

Announcements

- Homework 2 was posted.
-

Today's Lecture

- Ridge regression (finish)
- Notation comment
- Bayesian inference

NOTE: A COUPLE LINES HAVE BEEN ADDED TO pg-8
FOR CLARIFICATION.

RIDGE REGRESSION, PART 2

Subs. (2) & (3) \rightarrow (1): 1st term

$$\hat{\underline{w}} = \hat{\underline{\theta}} = \underset{\underline{w}}{\operatorname{argmax}} \left\{ \underbrace{\sum_{i=1}^N \ln N(y_i | w_0 + \underline{w}^T \underline{x}_i, \sigma^2)}_{\text{1st term}} + \underbrace{\sum_{j=1}^D \ln N(w_j | 0, \tau^2)}_{\text{2nd term}} \right\}$$

$$\hat{\underline{w}} = \underset{\underline{w}}{\operatorname{argmax}} \left\{ -\tilde{J}_R(\underline{w}, \sigma) \right\}$$

$\Rightarrow \tilde{J}_R$ IS OBJECTIVE FCN. (IN

UNSIMPLIFIED FORM)

1st term of $\tilde{J}_R(\underline{w}, \sigma)$:

$$-\sum_{i=1}^N \ln N(y_i | w_0 + \underline{w}^T \underline{x}_i, \sigma^2)$$

$$(4) \quad = -\cancel{\sum_{i=1}^N (\text{const. of } \underline{w})} + \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - (w_0 + \underline{w}^T \underline{x}_i))^2$$

2nd term of \tilde{J}_R :

$$- \sum_{j=1}^D \ln N(w_j | 0, \tau^2)$$

$$= - \sum_j \ln \left[\frac{1}{\sqrt{2\pi\tau^2}} \exp \left\{ -\frac{w_j^2}{2\tau^2} \right\} \right]$$

$$= + D \cdot \frac{1}{2} \ln(2\pi\tau^2) + \sum_{j=1}^D \frac{w_j^2}{2\tau^2}$$

$$(5) = \frac{1}{2\tau^2} \|\underline{w}\|_2^2$$

$$2\tau^2 [(4) + (5)] \Rightarrow$$

CORRECT VERSION OF
Eq. (7.32) \Leftarrow

$$\tilde{J}_R(\underline{w}, \sigma) = \sum_{i=1}^N \left[y_i - (\underline{w}_0 + \underline{w}^T \underline{x}_i) \right]^2 + \lambda \|\underline{w}\|_2^2$$

$\lambda = \frac{\sigma^2}{\tau^2}$

\uparrow TRAINING SAMPLE OUTPUT. $\hat{f}(\underline{x}_i)$

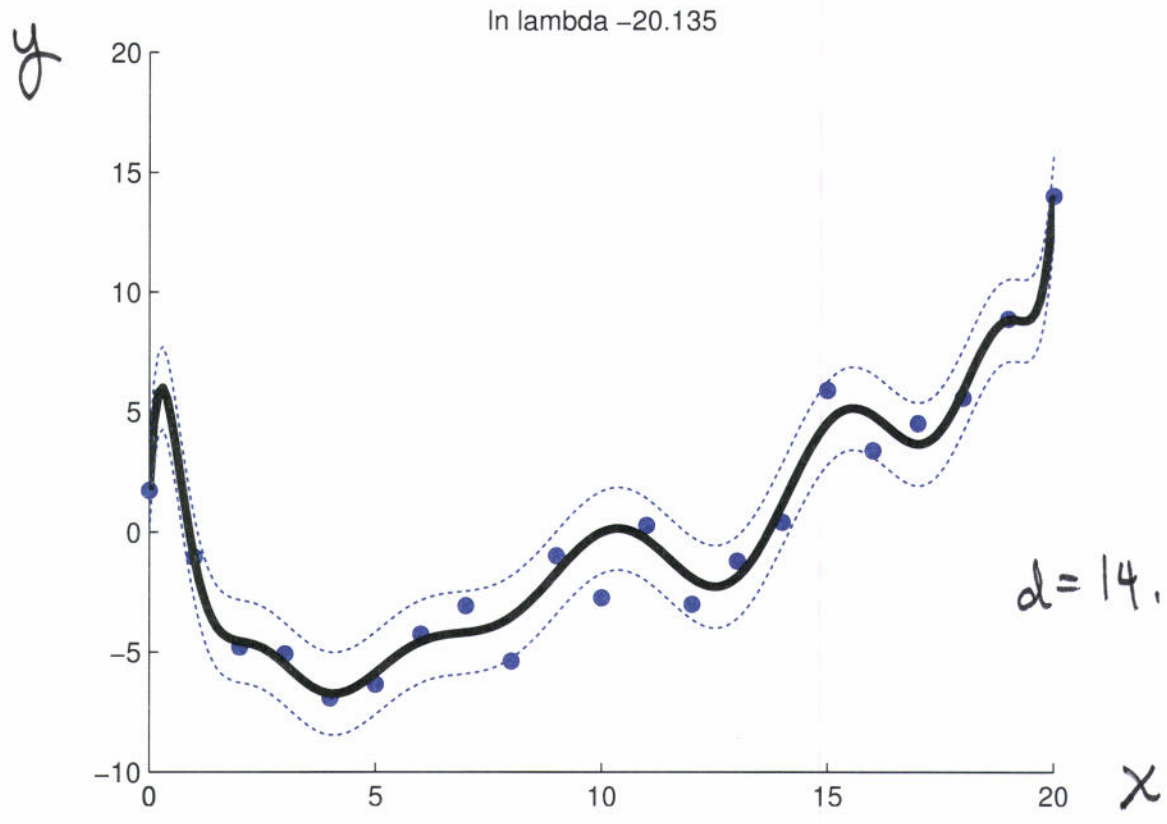
$\underbrace{\hspace{10em}}_{J_{MLE}} \quad \underbrace{\hspace{10em}}_{\text{REGULARIZER}} \quad \underbrace{\hspace{10em}}_{\text{new (prior term)}}$

