Working with genomics data

Andrew Gentles Medicine (BMIR) and Biomedical Data Sciences









Physics -> Biology



How do tumors develop?



Breast Cancer Development



5 year survival:

Stage 1: 95-100% Stage 2: 86% Stage 3: 57% Stage 4: 20%

BUMC.org

Tumors are complex cellular ecosystems

Malignant cells Myeloid cells Lymphoid cells: CD4 T cells CD8 T cells Tregs B cells



Working with data - systems biology

- What is genomic data?
- Regulatory networks
 - Transcriptional networks
- Connecting genomics with outcomes (survival)
- Dissecting tissues (e.g. tumors) at single cell resolution

Central dogma of molecular biology



DNA and RNA complementarity



Using DNA Microarrays to Measure Gene Expression

Fluorescent cDNA from muscle cell lights up myosin gene



Fluorescent cDNA from lymphocyte lights up immunoglobulin gene



We can assess activity of thousands of genes in a sample



If a gene is expressed, it will bind to its complementary probe on the array



http://learn.genetics.utah.edu/content/labs/microarray/

Next generation sequencing



Downstream analysis

https://en.wikipedia.org/wiki/RNA-Seq

Cost per Genome



Growth of gene expression omnibus at NCBI



Raw data – ballpark sizes

Data type	Typical size ~					
DNA microarray (Affymetrix cell file)	10 Mb					
TCGA histology slide scan (SVS)	300 Mb					
Exome sequence fastq (100X)	10 Gb					
Single RNA-seq experiment	10 Gb					
Single whole genome	600 Gb					
Complete TCGA dataset (11,000 tumors)	50 Pb (50,000 Tb)					

Sequencing is cheap(ish)

- What about the analysis ?
 - Storage
 - Compute costs
 - Interpretation
- It will be a long time before computational biologists are out of a job

Moderately large scale data

	J15	1 🗧 😣 🛇) (= fx			
1	A	В	С	D	E	
1	Acc ID	Exp1	Exp2	Exp3	Exp4	Exp5
2	NM_007818	67540.89	70924.09	80243.76	3501.2	5
3	NM_001105160	811.93	801.36	740.71	128.67	
4	NM_028089	190.41	211.06	236.19	9.05	
5	NM_016696	66.77	57.56	101.09	750.9	
6	NM_013459	3.3	11.29	1.89	735.82	
7	NM_007809	45.34	36.12	51.02	245.27	
8	NM_009999	103.04	370.21	200.29	17.09	
9	NM_133960	7708.78	6976.38	6569.04	1731	1
10	NM_027881	31.32	10.16	24.56	268.39	
11	NM_054053	31.32	24.83	19.84	323.68	
12	NM_007377	47.81	89.17	70.86	370.93	
13	NM_028064	703.95	689.62	662.29	214.11	
14	NM_008182	222.56	339.73	226.75	30.16	_
15	NM_013661	12.36	11.29	8.5	97.51	
16	NM_007815	20613.09	25218.13	31540.46	5209.07	

10-1000's columns (samples)

10,000's-100,000's rows (genes, transcripts)



Genomic profiles and clinical outcome



Alizadeh et al Nature 2000. Distinct types of diffuse large B-cell lymphoma identified by gene expression profiling



Sorlie et al PNAS 2001



Gentles et al JAMA 2010



Kratz et al. Lancet 2012

Associating survival with expression levels

Variable	Patient1	Patient2	Patient4	Patient5	Patient6	Patient7	•••
Time	12	3	8	35	14	22	
Status	1	1	1	0	0	1	
Gene1	1.45	0.15	-0.59	-1.88	-0.83	-0.26	•••
Gene2	0.94	-0.35	2.66	-0.23	2.09	-0.13	•••
Gene3	0.91	-0.32	-0.82	-0.35	0.86	0.32	•••

We could look at a fixed time e.g. alive/dead at 5 years – but this throws away information

Survival often assessed by hazard ratio – the increase in risk of an event (e.g. death) for each unit increase in some variable (e.g. age, expression level of a gene)

