



POLITECNICO DI TORINO

Corso di Laurea in Ingegneria Informatica

Tesi di Laurea

Predizione di difettosità nello sviluppo software attraverso machine learning

Artificial Intelligence applied to Software Engineering

Relatore

prof. Maurizio MORISIO

Candidato

Jacopo NASI [255320]

Supervisore Aziendale

dott. Davide PIAGNERI

APRILE 2020

Sommario

Ogni giorno migliaia di commit vengono eseguiti, ognuno di loro contiene molte informazioni: file modificati, modifiche, commenti, registri di test e molto altro. Una strutturata e corretta gestione delle piattaforme di concontrollo sorgente permette l'estrazione di dati utili analizzabili utilizzando modelli statistici di intelligenza artificiale.

Al fine di poter correttamente utilizzare questi dati sono necessari alcuni step preliminari: la prima fase riguarda l'analisi della struttura dati al fine di permettere l'estrazione di tutte le possibili informazioni, successivamente la pre-elaborazione per rimuovere informazioni di inutili e di disturbo, con i dati puliti è possibile procedere con l'estrazione di dati combinati, come la seniority degli sviluppatori, una lista di parole dei componenti modificati, la versione ed altre informazioni di carattere più matematico. L'ultima fase prevede la sostituzione dell'etichetta testuale relativa alla priorità con un valore numerico corrispondente al valor medio della distribuzione della durata di quella etichetta, questo valore prenderà il nome di severity. I dati verranno poi aggregati per settimana. Una volta generati i dati verranno utilizzati per allenare tre differenti modelli: Random Forest, Gradient Boosting e Reti Neurali. L'allenamento sarà gestito in tre differenti modalità: la prima allena e predice utilizzando lo stesso filone di dati, la seconda, cross-version, prevede che il modello venga allenato su dati relativi ad alcune versione del progetto per poi effettuare la predizione sulle successive, la terza, cross-project, allena il modello con dati relativi ad un progetto per poi prevedere l'andamento di uno

differenti.

Tutte le tipologie ottengono dei buoni risultati, il migliore è quello cross-project che riesce ad ottenere una precisione maggiore del 90% fino a quattro settimane e comunque maggiore del 70% fino a 20 settimane.

Ringraziamenti

Un ringraziamento speciale a Smirnuff ed i suoi cavalieri, luce della mia battaglia.

Indice

1	Introduzione	7
1.1	Problema Generale	7
1.2	Strumenti utilizzati	8
2	Stato dell'arte	11
2.1	Lavori correlati	11
3	Dati	13
3.1	SEOSS33	13
4	Apprendimento Automatico	19
4.1	Introduzione	19
4.2	Apprendimento ensamble	21
4.3	Reti Neurali	23
4.4	Metriche di valutazione	31
5	Pre-elaborazione	33
6	Predizione	35
7	Conclusioni	37
	Bibliografia	39

Capitolo 1

Introduzione

1.1 Problema Generale

Lo sviluppo software non si presenta molto differente dallo sviluppo di qualsiasi altro prodotto, dopo una fase iniziale di progettazione lo sviluppo del codice può avere inizio, durante esso emergeranno sistematicamente dei problemi che dovranno essere risolti prima della consegna della versione finale.

Ogni progetto software è costituito da diversi commit per giorno, ognuno di essi contiene innumerevoli informazioni le quali possono essere utilizzate per analisi statistiche. La predizione della difettosità può migliorare enormemente il processo di sviluppo, allocando un corretto numero di sviluppatori per risolvere le problematiche e riducendo quindi le tempistiche per la correzione. Anche il machine learning può essere utilizzato per la predizione dei difetti.

La predizione è uno strumento sempre più utilizzato a livello industriale, un corretto utilizzo può generare enormi benefici a livello produttivo, permettendo la riduzione di sprechi, l'ottimizzazione delle vendite e tante altri vantaggi. Lo sviluppo di progetti di natura informatica è sempre di più centrale all'interno della nostra società attuale, anche questo processo potrebbe trarre beneficio dai vantaggi della predizione. L'implementazione di tecniche statistiche viene in supporto, vista la natura

intellettuale della programmazione, nella generazione di predizioni utili.

1.2 Strumenti utilizzati

Lo sviluppo di questo progetto a richiesto l'utilizzo di diversi strumenti, di seguito una lista degli stessi:

Python Il linguaggio di programmazione principale, utilizzato per la gestione dei dati, l'estrazione di informazioni, l'applicazione di algoritmi matematici e l'interazione con altri software. Nello specifico la versione utilizzata è stata la v3.7.0

Pandas Libreria open source ad alte prestazioni, con semplici strutture e strumenti adatti all'analisi dati attraverso Python.

NumPy Libreria per il calcolo scientifico attraverso Python.

Matplotlib Libreria per il disegno di grafici 2D in Python.

Seaborn Libreria avanzata per il disegno 2D in Python.

Tensorflow Piattaforma per machine learning.

Keras API di alto livello per reti neurali.

SciKit-Learn Strumenti e librerie per machine learning.

GitLab Piattaforma di sourcing basata su Git. Utilizzata per il codice sorgente del progetto.

GitHub Piattaforma di sourcing basata su Git. Utilizzata per il calendario e elaborato testuale:

- Tesi: <https://github.com/Jacopx/Thesis>
- Calendario: <https://github.com/Jacopx/ThesisCalendar>

JetBrains IDEs IDE per lo sviluppo di diversi linguaggi di programmazione, gratuita per gli studenti:

- PyCharm: <https://www.jetbrains.com/pycharm/>
- DataGrip: <https://www.jetbrains.com/datagrip/>

Capitolo 2

Stato dell'arte

2.1 Lavori correlati

Parlando di altri lavori su simili tematiche.

