Proyek UTS PMDPM Gasal 2023/2024

Harap jangan menghapus cell dan komentar yang diberikan!

- Setiap cell markdown dan code akan berisi instruksi pengerjaan Proyek UTS PMDPM Gasal 2023/2024
- Dalam notebook ini Anda akan diminta untuk membuat sebuah proyek Pembelajaran Mesin menggunakan dataset yang sudah disediakan.
- Proyek akan terdiri dari proses inisialisasi, data loading, data cleansing dan encoding, modelling, dan evaluasi model.
- Pada bagian akhir silahkan berikan laporan singkat berupa jawaban pertanyaan dari proyek yang sudah dikerjakan.

Inisialisasi

Bagian berikut berisi import library yang dibutuhkan

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

from google.colab import drive
drive.mount('/content/drive')

Drive already mounted at /content/drive; to attempt to forcibly
remount, call drive.mount("/content/drive", force_remount=True).
```

Data Loading

 Bagian berikut berisi proses data loading (boleh dengan file upload atau dengan mount drive jika menggunakan Google Colab)

```
uts1 = pd.read csv("/content/drive/MyDrive/Mesin/Dataset UTS.csv")
uts1.head(10)
                   numberofrooms hasyard haspool
                                                     floors
                                                              citycode \
   squaremeters
                                3
0
           75523
                                        no
                                                yes
                                                          63
                                                                  9373
1
           55712
                               58
                                                          19
                                                                 34457
                                        no
                                                yes
2
           86929
                                                                 98155
                              100
                                                          11
                                       yes
                                                 no
3
           51522
                                3
                                                                  9047
                                                          61
                                        no
                                                 no
4
           96470
                               74
                                                          21
                                                                 92029
                                       yes
                                                 no
5
           79770
                                3
                                                ves
                                                          69
                                                                 54812
                                        no
6
           75985
                               60
                                                          67
                                                                  6517
                                       yes
                                                 no
7
                                                                 61711
           64169
                               88
                                                          6
                                        no
                                                yes
8
           92383
                               12
                                                          78
                                                                 71982
                                        no
                                                 no
9
           95121
                               46
                                                           3
                                                                  9382
                                        no
                                                yes
```

		rtrange \	numprevowners	made	isnewbuilt	hasstorr	mprotector
0	2	3	8	2005	old		yes
431 1	3	6	8	2021	old		no
293	7	3	4	2002	201		no
2 632	6	3	4	2003	new		no
3		8	3	2012	new		yes
632 4		4	2	2011	new		yes
541	4						-
5 887	1	10	5	2018	old		yes
6		6	9	2009	new		yes
487 7	8	3	9	2011	new		yes
305	4						,
8 750	7	3	7	2000	old		no
9		7	9	1994	old		no
615							
	attic		hasstorageroom	hasgı	uestroom		category
0 1	9005 8852	956 135	no yes			559081.5 574642.1	Luxury Middle
2	4748	654	no			696869.3	Luxury
3	5792	807	yes			154055.2	Middle
4	1172	716	yes			652258.1	Luxury
5	7117 281	240 384	no			986665.8 607322.9	Luxury
6 7	129	726	yes no			420823.1	Luxury Middle
8	9056	892				244344.0	
			-				•
9 sns df_	1221 .set_c uts =	328 ontext('	yes no 'poster") _csv('Dataset UT	S.csv	10 9	244344.0 515440.4	Luxury Luxury
u	u co i ile	au (IV)					

Data Cleansing & Encoding

- Bagian berikut berisi proses pembersihan data.
- Periksa apakah terdapat missing value dan data duplikat,
- Ubah data kategorik string menjadi numerik.
- Jika jumlah kelas pada data latih tidak seimbang, kalian dapat menggunakan metode oversampling.
- Untuk klasifikasi, pastikan Kategori menjadi target dan kolom Harga dihapus.

