



Università degli Studi di Salerno

Dipartimento di Informatica

Tesi di Laurea Triennale in

Informatica

**UN ALGORITMO PER L'INFERENZA DI
DIPENDENZE FUNZIONALI RILASSATE:
RICERCA DI PATTERN MINIMALI**

Relatore

Chiar.mo Prof. Vincenzo Deufemia

Candidato

Raffaele Ceruso

Secondo Relatore

Dott.sa Loredana Caruccio

Matr. 0512102927

Anno accademico 2016-2017

*Alla mia famiglia,
per essermi stata sempre vicino lungo questo percorso.*

*Alla mia famiglia,
per i sacrifici fatti affinché tutto questo potesse realizzarsi.*

*A Miriam,
per avermi supportato quando c'era il bisogno.*

Abstract

Le basi di dati sono utilizzate in larga scala in moltissimi aspetti dell'ambito tecnologico, per questo motivo durante la loro progettazione ci sono aspetti essenziali da prendere in considerazione per assicurare un servizio quanto più efficiente possibile. La qualità dei dati contenuti è un servizio che di certo una buona base di dati deve garantire, motivo per il quale la *data quality* è divenuta una materia estremamente interessante negli ultimi anni. Per ridurre anomalie ed inconsistenze ci vengono incontro le *Dipendenze funzionali*, utilizzate ampiamente per stabilire vincoli di integrità tra i dati. La grande mole di dati, però, ha reso necessario un riadattamento delle dipendenze funzionali rendendole in grado di catturare inconsistenze più ampie nei dati. Le *Dipendenze funzionali rilassate o approssimate (RFD)* sono da considerarsi come una naturale evoluzione o generalizzazione delle *dipendenze funzionali canoniche*. Il concetto più importante introdotto dalle RFD è quello della *similarità*. Nelle dipendenze funzionali classiche esisteva soltanto il concetto di uguaglianza tra dati, nelle RFD espandiamo questo concetto ad una similarità, questo ci permetterà di coprire una quantità di dati maggiore. Tuttavia le RFD possono fornire vantaggi solo se possono essere scoperte automati-

camente. Il lavoro di tesi si è basato su questo ultimo concetto di ottenere le RFD in seguito ad una procedura automatizzata. Durante le varie fasi di studio si è pensato ed implementato un algoritmo che permette, attraverso tre fasi intermedie, la scoperta di RFD di un dataset dato come input. Per questo lavoro di tesi mostreremo l'idea dell'algoritmo generale ed entreremo nel dettaglio della seconda fase di sviluppo (Minimality), mostrando, infine, i risultati della sperimentazione.

Indice

1	Introduzione	4
1.1	Incipit	4
1.2	Nozioni Preliminari	6
1.2.1	Schema di relazione	6
1.2.2	Dipendenze funzionali canoniche	7
1.2.3	Dipendenze funzionali rilassate	8
1.2.4	Scoperta RFD	12
1.2.5	Dominanza	13
1.3	Studi preliminari	13
2	Stato dell'arte	15

Elenco delle tabelle

1.1	Esempio di schema di relazione	6
1.2	Esempio di Relazione con anomalie	11

Snippet di codice

Introduzione

1.1 Incipit

Nella progettazione di una base di dati ci sono aspetti essenziali da prendere in considerazione per assicurare un servizio quanto più efficiente possibile. Uno di questi servizi è certamente la *qualità dei dati*, una base di dati con questa caratteristica farà sì che le inconsistenze tra i dati siano il minor numero possibile. Negli ultimi anni la crescita delle reti ha portato ad un aumento considerevole del flusso di dati rendendo la *data quality* una materia estremamente interessante vista la cospicua presenza di dati "sporchi" proveniente da fonti differenti. Per ridurre questo tipo di anomalie è impensabile tentare di eliminare le *inconsistenze* manualmente, una procedura di questo tipo può essere facilmente incline ad errori soprattutto con la quantità di dati precedentemente citata. In questo lavoro ci vengono incontro le *Dipendenze funzionali*, utilizzate ampiamente per stabilire vincoli di integrità tra i dati e ridurre anomalie e inconsistenze all'interno della nostra base di dati. La grande mole di

