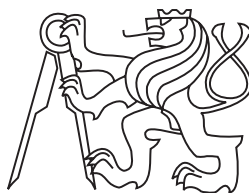


ČESKÉ VYSOKÉ UČENÍ TECHNICKÉ V PRAZE
FAKULTA JADERNÁ A FYZIKÁLNĚ INŽENÝRSKÁ

Katedra softwarového inženýrství v ekonomii
Obor: Inženýrská informatika
Zaměření: Softwarové inženýrství v ekonomii



Modely zátěže výpočetních serverů

Computational requirements modelling

BAKALÁŘSKÁ PRÁCE

Vypracoval: Dmitriy Burdin
Vedoucí práce: Ing. Jan Doubek
Rok: 2014

Před svázáním místo téhle stránky

 s podpisem děkana a do pdf verze oskenované zadání.

Prohlášení

Prohlašuji, že jsem svou bakalářskou práci vypracoval samostatně a použil jsem pouze podklady (literaturu, projekty, SW atd.) uvedené v příloženém seznamu.

V Praze dne

.....
Dmitriy Burdin

Poděkování

Děkuji Ing. Janu Doubkovi za vedení mé bakalářské práce a za podnětné návrhy, které ji obohatily.

Dmitriy Burdin

Název práce:

Modely zátěže výpočetních serverů

Autor: Dmitriy Burdin

Obor: Inženýrská informatika

Druh práce: BAKALÁŘSKÁ PRÁCE

Vedoucí práce: Ing. Jan Doubek
CISCO Systems s.r.o.

Abstrakt: Cílem této práce je efektivně modelovat časovou, výpočetní a paměťovou náročnost distribuovaných úloh. Implementace modelů na základě časových řad výsledků minulých úloh. Výsledkem práce bude studie použitelnosti ekonometrických metod v prostředí velkých data center. Dále návrh algoritmů pro implementaci studovaných metod. Výstupní modely budou použity pro lepší plánování rozvržení výpočetních zdrojů.

Klíčová slova: Časová řada, ekonometrie, predikce, analýza, server

Title:

Computational requirements modelling

Author: Dmitriy Burdin

Abstract: The goal of this thesis is to effectively make a model of temporal, computational and memory requirements of distributed tasks. Model implementation based on time series results of past tasks. The result of the thesis will be a study of applicability of econometric methods in large scale data centres. In addition, design of algorithms for implementation of the studied methods. Output models will be used for better planning of computational sources arrangement.

Key words: Times series, econometrics, forecasting, analysis, server

Obsah

Úvod	7
1. Teoretická část	8
1.1. Časové řady	8
1.1.1. Úvod a charakteristika	8
1.1.2. Analýza časových řad	9
1.1.3. Predikce	12
1.2. Regresní analýza	12
1.2.1. Úvod	12
1.2.2. Lineární regresní model	12
1.2.3. Nelineární regresní model	12
1.3. Softwarová aplikace	13
1.3.1. Základní předpoklady a nástroje	13
1.3.2. Způsob realizace	14
2. Praktická část	15
2.1. Chybné pokusy	15
2.2. Implementace	15
2.2.1. Datová vrstva	15
2.2.2. Aplikační vrstva	15
2.2.3. Prezentační vrstva	15
2.3. Testování aplikace	15
2.3.1. Testovací data	15
2.3.2. Zpracování a analýza dat	15
2.3.3. Predikce	15
2.3.4. Vizualizace	15
Závěr	16
Literatura	17
Přílohy	18
A. Název přílohy	18

Úvod

Obrovské a rostoucí množství informací zaplavují dnešní podniky. To se stavá hlavně kvůli rostoucímu počtu lidí, firem a zařízení připojených k internetu. Kolem třetiny světové populace má dnes přístup k internetu. Velké množství informace, jinak řečeno - velká data, představuje rychle rostoucí sféru, ve které se pro maximální efektivitu vyžaduje vhodná kombinace softwaru, hardwaru a navíc speciální úpravy s ohledem na oblast působnosti. Tyto fakty vedou k výraznému zaměření na velká data a na nové metody správy a analýzování tohoto proudu informací.

