

## Unit 2: Tourism Data Collection, Preprocessing and Exploratory Data Analysis

Course: AI and Tourism – MIT-AI @ Gandaki University

Bidur Devkota, PhD  
GCES Pokhara

December 17, 2025

## Unit 2 Overview

**Duration:** 8 Hours

### Learning Objectives:

- Understand tourism data sources and collection methods
- Learn to preprocess and transform raw tourism data
- Apply descriptive statistics and visualization techniques
- Conduct exploratory analysis to identify trends and patterns

### Structure:

- 1 Data Sources and APIs
- 2 Data Preprocessing
- 3 Descriptive Statistics & Visualization
- 4 Case Study Applications

## 2.1 Data Sources and APIs

### Modern Tourism Data Landscape

Traditional Sources	Digital Sources
<ul style="list-style-type: none"><li>• Government surveys</li><li>• Airport arrival records</li><li>• Hotel registration data</li><li>• Tourism board reports</li></ul>	<ul style="list-style-type: none"><li>• Online booking platforms</li><li>• Social media posts</li><li>• Review websites</li><li>• Mobile app data</li></ul>

### Key Technologies

VGI (Volunteered Geographic Information) + APIs = Real-time Tourism Insights

## Volunteered Geographic Information (VGI)

**Definition:** Geographic data created by non-experts via digital platforms

### Tourism VGI Examples:

- TripAdvisor reviews/ratings
- Google Maps location reviews
- Instagram geotagged posts
- Foursquare check-ins
- Travel blog posts

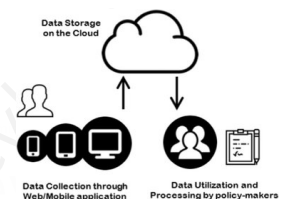


Figure: VGI Data Flow

<https://link.springer.com/article/10.1007/s13132-017-0504-y>

### Impact

Provides real-time, localized insights impossible with traditional surveys

## APIs for Tourism Data Collection

API	Data Provided
Google Places API	Location details, reviews, popularity times, photos
TripAdvisor API	Hotel/attraction reviews, ratings, rankings
Booking.com API	Availability, prices, property features
OpenWeather API	Weather conditions, forecasts
Skyscanner API	Flight prices, availability

[Table](#): Major Tourism Data APIs

## API Considerations

### Key Considerations:

- API rate limits
- Data usage terms
- Authentication
- Data pagination
- Error handling

## 2.2 Data Preprocessing

### Why Preprocess Tourism Data?

Common Issues	Preprocessing Steps
<ul style="list-style-type: none"><li>• Missing values</li><li>• Inconsistent formats</li><li>• Duplicate entries</li><li>• Outliers</li><li>• Unstructured text</li></ul>	<ol style="list-style-type: none"><li>1 Data cleaning</li><li>2 Handling missing data</li><li>3 Text processing</li><li>4 Feature engineering</li><li>5 Normalization</li></ol>

### Golden Rule

Garbage In = Garbage Out  
Quality preprocessing = Reliable insights

## Processing Unstructured Tourism Data

### Text Processing Pipeline:

- 1 Tokenization
- 2 Lowercasing
- 3 Stopword removal
- 4 Stemming/Lemmatization
- 5 Sentiment analysis
- 6 Entity recognition

### Example Review:

"Amazing hotel! Great location near the beach.  
Service was exceptional but rooms were small."

### Extracted Features:

- Sentiment: Positive (0.7)
- Keywords: hotel, beach, service, rooms
- Aspects: location, service, room size

## Tourism Feature Engineering

Raw Feature	Engineered Features
Arrival Date	Season, Day of Week, Holiday Flag, Peak Season Flag
Hotel Price	Price Category, Relative Price Index, Discount Percentage
Review Text	Sentiment Score, Topic Categories, Review Length
Location Coordinates	Distance to City Center, Proximity to Attractions
Booking Date	Advance Booking Days, Lead Time Category

Table: Feature Engineering Examples

## 2.3 Descriptive Statistics & Visualization

### Descriptive Statistics for Tourism

#### Central Tendency

- Average stay duration
- Mean tourist spending
- Median hotel price
- Mode travel package

