# Analysing SNP data

## PCA and GWAS

M.Timothy.Rabanus-Wallace 🐦 @mtrw85

# Population Genetic Diversity: An evolutionary legacy



Time →

M.Timothy.Rabanus-Wallace 🐦 @mtrw85

# Population Genetic Diversity:
# An evolutionary legacy



Population Size

Time →

M.Timothy.Rabanus-Wallace 🐦 @mtrw85

# Population Genetic Diversity: An evolutionary legacy



Population Size

Medium Population Size
Lots of interbreeding

A PCA plot
- Distances approximate the genetic similarity between individuals

Time →

# Population Genetic Diversity: An evolutionary legacy

Population Size

Population bottleneck
Causes a loss of diversity

Time →

M.Timothy.Rabanus-Wallace 🐦 @mtrw85

# Population Genetic Diversity:
# An evolutionary legacy



Population Size

Population separation
Causes genetic isolation

M.Timothy.Rabanus-Wallace  @mtrw85

# Population Genetic Diversity:
# An evolutionary legacy



Population Size

Large population gains diversity
Small population has less diversity
Long period of isolation

M.Timothy.Rabanus-Wallace 🐦 @mtrw85

# Population Genetic Diversity: An evolutionary legacy



Unknown

Known

M.Timothy.Rabanus-Wallace 🐦 @mtrw85
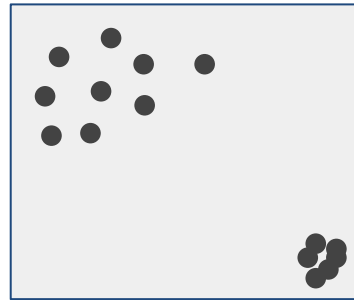
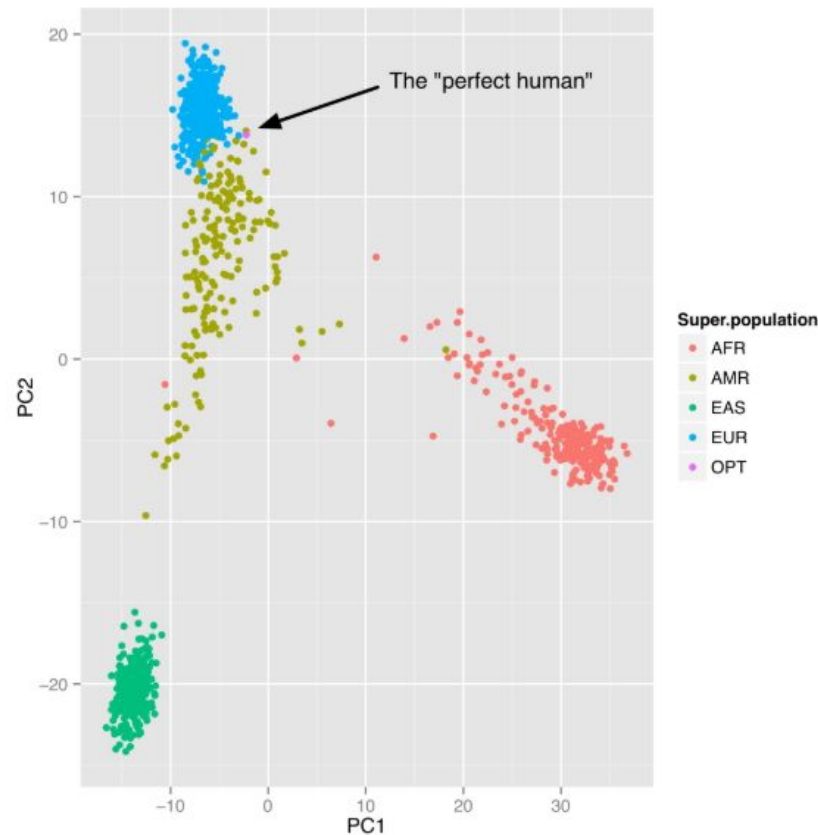# Population Genetic Diversity:
# A valuable resource

Each individual's genetic material is the product of millions of years of evolution to to tolerate different ...

- Temperatures
- Moisture levels
- Viruses
- Soil conditions
- Day lengths
- Cold periods or frost
- Wind conditions
- Bacteria
- Nematodes, insects
- Fungal pathogens
- Seasonal extremes
- Droughts
- Nutrient deficiencies
- Nutrient excesses
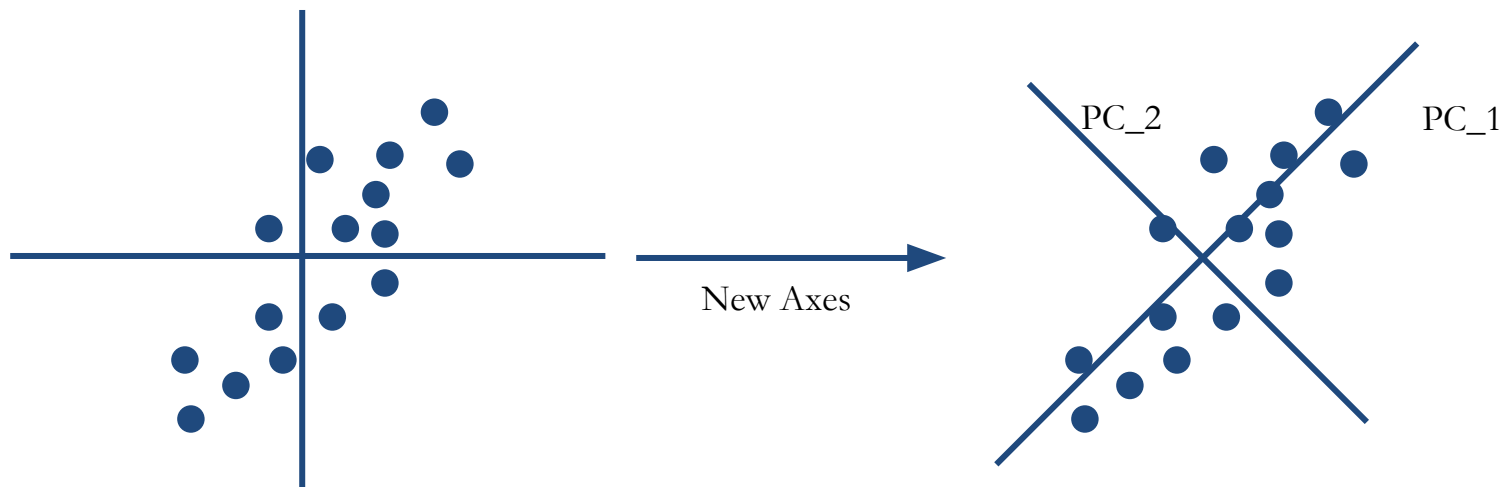- Salt levels
- Herbivores
- Symbioses
  - … etc etc etc

M.Timothy.Rabanus-Wallace 🐦 @mtrw85

# PCA with SNP data:

- An excellent way to visualise diversity, and
- An efficient mathematical way to specify population structure

