

Zaverecna sprava projektu

Projekt : Znackovani slov

Katarina Kejstova < 433 820 >

Pavel Kucera < 422 626 >

Simon Struk < 422 522 >

Jiri Vahala < 422 802 >

Zaverecna sprava projektu : Projekt : Znackovani slov

Katarina Kejstova , Pavel Kucera , Simon Struk a Jiri Vahala

Datum vydání 2016-06-07

Obsah

1. Popis projektu	1
Zadanie	1
2. Rozdelenie uloh	2
Pavel	2
Jirka	2
Simon	2
Katarina	2
3. Zaverene spravy	3
Simon	3
Pavel	3
Navrh databaze	3
MVC	3
Katarina	4
Jiri	4

Kapitola 1. Popis projektu

Zadanie

Pre mnohe ulohy strojoveho ucenia je potrebne ziskat oznackovane data. Vasou ulohou bude vytvorit system, ktory umozni ich znackovanie. Anotator odpoveda len na: otazka - odpoved - hodnotenie (vlastna akcia). Aby sa zabranilo typickym problemom pri znackovani je potrebne implementovat niekoľko akcií, ktore napr. pridavaju sum, opakuju otazky, ... Na zaklade vysledkov, by mali byt prezentovane zakladne statistiky.

Kapitola 2. Rozdelenie uloh

Pavel

vytvorenie databaze

spracovanie dat v csv-formate

Jirka

vytvorenie statistik na zaklade ziskanych dat

vybrat vhodne statistiky

Simon

Java backend

JSP server

spracovavanie poziadavkov

Katarina

html frontend, pomocou bootstrap

spracovanie wiki projektu

Kapitola 3. Zaverene spravy

Simon

V ramci mojho podielu na projekte som si nastudoval, ako sa vytvara webova aplikacia pomocou Java Server Pages a Java servletov. Nasledne som implementoval funkcionalitu zabezpecujucu prihlasovanie a odhlasovanie uzivatelov a ochranu pred neautorizovanim pristupom k privatnym castiam aplikacie. Mojou dalsou ulohou bolo vyriesenie nahravania vstupnych dat vo forme CSV suborov.

Pri prihlaseni uzivatela sa skontroluju predane prihlasovacie udaje voci databaze a pri platnych udajoch sa nastavi session atribut, ktory potom umoznuje pristup k privatnym castiam aplikacie. Tento atribut je kontrolovany filtrom a pri jeho neexistencii je uzivatel presmerovany na prihlasovaci stranku. Pri odhlaseni uzivatela sa tento atribut zneplatni.

Tato cast aplikacie bola implementovana ako prva a posluzila mojim kolegom ako referencia pri implementovani dalsich servletov a JSP stranok.

Nahravanie vstupnych dat je riesene pomocou servletu, v ramci ktoreho sa skontroluje format CSV suboru a vytvoria sa prislusne baliky po 1000 slovach. Tieto baliky sa potom uložia do databazy a su pripravené na dalsie spracovanie. Moj navrh bol potom mojimi kolegami prepracovany do jeho finalnej podoby.

Pavel

Navrh databaze

Aplikacni data byla potreba nekma ukladat a protoze jsem v tymu jediny s realnou zkusenosti pouzivani databazi, dobrovolne jsem se chopil vyberu databazove technologie a zpusobu, jakym data ukladat. Jelikož jsem nikdy predtim nepracoval s zadnou NoSQL databazi, zacal jsem se poohlizet prave timto smerem, abych se naucil neco noveho.

Po pocatecnim hledani padla volba na MongoDB [<https://www.mongodb.com/>], dokumentovou databazi s nekterymi charektistikami relacnich technologii. MongoDB mi prisla jako dobry mezikrok, kde bych mohl uplatnit sve znalosti a zaroven se naucit pracit v jinem paradigmatu. Pri navrhu entit samotnych jsem se snazil zohledit budouci pouziti, aby dotazy na data byly efektivni, pricemz jsem vysel ze schematu, jake bych navrhnul v relacni databazi. Pro zachovani efektivity jsem musel misty duplikovat data, coz by sice mela byt bezna praktika, ale stejne mi to nebylo uplne po chuti a velmi rad bych videl nejakou vetsi fungujici aplikaci nad NoSQL databazi, abych se z ni mohl poucit. Nerad znovu vymyslim kolo.

Navrzené schema slo, s mensimi upravami, pouzit pro veskerou zamyslenou funkcnost aplikace, coz beru jako pozitivum. Pri dalsim vyvoji bych se poohlizel po konzultaci s nekym zkusenejším, abych dokazal urcit vykonostni dopady.

