

Stability and Aggregation of Experimental Results

Cheng Soon Ong

Machine Learning Research Group
Data61 | CSIRO
Australian National University

10 August 2017
ICML workshop on Computational Biology

What is machine learning?

Machine learning is about prediction

Examples/features	$x_1, \dots, x_n \sim \mathcal{X}$
Labels/annotations	$y_1, \dots, y_n \sim \mathcal{Y}$
Predictor	$f_{\mathbf{w}}(x) : \mathcal{X} \rightarrow \mathcal{Y}$

Estimate best predictor = training

Given data $(x_1, y_1), \dots, (x_n, y_n)$, find a predictor $f_{\mathbf{w}}(\cdot)$.

- No mechanistic model of the phenomenon
- There is relatively large amounts of data (examples, x usually \mathbb{R}^d)
- The outcomes (labels, y usually binary) are well defined

Prediction \neq understanding

How can we use prediction to help with scientific research?

What are good features?

$$f_{\mathbf{w}}(x) : \mathcal{X} \rightarrow \mathcal{Y}$$

What to measure?

$$f_{\mathbf{w}}(x) : \mathcal{X} \rightarrow \mathcal{Y}$$

Use predictive model ...

$$f_{\mathbf{w}}(x) : \mathcal{X} \rightarrow \mathcal{Y}$$

... to discover good biomarkers

$$f_{\mathbf{w}}(x) : \mathcal{X} \rightarrow \mathcal{Y}$$

- Genome Wide Association Studies
- Spearman's Correlation
- Stability and Aggregation

Case-control studies

A cohort of sick individuals (**cases**) and healthy individuals (**controls**) are genotyped and their corresponding binary phenotype are recorded.

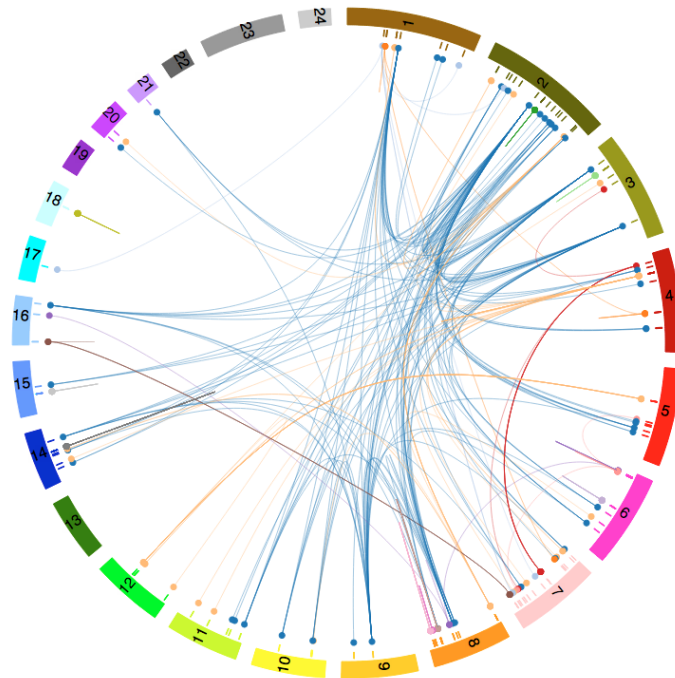
We use the framework of hypothesis testing

Hypothesis testing Given a case control study, test whether a particular SNP is associated with the phenotype.

Epistatic Interactions

- WTCCC data
- Need to tabulate 125 billion contingency tables
- Consider specificity and sensitivity
- Gain over univariate ROC
- CPU (\approx days) and GPU (\approx hours)
- Store the top 1 million pairs

Interacting with results



D3.js

- Circular plot
- Linear plot
- Manhattan plot
- Heat map

Interaction

- Filter
- Zoom
- Drill down
- Call out

Cristovao Freitas Iglesias Junior, Stefan Sevelde
github.com/chengsoon.org/rede

Interpreting p-values

Is 10^{-10} probability of association very significant?

Quote

... but a reliable method of procedure. In relation to the test of significance, we may say that a phenomenon is experimentally demonstrable when we know how to conduct an experiment which will rarely fail to give us a statistically significant result.

Fisher, *The Design of Experiments*, 1947, p. 14

Stability of scoring

We consider p-values as a score of association.

