

Alberi di Decisione e Potatura

Mihail Teodor Gurzu
5940299

13 aprile 2020

1 Intento

In questo elaborato si implementa in Python un algoritmo per l'apprendimento di Alberi di Decisione utilizzando *l'entropia* come misura dell'impurità, e successivamente un algoritmo per la potatura basato sull'errore sul *validation set*.

Il codice viene poi applicato a tre datasets e vengono ricavati e mostrati i risultati ottenuti (numero dei nodi, precisione sul test set) prima e dopo la potatura.

2 Implementazione

2.1 Alberi di Decisione

Gli Alberi di Decisione ricoprono un ruolo molto importante nell'ambito del *Machine Learning* e vengono utilizzati per costruire dei modelli per vari tipi di problemi come la *regressione* e la *classificazione*.

L'apprendimento di Alberi di Decisione consiste nel costruire un albero di decisione partendo da un dataset di dati in modo tale da poter predire certe caratteristiche sui dati futuri, con l'obiettivo di massimizzare la precisione e minimizzare lo spazio occupato in memoria.

Infatti uno dei principali problemi di questa struttura è che essa aumenta in modo esponenziale con l'aumento degli attributi e dei relativi possibili valori. Bisogna quindi adottare una certa tecnica nella costruzione di tali strutture e quella che è stata adottata in questo elaborato si basa sul concetto di *information gain* utilizzando *l'entropia* come misura delle impurità.

```
1 def ID3(examples, default):
2     ...
3     else:
4         best = chooseAttr(allAttr, examples, allAttr[last])
5         t = Node(best, {})
6     ...
7     return t
```

Come si vede ID3 usa la funzione *chooseAttr* per scegliere il miglior attributo in base al miglior valore dell' *information gain* che viene mostrato in seguito.

ID3 prende in ingresso l'insieme degli esempi ed un valore di default per l'attributo target e ritorna l'albero come istanza di *Node*.

```
1 class Node:
2     """ seen and pruned attributes are used in the pruning routine """
3     def __init__(self, label, children=None, seen=False, pruned=False):
4         self.label = label
5         self.children = {} if children == None else children
6         self.seen = seen
7         self.pruned = pruned
8     ...
```

2.1.1 Entropia

L'*entropia* definisce una misura comunemente usata nella teoria dell'informazione che caratterizza l'impurità di una collezione di esempi.

Data per esempio una certa collezione S dove l'attributo target può prendere c differenti valori, l'*entropia* di S relativa a questa classificazione è definita come:

$$Entropy(S) = \sum_{i=1}^c -p_i \log_2 p_i$$

dove p_i è la proporzione di S appartenente alla classe i .

```
1 def entropy(examples, attributes, target):
2     ...
3     i = 0
4     for entry in attributes:
5         if(target == entry):
6             break
7         i += 1
8     a = attributes[i]
9     for entry in examples:
10        if((entry[a]) in valfreq):
11            valfreq[entry[a]] += 1.0
12        else:
13            valfreq[entry[a]] = 1.0
14    for freq in valfreq.values():
15        entrop += (-freq/len(examples)) * math.log(freq / len(examples), 2)
16    return entrop
```

2.1.2 Information Gain

L'*Information Gain* di un particolare attributo è la riduzione attesa dell'entropia causata dal partizionare gli esempi relativamente a quel attributo.

Precisamente, l'*information gain* di un attributo A relativo ad una collezione di esempi S è definita come:

$$Gain(S,A) = Entropy(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} Entropy(S_v)$$

dove $Values(A)$ è l'insieme di tutti i possibili valori per l'attributo A e S_v è il sottoinsieme di S per il quale l'attributo A ha valore v .

```
1 def gain(examples, allAttr, a, target):
2     valfreq = {}
3     subs_ent = 0.0
4     for entry in examples:
5         if((entry[a]) in valfreq):
6             valfreq[entry[a]] += 1.0
7         else:
8             valfreq[entry[a]] = 1.0
9     for val in valfreq.keys():
10        valprob = valfreq[val] / sum(valfreq.values())
11        subs = [entry for entry in examples if entry[a] == val]
12        subs_ent += valprob * entropy(subs, allAttr, target)
13    return(entropy(examples, allAttr, target) - subs_ent)
```

