



UNIVERSITÀ
DEGLI STUDI
FIRENZE

Scuola di Scienze Matematiche, Fisiche e Naturali
Corso di Laurea in Informatica

Tesi di Laurea

TITOLO ITALIANO

TITOLO INGLESE

NOME CANDIDATO

Relatore: *Relatore*
Correlatore: *Correlatore*

Anno Accademico 2014-2015

Nome candidato: *Titolo italiano*, Corso di Laurea in Informatica, © Anno
Accademico 2014-2015

INDICE

1	Introduzione	7
2	Fondamentali	9
2.1	Scansione dell'ambiente	9
2.2	Raccoglimento dei dati	10
2.3	Onde Magnetiche	10
2.4	Elaborazione dei dati	11
2.4.1	Estrazione del magnitudo	11
2.4.2	Raggruppamento	11
2.4.3	Estrazione degli attributi	12
2.5	Costruzione di un modello	12
2.6	IA ed Apprendimento	13
2.7	Tipi di apprendimento	13
2.8	Alcune nozioni sull'apprendimento automatico	15
2.8.1	Verifica delle prestazioni	15
2.8.2	Rumore e sovradattamento	16
2.8.3	Cross validation	16
2.9	K Nearest Neighbour	17
2.9.1	Scelta del parametro k	18
2.10	Naive bayes	18
2.10.1	Teorema di bayes	18
2.10.2	Rete bayesiana	18
2.10.3	Naive Bayes	19
2.10.4	Costruzione di una rete bayesiana ingenua	19
2.10.5	Predizione di risultati	20
2.10.6	Modelli generativi e discriminativi	20
2.10.7	Calcolo della probabilita' a priori	20
2.10.8	Modelli ad eventi	20
2.10.9	Un piccolo esempio basato su naive bayes	21
2.11	Alberi di decisione	22
2.11.1	Costruzione di un albero di decisione	23
2.11.2	Classificazione dei nuovi elementi	23
2.11.3	Un esempio di albero di decisione	23
2.12	Ricerca della posizione	24
3	Struttura del software	27
3.1	Linguaggi e framework	27

2 Indice

3.2	Interfaccia grafica	27
3.3	Persistenza dei dati	28
3.4	Struttura del codice e design pattern	28
3.5	Analisi dei dati	29
3.6	Codice usato per l'analisi	29
4	Test	31
4.1	Base dei test	31
4.2	Caratteristiche	32
4.3	Classificatori a confronto	32
4.4	Un rimedio ingenuo al rumore	34
5	Miglioramenti	37
5.1	Possibili miglioramenti per la raccolta di dati	37

ELENCO DELLE FIGURE

Figura 1	Assi x, y, z centrati sul cellulare	10
Figura 2	Un esempio di mappa del campo magnetico	11
Figura 3	Un esempio grafico di insieme d'addestramento ed una sua classificazione	14
Figura 4	Classificazione e regressione a confronto	14
Figura 5	Apprendimento supervisiono e non supervisionato a confronto	15
Figura 6	L'apprendimento con rinforzo e' molto adatto ai giochi, come per esempio pacman	15
Figura 7		16
Figura 8	Rappresentazione grafica della Cross Validation	17
Figura 9	Esempio grafico dell'algoritmo KNN	17
Figura 10	Esempio di rete bayesiana	19
Figura 11	Esempio di rete bayesiana ingenua	19
Figura 12	Esempi di partite di tennis giocate in base alle condizioni meteorologiche e tabella di distribuzione delle probabilita'	22
Figura 13		22
Figura 14	Tabella degli attributi meteorologici con valore associato un booleano che stabilisce se la partita e' stata giocata o no	24
Figura 15	Albero di decisione ricavato dalla tabella precedente	24
Figura 16	Interfaccia grafica all'avvio dell'applicazione	28
Figura 17	Una piccola raffigurazione delle stanze usate per le prove con sopra scritto la <i>label</i> assegnata	31
Figura 18	Grafico in 2 dimensioni della media e varianza di tutte le onde magnetiche. I colori dei punti rappresentano le etichette	32
Figura 19		33
Figura 20		33
Figura 21	Numero di errori nella predizione per etichetta	34
Figura 22	Percentuale di errore nella predizione per etichetta	34

4 Elenco delle figure

Figura 23 35

"Inserire citazione"
— *Inserire autore citazione*

INTRODUZIONE

La seguente tesi e' basata su un tirocinio esterno svolto con l'azienda KeepUp in cui e' stata sviluppata la base di un'applicazione Android col compito di localizzare all'interno degli edifici la posizione dello *smartphone* sfruttando le distorsioni del campo magnetico terrestre.

Durante il primo capitolo, i fondamentali, approfondiremo il tipo di dato che dobbiamo trattare, le sue origini e la sua struttura per poi passare all'apprendimento automatico, usato per predire la posizione dell'utente all'interno dell'edificio elencandone i vari tipi e definendo alcuni termini gergali per concludere con la descrizione approfondita di alcuni classificatori molto conosciuti nel mondo dell'AI.

Nel secondo capitolo parleremo della struttura del software realizzato durante il tirocinio, dal codice mobile Android a quello per la verifica dei risultati.

Nel terzo capitolo vedremo i risultati ottenuti durante il tirocinio tramite vari grafici che mostreranno punti di forza e debolezza della nostra applicazione.

Nel quarto capitolo discuteremo dei possibili miglioramenti futuri all'applicazione.

FONDAMENTALI

L'identificazione della posizione all'interno di un edificio si svolge in 2 fasi:

1. Scansione dell'ambiente
2. Ricerca della posizione

2.1 SCANSIONE DELL'AMBIENTE

La scansione dell'ambiente consiste in un'analisi dell'ambiente per mezzo del magnetometro con lo scopo di registrare tutte le distorsioni del campo magnetico presenti in modo da permettere la ricerca della posizione.

2.2 RACCOGLIMENTO DEI DATI

Il raccoglimento dei dati avviene tramite il magnetometro del nostro *smartphone* da cui viene catturato diverse volte in un secondo il campo magnetico intorno ad esso. Ad ogni onda magnetica viene assegnata una *label*: un numero che identifica univocamente una parte dell'ambiente chiuso il cui uso verrebbe definito durante la fase di costruzione del modello.

