

**Python programming  
in medical data analysis**  
**FINAL PROJECT**

B323111023 鄭宇翔  
B313111055 洪宣蓉  
B313111056 潘湘昀

# OUTLINE

**01 Introduction**

**02 Method**

**03 Datasets**

**04 Analysis**

**05 Data processing**

**06 Model-RNN SVM**

**07 Selected Model PCA->KNN**

**08 Data splitting-Train/Test**

**09 Prediction and Result**

# INTRODUCTION

**Machine learning identifies girls with central precocious puberty based on multisource data**



## **Central precocious puberty(CPP)**

A disease caused by premature activation of the hypothalamic-pituitary-gonadal (HPG) axis with clinical pubertal symptoms among girls under 8 years old and boys under 9 years of age.



## **CPP diagnosis**

We usually use gonadotropin-releasing hormone (GnRH) stimulation test ,which is the gold standard for CPP diagnosis.



## **Objective**

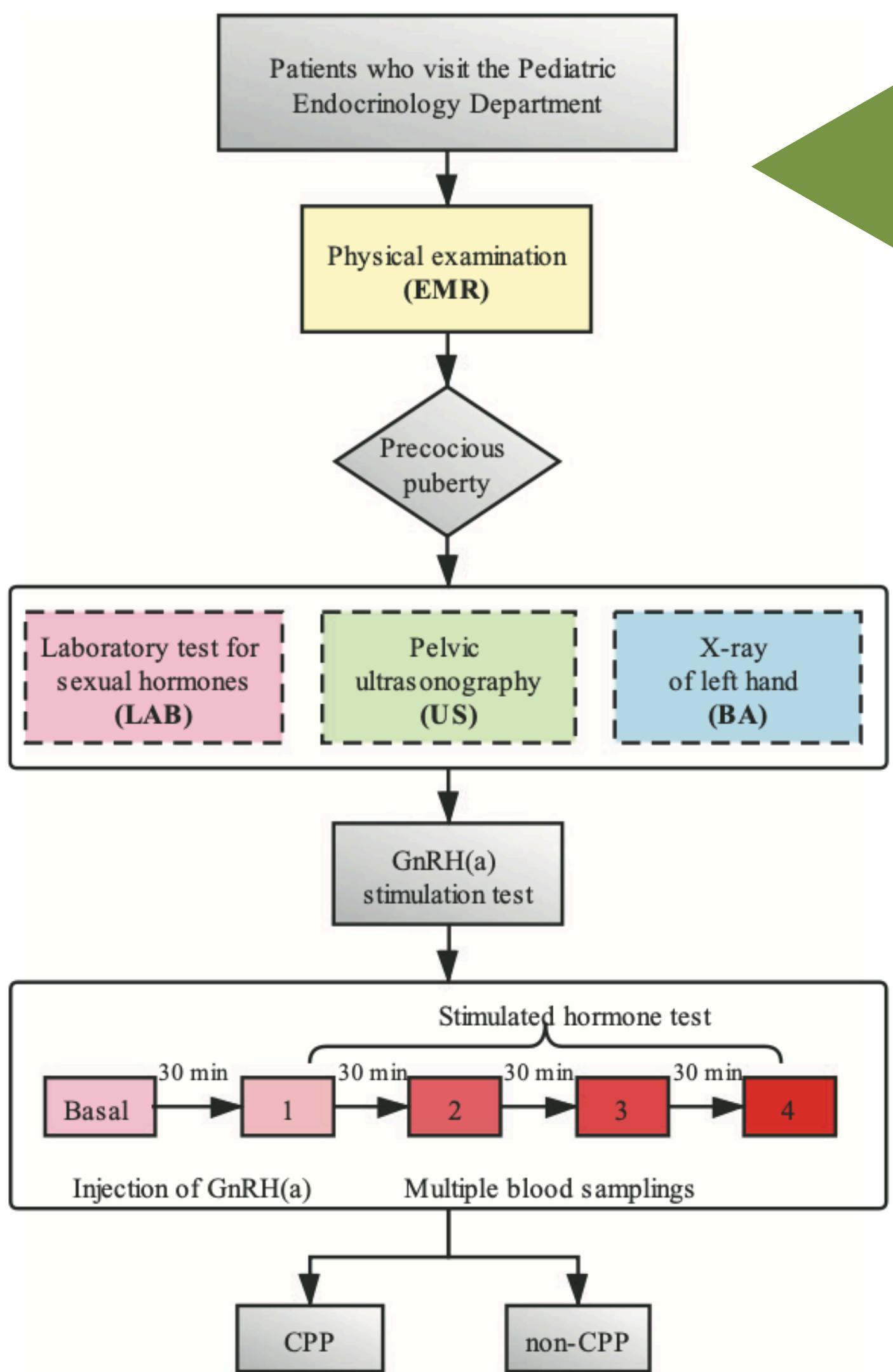
The study aimed to develop simplified diagnostic models for identifying girls with CPP, without GnRH stimulation test.

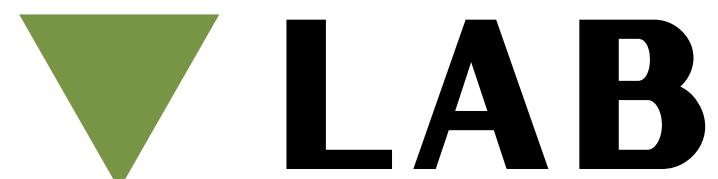
# METHOD

## CPP diagnosis pathway

the data obtained using these different approaches from the same patient were considered as data from different sources.

## Data sources:EMR, LAB, BA, and US





# Laboratory test for sexual hormones

Python-programming-in-medical-data-analysis / shared\_data / test.csv

shtthfkup Add files via upload

Preview Code Blame 2524 lines (2524 loc) · 95.8 KB Code 55% faster with GitHub Copilot

Search this file

1	LH	FSH	IGF-1	IGFBP-3	GH	TTE	PRL	E2	Age	IS_CPP
2	0.07	0.3	403.0	5.66	3.16	0.35	14.73	66.54	8	0
3	0.07	2.36				0.35	3.96	130.79	5	0
4	0.09	4.04	232.0	3.98	0.29	0.79	6.49	154.86	5	0
5	0.07	1.87	188.0	3.94	2.55	0.35	4.59	60.54	8	0
6	0.07	0.8	356.0	5.48	1.67	0.86	2.99	72.44	8	1
7	0.39	4.24				0.77	14.19	81.45	7	0
8	1.89	2.01	288.0	4.24	9.39	1.66	7.24	83.61	7	1
9	0.12	1.74				1.11	13.31	116.04	7	0
10	0.1	3.23	204.0		7.75	0.35	3.01	121.16	8	0

## OUTLINE

-10 variables

(LH, FSH, IGF-1,  
IGFBP-3, GH, TTE, PRL,  
E2, Age, IS\_CPP)

-2524 rows data

-CPP 1153  
Non-CPP 1370

<b>FEATURE</b>	<b>NON-CPP</b>	<b>CPP</b>	<b>P-VALUE</b>
<b>LH (IU/L)</b>	0.11	0.78	< 0.001
<b>FSH (IU/L)</b>	1.83	2.94	< 0.001
<b>IGF-1 (ng/L)</b>	232.25	295.19	< 0.001
<b>IGFBP-3 (μg/mL)</b>	4.69	0.78	< 0.001
<b>GH (ng/mL)</b>	3.18	4.13	< 0.001
<b>TTE (nmol/L)</b>	0.80	0.90	< 0.001
<b>PRL (ng/mL)</b>	9.42	8.765	< 0.001
<b>E2 (ng/mL)</b>	106.03	120.79	< 0.001
<b>AGE (years)</b>	7.05	7.47	< 0.001

# TEST\_BA

## X-ray of left hand

Python-programming-in-medical-data-analysis / shared\_data / test\_ba.csv

shthfkup Add files via upload

Preview Code Blame 1613 lines (1613 loc) · 81.8 KB Code 55% faster with GitHub Copilot

