

# Lending Club Predict Loan Payment

Fri. Aug. 25, 2017  
**Jeon Sohyun**

**Fast Campus  
Data Science School**

### 1. 프로젝트 요약

---

### 2. 전처리, 시각화, 변수 생성

---

### 3. 다중 공선성, PCA

---

### 4. 모델링

---

### 5. 결론

---

## 1. Summary of Project

# Lending Club Loan 데이터를 이용한 채무상환 예측

- Lending Club은 미국 최대의 P2P회사

- 데이터 출처: <https://www.kaggle.com/wendykan/lending-club-loan-data>

기본적이지만, 의사 결정과 문제 해결에 가장 핵심으로 사용되는

### ‘분류문제’

- 물건 구매예측
- 대출 상환예측
- 마케팅 성공여부

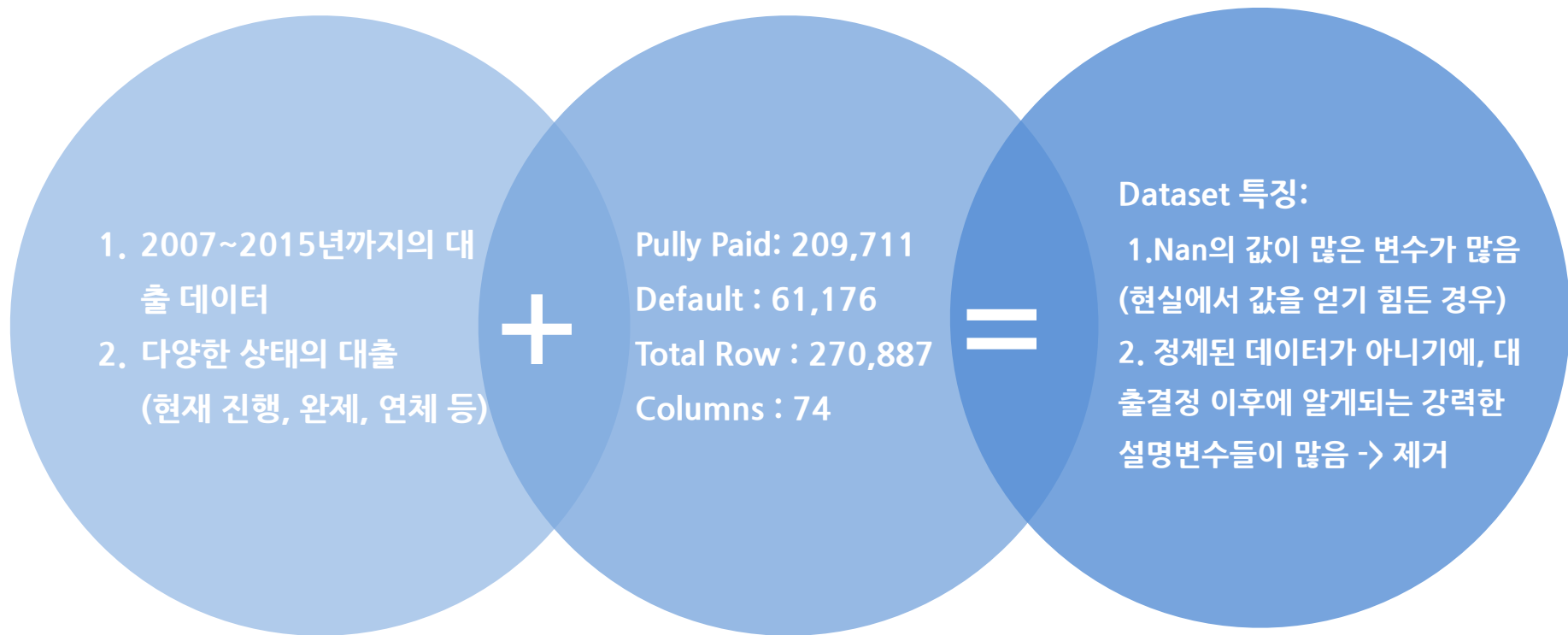
1. 대출 상환 여부는 금융권에 있어서 가장 중요한 문제.

