

## Theory Questions (hw2)

(1)

$$\Delta(y, \hat{y}) = \begin{cases} 0 & \hat{y} = y \\ \lambda & \hat{y} = L+1 \\ 1 & \text{otherwise} \end{cases}$$

(a)

$$L(h) = \sum_{x,y} P(Y=y \wedge X=x) \Delta(y, h(x)) =$$

for a given  $\hat{x} = \sum_y P(Y=y | X=\hat{x}) P(X=\hat{x}) \Delta(y, h(\hat{x}))$

$$= P(X=\hat{x}) \cdot \sum_i P(Y=i | X=\hat{x}) \Delta(i, h(\hat{x}))$$

$$= P(X=\hat{x}) \cdot \begin{cases} \lambda \cdot \sum_i P(Y=i | X=\hat{x}) & h(\hat{x}) = L+1 \\ \sum_{j \neq i} P(Y=j | X=\hat{x}) & h(\hat{x}) = i \end{cases}$$

$$\forall i \in [L] \quad \text{s.t. } i = \arg \max_{j \in [L]} P(Y=j | X=\hat{x})$$

$$\text{and } P(Y=i | X=\hat{x}) \geq 1 - \alpha$$

$$\sum_{j \neq i} P(Y=j | X=\hat{x}) = 1 - P(Y=i | X=\hat{x}) \leq 1 - (1-\alpha) = \alpha$$

$$\lambda \sum_i P(Y=i | X=\hat{x}) = \lambda$$

$$\Rightarrow \sum_{j \neq i} P(Y=j | X=\hat{x}) \leq \lambda \cdot \sum_i P(Y=i | X=\hat{x})$$

Also, for  $i = \arg \max_{j \in [L]} P(Y=j | X=\hat{x})$

$$\sum_{j \neq i} P(Y=j | X=\hat{x}) = \sum_j P(Y=j | X=\hat{x}) - P(Y=i | X=\hat{x})$$

is minimal

if  $P(Y=i | X=\hat{x}) < 1-\lambda$

$$\Rightarrow \sum_{j \neq i} P(Y=j | X=\hat{x}) = 1 - P(Y=i | X=\hat{x}) > 1 - (1-\lambda) > \lambda //$$

$$\Rightarrow \lambda = \lambda \sum_j P(Y=j | X=\hat{x}) < \sum_{j \neq i} P(Y=j | X=\hat{x})$$

$\Rightarrow$  we Always Reject! //

③ we note that for  $\lambda=0$   $L(h)$  is minimal only for  $h(\hat{x})=2+1 \forall \hat{x}$  and for  $\lambda=1$ :

$$\sum_{j \neq i} P(Y=j | X=\hat{x}) \leq \lambda \sum_j P(Y=j | X=\hat{x})$$

$\Rightarrow$  we never Reject!

$\lambda$  is a Parameter that controls the probability to reject.

when it increases from 0 to 1 so does the chance that  $h$  will reject.

(2) (a) a good algorithm for ERM shall be:

given a sample  $S$  of pairs  $(x_i, y_i)$ .

- if there exist a pair  $(x_j, 1)$

$\underline{\text{ERM}}(S) = h_{x_j}$ . This is because in the realizable case, given a pair

like this, we know the true

distribution  $P$ . for this case  $\ell_P(\text{ERM}(S))=0$ .

- if  $y_i=0 \quad \forall (x_i, y_i) \in S$ , then  $\text{ERM}(S) = h^-$  meaning we say that  $\forall x \in \mathcal{X} \quad h^-(x)=0$ .