Search this file

1	LH	FSH	IGF-1	IGFBP-3	GH	TTE	PRL	E2	Age	IS_CPP	BoneAge	Age_ratio
2	0.07	0.47	217.0	3.81	6.12		6.83	50.75	8	0	9.0	1.125
3	1.22	3.48	180.0		0.24				7	1	10.0	1.428571429
4	0.85	3.45	421.0			0.35			8	1	11.0	1.375
5	0.07	1.09	312.0		3.02	0.76			7	0	8.0	1.142857143
6	2.13	2.87	429.0	4.78	7.38	0.35			7	1	9.0	1.285714286
7	0.19	2.92	298.0	5.78	5.21		4.97	188.18	8	1	10.5	1.3125
8	0.07	0.97	278.0	4.96	0.09	0.48	6.45	116.07	6	0	9.0	1.5
9	0.94	0.63			3.4	0.35			8	1	11.0	1.375
10	0.07	0.96	238.0	4.72	2.74	0.81	9.75	58.58	6	0	8.25	1.375
11	0.07	0.85	285.0	4.98	3.18	0.92	3.82	81.04	8	1	11.0	1.375
12	0.07	1.19	231.0	5.49	1.91	0.96	11.2	119.28	8	0	8.25	1.03125
13	0.07	1.37	222.0	5.35	1.8	0.35			7	0		1.25

## OUTLINE

-2 more variables  
(BoneAge, Age\_ratio)

-Age\_ratio= BoneAge/Age

-1613 rows data

-CPP 715  
Non-CPP 897

<b>FEATURE</b>	<b>NON-CPP</b>	<b>CPP</b>	<b>P-VALUE</b>
<b>BoneAge (years)</b>	897	715	< 0.001
<b>Age Ratio</b>	1.25	1.32	< 0.001

# TEST\_EMR

## Physical Examination

[Python-programming-in-medical-data-analysis / shared\\_data / test\\_emr.csv](#)

shtthfkup Add files via upload

Preview Code Blame 1676 lines (1676 loc) · 114 KB Code 55% faster with GitHub Copilot

Search this file

1	LH	FSH	IGF-1	IGFBP-3	GH	TTE	PRL	E2	Age	IS_CPP	abnormal	height	weight	breast	core	vulva	pubes	pigmentation
2	0.13	1.18			3.35			135.0	7	0	6.0	126.7	25.0					
3	0.26	2.19				0.69		124.0	8	0	6.0	128.3	32.0					
4	0.32	2.55				0.69		259.0	8	0	3.0	126.3	21.0					
5	0.1	1.62				0.69		167.0	6	0	0.2	120.7	22.0					
6	0.1	0.75	341.0			3.31		137.0	7	0	3.0	138.8	38.0					
7	0.1	0.82				0.42	0.69	124.0	8	0	18.0	131.8	26.0					
8	0.1	1.52					0.69	127.0	6	0	3.0	128.9	23.0					
9	1.15	4.19				0.69		157.0	6	1	24.0	119.7	23.5					
10	0.18	1.36				0.69		136.0	8	1	3.0	136.5	34.0					
11	0.73	3.84				0.69		162.0	5	1	3.0	109.3	18.0					
12	0.99	3.55				0.69		232.0	8	1	6.0	142.0	32.0					
13	0.1	1.52				0.69		196.0	7	0	6.0	131.9	25.0					

## OUTLINE

- 8 more variables
- 1676 rows data
- CPP 794
- Non-CPP 881
- Abnormal =abnormal duration
- Breast formal: tanner stage 1-5
- Vulva/Pubes formal: tanner stage 1-3
- Core/Pigmentation formal: Yes/No

<b>FEATURE</b>	<b>NON-CPP</b>	<b>CPP</b>	<b>P-VALUE</b>
<b>Duration (years)</b>	8.15	9.95	< 0.001
<b>Height (cm)</b>	126.11	129.95	< 0.001
<b>Weight (kg)</b>	26.04	28.33	< 0.001
<b>BMI (kg/m^2)</b>	16.22	16.66	< 0.001
<b>Breast tanner stage 1</b>	72	47	
<b>stage 2</b>	418	197	
<b>stage 3</b>	347	402	< 0.001
<b>stage 4</b>	43	143	
<b>stage 5</b>	1	5	

FEATURE	NON-CPP	CPP	P-VALUE
<b>Core YES</b>	724	684	0.027
<b>Core NO</b>	157	110	
<b>Vulva tanner stage 1</b>	833	722	0.009
<b>stage 2</b>	44	61	
<b>stage 3</b>	4	11	0.002
<b>Pubes tanner stage 1</b>	837	773	
<b>stage 2</b>	43	17	0.137
<b>stage 3</b>	1	4	
<b>Pigmentation YES</b>	57	38	0.137
<b>Pigmentation NO</b>	824	756	

# TEST\_US

## Pelvic Ultras Onography

Python-programming-in-medical-data-analysis / shared\_data / test\_us.csv 

shtthfkup Add files via upload

Preview Code Blame 2008 lines (2008 loc) · 160 KB Code 55% faster with GitHub Copilot

Q Search this file

1	LH	FSH	IGF-1	IGFBP-3	GH	TTE	PRL	E2	Age	IS_CPP	uterus1	uterus2	uterus3	lovary1	lovary2	lovary3	rovary1	rovary2	rovary3
2	0.07	0.47	217.0	3.81	6.12		6.83	50.75	8	0	23.0	19.0	13.0	33.0	10.0	14.0	35.0	9.0	15.0
3	1.22	3.48	180.0			0.24			7	1	32.0	22.0	18.0	25.0	18.0	12.0	30.0	19.0	14.0
4	0.85	3.45	421.0			0.35			8	1	24.0	20.0	12.0	27.0	14.0	25.0	13.0		
5	0.07	1.09	312.0		3.02	0.76			7	0	15.0	14.0	10.0	19.0	9.0	10.0	20.0	11.0	11.0
6	0.07	1.74	180.0		2.43	0.53			6	0	20.0	9.0	10.0	26.0	15.0	15.0	19.0	11.0	13.0
7	2.13	2.87	429.0	4.78	7.38	0.35			7	1	31.0	19.0	14.0	16.0	24.0	15.0	26.0	11.0	25.0
8	0.19	2.92	298.0	5.78	5.21		4.97	188.18	8	1	20.5	13.0	9.3	26.0	15.0	14.0	16.0	11.0	10.5
9	0.07	0.97	278.0	4.96	0.09	0.48	6.45	116.07	6	0	21.0	7.0	13.0	25.0	13.0	15.0	23.0	9.0	11.0
10	0.94	0.63			3.4	0.35			8	1	12.0	20.0	17.0	18.0	14.0	15.0	17.0	13.0	16.0
11	0.07	0.96	238.0	4.72	2.74	0.81	9.75	58.58	6	0	17.0	6.0	4.0	23.0	12.0	18.0	25.0	12.0	21.0
12	0.07	0.85	285.0	4.98	3.18	0.92	3.82	81.04	8	1	19.0	9.0	14.0	25.0	13.0	18.0	20.0	11.0	15.0
13	0.07	1.19	231.0	5.49	1.91	0.96	11.2	119.28	8	0	15.0	8.0	13.0	16.0	10.0	17.0	23.0	11.0	20.0
14	0.07	1.37	222.0	5.35	1.8	0.35			7	0	18.0	8.7	5.0	38.0	13.0	7.0	11.8	8.3	5.3
15	0.07	11	299.0	5.1	0.62	0.69			7	0	26.0	15.0	10.0	26.0	10.0	11.0	34.0	10.0	11.0

## OUTLINE

-9 more variables

-2000 rows data

-CPP 923

Non-CPP 1084

-uterine volume

=uterus1 x uterus2 x uterus3  
(length x width x height)

-(L)ovary volume

=lovary1 x lovary2 x lovary3  
(length x width x height)

