

กลุ่ม Diabetes Prediction

- 6610402205 นายอภิชาติ รุจระนโษ พญ 1
- 6610402132 นายวรวิทย์ ลี้นกรราชโชติ พญ 1
- 6610401965 นายรัชชิต พญวัฒนา พญ 1

วัตถุประสงค์ของระบบต้นแบบ

กลุ่มของแพทย์และผู้สนใจมีความสำคัญทางด้านสุขภาพของประชาชนคนไทย ในปัจจุบันคนไทยประสบกับปัญหาสุขภาพมากมายแต่ยังขาดการรู้ตัวว่าตนเองมีสุขภาพ และคนไทยเป็นจำนวนมากที่ขาดความรู้โรคเบาหวาน พฤติกรรมการใช้ชีวิตของคนไทยในปัจจุบันยังจำเป็นต้องมีค่าใช้ชีวิตรวมของผู้ป่วยโรคเบาหวานมาช่วยพิจารณาเพื่อค้นหาว่าผู้ป่วยรายใดมีความเสี่ยงที่จะเป็นโรคเบาหวานหรือไม่เพื่อให้ได้การตรวจวัดและประเมินผลว่าผู้ป่วยใดที่มีความการใช้ชีวิตของเขาสอดคล้องกับแบบของเขาหรือไม่เป็นค่า 100 เมตรเช่นนี้ไม่สามารถใช้ได้เลยเพราะมันมีค่าเป็นจำนวนเต็มซึ่งเมื่อเราสุ่มสุ่มมาใช้จำนวนจะได้ค่าที่ความละเอียดทางค่าที่น้อยเกินไปซึ่งถ้าเราสุ่มสุ่มค่าออกมาแล้ว

สิ่งที่ไปยังข้อมูลที่จะใช้ในระบบต้นแบบ

Link to data: <https://www.kaggle.com/datasets/alextebboul/diabetes-health-indicators-dataset>

In [] :	<pre>import pandas as pd diabetes = pd.read_csv('diabetes_binary_health_indicators_BRFSS2015.csv')</pre>
Data Engineer	
ตรวจสอบข้อมูลสูญหาย (Missing Value)	
In [] :	<pre>diabetes.isna().sum()</pre>
Out[] :	<pre>Diabetes_binary 0 HighBP 0 HighChol 0 CholCheck 0 BMI 0 Smoker 0 Stroke 0 HeartDiseaseorAttack 0 PhysActivity 0 Fruits 0 Veggies 0 HvyAlcoholConsump 0 AnyHealthcare 0 NeDocCost 0 GenHlth 0 MenHlth 0 DiffWalk 0 PhysHlth 0 Sex 0 Age 0 Education 0 Income 0 dtype: int64</pre>
ตรวจสอบข้อมูลซ้ำซ้อน	
In [] :	<pre>diabetes.duplicated().sum()</pre>
Out[] :	<pre>24206</pre>
In [] :	<pre>diabetes.drop_duplicates(inplace=True)</pre>

อธิบายความ : ในไฟล์ข้อมูลซึ่งมีทั้งหมดทั้งหมด 5 ล้านชื่อที่มาจากจากทาง feature importance ซึ่งได้มาจากทางกลุ่มการวิเคราะห์ของทาง IBM GenHlth, BMI, HeartDiseaseorAttack, DiffWalk, HvyAlcoholConsump ซึ่งได้ดำเนินการใช้โปรแกรมด้วยโปรแกรมในลักษณะที่นำข้อมูลมาวิเคราะห์โดยใช้วิธีทางความสัมพันธ์ของแบบสถิติเพื่อ เพื่อสร้างการเชื่อมโยงว่ามีความสัมพันธ์กับสถิติเชิงพยากรณ์ โดยได้ทำการใช้วิธีการทางความสัมพันธ์กับเชิงวิเคราะห์จากทางกลุ่มของทาง IBM GenHlth, HighBP, HighChol, BMI, Age, DiffWalk และ DiffWalk

