

# Blessings of Dimensionality: Theoretical analysis of nearest-neighbor projected-distance methods for detecting interactions in high dimension

Bryan A. Dawkins<sup>1</sup>, Trang T. Le<sup>2</sup> and Brett A. McKinney<sup>1,3,\*</sup>

<sup>1</sup>Department of Mathematics, University of Tulsa, Tulsa, OK 74104, USA

<sup>2</sup>Department of Biostatistics, Epidemiology and Informatics, University of Pennsylvania, Philadelphia, PA 19104

<sup>3</sup>Tandy School of Computer Science, University of Tulsa, Tulsa, OK 74104, USA.

## Abstract

It is commonly known that high-throughput data has many inherent statistical challenges, such as multiple testing, sparsity and over fitting. Collectively these challenges are known as the Curse of Dimensionality. Here we highlight an important Blessing of Dimensionality: the ability to identify interactions with nearest neighborhoods. We review nearest-neighbor concepts for finding interactions, and we derive important distribution moments for distance metrics in high dimensional spaces. We use these theoretical results and simulated data to offer recommendations for computational approaches to find nearest neighbors in high dimension. We discuss ways to maximize the blessings and minimize the curses of dimensionality to reliably identify interactions.

## Author summary

## Introduction

Relief-based methods identify interacting attributes as important by using nearest-neighbor information in higher dimensions (the “blessings of dimensionality”). Myopic methods, such as univariate tests, that do not account for information from higher dimensions, are susceptible to false negatives when there are interactions. For example in the plot of variable A versus C in a three-variable simulation (Fig. 1a), variable A appears to show no difference between cases and controls (the marginal group means are the same). However, A is actually simulated to have a strong differential correlation with B, conditioned on the outcome variable (Fig. 2b). Current Relief-based methods determine the importance of an attribute by computing the average difference of a target instance (X) and its nearest instance from the same class (Hit) projected onto the attribute A dimension ( $d_{X,H}(A)$ ) subtracted from the projected difference of target X and its nearest instance from the opposite class (Miss) ( $d_{X,M}(A)$ ). When the inequality  $d_{X,M}(A) > d_{X,H}(A)$ , it suggests that attribute A is useful for discriminating between cases and controls.



**Fig 1. Imposters vs true neighbors in the presence of interactions with three variables.** Scatter plot of simulated irrelevant Attribute C with a functional Attribute A (a). None of the attributes has a main effect, but Attribute B and C interact through differential correlation (b). Computing nearest neighbors with irrelevant attributes (a) or lower dimensions leads to imposter nearest neighbors and degrades the ability of Relief-based methods to identify interaction effects. Computing distances in only these two dimensions leads to an imposter false miss (FM) for the nearest neighbor

from the opposite outcome class for target instance X. This imposter leads to attribute A predicting closer projected distances for misses than hits (H), which incorrectly indicates that A is a poor discriminator (yellow boxes in (a)). Computing nearest neighbors in higher dimensions (c-d) or with the correct interaction partner leads to imposter nearest neighbor (FM) being replaced by the true nearest miss neighbor (TM) for target instance X, which correctly leads to attribute A predicting closer projected distances for hits (H) than misses, which is an indication that attribute A is a good discriminator (yellow boxes (b)).

Relief-based methods use information from all attributes available to it (omnigenic) to estimate an attribute's importance. However, if relevant higher-dimensional information is not used, even Relief-based methods will miss the effect of A because "imposter" neighbors will be used in the attribute estimate (False Miss (FM) in Fig. 1, where  $d_{X,FM}(A) < d_{X,H}(A)$ ). If one were to compute nearest neighbors in the A-C plane (ignoring the B dimension), the nearest miss would be an imposter (FM), which leads to a negative contribution to the importance score for A. One might call this C attribute a type-I confounding attribute because it increases the chances of interacting attributes to be false negatives. When nearest neighbors are calculated based on higher dimensions with relevant information (Fig. 2c), it is clear that TM is closer to X than FM. The imposter (FM) is replaced by the true nearest miss (TM) and attribute A correctly shows a greater projected difference between misses than hits (Fig. 2d  $d_{X,TM}(A) > d_{X,H}(A)$ ), which is the signature of an important attribute. Univariate methods still cannot find the importance of A unless the interaction is explicitly modeled, but as long as functional variables A and B are in the space for nearest neighbor calculations (Fig. 2c-d), imposters can be excluded and Relief-based methods will find that A (and B) are important discriminators.

[Ideas: Relating the increasing  $k$  and myopic view to other distance-related method such as MDS/t-SNE vs PCA - local vs global distance - capturing non-linear manifold structure.

<https://www.kdnuggets.com/2018/08/introduction-t-sne-python.html>

Using same interaction, increase background noise genes to see degrading of A and B Relief importance because of curse of dimensionality (sparseness).



Fig 2. True neighbors

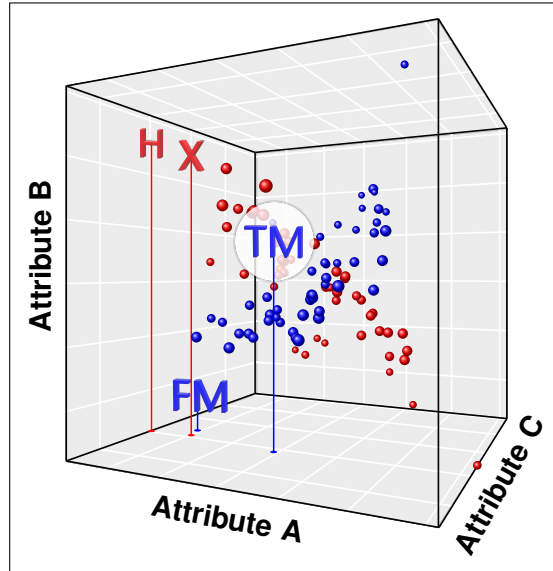


Fig 3. 3D AB view. Still working on this.

## 1 Neighborhood methods

NPD methods rely on a neighborhood algorithm for feature selection. One may specify a fixed- $k$  number of neighbors, an average radius SURF, a multiSURF radius that adapts for each instance [1], or a gene-wise adaptive- $k$ .

## 2 Derivation of expected $k$ for multiSURF neighborhoods

The multiSURF radius for an instance is the mean of its distances to all other instances subtracted by  $\alpha = 1/2$  of the standard deviation of this mean. Previously we showed



Fig 4. 3D AC view. Still working on this.

empirically for balanced case-control datasets that a good constant- $k$  approximation to the expected number of neighbors within the multiSURF radii is  $k = m/6$  [2], where  $m$  is the number of samples. Here we derive a more exact theoretical mean that shows the mathematical connection between neighbor-finding methods. This fixed- $k$  approximation to multi-SURF is independent of the type of data and the particular radii of each instance in the data.

