

Abschluss Arbeit

Neo4j und..

Asmaa Haja
348758

Ilhem Bouzir
350778

Jossep
330187

6. Juli 2015

Muss geändert werden: TU Berlin – Fakultät IV
Institut für Wirtschaftsinformatik und Quantitative Methoden
Fachgebiet Informatik und Gesellschaft
Prof. Dr.-Ing. Frank Pallas
Max Ulbricht

Inhaltsverzeichnis

1	Abstract	1
2	Einleitung	1
3	Aufgaben	1
4	Umsetzung	2
4.1	Knolwledge Graphen	2
4.2	Neo4j und Cypher	2
4.3	Sparql	3
4.4	Dbpedia	4
4.5	OWL, RDF, RDFS	4
4.6	Anfragen an die strukturierten Internetwebseiten über die Sparql-Interface	6
4.7	Automatizierung der Sparqlanfragen über Java	6
4.8	Benutzerschnittstelle(GUI)	6
5	Anleitung	9
6	Evaluation	10
7	Zusammenfassung	10
	Literatur	11

Abbildungsverzeichnis

1	Query Beispeil [3]	3
2	Ontologie	5
3	RDF	5
4	KlassenDiagram [3]	7
5	Hauptansicht	7
6	Datenbank wurde erzeugt	8
7	Auswahl von start und end City gleich gewählt	8
8	Ergebnis Bei der berechnung zwei verschiedene Städte	8
9	Hauptansicht	9
10	Hauptansicht	9
11	Hauptansicht	9

Listings

1 Abstract

Dieser Bericht beschreibt die Erstellung eines knowledgegraphs mit Hilfe von Neo4j. Dieser Graph repräsentiert verschiedene Deutsche Städte, deren Sehenswürdigkeiten und deren Entfernungen voneinander. Das Projekt fand im Rahmen der Veranstaltung DBPRO an der TU Berlin statt.

2 Einleitung

Im Bereich Big Data fallen zum Beispiel bei Texten sehr große Datenmengen an. Sinnvollerweise sollten solche Textinformationen auf einer übersichtlicheren und anschaulicherer Weise dargestellt werden können, so dass anhand einer bestimmten Repräsentation der Daten aus diesem Text gewonnen werden können. Diese erlaubt eine leichtere Suche nach bestimmten Informationen, ohne dass dabei wichtige Daten verloren gehen. Ziel des Projektes war es eine Knowledge Graph mit Hilfe der Graphendatenbank Neo4j zu erstellen basierend auf bestimmten rausgesuchten Daten. Mit der Graphenbasierte Abfragesprache sollen dann benötigte Informationen aus einer Wikipedia Informationen basierte Linked Open Data Version (DBpedia) extrahiert werden können und in einem bestimmten Format zu Neo4j übergeben werden, so dass diese die Daten in strukturierter Graphenform automatisch bringt. Das Projekt fand im Sommersemester 2015 an der TU Berlin im Rahmen der Veranstaltung DBPRO statt. Das Projektteam bestand aus Asmaa Haja, Ilhem Bouzir und Josseph, welche durch Dr. Holmer Hensen betreut wurden. Das Projekt umfasste drei Meilensteine. Zunächst fand im 1. Meilenstein eine Einarbeitung in das Thema Knowledge Graph, Neo4j und SPARQL statt. Als nächstes wurde dann konkrete Beispiele für Tourism in Deutschland mit SPARQL erstellt. Im 3. Meilenstein wurde aus der zuvor erstellten Abfragen mit Hilfe von Neo4j ein Knowledge Graph erstellt. Diese Arbeiten werden in diesem Bericht nach einer genauen Identifizierung der Aufgaben näher erläutert. Abschließend folgen ein Kapitel zur Evaluation und eine kurze Zusammenfassung der Projektergebnisse.