2. 부실은 과도한 충당금으로 건전성에 영향을 줄 수 있음.

3. 효과적인 상환예측모델은 비용감소를 기대할 수 있음

현재 국내의 금융은 엄격한 규제  
로 비식별 데이터를 구하기 힘든  
상황이기에  
Kaggle의 Dataset을 이용.  
(Competition은 아님)

## 1. Summary of Project



## 2. Preprocessing, Visualize

### 1. 데이터 전처리

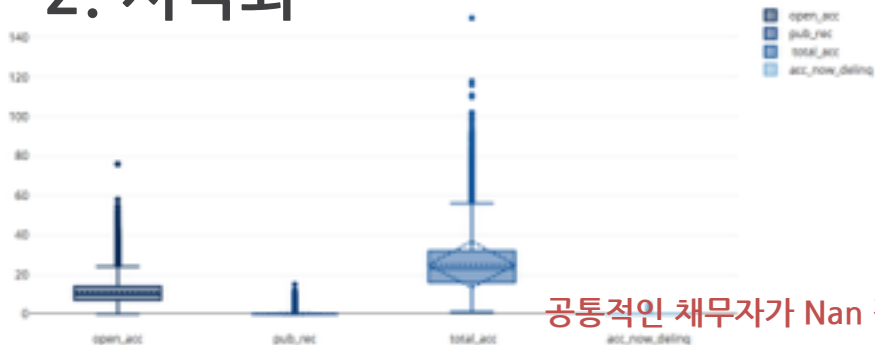
initial_list_status	0	id	0
collections_12_mths_ex_med	97	loan_amnt	0
mths_since_last_major_derog	164247	funded_amnt	0
policy_code	0	funded_amnt_inv	0
application_type	0	term	0
annual_inc_joint	203162	int_rate	0
dti_joint	203162	installment	0
verification_status_joint	203162	grade	0
acc_now_delinq	17	sub_grade	0
tot_coll_amt	49996	emp_title	11317
tot_cur_bal	49996	emp_length	0
open_acc_6m	203056	home_ownership	0
open_il_6m	203056	annual_inc	1
open_il_12m	203056	verification_status	0
open_il_24m	203056	issue_d	0
mths_since_rent_il	203059	pymnt_plan	0
total_bal_il	203056	url	0
il_util	203068	desc	133865
open_rv_12m	203056	purpose	0
open_rv_24m	203056	title	10
max_bal_bc	203056	zip_code	0
all_util	203056	addr_state	0
total_rev_hi_lim	49996	dti	0
inq_fi	203056	delinq_2yrs	17
total_cu_tl	203056	earliest_cr_line	17
inq_last_12m	203056	inq_last_6mths	17
loan_status	0	mths_since_last_delinq	111956
dtype: int64		mths_since_last_record	176929



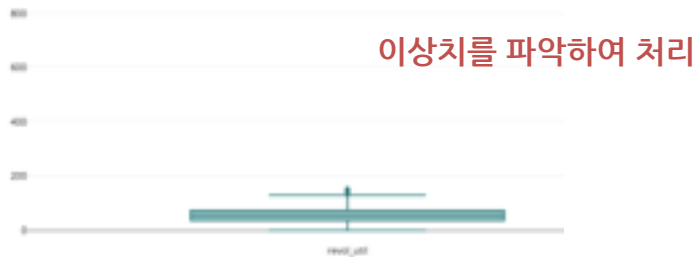
각 변수의 특징을 파악,  
추론하여 Nan 값을  
평균값, 0, 4분위수 등으로 처리

