

## SOME INEQUALITIES FOR INFORMATION DIVERGENCE AND RELATED MEASURES OF DISCRIMINATION

Flemming Topsøe

ABSTRACT. Inequalities which connect information divergence with other measures of discrimination or distance between probability distributions are used in information theory and its applications to mathematical statistics, ergodic theory and other scientific fields. We suggest new inequalities of this type. One of them is a sharpening of the well known Pinsker inequality. Our results depend on two measures of discrimination, called capacity discrimination and triangular discrimination. The discussion contains references to related research. This concerns, in particular, Csiszár's f-divergences and the Hellinger distance.

### 1 INTRODUCTION

We shall study two probability measures  $P$  and  $Q$  over a finite alphabet (for a more general set-up, see the discussion, Section VI). The point probabilities corresponding to  $P$  and  $Q$  are denoted by  $(p_i)$ , respectively  $(q_i)$ . Apart from information divergence  $D(P||Q)$  and variational distance  $V(P, Q)$ , we consider two other measures of discrimination between  $P$  and  $Q$ , viz. *capacity discrimination* given by  $C(P, Q) = D(P||M) + D(Q||M)$  where  $M = \frac{1}{2}(P + Q)$  and *triangular discrimination* given by

$$\Delta(P, Q) = \sum \frac{|p_i - q_i|^2}{p_i + q_i}.$$

We also consider a "directional" version of  $\Delta$ , denoted  $\Delta^*$ , and defined by

$$\Delta^*(P||Q) = \sum_{k=0}^{\infty} 2^k \Delta(M_k, Q).$$

Here,  $M_k = 2^{-k}P + (1 - 2^{-k})Q$ .

As we shall see below, capacity and triangular discrimination behave similarly. Indeed,

$$\frac{1}{2}\Delta(P, Q) \leq C(P, Q) \leq \log 2 \cdot \Delta(P, Q).$$

Via a general identity, this implies that

$$\frac{1}{2}\Delta^*(P||Q) \leq D(P||Q) \leq \log 2 \cdot \Delta^*(P||Q).$$

Note that the inequality  $D \geq \frac{1}{2}\Delta^*$  is a strengthening of Pinsker's inequality since  $\Delta^* \geq V^2$ .

The importance of inequalities for information divergence in mathematical statistics, information theory proper and other fields is well recognized. In particular, this is true for Pinsker's inequality.

### 2 DEFINITIONS AND AUXILIARY RESULTS

By  $A = \{a_i \mid 1 \leq i \leq n\}$  we denote an  $n$ -set, the *alphabet*. Distributions over  $A$ , always assumed to be probability distributions, are, typically, denoted by  $P, Q, \Pi$ , and the associated point probabilities are denoted by  $p_i, q_i$  and  $\pi_i$ , respectively. *Variational distance* ( $\ell_1$ -distance) and *information divergence* (Kullback-Leibler divergence) are defined as usual:

$$V(P, Q) = \sum_{i=1}^n |p_i - q_i|,$$
$$D(P||Q) = \sum_{i=1}^n p_i \log \frac{p_i}{q_i}.$$

---

*Key words and phrases.* Information divergence, Pinsker's inequality, capacity discrimination, triangular discrimination, Hellinger distance.

Here,  $\log$  denotes natural logarithm.

Information divergence is the basis for a new measure of discrimination, called *capacitory discrimination*, which is defined by

$$(2.1) \quad C(P, Q) = D(P\|M) + D(Q\|M),$$

where  $M = \frac{1}{2}(P + Q)$ . Furthermore, we define *triangular discrimination* between  $P$  and  $Q$  by

$$(2.2) \quad \Delta(P, Q) = \sum_{i=1}^n \frac{|p_i - q_i|^2}{p_i + q_i}.$$

A variant of this measure, depending on a natural number  $\nu$  as parameter, and called *triangular discrimination of order  $\nu$* , will also be needed. It is defined by the equation

$$(2.3) \quad \Delta_\nu(P, Q) = \sum_{i=1}^n \frac{|p_i - q_i|^{2\nu}}{(p_i + q_i)^{2\nu-1}}.$$

For  $\nu = 1$  we are back to (2.2), i.e.  $\Delta_1 = \Delta$ . In the discussion we point out that  $\Delta_\nu$  is a power of a distance (a metric), but for the main results we do not need this fact.

Fix two probability distributions  $P$  and  $Q$  over  $A$  and consider the communication channel with two input letters, say 0 and 1, and with  $P$  and  $Q$  as the conditional distributions over  $A$  given 0 and 1, respectively. If  $(\alpha, \beta)$  defines an input distribution ( $\alpha, \beta \geq 0, \alpha + \beta = 1$ ), then the *information transmission rate*, here denoted  $I(\alpha, \beta)$ , is defined as usual, i.e.