```
uts2 = uts1.drop('price', axis=1)
uts2.head(10)
                   numberofrooms hasyard haspool
                                                      floors
                                                                citycode \
   squaremeters
0
           75523
                                 3
                                                           63
                                                                    9373
                                         no
                                                 yes
                                58
1
           55712
                                                           19
                                                                   34457
                                         no
                                                 yes
2
           86929
                               100
                                                           11
                                                                   98155
                                        yes
                                                  no
3
           51522
                                 3
                                                           61
                                                                    9047
                                                  no
                                         no
4
           96470
                                74
                                                           21
                                                                   92029
                                        yes
                                                  no
5
                                 3
           79770
                                                           69
                                                                   54812
                                         no
                                                 yes
6
           75985
                                                           67
                                                                    6517
                                60
                                        yes
                                                  no
7
           64169
                                88
                                                            6
                                                                   61711
                                         no
                                                 yes
8
                                12
           92383
                                                           78
                                                                   71982
                                         no
                                                  no
9
           95121
                                46
                                                            3
                                                                    9382
                                         no
                                                 yes
                                     made isnewbuilt hasstormprotector
   citypartrange
                    numprevowners
basement \
                 3
                                  8
                                     2005
                                                   old
0
                                                                        yes
4313
1
                 6
                                  8
                                     2021
                                                   old
                                                                         no
2937
                 3
                                  4
                                     2003
                                                   new
                                                                         no
6326
                 8
                                  3
                                     2012
3
                                                   new
                                                                        yes
632
                 4
                                  2
                                     2011
                                                   new
                                                                        yes
5414
                10
5
                                     2018
                                                   old
                                                                        yes
8871
                 6
                                  9
                                     2009
                                                   new
6
                                                                        yes
4878
7
                 3
                                     2011
                                                   new
                                                                        yes
3054
                 3
                                     2000
                                                   old
8
                                                                         no
7507
                 7
                                  9
                                     1994
                                                   old
                                                                         no
615
   attic
           garage hasstorageroom
                                     hasguestroom category
0
    9005
              956
                                                      Luxury
                                 no
1
    8852
              135
                                                  9
                                                      Middle
                                yes
2
    4748
              654
                                 no
                                                 10
                                                      Luxury
3
    5792
              807
                                                  5
                                                      Middle
                                yes
4
                                                  9
    1172
              716
                                                      Luxury
                                yes
5
    7117
              240
                                                  7
                                 no
                                                      Luxury
6
     281
              384
                                                  5
                                                      Luxury
                                yes
7
     129
                                                  9
              726
                                                      Middle
                                 no
8
    9056
              892
                                                  1
                                                      Luxury
                                yes
9
               328
    1221
                                                 10
                                 no
                                                      Luxury
```

```
print("data null \n",uts2.isnull().sum())
print("data kosong \n",uts2.empty)
print("data nan \n",uts2.isna().sum())
data null
 squaremeters
                       0
numberofrooms
                      0
                      0
hasyard
                      0
haspool
                      0
floors
                      0
citycode
citypartrange
                      0
                      0
numprevowners
                      0
made
isnewbuilt
                      0
                      0
hasstormprotector
basement
                      0
                      0
attic
                      0
garage
                      0
hasstorageroom
                      0
hasquestroom
category
                      0
dtype: int64
data kosong
False
data nan
                       0
 squaremeters
numberofrooms
                      0
hasyard
                      0
haspool
                      0
floors
                      0
citycode
                      0
citypartrange
                      0
                      0
numprevowners
                      0
made
isnewbuilt
                      0
hasstormprotector
                      0
basement
                      0
                      0
attic
                      0
garage
hasstorageroom
                      0
                      0
hasguestroom
category
dtype: int64
print("Sebelum pengecekan data duplikat",uts2.shape)
cols_to_check = ["squaremeters", "numberofrooms", "hasyard", "haspool",
                  "floors", "citycode", "citypartrange",
"numprevowners", "made", "isnewbuilt", "hasstormprotector",
```

Train-test split

- Untuk nilai parameter random_state, sesuaikan dengan dua digit terakhir nomor pegawai terbesar (red: dua digit terakhir NPM terbesar).
- Silahkan memodifikasi persentase train-test split terbaik antara 80:20, 75:25, atau 70:30.

