

# **Intelligenza Artificiale e Machine Learning**

Appunti delle Lezioni di Intelligenza Artificiale e Machine Learning

*Anno Accademico: 2024/25*

*Giacomo Sturm*

*Dipartimento di Ingegneria Civile, Informatica e delle Tecnologie Aeronautiche  
Università degli Studi “Roma Tre”*

Sorgente del file LaTeX disponibile al seguente link:

<https://github.com/00Darxk/Intelligenza-Artificiale-e-Machine-Learning>

## Indice

<b>1</b>	<b>Introduzione: Intelligenza Artificiale e Machine Learning</b>	<b>1</b>
1.1	Storia dell'IA . . . . .	1
1.2	Il Machine Learning . . . . .	3
<b>2</b>	<b>Risoluzione dei Problemi e Ricerca</b>	<b>6</b>
2.1	Algoritmo di Ricerca Generale: Tree Search . . . . .	7
2.1.1	Operazioni su Frontiera, Nodi e Problemi . . . . .	7
2.1.2	Implementazione . . . . .	8
2.2	Criteri di Valutazione . . . . .	9
2.3	Ricerca non Informata o Cieca . . . . .	9
2.3.1	Algoritmo di Ricerca in Ampiezza: Breadth First Search (BFS) . . . . .	9
2.3.2	Algoritmo di Ricerca Guidata dal Costo: Dijkstra . . . . .	10
2.3.3	Algoritmo di Ricerca in Profondità: Depth First Search (DFS) . . . . .	10
2.3.4	Algoritmo di Ricerca in Profondità Limitata . . . . .	11
2.3.5	Algoritmo di Ricerca Iterative-Deepening Search . . . . .	11
2.4	Problema della Ripetizione degli Stati: Graph-Search . . . . .	12
2.5	Algoritmo di Ricerca Informata o Euristica: Best First Search . . . . .	13
2.5.1	Algoritmo di Ricerca Greedy . . . . .	14
2.5.2	Algoritmo A* . . . . .	15
2.6	Algoritmo di Ricerca Locale: Hill-Climbing . . . . .	16
2.6.1	Algoritmo Steepest Ascent Hill-Climbing . . . . .	17
2.6.2	Algoritmo Random-Restart-Hill Climbing . . . . .	17
2.6.3	Algoritmo Stochastic Hill-Climbing . . . . .	18
2.6.4	Algoritmo di Simulated Annealing . . . . .	19
<b>3</b>	<b>Introduzione a Python</b>	<b>21</b>
3.1	Operatori . . . . .	21
3.2	Istruzioni Condizionali . . . . .	21
3.3	Funzioni Built-In, Moduli e Definizione di Funzioni . . . . .	22
3.4	Cicli . . . . .	23
3.5	Stringhe . . . . .	23
3.6	Liste . . . . .	25
3.7	Dizionari e Set . . . . .	27
3.8	Ambienti di Sviluppo . . . . .	29
3.9	Classi e oggetti . . . . .	31
3.10	Liste Concatenate . . . . .	33

## 1 Introduzione: Intelligenza Artificiale e Machine Learning

L'intelligenza artificiale è l'area di studio che analizza come permettere ai computer di compiere operazioni che eseguite da un essere umano richiederebbero intelligenza. Rappresenta uno dei campi di ricerca più interessanti recentemente ed ha avuto un'enorme crescita dal punto di vista economico e tecnologico.

Un modo per poter riepilogare la storia dell'intelligenza artificiale consiste nell'elencare i vincitori del Turing Award, in questo campo. Questo premio rappresenta l'analogo del premio Nobel per l'informatica.

### 1.1 Storia dell'IA

La prima fase di questa storia si può attribuire a personaggi come Alan Turing, che introdusse il concetto di "test di Turing" per determinare se una macchina si può considerare intelligente nel 1950. Una macchina secondo questo criterio si può considerare intelligente se un essere umano interagendo con essa, senza saperlo a priori, non sia in grado di sapere se si tratta di una macchina o di un altro essere umano.

Un altro di questi personaggi fondamentali è John McCarthy, ha coniato il termine "intelligenza artificiale", nel 1955 ed ha proposto la creazione di un gruppo di lavoro l'anno successivo al collage di Dartmouth. Al convegno di Dartmouth del '56 parteciparono tutti i personaggi che vengono chiamati i padri fondatori dell'IA. L'obiettivo del convegno era lo studio della costruzione di macchine in grado di formare astrazioni e concetti, per risolvere problemi, al tempo riservati agli umani. In questo convegno si teorizzò che queste macchine potranno svolgere funzioni umane, ritenute intelligenti.

In questo convegno Cliff Shaw, Allen Newell e Herbert Simon dimostrarono il primo programma di intelligenza artificiale, chiamato "Logic Theorist", in grado di dimostrare i teoremi dei Principia Mathematica. Proposero l'idea di "General Problem Solver" per risolvere una grande varietà di problemi simulando i processi mentali umani. Questo sistema riusciva a dimostrare teoremi, calcolare funzioni matematica e risolvere problemi logici. Da sistemi come questi si crearono grandi aspettative per l'intelligenza artificiale, in questa seconda fase tra il 1952 ed il 1969. In questo periodo Nathaniel Rochester all'IBM sviluppò tra i primi programmi di IA. Nel 1959 Herbert Gelernter scrisse il "Geometry Theorem Prover", capace di dimostrare teoremi matematici complessi. John McCarty nel 1958 all'MIT definì il linguaggio di lato livello Lisp, destinato a diventare il linguaggio di programmazione per eccellenza dell'IA per i seguenti trent'anni. Nel 1963 fondò il laboratorio dell'IA a Stanford, per costruire una versione definita dell'"Advice Taker".

Nonostante queste grandi aspettative, i sistemi di IA realizzati per problemi semplici non dimostravano un comportamento altrettanto soddisfacente per problemi di natura più complessa. Uno dei problemi per questo mancato successo dell'IA in questo periodo tra il '66 ed il '73 fu l'incapacità di comprendere l'intrattabilità di molti problemi, su cui l'IA stava cercando di operare. Nel 1973 venne stilato il rapporto Lighthill, dove venne criticata l'intelligenza artificiale per la mancata soluzione al problema dell'"esplosione combinatoria", ovvero sull'utilizzo dell'IA per problemi di natura reale.

Un'altra critica giunse dal libro "Perceptrons" scritto nel 1969 da Marvin Minsky e Seymour Papert che, come altri, discusse il limite fondamentale della struttura di base per generare un comportamento intelligente. In questo periodo il clima generale attorno all'IA era pessimista e libri come questo aiutarono a maturare questo clima di diffidenza nei confronti delle potenzialità dell'IA.

Negli anni successivi per risolvere il problema dell'esplosione combinatoria si realizzarono sistemi basati sulla conoscenza, riservati ad esperti umani del settore. Questi programmi erano estremamente utili su settori molto ristretti, costituiti da un motore di inferenza. Il collo di bottiglia per questi sistemi era la conoscenza, per cui esistevano intere posizioni dedicate a trovare un modo per rappresentare ed astrarre il dominio di interesse su cui opera l'IA. Queste persone venivano chiamate scienziati della conoscenza, ed erano tipicamente professionisti ed esperti di quel settore specifico.

Nel 1981 il governo giapponese avviò un progetto chiamato "Quinta Generazione" per realizzare computer intelligenti utilizzando il linguaggio Prolog. Analogamente sia gli Stati Uniti che l'Inghilterra finanziarono progetti simili.

Tuttavia questi progetti non riuscirono a mantenere le loro promesse sulle capacità dell'IA o di impatto commerciale. I sistemi erano comunque limitati da una conoscenza limitata ed una difficile capacità di apprendere dall'esperienza e verificare la correttezza. Erano sistemi poco flessibili e robusti rispetto agli obiettivi promessi da questi progetti.

Nonostante questo tra l'inizio e la fine degli anni '80 l'industria dell'IA conobbe un boom economico con centinaia di aziende costruttrici di sistemi esperti, di visione artificiale, di robot, di software ed hardware specializzati per questi scopi.

Molte di queste aziende fallirono per l'impossibilità di mantenere le loro promesse e seguì un periodo chiamato "inverno dell'IA".

A metà degli anni '80 almeno quattro gruppi diversi reinventarono l'algoritmo di apprendimento basato sulla retropropagazione, sviluppato negli anni '60. Questi modelli furono considerati in diretta opposizione ai modelli simbolici di Newell e Simon e dell'approccio logicista di McCarthy. Geoff Hinton è una delle figure in primo piano nel risorgere delle reti neurali durante gli anni '80 e 2010, descrisse i simboli come l'"etere luminifero dell'IA".

Dai limiti dei sistemi esperti si introdusse un nuovo approccio basato sulla probabilità invece che sulla logica booleana, sul "Machine Learning" più che sulla programmazione manuale, su risultati sperimentali più che su affermazioni filosofiche. Questo permette a questi sistemi con logiche non cablate di poter apprendere dall'esperienza e migliorarsi nel tempo.

Nei primi anni del 2000 con lo sviluppo e l'espansione del "World Wide Web" ed il progresso sulla potenza di calcolo ha permesso di realizzare data set molto grandi, fenomeno indicato con il termine "big data". Questo ha portato allo sviluppo di algoritmi di machine learning progettati per trarre vantaggio da questi grandi insiemi di dati. La loro disponibilità ed il cambiamento di approccio verso il ML ha permesso all'IA di recuperare appetibilità commerciale.

Nel 2006 Geoffrey Hinton, introduce un algoritmo di apprendimento veloce per reti neurali, dando il via alla rivoluzione del "Deep Learning". Utilizzando molteplici livelli di elementi computazionali il ML diventa Deep Learning. Questo sistema ha ottenuto successi in molti domini applicativi, aumentando l'interesse verso l'IA.

Attualmente l'IA è in grado di lavorare in molti settori:

- Può guidare veicoli robotizzati come automobili o droni autonomi. Nel 2004 DARPA introduce una sfida per la guida autonoma di veicoli;
- Permette di eseguire locomozione su arti con robot umanoidi;
- Pianificazione e scheduling autonomo, per la navigazione autonoma ed il sistema GPS globale;
- Traduzione automatica;
- Riconoscimento vocale;
- Raccomandazioni su social network per gli utenti. Nel 2006 Netflix ha lanciato una borsa per aumentare la precisione del loro sistema di raccomandazione tramite ML, vinta nel 2009;
- Può battere i migliori giocatori umani in molti giochi. Nel 1997 Deep Blue batte il campione del mondo di scacchi Garry Kasparov, e nel 2011 IBM Watson batte i campioni del gioco "Jeopardy!", ora questo sistema è utilizzato in molti sistemi. Tra il 2015 ed il 2017 AlphaGo di DeepMind batte i campioni del mondo nel gioco di Go, lo stesso team ha sviluppato AlphaZero in grado di giocare sia scacchi, shogi e Go;
- Può interpretare e riconoscere immagini. Nel 2021 ImageNet comincia il suo concorso annuale per rilevare e classificare correttamente oggetti in un set di immagini ampio e ben accurato, con un incremento nella precisione dovuto ai progressi nelle "deep convolutional neural networks". Nel 2014 Facebook pubblica un lavoro sul DeepFace in grado di identificare volti con un'accuratezza del 97%;
- Permette di diagnosticare malattie, raggiungendo o superando la diagnosi di medici esperti;
- Recentemente ha accelerato la ricerca per il vaccino anti-Covid, nell'individuazione delle proteine utilizzate;
- Permette di rilevare informazioni dettagliate su eventi climatici estremi.

Nel 2014 il consumo di energia per raffreddare i Data Center è stato ridotto del 40% con un modello di ML, analogamente nel 2017 il software per analizzare le immagini di galassie sotto lenti gravitazionali è stato velocizzato di un fattore di  $10^7$ .

Nel 2022 venne rilasciato ChatGPT, "Generative Pre-trained Transformer", il chatbot di OpenAI progettato per rendere efficace la comunicazione con un utente umano.

## **1.2 Il Machine Learning**

L'intelligenza artificiale è una disciplina molto vasta, non copre solo tematiche sulla ricerca di una soluzione ottima. Esistono sistemi esperti per effettuare decisioni, offrire raccomandazioni, aiutare nella programmazione.

