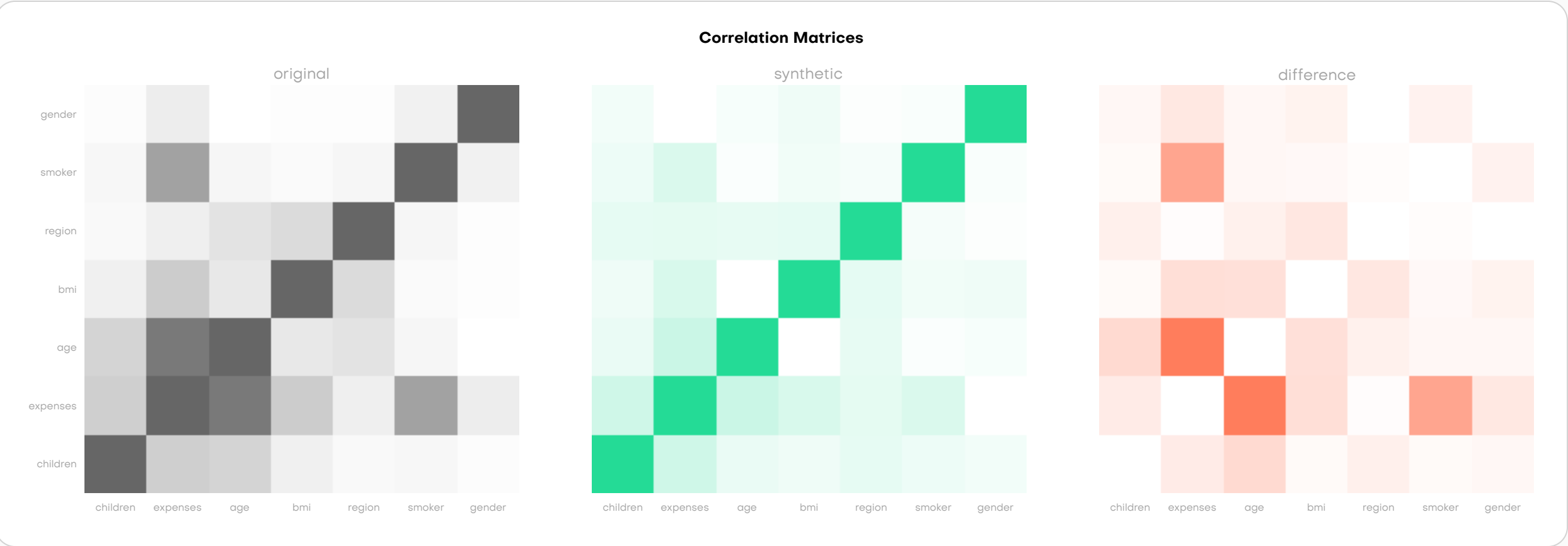


Model Report for `insurance_dataset`

Generated on 24 Jul 2024, 23:00 • Generator: [e9803b10-c368-4d5b-958f-22dfa1c13ce3](#)