## 2. Preprocessing, Visualize

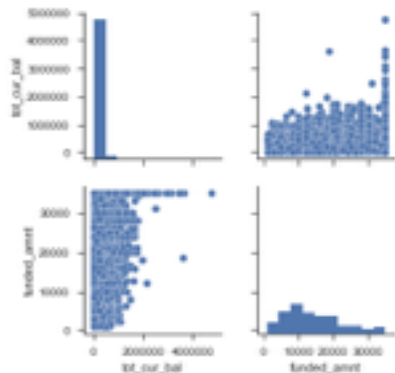
## 2. 시각화



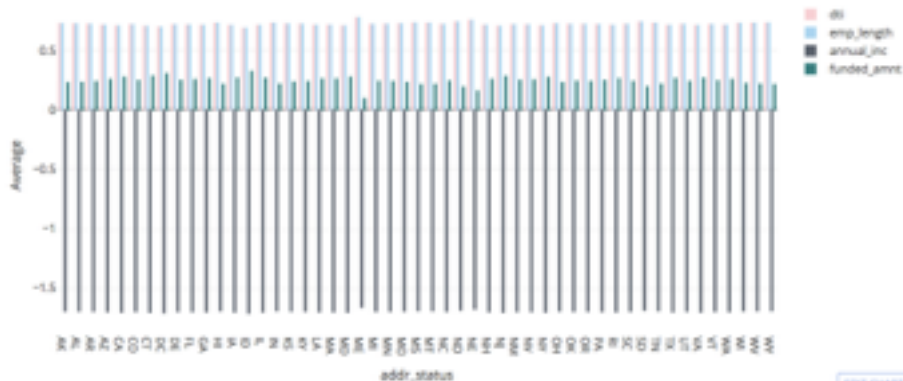
공통적인 채무자가 Nan 값을 가지는 변수의 분포를 파악하여 값을 추정



이상치를 파악하여 처리



두 변수간의 상관관계를 파악



## 2. Preprocessing, Visualize

### 3. 변수 분석 및 생성

```
loan_check.head()
```

Out[78]:

loan_status	sub_grade	paid	default
0	A1	4291	177
1	A2	4462	274
2	A3	5037	364
3	A4	7198	720
4	A5	8299	978

```
In [79]: loan_check['grade_percentage'] = loan_check.paid / (loan_check.paid + loan_check.default)
```

Out[79]:

loan_status	sub_grade	paid	default	grade_percentage
0	A1	4291	177	0.96090
1	A2	4462	274	0.94046
2	A3	5037	364	0.93365
3	A4	7198	720	0.91174
4	A5	8299	978	0.89478
5	B1	9202	1129	0.87905
6	B2	9739	1336	0.88248

Stringtype의 신용등급에 대한 가중치 계산하여 변수값 할당

	emp_length	annual_inc	norm
0	0	47968.955389	0.205355
1	1	67266.308913	0.287961
2	2	70006.602745	0.299701
3	3	70812.048165	0.303140
4	4	69765.991235	0.298662
5	5	71294.267054	0.305205
6	6	71431.280815	0.305791
7	7	72895.635531	0.312060
8	8	74238.551182	0.317809
9	9	74638.035644	0.319519
10	10	79853.937696	0.341848

재직기간에 따라 One-hot-encoding 외에 연소득과의 관계로 가중치 할당

	issue_d	earliest_cr_line	credit_period
0	24124	23983	141
1	24164	24008	150
2	24158	24003	155
3	24161	23883	278
4	24158	24000	158
5	24184	23982	202
6	24156	24000	156
7	24176	23997	179
8	24170	24005	165
9	24154	23743	411
10	24185	24008	179

