

IV. Group Law

(Ch. 4 in PB&vdG (70-11))

$$Y = X\beta^0 + \varepsilon$$

estimated (or true) β is a high-dim. parameter vector
with a group structure

$$\beta = (\beta_{g_1}, \beta_{g_2}, \dots, \beta_{g_p}), \quad \beta_{g_j} = \{\beta_r : r \in g_j\}$$
$$\bigcup_{j=1}^p g_j = \{1, \dots, p\}, \quad g_i \cap g_j = \emptyset \quad (i \neq j)$$

Example: (Ch 4.3 in PB & v.6)

$\gamma \in \mathbb{R}$

$$X = (X^{(1)}, \dots, X^{(p)})$$

$H_j: X^{(j)} \in \{0, 1, 2, 3\}$ a factor with 4 levels

e.g.: labels A, C, G, T in DNA sequences

for example $p=2$

$$y_i = \mu + \underbrace{\sum_{h=0}^3 \gamma_h I(X_i^{(2)} = h)}_{\text{main effects}} + \sum_{h=0}^3 \delta_h I(X_i^{(2)} = h)$$

$$+ \underbrace{\sum_{h,l=0}^3 \kappa_{h,l} I(X_i^{(2)} = h, X_i^{(1)} = l)}_{\text{interaction terms}} + \epsilon_i$$

overparameterized $\rightarrow \sum_{h=0}^3 \gamma_h = 0, \sum_{h=0}^3 \delta_h = 0$

$$\forall l \sum_{h=0}^3 \kappa_{h,l} = 0; \sum_l \kappa_{h,l} = 0 \quad \forall h$$

$$\left(\rightarrow \sum_{h,l} \kappa_{h,l} = 0 \right)$$

$\underbrace{1 + 1 + 4 + (4 - 1)}_{9 \text{ constraints}}$ constraints

no free parameters: $1 + 4 + 4 + 16 - 9 = 16$

can write this as

$$Y = X\beta + \epsilon$$

Group parameters: $\beta_1 = \mu$ / main eff. of $X^{(2)}$

$$g_1 = \{2, 3, 4\}, g_2 = \{5, 6, 7\}$$

Contr. to main eff. of $X^{(1)}$ $g_3 = \{8, 9, \dots, 16\}$ — interaction effects

aim: we want sparsity in terms of whole groups

either $\hat{\beta}_{g_j} = 0$ or $(\hat{\beta}_{g_j})_r \neq 0 \quad \forall r \in g_j$

this can be achieved by the Group Lasso penalty

(Yuan and Lin, 2006)

$$\hat{\beta}(\lambda) = \arg \min_{\beta} \left(\|Y - X\beta\|_2^2 / n + \lambda \sum_{j=1}^q m_j \|\beta_{g_j}\|_2 \right)$$

$$m_j = \sqrt{|g_j|}$$

Lemma: $g_j = \{j\}$

$$m_j = 1; \quad \|\beta g_j\|_2 = \sqrt{\beta_j^2} = |\beta_j|$$

~ Group basis equals basis for $g_j = \{j\}$

Why $m_j = \sqrt{|g_j|}$?

for basis: $(\beta g_j)_r \neq 0 \quad \forall r \in g_j$

$$\sum_{r \in g_j} |\beta g_j)_r| = \sqrt{|g_j|}$$

with Group basis: $\|\beta g_j\|_2 = \sqrt{|g_j|} \approx m_j \|\beta g_j\|_2 = |g_j|$

why does sparsity hold in terms of

$$\hat{\beta}_{g_j} = 0 \quad \text{or} \quad (\hat{\beta}_{g_j})_r \neq 0 \quad \forall r \in g_j$$

Reason:

sparsity happens at points of non-differentiability
(if derivative is "undefined" "not determined" it goes off
to "set to zero")

$$\text{here: } \|\hat{\beta}_{g_j}\|_2 = \sqrt{\|\beta_{g_j}\|_2^2}$$

$$\frac{\partial}{\partial \beta_{qj}} : \text{ if } \|\beta_{qj}\|_2 \neq 0$$

$$\rightsquigarrow = \frac{1}{2} (\|\beta_{qj}\|_2^2)^{-1/2} \cdot 2 \beta_{qj}$$

$$= \frac{\beta_{qj}}{\|\beta_{qj}\|_2}$$

if $\|\beta_{qj}\|_2 = 0 \rightsquigarrow$ not differentiable

\rightsquigarrow either $\hat{\beta}_{qj} = 0$ or $(\hat{\beta}_{qj})_r \neq 0 \forall r \in q_j$
 because Ridge regression is not sparse

IV. Additive models and many smooth univariate

functions

(Ch. 5 in PB8 ud6)

IV. 1. Model and estimation (Ch. 5.2 in PB8 ud6)

additive model:

$$Y_i = \mu + \sum_{j=1}^p f_j^0(X_i^{(j)}) + \varepsilon_i$$

$E[\varepsilon_i] = 0$, ε_i indep. of X_i (if X_i is random)

identifikation: $\sum_{j=1}^p f_j^0(X_i^{(j)}) = 0 \quad \forall j$

for $p < n$: e.g. Hastie & Tibshirani (1990)

f_j^0 : smooth $\mathbb{R} \rightarrow \mathbb{R}$

regularize w.r.t. smoothness

for $p \gg n$: need to regularize further

approach: use basis expansion

each f_j^0 is approximated with K basis functions

\rightsquigarrow # param: $1 + pK$

still OK if problem is sparse since

$\log(\text{dim})/n = \log(1 + pK)/n$ small