2.2 Potatura

Un' altro problema caratteristico degli Alberi di Decisione riguarda l'*overfitting* dei dati; ciò avviene quando l'Albero di Decisione si specializza troppo sul *training set* a scapito della generalizzazione sul problema in questione. Ciò si manifesta in una precisione altissima sul

training set che cala drasticamente quando consideriamo il *test set*, cioè dati nuovi, mai incontrati prima.

Per ovviare a questo problema, applichiamo la '*potatura*' all'albero con l'intento di despecializzarlo riguardo al *training set* aumentandone la generalizzazione e di conseguenza la precisione sul *test set*.

La tecnica utilizzata in questo elaborato si chiama *REDUCED ERROR PRUNNING* che descriviamo brevemente di seguito.

2.2.1 Reduced Error Pruning

Per applicare questa tecnica, dobbiamo innanzitutto suddividere il nostro *training set* in due parti: il *grow set* utilizzato per creare l'albero e il *validation set*, usato per decidere se potare un sottoalbero avente radice un certo nodo oppure no. Questa tecnica richiede avere a disposizione un numero abbastanza elevato di esempi.

La potatura, quindi consiste nello considerare, uno per uno, ciascun nodo non foglia e sostituire al relativo sottoalbero la più comune classificazione sul *grow set* relativo a quel nodo. Il nodo diventa quindi un candidato per la potatura se la precisione sul *validation set* dell'albero potato non è inferiore a quella dell'albero originale sul *validation set*. Dopo aver applicato questa procedura a tutti i nodi non foglia, si procede con la potatura per il nodo che massimizza l'aumento della precisione sul *validation set*. A questo punto, l'intero procedimento viene ripetuto finché c'è un qualche nodo che si candida per la potatura.

Adesso, prima di far vedere l'implementazione precisiamo che la classe *Data Manager* si occupa di suddividere il dataset nelle tre parti che ci servono così composte:

$$growset(70\%) + validationset(30\%) = trainingset$$

$$trainingset(70\%) + testset(30\%) = dataset$$

Inoltre *Data Manager* si occupa anche dell'introdurre rumore nei dati del *grow set* in base ad una certa percentuale, attraverso il metodo *introduce_noise()*

Per quanto riguarda l'implementazione della potatura, ne sono state fatte due, dove la prima è quella descritta sopra e che come si può intuire, quando l'albero diventa grande ha un tempo di esecuzione molto grande. La seconda implementazione si differenzia dalla prima col fatto che la procedura pota il primo nodo che incontra e che aumenta la precisione sul validation set (partendo dai nodi a profondità maggiore) e poi ricomincia da capo la ricerca dei nodi da potare, terminando quando non ne trova alcuno che migliori la precisione sul validation set con la conseguenza che il tempo di esecuzione viene così estremamente ridotto. La differenza principale nel codice sta nel metodo *managePrune()*.

Per primo facciamo vedere l'implementazione originale:

```

1 def _manage_prune(node, validation_set, grow_set):
2     pruned_nodes = False
3     iteration = True
4     tmp_best_tree = copy.deepcopy(node)
5     best_tree = node
6     best_performance = test(node, validation_set)
7     while iteration:
8         iteration = False
9         new_tree = copy.deepcopy(tmp_best_tree)
10        iteration = _prune(new_tree, tmp_best_tree, grow_set, iteration)
11        if iteration:
12            new_performance = test(new_tree, validation_set)
13            if new_performance >= best_performance:
14                best_tree = new_tree
15                best_performance = new_performance
16                pruned_nodes = True
17    reset_tree(best_tree)