Capitolo 3

Dati

Le seguenti sezioni analizzeranno le basi di dati utilizzate in questo progetto.

3.1 SEOSS33

SEOSS33[1] è una [base dati](#) collezionante errori, issue e tante altre informazioni a proposito di 33 progetti open source. I dati sono stati tutti collezionati estraendo le informazioni dalle piattaforme di controllo del codice sorgente, Version Control System (VCS), come GitHub e dalle piattaforme per la gestione dello sviluppo, Issue Tracking System (ITS), come Jira di Atlassian.

Ad oggi nessun altro progetto di ricerca, su questi dati, è stato effettuato.

Ogni progetto prevede una propria linea durante la fase di sviluppo, tutte le metodologie e linee guida sono alla base degli studi di ingegneria del software. Tuttavia è possibile unificare ed accorpare secondo una categorizzazione standard molte delle differenze specifiche. Lo sviluppo della base dati SEOSS33 mira proprio alla creazione di una serie di dati generalizzati e fruibili attraverso medesime procedure senza la necessità di adattarsi alle specifiche caratteristiche di ogni singolo progetto.

Il mondo open source presenta una quantità pressoché infinita di differenti software, parte del progetto in questione è stata dedicata alla selezione dei software da

analizzare per l’inserimento nella base dati condivisa, per questo motivo sono state definite alcune caratteristiche che accomunassero i vari progetti in modo da costituire una base: discretamente omogenea a livello di dimensionalità, ma con differenze strutturali utili per successive analisi come quella relativa a questo progetto. Il requisito principale riguardava il linguaggio di programmazione, considerare progetti sviluppati per la maggior parte in un singolo linguaggio di programmazione permette di ridurre la variabilità interna ad ogni singolo progetto. Vista la natura di analisi attraverso il machine learning, un’altra importante caratteristica riguardava il numero di issue, il quale doveva essere sufficientemente elevato. I progetti, oltre a dover essere attualmente in sviluppo, dovevano presentare un età di almeno 3 anni. La definizione di tutti questi parametri a permesso di generare una base dati contenente 33 progetti simili come struttura ma con caratteristiche differenti.

Lo sviluppo del presente progetto di tesi si è concentrato solamente su cinque di questi schemi, sono stati scelti i progetti più grossi e quelli in sviluppo dal maggior tempo, nello specifico i selezionati sono riportati in tabella 3.1:

Tabella 3.1. Distribuzione dati

Progetto	Mesi	Issue
Hadoop	150	39086
Hbase	131	19247
Maven	183	18025
Cassandra	106	13965
Hive	113	18025

Al fine di generalizzare le specifiche differenze, le varie issue: *New Feature*, *Bug Report*, ecc... Sono state mappate su cinque categorie:

- Bug: Un problema che previene il funzionamento del prodotto
- Feature: Una nuova funzionalità del prodotto
- Improvement: Un miglioramento di una funzionalità già esistente

- Task: Un compito necessario
- Other: Vario

La figure 3.1 visualizza la distribuzione, rispetto le varie categorie, del numero di issue per ogni progetto.

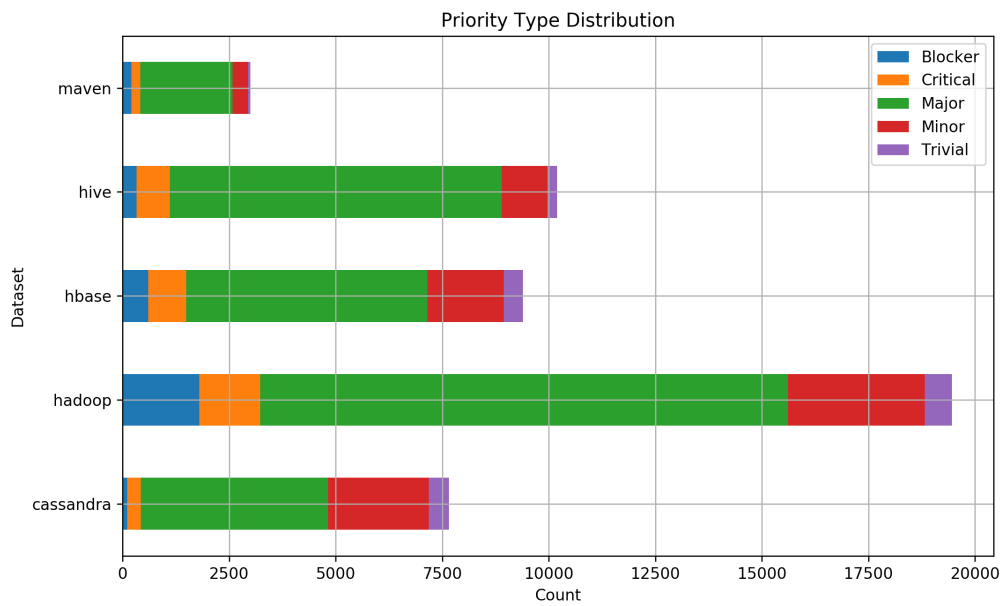


Figura 3.1. Distribuzione issue per progetto

Per poter estrarre ed utilizzare al meglio i dati è necessario conoscere al meglio la struttura contenitore. I dati relativi ad ogni software sono salvati in un file SQLITE, un database SQL offline che permette l'accesso sfruttando le potenzialità delle query, senza la necessità di un server vero e proprio. La figura 3.2 riporta lo schema integrale della struttura.

