

Primo Modell

In questo notebook sono usati `LogisticRegression` e la **selezione delle feature**.

In prima istanza è allenato un

Per entrambi i modelli si prevede un algoritmo per fare **tuning degli iperparametri**: rispettivamente il **massimo numero di foglie** con l'iperparametro `max_leaf_nodes` per l'albero di decisione e il **numero di modelli** implementati con l'iperparametro `n_estimators` per il

l'iperparametro `max_leaves` per l'albero di decisione e il **numero di modelli impiegati** con l'iperparametro `n_estimators` per il modello di boosting. Per fare tuning dei parametri si fa uso di un **dataset di validation**: si mantiene l'iperparametro che produce il minimo errore su validation.

```
# 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.tree import DecisionTreeRegressor
```

```
from sklearn.preprocessing import Pipeline
from sklearn.utils import resample
from sklearn.ensemble import AdaBoostRegressor

warnings.filterwarnings('ignore')
```

Letture dei dati

```
# local file paths

dir_name = 'selezione'
region_names = np.array(['A', 'B', 'C'])

fp_x_train = []
fp_x_val = []
fp_x_test = []
fp_y_train = []
fp_y_val = []
fp_y_test = []

for i in range(3):
    fp_x_train.append(dir_name + f'/X_train{region_names[i]}.csv')
    fp_x_val.append(dir_name + f'/X_val{region_names[i]}.csv')
    fp_x_test.append(dir_name + f'/X_test{region_names[i]}.csv')
    fp_y_train.append(dir_name + f'/y_train{region_names[i]}.csv')
    fp_y_val.append(dir_name + f'/y_val{region_names[i]}.csv')
    fp_y_test.append(dir_name + f'/y_test{region_names[i]}.csv')
```

Letture dei dataset su cui è stata fatta la selezione delle feature.

```
X_train = []
X_val   = []
X_test  = []
y_train = []
y_val   = []
```

```

x_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_val[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.append(pd.read_csv(fp_ytrain[i], low_memory=False))
    y_val.append(pd.read_csv(fp_val[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)

def dimensionality(y=False):
    for i in range(3):
        print('X_train(region_names[i]): (X_train[i].shape)')
        print('X_val(region_names[i]): (X_val[i].shape)')
        print('X_test(region_names[i]): (X_test[i].shape)')
        if y:
            print('y_train(region_names[i]): (y_train[i].shape)')
            print('y_val(region_names[i]): (y_val[i].shape)')
            print('y_test(region_names[i]): (y_test[i].shape)')
        print()

dimensionality(y=True)

X_traina: (26819, 55)

```

```
X_testA: (9085, 55)
y_trainA: (26819, 1)
y_testA: (9006, 1)
y_testA: (9085, 1)

X_trainB: (8119, 32)
X_valB: (2658, 32)
X_testB: (2606, 32)
y_trainB: (8119, 1)
y_valB: (2658, 1)
y_testB: (2606, 1)

X_trainC: (64771, 33)
X_valC: (21908, 33)
X_testC: (21876, 33)
y_trainC: (64771, 1)
y_valC: (21908, 1)
y_testC: (21876, 1)

Variabili globali

# NomI delle regioni
region_ida = np.array(['1206', '2061', '3101'])

# Riduzione del dataset C
C_PERC = 2/5
C_THD = 2

# Percentuale per la sottoinsieme su cui costruire il modello
SUB_PERC = [ 1/3, 1, 1/3] # A: 8000, B: 100 C: 100 Test rapido - 5 min
SUB_PERC = [ 1/3, 1, 1/3] # A: 8000, B: 8000, C: 8000 Test medio - 3 ore
```

```

train_PERC = [ 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100]
# Decision Tree Regressor
DTR_START = 10
DTR_END = 300
DTR_STEP = 10

# AdaBoost Regressor
BOOST_START = 20
BOOST_END = 500
BOOST_STEP = 20
BOOST_NTEST = 20
#BOOST_SAMPLE_PERC = 1/3 # 0.33
BOOST_SAMPLE_PERC = 2/3 # 0.67

```

```
dimensionality(y=True)
X_trainA: (26819, 55)
X_valA:   (9006, 55)
```

```
X_testA: (9085, 55)
y_trainA: (26819, 1)
y_valA: (9006, 1)
y_testA: (9085, 1)
```

```
X_trainB: (8119, 32)
X_valB: (2658, 32)
X_testB: (2606, 32)
y_trainB: (8119, 1)
y_valB: (2658, 1)
y_testB: (2606, 1)

X_trainC: (25908, 33)
X_valC: (8763, 33)
```

```
x_trainC: (8750, 33)
y_trainC: (25908, 1)
y_valC: (8763, 1)
y_testC: (8750, 1)
```

Gli algoritmi usati in questo notebook hanno un costo computazionale elevato, per questo è definito un sottoinsieme dei dataset originali su cui far girare gli algoritmi; questo per facilitare la fase di creazione e di testing e ottenere risultati verosimili in tempi utili a verificare il corretto funzionamento del processo.

Nella versione finale gli algoritmi useranno circa 8000 istanze del dataset originale, poiché il dataset allenato nella sua totalità richiede per il boosting tempi non sostenibili.

```
for i in range(3):
    print(f'Region {region_names[i]}')
    print(int(len(X_train[i])*SUB_PERC[i]))
    print(int(len(X_val[i])*SUB_PERC[i]))
    print()
```

```
Region A
8939
3902

Region B
8119
2658

Region C
8636
2921
```

```
X_train_sub = []
y_train_sub = []
X_val_sub = []
y_val_sub = []
```

```
for i in range(3):
    Xt_sub, yt_sub = resample(X_train[i], y_train[i], n_samples = int(SUB_PERCENT * X_train[i].shape[0]))
    Xv_sub, yv_sub = resample(X_val [i], y_val [i], n_samples = int(SUB_PERCENT * X_val [i].shape[0]))
    X_train_sub.append(Xt_sub)
    y_train_sub.append(yt_sub)
    X_val_sub.append(Xv_sub)
    y_val_sub.append(yv_sub)
```

```

        y_val_sub = append(yv_sub)

def dimensionality_sub(y=False):
    for i in range(3):
        print(f'X_train_sub(region_names[1]): {X_train_sub[1].shape}')
        if y:
            print(f'X_val_sub(region_names[1]): {X_val_sub [1].shape}')
        print(f'y_train_sub(region_names[1]): {y_train_sub[1].shape}')
        print(f'y_val_sub(region_names[1]): {y_val_sub [1].shape}')
    print()

