

I risultati ottenuti sono analoghi per tutti e tre le regioni: 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 tenda 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 la baseline del modello banale.

```
In [27]: def print_stats(X, y, models):
    for i in range(3):
        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()
    print(f'bias': ((y - np.mean(y_pred))**2).mean(), \
          'var': (np.var(y_pred).mean(), \
          'mse': ((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():
    print_stats_trivial(X_train):
    for i in range(len(X_train)):
        print(region_names[i])
        print_stats_trivial(y_train[i].values.ravel(), y_test[i].values.ravel())
        print()
```

Statistiche per le tre regioni 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': 0.007536264902483718, 'var': 0.000345142012373853424, 'mse': 0.007881369615220221)
C: ('bias': 0.0116745081204912, 'var': 0.0005817505209895822, 'mse': 0.012256286641480767)

Validation
A: ('bias': 0.0276636307669438, 'var': 0.0005273470318553495, 'mse': 0.027793701107504842)
B: ('bias': 0.019614360581975956, 'var': 0.000580683812700375, 'mse': 0.020195044392324598)
C: ('bias': 0.02217266495119075, 'var': 0.0005495066299329632, 'mse': 0.022722171951123617)

Test
A: ('bias': 0.02864172421576377, 'var': 0.00030789391837899747, 'mse': 0.02896961813395547)
B: ('bias': 0.024158109782064934, 'var': 0.0002547132513325053, 'mse': 0.024412823033397434)
C: ('bias': 0.03776212365272005, 'var': 0.0004758397461138191, 'mse': 0.03823796339883999)

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

A
logerror
count 26819.000000
mean 0.012913
std 0.033157
std 0.086804
min -0.470000
5% -0.002000
50% 0.006000
75% 0.036300
max 0.756600

B
logerror
count 8119.000000
mean 0.013157
std 0.086804
min -0.470000
5% -0.002000
50% 0.006000
75% 0.036300
max 0.756600

C
logerror
count 25908.000000
mean 0.012105
std 0.030851
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': 0.028637396765634218, 'var': 1.20370821524020224e-35, 'mse': 0.02863739676563433)

B
('bias': 0.0241624624499493516, 'var': 3.009265538105056e-36, 'mse': 0.02416246449493482)

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
```

Scelta 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 proviene principalmente dal bias
- è capace di raggiungere un buono score anche partendo da modelli relativamente semplici, come quelli selezionati dall'albero di decisione

Uso come modelli base gli alberi individuati ai punti precedenti

```
In [34]: base_model = []
        #n_estimators = 500
        for dt in dt_models:
            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=''):
    def get_adaboost_regressor(Xs, ys, estimators):
        adaboost = AdaBoostRegressor(
            baseModel,
            n_estimators = estimators
        )
        adaboost.fit(Xs,ys)
        return adaboost

    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(f'({name}): TUNING DEL NUMERO DI STIMATORI')
        print()
        for n in ['mse', 'bias', 'var']:
            min_ = min(stats[n])
            best = (np.argmin(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'({name}): 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):
                if debug:
                    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_adaboost_regressor(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'({train MSE:      (trn_mse_s) - Val MSE:      (val_mse_s)}) * \n'
                    f'({trn_bias_s) - Val Bias:      (val_bias_s)}) * \n'
                    f'({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')
            min_ = min(train_stats, 'Validation')
            plot_mse(val_stats, 'Validation')

            return get_adaboost_regressor(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 Regione 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: 0.016299857197165746 (500) stimatori
Best mse: 0.008570762096412658
Best number of Estimators: 20

Punteggio finale: 0.00861629437220689 (500) stimatori
Best bias: 0.007319919202575296
Best number of Estimators: 20

Punteggio finale: 0.0076392207759945063 (500) stimatori
Best var: 0.001450849306838356
Best number of Estimators: 20

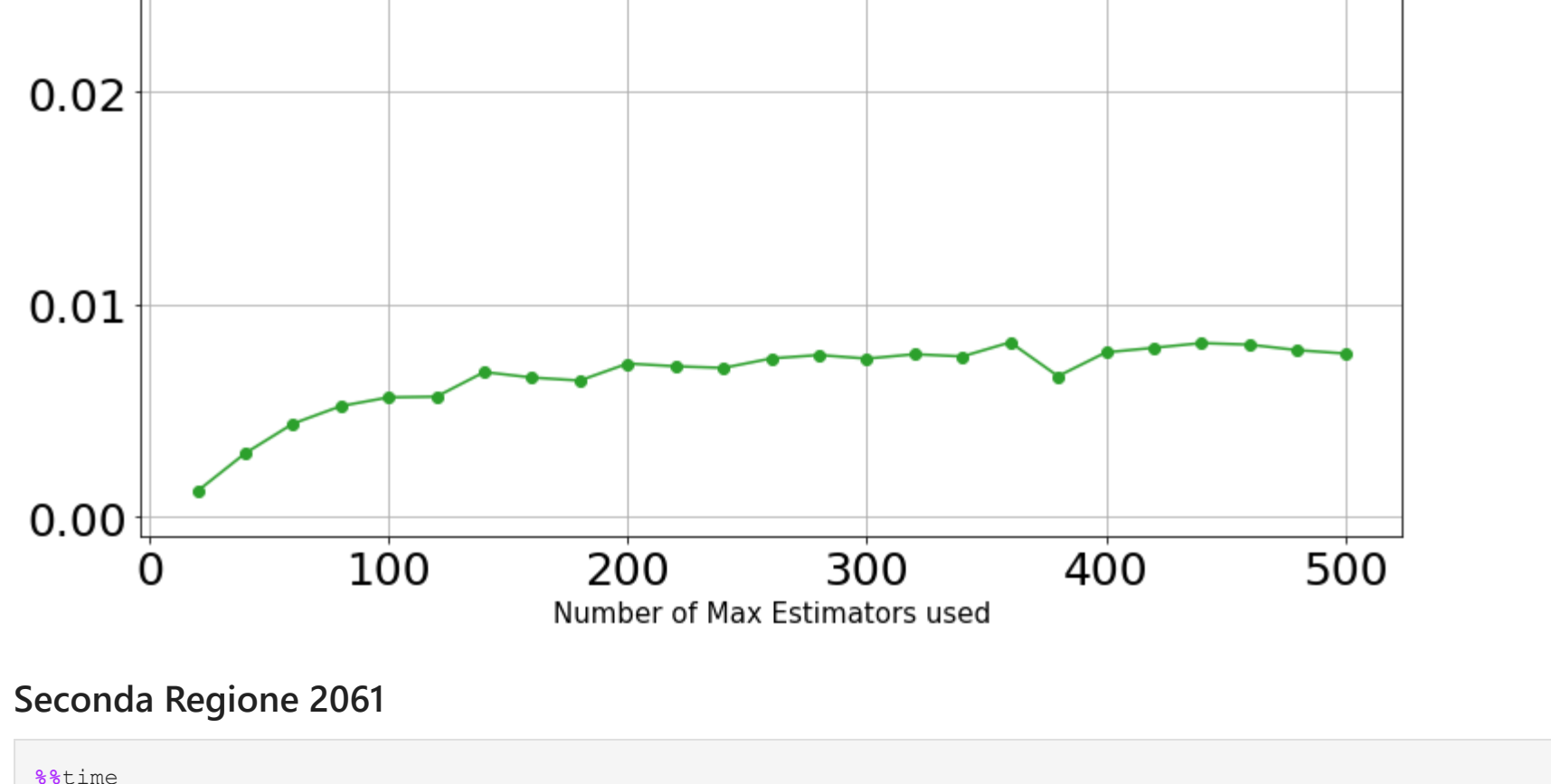
Validation: TUNING DEL NUMERO DI STIMATORI
Best mse: 0.03505686725975199
Best number of Estimators: 20

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

Punteggio finale: 0.007861629437220689 (500) stimatori
Best var: 0.001250562282222297
Best number of Estimators: 20

Wall time: 45min 32s

Train: MSE, BIAS, VARIANCE on different Estimators
```



Seconda Regione 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: 0.013847338941855728 (500) stimatori
Best mse: 0.00911286829265246
Best number of Estimators: 20

Punteggio finale: 0.009118869421700655 (500) stimatori
Best bias: 0.007913888051394728
Best number of Estimators: 20

Punteggio finale: 0.004726469520155089 (500) stimatori
Best var: 0.001198980241616206
Best number of Estimators: 20

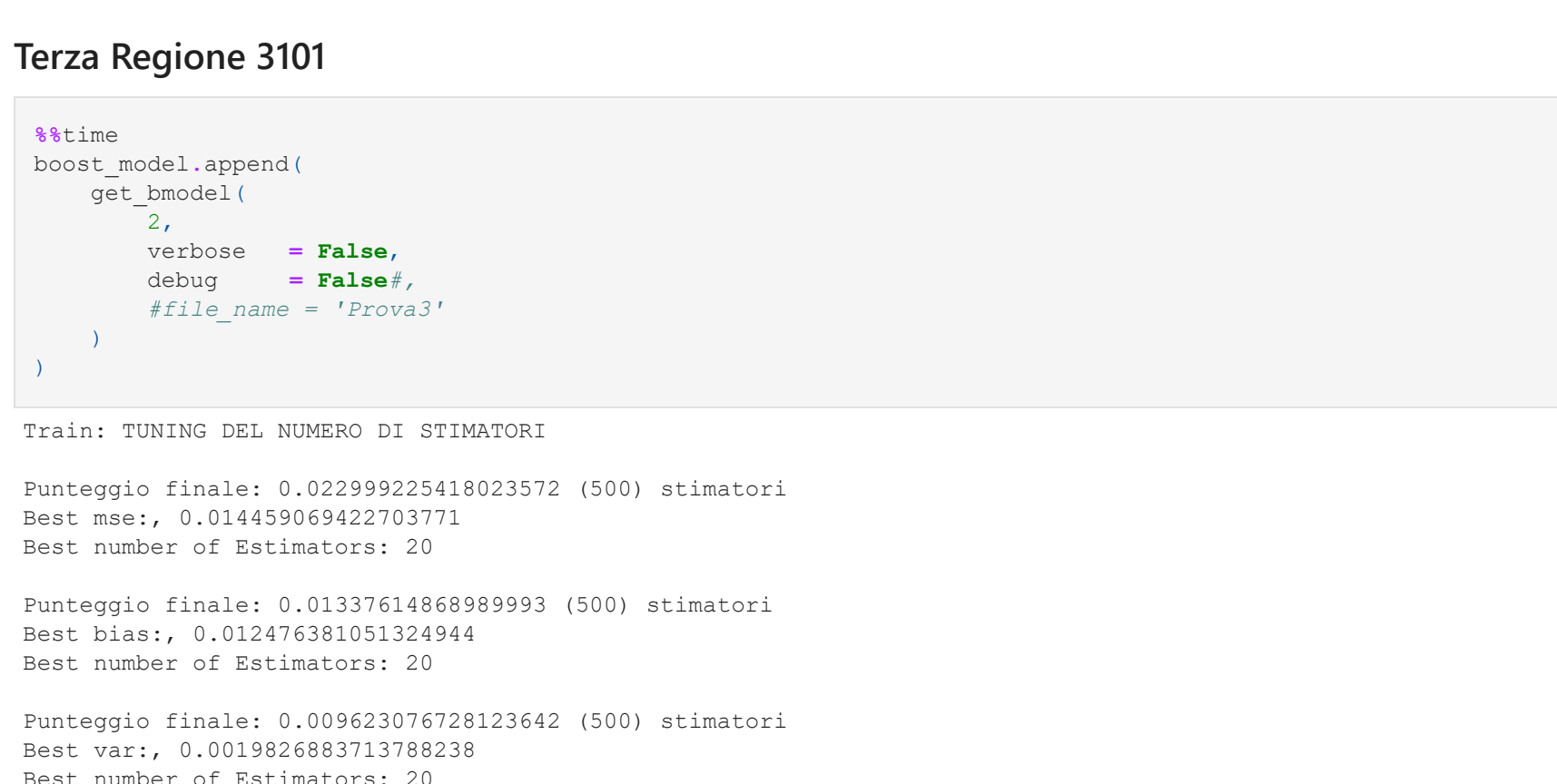
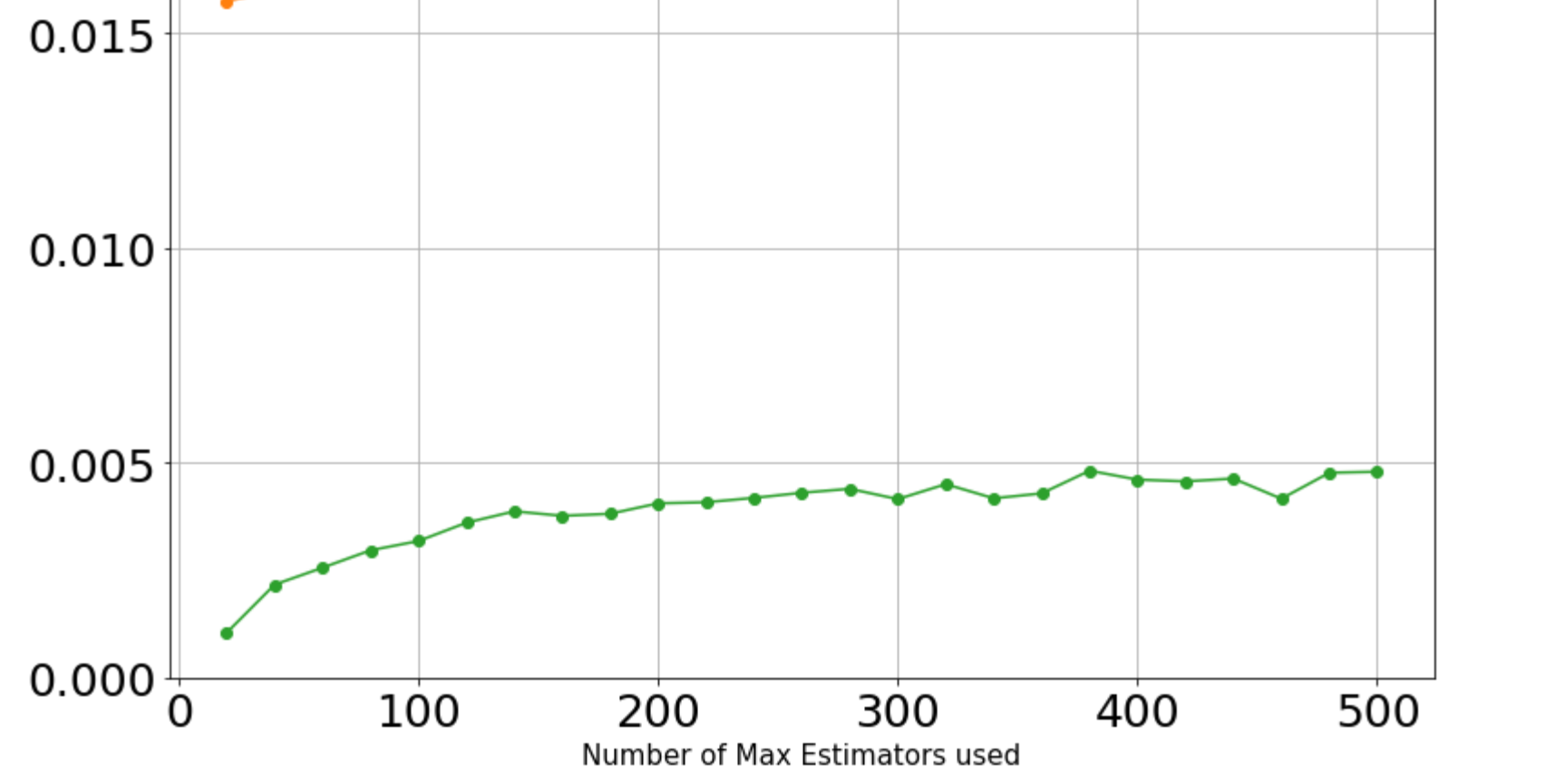
Validation: TUNING DEL NUMERO DI STIMATORI
Punteggio finale: 0.021694789718708297 (500) stimatori
Best mse: 0.0164744749110757
Best number of Estimators: 20

Punteggio finale: 0.016914292931888616 (500) stimatori
Best bias: 0.01573298230553015
Best number of Estimators: 20

Punteggio finale: 0.00478049676819669 (500) stimatori
Best var: 0.001031732291132495
Best number of Estimators: 20

Wall time: 32min 11s

Train: MSE, BIAS, VARIANCE on different Estimators
```



Terza Regione 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: 0.02399225418023572 (500) stimatori
Best mse: 0.01459084822703771
Best number of Estimators: 20

Punteggio finale: 0.01337614669899993 (500) stimatori
Best bias: 0.012476381051324944
Best number of Estimators: 20

Punteggio finale: 0.009623076728123642 (500) stimatori
Best var: 0.001982688713788238
Best number of Estimators: 20

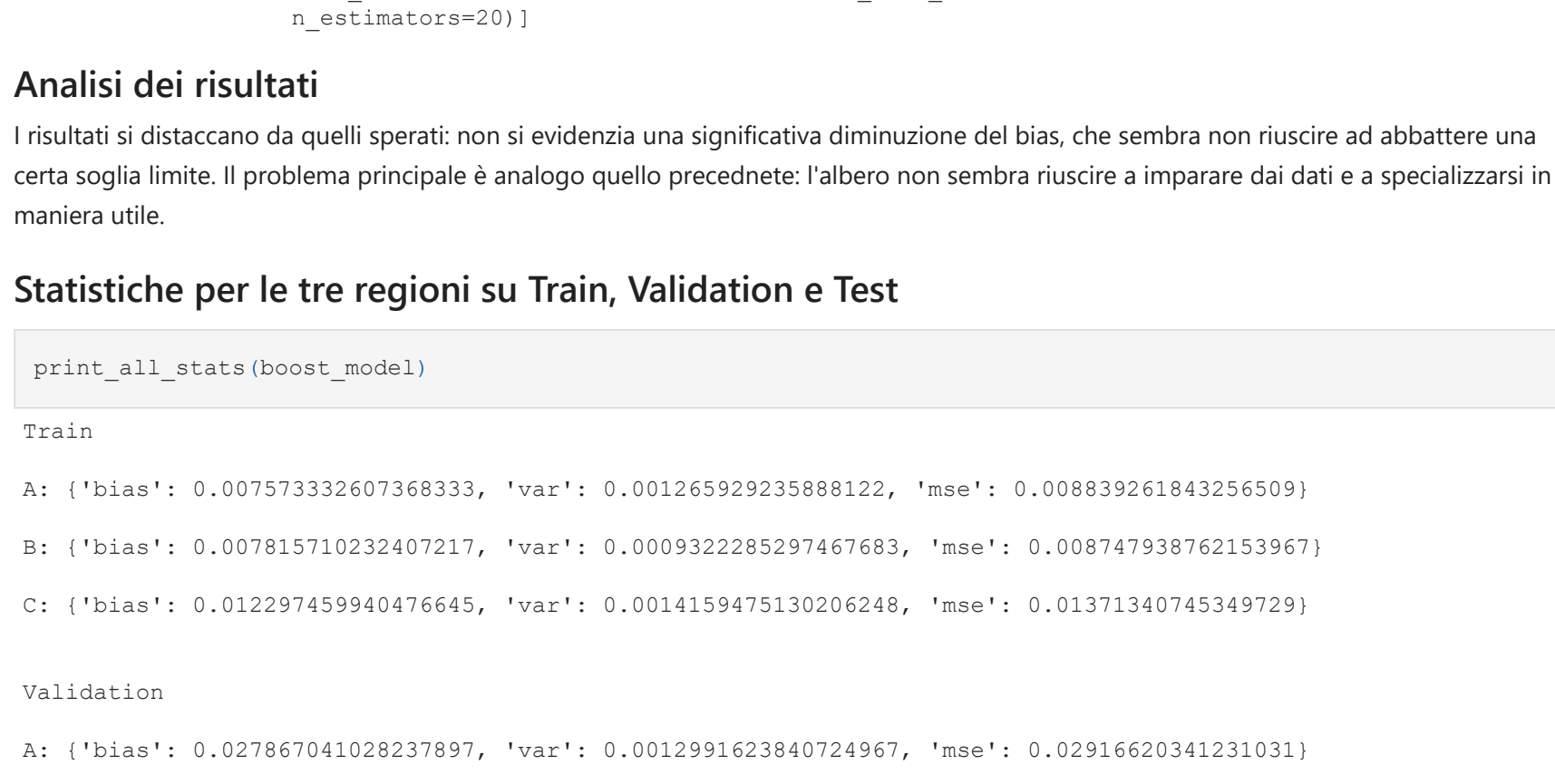
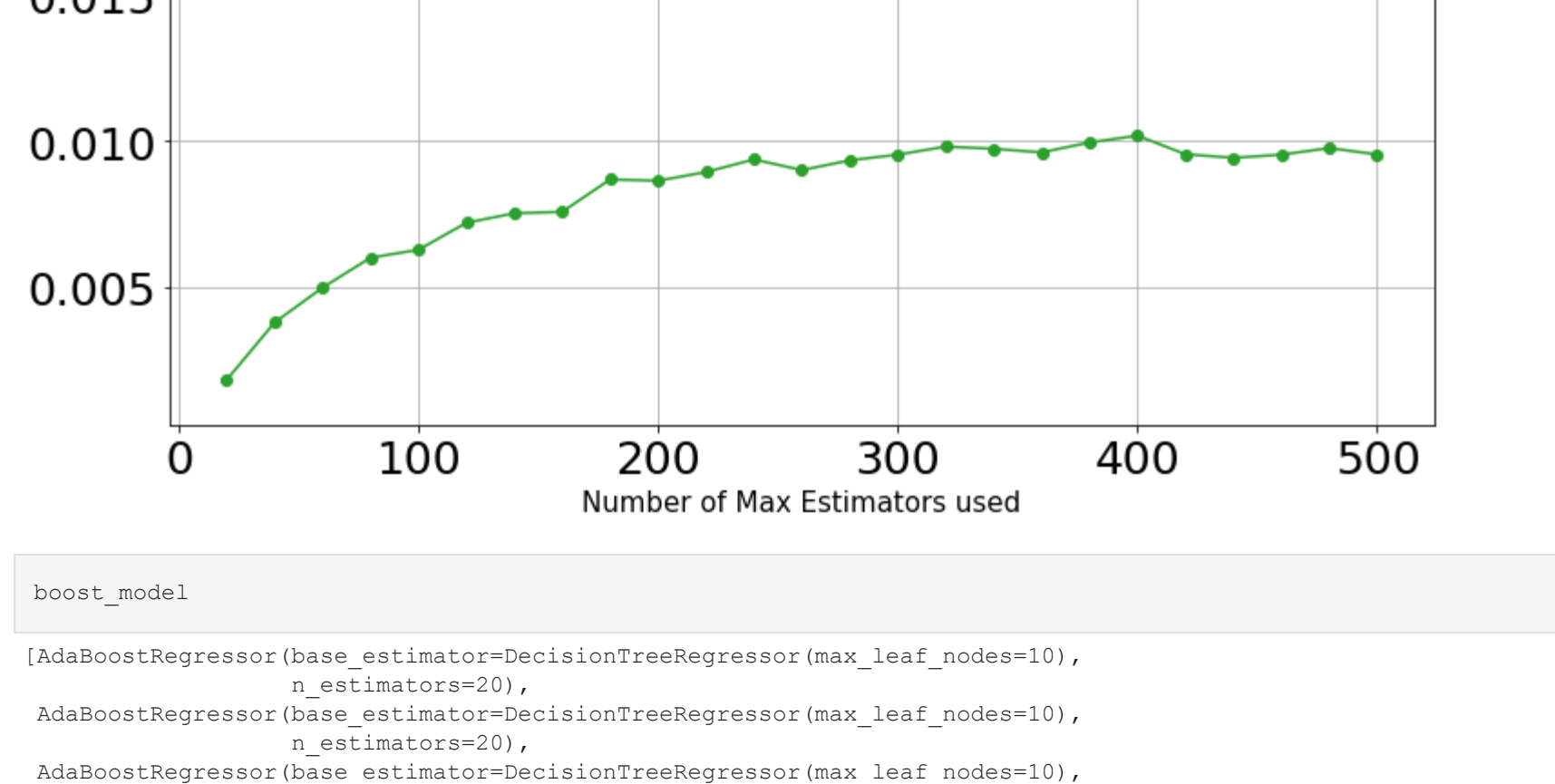
Validation: TUNING DEL NUMERO DI STIMATORI
Punteggio finale: 0.03147336445053434 (500) stimatori
Best mse: 0.022620905013911113
Best number of Estimators: 20

Punteggio finale: 0.02193619918809712 (500) stimatori
Best bias: 0.02082657073149757
Best number of Estimators: 20

Punteggio finale: 0.009531746531724628 (500) stimatori
Best var: 0.001794334276013548
Best number of Estimators: 20

Wall time: 36min 58s

Train: MSE, BIAS, VARIANCE on different Estimators
```



```
In [41]: boost_model

Out[41]: [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 distaccano da quelli sperati: 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 regioni su Train, Validation e Test

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

Train
A: ('bias': 0.00757333607368333, 'var': 0.001265929233888122, 'mse': 0.008639261843256509)
B: ('bias': 0.0078517302340721, 'var': 0.000932285297467683, 'mse': 0.008747938762306486)
C: ('bias': 0.012293745894076645, 'var': 0.0014159475130206248, 'mse': 0.01371340745349729)

Validation
A: ('bias': 0.027867041026237897, 'var': 0.0012991623840724967, 'mse': 0.02916620341231031)
B: ('bias': 0.019782689580296664, 'var': 0.0007790469655768118, 'mse': 0.02051327823606486)
C: ('bias': 0.02293935151010224, 'var': 0.0014057924039759516, 'mse': 0.0242399148559860275)

Test
A: ('bias': 0.029037958845038705, 'var': 0.001252205334630516, 'mse': 0.030290164179669298)
B: ('bias': 0.02414536449436275, 'var': 0.0007621476849189092, 'mse': 0.02490751217928166)
C: ('bias': 0.0382179434829629, 'var': 0.001389062739440926, 'mse': 0.03960706224044005)

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 predittore banale.
```