Dopo il 2006 le aziende hanno capito l'impatto che le raccomandazioni hanno nel mantenere l'attenzione del cliente e migliorare il loro servizio, quando Netflix ha annunciato un concorso di un milione di dollari per migliorare il suo programma di raccomandazioni del 10%. Da allora tutte le

principali aziende di investimenti, di viaggi, di intrattenimento, etc. Tutte le aziende che hanno di interesse l'interazione degli utenti.

Ora è veramente di interesse, infatti la massima conferenza sulla ricerca SIGIR "Special Interest Group on Information Retrieval", prima trattava principalmente problemi di ricerca, mentre recentemente la sua attenzione si è rivolta su sistemi di raccomandazione.

Un sistema di Machine Learning ML, all'interno di un dominio, dove ci sono dati validi con una loro consistenza, rappresenta un processo di apprendimento supervisionato o meno su questi dati per poter gestire e generalizzare i nuovi dati nello stesso dominio applicativo.

Uno dei grandi passi avanti nel campo delle immagini si deve al grande dataset ImageNet, grazie ad una ricercatrice e professoressa di Stanford, che ha investito il suo tempo nel progetto. Questo insieme contiene milioni di immagini descritte e la loro classe di appartenenza dell'oggetto raffigurato, sono infatti principalmente immagini in primo piano. Tramite insiemi di dati come questi, il sistema apprende in un approccio supervisionato, poiché ad ogni immagine viene indicata la classe di appartenenza, ed in alcuni casi si indica dov'è presente l'oggetto nell'immagine. Date queste due indicazioni si possono creare sistemi di identificazione o individuazione di un'immagine. Questo approccio può realizzare sistemi di riconoscimento facciale, individuando le facce da un'immagine ed in seguito riconoscere le facce dall'immagine.

La maggior parte dei dataset saranno supervisionati, ma esistono moltissimi dataset i cui data-point non sono indicati o descritti. Questo tipo di apprendimento prevede l'assenza delle classi dell'appartenenza, questo viene utilizzato per creare sistemi che possono essere in grado di identificare il rischio di una certa malattia, tramite delle caratteristiche comportamentali o genetiche. Oppure si assegnano a degli algoritmi di clustering un insieme di utente, che sulla base delle caratteristiche di quest'ultimi individua delle classi di utenti. Nel diagramma di dispersione dove sono presenti questi data-point si possono quindi individuare cluster diverse. Si utilizza questo approccio nelle comunità social per identificare comunità distinte.

Se si trattano fenomeni fisici che non sono caratterizzabili a priori si effettuano degli studi per vedere se sono presenti correlazioni tra i data-point e quindi dei cluster, ed eventualmente individuare informazioni da questi dati.

L'apprendimento per rinforzo è molto comune, ed è un sistema dove un agente interagisce con un ambiente, senza un dataset, per permettere all'agente di operare sull'ambiente, questo fornisce delle ricompense all'agente in caso effettui un certo tipo di azioni. Queste ricompense possono essere positive o negative, e sulla base di queste, si realizza il meccanismo di apprendimento. Alcune aree di dominio che si predispongono a questo tipo di apprendimento, dove sulla base delle ricompense si possono suscitare reazioni e quindi comportamenti desiderati dall'agente, sono i videogiochi.

Grazie al DL si ha una sistema di "Deep Reinforcement Learning", dove il sistema è in grado di imparare direttamente da un'immagine, un video, uno spettrogramma, etc.

Nel 2013 la DeepMind, di proprietà di Google ha dimostrato la possibilità di un agente di apprendere tramite DRL, su numerose partite di Breakout un comportamento "super-umano", e diversi altri giochi arcade della Atari. Un altro grande dominio dove questo tipo di approccio è utile è la robotica, dimostrato dalla Boston Dynamics e dal loro robot umanoide Atlas, in grado di camminare, saltare, portare pesi, su diversi tipi di terreni autonomamente. Il loro successivo robot Spot, ha risolto alcuni dei problemi, non dovendo imparare una locomozione su due piedi molto più complessa, basandosi invece sullo scheletro di un animale quadrupede come un cane.

Il reinforcement learning venne introdotto per la prima volta una ventina di anni fa da Richard Sutton e Andrew Barto, degli psicologi cognitivi, senza tutta la serie di risorse di calcolo e di data per favorire questo approccio. Quindi non si è affermato per molto tempo.

Nel 1997 venne proposta l'idea della RoboCup, un campionato di calcio composto da squadre di robot autonomi, da parte di diverse università entro il 2050, in grado di confrontarsi con la squadra di calcio campione al mondo.

Il ML quindi ha l'obiettivo di generalizzare ed effettuare decisioni, il quanto più accurate, su un relativo dominio i cui dati non ha ancora visto.

L'apprendimento è una componente chiave del ragionamento, ma non è sufficiente a generare IA in grado di generare nuova conoscenza. Tecniche di generazione del testo, delle immagini e musica, etc., sono note già da anni, ma solo recentemente con ChatGPT utilizzando una serie di tecniche che consentono un addestramento molto accurato ed efficiente su di un patrimonio di dati che è sterminato. Tutte le aziende informatiche, capiscono che questo è una minaccia, poiché permetterebbe di catturare oltre alla frequenza dei termini anche il loro contesto. Dal contesto è quindi in grado di determinare la probabilità che una certa parola appaia. Infatti ChatGPT è uno strumento statistico, poiché non ha alcun tipo di ragionamento, per questo si parla anche di allucinazioni, poiché non può creare nuova conoscenza, ma offre esempi dove lo strumento non rinuncia di predire, ma produce un risultato scorretto.

Il processo di ML permetterebbe di realizzare software 2.0, tramite una rete neurale in grado di identificare le istruzioni da dover eseguire per effettuare una data azione, senza l'appoggio di un programmatore esterno. Questo solleva una serie di problematiche relative al pericolo del ML.

Jeremy Howard ha esposto su una "TED talk" un problema analogo alla rivoluzione industriale in termini di occupazione, dove l'80% delle persone lavora sui servizi, un dominio di interesse dove l'intelligenza artificiale è in grado di lavorare molto bene. Per cui quando questi sistemi entreranno a regime e le aziende si renderanno conto di questi strumenti così potenti a disposizione, allora sostituiranno questa occupazione umana, creando problemi occupazionali. Già nel 2015 nel concorso ImageNet, sullo stesso dataset, il tasso di errore di agenti intelligenti è sceso al disotto del tasso di errore umano nel riconoscimento delle immagini. Anche se l'essere umano è in grado di identificare in più classi distinte rispetto all'agente artificiale ed è in grado di contestualizzare queste immagini. Ma in alcune professioni quello che le macchine possono offrire, in termini economici, è più conveniente rispetto a quello che lavoratori umani possono offrire. Nel 2023 c'è stato lo sciopero più lungo nell'ambito di Hollywood, per l'introduzione dell'IA nel cinema, poiché questo avrebbe messo a rischio il loro lavoro e la loro occupazione. In seguito a questo sciopero le case cinematografiche hanno imposto dei limiti per evitare questo caso peggiore.

Il ML rappresenta solo una piccola parte dell'intelligenza artificiale, ma comprende parti da altre discipline come il "neurocomputing", la statistica come classificatori bayesiani, il riconoscimento di andamenti.

## 2 Risoluzione dei Problemi e Ricerca

Si definisce un agente risolutore di problemi un agente con uno specifico obiettivo da raggiungere e che deve identificare una sequenza di azioni per raggiungerlo.

Bisogna determinare l'obiettivo, un insieme degli stati del mondo dove ci si trova, che si vuole raggiungere. Inoltre bisogna formulare il problema, ovvero le azioni e gli stati considerati dall'agente.

Un agente che ha a disposizione diverse opzioni immediate di valore sconosciuto, può decidere quale scegliere la sua azione analizzando le diverse possibili sequenze di azioni, che portano a stati di valore conosciuto, per scegliere la sequenza di costo migliore.

Questo processo di selezione viene definito ricerca, un algoritmo di ricerca quindi prende un problema e restituisce una soluzione costituita da una sequenza di azioni.

Nella formulazione si definisce uno stato iniziale e dell'obiettivo, come un insieme di stati e si definiscono le azioni come transizioni tra stati. Dopo aver trovato la sequenza di azioni corrispondente alla soluzione, la esegue.

Si possono distinguere due tipi di problemi, i problemi giocattolo o "toy problems", sono ideati come illustrazione o esercitazione dei metodi risolutivi. Rappresentano delle astrazioni, anche semplificate, dei problemi del mondo reale in generale più difficili a cui si è effettivamente interessati.

I problemi del mondo reale possono essere la configurazione VLSI, la navigazione dei robot, la sequenza di montaggio, la ricerca dell'itinerario e generali problemi di viaggio come il commesso viaggiatore.

Si possono utilizzare due tipi diversi di formulazione di un problema. Si può partire da uno stato vuoto, ed utilizzando operatori si può estendere progressivamente la descrizione dello stato. Invece nella formulazione a stato completo consiste nel partire da uno stato iniziale completo, ed ogni operatore altera questo stato per cercare una soluzione.

Lo spazio degli stati è un grafo che rappresenta tutti i possibili stati, come nodi, collegati tra archi che rappresentano le possibili azioni. Il problema consiste nel trovare un percorso in questo spazio degli stati, dallo stato iniziale ad uno dei possibili stati soluzione. L'algoritmo deve decidere ad ogni stato quale azione prendere e quindi a quale nodo del grafo spostarsi. Per definire il costo della soluzione, si considera il costo del cammino delle azioni intraprese dall'agente.

Per definire formalmente un problema, sono necessarie quattro componenti. Uno stato iniziale in cui si trova l'agente. Una descrizione delle azioni possibili, questa può utilizzare una funzione successore che dato uno stato restituisce i suoi possibili successori. Tramite la funzione successore e lo stato iniziale si può costruire lo spazio degli stati. Oppure si può utilizzare un insieme di operatori. Un test obiettivo per determinare se un particolare stato è uno stato obiettivo. Come ultimo componente è necessaria una funzione di costo del cammino per determinare il costo di una data soluzione, assegnando un valore numerico ad ogni cammino.

Il tipo di dato problema è rappresentato da questi quattro componenti, le istanze di questo tipo di dato rappresentano gli input degli algoritmi di ricerca.

Il mondo reale è estremamente complesso, e lo spazio degli stati deve essere creato mediante un processo di astrazione. In questo processo, lo stato astratto rappresenta un insieme di stati reali più complessi. Analogamente per le azioni astratte, queste rappresentano combinazioni di azioni reali. Nei problemi giocattolo non è necessario effettuare questo processo di astrazione, poiché rappresenta un problema semplice.



Per individuare queste possibili sequenze di azioni l'algoritmo si può costruire un albero di ricerca, dove il nodo radice corrisponde allo stato iniziale, ed i rami rappresentano le azioni possibili, ed i nodi figli rappresentano gli stati successori di un certo stato. Tuttavia i nodi dell'albero di ricerca e gli stati dello spazio degli stati sono differenti, poiché è possibile che più nodi condividano lo stesso stato, mentre ogni stato nello spazio è univoco. Generalmente si vuole evitare di ripetere stati all'interno di un cammino, questo rappresenta il problema degli stati ripetuti, e sarà analizzato successivamente.

Ad ogni nodo si possono inserire altre informazioni utili, oltre allo stato, un riferimento al genitore, l'operatore che ha generato lo stato, la profondità, il costo del cammino parziale fino a questo stato, etc.

`NODO = <stato, genitore, operatore, profondità, costo parziale, ...>`

Il processo di ricerca comporta la stessa sequenza di azioni. Deve scegliere tra le foglie dell'albero corrente un nodo da "espandere", secondo un certo criterio o strategia. In seguito bisogna determinare se questo nodo rappresenta un'obiettivo del problema, altrimenti vengono generati i suoi nodi figli, ed i corrispondenti stati successori, e tutte le componenti dei nodi figli. La collezione dei nodi in attesa di essere espansi viene chiamata in vari modi: confine, frontiera, frangia o lista aperta.

## 2.1 Algoritmo di Ricerca Generale: Tree Search

Un algoritmo generale di ricerca può essere chiamato **TREE-SEARCH**. Quest'algoritmo prende un problema ed una strategia come input e restituisce una soluzione oppure un fallimento. Viene realizzato semplicemente ad un ciclo che ripete le operazioni precedentemente descritte, fino a quando non identifica una soluzione o viene sollevato un problema, e quindi restituisce un fallimento. Il primo passo è la generazione dell'albero di ricerca del problema, se la frontiera è vuota, ovvero non esistono nodi candidati per l'espansione viene riportato un fallimento, poiché non è stato ancora trovato uno stato obiettivo. Se si arriva ad un nodo corrispondente ad uno stato obiettivo viene restituita rappresenta la sequenza di nodi ottenuta come soluzione.

