# Capstone Project 2
## FakeNewsNet Database Design and Implementation

Siu Chun Anson Chan

62006049

May 10, 2025

# Contents

# Chapter 1

# Introduction

This report presents the data modelling and database implementation for the FakeNews-Net dataset based on the dataset collected in Capstone Project 1. The project involves designing a comprehensive database schema to efficiently store and query fake news data, including news articles, social media interactions, and user information.

# Chapter 2

# Step 1: Understanding the Dataset

## 2.1 Dataset Summary

The FakeNewsNet dataset consists of news articles and their associated social media data collected from two fact-checking sources: PolitiFact and GossipCop. The dataset structure includes:

### 2.1.1 Main Data Components

- **News Articles**: Contains news content with unique IDs, URLs, titles, and labels (fake/real)
- **News Sources**: PolitiFact and GossipCop as fact-checking sources
- **Tweets**: Social media posts sharing the news articles
- **Users**: Twitter users who posted tweets about the news
- **Retweets**: Reshares of original tweets
- **User Relationships**: Follower/following connections
- **User Timelines**: Historical tweets from users
- **News Content**: Full article text, images, and publication dates

### 2.1.2 Data Attributes

- News: id, url, title, text, publish_date, label (fake/real)
- Tweet: tweet_id, user_id, content, timestamp
- User: user_id, username, profile_data, follower_count
- Relationships: Multi-dimensional social network data

### 2.1.3  Database Type Selection

A **Relational Database Management System (RDBMS)** is chosen for this project because:

- The data has clear relationships (news-tweets, tweets-users, users-followers)

- Need for ACID compliance and data integrity

- Complex queries involving multiple joins

- Well-defined schema with structured data

# Chapter 3

# Step 2: Conceptual Data Modelling

## 3.1   Entity Identification

Based on the dataset analysis, we identify the following key entities:

1. **NewsSource**: Fact-checking sources (PolitiFact, GossipCop)

2. **NewsArticle**: Individual news articles

3. **NewsContent**: Detailed content of news articles

4. **NewsImage**: Images associated with news articles

5. **Tweet**: Social media posts about news

6. **User**: Twitter users

7. **Retweet**: Retweet relationships

8. **UserFollower**: Follower relationships between users

9. **UserTimeline**: Historical tweets from users

10. **NewsCategory**: Categories/topics of news articles

## 3.2   Relationships

- NewsSource (1) — publishes — (M) NewsArticle

- NewsArticle (1) — has — (1) NewsContent

- NewsArticle (1) — contains — (M) NewsImage

- NewsArticle (1) — shared_by — (M) Tweet

- User (1) — posts — (M) Tweet

- Tweet (1) — has — (M) Retweet

- User (1) — retweets — (M) Retweet

- User (M) — follows — (M) User (UserFollower)

- User (1) — has — (M) UserTimeline

- NewsArticle (M) — belongs_to — (M) NewsCategory

## 3.3   Model Suitability

This conceptual model is suitable for the FakeNewsNet case study because:

- It captures all dimensions of fake news data (content, social context, temporal aspects)

- Supports tracking of information dissemination through social networks

- Enables analysis of user behaviour patterns

- Maintains data integrity through proper relationship definitions

- Allows for complex queries about news credibility and social impact

# Chapter 4

# Step 3: Logical Data Modelling

## 4.1 Entity Attributes and Data Types

### 4.1.1 NewsSource

- source_id (INT) - Primary Key

- source_name (VARCHAR(50)) - Not Null, Unique

- source_url (VARCHAR(255))

- credibility_rating (DECIMAL(3,2))

### 4.1.2 NewsArticle

- article_id (VARCHAR(50)) - Primary Key

- source_id (INT) - Foreign Key

- url (VARCHAR(500)) - Not Null

- title (VARCHAR(500)) - Not Null

- label (ENUM('fake', 'real')) - Not Null

- created_at (TIMESTAMP)