The distance between instances  $i$  and  $j$  in the data set  $X^{m \times p}$  of  $m$  instances and  $p$  attributes is calculated in the space of all attributes ( $a \in A$ ,  $|A| = p$ ) using a metric such as

$$D_{ij}^{(q)} = \left( \sum_{a \in A} |d_{ij}(a)|^q \right)^{1/q}, \quad (1)$$

which is typically Manhattan ( $q = 1$ ) but may also be Euclidean ( $q = 2$ ). The quantity  $d_{ij}(a)$ , known as a “diff” in Relief literature, is the projection of the distance between instances  $i$  and  $j$  onto the attribute  $a$  dimension. The function  $d_{ij}(a)$  supports any type of attributes (e.g., numeric continuous versus categorical). For example, the projected difference between two instances  $i$  and  $j$  for a continuous numeric ( $d^{\text{num}}$ ) attribute  $a$  may be

$$\begin{aligned} d_{ij}^{\text{num}}(a) &= \text{diff}(a, (i, j)) \\ &= |\hat{X}_{ia} - \hat{X}_{ja}|, \end{aligned} \quad (2)$$

where  $\hat{X}$  represents the standardized data matrix  $X$ . We use a simplified  $d_{ij}(a)$  notation in place of the  $\text{diff}(a, (i, j))$  notation that is customary in Relief-based methods. We omit the division by  $\max(a) - \min(a)$  used by Relief to constrain scores to the interval from  $-1$  to  $1$ . As we show in subsequent sections, NPDR scores are [standardized] regression coefficients with corresponding P values, so any scaling operation at this stage is unnecessary for comparing attribute scores. The numeric  $d_{ij}^{\text{num}}(a)$  projection is simply the absolute difference between row elements  $i$  and  $j$  of the data matrix  $X^{m \times p}$  for the attribute column  $a$ .

We define the NPDR neighborhood set  $\mathcal{N}$  of ordered pair indices as follows. Instance  $i$  is a point in  $p$  dimensions, and we designate the topological neighborhood of  $i$  as  $N_i$ .

This neighborhood is a set of other instances trained on the data  $X^{m \times p}$  and depends on the type of Relief neighborhood method (e.g., fixed- $k$  or adaptive radius) and the type of metric (e.g., Manhattan or Euclidean). If instance  $j$  is in the neighborhood of  $i$  ( $j \in N_i$ ), then the ordered pair  $(i, j) \in \mathcal{N}$  for the projected-distance regression analysis. The ordered pairs constituting the neighborhood can then be represented as nested sets:

$$\mathcal{N} = \{\{(i, j)\}_{i=1}^m\}_{\{j \neq i: j \in N_i\}}. \quad (3)$$

The cardinality of the set  $\{j \neq i: j \in N_i\}$  is  $k_i$ , the number of nearest neighbors for subject  $i$ .

## 2.1 Distribution of pairwise distances

Suppose that  $X_{ia}, X_{ja} \stackrel{iid}{\sim} \mathcal{F}_X(\mu_X, \sigma_X^2)$  for two fixed and distinct instances  $(i, j) \in \mathcal{N}$  and a fixed attribute  $a \in \mathcal{A}$ .

It is clear that  $|X_{ia} - X_{ja}|^q = |d_{ij}(a)|^q$  is another random variable. Let  $Z_a^q \sim \mathcal{F}_{Z^q}(\mu_{z^q}, \sigma_{z^q}^2)$  be the random variable such that

$$Z_a^q = |d_{ij}(a)|^q = |X_{ia} - X_{ja}|^q, \quad a \in \mathcal{A}. \quad (4)$$

Furthermore, the collection  $\{Z_a^q | a \in \mathcal{A}\}$  is a random sample of size  $p$  of mutually independent random variables. Hence, the sum of  $Z_a^q$  over all  $a \in \mathcal{A}$  is asymptotically normal by the Classical Central Limit Theorem (CCLT). More explicitly, this implies that

$$\left(D_{ij}^{(q)}\right)^q = \sum_{a \in \mathcal{A}} |d_{ij}(a)|^q = \sum_{a \in \mathcal{A}} |X_{ia} - X_{ja}|^q = \sum_{a \in \mathcal{A}} Z_a^q \sim \mathcal{N}(\mu_{z^q} p, \sigma_{z^q}^2 p) \quad (5)$$

Consider the smooth function  $g(z) = z^{1/q}$  that is continuously differentiable for  $z > 0$ . Assuming that  $\mu_{z^q} > 0$ , the Delta Method [3] can be applied to show that

$$\begin{aligned} g\left(\left(D_{ij}^{(q)}\right)^q\right) &= g\left(\sum_{a \in \mathcal{A}} Z_a^q\right) \\ &= \left(\sum_{a \in \mathcal{A}} |X_{ia} - X_{ja}|^q\right)^{1/q} \\ &= D_{ij}^{(q)} \sim \mathcal{N}\left(g(\mu_{z^q} p), [g'(\mu_{z^q} p)]^2 \sigma_{z^q}^2 p\right) \\ \Rightarrow D_{ij}^{(q)} &\sim \mathcal{N}\left((\mu_{z^q} p)^{1/q}, \frac{\sigma_{z^q}^2 p}{q^2 (\mu_{z^q} p)^{2(1-\frac{1}{q})}}\right) \end{aligned} \quad (6)$$

