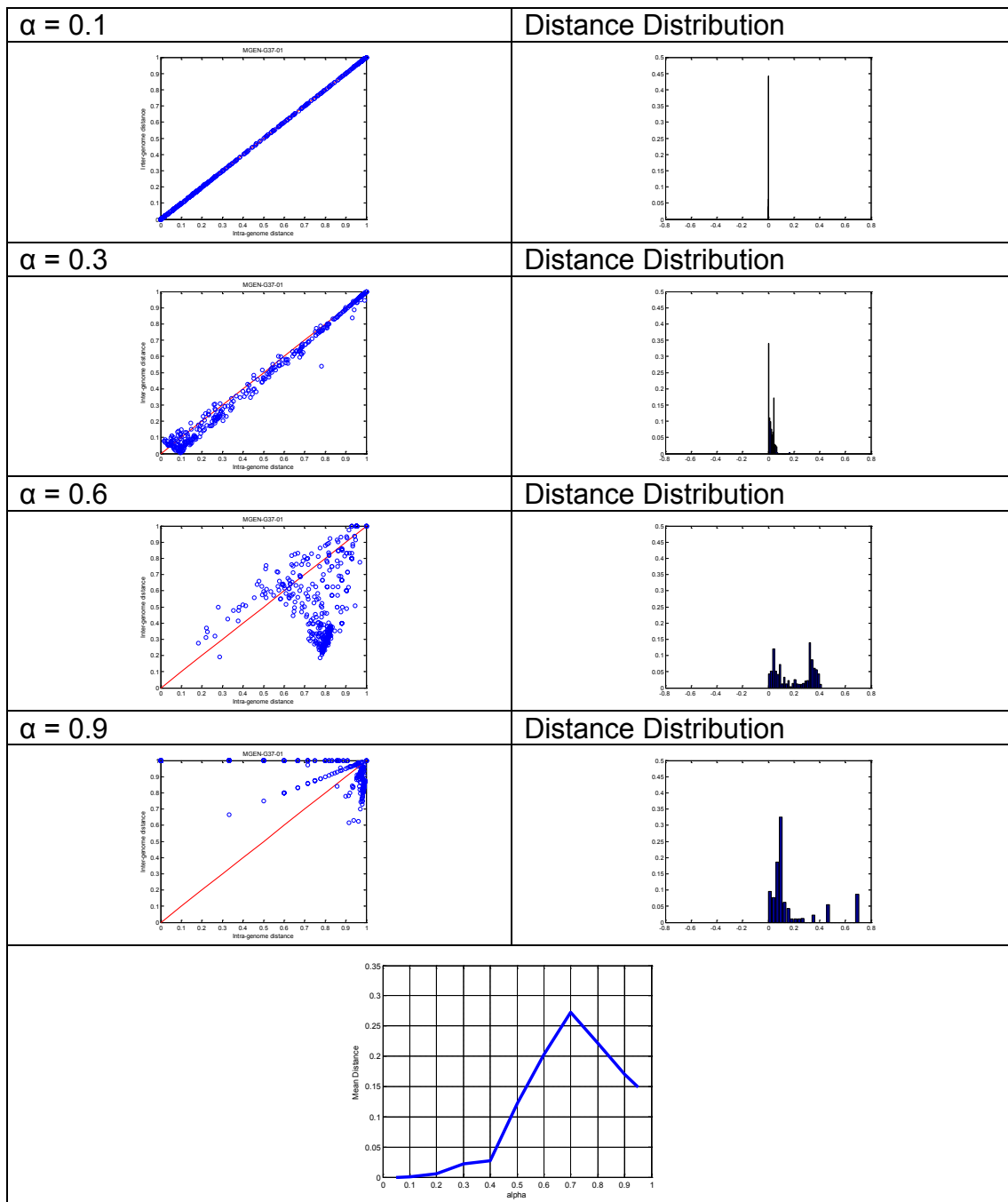