```
from sklearn.model_selection import train_test_split
X = uts2.drop(columns=['category'])
y = uts2.category
X train bf, X test, y train bf, y test = train test split(X, y,
test size=0.20, random state=50)
print(X train bf.shape)
print(X test.shape)
(8000, 16)
(2000, 16)
from sklearn.preprocessing import OneHotEncoder
from sklearn.compose import make column transformer
cat_cols=['hasyard', 'haspool', 'isnewbuilt',
          'hasstormprotector', 'hasstorageroom']
transformer = make column transformer(
    (OneHotEncoder(), cat cols),
    remainder='passthrough'
)
X train enc = transformer.fit transform(X train bf)
X test enc= transformer.transform(X test)
df train enc = pd.DataFrame(X train enc,
```

```
columns=transformer.get feature names out())
df test enc = pd.DataFrame(X test enc,
columns=transformer.get feature names out())
X_enc=df_train_enc
df train enc.head(10)
df test enc.head(10)
   onehotencoder hasyard no
                                onehotencoder__hasyard_yes
0
                           0.0
                                                          1.0
                           0.0
1
                                                          1.0
2
                           1.0
                                                          0.0
3
                           0.0
                                                          1.0
4
                           0.0
                                                          1.0
5
                           0.0
                                                          1.0
6
                           0.0
                                                          1.0
7
                           0.0
                                                          1.0
8
                           0.0
                                                          1.0
9
                           1.0
                                                          0.0
   onehotencoder__haspool_no
                                 onehotencoder__haspool_yes
0
                           \overline{1}.0
                                                          0.0
1
                           0.0
                                                          1.0
2
                           0.0
                                                          1.0
3
                           1.0
                                                          0.0
4
                           0.0
                                                          1.0
5
                           1.0
                                                          0.0
6
                           0.0
                                                          1.0
7
                           1.0
                                                          0.0
8
                           0.0
                                                          1.0
9
                           0.0
                                                          1.0
   onehotencoder isnewbuilt new
                                     onehotencoder isnewbuilt old \
0
                                0.0
                                                                  1.0
1
                                0.0
                                                                  1.0
2
                                1.0
                                                                  0.0
3
                                0.0
                                                                  1.0
4
                                0.0
                                                                  1.0
5
                                1.0
                                                                  0.0
6
                                0.0
                                                                  1.0
7
                                0.0
                                                                  1.0
8
                                0.0
                                                                  1.0
9
                                1.0
                                                                  0.0
   onehotencoder hasstormprotector no
onehotencoder__hasstormprotector_yes \
0
                                      1.0
0.0
                                      0.0
```

1.6	0 0.0		
1.6	0		
3 0.0	1.0		
4	0.0		
1.6	0 0.0	j	
1.6	0		
6 0.0	1.0 0		
7 0.0	1.0		
8	0.0		
1.6	0	j	
0.0			
	onehotencoderhasstorageroom_no o	nehotencoderhass	torageroom_yes
0	1.0		0.0
1	1.0		0.0
2	1.0		0.0
3	0.0		1.0
4	1.0		0.0
 5	1.0		0.0
6	1.0		0.0
7	0.0		1.0
8	1.0		0.0
9	0.0		1.0
0 1 2 3 4 5	remaindernumberofrooms remainder 69.0 35.0 97.0 84.0 61.0 15.0 84.0	floors remainde 4.0 20.0 44.0 96.0 63.0 11.0 95.0	rcitycode \ 28290.0 78664.0 48559.0 70430.0 17514.0 40862.0 62672.0

7 8 9		4. 26. 50.	9	14.0 82.0 96.0	9074	72870.0 90745.0 77787.0	
\	remainder_	_citypartrang	e remainde	rnumprevowner	s remainder	made	
0		1.	9	7.	9	1996.0	
1		1.	9	10.	9	1995.0	
2		8.	9	6.	9	2012.0	
3		8.	9	8.	9	1996.0	
4		8.	9	7.	9	2016.0	
5		2.	9	3.	9	2014.0	
6		7.	9	9.	9	2018.0	
7		7.	9	7.	9	2004.0	
8		7.	9	5.	9	1991.0	
9		9.	9	4.	9	2021.0	
0 1 2 3 4 5 6 7 8	remainder_	_basement re 8825.0 7528.0 9379.0 3856.0 8808.0 1832.0 6502.0 8764.0 4052.0 187.0	mainderat- 6894 1094 8734 344 6952 6730 5899 3150 3302	4.0 4.0 4.0 6.0 2.0 6.0 9.0 9.0	_garage \ 329.0 219.0 958.0 331.0 999.0 878.0 438.0 427.0 397.0 981.0		
0 1 2 3 4 5 6 7 8 9	remainder_	_hasguestroom 3.0 9.0 1.0 5.0 1.0 3.0 0.0 4.0 0.0 2.0					

