# COS 424 Homework 2

**Vladimir Feinberg**
Princeton University
`vyf@princeton.edu`

## Abstract

abstract

## 1 Introduction and Related Work

Imputation is a general statistical process involving replacing unknown values with a best guess based on a given context. Formally, we consider some vectors $\mathbf{x}_i \in S^n$ for some set $S$. If we take these $\mathbf{x}_i$ to be sampled from a distribution over $S^n$, we can make some predictions about the posterior of a $\mathbf{x}$ given some limited observation into a set of $\{\mathbf{x}_i\}$.

In our case, $S$ is a tuple over various biological information: chromosome number, genomic offset, strand type (3'-5' or 5'-3'), and methylation, the proportion of sites in a laboratory sample that had an additional methyl group attached to a cytosine molecule. The methylation count in our data set may be missing information.

Predicting whether the proportion of methylated sites from a sample is at least half has been done with over 90% accuracy by relying on methylation levels of nearby sites, DNA sequence properties such as ENCODE, and genomic position to the extent that the site is co-located (in the binary sense) with a DNA sequence for a particular protein or a sequence for a cis-regulartory element (CRE) [3]. Other studies have also found co-localization between CREs and CpGs (sites where the methylation occurs) [4].

## 2 Data

We focus our analysis mainly on chromosome 1, which has samples from 379551 sites. There are 33 samples for each site, each taken from an expensive WGBS procedure, this method is able to measure about 91% of sites [2]. Our test sample only has about 2% of the sites available from a cheaper procedure: methylation microarrays [3].

We are to use the limited information to impute the missing values in the sample. Because we can only use WGBS as the true value, there are some sites that we don't know the correct value for, even in testing. The imputation procedure may *not* be specified to only impute the sites we are testing for because this is a facet of our testing procedure. The methods constructed attempt to impute the entire genome. We are only able to provide an estimate for our testing error because we can only observe some of the true methylation values with the more expensive technique.

The sites available at test time are consistent - the technology used to assay the methylation values samples sites consistently - bare collaborative filtering solves a larger problem where the new sample may be missing values at different sites, so it cannot take advantage of learning patterns particular to the sites consistently tested by the methylation microarrays [1].

# 3 Methods

# 4 Results

# 5 Discussion and Conclusion

# References

[1] ILLUMINA. Ifinium HumanMethylation450 BeadChip. Data Sheet, March 2012.

[2] LAIRD, P. W. Principles and challenges of genome-wide dna methylation analysis. *Nature Reviews Genetics 11*, 3 (2010), 191–203.

[3] ZHANG, W., SPECTOR, T. D., DELOUKAS, P., BELL, J. T., AND ENGELHARDT, B. E. Predicting genome-wide dna methylation using methylation marks, genomic position, and dna regulatory elements. *Genome Biol 16* (2015), 14.

[4] ZILLER, M. J., GU, H., MÜLLER, F., DONAGHEY, J., TSAI, L. T.-Y., KOHLBACHER, O., DE JAGER, P. L., ROSEN, E. D., BENNETT, D. A., BERNSTEIN, B. E., ET AL. Charting a dynamic dna methylation landscape of the human genome. *Nature 500*, 7463 (2013), 477–481.