# Contents

**4   Probabilistic graphical models    95**

## II Inference 293

## 7 Inference algorithms: an overview 295

## 8 Message passing inference 307

# 9   Variational inference     365

xvi

9.3.5 Low-rank Gaussian VI 387
9.3.6 Example: Full-rank vs diagonal GVI on 1d linear regression 388
9.3.7 Sparse Gaussian VI 390
9.3.8 Automatic differentiation VI 392
9.3.9 Non-Gaussian reparameterized VI 394
9.3.10 Amortized inference 396
9.3.11 Online variational inference 398
9.4 More accurate variational posteriors 401
9.4.1 Structured mean field 401
9.4.2 Hierarchical (auxiliary variable) posteriors 401
9.4.3 Normalizing flow posteriors 402
9.4.4 Implicit posteriors 403
9.4.5 Combining VI with MCMC inference 403
9.5 Lower bounds 404
9.5.1 Multi-sample ELBO (IWAE bound) 404
9.5.2 The thermodynamic variational objective (TVO) 405
9.6 Upper bounds 406
9.6.1 Minimizing the $\chi$-divergence upper bound 407
9.6.2 Minimizing the evidence upper bound 408
9.7 Expectation propagation (EP) 408
9.7.1 Minimizing forwards vs reverse KL 408
9.7.2 EP as generalized ADF 410
9.7.3 Algorithm 410
9.7.4 Example 412
9.7.5 Optimization issues 412
9.7.6 Power EP and $\alpha$-divergence 413
9.7.7 Stochastic EP 413
9.7.8 Applications 414

**10 Monte Carlo inference** **415**
10.1 Introduction 415
10.2 Monte Carlo integration 415
10.2.1 Example: estimating $\pi$ by Monte Carlo integration 416
10.2.2 Accuracy of Monte Carlo integration 416
10.3 Generating random samples from simple distributions 418
10.3.1 Sampling using the inverse cdf 418
10.3.2 Sampling from a Gaussian (Box-Muller method) 419
10.4 Rejection sampling 419
10.4.1 Basic idea 420
10.4.2 Example 421
10.4.3 Adaptive rejection sampling 421
10.4.4 Rejection sampling in high dimensions 422
10.5 Importance sampling 422
10.5.1 Direct importance sampling 423

# III   Prediction   517

## IV    Generation    723

## 19 Generative models: an overview    725

# VI   Decision making     1081

## 33 Multi-step decision problems     1083