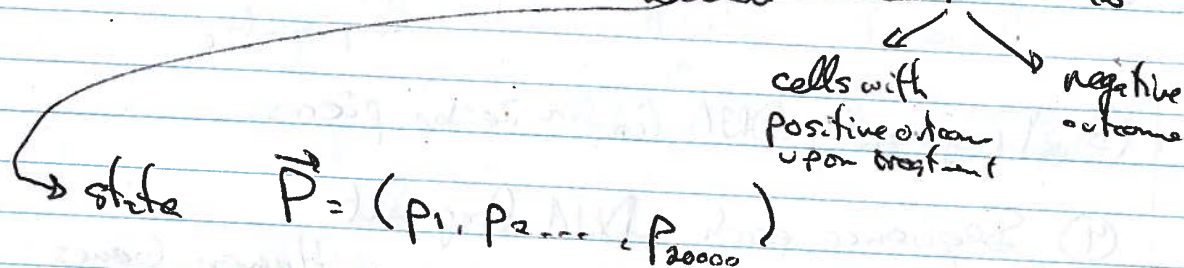


Gene Expression + Class comparison

Final Exam: Dec 14th ^{6pm} : open book, covers all topics

Goal: Capture and compare "state" of different cells



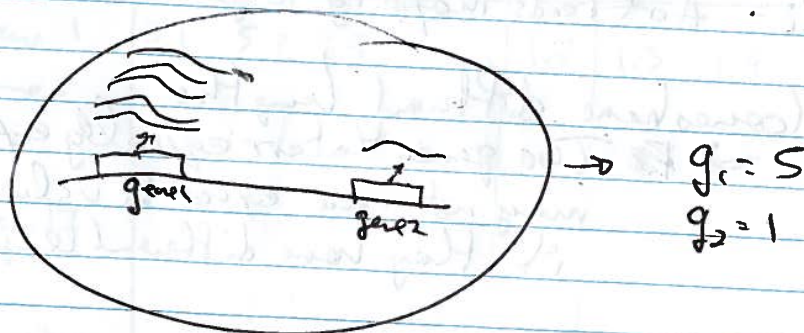
where p_i = abundance of protein p_i in cells.

Problem: Measuring protein abundance is hard (mass spectrometry)

Alternative: Measure mRNA abundance

$$\vec{G} = (g_1, g_2, \dots, g_{20000})$$

where g_i = abundance of mRNA from gene i in cells

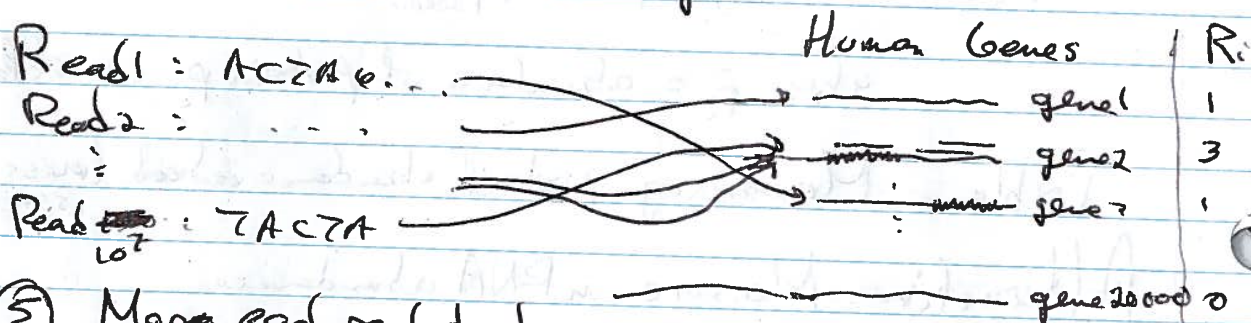


Note: $g_i \neq p_i$ because mRNA degradation
translation is regulated

RNA-sequencing (RNA-seq)

Goal: Measure gene expression levels $\vec{G} = (g_1, \dots, g_{20000})$

- ① Extract RNA from cells
- ② Reverse Transcribed RNA to cDNA
↑
complementary
- ③ Fragment cDNA in ~ 200 bp pieces
- ④ Sequence each cDNA fragment