$$(2.4) \quad I(\alpha, \beta) = \alpha D(P\|S) + \beta D(Q\|S)$$

where  $S = \alpha P + \beta Q$  is the *output distribution induced by  $(\alpha, \beta)$* . We see that

$$(2.5) \quad C(P, Q) = 2 \cdot I\left(\frac{1}{2}, \frac{1}{2}\right).$$

The *maximum transmission rate*, i.e. the capacity, is defined by

$$I_{\max} = \max_{(\alpha, \beta)} I(\alpha, \beta).$$

With any distribution  $\Pi$  over  $A$ , now conceived as a *predictor*, we associate a (*guaranteed*) *redundancy*  $R(\Pi)$ , defined by

$$(2.6) \quad R(\Pi) = \max(D(P\|\Pi), D(Q\|\Pi)).$$

The *minimum (guaranteed) redundancy* is the quantity

$$(2.7) \quad R_{\min} = \min_{\Pi} R(\Pi).$$

The above concepts are well known from the theory of universal coding and prediction. Here, we are dealing with a particularly simple instance of that theory. From the general theory we need a key fact, the *Gallager-Ryabko theorem*, cf. the discussion in Ryabko [24] for the history of this result, which tells us that

$$(2.8) \quad I_{\max} = R_{\min}.$$

We note that the “ $\leq$ -part” of (2.8), which is the part we need, follows immediately from a simple identity (called the *compensation identity*, in some planned publications). In the present context this identity states that for any input distribution  $(\alpha, \beta)$  and any predictor  $\Pi$ ,

$$(2.9) \quad I(\alpha, \beta) + D(S\|\Pi) = \alpha D(P\|\Pi) + \beta D(Q\|\Pi)$$

where  $S$  is the output distribution induced by  $(\alpha, \beta)$ .

The elementary inequalities

$$(2.10) \quad \frac{1}{2}C \leq I_{\max} \leq R_{\min} \leq R(M) \leq C$$

serve as motivation for choosing the name “capacitory discrimination” for  $C$ . *Note:* Here,  $C$  stands for  $C(P, Q)$ ; similar abbreviations are used in the sequel for  $\Delta(P, Q)$  and  $V(P, Q)$ .

Regarding triangular discrimination <sup>1</sup> note that

$$(2.11) \quad \frac{1}{2}V^2 \leq \Delta \leq V$$

The last inequality is trivial, and the first follows by writing  $\Delta$  as the square of an  $\ell_2$ -norm w.r.t. the probability measure  $\frac{1}{2}(P+Q)$  and using the fact that this  $\ell_2$ -norm dominates the corresponding  $\ell_1$ -norm (for a generalization, see (6.4)).

### 3 A CONNECTION BETWEEN CAPACITORY AND TRIANGULAR DISCRIMINATIONS

Our first observation displays a connection between capacity and triangular discrimination (of varying orders).

**Theorem 3.1.** *For any distributions  $P$  and  $Q$  over the alphabet  $A$ ,*

$$(3.1) \quad C(P, Q) = \sum_{\nu=1}^{\infty} \frac{1}{2\nu(2\nu-1)} \Delta_{\nu}(P, Q).$$

*Proof.* In addition to notation from Section II we introduce the following auxiliary quantities ( $1 \leq i \leq n$ ):

$$m_i = \frac{1}{2}(p_i + q_i), \quad \varepsilon_i = |p_i - m_i|, \quad k_i = m_i/\varepsilon_i.$$

To simplify the exposition we assume that, for all  $i$ ,  $p_i \neq 0$ ,  $q_i \neq 0$  and  $p_i \neq q_i$  hold. Note that  $1/k_i = |p_i - q_i|/(p_i + q_i)$  so that  $0 < 1/k_i < 1$ . We have:

$$\begin{aligned} C(P, Q) &= \sum_{i=1}^n p_i \log \frac{p_i}{m_i} + \sum_{i=1}^n q_i \log \frac{q_i}{m_i} \\ &= \sum_{i=1}^n \varepsilon_i \left( (k_i + 1) \log \frac{k_i + 1}{k_i} + (k_i - 1) \log \frac{k_i - 1}{k_i} \right) \\ &= \sum_{i=1}^n \varepsilon_i \left( k_i \log \left( 1 - \frac{1}{k_i^2} \right) \right. \\ &\quad \left. + \log \left( 1 + \frac{1}{k_i} \right) - \log \left( 1 - \frac{1}{k_i} \right) \right) \\ &= \sum_{i=1}^n \varepsilon_i \sum_{\nu=1}^{\infty} \frac{1}{\nu(2\nu-1)} k_i^{-(2\nu-1)}. \end{aligned}$$

Reversing the order of summation gives (3.1). ■

As a corollary we obtain:

**Theorem 3.2.** *Capacity and triangular discrimination behave similarly as they are connected by the inequalities*

$$(3.2) \quad \frac{1}{2}\Delta(P, Q) \leq C(P, Q) \leq \log 2 \cdot \Delta(P, Q).$$

*Proof.* The first inequality follows by considering only the first term in (3.1). Noting that  $\Delta_{\nu}(P, Q)$  decreases with increasing  $\nu$ , we also deduce from (3.1) that

$$C(P, Q) \leq \sum_{\nu=1}^{\infty} \frac{1}{2\nu(2\nu-1)} \Delta_1(P, Q) = \log 2 \cdot \Delta(P, Q),$$

as claimed. ■

---

<sup>1</sup> the motivation behind the chosen terminology is connected with the *triangular net* in  $\mathbb{R}_+^n$  which consists of all points of the form  $(t_{j_1}, \dots, t_{j_n})$  where the  $j_i$ 's are non-negative integers and  $t_j = \frac{1}{2}j(j+1)$ , the  $j$ 'th *triangular number*. Then note that any set of neighbouring points in the triangular net are unit-distance apart measured by  $\Delta$ , i.e.  $\Delta(x, y) = 1$  for any such set of points (this requires an extension of the definition (2.2) to arbitrary points in  $\mathbb{R}_+^n$ ).

Thus, numerically,  $0.50 \leq C/\Delta \leq 0.70$ .

