

# Polymorphism and Variant Analysis

Matt Hudson

University of Illinois

# Mutations and Variations

- In this class we will cover:
  - SNPs & SNVs
  - Mutant impact analysis (predicting when a mutation might be damaging)
  - Uses of machine learning in genetics
  - Genome-wide association studies (GWAS)

# What is a SNP ? And a SNV?

- Single nucleotide polymorphism
- Single nucleotide variant

```
I1 : AACGAGCTAGCGATCGATCGACTACGACTACGAGGT
I2 : AACGAGCTAGCGATCGATCGACAACGACTACGAGGT
I3 : AACGAGCTAGCGATCGATCGACTACGACTACGAGGT
I4 : AACGAGCTAGCGATCGATCGACTACGACTACGAGGT
I5 : AACGAGCTAGCGATCGATCGACAACGACTACGAGGT
I6 : AACGAGCTAGCGATCGATCGACTACGACTACGAGGT
I7 : AACGAGCTAGCGATCGATCGACTACGACTACGAGGT
I8 : AACGAGCTAGCGATCGATCGACTACGACTACGAGGT
```

Individuals I2 and I5 have a variation (T -> A). This position is both.

# Notes on SNPs and SNVs

- A SNV is any old change (e.g. could be a somatic mutation in an individual, or even an artifact)
- To be called a SNP, has to be **polymorphic**
- Lots of SNPs in databases, eg. the 1000 Genomes project recorded ~41 Million SNPs by sequencing ~1000 humans.

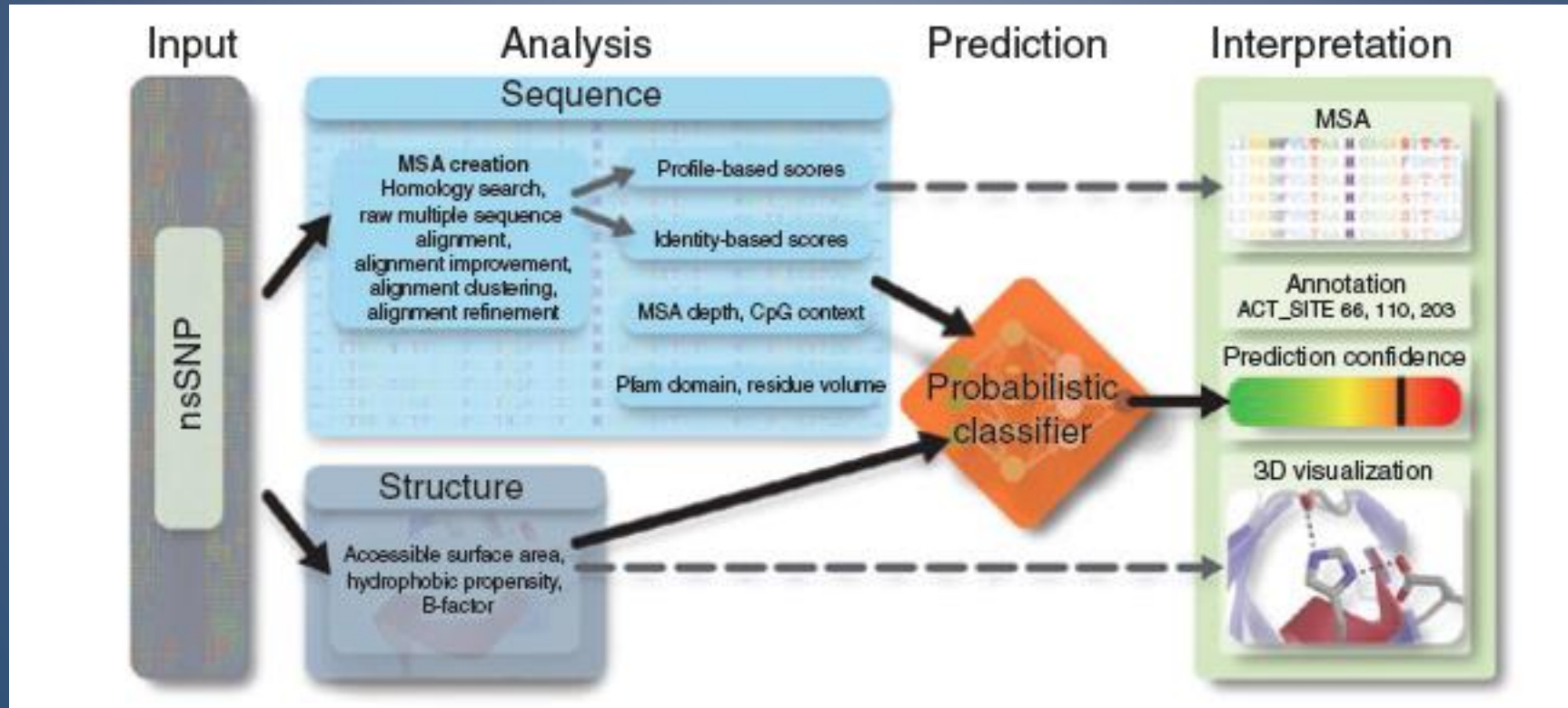
# Thus, your fields may differ

- If you are a population geneticist doing GWAS, you are generally only interested in SNPs
- If you are a cancer geneticist looking at sequence data from tumors, you are primarily interested in SNVs
- In non-human biology there can be other complications (e.g. polyploidy, HGT etc.).
- Definitions vary by field

# Predicting when a coding SNV (or SNP) is bad news

- Question:
  - *I found a SNP inside the coding sequence. Knowing how to translate the gene sequence to a protein sequence, I discovered that this is a non-synonymous change, i.e., the encoded amino acid changes. This is an nsSNP.*
  - *Will that impact the protein's function?*
  - *(And I don't quite know how the protein functions in the first place ...)*

# PolyPhen 2.0



# Data for training/evaluation

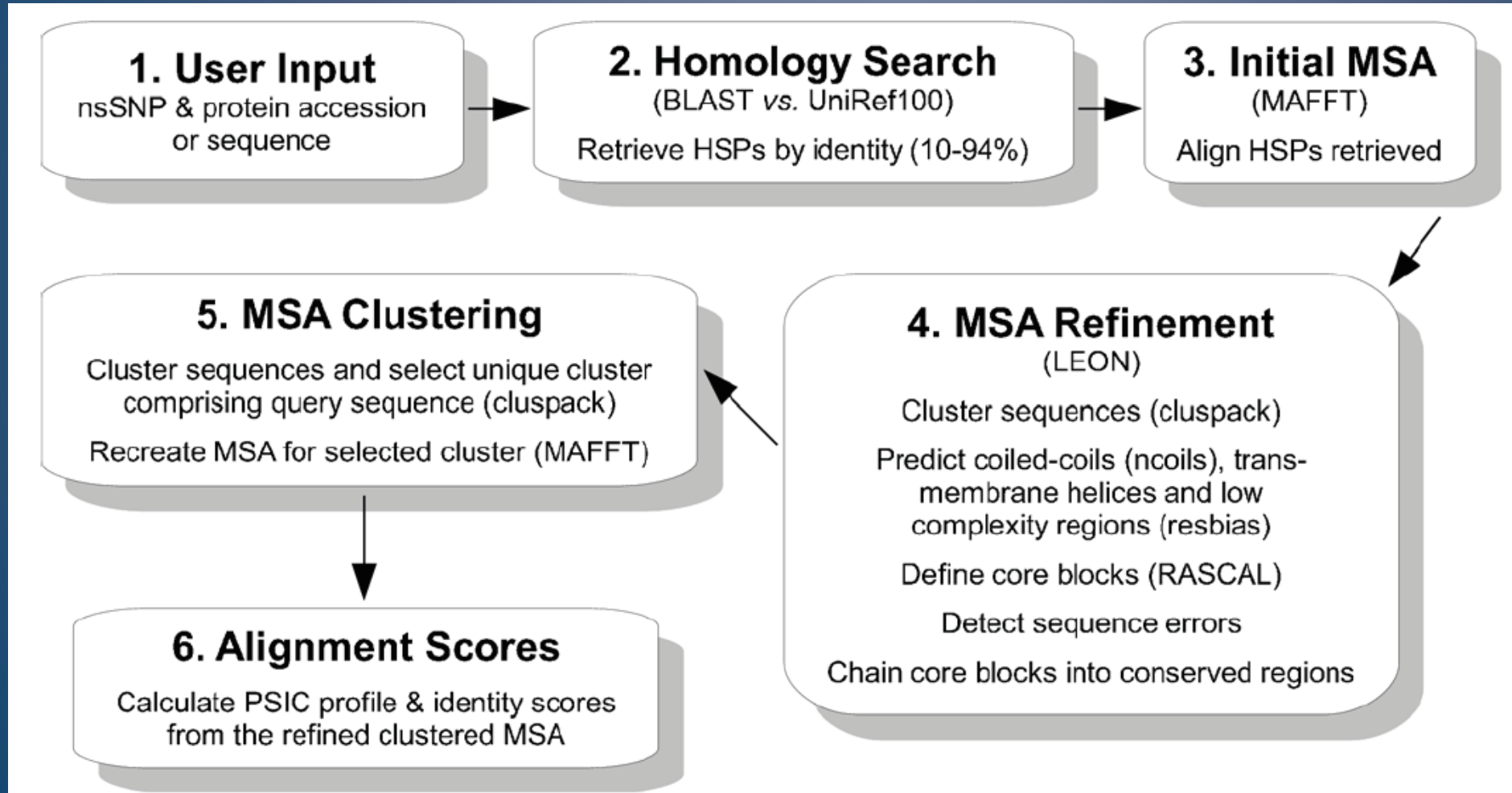
- HumDiv
  - Damaging mutations from UniProtKB. Look for annotations such as “complete loss of function”, “abolishes”, “no detectable activity”, etc.
  - Non-damaging mutations: differences in homologous proteins in closely related mammalian species



# “Features”

Name	Definition	Values with ranges in HumDiv
nt1	wild type allele nucleotide	A,C,G,T
nt2	mutation allele nucleotide	A,C,G,T
site	SITE annotation from UniProt/Swiss-Prot	Yes, No
region	REGION annotation from UniProt/Swiss-Prot	NO, PROPER, SIGNAL, TRANSMEM
phat	PHAT matrix element in the TRANSMEM region	[-8.0, 4.0], mean = -0.04
score1	PSIC score for the wild type allele	[-1.1], mean = 1.07
score2	PSIC score for the mutant allele	[-1.39, 2.64], mean = .166
score_delta	difference of PSIC scores (Score1-Score2)	[-3.23, 4.57], mean = .905
num_observ	number of residues observed at the position of the multiple alignment	[1, 432], mean 69.3
delta_volume	change in residue side chain volume	[-167, 167], mean = -1,93
transv	mutation origin by transversion or transition	Yes, No
CpG	mutation origin in the CpG hypermutable context	Yes, No
pfam_hit	position of the mutation within/outside a protein domain as defined by Pfam	Yes, No
id_p_max	congruency of the mutant allele to the multiple alignment	[0, 95.5], mean = 24
id_q_min	sequence identity with the closest homologue deviating from wild type allele	[1.56, 95.5], mean 68.76
cpgVar1 Var2	presence of the CpG context combined with wild type and mutant amino acid types	NO, AA1_AA2
cpg_transition	whether variant happened as transition in CpG context	No, Transition, Transversion
charge_change	change in electrostatic charge	0,1,2
hydroph_change	change in hydrophobicity	[0, 2.85], mean 0.80
ali_ide	sequence identity with the closest homolog with known 3D structure	[0, 1], mean 0.33
ali_len	alignment length with the closest homolog with known 3D structure	[0, 1213], mean 130.0
acc_normed	normalized accessible surface area of amino acid residue	[0, 1.55], mean .35
sec_str	secondary structure	HELIX, SHEET, OTHER
map_region	region of the Ramachandran map	ALPHA, BETA, OTHER
delta_prop	change in accessible surface area propensity	[-2.89, 2.89], mean -0.07
b_fact	crystallographic beta-factor	[-1.85, 5.17], mean 0.0
het_cont_ave_num	average number of contact with heteroatoms	Yes, No
het_cont_min_dist	minimal distance to a heteroatom	Yes, No
inter_cont_ave_num	average number of interchain contacts in a protein complex	Yes, No
inter_cont_min_dist	average minimal interchain distance	Yes, No
delta_volume_new	change in residue volume for buried residues	[-119, 138], mean -0.5
delta_prop_new	change in accessible surface area propensity for buried residues	[-1.83, 2.89], mean 0.0026

# The MSA part of the pipeline



# Position Specific Independent Count (PSIC)

- Reflects the amino acid's frequency at the specific position in sequence, given an MSA.



