

## **ABSTRACT**

**WildStat** is a comprehensive data-driven platform focused on biodiversity analysis within U.S. National Parks. Leveraging Python and powerful data science libraries, the project aims to uncover patterns in endangered species distributions, analyze ecological trends, and support informed decision-making for conservation efforts. This project represents a vital application of statistical analysis and data visualization in addressing ecological challenges and promoting sustainable biodiversity management.

The report begins with an overview of biodiversity threats and the critical role of data analysis in wildlife conservation. It reviews existing efforts in biodiversity monitoring and highlights the limitations of manual observation methods and fragmented datasets. The literature review explores previous research in ecological informatics, endangered species tracking, and public data utilization to identify the gap WildStat addresses.

The methodology section outlines the complete data science pipeline used in the project, including data acquisition from the National Park Service, preprocessing, exploratory data analysis, and visualization techniques. Using libraries such as pandas, NumPy, Matplotlib, and Seaborn, the project identifies species richness, conservation status distributions, and inter-park biodiversity comparisons.

The implementation chapter details the use of Jupyter Notebooks as the primary environment for analysis, explaining the design of graphs, statistical metrics, and interpretations. Key findings include hotspots of endangered species, correlations between park types and biodiversity levels, and insights into species protection status.

Finally, the report discusses future directions such as integrating geospatial data, building a dashboard for real-time biodiversity monitoring, and incorporating predictive modeling to forecast species population trends. The report concludes by emphasizing WildStat's contribution to data-backed ecological research and its potential impact in conservation policy and educational outreach.

## **SYNOPSIS OF THE PROJECT**

**(Topic – WildStat)**

### **IV)**

#### **Problem Statement:**

Biodiversity is declining at an alarming rate due to habitat destruction, climate change, and human activities. Many species are either endangered, threatened, or at risk of extinction, but the lack of proper data analysis and visualization makes it difficult to track and take action. National parks play a crucial role in conservation, yet we often fail to analyze trends in species populations and understand which species need immediate protection.

#### **Proposed Solution:**

To address the challenges of biodiversity loss and species conservation, this project proposes an analytical approach using Python and data visualization techniques to study species distribution in U.S. National Parks. By processing and analyzing biodiversity datasets, the project will uncover trends in species populations, conservation status, and park-wise diversity. Using libraries like Pandas, Matplotlib, and Seaborn, the data will be transformed into meaningful visualizations such as bar charts, pie charts, and heatmaps, helping researchers and conservationists identify species at risk. This solution not only provides data-driven insights for better conservation planning but also lays the groundwork for future AI/ML-based predictions of species population trends and extinction risks.

#### **Objectives:**

- Analyze species distribution across U.S. National Parks.
- conservation statuses of different species (Endangered, Threatened, Safe, etc).
- Visualize biodiversity trends using Python, Pandas, Matplotlib, and Seaborn.
- Provide insights that can assist in better conservation planning.

#### **Literature Review:**

Biodiversity conservation has been widely studied, with research emphasizing the impact of climate change, habitat destruction, and human activities on species extinction. Studies from organizations like the IUCN and GBIF highlight the need for data-driven approaches to track species populations and their conservation status. Traditional methods rely on field surveys and ecological modeling, but recent advancements in data science have enabled better analysis through visualization and predictive analytics. Existing biodiversity projects use GIS mapping, machine learning, and statistical analysis to study species distribution. However, there is a gap in easily accessible, interactive data visualization that helps conservationists make real-time decisions. This project builds on past research by using Python, Pandas, and visualization libraries to analyze biodiversity data from U.S. National Parks, providing insights that can support conservation efforts and future AI-based predictions.

## **Methodology:**

The Prototype Methodology is being used in this project, which enables ongoing improvement through user feedback and iterative development. The process includes:

1. Data Collection: Import biodiversity datasets (species information + observations in national parks).
2. Data Cleaning & Processing: Handle missing values, filter important data.
3. Exploratory Data Analysis (EDA): Find patterns, conservation status distribution, and park- wise species count.
4. Visualization: Use bar charts, pie charts, and heatmaps to represent findings.
5. Insights & Conclusions: Identify key patterns and suggest conservation efforts.

## **Tools:**

- Programming Language: Python
- Development Environment: Jupyter Notebook
- Data Processing & Analysis: Pandas, NumPy
- Data Visualization: Matplotlib, Seaborn
- Database : SQLite / CSV-based datasets
- Dataset Source: National Park Service Biodiversity Dataset

## **Research Gap:**

There is a lack of data-driven, interactive approaches that provide real-time insights into species distribution and conservation status. Existing projects primarily rely on static datasets and general visualizations, which may not fully capture the dynamic nature of biodiversity changes across different parks. Furthermore, there is limited use of machine learning to predict trends or provide actionable recommendations based on historical data.

This project aims to fill these gaps by combining real-time data analysis with interactive visualizations, and by exploring the potential of AI/ML for predicting species population trends in U.S. National Parks. The goal is to create a more dynamic, user-friendly tool that can assist conservationists with data-driven decision-making.

## **References:**

- Global Biodiversity Information Facility (GBIF): Official Website. <https://www.gbif.org>
- IUCN Red List: Official Website. <https://www.iucnredlist.org>
- Matplotlib: Official Documentation. <https://matplotlib.org/>
- Seaborn: Official Documentation. <https://seaborn.pydata.org/>

# TABLE OF CONTENTS

S NO.	TOPICS	PAGE NO.
I.	Self-Certificate	2
II.	Acknowledgement	3
III.	Certificate From the Guide	4
IV.	Synopsis	6
V.	Main Report i. Objective & Scope of the Project. ii. Theoretical Background Definition of Problem. iii. System Analysis & Design vis-a-vis User Requirements. iv. System Planning (PERT Chart). v. Methodology adopted, System Implementation & Details of Hardware & Software used System Maintenance & Evaluation. vi. Detailed Life Cycle of the Project a. ERD, DFD b. Screenshots of the project c. Process involved d. Methodology used testing e. Test Report	9 10 11 12 16 17 20
VI.	Outcome of the project	68
VII.	Conclusion and Future Scope	80
VIII.	References	82

# MAIN REPORT

V)

## i) Objective & Scope:

### **objective:**

The primary objective of this report is to:

- **Develop WildStat**, a data-driven tool to analyze and visualize biodiversity trends in U.S. National Parks using Python-based data science techniques.
- **Classify species** (flora and fauna) by conservation status (endangered, threatened, vulnerable, safe) and identify at-risk populations.
- **Compare biodiversity metrics** across parks to highlight conservation priorities.
- **Educate stakeholders** (researchers, policymakers, park visitors) through interactive visualizations on species distribution and threats.
- **Evaluate WildStat's potential** to support conservation efforts by providing actionable insights from open-source datasets.

### **Scope:**

This report will encompass the following:

#### **Focus Areas for Analysis:**

- **Species Categorization:**
  - Animals (mammals, birds, reptiles, amphibians) and plants.
  - Conservation status labels (e.g., "Endangered" per IUCN/NPS standards).
- **Park-Specific Trends:**
  - Geographic distribution of endangered species (e.g., Yellowstone vs. Everglades).
  - Correlation between human activity and species vulnerability.
- **Data Sources:**
  - National Parks Service (NPS) and Kaggle datasets.
  - Variables: Species names, categories, conservation status, park locations. Geographical Focus
- **Primarily U.S. National Parks due to data availability and standardization.**
- **Future expansion potential:** Global/Indian parks (if datasets are accessible).  
Target Audience
- **Students & Researchers:** For academic projects in ecology/data science.
- **Park Authorities:** To monitor species and allocate conservation resources.
- **General Public:** To raise awareness about biodiversity loss.

## **Technology Utilized**

- **Python Libraries:** Pandas (data cleaning), Matplotlib/Seaborn (visualizations).
- **Tools:** Jupyter Notebook, GitHub for code sharing.  
Limitations
- Relies on static datasets (no real-time updates).
- Excludes socio-economic factors affecting conservation (e.g., funding, poaching).
- Simplified assumptions (e.g., uniform data quality across parks).  
Future Directions
- Machine learning models to predict species at risk.
- Integration with real-time APIs (e.g., NPS wildlife tracking).
- Mobile app for park visitors to report species sightings.

## **ii) Theoretical Background Definition of Problem**

### **Theoretical Background**

Biodiversity is the foundation of healthy ecosystems, directly influencing ecological balance, climate regulation, food security, and the overall well-being of our planet. However, with increasing urbanization, habitat destruction, pollution, and climate change, the world is facing an unprecedented biodiversity crisis. Numerous species are rapidly declining, and many are on the brink of extinction, threatening the delicate balance of nature.

U.S. National Parks serve as critical biodiversity reservoirs, protecting unique ecosystems, endangered species, and natural habitats. Government and environmental organizations have long been collecting data on species distributions, conservation status, and ecological changes. However, this data is often underutilized, scattered across sources, or presented in formats not accessible to the general public.

Analyzing biodiversity requires understanding the complex relationships between species, habitats, and environmental threats. Traditional methods of ecological analysis involve manual fieldwork and static reporting, which, while valuable, are time-consuming and lack scalability. The rise of data science has opened new doors for efficient, large-scale ecological analysis by enabling rapid exploration, visualization, and interpretation of structured biodiversity datasets.

**WildStat** leverages this potential by providing a centralized, Python-based platform to analyze species data from U.S. National Parks. It simplifies biodiversity analysis using visualizations and summary statistics, allowing both experts and non-experts to understand patterns and derive actionable insights. Through tools like Jupyter Notebooks and libraries such as pandas and Seaborn, the platform presents an intuitive approach to conservation-driven analytics.

### **Definition of Problem**

Despite the availability of biodiversity data, several key challenges persist:

- **Lack of Accessibility:** Existing data on endangered species and biodiversity trends are often siloed in large databases or research papers, inaccessible to students, enthusiasts, and local policymakers.
- **Insufficient Visualization:** Raw ecological data is difficult to interpret without proper visual representation. Many datasets lack dynamic tools to help users explore species distribution or conservation status meaningfully.
- **Fragmented Analysis Tools:** While various ecological studies exist, there is no lightweight, open-source solution that aggregates and visualizes biodiversity metrics across multiple national parks in a user-friendly format.
- **Limited Awareness:** Without a clear, engaging way to explore biodiversity patterns, public awareness and interest in conservation remain limited, reducing the impact of available environmental efforts.

WildStat addresses these problems by providing a centralized and visual approach to biodiversity analysis. It enables quick, informative exploration of ecological data that supports education, conservation planning, and informed public engagement.

### **iii) System Analysis & Design vis-a-vis User Requirements**

#### **System Analysis and Design:**

The **WildStat** application is designed with a strong emphasis on user accessibility, data-driven insights, and educational value. Its architecture supports biodiversity analysis through clean data visualization, intuitive navigation, and dynamic content delivery. The system is built using Python and Jupyter Notebooks, integrated with popular data science libraries like pandas, NumPy, Matplotlib, and Seaborn, ensuring a powerful yet user-friendly experience.

#### **Modular Data Filtering and Input System:**

WildStat allows users to interactively select specific national parks or species categories (e.g., mammals, birds, endangered species) via modular dropdowns and filters. This design ensures users can tailor the analysis to their focus area, making the tool flexible for students, researchers, or conservationists. Each filter updates the underlying dataset dynamically without overwhelming the interface.

#### **Dynamic Data Processing Engine:**

Once a park or species group is selected, the backend Python engine processes biodiversity metrics in real time. It performs aggregation, filtering, and analysis using pandas, and outputs cleaned datasets for visualization. This engine ensures that all outputs are based on the latest user selections and provides meaningful summaries like species counts, threat status breakdowns, or species-per-park density.

#### **Interactive Visualizations and Dashboards:**

WildStat employs **Matplotlib**, **Seaborn**, and optionally **Plotly** to generate insightful charts, including:

- Species count comparisons by park
- Distribution of endangered vs. non-endangered species
- Pie charts showing threat levels or habitat types
- Time-based trends in conservation status (if temporal data is added)

These visualizations help users digest complex ecological data through clean, attractive graphs that support comparison and decision-making.

#### **Real-Time Summary Statistics:**

Key statistics like the number of endangered species, species diversity index, and species richness are calculated instantly and displayed using a structured card layout within the Jupyter Notebook. These summaries change as users modify their filters, creating a responsive analytical experience.

#### **User-Focused Interface and Experience:**

Although built within Jupyter Notebooks, the system is organized into clearly labeled sections with markdown headings, collapsible cells, and instructions to help guide users. This design reduces confusion, especially for non-technical users, making the notebook both a data tool and a learning aid.

#### **Session Persistence via Code Blocks:**

To improve the user experience, WildStat supports pseudo-session persistence by enabling users to re-run cells and update analysis as needed. Users can return to earlier inputs, tweak them, and re-process the outputs, simulating the effect of session management within a notebook-based interface.

#### **Clean Code and Documentation:**

The system includes detailed comments, markdown descriptions, and helper messages within the codebase. This not only assists in understanding the analysis pipeline but also promotes educational use. Users can learn both ecological concepts and coding practices simultaneously.

### **Scalable and Open-Source Architecture:**

As a lightweight, Python-powered solution, WildStat is scalable and open for future expansion. It can be extended to include more datasets (e.g., climate, invasive species), real-time data from APIs, or converted into a Streamlit web application for broader accessibility. The notebook structure allows developers to clone, modify, or integrate it into larger conservation platforms.

### **Responsive and Aesthetic Layout:**

The notebook is structured for readability, with consistent color themes in visualizations, structured sections, and intentional whitespace to avoid clutter. Every graph is labeled, every step is explained, and the overall design promotes an engaging user journey from data to insight.

This system design ensures that WildStat isn't just another data tool—it becomes a platform where biodiversity analysis is made accessible, interactive, and insightful for all kinds of users, from students and educators to policymakers and ecologists.

## **User Requirements Analysis:**

The WildStat application is designed to meet the following essential user requirements, with a strong emphasis on accessibility, usability, and educational value for biodiversity researchers, students, and conservation enthusiasts:

### **Intuitive and User-Friendly Interface:**

WildStat features a clean and organized Jupyter Notebook interface that is easy to navigate for both technical and non-technical users. Through well-structured sections, labeled code cells, and markdown instructions, users can interact with the system effortlessly without prior programming knowledge. The design philosophy prioritizes simplicity, clarity, and engagement.

### **Comprehensive Biodiversity Insights and Filtering:**

Users can filter data based on national parks, conservation status (e.g., endangered, threatened), and species categories (mammals, birds, etc.). This allows for personalized exploration of biodiversity metrics, enabling users to focus on areas or species of specific interest with ease and accuracy.

### **Clear and Understandable Visualizations:**

The system translates raw biodiversity data into easy-to-understand graphs and charts using **Matplotlib** and **Seaborn**. Users can view:

- Species counts per park
- Endangered species distribution
- Pie charts and bar graphs representing threat statuses, habitat types, and park-wise diversity

### **Personalized Feedback and Observations:**

Based on selected inputs, WildStat provides tailored summaries such as species richness, conservation risks, and diversity indicators for specific parks. Users are given contextual insights, such as which parks have the highest biodiversity or face the greatest conservation challenges, fostering a deeper understanding of regional ecological health.

### **Actionable and Educational Insights:**

While the tool is analytical in nature, it also serves as an educational platform. Users can draw meaningful conclusions from data, such as identifying hotspots for endangered species or recognizing patterns in species distribution, which may inform future conservation decisions, reports, or field research.

### **Progressive Expansion Potential:**

Although the current version of WildStat focuses on static datasets, future enhancements could include real-time data updates, integration with APIs (e.g., for environmental conditions), or additional analytical features like biodiversity trend tracking over time or invasive species mapping.

### **Educational Support and Explanations:**

The application includes markdown-based explanations and in-line documentation that break down ecological terms and statistical concepts used in the analysis. This makes the system a valuable learning resource for users unfamiliar with biodiversity science or Python-based data analysis.

### **Data Privacy and Transparency:**

WildStat does not collect any personal data. All data manipulation happens locally within the notebook using open-source datasets. This ensures that user interaction remains private, transparent, and fully under their control.

### **Cross-Platform Compatibility:**

As a Jupyter Notebook-based tool, WildStat can be accessed and executed on various platforms, including Windows, macOS, Linux, and cloud-based environments like Google Colab. This guarantees flexibility for users to explore biodiversity data across multiple devices.

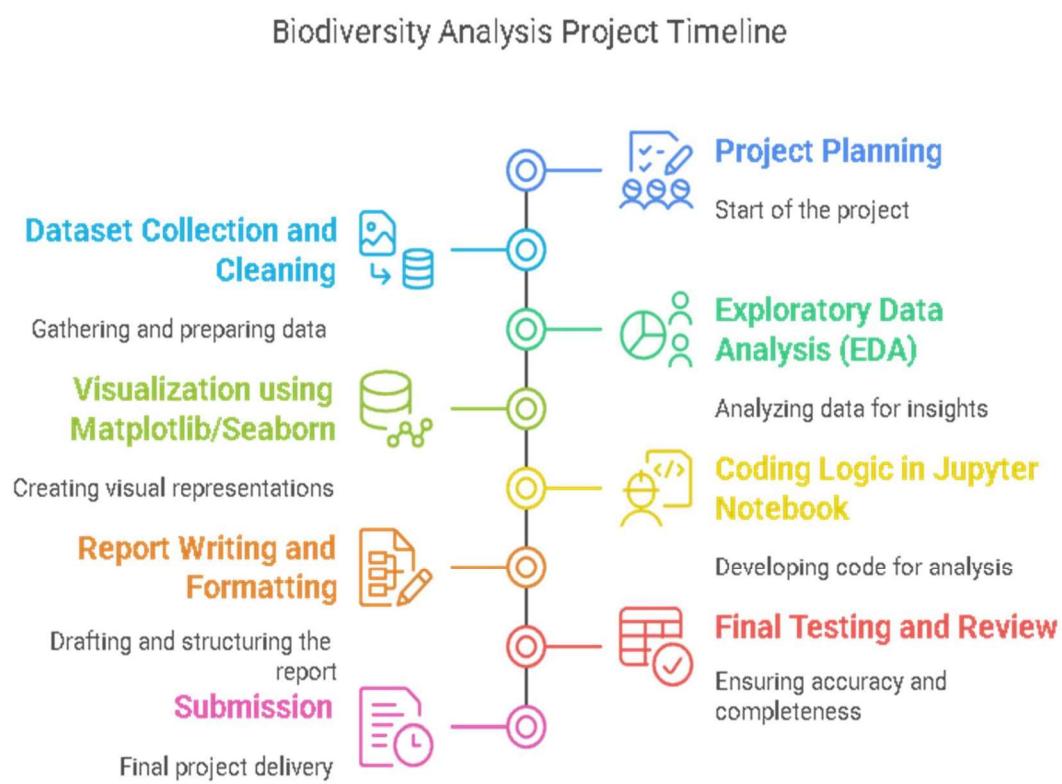
### **Flexible Unit Display and Interpretation:**

WildStat is designed to display data in intuitive, region-specific units (e.g., number of species, percentage of endangered species, etc.), with clear legends and labels on graphs to avoid confusion. This helps users interpret ecological statistics with ease, regardless of scientific background.

### **In-Notebook Help and Guidance:**

Inline markdown instructions and comments within the code serve as contextual help, guiding users through each analysis step. These tooltips and helper notes ensure clarity and reduce the chance of misinterpretation or errors during interaction.

#### iv) System Planning (PERT Chart)



**v) Methodology adopted; System Implementation & Details of Hardware & Software used System Maintenance & Evaluation**

**Methodology Adopted**

The development of the *WildStat – Biodiversity Analysis of U.S. National Parks* project followed a structured and agile-inspired methodology, ensuring flexibility, user relevance, and data accuracy throughout the system's lifecycle. The methodologies used include:

**Iterative Prototyping:** Development began with foundational tasks such as data loading, cleaning, and initial visualization. As the project progressed, new functionalities like filtering endangered species by park, category, or conservation status were added in successive iterations. This allowed for feedback-driven adjustments and continuous enhancement of data insights and visuals.

**User-Centered Design:**

The notebook and visualizations were created keeping in mind end users such as environmental researchers, students, or park services. Simplicity, clarity, and usability were emphasized through readable plots, structured dataframes, and clean layouts to support comprehension and decision-making.

**Data-Driven Approach:**

The project is entirely based on publicly available biodiversity datasets, including endangered species data from U.S. National Parks. Analytical decisions—such as species distribution, frequency of endangerment by category, and park-wise impact—were derived from real data using statistical and exploratory techniques.

**Descriptive and Analytical Approach:**

The project is both descriptive (explaining the structure, sources, and features of the dataset) and analytical (analyzing species risk distribution, park-specific conservation concerns, and visualization of key biodiversity patterns). Techniques such as bar plots, pie charts, and histograms were used to provide meaningful insights.

## **System Implementation & Details of Hardware & Software Used**

### **Software Used:**

Data Analysis Platform: The project was implemented in a Jupyter Notebook environment using Python for all stages of data handling and visualization.

### **Key Libraries Used:**

- pandas for data manipulation
- NumPy for numerical computations
- Matplotlib and Seaborn for generating visualizations

### **Data Used:**

- Endangered species dataset from the U.S. National Parks Service (available publicly or via Kaggle)

### **Core Functionalities:**

- Importing and cleaning biodiversity data
- Filtering and sorting endangered species by park and category
- Visualization of species risk trends across different locations
- Park-wise insights into conservation priorities

### **Hardware and Software Used:**

#### **Hardware:**

- Standard personal laptop or desktop with at least:
- 4GB RAM (8GB recommended)
- Intel i3 or equivalent processor
- Internet access for downloading datasets and libraries

#### **Software:**

- Python 3.x
- Jupyter Notebook (via Anaconda or standalone)
- Libraries: pandas, numpy, matplotlib, seaborn
- Operating System: Windows 10

## **System Maintenance & Evaluation**

### **Maintenance:**

- Regular updating of datasets if new conservation data becomes available.
- Code modularity ensures that future updates (new features or extended datasets) can be integrated with minimal rework. Clear documentation in cells and markdown helps with long-term maintainability and handovers.

### **Evaluation:**

- The system is evaluated based on:
- Accuracy of filtering and visual analysis
- Clarity and readability of visual outputs
- Responsiveness of plots and transformations
- User interpretation — ease of understanding biodiversity trends from visual insights

Future extensions could include interactive dashboards using Streamlit or Plotly Dash, integration of real-time data (e.g., endangered species updates), and machine learning-based species threat predictions.

## **System Evaluation:**

The effectiveness and reliability of the *WildStat – Biodiversity Analysis of U.S. National Parks* application can be evaluated using the following criteria:

### **Usability Testing:**

Although formal usability testing may not have been conducted in this academic phase, potential future deployment should include sessions with environmental science students, researchers, and conservation professionals. This would help assess the clarity of visualizations, ease of navigation in the notebook, and overall user comprehension of biodiversity trends across national parks.

### **User Engagement Metrics**

In a web-based or dashboard version of WildStat, engagement can be measured through user interaction data such as:

- Number of parks or species filtered/viewed
- Frequency of use over time
- Interaction with different visualization types (e.g., pie charts vs. bar graphs)
- These insights would help identify which areas of biodiversity data are of most interest and inform future improvements.

If shared as a public resource or interactive tool, embedding a simple feedback form or link can collect user opinions on:

- The usefulness of species/category insights
- Suggestions for additional filters or data visualizations
- Understanding of conservation status and biodiversity threats

WildStat is based on verified U.S. National Parks biodiversity data. However, evaluating how well the system reflects actual on-ground biodiversity trends could involve:

- Comparing the visualized results with government reports or field surveys
- Consulting subject matter experts for review of conclusions drawn from the visual data

While not resource-intensive in its current notebook format, performance can be evaluated by:

- Testing execution time for larger datasets
- Ensuring visualizations remain responsive with filters and grouped data
- Checking scalability for possible deployment as a web application (e.g., using Streamlit or Flask)
- Alignment with Objectives

The project can be evaluated by checking how well it meets its intended goals:

- Objective: Make biodiversity patterns in U.S. national parks clear and understandable — *Evaluation:* Effective use of pie/bar charts, filtering options, and park-wise breakdowns.
- Objective: Help identify endangered species and high-risk parks — *Evaluation:* Accurate categorization and species count summaries support this aim.

