Il software Weka Soluzioni degli esercizi

Prof. Matteo Golfarelli

Alma Mater Studiorum - Università di Bologna

Il data set Iris

- Il data set Iris modella le caratteristiche di una famiglia di piante
 - √ 150 istanze
 - ✓ Nessun dato mancante

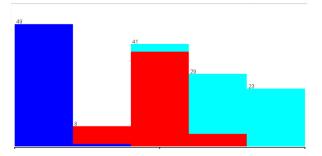
Attributo	Descrizione
SepalLength	Lunghezza del sepalo
SepalWidth	Larghezza del sepalo
PetalLength	Lunghezza del petalo
PetalWidth	Larghezza del petalo
Class	Sotto famiglia {setosa, virginica, versicolor}

Pre processing Iris dataset

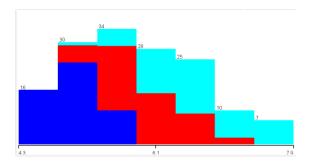
- Caricare il file
- 2. Effettuare un'analisi manuale dei dati mediante visualizzazione
 - ✓ Istogrammi attributo classe (da pagina Preprocess)
 - Quale attributo sembra essere maggiormente discriminante? Perché?
 - Quale attributo sembra essere meno discriminante? Perché?
 - ✓ Distribuzioni attributo attributo classe (da pagina Visualize)
 - Quale coppia di attributi sembra essere maggiormente discriminante? Perché?

Pre processing Iris dataset

- 1. Caricare il file
- 2. Effettuare un'analisi manuale dei dati mediante visualizzazione
 - ✓ Istogrammi attributo classe (da pagina Preprocess)
 - Quale attributo sembra essere maggiormente discriminante? Perché?
 - Quale attributo sembra essere meno discriminante? Perché?



L'attributo maggiormente discriminante sembra essere *PetalWidth* perché la sua misura determina una maggiore separazione tra gli elementi delle sottofamiglie. In particolare, i fiori di Setosa (classe blu) sono facilmente distinguibili



L'attributo meno discriminante sembra essere *SepalLength* perché i fiori delle diverse sottofamiglie presentano valori molto simili (sovrapposti) per la sua misura

Pre processing Iris dataset

- 1. Caricare il file
- 2. Effettuare un'analisi manuale dei dati mediante visualizzazione
 - ✓ Istogrammi attributo classe (da pagina Preprocess)
 - Quale attributo sembra essere maggiormente discriminante? Perché?
 - Quale attributo sembra essere meno discriminante? Perché?
 - ✓ Distribuzioni attributo attributo classe (da pagina Visualize)
 - Quale coppia di attributi sembra essere maggiormente discriminante? Perché?



PetalWidth+PetalLength istanze di classi diverse ben sepatate = alta capacità discriminante

SepalWidth+SepalLength istanze di classi diverse sovrapposte = bassa capacità discriminante

Classificazione Iris dataset

- Classificare il dataset utilizzando l'algoritmo trees->J48 e commentare il risultato
 - ✓ Valutare il risultato con Use training set / Cross-validation 10 folds / Percentage split
 - ✓ Rappresentare e discutere i decision boundary

Classificazione Iris dataset

- Classificare il dataset utilizzando l'algoritmo trees->J48 e commentare il risultato
 - ✓ Valutare il risultato con Use training set / Cross-validation 10 folds / Percentage split
 - Il dataset è semplice da classificare e le prestazioni sono sempre molto buone
 - Le istanze di Setosa sono più semplici da classificare (anche visivamente)

J48 – Training set 98%

```
a b c <-- classified as
50 0 0 | a = Iris-setosa
0 49 1 | b = Iris-versicolor
0 2 48 | c = Iris-virginica
```

J48 – Cross val 96%

```
a b c <-- classified as
49 1 0 | a = Iris-setosa
0 47 3 | b = Iris-versicolor
0 2 48 | c = Iris-virginica
```

