# ASSIGNMENT – II

**TITLE :**
Develop a machine learning model to predict customer churn based on historical data.

**NAME :** Shinde Shubham Dnyandev.      **ROLL NO :** 23107121.      **BATCH :** B