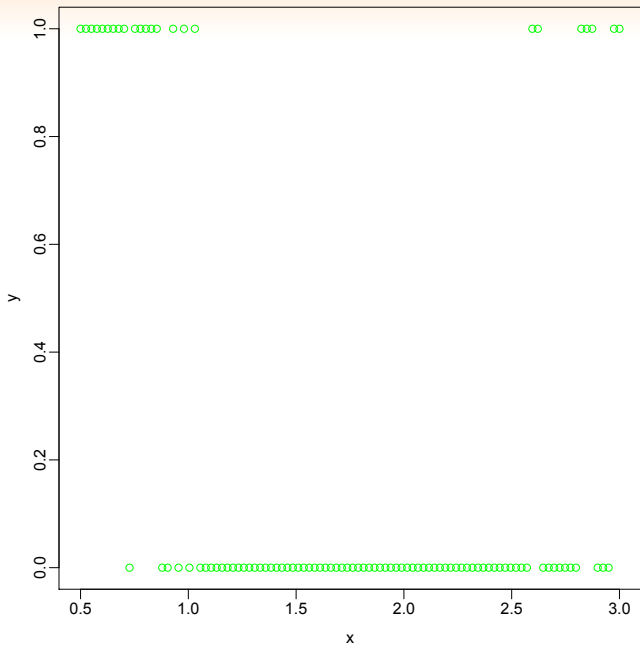


Alberi di classificazione

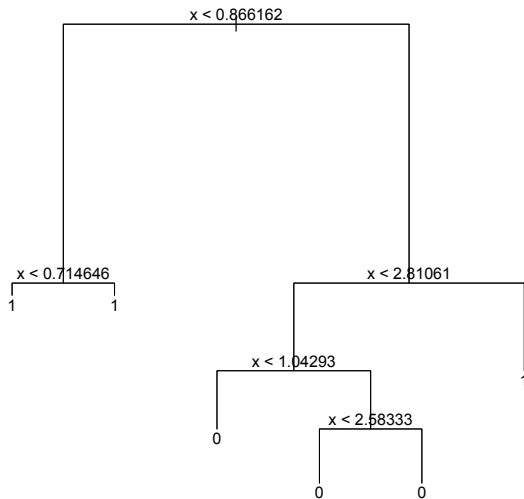


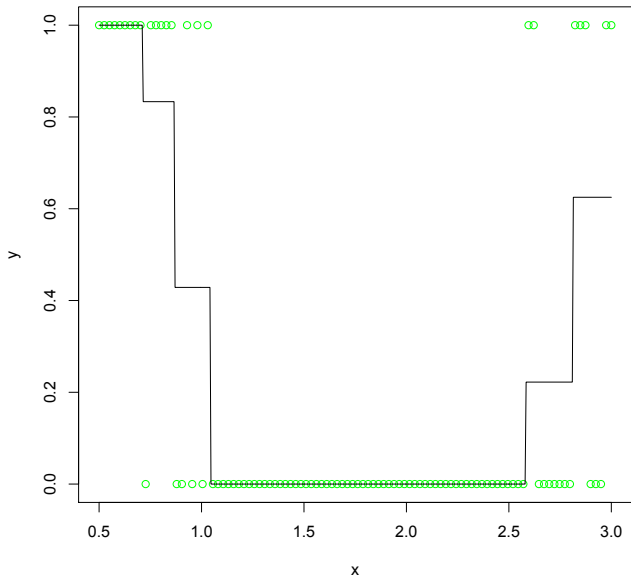
Se si indicano le due classi con 0 e 1 si potrà scrivere:

$$p(x) = P(y = 1|x)$$

$$\hat{p}(x) = \sum_{j=1}^J P_j I(x \in R_j)$$

con P_j che rappresenta la probabilità che $y = 1$ per la regione considerata.





In questo caso il valore associato alle foglie sarà 0 o 1. Quando si fa cadere un'osservazione lungo l'albero e raggiunge una foglia con la probabilità associata, questa verrà allocata alla classe $C(\hat{p}(x))$ con $C(p) = 0$ se $p \leq 1/2$ e $C(p) = 1$ altrimenti.

Per stimare i P_j si usa la media aritmetica:

$$\hat{P}_j = M(y_i : x_i \text{ in } R_j) = \frac{1}{n_j} \sum_{i \in R_j} I(y_i = 1)$$

che rappresenta la frequenza relativa nella regione considerata.

In questo caso il valore associato alle foglie sarà 0 o 1. Quando si fa cadere un'osservazione lungo l'albero e raggiunge una foglia con la probabilità associata, questa verrà allocata alla classe $C(\hat{p}(x))$ con $C(p) = 0$ se $p \leq 1/2$ e $C(p) = 1$ altrimenti.

Per stimare i P_j si usa la media aritmetica:

$$\hat{P}_j = M(y_i : x_i \text{ in } R_j) = \frac{1}{n_j} \sum_{i \in R_j} I(y_i = 1)$$

che rappresenta la frequenza relativa nella regione considerata.

Ovviamente ora non si considera più la classica devianza ma si utilizza la devianza connessa alla distribuzione binomiale.

$$D = -2 \sum_{i=1}^n (y_i \log \hat{p}_i + (1 - y_i) \log(1 - \hat{p}_i))$$

o anche

$$D = \sum_{j=1}^J -2n_j [\hat{P}_j \log \hat{P}_j + (1 - \hat{P}_j) \log(1 - \hat{P}_j)] = \sum_j D_j$$

Si dimostra inoltre che

$$D = 2n \sum_j (n_j/n) Q(\hat{P}_j)$$

che, a meno della costante $2n$, é la media delle entropie pesate con la numerosità delle foglie:

$$\begin{aligned} Q(P_j) &= - \sum_{k=0,1} P_{jk} \log P_{jk} = -[P_{j1} \log P_{j1} + P_{j0} \log P_{j0}] \\ &= -[P_{j1} \log P_{j1} + (1 - P_{j1}) \log(1 - P_{j1})] \end{aligned}$$

I termini Q rappresentano una misura di impurità, cioè un indicatore del fatto che gli elementi di una certa foglia sono disomogenei rispetto alla variabile di risposta. $Q(p) = 0$ infatti se $p = 0$ o $p = 1$, mentre aumenta quanto più ci sia avvicina al valore $1/2$.