```
[10 rows x 21 columns]
from imblearn.over_sampling import SMOTE
from collections import Counter

smt = SMOTE()
print(Counter(y_train_bf))

X_train,y_train = smt.fit_resample(X_train_enc,y_train_bf)
print(Counter(y_train))

Counter({'Basic': 3478, 'Luxury': 2436, 'Middle': 2086})
Counter({'Luxury': 3478, 'Basic': 3478, 'Middle': 3478})
```

Modelling

Tugas Klasifikasi

- Buatlah dua Pipeline untuk perbandingan model algoritme.
- Tahap pemodelan dimulai dari data scaling, feature selection, hingga algoritme classifier.
- Bandingkan dua metode penskalaan yaitu StandardScaler dan MinMaxScaler menggunakan parameter grid.
- Kalian perlu bereksperimen dengan membandingkan dua dari empat metode feature selection (SelectKBest, SelectPercentile, SelectFromModel, dan RFE) dan jumlah feature yang dipilih menggunakan parameter grid.
- Kalian wajib menyesuaikan parameter dari algoritme classifier yang kalian pilih agar model dapat bekerja dengan baik pada dataset.
- Dua algoritme classifier yang dibandingkan bebas.

```
from sklearn.model selection import GridSearchCV, StratifiedKFold
from sklearn.neighbors import KNeighborsClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import MinMaxScaler
from sklearn.feature selection import SelectFromModel, SelectKBest
from sklearn.metrics import classification report, confusion matrix
pipe KNN = Pipeline(steps=[
    ('scale',MinMaxScaler()),
    ('feat select', SelectKBest()),
    ('clf', KNeighborsClassifier())])
param grid KNN= [
    {'feat select k': np.arange(2,5),
     'clf__n_neighbors': [3, 5, 7, 9],
     'clf_weights': ['uniform', 'distance'],
     'clf p': [1, 2]},
```

```
'feat select' :
[SelectFromModel(estimator=DecisionTreeClassifier
       (random state=50, max depth=3), max features=4)],
      'clf n neighbors': [\overline{3}, 5, 7, 9],
      'clf weights': ['uniform', 'distance'],
      'clf p': [1, 2]
]
GSCV KNN = GridSearchCV(pipe KNN,
param grid KNN,cv=StratifiedKFold(n splits=5))
GSCV KNN.fit(X train enc,y train bf)
mask =
GSCV KNN.best estimator .named steps['feat select'].get support()
print("Best model:{}".format(GSCV KNN.best estimator ))
print("Selected features:{}".format(X enc.columns[mask]))
print("Best CV score: {:.2f}".format(GSCV KNN.best score ))
print("Test set score:
{:.2f}".format(GSCV_KNN.score(X_test_enc,y_test)))
Best model:Pipeline(steps=[('scale', MinMaxScaler()), ('feat select',
SelectKBest(k=4)),
                ('clf', KNeighborsClassifier(n neighbors=9, p=1))])
Selected features:Index(['onehotencoder hasyard yes',
'onehotencoder__haspool_no',
       'onehotencoder haspool yes', 'remainder squaremeters'],
      dtvpe='object')
Best CV score: 0.94
Test set score: 0.94
from sklearn.tree import DecisionTreeClassifier, plot tree
pipe DT = Pipeline(steps=[
    ('feat select', SelectKBest()),
    ('clf',DecisionTreeClassifier(random state=50))])
param grid DT = [
    {'feat select k': np.arange(2,5),
     'clf__max_depth': [2,3,4,5],
     'clf criterion': ['gini', 'entropy']
    },
    {
      'feat select' : [SelectFromModel
(estimator=DecisionTreeClassifier(random state=50), max features=4)],
      'clf max depth': [2,3,4,5],
      'clf__criterion': ['gini','entropy']
```

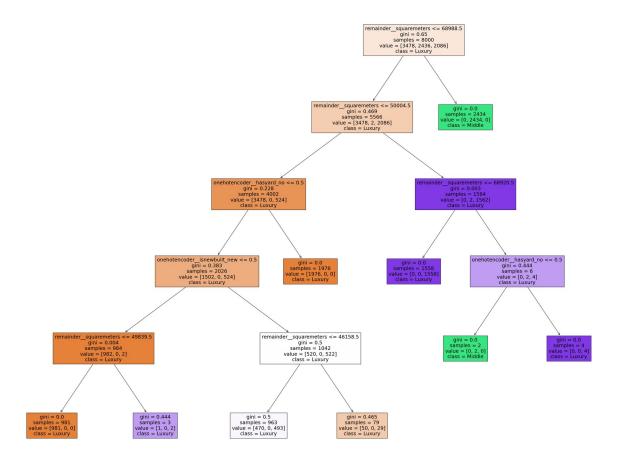
```
}
GSCV DT = GridSearchCV(pipe DT, param grid DT,
cv=StratifiedKFold(n splits=5))
GSCV DT.fit(X train enc, y train bf)
mask =
GSCV DT.best estimator .named steps['feat select'].get support()
print("Best model:{}".format(GSCV DT.best estimator ))
print("Selected features:{}".format(X_enc.columns[mask]))
print("Best CV score: {:.2f}".format(GSCV DT.best score ))
print("Test set score:
{:.2f}".format(GSCV DT.score(X test enc,y test)))
Best model:Pipeline(steps=[('feat select',
SelectFromModel(estimator=DecisionTreeClassifier(random state=50),
                                 max features=4)),
                ('clf', DecisionTreeClassifier(max_depth=5,
random state=50))])
Selected features:Index(['onehotencoder hasyard no',
'onehotencoder isnewbuilt new',
       'remainder squaremeters'],
      dtype='object')
Best CV score: 0.94
Test set score: 0.94
```

Evaluasi Model

- Evaluasi dilakukan dengan membuat masing-masing dua Grid Search Cross Validation dengan metode Stratified KFold Cross Validation.
- Untuk klasifikasi, tampilkan feature yang relevan dipilih oleh model, hasil pengukuran kinerja model klasifikasi dengan confusion matrix, serta metrik accuracy, precision, recall, dan F1-score. Untuk memudahkan pimpinan divisi membaca hasil, buatlah confusion matrix dalam bentuk representasi visual menggunakan fungsi ConfusionMatrixDisplay dari library scikit-learn.