	Fruits	-0.024805	-0.019467	-0.026257	-0.017860	-0.067424	-0.061731	-0.004486	...	-0.006946	0.125023	1.000000	...	0.022659	-0.032387	-0.071221	-0.052191	-0.024441	-0.029932	-0.088768	0.073515	0.084857	0.050907
	Veggies	-0.041734	-0.042994	-0.027399	-0.000553	-0.044054	-0.013744	-0.033029	...	-0.027180	0.135240	0.242941	...	0.020530	-0.019876	-0.094115	-0.042515	-0.045130	-0.063189	-0.066113	-0.003586	0.131624	0.125068
	HvyAlcoholConsump	-0.065650	-0.014178	-0.119057	-0.101957	-0.029759	-0.058420	-0.021347	...	-0.035561	0.023378	-0.028221	...	-0.006202	-0.001272	-0.057673	0.016852	-0.036860	-0.047655	-0.041018	0.039132	0.071863	
	PhysActivity	-0.100004	-0.104382	-0.063443	-0.004555	-0.127780	-0.066869	0.059306	...	0.073094	1.000000	0.125023	...	0.023959	-0.046440	-0.237511	-0.105914	-0.193007	-0.235719	0.033516	-0.087881	0.170931	
	Education	-0.102686	-0.112887	-0.050045	-0.009935	-0.074433	-0.136557	-0.064178	...	-0.082288	0.170931	0.084857	...	0.111367	-0.083260	-0.244752	-0.076122	-0.127687	-0.169350	0.015956	-0.092747	1.000000	
	Income	-0.140659	-0.140030	-0.062089	0.001989	-0.069097	-0.095314	-0.117108	...	-0.122728	0.165869	0.050907	...	0.146144	-0.187577	-0.331782	-0.185689	-0.240929	-0.299064	0.130997	-0.116361	0.419045	
	22 rows × 22 columns																						

```
In [ ] : # Features = ["GenHlth", "BMI", "HeartDiseaseorAttack", "DiffWalk", "HvyAlcoholConsump"] # Feature Importance
Features = ["GenHlth", "HighBP", "HighChol", "CholCheck", "Bmi", "Age", "DiffWalk", "Sex"] # Correlation
# Features = diabetes.drop("Diabetes_Binary", axis=1).columns # All
label = "Diabetes_Binary"

# Feature Importance and Imbalance Feature การทำ Under Sampling

from joblib import under_sampling
import RandomUnderSampler
```

```
In [ ] : # Features = ["GenHlth", "BMI", "HeartDiseaseorAttack", "DiffWalk", "HvyAlcoholConsump"] # Feature Importance
features = ["GenHlth", "HighBP", "HighChol", "BMI", "Age", "DiffWalk"] # Correlation
# Features = diabetes.drop("Diabetes_binary", axis=1).columns # All
label = "Diabetes_binary"
```

ตรวจสอบและจัดการกับ Imbalance ด้วยการทำ Undersampling

8	5.0	1.0	1.0	30.0	9.0	1.0
10	3.0	0.0	0.0	25.0	13.0	0.0
13	4.0	1.0	1.0	28.0	11.0	1.0
15	2.0	1.0	0.0	33.0	6.0	0.0
16	3.0	1.0	1.0	21.0	10.0	0.0

```
In [ ]: diabetes.Diabetes_binary.value_counts()

Out[ ]: Diabetes_binary
1.0    35097
0.0    35097
Name: count, dtype: int64
```

In [] : diabetes.Diabetes_binary.value_counts()

```
Out[ ] : Diabetes_binary
1.0    35987
0.0     35987
Name: count, dtype: int64
```

Traning Data

แบ่งชุดข้อมูลเป็น Train, ชุดข้อมูลทดสอบ(Test), และชุดข้อมูลตรวจสอบ (Validation)

In [] :	<pre>from sklearn.model_selection import train_test_split train, test = train_test_split(diabetes, test_size = 0.3, stratify = diabetes.Diabetes_binary, random_state = 1234) test, validation = train_test_split(test, test_size = 0.3, stratify = test.Diabetes_binary, random_state = 1234)</pre>
In [] :	<pre>len(train), len(test), len(validation)</pre>
Out[] :	<pre>(49139, 14741, 6518)</pre>

ตรวจสอบความสมดุลของชุดข้อมูล Train

In [] :	<pre>train.Diabetes_binary.value_counts()</pre>
Out[] :	<pre>Diabetes_binary 0.0 24568 1.0 24567 Name: count, dtype: int64</pre>

Model Development

สร้าง Pipeline เข้ามาเอา Logistic Regression มาเข้า Feature Engineer ใน Pipeline

In [] :	<pre>from sklearn.pipeline import make_pipeline from sklearn.compose import make_column_transformer from sklearn.preprocessing import KBinsDiscretizer from sklearn.model_selection import LogisticRegression def create_pipeline_logistic(): return make_pipeline(make_column_transformer((KBinsDiscretizer(encode="ordinal", strategy="quantile"), ["BMI"]), remainder="passthrough"), LogisticRegression(max_iter=500)) logistic_pipeline = create_pipeline_logistic() logistic_pipeline</pre>
Out[] :	<div><div>Pipeline</div><div><div>columntransformer: ColumnTransformer</div><div><div>kbinsdiscretizer</div><div>remainder</div><div>KBinsDiscretizer</div><div>passthrough</div></div></div><div>LogisticRegression</div></div>

