# Advances in Financial Machine Learning

MARCOS LOPEZ DE PRADO

**WILEY**

# Contents

**Labeling**

**Sample Weights**