

Lecture 3 - Regression

Instructor: Prof. Ganesh Ramakrishnan

The Simplest ML Problem: Least Square Regression

- Curve Fitting: Motivation
 - ▶ Error measurement
 - ▶ Minimizing Error
- Method of Least Squares

linear algebra, optimization

Simplest is least squares,
bayesian linear reg,
regularized error,
parametric vs non-param

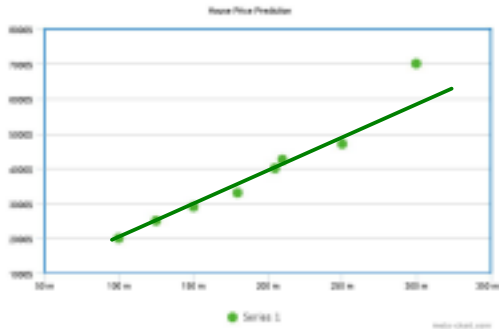
Curve Fitting: Motivation



- Example scenarios:
 - ▶ Prices of house to be fitted as a function of **the area of the house**
 - ▶ Temperature of a place to be fitted as a function of **its latitude and longitude and time of the year**
 - ▶ Stock Price (or BSE/Nifty value) to be fitted as a function of **Company Earnings** *Multivariate regression*
 - ▶ Height of students to be fitted as a function of **their weight**
- One or more **observations/parameters in the data** are expected to represent the output in the future

Higher you go, the more expensive the house!

- Consider the variation of price (in \$) of house with variations in its height (in m) above the ground level (Mumbai)
- These are specified as coordinates of the 8 points:
 $(x_1, y_1), \dots, (x_8, y_8)$
- Desired: Find a pattern or curve that characterizes the price as a function of the height



Errors and Causes

- (Observable) Data is generally collected through measurements or surveys
 - ▶ Surveys can have random human **errors** → *y axis error*
 - ▶ Measurements are subject to **imprecision** of the measuring or recording instrument → *x-axis error*
 - ▶ **Outliers** due to variability in the measurement or due to some experimental error;
- **Robustness to Errors:** Minimize the effect of error in predicted model
- **Data cleansing:** Outlier handling in a pre-processing step
 - ↳ *most often, you want the model building to do implicit data cleansing.*

Curve Fitting: The Process

- *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

Curve Fitting: The Process

- *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
- Need quantitative criteria to find the **best fit**
- **Error function** E : curve $f \times$ dataset $\mathcal{D} \rightarrow \mathfrak{R}$
- Error function must capture the deviation of prediction from expected value

Example

- Consider the two candidate prediction curves in blue and red respectively respectively. Which is the better fit?

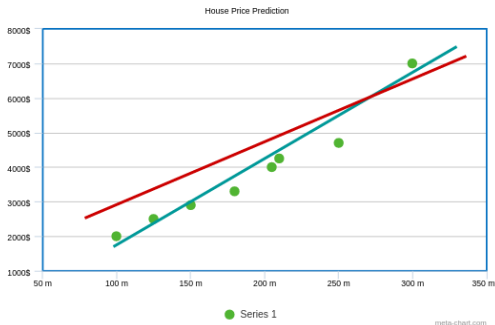
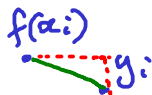


Figure: Price of house vs. its height - for illustration purpose only

Question



What are some options for error function $E(f, D)$ that measure the deviation of prediction from expected value?

$$\sum_i (f(x_i) - y_i)^2 : \text{Euclidean distance}$$

$$\sum_i |f(x_i) - y_i| : \text{Manhattan distance}$$

$$\sum_i f(x_i) - y_i : \text{Unsigned distance} \\ (\text{when bias is desired})$$

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

} In 1-d case, i.e. $y_i, f(x_i) \in \mathbb{R}$
these behave similarly except
that $(f(x_i) - y_i)^2$ discourages
really far away pts
So $(f(x_i) - y_i)^2$ might be
less robust to outliers than
 $|f(x_i) - y_i|$

Fixing outlier sensitivity
can be through ① Fixing error fn OR
② Regularization / bayesian estimation

Question

Which choice F do you think can give us best fit curve and why?

Hint: Think of these errors as distances.

Squared Error

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

- One best fit curve corresponds to f that minimizes the above function. It..
 - ① Is continuous and differentiable
 - ② Can be visualized as square of Euclidean distance between predicted and observed values
- Mathematical optimization of this function: Topic of following lectures.
- This is the 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 $\langle x_i, y_i \rangle$,
 - ★ 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 $\langle x_i, y_i \rangle$,
 - ★ 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

$$y = \beta_0 + \beta_1 * x \quad (1)$$

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 β_0 and β_1 from the past observations $\beta_0, \beta_1 = \operatorname{argmin} \sum_{x \in D} (y - f(x))^2$

Linear Regression with Illustration

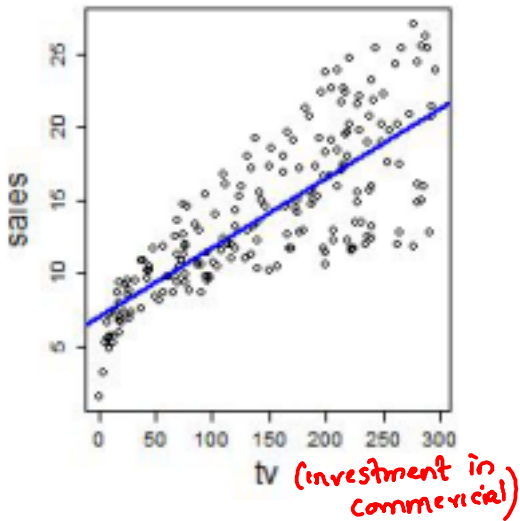


Figure: Linear regression on T.V advertising vs sales figure

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

$y \approx f(\phi_1(x), \phi_2(x) \dots \phi_n(x))$

ϕ_i : value

$\phi_2(x)$ \rightarrow (value)²

$\phi_n(x)$ \rightarrow log(value) or time of day of advertisement

x is an object representing specific investments

$\phi_1 \dots \phi_n$ are basis fns or features with hope that linear combination of ϕ 's as f is a good approx to y