MVC

Aby aplikace zvladla i „trvale udrzitelny rust“ jsem se snazil veskery napsany kod uklizet do typickych kategori model, view a controller. To trochu komplikoval fakt, ze jsme nepouzili zadny framework, ktory by podobne rozdeleni podporoval implicitne, ale po pridani malych komponent a nekterych knihoven, se nam podarilo aplikaci rozdelit.

Business logika je nyní kompletne oddelena a k predavani zavislosti uouziva (skoro) vyhradne konstruktory, takže je i lehce prenositelna. Servlety predstavuji jednotlivé controllery a jsp soubory pouzivame jako

sablony. V pripadech stranek bez jakékoli logiky jsme jsp soubory pouzili i bez vytvoreni odpovidajiciho servletu.

Katarina

Mojou ulohou bolo vytvorit frontend pre nasu aplikaciju. Kedze bola poziadavka na responzivnu aplikaciju, zvolili sme framework Bootstrap (<http://getbootstrap.com/>), popularny HTML, CSS a JS framework na vyvoj rezponzivnych mobilnych projektov na webe.

Cely navrh je vo jsp formate, a teda HTML + CSS, vyuzivajuca prvky javy. Vsetka funkcionalita je teda zabezpecena funkciami a servletom, pisanim v jave. Kazdy request je teda najprv spracovany servletom, a nasledne presmerovany na jsp stranku. Servlety zabezpecuju ako prihlasovanie a odhlasovanie, tak aj komunikaciju a zapisovanie do databazy, a taktiez vypis dostupnych balickov a nasledne znackovanie slov, ktore som riesila ja.

Navrh som sa pokusala urobiť čo najviac uzivatelsky privetivy, a jednoducho pochopitelny. Po spusteni aplikacie bolo kazdemu uzivateli umožnene prihlasiť sa, alebo registrovať. Po naslednom prihlásení uzivatel obdržal kratku spravu, vysvetlujucu, čo moze robit, a kde čo najde.

Menu je situovane v hornej casti, a pocas celej doby, s vynimkou samotneho znackovania, sa nemeni. Ponuka moznosti znackovania, a teda vypis aktualne dostupnych balickov, export roznych statistik, nahranie noveho balicka slov, a odhlasenie.

Samotna znackovacia cast je urobena tak, aby bola pohodlne obsluzitelna aj na tablete, a teda obsahuje velku plochu na odpovedanie "ano" vpravo zelenej farby, a "nie" vľavo cervenej farby. Odpoved je zaznamenana kliknutim na danu plochu.

Wiki stranku projektu sme spracovali priebežne, obsahuje ako rozdelenie uloh, tak dokumentaciju vsetkych casti : databaze, struktury, prihlasovanie a autentizacia, forma odpovedania, a nasledne spracovanie, a teda popis sumu a statistik.

Jiri

Mym ukolem bylo zpracovani nashromazdenych dat v databazi a jejich transformace do citelne podoby.

Nashromasdena data slouzi hlavne pro vypocitani statistik nad ruznymi uzivateli, ale i baliky slov nebo jednotlivymi slovy. Takova data mohou byt dale ruzne interpretovana a vyhodnocovana spravcem aplikace.

Jelikoz se jedna o aplikaci anotujici data, není možné pocítat statistiky nad vším stejne. Proto se rozlišují dve hlasovani uzivatele. Prvni typ hlasovani probiha nad slovy, u kterych není známy jejich typ. V takovem pripade muzeme generovat statistiky trueRatio a duration nad jednotlivymi slovy a baliky. Druhe hlasovani je nad slovy, u kterych dopredu typ známe. Tato slova jsou pridana do baliku slov schvalne kvuli overovani kvality hlasovani jednotlivych uzivatele.

Pokud nejaky uzivatel nema dostatecne kvalitni odpovedi na nami známých slovech, je možné jej z hlasovani vyradit nebo nebrat jeho hlasy v potaz. Pro pocitani uspesnosti jednoho uzivatele jsem zvolil vypocet confusion matrix a nad ni její známe metriky, zahrnujici: positive a negative F1-score, Recall, Precision, Random Agreement Probability, Cohen's Kappa coefficient, overall, average a mean classification accuracy.

Zalezi jen na spravci, jak se k daným statistikám postavi a ktere si vybere ke svým kriteriim.

Statistiky jsou ulozeny do formatu XML pres pouziti XStream knihovny. Pouziti je velice jednoduche a primocare.

Jednotlive dotazy do databaze jsou realizovany pouzitim MongoDB knihovny pro jazyk Java.

Do budoucna by se pro vyssi efektivitu daly statistiky implementovat do databaze, která by je on-the-fly prepocitavala. Tento krok by zajistil vyssi efektivitu. Ovsem ani v soucasne verzi neni s efektivitou zadny problem.