

Implementation and Evaluation of a Docking-Rescoring Method using Molecular Footprint Comparisons

Trent E. Balius, Sudipto Mukherjee, and Robert C. Rizzo

Supplemental Material

The following corollary and theorem prove the strong functional relationship (see Figure 6) between normalized Euclidean distance and standard Pearson Correlation methods used to calculate footprint similarity described in the text. Note that the relationship is strong because the mean of the footprint vectors is usually close to zero. If the mean is not close to zero (i.e. when using a threshold-based footprint) then the functional relationship will be weaker.

Corollary 1: if \vec{u} and \vec{v} are unit vectors, then $\|\vec{u} - \vec{v}\| = \sqrt{2(1 - \cos(\theta))}$, where θ is the angle between \vec{u} and \vec{v} .

Proof:

Let \vec{u} and \vec{v} be unit vectors. Then,

$$\begin{aligned}\cos(\theta) &= \frac{\vec{u} \cdot \vec{v}}{\|\vec{u}\| \|\vec{v}\|} = \vec{u} \cdot \vec{v} \quad (\because \|\vec{u}\| = 1) \\ \|\vec{u} - \vec{v}\| &= \sqrt{\sum (u_i - v_i)^2} = \sqrt{\sum u_i^2 + \sum v_i^2 - 2\sum u_i v_i} \\ &= \sqrt{1 + 1 - 2\sum u_i v_i} \quad (\because \|\vec{u}\|^2 = \sum u_i^2 = 1) \\ &= \sqrt{2(1 - \vec{u} \cdot \vec{v})}\end{aligned}$$

Therefore, $\|\vec{u} - \vec{v}\| = \sqrt{2(1 - \cos(\theta))}$

Theorem: $d_{norm} \approx \sqrt{2(1 - r)}$ for two vectors x and y , whose means are close to zero, where d_{norm} is the normalized Euclidean distance and r is the Pearson Correlation Coefficient (both calculated between the two vectors).

Proof:

$$r = \frac{\text{cov}(\vec{x}, \vec{y})}{\sqrt{\text{var}(\vec{x})} \sqrt{\text{var}(\vec{y})}}$$

The correlation coefficient can also be thought of as the cosine of the angle formed between the mean-modulated vectors \vec{x}^μ and \vec{y}^μ where $\vec{x}^\mu = [x_i - \mu_x]$ and $\vec{y}^\mu = [y_i - \mu_y]$, and μ represents the mean of each vector.

$$r = \frac{\vec{x}^\mu \cdot \vec{y}^\mu}{\|\vec{x}^\mu\| \|\vec{y}^\mu\|} = \cos(\theta)$$

Note that

$$r = \cos(\theta) = \frac{\bar{x}^\mu \cdot \bar{y}^\mu}{\|\bar{x}^\mu\| \|\bar{y}^\mu\|} \approx \frac{\bar{x} \cdot \bar{y}}{\|\bar{x}\| \|\bar{y}\|} = \cos(\theta^*)$$

The mean of a footprint is normally close to zero since most footprint entries are close to zero so this is a reasonable approximation.

The normalized Euclidian distance (d_{norm}) is defined as:

$$d_{norm} = \|\vec{\chi} - \vec{\gamma}\|$$

where $\vec{\chi} = \frac{\bar{x}}{\|\bar{x}\|}$, and $\vec{\gamma} = \frac{\bar{y}}{\|\bar{y}\|}$

$$d_{norm} = \sqrt{2(1 - \cos(\theta^*))} \approx \sqrt{2(1 - r)}$$

There is a relationship between the $\cos(\theta^*)$ and normalized Euclidean (d_{norm}) by **Corollary 1** and because the angle between two vectors is the same as that between their unit vectors.

Therefore, the approximate relationship is demonstrated.