

Statistical methods for high-dimensional biology

STAT 540 2015

Data exploration

Paul Pavlidis

Our prototype data set

Each value =
Measurement of one gene one sample

Samples

Assays

Case type B

gene, in a sample

Assays

assay

One row per gene (200000)

What you usually get

Data file: contains the measured data

Metadata file: describes the samples

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q
1	probe		8.1	54.1	36.1	23.1	17.1	40.1	45.1	55.1	11.1	38.1	26.1	41.1	31.1	45.	
2	1007_s_at	7.926	7.0559	7.9411	7.7977	7.4689	7.5711	7.4303	7.2057	7.4715	7.1906	8.2236	7.2317	8.0644	7.73	8.0763	7.642
3	1053_at	6.6365	7.1358	7.5862	7.8645	7.4006	7.3865	7.9948	7.4433	5.6344	7.903	7.026	7.057	8.3421	7.6931	7.7	
4	117_at	8.959	4.6803	10.318	10.271	7.9707	10.321	10.456	10.382	9.5865	9.6003	9.3467	7.9766	9.0111	10.051	9.8787	10.35
5	121_at	8.855	4.855	10.318	10.271	7.9707	10.321	10.456	10.382	9.5865	9.6003	9.3467	7.9766	9.0111	10.051	9.8787	10.35
6	1255_g_at	8.9677	8.3877	8.9665	8.9394	8.9049	8.3048	3.7421	3.4679	8.9584	8.3263	8.9448	8.3905	4.024	4.0581	4.0189	3.869
7	1294_at	8.0786	9.0139	8.9042	8.8864	8.0752	3.1849	9.0646	9.3711	8.0459	9.3193	8.1166	8.7887	9.2214	9.0581	8.9018	9.332
8	1302_at	8.959	5.7948	7.5862	7.8645	7.4006	7.3865	7.9948	7.4433	5.6344	7.903	7.026	7.057	8.3421	7.6931	7.7	
9	1320_at	8.147	5.1389	5.2289	5.1778	5.2646	4.6999	4.9862	5.0379	5.1504	4.7234	5.1441	5.0783	5.2984	5.3087	5.1752	4.987
10	1405_i_at	6.2639	7.5861	6.7248	5.8717	6.1997	10.464	6.0142	5.6864	5.87	6.577	6.0243	6.3546	7.5059	7.5744	5.403	
11	1431_at	4.1273	3.8827	4.2963	4.437	4.5913	4.3935	4.4605	4.4572	4.4081	4.9116	4.4206	4.3816	4.7248	4.437	4.7811	4.4576
12	1438_at	6.7376	5.9663	6.7058	6.804	6.0482	5.1074	1.491	1.461	5.974	6.974	6.8844	6.6017	6.6978	6.6462	6.793	
13	1487_at	8.1541	7.9911	8.449	8.6151	8.849	8.6511	8.0206	8.5319	8.489	8.5778	8.712	8.712	8.712	8.712	8.712	
14	1494_f_at	5.6342	5.6148	6.5304	6.5367	6.3581	5.695	6.314	6.1736	6.8538	5.7593	6.241	6.3872	6.4009	6.3022	6.7324	6.033
15	1552256_e	8.2499	7.537	7.9888	8.3884	7.624	6.9916	7.9671	7.673	7.5446	8.2511	7.9856	7.6377	8.1844	8.3138	8.2106	7.464
16	1552257_e	7.8619	6.6159	7.7546	7.7172	7.1255	7.2768	7.2925	8.609	7.6841	7.7207	7.4003	7.17	7.9793	7.0081	7.6714	7.316
17	1552258_e	7.8622	6.6162	7.7546	7.7172	7.1255	7.2768	7.2925	8.609	7.6841	7.7207	7.4003	7.17	7.9793	7.0081	7.6714	7.316
18	1552259_e	7.2138	5.8483	5.2059	5.2098	5.2198	4.6997	4.798	5.3421	5.3708	5.4281	5.3454	5.3999	4.9337	5.1037	5.124	
19	1552260_e	8.1053	8.8003	8.5589	8.5178	8.2161	8.0552	9.1113	9.1421	8.3856	8.1584	8.0124	7.9257	9.142	8.9135	8.6977	8.973
20	1552261_e	8.1056	8.8006	8.5592	8.5181	8.2164	8.0555	9.1116	9.1424	8.3859	8.1587	8.0127	7.9257	9.142	8.9136	8.6978	8.974
21	1552262_e	8.1056	8.8006	8.5592	8.5181	8.2164	8.0555	9.1116	9.1424	8.3859	8.1587	8.0127	7.9257	9.142	8.9136	8.6978	8.974
22	1552263_e	4.4192	4.2144	4.437	4.5913	4.3935	4.4655	4.4572	4.4081	4.9116	4.4206	4.3816	4.7248	4.437	4.7811	4.4576	4.503
23	1552264_e	6.3119	6.545	6.3908	6.2609	6.2344	6.3756	6.2433	6.3708	6.2118	6.1392	6.2828	7.1272	6.4115	6.3586	6.6803	6.58
24	1552272_e	8.3128	8.3465	7.7187	7.5153	7.7277	7.4438	7.2911	7.5992	7.7983	7.7883	7.178	7.8233	7.6738	7.5345	7.195	
25	1552273_e	8.3128	8.3465	7.7187	7.5153	7.7277	7.4438	7.2911	7.5992	7.7983	7.7883	7.178	7.8233	7.6738	7.5345	7.195	
26	1552274_e	8.3128	8.3465	7.7187	7.5153	7.7277	7.4438	7.2911	7.5992	7.7983	7.7883	7.178	7.8233	7.6738	7.5345	7.195	
27	1552275_e	8.7279	6.6933	6.6935	6.6933	6.2294	6.6975	6.5971	6.4809	6.2384	6.5713	6.4075	6.0216	6.749	6.6996	6.7675	6.633
28	1552276_e	8.7279	6.6933	6.6935	6.6933	6.2294	6.6975	6.5971	6.4809	6.2384	6.5713	6.4075	6.0216	6.749	6.6996	6.7675	6.633