```
Classification report KNN:
                             recall f1-score
               precision
                                                support
       Basic
                   0.92
                             0.93
                                        0.93
                                                   866
                   1.00
                              1.00
                                        1.00
                                                   629
      Luxury
      Middle
                   0.88
                             0.87
                                        0.87
                                                   505
                                        0.94
                                                  2000
    accuracy
   macro avq
                   0.93
                              0.93
                                        0.93
                                                  2000
                   0.93
                              0.94
                                        0.93
                                                  2000
weighted avg
pred = GSCV DT.predict(X test enc)
print("Confusion matrix DT:\n",confusion matrix(y test,pred))
print("Classification report DT:\
n",classification report(y test,pred,zero division=0))
Confusion matrix DT:
 [[749
        0 117]
   0 629
            01
   7
        1 49711
Classification report DT:
               precision
                             recall f1-score
                                                support
                   0.99
                             0.86
       Basic
                                        0.92
                                                   866
      Luxury
                   1.00
                              1.00
                                        1.00
                                                   629
      Middle
                   0.81
                              0.98
                                        0.89
                                                   505
                                        0.94
    accuracy
                                                  2000
   macro avg
                   0.93
                             0.95
                                        0.94
                                                  2000
weighted avg
                   0.95
                             0.94
                                        0.94
                                                  2000
from matplotlib import pyplot as plt
DT model = GSCV DT.best estimator .named steps['clf']
plt.figure(figsize = (40,30))
pict = plot_tree(DT_model, filled=True, feature names =
```

X_enc.columns[mask], class_names=uts2.category)



Data Cleansing & Encoding Regresi

- Load kembali dataset yang akan digunakan.
- Ubah data kategorik string menjadi numerik.
- Untuk regresi, pastikan Harga menjadi target dan kolom Kategori dihapus.

