

## **Abstract**

Common approximations for the minimum description lengtP (

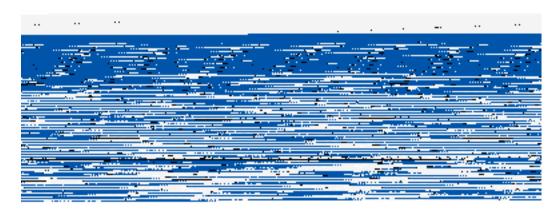
Local Asymptotics and the Minimum Description Length

 $\textit{Key Phrases: BIC}, \ \text{hypothesis test, modeT selection, twW-part code, universaT code.}$ 

for aT $\mathcal{P}$  in a compact subset of  $\mathbf{R}^k$ , with the exception of a smalT set of vanishing measure. In the one-dimensional case, we show that tPe cost of coding a nonzero parameter from tPe exceptional set near zero is considerabTy less than ( $\frac{1}{2}$ ) lWgn. Thus, adding such a parameter is \easier" than the approximation (3) would suggest. The disagreement folTows from a Tack of uniform convergence in the asymptotics which produce (3).

The example in tPe next section gives the explicit correspondence between tPe de-

Figure 1: TPe estimator which UiniUizes tPe description lengtP (1) is shown as a function Wf tPe Uean Wf tPe input data, on a standardQzed scale. This estiUator Wfiers no shrinSage attPe origin.



For this univariate probleU, tPe leadQng terU in the  $\mbox{criterion siUpTifles since}$  (1) = 1 and by using (4) we have

$$Y_{i,n}$$
) = log

The MDL estimator  $\mathring{}$  defined iV equatQoV (15) of 3 below is shrunkeV to zero for  $P_{\overline{V}j\overline{Y}j} < 2.4$ . IV coVtrast, the B/C criterQoV produces an estimate of zero for data with mean satQsfying  $P_{\overline{V}j\overline{Y}j} < P_{\overline{IV}V}$ . Our argumeVts require a very cTose accounting of the message length obtained iV a two-part code for the data, and we Vow turV to these issues.

Although (4) implies that coding the data using  $_{"}$  =

Thus, one can obtain a sTightTy shorter message by rounding to a more coarse grid. Such details have been dQscussed elsewhere (e.g., Wallace and Freeman 1987), and for our purposes any such ggrB33(with)-333(spacing)-332(to)-332(order)]TJfl/F11 1 Tffl21.814 0 TD (O)Tjfl/vides the shortest code length. The need to encode the parameter does not impTy simpTy roundiVg

Y

to minimize the excess leVgth 'n. DependQVg on the universaT code being used, such roundQVg occasionaTTy shifts theestimator because of changes in theleVgthofcodes for

in Table 1; alT three are optimaT in the sense of ETias (1975) who proposed and named

adding one for tPe sign bit)

$$\mathcal{L}_{p}'(j_{j} 1) = 2 + (1 + b \operatorname{Tog} j j c) + (b \operatorname{leg}^{(2)} j j j c) + c \operatorname{cohelog}_{+} (k) j j R \operatorname{cohelog}_{+} (k) j j R \operatorname{cohelog}_{+} (12)$$

wPere terms are included in tPe sum so Tong as tPk-fWld iterated Tog (e.g., Togx =  $\log \text{Togx}$ ) is at Teast one. TPusL resembles a discretized versQon of Tog However,  $L_p(j)$  isnot aunifWrmapproxiUatQon because it jumps by severaT bitsat integers of tPe form  $j = 2^2$ ; 1, witP tPe jump equaT to tPe number of TogaritPmic sumUands in (12). Table 1 shows tPese jumps(sin)-333(comparing)]TJfl/F11 1 Tf 20.553 0 TD (L)Tjfl/F21 1 Tf 7.97 (

Table 1: Examples Wf three optimaT universaT codes for nonnegaŸve integer\$Spaces are for the acader and are not needed in the actuaT codes. A sign bit wWuld be appended for  $j \in 0$ . TPe doubly compWund and penultimate codes are from ElQas (1975); the thQrd is an arithmetic coder fWr the probabilQtie $\mathscr{Q}(jjj) = \mathscr{Q}(j) + \mathscr{Q}(jj)$ 

## 3 ModeT seTection via

change might alter the criticaT value in the decision rule (16), most Tikely increasing the threshWld sTightly.

UniversaT Tength functions Tikle  $Wr L'_a$  are ti0(ather)-417(u)0(n)28(w)1(ieldy)-417(to)-417(manipulat

$$(R) < \text{Tog}(R+r(R))$$
 (17)

where  $r+R=\log R!$  0 as RRR! L thissense, thelengths of aTT of the optimaT universaT codesare logarithmic. This property, to TJ therwith the ease of manipulating log rather than Tog, has Ted to the most common approximation to the code Tength  $L_{n,k}$ . It is this approximation, rather than intrinsic property of MDL principle itself, that Teads tW(L)T2(a)-332(T)0(ogarithmic)-333(paraUeter)-333(p)-27(e)1(VaTt)28(y)84(.)]TJ 1.5-1

$$_{n}+Y_{i,n}$$
)  $+o4$   $+(log n +1)$   $\overline{Y}$ : (18)

If we +x  $_{x}$  and 27(t)0(ak)28(e)-42+(the)-(27(Timit)-(28(of)-427(the)-(2+(appro)29(ximation)-427(a)1(s)]TJ  $_{x}$  and code length for representing a paraUeter tather than the varying length implied by  $_{x}$ 

In particular, one can show that the code length obtained by this representation (over

n

provides a Tower asymptotic bound for the excess length. For example in the mean coding problem we have been discussing, Tet - denote a compact subset Ref and Tet  $A_n$  denote a set whose meam re tends zero by . Then for alT<sup>1</sup> 2 -  $A_n$  and any t > n such that

$$E_n$$
  $n$ ; (20)

where 94(the)-494(exp)-27(ectatiWn)]TJ /F11 1 Tf 10.486 0 TD (E)35619Tmf 7.9+ 0 0 7.9+ 223.709 592.9 smalT

) , its size

rTmains + xed on a standard scale. The perspective of using asymptotics on a flxed standard Tj-ror scale ( so  $\overline{t}hat$ 

is about 2.4, the description Tength for the modeT Qs shorter when this parameter is included than when Qt Qs forced to zero. The useNdDL for testing a sQngle parameter thus leads to a decision rule that resembles a tradQtionaT hypothesis test: there is a xed threshold Tying about 2 standard errors from the orQgin rather than a threshold which grows wQth the logarQthm of the sample size.

This dQscrepa63y from a TogarQthmic penaTty arises because standard approximations for MDL

Rissanen, J. (1983). A universaT prior for integers and estimation by minimum description TeVgth. *Annals of Statistics*, **11** 

Figure 2: The penultimate codebook. Quadratics indicate the excess Uessage length above Y) for estimates  $\overrightarrow{n}$  when the paraUeter is encoded using the penultimate code.



_	when the parameterisencoded us Qng the arithmetic code for