- How stable is this score if we repeat the experiment?
- How do we combine scores?

Challenges

- Scores available for only the top-k examples
- Scores from different sources not calibrated

Use predictive model ...

$$f_{\mathbf{w}}(x) : \mathcal{X} \rightarrow \mathcal{Y}$$

... to discover good biomarkers

$$f_{\mathbf{w}}(x) : \mathcal{X} \rightarrow \mathcal{Y}$$

Stability of feature selection

How to measure overlap?

Rank aggregation

How to combine different sources of information?

Modeling using Spearman's correlation

How to represent ranks?



Multiple ways to represent ranks

- Ordered list of n objects selected from Ω
- List of values $[1, \dots, n]$ (the ranks of the object)
- Normalised ranks $\in (0, 1)$
- Permutation mapping $R : \Omega \rightarrow (0, 1)$

Motivation

Given a set of replicated experiments, how do we measure overlap?

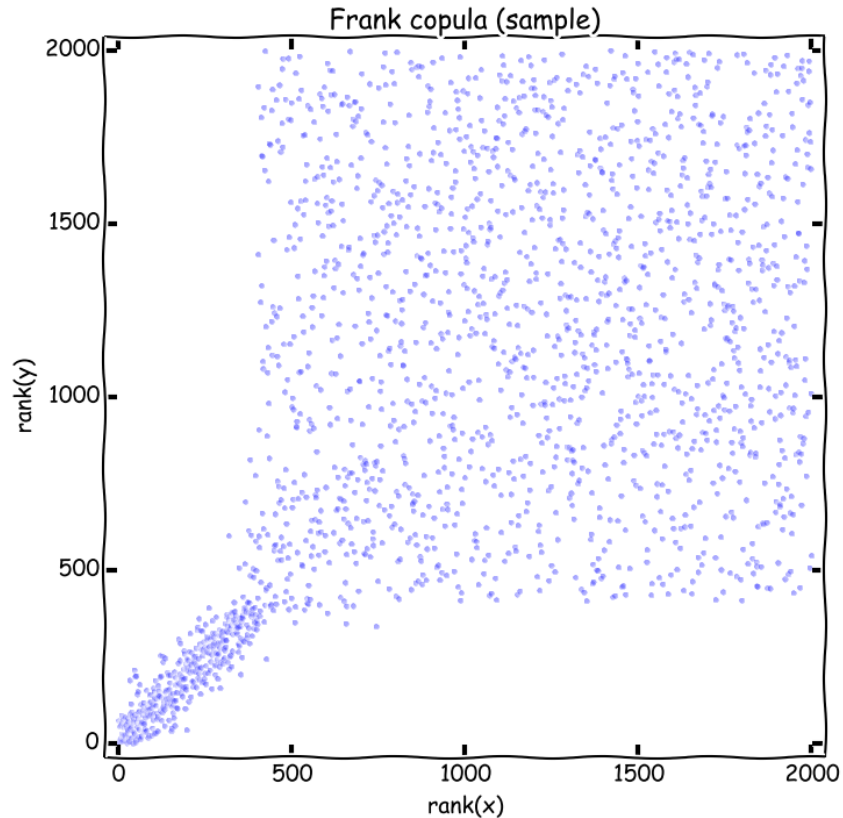
Examples

- Perform repeated splits of the data
- Experiments on different cohorts
- Multiple sources of information

Challenges

- Scores available for only the top-k examples
- Scores from different sources not calibrated

Signal and Noise



Running example (6 objects)

$$A = [a, b, c, d, e, f]$$

$$B = [a, b, e, f, c, d]$$

Jaccard Index

$$\text{overlap} = \frac{|A \cap B|}{|A \cup B|}$$

Measuring stability

- Easy to compute
- Works for top-k lists

Consider the top-3 lists from above:

$$\text{Jaccard index} = \frac{|\{a, b\}|}{|\{a, b, c, e\}|} = \frac{1}{2}$$

- Ignores the order given by scores

- Similar to Pearson's correlation for the measure of dependence
- Spearman's ρ is a correlation measure between ranked lists

$$\rho(A, B) := \frac{\sum_i (r_A^{(i)} - \bar{r}_A)(r_B^{(i)} - \bar{r}_B)}{\sqrt{\sum_i (r_A^{(i)} - \bar{r}_A)^2 \sum_i (r_B^{(i)} - \bar{r}_B)^2}},$$