최초 신용개설일과 대출시작시점을 계산하여 신용 기간 추정

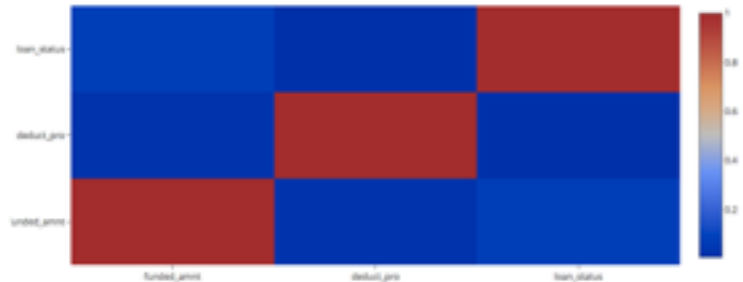
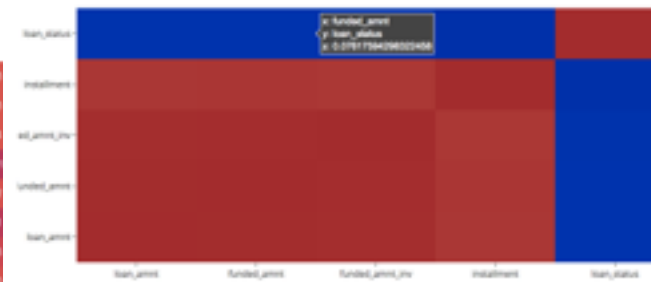
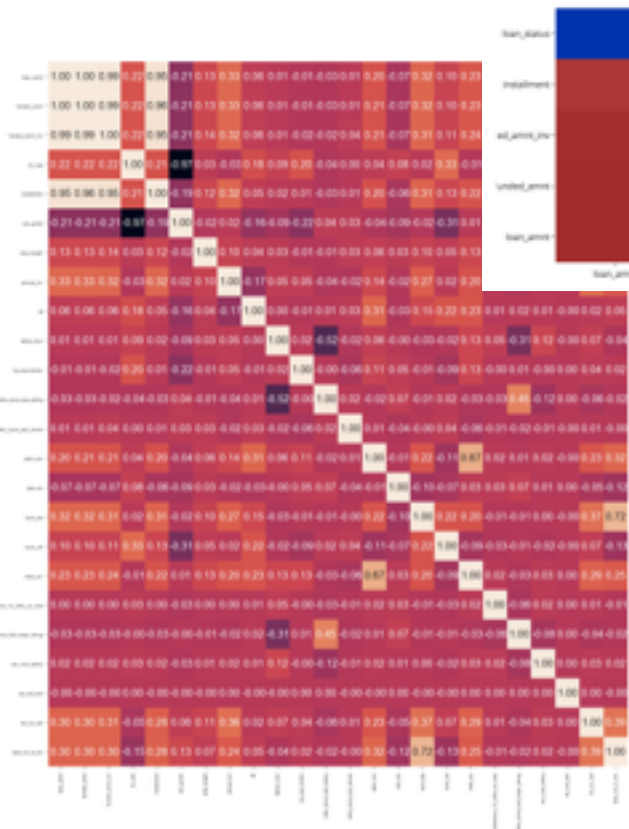
그 이외의 변수 생성과  
질적변수의 더미화



columns : 74 -> 142

### 3. Correlation, PCA

## 1. 다중 공선성 확인



loan\_amnt - funded\_amnt : 1.0  
 loan\_amnt - funded\_amnt\_inv : 0.99  
 loan\_amnt - installment : 0.95

funded\_amnt - funded\_amnt\_inv : 0.99  
 funded\_amnt - installment : 0.95

funded\_amnt\_inv - installment : 0.95

open\_acc - total\_acc : 0.67  
 revol\_bal - total\_rev\_hi\_lim : 0.72

mths\_since\_last\_delinq -  
 mths\_since\_last\_major\_derog : 0.45

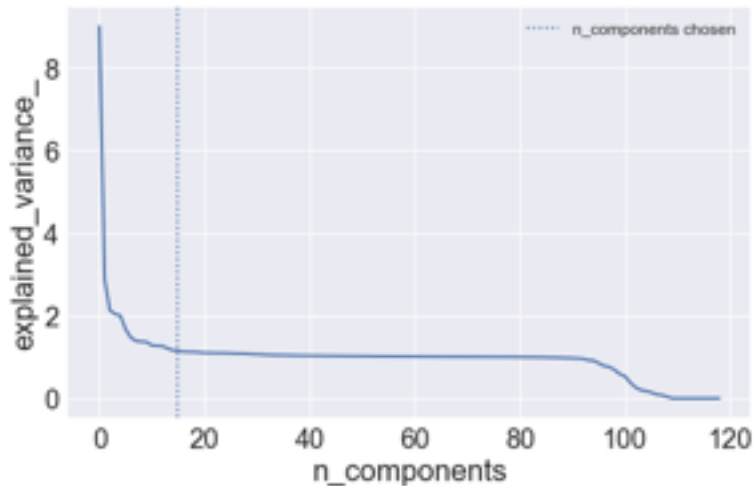
상관관계가 큰 변수는 제거



### 3. Correlation, PCA

## 2. PCA

- 0의 값이 전체의 80%를 초과하는 변수와 더미화한 변수를 PCA로 차원축소 한다 (과적합 우려)



