

4. Algoritmi

One rule

- Classificatore estremamente semplice che costruisce una singola regola decisionale basata su un solo attributo
- La sua efficacia, sorprendentemente elevata su dataset semplici, lo rende una solida baseline per confronti
- Mecanismo**
 - Per ciascun attributo, si costruisce una regola che associa ad ogni valore la classe più frequente
 - Si valuta l'accuratezza di ciascuna regola
 - Si sceglie l'attributo con la regola più accurata (minimo errore)
- Assunzioni**
 - L'attributo scelto è sufficiente per una buona predizione
 - I dati sono etichettati correttamente e senza ambiguità
- Vantaggi**
 - Semplice da implementare e spiegare
 - Robusto su dataset con forte rilevanza di un singolo attributo
 - Utile come baseline nei confronti sperimentali
 - Ignora tutte le interazioni tra attributi
- Svantaggi**
 - Rischio elevato di overfitting con dati numerosi o molti valori distinti (uso di numeri per descrivere una condizione)

Soluzione

Imporre un numero minimo di istanze (ex. 3) in una classe

64	65	68	69	70	71	72	72	75	75	80	81	83	85	
Tues	S	Mon	F	Tues	Tues	Tues	Tues	Mon	Tues	Mon	S	Tues	F	Mon

64	65	68	69	70	71	72	72	75	75	80	81	83	85
Tues	Mon	Tues	Tues	Tues	S	Mon	Tues	Tues	Tues	Mon	Tues	F	Mon

Ogni intervallo prende il valore della classe di maggioranza (es. se ci sono da 64 a 70 gradi fahrenheit si gioca (S o 2 no))

- Variazioni**
 - Hyper-pipes**
 - Si costruisce una regola per ogni classe come congiunzione di test, una per ogni attributo
 - Il test verifica se il valore dell'istanza si trova all'interno di un intervallo
 - Il test verifica se il valore è uno dei sottoinsiemi di valori dell'attributo
 - Dato del minimo e del massimo osservati nei dati di addestramento
 - Sottoinsieme dato da tutti i possibili valori osservati nell'addestramento
- La classe con il maggior numero di test corrispondenti viene predetta

Modelli statistici

- Due assunzioni principali**
 - Gli attributi sono ugualmente importanti
 - Gli attributi sono statisticamente indipendenti
- Conoscere il valore di un attributo non dice nulla sul valore di un altro
- Questa assunzione non è praticamente mai ponibile nella realtà
- Ma questo schema funziona comunque bene nella pratica
- | Outlook | Temp. | Humidity | Windy | Play |
|---------|-------|----------|-------|------|
| Sunny | Cool | High | True | ? |

Likelihood of the two classes

For "yes" = $2/9 \times 3/9 \times 3/9 \times 3/9 \times 9/14 = 0.0053$

For "no" = $3/5 \times 1/5 \times 4/5 \times 3/5 \times 5/14 = 0.0206$

Conversion into a probability by normalization:

$P(\text{"yes"}) = 0.0053 / (0.0053 + 0.0206) = 0.205$

$P(\text{"no"}) = 0.0206 / (0.0053 + 0.0206) = 0.795$
- Es. 9/14 e 5/14
- Pr(H|E) è la probabilità di un evento H data un'evidenza E
- Pr(E|H) probabilità dell'evidenza dato un certo evento
- Pr(H) è la probabilità che ci sia un evento prima che una evidenza sia visibile
- Pr(E) è la probabilità dell'evidenza
- Difficile da calcolare
- Se gli eventi sono indipendenti la probabilità congiunta è il prodotto delle singole probabilità
- Rappresentabile tutto mediante la probabilità di Bayes
- Naive Bayes**
 - Assunzione Naive su Bayes: L'evidenza si divide in parti (cioè attributi) che sono indipendenti
 - Dato che: La probabilità congiunta è il prodotto delle singole probabilità
 - Siamo nel caso in cui gli eventi sono statisticamente indipendenti
 - Alloca: Possiamo contare la probabilità di E dato H come la sua proiezione rispetto ai singoli attributi
 - $$Pr(H|E) = \frac{Pr(H_1|H)Pr(H_2|H) \dots Pr(H_n|H)Pr(H)}{Pr(E)}$$

← Evidenza E

Outlook	Temp.	Humidity	Windy	Play
Sunny	Cool	High	True	?

Probability of class "yes"

$Pr(\text{yes}|E) = Pr(\text{Outlook} = \text{Sunny} | \text{yes}) \times Pr(\text{Temp} = \text{Cool} | \text{yes}) \times Pr(\text{Humidity} = \text{High} | \text{yes}) \times Pr(\text{Windy} = \text{True} | \text{yes}) \times \frac{Pr(\text{yes})}{Pr(H)}$

$\frac{2}{9} \times \frac{3}{9} \times \frac{3}{9} \times \frac{3}{9} \times \frac{9}{14}$
 - Nell'esempio: Se abbiamo 0 in uno dei valori è un problema
 - Cosa succede se il valore di un attributo non è presente in ogni valore della classe? Tutto quanto va a 0 nel prodotto. Ex. P(Outlook = Overcast | No) = 0/5
 - Preprocessing necessario: Ex. Aumentiamo tutti quanti di 1 → Smoothing
 - Attenzione all'overfitting: Il modello è più complesso

