

Experimental Design – Part II

Yi-Ju Li, Ph.D.

Department of Biostatistics & Bioinformatics
Duke University Medical Center

July 22, 2015

- Technical variability in RNA-Seq
- Experimental Designs in RNA-Seq
- Factor of population structure

Main References:

1. Marioni et al. Genome Research 2008, 18:1509-17
2. McIntyre et al. BMC Genomics 2011, 12:293
3. Auer et al. Genetics 2010, 185:405-416

Principles of Design of Experiments

Four commonly considered principles in the design of experiment¹

- **Representativeness:** Are the experimental units used in the experiment sufficient to represent the conclusion to be made?
- **Randomization:** Help to avoid unknown bias.
- **Replication:** Increase the precision of the data.
- **Error control or blocking:** Help to reduce known bias (e.g. batch effect).

Experiment needs to be comparative.

¹Fisher R.A., 1935 The Design of Experiment.

Section 1

Variability in RNA-Seq data

Steps of a RNA-Seq experiment²

1. RNA is isolated from cells, fragmented at random positions, and copied into complementary DNA (cDNA)
2. Fragments meeting a certain specified size (e.g. 200 – 300 bp) are retained for PCR
3. Sequencing
4. Sequence alignment to generate sequence reads at each position
5. **Data:** Counts of sequence reads or **digital gene expression (DGE)**
6. **Types of reads:** junction reads, exonic reads, polyA reads

²Auer et al. **Genetics** 2010.

Types of variability applying to any experiments

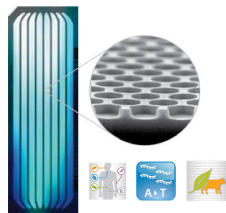
- Technical variability
- Biological variability
 - ▶ Variability between experimental units (samples)
 - ▶ Variability between factors of interest (treatment groups)
 - ▶ Biological variability is not affected by technical variability.

These sources of variability need to be considered in the experimental design.

Types of technical variability:

- between sequencing platforms
- between library construction
- between flow cells (different runs)
- between lanes

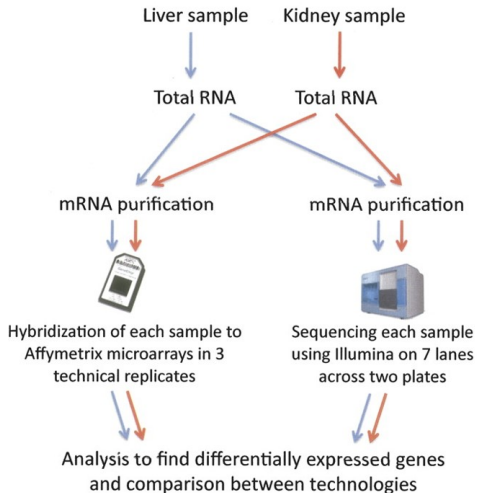
Flow cells: A glass slide with 1, 2, or 8 separate lanes (Illumina RNA-Seq)



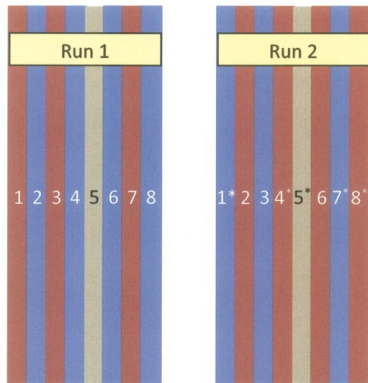
Example I³

- **Objective:** Assess the technical reproducibility of Illumina RNA-Seq
 - ▶ Comparison between platforms
 - ▶ Evaluate technical variability of RNA-Seq
- **Outline of the experiment**
 - ▶ Two sequencing platforms: Illumina RNA-Seq (8 lanes) and Affymetrix microarray
 - ▶ Two samples: a liver and a kidney samples
 - ▶ Two cDNA concentration (3pM and 1.5pM)
 - ▶ One lane for a control sample.
 - ▶ Each sample were sequenced 7 times total in two flow-cell runs.

³Marioni et al. 2008.

