# Data-Driven Exploration of Chronic Kidney Disease

## Introduction

The dataset for this study is sourced from Kaggle, a leading platform for data science and machine learning resources. The "Chronic Kidney Disease" dataset is hosted on [Kaggle](), a premier platform for open data and machine learning resources. Kaggle is renowned for its extensive repository of datasets, community-driven support, and robust tools for data analysis. The dataset is available for free under Kaggle's licensing terms, which allow for non-commercial academic and research purposes. The dataset is easily downloadable in CSV format and accessible through a straightforward interface that supports data exploration through discussions, kernels (code notebooks), and user reviews. This specific dataset in question is titled **"Chronic Kidney Disease"** and can be accessed through the Kaggle website: [CKD Dataset]() Chronic Kidney Disease (CKD) is a significant global health issue, impacting millions of lives worldwide. Early detection and treatment are crucial to improve patient outcomes and reduce the burden on healthcare systems. This study utilizes the "Chronic Kidney Disease" dataset from Kaggle, which includes comprehensive medical records of patients diagnosed with CKD. The objective is to explore the various dimensions of CKD by developing and implementing a MySQL database and a Node.js application. This project aims to analyze the data for a better understanding of CKD epidemiology, risk factors, and potential predictive modelling, thereby contributing to more effective management strategies and improved patient care.

This report presents an in-depth analysis of a publicly available dataset from Kaggle and offers an interactive platform for detailed exploration of CKD data. The significance of this analysis lies in its ability to uncover patterns and correlations within the dataset, which can lead to more accurate predictions of CKD and better management of the disease.

The study is structured around the following key areas:

- **Dataset Analysis**: Focuses on the "Chronic Kidney Disease" dataset, which includes various medical attributes such as patient demographics, clinical history, and diagnostic results.
- **Database Design**: Involves the creation of an Entity-Relationship (E/R) model to represent the dataset's structure, ensuring efficient data retrieval and support for complex queries.

- **Application Development**: Entails the implementation of a Node.js application to facilitate data interaction, providing functionalities for data retrieval, manipulation, and visualization.
- **Research Questions**: Aims to investigate demographic factors associated with CKD, the impact of various medical attributes on CKD diagnosis, and the potential for predictive modelling.
- **Performance Evaluation**: Assesses the performance of the MySQL database and the Node.js application in handling the dataset and supporting research activities.

**Strengths:**

- **Free Access:** No subscription or purchase is required, making it accessible for educational and non-commercial research.
- **Community Features:** Active community discussions and shared code notebooks that facilitate learning and problem-solving.
- **User-Friendly Interface:** Simple download options and an easy-to-navigate dataset page.

**Weaknesses:**

- **Commercial Use Limitations:** The dataset is not available for commercial use, which may restrict applications in for-profit ventures or product development.

By defining these parameters, the study ensures a comprehensive approach to data exploration, focusing on key areas that contribute to a deeper understanding of CKD.

## Dataset Selection and Critique

### Dataset Source and Accessibility

The "Chronic Kidney Disease" dataset is obtained from Kaggle, a renowned platform for data science and machine learning resources. Kaggle provides an intuitive interface for dataset downloads and supports community interaction through discussions, code snippets, and user reviews. The dataset is freely available for non-commercial purposes, making it highly accessible for academic research. Kaggle's user-friendly design and robust community support enhance the dataset's accessibility.

### Quality and Detail

The dataset is of high quality, featuring extensive information on various aspects of CKD, including patient demographics, medical history, and diagnostic results. The data is well-organized and clean, though it contains some missing values that need to be addressed during data preprocessing. The detailed attributes provided in the dataset enable in-depth analyses across multiple dimensions of CKD. The reliability of the dataset assumes thorough

and systematic data collection, although specific details on the collection methods are not explicitly provided.

The dataset provides a comprehensive collection of medical and demographic information relevant to Chronic Kidney Disease (CKD). It includes a variety of attributes such as age, blood pressure, specific gravity, and serum levels that are instrumental for CKD analysis.

**Quality Analysis:**

- **Completeness:** The dataset is relatively complete but does contain missing values, which may affect the accuracy of analyses. Strategies for handling missing data include imputation, deletion, or substitution based on domain knowledge.
- **Consistency:** The dataset features consistent attribute definitions, though data entry errors or outliers might be present. For instance, inconsistencies in data units or missing values for certain attributes could be scrutinized.
- **Accuracy:** The dataset's accuracy relies on the credibility of the original data sources, which are not explicitly detailed in the dataset description. Validation against external sources or comparison with similar datasets might help assess accuracy.