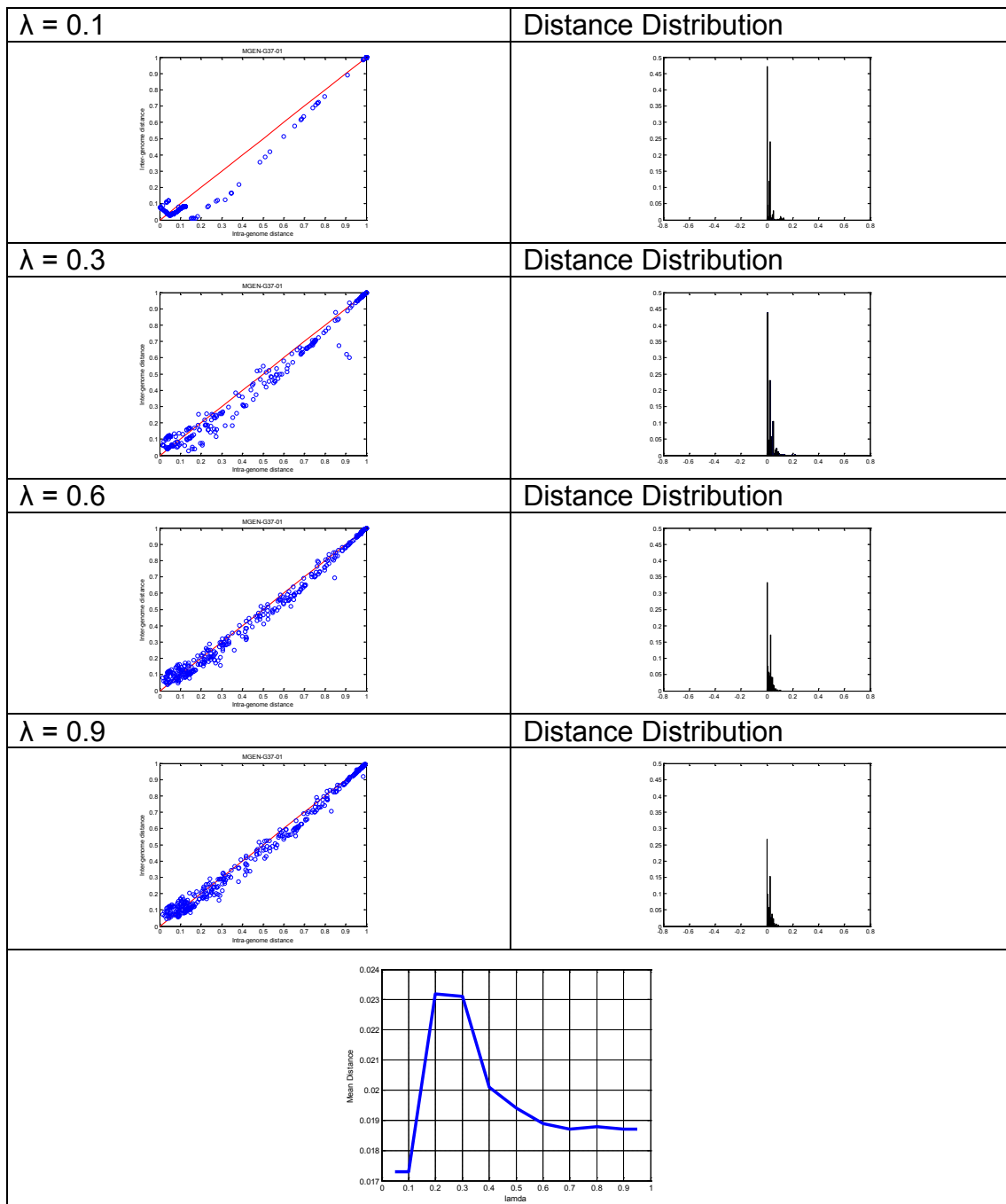