- ⑤ Map each read to its gene
↑

Find the gene to which the read aligns

- ⑥ Count $R_i = \# \text{ of reads mapping to gene } i$

Problem: Genes have different lengths
 \Rightarrow Two genes that are equally expressed
 may not have equal R values
 if they have different lengths

- ⑦ Normalize:

$$\text{FPKM}(g_i) = \frac{R_i}{(\text{length of } g_i) \cdot (\text{total reads})}$$

↓
 Fragment per kilobase
 per million reads

(in kb) (in Millions)

Class Comparison Problem

Given: Normalized gene expression data from two sets of samples

A: (control) with $N_A = 20$ samples

$$\vec{A}_i = (A_{i,1}, A_{i,2}, \dots, A_{i,20000})$$

$A_{i,j}$ \uparrow FPKM of gene i in sample j

$$\vec{A}_2 = (\quad)$$

$$\vec{A}_{N_A} = (\quad)$$

B: (treatment) with $N_B = 25$ samples

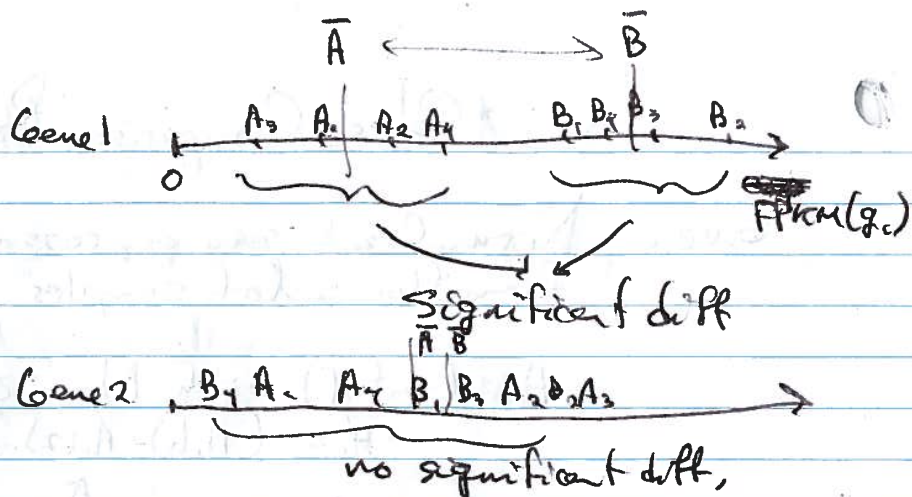
$$\vec{B}_1 = (\quad)$$

$$\vec{B}_{N_B} = (\quad)$$

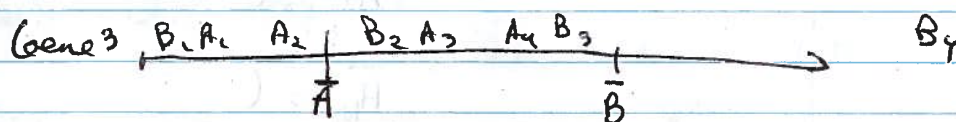
	20				25				t-stat	p-value
	A_1	$A_2 \dots$	A_{N_A}		B_1	B_{N_B}				
Gene 1	5.1	5.7	3.7		1.7	1.3	1.9	2.1	7.1	0.001
Gene 2									1.3	0.25
\vdots										
Gene 17									best genes	0.0001 = 1/10000
Gene 20000										0.35

Goal: Find genes that are "differentially expressed" b/w A, B

Example



$$d(g_i) = \text{diff. of means} = \bar{A}(g_i) - \bar{B}(g_i)$$



Student t-test (performed separately for each gene)

H_0 : Expression values from samples A and B come from the same normal distribution

$$\mu_A = \mu_B$$

$$\sigma_A = \sigma_B$$

Assumption

H_1 : $\mu_A \neq \mu_B$

$$\textcircled{1} \text{ Calculate } t(g_i) = \frac{\bar{A}(g_i) - \bar{B}(g_i)}{\sqrt{\frac{S_A^2}{N_A} + \frac{S_B^2}{N_B}}}$$

