

Model Stacking

In questo notebook sono usati i dataset su cui è stata operata la fase preliminare di **preparazione dei dati**, la **divisione** in base alla contea e la **selezione delle feature**.

È costruito un **metamodello** che fa utilizzo di un dataset le cui feature sono le predizioni dei tre modelli costruiti, rispettivamente un **DecisionTreeRegressor**, un **AdaBoostRegressor** e un **RandomForestRegressor**.

Il metamodello è allenato e allenato con un **LinearRegressor** sulle stesse istanze di train usate per generare le predizioni del nuovo dataset.

```
In [2]: # libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import warnings

from sklearn.metrics import mean_squared_error
from sklearn.ensemble import RandomForestRegressor
from sklearn.ensemble import AdaBoostRegressor
from sklearn.tree import DecisionTreeRegressor
from sklearn.linear_model import LinearRegression
# Lettura dei dati ##

warnings.filterwarnings('ignore')
```

Lettura dei dati

```
In [3]: # local file paths
dir_name = 'selezioni'
region_names = np.array(['A', 'B', 'C'])

fp_Xtrain = []
fp_Xval = []
fp_Xtest = []
fp_Ytrain = []
fp_Yval = []
fp_Ytest = []

for i in range(3):
    fp_Xtrain.append(dir_name + f'/{X_train(region_names[i]).csv}')
    fp_Xval.append(dir_name + f'/{X_val(region_names[i]).csv}')
    fp_Xtest.append(dir_name + f'/{X_test(region_names[i]).csv}')
    fp_Ytrain.append(dir_name + f'/{Y_train(region_names[i]).csv}')
    fp_Yval.append(dir_name + f'/{Y_val(region_names[i]).csv}')
    fp_Ytest.append(dir_name + f'/{Y_test(region_names[i]).csv}')
```

```
In [4]: # Lettura dei dati

X_train = []
X_val = []
X_test = []
Y_train = []
Y_val = []
Y_test = []

for i in range(3):
    X_train.append(pd.read_csv(fp_Xtrain[i], low_memory=False))
    X_val.append(pd.read_csv(fp_Xval[i], low_memory=False))
    X_test.append(pd.read_csv(fp_Xtest[i], low_memory=False))
    Y_train.append(pd.read_csv(fp_Ytrain[i], low_memory=False))
    Y_val.append(pd.read_csv(fp_Yval[i], low_memory=False))
    Y_test.append(pd.read_csv(fp_Ytest[i], low_memory=False))

X_train = np.array(X_train, dtype=object)
X_val = np.array(X_val, dtype=object)
X_test = np.array(X_test, dtype=object)
Y_train = np.array(Y_train, dtype=object)
Y_val = np.array(Y_val, dtype=object)
Y_test = np.array(Y_test, dtype=object)
```

```
In [5]: # dimensionality (y=False)
def i in range(3):
    print(f'X_train(region_names[i]): {X_train[i].shape}')
    print(f'X_val(region_names[i]): {X_val[i].shape}')
    print(f'X_test(region_names[i]): {X_test[i].shape}')
    if y:
        print(f'Y_train(region_names[i]): {Y_train[i].shape}')
        print(f'Y_val(region_names[i]): {Y_val[i].shape}')
        print(f'Y_test(region_names[i]): {Y_test[i].shape}')
    print()
```

```
In [6]: dimensionality(y=True)

X_train: (26819, 55)
X_val: (9085, 55)
X_test: (26819, 1)
Y_train: (26819, 1)
Y_val: (9085, 1)
Y_test: (9085, 1)
```

```
X_train: (8119, 32)
X_val: (26819, 32)
X_test: (26819, 32)
Y_train: (8119, 1)
Y_val: (26819, 1)
Y_test: (26819, 1)
```

```
X_trainC: (64771, 33)
X_valC: (21876, 33)
X_testC: (21876, 33)
Y_trainC: (64771, 1)
Y_valC: (21876, 1)
Y_testC: (21876, 1)
```

Modelli precedentemente individuati

Albero di decisione

```
In [7]: dt_model = np.array([
    DecisionTreeRegressor(max_leaf_nodes = 10),
    DecisionTreeRegressor(max_leaf_nodes = 10),
    DecisionTreeRegressor(max_leaf_nodes = 10)
])
```

