## Announcements

- HW7 (project proposal) is due tomorrow.

- HW8 will be posted.

## Today's Lecture

- Bayesian feature selection

  • Spike and slab model

  • Bernoulli-Gaussian model

  • L0 regularization

- p-norm and Bridge Regression

THEN:

$$\log p(\underline{\gamma} \mid \pi_0) = \|\underline{\gamma}\|_0 \log \pi_0 + (D - \|\underline{\gamma}\|_0)$$
$$\cdot \log (1 - \pi_0)$$

$$= \|\underline{\gamma}\|_0 \left[ \log \pi_0 - \log (1 - \pi_0) \right]$$
$$+ D \log (1 - \pi_0)$$

$$= -\lambda \|\underline{\gamma}\|_0 + C', \qquad \lambda \triangleq \log \left( \frac{1 - \pi_0}{\pi_0} \right)$$

ASSUME: $\bar{x}_j = 0, \quad \bar{y} = 0$

$$p(\underline{\mathscr{D}} \mid \underline{\gamma}) = p(\underline{y} \mid \underline{\underline{X}}, \underline{\gamma})$$

$$= \int p(\underline{y} \mid \underline{\underline{X}}, \underline{w}, \underline{\gamma}) \underbrace{p(\underline{w} \mid \underline{\gamma})}_{\text{prior on } \underline{w}.} d\underline{w}$$

$$p(\underline{w} \mid \underline{\gamma}) = \prod_{j=1}^{D} p(w_j \mid \gamma_j)$$

$$= \prod_{j=1}^{D} \begin{cases} \delta_0(w_j), & \text{IF } \gamma_j = 0 \\ N(w_j \mid 0, \sigma^2 \sigma_w^2), & \text{IF } \gamma_j = 1. \end{cases}$$

$\sigma^2$ = VARIANCE OF $y$.

$\sigma_w^2$ = VARIANCE OF RELEVANT FEATURE $(w_j)$
RELATIVE TO $\sigma^2$.



$p(w_j \mid \gamma_j = 0)$   $p(w_j \mid \gamma_j = 1)$

$w_j$   $w_j$

$\mathcal{C}$ "SPIKE AND SLAB" MODEL

— [MURPHY FIG. 13.1 (a), (c)] —

TYPICALLY $\sigma_w^2$ = LARGE $(\gg 1)$, BUT DEPENDS
ON REGULARIZATION DESIRED ON $\gamma_j = 1$ (RELEVANT)
FEATURES.

SINCE THIS APPROACH EFFECTIVELY COMPUTES THE
FULL $p(\underline{\gamma} \mid \mathcal{D})$, IT IS:

    - THOROUGH

      - COMPUTATIONALLY DEMANDING, OR
               APPROXIMATE / SUB-
               OPTIMAL.

# $l_0$ REGULARIZATION

SLIGHTLY DIFFERENT VIEW:

COMBINE $\underline{\gamma}$ AND $\underline{w}$ INTO 1 VECTOR.

MODEL:

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

WITH: $\underline{w}_\gamma = \begin{bmatrix} \gamma_1 w_1 \\ \gamma_2 w_2 \\ \vdots \\ \gamma_D w_D \end{bmatrix}$

($\underline{\gamma}$ IS A BINARY MASK FOR $\underline{w}$).

$$p(\underline{w}, \underline{\gamma} \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \underline{w}, \underline{\gamma}) \, p(\underline{w}, \underline{\gamma})}{p(\mathcal{D})}$$

LIKELIHOOD $p(\mathcal{D} \mid \underline{w}, \underline{\gamma})$

$$= N(\underline{y} \mid \underline{\underline{X}} \, \underline{w}_\gamma, \sigma^2 \underline{\underline{I}})$$

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

Joint prior $p(\underline{w}, \underline{\gamma}) = ?$

$$p(\underline{w}, \underline{\gamma}) = p(\underline{w}) \, p(\underline{\gamma})$$

$$= N(\underline{w} \mid \underline{0}, \sigma_w^2 \underline{\underline{I}}) \, \pi_0^{\|\underline{\gamma}\|_0} (1-\pi_0)^{D-\|\underline{\gamma}\|_0}$$

$\hookleftarrow$ <u>Bernouli-Gaussian model</u>

$$f_{obj}(\underline{w}, \underline{\gamma}) \propto -\log p(\underline{w}, \underline{\gamma} \mid \mathcal{D})$$

(drop const.'s of $\underline{w}, \underline{\gamma}$).

Take $\sigma_w^2 \to \infty$, so don't regularize $\gamma_j = 1$

weights

<u>Can show:</u>

$$f_{obj}(\underline{w}, \underline{\gamma}) = \|\underline{y} - \underline{\underline{X}} \underline{w}_\gamma\|_2^2 + \lambda \|\underline{\gamma}\|_0$$

$$\downarrow f_{obj}(\underline{w}_\gamma) = \|\underline{y} - \underline{\underline{X}} \underline{w}_\gamma\|_2^2 + \lambda \|\underline{w}_\gamma\|_0$$

Redefine $\underline{w} \stackrel{\triangle}{=} \underline{w}_\gamma$.

