

Australian Company Data Pipeline Documentation

Susanta Baidya

March 2, 2025

Abstract

This document describes the implementation of a data pipeline that integrates Australian company website data from Common Crawl with business information from the Australian Business Register (ABR). The system collects, cleans, and stores data in a MySQL database, enabling comprehensive business analysis.

1 Pipeline Architecture

The data processing workflow follows three key stages:

- **Data Ingestion:** Python scripts extract raw data from Common Crawl (200k+ URLs) and ABR
- **Transformation:** DBT models clean, normalize, and merge datasets
- **Storage:** Processed data stored in MySQL with optimized schema
- **Validation:** Automated data quality checks using DBT tests

2 Implementation Details

2.1 Data Collection

```
1 # Common Crawl URL processing
2 def process_url(url):
3     try:
4         response = requests.get(url, timeout=10)
5         soup = BeautifulSoup(response.text, 'html.parser')
6         return extract_metadata(soup)
7     except Exception as e:
8         log_error(url, str(e))
9         return None
```

2.2 Data Cleaning

DBT staging model for ABR data:

```
1 {{ config(materialized='view') }}
2
3 SELECT
4     abn,
5     TRIM(company_name) AS company_name,
6     entity_type,
7     state,
8     postcode,
9     registration_date
10 FROM {{ source('raw', 'abr_raw') }}
11 WHERE LENGTH(abn) = 11
```

3 Database Schema

Optimized MySQL table structure:

```
1 CREATE TABLE companies (  
2     company_id INT AUTO_INCREMENT PRIMARY KEY,  
3     url VARCHAR(512) NOT NULL UNIQUE,  
4     company_name VARCHAR(255) NOT NULL,  
5     industry VARCHAR(255),  
6     abn CHAR(11),  
7     entity_type VARCHAR(50),  
8     state CHAR(3),  
9     postcode CHAR(4),  
10    registration_date DATE  
11 );
```

4 Data Validation

Implemented quality checks:

- **Uniqueness:** Enforced through database constraints
- **Completeness:** NULL checks for mandatory fields
- **Consistency:** ABN format validation (11 digits)
- **Referential Integrity:** Relationship between tables

5 Example Queries

Top industries analysis:

```
1 SELECT industry, COUNT(*) AS company_count  
2 FROM companies  
3 GROUP BY industry  
4 ORDER BY company_count DESC  
5 LIMIT 10;
```

Recent registrations:

```
1 SELECT company_name, registration_date  
2 FROM companies  
3 WHERE registration_date > '2023-01-01'  
4 ORDER BY registration_date DESC;
```

References

- Common Crawl: <https://commoncrawl.org>
- ABR Dataset: <https://data.gov.au>
- DBT Documentation: <https://docs.getdbt.com>
- MySQL Documentation: <https://dev.mysql.com/doc>