# Why 'stratify' our fold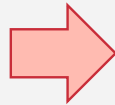s?