

Lecture 14

X) Regression:

Given $x_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}^p$,

for a test point x_{test} , we want to predict

$$y_{\text{test}} = \hat{f}(x_{\text{test}})$$

$$\hat{f}(x) = w_1 x + w_0$$

$$w_2 x^2 + w_1 x + w_0$$

$$d=1 \quad p=1$$

$$X = \begin{bmatrix} x_1 & x_2 & \dots & x_n \\ 1 & 1 & \dots & 1 \end{bmatrix}$$

$$y \in \mathbb{R}^{n \times 1}$$

$$W = (X'X)^{-1} X' y$$

What if you take several sub-datasets where samples
Come from same distribution and then compute f^{avg} ?

$$\hat{f}_1(x) \leftarrow w_2' x^2 + w_1' x + w_0' \text{ learnt } \mathcal{D}_1$$

$$\hat{f}_2(x) \leftarrow w_2^2 x^2 + w_1^2 x + w_0^2 \text{ learnt } \mathcal{D}_2$$

\vdots

$$\hat{f}_n(x) \leftarrow w_2^n x^2 + w_1^n x + w_0^n \leftarrow \mathcal{D}_n$$

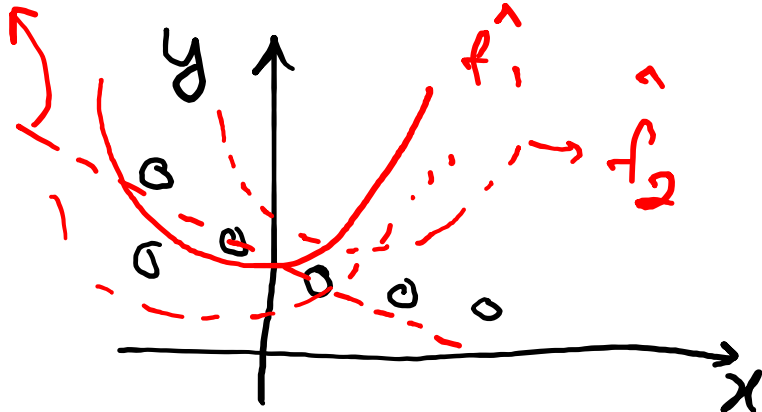
$$\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_n \leftarrow \mathcal{D}$$

$$\hat{f}_{\text{avg}}(x) = x^2 E(w_2) + x E(w_1) + E(w_0)$$

Bias-Variance Tradeoff

$$E(\hat{f}) \approx f$$

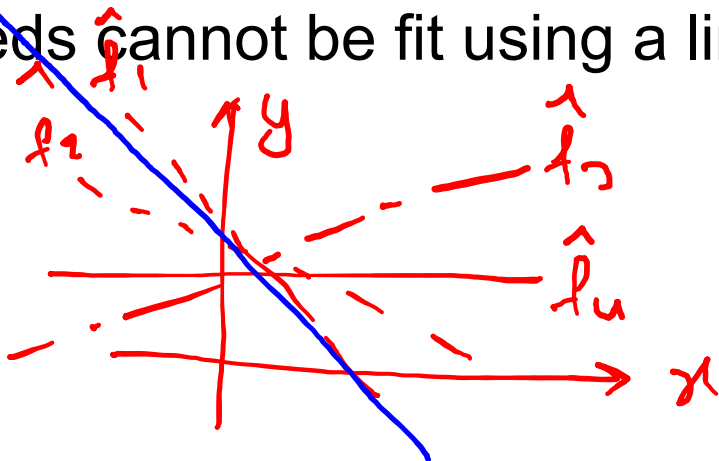
- For higher M : \hat{f} avg of $E(\hat{f})$ can fit into data that requires degree M or less.
- Variance: $(1/n-1)\{\hat{f}_i(x) - E(\hat{f})\}^2$ is **large**
- For eg. $\hat{f}_i(x)$ may be quadratic, whereas, expectation may be quad. Linear or constant.
- Bias: Will be **small** as $E(\hat{f})$ approx. equal to true f .



True $f \rightarrow$ was a straight line

$E(\hat{f})$ can very well model straight line.

