

Article

A Kernel-Based Calculation of Information on a Metric Space

R. Joshua Tobin ¹ and Conor J. Houghton ^{2,*}

¹ School of Mathematics, Trinity College Dublin, Dublin 2, Ireland

² Department of Computer Science, University of Bristol, Merchant Venturers Building, Woodland Road, Bristol BS8 1UB, England

* Author to whom correspondence should be addressed; conor.houghton@bristol.ac.uk, +44 (0) 117 954 5140

Version October 13, 2013 submitted to *Entropy*. Typeset by *LaTeX* using class file *mdpi.cls*

1 **Abstract:** Kernel density estimation is a technique for approximating probability
2 distributions. Here, it is applied to the calculation of mutual information on a metric space.
3 This is motivated by the problem in neuroscience of calculating the mutual information
4 between stimuli and spiking responses; the space of these responses is a metric space. It
5 is shown that kernel density estimation on a metric space resembles the k -nearest-neighbor
6 approach. This approach is applied to a toy data-set designed to mimic electrophysiological
7 data.

8 **Keywords:** mutual information, neuroscience, electrophysiology, metric spaces, kernel
9 density estimation

10 1. Introduction

11 This paper is concerned with the calculation of mutual information for spike trains using the data
12 that are available in a typical *in vivo* electrophysiology experiment in the sensory system. It uses a
13 kernel-based estimation of probability distributions.

14 In particular, this paper is concerned with computing the mutual information $I(R; S)$ between
15 two random variables, R and S . The motivating neuroscience example is a typical sensory pathway
16 electrophysiology experiment in which the corpus of sensory stimuli are presented over multiple trials,
17 so there is a set of recorded responses for each of a number of stimuli. The stimuli are drawn from a
18 discrete space, the corpus, but the responses are spike trains. The space of spike trains is peculiar; locally
19 it is like a smooth manifold with the spike times behaving like coordinates, but globally it is foliated into

subspaces, each with a different number of spikes. The space of spike trains does, however, have a metric. As such, S takes values in a discrete set, \mathcal{S} , and models the stimulus, and R takes values in a metric space, \mathcal{R} , and models the response.

\mathcal{R} is not a discrete space and so, to calculate the mutual information between S and R , it is necessary to either discretize \mathcal{R} or to use differential mutual information. In the application of information theory to electrophysiological data, it is common to take the former route and discretize the data. Here the latter alternative is chosen and the differential mutual information is estimated.

The mutual information between two random variables R and S is a measure of the average amount of information that is gained about S from knowing the value of R . With S a discrete random variable taking values in \mathcal{S} and R a continuous random variable, the mutual information is

$$I(R; S) = \sum_{s \in \mathcal{S}} \int_{\mathcal{R}} p(r, s) \log_2 \frac{p(r, s)}{p(r)p(s)} dr \quad (1)$$

where dr is the measure on \mathcal{R} : computing the differential mutual information between R and S requires integration over \mathcal{R} . Integration requires a measure, and when there are coordinates on a space, it is common to use the integration measure derived from these coordinates.

The space of spike trains has no system of coordinates and so there is no coordinate-based measure. This does not mean that the space has no measure, as a sample space it has an intrinsic measure corresponding to the probability distribution; thus, there is a measure, just not one derived from coordinates. The probability of an event occurring in a region of sample space gives a volume for that region. In other words, the volume of a region \mathcal{D} can be identified with $P(\mathbf{x} \in \mathcal{D})$. This is the measure that will be used throughout this paper; it does not rely on coordinates and so can be applied to the case of interest here.

Of course, in practice, the probability density is not usually known on the space of spike trains, but $P(\mathbf{x} \in \mathcal{D})$ can be estimated from the set of experimental data. A Monte-Carlo-like approach is used: the volume of a region is estimated by counting the fraction of data points that lie within it

$$\text{vol}(\mathcal{D}) = P(\mathbf{x} \in \mathcal{D}) \approx \frac{\text{number of data points in } \mathcal{D}}{\text{total number of points}}. \quad (2)$$

This is exploited in this paper to estimate the volume of square kernels, making it possible to estimate conditional probabilities using kernel density estimation.

The classical approach to the problem of estimating $I(R; S)$ is to map the spike trains to binary words using temporal binning [1,2] giving a histogram approximation for $p(r, s)$. This approach is very successful, particularly when supplemented with a strategically chosen prior distribution for the underlying probability distribution of words [3,4]. This is sometimes called the plug-in method and that term is adopted here. One advantage of the plug-in method is that the mutual information it calculates is correct in the limit: in the limit of an infinite amount of data and an infinitesimal bin size it gives the differential mutual information.

Nonetheless, it is interesting to consider other approaches, and in this spirit, an alternate approach is presented here. This new method exploits the inherent metric structure of the space of spike trains, it is very natural and gives an easily implemented algorithm which is accurate on comparatively small data sets.

56 2. Methods

57 This section describes the proposed method for calculating mutual information. Roughly, the
58 conditional probability is approximated using kernel density estimation and, by using the unconditioned
59 probability distribution as a measure, integration is approximated by the Monte-Carlo method of
60 summing over data points.

61 Since this is a kernel-based approach, a review of kernel density estimation is given in section 2.1.
62 This also serves to establish notation. The two key steps used to derive the kernel-based estimate are a
63 change of measure and a Monte-Carlo estimate. The change of measure, described in section 2.2, permits
64 the estimation of probabilities by a simple Monte-Carlo method. The new measure also simplifies the
65 calculation of $I(R; S)$, resulting in a formula involving a single conditional distribution. This conditional
66 distribution is estimated using a Monte-Carlo estimate in section 2.3.

67 2.1. Kernel Density Estimation

68 The non-parametric kernel density estimation (KDE) method [6–8] is an approach to estimating
69 probability densities. In KDE, a probability density is estimated by filtering the data with a kernel. This
70 kernel is normalized with integral one and is usually symmetric and localized. For an n -dimensional
71 distribution with outcome vectors $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$ and a kernel $k(\mathbf{x})$ the estimated distribution is
72 usually written

$$\tilde{p}(\mathbf{x}) = \frac{1}{m} \sum_i k(\mathbf{x} - \mathbf{x}_i) \quad (3)$$

73 where, because the argument is $\mathbf{x} - \mathbf{x}_i$, there is a copy of the kernel centered at each data point. In
74 fact, this relies on the vector-space structure of n -dimensional space; in the application considered here
75 a more general notation is required, with $k(\mathbf{x}; \mathbf{y})$ denoting the value at \mathbf{x} of the kernel when it is centered
76 on \mathbf{y} . In this situation the estimate becomes

$$\tilde{p}(\mathbf{x}) = \frac{1}{m} \sum_i k(\mathbf{x}; \mathbf{x}_i). \quad (4)$$

77 The square kernel is a common choice: for a vector space this is

$$k(\mathbf{x}; \mathbf{y}) = \begin{cases} \frac{1}{V} & \|\mathbf{x} - \mathbf{y}\| < 1 \\ 0 & \text{otherwise} \end{cases} \quad (5)$$

78 where V is chosen so that the kernel integrates to one. The kernel is usually scaled to give it a bandwidth:

$$k(\mathbf{x}; \mathbf{y}, h) = \begin{cases} \frac{1}{h^V} & \|\mathbf{x} - \mathbf{y}\| < h \\ 0 & \text{otherwise} \end{cases}. \quad (6)$$

79 This bandwidth h specifies the amount of smoothing. The square kernel is the most straight-forward
80 choice of kernel mathematically and so in the construction presented here a square kernel is used.

81 In the case which will be of interest here, where \mathbf{x} and \mathbf{y} are not elements of a vector space, the
82 condition $\|\mathbf{x} - \mathbf{y}\| < h$ must be replaced by $d(\mathbf{x}, \mathbf{y}) < h$ where $d(\mathbf{x}, \mathbf{y})$ is a metric measuring the
83 distance between \mathbf{x} and \mathbf{y} . Calculating the normalization factor, V , is more difficult since this requires
84 integration. This problem is discussed in the next subsection.

85 2.2. *Change of Measure*

86 Calculating the differential mutual information using KDE requires integration, both the integration
87 required by the definition of the mutual information, and the integration needed to normalize the kernel.
88 As outlined above, these integrals are estimated using a Monte-Carlo approach; this relies on a change
89 of measure which is described in this section.

90 For definiteness, the notation used here is based on the intended application to spike trains. The
91 number of stimuli is n_s , and each stimulus is presented for n_t trials. The total number of responses
92 n_r is then $n_r = n_s n_t$. Points in the set of stimuli are called s and in the response space, r ; the actual
93 data points are indexed, r_i , and (r_i, s_i) is a response-stimulus pair. As above, the random variables for
94 stimulus and response are S and R whereas the set of stimuli and the space of responses are denoted by
95 a calligraphic \mathcal{S} and \mathcal{R} respectively. It is intended that when the method is applied the responses, $r \in \mathcal{R}$,
96 will be spike trains.

97 The goal is to calculate the mutual information between the stimulus and the response. Using the
98 Bayes theorem, this is

$$I(R; S) = \sum_{s \in \mathcal{S}} \int_{\mathcal{R}} p(r, s) \log_2 \frac{p(r|s)}{p(r)} dr. \quad (7)$$

99 Unlike the differential entropy, the differential mutual information is invariant under the choice of
100 measure. Typically, differential information theory is applied to examples where there are coordinates
101 (x_1, x_2, \dots, x_n) on the response space and the measure is given by $dr = dx_1 dx_2 \dots dx_n$. However, here
102 it is intended to use the measure provided by the probability distribution $p(r)$. Thus, for a region $\mathcal{D} \subset \mathcal{R}$
103 the change of measure is

$$\text{vol}(\mathcal{D}) = \int_{\mathcal{D}} p(r) dr = \int_{\mathcal{D}} d\beta \quad (8)$$

104 SO

$$d\beta = p(r) dr. \quad (9)$$

105 The new probability density relative to the new measure, $p_\beta(r)$, is now one:

$$p_\beta(r) = \frac{p(r)}{d\beta/dr} = 1. \quad (10)$$

106 Furthermore, since $p(r|s)$ and $p(r)$ are both densities, $p(r|s)/p(r)$ is invariant under a change of measure
107 and

$$I(R; S) = \sum_s \int_{\mathcal{R}} p_\beta(r, s) \log_2 \frac{p_\beta(r|s)}{p_\beta(r)} d\beta = \sum_s \int_{\mathcal{R}} p_\beta(r, s) \log_2 p_\beta(r|s) d\beta \quad (11)$$

108 where, again, $p_\beta(r, s)$ and $p_\beta(r|s)$ are the values of the densities $p(r, s)$ and $p(r|s)$ after the change of
109 measure.

110 The expected value of any function $f(R, S)$ of random variables R and S is

$$\langle f \rangle = \sum_{s \in \mathcal{S}} \int_{\mathcal{R}} p_\beta(r, s) f(r, s) d\beta \quad (12)$$

111 and this can be estimated on a set of outcomes $\{(r_i, s_i)\}$ as

$$\langle f \rangle \approx \frac{1}{n_r} \sum_i f(r_i, s_i). \quad (13)$$

112 For the mutual information this gives

$$I(R; S) \approx \frac{1}{n_r} \sum_i \log_2 p_\beta(r_i | s_i). \quad (14)$$

113 Now, an estimate for $p_\beta(r_i | s_i)$ is needed; this is approximated using KDE.

114 2.3. A Monte-Carlo Estimate

115 One advantage to using $d\beta$ as the measure is that $p_\beta(r) = 1$ and this simplifies the expression
 116 for $I(R; S)$. However, the most significant advantage is that under this new measure volumes can be
 117 estimated by simply counting data points. This is used to normalize the kernel. It is useful to define
 118 the support of a function: if $f(r)$ is a function then the support of $f(r)$, $\text{supp}[f(r)]$, is the region of its
 119 domain where it has a non-zero value,

$$\text{supp}[f(r)] = \{r : f(r) \neq 0\}. \quad (15)$$

120 Typically the size of a square kernel is specified by the radius of the support. Here, however, it is
 121 specified by volume. In a vector space where the volume measure is derived from the coordinates, there
 122 is a simple formula relating radius and volume. That is not the case here and specifying the size of
 123 a kernel by volume is not equivalent to specifying it by radius. Choosing the volume over the radius
 124 simplifies subsequent calculations and also has the advantage that the size of the kernel is related to the
 125 number of data points. This also means that the radius of the kernel varies across \mathcal{R} .

126 The term bandwidth will be used to describe the size of the kernel even though here the bandwidth is
 127 a volume rather than a radius. Since $d\beta$ is a probability measure, all volumes are between zero and one:
 128 let h be a bandwidth in this range. If $k(r'; r, h)$ is the value at r' of a square kernel with bandwidth h
 129 centered on r , the support will be denoted as $\mathcal{S}(r; h)$:

$$\mathcal{S}(r; h) = \text{supp}[k(r'; r, h)] \quad (16)$$

130 and the volume of the support of the kernel is $\text{vol}[\mathcal{S}(r; h)]$. The value of the integral is set at one,

$$\int_{\mathcal{S}(r; h)} k(r'; r, h) d\beta = 1, \quad (17)$$

131 and so, since the square kernel is being used, $k(r'; r, h)$ has a constant value of $1/\text{vol}[\mathcal{S}(r; h)]$ throughout
 132 $\mathcal{S}(r; h)$.

133 Thus, volumes are calculated using the measure $d\beta$ based on the probability density. However, this
 134 density is unknown and so volumes need to be estimated. As described above, using $d\beta$, the volume
 135 of a region is estimated by the fraction of data points that lie within it. In other words, the change of
 136 measure leads to a Monte-Carlo approach to calculating the volume of any region. In the Monte-Carlo
 137 calculation the volume of the support of a kernel is estimated as the fraction of data points that lie within
 138 it. A choice of convention has to be made between defining the kernel as containing $\lfloor hn_r \rfloor$ or $\lceil hn_r \rceil$
 139 points, that is, on whether to round hn_r down or up. The former choice is used, so, the kernel around a

140 point r is estimated as the region containing the nearest $n_h = \lfloor hn_r \rfloor$ points to r , including r itself. Thus,
 141 the kernel around a point r_i is defined as

$$k(r; r_i, n_h) = \begin{cases} \frac{1}{n_h} & r \text{ is one of the } n_h \text{ closest points to } r_i \\ 0 & \text{otherwise} \end{cases} \quad (18)$$

142 and the support $\mathcal{S}(r_i; n_h)$ has $r_j \in \mathcal{S}(r_i; n_h)$ if $k(r_j; r_i, n_h) = 1/n_h$, or, put another way, r_j is one of
 143 the n_h nearest data points. In practice, rather than rounding hn_r up or down, the kernel volume in a
 144 particular example can be specified using n_h rather than h .

145 Typically, kernels are balls: regions defined by a constant radius. As such, the kernel described here
 146 makes an implicit assumption about the isotropic distribution of the data points. However, in the normal
 147 application of KDE special provision must be made near boundaries, where the distribution of data points
 148 is not isotropic [14]. Here these cases are dealt with automatically.

149 Since $p_\beta(r_i|s_i) = n_s p_\beta(r_i, s_i)$ here, the conditional distribution $p_\beta(r_i|s_i)$ is estimated by first
 150 estimating $p_\beta(r_i, s_i)$. As described above, a kernel has a fixed volume relative to the measure based
 151 on $p_\beta(r)$. Here the kernel is being used to estimate $p_\beta(r_i, s_i)$:

$$\tilde{p}_\beta(r_i, s_i) = \frac{c(r_i, s_i; n_h)}{n_h} \quad (19)$$

152 where $c(r_i, s_i; n_h)$ is the number of data points evoked to stimulus s_i for which r_i is one of the n_h closest
 153 points

$$c(r_i, s_i; n_h) = |\{(r_j, s_i) : r_j \in \mathcal{S}(r_i; n_h)\}| \quad (20)$$

154 This gives the estimated mutual information

$$I(R; S) \approx I(R, S; n_h) = \frac{1}{n_r} \sum_i \log_2 \frac{n_s c(r_i, s_i; n_h)}{n_h} \quad (21)$$

155 Remarkably, although this is a KDE estimator it resembles a k -, or here n_h -, nearest-neighbors estimator.
 156 Basing KDE on the data available for spike trains appears to lead naturally to nearest neighbor estimation.

157 The formula for $I(R, S; n_h)$ behaves well in the extreme cases. If the responses to each stimulus
 158 are close to each other, but distant from responses to all other stimuli, then $c(r_i, s_i; n_h) = n_h$ for all
 159 stimulus-response pairs (r_i, s_i) . That is, for each data point, all nearby data points are from the same
 160 stimulus. This means that the estimate will be

$$I(R, S; n_h) = \log_2 n_s. \quad (22)$$

161 This is the correct value because, in this case, the response completely determines the stimulus, and so
 162 the mutual information is exactly the entropy of the stimulus. On the other hand, if the responses to each
 163 stimulus have the same distribution then $c(r_i, s_i; n_h)/n_h \approx 1/n_s$, so the estimated mutual information
 164 will be close to zero. This is again the correct value, because in this case the response is independent of
 165 the stimulus.

166 3. Results

167 As a test, this method has been applied to a toy data set modelled on the behavior of real spike trains. It
 168 is important that the method is applied to toy data that resemble the data type, electrophysiological data,
 169 that the method is intended to perform well on. As such, the toy model is selected to mimic the behavior
 170 of sets of spike trains. The formula derived above acts on the matrix of inter-data-point distances rather
 171 than the points themselves, and so the data set is designed to match the distance distribution observed in
 172 real spike trains [5]. The test data set is also designed to present a stiff challenge to any algorithm for
 173 estimating information.

174 The toy data are produced by varying the components of one of a set of source vectors. More
 175 precisely, to produce a test data set a variance σ^2 is chosen uniformly from $[0, 1]$ and n_s sources are
 176 chosen uniformly in a n_d -dimensional box centered at the origin with unit sides parallel to the Cartesian
 177 axes. Thus, the sources are all n_d -dimensional vectors. The data points are also n_d -dimensional vectors,
 178 they are generated by drawing each component from a normal distribution about the corresponding
 179 component of the source. Thus, data points with a source $\mathbf{s} = (s_1, s_2, \dots, s_{n_d})$ are chosen as
 180 $\mathbf{r} = (r_1, r_2, \dots, r_{n_d})$ where the r_i are all drawn from normal distributions with variance σ^2 centered
 181 at the corresponding s_i :

$$r_i \sim \mathcal{N}(s_i, \sigma^2). \quad (23)$$

182 n_t data points are chosen for each source giving $n_r = n_s n_t$ data points in all.

183 Each test uses 200 different data sets; random pruning is used to ensure the values of mutual
 184 information are evenly distributed over the whole range from zero to $\log_2 n_s$, otherwise there tends to be
 185 an excess of data sets with a low value. The true mutual information is calculated using a Monte-Carlo
 186 estimate sampled over 10,000 points. The actual probability distributions are known: the probability of
 187 finding a point \mathbf{r} generated by a source \mathbf{s} depends only on the distance $d = |\mathbf{r} - \mathbf{s}|$ and is given by the
 188 χ -distribution

$$p(d) = \frac{2^{1-n_d/2}}{\Gamma(n_d/2)} \left(\frac{d}{\sigma}\right)^{n_d-1} e^{-d^2/2\sigma^2}. \quad (24)$$

189 There is a bias in estimating the mutual information, in fact, bias is common to any approach to
 190 estimating mutual information [15]. The problem of reducing bias, or defining the mutual information so
 191 that the amount of bias is low, is well studied and has produced a number of sophisticated approaches [4,
 192 15–19]. One of these, quadratic estimation, due to [16,18], is adapted to the current situation. Basically,
 193 it is assumed that for large numbers of data points n_t the estimated information $\tilde{I}(R; S)$ is related to the
 194 true mutual information $I(R; S)$ by

$$\tilde{I}(R; S) = I(R; S) + \frac{A}{n_t} + \frac{B}{n_t^2} + O(1/n_t^3). \quad (25)$$

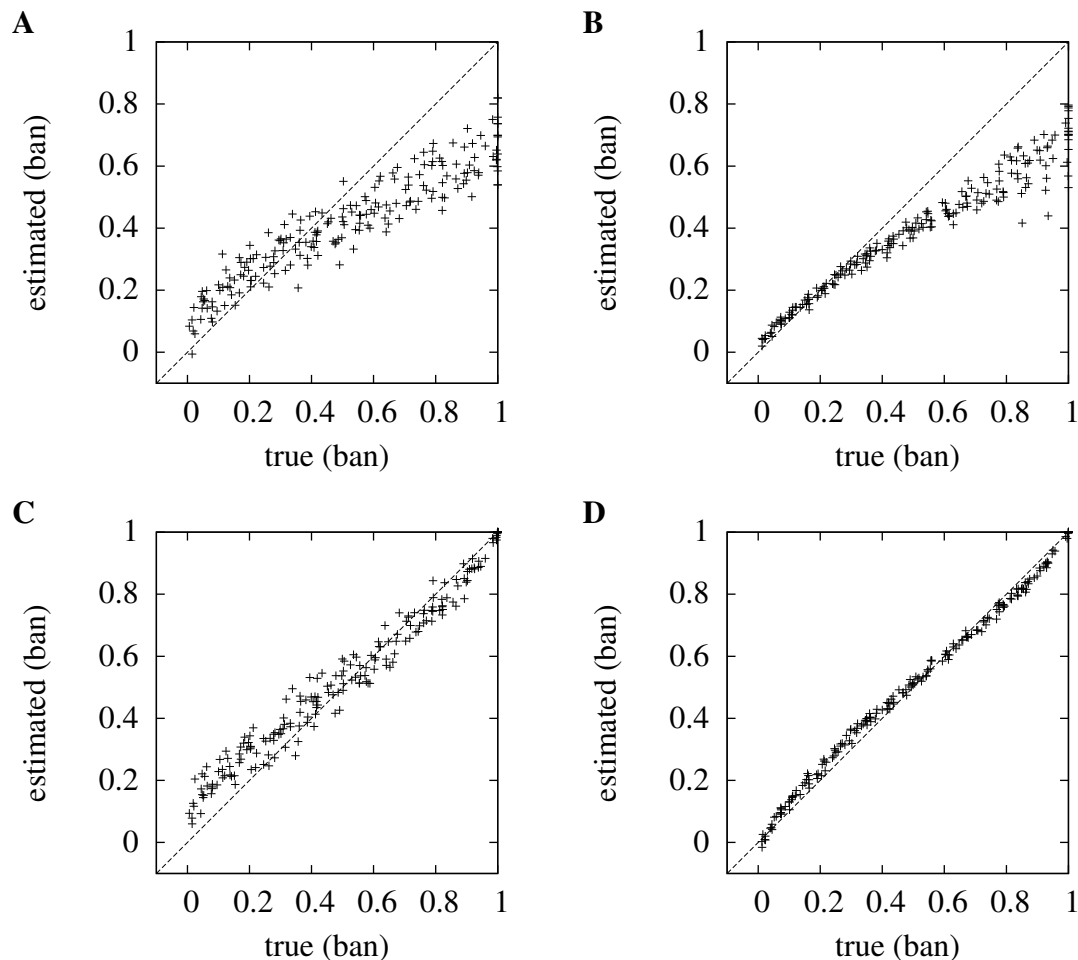
195 This asymptotic expansion is well-motivated in the case of the plug-in approach to spike train information
 196 [15,16,20–22] and it is assumed the same expansion applies. To extract $I(R; S)$ the estimate $I(R; S; n_h)$
 197 is calculated for λn_r with λ taking values from 0.1 to one in 0.1 increments. Least squares fit is used to
 198 estimate $I(R; S)$ from these ten values.

199 The new method works well on these toy data. It is compared to a histogram approach where the
 200 n_d -dimensional space is discretized into bins and counting is used to estimate the probability of each

201 bin. This is an analog of the plug-in method and the same quadratic estimation technique is used to
 202 reduce bias.

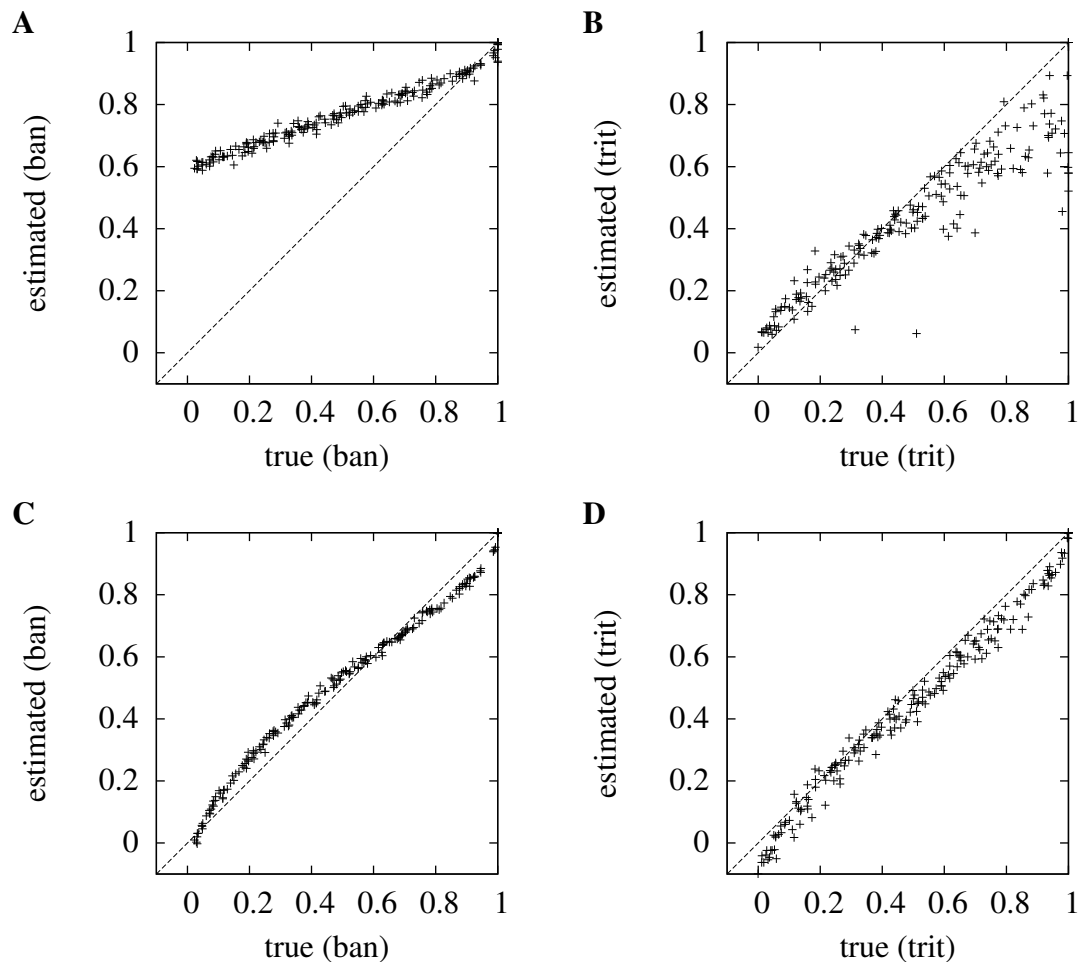
203 In Fig. 1 the new method is compared to the histogram method when $n_s = 10$ and $n_d = 3$ and for
 204 both low and high numbers of trials, $n_t = 10$ and $n_t = 200$. For the histogram method the optimum
 205 discretization width is used. This optimal width is large, $h = 5$ in each case; this roughly corresponds
 206 to a different bin for each octant of the three-dimensional space containing the data. In the new method
 207 the bandwidth is not optimized on a case by case basis, instead, the kernel bandwidth n_h is chosen as
 208 being equal to the number of trials n_t . It can be seen that the new method is better at estimating the
 209 information; for $n_t = 10$ it has an average absolute error of 0.189 bits, compared to 0.481 bits for the
 210 histogram method, for $n_t = 200$ the average absolute error is 0.083 bits, compared to 0.442 bits for the
 211 histogram approach.

Figure 1. Comparing the KDE to the histogram method for ten sources, $n_s = 10$, and three dimensions $n_d = 3$. In each case the true information is plotted against the estimated information; the line $y = x$ which represents perfect estimation is plotted for clarity. For convenience, the mutual information has been normalized, so in each case the value plotted is the estimate of $I(R; S) / \log_2 n_s$, with a maximum value of one; in the cases plotted here that means the information is measured in ban. **A** and **B** show the distribution for the histogram method for $n_t = 10$ and $n_t = 200$, **C** and **D** show the kernel method.



212 In Fig. 2 the histogram and kernel methods are compared for $n_s = 10$ and $n_d = 10$ and for $n_s = 3$
 213 and $n_d = 3$; the number of trials is $n_t = 200$ in each case. The kernel method outperforms the histogram
 214 method. When $n_s = 10$ and $n_d = 10$ the average absolute error for the kernel method is 0.139 bits,
 215 compared to 0.876 bits for the histogram method; for $n_s = 3$ and $n_d = 3$ its average absolute error
 216 is 0.076 bits compared to 0.141 bits for the histogram. Furthermore, the errors for the kernel method
 217 are less clearly modulated by the actual information, which makes the method less prone to producing
 218 misleading results.

Figure 2. Comparing the KDE to the histogram method for high and low numbers of sources and dimensions. The true information is plotted against the estimated information; in **A** and **C** $n_s = 10$ and $n_d = 10$, in **B** and **D** $n_s = 3$ and $n_d = 3$. The top row, **A** and **B**, are for the histogram method, the bottom row, **C** and **D**, are for the kernel method. As before, the normalized information, $I(S; R)/\log_2 n_s$ is plotted, so for $n_s = 10$ the information is in ban, for $n_s = 3$ it is in trit and in each case the maximum mutual information is one. $n_r = 200$ for all graphs.



220 Although the actual method presented here is very different, it was inspired in part by the transmitted
 221 information method for calculating mutual information using metric-based clustering described in [26]
 222 and by the novel approach introduced in [11] where a kernel-like approach to mutual information
 223 is developed. Another significant motivation was the interesting technique given in [12] where the
 224 information is estimated by measuring how large a sphere could be placed around each data point without
 225 it touching another data point. In [12], the actual volume of the sphere is required, or rather the rate the
 226 volume changes with diameter. This is calculated by foliating the space of spike trains into subspaces
 227 with fixed spike number and interpreting the spike times as coordinates. This is avoided here by using the
 228 Monte-Carlo estimate of volumes. Finally, the copula construction is related to the approach described
 229 here. In fact, the construction here can be thought of as a reverse copula construction [13].

230 An important part of the derivation of the kernel method is the change of measure to one based on
 231 the distribution. Since the kernel size is defined using a volume based on this measure, the radius of
 232 the kernel adapts to the density of data points. This is similar to the adaptive partitioning described for
 233 example in [24]. Like the plug-in method of computing mutual information for spike trains, adaptive
 234 partitioning is a discretization approach. However, rather than breaking the space into regions of fixed
 235 width, the discrete regions are chosen dynamically, using estimates of the cumulative distribution, similar
 236 to what is proposed here.

237 One striking aspect of KDE seen here is that it reduces to a k th nearest-neighbor (kNN) estimator.
 238 The kNN approach to estimating the mutual information of variables lying in metric spaces has been
 239 studied directly in [23]. Rather than using a KDE of the probability distribution, a Kozachenko-Leonenko
 240 estimator [25] is used. To estimate $I(X; Y)$ where X and Y are both continuous random variables taking
 241 values in \mathcal{X} and \mathcal{Y} , Kozachenko-Leonenko estimates are calculated for $H(X)$, $H(Y)$ and $H(X, Y)$; by
 242 using different values of k in each space the terms that would otherwise depend on the dimension of \mathcal{X}
 243 and \mathcal{Y} cancel.

244 This approach can be modified to estimate $I(R; S)$ where S is a discrete random variable. Using the
 245 approach described in [23] to estimate $H(R)$ and $H(R|S)$ gives

$$I_e(R; S) \approx F(n_k) + F(n_t n_s) - F(n_t) - \frac{1}{n_r} \sum_i F[C(r_i, s_i; n_k)] \quad (26)$$

246 where $F(x)$ is the digamma function, n_k is an integer parameter and $C(r_i, s_i; n_k)$ is similar to $c(r_i, s_i; n_h)$
 247 above. Whereas $c_k(r_i, s_i; n_h)$ counts the number of responses to s_i for which r_i is one of the n_h closest
 248 data points, $C(r_i, s_i; n_k)$ is computed by first finding the distance d from r_i to the n_k th nearest spike-train
 249 response to stimulus s_i ; then $C(r_i, s_i; n_k)$ counts the number of spike trains, from any stimulus, that are at
 250 most a distance of d from r_i . $I_e(R; S)$ is the mutual information with base e , so $I(R; S) = I_e(R; S)/\ln 2$.
 251 During the derivation of this formula, expressions involving the dimension of \mathcal{S} appear, but ultimately
 252 they all cancel, leaving an estimate which can be applied in the case of interest here, where \mathcal{S} has no
 253 dimension. Since the digamma function can be approximated as

$$F(x) = \ln x - \frac{1}{2x} \quad (27)$$

254 for large x this kNN approach and the kernel method produce very similar estimates. The similarity
 255 between the two formulas, despite the different routes taken to them, lends credibility to both estimators.

256 Other versions of the kernel method can be envisaged. A kernel with a different shape could be used
257 or the kernel could be defined by the radius rather than by the volume of the support. The volume of the
258 support and therefore the normalization would then vary from data point to data point. This volume could
259 be estimated by counting, as it was here. However, as mentioned above, the volume based bandwidth
260 has the advantage that it gives a kernel which is adaptive, the radius varies as the density of data points
261 changes. Another intriguing possibility is to investigate if it would be possible to follow [12] and [23]
262 more closely than has been done here and use a Monte-Carlo volume estimate to derive a Kozachenko
263 and Leonenko estimator. Finally, KDE applied to two continuous random variables could be used to
264 derive an estimate for the mutual information between two sets of spike trains, or between a set of spike
265 trains and a non-discrete stimulus such as position in a maze.

266 There is no general, principled, approach to choosing bandwidths for KDE methods. There are
267 heuristic methods, such as cross-validation [9,10], but these include implicit assumptions about how the
268 distribution of the data is itself drawn from a family of distributions, assumptions which may not apply
269 to a particular experimental situation. The KDE approach developed here includes a term analogous to
270 bandwidth and, although a simple choice of this bandwidth is suggested and gives accurate estimates,
271 the problem of optimal bandwidth selection will require further study.

272 Applying the KDE approach to spike trains means it is necessary to specify a spike train metric
273 [26–28]. Although the metric is only used to arrange points in order of proximity, the dependence
274 on a metric does mean that the estimated mutual information will only include mutual information
275 encoded in features of the spike train that affect the metric. As described in [12], in the context of
276 another metric-dependent estimator of mutual information, this means the mutual information may
277 under-estimate the true mutual information, but it does allow the coding structure of spike trains to
278 be probed by manipulating the spike train metrics.

279 It is becoming increasingly possible to measure large number spike trains from large numbers of spike
280 trains simultaneously. There are metrics for measuring distances between sets of multi-neuron responses
281 [29–31] and so the approach described here can also be applied to multi-neuronal data.

282 Acknowledgements

283 RJT is grateful to the Irish Research Council in Science, Engineering and Technology for financial
284 support. CJH is grateful to the James S McDonnell Foundation for financial support through a Scholar
285 Award in Human Cognition.

286 Conflicts of Interest

287 The authors declare no conflicts of interest.

288 References

- 289 1. de Ruyter van Steveninck, R.R.; Lewen, G.D.; Strong, S.P.; Koberle, R.; Bialek, W.
290 Reproducibility and variability in neural spike trains. *Science* **1997**, *275*, 1805–1808.
- 291 2. Strong, S.; Koberle, R.; de Ruyter van Steveninck, R.R.; Bialek, W. Entropy and information in
292 neural spike trains. *Physics Review Letters* **1998**, *80*, 197–200.

- 293 3. Nemenman, I.; Bialek, W.; de Ruyter van Steveninck, R. Entropy and information in neural spike
294 trains: Progress on the sampling problem. *Physical Review E* **2004**, *69*, 056111.
- 295 4. Nemenman, I.; Lewen, G.; Bialek, W.; de Ruyter van Steveninck, R.R. Neural coding of natural
296 stimuli: information at sub-millisecond resolution. *BMC Neuroscience* **2007**, *8*, S7.
- 297 5. Gillespie, J.; Houghton, C. A metric space approach to the information capacity of spike trains.
298 *Journal of Computational Neuroscience* **2011**, *30*, 201–209.
- 299 6. Rosenblatt, M. Remarks on some nonparametric estimates of a density function. *The Annals of*
300 *Mathematical Statistics* **1956**, pp. 832–837.
- 301 7. Parzen, E. On estimation of a probability density function and mode. *The Annals of Mathematical*
302 *Statistics* **1962**, *33*, 1065–1076.
- 303 8. Silverman, B. *Density Estimation*; Chapman and Hall: London, 1986.
- 304 9. Rudemo, M. Empirical choice of histograms and kernel density estimators. *Scandinavian*
305 *Journal of Statistics* **1982**, pp. 65–78.
- 306 10. Hall, P. Large sample optimality of least squares cross-validation in density estimation. *The*
307 *Annals of Statistics* **1983**, pp. 1156–1174.
- 308 11. Brasselet, R.; Johansson, R.S.; Arleo, A. Quantifying neurotransmission reliability through
309 metrics-based information analysis. *Neural Computation* **2011**, *23*, 852–881.
- 310 12. Victor, J.D. Binless strategies for estimation of information from neural data. *Physical Review E*
311 **2002**, *66*, 051903.
- 312 13. Calsaverini, R.S.; Vicente, R. An information-theoretic approach to statistical dependence:
313 Copula information. *EPL (Europhysics Letters)* **2009**, *88*, 68003.
- 314 14. Jones, M.C. Simple boundary correction for kernel density estimation. *Statistics and Computing*
315 **1993**, *3*, 135–146.
- 316 15. Paninski, L. Estimation of entropy and mutual information. *Neural Computation* **2003**, *15*, 1191–
317 1253.
- 318 16. Treves, A.; Panzeri, S. The upward bias in measures of information derived from limited data
319 samples. *Neural Computation* **1995**, *7*, 399–407.
- 320 17. Panzeri, S.; Treves, A. Analytical estimates of limited sampling biases in different information
321 measures. *Network* **1996**, *7*, 87–107.
- 322 18. Panzeri, S.; Senatore, R.; Montemurro, M.A.; Petersen, R.S. Correcting for the sampling bias
323 problem in spike train information measures. *Journal of Neurophysiology* **2007**, *98*, 1064–1072.
- 324 19. Montemurro, M.; Senatore, R.; Panzeri, S. Tight data-robust bounds to mutual information
325 combining shuffling and model selection techniques. *Neural Computation* **2007**, *19*, 2913–2957.
- 326 20. Miller, G. Note on the bias of information estimates. In *Information theory in psychology II-B*;
327 Quastler, H. Ed.; Free Press: Glencoe, IL, 1955; pp. 95–100.
- 328 21. Carlton, A. On the bias of information estimates. *Psychological Bulletin* **1969**, *71*, 108–109.
- 329 22. Victor, J.D. Asymptotic bias in information estimates and the exponential (Bell) polynomials.
330 *Neural Computation* **2000**, *12*, 2797–2804.
- 331 23. Kraskov, A.; Stögbauer H.; Grassberger, P. Estimating mutual information *Physical Review E*
332 **2004**, *69*, 066138

- 333 24. Darbellay, G.A.; Vajda I. Estimation of the information by an adaptive partitioning of the
334 observation space *IEEE Transactions on Information Theory* **1999**, *45*, 1315–1321
- 335 25. Kozachenko, L.; Leonenko, N. On statistical estimation of entropy of a random vector *Problems*
336 *of Information Transmission* **1987**, *23*, 9–16
- 337 26. Victor, J.D.; Purpura, K.P. Nature and precision of temporal coding in visual cortex: a metric-
338 space analysis. *Journal of Neurophysiology* **1996**, *76*, 1310–1326.
- 339 27. van Rossum, M. A novel spike distance. *Neural Computation* **2001**, *13*, 751–763.
- 340 28. Houghton, C.; Victor, J.D. Spike rates and spike metrics. In *Visual Population Codes: Toward a*
341 *Common Multivariate Framework for Cell Recording and Functional Imaging*; Kriegeskorte N.,
342 Kreiman, G. Ed.; MIT Press: Cambridge, MA, 2012; chap. 8.
- 343 29. Aronov, D.; Reich, D.S.; Mechler, F.; Victor, J.D. Neural coding of spatial phase in v1 of the
344 macaque monkey. *Journal of Neurophysiology* **2003**, *89*, 3304–3327.
- 345 30. Houghton, C.; Sen, K. A new multi-neuron spike-train metric. *Neural Computation* **2008**,
346 *20*, 1495–1511.
- 347 31. Kreuz, T.; Chicharro, D.; Houghton, C.; Andrzejka, R.G.; Mormann, F. Monitoring spike train
348 synchrony. *Journal of Neurophysiology* **2013**, *109*, 1457–1472.

349 © October 13, 2013 by the authors; submitted to *Entropy* for possible open access
350 publication under the terms and conditions of the Creative Commons Attribution license
351 <http://creativecommons.org/licenses/by/3.0/>.