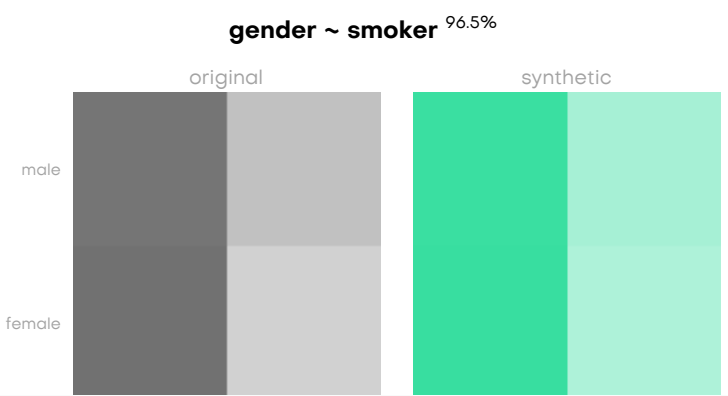
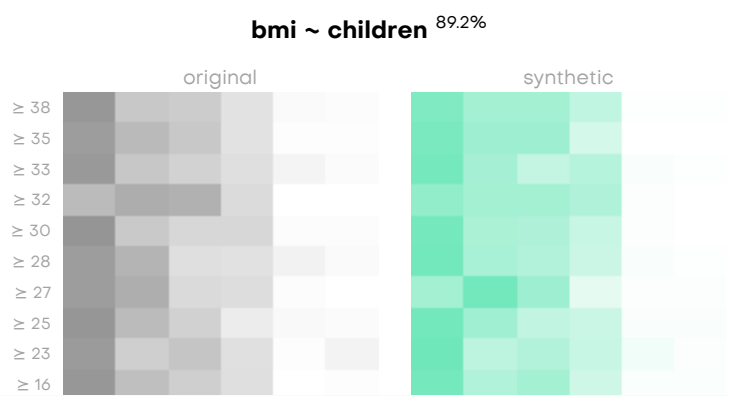
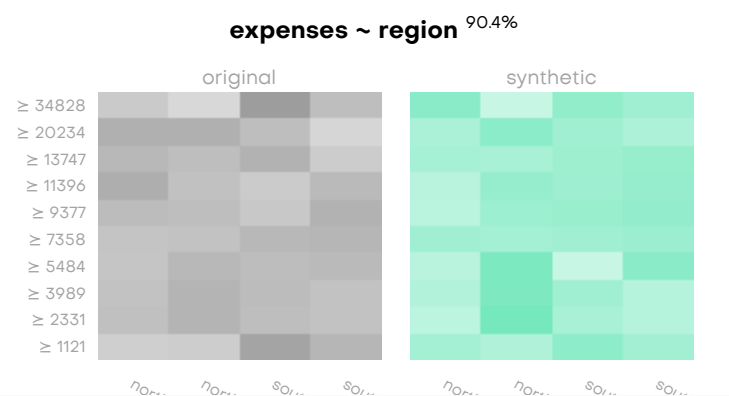
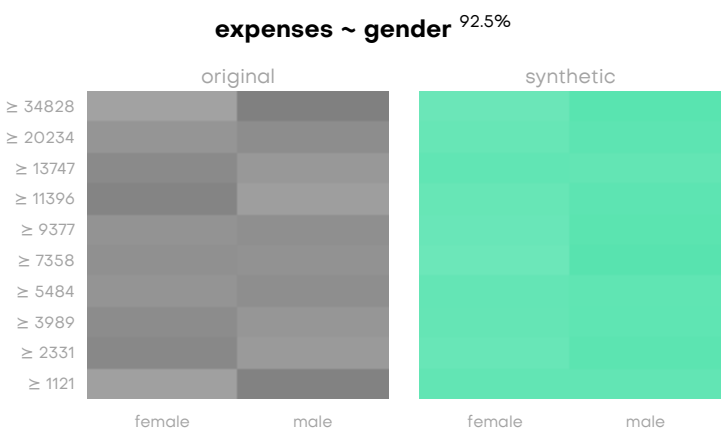
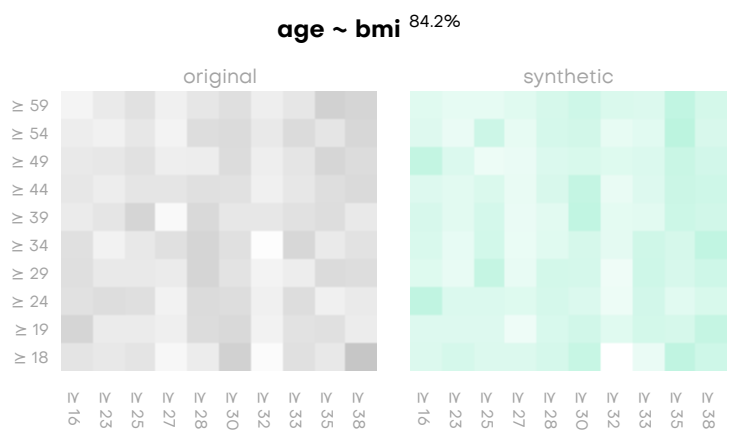
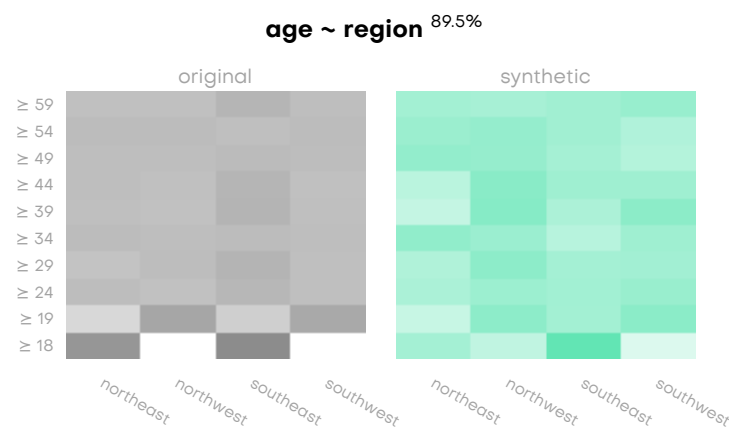
