

MACHINE LEARNING BOOK

$E_{\rho_0}[\chi(p)]$ -- folu) hgpovir.



$$w_{L+1} = w_0 - \eta \nabla w_L$$

HEISENBERG

1. Math Basic

- 1.1 Gaussian Distribution - MLE — — — — — — — P₁
- 1.2 Gaussian Distribution - MLE (unbiased vs. biased) — — — — — — — P₂
- 1.3 Gaussian Distribution from pdf perspective — — — — — — — P₂
- 1.4 Gaussian Distribution - Limitation — — — — — — — P₃
- 1.5 Gaussian Distribution - Marginal & Conditional probability — — — — — — — P₃
- 1.6 Gaussian Distribution - Joint Probability — — — — — — — P₄

2. Linear Regression

- 2.1 Least Square Method — — — — — — — — — — — — — P₅
- 2.2 Least Square Method - from probabilistic perspective — — — — — — — P₅
- 2.3 Regularization - Ridge Regression - frequentist — — — — — — — P₆
- 2.4 Regularization - Ridge Regression - Bayesians — — — — — — — P₇

3. Linear classification

- 3.1 Linear classification Background — — — — — — — — — — — — — P₇
- 3.2 Linear Classification perceptron Algorithm — — — — — — — — — — — — — P₈
- 3.3 Fisher - Linear Discriminant Analysis — — — — — — — — — — — — — P₈
- 3.4 Logistic Regression — — — — — — — — — — — — — P₁₀
- 3.5 Gaussian Discriminant Analysis — — — — — — — — — — — — — P₁₀
- 3.6 Learning of the parameter about GDA — — — — — — — — — — — — — P₁₁
- 3.7 Naive Bayes classifier — — — — — — — — — — — — — P₁₂

4. Dimensionality Reduction

- 4.1 Background — — — — — — — — — — — — — P₁₃
- 4.2 Sample Mean & Variance Matrix — — — — — — — — — — — — — P₁₃
- 4.3 PCA - maximum variance perspective — — — — — — — — — — — — — P₁₄
- 4.4 PCA - minimum error perspective — — — — — — — — — — — — — P₁₄
- 4.5 PCA - SVD perspective — — — — — — — — — — — — — P₁₅
- 4.6 PPCA - probabilistic PCA — — — — — — — — — — — — — P₁₅

5. Support Vector Machine

| | | |
|--|-----------|-----|
| 5.1 hard - margin SVM - model definition | — — — — — | P17 |
| 5.2 hard - margin SVM - model Solution | — — — — — | P17 |
| 5.3 Soft - margin SVM | — — — — — | P19 |
| 5.4 constraint Optimization - weak Duality | — — — — — | P19 |
| 5.5 Geometric interpretation of duality | — — — — — | P20 |
| 5.6 slater Condition - Duality problem | — — — — — | P20 |
| 5.7 karush - Kuhn - Tucker Condition | — — — — — | P20 |

6. kernel Method

| | | |
|--------------------------------|-----------|-----|
| 6.1 kernel Method Introduction | — — — — — | P21 |
| 6.2 Positive Definite Kernel | — — — — — | P22 |

7. Exponential Family Distribution

| | | |
|---------------------------------|-----------|-----|
| 7.1 Background | — — — — — | P23 |
| 7.2 Gaussian Distribution | — — — — — | P23 |
| 7.3 Log Partition Function | — — — — — | P23 |
| 7.4 MLE & Sufficient Statistics | — — — — — | P24 |
| 7.5 Maximum Entropy perspective | — — — — — | P24 |

8. Probabilistic Graphical Medels

| | | |
|---|-----------|-----|
| 8.1 PGM Background | — — — — — | P25 |
| 8.2 Bayesian Network - condition independence | — — — — — | P26 |
| 8.3 Bayesian Network - D- Separation | — — — — — | P26 |
| 8.4 Bayesian Network - example | — — — — — | P26 |
| 8.5 Markov Random Field - Representation - Conditional Independence | — — — — — | P27 |
| 8.6 Markov Random Field - Representation - Factorization | — — — — — | P27 |
| 8.7 Inference - Introduction | — — — — — | P27 |
| 8.8 Inference - Variable Elimination | — — — — — | P28 |
| 8.9 Inference - Belief propagation | — — — — — | P28 |

| | | | |
|------|-----------------------------------|-----------|-----|
| 8.10 | Inference - Max Product Algorithm | - - - - - | P29 |
| 8.11 | Supplement - Moral Graph | - - - - - | P30 |
| 8.12 | Supplement - Factor Graph | - - - - - | P30 |
| 9. | EM Algorithm | | |
| 9.1 | EM Introduction | - - - - - | P30 |
| 9.2 | ELBO + kL Divergence | - - - - - | P31 |
| 9.3 | ELBO + Jensen's Inequality | - - - - - | P32 |
| 9.4 | EM Review | - - - - - | P32 |
| 9.5 | Generalized EM | - - - - - | P32 |
| 10. | Gaussian Mixture Model | | |
| 10.1 | Introduction | - - - - - | P33 |
| 10.2 | MLE Not Applicable | - - - - - | P33 |
| 10.3 | EM - E-Step | - - - - - | P33 |
| 10.4 | EM - M-Step | - - - - - | P34 |
| 11. | Variational Inference | | |
| 11.1 | Introduction | - - - - - | P35 |
| 11.2 | Formula Deduction | - - - - - | P35 |
| 11.3 | Review | - - - - - | P36 |
| 11.4 | SGVI | - - - - - | P36 |
| 12. | Markov Chain & Monte Carlo | | |
| 12.1 | Sampling Method Introduction | - - - - - | P38 |
| 12.2 | Markov Chain | - - - - - | P38 |
| 12.3 | MH Algorithm | - - - - - | P39 |
| 12.4 | Gibbs Algorithm | - - - - - | P39 |
| 12.5 | Review | - - - - - | P40 |
| 12.6 | Stationary Distribution | - - - - - | P40 |

| | | |
|---|-----------|-----|
| 12.7 Problem & Thinking | - - - - - | P41 |
| 13. Hidden Markov Model | | |
| 13.1 Background | - - - - - | P41 |
| 13.2 Forward Algorithm | - - - - - | P42 |
| 13.3 Backward Algorithm | - - - - - | D42 |
| 13.4 Baum Welch Algorithm-EM | - - - - - | P43 |
| 13.5 Viterbi Algorithm | - - - - - | P44 |
| 13.6 Summary | - - - - - | P44 |
| 14. Linear Dynamic System-Kalman Filter | | |
| 14.1 Background | - - - - - | P45 |
| 14.2 Filtering | - - - - - | P45 |
| 15. Particle Filter | | |
| 15.1 Background | - - - - - | P46 |
| 15.2 Importance Sampling & SIS | - - - - - | P46 |
| 15.3 Particle Filter↔Re-Sampling | - - - - - | P47 |
| 15.4 SIR Filter | - - - - - | P47 |
| 16. Conditional Random Field | | |
| 16.1 Background | - - - - - | P48 |
| 16.2 HMM vs. MEMM | - - - - - | P48 |
| 16.3 MEMM vs. CRF | - - - - - | P49 |
| 16.4 Pdf - Parameter perspective | - - - - - | P49 |
| 16.5 Pdf - Vector perspective | - - - - - | P50 |
| 16.6 Model Problems | - - - - - | P50 |
| 16.7 Marginal probability | - - - - - | P50 |

| | | |
|--|-----------|-----|
| 16.8 Learning | - - - - - | P52 |
| 17. Gaussian Network | | |
| 17.1 Background | - - - - - | P53 |
| 17.2 Gaussian Bayesian Network | - - - - - | P53 |
| 17.3 Gaussian MRF | - - - - - | P54 |
| 18. Bayesian Linear Regression | | |
| 18.1 Background | - - - - - | P55 |
| 18.2 Inference | - - - - - | P55 |
| 18.3 Prediction | - - - - - | P56 |
| 19. Gauss Process | | |
| 19.1 Background | - - - - - | P56 |
| 19.2 Weight Space Perspective | - - - - - | P57 |
| 19.3 From Weight Space to Function Space | - - - - - | P59 |
| 19.4 Function Space Perspective | - - - - - | P59 |
| 20. Restricted Boltzmann Machine | | |
| 20.1 Background | - - - - - | P60 |
| 20.2 Representation | - - - - - | P61 |
| 20.3 Representation - review | - - - - - | P61 |
| 20.4 Inference - Posterior Distribution | - - - - - | P62 |
| 20.5 Inference - Marginal Distribution | - - - - - | P63 |
| 21. Spectral Clustering | | |
| 21.1 Background | - - - - - | P63 |
| 21.2 Model Introduction | - - - - - | P63 |
| 21.3 Matrix - Indicator vector / Diagonal Matrix | - - - - - | P64 |
| 22. Feedforward Neural Network | | |
| 22.1 From ML to DL | - - - - - | P65 |

| | | |
|---|-----------|-----|
| 22.2 From perceptron to deep learning | — — — — — | P66 |
| 22.3 Non-Linear Transformation | — — — — — | P66 |
| 23. Confronting Partition Function | | |
| 23.1 The Log-Likelihood gradient | — — — — — | P67 |
| 23.2 Stochastic Maximum Likelihood | — — — — — | P67 |
| 23.3 What is Contrastive Divergence | — — — — — | P68 |
| 23.4 The name of Contrastive Divergence | — — — — — | P68 |
| 23.5 Log-Likelihood Gradient of Energy-Based Model/RBM Learning | — — | P68 |
| 23.6 Log-Likelihood Gradient of RBM | — — — — — | P69 |
| 23.7 CD-K Algorithm for RBM | — — — — — | P70 |
| 24. Approximate Inference | | |
| 24.1 Introduction | — — — — — | P70 |
| 24.2 Inference is Optimization | — — — — — | P70 |
| 25. Sigmoid Belief Network | | |
| 25.1 Introduction | — — — — — | P71 |
| 25.2 Gradient of Log-Likelihood | — — — — — | P71 |
| 25.3 Wake-Sleep Algorithm-Introduction | — — — — — | P72 |
| 25.4 Wake-Sleep Algorithm - KL Divergence | — — — — — | P72 |
| 26. Deep Belief Network | | |
| 26.1 Introduction | — — — — — | P73 |
| 26.2 Stacking RBM | — — — — — | P73 |
| 26.3 Extra Layers improve ELBO | — — — — — | P73 |
| 26.4 Pre-Training | — — — — — | P74 |

27. Boltzmann Machine

| | | |
|------------------------------------|---------------------|-----|
| 27.1 Introduction | — — — — — — — — — — | P74 |
| 27.2 Gradient of Log-Likelihood | — — — — — — — — — — | P75 |
| 27.3 Gradient Ascend Based on MCMC | — — — — — — — — — — | P75 |
| 27.4 Conditional Probability | — — — — — — — — — — | P76 |
| 27.5 Variational Inference | — — — — — — — — — — | P77 |

28. Deep Boltzmann Machine

| | | |
|------------------------------|---------------------|-----|
| 28.1 Introduction | — — — — — — — — — — | P78 |
| 28.2 Pre-Training | — — — — — — — — — — | P78 |
| 28.3 Double Counting Problem | — — — — — — — — — — | P79 |
| 28.4 Pre-Training Continue | — — — — — — — — — — | P79 |

29. Generative Model

| | | |
|--|---------------------|-----|
| 29.1 Supervised Learning VS. Unsupervised Learning | — — — — — — — — — — | P80 |
| 29.2 Presentation & Inference & Learning | — — — — — — — — — — | P80 |
| 29.3 Model Classification | — — — — — — — — — — | P81 |
| 29.4 Probabilistic Graph vs. Neural Network | — — — — — — — — — — | P81 |
| 29.5 Reparameterization Trick | — — — — — — — — — — | P81 |

30. Generative Adversarial Network

| | | |
|---------------------------|---------------------|-----|
| 30.1 One Example | — — — — — — — — — — | P82 |
| 30.2 Model representation | — — — — — — — — — — | P82 |
| 30.3 Global Optimality | — — — — — — — — — — | P82 |

31. Variational Autoencoder

| | | |
|---------------------|---------------------|-----|
| 31.1 Representation | — — — — — — — — — — | P83 |
| 31.2 Learning | — — — — — — — — — — | P83 |

32. Normalizing Flow

32.1 Model Representation — — — — — P84

1-1 Gaussian Distribution-MLE
 Data: $\mathbf{x} = (x_1, \dots, x_N) = \begin{pmatrix} x_1 \\ \vdots \\ x_N \end{pmatrix}_{N \times p}$ Machine Learning
 $x_i \in \mathbb{R}^p$, $x_i \sim N(\mu, \Sigma)$, $\theta = (\mu, \Sigma)$ MLE (Maximum Likelihood Estimation)

MLE: $\theta_{MLE} = \arg\max_{\theta} P(x|\theta)$. $P(x) = \frac{1}{(2\pi)^{\frac{N}{2}} |\Sigma|^{\frac{1}{2}}} \exp(-\frac{1}{2}\mathbf{x}^T \Sigma^{-1} (\mathbf{x} - \mu))$.
 令 $\mu = (\mu, \Sigma) = (\mu, \sigma^2)$ 以 2 维为例 $x_i \sim N(\mu, \sigma^2)$
 $\log P(x|\theta) = \log \prod_{i=1}^N P(x_i|\theta) = \sum_{i=1}^N \log P(x_i|\theta)$
 $= \sum_{i=1}^N \left[\underbrace{\log \frac{1}{\sqrt{2\pi\sigma^2}}}_{\text{无关}} + \underbrace{\log \frac{1}{\sigma}}_{\text{无关}} - \frac{(x_i - \mu)^2}{2\sigma^2} \right]$

$\hat{\mu}_{MLE} = \arg\max_{\mu} \log P(x|\theta) = \arg\max_{\mu} \sum_{i=1}^N -\frac{(x_i - \mu)^2}{2\sigma^2} = \arg\min_{\mu} \sum_{i=1}^N (x_i - \mu)^2$
 $\frac{\partial}{\partial \mu} \sum_{i=1}^N (x_i - \mu)^2 \approx \sum_{i=1}^N 2(x_i - \mu)(-1)$,
 令 $2 \sum_{i=1}^N (x_i - \mu)(-1) = 0$
 $\sum_{i=1}^N (x_i - \mu) = 0$
 $\sum_{i=1}^N x_i - N\mu = 0$
 $\boxed{\hat{\mu}_{MLE} = \frac{1}{N} \cdot \sum_{i=1}^N x_i}$ 无偏

$E[\hat{\mu}_{MLE}] = \frac{1}{N} \sum_{i=1}^N E[x_i] = \frac{1}{N} \sum_{i=1}^N \mu = \frac{1}{N} \cdot N \cdot \mu = \mu$

$\hat{\sigma}^2_{MLE} = \arg\max_{\sigma^2} \log P(x|\theta) = \arg\max_{\sigma^2} \sum_{i=1}^N \left[\log \frac{1}{\sigma} - \frac{(x_i - \mu)^2}{2\sigma^2} \right] = \arg\max_{\sigma^2} \sum_{i=1}^N \underbrace{\left[-\log \sigma - \frac{(x_i - \mu)^2}{2\sigma^2} \right]}_{d(\sigma)}$
 $\frac{\partial d}{\partial \sigma} = \sum_{i=1}^N \left[-\frac{1}{\sigma} - \frac{1}{2} (x_i - \mu)^2 \cdot (-2) \frac{1}{\sigma^3} \right] = \sum_{i=1}^N \left[\frac{1}{\sigma} + \frac{1}{\sigma^3} (x_i - \mu)^2 \right]$
 令 $\frac{\partial d}{\partial \sigma} = 0$, 则:
 $\sum_{i=1}^N \left[-\frac{1}{\sigma} + \frac{1}{\sigma^3} (x_i - \mu)^2 \right] = 0$.
 $\sum_{i=1}^N \left[-\sigma^2 + (x_i - \mu)^2 \right] = 0$.
 $\sum_{i=1}^N \sigma^2 = \sum_{i=1}^N (x_i - \mu)^2$
 $\boxed{\hat{\sigma}^2_{MLE} = \frac{1}{N} \cdot \sum_{i=1}^N (x_i - \mu)^2}$ 有偏估计: $E[\hat{\sigma}^2_{MLE}] = \frac{N-1}{N} \sigma^2$
 无偏估计: $\hat{\sigma}^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \hat{\mu}_{MLE})^2$

$$E[\hat{\sigma}^2_{MLE}] = \sigma^2, \text{ 相等则无偏 } \text{Var}[M_{MLE}] = \text{Var}\left[\frac{1}{N} \sum_{i=1}^N x_i\right] = \frac{1}{N^2} \sum_{i=1}^N \text{Var}(x_i) = \frac{1}{N^2} \sum_{i=1}^N \sigma^2 = \frac{1}{N} \sigma^2$$

$$\begin{aligned} \hat{\sigma}^2_{MLE} &= \frac{1}{N} \sum_{i=1}^N (x_i - M_{MLE})^2 = \frac{1}{N} \sum_{i=1}^N (x_i^2 - 2x_i M_{MLE} + M_{MLE}^2) = \frac{1}{N} \sum_{i=1}^N x_i^2 - \frac{1}{N} \sum_{i=1}^N 2x_i M_{MLE} + \frac{1}{N} \sum_{i=1}^N M_{MLE}^2 \\ &= \frac{1}{N} \sum_{i=1}^N x_i^2 - 2M_{MLE} + M_{MLE}^2 = \frac{1}{N} \sum_{i=1}^N x_i^2 - M_{MLE}^2 \end{aligned}$$

则有: $E[\hat{\sigma}^2_{MLE}] = E\left[\frac{1}{N} \cdot \sum_{i=1}^N x_i^2 - M_{MLE}^2\right] = E\left[\left(\frac{1}{N} \sum_{i=1}^N x_i^2 - \mu^2\right) - (M_{MLE}^2 - \mu^2)\right]$

$$\begin{aligned} &= E\left(\frac{1}{N} \sum_{i=1}^N x_i^2 - \mu^2\right) - E(M_{MLE}^2 - \mu^2) = \sigma^2 - \frac{1}{N} \sigma^2 = \frac{N-1}{N} \sigma^2 \\ &= E\left(\frac{1}{N} \sum_{i=1}^N (x_i^2 - \mu^2)\right) - E(M_{MLE}) - E(\mu^2) \\ &\quad \frac{1}{N} \sum_{i=1}^N E(x_i^2 - \mu^2) \quad E(M_{MLE}) - \mu^2 \\ &\quad \frac{1}{N} \sum_{i=1}^N (E(x_i^2) - \mu^2) \quad E(M_{MLE}) - E(M_{MLE}) \rightarrow E(w) = E(M_{MLE}) \text{ 无偏} \\ &\quad \underbrace{\text{Var}(x_i)}_{E(x_i) = \mu} \quad \underbrace{\text{Var}(M_{MLE})}_{E(M_{MLE})} \\ &\quad \underbrace{\frac{1}{N} \sum_{i=1}^N \sigma^2}_{\sigma^2} \quad \underbrace{\frac{1}{N} \sigma^2}_{\sigma^2} \end{aligned}$$

1.3 Gaussian Distribution from pdf perspective

1.2 Gaussian Distribution-MLE (unbiased vs biased)

极大似然估计中, $\hat{\sigma}^2$ 有偏差

差: $E(x - \bar{x})$ 而非 $E|x - \mu|$

高维情况下: 高斯分布

$$x \sim N(\mu, \Sigma) = \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right), x \in \mathbb{R}^p, \text{ random vector} \rightarrow \text{r.v.}$$

$$x = \begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix}; \mu = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_p \end{pmatrix}; \Sigma = \begin{pmatrix} \sigma_{11} & \sigma_{12} & \cdots & \sigma_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{p1} & \sigma_{p2} & \cdots & \sigma_{pp} \end{pmatrix}$$

正定的 (以下讨论)

半正定 (一般情况) 且对称

协方差矩阵

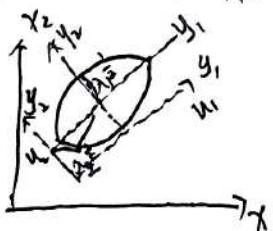
$(x - \mu)^T \Sigma^{-1} (x - \mu)$: 马氏距离 (x 与 μ 之间) 表示点与一个分布之间的距离。它是一种有效计算两个样本集相似度的方法。与欧氏距离不同的是, 它考虑到各种特征之间的联系

$\Sigma = I$, 马氏距离 = 欧氏距离

$$\text{e.g. } z_1 = \begin{pmatrix} z_{11} \\ z_{12} \end{pmatrix}; z_2 = \begin{pmatrix} z_{21} \\ z_{22} \end{pmatrix} \text{ 有 } (z_1 - z_2)^T \Sigma^{-1} (z_1 - z_2) = (z_{11} - z_{21}, z_{12} - z_{22}) \cdot \Sigma^{-1} \begin{pmatrix} z_{11} - z_{21} \\ z_{12} - z_{22} \end{pmatrix}$$

$$\text{对角化: } \sum = \sum_{i=1}^p \frac{y_i^2}{\lambda_i}, \text{ 令 } \lambda = 2,$$

$$\text{则 } p=2, \Delta = \frac{y_1^2}{\lambda_1} + \frac{y_2^2}{\lambda_2}, \text{ 但 } \Delta \neq 1, \text{ 這是椭圆}$$



$\Delta = \Sigma \rightarrow \Sigma \text{ 对称}$
取不同的
相当不同的
等高线

$$\Delta = \Sigma = U \Lambda U^T, VV^T = V^T V = I, \Lambda = \text{diag}(\lambda_i) \text{ } i=1, \dots, p, \text{ 令 } V = (v_1, \dots, v_p) \text{ } \text{op}$$

$$\begin{aligned} \Delta &= (U_1, \dots, U_p) \begin{pmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \lambda_p \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_p \end{pmatrix} = (U_1 \lambda_1, U_2 \lambda_2, \dots, U_p \lambda_p) \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_p \end{pmatrix} = \sum_{i=1}^p U_i \lambda_i u_i^T \\ \Sigma^{-1} &= (V^T \Lambda V)^{-1} = (V^T)^T \Lambda^{-1} (V^T) = (V^T)^{-1} \Lambda^{-1} (V^T) = V \Lambda^{-1} V^T \end{aligned}$$

$$\Delta = \sum_{i=1}^p \frac{1}{\lambda_i} U_i \frac{1}{\lambda_i} U_i^T = \sum_{i=1}^p U_i \frac{1}{\lambda_i} U_i^T \xrightarrow{\text{diag}} \text{diag}\left(\frac{1}{\lambda_i}\right), i=1, \dots, p$$

$$\Delta = (x - \mu)^T \Sigma^{-1} (x - \mu) = (x - \mu)^T \sum_{i=1}^p U_i \frac{1}{\lambda_i} U_i^T \cdot (x - \mu) \quad P2$$

$$= \sum_{i=1}^p (x - \mu)^T u_i \cdot \frac{1}{\lambda_i} u_i^T (x - \mu)$$

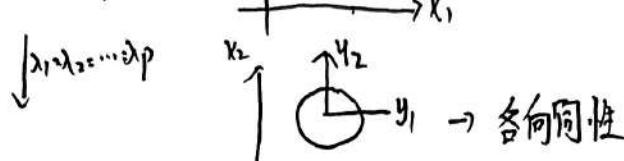
$$y = \begin{pmatrix} y_1 \\ y_p \end{pmatrix}; y_i = (x - \mu)^T u_i, y_i \text{ 是 } i \text{ 维的}$$

1.4 Gaussian Distribution - ~~from p. 10~~ ^{Limitation} $= \sum_{i=1}^p y_i \frac{1}{\lambda_i} y_i^T = \boxed{\sum_{i=1}^p \frac{y_i^2}{\lambda_i}}$

局限性：

①. $\Sigma_{pp} \rightarrow$ 参数个数: $C_p + p = \frac{p(p+1)}{2} + p = \frac{p(p+3)}{2} \rightarrow O(p^2)$

$\Sigma \rightarrow$ 对角矩阵 $(\lambda_1 \lambda_2 \dots \lambda_p) \rightarrow$



Gaussian Distribution - Marginal & Conditional probability

factor analysis (因子分析) \rightarrow 对角矩阵

P - PCA \rightarrow 各向同性

1.4.5 $x = \begin{pmatrix} x_a \\ x_b \end{pmatrix} \rightarrow m+n=p$ $u = \begin{pmatrix} u_a \\ u_b \end{pmatrix}$ $\Sigma = \begin{pmatrix} \sum_{aa} & \sum_{ab} \\ \sum_{ba} & \sum_{bb} \end{pmatrix}$ \rightarrow 对角 Gaussian 分布定理 q维 qp pm

求 $P(x_a), P(x_b|x_a); P(x_b), P(x_a|x_b)$ 对称

$$x_a = \underbrace{(I_m 0)}_A \underbrace{\begin{pmatrix} x_a \\ x_b \end{pmatrix}}_x + \underbrace{B}_0.$$

$$E[x_a] = Ax + B = (I_m 0) \begin{pmatrix} u_a \\ u_b \end{pmatrix} + 0 = u_a$$

$$\text{Var}[x_a] = (I_m 0) \left(\begin{pmatrix} \sum_{aa} & \sum_{ab} \\ \sum_{ba} & \sum_{bb} \end{pmatrix} \right) \left(\begin{pmatrix} I_m \\ 0 \end{pmatrix} \right) = (\sum_{aa} \sum_{ab}) \begin{pmatrix} I_m \\ 0 \end{pmatrix} = \sum_{aa}$$

$\therefore x_a \sim N(u_a, \sum_{aa})$

$$x_b|x_a, x_b-a = x_b - \sum_{ba} \sum_{aa}^{-1} x_a \rightarrow *$$

$$\text{令 } M_b-a = M_b - \sum_{ba} \sum_{aa}^{-1} M_a \quad \text{结果定义}$$

$$\sum_{bb-a} = \sum_{bb} - \sum_{ba} \sum_{aa}^{-1} \sum_{ab} \rightarrow \sum_{aa}^{-1} \text{ Schur complementary}$$

$$x_{b-a} = (- \sum_{ba} \sum_{aa}^{-1} I) \begin{pmatrix} x_a \\ x_b \end{pmatrix}$$

$$E[x_{b-a}] = (- \sum_{ba} \sum_{aa}^{-1} I) \begin{pmatrix} u_a \\ u_b \end{pmatrix} + \frac{B}{0} = M_b - \sum_{ba} \sum_{aa}^{-1} M_a = M_b-a.$$

$$\text{Var}[x_{b-a}] = (- \sum_{ba} \sum_{aa}^{-1} I) \left(\begin{pmatrix} \sum_{aa} & \sum_{ab} \\ \sum_{ba} & \sum_{bb} \end{pmatrix} \right) \left(\begin{pmatrix} -\sum_{aa}^{-1} \sum_{ba} \\ I \end{pmatrix} \right) = \begin{pmatrix} -\sum_{ba} + \sum_{ba} & -\sum_{ba} \sum_{aa}^{-1} \sum_{ab} + \sum_{bb} \\ -\sum_{ba} + \sum_{ba} & -\sum_{ba} \sum_{aa}^{-1} \sum_{ab} \end{pmatrix} \left(\begin{pmatrix} -\sum_{aa}^{-1} \sum_{ba} \\ I \end{pmatrix} \right) = (0 - \sum_{ba} \sum_{aa}^{-1} \sum_{ab} + \sum_{bb}) \left(\begin{pmatrix} -\sum_{aa}^{-1} \sum_{ba} \\ I \end{pmatrix} \right) = -\sum_{ba} \sum_{aa}^{-1} \sum_{ab} + \sum_{bb} = \sum_{bb-a}$$

$x_{b-a} \sim N(M_b-a, \sum_{bb-a})$

$$x_b | x_a$$

$$x_b = x_b \cdot a + \sum_{ba} \sum_{aa} x_a = \underbrace{x_b \cdot a}_{\text{固定映射}} + \underbrace{\sum_{ba} \sum_{aa} x_a}_B$$

$$E[x_b | x_a] = \underbrace{x_b \cdot a}_{\text{固定映射}} + \sum_{ba} \sum_{aa} x_a$$

$$\text{Var}[x_b | x_a] = \text{Var}[x_{b \cdot a}] = \sum_{bb} a^2.$$

$$x_b | x_a \sim N(\mu_{b \cdot a} + \sum_{ba} \sum_{aa} x_a, \sum_{bb} a)$$

$$x_a \sim N(\mu_a, \sum_{aa})$$

$$x_b \sim N(\mu_b, \sum_{bb})$$

$$x_a | x_b \sim N(\mu_{a \cdot b} + \sum_{ab} \sum_{bb} x_b, \sum_{aa} b)$$

x_b 与 x_a 相关的那个样子，体现了 x_a 与 x_b 的一致性关系，对给定的 x_a ，有固定映射的 x_b ，所以此时 $E(x_b)$ 为给定条件的 $E(x_b)$ ，亦即 $E(x_b | x_a)$

$$x_b \cdot a \sim N(\mu_{b \cdot a}, \sum_{bb} a)$$

在 $x_b | x_a$ 中， x_b, x_a 相互独立。

$$E[x_b | x_a] = E[x_b] = E[x_b \cdot a + \sum_{ba} \sum_{aa} x_a]$$

$$= E(x_{b \cdot a}) + \sum_{ba} \sum_{aa} x_a.$$

$$\because A = I$$

$$\therefore E[x_b | x_a] = \mu_{b \cdot a} + \sum_{ba} \sum_{aa} x_a.$$

$$\boxed{\text{Var}[x_b | x_a] = \text{Var}[x_b] = \text{Var}[x_{b \cdot a} + \sum_{ba} \sum_{aa} x_a]}$$

$$\text{常数} \text{Var}[b] = 0, x | \text{Var}[x_b | x_a] = \text{Var}[x_{b \cdot a}] = \sum_{bb} a$$

1.5.6 Gaussian Distribution - Joint Probability

在此基础上，依据Gauss分布定理，求：

$$\text{已知: } p(x) = N(x | \mu, A^{-1}) \quad \begin{matrix} \text{precision matrix} \\ \text{精度} \end{matrix}$$

$$p(y|x) = N(y | Ax + b, L^{-1}) \quad \begin{matrix} \text{covariance matrix} \\ \text{协方差} \end{matrix}$$

$$\text{求: } p(y), p(x|y)$$

$$\text{解: } (1) y = Ax + b \quad x, y, \varepsilon, \text{随机变量, } A, b \text{ 系数}$$

$$\varepsilon \sim N(0, L^{-1}) \quad \varepsilon \perp x, \varepsilon \text{与 } x \text{ 独立.}$$

$$E[y] = E[Ax + b + \varepsilon] = E[Ax] + E[\varepsilon] = A\mu + b$$

$$\text{Var}[y] = \text{Var}[Ax + b + \varepsilon] = \text{Var}[Ax] + \text{Var}[b] + \text{Var}[\varepsilon] = A \cdot A^{-1} A^T + L^{-1}$$

$$\therefore y \sim N(A\mu + b, L^{-1} + A \cdot A^{-1} A^T) \quad \Delta = A^{-1} A^T \text{ 对称.}$$

$$(2) z = (y) \sim N\left(\underbrace{\mu}_{E[z]}, \underbrace{\begin{bmatrix} A^{-1} & 0 \\ 0 & L^{-1} + A \cdot A^{-1} A^T \end{bmatrix}}_{\text{Var}[z]}\right)$$

$$p(z) = N\left(\underbrace{\mu}_{\mu}, \underbrace{\begin{pmatrix} A^{-1} & A^T \\ A \cdot A^{-1} & L^{-1} + A \cdot A^{-1} A^T \end{pmatrix}}_{\Sigma}\right).$$

$$E[x|y] = \mu + A^{-1} A^T (L^{-1} + A \cdot A^{-1} A^T)^{-1} (y - A\mu - b)$$

$$\text{Var}[x|y] = A^{-1} - A^{-1} A^T (L^{-1} + A \cdot A^{-1} A^T)^{-1} A \cdot A^{-1}$$

