





Jimmy Lin and Alek Kolcz

Twitter, Inc.



Image source:google.com/images



## Outline

# Outline

- •Is twitter big data?
- How can machine learning help twitter?
- Existing challenges?
- Existing literature of large-scale learning
- Overview of machine learning
- •Twitter analytic stack
- •Extending pig
- •Scalable machine learning
- Sentiment analysis application



## Focus of talk..

### What we will not talk about :

- Different "useful" application of twitter
- Why Twitter is a great product and one of its kind

### What we will talk about :

- Challenges faced while making it a good product
- Solution approach by "Insiders"



## Some twitter bragging ..

# The Scale of Twitter

- •Twitter has more than 280 million active users
- •500 million Tweets are sent per day
- •50 million people log into Twitter every day
- •Over 600 million monthly unique visitors to twitter.com

## Large scale infrastructure of information delivery

- •Users interact via web-ui, sms, and various apps
- •Over 70% of our active users are mobile users
- •Real-time redistribution of content
- At Twitter HQ we consume 1,440 hard boiled eggs weekly
- We also drink 585 gallons of coffee per week



Problems in hand ..

Support for user interaction

(other) problems we are trying to solve

- •Search
- -Relevance ranking
- User recommendation
- WTF or Who To Follow
- •Content recommendation
- -Relevant news, media, trends

- •Trending topics
- Language detection
- Anti-spam
- Revenue optimization
- User interest modeling
- •Growth optimization





## To put learning formally ..

# Supervised classification in a nutshell

Given  $D = \{ (x_i, y_i) \}_i^n$  (sparse) feature vector Induce  $f: X \to Y$  s.t. loss is minimized empirical loss =  $\frac{1}{n} \sum_{i=0}^n \ell(f(x_i), y_i)$  loss function

Consider functions of a parametric form:

$$\arg\min_{\theta} \frac{1}{n} \sum_{i=0}^{n} \ell(f(\mathbf{x}_{i}; \boldsymbol{\theta}), y_{i})$$

model parameters

Key insight: machine learning as an optimization problem! (closed form solutions generally not possible)



### Literature..

### Literature

•Traditionally, the machine learning community has assumed sequential algorithms on data fit in memory (which is no longer realistic)

•Few publication on machine learning work-flow and tool integration with data management platform

Google – adversarial advertisement detection Predictive analytic into traditional RDBMSes Facebook – business intelligence tasks LinkedIn – Hadoop based offline data processing But they are not for machine learning specificly. Spark ScalOps

But they result in end-to-end pipeline.



# What is author's contribution ..

## Contribution

- Provided an overview of Twitter's analytic stack
- Describe pig extension that allow seamless integration of machine learning capability into production platform
  Identify stochastic gradient descent and ensemble methods as being particularly amenable to large-scale machine learning

Note that, No fundamental contributions to machine learning



Scalable Machine Learning

# Scalable Machine learning

- Techniques for large-scale machine learning
- Stochastic gradient descent
- Ensemble method



## Gradient Descent..





ŵ

### Gradient Descent..

General method for nonlinear optimization

Start at  $\mathbf{w}(0)$ ; take a step along steepest slope

Fixed step size:

$$\mathbf{w}(1) = \mathbf{w}(0) + \eta$$

Move =

ght = Current Weight + mov = Step Size \* Unit Vecor

What is the direction  $\hat{\mathbf{v}}?$ 





Gradient Descent..