For  $q = 2$ , a more accurate approximation of the sample mean is given by

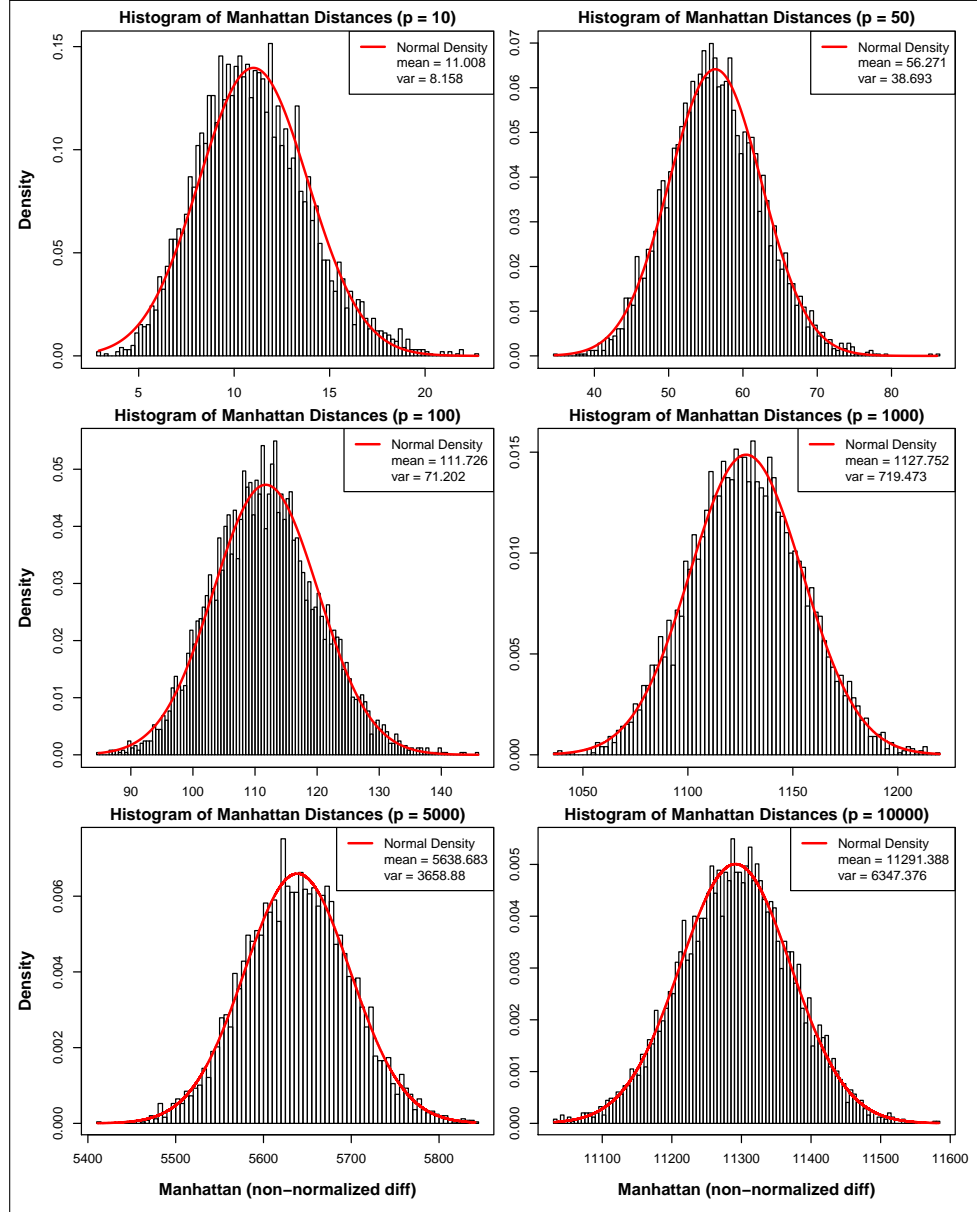
$$\mathbb{E}\left(D_{ij}^{(q)}\right) = \left(\mu_{z^q} p - \frac{\sigma_{z^q}^2 p}{q^2 (\mu_{z^q} p)^{2(1-\frac{1}{q})}}\right)^{1/q} \quad (7)$$

This is because

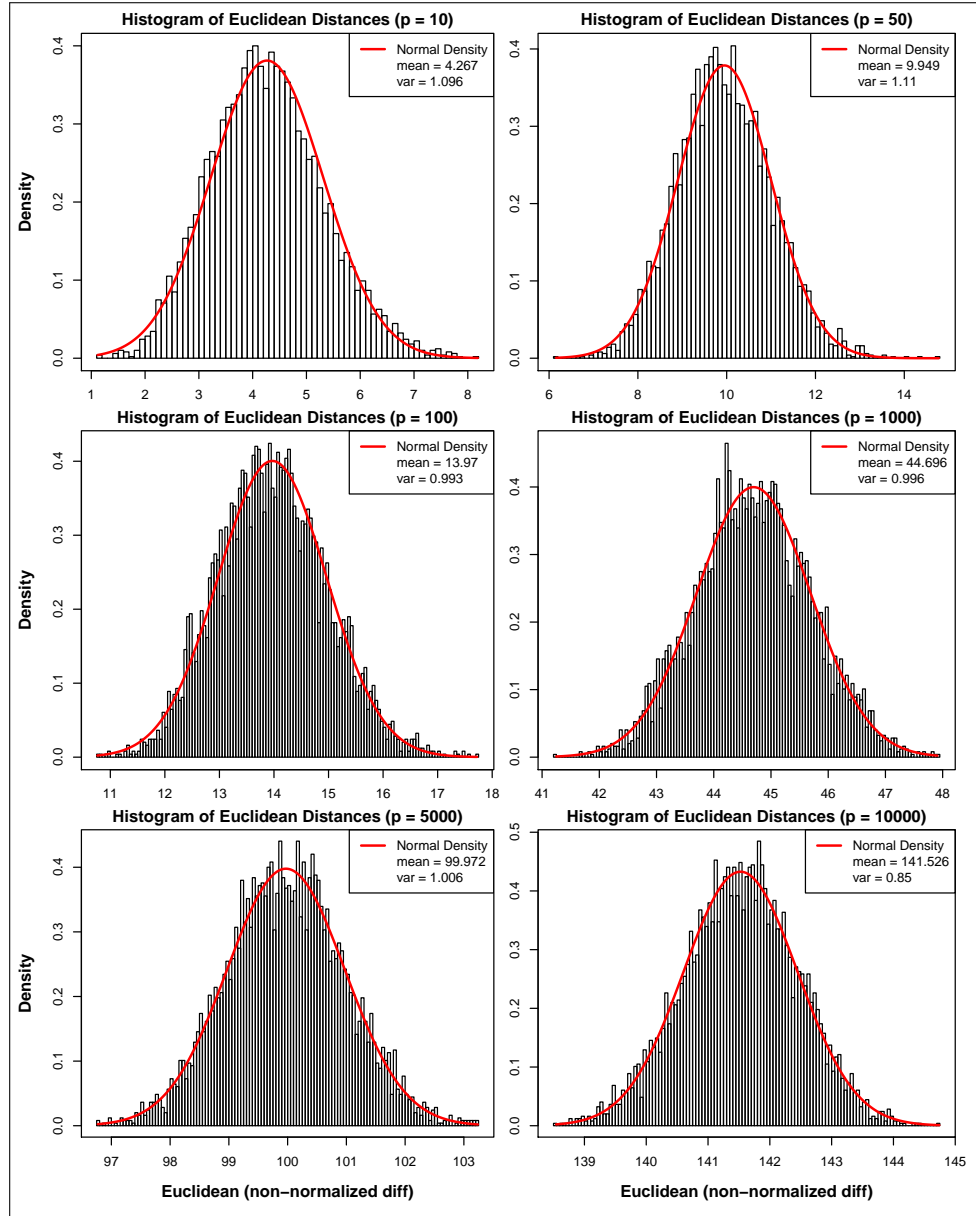
$$\begin{aligned} \mathbb{E}\left(D_{ij}^{(2)}\right) &= \sqrt{\mathbb{E}\left[\left(D_{ij}^{(2)}\right)^2\right] - \text{Var}\left(D_{ij}^{(2)}\right)} \\ &= \sqrt{\mu_{z^2} p - \frac{\sigma_{z^2}^2}{4\mu_{z^2}}} \end{aligned} \quad (8)$$

Therefore, the distance between two fixed and distinct instances  $i$  and  $j$  given by Eq. 1 is asymptotically normal.