29	1552277_e	8.7279	6.6933	6.6935	6.6933	6.2294	6.6975	6.5971	6.4809	6.2384	6.5713	6.4075	6.0216	6.749	6.6996	6.7675	6.633
30	1552278_e	8.7279	6.6933	6.6935	6.6933	6.2294	6.6975	6.5971	6.4809	6.2384	6.5713	6.4075	6.0216	6.749	6.6996	6.7675	6.633
31	1552279_e	8.7279	6.6933	6.6935	6.6933	6.2294	6.6975	6.5971	6.4809	6.2384	6.5713	6.4075	6.0216	6.749	6.6996	6.7675	6.633
32	1552280_e	8.7279	6.6933	6.6935	6.6933	6.2294	6.6975	6.5971	6.4809	6.2384	6.5713	6.4075	6.0216	6.749	6.6996	6.7675	6.633
33	1552281_e	8.7279	6.6933	6.6935	6.6933	6.2294	6.6975	6.5971	6.4809	6.2384	6.5713	6.4075	6.0216	6.749	6.6996	6.7675	6.633
34	1552282_e	8.7279	6.6933	6.6935	6.6933	6.2294	6.6975	6.5971	6.4809	6.2384	6.5713	6.4075	6.0216	6.749	6.6996	6.7675	6.633
35	1552283_e	8.7279	6.6933	6.6935	6.6933	6.2294	6.6975	6.5971	6.4809	6.2384	6.5713	6.4075	6.0216	6.749	6.6996	6.7675	6.633
36	1552284_e	8.7279	6.6933	6.6935	6.6933	6.2294	6.6975	6.5971	6.4809	6.2384	6.5713	6.4075	6.0216	6.749	6.6996	6.7675	6.633
37	1552285_e	8.7279	6.6933	6.6935	6.6933	6.2294	6.6975	6.5971	6.4809	6.2384	6.5713	6.4075	6.0216	6.749	6.6996	6.7675	6.633
38	1552286_e	8.7279	6.6933	6.6935	6.6933	6.2294	6.6975	6.5971	6.4809	6.2384	6.5713	6.4075	6.0216	6.749	6.6996	6.7675	6.633
39	1552287_e	8.7279	6.6933	6.6935	6.6933	6.2294	6.6975	6.5971	6.4809	6.2384	6.5713	6.4075	6.0216	6.749	6.6996	6.7675	6.633
40	1552288_e	8.7279	6.6933	6.6935	6.6933	6.2294	6.6975	6.5971	6.4809	6.2384	6.5713	6.4075	6.0216	6.749	6.6996	6.7675	6.633
41	1552289_e	8.7279	6.6933	6.6935	6.6933	6.2294	6.6975	6.5971	6.4809	6.2384	6.5713	6.4075	6.0216	6.749	6.6996	6.7675	6.633
42	1552290_e	8.7279	6.6933	6.6935	6.6933	6.2294	6.6975	6.5971	6.4809	6.2384	6.5713	6.4075	6.0216	6.749	6.6996	6.7675	6.633
43	1552291_e	8.7279	6.6933	6.6935	6.6933	6.2294	6.6975	6.5971	6.4809	6.2384	6.5713	6.4075	6.0216	6.749	6.6996	6.7675	6.633
44	1552292_e	8.7279	6.6933	6.6935	6.6933	6.2294	6.6975	6.5971	6.4809	6.2384	6.5713	6.4075	6.0216	6.749	6.6996	6.7675	6.633
45	1552293_e	8.7279	6.6933	6.6935	6.6933	6.2294	6.6975	6.5971	6.4809	6.2384	6.5713	6.4075	6.0216	6.749	6.6996	6.7675	6.633
46	1552294_e	8.7279	6.6933	6.6935	6.6933	6.2294	6.6975	6.5971	6.4809	6.2384	6.5713	6.4075	6.0216	6.749	6.6996	6.7675	6.633
47	1552295_e	8.7279	6.6933	6.6935	6.6933	6.2294	6.6975	6.5971	6.4809	6.2384	6.5713	6.4075	6.0216	6.749	6.6996	6.7675	6.633
48	1552296_e	8.7279	6.6933	6.6935	6.6933	6.2294	6.6975	6.5971	6.4809	6.2384	6.5713	6.4075	6.0216	6.749	6.6996	6.7675	6.633
49	1552297_e	8.7279	6.6933	6.6935	6.6933	6.2294	6.6975	6.5971	6.4809	6.2384	6.5713	6.4075	6.0216	6.749	6.6996	6.7675	6.633
50	1552298_e	8.7279	6.6933	6.6935	6.6933	6.2294	6.6975	6.5971	6.4809	6.2384	6.5713	6.4075	6.0216	6.749	6.6996	6.7675	6.633
51	1552299_e	8.7279	6.6933	6.6935	6.6933	6.2294	6.6975	6.5971	6.4809	6.2384	6.5713	6.4075	6.0216	6.749	6.6996	6.7675	6.633
52	1552299_e	8.7279	6.6933	6.6935	6.6933	6.2294	6.6975	6.5971	6.4809	6.2384	6.5713	6.4075	6.0216	6.749	6.6996	6.7675	6.633
53	1552299_e	8.7279	6.6933	6.6935	6.6933	6.2294	6.6975	6.5971	6.4809	6.2384	6.5713	6.4075	6.0216	6.749	6.6996	6.7675	6.633
54	1552299_e	8.7279	6.6933	6.6935	6.6933	6.2294	6.6975	6.5971	6.4809	6.2384	6.5713	6.4075	6.0216	6.749	6.6996	6.7675	6.633
55	1552299_e	8.7279	6.6933	6.6935	6.6933	6.2294	6.6975	6.5971	6.4809	6.2384	6.5713	6.4075	6.0216	6.749	6.6996	6.7675	6.633
56	1552299_e	8.7279	6.6933	6.6935	6.6933	6.2294	6.6975	6.5971	6.4809	6.2384	6.5713	6.4075	6.0216	6			

Ready for exploration

- Understand/get a feel for the data
- Formulate hypotheses / develop models
- Identify problems

The biggest mistake in data analysis

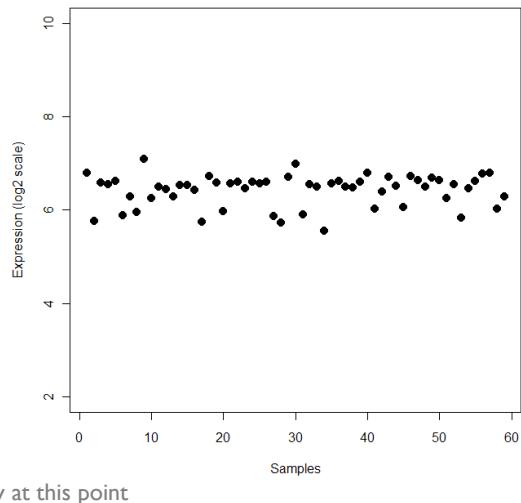
If you don't look at the data ... you are going to get in trouble.

