Epigenomics: Some Statistical Applications



Rafael A. Irizarry Department of Biostatistics

John Hopkins Bloomberg School of Public Health

Acknowledgements

- Tom Albert, Nimblegen
- Benilton Carvalho, JHU Biostatistics
- Andy Feinberg Lab, JHU Medicine
- Todd Richmond, Nimblegen
- Hao Wu, JHU Biostatistics
- Jean Wu, Brown Biostat
- Vasan Yegnasubramanian, JHU Oncology

Outline

- Quick Introduction to Epigenetics
- Introduction to Methylation
- Overview of competing technologies
- Review: Expression arrays lessons
- Comparison
- Role of statisticians

Genetics: the alphabet of life

 Letters of DNA sequence carry the information



Epigenetics



(Klug & Cummings 1997)

 $(3.4x10^{-10} \text{ meters/bp}) \times (6x10^9 \text{ bp/genome}) = \sim 2 \text{ meters/genome}$

Radius of the nucleus is ~ 10 μ M !!!

Klug and Cummings, 1997



[(6 x 10⁹ bp/genome) / (195 bp/nucleosome)] = ~ 30.8 x 10⁶ nucleosomes/genome ~ 5 % of nuclear volume

Nucleosome, Solenoid Model of Chromatin and Chromosome





http://www.albany.edu/~achm110/solenoidchriomatin.html

Epigenetics: the grammar of life



DNA methylation







Observed to expected = Pr(CG) / { Pr(C) Pr(G) }

DNA methylation can lead to silencing of gene expression



Robertson and Wolffe, Nat Rev Genet, 2000

ENCODE Track





Expression Array Lessons

Normalization



Probe effect



Intensity = Background + Probe Effect x Quantity x Error

Sequence effect for BG

Wu et al. (2004) JASA 99(468) 909



Back to Methylation

High throughput of course....

Densities for three methods



HCT116 lots of methylation **DKO** very little methylation

Hunh?



MeDIP (like ChIPchip)



Some Data



Problem: Not specific



HELP: Two enzymesCuts at CCGGCuts at C^MCGG



HELP after PCR



HELP



HELP





The Problem



Obsered to expeced = Pr(CG) / { Pr(C) Pr(G) }

Proportion of neighboring CpG also methylated/not methylated



McRBC on Tiling array



ROC now



ENCODE Track



Problems for Statisticians

- Background Correction + Normalization
- Probability Model for Segments
- Use these to from null and alternative models... we need power!
- Use these to create bump finding algorithms
- Adapt to high-throughput sequencing

Supplemental Slides

McRBC: One enzyme

Input

Cuts at A^mCG or G^mCG



McRBC after Gel



McRBC after Gel



McRBC



McRBC after **GEL**



McRBC after GEL





