# REGRESSION

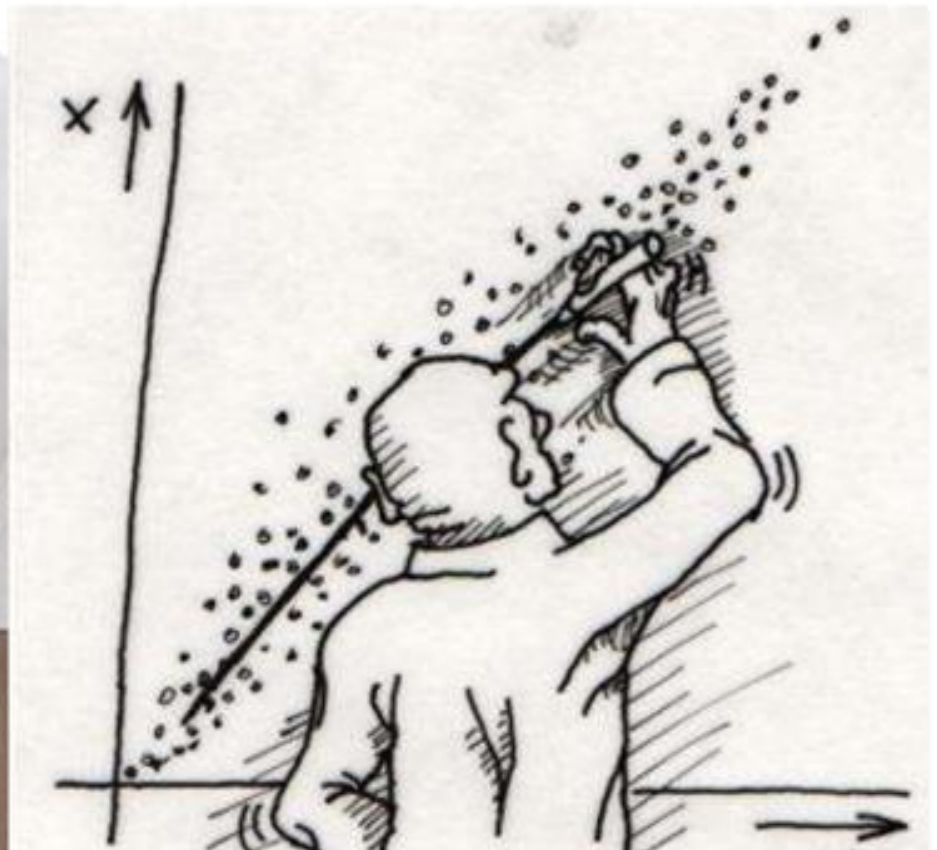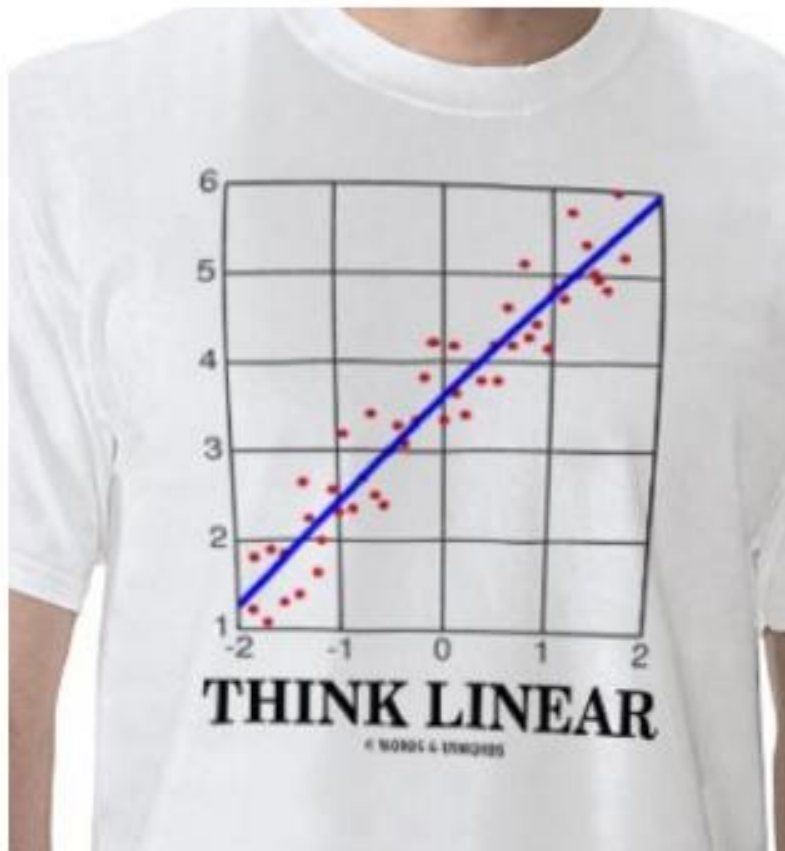**Presented By:**