\Leftarrow RIDGE REGRESSION OBJECTIVE FCN.

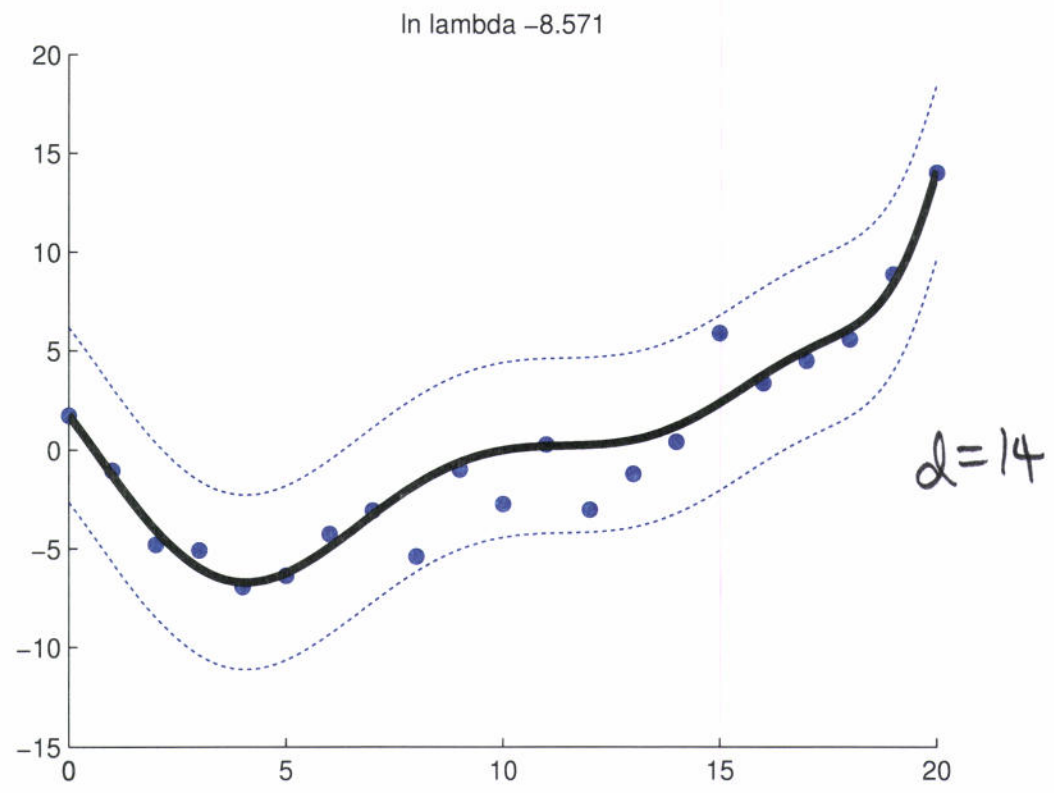
ASSUMED: NORMAL DENSITIES $p(y|\underline{x})$
 $p(\underline{w})$

& INDEP. OF w_j .

POINTS IN \mathcal{D} ARE i.i.d.



