

## Announcements

- HW7 (project proposals) are due this Friday.
  - Collaborative Kaggle topic is posted on piazza, with comments
- Midterm is two weeks from today

---

## Today's Lecture

- Bayesian and MAP feature selection
- Sparsity and regularization

# BAYESIAN AND MAP FEATURE SELECTION

## WHY FEATURE SELECTION?

- TO BETTER BALANCE  $N$  WITH  $\mathcal{H}$  COMPLEXITY  
(dvc, or  $\sim$  #d.o.f.)
- DISCOVER WHICH FEATURES ARE MOST  
IMPORTANT FOR PREDICTION.

1. USE CONTINUOUS VARIABLES  $\underline{w}$ .

EX: APARTMENT RENTS

$x_1 = \overset{\text{AREA}}{\text{LIVING ROOM (sq.f.)}}: 500 \leq x_1 \leq 3500$

$x_2 = \# \text{ of ROOMS}: 1 \leq x_2 \leq 12$

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

$$\nwarrow \quad \nearrow$$

$$|w_1| \Leftrightarrow |w_2|$$

CAN USE  $|w_i| \Leftrightarrow |w_j|$  ONLY IF

DATA ( $x_i$ ) ARE APPROPRIATELY

NORMALIZED (e.g., STANDARDIZED)

POSE AS A MAP ESTIMATION PROBLEM.

$$\begin{aligned}
 p(\underline{w} | \mathcal{D}) &= \frac{1}{K} p(\mathcal{D} | \underline{w}) p(\underline{w}) \\
 &= \frac{1}{K} \left[ \prod_{i=1}^N p(y_i | \underline{x}_i, \underline{w}) \right] p(\underline{w}) \\
 (1) \left\{ \begin{aligned} &= \frac{1}{K} p(\mathcal{D} | \underline{X}, \underline{w}) p(\underline{w}) \end{aligned} \right.
 \end{aligned}$$