```

```
dimensionality_sub(y=True)

X_train_subA: (8939, 55)
X_val_subA: (3002, 55)
y_train_subA: (8939, 1)
y_val_subA: (3002, 1)

X_train_subB: (8119, 32)
```

```

y_val_subB: (2658, 32)
y_train_subB: (8119, 1)
y_val_subC: (2658, 1)

X_train_subC: (8636, 33)
X_train_subC: (2921, 33)
y_train_subC: (8636, 1)
y_val_subC: (2921, 1)

```

DecisionTreeRegressor

Tuning del massimo numero di foglie attraverso il parametro `max_leaf_nodes`. Per studiare l'andamento la produzione vuol grafi che illustrano l'andamento di *bias*, *varianza* e *mean squared error* in funzione del massimo numero di foglie sia per l'insieme di Train che per quello di Validation.

```

plt.rcParams.update({'font.size': 35})

```

```

def get_bias_var_mse(X, y, model):
    y_pred = model.predict(X)
    return {
        'bias': ((y - np.mean(y_pred))**2).mean(), \
        'var': np.var(y_pred).mean(), \
        'mse': ((y_pred - y.reshape(-1,1))**2).mean()
    }

```

```

# Costruzione Decision TreeRegressor
def DecisionTreeRegressor_validation(X_train, y_train, X_val, y_val, verbose=False, debug=False, file_name=''):
    def get_dec_tree_rec(max_leaf):

```

```

dt = DecisionTreeRegressor(
    max_leaf_nodes = max_leaf,
    criterion = 'labeling_error'
)

dt.fit(X_train, y_train)

return dt

def bias_var_mse(X, y, model):
    stats = get_bias_var_mse(X, y, model)
    return stats['bias'],\
           stats['var'],\
           stats['mse']

def plot_mse(stats, name):

    print ("{}: TUNING DEL MASSIMO NUMERO DI FOGLIE".format(name))
    print()

    for n in ['mse', 'bias', 'var']:

        min_ = min(stats[n])
        best = (np.argmax(stats[n]) + DTR_STEP) + DTR_START

        print ("Punteggio finale: {} (stats[n] - 1) (DTR_END) MaxLeaves)".format(
            best - 1, min_))
        print ("Best number of MaxLeaves: (best)")
        print()

    fig, ax = plt.subplots(figsize=(len(stats['mse'])/2, 10))

    ax.tick_params(axis='both', which='major', labelsize=25)
    ax.tick_params(axis='both', which='minor', labelsize=15)

    ax.plot_range(DTR_START, DTR_END+1, DTR_STEP), stats['mse'], 'o-', label='MSE')

```

```
ax.plot(range(DTR_START, DTR_END+1, DTR_STEP), stats['bias'], 'o-', label='BIAS')
ax.plot(range(DTR_START, DTR_END+1, DTR_STEP), stats['var'], 'o-', label='VARIANCE')

ax.set_title(f'({name} MSE, BIAS, VARIANCE on different MaxLeaves', fontsize=15)
ax.set_xlabel('Number of Max Leaves used', fontsize=15)
ax.grid()
ax.legend(prop={'size': 12})

if file_name != '':
    fig.savefig('Images/' + file_name + '_DecisionTreeRegressor_' + name + '.jpg')

y_train = y_train.values.ravel()
y_val = y_val.values.ravel()

first = True

info = {}

train_stats = {
    'bias': [],
    'var': [],
    'mse': []
}

val_stats = {
    'bias': [],
    'var': [],
    'mse': []
}

for max_leaf in range(DTR_START, DTR_END+1, DTR_STEP):

    if debug:
        print(f'({max_leaf}/{DTR_END})')
```

```

model = get_dec_tree_max_leaf_train

trn_bias, trn_var, trn_mse = bias_var_mse(X_train, y_train, model)
val_bias, val_var, val_mse = bias_var_mse(X_val, y_val, model)

trn_stats["bias"].append(trn_bias)
trn_stats["var"].append(trn_var)
trn_stats["mse"].append(trn_mse)

val_stats["bias"].append(val_bias)
val_stats["var"].append(val_var)
val_stats["mse"].append(val_mse)

info.append('Max Leaf: %max_leaf' % \
            f'\n(train MSE: (trn_mse) - Val MSE: (val_mse) )' % \
            f'\n(Train Bias: (trn_bias) - Val Bias: (val_bias))' % \
            f'\n(Train Variance: (trn_var) - Val Variance: (val_var) )' )

if debug:
    print(f' %max_leaf' / DTR_END')

if (first or val_mse < best_mse):
    first = False
    best_mse = val_mse
    best_max_leaf = max_leaf
    best_model = model

if verbose:
    print()
    print(f'(name) MSE, BIAS, VARIANCE Train e Validation')
    print(f'Info, sep='\n')
    print()

```

```

    plot_mse_val_stats, train
    plot_mse_val_stats, "Validation")

    return best_model

dt_model = {}

def get_dt(index, verbose=False, debug=False, file_name=''):
    if file_name == '':
        file_name = region_ids[index]
    return DecisionTreeRegressorValidation(
        X_train_sub(index),
        y_train_sub(index),
        X_val_sub(index),
        y_val_sub(index),
        verbose = verbose,
        debug = debug,
        file_name = file_name
    )

Prima Contea 1286

%time
dt_model.append(
    get_dt(
        0,
        verbose = False,
        debug = False,
        file_name = 'prova'
    )
)

