

Are you ready for the Era of Big Data?

CEDT Meeting
Chiranjib Bhattacharyya
Dept of CSA, IISc

Machine Learning lab
Dept of CSA, IISc
`chiru@csa.iisc.ernet.in`
`http://drona.csa.iisc.ernet.in/~chiru`

Aug 08, 2012

Introduction

custard powder

Introduction

custard powder

Google translate

From: English ▾



To: Swedish ▾

Translate

English to Swedish translation

custard powder

vaniljsås pulver

How does Google Translate work?

- Does not use Rules
- Uses **Statistical Machine Translation**
- Statistical models are trained from **large corpora**

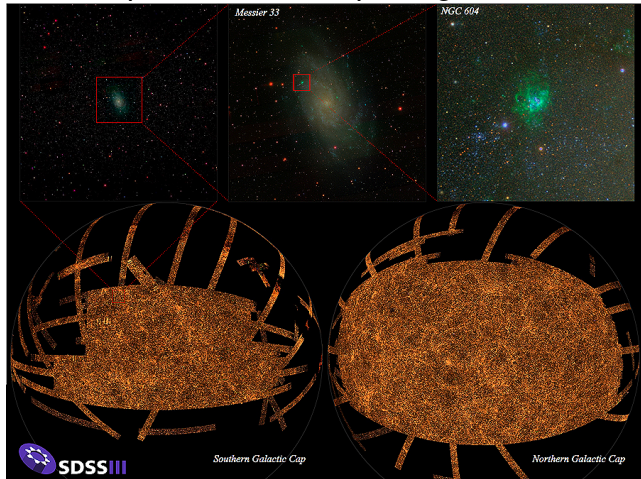
Outline

- 1 What is Big Data?
- 2 Big Data needs Data Scientist
- 3 What is Machine Learning
- 4 Understanding Application Workloads from NFS packet traces
- 5 Mining reviews

What is BigData

(Source: Wikipedia)

Sloan Digital Sky Survey has amassed data of 1.4 TeraBytes.
Need to process 200GB per night.



Examples of BigData

- **Large Hadron Collider** 13000 Terabytes
- **Walmart** handles 1 million customer transactions every hour, 2.5 petabytes
- **Biology** Human genome project took 10 years, now it can be done in one week

What is Big Data?

No precise definition

Working Definition

Datasets so large and complex that they become awkward using on-hand database management tools

Challenges

- capture
- storage
- sharing
- analysis
- visualization

The era of Big Data is upon us

Message from the president(June 2011)
(International Society for Bayesian Analysis)
Prof. M. I. Jordan, Dept of Statistics, UC Berkeley

The era of Big Data is upon us

Message from the president(June 2011)
(International Society for Bayesian Analysis)
Prof. M. I. Jordan, Dept of Statistics, UC Berkeley

In science, massive streams of data have become the norm in areas such as astronomy, high-energy physics, ecology, genetics and molecular biology.

Companies are increasingly looking to hire people who have expertise in data analysis. While the job descriptions sometimes refer to **statistics** they often refer to **data analytics**, **data mining** and **machine learning**.

Are you ready for big data?

Mckinsey Report

McKinsey Global Institute

Research ▾

People

In the news

Contact us

Article | *McKinsey Quarterly*

Are you ready for the era of 'big data'?

October 2011 | by Brad Brown, Michael Chui, and James Manyika

Businesses can gain more if they exploit and understand BigData

Huge Demand for Engineers who can effectively address issues related to Big Data.

Big Data needs Data Scientist

Posted on [LinkedIn](#), 25th June 2012

Data Scientist- "Machine Learning"

Amazon - Bengaluru Area, India



Job Description

Online advertising is set to become a \$100 billion market. Our Advertising Technologies group is building the infrastructure to meet the needs of advertisers and publishers in this rapidly changing industry. We power advanced online advertising programs for some of the world's largest websites, including Amazon.com and other prime online properties. We supply the technology to show the right ad to the right customer at the right time.

The Traffic Validation team in Bangalore is part of Advertising technology group and is completely responsible for ensuring high traffic quality for multiple ad-programs within Amazon. We build systems that process vast amounts of data of the order of tera-bytes on daily basis, using advanced algorithms in data mining, machine learning, and statistical analysis. Our distributed systems are built on cutting edge scalable technologies such as Hadoop and Amazon's EC2 and S3 cloud services.

The data scientist responsible for solving complex large-data problems in online advertising fraud space using data mining, machine learning and statistical analysis. We are looking for a highly motivated individuals who are passionate to apply research to solve actual business problems in fraud and spam detection space. The role involves analyzing large datasets which run into billions of records per day, mine patterns in the data and build models that can detect fraudulent and spam traffic. An ideal candidate should be able to analyze traffic patterns on platforms like Hadoop and develop scalable and low latency models that can detect the patterns.

PhD/MTech/MS or equivalent degree in Computer Science or Mathematics or Statistics
2-5 year relevant industry or research experience

Skill Set

Responsibilities

Solve complex large-data problems in online advertising fraud space using **data mining, machine learning** and **statistical analysis**. Design algorithms which can handle **Tera-Bytes** of data on a daily basis in a **cloud computing**.