One can readily verify the normality of distances between independent instances through sampling from any data distribution and plotting the histogram of pairwise distances. Histograms for the case of standard normal data and Manhattan ( $q=1$ ) and Euclidean ( $q=2$ ) metrics are shown in Figs. 5 and 6, respectively. For these simulated distances, we fixed  $m = 100$  and varied  $p$  from 10 to 10000. For even a rather small number of predictors, as in the case of  $p = 10$ , the sample distances are approximately normal. As  $p$  increases, the normality becomes stronger.



**Fig 5.** Convergence to normality of Manhattan distances between iid random normal instances. For each simulated distance distribution, we fixed  $m = 100$  instances but varied  $p$  from 10 to 10000. It is clear that convergence is rapid, and approximate normality can be safely assumed for even  $p = 10$ .



**Fig 6.** Convergence to normality of Euclidean distances between iid random normal instances. For each simulated distance distribution, we fixed  $m = 100$  instances but varied  $p$  from 10 to 10000. It is clear that convergence is rapid, and approximate normality can be safely assumed for even  $p = 10$ .

For distance based learning methods, all pairwise distances are used to determine relative importances for attributes. The collection of all distances above the diagonal in an  $m \times m$  distance matrix does not satisfy the independence assumption used in the previous derivations. This is because of the redundancy that is inherent to the distance matrix calculation. However, this collection is still asymptotically normal with mean and variance approximately equal to those given in Eq. 6. Hence, all fixed-radius methods will use a fixed radius that is some fraction of the expected pairwise distance for a given metric and data type. This implies that the probability of a fixed instance  $j$  being within a fixed radius of a given instance  $i$  can be parameterized by the expected

101  
102  
103  
104  
105  
106  
107  
108  
109

pairwise distance and the variance of the pairwise distance. This probability is obtained by evaluating the normal cumulative distribution function (CDF), with corresponding mean and variance, at the quantile given by some function of the fixed radius. Therefore, we can derive the expected number of neighbors in the neighborhood of a fixed instance  $i$ . In other words, for sufficiently large data sets, the sample mean of the number of neighbors in a given neighborhood is well approximated by the product between the total number of possible neighbors and the expected probability of an instance being in a given neighborhood. The total number of possible neighbors for a fixed instance  $i$  is always  $m - 1$ , but this becomes approximately  $\lfloor \frac{m-1}{2} \rfloor$  when delineating between possible hits and misses for balanced data.

## 2.2 Predicted number of neighbors in the multiSURF alpha neighborhood

Regardless of the predictor data type (numeric or categorical), the distribution of the  $p$  predictors (uniform, Gaussian, or binomial), or the metric used to compute distances (Manhattan or Euclidean), the  $m(m - 1)/2$  pairwise distances in the  $p$ -dimensional space are well approximated by a normal distribution. An instance  $j$  is in the adaptive  $\alpha$ -radius neighborhood of  $i$  ( $j \in N_i^\alpha$ ) under the condition

$$D_{ij} \leq R_i^\alpha \implies j \in N_i^\alpha, \quad (9)$$

where the threshold radius for instance  $i$  is

$$R_i^\alpha = \bar{D}_i - \alpha \sigma_{\bar{D}_i} \quad (10)$$

and

$$\bar{D}_i = \frac{1}{m-1} \sum_{j \neq i} D_{ij}^{(\cdot)} \quad (11)$$

is the average of instance  $i$ 's pairwise distances (using Eq. D Equation) with standard deviation  $\sigma_{\bar{D}_i}$ . MultiSURF uses  $\alpha = 1/2$  [4].

The probability of the remaining  $m - 1$  instances being inside the  $\alpha$ -radius of instance  $i$  ( $R_i^\alpha$ ) can be viewed as  $m - 1$  Bernoulli trials each with a probability of success  $q_\alpha$ . Then the average average number of neighbors is given by

$$\bar{k}_\alpha = (m - 1)q_\alpha, \quad (12)$$

from the mean of a binomial random variable. To calculate  $q_\alpha$ , we assume the distribution of distances ( $\{D_{ij}\}_{j \neq i}$ ) of neighbors of instance  $i$  is normal  $N(\bar{D}_i, \sigma_{\bar{D}_i})$ . Our empirical studies confirm a normal distribution and that it is robust to data type and metric. Extreme violations of independence of attributes (extreme correlations or interactions) will cause the distribution to be right skewed, but this effect is difficult to observe in real data. Thus, for a Gaussian pairwise distance distribution, the probability  $q_\alpha$  for one instance  $j \neq i$  to be in the neighborhood of  $i$  ( $j \in N_i^\alpha$ ) is given by the area under the mean-centered ( $\bar{D}_i$ ) Gaussian from  $-\infty$  to  $R_i^\alpha$ . **show Gaussian plot illustration?** An illustration of the area computed to estimate  $q_\alpha$  is given by Fig. 7. This integral can be written in terms of the error function (erf):

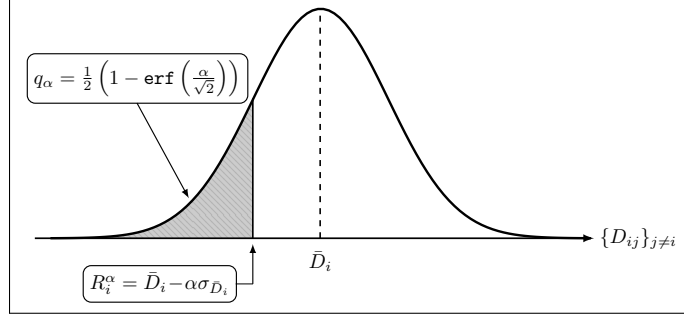
$$q_\alpha = \frac{1}{2} \left( 1 - \operatorname{erf} \left( \frac{\alpha}{\sqrt{2}} \right) \right). \quad (13)$$

And finally using Eqs. (12 and 13) we find

$$\bar{k}_\alpha = \left\lfloor \frac{m-1}{2} \left( 1 - \operatorname{erf} \left( \frac{\alpha}{\sqrt{2}} \right) \right) \right\rfloor, \quad (14)$$



where we apply the floor to ensure the number of neighbors is integer. For data with balanced hits and misses in standard fixed- $k$  Relief, one divides this formula by 2. For multiSURF ( $\alpha = 1/2$ ), this formula gives  $\bar{k}_{1/2}^{\text{hit/miss}} = \frac{1}{2}\bar{k}_{1/2} = .154(m-1)$ , which is very close to our previous empirical estimate  $m/6$ . When we compare multiSURF neighborhood methods with fixed- $k$  neighborhoods, we use  $\bar{k}_{1/2}$ . Using this  $\alpha = 1/2$  value has been shown to give good performance for simulated data sets. However, the best value for  $\alpha$  is likely data-specific and may be determined through nested cross-validation and other parameter tuning methods.



