

Linear Regression & Least Squares



• Learning Problem :-

- Given input-output pairs, we need to learn a function, f , that maps input to the output.

• Components & Terminology:-

X : input matrix with multiple variables
 $X_i \in \mathbb{R}^p$

y : output vector, dependent variable.

$\text{Col}(X)$: is the vector space that contains the input vectors

B : Matrix of the trainable parameters

f : the function that can be used to estimate the values of y based on X & B .

• Linear Regression:-

- Fit a line (represented by f) on the input data.
- is a common technique for many apps in econometrics, Genomics, Ecology, engineering, etc.

$$y = mX + b$$

Diagram illustrating the components of the linear regression equation $y = mX + b$:

- y : predicted
- m : slope
- X : input vars
- b : bias / Intercept

- In the training time, (X, y) pairs are available, so we estimate parameters

- In the inference time, we estimate \hat{y} based on input x & parameters $\{m, b\}$

Q: How to estimate parameters β , where they represent the best model fitting data?

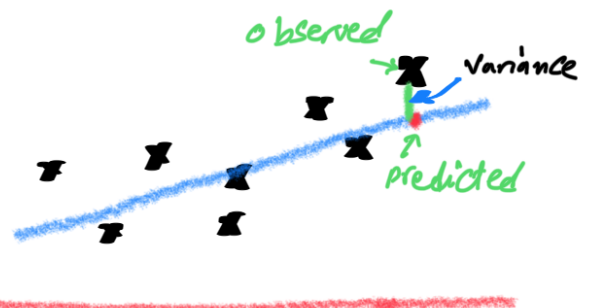
A: We use optimization techniques, such as Least Squares or Gradient Descent.

- understanding Least Squares Method

- objective:-

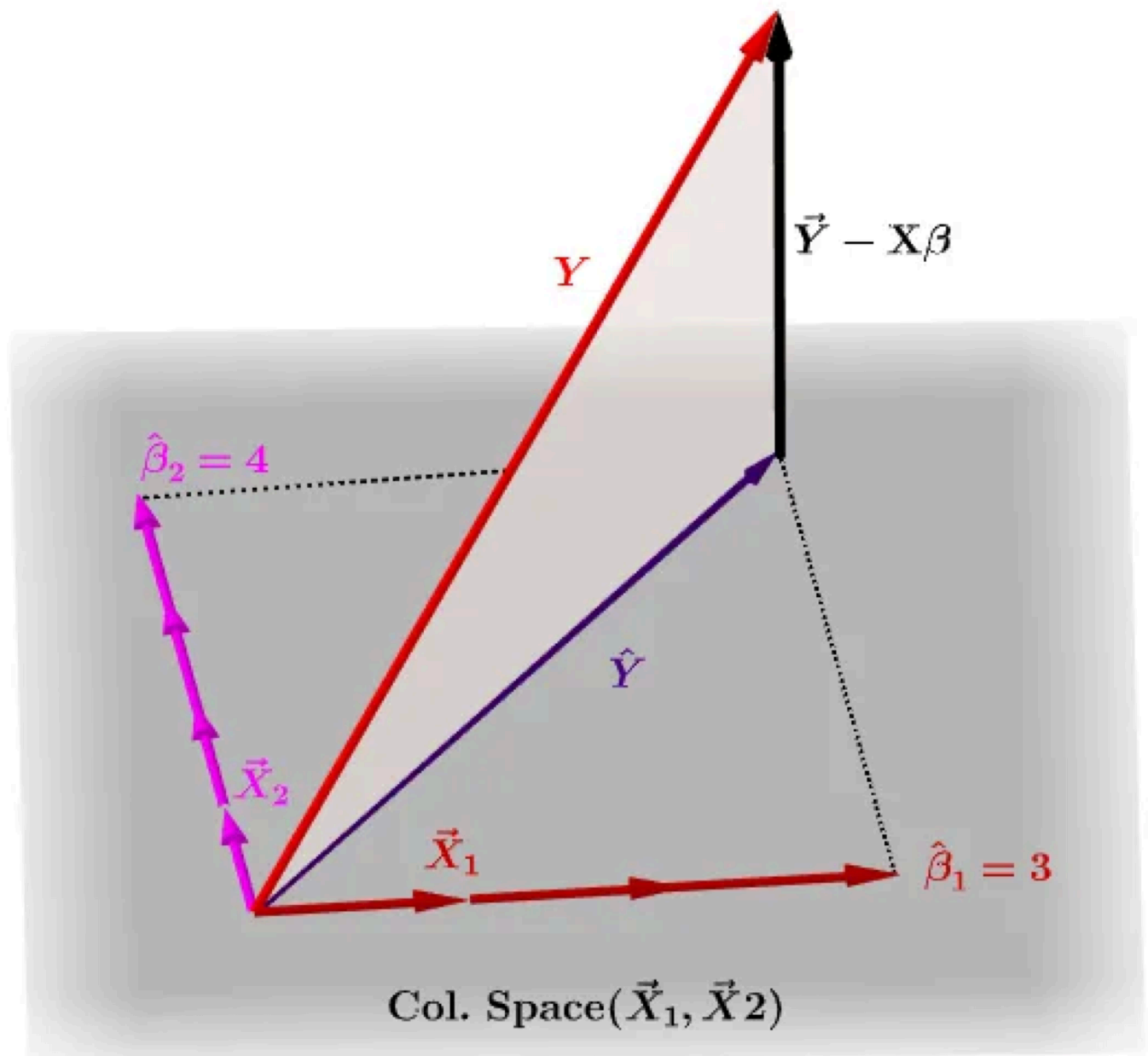
Find the function, (in this case line), that best fits the input data, and enable good estimate

of output variables. This can be done by minimizing the error / variance between estimated / predicted values and observed y values.



- using "Least Squares"

- Geometric Interpretation



The idea here is that the error vector is orthogonal onto the column space of the input.

Thus estimating the optimum values of $\hat{\beta}$ is possible. ← Good news !!!

- Mathematical Representation.

Given input $X = (X_1, X_2, \dots, X_N)^T$, and output $y = (y_1, y_2, \dots, y_N)$, we predict \hat{y} as:

$$\hat{y} = \hat{B}_0 + \sum_{j=1}^N X_j \hat{B}_j$$

Predicted
Parameters
Predictors

- By including B_0 into the matrix of the

learned parameters, we have

$$\hat{y} = X^T \hat{B}$$

- Given Pairs of X & y , we can estimate \hat{B} as follows

$$RSS(\hat{B}) = \sum_{i=1}^N (y_i - X_i^T \hat{B})^2$$

(Residual sum of squares) →

We minimize this function by taking the Partial derivatives w.r.t. the individual parameters, B_0, B_1, B_2, \dots

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