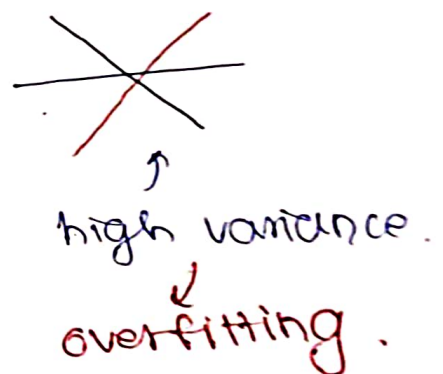
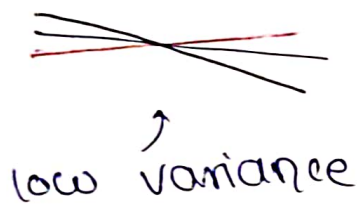


\* Regularization part 1 :-

#### 4 Bias - Var Trade off +

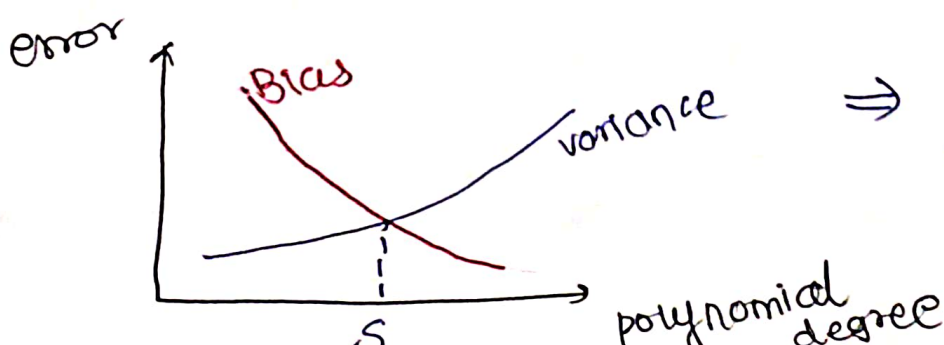
Reducible error  $(f(x) - f'(x))$

- Bias :- the inability of a ML model to fit the training data.  
(Deviation from actual relation)
- Variance :- how much the ML model changes upon changing the training data.



→ high bias  $\Rightarrow$  underfitting

Bias  $\propto \frac{1}{\text{Variance}}$



⇒ target is to find this S.

- Expected value :- (pop. mean)  
 It represents the average outcome of a random variable over a large number of trials.

$$E[X] = x_1 p(x_1) + x_2 p(x_2) + \dots + x_n p(x_n)$$

discrete RV

Cont. RV  $\rightarrow E[X] = \int x_i p(x_i) dx$

pop. of variance.

$$\text{Var}[X] = E[X^2] - (E[X])^2$$

$$= E[(X - E[X])^2]$$

- what are bias & variance mathematically:

$\rightarrow$  Bias :- It refers to systematic error that a model introduces bcoz it cannot capture the true relationship in the data.

It represents the diff. b/w expected prediction of our model and the correct value which we are trying to predict.

$$\text{Actual relation} = f(x)$$

$$\text{model relation} = f'(x) \rightarrow \text{for a sample}$$

Suppose, we take multiple samples and for each we have a specific  $f'(x)$ , so we take the expected value  $E[f'(x)]$

Note,

$$\text{Bias} = E[f'(x)] - f(x)$$

→ This can never be calculated as we never know  $f(x)$ , it's just an interpretation.

→ Variance :-

↳ refers to amount by which the prediction of our model will change if we used a diff. training data.

$$\text{Var} = E[(f'(x) - E[f'(x)])^2]$$

→ using Ridge:

```
from mixend.evaluate import  
    bias-variance-decomp  
alphas = np.linspace(0, 30, 100)  
loss = []  
bias = []  
variance = []  
for i in alphas:  
    reg = Ridge(alpha=i)  
    avg-expected-loss, avg-bias, avg-var =  
    bias-variance-decomp(reg, X-train, y-train,  
        X-test, y-test, loss = 'mse',  
        random-seed = 123)
```

