

10/04/16 CS535 Big Data - Fall 2016 W7.A.1

CS535 BIG DATA

PART 1. BATCH COMPUTING MODELS FOR BIG DATA ANALYTICS  
**2. LARGE SCALE DATA ANALYSIS USING SPARK WITH CASE STUDY**

Sangmi Lee Pallickara  
Computer Science, Colorado State University  
<http://www.cs.colostate.edu/~cs535>

10/04/16 CS535 Big Data - Fall 2016 W7.A.1

FAQs

- PA1 demo is going on
- Term Project
  - 10/7 5:00PM
  - Description available on the course web
  - Presentation schedule
  - Google computing cluster credit is available

10/04/16 CS535 Big Data - Fall 2016 W7.A.2

Proposal

- Title of your project
- Problem formulation
- Your strategy to solve the problem
- Functions targeted by your software
- Plan for testing
- Evaluation method
- Project timeline (weekly plan)
- Bibliography
- Submission (1,200-1,800 words)

10/04/16 CS535 Big Data - Fall 2016 W7.A.3

Presentations (10/11, and 10/13)

- Slide 1: Title (with the team info)
- Slide 2: Problem statement
- Slide 3: Your approach
- Slide 4: Your software
- Slide 5: Plan for software testing
- Slide 6: Evaluation Method
- Presentation should be no longer than **12 minutes** including the Q&A session. (10 minutes: presentation, 2 minutes: Q&A)
- All of the team members should present
- Audience will get 2% of participation score based on their questions and attendance
- Please send me your slides **at least 2 hours** before your presentation

10/04/16 CS535 Big Data - Fall 2016 W7.A.4

Objectives

- Recommendation systems
  - Collaborative filtering
  - Latent factor approach
- Linear Methods

10/04/16 CS535 Big Data - Fall 2016 W7.A.5

Large scale data analysis using Spark  
CASE STUDY: Recommending Music and  
the Audioscrobbler Dataset  
**Evaluating the Recommendation Model**

10/04/16 CS535 Big Data - Fall 2016 W7.A.6

## What is a “good” recommendation?

- “a popular artist”?
- “artists the user has listened to”?
- “artists the user will listen to”?

10/04/16 CS535 Big Data - Fall 2016 W7.A.7

## Preparing data for evaluation

- To perform a meaningful evaluation, some of the artist play data can be set aside
  - Hidden from the ALS model building process
- The held-out data can be used as a collection of good recommendations for each user
- Compute the recommender's score

For building model

For testing model

10/04/16 CS535 Big Data - Fall 2016 W7.A.8

## AUC metric

- Rank 1.0 is perfect, 0.0 is the worst
- Receiver Operating Characteristic (ROC)
  - Based on the rank used to decide final recommendations
- Area Under the Curve (AUC) of ROC may be used as the probability that a randomly chosen good recommendation ranks above a randomly chosen bad recommendation
- Spark's `BinaryClassificationMetrics`
  - Computes AUC per users and averages the result
  - Generating mean AUC

10/04/16 CS535 Big Data - Fall 2016 W7.A.9

## MAP metric

- Mean average precision
  - Focuses on the top recommendations

10/04/16 CS535 Big Data - Fall 2016 W7.A.10

## Computing AUC

- 90% of the data is used for training and the remaining 10% for validation

```
import org.apache.spark.rdd._

def areaUnderCurve(
  positiveData: RDD[ Rating],
  bAllItemIDs: Broadcast[ Array[ Int]],
  predictFunction: (RDD[( Int, Int)] => RDD[Rating])) = {
  ...
}

val allData = buildRatings( rawUserArtistData, bArtistAlias)
val Array( trainData, cvData) =
  allData.randomSplit(Array( 0.9, 0.1))
```

10/04/16 CS535 Big Data - Fall 2016 W7.A.11

## Computing AUC

- continued

```
trainData.cache()
cvData.cache()

val allItemIDs = allData.map(_._2.product).distinct().collect()

val bAllItemIDs = sc.broadcast( allItemIDs)

val model = ALS.trainImplicit( trainData, 10, 5, 0.01, 1.0)

val auc = areaUnderCurve( cvData, bAllItemIDs, model.predict)
```

