# Lecture 2 - Regression

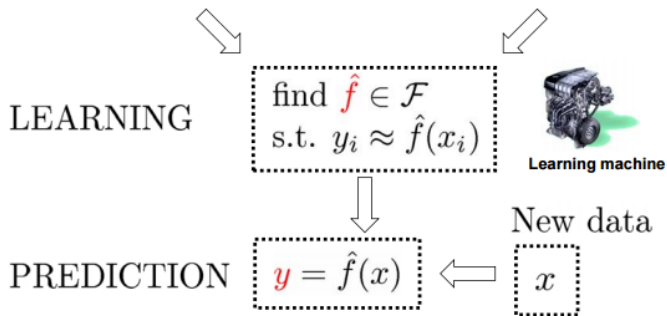Instructor: Prof. Ganesh Ramakrishnan

# Supervised Learning

Functions $F$      Training Data

$$f : X \to Y \quad \{ (x^i, y^i) \in X * Y \}$$



LEARNING    find $\hat{f} \in \mathcal{F}$
s.t. $y_i \approx \hat{f}(x_i)$

Learning machine

New data

PREDICTION    $y = \hat{f}(x) \Longleftarrow x$

# Next ....

We will start with linear regression and least square method to calculate parameters for linear regression problems.

# Recap

- **Machine Learning in general**
  - ▶ Supervised Learning
  - ▶ Unsupervised Learning
  - ▶ Applications and examples

- **Canonical Learning Problems**
  - ▶ Regression Supervised
  - ▶ Classification Supervised
  - ▶ Unsupervised modeling of data

# Agenda

- What is data?
  - Noise in data
- How to predict?
  - Fitting a curve
  - Error measurement
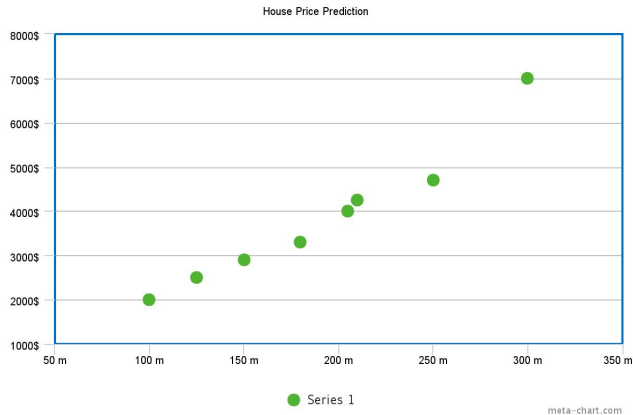  - Minimizing Error
- Method of Least Squares

# What is data?

- For us, data is the information about the problem, you are solving using ML, in quantized form
- This data can be from any source, some examples are
  - Prices of stock and stock indexes such as BSE or Nifty
  - Prices of house, area and size of the house
  - Temperature of a place, latitude, longitude and time of year
- The objective of ML is to predict or classify something using the given data
- Hence, one or more than one parameters of the data must also represent the output of our program

# Noise in Data

- Data in real life problems are generally collected through surveys
- Surveys may have random human errors
- Most methods we will be using deal with expectations as they minimize the effect of error in our predictions
- Data Cleansing by finding outliers

# Example dataset for this lecture

- For this lecture we will consider variation of cost of the house with the area of the house
- In this example we want to find a pattern or curve which this dataset follows, hence predict the price for any value of area



House Price Prediction

# How to predict?

- Curve fitting is the process of constructing a curve, or mathematical function, that has the best fit to a series of data points, possibly subject to constraints. - Wikipedia
- Thus we need a critera to compare two curves on a dataset
- We describe an error function $E(f, D)$ which takes a curve $f$ and dataset $D$ as input and returns a real number
- Error function must be such that it can capture how bad the prediction is

# Example

- Consider the example below where we have two curves on our dataset defined by blue($f_b$) and red($f_r$) line respectively. We want to find which is the better fit.



Figure: House purchase data curve fit

# Question

What are some options for $E(f, D)$?
Hint: Measurement of difference from original value.