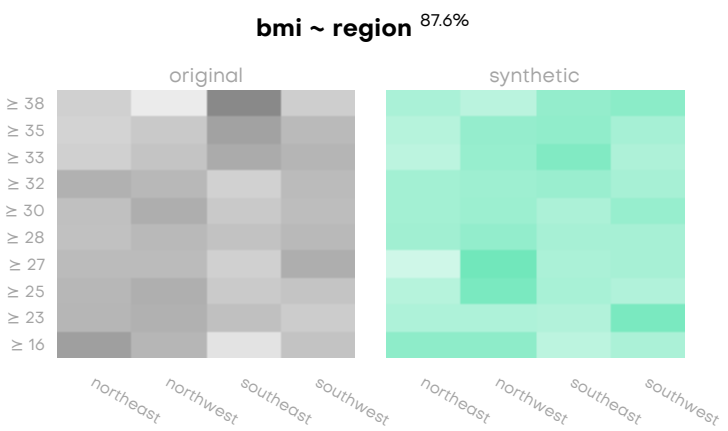
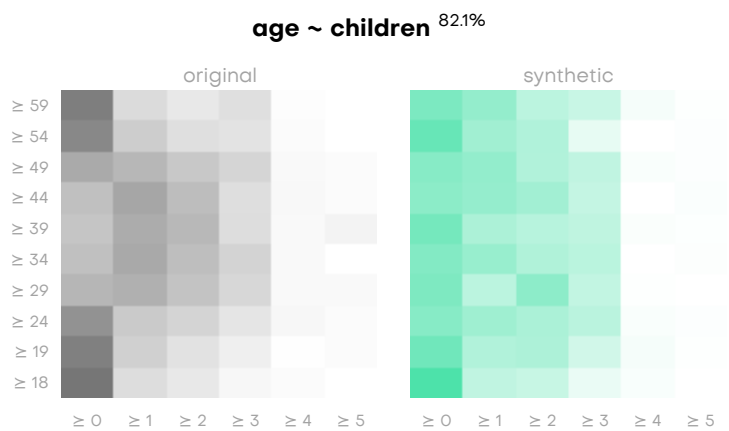
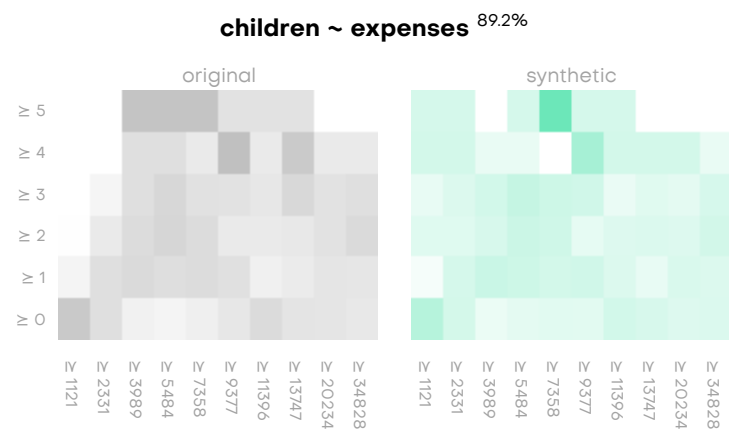
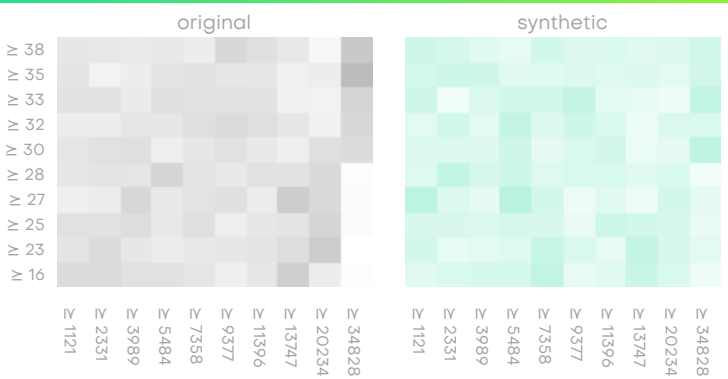
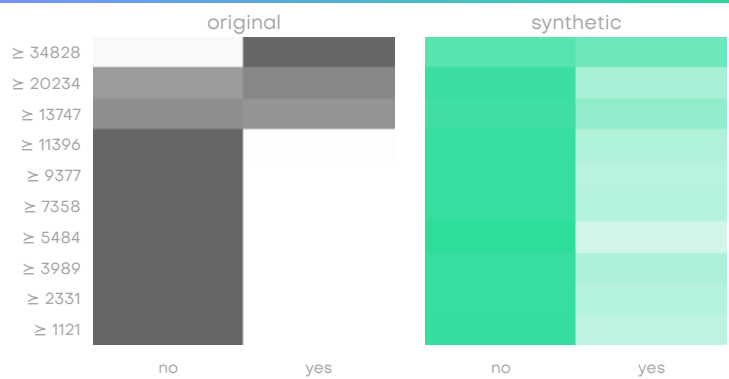
