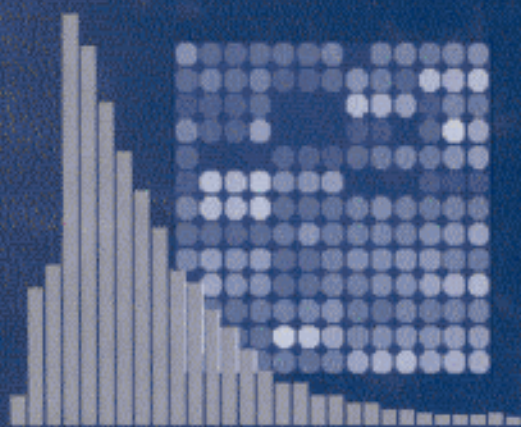# EXPLORATION AND ANALYSIS OF DNA MICROARRAY AND PROTEIN ARRAY DATA

## DHAMMIKA AMARATUNGA
## JAVIER CABRERA

# Contents