#### Dispersion

- Range of tour prices
- Variance in monthly arrivals
- Standard deviation of ratings
- IQR of spending patterns

### Application Example

Calculating average daily rate (ADR) and revenue per available room (RevPAR) for hotel performance analysis

## Tourism Data Visualization Techniques

### Temporal Patterns

- Line charts (arrivals over time)
- Heatmaps (seasonal demand)
- Calendar plots (daily occupancy)

### Spatial Patterns

- Choropleth maps
- Point density maps
- Flow maps (tourist movement)

### Comparative Analysis

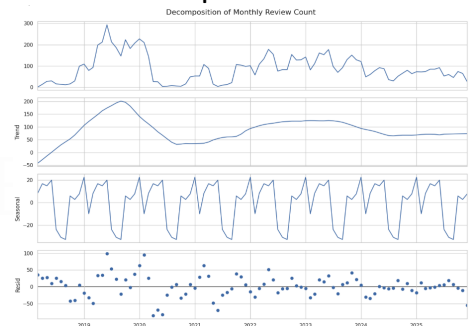
- Bar charts (spending by category)
- Pie charts (visitor demographics)
- Radar charts (destination attributes)

### Relationship Analysis

- Scatter plots (price vs. rating)
- Correlation matrices
- Bubble charts

## Visualization Examples

### Time Series Decomposition

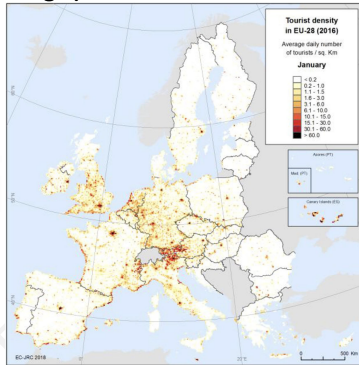


timeseries decomposition(Trend + Seasonality + Residuals)  
of reviews of Bindabasini Temple

**Tools:** matplotlib, seaborn, plotly, Tableau, Power BI

## Visualization Examples

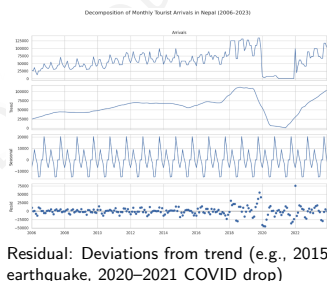
### Geographic Distribution



## Nepal Tourism Seasonality Analysis

### Observed:

- raw monthly data
- Trend — Long-term growth (rise until 2019, collapse in 2020–2021 due to COVID, strong recovery in 2022–2023).
- Seasonal — Consistent yearly cycle: peaks in October–November (post-monsoon/autumn trekking) and March (spring), lows in June–August (monsoon) and January.
- Residual — Irregular variations/noise, with large negative residuals in 2020–2021

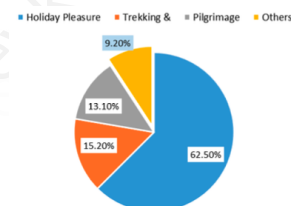


© 2025 Bidur Devkota, PhD. CC BY-NC.

## Nepal Tourism Market Segmentation

### By Purpose of Visit (2023):

- **Holiday/Pleasure:** 634,301 (62.5%)
- **Trekking & Mountaineering:** 154,262 (15.2%)
  - Expedition teams: 2,253
  - Climbers: 9,398
  - Royalty: Rs. 1,013.6 million
- **Pilgrimage:** 132,949 (13.1%)
  - Lumbini visitors: 903,883 total
  - Pashupatinath: 118,531 foreign visitors
- **Other:** 93,370 (9.2%)



### Adventure Tourism Highlights:

- Popular peaks: Everest, Ama Dablam, Manaslu
- **TIMS trekkers: 28,691 in 2023**

© 2025 Bidur Devkota, PhD. CC BY-NC.