**Fig 7.** Illustration of the expected probability of a fixed instance  $j$  being in the fixed radius neighborhood of another instance  $i$ . The fixed radius is parameterized by a fraction  $\alpha$  of the standard deviation of all pairwise distances measured from instance  $i$  to all possible neighbors.

### 3 Derivation of means and standard deviations for metrics and data distributions

#### 3.1 Distribution of $|d_{ij}(a)|^q = |X_{ia} - X_{ja}|^q$

Suppose that  $X_{ia}, X_{ja} \stackrel{iid}{\sim} \mathcal{F}_X(\mu_x, \sigma_x^2)$  and define  $Z_a^q = |d_{ij}(a)|^q = |X_{ia} - X_{ja}|^q$ , where  $a \in \mathcal{A}$  and  $|\mathcal{A}| = p$ .

- (i) Suppose that  $X_{ja} = X_{ia} - (Z_a^q)^{1/q}$ . Based on the iid assumption for  $X_{ia}$  and  $X_{ja}$ , it follows that the joint density function  $g^{(1)}$  of  $X_{ia}$  and  $Z_a^q$  is given by

$$\begin{aligned} g^{(1)}(x_{ia}, z_a) &= f_X(x_{ia}, x_{ja}) \left| \frac{\partial x_{ja}}{\partial z_a} \right| \\ &= f_X(x_{ia}) f_X(x_{ja}) \left| \frac{-1}{q} (z_a^q)^{\frac{1}{q}-1} \right| \\ &= \frac{1}{q (z_a^q)^{1-\frac{1}{q}}} f_X(x_{ia}) f_X \left( x_{ia} - (z_a^q)^{1/q} \right), \quad z_a > 0 \end{aligned} \tag{15}$$

The density function  $f_{Z_a^q}^{(1)}$  of  $Z_a^q$  is then defined as

$$\begin{aligned} f_{Z_a^q}^{(1)}(z_a^q) &= \int_{-\infty}^{\infty} g^{(1)}(x_{ia}, z_a^q) dx_{ia} \\ &= \frac{1}{q (z_a^q)^{1-\frac{1}{q}}} \int_{-\infty}^{\infty} f_X(x_{ia}) f_X \left( x_{ia} - (z_a^q)^{1/q} \right) dx_{ia}, \quad z_a > 0 \end{aligned} \tag{16}$$

- (ii) Suppose that  $X_{ja} = X_{ia} + (Z_a^q)^{1/q}$ . Based on the iid assumption for  $X_{ia}$  and  $X_{ja}$ , it follows that the joint density function  $g^{(2)}$  of  $X_{ia}$  and  $Z_a^q$  is given by

$$\begin{aligned}
g^{(2)}(x_{ia}, z_a) &= f_X(x_{ia}, x_{ja}) \left| \frac{\partial x_{ja}}{\partial z_a} \right| \\
&= f_X(x_{ia}) f_X(x_{ja}) \left| \frac{1}{q} (z_a^q)^{\frac{1}{q}-1} \right| \\
&= \frac{1}{q (z_a^q)^{1-\frac{1}{q}}} f_X(x_{ia}) f_X \left( x_{ia} - (z_a^q)^{1/q} \right), \quad z_a > 0
\end{aligned} \tag{17}$$

The density function  $f_{Z_a^q}^{(2)}$  of  $Z_a^q$  is then defined as

$$\begin{aligned}
f_{Z_a^q}^{(2)}(z_a^q) &= \int_{-\infty}^{\infty} g^{(2)}(x_{ia}, z_a^q) dx_{ia} \\
&= \frac{1}{q (z_a^q)^{1-\frac{1}{q}}} \int_{-\infty}^{\infty} f_X(x_{ia}) f_X \left( x_{ia} + (z_a^q)^{1/q} \right) dx_{ia}, \quad z_a > 0
\end{aligned} \tag{18}$$

Let  $F_{Z_a^q}$  denote the distribution function of the random variable  $Z_a^q$ . Furthermore, we define the events  $E^{(1)}$  and  $E^{(2)}$  as

$$E^{(1)} = \{|X_{ia} - X_{ja}|^q \leq z_a^q | X_{ja} = X_{ia} - (Z_a^q)^{1/q}\} \tag{19}$$

and

$$E^{(2)} = \{|X_{ia} - X_{ja}|^q \leq z_a^q | X_{ja} = X_{ia} + (Z_a^q)^{1/q}\}. \tag{20}$$

Then it follows from fundamental rules of probability that

$$\begin{aligned}
F_{Z_a^q}(z_a^q) &= P[Z_a^q \leq z_a^q] \\
&= P[|X_{ia} - X_{ja}|^q \leq z_a^q] \\
&= P[E^{(1)} \cup E^{(2)}] \\
&= P[E^{(1)}] + P[E^{(2)}] - P[E^{(1)} \cap E^{(2)}] \\
&= P[E^{(1)}] + P[E^{(2)}] \\
&= \int_{-\infty}^{z_a^q} f_{Z_a^q}^{(1)}(t) dt + \int_{-\infty}^{z_a^q} f_{Z_a^q}^{(2)}(t) dt \\
&= \int_{-\infty}^{z_a^q} \left( f_{Z_a^q}^{(1)}(t) + f_{Z_a^q}^{(2)}(t) \right) dt \\
&= \frac{1}{q (z_a^q)^{1-\frac{1}{q}}} \int_{-\infty}^{z_a^q} \left( \int_{-\infty}^{\infty} f_X(x_{ia}) [f_X(x_{ia} - t) + f_X(x_{ia} + t)] dx_{ia} \right) dt, \quad z_a > 0
\end{aligned} \tag{21}$$

It follows directly from the result in Eq. 21 that the density function of the random variable  $Z_a^q$  is given by

$$\begin{aligned}
f_{Z_a^q}(z_a^q) &= \frac{\partial}{\partial z_a^q} F_{Z_a^q}(z_a^q) \\
&= \frac{1}{q (z_a^q)^{1-\frac{1}{q}}} \int_{-\infty}^{\infty} f_X(x_{ia}) \left[ f_X \left( x_{ia} - (z_a^q)^{1/q} \right) + f_X \left( x_{ia} + (z_a^q)^{1/q} \right) \right] dx_{ia}
\end{aligned} \tag{22}$$

where  $z_a > 0$ .

