

## A theory of linear models:

Consider, a set of observations  
of a system,  $(x_1, y_1), (x_2, y_2)$   
- - -  $(x_n, y_n)$

The underlying function, generating  
these observations, is non-linear,  
complicated & unknown:

$$y_n = g(x_n)$$

## Page 2

A linear model seeks a linear approximation to the true ' $g$ ' model (which can be generally non linear)

So, we assume that there exists, a set of coefficients  $\{\alpha_1, \alpha_2\}$  that describes the linear model.

$$y_1 = \alpha_1 x_1 + \alpha_2$$

$$y_2 = \alpha_1 x_2 + \alpha_2$$

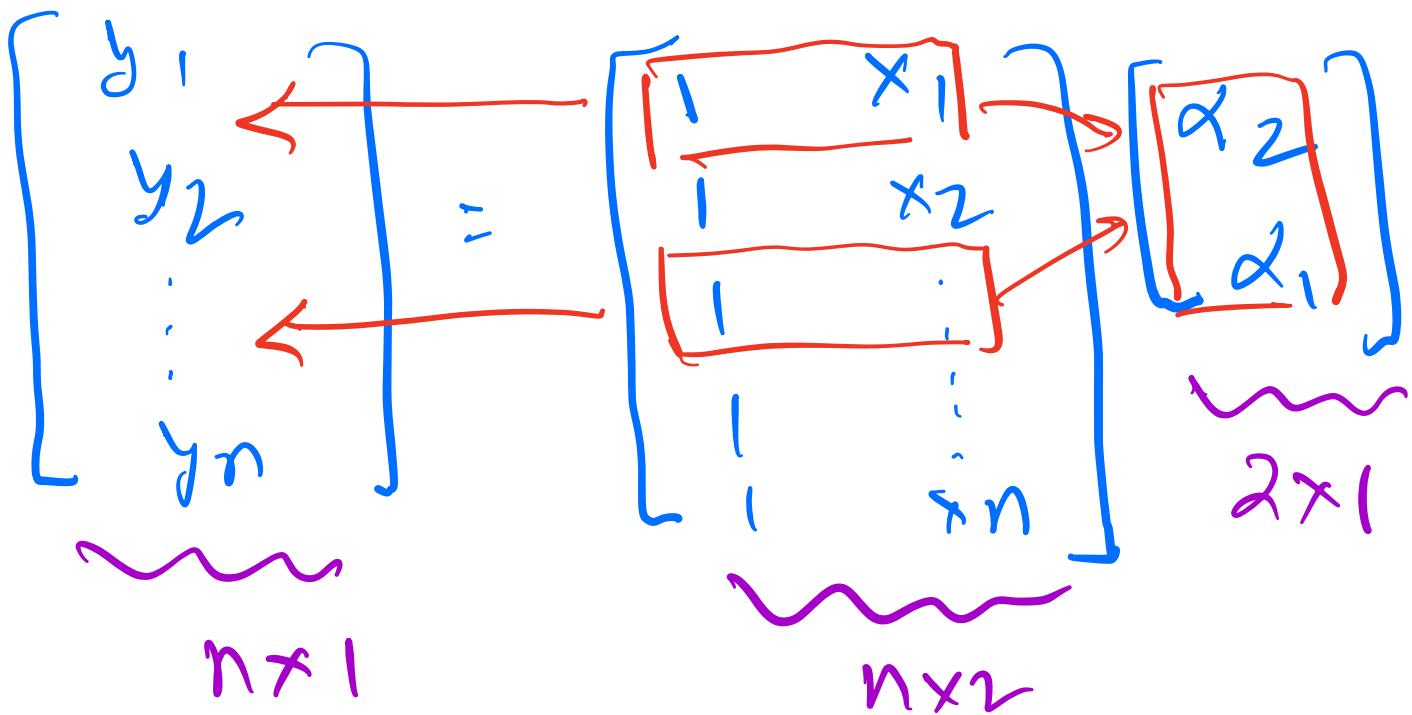
⋮

$$y_n = \alpha_1 x_n + \alpha_2$$

We can write this in matrix fashion.