## \* Regularization Part 2 ∴

### • Bias Variance Decomposition:-

$$\text{Loss} = \text{bias} + \text{variance} + \text{irreducible error}$$

→ Derivation:-

$$\text{mse} = \underbrace{\frac{\sum (y - \hat{y})^2}{n}}_{\text{Expected value.}} \quad \left\{ \begin{array}{l} \text{Assumption} = E(\epsilon) = 0 \\ V[\epsilon] = \sigma^2 \end{array} \right.$$

$$y = f(x) + \epsilon = \theta + \epsilon$$

$$\hat{y} = f'(x) = \hat{\theta}$$

$$\text{mse} = E[(y - \hat{y})^2]$$

$$= E[(\theta + \epsilon - \hat{\theta})^2]$$

(considering  $\epsilon = a$   
 $\theta - \hat{\theta} = b$ )

$$= E[(\theta - \hat{\theta})^2 + \epsilon^2 + 2\epsilon(\theta - \hat{\theta})]$$

$$= E[(\theta - \hat{\theta})^2] + E[\epsilon^2] + 2E[\epsilon](\theta - \hat{\theta})$$

$$\text{mse} = E[(\theta - \hat{\theta})^2] + E[\epsilon^2]$$

we know,  $V[\epsilon] = \sigma^2 = E[(\epsilon - E[\epsilon])^2]$   
 $= E[\epsilon^2]$

$$\therefore \text{mse} = E[(\theta - \hat{\theta})^2] + \text{Var}(\epsilon)$$



$$\begin{aligned}
 E[(\theta - \hat{\theta})^2] &= E[(\theta - E[\hat{\theta}] + E[\hat{\theta}] - \hat{\theta})^2] \\
 &= E[(\theta - E[\hat{\theta}])^2 + (E[\hat{\theta}] - \hat{\theta})^2 \\
 &\quad + 2(\theta - E[\hat{\theta}])(E[\hat{\theta}] - \hat{\theta})] \\
 &= E[(\theta - E[\hat{\theta}])^2] + E[(E[\hat{\theta}] - \hat{\theta})^2] \\
 &\quad + E[2(\theta - E[\hat{\theta}])(E[\hat{\theta}] - \hat{\theta})] \\
 &\quad \quad \quad \downarrow \quad \quad \quad \downarrow \quad \quad \quad \downarrow \\
 &\quad \quad \quad 2 \quad \quad \quad \theta - E[\hat{\theta}] \quad \quad \quad 0 \\
 &\quad \quad \quad \downarrow \\
 &(\theta - E[\hat{\theta}])^2 \rightarrow \text{Bias}^2
 \end{aligned}$$

$\text{Var}(\hat{\theta})$

$$\therefore \text{Loss} = \underbrace{\text{Bias}^2 + \text{Variance}}_{\text{reducible}} + \underbrace{\text{Var}(\epsilon)}_{\text{irreducible}}$$

- Regularization  $\rightarrow$  Reduce variance.  
 $\downarrow$   
 reduce overfitting.

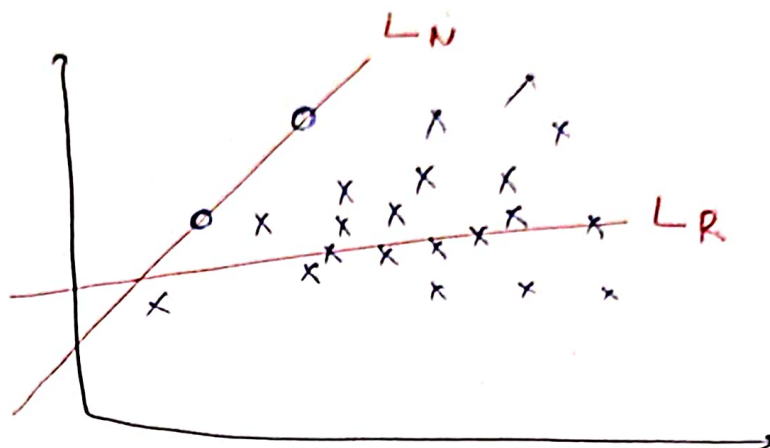
$\rightarrow$  Complex models are good way to reduce bias.

