## Announcements

- Homework 2 is due tomorrow

- Homework 3 will be posted

## Today's Lecture

- Bayesian inference (regression)

- Summary of estimation techniques

页2

# BAYESIAN INFERENCE

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

INSTEAD OF FINDING A POINT ESTIMATE $\hat{\underline{\theta}}$,

LET'S ESTIMATE THE DENSITY:

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

WE HAVE A MODEL:

(i) $\quad p(y \mid \underline{x}, \underline{\theta})$    US: [REGRESSION] [DISCRIMINATIVE]

(ii) OR $\quad p(\underline{x} \mid y, \underline{\theta})$    US: [CLASSIFICATION]. [GENERATIVE]

---

EE559: $\quad p(\underline{x} \mid S_i) \Leftarrow$ ASSUME A MODEL.

(i) DISCRIMINATIVE APPROACH

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

(ii) GENERATIVE APPROACH.

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

     NOTE: MODELING $p(\underline{x} \mid y = c, \underline{\theta})$ IN

     CLASSIFICATION, WE CAN:

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

AND:

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

(b)
$$OR = p(\underline{x} \mid y=c, \underline{\theta}) \, p(y=c)$$

FOR (a), WE CAN USE $p(\underline{x} \mid \underline{\theta}) = \sum_{c=1}^{C} p(\underline{x} \mid y=c, \underline{\theta}) \, p(y \mid \underline{\theta})$

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

## BAYESIAN INFERENCE (part 2)

$$\underline{w} = \underline{w}^{(+)}. \quad \text{[not always true in Murphy 7.6]}.$$

GIVEN OUR MODEL, WE CAN COMPUTE THE LIKELIHOOD:

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

USE BAYES THEOREM:

(1)
$$p(\underline{\theta}|\mathcal{D}) = \frac{p(\mathcal{D}|\underline{\theta}) \, p(\underline{\theta})}{p(\mathcal{D})}$$

IN WHICH: $p(\mathcal{D}) = \int p(\mathcal{D}|\underline{\theta}) \, p(\underline{\theta}) \, d\underline{\theta}$

(OR SUM IF $\underline{\theta}$ IS DISCRETE)

IN Ch.3 READING, $\underline{\theta} = h = $ HYPOTHESIS.

IN REGRESSION, $\underline{\theta} = \underline{w}$ (AND MAYBE $\sigma^2$).

EXAMPLES OF (1):

- NUMBERS GAME (HW2, M Ch.3)

$$h \in \{h_2, h_4\}$$

$$p(h \mid \not{o}) = \frac{p(\not{o} \mid h) \, p(h)}{p(\not{o})}$$

LIKELIHOOD
(use strong
sampling
assumpt.)

$$p(\not{o}) = \sum_{h} p(\not{o} \mid h) \, p(h)$$

PRIOR.

- REGRESSION (APT. RENT EX.)

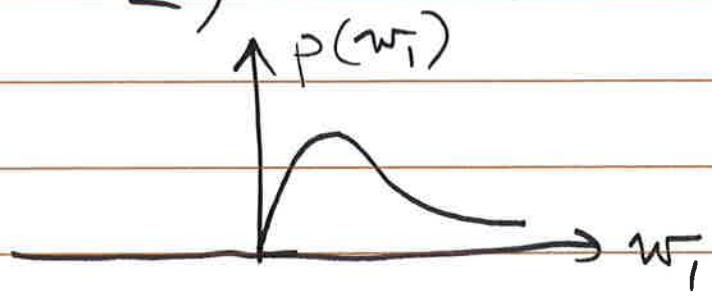  $\underline{\theta}$ = COEFFICIENTS (WEIGHTS) ON INPUTS
  (LIVING AREA, #ROOMS, ETC.)

$|\theta_i|$ = IMPORTANCE OF FEATURE $x_i$. IF
ALL ~~X-SCALE~~ $x_j$'S ARE NORMALIZED
TO SAME RANGE, (e.g., STANDARIZED
TO $\mu_j = 0$, $\sigma_j^2 = 1$), $j = 1, \cdots, D$.

$p(\underline{\theta} \mid \not{o})$ = DENS. OF $w_j$ COEFF. OF $x_j$.

$$= \left[ \prod_{i=1}^{N} p(y_i \mid \underline{x}_i, \underline{\theta}) \right] p(\underline{\theta}) \Big/ k$$

USE OUR MODEL, E.G.: $N(y_i \mid \underline{w}^T \underline{x}_i, \sigma^2)$.

$$p(\underline{\theta}) = \text{PRIOR ON } \underline{\theta}, \quad \text{E.G. MAYBE:}$$



THEN GET POSTERIOR PREDICTIVE

$$p(y \mid \underline{x}, \mathcal{D})$$

FROM: $p(y) = \int p(y \mid \underline{\theta}) p(\underline{\theta}) \, d\underline{\theta}$

WE GET!

(2)
$$p(y \mid \underline{x}, \mathcal{D}) = \int p(y \mid \underline{x}, \underline{\theta}, \mathcal{D}) \, p(\underline{\theta} \mid \underline{x}, \mathcal{D}) \, d\underline{\theta}$$

$$\text{USUALLY!} = \int p(y \mid \underline{x}, \underline{\theta}) \, p(\underline{\theta} \mid \mathcal{D}) \, d\underline{\theta}$$

OUR MODEL
(REGR, DISCRIMINATIVE)

$p(\underline{w} \mid \mathcal{D})$
— GET FROM (1).