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \alpha_1 \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} + \alpha_2 \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

# Page 4



Note: Each row is a constraint.

$(d_1, d_2) \rightarrow$  is a tuple of parameters

We want to find **optimal**

$$(d_1, d_2) = \theta$$

In machine learning, the matrix  $\begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}$  is called

The **feature matrix** where each column is a **feature**.

This is also the **design matrix**. We will call it  $A$ .

$\therefore y = A\theta \rightarrow$  linear problem

# Page 6

Now, we could had written  
this as:

$$\text{For } y_n, \quad y_n = \alpha_1 x_n + \alpha_2 x_n^2 + \alpha_3$$

Then:

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_1 & x_1^2 \\ \vdots & x_2 & x_2^2 \\ \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 \end{bmatrix} \begin{bmatrix} \alpha_3 \\ \alpha_2 \\ \alpha_1 \\ \hline 3 \times 1 \end{bmatrix}$$

$n \times 1$        $n \times 3$

# Page 7

We can keep adding more columns like these,  $x_n^2, x_n^3, \dots, e^{x_n}, \log(x_n), \dots$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_1 & x_1^2 & x_1^3 & \dots & x_L^p \end{bmatrix} \begin{bmatrix} d_p \\ d_{p-1} \\ \vdots \\ d_1 \end{bmatrix}$$

$n \times 1$        $n \times p$

new features.

# Page 8

This is still a linear model.

Because it is linear in  $\theta$

Note! Does this remind you  
of the basis functions in  
 $\sin \theta$ ?

So, every linear model is

; we can have

$$y = A \theta$$

$\underbrace{\quad}_{n \times 1} \quad \underbrace{\quad}_{n \times p} \quad \underbrace{\quad}_{p \times 1}$

$\left\{ \begin{array}{l} n > p \\ \text{or } n < p \end{array} \right.$

Over-determined

When  $n \gg p$ , we have a tall skinny matrix. In such a situation, "generally", we have no unique solution. This is because we have more equations than unknowns. So, instead, we search for a least square solution:

$$L(\theta) = \|A\theta - y\|_2^2$$

$$\theta^* = \underset{\theta}{\operatorname{arg\,min}} L(\theta).$$

## Page 10

In the under-determined case, where  $n \ll p$ , we have a short & fat matrix. There are infinitely many solutions to  $A\theta = y$ . We find the solution, where  $\|\theta\|_2$  is minimized.

Role of SVD in computing  $\theta^*$

$A = U\Sigma V^*$ . We can invert each of these matrices to compute the moore-penrose inverse of rectangular  $A$ .

Page 11

$$A^{-1} \equiv A^+ = V \Sigma^{-1} U^*$$

$$A^+ A = I_{(m \times m)}$$

$$A^+ A \theta = A^+ y \Rightarrow$$

$$\theta = \boxed{V \Sigma^{-1} U^* y}$$

Another way of solving this:

$$y = A\theta \text{ for } n \geq p$$

The next few pages assumes some knowledge on Random Matrix Theory

$$y A^T = (A^T A) \theta$$

$$\theta^* = \underbrace{(A^T A)^{-1}}_{PAP^{-1}} (A^T y)$$

normal equation

Suppose  $y_n = g(x_n)$  is an actual linear process with noise:

$$[y_n] = A \hat{\theta} + \eta \quad \eta \sim N(0, \sigma^2 I)$$

Then:

$$\|\hat{\theta} - \theta^*\|^2$$

$$= \|\hat{\theta} - (A^T A)^{-1} A^T y\|^2$$

$$= \|\hat{\theta} - (A^T A)^{-1} A^T (A \hat{\theta} + \eta)\|^2$$

$$= \|\cancel{\hat{\theta}} - \underbrace{(A^T A)^{-1} (A^T A)}_{I} \hat{\theta} - (A^T A)^{-1} A^T \eta\|^2$$

$$= \|(A^T A)^{-1} A^T \eta\|^2$$

$$= (A^T A)^{-1} A^T \eta \eta^T A (A^T A)^{-1}$$

We can re-write this as:

$$\|\hat{\theta} - \theta^*\|^2 = (\eta^T A (A^T A)^{-2} A^T \eta)$$

For multiple samples of  $\eta \sim N(0, \sigma^2 I)$

$$E_N [\|\hat{\theta} - \theta^*\|^2]$$

$$= E_N [\eta^T A (A^T A)^{-2} A^T \eta]$$

Now, we will use two results from classical statistics:

$$E_x[x^T B x] = \text{Tr}(B)$$

↳ Trace

$$\tilde{\eta} \sim N(0, I)$$

$$S_0: E_N[\eta^T A (A^T A)^{-2} A^T \eta]$$

$$= \sigma^2 \text{Tr}(A (A^T A)^{-2} A^T)$$

Page 16

$$= \sigma^2 \text{Tr} \left( (\mathbf{A}^\top \mathbf{A})^{-1} \right)$$

$$= \sigma^2 \frac{P}{n} \text{tr} \left( \frac{\mathbf{A}^\top \mathbf{A}}{n} \right)^{-1}$$

$\text{tr}$ : normalized trace

$$= \frac{\text{Tr}}{n}$$

So we have:

$$E_N \left[ \|\hat{\theta} - \theta^*\|^2 \right]$$

$$= \frac{\sigma^2}{n} \text{tr} \left( \left( \frac{1}{n} A^T A \right)^{-1} \right)$$

At this point, we need  
a primer on random  
matrix theory

Wishart matrices: A

matrix  $\hat{\Sigma}$  is a wishart  
matrix  $\hat{\Sigma} = \frac{1}{n} X^T X$

where each row of the  
matrix  $X$  is drawn  
from a Gaussian

$$N(0, \Sigma_{P \times P})$$

For Wishart matrices,  
the histogram of the  
eigenvalues of  $\Sigma = \frac{1}{n} X^T X$   
converges to the  
Marchenko-Pastur  
distribution.

If  $P_n \rightarrow \gamma \in [0, 1]$

as  $n \rightarrow \infty$

$\varphi_{MP}(t) \rightarrow \text{Pdf of}$   
markovian process

$$= \frac{1}{2\pi\gamma t} \sqrt{(x_f - t)(t - x_i)}$$

on  $[x_i, x_f]$

$$\gamma_- = (1 - \sqrt{\gamma})^2$$

$$\gamma_+ = (1 + \sqrt{\gamma})^2$$

Stieljes Transform:

---

For Wishart matrices,  $Z$

$$S_n(Z) = \sum \left[ \text{tr} \left[ \left( \hat{\Sigma} - Z^\top Z \right)^{-1} \right] \right]$$

For the marchenko-pastur distribution, the Stieltjes transform is given as:

$$S_{MP}(z) = \frac{1}{(t-z)} \varphi_{MP} dt$$
$$\forall z \in \mathbb{C} \setminus \mathbb{R}$$

We come back to our estimated  $\theta^*$  in linear regression

$$E_N \left[ \|\hat{\theta} - \theta^*\|^2 \right]$$

$$= \sigma^2 \frac{1}{n} \text{tr} \left( \left( \frac{1}{n} A^T A \right)^{-1} \right)$$

If  $A$  is a standard Gaussian random matrix, then:

$\hat{\Sigma} = \frac{1}{n} A^T A$  is a Wishart matrix, which as  $n \rightarrow \infty$  converges to  $\Sigma_{\text{GFS}}(0)$

We know:

$$S(z) = \int \frac{1}{t-z} \Psi_M dt$$

and

$$1 + z S(z) = \frac{S(z)}{1 + \gamma S(z)}$$

of  $\hat{\theta} = \theta$ ,

$$S(\theta) = \frac{1}{1-\gamma}$$

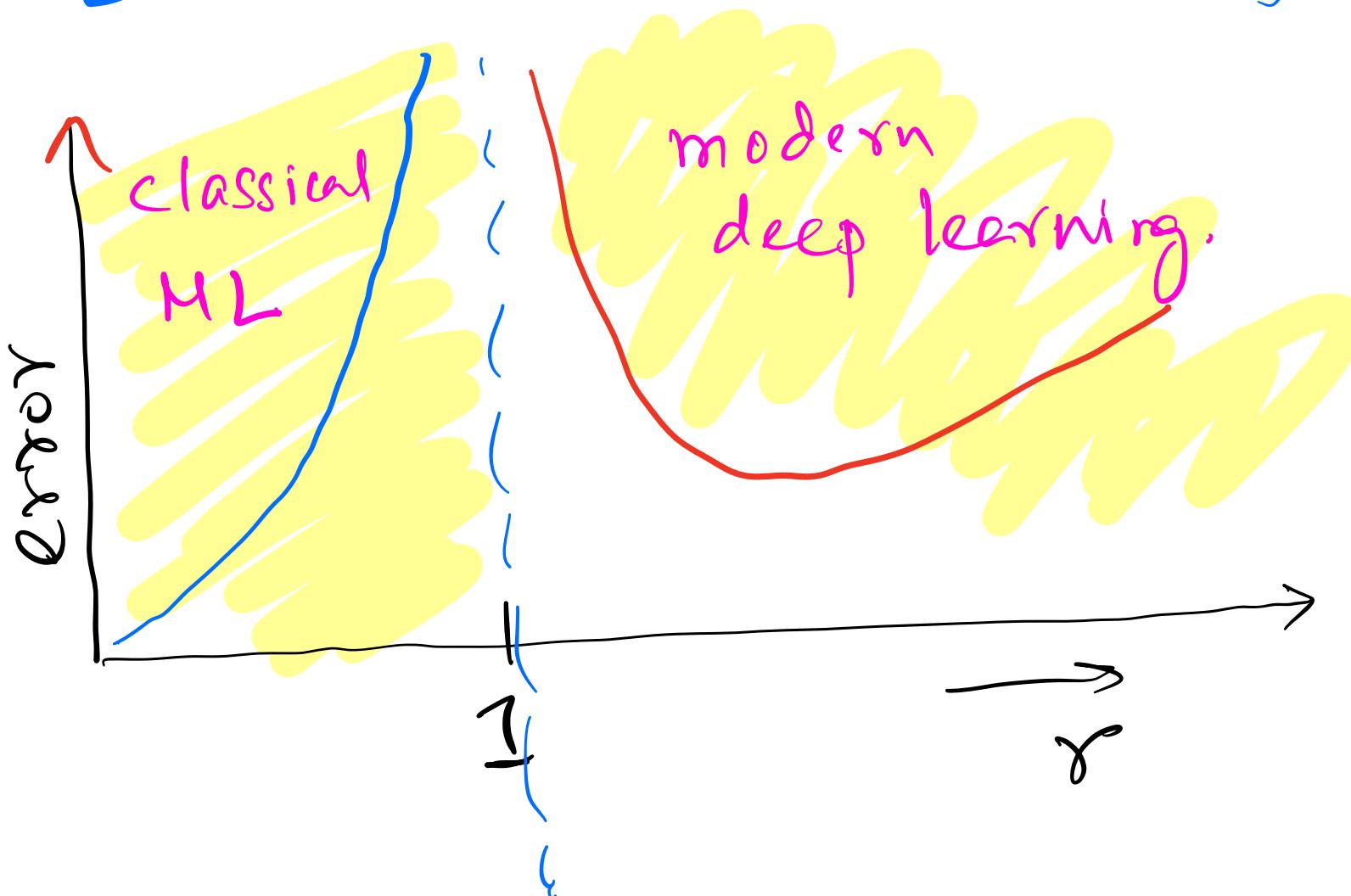
$\nearrow n \gg p$   
range

$\therefore E[\|\hat{\theta} - \theta^*\|^2] = \sigma^2 \frac{\gamma}{1-\gamma}$

→ A lot of analysis to prove that the risk in estimated parameters is a function of parameters vs samples ( $\gamma$ )

For  $p \gg n$ , we have.

$$E[\|\hat{\theta} - \theta^*\|^2] = \|\hat{\theta}\|^2 \left(1 - \frac{1}{\gamma}\right)$$