Tutto il modello si basa sulla sua entità centrale, la issue, ovvero l'attività di segnalazione che è stata creata da uno sviluppatore per gestire una problematica. Ognuna di queste issue è caratterizzata dal proprio *issue_id* il quale ne rappresenta la chiave primaria ed univoca, normalmente è strutturata con il nome del

progetto seguito da un numero progressivo. La tabella relativa alle issue contiene ulteriori informazioni direttamente correlate, la tipologia, la priorità, le informazioni temporali di apertura, aggiornamento e chiusura della stessa, un breve riassunto della problematica, lo stato e le informazioni relative allo sviluppatore che l'ha aperta. Direttamente collegate, tramite la chiave primaria, vi sono le tabelle contenenti i commenti *issue_comment*, la versione *issue_fix_version*, il componente modificato *issue_component* ed la tabella *change_set_link* la quale collega i vari commit alle issue. Durante l'estrazione delle varie informazioni sono state utilizzate tutte le tabelle ad esclusione di *issue_link* la quale viene utilizzare per correlare le differenti issue tra di loro.

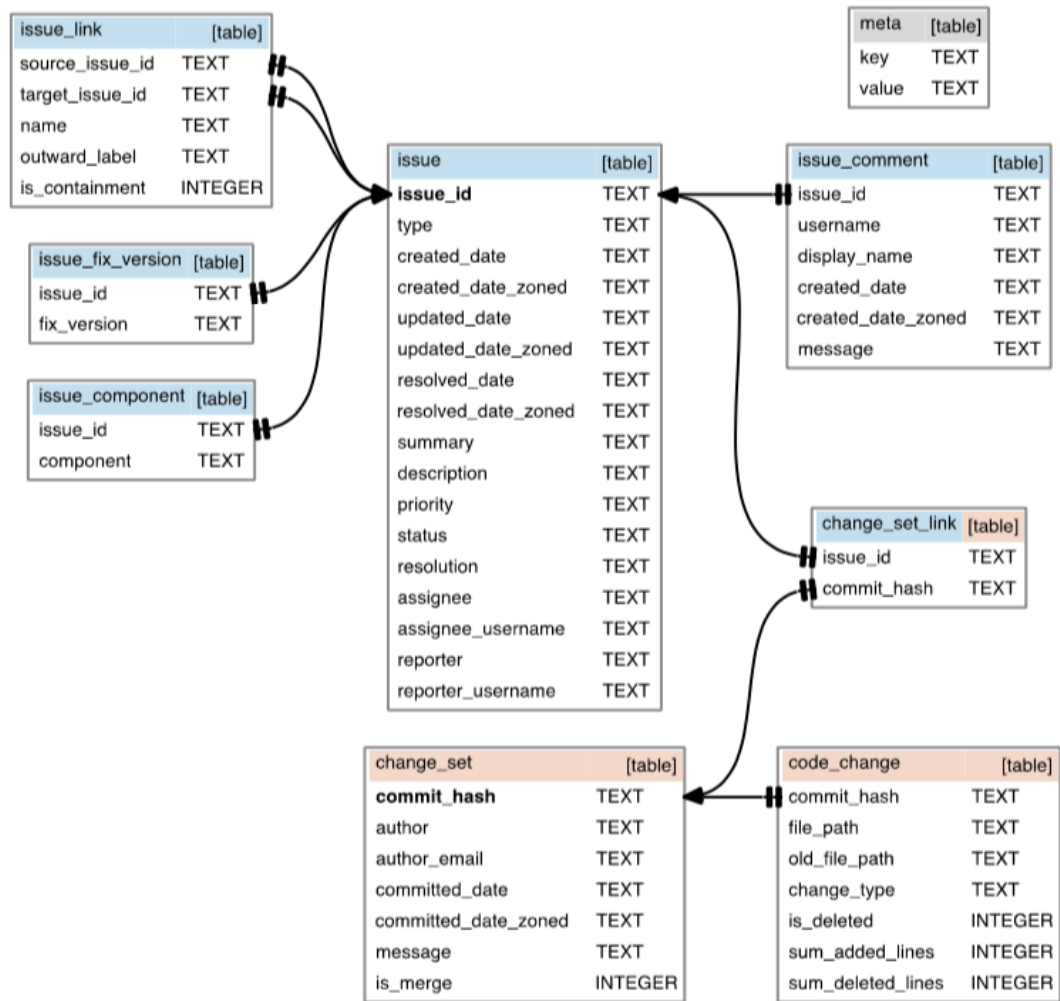


Figura 3.2. Struttura dati di SEOSS33

Capitolo 4

Apprendimento Automatico

4.1 Introduzione

L'apprendimento automatico, meglio conosciuto come Machine Learning (ML), è una branca dell'intelligenza artificiale basata sullo studio di algoritmi e modelli statistici utilizzabili dai calcolatori per svolgere determinati compiti senza essere esplicitamente istruiti per farlo. Questo settore è ormai diventato di dominio pubblico, solo negli ultimi decenni l'utilizzo di queste tecniche è cresciuto enormemente nonostante la maggior parte di esse furono teorizzate già molti anni prima. La motivazione principale di questo ritardo è da ricercare nella natura stessa di queste strategie, la capacità computazionale diventa rilevante e fondamentale all'applicazione degli stessi, grazie alla crescita di essa è ora possibile sfruttare questi algoritmi anche nei computer di casa.