this is of course consistent with

$$h(x) = \arg \max_{h \in \mathcal{H}} P(Y=1 | X=x);$$

say we choose  $h_{\bar{x}}$  for a random  $\bar{x}$  instead of  $h_{\bar{x} \notin S}$

$$L(h) = \sum_{\substack{x \in S \\ x \neq \bar{x}}} P(Y=1 | X=x) P(x) + P(Y=0 | X=\bar{x}) P(\bar{x}) \quad h=h_{\bar{x}}$$

$$\sum_{\bar{x}} P(Y=1 | X=\bar{x}) P(\bar{x}) \quad h=h^-$$

$\Rightarrow$  choosing  $h=h^-$  if  $P(Y=0 | X=\bar{x}) \geq P(Y=1 | X=\bar{x})$

but for  $|X| > |S|$  choosing  $\bar{x}$  for which  $y=0$ ,

$$P(X=\bar{x}, y=1) = \frac{1}{|\mathcal{X}| - |S|}; \quad P(X=\bar{x}, y=0) \geq \frac{|X|-|S|-1}{|\mathcal{X}| - |S|}$$

$$\Rightarrow P(Y=1 | X=\bar{x}) = \frac{1}{P(\bar{x})(|\mathcal{X}| - |S|)} \leq \frac{|\mathcal{X}| - |S| - 1}{P(\bar{x})(|\mathcal{X}| - |S|)} \leq P(Y=0 | X=\bar{x})$$

(2) (b) if we had some  $(x_i, 1)$  in our sample - Then we know the distribution and  $e_p(h) = 0$ .

otherwise, according to our ERM we choose  $h = h^*$ ,  $h^*(x) = 0 \forall x$ .

Assume there exist  $x^* \notin S$  for which  $y(x^*) = 1$  and  $P(x^*) \geq \varepsilon$ .

(if  $P(x^*) < \varepsilon$  then  $e_p(h^*) = P(Y=1|X=x^*)P(x^*) < \varepsilon$ )  
and we are done.

Then:

$$= 0 *$$

$$\begin{aligned} P(e_p(h) > \varepsilon) &= \underbrace{P(e_p(h) > \varepsilon | x^* \in S)}_{\leq 1} \cdot P(x^* \in S) \\ &\quad + \underbrace{P(e_p(h) > \varepsilon | x^* \notin S)}_{-} P(x^* \notin S) \end{aligned}$$

$$\Rightarrow P(e_p(h) > \varepsilon) \leq P(x^* \notin S) = P(\bigcup_{x_i \in S} x_i \neq x^*)$$

$$= \prod_i P(x_i \neq x^*) = \prod_i (1 - P(x^*)) \leq \prod_i (1 - \varepsilon)$$

$$= (1 - \varepsilon)^n \leq e^{-\varepsilon n} \leq \sigma$$

$$\Rightarrow n \geq \frac{1}{\varepsilon} \ln \frac{1}{\sigma}$$

PAC learnable!

(3) PAC learnable in expectation

def: for Algorithm A and  $N(a): (0, 1) \rightarrow N$   
s.t.  $\forall a \in (0, 1)$ ,  $\forall \delta > 0$  and sample S  
so that  $|S| \geq N(a)$   
 $\Rightarrow E[\ell_p(A(S))] \leq \alpha$

$H$  is PAC learnable  $\iff H$  is PAC learnable with expectation

Proof

$H$  is PAC learnable  $\Rightarrow$

$\forall n \geq N(\epsilon, \delta)$  and samp. S. from P

Denote  $\ell_p(A(S)) = \hat{\ell}_p$

$$E[\hat{\ell}_p] = E[\hat{\ell}_p | \hat{\ell}_p \geq \epsilon] P(\hat{\ell}_p \geq \epsilon) + E[\hat{\ell}_p | \hat{\ell}_p < \epsilon] P(\hat{\ell}_p < \epsilon)$$

total expectation  $\leq 1 \leq \delta \leq \epsilon \leq 1$

$H$  is PAC learnable

$$\Rightarrow E(\hat{\ell}_p) \leq \delta + \epsilon$$

Specifically there exist  $\epsilon = \delta = \frac{\alpha}{2} \cdot a \in (0, 1)$

$\Rightarrow E(\hat{\ell}_p) \leq \alpha \Rightarrow H$  is PAC learnable in expectation

Conversely, if  $H$  is PAC learnable in expectation

$$\text{then: } P(\hat{\ell}_p \geq \epsilon) \leq \frac{E(\hat{\ell}_p)}{\epsilon} \leq \frac{\alpha}{\epsilon} = \frac{\alpha}{\epsilon} = \delta$$

markov  $a = \delta \cdot \epsilon$   
 $(0 \leq \hat{\ell}_p \leq 1)$   $\delta, \epsilon \in (0, 1)$

$\Rightarrow H$  is PAC learnable.