La frontiera può contenere un nodo relativo ad uno stato obiettivo, ma l'algoritmo non termina fino a quando non viene scelto per essere espanso.

Una strategia di ricerca rappresenta un criterio per decidere quale nodo da espandere. Può essere definita come una funzione per la scelta di un elemento tra un insieme di nodi, la frontiera. Oppure si può considerare come una funzione di inserimento di un elemento in una sequenza. Se già nella fase di inserimento si analizza tramite una metrica il valore di ognuno di questi stati, allora l'elemento in prima posizione in questa struttura dati rappresenta il nodo migliore.

### 2.1.1 Operazioni su Frontiera, Nodi e Problemi

La frontiera viene implementata con una struttura dati chiamata coda, ma non necessariamente segue la disciplina FIFO. Su questa coda sono definite una serie di operazioni:

- **MAKE-QUEUE**: prende come input un nodo **n** e restituisce una coda **q**, realizza una coda contenente solo il nodo **n**;

- **EMPTY**: prende come input una coda *q* e restituisce un booleano, verifica se la coda è vuota;
- **REMOVE-FRONT**: prende come input una coda *q* e restituisce il primo nodo *n* della lista;
- **QUEUING-FN**: prende come input una coda *q* ed una lista di nodi *n*, e restituisce la coda con aggiunti tutti questi nodi.

L'ultima funzione si utilizza quando si producono tutti i nodi successori e si vogliono aggiungere alla lista. Questa funzione non inserisce generalmente in coda, ma dipende dalla strategia di ricerca utilizzata.

Assumendo che esistano le operazioni sul tipo di dato problema, e sul tipo di dato nodo. Le operazioni sul tipo di dato problema sono:

- **INITIAL-STATE**: prende come input un problema *p* e restituisce lo stato iniziale del problema *n*;
- **GOAL-TEST**: prende come input un problema *p* ed uno stato *n* e verifica se questo rappresenta una soluzione, restituendo un booleano;
- **OPERATORS**: prende come input un problema *p* e restituisce una lista con tutti gli operatori del problema *ops*. Ogni operatore *op* applicato ad uno stato *n* restituisce una lista di stati *ss*.

Le operazioni sul tipo di dato nodo sono:

- **MAKE-NODE**: prende come input uno stato *s* e costruisce un nodo su di esso *n*;
- **STATE**: prende come input un nodo *n* e ne restituisce lo stato contenuto *s*;
- **EXPAND**: prende come input un nodo *n* ed una lista di operatori *ops* e restituisce una lista di nodi successori *ns*.

### 2.1.2 Implementazione

Date queste operazioni, si può rappresentare in pseudocodice l'algoritmo di **TREE-SEARCH** in modo più semplice:

```
function TREE-SEARCH(problem) returns a solution or failure

fringe <- MAKE-QUEUE(MAKE-NODE(INITIAL-STATE(problem)))
loop do
  if EMPTY(fringe) then return failure
  node <- REMOVE-FRONT(fringe)
  if GOAL-TEST(problem, STATE(node)) then return SOLUTION(node)
  fringe <- QUEUING-FN(fringe, EXPAND(node, OPERATOR(problem)))
end
```

In questa variazione non viene conservato l'intero albero, ma solamente la coda con i nodi della frontiera.

## 2.2 Criteri di Valutazione

Per valutare questi algoritmi oltre alla complessità temporale e spaziale, si utilizzano altre due criteri, la completezza e l'ottimalità. Un'algoritmo si definisce completo, se quando esiste una soluzione è garantito sia in grado di trovarla. Un algoritmo si dice ottimo se dato un problema con diverse soluzioni, individua sempre la migliore, quella a costo minimo. La complessità dell'algoritmo dipende dal fattore di ramificazione  $b$  dello spazio degli stati e dalla profondità  $d$  della soluzione più superficiale. Il fattore di ramificazione  $b$  rappresenta il massimo numero di figli che un nodo può avere. Mentre la profondità  $d$  è la minima lunghezza di un cammino dal nodo iniziale alla radice.

## 2.3 Ricerca non Informata o Cieca

Questi algoritmi non sono molto efficienti in generale, ma sono utili per comprendere il comportamento gli algoritmi di ricerca informata o di euristica, che si avvalgono della conoscenza sul dominio dello spazio degli stati e dalla creazione dell'albero di ricerca per scegliere il percorso più promettente. Nel caso medio quest'ultimi sono certamente più efficienti degli algoritmi trattati in questa sezione.

### 2.3.1 Algoritmo di Ricerca in Ampiezza: Breadth First Search (BFS)

Nella ricerca in ampiezza si espande il nodo radice, e si espandono i nodi generati dalla radice, e si ripete per ogni nodo successore. Per implementare un algoritmo che utilizza una strategia di ricerca non informata in ampiezza, la "Queuing Function" inserisce i nodi appena generati in coda. Questa funzione quindi rappresenta sempre un inserimento in coda e si può chiamare "Enqueue at the End": **ENQUEUE-AT-END**.

Un algoritmo che utilizza questo tipo di strategia viene chiamato "Breadth First Search" o in ampiezza, e si può implementare in modo analogo all'algoritmo di ricerca generale trattato precedentemente:

```
function BREADTH-FIRST-SEARCH(problem) returns a solution or failure

fringe <- MAKE-QUEUE(MAKE-NODE(INITIAL-STATE(problem)))
loop do
  if EMPTY(fringe) then return failure
  node <- REMOVE-FRONT(fringe)
  if GOAL-TEST(problem, STATE(node)) then return SOLUTION(node)
  fringe <- ENQUEUE-AT-END(fringe, EXPAND(node, OPERATOR(problem)))
end
```

In questo approccio, tutti i nodi di profondità  $d$  vengono espansi prima dei nodi di profondità  $d+1$ . Rappresenta una strategia sistematica, ma permette di individuare solamente i nodi obiettivi più superficiali, non è garantito che questo rappresenta la soluzione ottima. Questo algoritmo è quindi completo, ma non è ottimo. Invece è ottimale se il costo del cammino  $g(n)$  è una funzione monotona non decrescente della profondità del nodo  $p(n)$ :

$$p(n) < p(m) \implies g(n) \leq g(m)$$

$$p(n) = p(m) \implies g(n) = g(m)$$

Ovvero se due nodi  $n$  ed  $m$  sono a profondità diverse, dove il nodo  $m$  è a profondità maggiore, il costo del cammino dalla radice al nodo  $n$  è al massimo uguale al costo del cammino dalla radice al nodo  $m$ . Inoltre se i nodi sono alla stessa profondità, allora i costi dei loro cammini dalla radice sono uguali.

Utilizzando questo algoritmo, bisogna generare un numero di nodi, prima di trovare una soluzione, almeno pari a tutti i nodi precedenti al nodo soluzione. Nel caso peggiore questo nodo è l'ultimo nodo espanso alla profondità  $d$ , e per ogni nodo vengono generati esattamente  $b$  figli, quindi bisogna espandere al massimo un numero di nodi  $N$  pari a:

$$N = \left( \sum_{i=1}^{d+1} b^i \right) - b$$

Supponendo che ogni generazione rappresenta un'operazione semplice allora la complessità temporale di questa ricerca è  $O(N) = O(b^d)$ . La complessità temporale è analogamente  $O(b^d)$  poiché bisogna memorizzare tutte le foglie generate.

### 2.3.2 Algoritmo di Ricerca Guidata dal Costo: Dijkstra

Modificando la ricerca in ampiezza espandendo il nodo della frontiera di costo più basso, si può aumentare l'efficienza dell'algoritmo precedente.

Si definisce con  $g(n)$  il costo del cammino dalla radice al nodo  $n$ . Questo valore viene salvato nella struttura dati nodo, e viene scelto il nodo di costo  $g(n)$  minore per essere espanso. Se il costo del cammino corrisponde alla funzione di profondità, si ha la ricerca in ampiezza. Questo algoritmo è completo e ottimale quando il costo di ogni step è sempre maggiore o uguale ad una costante positiva  $\varepsilon$ . Per cui è garantito che non attraversi costantemente lo stesso cammino, senza espandere altri nodi di profondità minore.

I costi di ogni nodo vengono salvati in un campo etichetta nella struttura dati nodo. Dalla frontiera si estrae sempre il nodo a costo minore, questa collezione viene quindi ordinata in base al costo delle etichette di ogni nodo.

### 2.3.3 Algoritmo di Ricerca in Profondità: Depth First Search (DFS)

La ricerca in profondità consiste nell'espansione del nodo più profondo, dopo aver espanso la radice. Per implementare questa funzione, si utilizza una queuing function che inserisce i nodi appena espansi all'inizio della lista, con una "Enqueue at the Front": **ENQUEUE-AT-FRONT**:

```
function DEPTH-FIRST-SEARCH(problem) returns a solution or failure

fringe <- MAKE-QUEUE(MAKE-NODE(INITIAL-STATE(problem)))
loop do
  if EMPTY(fringe) then return failure
  node <- REMOVE-FRONT(fringe)
```

```

    if GOAL-TEST(problem, STATE(node)) then return SOLUTION(node)
    fringe <- ENQUEUE-AT-FRONT(fringe, EXPAND(node, OPERATOR(problem)))
end

```

Questa funzione si può implementare mediante una funzione ricorsiva, dove viene passata una versione del problema, dove i nodi appena generati rappresentano i nuovi nodi radice. Quindi per ogni espansione vengono generati al massimo  $b$  sotto-problemi, risolti dallo stesso algoritmo. La lista dei nodi da visitare viene conservata implicitamente nello stack dei record di attivazione delle varie chiamate ricorsive. Quando si raggiunge un nodo foglia non obiettivo, si effettua il backtracking, ovvero si risale l'albero fino a trovare il nodo a profondità maggiore non ancora espanso su cui è possibile effettuare una scelta.

Questa ricerca non è né completa né ottimale, ha una complessità temporale di  $O(b^m)$ , dove  $m$  rappresenta la profondità massima dell'albero di ricerca. Se una soluzione è presente a profondità minore nel sotto-albero di destra, non verrà mai individuata se non è stato già espanso tutto il sotto-albero di sinistra, senza aver trovato una soluzione. Quindi se individua una soluzione la restituisce indipendentemente dalla sua ottimalità.

Si guadagna rispetto alla ricerca in ampiezza nella complessità spaziale. Infatti non bisogna memorizzare l'intero albero, ma solamente il cammino dalla radice alla foglia, ed i fratelli non espansi di ciascun nodo del cammino di profondità  $m$ :  $O(b \cdot m)$ . Se l'albero ha rami infiniti, allora la ricerca non termina.

### 2.3.4 Algoritmo di Ricerca in Profondità Limitata

Nella ricerca in profondità limitata si impone un limite alla profondità massima, per impedire di proseguire all'infinito su uno stesso cammino. Un nodo viene espanso solo se la lunghezza del cammino dalla radice al nodo è minore del massimo stabilito. Se non viene trovata alcuna soluzione restituisce il valore speciale taglio se alcuni nodi non sono stati espansi, altrimenti fallisce.

Si possono utilizzare conoscenze specifiche al problema per fissare questo limite. Se si lavora su di un grafo si potrebbe utilizzare il diametro del grafo come la profondità. Se non si sceglie un valore adeguato per questo limite allora l'algoritmo non funzionerà correttamente. L'algoritmo è completo, se la soluzione è ad una profondità minore della lunghezza  $l$  imposta, mentre non è ottimale. La complessità temporale e spaziale è rispettivamente  $O(b^l)$  e  $O(b \cdot l)$ . Risolve il problema della completezza, ma non risolve l'ottimale.

### 2.3.5 Algoritmo di Ricerca Iterative-Deepening Search

Questo algoritmo risolve il problema dell'ottimalità sugli algoritmi di ricerca in profondità senza conoscere un limite adeguato. Evita il problema della scelta del limite provando iterativamente tutti i limiti possibili fino a quando non individua una soluzione:

```

function ITERATIVE-DEEPENING-SEARCH(problem) returns solution or failure
    for depth = 0 to  $\infty$  do
        if DEPTH-LIMITED-SEARCH(problem, depth) succeeds

```

```

    then return ist result
end

```

Questo approccio combina i benefici di una ricerca in ampiezza con i benefici di una ricerca in profondità, poiché per ogni profondità vengono analizzati tutti i nodi.