Combining with (2.10) and (2.11) we get:

$$\begin{aligned} \frac{1}{8}V^2 &\leq \frac{1}{4}\Delta \leq \frac{1}{2}C \leq I_{\max} \leq R_{\min} \\ &\leq C \leq \log 2 \cdot \Delta \leq \log 2 \cdot V \end{aligned}$$

#### 4 INEQUALITIES FOR INFORMATION DIVERGENCE

As in the previous sections we study two distributions  $P$  and  $Q$  over the  $n$ -letter alphabet  $A$ . With  $P$  and  $Q$  we associate the *successive midpoints*  $(M_\nu)_{\nu \geq 0}$  in the direction  $Q$  which are given by  $M_0 = P$  and  $M_\nu = \frac{1}{2}(M_{\nu-1} + Q)$ , i.e.

$$M_\nu = 2^{-\nu}P + (1 - 2^{-\nu})Q; \quad \nu \geq 0.$$

We intend to show that the following quantity:

$$\Delta^*(P\|Q) = \sum_{k=0}^{\infty} 2^k \Delta(M_k, Q).$$

behaves much like  $D(P\|Q)$ . Note that

$$(4.1) \quad \Delta^*(P\|Q) = \sum_{\nu=1}^{\infty} \sum_{i=1}^n \frac{|p_i - q_i|^2}{p_i + (2^\nu - 1)q_i}.$$

**Theorem 4.1.** *Let  $P$  and  $Q$  be distributions over the alphabet  $A$  and consider the successive midpoints  $(M_\nu)_{\nu \geq 0}$  in the direction  $Q$ . The following identity and inequalities hold:*

$$(4.2) \quad D(P\|Q) = \sum_{\nu=0}^{\infty} 2^\nu C(M_\nu, Q),$$

$$(4.3) \quad \frac{1}{2}\Delta^*(P\|Q) \leq D(P\|Q) \leq \log 2 \cdot \Delta^*(P\|Q).$$

*Proof.* If  $D(P\|Q) = \infty$ , there exists  $i$  with  $q_i = 0$  and  $p_i > 0$ . Then, by (3.2) and (4.1),  $\sum_{\nu=0}^{\infty} 2^\nu C(M_\nu, Q) \geq \frac{1}{2} \sum_{\nu=0}^{\infty} 2^\nu \Delta(M_\nu, Q) = \infty$ , hence (4.2) and (4.3) hold in this case.

Now assume that  $D(P\|Q) < \infty$ . From the compensation identity (2.9) with  $\Pi = Q$  and  $\alpha = \beta = \frac{1}{2}$  we get

$$D(P\|Q) = C(P, Q) + 2D(M\|Q).$$

Iterating this, we find, for every  $k \geq 0$ ,

$$(4.4) \quad D(P\|Q) = \sum_{\nu=0}^k 2^\nu C(M_\nu, Q) + 2^{k+1} D(M_{k+1}\|Q).$$

A direct calculation shows that here, the last term tends to 0 as  $k \rightarrow \infty$ . Indeed,

$$\begin{aligned} 2^k D(M_k\|Q) &= \sum_{i=1}^n \left( \frac{(p_i - q_i)^2}{q_i} 2^{-k} + (p_i - q_i) \right) \\ &\quad \cdot \frac{\log(1 + 2^{-k} \frac{p_i - q_i}{q_i})}{2^{-k} \frac{p_i - q_i}{q_i}} \\ &\rightarrow \sum_{i=1}^n (p_i - q_i) = 0. \end{aligned}$$

With this auxiliary result at hand, (4.2) follows readily from (4.4) and then, (4.3) follows from (3.2). ■

Pinsker's inequality  $D \geq \frac{1}{2}V^2$  is an immediate consequence of the left hand side inequalities of (2.11) and (4.3). Indeed, as  $\Delta \geq \frac{1}{2}V^2$ ,

$$(4.5) \quad D(P\|Q) \geq \frac{1}{2} \sum_{\nu=0}^{\infty} 2^{\nu} \frac{1}{2} (V(M_{\nu}, Q))^2 = \frac{1}{2} V(P, Q)^2.$$

Due to the fact that  $D$  may be large, indeed infinite, while  $V$  is small, upper bounds for information divergence in terms of discrimination measures such as variational distance cannot hold without introducing extra terms.

**Theorem 4.2.** *Let  $P$  and  $Q$  be distributions over  $A$  and put  $c = \max_i(p_i/q_i)$ . Then*

$$(4.6) \quad \begin{aligned} D(P\|Q) &\leq \log 2 \cdot V(P, Q) + \log \left(\frac{1}{2}(1+c)\right) \\ &\leq \log 2 \cdot V(P, Q) + \log c. \end{aligned}$$

*Proof.* Put  $d = \frac{1}{2}(1+c)$ . Then  $p_i/q_i \leq dp_i/m_i$  for all  $i$  (again,  $m_i = \frac{1}{2}(p_i + q_i)$ ). Therefore, by (2.10) and Theorem 2,

$$\begin{aligned} D(P\|Q) &\leq \sum_{i=1}^n p_i \log \frac{dp_i}{m_i} = D(P\|M) + \log d \\ &\leq C(P, Q) + \log d \end{aligned}$$

which, in view of (3.2) and (2.11), is even stronger than (4.6). ■

## 5 SOME REFINEMENTS

We shall focus on the lower bounds  $C \geq \frac{1}{4}V^2$  and  $D \geq \frac{1}{2}V^2$  (Pinsker's inequality). These inequalities can be conceived as giving only the first term in an infinite expansion. It appears to be more natural to use half the variational distance rather than variational distance itself as parameter. Choosing the parameter this way, it varies in the unit interval.