Důležitým prvkem řešení této problematiky je datacentrum. Jsou to specializované prostory pro umístění a zajištění stabilního provozu výpočetní a serverové techniky. Většinou se tato technika skládá z počítačových clusterů, což je několik spolupracujících počítačů propojených počítačovou sítí. Clustery se využívají pro výpočet komplikovaných početních úloh. Jeden z hlavních parametrů kvalitních datacenter je správné plánování distribuovaných úloh pro rychlé zpracování dat. Při stále se zvětšujícím množství dat se tento problém komplikuje. Proto je potřeba vědět v jakém pořadí zpracovávat určité úlohy a jak dlouho to bude trvat.

Cílem této bakalářské práce bylo navrhnout a vytvořit softwarovou aplikaci pro analytické zpracování dat, se kterými pracují výpočetní servery. Podstatou aplikace je predikce budoucího chování určitých hodnont na základě analýzy minulých výsledků. Protože velká data obvykle obsahují dynamické systémy dat, které se mění s časem, zkoumaná data se předpokládají být ve formě časových řad. Proto se v této práci nejdříve věnuje způsobům zpracování, analýzy a predikce časových řad. Poté následuje popis algoritmu používání vybraných postupů pro vhodnou analýzu a implementace příslušného algoritmu. Na závěr je demonstrováno využití aplikace na příkladových datech.

1. Teoretická část

1.1. Časové řady

1.1.1. Úvod a charakteristika

Jak již bylo zmíněno, časové řady jsou základním zkoumaným prvkem při analýze různých dynamických systémů obsahující chronologicky uspořádaná data. Časové řady jsou vlastně soubory jednoznačně uspořádaných podle času pozorování v příslušném systému. Data ve formě časových řad vznikají v úplně různých odvětvích buď to fyzikální věda, biologie, společenská věda, medicína atd.

Například ve společenských vědách, časové řady jsou užitečné při popsání počtu obyvatel, porodnosti, nemocnosti. V ekonomii teorie časových řad je jednou z nejdůležitějších metod při analýze ekonomických procesů. Mohou popisovat vývoj určitého ukazatele jako objem výroby, produktivita práce, nezaměstnanost nebo spotřeba surovin. V technice časové řady mohou představovat průběh signálu, spolehlivost nebo intenzitu zatížení elektrického zařízení.

Pro lepší porozumění určitého mechanismu nebo procesu, popsaného časovou řadou, existuje analýza časových řad skládající se z různých metod. Tyto metody pomáhají vytvořit vhodný model popisující chování pozorovaných hodnot. Znalost takového modelu umožňuje kontrolovat činnost a sledovat vývoj určitého systému. Jako důsledek správné provedené analýzy stavá se možným predikovat budoucí chování systému.

Časové řady se skládají ze dvou prvků:

- časový úsek, během kterého byly udělány pozorování
- hodnoty příslušných ukazatelů časové řady

Podle časového úseku časové řady se člení do dvou typů. Jedním je okamžiková časová řada, jejíž pozorování jsou naměřena v jisté časové okamžiky a sčítání hodnot příslušných pozorování nedává smysl. Příkladem může být řada udávající počet zaměstnanců ve firmě na začátku roku. Druhým typem je intervalová časová řada. Pozorování v intervalových časových řadách jsou závislé na délce časového intervalu sledování a v tomto případě už sčítání hodnot ukazatele časové řady dává smysl, neboť hodnotu ukazatele za větší interval lze získat sčítáním hodnot za jednotlivé části příslušného intervalu. Například porodnost ve státě za rok.

Časové úseky jsou obvykle stejnoměrně rozděleny, což znamená, že mezi jednotlivými pozorování stejné časové intervaly. V opačném případě pozorování jsou rozděleny různorodě a tím se analýza časových řad komplikuje.

Obvykle rozlišujeme dva základní modely časových řad:

- Determenistický - model, ve kterém časová řada nemá náhodné prvky a je generována známou matematickou funkcí. Důsledkem je srovnatelně jednoduchá analýza časové řady.
- Stochastický - model popisuje náhodný proces a časová řada obsahuje náhodný prvek. Většina běžně se vyskytujících časových řad v praxi mají stochastické modely.

Jedna z dalších důležitých charakteristik časových řad, vyplývající ze stochastického modelu, je jejich stacionarita, případně nestacionarita. Při stacionaritě střední hodnota a rozptyl časové řady se v čase nemění, což nelze říct o nestacionárních časových řadách, ve kterých se objevují změny ve střední hodnotě či rozptylu. Říká se proto, že nestacionární časové řady mají určitý trend vývoje. I když stacionarita je běžným předpokladem většiny metod analýzy časových řad, některé modely se omezují na modely stacionárních procesů.