(3)

$$H_{1k} = \{h_I \mid I = \{(l_1, u_1) \dots (l_k, u_k)\} \text{ (disjoint intervals)}\}$$

First, we prove that there exist a sample of size  $2k$  that can be shattered meaning  $\sqrt{C} \dim(H_{1k}) \geq 2k$

Proof: for a set of unique  $x_i$  (sorted by the index  $i$ ) and corresponding  $y_i \in \{0, 1\}$  we can choose  $h(x_i)$  like this:

- for a sequence of  $j = m, m+1 \dots m+l$  such that  $y_j = 1 \forall j$ , we define an interval  $[x_m, x_{m+l}]$  in which  $h(x_j) = 1$
- we take  $l, m$  such that  $y_{m-1} = y_{m+l+1} = 0$  meaning, we consider the longest possible sequence of  $j$ 's and not subsequences of it.

We are expected to use at most  $k$  intervals this is for the case of  $y_1 = 0, y_i = |y_{i-1} - 1|$  meaning all labels are different from the labels of the two nearest neighbors and we have  $k$  sequences of  $j$ 's, each with length 1.

Also, if we use less than  $k$  intervals to cover all sequences of  $j$ 's, we can always create two different intervals from one.  $[x_m, x_{m+l}] \rightarrow [x_m, x_p] \cup [x_p, x_{m+l}]$  until we have  $k$  intervals to describe  $h$  such that  $h \in H_{1k}$ .

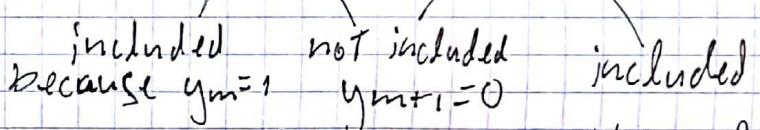
Now we shall prove that  $V\text{dim}(H_k) \leq 2k$

Consider  $2k+1$  points,  $x_i$  and the set of corresponding labels  $y_i = i \% 2$ .

( $1, 0, 1, \dots, 1$ ). ( $x_i$  does not have to be unique points)

for this set of labels we have to define  $k+1$  disjoint intervals, each include exactly 1 point  $x_i$  for which  $y_i = 1$ . but we only have  $k$  intervals to use, so if two adjacent

intervals are stretched  $[x_m, x_{m+1}], [x_{m+1}, x_{m+2}]$



$\rightarrow [x_m, x_{m+2}]$  we now have  $k$  intervals

but for one of the points,  $x_{m+1}$  for which  $y_{m+1} = 0$  we get  $l(x_{m+1}) = 1$

$\Rightarrow$  the given labels sequence could not be achieved using any hypothesis

$\Rightarrow V\text{dim} \leq 2m$

$\Rightarrow V\text{dim} = 2k$

using first part as well

⑤

$$\sum_{h \in H} w(h) \leq 1 \quad H = \cup_{i \in N} \{h_i\}$$

$$w(h) \in [0, 1] , \quad \sum_{h \in H} w(h) \leq 1$$

we start with the fact that

$$P(|e_p(h) - e_s(h)| > \varepsilon) \leq 2|H|e^{-2n\varepsilon^2} \leq \delta$$

$$\Rightarrow \varepsilon \leq \sqrt{\frac{1}{2n} \ln \frac{2|H|}{\delta}}$$

which means that

$$P(e_p > e_s + \sqrt{\frac{1}{2n} \ln \frac{2|H|}{\delta}}) \leq \delta$$

also, applying uniform convergence for each  $h_i$  separately with

$$\delta_i = \delta \cdot w(h_i)$$

$$P(e_p(h_i) > e_s(h_i) + \sqrt{\frac{1}{2n} \ln \frac{2|H|}{w(h_i)\delta}}) \leq \delta_i$$