Quindi questo algoritmo è ottimale e completo per le condizioni della ricerca in ampiezza. La complessità spaziale è  $O(b \cdot d)$ , quindi non è esponenziale. Mentre la complessità temporale è simile a quella a quella della ricerca in ampiezza. L'algoritmo viene richiamato ogni volta che si aumenta il limite, quindi i nodi a profondità minore vengono generati ogni volta che viene eseguito nuovamente l'algoritmo. Quindi vengono generati in totale  $N$  nodi:

$$N = \sum_{i=1}^d b^i \cdot (d + 1 - i)$$

I primi  $b$  nodi a profondità 1 vengono generati  $d$  volte, fino ai nodi al livello  $d$  generati una sola volta. Questo algoritmo è quindi circa l'11% meno efficiente rispetto alla ricerca in ampiezza. Ma non rappresenta un incremento considerevole rispetto alla ricerca in ampiezza, quindi è accettabile.

## 2.4 Problema della Ripetizione degli Stati: Graph-Search

Il problema della ripetizioni degli stati può provocare gravi complicazioni nel processo di ricerca. Questo problema sorge soprattutto quando sono possibili azioni bidirezionali ed in questo caso è possibile che gli alberi di ricerca siano infiniti. Si vuole quindi evitare quando è possibile di ripetere gli stessi stati in più nodi dell'albero di ricerca.

Questi stati ripetuti possono in certi casi rendere il problema irrisolvibile, è conveniente controllare se uno stato è replicato. Se un algoritmo arriva ad uno stesso stato attraverso due cammini differenti, allora ha individuato uno stato ripetuto e deve scartare uno di questi due cammini, per determinare quale scartare si sceglie generalmente l'ultimo cammino ottenuto. Si scarta anche se questo cammino è migliore del cammino precedente. Si utilizza questo approccio poiché negli algoritmi di ricerca euristica, sotto alcune condizioni, quando trova un percorso questo è ottimo, quindi non sorgono problemi nello scartare cammini che portano allo stesso stato.

In altre implementazioni dove non è garantito che il primo percorso trovato sia il migliore, bisogna controllare quale dei due cammini presenti il costo migliore. Per evitare la ripetizione bisogna contenere gli stati già visitati in memoria, tramite un'altra struttura dati chiamata insieme esplorato o lista chiusa, contenente ogni nodo espanso.

Si modifica l'algoritmo Tree Search nell'aggiunta alla frontiera per verificare la ripetizione degli stati:

```

function GRAPH-SEARCH(problem) returns a solution or failure

  close <- empty set
  fringe <- MAKE-QUEUE(MAKE-NODE(INITIAL-STATE(problem)))
  loop do
    if EMPTY(fringe) then return failure

```

```

node <- REMOVE-FRONT(fringe)
if GOAL-TEST(problem, STATE(node)) then return SOLUTION(node)
if STATE(node) not in close
  then ADD(close, node)
  child_list <- EXPAND(node, OPERATOR(problem))
  for child_node in child_list
    if STATE(child_node) not in close then
      fringe <- QUEUING-FN(fringe, child_node)
end

```

Questo approccio si chiama Graph Search, dove prima di aggiungere un nodo alla frontiera, si controlla se il suo stato è già stato avvistato. Si suppone che il primo cammino che raggiunge uno stato  $s$  è il più conveniente. Questo algoritmo realizza un albero direttamente sul grafo dello spazio degli stati, poiché è presente al massimo una singola copia di ogni stato. La frontiera separa nel grafo dello spazio degli stati in due regioni, una esplorata, ed una da esplorare. In questo modo ogni cammino dallo stato iniziale ad uno stato inesplorato deve passare attraverso uno stato sulla frontiera.

L'algoritmo scarta sempre il cammino appena trovato, se lo stato raggiunto è ripetuto, quindi potrebbe scartare un cammino corrispondente ad una soluzione migliore. Potrebbe quindi non essere un algoritmo ottimale.

Inoltre l'uso della lista chiusa significa che la ricerca in profondità e quella ad approfondimento iterativo non richiedono requisiti spaziali lineari.

## 2.5 Algoritmo di Ricerca Informata o Euristica: Best First Search

Quando si parla di ricerca euristica, l'algoritmo può sfruttare conoscenze specifiche sul problema in questione, aiutandolo nella scoperta della soluzione. Questa modifica si inserisce nella funzione di inserimento in coda **QUEUING-FN**. Questa conoscenza sul dominio del problema viene implementata tramite una funzione di valutazione  $f$  applicata ai nodi dell'albero di ricerca. Questa funzione stima quanto un nodo sia più o meno "promettente", ovvero stima della desiderabilità di espandere il nodo associato. Generalmente la funzione  $f$  è una funzione di stima del costo della soluzione, per cui si considera il nodo  $n$  appartenente alla frontiera con  $f(n)$  minore.

L'algoritmo "Best First Search" ordina i nodi inseriti nella frontiera dal migliore al peggiore secondo una data funzione di valutazione  $f$ . Il nodo scelto da espandere è quello con valutazione migliore, in genere con  $f(n)$  minore, per cui a differenza di funzioni di valutazione  $f$  si hanno diverse versioni di questo algoritmo di ricerca.

La ricerca guidata dal costo (2.3.2) si può considerare un caso particolare della ricerca best first dove la funzione di valutazione  $f(n)$  coincide con la funzione di costo del cammino  $g(n)$  dallo stato iniziale al nodo  $n$ . Questo rappresenta un'informazione di euristica nulla.

```

function BEST-FIRST-SEARCH(problem, EVAL-FN) returns a solution or failure

QUEUING-FN <- function that orders nodes by EVAL-FN
fringe <- MAKE-QUEUE(MAKE-NODE(INITIAL-STATE(problem)))

```

```

loop do
  if EMPTY(fringe) then return failure
  node <- REMOVE-FRONT(fringe)
  if GOAL-TEST(problem, STATE(node)) then return SOLUTION(node)
  fringe <- QUEUING-FN(fringe, EXPAND(node, OPERATOR(problem)))
end

```

Dove EVAL-FN rappresenta la funzione di valutazione.

Per favorire la comprensione si utilizza una serie di notazioni:

- $s, s_0, s_1, \dots, s_i$ : stati del problema;
- $s_0$ : stato iniziale;
- $k^*(s_i, s_j)$ : costo del cammino minimo da  $s_i$  a  $s_j$ , se esiste;
- $g^*(s_i) = k^*(s_0, s_i)$ : costo del cammino minimo dallo stato iniziale a  $s_i$ ;
- $h^*(s_i)$ : costo effettivo, non una stima, di un cammino da uno stato  $s_i$  ad uno stato obiettivo;
- $f^*(s_i) = g^*(s_i) + h^*(s_i)$ : il costo minimo di una soluzione vincolata a passare per  $s_i$ .

Si introducono ulteriori notazioni per l'albero di ricerca:

- $n, n_1, n_2, \dots, n_i$ : nodi dell'albero di ricerca;
- $g^*(n)$ : costo del cammino dallo stato iniziale allo stato associato al nodo  $n$ ;
- $h^*(n)$ : costo effettivo di un cammino dallo stato associato al nodo  $n$  ad uno stato obiettivo;
- $g(n)$ : costo del cammino dallo stato iniziale allo stato di  $n$ ,  $g(n) \leq g^*(n)$ ;
- $h(n)$ : stima di  $h^*(n)$ .

Le funzioni senza asterisco all'apice si riferiscono a funzioni di stima, altrimenti sono funzioni di costi effettivi. La funzione  $h$  si dice funzione euristica.

### 2.5.1 Algoritmo di Ricerca Greedy

Questo algoritmo "goloso" utilizza la funzione  $h$  come funzione di valutazione, si espande quindi il primo nodo che si ritiene sia vicino all'obiettivo. Se  $n$  corrisponde allo stato obiettivo, deve essere  $h(n) = 0$ . Minimizza il costo stimato per raggiungere l'obiettivo.

Una possibile funzione euristica è la distanza a linea d'aria "Straight Line Distance" SLD:  $h_{SLD}(n)$ , per problemi di viaggio. Questo algoritmo non espande nodi inutilmente, ma la soluzione trovata dall'algoritmo potrebbe non essere la soluzione ottima del problema. In generale la ricerca golosa non è ottimale. Inoltre non è neanche completa, poiché se una soluzione richiedesse allontanarsi dall'obiettivo e quindi andare verso stati con una valutazione peggiore, non sarebbe mai scoperta da questo algoritmo.



Nel caso peggiore è anche esponenziale spazialmente e temporalmente con una complessità asintotica di  $O(b^m)$ . Nonostante questi risultati con una buona euristica è possibile ottenere buoni risultati anche con un algoritmo greedy. Uno dei difetti è che la scelta del nuovo stato è effettuata interamente dalla stima della distanza, senza considerare il percorso parziale.

### 2.5.2 Algoritmo A\*

L'algoritmo A\* risolve il problema dell'algoritmo greedy, considerando anche il percorso parziale nella sua funzione di valutazione:  $f(n) = g(n) + h(n)$ , fornisce quindi una stima del costo del cammino dallo stato iniziale ad uno stato obiettivo, vincolato a passare per il nodo  $n$ .

Con questa semplice aggiunta è possibile migliorare notevolmente l'algoritmo precedente, inoltre sotto certe condizioni è possibile recuperare l'ottimalità e la completezza. Date queste due ipotesi sul grafo di partenza:

- Ogni nodo del grafo abbia un numero finito di successori;
- Tutti i costi abbiano costi maggiori di una quantità positiva  $\delta$ .

Se  $h$  è un'euristica ammissibile, ovvero se per ogni nodo  $n$  del grafo vale la condizione  $h(n) \leq h^*(n)$ , allora l'algoritmo A\* è ottimale e completo:

$$\forall n \text{ t.c. } h(n) \leq h^*(n) \implies \text{A*}: \text{completo ed ottimale} \quad (2.5.1)$$

In generale per un dato problema è possibile identificare diverse funzioni di euristica che soddisfano la condizione di limite inferiore. Nei problemi della ricerca di itinerari, e non solo, un'altra possibile alternativa è la funzione di Manhattan che considera ogni mossa, o spostamento, sicuramente ammissibile. La ricerca guidata dal costo è un approccio molto particolare della ricerca A\* [ $f(n) = g(n)$ ], dove la funzione di euristica  $h$  è nulla per ogni nodo, è quindi uno stimatore super ottimistico.

Date due versioni dell'algoritmo A\* con due euristiche diverse  $h_1 < h_2$ , per tutti i nodi non obiettivo, si dice che l'algoritmo A\*<sub>2</sub> è più informato dell'algoritmo A\*<sub>1</sub>. Vale allora il teorema secondo cui al termine delle loro ricerche su un qualsiasi grafo con un percorso dal nodo iniziale  $n_o$  al nodo obiettivo, allora ogni nodo espanso da A\*<sub>2</sub> sarà anche espanso da A\*<sub>1</sub>. Quindi il primo algoritmo espande almeno tanti nodi quanto il secondo algoritmo, quindi l'algoritmo più informato A\*<sub>2</sub> è più efficiente di A\*<sub>1</sub>.

In generale è più conveniente scegliere un'euristica per cui l'algoritmo è più informato. Per determinare una funzione euristica, un buon approccio consiste nel determinare la funzione di costo effettivo per un problema simile all'originale con minori restrizioni sugli operatori. Infatti il costo effettivo di una soluzione in questi "Relaxed Problems" è una buona euristica per il problema originale. Ma il calcolo della funzione euristica potrebbe avere un calcolo computazionale elevato, quindi bisogna scegliere l'euristica considerando anche la loro complessità.

L'algoritmo A\* può essere eseguito sia in modalità Tree Search che Graph Search, ma in queste due modalità potrebbe trovare due soluzioni differenti allo stesso problema.

