Cho. BMI Smol	cer 0 cke 0 ctDiseaseorAttack 0 csActivity 0 cts 0 cyies 0 AlcoholConsump 0
Anyl NoDe Geni Men Phy: Dif: Sex Age Educ Inco	Mealthcare 0 DocbcCost 0 Hith 0 Hith 0 Walk 0 Contain 0 Contain 0 Contain 0 Contain 0
อธิบ ควา: ฟีเจฮ	etes.drop_duplicates(inplace= True) ายเพิ่ม: ในครั้งก่อนๆฟีเจอร์ทั้งหมดที่นำมาใช้นั้นเป็นฟีเจอร์ 5 ฟีเจอร์ที่หามาจากการทำ feature importance ที่มีค่าความสำคัญมากที่สุดเรียงลงมา ได้แก่ GenHlth, BMI, HeartDiseaseorAttack, DiffWalk, HvyAlcoholConsump ซึ่งเมื่อนำไปให้โมเดลเรียนรู้ ผลลัพธ์ที่ได้คือโม แม่นยำไม่มากพอที่จะนำไปใช้งานจริง ในครั้งนี้จึงเปลี่ยนมาใช้วิธีการดูความสัมพันธ์ของแต่ละฟีเจอร์ เพื่อดูว่าฟีเจอร์ใดบ้างที่มีความสัมพันธ์ต่อฟีเจอร์อื่นๆ โดยได้ฟีเจอร์อื่นๆมากที่สุดลงมา ได้แก่ GenHlth, HighBP, HighChol, BMI, Age, DiffWalk แล ร์เหล่านี้ไปให้โมเดลเรียนรู้ผลปรากฏว่าโมเดลมีความแม่นยำเพิ่มขึ้นในระดับนึง
diak	Diabetes_binary HighBP HighChol CholCheck BMI Smoker Stroke HeartDiseaseorAttack PhysActivity Fruits AnyHealthcare NoDocbcCost GenHth MentHith PhysHith DiffWalk Sex Age Education Income Diabetes_binary 1.000000 0.254318 0.194944 0.072523 0.205086 0.045504 0.099193 0.168213 -0.100404 -0.024805 0.025331 0.020048 0.276940 0.054153 0.156211 0.205302 0.032724 0.177263 -0.102686 -0.140659 -0.140659 -0.140659 -0.140659 -0.140659 -0.140659 -0.140659 -0.254318 1.00000 0.283813 0.00000 0.283813 0.111259 0.194218 0.074264 0.124558 0.201443 -0.104382 -0.019467 0.052044 0.002292 0.272784 0.037482 0.146666 0.211759 0.047119 0.339802 -0.112887 -0.140030 -0.140030 -0.140030 -0.140659 -0.1406
	DiffWalk 0.205302 0.211759 0.136045 0.049107 0.182556 0.108144 0.1069399 0.202657 -0.235719 -0.029932 0.017714 0.106225 0.446696 0.218733 0.466852 1.00000 -0.073405 0.209064 -0.169350 -0.299064 BMI 0.205086 0.194218 0.089734 0.042487 1.00000 -0.009294 0.011006 0.039820 -0.127780 -0.067424 -0.008519 0.045795 0.208351 0.068569 0.102768 0.182556 0.03902 -0.074433 -0.069097 HighChol 0.194944 0.283963 1.000000 0.094772 0.089734 0.074583 0.089375 0.176446 -0.063443 -0.026257 0.052363 0.003020 0.188139 0.050346 0.111008 0.136045 0.02859 0.263841 -0.050045 -0.062089 Age 0.177263 0.339802 0.263841 0.076952 0.128209 0.128209 0.023912 -0.087881 0.073515 0.14746
Hear	EdiseaseorAttack 0.168213 0.201443 0.176446 0.050086 0.039820 0.198814 1.000000 -0.073094 -0.006946 0.025987 0.021971 0.246328 0.052601 0.170335 0.202657 0.089828 0.223912 -0.082288 -0.122728 PhysHith 0.156211 0.144656 0.111008 0.040758 0.102768 0.140806 0.170335 -0.199307 -0.024441 0.002924 0.136421 0.516476 0.340191 1.000000 0.466852 -0.044433 0.095483 -0.127687 -0.240929 Stroke 0.099193 0.124558 0.089375 0.027955 0.011006 0.054414 1.000000 0.198814 -0.059306 -0.004486 0.013627 0.028613 0.169809 0.061996 0.140806 0.169339 0.003626 0.128209 -0.064178 -0.117108 CholCheck 0.072523 0.111259 0.094772 1.000000 0.042487 -0.003721 0.051965 0.052601 -0.105914 -0.052191