define  $A_{h_i}$  as the event

$$e_p(h_i) > e_s(h_i) + \sqrt{\frac{1}{2n} \ln \frac{2|H|}{w(h_i)\delta}}, \text{ then}$$

$$P(\cup A_i) \leq \sum_i P(A_i) \leq \sum_i \delta \cdot w(h_i) \leq \delta$$

meaning, with confidence  $1-\delta$ ,  $\forall h \in H$

$$e_p(h) \leq e_s(h) + \sqrt{\frac{1}{2n} \ln \frac{2|H|}{\delta w(h)}}$$

## Programming Assignment

(b) in this case, we know the true prob.  
and we would like to choose that is  
consistent with  $\hat{h} = \arg \max_{y \in \{0,1\}} P(Y=y | X=x)$

because we have intervals  $I$  for which  
the prob for  $y=1$  is  $> 0.5$  for  $x \in I$   
and  $< 0.5$  for  $x \notin I$  we can choose

$$h^*(x) = \begin{cases} 1 & x \in I \\ 0 & x \notin I \end{cases} \quad I = [0, 0.2] \cup [0.5, 0.6] \cup [0.8, 1]$$

the true error for such  $h$  is

$$\begin{aligned} \epsilon_p(h^*) &= P(Y=0 \wedge h^*(x)=1) + P(Y=1 \wedge h^*(x)=0) \\ &= P(Y=0 | h^*(x)=1) P(h^*(x)=1) + P(Y=1 | h^*(x)=0) P(h^*(x)=0) \\ &= \left( \begin{array}{l} x \sim U(0, 1) \text{ and } P(h(x)=1) = P(x \in I) \\ P(h(x)=0) = P(x \notin I) \end{array} \right) \end{aligned}$$

so putting into the eq the numbers from  $P$ ,

$$= 0.2 \cdot \frac{\text{length}(I)}{1} + 0.9 \cdot \frac{\text{length}(I^c)}{1} =$$

$$0.2 \cdot 0.6 + 0.1 \cdot 0.5 = 0.16 //$$

this is of course the minimal true error

① from the results we can see that the empirical error increases with  $m$  - it is easy to find 3 intervals for correct classification of several points, but with an increasing sample size we expect some points that will not be consistent with our hope to classify many neighboring points with the same label using an interval. (see plot of section a). still, in any interval,  $P$  can generate either 0 or 1 labels. The true error however decreases with  $m$ . the more samples we have, the more likely we are to choose  $h$  that classifies with condition similar (Although not identical) to  $h^*$  (the best hypothesis)

② As we expect, the empirical error drops with  $k$  - more degrees of freedom in choosing intervals that will classify more sample points correctly. The true error however, increases with  $k$ . This is of course consistent with the fact that it should increase with  $|H|$  as we seen in class.

$k^* = k$  with the best empirical error. we got  $k^* = 10$  but it is of course not the best  $k$ , but merely the best for describing the sample. (with  $k=0 \dots 10$ )

if we had ran the algorithm with largest  $K$ 's, we would get larger  $K^*$  (expectedly, the largest  $K$ ) because we had more degrees of freedom in the sample classification. it would not, however decrease the true error that after  $K=3$  will increase due to poorer representation of the true distribution.

(e) Qn. 5 in theoretical part:

we got  $\text{VC dim} = 2k$ .

Also we have

$$P\left[\sup_{h \in H} |\ell_s(h) - \ell_p(h)| \geq \varepsilon\right] \leq n \prod_H(n) e^{-\frac{n\varepsilon^2}{8}}$$

$$\text{where } \prod_H(n) \leq \left(\frac{en}{d}\right)^d \quad (d = \text{VC dim})$$

$$\Rightarrow P\left[\sup_{h \in H} |\ell_s(h) - \ell_p(h)| \geq \varepsilon\right] \leq n \cdot \left(\frac{2en}{2k}\right)^{2k} e^{-\frac{n\varepsilon^2}{8}}$$

