

Announcements

- Homework 1 was due today.
 - Homework 2 will be posted.
 - Students on waiting list can register soon
-

Today's Lecture

- Notation (data; augmented and non-augmented vectors)
- Comment on definition of dataset D ($y|x$ or y,x)
- Regression
 - Introduction
 - Based on MLE
 - Ridge regression (start)

Notation for augmented & unaugmented quantities

Non-augmented space

$$\underline{w} = \underline{w}^{(0)} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_D \end{bmatrix}$$

$$\underline{x} = \underline{x}^{(0)} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix}$$

Linear $\hat{f}(\underline{x}) = w_0 + \underline{w}^T \underline{x}$

Augmented space

$$\underline{w} = \underline{w}^{(+)} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_D \end{bmatrix}$$

$$\underline{x} = \underline{x}^{(+)} = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_D \end{bmatrix} = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_D \end{bmatrix}$$

Linear $\hat{f}(\underline{x}) = \underline{w}^T \underline{x}$

(dropping superscripts)

Similarly for $\underline{\phi}^{(0)}(\underline{x})$, $\underline{\phi}^{(+)}(\underline{x})$, and $\underline{\phi}(\underline{x})$.

REGRESSION [MURPHY Ch. 7]

STOCK PRICE EX:

LET $d=2$ AND $\underline{w} = \underline{w}^{(+)}$ (also $\Rightarrow \underline{x} = \underline{x}^{(+)}$,
 $\underline{\phi} = \underline{\phi}^{(+)}$).

$$\hat{f}(\underline{x}) = w_0 + w_1 x + w_2 x^2$$

$$\text{LET } \underline{\phi}(\underline{x}) = \begin{bmatrix} \phi_0(x) \\ \phi_1(x) \\ \phi_2(x) \end{bmatrix} = \begin{bmatrix} 1 \\ x \\ x^2 \end{bmatrix}$$

$$\begin{aligned} \hat{f}(\underline{x}) &= w_0 \phi_0(x) + w_1 \phi_1(x) + w_2 \phi_2(x) \\ &= \underline{w}^T \underline{\phi}(\underline{x}) \end{aligned}$$

"BASIS ~~SET~~ FUNCTION EXPANSION"

"NONLINEAR MAPPING"

" ϕ MACHINE"

4
Ex2: RENT OF APARTMENTS NEAR BEACH

$$\text{LET } \underline{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} \text{LIVING AREA} \\ \text{NO. OF ROOMS} \end{bmatrix}$$

$$\underline{w} = \underline{w}^{(+)} \text{ THROUGHOUT.}$$

LINEAR MODEL:

$$\hat{f}(\underline{x}) = w_0 + w_1 x_1 + w_2 x_2 = \underline{w}^T \underline{x}$$

NONLINEAR MODEL: (NONLINEAR IN \underline{x}).

$$\hat{f}(\underline{x}) = \underline{w}^T \underline{\phi}(\underline{x})$$

$$\text{Ex: QUADRATIC: } \hat{f}(\underline{x}) = \underline{x}^T \underline{W} \underline{x},$$

\underline{W} IS UPPER TRIANGULAR.

5
KEY ISSUE: HYPOTHESIS SET, OR MODEL [M7.2]

MODEL IS STATISTICAL.

$\underline{w} = \underline{w}^{(+)}$ THROUGHOUT

MODEL y AS:

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

↑ UNKNOWN PARAMETERS,
TO BE ESTIMATED
FROM \mathcal{D} .

→ MAKE ASSUMPTION ✖

(i) HERE: $p(y | \underline{x}, \underline{\theta}) = N(y | \underline{w}^T \underline{x}, \sigma^2)$
[LINEAR
IN \underline{x}]

~~(i)~~ (i)' OR $= N(y | \underline{w}^T \phi(\underline{x}), \sigma^2)$
[NONLINEAR IN \underline{x}]

EQUIVALENT TO:

$$y(\underline{x}) = \underline{w}^T \underline{x} + n,$$

$$n \sim N(s_y | 0, \sigma^2).$$

OR $\phi(\underline{x})$

REGRESSION USING MAXIMUM LIKELIHOOD ESTIMATE

$\underline{w} = \underline{w}^{(+)}$ THROUGHOUT

(MLE)

[M 7.3]

$p(\mathcal{D} | \underline{\theta}) = \text{LIKELIHOOD OF } \underline{\theta}$.

