



UNIVERSITÀ DEGLI
STUDI DI FIRENZE

Previsione: Random Forest vs BART

Statistical Learning – Contest

Arianna Russo e Simone Capucci

Laurea Magistrale in Data Science, Calcolo Scientifico e Intelligenza Artificiale

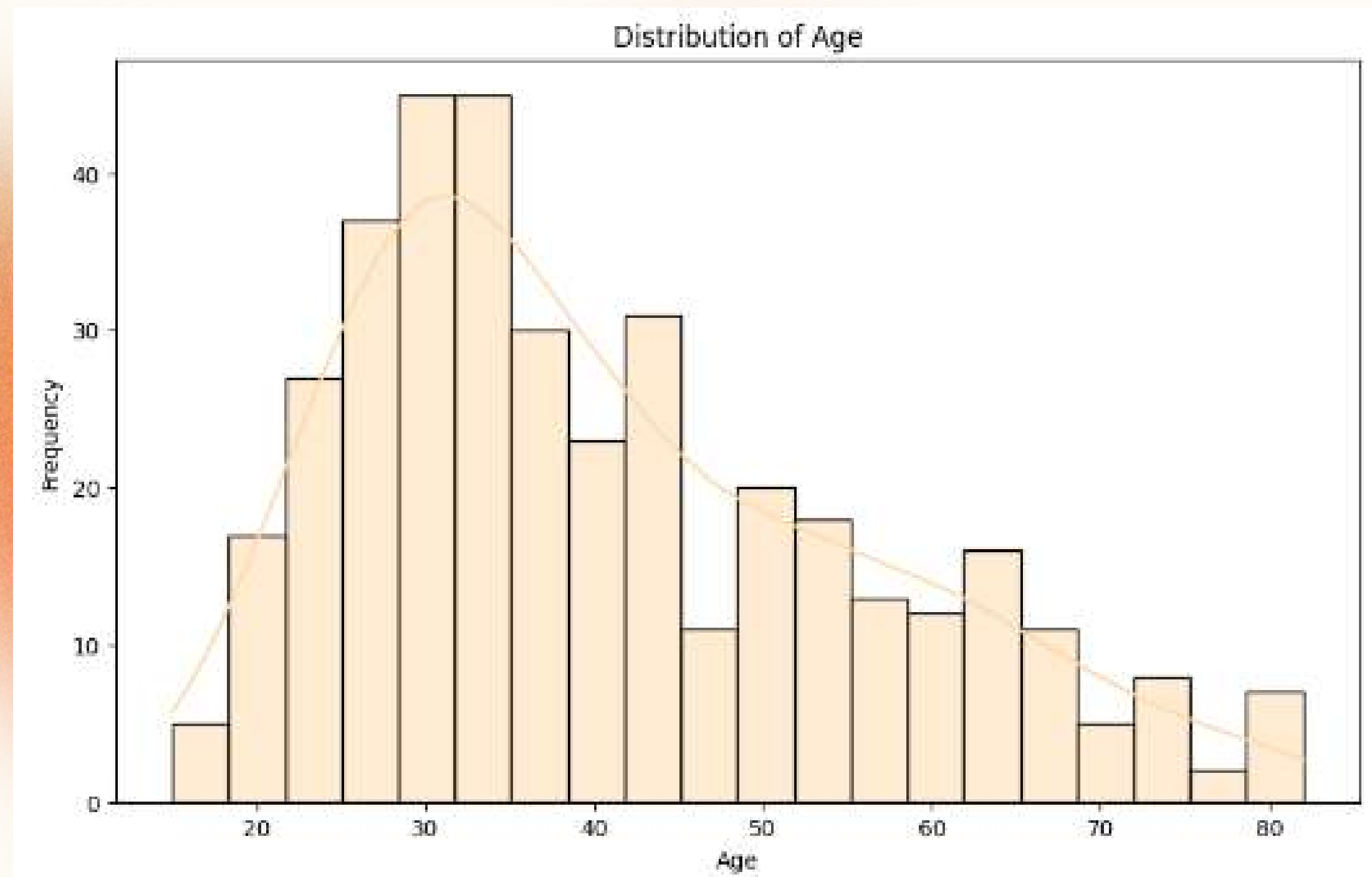


Recidiva del cancro alla tiroide

Obiettivo: Prevedere la recidiva del cancro alla tiroide

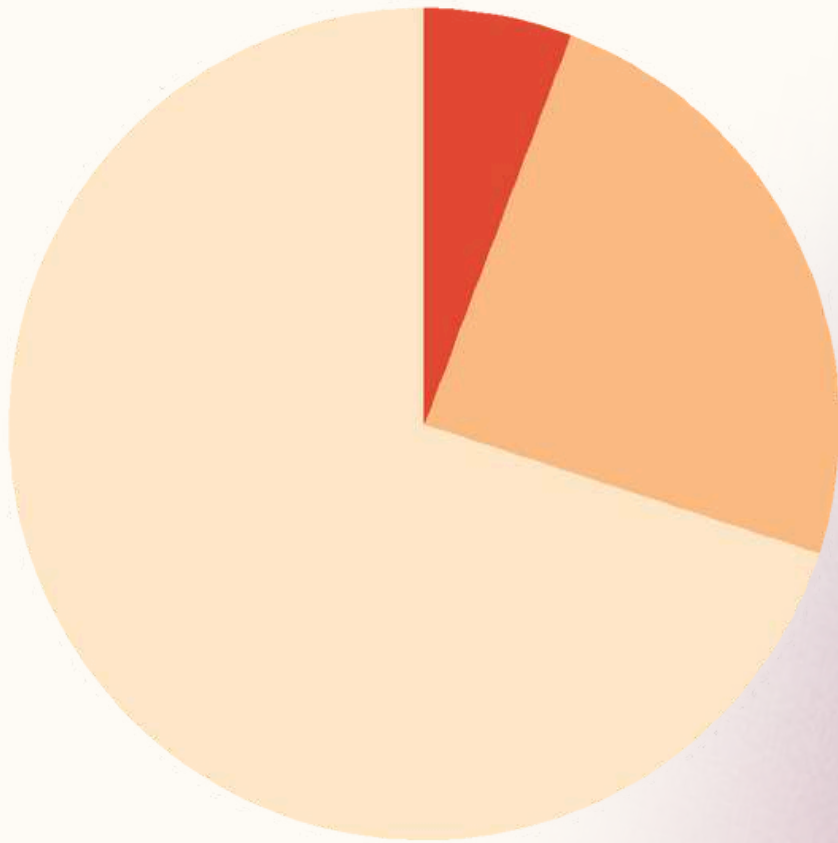
Il dataset:

Il dataset si compone di **383 osservazioni su 18 variabili** (collegate o possibilmente collegate al tema trattato), raccolte da soggetti differenti. Le variabili sono per lo più **categoriche** (a parte la variabile Age).





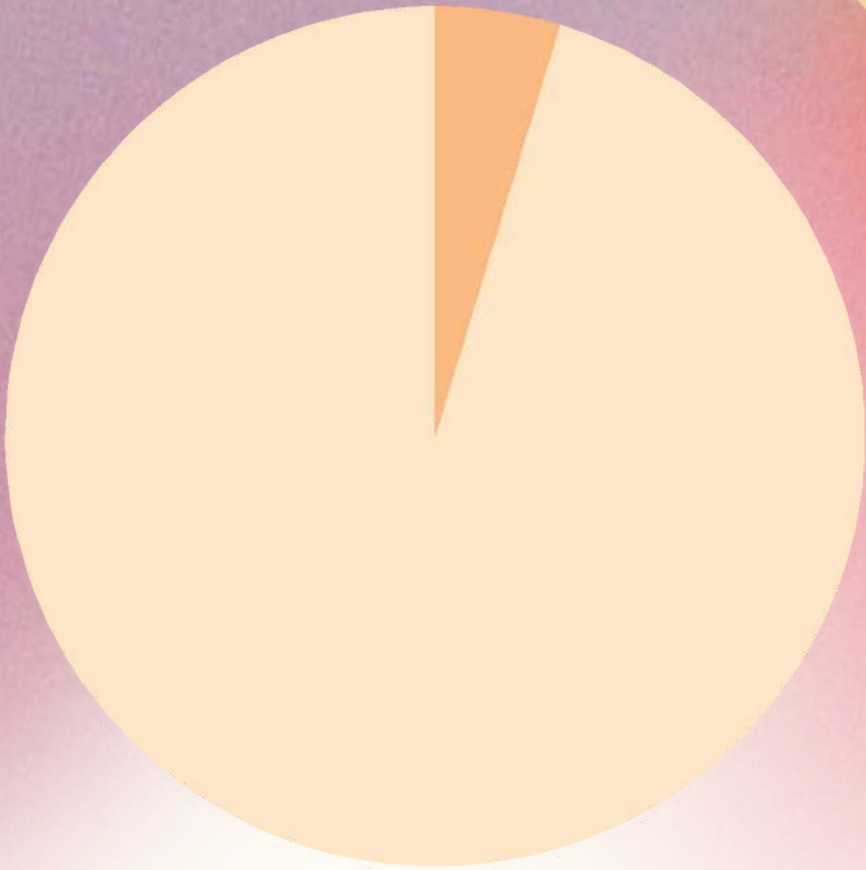
Lymph Nodes



Lymph_Nodes

- no evidence of regional lymph node metastasis
- regional lymph node metastasis in the central of the neck
- regional lymph node metastasis in the lateral of the neck