J48 – Split 98%

```
a b c <-- classified as

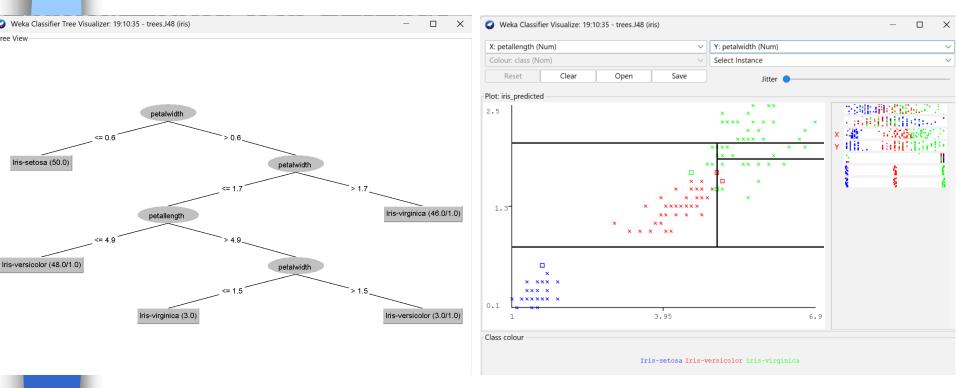
15 0 0 | a = Iris-setosa

0 19 0 | b = Iris-versicolor

0 2 15 | c = Iris-virginica
```

Classificazione Iris dataset

- Classificare il dataset utilizzando l'algoritmo trees->J48 e commentare il risultato
 - ✓ Rappresentare e discutere i decision boundary
 - Sono utilizzati due soli attributi: le informazioni portate dalle misure del Sepalo sono molto limitate ai fini della classificazione e l'algoritmo ha scelto di non utilizzarle

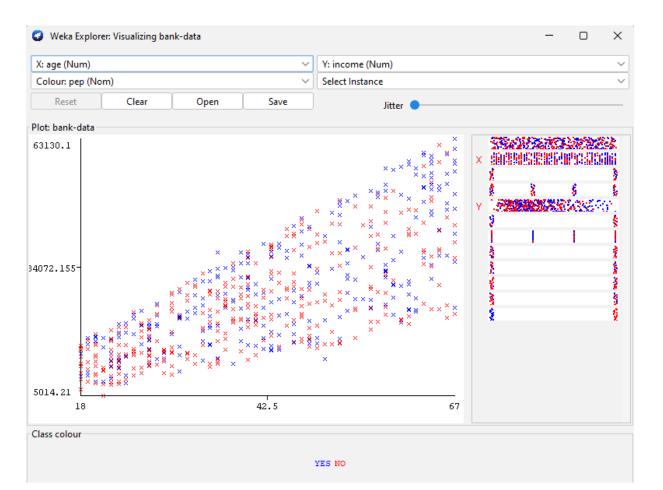


- II date set bank-data.csv
 - √ 600 istanze
 - ✓ Nessun dato missing

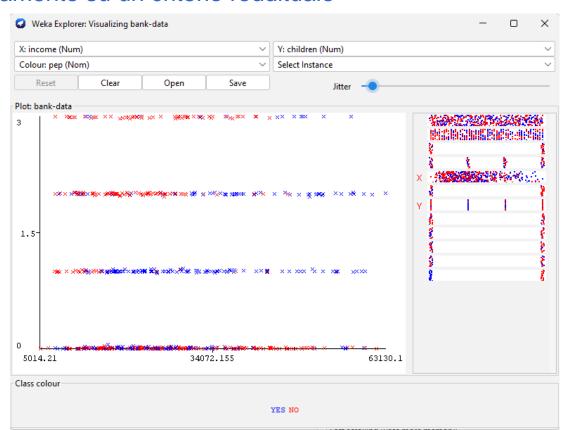
Attributo	Descrizione
ld	Identificatore unico
Age	età del cliente in anni (numeric)
Sex	MALE / FEMALE
Region	inner_city/rural/suburban/town
Income	reddito del cliente (numeric)
children	sposato? (YES/NO)
car	possiede un'automobile? (YES/NO)
save_acct	ha un conto di risparmio? (YES/NO)
current_acct	ha un conto corrente? (YES/NO)
Mortgage	ha un mutuo (YES/NO)
pep	ha acquistato un PEP (Personal Equity Plan) dopo l'ultimo invio postale? (YES/NO)