-(R)ovary volume

=rovary1 x rovary2 x rovary3  
(length x width x height)

```

3 merged_df['uterus_vol'] = merged_df['uterus1'] * merged_df['uterus2'] * merged_df['uterus3']
4 merged_df['lovary_vol'] = merged_df['lovary1'] * merged_df['lovary2'] * merged_df['lovary3']
5 merged_df['rovary_vol'] = merged_df['rovary1'] * merged_df['rovary2'] * merged_df['rovary3']

```

<b>FEATURE</b>	<b>NON-CPP</b>	<b>CPP</b>	<b>P-VALUE</b>
<b>Uterine Volume (mL)</b>	1.57	2.69	< 0.001
<b>Left Ovarian Volume (mL)</b>	2.04	2.46	< 0.001
<b>Right Ovarian Volume (mL)</b>	1.94	2.29	< 0.001

# ANALYSIS

## DATA PROCESSING

-Take the logarithm of 8 variables

(LH,FSH,IGF-1,GH,TTE,PRL,E2)

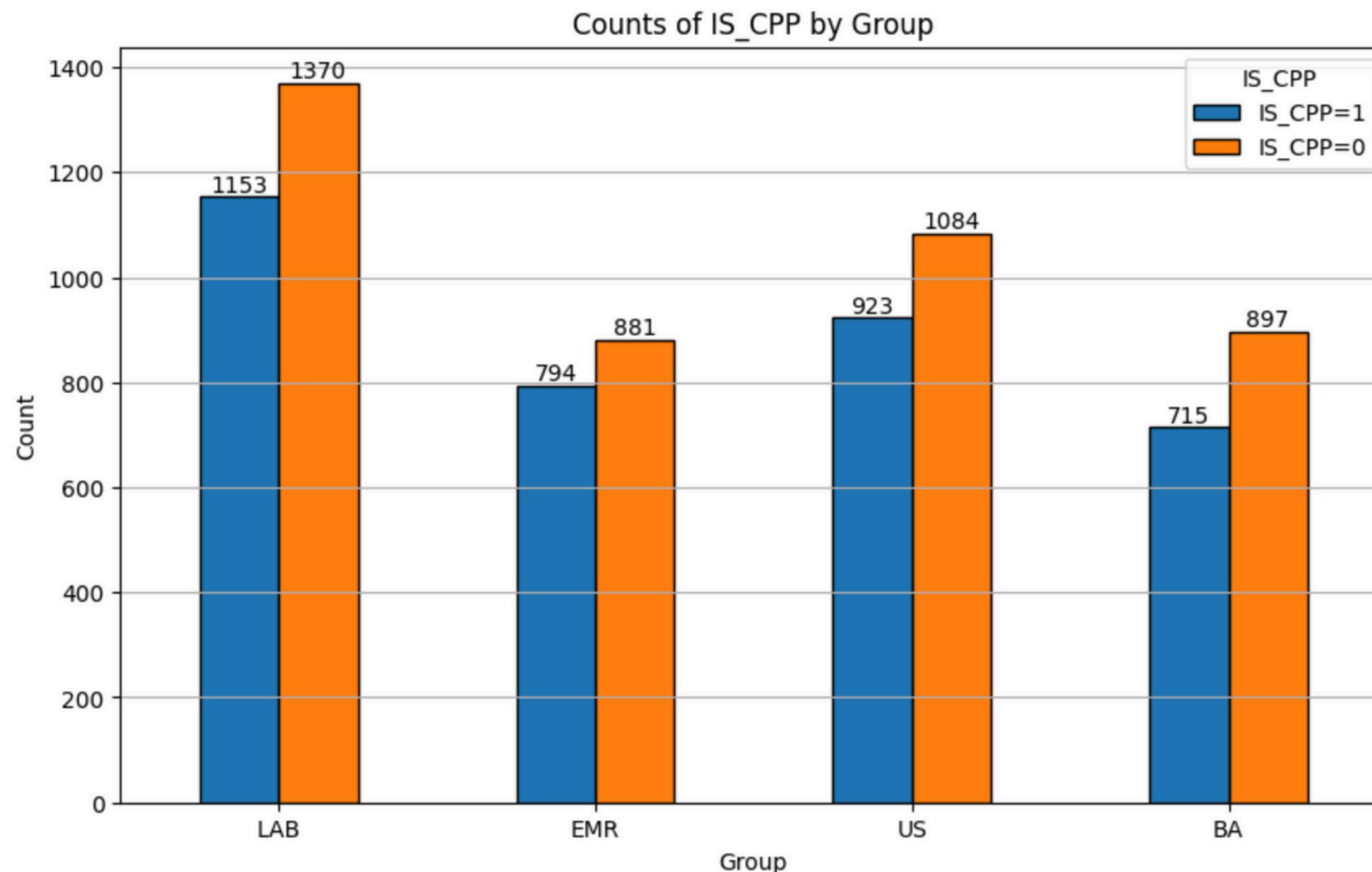
-Get BMI,Bone/Age difference, Uterus/(L)Ovary/(R)Ovary Volume

```
[17] 1 merged_df['LH_log'] = merged_df['LH'].apply(np.log)
    2 merged_df['FSH_log'] = merged_df['FSH'].apply(np.log)
    3 merged_df['IGF-1_log'] = merged_df['IGF-1'].apply(np.log)
    4 merged_df['GH_log'] = merged_df['GH'].apply(np.log)
    5 merged_df['TTE_log'] = merged_df['TTE'].apply(np.log)
    6 merged_df['PRL_log'] = merged_df['PRL'].apply(np.log)
    7 merged_df['E2_log'] = merged_df['E2'].apply(np.log)
```

```
[18] 1 merged_df['BMI'] = merged_df['height']**2/merged_df['weight']
    2 merged_df['Bone_age_diff'] = merged_df['BoneAge'] - merged_df['Age']
    3 merged_df['uterus_vol'] = merged_df['uterus1'] * merged_df['uterus2'] * merged_df['uterus3']
    4 merged_df['lovary_vol'] = merged_df['lovary1'] * merged_df['lovary2'] * merged_df['lovary3']
    5 merged_df['rovary_vol'] = merged_df['rovary1'] * merged_df['rovary2'] * merged_df['rovary3']
```