La misura di entropia si può quindi sostituire con altre misure di impurità, ad esempio con l'indice di Gini.

$$Q(P_j) = \sum_{k=0,1} P_{jk}(1 - P_{jk})$$

Con questo adattamento, l'algoritmo di crescita è analogo a quello degli alberi di regressione. Lo stesso vale per la procedura di potatura. In alcuni casi si utilizza anche il semplice tasso di errata classificazione come misura di discrepanza.

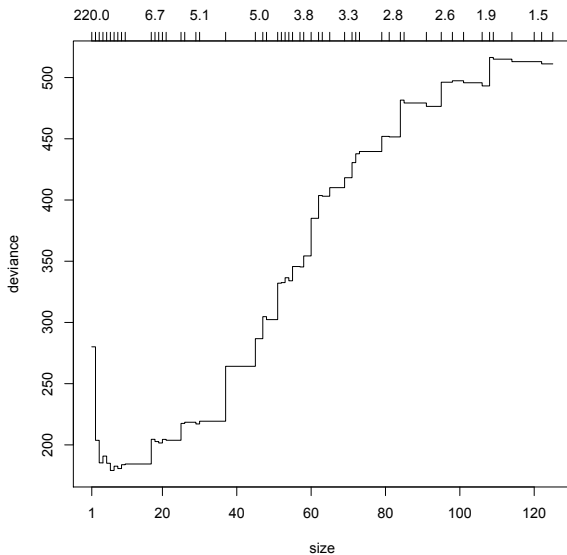
I termini Q rappresentano una misura di impurità, cioè un indicatore del fatto che gli elementi di una certa foglia sono disomogenei rispetto alla variabile di risposta. $Q(p) = 0$ infatti se $p = 0$ o $p = 1$, mentre aumenta quanto più ci sia avvicina al valore $1/2$.

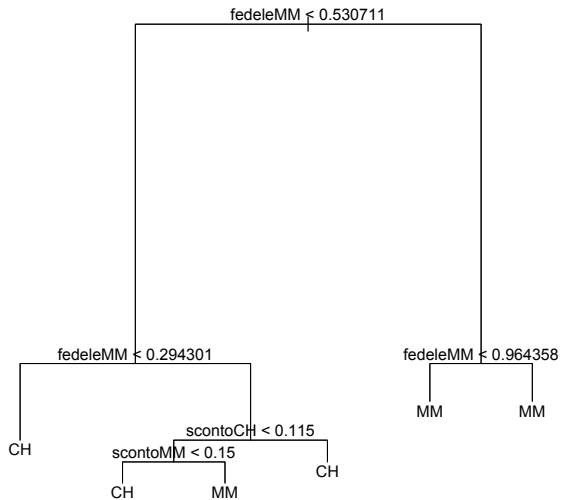
La misura di entropia si può quindi sostituire con altre misure di impurità, ad esempio con l'indice di Gini.

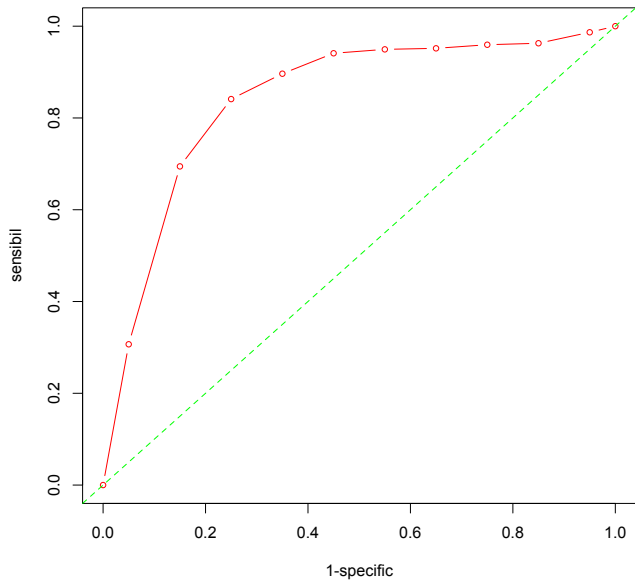
$$Q(P_j) = \sum_{k=0,1} P_{jk}(1 - P_{jk})$$

Con questo adattamento, l'algoritmo di crescita è analogo a quello degli alberi di regressione. Lo stesso vale per la procedura di potatura. In alcuni casi si utilizza anche il semplice tasso di errata classificazione come misura di discrepanza.

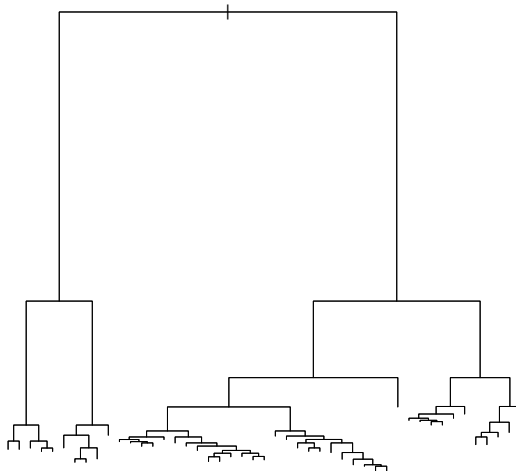
Succhi di frutta: 600 osservazioni per la stima e 202 per la potatura.

