**The University of Chicago**

**Department of Statistics**

**Ph.D. Seminar**

**OMAR DE LA CRUZ**

Department of Statistics

The University of Chicago

## Geometric Approaches in the Analysis of Genetic Data

**WEDNESDAY, July 30, 2008 at 3:00 PM**
**110 Eckhart Hall, 5734 S. University Avenue**

## ABSTRACT

We propose a method for detecting cell-cycle-regulated genes by studying the geometric structure of gene expression data obtained by assaying individual cells from a growing population: under reasonable assumptions, the data points will cluster around a closed curve that represents the ideal evolution of gene expression during the cell cycle. We describe a statistical model as well as a general strategy for fitting the data, divided in two main steps: first, using robust local orthogonal regression to obtain an initial estimate of the curve; and second, improving the likelihood of the estimate by direct search.

The method sketched above leads in a natural way to the notion of geometric learning. We will present some theoretical contributions to this field, by showing that the robust local orthogonal regression approach corresponds, when the sample size goes to infinity, to the detection of what we call modal ridges in the density of the distribution induced by the model. However, these theoretical results are not directly applicable to the case of cycle-regulated gene expression, since they depend on large sample sizes.

We also discuss how to integrate geometric learning with more traditional statistical procedures, like regression, as a way to obtain more easily interpretable inferences. This makes it possible to establish which coordinates contribute significantly to the geometric structure (e.g., which genes are cycle-regulated). Also, it makes it possible to adjust for the influence of the geometric structure, in order to more accurately measure other properties of the individuals; for example, in studies of gene expression in single cells, adjusting for the cell cycle allows a more accurate estimation of other characteristics, like membership in cell subpopulations or treatment effects.

Information about building access for persons with disabilities may be obtained in advance by calling Jewanna Carver at 773.834-5169 or by email (carver@galton.uchicago.edu).