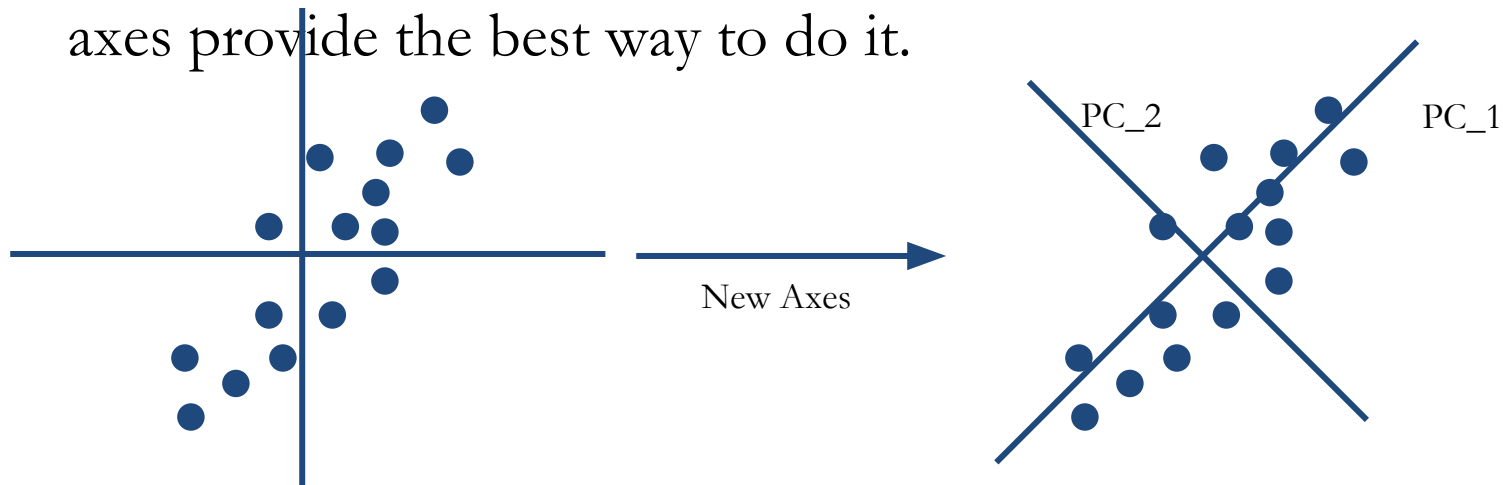M.Timothy.Rabanus-Wallace 🐦 @mtrw85

# PCA

- A way of specifying new "axes" to the data, so that the new axes express the variation in the data best.



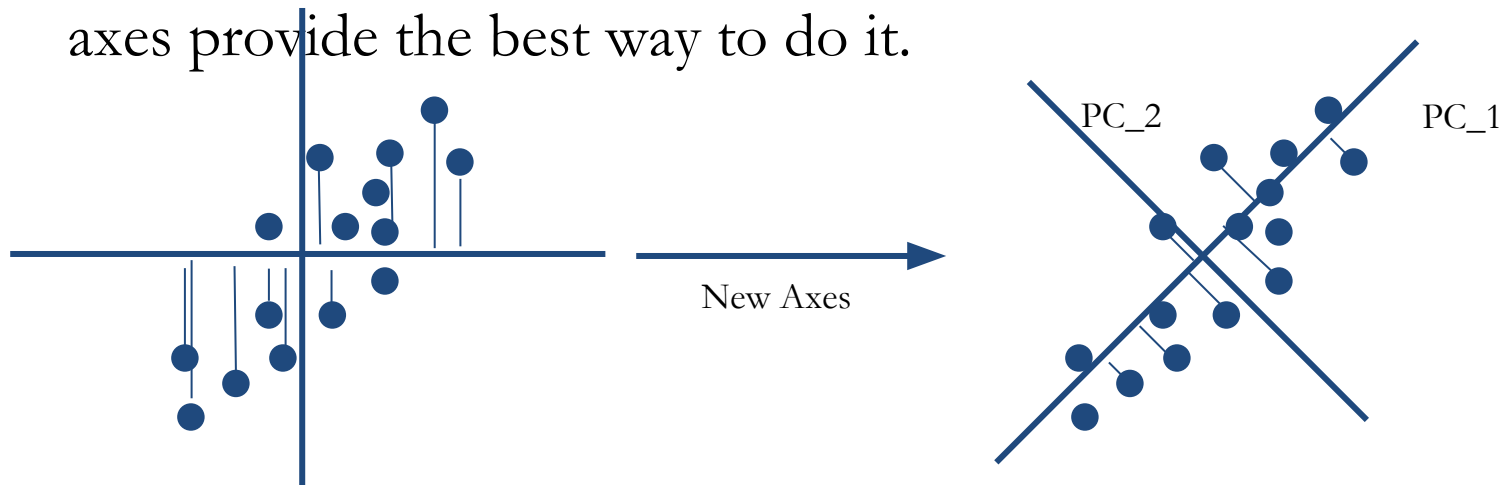New Axes

PC_2    PC_1
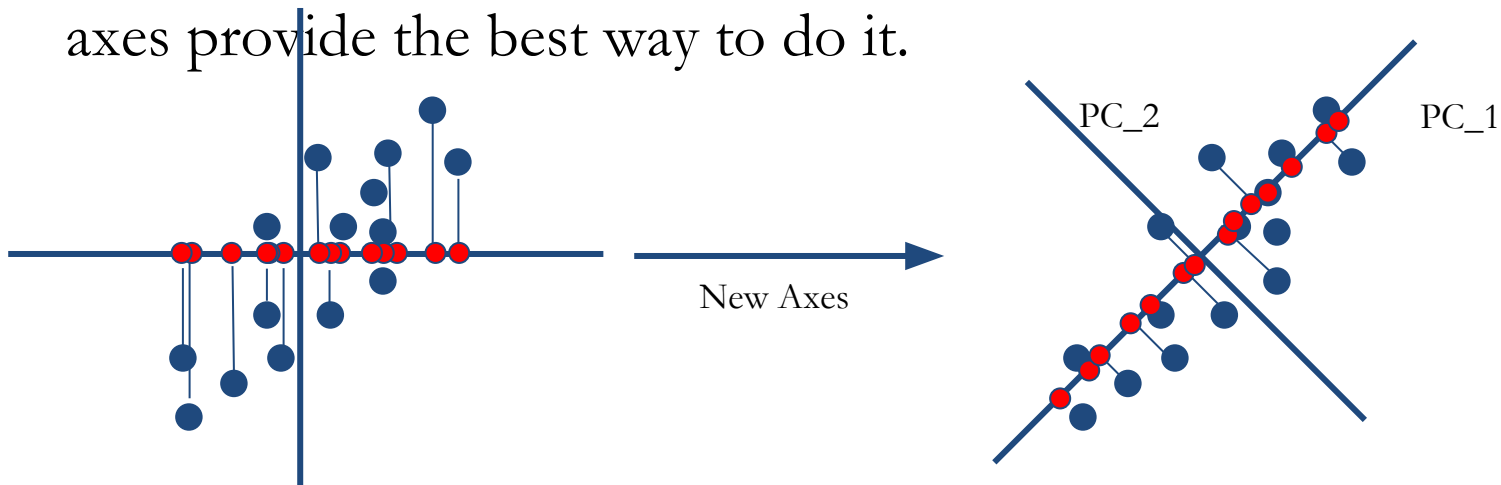
M.Timothy.Rabanus-Wallace 🐦 @mtrw85

# PCA

- A way of specifying "new axes" for the data, so that the new axes (or Principal Components) capture the highest possible variation in the data.
- If we need to describe the data using fewer dimensions, the new axes provide the best way to do it.