# Bias - Variance Trade-Off

---

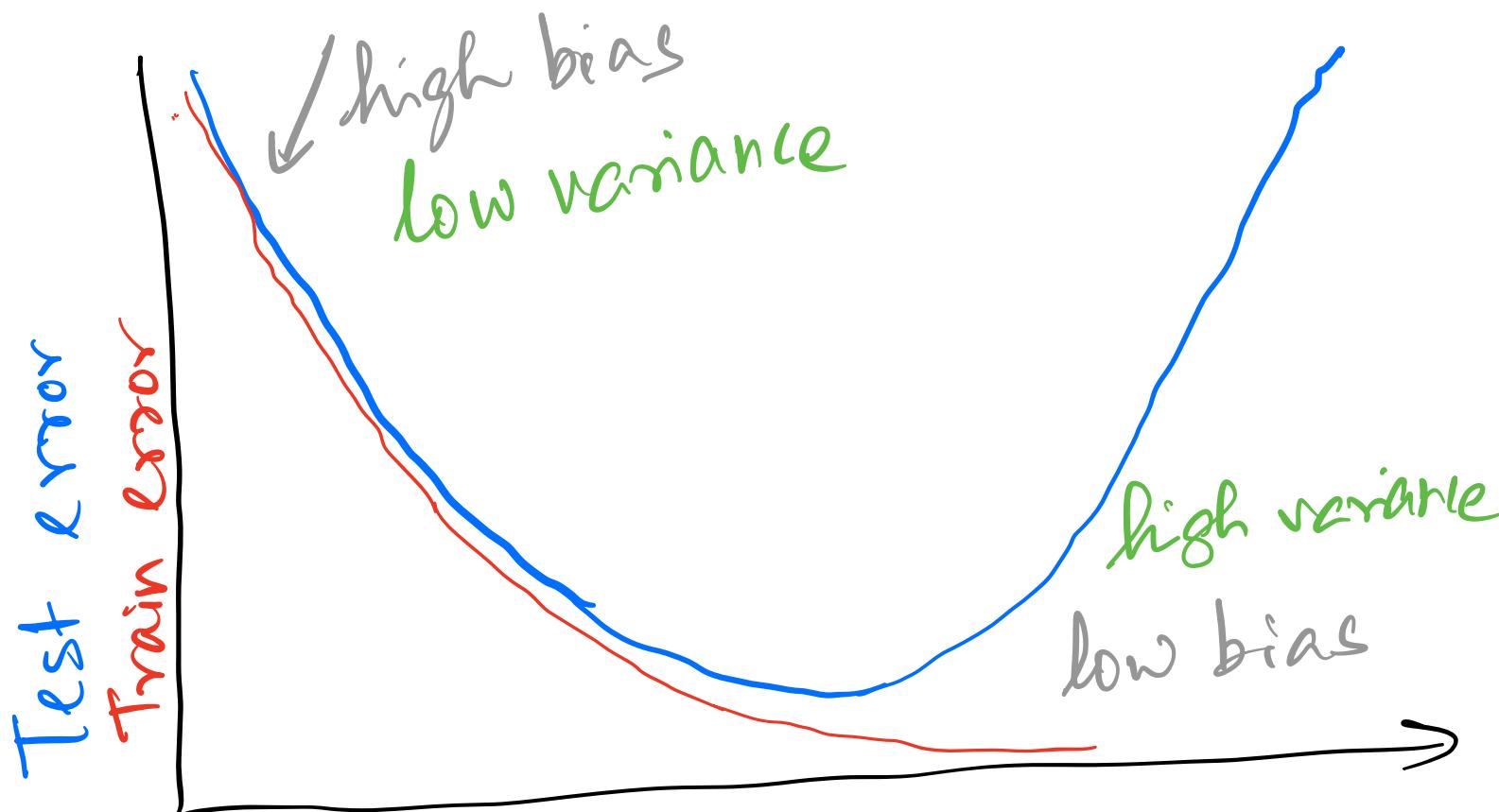
How do we design a model?

How many parameters should we have?

## Classical Wisdom:

- Start with a small model
- Keep increasing the model as test error decreases.
- Stop, when the test error increases again.

See the slide deck to understand visually how bias decreases and variance increases



We will prove that :

$$\text{MSE error} = (\text{bias})^2 + \text{variance} + \text{irreducible error}$$