```

```

Train: TUNING DEL MASSIMO NUMERO DI FOGLIE
Punteggio finale: 0.01078387673563775 (300) MaxLeaves
Best mse: 0.007326822203794711
Best number of MaxLeaves: 10

Punteggio finale: 0.0069320060368185344 (300) MaxLeaves
Best bias: 0.0069320060368185344
Best number of MaxLeaves: 10

Punteggio finale: 0.0038518706988191907 (300) MaxLeaves
Best var: 0.00840427616897616185
Best number of MaxLeaves: 10

Validation: TUNING DEL MASSIMO NUMERO DI FOGLIE

Punteggio finale: 0.03821001828921705 (300) MaxLeaves
Best mse: 0.03863314023928795
Best number of MaxLeaves: 10

Punteggio finale: 0.03364000821357107 (300) MaxLeaves
Best bias: 0.03361805250079998
Best number of MaxLeaves: 10

Punteggio finale: 0.0045700106493506504 (300) MaxLeaves
Best var: 0.00023888533095541306
Best number of MaxLeaves: 10

Wall time: 25.2 s

```

0.010

The figure consists of two line plots. The top plot shows the validation MSE, BIAS, and VARIANCE on the y-axis (ranging from 0.000 to 0.008) against the Number of Max Leaves on the x-axis (ranging from 0 to 300). The bottom plot shows the validation MSE, BIAS, and VARIANCE on the y-axis (ranging from 0.030 to 0.040) against the Number of Max Leaves on the x-axis (ranging from 0 to 300). Both plots include a legend for MSE (blue line with circles), BIAS (orange line with circles), and VARIANCE (green line with circles).

Top Plot Data (Approximate):

Number of Max Leaves	MSE	BIAS	VARIANCE
10	0.0075	0.0070	0.0005
50	0.0080	0.0070	0.0015
100	0.0078	0.0070	0.0022
150	0.0075	0.0070	0.0026
200	0.0072	0.0070	0.0031
250	0.0070	0.0070	0.0036
300	0.0068	0.0070	0.0039

Bottom Plot Data (Approximate):

Number of Max Leaves	MSE	BIAS	VARIANCE
10	0.0345	0.0340	0.0345
50	0.0355	0.0340	0.0345
100	0.0360	0.0340	0.0345
150	0.0362	0.0340	0.0345
200	0.0365	0.0340	0.0345
250	0.0368	0.0340	0.0345
300	0.0370	0.0340	0.0345

Seconda Contea 2061

Number of Max Leaves used	Value
0	0.000
10	0.000
20	0.001
30	0.000
40	0.001
50	0.002
60	0.003
70	0.002
80	0.003
90	0.003
100	0.003
110	0.003
120	0.003
130	0.003
140	0.003
150	0.003
160	0.004
170	0.004
180	0.004
190	0.003
200	0.003
210	0.004
220	0.004
230	0.005
240	0.004
250	0.005
260	0.004
270	0.005
280	0.005
290	0.004
300	0.005

```
get_at(
    1,
    verbose =
    debug =
    #file_name
)
```

```

Train: TUNING DEL MASSIMO NUMERO DI FOGLIE
Punteggio finale: 0.0121171755678008781 (300) MaxLeaves
Best mae: 0.07864440748399235
Best number of MaxLeaves: 10

Punteggio finale: 0.00740842968893943 (300) MaxLeaves
Best bias: 0.00740842968893943
Best number of MaxLeaves: 60

Punteggio finale: 0.004703325989614886 (300) MaxLeaves
Best var: 0.000456010580042893
Best number of MaxLeaves: 10

Validation: TUNING DEL MASSIMO NUMERO DI FOGLIE

Punteggio finale: 0.022142948201747773 (300) MaxLeaves
Best mae: 0.018816528962922508
Best number of MaxLeaves: 10

Punteggio finale: 0.015206957949504602 (300) MaxLeaves
Best bias: 0.015205301853708214
Best number of MaxLeaves: 220

Punteggio finale: 0.00695601025243171 (300) MaxLeaves
Best var: 0.0006105841246927085
Best number of MaxLeaves: 10

Wall time: 20.4 s

```

The graph displays two metrics, Bias and Variance, plotted against the Number of Max Leaves used (ranging from 0 to 300). The Y-axis represents the magnitude of these metrics, ranging from 0.000 to 0.012.

Bias (Orange line with 'x' markers): The bias remains constant across all values of Max Leaves, starting at approximately 0.0075 and ending at approximately 0.0075.

Variance (Green line with '+' markers): The variance starts at approximately 0.0005 for 10 Max Leaves and increases steadily as the number of Max Leaves increases, reaching approximately 0.0045 at 300 Max Leaves.

Number of Max Leaves used	Bias	Variance
10	0.0075	0.0005
50	0.0075	0.0018
100	0.0075	0.0028
150	0.0075	0.0035
200	0.0075	0.0040
250	0.0075	0.0043
300	0.0075	0.0045

Validation MSE, BIAS, VARIANCE on different MaxLeaves

MaxLeaves	MSE	BIAS	VARIANCE
100	0.0158	0.015	0.0005
110	0.0160	0.015	0.0010
120	0.0165	0.015	0.0015
130	0.0175	0.015	0.0025
140	0.0185	0.015	0.0035
150	0.0175	0.015	0.0025
160	0.0185	0.015	0.0035
170	0.0190	0.015	0.0040
180	0.0190	0.015	0.0045
190	0.0195	0.015	0.0045
200	0.0190	0.015	0.0040
210	0.0205	0.015	0.0050
220	0.0215	0.015	0.0055
230	0.0210	0.015	0.0065
240	0.0220	0.015	0.0060
250	0.0205	0.015	0.0055
260	0.0210	0.015	0.0060
270	0.0225	0.015	0.0075
280	0.0225	0.015	0.0070
290	0.0225	0.015	0.0070
300	0.0220	0.015	0.0065
310	0.0225	0.015	0.0070

```

0          50          100          150          200          250          300
Number of Max Leaves used

Terza Conta 3101

%%time
dt_model.append(
    get_dt(
        2,
        verbose      = False,
        debug         = False,
        file_name     = 'prova2'
    )
)

Train: TUNING DEL MASSIMO NUMERO DI FOGLIE

Punteggio finale: 0.018429791191878 (300) MaxLeaves
Best mso: 0.012663186996844945
Best number of MaxLeaves: 10

Punteggio finale: 0.012049414522519223 (300) MaxLeaves
Best bias: 0.012049414522519233
Best number of MaxLeaves: 10

Punteggio finale: 0.0063803766693586355 (300) MaxLeaves
Best var: 0.000613772474325895
Best number of MaxLeaves: 10

Validation: TUNING DEL MASSIMO NUMERO DI FOGLIE

Punteggio finale: 0.02857894831510075 (300) MaxLeaves