3 Aufgaben

Als nächstes gliedern wir die Erarbeitung des Projektes in drei Meilensteine, welche die zu erreichenden Ziele im Laufe des Projektes beinhalten:

- Meilstein 1
 - Eingewöhnung Knowledge Graphen
 - Eingewöhnung Neo4j und Cyper-Sprache
 - Eingewöhnung Sparql
- Meilstein 2
 - Eingewöhnung Dbpedia

- Eingewöhnung OWL, RFD, RFDS
- Anfragen an die strukturierten Internetwebseiten über die Sparql-Interface
- Meilstein 3
 - Automatisierung der Sparqlanfragen über Java.
 - Automazierte Erstellung des Datenbakes über Java.
 - Erstellung einer Schnittstelle(GUI) für die Benutzer um Anfragen an die Datenbank zu ermöglichen

4 Umsetzung

Die bearbeitung der einzelnen Ziele werden demnächst näher erläutert.

4.1 Knowledge Graphen

knowledge Graph wurde im 1982 von Hoede and Stokman entwickelt. Sie hatten den Ziel Wissen aus medizinischen und soziologischen Texte zu extrahieren um Expertensysteme zu erhalten und dadurch neue Methoden um das Wissen zu repräsentieren. Knowledge Graphen gehört zu der Kategorie semantischen Netze. [6]

Wissen kommt im Netzwek im strukturierter als auch im unstrukturierter Form vor. Außerdem kann das in deutlicher als auch undeutlicher Form gestaltet sein. Das führt dazu, dass der nach Information suchenden Netzwerknutzer zu sowohl sinnvollem als auch unbrauchbarem Wissen kommt. Um der Neztwerknutzer einen sinnvollen Überblick über das Wissen im Nezt zu erschaffen, Wird das in sctrukturierter Form aufgebaut. Auf diese Weise kommen wir zur Knowledge Graphen. Was nicht anderes ist als eine Art von semantischen Netzwerke, wo aus der Analysephase von texten zu einer Liste von Konzepte führt. Aus diese Listen werden bennaten Knoten erzeugt und letztendlich verbindet man sie anhand von Relationen, welche die Beziehungen zwieschen Konzepte(Knoten) darstellt und zu der Strukturierung des Wissen hinführt.[4]

4.2 Neo4j und Cypher

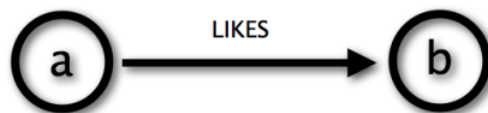
Neo4j ist ein open-source NoSQL graph Database, welche im Java und Scala implementiert wurde und von Neo Technology unterstützt wird. Die entwicklug began im 2003 und der freie Zugang ist seit 2007 veröffentlicht worden. Diese Technologien findet Anwendungen im folgenden Bereiche: Matchmaking, Netzwerk-Management, Software-Analytik, wissenschaftliche Forschung , soziale Netzwerke und mehr. Um die Arbeitsweise von Neo4j zu verstehen brauchen wir die Begriffe "Property Graph Database und "NoSQL"kennen lernen. Ein Property Graph ist, wie der Name einwenig verratet, ein aus Knoten und kanten Bestehendes Graph.welches sich von einem normalen Graph dadurch unterscheidet, dass sowohl Knoten als auch Kanten Eigenschaften besitzen können. Und der Begriff "NoSQList die Abkürzung für No Only SQL,diese Spracheanfrage umgeht

die Leistungsproblemen bei der Anfrage an relationalen Datenbanken in dem durch vernachlässigen von Consistency eine Verbesserung von Availability und Partitioning statt findet. Das führt dazu, dass die NoSQL Database sich besser skalieren lassen.(Die Abbildung 1 zeigt wie Daten mit Neo4j dargestellt werden

Cypher ist eine deklarative Sprache, welche von Neo4j benutzt wird, Um Anfragen an die Datenbank zu ermöglichen. Ein wichtiger Vorteil ist, dass beim Stellen einer Frage an die Datenbank nur angegeben wird, wonach gesucht wird und nicht wie man es finden soll.[3]

Die Abbildung 1 zeigt ein kleines Beispiel,wie Daten im Neo4j dargestellt werden und wie man mit der untenliegenden Cypher-Anfrage auf die Knoten und die dazwischen liegende Relation zugreifen kann. Beide Knoten sind von Type City Und die Relation ist von Type Distance

Cypher using relationship 'likes'



Cypher

(a) -[:LIKES]-> (b)

© All Rights Reserved 2013 | Neo Technology, Inc.