- Caricare il file e salvarlo in formato ARFF con il nome "bank data.arff"
- 2. Effettuare un'analisi manuale dei dati mediante visualizzazione
 - ✓ Istogrammi attributo attributo (da pagina Visualize)
 - ✓ Distribuzioni attributo attributo classe (da pagina Visualize)
 - Commentare la distribuzione Income –Children PEP
- 3. Eliminare l'attributo ID e salvare nuovamente con il nome "bank data.arff"
- 4. Discretizzare l'attributo AGE (Equal- frequency 10 bin) e Children (Manualmente) e salvare il dataset con il nome "bank data1.arff"

- 1. Discutere le relazioni e la capacità discriminante degli attributi
 - Quando si inizia il percorso lavorativo i redditi sono più bassi e più simili, con il procedere degli anni i redditi si differenziano anche in funzione delle carriere



- 1. Discutere le relazioni e la capacità discriminante degli attributi
 - ✓ La propensione all'acquisto del PEP (punti blu) è strettamente legata al reddito in funzione del numero (ossia al costo) dei figli.
 - ✓ Per chi non ha figli la correlazione propensione d'acquisto reddito è non si evidenzia: ossia la propensione all'acquisto non è basata strettamente su un criterio reddituale



- Utilizzare i seguenti algoritmi per eseguire la classificazione e commentare il risultato
 - ✓ Valutare il risultato con Cross-validation 10 folds / Use training set / Percentage split
 - ✓ Utilizzare sia il data set discretizzato (bank data1.arff), sia quello non discretizzato (bank data.arff)
 - 1. J48
 - Rappresentare e discutere i decision boundary
 - 2. J48 senza post-pruning (unpruned = True)
 - Visualizzare l'albero di decisione
 - 3. JRip
 - 4. IBk con k=1 e k=5 sul data set non discretizzato (BankData.arff)
 - Dopo aver normalizzato i rimanenti attributi numerici
 - 5. Ripetere la classificazione utilizzando BankData1.arff dopo aver discretizzato anche l'attributo income (equal frequency 10 bins)
 - Salvare il training set nel file BankData2.arff

Bank data1.arff: Discretizzato

J48 – Cross val 91.3%

J48 – Split 88.2%

Bank data.arff: Non discretizzato

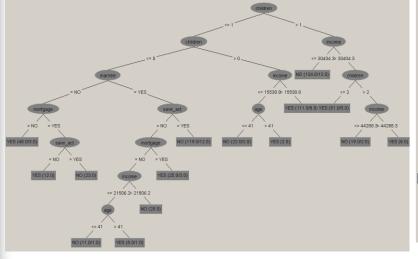
J48 – Training set 92.3%

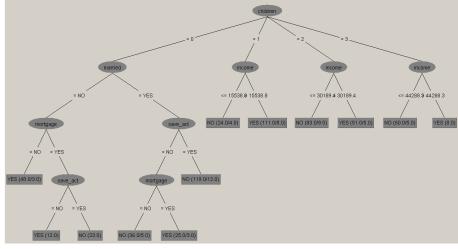
J48 – Cross val 89.8%

J48 - Split 86.8%

Bank data.arff: Non discretizzato + Split

Bank data1.arff: Discretizzato + Split

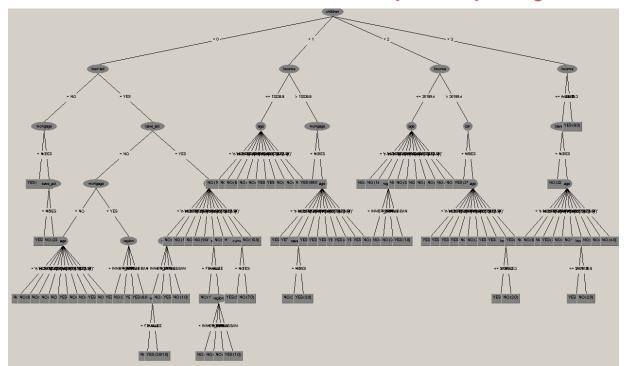




Bank data1.arff: Discretizzato + No pruning

Migliora la classificazione sul training set, peggiora negli altri casi Fatica a generalizzare! → Overfitting