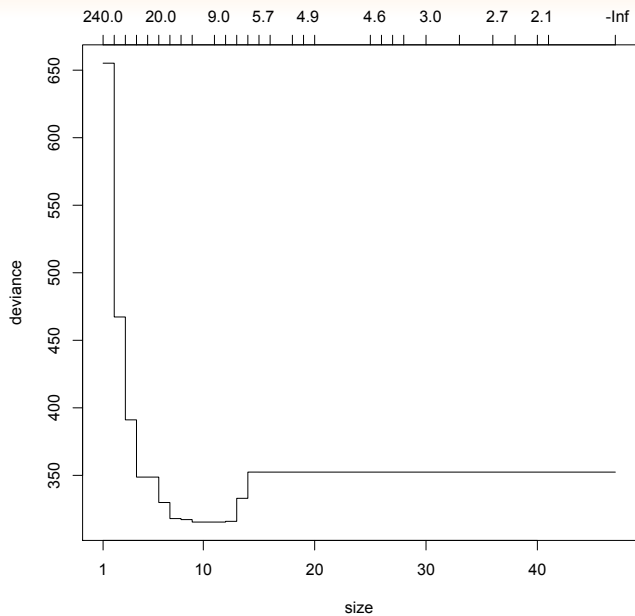


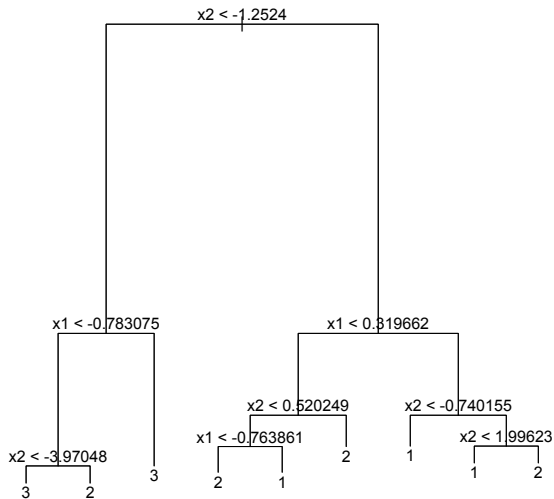


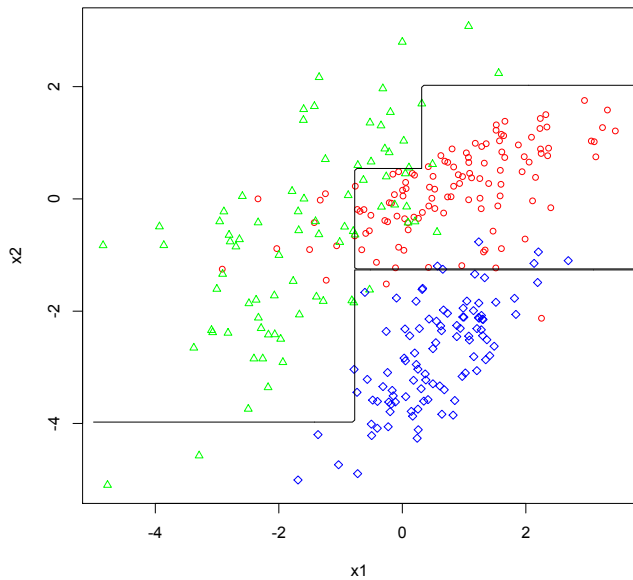


Se K é maggiore di 2 bisogna considerare degli adattamenti, ma il meccanismo sostanzialmente non cambia:







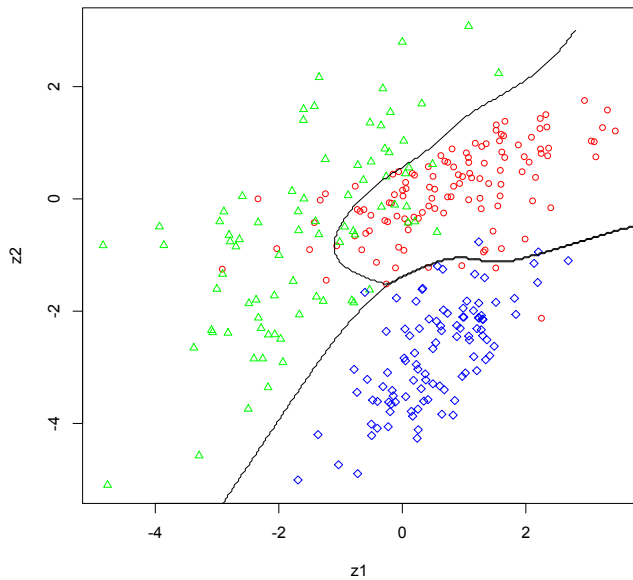


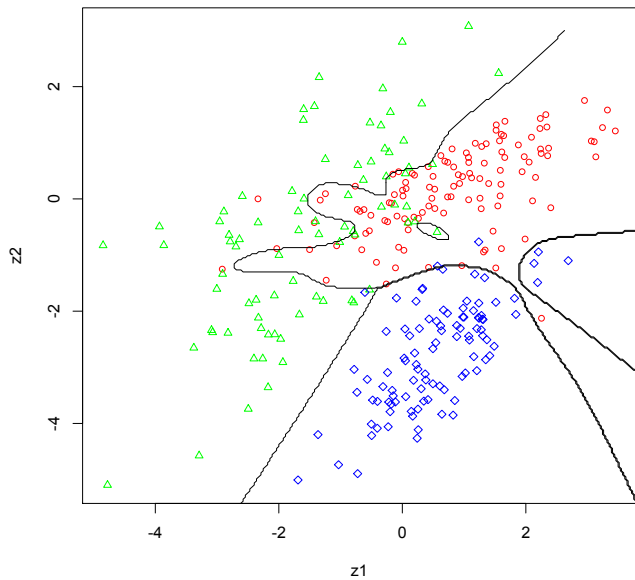
Reti neurali

Anche in questo caso si può estendere la procedura proposta per le variabili quantitative. L'adattamento da introdurre é nella funzione di attivazione f_1 che deve avere come codominio l'intervallo $0, 1$. Per tali motivi in genere la piú utilizzata é la funzione logistica. In alcuni casi si possono configurare le due classi con -1 e 1 e la funzione utilizzata sar :

