# Regression
Instructor: Prof. Ganesh Ramakrishnan

# Recap

- Supervised (Classification and Regression) vs Unsupervised Learning
    - Three canonical learning problems
- What is data and how to predict
    - More on this today in the context of regression
- Squared Error

# Agenda

- What is Regression
- Formal Defintion
- Types of Regression
- Least Square Solution
- Geometric Interpretation of least square solution

# Regression

- Finding correlation between a <u>set of output variables</u> and a set of <u>input variables</u> *(x)*

  *so far single output variable (y)*

- Input variables are called *independent variables*

- Output variables are called *dependent variables*

$$y = f(x) \ldots \text{Linear regression, } f \text{ is linear}$$

# Examples

- A company wants to how much money they need to spend on T.V advertising to increase sales to a desired level, say y*
- They have <u>previous data of form</u> $<x_i, y_i>$, where $x_i$ is money spent on advertising and $y_i$ are sale figures
- They now fit the data with a function, lets say linear function

$\therefore$ To get $\Delta y$ increase
in y, you need $\Delta y / \beta_1$ ← $y = \beta_0 + \beta_1 * x$ → so that $y_i \approx \beta_0 + \beta_1 x_i$ $\forall i$ (1)
increase in x

and then find the money they need to spend using this function

- Regression problem is to find the appropriate function and its coefficients
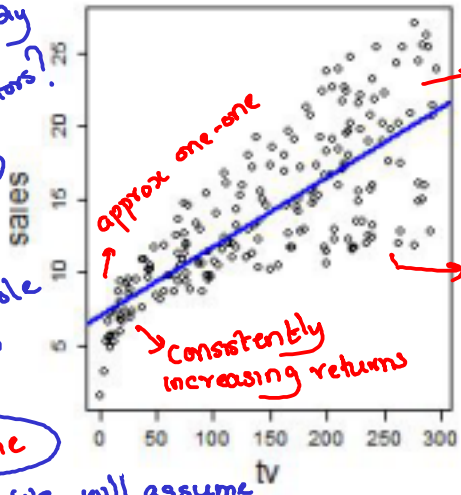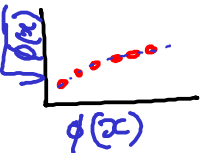
We will replace $x$ by $\phi(x)$

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

# What if sales is a non-linear function of advertising?

$$y = \beta_0 + \beta_1 \phi(x)$$

sales ___ investment in advertisement.

$x =$ date for TV company



If $y$ varies approximately as $\sqrt{\phi(x)}$

$$y = \beta_0 + \beta_1 \sqrt{\phi(x)}$$

OR $y$ varies approximately as $\phi^2(x)$

$$y = \beta_0 + \beta_1 \phi^2(x)$$

In reality, could be a combination!

$$y = \beta_0 + \beta_1 \phi(x) + \beta_2 \sqrt{\phi(x)} + \beta_3 \phi^2(x) + \cdots$$

we want to learn $\beta_0, \beta_1, \beta_2 \cdots$

We assume 1-1 mapping between $y$ & $x$ (approx mapping)

Given n observations:

$$\begin{bmatrix} y_1, \phi_1(x_1) \ \phi_2(x_1) \cdots \phi_p(x_1) \\ y_2, \phi_1(x_2), \phi_2(x_2) \cdots \phi_p(x_2) \\ \vdots \\ y_n, \phi_1(x_m), \phi_2(x_m) \cdots \phi_p(x_m) \end{bmatrix}$$

$x_i$ = corresponds to day $i$ for the TV company

$\phi_1(x_i), \phi_2(x_i) \cdots \phi_p(x_i)$ are different attributes of the T.V company on day $i$

$y_i$ is the amount of sales on day $i$

We need to estimate $w_0, w_1, w_2 \cdots w_k$ ($\beta_1, \beta_2 \cdots$ on previous slide) such that

$$y_i \approx w_0 + w_1 \phi_1(x_i) + w_2 \phi_2(x_i) + \cdots + w_k \phi_k(x_i)$$