**Theorem 5.1.** *For any distributions  $P$  and  $Q$  over the alphabet  $A$ ,*

$$(5.1) \quad C(P, Q) \geq \sum_{\nu=1}^{\infty} a_{\nu} \left(\frac{1}{2}V(P, Q)\right)^{2\nu}$$

where the constants  $a_{\nu}$ , all positive, are given by

$$(5.2) \quad a_{\nu} = 2 \frac{1}{2\nu(2\nu-1)}; \quad \nu \geq 1.$$

*Equality holds in (5.1) if and only if, either  $P = Q$  or else,  $P$  and  $Q$  are supported by the same 2-set and are symmetrically placed (in the sense that  $\frac{1}{2}(P+Q)$  is the uniform distribution over the 2-set in question).*

*The constants  $a_{\nu}$  are best possible in the natural sense (cf. below).*

*Proof.* We shall apply the data reduction inequality for information divergence (cf. Csiszár [6], Csiszár and Körner [8]). For another proof, see the discussion.

Let  $\partial$  denote the reduction defined by decomposing  $A$  into  $A^+ = \{i|p_i \geq q_i\}$ , and  $A^- = \{i|p_i < q_i\}$ .

By this reduction,  $P$  and  $Q$  are replaced by  $\partial P$  and  $\partial Q$ , respectively, where these measures are concentrated on the 2-set  $\partial A = \{A^+, A^-\}$  and have point masses  $P(A^+)$ ,  $P(A^-)$  and  $Q(A^+)$ ,  $Q(A^-)$ , respectively. Note, firstly, that this reduction leaves the variational distance unchanged:  $V(P, Q) = V(\partial P, \partial Q)$  and, secondly, that the two pairs  $(P, M)$  and  $(M, Q)$  (with  $M = \frac{1}{2}(P+Q)$  as usual) lead to the same reduction as that defined by the pair  $(P, Q)$ .

From Theorem 1 we can now conclude that

$$\begin{aligned} C(P, Q) &\geq C(\partial P, \partial Q) = \sum_{\nu=1}^{\infty} \frac{1}{2\nu(2\nu-1)} \Delta_{\nu}(\partial P, \partial Q) \\ &= \sum_{\nu=1}^{\infty} \frac{(\frac{1}{2}V(P, Q))^{2\nu}}{2\nu(2\nu-1)} \left( \frac{1}{(P(A^+) + Q(A^+))^{2\nu-1}} \right. \\ &\quad \left. + \frac{1}{(P(A^-) + Q(A^-))^{2\nu-1}} \right) \\ &\geq 2 \sum_{\nu=1}^{\infty} \frac{(\frac{1}{2}V(P, Q))^{2\nu}}{2\nu(2\nu-1)} \end{aligned}$$

which is the result stated. Note that the last inequality above follows from elementary considerations as it is easy to show that, for  $n \in \mathbb{N}$ , the minimal value of  $x^{-n} + y^{-n}$ , where  $x$  and  $y$  are non-negative numbers with sum 2, is 2, which is attained for  $x = y = 1$  (and for no other values of  $x$  and  $y$ ).

Really, the fact stated concerning equality in (5.1), follows by inspection of the above proof.

Concerning the last statement of the theorem, we first define the *best constants*  $a_{\nu}^{\max}$  for an inequality like (5.1). This is done by recursion: For each  $\nu_0 \geq 1$ ,  $a_{\nu_0}^{\max}$  is defined as the maximum over all  $a_{\nu_0}$  where  $a_{\nu_0}$  fits into an inequality of the form (5.1) with  $a_{\nu} = a_{\nu}^{\max}$  for all  $\nu < \nu_0$ .

The proof that  $a_{\nu} = a_{\nu}^{\max}$  for all  $\nu \geq 1$  with  $a_{\nu}$  given by (5.2) is then carried out by induction: Let  $\nu_0 \geq 1$  and assume that  $a_{\nu} = a_{\nu}^{\max}$  for  $\nu < \nu_0$ . By (5.1),  $a_{\nu_0}^{\max} \geq a_{\nu_0}$ . To prove the reverse inequality, note that by the result proved regarding instances for which equality holds in (5.1) it follows that for all  $0 \leq x \leq 1$ ,

$$\sum_{\nu=1}^{\infty} a_{\nu} x^{2\nu} \geq \sum_{\nu=1}^{\nu_0} a_{\nu}^{\max} x^{2\nu}.$$

This implies that for all  $0 < x \leq 1$ ,

$$(a_{\nu_0} - a_{\nu_0}^{\max}) + \sum_{\nu=\nu_0}^{\infty} a_{\nu} x^{2(\nu-\nu_0)} \geq 0,$$

and  $a_{\nu_0} \geq a_{\nu_0}^{\max}$  follows. Thus  $a_{\nu_0} = a_{\nu_0}^{\max}$  must hold. ■

By summing the infinite series in (5.1) we obtain a corollary which in essence is equivalent to (5.1). Indeed, putting  $x = \frac{1}{2}V$ , one finds that

$$(5.3) \quad \begin{aligned} C(P, Q) &\geq (1+x) \log(1+x) + (1-x) \log(1-x) \\ &= 2 \cdot D((\alpha, \beta) \| (\frac{1}{2}, \frac{1}{2})), \end{aligned}$$

with  $\alpha = \frac{1}{2}(1+x) = \frac{1}{2} + \frac{1}{4}V$  and  $\beta = \frac{1}{2}(1-x) = \frac{1}{2} - \frac{1}{4}V$ . By combining with (4.2) one further finds that

$$(5.4) \quad D(P \| Q) \geq \sum_{\nu=1}^{\infty} b_{\nu} (\frac{1}{2}V(P, Q))^{2\nu}$$

with

$$(5.5) \quad b_{\nu} = \frac{2^{2\nu}}{2^{2\nu-1} - 1} \cdot \frac{1}{2\nu(2\nu-1)}; \quad \nu \geq 1.$$

