

Table of contents

| | | |
|----------|--|----------|
| 1 | Evaluation - TravelHunters Project | 2 |
| 1.1 | Evaluation Workshop Results | 2 |
| 1.1.1 | Assessment of Project Results | 2 |
| 1.1.2 | Technical Evaluation | 3 |
| 1.2 | Decisions | 3 |
| 1.2.1 | 1. Project Continuation: YES | 3 |
| 1.2.2 | 2. Next Development Steps | 3 |
| 1.2.3 | 3. Additional Data Mining Iteration: YES | 4 |
| 1.3 | Lessons Learned | 5 |
| 1.3.1 | Technical Insights | 5 |
| 1.3.2 | Project Management | 5 |
| 1.4 | Risks and Mitigation | 6 |
| 1.4.1 | Technical Risks | 6 |
| 1.4.2 | Legal Risks | 6 |
| 1.4.3 | Business Risks | 6 |
| 1.5 | Future Recommendations | 6 |
| 1.5.1 | Short-term (1-3 months) | 6 |
| 1.5.2 | Medium-term (3-6 months) | 6 |
| 1.5.3 | Long-term (6-12 months) | 7 |
| 1.6 | Conclusion | 7 |

1 Evaluation - TravelHunters Project

This report summarizes the evaluation results of the TravelHunters project and documents decisions for further development and deployment of the travel recommendation system.

1.1 Evaluation Workshop Results

Date: July 3, 2025

Participants: Leona Kryeziu, Evan Blazo, Joan Felber, Jakub Baranec

1.1.1 Assessment of Project Results

1.1.1.1 Fulfilled Requirements

- **Data Collection:** Successfully collected and structured >2,800 travel datasets
 - 2,047 hotels from Booking.com
 - 90 activities from GetYourGuide
 - 671 destinations from various sources
- **Global Coverage:** Hotels and activities available from different continents
- **Data Quality:** >95% completeness in critical fields (name, location, rating)
- **Automation:** Fully automated web scraping with pagination implemented
- **System Architecture:** Scalable SQLite database with structured schema

1.1.1.2 Partially Fulfilled Requirements

- **Geographic Balance:** Europe bias in hotel data (predominantly European cities)
- **Image URLs:** Not all entries have working image URLs
- **Price Formats:** Inconsistent price representation across different sources

1.1.1.3 Unfulfilled Requirements

- **Real-time Data:** Static data collection without live updates
- **User Reviews:** Limited number of user reviews per entry
- **Multi-Language Support:** Only English/German content captured

1.1.2 Technical Evaluation

Web-Scraping Performance:

- **Success Rate:** 95% of requests successful
- **Throughput:** ~50 items per minute
- **Error Handling:** Robust error handling implemented
- **Duplicate Detection:** 6,200+ duplicates successfully removed

Data Processing Quality:

- **Cleaning:** Automatic text normalization and URL validation
- **Integration:** Successful merging of all data sources
- **Consistency:** Uniform schema across all categories

1.2 Decisions

1.2.1 1. Project Continuation: YES

Rationale:

- Proof-of-concept successfully demonstrated
- Technical feasibility for larger datasets confirmed
- Solid foundation for advanced features available
- Summer School learning objectives fully achieved

1.2.2 2. Next Development Steps

1.2.2.1 Phase 1: Optimization (Immediately actionable)

- **Improve Data Quality:**
 - Implement better image URL extraction
 - Normalize price formats
 - Expand geographic coverage

- **Performance Tuning:**
 - Implement parallel scraping
 - Introduce caching mechanisms
 - Optimize database indexing

1.2.2.2 Phase 2: Feature Extension (Medium-term)

- **Develop Recommendation System:**
 - Implement content-based filtering
 - Similarity algorithms for hotels/activities
 - User preference engine
- **API Development:**
 - RESTful API for data access
 - Enable real-time queries
 - Implement rate limiting

1.2.2.3 Phase 3: Production System (Long-term)

- **Develop Web Interface:**
 - User-friendly search interface
 - Interactive map integration
 - Mobile-responsive design
- **Commercialization:**
 - Affiliate links to booking platforms
 - Premium features for users
 - Business intelligence dashboard

1.2.3 3. Additional Data Mining Iteration: YES

1.2.3.1 Priority Improvements

1. **Expand Data Acquisition:**
 - Airbnb and alternative accommodation platforms
 - TripAdvisor for extended reviews
 - Local tourism websites for better regional coverage
 - Social media sentiment analysis

2. Improve Data Quality:

- Image processing for better image URLs
- Geocoding for precise location data
- Automate categorization and tagging
- Improve duplicate detection through ML

3. Analytical Improvements:

- Price prediction models
- Seasonality analysis
- Trend detection in destinations
- User behavior modeling

1.3 Lessons Learned

1.3.1 Technical Insights

1. Web-Scraping Challenges:

- Rate-limiting is critical for website stability
- Robust error-handling prevents data loss
- Pagination-handling complex but important for completeness

2. Data Processing:

- Early duplicate detection saves resources
- Uniform data models essential for integration
- Automated validation reduces manual rework

3. Infrastructure:

- SQLite sufficient for prototyping, PostgreSQL for production
- Documentation with Quarto very effective
- Version control for data pipelines important

1.3.2 Project Management

1. **Agile Development** works well for Data Science projects
2. **Continuous Evaluation** prevents direction changes
3. **Stakeholder Feedback** collect early and regularly

1.4 Risks and Mitigation

1.4.1 Technical Risks

- **Website Changes:** Plan regular spider updates
- **Scaling Problems:** Early architecture planning
- **Data Quality:** Implement automated monitoring tools

1.4.2 Legal Risks

- **robots.txt Compliance:** Automated verification
- **Rate-Limiting:** Respectful scraping practices
- **Copyright:** Only use publicly available data

1.4.3 Business Risks

- **Competition:** Focus on unique features and better UX
- **Data Currency:** Establish regular update cycles
- **User Acceptance:** Continuous user testing

1.5 Future Recommendations

1.5.1 Short-term (1-3 months)

1. **Clean up codebase** and documentation
2. **Automated tests** for all spiders implemented
3. **CI/CD Pipeline** for continuous data updates
4. **Monitoring Dashboard** for data quality

1.5.2 Medium-term (3-6 months)

1. **Machine Learning Pipeline** for recommendations development
2. **API Gateway** for external data access
3. **Performance Optimization** for larger datasets
4. **International Expansion** of data sources

1.5.3 Long-term (6-12 months)

1. **Production Web Application** development
2. **Mobile Apps** for iOS and Android
3. **AI-powered Travel Planning** implementation
4. **Partnerships** with travel providers establishment

1.6 Conclusion

The TravelHunters project has successfully demonstrated how web scraping, data processing, and machine learning can be applied to a real problem in the tourism sector.

Key Achievements:

- Automated data collection of >2,800 travel datasets
- Robust data pipeline with 95% success rate
- Scalable architecture for future extensions
- Practical application of Data Science methods

Project Status: SUCCESSFULLY COMPLETED

The project forms a solid foundation for further developments and has achieved or exceeded all originally set goals.

Evaluated by: Data Science Summer School Team

Approval: Project Leadership

Next Review: Upon Phase 2 Implementation