

Cross-Validation Approach

This strategy will be used primarily by the groups conducting voxel-based analyses. However, the partitioning assignments will be made available for others who wish to use this strategy. For example, to ensure that a region of interest drawn during the data summary process is selected in an unbiased fashion, we will use a cross-validation approach. Here we describe a pre-set partitioning scheme that will allow different groups working with the data to use the same test and validation subsets.

Training-test: 40% of all subjects will be chosen for a training dataset A and 60% for a test dataset B . The scheme for choosing the subjects will be stratified by a) diagnosis: NC, MCI, AD, b) study arm (1.5T only, PET+1.5T, 1.5T and 3T) and c) young (<76) vs old (>76). This will ensure that the correct proportion of age group and each diagnosis group is in each set and that the PET data will *also* be split 40-60 to give sets A_{PET} and B_{PET} . Note that the partitioning scheme will be such that all subjects in the set A_{PET} will also be contained in the set A and similarly for B_{PET} and B . We believe providing a standard split will foster homogeneity of methods and comparability across different ADNI manuscripts. Some investigators may wish, however, to split using different proportions. For increased precision of estimation, we will also provide access to a pre-specified leave-k-out cross-validation scheme.

Leave k-out cross-validation:

Set A will be divided into 4 equal parts and set B into 6 equal parts (again stratified by diagnosis, study arm, and young vs old) to give 10 sets each consisting of a stratified 10% of the ADNI sample.

Implementation: Randomization to training-test and leave-k-out datasets will be designed by the biostatistics core and implemented in coordination with the UCSD clinical core and UCLA data management core. The randomization will be carried out for individuals before data for analysis are posted to the LONI database, using the following steps:

1. Determination of stratum: Based on clinical data, study arm, and age, each individual will be classified into one of eighteen strata according to baseline clinical diagnosis (NC, MCI, AD), study arm (1.5T only, 1.5T and PET, 1.5T and 3T) and age (<76 vs >76).
2. Assignment of a sequence number within stratum: Individuals within a stratum will be assigned consecutive integers according to the order in which they were processed into the stratum.
3. Assignment to cross-validation partition subset: Each individual will be assigned to a cross-validation partition subset using a table-lookup scheme. The table, prepared by the Biostatistics Core, will consist of random block permutations to ensure balance over time. Each individual will thus be identified in one of 10 approximately equal-sized blocks for leave-k-out cross validation. Four of these blocks, chosen at random, will comprise the training set A and the remaining 6 the validation set B , and individuals will also be labeled as belonging to the training set or validation set.
4. Posting of labels: Each individual's training set/ validation set and k-cross-validation set label will be posted on adni-info.org.
5. Instructions for appropriate statistical use of the training set/ validation set, and of the k-cross-validation sets, will be posted to the LONI biostatistics website.

Whenever comparisons between MRI and PET are made, only the appropriate PET subsets should be used. The above design ensures that *each* (PET) set is a subset of the corresponding MRI set, which will simplify things when MRI groups modify analyses for PET comparison.