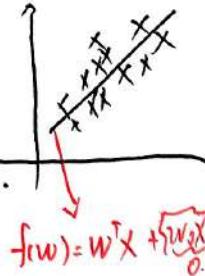
Gauss分布: $x \sim N(\mu, \Sigma), y = Ax + b$.

$$y \sim N(A\mu + b, A\Sigma A^T)$$

$$\begin{aligned} A \cdot \text{cov}(x, y) &= E[(x - E[x]) \cdot (y - E[y])^T] \\ &= E[(x - \mu) \cdot (y - A\mu - b)^T] \\ &= E[(x - \mu)(Ax + b + \varepsilon - A\mu - b)^T] \\ &= E[(x - \mu)(Ax - A\mu + \varepsilon)^T] \\ &= E[(x - \mu)(Ax - A\mu)^T + (x - \mu)\varepsilon^T] \\ &= E[(x - \mu)(Ax - A\mu)^T] + E[(x - \mu)\varepsilon^T] \\ &= E[(x - \mu)(Ax - A\mu)^T] \quad \varepsilon \perp \varepsilon \\ &= E[(x - \mu)(x - \mu)^T] \quad x - \mu \perp \varepsilon \\ &= E[(x - \mu)(x - \mu)^T] \quad E(x - \mu) \cdot E(\varepsilon) \\ &= E[(x - \mu)(x - \mu)^T] \cdot A^T \\ &= \text{Var}[x] \cdot A^T = A^{-1} \cdot A^T \end{aligned}$$

2.1.2.1 - Least Square Method

2 线性回归 (Linear Regression)



$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$$

$$x_i \in \mathbb{R}^P, y_i \in \mathbb{R}, i=1, 2, \dots, N$$

$$x = (x_1, x_2, \dots, x_N)^T = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{pmatrix} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1P} \\ x_{21} & x_{22} & \dots & x_{2P} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \dots & x_{NP} \end{pmatrix}_{N \times P}$$

$$Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}_{N \times 1}$$

最小2乘估计：

$$\begin{aligned} L(w) &= \sum_{i=1}^N \|w^T x_i - y_i\|^2 = \sum_{i=1}^N (w^T x_i - y_i)^2 = (w^T x_1 - y_1, w^T x_2 - y_2, \dots, w^T x_N - y_N) \begin{pmatrix} w^T x_1 - y_1 \\ w^T x_2 - y_2 \\ \vdots \\ w^T x_N - y_N \end{pmatrix} \\ &= [w^T(x_1, x_2, \dots, x_N) - (y_1, y_2, \dots, y_N)] \begin{pmatrix} w^T x_1 - y_1 \\ w^T x_2 - y_2 \\ \vdots \\ w^T x_N - y_N \end{pmatrix} = (w^T x^T - Y^T) \cdot (xw - Y) \\ &= w^T x^T x w - w^T x^T Y - Y^T x w + Y^T Y = w^T x^T x w - 2w^T x^T Y + Y^T Y \end{aligned}$$

$$\hat{w} = \arg \min L(w)$$

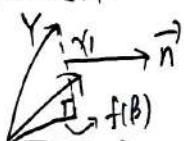
$$\frac{\partial L(w)}{\partial w} = 2x^T x w - 2x^T Y = 0.$$

$$\therefore x^T x w = x^T Y$$

$$w = (x^T x)^{-1} x^T Y \quad \text{解析解}$$

x⁺ 伪逆

几何理解。



$\vec{n} \perp$ 于 P 维空间

$$\vec{n} = (Y - X\beta)$$

$$\text{由 } \vec{a} \perp \vec{b} \Leftrightarrow \vec{a}^T \vec{b} = 0$$

Y 在 P 维空间中 (噪声 noise, random)

Y 到 P 维空间最短距离 \vec{n} (垂线法)

$f(\beta)$ 这一个投影就是 (x_1, \dots, x_p) 的线性组合 $\Leftrightarrow x\beta$

则有 $\vec{n} = (Y - X\beta)$, 而由 $\vec{n} \perp P$ 维空间, 则有 $x^T(Y - X\beta) = 0 \Leftrightarrow \boxed{\beta = (x^T x)^{-1} x^T Y}$.

$$f(w) = w^T x = x^T \beta$$

$$\beta = \begin{pmatrix} b_1 \\ b_p \end{pmatrix}_{p \times 1}, x^T = (\text{系数})_{p \times N}$$

$$(Y - X\beta) \perp x \Leftrightarrow x^T(Y - X\beta) = 0 = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}_{p \times 1} \Leftrightarrow x^T Y = x^T X \beta$$

$$\boxed{\beta = (x^T x)^{-1} x^T Y}$$

2.2. Least Square Method - from probabilistic perspective

概率视角

$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$$

$$x_i \in \mathbb{R}^P, y_i \in \mathbb{R}, i=1, 2, \dots, N$$

$$x = (x_1, x_2, \dots, x_N)^T = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{pmatrix} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1P} \\ x_{21} & x_{22} & \dots & x_{2P} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \dots & x_{NP} \end{pmatrix}_{N \times P}$$

$$Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}$$

x : 样本 y 值

$$P(y|x; w) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y - w^T x)^2}{2\sigma^2}\right)$$

最小二乘估计：

$$\begin{aligned} L(w) &= \sum_{i=1}^N \|w^T x_i - y_i\|^2 \\ \hat{w} &= \arg \min_w L(w) \\ \hat{w} &= (x^T x)^{-1} x^T Y \end{aligned}$$

噪声: $\epsilon \sim N(0, \sigma^2)$

$$y = f(w) + \epsilon$$

$$f(w) = w^T x$$

$$y = w^T x + \epsilon$$

$$y|x, w \sim N(w^T x, \sigma^2)$$

$y|x, w$ 是给定 x, w 的条件分布, 方差来源是误差
这里 x, w 看作 constant

给定特征 x 和 w 的条件下, 标签 y 的概率分布

w 是模型参数

给定 w 时, 观测到 y 的概率且套上 \log , 称为 log-likelihood

$$\begin{aligned} MLE: \quad \underset{\text{log-likelihood}}{L(w)} &= \log P(Y|x, w) = \log \prod_{i=1}^N P(y_i|x_i, w) = \sum_{i=1}^N \log P(y_i|x_i, w) = \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi}\sigma} + \log \exp\left(-\frac{(y_i - w^T x_i)^2}{2\sigma^2}\right) \\ MLE &= \sum_{i=1}^N \left[\log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{2\sigma^2} (y_i - w^T x_i)^2 \right] \end{aligned}$$

定义相似.

$$\hat{w} = \underset{\substack{\text{最小二乘} \\ \text{极似然}}}{\arg \max_w} L(w) = \arg \max_w \left[-\frac{1}{2\sigma^2} (y - w^T x)^2 \right] = \underset{w}{\arg \min} (y - w^T x)^2$$

则有 $\hat{w} = (x^T x)^{-1} x^T Y$, x : 样本 y : 观测值

$LSE \leq MLE$

条件 (noise is Gaussian Dist)

正则化-岭回归-频率角度

w is constant.

$$\text{Loss Function: } L(w) = \sum_{i=1}^N \|w^T x_i - y_i\|^2 \rightarrow \text{目标函数}$$

$$\hat{w} = (x^T x)^{-1} x^T Y$$

$X \in \mathbb{R}^{N \times p}$, N 个样本, $x_i \in \mathbb{R}^p$, 正常来说 $N > p$, 但也会存在 $N < p$ 甚至 $N \ll p$, 无法求 x^{-1} , 无法得到解析解的情况, 此时引入正则化

样本少, 维度高

正则化框架

$$\underset{w}{\arg \min} \left[L(w) + \lambda P(w) \right] \quad \cdots \quad (1)$$

loss Penalty (惩罚项)

① 加数据

且容易过拟合

② 解析解 / 特征提取 PCA

③ 正则化

L_1 : Lasso, $P(w) = \|w\|_1$

L_2 : Ridge, 岭回归, $P(w) = \|w\|_2^2 = w^T w$

权值衰减

$$\begin{aligned} \text{对于 } (1) \rightarrow L_2 \text{ 有: } (1) &= \sum_{i=1}^N \|w^T x_i - y_i\|^2 + \lambda w^T w = (w^T x^T - Y^T)(x w - Y) + \lambda w^T w = w^T x^T x w - 2w^T x^T Y + Y^T Y + \lambda w^T w \\ &= w^T x^T x w + \underbrace{w^T \lambda w}_{\text{实数可以写到中间}} - 2w^T x^T Y + Y^T Y = w^T (x^T x + \lambda I) w - 2w^T x^T Y + Y^T Y \end{aligned}$$

PC 推导

$$\underset{w}{\operatorname{argmin}} J(w) = \hat{w}$$

$$\frac{\partial J(w)}{\partial w} = 2(x^T x + \lambda I)w - 2x^T Y \triangleq 0$$

$$W = (x^T x + \lambda I)^{-1} x^T Y$$

$x^T x$ 半正定 λ 正定 + 对角 \Rightarrow 可逆
 $I \rightarrow \text{对角}$

LSE

$$W = (x^T x)^{-1} x^T Y$$

2.4. Regularization - Ridge Regression - Bayesians

Bayesian Perspective: w is r.v.

$$w \sim N(0, \sigma^2 I), P(y|w) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(y-w^T x)^2}{2\sigma^2}\right], P(w) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{\|w\|^2}{2\sigma^2}\right]$$

$$P(w|y) = \frac{P(y|w)P(w)}{P(y)}$$

$$\text{MAP: } \hat{w} = \underset{w}{\operatorname{argmax}} P(w|y) = \underset{w}{\operatorname{argmax}} P(y|w)P(w)$$

$$= \underset{w}{\operatorname{argmax}} \log [P(y|w)P(w)] = \underset{w}{\operatorname{argmax}} \log \left(\frac{1}{\sqrt{2\pi}\sigma} \cdot \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(y-w^T x)^2}{2\sigma^2} - \frac{\|w\|^2}{2\sigma^2}\right] \right)$$

$$= \underset{w}{\operatorname{argmin}} \left(\frac{(y-w^T x)^2}{2\sigma^2} + \frac{\|w\|^2}{2\sigma^2} \right) \stackrel{\text{同乘 } 2\sigma^2}{=} \underset{w}{\operatorname{argmin}} \left\{ (y-w^T x)^2 + \frac{\sigma^2}{2} \|w\|^2 \right\} \dots \textcircled{1}$$

$$\text{为了简化运算,省略了} = \underset{w}{\operatorname{argmax}} \sum_{i=1}^N P(y_i|w)P(w)$$

$$\downarrow \text{上一节的 } J(w). \text{ P6.} \quad \lambda = \frac{\sigma^2}{2}$$

$$\text{则代入有: } \hat{w}_{\text{MAP}} = \underset{w}{\operatorname{argmin}} \sum_{i=1}^N \left[(y_i - w^T x_i)^2 + \frac{\sigma^2}{2} \|w\|^2 \right] = \underset{w}{\operatorname{argmin}} \sum_{i=1}^N \left[(w^T x_i - y_i)^2 + \lambda w^T w \right] = \underset{w}{\operatorname{argmin}} [L(w) + \lambda P(w)] \dots \textcircled{1}$$

Regularized LSE \Leftrightarrow MAP (noise为 Gaussian Dist)
(Prior 也是 GD)
 \star

3 Linear classification 线性分类

3.1 Linear classification Background

频率派 \rightarrow 统计机器学习
贝叶斯派 \rightarrow 概率图模型

统计机器学习 \rightarrow Linear Regression
 $y = f(w, b)$ $f(w, b) = w^T x + b$
 $x \in \mathbb{R}^p$
 \uparrow 线性判别分析 (fisher)

线性分类
硬分类: 感知机
软分类: 生成式: Gaussian Discriminal Analysis
 $y \in \{0, 1\}$
半判别式: Logistic Regression

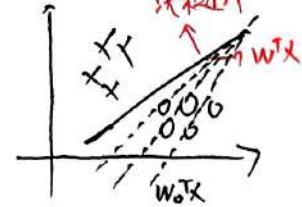
① 线性 \rightarrow 属性非线性: 特征转换 (多项式回归)
属性非线性: 线性分类 (数据函数是非线性)
属性非线性: 神经网络, 感知机

② 全局性与线性样条回归, 决策树
③ 数据未加工 \rightarrow PCA, 流形.

activation function
 \uparrow
 $y = f(w^T x + b), y \in \{0, 1\}$
 $f^{-1}: \text{link function}$
 $f: w^T x + b \mapsto \{0, 1\}$
 $f^{-1}: \{0, 1\} \mapsto w^T x + b$

线性回归 $\xrightarrow{\text{激活函数}}$ 线性分类
 \downarrow 降维
 $\{0, 1\}$

3.2. Linear Classification Perceptron Algorithm 1957年



D: {被错误分类的样本}

一次一次移动 $w_i^T x$

样本集: $\{(x_i, y_i)\}_{i=1}^N$

模型: $f(x) = \text{sign}(w^T x)$, $x \in \mathbb{R}^P$, $w \in \mathbb{R}^P$

$$\text{sign}(a) = \begin{cases} +1, & a > 0 \\ -1, & a < 0 \end{cases}$$

策略: loss function:

$$L(w) = \sum_{i=1}^N I\{y_i w^T x_i < 0\}$$

NP hard 不可导, 不连续, w 轻易变动, 且有可能变为0或1

$$\left. \begin{array}{l} x_i, y_i \\ w^T x_i > 0, y_i = +1 \\ w^T x_i < 0, y_i = -1 \end{array} \right\} \Rightarrow \left\{ \begin{array}{l} y_i w^T x_i > 0 \text{ 正确分类} \\ y_i w^T x_i < 0 \text{ 错误分类} \end{array} \right.$$

关于 w 连续函数

新 loss function:

$$L(w) = \sum_{i \in D} -y_i w^T x_i$$

$$\nabla_w L = -y_i x_i$$

算法: SGD: $w^{(t+1)} \leftarrow w^{(t)} - \lambda \nabla_w L = w^{(t)} + y_i x_i \cdot \lambda$ 学习率
随机梯度

3.3 Fisher - Linear Discrimination Analysis

线性判别分析

$$x = (x_1, x_2, \dots, x_N)^T = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{pmatrix}_{n \times p}, Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}_{n \times 1}$$

$$\{(x_i, y_i)\}_{i=1}^N, x_i \in \mathbb{R}^P, y_i \in \{+1, -1\}$$

$$X_{C_1} = \{x_i | y_i = +1\}, X_{C_2} = \{x_i | y_i = -1\}$$

$$|X_{C_1}| = N_1, |X_{C_2}| = N_2, N_1 + N_2 = N$$

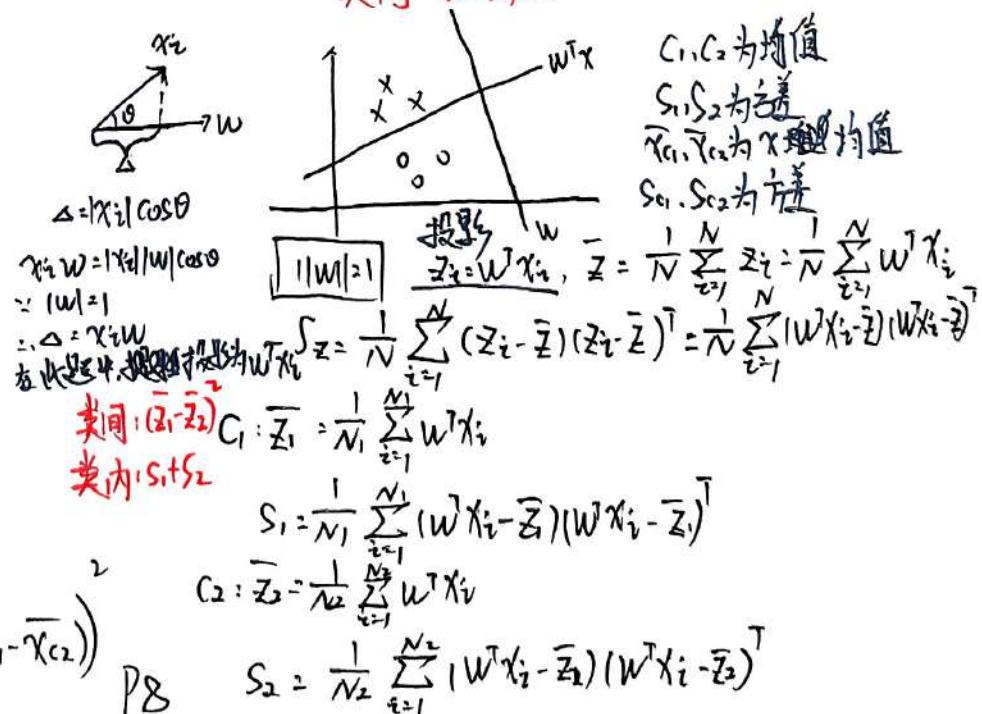
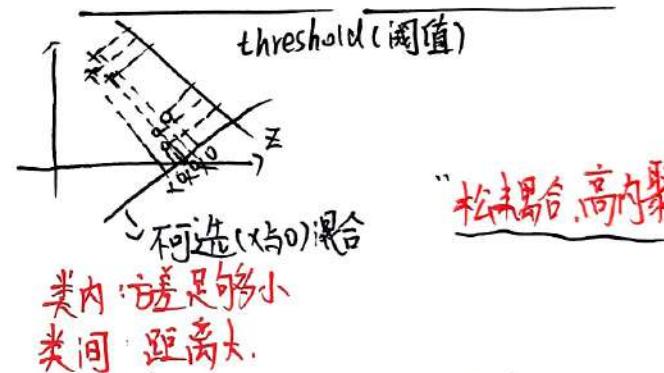
思想: 类内小, 类间大 \rightarrow $\frac{(z_1 - z_2)^2}{S_1 + S_2} \rightarrow$ 类间
目标函数: $J(w) = \frac{(z_1 - z_2)^2}{S_1 + S_2} \rightarrow$ 类内

$$w = \arg \max J(w)$$

$$J(w) = \frac{(\bar{z}_1 - \bar{z}_2)^2}{S_1 + S_2}$$

$$\text{分子: } (\frac{1}{N_1} \sum_{i=1}^{N_1} w^T x_i - \frac{1}{N_2} \sum_{i=1}^{N_2} w^T x_i)^2$$

$$= (w^T (\frac{1}{N_1} \sum_{i=1}^{N_1} x_i - \frac{1}{N_2} \sum_{i=1}^{N_2} x_i)) (w^T (\bar{x}_{C_1} - \bar{x}_{C_2}))$$



分子 = $S_1 + S_2$

$$S_1 = \frac{1}{N_1} \sum_{i=1}^{N_1} (W^T \chi_i - \frac{1}{N_1} \sum_{j=1}^{N_1} W^T \chi_j) (W^T \chi_i - \frac{1}{N_1} \sum_{j=1}^{N_1} W^T \chi_j)^T$$

$$= \frac{1}{N_1} \sum_{i=1}^{N_1} W^T (\chi_i - \bar{\chi}_{C_1}) (\chi_i - \bar{\chi}_{C_1})^T W$$

$$= W^T \left[\frac{1}{N_1} \sum_{i=1}^{N_1} (\chi_i - \bar{\chi}_{C_1}) (\chi_i - \bar{\chi}_{C_1})^T \right] W$$

S_{C_1} 表示 χ_i 与 $\bar{\chi}_{C_1}$ 的

$$= W^T \cdot S_{C_1} \cdot W$$

分母 = $W^T S_{C_1} W + W^T S_{C_2} W = W^T (S_{C_1} + S_{C_2}) W$

$$\therefore J(w) = \frac{W^T (\bar{\chi}_{C_1} - \bar{\chi}_{C_2}) (\bar{\chi}_{C_1} - \bar{\chi}_{C_2})^T W}{W^T (S_{C_1} + S_{C_2}) W}$$

$$= \frac{W^T S_b W}{W^T S_w W} \quad (\bar{\chi}_{C_1} - \bar{\chi}_{C_2}) (\bar{\chi}_{C_1} - \bar{\chi}_{C_2})^T$$

S_b : between-class 类间方差 $\Rightarrow S_b =$

S_w : within-class 类内方差 $\Rightarrow S_w = S_{C_1} + S_{C_2}$

$$J(w) = \frac{W^T S_b W}{W^T S_w W} = W^T S_b W \cdot (W^T S_w W)^{-1}$$

$$\nabla J(w) / \nabla w = 2 S_b W \cdot (W^T S_w W)^{-1} + W^T S_b W \cdot (-1) \cdot (W^T S_w W)^{-2} \cdot 2 S_w W \triangleq 0$$

$$\Rightarrow S_b W \cdot (W^T S_w W) - W^T S_b W S_w W = 0$$

$$\underbrace{W^T S_b W S_w W}_{\text{实数 GR}} = S_b W \cdot \underbrace{(W^T S_w W)}_{\substack{\text{实数} \\ \text{GR}}} \quad \begin{matrix} \text{实数} \\ \text{GR} \end{matrix} \quad \begin{matrix} \text{实数} \\ \text{GR} \end{matrix}$$

$$W^T = P \times P$$

$$S_w = P \times P$$

$$W = P \times 1$$

$$S_w W = \frac{W^T S_w W}{W^T S_b W} \cdot S_b W$$

$$W = \frac{W^T S_w W}{W^T S_b W} \cdot S_b^{-1} S_b W \quad \begin{matrix} \text{不关心 } W \text{ 大小} \\ \rightarrow \text{只关心 } W \text{ 方向.} \end{matrix}$$

$$\propto S_w^{-1} S_b W$$

$$\propto S_w^{-1} (\bar{\chi}_{C_1} - \bar{\chi}_{C_2}) \underbrace{(\bar{\chi}_{C_1} - \bar{\chi}_{C_2})^T}_{P \times P} \cdot \underbrace{W}_{P \times 1}$$

$$W \propto S_w^{-1} (\bar{\chi}_{C_1} - \bar{\chi}_{C_2}) \quad \begin{matrix} \text{GIR} \\ - \text{快矩阵} \times \text{向量.} \end{matrix} \rightarrow \text{方向.}$$

若 S_w 是对角阵且各向同性, $S_w^{-1} \propto 1$, 则 $W \propto (\bar{\chi}_{C_1} - \bar{\chi}_{C_2})$

3.4. Logistic Regression

Data: $\{(x_i, y_i)\}_{i=1}^N$

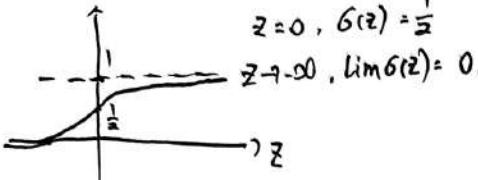
$x_i \in \mathbb{R}^P, y_i \in \{0,1\}$

Sigmoid function: $G(z) = \frac{1}{1+e^{-z}}$

$$z \rightarrow \infty, \lim G(z) = 1$$

$$z = 0, G(z) = \frac{1}{2}$$

$$z \rightarrow -\infty, \lim G(z) = 0.$$



有时我们只要得到一模型的概率，那么我们需要一种能输出 $[0,1]$ 区间的值的函数，考虑二分类模型，我们利用判别模型，希望对 $P(C|x)$ 建模，利用 Bayes.

$$P(C|x) = \frac{P(x|C_1)P(C_1)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)}$$

$$\text{取 } \alpha = \ln \frac{P(x|C_1)P(C_1)}{P(x|C_2)P(C_2)}, \text{ 于是 } P(C_1|x) = \frac{1}{1+e^{\alpha}} = \frac{1}{1+e^{w^T x}}$$

上面式子叫 Logistic Sigmoid 函数

$$\alpha = w^T x$$

$$G: \mathbb{R} \mapsto (0,1)$$

$$w^T x \mapsto P$$

$$\varphi(x, w)$$

$$\begin{aligned} p_1 &= P(y=1|x) = G(w^T x) = \frac{1}{1+e^{-w^T x}}, y=1 \\ p_0 &= P(y=0|x) = 1 - P(y=1|x) = \frac{e^{-w^T x}}{1+e^{-w^T x}}, y=0 \end{aligned}$$

$$\begin{aligned} \text{MLE: } \hat{w} &= \underset{w}{\operatorname{argmax}} \log P(Y|x) = \underset{w}{\operatorname{argmax}} \times \log \prod_{i=1}^N P(y_i|x_i) = \underset{w}{\operatorname{argmax}} \sum_{i=1}^N \log P(y_i|x_i) = \underset{w}{\operatorname{argmax}} \sum_{i=1}^N (y_i \log p_1 + (1-y_i) \log p_0) \\ &= \underset{w}{\operatorname{argmax}} \sum_{i=1}^N [y_i \log \varphi(x_i, w) + (1-y_i) \log (1-\varphi(x_i, w))] \end{aligned}$$

$$\begin{matrix} (\max) & (\min) & -\text{cross Entropy} \end{matrix}$$

MLE \Rightarrow loss function (min cross Entropy)

$$J(w) = \sum_{i=1}^N (y_i \log p_1 + (1-y_i) \log p_0), \quad p_1' = \left(\frac{1}{1+e^{-w^T x}} \right)' = p_1(1-p_1) \quad \text{损失函数是凸函数，但非二次函数，无法直接求闭式解}$$

$$\text{则 } J(w) = \sum_{i=1}^N y_i \cdot \frac{p_1(1-p_1)}{p_1} \cdot x_i + y_i p_1 x_i - p_1 x_i = \sum_{i=1}^N (y_i - p_1) x_i \rightarrow \text{含参数等10求出闭式解}$$

由于概率值非线性，无法直接求（与 Perceptron 类似），引入随机梯度上升（下降）求最大最小值

3.5. Gaussian Discriminant Analysis

Data: $\{(x_i, y_i)\}_{i=1}^N, x_i \in \mathbb{R}^P, y_i \in \{0,1\}$

$$\hat{y} = \underset{y \in \{0,1\}}{\operatorname{argmax}} P(y|x) = \underset{y \in \{0,1\}}{\operatorname{argmax}} P(y) P(x|y) \rightarrow \text{生成模型的基本原理}$$

$$y \sim \text{Bernoulli}(\phi) \quad \frac{y}{P} \mid \begin{array}{l} 1 \\ 0 \end{array} \xrightarrow{\phi^y, y=1} \frac{1-(1-\phi)^y, y=0}{\phi^y, y=1} \Rightarrow \phi^y \phi^{1-y} \text{ 合并为一项} \quad y \in \{0,1\}$$

$$\begin{cases} x|y=1 \sim N(\mu_1, \Sigma) \\ x|y=0 \sim N(\mu_2, \Sigma) \end{cases} \xrightarrow{\text{GDA}} N(\mu_1, \Sigma)^{y=1} \cdot N(\mu_2, \Sigma)^{y=0}$$

$$\begin{aligned} \log -\text{likelihood} &= \log \prod_{i=1}^N P(x_i, y_i) = \sum_{i=1}^N [\log P(x_i|y_i) P(y_i)] = \sum_{i=1}^N [\log P(x_i|y_i)] + \sum_{i=1}^N [\log P(y_i)] = \sum_{i=1}^N [\log P(x_i|y_i)] + \log P(y_i) \\ \theta &= (\mu_1, \mu_2, \Sigma, \phi) = \sum_{i=1}^N [\log N(\mu_1, \Sigma)^{y_i} \cdot N(\mu_2, \Sigma)^{1-y_i} + \log \phi^{y_i} (1-\phi)^{1-y_i}] \end{aligned}$$

$$\begin{aligned} \hat{\theta} &= \operatorname{argmax} \mathcal{L}(\theta) = \sum_{i=1}^N [\log N(\mu_1, \Sigma)^{y_i} + \log N(\mu_2, \Sigma)^{1-y_i} + \log \phi^{y_i} (1-\phi)^{1-y_i}] = \sum_{i=1}^N [\underbrace{y_i \log N(\mu_1, \Sigma)}_{①} + \underbrace{(1-y_i) \log N(\mu_2, \Sigma)}_{②} + \underbrace{\log \phi^{y_i} (1-\phi)^{1-y_i}}_{③}] \end{aligned}$$

$$\begin{aligned} y=1 &: N_1 \\ y=0 &: N_2 \\ N &= N_1 + N_2 \end{aligned}$$

3. b. 求得GPA参数(学习过程)

$$\text{求 } \phi: \textcircled{3} = \sum_{i=1}^N \log \phi \frac{y_i}{\phi} + (1-y_i) \log (1-\phi) > \sum_{i=1}^N y_i \log \phi$$

$$\frac{\partial \textcircled{3}}{\partial \phi} = \sum_{i=1}^N \left[y_i \log \phi + (1-y_i) \log (1-\phi) \right]$$

$$\frac{\partial \textcircled{3}}{\partial \phi} = \sum_{i=1}^N \left[y_i \frac{1}{\phi} + (1-y_i) \cdot \frac{1}{1-\phi} (-1) \right] = \sum_{i=1}^N \left[y_i \frac{1}{\phi} - (1-y_i) \frac{1}{1-\phi} \right] \triangleq 0$$

$$\Rightarrow \sum_{i=1}^N [(1-y_i)\phi y_i - (1-y_i)\phi] = 0 \Leftrightarrow \sum_{i=1}^N [y_i - \phi y_i - \phi + \phi y_i] = 0$$

$$\sum_{i=1}^N [y_i - \phi] = 0 \quad \begin{cases} y=1, N_1 \\ y=0, N_2 \end{cases} \quad \sum_{i=1}^N y_i = \sum_{i=1}^{N_1} y_i + \sum_{i=1}^{N_2} y_i = N_1 \cdot 1 + N_2 \cdot 0 = N,$$

$$\therefore \hat{\phi} = \frac{1}{N} \sum_{i=1}^N y_i = \frac{N_1}{N}$$

$$\text{求 } \mu_1^0 = \sum_{i=1}^N \log N(\boldsymbol{x}_i, \boldsymbol{\Sigma})^y_i = \sum_{i=1}^N y_i \log \frac{1}{(2\pi)^{\frac{D}{2}} |\boldsymbol{\Sigma}|} \exp \left\{ -\frac{1}{2} (\boldsymbol{x}_i - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{x}_i - \boldsymbol{\mu}_1) \right\}$$

$$\begin{aligned} \mu_1 &= \arg \max \textcircled{1} = \arg \max_{i=1}^N y_i \cdot \left[-\frac{1}{2} (\boldsymbol{x}_i - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{x}_i - \boldsymbol{\mu}_1) \right] = \arg \max \frac{1}{2} \sum_{i=1}^N y_i [\boldsymbol{x}_i - \boldsymbol{\mu}_1]^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{x}_i - \boldsymbol{\mu}_1) \\ &= \arg \max \frac{1}{2} \sum_{i=1}^N y_i (\boldsymbol{x}_i^T \boldsymbol{\Sigma}^{-1} - \boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1}) (\boldsymbol{x}_i - \boldsymbol{\mu}_1) = \arg \max \sum_{i=1}^N y_i (\underbrace{\boldsymbol{x}_i^T \boldsymbol{\Sigma}^{-1} \boldsymbol{x}_i}_{\text{常数}} - \underbrace{\boldsymbol{x}_i^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1}_{\text{常数}} - \underbrace{\boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{x}_i}_{\text{常数}} + \underbrace{\boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1}_{\text{常数}}) \\ &= \arg \max \underbrace{\sum_{i=1}^N y_i}_{\Delta} \underbrace{(\boldsymbol{x}_i^T \boldsymbol{\Sigma}^{-1} \boldsymbol{x}_i - 2 \boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{x}_i + \boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1)}_{\Delta} \end{aligned}$$

$$\frac{\partial \Delta}{\partial \boldsymbol{\mu}_1} = -\frac{1}{2} \sum_{i=1}^N y_i (-2 \boldsymbol{\Sigma}^{-1} \boldsymbol{x}_i + 2 \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1) \triangleq 0 \quad \text{常数舍弃}$$

$$\sum_{i=1}^N y_i (\boldsymbol{\Sigma}^{-1} \boldsymbol{x}_i - \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1) \triangleq 0$$

$$\sum_{i=1}^N y_i (\boldsymbol{\mu}_1 - \boldsymbol{x}_i) = 0$$

$$\sum_{i=1}^N y_i \boldsymbol{\mu}_1 = \sum_{i=1}^N y_i \boldsymbol{x}_i$$

$$\hat{\boldsymbol{\mu}}_1 = \frac{\sum_{i=1}^N y_i \boldsymbol{x}_i}{\sum_{i=1}^N y_i} = \frac{\sum_{i=1}^N y_i \boldsymbol{x}_i}{N_1}$$

