

PCC	$r = \frac{1}{n} \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{\sigma_x} \right) \left( \frac{y_i - \bar{y}}{\sigma_y} \right)$
Uncentered Correlation Coefficient	(same as PCC, except that sample means are set to 0)
Frequency Dot Product (FDP) (note: we've been calling this idf_unnorm)	$FDP = \sum_{i=1}^n \left[ (x_i)(y_i) \left( \log_2 \frac{D}{count_i} \right)^2 \right]$ <ul style="list-style-type: none"> <li>• <math>x_i</math> and <math>y_i</math> are the <math>i^{\text{th}}</math> phenotypes in the binary phenotype vectors <math>\vec{x}</math> and <math>\vec{y}</math>.</li> <li>• <math>n</math> = total number of phenotypes</li> <li>• <math>D</math> = total number of binary phenotype vectors</li> <li>• <math>count_i</math> = count of number of times the phenotype <math>i</math> appears in the data</li> </ul>
1.	
Inverse Document Frequency (IDF)	$IDF = \frac{FDP}{ \vec{x}  *  \vec{y} }$ <ul style="list-style-type: none"> <li>• IDF is the FDP normalized by the lengths of the two vectors.</li> </ul>
Euclidean Distance	$d = \sum_{i=1}^n (x_i - y_i)^2$
Jaccard Similarity Coefficient	$J = \frac{M_{11}}{M_{01} + M_{10} + M_{11}}$ <ul style="list-style-type: none"> <li>• <math>M_{11}</math> = # of times 1 is observed in both vectors</li> <li>• <math>M_{01} + M_{10}</math> = # of times 1 is observed in exactly one of the vectors in the pair.</li> </ul>
Mutual Information	$I(X; Y) = \sum_{y \in Y} \sum_{x \in X} p(x, y) \log \left( \frac{p(x, y)}{p(x)p(y)} \right)$
Residual IDF (RIDF)	$RIDF = IDF - \log_2 \frac{1}{1 - Poisson(0; \lambda_i)}$ <ul style="list-style-type: none"> <li>• RIDF is the difference between the actual IDF and the inverse document frequency predicted by a Poisson distribution.</li> <li>• <math>\lambda_i</math> is the Poisson parameter, the average number of occurrences of the phenotype in each phenotype vector.</li> </ul>

Table of measures tested and their mathematical formulas.