# Examples of $E$

- $\sum_D f(x_i) - y_i$
- $\sum_D |f(x_i) - y_i|$
- $\sum_D (f(x_i) - y_i)^2$
- $\sum_D (f(x_i) - y_i)^3$
- and many more

# Question

What $E$ do you think can give us best fit curve and why?
Hint: Intuition of distances.

# Squared Error

$$\sum_D (f(x_i) - y_i)^2$$

- To find the best fit curve we try to minimize the above function
- It is continuous and differentiable
- It can be visualized as square of Euclidean distance between predicted points and actual points
- How we can perform mathematical treatment over this function will be covered in further lectures.
- This mathematical treatment is known as method of least squares.

# Regression, More Formally

- Formal Definition
- Types of Regression
- Geometric Interpretation of least square solution

Linear Regression as a canonical example

- **Optimization** (Formally deriving least Square Solution)
- **Regularization** (Ridge Regression, Lasso), **Bayesian Interpretation** (Bayesian Linear Regression)
- **Non-parametric estimation** (Local linear regression),
- **Non-linearity through Kernels** (Support Vector Regression)

# Linear Regression with Illustration

- Regression is about learning to predict a set of output variables (*dependent variables*) as a function of a set of input variables (*independent variables*)
- Example
  - A company wants to determine how much it should spend on T.V commercials to increase sales to a desired level $y^*$
  - **Basis?**

# Linear Regression with Illustration

- Regression is about learning to predict a set of output variables (*dependent variables*) as a function of a set of input variables (*independent variables*)
- Example
  - A company wants to determine how much it should spend on T.V commercials to increase sales to a desired level $y^*$
  - **Basis?** It has previous observations of the form $<x_i, y_i>$,
    - $x_i$ is an instance of money spent on advertisements and $y_i$ was the corresponding observed sale figure

# Linear Regression with Illustration

- Regression is about learning to predict a set of output variables (*dependent variables*) as a function of a set of input variables (*independent variables*)
- Example
  - A company wants to determine how much it should spend on T.V commercials to increase sales to a desired level $y^*$
  - **Basis?** It has previous observations of the form $<x_i, y_i>$,
    - $x_i$ is an instance of money spent on advertisements and $y_i$ was the corresponding observed sale figure
  - Suppose the observations support the following linear approximation

$$y = \beta_0 + \beta_1 * x \tag{1}$$

  Then $x^* = \frac{y^* - \beta_0}{\beta_1}$ can be used to determine the money to be spent

- **Estimation** for Regression: Determine appropriate value for $\beta_0$ and $\beta_1$ from the past observations
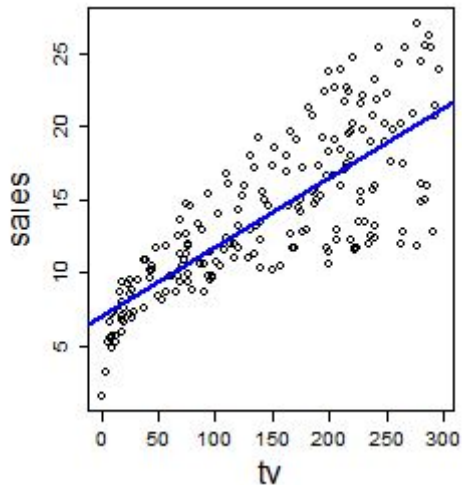
# Linear Regression with Illustration



Figure: Linear regression on T.V advertising vs sales figure

What will it mean to have sales as a non-linear function of investment in advertising?