## 6 DISCUSSION

Capacitory discrimination as introduced here can be viewed as a symmetrized and smoothed variant of information divergence. Furthermore, it has an interpretation as twice an information transmission rate, cf. (2.5). It is natural to compare capacitory discrimination with Jeffrey's measure of distance between two distributions given by  $J(P, Q) = D(P \| Q) + D(Q \| P)$ , cf. e.g. Csiszár and Körner [8] (Problem I.3.16). It is intuitively clear that  $C \leq J$ . In fact, even  $C \leq \frac{1}{2}J$  holds. The proof of Theorem 3 actually shows that

$$(6.1) \quad \frac{1}{2}J(P, Q) - C(P, Q) = D(M \| P) + D(M \| Q).$$

For many of our results, we introduced some extra distributions, typically midpoints, in addition to the basic measures under study. Though related only weakly to results presented here, the reader may wish to consult Dembo [11], Proposition 1, where wider instances of qualitatively similar inequalities are proved in order to study the measure concentration problem.

Triangular discrimination arose in an attempt to sharpen a lower bound of  $\varepsilon$ -capacity given in Ryabko and Topsøe [25]. This theme may be pursued in a later publication. Here, the significance of the triangular discrimination measures lies in the strong connection with capacity discrimination and information divergence (Theorems 1 and 3).

In fact, triangular discrimination has occurred before in the literature, perhaps for the first time in Vincze [28] where its statistical significance is briefly indicated. We also find triangular discrimination in LeCam [19] (cf. p.xvii and pp. 47–48) who refers to  $\Delta$  as a “chi-square like distance”, and LeCam indicates that this measure of discrimination was already used by Hellinger. However, the author has not been able to check up on this.

The importance for statistical research of discrimination measures like  $\Delta_\nu$  and, more generally, of arbitrary  $f$ -divergences (see below) is further discussed in Österreicher and Vajda [22]. For our purposes, the observation by Kafka, Österreicher and Vincze [13] (cf. also the announcement in Feldman and Österreicher [12]) that the triangular discrimination measures  $\Delta_\nu$  are all powers of metrics (indeed,  $(\Delta_\nu)^{\frac{1}{2\nu}}$  is a metric, and the exponent here is the largest possible one) is most noteworthy. In particular,  $\Delta$  is the square of a metric. This was in fact already noticed by LeCam [19]. Using identities and inequalities of Theorems 1–3, this makes it possible to relate information divergence to true metrics. This is of importance for (parts of) research as that of LeCam [19], Birgé [2], Birgé and Massart [3] and Yang and Barron [29], where the lack of metric properties for information divergence is compensated for by researching relations to true distances (typically then, the Hellinger distance is the preferred model metric).

The most important inequality discussed is without doubt Pinsker's inequality which has numerous applications. Some general references are Ahlswede and Wegener [1], Csiszár and Körner [8], Kullback [18] and Pinsker [23]. For more specific applications of Pinsker's inequality we can, for instance, point to Csiszár [9], Marton [20, 21], Ryabko and Topsøe [25] and Topsøe [26] for five different applications. The inequality originated with Pinsker [23] and was refined in Csiszár [6], cf. also Csiszár and Körner [8], Problem I.3.17. Note that our approach offers an alternative and rather elementary proof of Pinsker's inequality (for this remark note that the first inequality of (4.3) depends on (4.4) but does not need (4.2) for its proof).

The refined inequality for capacity discrimination given in Theorem 5 is in a satisfactory form, but the corresponding result quoted for information divergence ((23) with constants as in (5.4)) is not, as the constants are not best possible. It is known that the best two-term inequality is  $D \geq \frac{1}{2}V^2 + \frac{1}{36}V^4$  (Krafft [15]) and not  $D \geq \frac{1}{2}V^2 + \frac{1}{84}V^4$  as one might think, considering (23). Let us briefly indicate a proof of this (modifying the approach of Krafft [15]). We need only consider probability measures  $P$  and  $Q$  of the form  $P = (\frac{1-\alpha}{2}, \frac{1+\alpha}{2})$ ,  $Q = (\frac{1+\beta}{2}, \frac{1-\beta}{2})$ . A straight forward expansion, cf. Kambo and Kotz [14], shows that

$$(6.2) \quad D(P\|Q) = \sum_0^{\infty} \frac{1}{2\nu(2\nu-1)} T_\nu(\alpha, \beta)$$

with the polynomials  $T_\nu$  given by

$$(6.3) \quad T_\nu(\alpha, \beta) = \alpha^{2\nu} + 2\nu\alpha\beta^{2\nu-1} + (2\nu-1)\beta^{2\nu}.$$

Then  $T_\nu \geq 0$  (evident if  $\alpha$  and  $\beta$  are of the same sign, and, if not, factor out  $(\alpha + \beta)^2$ ). Noting that  $T_1 = (\alpha + \beta)^2 = V^2$  and that  $T_2 = \frac{1}{3}V^2(V^2 + 2(\alpha - 2\beta)^2)$ , it is easy to see that the two first best constants in an inequality like  $D \geq \sum_1^{\infty} c_\nu V^{2\nu}$  are  $c_1 = \frac{1}{2}$  and  $c_2 = \frac{1}{36}$ . For further details along these lines, we refer to Vajda [27], and Krafft and Schmitz [16]. The last mentioned reference adds the term  $cV^6$  with  $c = \frac{1}{288}$ . However, the best constant here is  $c = \frac{1}{270}$ , as the author may return to elsewhere.