Si introduce quindi la condizione di consistenza, e si dice che una certa funzione euristica  $h$  obbedisce a questa condizione se per tutte le coppie di nodi  $n_j$ , successore di  $n_i$  nel grafo di ricerca

si ha:

$$h(n_i) \leq c(n_i, n_j) + h(n_j) \quad (2.5.2)$$

Dove  $c(n_i, n_j)$  è il costo dell'arco che congiunge i due nodi. Analogamente:

$$h(n_j) \geq h(n_i) - c(n_i, n_j)$$

Ovvero su un qualsiasi percorso, la stima del costo ottimo non può diminuire più del costo di un arco lungo quel percorso. Si può considerare anche come un tipo di disequaglianza triangolare:

$$h(n_i) \leq c(n_i, n_j) + h(n_j)$$

La condizione di consistenza impone che i valori della funzione di valutazione  $f$  dei nodi nell'albero di ricerca siano strettamente non-decrescenti all'allontanarsi dal nodo di partenza. Dati due nodi  $n_j$ , successore di  $n_i$ , se è soddisfatta si ha:

$$f(n_j) \geq f(n_i)$$

Dalla condizione di consistenza si aggiunge ad entrambi i membri  $g(n_j) = g(n_i) + c(n_i, n_j)$ :

$$\begin{aligned} h(n_j) &\geq h(n_i) - c(n_i, n_j) \\ h(n_j) + g(n_j) &\geq h(n_i) - c(n_i, n_j) + g(n_j) = h(n_i) - \cancel{c(n_i, n_j)} + g(n_i) + \cancel{c(n_i, n_j)} \\ f(n_j) &\geq f(n_i) \end{aligned} \quad (2.5.3)$$

Per cui spesso la condizione di consistenza sulla funzione euristica  $h$  viene spesso chiamata condizione monotona su  $f$ .

Se la condizione di consistenza è soddisfatta su  $h$ , allora quando l'algoritmo A\* espande un nodo  $n$  ha già trovato il percorso ottimo per  $n$ . In questo modo la ricerca su grafo non è differente dalla ricerca su un albero per quanto riguarda l'ottimalità della soluzione.

## 2.6 Algoritmo di Ricerca Locale: Hill-Climbing

Nei problemi precedenti quando l'algoritmo risolutivo raggiungeva uno stato obiettivo, il cammino verso quello stato rappresenta una soluzione del problema. Tuttavia in alcuni problemi lo stato obiettivo contiene tutte le informazioni rilevanti per la soluzione, dove il cammino è irrilevante. Come esempio si consideri il problema dell'otto regine, è indifferente il cammino attraverso gli stadi intermedi, solamente lo stato finale, la disposizione delle regine nello stato finale.

Algoritmi di ricerca locale si utilizzano per risolvere questo tipo di problemi. In questi problemi è sempre presente uno spazio degli stati ed uno spazio degli stati aventi ciascuno una sua valutazione. Si può immaginare questi stadi su una superficie del territorio, uno spazio dove l'altezza di questo stato rappresenta la sua valutazione. L'algoritmo quindi itera su ognuno di questi stati per cercare quello di altezza maggiore, o minore, identificando quindi la soluzione al problema, indipendentemente dal cammino preso per raggiungerla. Questi punti di massimo rappresentano dei picchi, i cui punti adiacenti sono strettamente minori dello stato di massimo. Quindi l'algoritmo che parte da uno stato iniziale deve cercare un massimo globale in questo spazio, determinando quale sia tra

i vari massimi locali ed i massimi locali piatto, e le “spalle” massimi locali “piatti”, prima di un massimo globale.

Questi algoritmi chiamati anche di miglioramento iterativo, si muovono sulla superficie cercando questi picchi, senza tenere traccia del cammino effettuato, tenendo solamente traccia dello stato attuale e dei suoi vicini o successori, gli stati immediatamente adiacenti. Bisogna formulare il problema in modo che l'algoritmo non rimanga bloccato tra due massimi locali.

Questo algoritmo segue sempre le colline più ripide, si muove sempre verso l'alto nella direzione dei valori crescenti, e termina quando raggiunge uno stato per il quale si ha un picco che non ha vicino stati di valore maggiore. Tuttavia questo algoritmo può rimanere intrappolato su massimi globali.

Non viene memorizzato lo stato corrente, solamente il valore attraverso nodi che contengono lo stato ed il suo valore.

Esistono diversi tipi di algoritmi di questo genere per evitare di rimanere bloccati su picchi locali, utilizzando diverse tecniche.

- Steepest Ascent Hill-Climbing;
- First-Choice Hill-Climbing;
- Random-Restart Hill-Climbing (Iterated Hill-Climbing);
- Stochastic Hill-Climbing.

### 2.6.1 Algoritmo Steepest Ascent Hill-Climbing

Si considera il seguente pseudocodice dell'algoritmo Hill-Climbing di tipo “Steepest Ascent”:

```
function HILL-CLIMBING(problem) returns a state local max
  current <- MAKE-NODE(INITIAL-STATE(problem))
  loop do
    next <- MAX(EXPAND(current))
    if VALUE(next) < VALUE(current) then return STATE(current)
    current <- next
  end
```

Dallo stato corrente si ricava il suo valore, in seguito comincia un ciclo che prende in considerazione tutti i successori e si sceglie come **next** il nodo di valore più alto. Se tutti i nodi adiacenti hanno un valore minore di **next** allora questo rappresenta la soluzione dell'algoritmo e l'algoritmo termina, tuttavia questo stato potrebbe corrispondere ad un massimo locale, invece se esiste uno stato adiacente di valore maggiore, questo diventa **next** e si passa alla nuova iterazione.

### 2.6.2 Algoritmo Random-Restart-Hill Climbing

L'algoritmo contiene una componente di ripartenza casuale, questo infatti conduce una serie di ricerche di Hill-Climbing partendo da stati generati casualmente. Questo algoritmo da un punto di

vista teorico è completo, poiché con una serie infinita di ripartenze, sicuramente l'algoritmo visita tutti gli stati del sistema, trovando sicuramente la soluzione ottima del problema.

```
function RANDOM-RESTART-HILL-CLIMBING(problem) returns a state solution

  t <- 0
  best <- MAKE-NODE(NULL)

  repeat
    local <- false
    current <- RANDOM(problem)
    repeat
      next <- MAX(EXPAND(current))
      if VALUE(next) > VALUE(current)
        then current <- next
      else local <- true
    until local
    t <- t + 1
    if VALUE(current) > VALUE(best)
      then best <- current
  until t = MAX
  return STATE(best)
end
```

Il ciclo interno ad ogni iterazione genera un ottimo locale, e prova ad evitare ottimi locali effettuando una nuova ricerca da un nuovo stato scelto casualmente. L'algoritmo è completo con probabilità tendente ad uno, poiché è possibile generi come stato iniziale uno stato obiettivo.

### 2.6.3 Algoritmo Stochastic Hill-Climbing

L'algoritmo Stochastic Hill-Climbing si ottiene modificando la procedura normale dell'algoritmo. Invece di valutare tutti i vicini dallo stato corrente, l'algoritmo sceglie casualmente uno solo dei suoi successori da valutare per determinare se si tratta il successore, ed in caso diventa il nuovo stato corrente *next*, questo viene accettato con una probabilità che dipende dalla differenza della valutazione tra i due punti:  $\Delta E = \text{VALUE}(\text{current}) - \text{VALUE}(\text{next})$ .

```
function STOCHASTIC-HILL-CLIMBING(problem) returns a state solution

  t <- 0
  current <- RANDOM(problem)
  best <- MAKE-NODE(NULL)

  repeat
    next <- RANDOM(EXPAND(current))
```

```

    if  $p = 1/(1 + e^{\Delta E/T})$ 
      then current <- next
      if VALUE(current) > VALUE(best)
        then best <- current
    t <- t + 1
until t = MAX
return STATE(best)
end

```

Il nuovo stato viene scelto con una probabilità  $p$ , calcolata come:

$$p = \frac{1}{1 + e^{\Delta E/T}} \quad (2.6.1)$$

In seguito dopo una serie di iterazioni l'algoritmo restituisce uno stato ottimo. L'algoritmo ha quindi un solo ciclo, e può scegliere un nuovo punto con una probabilità  $p$ , quindi anche di valore minore. Questa probabilità dipende da un parametro  $T$  costante durante l'esecuzione dell'algoritmo. Se vale 1, la probabilità di accettazione è sostanzialmente pari al 100%.

All'aumentare del valore di  $T$  la probabilità di accettazione tende al 50%, diventa quindi sempre meno importante la differenza della valutazione tra i due punti, effettivamente comporta una ricerca casuale, mentre al diminuire di  $T$ , la procedura rappresenta un semplice algoritmo Hill-Climbing.

In caso di stati di valore uguale, la probabilità è del 50%, se il valore dello stato **next** è minore, la probabilità diminuisce, mentre se il valore di **next** è maggiore dello stato corrente, la probabilità aumenta.

Bisogna trovare una "link function" tra l'intervallo  $\Delta E/T$  e la probabilità  $p$ .

La caratteristica di poter scegliere come passo uno stato peggiore questo algoritmo potrebbe evitare massimi locali.

#### 2.6.4 Algoritmo di Simulated Annealing

L'algoritmo di Simulated Annealing, introdotto nel 1983 da S. Kirkpatrick, C.D. Gelatt, Jr. e M.P. Vecchi nella rivista Science, è un algoritmo che migliora considerevolmente l'approccio dell'algoritmo precedente. Ha causato una vera e propria rivoluzione in termini di ottimizzazione, ebbe un enorme successo in ogni settore dell'informatica come l'algoritmo migliore per problemi di ricerca locale, anche solo nella sua versione base. Questo algoritmo è talmente importante che ogni anno vengono riuniti congressi annuali internazionali per discutere possibili ottimizzazioni.

Questo algoritmo prende il nome dall'analogia con il processo di metallurgia per temperare un materiale, questo processo infatti raggiungere uno stato di struttura cristallina ad energia minima. La differenza principale con l'algoritmo stocastico, è la possibilità di variare il valore di  $T$  diminuisce gradualmente durante l'esecuzione dell'algoritmo. Il valore di  $T$  parte da un valore elevato, per poi diminuire nel tempo, come se fosse la temperatura durante un processo di temperatura, per cui l'algoritmo si comporta in modo molto simile ad un normale hill-climber. Inoltre sceglie sempre uno stato se è migliore del punto corrente.

```

function SIMULATED-ANNEALING(problem) returns a state solution

  t <- 0
  current <- RANDOM(problem)
  best <- MAKE-NODE(NULL)

  repeat
    repeat
      next <- RANDOM(EXPAND(current))
      if VALUE(next) > VALUE(current)
        then current <- next
      if VALUE(current) > VALUE(best)
        then best <- current
      else if RANDOM[0,1) <  $e^{-\Delta E/T}$ 
        then current <- next
    until termination-condition
    T = g(T, t)
    t <- t + 1
  until halting-condition
  return STATE(best)
end

```

La probabilità di accettazione è leggermente diversa rispetto all'algoritmo precedente, poiché nel simulated annealing si considera solamente un semiasse dell'ascissa, dato che in caso **next** sia migliore non viene calcolata la probabilità di accettazione. Si accetta sempre uno stato di valore migliore, mentre l'accettazione di uno stato peggiore dipende da una probabilità. Quindi il mapping della link function non viene realizzata sull'intero intervallo  $(-\infty, +\infty)$ , ma solo su metà dell'ascissa, su  $[0, +\infty)$ , per cui è sufficiente utilizzare la funzione  $e^{-\Delta E/T}$  avente come dominio il semiasse  $[0, +\infty)$  e come codominio  $(0, 1]$ .

Questo ciclo interno viene effettuato un certo numero di volte fino ad una condizione di terminazione, che verrà trattata nelle implementazioni future. Finito questo ciclo si abbassa leggermente la temperatura tramite una funzione  $g$  e si incrementa il contatore delle iterazioni **t**. Anche la condizione di terminazione dell'algoritmo dipende dallo specifico problema e verranno trattate in futuro.

Molte implementazioni dell'algoritmo seguono la stessa sequenza di passi. Si assegna la variabile **T** alla temperatura massima, e si sceglie uno stato corrente casuale al primo passo. Si determina un successore assegnandolo direttamente se è migliore oppure tramite la funzione di probabilità e si ripete per un certo numero di cicli, e come passo finale si diminuisce la temperatura e si ripete dal secondo passo se la temperatura non ha raggiunto la temperatura minima. Quando la temperatura ha raggiunto la temperatura minima, si può scegliere se terminare l'algoritmo o ripeterlo un certo numero di volte, ripartendo dal primo passo.

## 3 Introduzione a Python

Python è un linguaggio di programmazione vastamente utilizzato nell'area dell'intelligenza artificiale e nel machine learning. Recentemente è diventato il linguaggio di programmazione più diffuso al mondo. Python è un linguaggio general-purpose, ideato da Guido van Rossum nel 1989, a più alto livello del C, poiché gestisce automaticamente le più fondamentali operazioni. Per cui è molto semplice e spesso utilizzato a fine didattico tra i primi linguaggi di programmazione insegnati.

