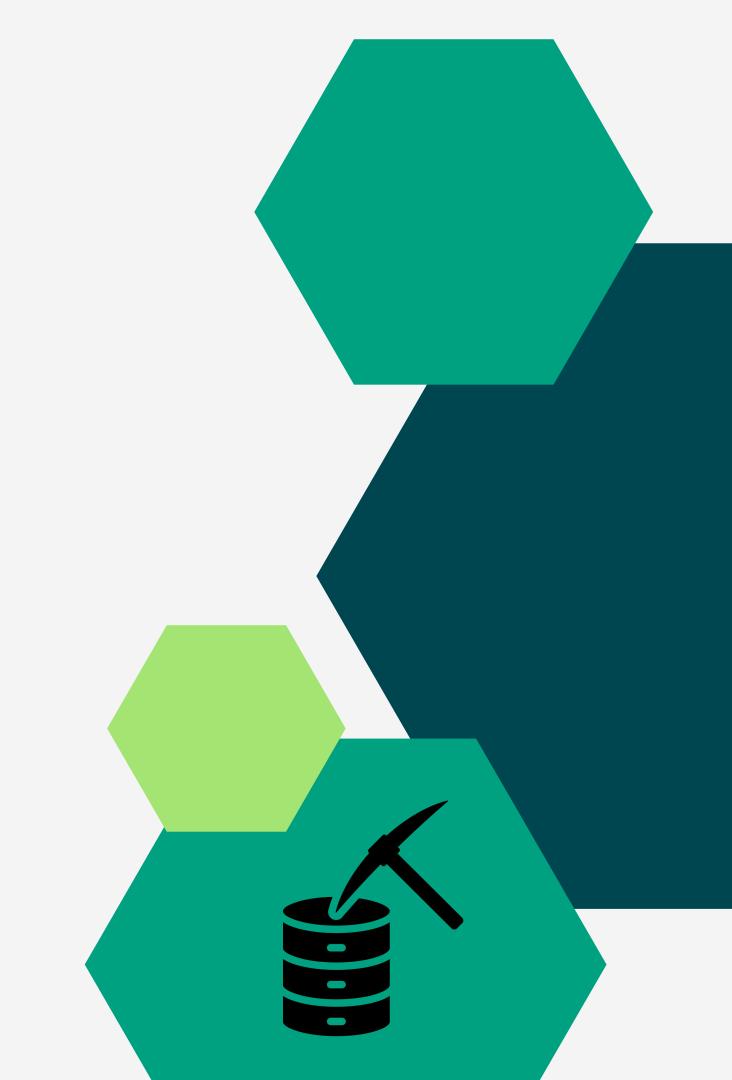
Binary Classification with Convolutional Neural Network Algorithm



### (1) Business Understanding

Proyek ini menggunakan data BPJS Hackathon (Fraud Detection) untuk melakukan prediksi dengan data yang telah kita miliki pendekatan yang digunakan yaitu metode Binary Classification with Convolutional Neural Network Algorithm.



# 2 Data Understanding

### Collecting data

Collecting data merupakan proses pengumpulan, pengukuran serta analisis data yang digunakan dalam penelitian.

#### Describe data

potensi terjadinya *fraud* pada klaim pelayanan Rumah Sakit maka set data yang digunakan dengan menggunakan algoritma *Convolutional Neural Network* (CNN) adalah *fraud\_detection\_train dataset. Dataset* ini terdiri dari 53 *variable* dengan total 200217

observasi.