- Not just at the beginning, but at every stage.
- That could mean making graphs, or staring at numbers. Probably both.
- “Sanity checks” should make up a lot of your effort.

What are some “first questions” you might ask?

Data for one gene

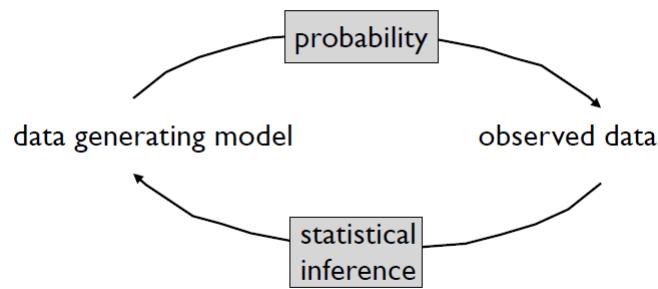
Why do the values vary?



Order of samples is arbitrary at this point

Making sense of the data

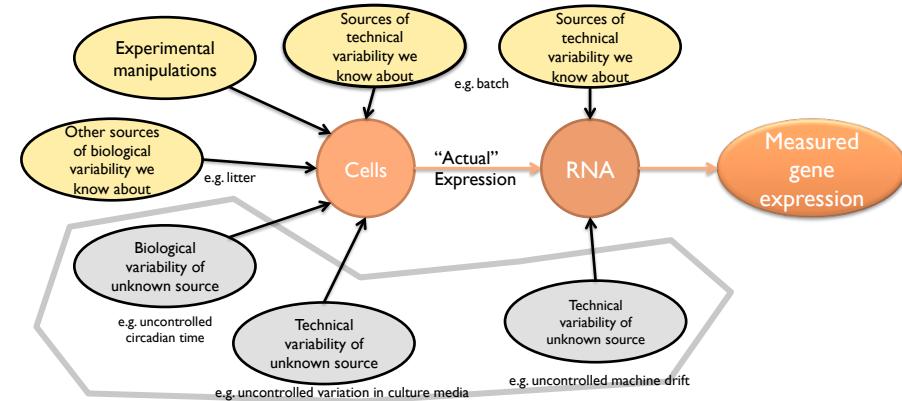
- Review: Data is what we observe, we want to infer something about “where it came from”



Jenny Bryan

Model of where expression data comes from

- The measured expression level of gene_i is the combination of many effects.
- Analysis goal is often to determine relative role of effects, separate interesting from “uninteresting”
- One person’s noise is another’s signal.



Variability: friend and foe

- First line of defense: Know the enemy
 - You can only “correct” for things you know about.
 - Keep track of potential sources of variance: Batches of reagents, slides, personnel
- Design experiments to minimize impact of technical variability.
 - Avoid batches / minimize batch differences
 - Randomize design with respect to batches
- Replication
 - Biological (important)
 - Technical (*usually* less important but might need to convince yourself)
- Much more later in course

Exploratory numbers

Use as a rough description of the data

- Range
- Mean, Mode, Variance
- Number of missing values

Examine these values per-sample and per-element (e.g. gene)!

Exploratory graphics

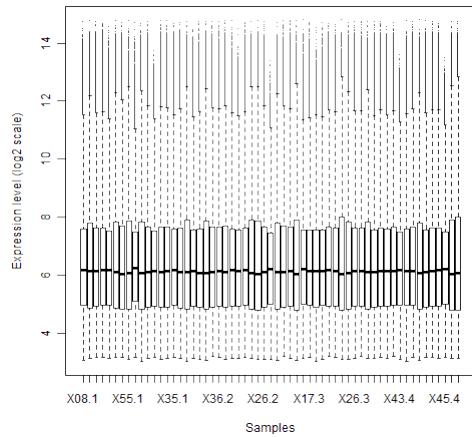
- “Exploratory analysis” is “compute a little and graph a lot”
- I’ll show a few simple approaches that are common in genomics
- Basis for examples: a data set of 59 Affymetrix microarrays, each of which have 54675 measurements. RNA was from blood cells of patients being treated for periodontal infection.

Boxplots to compare samples

Quick and dirty; Reasonable tool for summarizing large amounts of data.

Not so great if your distribution is multimodal.

Don't use boxplots if you don't have to
– show the actual data (esp. for small sample sizes!)

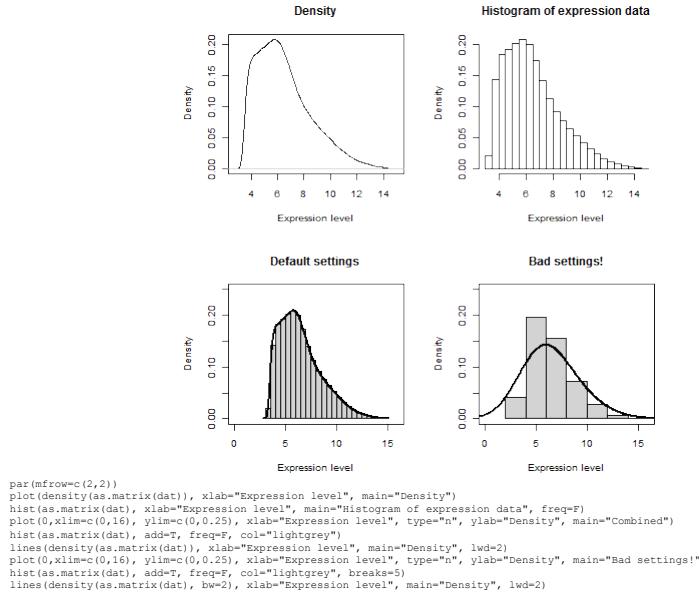


```
boxplot(dat, pch='.', xlab="Samples", ylab="Expression level (log2 scale)")
dev.print("basic.boxplot.png", device=png, width=500)
boxplot(dat[,1], pch='.', xlab=names(dat)[1], ylab="Expression level (log2 scale)")
```