**Rao Muhammad Umer**

# REGRESSION

# REGRESSION

- **Section-I:**
  - **Motivation**
- **Section-II:**
  - **Theoretical Bases**
- **Section-III:**
  - **Implementation**

# REGRESSION

## PART-I:
## Motivation

# Outline

- **What?**

- **Why?**

- **Who?**

- **How?**

# Outline

- **What?**

- **Why?**

- **Who?**

- **How?**

# Regression

- What about **Machine Learning**?

- What about **Deep Learning**?

- What about **Data Science**?

- What about **Big Data Analysis**?

- What about **Predictive Analytics** ?

# **Regression**

- Predicting a real numeric value for an entity

  with a given set of features

# Deterministic Vs Probabilistic Modeling

- $Y_i = a + b\,X_i$
  - Substituting a value of X in the equation, we can completely determine a unique value of Y.(Exact relationship)
  - Examples: **F = 32 + (9/5) C** ,  $A = \pi r^2$

- $Y_i = a + b\,X_i + e_i$
  - Inexact relationship b/w variables
  - $e_i$ is known as random error

# Types of Regression

- Simple Linear Regression

- Multiple Regression

- Polynomial Regression

- Logistic Regression

# **Outline**

- **What?**

- **Why?**

- **Who?**

- **How?**

# **Predicting House Price**

# **Predicting House Price**

# Predicting which Television Show will have more viewers for next week
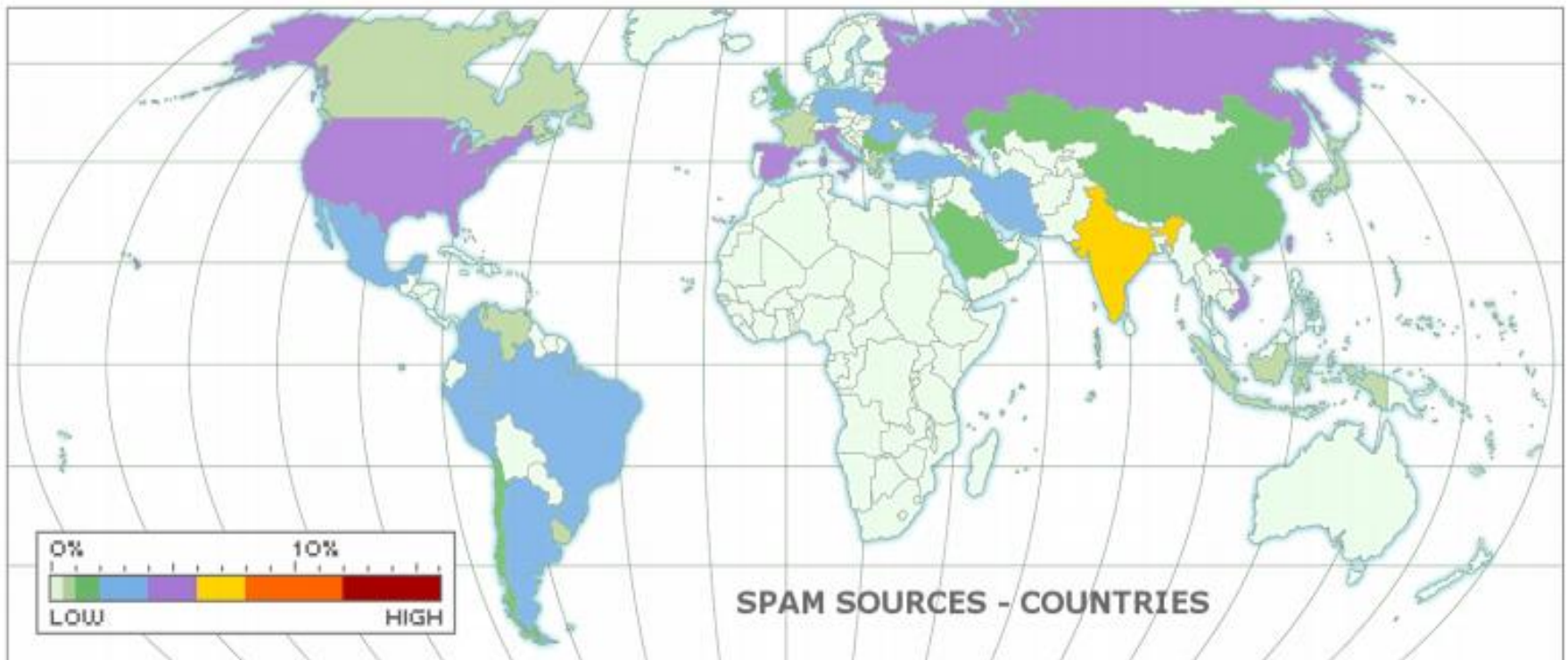
# What Causes Retweets?