La versione di Python utilizzata nel corso è la versione 3, nell'ambiente Anaconda. Può essere avviato tramite un interprete.

### 3.1 Ambienti di Sviluppo

I principali ambienti di programmazione di Python sono Anaconda, Eclipse e Google Colab. Tutte e tre permettono di programmare in modo semplificato. Esistono due approcci alla traduzione ed esecuzione dei programmi per i linguaggi di alto livello. Possono essere trasformati in un programma in linguaggio macchina ed eseguiti, compilazione, oppure ciascuna istruzione di alto livello può essere trasformata in almeno un' istruzione di linguaggio macchina ed eseguita direttamente, interpretazione. Storicamente la modalità interpretata viene associata alla programmazione su linea di comando, e solamente su rari casi viene utilizzata, anche se con Python ha cominciato ad essere riutilizzata.

Quindi non è presente una fase intermedia di compilazione o un programma compilato da eseguire, ma il programma viene eseguito direttamente istruzione per istruzione tramite un interprete. I linguaggi compilati sono specifici per ogni piattaforma, ma estremamente ottimizzati per quella esatta piattaforma. Mentre i linguaggi interpretati prevedono la distribuzione diretta del file sorgente, uguale per tutte le piattaforme, ma deve essere disponibile il programma interprete e sono generalmente più lenti poiché l'esecuzione attende ogni nuova istruzione da eseguire. In assenza di vincoli sul tempo dell'esecuzione il ritardo introdotto dall'interpretazione è accettabile.

Linguaggi come C e C++ sono linguaggi compilati, mentre Python è un linguaggio interpretato. Java ha un approccio ibrido generando il bytecode, una versione di linguaggio macchina eseguibile dalla JVM, elaboratore virtuale simile ad un hardware tradizionale, con una conversione quasi uno ad uno delle istruzioni.

Python è un linguaggio interpretato ad alto livello simbolico. Un programma Python può essere eseguito in modalità interattiva, dove ogni istruzione viene eseguita singolarmente, ed il suo output viene visualizzato sulla console in sequenza. Questo programma può essere interrotto, si possono rieseguire certe istruzioni, etc. Si può quindi modificare il flusso e le istruzioni sulla base dell'output ottenuto. Questo tipo di codice viene utilizzato per data analytics e exploration. Altrimenti all'interprete può essere passato l'intero programma come un file di estensione `.py`. L'output viene visualizzato solamente se si invoca la funzione di stampa.

Per utilizzare la modalità interattiva serve un ambiente di interazione, l'interprete più diffuso è Cpython.

Può essere scaricato su tutte le piattaforme, aprendo una console si può entrare nella modalità interattiva con il comando `python3`. Per uscire da questa modalità è sufficiente invocare la funzione `quit()`.

Anaconda è una distribuzione di Python che semplifica la gestione ed il deployment delle librerie. Permette di creare ambiente in cui installare versioni specifiche delle librerie compatibili con certi software. Il gestore di queste librerie o package si chiama conda e si usa principalmente da linea di comando e gestisce queste librerie e le mantiene aggiornate. Può essere installato insieme a Cpython.

I principali comandi in Anaconda sono:

```
# crea un nuovo environment di versione specificata:
conda create -n nomeEnv python=versionePython
# visualizza la lista di environment creati:
conda info -envs
# attiva l'ambiente di sviluppo
source activate nomeEnv
# installa il package specificato all'interno dell'ambiente corrente
conda install nomePackage=versionePackage
```

Se non viene specificata la versione, sia per Python che per i pacchetti da installare, conda utilizzerà l'ultima versione come default.

Le principali librerie utilizzate in questo corso sono pandas, matplotlib, NumPy, SciPy, IPython e scikit-learn. Dentro Anaconda c'è un'interfaccia browser per avere un ambiente di programmazione per programmare chiamato Jupiter, dove è possibile attivare la linea di comando per programmare in modalità interattiva. Si può attivare con il comando `jupyter notebook`, che mostra a schermo l'indirizzo a cui collegarsi per accedere, generalmente 127.0.0.1. Si analizzerà in seguito insieme all'ambiente Google Colab, anch'esso su un'interfaccia browser.

Sui notebook è possibile inserire del testo, ed accetta macro e comandi in formato Markdown, e funzioni scritte in LaTeX.

Eclipse è un ambiente di sviluppo storico, multi-linguaggio, installando dal marketplace il plugin pydev per poter programmare in Python.

Google Colab è un ambiente di sviluppo fornito gratuitamente da Google, utilizzando risorse cloud, su di un singolo notebook, direttamente collegato all'account google. Ha un'interfaccia a blocchi interattiva, da data scientist, e permette di scegliere il runtime e la GPU da utilizzare per eseguire il programma. Il codice viene salvato su un notebook, quindi su un file di estensione `.ipynb`, ed è possibile passare al Python tradizionale e viceversa. Su dispositivi mobile è possibile programmare con un'interfaccia molto simile ad un computer.

Si divide il programma in celle per dividere il codice e permettere di elaborare singoli blocchi, senza dover ricalcolare tutti i blocchi già eseguiti. Questo inoltre rende il codice più chiaro e comprensibile.

JupyterLab è un estensione di Jupyter che aggiunge funzionalità per programmare a progetto per visualizzare e gestire diverse tipologie di file.

Si possono utilizzare anche ambienti virtuali come Docker per eseguire il codice su container isolati. Inoltre esistono ambienti dedicati all'ambito scientifico con un'analisi più dettagliata dei dati prodotti come spyder, "The Scientific Python Development Environment".



## 3.2 Operatori

Alla creazione di una variabile non è necessario definirne il tipo, il nome identificativo è arbitrario e può contenere numeri, ma non cominciare con un numero, viene consigliato di utilizzare un carattere minuscolo come primo carattere del nome. Esistono 33 parole chiave, non utilizzabili come nomi di variabili. Si può assegnare un valore ad una variabile tramite l'operatore `=`, senza specificarne il tipo.

Esistono una serie di operatori aritmetici come `+`, `-`, `*`, `\`, `**`, per l'elevamento a potenza, `%` per il modulo. Una differenza tra Python 2 consiste nella gestione della divisione, infatti in Python 2 viene considerata solo la parte intera dell'operazione. Per ottenere lo stesso risultato esiste l'operatore `//`, chiamato "floor division".

Gli operatori seguono un ordine di precedenza naturale, come la sintassi moderna matematica:

1. Parentesi;
2. Elevamento a potenza;
3. Moltiplicazione e divisione;
4. Addizione e sottrazione;
5. Operatori con lo stesso ordine valutati da sinistra verso destra.

In Python sono presenti tutti gli operatori booleani del C come `==` ed operatori di confronto come `<`, `<=`, `>`, `>=`, in aggiunta sono presenti altri operatori `is` ed `is not`. Inoltre sono presenti due versioni degli operatori logici `&&` e `and`, `—` e `or` e `!=` e `not`. Come in C un qualsiasi valore diverso da zero corrisponde al booleano `True`, inoltre associa numericamente questo valore ad 1, mentre `False` a 0. È possibile convertire un tipo di dato booleano ad un tipo di dato numerico come un intero `int` o un numero reale `float`.

Si possono inserire dati dall'utente tramite la funzione `input()` e la funzione `raw_input()` per lo stesso comportamento di Python 2, e si può convertire in un tipo specifico con `tipo(var)`. I commenti vengono realizzati tramite il carattere `#`.

## 3.3 Istruzioni Condizionali

In Python per identificare funzioni o istruzioni condizionali non si usano parentesi, ma si indenta di quattro posizioni. Dopo la condizione dell'istruzione condizionale vanno inserite dei due punti `::`

```
if condizione:
    # corpo dell'if
else:
    # corpo dell'else
```

Si possono gestire le eccezioni con il costrutto `try` ed `except`:

```
try:
    # corpo del try
except:
    # corpo dell'except
```

### 3.4 Funzioni Built-In, Moduli e Definizione di Funzioni

In Python sono integrate tantissime funzioni utili, per svolgere attività comuni, utilizzabili senza doverle definire, queste sono funzioni “built-in”. Per invocare funzioni presenti in un certo modulo e non built-in si utilizza la notazione puntata `nomeModulo.nomeFunzione()`. Alcune tra le funzioni built-in più utili sono `max()` e `min()` che restituiscono il carattere più grande e più piccolo in una stringa; la funzione `len()` che restituisce la lunghezza della stringa. Per convertire variabili in certi tipi è già stata mostrata la funzione `tipo()`, dove al posto di `tipo` si inserisce il tipo specifico, si usa `str` per convertire in una stringa.

Per importare moduli contenenti altre funzioni si utilizza `importa` seguito dal nome del modulo da scaricare, che crea un object module con quel nome, si può rinominare seguendo questa istruzione con `as` seguito dall’alias del modulo. Si utilizza un alias per semplificare la notazione puntata.

Dati gli algoritmi analizzati precedentemente, si nota la necessità di introdurre generatori di numeri casuali. La maggior parte di generatori casuali, sono deterministici, ovvero dato lo stesso input, generano la stessa sequenza di numeri casuali. Si utilizzano quindi numeri pseudo-casuali, generati da un calcolo deterministico, ma non è quasi possibile distinguerli da numeri generati casualmente. In Python esiste il modulo `random` contenente funzioni pertinenti alla generazione di numeri casuali. La funzione `random()` genera un numero casuale tra 0.0, compreso e 1.0, non compreso. Un’altra funzione `randint()` accetta due parametri, estremi dell’intervallo, inclusi, per generare un numero intero tra loro compreso. Con la funzione `choice()` si può scegliere un elemento casualmente da una sequenza passata come argomento.

Per definire nuove funzioni si utilizza la parola chiave `def`, specificando il nome, tra parentesi tonde gli argomenti ed i due punti, indentando di quattro posizioni per scrivere il corpo della funzione:

```
def nomeFunzione(listaArgomenti):
    # corpo della funzione
    # resto del codice
```

Dopo aver passato degli argomenti ad una funzione, questi vengono assegnati a delle variabili locali. Si può utilizzare anche una variabile come argomento. Inoltre tutte le aggiunte possibili alle funzioni built-in, si possono effettuare sulle funzioni definite dall’utente.

Si dividono le funzioni in due tipi “fruitful function”, funzioni produttive, e “void function”, funzioni vuote, le prime restituiscono un valore, le seconde non restituiscono valore. Le prime quindi vengono usate per assegnare o inizializzare variabili. Se si tenta di assegnare il risultato di una void function ad una variabile, viene ottenuto un valore chiamato `None`. Questo valore ha un suo proprio tipo. Per definire una funzione produttiva, nel corpo si inserisce la parola chiave `return` seguita dai parametri da restituire come risultato della funzione.

### 3.5 Cicli

Si possono realizzare cicli tramite il costrutto `while` o `for`, seguito da una condizione booleana e dai due punti ::

```
while condizione:
    # corpo del ciclo
```

Si può interrompere il ciclo con `break`, e si può saltare l'iterazione corrente con `continue`. Quando bisogna iterare su una collezione, un insieme di elementi, si può realizzare un ciclo "for-each":

```
for elemento in collezione:
    # corpo del ciclo
```

### 3.6 Stringhe

Le stringhe sono sequenze di caratteri, indicizzati come fosse un array:

```
stringa[i] # (i+1)-esimo carattere
```

La funzione già discussa `len()` restituisce il numero di caratteri di una stringa, anche se può essere utilizzata per altri tipi di dati come dizionari. Poiché è strutturata come un array è possibile scandire ogni carattere della stringa individualmente con un ciclo for-each.

Talvolta è comodo accedere ad una sottostringa, o "slice", della stringa di partenza. La selezione di una sottostringa è simile alla selezione di un carattere, utilizzando due indici divisi da due punti per indicare l'inizio e la fine della sottostringa, il primo estremo è compreso, mentre il secondo no:

```
stringa[i:j] # slice contenente i caratteri da i a j-1
```

A volte si ha la necessità di realizzare una sottostringa che parte dall'inizio o la fine della stringa originaria, per effettuarlo si può omettere l'estremo corrispondente:

```
stringa[:j] # slice dall'inizio della stringa
stringa[i:] # slice fino alla fine della stringa
```

I valori di una stringa sono immutabili una volta definiti, per cui non è possibile modificarne il valore accedendo tramite indice, verrà sollevato un messaggio di errore. Si può modificare una stringa realizzando una nuova stringa, come variante, tramite l'operatore di concatenazione `+`:

```
stringa = stringa_1 + stringa_2
```

L'operatore `*` applicato su una lista, la replica un certo numero di volte specificato.

