

$L_{\alpha,10} = \text{soft-threshold estimator if } X^T X/n = I$

Proof:

$$\|y - X\beta\|_2^2/n + \lambda \|\beta\|_1$$
$$= \|y\|_2^2/n + \sum_{j=1}^p \underbrace{\left\{ -2 \underbrace{(X^T y)_j/n}_{z_j} \cdot \beta_j + \beta_j^2 + \lambda |\beta_j| \right\}}_{\text{minimize w.r.t. } \beta_j} \quad \forall j$$

decoupled problem!

$$-2z_j \beta_j + \beta_j^2 + 1 |\beta_j| \stackrel{!}{=} \min \beta_j$$

$$\text{if } 0 \leq z_j \leq \frac{1}{2} \Rightarrow \hat{\beta}_j = 0$$

$$\uparrow \text{ because } A |\beta_j| - 2z_j \beta_j \geq 0 \quad \forall \beta_j$$

$$= 0 \quad \beta_j = 0$$

$$\beta_j^2 \geq 0 \quad \forall \beta_j \quad \beta_j = 0$$

$$\text{likewise: } -\frac{1}{2} \leq z_j \leq 0 \Rightarrow \hat{\beta}_j = 0$$

for $z_j > 1/2 \rightsquigarrow \hat{\beta}_j > 0$

and hence objective fct. is differentiable for > 0

$$\frac{\partial}{\partial \beta_j} : -2z_j + 2\beta_j + 1 \stackrel{!}{=} 0$$

$$\rightsquigarrow \hat{\beta}_j = z_j - 1/2$$

likewise for $z_j \leq -1/2$.

□

in general: sub-differential calculus

II. 4. Prediction with the Lasso

goal: estimation of regression function

$$f(x) = E[Y|X=x] = \sum_{j=1}^p \beta_j x_j = \beta^T x = x^T \beta$$

II. 4.1 Practical aspects

we $\hat{f}(x) = \hat{\beta}(A)^T x$

and choose λ via cross-validation

Why prediction?

