

* Bias Variance Tradeoff

Source: Coursera Andrew Ng
Mitesh M Khapra -
Deep Learning NPTEL Week 8

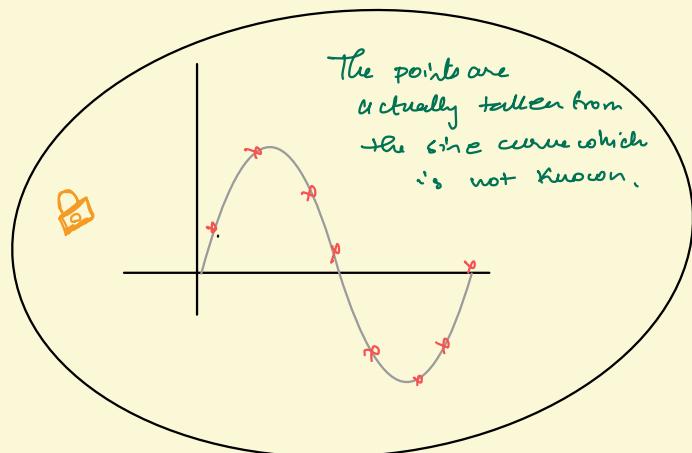
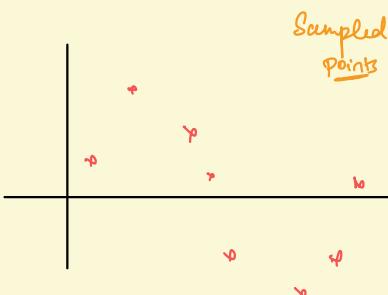
Overview of Bias - Variance:

Sundarayya Sarathi

Usual prob. of fitting a curve through a given set of points.

if a true relation b/w $x \& y$ which is unknown.

Consider the points ...



Two models

Linear model
Simple (deg 1)

i)

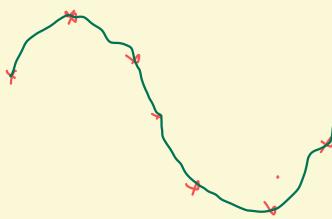
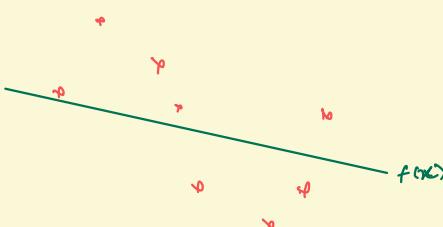
$$f(x) = w_1 x + w_0$$

How do we get this line?: By loss function

Polynomial model.
Complex (deg 25)

$$i) f(x) = \sum_{i=1}^{25} w_i x^i + w_0$$

a degree 25 polynomial obtained through gradient descent.



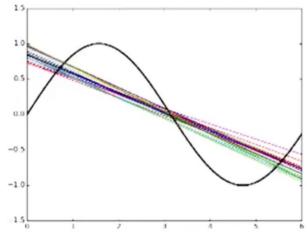
2) On average trying to minimize the distance of f from all points.

2) It's a complex model.

That's how we get f .

We sample 25 points from training data & train a simple & complex model and repeat the process k times to train multiple models.

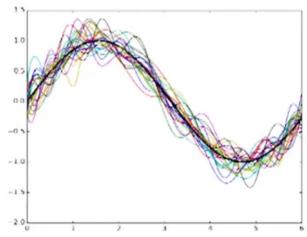
We won't get the same curve everytime since the training data is different each time, we get a different curve



(a)

- Simple models trained on different samples of the data do not differ much from each other
- However they are very far from the true sinusoidal curve (under fitting)
- On the other hand, complex models trained on different samples of the data are very different from each other (high variance)
- But they are closer to the true sinusoidal curve (Over fitting)

the different colors correspond to different curves obtained using different training samples



(b)

Bias = "has a very strong wrong preconception about the data".

① Bias: let $f(x)$ = true model
 $\hat{f}(x)$ = estimate of the model using a set of training points

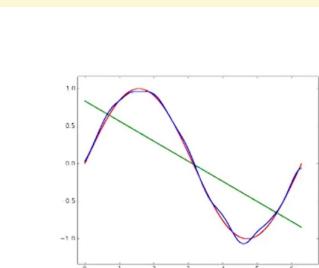
$$\text{Bias}(\hat{f}(x)) = E(\hat{f}(x)) - f(x)$$

$E(\hat{f}(x))$ = expectation over all the approximate functions obtained by training different data samples. (here, expectation of all those diff. colored lines)

i.e. It is the difference b/w true model and $E(\hat{f}(x))$

In (a) this difference is very high.