# ANALYSIS

## COUNTS OF IS\_CPP



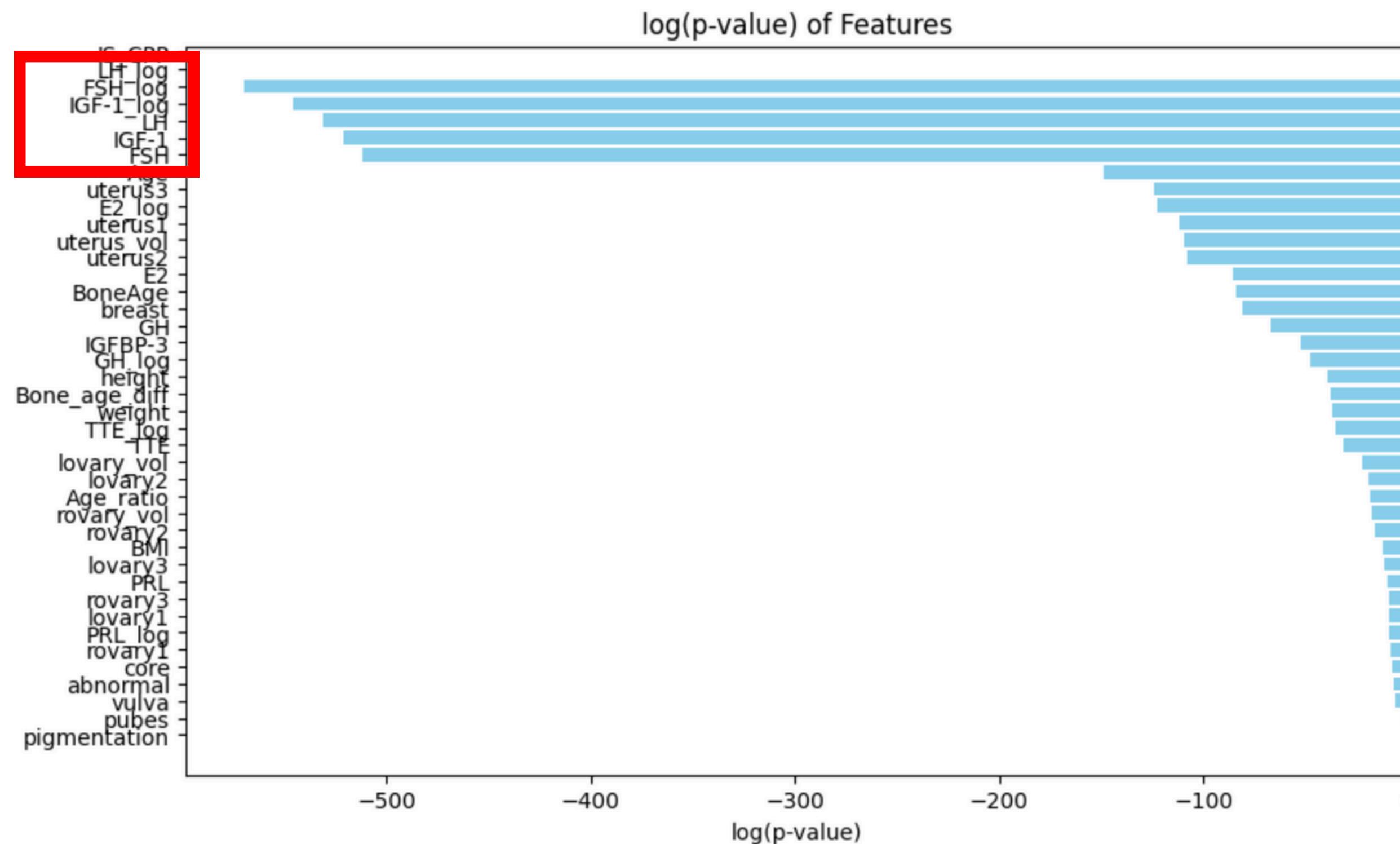
# ANALYSIS

## COUNTS OF IS\_CPP

```
[45] 1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 counts = {
5     "LAB": [len(df_LAB[df_LAB["IS_CPP"] == 1]), len(df_LAB[df_LAB["IS_CPP"] == 0])],
6     "EMR": [len(df_EMR[df_EMR["IS_CPP"] == 1]), len(df_EMR[df_EMR["IS_CPP"] == 0])],
7     "US": [len(df_US[df_US["IS_CPP"] == 1]), len(df_US[df_US["IS_CPP"] == 0])],
8     "BA": [len(df_BA[df_BA["IS_CPP"] == 1]), len(df_BA[df_BA["IS_CPP"] == 0])]
9 }
10
11
12 counts_df = pd.DataFrame(counts, index=["IS_CPP=1", "IS_CPP=0"])
13 counts_df = counts_df.T
14
15 ax = counts_df.plot(kind='bar', figsize=(10, 6), edgecolor='black', rot=0)
16 plt.title('Counts of IS_CPP by Group')
17 plt.xlabel('Group')
18 plt.ylabel('Count')
19 plt.legend(title='IS_CPP')
20 plt.grid(axis='y')
21
22 for container in ax.containers:
23     ax.bar_label(container, label_type='edge')
24
25 plt.show()
```

# ANALYSIS

## LOG(P-VALUE) OF FEATURES



# ▼ DATA PROCESSING

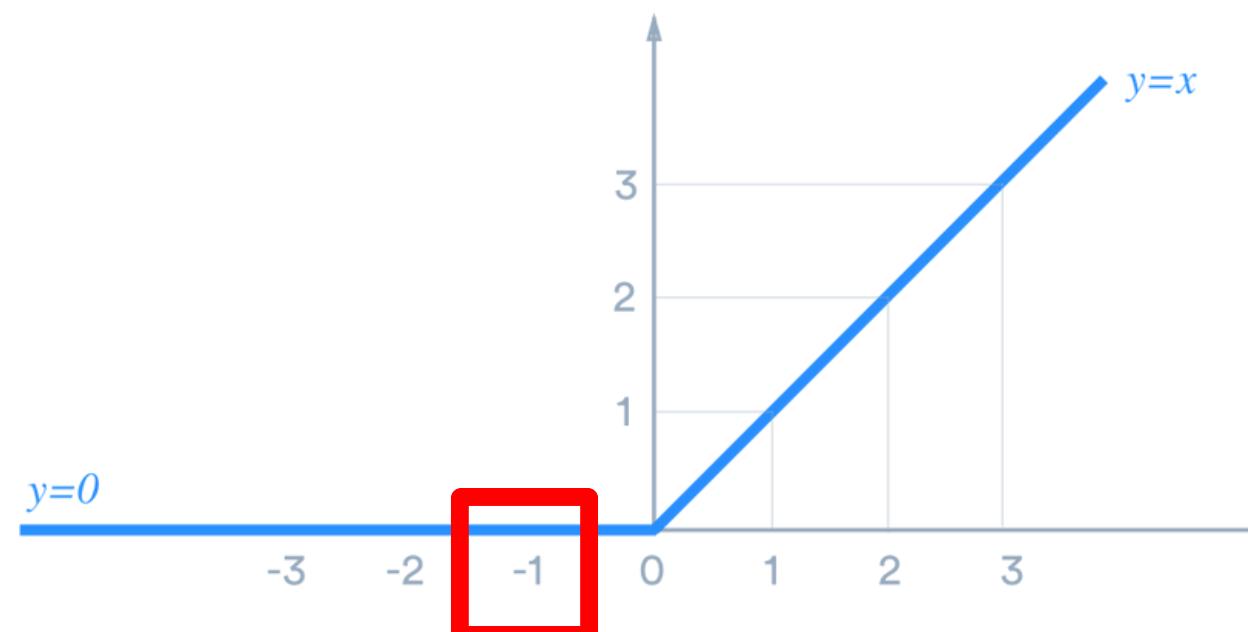
## Min Max Normalize (0~1)

```
[19] 1 numeric_cols = merged_df.select_dtypes(include=np.number)
2 min_vals = numeric_cols.min()
3 max_vals = numeric_cols.max()
4 normalized_df = (numeric_cols - min_vals) / (max_vals - min_vals)
5 normalized_df = pd.concat([normalized_df, merged_df["data_class"]], axis=1)
```