Tutti gli algoritmi possono essere catalogati in una delle seguenti categorie:

- Knowledge-based: Acquisizione e modellazione di leggi conosciute (dalle regole ai fatti)
- Learning: Estrazione della conoscenza e delle regole attraverso esempi ed esperienza (dai fatti alle regole)

Tutti gli algoritmi in ambito ML fanno parte della seconda categoria. A loro volta questi algoritmi possono essere divisi in tre principali sotto-categorie: apprendimento supervisionato, apprendimento non supervisionato e apprendimento per rinforzo. Nel primo vengono forniti modelli di dati in ingresso e i dati desiderati in uscita e lo scopo è quello di definire una regola che associ i due parametri. Nel modello non supervisionato il modello ha il compito di trovare una struttura ai dati in ingresso, senza che essi siano precedentemente etichettati in alcun modo. L'ultimo invece viene allenato per un compito, senza che gli venga insegnato come fare ma solamente conoscendo il risultato finale delle proprie azioni.

Esistono una varietà enorme di modelli di questa tipologia, i successivi paragrafi tratteranno quelli utilizzati in questo progetto.

Apprendimento Supervisionato è una tecnica che prevede il processamento dei dati in ingresso con associati i valori desiderati in uscita, lo scopo del modello è quindi quello di sviluppare una correzione matematica tra tutte le informazioni che riceve in ingresso ed i valori desiderati in uscita. Una volta terminata la fase di allenamento il modello potrà essere utilizzato per prevedere il valore di uscita dati i valori in ingresso. Questa metodologia può essere applicata nella risoluzione di problemi di due differenti categorie, quelli della classificazione e quelli della regressione lineare. Lo scopo del primo è quello di assegnare una etichetta ai dati per classificarli in diverse categorie, per esempio le transazioni sane o fraudolente di una banca. L'assegnazione può essere binaria, quindi solo due etichette, o multi-etichetta. La regressione invece ha l'obiettivo di predire un valore continuo di uscita, cercare di sviluppare una relazione matematica tra tutte le variabili in ingresso, cercando di prevedere, con il miglior livello di approssimazione il valore finale. Nel nostro progetto verranno solo impiegati questi ultimi, la classificazione non verrà ulteriormente trattata.

4.2 Apprendimento ensemble

L'apprendimento di insieme raggruppa una serie di tecniche sviluppate al fine di migliorare i risultati dei singoli predittori. Invece che utilizzare un singolo modello, nella fase di apprendimento, vengono simultaneamente allenate diverse copie dello stesso modello con parametri differenti, ciò porterà ad una differenziazione delle risultato di previsione. L'aggregazione, attraverso tecniche come bagging, boosting o stacking, permetterà di produrre un risultato più accurato e meno dipendente dalla rumorosità dei dati.

Foresta casuale conosciuta anche come Random Forest (RF) è un algoritmo di apprendimento supervisionato, basato sulle metodologie d'insieme, per la classificazione e la regressione. È costituito combinando la predizione di diversi alberi, ognuno allenato separatamente, tramite media [2]. In figura 4.1 una visualizzazione del modello.

La definizione di una foresta richiede tre parametri principali: (1) la metodologia per la divisione in foglie, (2) il tipo di predittore da usare in ciascuna foglia e (3) il metodo per garantire la randomicità. La divisione in foglia richiede la selezione della forma e metodologia per la valutazione di ogni candidato. Una tipica scelta è quella chiamata *axis-aligned*, dove i dati vengono diretti nei vari sotto alberi in base al passaggio o meno di un valore soglia, il quale viene scelto casualmente o dalla funzione di ottimizzazione della foglia. Al fine di dividere una foglia vengono generati diversi candidati e viene definito un criterio per scegliere tra essi. Un primo approccio potrebbe essere quello della selezione causuale uniforme, altrimenti la scelta può essere guidata da una funzione di purezza, cercandone la massimizzazione.

Possono essere utilizzate diverse tecniche per generare casualità nella foresta, attraverso la definizione di soglie senza l'utilizzo di funzioni oppure effettuando con

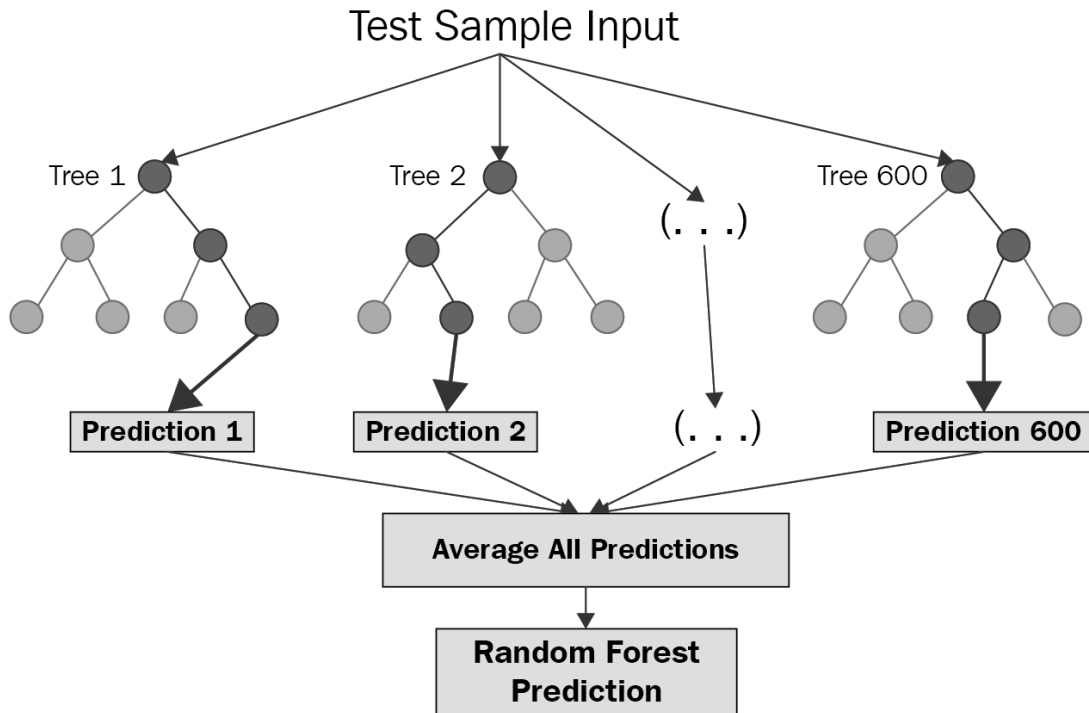


Figura 4.1. Schema semplificato di foresta casuale [3]

l'allenamento di ogni albero su una selezione di dati ristretta in modo da diversificare direttamente il risultato a livello di insieme.