The refined inequalities (5.1) and (23) can also be derived in a direct way from (3.1) and (4.2) by using the following inequalities which are of some independent interest:

$$(6.4) \quad \begin{aligned} 2^{-2\nu+1}V^{2\nu} &\leq 2^{-\nu+1}\Delta^\nu \\ &\leq \Delta_\nu \leq \Delta \leq V. \end{aligned}$$

The non-trivial parts follow as the numbers

$$\left(\frac{1}{2}\Delta_\nu(P, Q)\right)^{\frac{1}{2\nu}} = \left(\sum \left(\frac{|p_i - q_i|}{p_i + q_i}\right)^{2\nu} \left(\frac{1}{2}p_i + \frac{1}{2}q_i\right)\right)^{\frac{1}{2\nu}},$$

recognized as  $\ell_{2\nu}$ -norms w.r.t. a probability measure, are weakly increasing in  $\nu$  (allowing also the value  $\nu = \frac{1}{2}$  for which  $\Delta_\nu = V$ ). Note that the constants in (6.4) are best possible as equality holds throughout for  $P = (1, 0)$ ,  $Q = (0, 1)$ .

The inequalities (6.4) can be used to derive a strengthening of (5.1). Indeed, one finds that

$$(6.5) \quad C(P, Q) \geq \sum_{\nu=1}^{\infty} \frac{2^{-(\nu-1)}}{2\nu(2\nu-1)} \Delta(P, Q)^\nu,$$

again with best constants as coefficients (for this, observe that the discussion of equality in (6.5) is similar to the one in Theorem 5).

One may also derive a variant of (23) of the form  $D \geq \sum c_\nu \cdot \Delta^\nu$ , but a straight forward derivation gives constants far from the best possible ones. The most interesting inequality of this type is of the form  $D \geq c \cdot \Delta$ . Let  $c^{max}$  denote the largest such constant. From (4.3) and from (2.11) and Pinsker's inequality it follows that  $\frac{1}{2} \leq c^{max} \leq 1$ . In fact, we have  $0.8033 \leq c^{max} \leq 0.8961$ . The first inequality follows as one can show, using our identities and inequalities, that  $D \geq \frac{1}{2} \sum_1^\infty \frac{1}{2^{\nu-1}} \cdot \Delta$ , and the second follows from consideration of probability measures of the form  $P = (\varepsilon, 1 - \varepsilon)$ ,  $Q = (\rho\varepsilon, 1 - \rho\varepsilon)$  for small values of  $\varepsilon$  (and the appropriate value of  $\rho$ ,  $\rho \approx 4.48$ ). Details are left to the interested reader.<sup>2</sup>

The "star-construction" used for  $\Delta$  may also be introduced for  $C$  and  $\Delta_\nu$ :

$$(6.6) \quad \begin{aligned} C^*(P\|Q) &= \sum_0^\infty 2^k C(M_k, Q), \\ \Delta_\nu^*(P\|Q) &= \sum_0^\infty 2^k \Delta_\nu(M_k, Q) \end{aligned}$$

with  $M_0, M_1, \dots$  denoting the successive midpoints in direction  $Q$ . Then, from (3.1) and (4.2), we get:

$$(6.7) \quad D(P\|Q) = C^*(P\|Q) = \sum_1^\infty \frac{1}{2\nu(2\nu-1)} \Delta_\nu^*(P\|Q).$$

The two inequalities of (3.2) are best possible in a natural sense: By considering  $P = (1, 0)$  and  $Q = (0, 1)$  it follows that  $\log 2$  is the smallest possible constant in an inequality like  $C \leq \alpha \cdot \Delta$ . And, regarding the other inequality, we may either appeal to Pinsker's inequality or, more directly, we may consider the two functions  $f(\varepsilon) = C(P, Q_\varepsilon)$  and  $g(\varepsilon) = \Delta(P, Q_\varepsilon)$  with  $P = (\frac{1}{2}, \frac{1}{2})$  and  $Q_\varepsilon = (\frac{1}{2} + \varepsilon, \frac{1}{2} - \varepsilon)$  and note that  $f(0) = f'(0) = g(0) = g'(0) = 0$  and that  $f''(0) = 2$ ,  $g''(0) = 4$ . Hence  $C(P, Q_\varepsilon)/\Delta(P, Q_\varepsilon) \rightarrow \frac{1}{2}$  as  $\varepsilon \rightarrow 0$ , and it follows that  $\beta = \frac{1}{2}$  is the largest possible constant in an inequality of the form  $\beta \cdot \Delta \leq C$ .

For the left hand inequality of (4.3), i.e.  $\frac{1}{2}\Delta^* \leq D$ ,  $\frac{1}{2}$  is the largest possible constant since any larger constant would give rise to a strengthening of Pinsker's inequality beyond what is possible, cf. (4.5).

All measures of discrepancy considered in this correspondence are particular instances of Csiszár  $f$ -divergences. Recall, cf. Csiszár [5, 6, 7], that for a convex function  $f : [0, \infty[ \rightarrow \mathbb{R}$ , the Csiszár  $f$ -divergence between  $P$  and  $Q$  is defined by

$$(6.8) \quad I_f(P, Q) = \sum q_i f\left(\frac{p_i}{q_i}\right)$$

(we need only instances with  $f(0)$  finite; below we typically normalize so that  $f(0) = 1$ ). The family of functions  $(f_s)_{s \geq 1}$  with  $f_s(u) = |u - 1|^s \cdot (u + 1)^{-(s-1)}$  gives rise to variational distance  $V$  ( $s = 1$ ), triangular discrimination  $\Delta$  ( $s = 2$ ) and triangular discrimination of order  $\nu$ ,  $\Delta_\nu$  ( $s = 2\nu$ ). And if we take  $f(u) = (u - 1)^2 / (u + 2^{\nu+1} - 1)$ , we are led to the map  $(P, Q) \rightsquigarrow \Delta(M_\nu, Q)$  where, as usual,  $M_\nu = 2^{-\nu}P + (1 - 2^{-\nu})Q$ .