PMID 10360979

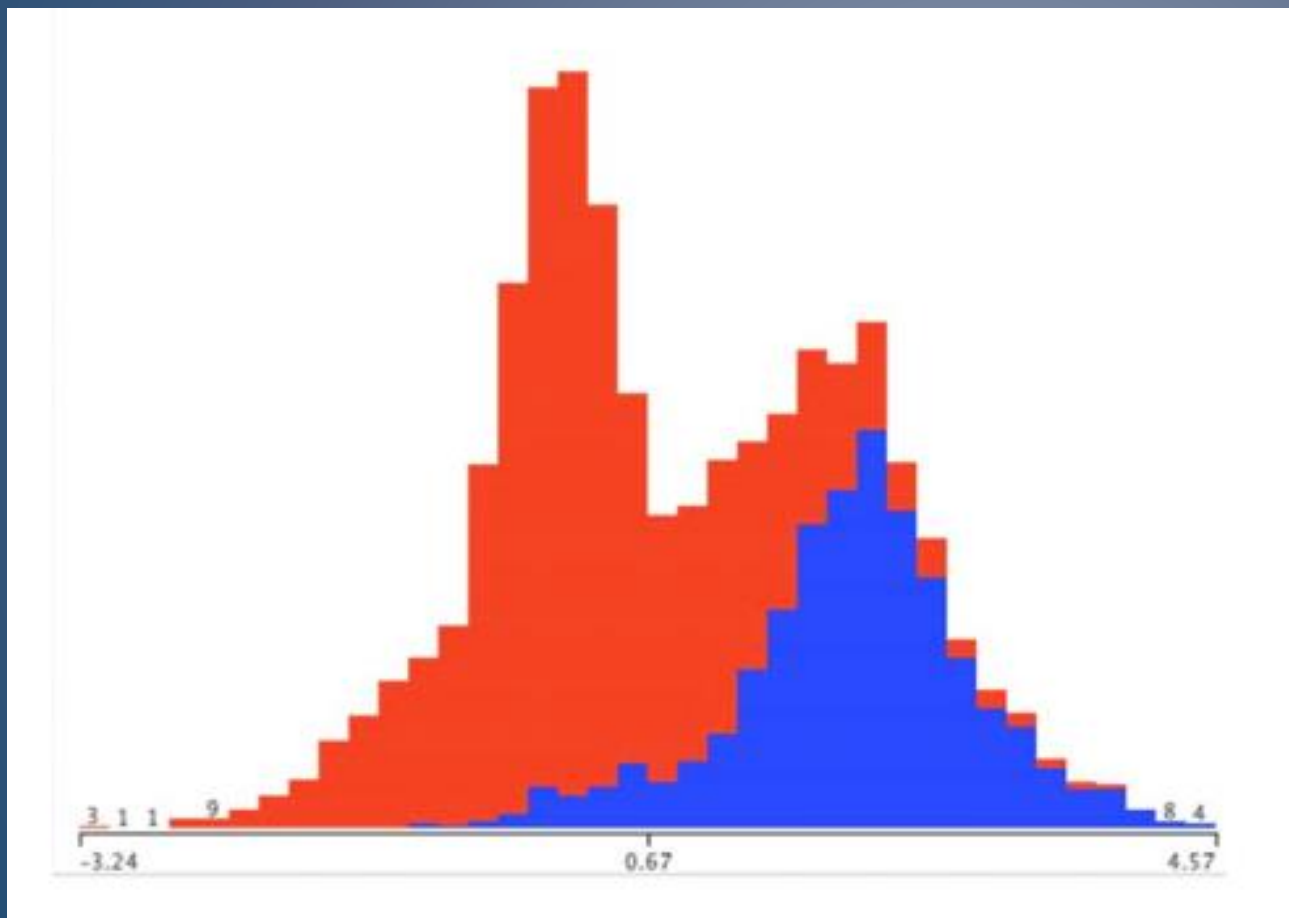
# PSIC Score

- For each column, calculate frequency of each amino acid:

$$p(a,i) = \frac{n(a,i)_{eff}}{\sum_b n(b,i)_{eff}}$$

Q5E940_BOVIN	-----	-MPREDRATWKSNYFLKIIQLDDYPKCFIVGADNVGSKOMQOIRMSLRGK-AVVLMGKNTMMRKAIRGHLENN--PALE
RLA0_HUMAN	-----	-MPREDRATWKSNYFLKIIQLDDYPKCFIVGADNVGSKOMQOIRMSLRGK-AVVLMGKNTMMRKAIRGHLENN--PALE
RLA0_MOUSE	-----	-MPREDRATWKSNYFLKIIQLDDYPKCFIVGADNVGSKOMQOIRMSLRGK-AVVLMGKNTMMRKAIRGHLENN--PALE
RLA0_RAT	-----	-MPREDRATWKSNYFLKIIQLDDYPKCFIVGADNVGSKOMQOIRMSLRGK-AVVLMGKNTMMRKAIRGHLENN--PALE
RLA0_CHICK	-----	-MPREDRATWKSNYFMKIIQLDDYPKCFVVGADNVGSKOMQOIRMSLRGK-AVVLMGKNTMMRKAIRGHLENN--PALE
RLA0_RANSY	-----	-MPREDRATWKSNYFLKIIQLDDYPKCFIVGADNVGSKOMQOIRMSLRGK-AVVLMGKNTMMRKAIRGHLENN--SALE
Q7ZUG3_BRARE	-----	-MPREDRATWKSNYFLKIIQLDDYPKCFIVGADNVGSKOMQOIRMSLRGK-AVVLMGKNTMMRKAIRGHLENN--PALE
RLA0_ICTPU	-----	-MPREDRATWKSNYFLKIIQLDDYPKCFIVGADNVGSKOMQOIRMSLRGK-AVVLMGKNTMMRKAIRGHLENN--PALE
RLA0_DROME	-----	-MVRENKAAWKAQYFIKVVLFDEFPKCFIVGADNVGSKOMONIRTSLRGL-AVVLMGKNTMMRKAIRGHLENN--POLE
RLA0_DICDI	-----	-MSGAG-SKRKKLFIKATKLFYTDKMIYAEADYVGSOLOKIRKSRGI-GAVLMGKNTMIRKVVIRDLADSK--PELD
Q54LP0_DICDI	-----	-MSGAG-SKRKNVFIKATKLFYTDKMIYAEADYVGSOLOKIRKSRGI-GAVLMGKNTMIRKVVIRDLADSK--PELD
RLA0_PLAF8	-----	-MAKLSKQKKQMYIEKLSSLIQQYSKILIVHVDNVGSKOMASVYRKS LRGK-AVVLMGKNTMIRKVVIRDLADSK--PELD
RLA0_SULAC	-----	-HIGLAVTTTKKIAKWKVDEVAELTSLKTHKTHIIIANIEGFPADKLHEIRKRLRGK-ADIKVTIKNNLFIKALKNAG--YDIX
RLA0_SULTO	----	-MRIMAVITQERKIAKWKIEEVKELQLREYHTIIIANIEGFPADKLHDIRKKMRGM-AEIKVTIKNTLFIKALKNAG--LDVS
RLA0_SULSO	----	-MKRLALALKQRKVASWKLEEVKELTELKNSNTILIGNLEGFPADKLHEIRKRLRGK-ADIKVTIKNTLFIKALKNAG--LDVS
RLA0_AERPE	MSVYSLVQ	QMYKREKPIPEWKTLMLELELFSKHVVVFLADLTGTFVYQYRVRKKLWKKYPMVAKKRIILRAMKAAGLE--LDDN
RLA0_PYRAE	-MMLAIG	KRRYVTRQYPAKVKIVSEATLLQKYPYVFLFDLHGLSRILHEVRYRLRRY-GVIKTIKPTLFIKALKNAG--IPAE
RLA0_METAC	-----	-MAEERNHTEHIPQWKDEIENIKELIQSHKVFQMVIEGILATKMKIIRDLKDY-AVLKYSRNTLTERALNQLG--ETIP
RLA0_METMA	-----	-MAEERNHTEHIPQWKDEIENIKELIQSHKVFQMVIEGILATKMKIIRDLKDY-AVLKYSRNTLTERALNQLG--ESIP
RLA0_ARCFU	-----	-MAAVRGS--PPEYKVRAVEEIKRMISKPVVAIVSFRNVPAGOMQIRREFRGK-AEIKVVKNTLLEALDALG--GDYL
RLA0_METKA	MAVKAKGQPP	SCYEKVAENKRREYKELKELMDEYENVGLVDLEGIPAPOLQEIIRAKLRERDTIIRMSRNTLMRITALEEKLDER--PELE
RLA0_METTH	-----	-MAHVAEWKKEVQELHDLIKGYEVVGIANLADIPAROLOKMRQTLRDS-ALIRMSKKTLLISLALAKAGREL--ENVD
RLA0_METTL	-----	-MITAESEHKIAPWKIEEVNKLKELLNKQIVALVDMMEVPAROLOEIRDKIR-GTMTLKMSRNTLIERAIKEVAEETGNPEFA
RLA0_METVA	-----	-MIDAKSEHKIAPWKIEEVNALKELLNKSNVIALIDMMEVPAROLOEIRDKIR-GTMTLKMSRNTLIERAIKEVAEETGNPEFA
RLA0_METJA	-----	-METKVAHVAPWKIEEVNKLKELIKSKPVVAIVDMMDYPAPOLQEIIRDKIR-DKVKLRMSRNTLIERAIKEVAEELNPKLA

# PSIC Score histogram from HumDiv



# Classification

- Naive Bayes method
- What is a Naive Bayes method/classifier?

# Naive Bayes Classifier

“Training Data”



$$\begin{aligned} &\Pr(x_1 | +), \Pr(x_1 | -), \\ &\Pr(x_2 | +), \Pr(x_2 | -), \dots \\ &\Pr(x_n | +), \Pr(x_n | -), \end{aligned}$$

Bayesian inference:

Expresses how a subjective assessment of likelihood should rationally change to account for evidence

$$\Pr(+ | x_1, x_2, \dots, x_n) \propto \Pr(x_1 | +) \Pr(x_2 | +) \dots \Pr(x_n | +) \Pr(+)$$

$$\Pr(- | x_1, x_2, \dots, x_n) \propto \Pr(x_1 | -) \Pr(x_2 | -) \dots \Pr(x_n | -) \Pr(-)$$



+ or -

# Bayesian probability

- In statistics, *frequentists* and *Bayesians* often disagree.
- A *frequentist* is a person whose long-run ambition is to be wrong 5% of the time.
- A *Bayesian* is one who, vaguely expecting a horse, and catching a glimpse of a donkey, strongly believes he has seen a mule.

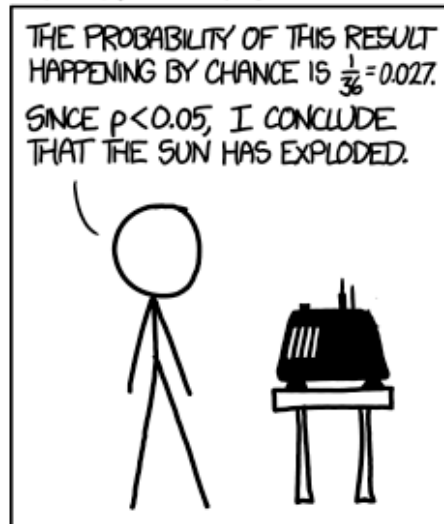


Or...

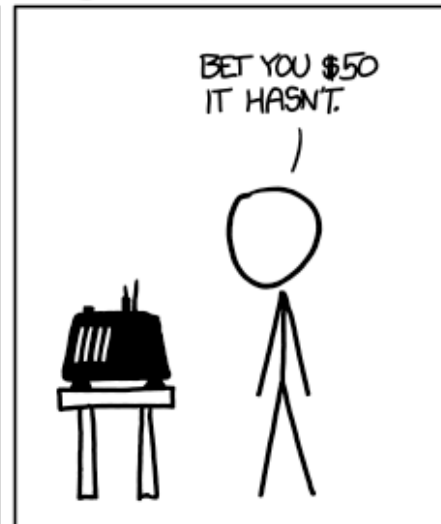
# DID THE SUN JUST EXPLODE? (IT'S NIGHT, SO WE'RE NOT SURE.)



