

Lecture Notes V – Model selection

Marina Meilă
mmp@stat.washington.edu

Department of Statistics
University of Washington

November, 2021

Cross-validation

AIC and BIC

Structural risk minimization and VC dimension

Reading HTF Ch.: Ch. 7, Murphy Ch.: BIC, AIC 8.4.2 (pp 255), SRM 6.5 (pp 204)

Crossvalidation

K-fold Crossvalidation

AIC and BIC

Hold for

- ▶ parametric \mathcal{F}
- ▶ log-likelihood loss $L(y, f(x)) = -\ln P(y|x, f)$
Note that $-n\hat{L}(f) = \ln P(y^{1:n}|x^{1:n}, f)$ data log likelihood
- ▶ $\hat{f} \in \mathcal{F}$ estimated by Maximum Likelihood
- ▶ (for BIC: $\frac{\partial^2 L}{\partial \text{parameters}}$ non-singular at \hat{f})

Akaike's Information Criterion (AIC)

$$AIC(\hat{f}) = -n\hat{L}(\hat{f}) - d, \quad (1)$$

where $d = \#\text{parameters}(f)$, and $n =$ the size of \mathcal{D} .

The Bayesian Information Criterion (BIC)

$$BIC(\hat{f}) = -n\hat{L}(\hat{f}) - \frac{d}{2} \ln n, \quad (2)$$

with $d = \#\text{parameters}(f)$

AIC and BIC

VC dimension

\mathcal{F} shatters $\mathcal{D}_h = \{x^1, \dots, x^h\}$

iff, for every possible labeling $y^{1:h} \in \{\pm 1\}$ of \mathcal{D}_h , there is a function $f \in \mathcal{F}$ that achieves that labeling, i.e. $\text{sgn}f(x^i) = y^i$ for all $i = 1 : m$.

VC dimension

\mathcal{F} shatters $\mathcal{D}_h = \{x^1, \dots, x^h\}$

iff, for every possible labeling $y^{1:h} \in \{\pm 1\}$ of \mathcal{D}_h , there is a function $f \in \mathcal{F}$ that achieves that labeling, i.e. $\text{sgn}f(x^i) = y^i$ for all $i = 1 : m$.

VC dimension of \mathcal{F}

A model class \mathcal{F} over \mathbb{R}^d has VC dimension h iff h is the maximum positive integer so that there exists a set of h points in \mathbb{R}^d that is shattered by \mathcal{F} .

Structural risk minimization

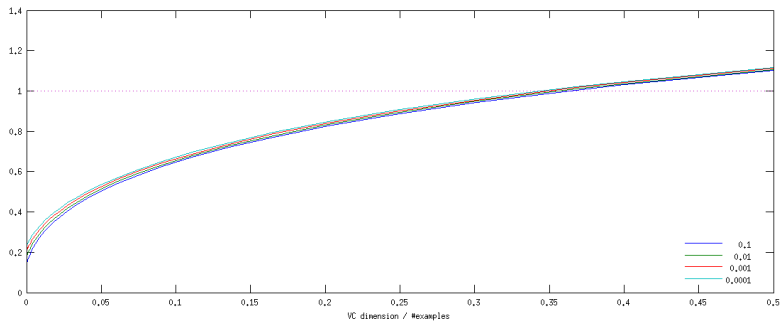
Theorem

Let \mathcal{F} be a model class of VC-dimension h and f a classifier in \mathcal{F} . Then, with probability w.p. $> 1 - \delta$ over training sets

$$L_{01}(f) \leq \hat{L}_{01}(f) + \sqrt{\frac{h[1 + \log(2n/h)] + \log(4/\delta)}{n}}. \quad (3)$$

Structural risk minimization

$$\sqrt{\frac{h[1 + \log(2n/h)] + \log(4/\delta)}{n}}$$



Structural risk minimization

Theorem

Let \mathcal{F} be a model class of VC-dimension h , with $f(x) \in [-1, 1]$ for all x and for all $f \in \mathcal{F}$. Let $\delta > 0$ and $\zeta \in (0, 1)$. Denote $\mathcal{D} = \{(x^i, y^i), i = 1 : n\}$ the current training set. Then, with probability w.p. $> 1 - \delta$ over training sets

$$L_{01}(f) \leq \hat{L}_{01,\zeta}(f) + \tilde{O} \left(\sqrt{\frac{h}{n\zeta^2}} \right) \quad (4)$$

for any $f \in \mathcal{F}$.

The test set method of bounding the classification error

[after <https://www.jmlr.org/papers/volume6/langford05a/langford05a.pdf>]

Given a classifier f and a data set $\mathcal{D}^{\text{test}}$ of size n .

$\hat{L}_{01}(f) \sim \text{Binomial}(L_{01}(f), n)$ denote $\bar{b}(m, L_{01}, \delta) = \max\{L_{01} \mid \Pr[m \mid L_{01}, n] \geq \delta\}$

$$L_{01}(f) \leq \hat{L}_{01}(f) + \sqrt{\frac{\ln 1/\delta}{2n}} \quad \text{w.p. } 1 - \delta \quad (5)$$

$$L_{01}(f) \leq \hat{L}_{01}(f) + \sqrt{\frac{2\hat{L}_{01} \ln 1/\delta}{n}} + \frac{2 \ln(1/\delta)}{n} \quad \text{w.p. } 1 - \delta \quad (6)$$

$$L_{01}(f) \leq \frac{\ln 1/\delta}{n} \quad \text{w.p. } 1 - \delta \quad \text{when } \hat{L}_{01}(f) = 0 \quad (7)$$

$$|L_{01}(f) - \hat{L}_{01}(f)| \leq \sqrt{\frac{\ln 1/\delta}{2n}} \quad \text{w.p. } 1 - \delta \quad (8)$$

The test set method of bounding the classification error

$$|L_{01}(f) - \hat{L}_{01}(f)| \leq \sqrt{\frac{\ln 1/\delta}{2n}} \quad \text{w.p. } 1 - \delta$$