	Smoker 0.045504 0.074264 0.074264 0.074583 -0.003721 -0.009294 1.00000 0.054414 0.105169 -0.066869 -0.061731 -0.013963 0.037335 0.134894 0.077641 0.100447 0.108144 0.09650 0.135657 -0.095314 Sex 0.032724 0.047119 0.022859 -0.024332 0.03565 0.03626 0.089828 0.03516 -0.087688 -0.021221 -0.046507 -0.01519 -0.03386 -0.044433 -0.073405 1.00000 -0.031862 -0.04443 0.015956 0.113667 0.136997 AnyHealthcare 0.025331 0.052044 0.052363 0.115498 -0.0031963 0.013627 0.025987 0.025987 0.023659 1.000000 -0.227469 -0.03286 -0.043582 0.007744 0.01714 -0.021221 0.147465 0.11367 0.146144 NoDocbcCost 0.024805 -0.01467 -0.054198 0.037335 0.028613 0.024599 0.0227469 1.000000 0.149690
Hvy	Veggies -0.041734 -0.042994 -0.027399 -0.04054 -0.033029 -0.027180 0.135240 0.242941 0.020530 -0.04115 -0.04215 -0.045130 -0.063189 -0.063189 -0.066113 -0.003856 0.131624 0.125068 AlcoholConsump -0.065950 -0.014178 -0.019057 -0.02975 -0.058420 0.096048 -0.021347 -0.035561 0.023378 -0.028221 -0.006202 -0.001272 -0.055783 0.016852 -0.036860 -0.047655 0.00435 -0.041018 0.039132 0.071863 PhysActivity -0.100404 -0.104382 -0.063443 -0.027780 -0.06869 -0.059306 -0.073094 1.000000 0.125023 0.023959 -0.046440 -0.237511 -0.105914 -0.199307 -0.235719 0.033516 -0.087881 0.170931 0.165869 Education -0.140659 -0.140659 -0.0140659 -0.0140659 -0.062089 -0.064178 -0.064178 -0.082288 0.170931 0.084857
# fe feat # fe	ws × 22 columns **atures = ["GenHlth", "BMI", "HeartDiseaseorAttack", "DiffWalk", "HvyAlcoholConsump"] # Feature Importance ures = ["GenHlth", "BMI", "DiffWalk", "Age", "Smoker", "Sex"] # Correlation **atures = diabetes.drop("Diabetes_binary", axis=1).columns # All 1 = "Diabetes_binary"
from resa	สสอนและจัดการกับ Imbalance ด้วยการทำ Undersampling imblearn.under_sampling import RandomUnderSampler mpler = RandomUnderSampler(random_state=1234) etes[features], diabetes[label] = resampler.fit_resample(diabetes[features], diabetes[label]) etes.dropna(inplace=True) etes[features].head()
8 10 13	GenHth BMI DiffWalk Age Smoker Sex 5.0 30.0 1.0 9.0 1.0 0.0 3.0 25.0 0.0 13.0 1.0 1.0 4.0 28.0 1.0 11.0 0.0 0.0
Dial	2.0 33.0 0.0 6.0 0.0 0.0 3.0 21.0 0.0 10.0 0.0 10.0 0.0 Detes_binary .value_counts() Detes_binary .35097 35097
Trar แบ่ง	รรบอท e: count, dtype: int64 กักทู Data สุดข้อมูลเรียนรู้(Train), ชุดข้อมูลทดสอบ(Test), และชุดข้อมูลตรวจสอบ (Validation)
len (49)	n, test = train_test_split(diabetes, test_size = 0.3, stratify = diabetes_Diabetes_binary, random_state = 1234) , validation = train_test_split(test, test_size = 0.3, stratify = test_Diabetes_binary, random_state = 1234) train), len(test), len(validation) 135, 14741, 6318) สอบความสมดุลของชุดข้อมูลเรียนรู้
Dial 0.0 1.0 Name	
from from from	Pipeline ของโมเดล Logistic Regression และทำ Feature Engineer ใน Pipeline sklearn.pipeline import make_pipeline sklearn.compose import make_column_transformer sklearn.preprocessing import KBinsDiscretizer sklearn.linear_model import LogisticRegression create_pipline_logistic():
logi	<pre>return make_pipeline(make_column_transformer((KBinsDiscretizer(encode="ordinal", strategy="quantile"), ["BMI"]), remainder="passthrough"), LogisticRegression(max_iter=500)) stic_pipeline = create_pipline_logistic()</pre>
	Pipeline Pipeline Columntransformer: ColumnTransformer () kbinsdiscretizer remainder KBinsDiscretizer passthrough
