

23. Model Selection and Validation

Tuesday, March 5, 2024 6:36 AM

The issues we want to solve:

- How to set the **hyperparameters** in your algorithm?
(e.g. # rands T of boosting?)
- How to choose which **algorithm** or hypothesis class \mathcal{H} ?
- Go beyond theory (w , unknown constants...) — we want something practical!
- [SS] presents 2 methods, we'll add in classical methods
 - Correct method depends on the data **regime**: lots of data?
limited data?

Generic Setups

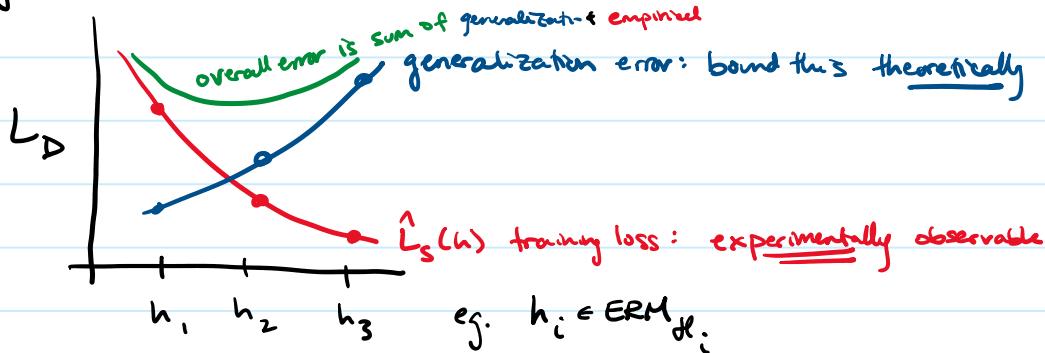
① $\mathcal{H}_1 \subseteq \mathcal{H}_2 \subseteq \mathcal{H}_3 \subseteq \dots$ e.g. $\mathcal{H}_n = \{ \text{polynomials of degree } n \text{ or less} \}$. Choose $n \in \mathbb{N}$

② algo is Tikhonov regularized least squares $\min_w \frac{1}{2} \|X \cdot w - y\|^2 + \alpha \|w\|^2$.
aka Ridge Regression Choose $\alpha \geq 0$.

Technique #1: Model Selection via SRM (= Structured Risk Minimization) [§11.1]

mostly of theoretical interest.

Idea:



Similarly:

error decomp.: $\xrightarrow{\text{complexity}}$

$$L_D(h) = \left(L_D(h) - \min_{h' \in \mathcal{H}} L_D(h') \right) + \left(\min_{h' \in \mathcal{H}} L_D(h') \right)$$

estimation error approx. error or bias

SRM based on a (confusing) weight for obviously this method only works well if theoretical bound captures right behaviour

For theoretical bound, use something like:

- VC dimension
- Minimum Description Length (MDL) (we'll revisit)

- Rademacher Complexity $L_D(h) = \hat{L}_s(h) + R_m(\mathcal{H}) + \sqrt{\log \frac{1}{\delta}} \underbrace{\text{use in SRM}}$

23a. Model Selection and Validation

Tuesday, March 5, 2024 1:35 PM

Technique #2: Validation set i.e. train/validation/test data split

Method of choice (simple, rigorous) if data is plentiful

Let D be any distribution, $S = (z_i)_{i=1}^m$ training set, $h = h_S$,

$V = (z_i)_{i=1}^{m_V}$ a validation or test set, $V \cap S = \emptyset$, all iid, then