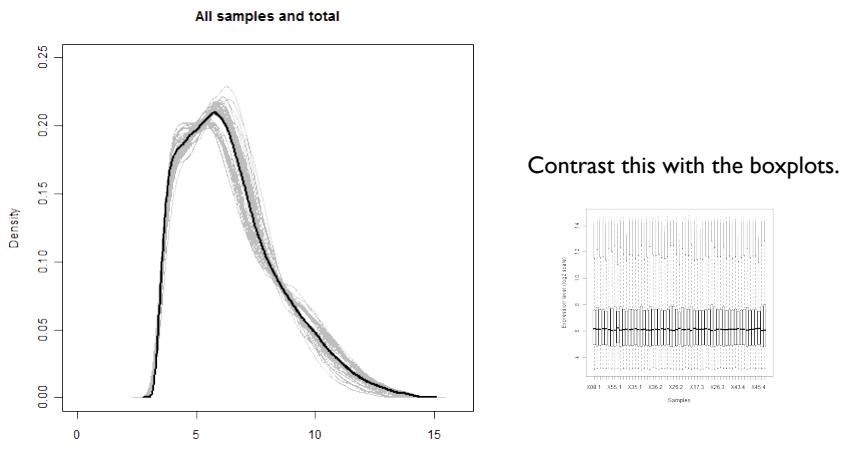
Histograms

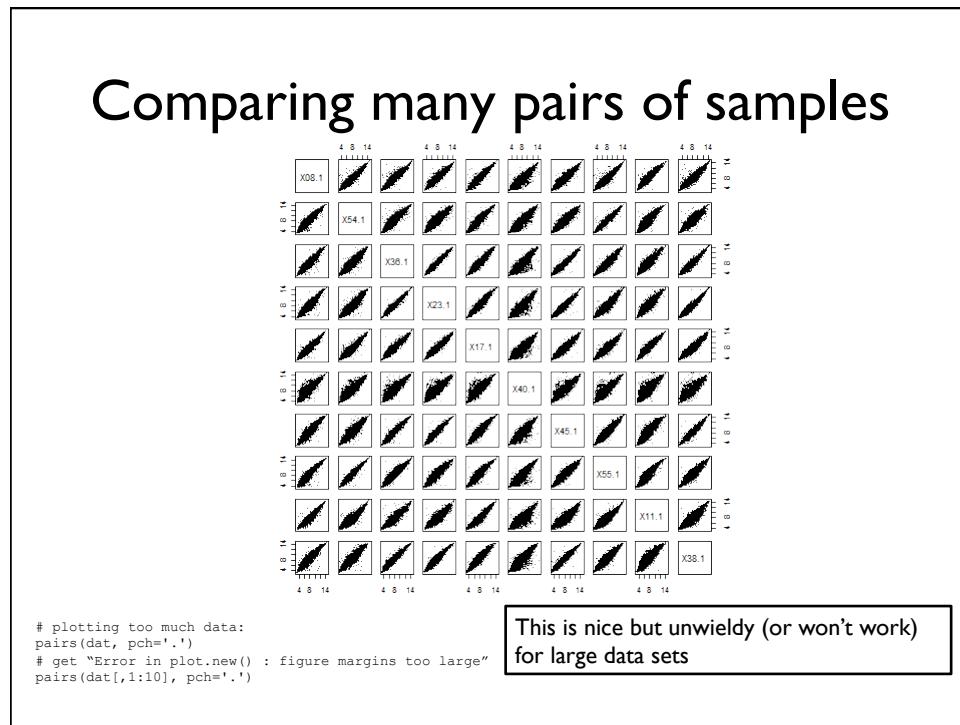
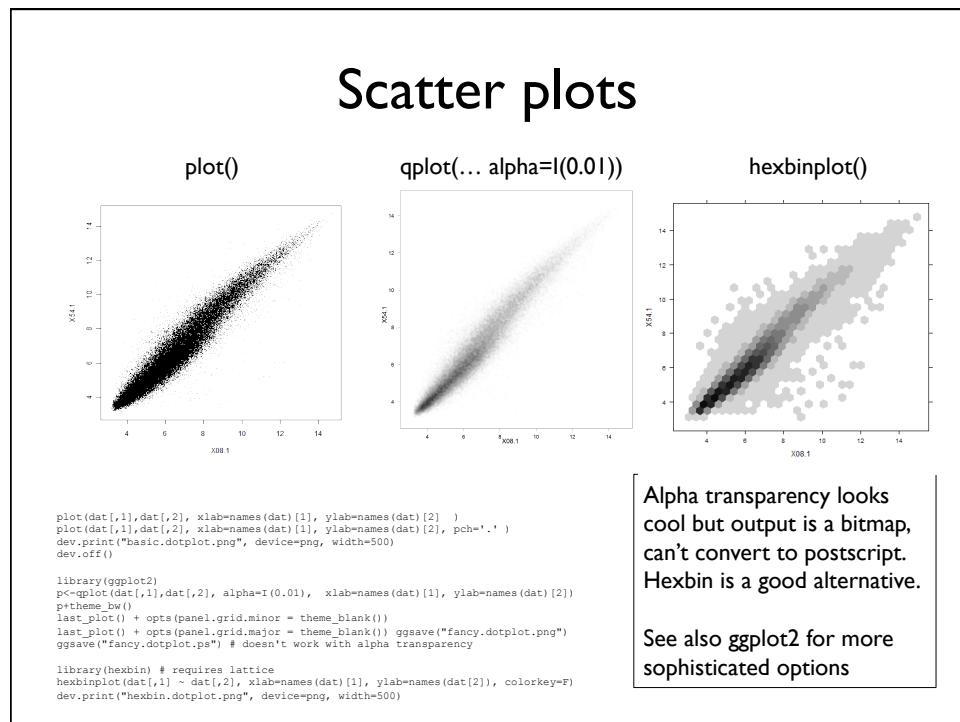
- Box plot's big brother but takes up more space.
- Consider using “density()” instead
- Choose sensible bin sizes and bandwidth