<sup>2</sup> the sum  $\sum \frac{1}{2^{\nu-1}}$  is of interest in analytic number theory. Indeed,  $\sum \frac{x^n}{1-x^n} = \sum d_n x^n$ , where  $d_n$  is the number of divisors of  $n$ . The interested reader may consult Borwein [4] from which we see that, typically, these sums are irrational for rational  $x$ , in particular for  $x = \frac{1}{2}$ .

Among the most popular choices we mention  $f(u) = \frac{1}{2}(\sqrt{u} - 1)^2$  which gives rise to the *Hellinger discrimination*  $h^2$ :

$$h^2(P, Q) = \frac{1}{2} \sum (\sqrt{p_i} - \sqrt{q_i})^2.$$

The square-root  $h = \sqrt{h^2}$  is a true distance.

The basic relations between  $V$ ,  $\Delta$  and  $h^2$  are the following three:  $\frac{1}{2}V^2 \leq \Delta \leq V$ , cf. (2.11), then the relation  $2h^2 \leq \Delta \leq 4h^2$  (derived e.g. by comparing the corresponding  $f$ -functions, cf. also LeCam [19] and Dacunha–Castelle [10]) and lastly, the relation  $\frac{1}{8}V^2 \leq h^2 \leq \frac{1}{2}V$  (follows from the two first). The occurring coefficients are best possible as may be seen by considering the case  $P = (1, 0)$ ,  $Q = (0, 1)$  or the case  $P = (\frac{1}{2} - \varepsilon, \frac{1}{2} + \varepsilon)$ ,  $Q = (\frac{1}{2}, \frac{1}{2})$  for small values of  $\varepsilon$ .

Kraft [17] improved part of this by pointing out that  $\frac{1}{8}V^2 \leq h^2(1 - \frac{1}{2}h^2)$  (follows by applying Cauchy–Schwartz' inequality). The inequalities between  $h^2$  and  $\Delta$  were also noted by LeCam [19].

In Dacunha–Castelle [10] we find the important inequality  $D \geq -2 \log(1 - h^2)$  (follows by Jensen's inequality), in particular,  $D \geq 2h^2$ . This is best possible as an inequality  $D \geq ch^2$  cannot hold for any  $c^2$  (consider a large  $\rho$  and  $P = (\varepsilon, 1 - \varepsilon)$ ,  $Q = (\rho\varepsilon, 1 - \rho\varepsilon)$  for small  $\varepsilon$ 's).

The above comments indicate that  $\Delta$  and  $h^2$  have similar properties. It may be true that  $\Delta$  is closer to  $D$  in some sense (as our results have shown), but, on the other hand,  $h^2$  has nice structural properties (discussed elsewhere) which are not shared by  $\Delta$ . In conclusion then, triangular discrimination cannot offer a replacement of Hellinger distance, but does appear to provide a convenient supplement.

Generalizations of identities and inequalities here presented to a countably infinite alphabet or to distributions defined with reference to general measure spaces are, basically, straight forward and a matter of routine (this is important, e.g. if  $\Delta$  will replace  $h^2$  for certain investigations in statistics). However, one issue deserves special mention, viz. regarding the general validity of (4.2) under the condition  $D(P\|Q) < \infty$ . So assume that  $P$  and  $Q$  are probability measures on a general measurable space and that  $D(P\|Q) < \infty$ . In particular then,  $P$  is absolutely continuous w.r.t.  $Q$ . We put  $M_\varepsilon = \varepsilon P + (1 - \varepsilon)Q$ ;  $0 < \varepsilon < 1$ , and shall prove that

$$(6.9) \quad \lim_{\varepsilon \rightarrow 0} \frac{1}{\varepsilon} D(M_\varepsilon \| Q) = 0.$$

Put  $\varphi = \frac{dP}{dQ}$ . A simple computation shows that

$$\frac{1}{\varepsilon} D(M_\varepsilon \| Q) = \frac{1}{\varepsilon} \int \log \frac{dM_\varepsilon}{dQ} dM_\varepsilon = A + B + C$$

with

$$A = \frac{1}{\varepsilon} \log(1 - \varepsilon), \quad B = \int \log\left(1 + \frac{\varepsilon}{1 - \varepsilon} \varphi\right) dP,$$

$$C = \int \frac{1 - \varepsilon}{\varepsilon} \log\left(1 + \frac{\varepsilon}{1 - \varepsilon} \varphi\right) dQ.$$

Clearly,  $A \rightarrow -1$  as  $\varepsilon \rightarrow 0$ . Concerning  $B$ , call the integrand  $f_\varepsilon$ , and note that  $f_\varepsilon$  decreases pointwise to 0 as  $\varepsilon$  decreases to 0 and that  $f_{\frac{1}{2}} = \log(1 + \varphi)$  is  $P$ -integrable (as  $\log(1 + x) \leq 2 \log x$  for  $x \geq 2$  and as  $D(P\|Q) < \infty$ ). It follows that  $B \rightarrow 0$  as  $\varepsilon \rightarrow 0$ . Finally, concerning  $C$ , we observe that  $C$  is of the form  $\int g_\varepsilon dQ$  with all  $g_\varepsilon \leq \varphi$  and that  $g_\varepsilon \rightarrow \varphi$  hence, again by dominated convergence,  $C \rightarrow 1$  as  $\varepsilon \rightarrow 0$ . All in all,  $A + B + C \rightarrow 0$  as  $\varepsilon \rightarrow 0$  and (6.9) follows.