```

```

18     if pruned_nodes:
19         best_tree, best_performance, pruned_nodes = _manage_prune(best_tree,
20             validation_set, grow_set)
21     return best_tree, best_performance, pruned_nodes

```

Come si può notare dal codice, la procedura nel caso peggiore ha un andamento asintotico pari a $O(n^2)$ dove n è il numero dei nodi dell'albero.

```

1 def manage_prune(node, validation_set, grow_set):
2     nr_pruned_nodes = 0
3     iteration = True
4     best_tree = node
5     best_performance = test(node, validation_set)
6     while iteration:
7         iteration = False
8         new_tree = copy.deepcopy(best_tree)
9         iteration = prune(new_tree, best_tree, grow_set, iteration)
10        new_performance = test(new_tree, validation_set)
11        if iteration and new_performance >= best_performance:
12            best_tree = new_tree
13            best_performance = new_performance
14            nr_pruned_nodes += 1
15    return best_tree, best_performance, nr_pruned_nodes

```

Questa versione invece ha un andamento asintotico pari a $\Theta(n)$ con n il numero dei nodi.

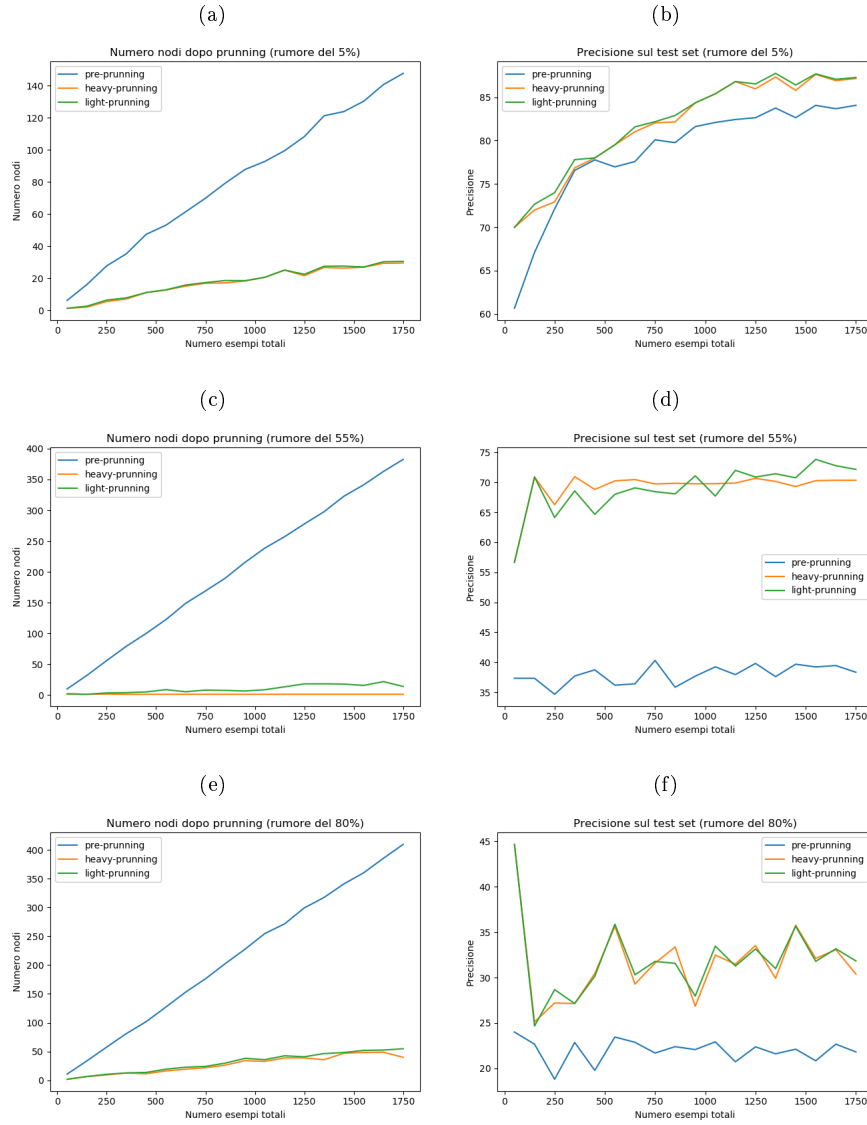
Entrambe le implementazioni fanno uso di *prune()*:

```

1 def prune(node, original_node, examples, pruned):
2
3     if len(node.children) == 0:
4         return pruned
5     if node.pruned == True:
6         return pruned
7
8     if len(examples) == 0:
9         node.children = {}
10        node.label = mode(examples, T)
11        return pruned
12
13    subsets = subset(examples, node.label)
14
15    for child in node.children.keys():
16        if pruned == True:
17            node.seen = False
18            return pruned
19        if node.children[child].children != {} and node.children[child].pruned !=
20            True:
21            node.seen = True
22            if child in subsets.keys():
23                pruned = prune(node.children[child], original_node.children[child],
24                    subsets[child], pruned)
25            else:
26                pruned = prune(node.children[child], original_node.children[child],
27                    {}, pruned)
28        if node.seen:
29            node.seen = False
30            pruned = True
31            return pruned
32
33    node.children = {}
34    ex_label = node.label
35    node.label = mode(examples, T)
36    node.pruned = True
37    original_node.pruned = True
38    pruned = True
39    return pruned