In (b) this difference is low.



Green Line: Average value of $\hat{f}(x)$ for the simple model

Blue Curve: Average value of $\hat{f}(x)$ for the complex model

Red Curve: True model ($f(x)$)

② Variance:

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

describes the spread of the different approx. models.

In (a) All models obtained using different data points are close to each other and not very far from each other. Hence variance is low

In (b) all the models are very different from each other. Spread is very high.
∴ Variance is high

In (a) : Simple model: low variance, high bias

In (b) : Complex : High variance, low bias.

Is it always the case?

• Train Error vs Test Error

Consider a new point (x, y) not seen during training.

Mean square error = $E[(y - \hat{f}(x))^2]$ (Averaging square error in predicting y for many such unseen points)

Why is $(y - \hat{f}(x))^2$ a Rand. variable?

$\because x$ feeding during test time is going to vary \therefore it's a rand. variable
D/p x is changing.

$$\star E[(y - \hat{f}(x))^2] = \text{Bias}^2 + \text{Variance} + \sigma^2 \text{ (lmm. error)}$$

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

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

$$\text{Variance} = \frac{1}{m} \sum_{i=1}^m (x_i - \bar{x})^2$$

$$\begin{aligned} \text{Var} &= E((x - E(x))^2) = E(x^2 + (E(x))^2 - 2x E(x)) \\ &= E(x^2) + (E(x))^2 - 2(E(x))^2 = E(x^2) - (E(x))^2 \end{aligned}$$

$$\begin{aligned} \text{Error} &= E((\hat{f} - f)^2) = E((f - E(\hat{f})) - (\hat{f} - E(\hat{f})))^2 \\ &= E((f - E(\hat{f}))^2) + E((\hat{f} - E(\hat{f}))^2) - 2 E((f - E(\hat{f}))(f - E(\hat{f}))) \end{aligned}$$

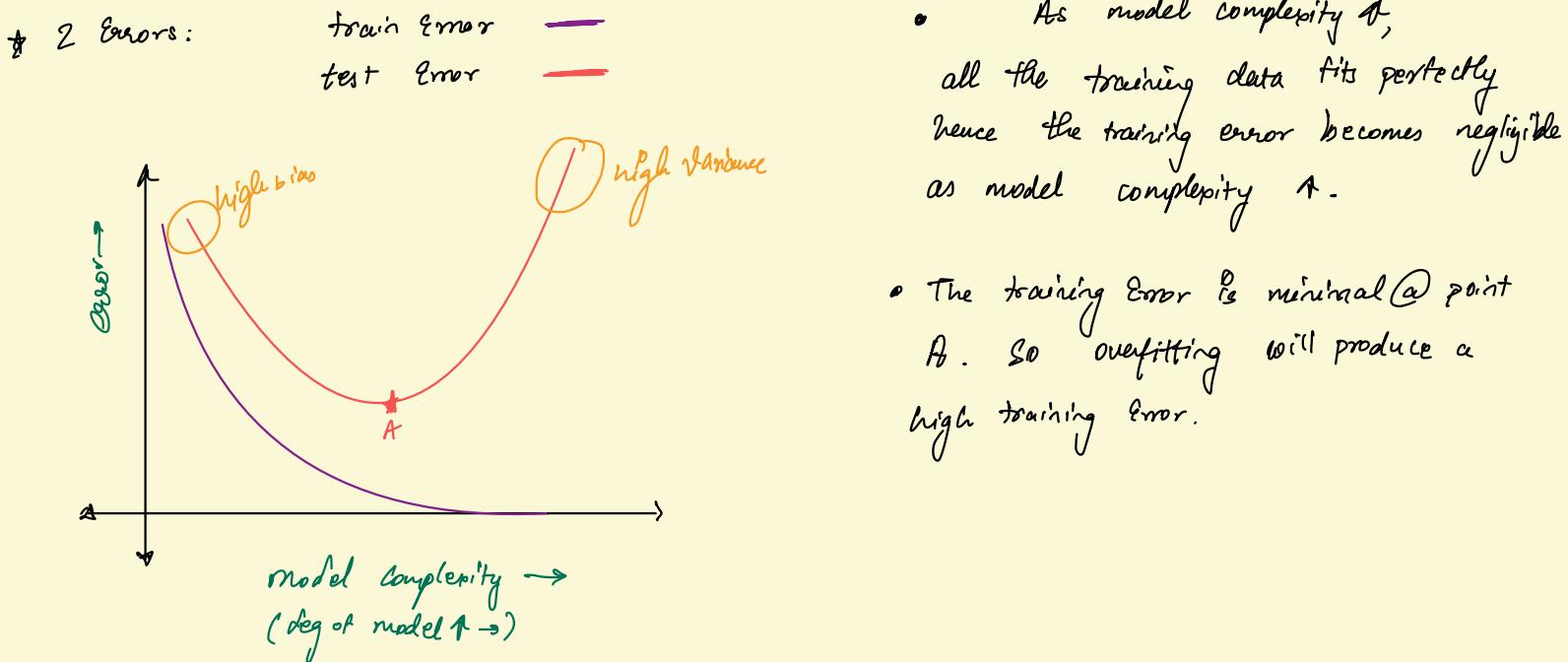
(Refer Pg 49 TB Hui Jiang)