EST. $\underline{\theta}$ USING MLE:

$$\hat{\underline{\theta}}_{MLE} = \underset{\underline{\theta}}{\operatorname{argmax}} \ln p(\mathcal{D} | \underline{\theta})$$

IF WE ASSUME (i) OR (i)', WITH σ, ϕ KNOWN,

THEN:

$$\underline{\theta} = \underline{w}$$

AND: $\hat{\underline{w}}_{MLE} = \underset{\underline{w}}{\operatorname{argmax}} \ln p(\mathcal{D} | \underline{w})$

$$\sum_{i=1}^N \ln p(y_i | x_i, \underline{w})$$

ASSUMING POINTS IN \mathcal{D} ~~ARE~~

ARE INDEPENDENTLY DISTRIBUTED (i.d.).

2. OBJECTIVE FUNCTION

$$J(\underline{w}, \theta) = -\ln p(\mathcal{D}|\underline{w}) = \text{NLL}(\underline{w}).$$

(or $-p(\mathcal{D}|\underline{w})$)

LET:

$$\underline{X} = \begin{bmatrix} -\underline{x}_1^T - \\ -\underline{x}_2^T - \\ \vdots \\ -\underline{x}_N^T - \end{bmatrix} \text{ and } \underline{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

COMPRISE θ .

Using probability relations we have:

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

where for the last step we have dropped the last condition on $\underline{\theta}$ because it tells us nothing useful about $p(\underline{x})$. If we are interested in maximizing (or minimizing) the likelihood, we will take:

$$\begin{aligned} \arg \max_{\underline{\theta}} p(\mathcal{D}|\underline{\theta}) &= \arg \max_{\underline{\theta}} \left\{ \prod_{i=1}^N p(y_i, \underline{x}_i|\underline{\theta}) \right\} \\ &= \arg \max_{\underline{\theta}} \left\{ \prod_{i=1}^N p(y_i|\underline{x}_i, \underline{\theta}) p(\underline{x}_i) \right\} = \arg \max_{\underline{\theta}} \left\{ \left(\prod_{i=1}^N p(\underline{x}_i) \right) \prod_{i=1}^N p(y_i|\underline{x}_i, \underline{\theta}) \right\} \\ &= \arg \max_{\underline{\theta}} \left\{ \prod_{i=1}^N p(y_i|\underline{x}_i, \underline{\theta}) \right\} \end{aligned}$$

and to obtain the last line, $\prod_{i=1}^N p(\underline{x}_i)$ was dropped because it is a positive multiplicative term that is a constant of $\underline{\theta}$. This can equivalently be seen by using the log likelihood instead:

$$\begin{aligned} \arg \max_{\underline{\theta}} p(\mathcal{D}|\underline{\theta}) &= \arg \max_{\underline{\theta}} \left\{ \ln \prod_{i=1}^N p(y_i, \underline{x}_i|\underline{\theta}) \right\} \\ &= \arg \max_{\underline{\theta}} \left\{ \ln \prod_{i=1}^N p(y_i|\underline{x}_i, \underline{\theta}) p(\underline{x}_i) \right\} \\ &= \arg \max_{\underline{\theta}} \left\{ \sum_{i=1}^N \ln [p(y_i|\underline{x}_i, \underline{\theta}) p(\underline{x}_i)] \right\} \\ &= \arg \max_{\underline{\theta}} \left\{ \sum_{i=1}^N [\ln p(y_i|\underline{x}_i, \underline{\theta}) + \ln p(\underline{x}_i)] \right\} \\ &= \arg \max_{\underline{\theta}} \left\{ \sum_{i=1}^N \ln p(y_i|\underline{x}_i, \underline{\theta}) \right\} \end{aligned}$$

and to obtain the last line, the additive terms that don't depend on $\underline{\theta}$ have been dropped.