Then:

$$\boxed{f_{obj}(\underline{w}) = \|\underline{y} - \underline{\underline{X}} \underline{w}\|_2^2 + \lambda \|\underline{w}\|_0}$$
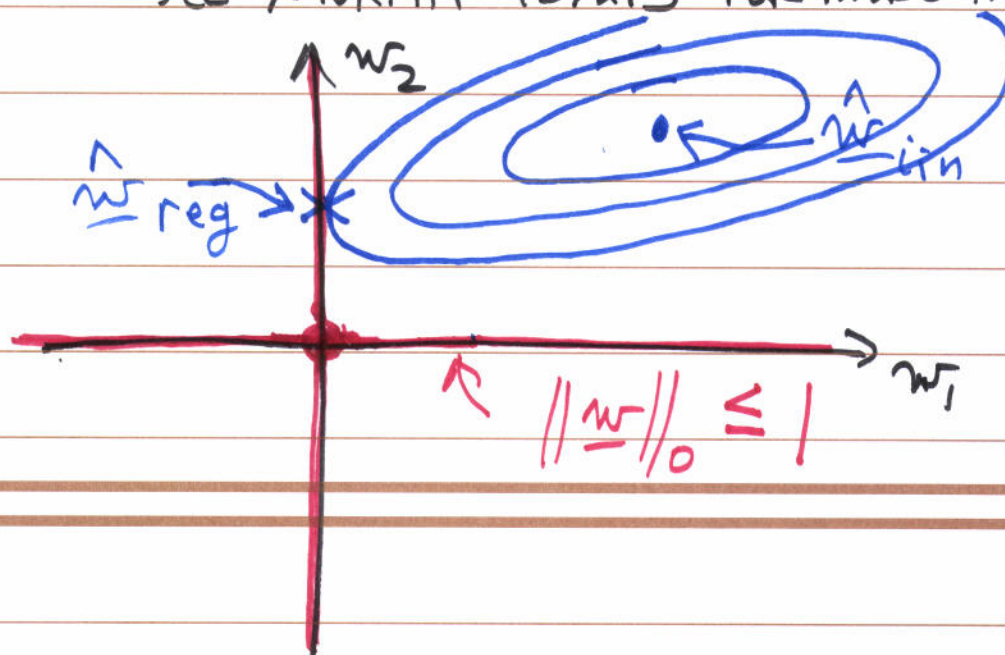
COMMENTS:

1. MIN. OVER CONTINUOUS VARIABLES $\underline{w}$ NOW.

2. Is $f_{obj.}$ CONVEX?

→ No, $||\underline{w}||_0$ IS NOT CONVEX.

→ SEE MURPHY 13.2.3 FOR MORE INFO [N.R.F.]

$$\hat{\underline{w}}_{reg} \qquad \hat{\underline{w}}_{lin}$$

$w_2$

$w_1$

$||\underline{w}||_0 \leq 1$

# BRIDGE REGRESSION

GENERALIZE $l_0, l_1, l_2$ REGULARIZERS / CONSTRAINTS.

EST. $\underline{w}$ USING $MAP$ WITH AN

EXPONENTIAL POWER DISTRIBUTION AS PRIOR:

$$ExpPwr\left(w_j \,|\, \mu_j, a, b\right)$$

$$= \frac{b}{2a\,\Gamma\left(\frac{1}{b}\right)} \exp\left\{- \frac{|x_j - \mu_j|^b}{a^b}\right\}$$

$\longrightarrow$ Murphy Eq. (13.132) may have errors.

LEADS TO AN ESTIMATE:

$$\hat{\underline{w}} = \underset{\underline{w}}{argmin}\left\{ NLL(\underline{w}) + \lambda \sum_{j=1}^{D} |w_j|^b \right\}, \quad b \geq 0.$$

LIKE AN "$l_b$" REGULARIZER,
OR "$p-norm$" BASED
REGULARIZER.

EX: $b=0$: $\ell_0$ REG., (IF WE DEFINE $|0|^0 = 0$).

$b=1$: $\ell_1$ REG.,

$b=2$: $\ell_2$ REG.

$$p\text{-norm} = \left[\sum_{j=1}^{D} |w_j|^p\right]^{1/p} \quad \text{(IS A NORM FOR } p \geq 1\text{)}.$$

$$\sum_{j=1}^{D} |w_j|^b \quad \text{IS CONVEX FOR } b \geq 1.$$

EX: $\quad f_{obj}(\underline{w}) = -\log p(\underline{w} \mid \mathcal{D})$

$$\text{or} \quad \|\underline{y} - \underline{X}\,\underline{w}\|_2^2 + \lambda \sum_{j=1}^{D} |w_j|^b$$

$\mathcal{D}$: ONE POINT $\underline{x}_1, y_1$

$$\Rightarrow \quad f_{obj}(\underline{w}) \propto \left[y_1 - (w_1 x_{11} + w_2 x_{12})\right]^2$$

given data point

$$+ \lambda \left[|w_1|^b + |w_2|^b\right].$$

$b=2$　　　　　$b=1$　　　　　$0<b<1$



PRIORS

$w_2$

$w_1$

$\hat{\underline{w}}_{Lin}$

$\hat{\underline{w}}_{Lin}$

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

Murphy  Fig. 13.7