dati, però, ha reso necessario un riadattamento delle dipendenze funzionali rendendole in grado di catturare inconsistenze più ampie nei dati. Le *Dipendenze funzionali rilassate o approssimate* (**RFD**) sono da considerarsi come una naturale evoluzione o generalizzazione delle *dipendenze funzionali canoniche*. Questo nuovo strumento ci permette di adattare le semplici dipendenze funzionali a diversi contesti applicativi, infatti, le RFD possono applicarsi anche solo ad una porzione di database. Il concetto più importante introdotto dalle RFD, però, è quello della *similarità*. Nelle dipendenze funzionali classiche esisteva soltanto il concetto di uguaglianza tra dati, nelle RFD espandiamo questo concetto ad una similarità, questo ci permetterà di coprire una quantità di dati maggiore e sfruttare le RFD appena scoperte per effettuare una operazione di *cleaning* sulla base di dati. Tuttavia le RFD possono fornire vantaggi solo se possono essere scoperte automaticamente. Il lavoro di tesi si è basato su questo ultimo concetto di ottenere le RFD in seguito ad una procedura automatizzata. Durante le varie fasi di studio si è pensato ed implementato un algoritmo che permette, attraverso tre fasi intermedie, la scoperta di RFD di un dataset dato come input. Le tre fasi di questo algoritmo sono: *Feasibility, Minimality, RFD Discovery* . Per questo lavoro di tesi mostreremo l'idea dell'algoritmo generale ed entreremo nel dettaglio della prima fase di sviluppo(Feasibility), mostrando, infine, i risultati della sperimentazione. Per questo algoritmo, particolare attenzione è stata posta sull'efficienza, oltre che sull'efficacia, studiando un'implementazione basata sul multithreading e predisponendola ad eventuale adattamento parallelo.

1.2 Nozioni Preliminari

E' necessario, prima di cominciare con lo studio del nostro algoritmo, introdurre alcuni concetti preliminari volti alla comprensione della logica dietro le RFD.

1.2.1 Schema di relazione

Uno schema di relazione è costituito da un simbolo R , detto nome della relazione, e da un insieme di attributi $X = \{A_1, A_2, \dots, A_n\}$, di solito indicato con $R(X)$. A ciascun attributo $A \in X$ è associato un dominio $dom(A)$. Uno schema di base di dati è un insieme di schemi di relazione con nomi diversi:

$$R = \{R_1(X_1), R_2(X_2), \dots, R_n(X_n)\}.$$

Una relazione su uno schema $R(X)$ è un insieme r di tuple su X . Per ogni istanza $r \in R(X)$, per ogni tupla $t \in r$ e per ogni attributo $A \in X$, $t[A]$ rappresenta la proiezione di A su t . In modo analogo, dato un insieme di attributi $Y \subseteq X$, $t[Y]$ rappresenta la proiezione di Y su t . [1]

Matricola	Cognome	Nome	Data di nascita
123456	Rossi	Mario	25/11/1991
567891	Neri	Anna	23/04/1992

Tabella 1.1: Esempio di schema di relazione

1.2.2 Dipendenze funzionali canoniche

Una *dipendenza funzionale*, abbreviata in FD, è un vincolo di integrità semantico per il modello relazionale che descrive i legami di tipo funzionale tra gli attributi di una relazione.

Data una relazione r su uno schema $R(X)$ e due sottoinsiemi di attributi non vuoti Y e Z di X , diremo che esiste su r una dipendenza funzionale tra Y e Z , se, per ogni coppia di tuple t_1 e t_2 di r aventi gli stessi valori sugli attributi Y , risulta che t_1 e t_2 hanno gli stessi valori sugli attributi Z :

$$\forall t_1, t_2 \in r, t_1[Y] = t_2[Y] \implies t_1[Z] = t_2[Z] \quad (1.1)$$

Una dipendenza funzionale tra gli attributi Y e Z viene indicata con la notazione $Y \rightarrow Z$ e viene associata ad uno schema.