A**B**

Illumina study design



Kidney
Liver

* Sequenced at a concentration of 1.5 pM

Example I: Marioni et al.⁴:

● Outline of the experiment

- ▶ Two sequencing platforms: Illumina RNA-Seq (8 lanes) and Affymetrix microarray
- ▶ Two samples: a liver and a kidney samples
- ▶ One lane for a control sample.
- ▶ Each sample were sequenced 7 times total in two flow-cell runs.

● What can they compare with this design?

⁴Marioni et al. 2008.

Example I: Marioni et al.⁴:

● Outline of the experiment

- ▶ Two sequencing platforms: Illumina RNA-Seq (8 lanes) and Affymetrix microarray
- ▶ Two samples: a liver and a kidney samples
- ▶ One lane for a control sample.
- ▶ Each sample were sequenced 7 times total in two flow-cell runs.

● What can they compare with this design?

- ▶ **Platform differences:** Two methods for gene expression

⁴Marioni et al. 2008.

Example I: Marioni et al.⁴:

● Outline of the experiment

- ▶ Two sequencing platforms: Illumina RNA-Seq (8 lanes) and Affymetrix microarray
- ▶ Two samples: a liver and a kidney samples
- ▶ One lane for a control sample.
- ▶ Each sample were sequenced 7 times total in two flow-cell runs.

● What can they compare with this design?

- ▶ **Platform differences:** Two methods for gene expression
- ▶ **Technical variability:** Same sample sequenced in different lanes, two separate runs

⁴Marioni et al. 2008.

Example I: Marioni et al.⁴:

● Outline of the experiment

- ▶ Two sequencing platforms: Illumina RNA-Seq (8 lanes) and Affymetrix microarray
- ▶ Two samples: a liver and a kidney samples
- ▶ One lane for a control sample.
- ▶ Each sample were sequenced 7 times total in two flow-cell runs.

● What can they compare with this design?

- ▶ **Platform differences:** Two methods for gene expression
- ▶ **Technical variability:** Same sample sequenced in different lanes, two separate runs
- ▶ **Effect of cDNA concentration:** two concentrations (3 and 1.5pM)

⁴Marioni et al. 2008.

Example I: Marioni et al.⁴:

● Outline of the experiment

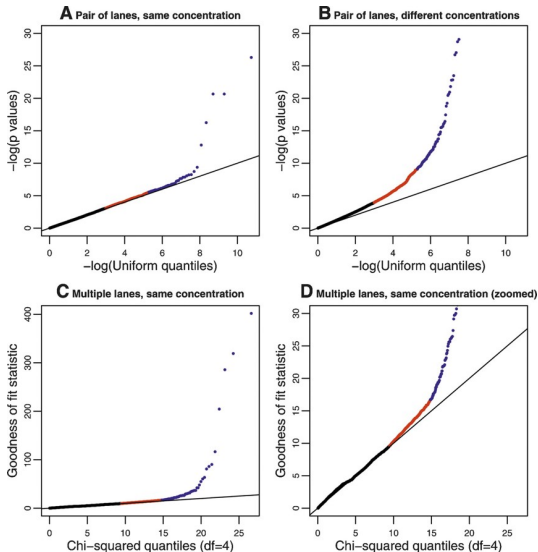
- ▶ Two sequencing platforms: Illumina RNA-Seq (8 lanes) and Affymetrix microarray
- ▶ Two samples: a liver and a kidney samples
- ▶ One lane for a control sample.
- ▶ Each sample were sequenced 7 times total in two flow-cell runs.

● What can they compare with this design?

- ▶ **Platform differences:** Two methods for gene expression
- ▶ **Technical variability:** Same sample sequenced in different lanes, two separate runs
- ▶ **Effect of cDNA concentration:** two concentrations (3 and 1.5pM)
- ▶ **Differential expression between liver and kidney tissues:** Two tissue samples