```
df uts1 = df uts[~df uts[cols to check].duplicated(keep='last')]
print("Sesudah pengecekan data duplikat",df utsl.shape)
df uts1.price.round(1)
df uts1.price.plot(kind='box')
from sklearn.model selection import train test split
X=df uts1.drop(columns=['category','price'],axis=1)
v=df uts1.price
X train, X test,y train,y test =
train test split(X,y, test size=0.20, random state=85)
print(X train.shape)
print(X test.shape)
from sklearn.preprocessing import OneHotEncoder
from sklearn.compose import make column transformer
cat_cols=['hasyard', 'haspool','isnewbuilt',
          'hasstormprotector', 'hasstorageroom']
transformer = make column transformer(
    (OneHotEncoder(), cat cols),
    remainder='passthrough'
)
X train enc = transformer.fit transform(X train)
X test enc= transformer.transform(X test)
df train enc = pd.DataFrame(X train enc,
columns=transformer.get feature names out())
df test enc = pd.DataFrame(X test enc,
columns=transformer.get feature names out())
df train enc.head(10)
df test enc.head(10)
```

Tugas Regresi

- Buatlah dua Pipeline untuk perbandingan model algoritme.
- Bandingkan dua metode penskalaan yaitu StandardScaler dan MinMaxScaler menggunakan parameter grid.
- Kalian juga perlu menyesuaikan parameter dari algoritme regressor yang kalian pilih agar model dapat bekerja dengan baik pada dataset.

• Dua algoritme regressor yang dibandingkan bebas.

```
from sklearn.linear model import LinearRegression
from sklearn.model selection import GridSearchCV
from sklearn.preprocessing import StandardScaler
from sklearn.preprocessing import MinMaxScaler
from sklearn.metrics import mean absolute error, mean squared error
from sklearn.pipeline import Pipeline
pipe LR = Pipeline(steps=[
    ('scale', MinMaxScaler()),
    ('reg', LinearRegression())
param grid LR = {
}
GSCV LR = GridSearchCV(pipe LR, param grid LR, cv=5,
scoring='neg mean squared error')
GSCV_LR.fit(X train enc, y train)
print("Best model:{} ".format(GSCV LR.best estimator ))
print('Koefisien/bobot:
{}'.format(GSCV LR.best estimator .named steps['reg'].coef ))
print('Intercept/Bias:
{}'.format(GSCV LR.best estimator .named steps['reg'].intercept ))
LR pred = GSCV LR.predict(X test enc)
mae_LR = mean_absolute_error(y_test, LR_pred)
mse LR = mean squared error(y test, LR pred)
print("LR MAE: ", mae_LR)
print("LR MSE: ", mse_LR)
print("LR Root Mean Square Error: ", np.sqrt(mse LR))
df result = pd.DataFrame(y test)
df result['LR pred'] = LR pred.round(1)
df result.head(10)
from sklearn.linear model import Ridge
pipe Ridge = Pipeline(steps=[
    ('scale', StandardScaler()),
    ('reg', Ridge())
    1)
param grid Ridge = {
    'reg alpha': [0.1, 1, 10, 100, 1000],
    'reg tol': [0.1, 1, 10, 100, 1000],
```