No.	Nama atribut	Tipe	Deskripsi
	(variabel)	atribut	
1.	visit_id	Nominal	ID kunjungan
2.	kdkc	Nominal	Kode wilayah kantor cabang BPJS
			Kesehatan
3.	dati2	Nominal	Kode kabupaten/kota
4.	typeppk	Nominal	Kode tipe Rumah Sakit
5.	jkpst	Binary	Jenis kelamin peserta JKN-KIS
6.	umur	Nominal	Umur peserta saat mendapatkan
			pelayanan rumah sakit
7.	jnspelsep	Binary	Kode tingkat pelayanan
			1 : rawat inap
			2 : rawat jalan
8.	los	Numerik	Lama pasien dirawat di rumah sakit
9.	cmg	Numerik	Klasifikasi CMG (Case Mix Group)
10.	severitylevel	Numerik	Tingkat urgensi
11.	diagprimer	Nominal	Diagnosa primer
12.	dx2	Numerik	Diagnosa sekunder
13.	proc	Nominal	Kode kelompok procedure
14.	label	Binary	Flag Fraud
			1 : fraud
			0: tidak fraud.

#### Validation Data

Pada subbab ini berisi tahapan evaluasi, kelengkapan data dan kualitas data yang digunakan dalam mengerjakan proyek. Terjadinya missing value maupun noise pada data diakibatkan karena terjadinya kesalahan maupun error pada saat melakukan penginputan data.



Pada tahapan data preparation berikut akan dijabarkan proses menyiapkan data, pemilahan variabel yang akan dianalisis, serta pembersihan data.

#### Data Selection

Data selection atau feature selection digunakan untuk memilih beberapa feature untuk membangun model klasifikasi. Proses seleksi dilakukan dengan melakukan penggabungan terhadap feature yang terkait menjadi satu selanjutnya memilih feature yang akan digunakan sebagai input feature

1 (	df_drop	na['di	lagprime	r'] =	df_dr	opna['dx2	_a00	а_699	] + df_dropna[	'dx2_c	0_d48'] + df_dropna['dx2_d50_d89'] + df_dropna[	'dx2_e
	4											-
Code	diatas d	igunak	an untuk r	menyat	ukan v	alue dari ko	olom	diagp	rimer yang sumbe	r dari ma	sing masing tabel value tersebut	
2 (		na.dro				b99', 'dx =True, ax			8', 'dx2_d50_d8	9', 'd	2_e00_e90', 'dx2_f00_f99', 'dx2_g00_g99', 'dx2_l	h00_h5
	4											-
	kdkc	dati2	typeppk	jkpst	umur	jnspelsep	los	cmg	severitylevel diag	primer	abel	
	0 1107	150	SB	1	64	2	0	F	0	0	1	
	1 1303	200	С	0	45	1	9	E	3	5	1	
	2 1114	172	В	1	34	2	0	Q	0	0	1	
	3 601	90	SC	0	34	2	0	Q	0	0	1	
	4 1006	130	В	0	27	2	0	F	0	0	1	
		***	***	***	***	***				***	No.	
20021	12 2102	353	В	1	48	2	0	z	0	0	0	
20021	13 1308	212	SD	0	1	2	0	Q	0	0	0	
	14 201	38	SB	1	3	2	0	Q	0	0	0	
20021		400	В	1	52	1	1	J	1	0	0	
	15 1008	128										

#### Data Cleaning

Data Cleaning merupakan proses persiapan data dengan cara menghapus atau memodifikasi data yang salah, tidak akurat, tidak terformat maupun duplikat. Data yang rusak tentunya akan

berpengaruh pada kinerja pada sistem.

```
1 df_dropna = df.dropna()
2 print(df_dropna.shape)

(200217, 53)

df_dropna.duplicated().sum()
```

```
df dropna.isna().sum()
[19]: visit_id
      typeppk
      jkpst
      jnspelsep
      severitylevel
      diagprimer
      dx2_a00_b99
      dx2 c00 d48
      dx2 d50 d89
      dx2 e00 e90
      dx2 f00 f99
      dx2_g00_g99
      dx2_h00_h59
      dx2_h60_h95
      dx2 i00 i99
      dx2_j00_j99
      dx2_koo_k93
      dx2_100_199
      dx2_m00 m99
      dx2 n00 n99
      dx2_o00_o99
      dx2 p00 p96
      dx2 q00 q99
      dx2_r00_r99
      dx2_s00_t98
      dx2_u00_u99
      dx2 v01 y98
      dx2_z00_z99
      proc00_13
```