Se l'insieme Z è composto da attributi A_1, A_2, \dots, A_k , allora una relazione soddisfa $Y \rightarrow Z$ se e solo se essa soddisfa tutte le k dipendenze $Y \rightarrow A_1, Y \rightarrow A_2, \dots, Y \rightarrow A_k$. Di conseguenza, quando opportuno, possiamo assumere che le dipendenze abbiano la forma $Y \rightarrow A$, con A singolo attributo.

Una relazione funzionale è *non banale* se A non compare tra gli attributi di Y .

Data una chiave K di una relazione r , si può facilmente notare che esiste una dipendenza funzionale tra K ed ogni altro attributo dello schema di r . Quindi una dipendenza funzionale $Y \rightarrow Z$ su uno schema $R(X)$ degenera nel vincolo di chiave se l'unione di Y e Z è pari a X . In tal caso Y è superchiave per lo schema $R(X)$.

Con la notazione $\langle R(X), F \rangle$ indicheremo uno schema $R(X)$ su cui è definito un

insieme di dipendenze funzionali F . Un'istanza r di $R(X)$ viene detta *istanza legale* di $\langle R(X), F \rangle$ se soddisfa tutte le dipendenze funzionali in F . Infine, data una relazione funzionale $Y \rightarrow Z$, se ogni istanza legale r di $\langle R(X), F \rangle$ soddisfa anche $Y \rightarrow Z$, allora diremo che F *implica logicamente* $Y \rightarrow Z$, indicato come $F \models Y \rightarrow Z$.

1.2.3 Dipendenze funzionali rilassate

In alcuni casi per risolvere dei problemi in alcuni di domini di applicazioni, come l'identificazione di inconsistenze tra i dati, o la rilevazione di relazioni semantiche fra i dati, è necessario rilassare la definizione di dipendenza funzionale, introducendo delle approssimazioni nel confronto dei dati. Invece di effettuare dei controlli di uguaglianza, si utilizzano dei controlli di similarità. Inoltre spesso si potrebbe desiderare che una certa dipendenza valga solo su un sottoinsieme di tuple che su tutte. Per questo motivo sono nate delle dipendenze funzionali che rilassano alcuni dei vincoli delle FD, prendono il nome di Dipendenze Funzionali Rilassate o Approssimate ¹. Esistono differenti tipi di RFD, ciascuna di esse rilassa uno o più vincoli delle FD, si possono dividere in due macro aree:

1. Confronto di attributi: La funzione di uguaglianza delle FD canoniche viene sostituita da una funzione di similarità , ciò implica che l'AFD deve descrivere una soglia di rilassamento per ogni attributo.

¹RFD abbreviazione di Relaxed Functional Dependency.

2. Estensione: Permette che il vincolo non sia valido su tutte le tuple, ma solo su di un sottoinsieme di esse.

Le RFD sono utilizzate in attività di: data cleaning, record matching e di rilassamento delle query.

Le definizione formale di una RFD è la seguente:

Teorema 1. *Sia R uno schema relazionale definito su di un insieme di attributi finito, e sia $R = (A_1, A_2, \dots, A_k)$ una relazione definita su R . Una RFD φ su R viene rappresentata come:*