Histogram and density examples

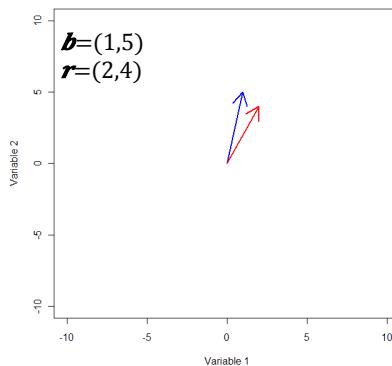


Overlay densities of all samples

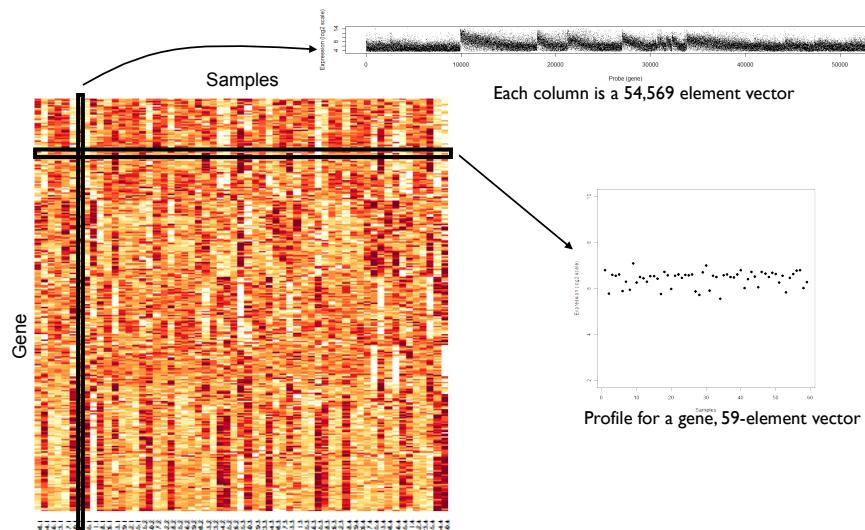




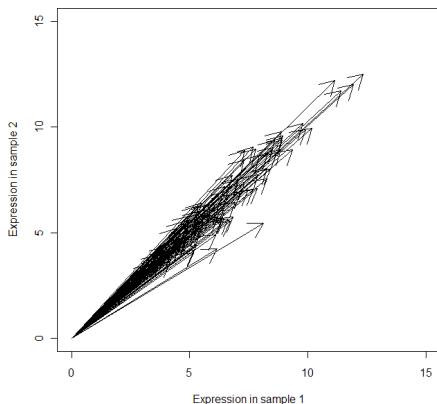
Digression into vectors



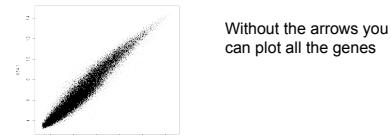
Rows and columns of your data matrix are vectors



First two dimensions of a bunch of gene vectors



Each vector is a gene (just 100 shown).
Variables/Coordinates are “assayed samples”

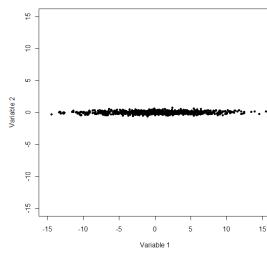


Without the arrows you
can plot all the genes

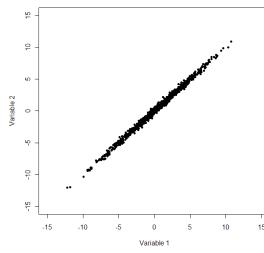
Do the genes form “interesting patterns” in “sample space”?

Data as points in space

- Distributed with a certain shape (and orientation.)
- Want to find “interesting patterns” in these “clouds”
- Do we need all 59 (or 54569) dimensions?
 - Often we are interested in finding interesting lower-dimensional structure in the (or representations of) data (regression, clustering, PCA)
- Motivating examples:



One of these dimensions is useless

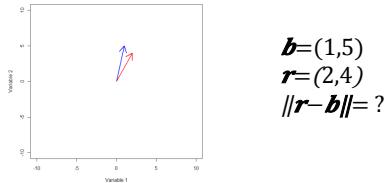


Dimensions are redundant

Comparing vectors: Euclidean distance

- The length of the difference between two vectors.
- In high dimensions, distances get big.

$$d(x, y) = \|x - y\| = \sqrt{\sum (x_i - y_i)^2}$$

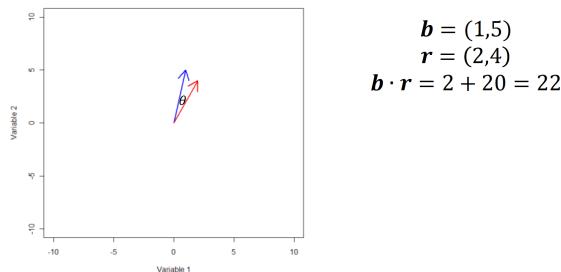


For matrices in R: `dist()`

Comparing vectors: Dot product

- Also known as **inner product**
- If two vectors are “nearer” each other, it’s bigger
- Can be negative, unlike Euclidean distance

$$\mathbf{x} \cdot \mathbf{y} = \sum x_i y_i = \cos \theta \|\mathbf{x}\| \|\mathbf{y}\|$$



In R: `%*%`

Covariance

Recall: Sample variance $s^2(\mathbf{x}) = \frac{\sum_i^N (x_i - \bar{x})^2}{N - 1}$

Sample covariance $\text{Cov}(\mathbf{x}, \mathbf{y}) = \frac{\sum_i^N (x_i - \bar{x})(y_i - \bar{y})}{N - 1}$

“Are you far from your mean when
I am far from mine?”