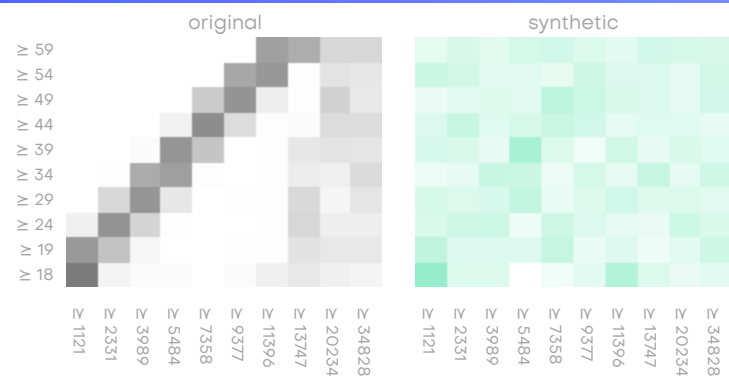
Dataset ⓘ	Original Samples	1,338
	Synthetic Samples	1,338
	Target Columns	7
Accuracy ⓘ	Univariate	95.8%
	Bivariate	88.2%
92.0% (94.6%)		
Distances ⓘ	Identical Matches	0.0% (0.1%)
	Average Distances	2.19 (1.28)