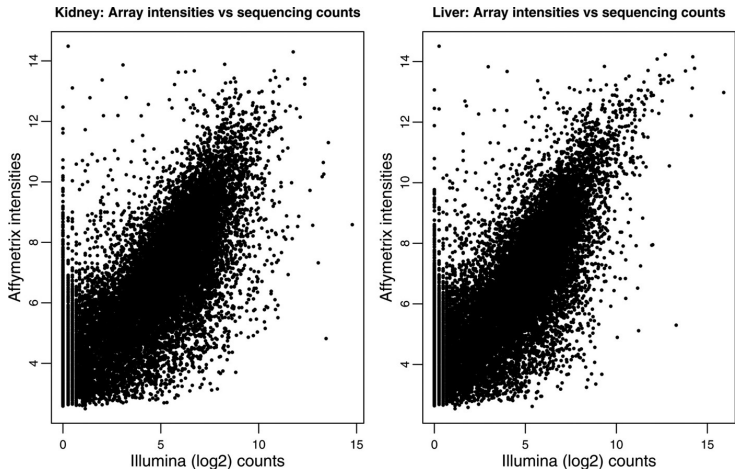
⁴Marioni et al. 2008.

Plots for assessing lane effect



A: Same sample, same concentration; **B:** Sample sample, different concentration; **C&D:** Goodness-of-fit for Poisson distribution – kidney samples

Comparison between platforms Comparing counts from Illumina sequencing with normalized intensities from the array, for kidney (left) and liver (right).



Spearman correlation = 0.73 for liver, 0.75 for kidney

Conclusion and issues in the design

- Summary of Marioni et al. 2008
 - ▶ Illumina RNA-Seq is replicable and has advantage over microarray
 - ▶ **Lane effect is small. (??)**
 - ▶ Larger difference between runs (**Batch effect**)
 - ▶ larger difference between cDNA concentration
 - ▶ Suggested that it is OK to run one sample per lane
- **Issues in the Design: No replicates, one sample only**
 - ▶ Is it sufficient to use one sample per tissue type to conclude low lane effect?
 - ▶ Can we partition biological variation (e.g. liver vs. kidney) from technical variation?

Example II: McIntyre et al.⁵ *RNA-seq: technical variability and sampling*

- **Objective:**

- ▶ Does technical variability exist?
- ▶ Is the impact of technical variability the same for all levels of coverage?

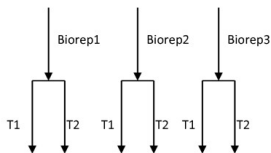
- **Experiments:**

1. Three independent samples (*D. melanogaster* female), two technical replicates per sample, run on two lanes of a Solex/Illumina flow cell.
2. Three independent samples (*D. simulans* male), two technical replicates per sample, run on two lanes of a flow cell.
3. One sample (*D. melanogaster* cell lines), 5 replicates, run on 5 lanes of a flow cell.

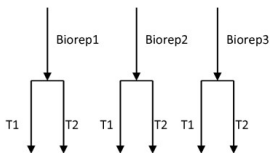
All are 36 base-paired end. The relationship of lanes for cell lines (same or independent flow cells) is unknown.

⁵McIntyre et al. *BMC Genomics* 2011.

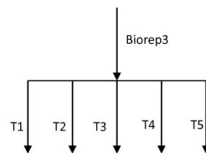
D. melanogaster



D. simulans



Cell line "c167"



D. melanogaster and *D. simulans* are single library run on multiple lanes.
D. melanogaster c167 cell lines are not exactly the same library run.

Variation of sequence reads

Partial data from Table 1 (McIntyre et al. 2011)

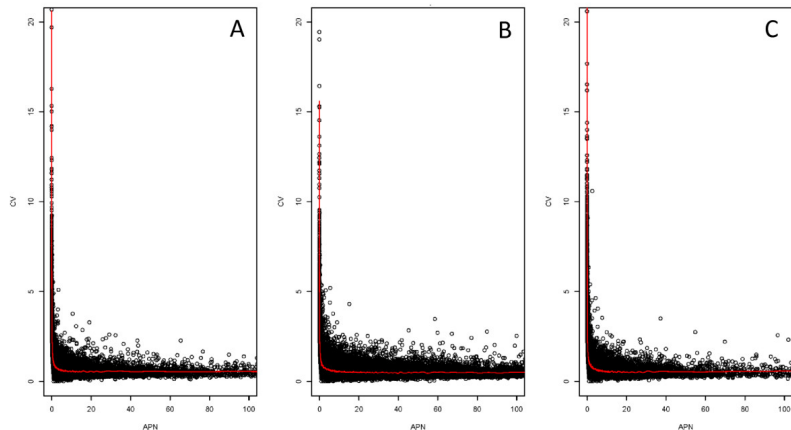
Table 1 Mappable reads per lane in each of the three experiments

