

Sem vložte zadání Vaší práce.

ČESKÉ VYSOKÉ UČENÍ TECHNICKÉ V PRAZE
FAKULTA INFORMAČNÍCH TECHNOLOGIÍ
KATEDRA . . . (SOFTWAREVÉHO INŽENÝRSTVÍ)



Bakalářská práce

InfoWeb - Nástroj získávání informací z webů

Vedoucí práce: Ing. Jiří Hunka

11. dubna 2017

Poděkování

Chtěl bych poděkovat za trpělivost vedoucímu, Ing. Jiřímu Hunkovi.

Prohlášení

Prohlašuji, že jsem předloženou práci vypracoval(a) samostatně a že jsem uvedl(a) veškeré použité informační zdroje v souladu s Metodickým pokynem o etické přípravě vysokoškolských závěrečných prací.

Beru na vědomí, že se na moji práci vztahují práva a povinnosti vyplývající ze zákona č. 121/2000 Sb., autorského zákona, ve znění pozdějších předpisů. V souladu s ust. § 46 odst. 6 tohoto zákona tímto uděluji nevýhradní oprávnění (licenci) k užití této mojí práce, a to včetně všech počítačových programů, jež jsou její součástí či přílohou, a veškeré jejich dokumentace (dále souhrnně jen „Dílo“), a to všem osobám, které si přejí Dílo užít. Tyto osoby jsou oprávněny Dílo užít jakýmkoli způsobem, který nesnižuje hodnotu Díla, a za jakýmkoli účelem (včetně užití k výdělečným účelům). Toto oprávnění je časově, teritoriálně i množstevně neomezené. Každá osoba, která využije výše uvedenou licenci, se však zavazuje udělit ke každému dílu, které vznikne (byť jen zčásti) na základě Díla, úpravou Díla, spojením Díla s jiným dílem, zařazením Díla do díla souborného či zpracováním Díla (včetně překladu), licenci alespoň ve výše uvedeném rozsahu a zároveň zpřístupnit zdrojový kód takového díla alespoň srovnatelným způsobem a ve srovnatelném rozsahu, jako je zpřístupněn zdrojový kód Díla.

V Praze dne 11. dubna 2017

.....

České vysoké učení technické v Praze
Fakulta informačních technologií

© 2017 Jakub Tuček. Všechna práva vyhrazena.

Tato práce vznikla jako školní dílo na Českém vysokém učení technickém v Praze, Fakultě informačních technologií. Práce je chráněna právními předpisy a mezinárodními úmluvami o právu autorském a právech souvisejících s právem autorským. K jejímu užití, s výjimkou bezúplatných zákonných licencí, je nezbytný souhlas autora.

Odkaz na tuto práci

Tuček, Jakub. *InfoWeb - Nástroj získávání informací z webů*. Bakalářská práce. Praha: České vysoké učení technické v Praze, Fakulta informačních technologií, 2017.

Abstrakt

V několika větách shrňte obsah a přínos této práce v češtině. Po přečtení abstraktu by se čtenář měl mít čtenář dost informací pro rozhodnutí, zda chce Vaši práci číst.

Klíčová slova Nahradte seznamem klíčových slov v češtině oddělených čárkou.

Abstract

Sem doplňte ekvivalent abstraktu Vaší práce v angličtině.

Keywords Nahradte seznamem klíčových slov v angličtině oddělených čárkou.

Obsah

Úvod	1
1 Cíl práce	3
1.1 Analytické cíle	3
1.2 Praktické cíle	3
2 Analýza a návrh	5
2.1 Získávání informací z webů	5
2.2 Současný stav řešení potřeb internetových obchodů	6
2.3 Cíl týmového projektu	8
2.4 Vývoj a implementace	9
2.5 Zhodnocení současného stavu projektu	11
2.6 Návrh na vylepšení	11
2.7 Analýza nového řešení	11
3 Realizace	13
3.1 Způsob realizace stávajícího řešení	13
3.2 Implementace vylepšení	13
Závěr	15
Literatura	17
A Seznam použitých zkratk	19
B Obsah příloženého CD	21

Seznam obrázků

Úvod

V předmětech BI-SP1 a BI-SP2 v prostředí FIT ČVUT byl realizován týmový projekt umožňující získávání informací z webů s primárním zaměřením na potřeby obchodů. Projekt řešil problém automatizace získávání dat z webů, jelikož stávající služby neposkytují veřejné rozhraní nebo mají velkou chybovost dat.