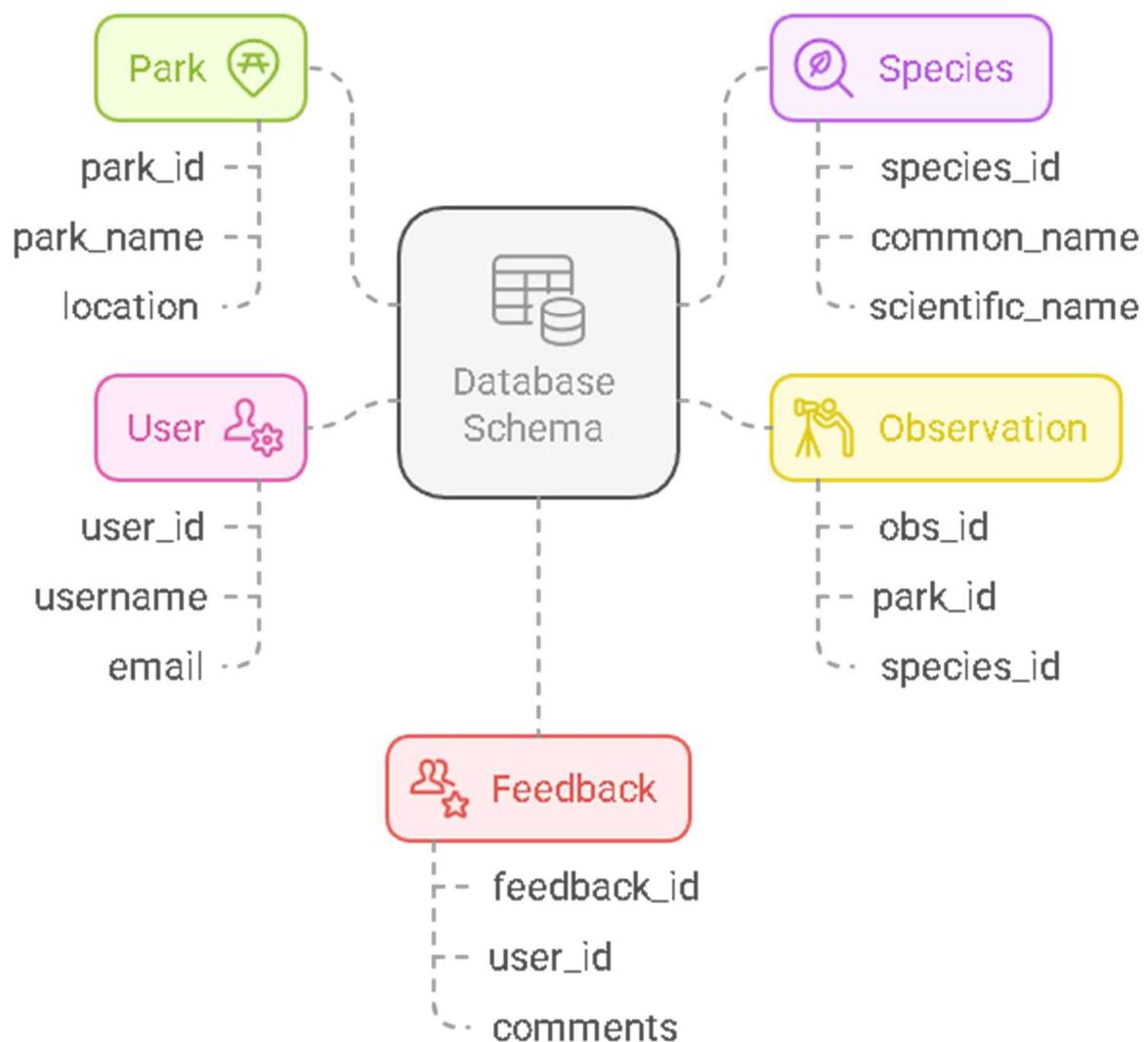
Visualization quality is critical for comprehension. Effectiveness is determined by:

- Clear labeling of axes, legends, and titles
- Logical color schemes for conservation status or category

## vi) Detailed Life Cycle of the Project

### a) ERD

Database Schema for Park and Wildlife Management



## **Detailed Explanation of the ERD:**

This ERD models the data structure of WildStat, a biodiversity and conservation tracking platform focused on analyzing species data in U.S. National Parks. It supports features like endangered species monitoring, park-wise insights, and user interaction. Below is the entity-wise breakdown and relationship mapping.

### **1.) parks**

- id (string, pk): Unique identifier for each national park.
- name (string): Name of the national park.
- location (string): General location or address of the park.
- state (string): U.S. state where the park is located.
- area\_acres (float): Total area of the park in acres.

**Purpose:** Stores geographic and descriptive information about the national parks being analyzed.

### **2.) species**

- id (string, pk): Unique identifier for each species.
- common\_name (string): Commonly used name of the species.
- scientific\_name (string): Scientific/Latin name of the species.
- category (string): Type or group (e.g., Mammal, Bird, Reptile).
- conservation\_status (string): Status like "Endangered", "Threatened", etc.

**Purpose:** Tracks species present in the national parks, along with their classification and conservation info.

### **3.) observations**

- id (string, pk): Unique observation record ID.
- park\_id (string, fk): References the park where the species was observed.
- species\_id (string, fk): References the species being observed.
- presence (boolean): Whether the species was observed or not.
- abundance (string): Estimate of how commonly the species was seen (e.g., "Common", "Rare").

**Purpose:** Records species presence data per park, allowing biodiversity analysis.

### **4.) users**

- id (string, pk): Unique identifier for the user.
- name (string): Full name of the user.
- email (string): Email for login/communication.
- role (string): Role like "admin", "viewer", "researcher".
- created\_at (timestamp): Date the user account was created.

**Purpose:** Stores credentials and access control for users interacting with the system.

## 5) Feedback

- id (string, pk): Unique ID for the feedback entry.
- user\_id (string, fk): References the user giving feedback.
- comments (text): Feedback or suggestion provided by the user.
- rating (int): Rating (e.g., 1–5 stars).
- date (timestamp): Date of the feedback.

Purpose: Collects user suggestions and evaluations for system improvement.

## 6) park\_statistics

- id (string, pk): Unique record ID for each park's statistical snapshot.
- park\_id (string, fk): References the park.
- timestamp (timestamp): When the data was recorded.
- endangered\_species\_count (int): Number of endangered species in the park.
- total\_species\_count (int): Total species observed in the park.
- diversity\_index (float): Biodiversity index (e.g., Simpson's or Shannon Index).

Purpose: Stores analytical and ecological metrics for each park to track biodiversity over time.

## 7) climate\_data

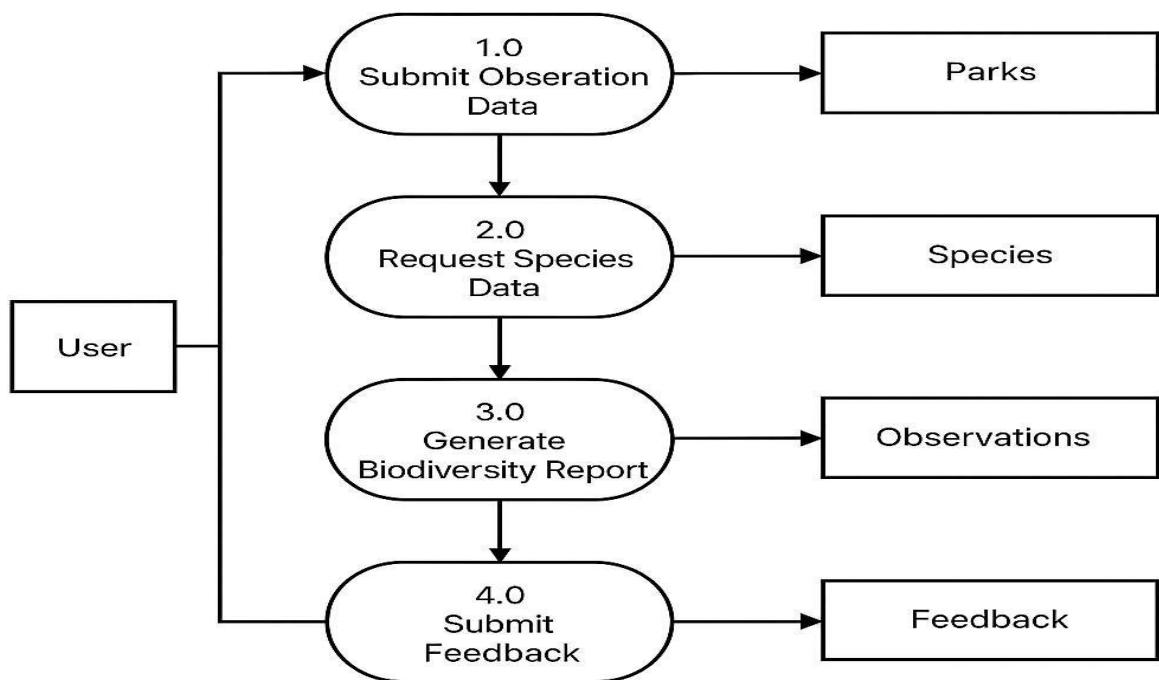
- id (string, pk): Unique ID for each climate record.
- park\_id (string, fk): References the park the data belongs to.
- timestamp (timestamp): Date and time of climate measurement.
- temperature\_celsius (float): Recorded temperature.
- humidity\_percent (float): Humidity level.
- precipitation\_mm (float): Rainfall data.

Purpose: Provides climate context to analyze environmental impact on species diversity and behavior.

## Relationships

Relationship	Type	Description
parks → observations	1-to-many	A park can have many observations.
species → observations	1-to-many	A species can be observed in many parks.
users → feedback	1-to-many	A user can submit multiple feedback entries.
parks → park_statistics	1-to-many	Statistics are tracked over time per park.
parks → climate_data	1-to-many	Multiple climate readings per park.

## DFD



## **Detailed Explanation of the DFD**

The Data Flow Diagram (DFD) for WildStat illustrates how biodiversity data is collected, processed, analyzed, and presented to end users. The system follows a structured workflow to transform raw species data into meaningful conservation insights for U.S. National Parks.

The process begins when researchers or park officials input raw biodiversity data into the system. This data includes species names, conservation statuses, population numbers, and geographic locations within various national parks. The frontend interface provides user-friendly forms for data entry, with validation checks to ensure data quality and consistency. Users can either upload bulk datasets in CSV format or enter information manually through structured forms.

Once received, the backend system processes this raw data through several cleaning and standardization steps. The data processing module handles missing values, corrects taxonomic inconsistencies, and verifies conservation status classifications against official IUCN and NPS standards. The system cross-references species names with authoritative databases to ensure accuracy. This cleaned data is then stored in a structured database, organized by park location, species category, and conservation priority.

For analytical processing, the system employs statistical models and machine learning algorithms. The analytics engine calculates key biodiversity metrics such as species richness, endemism rates, and threat indices. Advanced algorithms identify concerning population trends and flag species that may require conservation attention. These analyses generate processed datasets that serve as the foundation for visualization and reporting.

The visualization module transforms these analytical outputs into interactive dashboards and reports. Using libraries like Matplotlib and Plotly, the system creates dynamic charts that show species distribution patterns, conservation status breakdowns, and temporal trends. Park managers can filter these visualizations by specific taxa, geographic regions, or time periods to focus on areas of concern.

User interaction forms a critical component of the data flow. Authorized personnel can access the system through secure login and query specific information through search interfaces. The system supports various output formats, enabling users to download customized reports containing statistical summaries, visualizations, and conservation recommendations. These reports help inform management decisions and conservation strategies.

For decision support, WildStat incorporates a recommendation engine that suggests priority conservation actions based on the analyzed data. This module considers factors like species vulnerability, habitat quality, and available resources to generate actionable guidance. The system can highlight which parks require immediate intervention and which conservation strategies might be most effective.

Data security and integrity measures are implemented throughout this workflow. The system maintains audit logs of all data modifications and implements role-based access controls to protect sensitive information. Regular backups ensure data preservation, while validation checks at each processing stage maintain data quality.

The DFD demonstrates how WildStat integrates multiple components - from data collection to decision support - into a cohesive system. By streamlining the flow of biodiversity information, the platform enables more efficient monitoring and protection of species across U.S. National Parks. The end-to-end processing transforms disconnected raw data into organized knowledge that directly supports conservation efforts.

Future enhancements could expand this data flow to include real-time monitoring inputs from sensor networks or citizen science reports. The modular design allows for integration of additional data sources

and more sophisticated analytical models as they become available. This scalable architecture ensures WildStat can grow to meet evolving conservation needs while maintaining its core functionality of transforming biodiversity data into actionable insights.

### **System Architecture:**

The architecture of the *WildStat* system consists of a multi-layered structure:

- **Frontend:** The user-facing interface, where data is entered, and results are visualized.
- **Backend:** Coordinates data flow, validates inputs, processes data, and retrieves analysis results.
- **Database:** Stores and manages large sets of biodiversity and environmental data.
- **Analytics Engine:** Handles the statistical and analytical tasks required to process data and generate insights.
- **Machine Learning Model:** Predicts future trends and assesses potential threats to species.

## b) Input and Output

The screenshot shows a Google Colab interface with two tabs open. The top tab, titled 'Biodiversity\_Analysis.ipynb', displays an analysis overview for National Parks. It discusses datasets like 'observations.csv' and 'species\_info.csv' and provides code snippets for importing libraries and reading CSV files. The bottom tab, also titled 'Biodiversity\_Analysis.ipynb', shows a section on understanding data, listing steps for importing libraries and loading datasets. Both tabs show a connection to a Python 3 Google Compute Engine backend.

**Biodiversity in National Parks - Data Analysis**

**Analysis Overview -**

In this analysis, we will be analyzing data from two datasets containing information on National Parks and the species who inhabit them. This analysis is 100% self-guided. Other than the datasets below, no other information was provided for the analysis.

The datasets being analyzed are:

- "observations.csv" - includes data on the national park being observed, what species live there, and how many times each species has been recorded being spotted within the last 7 days.
  - **scientific\_name** - The scientific name of each species
  - **park\_name** - The name of the national park
  - **observations** - The number of observations in the past 7 days
- "species\_info.csv" - includes data on a large number of different species. It features the species' scientific name, common name, category of species, and conservation status.
  - **category** - The category that each species falls under
  - **scientific\_name** - The scientific name of each species
  - **common\_names** - The common names for each species
  - **conservation\_status** - The species conservation status

Both datasets were provided by [Codecademy.com](#)

**Connected to Python 3 Google Compute Engine backend**

**Section 1 - Understanding the Data**

In order to work with any data, you first need to understand what you'll be working with.

1. Import the Python libraries you'll need to conduct your analysis
2. Load the two CSV datasets into two separate Pandas DataFrames, parks and species
3. View the first 10 rows of the parks DataFrame
4. View the first 10 rows of the species DataFrame

```
[3] # 1
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import scipy.stats as ss

[4] # 2
parks = pd.read_csv('observations.csv')
species = pd.read_csv('species_info.csv')
```

**Connected to Python 3 Google Compute Engine backend**

Biodiversity\_Analysis.ipynb - Colab

[colab.research.google.com/drive/163FrkegS4Zq9JryY1chNC77CzBpym0m#scrollTo=qwfNiDae3cnT](https://colab.research.google.com/drive/163FrkegS4Zq9JryY1chNC77CzBpym0m#scrollTo=qwfNiDae3cnT)

Gmail YouTube My Perfect Hotel Napkin AI All Bookmarks

Biodiversity\_Analysis.ipynb File Edit View Insert Runtime Tools Help

Commands + Code + Text

```
# 2
parks = pd.read_csv('observations.csv')
species = pd.read_csv('species_info.csv')

# 3
{x}
parks.head(10)
```

	scientific_name	park_name	observations
0	Vicia benghalensis	Great Smoky Mountains National Park	68
1	Neovison vison	Great Smoky Mountains National Park	77
2	Prunus subcordata	Yosemite National Park	138
3	Abutilon theophrasti	Bryce National Park	84
4	Gilhpopsis specularioides	Great Smoky Mountains National Park	85
5	Elymus virginicus var. virginicus	Yosemite National Park	112
6	Spizella pusilla	Yellowstone National Park	228
7	Elymus multiseta	Great Smoky Mountains National Park	39
8	Lysimachia quadrifolia	Yosemite National Park	168
9	Diphyscium cumberlandianum	Yellowstone National Park	250

Connected to Python 3 Google Compute Engine backend

Biodiversity\_Analysis.ipynb - Colab

[colab.research.google.com/drive/163FrkegS4Zq9JryY1chNC77CzBpym0m#scrollTo=4jPoftVN3cnV](https://colab.research.google.com/drive/163FrkegS4Zq9JryY1chNC77CzBpym0m#scrollTo=4jPoftVN3cnV)

Gmail YouTube My Perfect Hotel Napkin AI All Bookmarks

Biodiversity\_Analysis.ipynb File Edit View Insert Runtime Tools Help

Commands + Code + Text

```
# 4
{x}
species.head(10)
```

	category	scientific_name	common_names	conservation_status
0	Mammal	Clethrionomys gapperi gapperi	Gapper's Red-Backed Vole	NaN
1	Mammal	Bos bison	American Bison, Bison	NaN
2	Mammal	Bos taurus	Aurochs, Aurochs, Domestic Cattle (Feral), Dom...	NaN
3	Mammal	Ovis aries	Domestic Sheep, Mouflon, Red Sheep, Sheep (Feral)	NaN
4	Mammal	Cervus elaphus	Wapiti Or Elk	NaN
5	Mammal	Odocoileus virginianus	White-Tailed Deer	NaN
6	Mammal	Sus scrofa	Feral Hog, Wild Pig	NaN
7	Mammal	Canis latrans	Coyote	Species of Concern
8	Mammal	Canis lupus	Gray Wolf	Endangered
9	Mammal	Canis rufus	Red Wolf	Endangered

Connected to Python 3 Google Compute Engine backend

Biodiversity\_Analysis.ipynb - Colab

<https://colab.research.google.com/drive/163FrkegS4Zq9JryY1chNC77CzBpym0m#scrollTo=NSa-W04j3cnY>

Gmail YouTube My Perfect Hotel Napkin AI All Bookmarks

Biodiversity\_Analysis.ipynb Share Gemini RAM Disk

File Edit View Insert Runtime Tools Help

Commands + Code + Text

```
# 6
In [6]: parks.describe(include='all')

Out[6]:
   scientific_name      park_name  observations
  count            23296          23296.000000
  unique           5541             4.000000    NaN
  top  Hypochaeris radicata  Great Smoky Mountains National Park    NaN
  freq              12             5824.000000    NaN
  mean             NaN             NaN  142.287904
  std              NaN             NaN  69.890532
  min              NaN             NaN  9.000000
  25%             NaN             NaN  86.000000
  50%             NaN             NaN 124.000000
  75%             NaN             NaN 195.000000
  max              NaN             NaN 321.000000
```

[9] # 7
In [9]: parks.park\_name.value\_counts()

Connected to Python 3 Google Compute Engine backend

Biodiversity\_Analysis.ipynb - Colab

[https://colab.research.google.com/drive/163FrkegS4Zq9JryY1chNC77CzBpym0m#scrollTo=q2\\_TV87U3cnZ](https://colab.research.google.com/drive/163FrkegS4Zq9JryY1chNC77CzBpym0m#scrollTo=q2_TV87U3cnZ)

Gmail YouTube My Perfect Hotel Napkin AI All Bookmarks

Biodiversity\_Analysis.ipynb Share Gemini RAM Disk

File Edit View Insert Runtime Tools Help

Commands + Code + Text

```
# 8
In [8]:
for num in range(10):
    print(f"Sample {num + 1}")
    parks_sample = parks.scientific_name.sample(10) # sample size of 10 used to conserve space but 20-25 preferred
    print(parks_sample)
    print()

Sample 1
9844      Fraxinus americana
19205     Botrychium lanceolatum
13737     Pohlia wahlenbergii
1625      Gilia capitata ssp. abrotanifolia
15821     Troglodytes pacificus
4085      Carex hirsutella
10490     Phlox subulata
11146      Pinicola enucleator
3708      Linanthus bicolor
12134     Cuscuta californica var. breviflora
Name: scientific_name, dtype: object

Sample 2
10197     Poncirus trifoliata
5837      Callitricha hermaphroditica
4519      Atriplex canescens var. canescens
9246      Sturnella magna
20457     Prunus pensylvanica
22916     Chamaesyce serpyllifolia ssp. serpyllifolia
19751     Carex albicans var. australis
6452      Eriogonum umbellatum var. majus
22556     Bolandra californica
23251     Lolium temulentum
```

Connected to Python 3 Google Compute Engine backend

Biodiversity\_Analysis.ipynb - Colab

File Edit View Insert Runtime Tools Help

Commands + Code + Text

```

18189      Sidalcea reptans
13185      Lepidium virginicum

```

According to the samples, the data is good to go!

Time to check out the species dataframe!

9. Check to see if the species dataframe contains any missing values. Check column names and data types as well

10. Take a look at the summary statistics

11. Find out how many species there are within each category

12. Find out how many species there are within each conservation\_status

13. Find out the number of NaN values in the conservation\_status column, and the percentage

```
# 9
species.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5824 entries, 0 to 5823
Data columns (total 4 columns):
 #   Column       Non-Null Count  Dtype  
--- 
 0   category     5824 non-null   object  
 1   scientific_name  5824 non-null  object  
 2   common_names   5824 non-null   object  
 3   conservation_status 191 non-null   object  
dtypes: object(4)
memory usage: 182.1+ KB
```

✓ Connected to Python 3 Google Compute Engine backend

[11] # 9
species.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5824 entries, 0 to 5823
Data columns (total 4 columns):
 # Column Non-Null Count Dtype 
--- 
 0 category 5824 non-null object 
 1 scientific\_name 5824 non-null object 
 2 common\_names 5824 non-null object 
 3 conservation\_status 191 non-null object 
dtypes: object(4)
memory usage: 182.1+ KB

# 10
species.describe(include='all')

	category	scientific_name	common_names	conservation_status
count	5824	5824	5824	191
unique	7	5541	5504	4
top	Vascular Plant	Columba livia	Brachythecium Moss	Species of Concern
freq	4470	3	7	161

✓ Connected to Python 3 Google Compute Engine backend

Biodiversity\_Analysis.ipynb - Colab

[colab.research.google.com/drive/163FrkegSI4Zq9JryY1chNC77CzBpym0m#scrollTo=RkVIZgLk3cnc](https://colab.research.google.com/drive/163FrkegSI4Zq9JryY1chNC77CzBpym0m#scrollTo=RkVIZgLk3cnc)

Gmail YouTube My Perfect Hotel Napkin AI All Bookmarks

Biodiversity\_Analysis.ipynb File Edit View Insert Runtime Tools Help

Commands + Code + Text

```
[14]: # 12
species.conervation_status.value_counts()
   count
Species of Concern
  Endangered      16
 Threatened       10
 In Recovery       4
dtype: int64
```

```
# 13
nan_count = len(species[species.conervation_status.isnull() == True])
nan_percent = (len(species[species.conervation_status.isnull() == True]) / len(species)) * 100
print(f'NaN values in \'conservation_status\' column: {nan_count} out of {len(species)}')
print(f'Percentage of NaN values in \'conservation_status\' column: {round(nan_percent, 2)}%')

NaN values in 'conservation_status' column: 5633 out of 5824
Percentage of NaN values in 'conservation_status' column: 96.72%
```

Connected to Python 3 Google Compute Engine backend

Biodiversity\_Analysis.ipynb - Colab

[colab.research.google.com/drive/163FrkegSI4Zq9JryY1chNC77CzBpym0m#scrollTo=iB9G2ZoD3cnc](https://colab.research.google.com/drive/163FrkegSI4Zq9JryY1chNC77CzBpym0m#scrollTo=iB9G2ZoD3cnc)

Gmail YouTube My Perfect Hotel Napkin AI All Bookmarks

Biodiversity\_Analysis.ipynb File Edit View Insert Runtime Tools Help

Commands + Code + Text

## Section 2 - Cleaning the Data

Before creating any visualizations of the data, you first need to clean the data to make sure that our visualizations accurately represent what is happening within the dataset.

14. Fill in the NaN values in the conservation\_status column with the string value "Not at risk"

```
[16]: # 14
species.conervation_status.fillna('Not at Risk', inplace=True)
species.conervation_status.value_counts()

<ipython-input-16-5bb5c9ebf8d6>:3: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method. The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves as a copy.

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method([col: value], inplace=True)' or 'df[col] = df[col].method(value)' instead, to perform the operation on the original object.

species.conervation_status.fillna('Not at Risk', inplace=True)
   count
Not at Risk      5633
Species of Concern      161
```

Connected to Python 3 Google Compute Engine backend

After doing this, it wouldn't hurt to check the scientific name column for any duplicate values. If you refer back to code block #10, you can see that the column does have some repeating values. This should be looked into.

15. Figure out how many values are duplicates within the `scientific_name` column  
 16. Find out how many duplicate rows exist within the `species` dataframe, if any

```
[17] # 15
dup = species.scientific_name.duplicated().sum()
print(f'Duplicated values in \'scientific_name\' column: {dup}')

Duplicated values in 'scientific_name' column: 283

[18] # 16
dup = species.duplicated().sum()
print(f'Duplicated rows in \'species\' dataframe: {dup}')

Duplicated rows in 'species' dataframe: 0
```

So there aren't any duplicate rows within the `species` dataframe, but there are 283 duplicates within the `scientific_name` column of the `species` dataframe. This may be due to the fact that more than one person recorded the same species and they recorded different numbers observed in the observations column. Let's find out.

