

Generalized vs Regularized Regression

G.L.R. :- Refers to conventional linear models for a continuous response variable given continuous/categorical predictors.
 e.g. → Simple Linear Regression,
 → Binary Logistic Regression.

Regularized Regression :-

Are type of regression where the coefficient estimates are constrained to zero. The size of Go efficient and the size of less terms are penalized.
 Complex modes are unwarranted due to overfitting.

Types of Regularizers

- ① L1
- ② L2
- ③ Elastic Net

Recap!

Bias :- Biases are the underlying assumptions that our model makes.

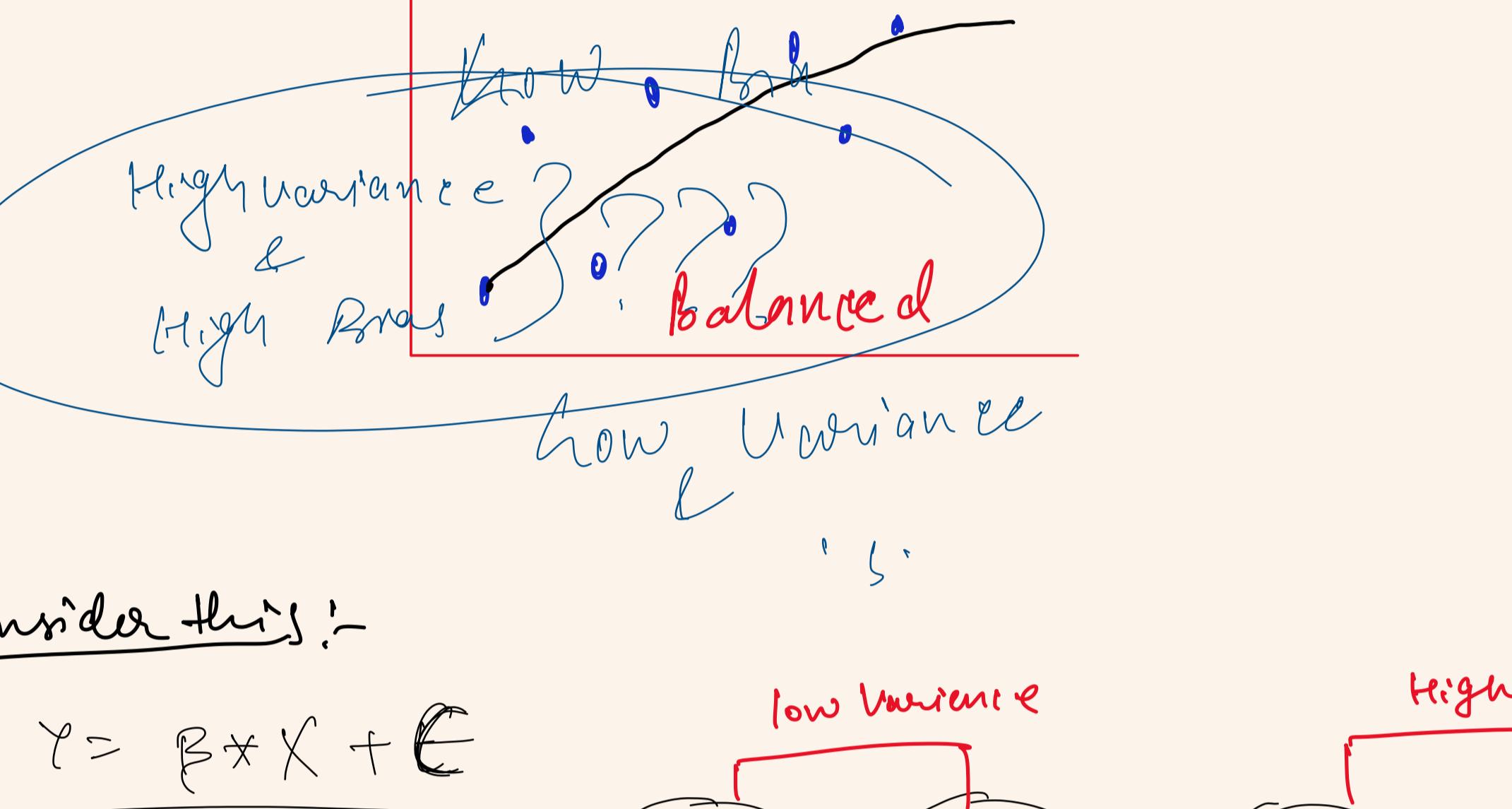
↓ Training time due to low sensitivity
 ↓ Training time due to low sensitivity

May lead to underfitting

Variance :-

→ Happens due to model's sensitivity to small fluctuations of data.