- For lower M : \hat{f}^{avg} of $E(\hat{f})$ will under fit the data that requires a higher degree.
- A data that needs cannot be fit using a line.
- Variance: $(1/n-1)\{\hat{f}_i(x) - E(\hat{f})\}^2$ is **small**
- For eg. $\hat{f}_i(x)$ may be linear, whereas, expectation may be, Linear or constant.
- Bias: Will be **large** as $E(\hat{f})$ wont be approx. equal to true f . A data that needs cannot be fit using a line.



$E(\hat{f})$ will be farther away from true f

Cross-validation

- Consider degree $m = 1, 2, \dots, M$
 - Divide data into K folds. $K = 5, 10$.
 - Hold one fold and use “remaining folds”.
 - Learn W from “remaining folds”.
 - Apply W to compute error on “remaining folds” – call this training error (k)
 - Apply W to compute error on “held out folds” – call this validation error (k)
 - end
 - avgTrainErr(m) = mean{train error}
 - avgValErr(m) = mean{val error}
- end

$$D = \{(x_1, y_1), (x_2, y_2) \dots (x_{10}, y_{10})\}$$

$D \rightarrow$ into 5 folds.

$$1^{st} \text{ fold} \rightarrow (x_1, y_1) \quad (x_2, y_2)$$

$$2-5^{th} \text{ fold} \rightarrow (x_3, y_3) \dots (x_{10}, y_{10})$$

$$W \leftarrow (X'X)^{-1} X'Y$$

$$X = \begin{bmatrix} x_1 & \dots & x_{10} \\ 1 & & 1 \end{bmatrix} \quad Y = \begin{bmatrix} y_1 \\ \vdots \\ y_{10} \end{bmatrix}$$

apply W on 1st fold.

$$(x_1, y_1) (x_2, y_2)$$

val. error (1)

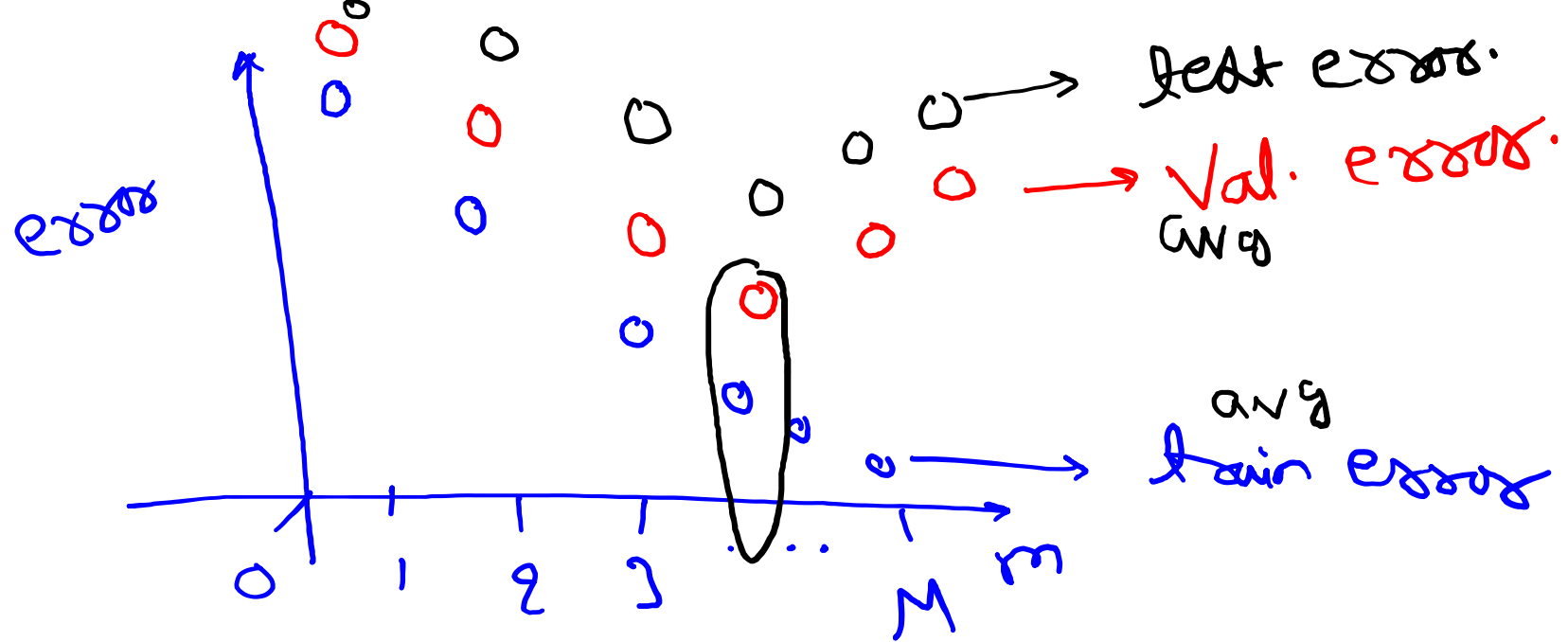
apply W on 2nd - 5th fold.