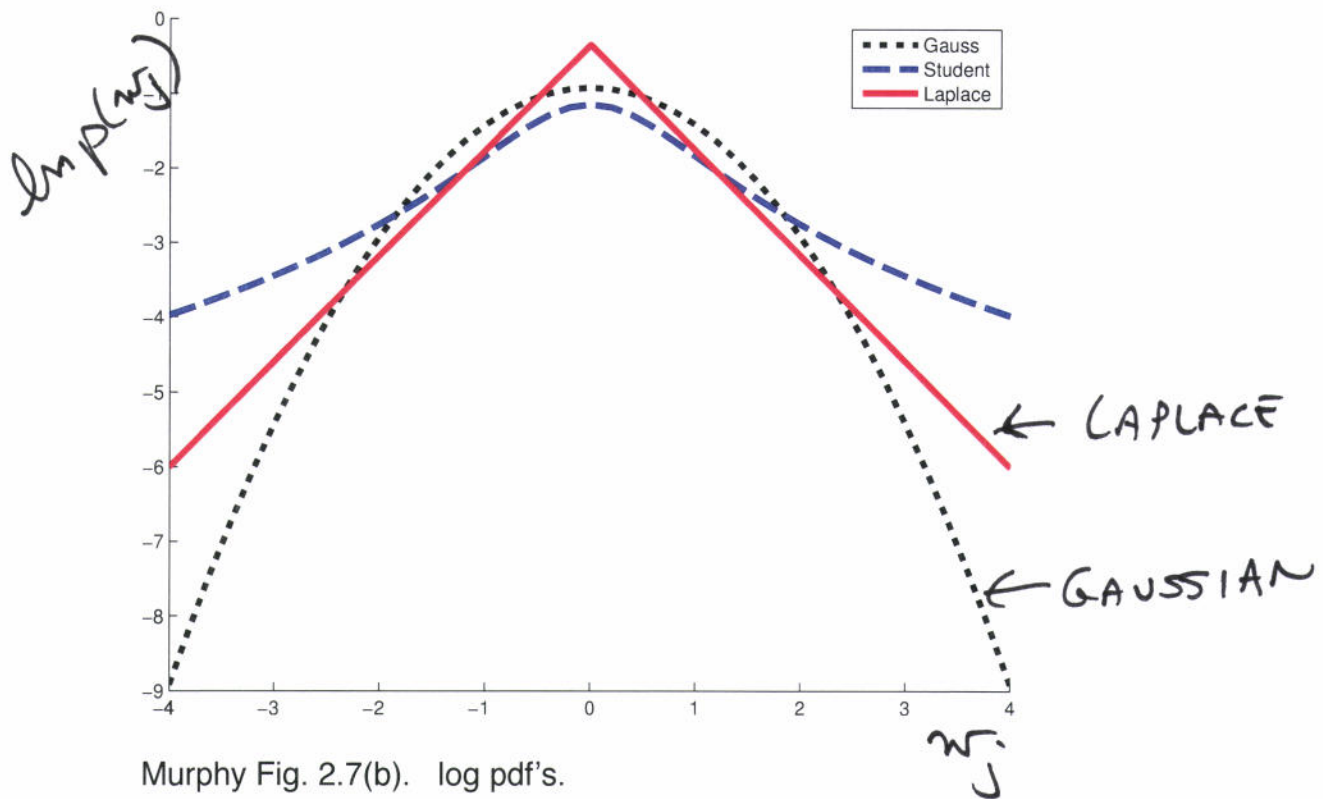
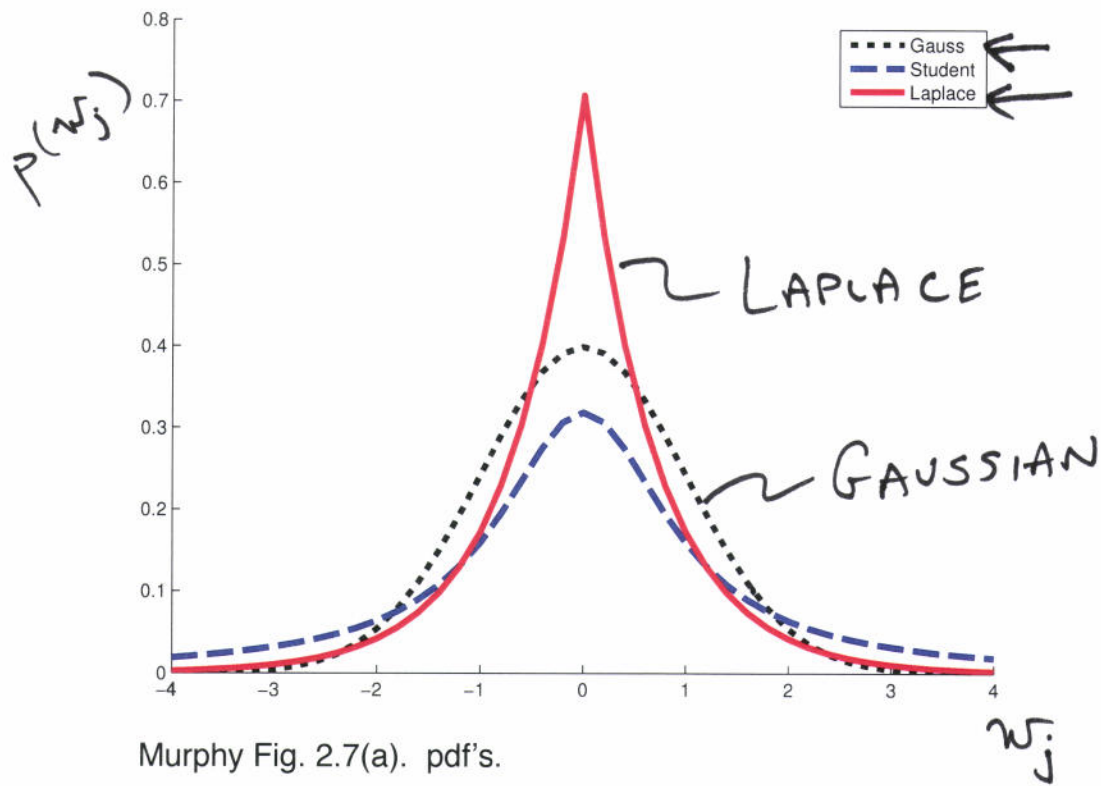
IF LINEAR REGRESSION w/ USUAL ASSUMPTION:

$$p(y | \underline{x}, \underline{w}) = N(y | \underline{w}^T \underline{x}, \sigma^2)$$

