# Contents

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

## II   Inference   299

## 7   Inference algorithms: an overview   301

## 8   Message passing inference   313

## 12 Sequential Monte Carlo inference    485

# III   Prediction      507

## 13 Predictive models: an overview      509

## 14 Generalized linear models      525

## IV Generation     713

## 19 Generative models: an overview     715

# V   Discovery    865

# 26 Discovery methods: an overview    867

# 27 Latent factor models    869

# 28 State-space models    917

# VI    Decision making     1051

# 33 Multi-step decision problems     1053

# 34 Reinforcement learning     1079

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47