1.1.2. Analýza časových řad

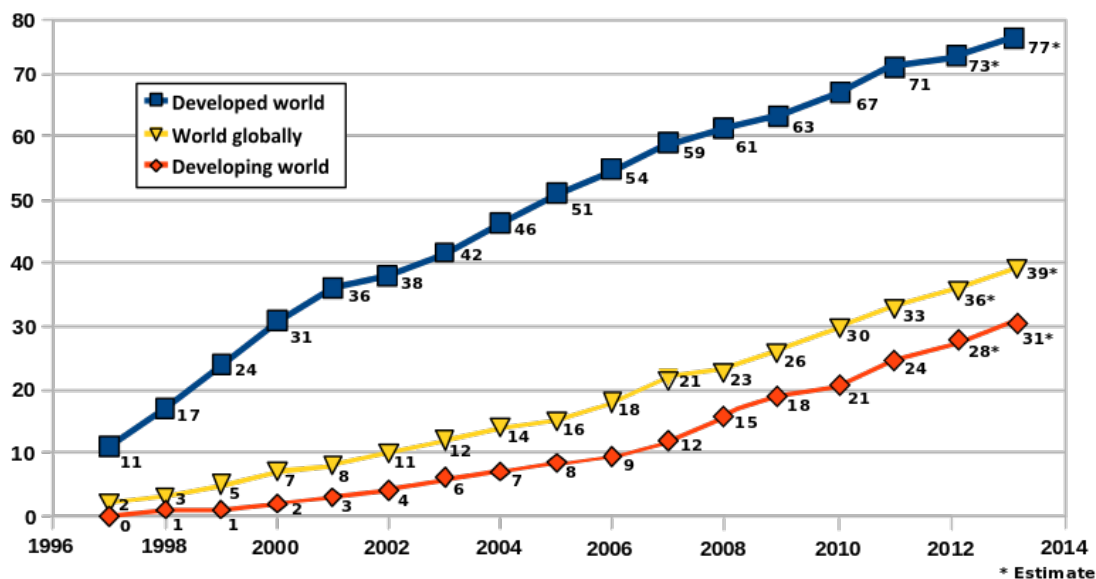
Grafická analýza časových řad

Jedním ze základních prostředků prezentace časových řad je jejich graf. Nejčastěji se graficky znázorňují původní hodnoty časové řady, nebo kumulativní časové řady, které vznikají postupným načítáním (kumulováním) jednotlivých hodnot (u okamžikových časových řad nemají smysl, neboť výše jejich hodnot nezávisí na daném časovém intervalu). Často se ale časové řady zobrazují tak, aby více vynikly jejich charakteristické vlastnosti a rysy. K tomu slouží speciální typy grafů jako jsou například spojnicové grafy jedné nebo více časových řad, krabičkové grafy, grafy sezónních hodnot.

Prvotní informace pro analýzu časových řad získáme ze **spojnicových grafů**. Jejich princip spočívá v zakreslení jednotlivých hodnot časové řady do souřadných os, na kterých jsou vyznačeny příslušné stupnice. Na osu horizontální se vynáší časová proměnná a na osu vertikální hodnoty časové řady nebo její funkce.

Do spojnicového grafu můžeme zakreslit i více časových řad. V případě, že zobrazujeme např. dvě časové řady lišící se měřítkem, je možné použít kromě levé i pravou vertikální osu.

V některých případech je užitečné provést detailnější pohled na časovou řadu. *Krabičkový graf* na rozdíl od jiných grafů obsahuje souhrnné charakteristiky zkoumané časové řady. Tento graf umožní odhalit některé důležité vlastnosti řady, které z jiných grafů nejsou zřetelné. Jeho základním prvkem je krabička, jejíž dolní a horní hrana je tvořena 25% a 75% kvantilem, uvnitř je vyznačen medián a symbolem „+“ aritmetický průměr. Na koncích svislých čar vycházejících z krabičky leží hodnoty minima a maxima. Protože délka této svislé čáry může být maximálně 1,5x delší než krabička, jsou hodnoty přesahující tyto hranice označovány jako odlehlé a jsou zakresleny jako samostatné body.