Experiment	BR	TR	Mappable Reads	Exons detected	Exons with an average coverage of more than 5 reads per nucleotide	Contigs present in all samples of each experiment
c167	1	1	5888686	39156	13432	19248
c167	1	2	5951769	39202	13517	19248
c167	1	3	7146461	39954	15684	19248
c167	1	4	7544117	40201	16355	19248
c167	1	5	7377032	40120	16089	19248
D. sim.	1	1	5174398	45878	14517	20339
D. sim.	1	2	4979485	45808	13912	20339
D. sim.	2	1	27595266	51701	35303	20339
D. sim.	2	2	28691914	51857	35942	20339
D. sim.	3	1	27601233	51834	34968	20339
D. sim.	3	2	27748704	51822	35008	20339

Data variation seen between technical replicates and between biological replicates.

Visualization of the data

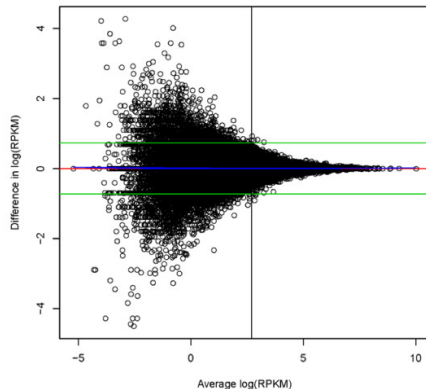
Coefficient of variation (CV) vs. average depth per nucleotide (APN); APN: within each lane, average number of reads per exon



Lower coverage has higher variation

Visualization of the data

MA plot: Minus vs. average, Bland and Altman plot



Green line: One standard deviation region

Low expression level has higher disagreement

- **Technical variation exists:**

- ▶ Mappable reads per lane varies among the technical replicates
- ▶ Inconsistent detection of exons between technical replicates: The number of exons detected increases with the number of mappable reads.
- ▶ Agreement between technical replicates varies: kappa ranges from 0.63-0.81
- ▶ Higher variability for those with low coverage (< 5 reads per nucleotide) or low expression level
- ▶ Random sampling of total RNA reads (e.g. 0.0013% of 30 millions reads) may contribute to the variability.

- **Biological variation is larger than technical variation.**

- **Suggestion:** Inclusion of technical replicates is as important as biological replicates. Multiplexing design can eliminate the lane effect for a small experiment.

Section 2

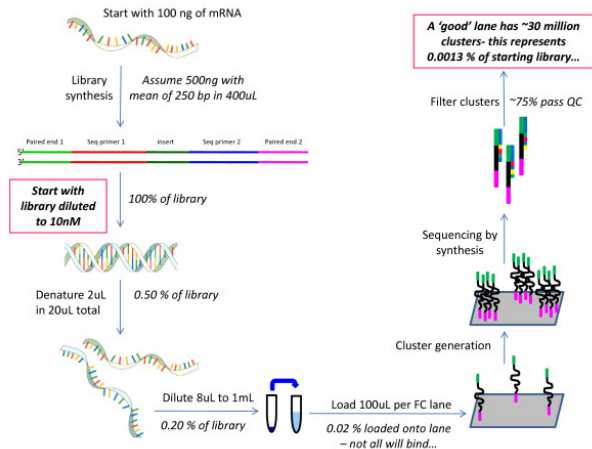
Designs for RNA-Seq

- **Three levels of sampling in RNA-Seq:**
 - ▶ Subject sampling
 - ▶ RNA sampling
 - ▶ Fragments sampling
- **Unreplicated data:**
 - ▶ Mostly from observational studies
 - ▶ No biological replicates
 - ▶ One sample per treatment group
- **Problem:** This design can investigate only the differences derived from RNA and fragment-level sampling, but not from subject sampling.