10/04/16 CS535 Big Data - Fall 2016 W7.A.12

## k-Fold Cross-validation

- Create a  $k$ -fold partition of the dataset
  - For each of the  $k$  experiments use  $K-1$  folds for training
    - The remaining fold for testing

Test example

Experiment 1  
Experiment 2  
Experiment 3  
Experiment 4

← Total number of examples →

10/04/16 CS535 Big Data - Fall 2016 W7.A.13

## True error estimate

- $k$ -fold cross validation is similar to random subsampling
- The advantage of  $k$ -Fold Cross validation
  - All the examples in the dataset are eventually used for both training and testing
- The true error is estimated as the average error rate

$$E = \frac{1}{K} \sum_{i=1}^K E_i$$

10/04/16 CS535 Big Data - Fall 2016 W7.A.14

## k-Fold Cross-validation with Spark

```

MLUtils.kFold()

def predictMostListened(
  sc: SparkContext,
  train: RDD[Rating])(allData:
  RDD[(Int, Int)]) = {

  val bListenCount = sc.broadcast(
    train.map(r => (r.product, r.rating)).
    reduceByKey(_ + _).collectAsMap()
  )

  allData.map { case (user, product) =>
    Rating(
      user, product,
      bListenCount.value.getOrElse(product, 0.0) )
  }
}

val auc = areaUnderCurve(cvData, bAllItemIDs,
  predictMostListened(sc, trainData))

```

10/04/16 CS535 Big Data - Fall 2016 W7.A.15

## Hyperparameter selection

- MatrixFactorizationModel
- ALS.trainImplicit()
  - rank = 10
    - The number of latent factors in the model
    - The number of columns,  $k$
  - iterations = 5
    - The number of iterations that the factorization runs
  - lambda = 0.1
    - A standard overfitting parameter
    - Higher value guards against overfitting
    - Values that are too high will decrease the factorization's accuracy
  - alpha = 1.0
    - Controls the relative weight of observed versus unobserved user-product interactions in the factorization

10/04/16 CS535 Big Data - Fall 2016 W7.A.16

Large scale data analysis using Spark  
CASE STUDY: Recommendation Systems  
**Amazon.com : Item-to-item collaborative filtering**

10/04/16 CS535 Big Data - Fall 2016 W7.A.17

## This material is built based on,

- Greg Linden, Brent Smith, and Jeremy York, "Amazon.com Recommendations, Item-to-Item Collaborative Filtering" IEEE Internet Computing, 2003

10/04/16 CS535 Big Data - Fall 2016 W7.A.18

- Amazon.com uses recommendations as a targeted marketing tool
  - Email campaigns
  - Most of their web pages

10/04/16 CS535 Big Data - Fall 2016 W7.A.19

### Item-to-item collaborative filtering

- It does NOT match the user to similar customers
- Item-to-item collaborative filtering
  - Matches each of the user's purchased and rated items to similar items
  - Combines those similar items into a recommendation list

10/04/16 CS535 Big Data - Fall 2016 W7.A.20

### Determining the most-similar match

- The algorithm builds a similar-items table
  - By finding items that customers tend to purchase together
- How about building a product-to-product matrix by iterating through all item pairs and computing a similarity metric for each pair?
- Many product pairs have no common customer
  - If you already bought a TV today, will you buy another TV again today?

10/04/16 CS535 Big Data - Fall 2016 W7.A.21

- Calculating the similarity between a single product and all related products:

```
For each item in product catalog, I1
  For each customer C who purchased I1
    For each item I2 purchased by customer C
      Record that a customer purchased I1 and I2
    For each item I2
      Compute the similarity between I1 and I2
```

10/04/16 CS535 Big Data - Fall 2016 W7.A.22

### Computing similarity

- Using cosine measure
  - Each vector corresponds to an item rather than a customer
  - M dimensions correspond to customers who have purchased that item

10/04/16 CS535 Big Data - Fall 2016 W7.A.23

### Creating a similar-item table

- Similar-items table is extremely computing intensive
  - Offline computation
  - $O(N^2M)$  in the worst case
    - Where N is the number of items and M is the number of users
  - Average case is closer to  $O(NM)$ 
    - Most customers have very few purchases
  - Sampling customers who purchase best-selling titles reduces runtime even more
    - With little reduction in quality

10/04/16 CS535 Big Data - Fall 2016 W7.A.24

## Scalability

- Amazon.com has around 110 million active customers(244 million total customers) and several million catalog items
- Traditional collaborative filtering does little or no offline computation
- Online computation scales with the number of customers and catalog items.