17. Experiment with the `species` dataframe and determine if these duplicate values should be dropped or not.

✓ Connected to Python 3 Google Compute Engine backend

```
[17] # 17
species[species.scientific_name.duplicated()].head()

category      scientific_name      common_names  conservation_status
3017    Mammal        Cervus elaphus  Rocky Mountain Elk      Not at Risk
3019    Mammal  Odocoileus virginianus  White-Tailed Deer, White-Tailed Deer      Not at Risk
3020    Mammal          Canis lupus       Gray Wolf, Wolf     In Recovery
3022    Mammal        Puma concolor      Cougar, Mountain Lion, Puma      Not at Risk
3025    Mammal        Lutra canadensis        River Otter      Not at Risk

[20] # 17 (Cont.)
species[species.scientific_name == 'Canis lupus']

category      scientific_name      common_names  conservation_status
  8    Mammal          Canis lupus       Gray Wolf      Endangered
3020    Mammal          Canis lupus       Gray Wolf, Wolf     In Recovery
4448    Mammal          Canis lupus       Gray Wolf, Wolf      Endangered
```

This doesn't seem right. This is the same mammal, recorded 3 times, and once with a different conservation status. The only sensible

✓ Connected to Python 3 Google Compute Engine backend

Biodiversity\_Analysis.ipynb - Colab

[colab.research.google.com/drive/163FrkegS4Zq9JryY1chNC77CzBpym0m#scrollTo=iB9G2ZoD3onc](https://colab.research.google.com/drive/163FrkegS4Zq9JryY1chNC77CzBpym0m#scrollTo=iB9G2ZoD3onc)

Gmail YouTube My Perfect Hotel Napkin AI All Bookmarks

Biodiversity\_Analysis.ipynb File Edit View Insert Runtime Tools Help

Commands + Code + Text

20. The same `scientific_name` shouldn't be recorded at the same `park_name` more than once. Find out how many duplicates you have when you check the duplicates between `scientific_name` and `park_name`

21. There is a lot of duplicates. Experiment with the `parks` datafram and determine if these duplicate values should be dropped or not

```
[26] # 20
dup = parks[['scientific_name', 'park_name']].duplicated().sum()
print(f'Duplicates: {dup}')

Duplicates: 1117
```

```
[27] # 21
parks[parks[['scientific_name', 'park_name']].duplicated()]

  scientific_name          park_name  observations
483      Agrostis gigantea  Yellowstone National Park      235
490      Agrostis mertensii    Yosemite National Park      128
945       Rumex crispus  Yellowstone National Park      255
1213     Dianthus barbatus      Bryce National Park      110
1259      Riparia riparia      Bryce National Park       91
...           ...           ...

```

Connected to Python 3 Google Compute Engine backend

Biodiversity\_Analysis.ipynb - Colab

[colab.research.google.com/drive/163FrkegS4Zq9JryY1chNC77CzBpym0m#scrollTo=iB9G2ZoD3onc](https://colab.research.google.com/drive/163FrkegS4Zq9JryY1chNC77CzBpym0m#scrollTo=iB9G2ZoD3onc)

Gmail YouTube My Perfect Hotel Napkin AI All Bookmarks

Biodiversity\_Analysis.ipynb File Edit View Insert Runtime Tools Help

Commands + Code + Text

23267	Oxalis corniculata	Yosemite National Park	164
23273	Dactylis glomerata	Bryce National Park	89
23280	Botrychium simplex	Yellowstone National Park	241

1117 rows × 4 columns

```
[28] # 21 (Cont.)
parks[parks.scientific_name == 'Agrostis gigantea']

  scientific_name          park_name  observations
449      Agrostis gigantea  Yellowstone National Park      253
483      Agrostis gigantea  Yellowstone National Park      235
6824     Agrostis gigantea      Bryce National Park      104
7763     Agrostis gigantea  Great Smoky Mountains National Park      93
8676     Agrostis gigantea      Bryce National Park      116
11602    Agrostis gigantea  Great Smoky Mountains National Park      57
13907    Agrostis gigantea    Yosemite National Park      148
17535    Agrostis gigantea    Yosemite National Park      128

```

```
[29] # 21 (Cont.)
```

Connected to Python 3 Google Compute Engine backend

Biodiversity\_Analysis.ipynb - Colab

File Edit View Insert Runtime Tools Help

Commands + Code + Text

```
1259 Riparia riparia Bryce National Park 89
```

```
# 22 (Cont.)
parks[(parks.scientific_name == 'Agrostis gigantea') & (parks.park_name == 'Yellowstone National Park')]
```

	scientific_name	park_name	observations
449	Agrostis gigantea	Yellowstone National Park	244
483	Agrostis gigantea	Yellowstone National Park	244

```
[33] # 22 (Cont.)
parks[(parks.scientific_name == 'Riparia riparia') & (parks.park_name == 'Bryce National Park')]
```

	scientific_name	park_name	observations
872	Riparia riparia	Bryce National Park	89
1259	Riparia riparia	Bryce National Park	89

Everything checks out! The duplicate `scientific_name` and `park_name` values now have duplicate `observations` values as well, which are equal to the average of the observations between the duplicates. Now we can simply just drop the duplicate rows of the dataframe since they are equal.

23. Now that the values have been correctly updated, you can drop the duplicates

The species and parks dataframes cleaned up and ready for further analysis!

24. Do a final summary statistics check on the species dataframe  
 25. Do a final summary statistics check on the parks dataframe

```
# 24
species.info()
```

	Column	Non-Null Count	Dtype
0	category	5541	non-null object
1	scientific_name	5541	non-null object
2	common_names	5541	non-null object
3	conservation_status	5541	non-null object

```
[36] # 24 (Cont.)
species.describe(include='all')
```

	category	scientific_name	common_names	conservation_status
count	5541	5541	5541	5541
unique	7	5541	5229	5

Biodiversity\_Analysis.ipynb - Colab

[colab.research.google.com/drive/163FrkegS4Zq9JryY1chNC77CzBpym0m#scrollTo=iB9G2ZoD3onc](https://colab.research.google.com/drive/163FrkegS4Zq9JryY1chNC77CzBpym0m#scrollTo=iB9G2ZoD3onc)

Gmail YouTube My Perfect Hotel Napkin AI All Bookmarks

## Biodiversity\_Analysis.ipynb

File Edit View Insert Runtime Tools Help

Commands + Code + Text

max NaN NaN 321.000000

The dataframes look great! Each dataframe contains 0 null values and 5541 unique values for the `scientific_name` column. The only time there will be a duplicate `scientific_name` value in the parks dataframe is when the scientific name corresponds to a different `park_name` value. Therefore, since there are four unique `park_name` values, there can be a maximum of four identical `scientific_name` values corresponding to a different unique `park_name` value with each occurrence.

The final step you need to take before looking further into the data is to merge the two dataframes. This will make plotting our findings much easier.

26. Merge the species and parks dataframes by inner merging species onto parks

```
[39] # 26
parks_species = parks.merge(species)

[40] # 26
parks_species.head()
```

	scientific_name	park_name	observations	category	common_names	conservation_status
0	Vicia benghalensis	Great Smoky Mountains National Park	68	Vascular Plant	Purple Vetch, Reddish Tufted Vetch	Not at Risk
1	Neovison vison	Great Smoky Mountains National Park	77	Mammal	American Mink	Not at Risk

Connected to Python 3 Google Compute Engine backend

Biodiversity\_Analysis.ipynb - Colab

[colab.research.google.com/drive/163FrkegS4Zq9JryY1chNC77CzBpym0m#scrollTo=iB9G2ZoD3onc](https://colab.research.google.com/drive/163FrkegS4Zq9JryY1chNC77CzBpym0m#scrollTo=iB9G2ZoD3onc)

Gmail YouTube My Perfect Hotel Napkin AI All Bookmarks

## Biodiversity\_Analysis.ipynb

File Edit View Insert Runtime Tools Help

Commands + Code + Text

TIME TO MOVE INTO FURTHER EXPLORATION OF THE DATA:

Section 3 - Exploring the Data - Part 1

### Conservation Statuses within Different Parks

With a 'tidy' pair of dataframes, you can now dive into deeper exploration of the data.

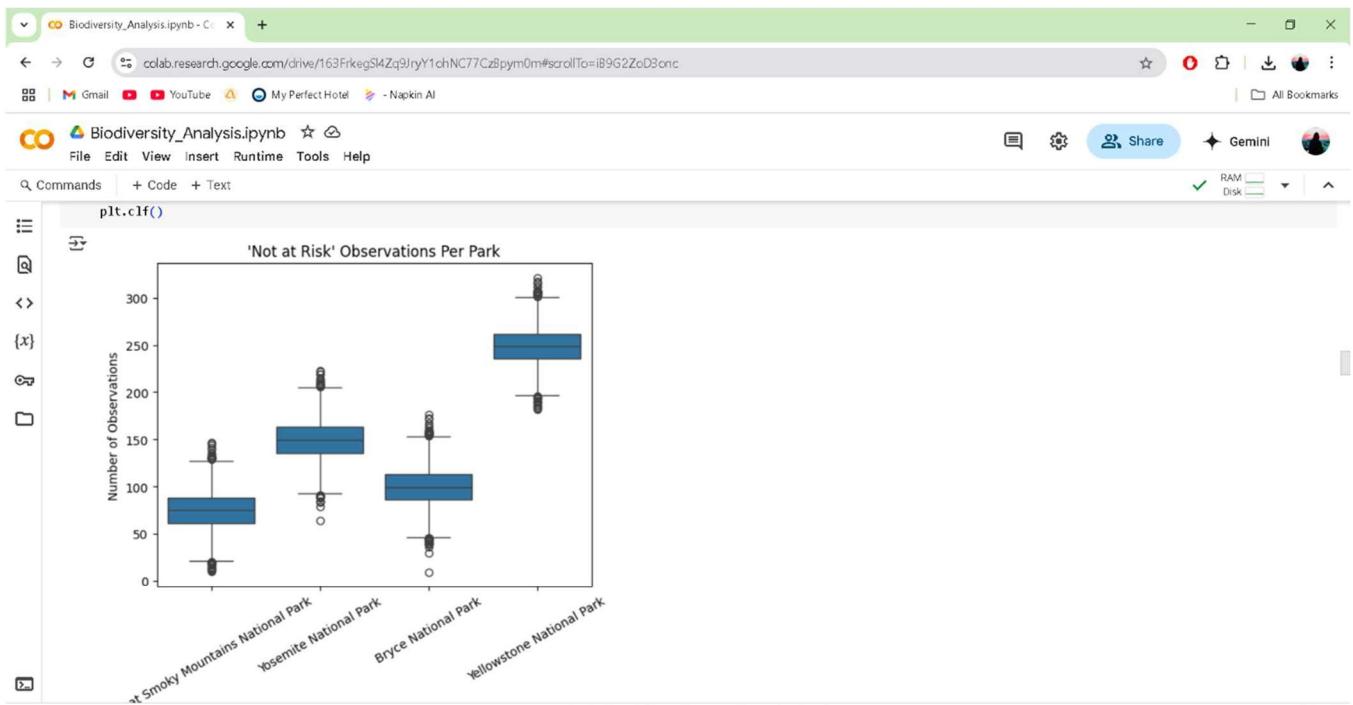
In this analysis, you want to find out if there are any common factors amongst endangered species in order to prevent other species becoming endangered in the future. For this section, the focus will be on the observations of different conservation statuses within different parks.

27. Find out which parks have the highest amount of 'Not at Risk' species  
 28. Look at the total observations for all `conservation_status` values within every park

```
[41] # 27
not_at_risk = parks_species[parks_species.conervation_status == 'Not at Risk']

[42] # 27
sns.boxplot(data=not_at_risk, x='park_name', y='observations')
ax = plt.subplot()
```

Connected to Python 3 Google Compute Engine backend



```
✓ Connected to Python 3 Google Compute Engine backend
```

Great Smoky Mountains National Park  
<Figure size 640x480 with 0 Axes>

You can see that the highest number of 'Not at Risk' species observations is in Yellowstone National Park. This chart is really just for curiosity because this doesn't necessarily mean that Yellowstone is a 'safer' park for a particular species. This could be for a variety of reasons such as the park could be larger than the other parks or the park having a larger species population, etc.

```
[43] # 28
observations_per_park = parks_species.groupby(['park_name', 'conservation_status']).observations.sum().reset_index()

[44] # 28 (Cont.)
observations_per_park_pivot = observations_per_park.pivot(
    index='park_name',
    columns='conservation_status',
    values='observations'
)

[45] # 28 (Cont.)
observations_per_park_pivot
```

	Endangered	In Recovery	Not at Risk	Species of Concern	Threatened
conservation_status					

```
✓ Connected to Python 3 Google Compute Engine backend
```

Biodiversity\_Analysis.ipynb - Colab

File Edit View Insert Runtime Tools Help

Commands + Code + Text

	Yellowstone National Park	1008	559	1337313	33569	1087
	Yosemite National Park	616	386	799611	20187	672

```
# 28 (Cont.)
# Using .sum() and specifying axis 1 will add up all the values across the rows instead of down the columns
observations_per_park_pivot.sum(axis=1)
```

park_name	
Bryce National Park	548159
Great Smoky Mountains National Park	410825
Yellowstone National Park	1373536
Yosemite National Park	821472

**dtype:** int64

Now that you have a pivot table, you can see the total population for every conservation status in every national park. In order to learn more from the pivot table, you can calculate the total observations for every conservation status and then find the percentages of the different conservation statuses within each park.

You don't just want to make calculations; it would be nice to save the calculations into a variable as well. You can update the entire pivot

Connected to Python 3 Google Compute Engine backend

```
30. Update the new pivot table variable with percentages for every conservation status in each park
```

```
# 29
observations_per_park_pivot_percent = observations_per_park_pivot.copy()
```

```
[48] # 30
bryce_park_total = [observations_per_park_pivot.sum(axis=1).index[0], observations_per_park_pivot.sum(axis=1)[0]]
smokies_park_total = [observations_per_park_pivot.sum(axis=1).index[1], observations_per_park_pivot.sum(axis=1)[1]]
yellowstone_park_total = [observations_per_park_pivot.sum(axis=1).index[2], observations_per_park_pivot.sum(axis=1)[2]]
yosemite_park_total = [observations_per_park_pivot.sum(axis=1).index[3], observations_per_park_pivot.sum(axis=1)[3]]

print(bryce_park_total)
print(smokies_park_total)
print(yellowstone_park_total)
print(yosemite_park_total)

park_list = [bryce_park_total, smokies_park_total, yellowstone_park_total, yosemite_park_total]
```

```
['Bryce National Park', np.int64(548159)]
['Great Smoky Mountains National Park', np.int64(410825)]
['Yellowstone National Park', np.int64(1373536)]
['Yosemite National Park', np.int64(821472)]
<ipython-input-48-8b7647848d08>:3: FutureWarning: Series._getitem_ treating keys as positions is deprecated. In a future version, integer keys will always be treated as labels
bryce_park_total = [observations_per_park_pivot.sum(axis=1).index[0], observations_per_park_pivot.sum(axis=1)[0]]
<ipython-input-48-8b7647848d08>:4: FutureWarning: Series._getitem_ treating keys as positions is deprecated. In a future version, integer keys will always be treated as labels
smokies_park_total = [observations_per_park_pivot.sum(axis=1).index[1], observations_per_park_pivot.sum(axis=1)[1]]
<ipython-input-48-8b7647848d08>:5: FutureWarning: Series._getitem_ treating keys as positions is deprecated. In a future version, integer keys will always be treated as labels
yellowstone_park_total = [observations_per_park_pivot.sum(axis=1).index[2], observations_per_park_pivot.sum(axis=1)[2]]
```

Connected to Python 3 Google Compute Engine backend

Biodiversity\_Analysis.ipynb - Colab

File Edit View Insert Runtime Tools Help

Commands + Code + Text

Yosemite National Park    0.075%    0.047%    97.339%    2.457%    0.082%

# 30 (cont.)

observations\_per\_park\_pivot

	conservation_status	Endangered	In Recovery	Not at Risk	Species of Concern	Threatened
park_name						
Bryce National Park	402	258	533093	13979	427	
Great Smoky Mountains National Park	294	189	400343	9659	340	
Yellowstone National Park	1008	559	1337313	33569	1087	
Yosemite National Park	616	386	799611	20187	672	

Now you have two separate pivot tables, one with the number of observations per park, and one with the percentage of observations per park!

What is the takeaway from doing this though?

Well, take a look at the pivot table directly above, with the number of observations. It looks as if the 'safest' park for wildlife is Yellowstone. It also looks like that is the 'most dangerous' park for wildlife as well. It's hard to interpret what you're really looking at. This is why you created the pivot table containing percentages instead:

Look at the pivot table above the one with the number of observations, the one that contains the decimal numbers as values (The

Connected to Python 3 Google Compute Engine backend

Section 3 - Exploring the Data - Part 2

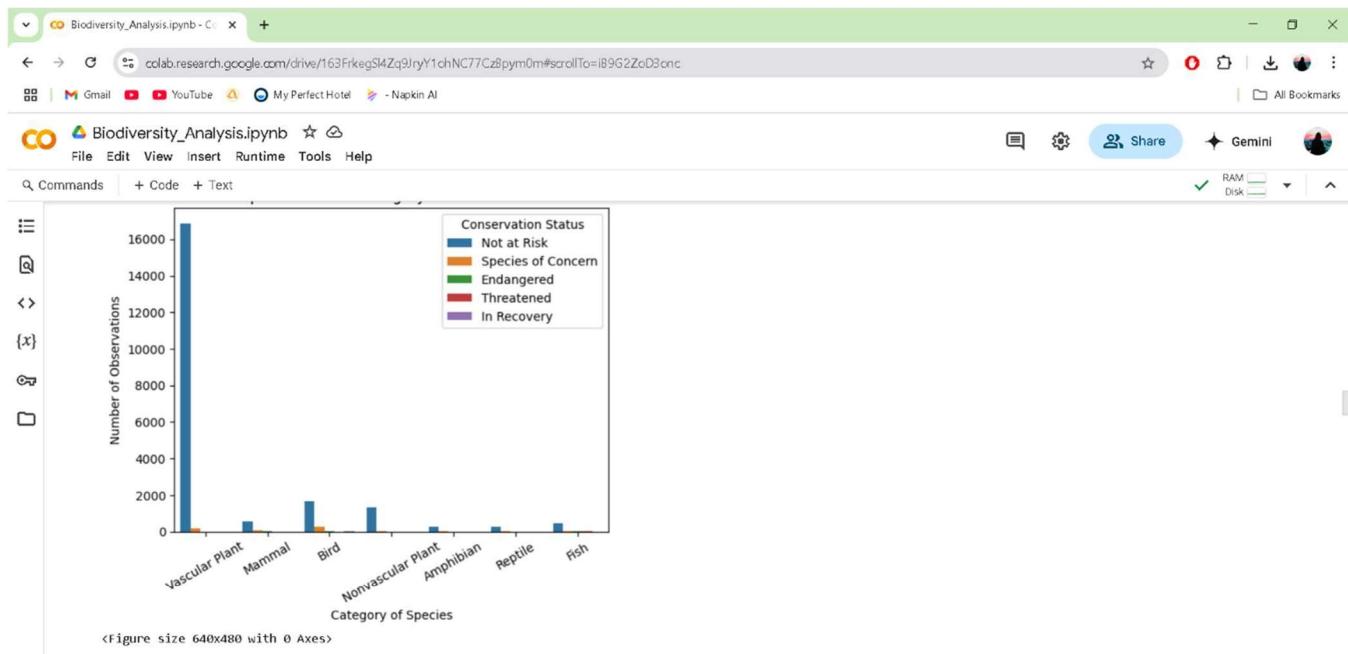
*Plotting the Proportions*

For most of the plots you'll make throughout the analysis you'll need to use a dataframe that excludes the 'Not at Risk' conservation\_status value. This is due to the fact that the proportion of species that are 'Not at Risk' exceeds the other conservation\_status values by so much that the graphs/charts can't be read. This is fine, as long as we understand the actual proportions of the data with the 'Not at Risk' values included.

31. Plot an example of keeping the 'Not at Risk' values in your graphs/charts
32. Create a dataframe that excludes the 'Not at Risk' values

```
[52] # 31
sns.countplot(data=parks_species, x='category', hue='conservation_status')
ax = plt.subplot()
ax.set_xticks(range(len(parks_species.category.unique())))
ax.set_xticklabels(parks_species.category.unique(), rotation=30)
plt.title('# of obsv\'s of species in each category based off conservation status')
plt.ylabel('Number of Observations')
plt.xlabel('Category of Species')
plt.legend(title='Conservation Status')
plt.show()
plt.clf()
```

Connected to Python 3 Google Compute Engine backend



<Figure size 640x480 with 0 Axes>

Not much to see here! Let's fix that!

You now have a dataframe that will be easier to plot with!

Time for some visualizations!

33. Make a figure showing the proportions of observations for each conservation status within each park

```
# 33
plt.figure(figsize=(10, 10))

for i in range(len(pk_sp_nar_excluded.park_name.unique())):
    pk_df = pk_sp_nar_excluded[pk_sp_nar_excluded.park_name == pk_sp_nar_excluded.park_name.unique()[i]]
    pk_df = pk_df[['conservation_status', 'observations']]
    pk_df = pk_df.groupby('conservation_status', sort=False).observations.sum()
    pie_labels = pk_df.index

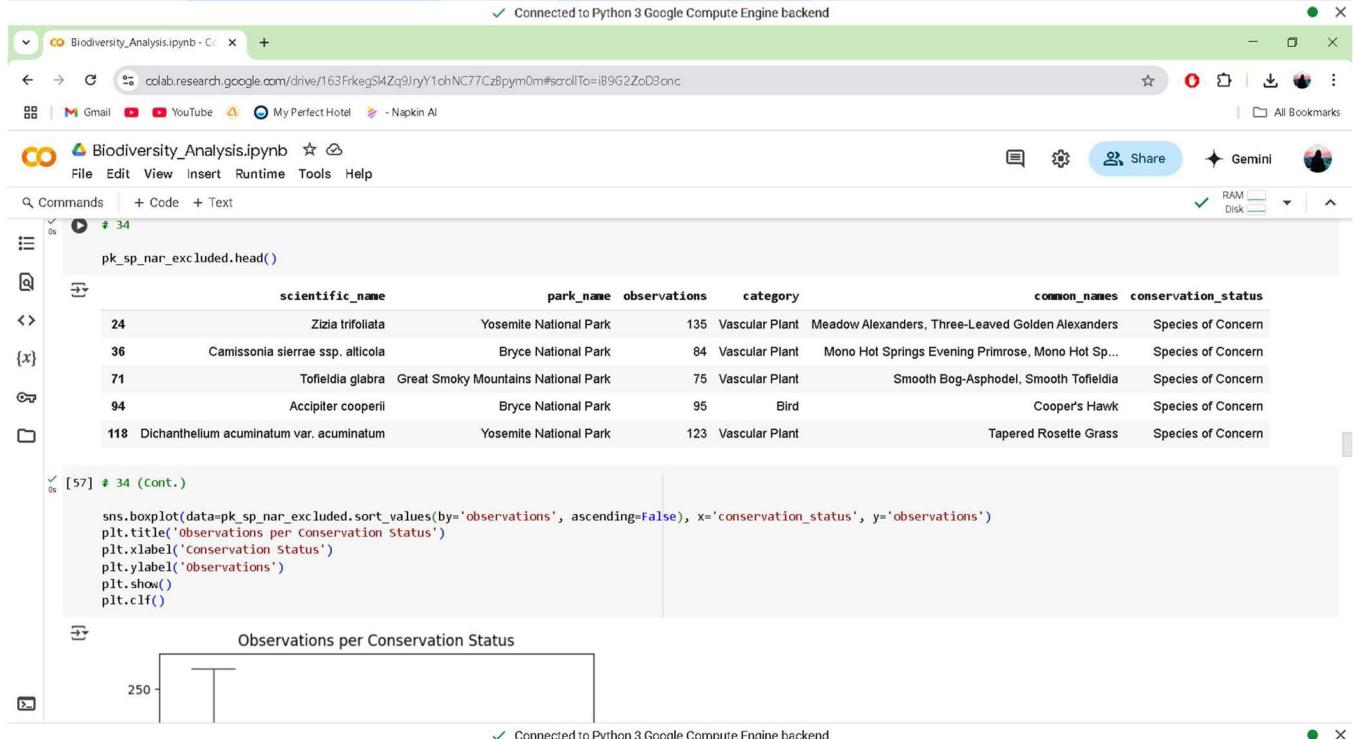
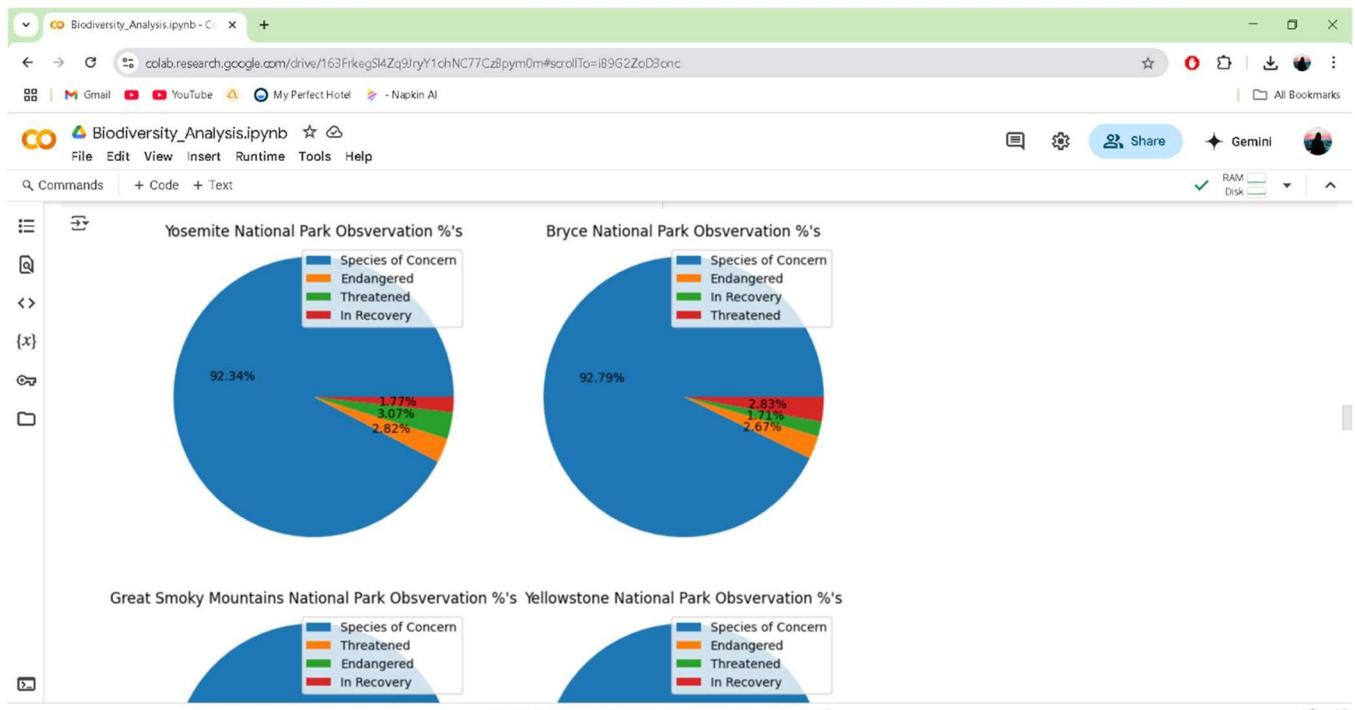
    plt.subplot(2, 2, i + 1)
    plt.pie(pk_df, autopct='%.2f%%')
    plt.axis('equal')
    plt.title(f'(pk_sp_nar_excluded.park_name.unique()[{i}]) Observation %\'s')
    plt.legend(pie_labels)

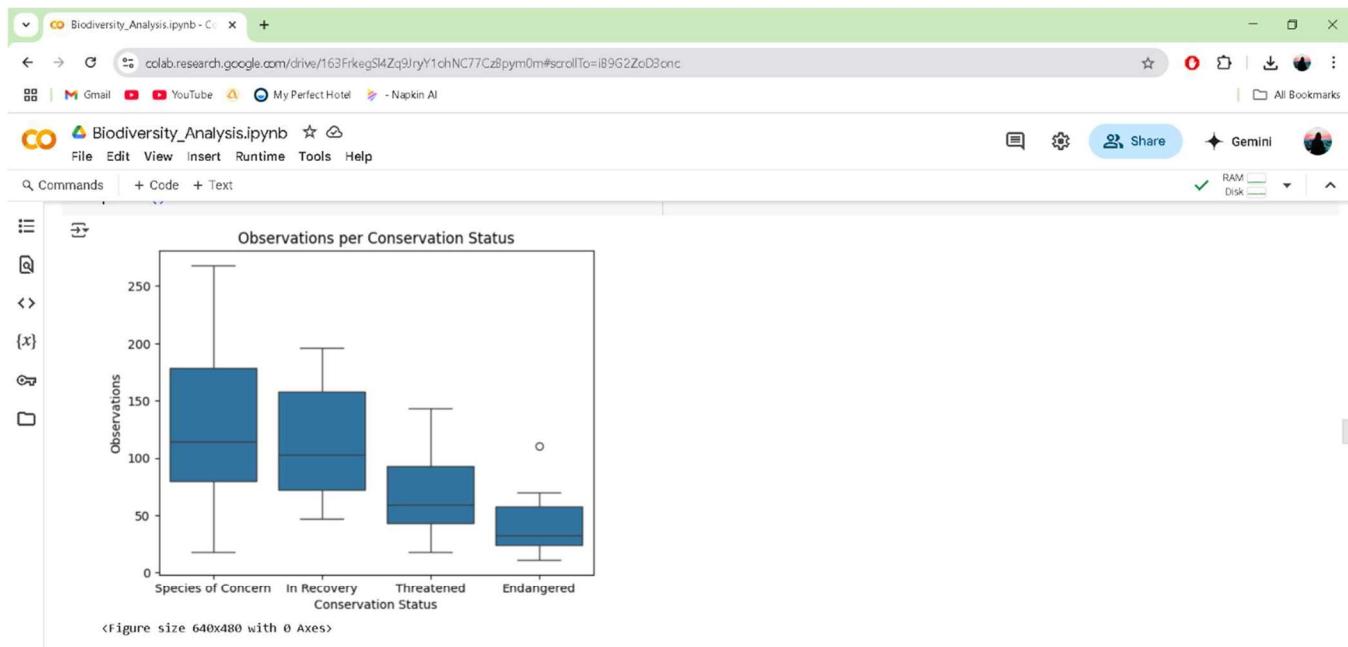
plt.show()
plt.clf()
```