- PhD/Masters in CS/Maths/Statistics
- Knowledge of Hadoop and other distributed computing platform
- Experience with Analysis on large scale datasets

Takeaway

Need firm grounding in

- Machine Learning
- Statistics
- Distributed Optimization

Ability to build systems

Takeaway

Need firm grounding in

- Machine Learning
- Statistics
- Distributed Optimization

Ability to build systems **Why not come to IISc?**

What is Machine Learning?

Machine Learning Department - Carnegie Mellon University

http://www.ml.cmu.edu/ Reader

Apple Yahoo! YouTube Wikipedia News Popular

Information for Current Students
Contact Us

What is the Machine Learning Department?

The Machine Learning Department is an academic department within Carnegie Mellon University's School of Computer Science. We focus on research and education in all areas of statistical machine learning. Watch an interview with Tom Mitchell, Department Head:



[Interview with Tom Mitchell](#)

What is Machine Learning?

Machine Learning is a scientific field addressing the question "How can we program systems to automatically learn and to improve with experience?" We study learning from many kinds of experience, such as learning to predict which medical patients will respond to which treatments, by analyzing experience captured in databases of online medical records. We also study mobile robots that learn how to successfully navigate based on experience they gather from sensors as they roam their environment, and computer aids for scientific discovery that combine initial scientific hypotheses with new experimental data to automatically produce refined scientific hypotheses that better fit observed data.

ATTEND

[Machine Learning Lu
Seminar](#)

[ML/Google Distinguis
Lecture Series](#)

[SCS Seminars](#)

[Statistics Seminar](#)

CALENDAR O EVENTS

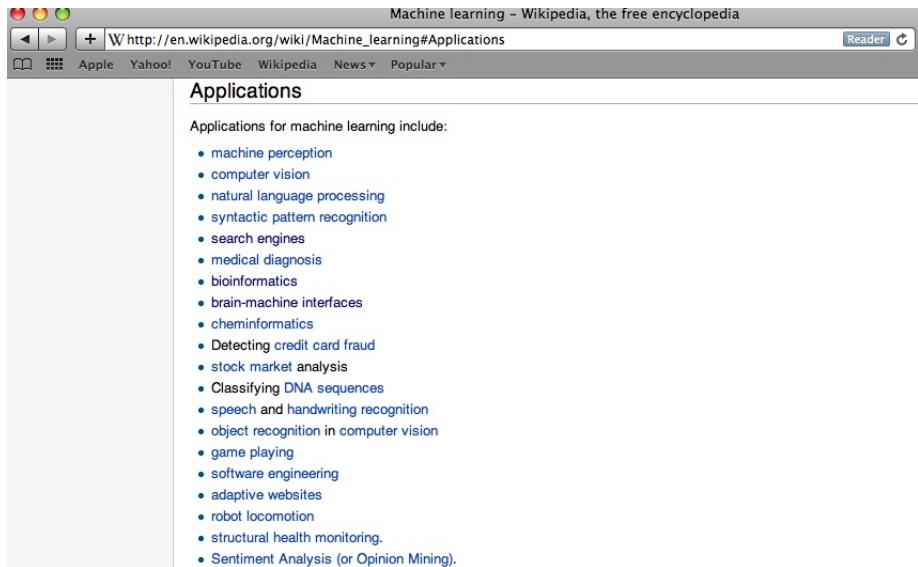
[ML IC Schedul](#)

What is Machine Learning?(CMU ML Dept)

Machine Learning is a scientific field addressing the question

How can we program systems to automatically learn and to improve with experience?

Scope of Machine Learning



The image shows a screenshot of a web browser window. The title bar reads "Machine learning - Wikipedia, the free encyclopedia". The address bar shows the URL "http://en.wikipedia.org/wiki/Machine_learning#Applications". The browser's search bar contains "W" and the address bar has a "Reader" button. Below the browser interface, the Wikipedia page content is visible, starting with the heading "Applications" and a list of applications for machine learning.

Applications

Applications for machine learning include:

- machine perception
- computer vision
- natural language processing
- syntactic pattern recognition
- search engines
- medical diagnosis
- bioinformatics
- brain-machine interfaces
- cheminformatics
- Detecting credit card fraud
- stock market analysis
- Classifying DNA sequences
- speech and handwriting recognition
- object recognition in computer vision
- game playing
- software engineering
- adaptive websites
- robot locomotion
- structural health monitoring.
- Sentiment Analysis (or Opinion Mining).

Predictive models

- $y = \{1, -1\}$ binary Classification

Predictive models

- $y = \{1, -1\}$ binary Classification
- $y = \{1, \dots, k\}$ multi-category classification

Predictive models

- $y = \{1, -1\}$ binary Classification
- $y = \{1, \dots, k\}$ multi-category classification
- Ordinal regression

Predictive models

- $y = \{1, -1\}$ binary Classification
- $y = \{1, \dots, k\}$ multi-category classification
- Ordinal regression
- Regression

Predictive models

- $y = \{1, -1\}$ binary Classification
- $y = \{1, \dots, k\}$ multi-category classification
- Ordinal regression
- Regression
- Reinforcement learning