- How can we help marketers use Twitter data to influence the number of times a tweet is retweeted?
- With over 555 MM registered users who tweet 58 MM times per day, Tweeting is potentially big business.
- We selected the top 30 entertainment industry related accounts to find out, ranging from rock stars, actors, to famous sports figures.
- In particular, we are attempting to understand the factors that contribute to a "retweet", i.e. an instance where a Twitter account holder's tweet is shared by another Twitter user.
- This has important implications for marketers.
- By understanding the factors that contribute to a retweet a marketer can maximize the potential audience and can better engage with target audiences.

# SPAM CLASSIFICATION

>67%
Of all inbound email is spam

"If you are tired of spam email, imagine if there were no spam filters"

SPAM SOURCES - COUNTRIES

# SPAM CLASSIFICATION

# Tagging a Photo!!!

# Self Customizing Programs

- Self customizing programs
  - Netflix
  - Amazon

# Stock Exchange Price Prediction

• Predicting a continuous response like stock price

# Detecting Fraudulent Transaction

- Online Transactions:

  - Fraudulent((Yes(/(No)?

- Credit card fraud

- Online payment fraud

- Spam instant messages etc.

# Clustering

# **Clustering**

# Outline

- **What?**

- **Why?**

- **Who?**

- **How?**

# **My Background**

- **Rao Muhammad Umer**

- MS Computer Science (Continue…)

- Email address: engr.raoumer943@gmail.com

- LinkedIn Profile: https://pk.linkedin.com/in/raomumer

- Github Repository Link: https://github.com/RaoUmer/GPU-Workshop-PIEAS-2016

# About You

# **Outline**

- **What?**

- **Why?**

- **Who?**

- **How?**

Abstractions...

...and Tools

# IPython Notebooks

# Is this possible for me ???

# **Prerequisites**

- Programming experience

    **C**, C++, Java, Python, **Cuda C** etc.

- Basic statistical knowledge

    Mean, Standard Deviation, Probability etc.

- Willingness to learn new software & tools
  - This can be time consuming
  - You will need to read online documentation

Be Patient

Be Flexible

Be Constructive

# REGRESSION

# PART-II:
# Theoretical Basics

# Supervised Learning

# Unsupervised Learning

# Linear Regression



Housing price prediction.

# **Housing Price Prediction**

- Supervised Learning

  - "right answers" given

- Regression: Predict continuous valued output (price)

# Training set of Housing Prices

| Size in feet² (x) | Price ($) in 1000's (y) |
|---|---|
| 2104 | 460 |
| 1416 | 232 |
| 1534 | 315 |
| 852 | 178 |
| ... | ... |

$m = 47$

# Simple Linear Regression

- **Housing Prices**



- **Supervised Learning**
  - **Given the "right answer" for each example in the data**
- **Regression Problem**
  - **Predict real-valued output**