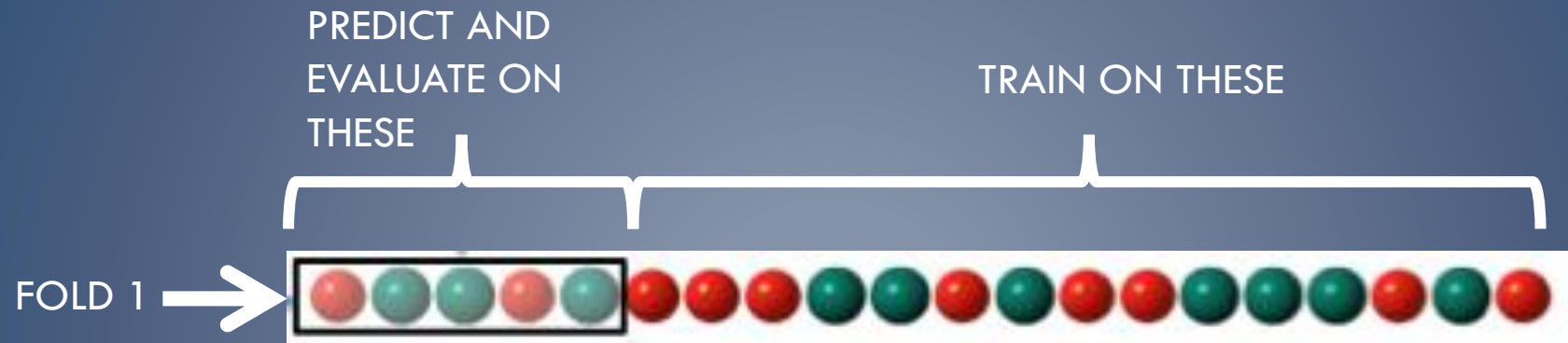
## FREQUENTIST STATISTICIAN:



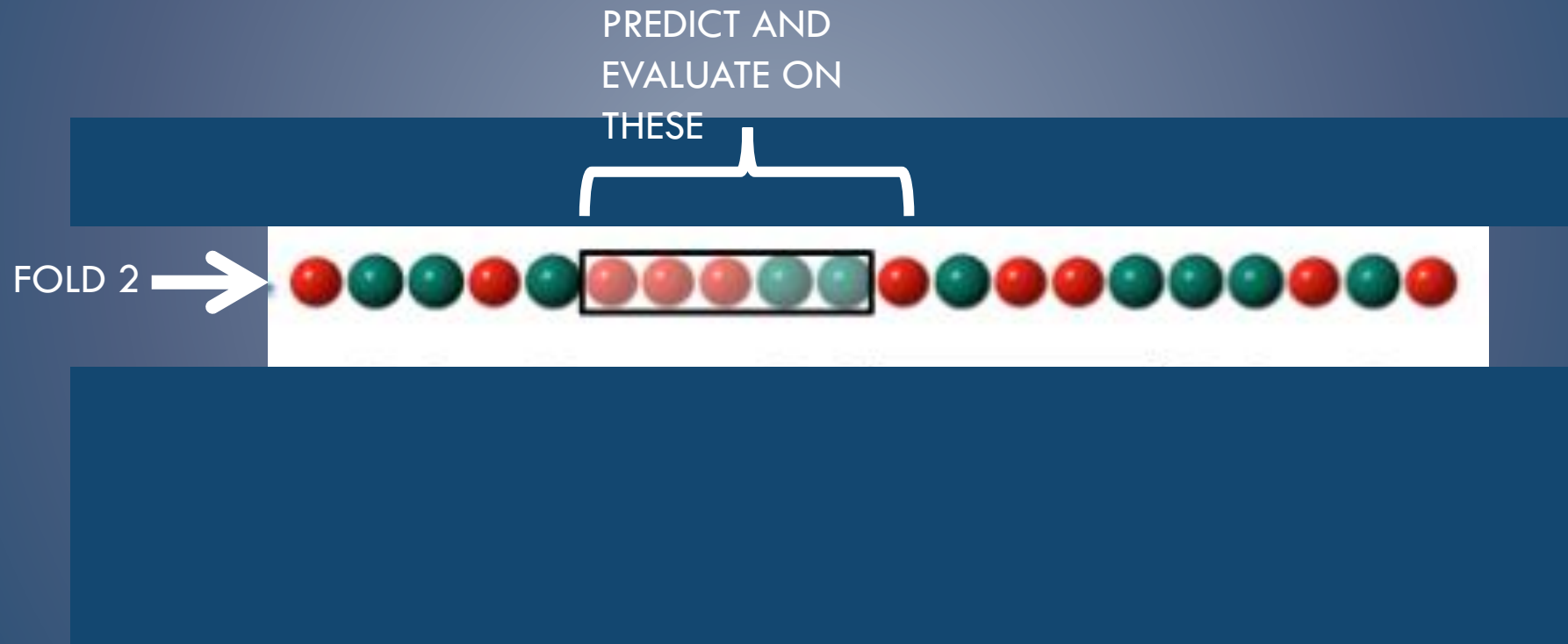
## BAYESIAN STATISTICIAN:



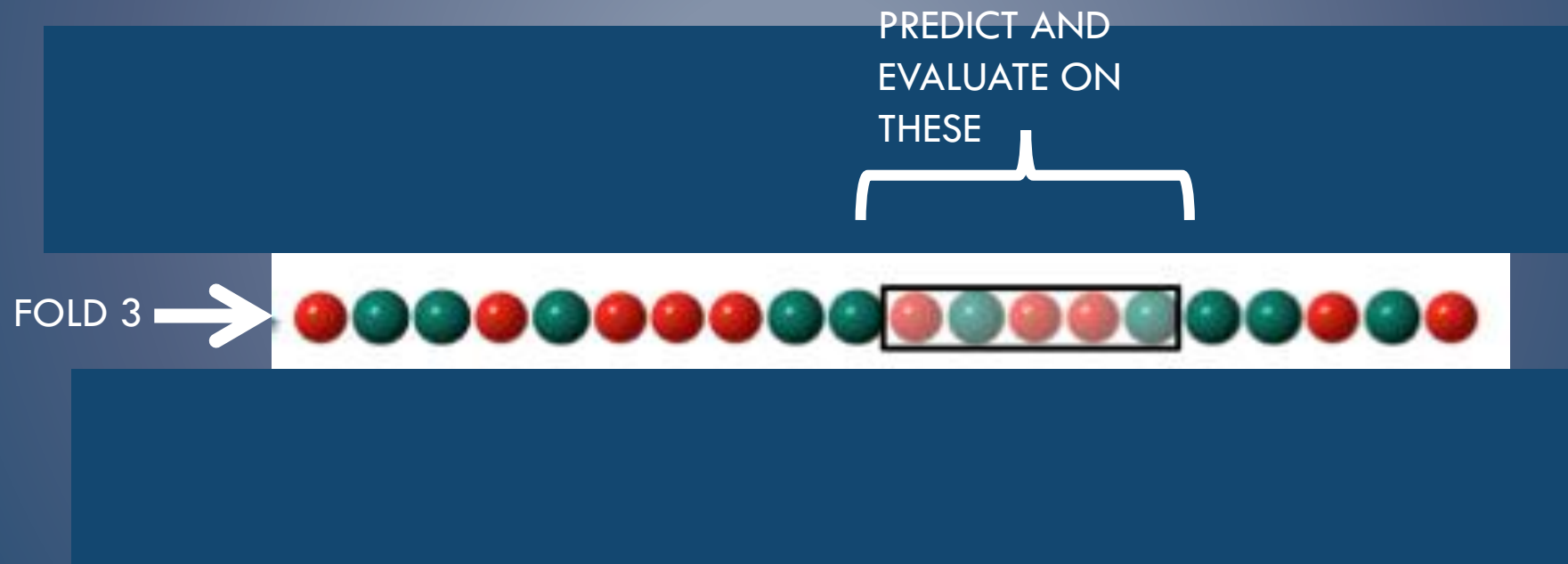
# Evaluating a classifier: Cross-validation



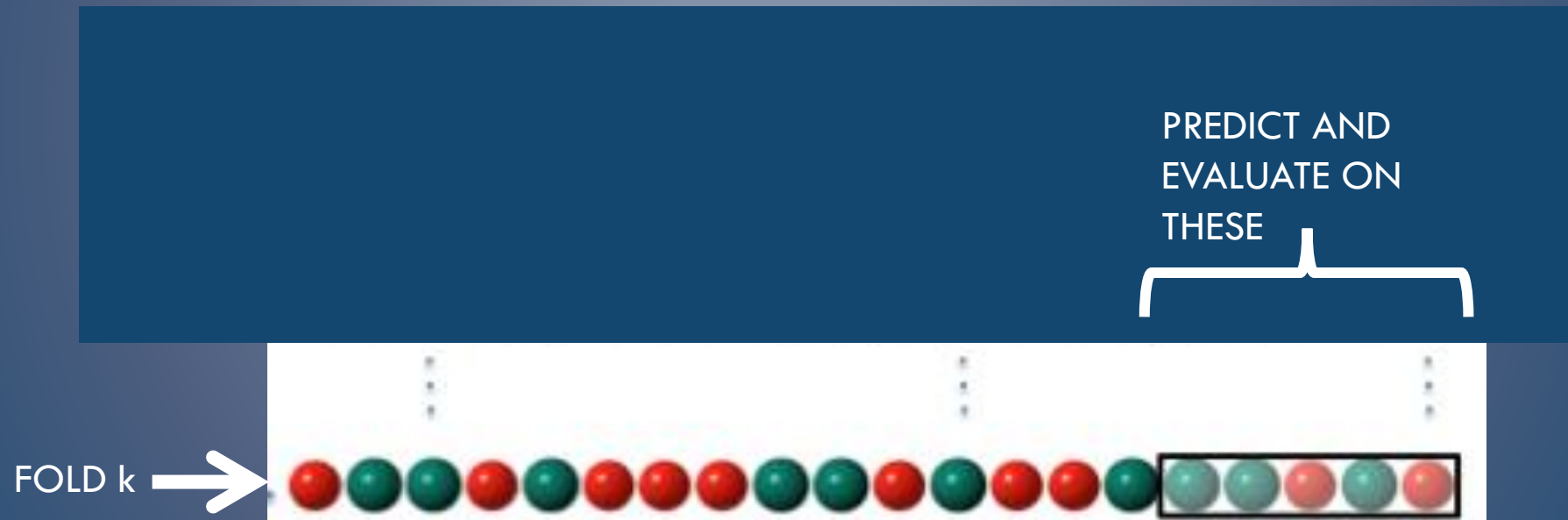
# Evaluating a classifier: Cross-validation



# Evaluating a classifier: Cross-validation

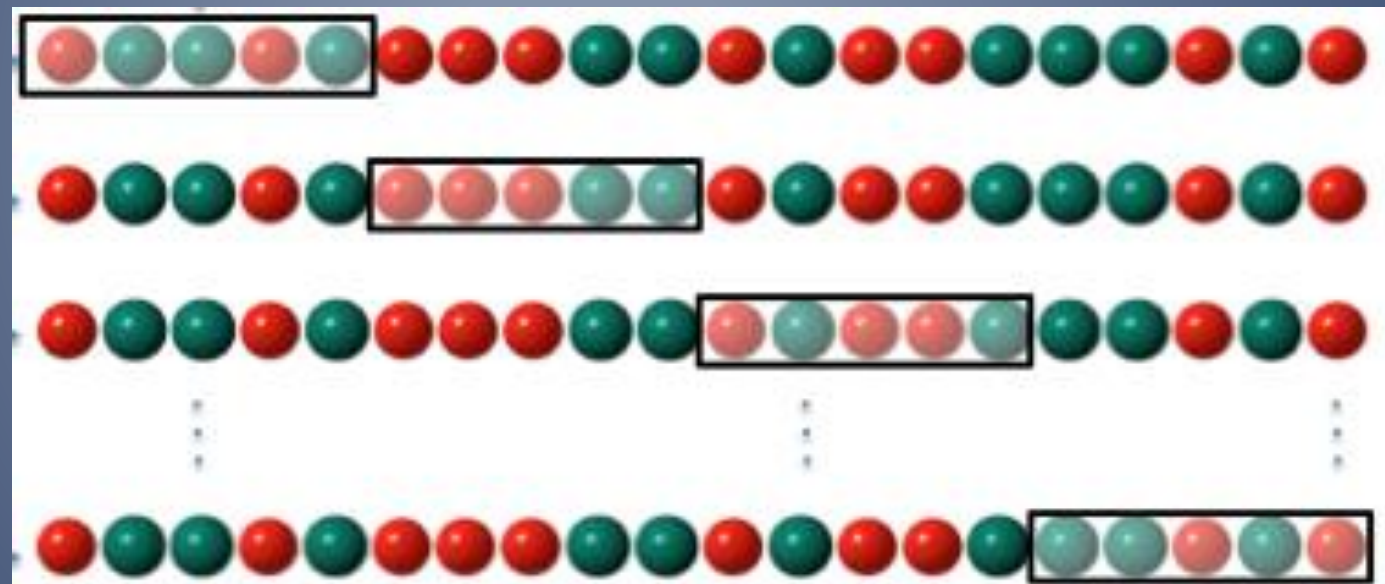


# Evaluating a classifier: Cross-validation



# Evaluating a classifier: Cross-validation

Collect all evaluation results (from k "FOLD"s)