Using Eq. 22, we can compute the mean and variance of the magnitude difference as 171

$$\mu_{z^q} = \int_{-\infty}^{\infty} z_a^q f_{Z^q}(z_a^q) dz_a^q \quad (23)$$

and 172

$$\sigma_{z^q}^2 = \int_{-\infty}^{\infty} (z_a^q)^2 f_{Z^q}(z_a^q) dz_a^q - \mu_{z^q}^2. \quad (24)$$

It follows immediately from Eqs. 23 and 24 and the Classical Central Limit Theorem (CCLT) that 173  
174

$$\left(D_{ij}^{(q)}\right)^q = \sum_{a \in \mathcal{A}} Z_a^q = \sum_{a \in \mathcal{A}} |X_{ia} - X_{ja}|^q \sim \mathcal{N}(\mu_{z^q} p, \sigma_{z^q}^2 p) \quad (25)$$

Applying the result given in Eq. 6, the distribution of  $D_{ij}^{(q)}$  is given by 175

$$D_{ij}^{(q)} \sim \mathcal{N}\left((\mu_{z^q} p)^{1/q}, \frac{\sigma_{z^q}^2 p}{q^2 (\mu_{z^q} p)^{2(1-\frac{1}{q})}}\right), \quad \mu_{z^q} > 0 \quad (26)$$

with improved estimate of the mean for  $q = 2$  given by Eq. 7. 176

### 3.1.1 Standard normal data 177

If  $X_{ia}, X_{ja} \stackrel{iid}{\sim} \mathcal{N}(0, 1)$ , then the marginal density functions with respect to  $X$  for  $X_{ia}$ ,  $X_{ia} - (Z_a^q)^{1/q}$ , and  $X_{ia} + (Z_a^q)^{1/q}$  are defined as 178  
179

$$f_X(x_{ia}) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x_{ia}^2}, \quad (27)$$

$$f_X\left(x_{ia} - (z_a^q)^{1/q}\right) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x_{ia} - (z_a^q)^{1/q})^2}, \quad z_a > 0, \text{ and} \quad (28)$$

$$f_X\left(x_{ia} + (z_a^q)^{1/q}\right) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x_{ia} + (z_a^q)^{1/q})^2}, \quad z_a > 0 \quad (29)$$

Substituting the results given by Eqs. 27-29 into Eq. 22 and completing the square on  $x_{ia}$  in the exponents, we have 180  
181

$$f_{Z^q}(z_a^q) = \frac{1}{2q\pi (z_a^q)^{1-\frac{1}{q}}} e^{-\frac{1}{4}(z_a^q)^{2/q}} \int_{-\infty}^{\infty} \left( e^{-\frac{1}{2}[\sqrt{2}x_{ia} - \frac{\sqrt{2}}{2}(z_a^q)^{1/q}]^2} + e^{-\frac{1}{2}[\sqrt{2}x_{ia} + \frac{\sqrt{2}}{2}(z_a^q)^{1/q}]^2} \right) dx_{ia} \quad (30)$$

$$= \frac{1}{2q\sqrt{\pi} (z_a^q)^{1-\frac{1}{q}}} e^{-\frac{1}{4}(z_a^q)^{2/q}} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \left( e^{-\frac{1}{2}u^2} + e^{-\frac{1}{2}u^2} \right) du \quad (31)$$

$$= \frac{1}{2q\sqrt{\pi} (z_a^q)^{1-\frac{1}{q}}} e^{-\frac{1}{4}(z_a^q)^{2/q}} (1 + 1) \quad (32)$$

$$= \frac{1}{q\sqrt{\pi}} (z_a^q)^{\frac{1}{q}-1} e^{-\frac{1}{4}(z_a^q)^{2/q}} \quad (33)$$

$$= \frac{\frac{2}{q}}{(2q)^{1/q} \Gamma\left(\frac{1}{q}\right)} (z_a^q)^{\frac{1}{q}-1} e^{-\left(\frac{z_a^q}{2q}\right)^{2/q}} \quad (34)$$

The density function given by Eq. 30 is a Generalized Gamma density with parameters  $b = \frac{2}{q}$ ,  $c = 2^q$ , and  $d = \frac{1}{q}$ . This distribution has mean and variance given by

$$\begin{aligned}\mu_{z^q} &= \frac{c\Gamma\left(\frac{d+1}{b}\right)}{\Gamma\left(\frac{d}{b}\right)} \\ &= \frac{2^q\Gamma\left(\frac{q+1}{2}\right)}{\sqrt{\pi}}\end{aligned}\quad (35)$$

and

$$\begin{aligned}\sigma_{z^q}^2 &= c^2 \left[ \frac{\Gamma\left(\frac{d+2}{b}\right)}{\Gamma\left(\frac{d}{b}\right)} - \left( \frac{\Gamma\left(\frac{d+1}{b}\right)}{\Gamma\left(\frac{d}{b}\right)} \right)^2 \right] \\ &= 4^q \left[ \frac{\Gamma\left(q + \frac{1}{2}\right)}{\sqrt{\pi}} - \frac{\Gamma^2\left(\frac{1}{2}q + \frac{1}{2}\right)}{\pi} \right]\end{aligned}\quad (36)$$

By linearity of the expected value and variance operators under the iid assumption, Eqs. 35 and 36 allow the  $p$ -dimensional mean and variance of the  $D_{ij}^{(q)}$  distribution to be computed directly as

$$\mu_{(D_{ij}^{(q)})^q} = \mathbb{E} \left[ (D_{ij}^{(q)})^q \right] = \mathbb{E} \left( \sum_{a \in \mathcal{A}} Z_a^q \right) = \sum_{a \in \mathcal{A}} \mathbb{E} (Z_a^q) = \sum_{a \in \mathcal{A}} \frac{2^q\Gamma\left(\frac{q+1}{2}\right)}{\sqrt{\pi}} = \frac{2^q\Gamma\left(\frac{q+1}{2}\right)}{\sqrt{\pi}} p \quad (37)$$

and

$$\begin{aligned}\sigma_{(D_{ij}^{(q)})^q}^2 &= \text{Var} \left[ (D_{ij}^{(q)})^q \right] = \text{Var} \left( \sum_{a \in \mathcal{A}} Z_a^q \right) \\ &= \sum_{a \in \mathcal{A}} \text{Var} (Z_a^q) \\ &= \sum_{a \in \mathcal{A}} 4^q \left[ \frac{\Gamma\left(q + \frac{1}{2}\right)}{\sqrt{\pi}} - \frac{\Gamma^2\left(\frac{1}{2}q + \frac{1}{2}\right)}{\pi} \right] \\ &= 4^q \left[ \frac{\Gamma\left(q + \frac{1}{2}\right)}{\sqrt{\pi}} - \frac{\Gamma^2\left(\frac{1}{2}q + \frac{1}{2}\right)}{\pi} \right] p\end{aligned}\quad (38)$$

Therefore, the asymptotic distribution of  $D_{ij}^{(q)}$  for standard normal data is

$$\mathcal{N} \left( \left( 2 \frac{\Gamma\left(\frac{q+1}{2}\right)}{\sqrt{\pi}} p \right)^{1/q}, \frac{4^q p}{q^2 \left( \frac{2^q\Gamma\left(\frac{1}{2}q + \frac{1}{2}\right)}{\sqrt{\pi}} p \right)^{2(1-\frac{1}{q})}} \left[ \frac{\Gamma\left(q + \frac{1}{2}\right)}{\sqrt{\pi}} - \frac{\Gamma^2\left(\frac{1}{2}q + \frac{1}{2}\right)}{\pi} \right] \right) \quad (39)$$