New Axes

PC_2    PC_1

# PCA

- A way of specifying "new axes" for the data, so that the new axes (or Principal Components) capture the highest possible variation in the data.
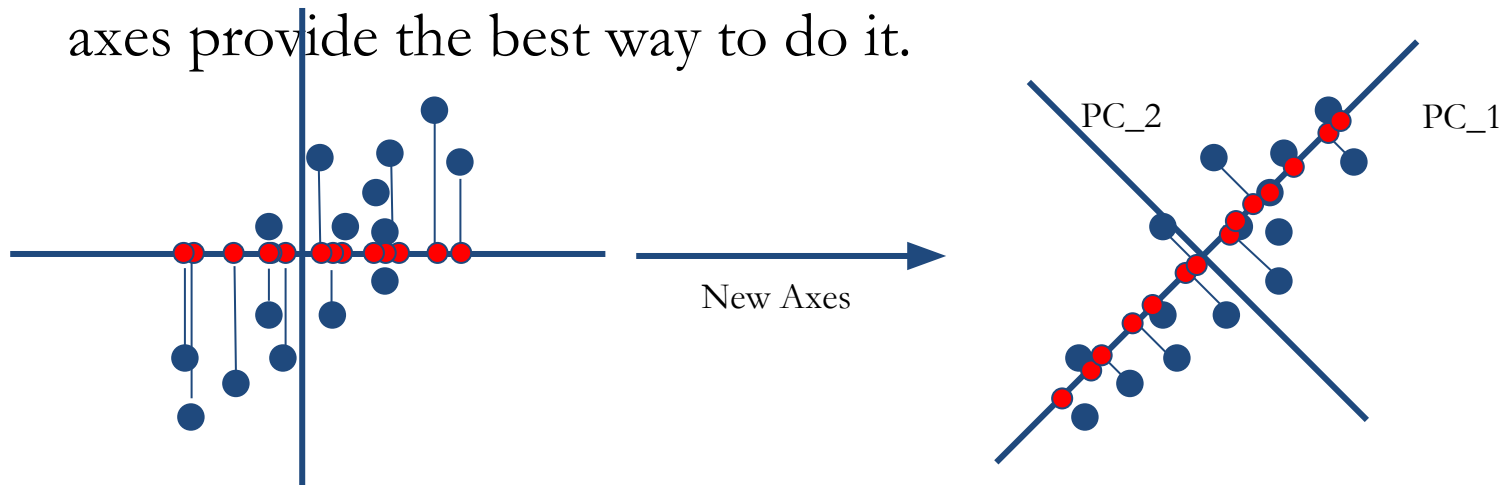- If we need to describe the data using fewer dimensions, the new axes provide the best way to do it.

New Axes

PC_2    PC_1

From 2D to 1D:
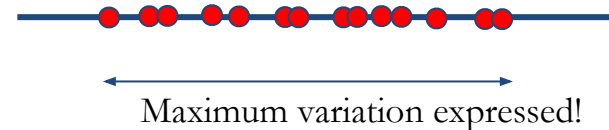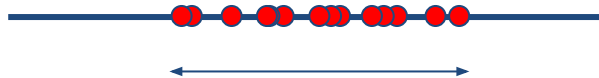
M.Timothy.Rabanus-Wallace 🐦 @mtrw85

# PCA

- A way of specifying "new axes" for the data, so that the new axes (or Principal Components) capture the highest possible variation in the data.
- If we need to describe the data using fewer dimensions, the new axes provide the best way to do it.



New Axes

PC_2    PC_1

From 2D to 1D:

# PCA

- A way of specifying "new axes" for the data, so that the new axes (or Principal Components) capture the highest possible variation in the data.
- If we need to describe the data using fewer dimensions, the new axes provide the best way to do it.

New Axes →

PC_2    PC_1

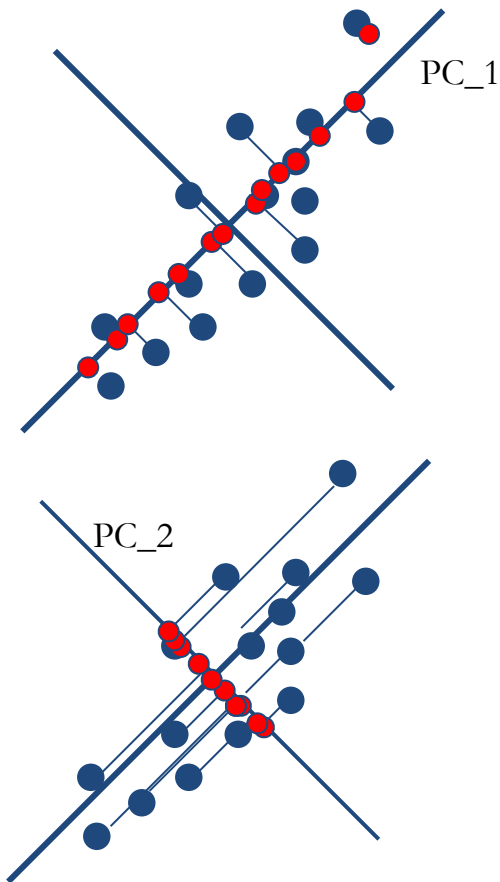From 2D to 1D:

Maximum variation expressed!

# PCA

- A way of specifying "new axes" for the data, so that the new axes (or Principal Components) capture the highest possible variation in the data.
- If we need to describe the data using fewer dimensions, the new axes provide the best way to do it.
- From 3D to 2D:

https://www.google.com/search?biw=1440&bih=767&tbm=isch&sa=1&ei=cw2cXLOzAZGblwT-o6TQCQ&q=eigenvectors+3d+gif&oq=eigenvectors+3d+gif&gs_l=img.3...551553.556272..556385...0.0..0.84.386.6......1....1..gws-wiz-img.2TdLFyVmTjY#imgrc=30k1xua4SqZHkM:

M.Timothy.Rabanus-Wallace 🐦 @mtrw85

# PCA

- A way of specifying "new axes" for the data, so that the new axes (or Principal Components) capture the highest possible variation in the data.
- If we need to describe the data using fewer dimensions, the new axes provide the best way to do it.
- From 3D to 2D
- A SNP dataset has as many "dimensions" as there are SNPs.
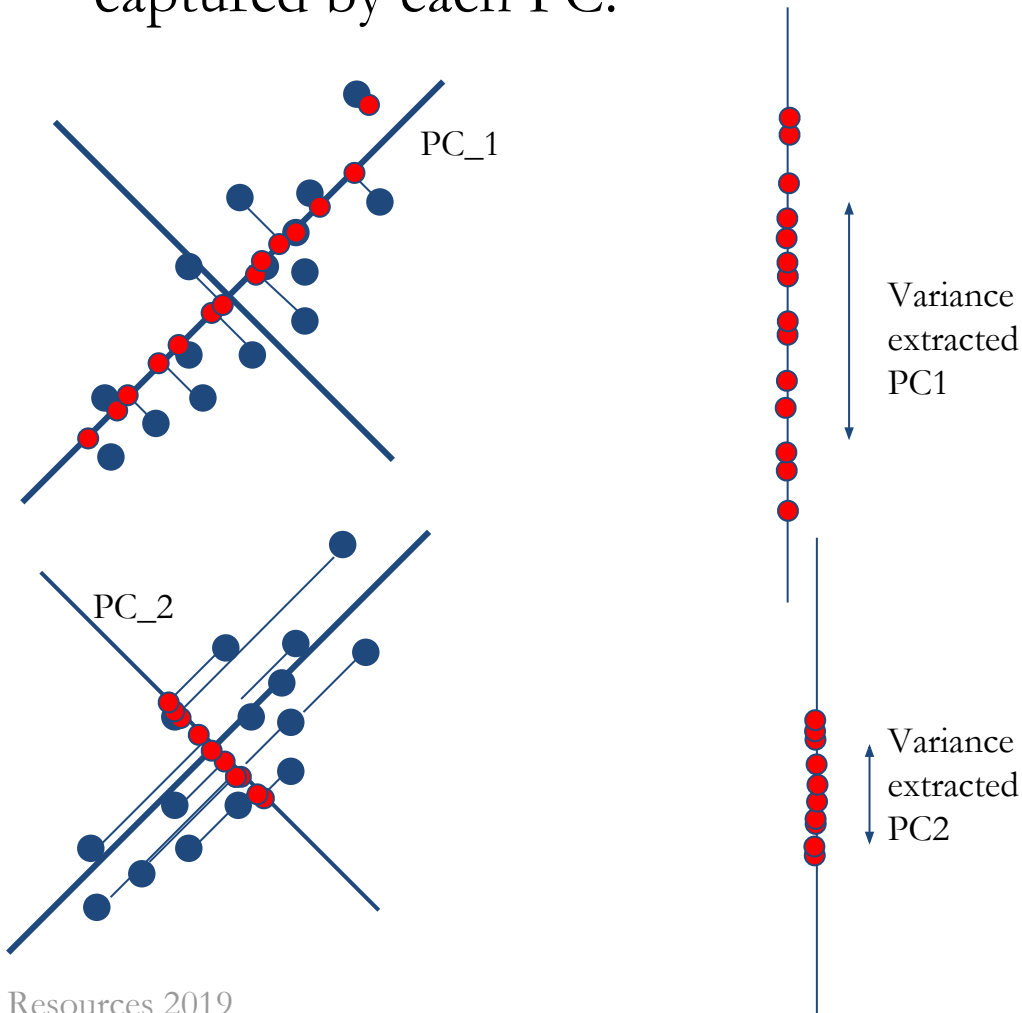  - And the values in it are discrete (e.g., 0, 1, and 2)
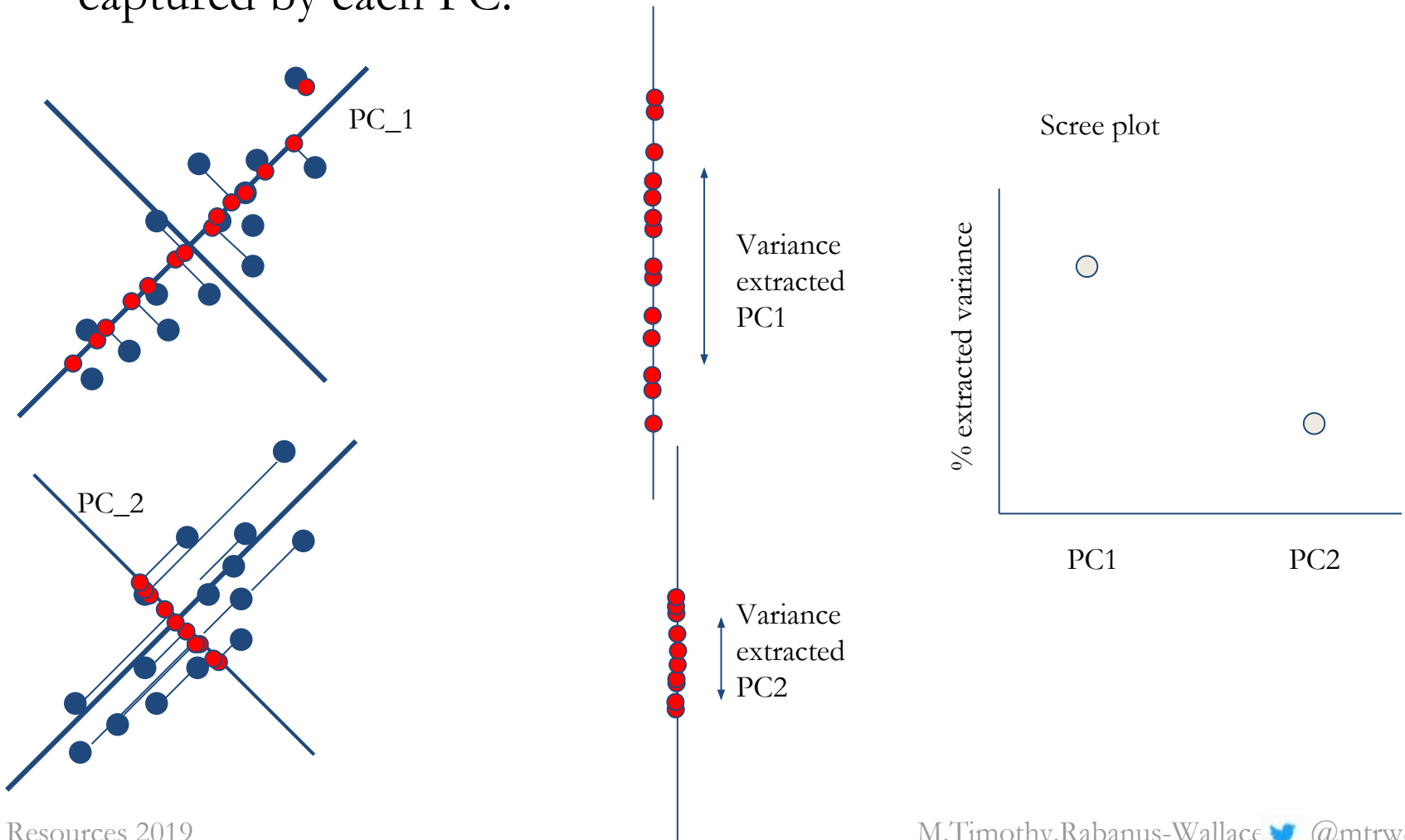
M.Timothy.Rabanus-Wallace 🐦 @mtrw85

# PCA

- A "scree plot" shows how much of the variation in the data is captured by each PC.