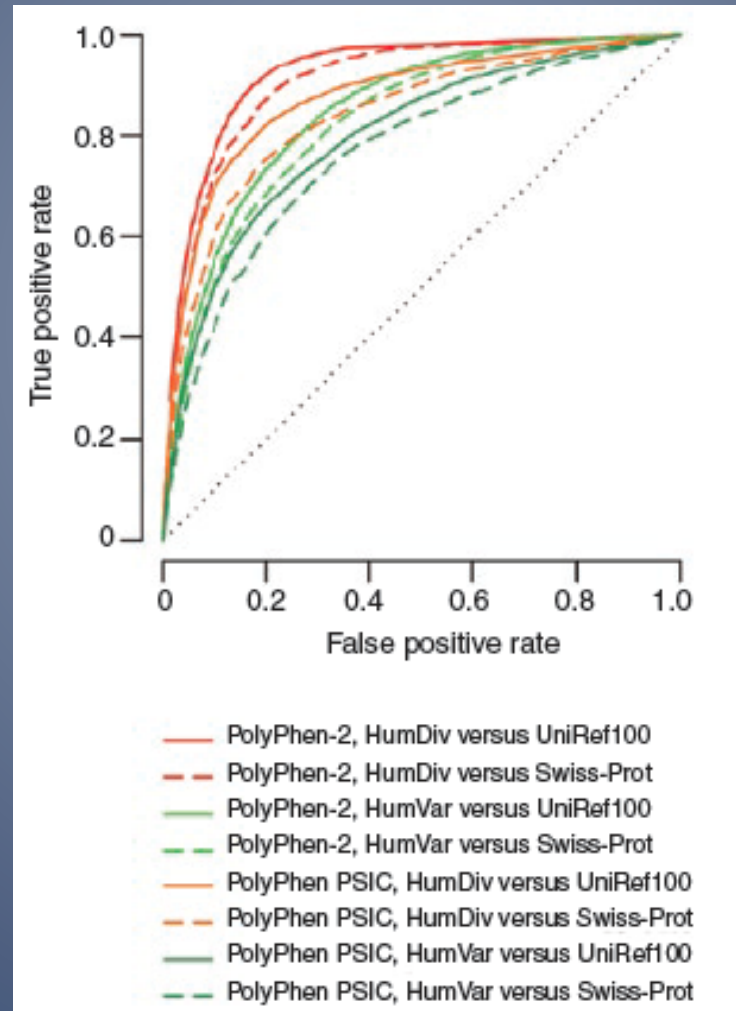


# Evaluating classification performance

		Patients with <b>bowel cancer</b> (as confirmed on <b>endoscopy</b> )		
		Condition Positive	Condition Negative	
Fecal Occult Blood Screen Test Outcome	Test Outcome Positive	<b>True Positive</b> (TP) = 20	<b>False Positive</b> (FP) = 180	<b>Positive predictive value</b> = $TP / (TP + FP)$ = $20 / (20 + 180)$ = <b>10%</b>
	Test Outcome Negative	<b>False Negative</b> (FN) = 10	<b>True Negative</b> (TN) = 1820	<b>Negative predictive value</b> = $TN / (FN + TN)$ = $1820 / (10 + 1820)$ ≈ <b>99.5%</b>
		<b>Sensitivity</b> = $TP / (TP + FN)$ = $20 / (20 + 10)$ ≈ <b>67%</b>	<b>Specificity</b> = $TN / (FP + TN)$ = $1820 / (180 + 1820)$ = <b>91%</b>	

# ROC of PolyPhen 2.0 on HumDiv

The Receiver Operating Characteristic (ROC) curve: True +ve vs False +ve





# What about SNPs outside coding regions?

- Generally hard enough to predict within coding regions – regulatory sequences notoriously hard to pin down (see ENCODE controversy)
- One interesting approach uses Support Vector Machine (SVM) classifiers to describe damage to cell-specific regulatory motif vocabularies.

## A method to predict the impact of regulatory variants from DNA sequence

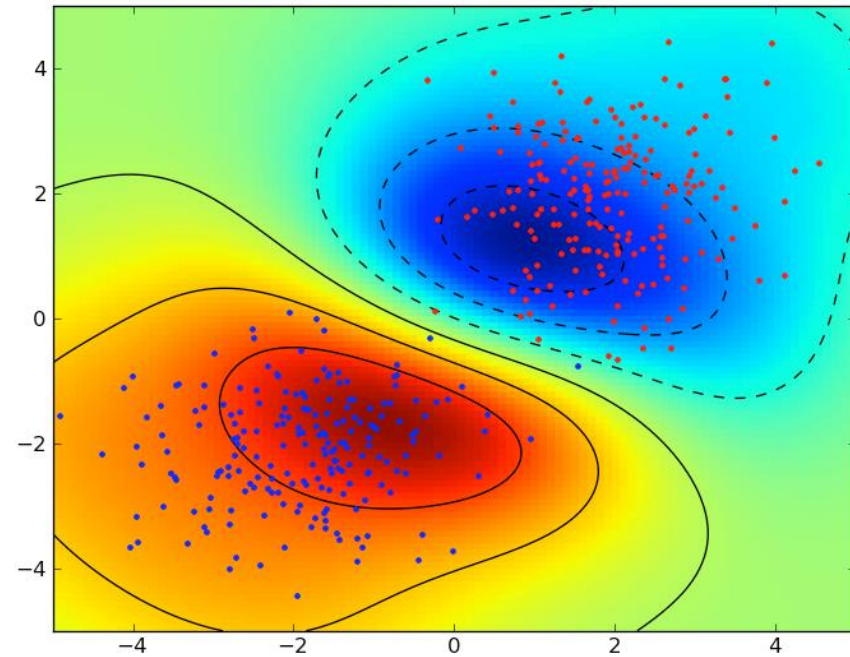
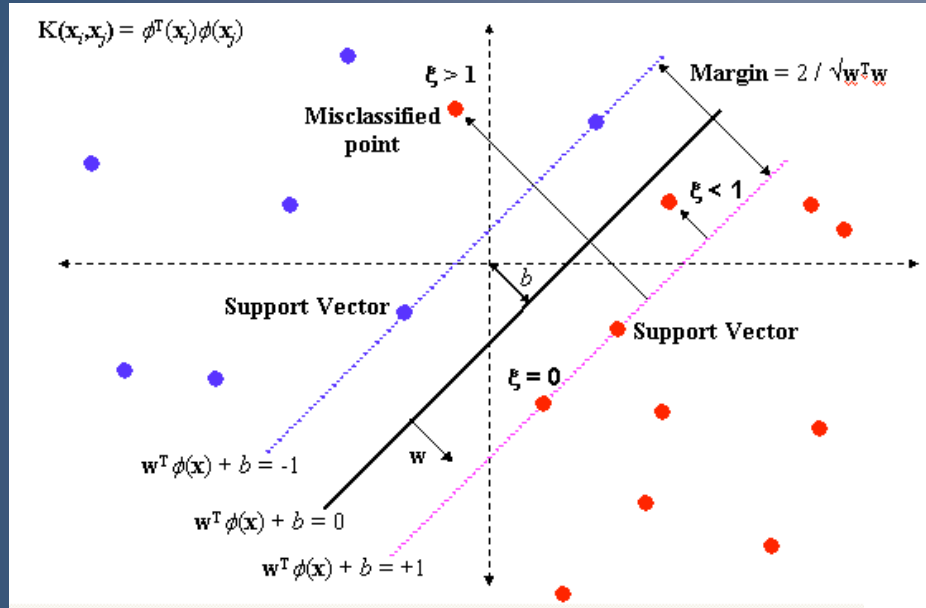
**Dongwon Lee, David U Gorkin, Maggie Baker, Benjamin J Strober, Alessandro L Asoni, Andrew S McCallion & Michael A Beer**

[Affiliations](#) | [Contributions](#) | [Corresponding authors](#)

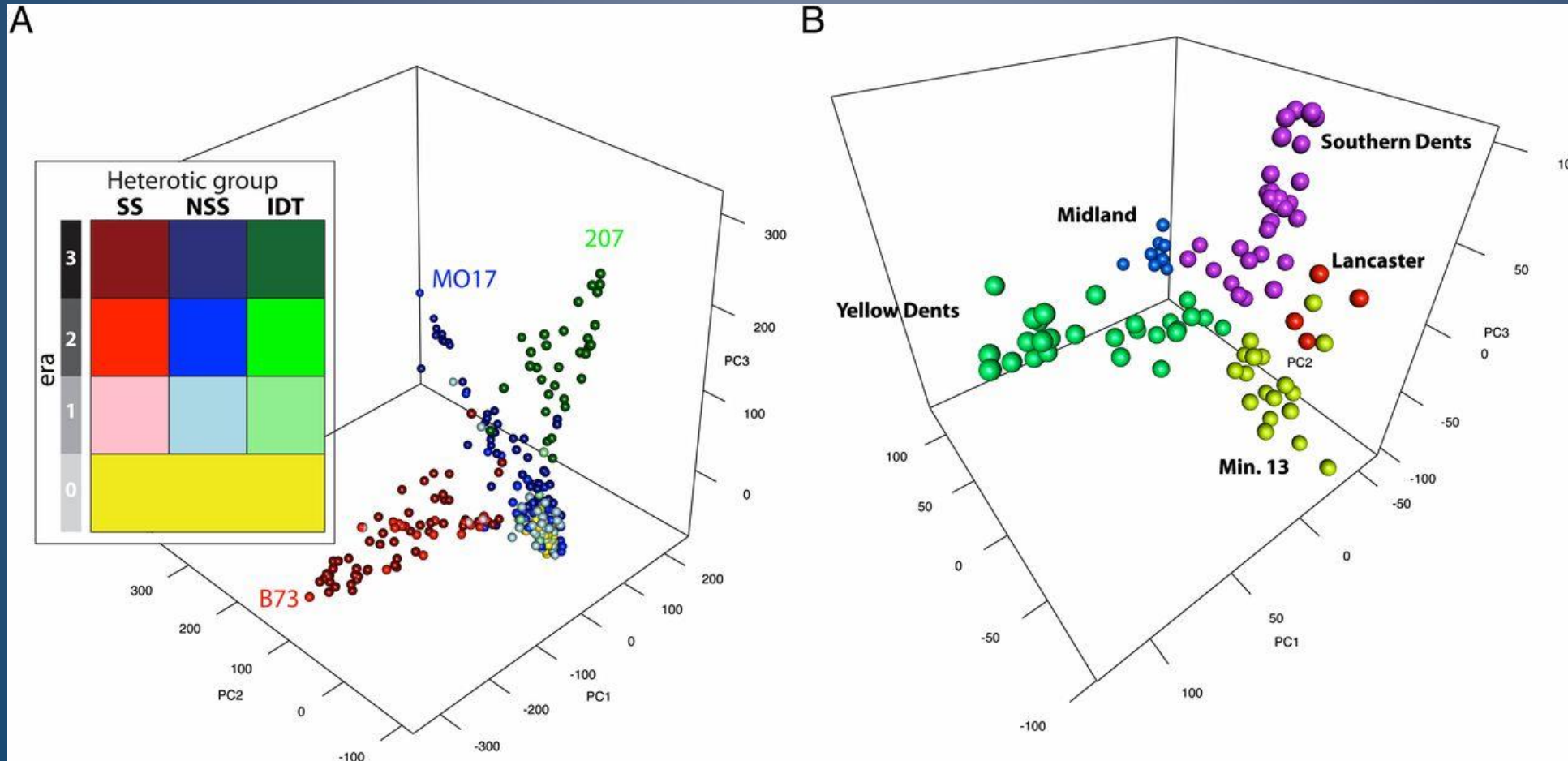
*Nature Genetics* (2015) | doi:10.1038/ng.3331

Received 12 February 2015 | Accepted 08 March 2015 | Published online 15 June 2015

# Support Vector Machines



# Populations and genetics



# Matrices of SNPs

- A typical genotype matrix looks like this

SNP ID	Line1	Line2	Line3	Line4
C1001	AA	AC	AC	CC
C1002	GG	GG	GG	GA
C1003	CC	AC	CC	CC
C1004	AA	AT	TT	TT
C1005	AG	AA	AA	GG

How do we do any kind of math on this?