Práce pojednává o požadavcích internetových obchodů, které jsou především tvořeny nutností držet krok s trhem a sledovat vývoj cen prodáváných produktů u konkurence.

Cílem této práce je popsat požadavky internetových obchodů, stávající stav a možná řešení. Dále na základě těchto poznatků zhodnotit vytvořené řešení a včetně korektnosti zvolených postupů navrhnout vylepšení. Ty implementovat, řádně vylepšení otestovat a zhodnotit výsledný stav projektu.

Cíl práce

1.1 Analytické cíle

1. Rešerše aktuálního stavu získávání dat pro potřeby internetových obchodů
2. Analýza vzniklého řešení týmového projektu vzniklého v prostředí ČVUT FIT, včetně důrazu na použité postupy při softwarovém vývoji
3. Návrh a zhodnocení implementovaných vylepšení

1.2 Praktické cíle

1. Implementace vylepšení systému

Analýza a návrh

V této kapitole se budu nejprve zabývat samotnou problematikou získávání informací z webů s důrazem na internetové obchody. Jelikož je tato problematika již řešena existujícími službami, je nutné je zhodnotit a zanalyzovat jejich funkcionalitu. Dále zhodnotím současný stav projektu, zvolené postupy při vývoji a výslednou funkcionalitu. Nakonec navrhnu vylepšení vzniklého systému, ty nejdůležitější implementuji a výsledný stav zhodnotím.

2.1 Získávání informací z webů

2.1.1 Problematika

Získávání informací z webů je efektivní možnost jak získat databázi informací, které se na internetu vyskytují. Tato činnost však stojí na problematice data získávat a uchovávat v potřebné struktuře, jelikož jinak z dat nejsme schopni vyčíst potřebné informace. Vzhledem k specifitě dat, které jsou v kontextu činnosti zajímavá a dále kvůli unikátnosti webových stránek není možné jednoznačně určit jednotný a zcela automatizovaný postup, jak data získat v požadovaném formátu.

2.1.2 Výběr dat

Nejčastější řešení je kombinace automatizace a prvku lidské inteligence. To je obvykle dosaženo roboty, kteří data stahují a lidské práce určující jaké informace nás ve stažených datech zajímají.

Získávání informací ze stažených stránek lze poté zjednodušit na problematiku určení elementů v HTML, které jsou pro nás zajímavé. Lokaci elementu v HTML se kterým je potřeba pracovat lze poté jednoznačně určit pomocí dvou možností:

1. XPath

2. CSS Selector

2.1.3 XML Path Language

XML Path Language[1] nazývaný zkráceně XPath je jazyk, který slouží k výběru elementu v dokumentu ve formátu XML[2] .

XML chápeme jako jazyk popisující strukturu dat, které jsou strojově i lidsky čitelné. HTML lze chápat jako strukturu podobnou XML, ačkoliv se přímo o XML dokument nejedná [3]. HTML popisuje obsah dat pro prezentaci ve webovém prohlížeči pomocí předem definované struktury, které prohlížeč rozumí. Díky této vlastnosti lze použít XPath pro definování cesty k prvku (a jeho obsahu), který uchovává potřebnou informaci na webové stránce.

2.1.4 CSS Selector

Jazyk CSS je používán pro vizuální popis prezentace webové stránky definované v HTML. Jazyk k určení prvků se kterými chce pracovat používá selektory, které označují prvek v HTML. Jako selektor může být použit jak samotný název prvku, tak vlastní definované třídy.[4]

Pomocí řetězení těchto selektorů jsme schopni jednoznačně získat element v HTML.

2.2 Současný stav řešení potřeb internetových obchodů

I v kontextu malého trhu jako Česká republika se lze bavit o velké konkurenci na poli maloobchodů prodávající své zboží na internetu. Internetové obchody potřebují monitorovat konkurenci a trh. Vzhledem k jejich zaměření je tedy nejvíce zajímaví obchody prodávající stejné zboží. Potřebné informace o prodáváných produktech konkurencí se skládají z následujících hlavních atributů:

1. Cena
2. Inzerovaný název
3. Dostupnost

2.2.1 Srovnávače cen

Data lze získat pomocí srovnávačů cen jako *zbozi.cz*[5] nebo *heureka.cz*[6]. Problém u těchto služeb je však že jsou určeny koncovým zákazníkům pro nalezení nejlepší ceny na trhu pro určitý produkt. S tím souvisí to, že největší srovnávače cen neposkytují data nebo rozhraní přes která by je bylo možné jednoduše získat.