S_A^2 = Variance in A

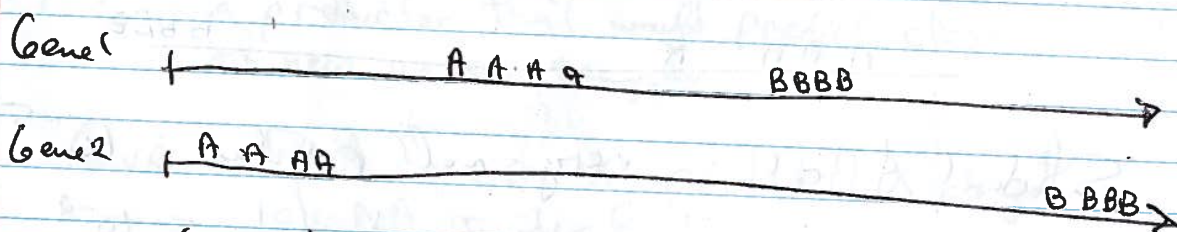
$\textcircled{2}$ Calculate [p-value for $t(g_i)$] = Prob that two random samples drawn from same distrib. would have t-statistics $t(g_i)$

Under H_0 : t follow a Student (m) ^{degrees of freedom}

$$m = \frac{\left(\frac{S_A^2}{N_A} + \frac{S_B^2}{N_B} \right)^2}{\frac{\left(\frac{S_A^2}{N_A} \right)^2}{N_A - 1} + \frac{\left(\frac{S_B^2}{N_B} \right)^2}{N_B - 1}}$$

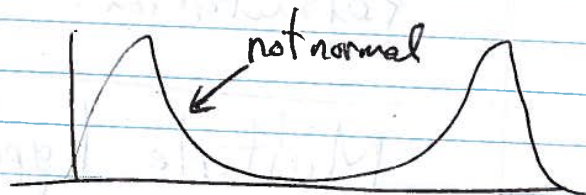
If $p\text{-value}(g_i) \leq 0.05$, then call g_i "diff. expressed"
 > 0.05 , then g_i is not diff. exp.

Issue: Violation of assumptions: - Normality of data



$$p\text{-value}(g_2) < p\text{-value}(g_1)$$

Suppose distribution under H_0



Permutation test: Estimate $p\text{-value}$ without assumption about underlying distrib. of data

For each gene i

Real	A_1	A_2	A_4	A_3
Shuffled	B_4	A_1	A_2	B_3

① Calculate $t(g_i)$

② Repeat $K=1000$ times

2.1 Randomly reshuffle class to obtain \tilde{A}, \tilde{B}

2.2 Calculate $\tilde{t}(g_i)$ from

2.3 If $\tilde{t}(g_i) \geq t(g_i)$ then

③ Report $p\text{-value}(g_i) = \frac{\text{success}}{K}$

A A A A B B B B

Student t-test : very small p-value : 1e

Permutation : p-value = $\frac{1}{\binom{8}{4}} = 10^{-4}$

Multiple hypothesis testing

We've done 20,000 tests

~~If~~ If all genes come from H_0 , then best value we would expect to observe would be: $1/20,000$

Bonferroni correction : $\text{corrected } p\text{-value}(g_i) = p\text{-value}(g_i) \cdot N$

Class Prediction Problem

	A_1	A_2	...	A_{N_A}	B_1	...	B_{N_B}	X
Gene 1	RNA-seq data							.
2								.
3								.
Gene 1000								.