Obrázek 1.1.: Počet uživatelů Internetu na 100 obyvatel ve světě.
(Příklad spojnicového grafu se třemi časovými řadami)

Dekompozice časových řad

Abychom mohli s řadou pracovat, ve smyslu něco předpovídat, dělat odhady (ať už intervalové či bodové) apod., potřebujeme mít alespoň stacionární časovou řadu. Respektivě alespoň slabě stacionární časovou řadu.

Většina pozorovaných časových řad není řadou nezávislých stejně rozdělených náhodných veličin ani stacionární časovou řadou. Ale mnoho kdy se nám může podařit vhodným očištěním řady, získat stacionární řadu. Eliminace jednotlivých složek časové řady, která se při dekompozici provádí, může mít různé cíle:

1. Zajímavé poznatky může přinést samotné studium oddělených (eliminovaných) složek časové řady, neboť tak lze objevit některé zákonitosti chování sledované řady, rozpoznat vnější vlivy působící na její průběh a provést účinné srovnání průběhu několika časových řad. Např. pomocí trendové složky lze posuzovat tempo rozvoje dvou příbuzných průmyslových odvětví nebo pomocí sezónní složky lze posoudit průběh poptávky po určitém druhu zboží během roku.
2. Podle bodu 1 je úkolem dekompozice časové řady proniknout hlouběji do podstaty historického průběhu řady. Neméně důležitým cílem dekompozice je ale také cíl extrapolační, kdy nás zajímá budoucí vývoj jednotlivých složek časové řady (např. jaké bude budoucí tempo růstu letecké přepravy) a nebo konstruujeme předpověď v celé časové řadě tak, že ji složíme z předpovědí v jednotlivých složkách, které se obvykle sestojí poměrně jednoduše a přesně.
3. Často je vzhledem k podstatě řešeného problému výhodné znát chování časové řady "očištěné" od některých jejích složek. V ekonomických časových řadách se například velice často provádí tzv. sezónní očišťování, kdy z řady odstraníme sezónní fluktuace.

Klasická analýza ekonomických časových řad vychází z předpokladu, že časovou řadu y_t pro $t = 1, 2, \dots, T$ je možné rozložit na čtyři složky: trendovou, cyklickou, sezónní a nesystematickou.

Trendová složka (T_t) vyjadřuje dlouhodobou tendenci vývoje zkoumaného jevu. Je výsledkem faktorů, které dlouhodobě působí stejným směrem např. technologie výroby, demografické podmínky, podmínky na trhu apod.

Cyklická složka (C_t) vyjadřuje kolísání okolo trendu, ve kterém se střídají fáze růstu a poklesu. Jednotlivé cykly (periody) se vytvářejí za období delší než jeden rok a mají nepravidelný charakter, tj. různou délku a amplitudu. Cykly jsou v ekonomických časových řadách způsobeny ekonomickými a neekonomickými faktory. V posledních letech se věnuje pozornost zejména technologickým, inovačním či demografickým cyklům.

Sezónní složka (S_t) vyjadřuje pravidelné kolísání okolo trendu v rámci kalendářního roku. Sezónní výkyvy se opakují každoročně ve stejných obdobích (délka periody je jeden rok) a vznikají v důsledku střídání ročních období nebo vlivem různých institucionalizovaných zvyků, jako jsou např. svátky, dovolené apod.

Poslední složkou časové řady je *nesystematická složka* (I_t nebo a_t). Tato složka vyjadřuje nahodilé a jiné nesystematické výkyvy, ale také chyby měření apod.

Dekompozice časové řady může být:

1. aditivní, hodnoty časové řady se dají určit jako součet hodnot jednotlivých složek, tj.

$$y_t = T_t + C_t + S_t + I_t \quad (1.1)$$

2. multiplikativní, hodnoty časové řady se dají určit jako součin hodnot jednotlivých složek

$$y_t = T_t \cdot C_t \cdot S_t \cdot I_t \quad (1.2)$$

Po aditivní dekompozici jsou jednotlivé složky časové řady ve stejných měrných jednotkách jako původní časová řada. Aditivní dekompozice se používá v případě, že variabilita hodnot časové řady je přibližně konstantní v čase.