PC_1

PC_2

# PCA

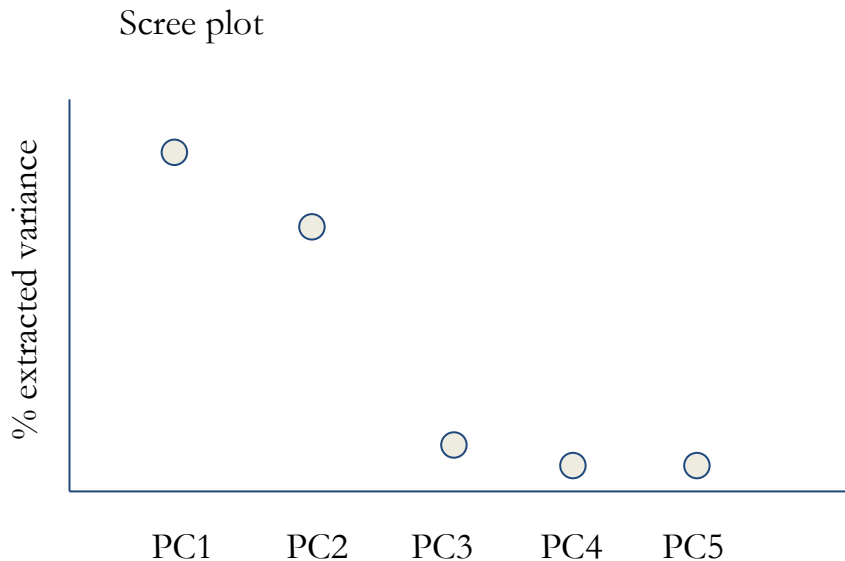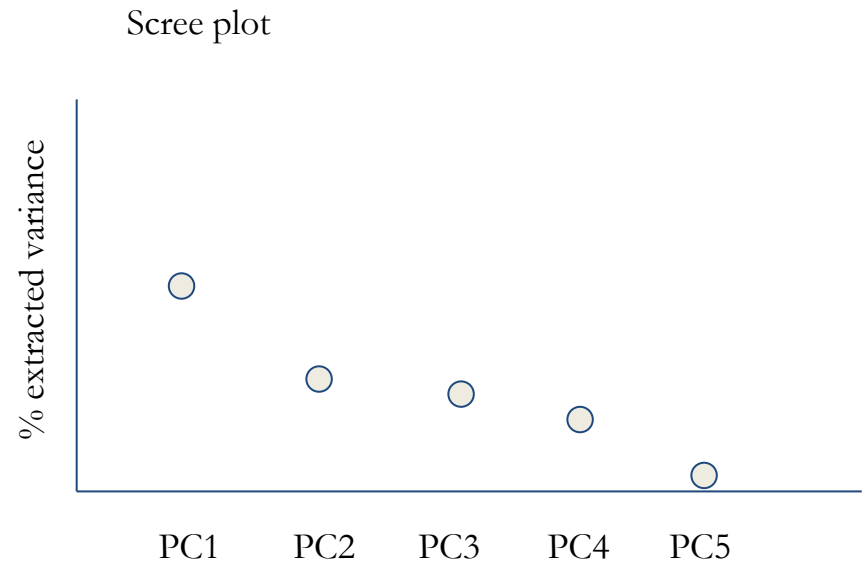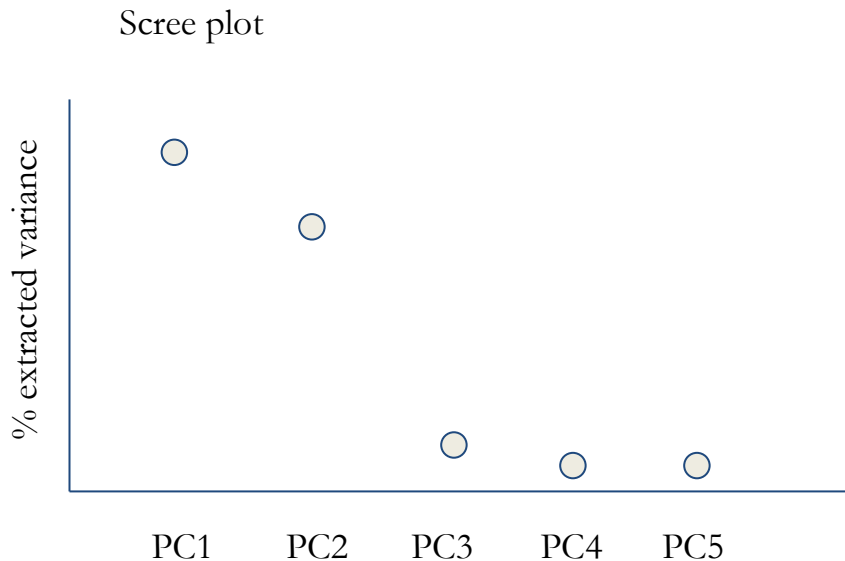- A "scree plot" shows how much of the variation in the data is captured by each PC.



PC_1

PC_2

Variance
extracted
PC1

Variance
extracted
PC2

# PCA

- A "scree plot" shows how much of the variation in the data is captured by each PC.



PC_1

PC_2

Variance extracted PC1

Variance extracted PC2

Scree plot

% extracted variance

PC1          PC2

# PCA

- A "scree plot" shows how much of the variation in the data is captured by each PC.
  - Interpreting scree plots



Scree plot

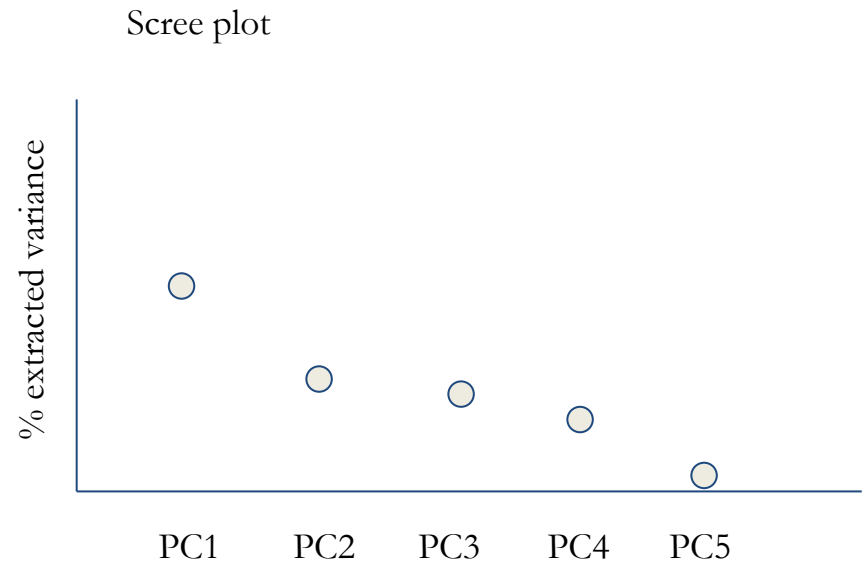PCs 1 and 2 can be used to describe the dataset very well, because they explain most of the variance!

Need more PCs to describe most of the variation in the data. Use PCs 1--4?

# PCA

- A "scree plot" shows how much of the variation in the data is captured by each PC.
  - Interpreting scree plots



Scree plot

PCs 1 and 2 can be used to describe the dataset very well, because they explain most of the variance!

Need more PCs to describe most of the variation in the data. Use PCs 1--4?

A geneticist might use a rule like "use as many PCs as needed to describe 95% of the variation".

M.Timothy.Rabanus-Wallace 🐦 @mtrw85

# Genetic Associations

- Testing whether genetic markers (SNPs in our case) are associated with a phenotype.
  - That's 'associated with' … not 'cause'!

M.Timothy.Rabanus-Wallace 🐦 @mtrw85

# Genetic Associations

- Testing whether genetic markers (SNPs in our case) are associated with a phenotype.
  - That's 'associated with' … not 'cause'!
  - Association suggests *linkage to* the cause of the phenotype

M.Timothy.Rabanus-Wallace @mtrw85

# Genetic Associations

- Testing whether genetic markers (SNPs in our case) are associated with a phenotype.
  - That's 'associated with' … not 'cause'!
  - Association suggests *linkage to* the cause of the phenotype
  - Means the SNP *may* be physically near a gene that causes the phenotype

M.Timothy.Rabanus-Wallace 🐦 @mtrw85

# Genetic Associations

- Testing whether genetic markers (SNPs in our case) are associated with a phenotype.
  - That's 'associated with' … not 'cause'!
  - Association suggests *linkage to* the cause of the phenotype
  - Means the SNP **may** be genetically near a gene that causes the phenotype
  - Means the SNP allele can be used as a marker to screen for the (probable) presence of the causal gene(s)

# Genetic Associations

Intuition

| Phenotype | Genotype |
|-----------|----------|
| 🌶️ | 2 |
| 🌶️ | 0 |
| 🌶️ | 2 |
| 🌶️ | 2 |
| 🌶️ | 0 |
| 🌶️ | 0 |
| 🌶️ | 0 |

Strong association!