Predictive models

- Need to quantify the fit between the target y and the prediction $f(x)$ on D

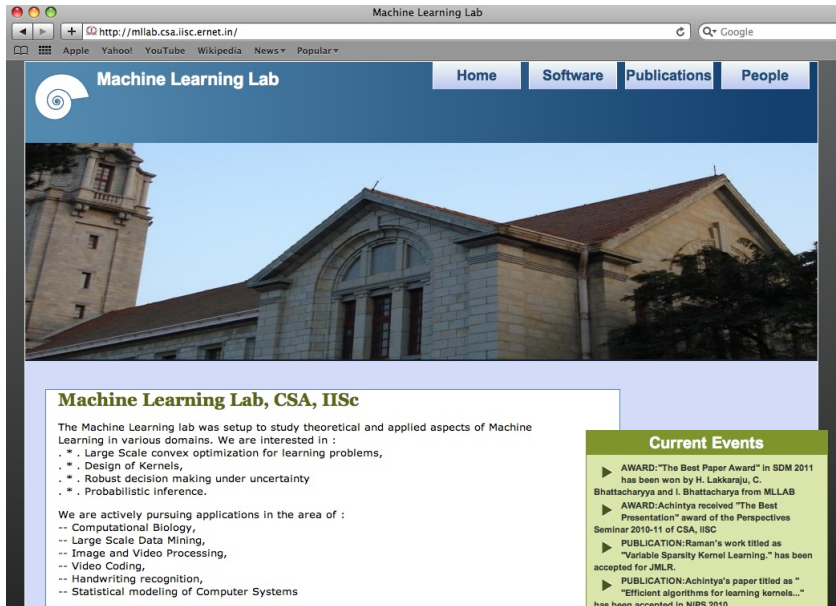
Predictive models

- Need to quantify the fit between the target y and the prediction $f(x)$ on D
- the model should hold for all x , even on data not present in D sometimes called *generalization ability*.

Predictive models

- Need to quantify the fit between the target y and the prediction $f(x)$ on D
- the model should hold for all x , even on data not present in D sometimes called *generalization ability*.
- The problem now becomes finding a model f , which *generalizes well* and fits the training data

What we do at IISc



Machine Learning Lab

Home Software Publications People

Machine Learning Lab, CSA, IISc

The Machine Learning lab was setup to study theoretical and applied aspects of Machine Learning in various domains. We are interested in :

- . * . Large Scale convex optimization for learning problems,
- . * . Design of Kernels,
- . * . Robust decision making under uncertainty
- . * . Probabilistic inference.

We are actively pursuing applications in the area of :

- Computational Biology,
- Large Scale Data Mining,
- Image and Video Processing,
- Video Coding,
- Handwriting recognition,
- Statistical modeling of Computer Systems

Current Events

- ▶ AWARD: "The Best Paper Award" in SDM 2011 has been won by H. Lakkaraju, C. Bhattacharyya and I. Bhattacharya from MLLAB
- ▶ AWARD: Achintya received "The Best Presentation" award of the Perspectives Seminar 2010-11 of CSA, IISc
- ▶ PUBLICATION: Raman's work titled as "Variable Sparsity Kernel Learning." has been accepted for JMLR.
- ▶ PUBLICATION: Achintya's paper titled as "Efficient algorithms for learning kernels..." has been accepted in NIPS 2010.

What we do at IISc

- Large Scale convex optimization for learning problems,
- Design of Kernels
- Robust decision making under uncertainty
- Probabilistic inference
- Computational Biology,
- Image and Video Processing,
- Statistical modeling of Computer Systems
- Large Scale Text Mining

What we do at IISc

- Large Scale convex optimization for learning problems,
- Design of Kernels
- Robust decision making under uncertainty
- Probabilistic Inference
- Computational Biology,
- Image and Video Processing,
- Statistical modeling of Computer Systems
- Large Scale Text Mining

Discovery of Application Workloads from Network File Traces
Yadwadkar N. , Bhattacharyya C., K. Gopinath, N. Thirumale, S.
Susarla
FAST 2010