Yosemite National Park Observation %'s

Bryce National Park Observation %'s

Connected to Python 3 Google Compute Engine backend





You can see a clear trend in the boxplot above. The conservation statuses are ordered from the least severe to most severe situation. The graph shows that, on average, as the severity of conservation status increases, the average observations of a particular species decreases.

Connected to Python 3 Google Compute Engine backend

Biodiversity\_Analysis.ipynb - Colab

[colab.research.google.com/drive/163FrkegS4Zq9JryY1chNC77CzBpym0m#scrollTo=iB9G2ZoD3onc](https://colab.research.google.com/drive/163FrkegS4Zq9JryY1chNC77CzBpym0m#scrollTo=iB9G2ZoD3onc)

Gmail YouTube My Perfect Hotel Napkin AI All Bookmarks

Biodiversity\_Analysis.ipynb File Edit View Insert Runtime Tools Help

Commands + Code + Text

Section 3 - Exploring the Data - Part 3

Categories of species

The focus of this section will be directed towards the different categories of species and the conservation statuses of those species.

35. First, create a new variable that contains only the parts of the dataframe that you'll need to use for this section of the analysis

```
[58] # 35
pk_sp_nar_excluded.head()
```

	scientific_name	park_name	observations	category	common_names	conservation_status
24	Zizia trifoliata	Yosemite National Park	135	Vascular Plant	Meadow Alexanders, Three-Leaved Golden Alexanders	Species of Concern
36	Camissonia sierrae ssp. alticola	Bryce National Park	84	Vascular Plant	Mono Hot Springs Evening Primrose, Mono Hot Sp...	Species of Concern
71	Tofieldia glabra	Great Smoky Mountains National Park	75	Vascular Plant	Smooth Bog-Asphodel, Smooth Tofieldia	Species of Concern
94	Accipiter cooperii	Bryce National Park	95	Bird	Cooper's Hawk	Species of Concern
118	Dichanthelium acuminatum var. acuminatum	Yosemite National Park	123	Vascular Plant	Tapered Rosette Grass	Species of Concern

```
[59] # 35 (cont.)
```

Connected to Python 3 Google Compute Engine backend

Biodiversity\_Analysis.ipynb - Colab

[colab.research.google.com/drive/163FrkegSI4Zq9JryY1chNC77Cz8pym0#scrollTo=iB9G2z0D3nc](https://colab.research.google.com/drive/163FrkegSI4Zq9JryY1chNC77Cz8pym0#scrollTo=iB9G2z0D3nc)

Gmail YouTube My Perfect Hotel Napkin AI All Bookmarks

**Biodiversity\_Analysis.ipynb**

File Edit View Insert Runtime Tools Help

Commands + Code + Text

	scientific_name	park_name	category	conservation_status
24	Zizia trifolia	Yosemite National Park	Vascular Plant	Species of Concern
36	Camissonia sierrae ssp. alticola	Bryce National Park	Vascular Plant	Species of Concern
71	Tofieldia glabra	Great Smoky Mountains National Park	Vascular Plant	Species of Concern
94	Accipiter cooperii	Bryce National Park	Bird	Species of Concern
118	Dichanthelium acuminatum var. acuminatum	Yosemite National Park	Vascular Plant	Species of Concern

Great, you won't need the common\_names or observations columns moving forward.

Let's better understand the relationship between category of species and conservation status with some visualizations!

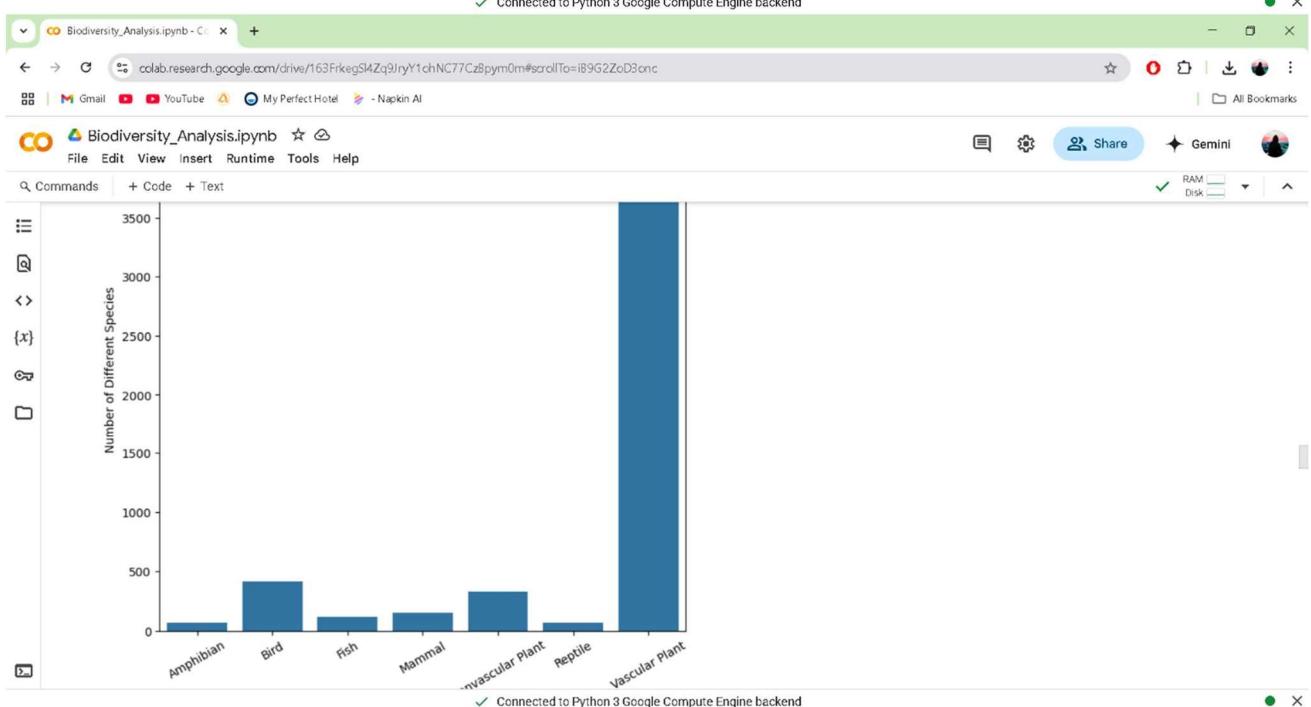
36. Plot the number of unique 'At Risk' species for each conservation status in each category

```
[60]: # Have to drop duplicates or each scientific name will be counted 4 times each
category_status_species = category_df.drop_duplicates(subset=['scientific_name'])

category_status_species = category_status_species.groupby(['category', 'conservation_status']).scientific_name.count().reset_index()

plt.figure(figsize=(8, 8))
sns.barplot(data=category_status_species, x='category', y='scientific_name', hue='conservation_status')
ax = plt.subplot()
ax.set_xticks(range(len(category_status_species.category.unique())))
ax.set_xticklabels(category_status_species.category.unique(), rotation=30)
plt.title('Species per Category based off Conservation Status')
```

Connected to Python 3 Google Compute Engine backend



Biodiversity\_Analysis.ipynb - Colab

[colab.research.google.com/drive/163FrkegS4Zq9JryY1chNC77Cz8pym0m#scrollTo=iB9G2ZoD3onc](https://colab.research.google.com/drive/163FrkegS4Zq9JryY1chNC77Cz8pym0m#scrollTo=iB9G2ZoD3onc)

Gmail YouTube My Perfect Hotel Napkin AI All Bookmarks

Biodiversity\_Analysis.ipynb Share Gemini RAM Disk

```

File Edit View Insert Runtime Tools Help
Commands + Code + Text

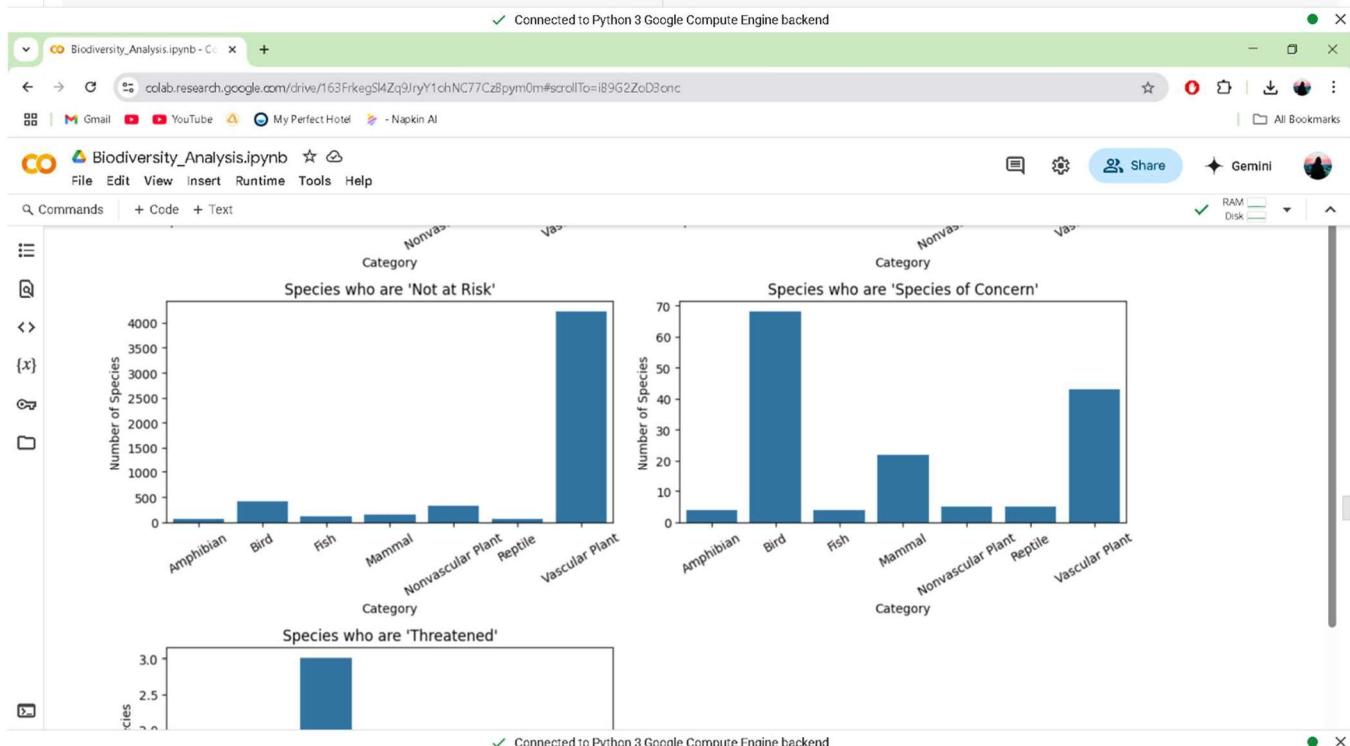
In [63]: parks_species_cat_pivot.fillna(0.0, inplace=True)
parks_species_cat_pivot

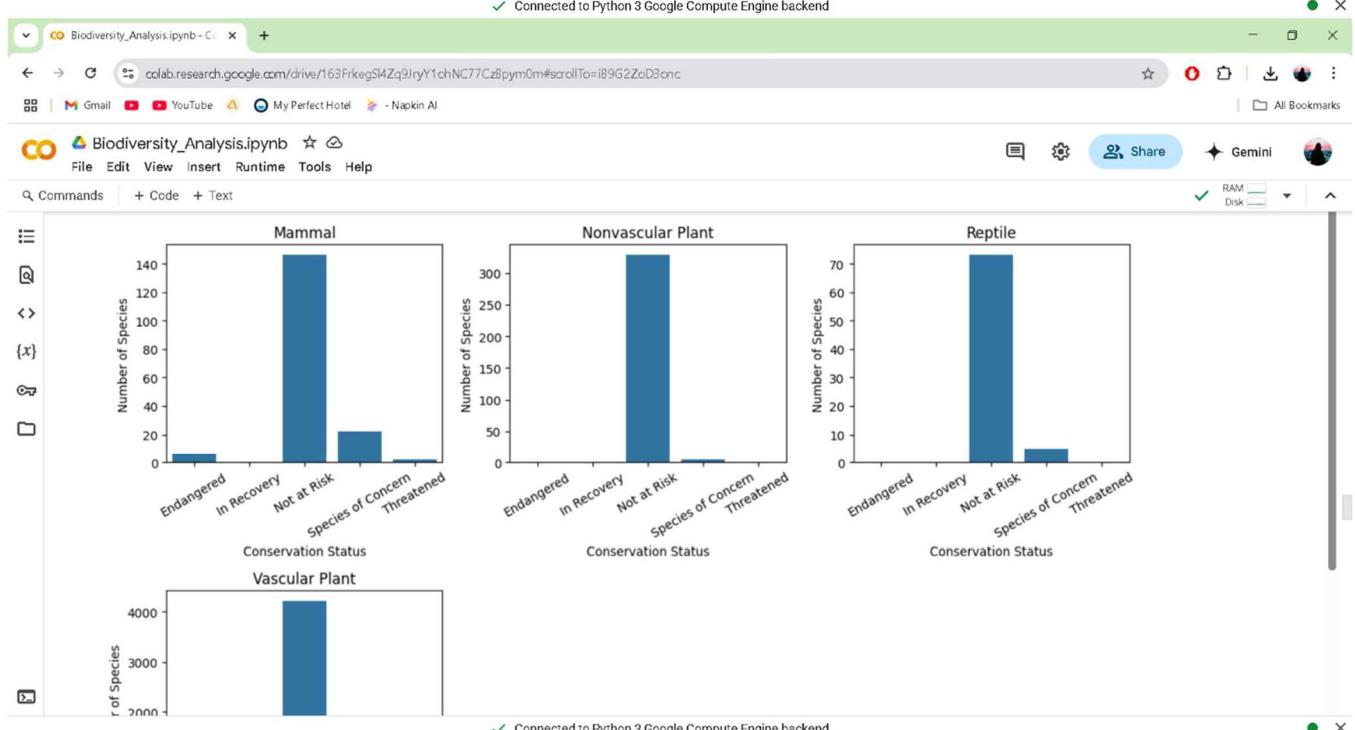
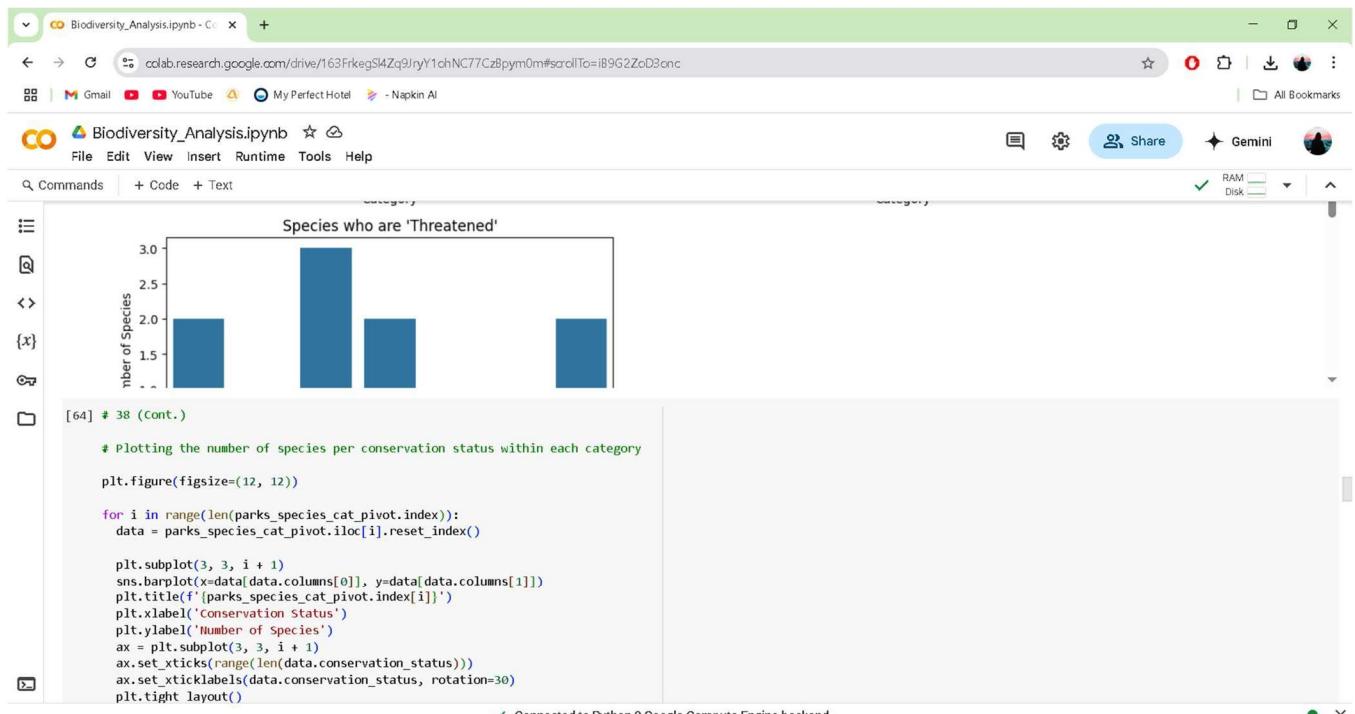
Out[63]:
   conservation_status  Endangered  In Recovery  Not at Risk  Species of Concern  Threatened
category
{x}
  Amphibian           1.0         0.0       72.0        4.0            2.0
  Bird                4.0         3.0      413.0       68.0            0.0
  Fish                3.0         0.0      115.0        4.0            3.0
  Mammal              6.0         0.0      146.0       22.0            2.0
  Nonvascular Plant    0.0         0.0      328.0        5.0            0.0
  Reptile              0.0         0.0      73.0        5.0            0.0
  Vascular Plant       1.0         0.0     4216.0       43.0            2.0

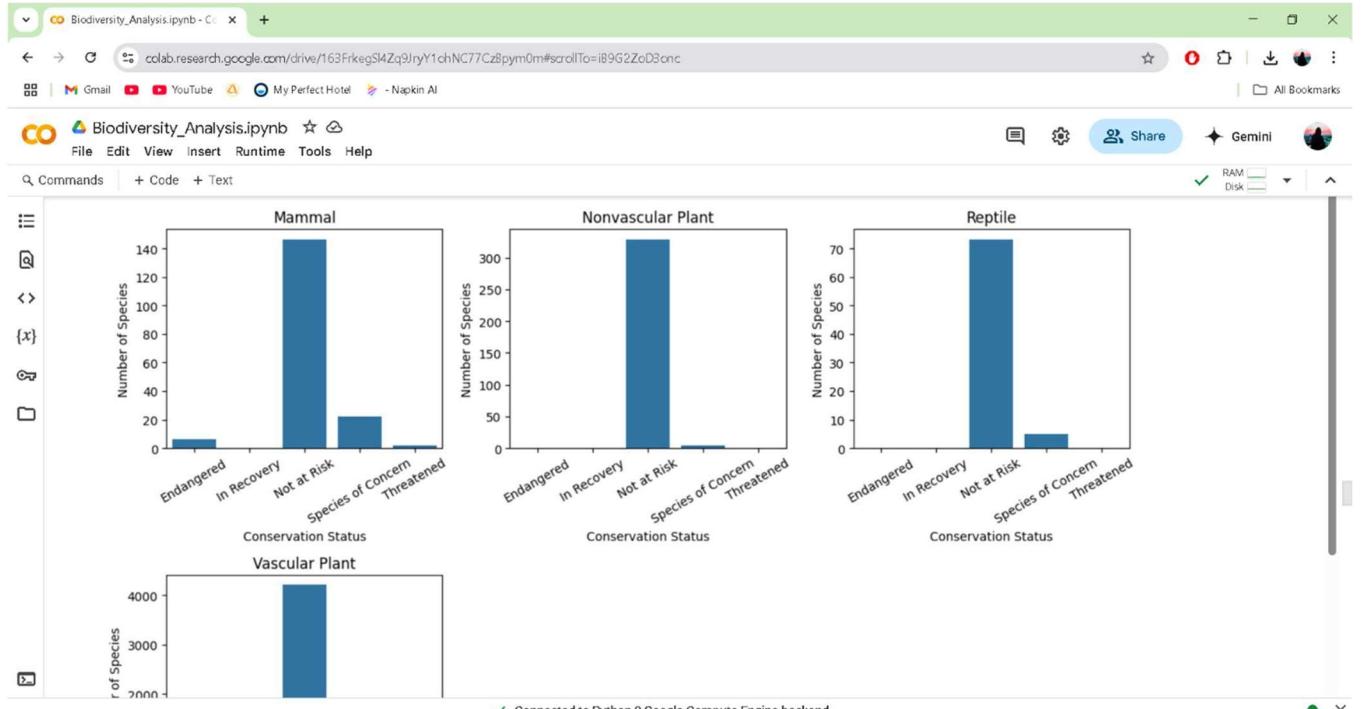
[63]: # 38 (Cont.)
      # Plotting the number of species per category within each conservation status
      plt.figure(figsize=(12, 12))

      for i in range(len(parks_species_cat_pivot.columns)):
          data = parks_species_cat_pivot[parks_species_cat_pivot.columns[i]].reset_index()

