

Realizzazione di un Add-on per Google Docs per estrazione interattiva di pattern sintattici



Relatore

Alberto Bartoli

Correlatore

Eric Medvet

Candidato

Lorenzo Gasparini



Descrizione del problema

- Dato un documento di testo si vogliono *individuare* ed *estrarre* tutte le occorrenze di un determinato *pattern sintattico*
- **Esempi:** Indirizzi IP (XXX.XXX.XXX.XXX), Date (DD/MM/YYYY), E-mail (alias@domain.ext)

Motivazione di base

```
/(?:[0-9]{1,3}\.){3}[0-9]{1,3}/g
```

Figura: RegEx (semplificata) per indirizzo IP

- Gli strumenti per specificare ed estrarre pattern da documenti testuali esistono già (e.g. RegEx), perchè un nuovo approccio?
- La curva di apprendimento di tali strumenti è ripida, non sono alla portata degli utenti comuni
- **Idea:** l'utente fornisce degli **esempi** di entità da estrarre e da non estrarre, l'algoritmo **deduce** il pattern e lo estrae dal testo



Obiettivo della tesi

- Implementare un algoritmo di estrazione delle entità basato sulla sintassi, sotto forma di Add-on per Google Docs
- L'algoritmo è stato sviluppato nel laboratorio di Machine Learning, ed è risultato in media il migliore in un confronto con altri algoritmi su 10 dataset
- L'algoritmo si basa sull'**active learning**: genera un estrattore sulla base degli esempi forniti dall'utente; sceglie l'esempio che deve aggiungere l'utente per migliorare l'estrattore

Google Docs



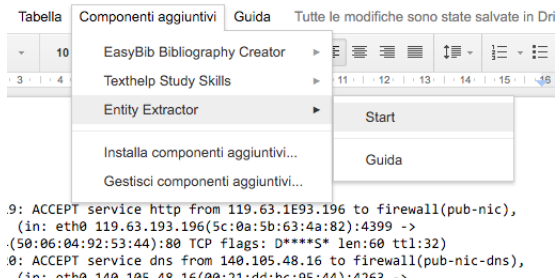
- Nasce nel 2006, è una piattaforma online di elaborazione testi
- Consente agli utenti di creare e modificare documenti direttamente nel browser e di collaborare con altri utenti in tempo reale



Add-on per Google Docs

- Dal 2014 è possibile sviluppare componenti aggiuntivi che permettono di ampliare le funzionalità della piattaforma
- Un Add-on è composto da un insieme di file HTML/Javascript/CSS e di script *Google Apps Script*, il quale:
 - È basato su Javascript
 - Viene eseguito dai server di Google

Interfaccia Add-on



- Per aprire l'Add-on, si apre un documento di testo Google Docs e si preme *Start* nella relativa voce del menù

Interfaccia Add-on

Modifica Visualizza Inserisci Formato Strumenti Tabella Componenti aggiuntivi Guida Tutte le modifiche s... Visualizza Consola

100% Testo norm... Consolas 10 Altre

Entity Extractor

☒ Desired extractions: Add

119.63.193.196 ×

☒ Desired unextractions: Add

Jan 12 06:26:19: ACCEPT service http from m ×

Extract Reset EXPORT

```
Jan 12 06:26:19: ACCEPT service http from 119.63.193.196 to firewall(pub-nic),
prefix: "none" (in: eth0 119.63.193.196(5c:0a:5b:63:4a:82):4399 ->
140.105.63.164(50:06:04:92:33:41):80 TCP flags: D****S* len:60 ttl:32)
Jan 12 06:26:20: ACCEPT service dns from 140.105.48.16 to firewall(pub-nic-dns),
prefix: "none" (in: eth0 140.105.48.16(00:21:00:00:00:00):4263 ->
140.105.63.158(00:14:31:83:c6:8d):53 UDP len:76 ttl:62)
Jan 12 06:27:09: DROP service 68->67(udp) from 216.34.211.83 to 217.70.158.154,
prefix: "spoof iana-0/8" (in: eth0 213.92.153.78(00:1f:d6:19:0a:80):68 ->
69.43.177.110(00:30:fe:fd:d6:51):67 UDP len:576 ttl:64)
Jan 12 06:27:13: DROP service 68->67(udp) from 213.92.39.37 to 216.34.41.186,
prefix: "spoof iana-0/8" (in: eth0 216.34.190.233(00:00:5a:49:61:ab):68 ->
69.43.93.45(00:26:32:9d:8d:35):67 UDP len:576 ttl:64)
Jan 12 06:27:16: DROP service 68->67(udp) from 69.43.127.88 to 69.43.242.47,
prefix: "spoof iana-0/8" (in: eth0 172.45.240.101(00:1d:d8:00:1c:c0):68 ->
216.34.219.223(00:16:a5:65:db:c0):67 UDP len:576 ttl:64)
Jan 12 06:28:18: DROP service 68->67(udp) from 213.92.89.154 to 217.70.194.196,
prefix: "spoof iana-0/8" (in: eth0 217.70.158.16(30:5d:38:86:35:7a):68 ->
217.70.75.13(00:0b:7a:c3:6a:61):67 UDP len:576 ttl:64)
```

