

Abdelrahman Adel  
17012296

## Sheet 2

### Lecture 1

①. Chromosomes have genes inside of them, chromosomes are made out of genes

• Humans have 23 pairs (46) chromosomes in all cells except reproductive cells where they have only 23 chromosomes

②. Genome is the whole DNA which is composed of multiple genes & Gene is only a small section of the Genome  
• and approximately there are around 20,000 Genes

③ Central Dogma says

DNA  $\xrightarrow{\text{Transcription}}$  RNA  $\xrightarrow{\text{Translation}}$  Protein

④ Before bioinformatics, processing data took a long time and a lot of effort but after it, it was more efficient to process huge amounts of data and even further efficiency is obtained using parallel processing

⑤ Late 90's, early 2000's due to the human genome project

⑥ The Human Genome Project was the first large-scale international effort to map and sequence the entire human genome. And it was the most significant project ever since the start of research of the human genome.

⑦ Phylogeny: it's the arrangement of species or sequences into groups based on distance and visualizing it is mainly using Phylogenetic trees & Heat maps  
more visual

1



- ③ Pairwise: Only aligning two sequences to achieve the optimal pairing  
Multiple: Aligning three or more sequences to achieve the optimal pairing

- ④ Scoring matrix: two dimensional matrix where the two sequences to be aligned are written along the axes

Scoring scheme: gives match, mismatch, and gap scores that can be used in the scoring matrix

Traceback: after finishing the scoring matrix you trace back to create the sequences.

- ⑤ Inputs: 2 sequences, scoring scheme

Outputs: 2 aligned sequences

Methods: global alignment & local alignment.

- ⑥ Genomic Variation is when at least 2 sequences have a difference between them. A SNP is single Nucleotide Polymorphism where it only looks at a singular base, and an allele has 2 types, Major allele and minor allele and they are the two forms of SNPs

- ⑦ Major alleles, and minor alleles, we identify them statistically if the main SNP is majority a specified base then it's a major allele else it's a minor allele

- ⑧ Whole Genome Sequence: sequencing ~300,000,000 bases  
Whole Exome Sequence: sequencing ~60,000,000 bases  
Large Scale Genotyping: sequencing ~1,000,000 bases
- ↓ less expensive  
↑ more expensive

- ⑨ They indicate the start and end of the coding region and without them they can't know the genes of coding so they are not junk DNA

- ⑩ BLAST breaks the query into "words" of fixed length and searches for exact matches while FASTA identifies short runs of identical residues as starting points.



### Lecture 3

- ① GWAS refers to the study in which hundreds of thousands of SNPs are genotyped across the genome and tested for association with the phenotype of interest

Input: Thousands of SNPs

Output: Association with phenotype

Methods: using the p-value & null hypothesis

- ② Phenotypes are like features, for example, phenotype of diabetes if you have it or not and these phenotypes can be influenced by genomic expressions, so much so that only one SNP can make it happen.

- ③ Manhattan plots are mainly used for GWAS, and it represents the p-values of each chromosome, and sets a threshold that anything above it is rejected by the null hypothesis

- ④ We use  $-\log p\text{-value}$  to make it easier to start from zero and increase when the p-value decreases. and the coding regions of chromosomes in the x-axis

- ⑤ Minor Allele Frequency: if a minor allele frequency under a specific percentage the subject is excluded from the Frequency output

Genotyping rate: if a SNPs with a genotyping rate missing more than 2% exclude them

Hardy-Weinberg Equilibrium: if SNP fail HWE then exclude

- ⑥ If the missing Rate of a genotype is low (less than 95% are not missing) (more than 2-5% missing) then the study of that SNP would only give incorrect results.

- ⑦ PED → contains: family ID, individual ID, paternal ID, maternal ID, sex, phenotype & genotype

MAP → contains: chromosome, marker ID, genetic distance, & physical position

BIM → contains: chromosome, variant ID, genetic position, & physical position

BED → Binary file requires .fam & .bim

- ④ FAM → contains: family ID, individual ID, paternal ID, maternal ID, sex & phenotype

- ⑧ False positive can occur as typically hundreds of thousands to millions of SNPs are tested simultaneously for association with the phenotype  
→ solution: multiple hypothesis testing is performed & confounders presence.