PRIOR TERM:

$$\begin{aligned}
 p(\underline{w}) &= p(\underline{w} | \lambda) \\
 &= \prod_{j=1}^D \text{Lap}(w_j | 0, \frac{1}{\lambda}) \\
 &= \prod_{j=1}^D \frac{\lambda}{2} e^{-\lambda |w_j|} \propto \prod_{j=1}^D e^{-\lambda |w_j|}
 \end{aligned}$$

$$p(w_0) \propto 1 \quad (\text{UNIFORM}).$$





(1)  $\Rightarrow$ 

$$J(\underline{w}) = -\ln p(\underline{y} | \underline{X}, \underline{w}) - \ln p(\underline{w} | \lambda)$$

$$= \text{NLL}(\underline{w}) + ~~\lambda \|\underline{w}\|_1~~ \lambda \|\underline{w}\|_1,$$

$$\|\underline{w}\|_1 \triangleq \sum_{j=1}^D |w_j|.$$

FOR LINEAR MODEL:

$$= \sum_{i=1}^N \frac{1}{2\sigma^2} [y_i - (w_0 + \underline{w}^T \underline{x}_i)]^2 + \lambda \|\underline{w}\|_1,$$

$$J(\underline{w}) \propto \text{RSS}(\underline{w}) + \lambda' \|\underline{w}\|_1, \quad \lambda' = 2\lambda\sigma^2.$$

↖ BASIS PURSUIT DENOISING (BPDN) [M13.3.0]

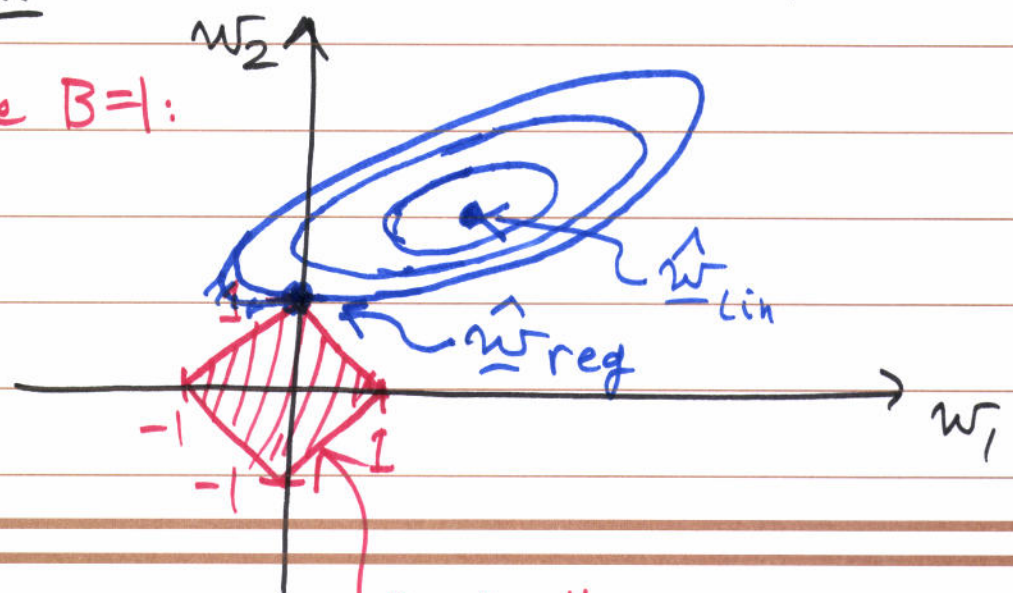
OR LASSO MINIMIZES THIS  $J(\underline{w})$  USING  
SUBGRADIENTS (M 13.3.2)