```







```
# 39
parks_species_cat_pivot_props = parks_species_cat_pivot.copy()

for cat in parks_species_cat_pivot_props.index:
    cat_total = parks_species_cat_pivot_props.loc[cat].sum(axis=0)
    print(f'Total {cat} Species: {cat_total}')

    parks_species_cat_pivot_props.loc[cat] = parks_species_cat_pivot_props.loc[cat].apply(lambda x: str(round((x / cat_total) * 100, 2)) + '%')

parks_species_cat_pivot_props
```

Total Amphibian Species: 79.0  
Total Bird Species: 488.0  
Total Fish Species: 125.0  
Total Mammal Species: 176.0  
Total Nonvascular Plant Species: 333.0  
Total Reptile Species: 78.0  
Total Vascular Plant Species: 4262.0

<ipython-input-65-7ed4ffd7b90>:10: FutureWarning: Setting an item of incompatible dtype is deprecated and will raise an error in a future version of pandas. Value '1.27%' has c  
parks\_species\_cat\_pivot\_props.loc[cat] = parks\_species\_cat\_pivot\_props.loc[cat].apply(lambda x: str(round((x / cat\_total) \* 100, 2)) + '%')  
<ipython-input-65-7ed4ffd7b90>:10: FutureWarning: Setting an item of incompatible dtype is deprecated and will raise an error in a future version of pandas. Value '0.0%' has dt  
parks\_species\_cat\_pivot\_props.loc[cat] = parks\_species\_cat\_pivot\_props.loc[cat].apply(lambda x: str(round((x / cat\_total) \* 100, 2)) + '%')  
<ipython-input-65-7ed4ffd7b90>:10: FutureWarning: Setting an item of incompatible dtype is deprecated and will raise an error in a future version of pandas. Value '91.14%' has  
parks\_species\_cat\_pivot\_props.loc[cat] = parks\_species\_cat\_pivot\_props.loc[cat].apply(lambda x: str(round((x / cat\_total) \* 100, 2)) + '%')  
<ipython-input-65-7ed4ffd7b90>:10: FutureWarning: Setting an item of incompatible dtype is deprecated and will raise an error in a future version of pandas. Value '5.06%' has c  
parks\_species\_cat\_pivot\_props.loc[cat] = parks\_species\_cat\_pivot\_props.loc[cat].apply(lambda x: str(round((x / cat\_total) \* 100, 2)) + '%')

Biodiversity\_Analysis.ipynb - Colab

<https://colab.research.google.com/drive/163FrkegS4Zq9JryY1chNC77CzBpym0m#scrollTo=iB9G2ZoD3onc>

Gmail YouTube My Perfect Hotel Napkin AI All Bookmarks

Biodiversity\_Analysis.ipynb

File Edit View Insert Runtime Tools Help

Commands + Code + Text

```

parks_species_cat_pivot_props.loc[cat] = parks_species_cat_pivot_props.loc[cat].apply(lambda x: str(round((x / cat_total) * 100, 2)) + '%')
<ipython-input-65-7ed4fd7b90>:10: FutureWarning: Setting an item of incompatible dtype is deprecated and will raise an error in a future version of pandas. Value '2.53%' has c
parks_species_cat_pivot_props.loc[cat] = parks_species_cat_pivot_props.loc[cat].apply(lambda x: str(round((x / cat_total) * 100, 2)) + '%')

conservation_status Endangered In Recovery Not at Risk Species of Concern Threatened
```

category	Amphibian	Bird	Fish	Mammal	Nonvascular Plant	Reptile	Vascular Plant
Amphibian	1.27%	0.0%	91.14%	5.06%	2.53%		
Bird	0.82%	0.61%	84.63%	13.93%	0.0%		
Fish	2.4%	0.0%	92.0%	3.2%	2.4%		
Mammal	3.41%	0.0%	82.95%	12.5%	1.14%		
Nonvascular Plant	0.0%	0.0%	98.5%	1.5%	0.0%		
Reptile	0.0%	0.0%	93.59%	6.41%	0.0%		
Vascular Plant	0.02%	0.0%	98.92%	1.01%	0.05%		

[66] # 39 (Cont.)

```
# Remove the '%' off of each value in the pivot table and convert values to floats for graphing
for cat in parks_species_cat_pivot_props.index:
    parks_species_cat_pivot_props.loc[cat] = parks_species_cat_pivot_props.loc[cat].apply(lambda x: float(x.strip('%')))
```

[67] # 39 (Cont.)

Connected to Python 3 Google Compute Engine backend

Biodiversity\_Analysis.ipynb - Colab

<https://colab.research.google.com/drive/163FrkegS4Zq9JryY1chNC77CzBpym0m#scrollTo=iB9G2ZoD3onc>

Gmail YouTube My Perfect Hotel Napkin AI All Bookmarks

Biodiversity\_Analysis.ipynb

File Edit View Insert Runtime Tools Help

Commands + Code + Text

```
parks_species_cat_pivot_props.loc[cat] = parks_species_cat_pivot_props.loc[cat].apply(lambda x: float(x.strip('%')))

# 39 (Cont.)
```

# Plotting the number of species per category within each conservation status

```
plt.figure(figsize=(12, 12))

for i in range(len(parks_species_cat_pivot_props.columns)):
    data = parks_species_cat_pivot_props[parks_species_cat_pivot_props.columns[i]].reset_index()

    plt.subplot(3, 2, i + 1)
    sns.barplot(x=data[data.columns[0]], y=data[data.columns[1]])
    plt.title(f'Species who are \'{data.columns[1]}\'')
    plt.xlabel('category')
    plt.ylabel('Percent of Species')
    ax = plt.subplot(3, 2, i + 1)
    ax.set_xticks(range(len(parks_species_cat_pivot.index)))
    ax.set_xticklabels(parks_species_cat_pivot.index, rotation=30)
    ax.set_yticks(range(0, 110, 10))
    ax.set_yticklabels(range(0, 110, 10))
    plt.tight_layout()

plt.show()
plt.clf()
```

Species who are 'Endangered'

Species who are 'In Recovery'

Connected to Python 3 Google Compute Engine backend

```

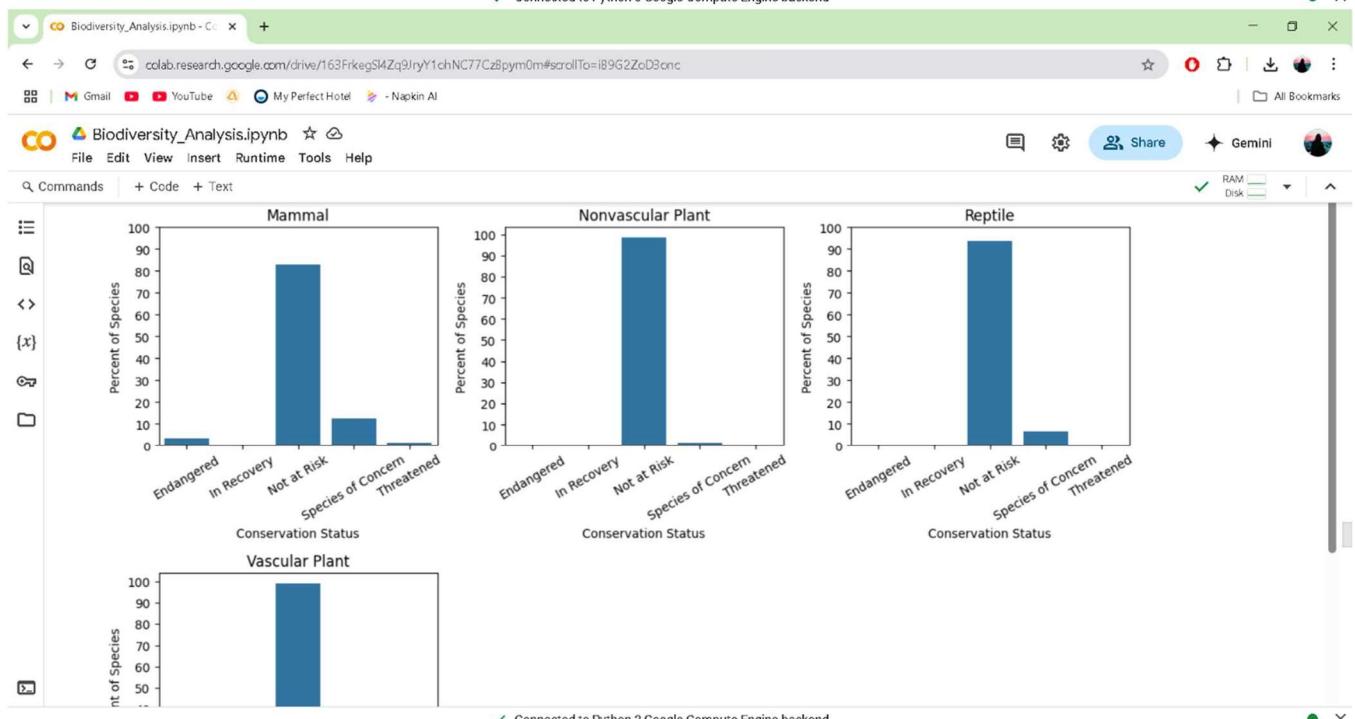
Biodiversity_Analysis.ipynb - Colab Research
File Edit View Insert Runtime Tools Help
File Edit View Insert Runtime Tools Help
# Plotting the number of species per conservation status within each category
plt.figure(figsize=(12, 12))

for i in range(len(parks_species_cat_pivot_props.index)):
    data = parks_species_cat_pivot_props.iloc[i].reset_index()

    plt.subplot(3, 3, i + 1)
    sns.barplot(x=data[data.columns[0]], y=data[data.columns[1]])
    plt.title(f'(parks_species_cat_pivot.index[i])')
    plt.xlabel('Conservation Status')
    plt.ylabel('Percent of Species')
    ax = plt.subplot(3, 3, i + 1)
    ax.set_xticks(range(len(data.conervation_status)))
    ax.set_xticklabels(data.conervation_status, rotation=30)
    ax.set_yticks(range(0, 110, 10))
    ax.set_yticklabels(range(0, 110, 10))
    plt.tight_layout()

plt.show()
plt.clf()

```



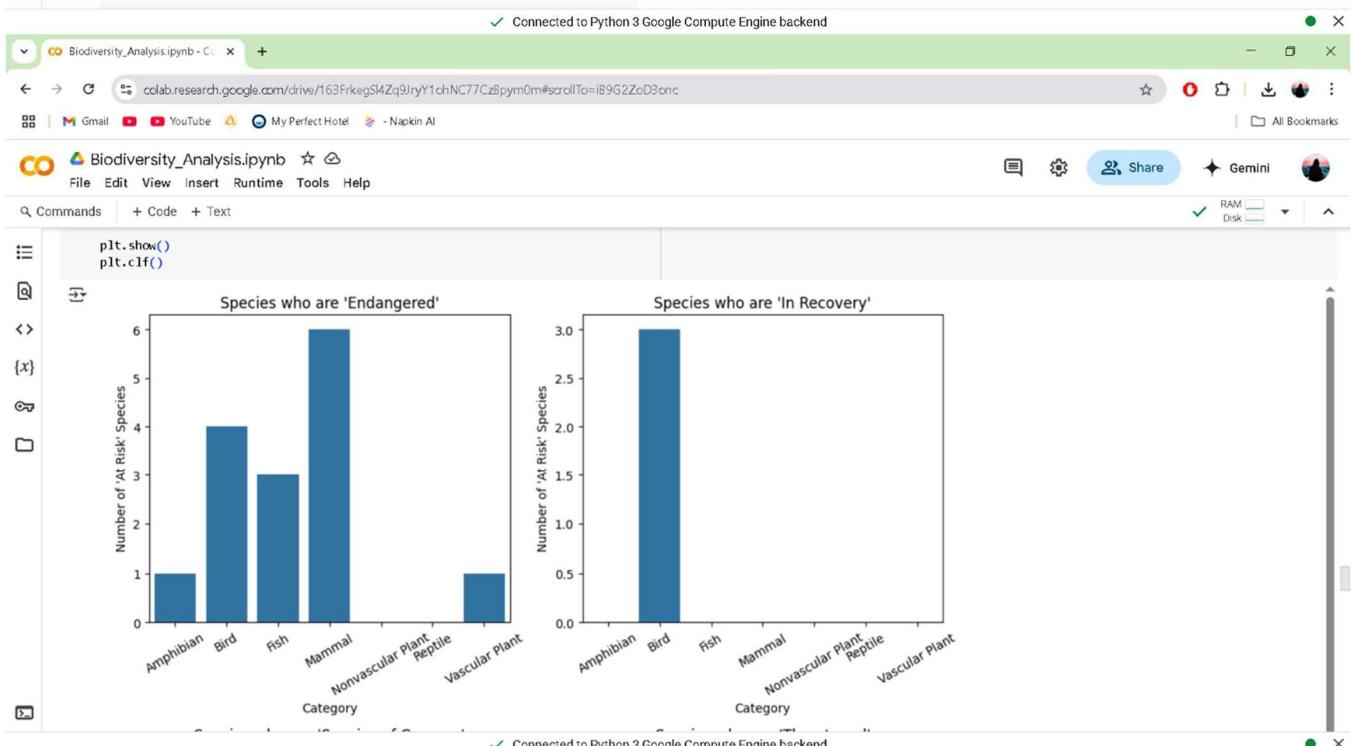
Biodiversity\_Analysis.ipynb - Colab

```

pk_sp_cat_nar_excluded = pk_sp_nar_excluded[['scientific_name', 'park_name', 'category', 'conservation_status']]
pk_sp_cat_nar_excluded = pk_sp_cat_nar_excluded.drop_duplicates(subset=['scientific_name'])
pk_sp_cat_nar_excluded = pk_sp_cat_nar_excluded.groupby(['category', 'conservation_status']).scientific_name.count().reset_index()
pk_sp_cat_nar_excluded_pivot = pk_sp_cat_nar_excluded.pivot(
    index='category',
    columns='conservation_status',
    values='scientific_name'
)
pk_sp_cat_nar_excluded_pivot.fillna(0.0, inplace=True)
pk_sp_cat_nar_excluded_pivot

```

	conservation_status	Endangered	In Recovery	Species of Concern	Threatened
category					
Amphibian	1.0	0.0	4.0	2.0	
Bird	4.0	3.0	68.0	0.0	
Fish	3.0	0.0	4.0	3.0	
Mammal	6.0	0.0	22.0	2.0	



```
# 40 (Cont.)

# Plotting the number of species per conservation status within each category
plt.figure(figsize=(10, 10))

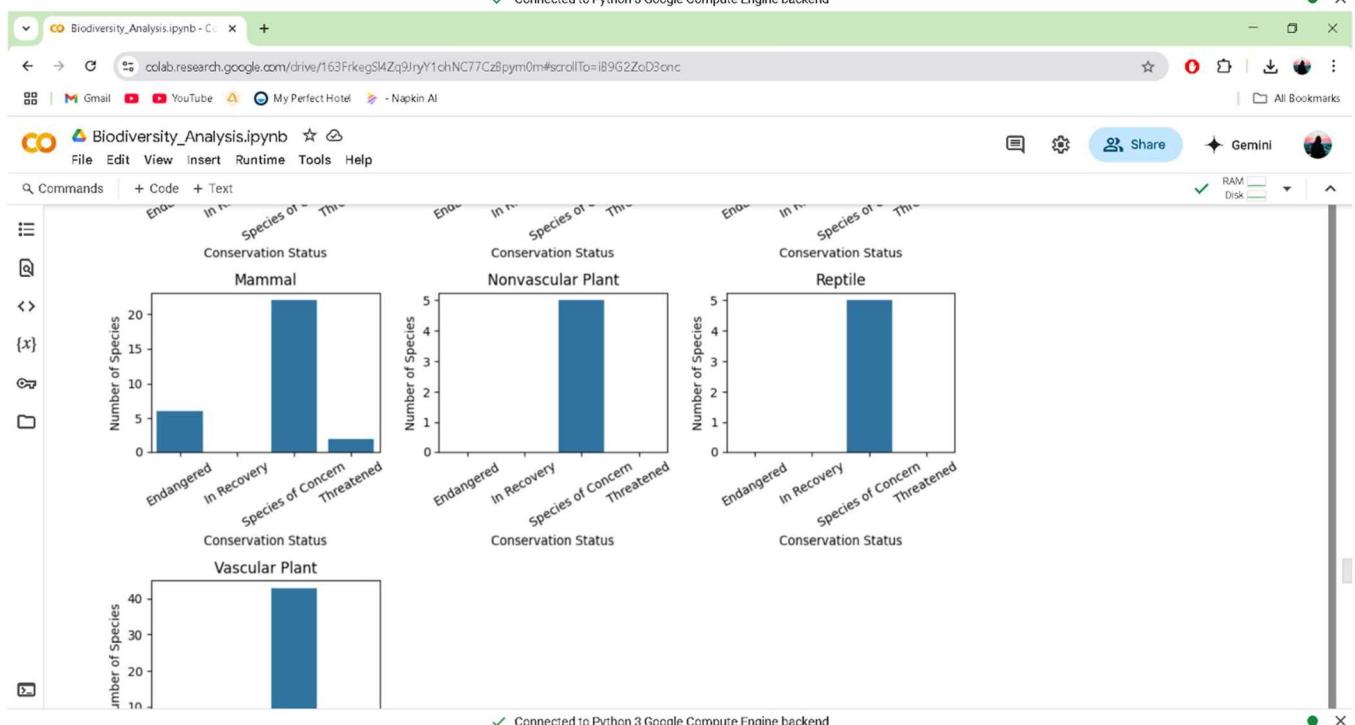
for i in range(len(pk_sp_cat_nar_excluded_pivot.index)):
    data = pk_sp_cat_nar_excluded_pivot.iloc[i].reset_index()

    plt.subplot(3, 3, i + 1)
    sns.barplot(x=data[data.columns[0]], y=data[data.columns[1]])
    plt.title(f'{pk_sp_cat_nar_excluded_pivot.index[i]}')
    plt.xlabel('Conservation Status')
    plt.ylabel('Number of Species')
    ax = plt.subplot(3, 3, i + 1)
    ax.set_xticks(range(len(data.conervation_status)))
    ax.set_xticklabels(data.conervation_status, rotation=30)
    plt.tight_layout()

plt.show()
plt.clf()
```

The code generates three bar charts side-by-side. Each chart has 'Conservation Status' on the x-axis (Endangered, In Recovery, Species of Concern, Threatened) and 'Number of Species' on the y-axis.

- Amphibian:** Number of Species: ~45, ~45, ~10, ~5
- Bird:** Number of Species: ~60, ~45, ~5, ~5
- Fish:** Number of Species: ~45, ~45, ~5, ~5



# 40 (Cont.)

```
pk_sp_cat_nar_excluded_pivot_props = pk_sp_cat_nar_excluded_pivot.copy()

for cat in pk_sp_cat_nar_excluded_pivot_props.index:

    cat_total = pk_sp_cat_nar_excluded_pivot_props.loc[cat].sum(axis=0)
    print(f'Total {cat} Species: {cat_total}')

    pk_sp_cat_nar_excluded_pivot_props.loc[cat] = pk_sp_cat_nar_excluded_pivot_props.loc[cat].apply(lambda x: str(round((x / cat_total) * 100, 2)) + '%')

pk_sp_cat_nar_excluded_pivot_props
```

Total Amphibian Species: 7.0  
Total Bird Species: 75.0  
Total Fish Species: 10.0  
Total Mammal Species: 30.0  
Total Nonvascular Plant Species: 5.0  
Total Reptile Species: 5.0  
Total Vascular Plant Species: 46.0

<ipython-input-72-9e9ed62e09bb:10> FutureWarning: Setting an item of incompatible dtype is deprecated and will raise an error in a future version of pandas. Value '14.29%' has dt\_pk\_sp\_cat\_nar\_excluded\_pivot\_props.loc[cat] = pk\_sp\_cat\_nar\_excluded\_pivot\_props.loc[cat].apply(lambda x: str(round((x / cat\_total) \* 100, 2)) + '%')

<ipython-input-72-9e9ed62e09bb:10> FutureWarning: Setting an item of incompatible dtype is deprecated and will raise an error in a future version of pandas. Value '0.0%' has dt\_pk\_sp\_cat\_nar\_excluded\_pivot\_props.loc[cat] = pk\_sp\_cat\_nar\_excluded\_pivot\_props.loc[cat].apply(lambda x: str(round((x / cat\_total) \* 100, 2)) + '%')

<ipython-input-72-9e9ed62e09bb:10> FutureWarning: Setting an item of incompatible dtype is deprecated and will raise an error in a future version of pandas. Value '57.14%' has dt\_pk\_sp\_cat\_nar\_excluded\_pivot\_props.loc[cat] = pk\_sp\_cat\_nar\_excluded\_pivot\_props.loc[cat].apply(lambda x: str(round((x / cat\_total) \* 100, 2)) + '%')

<ipython-input-72-9e9ed62e09bb:10> FutureWarning: Setting an item of incompatible dtype is deprecated and will raise an error in a future version of pandas. Value '28.57%' has dt\_pk\_sp\_cat\_nar\_excluded\_pivot\_props.loc[cat] = pk\_sp\_cat\_nar\_excluded\_pivot\_props.loc[cat].apply(lambda x: str(round((x / cat\_total) \* 100, 2)) + '%')

conservation status Endangered In Recovery Species of Concern Threatened



Biodiversity\_Analysis.ipynb - Colab

[colab.research.google.com/drive/163FrkegS4Zq9JryY1chNC77CzBpym0m#scrollTo=iB9G2ZoD3onc](https://colab.research.google.com/drive/163FrkegS4Zq9JryY1chNC77CzBpym0m#scrollTo=iB9G2ZoD3onc)

Gmail YouTube My Perfect Hotel Napkin AI All Bookmarks

Biodiversity\_Analysis.ipynb File Edit View Insert Runtime Tools Help

Commands + Code + Text

```
# 40 (cont.)
# Plotting the proportions of individual categories conservation statuses
plt.figure(figsize=(10, 10))
for i in range(len(pk_sp_cat_nar_excluded_pivot_props.index)):
    data = pk_sp_cat_nar_excluded_pivot_props.iloc[i].reset_index()
    plt.subplot(3, 3, i + 1)
    sns.barplot(x=data[data.columns[0]], y=data[data.columns[1]])
    plt.title(f'{pk_sp_cat_nar_excluded_pivot_props.index[i]}')
    plt.xlabel('Conservation Status')
    plt.ylabel('Percent of Species')
    ax = plt.subplot(3, 3, i + 1)
    ax.set_xticks(range(len(data.conervation_status)))
    ax.set_xticklabels(data.conervation_status, rotation=30)
    ax.set_yticks(range(0, 110, 10))
    ax.set_yticklabels(range(0, 110, 10))
    plt.tight_layout()
plt.show()
plt.clf()
```

Connected to Python 3 Google Compute Engine backend

Biodiversity\_Analysis.ipynb - Colab

[colab.research.google.com/drive/163FrkegS4Zq9JryY1chNC77CzBpym0m#scrollTo=iB9G2ZoD3onc](https://colab.research.google.com/drive/163FrkegS4Zq9JryY1chNC77CzBpym0m#scrollTo=iB9G2ZoD3onc)

Gmail YouTube My Perfect Hotel Napkin AI All Bookmarks

Biodiversity\_Analysis.ipynb File Edit View Insert Runtime Tools Help

Commands + Code + Text

## Section 3 - Exploring the Data - Part 4

### Categories, Parks, and Observations

The next section of the analysis will look into the category observations and see how they differ from park to park, if they differ at all.

For this section, you won't be using the dataframe that excluded the 'Not at Risk' conservation\_status value since you won't be using that column of the data.