求 $\boldsymbol{\mu}_2$:

$$\frac{\partial \Delta}{\partial \boldsymbol{\mu}_2} = \sum_{i=1}^N (1-y_i)(\boldsymbol{\mu}_2 - \boldsymbol{x}_i) = 0$$

$$\hat{\boldsymbol{\mu}}_2 = \frac{\sum_{i=1}^N \boldsymbol{x}_i (1-y_i)}{\sum_{i=1}^N (1-y_i)} = \frac{\sum_{i=1}^N \boldsymbol{x}_i (1-y_i)}{N_1 + N_2 - N_1} = \frac{\sum_{i=1}^N \boldsymbol{x}_i (1-y_i)}{N_2}$$

求 $\boldsymbol{\Sigma}$:

$$C_1 = \{\boldsymbol{x}_i | y_i = 1, i=1, 2, \dots, n\} \quad |C_1| = N_1, |C_2| = N_2$$

$$C_2 = \{\boldsymbol{x}_i | y_i = 0, i=1, 2, \dots, n\} \quad N_1 + N_2 = N$$

$$\hat{\boldsymbol{\Sigma}} = \arg \max \textcircled{1} + \textcircled{2}, \quad \textcircled{1} + \textcircled{2} = \sum_{\boldsymbol{x}_i \in C_1} \log N(\boldsymbol{x}_i, \boldsymbol{\Sigma}) + \sum_{\boldsymbol{x}_i \in C_2} \log N(\boldsymbol{x}_i, \boldsymbol{\Sigma})$$

$$\begin{aligned}\frac{\partial \text{tr}(AB)}{\partial A} &= B^T \\ \frac{\partial |A|}{\partial A} &= |A| \cdot A^{-1} \\ \text{tr}(AB) &= \text{tr}(BA) \\ \text{tr}(ABC) &= \text{tr}(B \cdot CA) = \text{tr}(CAB)\end{aligned}$$

$$S = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}) (\mathbf{x}_i - \boldsymbol{\mu})^T$$

$$\begin{aligned}&= \sum_{i=1}^N \text{tr}((\mathbf{x}_i - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x}_i - \boldsymbol{\mu})) \\ &= \sum_{i=1}^N \text{tr}((\mathbf{x}_i - \boldsymbol{\mu}) (\mathbf{x}_i - \boldsymbol{\mu})^T \Sigma^{-1}) \\ &= \text{tr}(\underbrace{\sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}) (\mathbf{x}_i - \boldsymbol{\mu})^T}_{NS} \Sigma^{-1}) \\ &= N \text{tr}(S \Sigma^{-1})\end{aligned}$$

$$\frac{\partial \text{tr}(S_1 \Sigma^{-1})}{\partial \Sigma^{-1}} = S_1^T \Sigma^{-2} = -S_1^T \Sigma^{-2}$$

$$\begin{aligned}\therefore (1)+(2) &= -\frac{1}{2} N_1 \log |\Sigma| - \frac{1}{2} N_1 \cdot \text{tr}(S_1 \cdot \Sigma^{-1}) - \frac{1}{2} N_2 \log |\Sigma| - \frac{1}{2} N_2 \cdot \text{tr}(S_2 \cdot \Sigma^{-1}) + C \\ &= -\frac{1}{2} \underbrace{N \log |\Sigma|}_{N=N_1+N_2} - \frac{1}{2} N_1 \cdot \text{tr}(S_1 \cdot \Sigma^{-1}) - \frac{1}{2} N_2 \cdot \text{tr}(S_2 \cdot \Sigma^{-1}) + C \\ &= -\frac{1}{2} (N \log |\Sigma| + N_1 \text{tr}(S_1 \cdot \Sigma^{-1}) + N_2 \cdot \text{tr}(S_2 \cdot \Sigma^{-1})) + C\end{aligned}$$

$$\begin{aligned}\frac{\partial (1)+(2)}{\partial \Sigma} &= -\frac{1}{2} \left(N \cdot \frac{|\Sigma| \cdot \Sigma^{-1}}{|\Sigma|} - N_1 \cdot S_1^T \Sigma^{-2} - N_2 \cdot S_2^T \Sigma^{-2} \right) \\ &= -\frac{1}{2} (N \Sigma^{-1} - N_1 S_1^T \Sigma^{-2} - N_2 S_2^T \Sigma^{-2}) \triangleq 0.\end{aligned}$$

则有:

$$\begin{aligned}N \Sigma - N_1 S_1^T - N_2 S_2^T &= 0 \\ \Sigma &= \frac{1}{N} (N_1 S_1 + N_2 S_2)\end{aligned}$$

故有

$$\begin{cases} \hat{\boldsymbol{\mu}}_1 = \frac{\sum_{i=1}^N y_i \mathbf{x}_i}{N_1} \\ \hat{\boldsymbol{\mu}}_2 = \frac{\sum_{i=1}^N (1-y_i) \mathbf{x}_i}{N_2} \\ \Sigma = \frac{1}{N} (N_1 S_1 + N_2 S_2) \\ \hat{\phi} = \frac{N_1}{N} \end{cases}$$

3.7 Naive Bayes classifier

Data: $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$

$\mathbf{x}_i \in \mathbb{R}^p$, $y_i \in \{0, 1\}$

给定 $\mathbf{x}, y=?$ 0/1

$\hat{y} = \arg \max_y P(y|\mathbf{x})$

$= \arg \max_{y \in \{0, 1\}} \frac{P(x, y)}{P(x)}$

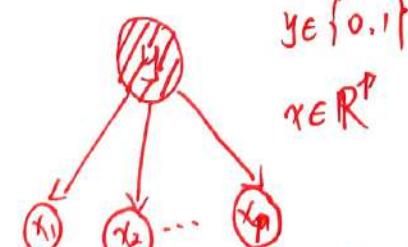
$= \arg \max_y P(y) P(x|y)$

思想: 朴素贝叶斯假设

条件独立性假设

最简单的概率图模型

有何图



且 $x_i \perp x_j | y$ ($i \neq j$) 相互独立

二分类 0/1: $y \sim \text{Bernoulli Dist}$

多分类: $y \sim \text{Categorical Dist}$

$$P(x|y) = \prod_{j=1}^p P(x_j|y)$$

动机: 简化运算

$$P(y|x) = \frac{P(x,y)}{P(x)} = \frac{P(y) \cdot P(x|y)}{P(x)} \propto P(y) \cdot P_{\text{prior}}(x)$$

做 n 次 \rightarrow Binomial 分布

Bernoulli \rightarrow Binomial

Categorical \rightarrow Multinomial

做 n 次

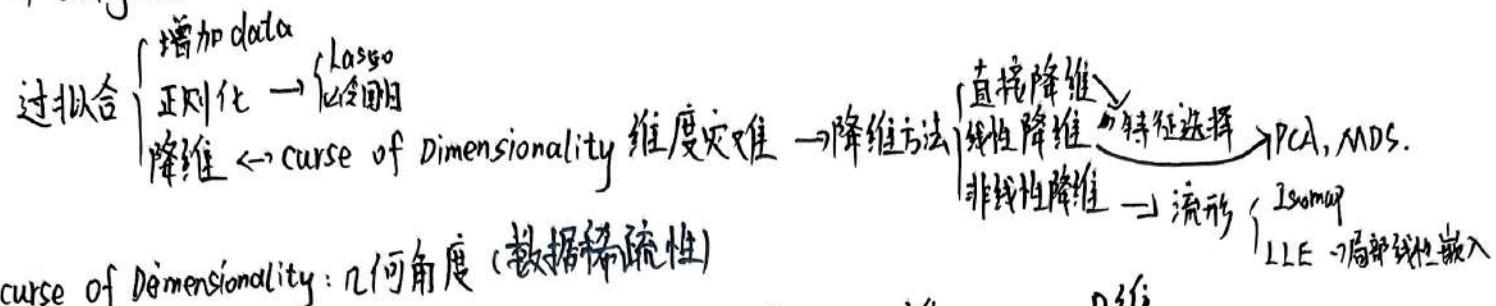
多项分布

x 离散 $\rightarrow x_j \sim \text{Categorical Dist}$

x 连续 $\rightarrow x_j \sim N(\mu_j, \sigma_j^2)$

4. Dimensionality Reduction 降维

4.1 Background.



D²维:



$$V_{超立方体} = 1$$

$$V_{超球体} = k \cdot 0.5^D \quad D \rightarrow \infty$$

说明数据都在球内, 分散的

2维

3维

... D维

1

$$k \cdot 0.5^D \rightarrow 0$$

正: 1

1

1

圆: $\pi \cdot 0.5^2$

$\frac{4}{3}\pi \cdot 0.5^3$

$k \cdot 0.5^D$

②



r=1, D维

$$V_{超球体} = k \cdot 1^D = k$$

$$V_{环形带} = V_{超球体} - V_{内球} = k - k(r-\epsilon)^D$$

$$\text{则有: } \frac{V_{超球体}}{V_{超球体}} = 1 - (r-\epsilon)^D, \quad r=1, 0 < \epsilon < 1$$

$$= 1 - (1-\epsilon)^D$$

$$\lim_{D \rightarrow \infty} (1-\epsilon)^D = 0, \text{ 则有 } \lim_{D \rightarrow \infty} \frac{V_{超球体}}{V_{超球体}} = 1$$

4.2 Sample Mean & Variance Matrix

$$\text{Data: } X = (X_1, X_2, \dots, X_N)_{N \times P}^T = \begin{pmatrix} X_1^T \\ X_2^T \\ \vdots \\ X_N^T \end{pmatrix} = \begin{pmatrix} X_{11} & X_{12} & \dots & X_{1P} \\ X_{21} & X_{22} & \dots & X_{2P} \\ \vdots & \vdots & \ddots & \vdots \\ X_{N1} & X_{N2} & \dots & X_{NP} \end{pmatrix}_{N \times P}$$

$$I_N = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}_{N \times 1}$$

$$\text{Sample Mean: } \bar{X} = \frac{1}{N} \sum_{i=1}^N X_i = \frac{1}{N} X^T I_N$$

$$\text{Sample Covariance: } S = \frac{1}{N} \sum_{i=1}^N (X_i - \bar{X})(X_i - \bar{X})^T = \frac{1}{N} X^T H X, \quad (H = I_N - \frac{1}{N} I_N I_N^T)$$

$$\bar{X} = \frac{1}{N} \sum_{i=1}^N X_i = \frac{1}{N} (X_1, X_2, \dots, X_N) \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}_{N \times 1} = \frac{1}{N} \cdot X^T \cdot I_N$$

$$H = (I_N - \frac{1}{N} I_N I_N^T)$$

$$H^T = (I_N - \frac{1}{N} I_N I_N^T) = H$$

$$H^2 = H \cdot H = (I_N - \frac{1}{N} I_N I_N^T) \cdot (I_N - \frac{1}{N} I_N I_N^T)$$

$$= I_N - \frac{2}{N} I_N \cdot I_N^T + \frac{1}{N^2} I_N I_N^T I_N I_N^T$$

$$= I_N - \frac{2}{N} I_N I_N^T + \frac{1}{N^2} \cdot N I_N I_N^T$$

$$= I_N - \frac{1}{N} I_N I_N^T$$

$$= H$$

$$\text{则有 } H^n = H$$

$$\text{即 } H = H^T; H^2 = H$$

$$= \frac{1}{N} X^T \underbrace{(I_N - \frac{1}{N} I_N I_N^T)}_{H_w} \cdot \underbrace{(I_N - \frac{1}{N} I_N I_N^T)}_{H^T} \cdot X$$

$$= \frac{1}{N} X^T H_w H^T X = \frac{1}{N} X^T H X$$

4.3. PCA - maximum variance perspective

Principal Component Analysis, classical)

一个中心：原始特征空间的重构（相互正交的基）
相关 → 无关

两焦点：最大投影方差 ★ 本节重点
最小重构距离 ★ 4.4 重点

这节课的一个点

$$\begin{cases} \hat{u}_1 = \operatorname{argmax}_{\hat{u}_1} \hat{u}_1^T S \hat{u}_1 \\ \text{s.t. } \hat{u}_1^T \hat{u}_1 = 1 \end{cases}$$

$$L(u_1, \lambda) = u_1^T S u_1 + \lambda (u_1^T u_1 - 1)$$

$$\frac{\partial L}{\partial u_1} = 2S u_1 - 2\lambda u_1 \stackrel{\Delta}{=} 0$$

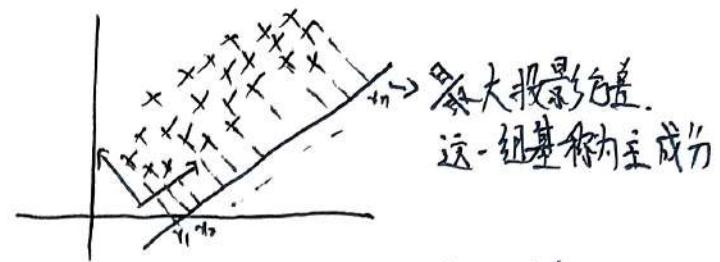
$$S u_1 - \lambda u_1 = 0$$

$$\boxed{\frac{S u_1}{\lambda} = u_1}$$

eigen-value (特征值)
eigen-vector.

$$\begin{cases} \hat{u}_j = \operatorname{argmax}_{\hat{u}_j} \sum_{i=1}^N u_i^T S u_j \\ \text{s.t. } \hat{u}_j^T \hat{u}_j = 1 \end{cases}$$

$$\Rightarrow J = \sum_{i=1}^N \lambda_i (\text{大的前})$$



而把 $\hat{x}_1, \hat{x}_2, \dots, \hat{x}_N$ 这些投影点重构回去，
恢复成原样，代价就叫最小重构距离

① 中心化 (平移)

$$x_i - \bar{x} \quad u_1 \quad \|u_1\|=1$$

② 投影误差：(已中心化，均值已减去)

$$J = \frac{1}{N} \sum_{i=1}^N ((x_i - \bar{x})^T u_1)^2$$

$$= \sum_{i=1}^N \frac{1}{N} u_1^T (x_i - \bar{x}) \cdot (x_i - \bar{x})^T u_1$$

$$= u_1^T \left(\sum_{i=1}^N \frac{1}{N} (x_i - \bar{x})(x_i - \bar{x})^T \right) u_1$$

$$J = u_1^T S u_1$$

4.4. PCA - minimum error perspective

$$J = \frac{1}{N} \sum_{i=1}^N \|x_i - \hat{x}_i\|^2 \rightarrow \text{前 } q \text{ 个相同}$$

$$= \frac{1}{N} \sum_{i=1}^N \left\| \sum_{k=1}^p (x_i^T u_k) u_k \right\|^2$$

有：坐标系 $u_{q+1}, u_{q+2}, \dots, u_p$
坐标 $x_i^T u_{q+1}, \dots, x_i^T u_p$

$$\alpha = (a_1, a_2, \dots, a_n)$$

$$\|\alpha\|^2 = \sum_{i=1}^n a_i^2$$

中心化：

$$\hat{x}_i = \frac{1}{N} \sum_{i=1}^N \sum_{k=q+1}^p ((x_i - \bar{x})^T u_k) u_k$$

$$= \sum_{i=1}^N \frac{1}{N} \sum_{k=q+1}^p ((x_i - \bar{x})^T u_k)^2 = \sum_{k=q+1}^p \sum_{i=1}^N \frac{1}{N} ((x_i - \bar{x})^T u_k)^2$$

$$u_k^T S u_k$$

$$= \sum_{k=q+1}^p u_k^T S u_k$$

$$\text{s.t. } u_k^T u_k = 1$$

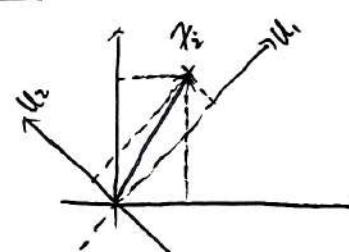
这节课的所有点

$$u_k = \operatorname{argmin}_{u_k} \sum_{i=1}^N u_k^T S u_k \quad (u_{q+1}, \dots, u_p) \text{ 线性无关，可以一个一个求，参考 4.3 求最小较小的 } \lambda_i$$

重构余量中的 λ_i

$$\text{s.t. } u_k^T u_k = 1$$

$$J \propto = \sum_{k=q+1}^p \lambda_k \quad (\text{余量中最小的这些 } \lambda_i)$$



$$x_i = (x_i^T u_1) u_1 + (x_i^T u_2) u_2$$

提

则对于 x_1, x_2, \dots, x_N 经过去中心化。
重构有 u_1, u_2, \dots, u_p

$$x_i = \sum_{k=1}^p (x_i^T u_k) u_k$$

$$\hat{x}_i = \frac{1}{N} \sum_{k=1}^p (x_i^T u_k) u_k \rightarrow \text{重构后第 } q \text{ 组}$$

$$J = \frac{1}{N} \sum_{i=1}^N \|x_i - \hat{x}_i\|^2$$

N 样本点的最小重构代价

方差矩阵 S 进行特征值分解 $= S = G K G^T; G^T G = I$

$$K = \begin{bmatrix} k_1 & k_2 & \dots & k_p \end{bmatrix} \quad k_1 \geq k_2 \geq \dots \geq k_p$$

4.5. PCA - SVD perspective (SVD: 奇异值分解)

参考 4.2 中数据格式

$$\text{中心化: } Hx = V \Sigma V^T$$

$$S = \frac{1}{n} X^T H X$$

为了方便推导, 先省略方

$$\left| \begin{array}{l} V^T V = I \\ V^T V = V V^T = I \end{array} \right.$$

$$\left| \begin{array}{l} H^T = H \\ H^T = H \end{array} \right.$$

V 是特征向量 代表方向
方向, 坐标为摄影: $Hx \cdot V$

$$\text{则 } S = X^T H X = X^T H^T H X = V \Sigma U^T \cdot V \Sigma V^T = V \cdot \Sigma^2 \cdot V^T$$

由上节 4.4. $S = G K G^T$ 得: $G = V$, $K = \Sigma^2$ \star

①先进行中心化; ②再进行奇异值分解; 效果相同

$$T \triangleq H X X^T H^T = V \Sigma V^T \cdot V \Sigma U^T = V \Sigma^2 U^T$$

T 和 S 有相同的 eigenvalue

eigen-vector

坐标矩阵

特征向量的分解

S : 特征分解: 得到方向(主成分), 然后 $Hx \cdot V \rightarrow$ 坐标: $V \Sigma V^T \cdot V = V \Sigma$

T : 特征分解: 直接得到坐标. 称为主坐标分析 (Principle coordinate analysis) PCA

$$\begin{aligned} T &= U \Sigma^2 U^T \\ T \Sigma &= U \Sigma^2 U^T \cdot V \Sigma \\ &= U \Sigma^3 \\ &= U \Sigma \cdot \Sigma^2 \end{aligned}$$

$$\begin{aligned} T \Sigma &= U \Sigma^2 U^T \cdot V \Sigma \\ &= U \Sigma \cdot \Sigma^2 \end{aligned}$$

Σ^2 : 特征值

$U \Sigma$: 特征向量

, 而有: 坐标: $Hx \cdot V = U \Sigma$ 知
此方式下的特征向量就是坐标

特征值

根据 n 与 p 的大小适当选取 PCA 或 PCoA. 维度过高用 PCoA.

4.6 PPCA - probabilistic PCA

$$x \in \mathbb{R}^p, z \in \mathbb{R}^q, q < p$$

\downarrow latent variable (无法观测的变量)

observed data (观察数据)

$$\begin{bmatrix} \sigma^2 & \sigma^2 & \dots & \sigma^2 \end{bmatrix} \rightarrow \text{各向同性}$$

$$z \sim N(0_{q \times 1}, I_{q \times q}), x = wz + \mu + \epsilon, \epsilon \sim N(0, \sigma^2 I_p)$$

Linear Gaussian Model

$$z, x | z, x, z | x$$

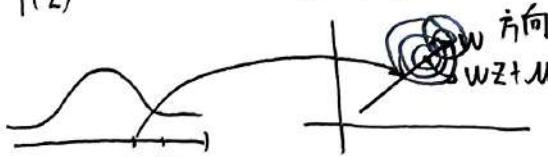
$$\text{P-PCA} \left\{ \begin{array}{l} \text{Inference: } P(z|x) \\ \text{Learning: } W, \mu, \sigma^2 \rightarrow \text{EM (期望最大)} \end{array} \right. \rightarrow \text{① 求 } P(z|x)$$

$$\text{Learning: } W, \mu, \sigma^2 \rightarrow \text{EM (期望最大)} \rightarrow \text{② 计算参数}$$

$$P(z)$$

$$P(x|z)$$

$$P(x)$$



采样点越多, 形成更多高斯分布

$$z \sim N(0, I)$$

$$\chi = Wz + \mu + \varepsilon$$

$$\varepsilon \sim N(0, \sigma^2 I)$$

$$\varepsilon \perp z$$

$$E[\chi|z] = E[Wz + \mu + \varepsilon] = Wz + \mu$$

$$\text{Var}[\chi|z] = \text{Var}[Wz + \mu + \varepsilon] = \text{Var}[\varepsilon] = \sigma^2 I$$

$$\chi|z \sim N(Wz + \mu, \sigma^2 I)$$

$$E[\chi] = E[Wz + \mu + \varepsilon] = E[Wz + \mu] + E[\varepsilon] = \mu, \quad \text{Var}[\chi] = \text{Var}[Wz + \mu + \varepsilon] = \text{Var}[Wz + \mu] + \text{Var}[\varepsilon] \\ = \text{Var}[Wz] + \text{Var}[\varepsilon] = W[I]W^T + \sigma^2 I = WW^T + \sigma^2 I$$

$$\boxed{\chi \sim N(\mu, WW^T + \sigma^2 I)} \rightarrow \begin{pmatrix} \chi \\ z \end{pmatrix} \sim N\left(\begin{bmatrix} \mu \\ 0 \end{bmatrix}, \begin{bmatrix} WW^T + \sigma^2 I \\ 0 \end{bmatrix}\right) \text{ 求 } \Delta$$

$$\chi = \begin{pmatrix} \chi_a \\ \chi_b \end{pmatrix}, \quad \chi \sim N\left(\begin{bmatrix} \mu_a \\ \mu_b \end{bmatrix}, \begin{bmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{bmatrix}\right) \text{ 下面为P4公式回忆}$$

$$\chi_{b-a} = \chi_b - \Sigma_{ba} \Sigma_{aa}^{-1} \chi_a$$

$$\mu_{b-a} = \mu_b - \Sigma_{ba} \Sigma_{aa}^{-1} \mu_a$$

$$\Sigma_{bba} = \Sigma_{bb} - \Sigma_{ba} \Sigma_{aa}^{-1} \Sigma_{ab} \quad \text{Schur Complementary}$$

$$\chi_{b-a} \sim N(\mu_{b-a}, \Sigma_{bba})$$

$$z: \chi_b = \chi_{b-a} + \Sigma_{ba} \Sigma_{aa}^{-1} \chi_a$$

$$\begin{cases} E[\chi_b|\chi_a] = E[\chi_{b-a}] + \Sigma_{ba} \Sigma_{aa}^{-1} \chi_a = \mu_{b-a} + \Sigma_{ba} \Sigma_{aa}^{-1} \mu_a = \mu_b + \Sigma_{ba} \Sigma_{aa}^{-1} (\chi_a - \mu_a) \end{cases}$$

$$\text{Var}[\chi_b|\chi_a] = \text{Var}[\chi_{b-a}] = \Sigma_{bba}$$

$$\underbrace{\chi_b|\chi_a \sim N(\mu_b + \Sigma_{ba} \Sigma_{aa}^{-1} (\chi_a - \mu_a), \Sigma_{bba})}_{= \mu_b} \quad (z) \sim N\left(\begin{bmatrix} \mu \\ 0 \end{bmatrix}, \begin{bmatrix} \Sigma_{bba} & 0 \\ 0 & I \end{bmatrix}\right) \quad (2)$$

P4. 例 13

$$\text{Corr} = E[(\chi - \mu)(z - 0)^T] = E[(\chi - \mu) \cdot z^T] = E[(Wz + \varepsilon) \cdot z^T] = E[Wz \cdot z^T] + \underbrace{E[\varepsilon z^T]}_{=0} = E[Wz \cdot z^T]$$

$$= W E[z \cdot z^T] = W \cdot I = W$$

$$\underbrace{E[\varepsilon] \cdot E(z^T)}_{=0}$$

$$\text{故 } (z) \sim N\left(\begin{bmatrix} \mu \\ 0 \end{bmatrix}, \begin{bmatrix} WW^T + \sigma^2 I & W \\ W^T & I \end{bmatrix}\right)$$

利用①②即可求出 $\chi_b|\chi_a$ 及 $\chi|z$

5. Support Vector Machine

5.1 hard-margin SVM - model definition

SVM有三宝，间隔，对偶，核技巧

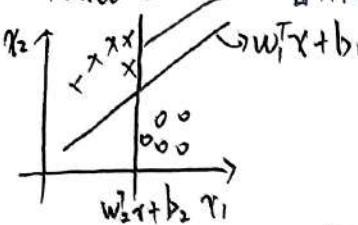
hard-margin SVM → 最大间隔分类器 $\rightarrow \max_{w,b} \text{margin}(w,b)$ s.t. $\begin{cases} w^T x_i + b > 0, y_i = +1 \\ w^T x_i + b < 0, y_i = -1 \end{cases} \Rightarrow y_i(w^T x_i + b) > 0$

soft-margin SVM

kernel SVM

鲁棒性，泛化误差不好，对噪声很敏感

还很弱



$$f(w) = \text{sign}(w^T x + b) \rightarrow \text{判别模型}$$

Data: $\{(x_i, y_i)\}_{i=1}^N, x_i \in \mathbb{R}^p, y_i \in \{-1, 1\}$

$$\text{margin}(w, b) = \min_{i=1,2,\dots,n} \frac{1}{\|w\|} |y_i(w^T x_i + b)|$$

$$= \min_{w,b,x_i} \frac{1}{\|w\|} |w^T x_i + b|$$

$$\Rightarrow \max_{w,b} \min_{i=1,2,\dots,n} \frac{1}{\|w\|} |w^T x_i + b| \Rightarrow \max_{w,b} \min_{i=1,2,\dots,n} \frac{1}{\|w\|} y_i(w^T x_i + b)$$

约束 $\rightarrow y_i(w^T x_i + b) > 0$ $r=1$ 不造成影响 \rightarrow

$$\Rightarrow \max_{w,b} \min_{i=1,2,\dots,n} y_i(w^T x_i + b)$$

$$| y_i(w^T x_i + b) > 0 = \Rightarrow Y > 0, \text{s.t. } \min_{i=1,2,\dots,n} y_i(w^T x_i + b) = Y$$

$$\text{distance} = \frac{1}{\|w\|} |w^T x_i + b|$$

总可以调整 w 使 $r=1$, $w^T x_i + b = 0$

w 与 x 相同，限制 $r=1$ 只是为了限制超平面的个数

$$\Rightarrow \begin{cases} \max_{w,b} \frac{1}{\|w\|} \\ \text{s.t. } \min_{i=1,2,\dots,n} y_i(w^T x_i + b) = 1 \\ \quad \Rightarrow y_i(w^T x_i + b) \geq 1, i=1, \dots, n \end{cases}$$

$$\Rightarrow \begin{cases} \min_{w,b} \|w\| \\ \text{s.t. } y_i(w^T x_i + b) \geq 1, \text{ for } i=1, \dots, n \end{cases}$$

N个约束

5.2 hard-margin SVM - model Solution

Data: $\{(x_i, y_i)\}_{i=1}^N, x_i \in \mathbb{R}^p, y_i \in \{-1, 1\}, \lambda = (\lambda_1, \dots, \lambda_N)$

$$\min_{w,b} \frac{1}{2} w^T w$$

①② \rightarrow Primal problem

$$\text{s.t. } y_i(w^T x_i + b) \geq 1 \Leftrightarrow \underbrace{1 - y_i(w^T x_i + b)}_{\text{带约束 for } i=1,2,\dots,N} \leq 0$$

$$L(w, b, \lambda) = \frac{1}{2} w^T w + \sum_{i=1}^N \lambda_i [1 - y_i(w^T x_i + b)]$$

约束

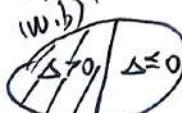
$$\min_{w,b} \max_{\lambda} L(w, b, \lambda)$$

这是一种化简优化问题的思想

$$\text{s.t. } \lambda_i \geq 0$$

①②是等价的 在 $\lambda_i \geq 0$ 时。

解释:



如果 $1 - y_i(w^T x_i + b) > 0$, $\max_{\lambda} L(w, b, \lambda) = \frac{1}{2} w^T w + \infty = \infty$

如果 $1 - y_i(w^T x_i + b) \leq 0$, $\max_{\lambda} L(w, b, \lambda) = \frac{1}{2} w^T w$, $\Leftrightarrow \min_{w,b} \max_{\lambda} L(w, b, \lambda) = \min_{w,b} \frac{1}{2} w^T w$

在优化问题②中, $\Delta > 0$ 中的 (w, b) 解全被舍弃, 且 w^*, b^* 必出自于 $\Delta \leq 0$.

$$\begin{aligned} \text{利用 } & \max_{w,b} \min_{\lambda} \mathcal{L}(w,b,\lambda) \\ \text{强对偶} & \text{s.t. } \lambda_i \geq 0. \end{aligned}$$

dual problem

$$\min_{w,b} \mathcal{L} \geq \max_{\lambda} \min_{w,b} \mathcal{L}$$

弱对偶关系

$$\min_{w,b} \mathcal{L} = \max_{\lambda} \min_{w,b} \mathcal{L}$$

强对偶关系

梯中. 凸优化, $\frac{1}{2} w^T w$ 二次, 凸二次
优化问题, 天生满足强对偶关系
约束线性, 需证明, 参看优化

$$f_0(w,b,\lambda) = \frac{1}{2} w^T w + \sum_{i=1}^N \lambda_i [1 - y_i (w^T x_i + b)]$$

$$\frac{\partial \mathcal{L}}{\partial b} = - \sum_{i=1}^N \lambda_i y_i \triangleq 0 \Rightarrow \boxed{\sum_{i=1}^N \lambda_i y_i = 0}. \text{ 将其代入 } \mathcal{L}(w,b,\lambda)$$

$$\mathcal{L}(w,b,\lambda) = \frac{1}{2} w^T w + \sum_{i=1}^N \lambda_i - \sum_{i=1}^N \lambda_i y_i w^T x_i - \underbrace{\sum_{i=1}^N \lambda_i b y_i}_{\approx 0 \therefore = 0}.$$

$$f_0(w,b,\lambda) = \frac{1}{2} w^T w + \sum_{i=1}^N \lambda_i - \sum_{i=1}^N \lambda_i y_i w^T x_i$$

$$\frac{\partial \mathcal{L}}{\partial w} = \frac{1}{2} \cdot 2 \cdot w - \sum_{i=1}^N \lambda_i y_i x_i \triangleq 0 \Rightarrow \boxed{w = \sum_{i=1}^N \lambda_i y_i x_i}. \text{ 将其代入 } \mathcal{L}(w,b,\lambda)$$

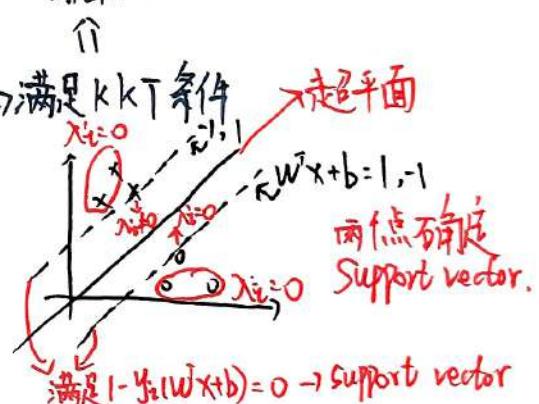
$$\mathcal{L}(w,b,\lambda) = \frac{1}{2} \left(\sum_{i=1}^N \lambda_i y_i x_i \right)^T \left(\sum_{j=1}^N \lambda_j y_j x_j \right) - \sum_{i=1}^N \lambda_i y_i \left(\sum_{j=1}^N \lambda_j y_j x_j \right)^T \cdot x_i + \sum_{i=1}^N \lambda_i$$

$$= \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j x_i^T x_j - \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j x_j^T x_i + \sum_{i=1}^N \lambda_i$$

$$= -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j x_i^T x_j + \sum_{i=1}^N \lambda_i$$

$$\begin{cases} \max \mathcal{L}_0(w,b,\lambda) \\ \text{s.t. } \lambda_i \geq 0 \end{cases} \Rightarrow \begin{cases} \min \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j x_i^T x_j - \sum_{i=1}^N \lambda_i \\ \quad \text{s.t. } \lambda_i \geq 0 \\ \quad \sum_{i=1}^N \lambda_i y_i = 0 \end{cases}$$

解出 w^*, b^*



KKT 条件: (Karush-Kuhn-Tucker) \rightarrow 原对偶问题是具有强对偶关系 \Leftrightarrow 满足 KKT 条件

$\begin{cases} \frac{\partial \mathcal{L}}{\partial w} = 0, \frac{\partial \mathcal{L}}{\partial b} = 0 \\ \lambda_i [1 - y_i (w^T x_i + b)] = 0 \\ \lambda_i \geq 0 \\ 1 - y_i (w^T x_i + b) \leq 0 \end{cases}$ slackness complementary 松弛互补

$\lambda_i [1 - y_i (w^T x_i + b)] = 0 \rightarrow$ 当 $1 - y_i (w^T x_i + b) = 0$ 时, λ_i 不等于 0.
 $y_i = \pm 1$
 $w^* = \sum_{i=1}^N \lambda_i x_i y_i, b^* = y_k - \sum_{i=1}^N \lambda_i x_i^T y_i x_k$ 对偶最优解、最优拉格朗日乘子

$$\exists (y_k, x_k), \text{s.t. } 1 - y_k (w^T x_k + b) = 0 \Leftrightarrow y_k (w^T x_k + b) = 1 \Leftrightarrow y_k^2 (w^T x_k + b) = y_k \Leftrightarrow w^T x_k + b = y_k$$

$$\therefore b^* = y_k - w^T x_k = y_k - \sum_{i=1}^N \lambda_i x_i^T y_i x_k$$

$$\therefore f(x) = \text{sign}(w^T x + b^*) \text{, 决策平面 } w^T x + b^*$$

w^* : 是 data 的线性组合, 只有在 Support vector 上的点才会对 w^* 有影响, 其他的都不让 $\lambda_i = 0$, 这是 KKT 条件.