### 3.1.2 Standard uniform data

If  $X_{ia}, X_{ja} \stackrel{iid}{\sim} \mathcal{U}(0, 1)$ , then the marginal density functions with respect to  $X$  for  $X_{ia}$ ,  $X_{ia} - (Z_a^q)^{1/q}$ , and  $X_{ia} + (Z_a^q)^{1/q}$  are defined as

$$f_X(x_{ia}) = 1, \quad 0 \leq x_{ia} \leq 1 \quad (40)$$

$$f_X \left( x_{ia} - (z_a^q)^{1/q} \right) = 1, \quad 0 \leq x_{ia} - (z_a^q)^{1/q} \leq 1, \text{ and} \quad (41)$$

$$f_X \left( x_{ia} + (z_a^q)^{1/q} \right) = 1, \quad 0 \leq x_{ia} + (z_a^q)^{1/q} \leq 1. \quad (42)$$

Substituting the results given by Eqs. 40-42 into Eq. 22, we have

$$\begin{aligned} f_{Z^q}(z_a^q) &= \frac{1}{q(z_a^q)^{1-\frac{1}{q}}} \int_{-\infty}^{\infty} f_X(x_{ia}) \left[ f_X \left( x_{ia} - (z_a^q)^{1/q} \right) + f_X \left( x_{ia} + (z_a^q)^{1/q} \right) \right] dx_{ia}, \\ & \quad 0 < z_a \leq 1 \\ &= \frac{1}{q(z_a^q)^{1-\frac{1}{q}}} \int_0^1 [f_X(x_{ia} - (z_a^q)^{1/q}) + f_X(x_{ia} + (z_a^q)^{1/q})] dx_{ia}, \quad 0 < z_a \leq 1 \\ &= \frac{1}{q(z_a^q)^{1-\frac{1}{q}}} \int_{(z_a^q)}^1 1 dx_{ia} + \int_0^{1-(z_a^q)} 1 dx_{ia}, \quad 0 < z_a \leq 1 \\ &= \frac{1}{q(z_a^q)^{1-\frac{1}{q}}} [(1 - (z_a^q)) + (1 - (z_a^q))], \quad 0 < z_a \leq 1 \\ &= \frac{1}{q} \cdot 2(z_a^q)^{\frac{1}{q}-1} [1 - (z_a^q)^{1/q}]^{2-1}, \quad 0 < z_a \leq 1 \end{aligned} \quad (43)$$

The density given by Eq. 43 is a Kumaraswamy density with parameters  $b = \frac{1}{q}$  and  $c = 2$  with moment generating function (MGF) given by

$$\begin{aligned} M_n &= \frac{c\Gamma(1 + \frac{n}{b})\Gamma(c)}{\Gamma(1 + c + \frac{n}{b})} \\ &= \frac{2}{(nq + 2)(nq + 1)} \end{aligned} \quad (44)$$

Using the MGF given by Eq. 44, the mean and variance of  $Z_a^q$  are computed as

$$\mu_{z^q} = \frac{2}{(q + 2)(q + 1)} \quad (45)$$

and

$$\sigma_{z^q}^2 = \frac{1}{(q + 1)(2q + 1)} - \left( \frac{2}{(q + 2)(q + 1)} \right)^2 \quad (46)$$

By linearity of the expected value and variance operators under the iid assumption, Eqs. 47 and 48 allow the  $p$ -dimensional mean and variance of the  $\left( D_{ij}^{(q)} \right)^q$  distribution to be computed directly as

$$\begin{aligned} \mu_{\left( D_{ij}^{(q)} \right)^q} &= \mathbb{E} \left[ \left( D_{ij}^{(q)} \right)^q \right] = \mathbb{E} \left( \sum_{a \in \mathcal{A}} Z_a^q \right) \\ &= \sum_{a \in \mathcal{A}} \mathbb{E}(Z_a^q) \\ &= \sum_{a \in \mathcal{A}} \frac{2}{(q + 2)(q + 1)} \\ &= \frac{2p}{(q + 2)(q + 1)} \end{aligned} \quad (47)$$

and

201

$$\begin{aligned}
\sigma^2_{\left(D_{ij}^{(q)}\right)^q} &= \text{Var} \left[ \left( D_{ij}^{(q)} \right)^q \right] = \text{Var} \left( \sum_{a \in \mathcal{A}} Z_a^q \right) \\
&= \sum_{a \in \mathcal{A}} \text{Var} (Z_a^q) \\
&= \sum_{a \in \mathcal{A}} \left[ \frac{1}{(q+1)(2q+1)} - \left( \frac{2}{(q+2)(q+1)} \right)^2 \right] \\
&= \left[ \frac{1}{(q+1)(2q+1)} - \left( \frac{2}{(q+2)(q+1)} \right)^2 \right] p
\end{aligned} \tag{48}$$

Therefore, the asymptotic distribution of  $D_{ij}^{(q)}$  for standard uniform data is

202

$$\begin{aligned}
&\mathcal{N} \left( \left( \frac{2p}{(q+2)(q+1)} \right)^{1/q}, \right. \\
&\quad \left. \frac{p}{q^2 \left( \frac{2p}{(q+2)(q+1)} \right)^{2(1-\frac{1}{q})}} \left[ \frac{1}{(q+1)(2q+1)} - \left( \frac{2}{(q+2)(q+1)} \right)^2 \right] \right).
\end{aligned} \tag{49}$$

**Table 1.** Summary of asymptotic distance distributions for common data types. Metrics with subscripts M and E represent Manhattan and Euclidean, respectively. Metrics with superscript \* represent a deviation from the standard metric by attribute range normalization. The function  $\Phi^{-1}(x)$  denotes the standard normal quantile function, where  $x \in (0, 1)$ .