So, when the goal of using the likelihood is to find its argmax or argmin w.r.t. $\underline{\theta}$, we can replace $p(y_i, \underline{x}_i|\underline{\theta})$ directly with $p(y_i|\underline{x}_i, \underline{\theta})$.

$$\text{Using } p(\mathcal{Y}/\underline{\theta}) = \prod_{i=1}^N p(y_i | \underline{x}_i, \underline{\theta})$$

$$= \prod_{i=1}^N N(y_i | \underline{w}^T \underline{x}_i, \sigma^2)$$

↑ plug in

$$\Rightarrow J_1(\underline{w}) = \underbrace{\frac{1}{2\sigma^2}}_{\text{Drop}} \sum_{i=1}^N (y_i - \underline{w}^T \underline{x}_i)^2 + \underbrace{\frac{N}{2} \log(2\pi\sigma^2)}_{\text{const. of } \underline{w}}$$

$$\text{Let } J(\underline{w}) = \frac{1}{2} \text{RSS}(\underline{w})$$

$$= \frac{1}{2} \sum_{i=1}^N (y_i - \underline{w}^T \underline{x}_i)^2$$

$$= \frac{1}{2} \|\underline{y} - \underline{X} \underline{w}\|_2^2$$

$$= \frac{1}{2} (\underline{y} - \underline{X} \underline{w})^T (\underline{y} - \underline{X} \underline{w})$$

3. OPTIMIZATION METHOD

$$\nabla_{\underline{w}} J(\underline{w}, \underline{y}) = \underline{0} \quad \& \text{ SOLVE ALGEBRAICALLY.}$$

SOLVING GIVES $\hat{\underline{w}}$:

$$\underline{\underline{X}}^T \underline{\underline{X}} \hat{\underline{w}} = \underline{\underline{X}}^T \underline{\underline{y}}$$

IF $(\underline{\underline{X}}^T \underline{\underline{X}})$ IS INVERTABLE, THEN

$$\hat{\underline{w}}_{\text{OLS}} = \underline{\underline{X}}^- \underline{\underline{y}} = (\underline{\underline{X}}^T \underline{\underline{X}})^{-1} \underline{\underline{X}}^T \underline{\underline{y}}$$

= ~~ORD~~ ORDINARY LEAST SQUARES
SOLUTION.

$\underline{\underline{X}}^-$ = PSEUDOINVERSE OF $\underline{\underline{X}}$.

~~OR~~ $\hat{\underline{w}}_{\text{OLS}}$ FOLLOWS FROM THE NORMAL
DENSITY ASSUMPTION (AND i.i.d.)

RIDGE REGRESSION

($\underline{w} = \underline{w}^{(0)}$ THROUGHOUT)

→ USE MAP ESTIMATE

$$\hat{\underline{\theta}}_{\text{MAP}} = \underset{\underline{\theta}}{\operatorname{argmax}} p(\underline{\theta} | \mathcal{D})$$

$$= \underset{\underline{\theta}}{\operatorname{argmax}} \left\{ \frac{p(\mathcal{D} | \underline{\theta}) p(\underline{\theta})}{p(\mathcal{D})} \right\}$$

DROP.

$$= \underset{\underline{\theta}}{\operatorname{argmax}} \{ p(\mathcal{D} | \underline{\theta}) p(\underline{\theta}) \}$$

$$(1) \hat{\underline{\theta}}_{\text{MAP}} = \underset{\underline{\theta}}{\operatorname{argmax}} \left\{ \underbrace{\ln p(\mathcal{D} | \underline{\theta})}_{\text{likelihood of } \underline{\theta}} + \underbrace{\ln p(\underline{\theta})}_{\text{prior}} \right\}$$

Assume (x_i, y_i) ARE i.i.d., FROM $N(y | w_0 + \underline{w}^T \underline{x}, \sigma^2)$

$$(2) \text{ so: } \ln p(\mathcal{D} | \underline{\theta}) = \ln \prod_i N(y_i | w_0 + \underline{w}^T \underline{x}_i, \sigma^2)$$

→ MLE TERM

FOR PRIOR, CHOOSE:

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

$$(3) \ln p(\underline{\theta}) = \sum_{j=1}^D \ln N(w_j | 0, \tau^2)$$

WHY A GAUSSIAN PRIOR?

1. ALGEBRAICALLY CONVENIENT

2. CAN CHOOSE τ (NARROW OR WIDE GAUSSIAN)

3. WILL PREFER w_j VALUES CLOSE TO 0.

BUT NOT NECESSARILY THE BEST CHOICE.