5.3 Soft-Margin SVM

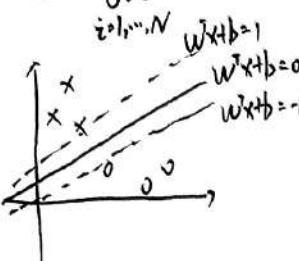
soft: 允许一点点错误: $\min_{w,b} \frac{1}{2} w^T w + \text{loss}$

Hard-Margin SVM:

Data: $\{(x_i, y_i)\}_{i=1}^N, x_i \in \mathbb{R}^p, y_i \in \{-1, 1\}$

$$\min_{w,b} \frac{1}{2} w^T w$$

$$\text{s.t. } y_i(w^T x_i + b) \geq 1$$



Loss definition: 指示函数

$$\text{① loss} = \sum_{i=1}^N I[y_i(w^T x_i + b) < 1]$$

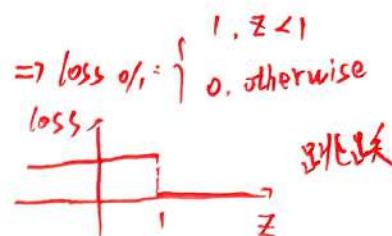
令 $z = y_i(w^T x_i + b)$ 关于 w 不连续

$$\text{② loss: 距离 (hinge loss) 铰链损失}$$

$\begin{cases} \text{if } y_i(w^T x_i + b) \geq 1, \text{ loss} = 0 \\ \text{if } y_i(w^T x_i + b) < 1, \text{ loss} = 1 - y_i(w^T x_i + b) \end{cases}$

$$\rightarrow \text{loss} = \max \{0, 1 - y_i(w^T x_i + b)\}$$

$$\text{loss}_{\max} = \max \{0, 1 - z\} \rightarrow$$



连续

Soft-Margin SVM (Hinge loss):

$$\min_{w,b} \frac{1}{2} w^T w + C \sum_{i=1}^N \max \{0, 1 - y_i(w^T x_i + b)\}$$

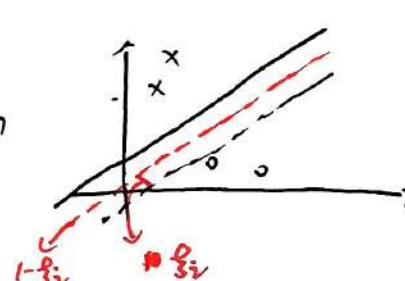
$$\text{s.t. } y_i(w^T x_i + b) \geq 1$$

$$z_i \geq 0$$

$$\exists j \lambda_j z_j = 1 - y_j(w^T x_j + b), z_j \geq 0,$$

$$\min_{w,b} \frac{1}{2} w^T w + C \sum_{i=1}^N z_i$$

$$\text{s.t. } y_i(w^T x_i + b) \geq 1 - z_i$$



求解过程与 H-SVM 雷同

① Lagrange function

② Dual function

③ KKT condition

④ Lagrange Multiplier

5.4. Constraint Optimization - Weak Duality

Primal Problem:

$$\min_{x \in \mathbb{R}^n} f(x)$$

$$\text{s.t. } m_i(x) \leq 0, i=1, \dots, M$$

$$n_j(x) = 0, j=1, \dots, N$$

Lagrangian function:

$$L(\gamma, \lambda, \eta) = f(x) + \sum_{i=1}^M \lambda_i m_i + \sum_{j=1}^N \eta_j n_j$$

Dual problem (原问题的对偶形式) (关于 x 的函数)

$$\min_{\gamma, \lambda, \eta} L(\gamma, \lambda, \eta) \Rightarrow x \in \{\text{好的 } x \text{ 集合}\}$$

$$\text{s.t. } \lambda_i \geq 0$$

如果 x 违反了约束 $m_i(x)$, 则 $m_i(x) > 0$, $\max_x L \rightarrow \infty$

如果才符合 $m_i(x) \leq 0$, $\max_x L \rightarrow \neq +\infty$, $\max_x L = f(x)$

$$\min_{\gamma, \lambda} \max_x L = \min_{\gamma, \lambda} \left\{ \max_x L, +\infty \right\} = \min_{\gamma, \lambda} \max_x L$$

Dual problem. (About $f(x, \lambda, \eta)$)

$$d \leftarrow \max_{\lambda, \eta} \min_x L(x, \lambda, \eta) \quad (\text{Dual problem is about } \lambda, \eta \text{ 's function})$$

$$\text{s.t. } \lambda_i \geq 0$$

Weak Duality:

Dual problem \leq Primal problem

$$d \leq p$$

Proof: $\max_{\lambda, \eta} \min_x L \leq \min_x \max_{\lambda, \eta} L$

$$\text{RP } \max_{\lambda, \eta} \min_x L(x, \lambda, \eta) \leq \min_x \max_{\lambda, \eta} L(x, \lambda, \eta)$$

$$\min_x L(x, \lambda, \eta) \leq L(x, \lambda, \eta) \leq \max_{\lambda, \eta} L(x, \lambda, \eta)$$

则有 $A(\lambda, \eta) \leq B(x)$

则有 $\max A(\lambda, \eta) \leq \min B(x)$

$$\max_{\lambda, \eta} \min_x L(x, \lambda, \eta) \leq \min_x \max_{\lambda, \eta} L(x, \lambda, \eta)$$

This complete the proof

5.5. 对偶性的几何解释:

$$\begin{cases} \min_{x \in \mathbb{R}^n} f(x) \\ \text{s.t. } m_i(x) \leq 0 \end{cases} \quad D \text{ 定义域}, D = \text{dom } f \cap \text{dom } m_i, \quad p^* = \inf_{t \in \mathbb{R}} \{ t | (u, t) \in G, u \leq 0 \}$$

便利于时, 取不等式约束, 且只取一个

$$L(x, \lambda) = f(x) + \lambda m_i(x), \lambda \geq 0.$$

$$p^* = \min_{x \in \mathbb{R}^n} f(x) \quad (\text{原问题的最优解})$$

$$d^* = \max_{\lambda \geq 0} \min_{x \in \mathbb{R}^n} L(x, \lambda) \quad (\text{对偶最优解})$$

$$G = \{(m_i(x), f(x)) | x \in D\} = \{(u, t) | x \in D\}$$

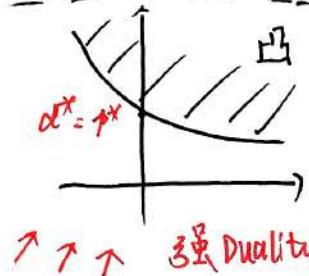
$$p^* = \inf_{t \in \mathbb{R}} \{ t | (u, t) \in G, u \leq 0 \}, d^* = \max_{\lambda \geq 0} g(\lambda)$$

$$d^* = \max_{\lambda \geq 0} \min_{x \in \mathbb{R}^n} L(x, \lambda) = \max_{\lambda \geq 0} \min_{x \in \mathbb{R}^n} (t + \lambda u) = \max_{\lambda \geq 0} g(\lambda),$$

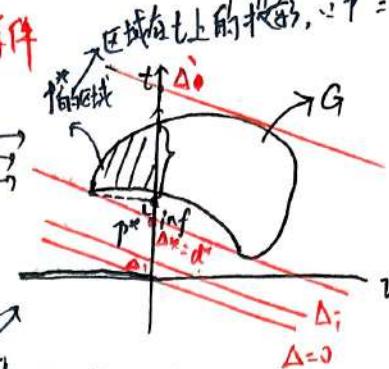
$$g(\lambda) = \inf_{t \in \mathbb{R}} \{ t + \lambda u | (u, t) \in G \} \Rightarrow \text{一次凸数与 } G \text{ 的交集的 inf.}$$

可以任取 t , 去平移, 相交

-次凸数
令 $t + \lambda u = \Delta$



强 Duality



Δ 在 $g(\lambda)$ 中, 它不是该入 t 组合的 $\inf \Delta$.

由于 $t + \lambda u = \Delta$, 可知 $g(\lambda)$ 就是 $\inf \{\Delta | (u, t) \in G\}$

且 Δ 是有限的, 由图上知, 不同的 λ, u 直线会平移产生不同的 Δ ,

即同一个斜率下得到的第一个切线, 就是 $g(\lambda)$ 的 - 一个元素, 而不同斜率所得的可数的 $\{\Delta_1, \dots, \Delta_n\}$ 就是 $g(\lambda)$ 的 t 而 $d^* = \max_{\lambda \geq 0} g(\lambda)$, 可知, 由图 $\Delta = \infty$ 两点确定一条直线时, 这就是 Weak Duality.

5.6. Slater Condition - Duality Problem

slater condition:

$$\exists \bar{x} \in \text{relint } D, \text{ s.t. } \forall i=1, \dots, M, m_i(\bar{x}) < 0$$

relint: relative interior: explain.

有边界就用法边界, 没有边界更好, 之后取内部区域

① 对于大多数凸优化, slater 成立

② 放松的 slater: if M 中有 k 个仿射函数, 则只需检验剩下的 $M-k$ 个是否满足 $m_i(x) \leq 0$

凸 + slater \Rightarrow 强 Duality

还可以是: else conditions.

凸二次规划: f 是凸, m_i 仿射

n_j 仿射
天然满足②

且或左半边不包含边界
的时或中有在一点

5.7 Karush-Kuhn-Tucker (KKT) Condition

$$\begin{cases} \min_{x \in \mathbb{R}^n} f(x), f(x) \text{ 可微} \\ \text{s.t. } m_i(x) \leq 0, i=1, \dots, M \\ n_j(x) = 0, j=1, \dots, N \end{cases} \quad \begin{array}{l} p^* \rightarrow x^* \\ \text{或} \\ \lambda = \begin{pmatrix} \lambda_1 \\ \vdots \\ \lambda_M \end{pmatrix} \\ \text{或} \\ y = \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix} \end{array}$$

$$L(x, \lambda, y) = f(x) + \sum_{i=1}^M \lambda_i m_i + \sum_{j=1}^N y_j n_j$$

$$g(x, y) = \min_x L(x, \lambda, y)$$

Dual problem:

$$\begin{cases} \max_{\lambda, y} g(x, \lambda, y) \\ \text{s.t. } \lambda \geq 0 \end{cases} \quad \begin{array}{l} \lambda^* \rightarrow \lambda^*, y^* \\ \text{或} \end{array}$$

convex + slater \Rightarrow strong duality \Leftrightarrow KKT (To solve the λ^*, y^*, x^*)

$$(d^* = p^*)$$

KKT:

| | |
|-------|---|
| 可行条件: | $\begin{cases} m_i(x^*) \leq 0 & (1) \\ n_j(x^*) = 0 & (2) \\ \lambda^* \geq 0 & (3) \end{cases}$ |
| 互补松弛: | $\lambda_i m_i = 0 \quad (\forall i=1, \dots, M) \quad (4)$ |
| 梯度为零: | $\frac{\partial f(x, \lambda^*, \eta^*)}{\partial x} \Big _{x=x^*} = 0 \quad (5)$ |

互补松弛与梯度为零都是从等号得出

$$\begin{aligned}
 d^* &= \max_{\lambda, \eta} g(\lambda, \eta) = g(\lambda^*, \eta^*) \\
 &= \min_{\lambda, \eta} L(x, \lambda^*, \eta^*) \\
 \lambda^* &\leq \underbrace{f(x, \lambda^*, \eta^*)}_{\text{由(5)得}} \quad \forall x \in \text{dom} f, \Rightarrow x^* \text{也成立} \rightarrow \text{取等} \\
 \text{slack} &= L(x^*, \lambda^*, \eta^*) \\
 \text{penalty} &= f(x^*) + \sum_{i=1}^M \lambda_i^* m_i + \sum_{j=1}^N \eta_j^* n_j \\
 &\stackrel{\lambda_i^* \geq 0}{\leq} f(x^*) \quad \text{这里只能取等} \rightarrow \sum_{i=1}^M \lambda_i^* m_i = 0 \Rightarrow \lambda_i m_i = 0 \quad (\forall i=1, \dots, M)
 \end{aligned}$$

6. Kernel Method

从思想角度
kernel Method

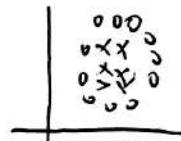
从计算角度
kernel Trick

(1) 非线性带来高维转换 (从模型角度) $x \rightarrow \phi(x)$

(2) 对偶表示带来内积 (从优化角度) $x_i^T x_j$

看①:

关于这类:



① PLA → 多层感知机 (神经网络), DP

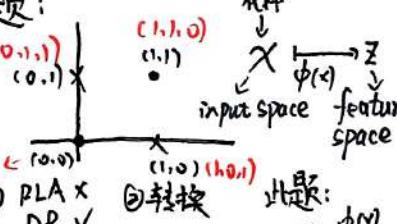
② 非线性可分 → 线性可分

非线性转换

可分超平面

看②:

异或问题:



cover Theorem: 高维比低维更易线性可分

困难之处: ① $\phi(x)$ 取值大; ② $\phi(x)$ 中内积计算大.

解决办法: 降低维数

看③: Hard-Margin SVM.

Primal problem:

$$\min_{w, b} \frac{1}{2} w^T w$$

凸优化

$$\text{s.t. } y_i(w^T x_i + b) \geq 1 \quad (i \text{ 的约束})$$

Dual problem:

$$\min_{\lambda} \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j x_i^T x_j - \sum_{i=1}^N \lambda_i$$

$$\text{s.t. } \lambda_i \geq 0, \text{ for all } i=1, 2, \dots, N, \sum_{i=1}^N \lambda_i y_i = 0$$

由此引出了 kernel function, 不去计算 $\phi(x)$, 直接得到 $\phi(x)^T \phi(x)$

左侧

线性可分 | 一点点错误 | 严格非线性 |

PLA (当初始值相关) | Pocket Algorithm | $\frac{\phi(x) + PLA}{\phi(x)}$ |

Hard-Margin SVM | Soft-Margin SVM | $\frac{-\phi(x) + Hard-Margin \text{ 转换举例}}{\phi(x)}$ |

Kernel Function: $\phi(x, x') = \phi^T(x) \phi(x') = \langle \phi(x), \phi(x') \rangle$ 利用 kernel function

$\forall x, x' \in X, \exists \phi: X \mapsto Z$ s.t. $K(x, x') = \phi(x)^T \phi(x')$ $\left[\text{求 } K(x, x'), \text{ 先求 } \phi(x), \phi(x') \right]$

对称 $K(x, x')$ 是一个核函数

e.g. $K(x, x') = \exp\left\{-\frac{(x-x')^2}{2\sigma^2}\right\}$ x, x' 为样本 "Gauss kernel function"

技巧: kernel Trick, 省去找 $\phi(x)$, 求 $\phi(x)^T \phi(x')$ 或径向基函数(RBF)

kernel function 适用场景: 线性不可分, 需要通过低维到高维变换, 但 $\phi(x)$ 难求, 直接使用 kernel Function, e.g.: 上例中的 Gauss Function

线性核: $K(x_i, x_j) = x_i^T x_j$

多项式核: $K(x_i, x_j) = (x_i^T x_j)^d$

高斯核: $K(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right)$

拉普拉斯核: $K(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|}{\sigma}\right)$

Sigmoid 核 $K(x_i, x_j) = \tanh(B x_i^T x_j + \theta)$

d 为多项式的次数

$\sigma > 0$ 为高斯核的带宽 (width)

$\sigma > 0$

\tanh 为双曲正切函数, $B > 0, \theta < 0$

6.2 Positive Definite kernel

kernel function: $K: X \times X \rightarrow \mathbb{R}, \forall x, z \in X$, 则称

$K(x, z)$ 为 kernel function

$X \rightarrow \text{input space}$

positive definite kernel: $K: X \times X \rightarrow \mathbb{R}, \forall x, z \in X$,

① 定义 有 $k(x, z)$, if $\exists: \phi: X \rightarrow \mathbb{H}$, $\phi \in \mathbb{H} \rightarrow (\text{Hilbert space})$
s.t. $k(x, z) = \langle \phi(x), \phi(z) \rangle$, 那么则称 $k(x, z)$ 为正定核函数

Positive Definite kernel: $K: X \times X \rightarrow \mathbb{R}, \forall x, z \in X$

② 定义 有 $k(x, z)$, if $k(x, z)$ 满足如下两条性质:

① 对称性

② 正定性

那么则称 $k(x, z)$ 为正定核函数

① 对称性 $\Leftrightarrow \phi(x, z) = \phi(z, x)$

② 正定性 \Leftrightarrow 任取 N 个元素, $x_1, x_2, \dots, x_N \in X$
对应的 Gram matrix 是半正定的

$$K = [k(x_i, x_j)]$$

$$\lim_{n \rightarrow \infty} k_n = K \in \mathbb{H}$$

对极限操作是封闭的

元素为函数

introduce Hilbert space: 完备的, 可能是无限维的, 被賦予内积运算的
 $f, g \in \mathbb{H}$ 线性空间

向量空间

(保加法, 保数乘)

$$\begin{cases} \langle f, g \rangle = \langle g, f \rangle \\ \langle f + g, h \rangle = \langle f, h \rangle + \langle g, h \rangle \\ \langle \alpha f, g \rangle = \alpha \langle f, g \rangle \end{cases}$$

$$= \langle f_1, g_1 \rangle + \langle f_2, g_2 \rangle$$

$$= \langle f_1, g_1 \rangle + \langle f_2, g_2 \rangle$$

证明①②等价: 由于对称性明显, 这里讨论正定性:

要证: $k(x, z) = \langle \phi(x), \phi(z) \rangle \Leftrightarrow \text{Gram matrix 半正定}$

\Rightarrow

已知 $k(x, z) = \langle \phi(x), \phi(z) \rangle$, 证 Gram matrix 半正定, 且 $k(x, z)$ 对称

证: $k(x, z) = \langle \phi(x), \phi(z) \rangle, k(z, x) = \langle \phi(z), \phi(x) \rangle$

又: 内积拥有对称性质, 即 $\langle \phi(x), \phi(z) \rangle = \langle \phi(z), \phi(x) \rangle$

$$\therefore k(x, z) = k(z, x)$$

$\therefore k(x, z)$ 满足对称性质

欲证 Gram matrix 半正定:

Gram matrix: $K = [k(x_i, x_j)]_{N \times N}$

即证: $\forall \alpha \in \mathbb{R}^N, \alpha^T K \alpha \geq 0$

$$k_{ij} = k(x_i, x_j)$$

$$\begin{aligned} \alpha^T \cdot K \cdot \alpha &= (\alpha_1, \alpha_2, \dots, \alpha_N) \begin{pmatrix} k_{11} & k_{12} & \cdots & k_{1N} \\ k_{21} & k_{22} & \cdots & k_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ k_{N1} & k_{N2} & \cdots & k_{NN} \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_N \end{pmatrix} = \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j k_{ij} = \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j \langle \phi(x_i), \phi(x_j) \rangle = \sum_{i=1}^N \alpha_i \phi(x_i)^T \sum_{j=1}^N \alpha_j \phi(x_j) \\ &= \left[\sum_{i=1}^N \alpha_i \phi(x_i) \right]^T \sum_{j=1}^N \alpha_j \phi(x_j) = \langle \sum_{i=1}^N \alpha_i \phi(x_i), \sum_{j=1}^N \alpha_j \phi(x_j) \rangle = \left\| \sum_{i=1}^N \alpha_i \phi(x_i) \right\|^2 \geq 0 \end{aligned}$$

$\therefore K$ 是半正定的

此时 Δ 为对角矩阵

\Leftarrow 对 K 进行分解, 对于对称矩阵 $K = V \Delta V^T$, 那么令 $\phi(x_i) = \sqrt{\lambda_i} v_i$, 其中 v_i 是特征向量, 于是就构造了 $k(x, z) = \sqrt{\lambda_i} v_i^T v_j$

$\langle \phi(x_i), \phi(x_j) \rangle$

小结: 对于严格可分的 Data, Hard-Margin SVM 选定一个超平面, 保证所有数据到这个超平面距离最大, 对这个平面施加约束, 固定 $y_i(w^T x_i + b) = 1$, 得到一个 convex optimization prob, 并且所有约束 conditions 都是仿射 function, 已经满足 Slater 条件, 将 primal prob \Rightarrow Dual prob, 得到等价解, 并求出约束数: $\max_w -\frac{1}{2} \sum_{i,j} \lambda_i y_i y_j w_i^T w_j + \sum_i \lambda_i$, s.t. $\lambda_i \geq 0$.

对超平面参数的求解用 strong duality's KKT condition.

$$\begin{cases} \frac{\partial L}{\partial w} = 0, \frac{\partial L}{\partial b} = 0 \\ \lambda_k (1 - y_k w^T x_k + b) = 0 \quad \text{slackness complementary} \end{cases} \Rightarrow \begin{cases} \hat{w} = \sum_{i=1}^N \lambda_i y_i x_i \\ \hat{b} = y_k - w^T x_k = y_k - \sum_{i=1}^N \lambda_i y_i^T x_i \quad \text{s.t. } k, 1 - y_k w^T (x_k + b) = 0 \\ \lambda_i \geq 0 \\ 1 - y_k w^T (x_k + b) \leq 0 \end{cases}$$

当出现一点错误, Soft-Margin SVM, 加入 Hinge Function 错误大, $\Rightarrow \arg \min_w \frac{1}{2} w^T w + C \sum_{i=1}^N \xi_i$, s.t. $y_i(w^T x_i + b) \geq 1 - \xi_i, \xi_i \geq 0, i=1 \dots n$
对完全不可分, 采用特征转换, 在 SVM 中, 以 Positive Definite Kernel 对内积进行变换, 只要满足对称、正定, 就可以做核映射

充分统计量 sufficient statistics Exponential Family Distribution

7.1. Background

↑ online learning

$$P(x|\theta) = h(x) \exp(\eta^T \phi(x) - A(\eta))$$

参数向量, $\theta \in \mathbb{R}^n$ $y = h(x)$ $\theta \rightarrow \eta$ $P(x|\eta)$ 是概率
 $A(\eta) = A(h(x))$ 密度函数

指数分布
指数族分布

$$P(x|\theta) = \frac{1}{Z} P(x|\theta)$$

归一化因子: $Z = \int_x P(x|\theta) dx$
配分函数

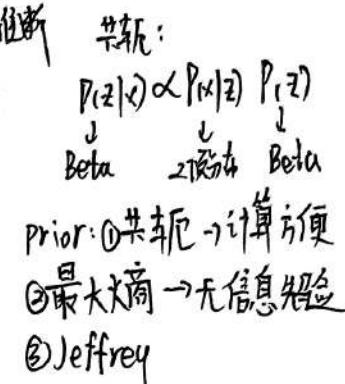
$$\int P(x|\theta) dx = \int \frac{1}{Z} P(x|\theta) dx$$

线性模型
线性组合 wx
Link function, w (数位系数)
指数族分布: $y|x$ 指数族分布

$$Z = \int P(x|\theta) dx$$

$$P(x|\eta) = h(x) \cdot \exp(\eta^T \phi(x)) \cdot \exp(-A(\eta)) = \frac{1}{\exp(A(\eta))} h(x) \cdot \exp(\eta^T \phi(x)) = \frac{1}{Z} P(x|\eta)$$

$$\exp(A(\eta)) = Z \Rightarrow A(\eta) = \log Z \Rightarrow \text{log partition function}$$



7.2. Gaussian Distribution

$$P(x|\theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\} \sim \mathcal{N}(\mu, \sigma^2)$$

$$= \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left\{-\frac{1}{2\sigma^2}(x^2 - 2\mu x + \mu^2)\right\} = \exp\left(\log(2\pi\sigma^2)^{-\frac{1}{2}}\right) \exp\left\{-\frac{1}{2\sigma^2}(x^2 - 2\mu x - \frac{\mu^2}{\sigma^2})\right\}$$

$$= \exp\left(\log(2\pi\sigma^2)^{-\frac{1}{2}}\right) \exp\left\{-\frac{1}{2\sigma^2}(-2\mu - 1)\left(\begin{pmatrix} x \\ x^2 \end{pmatrix} - \frac{\mu}{\sigma^2}\right)^T \left(\begin{pmatrix} x \\ x^2 \end{pmatrix} - \frac{\mu}{\sigma^2}\right)\right\} = \exp\left\{\underbrace{\left(\frac{1}{\sigma^2} - \frac{1}{2\sigma^2}\right)}_{\eta^T} \underbrace{\left(\begin{pmatrix} x \\ x^2 \end{pmatrix} - \frac{\mu}{\sigma^2}\right)^T \left(\begin{pmatrix} x \\ x^2 \end{pmatrix} - \frac{\mu}{\sigma^2}\right)}_{\phi(x)} - \underbrace{\left(\frac{\mu^2}{\sigma^2} + \frac{1}{2}\log(2\pi\sigma^2)^2\right)}_{A(\eta)}\right\}$$

设 $\eta = \begin{pmatrix} \eta_1 \\ \eta_2 \end{pmatrix} = \begin{pmatrix} \frac{\mu}{\sigma^2} \\ -\frac{1}{2\sigma^2} \end{pmatrix} \Rightarrow \begin{cases} \mu = -\frac{\eta_1}{2\eta_2} \\ \sigma^2 = -\frac{1}{2\eta_2} \end{cases} \Rightarrow \frac{\mu^2}{\sigma^2} = \frac{-\eta_1^2}{4\eta_2^2}$

$$A(\eta) = \frac{-\eta_1^2}{4\eta_2^2} + \frac{1}{2} \log(-\frac{\pi}{\eta_2}) ; h(x) = 1; \eta^T = \left(\frac{\mu}{\sigma^2} - \frac{1}{2\sigma^2}\right); \phi(x) = \begin{pmatrix} x \\ x^2 \end{pmatrix}$$

故 $P(x|\theta) = \exp\{\eta^T \phi(x) - A(\eta)\}$, Gauss Distribution is Exponential family Distribution

7.3 log partition function

$$P(x|\theta) = h(x) \exp(\eta^T \phi(x)) \cdot \exp(-A(\eta)) = \frac{1}{\exp(A(\eta))} h(x) \cdot \exp(\eta^T \phi(x))$$

$$\int_{-\infty}^{\infty} P(x|\theta) dx = 1, \text{ 则两边同时积分} \Rightarrow \exp(A(\eta)) = \int h(x) \exp(\eta^T \phi(x)) dx \Rightarrow \text{两边对} \eta \text{求偏导}$$

$$\Rightarrow \exp(A(\eta)) \cdot A'(\eta) = \frac{\partial}{\partial \eta} \left(\int h(x) \exp(\eta^T \phi(x)) dx \right)$$

$$= \int h(x) \exp(\eta^T \phi(x)) \phi(x) dx$$

$$A'(\eta) = \frac{\int h(x) \exp(\eta^T \phi(x)) \phi(x) dx}{\exp(A(\eta))} = \int h(x) \exp(\eta^T \phi(x) - A(\eta)) \cdot \phi(x) dx = \int P(x|\eta) \cdot \phi(x) dx$$

$$= E_{P(x|\eta)}[\phi(x)]$$

$$\therefore A'(\eta) = E_{P(x|\eta)}[\phi(x)]; A''(\eta) = \text{Var}[\phi(x)] \geq 0 \quad A(\eta) \text{ is convex function}$$

$\phi(x)$ 的期望

验证: $E[\phi(x)] = \left(\frac{E[x]}{E[x^2]}\right)$
 $E[x]$ 在 $\mathcal{N}(\mu, \sigma^2)$, 则 $E[x] = \mu$.
 $E[x^2]$ 在 $\mathcal{N}(\mu, \sigma^2)$, 则 $E[x^2] = \mu^2 + \sigma^2$.
 $A'(\eta) = -\frac{\eta_1^2}{4\eta_2^2} + \frac{1}{2} \log(-\frac{\pi}{\eta_2})$
 $A'(\eta) = \frac{G(\eta)}{\sigma^2 \eta_1} = \frac{\eta_1}{2\eta_2}, \text{ 令 } \begin{cases} \eta_1 = \frac{\mu}{\sigma^2} \\ \eta_2 = -\frac{1}{2\sigma^2} \end{cases} \Rightarrow A'(\eta) = \mu$

7.4 MLE & Sufficient Statistics

$$D = \{x_1, x_2, \dots, x_N\}$$

$$\begin{aligned} \hat{\eta}_{MLE} &= \underset{\eta}{\operatorname{argmax}} \log P(x|y) = \underset{\eta}{\operatorname{argmax}} \log \prod_{i=1}^N P(x_i|\eta) = \underset{\eta}{\operatorname{argmax}} \sum_{i=1}^N \log P(x_i|\eta) \\ &= \underset{\eta}{\operatorname{argmax}} \sum_{i=1}^N \log [h(x_i) \cdot \exp(y^T \phi(x_i) - A(\eta))] = \underset{\eta}{\operatorname{argmax}} \sum_{i=1}^N [\log h(x_i) + y^T \phi(x_i) - A(\eta)] \\ &= \underset{\eta}{\operatorname{argmax}} \sum_{i=1}^N (y^T \phi(x_i) - A(\eta)) \end{aligned}$$

$$\frac{\partial}{\partial \eta} \left(\sum_{i=1}^N (y^T \phi(x_i) - A(\eta)) \right) = \sum_{i=1}^N \frac{\partial}{\partial \eta} (y^T \phi(x_i) - A(\eta)) = \sum_{i=1}^N [\phi(x_i) - A'(\eta)] = \sum_{i=1}^N \phi(x_i) - N A'(\eta) \triangleq 0$$

$$A'(\eta) = \frac{1}{N} \sum_{i=1}^N \underbrace{\phi(x_i)}_{\text{sufficient statistics}} \Rightarrow A(\eta) \text{关于 } \eta \text{ 的函数, 则 } \eta = (A^{-1})^{-1} \text{ 反函数, 解出 } \eta_{MLE} = A^{-1}(\eta)$$