Boosting

```
In [8]: boost_model = np.array([
    AdaBoostRegressor(dt_model[0], n_estimators = 20),
    AdaBoostRegressor(dt_model[1], n_estimators = 20),
    AdaBoostRegressor(dt_model[2], n_estimators = 20)
])
```

Foresta

```
In [9]: rf_model = np.array([
    RandomForestRegressor(n_estimators = 850, n_jobs = -1),
    RandomForestRegressor(n_estimators = 700, n_jobs = -1),
    RandomForestRegressor(n_estimators = 700, n_jobs = -1)
])
```

```
In [10]: models = np.array([
    dt_model,
    boost_model,
    rf_model
])
```

Costruzione del Metamodello

Costruzione di un dataset **Y** cui ciascuna colonna contiene le **predizioni dei modelli selezionati per le istanze di Train**. Le colonne hanno il nome del tipo di modello usato, l'operazione è usata per ottenere un dataset **Y** sia per il **Train** che per il **Test**.

```
In [11]: def gen_Y(X, y, models, col_names, index):
    Y = pd.DataFrame()
    for i in range(len(col_names)):
        models[i].fit(X[index], y[index])
        Y[col_names[i]] = models[i][index].predict(X[index])
    return Y
```

```
In [12]: def gen_arrY(X, y, models, col_names):
    Y = []
    for i in range(len(X)):
        Y.append(gen_Y(X, y, models, col_names, i))
    Y = np.array(Y, dtype=object)
    return Y
```

```
In [13]: col_names = ['DecisionTree', 'Boosting', 'RandomForest']

In [14]: Y_train = gen_arrY(X_train, y_train, models, col_names)
Y_test = gen_arrY(X_test, y_test, models, col_names)
```

```
In [15]: Y_train[0].head(20)

Out[15]:
```

	DecisionTree	Boosting	RandomForest
0	0.050485	0.266712	0.055412
1	0.007182	0.133440	0.025987
2	0.007182	0.056837	-0.006864
3	0.007182	0.032949	-0.019451
4	0.017906	0.095386	0.021818
5	0.007182	0.091944	0.023328
6	0.017906	0.201193	0.123257
7	0.007182	0.182836	-0.028503
8	0.017906	0.182836	0.026966
9	0.007182	0.071312	-0.002090
10	0.007182	0.188450	0.006948
11	0.007182	0.032391	0.015790
12	0.007182	0.201981	-0.004285
13	0.007182	0.198270	-0.027130
14	0.007182	0.103031	-0.000350
15	0.017906	0.061509	0.035279
16	0.017906	-0.023329	0.004317
17	0.007182	0.058858	0.003308
18	0.007182	0.023251	-0.000582
19	0.017906	0.112282	0.097573

```
In [16]: Y_test[0].head(20)

Out[16]:
```

	DecisionTree	Boosting	RandomForest
0	0.007520	0.034881	0.015903
1	0.007520	0.059259	0.013394
2	0.007520	-0.89670	0.017264
3	0.007520	-0.025480	-0.052040
4	0.040975	-0.861412	0.000268
5	0.040975	-0.921076	-0.021154
6	0.007520	-1.097865	-0.000320
7	0.007520	0.007690	-0.010995
8	0.007520	0.071312	-0.033108
9	0.007520	0.007690	-0.044424
10	0.010141	0.137659	0.010438
11	0.007520	-0.505645	0.010550
12	0.010141	0.006302	0.023889
13	0.007520	0.592034	0.029549
14	0.007520	0.051201	0.028119
15	0.007520	0.005953	-0.008425
16	0.010141	0.032760	0.070837
17	0.007520	0.673133	0.065244
18	0.010141	0.710624	0.019638
19	0.007520	0.420446	0.014377

```
In [17]: def dimensionalityY(i):
    for i in range(3):
        print(region_names[i])
        print(f'X_train(region_names[i]): {X_train[i].shape}')
        print(f'Y_train(region_names[i]): {Y_train[i].shape}')
        print(f'X_test(region_names[i]): {X_test[i].shape}')
        print(f'Y_test(region_names[i]): {Y_test[i].shape}')
        print()
```

```
In [18]: dimensionalityY()

A
X_train: (26819, 55)
Y_train: (26819, 3)
Y_test: (9085, 3)
X_test: (9085, 55)
Y_test: (9085, 1)

B
X_train: (8119, 32)
Y_train: (8119, 3)
Y_test: (26819, 32)
X_test: (26819, 32)
Y_test: (26819, 1)

C
X_train: (64771, 33)
Y_train: (64771, 3)
Y_test: (21876, 33)
X_test: (21876, 33)
Y_test: (21876, 1)
```