2.2.2 Existující služby

Problematiku sledování trhu s důrazem na firemní klientelu, řeší aktuálně několik existujících služeb.

Služby mají v zásadě velmi podobnou povahu služeb. Rámcově se jedná o porovnávání cen včetně historie na různých internetových obchodem či na srovnávacích cen. Uživatel si zadá okruh či seznam produktů, buďto formou manuální či vstupem ze souboru, případně přímých napojením na e-shop. Následně je možné konkrétní data zobrazit v grafech označující vývoj cen, trendů či náhlých změn. Dále umožňují externí výstup do souboru v dostupných formátech.

Největší rozdíl služeb je zda jsou data získávána přímo z obchodů nebo ze srovnávacích. Další odlišností je možnost zda služba dokáže sledovat i zahraniční trh.

Cena služeb se nejvíce odvíjí od počtu sledovaných produktů a četnosti aktualizací. Proto se měsíční platby mohou pohybovat od stovek korun po desítek tisíc korun.

2.2.2.1 Price checking

2.2.2.2 Pricing intelligence

2.2.2.3 Sledování trhu

2.2.2.4 Pricebot

2.2.2.5 Zahraniční nástroje

Tyto nástroje jsou obecněji zaměřené a obvykle požadují od uživatele techničtější zaměření, jelikož je nutné přesně specifikovat kde, co a jak chce sledovat. Vzhledem k tomuto omezení nejsou přímo pro provozovatele e-shopů vhodné kvůli nedostatečným technickým kapacitám a pro tuto práci důležité.

Bodový seznam zahraničních nástrojů:

1. Screen scraper [7]
 - Webová služba
 - procházení web skrz odkazy
 - potvrzování formulářů
 - využití interního vyhledávání
 - export do širokého množství formátů souborů
 - cena: \$549 - \$2,799 za měsíc
2. Web extractor [8]
 - Windows Aplikace

- procházení zadaných stránek
- hledání stránek pomocí klíčových slov
- export do csv formátu
- cena: \$99 - \$199 jednorázově

3. Web Scraper [9]

2.3 Cíl týmového projektu

V předmětech BI-SP1 a BI-SP2 byl realizován týmový projekt, v souladu s osnovami těchto předmětů byl nejdříve v BI-SP1 vytvořen návrh systému a v BI-SP2 implementován.

Systém je navržen, aby systém umožňoval automatizované získávání informace o produktech prodáváných konkurencí. Důraz je především kladem na automatizaci maximálního počtu procesů a zbylé nechat v rukách administrátora u kterého se předpokládá minimální technické vzdělání. Jediná nutná problematika, co musí administrátor znát je parsování HTML stránek.

Návrh popisuje rozdělení aplikace na část poskytující veškeré webový rozhraní a na část zpracovávající všechny interní procesy. Vzhledem k požadavkům na škálovatelnost aplikace je druhá část z výše zmiňovaných složena z více samostatných menších služeb - modulů komunikující spolu pomocí front. Díky tomu, že každý modul zajišťuje určitou funkcionalitu je možné vytvářet více instancí modulu, čímž je možné procesy zpracovávat paralelně. Uživatelská a interní část spolu sdílejí data pomocí relační databáze[10].

2.3.1 Webové rozhraní

Webové rozhraní lze rozdělit na dvě části. První je uživatelská část, což je množina podstránek určených pouze pro konečné uživatele služby. Uživatelská část umožňuje vytvořit kampaň. Kampaň je proces trvající určitý časový úsek, který sleduje vložené produkty u konkurence. V těchto běžících kampaních má poté možnost uživatel vidět vizualizaci získaných dat, případně je mu umožněn export dat do formátu csv nebo xlsx. Získaná data obsahují, kde se sledované produkty prodávají a za jakou cenu na těchto obchodech.

Druhá část je určena pouze pro administrátory a slouží k monitorování kampaní uživatelů a řešení chyb, které systém není schopný vyřešit. Chyby jsou typicky problémy s parsováním webových stránek, párování produktu ke stránce nebo potvrzení zda jsou získaná data validní.