**Detail Analysis:**

- **Attribute Range:** The dataset contains diverse attributes, including demographic details (age, gender), medical history (blood pressure, diabetes), and diagnostic results (haemoglobin, urine glucose). This range allows for detailed explorations into CKD-related factors.
- **Attribute Documentation:** Clear explanations for each attribute are provided, though detailed methodologies for data collection are absent. Improved documentation might include information about data sources, collection methods, and validation processes.

**Documentation and Use**

The dataset includes clear and detailed documentation, explaining the meaning of each attribute and the overall data structure. This documentation is essential for understanding and correctly interpreting the data and is easily accessible on the Kaggle dataset page. The dataset can be utilised for various research purposes, such as studying CKD prevalence, evaluating diagnostic criteria, and developing predictive models. However, certain uses may be limited by the granularity of the data or missing values in some attributes.

The dataset includes basic documentation that outlines each attribute's meaning and purpose. This documentation is essential for understanding the data structure and ensuring correct analysis.

**Strengths:**

- **Clear Attribute Descriptions:** The dataset documentation explains the significance of each attribute.

- **Usage Examples:** Kaggle's platform includes examples of how the data has been used in previous projects, which can guide new analyses.

**Weaknesses:**

- **Limited Details:** The dataset lacks comprehensive methodological details, which limits understanding of the data collection processes and potential biases.

**Potential Uses:**

- **Research and Analysis:** The dataset is suitable for studying CKD's epidemiology, risk factors, and diagnostic criteria.
- **Predictive Modeling:** It supports the development of models to predict CKD based on medical and demographic variables.

**Interrelation and Discoverability**

Integrating the CKD dataset with other data sources could yield more comprehensive insights.

**Potential Integrations:**

- **Environmental Data:** Data on air and water quality could reveal environmental factors influencing CKD risk.
- **Socioeconomic Data:** Information on income, education, and occupation could help explore socioeconomic factors affecting CKD prevalence.
- **Health Records:** Combining CKD data with records of other chronic diseases could uncover comorbidity patterns.

**Discoverability:**

- **Kaggle Community:** Kaggle facilitates finding related datasets through user suggestions and shared resources.
- **External Repositories:** Government health databases and scientific research repositories offer additional datasets for broader research.

The dataset could be enhanced by integrating it with other related datasets, such as environmental data (e.g., air and water quality indices), socioeconomic data, or health records from other chronic diseases. Such integration could provide a more comprehensive understanding of CKD triggers and the socioeconomic factors affecting CKD management. Integrating related datasets can be facilitated by using Kaggle, which offers a wide range of data sources. Other potential sources include government health databases and scientific research repositories, which provide alternative datasets for broader analyses.

**Interest and Research Questions**

This dataset is particularly intriguing due to its potential to reveal valuable insights into CKD—a prevalent and impactful health issue. Analysis of this dataset can uncover patterns and correlations that inform better treatment strategies, policy formulation, and public health initiatives. Some SQL-based research questions that can be explored with this dataset include:

- What demographic factors are most strongly associated with CKD?
- What are the common medical attributes associated with CKD diagnosis?
- Can predictive models be developed to forecast CKD based on available data?
- What is the geographical distribution of CKD cases, and how does it relate to environmental factors?

A database application would be useful to answer these questions by providing tools for complex queries, data visualization, and statistical analysis. This approach aims to yield valuable insights that can inform better management and intervention strategies for CKD patients.

The dataset's potential to provide insights into CKD makes it highly valuable for various research questions and database applications.

**Research Questions:**

1. **Demographic Factors:** What demographic characteristics (age, gender, occupation) are most strongly associated with CKD?
   - *SQL Query Example:* `SELECT age, gender, COUNT(*) FROM ckd_data GROUP BY age, gender HAVING CKD = 'positive';`
2. **Medical Attributes:** What medical attributes (blood pressure, serum levels) are the most significant predictors of CKD?
   - *SQL Query Example:* `SELECT attribute, AVG(value) FROM ckd_data WHERE CKD = 'positive' GROUP BY attribute;`
3. **Predictive Modeling:** Can we build a predictive model for CKD diagnosis based on the dataset's attributes?
   - *SQL Query Example:* `SELECT * FROM ckd_data WHERE CKD = 'positive' OR CKD = 'negative';`
   - This question would involve data preprocessing and applying machine learning algorithms.
4. **Geographical Distribution:** What is the geographical distribution of CKD cases, and how might it relate to environmental and socioeconomic factors?
   - *SQL Query Example:* `SELECT location, COUNT(*) FROM ckd_data GROUP BY location HAVING CKD = 'positive';`

**Application of a Database:** A well-structured MySQL database could handle complex queries, perform data manipulation, and generate visualizations. A Node.js application could be developed for interactive data exploration and visualization, helping researchers and practitioners access insights from the data.