7.5 Maximum Entropy Perspective

信息量: $-\log p$

熵: $E[I - \log p] = \int p(x) \cdot \log p(x) dx$
 若离散: $= - \sum p(x) \cdot \log p(x)$

$$H[P] = - \sum p(x) \log p(x)$$

假设 x 是离散的

最大熵 \Leftrightarrow 等可能 \rightarrow 对某事件求导 (因为 p_i 是离散型)

$$\frac{\partial L}{\partial p_i} = \log p_i + 1 - \lambda \triangleq 0$$

$$\hat{p}_i = \exp(\lambda - 1)$$

$\therefore \hat{p}_i = \text{constant}$

$$\therefore \hat{p}_1 = \hat{p}_2 = \dots = \hat{p}_k = \frac{1}{k}$$

$\therefore p(x)$ 是均匀分布

在没有任何已知情况下, 熵达到最大的分布是均匀分布

| | | | | |
|-----|------------------------|---|---------|-----|
| x | 1 | 2 | \dots | k |
| p | p_1, p_2, \dots, p_k | | | |

$$\sum_{i=1}^k p_i = 1, \min \sum_{i=1}^k p_i \log p_i$$

$$\max H[P] = \max - \sum_{i=1}^k p_i \log p_i, p = \left(\frac{p_1}{p_2}, \frac{p_2}{p_3}, \dots, \frac{p_k}{p_1} \right)$$

$$\therefore p(x)$$

$$\hat{p}_i = \operatorname{argmax}_k H[P], \text{ Lagrangian Function}$$

$$\hat{p}_i = \operatorname{argmin} \sum_{i=1}^k p_i \log p_i$$

$$\therefore L(\eta, \lambda) = \sum_{i=1}^k p_i \log p_i + \lambda \left(1 - \sum_{i=1}^k p_i \right)$$

7.6 Maximum Entropy Perspective

$$P(x|y) = h(x) \exp\{y^T \phi(x) - A(y)\} = \frac{1}{Z(y)} h(x) \exp\{y^T \phi(x)\}$$

Data: $\{x_1, x_2, \dots, x_N\}$ 既定事实

$$\text{Empirical Distribution (经验分布)} \quad \hat{P}(x=x) = \hat{p}(x) = \frac{\text{count}(x)}{N}, \text{ e.g. } \hat{p}(x=x_1) = \frac{\text{count}(x_1)}{N}$$

$$E_{\hat{p}}[x], \text{Var}_{\hat{p}}[x]$$

$$f(x) \text{ 是任意关于 } x \text{ 的函数向量} \quad f(x) = \begin{pmatrix} f_1(x) \\ f_2(x) \\ \vdots \\ f_K(x) \end{pmatrix}, A = \begin{pmatrix} A_1 \\ A_2 \\ \vdots \\ A_K \end{pmatrix}$$

$$E_{\hat{p}}[f(x)] = A \rightarrow (\text{改写})$$

最大熵原理
满足已知事实约束

最大熵 \Leftrightarrow 等可能

$$P(x) \rightarrow P(x)$$

$$H[P] = - \sum_{x \in X} P(x) \log P(x)$$

则最大熵原理：

$$\left\{ \begin{array}{l} \min \sum_{x \in X} P(x) \log P(x) \\ \text{s.t. } \sum_{x \in X} P(x) = 1, E_p[f(x)] = E_p[f(x)] = \Delta^T f(x) \end{array} \right. \rightarrow P(x) \text{ 指数族分布}$$

这7.5多出的一件事，归宿事件

$$L(P, \lambda_0, \lambda) = \sum_x P(x) \log P(x) + \lambda_0 (1 - \sum_x P(x)) + \lambda^T (\Delta - E_p[f(x)]) \rightarrow \sum_x P(x) f(x)$$

$$\frac{\partial L}{\partial P(x)} = \sum_x [\log P(x) + 1] + -\sum_x \lambda_0 + -\sum_x \lambda^T f(x) \triangleq 0$$

$$\underbrace{\sum_x [\log P(x) + 1]}_0 - \lambda_0 - \lambda^T f(x) = 0$$

$$\log P(x) + 1 - \lambda_0 - \lambda^T f(x) = 0$$

$$\log P(x) = \lambda^T f(x) + \lambda_0 - 1$$

这就是指数族分布

$$\gamma = \begin{pmatrix} \lambda_0 \\ \lambda \end{pmatrix} \rightarrow \text{自然参数}$$

$$\phi(x) = \begin{pmatrix} 1 \\ f(x) \end{pmatrix} \rightarrow \text{sufficient statistic}$$

$$A(\gamma) = \dots \rightarrow \text{partition function}$$

$\therefore P(x)$ 是满足指数族分布

作用：「数据集不知道是什么分布时，用指数族分布 (Gauss Distribution)」

与 7.5 的区别：约束件不同，
最大熵约束件

均匀分布 仅限制取值范围 ($x \in [a, b]$)，
其它统计量 (μ, σ^2) 约束

指数族分布 存在矩约束 (如 $E[x] = \mu$ 或
更一般 $E[f_i(x)] = \mu_i$)

“在信息约束下，做最无偏的概率推断”

“有约束但信息不足，且需保持中立、不添假设”时用最大熵

e.g. 已知随机变量 均值 / 方差，但不知分布类型，正态分布是最大熵时。

8. Probabilistic Graphical Models

Bayesian Network
Gaussian BN

Gaussian MN

Markov Network

精确推断

Monte Carlo (MC) MC

随机近似

完备数据

无向

隐变量 EM

结构学习

Inference — 推断

Learning — 学习

Learning — 学习

Learning — 学习

Learning — 学习

8.1 PGM Background

高维随机变量 $P(x_1, x_2, \dots, x_p)$

$$\text{sum Rule: } P(x_1) = \int P(x_1, x_2) dx_2$$

$$\text{Product Rule: } P(x_1, x_2) = P(x_1) \cdot P(x_2 | x_1) = P(x_1) \cdot P(x_2 | x_1)$$

$$\text{chain Rule: } P(x_1, x_2, \dots, x_p) = \prod_{i=1}^p P(x_i | x_1, x_2, \dots, x_{i-1})$$

$$\text{Bayesian Rule: } P(x_1 | x_2) = \frac{P(x_1, x_2)}{P(x_2)} = \frac{P(x_1, x_2)}{\int P(x_1, x_2) dx_2} = \frac{P(x_2 | x_1) P(x_1)}{\int P(x_2) P(x_1 | x_2) dx_2}$$

困境：维度高，计算复杂， $P(x_1, x_2, \dots, x_p)$ 计算太大

$$\downarrow \text{(naive Bayes)} P(y|x) = \prod_{i=1}^p P(x_i | y)$$

$$\text{简化相加独立, } P(x_1, x_2, \dots, x_p) = \prod_{i=1}^p P(x_i)$$

独立现在，将来过去互相独立

Markov Property $x_j \perp x_{i+1} | x_i, j < i$

条件独立性

$x_A \perp x_B \perp x_C$
 x_A, x_B, x_C 是集合且不相交

P25 (HMM) 部分 Markov 假设

8.2. Bayesian Network - condition independence

$$P(x_1, x_2, \dots, x_p) = P(x_1) \cdot \prod_{i=2}^p P(x_i | x_{1:i-1}) \rightarrow \text{chain Rule}$$

条件独立性: $x_A \perp x_C | x_B$

$x_{pa(i)}$ 是 x_i 的父集

$$\text{因子分解: } P(x_1, x_2, \dots, x_p) = \prod_{i=1}^p P(x_i | x_{pa(i)})$$

有向图

检验图中的关系时, 使用 chain Rule
以及因子分解

拓扑排序



因子分解

$$P(a, b, c) = P(a) P(b|a) P(c|a)$$

$$P(a, b, c) = \underbrace{P(a) \cdot P(b|a)}_{\text{chain Rule}} \cdot P(c|a, b)$$

$$\Rightarrow P(c|a) = P(c|a, b)$$

$$\Rightarrow c \perp b | a$$

即若 a 被观测, 则路径被阻塞, 即

② $\overbrace{\quad \quad \quad}^{tail} \quad \overbrace{\quad \quad \quad}^{head}$ b, c 独立

$$\stackrel{\text{head}}{a} \rightarrow \stackrel{\text{tail}}{b} \rightarrow \stackrel{\text{head}}{c} \Rightarrow \text{head to tail}$$

$a \perp c | b$ 若 b 被观测, 则路径被阻塞

③ $\overbrace{\quad \quad \quad}^{head} \quad \overbrace{\quad \quad \quad}^{head}$ $\Rightarrow \text{head to head}$

默认情况下, $a \perp b$, 若 c 被
路径是阻塞的, 若 c 被观测, 则路径是通的

$$P(a, b, c) = P(a) P(b) \cdot P(c|a, b) \text{ 因子分解} \Rightarrow P(b) = P(b|a)$$

$$P(a, b, c) = P(a) P(b|a) P(c|a, b) \text{ chain Rule} \Rightarrow a \perp b$$

$$P(x_i | x_{-i}) = \frac{P(x_i, x_{-i})}{P(x_{-i})} = \frac{P(x)}{\int_{x_{-i}} P(x) dx_{-i}} = \frac{\prod_{j \neq i} P(x_j | x_{pa(j)})}{\int_{x_{-i}} \prod_{j \neq i} P(x_j | x_{pa(j)}) dx_{-i}}$$

与 x_i 相关 Δ \rightarrow 在分母中, 与 x_i 无关的 Δ
与 x_i 无关 Δ 可以被提出积分号, 而前部分也可以提出一个 Δ , 二者
相消, 最后就变为与 x_i 相关的函数, 即

$$P(x_i | x_{-i}) = f(\bar{x})$$

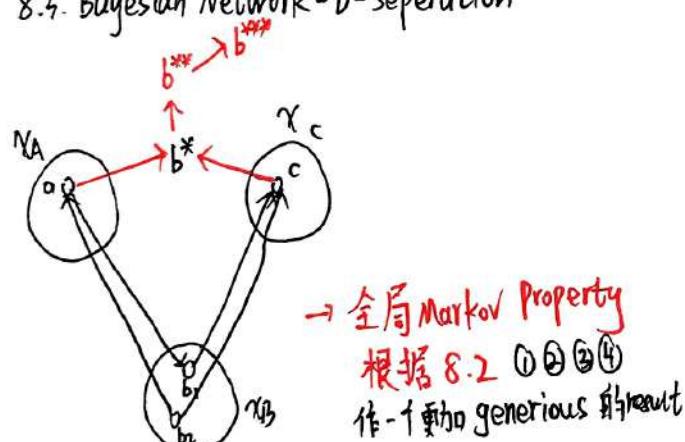
$$\frac{P(x_i | x_{pa(i)})}{P(x_{child(i)} | x_i, x_{parent(child(i))})}$$

双亲

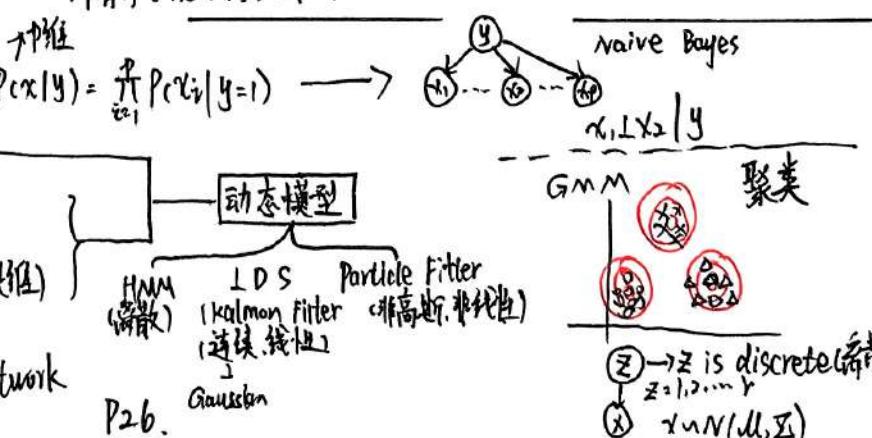
Markov Blanket

即在 Marcor Network 中, 一个元素 x_i 与所有网络元素之间的关系, 第
所有与它相关的, 即 x_i 的 Marcor Blanket.

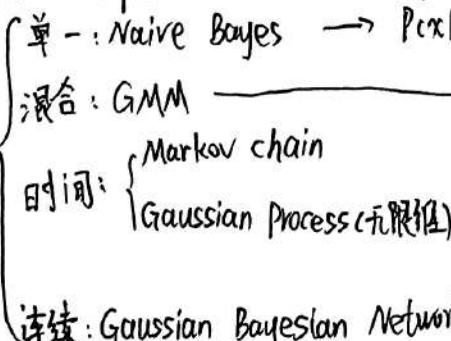
8.3. Bayesian Network - D-Separation



→ 全局 Markov Property
根据 8.2 ① ② ③ ④
做一个更通用的 result



8.4. Bayesian Network - example



8.5. Markov Random Field - Representation - 条件独立性

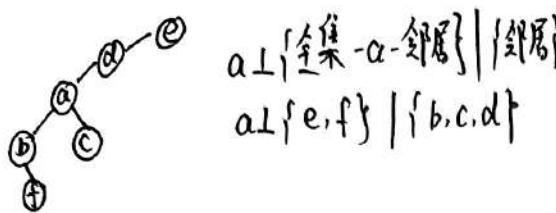
pairwise

① $X_A \perp X_C | X_B$
Global Markov

② 成对 Markov

$$X_i \perp X_j | X_{-i-j} \quad (i \neq j)$$

③ Local Markov



factorization:

图：
团：一个关于某点的集合，集合中的节点之间相互都是连通的
最大团：无法加入新节点，加入后就不是团

条件独立性体现在三方面：①全局 ②局部 ③成对
且 ① \Leftrightarrow ② \Leftrightarrow ③ 相互等价

概率分布、归一化因子

$$P(x) = \frac{1}{Z} \prod_{i=1}^K \phi(X_{ci})$$

$$Z = \sum_x \prod_{i=1}^K \phi(X_{ci})$$

$$= \sum_{\forall x_1 \dots x_K} \prod_{i=1}^K \phi(X_{ci})$$

需证明二者等价
才能说明这是一个概率率图的 factorization

8.6. Markov Random Field - Representation - Factorization

D-separation

基于：条件独立性
① Global Markov Property: $X_A \perp X_C | X_B$ 如果 A, B, C , sep $(A, C | B)$, 那么 $X_A \perp X_C | X_B$

② Local Markov Property: $X_i \perp X_{-i-nb(i)} | X_{nb(i)}$, $nb(i)$: neighborhood of node i

③ Pairwise Markov property: $X_i \perp X_j | X_{-i-j}$

① \Leftrightarrow ② \Leftrightarrow ③

factorization: 无向图，有向图 P26.

$$P(x) = \frac{1}{Z} \prod_{i=1}^K \psi(X_{ci}),$$

C_i : 最大团

X_{ci} : 最大团随机变量集合

$$Z = \sum_x \prod_{i=1}^K \psi(X_{ci}) = \sum_{x_1} \sum_{x_2} \dots \sum_{x_K} \prod_{i=1}^K \psi(X_{ci})$$

$\psi(X_{ci})$: 势函数，必须为正

如果一个概率分布，可以写成基于团上的因子分解，根据 Hammettley-clifford 可以证明它是马尔可夫随机场：① \Leftrightarrow ② \Leftrightarrow ③ \Leftrightarrow 因子分解(基于最大团)

$$\text{势 } \psi(X_{ci}) = \exp\{-E(X_{ci})\} > 0$$

Energy Function

$\rightarrow P(x)$ 称为 Gibbs Distribution / Boltzmann Distribution

形式上与指数族分布相似

$$P(x) = \frac{1}{Z} \prod_{i=1}^K \psi(X_{ci}) = \frac{1}{Z} \prod_{i=1}^K \exp\{-E(X_{ci})\}$$

$$= \frac{1}{Z} \exp\left\{-\sum_{i=1}^K E(X_{ci})\right\}$$

指数族分布 $\rightarrow P(x) = h(x) \cdot \exp\{y^T \phi(x) - A(y)\}$

$$= \frac{1}{Z(y)} h(x) \exp\{y^T \phi(x)\}$$

最大熵原理 \Rightarrow 指数族分布 (Gibbs 分布)

Markov Random Field \Leftrightarrow Gibbs Distribution

8.7 Inference - Introduction

Inference: 求概率: $P(x) = P(X_1, X_2, \dots, X_p)$

边缘概率: $P(x_i) = \sum_{X_1} \sum_{X_2} \dots \sum_{X_{i-1}} \sum_{X_{i+1}} \sum_{X_p} P(x)$, 除以对其余变量的所有取值求和

条件概率: $P(X_A | X_B)$

$$Y = X_A \cup X_B$$

MAP Inference: $P(z|x)$

$$\hat{z} = \arg \max_z P(z|x) \propto \arg \max_z P(z, x)$$

精确推断 {
 variable elimination (VE) ↗
 Belief Propagation (BP) ↗ Sum-Product Algorithm (针对树结构)
 Junction Tree Algorithm ↗ (普通图结构)

Inference {
 近似推断 {
 Loop Belief Propagation (有环图)
 Monte Carlo Inference: Importance Sampling, MCMC
 variational Inference:

8.8. Inference - Variable Elimination (乘法分配律)

Tasks: Inference (给定 $P(x) = P(x_1, x_2, \dots, x_p)$)

$$\text{边缘概率: } P(x_i) = \sum_{x_1, x_2, \dots, x_i, \dots, x_p} P(x_1, x_2, \dots, x_p)$$

$$\text{条件概率: } P(x_A | x_B)$$

$$\text{MAP } \hat{x}_A = \arg \max_{x_A} P(x_A | x_B) = \arg \max_{x_A} P(x_A, x_B)$$

$$P(x) = \prod_{x_i} \phi_i(x_i)$$

缺点: ① 重复计算 ② ordering \rightarrow NP-hard

(假设 a, b, c, d 均是离散二值变量)
 $r.v. a, b, c, d \in \{0, 1\}$

$$P(d) = \sum_{a,b,c} P(a, b, c, d) = \sum_{a,b,c} P(a) \cdot P(b|a) \cdot P(c|b) \cdot P(d|c)$$

$$= P(a=0) \cdot P(b=0|a=0) \cdot P(c=0|b=0) \cdot P(d=0|c=0)$$

$$+ P(a=1) \cdot P(b=0|a=1) \cdot P(c=0|b=0) \cdot P(d=0|c=0)$$

+

$$+ P(a=1) \cdot P(b=1|a=1) \cdot P(c=1|b=1) \cdot P(d=0|c=1)$$

= &·因子积

$\Rightarrow d=1$ 时 = &·因子积

$$= \sum_{b,c} P(c|b) \cdot P(d|c) \cdot \sum_a P(a) \cdot P(b|a) \xrightarrow{\sum_a \phi_a(b)} P(b)$$

$$= \sum_c P(d|c) \cdot \sum_b P(c|b) \cdot \underbrace{\phi_a(b)}_{\phi_b(c)}$$

$$= \phi_c(d) \quad (\text{乘法对加法分配律})$$

8.9. Inference - Belief Propagation

variable Elimination: (乘法分配律)



$$P(e) = P(a, b, c, d, e) = P(a) \cdot P(b|a) \cdot P(c|b) \cdot P(d|c) \cdot P(e|d)$$

$$P(e) = \sum_{a,b,c,d} P(a, b, c, d, e) = \sum_d P(e|d) \cdot \sum_c P(d|c) \cdot \sum_b P(c|b) \cdot \sum_a P(b|a) P(a)$$

$$P(e) = \sum_{a,b,d,e} P(a, b, d, e) = (\sum_b P(c|b) \cdot \sum_a P(b|a) P(a)) \cdot (\sum_d P(d|e) \sum_e P(e|d))$$

为解决重复计算, 引入 Belief Propagation:

$$P(a) = \sum_{b,c,d} P(a, b, c, d) \quad \text{结合①与图①, 从 c, d 往上}$$

$$= \sum_c \psi_c \cdot \psi_{bc}$$

$$= \sum_d \psi_d \cdot \psi_{bd}$$

$$= \psi_{ab}$$

$$= P(a)$$

Forward Algorithm

BP: Chain-Tree:

有向-无向

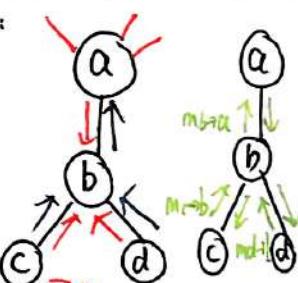
$$\nexists P(a)$$

(4个点, 3条边)

$$P(a, b, c, d) = \frac{1}{4} \psi_a(a) \cdot \psi_b(b) \cdot \psi_c(c) \cdot \psi_d(d) \cdot \phi_{ab}(b)$$

$$\psi_{bc}(b, c) \cdot \psi_{cd}(c, d) \dots \text{①}$$

$$P(a) = \sum_{b,c,d} P(a, b, c, d), P(b) = \sum_{a,c,d} P(a, b, c, d)$$



若把绿线全求出, $P(a), P(b), P(c), P(d)$ 自然也就出来了

$m_{mb \rightarrow a}(a) = \sum_b \psi_{ab} \psi_b \cdot m_{c \rightarrow b}(b) \cdot m_{d \rightarrow b}(b)$
 b 的除 a 外的 neighbor

$$P(a) = \psi_a m_{b \rightarrow a}(a) \quad \text{NB: neighbor}$$

$$m_{ij \rightarrow i}(i) = \sum_j \psi_{ij} \psi_j \prod_{k \in N(i) \setminus j} m_{k \rightarrow j}(j) \quad m_{k \rightarrow j}(j)$$

$$P(x_i) = \psi_i \prod_{k \in N(i)} m_{k \rightarrow i}(x_i)$$

如图中 a 可能有 neighbor
 Belief propagation

8.109 Inference - Belief propagation

$$M_{j \rightarrow i} = \sum_j \psi_{ij} \psi_j^T M_{k \rightarrow j}$$

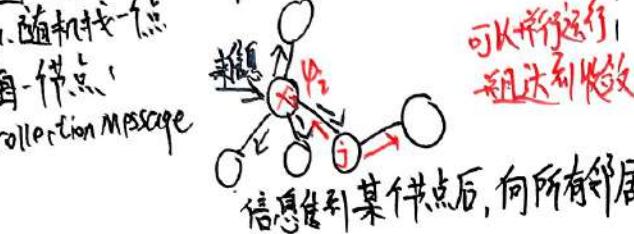
$$m_{b \rightarrow a} = \sum_b \psi_{ab} \psi_b M_{c \rightarrow b} M_{d \rightarrow b}$$

self children
belief (b)

$$\text{belief}(b) = \psi_b \cdot \text{children}$$

$$m_{b \rightarrow a} = \sum_b \psi_{ab} \cdot \text{belief}(b)$$

② Parallel Implementation 分布式



方法 II BP = VE + Caching

直接求 $M_{ij} \Rightarrow P(x_i)$
图的遍历

BP (Sequential Implementation)

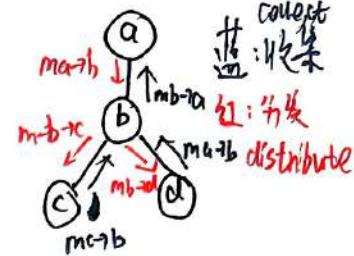
① Get root, assume a is root

② collect Message \rightarrow for x_i in $NB(a)$:
collect $\text{Msg}(x_i)$

③ distribute Message \rightarrow for x_j in $NB(a)$:
 $\text{distribute}(x_j)$

可得 M_{ij} for all $i, j \in V$

从而 $P(x_k), k \in V$



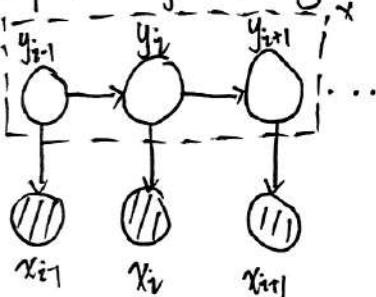
8.110 Inference - Max Product Algorithm

Graph: (X, E)

(I) 边缘概率: $E = \{e_1, e_2, \dots, e_k\}, P(E)$ (likelihood)

(II) 条件概率: $P(Y|E)$ (Posterior) $X = (Y, Z) P(Y|E) = \sum_Z P(X|E)$

(III) MAP (Decoding): $\hat{Y} = \arg \max_X P(X|E) \quad \hat{Y} = \arg \max_Y P(Y|E) = \arg \max_Y \sum_Z P(X|E)$



Decoding: $\hat{Y} = \arg \max_Y P(Y|X)$

viterbi Algorithm: 动态规划问题

Belief Propagation (Tree) Sum-product

$$m_{j \rightarrow i}(\gamma_i) = \prod_{k \in NB(i)-j} \psi_{kj}(\gamma_k) \cdot \psi_{ij}(\gamma_i, \gamma_j) \prod_{k \in NB(j)} m_{k \rightarrow j}(\gamma_j)$$

$$P(x_i) = \psi_i(x_i) \cdot \prod_{k \in NB(i)} m_{k \rightarrow i}(\gamma_i)$$

Max-product:

① BP 改进; ② viterbi 的推广

此处
 $m_{b \rightarrow a}$ 是经过了到
 $m_{c \rightarrow b}$ 下一个点的最
大概率

如 $m_{b \rightarrow c}$ 路过 c, d 后,
从 b 到 a 的最大概率
即 c, b, d 联合概率最大

MAX-product

$$m_{j \rightarrow i} = \max_{\gamma_j} \psi_j \cdot \psi_{ij} \cdot \prod_{k \in NB(i)-j} m_{k \rightarrow j}$$

$$m_{c \rightarrow b} = \max_{\gamma_c} \psi_c \cdot \psi_{cb}$$

$$m_{d \rightarrow b} = \max_{\gamma_d} \psi_d \cdot \psi_{bd}$$

与 c 相关的所有概率
无关圆无用
关于 c 的函数

与 d 相关的所有概率
无关圆无用
关于 d 的函数

关于 b 的函数

从底向上可求出 $\max_{\gamma_b} P(x_a, x_b, x_c, x_d | E)$

从上到底回溯可求出 $x_a^*, x_b^*, x_c^*, x_d^*$

$$(x_a^*, x_b^*, x_c^*, x_d^*) = \arg \max_{x_a, x_b, x_c, x_d} P(x_a, x_b, x_c, x_d | E)$$

x_a, x_b, x_c, x_d

$$x_d^* = \arg \max_{x_d} \psi_d \cdot \psi_{bd}$$

$$x_c^* = \arg \max_{x_c} \psi_c \cdot \psi_{bc}$$

$$x_b^* = \arg \max_{x_b} \psi_b \cdot \psi_{ab} \cdot m_{c \rightarrow b} \cdot m_{d \rightarrow b}$$

$$x_a^* = \arg \max_{x_a} \psi_a \cdot m_{b \rightarrow a}$$

$$\max_{x_a} P(x_a | x_b, x_c, x_d) = \max_{x_a} \psi_a \cdot m_{b \rightarrow a}$$

关于 a 的函数

回溯:

$(x_a^*, x_b^*, x_c^*, x_d^*)$

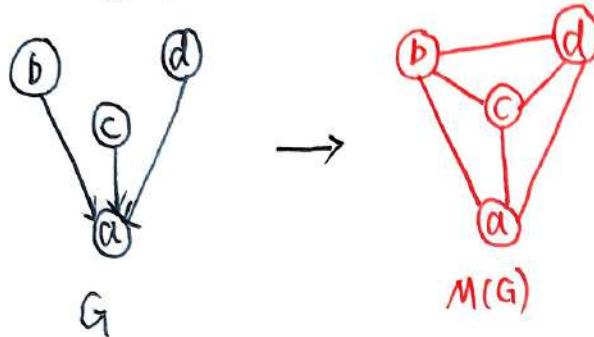
8.1. Supplement - Graph

Moral

原因：想把有向图转为无向图来研究。

有向图：head-to-head (V结构)

作法：① $\forall X_i \in G$, 将 Parent(X_i) 两两连接
② 将 G 中的有向别替换成无向边



$$\text{Sep}(A, B | C) \Leftrightarrow \text{D-Sep}(A, B | C)$$

无向图中

有向图

① head to tail

$$a \rightarrow b \rightarrow c \rightarrow \text{tail} \quad \rightarrow \quad \text{团 } \phi(a,b) \quad \text{团 } \phi(b,c)$$

$$P(a, b, c) = P(a) \cdot P_{cb}(a) \cdot P_{cc}(b)$$

$$\phi(a, b) \quad \phi(b, c)$$

② tail to tail

$$a \leftarrow b \leftarrow c \rightarrow \text{tail} \quad \rightarrow \quad \text{团 } \phi(a,b) \quad \text{团 } \phi(b,c)$$

$$P(a, b, c) = P(a) \cdot P_{cb}(a) \cdot P_{cc}(a)$$

$$\phi(a, b) \quad \phi(a, c)$$

③ head to head

$$a \leftarrow c \leftarrow b \rightarrow \text{tail} \quad \rightarrow \quad \text{团 } \phi(a,b,c)$$

$$P(a, b, c) = P(a) P(b) \cdot P_{cc}(a) \cdot P_{cb}(b) = P(a) P(b) \underbrace{P_{cc}(a, b)}_{\phi(a, b, c)}$$

按①②做法, (a, b, c) 显然不为一个团

8.2. Supplement - Factor Graph

有向图： $P(x) = \prod P(x_i | \pi(x_i))$

无向图： $P(x) = \frac{1}{Z} \prod_{i=1}^k \phi_{ci}(x_{ci})$, 最大团集合

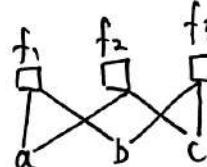
道德图：有向图 \rightarrow 无向图 (Tree-like Graph)
(树) (树) (引入环)

① BP 只能对树操作, 因子图可以去环

② 简便

因式分解本身对应一个特殊的因子图

因子图：看作是对因式分解的进一步分解



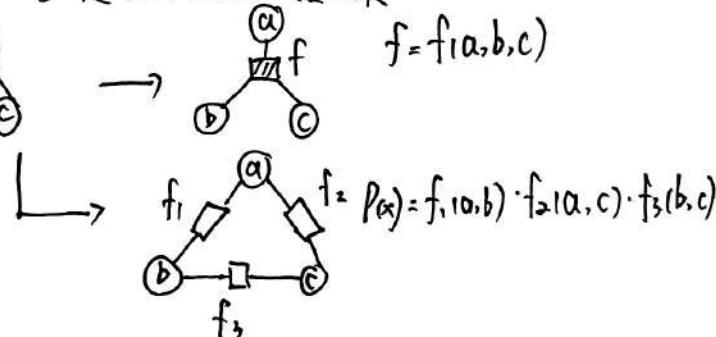
因子图： $x = x_i; \dots; x_p$

$$P(x) = \prod f_s(x_s)$$

S : 图的节点子集

x_s 是 x 的随机变量子集

$$f = f(a, b, c)$$



9. EM Algorithm

9.1 EM Introduction

可以想到解的例子：MLE: $P(x|\theta)$

$$\theta_{MLE} = \arg \max_{\theta} \log P(x|\theta)$$

log-likelihood

$$\text{EM: } \theta^{(t+1)} = \arg \max_{\theta} f_Z \underbrace{\log P(x, z|\theta)}_{E_{Z|x, \theta^{(t)}} [\log P(x, z|\theta)]} \cdot P(z|x, \theta^{(t)}) dz$$