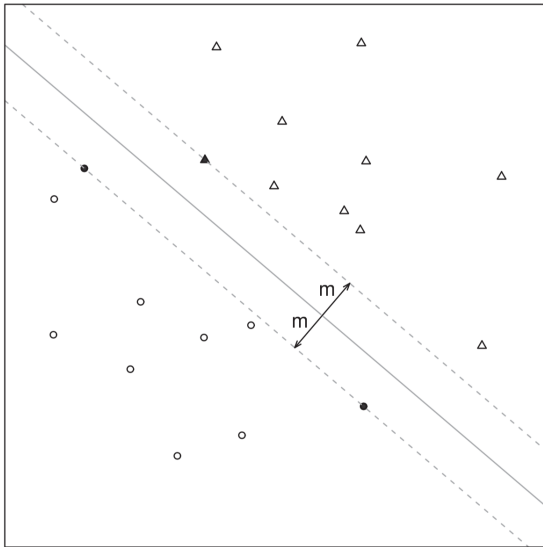
$$2I(x) - 1 = \frac{e^x - 1}{e^x + 1} = \tanh(x/2)$$

Se le classi sono piú di due si creano K variabili risposta di tipo $0 - 1$. Inoltre il termine che entra nella funzione obiettivo non é piú la distanza euclidea, ma l'entropia.





Support vector machines



L'equazione

$$\beta_0 + x'\beta = 0$$

individua un generico iperpiano candidato a separare le due classi, alle quali vengono convenzionalmente assegnati i valori -1 e 1.

L'iperpiano ottimale sarà ovviamente quello che massimizza la distanza m . Il problema di ottimizzazione può essere formalizzato come segue:

$$\max_{\beta_0, \beta} m$$

con i seguenti vincoli:

$$||\beta|| = 1$$

$$y_i(\beta_0 + \tilde{x}_i'\beta) \geq m \quad (i = 1, \dots, n)$$

dove (\tilde{x}, y) rappresenta una generica unità da classificare.

Per semplificare la condizione $\|\beta\| = 1$ i vincoli vengono riscritti nella seguente maniera

$$y_i(\beta_0 + \tilde{x}_i'\beta) \geq m\|\beta\|$$

e il problema diventa:

$$\min_{\beta_0, \beta} \frac{1}{2} \|\beta\|^2$$

con vincoli

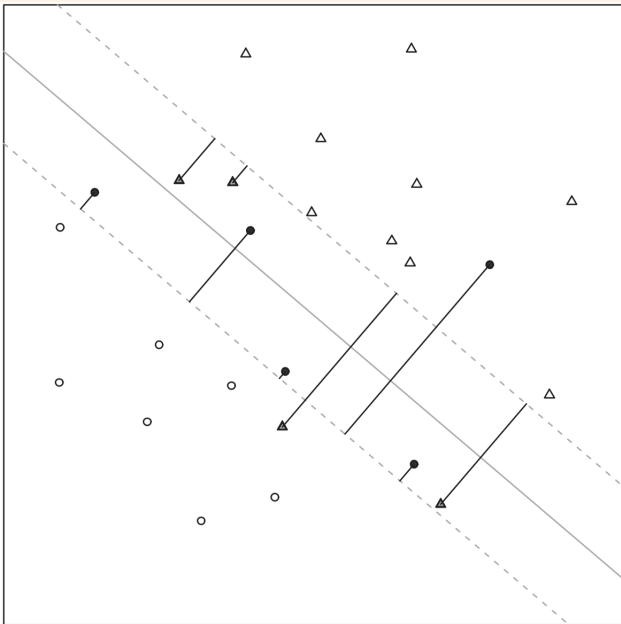
$$y_i(\beta_0 + \tilde{x}_i'\beta) \geq 1$$

Per la soluzione del problema si utilizzano elementi di geometria, ma ovviamente, in situazioni reali, é impossibile trovare una retta che separi perfettamente le due classi. É necessario quindi introdurre un requisito meno stringente, tollerando che alcuni punti saranno comunque mal classificati.

A tal fine vengono introdotte delle variabili ausiliarie a componenti non negative $\xi = (\xi_1, \dots, \xi_n)$ che esprimono di quanto i vari punti stanno al di la' della linea di margine della loro classe.

Per la soluzione del problema si utilizzano elementi di geometria, ma ovviamente, in situazioni reali, é impossibile trovare una retta che separi perfettamente le due classi. É necessario quindi introdurre un requisito meno stringente, tollerando che alcuni punti saranno comunque mal classificati.

A tal fine vengono introdotte delle variabili ausiliarie a componenti non negative $\xi = (\xi_1, \dots, \xi_n)$ che esprimono di quanto i vari punti stanno al di la' della linea di margine della loro classe.



In questo caso il problema di ottimizzazione viene adattato sostituendo i vincoli precedenti con la seguente forma

$$y_i(\beta_0 + \tilde{x}_i' \beta) \geq m(1 - \xi)$$

e l'ottimizzazione diventa:

$$\min_{\beta_0, \beta} \frac{1}{2} \|\beta\|^2 + \gamma \sum_{i=1}^n \xi_i$$

con vincoli

$$y_i(\beta_0 + \tilde{x}_i' \beta) \geq 1 - \xi_i$$

$$\xi_i \geq 0$$

con γ costante positiva che é un parametro di regolazione e rappresenta il costo di violazione delle barriere.

Si dimostra che:

$$\hat{\beta} = \sum_{i=1}^n a_i y_i \tilde{x}_i$$

dove solo alcuni a_i sono non nulli, quindi la soluzione é esprimibile solo con alcune osservazioni, dette support vectors.