More on sampling in RNA-Seq

- **Subject sampling:** Subjects (e.g. organisms or individuals) are ideally drawn from a large population to which the results can be generalized.
- **RNA sampling:** occurs during the experimental procedure when RNA is isolated from the cell(s).
- **Fragment sampling:** Only certain fragmented RNAs are retained for amplification. The sequencing reads do not represent 100% of the fragments loaded into a flow cell resulted in fragment sampling.

More on RNA and fragment primer sampling

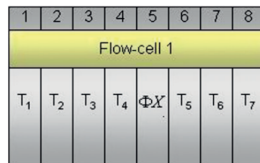


Library concentration 10nM with $400\mu\text{L} \rightarrow 4\text{pM} \rightarrow \frac{4}{10^{12}} \times 6.02 \times 10^{23} = 2.408 \times 10^{12}$ total molecules \rightarrow

$\frac{30,000,000}{2.408 \times 10^{12}} = 0.0013\%$ of molecules to be analyzed.

Outline of experiment:

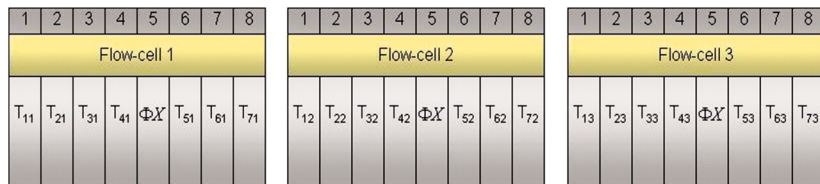
- mRNA isolated from subjects within different treatment group (T_1, \dots, T_7).
- a ΦX genomic sample is loaded to lane 5 as a control
- ΦX can be used to recalibrate the quality score of sequencing reads from other lane.



Problems:

- Lack of knowledge about biological variation
- Unable to estimate within treatment variation leading to no basis for inference of between treatment effect.
- Results are specific to the subjects in the study and can't be generalized.

Replicated data: Multiple flow-cell design



- **Exp Design:** One sample per treatment group, two additional biological replicates. T_{ij} for i^{th} treatment group and j^{th} replicate. $i = 1, \dots, 7$ and $j = 1 - 3$.
- **Factor of consideration:** treatment effect (τ_{ik})

$$(\text{Dependent variable})_{ijk} = \alpha_k + \tau_{ik} + \epsilon_{ijk}$$

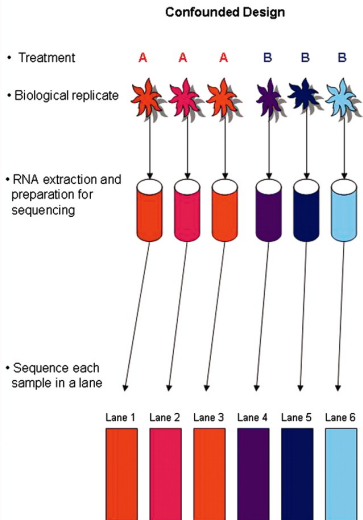
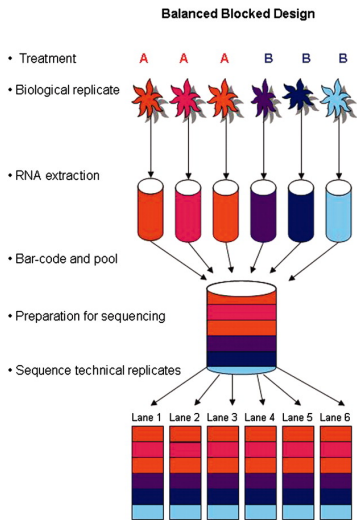
- **Problem:** Cannot separate treatment effect from technical effect since biological replicates are run in different flow-cells.

- **Objective:** To control two sources of technical variation: batch effect and lane effect.
- **Multiplexing:** All samples are pooled to be run within the same lane.
 - ▶ Take the advantage of bar coding of RNA fragments.
 - ▶ To keep the same sequence depth, divide the amplification product to run in multiple lanes
 - ▶ If # of lanes = # of samples, it produces the same sequence depth as running one sample per lane.
 - ▶ Each lane has the same set of samples – **eliminate the lane effect**

Example:

- Three biological replicates per treatment ($j = 1, \dots, 3$)
- treatment group (A and B) ($i = 1, \dots, 2$)
- RNA are bar-coded and pooled
- Divide the pool to six equal subset to run on 6 lanes (six technical replicates, $t = 1, \dots, 6$)
- Single flow cell run

Balanced Block design - I



Analysis model for BBD I

- **Dependent variable:** DGE measures, defined by the distribution you assumed for the sequence reads. For example,
 - ▶ Auer et al. assumed $y_{ijk} \sim \text{Poisson}(\mu_{ijk})$.
 - ▶ DESeq2 uses Negative Binomial model.