$\forall h$ (either deterministic or independent of V)

$$|\hat{L}_V(h) - L_D(h)| \leq \sqrt{\frac{\log(2/\delta)}{2m_V}} \text{ w.p. } \geq 1-\delta$$

for loss function
w/ range in $[0, 1]$

proof: Hoeffding

$$\text{aka } m_V = \frac{1}{2\varepsilon^2} \log(2/\delta)$$

How to use? e.g. model selection:

if we choose among $\{h_1, \dots, h_K\}$ then we have

$$(\text{Simultaneously}) \forall h \in \{h_1, \dots, h_K\} \quad |\hat{L}_V(h) - L_D(h)| \leq \sqrt{\frac{\log(2K/\delta)}{2m_V}}$$

via union bound.

i.e. Select $h \in \arg\min_{i \in [K]} \hat{L}_{V_1}(h_i)$. V_1 = validation

K too large is a bit like p-hacking... hence always b.
test err also.

then... re-estimate L_D using \hat{L}_{V_2} V_2 = testing data

to get unbiased estimate (otherwise we're returning an order statistic)

23b. Model Selection and Validation

Tuesday, March 5, 2024 1:56 PM

Techniques not in [SS]

Technique #3 Mallows's C_p / UPRE

cf. Vogel's '02 "Comp. Methods for Inverse Problems"

and

§ 7.5 in

Hastie, Tibshirani, Friedman's
"Elements of Stat-Learning"

Model: $y = X \cdot w_{\text{true}} + \eta$ ← noise (i.e. agnostic case)

$m \geq d$

We'll estimate w_{true} via \hat{w}

For a given \hat{w} , the predictive error is $p = X \cdot (\hat{w} - w_{\text{true}})$
and the predictive risk is $\frac{1}{m} \|p\|^2$

Restrict our attention to estimates \hat{w} that are linear functions of the data y ,
and parameterized by α , i.e. $\hat{w} = w_\alpha := R_\alpha^{-1} y$ matrix
(\neq Rademacher complexity!)

Main ex.: ridge regression $w_\alpha := \arg \min \frac{1}{2} \|X \cdot w - y\|^2 + \alpha \|w\|^2$ variant: $\|w - w_0\|^2$
 $= (X^T X + \alpha I)^{-1} X^T y$ always unique sol'n if $\alpha > 0$

Question: What value of α to use?

α too small and we might overfit ($m \gg d$) or no unique sol'n ($m < d$)

α too large and we bias \hat{w} too much toward 0. (more bias, less variance)

UPRE looks at $E \frac{1}{m} \|p_\alpha\|^2$,

$$\begin{aligned} p_\alpha &:= X \cdot (w_\alpha - w_{\text{true}}) \\ &= r_\alpha + \eta \\ &\quad \begin{matrix} \nearrow \text{observable} \\ \uparrow \text{unobservable} \end{matrix} \end{aligned}$$

$$\begin{aligned} r_\alpha &= X \cdot w_\alpha - y \quad // \text{residual} \\ &= X \cdot w_\alpha - (X \cdot w_{\text{true}} + \eta) \end{aligned}$$

$$\text{so } E \frac{1}{m} \|p_\alpha\|^2 = \frac{1}{m} E \|r_\alpha\|^2 + \frac{2}{m} E \langle r_\alpha, \eta \rangle + \frac{1}{m} E \|\eta\|^2 \quad \begin{matrix} \underbrace{}_A \\ = \sigma^2 \end{matrix}$$

Assume $E[\eta] = 0$

$\text{Var}(\eta) = \sigma^2 I_{m \times m}$
(“white” noise)

$$\begin{aligned} r_\alpha &= X w_\alpha - y \\ &= (X R_\alpha - I) y \quad \text{let } P := X R_\alpha \quad (\text{uppercase, } \neq p) \\ &= (X R_\alpha - I)(X w_{\text{true}} + \eta) = (P - I)\eta + (P - I) \cdot X w_{\text{true}} \end{aligned}$$

$$\begin{aligned} A &= \frac{2}{m} E \langle (P - I)\eta, \eta \rangle + \frac{2}{m} E \langle (P - I) \cdot X w_{\text{true}}, \eta \rangle \quad = 0 \text{ since } E[\eta] = 0 \\ &= \frac{2}{m} \left(E \eta^T P \eta - E \eta^T \eta \right) \quad = \frac{2}{m} E \eta^T P \eta - 2 \\ &\quad \underbrace{}_{m \cdot \sigma^2} \end{aligned}$$