<http://www.fool.com/investing/general/2014/05/24/how-many-customers-does-amazon-have.aspx>

10/04/16 CS535 Big Data - Fall 2016 W7.A.25

## Key scalability strategy for amazon recommendations

- Creating the expensive similar-items table offline
- Online component
  - Looking up similar items for the user's purchases and ratings
  - Scales independently of the catalog size or the total number of customers
- It is dependent only on how many titles the user has purchased or rated

10/04/16 CS535 Big Data - Fall 2016 W7.A.26

## Recommendation quality

- The algorithm recommends highly correlated similar items
  - Recommendation quality is excellent
  - Algorithm performs well with limited user data

10/04/16 CS535 Big Data - Fall 2016 W7.A.27

## Large scale data analysis using Spark

### CASE STUDY: Linear Models

10/04/16 CS535 Big Data - Fall 2016 W7.A.28

## Spark's linear models

- Classification
  - Linear Support Vector Machines (SVMs)
  - Logistic regression
- Regression
  - Linear least squares, Lasso, and ridge regression

10/04/16 CS535 Big Data - Fall 2016 W7.A.29

## Implementation

- SVMWithSGD
- LogisticRegressionWithBFGS
- LogisticRegressionWithSGD
- LinearRegressionWithSGD
- RidgeRegressionWithSGD
- LassoWithSGD

10/04/16 CS535 Big Data - Fall 2016 W7.A.30

## Large scale data analysis using Spark

### CASE STUDY: Linear Models

## Linear Regression Model to a Large Dataset

10/04/16 CS535 Big Data - Fall 2016 W7.A.31

## The linear regression model

- The structure is exactly the same as for the linear discriminant function  

$$f(x) = w_0 + w_1x_1 + w_2x_2 + \dots$$

	A	B	C	D	E	F	G	H
Math score	85	85	62	67	89	65	80	67
Physics score	89	92	70	95	95	80	75	80

- How big is the error of the fitted model?
  - We would like to minimize this error
- The model that fits the data best
  - The model with the **minimum sum of errors** on the training data
  - e.g. The sum or mean of the squares of the errors
  - Least squares regression

10/04/16 CS535 Big Data - Fall 2016 W7.A.32

## Squared error

- Squared error
  - Strongly penalizes very large errors
  - Drawback
    - Is very **sensitive to the data**
    - Erroneous or outlying data points can severely skew the resultant linear function
- We should choose the objective function to optimize

10/04/16 CS535 Big Data - Fall 2016 W7.A.33

## Root Mean Squared Error

- Measures the differences between values predicted by model estimator and the values actually observed

$$RMSD = \sqrt{\frac{\sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2}{n}}$$

10/04/16 CS535 Big Data - Fall 2016 W7.A.34

## Linear Regression

- Calculating a linear regression
  - Using the least square criterion

	A	B	C	D	E	F	G	H
Math score	85	85	62	67	89	65	80	67
Physics score	89	92	70	95	95	80	75	80

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1$$

where, the  $\theta$ s are the **parameters** (Math scores, or slopes)