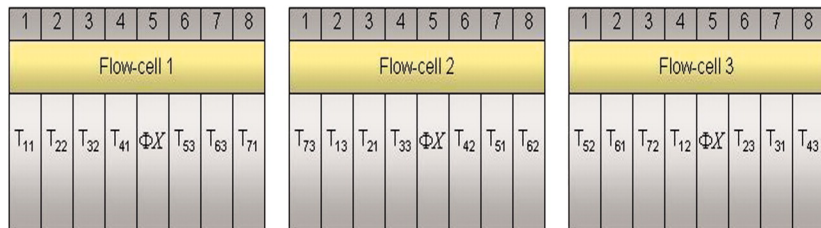
In Auer et al, $y_{ijk} = \sum_t y_{ijkt}$, where i for treatment, j for sample, k for gene, and t for the 6 technical replicates

- **Factors considered in the GLM:** treatment effect (τ_{ik}).

$$(\text{Dependent variable})_{ijk} = \alpha_k + \tau_{ik}$$

- No lane effect was included in this model as they considered lane effects were balanced across treatment groups.
- No batch effect in this case since it is only one flow-cell run.

Balanced block design II-without multiplexing



- A design that can run one sample per lane but also has good randomization of samples within each flow-cell.
- Three biological replicates within seven treatment groups. T_{ij} , where $i = 1, \dots, 7$ for treatment groups and $j = 1, \dots, 3$ for samples.
- **Two block effects:** flow cells and lanes.

- **Dependent variable:** Same as before, but it is coded to indicate treatment (i), flow-cell (f), lane (l), and gene (k).
- **Factors to consider:** treatment effect (τ_{ik}), flow-cell effect (ν_{fk}), and lane effect (ω_{lk}).

$$\text{(Dependent variable)}_{ijflk} = \alpha_k + \tau_{ik} + \nu_{fk} + \omega_{lk} + \epsilon_{ijflk}$$

ϵ_{ijflk} is the error term.

Summary for Balanced block design

- The feature of unique bar-code for RNA fragments in RNA-Seq makes blocking design possible.
- Can control batch and lane effects
- Multiplex design illustrated here requires the number of unique bar-codes equal or greater than the samples in each lane.

Balanced incomplete block design (BIBD)

Assume:

- number of treatment (I)
- number of biological replicates per treatment (J)
- number of unique barcodes (s) that can be included in one lane
- number of lanes available for sequencing (L)

If the number of unique bar codes (s) in one lane is less than the number of treatments ($s < I$), balanced block design is impossible.

- For a given number of treatment groups (I), sample per treatment group (J), unique barcodes (s), and number of available lanes (L), the total number of technical replicates (T) in BIBD is $T = \frac{sL}{IJ}$.
- **Example of BIBD:**
 - ▶ Assume 3 treatment group ($i = 3$), one subject per treatment group ($j = 1$), two unique barcodes ($s = 2$), and three available lanes ($L = 3$).
 - ▶ The total number of technical replicates is $T = \frac{2 \times 3}{3 \times 1} = 2$.