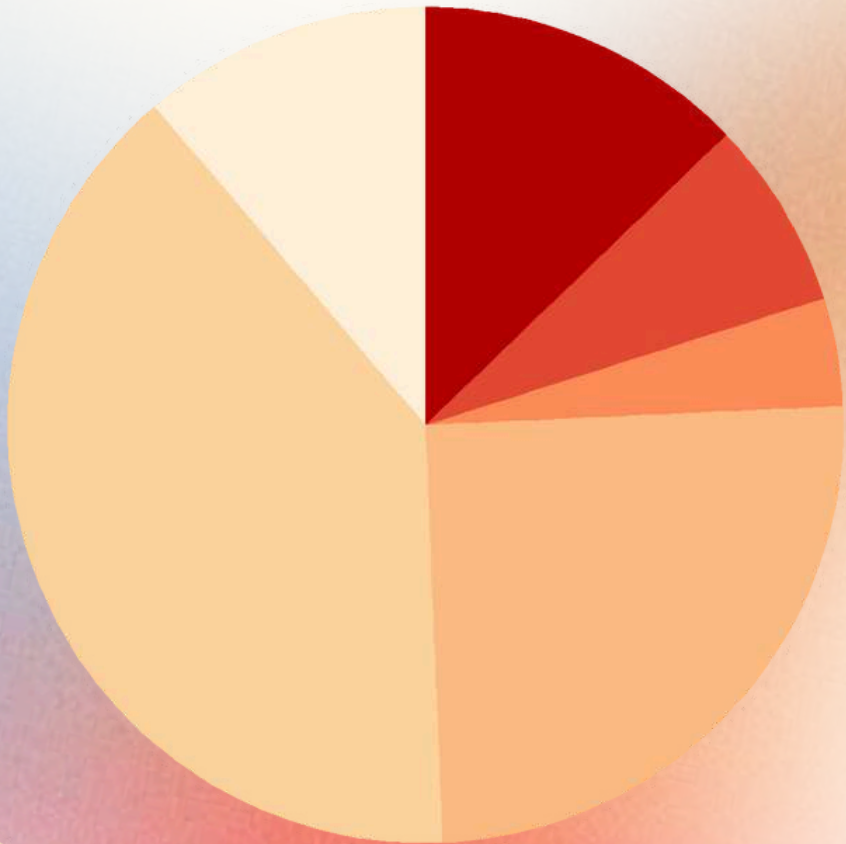
Cancer Metastasis



Cancer_Metastasis

- no evidence of distant metastasis
- the presence of distant metastasis

Tumor



Tumor

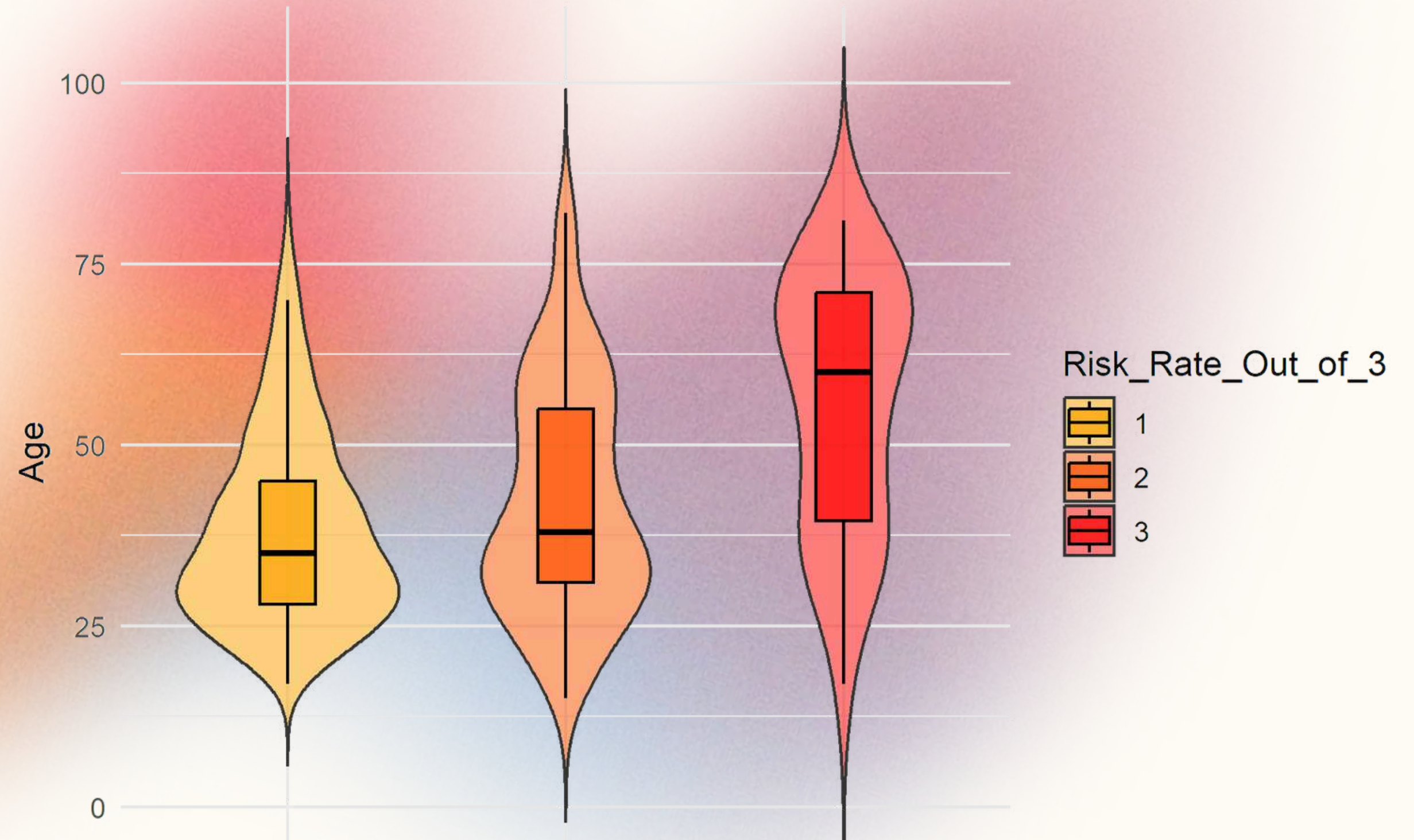
- tumor larger than 1 cm but not larger than 2 cm
- tumor larger than 2 cm but not larger than 4 cm
- tumor larger than 4 cm
- tumor that has grown outside the thyroid
- tumor that has invaded nearby structures
- tumor that is 1 cm or smaller