#### ACKNOWLEDGEMENTS

The research presented arose in discussions with Boris Ryabko and I am much indebted to his suggestions. In particular, Theorem 1 is due to Ryabko. Just before finishing the manuscript, the author had the opportunity to discuss the findings at the workshop "Information theory, statistical learning and pattern recognition", CIRM, Marseille, december 1998. The discussions at the workshop with Andrew Barron, Lucien Birgé, Amir Dembo, Katalin Marton and Ofer Zeitouni resulted in improvements of the discussion. In particular, the proof of (6.9) was suggested by Ofer Zeitouni.

## REFERENCES

- [1] R. Ahlswede, and I. Wegener, *Search problems*, Chichester: Wiley, 1987. (German original by Teubner, 1979).
- [2] L. Birgé, *Approximation dans les espaces métriques et théorie de l'estimation*, Z. Wahrscheinlichkeitstheorie verw.Geb., vol. 65, pp. 181–237, 1983.
- [3] L. Birgé, and P. Massart, *Rates of convergence for minimum contrast estimator*, Probab. Theory Relat. Fields, vol. 97, pp. 113–150, 1993.
- [4] P. B. Borwein, *On the irrationality of certain series*, Math. Proc. Camb. Phil. Soc. Probab., vol. 112, pp. 141–146, 1992.
- [5] I. Csiszár, *A note on Jensen's inequality*, Studia Sci. Math. Hungar., vol. 1, pp. 185–188, 1966.
- [6] I. Csiszár, *Information-type measures of difference of probability distributions and indirect observation*, Studia Sci. Math. Hungar., vol. 2, pp. 299–318, 1967.
- [7] I. Csiszár, *On topological properties of  $f$ -divergences*, Studia Sci. Math. Hungar., vol. 2, pp. 329–339, 1967.
- [8] I. Csiszár, and J. Körner, *Information Theory: Coding Theorems for discrete memoryless systems*, New York: Academic, 1981.
- [9] I. Csiszár, *Sanov property, generalized  $I$ -projection and a conditional limit theorem*, Ann. Prob. vol. 12, pp. 768–793, 1984.
- [10] D. Dacunha-Castelle, *Ecole d'Ete de Probabilités de Saint-Flour VII-1977*, Berlin, Heidelberg, New York: Springer, 1978.
- [11] A. Dembo, *Information inequalities and concentration of measure*, Ann. Prob., vol. 25, pp. 927–939, 1997.
- [12] D. Feldman, and F. Österreicher, *A note on  $f$ -divergences*, Studia Sci. Math. Hungar., vol. 24, pp. 191–200, 1989.
- [13] P. Kafka, F. Österreicher, and I. Vincze, *On powers of  $f$ -divergences defining a distance*, Studia Sci. Math. Hungar., vol. 26, pp. 415–422, 1991.
- [14] N. S. Kambo, and S. Kotz, *On exponential bounds for binomial probabilities*, Ann. Inst. Stat. Math., vol. 18, pp. 277–287, 1966.
- [15] O. Krafft, *A note on exponential bounds for binomial probabilities*, Ann. Inst. Stat. Math., vol. 21, pp. 219–220, 1969.
- [16] O. Krafft, and N. Schmitz, *A note on Hoeffding's inequality*, J. Amer. Statist. Assoc., vol. 64, pp. 907–912, 1969.
- [17] C. Kraft, *Some conditions for consistency and uniform consistency of statistical procedures*, Univ. of California Publ. in Statistics, vol. 1, pp. 125–142, 1955.
- [18] S. Kullback, *Information Theory and Statistics*, New York: Wiley, 1959.
- [19] L. LeCam, *Asymptotic Methods in Statistical Decision Theory*, New York: Springer, 1986.
- [20] K. Marton, *Bounding  $\bar{d}$ -distance by informational divergence: A method to prove measure concentration*, Ann. Prob., vol. 24, pp. 857–866, 1966.
- [21] K. Marton, *A simple proof of the blowing-up lemma*, IEEE Trans. Inform. Theory, vol. 32, pp. 445–446, 1986.
- [22] F. Österreicher, and I. Vajda, *Statistical information and discrimination*, IEEE Trans. Inform. Theory, vol. 39, pp. 1036–1039, 1993.
- [23] M.S. Pinsker, *Information and information stability of random variables and processes*, San Francisco: Holden-Day, 1964 (Russian original 1960).
- [24] B.Ya. Ryabko, *Comments on "A source matching approach to finding minimax codes"*, IEEE Trans. Inform. Theory, vol. 27, pp. 780–781, 1981. (Including also the ensuing Editor's Note).
- [25] B.Ya. Ryabko, and F. Topsøe, *On asymptotically optimal methods of prediction and adaptive coding*, manuscript submitted for publication.
- [26] F. Topsøe, *Information theoretical optimization techniques*, Kybernetika, vol. 15, pp. 8–27, 1979.

- [27] I. Vajda, *Note on discrimination information and variation*, IEEE Trans. Inform. Theory, vol. 16, pp. 771–773, 1970.
- [28] I. Vincze, *On the concept and measure of information contained in an observation*, Contributions to probability (ed. by J. Gani and V.K. Rohatgi), pp. 207–214, New York: Academic, 1981.
- [29] Y. Yang, and A. Barron, *Information-theoretic determination of minimax rates of convergence*, under publication.

DEPARTMENT OF MATHEMATICS  
UNIVERSITY OF COPENHAGEN