$$(x_3, y_3), \dots, (x_{10}, y_{10})$$

train error (1)

for $k=2$, $(x_1, y_1) (x_2, y_2), (x_5, y_5) \dots (x_{10}, y_{10})$

\nwarrow
 W apply them on $(x_3, y_3) (x_4, y_4) \rightarrow \text{val error}(2)$



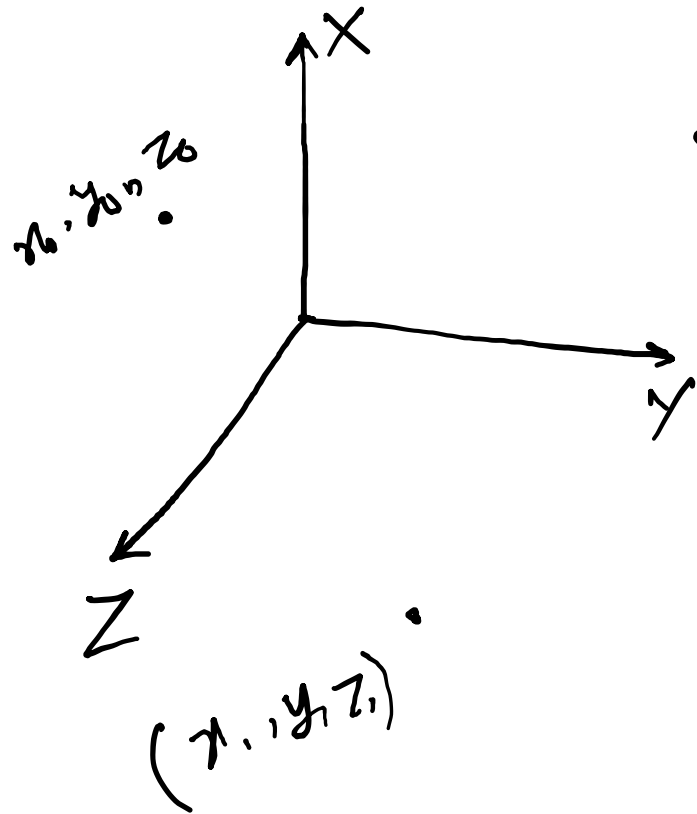
'm' that gives the least val. error. will be
'm' that you will use further.

for eg. $m=4$, gives least val. error.

$m=4$, & use the full training set.

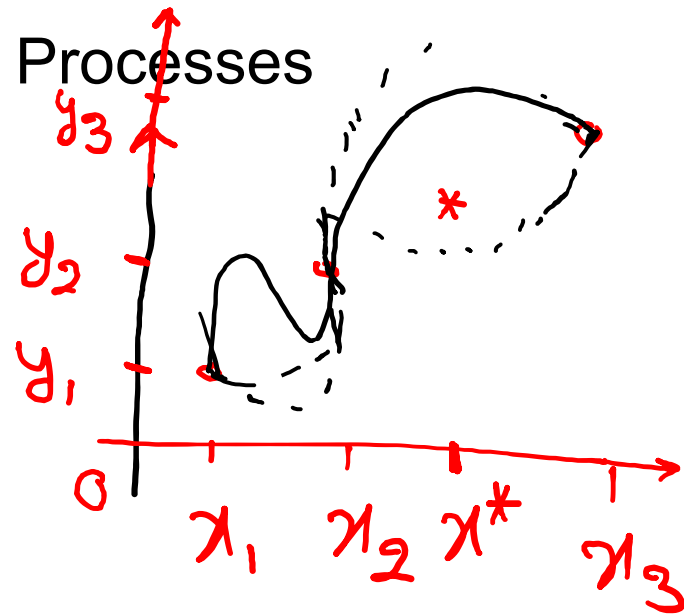
→ select $W \rightarrow$ testing set.