Data Visualization

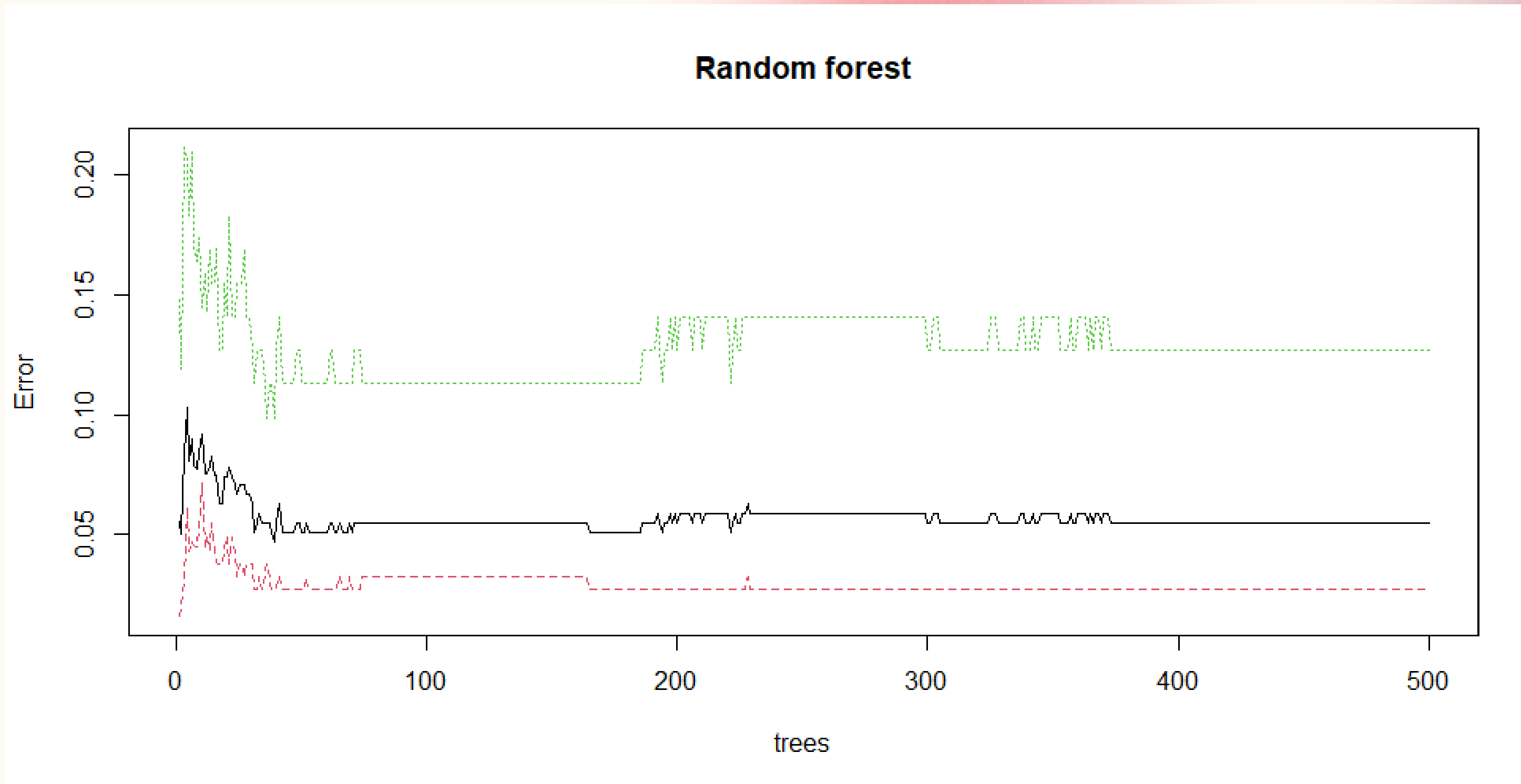
All'aumentare
dell'età **aumenta**
anche il rischio
per il paziente
associato al tumore