$$= n \left(\frac{n}{2k}\right)^{2k} e^{2k - n\varepsilon^2/8} \stackrel{!}{\leq} \sigma$$

$$e^{2k - n\varepsilon^2/8} \leq \frac{\sigma}{n} \left(\frac{n}{2k}\right)^{2k}$$

$$2k - n\varepsilon^2/8 \leq \ln\left(\frac{\sigma}{n}\right) + 2k \ln\left(\frac{n}{2k}\right)$$

$$\varepsilon^2 \geq \frac{16k}{n} + \frac{8}{n} \ln\left(\frac{n}{\sigma}\right) + \frac{16k}{n} \ln\left(\frac{n}{2k}\right)$$

$$\varepsilon \geq \sqrt{\frac{8}{n} \left(2k + 2 \ln\left(\frac{n}{2k}\right) + \ln\left(\frac{n}{\sigma}\right)\right)}$$

$$\sigma = 0.1$$

$$\Rightarrow \ell_p(h) \leq \ell_s(h) + \sqrt{\frac{8}{n} \left(2k + 2 \ln\left(\frac{n}{2k}\right) + \ln n\right)}$$

we can see that the penalty  $\sqrt{k - \ln k}$   
which increases with  $k$ . smallest penalty for  
 $k=1$ . The last inequality will hold

in prob  $\geq 1 - \sigma$

in section d we had  $k_{\text{best}} = 3$  of course,  
as explained in section b.

see plot for all errors wrt  $k$

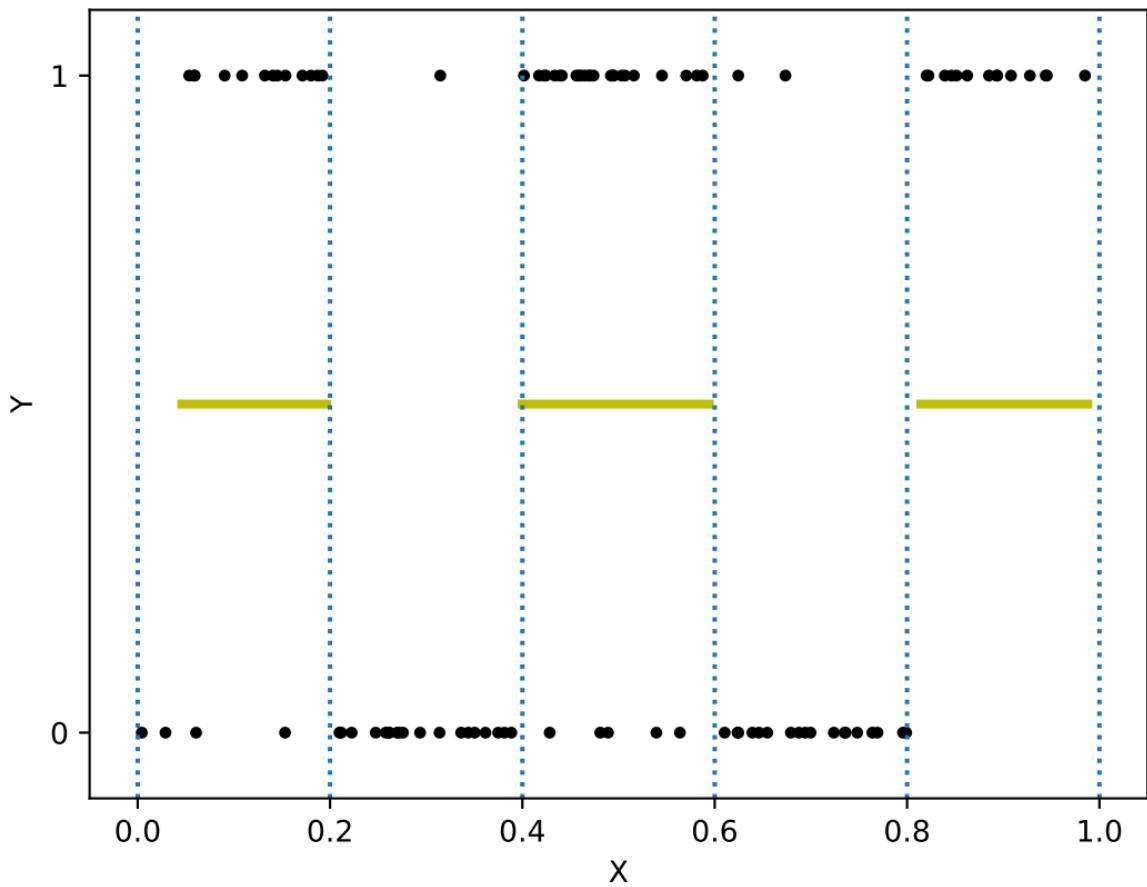
(t)

