







## Dealing With Missing Data • Some learning methods can handle missing values • Throw away cases with missing values • in some data sets, most cases get thrown away • if not missing at random, throwing away cases can bias sample towards certain kinds of cases • Treat "missing" as a new attribute value • what value should we use to code for missing with continuous or ordinal attributes? • if missing causally related to what is being predicted? • Impute (fill-in) missing values • once filled in, data set is easy to use • if missing values poorly predicted, may hurt performance of subsequent uses of data set

## Potential Problems Imputed values may be inappropriate: in medical databases, if missing values not imputed separately for male and female patients, may end up with male patients with 1.3 prior pregnancies, and female patients with low sperm counts many of these situations will not be so humorous/obvious! If some attributes are difficult to predict, filled-in values may be random (or worse) Some of the best performing machine learning methods are impractical to use for filling in missing values (neural nets) Beware of coding - reliably detect missing cases can be difficult

## Imputing Missing Values • Fill-in with mean, median, or most common value • Predict missing values using machine learning • Expectation Maximization (EM): Build model of data values (ignore missing vals) Use model to estimate missing values Build new model of data values (including estimated values from previous step) Use new model to re-estimate missing values Re-estimate model Repeat until convergence





