119



15

### 3. Correlation, PCA

## 3. 결과

#### # 그라디언트 부스트

```
from sklearn.ensemble import GradientBoostingClassifier
gbdt=GradientBoostingClassifier(n_estimators=200)

gbdt.fit(X_train, y_train)

print (gbdt.score(X_train, y_train))
print (gbdt.score(X_val, y_val))

0.777178371496
0.734594846492

y_pred = gbdt.predict_proba(X_val)

#y_pred = xgb.predict_proba(X_val)
y_true = y_val
pred_random=[]
THRESHOLD = 0.23

for val in y_pred:
    if val[1] > THRESHOLD:
        pred_random.append(1)
    elif val[1] <= THRESHOLD:
        pred_random.append(0)

score = f1_score(y_true, pred_random)
print ("f1 score: {}".format(score))

f1 score: 0.468144925342338
```

다중공선성 변수제외, PCA 적용 후

#### Light gbm

```
123> gbm = lgb.LGBMRegressor(objective='binary',
                             num_leaves=31,
                             learning_rate=0.05,
                             n_estimators=100)

124> gbm.fit(X_train, y_train,
            eval_set=[(X_val, y_val)],
            eval_metric='log_loss',
            early_stopping_rounds=5)

125> print (gbm.score(X_train, y_train))
print (gbm.score(X_val, y_val))

0.131354963325
0.105339628872

126> y_pred = gbm.predict(X_val, num_classes=2)
y_true = y_val
pred_random=[]
THRESHOLD = 0.24

for val in y_pred:
    if val > THRESHOLD:
        pred_random.append(1)
    elif val <= THRESHOLD:
        pred_random.append(0)

score = f1_score(y_true, pred_random)
print ("f1 score: {}".format(score))

f1 score: 0.46600651713119686
```

#### Random Forest

```
161> model_random = RandomForestClassifier(n_estimators=500)

171> model_random.fit(X_train, y_train)

172> RandomForestClassifier(bootstrap=True, class_weight=None,
                           max_depth=None, max_features='sqrt', max_leaf
                           min_impurity_split=0.07, min_samples_leaf=1,
                           min_samples_split=2, min_weight_fraction=0.1,
                           n_estimators=500, n_jobs=-1, oob_score=False,
                           verbose=0, warm_start=False)

181> y_pred = model_random.predict_proba(X_val)
y_true = y_val
pred_random=[]
THRESHOLD = 0.25

for val in y_pred:
    if val[1] > THRESHOLD:
        pred_random.append(1)
    elif val[1] <= THRESHOLD:
        pred_random.append(0)

random_f1 = f1_score(y_true, pred_random)
random_train_score = model_random.score(X_train, y_train)
random_val_score = model_random.score(X_val, y_val)

191> print("train score: {}".format(random_train_score))
print("val score: {}".format(random_val_score))
print("f1 score: {}".format(random_f1))

train score: 1.0
val score: 0.7779073823279256
f1 score: 0.47182413998957504
```

#### 그라디언트 부스트

```
11 from sklearn.ensemble import GradientBoostingClassifier
gbdt=GradientBoostingClassifier(n_estimators=200)

12 gbdt.fit(X_train_scaled, y_train_scaled)

13 GradientBoostingClassifier(criterion='friedman_mse', learning_rate=0.01, loss='deviance',
                             max_features=None, max_leaf_nodes=None,
                             min_impurity_split=0.7, min_samples_leaf=1,
                             min_samples_split=2, min_weight_fraction=0.1,
                             n_estimators=200, presort='auto',
                             subsample=1.0, verbose=0, warm_start=False)

14 print (gbdt.score(X_train_scaled, y_train_scaled))
print (gbdt.score(X_val_scaled, y_val_scaled))

0.777684954259
0.73464955111

15 y_pred = xgb.predict_proba(X_val_scaled)
y_true = y_val_scaled
pred_random=[]
THRESHOLD = 0.24

for val in y_pred:
    if val[1] > THRESHOLD:
        pred_random.append(1)
    elif val[1] <= THRESHOLD:
        pred_random.append(0)

score = f1_score(y_true, pred_random)
print ("f1 score: {}".format(score))

f1 score: 0.4691500997233482
```