Gaussian Processes



\mathcal{S}_0

\mathcal{S}_1



*) (x_1, x_2) are closer \rightarrow

(y_1, y_2) similar

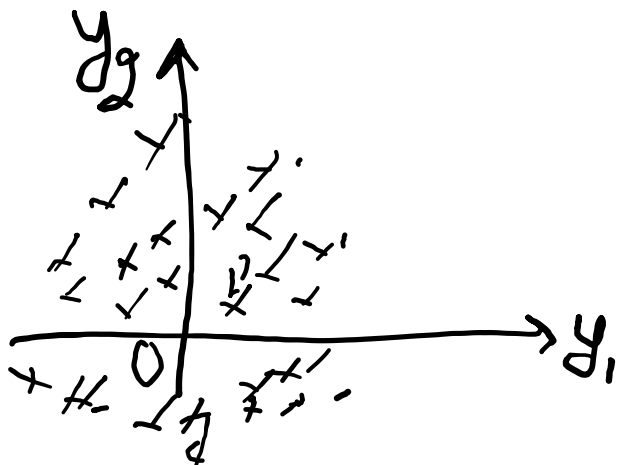
*) If (x_1, x_2) are farther

(y_1, y_2) dis-similar

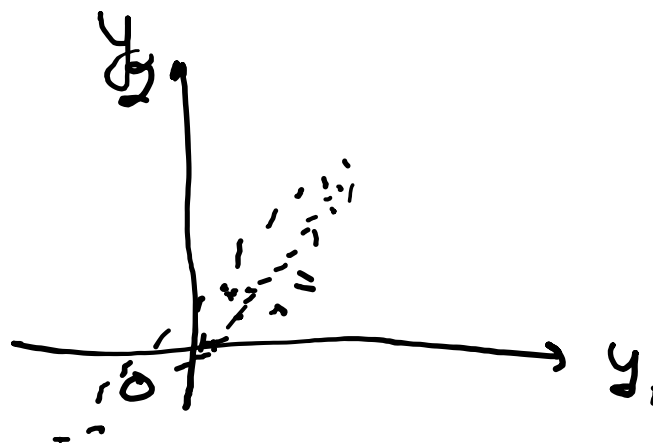
$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$ assume coming M.V.G \sim

$$N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}\right)$$

$$\text{Cov.} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



$$\text{Cov.} \begin{bmatrix} 1 & 0.99 \\ 0.99 & 1 \end{bmatrix}$$



$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} \sim N \left(\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0.8 & 0.2 \\ 0.8 & 1 & 0.5 \\ 0.2 & 0.5 & 1 \end{bmatrix} \right)$$

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} \sim N \left(\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} K(x_1, x_1) & K(x_1, x_2) & K(x_1, x_3) \\ K_{21} & K_{22} & K_{23} \\ K_{31} & K_{32} & K_{33} \end{bmatrix} \right)$$

Sim. function. $K(a, b)$

$$K(x_i, x_j) = K_{ij}$$

$$K(x_i, x_j) = \sigma^2 e^{-\|x_i - x_j\|_2^2 / 2l^2}$$

$l \rightarrow$ width of Gaussian kernel

$\sigma \rightarrow$ max of kernel.