Focusing on the Opcode Sequence

Time Stamp	Source IP	Destination IP	OPCODE	Parameters
01:39:39.146061	10.192.25.47	10.192.25.18	NFS V2 LOOKUP	Call, DH:0x7c58939e/.Trash
01:39:39.146099	10.192.25.18	10.192.25.47	NFS V2 LOOKUP	Reply (Call In 12) Error:NFS3ERR_NOENT
01:39:39.146667	10.192.25.47	10.192.25.18	NFS V2 LOOKUP	Call, DH:0x7c58939e/.Trash-1003
01:39:39.146701	10.192.25.18	10.192.25.47	NFS V2 LOOKUP	Reply (Call In 14) Error:NFS3ERR_NOENT
01:39:39.151332	10.192.25.47	10.192.25.18	NFS V2 ACCESS	Call, FH:0x7c58939e
01:39:39.151381	10.192.25.18	10.192.25.47	NFS V2 ACCESS	Reply (Call In 16)
01:39:39.152028	10.192.25.47	10.192.25.18	NFS V2 LOOKUP	Call, DH:0x7c58939e/.Trash-1001
01:39:39.152072	10.192.25.18	10.192.25.47	NFS V2 LOOKUP	Reply (Call In 18) Error:NFS3ERR_NOENT
01:39:42.115414	10.192.25.47	10.192.25.18	NFS V2 GETATTR	Call, FH:0xe21586f0
01:39:42.115481	10.192.25.18	10.192.25.47	NFS V2 GETATTR	Reply (Call In 69) Regular File mode:0755 uid:1000 gid:1000
01:39:42.115999	10.192.25.47	10.192.25.18	NFS V2 READ	FH:0xe21586f0 Offset:0 Len:16384
01:39:42.116050	10.192.25.18	10.192.25.47	NFS V2 READ	Reply (Call In 72) Len:16384[Unreassembled Packet [incorrect TCP checksum]]
01:39:42.118382	10.192.25.47	10.192.25.18	NFS V2 GETATTR	Call, FH:0xe21586f0
01:39:42.118423	10.192.25.18	10.192.25.47	NFS V2 GETATTR	Reply (Call In 85) Regular File mode:0755 uid:1000 gid:1000
01:39:42.119314	10.192.25.47	10.192.25.18	NFS V2 LOOKUP	Call, DH:0xf21596f0/conf29395.sh
01:39:42.119363	10.192.25.18	10.192.25.47	NFS V2 LOOKUP	Reply (Call In 87) Error:NFS3ERR_NOENT
01:39:42.119892	10.192.25.47	10.192.25.18	NFS V2 CREATE	Call, DH:0xf21596f0/conf29395.sh Mode:UNCHECKED
01:39:42.121303	10.192.25.18	10.192.25.47	NFS V2 CREATE	Reply (Call In 89)
01:39:42.122465	10.192.25.47	10.192.25.18	NFS V2 WRITE	Call, FH:0xe53f81da Offset:0 Len:11 UNSTABLE
01:39:42.122526	10.192.25.18	10.192.25.47	NFS V2 WRITE	Reply (Call In 93) Len:11 UNSTABLE
01:39:42.123036	10.192.25.47	10.192.25.18	NFS V2 COMMIT	Call, FH:0xe53f81da
01:39:42.128024	10.192.25.18	10.192.25.47	NFS V2 COMMIT	Reply (Call In 95)
01:39:42.128539	10.192.25.47	10.192.25.18	NFS V2 GETATTR	Call, FH:0xf21596f0
01:39:42.128575	10.192.25.18	10.192.25.47	NFS V2 GETATTR	Reply (Call In 97) Directory mode:0777 uid:1000 gid:1000
01:39:42.129092	10.192.25.47	10.192.25.18	NFS V2 GETATTR	Call, FH:0xe53f81da
01:39:42.129151	10.192.25.18	10.192.25.47	NFS V2 GETATTR	Reply (Call In 99) Regular File mode:0644 uid:1000 gid:1000
01:39:42.129662	10.192.25.47	10.192.25.18	NFS V2 ACCESS	Call, FH:0xe53f81da
01:39:42.129695	10.192.25.18	10.192.25.47	NFS V2 ACCESS	Reply (Call In 101)
01:39:42.130219	10.192.25.47	10.192.25.18	NFS V2 WRITE	Call, FH:0xe53f81da Offset:0 Len:18 UNSTABLE
01:39:42.130271	10.192.25.18	10.192.25.47	NFS V2 WRITE	Reply (Call In 103) Len:18 UNSTABLE
01:39:42.130772	10.192.25.47	10.192.25.18	NFS V2 COMMIT	Call, FH:0xe53f81da
01:39:42.135557	10.192.25.18	10.192.25.47	NFS V2 COMMIT	Reply (Call In 105)
01:39:42.137001	10.192.25.47	10.192.25.18	NFS V2 SETATTR	Call, FH:0xe53f81da
01:39:42.141261	10.192.25.47	10.192.25.18	NFS V2 SETATTR	Reply (Call In 107)

Workload Identification

ACCESS
ACCESS
LOOKUP
LOOKUP
LOOKUP
LOOKUP
LOOKUP

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

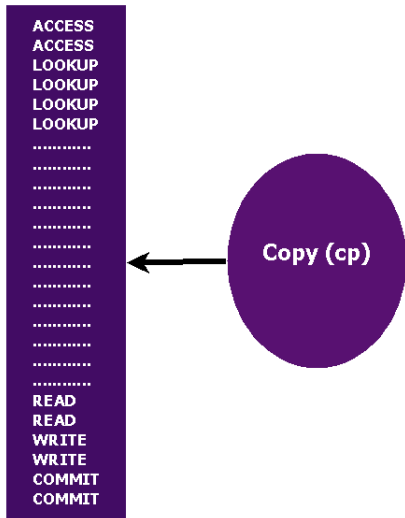
.....

.....

READ
READ
WRITE
WRITE
COMMIT
COMMIT



Workload Identification



Trace Annotation

ACCESS
ACCESS
LOOKUP
LOOKUP

.....

.....

.....

.....

.....

.....

.....

.....

.....

READ
WRITE
WRITE
COMMIT
COMMIT
ACCESS
ACCESS
LOOKUP

.....

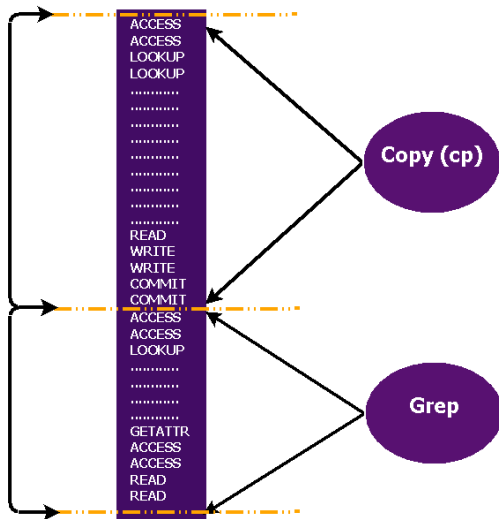
.....

.....

.....