Distribution of Ages by Risk Category





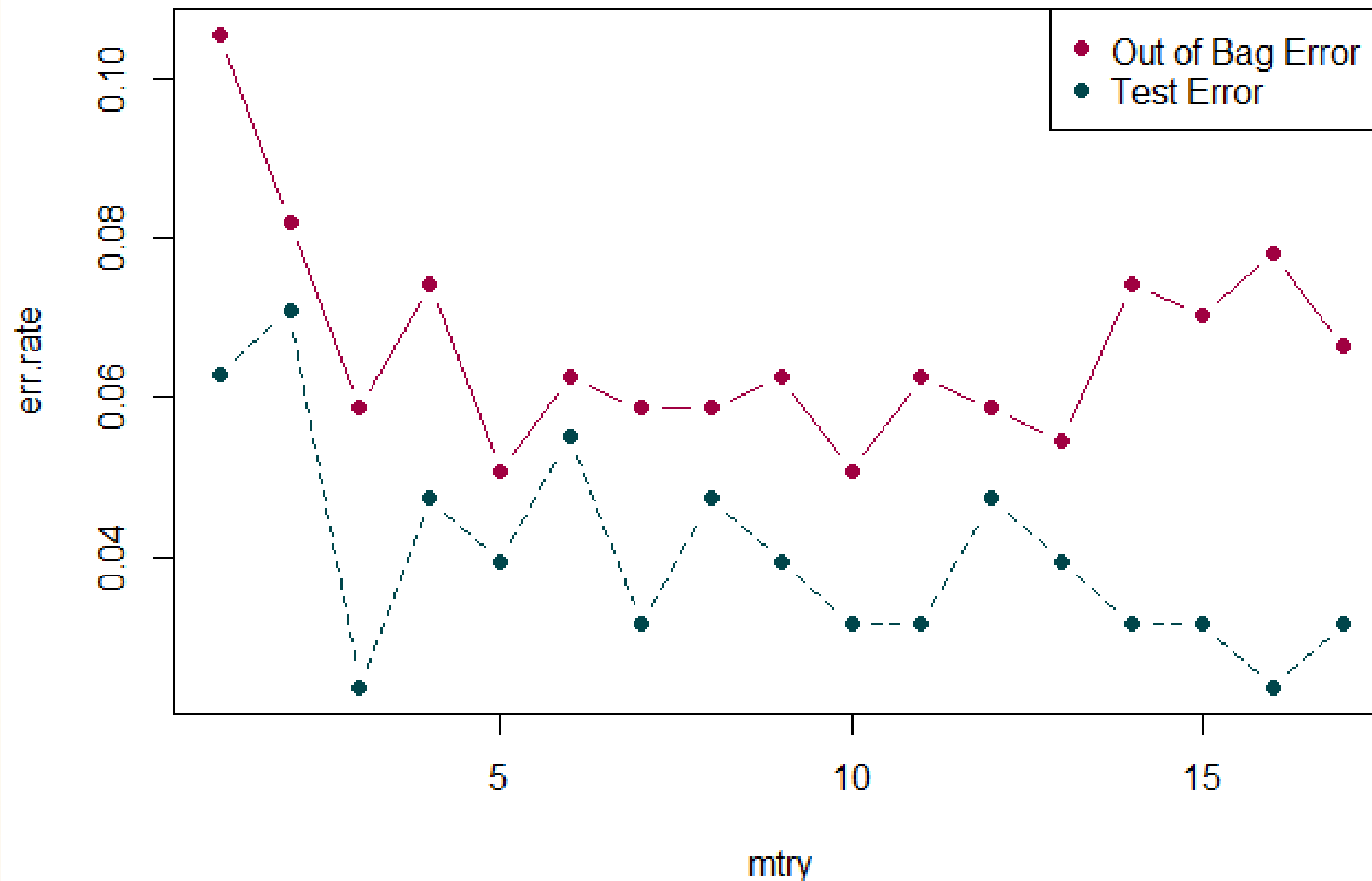
Random Forest



L'error rate
dell'oob si
stabilizza
con un
numero di
alberi pari
circa a **50**



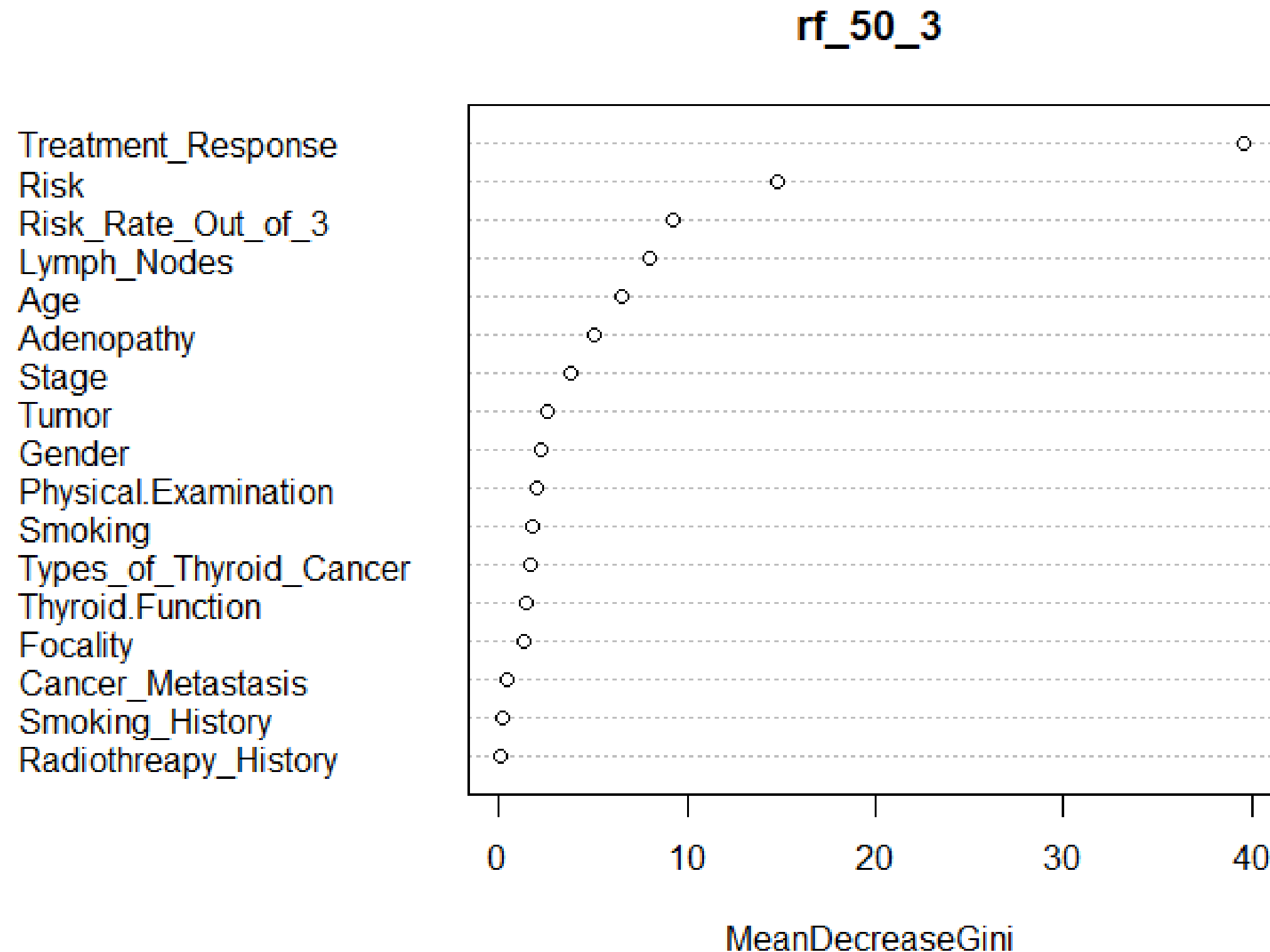
Random Forest



Il grafico mostra che il valore migliore di mtry (con ntree impostato a 50) è circa **5**



Random Forest



Dal grafico si può notare che la variabile **Treatment response** è la più importante del modello



UNIVERSITÀ DEGLI
STUDI DI FIRENZE

Bayesian Additive Regression Trees

Passi inferenziali nell'analisi bayesiana

- Sulla base delle conoscenze a priori si specifica un modello di probabilità iniziale (prior);
- Si aggiorna il modello sulla base dei dati osservati (posterior);
- Infine, si valuta la bontà del modello e la sensibilità delle conclusioni.