2.3.2 Interní část

Interní část je rozdělena do samostatných modulů, které spolu komunikují pomocí front. Moduly jsou detailně popsány v následujících pod sekcích.

2.3.2.1 Manager

Manager je hlavní modul, který jako jediný má možnost připojení přímo do sdílené databáze a jeho instance může existovat pouze jednou. Manager má za úkol plánování práce pro ostatní části systému a zpracování vstupů a výstupů z front.

2.3.2.2 Finder

Finder je modul, který má za úkol získávat URL adresy internetových obchodů a na nich vyhledávat URL adresy vedoucí na požadované detaily produktů. Detailem produktu je myšlena webová stránka, kde jsou obsaženy podrobné informace o prodávaném produktu. Typicky je na takové stránce pouze jeden produkt, případně odkazy na jiné detaily prodáváného zboží.

2.3.2.3 DataProvider

DataProvider je module, který zpracovává adresy vedoucí na detaily produktu. Po stažení stránky, se z ní pokusí získat požadované hodnoty. V případě neúspěchu odešle chybovou hlášku, v opačném případě data zanalyzuje korektnost získaných cen vůči historickým datům pokud existují. Výsledek je poté odeslán k zpracování „Managerem“

2.4 Vývoj a implementace

2.4.1 Pojmy

2.4.1.1 Verzovací systém Git

Git je pro sdílení jednotlivých verzí každého vyvojáře, umožňující jednoduchý přehled nad rozpracovanými částmi každého vyvojáře. Úložiště systému se nazývá repozitář, který obsahuje veškerý kód. Základní jednotkou tvoří verze, které jsou postupně vytvářeny vyvojáře po vytvoření každé malé funkcionality. Verze jsou poté uchovávány v jednotlivých větví programu. Vedlejší větve slouží pro práci vyvojáře a oddělení práce vyvojářů. Hlavních větví poté tento kód spojují. Git zajišťuje základní nástroje pro slučování kódu v případě spojování nebo slučování vedlejších větví.

2.4.1.2 Jednotkové a integrační testy

Jednotkovými testy se rozumí sada kladných a záporných testů ověřující funkcionalitu jedné třídy. Integrační testy pokrývají komunikaci více tříd a funkcionalitu závislou na externích službách, např. stažení webové stránky.

2.4.1.3 Statická analýza kódu

Statická analýza kódu je analýza softwerového produktu, která běží bez spuštění samotné aplikace. Kontroluje tedy pouze samotný kód. Označuje kritické konstrukce vedoucí k chybám nebo nedodržení programátorských konvencí daného jazyka.

2.4.1.4 Průběžná integrace

Průběžnou integrací se rozumí sada nástrojů sloužící k rychlému nalezení případných chyb. Základní součástí je vzdálený repozitář, kde na základě každé jeho změny se vytváří nový build verzovaného systému. Při buildy jsou poté spuštěny jednotkové a integrační testy. Na jejich základě je poté možné zjistit případné chyby.

2.4.2 Vývoj

Vývoj byl rozdělen do 5 iterací, z nichž každá obsahovala 10 sprintů. Před začátkem vývoje byla rozdělena práce do těchto částí, s tím, že na konci každé iterace probíhala prezentace vyučujícímu. Každý sprint byl poté průběžně rozdělen na jednotlivé úkoly, které byly přiřazeny členům týmu. Stav úkolů byl uchovávan na systému Redmine, ten umožňuje přehledné řízení projektu a možnost sledování objevených chyb.

Jako verzovací systém byl zvolen systém Git, zajišťován službou Gitlab. Gitlab umožňuje možnost uchovávat repozitář na vzdáleném serveru a poskytuje webové rozhraní pro snadnou správu. Projekt byl v rámci repozitáře rozdělen na 4 části (větvě):

- Master - hlavní větev uchovávající verze určené k nasazení na produkční server
- Develop - vývojová větev uchovávající aktuální stav vývoje
- Feature - větev vytvořená pro konkrétní úkol přidávající novou funkcionality
- Fix - větev určená pro úkoly opravující naleznou chybu

Jelikož přístup k přidání verze do větví Master a Develop měl pouze vedoucí projektu, musel být pro Feature a Fix větvě vytvořen požadavek o zařazení. Ten po kontrole vedoucí projektu zařadil nebo vrátil k opravení.

