# Index

## A

# J

Jaccard distance, 107

Jaccard index/coefficient, 107–108

joint cumulative distribution function, 135

joint distribution, 120, 135

joint probability, 120, 165, 166, 167

joint probability density functions, 136

joint probability mass functions, 135–136

Julia programming language, 23

# K

Kappa value, 77

kernel trick, 207–208

kernels, 207

*k*-fold cross-validation method, 68–70, 335, 374–375

*k*-means algorithm, 67, 247–255, 349

appropriate number of clusters, 249

elbow method, 249–254

strengths and weaknesses, 248

*k*-medoids algorithm, 255–257

## O