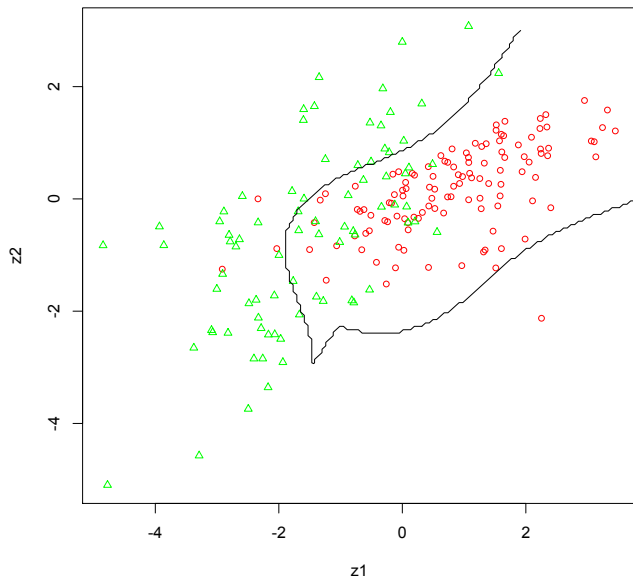
In alcuni casi é opportuno considerare delle trasformazioni delle variabili esplicative, sfruttando diverse funzioni nucleo che portano a risultati differenti.

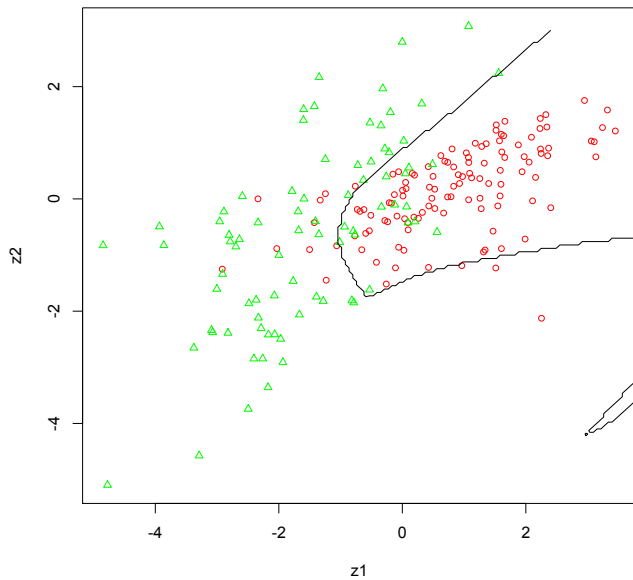
Si dimostra che:

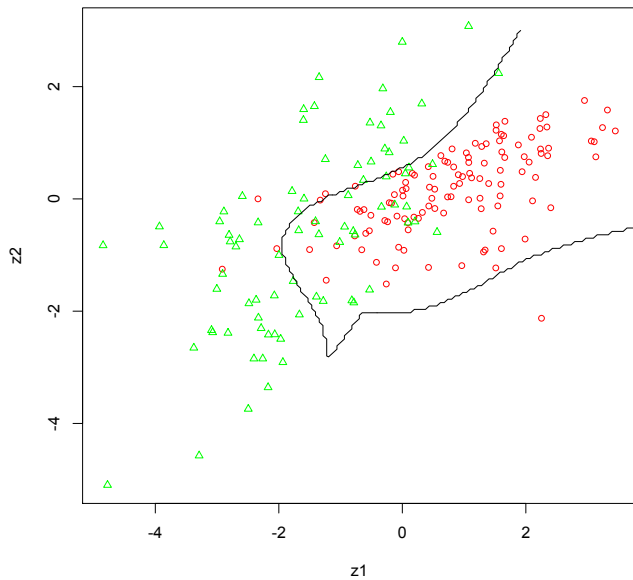
$$\hat{\beta} = \sum_{i=1}^n a_i y_i \tilde{x}_i$$

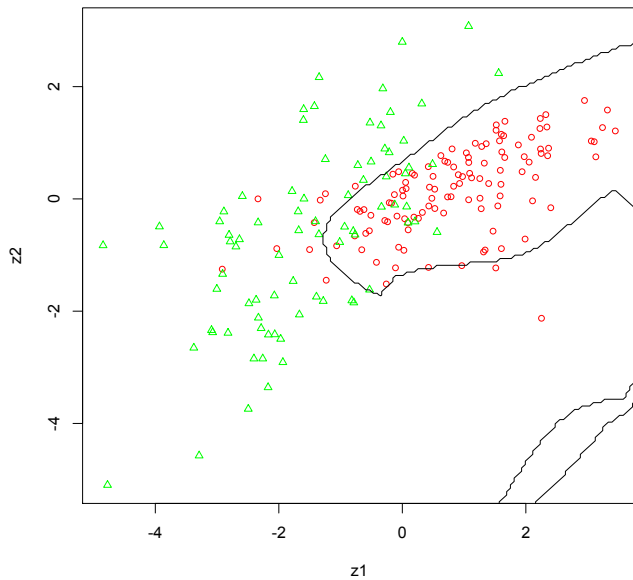
dove solo alcuni a_i sono non nulli, quindi la soluzione é esprimibile solo con alcune osservazioni, dette support vectors.

In alcuni casi é opportuno considerare delle trasformazioni delle variabili esplicative, sfruttando diverse funzioni nucleo che portano a risultati differenti.









Combinazioni di classificatori

In molte situazioni diversi modelli sembrano adattarsi bene ai dati e non sembra esserci uno preferito agli altri.

Ad esempio, se abbiamo 30 variabili esplicative e costruiamo un modello di regressione logistica con 5 variabili, ci saranno circa 140mila possibili gruppi che si possono formare.

Andando a calcolare l'errore di previsione su un insieme di verifica, sarà di poco diverso considerando tutti i modelli costruiti. Quindi i modelli saranno praticamente equivalenti per effettuare la previsione, ma ognuno ci spiegherà in maniera diversa il problema, a seconda del tipo di variabili che entreranno in gioco.

Combinazioni di classificatori

In molte situazioni diversi modelli sembrano adattarsi bene ai dati e non sembra esserci uno preferito agli altri.

Ad esempio, se abbiamo 30 variabili esplicative e costruiamo un modello di regressione logistica con 5 variabili, ci saranno circa 140mila possibili gruppi che si possono formare.

Andando a calcolare l'errore di previsione su un insieme di verifica, sarà di poco diverso considerando tutti i modelli costruiti. Quindi i modelli saranno praticamente equivalenti per effettuare la previsione, ma ognuno ci spiegherà in maniera diversa il problema, a seconda del tipo di variabili che entreranno in gioco.