Na konci každé iterace byla poslední verze vždy označena ve verzovacím systému, ta byla poté prezentována vedoucímu. Označení bylo zvoleno na základě číslování iterací. Například 1. iterace je označena verzí „0.1“.

Pro vývoj se využil princip průběžné integrace. Každá verze byla zkompilována, otestována a zanalyzována na vzdáleném serveru. Tyto činnosti zajišťovaly systém Jenkins. Jenkins aplikaci zkompiloval, pustil testy a statickou analýzu kódu zajištěnou aplikací SonarQube. Výsledky poté publikoval ve svém webovém rozhraní a zároveň v rozhraní Gitlab.

2.4.3 Implementace

2.4.3.1 Webové rozhraní

Webové rozhraní je implemetováno v jazyce PHP verze 7. Základní kamenem aplikace zajišťuje aplikační rámec Nette, poskytující již napsané nástroje pro automatickou správu závislostí, komunikaci s databází, vytváření bezpečných formulářů, zabezpečení aplikace a rozhraní pro tvorbu jednotkových testů. Dále automaticky vynucuje MVC architekturu, která odděluje prezenční a logickou vrstvu.

Snadnou správu závislostí umožňuje poté balíčkovací systém Composer. Na základě definujícího souboru obsahující požadované knihovny a jejich verze, knihovny automaticky stáhne z centrálního repozitáře. Tím je mimo automatického stažení zajištěno, že každý člen týmu pracuje se stejnými verzemi knihoven.

2.4.3.2 Interní část

Interní část je implemetována v jazyce Java verze 8. Aplikace

2.5 Zhodnocení současného stavu projektu

TODO

2.6 Návrh na vylepšení

TODO

2.7 Analýza nového řešení

TODO

Realizace

3.1 Způsob realizace stávajícího řešení

TODO

3.2 Implementace vylepšení

TODO

Závěr

Literatura

- [1] Extensible Markup Language (XML) 1.0 (Fifth Edition) [online]. 2008 [cit. 2017-04-10]. Dostupné z: <https://www.w3.org/TR/2008/REC-xml-20081126/#sec-intro/>
- [2] PERUGINI, Saverio. HTML versus XML. In: Virginia Tech - College of engineering: Department of computer science [online]. Virginia Tech: Virginia Tech, 2002 [cit. 2017-04-01]. Dostupné z: <http://courses.cs.vt.edu/~cs1204/XML/htmlVxml.html>
- [3] HTML5: A vocabulary and associated APIs for HTML and XHTML [online]. 2014 [cit. 2017-04-10]. Dostupné z: <https://www.w3.org/TR/html5/introduction.html#html-vs-xhtml>
- [4] Cascading Style Sheets Level 2 Revision 1 (CSS 2.1) Specification [online]. 2011 [cit. 2017-04-10]. Dostupné z: <https://www.w3.org/TR/CSS21/selector.html#q5.0/>
- [5] Heuréka [online]. [cit. 2017-04-01]. Dostupné z: <http://www.heureka.cz/>
- [6] Zboží [online]. [cit. 2017-04-01]. Dostupné z: <http://www.zbozi.cz/>
- [7] Screen scraper [online]. [cit. 2017-04-01]. Dostupné z: <http://www.screen-scraper.com/>
- [8] Web extractor [online]. [cit. 2017-04-01]. Dostupné z: <http://www.webextractor.com/>
- [9] Web Scraper [online]. [cit. 2017-04-01]. Dostupné z: <http://http://webscraper.io/>
- [10] A Relational Database Overview [online]. = [cit. 2017-04-10]. Dostupné z: <https://docs.oracle.com/javase/tutorial/jdbc/overview/database.html>

Seznam použitých zkratk

XML Extensible markup language

HTML Hypertext Markup Language

CSS Cascading style sheets

Obsah přiloženého CD

	readme.txt.....	stručný popis obsahu CD
	exe	adresář se spustitelnou formou implementace
	src	
	impl.....	zdrojové kódy implementace
	thesis	zdrojová forma práce ve formátu L ^A T _E X
	text	text práce
	thesis.pdf	text práce ve formátu PDF
	thesis.ps	text práce ve formátu PS