

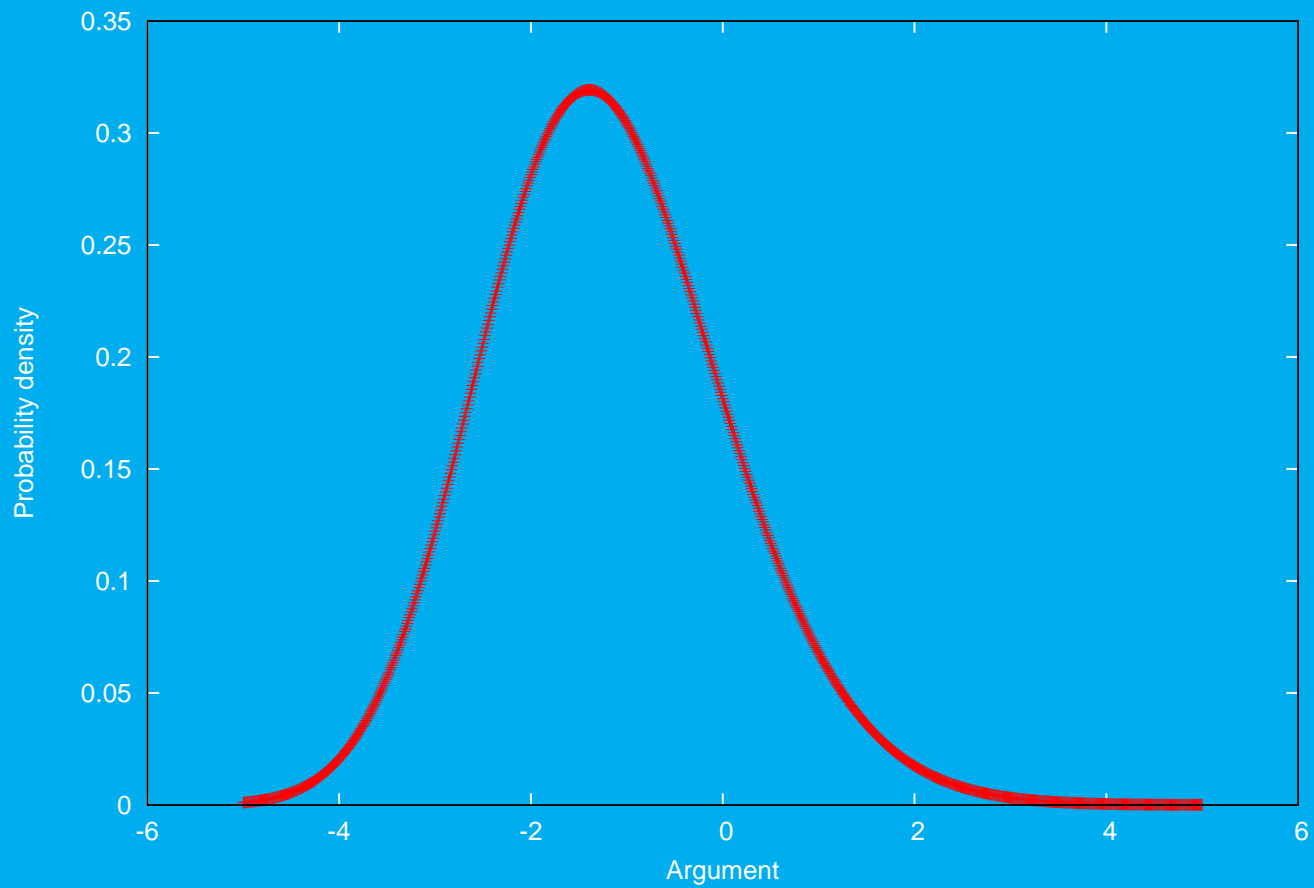
Given genotype data:
is it from homogeneous population?

Need statistic and formal test.

