

Fintech App Data Collection and Analysis: Research Methodology

Executive Summary

This document presents the comprehensive methodology employed in the ISB Fintech App Research Project, detailing the systematic approach to collecting, processing, and analyzing historical data for Indian fintech applications. The research combines advanced web scraping techniques, machine learning algorithms, and data validation strategies to create a robust dataset for market analysis.

1. Research Design and Objectives

1.1 Research Questions

- How has the Indian fintech app ecosystem evolved over time?
- What are the key performance indicators and trends in fintech app adoption?
- How can we systematically map applications to their parent companies?
- What temporal patterns exist in user ratings, downloads, and engagement metrics?

1.2 Data Requirements

- **Temporal Coverage:** Multi-year historical data spanning fintech app lifecycle
- **Breadth:** Comprehensive coverage of Indian fintech applications
- **Depth:** Rich metadata including ratings, downloads, reviews, and descriptions
- **Quality:** High accuracy and validation for research-grade analysis

2. Phase 1: Application Discovery and Collection

2.1 Search Strategy

Objective: Identify comprehensive set of fintech applications in Indian market

Implementation: `apk_mirror_app_scraper.py`, `wayback_single_app_scraper.py`

Methodology:

1. Keyword-Based Discovery

- 2. Fintech terminology: payments, banking, lending, investment, crypto
- 3. Regional terms: UPI, NEFT, RTGS, digital wallet
- 4. Company-specific searches for known fintech entities
- 5. **Systematic Expansion**
- 6. Initial seed set of 50 known fintech apps
- 7. Iterative expansion through related app suggestions
- 8. Developer-based discovery (apps from same fintech companies)
- 9. Category-based exploration within financial services

10. Data	Capture	Protocol
<pre>python # Key data points collected per application app_data = { 'app_name': str, 'package_id': str, 'developer': str, 'category': str, 'description': str, 'initial_discovery_date': datetime }</pre>		

2.2 Quality Assurance

- **Relevance Filtering:** Manual verification of fintech classification
- **Duplicate Detection:** Package ID and name-based deduplication
- **Completeness Validation:** Ensuring minimum required metadata present

3. Phase 2: Company Database Integration and Matching

3.1 Company Data Source

Source: Traxcn Indian Fintech Company Database **Coverage:** 500+ verified fintech companies **Fields:** Company name, sector, business model, domicile, funding status

3.2 Multi-Algorithm Matching Strategy

3.2.1 Fuzzy String Matching

Implementation: fuzzy_string_app_company_matcher.py

Algorithm: Levenshtein distance with threshold optimization

```

from fuzzywuzzy import fuzz, process

def fuzzy_match(app_name, company_names, threshold=85):
    matches = process.extractBests(app_name, company_names,
                                   scorer=fuzz.ratio, score_cutoff=threshold)

    return matches

```

Validation Strategy:

- Developer name to company name matching
- App name to company name/brand matching
- Cross-validation with multiple scoring methods

3.2.2 NLP Semantic Matching

Implementation: `nlp_semantic_app_company_matcher.py`

Technology: Sentence Transformers with 'all-MiniLM-L6-v2' model

Process:

1. Text Preprocessing

```

python app_text = app_name + ' ' + developer + ' ' + description
company_text = company_name + ' ' + sector + ' ' + business_model + ' ' + overview

```

1. Embedding Generation

```

python model = SentenceTransformer('all-MiniLM-L6-v2')
app_embeddings = model.encode(app_texts, convert_to_tensor=True)
company_embeddings = model.encode(company_texts, convert_to_tensor=True)

```

1. Similarity Computation

```

python similarity_matrix = util.cos_sim(app_embeddings,
company_embeddings)

```

Threshold Configuration:

- Primary threshold: 0.50 (moderate confidence)
- High confidence: 0.70+ (automatic acceptance)
- Manual review: 0.40-0.50 (human verification required)

3.2.3 Rule-Based Validation

Implementation: Custom validation logic

Rules Applied:

1. **Direct Name Matching:** Exact or near-exact company name in app name/developer
2. **Domain Matching:** Email domains matching company websites
3. **Cross-Reference Validation:** Known subsidiaries and brand relationships
4. **Geographic Constraints:** Domicile consistency checks