$$\text{Error} = \text{Bias}^2 + \text{Variance} + \sigma^2 \text{ (lmm. error)}$$

• If Bias is high, Error \uparrow and also if Variance is high, error \uparrow .

• That's why we don't want a very high bias or a very high variance.

Hence we want something in b/w. high bias & high variance.



- As model complexity ↑, all the training data fits perfectly hence the training error becomes negligible as model complexity ↑.
- The training error is minimal @ point A. So overfitting will produce a high training error.

b) Let there be n training points & m test points.

$$\text{train err} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2$$

$$\text{test err} = \frac{1}{m} \sum_{i=n+1}^{n+m} (y_i - \hat{f}(x_i))^2$$

As model complexity ↑ train err becomes overly optimistic and gives a wrong picture of how close \hat{f} is to f .

$$y_i = f(x_i) + \varepsilon_i$$

means, y_i is related to x_i by some true relation but there is also some noise ε in the relation.

Assume $\varepsilon \sim N(0, \sigma^2)$ for simplicity.

We use \hat{f} to approximate f and the parameters using T & D s.t

$$y_i = \hat{f}(x_i)$$

As the model complexity ↑, train err becomes overly optimistic & gives a wrong pic of how close \hat{f} is to f

We are interested in knowing $E[(\hat{f}(x_i) - f(x_i))^2]$

$$E[(\hat{y}_i - y_i)^2] = E[(\hat{f}(x_i) - f(x_i) - \varepsilon_i)^2] = E[(\hat{f}(x_i) - f(x_i))^2] - 2E[\varepsilon_i(\hat{f}(x_i) - f(x_i))] + E[\varepsilon_i^2]$$

$$\therefore E[(\hat{f}(x_i) - f(x_i))^2] = E[(\hat{y}_i - y_i)^2] - E[\varepsilon_i^2] + 2E[\varepsilon_i(\hat{f}(x_i) - f(x_i))]$$

$$E[(\hat{y}_i - y_i)^2] = \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i)^2$$

We can empirically evaluate R^2 by using training or test observations.

Case I

Using test observations

$$E((\hat{f}(x_i) - f(x_i))^2) = \frac{1}{n} \sum_{i=n+1}^{n+m} (y_p - \hat{y}_p)^2 - \frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i^2 + 2 E(\hat{\epsilon}_i (\hat{f}(x_i) - f(x_i)))$$

$\underbrace{\quad}_{\text{true error}}$ $\underbrace{\sum_{i=n+1}^{n+m} (y_p - \hat{y}_p)^2}_{\text{Empirical estimation of error}}$ $\underbrace{\frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i^2}_{\text{small constant}}$ $\underbrace{2 E(\hat{\epsilon}_i (\hat{f}(x_i) - f(x_i)))}_{\text{Covariance } \hat{\epsilon}_i (\hat{f}(x_i) - f(x_i))}$

When we have covariance b/w 2 things the two things must be random variables.

$$\begin{aligned} \text{Cov}(X, Y) &= E((\hat{x} - \mu_x)(\hat{y} - \mu_y)) \\ &= E(XY) - E(X)\mu_y - \mu_x E(Y) \\ &= E(XY) - E(X)E(Y) = E(XY) - \mu_y \mu_x \end{aligned}$$

Parameters of \hat{f} were known from the training data $\hat{\epsilon}_i$'s are coming from test data. That's why $\hat{\epsilon}_i$ & $(\hat{f}(x_i) - f(x_i))$ are independent.

$$\therefore E(\hat{\epsilon}_i (\hat{f}(x_i) - f(x_i))) = E(\hat{\epsilon}_i) E(\hat{f}(x_i) - f(x_i)) = 0 \cdot E(\hat{f}(x_i) - f(x_i)) = 0$$

∴ True Error = empirical test + small const. error

This tells if a model is learned & then I take the estimate of the error, on some data which was not used for the training, then that error is actually very close to the true error. It only differs by a constant.