Among 3 experiments, each for different values of  $K$  and a different subset of  $m$  points as a holdout set, we got that the best  $K$  (in most exp.) is  $\underline{K=3}$  (sometimes one of the exp yields best  $K=4$ ). This is not surprising of course as we expect that the holdout set would simulate the true distribution and therefore the error on it shall be close to the true error given the  $K$  intervals returned by the ERM on the train data.

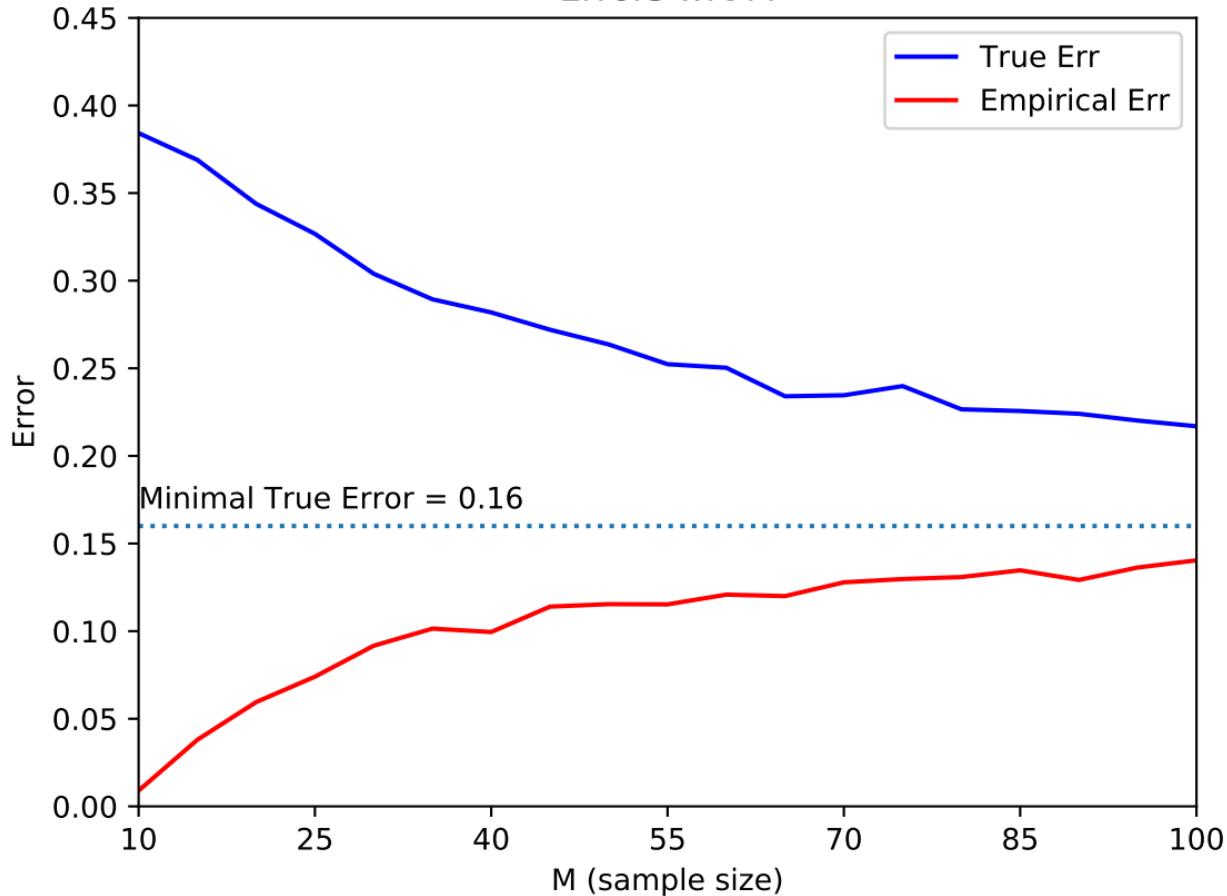
As discussed in section 2e, The best hypothesis shall be 3 disjoint intervals identical to those with the higher prob. for  $y=1$  in the true distribution, in which all points are classified as  $h(x)=1$ . The holdout set helps us find a hypothesis close to it because it represents the true distribution in nature. (see student for the intervals of each  $K$ )