What happens if the data are mean centered?

Comparing vectors: Pearson correlation

- Two vectors are pos. correlated if they point “in the same direction”
 - Equivalently, if the angle between them is small
- Always between 1 and -1 , no matter how many dimensions.
- What happens if the data is **standardized** (mean 0, var 1)?

Covariance:

$$r(\mathbf{x}, \mathbf{y}) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

Variances: “normalizes” the covariance

In R: `cov()` and `cor()`; `scale()` for standardizing

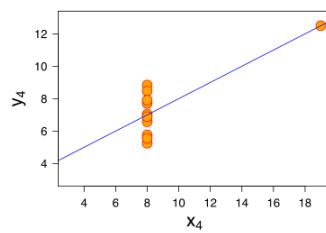
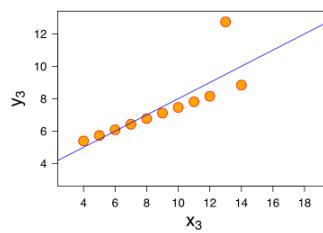
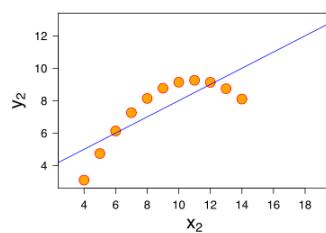
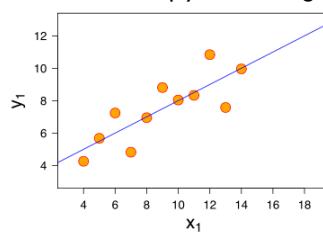
Self-study exercise: think about the metrics on the last few slides for cases of:

- orthogonal vectors
- vectors with opposite direction
- vectors with same direction
- High-dimensional data vs. lower-dimensional data
- What happens if variance of one vector is zero
- etc...

Supplement with sketches and/or plots in R

Look at the data!

Correlation does not imply “interesting”, and lack of correlation does not mean “boring”

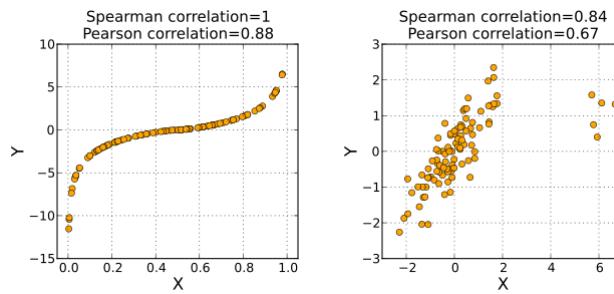


http://en.wikipedia.org/wiki/File:Anscombe%27s_quartet_3.svg

Remedy for some problems: ranks

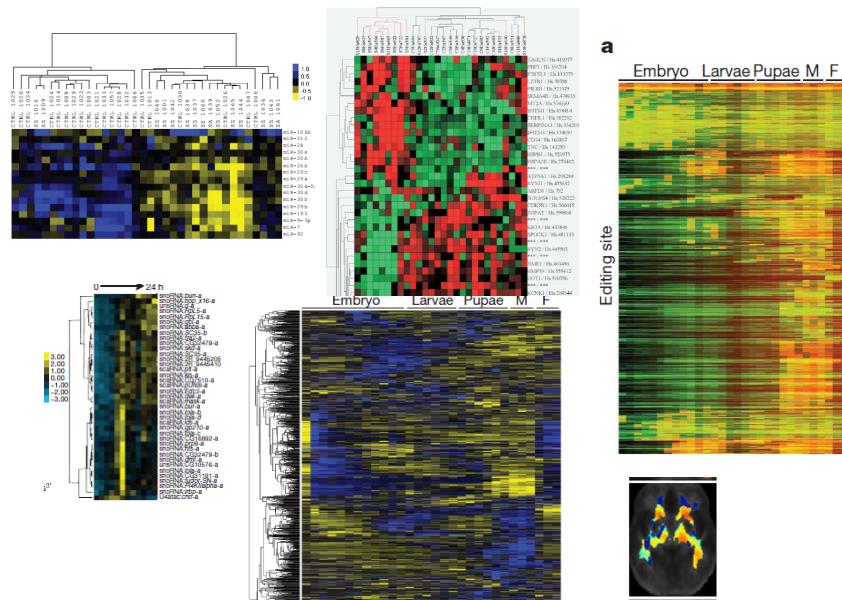
The correlation of the ranks is robust to “non-linear” relationships and outliers.

- Annoyance: dealing with ties



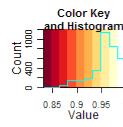
Wikipedia

Heat maps



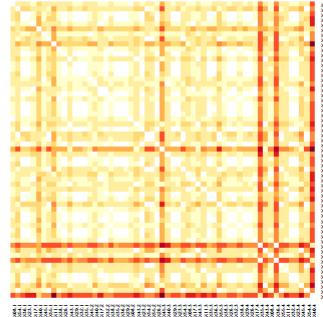
Using heat maps to compare assays

- Pairwise Pearson correlations of entire sample expression profiles



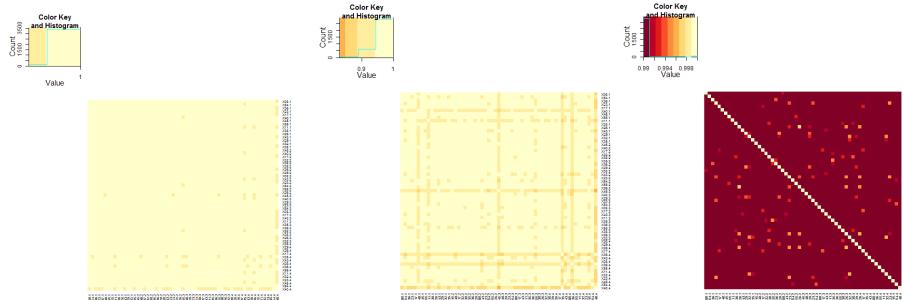
- This method loses some information, but is suitable for quite large data sets.
- Not for comparing features (genes)!

```
library(gplots)
library(RColorBrewer)
cols<-c(rev(brewer.pal(9,"YlOrRd")), "#FFFFFF")
heatmap.2(cor(dat), Rowv=NA, Colv=NA, symm=T,
trace="none", dendrogram="none", col=cols, cexCol=0.5,
cexRow=0.5)
dev.print("heatmap.png", device=png, width=500)
```