# Replace with numeric data

- Use the **number of copies of the reference allele**

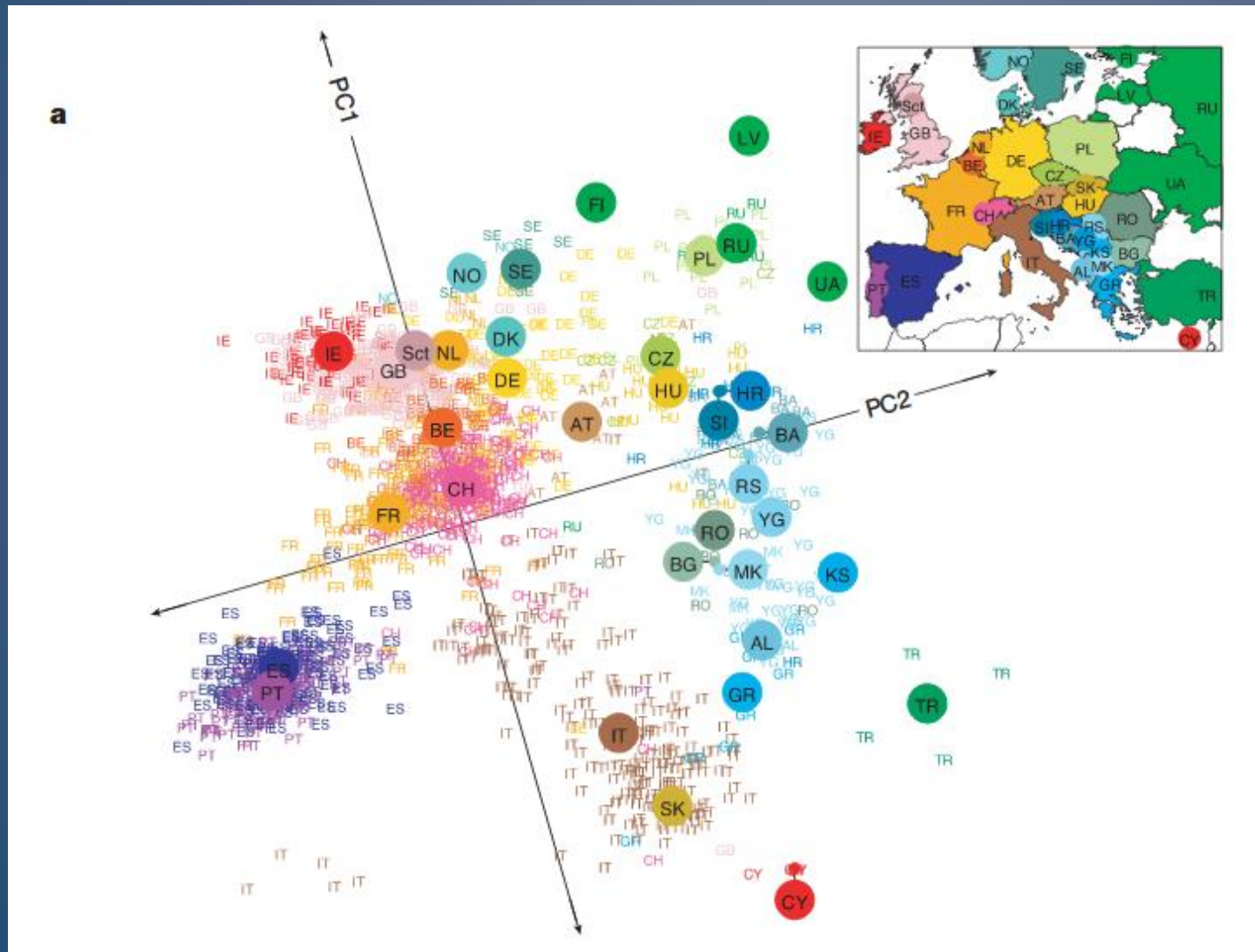
- so:

becomes:

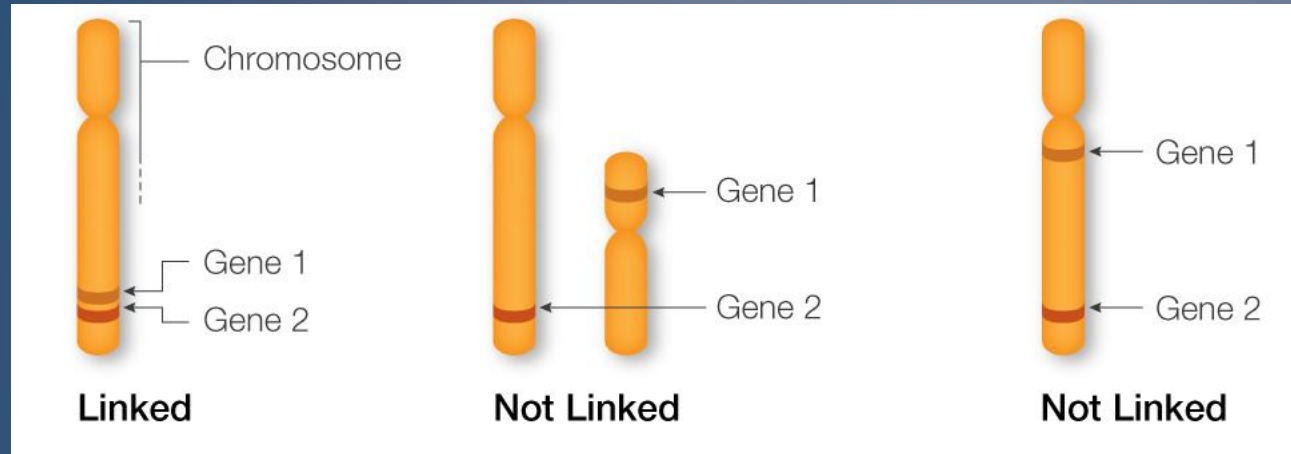
SNP ID	Line1	Line2	Line3	Line4
C1001	AA	AC	AC	CC
C1002	GG	GG	GG	GA
C1003	CC	AC	CC	CC
C1004	AA	AT	TT	TT
C1005	AG	AA	AA	GG

SNP ID	Line1	Line2	Line3	Line4
C1001	2	1	1	0
C1002	0	0	0	1
C1003	2	1	2	2
C1004	0	1	2	2
C1005	1	0	0	2

# Can then use numeric vectors for PCA

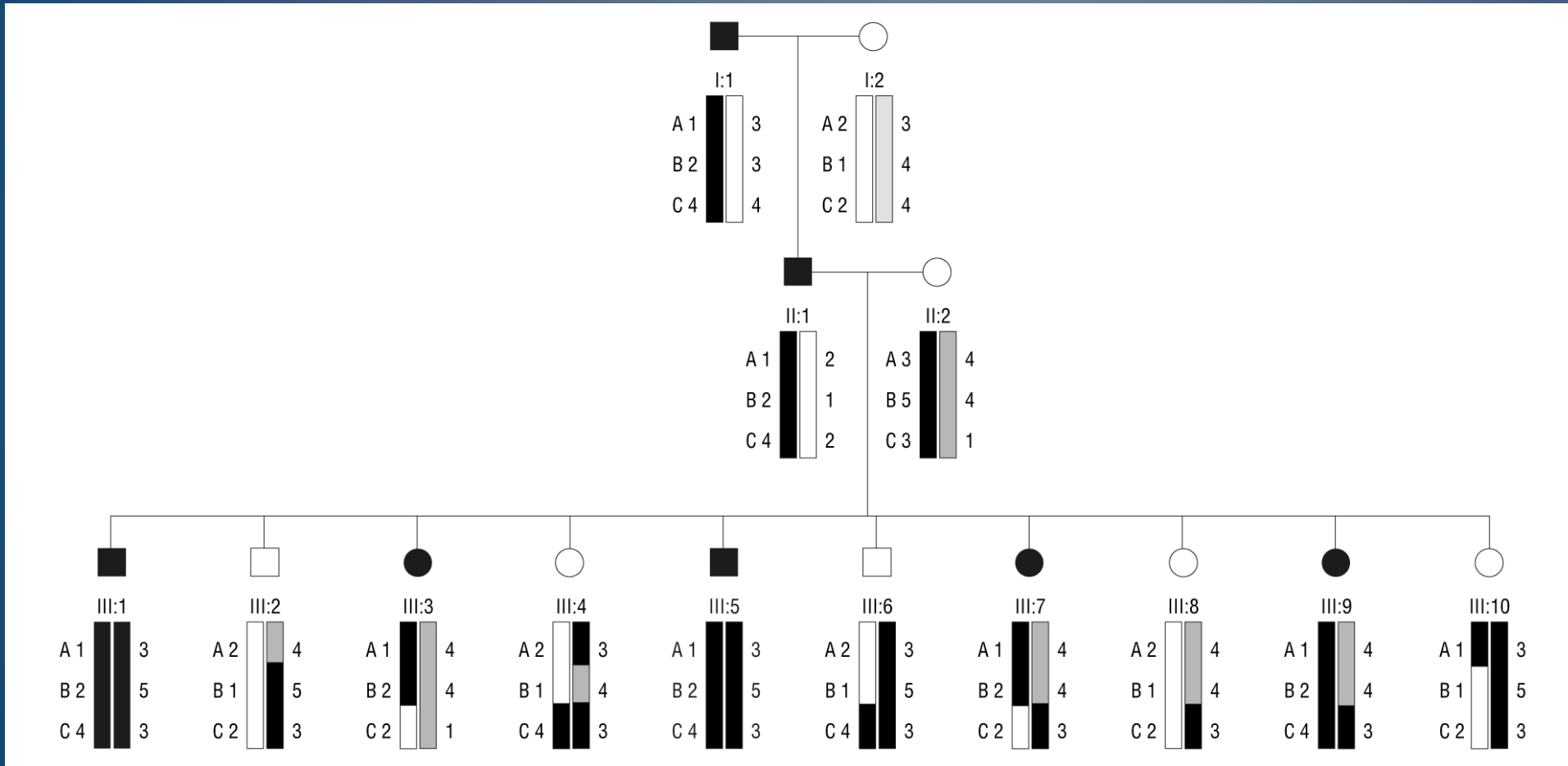


# Genetic linkage



**When a marker is correlated with a trait, it is likely to be genetically linked to the locus in a family analysis**

# Genetic linkage analysis





# Genetic linkage analysis

- Cystic Fibrosis and the CFTR gene mutations.
- “Linkage analysis”
  - Genotype members of a family (with some individuals carrying the disease)
  - Find a genetic marker that correlates with disease
  - Disease gene lies close to this marker.

# Limits of genetic linkage analysis

- Requires data from entire families, preferably large ones, where the trait is segregating – easy in plants, hard in humans
- Linkage analysis less successful with common diseases, e.g., heart disease or cancers.
- Requires single, large effect loci