10/04/16 CS535 Big Data - Fall 2016 W7.A.35

## Linear Regression -- continued

$$h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \dots$$

- To simplify the notation,

$$h(x) = \sum_{i=0}^n \theta_i x_i$$

10/04/16 CS535 Big Data - Fall 2016 W7.A.36

## Objective function (Cost function)

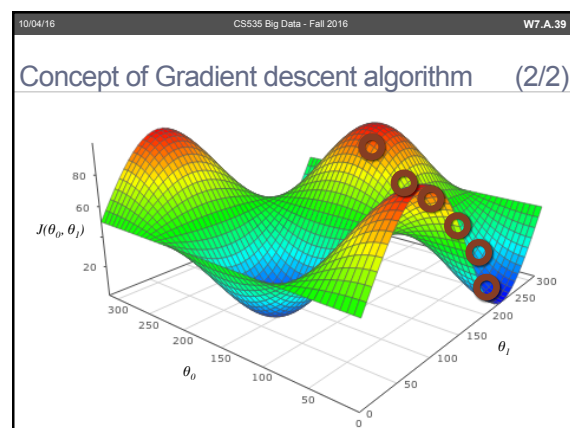
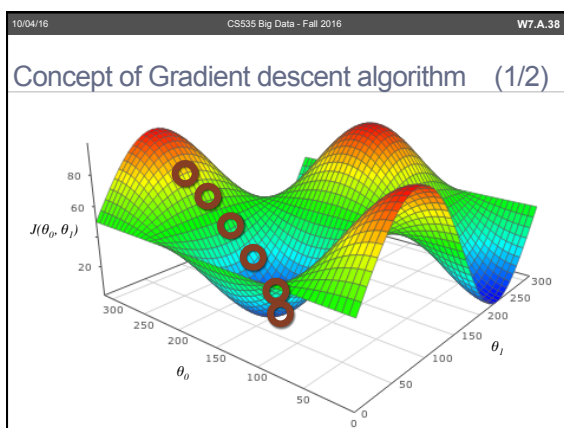
- For a given training set, how do we pick, or learn, the parameter  $\theta$ ?
- Make  $h(x)$  close to  $y$ 
  - Make your prediction close to the real observation
- We define the **objective (cost) function**

$$J(\theta) = \frac{1}{2} \sum_{i=0}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

10/04/16 CS535 Big Data - Fall 2016 W7.A.37

## Minimization problem

- We have a function  $J(\theta_0, \theta_1)$
- We want to find  $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$
- Goal: Find parameters to minimize the cost (output of the objective function)**
- Outline of our approach:
  - Start with some  $\theta_0, \theta_1$
  - Keep changing  $\theta_0, \theta_1$  to reduce  $J(\theta_0, \theta_1)$  until we end up at a minimum



10/04/16 CS535 Big Data - Fall 2016 W7.A.40

## Stochastic Gradient descent algorithm

Repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

(for  $j=0$  and  $j=1$ )