Se invece utilizziamo per esempio le reti neurali o gli alberi di classificazione, la stima sarà fortemente influenzata dall'insieme di dati usato per la stima. Si é provato per esempio che se si prende un insieme di stima e poi si elimina una piccola percentuale di osservazioni (2 – 3%) il modello può cambiare completamente anche se non varia l'errore di previsione.

Per superare tali inconvenienti, una possibilità é quella di combinare le previsioni ottenute da metodi diversi. Sono state in realtà seguite diverse strade, ma ciascuna cerca di raccogliere le qualità delle singole componenti portando spesso a previsioni più accurate.

Se invece utilizziamo per esempio le reti neurali o gli alberi di classificazione, la stima sarà fortemente influenzata dall'insieme di dati usato per la stima. Si é provato per esempio che se si prende un insieme di stima e poi si elimina una piccola percentuale di osservazioni (2 – 3%) il modello può cambiare completamente anche se non varia l'errore di previsione.

Per superare tali inconvenienti, una possibilità é quella di combinare le previsioni ottenute da metodi diversi. Sono state in realtà seguite diverse strade, ma ciascuna cerca di raccogliere le qualità delle singole componenti portando spesso a previsioni più accurate.

Bagging

- Si indica con $C(x)$ un classificatore determinato con uno dei metodi introdotti finora su un insieme di stima $Z = \{(x_1, y_1), \dots, (x_n, y_n)\}$
- Si utilizza poi la procedura bootstrap, considerando il campione Z_1^* ottenuto estraendo n volte con ripetizione gli elementi dell'insieme di stima Z .
- Utilizzando questo nuovo campione si ottiene un nuovo classificatore $C_1^*(x)$ che sarà probabilmente diverso dal precedente.

- Si ripete la procedura B volte, ottenendo i campioni $Z_b^*(b = 1, \dots, B)$ e da questa i nuovi classificatori $C_b^*(x)$.
- Si introduce quindi un nuovo classificatore che é una sintesi dei risultati ottenuti sui singoli campioni. In genere si utilizza la media aritmetica e si ha

$$C_{bag}(x) = \frac{1}{B} \sum_{b=1}^B C_b^*(x)$$

- Si classifica l'unità nella classe associata al valore 1 se $C_{bag}(x) > 1/2$ e a quella associata al valore 0 altrimenti.

L'errore di classificazione del nuovo modello risulta mediamente inferiore di quello calcolato sui singoli modelli "originali".

Il modello viene chiamato di Bootstrap AGGRegatING, da cui BAGGING.

La stessa procedura si può utilizzare per variabili quantitative, dove non si avrà più il classificatore ma quello che viene definito il previsore, cioè il valore assunto dalla variabile di risposta. In questo caso il nuovo modello avrà generalmente varianza inferiore rispetto al modello originario.

La combinazione dei risultati comporta però la perdita di ogni semplice struttura nei modelli originari, da cui deriva la difficoltà di interpretazione dei risultati.

Una procedura alternativa è il bumping che individua come classificatore il modello che ha il minor errore di previsione fra tutti i modelli ottenuti con i ricampionamenti bootstrap.

L'utilizzo di campioni casuali permette di utilizzare anche la procedura out-of-bag. La parte di dati che viene esclusa dal campione bootstrap viene praticamente usata come insieme di verifica

La stessa procedura si può utilizzare per variabili quantitative, dove non si avrà più il classificatore ma quello che viene definito il previsore, cioè il valore assunto dalla variabile di risposta. In questo caso il nuovo modello avrà generalmente varianza inferiore rispetto al modello originario.

La combinazione dei risultati comporta però la perdita di ogni semplice struttura nei modelli originari, da cui deriva la difficoltà di interpretazione dei risultati.

Una procedura alternativa è il bumping che individua come classificatore il modello che ha il minor errore di previsione fra tutti i modelli ottenuti con i ricampionamenti bootstrap.

L'utilizzo di campioni casuali permette di utilizzare anche la procedura out-of-bag. La parte di dati che viene esclusa dal campione bootstrap viene praticamente usata come insieme di verifica

La stessa procedura si può utilizzare per variabili quantitative, dove non si avrà più il classificatore ma quello che viene definito il previsore, cioè il valore assunto dalla variabile di risposta. In questo caso il nuovo modello avrà generalmente varianza inferiore rispetto al modello originario.

La combinazione dei risultati comporta però la perdita di ogni semplice struttura nei modelli originari, da cui deriva la difficoltà di interpretazione dei risultati.

Una procedura alternativa è il bumping che individua come classificatore il modello che ha il minor errore di previsione fra tutti i modelli ottenuti con i ricampionamenti bootstrap.

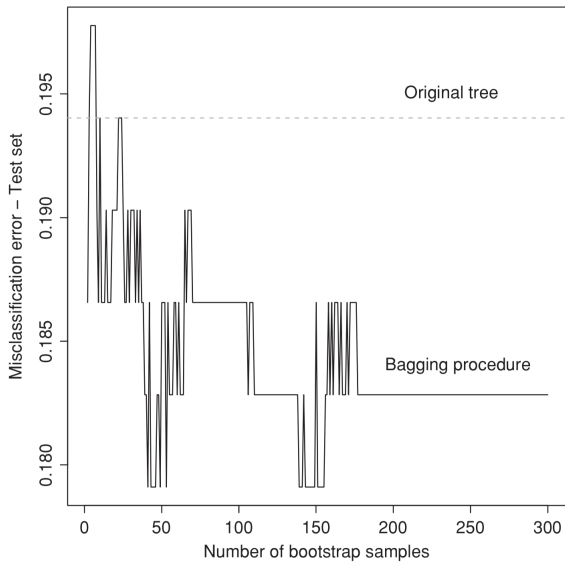
L'utilizzo di campioni casuali permette di utilizzare anche la procedura out-of-bag. La parte di dati che viene esclusa dal campione bootstrap viene praticamente usata come insieme di verifica