41. Make sure you have the correct dataframe variable to use
42. Manipulate it to be workable for this problem and create a pivot table

```
[76] # 41
parks_species.head()
```

	scientific_name	park_name	observations	category	common_names	conservation_status
0	Vicia benghalensis	Great Smoky Mountains National Park	68	Vascular Plant	Purple Vetch, Reddish Tufted Vetch	Not at Risk
1	Neovison vison	Great Smoky Mountains National Park	77	Mammal	American Mink	Not at Risk
2	Prunus subcordata	Yosemite National Park	138	Vascular Plant	Klamath Plum	Not at Risk
3	Abutilon theophrasti	Bryce National Park	84	Vascular Plant	Velvelleaf	Not at Risk
4	Githopsis specularioides	Great Smoky Mountains National Park	85	Vascular Plant	Common Bluecup	Not at Risk

Connected to Python 3 Google Compute Engine backend

Biodiversity\_Analysis.ipynb - Colab

<https://colab.research.google.com/drive/163FrkegS4Zq9JryY1chNC77Czbpym0m#scrollTo=iB9G2ZoD3cnc>

Gmail YouTube My Perfect Hotel Napkin AI All Bookmarks

**Biodiversity\_Analysis.ipynb** Share Gemini RAM Disk

File Edit View Insert Runtime Tools Help

Commands + Code + Text

```
# 3 Abutilon theophrasti Bryce National Park 84 Vascular Plant Velvetleaf Not at Risk
# 4 Githopsis specularioides Great Smoky Mountains National Park 85 Vascular Plant Common Bluecup Not at Risk
```

# 42

```
parks_cats = parks_species[['category', 'park_name', 'observations']]
parks_cats = parks_cats.groupby(['category', 'park_name']).observations.sum().reset_index()
parks_cats_pivot = parks_cats.pivot(
    index='park_name',
    columns='category',
    values='observations')
parks_cats_pivot
```

park_name	category	Amphibian	Bird	Fish	Mammal	Nonvascular Plant	Reptile	Vascular Plant
Bryce National Park		7299	48383	12223	16823	32992	7864	422686
Great Smoky Mountains National Park		5622	35290	9068	12301	24857	5616	318071
Yellowstone National Park		19191	119219	30131	41905	83021	19300	1060769
Yosemite National Park		11309	71290	18353	24887	49783	11335	634515

Connected to Python 3 Google Compute Engine backend

Biodiversity\_Analysis.ipynb - Colab

<https://colab.research.google.com/drive/163FrkegS4Zq9JryY1chNC77Czbpym0m#scrollTo=iB9G2ZoD3cnc>

Gmail YouTube My Perfect Hotel Napkin AI All Bookmarks

**Biodiversity\_Analysis.ipynb** Share Gemini RAM Disk

File Edit View Insert Runtime Tools Help

Commands + Code + Text

```
parks_cats_pivot_props = parks_cats_pivot.copy()
for park in parks_cats_pivot_props.index:
    park_total = parks_cats_pivot_props.loc[park].sum(axis=0)
    print(f'Total {park} Observations: {park_total}')
    parks_cats_pivot_props.loc[park] = parks_cats_pivot_props.loc[park].apply(lambda x: str(round((x / park_total) * 100, 2)) + '%')
parks_cats_pivot_props
```

Total Bryce National Park Observations: 548159  
Total Great Smoky Mountains National Park Observations: 410825  
Total Yellowstone National Park Observations: 1373536  
Total Yosemite National Park Observations: 821472  
<ipython-input-78-701992344311>:10: FutureWarning: Setting an item of incompatible dtype is deprecated and will raise an error in a future version of pandas. Value '1.33%' has c  
 parks\_cats\_pivot\_props.loc[park] = parks\_cats\_pivot\_props.loc[park].apply(lambda x: str(round((x / park\_total) \* 100, 2)) + '%')  
<ipython-input-78-701992344311>:10: FutureWarning: Setting an item of incompatible dtype is deprecated and will raise an error in a future version of pandas. Value '8.83%' has c  
 parks\_cats\_pivot\_props.loc[park] = parks\_cats\_pivot\_props.loc[park].apply(lambda x: str(round((x / park\_total) \* 100, 2)) + '%')  
<ipython-input-78-701992344311>:10: FutureWarning: Setting an item of incompatible dtype is deprecated and will raise an error in a future version of pandas. Value '2.23%' has c  
 parks\_cats\_pivot\_props.loc[park] = parks\_cats\_pivot\_props.loc[park].apply(lambda x: str(round((x / park\_total) \* 100, 2)) + '%')  
<ipython-input-78-701992344311>:10: FutureWarning: Setting an item of incompatible dtype is deprecated and will raise an error in a future version of pandas. Value '3.07%' has c  
 parks\_cats\_pivot\_props.loc[park] = parks\_cats\_pivot\_props.loc[park].apply(lambda x: str(round((x / park\_total) \* 100, 2)) + '%')  
<ipython-input-78-701992344311>:10: FutureWarning: Setting an item of incompatible dtype is deprecated and will raise an error in a future version of pandas. Value '6.02%' has c  
 parks\_cats\_pivot\_props.loc[park] = parks\_cats\_pivot\_props.loc[park].apply(lambda x: str(round((x / park\_total) \* 100, 2)) + '%')  
<ipython-input-78-701992344311>:10: FutureWarning: Setting an item of incompatible dtype is deprecated and will raise an error in a future version of pandas. Value '1.43%' has c  
 parks\_cats\_pivot\_props.loc[park] = parks\_cats\_pivot\_props.loc[park].apply(lambda x: str(round((x / park\_total) \* 100, 2)) + '%')  
<ipython-input-78-701992344311>:10: FutureWarning: Setting an item of incompatible dtype is deprecated and will raise an error in a future version of pandas. Value '77.09%' has c  
 parks\_cats\_pivot\_props.loc[park] = parks\_cats\_pivot\_props.loc[park].apply(lambda x: str(round((x / park\_total) \* 100, 2)) + '%')

Connected to Python 3 Google Compute Engine backend

Well, once again, the proportions show that each park is nearly identical in percentage of observations for each category.

For fun, make one last figure of graphs to visualize or pivot table findings.

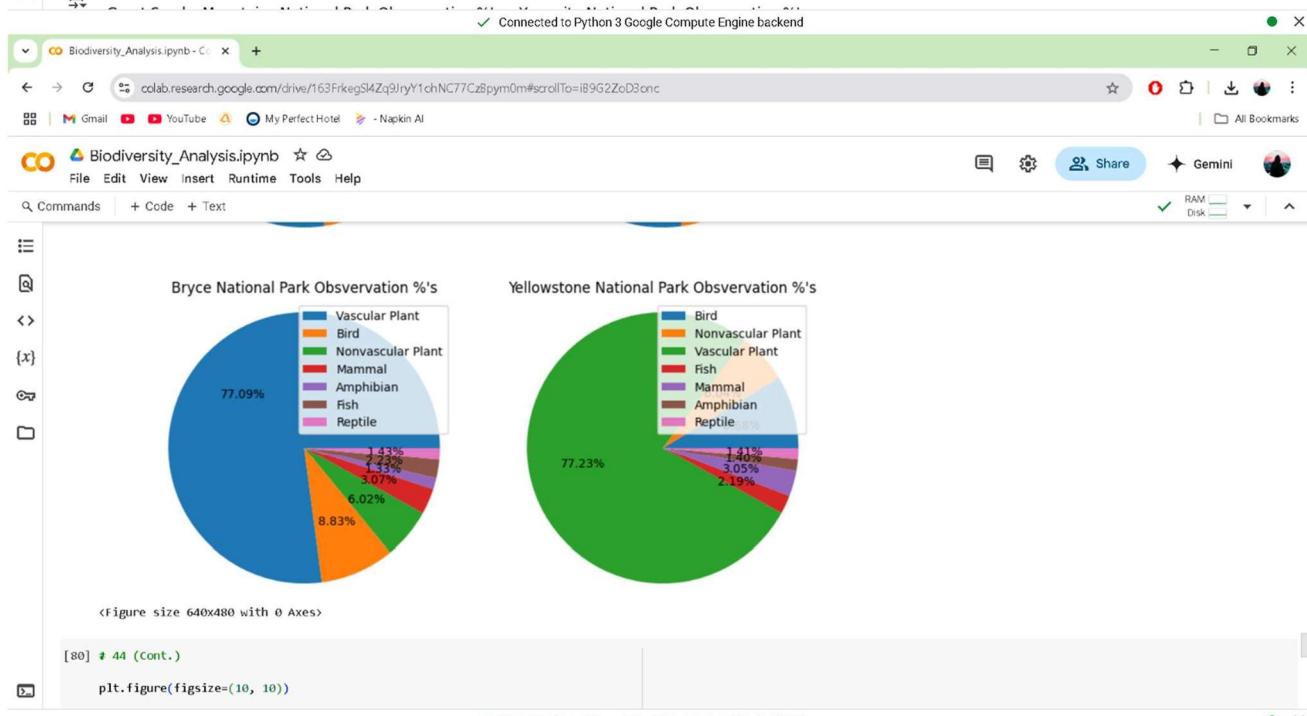
44. Create a visualization to show the observations of each category in each park

```
[79] # 44
plt.figure(figsize=(10, 10))

for i in range(len(parks_species.park_name.unique())):
    cats_parks_df = parks_species[parks_species.park_name == parks_species.park_name.unique()[i]]
    cats_parks_df = cats_parks_df[['category', 'observations']]
    cats_parks_df = cats_parks_df.groupby('category', sort=False).observations.sum()
    labels = cats_parks_df.index

    plt.subplot(2, 2, i + 1)
    plt.pie(cats_parks_df, autopct='%.2f%%')
    plt.axis('equal')
    plt.title(f'{parks_species.park_name.unique()[i]} Observation %'s')
    plt.legend(labels, loc='upper right')

plt.show()
plt.clf()
```



Biodiversity\_Analysis.ipynb - Colab

File Edit View Insert Runtime Tools Help

Commands + Code + Text

```
<Figure size 640x480 with 0 Axes>

# 44 (cont.)

plt.figure(figsize=(10, 10))

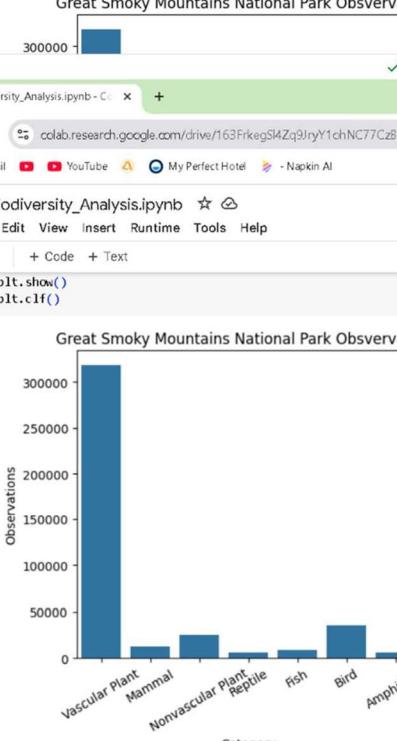
for i in range(len(parks_species.park_name.unique())):

    cats_parks_df = parks_species[parks_species.park_name == parks_species.park_name.unique()[i]]
    cats_parks_df = cats_parks_df[['category', 'observations']]
    cats_parks_df = cats_parks_df.groupby('category', sort=False).observations.sum()

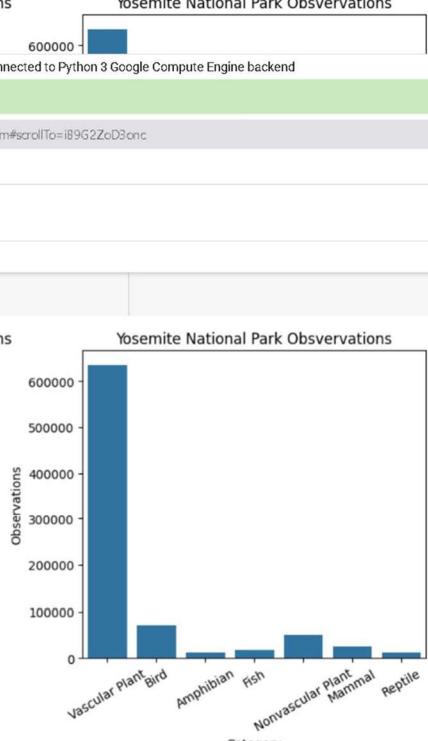
    plt.subplot(2, 2, i + 1)
    sns.barplot(x=cats_parks_df.index, y=cats_parks_df['observations'])
    plt.title(f'{parks_species.park_name.unique()[i]} Observations')
    plt.xlabel('Category')
    plt.ylabel('Observations')
    ax = plt.subplot(2, 2, i + 1)
    ax.set_xticks(range(len(cats_parks_df.index)))
    ax.set_xticklabels(cats_parks_df.index, rotation=30)
    plt.tight_layout()

plt.show()
plt.clf()
```

Great Smoky Mountains National Park Observations

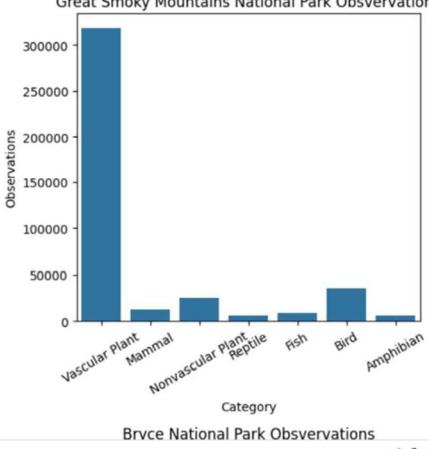


Yosemite National Park Observations

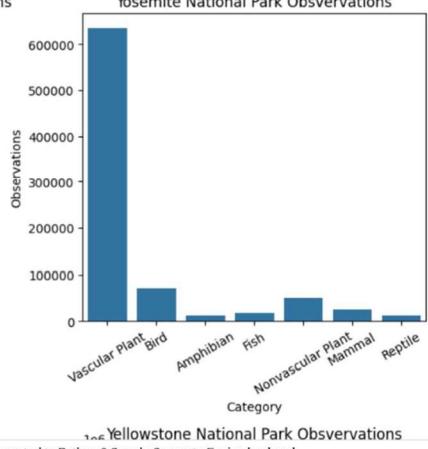


Connected to Python 3 Google Compute Engine backend

Great Smoky Mountains National Park Observations



Yosemite National Park Observations



Brvce National Park Observations

Yellowstone National Park Observations

Connected to Python 3 Google Compute Engine backend

Biodiversity\_Analysis.ipynb - Colab

File Edit View Insert Runtime Tools Help

Commands + Code + Text

## Section 4 - Statistical Significance Among Endangered Species

For the final part of the analysis, you'll refer back to some of the pivot tables from earlier and you'll test for statistical significance between categories of endangered species.

45. Take a look at the dataframe you'll be using for the statistical analysis

```
[81]: # 45
parks_species_cat_pivot
```

category	conservation_status	Endangered	In Recovery	Not at Risk	Species of Concern	Threatened
Amphibian	1.0	0.0	72.0	4.0	2.0	
Bird	4.0	3.0	413.0	68.0	0.0	
Fish	3.0	0.0	115.0	4.0	3.0	
Mammal	6.0	0.0	146.0	22.0	2.0	
Nonvascular Plant	0.0	0.0	328.0	5.0	0.0	
Reptile	0.0	0.0	73.0	5.0	0.0	

✓ Connected to Python 3 Google Compute Engine backend

46. Add a column to parks\_species\_cat\_pivot that has the total species for each category that is not endangered

```
# 46
parks_species_cat_pivot['Total (excluding Endangered)'] = parks_species_cat_pivot['In Recovery'] + \
parks_species_cat_pivot['Not at Risk'] + \
parks_species_cat_pivot['Species of Concern'] + \
parks_species_cat_pivot['Threatened']
```

category	conservation_status	Endangered	In Recovery	Not at Risk	Species of Concern	Threatened	Total (excluding Endangered)
Amphibian	1.0	0.0	72.0	4.0	2.0		78.0
Bird	4.0	3.0	413.0	68.0	0.0		484.0
Fish	3.0	0.0	115.0	4.0	3.0		122.0
Mammal	6.0	0.0	146.0	22.0	2.0		170.0
Nonvascular Plant	0.0	0.0	328.0	5.0	0.0		333.0
Reptile	0.0	0.0	73.0	5.0	0.0		78.0
Vascular Plant	1.0	0.0	4216.0	43.0	2.0		4261.0

Now, with the number of endangered species and the number of total species, you can run a chi-squared test by using the `chi2` statistic.

✓ Connected to Python 3 Google Compute Engine backend

Biodiversity\_Analysis.ipynb - Colab

File Edit View Insert Runtime Tools Help

# 48

```
for cat in parks_species_cat_pivot.index:
    if cat == 'Mammal':
        continue

    contingency = [[6, 170],
                   [parks_species_cat_pivot.loc[cat].Endangered, parks_species_cat_pivot.loc[cat]['Total (excluding Endangered)']]]
    p_val = ss.chisq_contingency(contingency)[1]
    print(f'{cat} P-value: {round(p_val, 4)}')
    if p_val < 0.05:
        print(f'The difference in the rate of endangerment between Mammals and {cat}(s) is statistically significant\n')
    else:
        print(f'The difference in the rate of endangerment between Mammals and {cat}(s) is NOT statistically significant\n')

Amphibian P-Value: 0.5795
The difference in the rate of endangerment between Mammals and Amphibian(s) is NOT statistically significant

Bird P-Value: 0.0397
The difference in the rate of endangerment between Mammals and Bird(s) is statistically significant

Fish P-Value: 0.8704
The difference in the rate of endangerment between Mammals and Fish(s) is NOT statistically significant

Nonvascular Plant P-Value: 0.0031
The difference in the rate of endangerment between Mammals and Nonvascular Plant(s) is statistically significant

Reptile P-Value: 0.2292
The difference in the rate of endangerment between Mammals and Reptile(s) is NOT statistically significant
```

Connected to Python 3 Google Compute Engine backend

The difference in the rate of endangerment between Mammals and Vascular Plant(s) is statistically significant

After running the chi-squared tests, we have some results, but what do these results mean?

You have three categories that have a statistically significant difference in endangerment rates when compared to Mammals:

- Birds
- Nonvascular Plants
- Vascular Plants

This means that, when you compare the endangerment rates of mammals to either of these three categories of species, mammals show a significantly higher rate of endangerment than the other three categories.

This is determined with a 5% threshold (lower than 0.05 p-value). This means that less than 5% percent of the time, the endangerment rate of mammals among a random sample population will NOT have a higher rate of endangerment than the other three categories of a random sample population.

Inversely, more than 95% of the time, the endangerment rate amongst mammals will be higher than the endangerment rate of the other three categories.

---

Section 5 - Conclusions

Analysis Findings

Connected to Python 3 Google Compute Engine backend

Biodiversity\_Analysis.ipynb - Colab

File Edit View Insert Runtime Tools Help

Analysis Findings -

*General analysis Findings:*

- Though there is a huge difference in observation counts for each park, there are near equal proportions of observations for each conservation status within each park (Ex: The endangered conservation status has nearly the same percent of observations in each park)
- Though there is a huge difference in observation counts for each park, there are near equal proportions of observations for each category of species within each park (Ex: The mammal category has nearly the same percent of observations in each park)
- On average, as the severity of the conservation status increases for a specific species, the sightings of that specific species decreases (Approaches Endangerment)
- On average, the observation count of a species can be a predictor of conservation status (higher count = lower risk status | lower count = higher risk status)

*Findings among 'At Risk' and 'Not at Risk' species:*

- Both Plant categories have the lowest proportion of species that are 'At risk'
- Birds have the highest proportion of species that are labeled a 'Species of Concern'
- Reptiles and Nonvascular Plants have no endangered species within their categories
- Mammals have the highest proportion of species that are at 'At Risk' (Lowest proportion of species that are 'Not at Risk')

✓ Connected to Python 3 Google Compute Engine backend

Biodiversity\_Analysis.ipynb - Colab

File Edit View Insert Runtime Tools Help

*Findings among 'At Risk' species:*

- Birds have the highest proportion of species that are labeled a 'Species of Concern'
- Reptiles and Nonvascular Plants have no endangered species within their categories
- Mammals have the highest proportion of species that are at 'At Risk' (Lowest proportion of species that are 'Not at Risk')
- Mammals have the highest proportion of species that are endangered

*Findings among 'At Risk' species:*

- Fish have the highest proportion of species that are endangered
- Species that live primarily in the water (Amphibians and Fish) have higher proportions of endangerment (Amphibians are 3rd highest behind Mammals) and are more threatened than any other category
- Birds, Reptiles, Vascular Plants and Nonvascular Plants each have a very high proportion of species labeled as a 'Species of Concern', and a low proportion of species that are endangered, relative to the other categories (This finding could be used for preventative measures)

*Findings among 'Not at Risk' species:*

- Vascular plants take up roughly 77% of all 'Not at Risk' species (4216 out of 5363 species)

*Statistical findings:*

- Mammals have statistically significant higher rates of endangerment when compared to Birds, Nonvascular Plants, and Vascular Plants (More than 95% of the time)

✓ Connected to Python 3 Google Compute Engine backend

### c) Process Involved

#### **Project Conception and Justification**

The WildStat project was conceived in response to the increasing global need for biodiversity conservation and ecological awareness. With accelerating habitat loss, climate change, and human encroachment, it's crucial to monitor and analyse biodiversity patterns in protected ecosystems such as U.S. National Parks.

WildStat aims to provide an intuitive data analysis platform that leverages real-world species data, visualizations, and statistical insights to highlight biodiversity trends, endangered species distribution, and park-wise ecological summaries. The project is implemented using Python with libraries such as pandas, NumPy, Matplotlib, and Seaborn within a Jupyter Notebook environment.

The core goal is to empower students, researchers, and environmental policymakers with meaningful, data-driven insights that can support conservation planning and ecological studies.

#### **The need for WildStat is justified by several key factors:**

- Biodiversity Loss and Conservation Needs: With numerous species facing extinction, there's an urgent requirement to study patterns and take preventive action.
- Lack of Publicly Engaging Tools: There are few accessible tools that allow non-specialists to analyze national park species data in an understandable way.
- Educational and Research Value: WildStat helps students and enthusiasts understand how data science can be applied to real-world environmental challenges.
- Data-Driven Environmental Planning: Through visualizations and summaries, WildStat enables better comprehension of ecological richness and threats.
- Scalability: The platform can be scaled to include more datasets (e.g., pollution levels, climate change impact, migratory patterns) across global protected areas.

#### **Requirements Gathering and Analysis**

##### **User Needs Assessment**

**Identifying Target Users:** Primary users include ecology students, data science learners, environmental researchers, and educators.

##### **User Goals:**

- Analyze biodiversity metrics for different U.S. National Parks.
- Identify the number and types of endangered species in each park.
- Visualize species distribution across categories (e.g., mammals, birds, plants). Understand data trends using statistical and graphical methods.
- Access an educational tool that bridges ecology with data science.

##### **Functional Requirements**

- Data cleaning and wrangling from CSV files.
- Summary statistics (mean, count, distribution) for species per park.
- Category-wise filtering (e.g., endangered, threatened).
- Graphical visualizations (bar plots, pie charts, line graphs).
- Correlation and comparison of park-wise biodiversity.
- Exporting or saving insights for reports or presentations.

##### **Non-Functional Requirements**

- Usability: Clean, readable Jupyter Notebook structure with markdown explanations.
- Performance: Efficient processing of moderately sized CSV datasets.
- Accessibility: The notebook can be shared or run on platforms like Google Colab.

- Reusability: Modular code and functions to allow updates with new data.
- Maintainability: Well-commented code and structured sections for easy future enhancements.

## System Analysis and Design

### Architecture Overview

WildStat follows a simple but effective architecture:

- Data Tier: CSV dataset containing records of species, categories, conservation status, and park locations.
- Application Tier: Jupyter Notebook powered by Python scripts and visualization libraries.
- Presentation Tier: Outputs are shown as readable tables, graphs, and markdown text within the notebook itself.

### Component Design

- Data Loading Module: Loads and reads the biodiversity CSV dataset.
- Data Cleaning & Preprocessing Module: Handles missing values, standardizes entries, and structures the data.
- Analysis Module:
  - Computes statistics like species count, endangered species count per park, etc.
  - Identifies biodiversity hotspots.
- Visualization Module:
  - Bar plots for park-wise species counts.
  - Pie charts for conservation statuses.
  - Line graphs to represent comparative trends.
- User Interaction (via code cells): Users can run and modify code blocks to analyze specific cases or customize plots.

### Data Flow

- Dataset is loaded from CSV.
- Data is cleaned and structured.
- Users run analysis functions to generate insights.
- Visualizations are rendered using Matplotlib and Seaborn.
- Insights are presented alongside markdown interpretations.

### UI and Output Design

- Clear code cell structure with headings.
- Inline comments for understanding logic.
- Graphs and charts embedded in the notebook.
- Text-based summaries after visualizations for interpretation.

## **Project Planning and Management:**

The **WildStat** project was strategically planned with a focus on delivering an insightful and educational data science tool for biodiversity analysis in U.S. National Parks. The project scope involved the use of Python and Jupyter Notebooks, leveraging libraries such as pandas, NumPy, Matplotlib, and Seaborn for data handling, analysis, and visualization. The phases of development included problem identification, data collection, preprocessing, exploratory data analysis, visualization, interpretation of results, and documentation.

### **Key milestones were mapped to:**

- Dataset exploration and cleaning,
- Implementation of biodiversity metrics and filtering logic,
- Generation of insightful graphs,
- Interpretation of results for conservation relevance,
- Finalization of the report and project notebook.

Resources were allocated in alignment with the academic goals of the project, focusing on data science skills, research capabilities, and technical tools required for effective analysis.

Project management emphasized clarity, consistency, and adaptability. Risks such as data inconsistencies, misinterpretation of biological terms, or visualization complexity were identified early and mitigated through iterative testing and validation. The use of a structured Jupyter Notebook format ensured transparency and ease of debugging. Regular checkpoints were set to monitor progress against deliverables, and scope adjustments were made when necessary to maintain focus on educational impact and usability.

The planning followed an iterative, data-driven workflow that allowed for continuous refinement and learning, ensuring that the final output met both technical expectations and academic standards.