M.Timothy.Rabanus-Wallace 🐦 @mtrw85

# Genetic Associations

Intuition

| Phenotype | Genotype |
|-----------|----------|
| 🌶️ (red) | 0 |
| 🌶️ (green) | 0 |
| 🌶️ (red) | 2 |
| 🌶️ (red) | 0 |
| 🌶️ (green) | 0 |
| 🌶️ (green) | 2 |
| 🌶️ (green) | 0 |

Weak / no association

M.Timothy.Rabanus-Wallace @mtrw85

# Genetic Associations

Intuition

| Phenotype | Genotype |
|-----------|----------|
|  | 2 |
|  | 0 |
|  | 2 |
|  | 2 |
|  | 0 |
|  | 2 |
|  | 0 |

Maybe associated .... ?

But how do we judge what is significant?

M.Timothy.Rabanus-Wallace 🐦 @mtrw85

# Genetic Associations

Any appropriate statistical test can be used to judge whether a SNP is likely associated with a phenotype, e.g.:

## Chi squared

**Phenotype**

|  | 2 | 0 |
|---|---|---|
| 🌶️ (red) | 20 | 52 |
| 🌶️ (green) | 32 | 32 |

**Phenotype** (vertical, left axis)

Null hypothesis: The genotype and the phenotype are independent

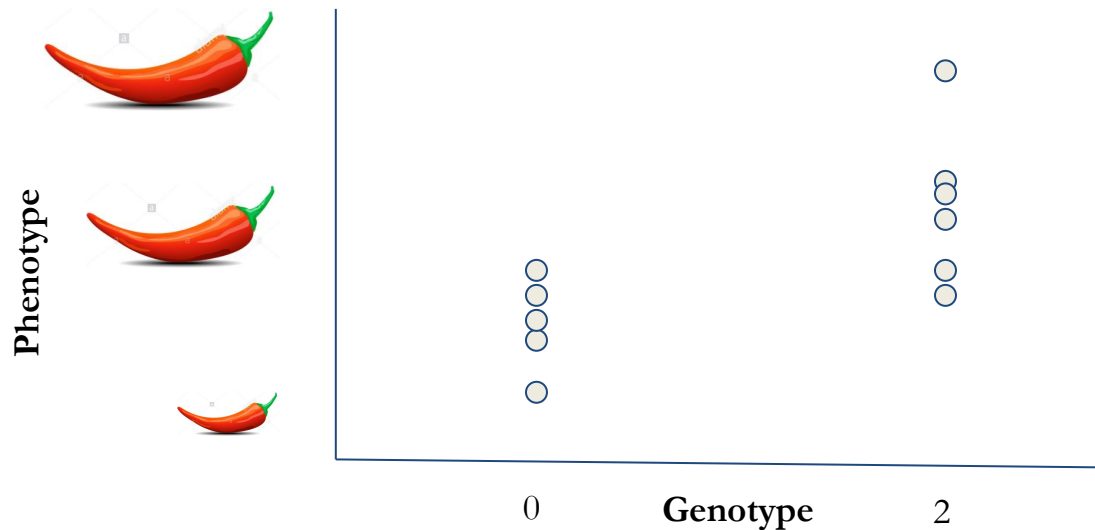… generate a p-value …

# Genetic Associations

Any appropriate statistical test can be used to judge whether a SNP is likely associated with a phenotype, e.g.:

## t-test



Null Hypothesis: The mean sizes of plants of the two genotypes are equal

… generate a p-value ….

M.Timothy.Rabanus-Wallace 🐦 @mtrw85

# Genetic Associations

Any appropriate statistical test can be used to judge whether a SNP is likely associated with a phenotype, e.g.:

**We will use a linear regression**
- Fast
- Helps account for population structure
- *You are not required to know details*, just:
  - Null hypothesis: "The SNP allele is required to explain the phenotype."

M.Timothy.Rabanus-Wallace 🐦 @mtrw85

# Genetic Associations

Any appropriate statistical test can be used to judge whether a SNP is likely associated with a phenotype, e.g.:

**We will use a linear regression**
- Fast
- Helps account for population structure
- Not required to know details, just:
  - Null hypothesis: "The SNP allele is required to explain the phenotype."
- The test is done for every SNP we chose in the genome … hence

Genome Wide Association Study (GWAS)

# Genetic Associations

Any appropriate statistical test can be used to judge whether a SNP is likely associated with a phenotype, e.g.:

**We will use a linear regression**
- Fast
- Helps account for population structure
- Not required to know details, just:
  - Null hypothesis: "The SNP allele is required to explain the phenotype."
- The test is done for every SNP we chose in the genome … hence

Genome Wide Association Study (GWAS)

(we just use chromosome 1H)

M.Timothy.Rabanus-Wallace 🐦 @mtrw85

# Associations and population structure



CGATTCG**T**GCGGGGCTCCTCTCAGGATGCTT**A**AA
CGATTCG**T**GCGGGGCTCCTCTCAGGATGCTT**A**AA
CGATTCG**T**GCGGGGCTCCTCTCAGGATGCTT**A**AA
CGATTCG**T**GCGGGGCTCCTCTCAGGATGCTT**G**AA
CGATTCG**T**GCGGGGCTCCTCTCAGGATGCTT**G**AA
CGATTCG**C**GCGGGGCTCCTCTCAGGATGCTT**G**AA
CGATTCG**C**GCGGGGCTCCTCTCAGGATGCTT**G**AA
CGATTCG**C**GCGGGGCTCCTCTCAGGATGCTT**G**AA
CGATTCG**C**GCGGGGCTCCTCTCAGGATGCTT**G**AA
CGATTCG**C**GCGGGGCTCCTCTCAGGATGCTT**G**AA

Does either SNP associate with fruit size?

M.Timothy.Rabanus-Wallace 🐦 @mtrw85

# Associations and population structure



What about now?

M.Timothy.Rabanus-Wallace 🐦 @mtrw85

# Associations and population structure

**A PCA plot would look a bit like this ...**



PC_1

PC_2

M.Timothy.Rabanus-Wallace ⚫ @mtrw85

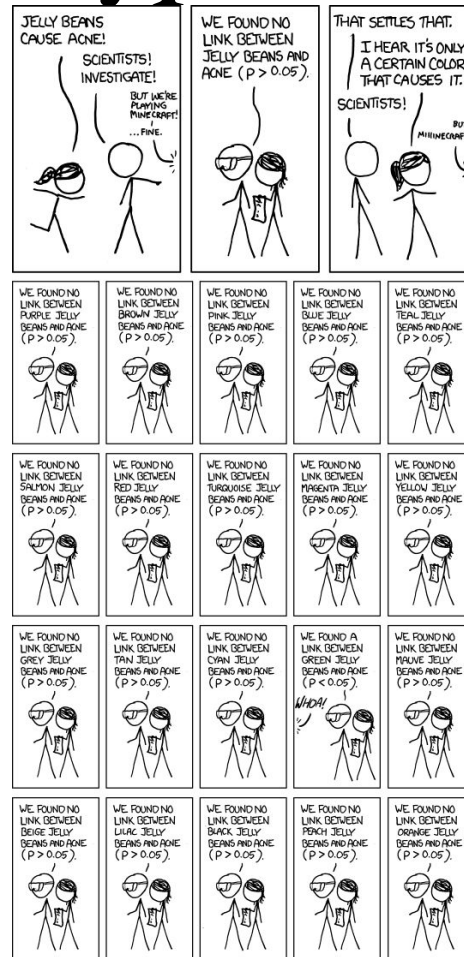# Associations and population structure

PC_1

PC_2

- The principal components are able to predict the phenotype!
  - i.e., population structure is a strong influence
- We can control for this in the linear model
  - The model will test how much extra predictive power the SNP allele gives us, when the population structure (summarised by the PCs) are also used to predict the phenotype.

