

* Bias Variance decomposition

We try to decompose the generalized error and understand what leads to the error. What terms exactly contribute to that error.

A parameter is some quantity about a distribution we are interested in. We estimate this parameter θ using an estimator $\hat{\theta}$.

$$(1) \text{ Bias } (\hat{\theta}, \theta) = E(\hat{\theta}) - \theta$$

What the bias measures is how off is the estimator from being a true estimator. This is the error we will still get, even if we train with infinite data. This is due to classifier being biased to a particular type of solution.

For eg. Linear classifiers have a bias that the features and labels are linearly related.

In other words, bias is inherent to the model & it does not depend on data as such.

$$(2) \text{ Variance } (\hat{\theta}) = E(\hat{\theta} - E(\hat{\theta}))^2$$

What variance computes is how much do the classifiers vary in general. Thus it measures how different our classifiers become if we change the training data.

For eg. Non linear classifiers get affected even by slight change of data. Thus, they depict high variance. Variance is independent of true prediction. It just indicates how much I will vary.

Thus,

D_i are different training data subsets drawn from $P(x, y)$.

ML
Algorithm
 \hat{A}

$E(\hat{A})$

$D_1 \quad D_2 \quad D_3 \quad D_4 \quad D_5 \quad \theta$

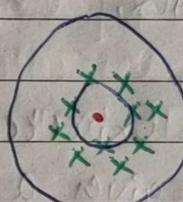
\leftarrow variance \rightarrow

bias

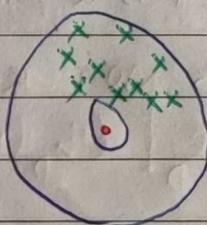
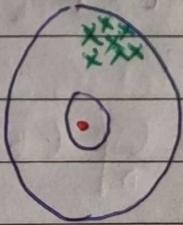
LOW
VARIANCE

HIGH
VARIANCE

LOW
BIAS



HIGH
BIAS



∴ VARIANCE MEASURES SPREAD while BIASES MEASURE ACCURACY.

This BIASES-VARIANCE TRADEOFF is what we need to consider while choosing a Model. This is what tells us that Non-linear models are NOT always better than linear models. It is the tradeoff over which a Model is selected.

Consider a distribution $P(x, y)$ where $y \in R$; i.e. Regression setting. For a given x ; we do NOT have a single y ; but probability distribution over y . This is because for e.g. two similar houses can be sold for different price. This is maybe because of dataset not capturing all features or inherent randomness.

$$\therefore P(x, y) = \underset{\text{Joint probability}}{\underset{\uparrow}{P(y|x)}} \underset{\text{Conditional probability}}{\underset{\uparrow}{P(x)}} \underset{\text{Marginal probability}}{\underset{\uparrow}{P(y)}}$$

Now,

$$(1) \text{ Expected label } (\bar{y}) : \bar{y}(x) = E_{y|x} [y]$$

(given a $x \in R^d$)

Note: If y is a random variable, $E[y]$

Or \bar{y} is NOT R.V.

$$= \int_y y P(y|x) dy$$

$$(2) \text{ Expected Test error} : E_{\substack{x \in \text{TEST} \\ y}} \left[(h_D(x) - y)^2 \right]$$

(given a classifier h_D)

test points

Squared loss

classifier obtained
from algorithm π
& dataset D

$$= \int_x \int_y (h_D(x) - y)^2 P(y|x) dy dx$$

$$h_D = \pi(D)$$

for all x ; you find the most possible y at expected y & then find test error over all x .

(3) Expected classifier ($\bar{h}(x)$) : $E_{D \sim P^n} [f(D)]$
 (given Algorithm f)

$$= \int h_D P(D) dD$$

This is nothing
but $E(\hat{\theta})$.

(4) Expected generalization : $E_{\substack{x, y \sim P \\ D \sim P^n}} [(h_D(x) - y)^2]$
 (error of algorithm f)

$$= \iiint_{\substack{D \\ Y \\ X}} [(h_D(x) - y)^2] \times P(y|x) \times P(D) dx dy dD$$

This is nothing but **expected square loss** of the **Algorithm f** . Hence, this is what we want to decompose & understand.

Consider the error at a fixed x , i.e.

(we will introduce)
 E_x in the end

$$= E_{Y,D} [(h_D(x) - y)^2]$$

Note: \bar{y} is NOT
a random variable

$$= E_{Y,D} [((h_D(x) - \bar{y}) + (\bar{y} - y))^2]$$

This term does
not depend
on D

$$= E_{\substack{Y \\ \text{fixed } D}} [(h_D(x) - \bar{y})^2] + E_Y [(\bar{y} - y)^2]$$

This term does NOT
depend on y

$$+ 2 E_{Y,D} [(h_D(x) - \bar{y})(\bar{y} - y)]$$

— (I)

However; if we look at third term;

$$E_{Y,D} \left[(h_D(x) - \bar{Y})(\bar{Y} - Y) \right] = E_D \left[(h_D(x) - \bar{Y}) \right] \times E_Y \left[(\bar{Y} - Y) \right]$$

$$\text{Now; } E_Y(\bar{Y} - Y) = \bar{Y} - E_Y(Y) = \bar{Y} - \bar{Y} = 0.$$

∴ Third term is $E(0) = 0$. — (II)