# Genome-wide Association Studies (GWAS)

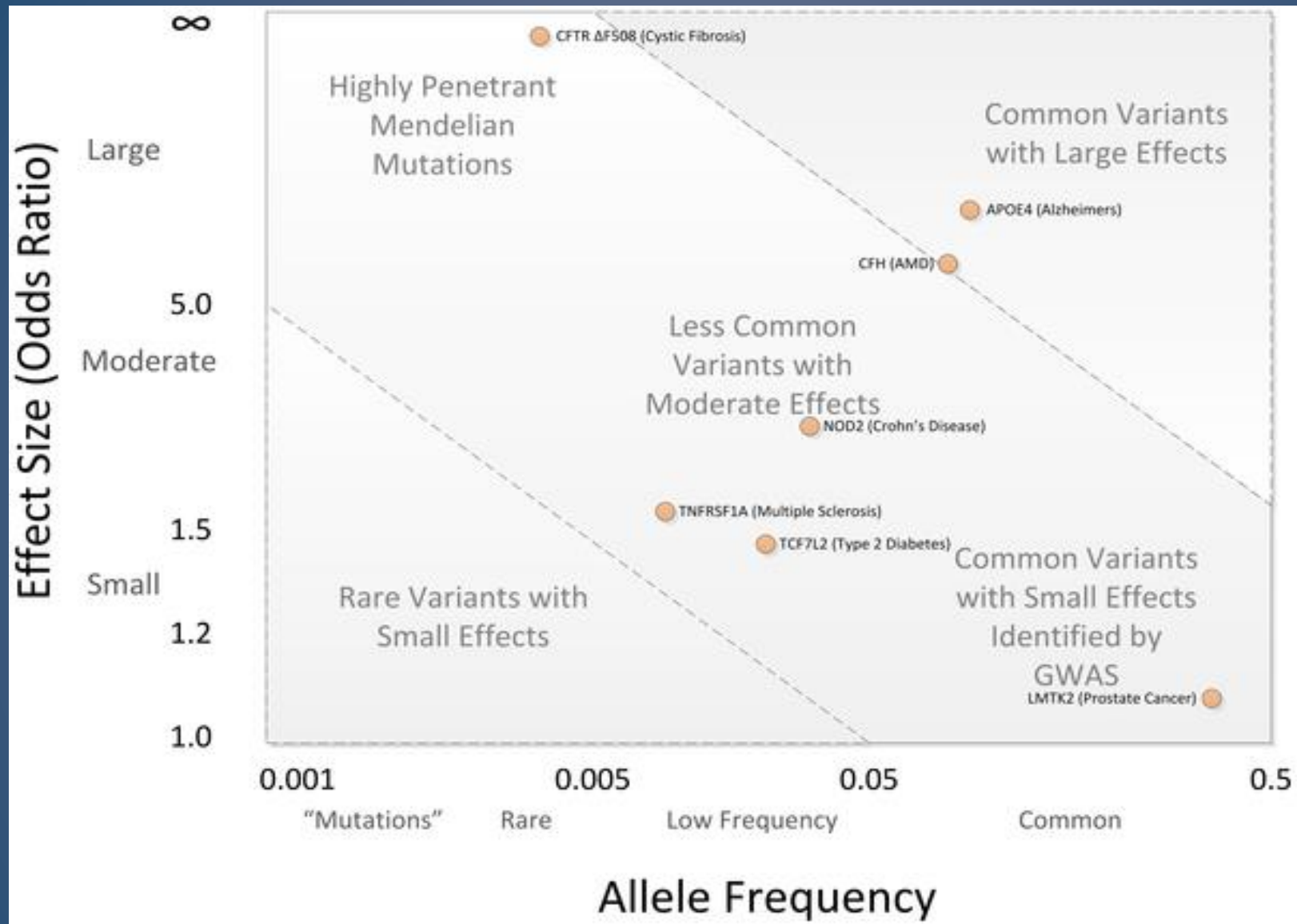
<https://doi.org/10.1371/journal.pcbi.1002828>

<https://doi.org/10.1371/journal.pcbi.1002822>

# Common disease common variant

- Hypothesis that common diseases are influenced by genetic variation that is “common” in the population
- Implications:
  - Any individual variation (SNP) will have relatively small correlation with disease
  - Multiple common alleles *together* influence the disease phenotype
- Argument for population-based studies versus family based studies. (Think about it!)

Figure 1. Spectrum of Disease Allele Effects.



Bush WS, Moore JH (2012) Chapter 11: Genome-Wide Association Studies. PLoS Comput Biol 8(12): e1002822.

doi:10.1371/journal.pcbi.1002822

<http://www.ploscompbiol.org/article/info:doi/10.1371/journal.pcbi.1002822>

# GWAS: Genotyping methodology

- Microarray technology to assay 0.5 - 1 million or more SNPs, e.g. Affymetrix and Illumina
- One population may need more SNPs to be put on the chip than another population
- Increasingly, people are using whole-genome sequencing. But LD limits utility, arrays still have advantages.

# GWAS: Phenotyping methodology

- Case/control vs. quantitative
  - Quantitative (e.g. blood pressure, LDL levels)
  - Case/control (qualitative, disease vs. no disease)
- Possible to look at more than one phenotype? Electronic medical records (EMR) for phenotyping?

# GWAS – a simple idea

## correlate genotype with phenotype

- Case/control:

		Disease?
I1:	AACGAGCTAGCGATCGATCGACTACGACTACGAGGT	-
I2:	AACGAGCTAGCGATCGATCGAC <b>A</b> ACGACTACGAGGT	+
I3:	AACGAGCTAGCGATCGATCGACTACGACTACGAGGT	-
I4:	AACGAGCTAGCGATCGATCGACTACGACTACGAGGT	-
I5:	AACGAGCTAGCGATCGATCGAC <b>A</b> ACGACTACGAGGT	+
I6:	AACGAGCTAGCGATCGATCGACTACGACTACGAGGT	-
I7:	AACGAGCTAGCGATCGATCGACTACGACTACGAGGT	-
I8:	AACGAGCTAGCGATCGATCGACTACGACTACGAGGT	-



# GWAS statistics: case vs control

- The Fisher Exact test

I1: AACGAGCTAGCGATCGATCGACTACGACTACGAGGT -  
I2: AACGAGCTAGCGATCGATCGACTACGACTACGAGGT -  
I3: AACGAGCTAGCGATCGATCGACTACGACTACGAGGT -  
I4: AACGAGCTAGCGATCGATCGACTACGACTACGAGGT -  
I5: AACGAGCTAGCGATCGATCGAC**A**ACGACTACGAGGT +  
I6: AACGAGCTAGCGATCGATCGACTACGACTACGAGGT -  
I7: AACGAGCTAGCGATCGATCGACTACGACTACGAGGT -  
I8: AACGAGCTAGCGATCGATCGAC**A**ACGACTACGAGGT +  
I9: AACGAGCTAGCGATCGATCGACTACGACTACGAGGT -  
I10: AACGAGCTAGCGATCGATCGACTACGACTACGAGGT +  
I11: AACGAGCTAGCGATCGATCGACTACGACTACGAGGT -  
I12: AACGAGCTAGCGATCGATCGACTACGACTACGAGGT -  
I13: AACGAGCTAGCGATCGATCGAC**A**ACGACTACGAGGT -  
I14: AACGAGCTAGCGATCGATCGAC**A**ACGACTACGAGGT +

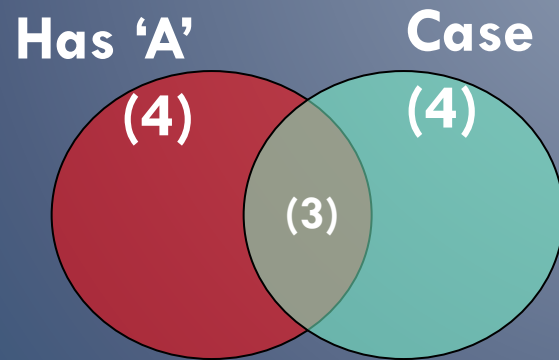
	Has 'A'	Has 'T'
Case	3	1
Control	1	9

# GWAS statistics: case vs control

- The Fisher Exact test

	Has 'A'	Has 'T'
Case	3	1
Control	1	9

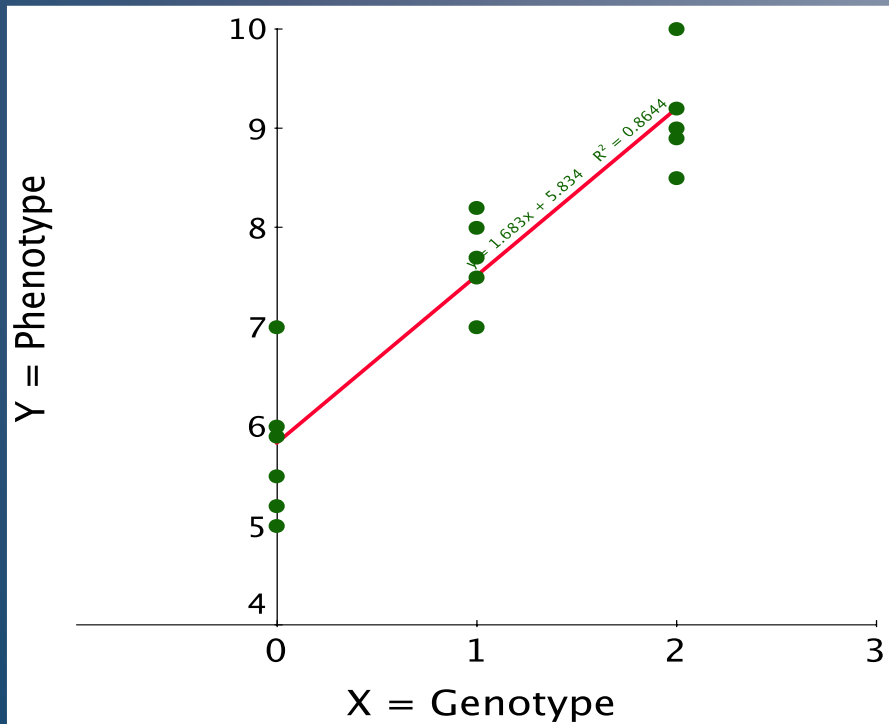
p-value < 0.05



All individuals (14)

# Quantitative phenotypes

- $Y_i$  = Phenotype value of Individual  $i$
- $X_i$  = Genotype value of Individual  $i$



Linear regression

$$Y = a + bX$$

If no association,  $b \approx 0$

The more  $b$  differs from 0, the stronger the association

# GWAS Gotchas

- Before we start on the stats, some gotchas:
  - Correlation is not causation
  - Population structure (see later)
  - Linkage disequilibrium (see later)
  - Phenotyping
- Also, even if it all works, can be hard to interpret
  - Say a SNP correlates well with heart disease
    - Could be a direct biochemical link
    - Could be behavioral (makes you like bacon...)

# GWAS statistics: case vs control

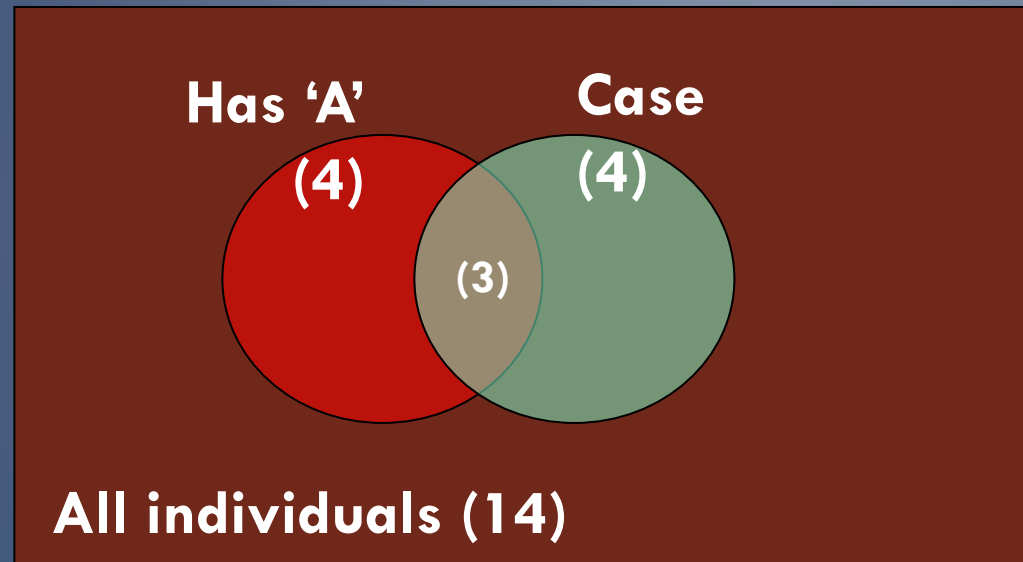
- The Fisher Exact test

I1: AACGAGCTAGCGATCGATCGACTACGACTACGAGGT -  
I2: AACGAGCTAGCGATCGATCGACTACGACTACGAGGT -  
I3: AACGAGCTAGCGATCGATCGACTACGACTACGAGGT -  
I4: AACGAGCTAGCGATCGATCGACTACGACTACGAGGT -  
I5: AACGAGCTAGCGATCGATCGAC**A**CGACTACGAGGT +  
I6: AACGAGCTAGCGATCGATCGACTACGACTACGAGGT -  
I7: AACGAGCTAGCGATCGATCGACTACGACTACGAGGT -  
I8: AACGAGCTAGCGATCGATCGAC**A**CGACTACGAGGT +  
I9: AACGAGCTAGCGATCGATCGACTACGACTACGAGGT -  
I10: AACGAGCTAGCGATCGATCGACTACGACTACGAGGT +  
I11: AACGAGCTAGCGATCGATCGACTACGACTACGAGGT -  
I12: AACGAGCTAGCGATCGATCGACTACGACTACGAGGT -  
I13: AACGAGCTAGCGATCGATCGAC**A**CGACTACGAGGT -  
I14: AACGAGCTAGCGATCGATCGAC**A**CGACTACGAGGT +

	Has 'A'	Has 'T'
Case	3	1
Control	1	9

# GWAS statistics: case vs control

- The Fisher Exact test



	Has 'A'	Has 'T'
Case	3	1
Control	1	9

p-value < 0.05

## GWAS statistics: case vs control

- Instead of the Fisher Exact test, can use the Chi Squared test.
- Do this test with EACH SNP separately. Get a p-value for each SNP.
- The smallest p-values point to the SNPs most associated with the disease

## Association tests: Allelic vs Genotypic

- What we saw was an “allelic association test”. Test if ‘A’ instead of ‘T’ at the position correlates with disease
- Genotypic association test: Each position is not one allele, it is two alleles (e.g, A & A, T & T, A & T).
- Correlate genotype at that position with phenotype of individual



# Genotypic association tests

- Various options
- Dominant model

	<b>AA or AT</b>	<b>TT</b>
<b>Case</b>	?	?
<b>Control</b>	?	?

