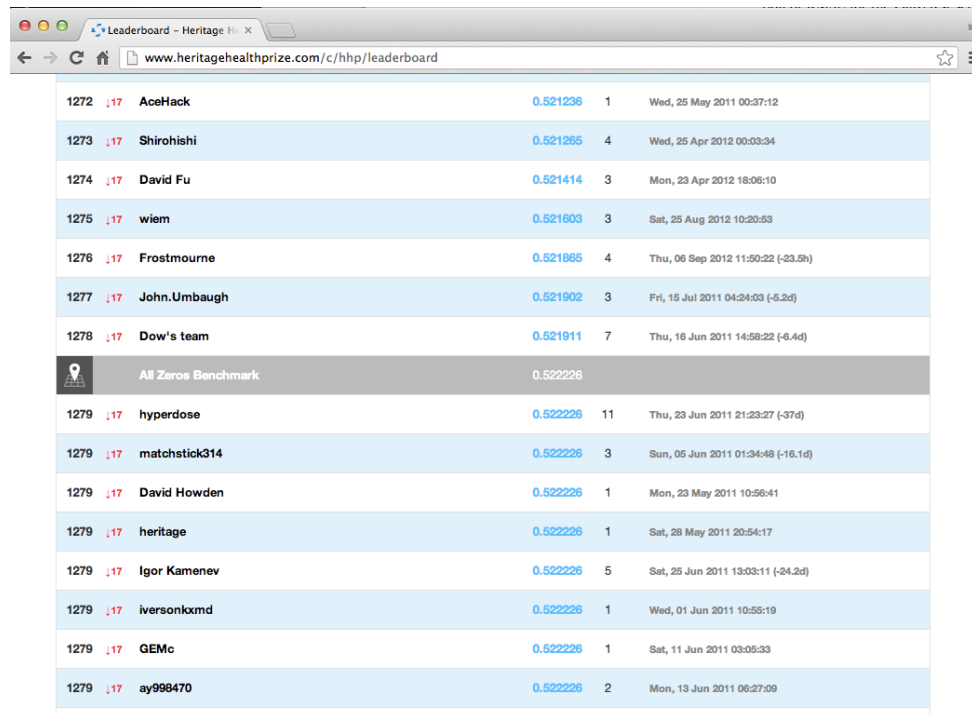# Prediction study design

Jeffrey Leek
Johns Hopkins Bloomberg School of Public Health

# Prediction study design

1. Define your error rate

2. Split data into:

   - Training, Testing, Validation (optional)

3. On the training set pick features

   - Use cross-validation

4. On the training set pick prediction function

   - Use cross-validation

5. If no validation

   - Apply 1x to test set

6. If validation

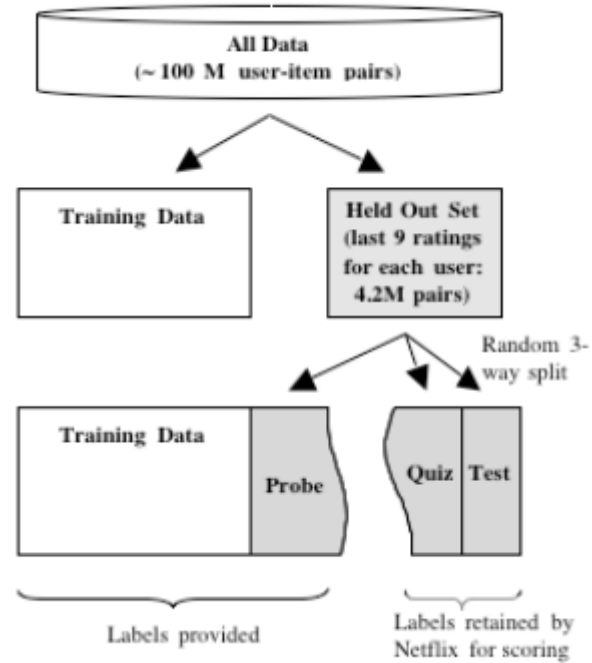   - Apply to test set and refine

   - Apply 1x to validation

# Know the benchmarks



http://www.heritagehealthprize.com/c/hhp/leaderboard

# Study design



http://www2.research.att.com/~volinsky/papers/ASAStatComp.pdf

# Used by the professionals



http://www.kaggle.com/

# Avoid small sample sizes

- Suppose you are predicting a binary outcome

    - Diseased/healthy

    - Click on ad/not click on ad

- One classifier is flipping a coin

- Probability of perfect classification is approximately:

    - $\left(\frac{1}{2}\right)^{sample\ size}$

    - $n = 1$ flipping coin 50% chance of 100% accuracy

    - $n = 2$ flipping coin 25% chance of 100% accuracy

    - $n = 10$ flipping coin 0.10% chance of 100% accuracy

# Rules of thumb for prediction study design

- If you have a large sample size

    - 60% training

    - 20% test

    - 20% validation

- If you have a medium sample size

    - 60% training

    - 40% testing

- If you have a small sample size

    - Do cross validation

    - Report caveat of small sample size

# Some principles to remember

· Set the test/validation set aside and *don't look at it*

· In general *randomly* sample training and test

· Your data sets must reflect structure of the problem

    - If predictions evolve with time split train/test in time chunks (called backtesting in finance)

· All subsets should reflect as much diversity as possible

    - Random assignment does this

    - You can also try to balance by features - but this is tricky