Regressione lineare

Modello di regressione lineare allentato sul nuovo dataset costruito con le predizioni per le **Y di train** e le stesse **y di train** su cui sono stati allenati i modelli. Le predizioni sono calcolate sulla base del nuovo dataset costruito con le predizioni per le **Y di test**.

```
In [19]: def linear_regression_preds(Y_train, y_train, Y_test):
    preds = []
    for Y_train, Y_test, y_train in zip(Y_train, Y_test, y_train):
        linear_regression(fit_intercept=True).fit(Y_train, y_train).predict(Y_test)
    preds = np.array(preds, dtype = object)
    return preds

In [20]: y_preds = linear_regression_preds(Y_train, y_train, Y_test)
```

Sulla base del confronto tra predizioni e valori attesi calcolo l'errore.

```
In [21]: def get_mse(y_trues, y_preds):
    mse = []
    for y_true, y_pred in zip(y_trues, y_preds):
        mse.append(mean_squared_error(y_true, y_pred))
    mse = np.array(mse, dtype = object)
    return mse

In [22]: mse = get_mse(y_test, y_preds)
```

```
In [23]: for i in range(3):
    print(region_names[i])
    print(mse[i])
    print()
```

```
A
0.06169797072856968

B
0.0021310291641643032

C
0.0022809873002168383
```

Confronto con l'errore del modello banale

```
In [24]: def mse_train(y_train, y):
    y_pred = np.repeat(y_train.mean(), len(y_train))
    return mse(y_train, y_pred)

In [25]: for i in range(3):
    print(region_names[i])
    print(mse_train(y_test[i].values.ravel(), y_preds[i]))
    print()
```

```
A
0.02757366586548814

B
0.01915225064940997

C
0.027546206852880815
```

Analisi dei risultati

Rispetto al modello banale gli errori del metamodello sono circa dieci volte più precise: seppur il modello non sia molto preciso, combinare le informazioni di più modelli sembra essere risultato vantaggioso.

Migliori Predizioni

```
In [46]: def best_prediction(X, y_trues, y_preds, index, nrow=10):
    err = abs(y_trues[index] - y_preds[index])
    err = err.values.ravel()
    args = err.argsort()
    for i in range(nrow):
        print(
            "Row : %d, logerror : %2.7f error: %2.7f" % (
                args[i],
                (y_trues[index].iloc[args[i]] - y_preds[index]).iloc[args[i]]
            )
        )
    return X[index].iloc[args[nrow]].transpose()
```

Prima Contea 1286

```
In [47]: best_prediction(X_test, y_test, y_preds, 0, nrow=5)

Row : 4204, logerror : 0.0469000 error: 0.0000018
Row : 7126, logerror : -0.0182000 error: 0.0000019
Row : 7055, logerror : -0.0050000 error: 0.0000040
Row : 7353, logerror : 0.0198000 error: 0.0000044
Row : 4263, logerror : 0.0046090 error: 0.0000059
```