La fase di training viene gestita indipendentemente da ogni albero attraverso punteggi di struttura e stima, i primi permettono la variazione della forma dello stesso, mentre i secondi guidano le funzioni di ottimizzazione delle singole foglie.

Un volta effettuato l'allenamento della rete è possibile utilizzare il modello per la predizione dei valori. Nella fase di stima, ogni singolo albero, genererà indipendentemente un proprio valore, la scelta finale avverrà calcolando la media aritmetica di tutti questi valori generati, il contributo è equamente ripartito tra tutti.

La nostra implementazione sfrutta le API per la Random Forest di SciKit-Learn v0.21:

```
from sklearn.ensemble import RandomForestRegressor
```

Gli specifici parametri utilizzati verranno illustrati durante il capitolo 6 sulla previsione.

Macchine ad aumento di gradiente – conosciute anche come Gradient Boosting Machines (GBM), sono una famiglia di potenti modelli statistici di apprendimento automatico in grado di ottenere ottimi risultati in una grande varietà di applicazioni. Una delle loro principali caratteristiche è la possibilità di personalizzare il modello in base alle caratteristiche dell'applicazione [4]. Tecniche come la foresta casuale, appena trattata, sono basate sulla semplice media dei risultati prodotti da ogni singolo componente. La famiglia dei metodi di aumento è basata su una differente strategia di unione dei pezzi per la formazione della modello finale. Il boosting aggiunge, sequenzialmente, nuovi parti all'insieme; durante la fase di allenamento vengono via via sviluppati nuovi piccoli modelli da aggiungere al fine di migliorare l'accuratezza nella previsione. Idealmente vengono costruiti nuovi modelli di base, come per esempio l'albero decisionale, per poter massimizzare la correlazione con il gradiente negativo della funzione di perdita (*loss*).

Vista l'alta flessibilità del modello, l'adattamento dello stesso a differenti ambienti non risulta difficoltoso, molte differenti sperimentazioni possono essere fatte.

Nel nostro progetto si è deciso di implementare il modello di Gradient Boosting Decision Tree (GBDT) sempre utilizzando la libreria SciKit-Learn v0.21:

```
from sklearn.ensemble import GradientBoostingRegressor
```

4.3 Reti Neurali

Le reti neurali, in inglese Neural Networks (NN), sono modelli di apprendimento automatico con diretta ispirazione al cervello umano e come esso procede alla fase di apprensione di un concetto, la rete è costituita dalla basilare unità di calcolo, il neurone (neuron), collegata ad altri neuroni attraverso le sinapsi (synapses). La

conoscenza è data alla rete, nella fase di allenamento, attraverso esempi, la forza delle connessioni inter neurali è la base per acquisire e mantenere la conoscenza.

La fase di apprendimento può essere sia supervisionata che non supervisionata. La modalità supervisionata viene utilizzata per il riconoscimento di schemi (pattern recognition) e regressione e viene effettuata sempre con i dati di input ed i desiderati dati di output. Invece, la modalità non supervisionata, è maggiormente utilizzata per sviluppare modelli adatti al raggruppamento (clustering) e l'allenamento viene effettuato senza il valore desiderato. Il nostro progetto farà uso di reti neurali per la regressione.

Questa tipologia di reti può essere di tre tipologie:

- Singolo livello flusso in avanti
- Multi livello flusso in avanti
- Ricorsiva

L'architettura standard è composta di tre diversi livelli, figura 4.2, lo strato di ingresso, le unità nascoste e il livello di uscita; tutti questi livelli sono correlati tra loro tramite le connessioni sviluppate durante la fase di apprendimento. Il neurone è l'unità basilare per il processamento all'interno della rete, si occupa di ricevere i dati in ingresso, gestirli e poi passarli ai successivi livelli. Ogni ingresso combina i dati con il proprio stato interno e la funzione di attivazione per poi procedere a passare il valore come ingresso del livello successivo. L'importante che ognuno di questi valori in ingresso avrà sarà determinata dal peso assegnato alla connessione durante la fase di allenamento della rete stessa. Ogni nodo ha la possibilità di ricevere più di un ingresso, per questo motivo, tutti i valori verranno aggregati in modo da consegnare un solo valore come ingresso del successivo strato, la formula per il calcolo della somma è:

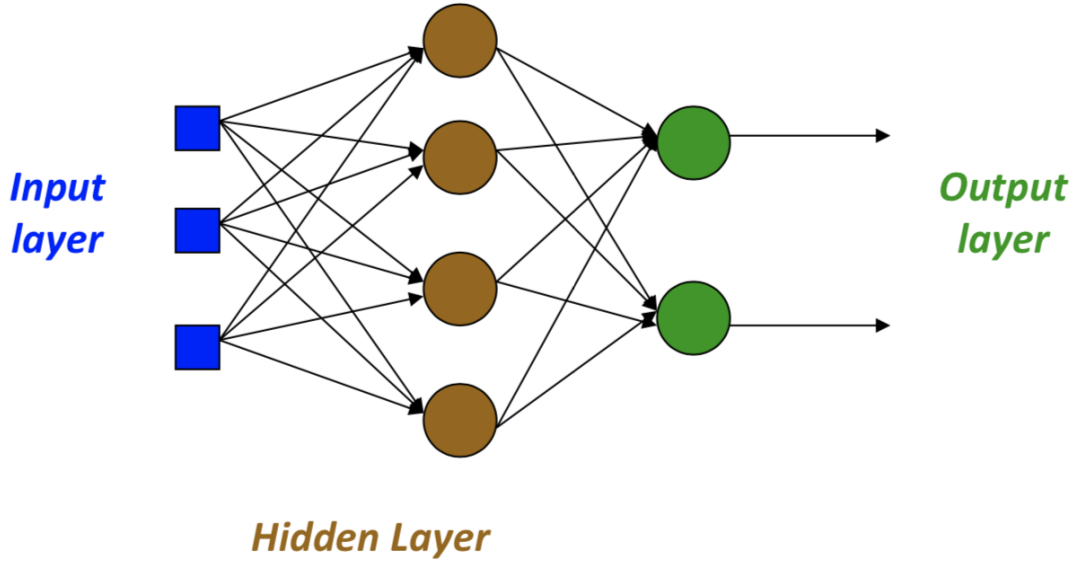


Figura 4.2. Rete a flusso avanti multi livello

$$u = \sum_{j=1}^m w_j x_j \quad (4.1)$$

Il valore calcolato viene scalato tramite una funzione di attivazione φ al fine di limitarne l'ampiezza:

$$y = \varphi(u + b) \quad (4.2)$$

La precedente funzione riporta il parametro b il quale rappresenta il bias, un parametro esterno del neurone. y rappresenta invece il valore di uscita dopo la computazione, il quale rappresenta il valore di ingresso del successivo livello gerarchico. Un esempio della struttura in questione si può trovare in figura 4.3.

Diverse funzioni di attivazione possono essere applicate al neurone, il loro compito è quello di emulare la tipica risposta biologica del sistema nervoso umano e le sue differenti metodologie di attivazione. Nel corso degli anni sono state definite

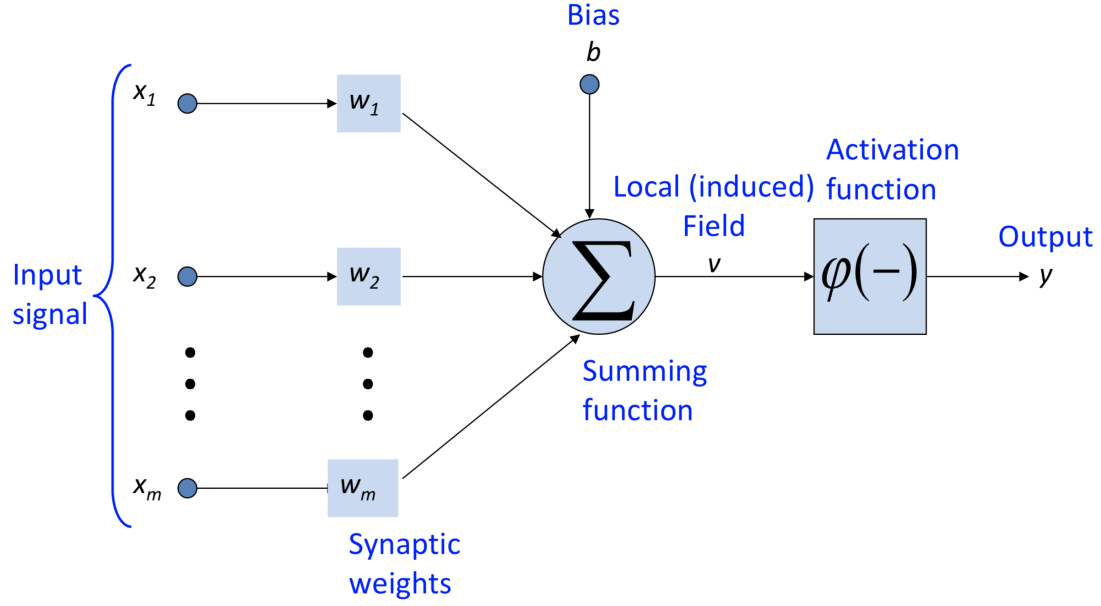


Figura 4.3. Visualizzazione di neurone

numerosi differenti funzioni, più o meno adatte a differenti contesti, con proprie peculiarità, problematiche e caratteristiche. Possono essere classificate in due categorie, lineari e non lineari, le più comuni sono: lineare, gradino, relu e sigmoide.

Gradino è una delle più comuni funzioni di attivazione, binaria, lineare e basata su soglia, in figura 4.4 la sua definizione. Quando il valore in ingresso è sopra o sotto la soglia definita, il neurone viene attivato e passa il valore in ingresso al successivo livello. La principale problematica correlata a questa funzione è la sua impossibilità di gestire valori multipli in uscita.

Lineare è una funzione di attivazione lineare della forma:

$$f(x) = x \quad (4.3)$$

La funzione, dato il valore in ingresso e moltiplicandolo per il peso del neurone, calcola il valore di uscita. Rispetto alla funzione gradito possono essere generati output multi valore, presenta comunque due problematiche: non sarà possibile utilizzare la propagazione indietro (trattata successivamente) per allenare la rete, vista la funzione derivata costante; l'altro problema riguarda il collasso di tutto i diversi livelli in uno solo, vista la sua natura lineare, il valore di uscita finale sarà in ogni caso la combinazione lineare tutti i livelli precedenti. La figura 4.5 visualizza la curva in questione.