→ Commonly referred as overfitting
 → The model basically learns every single point and doesn't do well on novel dataset that it hadn't seen before.

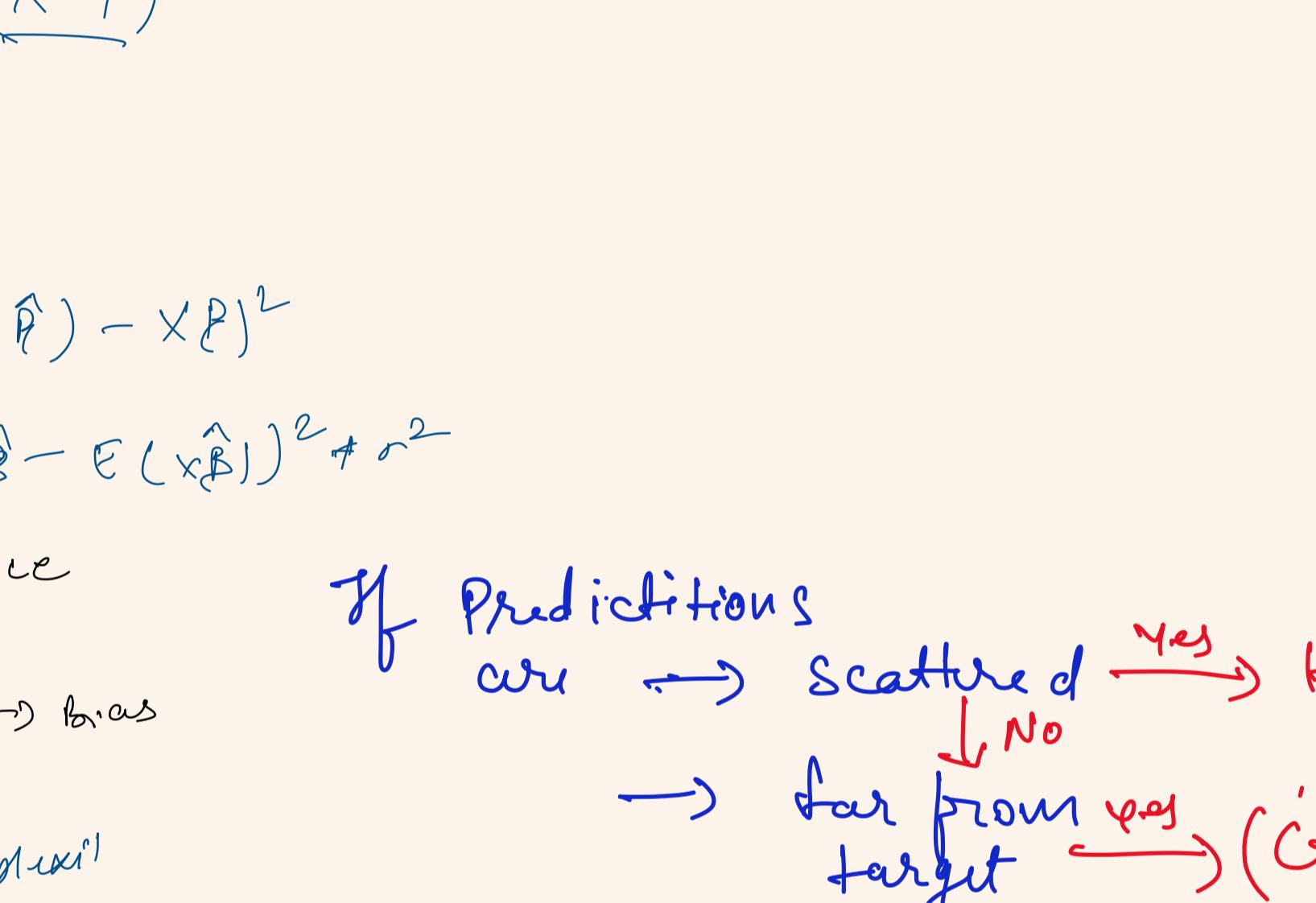


Consider this :-

$$Y = \beta * X + \epsilon$$

$$\hat{Y} = \hat{\beta} * X + \hat{\epsilon}$$

Goal is to get $\hat{\beta}$ that produces lower RSS.



