

1 Bounding Error on Mean

Given a set of examples $x_i = [0, 1]$ $i = \{1 \dots m\}$ drawn from a distribution D we would like to calculate the observed mean \hat{p}

$$\hat{p} = \frac{1}{m} \sum_{i=1}^m x_i \tag{1}$$

and compare that to the real mean of D , call it p . It is useful to define the quantity

$$q \equiv p + \epsilon. \tag{2}$$

We could weakly bound the probability of \hat{p} being greater than q using Markov's inequality:

$$\Pr[x \geq kEX] \leq 1/k \tag{3}$$

$$\Pr[\hat{p} \geq s] \leq p/s \tag{4}$$

$$\Pr[\hat{p} \geq q] \leq p/q < 1. \tag{5}$$

However, we can do better than that, in fact

THEOREM

$$\Pr[\hat{p} \geq q] \leq \exp(-\text{RE}(q|p)m) \tag{6}$$

where

$$\text{RE}(q|p) = p \ln \left(\frac{p}{q} \right) + (1-p) \ln \left(\frac{1-p}{1-q} \right) \tag{7}$$

PROOF

$$\Pr[\hat{p} \geq q] = \Pr[e^{\lambda \hat{p} m} \geq e^{\lambda q m}] \leq e^{-\lambda q m} E[e^{\lambda \hat{p} m}] \tag{8}$$

$$= e^{-\lambda q m} E[e^{\lambda \sum x_i}] = e^{-\lambda q m} E \left[\prod_{i=1}^m e^{\lambda x_i} \right] \tag{9}$$

now because each x_i is an independant measurement

$$e^{-\lambda q m} \prod_{i=1}^m E[e^{\lambda x_i}]. \tag{10}$$

Now note that

$$e^{\lambda x} \leq 1 - x + e^{\lambda} x \quad \forall x \in [0, 1] \tag{11}$$

which we can use to show that

$$e^{-\lambda q m} \prod_{i=1}^m E[e^{\lambda x_i}] \leq e^{-\lambda q m} \prod_{i=1}^m E[1 - x_i + e^{\lambda} x_i] \tag{12}$$

$$= e^{-\lambda q m} \prod_{i=1}^m (1 - p + e^{\lambda} p) = e^{-\lambda q m} (1 - p + e^{\lambda} p)^m \tag{13}$$

$$= \left(e^{-\lambda q} (1 - p + e^{\lambda} p) \right)^m. \tag{14}$$

Now if we minimize this probability with respect to lambda you get

$$\lambda_{min} = \ln \left(\frac{q(1-p)}{(1-q)p} \right). \quad (15)$$

Plugging that back into the probability you get the desired bound

$$\Pr[\hat{p} \geq q] \leq \exp(-\text{RE}(q|p)m) \quad (16)$$

Note that by defining $y_i = 1 - x_i$ we can prove the symmetric result that

$$\Pr[\hat{p} \geq p - \epsilon] \leq \exp(-\text{RE}(p - \epsilon|p)m) \quad (17)$$

This theorem can also be used to prove corollaries

$$\Pr[\hat{p} \geq p + \alpha p] \leq e^{-m\alpha^2/3} \quad (18)$$

$$\Pr[\hat{p} \leq p - \alpha p] \leq e^{-m\alpha^2/2} \quad (19)$$

This is done by plugging in $\epsilon = \alpha p$ and then bounding $\text{RE}(q|p)$.

Thinking back to the double-sample proof: $m(h)$ was the number of mistakes on S' . Thus $\hat{p} = m(h)/m$ and $p = \epsilon$ and so

$$\Pr[m(h) < m\epsilon/2] = \Pr[\hat{p} < p - p/2] \leq e^{-mp/8}. \quad (20)$$

If $m \geq 8/\epsilon$ then

$$\Pr[m(h) < m\epsilon/2] \leq e^{-1} < 1/2. \quad (21)$$

2 McDiarmid's Inequality

Let

$$f(x_1, \dots, x_m) \quad (22)$$

be any function such that for all $x_1, \dots, x_m; x'_i$,

$$|f(x_1, \dots, x_i, \dots, x_m) - f(x_1, \dots, x'_i, \dots, x_m)| \leq c_i. \quad (23)$$

In other words, changing x_i can never change f by more than c_i . Let $X_1 \dots X_m$ be independent but not necessarily identically distributed.

THEOREM

$$\Pr[f(x_1 \dots x_m) \geq \mathbb{E}[f(x_1 \dots x_m)] + \epsilon] \leq \exp\left(\frac{-2\epsilon^2}{\sum c_i^2}\right) \quad (24)$$

3 Hoeffding

Hoeffding's inequality

$$\Pr[\hat{p} \geq p + \epsilon] \leq e^{-2\epsilon^2 m} \quad (25)$$

is a special case of McDiarmid. Let

$$f(x_1 \dots x_m) = \frac{1}{m} \sum_{i=1}^m x_i. \quad (26)$$

$\mathbb{E}[f] = p$, $c_i = 1/m$ (for $0 \leq x \leq 1$). So making use of McDiarmid's Inequality,

$$\Pr[\hat{p} \geq p + \epsilon] \leq \exp\left(\frac{-2\epsilon^2}{\sum 1/m^2}\right) = \exp(-2\epsilon^2 m) \quad (27)$$

4 General Strategy

Let

$$err(h) = Pr_{x,y \sim D}[y \neq h(x)] \quad (28)$$

$$e\hat{r}r(h) = \frac{1}{m} |\{i : y_i \neq h(x_i)\}| \quad (29)$$

Minimize $e\hat{r}r(h)$. Show close to true error, thus effectively minimizing the true error.

