# 머신러닝을 활용한 대출자 연체 여부 예측

2 0 2 0 3 8 0 5 0 1  이 현 지
2 0 2 0 3 8 0 7 2 2  이 승 주

통계학과머신러닝 PBL

# 목차

# 데이터 설명

## 대출 연체 데이터 (Loan Data Set)
(614, 13)

| | Loan_ID | Gender | Married | Dependents | Education | Self_Employed | ApplicantIncome | CoapplicantIncome | LoanAmount | Loan_Amount_Term | Credit_History | Property_Area | Loan_Status |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | LP001002 | Male | No | 0 | Graduate | No | 5849 | 0.0 | NaN | 360.0 | 1.0 | Urban | Y |
| 1 | LP001003 | Male | Yes | 1 | Graduate | No | 4583 | 1508.0 | 128.0 | 360.0 | 1.0 | Rural | N |
| 2 | LP001005 | Male | Yes | 0 | Graduate | Yes | 3000 | 0.0 | 66.0 | 360.0 | 1.0 | Urban | Y |
| 3 | LP001006 | Male | Yes | 0 | Not Graduate | No | 2583 | 2358.0 | 120.0 | 360.0 | 1.0 | Urban | Y |
| 4 | LP001008 | Male | No | 0 | Graduate | No | 6000 | 0.0 | 141.0 | 360.0 | 1.0 | Urban | Y |

## 변수 ●

| Loan_ID | 아이디 | Self_Employed | 자영업 여부 | Credit_History | 과거 신용기록 |
|---|---|---|---|---|---|
| Gender | 성별 | ApplicantIncome | 대출신청자 소득 | Property_Area | 부동산 지역 |
| Married | 결혼 여부 | CoapplicantIncome | 공동대출신청자 소득 | Loan_Status | 연체 상태 |
| Dependents | 부양가족 수 | LoanAmount | 대출금 총액 | | |
| Education | 교육 수준 | Loan_Amount_Term | 대출 기간 | | |

상관분석

CORRELATION ANALYSIS



| | | |
|---|---|---|
| **LoanAmount** | **ApplicantIncome** | **0.57** |
| **Loan_Status** | **Credit_History** | **0.54** |
| **Married** | **Gender** | **0.36** |
| | **Dependents** | **0.33** |
| **Dependents** | **Gender** | **0.17** |

# 데이터 모델링

DATA MODELING

Logistic Regression

KNN

Decision Tree

Bayes Classifier

Support Vector Machine

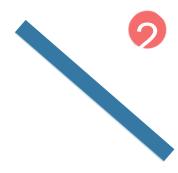Ensemble

# 데이터 모델링

DATA MODELING

AUC = 0.715517

cut-off = 0.676

# 데이터 모델링

DATA MODELING

|  | Train | Test |
|---|---|---|
| Logistic Regression | 0.807 | 0.822 |
| KNN (k=5) | 0.821 | 0.773 |
| Decision Tree (max_depth=5) | 0.844 | 0.789 |
| Bayes Classifier (Linear) | 0.808 | 0.816 |
| Bayes Classifier (Quadratic) | 0.808 | 0.816 |
| Bayes Classifier (Multinomial) | 0.808 | 0.816 |
| Bayes Classifier (Gaussian) | 0.809 | 0.800 |
| Support Vector Machine | 0.807 | 0.811 |

# 데이터 모델링

DATA MODELING

|  | Train | Test |
|---|---|---|
| **Logistic Regression** | **0.807** | **0.822** |
| KNN (k=5) | 0.821 | 0.773 |
| Decision Tree (max_depth=5) | 0.844 | 0.789 |
| **Bayes Classifier (Linear)** | **0.809** | **0.816** |
| Bayes Classifier (Quadratic) | 0.809 | 0.816 |
| Bayes Classifier (Multinomial) | 0.809 | 0.816 |
| Bayes Classifier (Gaussian) | 0.809 | 0.800 |
| **Support Vector Machine** | **0.807** | **0.811** |

# Voting Classifier

| | Train | Test |
|---|---|---|
| Logistic Regression | 0.807 | 0.822 |
| Bayes Classifier (Linear) | 0.809 | 0.816 |
| Support Vector Machine | 0.807 | 0.811 |

| | Train | Test |
|---|---|---|
| Voting Classifier (hard) | 0.809 | 0.816 |
| Voting Classifier (soft) | 0.809 | 0.816 |

# 앙상블 모형

ENSEMBLE MODEL

### Random Forest
Number of Trees = **181**
Maximum Depth of Individual Tree = **6**

### AdaBoost
Number of Trees = **3**
Maximum Depth of Individual Tree = **2**

### Gradient Boosting
Number of Trees = **5**
Maximum Depth of Individual Tree = **4**

### XGBoost
Number of Trees = **2**
Maximum Depth of Individual Tree = **3**

### LightGBM
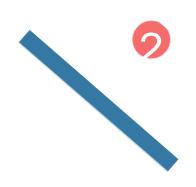Number of Trees = **30**
Maximum Depth of Individual Tree = **4**

### CatBoost

# 앙상블 모형

ENSEMBLE MODEL

|  | Train | Test |
|---|---|---|
| Random Forest | 0.848 | 0.821 |
| AdaBoost | 0.825 | 0.800 |
| Gradient Boosting | 0.822 | 0.805 |
| LightGBM | 0.828 | 0.816 |
| XGBoost | 0.823 | 0.806 |
| CatBoost | 0.828 | 0.822 |

# Random Forest

AUC = 0.896875

| Train | precision | recall | f1-score |
|-------|-----------|--------|----------|
| 0 (N) | 0.53 | 0.97 | 0.69 |
| 1 (Y) | 0.99 | 0.82 | 0.90 |

| Test | precision | recall | f1-score |
|------|-----------|--------|----------|
| 0 (N) | 0.43 | 1.00 | 0.60 |
| 1 (Y) | 1.00 | 0.79 | 0.89 |



Feature importances (Random Forest)

# LightGBM

| Train | precision | recall | f1-score |
|-------|-----------|--------|----------|
| 0 (N) | 0.50 | 0.91 | 0.64 |
| 1 (Y) | 0.98 | 0.81 | 0.89 |

| Test | precision | recall | f1-score |
|------|-----------|--------|----------|
| 0 (N) | 0.45 | 0.93 | 0.60 |
| 1 (Y) | 0.98 | 0.80 | 0.88 |

### Feature importances Top 20 (LightGBM)

| Feature | Importance |
|---------|-----------|
| LoanAmount | |
| ApplicantIncome | |
| CoapplicantIncome | |
| Credit_History | |
| Property_Area | |
| Dependents | |
| Married | |
| Self_Employed | |
| Education | |
| Loan_Amount_Term | |
| Gender | |

# CatBoost

AUC = 0.896875

| Train | precision | recall | f1-score |
|-------|-----------|--------|----------|
| 0 (N) | 0.53 | 0.97 | 0.69 |
| 1 (Y) | 0.99 | 0.82 | 0.90 |

| Test | precision | recall | f1-score |
|------|-----------|--------|----------|
| 0 (N) | 0.43 | 1.00 | 0.60 |
| 1 (Y) | 1.00 | 0.79 | 0.89 |



Feature importances (CatBoost)

# 변수 중요도

FEATURE IMPORTANCES



Feature importances Top 20 (LightGBM)

LightGBM

Feature importances (Random Forest)

Random Forest

Feature importances (CatBoost)

CatBoost

# CatBoost

# 결론

Random Forest

Gradient Boosting

LightGBM

AdaBoost

XGBoost

CatBoost