- Running example:

$$\rho([a, b, c, d, e, f], [a, b, e, f, c, d]) = 0.543$$

(Jaccard index = 1)

- **Need the same elements in A and B**

$$\rho([a, b, c], [a, b, e]) ?$$

Spearman's ρ on top k lists

Our idea

Define Spearman's ρ for top k lists

Key observation

Any elements in list A that do not appear in list B must have a rank higher than the number of elements in B

Running example (top-3)

$$A = [a, b, c, d, e, f] \quad \text{and} \quad B = [a, b, e, f, c, d]$$

$$A_3 = [a, b, c] \quad \text{and} \quad B_3 = [a, b, e]$$

$$A_3 \xrightarrow{B_3} = [a, b, c, e] \quad \text{and} \quad B_3 \xrightarrow{A_3} = [a, b, e, c]$$

$$\text{Spearman's } \rho = \rho(A_3 \xrightarrow{B_3}, B_3 \xrightarrow{A_3}) = 0.8$$

Extend the list

We expand lists A and B to complete rankings over the same set of elements, denoting them as $A \xrightarrow{B}$ and $B \xrightarrow{A}$ respectively.

The missing values in the extension are given the average rank.

Running example (top-4)

$$A_4 = [a, b, c, d] \quad \text{and} \quad B_4 = [a, b, e, f]$$

$$A_4 \xrightarrow{B_4} = [1, 2, 3, 4, 5.5, 5.5] \quad \text{and} \quad B_4 \xrightarrow{A_4} = [1, 2, 5.5, 5.5, 3, 4]$$

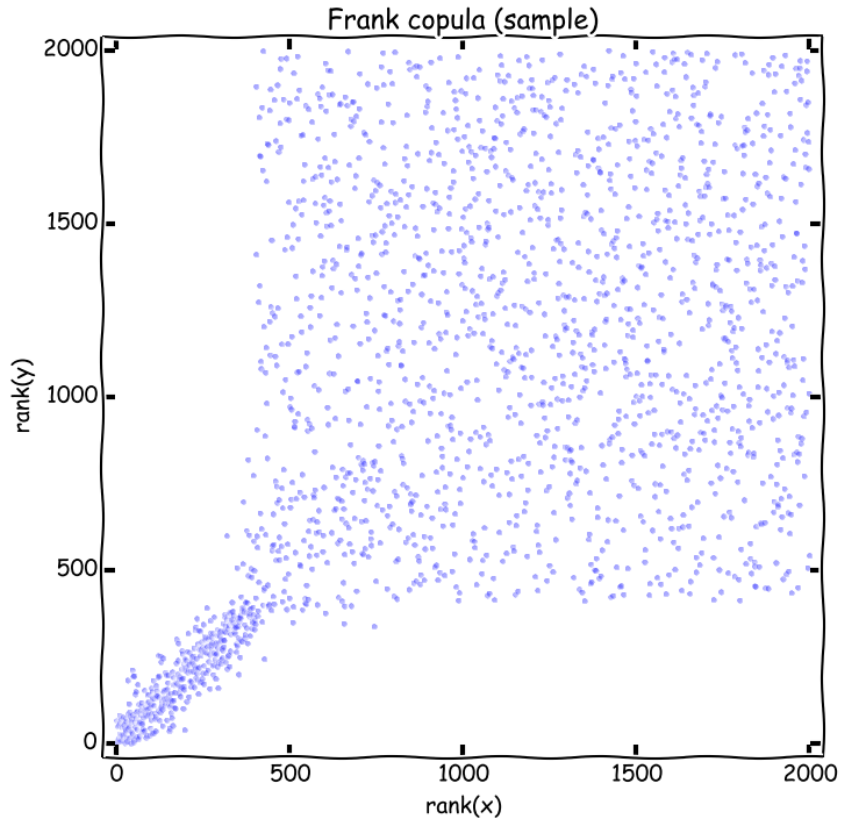
Makes no assumption about the order of the unranked objects