training data $(X_1, y_1), \dots, (X_n, y_n)$

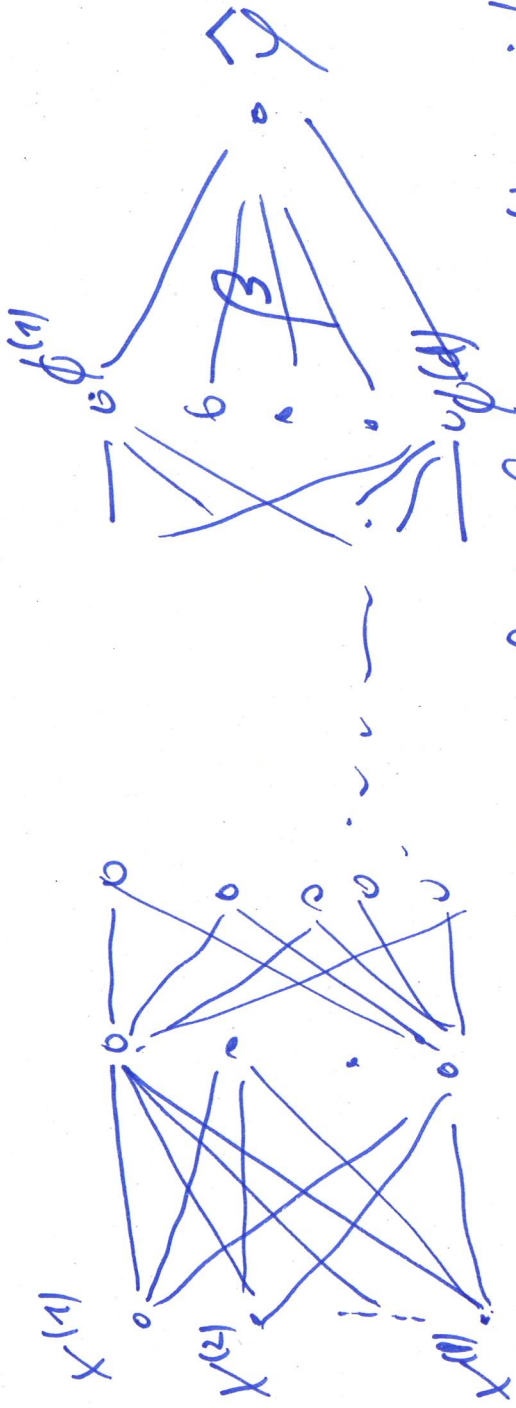
new data point: x_{new}, y_{new}

observe z to be predicted

$$\mathbb{E}[(y_{new} - \hat{\beta}^T x_{new})^2] = \sigma_{new}^2 + \underbrace{\mathbb{E}[(\hat{\beta}^T x_{new} - \beta^T x_{new})^2]}_{\text{indep. of training data}} + \mathbb{E}[\varepsilon_{new}^2]$$

mean squared error for estimating regression function

Comment: deep neural network



last layer with variables/features

$$\phi^{(1)}(x)$$
$$\phi^{(2)}(x)$$

$$d \gg n$$

$$\underbrace{\phi^{(d)}(x)}_{\text{learned from data}}$$

last step:

$$\vec{\beta}(\lambda) = \underset{\beta}{\operatorname{argmin}} \left(\|Y - \Phi\beta\|_2^2 + \lambda \|\beta\|_2^2 \right)$$

or $\lambda \|\beta\|_2^2$

$$\Phi_{n \times d} = \begin{bmatrix} \phi^{(1)}(X_1) & \dots & \phi^{(d)}(X_1) \\ \vdots & \ddots & \vdots \\ \phi^{(1)}(X_n) & \dots & \phi^{(d)}(X_n) \end{bmatrix} \quad d \gg n$$

Choice of λ : cross-validation \hat{I}_{CV}
(typically 10-fold CV; default in glmnet R-package)

good prediction and increased sparsity: "1se rule"
from glmnet

largest λ which is at most 1se away from minimum

II.4.2 Some results from asymptotic theory

goal: what are "some" properties of the Lasso in
terms of mathematical guarantees

consider asymptotics:

$$Y_{n,i} = \sum_{j=1}^{p_n} \beta_{n,j}^{(i)} X_{n,j,i} + \varepsilon_{n,i}$$

$$i = 1, \dots, n$$

$$P = P^n$$

$$n = 1, 2, 3, \dots$$

allow for $\underline{p_n} \geq n$ ($\frac{p_n}{n} \rightarrow \infty$ ($n \rightarrow \infty$))

What happens (to the lasso) as $n \rightarrow \infty$ (and $p_n \rightarrow \infty$)?

notationally: just write $\beta_{p \times 1}^0$

etc...

Consider (here) fixed design X

Theorem

Considers the Lasso $\hat{\beta}(\lambda)$.

Assumptions:

(1) $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$ i.i.d. $\sim \mathcal{N}(0, \sigma^2)$
(not crucial, but simplifies the proof)

$$(2) \frac{1}{n} \sum_{i=1}^n (X_i^{(j)})^2 \equiv 1 \quad \forall j = 1, \dots, p_n$$

Scaled covariates

$$(3) \|\beta^0\|_1 = \sigma \left(\sqrt{\frac{n}{\log(p_n)}} \right) \quad (n \rightarrow \infty)$$

(makes only sense if $\log(p_n) \ll n$) e.g. $p_n = n^\alpha$ works for any $\alpha < \infty$

choose $\lambda = \lambda_n = 4 \hat{\sigma}_n \sqrt{\frac{t_n^2 + 2 \log(p_n)}{n}}$ with

$t_n^2 \rightarrow \infty$, $t_n^2 = O(\log(p_n))$ ($t_n^2 = \log(p_n)$)

$\cdot P[C > \hat{\sigma}_n \geq \sigma] \rightarrow 1$

($\hat{\sigma}_n$ should overestimate σ)

Then: $\|X(\hat{\beta}(\lambda_n) - \beta^0)\|_2^2 / n \xrightarrow{P} 0$ ($n \rightarrow \infty$)
Convergence in probability

with $\|\beta^0\|_1$ is suffic. sparse

$$\lambda = \lambda_n \asymp \sqrt{\frac{\log(p_n)}{n}}$$

$\log(p_n) \ll n$