%tim	LogisticRegression e logistic_pipeline.fit(train[features], train[label]) stic_train_score = logistic_pipeline.score(train[features], train[label]) stic_validation_score = logistic_pipeline.score(validation[features], validation[label])
prir prir CPU t Wall Train Valid	stic_validation_score = logistic_pipeline.score(validation[features], validation[label]) t(f"Train Score: %.3f" % logistic_train_score) t(f"Validation Score: %.3f" % logistic_validation_score) imes: total: 0 ns time: 199 ms Score: 0.714 ation Score: 0.716 separameter Tuning
. { 'mo	stic_pipeline.get_params() emory': None, ceps': [('columntransformer', ColumnTransformer(remainder='passthrough',
'10' '10' '20' '20'	erbose': False, plumntransformer': ColumnTransformer (remainder='passthrough',
'ca'	plumntransformer_transformer_weights': None, plumntransformer_transformer_veights': [('kbinsdiscretizer', KBinsDiscretizer(encode='ordinal'), ['BMI'])], plumntransformer_verbose': False, plumntransformer_verbose_feature_names_out': True, plumntransformer_kbinsdiscretizer': KBinsDiscretizer(encode='ordinal'), plumntransformer_kbinsdiscretizer_dtype': None, plumntransformer_kbinsdiscretizer_encode': 'ordinal', plumntransformer_kbinsdiscretizer_encode': 'ordinal', plumntransformer_kbinsdiscretizer_n_bins': 5, plumntransformer
'co 'co 'lo 'lo 'lo 'lo 'lo 'lo 'lo 'lo 'lo 'l	clumntransformerkbinsdiscretizerrandom_state': None, columntransformerkbinsdiscretizerstrategy': 'quantile', columntransformerkbinsdiscretizersubsample': 200000, columntransformerkbinsdiscretizersubsample': 200000, columntransformerkbinsdiscretizersubsample': 200000, columntransformerkbinsdiscretizersubsample': 200000, columntransformerkbinsdiscretizersubsample': 200000, columntransformerkbinsdiscretizerstrategy': 'quantile', columntransformerkbinsdiscretizersubstantegy': 'long columntransformerkbinsdiscretizerstrategy': 'quantile', columntransformerkbinsdiscretizerstrategy': 'long columntransformerkbinsdiscretizersubstantegy': 'lo
'10 '10 '10 '10 '10 '10 '10 '10	ogisticregression_max_iter': 500, ogisticregression_multi_class': 'deprecated', ogisticregression_n_jobs': None, ogisticregression_penalty': '12', ogisticregression_random_state': None, ogisticregression_solver': 'lbfgs', ogisticregression_tol': 0.0001, ogisticregression_verbose': 0, ogisticregression_warm_start': False}
impo warr pipe	<pre>rt warnings ings.filterwarnings("ignore") line = create_pipline_logistic() m_dist = { 'columntransformerkbinsdiscretizerstrategy': ['uniform', 'quantile', 'kmeans'], 'columntransformerkbinsdiscretizerencode': ['ordinal', 'onehot-dense'], 'columntransformerkbinsdiscretizerencode': ['ordinal', 'onehot-dense'],</pre>
} sear	<pre>'columntransformerkbinsdiscretizern_bins': [5, 10, 20], 'logisticregressionC': [0.001, 0.01, 0.1, 1, 10, 100], 'logisticregressionpenalty': ['none', 'll', 'l2', 'elasticnet'], 'logisticregressionsolver': ['newton-cg', 'lbfgs', 'liblinear', 'sag', 'saga'], 'logisticregressionmax_iter': [100, 200, 300], 'logisticregressiontol': [1e-4, 1e-3, 1e-2]</pre> ch = RandomizedSearchCV(pipeline, param_dist, n_iter=20, random_state=0)
