# Chronic Kidney Disease (CKD) Prediction – ML Model Development Report

---

## 1. Problem Statement

The hospital's management has requested the development of a predictive machine learning model that can accurately identify **Chronic Kidney Disease (CKD)** in patients using various clinical and biochemical parameters. The goal is to assist doctors with early detection and diagnosis, potentially improving patient outcomes.

---

## 2. Dataset Overview

- **Source**: Provided by hospital

- **File Name**: `CKD.csv`

- **Total Records (Rows)**: 399

- **Total Features (Columns)**: 25

**Key Columns (Features):**

- **Demographic & Clinical**: `age`, `bp` (blood pressure), `sg` (specific gravity), `al` (albumin), `su` (sugar)

- **Lab Test Results**: `bgr`, `bu`, `sc`, `sod`, `pot`, `hemo`, `pcv`, `wc`, `rc`

- **Categorical Parameters**: `rbc`, `pc`, `pcc`, `ba`, `htn`, `dm`, `cad`, `appet`, `pe`, `ane`

- **Target Column**: `classification` (values: `ckd`, `notckd`)

---

# 3. Pre-processing Steps

The dataset required significant preprocessing due to the presence of:

- Missing values (?)

- Inconsistent categorical entries (e.g. yes, no, present, notpresent)

- Mixed data types

## Steps Performed:

| Step | Description |
|---|---|
| Data Cleaning | Removed whitespaces, converted all text to lowercase |
| Missing Values Handling | Replaced ? with NaN, then imputed numeric columns with **mean**, categorical columns using **LabelEncoder** |
| Label Encoding | Encoded categorical variables (e.g., yes/no, normal/abnormal) |
| Target Label Encoding | ckd → 1, notckd → 0 |
| Feature Scaling | Applied **StandardScaler** to normalize numerical features |
| Train-Test Split | Used 80/20 split (random_state=42) |

# 4. Model Development and Evaluation

We trained multiple models using the cleaned and scaled dataset. Below are their performance metrics:

**Model Comparison Table**

```
print(results_df.to_string(index=False))
```

```
Model Performance Summary:

             Model  Accuracy  Precision    Recall  F1 Score   ROC AUC
     Random Forest    1.0000   1.000000  1.000000  1.000000  1.000000
Logistic Regression    0.9875   1.000000  0.975610  0.987654  1.000000
           XGBoost    0.9875   0.976190  1.000000  0.987952  1.000000
     Decision Tree    0.9625   0.975000  0.951220  0.962963  0.962789
               SVM    0.9625   0.975000  0.951220  0.962963  0.998124
               KNN    0.9375   0.973684  0.902439  0.936709  0.997498
```

# 5. Research Results Summary

### (1) Logistic Regression
- Easy to interpret

- Performs well

- ROC AUC: **1.00**

### (2) Decision Tree
- Fast training

- Slightly overfits small datasets

### (3) Random Forest
- Best performance

- Robust to noise and overfitting

### (4) XGBoost
- Excellent accuracy

- Slightly more complex

- Also achieved perfect scores on this dataset

**(5) KNN & SVM**
- Decent performance

- SVM can struggle with large feature sets; KNN sensitive to scaling

---

# 6. Final Model Selection

**Chosen Model: Random Forest Classifier**

**Justification:**

- **Perfect scores** across all evaluation metrics

- **Handles both numerical and categorical data**

- **Robust** to overfitting and noise due to ensemble nature

- **Feature importance** is interpretable and provides insights to clinicians

---

# Conclusion

A highly accurate and robust **Random Forest model** was developed to predict Chronic Kidney Disease (CKD) from patient medical data. This model can now be integrated into hospital systems to support early diagnosis and improve clinical outcomes.