
⚠️ Skipped at index 766 | Name: | Address:
⚠️ Skipped at index 767 | Name: | Address:
⚠️ Skipped at index 768 | Name: | Address:
⚠️ Skipped at index 769 | Name: | Address:
⚠️ Skipped at index 770 | Name: | Address:
⚠️ Skipped at index 771 | Name: | Address:
⚠️ Skipped at index 772 | Name: | Address:
⚠️ Skipped at index 773 | Name: | Address:
⚠️ Skipped at index 774 | Name: | Address:
⚠️ Skipped at index 775 | Name: | Address:
⚠️ Skipped at index 776 | Name: | Address:
⚠️ Skipped at index 777 | Name: | Address:
⚠️ Skipped at index 778 | Name: | Address:
⚠️ Skipped at index 779 | Name: | Address:
⚠️ Skipped at index 780 | Name: | Address:
⚠️ Skipped at index 781 | Name: | Address:
⚠️ Skipped at index 782 | Name: | Address:
⚠️ Skipped at index 783 | Name: | Address:
⚠️ Skipped at index 784 | Name: | Address:
⚠️ Skipped at index 785 | Name: | Address:
⚠️ Skipped at index 786 | Name: | Address:
⚠️ Skipped at index 787 | Name: | Address:
⚠️ Skipped at index 788 | Name: | Address:
⚠️ Skipped at index 789 | Name: | Address:
⚠️ Skipped at index 790 | Name: | Address:
⚠️ Skipped at index 791 | Name: | Address:
⚠️ Skipped at index 792 | Name: | Address:
⚠️ Skipped at index 793 | Name: | Address:
⚠️ Skipped at index 794 | Name: | Address:
⚠️ Skipped at index 795 | Name: | Address:
⚠️ Skipped at index 796 | Name: | Address:
⚠️ Skipped at index 797 | Name: | Address:
⚠️ Skipped at index 798 | Name: | Address:
⚠️ Skipped at index 799 | Name: | Address:
⚠️ Skipped at index 800 | Name: | Address:
⚠️ Skipped at index 801 | Name: | Address:
⚠️ Skipped at index 802 | Name: | Address:
⚠️ Skipped at index 803 | Name: | Address:
⚠️ Skipped at index 804 | Name: | Address:
⚠️ Skipped at index 805 | Name: | Address:
⚠️ Skipped at index 806 | Name: | Address:
⚠️ Skipped at index 807 | Name: | Address:
⚠️ Skipped at index 808 | Name: | Address:
⚠️ Skipped at index 809 | Name: | Address:
⚠️ Skipped at index 810 | Name: | Address:
⚠️ Skipped at index 811 | Name: | Address:
⚠️ Skipped at index 812 | Name: | Address:
⚠️ Skipped at index 813 | Name: | Address:
⚠️ Skipped at index 814 | Name: | Address:
⚠️ Skipped at index 815 | Name: | Address:
⚠️ Skipped at index 816 | Name: | Address:
⚠️ Skipped at index 817 | Name: | Address:
⚠️ Skipped at index 818 | Name: | Address:
⚠️ Skipped at index 819 | Name: | Address:
⚠️ Skipped at index 820 | Name: | Address:
⚠️ Skipped at index 821 | Name: | Address:
⚠️ Skipped at index 822 | Name: | Address:
⚠️ Skipped at index 823 | Name: | Address:
⚠️ Skipped at index 824 | Name: | Address:
⚠️ Skipped at index 825 | Name: | Address:
⚠️ Skipped at index 826 | Name: | Address:
⚠️ Skipped at index 827 | Name: | Address:
⚠️ Skipped at index 828 | Name: | Address:
⚠️ Skipped at index 829 | Name: | Address:
⚠️ Skipped at index 830 | Name: | Address:
⚠️ Skipped at index 831 | Name: | Address:
⚠️ Skipped at index 832 | Name: | Address:
⚠️ Skipped at index 833 | Name: | Address:
⚠️ Skipped at index 834 | Name: | Address:
⚠️ Skipped at index 835 | Name: | Address:
⚠️ Skipped at index 836 | Name: | Address:
⚠️ Skipped at index 837 | Name: | Address:
⚠️ Skipped at index 838 | Name: | Address:
⚠️ Skipped at index 839 | Name: | Address:
⚠️ Skipped at index 840 | Name: | Address:
⚠️ Skipped at index 841 | Name: | Address:
⚠️ Skipped at index 842 | Name: | Address:
⚠️ Skipped at index 843 | Name: | Address:
⚠️ Skipped at index 844 | Name: | Address:
⚠️ Skipped at index 845 | Name: | Address:

⚠ Skipped at index 846 | Name: | Address:
⚠ Skipped at index 847 | Name: | Address:
⚠ Skipped at index 848 | Name: | Address:
⚠ Skipped at index 849 | Name: | Address:
⚠ Skipped at index 850 | Name: | Address:
⚠ Skipped at index 851 | Name: | Address:
⚠ Skipped at index 852 | Name: | Address:
⚠ Skipped at index 853 | Name: | Address:
⚠ Skipped at index 854 | Name: | Address:
⚠ Skipped at index 855 | Name: | Address:
⚠ Skipped at index 856 | Name: | Address:
⚠ Skipped at index 857 | Name: | Address:
⚠ Skipped at index 858 | Name: | Address:
⚠ Skipped at index 859 | Name: | Address:
⚠ Skipped at index 860 | Name: | Address:
⚠ Skipped at index 861 | Name: | Address:
⚠ Skipped at index 862 | Name: | Address:
⚠ Skipped at index 863 | Name: | Address:
⚠ Skipped at index 864 | Name: | Address:
⚠ Skipped at index 865 | Name: | Address:
⚠ Skipped at index 866 | Name: | Address:
⚠ Skipped at index 867 | Name: | Address:
⚠ Skipped at index 868 | Name: | Address:
⚠ Skipped at index 869 | Name: | Address:
⚠ Skipped at index 870 | Name: | Address:
⚠ Skipped at index 871 | Name: | Address:
⚠ Skipped at index 872 | Name: | Address:
⚠ Skipped at index 873 | Name: | Address:
⚠ Skipped at index 874 | Name: | Address:
⚠ Skipped at index 875 | Name: | Address:
⚠ Skipped at index 876 | Name: | Address:
⚠ Skipped at index 877 | Name: | Address:
⚠ Skipped at index 878 | Name: | Address:
⚠ Skipped at index 879 | Name: | Address:
⚠ Skipped at index 880 | Name: | Address:
⚠ Skipped at index 881 | Name: | Address:
⚠ Skipped at index 882 | Name: | Address:
⚠ Skipped at index 883 | Name: | Address:
⚠ Skipped at index 884 | Name: | Address:
⚠ Skipped at index 885 | Name: | Address:
⚠ Skipped at index 886 | Name: | Address:
⚠ Skipped at index 887 | Name: | Address:
⚠ Skipped at index 888 | Name: | Address:
⚠ Skipped at index 889 | Name: | Address:
⚠ Skipped at index 890 | Name: | Address:
⚠ Skipped at index 891 | Name: | Address:
⚠ Skipped at index 892 | Name: | Address:
⚠ Skipped at index 893 | Name: | Address:
⚠ Skipped at index 894 | Name: | Address:
⚠ Skipped at index 895 | Name: | Address:
⚠ Skipped at index 896 | Name: | Address:
⚠ Skipped at index 897 | Name: | Address:
⚠ Skipped at index 898 | Name: | Address:
⚠ Skipped at index 899 | Name: | Address:
⚠ Skipped at index 900 | Name: | Address:
⚠ Skipped at index 901 | Name: | Address:
⚠ Skipped at index 902 | Name: | Address:
⚠ Skipped at index 903 | Name: | Address:
⚠ Skipped at index 904 | Name: | Address:
⚠ Skipped at index 905 | Name: | Address:
⚠ Skipped at index 906 | Name: | Address:
⚠ Skipped at index 907 | Name: | Address:
⚠ Skipped at index 908 | Name: | Address:
⚠ Skipped at index 909 | Name: | Address:
⚠ Skipped at index 910 | Name: | Address:
⚠ Skipped at index 911 | Name: | Address:
⚠ Skipped at index 912 | Name: | Address:
⚠ Skipped at index 913 | Name: | Address:
⚠ Skipped at index 914 | Name: | Address:
⚠ Skipped at index 915 | Name: | Address:
⚠ Skipped at index 916 | Name: | Address:
⚠ Skipped at index 917 | Name: | Address:
⚠ Skipped at index 918 | Name: | Address:
⚠ Skipped at index 919 | Name: | Address:
⚠ Skipped at index 920 | Name: | Address:
⚠ Skipped at index 921 | Name: | Address:
⚠ Skipped at index 922 | Name: | Address:
⚠ Skipped at index 923 | Name: | Address:
⚠ Skipped at index 924 | Name: | Address:
⚠ Skipped at index 925 | Name: | Address:
⚠ Skipped at index 926 | Name: | Address:
⚠ Skipped at index 927 | Name: | Address:
⚠ Skipped at index 928 | Name: | Address:
⚠ Skipped at index 929 | Name: | Address:
⚠ Skipped at index 930 | Name: | Address:
⚠ Skipped at index 931 | Name: | Address:
⚠ Skipped at index 932 | Name: | Address:

```
⚠️ Skipped at index 933 | Name:  | Address:
⚠️ Skipped at index 934 | Name:  | Address:
⚠️ Skipped at index 935 | Name:  | Address:
⚠️ Skipped at index 936 | Name:  | Address:
⚠️ Skipped at index 937 | Name:  | Address:
⚠️ Skipped at index 938 | Name:  | Address:
⚠️ Skipped at index 939 | Name:  | Address:
⚠️ Skipped at index 940 | Name:  | Address:
⚠️ Skipped at index 941 | Name:  | Address:
⚠️ Skipped at index 942 | Name:  | Address:
⚠️ Skipped at index 943 | Name:  | Address:
⚠️ Skipped at index 944 | Name:  | Address:
⚠️ Skipped at index 945 | Name:  | Address:
⚠️ Skipped at index 946 | Name:  | Address:
⚠️ Skipped at index 947 | Name:  | Address:
⚠️ Skipped at index 948 | Name:  | Address:
⚠️ Skipped at index 949 | Name:  | Address:
⚠️ Skipped at index 950 | Name:  | Address:
⚠️ Skipped at index 951 | Name:  | Address:
⚠️ Skipped at index 952 | Name:  | Address:
⚠️ Skipped at index 953 | Name:  | Address:
⚠️ Skipped at index 954 | Name:  | Address:
⚠️ Skipped at index 955 | Name:  | Address:
⚠️ Skipped at index 956 | Name:  | Address:
⚠️ Skipped at index 957 | Name:  | Address:
⚠️ Skipped at index 958 | Name:  | Address:
⚠️ Skipped at index 959 | Name:  | Address:

📊 Preview of Top 5:
                                     Center Name  \
0                        Oasis Family Health NP RN
1   Integrative Rheumatology of Westchester | Dr. ...
2        Agile Infusion Services LLC - Hackensack, NJ
3          Thrivewell Infusion - Tarrytown/Elmsford
4                   VIVO Infusion - Valhalla, NY

                              Address
0        11 Medical Park Dr, Pomona, NY
1   838 Pelhamdale Ave, New Rochelle, NY
2           5 Summit Ave, Hackensack, NJ
3             555 Taxter Rd, Elmsford, NY
4          400 Columbus Ave, Valhalla, NY

📦 Total Valid Centers Extracted: 240
```

In [45]: `df.head()`

Out[45]:

| | Center Name | Address |
|---|---|---|
| 0 | Oasis Family Health NP RN | 11 Medical Park Dr, Pomona, NY |
| 1 | Integrative Rheumatology of Westchester \| Dr. ... | 838 Pelhamdale Ave, New Rochelle, NY |
| 2 | Agile Infusion Services LLC - Hackensack, NJ | 5 Summit Ave, Hackensack, NJ |
| 3 | Thrivewell Infusion - Tarrytown/Elmsford | 555 Taxter Rd, Elmsford, NY |
| 4 | VIVO Infusion - Valhalla, NY | 400 Columbus Ave, Valhalla, NY |

In [46]: `df.shape`

Out[46]: `(240, 2)`

In [47]: `df.to_csv("infusion_centers.csv")`

In [ ]: