# MEASURES OF INFORMATION IN PROBABILITY AND STATISTICS

Adolfo M. de Guzman
Associate Professor
Statistical Center
University of the Philippines

### Abstract

Applications of information measures in analysis, probability theory and statistics are pointed out. In testing the independence of random vectors under normality, the maximum likelihood criterion is shown to be equivalent to an entropy-based test. Two large sample tests based on entropic loss of information is presented.

Keywords and phrases: information measures, Fisher information, entropy, entropy-based tests.

## 1. INTRODUCTION

The idea of using measures of information to prove limit theorems is due to Linnik (1959). He gave an information-theoretic proof of the Central Limit Theorem on Lindeberg conditions. Renyi (1960) provides a measure-theoretic proof of a limit theorem for Markov chains. Only recently de Guzman (1989a) showed that Hadamard's Inequality follows almost trivially from the nonnegativity of the Kullback-Leibler information quantity. Jefferson, May and Ravi (1989) suggest the use of entropy to the scaling of some ordinal categorical data. Entropy-based goodness of fit tests have been developed by Vasicek (1976), Dudewicz and Van der Meulen (1981) and Gokhale (1983). In statistical pattern recognition some rules for feature evaluation are derived from information measures (Ben-Bassat, 1982). Srivastava (1973) attempts to extract the "intrinsic" dimensionality of a multivariate data set by exploiting Shannon's information function. In the last fifteen years, Akaike's Information criterion (AIC) which is based on Kullback-Leibler entropy has

found successful application in statistical model evaluation problems (Bozdogan, 1987). Hall (1987) shows that in kernel density estimation, if the kernel is chosen appropriately, likelihood cross-validation does result in asymptotic minimization of Kullback-Leibler loss (information measure). Malvestuto (1989) provides a computational method for obtaining the maximum entropy extension of given discrete probability distributions.

In this paper we shall introduce Fisher's information, give a measure-theoretic proof of Hadamard's inequality, and propose an entropy-based statistic for testing independence of two random vectors.

## 2. FISHER INFORMATION IN EQUIVALENCE/SINGULARITY DICHOTOMIES

Let $(\Omega, F, P)$ be an experiment i.e. $(\Omega, P)$ is a measurable space and $P$ a class of probability measure on $(\Omega, F)$. We say that an equivalence/singularity dichotomy holds if P, Q $\in P$ implies or P $\equiv$ Q or P $\perp$ Q.

Kakutani (1946) gave the following interesting example of dichotomy. Let $(\Omega, P) = (\prod_{i=1}^{\infty} \Omega_i, \sigma(\prod_{i=1}^{\infty} F_i))$ where $(\Omega_i, F_i)$, is a sequence of measurable spaces satisfying Kolmogorov consistency conditions. Let $\{Q_i\}$. be a sequence of measures, where $Q_i$ is on $(\Omega_i, F_i)$ for all i. If $P = \{\prod_{i=1}^{\infty} P_i : P_i \equiv Q_i$ for all i,$\}$ Kakutani showed that if $P = \prod_{i=1}^{\infty} P_i$ and $\tilde{P} = \prod_{i=1}^{\infty} \tilde{P}_i \in P$ then $P \equiv \tilde{P}$ or $P \perp \tilde{P}$ with the former holding iff $\sum_{i=1}^{\infty} H^2(P_i, \tilde{P}_i) < \infty$.

Here H is the Hellinger distance and is defined by

$$2 H^2(P_i, \tilde{P}_i) = \int \sqrt{(f_i - \sqrt{\tilde{f}_i})} \, dv_i$$

where $v_i = P_i + \tilde{P}_i$, $f_i \in dP_i/dv_i$ and $\tilde{f}_i \in d\tilde{P}_i/dv_i$.

In the case where $\Omega = R^\infty$, $F = B(R^\infty)$ and $P$ is the set of all

Gaussian measures, let $\overset{\infty}{\underset{i=1}{\pi}} N(0,1)$ and $\tilde{P} = \overset{\infty}{\pi} N(\mu_i, \sigma_i^2)$.

Feldman (1958) and Hajek (1958) showed that $P \equiv \tilde{P}$ if and

only if

$$\sum_{i=1}^{\infty} (\mu_i^2 + (1-\sigma_i)^2) < \infty .$$

In case $\sigma_i^2 = 1$, $P \equiv \tilde{P}$ if and only if

$\mu_i \in l^2$, the space of square summable sequences.

What other probability measures on $B(R)$ besides $N(0,1)$ satisfy
the property? Shepp (1965) gave the following answer: if $P$ is
a probability measure in $(R, B(R)$ and $P_t$ is the translate of

$P$ both then $\overset{\infty}{\underset{i=1}{\pi}} P \equiv \overset{\infty}{\underset{i=1}{\pi}} P_t$ for all $t_i \in l^2$ if and only if $P \equiv \lambda$

(Lebesque measure) and $P$ has finite information measure that
is, there exists a locally absolutely continous density f such
that

$$\int \frac{(f')^2}{f} \, d\lambda < \infty$$

Thus in a translation invariant experiment finiteness of Fisher
information as defined above is a necessary and sufficient
condition for an $l^2$-type dichotomy

## 3. AN INFORMATION THEORETIC PROOF OF HADAMARD'S INEQUALITY

**Theorem** (Hadamard's inequality) If A is an nxp real matrix of
rank p, then

$$|A^TA| \leq \prod_{i=1}^{p} \sum_{j=1}^{n} a_{ij}^2$$

where $A = (a_{ij})$.

**Proof.** Consider a p-variate random vector $\mathbf{X} = (X_1,\ldots,X_p)'$ which has p-variate normal density with dispersion matrix $\Sigma = \mathbf{A}^T\mathbf{A}$. The entropies $H(\mathbf{X})$ and $H(X_1)$ $i = 1,..,p$ are given by

$$H(\mathbf{X}) = \frac{p}{2} + \frac{p}{2}\log 2\pi + \frac{1}{2}\log |\Sigma| \;,\; H(X_1) = \frac{1}{2} + \frac{1}{2}\log 2\pi$$

$$+ \frac{1}{2}\log \sigma_1^2 \quad \text{where} \quad \sigma_1^2 = \sum_{j=1}^{n} a_{1j}^2 \quad i=1,\ldots,p$$

so that

$$\sum_{i=1}^{p} H(X_1) = \frac{p}{2} + \frac{p}{2}\log 2\pi + \frac{1}{2}\log \left[ \prod_{i=1}^{p} \sum_{j=1}^{n} a_{1j}^2 \right]$$

The Kullback-Leibler quantity

$$H(\mathbf{X}) - \sum_{i=1}^{p} H(X_1)$$

is nonnegative and the log function is monotonic yield

$$|\Sigma| = |\mathbf{A}^T\mathbf{A}| \leq \prod_{i=1}^{p} \sum_{j=1}^{n} a_{1j}^2 \;,$$

which is Hadamard's inequality. ∎

## 4. ENTROPY-BASED TESTS

Vasicek (1979) introduced a test on the composite hypothesis of normality based on sample estimate of entropy. The test was shown to be consistent against all alternatives without a singular continuous part. Its asymptotic normality was exhibited by Dudewicz and van der Meulen under the hypothesis of uniformity and also under a special class of alternative densities. A general form of a goodness-of-fit test statistic for families of maximum entropy distributions was given by Gokhale (1983). The proposed test is shown consistent against an appropriate class of alternatives and simulation and Monte Carlo results show favorable comparison with other goodness-of-fit tests.

To test the independence of two random vectors, de Guzman (1989b) showed that a test based on entropic loss of information is equivalent to the test based on the likelihood ratio criterion under normality assumptions.

Let $X\alpha$, $\alpha = 1,\ldots,N$ be a random sample from $n(\mu,\Sigma)$ and suppose $X' = (X^{(1)'} \ X^{(2)'})'$ and let $\mu$ and $\Sigma$ be partitioned correspondingly as

$$\mu = \begin{bmatrix} \mu^{(1)} \\ \mu^{(2)} \end{bmatrix} \qquad \Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$$

Also, let

$$S = \Sigma (X_\alpha - X)(X\alpha - X)'/N \quad \text{and} \quad A = NS \text{ with corresponding}$$

partitions $S = \begin{bmatrix} S_{11} & S_{12} \\ S_{21} & S_{22} \end{bmatrix}$ and $A = \begin{bmatrix} A_{11} & a_{12} \\ A_{21} & A_{22} \end{bmatrix}$

Consider the null hypothesis $H_0$: $X^{(1)}$ and $X^{(2)}$ are independent. Loss of information for discarding $X^{(2)}$ is given by

$$L(X) = H(X) - H(X^{(1)}) = c_{p-m} + \frac{1}{2} \log \frac{|\Sigma|}{|\Sigma_{11}|}$$

under normality assumption.

The null hypothesis is equivalent to

$$H_0: \Sigma_{12} = \Sigma_{21} = 0$$

so that $|\Sigma| = |\Sigma_{11}||\Sigma_{22}|$. Hence $L(X|H_0) = c_{p-m} + \frac{1}{2} \log |\Sigma_{22}|$.

Hence

$$W = L(X) - L(X|H_0) = \frac{1}{2} \log \frac{|\Sigma|}{|\Sigma_{11}||\Sigma_{22}|}$$

We therefore have the following theorem.

**THEOREM:** For testing the null hypothesis Ho: $X^{(1)}$ and $X^{(2)}$ are independent, the entropy based test statistic

$$W = \frac{1}{2} \log \frac{|A|}{|A_{11}||A_{22}|}$$

is equivalent to the likelihood criterion when **X** is multivariate normal.

We now remove the assumption that **X** is multivariate normal and construct a nonparametric large sample test based on entropic loss of information. Let **X** have density f and suppose $f_n$ is an estimator of f based on the observations $X_1, X_2, \ldots, X_n$.

Let $X = \begin{bmatrix} Y \\ Z \end{bmatrix}$ with g, h as the densities of Y and Z

Then $L(X) = -\int f(x) \log f(x) - \int g(y) \log g(y)$

and under the null hypotheses Ho: Y and X are independent

$$L(X|Ho) = -\int g(y)(z) \log[g(y)h(z)] dy dz - \int |g(y) \log(y) dx$$

$$= -\int g(y) \log g(y) dy - \int h(z) \log h(z) dz$$

$$+ \int g(y) \log g(y) dy = -\int h(z) \log h(z) dz,$$

so that

$$L(X)-L(X|Ho) = - \int f(x)\log f(x)\ dx + \int g(y)\ \log\ g(y)\ dy$$

$$+ \int h(z)|\log h(z)dz = H(X) - H(Y) - H(Z).$$

Let $f_n$, $g_n$ and $h_n$ be density estimates of f,g and h respectively and let

$$W = \hat{H}(X) - \hat{H}(Y) - \hat{H}(Z) = - \int f_n(x)\ \log\ f_n(x)\ dx$$

$$+ \int g_n(y)\ \log g_n(y)dy + \int h_n(z)\ \log\ h_n(z)dz)$$

If the null hypothesis Ho is true, W should be small. If the density estimates $f_n$, $g_n$ and $h_n$ are suitably chosen, W can be shown to be normal. Consider $H(X)/\{H(Z) + H(Y)\}$. Under Ho, $H(X)/H(Z) + H(Y)\} = 1$. Therefore the problem reduces to testing the null hypothesis Ho: $H(X)/\{H(Z) + H(Y)\} = 1$ against the alternative A: $H(X)/\{H(Z) + H(Y)\} < 1$. This suggest the test statistic

$$V = \hat{H}(X)/\{\hat{H}(Z) + \hat{H}(Y)\}.$$

The asymptotic distribution of V still has to be worked out.

# References

Ben-Bassat M. (1982). Use of distance measures, information measures and error bounds in feature evaluation. In: P.R. Krisnaiah and L.N. Kanal. *Handbook of Statistics* 2 North Holland Publishing Company 773-791.

Bozdogan, H. (1987). Model selection and Akaike's Information Criterion (AIC): the general theory and its analytical extension. *Psychemetrika* 50 3 345-370.

de Guzman, A. M. (1989a). Dimension, dimension reduction and loss of information. Dissertation submitted to the University of the Philippines.

_____, (1986b). An information theoretic proof of Hadamards' inequality. *Matimyas Matematika* (to appear).

Dudewicz, E. J. and van der Meulen (1981). Entropy based tests of uniformity. Journal of the ASA 76, 967-974.

Gokhale, D.V. (1983). On entropy based goodness-of-fit Test. *Computational Statistics and Data Analysis* 1 157-165.

Hall, P. (1987). On Kullback-Leibler loss and density estimation. *The Annals of Statistics.* 15 4 1491-1519.

Jefferson, R. R., May, J.H. and Ravi, N. (1989). An entropy approach to the scaling of ordinal categorical data. Psychometrika 54 2 203-215.

Linnik, Yu V. (1959). An information theoretical proof of the central limit theorem on Lindeberg conditions. *Teor. Veroyatnost. i Primenen*, 4 311-321. (In Russian).

Malvestuto, F.M. (1989). Computing the maximum-entropy extension of a given discrete probability distribution. *Computational Statistics and Data Analyses.* 83 299-311.

Neymann, J. (1960). *Proceedings of the fourth Berkeley symposium on mathematical statistics and probability* 1, 547-561.

Renyi, A. (1970). *Foundations of probability.* Holden-Day Inc. San Francisco.

Shore, J.E. et al (1984). Axiomatic derivation of the principle of maximum entropy and the principle of minimum cross entropy. *IEEE Information Theory* 26 26-37.

Srivastava, J.N. (1973). An information function approach to dimensionality analysis and curved maniforld clustering. In: *Multivariate Analysis III*. Academic Press, Inc.

Thelen, B.J. (1989). Fisher information and in equivalence/contiguity. *The Annals of Probability* **17** 4 1664-1690.

Vasicek, O. (1976). A test for normality based on sample entropy. *Journal of the Royal Statistical Society*. Series B **38** 54-59.