Bayesian Additive Regression Trees

Passi inferenziali nell'analisi bayesiana

- Sulla base delle conoscenze a priori si specifica un modello di probabilità iniziale (prior);
- Si aggiorna il modello sulla base dei dati osservati (posterior);
- Infine, si valuta la bontà del modello e la sensibilità delle conclusioni.

Modello assunto:

$$Y = f(\mathbf{x}) + \varepsilon_i \quad \varepsilon \sim N(0, \sigma^2)$$

$$f(\mathbf{x}) = \sum_{j=1}^m g(\mathbf{x}; \mathcal{T}_j, \mathcal{M}_j)$$

with \mathcal{T}_j = structure of the tree j , m = number of trees
 $\mathcal{M}_j = (\mu_1, \dots, \mu_{M_j})$: vector of parameters for the leaves



UNIVERSITÀ DEGLI
STUDI DI FIRENZE

Bayesian Additive Regression Trees

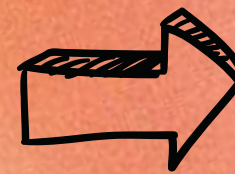
```
> summary(modello_BART)
bartMachine v1.3.4.1 for classification

training data size: n = 256 and p = 57
built in 1.8 secs on 1 core, 50 trees, 1000 burn-in and 2000 post. samples
```

confusion matrix:

	predicted 1	predicted 0	model errors
actual 1	69.000	7.000	0.092
actual 0	1.000	179.000	0.006
use errors	0.014	0.038	0.031

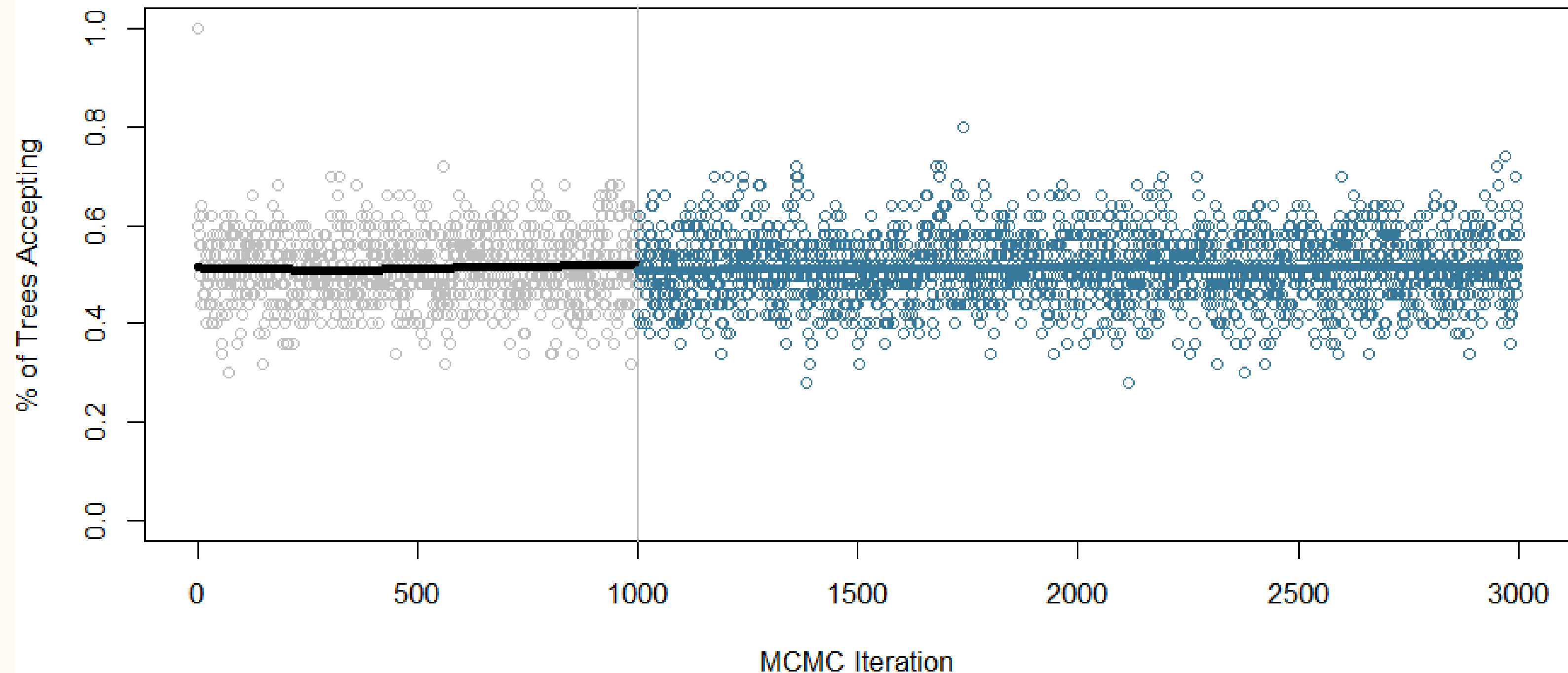
```
> table(predictions, Y_test)
      Y_test
predictions 0  1
      0 94  4
      1  1 28
```