2.3 ONDE MAGNETICHE

Le onde magnetiche sono un vettore di 3 elementi, quindi in \mathbb{R}^3 che classifica. Il primo valore rappresenta la forza del campo magnetico lungo l'asse X, il secondo lungo Y ed il terzo lungo Z.

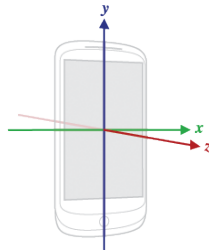


Figura 1: Assi x, y, z centrati sul cellulare

I tre valori sono espressi in μT (micro Tesla), unità di misura della densità di un flusso magnetico. Le onde magnetiche raccolte sono dati continui sia positivi che negativi. L'intensità dell'onda deriva dalla distorsione del campo magnetico generata dagli oggetti statici intorno al punto e dipende anche dalla velocità con cui ci muoviamo per cui, per semplicità, assumeremo d'ora in poi una velocità costante di 3 piedi al secondo.

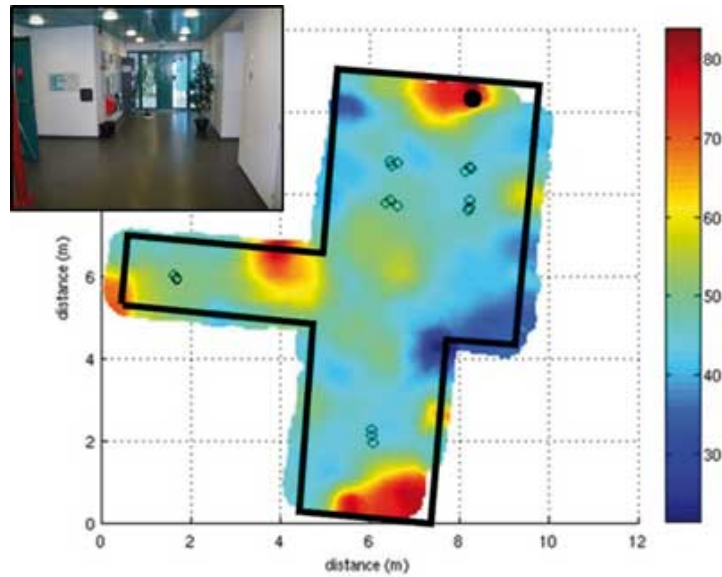


Figura 2: Mappa del campo magnetico all'università di Oulu: Discus Entrance hall

2.4 ELABORAZIONE DEI DATI

2.4.1 Estrazione del magnitudo

Per estrarre l'intensità di ogni onda magnetica eseguiamo semplicemente la norma euclidea di un vettore:

$$\sqrt{x^2 + y^2 + z^2}$$

2.4.2 Raggruppamento

Le onde magnetiche con la stessa *label* vengono raggruppate in *fingerprints*, insiemi di dimensione prefissata. A livello logico, ogni *fingerprint* cerca di identificare univocamente un punto all'interno di una zona, identificata con una *label*. L'insieme di *fingerprints* quindi, cerca di distinguere, tramite le caratteristiche dei campi elettromagnetici di ciascun punto, ogni *label* dall'altra.

2.4.3 Estrazione degli attributi

Per ogni *fingerprint*, l'estrazione delle *features* consiste nell'estrazione di variabili statistiche. In questo specifico caso sono:

- Media
- Varianza
- Deviazione standard
- Mediana
- Media troncata
- Coefficiente di variazione
- Massimo
- Minimo
- 1°, 5°, 95°, 99° percentile
- 1°, 2°, 3° quartile

2.5 COSTRUZIONE DI UN MODELLO

Dopo aver raccolto ed elaborato le onde magnetiche, un algoritmo creerà un classificatore in grado di assegnare un'etichetta ai nuovi input ricevuti durante l'uso dell'utente finale. Sono stati utilizzati vari classificatori per cercare il più preciso fra tutti. Qui di seguito introdurremo la teoria dietro ad ogni classificatore utilizzato.

2.6 IA ED APPRENDIMENTO

La definizione secondo wikipedia¹ di IA e' la seguente:

Definizioni specifiche possono essere date focalizzandosi o sui processi interni di ragionamento o sul comportamento esterno del sistema intelligente ed utilizzando come misura di efficacia o la somiglianza con il comportamento umano o con un comportamento ideale, detto razionale:

- Agire umanamente: il risultato dell'operazione compiuta dal sistema intelligente non e' distinguibile da quella svolta da un umano.
- Pensare umanamente: il processo che porta il sistema intelligente a risolvere un problema ricalca quello umano. Questo approccio Ã" associato alle scienze cognitive.
- Pensare razionalmente: il processo che porta il sistema intelligente a risolvere un problema Ã" un procedimento formale che si rifÃ alla logica.
- Agire razionalmente: il processo che porta il sistema

L'apprendimento automatico e' una branca dell'intelligenza artificiale che permette al computer di apprendere da insiemi di dati per generare conoscenza con lo scopo di effettuare previsioni.

2.7 TIPI DI APPRENDIMENTO

Esistono 3 tipi di apprendimento:

- Apprendimento supervisionato: al calcolatore vengono forniti esempi del tipo (x, y) per poter apprendere. L'obiettivo della predizione per nuovi input sara' quello di calcolare la variabile dipendente. Y puo' essere sia discreta che continua: nel primo caso si parla di classificazione, nel secondo di regressione.