1	2	3
T_{111}	T_{211}	T_{311}
T_{212}	T_{312}	T_{112}

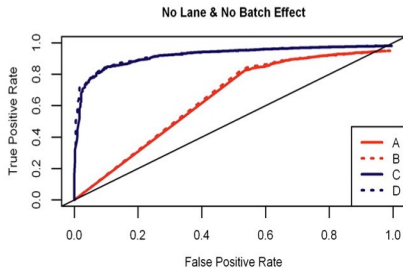
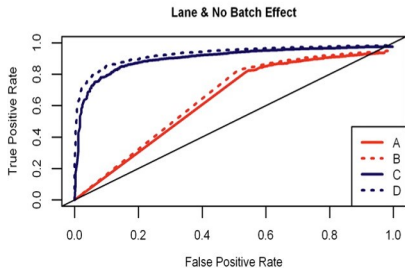
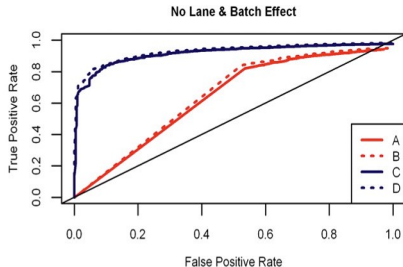
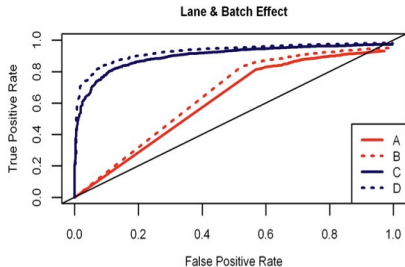
- T_{ijk} is for treatment i , subject j , and technical replicates t .
- For Illumina, a total of 12 unique barcodes can be used in one lane. Therefore, 96 samples can be multiplexed in one flow-cell run.

Performance comparison between designs by simulation studies

A		B		C						D					
1	2	1	2	1	2	3	4	5	6	1	2	3	4	5	6
T_{11}	T_{21}	T_{111}	T_{112}	T_{11}	T_{12}	T_{13}	T_{21}	T_{22}	T_{23}	T_{111}	T_{112}	T_{113}	T_{114}	T_{115}	T_{116}
		T_{211}	T_{212}							T_{121}	T_{122}	T_{123}	T_{124}	T_{125}	T_{126}
										T_{131}	T_{132}	T_{133}	T_{134}	T_{135}	T_{136}
										T_{211}	T_{212}	T_{213}	T_{214}	T_{215}	T_{216}
										T_{221}	T_{222}	T_{223}	T_{224}	T_{225}	T_{226}
										T_{231}	T_{232}	T_{233}	T_{234}	T_{235}	T_{236}

T_{ijk} : i for treatment, j for sample, k for technical replicates.

A: unreplicated data; **B**: no biological replicates, two technical replicates (BBD without biological replicates); **C**: no technical replicates (unblocked design); **D**: BBD with biological and technical replicates.



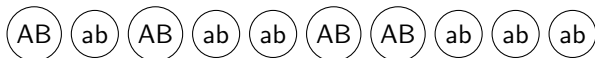
C&D always perform better than A&B. When simulation included lane and/or batch effects, **D (balanced block design)** performed better than **C (unblocked design)**.

Section 3

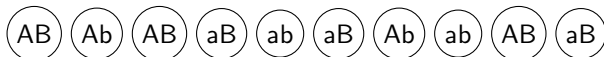
Factor of Population structure

- **Concept of allelic association:**

- ▶ Alleles A and B at two loci are **associated** if the event that a gamete carries A is not independent of the event that the gamete carries allele B.



- ▶ Alleles A and B are **not associated** if they occur in the gametes randomly.

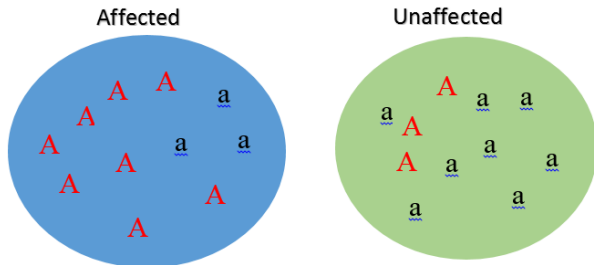


- ▶ Allelic association is population specific.

Background of disease genetic association

- **Disease/allele association**

- ▶ Look for the association between markers and disease phenotype.
- ▶ Allele **A** occurs more frequent in affected than unaffected subjects.



- ▶ Assume the marker is in allelic association with the causal allele.