```
'reg solver': ['auto', 'svd', 'cholesky', 'lsqr', 'sparse_cg',
'sag', 'saga'],
    'reg random state': [85]
GSCV Ridge = GridSearchCV(pipe Ridge, param grid Ridge, cv=5,
scoring='neg mean squared error')
GSCV Ridge.fit(X train enc, y train)
print("Best model:{} ".format(GSCV Ridge.best estimator ))
print('Koefisien/bobot:
{}'.format(GSCV Ridge.best estimator .named steps['reg'].coef ))
print('Intercept/Bias:
{}'.format(GSCV Ridge.best estimator .named steps['reg'].intercept ))
Ridge pred = GSCV Ridge.predict(X test enc)
mse_Ridge = mean_squared_error(y_test, Ridge_pred)
mae Ridge = mean absolute error(y test, Ridge pred)
print("Ridge MAE: ", mae_Ridge)
print("Ridge MSE: ", mse_Ridge)
print("Ridge Root Mean Square Error: ", np.sqrt(mse Ridge))
df result['Ridge pred'] = Ridge pred.round(1)
df result.head(10)
df final = df result.head(300)
plt.figure(figsize=(30,5))
data len = range(1,301)
plt.scatter(data len, df final.price, color='blue', label='Actual')
plt.plot(data len, df final['LR pred'], color='green',linestyle='--',
label='LR')
plt.legend()
plt.title('Perbadingan Data Actual dan Prediksi Linear Regression')
plt.show()
plt.figure(figsize=(30,5))
plt.scatter(data_len, df_final.price, color='blue', label='Actual')
plt.plot(data len, df final['Ridge pred'],
color='black',linestyle='-.', label='Ridge')
plt.legend()
plt.title('Perbadingan Data Actual dan Prediksi Ridge Regression')
plt.show()
```

Evaluasi Model

- Evaluasi dilakukan dengan membuat masing-masing dua Grid Search Cross Validation dengan metode Stratified KFold Cross Validation.
- Untuk regresi, tampilkan hasil pengukuran kinerja model regresi dengan Mean Absolute Error, Mean Squared Error, dan Root Mean Squared Error. Untuk memudahkan pimpinan divisi membaca hasil, buatlah tabel yang menampilkan perbandingan harga asli properti dan harga hasil prediksi dua model regresi beserta grafik visualnya.

Kesimpulan

Silahkan jawab pertanyaan berikut pada cell markdown yang sudah disediakan:

Klasifikasi

- 1. Apa saja feature-feature yang relevan untuk membedakan antara rumah basic, middle, dan luxury?
- Untuk model KNN yang digunakan adalah 'onehotencoder_hasyard_yes',
 'onehotencoder_haspool_no', 'onehotencoder_haspool_yes',
 'remainder_squaremeters'
- Untuk model Decision Tree yang digunakan adalah 'onehotencoder_hasyard_no',
 'onehotencoder_isnewbuilt_new', 'remainder_squaremeters'
- Model classifier apakah yang memiliki performa paling baik terhadap dataset?
 Apakah parameter setting yang paling optimal dari model tersebut? Kedua model memiliki performa yang sama
- 2. Dari hasil Confusion Matrix dan Classification Report, hal apa yang bisa kalian simpulkan dari performa model terbaik kalian? Kedua model memiliki hasil yang relatif sama

Regresi

- 1. Model regressor apakah yang memiliki performa paling baik terhadap dataset? Apakah parameter setting yang paling optimal dari model tersebut?
- Model Regressor terbaik adalah model Rigde Regression dengan memiliki nilai MAE dan MSE yang lebih kecil dari Linear Regression
- Untuk Ridge Regression, parameter setting yang paling optimal adalah alpha = 0.1, solve ='svd' dan tol = 0.1
- 1. Dari hasil metrik evaluasi, hasil prediksi, dan grafik regresi, hal apa yang bisa kalian simpulkan dari performa model terbaik kalian? -Model Algoritma Machine Learning yang dikembangkan memiliki prediksi yang akurat dengan pemyimpangan nilai yan kecil

tulis jawaban kalian di cell ini

Kelompok: 1 Gideon Bahtera A.P. / 21071125- Nama anggota 2 / NPM Nama anggota 3 / NPM Nama anggota 4 / NPM

Jawaban:

Klasifikasi

- Apa saja feature-feature yang relevan untuk membedakan antara rumah basic, middle, dan luxury?
- Untuk model KNN yang digunakan adalah 'onehotencoder_hasyard_yes',
 'onehotencoder_haspool_no', 'onehotencoder_haspool_yes',
 'remainder_squaremeters'
- Untuk model Decision Tree yang digunakan adalah 'onehotencoder_hasyard_no',
 'onehotencoder_isnewbuilt_new', 'remainder_squaremeters'
- 1. Model classifier apakah yang memiliki performa paling baik terhadap dataset? Apakah parameter setting yang paling optimal dari model tersebut? Kedua model memiliki performa yang sama
- 2. Dari hasil Confusion Matrix dan Classification Report, hal apa yang bisa kalian simpulkan dari performa model terbaik kalian? Kedua model memiliki hasil yang relatif sama

Regresi

- 1. Model regressor apakah yang memiliki performa paling baik terhadap dataset? Apakah parameter setting yang paling optimal dari model tersebut?
- Model Regressor terbaik adalah model Rigde Regression dengan memiliki nilai MAE dan MSE yang lebih kecil dari Linear Regression
- Untuk Ridge Regression, parameter setting yang paling optimal adalah alpha = 0.1, solve ='svd' dan tol = 0.1
- 1. Dari hasil metrik evaluasi, hasil prediksi, dan grafik regresi, hal apa yang bisa kalian simpulkan dari performa model terbaik kalian? -Model Algoritma Machine Learning yang dikembangkan memiliki prediksi yang akurat dengan pemyimpangan nilai yan kecil