

I like the idea and I've been playing with that too. I never tried clustering features using kmeans of the corr matrix though. Any motivation for kmeans of the corr matrix?

I did try clustering with some homegrown methods though. When I add mean, std, etc per group as additional features, the resulting models always show improvements OOS. What I did to cluster features was something like:

1. two features that correlate above some threshold are defined as neighbors
2. features groups are neighborhoods

You could also try linear regression coefficients instead of corr values to define neighbors, seems to give interesting results as well.