Population admixture can cause allelic association

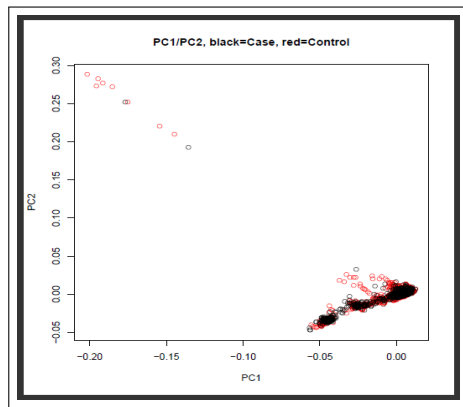
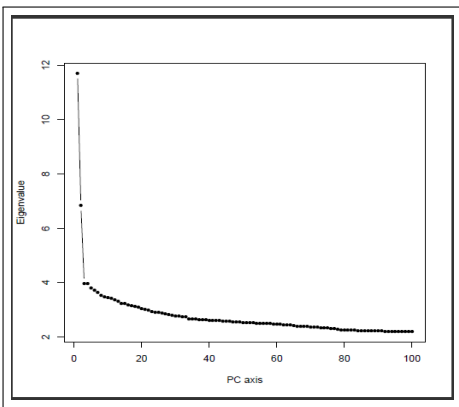
- Two (or more) mixing populations can lead to associations created due to the differences in allele frequencies in the mixing populations
 - ▶ **Population A:** **A** allele is VERY common; Disease allele (**D**) occurs randomly with **A** or **a** alleles.
 - ▶ **Population B:** **a** allele is VERY common; Disease allele (**D**) does not exist.
 - ▶ **Admixture population:** Assume equal mixed of populations A and B (allele **A** with frequency of 0.5), we will observed an association between A and D.
- Population structure will lead to false positive results for genetic association studies.
- How about gene expression studies? Does population structure have an effect on gene expression?

How to adjust population structure?

- **Family-based design**
- **Genomic control:** Estimate degree of population stratification by typing 20-60 unlinked markers on same cases and controls used for studying candidate gene association (Devlin and Roeder 1999)
- **Structure:** Alternative method based on explicit modeling of population structure (Pritchard and Rosenberg (1999))
- **Eigenstrat:** Use principal components (PCs) analysis to explicitly model ancestry differences between cases and controls. (Price et al. 2006)

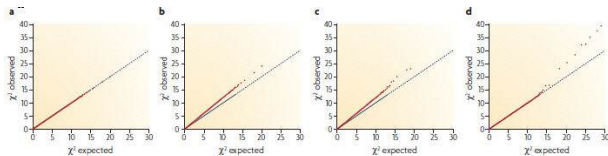
The idea of using principal components to adjust for population structure can also be applied to gene expression data.

Plots from EigenSoft



Make a decision: (1) incorporate PC(s) in the model; (2) exclude outliers.

Visualization: Q-Q plots



Quantile-Quantile Plots compare observed distribution of test statistics to that expected under the null hypothesis of no association (McCarthy et al. 2008).

- No association, observed = expected
- Probably mostly population substructure, deviations across distribution
- Possible true associations, but also population substructure
- True association, deviations at the highest end of the distribution

Useful R code: <http://www.broadinstitute.org/diabetes/scandinav/figures.html>

Example for gene expression vs. ethnicity groups

- Hicks et al.⁶ reported gene expression levels significantly differ between ethnic groups.
 - ▶ **Study Design:** Compare gene expression between four ethnic populations (Whites, Blacks, Hispanics, and Asians) using samples from B-Precursor acute lymphoblastic leukemia (B-ALL) patients.
 - ▶ Affymetrix microarray
 - ▶ Significant differential expression genes were found between ethnic populations.
- For gene expression studies in human, one should also pay attention on the effect of ethnicity (or population structure).
- Use QQ plot to examine whether the population structure could be a factor.

⁶Hicks et al. **Cancer Informatics 2013.**

Summary

- The classical principals of experimental design still apply to RNA-Seq
- Technical variation exists and should be taken into account.
 - ▶ Lane effect, batch effect
- RNA-Seq data consist of variation from subject sampling, RNA sampling, RNA fragment sampling
- Multiplexing in NGS allow us to implement randomization and blocking.
- Take advantages of visualization tools (e.g. scatter plots, MA plots, QQ plots) to learn your data.
- When you deal with human data for genetic study, make sure examining the effect of population structure.