- **How do we represent h?**



- Linear regression with one variable
- Univariate linear regression

# **Cost function**

**Hypothesis:**

$$h_\theta(x) = \theta_0 + \theta_1 x$$

**Parameters:**

$$\theta_0, \theta_1$$

**Cost Function:**

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

**Goal:** $\displaystyle \min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

# Gradient Descent

- Repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad (\text{for } j = 0 \text{ and } j = 1)$$

  }

# Gradient Descent

# Gradient Descent for Linear Regression

## Gradient descent algorithm

repeat until convergence {

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

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

}

## Linear Regression Model

$$h_\theta(x) = \theta_0 + \theta_1 x$$

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

# **Multiple Linear Regression**

• Linear Regression with Multiple variables

| Size (feet$^2$) | Number of bedrooms | Number of floors | Age of home (years) | Price ($1000) |
|---|---|---|---|---|
| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $y$ |
| 2104 | 5 | 1 | 45 | 460 |
| 1416 | 3 | 2 | 40 | 232 |
| 1534 | 3 | 2 | 30 | 315 |
| 852 | 2 | 1 | 36 | 178 |

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$$

# Polynomial Regression

- Housing Price Prediction

$$h_\theta(x) = \theta_0 + \theta_1 \times \boxed{frontage} + \theta_2 \times \boxed{depth}$$

$x_1$ $x_2$

# **Polynomial Regression**

# **Normal Equation**

- Gradient Descent

$$J(\theta)$$

$$\theta$$

- Normal equation
  - Method to solve for $\theta$ analytically

# Normal Equation

$$\theta = (X^T X)^{-1} X^T y$$

- **pinv(X' * X) *X' * y**

# Logistic Regression



Threshold classifier output $h_\theta(x)$ at 0.5:
If $h_\theta(x) \geq 0.5$, predict "y = 1"
If $h_\theta(x) < 0.5$, predict "y = 0"