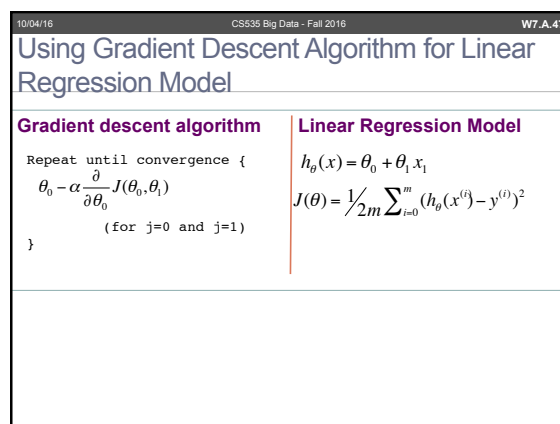
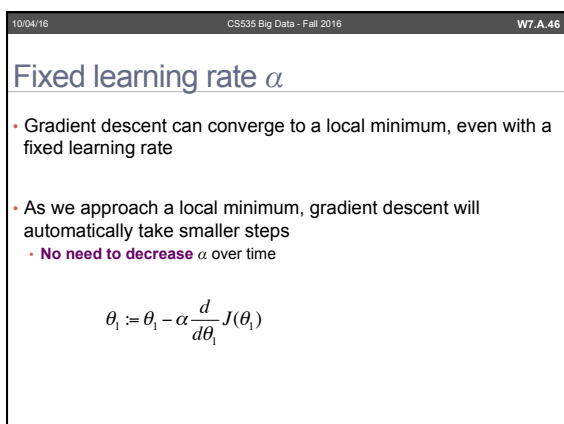
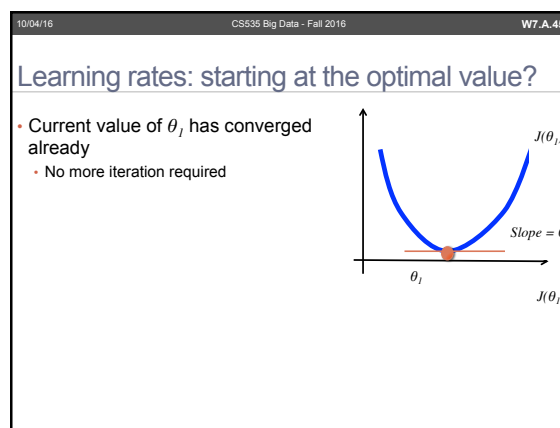
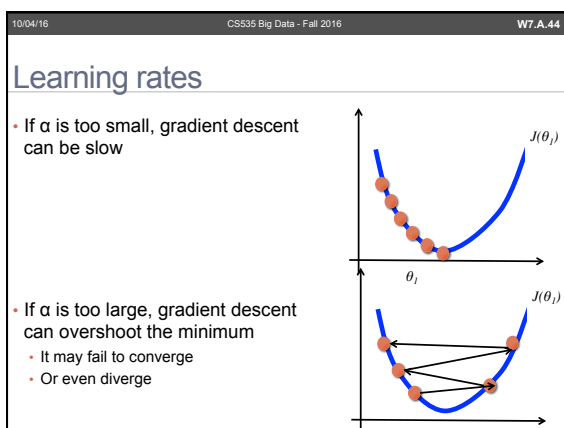
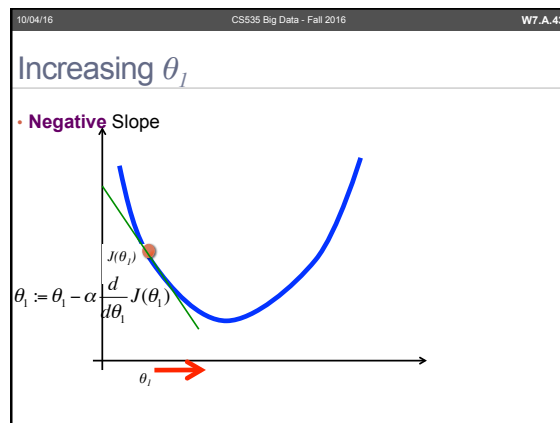
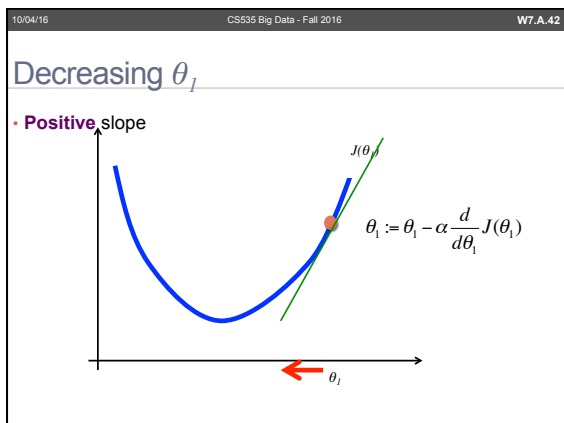
}

- Simultaneous update**
  - Your implementation should perform simultaneous update
  - See slide 41

10/04/16 CS535 Big Data - Fall 2016 W7.A.41

## Simultaneous update

Correct: Simultaneous update	Incorrect:
$\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$	$\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$
$\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$	$\theta_0 := \text{temp0}$
$\theta_0 := \text{temp0}$	$\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$
$\theta_1 := \text{temp1}$	$\theta_1 := \text{temp1}$





10/04/16 CS535 Big Data - Fall 2016 W7.A.48

$$\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$= \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2$$

Case 1,  $\theta_0$  ( $j = 0$ ) :

$$\frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_0} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})$$

Case 2,  $\theta_1$  ( $j = 1$ ) :

$$\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_1} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)}$$