Using (I) & (II);

$$E_D \left[(h_D(x) - \bar{Y})^2 \right] + E_Y \left[(\bar{Y} - Y)^2 \right] — (III)$$

Just looking at first term of (III);

$$E_D \left[(h_D(x) - \bar{Y})^2 \right] = E_D \left[((h_D(x) - \bar{h}(x)) + (\bar{h}(x) - \bar{Y}))^2 \right]$$

$$= E_D \left[(h_D(x) - \bar{h}(x))^2 \right] + E_D \left[(\bar{h}(x) - \bar{Y})^2 \right]$$

$$+ 2 E_D \left[(h_D(x) - \bar{h}(x)) (\bar{h}(x) - \bar{Y}) \right]$$

This is 0, similar to above

$$= E_D \left[(h_D(x) - \bar{h}(x))^2 \right] + E_D \left[(\bar{h}(x) - \bar{Y})^2 \right]$$

— (IV)

Resulting (IV) in (III);

Note: if x is R.V, $E(x)$ is NOT R.V. Thus; $E_D[(\bar{h}(x) - \bar{y})^2]$ is $(\bar{h}(x) - \bar{y})$.

$$= E_D[(h_D(x) - \bar{h}(x))^2] + E_D[(\bar{h}(x) - \bar{y})^2]$$

$$+ E_{\bar{y}}[(\bar{y} - y)^2]$$

$$= E_D[(h_D(x) - \bar{h}(x))^2] + (\bar{h}(x) - \bar{y})^2 + E_{\bar{y}}[(\bar{y} - y)^2]$$

VARIANCE + BIAS + NOISE

This term just tells us how varied our classifiers are from the mean classifier. It does NOT depend on y .

This term tells us how off our classifier is from the true prediction and is squared. It is not dependent on D ALREADY.

This term tells us how off our true prediction is from the true value. Even bayes optimal classifier suffers from noise.

NOW:

$$E_{x,y,D}[(h_D(x) - y)^2] = E_{x,D}[(h_D(x) - \bar{h}(x))^2] \text{ Variance}$$

$$+ E_x[(\bar{h}(x) - \bar{y})^2] \text{ Bias}^2$$

$$+ E_{x,y}[(\bar{y} - y)^2] \text{ Noise}$$

In other words :

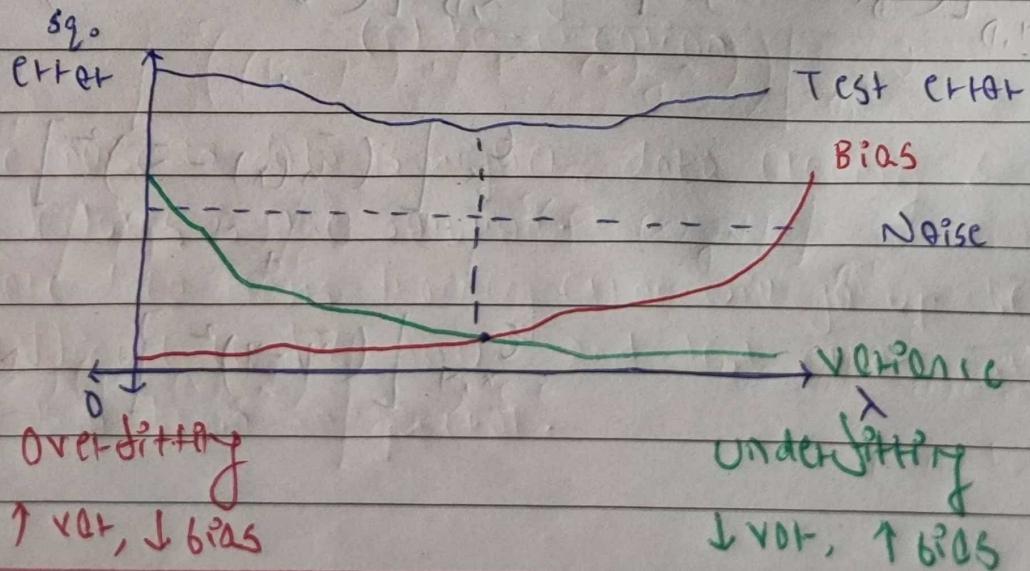
$$E_{x,y,D} \left[(h_D(x) - y)^2 \right] = E_x \left[\text{var}_{y|x}(y) + \text{var}_D(h_D(x)) \right. \\ \left. + \text{bias}(h_D(x))^2 \right]$$

The tradeoff here is that we want to choose h_D to balance between second & third term to reduce MSE. We can NOT just minimize one or other, but we can find a balance between these terms.

For eg. At times, increasing bias by little (for eg. regularization) allows us to substantially drop variance.

The variance term contains the wigginess, hence might prefer a simple function that yield predictions which are too varied. The bias term contains how close the average model is to truth. we need to pay attention to data to reduce the bias term.

Introducing structures like regularization allows us to tune bias-variance tradeoff.