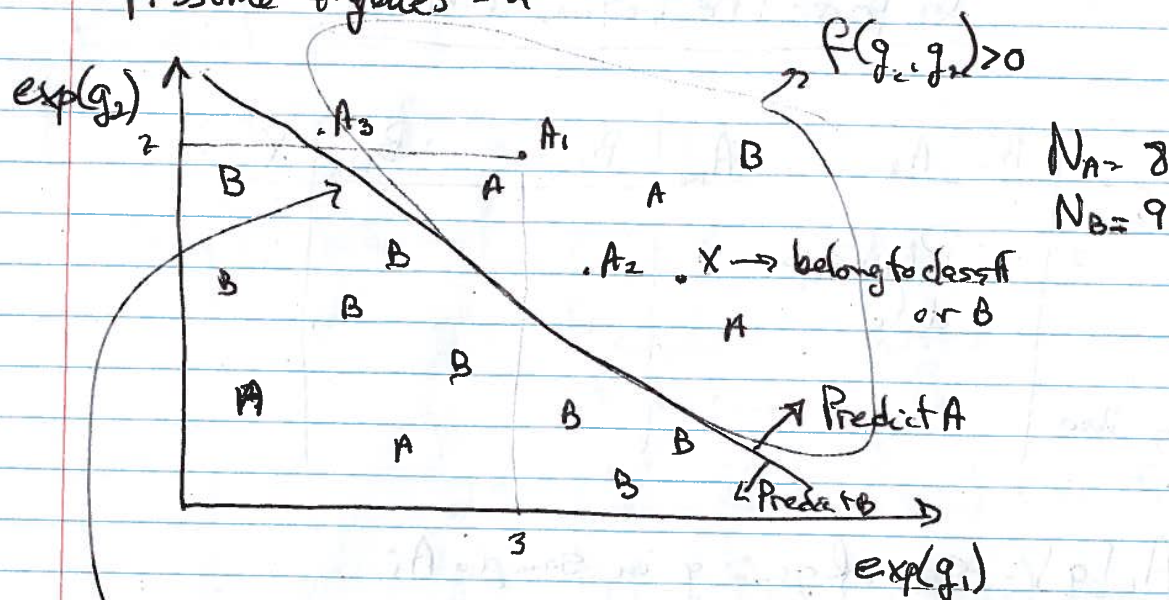
$A_i[g]$ = exp. of gene g in sample A_i

$B_j[g]$ =

Goal: From RNA-seq data from A, B , train a predictor that would predict class of new, unseen samples

Given: RNA-seq data X , predict if X belongs to class A or class B .

Assume #genes = 2



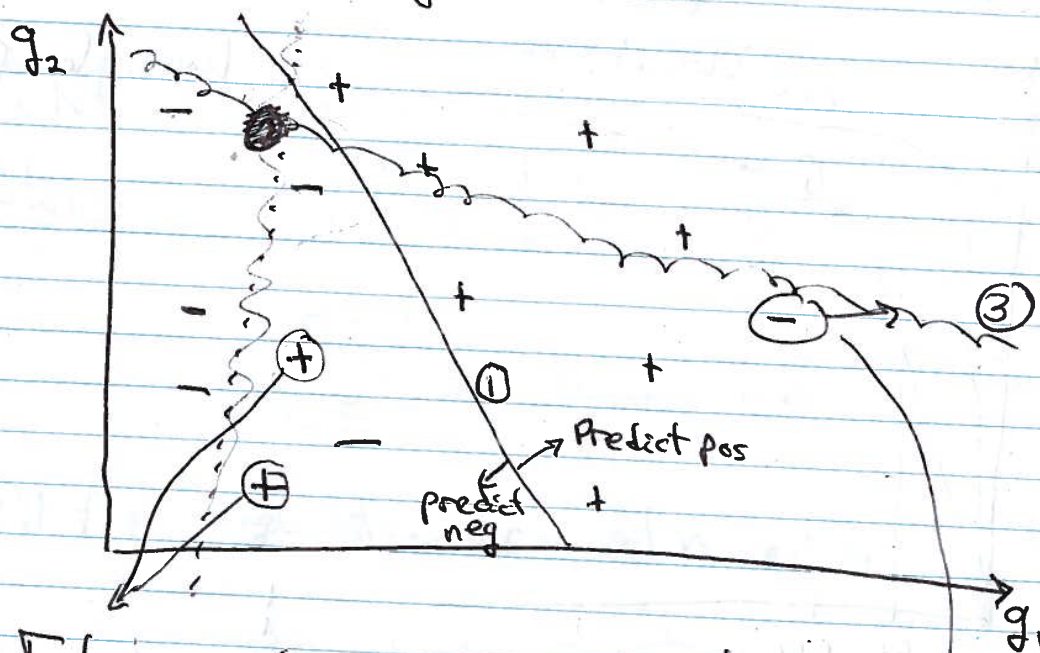
$g_0 = 10 - 2g_1$

$f(g_1, g_2) > 0 \rightarrow \text{predict A}$
 $< 0 \rightarrow \text{predict B}$