L'operatore `in` è estremamente importante, permette di individuare se una stringa è sottostringa di un'altra, restituisce un valore booleano vero o falso:

```
stringa_1 in stringa_2
```

Si possono inoltre utilizzare operatori di confronto tra stringhe all'uguaglianza con `==`, oppure con `<`, `>`, per confrontarle in ordine alfabetico. In Python le maiuscole vengono prima delle minuscole. Le stringhe sono degli oggetti che oltre alla sequenza di caratteri contengono oltre i dati anche i metodi disponibili per ogni istanza dell'oggetto stringa.

Con la funzione `type` si ha la possibilità di identificare il tipo dell'oggetto su cui viene operata e `dir` mostra i metodi disponibili. Per utilizzare un metodo si utilizza la notazione puntata con il nome dell'istanza dell'oggetto. Uno dei metodi sulle stringhe è `find()` che prende come argomento una sottostringa, ed un indice opzionale da cui cercare il carattere, e restituisce la prima posizione dell'occorrenza della sottostringa specificata.

Il metodo `strip()` rimuove lo spazio bianco prima e dopo la stringa. Il metodo `startswith()` restituisce un valore booleano se la stringa comincia con la sottostringa passata come argomento. Il metodo `capitalize()` consente di impostare a maiuscolo il primo carattere, per mettere tutti i caratteri in maiuscolo si utilizza il metodo `upper()`.

L'operatore "format" % permette di costruire stringhe formattando la stringa rispetto a dati contenuti in altre variabili. All'interno di una stringa chiamata "format string" si può inserire questo operatore seguito da una lettera per specificare il tipo di dato associato, chiamate "format sequences". Questa viene utilizzata come il primo argomento, mentre il secondo operando è la variabile da formattare, questo produce una stringa:

```
>>> x = 1
>>> 'il numero è %d' % x
'il numero è 1'
```

Se in una stringa compare una sequenza di format sequences si specifica l'operando come una sequenza, una tupla, tra parentesi tonde, separando gli elementi con virgole:

```
>>> x = 1
>>> 'il %s è %d' % ('numero', x)
'il numero è 1'
```

La tupla deve corrispondere in numero, ordine e tipo alle format sequences nella stringa. Si utilizza `%d` per numeri interi `%g` per floating point e `%s` per stringhe. Utilizzando numeri è possibile specificare il numero di cifre significative da utilizzare per rappresentare il numero, inserendo un numero subito dopo l'operatore format.

Esiste un nuovo operatore format, scrivendo tra parentesi graffe : seguito dal carattere corrispondente al tipo, dopo la stringa si invoca il metodo `format()` specificando nell'argomento il valore da inserire:

```
>>> 'il {:s} è {:d}' .format('numero', 1)
'il numero è 1'
```

Questa implementa tutte le caratteristiche già trattate per il format %.

### 3.7 Classi e oggetti

La fonte di ispirazione della programmazione orientata agli oggetti è dovuta a lavori realizzati dalla comunità sull'intelligenza artificiale, nell'ambito della rappresentazione della conoscenza quando

Marvin Minsky propose il formalismo dei frames nel 1975. Con tutte le caratteristiche dei moderni approcci di programmazione orientata agli oggetti.

Questo paradigma di programmazione consiste nel mantenere i dati e le funzioni operabili su di essi vicine, raggruppate in entità chiamate oggetti. Una classe è una descrizione di un insieme di oggetti che hanno lo stesso comportamento, ed un oggetto è una singola istanza di questa classe, dove è definita anche un insieme di metodi da poter usare su questi oggetti.

L'insieme di tutti i metodi resi disponibili da una classe, con la descrizione del loro comportamento, si chiama interfaccia pubblica della classe.

Una classe in Python viene definita tramite la parola chiave `class` seguita dal nome della classe, utilizzando la notazione dove ogni parola del nome ha la prima lettera maiuscola, "Camel-case":

```
class NomeClasse:
    pass
```

`pass` è un placeholder ed indica dove sarà scritto codice, questa rappresenta una definizione di una classe minimale, definendo solamente il nome e nessun altro comportamento.

Si può istanziare una nuova classe dato il suo nome:

```
NomeClasse()
```

Questo istanzia un oggetto appartenente alla classe specificata, allocata in una nuova zona di memoria, una variabile assegnata a questa nuova istanza mantiene il riferimento a questa zona di memoria, quindi si ha:

```
>>> a = NomeClasse()
>>> b = NomeClasse()
>>> a == b
False
```

Poiché quest'operazione di confronto effettuata sulle variabili ha azione sugli indirizzi di memoria contenuti. Le proprietà che una classe deve avere può essere definite utilizzando un metodo particolare costruttore chiamato `__init__()`. Questo inizializza una nuova istanza della classe ed assegna i valori al suo stato. Metodi con un doppio carattere `_` vengono chiamati "Dunder Methods", per "Data Underlined". Al costruttore vengono passati i parametri che si vuole, ma è sempre necessario passare la variabile chiamata `self` per definire nuovi attributi sull'oggetto. Poiché quando viene creata un'istanza di una classe, viene salvata automaticamente a questo parametro.

```
class NomeClasse:
    def __init__(self, attributo1, ..., attributoN):
        self.attributo1 = attributo1
        # ...
        self.attributoN = attributo N
```

Ognuna di queste assegnazioni in realtà crea anche un attributo chiamato con il nome utilizzato nella notazione puntata, riferita a `self`, e contenente il valore assegnatoli dal parametro passato come input al costruttore. `self` deve essere sempre il primo parametro.

Gli attributi creati in questo modo vengono chiamati “Instance Attributes”, attributi d’istanza, che sono attributi unici per ogni istanza della classe, chiamati al momento di costruzione dell’istanza. Ma possono essere definiti attributi comuni per tutte le classi, chiamati “Class Attributes”, all’interno del corpo della classe, senza inserirli all’interno del costruttore:

```
class NomeClasse:
    attributo = 'valore'
    def __init__(self, attributo1, ..., attributoN):
        self.attributo1 = attributo1
        # ...
        self.attributoN = attributo N
```

Ogni istanza avrà gli stessi class attributes, contenenti gli stessi valori. Per istanziare una classe, avendo un costruttore, bisogna passare i parametri richiesti al momento di creazione:

```
>>> istanza = NomeClasse('valore1', ..., 'valoreN')
```

Si passa un parametro in meno poiché **self** viene gestito internamente, al momento della creazione dell’oggetto. Si può accedere ad un certo attributo sempre con la notazione puntata, e si può modificare. Si può modificare sia un attributo di classe o istanza.

Gli instance methods sono funzioni definite all’interno di una classe, possono essere invocati da una qualsiasi istanza della classe.

È buona norma utilizzare un metodo come **description()** per fornire una stringa contenente informazioni utili per l’istanza della classe. Ma questo non è il modo migliore per ottenere. Si utilizza il metodo **.\_\_str\_\_()** per restituire queste informazioni in formato di stringa:

```
class NomeClasse:
    attributo = 'valore'

    def __init__(self, attributo1, ..., attributoN):
        self.attributo1 = attributo1
        # ...
        self.attributoN = attributo N

    def __str__():
        return f'{self.attributo1}, ..., {self.attributoN}'
```

Quando viene invocata una funzione di stampa su un’istanza di questa classe, viene stampata la stringa definita da questo metodo:

```
>>> istanza = NomeClasse('valore1', ..., 'valoreN')
>>> print(istanza)
valore1, ..., valoreN
```

## 3.8 Collezioni

### 3.8.1 Liste

Una lista è una sequenza di item di qualsiasi tipo, per cui a differenza di una stringa il tipo non è omogeneo. All'interno di una lista è possibile inserire una lista, creando una lista annidata. Una lista si realizza specificando gli elementi tra parentesi quadre:

```
[item_1, item_2, [item_3, item_4]]
```

Una lista che non contiene elementi si chiama lista vuota mediante [], oppure con il comando `list()`. Si può assegnare ad una variabile una lista allo stesso modo di un valore:

```
>>> lista = [item_1, item_2, [item_3, item_4]]
```

Le liste al contrario delle stringhe sono modificabili, con la stessa sintassi per le stringhe, specificando tra parentesi quadre l'indice corrispondente, ed è possibile aggiornare il valore contenuto in questo indice. Se ad un certo indice è presente una lista, questo elemento viene trattato come lista e quindi si possono utilizzare due coppie di parentesi quadre per indicizzare questi elementi annidati:

```
>>> lista[2][2]
item_4
```

Si può considerare una lista come un mapping tra indici ed elementi, come nelle stringhe possono essere variabili ed espressioni, se si prova a leggere o scrivere ad un indice non esistente si solleva un `IndexError`. Mentre se si utilizza un indice negativo, si conta all'indietro partendo dalla fine della stringa. L'operatore `in` ha un funzionamento analogo a quello per le stringhe.

Un modo per scandire la lista è con un ciclo for-each, come per le stringhe, oppure iterando manualmente su ogni elemento. Questo si può effettuare tramite la funzione `range()` che permette di iterare su un insieme di valori, fornito l'ultimo come argomento:

```
for i in range(len(lista)):
    # corpo del for
```

L'operatore di concatenazione vale anche per le liste, analogamente per l'operatore `*` che replica una lista un certo numero di volte. L'operatore slice può essere utilizzato anche sulle liste per ottenere una sottolista:

```
lista[i:j] # sottolista da i a j-1
lista[:j]  # sottolista da 0 a j-1
lista[i:]  # sottolista da i a len(lista)
lista[:]   # copia della lista
```

L'operatore slice può essere utilizzato per aggiornare una sequenza di elementi in una lista. La sottolista da aggiornare deve avere lo stesso numero dei valori della lista originaria.

Come per le stringhe, le liste sono oggetti e contengono metodi built-in. Il metodo `append()` aggiunge un nuovo elemento alla fine di una lista, il metodo `extend()` prende come argomento una

lista e la concatena alla lista su cui si è operato. Se si passa una lista al metodo `append()` si genera una lista annidata. Si può ordinare gli elementi di una lista dal minore al maggiore con `sort()` in ordine alfabetico. Molti dei metodi su liste, come questi, sono metodi void.

Per estrarre un elemento dalla lista si utilizza il metodo `pop()` che estrae dalla lista l'ultimo elemento, oppure l'elemento di indice specificato, lo restituisce rimuovendolo dalla lista:

```
>>> lista.pop(1)
item_2
>>> print(lista)
[item_1, [item_3, item_4]]
```

L'operatore `del` rimuove un elemento dalla lista:

```
>>> del lista[1]
>>> print(lista)
[item_1, [item_3, item_4]]
```

Utilizzando lo slice si può rimuovere una sottostringa allo stesso modo:

```
>>> del lista[:1]
>>> print(lista)
[[item_3, item_4]]
```

Se non si conosce l'indice dell'elemento da cancellare si può utilizzare il metodo `remove()`, ma non restituisce alcun valore:

```
>>> lista.remove(item_2)
>>> print(lista)
[item_1, [item_3, item_4]]
```

Alcuni metodi come `sum()` possono essere invocati solamente se la lista contiene solo numeri, altri metodi possono essere invocati se la lista contiene dati confrontabili tra di loro. Si può convertire una stringa in una lista tramite la funzione `list()`, suddividendo la stringa in caratteri singoli, mentre si può suddividere in parole singole, separate da spazi con la funzione `split()`. Si può specificare un delimitatore in questa funzione `split()` da utilizzare al posto del delimitatore di default spazio. La funzione inversa della `split()` è il metodo `join()`, permette di realizzare una stringa, applicato su una stringa contenente il delimitatore da utilizzare per separare gli elementi. Si può concatenare senza spazi con la stringa vuota `[]`.

In Python si possono confrontare due oggetti con l'operatore `is`, se si creano due stringhe di contenuto uguale, Python realizza un unico oggetto stringa, mentre se vengono realizzate due liste, contenenti gli stessi elementi, sono due oggetti distinti. Le due liste sono equivalenti poiché hanno lo stesso valore, ma non sono identiche, poiché non sono lo stesso oggetto. Due oggetti identici sono anche equivalenti.