## Drop NaN

```
[20] 1 df = normalized_df.fillna(-1)
2 df
```

## To fit Relu Activation



# MODEL SVM/RNN



## SVM of Age to LH\_log

Accuracy: 0.80



## RNN of Exclusive data

Batch size = 8

Epoch = 150

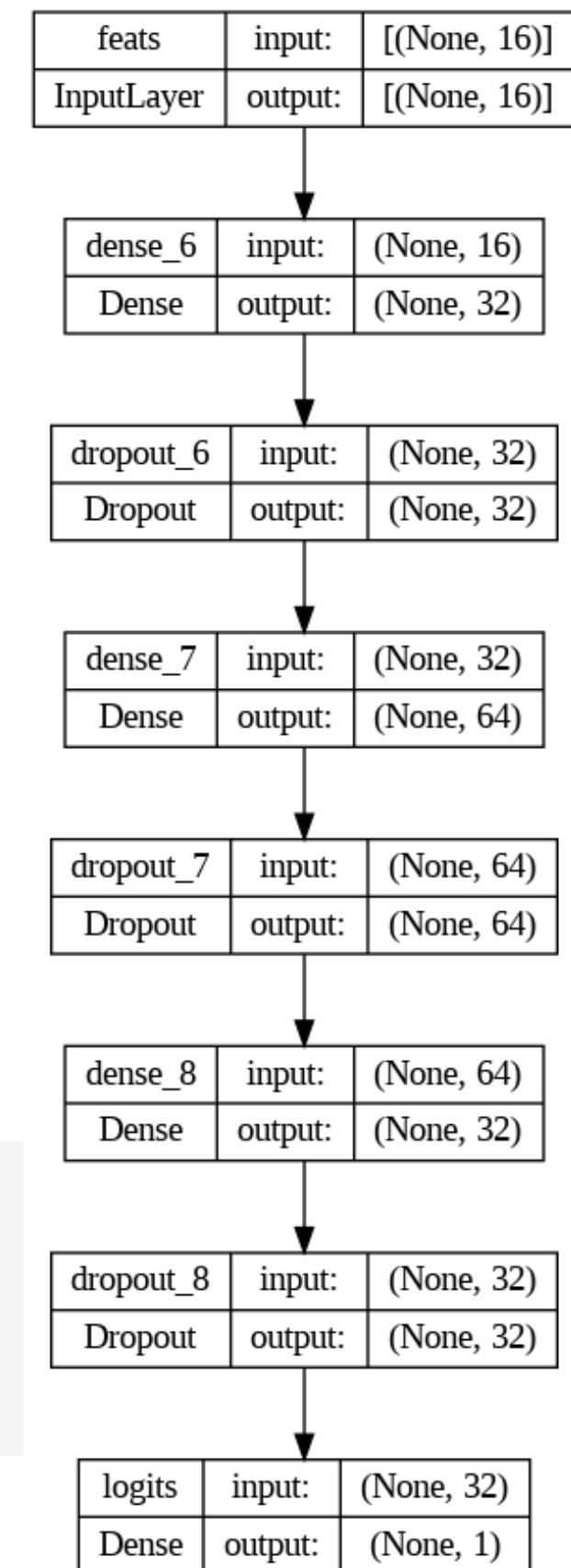
Valid. Data = 0.2

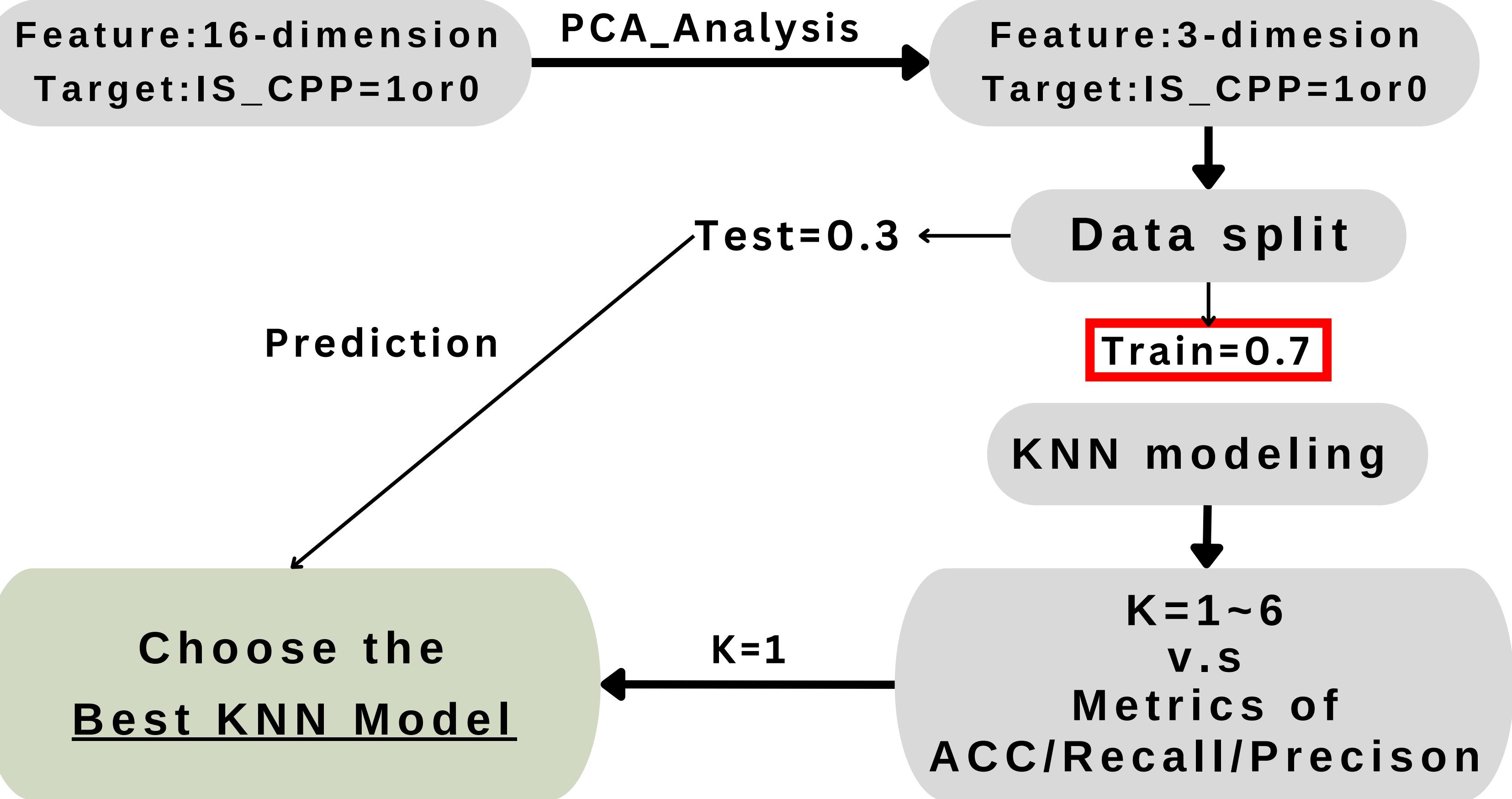
EMR Accuracy: 0.71

BA Accuracy: 0.84

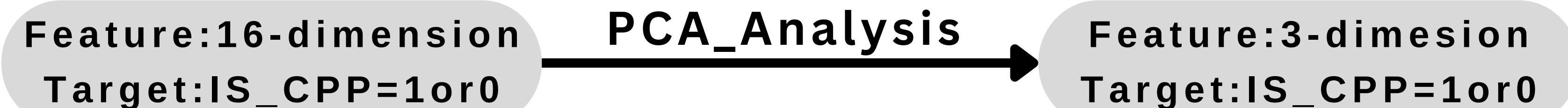
US Accuracy: 0.81

```
inputs = layers.Input(shape=(len(feature_columns), ), name="feats")
x = layers.Dense(32, activation="relu")(inputs)
x = layers.Dropout(0.1)(x)
x = layers.Dense(64, activation="relu")(x)
x = layers.Dropout(0.1)(x)
x = layers.Dense(32, activation="sigmoid")(inputs)
logits = layers.Dense(1, name="logits")(x)
```