Po multiplikativní dekompozici je trendová složka časové řady ve stejných měrných jednotkách jako původní časová řada, ale ostatní složky časové řady (cyklická, sezónní a nesystematická) jsou v relativním vyjádření. Multiplikativní dekompozice se používá v případě, že variabilita časové řady roste v čase, nebo se v čase mění. V praxi se dekompozice časových řad často používá z těchto důvodů:

- a) analýzou jednotlivých složek řady lze odhalit určité zákonitosti vývoje zkoumaného jevu,
- b) časové řady je možné očistit od sezónnosti, tj. z časové řady se odstraní sezónní složka, což umožňuje porovnávat trend několika časových řad současně,
- c) časové řady lze očistit od trendu, tj. z řady se odstraní trendová složka, což umožňuje lépe modelovat sezónnost, protože charakter sezónnosti je výraznější,

- d) často umožňuje přesněji určit předpovědi nejen jednotlivých složek časové řady, ale v konečném důsledku také samotné časové řady, v tom smyslu, že předpovědi jednotlivých složek se sečtou anebo vynásobí podle toho, který typ dekompozice jsme použili.

1.1.3. Predikce

1.2. Regresní analýza

1.2.1. Úvod

1.2.2. Lineární regresní model

1.2.3. Nelineární regresní model

1.3. Softwarová aplikace

1.3.1. Základní předpoklady a nástroje

Softwarová aplikace, která by měla být užitečnou pomůckou pro dosažení cíle práce, může být také považována za výsledek této práce. Její podstatou je možnost nejdřív načíst datový soubor, poté ho zanalyzovat regresní metodou, určit potřebnou informaci ze souboru, použít dosažené výsledky pro predikci budoucích hodnot a vykreslit vývoj nejlépe predikovaných hodnot ve tvaru grafu časových řad.

Jako prvek dat, se kterými pracují výpočetní servery, dostal jsem od svého vedoucího *log file* (nebo taky *žurnál*). Log je to většinou textový soubor s příponou *.log*, který obsahuje záznam činností nějakého programu nebo procesu. Tento záznam je ve tvaru posloupnosti jednotlivých kroků příslušného běhu programu, tj. skládá se z informace o tom jak, kdy, čím nebo kým byla využívána konkrétní aplikace, jaké zatížení bylo vyvoláno nebo kolik času trvalo aby se příslušná aplikace proběhla. Používají se logy pro kontrolu procesů a zjištění přesné příčiny nějaké chyby, pokud k takové došlo. Z jejich záznamů je možno určit co přivedlo k chybě, jaký měla vliv tato chyba a jaké pak měl běžící proces důsledky.