# Logistic Regression Model

- Sigmoid Function
- Logistic Function



$$h_\theta(x) = g(\theta^T x)$$

$$g(z) = \frac{1}{1+e^{-z}}$$

Suppose predict "$y = 1$" if $h_\theta(x) \geq 0.5$

predict "$y = 0$" if $h_\theta(x) < 0.5$

# Classification

# REGRESSION

## PART-III: Implementation

# Predicting House Price

# LINEAR REGRESSION EXAMPLE

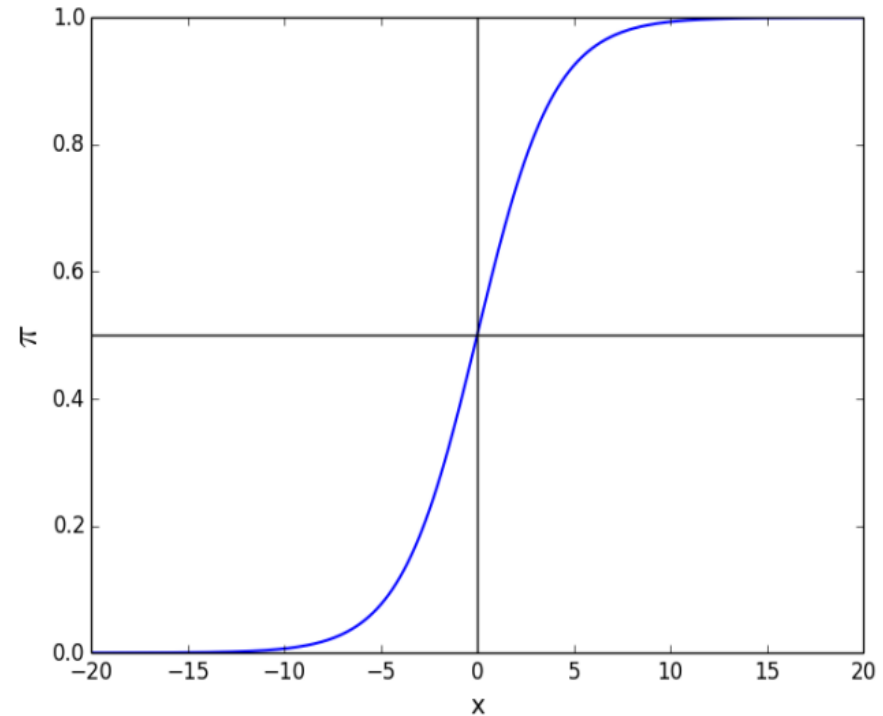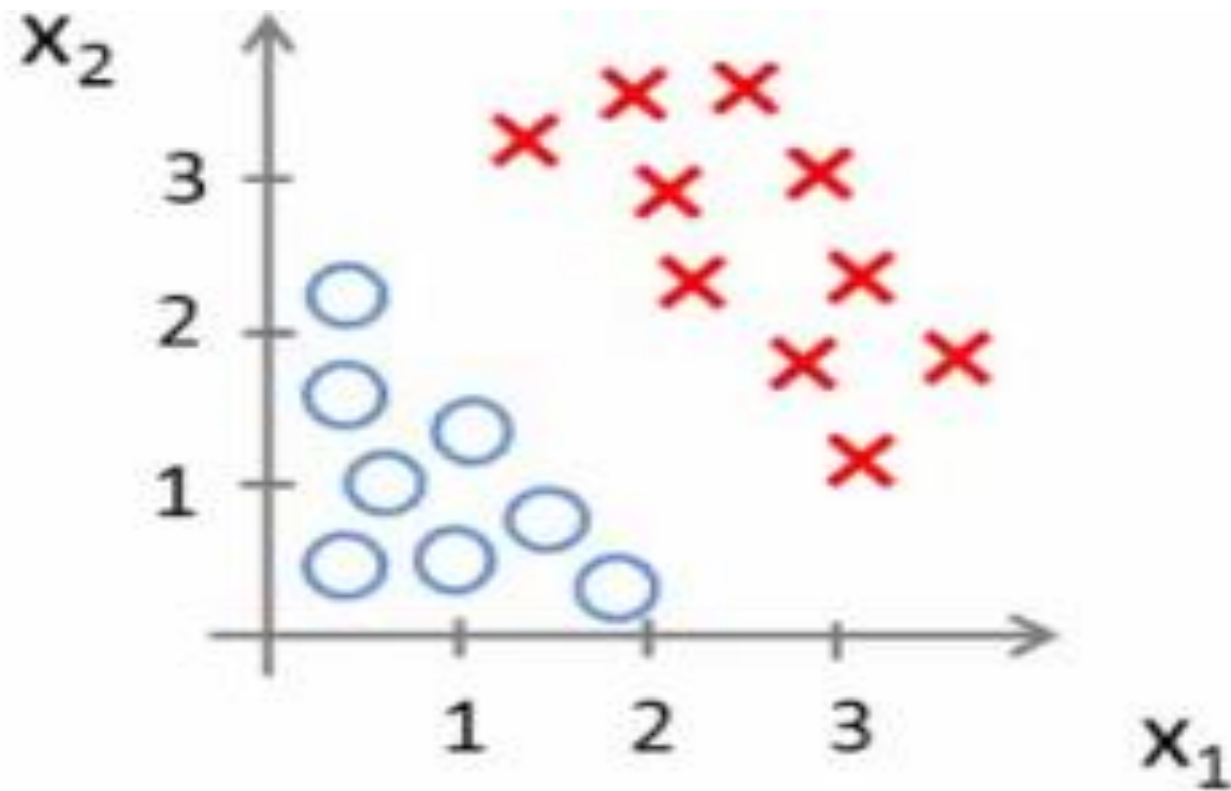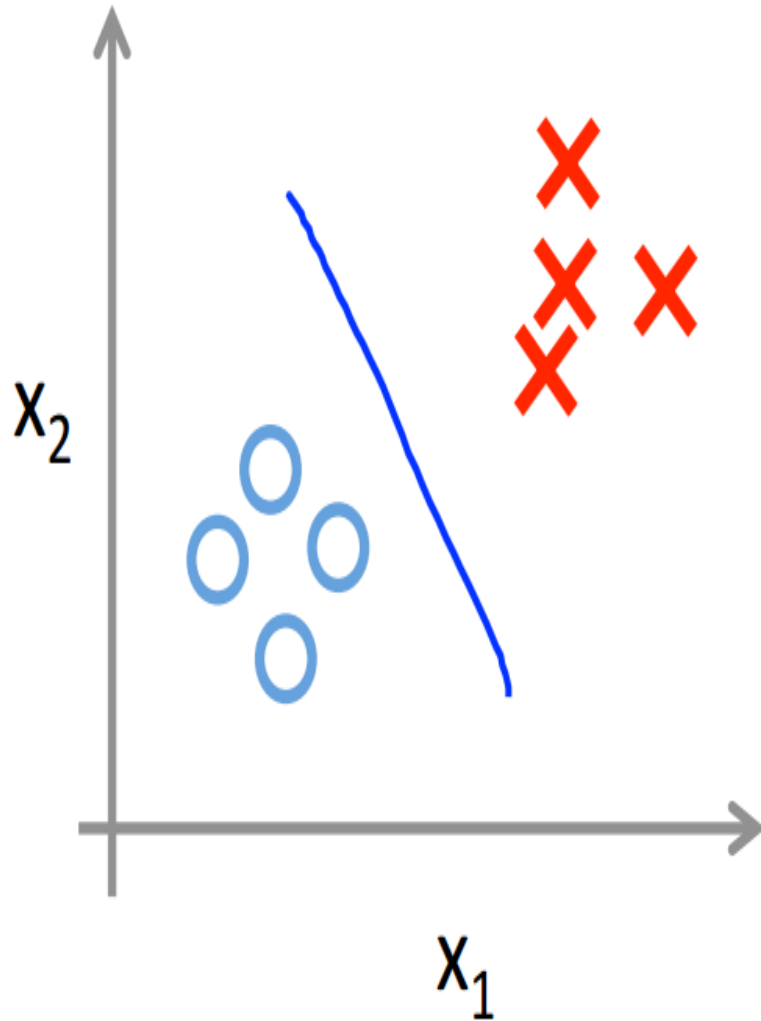| X | Y | XY | $X^2$ |
|---|---|---|---|
| 1125 | 7160 | 8055000 | 1265625 |
| 920 | 8804 | 8099680 | 846400 |
| 835 | 8108 | 6770180 | 697225 |
| 1000 | 6370 | 6370000 | 1000000 |
| 1150 | 6441 | 7407150 | 1322500 |
| 990 | 5154 | 5102460 | 980100 |
| 840 | 5896 | 4952640 | 705600 |
| 650 | 5336 | 3468400 | 422500 |
| 640 | 5041 | 3226240 | 409600 |
| 583 | 5012 | 2921996 | 339889 |
| ΣX | ΣY | ΣXY | $ΣX^2$ |
| 8733 | 63322 | 56373746 | 7989439 |

