# Contents

# II   Inference      339

# IV   Generation       763

## 20 Generative models: an overview       765

## 21 Variational autoencoders       781

# V   Discovery      915

## VI    Action      1091