# \* Ridge Regression :-

→ L2 Regularization.

→ Let's consider a case where training data has only 2 points.

o → training points  
x → test points.



we can see, the best fit line according to the training data is not best for test.

∴ we need to convince the model to take  $L_R$  as best fit.

Regularization term

$$\Rightarrow L = \sum (y_i - \hat{y}_i)^2 + \lambda (m^2)$$

Bias slightly increases.

for multiple IR  
 $= \lambda (\sum m_i^2)$

L2

from sklearn.linear-model import Ridge  
 $R = \text{Ridge}(\alpha = 0.0001)$

$R.\text{fit}(X_{\text{train}}, y_{\text{train}})$

$y = R.\text{predict}(X_{\text{test}})$

$\lambda$

$\lambda \rightarrow$  less value  $\Rightarrow$  less reduction in slope

$\rightarrow$  high value  $\Rightarrow$  high redn in slope  
 $\rightarrow$  may result in underfitting.

Similar to the formula of 'm' in linear regression,

M for Ridge regression,

$$m = \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2 + \lambda}$$

Extra term

→ for multiple linear regression:

$$L = (XB - Y)^T (XB - Y) + \lambda B^T B$$

$$\beta = (X^T X + \lambda I)^{-1} X^T Y$$

• for Ridge regression with gradient descent:

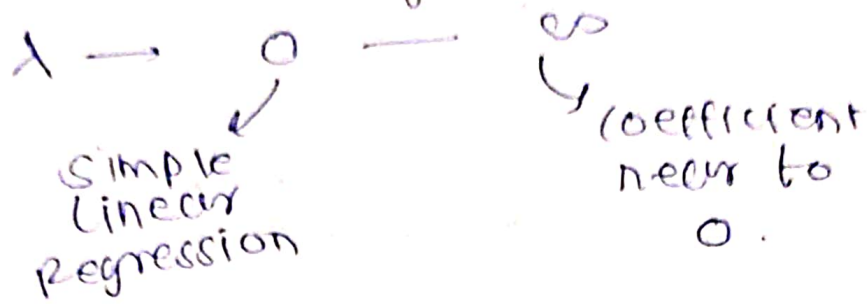
```
from sklearn.linear_model import SGDRegressor  
reg = SGDRegressor(penalty='l2', max_iter=5000,  
                    eta0=0.1, learning_rate='constant',  
                    alpha=0.001)
```

Now fit & predict

OR

use solver = 'sparse\_cg' in Ridge.

Q.1 How the coefficients get affected?



Higher values are affected more,  
at  $\alpha > 1000$ , all values come near  
to 0.

$\rightarrow$  for low  $\lambda \rightarrow$  Bias is low  
variance is high  
(overfit)

$\rightarrow$  for high  $\lambda \rightarrow$  Bias is high  
variance is low  
(underfit)

} Bias-variance trade off.

$\rightarrow$  Bias-var ke graph ke intersection se  
thode pehle wale point ki alpha value  
leni hai.



## \* Lasso Regression :-

✓

L1 Regularization

L1

$$L = \text{MSE} + \lambda \|m\|$$

In case of ridge regression, upon increasing the value of  $\lambda$ , corresponding coefficient goes near to zero, but in case of Lasso, the coefficient corresponding to variables which are not imp. for prediction becomes zero.  $\rightarrow$  feature selection.

from sklearn.linear-model import Lasso

reg = Lasso(alpha=)

reg.fit(X\_train, y\_train)

$$m = \frac{\sum (y_i - \bar{y})(x_i - \bar{x}) \pm \lambda}{\sum (x_i - \bar{x})^2}$$

for  $m > 0 \rightarrow -\lambda$

for  $m < 0 \rightarrow +\lambda$

for  $m = 0 \rightarrow$  no  $\lambda$  term

## \* ElasticNet Regression :-

↳ combination of L1 & L2 :

$$L = \sum (y_i - \hat{y}_i)^2 + a \|m\|^2 + b \|m\|$$

$$\lambda = a + b$$
$$l_1\text{-ratio} = \frac{a}{a+b}$$

} hyperparameters.

→ If there is multicollinearity in input cols then we should use ElasticNet