Type	Mean	Variance
$\mathcal{N}(0, 1) - \mathbf{d}_M$	$\frac{2p}{\sqrt{\pi}}$	$\frac{2p(\pi - 2)}{\pi}$
$\mathcal{N}(0, 1) - \mathbf{d}_M^*$	$\frac{p}{\sqrt{\pi}\mu(m)}$ where $\mu(m) = \frac{\log(\log(2))}{\Phi^{-1}(\frac{1}{m})} - \Phi^{-1}(\frac{1}{m})$	$\frac{p(\pi - 2)}{2\pi\mu^2(m)}$ where $\mu(m) = \frac{\log(\log(2))}{\Phi^{-1}(\frac{1}{m})} - \Phi^{-1}(\frac{1}{m})$
$\mathcal{N}(0, 1) - \mathbf{d}_E$	$\sqrt{2p - 1}$	1
$\mathcal{N}(0, 1) - \mathbf{d}_E^*$	$\frac{\sqrt{2p - 1}}{2\mu(m)}$ where $\mu(m) = \frac{\log(\log(2))}{\Phi^{-1}(\frac{1}{m})} - \Phi^{-1}(\frac{1}{m})$	$\frac{2\log(m)}{\pi^2 + 12\mu^2(m)\log(m)}$ where $\mu(m) = \frac{\log(\log(2))}{\Phi^{-1}(\frac{1}{m})} - \Phi^{-1}(\frac{1}{m})$
$\mathcal{U}(0, 1) - \mathbf{d}_M$	$\frac{p}{3}$	$\frac{p}{18}$
$\mathcal{U}(0, 1) - \mathbf{d}_M^*$	$\frac{(m+1)p}{3(m-1)}$	$\frac{(m^3 - 18m^2 - 5m + 2)p}{18(m^3 + m^2 + 2)(m-1)^2}$
$\mathcal{U}(0, 1) - \mathbf{d}_E$	$\sqrt{\frac{p}{6} - \frac{7}{120}}$	$\frac{7}{120}$
$\mathcal{U}(0, 1) - \mathbf{d}_E^*$	$\sqrt{\frac{p}{6} - \frac{7}{120}} \left( \frac{m+1}{m-1} \right)$	$\frac{7(m+1)^2(m+2)}{120(m^3 + m^2 + 2)}$

**Table 2.** Summary of asymptotic distance distributions for rs-fMRI and GWAS data. Metrics with superscript \* represent a deviation from the standard metric by attribute range normalization. The function  $\Phi^{-1}(x)$  denotes the standard normal quantile function, where  $x \in (0, 1)$ .

Type	Mean	Variance
rs-fMRI ( $\mathbf{d}_{\text{ROI}}$ )	$\frac{2p(p-1)}{\sqrt{\pi(p-3)}}$	$\frac{4(\pi-2)p(p-1)}{\pi(p-3)}$
rs-fMRI ( $\mathbf{d}_{\text{ROI}}^*$ )	$\frac{2p(p-1)}{\mu(m,p)\sqrt{\pi(p-3)}}$ where $\mu(m,p) = \frac{1}{\sqrt{p-3}}\Phi^{-1}\left(1 - \frac{1}{m(p-1)}\right)$	$\frac{2[6(p-3)\mu^2(m,p)\log[m(p-1)](\pi-2) - \pi^2]p(p-1)}{\pi(p-3)\mu^2(m,p)(\pi^2 + 12(p-3)\mu^2(m,p)\log[m(p-1)])}$ where $\mu(m,p) = \frac{1}{\sqrt{p-3}}\Phi^{-1}\left(1 - \frac{1}{m(p-1)}\right)$
GWAS ( $\mathbf{d}_{\text{GM}}$ )	$2 \sum_{a=1}^p F(a)$ where $F(a) = [2(1-f_a)^3 f_a + 2f_a^3(1-f_a) + (1-f_a)^2 f_a^2]$ , and $f_a$ is the probability of a minor allele at locus $a$ .	$2 \sum_{a=1}^p F(a)[1 - 2F(a)]$ where $F(a) = [2(1-f_a)^3 f_a + 2f_a^3(1-f_a) + (1-f_a)^2 f_a^2]$ , and $f_a$ is the probability of a minor allele at locus $a$ .
GWAS ( $\mathbf{d}_{\text{AM}}$ )	$2 \sum_{a=1}^p F(a)$ where $F(a) = [(1-f_a)^3 f_a + f_a^3(1-f_a) + (1-f_a)^2 f_a^2]$ , and $f_a$ is the probability of a minor allele at locus $a$ .	$\sum_{a=1}^p [G(a) - 4F^2(a)]$ where $F(a) = [(1-f_a)^3 f_a + f_a^3(1-f_a) + f_a^3(1-f_a) + (1-f_a)^2 f_a^2]$ , $G(a) = [(1-f_a)^3 f_a + f_a^3(1-f_a) + 2(1-f_a)^2 f_a^2]$ , and $f_a$ is the probability of a minor allele at locus $a$ .
GWAS ( $\mathbf{d}_{\text{TIV}}$ )	$(\gamma_0 + \gamma_2 + 2\gamma_1) \sum_{a=1}^p F(a) + \left[\frac{3}{2}(\gamma_0 + \gamma_2) + 2\gamma_1\right] \sum_{a=1}^p G(a)$ where $F(a) = [(1-f_a)^3 f_a + f_a^3(1-f_a)]$ and $G(a) = (1-f_a)^2 f_a^2$ , $f_a$ is the probability of a minor allele at locus $a$ , and $\gamma_0, \gamma_1$ , and $\gamma_2$ are probabilities of PuPu, PuPy, and PyPy, respectively, at locus $a$ .	$\left[\frac{1}{4}(\gamma_0 + \gamma_2) + \gamma_1\right] \sum_{a=1}^p F(a) + \left[\frac{9}{8}(\gamma_0 + \gamma_2) + 2\gamma_1\right] \sum_{a=1}^p G(a)$ $+ \sum_{a=1}^p \left[(\gamma_0 + \gamma_2 + 2\gamma_1)F(a) + \left[\frac{3}{2}(\gamma_0 + \gamma_2) + 2\gamma_1\right] G(a)\right]^2$ where $F(a) = [(1-f_a)^3 f_a + f_a^3(1-f_a)]$ and $G(a) = (1-f_a)^2 f_a^2$ , $f_a$ is the probability of a minor allele at locus $a$ , and $\gamma_0, \gamma_1$ , and $\gamma_2$ are probabilities of PuPu, PuPy, and PyPy, respectively, at locus $a$ .

## 4 Optimal neighborhood parameters for detecting effects

k or  $\alpha$ . Balancing blessing and curse of dimensionality.

## 5 ICA?

Using same interaction, increase background noise genes to see degrading of A and B Relief importance because of curse of dimensionality (sparseness).

## References

1. Ryan J. Urbanowicz, Randal S. Olson, Peter Schmitt, Melissa Meeker, and Jason H. Moore. Benchmarking relief-based feature selection methods for



- bioinformatics data mining. *Journal of Biomedical Informatics*, 85:168–188, 2018. 212
2. Trang T Le, Ryan J Urbanowicz, Jason H Moore, and Brett A McKinney. 213  
Statistical inference relief (stir) feature selection. *Bioinformatics*, page bty788, 214  
2018. 215
  3. Larry Wasserman. *All of Statistics: A Concise Course in Statistical Inference*. 216  
Springer, New York, NY, 2004. 217
  4. Delaney Granizo-Mackenzie and Jason H Moore. Multiple threshold spatially 218  
uniform relieff for the genetic analysis of complex human diseases. In *European* 219  
*Conference on Evolutionary Computation, Machine Learning and Data Mining in* 220  
*Bioinformatics*, pages 1–10. Springer, 2013. 221