**Terms of Use and Licensing**

The dataset is available under the terms specified by Kaggle, which allow for non-commercial use with proper attribution to the original contributor, Mansoor Iqbal. The terms do not significantly restrict academic and research use, making the dataset suitable for this project. Licensing and rights information is provided on the Kaggle dataset page, ensuring that users are aware of the usage conditions and can adhere to them.

By leveraging this dataset, the study aims to contribute to a deeper understanding of CKD and its various dimensions, ultimately supporting improved patient outcomes and healthcare strategies.

The dataset is available under Kaggle's terms of use, which allow non-commercial academic research and require attribution to the dataset creator, Mansoor Iqbal.

**Strengths:**

- **Academic and Research Use:** Suitable for educational and research purposes.
- **Clear Licensing:** Licensing terms are clearly stated, ensuring compliance with usage guidelines.

**Weaknesses:**

- **Commercial Restrictions:** Limitations on commercial use could hinder applications in for-profit scenarios.

**Entity-Relationship (E/R) Diagram**

The E/R diagram represents the structure of the dataset by identifying entities and their relationships. The list of entities and their respective attributes are shown below :

| Patients | |
|---|---|
| patient_id 🔑 | INT |
| age | INT |
| gender | VARCHAR(10) |
| blood_pressure | INT |
| specific_gravity | FLOAT |
| albumin | INT |
| sugar | INT |
| red_blood_cells | VARCHAR(20) |
| pus_cell | VARCHAR(20) |
| pus_cell_clumps | VARCHAR(20) |
| bacteria | VARCHAR(20) |
| blood_glucose_random | FLOAT |
| blood_urea | FLOAT |
| serum_creatinine | FLOAT |
| sodium | FLOAT |
| potassium | FLOAT |
| hemoglobin | FLOAT |
| packed_cell_volume | VARCHAR(10) |
| white_blood_cell_count | INT |
| red_blood_cell_count | FLOAT |
| hypertension | VARCHAR(10) |
| diabetes_mellitus | VARCHAR(10) |
| coronary_artery_disease | VARCHAR(10) |
| appetite | VARCHAR(10) |
| pedal_edema | VARCHAR(10) |
| anemia | VARCHAR(10) |
| classification | VARCHAR(10) |

| Medical_History | |
|---|---|
| history_id 🔑 | INT |
| patient_id | INT |
| hypertension | VARCHAR(10) |
| diabetes_mellitus | VARCHAR(10) |
| coronary_artery_disease | VARCHAR(10) |
| anemia | VARCHAR(10) |
| date_recorded | DATE |

| Diagnostic_Results | |
|---|---|
| result_id 🔑 | INT |
| patient_id | INT |
| blood_pressure | INT |
| specific_gravity | FLOAT |
| albumin | INT |
| sugar | INT |
| red_blood_cells | VARCHAR(20) |
| pus_cell | VARCHAR(20) |
| pus_cell_clumps | VARCHAR(20) |
| bacteria | VARCHAR(20) |
| blood_glucose_random | FLOAT |
| blood_urea | FLOAT |
| serum_creatinine | FLOAT |
| sodium | FLOAT |
| potassium | FLOAT |
| hemoglobin | FLOAT |
| packed_cell_volume | VARCHAR(10) |
| white_blood_cell_count | INT |
| red_blood_cell_count | FLOAT |
| date_recorded | DATE |

| Treatment_Plans | |
|---|---|
| treatment_id 🔑 | INT |
| patient_id | INT |
| medication | VARCHAR(255) |
| dosage | VARCHAR(100) |
| frequency | VARCHAR(50) |
| start_date | DATE |
| end_date | DATE |

## Entities and Attributes

- **Patients**
  - PatientID (Primary Key)
  - Name
  - Age
  - Gender
  - Address
- **Demographics**
  - DemographicID (Primary Key)
  - PatientID (Foreign Key)
  - Ethnicity
  - EducationLevel
  - Occupation
- **MedicalHistory**
  - MedicalHistoryID (Primary Key)
  - PatientID (Foreign Key)
  - Hypertension
  - Diabetes
  - CoronaryArteryDisease