Bank data1.arff: Discretizzato + Split + No pruning



Bank data1.arff: Discretizzato

JRip – Cross val 91.0%

JRip - Split 86.8%

Bank data.arff: Non discretizzato

JRip- Training set 91.2%

JRip- Cross val 90.8%

JRip-Split 85.8%

- Regole e alberi forniscono classificazioni molto simili
- La discretizzazione sembra migliorare lievemente le prestazioni anche se fissando arbitrariamente le soglie non consente agli algoritmi di learning di decidere i punti di split
 - ✓ Motivazione: alberi e regole più semplici? Nella versione del data set non discretizzata l'attributo income viene utilizzato più frequentemente

Bank data.arff2: Normalizzato + Discretizzato

IB1– Training set 100%

IB1– Cross val 74.7%

IB1- Split 70.6%

Bank data.arff: Normalizzato + Non discretizzato

IB1– Training set 100%

IB1– Cross val 64.2%

IB1- Split 61.3%

Bank data.arff2: Normalizzato + Discretizzato

IB5 – Training set 82.3% IB5– Cross val 76.3%

IB5- Split 75.5%

Bank data.arff: Normalizzato + Non discretizzato

IB5 – Training set 78.2% IB5 – Cross val 65.2%

IB5 – Split 60.3%

Bank data: alcune considerazioni

- Le tecniche k-nearest neighbor non sembrano fornire buoni risultati
 - ✓ Rumore nei dati? (l'utilizzo di un k=5 aumenta di poco l'efficacia)
 - ✓ Training set limitato?
 - ✓ Funzioni distanza non appropriate? (l'algoritmo non riesce a sfruttare gli attributi numerici)
 - ✓ Presenza di attributi ridondanti o non rilevanti?
- E' possibile fare una verifica...
 - ✓ Selezionando un sottoinsieme di attributi
 - ✓ Eseguendo su di esso la classificazione
 - Sul training set che ha dato i risultati migliori bank data1.arff (senza campo id con children e age discretizzati) si esegue la selezione degli attributi (CfsSubsetEval) che restituisce:

income + married + children

- Si esegue la classificazione sulla versione normalizzata di questi attributi
- In alternativa si possono utilizzare i primi attributi scelti dall'albero decisionale che ha fornito i risultati migliori (ossia quello basato bank data1.arff):

children+income+married+mortgage+save act

bank data1.arff: income + married + children

IB5- Split 88.9%

bank data1.arff: children + income + married + mortgage + save act

IB5– Training set 90.3% IB5– Cross val 87.0%

IB5- Split 84.8%

- Il date set bank-data.csv riporta dati del censimento USA (http://cps.ipums.org/)
 - √ 1000 istanze per il training
 - √ 31561 istanze per la validazione

Attributo	Descrizione
age	Età in anni
workclass	Classe di lavoro
fnlwgt	"Final sampling weight" peso dell'istanza (campione) rispetto alla popolazione
education	Titolo ottenuto
education-num	Numero di anni di studio
marital-status	Stato civile
occupation	Occupazione
relationship	Tipo di relazione con il capo famiglia
race	Razza
sex	Sesso
capital-gain	Utili da capitali (plus valenza)
capital-loss	Perdite da capitali (minus valenza)
hours-per-week	Ore di lavoro settimanali
native-country	Nazionalità
Total Income	L'individuo guadagna più o meno di 50K\$

- Obiettivo dello studio è trovare un modello che permetta di predire quali persone guadagnano più di 50K€
 - ✓ Ricerca di frodi fiscali
- Si proceda utilizzando la metodologia CRISP-DM
 - 1. Comprensione del dominio applicativo
 - 2. Comprensione dei dati
 - 3. Preparazione dei dati
 - 4. Creazione del modello
 - 5. Valutazione del modello e dei risultati
 - 6. Deployment

- C'è una forte correlazione tra i campi education ed education num
- Il campo fnlwgt sembra poco rilevante ai fini della determinazione della classe
- L'attributo native country è fortemente sbilanciato sul valore United States
 - ✓ Le selezioni fatte utilizzando questo attributo saranno sempre poco significative
- Alcuni campi presentano valori nulli ('?')
 - ✓ E' necessario valutare se calcolarli sostituendoli con media e moda
- Tra i campi che "a vista" meglio si prestano a fornire informazioni utili alla classificazione
 - ✓ Relationship, Race, Sex, Education Education num
 - ✓ Il profilo di chi guadagna >50k\$ è maschio, bianco, sposato con una educazione di livello superiore.
 - ✓ Dalla prima analisi non è possibile evincere se tra questi attributi esistano delle correlazioni e quindi debbano essere usati in alternativa