¹ https://it.wikipedia.org/wiki/Intelligenza_artificiale

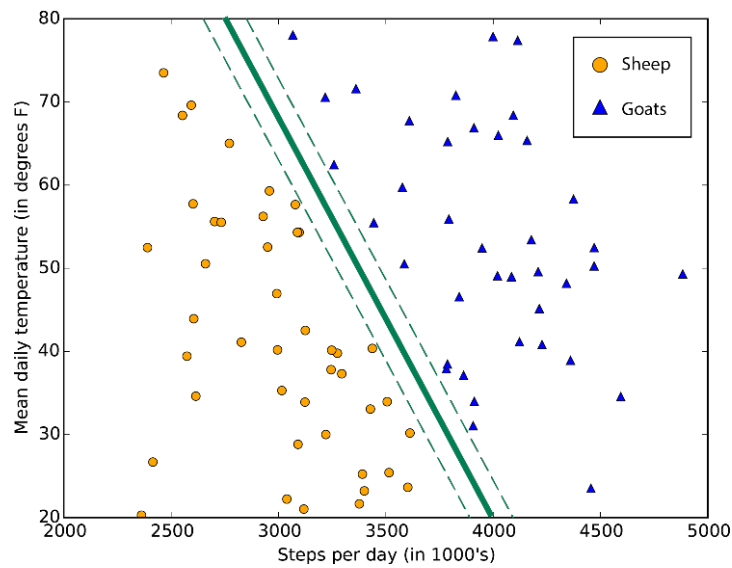


Figura 3: Un esempio grafico di insieme d'addestramento ed una sua classificazione

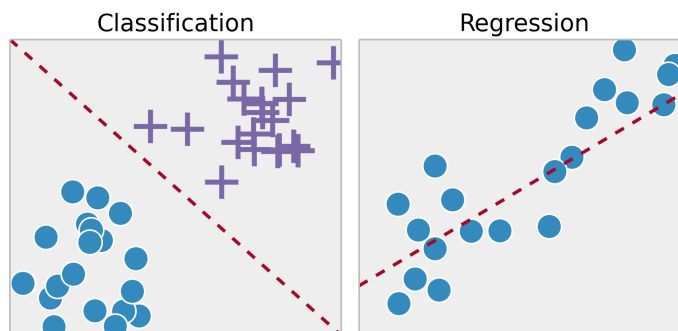


Figura 4: Classificazione e regressione a confronto

- **Apprendimento non supervisionato:** Gli esempi non contengono una variabile dipendente ma solo un insieme di attributi x . L'obiettivo è quello di inferire pattern nascosti dai dati non etichettati. Un'importante applicazione è il *clustering*: raggruppare i dati in base ad una similarità fra di essi.

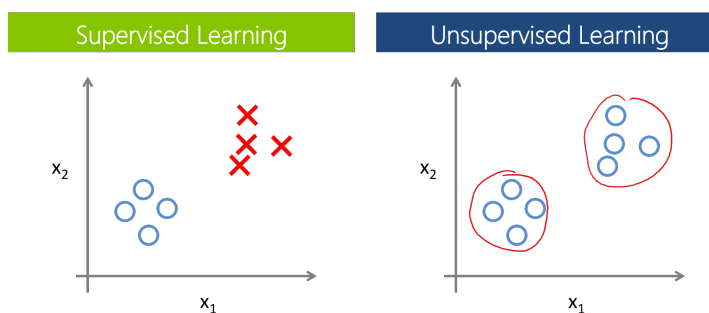


Figura 5: Apprendimento supervisionato e non supervisionato a confronto

- Apprendimento con rinforzo: per apprendere viene fornita una funzione ricompensa cioè una funzione che, data un'azione effettuata dall'agente, restituirà una ricompensa di tipo numerico.

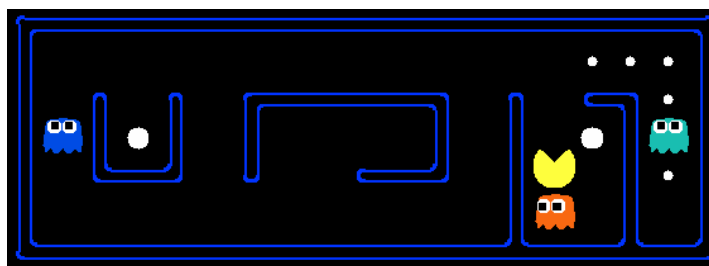


Figura 6: L'apprendimento con rinforzo è molto adatto ai giochi, come per esempio Pacman

2.8 ALCUNE NOZIONI SULL'APPRENDIMENTO AUTOMATICO

2.8.1 Verifica delle prestazioni

Per verificare la correttezza del classificatore viene diviso in 2 parti il *dataset* a nostra disposizione: il primo si chiama insieme di addestramento mentre il secondo insieme di test. Il nostro classificatore si allenerà sull'insieme di addestramento, cioè imparerà dagli esempi come classificare i nuovi input ed effettuerà le predizioni sull'insieme di test, su cui verranno misurate le prestazioni in base al rapporto fra predizioni corrette e totali.

2.8.2 Rumore e sovradattamento

Un problema molto comune durante l'addestramento dei classificatori e' il rumore: durante la discriminazione degli esempi, ci ritroviamo ad un punto in cui gli attributi rimanenti sono identici ma con etichetta differente. Una probabile causa potrebbe essere la presenza di errori nei dati mentre una possibile soluzione consisterebbe nel voto di maggioranza. Un'altro problema e' il sovradattamento: la costruzione di un classificatore consistente con tutti gli esempi a causa dell'utilizzo di attributi irrilevanti nella classificazione. Supponiamo di voler predire l'esito del lancio di un dado e fra gli attributi di avere ora, giorno, mese ed anno; ecco un esempio lampante di sovradatamento. Nei casi reali tuttavia non sono cosi' evidenti gli attributi insignificanti e, per esempio, una tecnica utilizzata negli alberi di decisione e' la potatura.

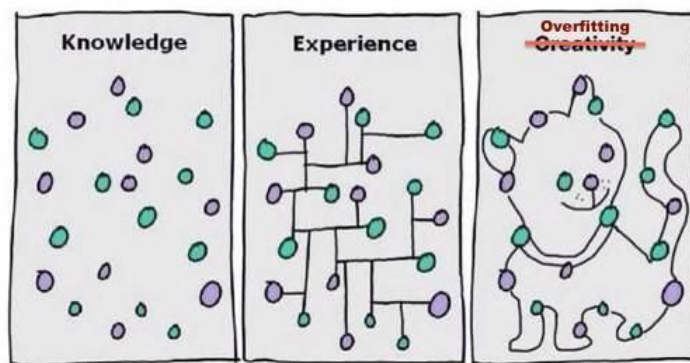


Figura 7

2.8.3 Cross validation

Per migliorare le performance di un classificatore si suddivide l'intero *dataset* in n parti uguali (di solito 10), ognuna delle quali svolgera' per una volta il ruolo di insieme di test mentre il resto sara' l'insieme di addestramento. Questa tecnica risolve vari problemi tra i quali l'*overfitting*.