### 4.1.3 NewsContent

- content_id (INT) - Primary Key, Auto Increment

- article_id (VARCHAR(50)) - Foreign Key, Unique

- text (TEXT)

- publish_date (DATETIME)

- author (VARCHAR(255))

- word_count (INT)

## 4.2 Normalisation to 3NF

The logical model is normalised to Third Normal Form (3NF):

1. **1NF**: All attributes contain atomic values, no repeating groups

2. **2NF**: All non-key attributes are fully functionally dependent on the primary key

3. **3NF**: No transitive dependencies exist

# Chapter 5

# Step 4: Physical Data Modelling & Database Schema

## 5.1 DBMS Selection

PostgreSQL is selected as the DBMS for the following reasons:

- Excellent support for complex queries and joins

- JSONB data type for semi-structured data

- Full-text search capabilities

- Open-source with strong community support

- Scalability for large datasets

## 5.2 Physical ERD Design

The physical ERD extends the logical model with implementation-specific details for PostgreSQL:

- **Storage Parameters**: HEAP storage with optimised FILLFACTOR settings

- **Indexing Strategy**: B-tree indexes for primary keys, GIN indexes for full-text search

- **Partitioning**: Date-based partitioning for high-volume tables (tweet, retweet)

- **Performance Optimisation**: Composite indexes for common query patterns

- **Constraint Implementation**: CHECK constraints, foreign key constraints, and unique constraints

Key physical design decisions:

- Tweet table partitioned by creation date for scalability

- Full-text search indexes on article titles and content

- Optimised follower relationship queries through strategic indexing

- Memory-efficient storage for large text fields using TOAST

## 5.3 SQL Scripts for Table Creation

Listing 5.1: Database and Initial Setup

```sql
-- Create database
CREATE DATABASE fakenewsnet_db;

-- Create ENUM type for news labels
CREATE TYPE news_label AS ENUM ('fake', 'real');
```

Listing 5.2: NewsSource Table

```sql
CREATE TABLE news_source (
    source_id SERIAL PRIMARY KEY,
    source_name VARCHAR(50) NOT NULL UNIQUE,
    source_url VARCHAR(255),
    credibility_rating DECIMAL(3,2)
        CHECK (credibility_rating >= 0 AND credibility_rating <=
            1)
);
```

Listing 5.3: NewsArticle Table

```sql
CREATE TABLE news_article (
    article_id VARCHAR(50) PRIMARY KEY,
    source_id INTEGER NOT NULL,
    url VARCHAR(500) NOT NULL,
    title VARCHAR(500) NOT NULL,
    label news_label NOT NULL,
    created_at TIMESTAMP DEFAULT CURRENT_TIMESTAMP,
    FOREIGN KEY (source_id) REFERENCES news_source(source_id)
);

-- Create indexes for faster queries
CREATE INDEX idx_news_article_source ON news_article(source_id);
CREATE INDEX idx_news_article_label ON news_article(label);
```

# Chapter 6

# Step 5: Querying the Database

## 6.1 Easy Queries

### 6.1.1 Query 1: Total Tweets by News Label

Listing 6.1: Easy Query 1

```sql
-- Get the total number of tweets for fake vs real news articles
SELECT
    na.label,
    COUNT(t.tweet_id) as total_tweets,
    AVG(t.retweet_count) as avg_retweets
FROM news_article na
INNER JOIN tweet t ON na.article_id = t.article_id
GROUP BY na.label
ORDER BY total_tweets DESC;
```

**Step-by-step Explanation**:

1. **SELECT clause**: Selects the news label, counts total tweets, and calculates average retweets

2. **FROM clause**: Starts with the news_article table as the main data source

3. **INNER JOIN**: Joins with the tweet table to connect articles with their associated tweets

4. **GROUP BY**: Groups results by news label (fake/real) to aggregate data for each category

5. **ORDER BY**: Sorts results by total tweet count in descending order

This query helps understand engagement patterns between different news types by analysing tweet volume and retweet behaviour.