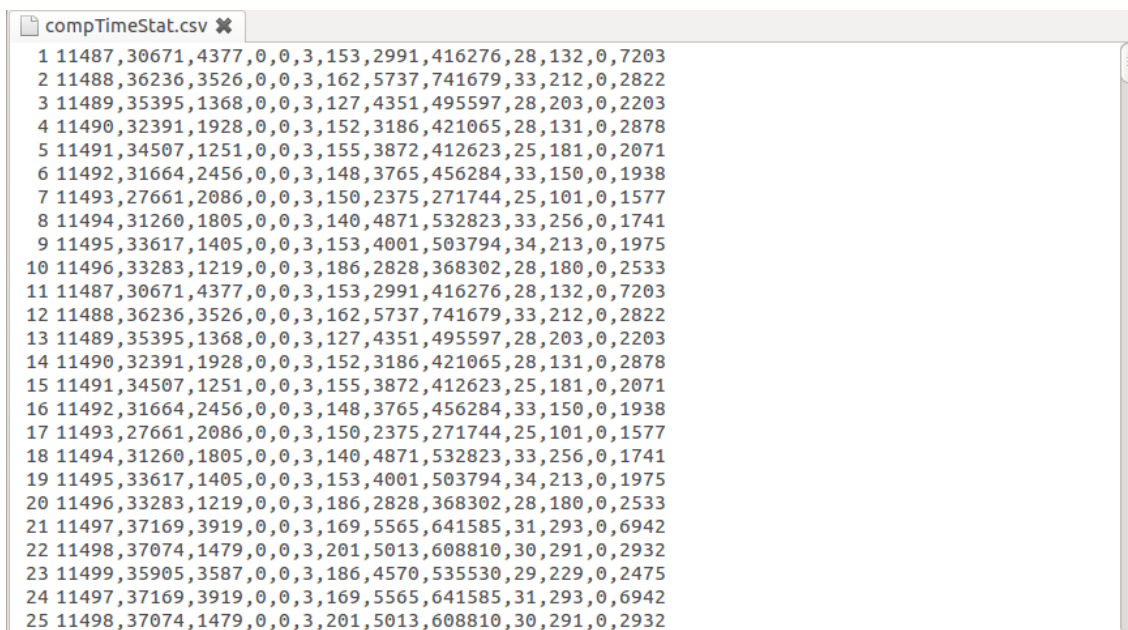
Daný soubor skládá se z dostatečně velkého množství řádků, v každém ze kterých zaznamenáváno několik parametrů. Tyto parametry obsahují svou hodnotu a/nebo svůj název.



```
compTimeStat.log
1 2014.30.03 08:44:52.369 [pool-2-thread-1] INFO 1278716032 l.c.c.a.g.m.GTWLListener - Received message: {headers=
{workerTaskId=1278716032_11487_507, __serialNumberBase=1396165942448, customerId=1278716032, computationTime=
{"dataSetId":"1278716032_http_20140330T0715p0000_11487", "dataSetNumber":11487, "computationTime":7203, "numFlows":30671, "flowLoadingTime":4377
{"numCollections":2, "collectionDuration":1937, "majorGCStats":
{"numCollections":0, "collectionDuration":0, "numPrimitiveIncidents":3, "numEvents":153, "numFlowsInEvents":2991, "resultSize":416276, "numArtif
serialNumber=408, messageType=DataSetInfoProcessingMessage}, body=}
2 2014.30.03 08:44:57.370 [pool-2-thread-4] INFO 1278716032 l.c.c.a.g.m.GTWLListener - Received message: {headers=
{workerTaskId=1278716032_11487_507, __serialNumberBase=1396165942448, customerId=1278716032, computationTime=
{"dataSetId":"1278716032_http_20140330T0720p0000_11488", "dataSetNumber":11488, "computationTime":2822, "numFlows":36236, "flowLoadingTime":3526
{"numCollections":2, "collectionDuration":304, "majorGCStats":
{"numCollections":0, "collectionDuration":0, "numPrimitiveIncidents":3, "numEvents":162, "numFlowsInEvents":5737, "resultSize":741679, "numArtif
serialNumber=409, messageType=DataSetInfoProcessingMessage}, body=}
3 2014.30.03 08:45:02.372 [pool-2-thread-1] INFO 1278716032 l.c.c.a.g.m.GTWLListener - Received message: {headers=
{workerTaskId=1278716032_11487_507, __serialNumberBase=1396165942448, customerId=1278716032, computationTime=
{"dataSetId":"1278716032_http_20140330T0725p0000_11489", "dataSetNumber":11489, "computationTime":2203, "numFlows":35395, "flowLoadingTime":1368
{"numCollections":2, "collectionDuration":253, "majorGCStats":
{"numCollections":0, "collectionDuration":0, "numPrimitiveIncidents":3, "numEvents":127, "numFlowsInEvents":4351, "resultSize":495597, "numArtif
serialNumber=410, messageType=DataSetInfoProcessingMessage}, body=}
4 2014.30.03 08:45:07.373 [pool-2-thread-5] INFO 1278716032 l.c.c.a.g.m.GTWLListener - Received message: {headers=
{workerTaskId=1278716032_11487_507, __serialNumberBase=1396165942448, customerId=1278716032, computationTime=
{"dataSetId":"1278716032_http_20140330T0730p0000_11490", "dataSetNumber":11490, "computationTime":2878, "numFlows":32391, "flowLoadingTime":1928
{"numCollections":1, "collectionDuration":150, "majorGCStats":
{"numCollections":0, "collectionDuration":0, "numPrimitiveIncidents":3, "numEvents":152, "numFlowsInEvents":3186, "resultSize":421065, "numArtif
serialNumber=413, messageType=DataSetInfoProcessingMessage}, body=}
5 2014.30.03 08:45:12.374 [pool-2-thread-4] INFO 1278716032 l.c.c.a.g.m.GTWLListener - Received message: {headers=
{workerTaskId=1278716032_11487_507, __serialNumberBase=1396165942448, customerId=1278716032, computationTime=
{"dataSetId":"1278716032_http_20140330T0735p0000_11491", "dataSetNumber":11491, "computationTime":2071, "numFlows":34507, "flowLoadingTime":1251
{"numCollections":2, "collectionDuration":243, "majorGCStats":
{"numCollections":0, "collectionDuration":0, "numPrimitiveIncidents":3, "numEvents":155, "numFlowsInEvents":3872, "resultSize":412623, "numArtif
serialNumber=415, messageType=DataSetInfoProcessingMessage}, body=}
6 2014.30.03 08:45:12.374 [pool-2-thread-3] INFO 1278716032 l.c.c.a.g.m.GTWLListener - Received message: {headers=
{workerTaskId=1278716032_11487_507, __serialNumberBase=1396165942448, customerId=1278716032, computationTime=
{"dataSetId":"1278716032_http_20140330T0740p0000_11492", "dataSetNumber":11492, "computationTime":1938, "numFlows":31664, "flowLoadingTime":2456
```