```

Figura 1: heavy pruning è l'implementazione originale, light pruning quella asintoticamente più veloce. I grafici si riferiscono al dataset *cardata*.



37

In Figura 1 è mostrata la differenza tra le due procedure riguardo il numero dei nodi dopo la potatura e la precisione sul test set (i risultati sono stati ottenuti applicando le due procedure allo stesso albero di decisione):

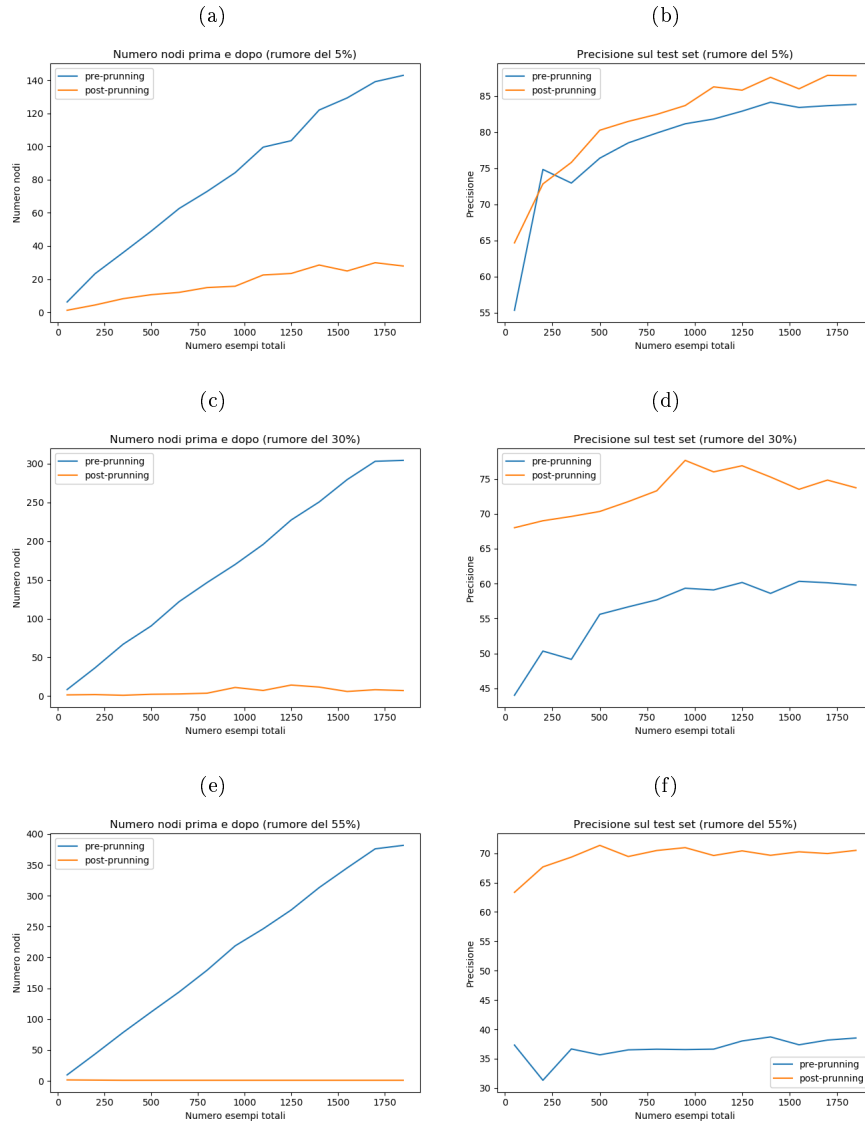
Come si può vedere, le due procedure sono pressoché equivalenti nei risultati ottenuti.