sear ({ '.	ch.fit(train[features], train[label]) ch.best_params_, f"(search.best_score_:.3f)" Logisticregressiontol': 0.01, Logisticregressionsolver': 'lbfgs', Logisticregressionpenalty': '12', Logisticregressionmax_iter': 300, Logisticregressionmax_iter': 300, Logisticregressionc': 0.1, Logisticregressionc': 0.1, Logisticregressionk': 'quantile', Logisticregressionk': 'quantile', Logisticregressionk': 'quantile', Logisticregressionk': 'quantile', Logisticregressionk': 'quantile', Logisticregressionk': 'quantile', Logisticregression_k': 'quantile', Logis
' 0 Mo สร้าง	columntransformer_kbinsdiscretizer_encode': 'onehot-dense'}, (715') del Evaluation Baseline มาเปรียบเทียบกับโมเดลที่มีอยู่ www. ยังไม่ได้ทำ feature engineer
# ba # base # base	seline = DummyClassifier(strategy="uniform", random_state=1234) seline = DummyClassifier(strategy="stratified", random_state=1234) line = DummyClassifier(strategy="most_frequent", random_state=1234) seline = DummyClassifier(strategy="most_frequent", random_state=1234) seline = DummyClassifier(strategy="constant", constant=0, random_state=1234) line.fit(train[features], train.Diabetes_binary)
logi logi prir prir	<pre>line_train_score = baseline.score(train[features], train.Diabetes_binary) line_validation_score = baseline.score(validation[features], validation.Diabetes_binary) stic_tuning_train_score = search.score(train[features], train[label]) stic_tuning_validation_score = search.score(validation[features], validation[label]) t(f"Baseline Train Score: {baseline_train_score:.3f}\nBaseline Validation Score: {baseline_validation_score:.3f}") t(f"Logistic Train Score: {logistic_tuning_train_score:.3f}\nLogistic Validation Score: {logistic_tuning_validation_score:.3f}") ine Train Score: 0.500</pre>
Logis Logis Clas from	ine Validation Score: 0.500 tic Train Score: 0.716 tic Validation Score: 0.716 sification Report a sklearn.metrics import classification_report ed = search.best_estimatorpredict(validation[features]) t(classification_report(validation[label], y_pred))
ma	precision recall f1-score support 0.0 0.73 0.69 0.71 3159 1.0 0.71 0.74 0.72 3159 ccuracy cro avg 0.72 0.72 0.72 6318 ted avg 0.72 0.72 0.72 6318
cros	s Validation sklearn.model_selection import cross_val_score s_val_logistic_pipeline = create_pipline_logistic() s_val_logistic_pipeline.set_params(**search.best_params_) see cv_scores = cross_val_score(cross_val_logistic_pipeline, train[features], train[label], cv = 5)
CPU t Wall cross	t(f"cross validation score: %.3f" % cv_scores.mean()) imes: total: 62.5 ms time: 288 ms validation score: 0.715 cl Calibration a sklearn.calibration import CalibratedClassifierCV
cali cali %tim trai vali	bration_logistic_pipeline = create_pipline_logistic() bration_logistic_pipeline.set_params(**search.best_params_) bration_logistic_pipeline = CalibratedClassifierCV(estimator=calibration_logistic_pipeline, cv = 5) calibration_logistic_pipeline.fit(train[features], train[label]) n_score = calibration_logistic_pipeline.score(train[features], train[label]) d_score = calibration_logistic_pipeline.score(validation[features], validation[label]) _score = calibration_logistic_pipeline.score(test[features], test[label])