Gene expression ~ survival



"Good" genes

"Bad" genes

Clustering of outcome matrix -> biological processes



http://precog.stanford.edu

Prognostic influence of immune infiltrates

										aled	é,											
									nach	Cive										6		
								.80	, RO) at			6						Jul 2	10 310	60,0	`
						1	0 0	orden	ord r	iello d	alland	Cinal	aleral	0g	.10	in l	2	.nº	in cit	munin	Jac. H	allo
			Nº of	ion alle	~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~	OA nai	Ame	Ameni	culai	ana	ulato.	astin.	mult	5 0	005	See.	es' d	alls	alls	unsu	sum	19 119
No. data	No.	JIS Nº	IIS Me	na lu	SCV 119	0,00	11SC	115101	11500	11510	ells	ells	one	opho	opho	opho	Aritic	ditto.	colle	colle	opt	Mohn
sets	analyzed	8° 8	20. 6192	1 CO	< con	< 0° <	0° < 0	\$ x d	\$ x c	10° - 24	- 14	10	c. Ws	N. 13	C. 18	or Oal	" Oor	Wa.	10	5 40	21. 46	50°
5	744	1.85 -1.53	1.53	2.62		2.67	-2.38	-1.61	2.28	-1.49	1.98				0.82	-2.03	-0.80	-2.10	0.50	3.17	015	AML
1	174	-1.61 1.44					-2.79			-1 16	149		1,94	-1.40		100	-0,85					B-ALL
1	107	2.08 11.00	2.44		1.39 -1	51 -2.06	1	-1,83		-2.39	0.54	-2.60			2.34		2.10	-0.84	-2.38			CLL
1	158	2.07		1.16	-1.	07 1 34	-2.21	-1-35	-1-11				-2.71	-1.35	-1.83							Burkitt's lymphoma
3	594	2.28			-2	47 1 99	-3.33	-1.70	-1.46	9.81		2.94	-4.03		3.56		1.51	-2.96				DLBCL
1	180	-0.92	1.50			-1.57	-1.62				132				1		-1,06	1.53				FL
2	189	-1.11	-0.44					1.67		1.72	1.82	-1.05	1.51		-0.93			-2.11	1.84			Multiple myeloma
2	70					100 -0 00	1.76	1.01	0.90	100		0.84	-		10 100	_			1.07		2.02	Astrocutoma
6	283	-100	-1.05		-1	71	3.11		1.60			-2.80	2.49					-1.62		-1.16		Glioblastoma
1	30						-0.75				-1.78			-1.20		-0.99	215	0 70		hatte		Meningioma
1	15	1023 1 11	1.91				-1.07			-0.62						1.43	-0.93				-1.63	Oligodendroglioma
-3																						
1	30	-1.16 -1.62		-0	1.97 2	20	-0.89	-2.07		1.79					-0.42		1.50		1.79			Bladder cancer
4	567	612	-1.62		0.76			-1.28			-1.51										3.91	Breast cancer
3	236	-1:02 2.49	-1.74	0:68	.75 -0.	91 -0.90		-1.80						-1:08								Colon cancer
1	20	0.001 (0.01	-1.86				-1.28				-1.27										2.58	Ewing sarcoma
1	18	0.00 1.90	-0.87								-1.04	2.03		-1.54	-1.72						-1.27	Gastric cancer
2	96	-2.21 1 18	-1.97 -	2.01 -1	43 2	34 -1.69	-2.70	-1.40				-1.48			2.05	2.27			3,27			Germ cell tumors
2	76	00-14		1.18		-1.39		-1:81			-1.46					-2.15						Head and neck cancer
9	902	-1.43 -1.40	-1.59		-2	67 3.00	-0.99					-1.09	2.62	2.04	2.36	-1.37	3.64	2.94				Lung adenocarcinoma
7	408	2.02	-	1.22		-1.75	-0.51	-1 23								45.54	-1.92	1 26				Lung squamous cell carcinoma
2	26	-1,25	2.71			-1 15		-2.04							-1.35	-1.70	0.45					Lung large cell carcinoma
1	19	-0.54	-1.35			-1.55	1:86	-1:29	-1.32											j l		Melanoma primary
2	62	-1.96		000-1	.54	55 0.5	-1.54				1.14											Melanoma metastasis
1	33	1 15		1.10	644	-1.56	1.30						-1.36									Osteosarcoma
6	745	1.42	-1.50	1:06	2	28 -2.69	-0.73															Ovarian cancer

Network types

- Protein-protein
- Protein-DNA
- miRNA-RNA
- Transcriptional (expression) networks
- Signaling networks



Sachs et al. http://www.sciencemag.org/content/308/5721/523.full

Gene regulatory networks

They are a "flowchart" of interactions, not a detailed model

Transcriptional regulatory networks

- Imperfect but powerful representation of the system
- Example: Glioblastoma subtypes
 - Hairball -> targets of transcription factors -> which are differentially expressed between GBM types

- How do these networks compare in size to ones in the physical world
 - E.g. electrical system, telecommunications
- What sort of qualitative differences are there?

• Do cats or dogs produce worse hairballs?

Heterogeneity of tissues (and tumors)

Single-cell RNA-seq

Profiling of sorted populations

Can you tell what's in a smoothie by taste?

• How many different flavours?

• How small a proportion of the mix?

• Is it easier to detect presence or absence ?

Averaging across cell types

Trapnell et al. 2015 Genome Res. 25: 1491-1498

Why single cell?

Trapnell et al. 2015 Genome Res. 25: 1491-1498

Single cell RNA-seq approaches

Papalexi & Satuja Nat Rev Imm 2018

Mass cytometry – single cell proteomics

Summary

- We have technologies to measure many things at the same time on large numbers of samples
- We can test which of thousands of genes are connected with things of importance such as how long patients survive
- Which genes can tell us about why some patients live longer
- The immune system is very important almost doesn't matter what treatment
- Possible treatment options