收敛性

$$\log P(x|\theta^{(t)}) \leq \log P(x|\theta^{(t+1)}) \dots \text{①}$$

EM: 最大化对数似然，保证下一步都比上一步大。

能达到极值

$$\text{证①式: } P(x, z|\theta) = P(z|x, \theta) \cdot P(x|z)$$

$$P(x|\theta) = \frac{P(x, z|\theta)}{P(z|x, \theta)}$$

$$\log P(x|\theta) = \log \left(\frac{P(x, z|\theta)}{P(z|x, \theta)} \right)$$

$$= \log P(x, z|\theta) - \log P(z|x, \theta)$$

$$\therefore \log P(x|\theta) = \log P(x, z|\theta) - \log P(z|x, \theta)$$

$$\text{左边} = \int_z P(z|x, \theta^{(t)}) \cdot \log P(x|\theta) dz$$

$$= \log P(x|\theta) \cdot \int_z P(z|x, \theta^{(t)}) dz$$

$$= \log P(x|\theta)$$

对任意随机变量 X (满足 $X > 0$ 且 $E[X]$ 存在), 若 f 是凸函数,
则: $f(E[X]) \geq E[f(X)]$

或由 Jensen 不等式:

concave 凸 (凹) convex (凸)

$\int_Z P(z|X, \theta^{(t)}) \cdot \log \frac{P(z|X, \theta^{(t)})}{P(z|X, \theta^{(t-1)})} dz$

$E[\log X] \leq \log E[X], f(E[X]) \geq E[f(X)]$

$f(x) = \log x$

$\leq \log \int_Z P(z|X, \theta^{(t+1)}) dz = \log 1 = 0$

$\therefore H(\theta^{(t+1)}, \theta^{(t)}) - H(\theta^{(t)}, \theta^{(t)}) \leq 0.$

$\therefore H(\theta^{(t+1)}, \theta^{(t)}) \leq H(\theta^{(t)}, \theta^{(t)})$

由 $P_3.0, \theta^{(t+1)} = \text{argmax}_{\theta} \dots$ 得出 $Q(\theta^{(t)}, \theta^{(t)})$

是变量 θ 行的依据, 欲证 $H(\theta^{(t+1)}, \theta^{(t)}) \leq H(\theta^{(t)}, \theta^{(t)})$ 自然成立 ①
然成立 只证: $H(\theta^{(t+1)}, \theta^{(t)}) - H(\theta^{(t)}, \theta^{(t)}) \leq 0$

也有帮助 $= \int_Z P(z|X, \theta^{(t)}) \log P(z|X, \theta^{(t+1)}) dz - \int_Z P(z|X, \theta^{(t)}) \log P(z|X, \theta^{(t)}) dz$

$= \int_Z P(z|X, \theta^{(t)}) \cdot \log \frac{P(z|X, \theta^{(t+1)})}{P(z|X, \theta^{(t)})} dz$, kullback-Leiber divergence

$= -KL(P(z|X, \theta^{(t)}) || P(z|X, \theta^{(t+1)}))$

$\leq 0 \quad \geq 0$

$\therefore H(\theta^{(t+1)}, \theta^{(t)}) \leq H(\theta^{(t)}, \theta^{(t)})$

注意变量替换: $\log P(x|\theta) \rightarrow \log P(x|\theta^{(t)}) = Q(\theta^{(t)}, \theta^{(t)}) - H(\theta^{(t)}, \theta^{(t)}) \leq \log P(x|\theta^{(t+1)})$

由 ① ② $\rightarrow Q(\theta^{(t)}, \theta^{(t)}) - H(\theta^{(t)}, \theta^{(t)}) \leq Q(\theta^{(t+1)}, \theta^{(t)}) - H(\theta^{(t+1)}, \theta^{(t)})$

故 $\log P(x|\theta^{(t)}) \leq \log P(x|\theta^{(t+1)})$ 收敛

上述推导并不严密: $P_3.0$ 中 $\log P(x|\theta)$ 中是取 θ , 因为我们证 ① 时写的是 $P(x, z|\theta) = P(z|\theta) \cdot P(x|z)$, 若把 θ 改为 $\theta^{(t)}$
则有 $\log P(x|\theta^{(t)}) = \log P(x, z|\theta^{(t)}) - \log P(z|\theta^{(t)})$ 之后证明关于变量替换部分就明显了, 我用粗笔画出

9.2 ELBO + KL Divergence

MLE: $\theta_{MLE} = \log P(x|\theta)$

x : observed data $x = \{x_1, \dots, x_n\}$ $x_i \sim P_{(x_i)}, P_{(x_i)} = \prod_{j=1}^n P_{(x_{ij})}$
 z : unobserved data (latent variable)

(x, z) : complete data

θ : parameter

E-step: $P(z|x, \theta^{(t)}) \rightarrow E_{z|x, \theta^{(t)}}[\log P(x, z|\theta)]$

M-step: $\theta^{(t+1)} = \arg \max_{\theta} E_{z|x, \theta^{(t)}}[\log P(x, z|\theta)]$

收敛性: $\log P(x|\theta^{(t)}) \leq \log P(x|\theta^{(t+1)})$

下面证明 EM:

$\log P(x|\theta) = \log P(x, z|\theta) - \log P(z|x, \theta)$

$= \log \frac{P(x, z|\theta)}{q(z)} - \log \frac{P(z|x, \theta)}{q(z)}, q(z) \neq 0$

同时对左右两边求期望
左边: $\int_Z q(z) \log P(x|\theta) dz = (\log P(x|\theta)) \cdot \int_Z q(z) dz = \log P(x|\theta)$

右边: $\int_Z q(z) \cdot \log \frac{P(x, z|\theta)}{q(z)} dz - \int_Z q(z) \log \frac{P(z|x, \theta)}{q(z)} dz$

ELBO = evidence lower bound

$KL(q(z) || P(z|x, \theta))$

对数似然的下界

故有: $\log P(x|\theta) = \text{ELBO} + KL(q||P) \rightarrow \text{Posterior}$
且 $KL(q||P) \geq 0 \Leftrightarrow q = P \quad KL(q(z) || P(z|x, \theta))$

即: $\log P(x|\theta) \geq \text{ELBO} \rightarrow \text{下界}$

注: ELBO 尽量大, 则 $\log P(x|\theta)$ 也会大.

$\hat{\theta} = \arg \max_{\theta} \text{ELBO}$

$= \arg \max_{\theta} \int_Z q(z) \log \frac{P(x, z|\theta)}{q(z)} dz$

上次选 $q(z) = P(z|x, \theta)$ 时, $\log P(x|\theta) = \text{ELBO}$. (θ 是自变量, parameter, 令为 $\theta^{(t)}$)

$= \arg \max_{\theta} \int_Z P(z|x, \theta^{(t)}) [\log P(x, z|\theta) - \log P(z|x, \theta^{(t)})] dz$

$\therefore \hat{\theta} = \arg \max_{\theta} \int_Z \log P(x, z|\theta) \cdot P(z|x, \theta^{(t)}) dz$

$\downarrow \theta^{(t+1)}$ Variable constant

9.3. ELBO + Jensen's Inequality

$$q = q(z)$$

$$P = P(z|x, \theta)$$

$$\log P(x|\theta) = \text{ELBO} + \text{KL}(q||P)$$

$$\log P(x|\theta) = \log \int_z P(x,z|\theta) dz$$

$$= \log \int_z \frac{P(x,z|\theta)}{q(z)} \cdot q(z) dz$$

$$= \log \mathbb{E}_{q(z)} \left[\frac{P(x,z|\theta)}{q(z)} \right]$$

$$\geq \mathbb{E}_{q(z)} \left[\log \frac{P(x,z|\theta)}{q(z)} \right] \rightarrow \text{ELBO}$$

$$\text{"="} \Leftrightarrow \frac{P(x,z|\theta)}{q(z)} = C$$

$$q(z) = \frac{1}{C} P(x,z|\theta)$$

$$1 = \int_z q(z) dz = \int_z \frac{1}{C} P(x,z|\theta) dz$$

$$= \frac{1}{C} \int_z P(x,z|\theta) dz$$

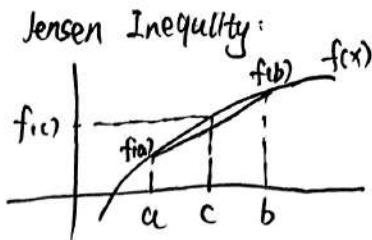
$$1 = \frac{1}{C} P(x|\theta)$$

$$C = P(x|\theta)$$

$$\text{则 } q(z) = \frac{1}{P(x|\theta)} \cdot P(x,z|\theta) = P(z|x, \theta) \rightarrow \text{Posterior}$$

我们随机抽入的 $q(z)$ 不是 $P(z|x, \theta)$

$$\text{则有 } \log P(x|\theta) \geq \mathbb{E}_{q(z)} \left[\log \frac{P(x,z|\theta)}{P(z|x, \theta)} \right] \Rightarrow \hat{\theta} = \arg \max_{\theta} \mathbb{E}_{q(z)} \left[\log \frac{P(x,z|\theta)}{P(z|x, \theta)} \right] \text{ 与 9.2 同理}$$



$f(x)$ = concave function

$$t \in [0, 1]$$

$$c = ta + (1-t)b$$

$$f(c) = f[ta + (1-t)b] \geq tf(a) + (1-t)f(b) \rightarrow \text{Jensen}$$

当 $t = \frac{1}{2}$ 时.

$$f\left[\frac{a+b}{2}\right] \geq \frac{f(a)}{2} + \frac{f(b)}{2}$$

$$f(E) \geq E[f]$$

9.4 EM Review 再回首

- ① 从狭义 EM \rightarrow 广义 EM
- ② 狹义 EM 是广义 EM 的一个特例
- ③ EM 变种

概率生成模型

$$P(x|\theta)$$

EM 就是求这个 θ

$$\hat{\theta} = \arg \max_{\theta} P(x|\theta)$$

$$x \sim P(x) \quad x = \{x_i\}_{i=1}^N$$

9.5. Generalized EM

$$\log P(x|\theta) = \underbrace{\text{ELBO}}_{\mathcal{L}(\theta, \theta)} + \text{KL}(q||P)$$

$$q = P(z)$$

$$\begin{cases} \text{ELBO} = \mathbb{E}_{q(z)} \left[\log \frac{P(x,z|\theta)}{q(z)} \right] \\ \text{KL}(q||P) = \int_z q(z) \cdot \log \frac{q(z)}{P(z|x, \theta)} dz \end{cases}$$

$$\text{熵: } \int_z q(z) \cdot \log \frac{1}{q(z)} dz$$

简化 θ

$$\begin{aligned} \mathcal{L}(\theta, \theta) &= \mathbb{E}_{q(z)} [\log P(x,z) - \log q(z)] = \mathbb{E}_q [\log P(x,z)] \\ &\quad - \underbrace{\mathbb{E}_q [\log q]}_{H[q]} \end{aligned}$$

$$\text{ELBO} = \mathbb{E}_{q(z)} [\log P(x,z|\theta)] + H[q]$$

故广义 EM 就是处理当 q 无法取到 P 的情况,

$q = q(z)$, $P = P(z|x, \theta)$, 而 q 越接近 P , 则 KL 就越小, ELBO 越大, 故共固定 θ 求最优 q , 在取了 q 求 $\hat{\theta}$, 并且都归到 $\mathcal{L}(\theta, \theta)$ 上进行计算

$$\text{固定 } \theta, \hat{\theta} = \arg \min_q \text{KL}(q||P) = \arg \max_q \mathcal{L}(\theta, \theta) \rightarrow \text{E-step}$$

$$\text{固定 } \hat{\theta}, \hat{\theta} = \arg \max_{\theta} \mathcal{L}(\hat{\theta}, \theta) \rightarrow \text{M-step}$$

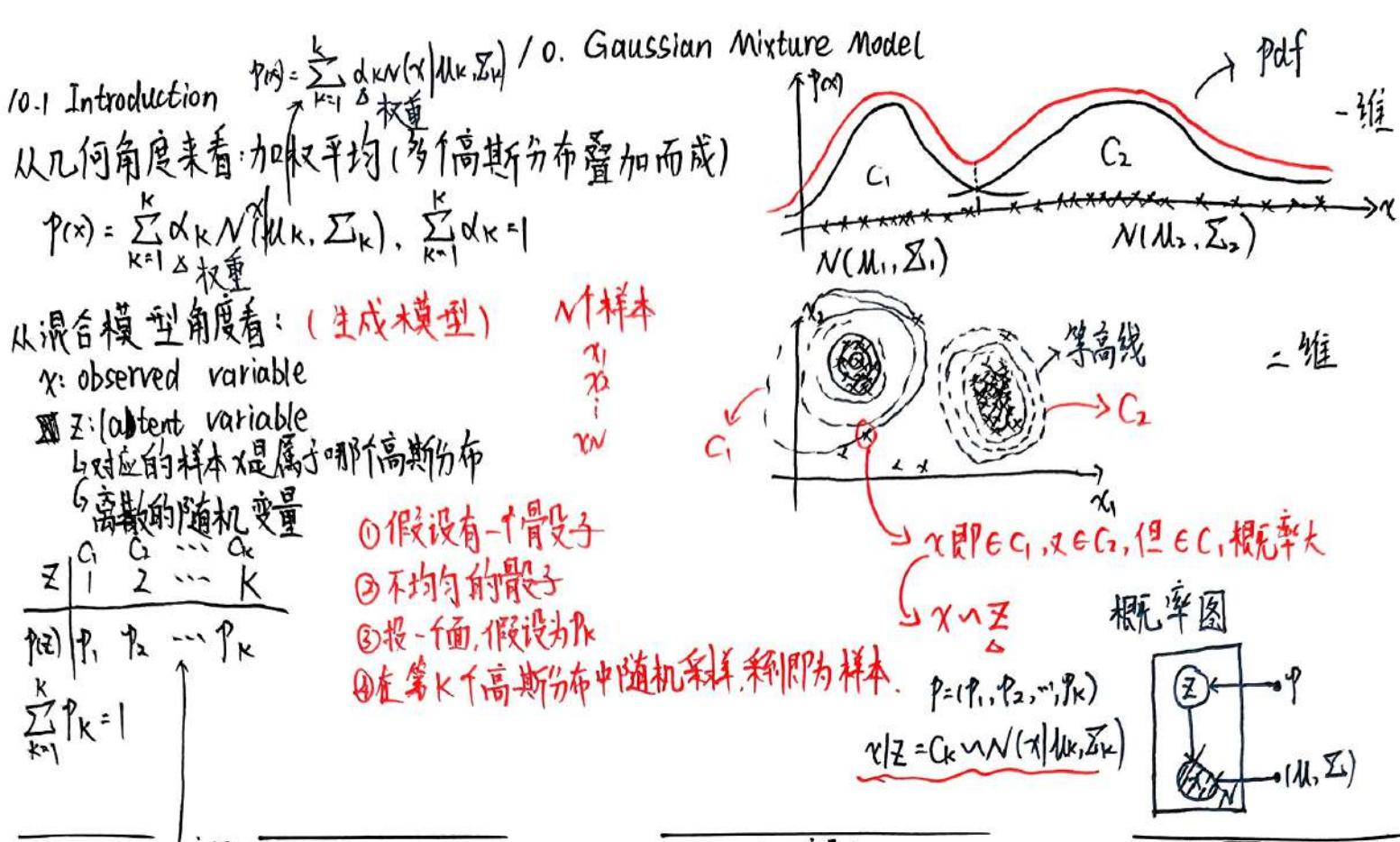
$$\arg \max_{\theta} \log P(x|\theta) = \arg \max_{\theta} \mathcal{L}(\theta, \theta)$$

$$\text{广义 EM: 不易计算 } \arg \max_{\theta} \mathcal{L}(\theta, \theta), \text{ 故 } \arg \max_{\theta} \text{KL}(q||P) = (\log P(x|\theta)) - \mathcal{L}(\theta, \theta)$$

$$\arg \max_{\theta} \mathcal{L}(\theta, \theta) \quad \arg \min_{\theta} \text{KL}(q||P) = \arg \max_{\theta} \mathcal{L}(\theta, \theta)$$

$$\begin{cases} \text{E-step: } q = \arg \max_q \mathcal{L}(\theta^{(t)}, \theta) \\ \text{M-step: } \theta^{(t+1)} = \arg \max_{\theta} \mathcal{L}(q^{(t)}, \theta) \end{cases}$$

梯度上升法
坐标上升法 (SMO)



10.2 MLE

由上节对于混合模型：

$$P(x) = \sum_z P(x, z) = \sum_{k=1}^K P(x, z=C_k) = \sum_{k=1}^K P(z=C_k) \cdot P(x | z=C_k)$$

Mixture model, z 下 x 的概率是 Gauss distribution

$$\sum_{k=1}^K P_k \cdot N(x | \mu_k, \Sigma_k)$$

$\hat{\theta}_{MLE} = \underset{\theta}{\operatorname{argmax}} \log P(x) = \underset{\theta}{\operatorname{argmax}} \log \prod_{i=1}^N P(x_i | \theta) = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^N \log P(x_i | \theta) = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^N \log \sum_{k=1}^K P_k \cdot N(x_i | \mu_k, \Sigma_k)$

直接用 MLE 求解 GMM 无法得出解析解
十分复杂，无法得到解析解
连加号十分困难

10.3 EM - E-Step

EM: $\theta^{(t+1)} = \operatorname{argmax}_{\theta} \underbrace{E_{z|x,\theta^{(t)}} [\log P(x, z|\theta)]}_{Q(\theta, \theta^{(t)})}$

$Q(\theta, \theta^{(t)}) = \sum_{z_1, z_2, \dots, z_N} \log \frac{1}{\prod_{i=1}^N P(x_i, z_i | \theta)} \cdot \prod_{i=1}^N P(z_i | x_i, \theta^{(t)}) dz$

$G(\theta, \theta^{(t)}) = \int_z \log P(x, z | \theta) \cdot P(z | x, \theta^{(t)}) dz$

$= \sum_z \underbrace{\log \frac{1}{\prod_{i=1}^N P(x_i, z_i | \theta)} \cdot \prod_{i=1}^N P(z_i | x_i, \theta^{(t)})}_{\text{原式: } Q(\theta, \theta^{(t)})} dz$

$= \sum_{z_1, z_2, \dots, z_N} \sum_{i=1}^N \log P(x_i, z_i | \theta) \cdot \prod_{i=1}^N P(z_i | x_i, \theta^{(t)})$

$= \sum_{z_1, z_2, \dots, z_N} [\log P(x_1, z_1 | \theta) + \log P(x_2, z_2 | \theta) + \dots + \log P(x_N, z_N | \theta)] \prod_{i=1}^N P(z_i | x_i, \theta^{(t)})$

对于: $\sum_{z_1, z_2, \dots, z_N} \log P(x_i, z_i | \theta) \cdot \prod_{i=1}^N P(z_i | x_i, \theta^{(t)})$ 第一项

$= \sum_{z_1, z_2, \dots, z_N} \log P(x_1, z_1 | \theta) \cdot P(z_1 | x_1, \theta^{(t)})$

$= \sum_{z_1, z_2, \dots, z_N} \log P(x_1, z_1 | \theta) \cdot \frac{1}{Q(\theta, \theta^{(t)})}$

P33

$$P(x) = \sum_{k=1}^K P_k \cdot N(x | \mu_k, \Sigma_k)$$

$$P(x, z) = P_z \cdot N(x | \mu_z, \Sigma_z)$$

$$P(z|x) = \frac{\sum_{k=1}^K P_k \cdot N(x | \mu_k, \Sigma_k)}{\sum_{k=1}^K P_k \cdot N(x | \mu_k, \Sigma_k)}$$

由 θ 是 parameter, θ 与 μ, Σ, P 都有关, 故代入式后 θ 可忽略
是第七个参数, 带上标 $P(z_i | x_i, \theta^{(t)})$

10.4 EM-M-Step

$$Q(\theta, \theta^{(t)}) = \sum_{z_i} \sum_{i=1}^N \log [P_{z_i} \cdot N(x_i | \mu_{z_i}, \Sigma_{z_i})] \cdot P(z_i | x_i, \theta^{(t)})$$

$$= \sum_{k=1}^K \sum_{i=1}^N \log [P_k \cdot N(x_i | \mu_k, \Sigma_k)] \cdot P(z_i = c_k | x_i, \theta^{(t)})$$

$$= \sum_{k=1}^K \sum_{i=1}^N [\log P_k + \log N(x_i | \mu_k, \Sigma_k)] \cdot P(z_i = c_k | x_i, \theta^{(t)})$$

$$\theta^{(t+1)} = \arg \max_{\theta} Q(\theta, \theta^{(t)}) - M\text{-step}$$

求 $P_k^{(t+1)}$: $P_k^{(t+1)} = \arg \max_{P_k} \sum_{k=1}^K \sum_{i=1}^N \log P_k \cdot P(z_i = c_k | x_i, \theta^{(t)})$, s.t. $\sum_{k=1}^K P_k = 1 \rightarrow$

$$\begin{cases} \max_{P_k} \sum_{k=1}^K \sum_{i=1}^N \log P_k \cdot P(z_i = c_k | x_i, \theta^{(t)}) \\ \text{s.t. } \sum_{k=1}^K P_k = 1 \end{cases}$$

$$P = (P_1, \dots, P_K)^T \quad L(P, \lambda) = \sum_{k=1}^K \sum_{i=1}^N \log P_k \cdot P(z_i = c_k | x_i, \theta^{(t)}) + \lambda (\sum_{k=1}^K P_k - 1).$$

对其中一事件 $\frac{\partial L}{\partial P_k} = \sum_{i=1}^N \frac{1}{P_k} \cdot P(z_i = c_k | x_i, \theta^{(t)}) + \lambda \stackrel{\text{与 } P_k \text{无关}}{=} 0$

Ck 求偏导

$$\Rightarrow \sum_{i=1}^N P(z_i = c_k | x_i, \theta^{(t)}) + P_k \lambda = 0 \rightarrow \text{通项}(k \text{ 事件})$$

有 $\sum_{i=1}^N P(z_i = c_k | x_i, \theta^{(t)}) + P_k \lambda = \dots = \sum_{i=1}^N P(z_i = c_k | x_k, \theta^{(t)}) + P_k \lambda = 0$

$$\Rightarrow \sum_{i=1}^N \sum_{k=1}^K P(z_i = c_k | x_i, \theta^{(t)}) + \sum_{k=1}^K P_k \lambda = 0$$

$$\because \sum_{k=1}^K P_k = 1, \sum_{k=1}^K P(z_i = c_k | x_i, \theta^{(t)}) = 1 \rightarrow \text{相当于积分}$$

$$\therefore \Rightarrow N + \lambda = 0$$

$$\lambda = -N \quad \text{这 } z_i = c_k \text{ 是 } P^{(t+1)} = (P_1^{(t+1)}, \dots, P_K^{(t+1)}) \text{ 时用}$$

$$P_k^{(t+1)} = \frac{1}{N} \sum_{i=1}^N P(z_i = c_k | x_i, \theta^{(t)}) \quad \begin{matrix} \uparrow \text{责任} \\ \text{只管这 } i \end{matrix} = \frac{1}{N} \sum_{i=1}^N \gamma_{ik}$$

$$\therefore P^{(t+1)} = (P_1^{(t+1)}, \dots, P_K^{(t+1)})$$

$$\text{②. } \mu_k^{(t+1)} = \frac{\sum_{i=1}^N P(z_i = c_k | x_i, \theta^{(t)}) \cdot x_i}{\sum_{i=1}^N P(z_i = c_k | x_i, \theta^{(t)})} ; \quad \text{③. } \Sigma_{ik}^{(t+1)} = \frac{\sum_{i=1}^N P(z_i = c_k | x_i, \theta^{(t)}) \cdot (x_i - \mu_k)(x_i - \mu_k)^T}{\sum_{i=1}^N P(z_i = c_k | x_i, \theta^{(t)})}$$

$$= \frac{\sum_{i=1}^N \gamma_{ik} \cdot x_i}{\sum_{i=1}^N \gamma_{ik}}$$

$$= \frac{\sum_{i=1}^N \gamma_{ik} \cdot (x_i - \mu_k)(x_i - \mu_k)^T}{\sum_{i=1}^N \gamma_{ik}}$$

11. Variational Inference

11.1 Introduction

频率角度 → 优化问题：

回归： $D = \{(x_i, y_i)\}_{i=1}^N, x_i \in \mathbb{R}^p, y_i \in \mathbb{R}$
 Model: $f(w) = w^T x$, loss function: $L(w) = \sum_{i=1}^N \|w^T x_i - y_i\|^2$
 解法：
 ① 解析解: $\frac{\partial L(w)}{\partial w} = 0 \Rightarrow w^* = (x^T x)^{-1} x^T y$
 ② 数值解: GD → Gradient Descent

SVM(分类): ① $f(w) = \text{sign}(w^T x + b)$; ② loss function: $\min_{w,b} \frac{1}{2} w^T w$ s.t. $y_i(w^T x_i + b) \geq 1 \quad i=1, \dots, N$ ⇒ 有约束, convex 优化
 ③ QP: Lagrange Dual

EM: $\hat{\theta} = \arg \max_{\theta} \log P(x|\theta) \Rightarrow \hat{\theta}^{(t+1)} = \arg \max_{\theta} \int_z \log P(x, z|\theta) \cdot P(z|x, \theta^{(t)}) dz$
 贝叶斯角度 → 积分问题:
 $P(\theta|x) = \frac{P(x|\theta) \cdot P(\theta)}{P(x)} \rightarrow \text{先验}$
 后验 $\int_{\theta} P(x|\theta) \cdot P(\theta) d\theta$
 看到 X 数据, 假设有隐变量 Z, 给定 P(Z), 求 P(Z|X)
 Inference: 精确推断 → 确定性近似 → VI
 近似推断 → 随机近似 → MCMC, MH, Gibbs

Bayes Inference; Bayes 决策: $x - N$ 个样本, \hat{x} 新的样本, 求 $P(\hat{x}|x)$
 $P(\hat{x}|x) = \int_{\theta} P(\hat{x}, \theta|x) d\theta = \int_{\theta} P(\hat{x}|\theta) \cdot P(\theta|x) d\theta = E_{\theta|x}[P(\hat{x}|\theta)]$

11.2. Formula Deduction

X : observed data
 Z : latent variable + parameter
 (X, Z) : complete data

$\log P(x) = \log P(x, z) - \log P(z|x)$
 $= \log \frac{P(x, z)}{q(z)} - \log \frac{P(z|x)}{q(z)}$

右边 = $\int_Z \log P(x) \cdot q(z) dz = [\log P(x)] = \text{constant}$

左边 = $\int_Z \log \frac{P(x, z)}{q(z)} dz = \underbrace{\int_Z \log \frac{P(x, z)}{q(z)} dz}_{\text{evidence lower bound ELBO}} - \underbrace{\int_Z \log \frac{P(z|x)}{q(z)} dz}_{KL(q||p)}$

$L(q) = \underbrace{\int_Z q(z) \log P(x, z) dz}_{①} - \underbrace{\int_Z q(z) \log q(z) dz}_{②}$ 这里的 Z 是大写, 代表所有
 隐变量

$① = \int_Z \prod_{i=1}^M q_i(z_i) \log P(x, z) dz_1 \dots dz_M$

$= \int_Z q_j(z_j) \underbrace{\left(\prod_{i \neq j}^M q_i(z_i) \log P(x, z) \cdot dz_1 \dots dz_{j-1} dz_{j+1} \dots d_M \right)}_{\log P(x, z) \cdot \prod_{i \neq j}^M q_i(z_i) dz_i} dz_j$

$\log P(x, z) \cdot \prod_{i \neq j}^M q_i(z_i) dz_i$ $\log P(x, z_j)$

$= \int_Z q_j(z_j) \cdot \underbrace{\left[\prod_{i \neq j}^M q_i(z_i) \log P(x, z) \right]}_{\text{关于 } z_j} dz_j$ (关于 z_j) $\int_Z q_j(z_j) dz_j$

$② = \int_Z q(z) \log q(z) dz = \int_Z \prod_{i=1}^M q_i(z_i) \cdot \log \prod_{i=1}^M q_i(z_i) dz$

$= \int_Z \prod_{i=1}^M q_i(z_i) \cdot \sum_{i=1}^M \log q_i(z_i) dz = \int_Z \prod_{i=1}^M q_i(z_i) \cdot [\log q_1(z_1) + \dots + \log q_M(z_M)] dz$

$\hat{q}(z) = \arg \max_{q(z)} L(q) \Rightarrow \hat{q}(z) \approx P(z|x)$ 其中 T: $\int_Z \prod_{i=1}^M q_i(z_i) \cdot \log q_i(z_i) dz = \int_Z q_1(z_1) \dots q_M(z_M) \cdot \log q_i(z_i) dz$

$q(z) = \prod_{i=1}^M q_i(z_i) \rightarrow \text{mean theory} \rightarrow q_j(z_j)$

P35 $= \int_Z q_1 \log q_1 dz_1 \cdot \int_Z q_2 dz_2 \dots \int_Z q_M dz_M = \int_Z q_i \log q_i dz_i$

$$\sum_{i=1}^M \int_{Z_i} q_i(z_i) \log q_i(z_i) dz_i = \boxed{\int_{Z_j} q_j(z_j) \log q_j(z_j) dz_j + C} \quad \text{其它项+常数项}$$

$$①-②: \int_{Z_j} q_j(z_j) \log \frac{\hat{P}(x, z_j)}{q_j(z_j)} dz_j + C = -\text{KL}(q_j || \hat{P}(x, z_j)) \leq 0, \text{ 当 } \frac{\hat{P}(x, z_j)}{q_j(z_j)} \text{ 为常数的话, 取等号}$$

$$\therefore q_j(z_j) = \hat{P}(x, z_j) \text{ 通过 argmax ELBO 作出, 它的核心就是 } P(x, z) \rightarrow \text{结合 P35. } \log q_j(z_j) = \log \hat{P}(x, z_j) = E_{\hat{P}} [\log P(x, z)]$$

11.3 Review (再回首)

$$x: \text{observed variable} \rightarrow X = \{x^{(i)}\}_{i=1}^N = \{x^{(i)}\}_{i=1}^N$$

$$z: \text{latent variable} \rightarrow Z = \{z^{(i)}\}_{i=1}^N = \{z^{(i)}\}_{i=1}^N$$

(X, z) : complete data

θ : model parameter

$$\text{ELBO} = E_{q(z)} [\log \frac{P(x, z | \theta)}{q(z)}] = E_{q(z)} [E[\log P_\theta(x | z)] + H[q(z)]] \text{ 故 VI. } \log P_\theta(x) = \underbrace{\text{ELBO}}_{f(\theta)} + \underbrace{\text{KL}(q || p)}_{\geq 0} \geq L(\theta)$$

$$\text{KL}(q || p) = \int q(z) \cdot \log \frac{q(z)}{P_\theta(z | x^{(i)})} dz$$

VI (mean field) \rightarrow classical VI

$$\text{Assumption: } q(z) = \prod_{i=1}^M q_i(z_i)$$

$$\log q_j(z_j) = E_{\prod_{i \neq j} q_i(z_i)} [\log P_\theta(x^{(i)}, z | \theta)] + \text{const}$$

$$= \int_{q_1} \int_{q_2} \dots \int_{q_{j-1}} \int_{q_{j+1}} \dots \int_{q_M} q_1 q_2 \dots q_{j-1} q_{j+1} \dots q_M [\log P_\theta(x^{(i)}, z)] \cdot dq_1 dq_2 \dots dq_{j-1} dq_{j+1} \dots dq_M$$

$$\text{则关于目标函数有: } \hat{q}_1(z_1) = \int_{q_2} \dots \int_{q_M} q_2 \dots q_M [\log P_\theta(x^{(i)}, z)] dq_2 \dots dq_M$$

$$\hat{q}_2(z_2) = \int_{q_1} \int_{q_3} \dots \int_{q_M} \hat{q}_1(q_1) \dots q_M [\log P_\theta(x^{(i)}, z)] dq_1 dq_3 \dots dq_M$$

$$\vdots$$

$$\hat{q}_M(z_M) = \int_{q_1} \int_{q_2} \dots \int_{q_{M-1}} \hat{q}_1(q_1) \dots \hat{q}_{M-1}(q_{M-1}) [\log P_\theta(x^{(i)}, z)] dq_1 dq_2 \dots dq_{M-1}$$