Other possible imputation approaches

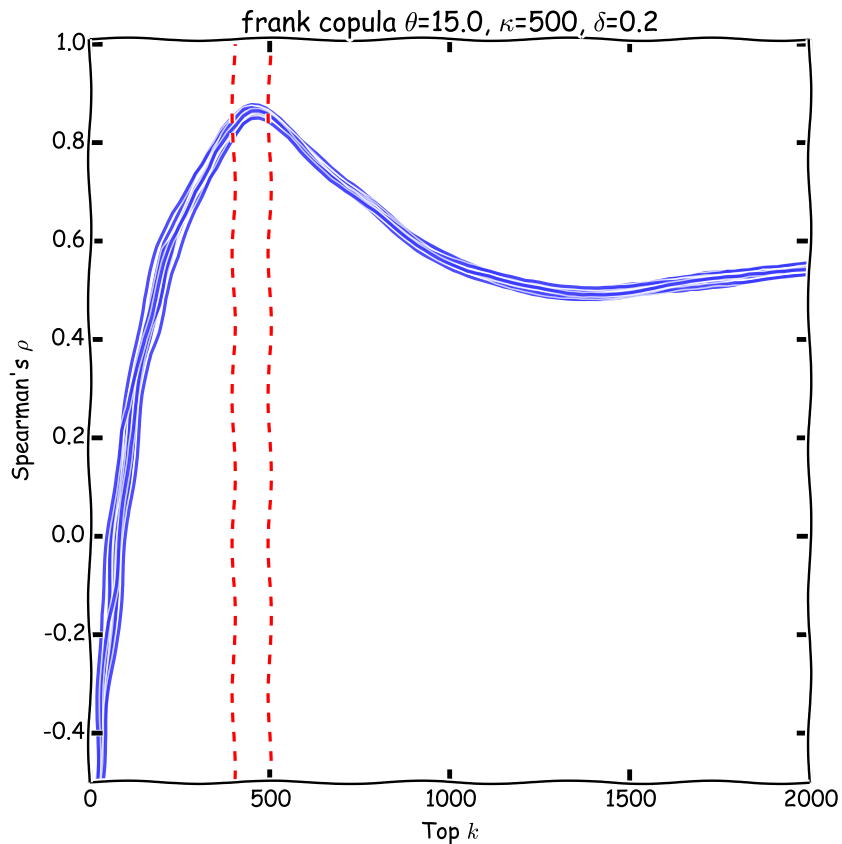
- Optimistic
- Worst case

Bedó, Rawlinson, Goudey, Ong, PLoS ONE, 2014

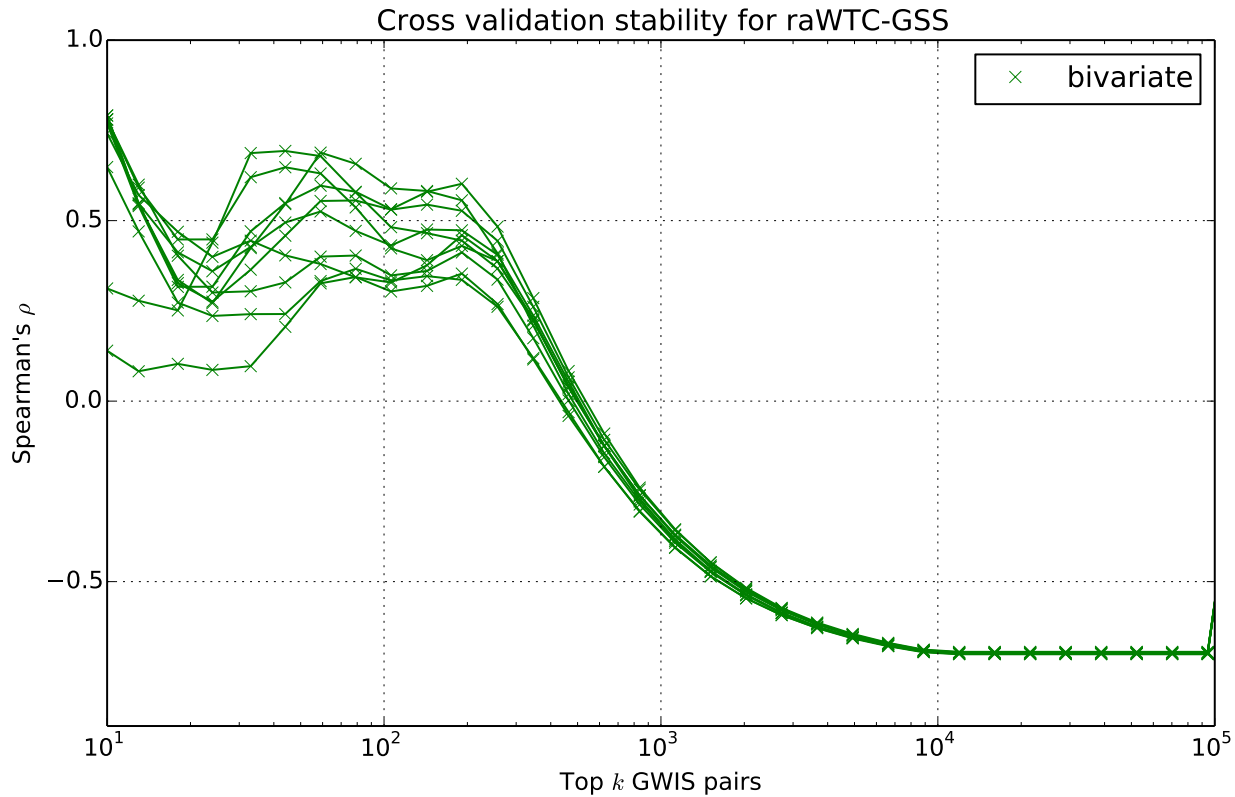
Signal and Noise



Spearman's ρ



Simulate two cohorts by splitting



Motivation

Given a set of replicated experiments, how do we measure overlap?

Challenges

- Scores available for only the top-k examples
- Scores from different sources not calibrated

Model

- **Ranked list** Instead of just using set intersection, we can use the scores from GWIS to order the results
- **top k** Traditional methods (Spearman's ρ) requires ranks for the whole list. We have incomplete information, but we know our ranks are the top ones.
- **Multivariate** Textbook Spearman's ρ is for computing correlation between two ranks. We want to compute the correlation between multiple ranked lists.

Intuition