| X-mean | Y-mean |
|---|---|
| 873.3 | 6332.2 |

| b0 | b1 |
|---|---|
| 3746.198 | 2.961183 |

$$\beta_0 = \frac{(\sum X^2)(\sum Y) - (\sum X)(\sum XY)}{n \sum X^2 - (\sum X)^2}$$

$$\beta_1 = \frac{n \sum XY - (\sum X)(\sum Y)}{n \sum X^2 - (\sum X)^2}$$

# LINEAR REGRESSION EXAMPLE

# **Machine Learning**

# **Machine Learning**



Figure shows how one feature can be used on a linear regression problem to predict new house prices

# **Machine Learning**



On Figure  how the complexity can grow easily from 2 dimensions linear to hundreds of dimensions polynomial

# **Gradient Descent**



Figure shows how a training data is plot and the error is calculated

# Large Scale Machine Learning



In Figure  we can see how the data is split into four parts and fed into four different processors

# **Large Scale Machine Learning**



Figure shows this configuration along with the parallelized part on the GPU cores
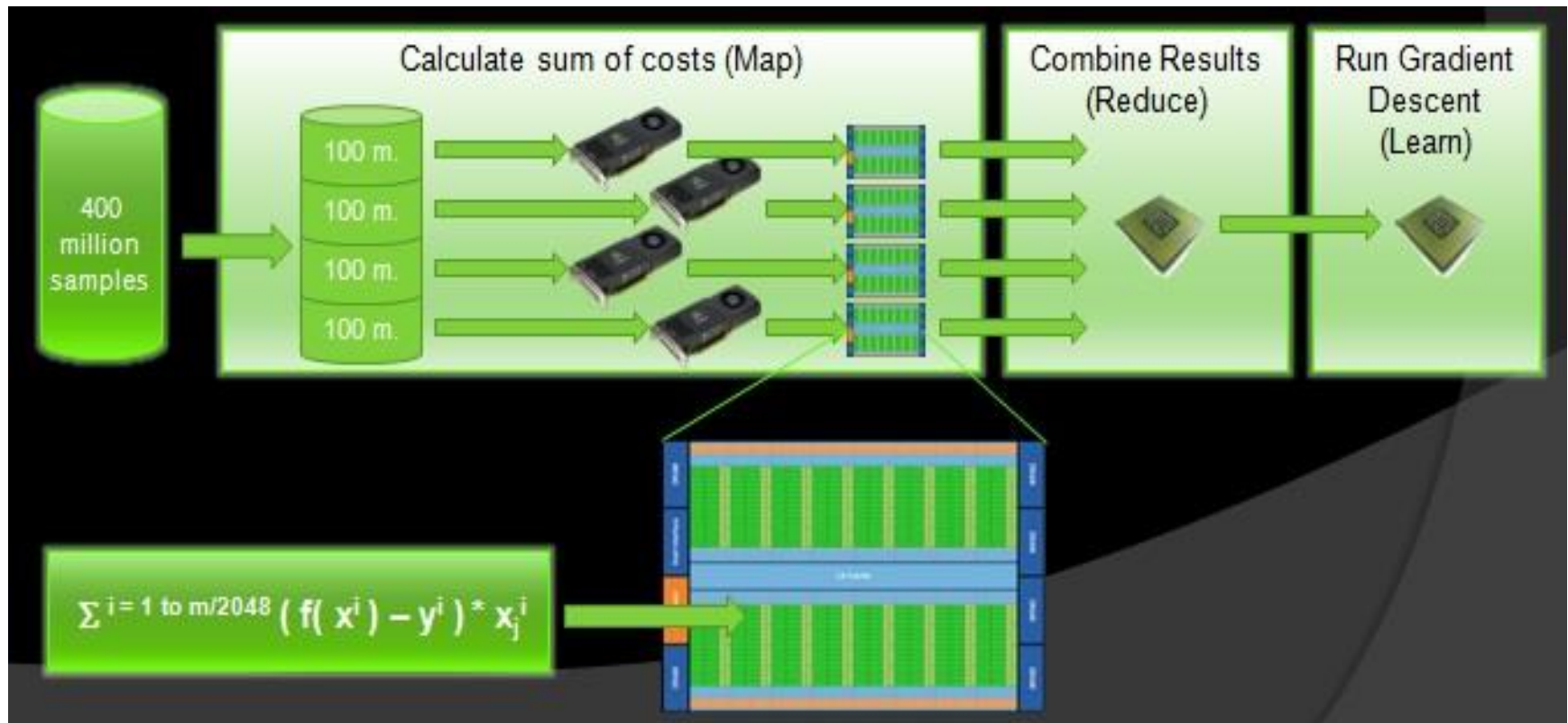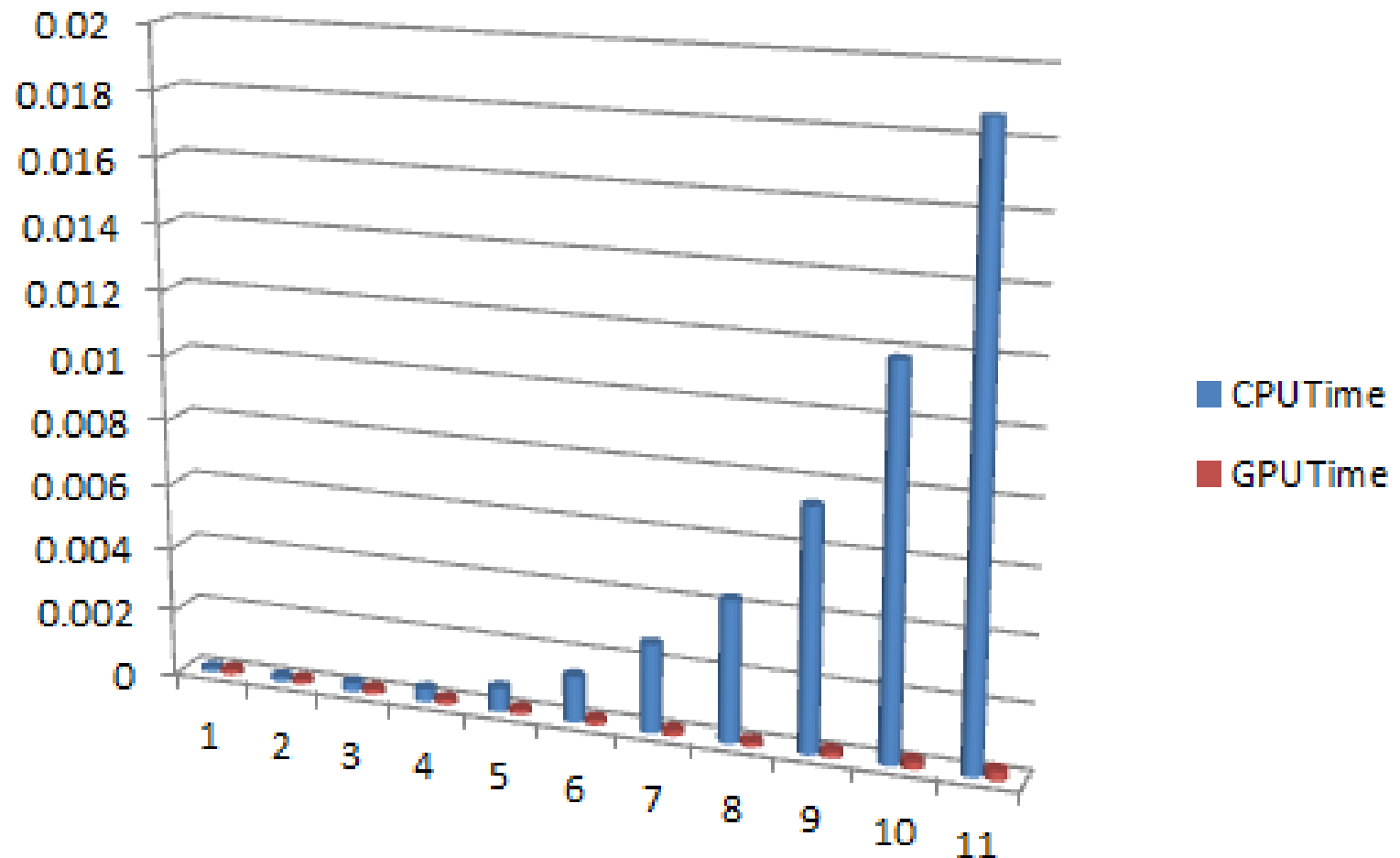
# GPGPU