## Geographic Distribution & Entry Points

### Top Entry Points (2023):

- **By Air:** 914,270 (90.09%) - TIA Kathmandu, Gautam Buddha Airport
- **By Land:** 100,612 (9.91%)
  - Belahiya (India border): 81,338 (81% of land arrivals)
  - Rasuwagadi (China border): 10,447
  - Kakarbhitta (India): 3,125
  - Birgunj (India): 2,006

### Top Source Markets Analysis:

Country	2023 Visitors	Change from 2022
India	319,936	+52.8%
USA	100,355	+30.2%
China	60,878	+534.1%
UK	52,865	+17.2%
Australia	38,798	+44.2%
Bangladesh	36,483	+43.8%

© 2025 Bidur Devkota, PhD. CC BY-NC.

## Economic Impact Analysis

### Foreign Exchange Earnings:

- **2023 Earnings:** \$548.2 million (Rs. 72.46 billion)
- **Growth:** 68% increase from 2022
- **Daily Spending:** \$41 per tourist per day
- **GDP Contribution:** 1.98% of GDP (hotel & restaurant sector)

### Tourism Infrastructure (2023):

- **Star Hotels:** 182
- **Non-star Hotels:** 1,416
- **Total Beds:** 54,370
- **Travel Agencies:** 4,845
- **Trekking Agencies:** 3,191
- **Tourist Guides:** 5,123
- **Trekking Guides:** 26,292

### Key Insight

Tourism recovery significantly contributes to foreign exchange reserves and employment generation

© 2025 Bidur Devkota, PhD. CC BY-NC.

## Hands-on Exercise: Nepal Tourism Data Analysis

**Dataset:** Nepal Tourism Statistics 2023 (provided PDF)

### Tasks:

- 1 Clean and structure tourist arrival data (2014-2023)
- 2 Analyze seasonality patterns for different nationalities
- 3 Calculate recovery rates post-COVID by market segment
- 4 Visualize trekking and expedition trends
- 5 Identify emerging markets and declining trends
- 6 Forecast 2024 arrivals based on historical patterns

### Expected Deliverables

- Time-series analysis of tourist arrivals
- Market segmentation dashboard
- Seasonality heatmaps
- Recovery trend analysis
- Strategic recommendations for tourism board

© 2025 Bidur Devkota, PhD. CC BY-NC.

## Recommended Tools & Libraries

### Python Ecosystem:

- pandas - Data manipulation
- numpy - Numerical computing
- matplotlib/seaborn - Visualization
- scikit-learn - Machine learning
- statsmodels - Statistical analysis
- geopandas - Spatial analysis
- nltk/spaCy - Text processing

### Other Tools:

- R with tidyverse
- Tableau/Power BI
- Google Data Studio
- QGIS (spatial analysis)
- SQL databases

### Data Sources:

- UNWTO datasets
- National tourism boards
- Kaggle tourism datasets
- API collections

© 2025 Bidur Devkota, PhD. CC BY-NC.

## Key Takeaways

- 1 **Quality preprocessing** is essential for reliable analysis
- 2 **Visualization** reveals patterns not obvious in raw data
- 3 **Seasonality and trends** drive strategic decisions
- 4 **Outlier detection** helps identify opportunities and risks

### Professional Competency

Ability to transform raw tourism data into actionable business intelligence

© 2025 Bidur Devkota, PhD. CC BY-NC.

## References & Further Reading

- 📄 Smith, J. (2022) *Tourism Analytics: Data Driven Decision Making*
- 📄 Chen, L., & Wang, Y. (2021) *Big Data in Tourism*
- 📄 UNWTO (2023) *Global Tourism Data Standards*
- 📄 Google (2023) *Travel Insights with Google*
- 📄 Kaggle Tourism Datasets <https://www.kaggle.com/tags/tourism>
- 📄 GitHub Tourism Analytics Projects <https://github.com/topics/tourism-analytics>
- 📄 Nepal Tourism Statistic 2023 [https://trade.ntb.gov.np/wp-content/uploads/2024/06/Nepal-Tourism-Statistic\\_2023-final.pdf](https://trade.ntb.gov.np/wp-content/uploads/2024/06/Nepal-Tourism-Statistic_2023-final.pdf)

Thank You

Questions?

© 2025 Bidur Devkota, PhD. CC BY-NC.