10/04/16 CS535 Big Data - Fall 2016 W7.A.49

### Gradient descent for Linear Regression

Repeat until convergence {

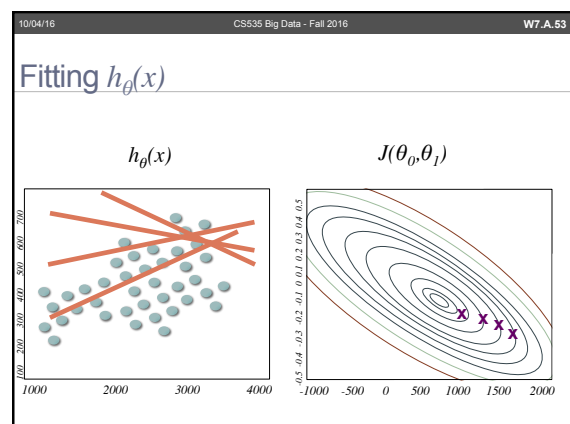
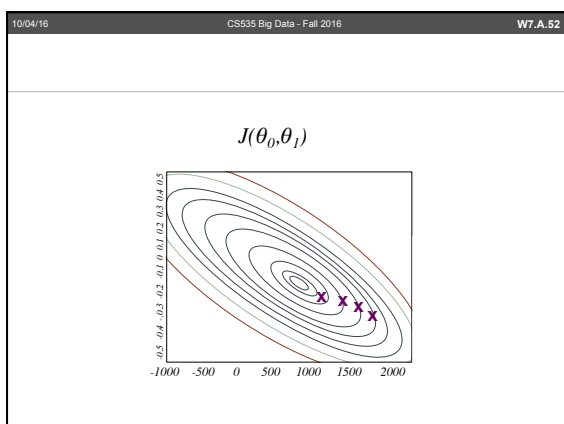
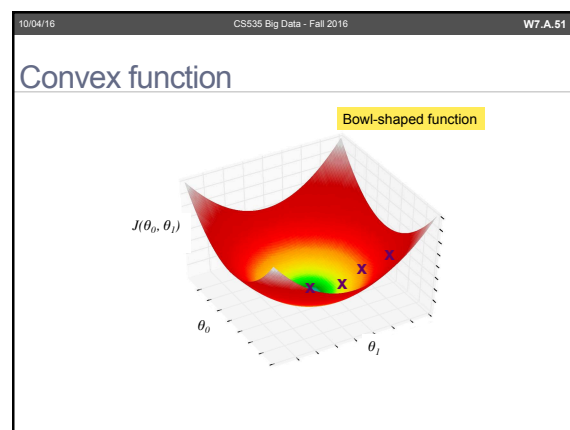
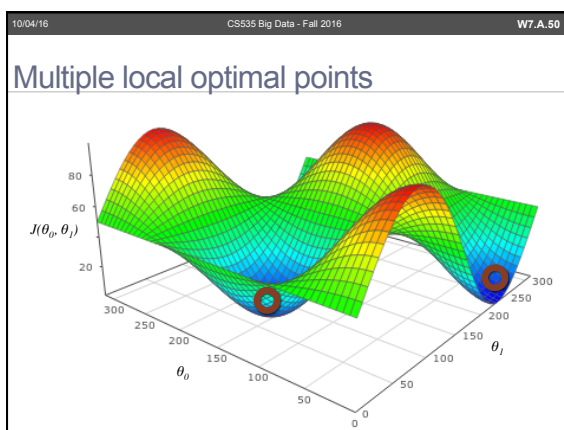
$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)}$$

(for  $j=0$  and  $j=1$ )

Update  $\theta_0$  and  $\theta_1$  simultaneously

}



10/04/16 CS535 Big Data - Fall 2016 W7.A.54

## "Batch" Gradient Descent

- Batch
  - Each step of gradient descent uses all of the training example

$$\theta_j := \theta_j + \alpha \frac{1}{m} \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

10/04/16 CS535 Big Data - Fall 2016 W7.A.55

## Running with Spark in parallel

- For the sample size 1,000 ( $m=1,000$ )
- Batch gradient descent:
 
$$\theta_j := \theta_j + \alpha \frac{1}{1,000} \sum_{i=1}^{1000} (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$
- Using 4 machines

10/04/16 CS535 Big Data - Fall 2016 W7.A.56

## continued

- Step 1. 4 input splits
  - $(x^{(1)}, y^{(1)}), \dots, (x^{(250)}, y^{(250)})$
  - $(x^{(251)}, y^{(251)}), \dots, (x^{(500)}, y^{(500)})$
  - $(x^{(501)}, y^{(501)}), \dots, (x^{(750)}, y^{(750)})$
  - $(x^{(751)}, y^{(751)}), \dots, (x^{(1000)}, y^{(1000)})$
- Step 2. Calculate temp1 ~ 4
 
$$temp1 = \sum_{i=1}^{250} (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

$$temp2 = \sum_{i=1}^{500} (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

$$temp3 = \sum_{i=1}^{750} (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

$$temp4 = \sum_{i=1}^{1000} (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$
- Step 3. Calculate final results