```

Best mse: 0.0201017095137600226
Best number of MaxLeaves: 10

Punteggio finale: 0.02030524671354402 (300) MaxLeaves
Best bias: 0.020273881054273456
Best number of MaxLeaves: 280

Punteggio finale: 0.008273701601756759 (300) MaxLeaves
Best var: 0.0007311905581820703
Best number of MaxLeaves: 10

Wall time: 22.4 s

Train MSE, BIAS, VARIANCE on different MaxLeaves

MaxLeaves	MSE	BIAS	VARIANCE
10	0.0125	0.0115	0.0115
20	0.0128	0.0115	0.0115
30	0.0132	0.0115	0.0115
40	0.0135	0.0115	0.0115
50	0.0138	0.0115	0.0115
60	0.0142	0.0115	0.0115
70	0.0145	0.0115	0.0115
80	0.0148	0.0115	0.0115
90	0.0152	0.0115	0.0115
100	0.0155	0.0115	0.0115
120	0.0158	0.0115	0.0115
140	0.0162	0.0115	0.0115
160	0.0165	0.0115	0.0115
180	0.0168	0.0115	0.0115
200	0.0170	0.0115	0.0115
220	0.0172	0.0115	0.0115
240	0.0174	0.0115	0.0115
260	0.0175	0.0115	0.0115
280	0.0175	0.0115	0.0115
300	0.0175	0.0115	0.0075

The figure consists of two line plots. The top plot shows the relationship between the Number of Max Leaves used (x-axis, 0 to 300) and the Validation MSE (y-axis, 0.0000 to 0.0050). The bottom plot shows the relationship between the Number of Max Leaves used (x-axis, 0 to 300) and the Validation BIAS (y-axis, 0.015 to 0.030). Both plots include a legend for MSE (blue line with circles), BIAS (orange line with circles), and VARIANCE (green line with circles).

Top Plot: Validation MSE vs. Number of Max Leaves used

Number of Max Leaves used	Validation MSE
20	0.0005
40	0.0010
60	0.0015
80	0.0020
100	0.0025
120	0.0030
140	0.0035
160	0.0040
180	0.0045
200	0.0050
220	0.0055
240	0.0060
260	0.0065
280	0.0070
300	0.0075

Bottom Plot: Validation BIAS vs. Number of Max Leaves used

Number of Max Leaves used	Validation BIAS
20	0.0210
40	0.0215
60	0.0220
80	0.0225
100	0.0230
120	0.0235
140	0.0240
160	0.0245
180	0.0250
200	0.0255
220	0.0260
240	0.0265
260	0.0270
280	0.0275
300	0.0280

Number of Max Leaves used	dt_model
10	0.0008
20	0.0012
30	0.0018
40	0.0022
50	0.0028
60	0.0032
70	0.0035
80	0.0032
90	0.0032
100	0.0042
110	0.0032
120	0.0055
130	0.0042
140	0.0050
150	0.0052
160	0.0055
170	0.0065
180	0.0055
190	0.0065
200	0.0062
210	0.0062
220	0.0065
230	0.0070
240	0.0060
250	0.0072
260	0.0068
270	0.0075
280	0.0065
290	0.0090
300	0.0085

```
dt_model

[DecisionTreeRegressor(max_leaf_nodes=10),
 DecisionTreeRegressor(max_leaf_nodes=10),
 DecisionTreeRegressor(max_leaf_nodes=10)]
```

Visualizzazione degli alberi

```
def plot_dt(index, file_name=''):
    if file_name == '':
        file_name = region_id[index]
    fig, ax = plt.subplots(figsize=(20,10))
    plot_tree(
        dt_model[index],
        ax = ax,
```

[illegible][illegible][illegible]

Analisi dei risultati

I risultati ottenuti sono analoghi per tutti e tre le contee: i risultati non sono per nulla soddisfacenti.

All'aumentare del numero di foglie, l'errore aumenta.

Decomponendo l'errore di più notare come questo errore provenga da una **varianza** che inevitabilmente tende a **crescere aumentando il numero di foglie**, e da un **bias** che al posto che diminuire specializzandosi nei vari nodi **rimane costante**.

Il modello non riesce a imparare dai dati specializzandosi nei nodi, ma probabilmente fornisce indipendentemente dal nodo una predizione banale pari circa alla media delle osservazioni.

Analizzando l'albero si riesce a visualizzare come la struttura tende a specializzarsi su una piccola fetta dei dati, lasciando un'unica foglia per la stragrande maggioranza delle osservazioni per cui la predizione è molto simile alla media della totalità del dataset.

Il modello non riesce così ad abbattere le baseline del modello banale.

```
In [27]: def print_stats(X, y, models):
    for i in range(len(models)):
        print(f"region_names[{i}]: (get_bias_var_mse(X[i], y[i].values.ravel(), models[i]))")
        print()
```

```
In [28]: def print_stats_trivial(y_train, y):
    y_pred = np.repeat(y_train.mean(), len(y_train))
    print(
        'bias': ((y - np.mean(y_pred))**2).mean(), \
        'mse': np.var(y_pred).mean(), \
        'var': ((y_pred - y.reshape(-1,1))**2).mean()
    )
```

```
In [29]: def print_all_stats(models):
    for X, y, name in zip(
        [X_train, X_val, X_test],
        [y_train, y_val, y_test],
        ['Train', 'Validation', 'Test']
    ):
        print(name)
        print()
        print_stats(X, y, models)
        print()
```

```
In [30]: def print_all_stats_trivial():
    for i in range(len(models)):
        print(region_names[i])
        print_stats_trivial(y_train[i].values.ravel(), y_test[i].values.ravel())
        print()
```

Statistiche per le tre contee su Train, Validation e Test

```
In [31]: print_all_stats(dt_model)
```