$$D_c : (X)_{\Phi_1} \xrightarrow{\Psi(X,Y) \leq \epsilon} (Y)_{\Phi_2}$$

dove

- $\mathbb{D}_c = \{t \in \text{dom}(R) \mid (\bigwedge_{i=1}^k c_i(t[A_i]))\}$, dove $c = (c_1, c_2, \dots, c_k)$ con c_i predicato sul $\text{dom}(A_i)$, utilizzato per filtrare le tuple a cui φ va applicata;
- $X, Y \subseteq \text{attr}(R)$ tali che $X \cap Y = \emptyset$;
- $\Phi_1 (\Phi_2 \text{ rispettivamente})$ è un insieme di vincoli $\phi[X](\phi[Y])$ definite sugli attributi di X (Y). Per qualsiasi coppia di tuple $(t_1, t_2 \in \mathbb{D}_c)$ il vincolo $\phi[X](\phi[Y])$ restituisce vero se la similarità fra $t_1[X]$ e $t_2[X]$ ($t_1[Y]$ e $t_2[Y]$) concordano con i vincoli specificati da $\phi[X](\phi[Y])$;
- $\Psi : \text{dom}(X) \times \text{dom}(Y) \rightarrow \mathbb{R}$ rappresenta una misura di copertura su \mathbb{D}_c e indica il numero di tuple che violano o soddisfano φ ;
- ϵ è la soglia che indica il limite superiore o inferiore per il risultato della misura di copertura;

Per questo algoritmo vengono trattate solo le RFD che rilassano il vincolo di uguaglianza, senza considerare l'estensione. Data RFD $X \rightarrow Y$ essa vale su una relazione r se e solo se la distanza fra due tuple t_1 e t_2 , i cui valori sui singoli attributi A_i non superano una certa soglia β_i , è inferiore ad una certa soglia a_A su ogni attributo $A \in X$, allora la distanza fra t_1 e t_2 su ogni attributo $B \in Y$ è minore di una certa soglia a_B .

La struttura delle RFD utilizzate è la seguente:

$$attr_1(\leq soglia_1), \dots, attr_n(\leq soglia_n) \rightarrow RHS$$

Gli attributi che si trovano a sinistra della freccia costituiscono la parte LHS², l'attributo che invece si trova dopo la freccia costituisce l'RHS³. È importante focalizzare l'attenzione su questo concetto in quanto le dipendenze funzionali hanno un verso, ed è quello indicato dalla freccia. Qualsiasi operazione effettuata con le RFD deve sempre tener conto del verso, le RFD non forniscono conoscenza nel verso opposto. Questa non è una proprietà riguardante solo le RFD, bensì riguarda qualsiasi tipo di dipendenza funzionale. Ad esempio consideriamo la relazione in questa tabella:

²Left Hand Side o lato sinistro.

³Right Hand Side o lato destro.

Impiegato	Stipendio	Progetto	Bilancio	Funzione
Rossi	20000	Sito web	2000	tecnico
Verdi	35000	App Mobile	15000	progettista
Verdi	35000	Server	15000	progettista
Neri	55000	Server	15000	direttore
Neri	55000	App Mobile	15000	consulente
Neri	55000	Sito web	2000	consulente
Mori	48000	Sito web	15000	direttore
Mori	48000	Server	15000	progettista
Bianchi	48000	Server	15000	progettista
Bianchi	48000	App Mobile	15000	direttore

Tabella 1.2: Esempio di Relazione con anomalie

Si può osservare che lo stipendio di ciascun impiegato è unico, quindi in ogni tupla in cui compare lo stesso impiegato verrà riportato lo stesso stipendio. Possiamo dire che esiste una Dipendenza Funzionale: $Impiegato \rightarrow Stipendio$. Si può fare lo stesso discorso tra gli attributi Progetto e Bilancio, quindi anche qui abbiamo una dipendenza funzionale $Progetto \rightarrow Bilancio$. Non si può dire che di conseguenza vale anche il verso opposto:

$$Impiegato \rightarrow Stipendio \neq Stipendio \rightarrow Impiegato$$

Infatti percepiscono 48000 di stipendio sia Mori che Bianchi.[1]

1.2.4 Scoperta RFD

Data una relazione r , la scoperta di una RFD è il problema di trovare un *minimal cover set* di RFD che si verificano per r . Questo problema rende ancor più complesso il problema della scoperta delle dipendenze dei dati visto l'ampio spazio di ricerca dei possibili vincoli di similarità. Dunque è necessario trovare algoritmi efficienti in grado di estrarre RFD con vincoli di similarità significativi.