For continuous random variables, copulas model the dependence component after discounting for univariate marginal effects

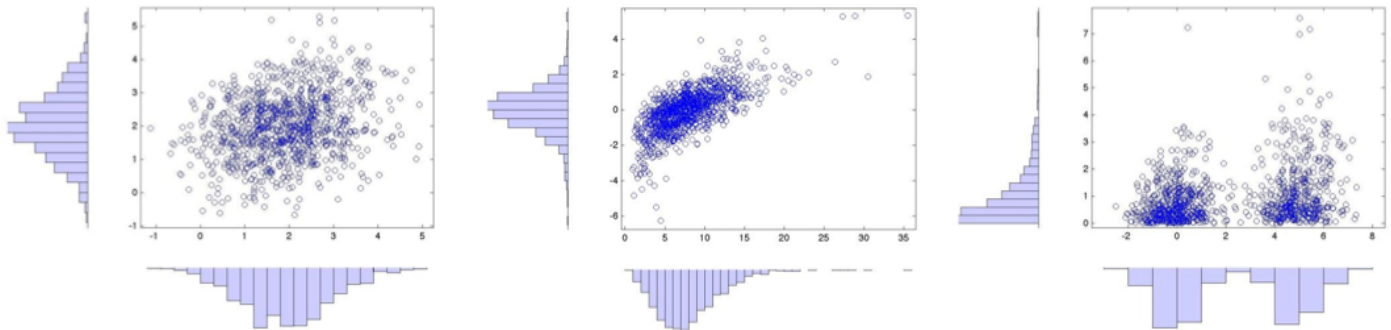
Probabilistic definition

Let U_1, \dots, U_d be real random variables $\sim U([0, 1])$.

A copula function $C : [0, 1]^d \rightarrow [0, 1]$ is a joint distribution

$$C_\theta(u_1, \dots, u_d) = P(U_1 \leq u_1, \dots, U_d \leq u_d)$$

The same Gaussian copula function



Spearman's ρ can be expressed in terms of the copula

$$\rho(A, B) = 12 \int_{[0,1]^2} C(u, v) du dv - 3$$

Proof

$$\begin{aligned} \rho(A, B) &= \frac{\sum_i (r_A^{(i)} - \bar{r}_A)(r_B^{(i)} - \bar{r}_B)}{\sqrt{\sum_i (r_A^{(i)} - \bar{r}_A)^2 \sum_i (r_B^{(i)} - \bar{r}_B)^2}} \\ &= \frac{\mathbb{E}[F(X)G(Y)] - \mathbb{E}[F(X)]\mathbb{E}[G(Y)]}{\text{STD}[F(X)]\text{STD}[G(Y)]} \\ &= \frac{\mathbb{E}[F(X)G(Y)] - \frac{1}{2}}{\frac{1}{12}} \\ &= 12\mathbb{E}[F(X)G(Y)] - 3 \\ &= 12 \int \int uvC(u, v) - 3 \\ &= 12 \int_{[0,1]^2} C(u, v) du dv - 3 \end{aligned}$$