다중공선성 변수제외, PCA 적용 전

적용전이 적용후 보다 점수가 좀더 높은 경향으로 보여 원 데이터를 살리기로 함

## 4. Modelling

### 1. y값의 비대칭 문제 발견

Model	scaled	hyperparameter	threshold	train score	val score	F1 score	ci
Random Forest	x	max_features='sqrt', n_estimators=500, n_jobs=-1, random_state=20	0.25	1.0	0.776	0.469	
	x	n_jobs=-1, max_features='sqrt', random_state=20, n_estimators=200	0.25	0.944	0.943	0.8298	
	x	n_jobs=-1, max_features='sqrt', random_state=20, n_estimators=200	0.32	0.944	0.945	0.85125	
	x	n_jobs=-1, max_features='sqrt', random_state=20, n_estimators=200	0.42	0.944	0.944	0.8648	
	x	n_jobs=-1, max_features='sqrt', random_state=10, n_estimators=200	0.42	0.944	0.944	0.8646	

하이퍼 파라미터를 튜닝하면서 획기적으로 높은 스코어를 얻음  
하지만 Test set을 적용 시 F1 Score가 0.21이 되는 기형적 현상을 발견

〈Training Set〉  
Pully Paid: 157,229  
Default : 45,936

상환과 부채의 비율이 7:3정도로  
모델이 불균형데이터에 학습,  
적절한 예측을 하지 못함  
최적화 예측 등 악기 노출

## 4. Modelling

# 2. 해결책

### (2) Dataframe sampling and shuffle

```
def shuffle_data(X_0, X_1): # select same number of sample with X_1 from X_0, and shuffle
    X_sample = X_0.sample(n=45936)
    X = pd.concat([X_sample, X_1])
    X = X.sample(frac=1).reset_index(drop=True)
    y = X['loan_status']
    X = X.drop('loan_status', axis=1)
    return X, y
```

Shuffle 함수

### (3) Cross Validation

```
def cross(model, X_0, X_1, thres):
    X, y = shuffle_data(X_0, X_1)
    train_list = []
    test_list = []
    fi_list = []
    auc_list = []

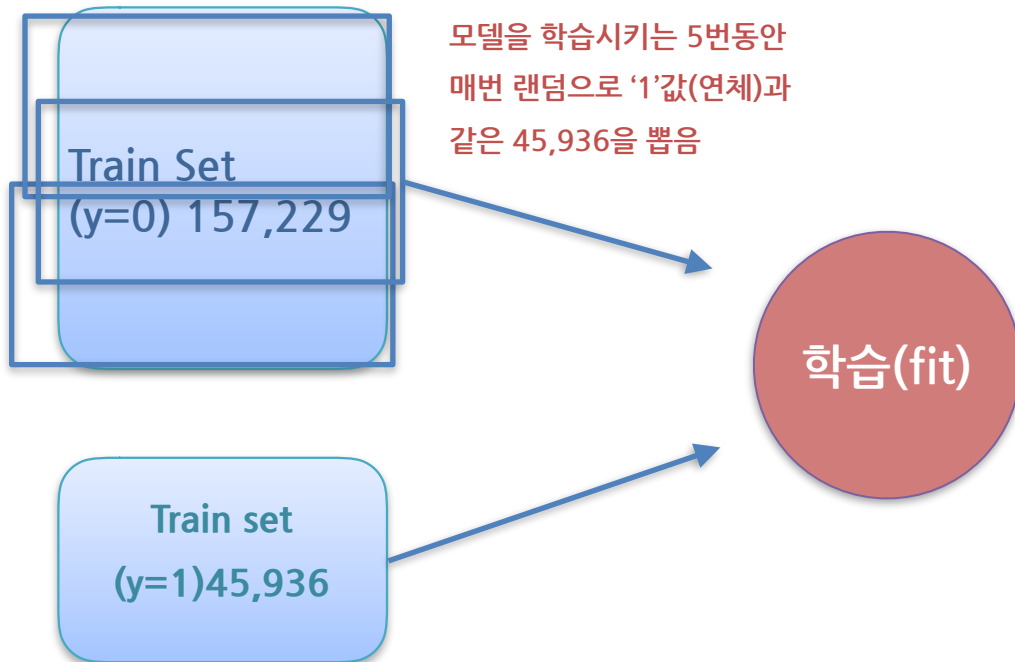
    # 5 times modeling and take score
    for i in range(5):
        X_train, X_val, y_train, y_val = train_test_split(X, y)

        # Training
        fitted_model = model.fit(X_train, y_train)

        # accuracy
        train_list.append(model.score(X_train, y_train))
        test_list.append(model.score(X_val, y_val))

        # AUC & ROC
        auc_list.append(roc_auc_score(y_val, model.predict(X_val)))
```