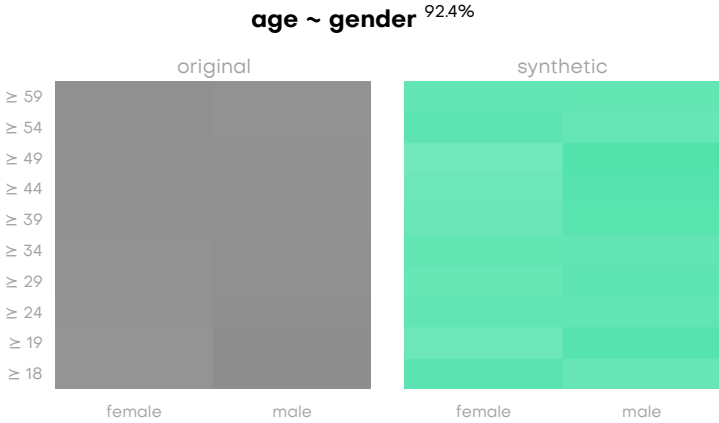
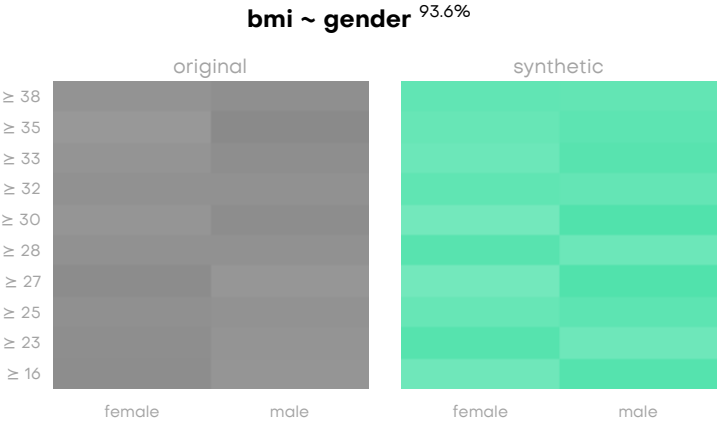
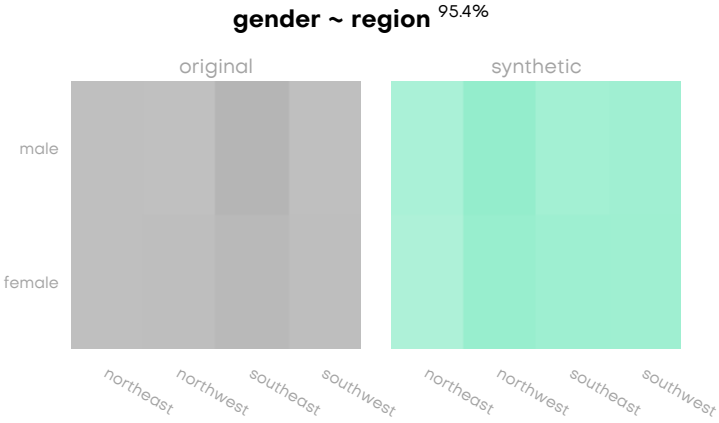
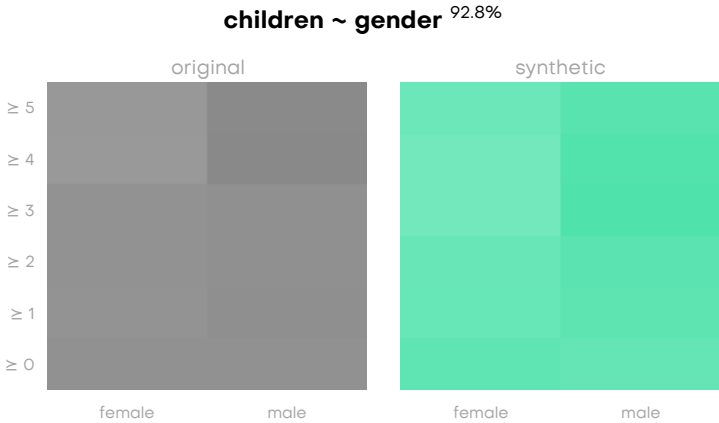
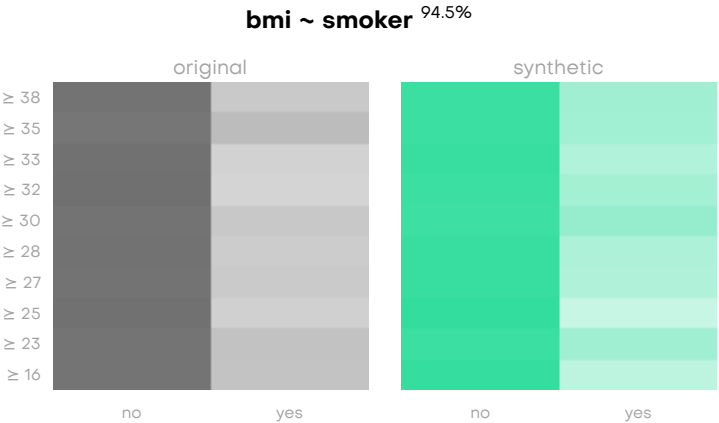