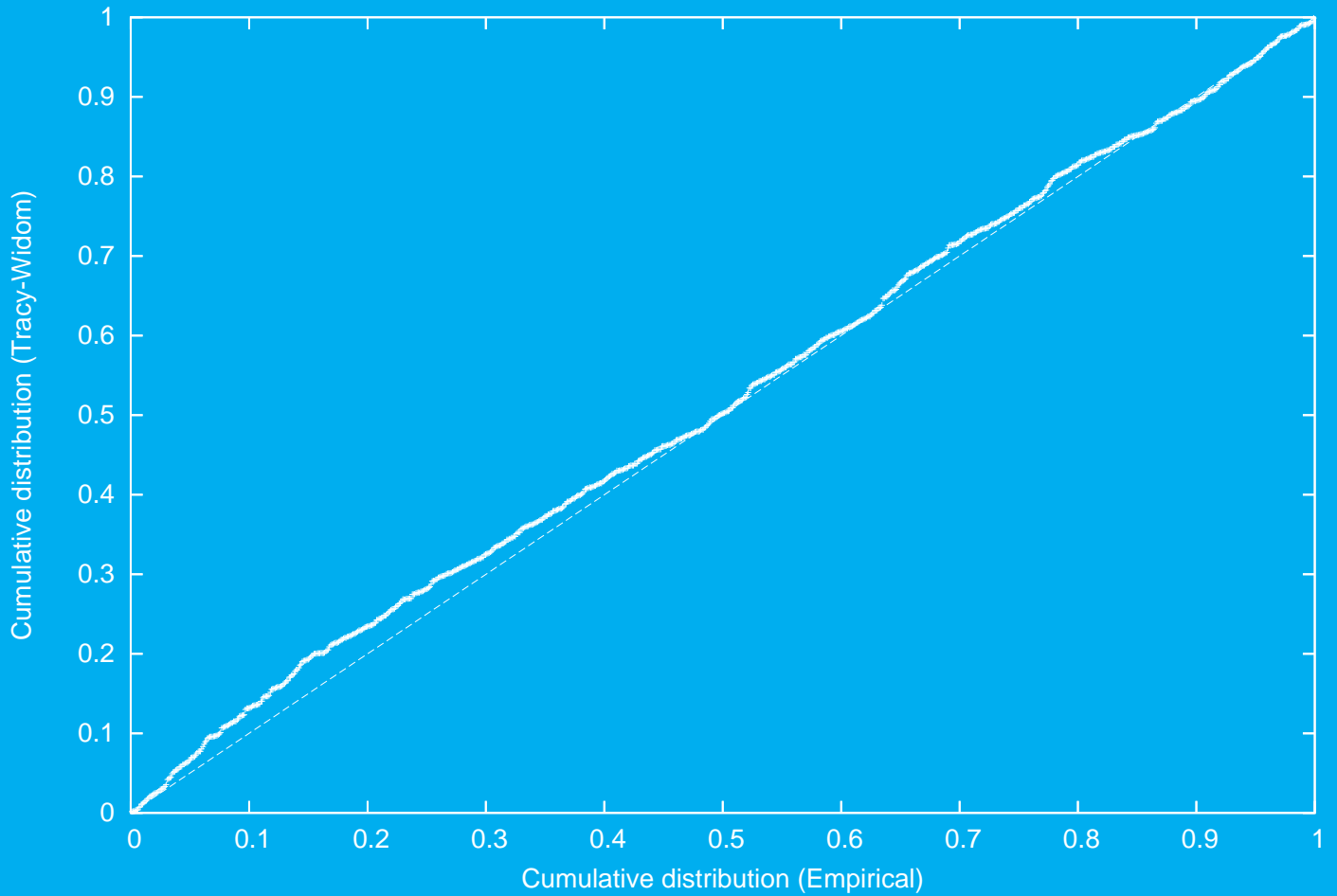
Want to test for additional structure when some found.

Should work when markers are in LD.

Tracy-Widom distribution



P=200, N=50000



Effective number of markers

Our matrices are *similar* to Wishart matrices.

(especially if markers have LD) *not Wishart*

Idea: Think of our covariance matrix as from a Wishart

BUT:

- Variance σ^2 .of Gaussian in each cell unknown
- Number of markers N unknown.

Estimate N , σ .

A moments estimator:

$$\hat{N} = \frac{(p+1)(\sum_i \lambda_i)^2}{((p-1)\sum_i \lambda_i^2) - (\sum_i \lambda_i)^2} \quad (1)$$
$$\hat{\sigma}^2 = \frac{\sum_i \lambda_i}{(p-1)\hat{N}}$$

This works better than max likelihood.