	4204	7126	7055	7353	4263
bedroomcnt	2.000000e+00	2.000000e+00	3.000000e+00	2.500000e+00	1.000000e+00
bedroom	4.000000e+00	3.000000e+00	4.000000e+00	3.000000e+00	4.000000e+00
buildingqualitytypeid	7.000000e+00	7.000000e+00	7.000000e+00	7.000000e+00	7.000000e+00
calculatedbathnr	2.000000e+00	2.000000e+00	3.000000e+00	2.500000e+00	1.000000e+00
calculatedfinishedsquarefeet	1.274000e+03	1.264000e+03	2.127000e+03	1.345000e+03	8.220000e+02
finishedsquarefeet12	1.274000e+03	1.264000e+03	2.127000e+03	1.345000e+03	8.220000e+02
fireplacecnt	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
latitude	3.36740e+07	3.37244e+07	3.361572e+07	3.358339e+07	3.377900e+07
longitude	1.178550e+08	1.180366e+08	1.176718e+08	1.177196e+08	1.179700e+08
lotsizesquarefeet	7.250000e+03	6.030000e+03	8.710000e+03	2.550000e+03	7.205000e+03
rawcensustractandblock	6.059003e+07	6.059100e+07	6.059032e+07	6.059034e+07	6.059088e+07
regionidcity	5.265000e+04	5.221800e+04	1.277300e+03	3.708600e+04	4.296700e+04
regionidzip	9.694700e+04	9.696700e+04	9.699500e+04	9.692400e+04	9.692300e+04
roomcnt	0.000000e+00	0.000000e+00	5.000000e+00	0.000000e+00	0.000000e+00
unitcnt	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00
yearbuilt	2.006000e+03	1.960000e+03	1.968000e+03	1.990000e+03	1.983000e+03
structuraevaluatedollarcnt	2.500940e+05	5.604800e+04	1.442240e+05	1.961390e+05	5.997000e+04
taxvaluatedollarcnt	4.617750e+05	4.641480e+05	3.262100e+05	3.583339e+05	3.377900e+05
landtaxvaluatedollarcnt	2.116810e+05	4.081000e+05	1.778400e+05	1.922440e+05	1.472380e+05
taxamount	5.697100e+03	5.336500e+03	3.316100e+03	3.246380e+03	2.674800e+03
int_transactiondate	2.630000e+02	2.030000e+02	2.630000e+02	1.810000e+02	8.600000e+01
period_mean_price	1.852401e-02	8.281449e-03	1.852401e-02	8.717589e-03	8.947002e-03
neighborhood_mean_price	1.916152e-02	1.168281e-02	1.227391e-02	1.051390e-02	1.159526e-02
living_area_prop	2.185046e-01	2.096168e-01	2.442021e-01	5.274510e-01	2.185004e-01
tax_ratio	7.106445e+01	6.897131e+01	9.713341e+01	6.918783e+01	6.847817e+01
tax_prop	1.181466e+00	1.373389e+01	8.107756e+01	6.580752e+01	4.079756e+01
buildingqualitytypeid_na_flag	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00
unitcnt_na_flag	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00
assessmentyear_2015.0	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00
heatingingsystemtypeid.1.0	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
heatingingsystemtypeid.6.0	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
heatingingsystemtypeid.7.0	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
heatingingsystemtypeid.10.0	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
heatingingsystemtypeid.11.0	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
heatingingsystemtypeid.12.0	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
heatingingsystemtypeid.13.0	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
heatingingsystemtypeid.18.0	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
heatingingsystemtypeid.24.0	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
poolcnt.1.0	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
propertycountylandusecode.0100	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
propertycountylandusecode.122	0.000000e+00	1.000000e+00	1.000000e+00	0.000000e+00	0.000000e+00
propertycountylandusecode.34	0.000000e+00	0.000000e+00	0.000000e+00	1.000000e+00	0.000000e+00
propertycountylandusecode.39	1.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
propertycountylandusecode.31.0	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
propertycountylandusecode.246.0	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
propertycountylandusecode.247.0	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
propertycountylandusecode.248.0	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
propertycountylandusecode.260.0	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
propertycountylandusecode.261.0	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
propertycountylandusecode.263.0	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
propertycountylandusecode.264.0	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
propertycountylandusecode.265.0	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
propertycountylandusecode.266.0	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
propertycountylandusecode.267.0	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
propertycountylandusecode.269.0	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00

Condiderazioni:

- Il valore assoluto del `logerror` è relativamente piccolo.
- la grandezza di 4 delle 5 case è simile come indicano `bedroomcnt`, `bathroom` e `lotsizesquarefeet`. Una casa ha un numero contenuto di stanze e una superficie decisamente più piccola rispetto alle altre.
- tutte le case hanno lo stesso `buildingqualitytypeid` pari a 7
- tutte le case sono di città e differenti, due appartengono alla stessa regione.
- tutte le case hanno lo stesso `rawcensustractandblock`.
- il `period_mean_price` è molto variabile, mentre la variabilità del `neighborhood_mean_price` è contenuta.
- tutte e cinque le case sono sprovviste di un impianto antincendio (`fireplacecnt` = 0) e sono sprovviste di piscina (<