Spearman's ρ can be expressed in terms of the copula

$$\rho(A, B) = 12 \int_{[0,1]^2} C(u, v) du dv - 3$$

Empirical copula

$$C_n(u, v) = \frac{1}{|\Omega|} \sum_{x \in \Omega} \mathbf{1}(R(x) \leq u, S(x) \leq v)$$

Why do the math?

- Unclear how to extend formula for Spearman's correlation.
- Multivariate distributions \Rightarrow multivariate copula.

A multivariate extension of Spearman's ρ

For a d dimensional set of random variables \mathbf{u} , the multivariate Spearman's ρ is given by

$$\rho(R_1, \dots, R_d) = Q(C, \pi) = h(d) \left(2^d \int_{[0,1]^d} \pi(\mathbf{u}) \, dC(\mathbf{u}) - 1 \right),$$

where

$$h(d) = \frac{d + 1}{2^d - (d + 1)}.$$

Empirical multivariate Spearman's correlation

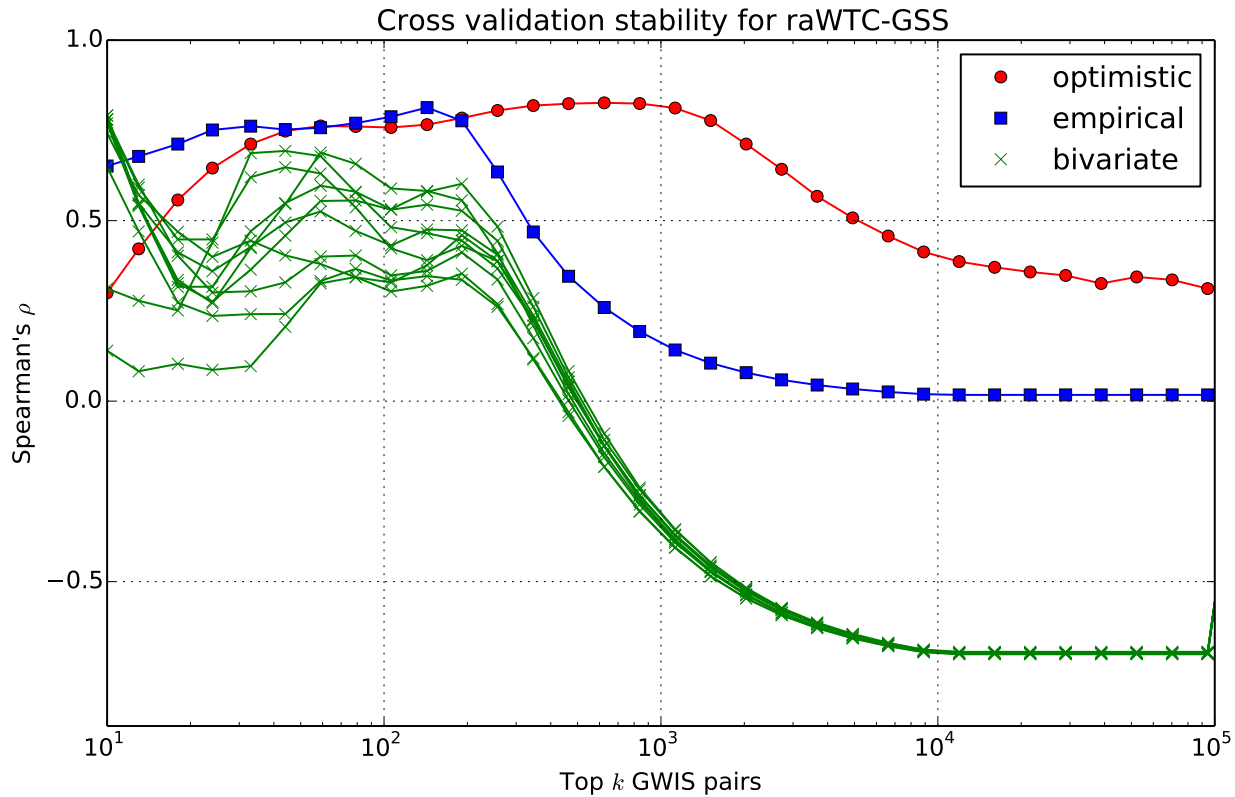
$$\rho_n(R_1, \dots, R_d) = h(d) \left[\frac{2^d}{n} \sum_x \prod_{j=1}^d R_j(x) - 1 \right].$$


```
n_rank /= n+1.  
prod = np.prod(1.-n_rank, axis=1)  
total = np.sum(prod)  
s_rho = h(d) * ((2**d/float(n)) * total - 1.)
```