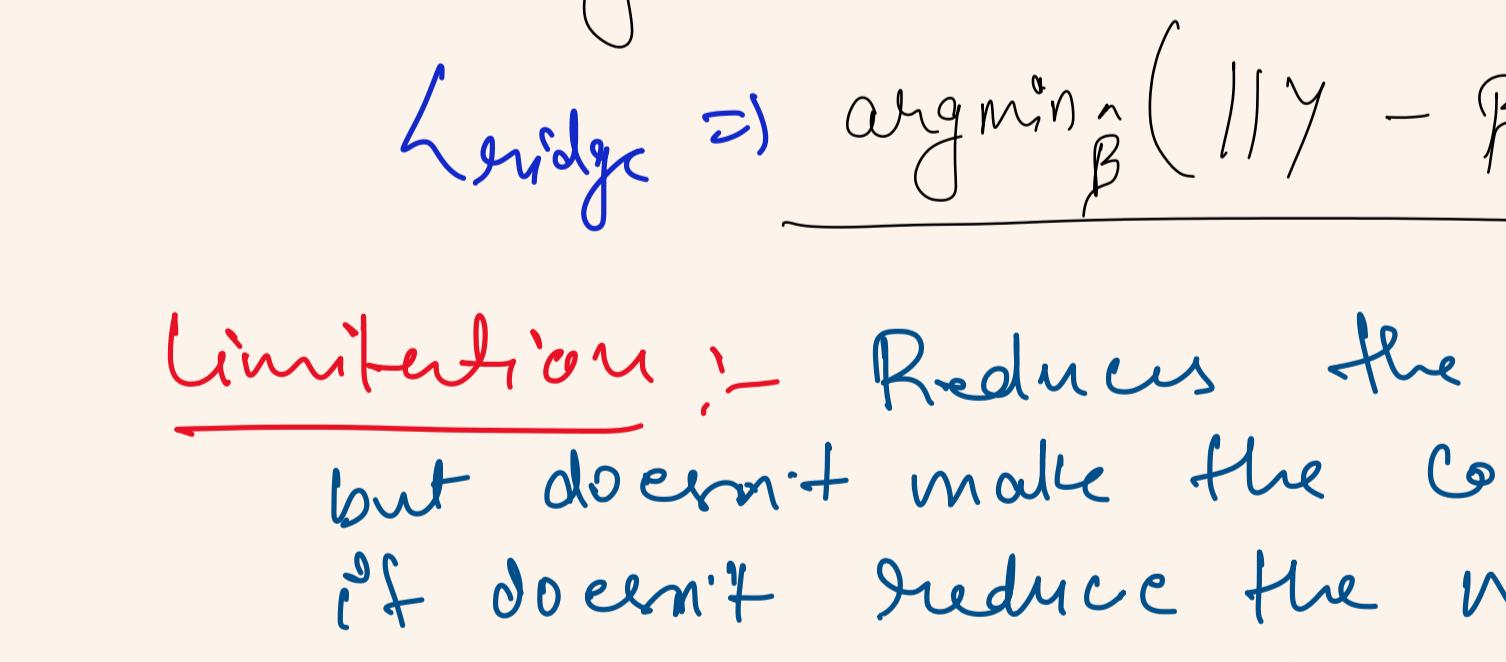
$$L(\hat{\beta}) = \sum_{i=0}^n || Y_i - \hat{\beta}_i * \hat{\beta} ||^2 = || Y - X\hat{\beta} ||^2$$

$$(\hat{\beta}) = (X^T X)^{-1} (X^T Y)$$

$$\text{Bias} = E(\hat{\beta}) - \beta$$

$$\text{Variance} = \sigma^2 (X^T X)^{-1}$$

$$\text{Error-term} = (E(X\hat{\beta}) - X\beta)^2$$



If Predictions are → scattered → High Variance
 ↓ No → far from target → High Bias
 (Crappy Model)

Ridge Regression :- (L^2)

→ We add penalty term which is equal to the square of the coefficient.

The L^2 term is equal to the square of the magnitude of the coefficients. we also add λ to control the penalty.

$$\text{Ridge} \Rightarrow \underset{\beta}{\operatorname{argmin}} (||Y - \beta * X||^2 + \lambda * \beta^2)$$

Limitations :- Reduces the model complexity but doesn't make the coeff. zero meaning if doesn't reduce the number of parameters. (Not good for feature reduction)

Lasso Regression :- (L_1)

(Least abs. Shrinkage and Selection Operator.)

→ It adds penalty term to the cost function.

→ This term is the absolute sum of coefficient.

$$L_{\text{lasso}} = \underset{\beta}{\operatorname{argmin}} (||Y - \beta * X||^2 + \lambda * \beta_1)$$

Can make coefficient '0' hence suitable for feature reduction.

Limitations :- Sometimes struggles with some specific type of data where no. of predictors > no. of observations.

→ If there are two highly collinear variables, lasso will select 1 randomly which is bad for interpretability.

Elastic Net :-

→ Sometimes, lasso produces bias in model where the prediction is too dependent on 1 variable.

In that case we use elastic Net because it contains the best of both.