Decision Tree

- Modelli di classificazione supervisionata che utilizzano una strategia divide-et-impera
- Innanzitutto si seleziona un attributo per il nodo radice
- Si crea un ramo per ogni possibile valore dell'attributo
- Ex. Per l'attributo Outlook: Sunny, Overcast, Rainy
- Poi si dividono le istanze in sottoinsiemi, uno per ogni ramo che si estende dal nodo
- Per ogni attributo scelgo un sottoinsieme di attributi
- Ex. Humidity o Windy
- Infine si ripete ricorsivamente per ogni ramo, usando solo le istanze che raggiungono il ramo
- Ci si ferma se tutte le istanze hanno la stessa classe
- Ex. Sunny -> Humidity porta sempre a Yes

Outlook	Humidity	Giocare?
Sunny	Normal	Yes
Sunny	Normal	Yes

- Ogni nodo interno rappresenta un test su un attributo
- Ogni ramo rappresenta un esito del test
- Ogni foglia rappresenta una classe di output
- Costruzione**
 - Voglio avere l'albero più piccolo → Più efficace da scorrere
 - Sciegliere l'attributo che produce il nodo più puro
 - Ovvero che genera quante più istanze della stessa classe
 - Etichetta che seguiamo: Criterio di purità popolare → Guadagno di informazioni
 - Aumenta con la purità media dei subset
 - Strategia: Sciegliere l'attributo con il massimo guadagno di informazioni
 - L'entropia misura la "purezza" di un insieme
 - entropia(p_1, p_2, \dots, p_n) = $-p_1 \log p_1 - p_2 \log p_2 - \dots - p_n \log p_n$
 - Dove p è la frazione di elementi che soddisfanno la condizione
 - Probabilità sul totale / totale misurazioni * quantitativo informativo
 - Si misura il guadagno di informazione come Entropia prima della divisione - Entropia dopo la divisione
 - Entropia (S) = $-9/14 \cdot \log(9/14) - 5/14 \cdot \log(5/14) \rightarrow 0.940$ bits
 - Quello che permette di minimizzare l'entropia residua e massimizzare l'informazione acquisita
 - Altra tecnica: Gain Ratio → Si valuta il numero di rami che genera il nodo → Se sono tanti rami lo penalizziamo altrimenti no
- Criteri di scelta degli attributi di divisione**
 - Convertire i decision tree in un set di regole
 - Ogni regola copre le istanze del dataset, di cui p appartengono alla classe target
 - Per questo si chiama approccio di copertura
 - Qualità della regola → Si misura con il rapporto p/n
 - Maggiore è questo valore, più la regola è considerata affidabile
 - Costruzione delle regole → Le regole si costruiscono cercando condizioni che identificano sottoinsiemi "puri"
 - Puri = alta concentrazione di positivi
 - Vantaggi**
 - Ogni regola è indipendente e interpretabile
 - Possiamo essere più semplici da leggere rispetto agli alberi
 - Svantaggi**
 - Possibili conflitti tra regole
 - Serve pruning per mantenere il modello semplice e preciso.
 - Esempio: PRISM**
 - L'obiettivo è massimizzare l'accuratezza della regola tenendo in considerazione
 - il numero totale di istanze coperte dalla regola (n)
 - Gli esempi positivi della classe coperti dalla regola (p)
 - Numero di errori fatti dalla regola 1-p
 - Ovvero il numero di elementi a cui si applica la regola
 - Abbiamo finito quando $p/n = 1$ o l'insieme delle istanze non può essere suddiviso ulteriormente

Algoritmi di copertura

- Convertire i decision tree in un set di regole
- Ogni regola copre le istanze del dataset, di cui p appartengono alla classe target
- Per questo si chiama approccio di copertura
- Qualità della regola → Si misura con il rapporto p/n
- Maggiore è questo valore, più la regola è considerata affidabile
- Costruzione delle regole → Le regole si costruiscono cercando condizioni che identificano sottoinsiemi "puri"
- Puri = alta concentrazione di positivi
- Vantaggi**
 - Ogni regola è indipendente e interpretabile
 - Possiamo essere più semplici da leggere rispetto agli alberi
- Svantaggi**
 - Possibili conflitti tra regole
 - Serve pruning per mantenere il modello semplice e preciso.
- Esempio: PRISM**
 - L'obiettivo è massimizzare l'accuratezza della regola tenendo in considerazione
 - il numero totale di istanze coperte dalla regola (n)
 - Gli esempi positivi della classe coperti dalla regola (p)
 - Numero di errori fatti dalla regola 1-p
 - Ovvero il numero di elementi a cui si applica la regola
 - Abbiamo finito quando $p/n = 1$ o l'insieme delle istanze non può essere suddiviso ulteriormente