- Utilizing GPUs to enable dramatic increases in computing performance of general purpose scientific and engineering computing is named GPGPU.

- NVIDIA is providing a parallel computing platform and programming model named CUDA to develop GPGPU software on C, C++ or Fortran which can run on any NVIDIA GPU.

# CPU vs GPU Time

# **Conclusion**

- **GPGPU**, **Machine Learning** and **Big Data** are three rising fields in the IT industry.

- As much as we get deeper into these fields ,we figure out how well they fit together.

- I hope this sample application gave you some basic idea and maybe just one perspective how you can use **NVIDIA CUDA** easily on machine learning problems.

# References

- **http://gpgpu.org/**
- **https://en.wikipedia.org/wiki/Linear_regression**
- **https://docs.nvidia.com/cuda/**
- **https://en.wikipedia.org/wiki/Polynomial_regression**
- **https://en.wikipedia.org/wiki/Gradient_descent**
- **https://en.wikipedia.org/wiki/Overfitting**
- **https://en.wikipedia.org/wiki/Regularization_(mathematics)**
- **http://www.nvidia.com/object/machine-learning.html**
- **https://developer.nvidia.com/nvidia-nsight-visual-studio-edition**
- **https://docs.nvidia.com/cuda/thrust/index.html**
- **https://adnanboz.wordpress.com/2012/02/25/large-scale-machine-learning-using-nvidia-cuda/**

# References

- **http://statsmodels.sourceforge.net/**
- **http://www.nature.com/ejhg/journal/v17/n8/full/ejhg20 095a.html**
- **http://www.alsharif.info/#!iom530/c21o7**