For eg, on prev slide: $\phi_1(x_i) = $ ad investment on $i^{th}$ day    $\sqrt{\phi_1(x_i)}$    $\phi_1^2(x_i)$
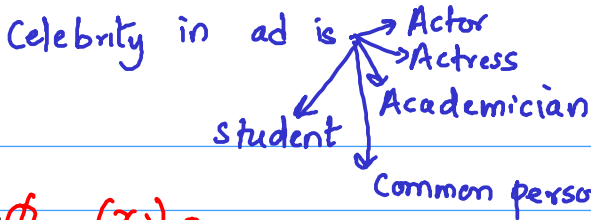
Some more notes:

① $\phi_1 \ldots \phi_K$ can be any attribute that has the <u>potential of</u> linearly influencing $y$

② $\phi_1(x_i), \phi_2(x_i) \ldots \phi_p(x_i)$ may or may not be functions of each other ...
But certain machine learning approaches are robust to interdependencies among $\phi_j$'s and certain others are NOT

③ You could use ML workbench such as WEKA (Waikato univ, New Zealand) to study $\phi_j$'s that "linearly influence" $y$.

④ How to deal with the attribute:
Celebrity in ad is → Actor
→ Actress
→ Academician
↙ Student
↓ Common person

Example of nominal attribute

$$\phi_{c=ar}(x_i)$$
$$\phi_{c=as}(x_i)$$
$$\phi_{c=cp}(x_i)$$
$$\phi_{c=s}(x_i)$$

For a given $i$ exactly one of these can take Value $=1$ and the rest take value $=0$

Nominal attribute: One that takes one value from a set of discrete possible values.

# Formal Definition

- Two sets of variables: $x \in \mathcal{R}^N$ (independent) and $y \in \mathcal{R}^k$ (dependent)
- D is a set of m data points: $<x_1, y_1>, <x_2, y_2>, ..., <x_m, y_m>$
- $\epsilon$ (f, D): An error function, designed to reflect the discrepancy between the predicted value $f(x_i)$ and $y_i$ $\forall i$    $f(x_i) = \omega_0 + \omega_1 \phi_1(x_i)$
  $+ \omega_2 \phi_2(x_i) \dots$
- Regression problem: Determine a function $f^*$ such that $f^*(x)$ is the best predictor for y, with respect to D,

$$f^* = \underset{f \in F}{\operatorname{argmin}} \; \epsilon(f, D) \qquad (2)$$

where, F denotes the class of functions over which the optimization is performed

# Types of Regression

- Depends on the <u>function class</u> and <u>error function</u>
- Linear Regression : establishes a relationship between dependent variable (Y) and one or more independent variables (X) using a best fit straight line, i.e

$$Y = a + b * X \tag{3}$$

  - Here F is of the form $\Sigma_{i=1}^{p} w_i \phi_i(x)$, where $\phi_i$ are called basis functions <span style="color:blue">(or attributes)</span>
  - Problem is to find $w^*$ where

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

- Ridge Regression : A shrinkage parameter (regularization parameter) is added in the <u>error function</u> to reduce discrepancies due to variance → *Linear regression with good generalization*
- Logistic Regression : Used to model conditional probability of dependent variable given independent variable and is extensively used in classification tasks

$$log\frac{p(y|x)}{1 - p(y|x)} = \beta_0 + \beta * x \quad (5)$$

→ *Different function class*

- <u>Lasso regression</u>, Stepwise regression and many more

↓

*Different error function*

# Least Square Solution

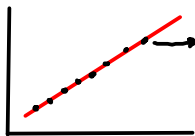- Form of $\epsilon$ plays a major role in the accuracy and tractability of the optimization problem
- The squared loss is a commonly used error/loss function. It is the sum of squares of the differences between the actual value and the predicted value