Abbildung 1: Query Beispiel [3]

4.3 Sparql

Sparql ist eine für relationale Datenbanken Anfragesprache, Welche im erste linie entwickelt worden ist um RDF Graphen anzufagen. Der Baustein für SPARQL Abfragen ist Basic Graph Patterns(BGP). SPARQL(BGP) ist eine Menge aus triple patterns, welche im Abschnitt RDF näher erläutert wird. Um komplexe Anfragen aufzubauen, benutzt man die Operatoren Select, Optional, Union, und Filter. [5]

Da RDF-Dokumente eine Menge von RDF-trippeln ist. Verwendet Spaqrl eine ähnliche strukturierung. d.h zum Aufbau eine Anfragemuster verwendet man Tripeln der Form (Subjekt,Predikat,Objekt), Wobei das Beenden des Tripeln duch einen Pukt gekennzeichnet wird. Außerdem kann man für die Teile eines Tripels variablen verwenden,

welsche mit "?" gekennzeichnet werden. In Abbildung 3 ist eine Beispielanfrage zu erkennen.

```
PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX dbpedia-owl: <http://dbpedia.org/ontology/>

SELECT DISTINCT ?city ?cprop ?country ?label
WHERE { ?city rdf:type dbpedia-owl:City ;
        rdfs:label ?label ;
        dbpedia-owl:country ?country ;
        dbpprop:countryCode ?cprop
}
```

Der SELECT Klausel enthält die Variablen, welche nach ausführung der Anfrage belegt worden sind. Im der WHERE Klausel wird der Aufbau des gewünschten Graphs beschrieben. Die komplette Sprachsyntax, sowie alle Sprachmittel, sind in [W3C SPARQL] zu finden.

4.4 Dbpedia

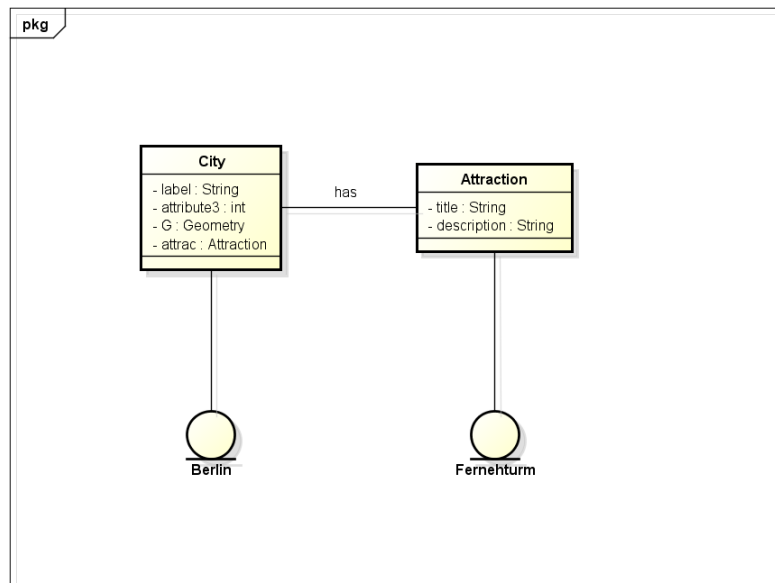
Im rahmen dieses Projekts benutzen wir Dbpedia als prinzipal Quelle um informationen aus Wikipedia zu extrahieren.