Se i vincoli di similarità e le soglie sono noti per ogni attributo del dataset, scoprire le RFD si riduce a trovare tutte le possibili dipendenze che soddisfano la seguente regola:

Lemma 1. *Le partizioni di tuple che sono simili sugli attributi contenuti nel lato sinistro o LHS della dipendenza, devono corrispondere a quelle che sono simili nel lato destro o RHS.*

Questo problema è simile a trovare le FD, dove bisogna trovare le partizioni di tuple che condividono lo stesso valore sull'RHS quando esse condividono lo stesso valore sull'LHS. Il problema viene reso più semplice dal fatto che, nel caso della scoperta delle FD, tali partizioni sono disgiunte, cosa che però non vale nelle RFD in quanto uno stesso valore può essere simile a valori differenti. Ciò impedisce quindi di sfruttare gli algoritmi utilizzati nella scoperta delle FD, nella scoperta delle RFD.

1.2.5 Dominanza

Nel corso del nostro lavoro, abbiamo applicato alcuni risultati dell'intelligenza artificiale al campo della discovery delle *Dipendenze Funzionali Rilassate*. In particolare, cercando di individuare le dipendenze funzionali rilassate ci siamo serviti di un importante risultato nella succitata materia: la dominanza stretta (*strict dominance*).

Teorema 2. *Dato un vettore di attributi $\mathbf{X} = X_1, X_2, \dots, X_n$, siano $\mathbf{x} = (x_1, x_2, \dots, x_n)$ e $\mathbf{y} = (y_1, y_2, \dots, y_n)$ due vettori di assegnamenti definiti sugli attributi di \mathbf{X} , dove l' i -esimo elemento x_i o y_i può essere sia un valore numerico sia un valore discreto con un assunto ordinamento su tali valori. Diremo che \mathbf{x} domina strettamente (o deterministicamente) \mathbf{y} se e solo se*

$$y_i \leq x_i \quad i = 1, 2, \dots, n,$$

ovvero

$$\mathbf{y} - \mathbf{x} \leq \mathbf{0}$$

1.3 Studi preliminari

Prima di iniziare lo sviluppo dell'algoritmo per la scoperta di RFD si è reso necessario uno studio approfondito di un algoritmo precedentemente sviluppato per un progetto di IA [2]. Tale algoritmo è stato sviluppato in Python, pertanto, abbiamo effettuato uno studio del linguaggio precedentemente citato. Oltre le principali caratteristiche di questo linguaggio, è stato fatto uno studio anche delle librerie utilizzate all'interno del progetto:

- ***Pandas***: È una libreria che include delle strutture dati e tool di analisi facili da usare e fortemente ottimizzate.
- ***Numpy***: È un package dedicato all'elaborazione scientifica sul linguaggio Python.

Una volta concluso questo tipo di studio si è cominciato a pensare allo sviluppo dell'algoritmo in un ambiente differente. La scelta è ricaduta su *Java*, tale scelta è dovuta, oltre che alla già piena conoscenza del team di questo linguaggio, alla potenza e versatilità che questo linguaggio ci offre, oltre che al gran numero di framework presenti per la gestione di parallelizzazione e concorrenza, essendo quest'ultimo un aspetto molto importante per l'efficienza dell'algoritmo. Le librerie esterne studiate ed utilizzate per la realizzazione di questo algoritmo saranno ben approfondite nel capitolo 4 di questo elaborato. Le sopracitate librerie esterne utilizzate sono:

- ***AKKA***: Framework per la gestione del parallelismo e concorrenza.
- ***Joinery Dataframe***: Struttura dati simile al dataframe presente in *Pandas* di Python.

All'infuori delle conoscenze legate ai linguaggi di programmazione, è stato necessario leggere e studiare vari documenti legati al mondo delle dipendenze funzionali.

Stato dell'arte

Questo è lo stato dell'arte

Bibliografia

- [1] F. P. P. S. T. R. Atzeni P., Ceri S., *Basi di dati: Modelli e linguaggi di programmazione*. McGraw Hill, 2013.
- [2] T. M. Altamura A., Di Pasquale D., “Relaxed functional dependencies discovery,” *Unisa*, 2017.