

Linear Regression / (Bias - Variance)

→ $y/x \sim N(\theta^T x, \sigma^2)$

→ $\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \operatorname{RSS}(\theta) \quad ; \quad \operatorname{RSS}(\theta) = \|y - X\theta\|^2$
$$= \sum_{i=1}^n (y^{(i)} - x^{(i)T} \theta)^2$$

⇒ Equivalent to maximum conditional likelihood Estimator

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \log(P(y^{(i)}/x^{(i)}))$$

→ If the Residual plot has a pattern, then the model is not right.
→ linear transformation may be used.

i.e. y is linear on x , but x could be non-linear transformation of original data x'

→ ^{LT} Linear Transformation will not impact the linear relationship b/w Variables.
$$\operatorname{Corr}(y, x) = \operatorname{Corr}(y, \operatorname{LT}(x))$$

Non-Linear Transformation may increase/decrease the linear relationship

→ Plot Residual plots to verify Non-linear Transformation helps given relationship.

→ Training Data

$$(X^{(i)}, y^{(i)}) \stackrel{iid}{\sim} p(x, y) \\ = \underbrace{p(x)}_{\text{arbitrary model}} \underbrace{p(y/x)}_{\text{linear regression model}}$$

→ Gradient Descent in Linear Regression - is used cause ~~it~~ its computationally cheaper.

Normal Eqⁿ - has matrix inversion $O(n^3)$
↳ very Expensive.

$$\begin{aligned} \rightarrow E(\hat{\theta}/x) &= (X^T X)^{-1} X^T E(\underbrace{y/x}_{X\theta}) \\ \hat{\theta} &= (X^T X)^{-1} X^T y \\ &= (X^T X)^{-1} (X^T X) \theta = \theta \end{aligned}$$

$$\text{Var}(\hat{\theta}/x) = \sigma^2 (X^T X)^{-1}$$

$$\left[\hat{\theta}/x \sim N(\theta, \sigma^2 (X^T X)^{-1}) \right]$$

Link [towardsdatascience.com/mse and bias-variance decomposition]
 $y = f(x) + \epsilon$

Bias-Variance Decomposition

$$MSE = E[(y - \hat{f}_S(x))^2]$$

Expected Value Squared ^{prediction} Errors

$$= \underbrace{\text{Var}(f(x) - \hat{f}_S(x))}_{\text{training data } S} + \underbrace{\text{Var}(\epsilon)}_{\text{Irreducible noise}} + (E[f(x)] - E[\hat{f}_S(x)])^2$$

Variances is the gap
 Has jumps is the
 the predicts real model & the
 on the training data
 S and test data (x, y)

Bias

Shows whether our
 model approximates the
 real model well.

i.e. if model passes through
 all the points then it has
 low Bias.

Irreducible
 noise

Watch
 Regularization
 DVA.

features \rightarrow samples.

\rightarrow if $D \ll N$, MLE overfits

\rightarrow MLE of linear regression Estimator $= 0$

$$E[\hat{\theta}] = 0$$

Bias/Intercept different from statistical bias (DL, 110 page)

$$E[\alpha] - \alpha$$

Bias in Linear Regression

- Adds flexibility to the model

- i.e. without bias the line has to pass through origin, limiting or constraining the model.

Ex: if all features $= 0$, then prediction must be zero irrespective of the data.

bias parameter of affine transformation

$$y \sim N(XW, \sigma^2 I)$$

for new (x_0, y_0) $y_0 \sim N(x_0^T W, \sigma^2)$

$$E \left[(y_0 - x_0^T \hat{w})^2 / X, x_0 \right]$$

Squared Bias

different from bias/intercept

Statistical Bias

$$= \underbrace{\sigma^2}_{\text{noise}} + x_0^T (W - E[\hat{w}]) (W - E[\hat{w}])^T x_0$$

Model Bias: how close to the solution we expect to be on average.

$$+ x_0^T \text{Var}(\hat{w}) x_0$$

Variance

how sensitive is our solution to the data.



LS Solution:-

unbiased, but potentially high variance

Ridge Regression:-

(RR)

biased, but lower variance than LS

Which is preferable?

Ultimately depends on how well our model generalizes the data.

→ Using proxy $E[(y_0 - x_0^T \hat{w})^2 / x_0]$

LS \hat{w} Bias = 0

In RR, we add Bias to lower variance.

Squared Bias = $x_0^T (E[\hat{w}] - E[w])$

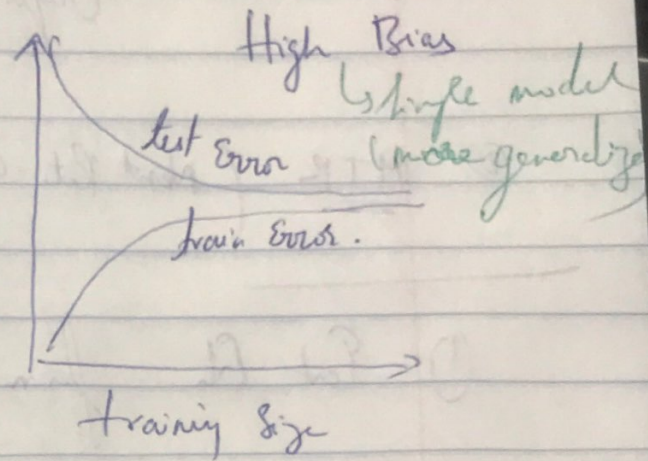
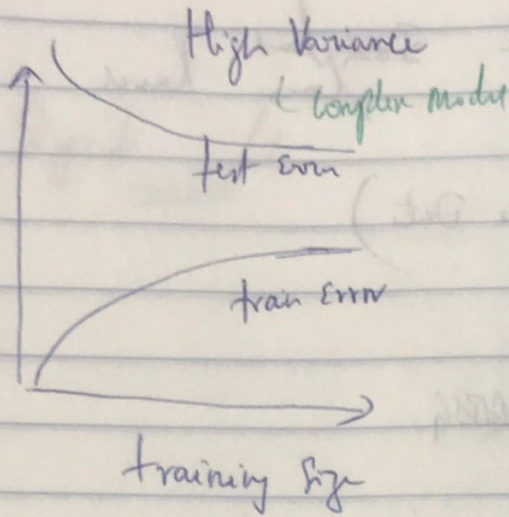
$$= x_0^T \left(\overset{w}{E[\hat{w}]} - E[w] \right) (w - E[w])^T x_0$$

In RR, we don't w (true weights), can't calculate Bias.

But we could make some statements about what w should be?

Ch ~~Learning Curves~~ (output)

→



↪

Difference b/w train error & test error
~~is~~