DBpedia ist eine Gemeinschaftsbemühung um strukturierter Information aus Wikipedia zu extrahieren. Es handelt sich um eine sammlung von Quellen, welche Wissen in verschiedenen Sprachen zur Verfügung stellt. Das zum extrahierenden Wissen kommt im struktuiertes Form vor, welche durch Ontologien definiert worden sind. Das DBpedia Project wurde im 2006 gestartet und hat mittlerweile große Interesse im Forschungsprojekte und in der Praxis geweckt[1]

4.5 OWL, RDF, RDFS

OWL, RDF und RDFS sind OntologieSprachen, die dazu dienen eine Ontologie zu beschreiben. Zuerst erläutern wir kurz was unter dem Begriff Ontologie zu verstehen ist. Dieser bedeutet eine Wissensvermittlung, die eine Bestimmte Domäne beschreibt. [2] Ontology ist ein Bestandteil von DBpedia. Diese kommt auch Komplexe Anfragen entgegen indem sie trotz der Komplexität genaue und Präzise Antworten liefert. DBpedia Ontology umfasst Klassen, Eigenschaften und Instanzen. Der Aufbau einer Ontology wird es die untere Abbildung zeigen. Diese hat 2 Klassen City und Attraction und die beiden dazugehörigen Instanzen Berlin und Fernsehturm. Die beiden Klassen verbindet außerdem die has Relation.

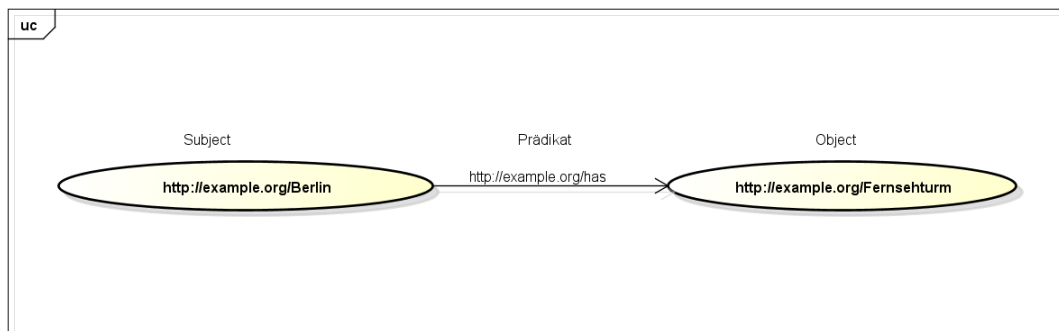
Damit eine Ontology von Maschinen erkennbar wird ist ein Formaler Aufbau benötigt. Deswegen wird Resource Description Framework oder kurz RDF verwendet. Dies ist eine



powered by Astah

Abbildung 2: Ontologie

Triple basierende Sprache und besteht aus Subject Prädikat und Object als Grundlage diese Elemente werden wie es in der Abbildung 3 zu sehen ist verwendet.



powered by Astah

Abbildung 3: RDF

RDF Schema oder kurz RDFS präsentiert eine Erweiterung bei dem Vokabular vom RDF das die verschiedenen Klassenhierarchien ausdrücken lässt. Mit der RDF Verwendung sind die Ressourcen minimal typisiert deswegen fügt RDFS zu den RDF Konstrukten sowie Klassen und Properties weitere Typen hinzu. Außerdem kann anhand des RDFS die Spezialisierung und Generalisierung spezifiziert zwischen den Klassen werden unter der Verwendung von Subklassen was garnicht mit dem RDF Vokabularien darstellbar ist. Letztendlich bietet sich noch eine weitere Sprache Web Ontology Language

ge kurz OWL die aus noch Komplexerer Vokabularien entsteht und damit komplexerer Sachverhalte modellieren erlaubt. Um die eine Ontology zu formalisieren basiert sich die OWL auf RDF und XML Syntax und verwendet auch RDFS Elemente. Zumal hat OWL andere Eigenschaften sowie Mengen Relation Operationen und Kardinalitätenbeschränkungen. Diese extra Funktionalitäten stärken die Ausdrucksmöglichkeiten und die Beschreibungsfähigkeiten von den zu formulierenden Aussagen.

4.6 Anfragen an die strukturierten Internetwebseiten über die Sparql-Interface

Um Daten aus Internetseiten wie DBpedia zu extrahieren benutzen wir die Snorql Interface, welche unter der Link <http://dbpedia.org/snorql> zu erreichen ist. Abbildung 2 zeigt eine Sparql-Anfrage, welche die einige Städte von Deutschland angibt.