坐标上升法: Coordinate Ascend

classical VI:

问题: ①假设太强: (mean field) 每个变量相互独立; ② intractable

11.4 SGVI

$$\text{目标函数: } \hat{q} = \arg \min_q \text{KL}(q || p) = \arg \max_q L(q) \quad \log P_\theta(x^{(i)}) = \underbrace{\text{ELBO}}_{L(\phi)} + \underbrace{\text{KL}(q || p)}_{\geq 0} \geq L(\phi)$$

假设 $q(z) \rightarrow q_\phi(z)$, $q_\phi(z)$ 的参数为 ϕ . 是一个指数族分布

$$\text{ELBO} = E_{q_\phi(z)} [\log P_\theta(x^{(i)}, z) - \log q_\phi(z)] = L(\phi)$$

$$\text{则关于 } \hat{\phi} = \arg \max_\phi L(\phi)$$

$$\nabla_\phi L(\phi) = \nabla_\phi E_{q_\phi} [\log P_\theta(x^{(i)}, z) - \log q_\phi] \dots ①$$

$$= \nabla_\phi \int q_\phi [\log P_\theta(x^{(i)}, z) - \log q_\phi] dz$$

$$= \int \nabla_{\phi} q_{\phi} \cdot [\log P_{\theta}(x^{(i)}, z) - \log q_{\phi}] dz + \int q_{\phi} \cdot \nabla_{\phi} [\log P_{\theta}(x^{(i)}, z) - \log q_{\phi}] dz$$

= ① + ②

$$\textcircled{2} = - \int q_{\phi} \cdot \nabla_{\phi} \log q_{\phi} dz$$

$$= - \int q_{\phi} \cdot \frac{1}{q_{\phi}} \nabla_{\phi} \cdot q_{\phi} dz$$

$$= - \int \nabla_{\phi} q_{\phi} dz$$

$$= - \nabla_{\phi} \underbrace{\int q_{\phi} dz}_{= 0}$$

variance Reduction

需要非常多样本

high variance

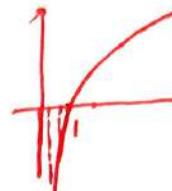
+ 分不稳定，并且在采样过程中
得到的 $\nabla_{\phi} L(\phi)$ 是近似的，之后
使用 SG 仍是近似，二者误差增加

$$\text{由: } \nabla_{\phi} q_{\phi} = \nabla_{\phi} \log q_{\phi} \cdot q_{\phi}$$

$$\text{得: } \int q_{\phi} \cdot \nabla_{\phi} \log q_{\phi} [\log P_{\theta}(x^{(i)}, z) - \log q_{\phi}] dz$$

$$= \mathbb{E}_{q_{\phi}} [\nabla_{\phi} \log q_{\phi} [\log P_{\theta}(x^{(i)}, z) - \log q_{\phi}]]$$

$$= 0$$



故: $\nabla_{\phi} L(\phi) = \mathbb{E}_{q_{\phi}} [\nabla_{\phi} \log q_{\phi} [\log P_{\theta}(x^{(i)}, z) - \log q_{\phi}]]$, $q_{\phi} = q_{\phi}(z)$ 可用 MC 采样.

$$z^{(i)} \sim q_{\phi}(z), i=1, 2, \dots, L$$

$$\approx \frac{1}{L} \sum_{l=1}^L \nabla_{\phi} \log q_{\phi}(z^{(i)}) [\log P_{\theta}(x^{(i)}, z^{(i)}) - \log q_{\phi}(z^{(i)})] \rightarrow$$

缺点:
重参数化技巧

降低 Variance: Reparameterization Trick (目的: 让 $q_{\phi} \rightarrow p_{\epsilon}$ 与 $\nabla_{\phi} L(\phi)$ 易求, 改进)

$$\nabla_{\phi} L(\phi) = \nabla_{\phi} \mathbb{E}_{q_{\phi}} [\log P_{\theta}(x^{(i)}, z) - \log q_{\phi}], p36 \text{ ①式}$$

假设 $z = g_{\phi}(\epsilon, x^{(i)})$, $\epsilon \sim p(\epsilon)$ 是一个确定的分布

$$z \sim q_{\phi}(z | x^{(i)}) \quad \text{可认为 } q_{\phi} \text{ 是一个变换}$$

$$\downarrow \quad \downarrow \quad \downarrow \quad \text{自己给定的一分布 } z \text{ 的随机性转换给 } \epsilon$$

$$\epsilon \sim p(\epsilon) \quad \int q_{\phi}(z | x^{(i)}) dz = 1$$

$$\text{有 } |q_{\phi}(z | x^{(i)})| dz = |p(\epsilon)| d\epsilon$$

$$\text{由 } \nabla_{\phi} L(\phi) = \nabla_{\phi} \mathbb{E}_{q_{\phi}} [\log P_{\theta}(x^{(i)}, z) - \log q_{\phi}]$$

$$= \nabla_{\phi} \int [\log P_{\theta}(x^{(i)}, z) - \log q_{\phi}] \cdot q_{\phi} dz$$

$$= \nabla_{\phi} \int [\log P_{\theta}(x^{(i)}, z) - \log q_{\phi}] \cdot p(\epsilon) d\epsilon$$

$$= \nabla_{\phi} \mathbb{E}_{p(\epsilon)} [\log P_{\theta}(x^{(i)}, z) - \log q_{\phi}] \quad \text{分布 } p(\epsilon) \text{ 与 } \phi \text{ 无关}$$

$$= \mathbb{E}_{p(\epsilon)} \left\{ \nabla_{\phi} [\log P_{\theta}(x^{(i)}, z) - \log q_{\phi}] \right\}$$

$$= \mathbb{E}_{p(\epsilon)} \left\{ \nabla_{\phi} [\log P_{\theta}(x^{(i)}, z) - \log q_{\phi}(z | x^{(i)})] \cdot \nabla_{\phi} g_{\phi}(\epsilon, x^{(i)}) \right\} \quad \text{链式法则}$$

$$= \mathbb{E}_{p(\epsilon)} \left\{ \nabla_{\phi} [\log P_{\theta}(x^{(i)}, z) - \log q_{\phi}(z | x^{(i)})] \cdot \nabla_{\phi} \cdot g_{\phi}(\epsilon, x^{(i)}) \right\}$$

$$\frac{\partial f}{\partial \phi} = \frac{\partial f}{\partial z} \cdot \frac{\partial z}{\partial \phi}, z = g(\phi)$$

此时采样:

$$\epsilon^{(i)} \sim p(\epsilon), i=1, 2, \dots, L$$

$$\nabla_{\phi} L(\phi) \approx \frac{1}{L} \sum_{l=1}^L \left\{ \nabla_{\phi} [\log P_{\theta}(x^{(i)}, z) - \log q_{\phi}(z | x^{(i)})] \cdot \nabla_{\phi} g_{\phi}(\epsilon^{(i)}, x^{(i)}) \right\}$$

SGVI:

$$\phi^{(t+1)} \leftarrow \phi^{(t)} + \lambda^{(t)} \cdot \nabla_{\phi} L(\phi)$$

12. MCMC (Markov chain & Monte Carlo)

12.1 Sampling Method Introduction

Inference $\begin{cases} \text{精确推断} \\ \text{近似推断} \end{cases} \rightarrow \begin{cases} \text{确定性} \\ \text{随机} \end{cases} \rightarrow \text{MCMC}$

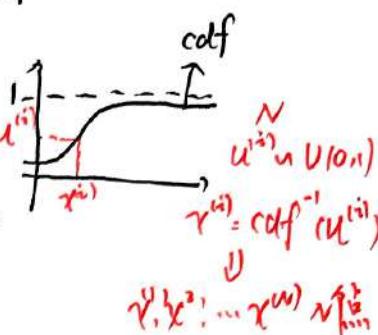
Monte Carlo Method: 基于采样的随机近似方法

$$P(z|x) \xrightarrow{\substack{\text{observed data} \\ \text{lateral variable}}} E_{z|x}[f(z)] = \int P(z|x) \cdot f(z) dz \approx \frac{1}{N} \sum_{i=1}^N f(z_i), \quad z^{(1)}, z^{(2)}, \dots, z^{(N)} \text{ 来自 } (n) P(z|x)$$

由 PDF 和 cdf 很困难

① 概率分布采样:

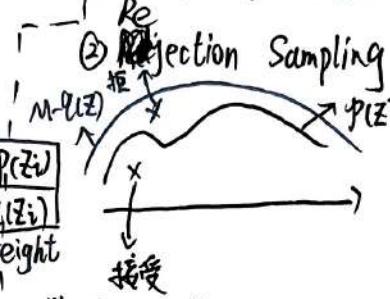
$P(z)$: PDF \rightarrow 求出 CDF



② Importance Sampling (期望采样)

$$\begin{aligned} E_{p(z)}[f(z)] &= \int p(z) f(z) dz = \int \frac{p(z)}{q(z)} \cdot q(z) f(z) dz \\ &= \int f(z) \frac{p(z)}{q(z)} \cdot q(z) dz \approx \frac{1}{N} \sum_{i=1}^N f(z_i) \frac{p(z_i)}{q(z_i)} \end{aligned}$$

$z_i \sim q(z), i=1, \dots, N$ 样本从 $q(z)$ 中选 weight



$$q(z) = \text{proposal distribution}$$

$$\forall z_i, \frac{p(z_i)}{q(z_i)} \geq 1$$

α : 接收率, $\alpha = \frac{p(z^{(i)})}{M q(z^{(i)})}, 0 \leq \alpha \leq 1$

取决于 $p(z)$ 与 $q(z)$ 相似程度

③ 的变形: Sampling - Importance - Resampling (重要思想)

Importance sampling 的错误: 乘上 (weight)

① 采样

② 采样, 从 $z^{(1)}, z^{(2)}, \dots, z^{(N)}$ 等样
以①中 weight 为概率直接去采样

借助 $U(0,1)$

i: $z^{(i)} \sim q(z)$

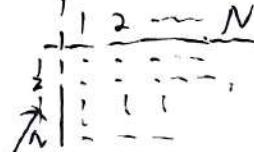
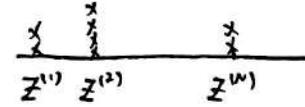
ii: $u \sim U(0,1)$

if $u \leq \alpha$, 接受 $z^{(i)}$

else 非拒

$p(z)$ 与 $q(z)$ 相似程度

采样效率不均, 很能有很高很低的



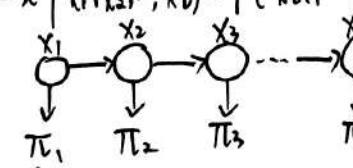
第N状态到第N状态的概率
每行之和为1, 即 $\sum p_{ij} = 1$

12.2 Markov Chain

Markov Chain: 时间和状态都是离散的 $\{X_t\}, P \rightarrow \text{转移矩阵 } P_{ij}$
未来只依赖于当前, 与过去无关

齐次 (-H) Markov Chain: $P(X_{t+1}=x|X_1, X_2, \dots, X_t) = P(X_{t+1}=x|X_t)$ 只依赖于前一个状态

$$P_{ij} = P(X_{t+1}=j|X_t=i)$$



每一个随机变量都有属于它的概率分布 π_i

平稳分布: $\{\pi_i\}$

$$\sum_{i=1}^{\infty} \pi_i = 1$$

若 $\{\pi_i\} = \{\pi^{(1)}, \pi^{(2)}, \dots, \pi^{(N)}, \dots\}$ 使得 $\pi(x^*) = \int \pi(x) \cdot P(x \rightarrow x^*) dx$ 成立

$$\begin{array}{c} P(x \rightarrow x^*) \\ t: \pi_t(x) \quad t+1: \pi_{t+1}(x^*) \end{array}$$

则 $\{\pi_i\}$ 是 $\{\pi_i\}$ 的平稳分布, 即每一个时刻 π_i 的分布是相同的

Detailed Balance: $\pi(x) \cdot P(x \rightarrow x^*) = \pi(x^*) \cdot P(x^* \rightarrow x)$ 这个结构的 Markov Chain

Detailed Balance \Leftrightarrow 平稳分布 (没有基于 $0 \rightarrow \dots \rightarrow 0 \rightarrow \dots$ 这个结构的 Markov Chain
讲法, 而是更 general 的)

证明: $\int \pi(x) \cdot P(x \rightarrow x^*) dx = \int \pi(x^*) \cdot P(x^* \rightarrow x) dx$

$$\begin{aligned} &= \pi(x^*) \int P(x^* \rightarrow x) dx \xrightarrow{\text{若为离散}} \sum_{j=1}^{\infty} P_{ij} = 1 \quad \text{转移矩阵的一行} \\ &= \pi(x^*) \end{aligned}$$

12.3 MH Algorithm 接受率

$$P(z) \cdot Q(z \rightarrow z^*) \cdot d(z, z^*) = P(z^*) \cdot Q(z^* \rightarrow z) \cdot d(z^*, z) \quad (1), \quad Q = [Q_{ij}] \text{ proposal matrix}$$

$$P(z) \cdot P(z \rightarrow z^*) = P(z^*) \cdot P(z^* \rightarrow z)$$

$$d(z, z^*) = \min(1, \frac{P(z^*) \cdot Q(z^* \rightarrow z)}{P(z) \cdot Q(z \rightarrow z^*)}) \cdot Q(z^* | z)$$

从 $P = [P_{ij}]$ 中随机选取的

①式中，若不加入 $d(z, z^*)$ 时不满足 Detailed Balance
因为 Q 是随机选取的，此时 $P(z)$ 并不满足 DB.

$$P(z) \cdot Q(z \rightarrow z^*) \neq P(z^*) \cdot Q(z^* \rightarrow z)$$

Q_{ij} 是任取的

从 Q 不可推后，加上 $d(z, z^*)$ 后可推后

$$\frac{P(z) \cdot Q(z \rightarrow z^*)}{P(z^*) \cdot P(z^* \rightarrow z)}$$

$d(z^*, z)$

$$\text{证明加入 } d(z, z^*) \text{ 后满足 DB: } P(z) \cdot P(z \rightarrow z^*)$$

$$\boxed{P(z) \cdot Q(z \rightarrow z^*) \cdot d(z, z^*)} = P(z) \cdot Q(z \rightarrow z^*) \cdot \min(1, \frac{P(z^*) \cdot Q(z^* \rightarrow z)}{P(z) \cdot Q(z \rightarrow z^*)})$$

$$= \min(P(z) \cdot Q(z \rightarrow z^*), P(z^*) \cdot Q(z^* \rightarrow z)) = \min(P(z^*) \cdot Q(z^* \rightarrow z) \cdot \min(1, \frac{P(z) \cdot Q(z \rightarrow z^*)}{P(z^*) \cdot Q(z^* \rightarrow z)}))$$

$$= P(z^*) \cdot Q(z^* \rightarrow z) \cdot d(z^*, z)$$

归一化因子，常数化因子
 $\frac{P(z)}{P(z^*)}$ 一般不可求

Metropolis-Hastings:

$$\begin{aligned} & \text{从 } Q \text{ 中采集} \\ & z^{(i)} \sim Q(z | z^{(i-1)}) \quad \uparrow \\ & d = \min(1, \frac{P(z^*) \cdot Q(z^* \rightarrow z)}{P(z) \cdot Q(z \rightarrow z^*)}) \quad \uparrow \\ & \text{if } U \leq d, z^{(i)} = z^* \\ & \text{else } z^{(i)} = z^{(i-1)} \end{aligned}$$

与 Rejection Sampling 区别: $U > d$ 时不拒绝，而是重复上一个样本，这样不会造成样本量减少

$$P(z) \rightarrow E_{P(z)}[f(z)] = \int_z P(z) f(z) dz \approx \frac{1}{N} \sum_{i=1}^N f(z^{(i)})$$

d : 接收率的作用，把 $Q \rightarrow P$
①从 Q 中采集；②和 d 判断是否接收

执行 N 次，得到 $z^{(1)}, z^{(2)}, \dots, z^{(N)}$ 个样本点

一维一维的采

12.4 Gibbs Algorithm (special MH Sampling)

$p(z) = p(z_1, \dots, z_n)$ 高维采样的数据是高维

$z_i \sim P(z_i | z_{-i})$, 固定其它样本

$p(z) = p(z_1, z_2, z_3)$

$$i: z_1^{(t+1)}, z_2^{(t+1)}, z_3^{(t+1)} \quad z_{-i} = z_{-i}^* \Rightarrow z_{-i} = z_{-i}^*$$

$i: t+1$:

$$z_1^{(t+1)} \sim P(z_1 | z_2^{(t+1)}, z_3^{(t+1)})$$

$$d = \min(1, \frac{P(z^*) \cdot Q(z^* \rightarrow z)}{P(z) \cdot Q(z \rightarrow z^*)})$$

$$z_2^{(t+1)} \sim P(z_2 | z_1^{(t+1)}, z_3^{(t+1)})$$

$$\text{设: } z_3^{(t+1)} \sim P(z_3 | z_1^{(t+1)}, z_2^{(t+1)}) = P(z^*) = P(z_i^* | z_{-i}^*)$$

$$\frac{P(z^*) \cdot Q(z^* \rightarrow z)}{P(z) \cdot Q(z \rightarrow z^*)} = \frac{P(z_i^* | z_{-i}^*) \cdot P(z_{-i}^*) \cdot P(z | z_{-i}^*)}{P(z_i | z_{-i}) \cdot P(z_{-i}) \cdot P(z^* | z)}$$

$$= \frac{P(z_i^* | z_{-i}^*) \cdot P(z_{-i}^*) \cdot P(z_i | z_{-i}^*)}{P(z_i | z_{-i}) \cdot P(z_{-i}) \cdot P(z_{-i}^* | z_{-i})}$$

$$\text{由 } z_{-i} = z_{-i}^*, \text{ 有: } \frac{P(z_i^* | z_{-i}^*) \cdot P(z_{-i}^*) \cdot P(z_i | z_{-i}^*)}{P(z_i | z_{-i}) \cdot P(z_{-i}) \cdot P(z_{-i}^* | z_{-i})} = 1$$

代入 MH 中的 $d = \min(1, \frac{P(z^*) \cdot Q(z^* \rightarrow z)}{P(z) \cdot Q(z \rightarrow z^*)})$ 中后， $d = 1$ ，效率更高

12.5 Review

采样的动机：①采样本身就是常见的任务

$$\text{②求和求积分, e.g. } \mathbb{E}_{p(x)}[f(x)] = \int p(x) f(x) dx \approx \frac{1}{N} \sum_{i=1}^N f(x^{(i)})$$

什么是好的样本:

①样本趋于高概率区域

②样本之间相互独立

采样是困难的,

(the curse of high dimensionality)

① partition function is intractable

② high dimension, e.g. $x \in \{1, 2, \dots, k\}^P$

状态空间 \mathbb{K}^P

$\frac{x_1 | 1 2 \dots k^P}{p | p_1 p_2 \dots p_{k^P}}$ * 考虑样本趋于高概率, 把每一个 $p_i \rightarrow p_{k^P}$ 比较不现实

因此有: Rejection Sampling, Importance Sampling, $P(x)$ 采样有困难, 通过 $q(x)$ 经过 $P(x)$

MCMC: $\begin{cases} MH \\ Gibbs \text{ MH 特例} \end{cases}$

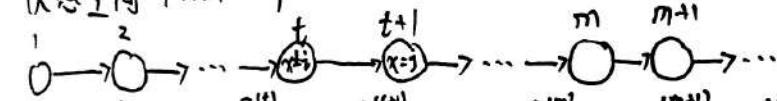
12.6 Stationary Distribution (平稳分布)

Rejection Sampling } ① $q(x)$ 和 $P(x)$ 接近

Importance Sampling } ③ $q(x)$ 简单

进入平稳分布

状态空间: $\{1, 2, \dots, k\}$



$q^{(t)}(x) = \frac{x_1 | 1 2 \dots k}{q_1^{(t)} q_2^{(t)} \dots q_k^{(t)}}$ -> 有 k 种状态

状态转移矩阵 Q : $Q = \begin{pmatrix} Q_{11} & Q_{12} & \dots & Q_{1k} \\ Q_{21} & Q_{22} & \dots & Q_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ Q_{k1} & Q_{k2} & \dots & Q_{kk} \end{pmatrix}$

随机矩阵

令 $q^{(t+1)} = (q^{(t+1)}_{x=1}, q^{(t+1)}_{x=2}, \dots, q^{(t+1)}_{x=k})_{1 \times k}$

而 $q^{(t+1)}_{x=j} = \sum_{i=1}^k q^{(t)}_{x=i} \cdot Q_{ij}$

$\therefore q^{(t+1)} = \left(\sum_{i=1}^k q^{(t)}_{x=i} \cdot Q_{ik} \right)_{1 \times k}$

$\sum_{i=1}^k q^{(t)}_{x=i} \cdot Q_{ik} = Q_{ik} \cdot (q^{(t)}_{x=1} + q^{(t)}_{x=2} + \dots + q^{(t)}_{x=k}) \cdot Q_{kk}$

$q^{(t)} = (q^{(t)}_{x=1}, q^{(t)}_{x=2}, \dots, q^{(t)}_{x=k})_{1 \times k}$

$q^{(t+1)} = q^{(t)} \cdot Q = q^{(t)} \cdot Q^2 = \dots = q^{(t)} \cdot Q^t$

R.H.S: $q^{(t+1)} = q^{(t)} \cdot (A \cdot \Delta A^{-1})^t = q^{(t)} \cdot A \cdot \Delta^t A^{-1}$

当 $t \rightarrow \infty$ 时, $\Delta^t \rightarrow \begin{pmatrix} 0 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{pmatrix} = A$

核心, 构造一条不可失链, 使其收敛的 $q(x)$ 与我们的 $P(x)$ 相近, 通过采样的方式, 在 m 之后的样本取出一些样本 ($1 \downarrow$), 则就实现了通过 $q^{(m)}$ 直接从 $P(x)$ 中采样

重要的点是转移矩阵 Q 的构造

$$q^{(t+1)}_{x=j} = \sum_{i=1}^k q^{(t)}_{x=i} \cdot Q_{ij}$$

$$\begin{aligned} \sum_{j=1}^k Q_{ij} &= 1 \quad (i=1, \dots, k) \\ \sum_{j=1}^k Q_{ji} &= 1 \quad (i=1, \dots, k) \end{aligned} \quad \text{每行每列根频率都为 1}$$

随机矩阵特征值的绝对值 ≤ 1 不妨设 $\lambda_i \neq 1$

$Q = A \cdot \Delta A^{-1}$

$$q^{(t+1)} = q^{(t)} \cdot A \cdot \Delta A^{-1}$$

$$q^{(m+1)} = q^{(m+1)} \cdot A \cdot \Delta A^{-1} = q^{(m+1)} \cdot A \cdot \Delta^{m+1} \cdot A^{-1} = q^{(m+1)}$$

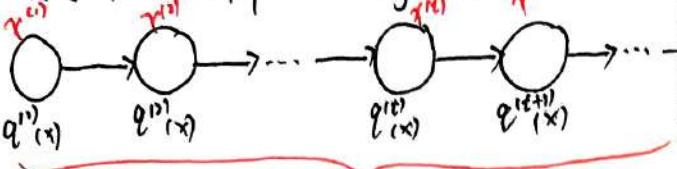
$$\Delta^{m+1} = \Delta^m \cdot \Delta = \begin{pmatrix} 0 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{pmatrix} = \begin{pmatrix} 0 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{pmatrix}$$

P40 $\therefore q^{(m+1)} = q^{(m+1)}$, 当 $t > m$ 时, $q^{(m+1)} = q^{(m+2)} = \dots = q^{(\infty)}$, 前提是同一个分布, 平稳

12.7 Problem & Thinking

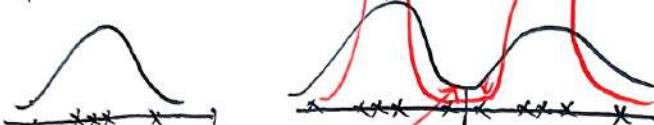
状态空间: $\{1, 2, \dots, k\}$

状态转移矩阵: $Q = [Q_{ij}]_{K \times K}$



burn-in
mixing time

单峰 vs 多峰 (1维)



$P_N = \frac{1}{Z} \exp(-E(\eta))$
partition function
能量函数
(2维)



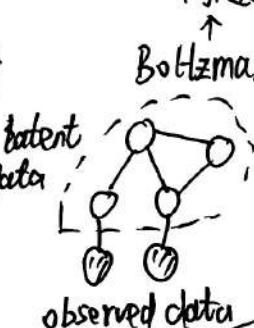
即有可能无法到达另一个
峰点, 无法超过这道低
沟, 因为算法都是趋
向于高概率区域的

MCMC: 利用Markov chain 收敛于平稳
分布设计 Q , 使得 平稳分布 目标分布
 $q(x)$ $P(x)$

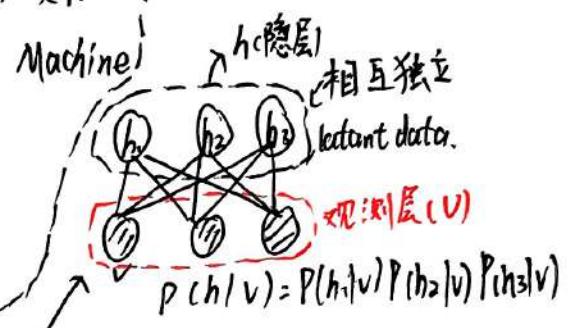
MCMC 问题:

- ① 理论上保证收敛性, 但无法知道何时收敛
- ② mixing time 过长 $P(x)$ 太复杂, 高维及相关性
- ③ 样本之间有一定的相关性 (每隔 100, 1000
样本取一个)

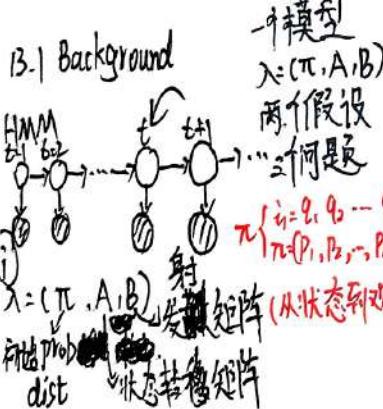
问题: 太复杂



RBM \rightarrow Restricted Boltzmann machine



13.1 Background



射线矩阵 (从状态到观测量)
状态转移矩阵

观测量
 $O, O_1, O_2, \dots, O_t, \dots, O_N$

状态
 $i, i_1, i_2, \dots, i_t, \dots, i_N$

观测量
 V, V_1, V_2, \dots, V_M

状态
 Q, Q_1, Q_2, \dots, Q_N

观测量
 b, b_1, b_2, \dots, b_K

两个假设:
①齐次Markov假设
②观测量独立假设

13. Hidden Markov Model

模型 $\lambda = (\pi, A, B)$
两个假设

概率图 $\left\{ \begin{array}{l} \text{有向 - Bayesian Network} \\ \text{无向 - Markov Random Field} \\ \text{(Markov Network)} \end{array} \right.$

Dynamic Model $\left\{ \begin{array}{l} \text{HMM} \\ \text{kalman Filter} \\ \text{Particle Filter} \end{array} \right.$

三个问题
time
观测量
动态
x_i 之间是 iid

Evaluation $P(o|\lambda) \rightarrow$ 前向后向
Learning \rightarrow 变向求
 $\lambda = \text{argmax } P(o|\lambda)$

Decoding: 找到一个状态序列
 $I = \text{argmax } P(I|o)$

Prediction $\rightarrow P(i_{t+1}|O, O_1, \dots, O_t)$

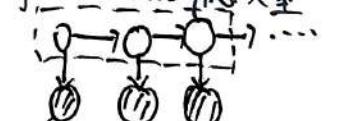
Filtration $\rightarrow P(i_t|O_1, O_2, \dots, O_b)$

GMM N, π_1, \dots, π_N

$\pi_i \sim \text{dir}(P_N|\theta)$



$P(y|z) \sim N(\mu, \Sigma)$
system state 隐变量



观测量

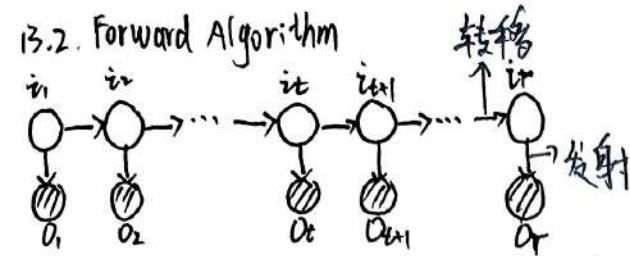
state \rightarrow 离散 HMM

$\lambda = \text{argmax } P(o|\lambda)$ \rightarrow EM
Bartel Welch \rightarrow 连续
Kalman Filter

$I = \text{argmax } P(I|o)$ \rightarrow Particle Filter

$P_{t+1} = P(i_{t+1}|O, O_1, \dots, O_t)$
滤波 $\rightarrow P(i_t|O_1, O_2, \dots, O_b)$

13.2. Forward Algorithm



$I = i_1, i_2, \dots, i_T \rightarrow$ 状态序列, $\mathcal{Q} = \{q_1, q_2, \dots, q_N\} \rightarrow$ 状态值集合

$O = o_1, o_2, \dots, o_T \rightarrow$ 观测序列, $\mathcal{V} = \{v_1, v_2, \dots, v_M\} \rightarrow$ 观测值集合

(i) $\lambda = (\pi, A, B)$

π : 初始概率分布, $\pi = (\pi_1, \dots, \pi_N)$, $\sum_i \pi_i = 1$, $\pi_i = P(i_1 = q_i)$
 $A = [A_{ij}]_{N \times N}$ 转移矩阵, $A_{ij} = P(i_{t+1} = q_j | i_t = q_i)$