Dato che le differenze tra le due procedure è piccola e dato che l'andamento asintotico dell'implementazione *'light'* è molto minore dell'altra, i risultati degli esperimenti sul dataset *nursery_admission_school* sono stati ottenuti attraverso l'applicazione della procedura *'light'*.

3 Risultati esperimenti

3.1 Cardata

Figura 2: Risultati sul dataset *cardata*.

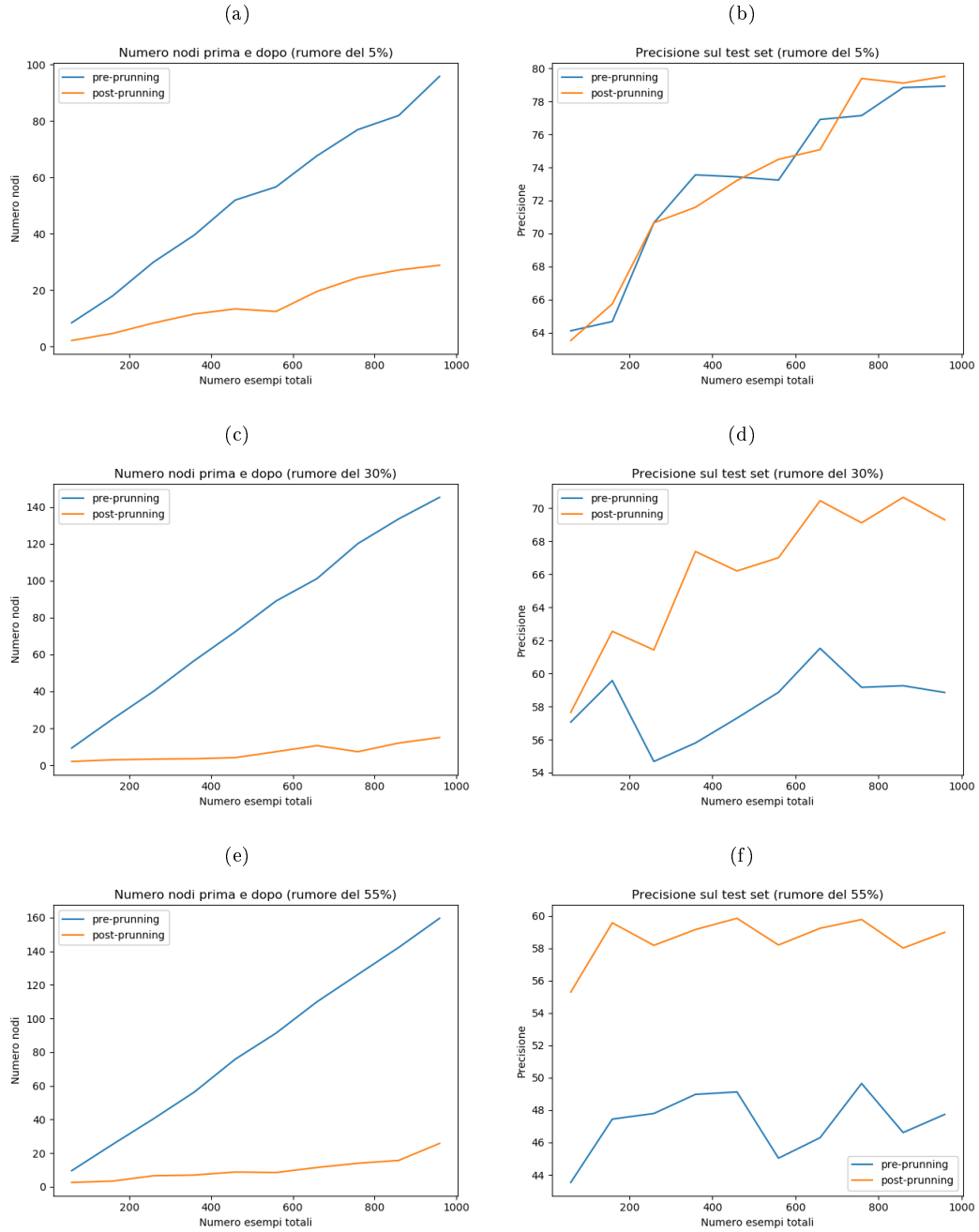


Come si evince dalla Figura 2, la potatura riduce al minimo il numero dei nodi aumentando la precisione sul test set.

3.2 tic-tac-toe

Nella Figura 3 sono mostrati i risultati sul dataset *tic-tac-toe*.

Figura 3: Risultati sul dataset *tic-tac-toe*.