In [] :	<pre>#time logistic_pipeline.fit(train[features], train[label]) logistic_train_score = logistic_pipeline.score(train[features], train[label]) logistic_validation_score = logistic_pipeline.score(validation[features], validation[label]) print(f"Train Score: %.3f" % logistic_train_score) print(f"Validation Score: %.3f" % logistic_validation_score)</pre>
Out[] :	<pre>CPU times: total: 31.2 ms wall time: 17.3 ms Train Score: 0.738 Validation Score: 0.728 Hyperparameter Tuning</pre>

In [] :	<pre>logistic_pipeline.get_params()</pre>
Out[] :	<pre>{'memory': None, 'steps': [(('columntransformer', ColumnTransformer(remainder='passthrough', transformer=[('kbinsdiscretizer', transform=KBinsDiscretizer(encode='ordinal'), ['BMI'])])), ('logisticregression', LogisticRegression(max_iter=500))), (verbose': False, 'columntransformer': ColumnTransformer(remainder='passthrough', transformer=[('kbinsdiscretizer', KBinsDiscretizer(encode='ordinal'), ['BMI'])])), 'logisticregression': LogisticRegression(max_iter=500), 'columntransformer__force_int_remainder_cols': True, 'columntransformer__n_jobs': None, 'columntransformer__remainder': 'passthrough', 'columntransformer__sparse_threshold': 0.3, 'columntransformer__transformer_weights': None, 'columntransformer__transformers': [(('kbinsdiscretizer', KBinsDiscretizer(encode='ordinal'), ['BMI']))], 'columntransformer__verbose': False, 'columntransformer__verbose_feature_names_out': True, 'columntransformer__kbinsdiscretizer': KBinsDiscretizer(encode='ordinal'), 'columntransformer__kbinsdiscretizer__dtype': None, 'columntransformer__kbinsdiscretizer__encode': 'ordinal', 'columntransformer__kbinsdiscretizer__n_bins': 5, 'columntransformer__kbinsdiscretizer__random_state': None, 'columntransformer__kbinsdiscretizer__strategy': 'quantile', 'columntransformer__kbinsdiscretizer__subsample': 290000, 'logisticregression__C': 1.8, 'logisticregression__class_weight': None, 'logisticregression__dual': False, 'logisticregression__fit_intercept': True, 'logisticregression__intercept_scaling': 1, 'logisticregression__l1_ratio': None, 'logisticregression__max_iter': 500, 'logisticregression__multi_class': 'deprecated', 'logisticregression__n_jobs': None, 'logisticregression__penalty': 'l2', 'logisticregression__random_state': None, 'logisticregression__solver': 'lbfgs', 'logisticregression__tol': 0.0001, 'logisticregression__verbose': 0, 'logisticregression__warm_start': False}</pre>

In [] :	<pre>from sklearn.model_selection import RandomizedSearchCV import warnings warnings.filterwarnings('ignore') pipeline = create_pipeline_logistic() param_dist = { 'columntransformer__kbinsdiscretizer__strategy': ['uniform', 'quantile', 'kmeans'], 'columntransformer__kbinsdiscretizer__encode': ['ordinal', 'onehot', 'onehot-dense'], 'columntransformer__kbinsdiscretizer__n_bins': [5, 10, 20], 'logisticregression__C': [0.001, 0.01, 0.1, 1, 10, 100], 'logisticregression__penalty': ['none', 'l1', 'l2', 'elasticnet'], 'logisticregression__solver': ['newton-cg', 'lbfgs', 'liblinear', 'sag', 'saga'], 'logisticregression__max_iter': [100, 200, 300], 'logisticregression__tol': [1e-4, 1e-3, 1e-2]} search = RandomizedSearchCV(pipeline, param_dist, n_iter=20, random_state=0) search.fit(train[features], train[label]) search.best_params_ # ("search.best_score_:.3f")</pre>
Out[] :	<pre>(('logisticregression__solver': 'lbfgs', 'logisticregression__penalty': 'l2', 'logisticregression__max_iter': 200, 'logisticregression__C': 0.1, 'columntransformer__kbinsdiscretizer__strategy': 'quantile', 'columntransformer__kbinsdiscretizer__n_bins': 10, 'columntransformer__kbinsdiscretizer__encode': 'onehot-dense', '0.729')</pre>