∴ We are trying to estimate the true error (which cannot be explicitly estimated) by using the data. If data used is test data then empirical estimation of test error is close to the true error.
Inference to estimate the error, we should always use a validation set (independent of the training set).

Case II

Using training observations

$$E((\hat{f}(x_i) - f(x_i))^2) = \frac{1}{n} \sum_{i=1}^n (y_p - \hat{y}_p)^2 - \frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i^2 + 2 E(\hat{\epsilon}_i (\hat{f}(x_i) - f(x_i)))$$

$\underbrace{\sum_{i=1}^n (y_p - \hat{y}_p)^2}_{\text{Empirical estimation of train error}} - \underbrace{\frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i^2}_{\text{small constant}} + 2 E(\hat{\epsilon}_i (\hat{f}(x_i) - f(x_i)))$
 \downarrow
 $+ 0$

$\varepsilon \propto \hat{f}(x)$ because it's used for estimating parameters of $\hat{f}(x)$

$$E(\varepsilon_i | f(x_i) - \hat{f}(x_i)) \neq E(\varepsilon_i) E(\hat{f}(x_i) - f(x_i)) \neq 0.$$

Empirical train error is small compared to the true error, & doesn't give a true picture of the error. In other words, it's overly optimistic.

Stein's lemma

$$\frac{1}{n} \sum_{i=1}^n \varepsilon_i (\hat{f}(x_i) - f(x_i)) = \frac{\sigma^2}{n} \sum_{i=1}^n \frac{\partial \hat{f}(x_i)}{\partial y_i}$$

The derivative is high if y changes a bit, the change in $\hat{f}(x_i)$ is going to be large.

↓
This happens for complex models (because Bias is high recall!)

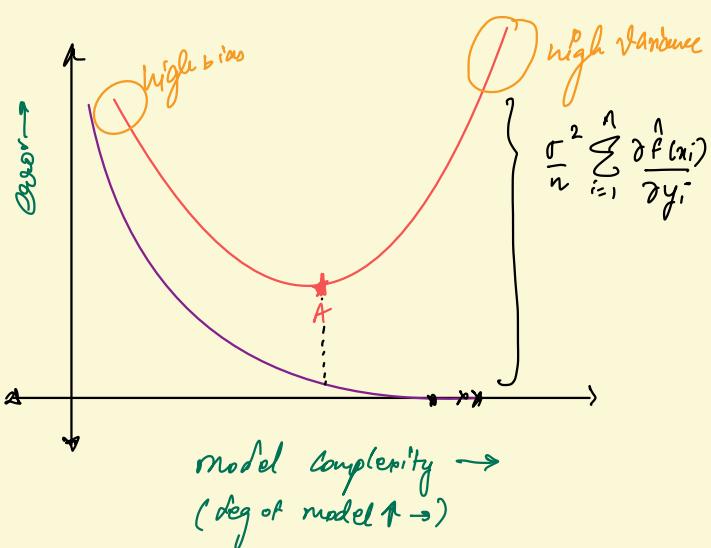
Hence,
true error = Empirical train + small const + Ω (model complexity)

$$\min_{\theta} J_{\text{train}}(\theta) + \Omega(\theta) = L(\theta)$$

where $J(\theta)$ would be high for complex models & small for simple models.
 $\Omega(\theta)$ acts as an approximate of $\frac{\sigma^2}{n} \sum_{i=1}^n \frac{\partial \hat{f}(x_i)}{\partial y_i}$.

2 Errors:

train Error test / validation	
--	---



we are aiming for the point A.
we should ensure using $\Omega(\theta)$ that the gap is also minimised.

$$\min_{\theta} J(\theta) + \Omega(\theta)$$

Hence model would generalize better on the test data.

* L^2 Regularization:

$$L(\omega) = L(\omega) + \frac{\lambda}{2} \|\omega\|^2$$

↑ $\Omega(\theta)$, relating to model complexity
by minimizing $\|\omega\|$, we are reducing the complexity of the model by ensuring that the weights don't take any possible value.