## ⑨ PLINK

Inputs: PED/MAP or BED/BIM/FAM files

Outputs: Frequency file & bed files

Filters: --geno, --hwe, --maf

## Lecture 4:

- ① Population structure → Systematic ancestry differences like geographic proximity or individuals sharing the same ethnicity
- ② False Positive → when it's a positive value but it's actually not positive  
False Negative → " " a negative " " " " not negative  
→ Happens when the statistical signal of the marker is not strong enough because of: a) too little evidence or b) the threshold is too low.
- ③ PCA creates principle components of the data and uses the first 2 or 3 PCAs (the heaviest and most important PCAs) as XYZ to make it possible to visualize population structure.
- ④ We can use box plots to show outliers and inliers or we can use logistical & linear Regression
- ⑤ Additive: for each minor allele you increase the value by 1, Range 0, 1, 2, 3  
Dominant: if the SNP have a minor allele then it's 1 and 0 if anything else  
Recessive: if the SNP is a homozygous minor allele then 1 else use it's 0
- ⑥ Two sample tests → creating a contingency table using case/control study without using any confounders.

(7) using the contingency table can be used to test of association using discrete test statistic, such as Fisher's exact test &  $\chi^2$  test or odds ratio

(8) All of them aim to find the p-value that would reject the null hypothesis with different equations

• Fisher's test: 
$$p = \frac{(a+b)!(c+d)!(a+c)!(b+d)!}{a!b!c!d!n!}$$

•  $\chi^2$  test:  $\chi^2 = \sum \frac{(O-E)^2}{E} \Rightarrow$  then get p from  $\chi^2$  chart

• Odds Ratio =  $\frac{\text{Odds}(A)}{\text{Odds}(B)}$ ,  $\text{Odds}(A) = \frac{P(\text{trait} | A)}{1 - P(\text{trait} | A)}$

(9) Linear models are like  $y = \beta_0 + \beta_1 x_G + \beta_2 x_C + \epsilon$

$\downarrow$  phenotype                       $\downarrow$  Genotype                       $\downarrow$  Covariates                       $\downarrow$  residuals

(10) Logistical models are used when the phenotype is binary not on a scale

(11) with LD we can figure out the correct SNPs if there was any missing SNPs.

### Lecture 5

(1) It's a way to hierarchical clustering using Distance means. Input is a list of sequences.

- (1) Compute distance between all sequences
- (2) identify the least distance
- (3) Make them on the same level
- (4) Repeat until grouping all sequences

We utilize Euclidean distance and outputs Phylogeny tree



## Lecture 6:-

- ① RNA polymerase  $\Rightarrow$  Creates an mRNA from DNA  
Reverse transcriptase  $\Rightarrow$  creates single strand DNA from mRNA  
(RT-PCR)  $\Rightarrow$  from mRNA we create cDNA
- ② cDNA represents the gene of interest DNA, not the whole DNA  
and it has a non-coding region to allow coding region to stick  
to the cDNA using the same coding region
- ③ Microarray assay:- are arrays that has many genes of the organism with  
single strand DNA probes tethered to a solid support  
when the cDNA is put in microarray it glows a specific  
color  
the color corresponds to the presence of the gene or not.
- ④ T-test: how significant is the difference between gene expression in two  
classes  
 $\rightarrow$  returns p-value  
and then we can use multiple testing correction of the p-value
- ⑤ Clustering  $\rightarrow$  unsupervised learning outputs a number of clusters given  
Classification  $\rightarrow$  supervised learning that outputs a best fit model to predict  
other new values
- ⑥ Regular K-Means is a special case of Fuzzy K-Means when  
$$P(\text{label } k | x_i, \mu_k) = \begin{cases} 1 & \text{if } x_i \text{ is closest to } \mu_k \\ 0 & \text{otherwise} \end{cases}$$
- ⑦ WGCNA can be used for finding clusters of highly correlated genes,  
for summarizing such clusters using module eigengene or intra module hub  
gene, default method of co-expression similarity  $s_{ij} = |Cor(x_i, x_j)|$