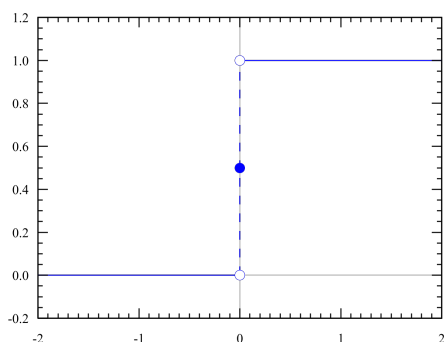


Figura 4.4. Funzione a gradino

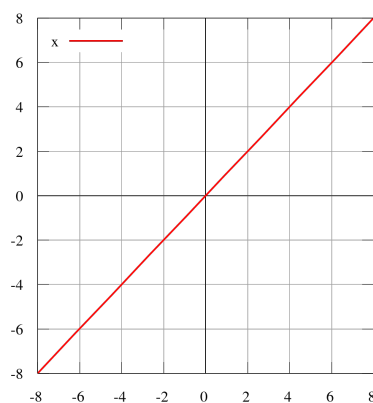


Figura 4.5. Funzione lineare

Sigmoide è la prima funzione di attivazione non lineare trattata, nello specifico è caratterizzata dalla seguente equazione:

$$f(x) = \frac{1}{1 + e^x} \quad (4.4)$$

La peculiare caratteristica di non linearità permette un più morbido gradiente in modo da prevenire valori vuoti scala, normalizzando il valore tra $[0, +1]$ si ottengono anche benefici a livello di pulizia dei dati in ingresso utilizzati successivamente per le previsioni. La funzione non si presenta esente da problematiche, la principale riguarda la vanificazione del gradiente, se da un lato permette di smorzare i

valori fuori scala, si tramuta in collo di bottiglia in altri casi, in caso di valori in ingresso molto elevati o molto bassi non vi sarà differenziazione nel valore di uscita. Inoltre l'applicazione del calcolo stesso è decisamente più impegnativa a livello computazione. La figure 4.6 descrive la curva in questione.

ReLU il quale acronimo sta per Rectified Linear Unit, unità lineare rettificata, è definita nella seguente maniera:

$$f(x) = \max(0, x) \quad (4.5)$$

Nonostante assomigli molto alla funzione di attivazione lineare, presenta una funzione derivate che permette la propagazione inversa e si presenta molto efficiente a livello computazione. Le problematiche si presentano in caso di valori in ingresso prossimi allo zero o addirittura negativi, il gradiente della funzione diventa nullo e la rete non potrà effettuare la propagazione inversa e conseguentemente non potrà portare avanti il processo di apprendimento.

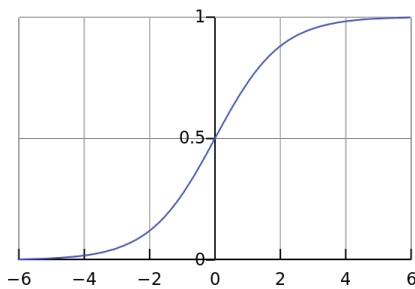


Figura 4.6. Sigmoide

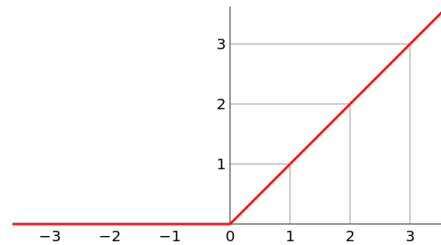


Figura 4.7. ReLU

Delta learning rule uses the difference between target output and the obtained activation to drive the learning. According to these rules, each time an output is computed, the weight of the neuron is adjusted according to an error function trying to reduce the difference between output and target.

Backpropagation short for "*backward propagation of errors*", is an algorithm for supervised learning of artificial neural networks using gradient descent. Given an artificial neural network and an error function, the method calculates the gradient of the error function concerning the neural network's weights. It is a generalization of the delta rule for perceptrons to multilayer feedforward neural networks. [5].

The main characteristic of this technique is that the gradient, proceeds backward through the network, with the gradient of the final layer calculated before the first. This solution allows an efficient calculation of the gradient for each layer. The algorithm is structured in this form:

1. Compute the error term for the output units
2. From the output layer, repeat until last hidden reached:
 - (a) Propagating the error term back to the previous layer
 - (b) Updating the weights between the two layers

BP suffers the Vanishing Gradient problem, more layers are embedded in the network and more difficult become to train it. Due to the nature of the backpropagation, when an output is generated, the weight of the neurons is fixed according to the rule, each layer back the correction value is smaller, following the derivate of the activation function, in case of shallow network this problem is solved easily, for deep network this can become a real issue. The ReLU activation function manages better this problem. Another solution is the batch normalization, that rescale the input value between $[-1,1]$ in order to keep the input value far from the activation function edges. Exist a kind of network, Long Short Term Memory, developed to mitigate the vanishing gradient problem, will be discussed later.

Recurrent Neural Networks (RNN) is a class of NN that keeps connections between nodes and a temporal sequence. The main difference, respect NN, is that it has feedback connections, this memory allows us to keep track of temporal

dynamic behavior, they can process a single data point or the entire sequence of data, like video or speech. They become useful as their intermediate values can store information about past inputs for a time that is not fixed a priori.

This kind of structures is used in a lot of different fields: image classification, captioning, sentiment analysis, machine translation, video classification, etc. . .