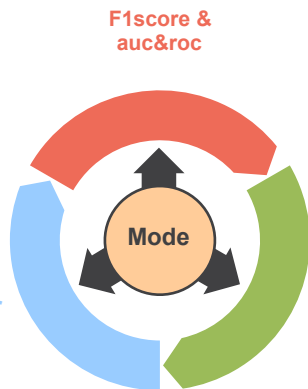
Cross Validation



## 4. Modelling

# 3. Hyper Parameter Tunning

- Grid Search를 이용하여 자동 Hyper Parameter를 뽑아 보았으나, 수기로 조정 시 보다 Performance가 좋지 않음을 확인함



Gradient Boosting	n_estimators=200, learning_rate=0.05, max_depth=5	0.34	0.66788	0.66025	0.6889		
	n_estimators=300, learning_rate=0.01, max_depth=4	0.34	0.66649	0.69041	0.6809		
	n_estimators=100, learning_rate=0.001, max_depth=3	0.34	0.6344218	0.6301	0.667854		
	n_estimators=100, learning_rate=0.001, max_depth=5	0.34	0.64030	0.6363	0.668889		
Neural network	alpha=0.001,max_iter= 300, hidden_layer_size=100	Random Forest	n_jobs=-1, max_features='sqrt', random_state=20, n_estimators=300	0.42	1.0	0.6544	0.691
	alpha=0.001,max_iter= 100, hidden_layer_size=100		n_jobs=-1, max_features='sqrt', random_state=20, n_estimators=500	0.29	1.0	0.65485	0.69637
	alpha=0.001,max_iter= 300, hidden_layer_size=200		n_jobs=-1, max_features='sqrt', random_state=2, n_estimators=300	0.29	1.0	0.6563	0.6926
	alpha=0.001,max_iter= 300, hidden_layer_size=100		n_jobs=-1, max_features='sqrt', random_state=2, n_estimators=300	0.36	1.0	0.65533	0.69706
	alpha=0.001,max_iter= 300, hidden_layer_size=100		n_jobs=-1, max_features='log2', random_state=25, n_estimators=300	0.33	1.0	0.63332	0.69604
	alpha=0.001,max_iter= 300, hidden_layer_size=300		random_state=20, n_estimators=200	0.33	1.0	0.6528	0.694783
Light GBM	objective='regression', num_leaves=40, learning_rate=0.001, n_estimators=1000	random_state=2, n_estimators=500, n_jobs=-1, max_features='sqrt'	0.33	1.0	0.654536	0.69787	0.65491
	objective='regression', num_leaves=40, learning_rate=0.001, n_estimators=2000, early_stopping_rounds	random_state=2, n_estimators=500, n_jobs=-1, max_features='sqrt'	0.35	1.0	0.6558	0.69857	0.6557
		random_state=2, n_estimators=500, n_jobs=-1, max_features='sqrt'	0.35	1.0	0.66009	0.70034	0.66019
	objective='regression', learning_rate=0.001, n_estimators=1000		0.65546	0.6481	0.6483		
	objective='regression', num_leaves=40, learning_rate=0.001, n_estimators=1000	0.36	0.676800	0.65848	0.69629	0.63855	
	objective='regression', num_leaves=50, learning_rate=0.001, n_estimators=1000	0.36	0.68080	0.657127	0.68805	0.65716	
	objective='regression', num_leaves=100, learning_rate=0.001, n_estimators=1000	0.36	0.70343	0.66066	0.70038	0.66047	
	objective='regression', num_leaves=200, learning_rate=0.001, n_estimators=1000	0.36	0.74348	0.6621	0.70188	0.66209	
	objective='regression', num_leaves=300, learning_rate=0.001, n_estimators=1000						