- ○ Anemia
- **ClinicalMeasurements**
  - ○ ClinicalMeasurementID (Primary Key)
  - ○ PatientID (Foreign Key)
  - ○ BloodPressure
  - ○ SpecificGravity
  - ○ Albumin
  - ○ Sugar
  - ○ BloodGlucoseRandom
  - ○ BloodUrea
  - ○ SerumCreatinine
  - ○ Sodium
  - ○ Potassium
  - ○ Hemoglobin
  - ○ PackedCellVolume
  - ○ WhiteBloodCellCount
  - ○ RedBloodCellCount
- **Symptoms**
  - ○ SymptomID (Primary Key)
  - ○ PatientID (Foreign Key)
  - ○ ShortnessOfBreath
  - ○ SwellingOfFeet
  - ○ Fatigue
  - ○ Nausea
- **Diagnosis**
  - ○ DiagnosisID (Primary Key)
  - ○ PatientID (Foreign Key)
  - ○ DiagnosisDate
  - ○ Stage
  - ○ TreatmentPlan

**Relationships**

- Demographics to Patients: One-to-One (A set of demographics is based on a single patient)
- Patients to MedicalHistory: One-to-One (Each patient has a medical history)
- ClinicalMeasurements to Patients: One-to-Many (Each patient can have multiple clinical measurements)
- Patients to Symptoms: One-to-Many (Each patient can have multiple symptoms)
- Diagnosis to Patients: One-to-One (Each patient can have a single diagnosis)

**Relational Schema**

Based on the E/R diagram, the entities and relationships are translated into a relational schema, which includes the tables and fields necessary to implement the database in a relational database management system (RDMS).

**Schema**

- **Patients**: The central table of the schema with `PatientID` as the primary key.
  - PatientID (Primary Key)
  - Name
  - Age
  - Gender
  - Address
- **Demographics**: Includes `DemographicID` as the primary key and contains demographic information about each patient.
  - DemographicID (Primary Key)
  - PatientID (Foreign Key)
  - Ethnicity
  - EducationLevel
  - Occupation
- **MedicalHistory**: Holds details about the patient's medical background.
  - MedicalHistoryID (Primary Key)
  - PatientID (Foreign Key)
  - Hypertension
  - Diabetes
  - CoronaryArteryDisease
  - Anemia
- **ClinicalMeasurements**: Includes `ClinicalMeasurementID` as the primary key and holds details about the patient's clinical measurements.
  - ClinicalMeasurementID (Primary Key)
  - PatientID (Foreign Key)
  - BloodPressure
  - SpecificGravity
  - Albumin
  - Sugar
  - BloodGlucoseRandom
  - BloodUrea
  - SerumCreatinine
  - Sodium
  - Potassium
  - Hemoglobin
  - PackedCellVolume
  - WhiteBloodCellCount
  - RedBloodCellCount
- **Symptoms**: Describes the symptoms experienced by patients.

- ○ SymptomID (Primary Key)
  - ○ PatientID (Foreign Key)
  - ○ ShortnessOfBreath
  - ○ SwellingOfFeet
  - ○ Fatigue
  - ○ Nausea
- **Diagnosis**: Determines the status of the patient's CKD diagnosis.
  - ○ DiagnosisID (Primary Key)
  - ○ PatientID (Foreign Key)
  - ○ DiagnosisDate
  - ○ Stage
  - ○ TreatmentPlan

## Normalization and Evaluation Against Normal Forms

This process involves organizing the fields and tables of the relational schema to minimize redundancy and dependency. The different types of normal forms are as follows:

- **First Normal Form (1NF)**: Ensures that each table has a primary key and that each column contains atomic, indivisible values.
- **Second Normal Form (2NF)**: Ensures that all non-key attributes are fully functionally dependent on the primary key.
- **Third Normal Form (3NF)**: Ensures that all non-key attributes are not only fully functionally dependent on the primary key but also non-transitively dependent.

**Tables Normal Form Justification**

- **Patients**: 3NF - All attributes are directly dependent on the primary key, PatientID.
- **Demographics**: 3NF - All attributes are directly dependent on the primary key, DemographicID, with a foreign key linking to PatientID.
- **MedicalHistory**: 3NF - All attributes are directly dependent on the primary key, MedicalHistoryID, with a foreign key linking to PatientID.
- **ClinicalMeasurements**: 3NF - All attributes are directly dependent on the primary key, ClinicalMeasurementID, with a foreign key linking to PatientID.
- **Symptoms**: 3NF - All attributes are directly dependent on the primary key, SymptomID, with a foreign key linking to PatientID.
- **Diagnosis**: 3NF - All attributes are directly dependent on the primary key, DiagnosisID, with a foreign key linking to PatientID.