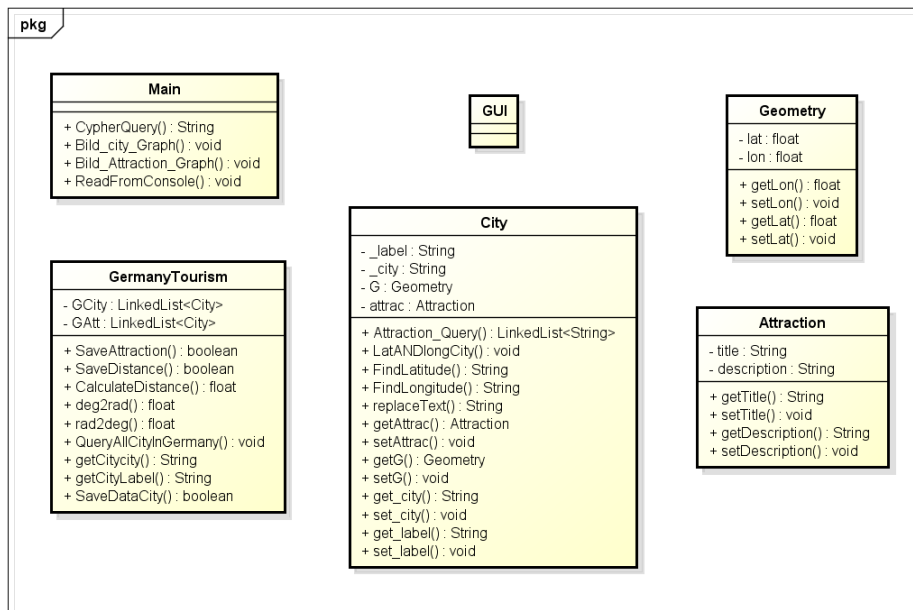
4.7 Automatisierung der Sparqlanfragen über Java

Damit wir unsere benötigte Tourism Informationen extrahieren und danach die automatisch an die Neo4j übergeben können, haben wir einen Java Code mit verschiedenen Funktionen implementiert. Diese sind bei der unteren Abbildung zu erkennen. Sodass am ende ein Knowledgegraph automatisch erstellt werden kann. Zuerst werden die Sparql Anfragen auf DBpedia ausgeführt sodass wir alle unsere Daten zugeliefert bekommen. Anhand der `Attraction_Query()` Methode werden dann die Attraktionen jeder Stadt in einer Liste gespeichert. Bei der Klasse `Germany_Tourism` werden dann diese erhaltene Daten geparkt und in CSV Dateien gespeichert anhand der beiden `saveAttraction()` und `saveDistance()` Methoden wobei die zweite wird erst aufgerufen nachdem die `calculateDistance()` Funktion bereits ausgeführt wurde. Zuallerletzt werden denn diese Datendateien an die Neo4j übermittelt. Dieser liest denn diese Informationen aus und stellt sie in Form eines Graphes. Soweit hat man der Tourism Knowledge Graph, der die verschiedenen Städten, deren zahlreichen Sehenswürdigkeiten und deren Entfernungen voneinander. Auf den Graph können dann beliebige Cypher basierende Anfragen angewendet werden um gesuchte Informationen zu erhalten. Diese kann man dann mit Hilfe der Entwickelten GUI abfragen.

4.8 Benutzerschnittstelle(GUI)

Damit einem Nutzer sich Informationen über die Tourismmöglichkeiten im Deutschland angeben lässt, haben wir einen Schnittstelle(GUI) erzeugt. Welche durch einfacheres klicken sich bedienen lässt. Die Bedienungsschritte werden wir in den nächsten Abchnitte beschreiben.