# Basic Notation

- Data set: $\mathcal{D} = <\mathbf{x_1}, \mathbf{y_1}>, .., <\mathbf{x_m}, \mathbf{y_m}>$
  - Notation (used throughout the course)
    - $m$ = number of training examples
    - $\mathbf{x}'s$ = input/independent variables
    - $\mathbf{y}'s$ = output/dependent/'target' variables
    - $(\mathbf{x}, \mathbf{y})$ - a single training example
    - $(\mathbf{x}_j, \mathbf{y}_j)$ - specific example ($j^{th}$ training example)
    - $j$ is an index into the training set

- $\phi_i$'s are the attribute/basis functions, and let

$$\Phi = \begin{bmatrix} \phi_1(\mathbf{x}_1) & \phi_2(\mathbf{x}_1) & ...... & \phi_p(\mathbf{x}_1) \\ . & & & \\ . & & & \\ \phi_1(\mathbf{x}_m) & \phi_2(\mathbf{x}_m) & ...... & \phi_p(\mathbf{x}_m) \end{bmatrix} \qquad (2)$$

-

$$\mathbf{y} = \begin{bmatrix} y_1 \\ . \\ y_m \end{bmatrix} \qquad (3)$$

# Formal Definition

- **General Regression problem**: Determine a function $f^*$ such that $f^*(x)$ is the best predictor for $y$, with respect to $\mathcal{D}$:

$$f^* = \operatorname*{argmin}_{f \in F} E(f, \mathcal{D})$$

Here, $F$ denotes the class of functions over which the error minimization is performed

- **Parametrized Regression problem**: Need to determine parameters $\mathbf{w}$ for the function $f(\phi(\mathbf{x}), \mathbf{w})$ which minimize our error function $E\big(f(\phi(\mathbf{x}), \mathbf{w}), \mathcal{D}\big)$

$$\mathbf{w}^* = \operatorname*{argmin}_{\mathbf{w}} \left\langle E\big(f(\phi(\mathbf{x}), \mathbf{w}), \mathcal{D}\big) \right\rangle$$

# Types of Regression

- Classified based on the function class and error function
- $F$ is space of linear functions $f(\phi(\mathbf{x}), \mathbf{w}) = \mathbf{w}^T \phi(\mathbf{x}) + b \implies$ Linear Regression
  - Problem is then to determine $\mathbf{w}^*$ such that,

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \; E(\mathbf{w}, \mathcal{D}) \tag{4}$$

# Types of Regression (contd.)

- **Ridge Regression:** A shrinkage parameter (regularization parameter) is added in the error function to reduce discrepancies due to variance
- **Logistic Regression:** Models conditional probability of dependent variable given independent variables and is extensively used in classification tasks

$$f(\phi(\mathbf{x}), \mathbf{w}) = \log \frac{\Pr(\mathbf{y}|\mathbf{x})}{1 - \Pr(\mathbf{y}|\mathbf{x})} = b + \mathbf{w}^T * \phi(\mathbf{x}) \tag{5}$$

- Lasso regression, Stepwise regression and several others

# Least Square Solution

- Form of $E()$ should lead to accuracy and tractability
- The squared loss is a commonly used error/loss function. It is the sum of squares of the differences between the actual value and the predicted value

$$E(f, \mathcal{D}) = \sum_{j=1}^{m} (f(x_j) - y_j)^2 \tag{6}$$

- The least square solution for linear regression is obtained as

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \ \sum_{j=1}^{m} (\sum_{i=1}^{p} (w_i \phi_i(x_j) - y_j)^2) \tag{7}$$

- The minimum value of the squared loss is zero
- If zero were attained at $\mathbf{w}^*$, we would have ...................

- The minimum value of the squared loss is zero
- If zero were attained at $\mathbf{w}^*$, we would have $\forall u, \phi^T(x_u)\mathbf{w}^* = y_u$, or equivalently $\Phi\mathbf{w}^* = \mathbf{y}$, where

$$\Phi = \begin{bmatrix} \phi_1(x_1) & ... & \phi_p(x_1) \\ ... & ... & ... \\ \phi_1(x_m) & ... & \phi_p(x_m) \end{bmatrix}$$

and

$$\mathbf{y} = \begin{bmatrix} y_1 \\ ... \\ y_m \end{bmatrix}$$

- It has a solution if $\mathbf{y}$ is in the column space (the subspace of $R^m$ formed by the column vectors) of $\Phi$

- The minimum value of the squared loss is zero
- If zero were NOT attainable at $\mathbf{w}^*$, what can be done?