Alternative: matrix2png (<http://www.chibi.ubc.ca/matrix2png>)

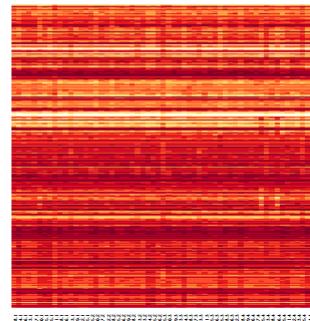
Choose an appropriate range



R sets “reasonable” ranges for you; but you can control it (`breaks`) and will often want to.

Setting up data heat maps

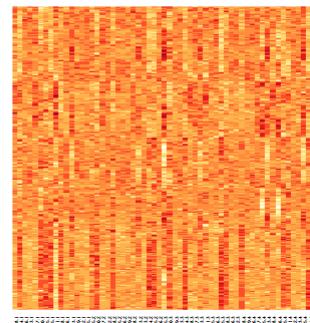
- I've taken 500 rows at random from my big data set
- The heat map here was created using this data "as is".
- I did cheat a little: The row ordering is not random.
- Note that this lecture doesn't get into dealing with missing data – how to evaluate it and visualize it



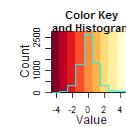
```
heatmap.2(bitOfData, Rowv=NA, Colv=NA, scale=NULL, trace="none", dendrogram="none", col=cols, cexCol=0.5, labRow=NA)
```

Revision I

- Same input data, but each row is scaled to have mean 0 and variance 1 (z-scores)
 - Subtract the mean; divide by the standard deviation. use `scale()` on the data rows.
- It is now easier to compare the rows and seem some structure.
- But looks kind of bland compared to ones you often see in the literature



Note: R heatmap scales rows by default. You must disable this to give yourself control over how it turns out.

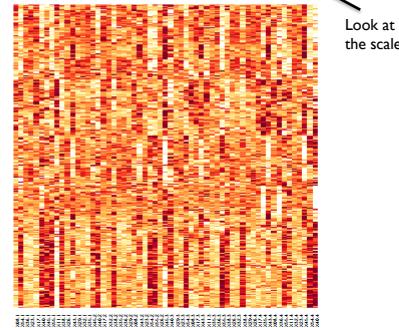
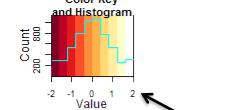


```
heatmap.2(scaledBit, Rowv=NA, Colv=NA, scale=NULL, trace="none", dendrogram="none", col=cols, cexCol=0.5, cexRow=0.5, labRow=NA)
```

Revision 2: Adjusting contrast

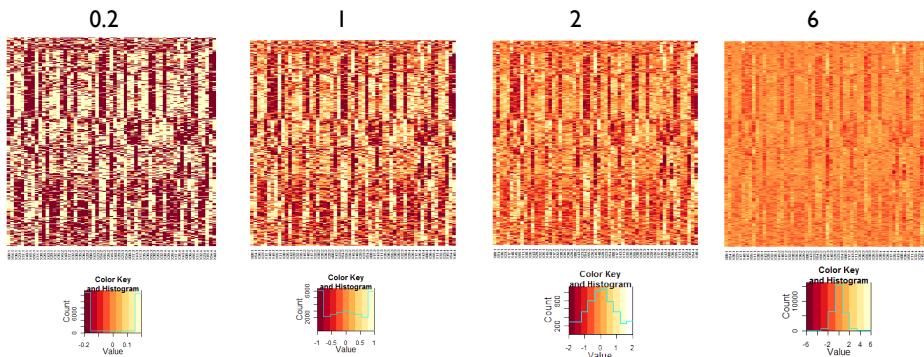
- Same input data again, but after scaling I clipped the range of values to $\{-2, 2\}$
- This means “anything more than two standard deviations from the row mean is set to 2”
- Now extremes have no effect; “Higher contrast”
- Limit values of the range limit of 2 or 3 are usually good.

Heat maps you see in the literature are almost always set up this way. Be aware of what's really going on.



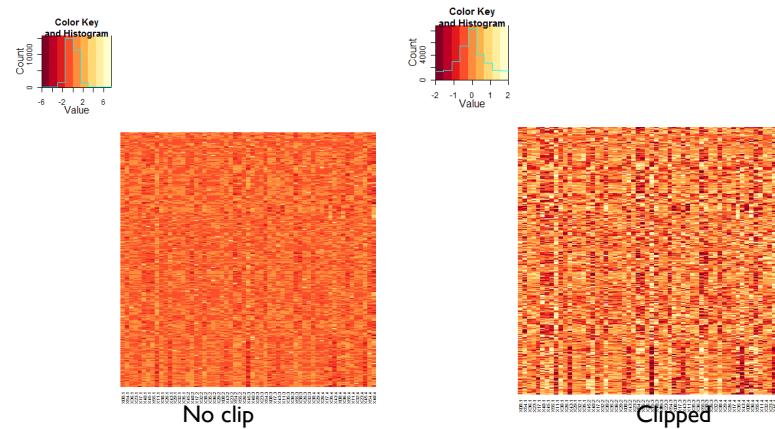
```
limit<-2
trim.scaledBit<-scaledBit
trim.scaledBit[which(trim.scaledBit < -limit)]<- -limit
trim.scaledBit[which(trim.scaledBit > limit)]<- limit
heatmap.2(trim.scaledBit, Rowv=NA, Colv=NA, scale=NULL, trace="none",
col=cols, cexCol=0.5, cexRow=0.5, labRow=NA)
```

Varying clip limit



Another reason why clipping is and isn't good

I added random “spiky” values to the data (I just multiplied them by 3)
 Clipping hides these “outliers” but allows us to see variation in the bulk of the data



Plotting too much data



The entire data set.