(see figure for 3 exp. result,  
holdout error wrt  $K$ )

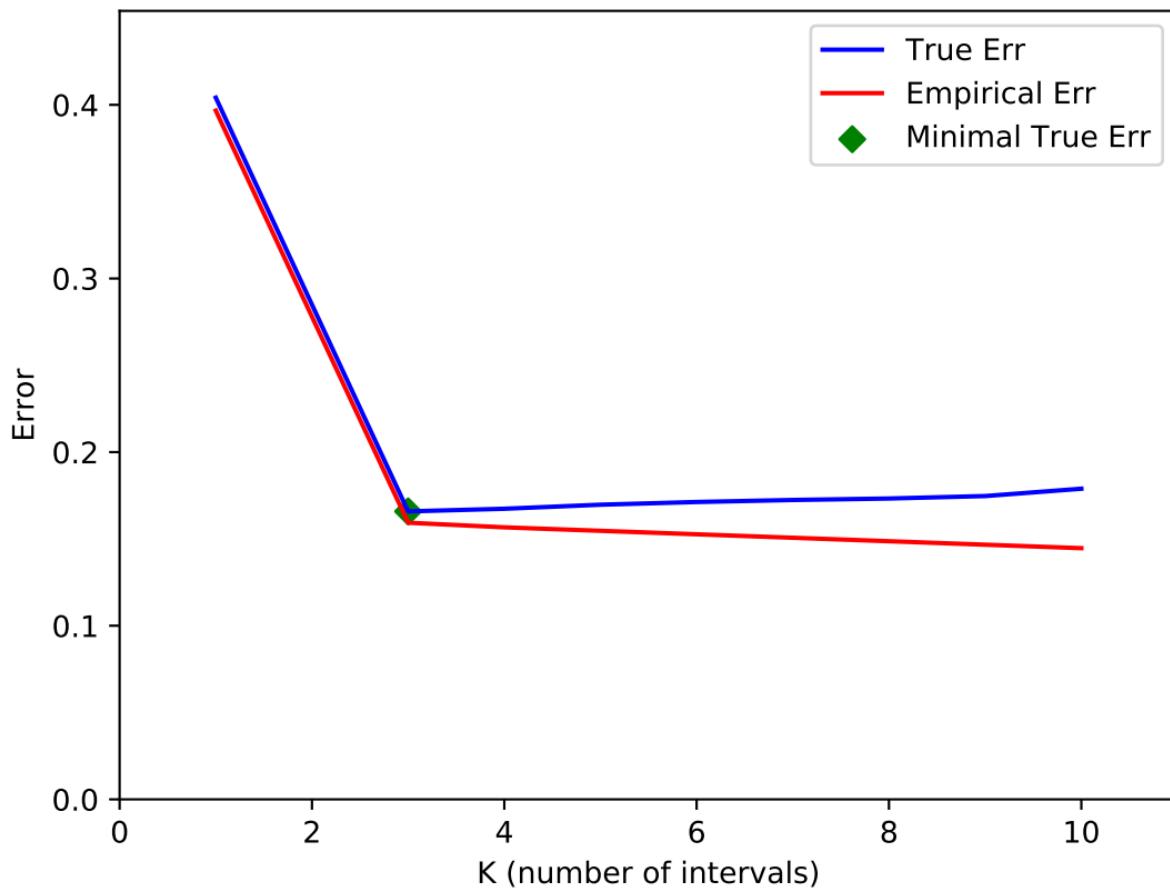
## Dataset Distribution and Best Intervals