GETATTR
ACCESS
ACCESS
READ
READ

Trace Annotation



Challenges

- Variability in the traces of the same workload

cp contacts.csv con.csv

```
ACCESS Call, FH:0xe003db8b
LOOKUP Call, DH:0xe003db8b/con.csv
LOOKUP Reply Error:NFS3ERR_NOENT
LOOKUP Call, DH:0xe003db8b/contacts.csv
LOOKUP Reply, FH:0x71d9fc7c
GETATTR Call, FH:0x71d9fc7c
ACCESS Call, FH:0x71d9fc7c
CREATE Call, DH:0xe003db8b/con.csv
SETATTR Call, FH:0x58d9d57c
GETACL Call
GETATTR Call, FH:0x58d9d57c
READ Call, FH:0x71d9fc7c ...
WRITE Call, FH:0x58d9d57c ...

COMMIT Call, FH:0x58d9d57c
```

cp contacts.csv dir/con.csv

```
LOOKUP Call, DH:0xe003db8b/dir
```

```
LOOKUP Reply, FH:0x0eb18814
```

```
ACCESS Call, FH:0x0eb18814
LOOKUP Call, DH:0x0eb18814/con.csv
LOOKUP Reply Error:NFS3ERR_NOENT
LOOKUP Call, DH:0xe003db8b/contacts.csv
LOOKUP Reply, FH:0x71d9fc7c
GETATTR Call, FH:0x71d9fc7c
ACCESS Call, FH:0x71d9fc7c
CREATE Call, DH:0x0eb18814/con.csv
SETATTR Call, FH:0x14b19214
GETACL Call
GETATTR Call, FH:0x14b19214
READ Call, FH:0x71d9fc7c ...
WRITE Call, FH:0x14b19214 ...

COMMIT Call, FH:0x14b19214
```

Key Contributions

- Identifying workloads from NFS opcodes
- Identifying transitions between workloads in a trace sequence
- Small snippets of traces are sufficient!
- Exploited the analogy with Biological sequence analysis problem
- Use of Profile Hidden Markov Models(Profile HMMs)

Analogy with Problem in Computational Biology

Unfortunately,

DP formulations for aligning r sequences, each of length n are expensive, $O(n^r)$, time and space complexity

Computational Biology

Conserved in critical regions

Diverge due to chance mutations

Problem at Hand

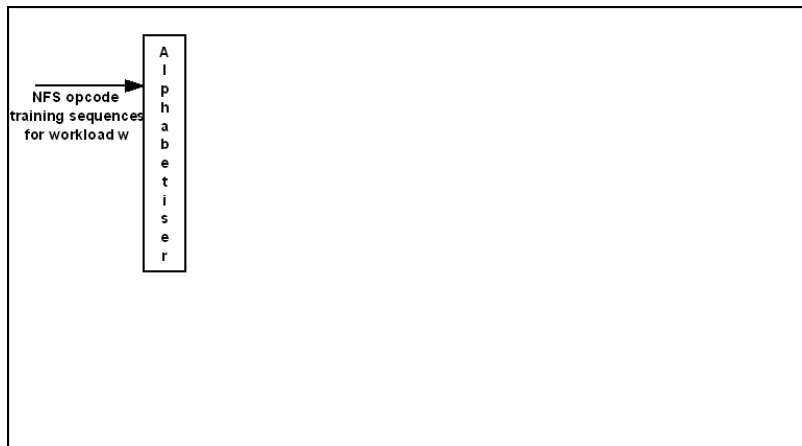
Similarity to a large extent

Additions, deletions and replacements of symbols observed

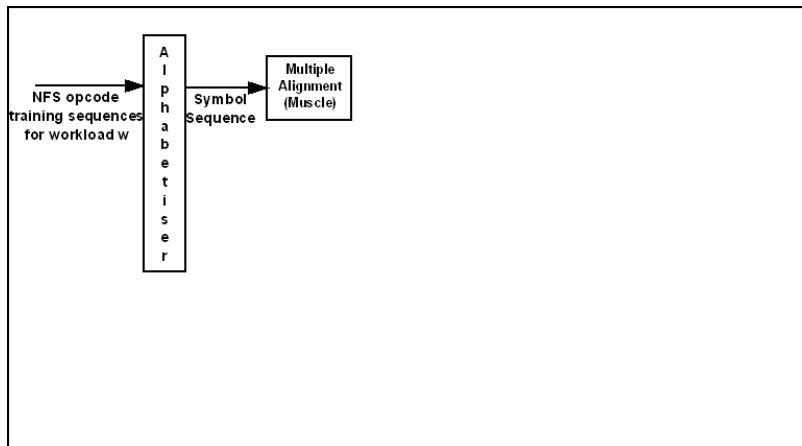
Proposal

To use Profile HMM for representing Profile of a workload

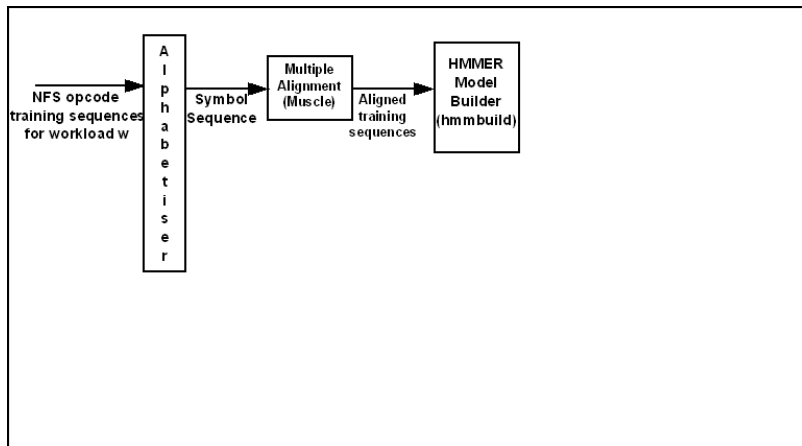
Profile HMMs Training and Usage Workflow: An Overview