Parameter: α



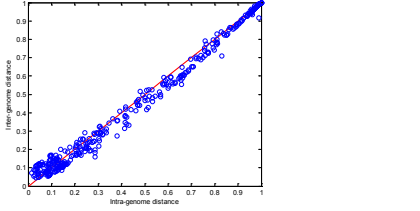
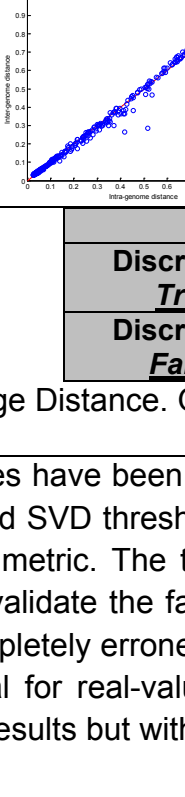
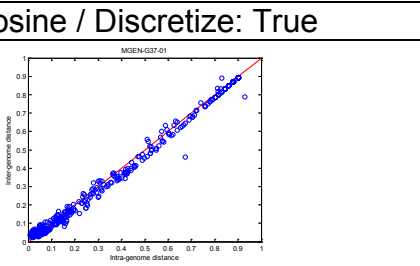
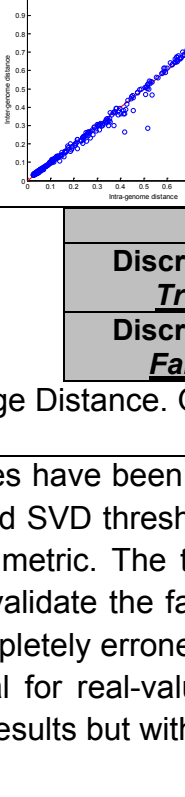
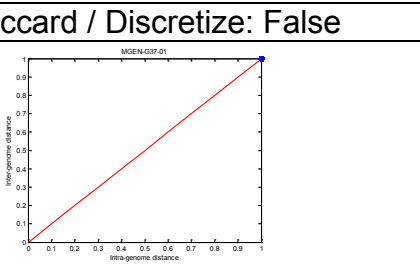
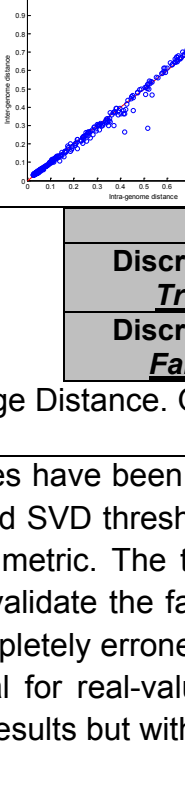
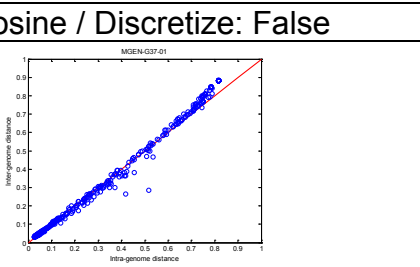
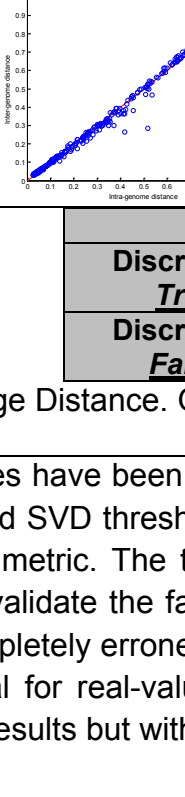
All Figures have been produced by keeping the parameters distance metric (*jaccard*) and SVD threshold λ (0.75) constant, and using different values for fuzzy threshold α . From the top figures, it is evident that the higher, more stringent values of α lead to a higher disparity of genes across the Distance Matrix; the expected presence of points on the main diagonal is lost. This is also evident by the average distance of each point from the main diagonal (bottom figure), where there is a sharp increase after 0.4. Therefore, aiming at the most flexible value of α , but without losing the on-diagonal presence of genes, the optimum range is between 0.3 and 0.4, hence the selection of 0.35 as our default α value.

Parameter: λ



All Figures have been produced by keeping the parameters distance metric (*jaccard*) and fuzzy threshold α (0.35) constant, and using different values for SVD threshold λ . From the top figures it is evident that the lower values of λ (i.e. low percentage of representation of the initial records) severely modify the expected distribution of points on the Distance matrix. This is also evident by the average distance of each point from the main diagonal (bottom figure), where the mean off-diagonal distance reaches a steady state around the range of 0.7 to 0.9, hence the selection of 0.75 as our default λ value (we need the lowest possible λ value in order to maximize the de-noising effect but also retain the distribution of the original records).

Parameter: Distance metric

<p>Metric: jaccard / Discretize: True</p> 	<p>Distance Distribution</p> 									
<p>Metric: cosine / Discretize: True</p> 	<p>Distance Distribution</p> 									
<p>Metric: jaccard / Discretize: False</p> 	<p>Distance Distribution</p> 									
<p>Metric: cosine / Discretize: False</p> 	<p>Distance Distribution</p> 									
<table border="1"> <thead> <tr> <th></th> <th>Jaccard</th> <th>Cosine</th> </tr> </thead> <tbody> <tr> <th>Discretize: <u>True</u></th> <td style="background-color: #90EE90;">0.0188</td> <td style="background-color: #C71585;">0.0149 (with NaN)</td> </tr> <tr> <th>Discretize: <u>False</u></th> <td style="background-color: #C71585;">0</td> <td style="background-color: #90EE90;">0.0137</td> </tr> </tbody> </table>			Jaccard	Cosine	Discretize: <u>True</u>	0.0188	0.0149 (with NaN)	Discretize: <u>False</u>	0	0.0137
	Jaccard	Cosine								
Discretize: <u>True</u>	0.0188	0.0149 (with NaN)								
Discretize: <u>False</u>	0	0.0137								
<p>Average Distance. Green color denotes optimal selection; the opposite is shown in red color.</p>										

All Figures have been produced by keeping the parameters fuzzy threshold α (0.35) and SVD threshold λ (0.75) constant, and using different values for the distance metric. The top figures, together with the summarizing table in the bottom, validate the fact the Jaccard metric is only applicable on binary data (it is completely erroneous on the real-value data), whereas the Cosine metric is optimal for real-value data and quite ill-suited for the binary data (it can provide results but with a fair amount of NaN values).