# Social Sequence Analysis

*Methods and Applications*

BENJAMIN CORNWELL
*Cornell University*

# Contents

# Figures

# Tables