#### Data Construct

200217 rows × 55 columns

Mengkonstruksi data merupakan bagian dari Data transformasi yang terdiri dari representasi fitur, menentukan korelasi dan mengintegrasikan data. representasi fitur digunakan untuk mengurangi kompleksitas, meningkatkan akurasi dan memilih fitur optimal.

```
1 df dropna['jkpst'].replace(to replace=['L','P'], value = [0,1], inplace = True)
Kode diatas digunakan untuk mengubah value kategorikal dari jenis kelamin berupa L dan P menjadi Binary
                 1 import numpy as np
                  2 # Numeric data type
                  3 data_num = df_dropna.select_dtypes(include=[np.number])
                  5 # Category data type
                  6 data_cat = df_dropna.select_dtypes(exclude=[np.number])
                  8 # Get dummies (data transformation)
                  9 transform_cat = pd.get_dummies(data_cat, prefix_sep='_', drop_first=True)
                  1 from numpy.core.defchararray import add
                  2 data_cat = transform_cat.assign(new=add('', np.arange(1, len(data_cat) + 1).astype(str)))
                  3 data_num = data_num.assign(new=add('', np.arange(1, len(data_num) + 1).astype(str)))
                 4 df_dropna = pd.concat([data_cat, data_num], axis=1)
                 5 df_dropna.drop(['new'], axis=1, inplace=True)
                Kode diatas digunakan untuk transformasi data pada semua data kategorikal sehingga semua fitur memiliki nilai numerik.
                 1 df_dropna
                 200212
                 200213
                 200214
                 200215
```

#### Labelling Data

Pada kasus Fraud Detection (Binary Classification) data dibagi menjadi data training dan data validation yang berbeda

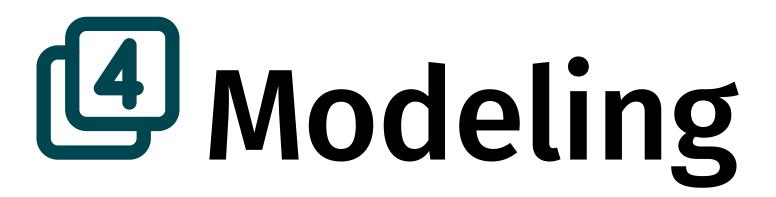
```
1 X = df_dropna.drop('label', axis = 1)
2 y = df_dropna['label']
```

Kode diatas merupakan Feature Selection utk menentukan Input maupun Target Features

#### Data Integration

Pada tahap mengintegrasi data dilakukan concatenation. Concatenation dapat dianggap sebagai sebuah pendekatan untuk menambahkan baris atau kolom ke data. Pendekatan ini dimungkinkan jika data terbagi menjadi beberapa bagian atau jika dilakukan perhitungan yang ingin ditambahkan ke set data yang sudah tersedia.

				version. Us d, ignore_i	e pandas. <mark>co</mark> ndex = True									
typ	peppk_B	typeppk_C	typeppk_D	typeppk_GD	typeppk_HD	typeppk_I1	typeppk_l2	typeppk_I3	typeppk_I4	typeppk_KB		cmg_Z	kdkc	dati2
0	0	0	0	0	0	0	0	0	0	0	- 200	0	1107	150
1	0	1	0	0	0	0	0	0	0	0	-	0	1303	200
2	1	0	0	0	0	0	0	0	0	0	***	0	1114	172
3	0	0	0	0	0	0	0	0	0	0		0	601	90
4	1	0	0	0	0	0	0	0	0	0		0	1006	130
***			440		***		***		444		100			
00505	0	0	0	0	0	0	0	0	0	0		0	1112	169
00506	1	0	0	0	0	0	0	0	0	0	100	0	105	3
00507	0	0	0	0	0	0	0	0	0	0		0	2102	354
00508	0	0	0	0	0	0	0	0	0	0		0	1112	168
00509	0	0	0	0	0	0	0	0	1	0		0	701	97