Si può assegnare ad una variabile ad un'altra in modo che entrambe si riferiscano allo stesso oggetto, e quindi anche se si tratta di liste il confronto `is` sarà verificato. Un oggetto che ha più di riferimenti ha un "alias", più di un nome, e si dice un oggetto "aliased". Se gli oggetti sono mutabili



questo può essere fonte di errore. Si consiglia di evitare di utilizzare alias quando si lavora con liste ed oggetti mutabili. Si distingue tra operazioni che creano nuove liste ed operazioni che modificano lista, come l'operatore di concatenazione `+` ed il metodo `append()`, quest'ultimo modifica la lista su cui è invocato analogamente di `extend()`, mentre l'operatore di concatenazione genera una nuova lista.

C'è uno speciale costrutto per accedere a tutti gli elementi di una lista effettuando la stessa operazione su di loro e memorizzare i nuovi elementi in un'altra lista:

```
>>> lista_1 = [1, 2, 3]
>>> lista_2 = [item * 2 for item in lista_1]
>>> print(lista_2)
[2, 4, 6]
```

Quando bisogna accedere contemporaneamente a più liste si può utilizzare la sintassi che utilizza la parola chiave `zip` per trattare la sequenza di liste passate contemporaneamente:

```
for item_1, item_2, ... in zip(lista_1, lista_2, ...):
    # corpo del ciclo, per ogni iterazione si ha
    # item_1 = lista_1[i], item_2 = lista_2[i], ...
```

In questo modo si può iterare su tutte le liste insieme, senza dover specificare per ognuna il loro indice. Questa funziona associa le liste e termina alla lista più corta.

### 3.8.2 Dizionari e Set

Un dizionario è simile ad una lista, ma più generale, permette di accedere agli elementi contenuti tramite una chiave definita a priori. Rappresenta un mapping tra due insiemi, uno di chiavi ed uno di valori, ogni chiave individua un valore, ogni item del dizionario è quindi una coppia tra una chiave ed un valore. Può essere costruito un dizionario vuoto tramite la funzione `dict()`, e se stampato produce `{}`.

```
>>> str2num = dict()
>>> print(str2num)
{}
>>>
```

Per aggiungere un item nel dizionario, si utilizzano parentesi quadre per specificare la chiave e si assegna il suo valore contenuto, come se fosse l'indice di una lista:

```
>>> str2num['uno'] = 1
>>> print(str2num)
{'uno' : 1}
```

Questo formato di output di un dizionario è un formato che può essere utilizzato per popolare un dizionario:

```
>>> str2num = {'uno' : 1, 'due' : 2, 'tre' : 3}
```

Ma l'ordine in cui sono stati inseriti gli elementi non viene mantenuto, questo non costituisce un problema, poiché gli elementi non vengono identificati dalla loro posizione, ma attraverso una chiave. Se si prova ad accedere ad un elemento la cui chiave non esiste viene sollevato un errore. Per accedere ad un valore nel dizionario si utilizza la chiave come fosse un indice:

```
>>> str2num['uno']  
1
```

Si può utilizzare la funzione `len()` per determinare il numero di elementi in un dizionario e la parola chiave `in` per determinare se una data chiave è presente in un dizionario, non controlla se è presente un valore:

```
>>> len(str2num)  
3
```

Per visualizzare i valori si può utilizzare il metodo `values()` che restituisce elementi su cui si può costruire una lista, su cui è possibile verificare la presenza di certi valori tramite `in`. Questo operatore utilizza diversi algoritmi per le liste e per i dizionari, utilizza una ricerca lineare per le liste ed una tabella di hash per i dizionari.

Una tupla è una sequenza di valori, simili ad una lista, di qualsiasi tipo, indicizzati da interi. Le tuple a differenza delle liste sono immutabili, ed è possibile confrontarle tra di loro, quindi è possibile siano utilizzate per realizzare chiavi di dizionari. Si definisce una tupla come una sequenza di valori divisi da virgole, ma convenzionalmente si tende ad includere questi valori in parentesi graffe per semplificare la leggibilità:

```
>>> t={1, 2, 3}
```

Per creare una tupla con un unico valore, bisogna comunque inserire la prima virgola, altrimenti si creerebbe o una stringa o un dato di tipo uguale al tipo del primo elemento.

Si può creare una tupla vuota con la funzione `tuple()`, invece è possibile inserire un argomento in questa funzione per generare una tupla dagli elementi della sequenza. Se si tratta di una stringa ogni carattere viene trattato come un elemento distinto, quindi viene scomposta. I vari operatori analizzati per le liste si possono utilizzare anche per le tuple:

```
t[i]          # valore all'indice i  
t[i:j]        # valori dall'indice i a j-1  
t[:j]         # valori dall'indice 0 a j-1  
t[i:]         # valori dall'indice i a len(t)  
t[:]          # copia della tupla
```

Se si provasse a modificare un elemento verrebbe sollevato un errore, è comunque possibile creare una copia sostituendo solamente certi valori, così come per le stringhe, effettuando una concatenazione.

Gli operatori di confronto funzionano sulle tuple e su altre sequenze su Python. Cominciano a confrontare dal primo elemento di ciascuna sequenza, se sono uguali passa all'elemento successivo fino a trovare una coppia di elementi diversi. Quindi gli elementi successivi non vengono presi in considerazione.

Si può utilizzare la funzione `sort()`, può essere utile con lo schema DSU, “Decorate, Sort, Undecorate”. Secondo questo approccio si decora una sequenza costruendo un elenco di tuple con una o più chiavi di ordinamento che precedono gli elementi della sequenza, si ordina tramite questa funzione ed in seguito si rimuove questa decorazione estraendo gli elementi della sequenza.

Una caratteristica di Python consente di assegnare più di una variabile alla volta, potendo inserire una tupla sul lato sinistro dell’assegnazione, senza parentesi tonde per convenzione. Un’applicazione interessante permette di scambiare i valori tra due variabili con una singola istruzione:

```
>>> a, b = b, a
```

Il numero di valori a destra e sinistra dell’assegnazione devono corrispondere per non sollevare errori. Un altro metodo chiamato `items()` permette di restituire una lista di tuple dove ogni tupla è una coppia chiave-valore. Questa lista non è ordinata, ma essendo le tuple comparabili, si può ordinare con il metodo `sort()`.

Gli insiemi sono collezioni non ordinate di altri oggetti. L’istruzione per creare un insieme è `set()`, inserendo una sequenza di elementi. In un insieme non possono esserci elementi ripetuti quindi vengono rimossi al momento della creazione.

Si possono effettuare le classi operazioni su insiemi, con metodi di unione `union()`, intersezione `intersection()`, differenza `difference()` e differenza simmetrica `symmetric_difference()`.

### 3.8.3 Liste Concatenate

Il tipo di dato astratto lista concatenata è diversa dal tipo di dato predefinito lista già presente su Python. Si parlerà di liste singolarmente concatenate. Questa lista è composta da un nodo, contenente il dato ed un puntatore all’elemento successivo di questa struttura concatenata. L’ultimo nodo della lista contiene un puntatore nullo. Per accedere agli elementi si hanno due puntatori di accesso, `head` e `tail`, per accedere in testa o in coda alla lista.

È comodo avere un accesso facilitato a queste due posizioni per effettuare facilmente operazioni di aggiunta, rimozione, modifica in testa o in coda. Un nodo è una classe contenente questi due unici attributi:

```
class Node:
    data = ''
    next = None

    def __init__(self, data, next):
        self.data = data
        self.next = next
```

L’utilizzo del paradigma orientato agli oggetti non è l’unico modo per definire una lista concatenata, ma è il modo utilizzato durante questo corso.

Una lista concatenata è una classe contenente i riferimenti alla testa ed alla coda della lista, inizialmente assegnati al valore nullo, per generare una lista vuota:

```
class SinglyLinkedList:
    def __init__(self):
        self.__head = None
        self.__tail = None
```

In realtà il campo dati è un riferimento ad un altro oggetto contenente i dati, ma per semplificare si considera il campo dati direttamente contenente questi.

Per scandire la lista elemento per elemento si può utilizzare un puntatore ausiliario, inizialmente assegnato al primo elemento della lista:

```
p = self.__head
```

Per passare all'elemento successivo bisogna assegnare questo puntatore al valore del puntatore a `next` contenuto nel nodo corrente:

```
p = p.next
```

Per scandire la lista quindi si utilizza un ciclo, fino a quando il valore contenuto in questo puntato `p` non corrisponde al valore `None`:

```
class SinglyLinkedList:
    def __init__(self):
        self.__head = None
        self.__tail = None

    def operazione():
        p = self.__head
        while p != None:
            # operazioni sul nodo
            p = p.next
```

Per cominciare a scandire la lista da un punto preciso si può inserire un parametro di classe `Node` a cui assegnare `p`. Si vogliono realizzare i seguenti metodi sulle liste singolarmente concatenate:

- `append`: Inserimento in coda;
- `insert_head`: Inserimento in testa;
- `insert_position`: Inserimento in una posizione intermedia;
- `delete`: Cancellazione di un elemento.

```
class SinglyLinkedList:
    def __init__(self):
        self.__head = None
        self.__tail = None
```

```
def append(self, newNode)
    if self.__tail != None:
        self.__tail.next = newNode
    if self.__head == None:
        self.__head = newNode
    self.__tail = newNode
    # newNode.next = None

def insert_head(self, newNode):
    if self.__head == None:
        self.__tail = self._head
    newNode.next = self.__head
    self.__head = newNode

def insert_position(self, newNode, position):
    p = self.__head
    i = 0
    while i < position - 1 and p != None:
        p = p.next
        i += 1
    newNode.next = p.next
    p.next = newNode

def delete(self, position):
    if self.__head == None:
        return None
    if position == 0:
        self.__head = self._head.next
        return None
    p = self.__head
    i = 0
    while i < position - 1 and p.next != None:
        p = p.next
        i += 1
    if i != position - 1:
        return None
    p.next = p.next.next
```

### 3.8.4 Liste Ordinate, Pile e Code

Una lista ordinata è una lista ai cui elementi sono ordinati rispetto ad una certa chiave. Ogni inserimento deve quindi rispettare l'ordine della lista.

Questa lista si può implementare tramite la liste concatenate:

```
class SortedLinkedList:
    def __init__(self):
        self.__head = None
        self.__tail = None
```

Si suppone di poter effettuare un inserimento in posizione intermedia. Bisogna scandire la lista fino a quando il dato da inserire non è maggiore del dato contenuto nel nuovo nodo da inserire. Si considera la chiave e la lista ordinata in modo crescente:

```
def add(self, newNode):
    if self.__head == None or self.__head.data > newNode.data:
        self.insert_head(newNode)
        return None
    elif self.__tail.data < newNode.data:
        self.append(newNode)
        return None

    p = self.__head
    while p.next != None and p.data < newNode.data:
        p = p.next
    newNode.next = p.next
    p.next = newNode
```

Una pila è un tipo di dato astratto basato sul modello delle liste che gestisce gli elementi, implementando la disciplina LIFO, "Last-In, First-Out". Le operazioni vengono effettuate su un estremo della lista chiamato testa **top**.

La classe **Node** è sempre la stessa, viene definita con un attributo **None** inizializzato a **None** e può essere comodo avere a disposizione una variabile che indica il numero di elementi nella pila:

```
class Stack:
    __top = None
    __size = 0
```

Ogni inserimento in una pila è un inserimento in testa, si chiama "push":

```
def push(self, newNode):
    newNode.next = self.__top
    self.__top = newNode
    self.__size += 1
```

Per effettuare un'estrazione si estrae un elemento dalla testa, e generalmente il dato estratto è di interesse:

```
def pop(self):
    p = self.__top
    if p != None
```

```
        self.__top = p.next
        p.next = None
    self.__size -= 1
    return p
```

Le code sono un altro tipo di dato astratto che utilizzano il modello delle liste, ma applicano la disciplina FIFO, “First-In, First-Out”:

```
class Queue:
    __head = None
    __tail = None
    __size = 0
```

L’inserimento di un elemento avviene sempre in coda, e si chiama “enqueue”:

```
def enqueue(self, newNode):
    if self.__head == None:
        self.__head = newNode
    else:
        self.__tail.next = newNode
    self.__tail = newNode
    self.__size += 1
```

Per rimuovere il primo elemento, si utilizza la funzione “dequeue”:

```
def dequeue(self):
    p = self.__head
    if p != None:
        self.__head = p.next
        p.next = None
    self.__size -= 1
    return p
```