# ▼ PCA->KNN

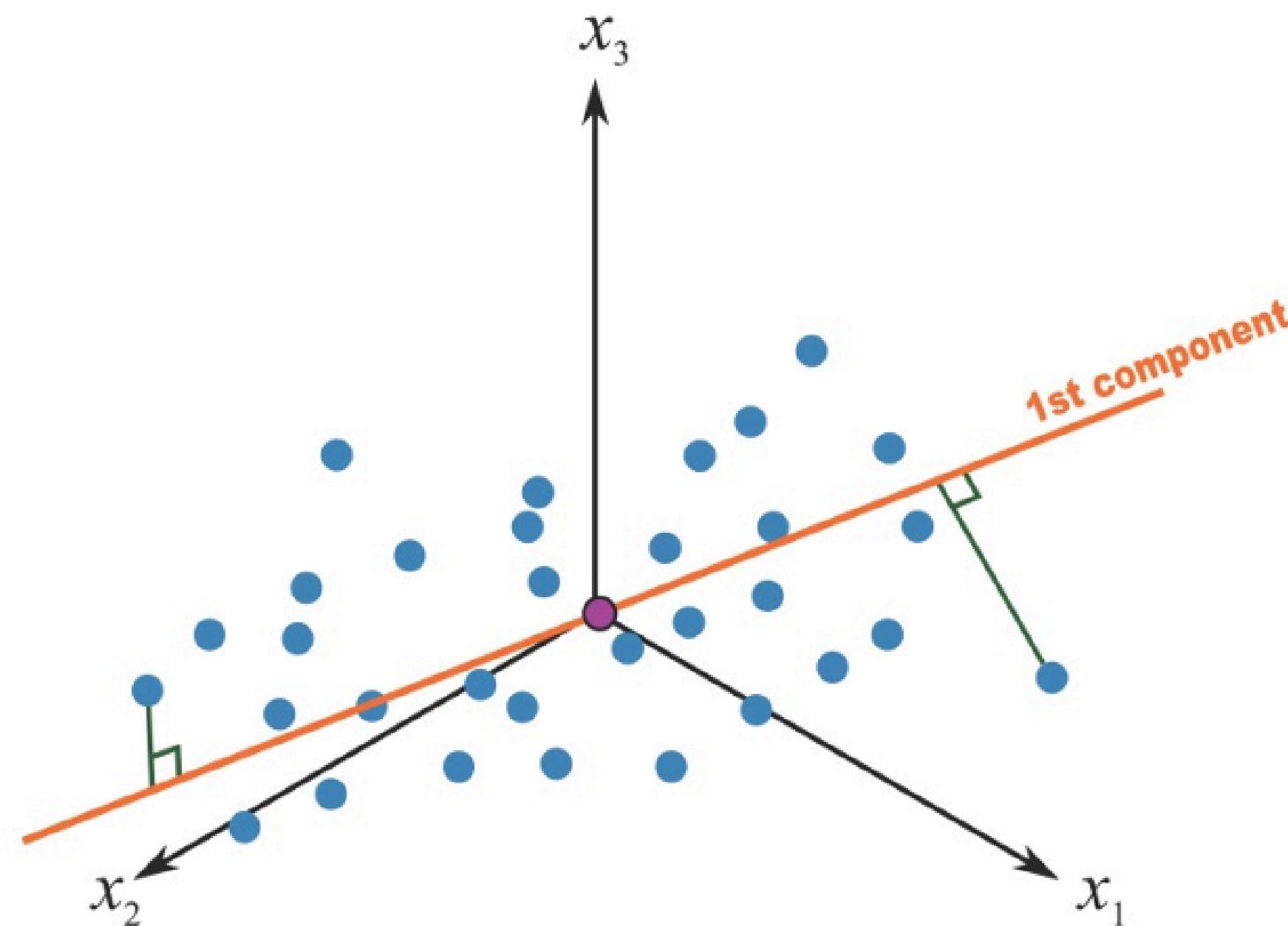


```
3 feature_columns = [  
4     "LH", "FSH", "IGF-1", "IGFBP-3", "GH", "TTE", "PRL", "E2", "Age",  
5     "LH_log", "FSH_log", "IGF-1_log",  
6     "GH_log", "TTE_log", "PRL_log", "E2_log"]
```