Train

A: ('bias': 0.0069294254023030055, 'var': 0.0004273882859140359, 'mse': 0.0073568136882170955)
B: ('bias': 0.007536254902481718, 'var': 0.000345142012739853424, 'mse': 0.007881396915220221)
C: ('bias': 0.0116745081204912, 'var': 0.0005817505209895822, 'mse': 0.012526528641480767)

Validation

A: ('bias': 0.0272663654076649438, 'var': 0.00052734703718553495, 'mse': 0.027793701107504842)
B: ('bias': 0.019614360581977596, 'var': 0.0005806838132700375, 'mse': 0.02019504439324598)
C: ('bias': 0.02217266495119075, 'var': 0.0005495066299329632, 'mse': 0.022722171581123617)

Test

A: ('bias': 0.028641724215576377, 'var': 0.00030789391837899747, 'mse': 0.02894961813395547)
B: ('bias': 0.0241581039782064934, 'var': 0.0002547132513325053, 'mse': 0.02441282303397434)
C: ('bias': 0.03776212365272005, 'var': 0.0004758397461138191, 'mse': 0.03823796339883399)

```
In [32]: for i in range(3):
    print(f"region_names[{i}]:")
    print(pd.DataFrame(y_train[i]).describe())
    print()
```

```
count      logerror
count  26819.000000
mean      0.012913
std       0.083244
min      -0.482335
25%      -0.017767
50%      0.006548
75%      0.034400
max       0.753300

B
count      logerror
count  8119.000000
mean      0.031357
std       0.086804
min      -0.470000
25%      -0.020200
50%      0.006000
75%      0.036300
max       0.756600

C
count      logerror
count  25908.000000
mean      0.012105
std       0.108051
min      -0.485174
25%      -0.029400
50%      0.005973
75%      0.041100
max       0.751627
```

L'errore quadratico medio è decisamente molto grande rispetto alla dimensionalità della risposta.

Statistiche per un modello banale sul test

```
In [33]: print_all_stats_trivial()
```

A: ('bias': 0.028637396765634219, 'var': 1.2037062152420224e-35, 'mse': 0.02863739676563433)
B: ('bias': 0.02412642449493516, 'var': 3.009265538105056e-36, 'mse': 0.024126424499493482)
C: ('bias': 0.03775923001036233, 'var': 3.009265538105056e-36, 'mse': 0.03775923001036172)

Gli errori dell'albero di decisione sono paragonabili a quelli del modello banale.

L'errore sembra provenire principalmente dal bias: si può pensare di sfruttare il **boosting** per ridurlo.

Boosting

Sceglia l'algoritmo di boosting principalmente per due motivi:

- è un processo che tende a ridurre il **bias** di un modello e l'**errore** degli alberi di decisione **provenienti** principalmente dal **bias**.
- è capace di raggiungere un buono score anche partendo da **weak learners**, come quelli selezionati dall'albero di decisione

Usa come modelli base gli alberi individuati ai punti precedenti

```
In [34]: base_model = []
#max_leaf = 500
base_model.append(DecisionTreeRegressor(max_leaf_nodes = dt_get_params()['max_leaf_nodes']))
#base_model.append(DecisionTreeRegressor(max_leaf_nodes = MAX_LEAF))
base_model
```

```
Out[34]: DecisionTreeRegressor(max_leaf_nodes=10),
DecisionTreeRegressor(max_leaf_nodes=10),
DecisionTreeRegressor(max_leaf_nodes=10)
```

```
In [35]: def boosting_train(X_train, y_train, X_val, y_val, baseModel, verbose=False, debug=False, file_name=''):
    adaBoost = AdaBoostRegressor(
        baseModel,
        n_estimators = estimators
    )
    adaBoost.fit(Xs,ys)
    return adaBoost

def bias_var_mse(y, y_model):
    stats = get_bias_var_mse(X, y, model)
    return stats['bias'],\
           stats['var'],\
           stats['mse']

def plot_mse(stats, name):
    print(f"({name}): TUNING DEL NUMERO DI STIMATORI")
    print()
    for n in ['mse', 'bias', 'var']:
        min_n = min(stats[n])
        best = (np.argmax(stats[n]) * BOOST_STEP) + BOOST_START

        print(f"Punteggio finale: {stats[n][1]} ((BOOST_END) stimatori)")
        print(f"Best (n):, (min)")
        print(f"Best number of Estimators: (best)")
        print()

    fig, ax = plt.subplots(figsize=(len(stats['mse'])/2, 10))
    ax.tick_params(axis='both', which='major', labelsize=25)
    ax.tick_params(axis='both', which='minor', labelsize=15)
    ax.plot(range(BOOST_START, BOOST_END+1, BOOST_STEP), stats['mse'], 'o-', label='MSE')
    ax.plot(range(BOOST_START, BOOST_END+1, BOOST_STEP), stats['bias'], 'o-', label='BIAS')
    ax.plot(range(BOOST_START, BOOST_END+1, BOOST_STEP), stats['var'], 'o-', label='VARIANCE')
    ax.set_title(f"Number: MSE, BIAS, VARIANCE on different Estimators, fontsize=15")
    ax.set_xlabel("Number of Max Estimators used", fontsize=15)
    ax.grid()
    ax.legend(prop={'size': 12})

    if file_name != '':
        fig.savefig('images/' + file_name + '_AdaBoostRegressor_' + name + '.jpg')

stats = np.array([])
boosts = range(BOOST_START, BOOST_END+1, BOOST_STEP)

y_train = y_train.values.ravel()
y_val = y_val.values.ravel()

first = True

info = []

train_stats = {
    'bias': [],
    'var': [],
    'mse': []
}

val_stats = {
    'bias': [],
    'var': [],
    'mse': []
}

for b in boosts:
    if debug:
        print(f"({b})/(BOOST_END)")

    train_stats_s = {
        'bias': [],\
        'var': [],\
        'mse': []
    }

    val_stats_s = {
        'bias': [],\
        'var': [],\
        'mse': []
    }

    # Resampling
    for i in range(BOOST_NTEST):
        print(f"({i+1})/(BOOST_NTEST)")

        X_sample, y_sample = resample(X_train, y_train, n_samples = int(BOOST_SAMPLE_PERC*len(y_train)))
        ada = get_adaBoostRegressor(X_sample, y_sample, b)

        trn_bias, trn_var, trn_mse = bias_var_mse(X_train, y_train, ada)
        val_bias, val_var, val_mse = bias_var_mse(X_val, y_val, ada)

        train_stats_s['bias'].append(trn_bias)
        train_stats_s['var'].append(trn_var)
        train_stats_s['mse'].append(trn_mse)

        val_stats_s['bias'].append(val_bias)
        val_stats_s['var'].append(val_var)
        val_stats_s['mse'].append(val_mse)

        trn_bias_s = np.array(train_stats_s['bias']).mean()
        trn_var_s = np.array(train_stats_s['var']).mean()
        trn_mse_s = np.array(train_stats_s['mse']).mean()

        val_bias_s = np.array(val_stats_s['bias']).mean()
        val_var_s = np.array(val_stats_s['var']).mean()
        val_mse_s = np.array(val_stats_s['mse']).mean()

        train_stats['bias'].append(trn_bias_s)
        train_stats['var'].append(trn_var_s)
        train_stats['mse'].append(trn_mse_s)

        val_stats['bias'].append(val_bias_s)
        val_stats['var'].append(val_var_s)
        val_stats['mse'].append(val_mse_s)

        info.append(
            f"({b})/({b})/(BOOST_END) (trn_mse_s - Val MSE: (val_mse_s))" + "\n" +
            f"({b})/({b})/(BOOST_END) (trn_bias_s - Val Bias: (val_bias_s))" + "\n" +
            f"({b})/({b})/(BOOST_END) (trn_var_s - Val Variance: (val_var_s))" + "\n"
        )

    if (first or val_mse_s < best_mse):
        first = False
        best_mse = val_mse_s
        best_estimators = b

    if verbose:
        print(f"Info, sep='\n'")

    plot_mse(train_stats, 'Train')
    plot_mse(val_stats, 'Validation')

    return get_adaBoostRegressor(X_train, y_train, best_estimators)
```

```
In [36]: boost_model = []
```

```
In [37]: def get_bmodel(index, verbose=False, debug=False, file_name=''):
    if file_name != '':
        file_name = region_ids[index]
    return boosting_train(
        X_train_sub[index],
        y_train_sub[index],
        X_val_sub[index],
        y_val_sub[index],
        base_model,
        verbose = verbose,
        debug = debug,
        file_name = file_name
    )
```

Prima Contea 1286

```
In [38]: %time
boost_model.append(
    get_bmodel(
        1,
        verbose = False,
        debug = False,
        file_name = "Prova1"
    )
)
```

Train: TUNING DEL NUMERO DI STIMATORI

Punteggio finale: 0.0162399857197165746 (500) stimatori
Best mse, 0.00850762096412658
Best number of Estimators: 20

Punteggio finale: 0.008661629437220689 (500) stimatori
Best bias, 0.007319919027574296
Best number of Estimators: 20

Punteggio finale: 0.00763927759945063 (500) stimatori
Best mse, 0.0012506460308356
Best number of Estimators: 20

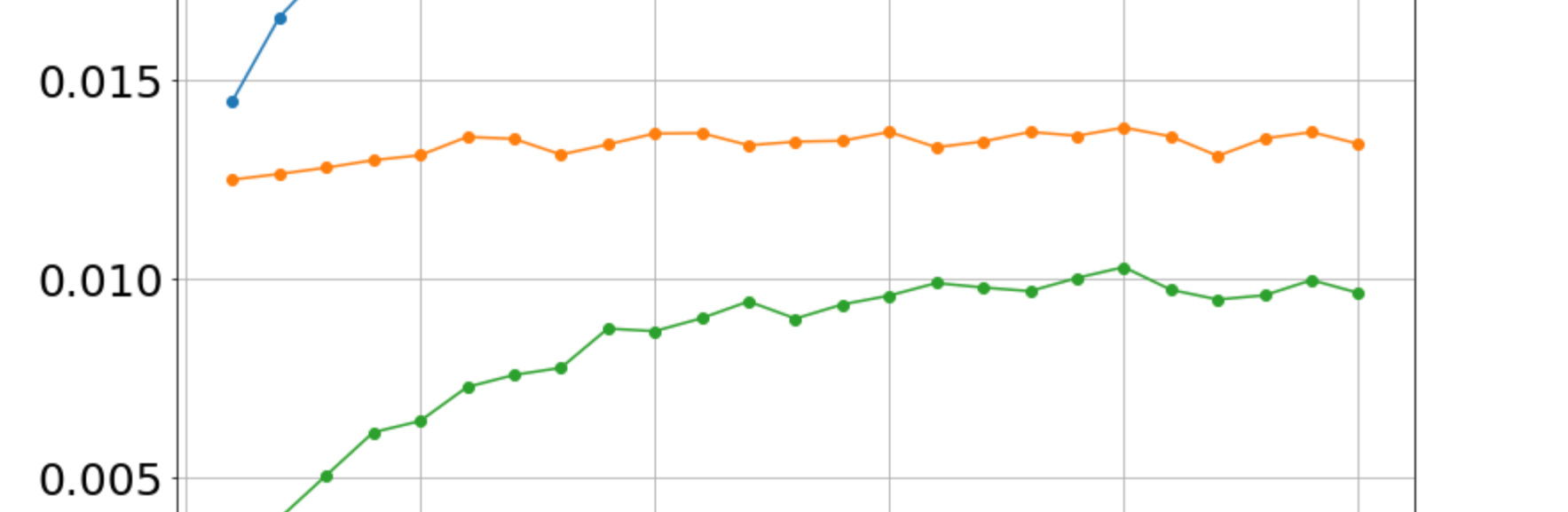
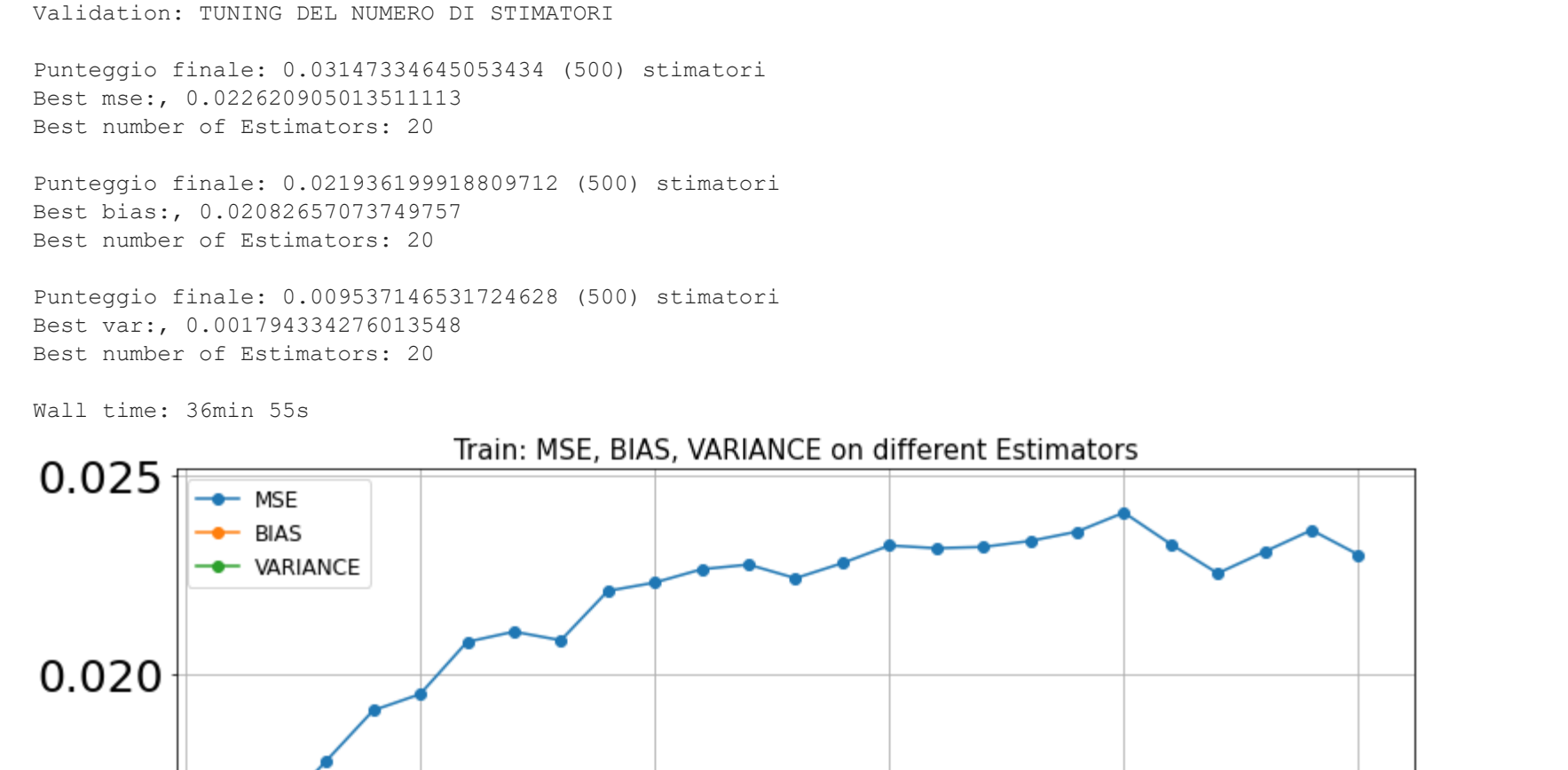
Validation: TUNING DEL NUMERO DI STIMATORI

Punteggio finale: 0.042728461646479128 (500) stimatori
Best mse, 0.0350568725915139
Best number of Estimators: 20

Punteggio finale: 0.03504045059819914 (500) stimatori
Best bias, 0.03385121053172969
Best number of Estimators: 20

Punteggio finale: 0.00768901048592151 (500) stimatori
Best var, 0.001205656728022977
Best number of Estimators: 20

Wall time: 45min 32s



Seconda Contea 2061

```
In [39]: %time
boost_model.append(
    get_bmodel(
        2,
        verbose = False,
        debug = False,
        file_name = "Prova2"
    )
)
```

Train: TUNING DEL NUMERO DI STIMATORI

Punteggio finale: 0.013847338931855728 (500) stimatori
Best mse, 0.00112868092056046
Best number of Estimators: 20

Punteggio finale: 0.009118669421700655 (500) stimatori
Best bias, 0.00791388051394728
Best number of Estimators: 20

Punteggio finale: 0.00478469520155089 (500) stimatori
Best var, 0.0011989802412615206
Best number of Estimators: 20