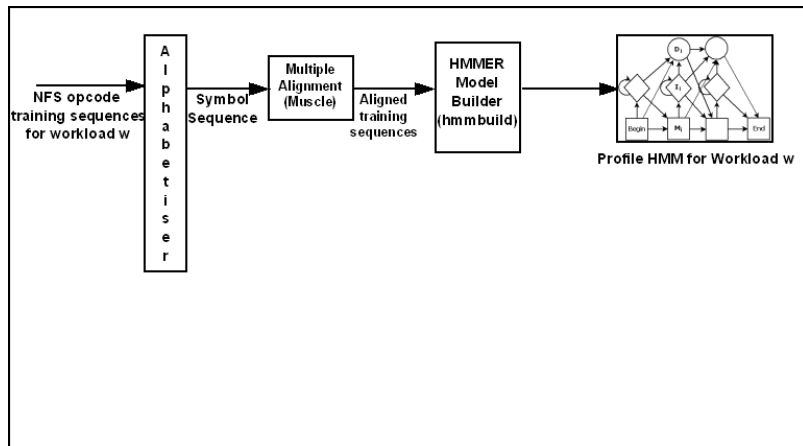
Profile HMMs Training and Usage Workflow: An Overview



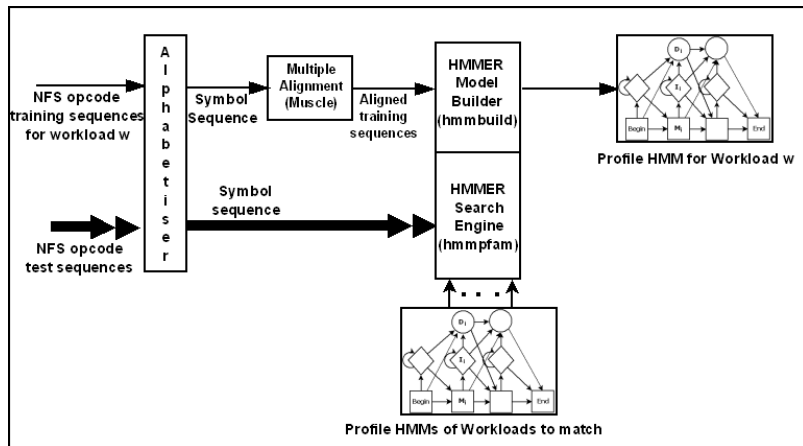
Profile HMMs Training and Usage Workflow: An Overview



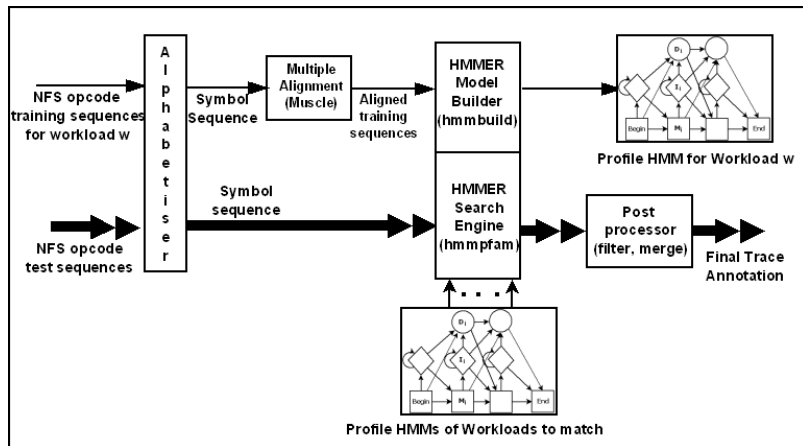
Profile HMMs Training and Usage Workflow: An Overview



Profile HMMs Training and Usage Workflow: An Overview



Profile HMMs Training and Usage Workflow: An Overview



Experimental Set-up

- Unix Commands
 - ▶ tar, untar, make, edit, copy, move, grep, find, compile
 - ▶ Accessing a subset of 14361 files and 1529 directories up to 7 levels deep
- TPC-C
 - ▶ 1 to 5 warehouses with 1 to 5 database clients per warehouse
- Postmark
 - ▶ a workload approximating a large internet email server

Experimental Results

Workload Identification Confusion Matrix

Trace Command	Models								
	make	find	grep	tar	untar	copy	move	edit	tpcc
make	91.7	1.2		1.2	2.4	3.6			
find		91.8	2.1			3.1	1		2.1
grep	1.3	1.3	85	1.3	11.3				
tar				100					
untar				1.2	98.8				
copy		1	1		6	82	1	9	
move		5.6	0.8	0.8		2.4	89.6	0.8	
edit								100	
tpcc									100

Best paper award

Exploiting Coherence in Reviews for Discovering Latent Facets and associated Sentiments

Himabindu L., Bhattacharyya C., Bhattacharya I., Merugu Srujana.
Siam Data Mining Conference 2011

Mining Customer Reviews

★★★★★ A great companion that you will want to take anywhere
I ordered my Acer Aspire One 10.1 Netbook (AOD150-1165) on Feb 28, 2009 from Amazon. Item shipped on March 2nd and arrived on March 14.

