

# Metodi di classificazione basati su alberi

Leonardo Michi

Laurea in Statistica

A.A. 2022/2023

# Alberi decisionali

Metodi di  
classificazione  
basati su  
alberi

Leonardo  
Michi

Gli **alberi decisionali** sono un particolare tipo di modelli predittivi non lineari. L'idea è quella di considerare una partizione ricorsiva dello spazio dei predittori e applicare un *modello semplice* per ogni regione della partizione. La partizione ricorsiva è rappresentata attraverso un albero.

Gli alberi di **classificazione** sono metodi che partizionano lo spazio dei predittori in regioni disgiunte e classificano le osservazioni sulla base della regione in cui l'osservazione cade.

# Un esempio di albero decisionale

Metodi di  
classificazione  
basati su  
alberi

Leonardo  
Michi

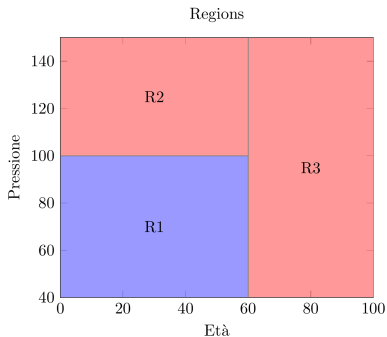
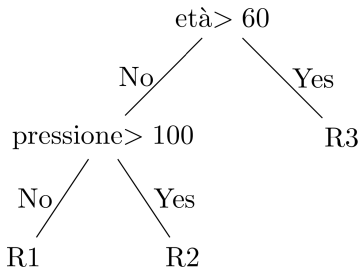


Figura 1: Esempio di albero di classificazione sinistra e partizione risultante destra

# Alberi di classificazione

Metodi di  
classificazione  
basati su  
alberi

Leonardo  
Michi

Dato il *training set*  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$ , con  $\mathbf{x}_i \in \mathbb{R}^p$  e  $y_i \in \{1, \dots, K\}$ , per un generico valore  $\mathbf{x} = (x_1, \dots, x_p)$  dello spazio dei predittori e un valore soglia  $s$  si individua la partizione binaria

$$R_1 = \{\mathbf{x} | x_j < s\}, \quad R_2 = \{\mathbf{x} | x_j > s\}$$

e si ripete il processo sulle sotto regioni finché non è soddisfatto lo *stopping criterion*.

Individuate le regioni nodi terminali  $R_1, \dots, R_M$ , ciascuna contenente  $n_m$  osservazioni,  $m \in \{1, \dots, M\}$ , si ottiene la stima della probabilità della classe  $k$

$$\hat{p}_{mk} = \frac{1}{n_m} \sum_{\mathbf{x}_i \in R_m} \mathbb{1}_{\{y_i=k\}}$$

e, in un problema di classificazione con  $K$  classi, si definisce la classe

$$k(m) = \operatorname{argmax}_{k \in \{1, \dots, K\}} \hat{p}_{mk}.$$

# Misure di impurità

La partizione binaria è definita minimizzando una misura di *impurità*. Possibili misure di impurità sono:

- *Mis-classification error*

$$E = 1 - \max_k (\hat{p}_{mk})$$

- Indice di *Gini* :

$$G = \sum_{k=1}^K \hat{p}_{mk}(1 - \hat{p}_{mk})$$

- Indice di *Entropia*:

$$D = - \sum_{k=1}^K \hat{p}_{mk} \log \hat{p}_{mk}$$

# Tree pruning

Metodi di  
classificazione  
basati su  
alberi

Leonardo  
Michi

Le procedure di tree pruning hanno lo scopo di evitare l'overfitting. Considerando l'albero più grande  $T_0$  il tree pruning mira a individuare dei sotto alberi adeguatamente selezionati. L'approccio che si basa sulla **complessità dei costi**, si sceglie l'albero  $T \subset T_0$  tale che sia minima la quantità:

$$C_\alpha(T) = \sum_{m=1}^{|T|} n_m Q_m(T) + \alpha |T|$$

Dove:

- $|T|$  è il numero di nodi terminali in  $T$
- $Q_m(T) = 1 - \frac{1}{n_m} \sum_{\mathbf{x}_i \in R_m} \mathbb{1}_{\{y_i = k(m)\}}$
- $\alpha$  è un parametro di tuning da scegliere opportunamente.

# Bagging e Random Forest

Metodi di  
classificazione  
basati su  
alberi

Leonardo  
Michi

Il **bagging** o bootstrap aggregation consiste in una procedura generale che consente di migliorare l'accuratezza della classificazione rispetto a un singolo albero. A partire dal training set, si ottengono  $B$  campioni bootstrap e su ciascuno di essi si costruisce un albero di classificazione. Indicata con  $p_k(\mathbf{x})$  la frazione di *alberi bagging* che classificano  $\mathbf{x}$  nella classe  $k$ , la classificazione finale si ottiene come:

$$\hat{k}_{\text{bag}}(\mathbf{x}) = \underset{k \in \{1, \dots, K\}}{\operatorname{argmax}} \{p_1(\mathbf{x}), \dots, p_K(\mathbf{x})\}$$

In presenza di predittori forti il bagging potrebbe non fornire un miglioramento sostanziale dell'accuratezza. Le **random forest** rappresentano una variante del bagging in cui ad ogni split si considera un sottoinsieme di  $q < p$  variabili ad ogni split. (usualmente  $q = \sqrt{p}$ )

# Boosting

Metodi di  
classificazione  
basati su  
alberi

Leonardo  
Michi

Il **Boosting** è una tecnica che si basa sulla costruzione *sequenziale* di alberi. L'obiettivo è quello di applicare sequenzialmente l'algoritmo di classificazione debole a versioni dei dati ripetutamente modificate in modo da avere una sequenza di classificatori deboli  $T^{(b)}$  con  $b = 1, \dots, B$ .