# Formula for the direction $\hat{\boldsymbol{v}}$

 $\Delta E_{\mathrm{in}} = E_{\mathrm{in}}(\mathbf{w}(0) + \eta \hat{\mathbf{v}}) - E_{\mathrm{in}}(\mathbf{w}(0))$ 

$$= \eta \nabla E_{\rm in}(\mathbf{w}(0))^{\rm T} \hat{\mathbf{v}} + O(\eta^2)$$

Using Taylor's series expansion

Because the surface non linear

$$\geq -\eta \| \nabla E_{\mathrm{in}}(\mathbf{w}(0)) \|$$

Since  $\hat{\mathbf{v}}$  is a unit vector,

$$\hat{\mathbf{v}} = -\frac{\nabla E_{\mathrm{in}}(\mathbf{w}(0))}{\|\nabla E_{\mathrm{in}}(\mathbf{w}(0))\|} \xrightarrow{\mathrm{Descent along gradies}} for the product of error.}$$



Stochastic Gradient Descent (SGD)

sto-chas-tic stə kastik/ *adjective* 1.randomly determined; having a random probability distribution or pattern that may be analyzed statistically but may not be predicted precisely.

Stochastic gradient descent

GD minimizes:

$$E_{\rm in}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^{N} \underbrace{\mathbf{e}\left(h(\mathbf{x}_n), y_n\right)}_{\ln\left(1+e^{-y_n \mathbf{w}^{\mathsf{T}} \mathbf{x}_n\right)} \leftarrow \text{ in logistic regression}}$$

by iterative steps along  $abla E_{
m in}$ :

$$\Delta \mathbf{w} = - \eta \ \nabla E_{\mathrm{in}}(\mathbf{w})$$

 $abla E_{ ext{in}}$  is based on all examples  $(\mathbf{x}_n, y_n)$ 

Slides from Yaser Abu Mostafa-Caltech

"batch" GD



Stochastic Gradient Descent (SGD)

# Stochastic gradient descent

# **Gradient Descent**

$$w^{(t+1)} = w^{(t)} + \gamma^{(t)} \frac{1}{n} \sum_{i=0}^{n} \nabla I\left(f(\mathbf{x}_{i}; \theta^{(t)}), y_{i}\right)$$

Compute the gradient in the loss function by optimizing value in dataset. This method will do the iteration for all the data in order to one a gradient value.

Inefficient and everything in the dataset must be considered.



Stochastic Gradient Descent (SGD)

# Stochastic gradient descent

Approximating gradient depends on the value of gradient for one instance.

$$w^{(t+1)} = w^{(t)} + \gamma^{(t)} \nabla I\left(f(\mathbf{x}; \boldsymbol{\theta}^{(t)}), y\right)$$

Solve the iteration problem and it does not need to go over the whole dataset again and again.

Stream the dataset through a single reduce even with limited memory resource.

But when a huge dataset stream goes through a single node in cluster, it will cause network congestion problem.



### Stochastic Gradient Descent (SGD)

# Benefits of SGD

1. cheaper computation  $E_{
m in}$ 2. randomization 3. simple Weights, w Rule of thumb: randomization helps A 4



Aggregation a.k.a Ensemble Learning

What is aggregation?

Combining different solutions  $h_1, h_2, \cdots, h_T$  that were trained on  $\mathcal{D}$ :



Regression: take an average

Classification: take a vote

a.k.a. ensemble learning and boosting



## Aggregation a.k.a Ensemble Learning

# Different from 2-layer learning

In a 2-layer model, all units learn jointly:



In aggregation, they learn **independently** then get combined:





Ensemble Learning..

# **Ensemble Methods**

Classifier ensembles: high performance learner

Performance: very well

Some rely mostly on randomization

-Each learner is trained over a subset of features and/or instances of the data

- Ensembles of linear classifiers
- Ensembles of decision trees (random forest)







Hoeffding's Inequality

In a big sample (large N), u is probably close to  $\mu$  (within  $\epsilon$ ).

Formally,

 $\mathbb{P}\left[\left|\nu-\mu\right| > \epsilon\right] \le 2e^{-2\epsilon^2 N}$ 

Sample frequency v is likely lose to bin frequency  $\mu$ .

This is called **Hoeffding's Inequality**.

Slide taken from Caltech's Learning from Data Course : Dr Yaser Abu Mostafa





Image Source: Apache Yarn Release



Hadoop Ecosystem at Twitter..