- L'aggiunta di nuovi esempi avviene mediante la selezione del testo ed il click sul relativo pulsante *Add*
- Gli esempi vengono evidenziati con colori diversi per un'indicazione visuale istantanea

Interfaccia Add-on

Formato: Consolas 10 Altro

```

Service http from 119.63.193.196 to firewall(pub-nic),
  119.63.193.196(5c:8a:5b:63:4a:82):4399 ->
  [92:53:44]:80 TCP flags: D****S len:68 ttl:32
Service dns from 140.105.48.16 to firewall(pub-nic-dns),
  140.105.48.16(00:21:dd:bc:95:44):4263 ->
  [83:c6:8d]:53 UDP len:78 ttl:62
Irrive 68->67(udp) from 216.34.211.83 to 216.34.253.94
  (in: eth0 213.92.153.78(00:1f:d0:19:8a:00):68 ->
  [d:d6:51]:67 UDP len:576 ttl:64)
Irrive 68->67(udp) from 213.92.39.37 to 216.34.41.186
  (in: eth0 216.34.199.233(00:80:5a:49:c1:ab):68 ->
  [8d:35]:67 UDP len:576 ttl:64)
Irrive 68->67(udp) from 69.43.127.88 to 69.43.242.47
  (in: eth0 172.45.240.101(00:1d:d0:00:1c:c8):68 ->
  [65:db:c0]:67 UDP len:576 ttl:64)
Irrive 68->67(udp) from 213.92.89.154 to 217.70.194.196
  (in: eth0 217.70.158.161(30:5d:38:86:35:7a):68 ->
  [16:6f:15]:67 UDP len:576 ttl:64)
Irrive 68->67(udp) from 69.43.36.247 to 172.45.119.91
  (in: eth0 213.92.156.146(00:01:f3:4d:0a:c1):68 ->
  [8:10:98):67 UDP len:576 ttl:64)
Irrive 68->67(udp) from 213.92.169.120 to 69.43.115.84
  (in: eth0 172.45.227.107(00:06:aa:c9:c4:b2):68 ->
  [18:51:04):67 UDP len:576 ttl:64)
Service dns from 66.249.75.180 to firewall(pub-nic-dns),
  66.249.75.180(00:1f:f3:bd:ba:aa):52378 ->
  [58:00:93]:53 UDP len:65 ttl:45)
Irrive 68->67(udp) from 172.45.186.97 to 217.70.191.141
  (in: eth0 213.92.59.210(fc:fa:f7:5c:c8:90):68 ->
  [ ]
  
```

Entity Extractor

☒ Desired extractions: Add

119.63.193.196 ✕

☒ Desired unextractions: Add

Jan 12 06:26:19: ACCEPT service http from ✕

Extract Reset EXPORT

☒ System extractions:

213.92.89.154

140.105.48.16(00:21:dd:bc:95:44):4263

to

Query

→

Extract Don't Extract Edit

- *Extract* avvia la costruzione dell'estrattore in base agli esempi forniti
- Il comportamento dell'estrattore è mostrato dal sistema
- Viene formulata una query allo scopo di ottenere un nuovo esempio (*active learning*)

Esportazione estrazioni

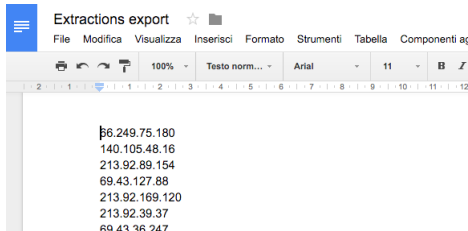
```
n: eth0 216.34.190.233(00:80:5a:49:61:ab):68 ->  
35).:67 UDP len:576 ttl:64)
```

System extractions export

The system extractions were exported. [Link](#)

```
ce 68->67(udp) from 213.92.169.120 to 69.43.115.84,  
n: eth0 172.45.227.107(00:06:ea:c9:c4:b2):68 ->
```

- Cliccando su Export è possibile esportare le attuali estrazioni suggerite in un nuovo documento
- Si aprirà una finestra modale con un link al nuovo documento contenente le estrazioni



Dettagli implementativi

Nome file	Righe	Contenuto
EntityExtractor.gs	561	Motore dell'algoritmo
Main.gs	507	Gestione interazione utente
Set.gs	54	Struttura dati insieme
Sidebar.css.html	67	CSS Sidebar
Sidebar.html	61	HTML Sidebar
Sidebar.js.html	322	Javascript Sidebar
Store.gs	53	Gestore memorizzazione dati server
TextRange.gs	34	Struttura dati annotazione



Limiti della piattaforma e sviluppi futuri

- Google Docs è una piattaforma proprietaria che presenta dei limiti intrinseci:
 - Le evidenziazioni sono permanenti, modificano la struttura del documento
 - Non è possibile gestire l'evento di chiusura dell'Add-on, impedendo l'esecuzione di azioni di pulizia del documento
 - Lo spazio di archiviazione lato server è ristretto, ciò rende difficile implementare meccanismi di cache atti a diminuire il carico computazionale dell'algoritmo
- Soluzione:
 - Migrazione a piattaforma web standalone o GUI desktop



Demo

Dimostrazione



Fine

Grazie per l'attenzione.