Model Evaluation

สร้าง Baseline มาเป็นต้นแบบใน Model Eval

In [] :	<pre># Dummy 0/1 label feature engineer from sklearn.dummy import DummyClassifier # baseline = DummyClassifier(strategy="uniform", random_state=1234) # baseline = DummyClassifier(strategy="stratified", random_state=1234) baseline = DummyClassifier(strategy="most_frequent", random_state=1234) # baseline = DummyClassifier(strategy="constant", constant=0, random_state=1234) baseline.fit(train[features], train.Diabetes_binary) baseline_train_score = baseline.score(train[features], train.Diabetes_binary) baseline_validation_score = baseline.score(validation[features], validation.Diabetes_binary) logistic_train_score = search.score(train[features], train[label]) logistic_tuning_validation_score = search.score(validation[features], validation[label]) print(f"Baseline Train Score: {baseline_train_score:.3f} Baseline Validation Score: {baseline_validation_score:.3f}") print(f"Logistic Train Score: {logistic_train_score:.3f} Logistic Validation Score: {logistic_tuning_validation_score:.3f}")</pre>
Out[] :	<pre>Baseline Train Score: 0.590 Baseline Validation Score: 0.590 Logistic Train Score: 0.728 Logistic Validation Score: 0.730 Classification Report</pre>

In [] :	<pre>from sklearn.metrics import classification_report y_pred = search.best_estimator_.predict(validation[features]) print(classification_report(validation[label], y_pred))</pre>
Out[] :	<pre>precision recall f1-score support 0.0 0.74 0.70 0.72 3159 1.0 0.72 0.75 0.74 3159 accuracy 0.73 6318 macro avg 0.73 0.73 0.73 6318 weighted avg 0.73 0.73 0.73 6318</pre>

Cross Validation

In [] :	<pre>from sklearn.model_selection import cross_val_score cross_val_logistic_pipeline = create_pipeline_logistic() cross_val_logistic_pipeline.set_params(**search.best_params_) #time cv_scores = cross_val_score(cross_val_logistic_pipeline, train[features], train[label], cv = 5) print(f"Cross validation score: %.3f" % cv_scores.mean())</pre>
Out[] :	<pre>CPU times: total: 125 ms wall time: 28 ms cross validation score: 0.729</pre>

Model Calibration

In [] :	<pre>from sklearn.calibration import CalibratedClassifierCV calibration_logistic_pipeline = create_pipeline_logistic() calibration_logistic_pipeline.set_params(**search.best_params_) calibration_logistic_pipeline = CalibratedClassifierCV(estimator=calibration_logistic_pipeline, cv = 5) #time calibration_logistic_pipeline.fit(train[features], train[label]) train_score = calibration_logistic_pipeline.score(train[features], train[label]) valid_score = calibration_logistic_pipeline.score(validation[features], validation[label]) test_score = calibration_logistic_pipeline.score(test[features], test[label]) print(f"Train Score: {train_score:.3f}") print(f"Validation Score: {valid_score:.3f}") print(f"Test Score: {test_score:.3f}")</pre>
Out[] :	<pre>CPU times: total: 199 ms wall time: 259 ms Train Score: 0.738 Validation Score: 0.729 Test Score: 0.731</pre>

In [] :	<pre>sample = test.sample(100) cnt = 0 y_pred = calibration_logistic_pipeline.predict(sample[features]) sample_label = sample[label] for i in range(len(sample)): if (y_pred[i] != sample_label.iloc[i]): cnt += 1 cnt</pre>
Out[] :	<pre>77</pre>

In [] :	<pre>!writefile requirements.txt pandas>2.2 scikit-learn>1.5 imbalanced-learn>12.3 gradient>4.42.0 Overwriting requirements.txt</pre>
Out[] :	<pre>test[features]</pre>

```
calibration_logistic_pipeline = CalibratedClassifierCV(estimator=calibration_logistic_pipeline, cv = 5)
# Use calibration_logistic_pipeline.fit(train[features], train[label])

train_score = calibration_logistic_pipeline.score(train[features], train[label])
valid_score = calibration_logistic_pipeline.score(validation[features], validation[label])
test_score = calibration_logistic_pipeline.score(test[features], test[label])

print(f"Train Score: {train_score:.3f}")
print(f"Validation Score: {valid_score:.3f}")
print(f"Test Score: {test_score:.3f}")

CPU times: total: 109 ms
wall time: 259 ms
Train Score: 0.730
Validation Score: 0.729
Test Score: 0.731
```

```
In [ ]: sample = test.sample(100)
```