3.3 Match Quality Assurance

- **Confidence Scoring:** Probabilistic confidence assignment
- **Manual Verification:** Human review for ambiguous cases
- **Cross-Validation:** Multiple algorithm consensus requirement
- **Iterative Refinement:** Continuous improvement based on validation results

4. Phase 3: Historical Data Collection via Wayback Machine

4.1 Wayback Machine Integration Strategy

Implementation: `wayback_bulk_historical_scraper.py`

Objectives:

- Comprehensive temporal coverage for each matched application
- High-quality HTML content preservation
- Scalable and respectful scraping practices

4.2 Technical Implementation

4.2.1 Snapshot Discovery

```
def get_wayback_snapshots(url):  
    """Retrieve all available snapshots using CDX API"""  
    cdx_api_url = f"https://web.archive.org/cdx/search/cdx?url={encoded_url}"  
    cdx_api_url += "&output=json&fl=timestamp,original,statuscode,digest,length"  
    cdx_api_url += "&filter=statuscode:200&collapse=timestamp:6"
```

Benefits:

- Efficient discovery of available snapshots
- Temporal distribution optimization

- Status code filtering for quality assurance

4.2.2 Advanced Scraping Infrastructure

Proxy Rotation System:

```
def get_free_proxies():  
    """Dynamic proxy discovery and rotation"""  
    # Scrapes free-proxy-list.net for HTTPS proxies  
    # Implements proxy health checking  
    # Maintains rotating proxy pool
```

Benefits:

- Distributed request sourcing
- Rate limiting mitigation
- Geographic diversification

Adaptive Rate Limiting:

```
def adaptive_delay():  
    """Intelligent delay calculation"""  
    base_delay = 10 # Base delay in seconds  
    jitter = random.uniform(1, 5) # Random component  
    backoff_factor = 2 ** retry_count # Exponential backoff  
    return base_delay + jitter + backoff_factor
```

Benefits:

- Respectful server interaction
- Reduced blocking probability
- Improved success rates

4.2.3 Error Handling and Resilience

Multi-Level Retry Strategy:

1. **Network Level:** Connection timeouts and retries
2. **HTTP Level:** Status code-based retry logic
3. **Content Level:** Validation of retrieved content
4. **Application Level:** Failed request queuing and reprocessing

Robust Error Recovery:

```
def get_snapshot_content(timestamp, url, max_retries=3):
    for attempt in range(max_retries):
        try:
            # Attempt request with current configuration
            response = session.get(wayback_url, proxies=proxy, timeout=30)
            if response.status_code == 200:
                return response.text
        except ConnectionError:
            # Implement exponential backoff
            wait_time = (2 ** attempt) * 60 + random.uniform(1, 30)
            time.sleep(wait_time)
```

4.3 Data Storage Strategy

- **Raw HTML Preservation:** Complete snapshot storage for reproducibility
- **Structured Metadata:** Extracted data points in CSV format
- **Temporal Organization:** Directory structure by app and time period
- **Compression:** Efficient storage for large datasets

5. Phase 4: HTML Processing and Data Extraction

5.1 Dual-Library Approach

Rationale: Different HTML parsers excel at different layouts and conditions

5.1.1 lxml Implementation

Implementation: `html_data_extractor_lxml.py`

Advantages:

- Fast C-based parsing engine
- Robust XPath expression support
- Memory efficient for large files
- Superior handling of malformed HTML

XPath Strategy:

```
def get_first_nonempty(tree, xpaths):
    """Cascade through multiple XPath expressions"""
    for xpath in xpaths:
        results = tree.xpath(xpath)
        if results and results[0].strip():
            return results[0].strip()
    return ""
```

5.1.2 BeautifulSoup Implementation

Implementation: `html_data_extractor_beautifulsoup.py`

Advantages:

- Intuitive CSS selector support
- Excellent handling of nested structures
- Robust text extraction capabilities
- Strong community and documentation

Selector Strategy:

```
def extract_with_fallbacks(soup, selectors):
    """Try multiple CSS selectors for robustness"""
    for selector in selectors:
        element = soup.select_one(selector)
        if element and element.text.strip():
            return element.text.strip()
    return ""
```