$f(g_1, g_2) = 2g_1 + g_2 - 10$

$2g_1 + g_2 - 10$

Assessing a classifier's accuracy



False-negative predictions (FN) : 2

False-positive prediction (FP) : 1

True-positive (TP) : 7

True-negative (TN) : 5

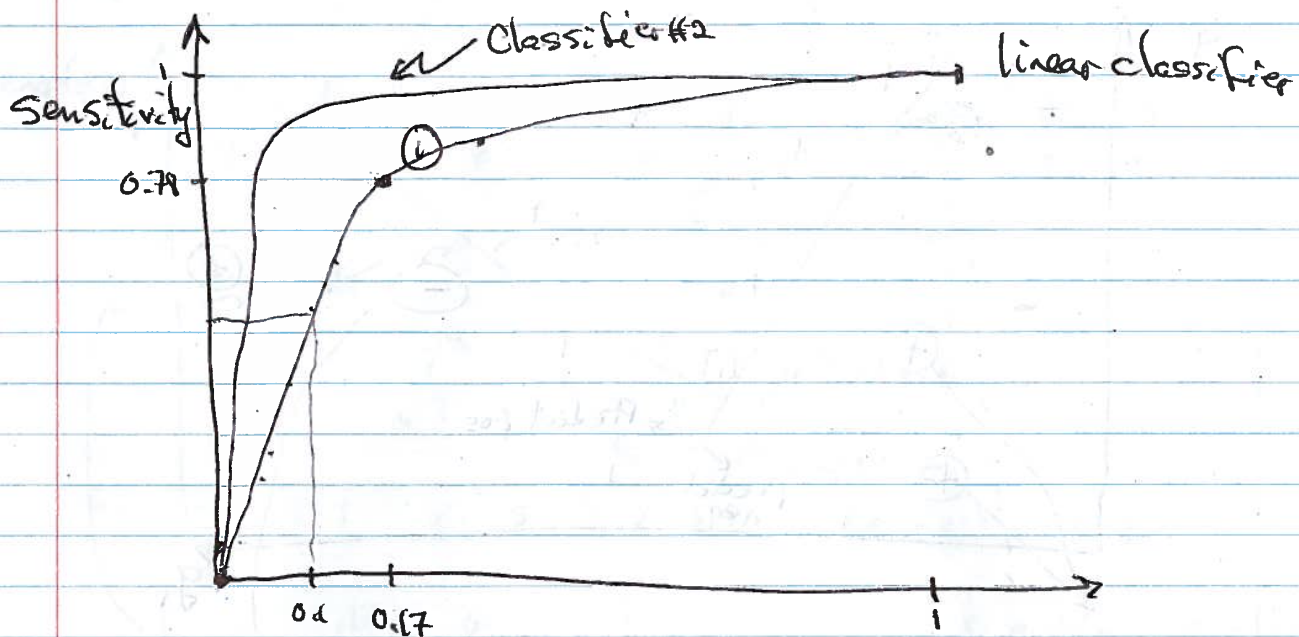
① Sensitivity: $\frac{TP}{TP+FN} = \frac{7}{7+2} = \frac{7}{9} \approx 78\%$

Specificity: $\frac{TN}{TN+FP} = \frac{5}{5+1} = \frac{5}{6} \approx 83\%$

② Sensitivity: 100%
Specificity: 50%

③ Sensitivity: $\frac{4}{9} \approx 55\%$
Specificity: 100%

Receiving-Operating Curve



1-specificity
 = False positive rate
 = Fraction of neg. examples
 that are predicted pos.