Big data management and predictive analytics for customer transactions & operations using Apache Spark and AWS

Large Scale Distributed Machine Learning

> BDA761 Big Data Mgmt in a Supercomputing environment

> > MUDIT UPPAL

# Problem with Big Data(s)

- Machine learning practices at scale for PB/TB data
- \* A framework which provides and computes models using virtual nodes with processors and memory getting cheaper every year
- \* Using GPU + multi-threading + make use of multiple cores
- Goal: Thinking in 'big data'; create a tool which can be used in any operations/Sales/customer analysis
- \* Parallelizing DATA and MODEL

## Case study

- \* *Rossmann*: operates over 3000 drug stores in 7 european countries
- \* Store sales are influenced by many factors, including promotions, competition, school and state holidays, seasonality, and locality.
- \* DECISION TASKS:
  - Forecast Sales for upto 6 weeks using stores data, customers, promotion data et cetera
  - Predicting Sales of ~1000 stores daily

Data Source: Rossmann, Walmart via Kagle

														TripType Visi	tNumber Weekday	Upc S	ScanCount	DepartmentDescription	FinelineNumber
Id	Store	DayOfWeek	Date	Open	Promo	State	Holiday	SchoolH	oliday					999	5 Friday	68113152929	-1	FINANCIAL SERVICES	1000
1	1	4	2015-09-17	1	1	0		0						30	7 Friday	60538815980	1	SHOES	8931
2	з	4	2015-09-17	1	1	0 Stor	re StoreType	Assortment	CompetitionDistance	CompetitionOpenSinceMonth	CompetitionOpenSinceYear	Promo2	Promo2SinceWeek	Promo2SinceYear	PromoInterval	7410811099	1	PERSONAL CARE	4504
-	5		2015 05 17	-	-	1	с	а	1270	9	2008	0				2238403510	2	PAINT AND ACCESSORIES	3565
3	7	4	2015-09-17	1	1	0 2	а	а	570	11	2007	1	13	2010	Jan,Apr,Jul,Oct	2006613744	2	PAINT AND ACCESSORIES	1017
4	8	4	2015-09-17	1	1	0 3	а	а	14130	12	2006	1	14	2011	Jan, Apr, Jul, Oct	2006618783	2	PAINT AND ACCESSORIES	1017
5	Q	4	2015-09-17	1	1	0 5	c	с э	29910	4	2009	0				2006613743	1	PAINT AND ACCESSORIES	1017
5	5	-	2015 05 17	-	-	6	a	a	310	12	2013	0				7004802737	1	PAINT AND ACCESSORIES	2802
6	10	4	2015-09-17	1	1	0 7	а	с	24000	4	2013	0				2238495318	1	PAINT AND ACCESSORIES	4501
7	11	4	2015-09-17	1	1	08	а	а	7520	10	2014	0				2238400200	-1	PAINT AND ACCESSORIES	3565
8	12	4	2015-09-17	1	1	0 9	а	c	2030	8	2000	0				5200010239	1	DSD GROCERY	4606
0	12		2013-09-17	1	1	10	а	а	3160	9	2009	0				38679300501	2	PAINT AND ACCESSORIES	3504
9	13	4	2015-09-17	1	1	0 11	а	c	960	11	2011	1	1	2012	Jan,Apr,Jul,Oct	22006000000	1	MEAT - FRESH & FROZEN	6009
10	14	4	2015-09-17	1	1	0 12	a	c	1070			1	13	2010	Jan, Apr, Jul, Oct	2236760452	1	PAINT AND ACCESSORIES	7
11	15	4	2015-00-17	4	1	0 14	a	a	1300	3	2014	1	40	2003	Jan.Apr.Jul.Oct	38679300501	-1	PAINT AND ACCESSORIES	3504
	15	4	2013-09-17	1	1	15	d	c	4110	3	2010	1	14	2011	Jan, Apr, Jul, Oct	2238400200	2	PAINT AND ACCESSORIES	3565
12	16	4	2015-09-17	1	1	0 16	а	с	3270			0				3019294203	1	PAINT AND ACCESSORIES	2801
13	19	4	2015-09-17	1	1	0 17	а	а	50	12	2005	1	26	2010	Jan,Apr,Jul,Oct	72450408840	1	PAINT AND ACCESSORIES	1028
14	20		2015 00 17			18	d	с	13840	6	2010	1	14	2012	Jan, Apr, Jul, Oct	25541500000	2	DAIBY	1305
14	20	4	2015-09-17	1	1	U 19	а	с	3240			1	22	2011	Mar,Jun,Sept,Dec	2310010776	- 1	PETS AND SUPPLIES	3300
15	21	4	2015-09-17	1	1	0 20	d	а	2340	5	2009	1	40	2014	Jan, Apr, Jul, Oct	72450402700	· 2		1019
16	22	4	2015-09-17	1	1	0 22	c	c	550	10	1999	1	45	2009	Jan,Apr,Jul,Oct	72430403700			707
10		•	2015 05 17	-	-	22	a d	a	4060	8	2005	0	22	2012	Jan, Apr, Jul, Oct	7874204967			707
17	23	4	2015-09-17	1	1	24	a	c	4590	3	2000	1	40	2011	Jan, Apr, Jul, Oct	5114139038	1	PAINT AND ACCESSORIES	4415
						25	c	а	430	4	2003	0				5114197561	1	PAINT AND ACCESSORIES	4415
~2.8 million rows data					26	d	а	2300			0				3270011053	3	PETS AND SUPPLIES	1001	
					27	а	а	60	1	2005	1	5	2011	Jan, Apr, Jul, Oct		1	NULL		
					28	а	а	1200	10	2014	1	6	2015	Mar,Jun,Sept,Dec					
					29	d	с	2170			0								
			-			30	а	а	40	2	2014	1	10	2014	Mar,Jun,Sept,Dec				
points						31	d	с	9800	7	2012	0							
1																			

Label: Sales

Features: Store, Sales, customers, Open, date, stateholiday, schoolHoliday, storetype, Assortment, competitionDistance, CompetitionSinceMonth, sinceYear, Promo, Promo2, PromoeInterval, Promo2SinceMonth, PromoSince year, DayOfWeek, etc...