Figura 8: Rappresentazione grafica della Cross Validation

2.9 K NEAREST NEIGHBOUR

Uno degli algoritmi piu' semplici di apprendimento automatico, e' il *k-nearest-neighbours* dove l'input consiste in k elementi presi dal *training set* piu' vicini in base ad un criterio scelto da chi utilizza l'algoritmo (per esempio la distanza euclidea o di Mahalanobis). Il KNN viene definito un algoritmo pigro(*lazy*) perche' non ha bisogno di apprendere dall'insieme di addestramento per poter creare un classificatore ma puo' usare direttamente i dati forniti per classificare i nuovi esempi. L'inconveniente che abbiamo ricade sul costo computazione: esso e' proporziale alla grandezza del dataset.

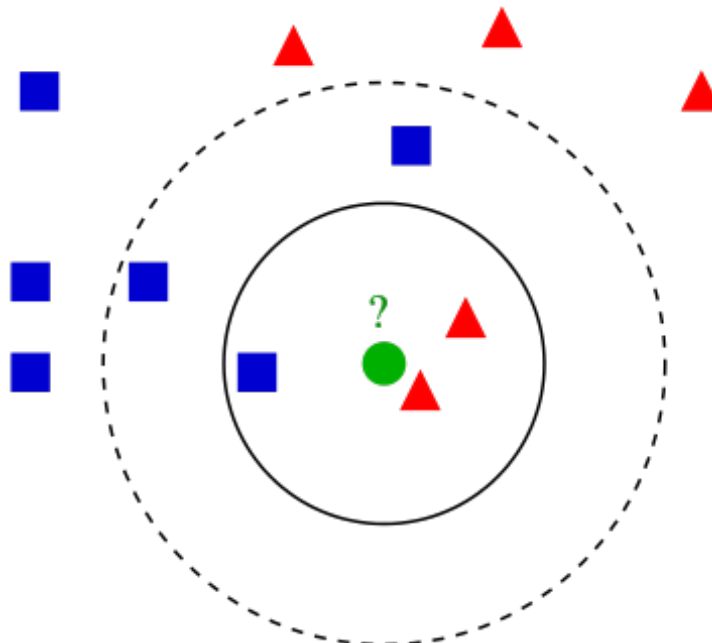


Figura 9: Esempio grafico dell'algoritmo KNN

2.9.1 Scelta del parametro k

La scelta del parametro dipende, ovviamente, dal tipo dei dati che abbiamo e dalla quantità, anche se in generale più è grande k meno rumore viene generato da questo algoritmo. Un buon metodo per trovare il giusto valore è l'uso di tecniche euristiche, come la *cross validation*. Un'altra fonte di rumore di cui bisogna stare attenti è la presenza di *features* insignificanti nella ricerca del vicino. Per porre rimedio possiamo, ad esempio, usare un algoritmo genetico per selezionare le *features* più significative

2.10 NAIVE BAYES

Un altro tipo di apprendimento usato per classificare le *label* è Naive bayes: un algoritmo di classificazione e regressione basato sulla statistica. Prima di spiegare in cosa consiste occorre spiegare un paio di concetti:

2.10.1 Teorema di bayes

Il teorema di bayes ci fornisce una relazione fra probabilità condizionate molto utile nel calcolo probabilistico ma anche per l'apprendimento automatico. L'equazione è la seguente:

$$P(B|A) = \frac{P(A|B)P(B)}{P(A)}$$

2.10.2 Rete bayesiana

Una rete bayesiana è un grafo diretto aciclico i cui nodi sono le variabili casuali del sistema mentre gli archi rappresentano la condizione di dipendenza fra nodi. Ad ogni nodo è associata una tabella di distribuzione delle probabilità la cui complessità è proporzionale al numero di archi entranti.

Per esempio se il nodo con variabile casuale Leggere ha un arco verso Istruito allora possiamo dire che Istruito è condizionalmente dipendente da Leggere. Qui sotto potete vedere una raffigurazione grafica di una semplice rete bayesiana formata da nodi padre ed un figlio:

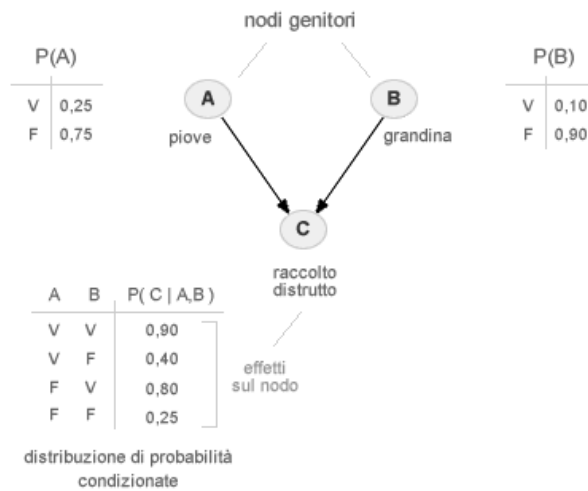


Figura 10: Esempio di rete bayesiana

2.10.3 Naive Bayes

Naive bayes e' una rete bayesiana in cui si assume l'indipendenza condizionale fra tutte le variabili casuali del sistema data la classe. Questa forte assunzione non mira a modellare esattamente la realta' ma fornisce delle buone performance sulla predizione di una classe.

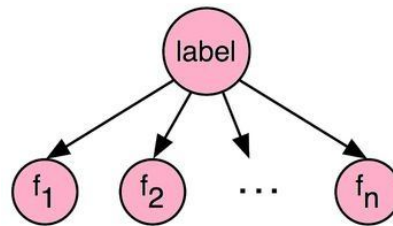


Figura 11: Esempio di rete bayesiana ingenua

2.10.4 Costruzione di una rete bayesiana ingenua

Preso l'insieme di esempi, ricavare la tabella di distribuzione delle probabilita' al variare di tutti i valori che la variabile dipendente puo' assumere. Successivamente costruire un grafo che ha come nodo radice la variabile dipendente e come figli tutte le variabili indipendenti.