# Multiple hypothesis testing ...

# Multiple hypothesis testing ...



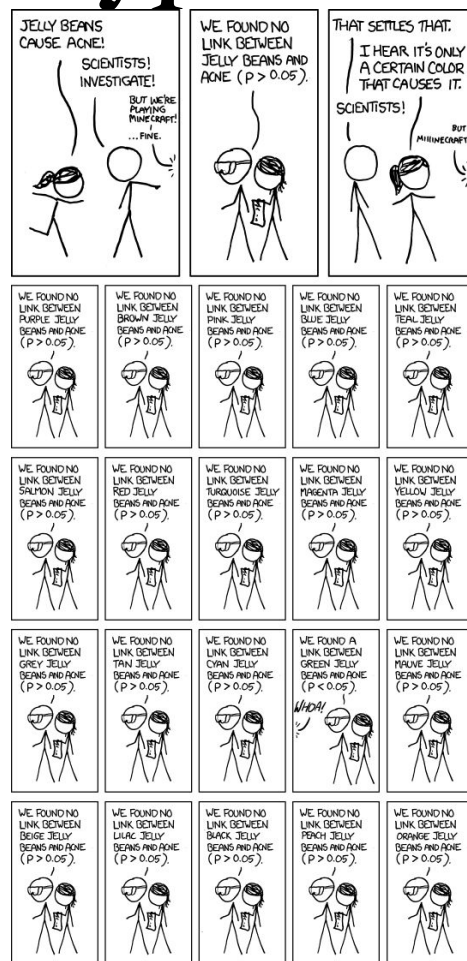- The normal "significance threshold" of p-value = 0.05 causes false positives when we test many hypotheses

# Multiple hypothesis testing ...



- The normal "significance threshold" of p-value < 0.05 causes false positives when we test many hypotheses.
- We correct for this by setting the threshold much more stringently.
- The "Bonferroni correction" involves simply dividing the p-value significance threshold (0.05) by the number of tests.
- The number of tests is the number of SNPs.
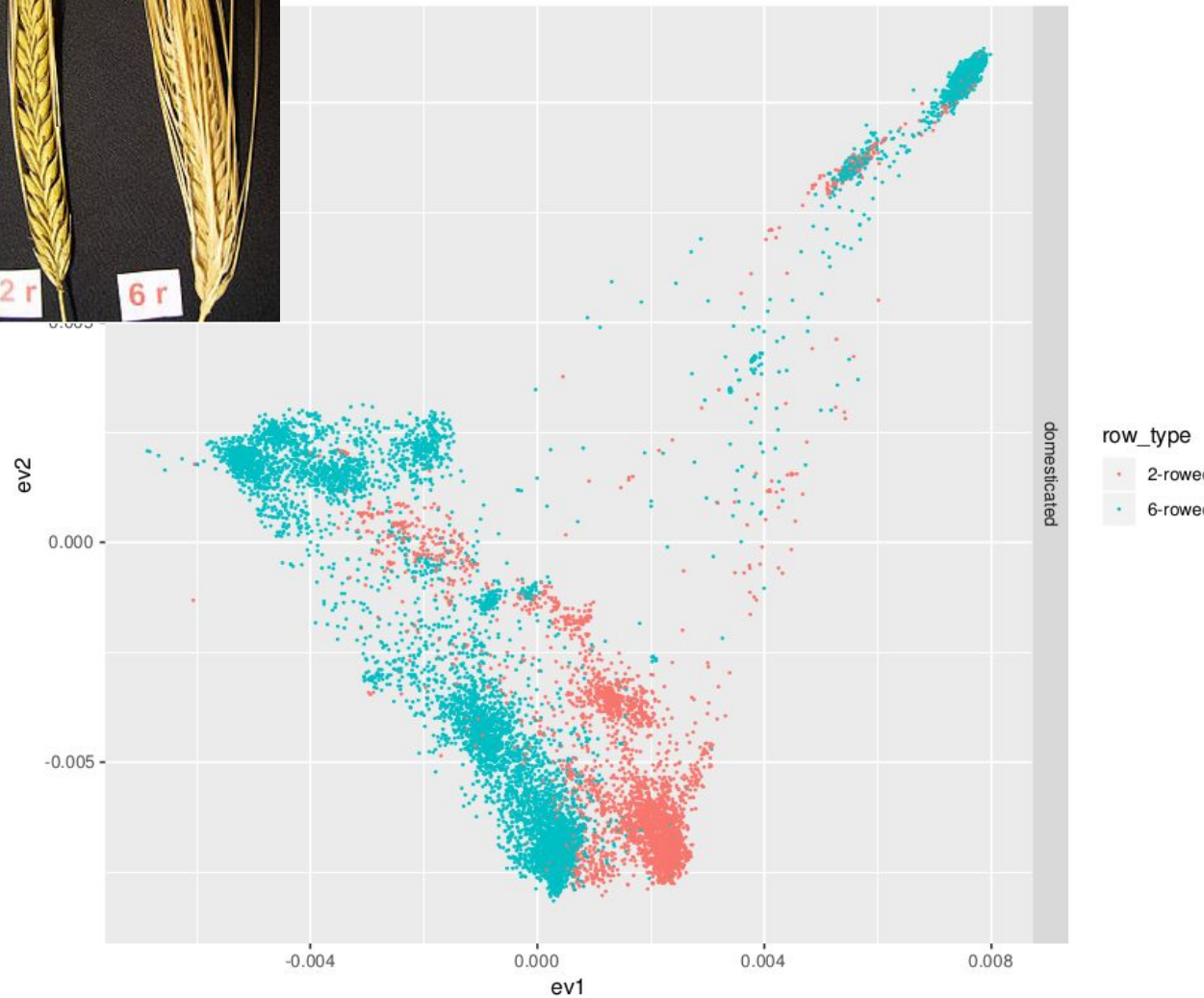
M.Timothy.Rabanus-Wallace  @mtrw85

# Association testing in this session

Test for differences in row type
(6-rowed vs 2-rowed)

M.Timothy.Rabanus-Wallace 🐦 @mtrw85

# Association testing in this session

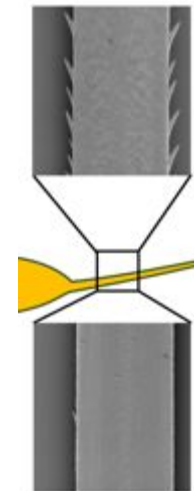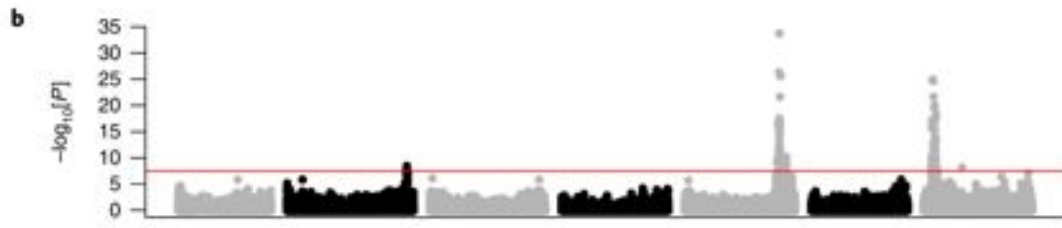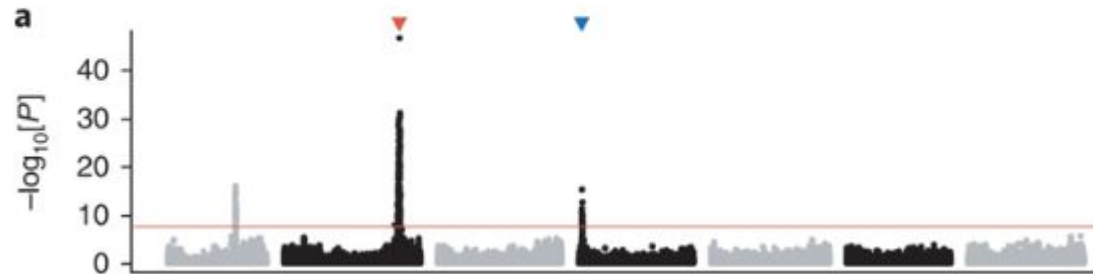Test for differences in row type
(6-rowed vs 2-rowed)



The collection our sample came from
shows some relationship between
row-type and population structure!