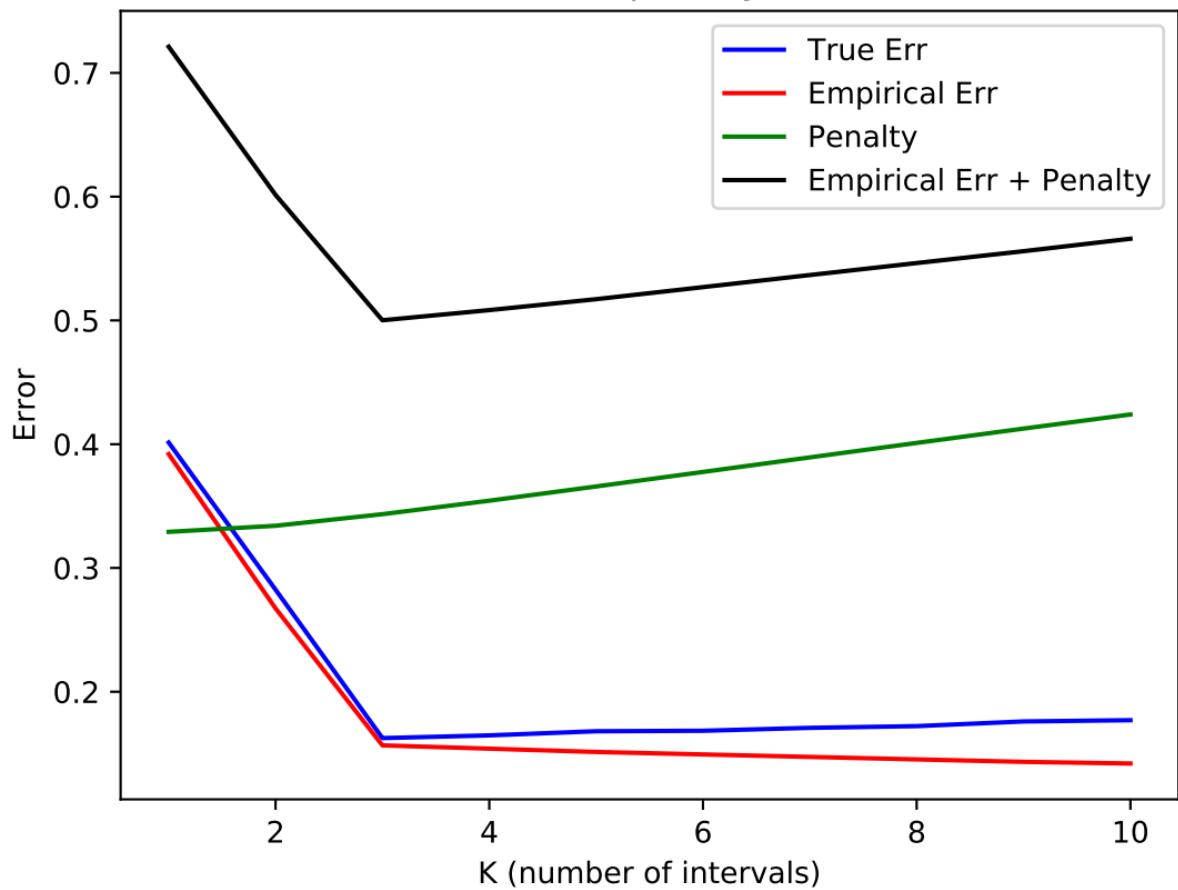
### Errors wrt M



### Errors wrt K



## Errors and penalty wrt K



### Hold out Errors wrt K