2.10.5 *Predizione di risultati*

Preso un insieme di attributi \mathbf{x} con valori $x_1 x_2 x_3 \dots x_n$ e tutti i valori della variabile dipendente $\mathbf{y} = y_1 y_2 y_3 \dots y_m$, per classificare l'insieme applicheremo il teorema di Bayes con un'approssimazione molto usata nel campo scientifico chiamata MAP (Massima ipotesi a posteriori)

$$\begin{aligned} \text{etichetta} &= y_j \text{ in } \max_j P(y_j | X_1 = x_1, X_2 = x_2 \dots X_N = x_n) \\ &= \max_j \frac{P(X_1 = x_1, X_2 = x_2 \dots X_N = x_n | y_j) P(y_j)}{P(X_1 = x_1, X_2 = x_2 \dots X_N = x_n)} \\ &= \max_j \frac{\prod_i P(x_i | y_j) P(y_j)}{P(X_1 = x_1, X_2 = x_2 \dots X_N = x_n)} \end{aligned}$$

i cui valori della precedente equazione si possono ricavare dalla tabella di distribuzione delle probabilità.

2.10.6 *Modelli generativi e discriminativi*

L'espressione $P(y_j | \mathbf{x})$ può essere calcolata in modi differenti in base al modello trattato: nel caso di modelli discriminativi, l'espressione viene ricavata direttamente dai dati sorgente mentre nei modelli generativi l'obiettivo è di costruire una distribuzione congiunta di probabilità per $P(\mathbf{x}, y_j)$ oppure da $P(\mathbf{x} | y_j)$ e $P(y_j)$.

2.10.7 *Calcolo della probabilità a priori*

Per calcolare la probabilità a priori esistono varie tecniche: l'equiprobabilità $\left(\frac{1}{|\mathbf{y}|}\right)$, rapporto fra esempi di classe j e il totale degli esempi dall'insieme di addestramento, le distribuzioni di probabilità oppure un modello non parametrico dall'insieme di addestramento.

2.10.8 *Modelli ad eventi*

La distribuzione viene anche chiamata modello ad eventi perché tratta gli attributi come probabilità di eventi. Ci sono 3 modelli usati con *Naive Bayes* che sono:

- Bernoulli: gli attributi sono di tipo *booleano* e l'attributo $x_i \in \mathbf{x} = (x_1, x_2, \dots, x_n)$ vale 1 se l'evento i è avvenuto, o altrimenti. Nel caso bernoulliano possiamo valutare $P(\mathbf{x} | y_j)$ come

$$P(\mathbf{x}|y_j) = \prod_{i=1}^n P_{ij}^{x_i} (1 - P_{ij})^{(1-x_i)}$$

Dove P_{ij} e' la probabilita' per y_j che x_i sia vero. Un esempio canonico e' quello della classificazione dei documenti, dove x_i rappresenta la presenza del termine w_i nei documenti di classe y_j e di conseguenza P_{ij} la probabilita' di trovarlo. Dobbiamo notar bene che Bernoulli a differenza della multinomiale valuta nella produttoria anche la probabilita' che l'evento i non avvenga.

- Multinomiale: Nella multinomiale gli esempi rappresentano la frequenza con il quale gli eventi sono stati generati dalla multinomiale $(p_1 \dots p_n)$ dove p_i e' la probabilita' che i occorra. Gli attributi $\mathbf{x} = (x_1, x_2 \dots x_n)$ contano quante volte l'evento i e' avvenuto. La probabilita' condizionata $P(\mathbf{x}|y_j)$ e' stimata come

$$P(\mathbf{x}|y_i) = \frac{(\sum_i x_i)!}{\prod_i x_i!} \prod_i p_{ki}^{x_i}$$

- Gaussiana: viene utilizzata con dati continui e si assume che siano distribuiti in base alla Gaussiana. Per ogni classe y_j , viene ricavata la media μ_j e la varianza σ_j^2 . Supponiamo di aver raccolto un insieme di valori v allora la probabilita' sara':

$$P(x = v|y_j) = \frac{1}{\sqrt{2\pi\sigma_j^2}} e^{-\frac{(v - \mu_j)^2}{2\sigma_j^2}}$$

Un'altra possibile opzione per trattare valori continui e' la discretizzazione da cui possiamo utilizzare Bernoulli/Multinomiale anche se bisogna stare attenti a non perdere informazioni discriminanti.

2.10.9 Un piccolo esempio basato su naive bayes

Supponiamo di dover usare Naive Bayes per predire se giocare una partita di tennis o no in base alle condizioni meteorologiche. Dati gli esempi qui sotto a sinistra, possiamo ricavare la tabella di distribuzione delle probabilita' generale come qui di seguito:

Weather	Play
Sunny	No
Overcast	Yes
Rainy	Yes
Sunny	Yes
Sunny	Yes
Overcast	Yes
Rainy	No
Rainy	No
Sunny	Yes
Rainy	Yes
Sunny	No
Overcast	Yes
Overcast	Yes
Rainy	No

Frequency Table		
Weather	No	Yes
Overcast		4
Rainy	3	2
Sunny	2	3
Grand Total	5	9

Likelihood table				
Weather	No	Yes		
Overcast		4	=4/14	0.29
Rainy	3	2	=5/14	0.36
Sunny	2	3	=5/14	0.36
All	5	9		
	=5/14	=9/14		
	0.36	0.64		

Figura 12: Esempi di partite di tennis giocate in base alle condizioni meteorologiche e tabella di distribuzione delle probabilit 

Ecco una rappresentazione grafica di *Naive Bayes*:

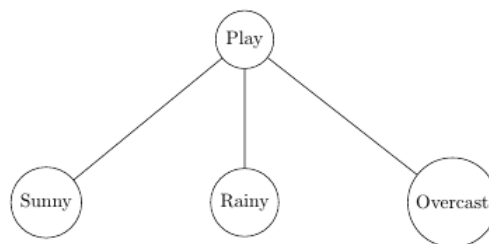


Figura 13

2.11 ALBERI DI DECISIONE

Gli alberi di decisione sono un altro tipo di apprendimento supervisionato, quindi dovremmo avere Da un punto di vista strutturale, l'albero di decisione   un albero (inteso come struttura dati) dove i nodi interni sono gli attributi, i rami tutti i possibili valori assumibili dall'attributo (oppure un range nel caso continuo) e le foglie sono la predizione da scegliere

D'ora in poi parleremo, per semplicit , solo di classificazione quando non espresso chiaramente.