Long Short Term Memory (LSTM) normal recurrent neural networks can connect short term event with the present, this feature can be useful in some case, another context could require more long term connection, this special RNN can handle this sort of correlations. Sometimes the recent information is enough to perform the present task, for example, trying to predict the last word of a sentence like "Sun bright in the *sky*" the previous words are enough to correctly predict the last one, no more context is required. In case of more complex sentences some additional information is needed, for the phrase: "I was born in Italy and I speak *italian*" the context becomes fundamental to solve the problem. In this kind of problem LSTM can be really helpful.

LSTM are a special kind of RNN, they were introduced by Hochreiter & Schmidhuber (1997) [6] and later studied and improved. It works really well on a large variety of problems. The special kind of network is explicitly designed for long term connections. The structure of a RNN is normally a single layer, in case of LSTM the structure is more complicated, is made upon four different layers, as shown in figure 4.8.

The first layer filter the input and select what keep and not keep, the second layer will manage what information store in the cell state. The third state is in charge of update the state based on the previous step. The last one decide the output.

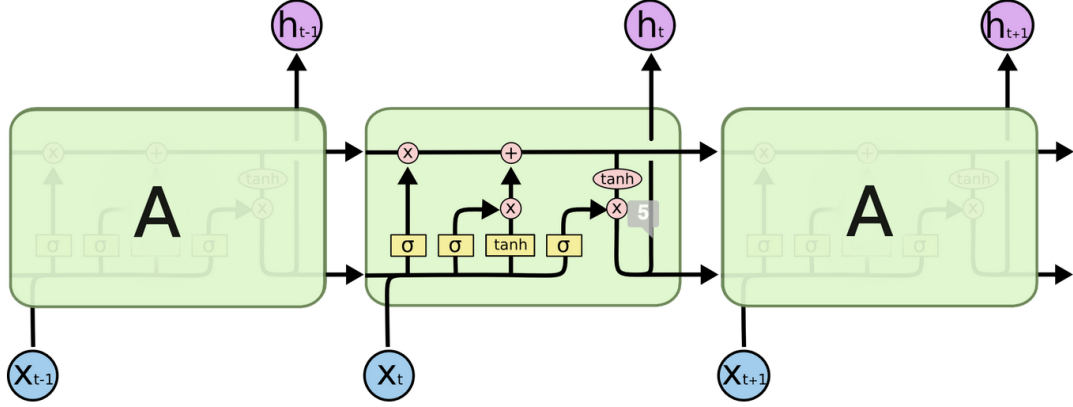


Figura 4.8. Long Short Term Memory layout [7]

4.4 Metriche di valutazione

Every model developed need to be evaluated, there are a lot of different methods that can be used to get metrics about the quality of the system behavior. The computation of these values is achieved via mathematical expression, for this reason, the results, are not always human-readable, some errors like mean absolute error, mean relative error and similar can be read and interpreted without any problem. For metrics like R2 or mean squared error the task becomes more complicated. The following is about the mathematical analysis and description of the different metrics used.

Most of metrics directly derive from the calculation of the simplest error, the *absolute error*, the difference between the target y and the prediction x :

$$\epsilon = |y - x| \quad (4.6)$$

Besides being the easiest to calculate is also the easiest to be understand because show directly the distance respect the desired value.

The first metric based on the previous metric is the *relative error*, that is computed dividing the absolute error by the target value:

$$\eta = \frac{\epsilon}{|x|} = \frac{|y - x|}{|x|} \quad (4.7)$$

This error rescale the output between $[0,1]$ in order to better understand the error in the estimation. For example, in case of a value target value of 530 and a prediction value of 570 the two error are:

$$\epsilon = |520 - 570| = 50 \quad (4.8)$$

$$\eta = \frac{50}{|520|} = 0.09 \quad (4.9)$$

the difference was 50, but respect the expected output the difference is really low, around 9%.

The previous errors are the basic ones, the following is the first based on aggregation of multiple errors calculation. The aggregation is fundamental to generate a single number to evaluate the quality of prediction of our model.

The *Mean Absolute Error* (MAE), is probably the easiest to be calculated and understand, is basically computed by calculating the average of all the absolute errors:

$$MAE = \frac{\sum_{i=1}^n |y_i - x_i|}{n} = \frac{\epsilon}{n} \quad (4.10)$$

Capitolo 5

Pre-elaborazione

Capitolo 6

Predizione

Capitolo 7

Conclusioni

Speaking about conclusion.

Bibliografia

- [1] M. Rath, P. Mäder, “The SEOSS 33 Dataset — Requirements, Bug Reports, Code History, and Trace Links for Entire Projects” in *Data in Brief*, v. 25, p. 104005, 05 2019. [Online]: <https://doi.org/10.7910/DVN/PDDZ4Q>
- [2] M. Denil, de Freitas, in “Narrowing the Gap: Random Forests In theory and In Prattice”
- [3] “Random Forest and its implementation.” [Online]: <https://towardsdatascience.com/random-forest-and-its-implementation-71824ced454f>
- [4] A. Natekin, A. Knoll, “Gradient boosting machines, a tutorial” in *Frontiers in Neurorobotics*, v. 7, p. 21, 2013. [Online]: <https://www.frontiersin.org/article/10.3389/fnbot.2013.00021>
- [5] C. W. John McGonagle, George Shaikouski, *et al.*, “Backpropagation.” [Online]: <https://brilliant.org/wiki/backpropagation/>
- [6] S. Hochreiter, J. Schmidhuber, “Long Short-term Memory” in *Neural computation*, v. 9, pp. 1735–80, 12 1997. [Online]: <https://dl.acm.org/doi/10.1162/neco.1997.9.8.1735>
- [7] C. Olah, in “Understanding LSTM Networks” 8 2015. [Online]: <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>