Store	StoreType	Assortment		Competiti	onDistance	CompetitionOpenSinceMonth			CompetitionOpenSinceYear		Promo2	Promo2SinceWeek	Promo2SinceYear	PromoInterval
1	с	а			1270	9			2008		0			
2	а	a		570		11			2007		1	13	2010	Jan,Apr,Jul,Oct
3	а	a			14130			12	2006		1	14	2011	Jan,Apr,Jul,Oct
Store	DavOfWee	k Date		Sales	Customers	Open	Promo	StateHoliday	SchoolHoliday					
					• • • • • • • •	••••		••••••						
1		5 2015-0	7-31	5263	555	1	1	0	1					
2		5 2015-0	7-31	6064	625	1	1	0	1					
3		5 2015-07-3		8314	821	1	1	0	1					

~

.

.....

.

### Goal

- \* Aim is to create a Predictive Analytics Framework
  - distributed machine learning for any size dataset
- 3 Demos in this presentation
  - \* Apache Spark
  - \* Exploratory DA with R
  - \* xgBoost(python) ML

### Data Source: Rossmann, Walmart via Kagle

### Main Demos

- Exploratory data analysis
- Apache Spark(sql spark context) Distributed ML across multiple machines/nodes
  - \* Linear Regression analysis
  - Gradient descent
- Ensemble and boosting algorithms (XGBoost)

### Tools

- R + python for exploratory analysis
- Apache Spark for implementation
- \* Hadoop
- Spark MLlib
- AWS cloud(m4 large instances)
- Ganglia(distributed monitoring system ti work with clusters)
- pyspark

# Data Pipeline



### Demo 1: exploratory analysis

- R+python
- Data and source code can be downloaded from: <u>http://muppal.com</u>





factor(DayOfWeek)

# ApacheSpark <---> xgboost

- Spark excels at distributing operations across a cluster while abstracting away many of the underlying implementation details.
- Thinking in terms of RDD(transformations and actions)
- Still under development

#### XGBoost:

- start off with a rough prediction and then building a series of decision trees; with each trying to correct the prediction error of the one before
- Large-scale and Distributed Gradient Boosting (GBDT, GBRT or GBM) Library, on single node, hadoop yarn etc
- XGBoost can also be distributed and scale to Terascale data
- You can define threads with "nthreads =.."
- https://github.com/dmlc/xgboost

$$\begin{bmatrix} 9 & 3 & 5 \\ 4 & 1 & 2 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 3 & -5 \\ 2 & 3 \end{bmatrix} = \begin{bmatrix} 28 & 18 \\ 11 & 9 \end{bmatrix}$$
$$\begin{bmatrix} 9 & 18 \\ 4 & 8 \end{bmatrix} + \begin{bmatrix} 9 & -15 \\ 3 & -5 \end{bmatrix} + \begin{bmatrix} 10 & 15 \\ 4 & 6 \end{bmatrix}$$



Example: n = 6; 3 workers



# Data pipeline/Tools



### Demo 2/3: xgboost + ML Spark on clusters

- xgboost + Apache Spark/pyspark
- \* Findings: [1199] train-rmspe:0.103954 eval-rmspe:0.094526
- Validating RMSPE: 0.094526
- \* Data and source code can be downloaded from: <u>http://muppal.com</u>

[TTOT]			Tiassie Free Data Science Apps Committation - Dear Mudit,
[1182]	train-rmspe:0.104195	eval-rmspe:0.094690	
[1183]	train-rmspe:0.104192	eval-rmspe:0.094686	
[1184]	train-rmspe:0.104186	eval-rmspe:0.094678	
[1185]	train-rmspe:0.104167	eval-rmspe:0.094676	
[1186]	train-rmspe:0.104151	eval-rmspe:0.094668	
[1187]	train-rmspe:0.104136	eval-rmspe:0.094649	
[1188]	train-rmspe:0.104137	eval-rmspe:0.094649	
[1189]	train-rmspe:0.104120	eval-rmspe:0.094639	
[1190]	train-rmspe:0.104109	eval-rmspe:0.094637	
[1191]	train-rmspe:0.104107	eval-rmspe:0.094637	
[1192]	train-rmspe:0.104065	eval-rmspe:0.094622	
[1193]	train-rmspe:0.104015	eval-rmspe:0.094553	
[1194]	train-rmspe:0.104018	eval-rmspe:0.094555	
[1195]	train-rmspe:0.104014	eval-rmspe:0.094551	
[1196]	train-rmspe:0.103998	eval-rmspe:0.094550	
[1197]	train-rmspe:0.103995	eval-rmspe:0.094550	
[1198]	train-rmspe:0.103958	eval-rmspe:0.094525	
[1199]	train-rmspe:0.103954	eval-rmspe:0.094526	
Validat	ting ou were in a vider		
RMSPE:	0.094526		

Thank you!

### References

- http://rcarneva.github.io/understanding-gradientboosting-part-1.html
- https://www.kaggle.com/c/walmart-recruiting-triptype-classification