prir prir CPU t Wall Train	t(f"Train Score: {train_score:.3f}") t(f"Validation Score: {valid_score:.3f}") t(f"Test Score: {test_score:.3f}") imes: total: 93.8 ms time: 297 ms Score: 0.716 ation Score: 0.715 Score: 0.715
cnt y_pr samp	<pre>le = test.sample(100)</pre>
pand	<pre>cnt += 1 itefile requirements.txt as=2.2.2 it-learn=1.5.1 cont += 1</pre>
grac	lanced-learn=0.12.3 io=4.42.0 riting requirements.txt [features] GenHith BMI DiffWalk Age Smoker Sex 43 3.0 37.0 0.0 8.0 1.0 1.0
2073 676 1214 776	78 2.0 24.0 0.0 5.0 1.0 1.0 61 4.0 28.0 1.0 13.0 1.0 1.0
1538 713 480 2074	27 4.0 32.0 1.0 11.0 1.0 1.0 17 2.0 24.0 0.0 8.0 0.0 0.0
Exp	97 4.0 28.0 0.0 7.0 1.0 1.0 I rows × 6 columns ort "Artifacts" ret joblib ib.dump(calibration_logistic_pipeline, "pipeline.joblib")
%%wrimpo	peline.joblib'] itefile app.py rt pandas as pd rt joblib rt numpy as np rt gradio as gr
feat pipe clas	ures = ["GenHlth", "BMI", "DiffWalk", "Age", "Smoker", "Sex"] line = joblib.load("pipeline.joblib") ses = ["Based on the assessment from the information you provided, we found that you have no risk of developing diabetes.", "Based on the assessment from the information you provided, we found that you have risk of developing diabetes.", "Based on the assessment from the information you provided, we found that you have risk of developing diabetes.", "Based on the assessment from the information you provided, we found that you have risk of developing diabetes.", "Based on the assessment from the information you provided, we found that you have risk of developing diabetes.", "Based on the assessment from the information you provided, we found that you have risk of developing diabetes.", "Based on the assessment from the information you provided, we found that you have risk of developing diabetes.", "Based on the assessment from the information you provided, we found that you have risk of developing diabetes.", "Based on the assessment from the information you provided, we found that you have risk of developing diabetes.", "Based on the assessment from the information you provided, we found that you have risk of developing diabetes.", "Based on the assessment from the information you provided, we found that you have risk of developing diabetes.", "Based on the assessment from the information you provided, we found that you have risk of developing diabetes.", "Based on the assessment from the information you provided, we found that you have risk of developing diabetes.", "Based on the assessment from the information you provided, we found that you have risk of developing diabetes.", "Based on the assessment from the information you provided, we found that you have risk of developing diabetes.", "Based on the assessment from the information you provided, we found that you have risk of developing diabetes.", "Based on the assessment from the information you provided, we found that you have risk of developing diabetes."