If the cells are less than 1 pixel,
 everything starts to turn to mush
 and can even be misleading.

(This won't work in R unless
 you print directly to a file.)

Choice of colours

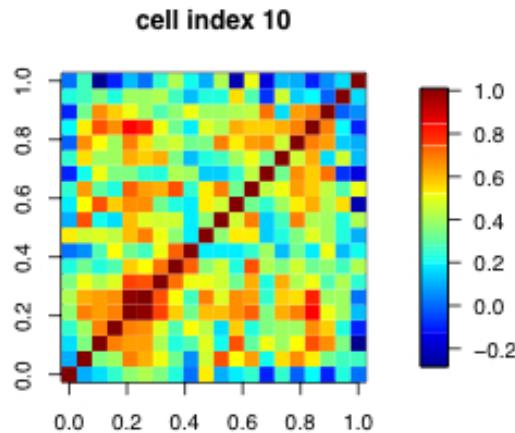
- R defaults: ketchup and mustard
- RColorBrewer: good maps based on work of visualization expert (matrix2png uses them too)
- Red-black-green you often see in papers is **bad choice** for colour blind individuals.
- Blue-black-yellow is better
- Greyscale: loss of dynamic range, but cheaper to publish!
- My favourite: matrix2png "black body"
- Humans can't really tell the difference between a 8 and 16 colour scale.

Divergent vs. sequential maps

Colours pass through black at (in this case) zero. Yellow="Below the mean", Blue="Above the mean" for the row. This can be the right thing to do, especially if your original data are naturally "symmetric" around zero (or some other value). Otherwise it might just be confusing.

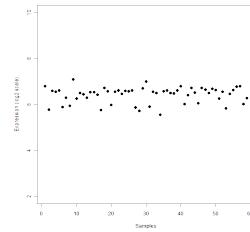
Colours go from light to dark. Darker colours mean "more".

A confusing heat map (at least for me; one can get used to anything)



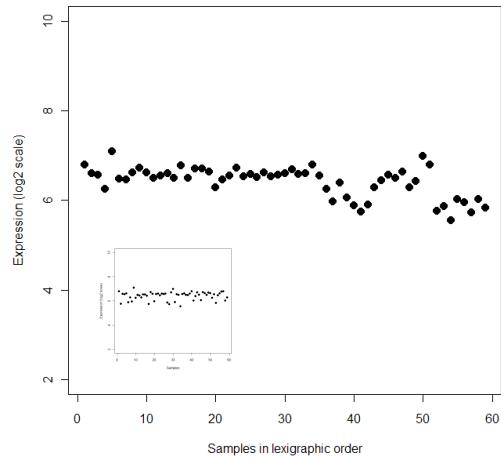
Slicing and dicing

- I mean: Viewing data arranged or grouped by factors of interest
- In my example, I have multiple “visits” for each “patient”
 - Any interesting trends we can spot?
 - Are intra-subject samples more similar than inter-subject samples?



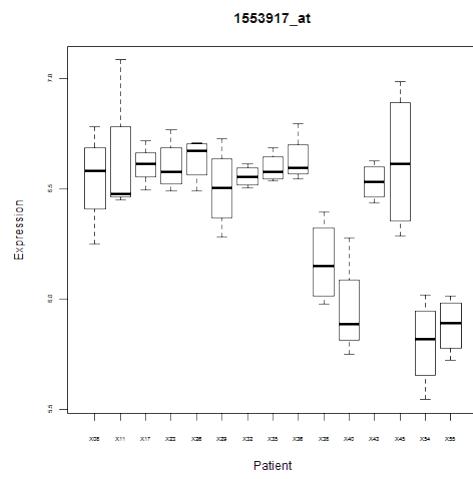
Playing around with one gene

- I sorted the samples by their names
- Is this structure meaningful? Is it interesting?
- Maybe we need more information.
- (And some statistics)



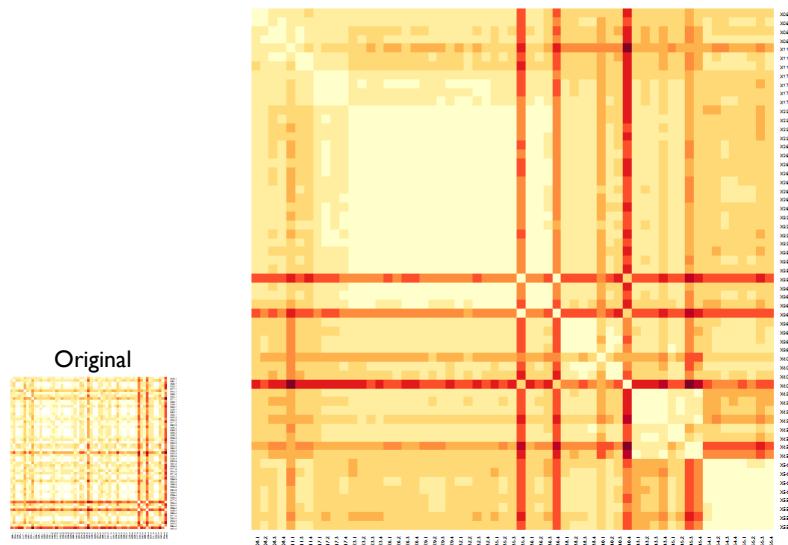
Explicitly group by patient

- We're developing the concept that we can look for "factors" that "explain" some of the variation.
- Is this interesting? We don't know, but it suggests that accounting for Patient differences might be reasonable.
- Need to get more comprehensive view ...



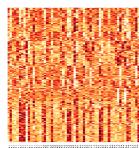
(Boxplots not ideal for this, but easy)

Arrange data by Patient



View of (part of) the data organized by Patient

(Again, row order
is not random
either)



Summary

- Sanity checks
- Use graphs to look at the data, slicing and dicing – build up a library of techniques that work for your data
- Additional exploratory techniques will be discussed later in the course
 - Clustering
 - PCA
- Also: missing values