# Contents

## II    Inference      289

## 7    Inference algorithms: an overview      291

## 8    State-space inference      301

# III    Prediction    537

# 14 Prediction models: an overview    539

# IV  Generation        767

# 21 Generative models: an overview        769

# 22 Variational autoencoders        783