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

- The least square solution for linear regression is given by

$$w^* = \underset{w}{\text{argmin}} \sum_{j=1}^{m}\left(\underbrace{\sum_{i=1}^{p} w_i \phi_i(x_j)}_{f(x_j)} - y_j\right)^2 \tag{7}$$

$$f(x_j) = w_1^* \phi_1(x_j) + \ldots w_p^* \phi_p(x_j) = y_j$$

assume that some $\phi_k(x_j) = 1$ to account for offset/bias of linear function

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

Suppose

$$\sum_{i=1}^{p} w_i^* \phi_i(x_j) = y_j \quad \forall j = 1 \ldots m$$

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

$$\phi = \begin{bmatrix} \phi_1(x_1) & \ldots & \phi_p(x_1) \\ \ldots & \ldots & \ldots \\ \phi_1(x_m) & \ldots & \phi_p(x_m) \end{bmatrix}$$
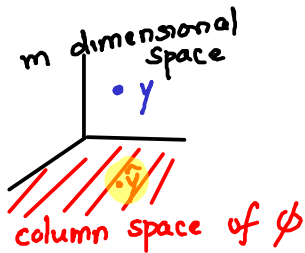
*each row is for an example*

and

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

- It has a solution if $\mathbf{y}$ is in the column space (the subspace of $R^n$ formed by the column vectors) of $\phi$ *: Obtain $\omega^*$ using Gaussian elimination*

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

Least squares :
- Consider all $\widehat{y}$'s in column space of $\phi$ matrix
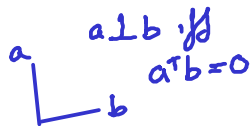- Find that $\widehat{y}$ which is as close as possible to observed $y$
- Thus, $\widehat{y}$ in column space st $(y-\widehat{y})$ is $\perp$ to column space is the LS soln

m dimensional space

• $y$

• $\widehat{y}$

column space of $\phi$

# 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............ $y^* = \phi w$

$$(y - y^*)^T \phi = 0$$

$a \perp b$

$a^T b = 0$

# 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

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

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

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

$$\phi^{\mathbf{T}}(\phi\omega)$$
$$= \phi^{\mathbf{T}}y^*$$
$$= \phi^{\mathbf{T}}y$$

$$(\phi^{\mathbf{T}}\phi)\,\omega = \phi^{\mathbf{T}}y \Rightarrow \omega = (\phi^{\mathbf{T}}\phi)^{-1}\phi^{\mathbf{T}}y$$

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

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

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

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

- Here $\phi^{T}\phi$ is invertible 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 $\phi\mathbf{x} = \mathbf{0} \Rightarrow \mathbf{x} = \mathbf{0}$

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

$$\Rightarrow \mathbf{x^T}\phi^{\mathbf{T}}\phi\mathbf{x} = \mathbf{0}$$
$$\Rightarrow (\phi\mathbf{x})^{\mathbf{T}}\phi\mathbf{x} = \mathbf{0}$$
$$\Rightarrow \phi\mathbf{x} = \mathbf{0}$$

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

If $\phi^T\phi$ is invertible then $\phi\mathbf{x} = \mathbf{0}$ implies $(\phi^T\phi\mathbf{x}) = \mathbf{0}$, which in turn implies $\mathbf{x} = \mathbf{0}$ , **This implies $\phi$ has full column rank if $\phi^T\phi$ is invertible. Hence, theorem proved**
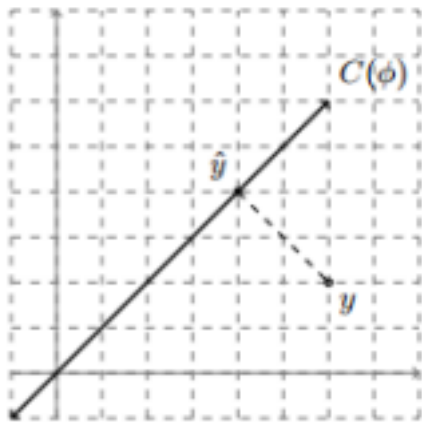
Figure: Least square solution $\mathbf{y}^*$ is the orthogonal projection of y onto column space of $\phi$