2.11.1 *Costruzione di un albero di decisione*

L'albero di decisione si potrebbe definire prendendo a caso un attributo ed iniziare a dividere gli esempi fino ad ottenere foglie, anche se così otterremo un albero non molto utile in tutti i casi in cui non abbiamo gli esempi. Quindi quale albero scegliere fra tutti quelli possibili? In questo caso ci aiuta il rasoio di Occam, che ci dice di scegliere quello più piccolo fra tutti. Per generare l'albero più piccolo dovremmo scegliere gli attributi più significativi per generare l'albero. Cosa intendiamo per significativo? Intendiamo l'attributo che genera figli con meno varietà di classi presenti all'interno di essi.

Un esempio: immaginiamo di avere 10 esempi con attributi A e B e come valore associato un booleano. se abbiamo l'attributo A che genera due figli con esempi aventi meta' valore vero e falso mentre se suddividiamo secondo B abbiamo figli con esempi esclusivamente veri oppure falsi. In questo caso indubbiamente l'attributo più significativo è B.

2.11.2 *Classificazione dei nuovi elementi*

Quando dovremmo predire l'etichetta di un nuovo insieme di attributi x basterà semplicemente scorrere l'albero di decisione dalla radice fino ad una foglia, che sarà l'etichetta da assegnare.

2.11.3 *Un esempio di albero di decisione*

Riprendiamo l'esempio della partita di tennis, questa volta con qualche attributo in più. La tabella è la seguente:

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Figura 14: Tabella degli attributi meteorologici con valore associato un booleano che stabilisce se la partita e' stata giocata o no

da cui possiamo ricavare il seguente albero di decisione:

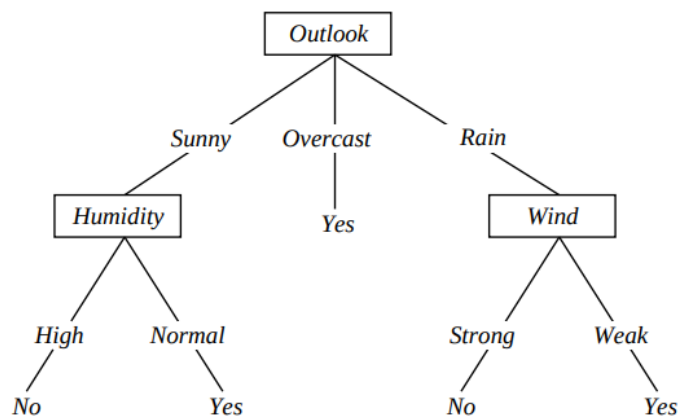


Figura 15: Albero di decisione ricavato dalla tabella precedente

2.12 RICERCA DELLA POSIZIONE

Dopo la costruzione del modello esso viene adoperato per predire la posizione corrente dell'utente che esegue l'applicazione. La ricerca consiste nel creare sul momento una *fingerprint* e poi, tramite un algoritmo di classificazione, cercare di inferire la *label* su cui ci troviamo. Gli algoritmi

utilizzati per analizzare il problema sono quelli descritti in precedenza quindi

1. *K Nearest Neighbour*
2. Bayes ingenuo
3. Alberi di decisione

STRUTTURA DEL SOFTWARE

3.1 LINGUAGGI E FRAMEWORK

L'applicazione ha come *target platform* Android perciò il linguaggio usato principalmente e' stato Java. Da notare che su Android e' presente nella versione 7, perciò non sono disponibili alcune funzionalita' di Java 8 come i metodi di default, *stream*, *lambda expression* ecc. Sono state usate delle librerie esterne fra i quali *Gson* che consente la serializzazione/deserializzazione dei dati, *lombok* che fornisce delle annotazioni che abbreviano il codice mantenendo una buona espressivita', varie librerie di supporto Android ed una libreria *Apache* che fornisce molte utilita' matematiche. Mosso dalla curiosita', ho provato ed alla fine utilizzato nel progetto *Kotlin*, un linguaggio che compila in *bytecode* per la JVM 100% interoperabile con *Java 6* (esuccessivi) che, fra le varie cose, fornisce le funzionalita' sopracitate mancanti a Java 7.

3.2 INTERFACCIA GRAFICA

L'interfaccia e' stata realizzata ai fini di test pratici delle funzionalita' finali dell'applicazione quindi non ha una grande cura da un punto di vista estetico come vedremo piu' avanti. La parte alta contiene delle label raffiguranti i 3 valori catturati tramite il magnetometro e presi tramite le API di Android.

Nel mezzo ci sono dei pulsanti per:

- Iniziare/Terminare la scansione dell'ambiente.
- Incrementare la *label* che verra' assegnata alla prossima *fingerprint* registrata
- Iniziare/Terminare la ricerca.

- Serializzare tutti i dati registrati finora
- De-serializzare i dati salvati in un JSON.

Nella parte bassa invece c'è una *textbox* contenente il log che verrà stampato durante l'esecuzione del programma.

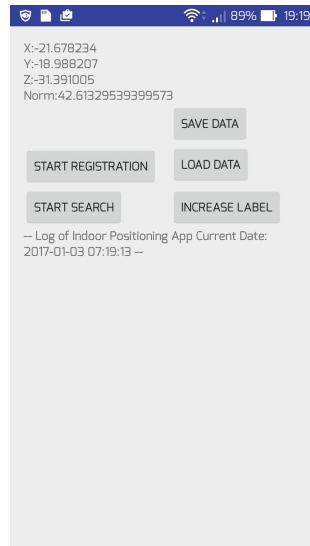


Figura 16: Interfaccia grafica all'avvio dell'applicazione

3.3 PERSISTENZA DEI DATI

L'applicazione consente anche di serializzare tutti i dati registrati fino a quel momento nel formato standard JSON tramite la libreria *gson* che fornisce delle funzioni per il linguaggio Java per la serializzazione/de-serializzazione di oggetti. Il file viene salvato nella cartella dati dell'applicazione non visibile all'utente.

3.4 STRUTTURA DEL CODICE E DESIGN PATTERN

Nello sviluppo del software sono stati applicati vari *design pattern* visti durante i vari corsi e principi di programmazione. Fra questi ultimi abbiamo il *dependency inversion principle*, *open closed principle*. Riguardo i *design pattern*, sono stati usati frequentemente l'*observer*, il *template* e *factory*.