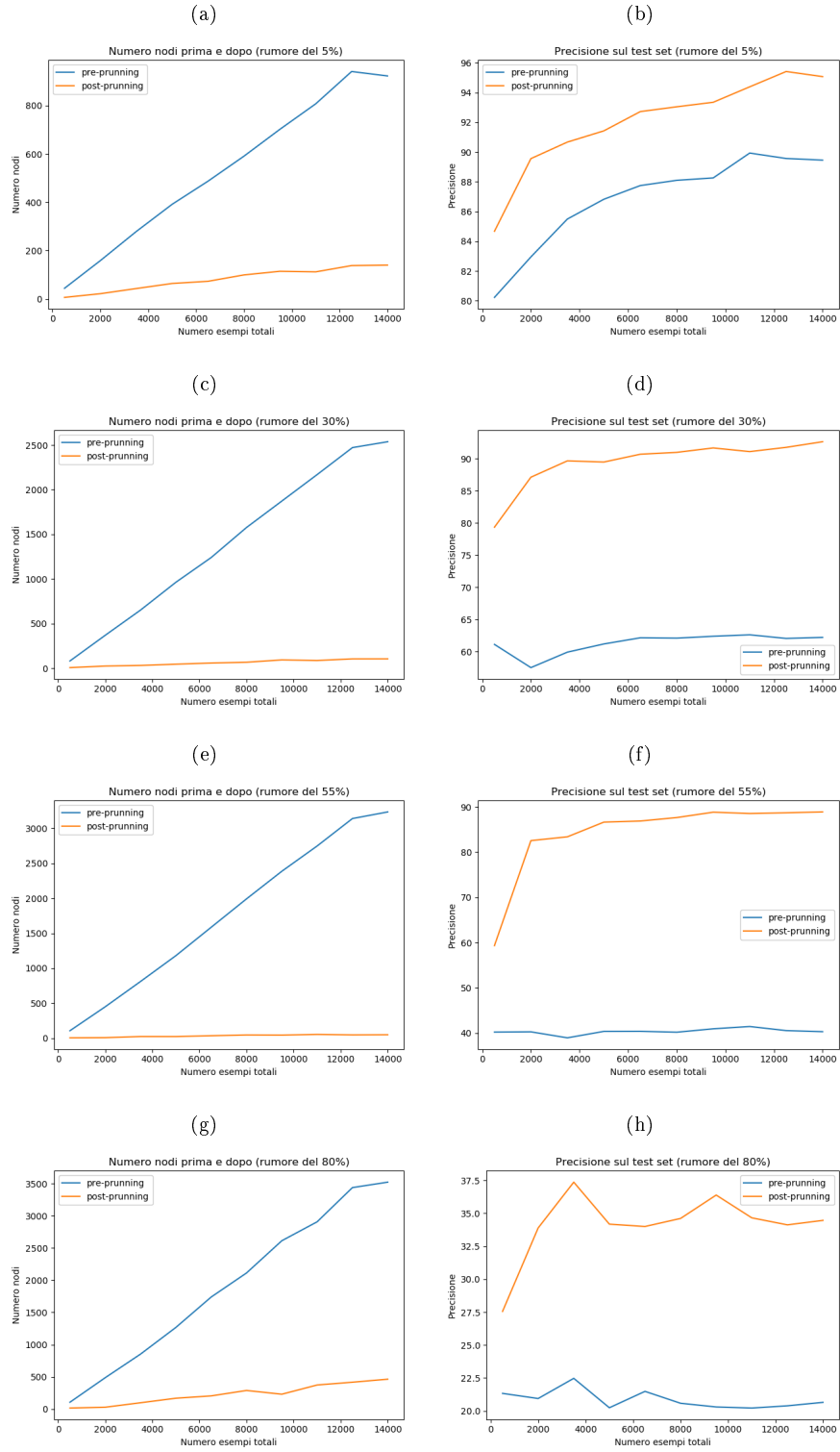


3.3 Nursery School Admission

Dalla Figura 4 notiamo che dopo la potatura il numero dei nodi viene ridotto al minimo, ed anche la precisione sul test set viene sostanzialmente aumentata.

Dai risultati di tutti i dataset si evince che la precisione massima si ha per un insieme di esempi elevato. Più sono gli esempi del dataset iniziale e di conseguenza più sono i nodi

Figura 4: Risultati sul dataset *nursery school admission*.



dell'albero risultante, maggiore risulta la precisione sul test set.

4 Come usare il codice

```
1 def main(dataset):
2     '''
3     elaborate() function generates a decision tree based on file dataset and then
4     pruns it plotting the results (True = light pruning; False = heavy pruning)
5     compare() function compares the two implementations of pruning (light and
6     heavy)
7     '''
8     elaborate(dataset, 5, 50, 3, True)
9     compare(dataset, 5, 50, 3)
10
11 main("cardata")
```

Per usare il codice, bisogna passare alla funzione `main()` il dataset su cui si desidera operare. Il dataset deve avere come prima riga gli attributi e l'attributo target deve occupare l'ultima colonna.

```
1 def elaborate(dataset, initial_error_rate, initial_examples_nr,
2               number_of_iterations, light)
3 ...
4 def compare(dataset, initial_error_rate, initial_examples_nr,
5             number_of_iterations)
```

I parametri delle due funzioni rappresentano:

- *file* - denota il dataset su cui operare;
- *initial_error_rate* - denota la percentuale dell'errore da introdurre nel grow set (la percentuale aumenta dello 25% ad ogni iterazione del ciclo principale delle rispettive funzioni);
- *initial_examples_nr* - denota il numero iniziale di esempi presi casualmente dall'intero dataset sul quale costruire i tre dataset (grow set, validation set e test set) usati dall'elaborazione. Il numero di esempi viene aumentato ad ogni iterazione interna di un numero prefissato;
- *number_of_iterations* - denota il numero delle iterazioni su cui fare la media (i grafici mostrati qui sono ottenuti come media su 10 iterazioni);
- *light* - se **True** verrà usata l'implementazione *'light'* della procedura di potatura; se **False** verrà usata quella *'heavy'*;