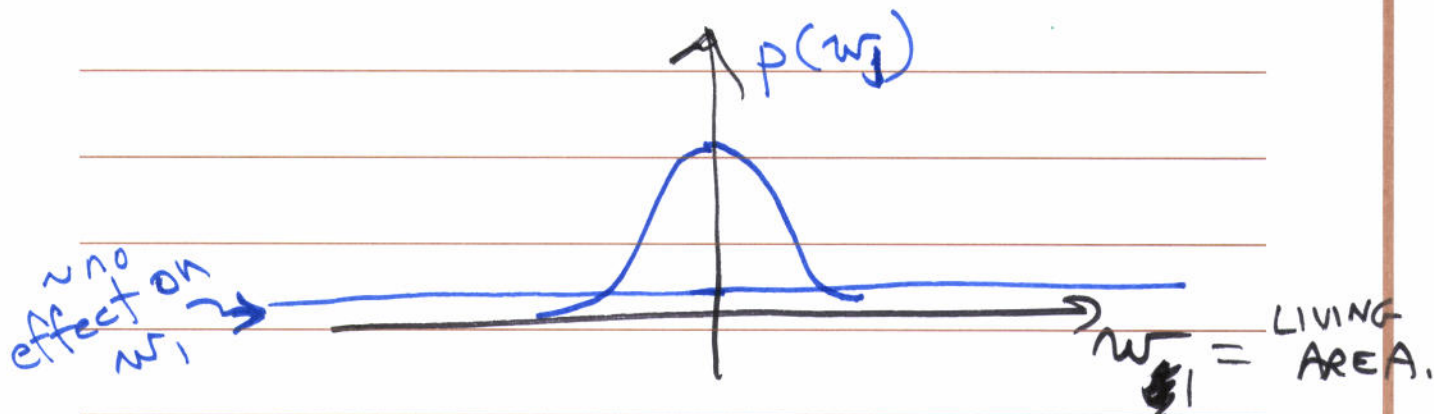
Murphy Fig. 7.7 (a)-(b). $N = 21$ data points, fit using regression function that is polynomial of degree 14, and differing amounts of L2 regularization.



5

$$\lambda = \frac{\text{VARIANCE OF DATA (FROM MODEL)}}{\text{VARIANCE OF PRIOR.}}$$

$$p(\underline{w}) = \prod_{j=1}^D N(w_j | 0, \tau^2)$$



[USING MORE DATA CAN HAVE A REGULARIZING-TYPE EFFECT — MURPHY FIG. 7.10]

NOTATION COMMENT

\underline{x} = GENERAL INPUT VARIABLE.

\underline{y} = " OUTPUT " (VALUE OR CLASS)

WHICH ARE DIFFERENT THAN:

VALUE OF DATAPPOINTS IN \mathcal{D} .

\underline{X} = INPUT VALUES OF DATASET \mathcal{D} (TRAINING DATA)

$$= \begin{bmatrix} \underline{x}_1^T \\ \underline{x}_2^T \\ \vdots \\ \underline{x}_N^T \end{bmatrix}$$

~~\underline{x}_i = INPUT VALUES OF~~

\underline{y}_i = OUTPUT VALUE FOR INPUT \underline{x}_i OF \mathcal{D} .

$$\underline{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

x_i = i^{th} COMPONENT OF FEATURE VECTOR \underline{x} .

BAYESIAN INFERENCE

WE WANT TO ESTIMATE $\underline{\theta}$

INSTEAD OF FINDING A POINT ESTIMATE $\hat{\underline{\theta}}$,
LET'S ESTIMATE THE DENSITY:

$$p(\underline{\theta} | \mathcal{D}).$$

WE HAVE A MODEL:

- (i) $p(y | \underline{x}, \underline{\theta})$ [REGRESSION]
(ii) OR $p(\underline{x} | y, \underline{\theta})$ [CLASSIFICATION].

TEST: $p(\underline{x} | S_i) \leftarrow$ ASSUME
A
MODEL.

(i) DISCRIMINATIVE APPROACH

MODELS $p(y | \underline{x}, \underline{\theta})$ DIRECTLY.

(ii) GENERATIVE APPROACH.

MODELS $p(y, \underline{x} | \underline{\theta})$

NOTE: MODELING $p(\underline{x} | y=c, \underline{\theta})$ IN
CLASSIFICATION, WE CAN:

$$p(y=c | \underline{x}, \underline{\theta}) = \frac{p(\underline{x} | y=c, \underline{\theta}) p(y=c)}{p(\underline{x})}$$

AND:

$$(a) \quad p(\underline{x}, y=c | \underline{\theta}) = p(y=c | \underline{x}, \underline{\theta}) p(\underline{x})$$

$$(b) \quad p_A = p(\underline{x} | y=c, \underline{\theta}) p(y=c)$$

$$\text{For (a), we can use } p(\underline{x} | \underline{\theta}) = \sum_{c=1}^C p(\underline{x} | y=c, \underline{\theta}) p(y=c | \underline{\theta})$$

$$\Rightarrow p(\underline{x}) = \sum_{c=1}^C p(\underline{x} | y=c, \underline{\theta}) p(y=c).$$