## **System Development and Implementation**

The implementation of the Biodiversity Analysis project involved developing a structured and functional data science pipeline using Python within the Jupyter Notebook environment. The focus was on effective data ingestion, cleaning, exploratory analysis, visualization, and basic interpretation. Key components of the system implementation are as follows:

### **Data Acquisition and Loading:**

The dataset used, sourced from Kaggle, includes information on species presence, conservation status, and park location across U.S. National Parks. It was loaded and managed using `pandas`, allowing for efficient manipulation and inspection of large tabular data.

### **Data Preprocessing and Cleaning:**

Initial preprocessing involved handling missing values, correcting inconsistent labels, and renaming columns for clarity. Functions were implemented to normalize textual data and convert it into formats suitable for grouped analysis.

### **Exploratory Data Analysis (EDA):**

Using `pandas`, `matplotlib`, and `seaborn`, various visualizations were generated to explore species distribution by category, conservation status, and park locations. These included bar plots, pie charts, and grouped bar graphs to represent biodiversity insights visually.

### **Custom Analysis:**

Specific queries and filters were implemented to identify the number of endangered species by park, observe the frequency of different species types, and analyze trends in protected vs. non-protected species. These were coded using `groupby` operations and condition-based filtering.

### **Visualization and Interpretation:**

Visual insights formed a core part of the analysis, helping interpret conservation challenges in national parks. Custom styling and annotations were added to ensure clarity in data storytelling. Seaborn's advanced features were used for aesthetic and informative plots.

### **Notebook Documentation and Structure:**

The notebook is divided into clearly labeled sections—data loading, cleaning, visualization, and conclusion—making it modular and readable. Markdown cells are used to explain the logic, insights, and significance of each step.

In the **Biodiversity Analysis of U.S. National Parks** project, testing and validation were conducted within the data analysis environment to ensure the accuracy of results and the reliability of insights derived from the dataset. Given that the project was implemented as a Jupyter Notebook for data analytics and visualization, the quality assurance process focused on data correctness, code reliability, and visual validation.

#### **Manual Code Validation:**

Each cell in the Jupyter Notebook was executed sequentially and validated to ensure that the code performed as expected without runtime errors. This iterative execution allowed for step-by-step verification of data transformations, calculations, and visualizations.

#### **Data Consistency Checks:**

Data quality was verified through checks for missing values, duplicate entries, and inconsistent labels. Functions were written to clean and normalize data, and outputs were visually inspected using summary statistics and `.info()` / `.describe()` methods in `pandas`.

#### **Visual Output Validation:**

All generated plots and visualizations (using `matplotlib` and `seaborn`) were evaluated manually to ensure they accurately represented the underlying data. Graphs were checked for correctness of axes, legends, labels, and overall readability.

#### **Exploratory Sanity Checks:**

Aggregated counts (e.g., total number of endangered species or species per park) were cross-verified using multiple methods to confirm consistency. These cross-checks helped prevent logical errors in groupings or filters.

#### **Iterative Refinement:**

Based on visual inspection and logical interpretation, code was modified and optimized for clarity and correctness. Insights drawn from the data were also evaluated to ensure they aligned with expectations based on the dataset's structure.

Though formal testing frameworks weren't used due to the nature of the project, the code was written with clarity and correctness in mind, with a focus on producing accurate and meaningful analytical results.

#### **Comprehensive Performance Monitoring:**

Continuous tracking of key application performance indicators (KPIs) such as response times, CPU utilization, memory usage, and network latency. Proactive identification of performance bottlenecks, slowdowns, or areas of inefficiency within the system. Implementation of optimization strategies, including code refactoring, database tuning, and infrastructure adjustments, to enhance application speed and responsiveness.

#### **Diligent Security Updates:**

Regular monitoring for security patches and updates released by the underlying operating systems, frameworks (Next.js, Node.js, Flask), and libraries used in the application. Prompt application of these security updates to mitigate potential vulnerabilities and protect the application and user data from known threats. Periodic security audits and vulnerability scanning to proactively identify and address potential weaknesses in the system.

### **Effective Technical Support:**

Providing timely and helpful assistance to users who encounter problems while using the application or have questions about its features and functionality. Troubleshooting user-reported issues through various channels (e.g., email, FAQs, potentially in-app support). Offering clear guidance and instructions to users on how to effectively utilize the application's features and resolve common issues.

### **Consistent System Administration:**

Regular maintenance of the server infrastructure hosting the application, including software updates, configuration management, and resource allocation. Implementation of robust database backup strategies and regular execution of backups to prevent data loss. Effective management of application logs for monitoring system behavior, diagnosing issues, and identifying potential security threats.

## **System Evolution and Enhancement:**

### **Detailed Requirements Gathering:**

Systematic collection of feedback from users through surveys, interviews, feedback forms, and analysis of usage patterns. Engagement with stakeholders to understand evolving business needs, new opportunities, and potential integrations. Thorough analysis of gathered feedback to identify valuable new features, desired improvements to existing functionality, and areas for enhanced user experience.

### **Structured Design and Development:**

Careful design of new features and enhancements, considering user needs, system architecture, and technical feasibility. Following a structured development process, including detailed design specifications, efficient coding practices, and thorough unit and integration testing. Adherence to coding standards and best practices to ensure maintainability and scalability of the codebase.

### **Managed Release Management:**

Careful planning and coordination of the deployment of updated application versions containing new features or enhancements. Thorough testing in staging environments before releasing to production to minimize the risk of introducing new issues. Clear communication of the changes and new features to users through release notes, in-app announcements, or other appropriate channels.

### **Strategic Technology Upgrades:**

Continuous evaluation of the underlying technologies used in the application to identify opportunities for improvement in performance, security, and maintainability. Planned upgrades to newer versions of frameworks (Next.js, Node.js, Flask), libraries, and the database system, ensuring compatibility and leveraging new features and security enhancements.

### **Thoughtful System Redesign:**

Periodic review of the application's overall architecture and design to identify fundamental limitations or areas that could benefit from significant improvements in scalability, performance, or maintainability.

Careful planning and execution of system redesign efforts, potentially involving significant modifications to the database schema, API architecture, or core components, to address evolving needs and ensure the long-term viability of the application.

## **d) Methodology used testing**

The **Biodiversity Analysis of U.S. National Parks** project was built using Python in a Jupyter Notebook environment. Testing was integrated directly into the development workflow to ensure accuracy, reliability, and consistency of the data analysis and visualizations. Although this project does not involve traditional software components like APIs or user interfaces, specific testing methodologies were applied relevant to data science projects.

### **Code Validation and Output Testing**

#### **Purpose:**

To verify that the data analysis code produces expected results and functions correctly on the given dataset.

#### **Description:**

Code blocks were tested iteratively to ensure expected outputs after each transformation (e.g., filtering endangered species, grouping by category, visualizing park-wise distribution). Pandas and NumPy operations were validated by checking data shapes, null values, and statistical summaries. Cross-verification was performed between calculated metrics (e.g., count of endangered species) and visualized plots to ensure consistency.

#### **Benefits:**

Ensures correctness of the analysis logic. Prevents misleading insights due to unnoticed coding mistakes.

### **Data Integrity Checks**

#### **Purpose:**

To validate that the dataset used in the analysis is clean, consistent, and suitable for deriving meaningful insights.

#### **Description:**

Null values were identified and handled using `isnull()` checks. Unnecessary columns were dropped, and relevant features were selected based on analysis goals. Categorical variables were reviewed for spelling inconsistencies or unexpected values.

#### **Benefits:**

Clean data ensures accurate visualizations and conclusions. Helps avoid false interpretations caused by data inconsistencies.

## **Visualization Verification**

### **Purpose:**

To confirm that plots and graphs accurately represent the underlying data and convey intended insights.

### **Description:**

All plots were manually verified by comparing them to summary statistics (e.g., bar plots vs. species counts). Color schemes and legends were checked to ensure clarity and interpretability. Axes labels and titles were reviewed to avoid confusion.

### **Benefits:**

Enhances readability and trustworthiness of the visualizations. Prevents misleading graphical representations.

## **Manual Peer Review**

### **Purpose:**

To get external feedback on the notebook's readability, correctness, and clarity.

### **Description:**

The project notebook was reviewed by peers and shared with mentors for informal validation. Feedback was implemented to improve explanations, markdown cells, and plot labeling.

### **Benefits:**

Provides fresh perspectives and helps identify overlooked issues. Ensures the project is understandable by others.

This testing approach, while lightweight compared to traditional software testing, was essential to maintain the **integrity, correctness, and educational value** of the analysis.

To ensure **WildStat** was reliable, insightful, and intuitive for end users, a **multi-layered testing methodology** was applied throughout its development. Given that WildStat assists users in analyzing biodiversity data and identifying patterns among endangered species, maintaining accuracy and system integrity was critical. The testing lifecycle encompassed everything from unit-level verifications to real-world user simulations.

## **e) Test Report**

### **Unit Testing**

#### **Purpose:**

To validate that individual functions and components performed correctly in isolation, especially those handling sensitive biodiversity data.

#### **Description:**

Key modules were tested with known inputs and expected outputs:

- **Species Data Parser:** Verified parsing and cleaning of raw biodiversity datasets.
- **ML API Client:** Tested communication logic between the backend and the Flask-based prediction service.
- **Data Visualization Logic:** Ensured chart rendering components displayed correct metrics (e.g., species count trends).
- **Authentication Handlers:** Validated secure login/logout, session tokens, and role management.

**Benefits:**

- Early bug detection.
- Higher code quality and maintainability.
- Safeguarded core features from regression after updates.

## Integration Testing

**Purpose:**

To confirm that individually verified modules worked correctly together across the full system stack.

**Description:**

Integration points were tested to ensure data flowed seamlessly between layers:

- **Frontend ↔ Backend:** Verified park/species queries and chart update triggers.
- **Backend ↔ ML API:** Ensured consistent request/response formats and prediction result handling.
- **Database ↔ Backend:** Confirmed accurate saving and retrieval of species records, park metadata, and predictions.

**Approach:**

Followed a **bottom-up** strategy — starting with lower-level services (e.g., API clients, DB handlers) and building up to full feature workflows.

**Benefits:**

- Detected interface mismatches and serialization bugs.
- Validated correct component orchestration.

## System Testing

**Purpose:**

To validate WildStat's complete functionality in a simulated production environment and test both **functional** and **non-functional** aspects.

**Description:**

Tested the full platform using actual park data and simulated usage patterns:

- **User Journey Tests:** From logging in to querying park data and receiving endangered species reports.

- **Performance Metrics:** Measured chart load times, ML prediction response time, and overall application responsiveness.
- **Security Checks:** Verified authentication integrity and secure API access.
- **UI Consistency:** Checked for layout stability and responsiveness across screen sizes.

### **System Testing Types:**

- **Functional Testing:** Feature-by-feature validation.
- **Usability Testing:** Informal reviews for navigation and clarity.
- **Performance Testing:** Ensured acceptable response times under sample data loads.
- **Security Testing:** Inspected for potential data leakage or weak authentication points.

### **Benefits:**

- Confidence in system behavior under realistic conditions.
- Verified stability and responsiveness across user flows.

## **User Acceptance Testing (UAT)**

### **Purpose:**

To validate that **WildStat** meets real-world expectations of its target users — students, conservation researchers, or wildlife analysts.

### **Description:**

A group of internal stakeholders tested the system under actual use cases:

- **Scenario Testing:** Entered real parks (e.g., Yellowstone) and reviewed predicted species at risk.
- **Feedback Gathering:** Focused on clarity of visualizations, trust in recommendations, and UI intuitiveness.
- **Defect Logging:** Reported bugs, UI glitches, or confusing workflows.

### **Acceptance Criteria:**

- Intuitive navigation from park selection to insight generation.
- Accurate and meaningful prediction results.
- Clear data visualizations with supporting metadata.

### **Benefits:**

- Aligned features with user expectations.
- Identified subtle usability or clarity issues.
- Built trust in insights generated by the platform.

## **Test Report Summary**

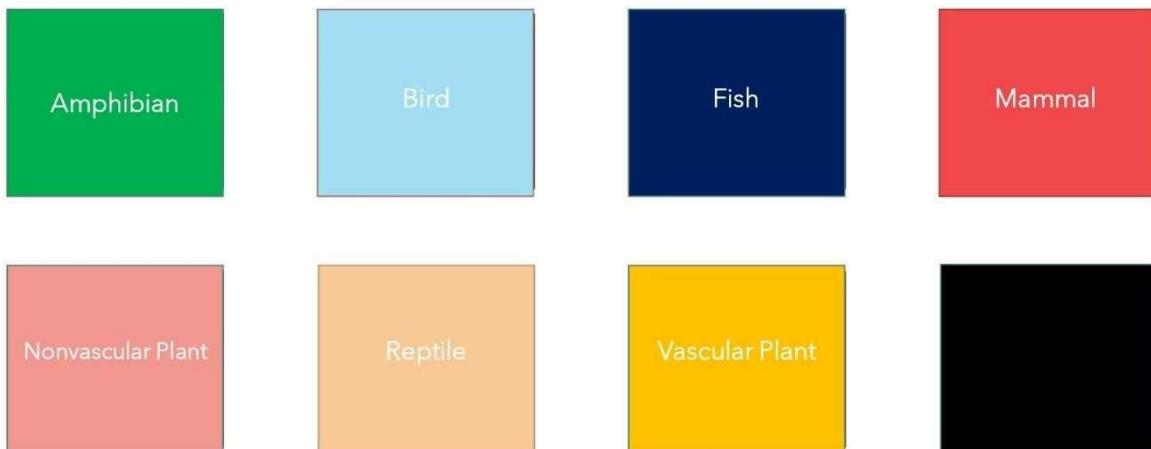
WildStat's quality assurance process followed a rigorous, multi-tiered testing approach:

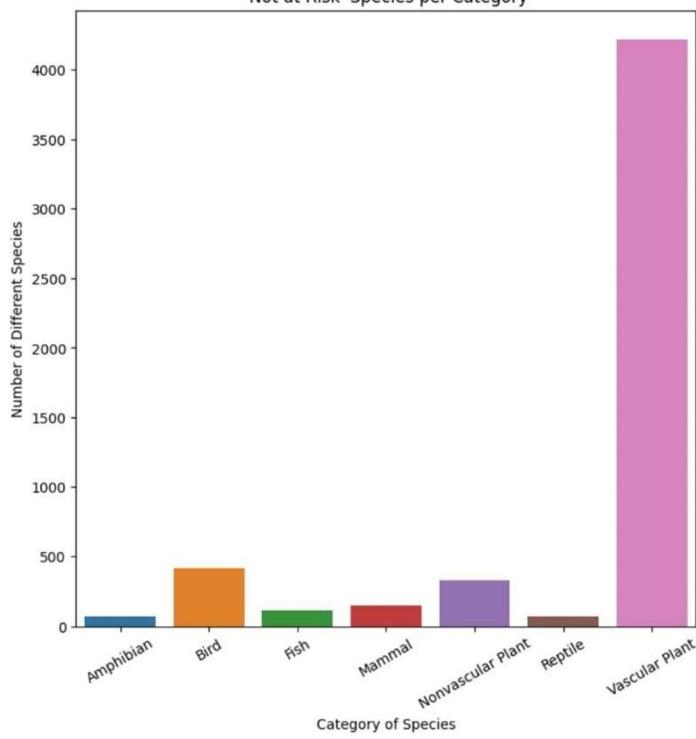
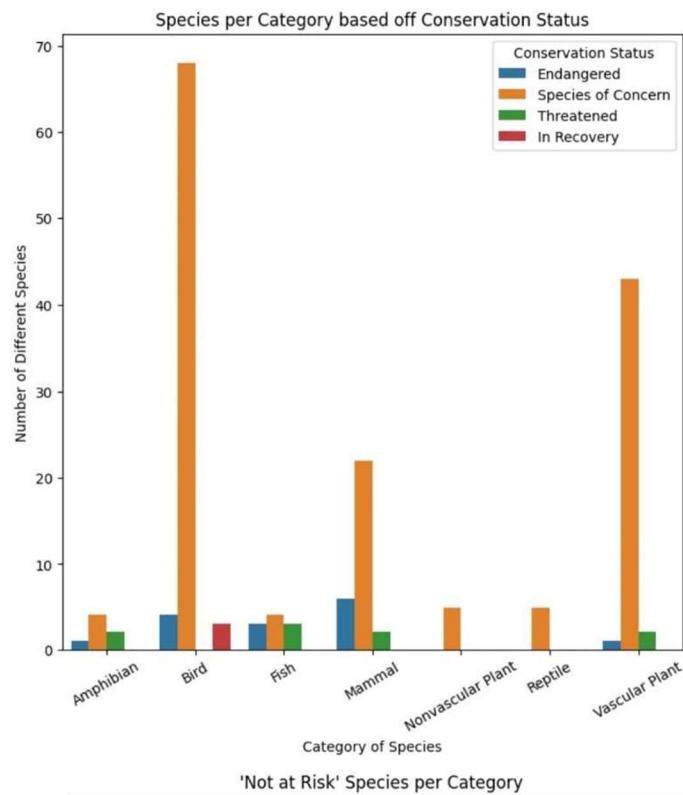
- **Unit Tests** ensured robust foundations across the backend, ML API client, and data visualizers.
- **Integration Testing** confirmed reliable communication among system modules (MongoDB, Flask ML, Node.js, React).
- **System Testing** validated platform performance and functionality using production-like scenarios.
- **User Acceptance Testing** refined the user experience and confirmed value delivery.

This structured process helped WildStat evolve into a stable, user-focused biodiversity analytics platform that delivers clear, actionable insights on species conservation.

## VI) OUTCOME OF THE PROJECT

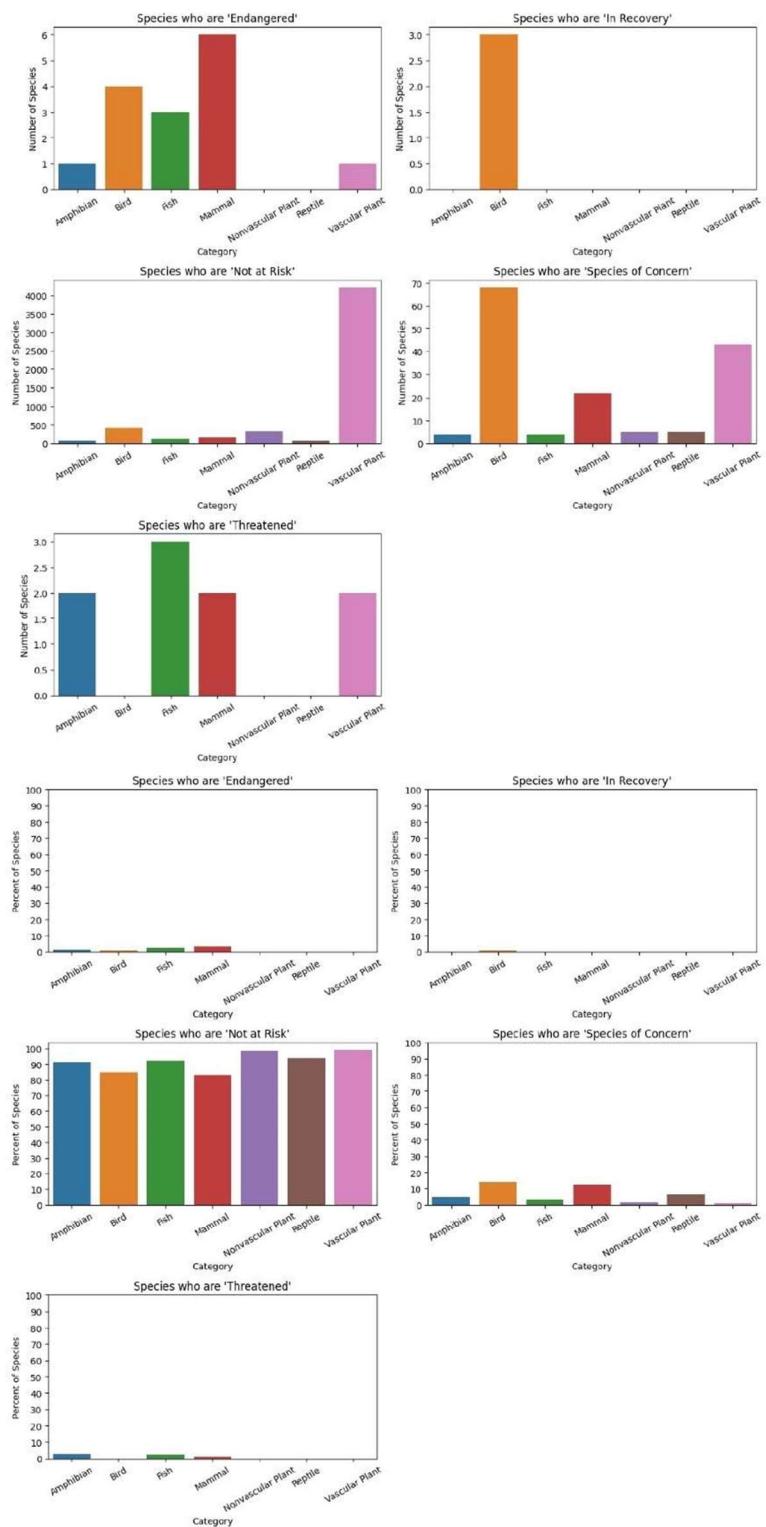
### CATEGORIES OF SPECIES

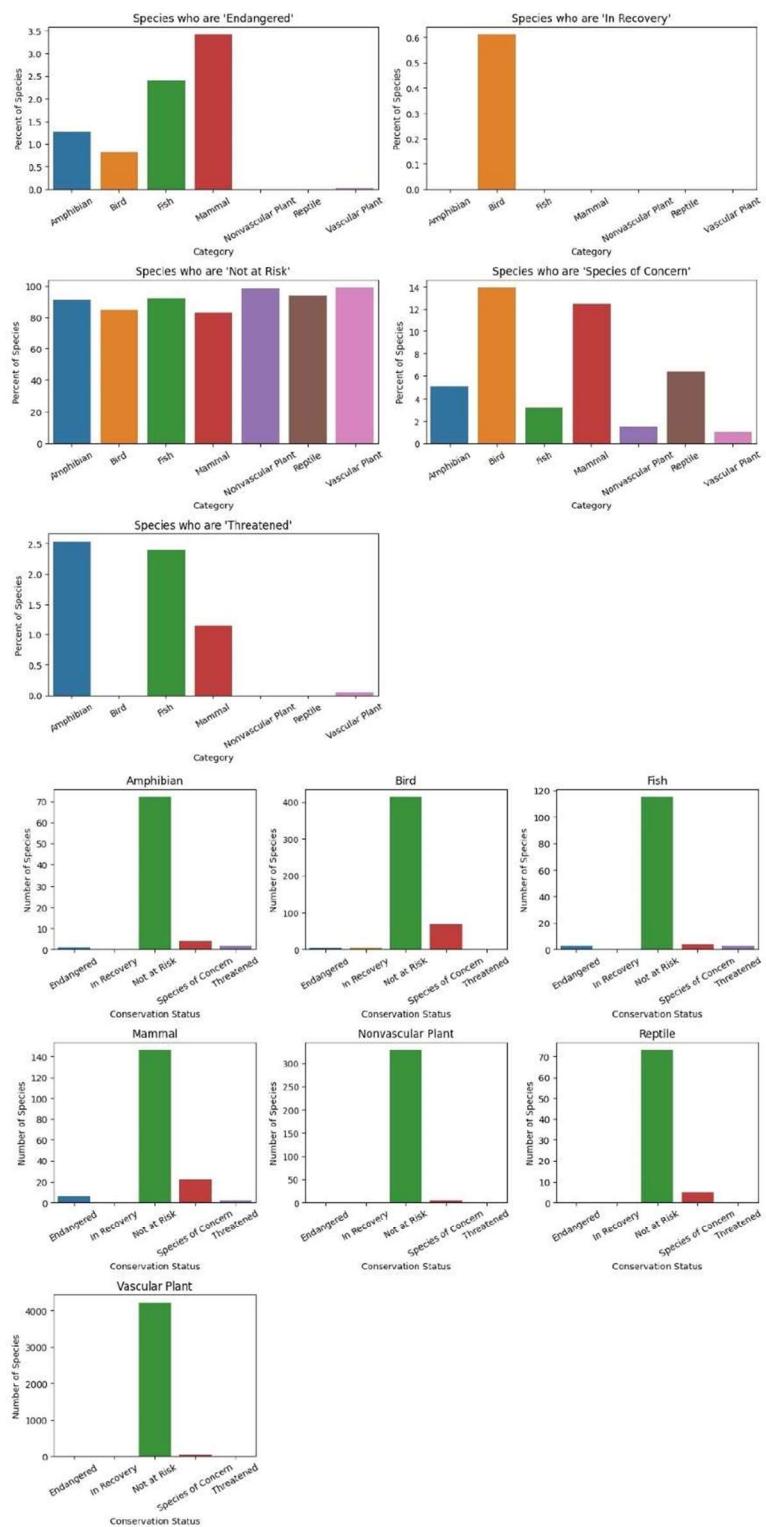


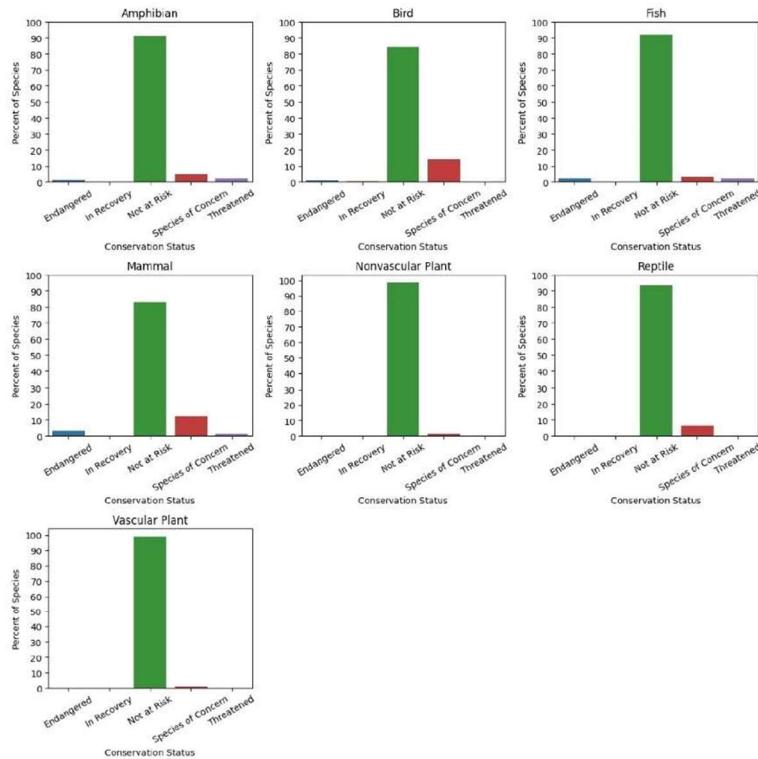


	Conservation Status	Endangered	In Recovery	Not at Risk	Species of Concern	Threatened	Total Species
Category							
Amphibian		1.0	0.0	72.0	4.0	2.0	79.0
Bird		4.0	3.0	413.0	68.0	0.0	488.0
Fish		3.0	0.0	115.0	4.0	3.0	125.0
Mammal		6.0	0.0	146.0	22.0	2.0	176.0
Nonvascular Plant		0.0	0.0	328.0	5.0	0.0	333.0
Reptile		0.0	0.0	73.0	5.0	0.0	78.0
Vascular Plant		1.0	0.0	4216.0	43.0	2.0	4262.0