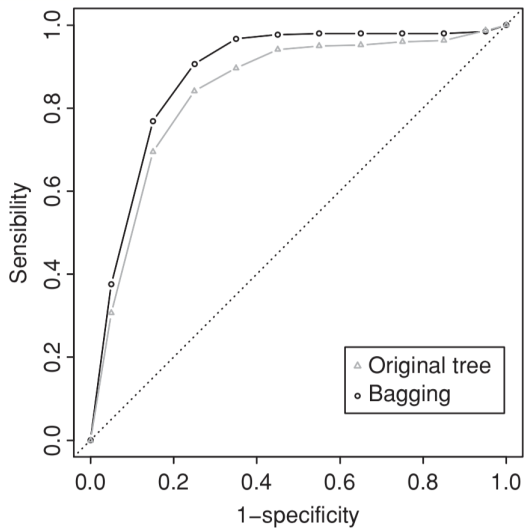
La stessa procedura si può utilizzare per variabili quantitative, dove non si avrà più il classificatore ma quello che viene definito il previsore, cioè il valore assunto dalla variabile di risposta. In questo caso il nuovo modello avrà generalmente varianza inferiore rispetto al modello originario.

La combinazione dei risultati comporta però la perdita di ogni semplice struttura nei modelli originari, da cui deriva la difficoltà di interpretazione dei risultati.

Una procedura alternativa è il bumping che individua come classificatore il modello che ha il minor errore di previsione fra tutti i modelli ottenuti con i ricampionamenti bootstrap.

L'utilizzo di campioni casuali permette di utilizzare anche la procedura out-of-bag. La parte di dati che viene esclusa dal campione bootstrap viene praticamente usata come insieme di verifica





Boosting

Anche il boosting consiste nel combinare i risultati di un modello adattato usando diversi insiemi di dati, ma selezionando le unità da inserire nel campione attraverso l'assegnazione a ciascuna unità di una diversa probabilità di entrata.

In particolare si assegna maggior peso alle osservazioni che in precedenza sono state classificate peggio: in questo modo si cerca di migliorare la prestazione del nuovo modello, agendo proprio in quelle aree dove il vecchio classificatore presentava maggiori difficoltà.

Boosting

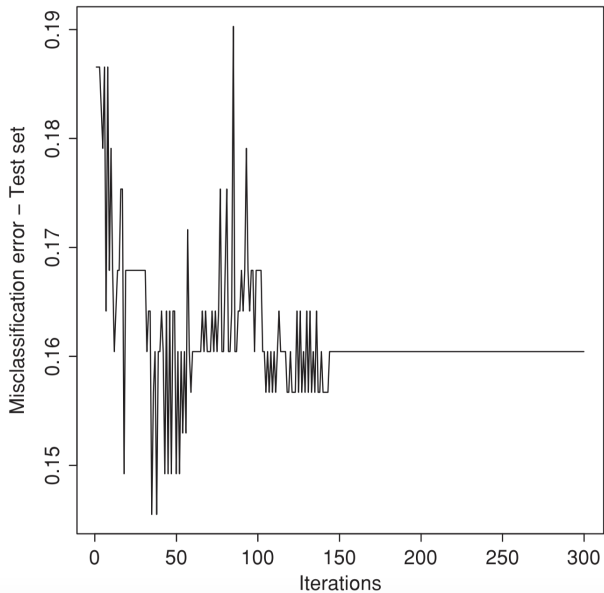
Anche il boosting consiste nel combinare i risultati di un modello adattato usando diversi insiemi di dati, ma selezionando le unità da inserire nel campione attraverso l'assegnazione a ciascuna unità di una diversa probabilità di entrata.

In particolare si assegna maggior peso alle osservazioni che in precedenza sono state classificate peggio: in questo modo si cerca di migliorare la prestazione del nuovo modello, agendo proprio in quelle aree dove il vecchio classificatore presentava maggiori difficoltà.

Si tratta di una procedura iterativa:

- classificatore iniziale costruito dando ad ogni osservazione lo stesso peso;
- si utilizza tale classificatore (definito debole) per la prima classificazione;
- l'insieme dei pesi viene aggiornato ad ogni iterazione;
- Alla fine della procedura si identifica il nuovo classificatore tenendo conto di tutti i modelli costruiti al variare delle iterazioni.

Tale logica é stata concretizzata in diversi modi, ma l'algoritmo piú utilizzato é quello definito AdaBoost.



Random forest (Foreste Casuali)

Bagging e boosting fanno variare ad ogni iterazione le unità statistiche che entrano in gioco nel modello, usando tutte le variabili esplicative disponibili. Le foreste casuali invece considerano diversi sottoinsiemi di variabili esplicative. Tali tecniche sono tipiche degli alberi di classificazione, ma usate anche per altri tipi di classificatori.

La procedura consiste nel selezionare in modo casuale, ad ogni nodo dell'albero, un piccolo gruppo di variabili esplicative che verranno ispezionate per individuare il punto di suddivisione ottimale.

In questo modo non si esplorerebbero tutti i punti di suddivisione, ma solo quelli relativi alle variabili scelte casualmente.

Random forest (Foreste Casuali)

Bagging e boosting fanno variare ad ogni iterazione le unità statistiche che entrano in gioco nel modello, usando tutte le variabili esplicative disponibili. Le foreste casuali invece considerano diversi sottoinsiemi di variabili esplicative. Tali tecniche sono tipiche degli alberi di classificazione, ma usate anche per altri tipi di classificatori.

La procedura consiste nel selezionare in modo casuale, ad ogni nodo dell'albero, un piccolo gruppo di variabili esplicative che verranno ispezionate per individuare il punto di suddivisione ottimale.

In questo modo non si esplorerebbero tutti i punti di suddivisione, ma solo quelli relativi alle variabili scelte casualmente.

Random forest (Foreste Casuali)

Bagging e boosting fanno variare ad ogni iterazione le unità statistiche che entrano in gioco nel modello, usando tutte le variabili esplicative disponibili. Le foreste casuali invece considerano diversi sottoinsiemi di variabili esplicative. Tali tecniche sono tipiche degli alberi di classificazione, ma usate anche per altri tipi di classificatori.