Testing additional values

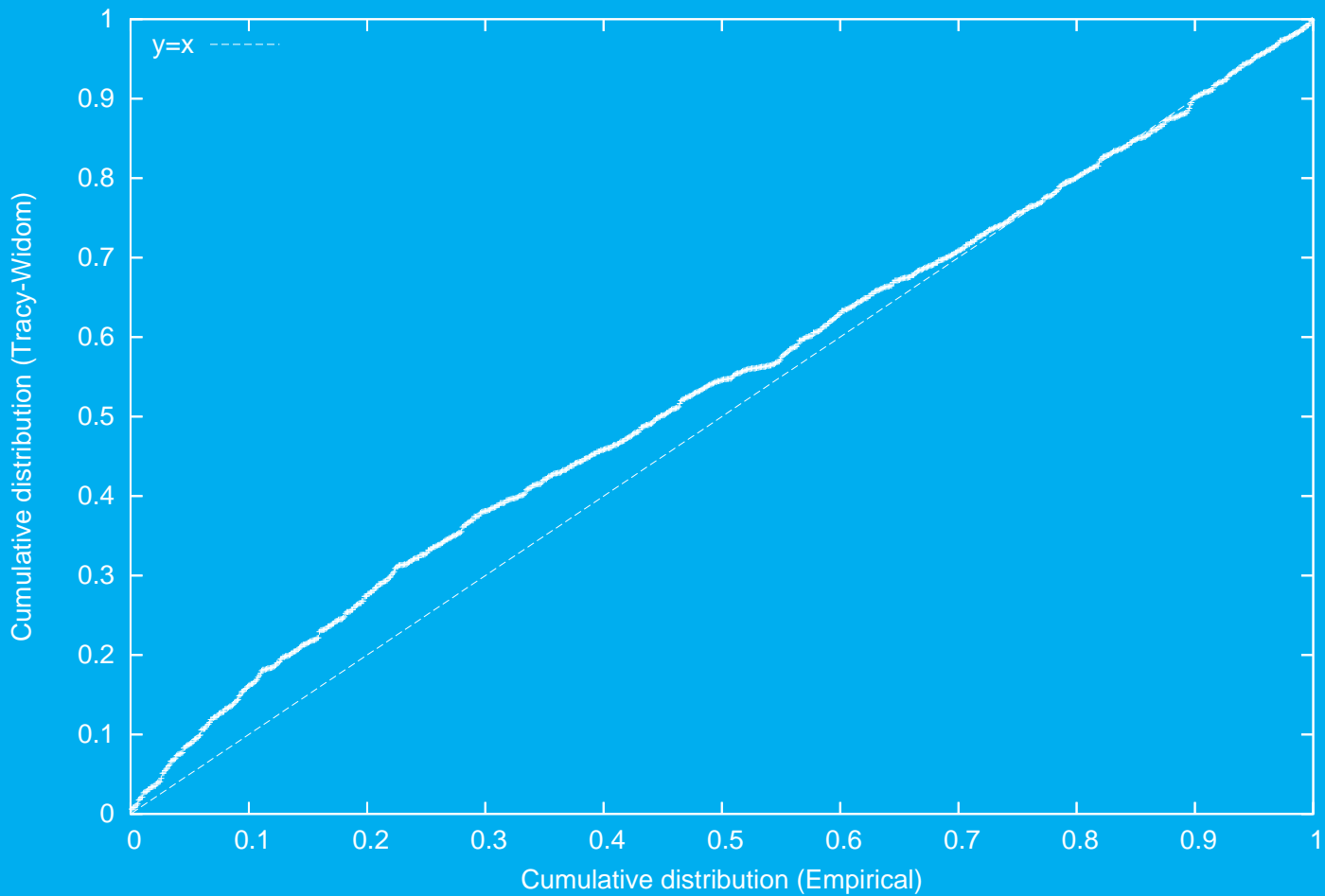
Just ignore eigenvalues already found to be significant

Run test on remaining eigenvalues

This is conservative (Johnstone, 2001)

Works *extremely well*

P-P plot. (Second eigenvalue) $p=100, n=5000$



The BBP Phase Change

How much data do we need to detect population structure?