5.2 Best-of-Both Integration

Algorithm: Intelligent selection of best extraction result

Selection Criteria:

1. **Completeness:** Result with most non-empty fields
2. **Consistency:** Cross-validation between parsers
3. **Data Quality:** Format validation and sanity checks
4. **Temporal Consistency:** Consistency with historical patterns

```
def select_best_extraction(lxml_result, bs_result):
    """Select optimal extraction result"""
    lxml_score = calculate_completeness_score(lxml_result)
    bs_score = calculate_completeness_score(bs_result)

    if abs(lxml_score - bs_score) < 0.1:
        # Scores similar, use cross-validation
        return cross_validate_results(lxml_result, bs_result)
    else:
        # Clear winner based on completeness
        return lxml_result if lxml_score > bs_score else bs_result
```

5.3 Extracted Data Schema

```
extracted_data = {
    'app_title': str,          # Application name
    'rating': float,          # Average user rating (1-5)
    'downloads': str,         # Download count range
    'reviews': int,           # Number of user reviews
    'description': str,       # Full application description
    'whats_new': str,         # Recent updates and changes
    'timestamp': str,         # Snapshot timestamp
    'version': str,           # App version number
    'last_updated': str,      # Last update date
    'size': str,              # Application size
    'content_rating': str     # Age rating classification
}
```

6. Phase 5: Data Cleaning and Standardization

6.1 Data Cleaning Pipeline

Implementation: `fintech_data_cleaner_standardizer.py`, `advanced_fintech_data_cleaner.py`

6.1.1 Rating Normalization

```
def clean_rating(rating_str):
    """Standardize rating formats"""
    if pd.isna(rating_str):
        return None

    # Extract numeric rating from various formats
    rating_patterns = [
        r'(\d+\.? \d*) out of 5',
        r'(\d+\.? \d*) stars',
        r'Reated (\d+\.? \d*)',
        r'^(\d+\.? \d*)$'
    ]

    for pattern in rating_patterns:
        match = re.search(pattern, str(rating_str))
        if match:
            return float(match.group(1))

    return None
```

6.1.2 Download Count Processing

```
def normalize_downloads(download_str):  
    """Convert download ranges to numerical values"""  
    if pd.isna(download_str):  
        return None  
  
    download_mappings = {  
        '1,000+': 1000,  
        '5,000+': 5000,  
        '10,000+': 10000,  
        '50,000+': 50000,  
        '100,000+': 100000,  
        '500,000+': 500000,  
        '1,000,000+': 1000000,  
        '5,000,000+': 5000000,  
        '10,000,000+': 10000000,  
        '50,000,000+': 50000000,  
        '100,000,000+': 100000000,  
        '500,000,000+': 500000000,  
        '1,000,000,000+': 1000000000  
    }  
  
    return download_mappings.get(download_str, None)
```

6.2 Temporal Data Processing

Objective: Create proper time-series structure

Implementation:

1. **Date Standardization:** Convert various date formats to ISO standard
2. **Temporal Sorting:** Chronological organization of data points
3. **Gap Identification:** Detection of missing time periods
4. **Interpolation Strategy:** Handling of missing data points

```
def create_timeseries(app_data):
    """Transform app data into time-series format"""
    # Sort by timestamp
    app_data['timestamp'] = pd.to_datetime(app_data['timestamp'])
    app_data = app_data.sort_values('timestamp')

    # Create temporal features
    app_data['year'] = app_data['timestamp'].dt.year
    app_data['month'] = app_data['timestamp'].dt.month
    app_data['quarter'] = app_data['timestamp'].dt.quarter

    # Calculate temporal deltas
    app_data['days_since_launch'] = (app_data['timestamp'] -
                                     app_data['timestamp'].iloc[0]).dt.days

    return app_data
```

7. Data Integration and Final Dataset Construction

7.1 Multi-Source Integration

Objective: Combine all data sources into unified dataset

Process:

1. **App-Company Mapping Integration:** Link app data with company information
2. **Temporal Data Merging:** Combine time-series data across applications
3. **Metadata Enhancement:** Add derived features and calculations
4. **Quality Validation:** Final data quality checks and validation