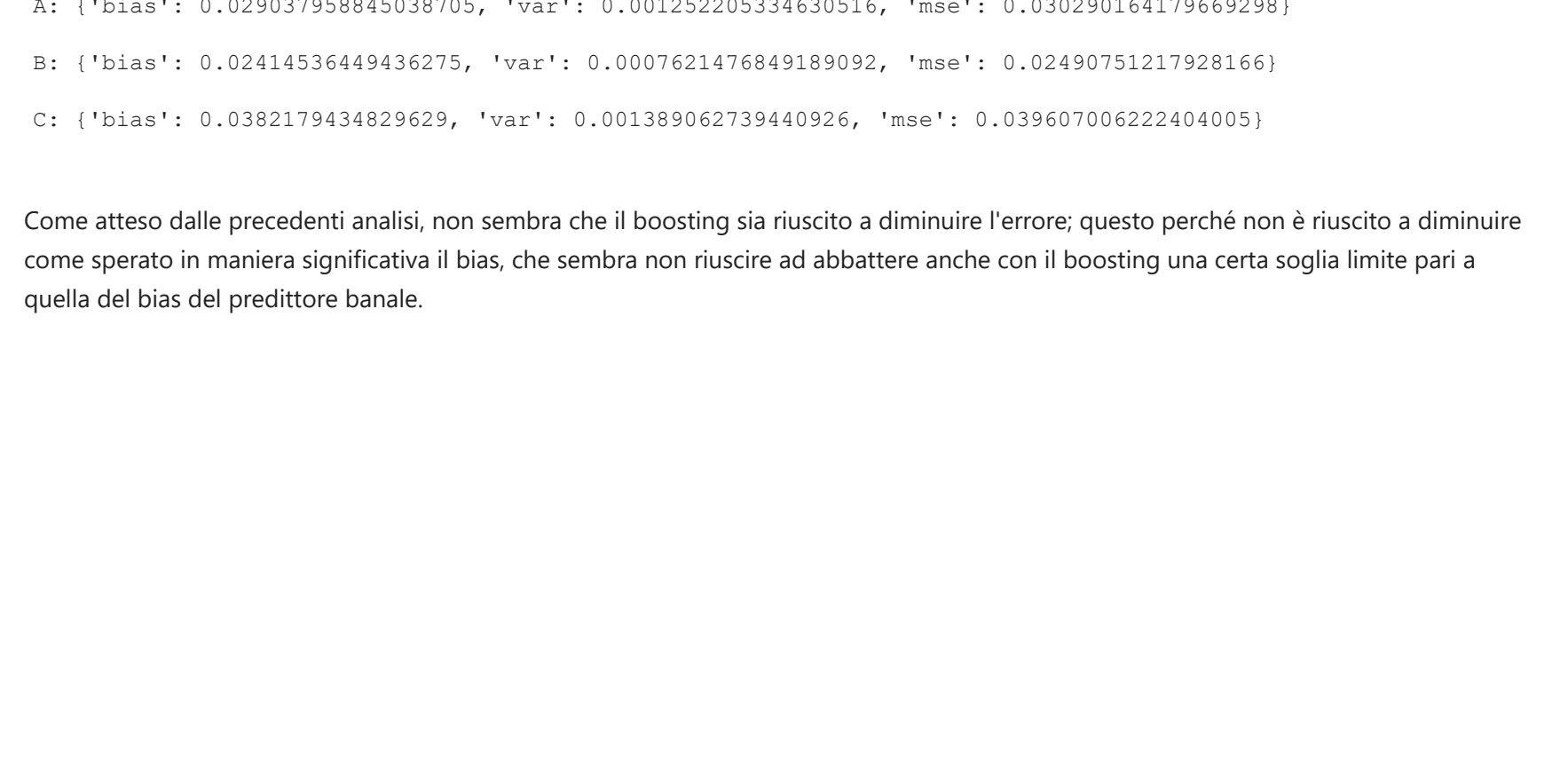
Validation: TUNING DEL NUMERO DI STIMATORI

Punteggio finale: 0.021694789718708297 (500) stimatori
Best mse, 0.01676477459170757
Best number of Estimators: 20

Punteggio finale: 0.016914292931888616 (500) stimatori
Best bias, 0.01573980100531015
Best number of Estimators: 20

Punteggio finale: 0.004780496784819669 (500) stimatori
Best var, 0.00103179229174565
Best number of Estimators: 20

Wall time: 32min 11s



Terza Contea 3101

```
In [40]: %time
boost_model.append(
    get_bmodel(
        3,
        verbose = False,
        debug = False,
        file_name = "Prova3"
    )
)
```

Train: TUNING DEL NUMERO DI STIMATORI

Punteggio finale: 0.022999225418023572 (500) stimatori
Best mse, 0.0148595069422102971
Best number of Estimators: 20

Punteggio finale: 0.01337614068989993 (500) stimatori
Best bias, 0.012476381051324944
Best number of Estimators: 20

Punteggio finale: 0.009623076728123642 (500) stimatori
Best var, 0.0019826883713788238
Best number of Estimators: 20

Validation: TUNING DEL NUMERO DI STIMATORI

Punteggio finale: 0.03147334645053434 (500) stimatori
Best mse, 0.02262090503511113
Best number of Estimators: 20

Punteggio finale: 0.021936199918809712 (500) stimatori
Best bias, 0.02082657373749757
Best number of Estimators: 20

Punteggio finale: 0.009537146531724628 (500) stimatori
Best var, 0.001794334276013548
Best number of Estimators: 20

Wall time: 36min 55s



```
In [41]: boost_model
```

AdaBoostRegressor(base_estimator=DecisionTreeRegressor(max_leaf_nodes=10),
n_estimators=20),
AdaBoostRegressor(base_estimator=DecisionTreeRegressor(max_leaf_nodes=10),
n_estimators=20),
AdaBoostRegressor(base_estimator=DecisionTreeRegressor(max_leaf_nodes=10),
n_estimators=20)

Analisi dei risultati

I risultati si distinguono da quelli precedenti: non si evidenzia una significativa diminuzione del bias, che sembra non riuscire ad abbattere una certa soglia limite. Il problema principale è analogo quello precedente: l'albero non sembra riuscire a imparare dai dati e a specializzarsi in maniera utile.

Statistiche per le tre contee su Train, Validation e Test

```
In [42]: print_all_stats(boost_model)
```

Train

A: ('bias': 0.007573332607368333, 'var': 0.00126592933888122, 'mse': 0.008639261843265099)
B: ('bias': 0.00785173023240721, 'var': 0.000932285297467683, 'mse': 0.00847937826360486)
C: ('bias': 0.012297458940476645, 'var': 0.0014159475130206248, 'mse': 0.01371340745349729)

Validation

A: ('bias': 0.027867041028237897, 'var': 0.00122991623840724967, 'mse': 0.02916620341231031)
B: ('bias': 0.01978228058029664, 'var': 0.000790469655768118, 'mse': 0.020495751217928166)
C: ('bias': 0.02299356156102243, 'var': 0.00144579224097579516, 'mse': 0.024399148595860275)

Test

A: ('bias': 0.02903795884038705, 'var': 0.001262205334630516, 'mse': 0.030290164179669298)
B: ('bias': 0.024153649436275, 'var': 0.0007822476849189092, 'mse': 0.02490751217928166)
C: ('bias': 0.0382179423829629, 'var': 0.001389062739440926, 'mse': 0.03960706222404005)

Come atteso dalle precedenti analisi, non sembra che il boosting sia riuscito a diminuire l'errore; questo perché non è riuscito a diminuire come sperato in maniera significativa il bias, che sembra non riuscire ad abbattere anche con il boosting una certa soglia limite pari a quella del bias del modello banale.