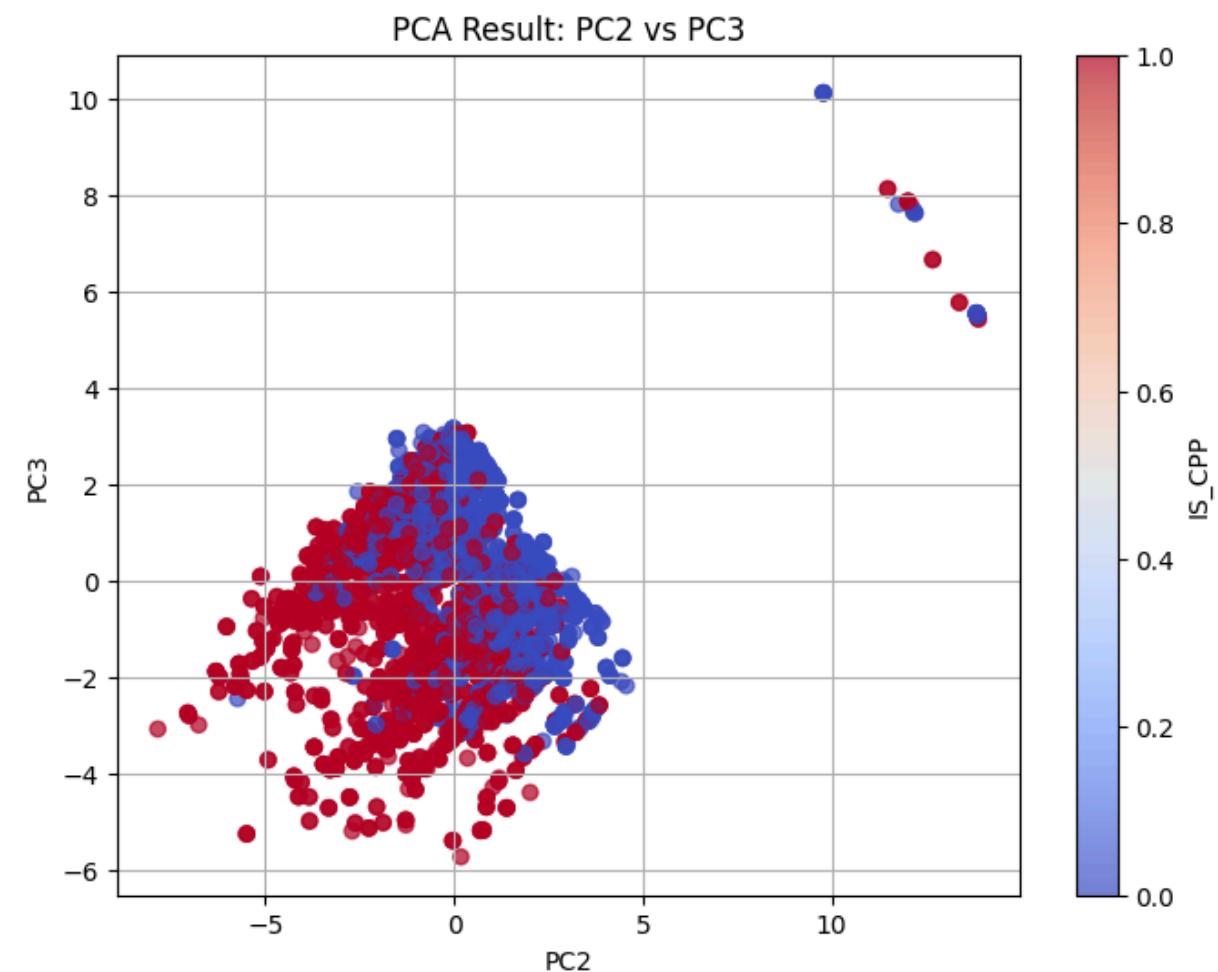
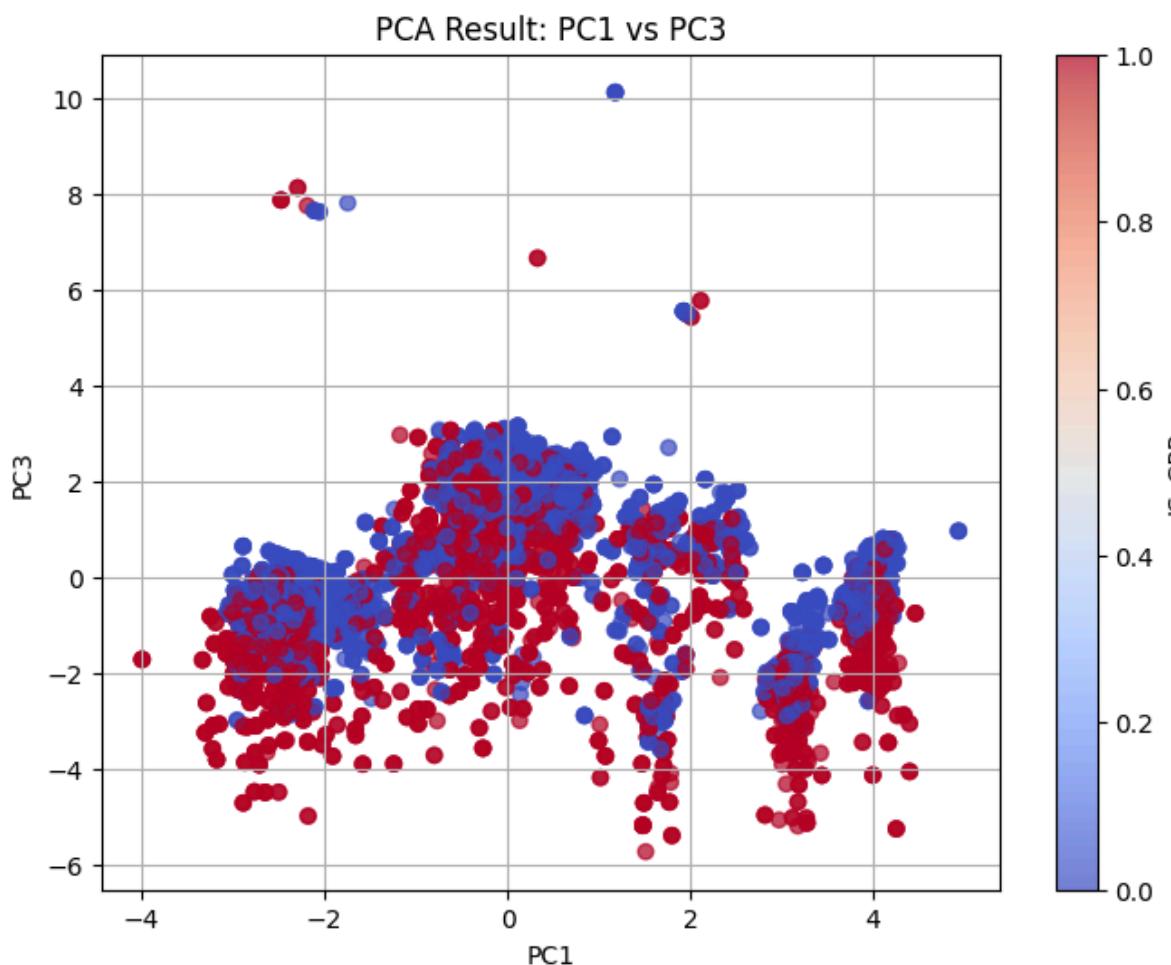
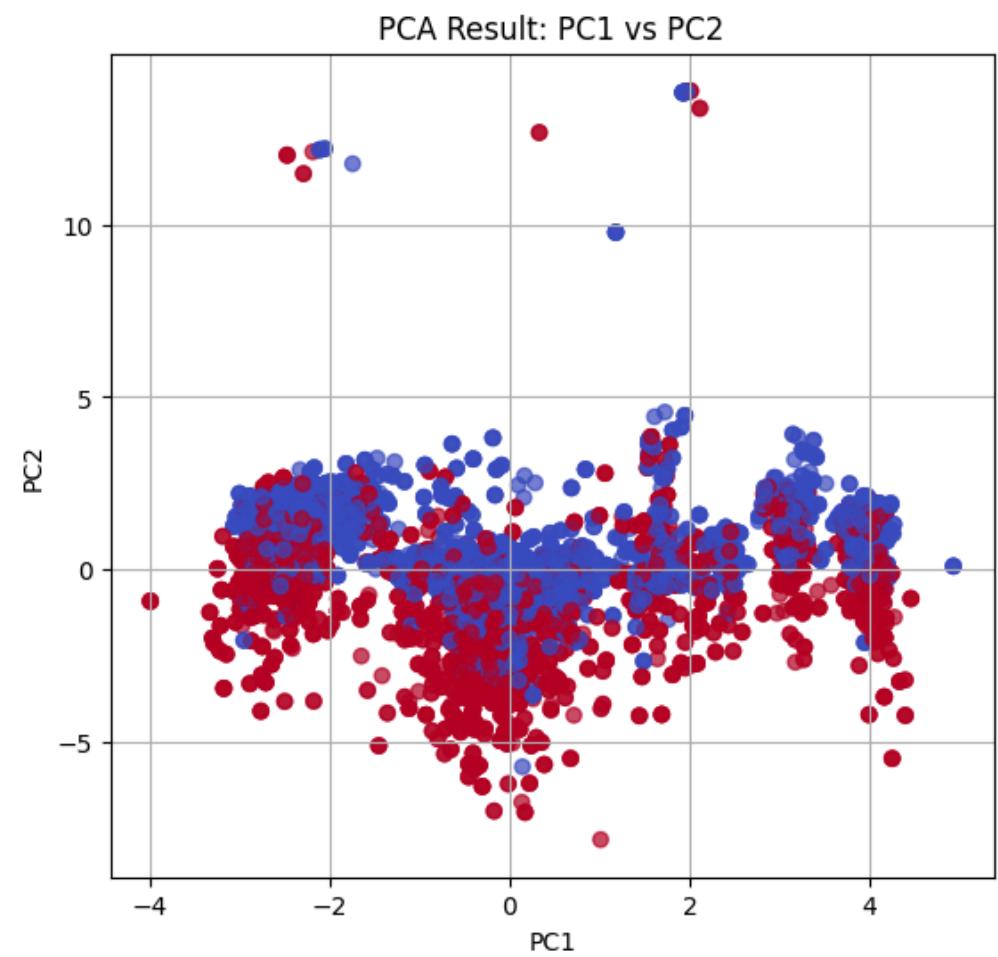
```
10 pca = PCA(n_components=3) #3-dimension  
11 pca_result = pca.fit_transform(X_scaled)  
12  
13 pca_df = pd.DataFrame(data=pca_result, columns=['PC1', 'PC2', 'PC3'])  
14
```

# ▼ PCA

a technique used to reduce the dimensionality of large datasets while retaining most of the variance in the data to reduce complexity and improve computational efficiency.



# ▼ PCA RESULT



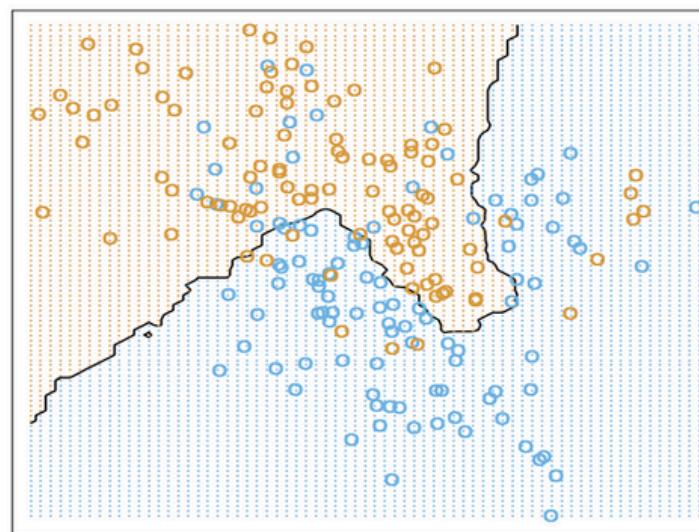
red spot: column of CPP  
blue spot :column of non-CPP