- Nach starten des Programs bekommt einen Benutzer das im Abbildung 3 gezeichnete GUI zu Gesicht. Durch klicken des Buttons Start wird der Datenbank generiert.



powered by Astah

Abbildung 4: KlassenDiagram [3]

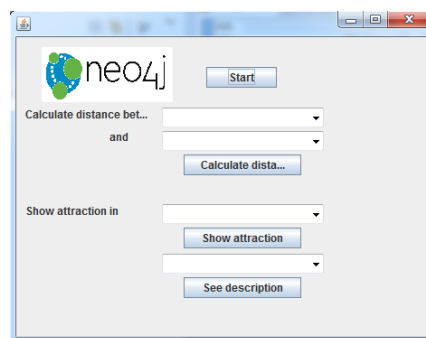


Abbildung 5: Hauptansicht

- Nach der Erstellung des Datenbanks werden Alle deutschen Städte im den Textfelder, die rechts von Calculate distance-,and- und Show attraction in-Labels aufgelistet, sodass der Nutzer per einfacheres klicken zwei städte auswählen kann und sich die Entfernung angeben lässt oder die Sehenswürdigkeiten eines Stadts aufgelistet angeben lässt.
- Falls für die Berechnung der Entfernung zwei Städte der selbe Stadt sowohl als start als auch als end Ziel gewählt wurde, erscheint einen nachrichten Fenster, Wodurch der Kunde angedeutet wird, dass er zwei verschiedene Städte auswählen soll. Im Abbildung 7 wird dies gezeigt.