$f(x) = w^T \phi(x)$ [Step 1 to non-linearity]

Basic Notation

- Data set: $\mathcal{D} = \langle \mathbf{x}_1, \mathbf{y}_1 \rangle, \dots, \langle \mathbf{x}_m, \mathbf{y}_m \rangle$
 - Notation (used throughout the course)
 - m = number of training examples
 - \mathbf{x}' 's = input/independent variables **objects**
 - \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
- ϕ_i 's are the attribute/basis functions, and let

$$\phi = \begin{bmatrix} \phi_1(\mathbf{x}_1) & \phi_2(\mathbf{x}_1) & \dots & \phi_p(\mathbf{x}_1) \\ \vdots & \vdots & & \\ \phi_1(\mathbf{x}_m) & \phi_2(\mathbf{x}_m) & \dots & \phi_p(\mathbf{x}_m) \end{bmatrix} \quad (2)$$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} \quad (3)$$

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 $E(f(\phi(\mathbf{x}), \mathbf{w}), \mathcal{D})$

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

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,

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

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

$$y \in \{0, 1\} \quad f(\phi(\mathbf{x}), \mathbf{w}) = \log \frac{\Pr(y|\mathbf{x})}{1 - \Pr(y|\mathbf{x})} = b + \mathbf{w}^T * \phi(\mathbf{x}) \quad (5)$$

- 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

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

$\downarrow f(x) = \omega^T \phi(x) + b$

- The least square solution for linear regression is obtained as

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

- The minimum value of the squared loss is zero
- If zero were attained at w^* , we would have

$$\Phi(w^*) = y$$

- 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

$$\phi = \begin{bmatrix} \phi_1(x_1) & \dots & \phi_p(x_1) \\ \dots & \dots & \dots \\ \phi_1(x_m) & \dots & \phi_p(x_m) \end{bmatrix}$$

and

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

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

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

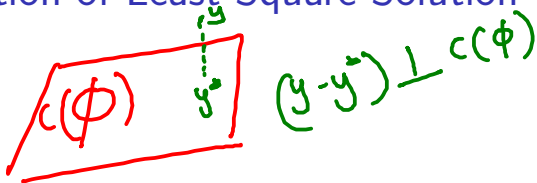
When?

ϕ is not full column rank

\Rightarrow ① Because $n > m$
(too many irrelevant attributes)

② Redundancy in ϕ

Geometric Interpretation of Least Square Solution



- Let y^* be a solution in the column space of ϕ
- The least squares solution is such that the distance between y^* and y is minimized
- Therefore.....

Geometric Interpretation of Least Square Solution

- Let \mathbf{y}^* be a solution in the column space of ϕ
- 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

$$\phi \mathbf{w} = \mathbf{y}^* \quad (8)$$

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

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

Recap: $\tilde{y} \in C(\phi) \dots$ Let $\tilde{y} = \phi w$

$$(\phi w)^T \phi = y^T \phi \quad (11)$$

$$w^T \phi^T \phi = y^T \phi \quad (12)$$

$$\phi^T \phi w = \phi^T y \quad (13)$$

Assumes

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

- Here $\phi^T \phi$ is invertible only if ϕ has full column rank

Proof?

Theorem : $\phi^T \phi$ is invertible if and only if ϕ is full column rank

Proof :

Given that ϕ 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 $\phi^T \phi \mathbf{x} = \mathbf{0}$

$\Rightarrow \phi^T \phi$ is not full row/column rank.

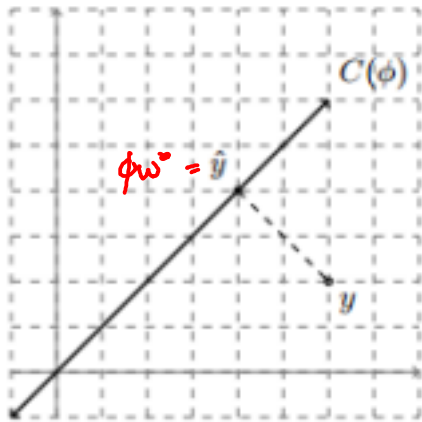
$$\Rightarrow \mathbf{x}^T \phi^T \phi \mathbf{x} = 0$$

$$\Rightarrow (\phi \mathbf{x})^T \phi \mathbf{x} = 0 \Rightarrow \|\phi \mathbf{x}\|_2^2 = 0$$

$$\Rightarrow \phi \mathbf{x} = \mathbf{0}$$

This is a contradiction. Hence $\phi^T \phi$ is invertible if ϕ 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 ϕ has full column rank if $\phi^T \phi$ is invertible. Hence, theorem proved**



if $(\phi^T \phi)^{-1}$
exists
then
 $w^* = (\phi^T \phi)^{-1} \phi^T y$

Figure: Least square solution y^* is the orthogonal projection of y onto column space of ϕ

What is Next?

- 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

Ridge regression
Lasso regression