### 6.1.2 Query 2: User Activity Summary

```
1  -- Get user activity summary showing tweet counts
2  SELECT
3      u.username,
4      u.verified,
5      COUNT(t.tweet_id) as tweet_count,
6      SUM(t.retweet_count) as total_retweets_received
7  FROM users u
8  INNER JOIN tweet t ON u.user_id = t.user_id
9  WHERE u.followers_count > 1000
10 GROUP BY u.user_id, u.username, u.verified
11 ORDER BY tweet_count DESC
12 LIMIT 10;
```

**Step-by-step Explanation**:

1. **SELECT clause**: Selects username, verification status, tweet count, and total retweets received

2. **FROM clause**: Starts with the users table to get user information

3. **INNER JOIN**: Joins with tweet table to connect users with their tweets

4. **WHERE clause**: Filters to include only users with more than 1000 followers

5. **GROUP BY**: Groups by user attributes to aggregate tweet statistics per user

6. **ORDER BY**: Sorts by tweet count in descending order to find most active users

7. **LIMIT**: Restricts results to top 10 most active users

This query identifies influential users based on follower count and activity level.

## 6.2 Medium Queries

### 6.2.1 Query 3: News Source Performance Analysis

Listing 6.3: Medium Query 1

```
1  -- Analyze news source performance with category breakdown
2  SELECT
3      ns.source_name,
4      nc.category_name,
5      na.label,
6      COUNT(DISTINCT na.article_id) as article_count,
7      COUNT(DISTINCT t.tweet_id) as total_tweets,
8      UPPER(SUBSTRING(na.title, 1, 50)) || '...' as sample_title
9  FROM news_source ns
10 INNER JOIN news_article na ON ns.source_id = na.source_id
11 INNER JOIN article_category ac ON na.article_id = ac.article_id
12 INNER JOIN news_category nc ON ac.category_id = nc.category_id
13 LEFT JOIN tweet t ON na.article_id = t.article_id
```

```
14  WHERE na.created_at >= CURRENT_DATE - INTERVAL '30 days'
15  GROUP BY ns.source_name, nc.category_name, na.label, na.title
16  ORDER BY total_tweets DESC;
```

**Step-by-step Explanation**:

1. **SELECT clause**: Selects source name, category, label, article count, tweet count, and sample title

2. **FROM clause**: Starts with news_source table as the primary data source

3. **INNER JOIN news_article**: Connects sources with their published articles

4. **INNER JOIN article_category**: Links articles to their categories through junction table

5. **INNER JOIN news_category**: Connects to category names and details

6. **LEFT JOIN tweet**: Includes tweet data where available (left join allows articles without tweets)

7. **WHERE clause**: Filters to articles created within the last 30 days

8. **GROUP BY**: Groups by source, category, label, and title to aggregate statistics

9. **ORDER BY**: Sorts by total tweets in descending order to show most engaging content first

This query provides comprehensive performance analysis of news sources across different categories and time periods.

## 6.3    Difficult Query

### 6.3.1    Query 5: Comprehensive Fake News Network Analysis

Listing 6.4: Difficult Query

```
1   -- Complex analysis of fake news dissemination networks
2   WITH user_influence AS (
3       SELECT
4           u.user_id,
5           u.username,
6           u.verified,
7           COUNT(DISTINCT t.tweet_id) as tweet_count,
8           RANK() OVER (ORDER BY u.followers_count *
9                       COUNT(DISTINCT t.tweet_id) DESC) as
                            influence_rank
10      FROM users u
11      INNER JOIN tweet t ON u.user_id = t.user_id
12      GROUP BY u.user_id, u.username, u.verified, u.followers_count
13  )
14  SELECT
15      ui.username as top_influencer,
```