When I ordered only Sapphire Blue was available but I am happy with Sapphire Blue. It looks very attractive. I would normally order either white or black model but I am glad they were not available :) The only thing I...

★☆☆☆☆ BEWARE - screen damages with even the slightest touch
It and it was working very well.. UNTIL tonight I touched the screen and it cracked. I know have a black mark on the screen but the size of two dimes and a crack mark... very unhappy Like I said I barely touched it. I subsequently found a full page of reviews about the screen being very cheap..

Facet	Sentiment
Memory	-
Screen	-
Appearance	Positive

Facet	Sentiment
Memory	-
Screen	Negative
Appearance	-

- **Central Problem:** Facet based sentiment analysis of customer reviews
- **Applications**
 - ▶ E-commerce : product recommendation for customers
 - ▶ Business Analytics : aiding product managers and decision makers in understanding the product's market standing

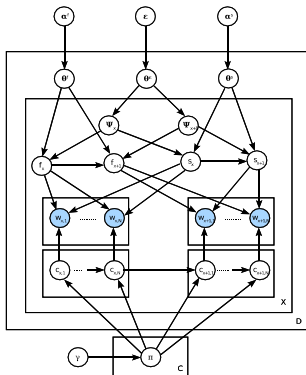
FACeT Sentiment extraction model (FACTS)

The *pictures* i took during my last trip with this camera were absolutely *great*. The *picture quality* is amazing and the pics come out *clear* and *sharp*. I am also very *impressed* with its *battery life*, unlike other cameras available in the market, the *charge* lasts *long* enough. However, I am *unhappy* with the *accessories*.

- FACTS aims at extracting both facets as well as associated sentiments from customer reviews
- Captures both the syntactic and semantic dependencies
- Loosely based on HMM LDA
- Facet and Sentiment classes comprise of topics

FACTS Model

Extends HMM-LDA to include topics within another syntactic class for **sentiments**



Choose $\theta_d^f \sim \text{Dir}(\alpha^f)$

Choose $\theta_d^s \sim \text{Dir}(\alpha^s)$

Choose $\theta_d^c \sim \text{Dir}(\epsilon)$

For each window $x \in \{1 \dots X_d\}$

a. Choose $\psi_{d,x} \sim \text{Mult}(\theta_d^s)$

b. if ($\psi_{d,x}=0$)

Choose $f_{d,x} = f_{d,x-1}$ and $s_{d,x} = s_{d,x-1}$

else if ($\psi_{d,x}=1$)

Choose $f_{d,x} \sim \theta_d^f$ and $s_{d,x} = s_{d,x-1}$

else Choose $f_{d,x} \sim \text{Mult}(\theta_d^f)$ and $s_{d,x} \sim \text{Mult}(\theta_d^s)$

c. For each word i in the window x

i. Choose $c_{d,x,i} \sim \text{Mult}(\pi^{c_{d,x,i-1}})$

ii. if $c_{d,x,i} = 1$, Choose $w_{d,x,i} \sim \text{Mult}(\phi_{f_{d,x}}^f)$

else if $c_{d,x,i} = 2$, Choose $w_{d,x,i} \sim \text{Mult}(\phi_{s_{d,x}}^s)$

else Choose $w_{d,x,i} \sim \text{Mult}(\phi_{c_{d,x,i}}^c)$

$c_{d,i} = 1 \Rightarrow$ **facet**

$c_{d,i} = 2 \Rightarrow$ **sentiment**

Coherence based FACTS model (CFACTS)

The *pictures* i took during my last trip with this camera were absolutely *great*. The *picture quality* is amazing and the pics come out *clear* and *sharp*. I am also very *impressed* with its *battery life*, unlike other cameras available in the market, the *charge* lasts *long* enough. However, I am *unhappy* with the *accessories*.

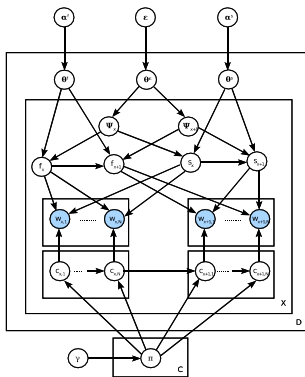
- **Coherence** is an important aspect of user generated content
- In case of reviews, *facet* and *sentiment* coherence are usually prevalent

Modeling Coherence

- Each review comprises of basic units of coherence windows
- Each window is associated with a single facet and sentiment
- Continuity of topics across windows governed by parameter ψ
 - ▶ $\psi = 0$: $f_{d,x} = f_{d,x-1}$ and $s_{d,x} = s_{d,x-1}$
 - ▶ $\psi = 1$: $f_{d,x} = \theta_d^f$ and $s_{d,x} = s_{d,x-1}$
 - ▶ $\psi = 2$: $f_{d,x} = \theta_d^f$ and $s_{d,x} = \theta_d^s$

CFACTS Model

- Extends FACTS to incorporate coherence in facets/sentiments
- Also, enables loose coupling of the facet and sentiment classes