- Si comparino le performance degli algoritmi studiati (J48, JRip e IBK) utilizzando come training set i file
 - ✓ AdultTraining.arff
 - ✓ AdultTrainingSmall.arff
- In entrambi i casi si utilizzi come modalità di validazione il test set
 - ✓ AdultTest.arff
- Si verifichi in seguito come cambiano i risultati ottenuti dai due data set utilizzando le tecniche
 - ✓ AdaBoost
 - ✓ Bagging
- Si commentino i risultati

L'utilità di calcolare i dati mancanti è verificata applicando gli algoritmi ai dati grezzi e ai dati modificati tramite il filtro ReplaceMissingValues

CensusTraining.arff

JRip- Supplied TS 82.6%

IB5- Supplied TS 80.7%

```
a b <-- classified as
3575 4034 | b = >50K
```

CensusTrainingNoMissing.arff

JRip- Supplied TS 83.0%

IB5- Supplied TS 80.1%

Si decide di mantenere i dati originali

Le performance degli algoritmi sul training set ridotto sono:

CensusTrainingSmall.arff

- Riducendo il numero di attributi considerati sulla base del filtro CfsSubsetEval la performance dell'algoritmo di k-nearest neighbor migliora leggermente
 - ✓ Age + education-num + marital-status + relationship + capital gain + capital loss + hours-per-week

CensusTraining.arff

Cosa succede discretizzando gli attributi numerici?

Se si introducono tecniche di boosting e bagging

CensusTraining.arff + bagging

CensusTraining.arff + AdaBoost

Il classificatore migliore risulta essere J48 utilizzato con la tecnica di bagging sul training set con valori mancanti

- Una ditta di collocamento possiede una banca dati contenente le informazioni (classe esclusa) relative a 50000 infermiere. Si vogliono contattare tramite posta tutte le infermiere appartenenti alla classe "priority". Tali infermiere saranno tutte inviate ad un insieme di ospedali che hanno fatto richiesta di nuove infermiere.
- L'operazione di invio delle lettere ha un costo fisso di 10,000 € e un costo individuale (per ogni infermiera contattata) pari a 5 €.
- Gli ospedali prenderanno in prova le infermiere selezionate dalla ditta e alla fine del periodo di prova pagheranno alla ditta 10 € per ogni infermiera che risulterà essere effettivamente una infermiera classificabile come appartenente alla classe "priority". Decidere quale modello, tra quelli studiati permette di ottenere potenzialmente il profitto maggiore e quante infermiere (in percentuale) devono essere contattate.
- Verificare quale delle tecniche di classificazione studiate fornisce il modello che fornisce il lift maggiore
 - ✓ Dataset di training: nursery.arff
 - ✓ Tecnica di validazione 10 folds cross-validation

Il lift è un indicatore che permette di compare più modelli di classificazione favorendo quelli che permettono di individuare un campione distorto della popolazione che massimizzi la probabilità di trovare istanze della classe desiderata Ci

$$Lift(ModelloX) = \frac{P(C_i \mid Campione)}{P(C_i \mid Popolazione)}$$

- Molto utilizzato in ambito marketing per selezionare i clienti su cui operare (campagne focalizzate).
 - ✓ La classe desiderata è quella degli utenti che risponderanno positivamente all'attività di marketing
- Calcolare il lift e il guadagno effettivo per J48 JRIP e IB1

J48

```
a b c d e <-- classified as

3333 0 0 0 0 | a = not_recom

0 0 2 0 0 | b = recommend

0 0 239 89 0 | c = very_recom

0 0 63 3785 103 | d = priority

0 0 54 2332 | e = spec_prior
```