To choose the best balance point between bias & variance, we use methods such as K-fold cross validation or early stopping.

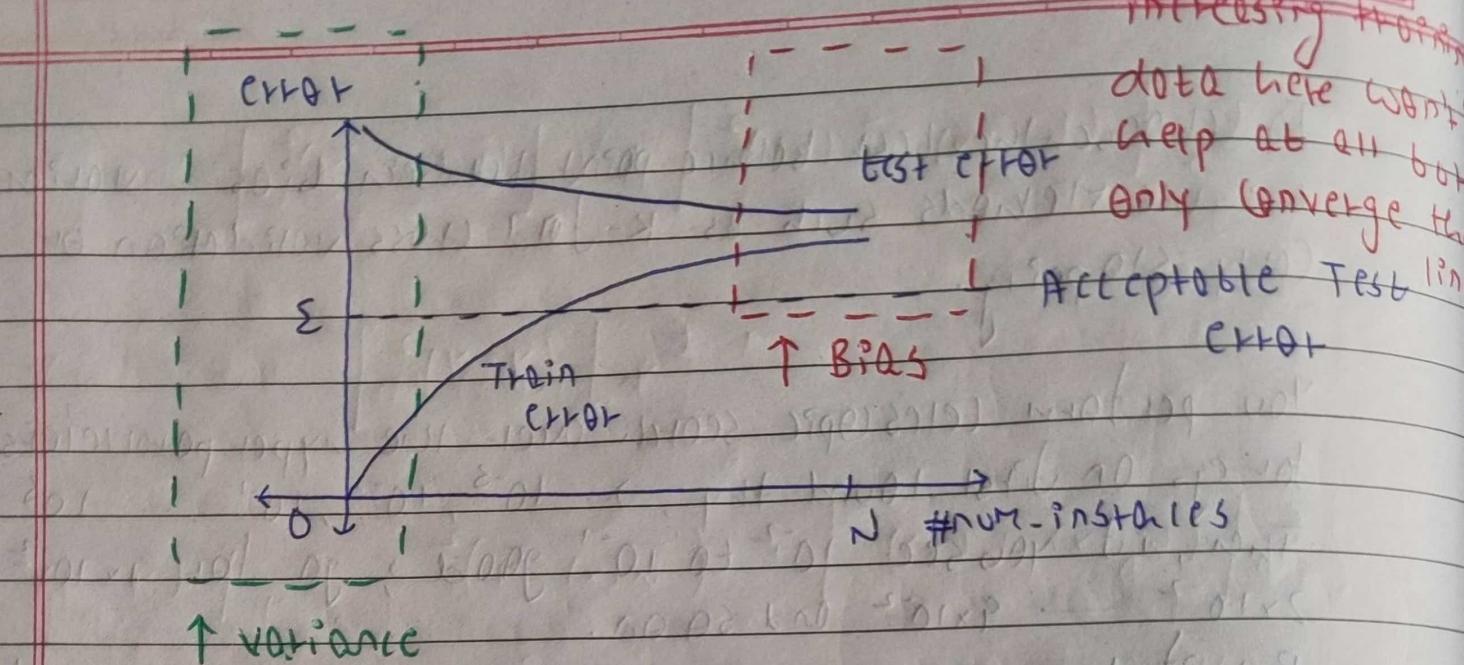
You perform telescopic search over the hyperparameters. For eg. On first round try, $\lambda = 10^{-3}, 10^{-2}, 10^{-1} \dots 10^3$. Then, if you see 10^{-2} to 10^{-1} good, go for $1 \times 10^{-2}, 2 \times 10^{-2} \dots 9 \times 10^{-2}$ and so on.

But if tune too much on validation set, we will overfit on validation set. This is because if you tune too hard to get accuracy from 80 to 90%, that would mean you tuned too hard to get 4-5 examples right, cause validation set is small! That's why we do **K-fold cross validation**.

Here; Training data is split in **K buckets** and for **K times** you can take **K-1 buckets to train & 1 bucket to validate**. The larger **K** is better, but at a time it takes more time. So if $K=N$; that is leave out one (**Cross Validation**).

Now when we get high error, we can estimate where the bias exactly is. If the variance of model is low, reducing variance won't help.

One idea for ML debugging is to **plot error vs num instances**. Here, we start from 1 instance in training data to all N instances. We train the model every time & plot **train error & test error**. Σ is the error limit we want by the model.



If (test error is $< \Sigma$) ; NO problem
 Else if (test error $> \Sigma$) & (train error $< \Sigma$) ;

This is ↑ variance OR Overfitting

- Solutions : (1) Add more training data
- (2) Bagging
- (3) Reduce model complexity
- (4) ↓ regularization

Else (test error $> \Sigma$) & (train error $> \Sigma$) ;

This is ↑ Bias OR Underfitting

- Solutions : (1) Use more complex model (non-linear)
- (2) Add more features
- (3) Boosting
- (4) ↑ regularization
- (5) Kernelization

If the noise in data is high, data cleaning & adding more features is the solution.

Also, more data means training error will go up because it's more difficult to fit new. And test error will be greater than train in general.