* Regularized Linear Regression:

$$L(\omega) = \underset{\omega}{\operatorname{arg\,min}} \left(\sum_{i=1}^m (\omega^T x_i - y_i)^2 + \lambda \|\omega\|_2^2 \right)$$

Closed form solution.

$$\sum_{i=1}^m (\omega^T x_i - y_i)^2 = \|x\omega - y\|^2, \quad \|\omega\|^2 = \omega^T \omega$$

$$L(\omega) = \underset{\omega}{\operatorname{arg\,min}} \|x\omega - y\|^2 + \lambda \|\omega\|^2$$

$$\frac{\partial L}{\partial \omega} = 0 \Rightarrow \frac{\partial}{\partial \omega} ((x\omega - y)^T (x\omega - y)) + \lambda \frac{\partial}{\partial \omega} (\omega^T \omega) = 0$$

\Downarrow
 $(x\omega^T x^T - y^T)(x\omega - y) + 2\lambda \omega = 0$

$$\Rightarrow \frac{\partial}{\partial \omega} (-y^T x \omega + y^T y + \omega^T x^T x \omega - \omega^T x^T y) + 2\lambda \omega = 0$$

$$\Rightarrow -y^T x + \frac{\partial}{\partial \omega} (x^T \omega x \omega) - x^T y + 2\lambda \omega = 0 \quad (A-B)^T = A^T - B^T$$

$$\Rightarrow -y^T x + \frac{\partial}{\partial \omega} (x x^T \omega^T \omega) - x^T y + 2\lambda \omega = 0$$

$$\Rightarrow -y^T x + 2x^T x \omega + 2\lambda \omega - x^T y = 0$$

$$\Rightarrow 2(x^T x + \lambda I)\omega = 2x^T y$$

$$\Rightarrow \omega = (x^T x + \lambda I)^{-1} (x^T y)$$

which is the closed form solution.

$$\frac{\partial}{\partial \omega} (\omega^T \omega) = 2\omega$$

$$\frac{\partial}{\partial x} (x^T y) = y$$

$$x = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_m^T \end{bmatrix}_{m \times n}$$

m samples
n features.

* Regularized linear regression

$$\min_{w,b} J(w,b) = \min_{w,b} \left(\frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2 \right)$$

m samples

n features

$$\frac{\partial}{\partial w_j} J(w,b) = \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} w_j$$

Gradient Descent

$$\frac{\partial}{\partial b} J(w,b) = \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})$$

repeat

$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(w,b)$$

$$b = b - \alpha \frac{\partial}{\partial b} J(w,b)$$

g simultaneous update

for implementation:

$$\begin{bmatrix} dw \\ dw_1 \\ \vdots \\ dw_n \end{bmatrix}^T$$

$$= \frac{1}{m} \begin{bmatrix} (f(x_1) - y_1), \dots, (f(x_m) - y_m) \end{bmatrix}^T \begin{bmatrix} x_1^1 & x_2^1 & \dots & x_n^1 \\ x_1^2 & x_2^2 & \dots & x_n^2 \\ \vdots & \vdots & \ddots & \vdots \\ x_1^m & x_2^m & \dots & x_n^m \end{bmatrix}$$



$$\begin{bmatrix} dw_1 \\ \vdots \\ dw_n \end{bmatrix} = f(x) \begin{bmatrix} s_1 & s_2 & \dots & s_m \\ x_1^1 & x_1^2 & \dots & x_1^m \\ x_2^1 & x_2^2 & \dots & x_2^m \\ \vdots & \vdots & \ddots & \vdots \\ x_n^1 & x_n^2 & \dots & x_n^m \end{bmatrix} \begin{bmatrix} f(x_1) - y_1 \\ f(x_2) - y_2 \\ \vdots \\ f(x_m) - y_m \end{bmatrix} + \frac{\lambda}{m} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}$$

$$+ \frac{\lambda}{m} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}^T$$

$$x^T (f - y) + \frac{\lambda}{m} w$$

$$w_j = w_j - \alpha \frac{\lambda}{m} w_j - \frac{\alpha}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$(1 - \frac{f(x)}{m})$ is a quantity smaller than 1 effectively reducing the size of $f(x)$,

* Regularized Logistic Regression

$$J(w,b) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(f_{w,b}(x^{(i)})) + (1-y^{(i)}) \log(1-f_{w,b}(x^{(i)})) + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

Gradient descent is same as the linear regression

Only $f_{w,b}(x)$ differs.

$$\frac{\partial}{\partial w_j} J(w,b) = \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} w_j$$

$$\frac{\partial}{\partial b} J(w,b) = \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})$$

$$\frac{1}{1 + e^{-x}}$$

