# 🚀 Web Scraping Workshop

## Extracting Insights in the Digital Age

**Presented by:**

Ch. Sindhura & S Jai Prakash (JP)

# 🤝 Workshop Guidelines

- ✋ Raise your hand if you need assistance
- 💭 Feel free to ask questions at any time
- 🎯 Stay focused and engaged
- 🤝 Respect others' opinions and questions
- 📱 Keep phones on silent
- 🤔 No question is too basic - we're here to learn!

# 🤔 Interactive Session: What Do You Know?

Let's start with some questions:

1. What comes to your mind when you hear "web scraping"?
2. Have you ever needed data from a website but couldn't easily download it?
3. What websites would you like to extract data from?
4. What challenges do you expect in web scraping?

*Raise your hand to share your thoughts!*

# 📜 The Evolution of Data Extraction

- 🌐 **Early Days**: Manual screen scraping and data collection
- 🤖 **Modern Era**: Automated scripts, intelligent APIs, AI-driven extraction
- 💡 **Transformation**: From time-consuming manual work to sophisticated, efficient data gathering

# 📊 Data: The New Digital Gold

Data Power

# 🎯 Quick Poll

Which of these have you used before?

- Python
- HTML/CSS
- APIs
- Browser Developer Tools

*Raise your hand for each tool you're familiar with!*

# 🤔 What is Web Scraping?

Web scraping is the automated extraction of information from websites, transforming raw web content into structured, actionable data.

## Key Technical Components

- **HTTP Requests**: Fetching web pages programmatically
- **HTML Parsing**: Extracting meaningful information
- **Data Transformation**: Converting raw data into usable formats

# 🤔 Discussion Time

## Share your experiences:

- Have you ever tried to copy data manually from websites?
- What was the largest amount of data you needed to collect?
- How long did it take?

*Let's discuss the inefficiencies of manual data collection!*

# 🌟 Diverse Applications of Web Scraping

- 📊 **Market Research**: Price tracking, competitor analysis
- 📈 **Business Intelligence**: Trend identification, market sentiment
- 🔍 **Academic & Research**: Data collection, trend analysis
- 🤖 **AI & Machine Learning**: Training data generation

# ⚖️ Ethical Web Scraping Guidelines

## Responsible Data Extraction

- **Respect Website Policies**: Always check robots.txt
- **Rate Limiting**: Prevent server overload
- **Legal Compliance**: Ensure authorized data access
- **Transparency**: Understand and respect data usage terms

# 🛠️ Python Scraping Toolkit

## Essential Libraries

- **Requests**: HTTP communication
- **BeautifulSoup**: HTML parsing
- **Selenium**: Dynamic content handling
- **Pandas**: Data manipulation

# 📝 Basic Web Scraping Example

```python
import requests
from bs4 import BeautifulSoup

def scrape_website(url):
    try:
        response = requests.get(url)
        soup = BeautifulSoup(response.text, 'html.parser')

        # Extract specific elements
        titles = soup.find_all('h2')
        return [title.text for title in titles]

    except requests.RequestException as e:
        print(f"Scraping error: {e}")
```

# 🌐 Understanding HTML Structure

```html
<div class="product">
    <h2 class="title">Product Name</h2>
    <span class="price">$99.99</span>
    <div class="description">
        Product details here...
    </div>
</div>
```

## How to Extract:

```python
# Find specific product
product = soup.find('div', class_='product')
title = product.find('h2').text
price = product.find('span', class_='price').text
```

13

# 🔄 Handling Dynamic Content

```python
from selenium import webdriver
from selenium.webdriver.common.by import By

driver = webdriver.Chrome()
driver.get(url)

# Wait for dynamic content to load
driver.implicitly_wait(10)

# Find elements after JavaScript execution
elements = driver.find_elements(By.CLASS_NAME, 'dynamic-content')
```

# 📊 Data Storage Best Practices

```python
import pandas as pd

# Create DataFrame
df = pd.DataFrame({
    'titles': titles,
    'prices': prices,
    'descriptions': descriptions
})

# Save to CSV
df.to_csv('scraped_data.csv', index=False)

# Export to JSON
df.to_json('scraped_data.json', orient='records')
```

# 🤖 AI-Enhanced Web Scraping

## Intelligent Data Extraction

- **Contextual Understanding**: Beyond simple text extraction

- **Adaptive Parsing**: Handling complex web structures

- **Machine Learning Integration**: Improving extraction accuracy

# 🎯 Workshop Practice Session

**Let's apply what we learned:**

1. Identify a website you want to scrape

2. Check its robots.txt

3. Plan your scraping strategy

4. Write basic code together

*Remember to ask for help if needed!*

# 🎯 Workshop Objectives

## Hands-on Projects

- Wikipedia Oscar Winners Data Extraction
- IMDB Movie Review Scraping
- Sentiment Analysis with ChatGPT
- Advanced Techniques with Selenium

# 🤝 Final Q&A Session

- What did you learn today?

- What would you like to scrape next?

- Any challenges you foresee?

- How will you use web scraping in your work?

# 💡 Best Practices & Tips

1. Always check website's Terms of Service

2. Implement proper error handling

3. Use appropriate delays between requests

4. Keep your code modular and documented

5. Regular testing and maintenance

6. Backup your scraped data

# 🚀 Resources for Further Learning

- Official Documentation:
  - Beautiful Soup
  - Selenium
  - Requests
- Online Courses and Tutorials
- Web Scraping Communities
- Practice Websites for Scraping

# 🚀 Happy Scraping! 🕸

## Connect with us:

- Ch. Sindhura: https://www.linkedin.com/in/sindhura-chinoori-710b5165/
- S Jai Prakash (JP): https://www.linkedin.com/in/s-jaiprakash/

# 📚 Additional Resources

## Recommended Reading:

- Web Scraping with Python by Ryan Mitchell
- Python Requests Documentation
- Beautiful Soup Documentation
- Selenium with Python Documentation

## Practice Websites:

- quotes.toscrape.com
- books.toscrape.com

23