T_EC: 0.03937



Percent Acceptance by MCMC Iteration

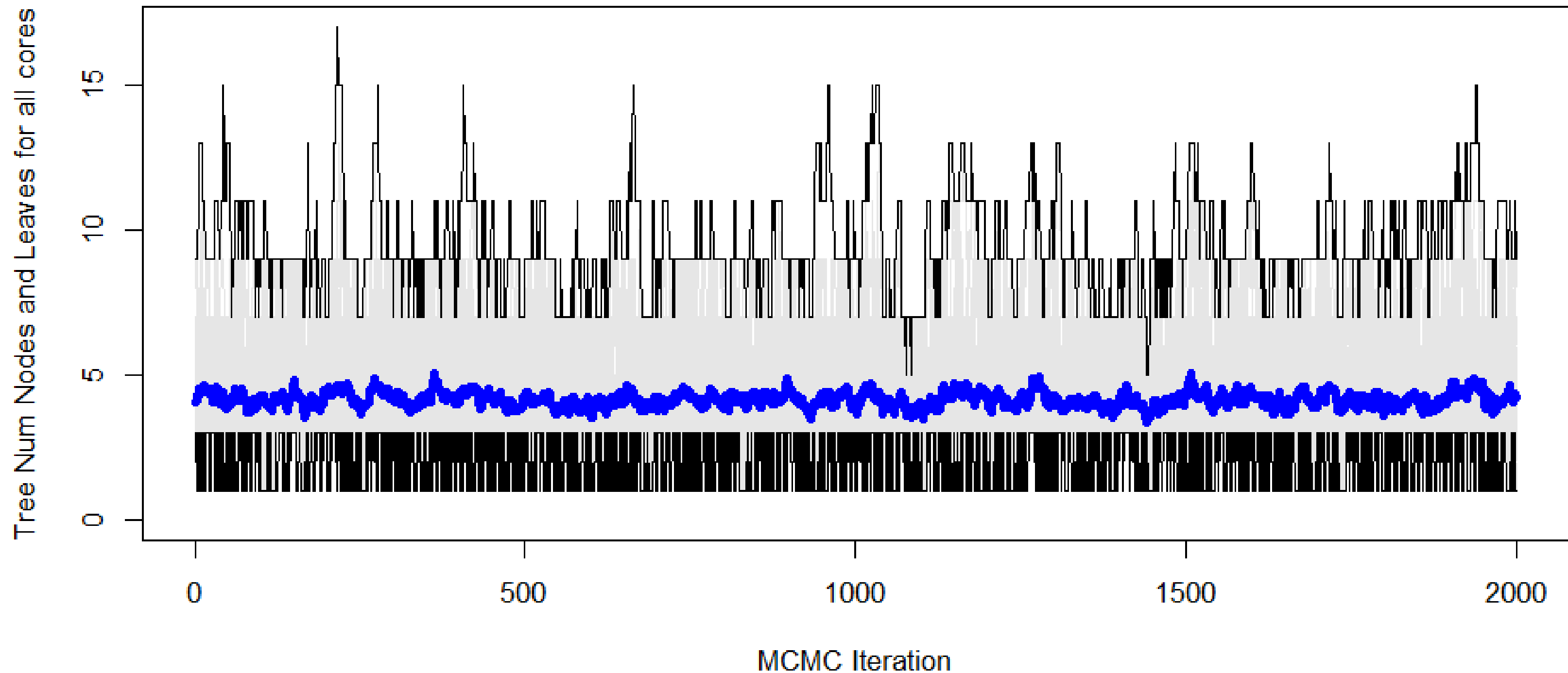




UNIVERSITÀ DEGLI
STUDI DI FIRENZE

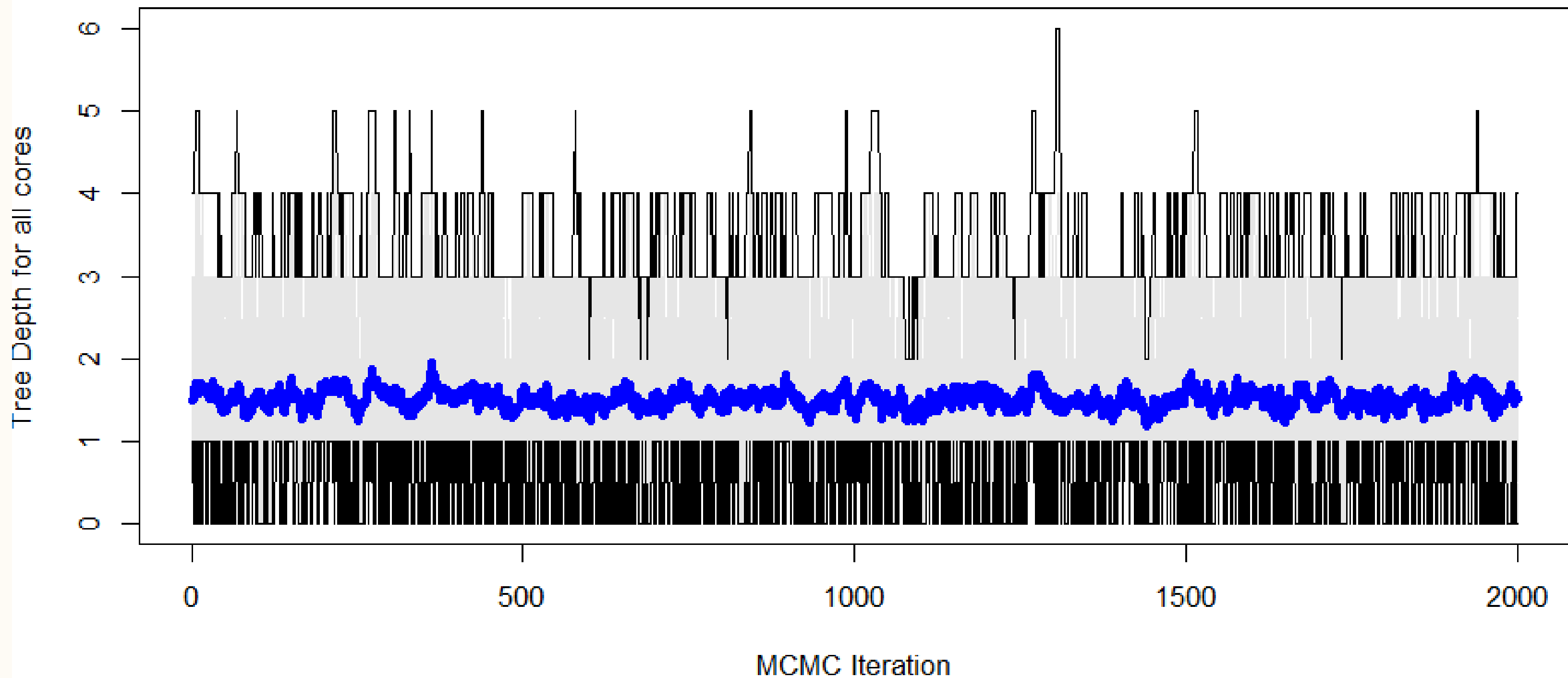
BART

Tree Num Nodes And Leaves by
MCMC Iteration After Burn-in



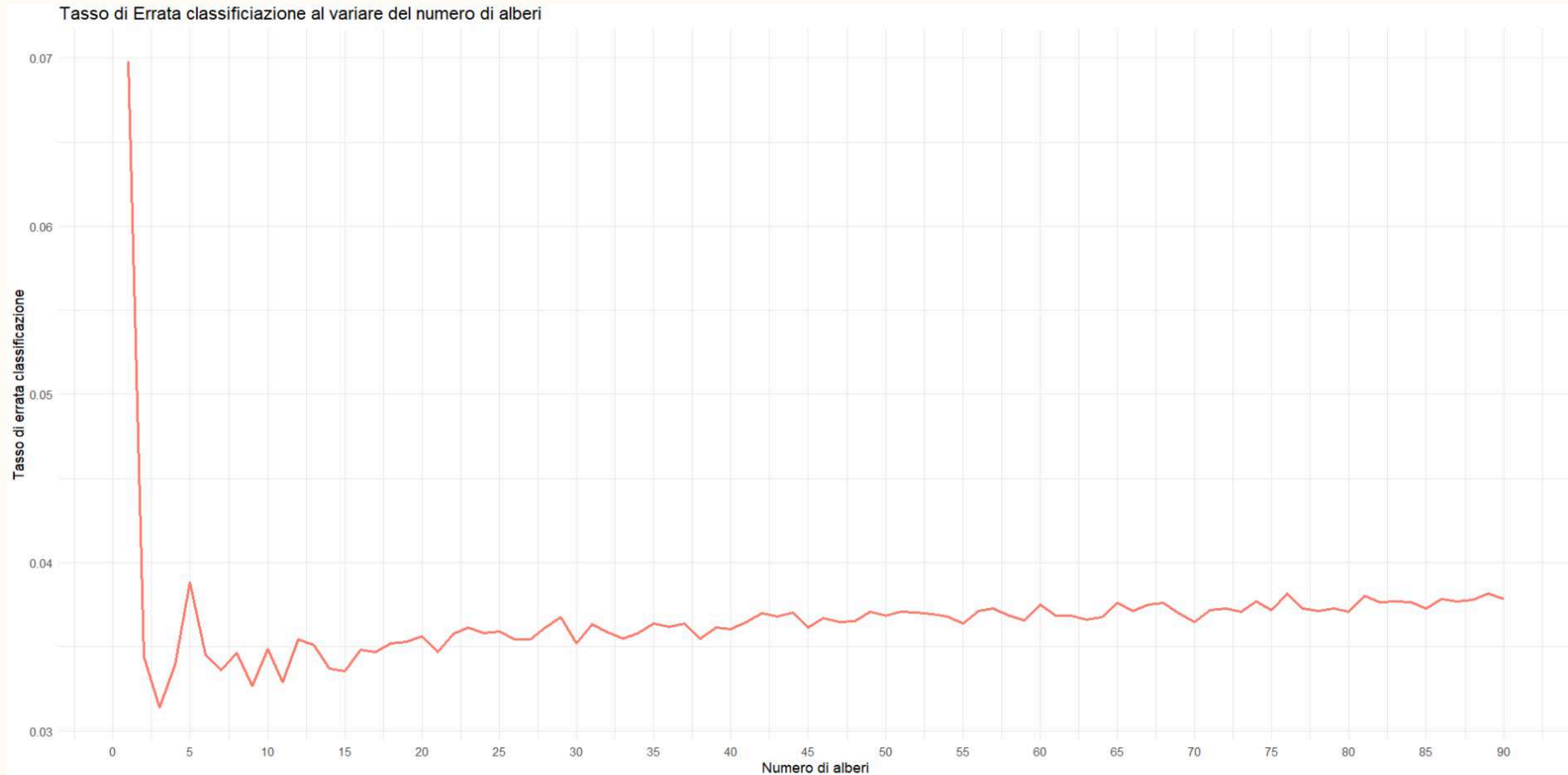


Tree Depth by MCMC Iteration After Burn-in





Cosa succede al variare del numero di alberi ?





Modello con il numero di alberi ottimale

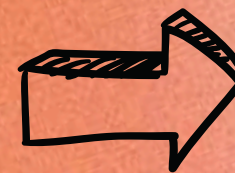
```
> summary(modello_BART)
bartMachine v1.3.4.1 for classification

training data size: n = 256 and p = 57
built in 0.4 secs on 1 core, 3 trees, 1000 burn-in and 2000 post. samples

confusion matrix:
```

	predicted 1	predicted 0	model errors
actual 1	69.000	7.000	0.092
actual 0	2.000	178.000	0.011
use errors	0.028	0.038	0.035

```
> table(predictions, Y_test)
      Y_test
predictions 0  1
      0 94  3
      1  1 29
```



T_EC: 0.03149