

Secondo Modello

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

È costruita una **foresta**.

È impiegato un algoritmo per fare **tuning dell'iperparametro $n_{estimators}$** : ossia il numero di alberi della foresta. Per fare tuning dei parametri si fa uso di un **dataset di validation**: si mantiene l'iperparametro che produce il minimo errore su validation.

```
In [3]: # 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.utils import resample

warnings.filterwarnings('ignore')
```

Letture dei dati

```
In [2]: # local file path
dir_name = 'selezione'
region_names = np.array(['A', 'B', 'C'])
region_ids = np.array(['1286', '2061', '3101'])

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}')
```

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

```
In [3]: # 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 [4]: def dimensionality(y=False):
    for 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(f'X_train: {X_train[i].shape}')
```

```
In [5]: dimensionality(y=True)

X_train: (26819, 55)
X_val: (9006, 55)
X_test: (9085, 55)
y_train: (26819, 1)
y_val: (9006, 1)
y_test: (9085, 1)

X_train: (8119, 32)
X_val: (2658, 32)
X_test: (2606, 32)
y_train: (8119, 1)
y_val: (2658, 1)
y_test: (2606, 1)

X_train: (46771, 33)
X_val: (21908, 33)
X_test: (21876, 33)
y_train: (46771, 1)
y_val: (21908, 1)
y_test: (21876, 1)
```

Variabili Globali

```
In [6]: # Nomi delle regioni
region_ids = np.array(['1286', '2061', '3101'])

# Riduzione del dataset C
C_IND = 2
C_IND = 2

# Percentuale per il sottoinsieme su cui costruire il modello
SUB_PERC = [1/1000, 1/100, 1/1000] # A: 100, B: 100, C: 100; Test rapido ~ 5 min
SUB_PERC = [1/3, 1, 1/3] # A: 8000, B: 8000, C: 8000; Test medio ~ 2 ore
SUB_PERC = [1, 1, 1] # A: 24000, B: 8000, C: 64000; Dataset Completo ~ 5 ore

# Random Forest Regressor
RF_START = 0
RF_END = 1000
RF_STEP = 50
```

Limite del numero di istanze del terzo dataset

Per limiti tecnici di memoria la gestione del dataset C risulta complicata in quanto incorre spesso in errori a run-time di categoria **MemoryError** la macchina su cui è eseguito il notebook non è in grado di allocare la quantità di memoria richiesta per lavorare sull'intero dataset di Train (circa 64000 righe per 33 colonne).

È dunque inevitabile dover ridurre l'insieme con un numero di righe il cui calcolo è supportato, cioè circa il 40% del dataset. Il dataset per la terza regione (**region_id**: 3101) è riscaltato su 2/5 delle osservazioni per Train, Validation e Test.

L'operazione introduce dunque una forte approssimazione per questa regione, che per limiti tecnici è però inevitabile.

```
In [7]: for X, y in zip((X_train, X_val, X_test), (y_train, y_val, y_test)):
    X[C_IND], y[C_IND] = resample(X[C_IND], y[C_IND], n_samples = int(C_PERC * len(X[2]))
    )
```

```
In [8]: dimensionality(y=True)

X_train: (26819, 55)
X_val: (9006, 55)
X_test: (9085, 55)
y_train: (26819, 1)
y_val: (9006, 1)
y_test: (9085, 1)

X_train: (8119, 32)
X_val: (2658, 32)
X_test: (2606, 32)
y_train: (8119, 1)
y_val: (2658, 1)
y_test: (2606, 1)

X_train: (46771, 33)
X_val: (21908, 33)
X_test: (21876, 33)
y_train: (46771, 1)
y_val: (21908, 1)
y_test: (21876, 1)
```

Definizione di un Subset per la costruzione del modello

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 la totalità delle istanze del dataset originale.

```
In [9]: for i in range(3):
    print(f'Region {region_names[i]}')
    print(f'len(X_train_sub[i]) * SUB_PERC[i]')
    print(f'len(X_val_sub[i]) * SUB_PERC[i]')
    print(f'len(X_test_sub[i]) * SUB_PERC[i]')
```

```
Region A
26819
9006

Region B
8119
2658

Region C
25908
8763
```

```
In [10]: X_train_sub = []
X_val_sub = []
X_test_sub = []
y_train_sub = []
y_val_sub = []
y_test_sub = []
```

```
In [11]: for i in range(3):
    X_train_sub.append(X_train[i], y_train[i], n_samples = int(SUB_PERC[i] * len(X_train[i])))
    X_val_sub.append(X_val[i], y_val[i], n_samples = int(SUB_PERC[i] * len(X_val[i])))
    X_test_sub.append(X_test[i], y_test[i], n_samples = int(SUB_PERC[i] * len(X_test[i])))
    y_train_sub.append(y_train[i])
    y_val_sub.append(y_val[i])
    y_test_sub.append(y_test[i])
```

```
In [12]: def dimensionality_sub(y=False):
    for i in range(3):
        print(f'X_train_sub(region_names[i]): {X_train_sub[i].shape}')
        print(f'X_val_sub(region_names[i]): {X_val_sub[i].shape}')
        if y:
            print(f'y_train_sub(region_names[i]): {y_train_sub[i].shape}')
            print(f'y_val_sub(region_names[i]): {y_val_sub[i].shape}')
            print(f'y_test_sub(region_names[i]): {y_test_sub[i].shape}')
```

```
In [13]: dimensionality_sub(y=True)

X_train_sub: (26819, 55)
X_val_sub: (9006, 55)
X_test_sub: (9085, 55)
y_train_sub: (26819, 1)
y_val_sub: (9006, 1)
y_test_sub: (9085, 1)

X_train_sub: (8119, 32)
X_val_sub: (2658, 32)
X_test_sub: (2606, 32)
y_train_sub: (8119, 1)
y_val_sub: (2658, 1)
y_test_sub: (2606, 1)

X_train_sub: (25908, 33)
X_val_sub: (8763, 33)
X_test_sub: (8750, 33)
y_train_sub: (25908, 1)
y_val_sub: (8763, 1)
y_test_sub: (8750, 1)
```

Costruzione Random Forest Regressor

Tuning numero di alberi attraverso il parametro **$n_{estimators}$** . Per studiare l'andamento la funzione produce due grafici che illustrano l'andamento di **bias**, **varianza** e **mean squared error** in funzione del numero di alberi sia per l'insieme di **Train** che per quello di **Validation**.

```
In [14]: plt.rcParams.update({'font.size': 35})
```

```
In [15]: 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 - y_pred)**2).mean()
```

```
In [16]: # Costruzione RandomForestRegressor
def RandomForestRegressor_validation(X_train, y_train, X_val, y_val, verbose=False, debug=False, file_name = ''):
    dt = RandomForestRegressor(
        n_estimators=1,
        criterion='squared_error',
        n_jobs=-1
    )
    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(f'({name}): TUNING DEL MASSIMO NUMERO DI ALBERI NELLA FORESTA')
    print()
    for n in ['mse', 'bias', 'var']:
        min = min(stats[n])
        best = np.argmax(stats[n]) * RF_STEP + RF_START

        print(f'Punteggio finale: {stats[n][best]} ((RF_END) alberi)')
        print(f'Best (n): {min}')
        print(f'Best number of Trees used: {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(RF_START, RF_END+1, RF_STEP), stats['mse'], 'o-', label='MSE')
        ax.plot(range(RF_START, RF_END+1, RF_STEP), stats['bias'], 'o-', label='BIAS')
        ax.plot(range(RF_START, RF_END+1, RF_STEP), stats['var'], 'o-', label='VARIANCE')

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

    if file_name != '':
        fig.savefig('images/' + file_name + '_RandomForestRegressor' + 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 estimator in range(RF_START, RF_END+1, RF_STEP):
        if debug:
            print(f'({estimator})/(RF_END)')

        model = get_rf_reg(estimator)

        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)

        train_stats['bias'].append(trn_bias)
        train_stats['var'].append(trn_var)
        train_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(f'Estimators: {estimator}') + \
            f'\nTrain MSE: (trn_mse) - Val MSE: (val_mse)' + \
            f'\nTrain Bias: (trn_bias) - Val Bias: (val_bias)' + \
            f'\nTrain Variance: (trn_var) - Val Variance: (val_var)'

    if (first or val_mse < best_mse):
        first = False
        best_mse = val_mse
        best_estimator = estimator
        best_model = model

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

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

    return best_model
```

```
In [17]: def get_rf(index, verbose=False, debug=False, file_name=''):
    if file_name == '':
        file_name = region_ids[index]
    return RandomForestRegressor_validation(
        X_train_sub[index],
        y_train_sub[index],
        X_val_sub[index],
        y_val_sub[index],
        verbose = verbose,
        debug = debug,
        file_name = file_name
    )
```

In [18]: rf_model = 1

Prima Contea 1286

```
In [19]: rf_model.append(
    get_rf(
        1,
        verbose = False,
        debug = False,
        file_name = 'Prova1'
    )
)
```

Train: TUNING DEL MASSIMO NUMERO DI ALBERI NELLA FORESTA

Punteggio finale: 0.01273302418721063 (1000) alberi
Best mse: 0.01270233531289177
Best number of Trees: 200

Punteggio finale: 0.006751486653928585 (1000) alberi
Best bias: 0.006751477646661825
Best number of Trees: 200

Punteggio finale: 0.004528183764792367 (1000) alberi
Best mse: 0.004518733251862547
Best number of Trees: 750

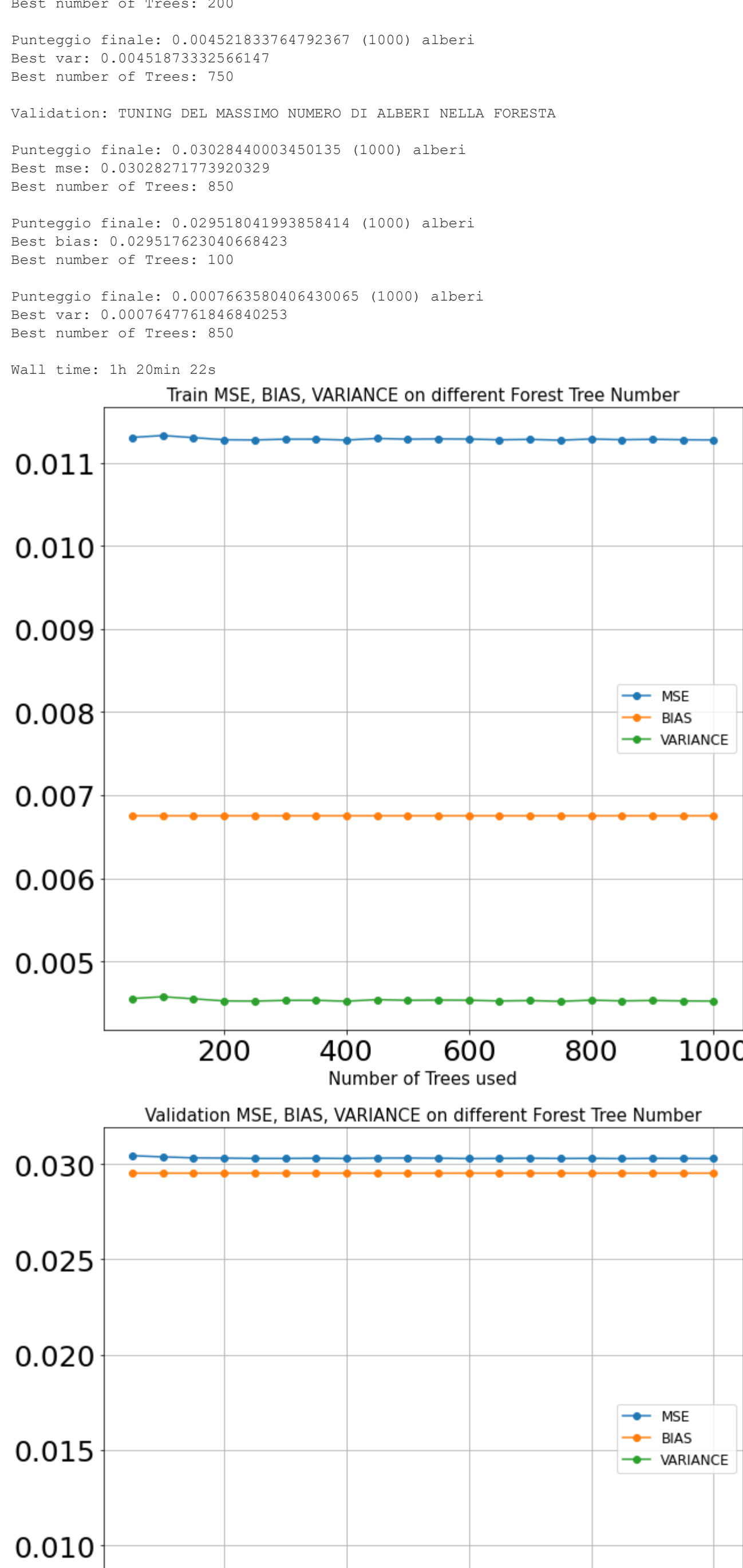
Validation: TUNING DEL MASSIMO NUMERO DI ALBERI NELLA FORESTA

Punteggio finale: 0.03028440003650135 (1000) alberi
Best mse: 0.0302827173920329
Best number of Trees: 850

Punteggio finale: 0.029517623040668423 (1000) alberi
Best bias: 0.029517623040668423
Best number of Trees: 200

Punteggio finale: 0.0075166359804633065 (1000) alberi
Best var: 0.00764761846840253
Best number of Trees: 850

Wall time: 1h 20min 22s



Seconda Contea 2061

```
In [20]: rf_model.append(
    get_rf(
        2,
        verbose = False,
        debug = False,
        file_name = 'Prova2'
    )
)
```

Train: TUNING DEL MASSIMO NUMERO DI ALBERI NELLA FORESTA

Punteggio finale: 0.012578312287415128 (1000) alberi
Best mse: 0.01861892215186235
Best number of Trees: 200

Punteggio finale: 0.0075187491384942226 (1000) alberi
Best bias: 0.0075187491384942226
Best number of Trees: 200

Punteggio finale: 0.005059528146636468 (1000) alberi
Best var: 0.00504244193897201
Best number of Trees: 200

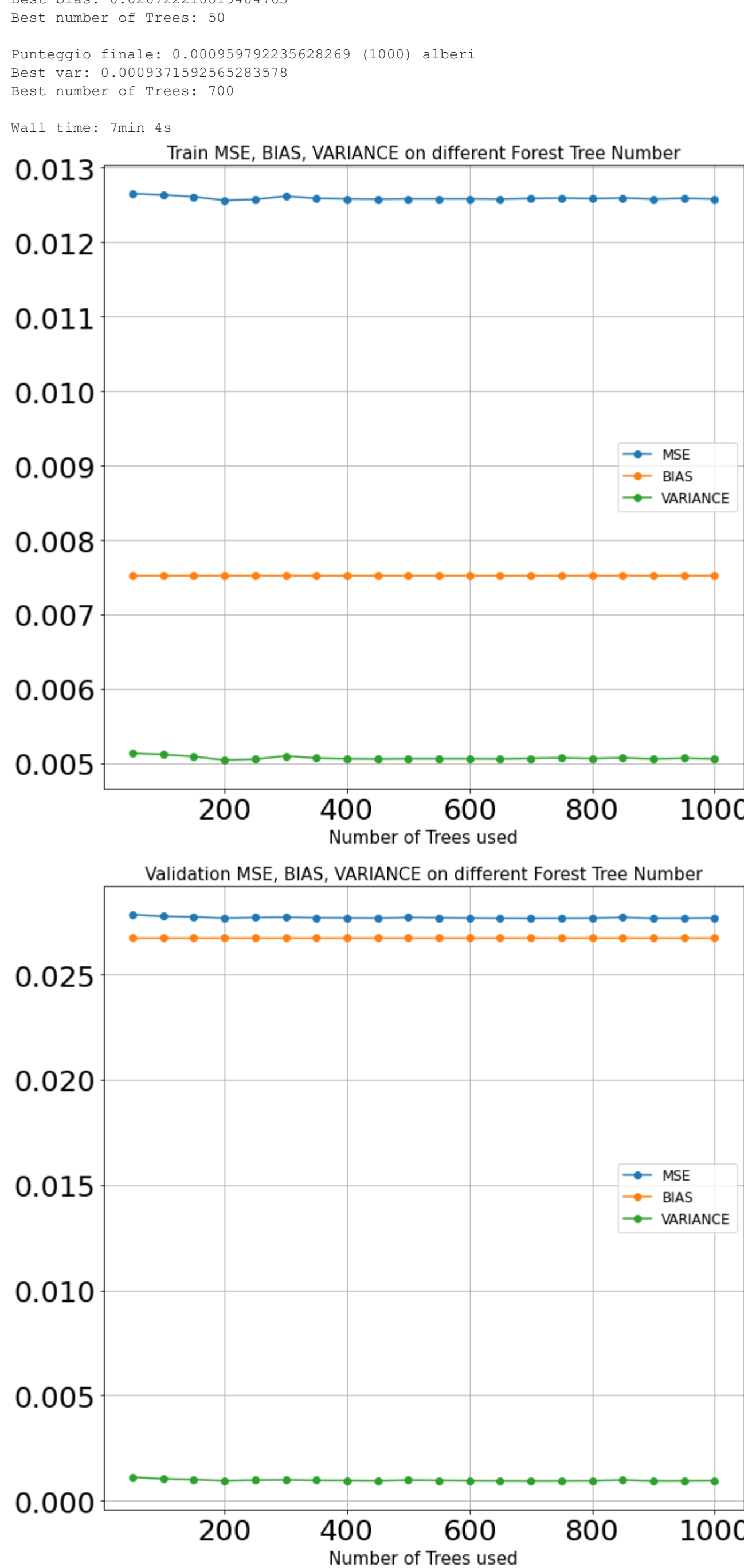
Validation: TUNING DEL MASSIMO NUMERO DI ALBERI NELLA FORESTA

Punteggio finale: 0.027682165249020445 (1000) alberi
Best mse: 0.027659561486424258
Best number of Trees: 700

Punteggio finale: 0.026727373013392222 (1000) alberi
Best bias: 0.026722210919404743
Best number of Trees: 50

Punteggio finale: 0.000957929235628269 (1000) alberi
Best var: 0.000937159256283578
Best number of Trees: 700

Wall time: 7min 4s



Seconda Contea 2061

```
In [21]: rf_model.append(
    get_rf(
        2,
        verbose = False,
        debug = False,
        file_name = 'Prova3'
    )
)
```

Train: TUNING DEL MASSIMO NUMERO DI ALBERI NELLA FORESTA

Punteggio finale: 0.018628124583239 (1000) alberi
Best mse: 0.01861892215186235
Best number of Trees: 200

Punteggio finale: 0.010770387924195961 (1000) alberi
Best bias: 0.010770387924195961
Best number of Trees: 50

Punteggio finale: 0.007858193722906867 (1000) alberi
Best var: 0.00784827874187144
Best number of Trees: 200

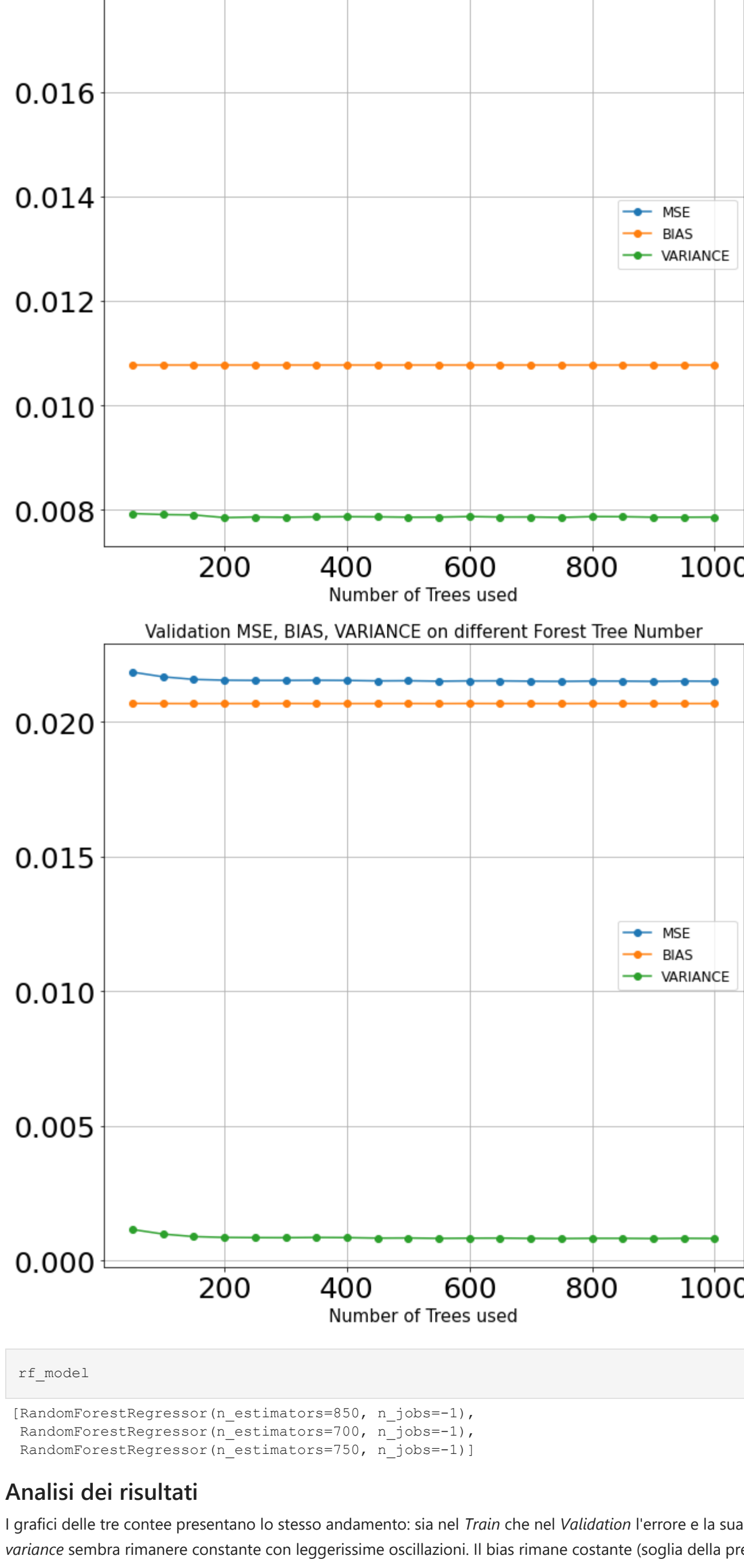
Validation: TUNING DEL MASSIMO NUMERO DI ALBERI NELLA FORESTA

Punteggio finale: 0.02151663792432453 (1000) alberi
Best mse: 0.02151663792432453
Best number of Trees: 750

Punteggio finale: 0.020868515472398165 (1000) alberi
Best bias: 0.020868515472398165
Best number of Trees: 200

Punteggio finale: 0.000802903899264212 (1000) alberi
Best var: 0.000802903899264212
Best number of Trees: 800

Wall time: 1h 1min 13s



Seconda Contea 2061

```
In [22]: rf_model.append(
    get_rf(
        2,
        verbose = False,
        debug = False,
        file_name = 'Prova4'
    )
)
```

Train: TUNING DEL MASSIMO NUMERO DI ALBERI NELLA FORESTA

Punteggio finale: 0.018628124583239 (1000) alberi
Best mse: 0.01861892215186235
Best number of Trees: 200

Punteggio finale: 0.010770387924195961 (1000) alberi
Best bias: 0.010770387924195961
Best number of Trees: 50

Punteggio finale: 0.007858193722906867 (1000) alberi
Best var: 0.00784827874187144
Best number of Trees: 200

Validation: TUNING DEL MASSIMO NUMERO DI ALBERI NELLA FORESTA

Punteggio finale: 0.02151663792432453 (1000) alberi
Best mse: 0.02151663792432453
Best number of Trees: 750

Punteggio finale: 0.020868515472398165 (1000) alberi
Best bias: 0.020868515472398165
Best number of Trees: 200

Punteggio finale: 0.000802903899264212 (1000) alberi
Best var: 0.000802903899264212
Best number of Trees: 800

Wall time: 1h 1min 13s



Seconda Contea 2061

```
In [23]: rf_model.append(
    get_rf(
        2,
        verbose = False,
        debug = False,
        file_name = 'Prova5'
    )
)
```

Train: TUNING DEL MASSIMO NUMERO DI ALBERI NELLA FORESTA

Punteggio finale: 0.018628124583239 (1000) alberi
Best mse: 0.01861892215186235
Best number of Trees: 200

Punteggio finale: 0.010770387924195961 (1000) alberi
Best bias: 0.010770387924195961
Best number of Trees: 50

Punteggio finale: 0.007858193722906867 (1000) alberi
Best var: 0.00784827874187144
Best number of Trees: 200

Validation: TUNING DEL MASSIMO NUMERO DI ALBERI NELLA FORESTA

Punteggio finale: 0.02151663792432453 (1000) alberi
Best mse: 0.02151663792432453
Best number of Trees: 750

Punteggio finale: 0.020868515472398165 (1000) alberi
Best bias: 0.020868515472398165
Best number of Trees: 200

Punteggio finale: 0.000802903899264212 (1000) alberi
Best var: 0.000802903899264212
Best number of Trees: 800

Wall time: 1h 1min 13s

Seconda Contea 2061

```
In [24]: rf_model.append(
    get_rf(
        2,
        verbose = False,
        debug = False,
        file_name = 'Prova6'
    )
)
```

Train: TUNING DEL MASSIMO NUMERO DI ALBERI NELLA FORESTA

Punteggio finale: 0.018628124583239 (1000) alberi
Best mse: 0.01861892215186235
Best number of Trees: 200

Punteggio finale: 0.010770387924195961 (1000) alberi
Best bias: 0.010770387924195961
Best number of Trees: 50

Punteggio finale: 0.007858193722906867 (1000) alberi
Best var: 0.00784827874187144
Best number of Trees: 200

Validation: TUNING DEL MASSIMO NUMERO DI ALBERI NELLA FORESTA

Punteggio finale: 0.02151663792432453 (1000) alberi
Best mse: 0.02151663792432453
Best number of Trees: 750

Punteggio finale: 0.020868515472398165 (1000) alberi
Best bias: 0.020868515472398165
Best number of Trees: 200

Punteggio finale: 0.000802903899264212 (1000) alberi
Best var: 0.000802903899264212
Best number of Trees: 800

Wall time: 1h 1min 13s

Seconda Contea 2061

```
In [25]: rf_model.append(
    get_rf(
        2,
        verbose = False,
        debug = False,
        file_name = 'Prova7'
    )
)
```

Train: TUNING DEL MASSIMO NUMERO DI ALBERI NELLA FORESTA

Punteggio finale: 0.018628124583239 (1000) alberi
Best mse: 0.01861892215186235
Best number of Trees: 200

Punteggio finale: 0.010770387924195961 (1000) alberi
Best bias: 0.010770387924195961
Best number of Trees: 50

Punteggio finale: 0.007858193722906867 (1000) alberi
Best var: 0.00784827874187144
Best number of Trees: 200

Validation: TUNING DEL MASSIMO NUMERO DI ALBERI NELLA FORESTA

Punteggio finale: 0.02151663792432453 (1000) alberi
Best mse: 0.02151663792432453
Best number of Trees: 750

Punteggio finale: 0.020868515472398165 (1000) alberi
Best bias: 0.020868515472398165
Best number of Trees: 200

Punteggio finale: 0.000802903899264212 (1000) alberi
Best var: 0.000802903899264212
Best number of Trees: 800

Wall time: 1h 1min 13s

Seconda Contea 2061

```
In [26]: rf_model.append(
    get_rf(
        2,
        verbose = False,
        debug = False,
        file_name = 'Prova8'
    )
)
```

Train: TUNING DEL MASSIMO NUMERO DI ALBERI NELLA FORESTA

Punteggio finale: 0.018628124583239 (1000) alberi
Best mse: 0.0186