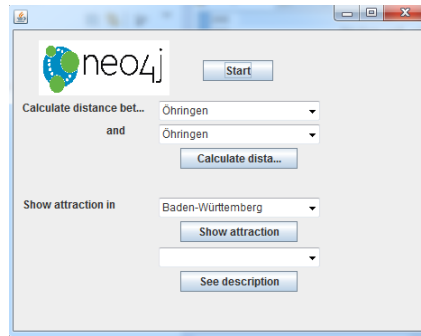


Abbildung 6: Datenbank wurde erzeugt

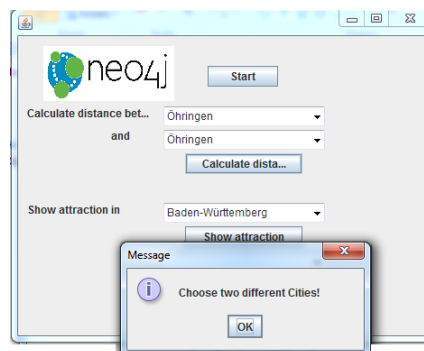


Abbildung 7: Auswahl von start und end City gleich gewählt

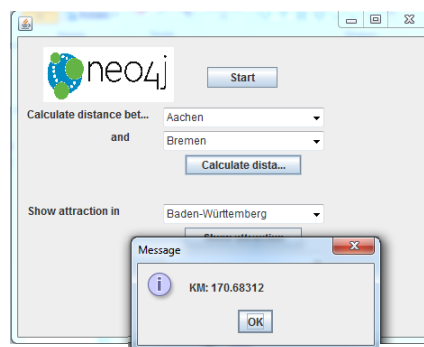


Abbildung 8: Ergebnis Bei der berechnung zwei verschiedene Städte

•

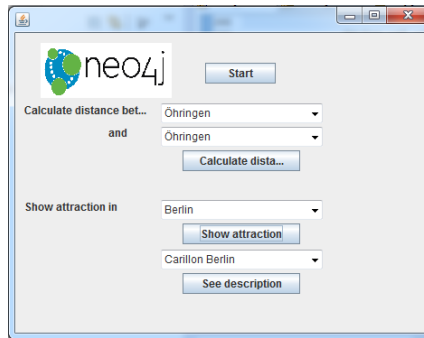


Abbildung 9: Hauptansicht

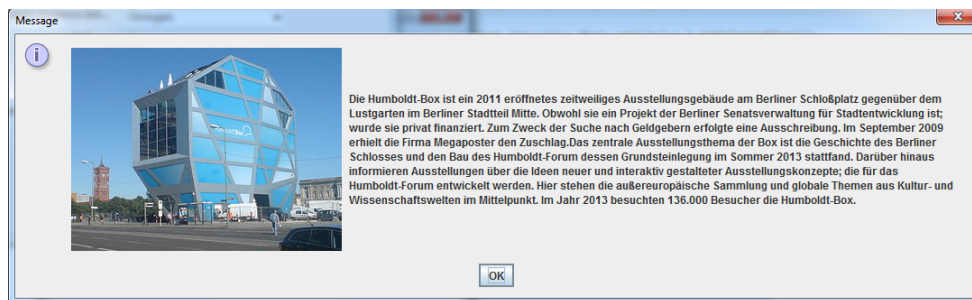


Abbildung 10: Hauptansicht

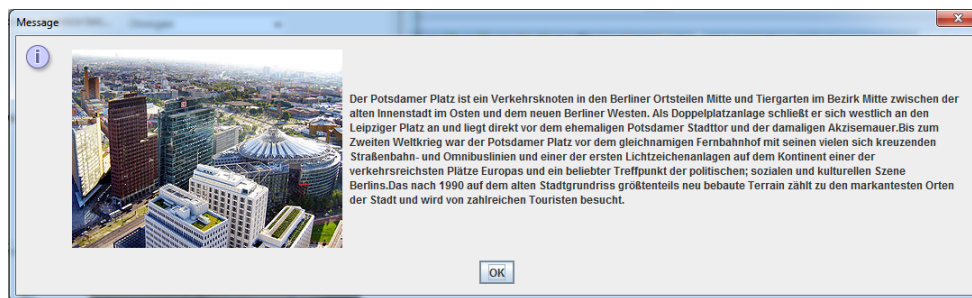


Abbildung 11: Hauptansicht

5 Anleitung

Dieser Abschnitt gibt einen kurzen Überblick über die Vorgehensweise beim Nutzen der entwickelten Knowledgegraphs und des implementierten Codes. Genauere Erläuterungen zu den einzelnen Abschnitten befinden sich im Kapitel 4.

6 Evaluation

Die im Abschnitt 3 gesetzten Ziele wurden erfüllt. Mit Hilfe der DBpedia und den Sparql Anfragen Deutsche Städte Namen und verschiedene dazugehörige Tourism Informationen. Mittles eines Java Codes und Neo4j wurden diese erhaltene Informationen, die Tourism in Deutschland inkorporieren im Form von Graphen dargestellt. Anschließend kann der Nutzer anhand eines Interfaces oder auch ohne nützliche gebrauchte Informationen Abfragen. Er kann sowohl kürzester Strecke Information als auch Attraktionen einer Stadt abfragen. Die Nutzung DBpedia erleichtert die Arbeit enorm, allerdings sind die Informationen nicht unbedingt Vollständig gespeichert wodurch Probleme entstanden sind vorallem beim Städten auflisten und beim Verknüpfen von dieser mit den entsprechenden Sehenswürdigkeiten. Das entwickelte Java-Programm wurde unter Windows mit Java 7 getestet.

7 Zusammenfassung

In der vorliegenden Arbeit wurde mit Hilfe der Neo4j Graphendatenbank durch Tourism in Deutschland zusammenhängende Informationen ein Knowledgegraph generiert. Anfragen an dieser Graph wurden mit Hilfe der Cypher Anfragesprache realisiert. Die Realisation erfolgte durch DBpedia basiertes Sparql Anfragen und einem Java-Modul, welches im Laufe des Projekts implementiert wurde.

Literatur

- [1] Sören Auer, Christian Bizer, Georgi Kobilarov, Jens Lehmann, Richard Cyganiak, and Zachary Ives. *Dbpedia: A nucleus for a web of open data*. Springer, 2007.
- [2] Jeff Heflin. Owl web ontology language-use cases and requirements. *W3C Recommendation*, 10:12, 2004.
- [3] Neo4j. cypher-query-language.
- [4] Roel Popping. Knowledge graphs and network text analysis, 2003.
- [5] Evren Sirin and Bijan Parsia. Sparql-dl: Sparql query for owl-dl. In *OWLED*, volume 258, page 2, 2007.
- [6] Lei Zhang. *Knowledge graph theory and structural parsing*. Twente University Press, 2002.