7.2 Final Dataset Structure

Output: `Fintech_App_History.csv`

Schema:

```
app_launch_date, app_title, company_name, domicile, file, timestamp,
rating, downloads, reviews, description, whats_new, Date,
normalized_downloads, rating_numeric, days_since_launch,
year, month, quarter
```

Characteristics:

- **Size:** 11MB+ of processed data
- **Coverage:** 200+ fintech applications
- **Temporal Span:** Multi-year historical coverage
- **Data Points:** Thousands of temporal observations

8. Validation and Quality Assurance

8.1 Data Quality Metrics

1. **Completeness:** Percentage of non-null values per field
2. **Consistency:** Cross-validation between different extraction methods
3. **Accuracy:** Manual verification of sample records
4. **Temporal Consistency:** Logical progression of metrics over time

8.2 Validation Strategies

- **Cross-Parser Validation:** Compare lxml vs BeautifulSoup results
- **Temporal Logic Checks:** Ensure ratings don't decrease dramatically
- **External Validation:** Cross-reference with known app information
- **Statistical Outlier Detection:** Identify and investigate anomalies

8.3 Error Analysis

- **Missing Data Patterns:** Analysis of systematic gaps
- **Extraction Failures:** Investigation of parser limitations
- **Temporal Inconsistencies:** Identification of problematic time periods
- **Quality Scoring:** Per-app and per-field quality metrics

9. Ethical Considerations and Compliance

9.1 Responsible Scraping Practices

- **Rate Limiting:** Respectful request frequency
- **Robot.txt Compliance:** Adherence to site policies
- **Terms of Service:** Compliance with platform guidelines
- **Data Attribution:** Proper acknowledgment of data sources

9.2 Privacy and Data Protection

- **No Personal Data:** Exclusion of user-identifiable information
- **Aggregated Metrics:** Focus on aggregate statistics only
- **Public Data:** Limitation to publicly available information
- **Research Use:** Academic and research-focused application

10. Technical Achievements and Innovations

10.1 Novel Methodological Contributions

1. **Hybrid Matching Algorithm:** Combination of fuzzy and semantic matching
2. **Best-of-Both Extraction:** Dual-parser approach with intelligent selection
3. **Adaptive Scraping:** Dynamic rate adjustment and proxy rotation
4. **Temporal Data Integration:** Historical data reconstruction methodology

10.2 Scalability Solutions

- **Modular Architecture:** Component-based design for maintainability
- **Batch Processing:** Efficient handling of large datasets
- **Memory Optimization:** Streaming data processing approaches
- **Error Recovery:** Robust failure handling and resumption

10.3 Performance Optimizations

- **C-based Parsing:** lxml for high-performance HTML processing
- **Tensor Computations:** GPU-accelerated similarity calculations
- **Caching Strategies:** Intelligent data caching for repeated operations
- **Parallel Processing:** Multi-threaded data processing where applicable

11. Limitations and Future Work

11.1 Current Limitations

- **Temporal Gaps:** Some historical periods may have limited data
- **Platform Dependency:** Focused primarily on Google Play Store
- **Manual Validation:** Some matching requires human verification

- **Dynamic Content:** Limited handling of JavaScript-generated content

11.2 Future Enhancement Opportunities

1. **Real-time Pipeline:** Live data collection and processing
2. **Multi-platform Support:** Integration with iOS App Store data
3. **Advanced NLP:** Integration of larger language models for matching
4. **Automated Validation:** Machine learning-based quality assessment

12. Conclusion

This methodology represents a comprehensive approach to fintech app data collection and analysis, combining advanced web scraping techniques, machine learning algorithms, and rigorous data validation processes. The resulting dataset provides unprecedented insight into the evolution of the Indian fintech app ecosystem and serves as a valuable resource for academic research and industry analysis.

The technical innovations introduced, particularly the hybrid matching algorithms and dual-parser extraction approach, contribute to the broader field of web data collection and processing. The emphasis on data quality, ethical scraping practices, and reproducible methodology ensures the research meets high academic standards while providing practical value for understanding fintech market dynamics.

Author: ISB Fintech Research Team

Institution: Indian School of Business (ISB)

Project Duration: [Research Period]

Last Updated: [Current Date]