Pada bab ini dijelaskan mengenai pemilihan teknik modeling, menghasilkan test design, membangun model atau membuat pemodelan, dan menilai model yang telah dibangun. Model yang digunakan adalah Binary Classification with Convolutional Neural Network Algorithm

#### Modelling

#### **Building Test Scenario**

Teknik pemodelan yang dilakukan pada penelitian melibatkan penerapan cnn dalam melakukan prediksi jumlah kasus dan unit cost pada sebuah daerah akibat penambahan Rumah Sakit dari 200217 observasi dan 53 variable. Adapun feature yang digunakan pada dataframe terdiri atas kdkc, dati2, typeppk, jkpst, umur, jnspelsep, los, cmg, severitylevel diagprimer untuk input feature serta label yang menjadi target feature.

```
1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.1, random_state=0)
 2 X train.shape, X test.shape
 1 # Buat objek scaler
 2 scaler = StandardScaler()
 3 # Sesuaikan scaler dengan data
 4 X_train = scaler.fit_transform(X_train)
 5 # Mengubah data train dan test
 6 X test = scaler.transform(X test)
 7 y_train = y_train.to_numpy()
 8 y test = y test.to numpy()
Scaling standardisasi berfokus pada mengubah data mentah menjadi informasi yang dapat digunakan sebelum dianalisis.
 1 X train.shape
(180459, 54)
 1 X_train = X_train.reshape(X_train.shape[0], X_train.shape[1], 1)
 2 X_test = X_test.reshape(X_test.shape[0], X_test.shape[1], 1)
```

#### Modelling

#### **Build Model**

Mendefenisikan model Convolutional Neural Network

```
model = Sequential()
model.add(Conv1D(32, 2, activation='relu', input_shape = (54, 1)))
model.add(BatchNormalization())
model.add(Dropout(0.1))

model.add(Conv1D(64, 2, activation='relu'))
model.add(BatchNormalization())
model.add(Dropout(0.2))

model.add(Flatten())
model.add(Dense(64, activation='relu'))
model.add(Dropout(0.4))

model.add(Dense(1, activation='sigmoid'))
model.summary()
```

Berikut kode untuk melakukan compile model dan fit model cnn

```
1 model.compile(optimizer='adam', loss = 'binary_crossentropy', metrics=['accuracy'])
1 history = model.fit(X_train, y_train, epochs=10, validation_data=(X_test, y_test), verbose=1)
```

Hasil akhir compile model dan fit model



Pada bab ini dilakukan tahap Evaluation (Evaluasi) dengan tujuan untuk memprediksi seberapa baik model akhir akan bekerja nantinya sehingga diketahui apakah model tersebut layak digunakan atau tidak dan untuk membantu menemukan model yang paling mewakili pelatihan data

#### Evaluation

#### Berikut ditampilkan hasil evaluasi terhadap model yang dikembangkan



```
print(classification_report(y_test, y_pred_cnn))
print('precision_score:',precision_score(y_test,y_pred_cnn))
print('accuracy_score:',accuracy_score(y_test,y_pred_cnn))
print('recall_score:',recall_score(y_test,y_pred_cnn))
```

	precision	recall	f1-score	support
0	0.64	0.58	0.61	10059
1	0.61	0.67	0.64	9992
accuracy			0.62	20051
macro avg	0.62	0.62	0.62	20051
weighted avg	0.62	0.62	0.62	20051

precision\_score: 0.6109589041095891 accuracy\_score: 0.6228617026582215 recall\_score: 0.6695356285028022

## © Deployment

Pada babi ini akan dijelaskan mengenai perencanaan fase penyebaran atau penggunaan model yang sudah dihasilkan, perencanaan pemantauan dan pemeliharaan

### Deployment

- Membuat rencana deployment model
- Melakukan deployment model
- Melakukan rencana pemeliharaan
- Melakukan Pemeliharaan



nant you