Clustering

- Le tecniche di clustering si usano quando non ci sono classi da predire
- L'obiettivo è di dividere le istanze in gruppi naturali
- Disjoint vs Overlapping**
 - Tipi di cluster
- Deterministic vs Probabilistic**
 - Tipi di cluster
- Fiat vs Hierarchical**
 - Un dataset adatto al clustering è una collezione di punti, che sono oggetti che appartengono a uno spazio
- Uno spazio è un insieme universale di punti, dai quale sono stati estratti i punti nel dataset
- O algoritmi aggregativi
- Ogni punto inizia come cluster separato
- I cluster vengono poi uniti iterativamente in base alla distanza
- Possiamo dividere gli algoritmi di clustering in due gruppi
- I punti vengono assegnati al cluster tramite un criterio di ottimizzazione
- Point assignment**
 - Una metrica di distanza (spazio euclideo o altro)
 - Un centroide è il punto centrale (o medio) di un cluster
 - Una regola di distanza inter-cluster (es. distanza tra centroidi)
 - Una condizione di stop
 - Ad ex. numero desiderato di cluster o soglia di dissimilarità
 - Richiede
 - Il risultato è rappresentato tramite un dendrogramma
 - Diagramma ad albero che mostra come i cluster vengono costruiti passo dopo passo
- Clustering gerarchico**
 - Algoritmo flat che richiede in input il numero k di cluster desiderati
 - Si inizializzano casualmente k centroidi
 - Si assegna ogni punto al centroide più vicino
 - Si ricalcolano i centroidi come media dei punti assegnati
 - Si ripete finché i centroidi non cambiano più (convergenza)
 - Fast
 - K-means
- Possiamo separare le istanze quasi perfettamente se introduciamo un confine che sia una linea retta, ma non perpendicolare agli assi
- È un classificatore lineare ed è una somma ponderata dei valori dei vari attributi

Classificazione con funzioni matematiche

- Funzioni discriminanti lineari**
 - Ce ne sono infinite che classificano potenzialmente il training set in maniera perfetta
 - Non è banale la scelta della retta migliore
 - Procedura
 - Rappresenti un obiettivo
 - Definire una funzione obiettivo che possa essere calcolata per un particolare set di pesi e un particolare set di dati
 - Linear regression
 - Logistic regression
 - SVM
 - Crearla è complesso
 - Ognuna usa una funzione obiettivo diversa
 - Modelli lineari che separano le classi con il massimo margine possibile
 - Costruisce una banda delimitata da due rette parallele
 - La retta centrale è la funzione decisionale
 - I punti che toccano la banda sono i support vectors e determinano la posizione dell'iperpiano
 - Penalizziamo un punto di addestramento perché si trova sul lato sbagliato del confine di decisione
 - Ovviamente avremo punti che usciranno dal lato sbagliato della banda
 - Due casi
 - Set linearmente separabili
 - Set non linearmente separabili
 - Addestramento migliore è un equilibrio tra un margine abbondante e un basso errore totale
 - Modello di classificazione binaria che calcola la probabilità di appartenenza alla classe positiva usando la funzione sigmoide
 - Classificazione
 - Un esempio viene classificato come positivo se $o(x) > 0.5$, altrimenti negativo
 - Linear Regression
 - Tipi
 - Least squares regression
 - Penalizza errori molto grandi
 - La regressione lineare standard usa i quadrati di questi errori
 - Inconveniente
 - Molto sensibile ai dati
 - In punti errati o fuori norma possono alterare gravemente la funzione lineare risultante
 - Minimizzazione dell'errore al quadrato
 - Si scelgono k+1 coefficienti che minimizzano l'errore al quadrato dai dati di addestramento
 - Gli esempi di training sono salvati da qualche parte
 - Non costruisco un modello, trovo l'esempio più vicino a me
 - Se l'attributo è categorico è più difficile trovare il valore più vicino, se invece è numerico si può usare la distanza euclidea

Instance Based Learning

- Differenti attributi sono misurati su scale diverse → I valori dei dati set vanno normalizzati
- Per classificare una nuova istanza, si considerano i k esempi più vicini nel dataset e si assegna la classe più frequente tra questi
- Funzionamento**
 - Si calcola la distanza (tipicamente euclidea) tra la nuova istanza e tutte le istanze del training set
 - Si selezionano i k esempi più vicini
 - Si applica una funzione di combinazione (es. voto di maggioranza per classificazione, media per regressione)
- Pesatura dei vicini**
 - $$peso = \frac{1}{distanza^k}$$
 - k=3 -> 3-NN
 - Il nome k-NN dipende dal numero di k
 - Numero dispari → Convenienti per evitare vogli di maggioranza in classificazioni in due classi
 - Quanti vicini (k) usare?
 - Maggiore è k → Maggiore sono i vicini con cui confrontarsi
 - Per classificazione → Si predice sempre la classe maggioritaria nell'intero set di dati
 - Per regressione → La media di tutti i valori target
 - Se aumentiamo k (-n) al massimo possibile, si usa l'intero dataset per ogni predizione
 - Due casi
 - Nessuna complessità → k=n
 - Modello estremamente complesso → k=1
 - Pone confini complicati
 - Ogni esempio si troverà in una regione etichettata della propria classe
 - K alto → Basso overfitting
- Overfitting