Empirical multivariate Spearman's correlation

$$\rho_n(R_1, \dots, R_d) = h(d) \left[\frac{2^d}{n} \sum_x \prod_{j=1}^d R_j(x) - 1 \right].$$

Multiple replicates



Wait... there's more

Modeling using Spearman's correlation

Use predictive model ...

$$f_{\mathbf{w}}(x) : \mathcal{X} \rightarrow \mathcal{Y}$$

... to discover good biomarkers

$$f_{\mathbf{w}}(x) : \mathcal{X} \rightarrow \mathcal{Y}$$

Stability of feature selection

How to measure overlap?

$$\rho(R_1, \dots, R_d)$$

Rank aggregation

How to combine different sources of information?

Macintyre, Yepes, Ong, Verspoor, PeerJ, 2014

How to combine different sources of information?

We maximise multivariate correlation

$$R^* = \arg \max_R \rho(R, R_1, R_2, \dots, R_d).$$

Theorem The aggregator that maximises multivariate Spearman's correlation is the product of the normalised ranks.

Use the geometric mean

NOT pairwise correlation

Instead of decomposing the association into a combination of pairwise similarities $\rho(R, R_1), \rho(R, R_2), \dots, \rho(R, R_d)$.

Method

1. Divide rank by number of items
2. return log average

Problem setting

We are given a ranking L of n objects which comprise our labels, and a set of d experts $\{R_j\}$.

Find a weighting of the experts.

Relaxed optimal aggregator

We solve the **least squares problem**

$$\min_{\omega} \sum_x \left(l(x) - \sum_{j=1}^d \omega_j r'_j(x) \right)^2,$$

where the outer sum is over the n examples x ,

$l(x)$ is the log scaled normalised labels,

$r'_j(x)$ is the log scaled normalised completed ranks.

LETOR 4.0

We (surprisingly) perform much better than state of the art

Bedó, Ong, JMLR 17(201):1–30, 2016

Why is ML for Comp Bio hard?

Deep data

- High dimensional low sample size data
- Finite domain (we want predictions genome wide)
- Difficult to satisfy i.i.d. assumption

Prior knowledge

- Hard and soft constraints
- Weak constraints
- Causal systems with feedback and delays

Different problems/communities

- Genomics, proteomics, medical imaging, health records
- Communication challenges

Keep up the great work!

Prediction \neq understanding

How can we use prediction to help with scientific research?

- Use predictive model to discover good features

Spearman's correlation - applied to GWAS

- Stability of scoring
- Rank aggregation
- Supervised learning to rank
- Imputation from top-k lists
- Multivariate correlation using copulas

The dream

Facilitating scientific knowledge discovery
through automated experimental design with machine learning

Prediction \neq understanding

How can we use prediction to help with scientific research?

- Use predictive model to discover good features

Spearman's correlation - applied to GWAS

- Stability of scoring
- Rank aggregation
- Supervised learning to rank
- Imputation from top-k lists
- Multivariate correlation using copulas

The dream

Facilitating scientific knowledge discovery
through automated experimental design with machine learning

Please make your research open

Active Learning

- Choose a particular example to label using heuristics
- Annotator assumed to provide ground truth

Bandits

- Select a choice from a set of actions
- Simple algorithms with theoretical guarantees
- Manage uncertainty with repeated sampling

Choice theory

- Aggregate set of ranks into one ordering
- Economics and social science, impossibility theorems

Designing Experiments

- Choose a set of trials to measure
- Optimisation algorithms with theoretical analysis
- Information theory, real random variables