

Lecture 22

..cv)
Regularization *

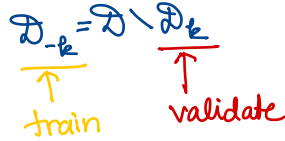
K-fold CV

1. Partition $\mathcal{D} = \mathcal{D}_1 \cup \dots \cup \mathcal{D}_k$ disjoint, $|\mathcal{D}_1| \approx \frac{n}{k}$

Large n
small n : k-fold

2. for $k=1:k$

train $f_{1:m}^{(k)}$ on \mathcal{D}_{-k}
compute $L(f_m^{(k)}; \mathcal{D}_k) = L_m^{(k)}$



$$\bar{L}_m = \frac{1}{k} \sum_k L_m^{(k)}$$

← avg reduces variance

$$m^* = \underset{m}{\operatorname{argmin}} \bar{L}_m$$

→ final predictor = model m^* trained on all \mathcal{D}

"Training set" $|\mathcal{D}_{-k}| = \frac{k-1}{k} n$
validation set $|\mathcal{D}_k| = \frac{1}{k} n$

$k \uparrow$ training time $\sim |\mathcal{D}_{-k}| \cdot k = \frac{k-1}{k} k \cdot n \propto k \cdot n$
needed only when $n \downarrow$

$k=n$ Leave-One-Out CV for $n < 100$

valid



$k=4$

Regularization - a bias for simple models

- weight decay = L2 Regularization
- lasso = L1 Regularization
- data augmentation
- early stopping
- dropout
- AE: denoising auto encoders, sparse autoencoders

Bagging DT \rightarrow Random Forest (RF) = $DT_1: M$ trained independently
 perturb $\left\{ \begin{array}{l} \text{data - e.g. subsample of } \mathcal{D} \text{ for } n'=n \text{ with replacement} \\ \text{algorithm - randomize splits} \end{array} \right.$

$M \approx 100$

TRAIN

PREDICT

$f(x) = \frac{1}{M} \sum_{k=1}^M DT_k(x)$ avg. over forest: Var \downarrow
 single $DT_k \rightarrow$ low bias, high var

TRAIN Model $y = \beta^T x + \epsilon$ $x, \beta \in \mathbb{R}^d$
 $\min_{\beta} L_{LS} + R(\beta) = \frac{1}{n} \sum_{i=1}^n (y_i - \beta^T x_i)^2 + \lambda$
 REGULARIZATION TERM

PREDICTION Want L_{LS} low

$L_2 \leftarrow$ closed form $\hat{\beta} = (X^T X + 2\lambda I)^{-1} X^T y$ **not sparse**
 $L_1 \leftarrow$ **CONVEX, sparse $\hat{\beta}$** , not closed form $\hat{\beta} = (X^T X)^{-1} X^T y$
 0 no regul. ariz. \leftarrow closed form
 $\#\{|\beta_j| \neq 0, j=1:d\} = s \leftarrow$ sparse reg \checkmark
NOT convex

$\lambda \uparrow$ bias \uparrow , var \downarrow
 $\lambda \downarrow$ bias \downarrow , var \uparrow
 $\lambda \geq 0$
 λ chosen by CV

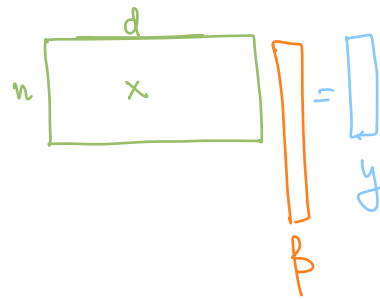
sparse Regression

Assumption: $\beta_{1:s} \neq 0, \beta_{s+1:d} = 0 \iff$ beta sparse
 $s \ll d$
 $y = \beta_1 x_1 + \dots + \beta_s x_s + \text{noise}$ ($x_{s+1:d} = \text{irrelevant}$)
relevant variables

sparse Regression

Assumption: $\beta_{1:s} \neq 0$, $\beta_{s+1:d} = 0 \Leftrightarrow \beta$ is sparse
 $s \ll d$

$y = \beta_1 x_1 + \dots + \beta_s x_s + \text{noise}$
 (relevant variables) $(x_{s+1:d} = \text{irrelevant})$



Convex Optimization for L1 regularization

$$\min_{\beta} \frac{1}{n} \sum_i (y_i - \beta^T x_i)^2 + \lambda \|\beta\|_1$$

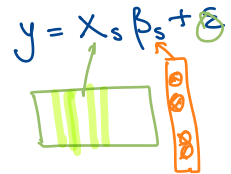
convex function of β

- No local optima
- algorithms \exists
- ——— tractable
- can be analyzed $\Rightarrow \hat{\beta}$ sparse

Theorem $y = X^T \beta^* + \varepsilon$, $\|\varepsilon\|$ small $\in \mathbb{R}^n$
 X "good", $n \sim s \log_2 \frac{d}{s}$
 s sparse

$$\Rightarrow \text{supp } \hat{\beta} \subseteq \text{supp } \beta^* = S$$

→ LS Regression Re-estimate β_s from



$$\text{supp } \beta^* = \{j, \beta_j^* \neq 0\}$$

Why $s \log_2 \frac{d}{s}$?

$P_0 =$ "find S " + estimate s β_j 's
 $\sim s$ data points
 #bits?

$N \in [0: 2^k - 1] \Rightarrow$ need k bits
 $\lceil \log_2 N \rceil$

$S =$ one of $\binom{d}{s}$ sets
 \downarrow
 N

$$\text{need } \log_2 \binom{d}{s} = \log_2 \frac{d(d-1) \dots (d-s+1)}{s!} \approx \log_2 \frac{d^s}{s^s} = s \log_2 \frac{d}{s} \text{ bits}$$