Obrázek 1.2.: Log soubor (malá část)

Pro pohodlnější práci se souborem a odstranění nepotřebných parametrů (např. čas, datum) rozhodl jsem se překonvertovat log soubor do souborového formátu *CSV* (comma-separated values, nebo hodnoty oddělené čárkami). Je to jeden z nejjednodušších formátů pro práci s tabulkovými daty. Považuji vybraný formát za vhodný, protože má podobnou strukturu jako u log souboru a je velice používán pro výměnu informací mezi různými systémy. Také dost důležité je, že se sloupce CSV souboru můžou tvářit jako časové řady, které se dobře využívají v regresní analýze.



Obrázek 1.3.: CSV soubor (malá část)

pokračování nasleduje..... popis nástrojů: Java, IntelliJIdea, Github, test-driven development.

1.3.2. Způsob realizace

Po důkladném seznámením s problematikou práce a se související s ní teorií byl stanoven způsob dosažení cílů práce, tj. způsob její praktické realizace. Bylo navrženo použít již existující nebo vytvořit svou knihovnu v jazyce Java pro implementaci metod regresní analýzy. Pomocí této analýzy zjistit určité parametry v datech, které nejméně nebo nejvíc ovlivňují ostatní parametry, a jestli mezi sebou vůbec mají nějakou souvislost. Potom na základě těchto údajů, s úvahou o výsledcích regresní analýzy udělat nejlepší predikci budoucích hodnot určitého parametru. V daném případě jde o parametr *Computation time* (doba trvání výpočtu), jelikož tento parametr vyvolává zatížení, snížení kterého je vlastně úkolem této práce.

možná ještě obrázek, ukazující tento proces realizace. 1 - 2 - 3 -...

2. Praktická část

2.1. Chybné pokusy

Při výběru programovacích nástrojů byla uvažována možnost využití cizích, svobodně distribuovaných knihoven nebo frameworků vhodných pro práci s daty a případně pro jejich analýzu. Takové pomocné nástroje se ovšem předpokládaly být použitelné spolu s programovacím jazykem Java.

Při hledání pomocných nástrojů neboli vhodných knihoven našel jsem dva zdroje, které jak jsem si myslel by mohly zlepšit průběh práce. Po hlubším zkoumání však nebyly tak vhodné než se zdálo na první pohled. Jsou to velká sbírka algoritmů pro práci s daty *Weka* a neurální framework *Encog*.

pokračování následuje..... popis a důvody proč nebyly vhodné weka a encog:

- Weka ...
- Encog ...

2.2. Implementace

2.2.1. Datová vrstva

2.2.2. Aplikační vrstva

2.2.3. Prezentační vrstva

2.3. Testování aplikace

2.3.1. Testovací data

2.3.2. Zpracování a analýza dat

2.3.3. Predikce

2.3.4. Vizualizace

Závěr

Literatura

- [1] Fumio Hayashi. *Econometrics*. Princeton. Princeton University Press. 2000.
- [2] Tomáš Cipra. *Analýza časových řad s aplikacemi v ekonomii*. Praha. SNTL - Nakladatelství technické literatury. 1986.
- [3] Petr Fiala. *Úvod do ekonometrie*. Praha. Nakladatelství ČVUT. 2008.

A. Název přílohy