# Association testing in this session

Previous GWAS results show genetic associations with awn barbs and row type



Fig. 3: Genome-wide association scans for morphological characters.

https://www.nature.com/articles/s41588-018-0266-x
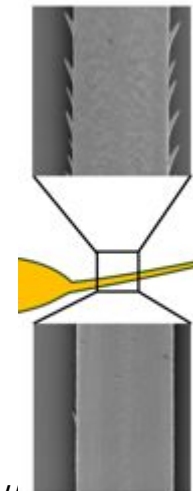
M.Timothy.Rabanus-Wallace 🐦 @mtrw85

# Association testing in this session

Previous GWAS results show genetic associations with row type (and also awn barbs, below)



Fig. 3: Genome-wide association scans for morphological characters.

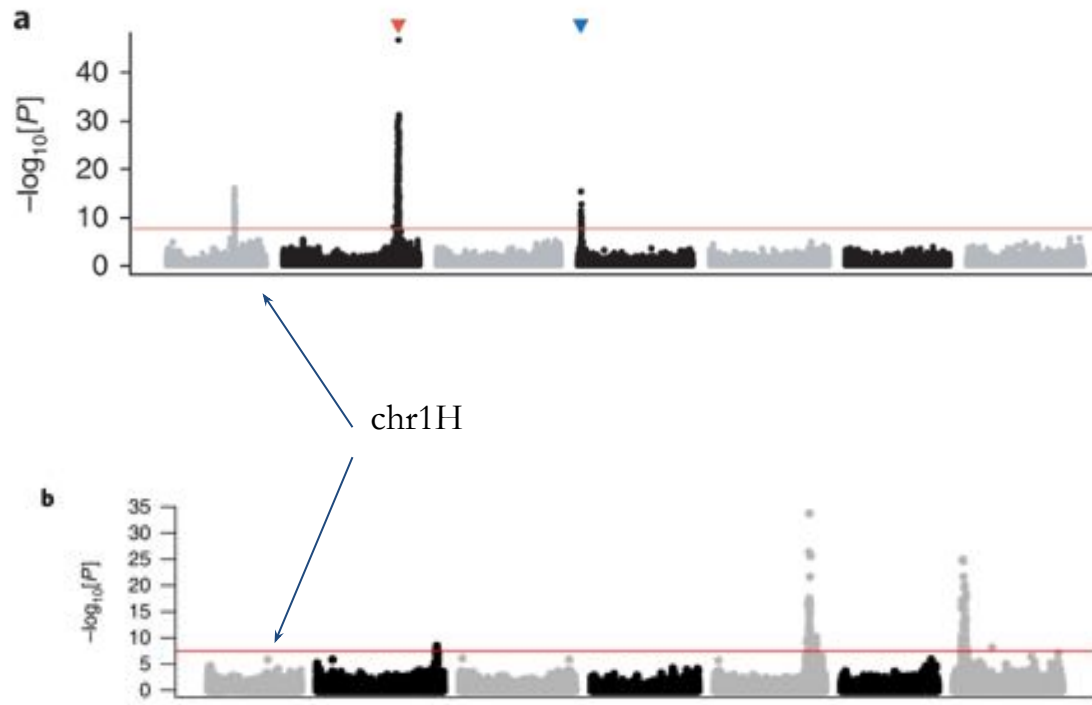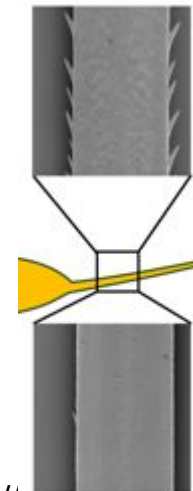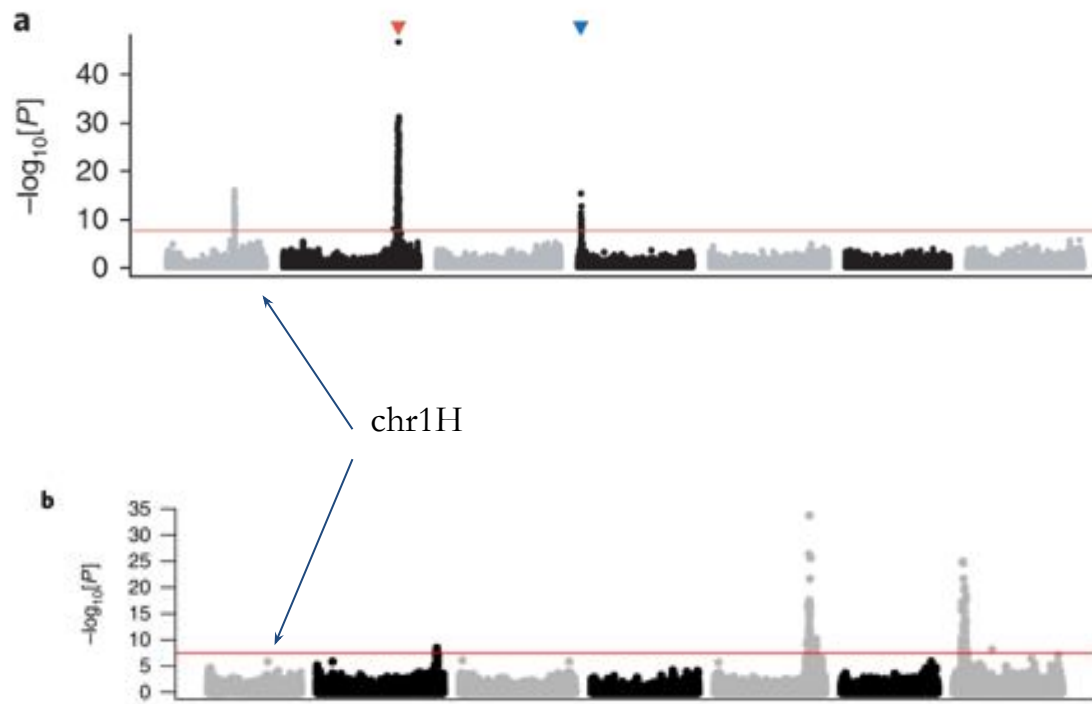chr1H

https://www.nature.com/articles/s41588-018-0266-x

# Association testing in this session

Previous GWAS results show genetic associations with row type (and also awn barbs, below)



Fig. 3: Genome-wide association scans for morphological characters.

chr1H

We are unlikely to find any strong correlations owing to small sample size ...

https://www.nature.com/articles/s41588-018-0266-x

# Association testing in this session

So … primary aims:

- Load our cleaned SNP data into R
- Create a PCA to summarise the population structure between samples
  - Plot a PCA plot!
- Link it to phenotype data from a database
  - Using a new data.table trick …
  - More PCA plots …
- Link three datasets: **The PCs** that summarise the population structure, **the SNPs**, and **the phenotype data**
- Run a GWAS to test associations at each SNP
- Plot the GWAS results as a Manhattan Plot with the significance threshold shown on the plot.