Daniel O. Stram

# Design, Analysis, and Interpretation of Genome-Wide Association Scans

2 Springer

Daniel O. Stram

# Design, Analysis, and Interpretation of Genome-Wide Association Scans

Springer

# Statistics for Biology and Health

Daniel O. Stram

# Design, Analysis, and Interpretation of Genome-Wide Association Scans

Daniel O. Stram
Department of Preventive Medicine
University of Southern California
   Keck School of Medicine
Los Angeles, CA, USA

*To Pavlova, Alex, and Douglas*

# Acknowledgments

# Contents