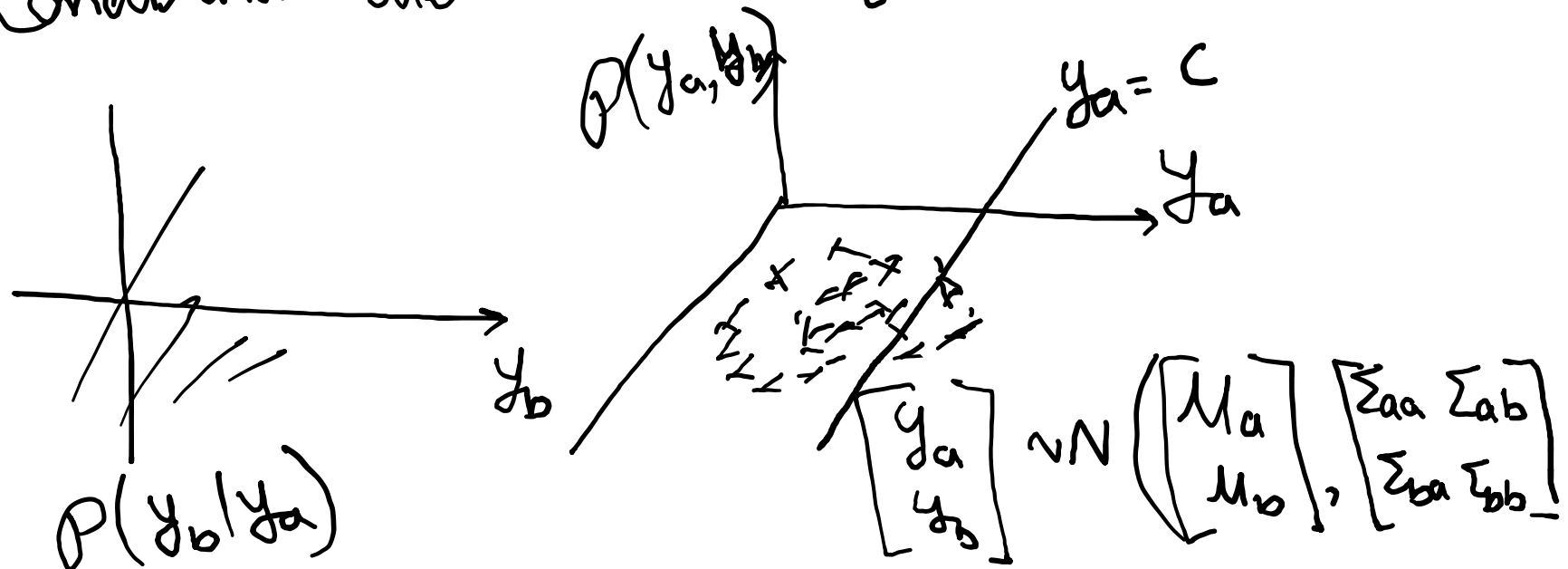
$\sigma, l \rightarrow$ Unknown.

$$\underbrace{y_a}_{\left[\begin{array}{c} y_1 \\ y_2 \\ y_3 \\ y_4 \end{array} \right]} \sim N \left(\underbrace{\left[\begin{array}{c} 0 \\ 0 \\ 0 \\ 0 \end{array} \right]}_{\Sigma_{aa}}, \underbrace{\left[\begin{array}{cccc} K_{11} & K_{12} & K_{13} & K_{14} \\ K_{21} & K_{22} & K_{23} & K_{24} \\ K_{31} & K_{32} & K_{33} & K_{34} \\ K_{41} & K_{42} & K_{43} & K_{44} \end{array} \right]}_{\substack{\Sigma_{aa} \quad \Sigma_{ab} \\ \downarrow \Sigma_{bb}}} \right)$$

$$\left. \begin{aligned} \mu_{b|a} &= \mu_b + \Sigma_{ba} \Sigma_{aa}^{-1} (y_a - \mu_a) \\ \Sigma_{b|a} &= \Sigma_{bb} - \Sigma_{ba} \Sigma_{aa}^{-1} \Sigma_{ab} \end{aligned} \right\} \text{theorem.}$$

$$P(y^* | x, y) \sim N(\mu^*, \sigma^*)$$

Conditional distribution on y^* given all other info



$$p(y^* | y_1, y_2, y_3) \sim$$

$$\mu^* = \sum_b a \sum_a^{-1} y_a$$

$$= \begin{bmatrix} K_{1x} & K_{2x} & K_{3x} \end{bmatrix} \begin{bmatrix} K_{11} & K_{12} & K_{13} \\ & \ddots & \\ K_{31} & & K_{33} \end{bmatrix}^{-1} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

$$\Sigma_{b|a} = K_{xx} - \begin{bmatrix} K_{1x} & K_{2x} & K_{3x} \end{bmatrix} \begin{bmatrix} & & \\ & \ddots & \\ & & \end{bmatrix}^{-1} \begin{bmatrix} K_{1x} \\ K_{2x} \\ K_{3x} \end{bmatrix}$$

C.I.

$$\mu^* \pm 1.96 \sqrt{\Sigma_{b|a}}$$

2-deg.