# ▼ DATA SPLITTING

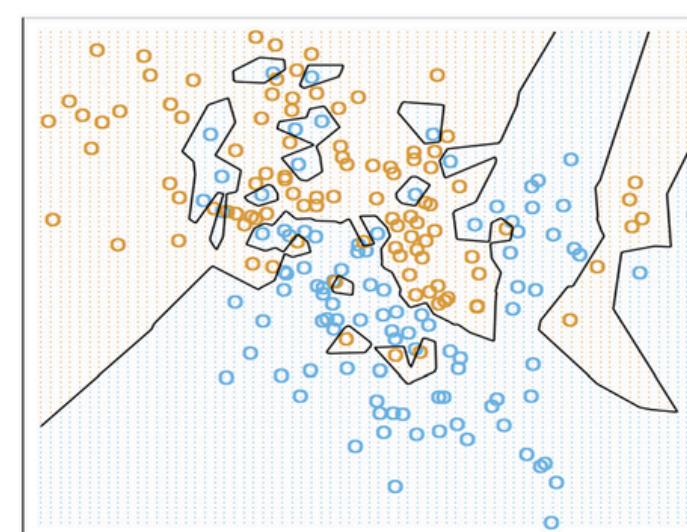
```
9 pca_df
10 X = pca_df.drop('IS_CPP', axis=1)
11 y = pca_df['IS_CPP']
12
13 X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.3, random_state=42)
14
15 accuracies = []
16 recalls = []
17 precisions = []
18
19 # 不同的 K 值
20 for k in range(1,6):
21     knn = KNeighborsClassifier(n_neighbors=k)
22     knn.fit(X_train, y_train)
23
24     y_pred = knn.predict(X_test)
25     accuracy = accuracy_score(y_test, y_pred)
26     recall = recall_score(y_test, y_pred)
27     precision = precision_score(y_test, y_pred)
28
29     accuracies.append(accuracy)
30     recalls.append(recall)
31     precisions.append(precision)
```

Train=0.7  
Test=0.3

# ▽ KNN MODELING



15-nearest neighbors

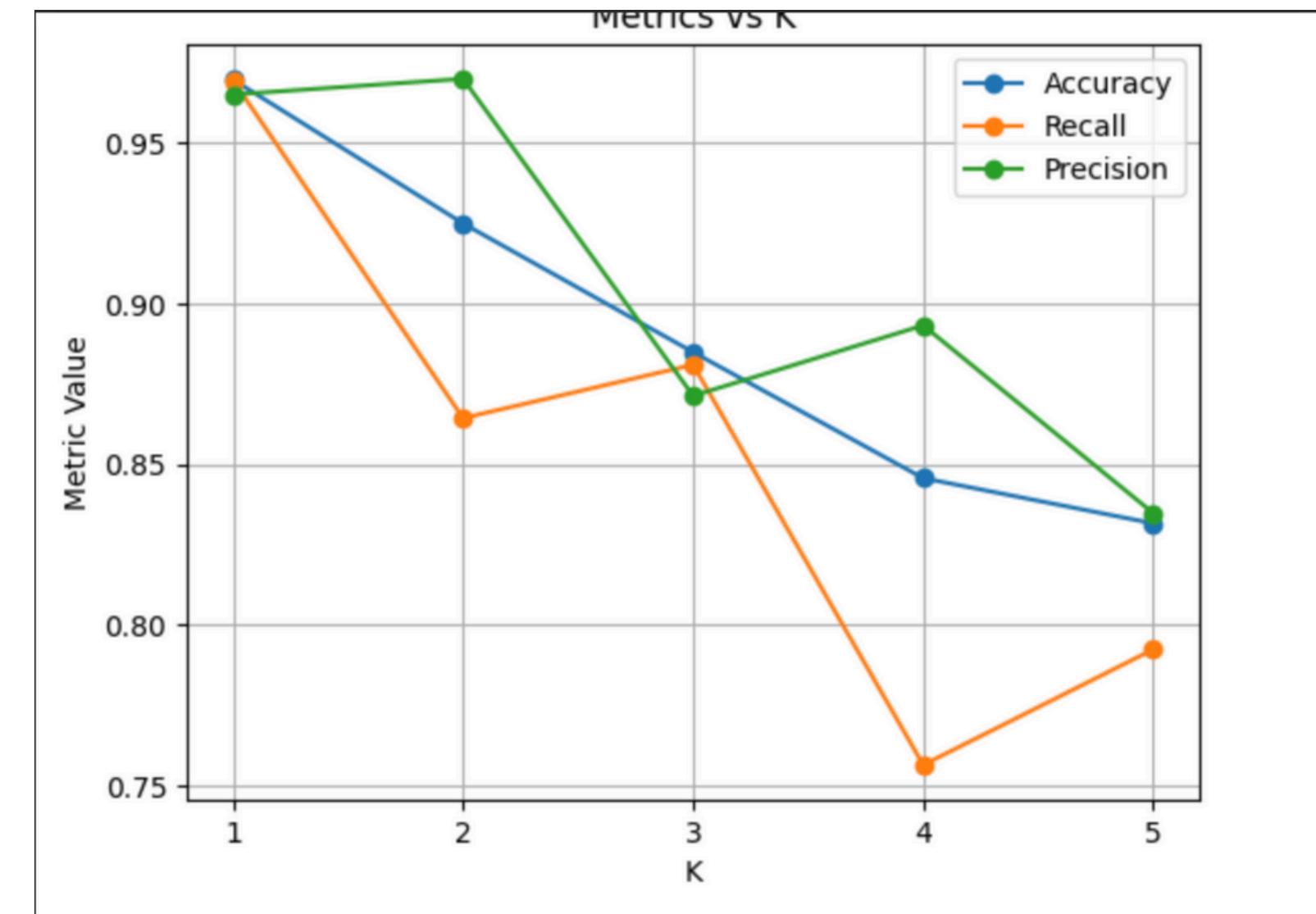


1-nearest neighbors

```
9 pca_df
10 X = pca_df.drop('IS_CPP', axis=1)
11 y = pca_df['IS_CPP']
12
13
14 X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.3, random_state=42)
15
16 accuracies = []
17 recalls = []
18 precisions = []
19
20 # Training KNN
21 for k in range(1,6):
22     knn = KNeighborsClassifier(n_neighbors=k)
23     knn.fit(X_train, y_train)
24
25     y_pred = knn.predict(X_test)
26     accuracy = accuracy_score(y_test, y_pred)
27     recall = recall_score(y_test, y_pred)
28     precision = precision_score(y_test, y_pred)
29
30     accuracies.append(accuracy)
31     recalls.append(recall)
32     precisions.append(precision)
```

# ▽ KNN PREDICTION

```
9 pca_df
10 X = pca_df.drop('IS_CPP', axis=1)
11 y = pca_df['IS_CPP']
12
13
14 X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.3, random_state=42)
15
16 accuracies = []
17 recalls = []
18 precisions = []
19
20 # 不同的 K 值
21 for k in range(1,6):
22     knn = KNeighborsClassifier(n_neighbors=k)
23     knn.fit(X_train, y_train)
24
25     y_pred = knn.predict(X_test)
26     accuracy = accuracy_score(y_test, y_pred)
27     recall = recall_score(y_test, y_pred)
28     precision = precision_score(y_test, y_pred)
29
30     accuracies.append(accuracy)
31     recalls.append(recall)
32     precisions.append(precision)
```



```
1 for k, (accuracy, recall, precision) in enumerate(zip(accuracies, recalls, precisions), 1):
2     print(f"K={k}: Accuracy = {accuracy:.4f}, Recall = {recall:.4f}, Precision = {precision:.4f}")
```

K=1: Accuracy = 0.9697, Recall = 0.9696, Precision = 0.9651  
K=2: Accuracy = 0.9250, Recall = 0.8644, Precision = 0.9700  
K=3: Accuracy = 0.8849, Recall = 0.8810, Precision = 0.8714  
K=4: Accuracy = 0.8457, Recall = 0.7565, Precision = 0.8932  
K=5: Accuracy = 0.8316, Recall = 0.7924, Precision = 0.8348

**Best model: k=1**

acc=0.9697

recall=0.9696

precision=0.9651



**THANK YOU**

**For Your Attention**