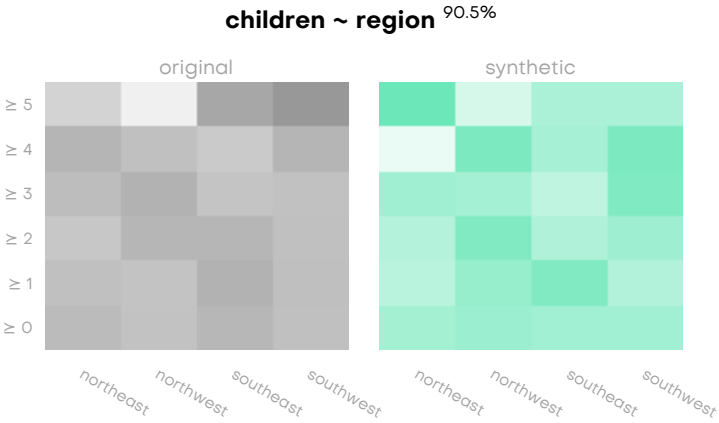
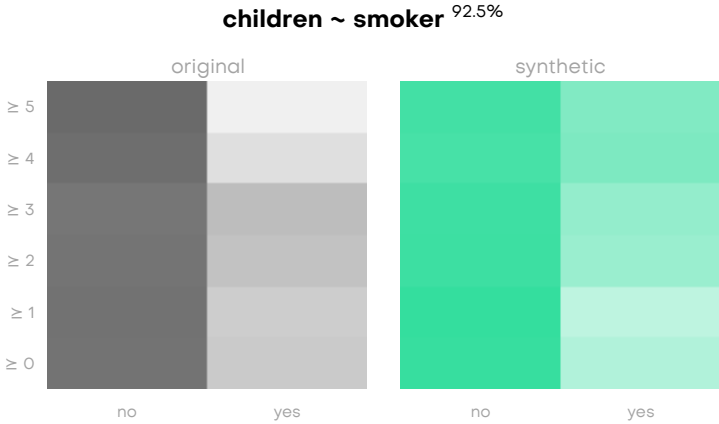
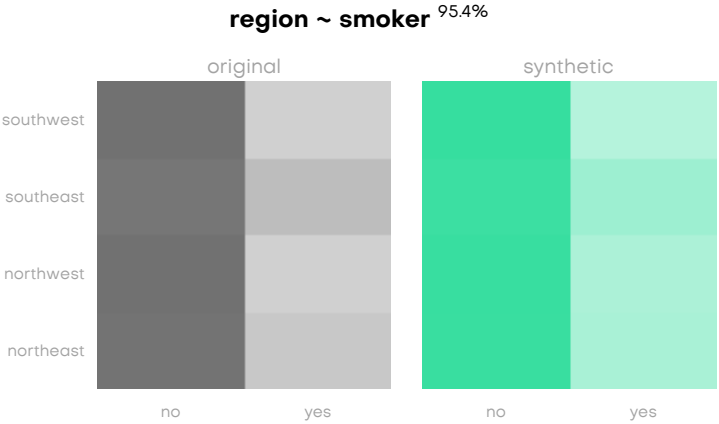
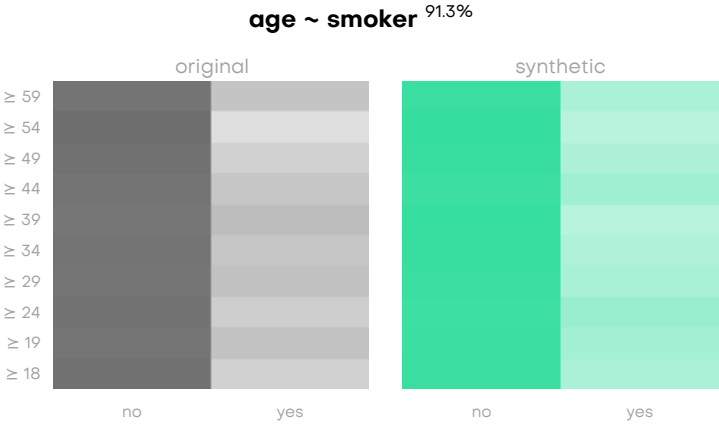
Correlations