La procedura consiste nel selezionare in modo casuale, ad ogni nodo dell'albero, un piccolo gruppo di variabili esplicative che verranno ispezionate per individuare il punto di suddivisione ottimale.

In questo modo non si esplorerebbero tutti i punti di suddivisione, ma solo quelli relativi alle variabili scelte casualmente.

L'albero viene fatto crescere fino alla sua massima dimensione e non viene però potato. Sarà infatti la combinazione di diversi alberi che permetterà di evitare i problemi di sovra-adattamento.

Il numero di variabili da selezionare é un parametro di regolazione da determinare e in genere viene mantenuto costante su tutti i nodi.

Spesso questo viene scelto considerando foreste costruite con un numero diverso di variabili esplicative e determinando quel valore che minimizza l'errore su un insieme di verifica.

Un altro parametro di regolazione da determinare é il numero di alberi che costituiscono la foresta. Si può mostrare che l'errore globale converge a una soglia inferiore al crescere del numero di alberi.

L'albero viene fatto crescere fino alla sua massima dimensione e non viene però potato. Sarà infatti la combinazione di diversi alberi che permetterà di evitare i problemi di sovra-adattamento.

Il numero di variabili da selezionare é un parametro di regolazione da determinare e in genere viene mantenuto costante su tutti i nodi.

Spesso questo viene scelto considerando foreste costruite con un numero diverso di variabili esplicative e determinando quel valore che minimizza l'errore su un insieme di verifica.

Un altro parametro di regolazione da determinare é il numero di alberi che costituiscono la foresta. Si può mostrare che l'errore globale converge a una soglia inferiore al crescere del numero di alberi.

L'albero viene fatto crescere fino alla sua massima dimensione e non viene però potato. Sarà infatti la combinazione di diversi alberi che permetterà di evitare i problemi di sovra-adattamento.

Il numero di variabili da selezionare é un parametro di regolazione da determinare e in genere viene mantenuto costante su tutti i nodi.

Spesso questo viene scelto considerando foreste costruite con un numero diverso di variabili esplicative e determinando quel valore che minimizza l'errore su un insieme di verifica.

Un altro parametro di regolazione da determinare é il numero di alberi che costituiscono la foresta. Si può mostrare che l'errore globale converge a una soglia inferiore al crescere del numero di alberi.

L'albero viene fatto crescere fino alla sua massima dimensione e non viene però potato. Sarà infatti la combinazione di diversi alberi che permetterà di evitare i problemi di sovra-adattamento.

Il numero di variabili da selezionare é un parametro di regolazione da determinare e in genere viene mantenuto costante su tutti i nodi.

Spesso questo viene scelto considerando foreste costruite con un numero diverso di variabili esplicative e determinando quel valore che minimizza l'errore su un insieme di verifica.

Un altro parametro di regolazione da determinare é il numero di alberi che costituiscono la foresta. Si può mostrare che l'errore globale converge a una soglia inferiore al crescere del numero di alberi.

In genere nella costruzione di una foresta si associa alla selezione casuale delle variabili anche una procedura di bagging.

Praticamente ciascun albero viene fatto crescere su un campione bootstrap e utilizzando per ciascun nodo un numero diverso di variabili.

Tale combinazione permette ancora una volta di utilizzare la procedura out-of-bag, utilizzata anche per valutare l'importanza di ogni variabile esplicativa.

In genere nella costruzione di una foresta si associa alla selezione casuale delle variabili anche una procedura di bagging.

Praticamente ciascun albero viene fatto crescere su un campione bootstrap e utilizzando per ciascun nodo un numero diverso di variabili.

Tale combinazione permette ancora una volta di utilizzare la procedura out-of-bag, utilizzata anche per valutare l'importanza di ogni variabile esplicativa.

Procedura:

- Si costruiscono gli alberi e si effettua la previsione sull'insieme out-of-bag e sullo stesso insieme con i valori della j -esima variabile esplicativa.
- Si misura la differenza tra l'errore di previsione nei due casi e al termine della procedura si considera la media delle differenze tra i vari alberi divisa per l'errore standard.
- Questo valore fornisce una misura di quanto la variabile influisce sulle previsioni.

L'accuratezza della previsione é simile a quella ottenuta col boosting ma soprattutto piú veloce perché ogni albero si basa su poche variabili e quindi la complessità computazionale é inferiore.

É inoltre facile ottenere un algoritmo di stima che sfrutta il calcolo parallelo che accelera ulteriormente la procedura.

Procedura:

- Si costruiscono gli alberi e si effettua la previsione sull'insieme out-of-bag e sullo stesso insieme con i valori della j -esima variabile esplicativa.
- Si misura la differenza tra l'errore di previsione nei due casi e al termine della procedura si considera la media delle differenze tra i vari alberi divisa per l'errore standard.
- Questo valore fornisce una misura di quanto la variabile influisce sulle previsioni.

L'accuratezza della previsione é simile a quella ottenuta col boosting ma soprattutto piú veloce perché ogni albero si basa su poche variabili e quindi la complessità computazionale é inferiore.

É inoltre facile ottenere un algoritmo di stima che sfrutta il calcolo parallelo che accelera ulteriormente la procedura.

Procedura:

- Si costruiscono gli alberi e si effettua la previsione sull'insieme out-of-bag e sullo stesso insieme con i valori della j -esima variabile esplicativa.
- Si misura la differenza tra l'errore di previsione nei due casi e al termine della procedura si considera la media delle differenze tra i vari alberi divisa per l'errore standard.
- Questo valore fornisce una misura di quanto la variabile influisce sulle previsioni.

L'accuratezza della previsione é simile a quella ottenuta col boosting ma soprattutto piú veloce perché ogni albero si basa su poche variabili e quindi la complessità computazionale é inferiore.

É inoltre facile ottenere un algoritmo di stima che sfrutta il calcolo parallelo che accelera ulteriormente la procedura.