Il risultato finale della classificazione si ottiene come :

$$\hat{k}_{boost}(x) = \underset{k \in \{1, \dots, K\}}{\operatorname{argmax}} \sum_{b=1}^B \alpha_b \mathbb{1}_{\{T^{(b)}(x)=k\}}$$

dove gli  $\alpha_b$  sono pesi costruiti ad ogni passo in modo da tener conto delle osservazioni classificate non correttamente.



# Misure di importanza delle variabili

Metodi di  
classificazione  
basati su  
alberi

Leonardo  
Michi

Alcune misure di importanza delle variabili sono:

- La *Mean decrease Gini* che si ottiene considerando una misura costruita come ammontare della decrescita dell'indice di Gini che risulta dagli split su una variabile (media sugli alberi ).
- La *Mean decrease accuracy* che si basa sull'ammontare della decrescita dell'accuratezza quando una certa variabile è permutata (media sugli alberi ).

# Model assessment

Metodi di  
classificazione  
basati su  
alberi

Leonardo  
Michi

		previsto	
		positivo	negativo
Effettivo	Popolazione totale (PT) positivo	Vero positivo (VP)	Falso negativo (FN)
	negativo	Falso positivo (FP)	Vero negativo (VN)

$$Acc = \frac{VP + VN}{PT} \quad spec = \frac{VN}{VN + FP} \quad sens = \frac{VP}{VP + FN}$$

$$prec = \frac{VP}{VP + FP} \quad F1 = \frac{2VP}{2VP + FP + FN}$$

# Classificazione dei tumori tramite alberi decisionali

Metodi di  
classificazione  
basati su  
alberi

Leonardo  
Michi

Il dataset è composto da 569 osservazioni di 30 variabili per caratterizzare forma e dimensioni del tumore. 357 osservazioni sono della classe benigni e 212 maligni. Si utilizzano 250 osservazioni come *training set* e le restanti come *test set*.

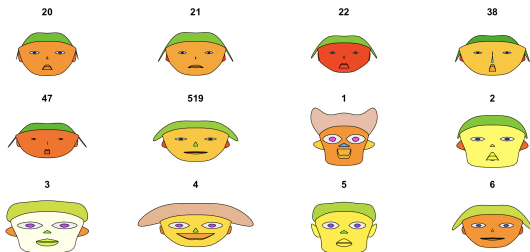


Figura 2: Facce di Chernoff

# Classificazione dei tumori tramite alberi decisionali 2

Metodi di  
classificazione  
basati su  
alberi

Leonardo  
Michi

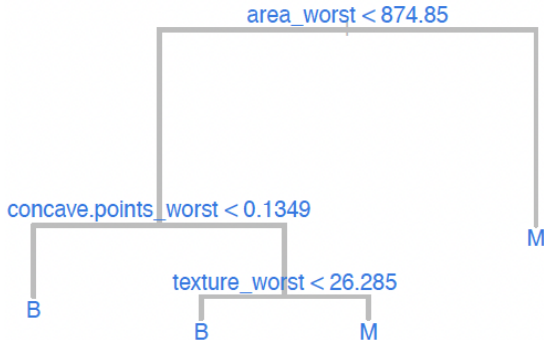


Figura 3: Esempio di singolo albero di classificazione dopo il pruning con accuratezza pari a 95.3%

# Classificazione dei tumori tramite alberi decisionali 3

Metodi di  
classificazione  
basati su  
alberi

Leonardo  
Michi

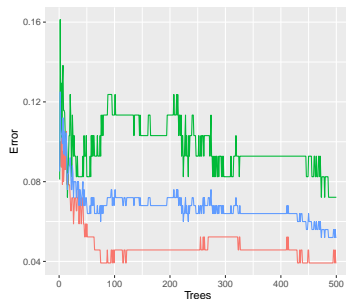
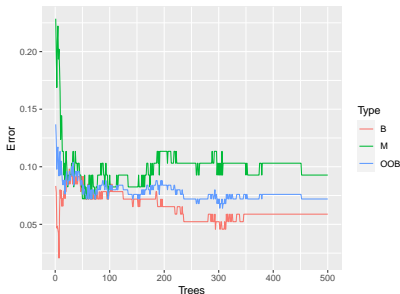


Figura 4: Errori di classificazione per bagging sinistra e random forest destra



# Some references

- James, G., Witten, D., Hastie, T., & Tibshirani, R. (2021). An introduction to statistical learning (Second Edition). New York: Springer.
- Breiman, L. Random Forests. Machine Learning 45, 5–32 (2001).
- Biau, G., Scornet, E. A random forest guided tour. TEST 25, 197–227 (2016).
- Introduzione ai Modelli Statistici Giovanni M. Marchetti Dipartimento di Statistica, Informatica, Applicazioni, Firenze 2013, 2019 rev. 1
- Wasserman, L. (2004). All of statistics. Springer-Verlag, New York.
- Friedman, J., Hastie, T., & Tibshirani, R. (2013). The elements of statistical learning. Second edition. Springer, Berlin: Springer series in statistics.