5 Finite Example

THEOREM Assume $|H|$ finite. Given m random examples, with probability $1 - \delta$

$$\forall h \in H : |err(h) - e\hat{r}r(h)| < \epsilon \quad (30)$$

$$\text{if } m \geq O\left(\frac{\ln |H| + \ln(1/\delta)}{\epsilon^2}\right) \quad (31)$$

PROOF

Fix an h .

$$X_i = \{1 \text{ if } h(x_i) \neq y_i, 0 \text{ else}\}. \quad (32)$$

$$E[X_i] = err(h) \quad (33)$$

$$\frac{1}{m} \sum_{i=1}^m X_i = e\hat{r}r(h) \quad (34)$$

$$Pr[\hat{p} \geq p + \epsilon] \leq e^{-2\epsilon^2 m} \quad (35)$$

$$Pr[e\hat{r}r(h) \geq err(h) + \epsilon] \leq e^{-2\epsilon^2 m} \quad (36)$$

$$Pr[|e\hat{r}r(h) - err(h)| \geq \epsilon] \leq 2e^{-2\epsilon^2 m} \quad (37)$$

Now using union bound we can bound the probability for all h .

$$Pr[\exists h \in H : |err(h) - e\hat{r}r(h)| \geq \epsilon] \leq 2|H|e^{-2\epsilon^2 m} \quad (38)$$

So by increasing m we can bound the probability. So in order to make the probability less than δ we need

$$m \geq \frac{\ln(2|H|) + \ln(1/\delta)}{\epsilon^2} \quad (39)$$

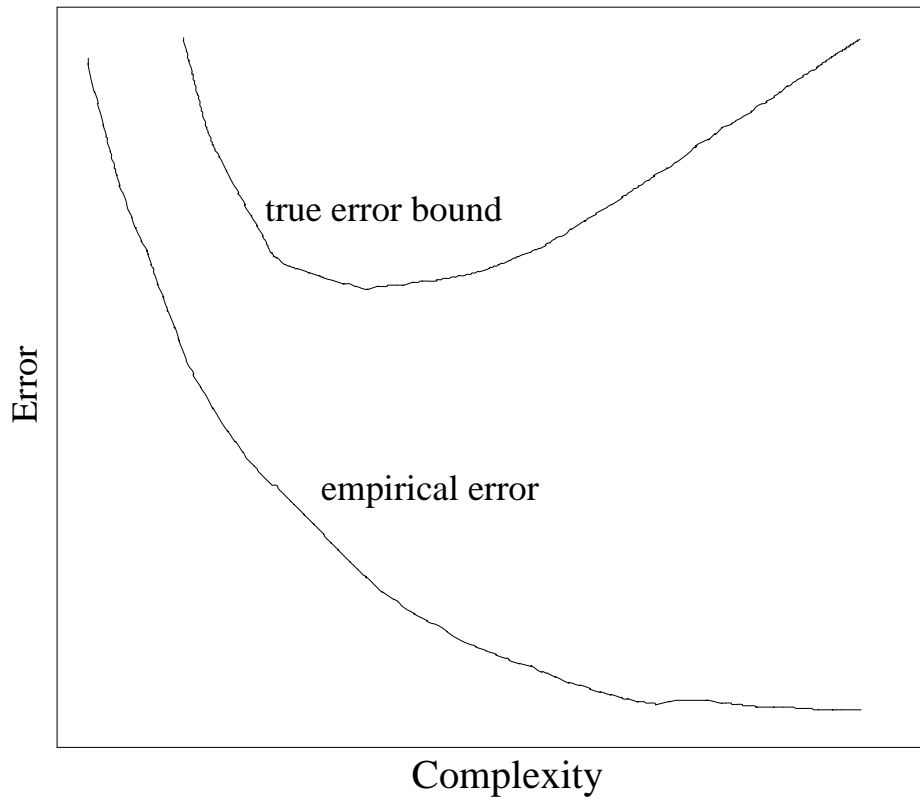
Note that this bound is weaker than before, it requires more examples for smaller errors (grows like $1/\epsilon^2$ vs $1/\epsilon$).

6 Overfitting

$$err(h) \leq e\hat{r}r(h) + O\left(\sqrt{\frac{\ln(|H|) + \ln(1/\delta)}{m}}\right) \quad (40)$$

$$err(h) \leq e\hat{r}r(h) + O\left(\sqrt{\frac{|h| + \ln(1/\delta)}{m}}\right) \quad (41)$$

As h gets more complex, the training error $e\hat{r}(h)$ tends to go down while the complexity $|h|$ increases, so the true error $err(h)$ may at first go down but then increase, as in the figure. This is called overfitting.



Solutions: (1) Cross Validation, which means, hold out some of the training data and use it to determine when to stop training. (2) Treat bound as real and optimize in $|h|$. This is called structural risk minimization. (3) Build algorithms that are resistant to overfitting.