# Geometric Interpretation of Least Square Solution

- Let $\mathbf{y}^*$ be a solution in the column space of $\Phi$
- The least squares solution is such that the distance between $\mathbf{y}^*$ and $\mathbf{y}$ is minimized
- Therefore............

# Geometric Interpretation of Least Square Solution

- Let $\mathbf{y}^*$ be a solution in the column space of $\Phi$
- The least squares solution is such that the distance between $\mathbf{y}^*$ and $\mathbf{y}$ is minimized
- Therefore, the line joining $\mathbf{y}^*$ to $\mathbf{y}$ should be orthogonal to the column space

$$\phi \mathbf{w} = \mathbf{y}^* \tag{8}$$

$$(\mathbf{y} - \mathbf{y}^*)^T \Phi = 0 \tag{9}$$

$$(\mathbf{y}^*)^T \Phi = (\mathbf{y})^T \phi \tag{10}$$

$$(\phi\mathbf{w})^T\Phi = \mathbf{y}^T\Phi \tag{11}$$

$$\mathbf{w}^T\Phi^T\Phi = \mathbf{y}^T\Phi \tag{12}$$

$$\Phi^T\Phi\mathbf{w} = \Phi^T\mathbf{y} \tag{13}$$

$$\mathbf{w} = (\Phi^T\Phi)^{-1}\Phi^T\mathbf{y} \tag{14}$$

- Here $\Phi^T\Phi$ is invertible if and only if $\Phi$ has full column rank

Proof?

**Theorem** : $\Phi^T\Phi$ is invertible if and only if $\Phi$ is full column rank

Proof :

Given that $\Phi$ has full column rank and hence columns are linearly independent, we have that $\Phi\mathbf{x} = 0 \Rightarrow \mathbf{x} = 0$

Assume on the contrary that $\Phi^T\Phi$ is non invertible. Then $\exists \mathbf{x} \neq 0$ such that $\Phi^T\Phi\mathbf{x} = 0$
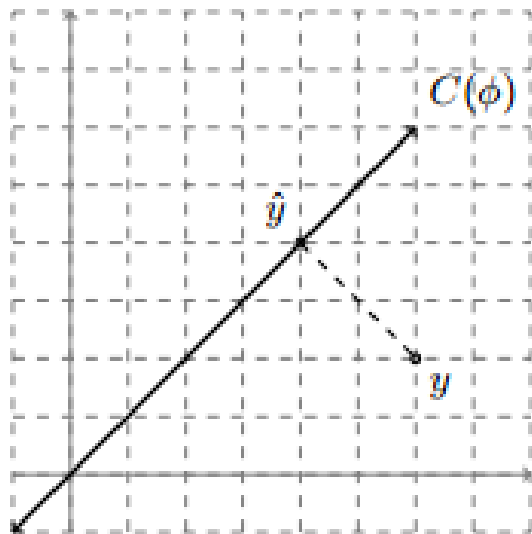
$$\Rightarrow \mathbf{x}^T\Phi^T\Phi\mathbf{x} = 0$$
$$\Rightarrow (\Phi\mathbf{x})^T\Phi\mathbf{x} = 0$$
$$\Rightarrow \Phi\mathbf{x} = 0$$

This is a contradiction. Hence $\Phi^T\Phi$ is invertible if $\Phi$ is full column rank

If $\Phi^T\Phi$ is invertible then $\Phi\mathbf{x} = 0$ implies $(\Phi^T\Phi\mathbf{x}) = 0$, which in turn implies $\mathbf{x} = 0$ , This implies $\Phi$ has full column rank if $\Phi^T\Phi$ is invertible. The converse can also be proved similarly.

# How about an Analytic Derivation?

- Some more questions on the Least Square Linear Regression Model
- More generally: How to minimize a function?
  - Level Curves and Surfaces
  - Gradient Vector
  - Directional Derivative
  - Hyperplane
  - Tangential Hyperplane
- Gradient Descent Algorithm