3.5 ANALISI DEI DATI

La predizione dei risultati e' stata implementata sia nel software *Android* sia sul computer. Nel codice mobile e' stato adoperato solamente il KNN per via della facilita' d'implementazione da zero e non e' stato utilizzato per testare la precisione, ma per verificare il corretto funzionamento dell'applicazione. Invece su computer, presi i dati serializzati dal software mobile, sono stati applicati tutti gli algoritmi di apprendimento elencati precedentemente e gia' tutti implementati da librerie di terze parti per verificare la precisione dei dati. Il linguaggio scelto su computer e' *Python* per via del suo buon supporto all' apprendimento automatico tramite la libreria *sklearn*.

3.6 CODICE USATO PER L'ANALISI

```
def train_model(dataset: list, model, cross_validation=True):
    """
    :param dataset: List of tuples (label, features).
    :param model: a model instance
    :param cross_validation: flag for cross validation
    :return: the model and the test error predictions. In case of
             cross validation the mean of all the errors
    """
    y, X = zip(*dataset)
    X_train, X_test, y_train, y_test = train_test_split(X, y,
                                                         test_size=0.3)
    model = model.fit(X_train, y_train)
    test = list(zip(X_test, y_test))
    return model, test_model(model, test, cross_validation)
```

```
add = lambda a, b: a + b
```

```
def test_model(model: {'predict'}, test: list, cross_validation):
    if cross_validation:
        X_test, y_test = zip(*test)
        scores = cross_val_score(model, X_test, y_test, cv=10,
                                  scoring='accuracy')
```

```

        return 1 - scores.mean()
    else:
        predictions = [(model.predict([X_i])[0], y_i) for X_i, y_i in
                        test]
        wrong_predictions = reduce(add, [1 for prediction, real_label
                                         in predictions if prediction != real_label], 0)
        return wrong_predictions / float(len(test))

```

Visto che ogni classificatore condivide i metodi per l'addestramento e la predizione con tutti gli altri, ho creato una funzione generica che dato in input il dataset, un'istanza del classificatore ed un flag per la *cross validation*, ritorni la percentuale d'errore sull'insieme di test.

Un semplice uso della funzione precedente e':

```

decision_tree, test_error = train_model(dataset, DecisionTreeClassifier())
print("Test error with decision tree = %d"% test_error)

```

TEST

4.1 BASE DEI TEST

Per testare l'effettivo funzionamento dell'applicazione ho usato alcune stanze di casa mia ed ho assegnato a ciascuna di esse una *label*. Qui di seguito una piccola piantina rappresentante le stanze utilizzate:

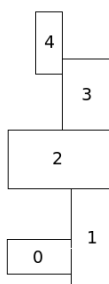


Figura 17: Una piccola raffigurazione delle stanze usate per le prove con sopra scritto la *label* assegnata

Sono stati raccolti circa 18000 campioni di onde magnetiche. La suddivisione fra addestramento e test e' 70/30.

4.2 CARATTERISTICHE

Dai grafici qui di seguito possiamo notare che c'è sovrapposizione fra i dati, quindi è presente del rumore nei dati.

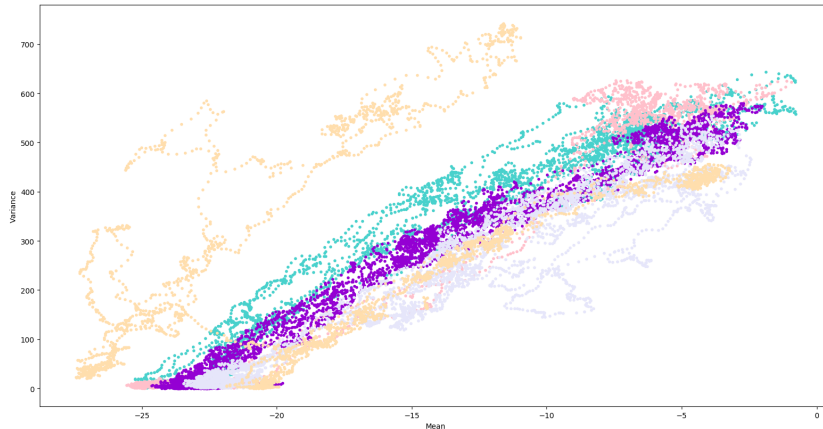


Figura 18: Grafico in 2 dimensioni della media e varianza di tutte le onde magnetiche. I colori dei punti rappresentano le etichette

La causa del rumore sono i sensori che offrono una misurazione non precisa. Un possibile riparo al problema è il *filtro di Kalman*

4.3 CLASSIFICATORI A CONFRONTO

Qui di seguito vediamo i risultati ottenuti da ciascun classificatore con un istogramma:

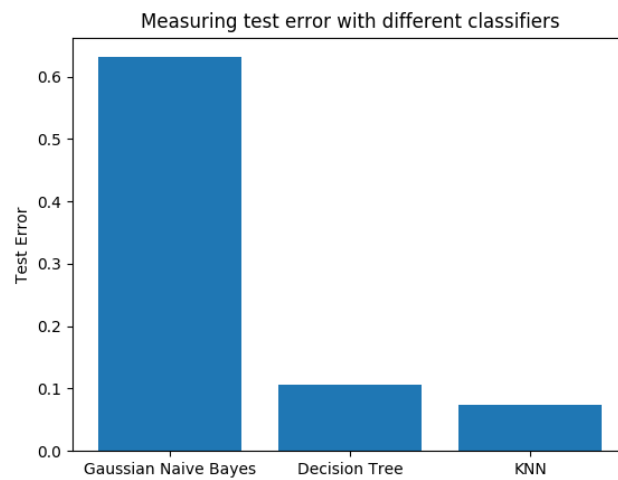


Figura 19

come possiamo notare *Gaussian Naive Bayes* e' totalmente inadatto alla classificazione di onde magnetiche mentre gli alberi di decisione e *K Nearest Neighbours* si comportano molto bene, con risultati leggermente migliori in quest'ultimo. A questo punto qualcuno potrebbe pero' pensare che gli ultimi 2 modelli si sono sovradattati agli esempi (*overfitting*) ed avrebbe ragione, perche' applicando la *cross validation* abbiamo risultati diversi da quelli precedenti.