23c. Model Selection and Validation

Tuesday, March 5, 2024 7:25 PM

what is $\mathbb{E} \eta^T P \eta$? (assuming $\mathbb{E}[\eta] = 0$, $\mathbb{E}[\eta \eta^T] = \text{Var}[\eta] = \sigma^2 I$)
 uncorrelated

$$\begin{aligned}\mathbb{E} \eta^T P \eta &= \mathbb{E} \operatorname{tr}(\eta^T P \eta) \stackrel{\text{cyclic property}}{=} \mathbb{E} \operatorname{tr}(P \cdot \eta \eta^T) \stackrel{\text{linearity}}{=} \operatorname{tr}(P \mathbb{E} \eta \eta^T) \\ &\stackrel{\substack{\text{trace of a } 1 \times 1 \text{ matrix} \\ \text{2 different ways to see it}}}{=} \mathbb{E} \left[\sum_i \eta_i^2 P_{ii} + \sum_{i \neq j} \eta_i \eta_j P_{ij} \right] \stackrel{=0}{\cancel{=}} \end{aligned}$$

Cyclic property of trace
 $\operatorname{tr}(ABC) = \operatorname{tr}(BCA)$
 $= \operatorname{tr}(CAB)$

So...

$$\mathbb{E} \frac{1}{m} \|P_\alpha\|^2 = \underbrace{\frac{1}{m} \mathbb{E} \|r_\alpha\|^2}_{\text{a realization is an unbiased estimate.}} + \sigma^2 \left(\frac{2}{m} \operatorname{tr}(P_\alpha) - 1 \right)$$

r_α is observable

$$\text{so the } C_p \text{ statistic is } U(\alpha) = \underbrace{\frac{1}{m} \|r_\alpha\|^2}_{\text{empirical risk}} + \underbrace{\frac{2\sigma^2}{m} \operatorname{tr}(P_\alpha)}_{\text{penalizes complexity}}$$

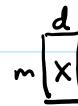
and UPRE chooses the α that minimizes this statistic.

Hastie calls $\frac{1}{m} P_\alpha$ the "effective number of parameters" d

Ex: ridge regression/Tikhonov $P_\alpha = X R_\alpha = X(X^T X + \alpha I)^{-1} X^T$

$$\begin{aligned}\operatorname{tr}(P_\alpha) &= \operatorname{tr}(X(X^T X + \alpha I)^{-1} X^T) \\ &= \operatorname{tr}(\underbrace{(X^T X + \alpha I)^{-1}(X^T X + \alpha I)}_{=I} - \alpha I) \\ &= \operatorname{tr}(I_{d \times d}) - \alpha \operatorname{tr}((X^T X + \alpha I)^{-1}) \\ &= d - \sum_{i=1}^d \frac{\alpha}{\alpha + \lambda_i(X^T X)} \end{aligned}$$

$\lambda_i(X^T X)$ = eigenvalues of $X^T X$



As $\alpha \rightarrow \infty$, $\operatorname{tr}(P_\alpha) = 0$ but empirical risk will be large
 $\alpha \rightarrow 0$, $\operatorname{tr}(P_\alpha) = d$ so "effective dimension" = "actual dimension"

Stein's Unbiased Risk Estimate "SURE" (Stein '81)

generalizes this to nonlinear estimates, $\omega_\alpha = R_\alpha(y)$

and replaces trace w/ divergence of a Jacobian.

$$\text{i.e. } \operatorname{tr}(P_\alpha) \rightarrow \operatorname{div}(R_\alpha(y)) := \sum_{i=1}^n \frac{\partial r_i}{\partial y_i} \cdots ?$$