## 4. Model Selection

\*Predicting overfitting, tuning hyper parameter

## <Random Forest>

F1: 0.699

auc&roc: 0.661

- Gradient Boosting takes too long time, gridsearch w

Non-scaled

## Gradient Boosting

F1: 0.690

auc&roc: 0.658

random\_states(2,alpha(0.0001 -> 1) :weight to be strong

<NN>

F1: 0.637

auc&roc: 0.555

```
xgbr = xgb.XGBClassifier(n_estimators=30, c
```

## <XGBoost>

F1: 0.696

auc&roc: 0.652

```
gbm = lgb.LGBMClassifier(objective='regression',
                           num_leaves=200,
                           learning_rate=0.001,
                           n_estimators=1000,
                           # early_stopping_rounds=
                           )
```

## <LightGBM>

F1: 0.698

auc&roc: 0.657

속도와 성능을 고려 시,  
LightGBM이 적절할  
것으로 판단

## 5. Conclusion

# 1. Test set Score

### 1) Random Forest

```
test(model_random_final, X_test, y_test, 0.35)
```

```
test score: 0.6466805383086  
auc&roc score : 0.656774332377025  
f1 score : 0.4235437819680892
```

### 2) Gradient Boosting

```
test(model_gradient_final, X_test, y_test, 0.25)
```

```
test score: 0.6570095389976669  
auc&roc score : 0.6582475146334387  
f1 score : 0.3997114978974729
```

### 3) Neural network model

```
test(model_nn_final, X_test, y_test, 0.33)
```

```
test score: 0.6462597888055449  
auc&roc score : 0.5806297067130928  
f1 score : 0.37103221255584296
```

### 4) XGBoost

```
test(model_xg_final, X_test, y_test, 0.33)
```

```
test score: 0.654957030211748  
auc&roc score : 0.6550374639133867  
f1 score : 0.42015710776136606
```

### 5) LightGBM

```
test(model_gbm_final, X_test, y_test, 0.36)
```

```
test score: 0.6529045214258291  
auc&roc score : 0.6562043019173878  
f1 score : 0.4350461047212538
```

Test set에서도 가장 좋은  
점수가 나옴을 알 수 있음

## Main Hyper Parameter

num\_leaves=200,  
earning\_rate=0.001,  
n\_estimators=3000,

## 상위 주요 변수 5위

annual\_inc 0.073342

dti 0.072496

revol\_util 0.066707

credit\_period 0.058479

tot\_cur\_bal 0.057720

log\_cml\_psi 0.021150

clsqir\_baupq 0.028118

## 5. Conclusion

# 2. Conclusion

하버드대학의 Dr. Asim Khwaja의 최근 논문에 따르면, 금융 정보 이외의 비금융 정보 (신청서 작성 시 자필 및 맞춤법, 통신 및 보험 정보, SNS의 평판 등)를 가지고 좀 더 정확하고 강력한 대출 상환 예측을 할 수 있다는 연구가 있었다.

현재 국내에서는 다양한 규제로 인해 비식별 데이터를 얻기 힘든 실정이긴 하지만 이런 비금융 정보를 이용하여 효과적인 모형을 만들어 보고 싶다



**Thank You**