

주제 소개

학회장팀

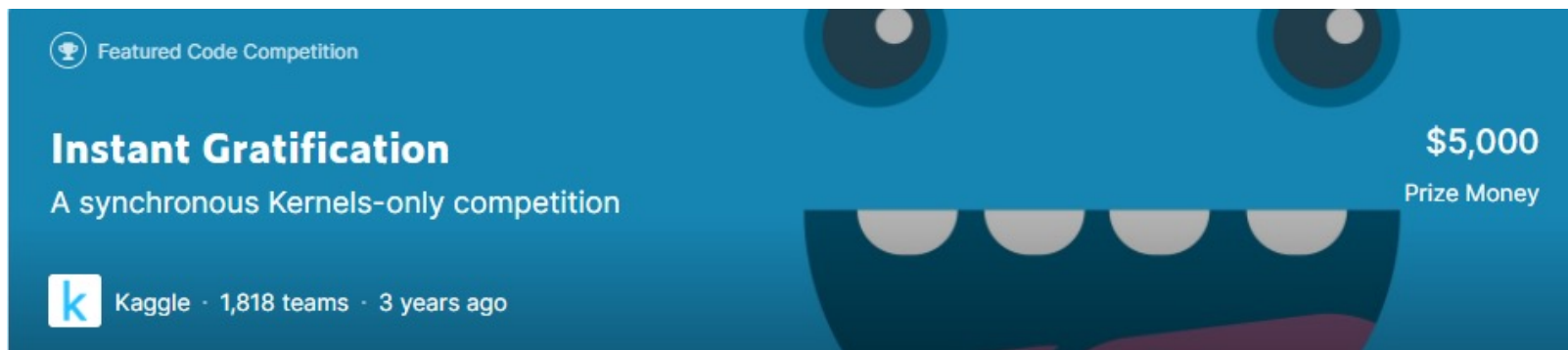


주혜인



이승우

| Reference



Kaggle의 Instant Gratification 대회를 레퍼런스로
데이터 생성 및 평가 지표를 결정

데이터 생성

```
1 NUM_SUB_DATASETS = 32
2 NUM_SAMPLES = 2048
3 NUM_FEATURES = 255
4 MAX_SEED = 2**32 - 1

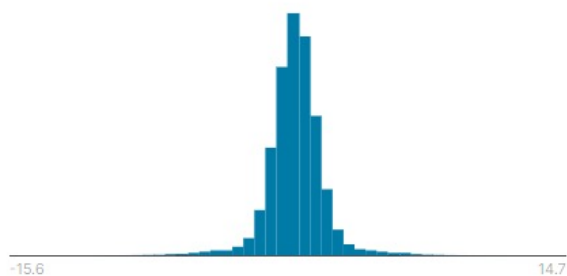
1 def create_dataset(random_seed):
2     random.seed(3 + random_seed)
3     X,y = make_classification(n_samples=NUM_SAMPLES,
4                             n_features=NUM_FEATURES,
5                             n_informative=random.randint(33,47),
6                             n_redundant=0,
7                             n_repeated=0,
8                             n_classes=2,
9                             n_clusters_per_class=3,
10                            weights=None,
11                            flip_y=0.05,
12                            class_sep=1.0,
13                            hypercube=True,
14                            shift=0.0,
15                            scale=1.0,
16                            shuffle=True,
17                            random_state=random_seed)
18     df = pd.DataFrame(X, columns=['V' + str(x) for x in range(0,len(X[0]))])
19     df['magic'] = random_seed
20     df['target'] = y
21     return df
```

아래 링크를 참고로 scikit-learn 라이브러리의
make_classification() 함수를 이용해 데이터 생성

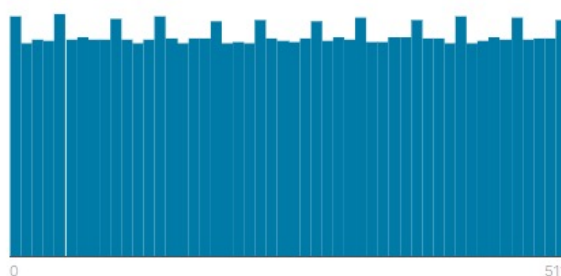
| 데이터 생성



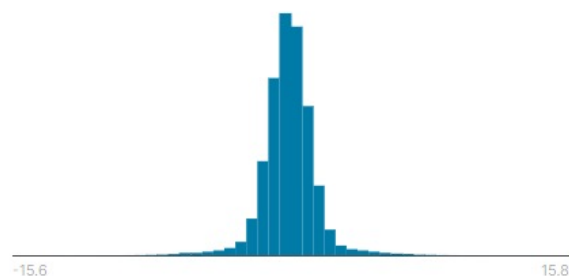
jumpy-thistle-discus-sorted



wheezy-copper-turtle-magic



muggy-smalt-axolotl-pembus



기존 대회와 달리 변수 이름이 복잡했지만 큰 의미가 없어
불필요한 데이터 탐색으로 인한 시간 낭비를 막고자 변수명 단순화

데이터 생성

```
1 seed_list = random.sample(range(1, MAX_SEED), NUM_SUB_DATASETS)
2 df = pd.concat([create_dataset(s) for s in seed_list], axis=0, sort = False).reset_index(drop=True)
3 df = df.sample(frac=1, random_state=9726).reset_index(drop=True)
4
5 le = preprocessing.LabelEncoder()
6 le.fit(df['magic'])
7 df['magic'] = le.transform(df['magic'])
```

```
1 train = df.groupby(['magic']).sample(frac=0.75, random_state=17)
2 test = df.drop(train.index, axis=0)
```

```
1 train = train.sample(frac=1).reset_index(drop=True)
2 test = test.sample(frac=1).reset_index(drop=True)
3
```

32개의 데이터셋을 이어 붙여 최종 데이터셋 생성



THANK YOU