Bivariate Distributions





Accuracy

Accuracy Matrix

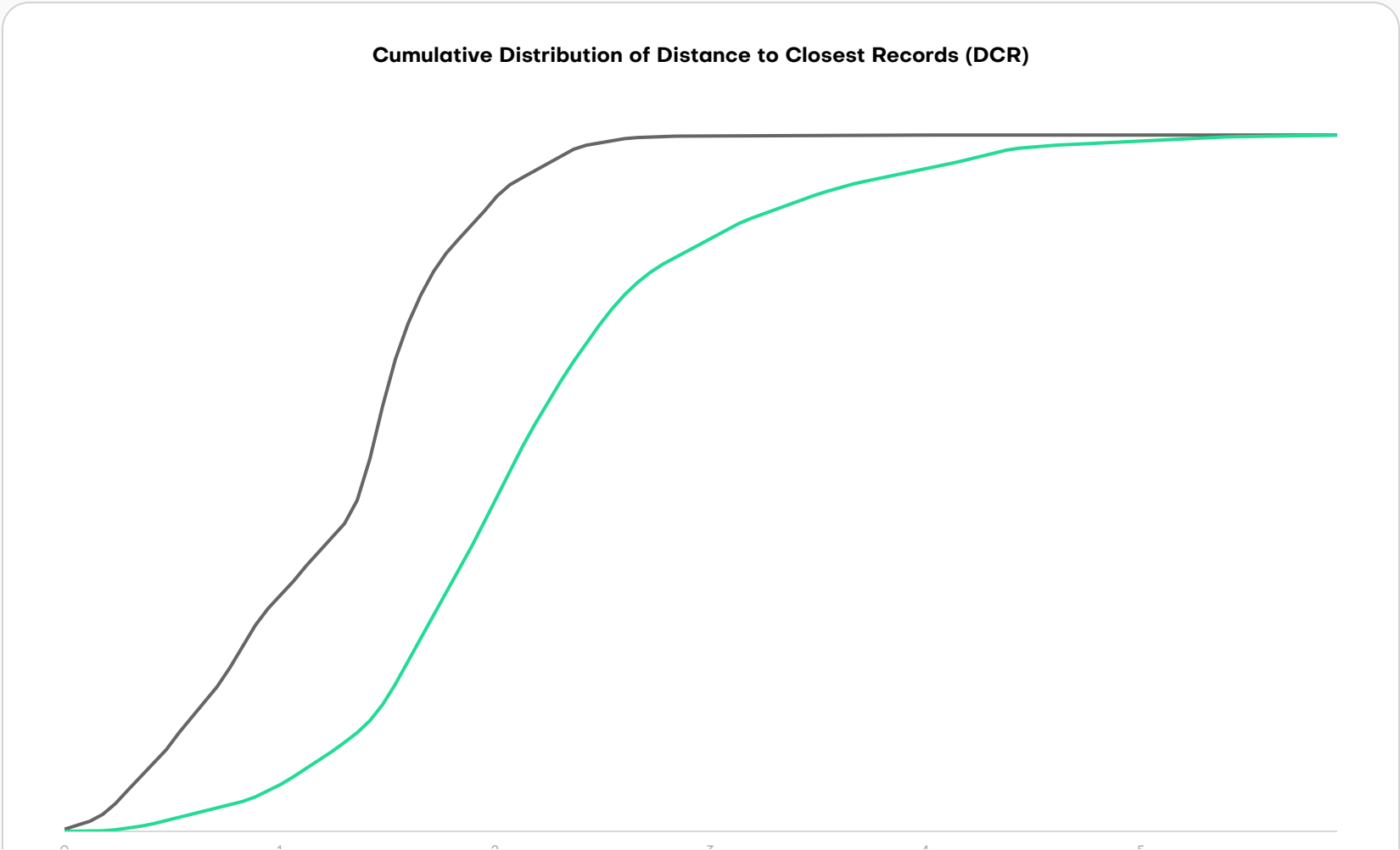
	age	bmi	children	expenses	gender	region	smoker
smoker	0.65	0.70	0.65	0.95	0.70	0.85	0.75
region	0.60	0.55	0.60	0.65	0.75	0.85	0.75
gender	0.70	0.75	0.65	0.65	0.85	0.85	0.75
expenses	0.95	0.55	0.60	0.85	0.65	0.60	0.95
children	0.55	0.70	0.70	0.60	0.75	0.65	0.70
bmi	0.55	0.85	0.65	0.55	0.75	0.55	0.75
age	0.75	0.55	0.50	0.95	0.70	0.65	0.65

Accuracy of synthetic data is assessed by comparing the distributions of the synthetic (shown in green) and the original data (shown in gray). For each distribution plot we sum up the deviations across all categories, to get the so-called total variation distance (TVD). The reported accuracy is then simply reported as 100% - TVD. These accuracies are calculated for all univariate and bivariate distributions. A final accuracy score is then calculated as the average across all of these.

Distances

Identical Matches: 0.0% (0.1%)

Average Distances: 2.19 (1.28)





Explainer

Synthetic data shall be close, but not "too close" to the original data in order to preserve the confidentiality of the original samples. This can be asserted empirically by measuring distances between synthetic records to their closest original records. These distances can then be compared against the observed distances within the original data itself, which serve as a reference for what it means to be "too close". For the visualization above, the distances for the synthetic data are displayed in green, and the distances for the original data displayed in gray. A green line that is significantly left of the gray line within the cumulative density plots would imply that the generated data is too close to the actual records.