	<pre>if (bmi >= 35): bmi = 5 elif (bmi >= 30): bmi = 4 elif (bmi >= 27): bmi = 3 elif (bmi >= 24): bmi = 2 elif (bmi >= 23):</pre>
	<pre>bmi = 1 sample = dict() sample["GenHlth"] = genhlth sample["BMI"] = bmi sample["DiffWalk"] = diffWalk sample["Age"] = age sample["Smoker"] = smoker sample["Smoker"] = sex</pre> sample = pd.DataFrame([sample])
def	<pre>y_pred = pipeline.predict_proba(sample)[0] y_pred = dict(zip(classes, y_pred)) return y_pred cast_string_to_float(input): result = float(input) return result gr.Blocks(theme=gr.themes.Soft()) as demo:</pre>
	options_genHealt = ["Excellent", "Very good", "Good", "Fair", "Poor"] values_genHealt = [1, 2, 3, 4, 5] options_age = ["18-24", "25-29", "30-34", "35-39", "40-44", "45-49", "50-54", "55-59", "60-64", "65-69", "70-74", "75-79", ">80"] values_age = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13] options_bmi = ["13-23", "24-26", "27-29", "30-34", ">35"] values_bmi = [1, 2, 3, 4, 5] options_yes_no = ["Yes", "No"] values_yes_no = [1, 0] options_sex = ["Male", "Female"]
	gr.Markdown("## Diabetes Indicator", elem_classes="markdown-label") gr.Markdown(""" <h3>The dataset used for training the model comes from https://www.kaggle.com/datasets/alexteboul/diabetes-health-indicators-dataset""") gr.Markdown(""" <div style="background-color: #f8d7da; border: 1px solid #f5c6cb; border-radius: 5px; padding: 15px; color: #721c24;"></div></h3>
	<pre></pre>
	<pre>style="color: #333333;">You feel very energetic and strong, with no health issues. You can perform physically demanding activities without any problems. style="color: #333333;">You exercise regularly and maintain a healthy diet. </pre>
	<pre> </pre> <pre> <pre></pre> <pre><pre></pre> <pre></pre> <</pre></pre>
	<pre><pre>style="color: #333333;">You have health issues that occasionally interfere with your daily activities.</pre> <pre></pre> <pre><pre></pre></pre></pre> <pre><pre><pre><pre><pre><pre><pre><</pre></pre></pre></pre></pre></pre></pre>
	<pre>""") with gr.Row(): smoker = gr.Radio(choices=list(zip(options_yes_no, values_yes_no)), label="Smoker", info="Have you smoked at least 100 cigarettes in your entire life? [Note: 5 packs = 100 cigarettes]", elem_classes="gr-radio") sex = gr.Radio(choices=list(zip(options_sex, values_sex)), label="Sex", info="What is your sex?", elem_classes="gr-radio") with gr.Row(): height = gr.Number(label="Height", info="What is your height?", value=1)</pre>
	<pre>height = gr.Number(label="Height", info="What is your height?", value=1) weight = gr.Number(label="Weight", info="What is your weight?", value=1) age = gr.Radio(choices=list(zip(options_age, values_age)), label="Age", info="Select your age", elem_classes="gr-radio") diffWalk = gr.Radio(choices=list(zip(options_yes_no, values_yes_no)), label="Difficult walk", info="Do you have serious difficulty walking or climbing stairs?", elem_classes="gr-radio") predict_btn = gr.Button("Predict", variant="primary", elem_classes="gr-button-primary") Diabetes_binary = gr.Label(label="Diabetes_binary", elem_classes="gr-label") inputs = [genhlth, diffWalk, age, smoker, sex, height, weight] output = [Diabetes_binary]</pre>
<pre>if _ Overw . %rur</pre>	<pre>predict_btn.click(predict, inputs=inputs, outputs=output)name == "main": demo.launch() riting app.py app.py ng on local URL: http://127.0.0.1:7908</pre>
	ng on local URL: http://127.0.0.1:7908 eate a public link, set `share=True` in `launch()`. ERROR The requested URL could not be retrieved
	lowing error was encountered while trying to retrieve the URL: http://127.0.0.1:7908/ Access Denied. It control configuration prevents your request from being allowed at this time. Please contact your service provider if you feel this is incorrect. In additional configuration is webmaster.
Your ca	d Sun, 08 Sep 2024 05:03:14 GMT by htmlpdf-proxy-us (squid/4.10)

กลุ่ม Diabetes Prediction

1. 6610402205 นายรักษิต รุ่งรัตนไชย หมู่ 1

วัตถุประสงค์ของระบบต้นแบบ

2. 6610402132 นายบวรรัตน์ ตั้งนรารัชชกิจ หมู่ 1 3. 6610401985 นายไชยวัตน์ หนูวัฒนา หมู่ 1