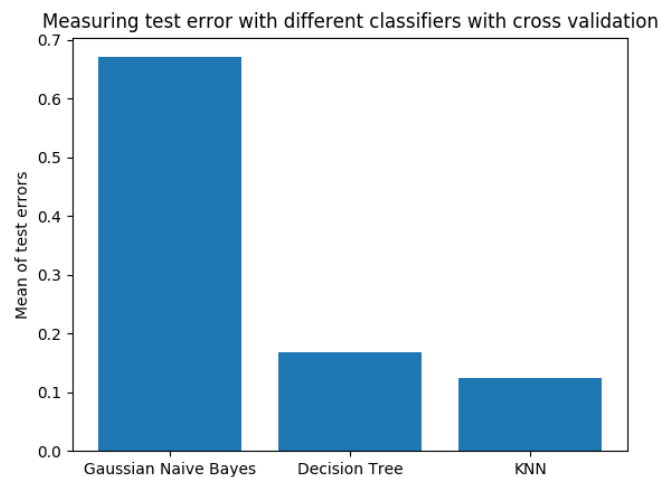


Figura 20

Adesso visualizziamo gli errori per etichetta:

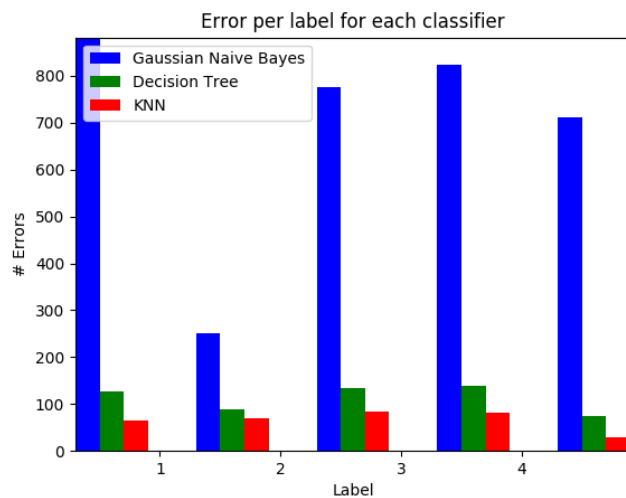


Figura 21: Numero di errori nella predizione per etichetta

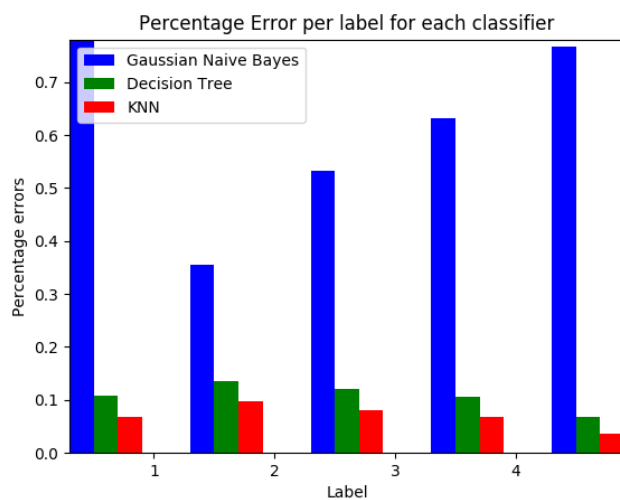


Figura 22: Percentuale di errore nella predizione per etichetta

4.4 UN RIMEDIO INGENUO AL RUMORE

Un approccio ingenuo per risolvere il problema al rumore potrebbe essere quello di prendere meno dati per etichetta. Il seguente grafico però ci mostra che ciò non è vero

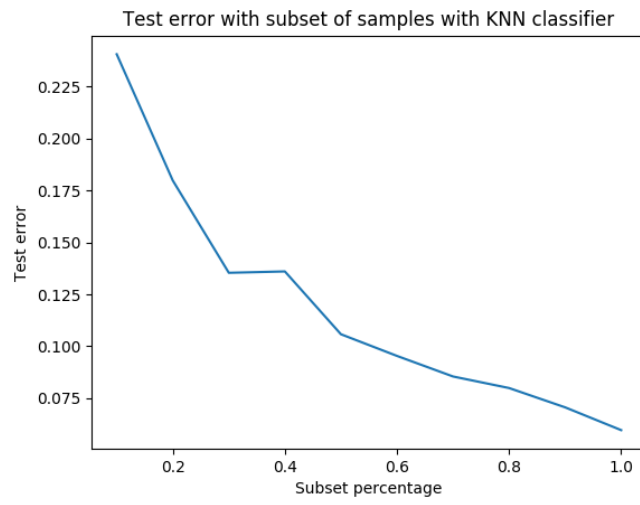


Figura 23

MIGLIORAMENTI

5.1 POSSIBILI MIGLIORAMENTI PER LA RACCOLTA DI DATI

- L'architettura *client-server* e' necessaria per avere un'applicazione che riesca a gestire con fluidita' la ricerca all'interno di grandi edifici. In riferimento alla memoria del telefono poiche' nella maggior parte dei casi possiede un quantitativo di *RAM* e memoria permanente molto limitati, ma anche rispetto al processore perche' la ricerca della posizione, assumendo di usare il *KNN*, ha bisogno di confrontare i propri attributi con tutti gli altri e va da se che piu' e' grande il dataset, piu' potenza di calcolo ci vorra' per ottenere una risposta in tempi umani.
- Durante la raccolta sarebbe opportuno applicare il *filtro di Kalman* per ridurre il rumore causato dall'imprecisione dei sensori.
- Per migliorare la precisione sarebbe opportuno appoggiarsi anche ad altri sensori presenti sullo *smartphone*. Prendiamo come esempio il Wi-Fi, se il dispositivo e' connesso alla rete dell'edificio potrebbe sfruttare la potenza di segnale per avere una precisione maggiore; sfruttare l'accelerometro per stimare la velocita' del telefono e quindi standardizzare tutte le rilevazioni effettuate col magnetometro. Questa tecnica viene chiamata *sensor fusion*([Link a qualche paper](#))