OR CAN BE OBTAINED FROM
$p(\underline{x} \mid y, \underline{\theta})$ (CLASS'N, GENER-
ATIVE)

$y = \$$

$\hat{f}(x)$ OR $E\{y\}$

o POINT IN $\mathcal{D}$.

$p(y | \tilde{x}, \mathcal{w})$

$x_1 = $ LIVING AREA

$\tilde{x}$

# BAYESIAN REGRESSION

1. MODEL IS $p(y | \underline{x}, \underline{\theta}) = p(y | \underline{x}, \underline{w}, \underline{\theta}')$

$\underline{\theta}'$ → any other unknowns.

EX: $\qquad = p(y | \underline{w}^T \underline{x})$ (LINEAR CASE)

$\qquad = p(y | \underline{w}^T \underline{\phi}(x))$

## 2. PARAMETER POSTERIOR

FROM EQ. (1):

$$p(\theta | \mathcal{D}) = p(\underline{w} | \underline{y}, \underline{X})$$

Likelihood

prior

$$= \frac{p(\underline{y} | \underline{w}, \underline{X}) \, p(\underline{w} | \underline{X})}{k}$$

$$k = \int p(\underline{y} | \underline{w}, \underline{X}) \, p(\underline{w} | \underline{X}) \, d\underline{w}$$

## 3. POSTERIOR PREDICTIVE

FROM EQ. (2):

$$p(y | \underline{x}, \mathcal{D}) = \int p(y | \underline{x}, \underline{w}, \mathcal{D}) \, p(\underline{w} | \underline{x}, \mathcal{D}) \, d\underline{w}$$

$$= \int p(y | \underline{x}, \underline{w}) \, p(\underline{w} | \mathcal{D}) \, d\underline{w}.$$

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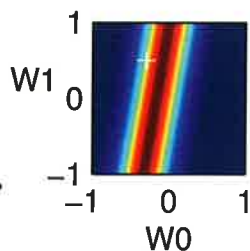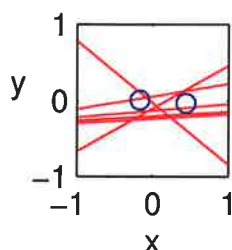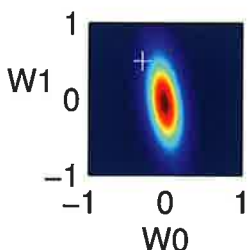


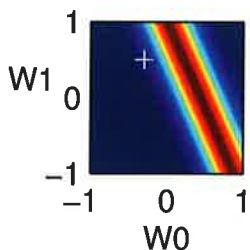likelihood          prior/posterior          data space

no data.

Prior $p(\underline{w})$

POSTERIOR $p(\underline{w}|\mathcal{D})$
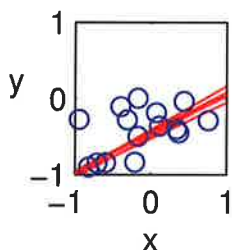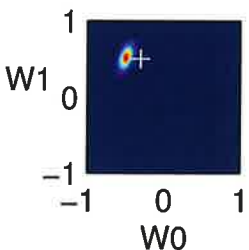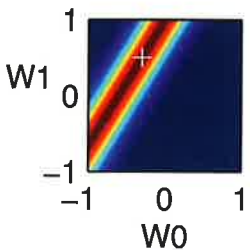
$N=1$ DATA POINT.

$N=2$

$N=20$

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

Murphy Fig. 7.11.

Linear regr. w/ scalar input $x_1$.

$$\underline{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$

$$y(x_1, \underline{w}) = w_0 + w_1 x_1 + \epsilon$$

$\hookleftarrow$ noise term.

# SUMMARY OF ESTIMATION TECHNIQUES

MLE: $\quad \hat{\underline{\theta}}_{MLE} = \underset{\underline{\theta}}{\text{argmax}} \left\{ \ln p(\mathcal{D}|\underline{\theta}) \right\}$

$\qquad$ GAUSSIAN CASE $\Longrightarrow$ $J_{MLE}(\underline{w}, \mathcal{D}) = MSE$

MAP: $\quad \hat{\underline{\theta}}_{MAP} = \underset{\underline{\theta}}{\text{argmax}} \left\{ \ln p(\mathcal{D}|\underline{\theta}) + \ln p(\underline{\theta}) \right\}$

$\qquad$ GAUSSIAN CASE $\Longrightarrow$ $J_{MAP}(\underline{w}, \mathcal{D}) =$
$\qquad\qquad\qquad\qquad N \cdot (MSE) + \lambda \|\underline{w}\|_2^2$

BAYESIAN: $\quad p(\underline{\theta}|\mathcal{D}) \propto p(\mathcal{D}|\underline{\theta}) p(\underline{\theta})$

$\Longrightarrow p(y|\underline{x}, \mathcal{D}) = \int p(y|\underline{x}, \underline{\theta}) p(\underline{\theta}|\mathcal{D}) d\underline{\theta}$

$\quad$ POSTERIOR PREDICTIVE