# Genotypic association tests

- Various options
- Recessive model

	AA	AT or TT
Case	?	?
Control	?	?

# Genotypic association tests

- Various options
- 2 x 3 table

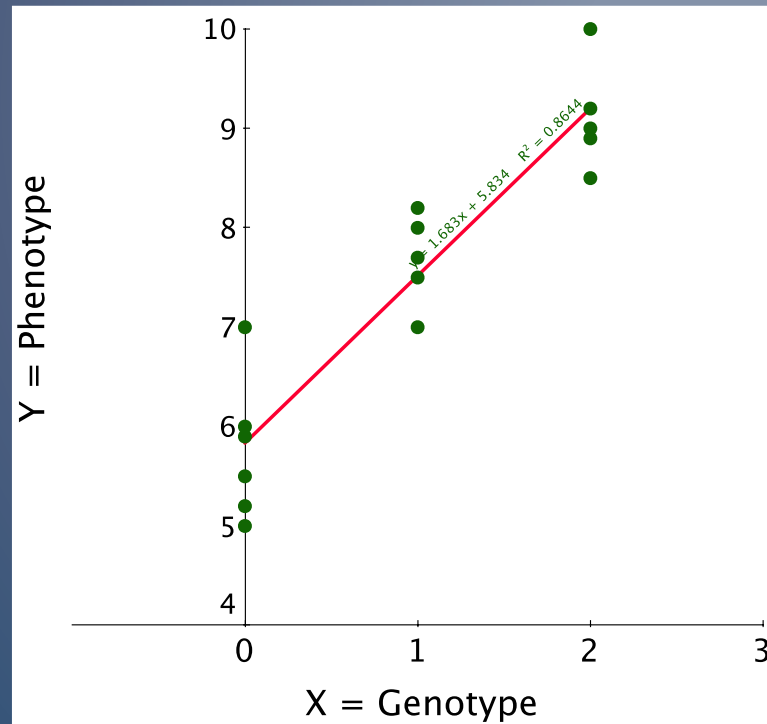
	AA	AT	TT
Case	$O_{11}$	$O_{12}$	$O_{13}$
Control	$O_{21}$	$O_{22}$	$O_{23}$

$$X^2 = \sum_i \sum_j \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

Chi-squared test

# Quantitative phenotypes

- $Y_i$  = Phenotype value of Individual  $i$
- $X_i$  = Genotype value of Individual  $i$



$$Y = a + bX$$

If no association,  $b \approx 0$

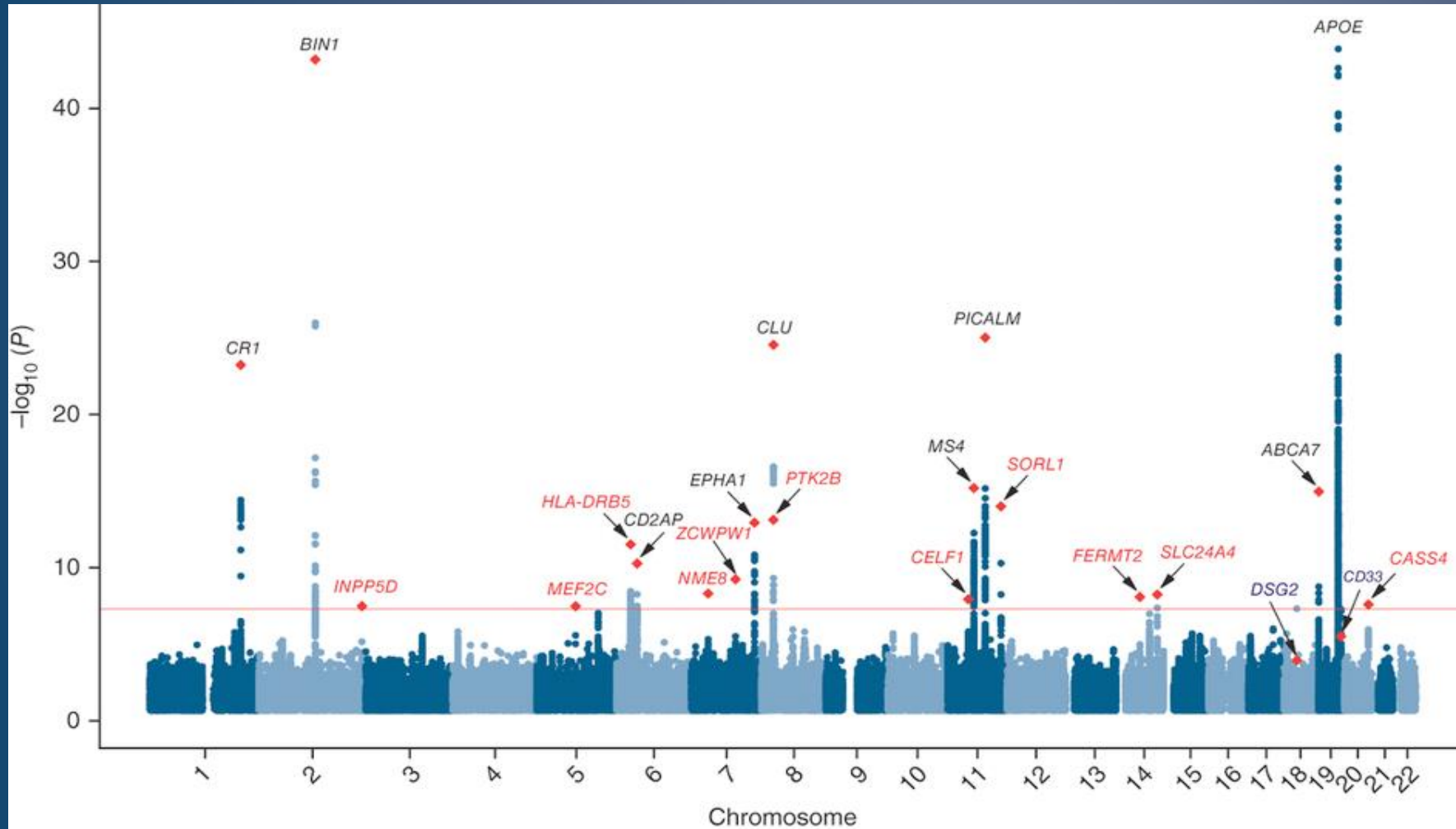
The more  $b$  differs from 0,  
the stronger the  
association

This is called “linear  
regression”

# Quantitative phenotypes

- Another statistical test commonly used on such GWAS matrices is “ANOVA” (Analysis of Variance)
- Statistical models for GWAS can get quite involved – can give refs on request.

# Manhattan plot



# Multiple hypothesis correction

- What does the “p-value of an association test = 0.01” mean ?
- It means that the observed correlation between genotype and phenotype has only 1% probability of happening just by chance. Pretty good?
- But if you repeat the test for 1 million SNPs, 1% of those tests, i.e., 10,000 SNPs will show this level of correlation, *just by chance (and by definition)*.
- <http://xkcd.com/882/>

**JELLY BEANS CAUSE ACNE!**

SCIENTISTS!  
INVESTIGATE!

BUT WE'RE PLAYING PRIME CRIB!  
... FINE.

WE FOUND NO LINK BETWEEN PURPLE JELLY BEANS AND ACNE ( $p > 0.05$ ).

WE FOUND NO LINK BETWEEN BROWN JELLY BEANS AND ACNE ( $p > 0.05$ ).

WE FOUND NO LINK BETWEEN PINK JELLY BEANS AND ACNE ( $p > 0.05$ ).

WE FOUND NO LINK BETWEEN BLUE JELLY BEANS AND ACNE ( $p > 0.05$ ).

WE FOUND NO LINK BETWEEN TEAL JELLY BEANS AND ACNE ( $p > 0.05$ ).

WE FOUND NO LINK BETWEEN SALMON JELLY BEANS AND ACNE ( $p > 0.05$ ).

WE FOUND NO LINK BETWEEN RED JELLY BEANS AND ACNE ( $p > 0.05$ ).

WE FOUND NO LINK BETWEEN TURQUOISE JELLY BEANS AND ACNE ( $p > 0.05$ ).

WE FOUND NO LINK BETWEEN MAGENTA JELLY BEANS AND ACNE ( $p > 0.05$ ).

WE FOUND NO LINK BETWEEN YELLOW JELLY BEANS AND ACNE ( $p > 0.05$ ).

WE FOUND NO LINK BETWEEN JELLY BEANS AND ACNE ( $p > 0.05$ ).

THAT SETTLES THAT.

I HEAR IT'S ONLY A CERTAIN COLOR THAT CAUSES IT.

SCIENTISTS!

BUT PRIME CRIB!

NO NO LINK BETWEEN TAN JELLY BEANS AND ACNE ( $p > 0.05$ ).

WE FOUND NO LINK BETWEEN TAN JELLY BEANS AND ACNE ( $p > 0.05$ ).

WE FOUND NO LINK BETWEEN CYAN JELLY BEANS AND ACNE ( $p > 0.05$ ).

WE FOUND A LINK BETWEEN GREEN JELLY BEANS AND ACNE ( $p < 0.05$ ).

WHOA!

WE FOUND NO LINK BETWEEN MALIVE JELLY BEANS AND ACNE ( $p > 0.05$ ).

NO NO LINK BETWEEN JELLY BEANS AND ACNE ( $p > 0.05$ ).

WE FOUND NO LINK BETWEEN LILAC JELLY BEANS AND ACNE ( $p > 0.05$ ).

WE FOUND NO LINK BETWEEN BLACK JELLY BEANS AND ACNE ( $p > 0.05$ ).

WE FOUND NO LINK BETWEEN PEACH JELLY BEANS AND ACNE ( $p > 0.05$ ).

WE FOUND NO LINK BETWEEN ORANGE JELLY BEANS AND ACNE ( $p > 0.05$ ).


**News**

**GREEN JELLY BEANS LINKED TO ACNE!**

**95% CONFIDENCE**

ONLY 5% CHANCE OF COINCIDENCE!

SCIENTISTS





# Bonferroni correction

- Multiply p-value by number of tests.
- So if the original test on a particular SNP gave a p-value of  $p$ , define the new p-value as  $p' = p \times N$ , where  $N$  is the number of SNPs tested (1 million ?)
- With  $N = 10^6$ , a p-value of  $10^{-9}$  is downgraded to  $p' = 10^{-9} \times 10^6 = 10^{-3}$ . This is quite good.

# False Discovery Rate

- Bonferroni correction will “kill” most reported associations (reduced statistical power)
- Too stringent for most applications (although good if it works). Need to balance false positive rate with false negative rate
- False Discovery Rate (FDR) is an alternative procedure to correct for multiple hypothesis testing, which is less stringent.

# False Discovery Rate

- Given a threshold  $\alpha$  (e.g., 0.05):
- Sort all p-values (N of them) in ascending order:
- $p_1 \leq p_2 \leq \dots \leq p_N$
- Count for each group of N, p from 1 to i:

$$p'_i = p_i \cdot \frac{N}{i}$$

- Require  $p' < \alpha$
- This ensures that the expected proportion of false positives in the reported associations is  $< \alpha$

# Beyond single locus associations?

- We tested each SNP separately
- Recall that our “common disease, common variant” hypothesis meant each individual SNP carries only a small effect.
- Maybe two SNPs together will correlate better with phenotype.
- So, methods for 2-locus association study.
- Main problem: Number of pairs  $\sim N^2$

# Beyond the probed SNPs?

- The SNP-chip has a large number of probes (e.g., 0.5 – 1 Million). But still, way fewer SNPs than WGS.
- But there are many more sites in the human genome where variation may exist. Are we going to miss any causal variant outside the panel of ~1 Million?
- Not necessarily.