Glorifying PIG

# Why use Pig?

 Suppose you have user data in one file, website data in another, and you need to find the top 5 most visited sites by users aged 18 - 25







### **Glorifying PIG**

# In Map-Reduce

ern abis arregtist The second terre are seare rating meret. Heartablemete. org-sparse-salong-sepred-solorosfi org-sparse-rationg-sepred-sepremeters/sepre-org-sparse-rationg-sepred-sepremeters 19112 118-18-1912 - 20102 - 20122 - 2010 - 2010 Alexan and a second sec name with styles ( the static state traditupes attacks residented as inplants regardingsridented. Tark. Tark. Tark ale shekir ekses tratavitikkarosare atkante mojaako Sepananka mojaretemperikatas, tark, tark, tark, tark And The Association of the Assoc iteing adva = test testating;;; the provide a star testating;; dot app = entropy - person testating; dot app = entropy - person testating; dot app = entropy - person testating; testa entropy = person testation; testa entropy testation; testa entropy testation; testation; ingustaria makes that attants reptationalizes patale wild enters (werk hay, starskerstuck flar, indepatient allow, tarks of allow a superior reporter, being allow the late term and in patent want, tipper out will the late term and shore it units (dier.mannet.)) ( tent t = ther.mat()) direct and the thermality () direct and the tent of tent o

agentian contraction and provide the the second product and statement the vester "DELET, affit the form of the second seco rapparenters, turk, turk, Longershaltan ( would mapp ten ber and been under wennen bererart. Der gestikteter of Begerker regerter bereis einen geben Begerker regerter bereis einen geben ing and the same hadrow the other that the burght and automa - au will char manatin () ( saarhong, max temperakataacaamoor sections. Longweiterne states theiltistes attants repartmeters enert a te le mais catera: ingentitata tar. Marthermath far. Marthermath far. Marthermath far. ( imay ordered free first and restricts ( ) of the state ais skatis min min minerativer () args; berns remeaghin ( scheme, b) o ma entrance () tags, () scheme, b) of mark tags, ()

Martin Control (Sector) (Secto CHARLES PROPERTY AND THE PROPERTY PROPERTY IN pakan / maara \* 1) / Independent mark andre Appel ander ja fra. Naar talek / / maar / galana, felg / flak anal, saarar \* 1) andre mark / / maar / galana, felg / flak anal, saarar \* 1) Andre mark - 1 haar fing / fing / fing / Full Australia in the finite Apply and and the state of the finite apply and the state of the state Solar - state apply in the state of the state of the state Solar - state apply in the state of the state of the state Solar - state apply in the state of the state of the state Solar - state apply in the state of the state of the state Solar - state apply in the state of the state of the state state of the state of the state of the state of the state state of the state of the state of the state of the state state of the state of the state of the state of the state state of the state of the state of the state of the state state of the state of the state of the state of the state state of the state of the state of the state of the state state of the state of the state of the state of the state state of the state of the state of the state of the state state of the state of the state of the state of the state state of the state of conserved prompt a feature of the second sec Plant and proved a simple rate ( prop. And Hand berger and an and the second second second Rin friter anter anter attagenter anter attagenter atta at hopens - new solo ter attan 11 trant addition in a star of the start te jaar andere beeren an en in ferder andere sterre beeren de seekere en de seekere besterre be bigett arter tel ber her riefte their atter Million and a second seco interfaces in a new minterfaces; what hep are allow for your 

#### 170 lines of code, 4 hours to write Credits : Hortonworks



## Glorifying PIG

# In Pig Latin

```
Users = load 'input/users' using PigStorage(',') as (name:chararray, age:int);
Fltrd = filter Users by age >= 18 and age <= 25;
Pages = load 'input/pages' using PigStorage(',') as (user:chararray,
url:chararray);
Jnd = join Fltrd by name, Pages by user;
Grpd = group Jnd by url;
Smmd = foreach Grpd generate group,COUNT(Jnd) as clicks;
Srtd = order Smmd by clicks desc;
Top5 = limit Srtd 5;
store Top5 into 'output/top5sites' using PigStorage(',');
```

#### 9 lines of code, 15 minutes to write

#### 170 lines to 9 lines of code

Credits : Hortonworks



2

### Maximizing the use of Hadoop ..

Maximizing the use of Hadoop

- •We cannot afford too many diverse computing environments
- Most of analytics job are run using Hadoop cluster
- -Hence, that's where the data live
- -It is natural to structure ML computation so that it takes advantage of the cluster and is performed close to the data



Integration into production workflows





What authors contributed technically ..