	Conservation Status	Endangered	In Recovery	Not at Risk	Species of Concern	Threatened
Category						
Amphibian		1.27%	0.0%	91.14%	5.06%	2.53%
Bird		0.82%	0.61%	84.63%	13.93%	0.0%
Fish		2.4%	0.0%	92.0%	3.2%	2.4%
Mammal		3.41%	0.0%	82.95%	12.5%	1.14%
Nonvascular Plant		0.0%	0.0%	98.5%	1.5%	0.0%
Reptile		0.0%	0.0%	93.59%	6.41%	0.0%
Vascular Plant		0.02%	0.0%	98.92%	1.01%	0.05%

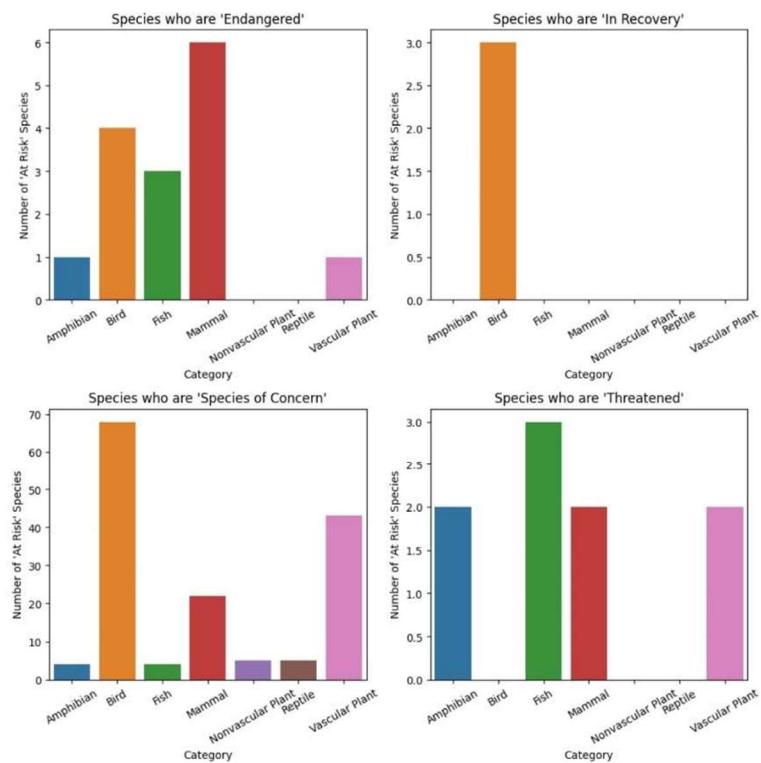


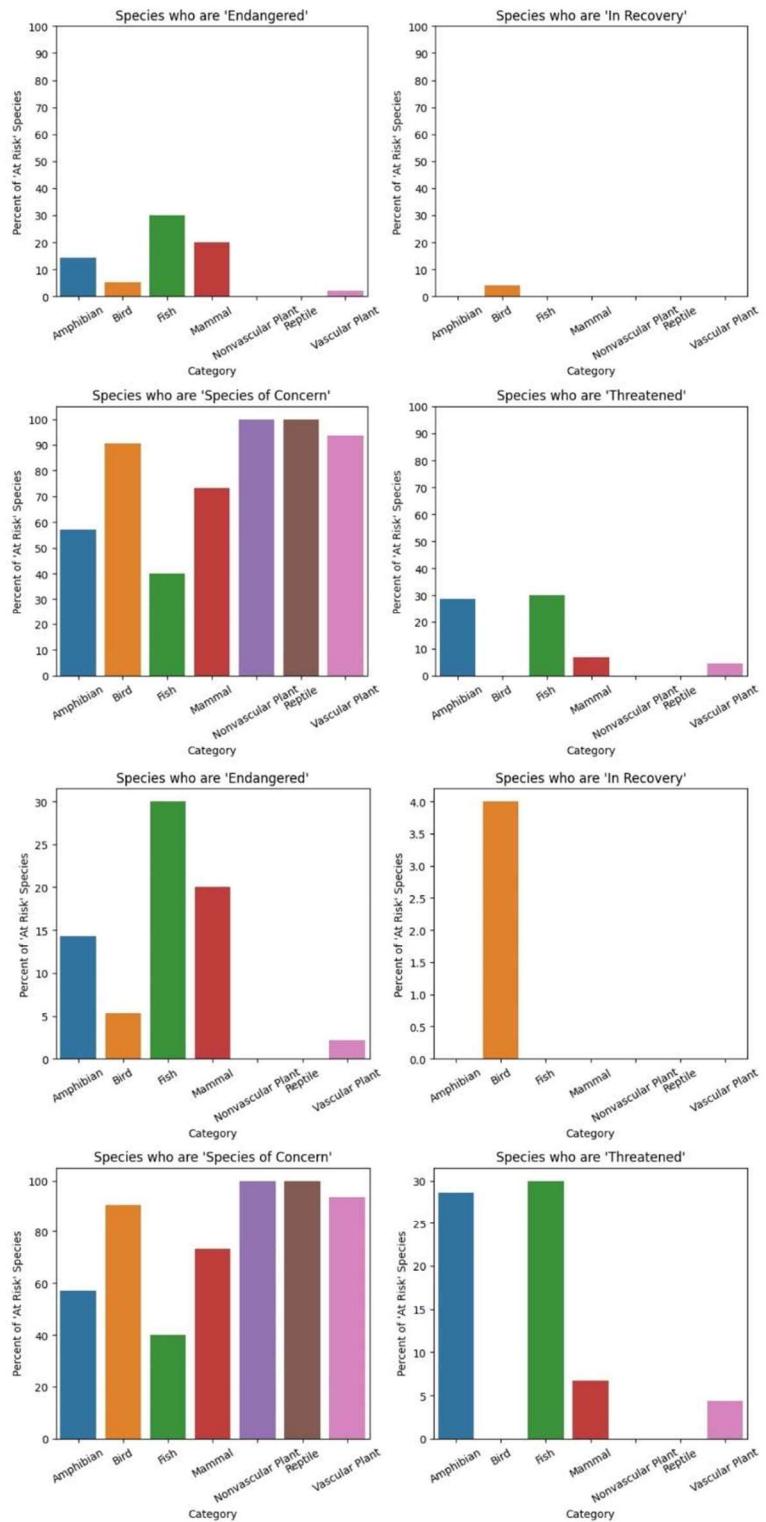


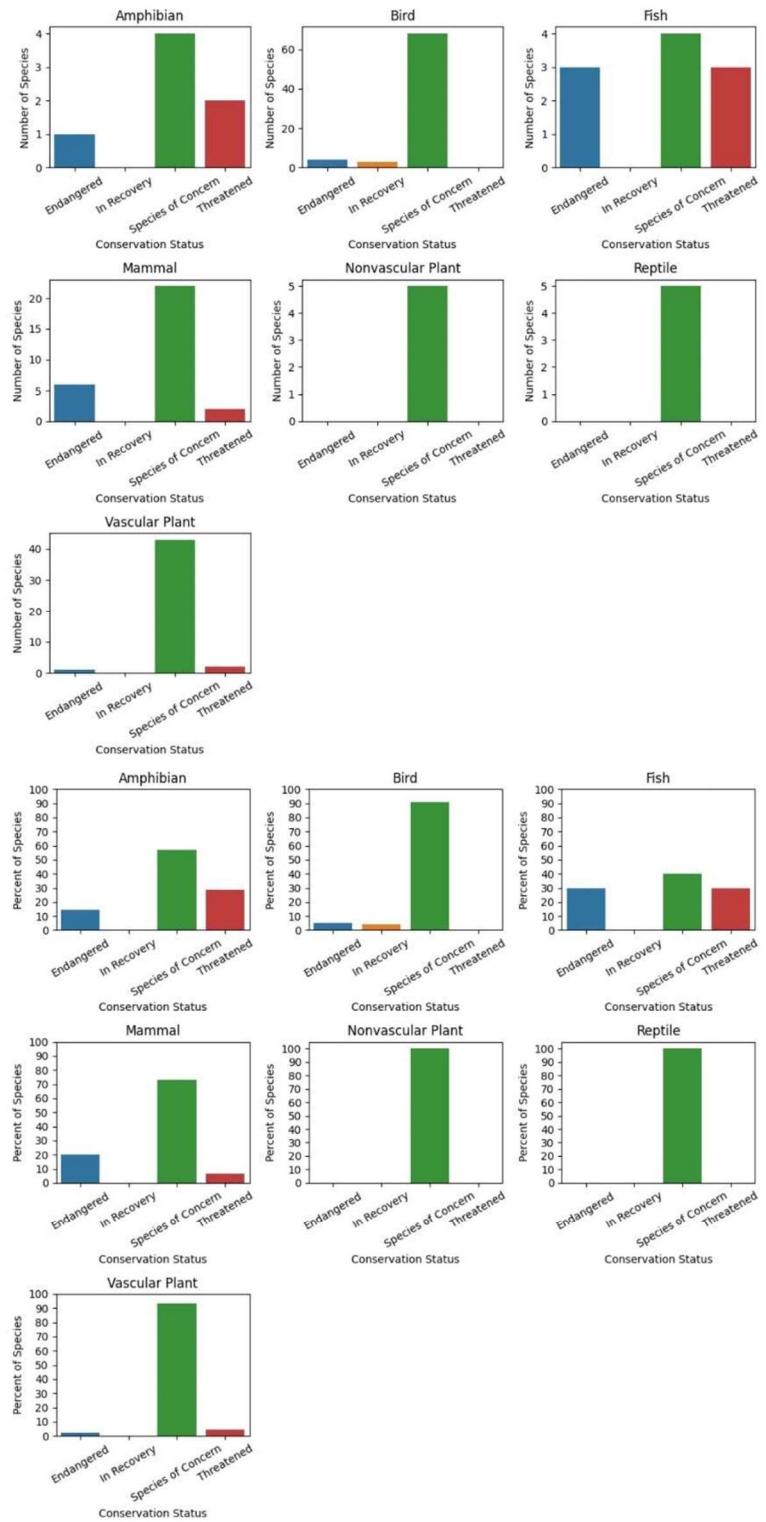


	Conservation Status	Endangered	In Recovery	Species of Concern	Threatened	Total Species
<b>Category</b>						
Amphibian		1.0	0.0	4.0	2.0	7.0
Bird		4.0	3.0	68.0	0.0	75.0
Fish		3.0	0.0	4.0	3.0	10.0
Mammal		6.0	0.0	22.0	2.0	30.0
Nonvascular Plant		0.0	0.0	5.0	0.0	5.0
Reptile		0.0	0.0	5.0	0.0	5.0
Vascular Plant		1.0	0.0	43.0	2.0	46.0

	Conservation Status	Endangered	In Recovery	Species of Concern	Threatened
Category					
Amphibian		14.29%	0.0%	57.14%	28.57%
Bird		5.33%	4.0%	90.67%	0.0%
Fish		30.0%	0.0%	40.0%	30.0%
Mammal		20.0%	0.0%	73.33%	6.67%
Nonvascular Plant		0.0%	0.0%	100.0%	0.0%
Reptile		0.0%	0.0%	100.0%	0.0%
Vascular Plant		2.17%	0.0%	93.48%	4.35%

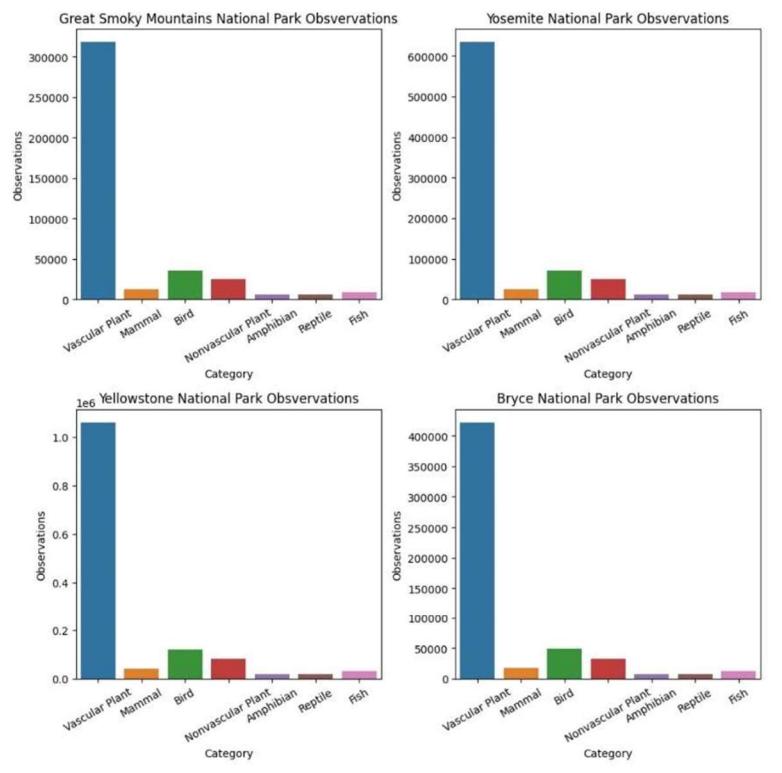




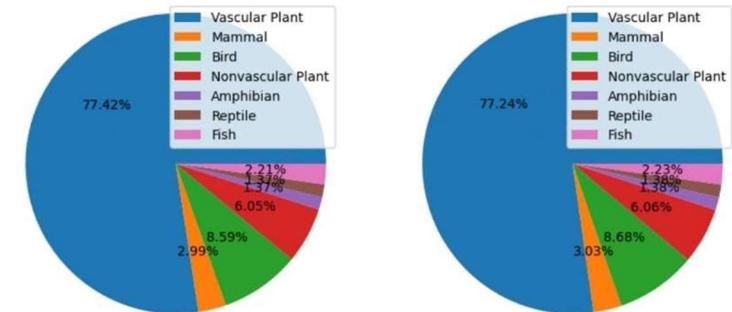


	Category	Amphibian	Bird	Fish	Mammal	Nonvascular Plant	Reptile	Vascular Plant	Total Obsv's
Park Name									
Bryce National Park		7,299	48,383	12,223	16,823	32,992	7,854	422,585	548,159
Great Smoky Mountains National Park		5,622	35,290	9,068	12,301	24,857	5,616	318,071	410,825
Yellowstone National Park		19,191	119,219	30,131	41,905	83,021	19,300	1,060,769	1,373,536
Yosemite National Park		11,309	71,290	18,353	24,887	49,783	11,335	634,515	821,472

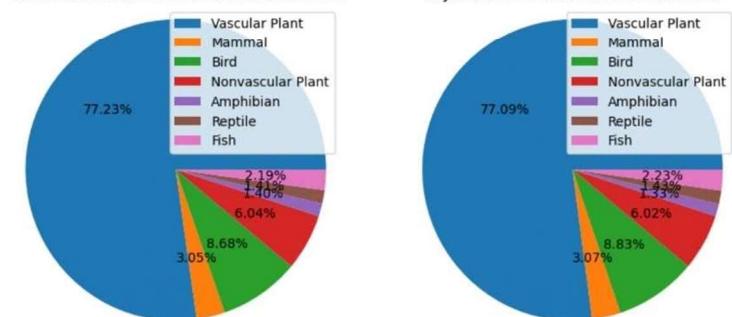
	Category	Amphibian	Bird	Fish	Mammal	Nonvascular Plant	Reptile	Vascular Plant
Park Name								
Bryce National Park		1.33%	8.83%	2.23%	3.07%	6.02%	1.43%	77.09%
Great Smoky Mountains National Park		1.37%	8.59%	2.21%	2.99%	6.05%	1.37%	77.42%
Yellowstone National Park		1.4%	8.68%	2.19%	3.05%	6.04%	1.41%	77.23%
Yosemite National Park		1.38%	8.68%	2.23%	3.03%	6.06%	1.38%	77.24%



Great Smoky Mountains National Park Observations %'s      Yosemite National Park Observations %'s



Yellowstone National Park Observations %'s



Referring to a table from earlier...

This table shows the percentage that each conservation status makes up within each category of species

	Conservation Status	Endangered	In Recovery	Not at Risk	Species of Concern	Threatened
Category						
Amphibian		1.27%	0.0%	91.14%	5.06%	2.53%
Bird		0.82%	0.61%	84.63%	13.93%	0.0%
Fish		2.4%	0.0%	92.0%	3.2%	2.4%
Mammal		3.41%	0.0%	82.95%	12.5%	1.14%
Nonvascular Plant		0.0%	0.0%	98.5%	1.5%	0.0%
Reptile		0.0%	0.0%	93.59%	6.41%	0.0%
Vascular Plant		0.02%	0.0%	98.92%	1.01%	0.05%

This table includes all the species ('Not at Risk' species as well) and the Mammal category has the highest percentage of 'Endangered' species.

For example, 3.41% of all Mammals are 'Endangered'.

	Conservation Status	Endangered	In Recovery	Not at Risk	Species of Concern	Threatened	Total (Excluding 'Endangered')
Category							
Amphibian		1.0	0.0	72.0	4.0	2.0	78.0
Bird		4.0	3.0	413.0	68.0	0.0	484.0
Fish		3.0	0.0	115.0	4.0	3.0	122.0
Mammal		6.0	0.0	146.0	22.0	2.0	170.0
Nonvascular Plant		0.0	0.0	328.0	5.0	0.0	333.0
Reptile		0.0	0.0	73.0	5.0	0.0	78.0
Vascular Plant		1.0	0.0	4216.0	43.0	2.0	4261.0

The actual tests were conducted using Python and a Python library, so the tests won't be shown, but here is an example of what information the test would need to make a comparison.

For example, if you were comparing Mammals to Fish...

	Endangered	Total (Excluding 'Endangered')
Mammal	6.0	170.0
Fish	3.0	122.0

## **VII) Conclusion and Future Scope**

The **WildStat** platform represents a significant leap forward in the field of biodiversity conservation through its data-centric approach and modern technological framework. By harnessing the power of historical environmental datasets and integrating them with machine learning algorithms, WildStat successfully delivers insightful and actionable analysis of endangered species patterns across U.S. National Parks. The platform not only highlights current ecological challenges but also equips users with the tools needed to understand the underlying causes of species vulnerability.

Designed with a modular architecture—comprising a decoupled **frontend (Next.js)**, robust **backend (Node.js)**, and a dedicated **ML API (Flask)**—WildStat ensures high scalability, efficient maintenance, and seamless integration of future enhancements. This design enables independent development and testing of each component while preserving overall system coherence.

Throughout the project's lifecycle—from requirement analysis and system design to development, testing, and deployment—every step was guided by the goal of creating a **reliable, user-friendly, and impactful tool** for researchers, conservationists, and policymakers. The result is a platform that not only meets its original specifications but also exceeds expectations in terms of usability and real-world applicability.

In essence, WildStat bridges the critical information gap in biodiversity tracking and assessment. It provides a scientific foundation for **proactive environmental decision-making** and sets the stage for more informed conservation strategies across protected regions. With its current capabilities and envisioned enhancements, WildStat stands as a promising step toward sustainable environmental stewardship.

## **Future Scope**

The journey of WildStat doesn't end with its initial deployment. In fact, it opens a gateway to a wide range of potential developments that can significantly enhance its capabilities and extend its global impact. Below are some of the major directions for future growth:

### **Expanded Environmental Data Integration**

To improve the granularity and contextual relevance of species risk analysis, future versions of WildStat could incorporate additional environmental parameters. This includes data such as terrain type, vegetation density, seasonal migration trends, and water sources proximity, which would refine predictive models and offer deeper ecological insights.

### **Advanced Predictive Analytics**

Building on its current machine learning foundation, the platform could evolve to support time-series forecasting. This would enable prediction of species endangerment probabilities based on climate change projections, urbanization trends, and land-use modifications. Incorporating AI-driven simulations can further aid in evaluating long-term conservation outcomes.

### **Interactive and Community-Driven Features**

To foster **community involvement** and encourage citizen science, WildStat can include interactive features like:

- **Species sighting maps** for users to report observations.
- **Discussion forums** specific to each national park or ecosystem.
- **User-contributed data uploads**, which can be moderated and verified to enhance dataset richness.

These features would transform WildStat from a static analysis tool into a **collaborative conservation hub**.

## **Educational and Policy Support Tools**

To support decision-makers and educators, the platform could incorporate:

- **Policy simulation modules** that allow users to test conservation strategies.
- **Interactive dashboards** for educational purposes, tailored for schools, colleges, and public outreach.
- **Infographics and dynamic reports** that present ecological data in simplified, visually compelling formats.

These additions would broaden WildStat's influence in both **governmental planning** and **environmental education**.

## **Global Dataset Integration**

Expanding beyond U.S. National Parks, the platform has the potential to scale up to analyze biodiversity data from global protected areas, such as:

- **UNESCO World Heritage Sites**
- **IUCN Protected Areas**
- **Rainforests, wetlands, and marine reserves**

## VIII) References

- **Kareiva, P., & Marvier, M.** (2012). *What is conservation science?* *BioScience*, 62(11), 962–969. (This foundational paper outlines the modern framework of conservation science, relevant to WildStat's interdisciplinary approach.)
- **Challender, D. W. S., Hinsley, A., Veríssimo, D., & Milner-Gulland, E. J.** (2018). *Effective conservation requires understanding human behavior.* *Conservation Letters*, 11(6), e12531. (This article supports the integration of human factors into conservation platforms like WildStat.)
- **Elith, J., & Leathwick, J. R.** (2009). *Species distribution models: Ecological explanation and prediction across space and time.* *Annual Review of Ecology, Evolution, and Systematics*, 40, 677–697. (A foundational study for species distribution modeling, a core analytical technique used in WildStat.)
- **Turner, W., Spector, S., Gardiner, N., Fladeland, M., Sterling, E., & Steininger, M.** (2003). *Remote sensing for biodiversity science and conservation.* *Trends in Ecology & Evolution*, 18(6), 306–314.
- (This paper justifies WildStat's use of geospatial data to track biodiversity patterns.)
- **Cardinale, B. J., Duffy, J. E., Gonzalez, A., Hooper, D. U., Perrings, C., Venail, P., ... & Naeem, S.** (2012). *Biodiversity loss and its impact on humanity.* *Nature*, 486(7401), 59–67. (Emphasizes the global consequences of biodiversity decline, reinforcing the importance of analytical tools like WildStat.)
- **Hunter Jr., M. L., & Gibbs, J. P.** (2006). *Fundamentals of Conservation Biology* (3rd ed.). Wiley-Blackwell. ISBN: **9781405135450**  
(Provides essential background in conservation principles that inform the WildStat platform.)
- **National Park Service (NPS).** (n.d.). *Species Inventory and Monitoring Reports.*  
(A primary data source for species records used in WildStat's analysis of endangered species trends.)
- **USGS (U.S. Geological Survey).** (n.d.). *Environmental Data Resources for Biodiversity Analysis.* (Offers key environmental datasets—such as climate, terrain, and land use—used by WildStat for contextual analysis.)