Let l_1 be the lead eigenvalue of *theoretical* covariance (rest of eigenvalues 1)

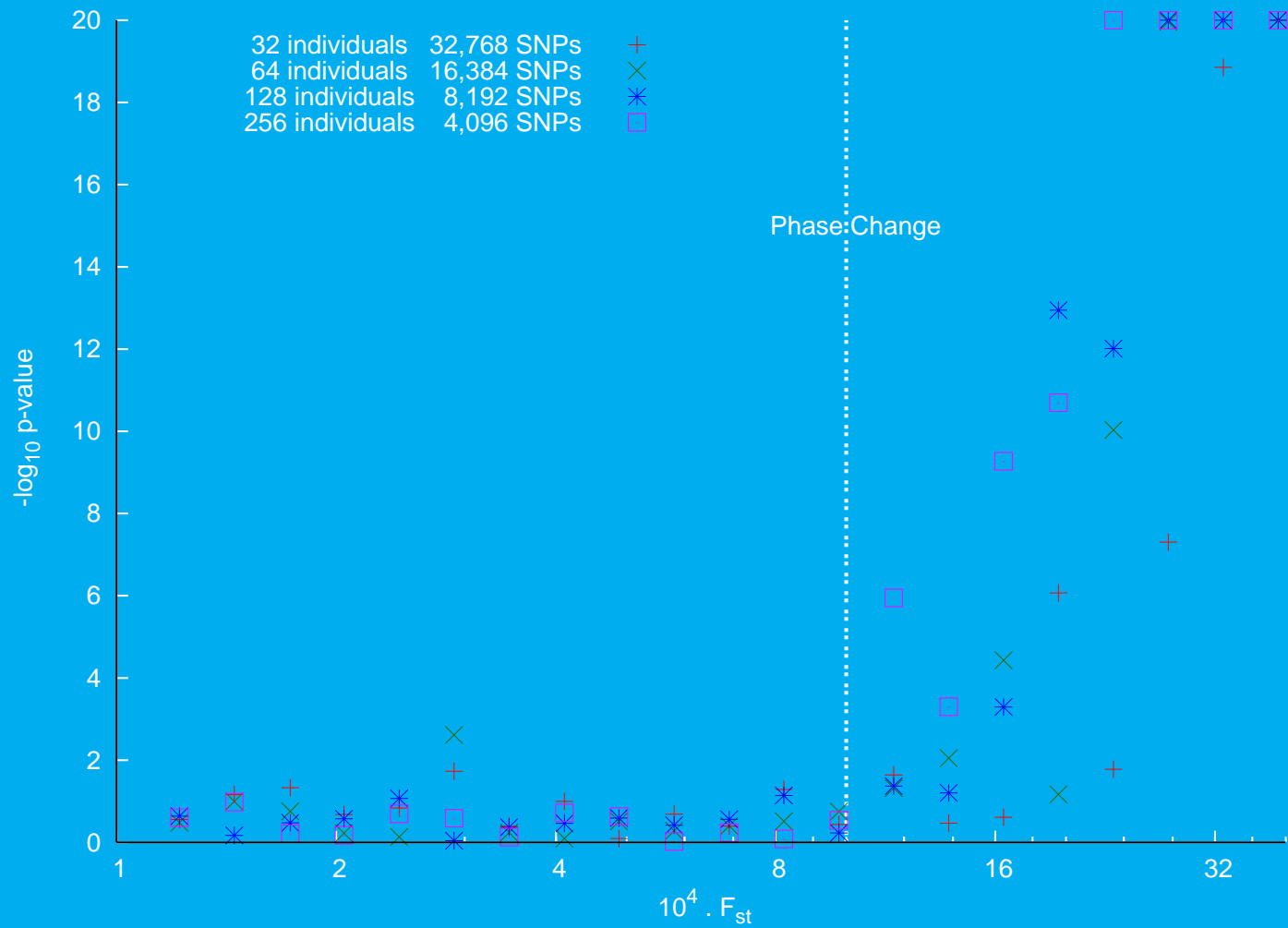
Baik et. al. (2005) prove that for a *complex* Wishart: Set $n/p = \gamma^2$. L_1 largest eigenvalue of sample covariance.

If $l_1 < 1 + 1/\gamma$ then as $p, n \rightarrow \infty$ L_1 tends in distribution to the same distribution as when $l_1 = 1$

If $l_1 > 1 + 1/\gamma$, then the TW-statistic becomes unbounded.

Change in character for small, large l_1 .

Conjecture (Baik, Ben Arous, P  ch   (2005)) : True for real Wishart too. Partially proved by Baik and Silverman (2006).