# **Core libraries:**

## **Core Java library**

Basic abstractions similar to existing packages (weka, mallet, mahout)

# Lightweight wrapper

Expose functionalities in Pig



**PIG Functions..** 

**Training models:** 



**PIG Functions..** 

Shuffling data:



**PIG Functions..** 

Using models:



My Scripts       Overy history         My scripts       Overy history         My scripts       Overy history         Pig script       Pig script         Settings       Ibatting = load 'Batting.csv' using PigStorage(','); 2 mons = FOREACH batting GENRATE \$9 as playerID, \$1 as year, \$9 as max_runs; 3 groupdate = GROUP mus by (year); 4 max_runs = FOREACH group as group, MAX(runs.runs) as max_runs; 5 join max_run = DOLM max_runs by (\$9, max_runs), runs by (year, runs); 7 join_data = FOREACH gin_max_run GENERATE \$9 as year, \$2 as playerID, \$1 as runs; 8 damp join_data;	÷		HortonWo	orks Way
My scripts         Image: DigScript         Settings         Image: DigScript         Image: DigScript: DigScript         Image: DigScript: DigScr	N 🐝 🐱 🚾 📃 🔍 🗭 My Scripts Query history	You have gone full screen.     Exit full screer	<u>n (F11)</u>	å hue ▼
	My scripts I pigScript Settings Email notification USER-DEFINED FUNCTIONS V Upload UDF Jar	<pre>Title: pigScript Pig script:  Pig script:  Pig script:  Pig script:  Pig helper  to batting = load 'Batting.csv' using PigStorage(','); runs = FOREACH batting GENERATE \$0 as playerID, \$1 a grp_data = GROUP runs by (year); max_runs = FOREACH grp_data GENERATE group as grp,MA  join_max_run = JOIN max_runs by (\$0, max_runs), runs join_data = FOREACH join_max_run GENERATE \$0 as year dump join_data; </pre>	as year, \$8 as runs; {X(runs.runs) as max_runs; ; by (year,runs); , \$2 as playerID, \$1 as runs;	



Final Model which works!!!

# Final Learning - Ensemble Methods





Use case..

# Example: Sentiment Analysis

Emotion Trick  $\odot$   $\otimes$ 

Test dataset: 1 million English tweets, minimum 20 letters-long

Training data: 1 million, 10 million and 100 million English training examples

Preparation: training and test sets contains equal number of positive and negative examples, removed all emoticons.



# Finally a graph ..





# Explaining a bit more of graph ..



- 1. The error bar denotes 95% confidence interval
- 2. The leftmost group of bars show accuracy when training a single logistic regression classifier on {1, 10, 100} million training examples.
- 3. 1-10 Change Sharp , 10 100 million : Not that sharp
- 4. The middle and right group of bars in Figure 2 show the results of learning ensembles
- 5. Ensembles lead to higher accuracy—and note that an ensemble trained with 10 million examples outperforms a single classifier trained on 100 million examples
- 6. No accurate running time reported as experiments were run on production clusters but informal observations are in sync with what the logical mind suggests ( ensemble takes shorter to train because models are learned in parallel )
- 7. In terms of applying the learned models, running time increases with the size of the ensembles—since an ensemble of n classifiers requires making n separate predictions.



Conclusion

What I loved about paper : I understood it  $\odot$  ?

"our goal has never been to make fundamental contributions to machine learning, we have taken the pragmatic approach of using off-the shelf toolkits where possible. Thus, the challenge becomes how to incorporate third-party software packages along with inhouse tools into an existing workflow"..