14

```
16      ui.verified as influencer_verified ,
17      ui.influence_rank ,
18      na.title ,
19      na.label ,
20      COUNT(DISTINCT t.tweet_id) as tweet_count
21  FROM user_influence ui
22  INNER JOIN tweet t ON ui.user_id = t.user_id
23  INNER JOIN news_article na ON t.article_id = na.article_id
24  WHERE ui.influence_rank <= 100
25    AND na.label = 'fake'
26  GROUP BY ui.username , ui.verified , ui.influence_rank , na.title ,
       na.label
27  ORDER BY ui.influence_rank
28  LIMIT 25;
```

**Step-by-step Explanation**:

1. **CTE (WITH clause)**: Creates a Common Table Expression called user_influence to pre-calculate user influence metrics

2. **CTE SELECT**: Calculates user influence based on followers count multiplied by tweet activity

3. **RANK() function**: Assigns influence rankings using a window function ordered by the influence metric

4. **Main SELECT**: Selects top influencer details, article information, and engagement metrics

5. **Multiple JOINs**: Connects user influence data with tweets and news articles

6. **WHERE clause**: Filters to top 100 influencers AND only fake news articles

7. **GROUP BY**: Groups results by user and article to aggregate tweet counts

8. **ORDER BY**: Sorts by influence rank to show most influential users first

9. **LIMIT**: Restricts to top 25 results for manageable output

This complex query identifies the most influential users spreading fake news by combining user metrics, social network analysis, and content classification. It demonstrates advanced SQL techniques including CTEs, window functions, and multi-table joins.

# Chapter 7

# Step 6: Database Implementation and Data Insertion

## 7.1 Sample Data Implementation

The database implementation includes comprehensive sample data insertion demonstrating:

- **News Sources**: PolitiFact, GossipCop, and Snopes with credibility ratings

- **News Categories**: Politics, Health, Technology, and Entertainment classifications

- **News Articles**: Mixed fake and real articles with proper labelling

- **User Profiles**: Verified and unverified users with realistic follower counts

- **Social Interactions**: Tweets, retweets, and follower relationships

- **Content Relationships**: Article-category associations and user timelines

## 7.2 Data Verification and Testing

Sample queries verify data integrity and relationship correctness:

- Foreign key constraint validation across all tables

- Data type conformance and constraint adherence

- Relationship cardinality verification

- Index performance testing for common query patterns

The implementation demonstrates successful data insertion with proper referential integrity maintenance and constraint enforcement throughout the database schema.

# Chapter 8

# Conclusion

This project successfully designed and implemented a comprehensive database schema for the FakeNewsNet dataset. The complete deliverables include:

## 8.1 Database Design Accomplishments

- **Conceptual ERD**: 10 entities with proper relationships and cardinalities

- **Logical ERD**: Normalised to 3NF with complete attribute specifications

- **Physical ERD**: PostgreSQL-specific implementation with performance optimisations

- **SQL Implementation**: Complete table creation scripts with constraints and indexes

## 8.2 Query Implementation

The project delivers 8 comprehensive SQL queries including:

- 2 Easy queries demonstrating basic aggregation and filtering

- 2 Medium queries showcasing complex joins and analytical functions

- 1 Difficult query utilising CTEs, window functions, and advanced SQL features

- 3 Additional queries for comprehensive database demonstration

## 8.3 Technical Achievements

The database design:

- Efficiently stores multi-dimensional fake news data across 11 related tables

- Supports complex queries for analysing news dissemination patterns

- Maintains referential integrity through comprehensive constraint implementation

- Provides optimised access patterns through strategic indexing

- Enables advanced analytics on social media engagement and user behaviour

- Implements physical optimisations including partitioning and full-text search

## 8.4  Project Impact

The implemented database system demonstrates capability to:

- Track fake news propagation through social networks

- Analyse user influence patterns and verification status

- Support real-time monitoring of news article engagement

- Enable research into misinformation spread dynamics

- Facilitate fact-checking organisation collaboration

The project provides a robust foundation for fake news research and detection systems, with scalable architecture supporting both analytical and operational workloads.