Seems to work very well. Note the sensitivity with large data sets:

If $n = 100,000$ (independent), $p = 1000$

$\tau = .001$ very easy to detect.

Detectable pop. structure will be present on most large datasets.

Other (Bayesian) models

Ancient frequency P .

Pop freq. in K populations:

$$\mathbf{p} = (p_1, p_2, \dots, p_K)$$

Conditional on P :

p_i has mean P

Covariance $P(1 - P)B$

Nicholson et al., STRUCTURE ...

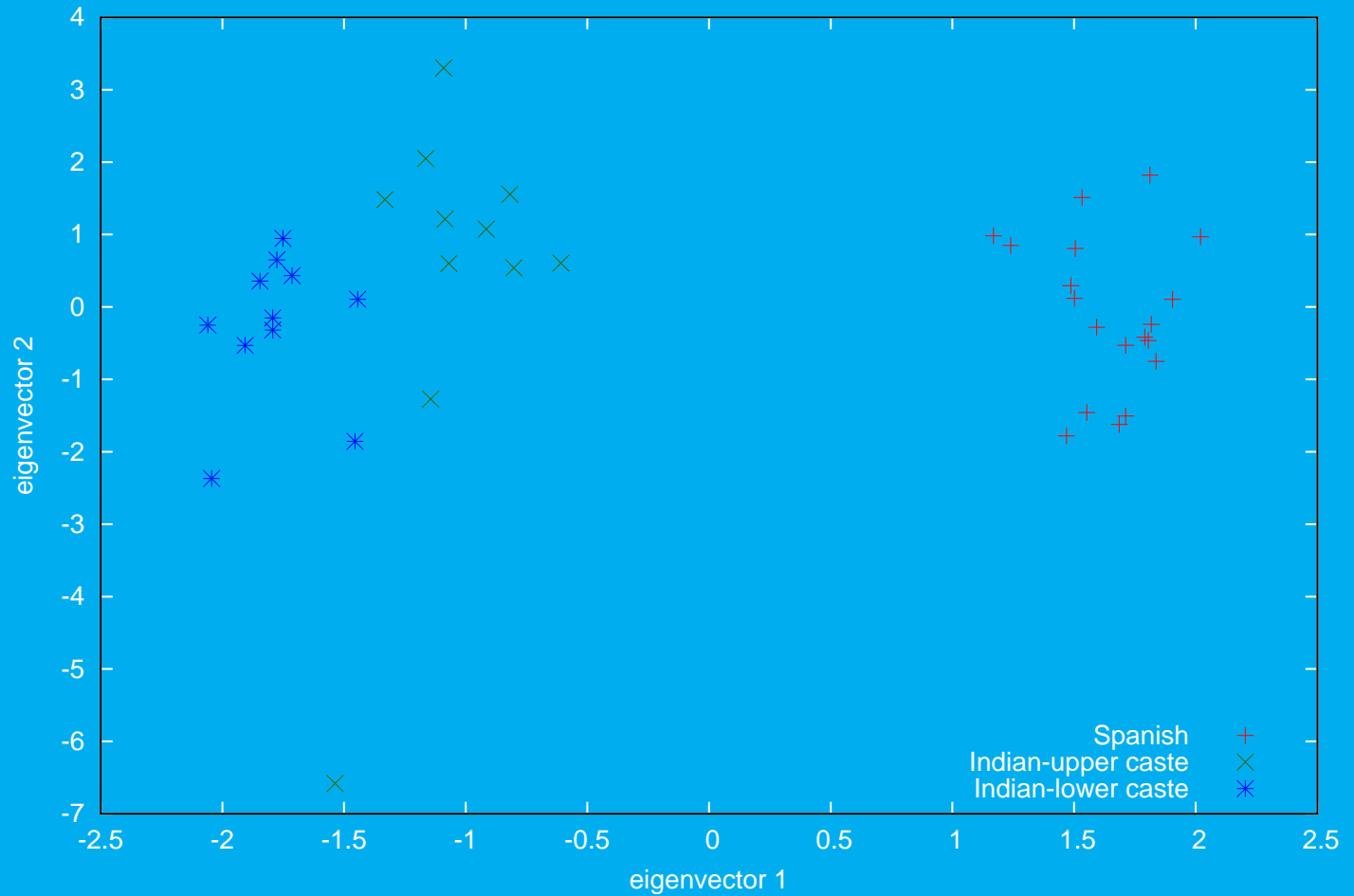
Sample covariance:

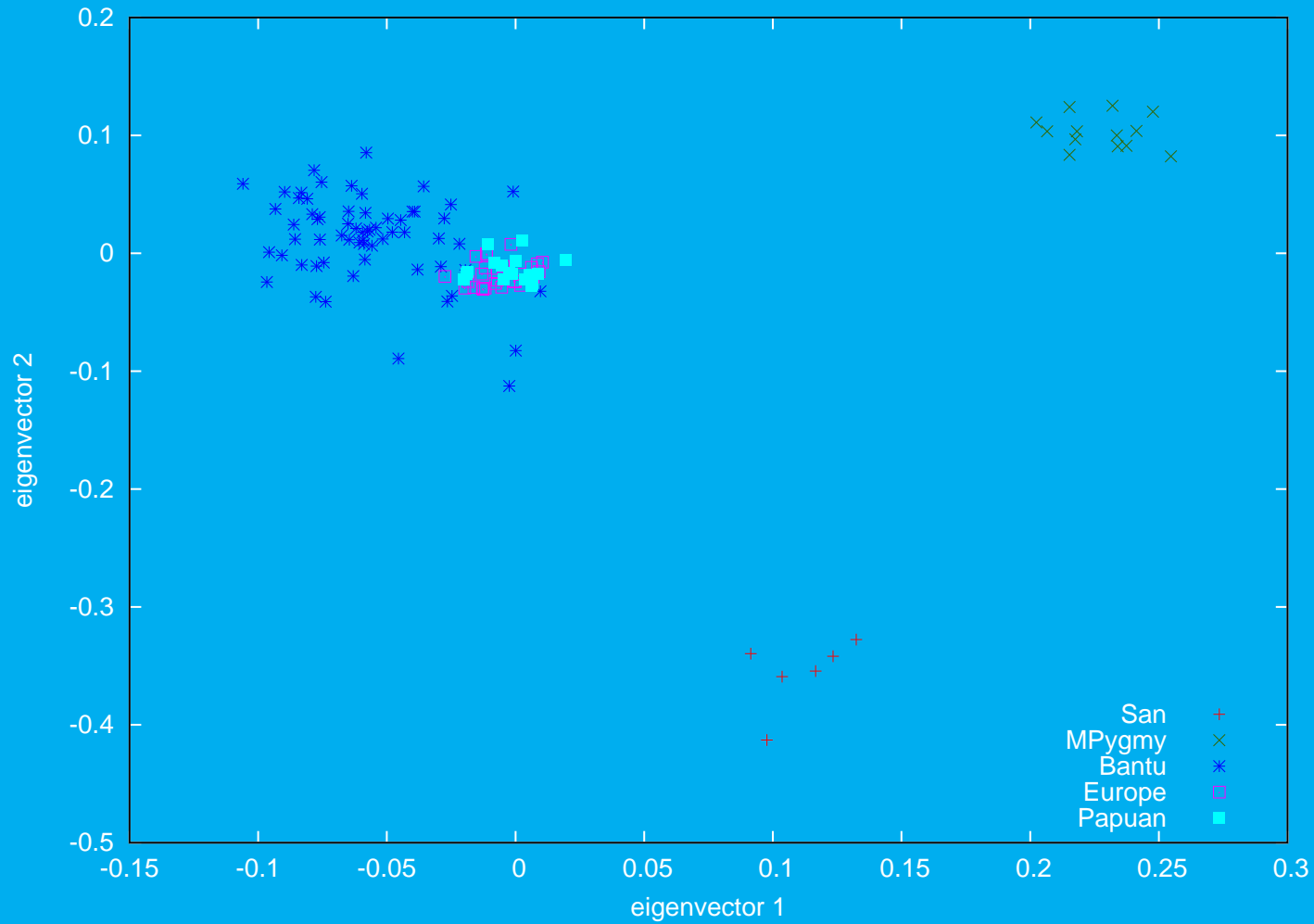
$(K - 1)$ 'large' eigenvalues

1 zero

Assumptions analyzed by Johnstone (2001)

A Spanish and two Indian populations





Missing data

Missing data is problematic

- Sample handling different for different pops.
- Some missing data genuine (pop. dependent deletions)
- Informative missingness (Clayton)
- In particular worse quality DNA \implies hets often called missing.