# $l_1$ REGULARIZATION AND SPARSITY [M 13.3.1]

→ LOOK AT AS CONSTRAINED MIN. PROBLEMS:

$$\min_{\underline{w}} \text{RSS}(\underline{w}) \quad \text{s.t.} \quad \|\underline{w}\|_1 \leq B \quad (\text{LASSO})$$

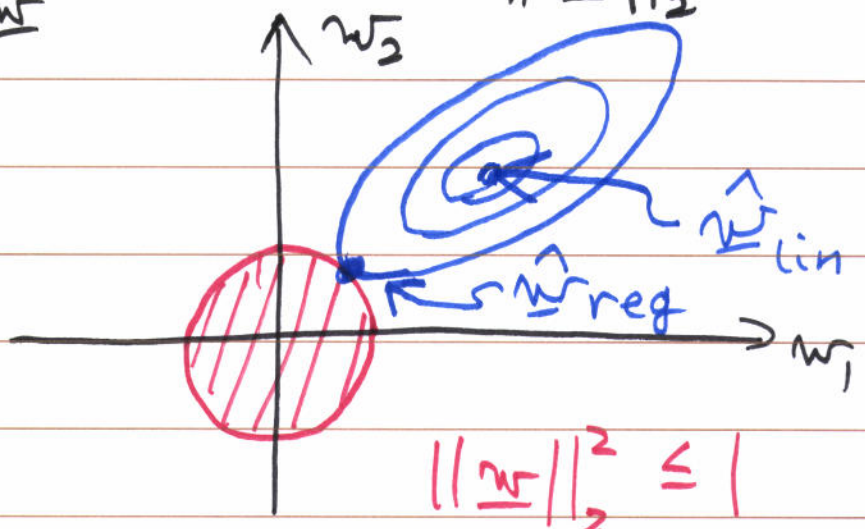
choose  $B=1$ :



$$\|\underline{w}\|_1 \leq 1$$

COMPARE WITH:

$$\min_{\underline{w}} \text{RSS}(\underline{w}) \quad \text{s.t.} \quad \|\underline{w}\|_2^2 \leq B \quad (\text{RIDGE})$$



$$\|\underline{w}\|_2^2 \leq 1$$

7  
EXAMPLE: PREDICTING PSA FROM  
VARIOUS PARAMETERS.

$\mathcal{D}$ : SET OF DATA FROM PROSTATE  
CANCER PATIENTS (97).

— MURPHY TABLE 13.1 —

— HASTIE et al, TABLE 3.3 —

2. BAYESIAN FEATURE (VARIABLE) SELECTION [M 13.2]

$\underline{x} \text{ }_{D \times 1}$ :  $x_j = \begin{cases} 1 & \text{FEATURE } j \text{ IS RELEVANT} \\ 0 & \text{OTHERWISE} \end{cases}$

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

$$p(\underline{x} | \mathcal{D}) = \frac{p(\mathcal{D} | \underline{x}) p(\underline{x})}{p(\mathcal{D})}$$

PRIOR TERM  $p(\underline{x})$ ?



8

ASSUME:  $p(\underline{x}) = \prod_{j=1}^D p(x_j)$

COMMON CHOICE:

$$p(x_j) = \text{Ber}(x_j | \pi_0) = \pi_0^{x_j} (1 - \pi_0)^{(1 - x_j)}$$

$\pi_0$  = PROBABILITY THAT A FEATURE IS RELEVANT.

$$\Rightarrow p(\underline{x}) = \pi_0^{\sum x_j} (1 - \pi_0)^{D - \sum x_j}$$

LET  $\|\underline{x}\|_0 \triangleq \sum_{j=1}^D x_j = \# \text{ OF NONZERO ELEMENTS IN } \underline{x}.$

=  $\ell_0$  PSEUDO-NORM.

$$p(\underline{x}) = \pi_0^{\|\underline{x}\|_0} (1 - \pi_0)^{D - \|\underline{x}\|_0}$$