- Error rate = 3.11%
- % TP stimata della popolazione = 39.51%
- % del campione rispetto alla popolazione = 3,928 / 10,000 = 39.28%
- % TP stimata del campione = 96.36%
- Lift(J48)=(3,785/3,928) / (3,951/10,000)=0.9636/0.3951=2.4389
- Costo su popolazione = -10,000€ 50,000 × 5€ + 50,000 × 39.51% × 10€ = 260,000€ + 197,550€ = -62,450 €
- Costo su campione = $-10,000 \in -50,000 \times 39.28\% \times 5\emptyset + 50,000 \times 39.28\% \times 96.36\% \times 10 = -108,200\emptyset + 189,251\emptyset = +81,051\emptyset$

JRIP

```
a b c d e <-- classified as

3333 0 0 0 0 | a = not_recom

0 0 2 0 0 | b = recommend

0 0 214 114 0 | c = very_recom

0 1 124 3719 107 | d = priority

0 0 35 2351 | e = spec prior
```

- Error rate = 3.81%
- % TP stimata della popolazione = 39.51%
- % del campione rispetto alla popolazione = 3,868 / 10,000 = 38.68%
- % TP stimata del campione = 96.15%
- Lift(JRIP)=(3,719/3,868) / (3,951/10,000)=0.9615/0.3951=2.4336
- Costo su popolazione = -10,000€ 50,000 × 5€ + 50,000 × 39.51% × 10€ = 260,000€ + 197,550€ = -62,450 €
- Costo su campione = $-10,000 \in -50,000 \times 38.68\% \times 5\emptyset + 50,000 \times 38.68\% \times 96.15\% \times 10 = -108,200\emptyset + 185,954\emptyset = +77,754\emptyset$

IB1

```
a b c d e <-- classified as

3333 0 0 0 0 | a = not_recom

0 0 2 0 0 | b = recommend

0 0 194 134 0 | c = very_recom

0 0 0 3916 35 | d = priority

0 0 89 2297 | e = spec prior
```

- Error rate = 2.60%
- % TP stimata della popolazione = 39.51%
- % del campione rispetto alla popolazione = 4,139 / 10,000 = 41.39%
- % TP stimata del campione = 94.61%
- Lift(IB1)=(3,916/4,139) / (3,951/10,000)=0.9461/0.3951=2.3946
- Costo su popolazione = -10,000€ 50,000 × 5€ + 50,000 × 39.51% × 10€ = 260,000€ + 197,550€ = -62,450 €
- Costo su campione = -10,000 € 50,000 × 41.39% × 5€ + 50,000 × 41.39% × 94.61% × 10 = -113,475€ + 195,795€ = + 82,320€

Considerazioni

- La tecnica di classificazione più idonea sembra essere quella dei k-mediani
- Le performance degli alberi decisionali non aumentano utilizzando tecniche di boosting
- Il risultato non dipende dalla componente variance del dataset ma piuttosto dalla componente bias
 - ✓ I decision boundary di questo data set sono difficilmente modellabili tramite segmenti paralleli agli assi
- Per il problema del lift non è detto che il classificatore migliore sia quello con accuracy o error rate minimi visto che siamo interessati agli errori commessi per una particolare classe
 - ✓ Classificazione cost based con costi più elevati per gli errori di interesse

Cost SensitiveClassifier +J48

```
Cost Matrix a b c d e <-- classified as 0 1 1 1 1 3333 0 0 0 0 0 | a = not_recom 1 0 1 1 1 0 0 0 2 0 0 | b = recommend 1 1 0 10 1 0 0 308 20 0 | c = very_recom 1 1 1 0 1 0 0 0 161 3642 148 | d = priority 1 1 1 10 0 0 0 11 2375 | e = spec prior
```

- Error rate = 3.42%
- % TP stimata della popolazione = 39.51%
- % del campione rispetto alla popolazione = 3,673 / 10,000 = 36.73%
- % TP stimata del campione = 99.16%
- Lift(CSC+J48)=(3,642/3,673) / (3,951/10,000)=0.9916/0.3951=2.5096
- Costo su popolazione = -10,000€ 50,000 × 5€ + 50,000 × 39.51% × 10€ = 260,000€ + 197,550€ = -62,450 €
- Costo su campione = $-10,000 \in -50,000 \times 36.73\% \times 5\emptyset + 50,000 \times 36.73\% \times 99.16\% \times 10 = -101,825\emptyset + 182,107\emptyset = +80,282\emptyset$

Gli errori si riducono, ma la riduzione del campione controbiliancia la maggior precisione