Choose $\theta_d^f \sim \text{Dir}(\alpha^f)$

Choose $\theta_d^s \sim \text{Dir}(\alpha^s)$

Choose $\theta_d^\epsilon \sim \text{Dir}(\epsilon)$

For each window $x \in \{1 \dots X_d\}$

a. Choose $\psi_{d,x} \sim \text{Mult}(\theta_d^\epsilon)$

b. if $(\psi_{d,x}=0)$

Choose $f_{d,x} = f_{d,x-1}$ and $s_{d,x} = s_{d,x-1}$

else if $(\psi_{d,x}=1)$

Choose $f_{d,x} \sim \theta_d^f$ and $s_{d,x} = s_{d,x-1}$

else Choose $f_{d,x} \sim \text{Mult}(\theta_d^f)$ and $s_{d,x} \sim \text{Mult}(\theta_d^s)$

c. For each word i in the window x

i. Choose $c_{d,x,i} \sim \text{Mult}(\pi^{c_{d,x,i-1}})$

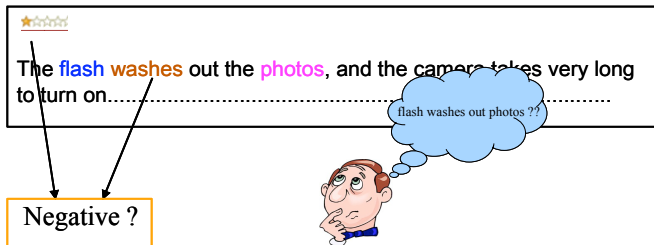
ii. if $c_{d,x,i} = 1$, Choose $w_{d,x,i} \sim \text{Mult}(\phi_{f_{d,x}}^f)$

else if $c_{d,x,i} = 2$, Choose $w_{d,x,i} \sim \text{Mult}(\phi_{s_{d,x}}^s)$

else Choose $w_{d,x,i} \sim \text{Mult}(\phi_{c_{d,x,i}}^c)$

Incorporating ratings - CFACTS-R

- Review ratings are valuable pointers to the sentiments expressed in reviews
- Does incorporating these review ratings help us extract sentiments better ?
 - ▶ Review ratings turn out be of immense help for 'ordering sentiment topics'



Experimental results

Qualitative Evaluation

Digital Camera Corpus

Model	Topic Label	Top Words
CFACTS-R (all topics)	Price Ease of use Picture quality Accessories Display Battery life Portability Features	fit, purse, pocket, pay, worth ease, weight, casing, <u>digicam</u> , travel shots, video, <u>camera</u> , images, pics charger, cable, battery, controls, button digital, viewfinder, shots, lens, clarity aa, batteries, life, <u>ease</u> , charge travel, ease, bags, portability, straps lens, memory, <u>point-and-shoot</u> , software
CFACTS (all topics)	Battery life Accessories Picture quality Features Ease of use Display Price Portability	battery, charge, <u>shutter</u> , aa batteries, alkaline charger, <u>camera</u> , cable, tripod, shutter button images, clarity, camera, brightness, focus zoom, <u>nikon</u> , face recognition, redevye, memory ease, use, design, <u>color</u> , grip <u>slr</u> , lcd, viewfinder, display, point-and-shoot price, worth, discount, warranty, fit ease, portability, size, lightweight, travel
FACTS-R (out of 8 topics)	Accessories Lens Portability - Picture quality	buttons, tripod, controls, batteries, purse shutter, <u>minolta</u> , <u>camera</u> , point-and-shoot range, size, weight, bag, design memory, quality, purchase, warranty, cams pictures, quality, images, resolution, sharp
FACTS (out of 8 topics)	Lens Portability - Picture quality Accessories	shutter, lens, <u>camera</u> , point-and-shoot range, size, weight, bag, design pics, shots, range, ease, straps pictures, quality, images, resolution, sharp buttons, controls, charger, tripod, purse
LDA (out of 9 topics)	Accessories - Picture quality -	replace, charger, reader, <u>digicam</u> , <u>easy</u> take, shoot, carry, great, <u>easy</u> images, <u>camera</u> , pics, <u>like</u> , <u>good</u> charger, lens, awful, camera, shutter

Quantitative Evaluation

Facet Coverage - the fraction of extracted facets that actually correspond to product attributes. Benchmarked against amazon's structured ratings facets

Facet Purity – the fraction of the top words in the facet that actually correspond to the product attribute

Corpus	Model	Facet Coverage(%)	Topic Purity(%)
Digital Cameras	CFACTS-R	100	80.18
	CFACTS	100	84
	FACTS-R	33	74.73
	FACTS	33	72.28
	LDA	16.67	44.37
Laptops	CFACTS-R	83.33	87.09
	CFACTS	83.33	87.09
	FACTS-R	33.33	74.19
	FACTS	33.33	77.41
	LDA	33.33	45.16
Mobile Phones	CFACTS-R	80	91.48
	CFACTS	80	89.36
	FACTS-R	40	74.46
	FACTS	40	80.85
	LDA	40	40.42
LCD TVs	CFACTS-R	80	78.94
	CFACTS	80	84.21
	FACTS-R	60	68.42
	FACTS	60	65.78
	LDA	40	36.84
Printers	CFACTS-R	100	79.31
	CFACTS	100	84.48
	FACTS-R	75	75.86
	FACTS	75	72.41
	LDA	75	36.76

Thanks !