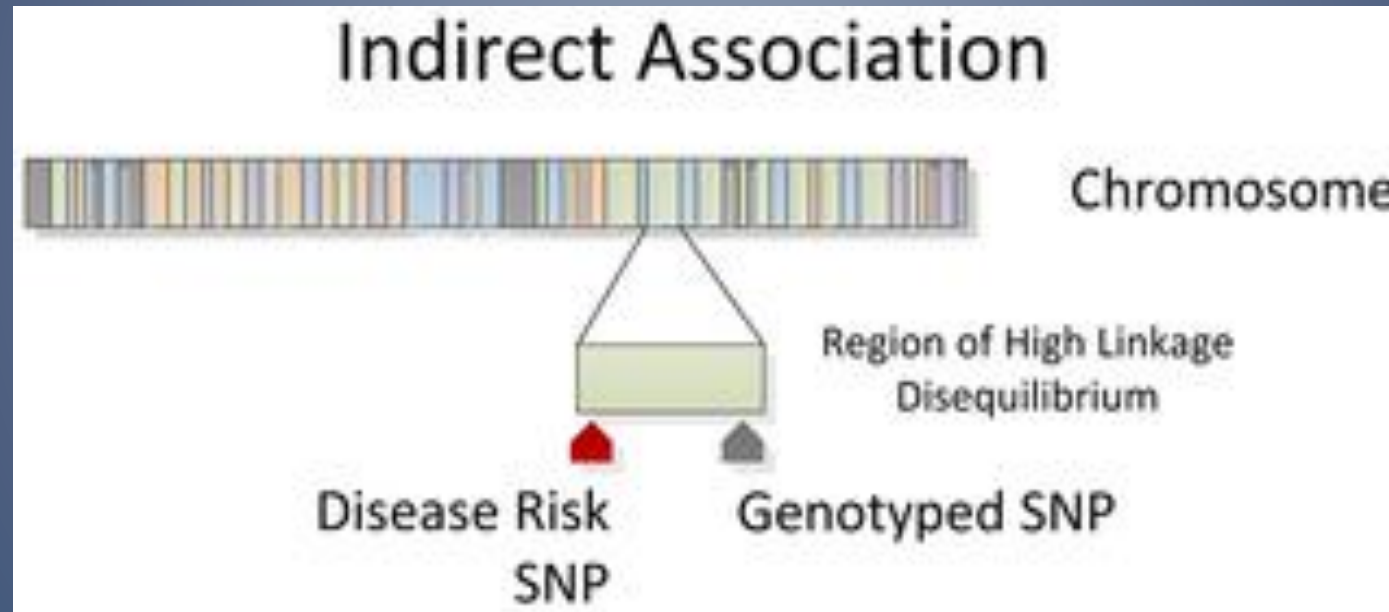
# Linkage disequilibrium

- Two sites close to each other may vary in a highly correlated manner. This is linkage disequilibrium (LD).
- Not enough recombination events have happened to make the inheritance of those two sites independent.
- If two sites are in a segment of high LD, then one site may serve as a “proxy” for the other.

# LD and its impact on GWAS

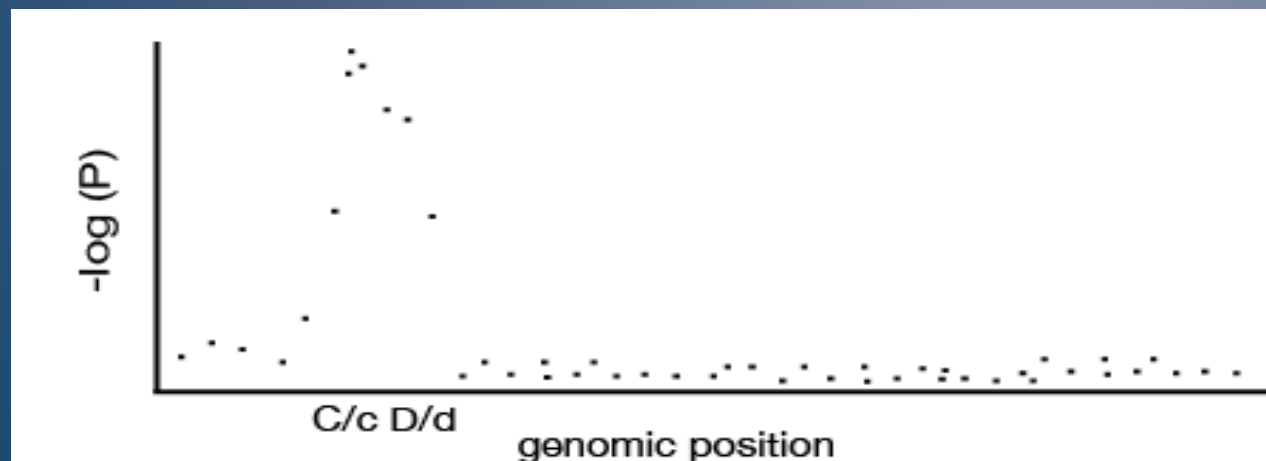
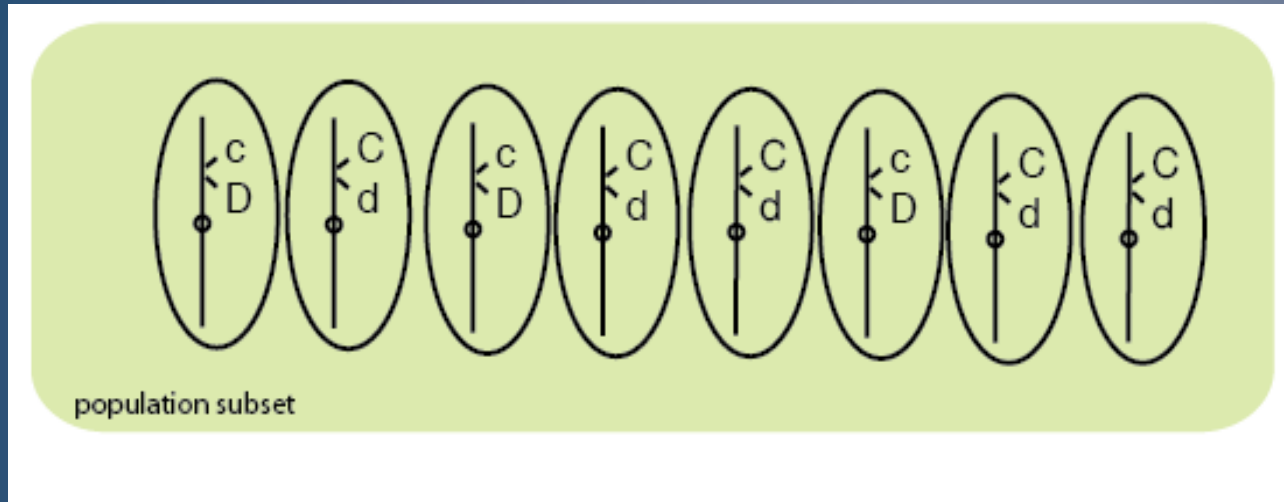
- If sites X & Y are in high LD, and X is on the SNP-chip, knowing the allelic form at X is highly informative of the allelic form at Y.
- So, a panel of 0.5 – 1 Million SNPs may represent a larger number, perhaps all of the common SNPs.
- But this also means: if X is found to have a high correlation with disease, the causal variant may be Y, and not X

# LD and its impact on GWAS

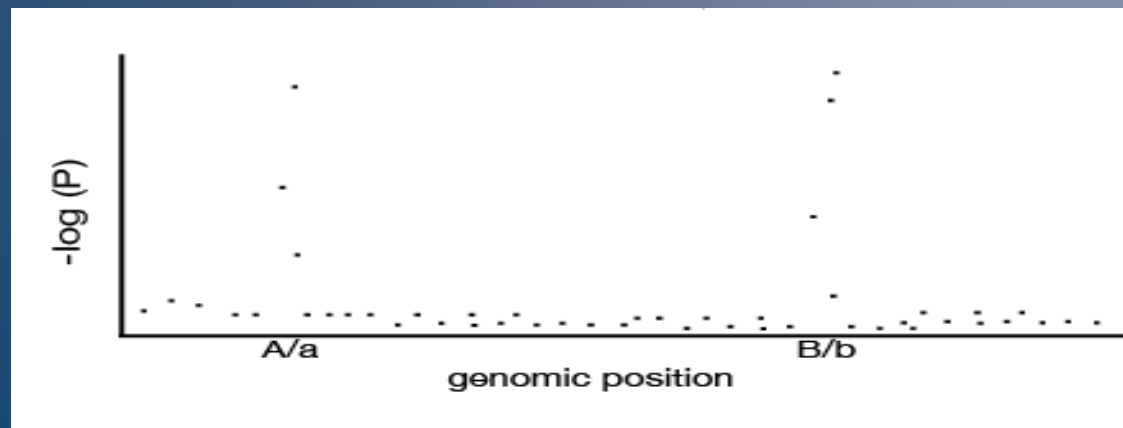
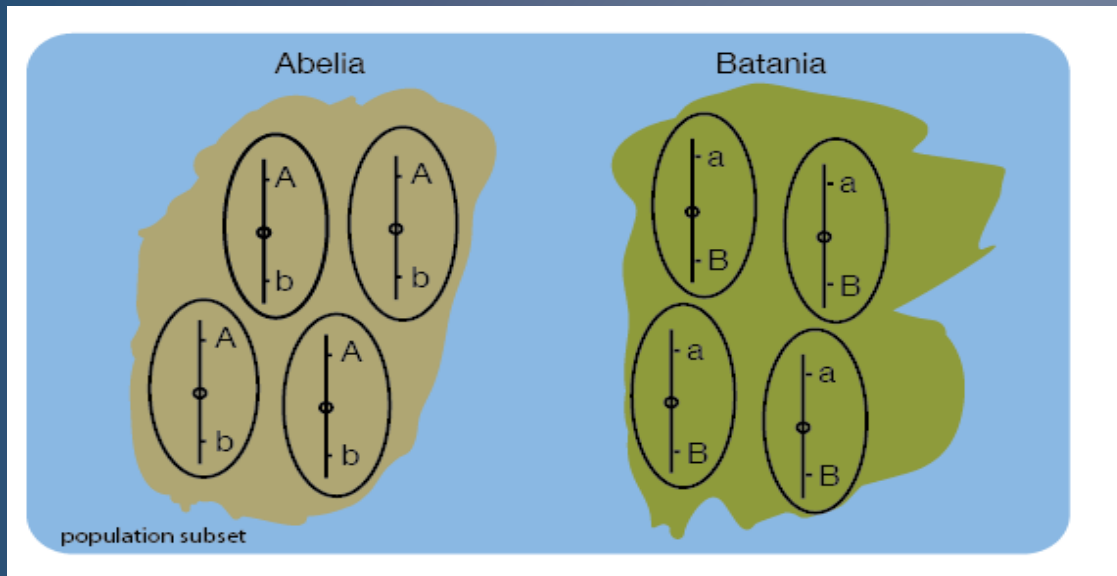




# LD impact



# Population structure



# Discussions

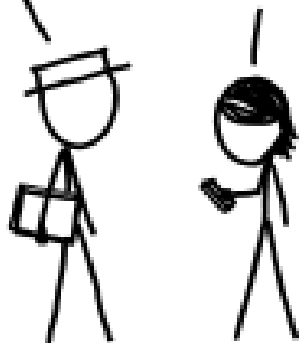
- In many cases, able to find SNPs that have significant association with disease. Risk factors, some mechanistic insights.
- GWAS Catalog : <http://www.genome.gov/26525384>
- Yet, final predictive power (ability to predict disease from genotype) is limited for complex diseases.
- “Finding the Missing Heritability of Complex Diseases”  
<http://www.genome.gov/27534229>

# Discussions

- Increasingly, whole-exome and even whole-genome sequencing used for variant detection
- Taking on the non-coding variants. Use functional genomics data as template
- Network-based analysis rather than single-site or site-pairs analysis
- Complement GWAS with family-based studies

BIOLOGY IS LARGELY SOLVED.  
DNA IS THE SOURCE CODE  
FOR OUR BODIES. NOW THAT  
GENE SEQUENCING IS EASY,  
WE JUST HAVE TO READ IT.

IT'S NOT JUST "SOURCE  
CODE." THERE'S A TON  
OF FEEDBACK AND  
EXTERNAL PROCESSING.



BUT EVEN IF IT WERE, DNA IS THE  
RESULT OF THE MOST AGGRESSIVE  
OPTIMIZATION PROCESS IN THE  
UNIVERSE, RUNNING IN PARALLEL  
AT EVERY LEVEL, IN EVERY LIVING  
THING, FOR FOUR BILLION YEARS.

IT'S STILL JUST CODE.



OK, TRY OPENING GOOGLE.COM  
AND CLICKING "VIEW SOURCE."

OK, I-... OH MY GOD.

THAT'S JUST A FEW YEARS OF  
OPTIMIZATION BY GOOGLE DEVS.  
DNA IS THOUSANDS OF TIMES  
LONGER AND WAY, WAY WORSE.

WOW, BIOLOGY  
IS IMPOSSIBLE.



# Polymorphism and Variant Analysis

Matt Hudson

University of Illinois