$B = [b_j(k)]_{M \times N}$ 发射矩阵, $b_j(k) = P(O_t = v_k | i_t = q_j)$

(ii) 两个假设: ① 齐次 Markov ② 观测独立 $i \rightarrow O$

$$1) P(i_{t+1} | i_1, \dots, i_t, o_1, \dots, o_t) = P(i_{t+1} | i_t)$$

$$2) P(O_t | i_1, \dots, i_t, o_1, \dots, o_t) = P(O_t | i_t)$$

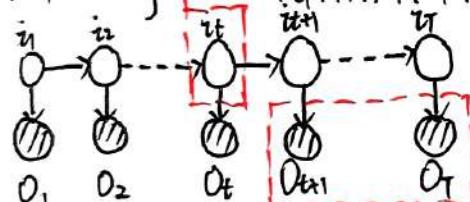
(iii) 三个问题

① Evaluation: Given λ , 求 $P(O|\lambda)$ Forward-Backward

② Learning: $\lambda_{MLE} = \arg \max P(O|\lambda)$ Baum-Welch

③ Decoding: $i = \arg \max P(i|O, \lambda)$ Viterbi

13.3 Backward Algorithm: 其实从后往前算一个状态和



记 $\beta_t(i) = P(O_{t+1}, \dots, O_T | i_t = q_i, \lambda)$

$$\beta_t(i) = P(O_2, \dots, O_T | i_t = q_i, \lambda)$$

$$P(O|\lambda) = P(O_1, \dots, O_T | \lambda) = \sum_{i=1}^N P(O_1, \dots, O_T, i_t = q_i | \lambda)$$

$$= \sum_{i=1}^N P(O_1, \dots, O_T | i_t = q_i, \lambda) \cdot P(i_t = q_i | \lambda)$$

$$= \sum_{i=1}^N P(O_1 | O_2, \dots, O_T, i_t = q_i, \lambda) \cdot P(O_2, \dots, O_T | i_t = q_i, \lambda) \cdot \pi_i$$

$$= \sum_{i=1}^N P(O_1 | i_t = q_i) \cdot \beta_t(i) \cdot \pi_i$$

$$= \sum_{i=1}^N b_i(O_1) \beta_t(i) \cdot \pi_i = P(O|\lambda)$$

$$\beta_t(i) = P(O_{t+1}, \dots, O_T | i_t = q_i, \lambda)$$

$$= \sum_{j=1}^N P(O_{t+1}, \dots, O_T, i_{t+1} = q_j | i_t = q_i, \lambda) \quad P(i_{t+1} = q_j, \dots, i_T = q_i, \lambda) | i_t = q_i, \lambda$$

$$= \sum_{j=1}^N P(O_{t+1}, \dots, O_T | i_{t+1} = q_j, i_t = q_i, \lambda) \cdot P(O_{t+1}, \dots, O_T, i_{t+1} = q_j | i_t = q_i, \lambda)$$

$$= \sum_{j=1}^N P(O_{t+1}, \dots, O_T | i_{t+1} = q_j, \lambda) A_{ij} = \sum_{j=1}^N P(O_{t+1}, \dots, O_T | i_{t+1} = q_j, \lambda) A_{ij}$$

head to tail, O_{t+1} 为观序列, $P_{it} \perp P_{it+1}$

Evaluation: Given λ , 求 $P(O|\lambda)$

$$P(O|\lambda) = \sum_I P(I, O|\lambda) = \sum_I P(O|I, \lambda) \cdot P(I|\lambda) = \sum_I P(i_1, i_2, \dots, i_T | \lambda)$$

$$\therefore P(I|\lambda) = P(i_1, i_2, \dots, i_T | \lambda) = P(i_T | i_1, i_2, \dots, i_{T-1}, \lambda) \cdot P(i_1, i_2, \dots, i_{T-1}, \lambda)$$

$$\text{由于太复杂, 放有 Forward, Backward} = \alpha_{i_{T-1}, i_T} \cdot \alpha_{i_{T-2}, i_{T-1}} \dots \alpha_{i_1, i_2} \cdot \pi(i_1)$$

$$= \pi(i_1) \cdot \prod_{t=2}^{T-1} \alpha_{i_{t-1}, i_t}$$

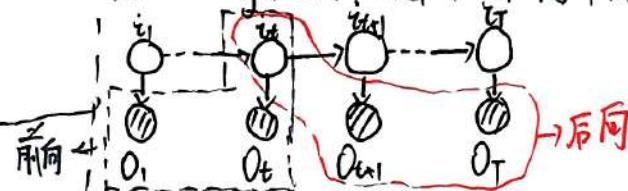
$$P(O|\lambda) = \prod_{t=1}^T b_{it}(O_t)$$

$$P(O|\lambda) = \sum_I \pi(i_1) \cdot \prod_{t=2}^T \alpha_{i_{t-1}, i_t} \prod_{t=1}^T b_{it}(O_t)$$

$$= \sum_{i_1} \sum_{i_2} \dots \sum_{i_T} \pi(i_1) \cdot \prod_{t=2}^T \alpha_{i_{t-1}, i_t} \cdot \prod_{t=1}^T b_{it}(O_t)$$

$N^T \Rightarrow O(N^T) \rightarrow$ 复杂度太高, 又难以计算

Forward Algorithm: 算出 d_t 再对每一个状态求和



$$d_t(i) = P(O_1, \dots, O_t, i_t = q_i | \lambda)$$

$$d_T(i) = P(O, i_T = q_i | \lambda)$$

$$P(O|\lambda) = \sum_{i=1}^N P(O, i_t = q_i | \lambda) = \sum_{i=1}^N d_T(i)$$

有 N 个状态值, 把它积掉

$$d_{t+1}(j) = P(O_1, \dots, O_t, O_{t+1}, i_{t+1} = q_j | \lambda) = \sum_{i=1}^N P(O_1, \dots, O_t, O_{t+1}, i_{t+1} = q_j, i_t = q_i | \lambda)$$

$$= \sum_{i=1}^N P(O_{t+1} | O_1, \dots, O_t, i_t = q_i, i_{t+1} = q_j, \lambda) \cdot P(O_1, \dots, O_t, i_t = q_i, i_{t+1} = q_j | \lambda)$$

$$= \sum_{i=1}^N P(O_{t+1} | i_{t+1} = q_j) \cdot P(i_{t+1} = q_j | O_1, \dots, O_t, i_t = q_i, \lambda)$$

$$= \sum_{i=1}^N P(O_{t+1} | i_{t+1} = q_j) \cdot P(i_{t+1} = q_j | i_t = q_i, \lambda) \cdot d_t(i)$$

$$= \sum_{i=1}^N b_j(O_{t+1}) \cdot a_{ij} \cdot d_t(i)$$

$$\therefore d_{t+1}(j) = \sum_{i=1}^N b_j(O_{t+1}) \cdot a_{ij} \cdot d_t(i) \text{ Forward}$$

$$\sum_{j=1}^N b_j(O_{t+1}) \cdot a_{ij} \cdot d_{t+1}(j) = \beta_t(i) \text{ Backward}$$

$$= \sum_{j=1}^N P(O_{t+1} | O_{t+2}, \dots, O_T, i_{t+1} = q_j, \lambda) \cdot P(O_{t+2}, \dots, O_T | i_{t+1} = q_j, \lambda)$$

$$= \sum_{j=1}^N P(O_{t+1} | i_{t+1} = q_j) \cdot P(i_{t+1} = q_j | i_t = q_i, \lambda)$$

$$= \sum_{j=1}^N b_j(O_{t+1}) \cdot a_{ij} \cdot \beta_{t+1}(j)$$

13.4 Baum Welch Algorithm - EM

$$\theta^{(t+1)} = \underset{\theta}{\operatorname{argmax}} \int_z \log P(x, z | \theta) \cdot P(z | x, \theta^{(t)}) dz$$

x : 观测
 \downarrow
 O

z : 隐变量
 \downarrow
 $I \rightarrow \text{离散}$

θ : 参数
 λ

$$P(O|\lambda) = \sum_I P(O, I | \lambda) = \sum_i \sum_{i_1}^I \pi(i_1) \prod_{t=2}^T a_{i_{t-1}, i_t} \prod_{t=1}^T b_{i_t}(O_t)$$

$$\pi_i = P(i_1 = q_i)$$

$$\lambda^{(t+1)} = \underset{\lambda}{\operatorname{argmax}} \sum_I \log P(O, I | \lambda) \cdot \underbrace{P(I | O, \lambda^{(t)})}_{\frac{P(I, O | \lambda^{(t)})}{P(O | \lambda^{(t)})}}$$

而 $\lambda^{(t)}$ 是上次迭代产生的一个常数，且 O 与 λ 无关，留下分子的 $\lambda^{(t)}$ 求解。

$$= \underset{\lambda}{\operatorname{argmax}} \sum_I \log P(O, I | \lambda) \cdot P(O, I | \lambda^{(t)}) \quad \lambda^{(t)} = (\pi^{(t)}, A^{(t)}, B^{(t)})$$

$$\begin{aligned} Q(\lambda, \lambda^{(t)}) &= \sum_I \log P(O, I | \lambda) \cdot P(O, I | \lambda^{(t)}) \\ &= \sum_I \log \left[\sum_{i_1} \sum_{i_2} \cdots \sum_{i_T} \pi(i_1) \cdot \prod_{t=2}^T a_{i_{t-1}, i_t} \cdot \prod_{t=1}^T b_{i_t}(O_t) \right] \cdot P(O, I | \lambda^{(t)}) \\ &= \sum_I \left\{ \log \pi(i_1) + \sum_{t=2}^T \log a_{i_{t-1}, i_t} + \sum_{t=1}^T \log b_{i_t}(O_t) \right\} \cdot P(O, I | \lambda^{(t)}) \end{aligned}$$

计算 $\pi^{(t+1)}$

$$\pi^{(t+1)} = \underset{\pi}{\operatorname{argmax}} Q(\lambda, \lambda^{(t)})$$

$$= \underset{\pi}{\operatorname{argmax}} \sum_I [\log \pi(i_1) \cdot P(O, I | \lambda^{(t)})]$$

$$= \underset{\pi}{\operatorname{argmax}} \sum_{i_1} \sum_{i_2} \cdots \sum_{i_T} [\log \pi(i_1) \cdot P(O, i_1, i_2, \dots, i_T | \lambda^{(t)})]$$

$$= \underset{\pi}{\operatorname{argmax}} \sum_{i_1} [\log \pi(i_1) \cdot P(O, i_1 | \lambda^{(t)})]$$

$$= \underset{\pi}{\operatorname{argmax}} \sum_{i_1} [\log \pi(i_1) \cdot P(O, i_1 = q_i | \lambda^{(t)})]$$

$$\text{s.t. } \sum_{i_1} \pi(i_1) = 1$$

$$L(\pi, \gamma) = \sum_{i_1}^N [\log \pi(i_1) \cdot P(O, i_1 = q_i | \lambda^{(t)})] + \gamma \left(\sum_{i_1}^N \pi(i_1) - 1 \right)$$

$$\frac{\partial L}{\partial \pi(i_1)} = \frac{1}{\pi(i_1)} P(O, i_1 = q_i | \lambda^{(t)}) + \gamma \triangleq 0$$

$$P(O, i_1 = q_i | \lambda^{(t)}) + \gamma \triangleq 0 \Rightarrow \pi(i_1) = \frac{1}{\gamma} P(O, i_1 = q_i | \lambda^{(t)})$$

$$\sum_{i_1}^N [P(O, i_1 = q_i | \lambda^{(t)})] + \sum_{i_1}^N \pi(i_1) = 0$$

$$P(O | \lambda^{(t)}) + \gamma = 0$$

$$\gamma = -P(O | \lambda^{(t)})$$

$$\pi_i^{(t+1)} = \frac{1}{P(O | \lambda^{(t)})} P(O, i_1 = q_i | \lambda^{(t)})$$

$$\pi^{(t+1)} = (\pi_1^{(t+1)}, \pi_2^{(t+1)}, \dots, \pi_N^{(t+1)})$$

$$\begin{aligned} d_t(i) \cdot \beta_t(i) &= P(i_t = q_i, O | \lambda) \\ \gamma_t(i) &= P(i_t = q_i | O, \lambda) = \frac{P(i_t = q_i, O | \lambda)}{P(O | \lambda)} \\ &= \frac{d_t(i) \cdot \beta_t(i)}{\sum_j d_t(j) \beta_t(j)} = \frac{\sum_i d_t(i) \beta_t(i)}{\sum_i \sum_j d_t(i) \beta_t(j)} \end{aligned}$$

$$\begin{aligned} \delta_t(i, j) &= P(i_t = q_i, i_{t+1} = q_j | O, \lambda) \\ \delta_t(i, j) &= \frac{P(i_t = q_i, i_{t+1} = q_j, O | \lambda)}{P(O | \lambda)} = \frac{P(i_t = q_i, i_{t+1} = q_j, O | \lambda)}{\sum_i \sum_j P(i_t = q_i, i_{t+1} = q_j, O | \lambda)} \end{aligned}$$

$$P(i_t = q_i, i_{t+1} = q_j, O | \lambda) = d_t(i) \alpha_{ij} b_j(O_{t+1}) \beta_{t+1}(j)$$

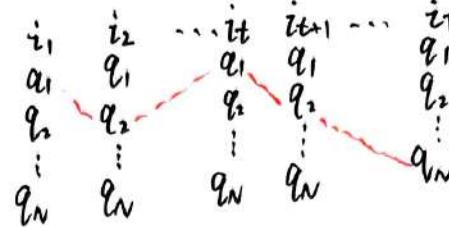
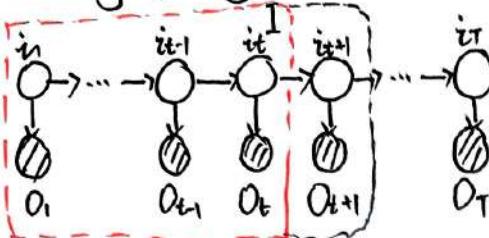
$$\delta_t(i, j) = \frac{d_t(i) \alpha_{ij} b_j(O_{t+1}) \beta_{t+1}(j)}{\sum_{i_1}^N \sum_{j_1}^M d_t(i_1) \alpha_{i_1 j_1} b_{j_1}(O_{t+1}) \beta_{t+1}(j_1)}$$

$$\text{设 } \alpha_{ij} = \frac{\sum_{t=1}^{T-1} \delta_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)}, \quad b_j(O) = \frac{\sum_{t=1}^T \delta_t(i, j)}{\sum_{t=1}^T \gamma_t(i)}, \quad \beta_{t+1}(j) = \frac{\sum_{t=1}^T \gamma_t(j)}{\sum_{t=1}^T \gamma_t(i)}$$

$$\pi(i) = \gamma_t(i)$$

13.5 Viterbi Algorithm

Decoding: $\hat{I} = \operatorname{argmax} P(I|O, \lambda) \rightarrow \text{Viterbi}$



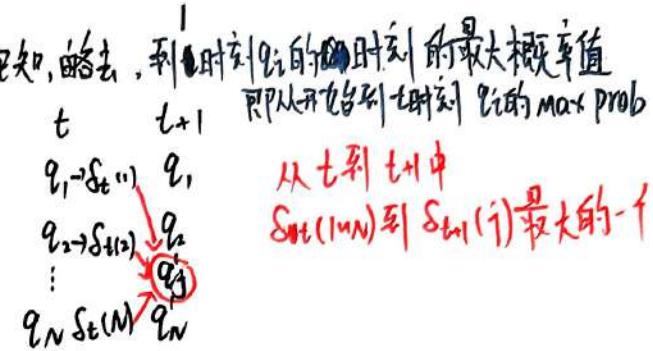
找一个概率最大的序列

$\delta_t(i) = \max_{i_1, i_2, \dots, i_{t-1}} P(O_1, O_2, \dots, O_t, i_1, i_2, \dots, i_{t-1}, i_t = q_i)$, condition 已知, 略去, 利用时刻 q_i 的前一时刻 i_{t-1} 的最大概率值
即从开始到时刻 t 的 $\max_{i_1, i_2, \dots, i_{t-1}} P(O_1, O_2, \dots, O_t, i_1, i_2, \dots, i_{t-1}, i_t = q_i)$

$\delta_{t+1}(j) = \max_{i_1, i_2, \dots, i_t} P(O_1, O_2, \dots, O_t, O_{t+1}, i_1, i_2, \dots, i_t, i_{t+1} = q_j)$

最大的概率值 $= \max_{1 \leq i \leq N} \delta_t(i) \alpha_{ij} b_j(O_{t+1})$

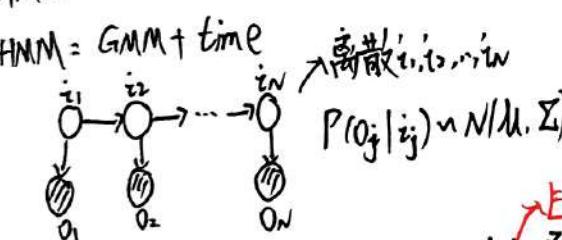
$\psi_{t+1}(j) = \operatorname{argmax}_{1 \leq i \leq N} \delta_t(i) \cdot \alpha_{ij}$ → 这个得到的是 t 时刻到 $t+1$ 时刻 q_j 处
P 最大的点 (时刻) 即 $\delta_t(i)$ 中的 i
 $\delta_t(i) \rightarrow \delta_{t+1}(j)$ 最大的路径的 i , 即 q_i 的位置



13.6 Summary

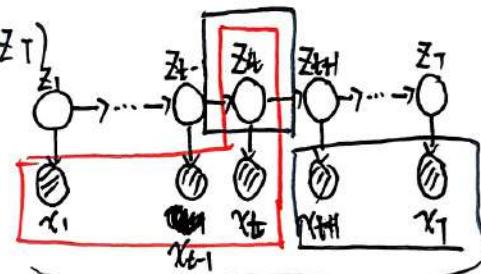
Dynamic Model:

HMM (mixture + time)

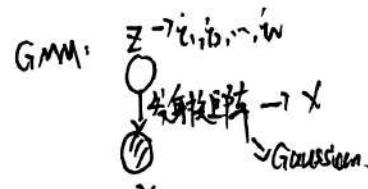


Forward: $d_t = P(x_1, \dots, x_t, z_t)$

Backward: $b_t = P(x_{t+1}, \dots, x_T | z_t)$



Dynamic Model → state space Model



{ learning: $\lambda_{MLE} = \operatorname{argmax}_{\lambda} P(x|\lambda)$ BW(CEM) }
decoding: $P(z_1, \dots, z_T | x_1, \dots, x_T)$ Viterbi Algorithm ②

{ Inference } prob of evidence: $P(x|\theta) = P(x_1, x_2, \dots, x_T | \theta)$ Forward Algorithm ③, Backward Algorithm ④

filtering: $P(z_t | x_1, x_2, \dots, x_t) \rightarrow$ Online → $P(z_1 | x_1) \rightarrow P(z_2 | x_1, x_2) \rightarrow \dots \rightarrow P(z_t | x_1, \dots, x_t)$ Forward Algorithm

smoothing: $P(z_t | x_1, x_2, \dots, x_T) \rightarrow$ Offline Forward-Backward Algorithm ⑤

prediction: $P(z_{t+1}, z_{t+2} | x_1, x_2, \dots, x_t)$ 预测下一段时序

$P(x_{t+1}, x_{t+2} | x_1, x_2, \dots, x_t)$

Filtering: $P(z_t | x_1, \dots, x_t) = \frac{P(z_t, x_1, \dots, x_t)}{P(x_1, \dots, x_t)} = \frac{P(x_1, \dots, x_t, z_t)}{\sum_{z_t} P(x_1, \dots, x_t, z_t)}$

Smoothing: $P(z_t | x_{1:T}) = \frac{P(x_1, \dots, x_t, z_t)}{P(x_1, \dots, x_T)}$

拆开前向后向

$P(x_{1:T}, z_t) = P(x_{1:T}, x_{t+1:T}, z_t) = P(x_{t+1:T} | x_{1:T}, z_t) \cdot P(x_{1:T} | z_t)$

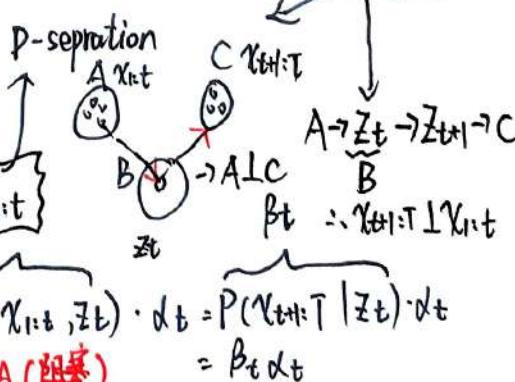
P44

x 观测值
z 隐变量

给定 z_t , A(阻塞)

$P(z_{t+1:T} | x_{1:T}, z_t) \cdot d_t = P(x_{t+1:T} | z_t) \cdot d_t = \beta_t d_t$

看图.



$$\therefore \text{Smoothing} \in P(z_t | x_{1:t}) \wedge P(x_t | z_t) = dt \beta t$$

转移矩阵

Prediction: (利用两个假设)

$$P(z_{t+1} | x_{1:t}) = \sum_{z_t} P(z_{t+1}, z_t | x_{1:t}) = \sum_{z_t} \underbrace{P(z_{t+1} | z_t, x_{1:t})}_{P(z_{t+1} | z_t)} \cdot P(z_t | x_{1:t}) = \sum_{z_t} P(z_{t+1} | z_t) \cdot P(z_t | x_{1:t})$$

$$P(x_{t+1} | x_{1:t}) = \sum_{z_{t+1}} P(x_{t+1}, z_{t+1} | x_{1:t}) = \sum_{z_{t+1}} \underbrace{P(x_{t+1} | z_{t+1}, x_{1:t})}_{P(x_{t+1} | z_{t+1})} \cdot P(z_{t+1} | x_{1:t}) = \sum_{z_{t+1}} P(x_{t+1} | z_{t+1}) \cdot P(z_{t+1} | x_{1:t})$$

已知发射矩阵

14. Linear Dynamic System (Kalman Filter)

14.1 Background

Dynamic Model (State space Model)

HMM: state is distance (离散转移矩阵, 发射矩阵)

Linear Dynamic System (Kalman filter): Linear Gaussian Model

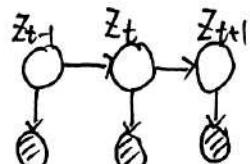
Partical Filter: Non-linear, NonGaussian → 连续, 状态和移维, 发射矩阵

Linear

$$\begin{aligned} \text{kalman filter: } z_t &= A \cdot z_{t-1} + b + \varepsilon \leftarrow \text{Gaussian Dist. } \quad z_{t-1} \quad z_t \quad z_{t+1} \\ y_t &= C \cdot z_t + d + \delta \leftarrow \text{Linear} \quad y_{t-1} \quad y_t \quad y_{t+1} \\ \varepsilon &\sim N(0, Q), \quad P(z_t | z_{t-1}) \sim N(A \cdot z_{t-1} + b, Q) \\ \delta &\sim N(0, R), \quad P(y_t | z_t) \sim N(C \cdot z_t + d, R) \\ z_t &\sim N(\mu_t, \Sigma_t) \\ \theta &= (A, B, C, D, Q, R, \mu_t, \Sigma_t) \end{aligned}$$

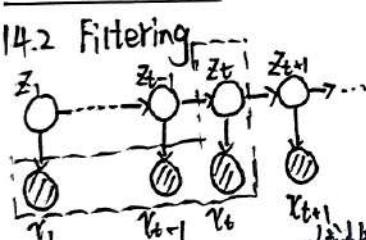
tasks
Learning
Inference

$$\text{HMM: } \lambda = (\pi, A, B)$$



$$A = [a_{ij}]_{N \times N}, a_{ij} = P(z_{t+1} = q_j | z_t = q_i)$$

$$B = [b_{j(k)}]_{N \times M}, b_{j(k)} = P(y_t = v_k | z_t = q_j)$$



$$\begin{cases} P(z_t | z_{t-1}) = N(A \cdot z_{t-1} + b, Q) \\ P(y_t | z_t) = N(C \cdot z_t + d, R) \\ P(z_1) = N(\mu_1, \Sigma_1) \end{cases}$$

Kalman Filter是Linear Gaussian Model

$$\begin{aligned} A \cdot z_t &= A \cdot z_{t-1} + b + \varepsilon, \quad \varepsilon \sim N(0, Q) \quad \text{误差 Gauss} \\ y_t &= C \cdot z_t + d + \delta, \quad \delta \sim N(0, R) \quad \text{Yt为常值} \end{aligned}$$

$$P(z_t | x_{1:t}, z_t) = \frac{P(x_1, \dots, x_t, z_t)}{P(x_1, \dots, x_t)} \propto P(x_t | x_{1:t}, z_t) = P(y_t | x_{1:t}, z_t) \cdot P(x_{1:t}, z_t) = P(x_t | z_t) \cdot P(x_{1:t}, z_t)$$

$$\begin{aligned} \text{Filtering} &= P(x_t | z_t) \cdot P(z_t | x_{1:t}), \quad P(x_{1:t}) \propto P(x_t | z_t) \cdot P(z_t | x_{1:t-1}) \\ &\quad P(x_t | z_t) \propto P(x_t | z_t) \cdot P(z_t | x_{1:t-1}) \end{aligned}$$

$$\text{故: } P(z_t | x_{1:t}, z_t) = \int_{z_{t-1}} P(z_t | z_{t-1}) \cdot P(z_{t-1} | x_{1:t-1}) dz_{t-1} = P(x_t | z_t) \cdot P(z_t | x_{1:t-1}) \int_{z_{t-1}} P(z_t | z_{t-1}, x_{1:t-1}) dz_{t-1}$$

$$\text{步骤: } t=1, P(z_1 | x_1) \rightarrow \text{update} \quad t=T, \int_{z_{t-1}} P(z_t | z_{t-1}) \cdot P(z_{t-1} | x_{1:t-1}) dz_{t-1}$$

$$\text{online } t=2, P(z_2 | x_1) \rightarrow \text{prediction} \quad \int_{z_{t-1}} P(z_t | z_{t-1}) \cdot P(z_{t-1} | x_{1:t-1}) dz_{t-1}$$

$$t=2, \begin{cases} P(z_2 | x_1, y_2) \rightarrow \text{update (correction)} \\ P(z_3 | x_1, y_2) \rightarrow \text{prediction} \end{cases}$$

Learning

Inference: $P(z_t | x)$

decoding → HMM (Viterbi)
prob of evidence → prob of

fitering: $P(z_t | x_{1:t}, z_{t-1})$
marginal posterior: $P(z_t | x_{1:t}, z_{t-1})$

prediction:
 $P(z_t | x_{1:t}, z_{t-1})$
 $P(x_t | x_{1:t}, z_{t-1})$

符合两个假设

update

Prediction

$$\begin{aligned} p(x) &= N(x|u, \Lambda^{-1}) \\ p(y|x) &= N(y|Ax+b, L^{-1}) \\ p(y) &= N(Ay|Au+b, L^{-1} + A\Lambda^{-1}A^T) \\ p(x|y) &= N(x|u + \Lambda^{-1}A^T(L^{-1} + A\Lambda^{-1}A^T)^{-1}(y - Ax - b), \Lambda^{-1} + A\Lambda^{-1}(L^{-1} + A\Lambda^{-1}A^T)^{-1}A\Lambda^{-1}) \end{aligned}$$

Filtering: (未看到 y_t) 代公式 ↑

step 1: Prediction: $p(y)$

$$P(z_t|y_1, \dots, y_{t-1}) = \int_{z_{t-1}} P(z_t|z_{t-1}) \cdot P(z_{t-1}|y_1, \dots, y_{t-1}) dz_{t-1}$$

$$p(y) \leftarrow N(u_t^*, \Sigma_t^*) = \int N(z_t|Az_{t-1}+B, Q) \cdot N(u_{t-1}, \Sigma_{t-1}) P(x)$$

step 2. update (看到 y_t)

$$P(z_t|y_1, \dots, y_t) \propto P(x_t|z_t) \cdot P(z_t|y_1, \dots, y_{t-1})$$

$$N(u_t, \Sigma_t) \propto N(y_t|Cz_t + D, R) \cdot N(u_t^*, \Sigma_t^*) P(x)$$

$$P(z_t|y_{1:t}) \propto P(y_t|z_t) \cdot P(z_t) \quad \text{在 step 1 中已经看到}$$

$$P(x|y) \quad P(y|x) \quad P(x)$$

先预测, 看到 y_t 后再做更新

这里有关于 Kalman Filter 的特性:

$$\begin{cases} \Lambda^{-1} = Z_{t-1}Z_{t-1}^T, u = U_{t-1}/U_{t-1}^T \\ L^{-1} = Q/R, A = A/C \\ b = B/D \\ y_t = z_t / \chi_t \\ P(z_t|z_{t-1}) = N(A \cdot z_{t-1} + B, Q) \\ P(y_t|z_t) = N(C \cdot z_t + D, R) \\ P(z_t) = N(u_t, \Sigma_t) \end{cases}$$

Kalman Filter 特性.

$$\begin{cases} u_t^* = A \cdot u_{t-1} + B \\ \Sigma_t^* = Q + A \cdot \Sigma_{t-1} \cdot A^T \end{cases} \Rightarrow \text{Step 1 中的参数}$$

由 z_t 与 y_t 的关系与 z_{t-1} 此时状态去测 z_t
且 (x_t, y_t) 为已观测到

$$\begin{cases} u_t = u_{t-1} + Z_{t-1}A^T(Q + A\Sigma_{t-1}A^T)^{-1}(z_t - Az_{t-1} - B) \\ \Sigma_t = \Sigma_{t-1} - \Sigma_{t-1}A^T(Q + A\Sigma_{t-1}A^T)^{-1}A\Sigma_{t-1} \end{cases}$$

看到 y_t 去根据 $z_t|y_{1:t-1}$ (t 时刻做 update)

$$\begin{cases} u_t = u_t^* + \Sigma_t^* A^T(R + \frac{C}{C} \Sigma_t^* C^T)^{-1}(y_t - A u_t^* - D) \\ \Sigma_t = \Sigma_t^* - \Sigma_t^* C^T(R + \frac{C}{C} \Sigma_t^* C^T)^{-1}C \Sigma_t^* \end{cases} \Rightarrow \text{Step 2 中的参数}$$

状态转移矩阵

$$\begin{cases} z_t = g(z_{t-1}, u, \varepsilon) \\ x_t = h(z_t, u, \delta) \end{cases} \quad \text{非线性, 噪声非高斯}$$

观测方程

15. Particle Filter

Filtering 问题: 求解: $P(z_t|y_1, \dots, y_t)$

step 1: Prediction:

$$P(z_t|y_1, \dots, y_{t-1}) = \int_{z_{t-1}} P(z_t|z_{t-1}) \cdot P(z_{t-1}|y_1, \dots, y_{t-1}) dz_{t-1} \rightarrow \text{P45. Prediction}$$

step 2. update:

$$P(z_t|y_1, \dots, y_t) \propto P(y_t|z_t) \cdot P(z_t|y_1, \dots, y_{t-1}) \rightarrow \text{P45. update}$$

15.2 Importance Sampling & SIS

Monte Carlo Method.

$$P(z|x) E_{z|x}[f(z)] = \int f(z) \cdot P(z|x) dz \approx \frac{1}{N} \sum_{i=1}^N f(z^{(i)})$$

若我们能均匀地从 $P(z|x)$ 中采样, 则有 N 个样本 $z^{(1)}, z^{(2)}, \dots, z^{(N)}$
若无法直接采样, 可通过 $q(z)$ 提议分布 proposal distribution

Importance Sampling:

$$E[f(z)] = \int f(z) P(z) dz = \int f(z) \frac{P(z)}{q(z)} \cdot q(z) dz \quad \text{①}$$

$$q(z) \text{ 是一个已知的采样的分布, } z^{(1)}, z^{(2)}, \dots, z^{(N)} \quad i=1, 2, \dots, N$$

$$\text{则 ① 式} = \frac{1}{N} \sum_{i=1}^N f(z^{(i)}) \cdot \frac{P(z^{(i)})}{q(z^{(i)})} = \frac{1}{N} \sum_{i=1}^N f(z^{(i)}) w_i^{(i)} = \sum_{i=1}^N f(z^{(i)}) \hat{w}_i^{(i)}$$

Filtering

$$P(z_t|y_1, \dots, y_t), w_t^{(i)} = \frac{P(z_t^{(i)}|y_1, \dots, y_t)}{q(z_t^{(i)}|y_1, \dots, y_t)}$$

$$t=1: w_1^{(i)}, i=1, \dots, N \rightarrow w_1^1, w_1^2, \dots, w_1^N$$

$$t=2: w_2^{(i)}, i=1, \dots, N \rightarrow w_2^1, w_2^2, \dots, w_2^N$$

由于不易求的问题, 是否找得到 $w_t^{(i)} \rightarrow w_t^{(i)}$ 有直接关系
引出:

Sequential Importance Sampling: SIS

此模型关注的是 $P(z_1, \dots, z_t|y_1, \dots, y_t)$

$$w_t^{(i)} \propto \frac{P(z_1, \dots, z_t|y_1, \dots, y_t)}{q(z_1, \dots, z_t|y_1, \dots, y_t)}$$

$$\begin{aligned} P(z_1, \dots, z_t|y_1, \dots, y_t) &= \frac{P(z_1|t, y_{1:t})}{P(x_{1:t})} = \frac{1}{C} P(z_1|t, y_{1:t}) \quad C \text{ constant} \\ &= \frac{1}{C} \cdot P(x_t|z_{1:t}, x_{1:t-1}) \cdot P(z_{1:t}, x_{1:t-1}) \\ &= \frac{1}{C} P(x_t|z_t) \cdot P(z_{1:t}, x_{1:t-1}) \\ &= \frac{1}{C} P(y_t|z_t) \cdot P(z_{1:t-1}, x_{1:t-1}) P(z_{1:t-1}, y_{1:t-1}) \\ &= \frac{1}{C} P(x_t|z_t) \cdot P(z_t|z_{t-1}) \cdot P(z_{1:t-1}|y_{1:t-1}) \cdot P(x_{1:t-1}) \\ &= \frac{P}{C} P(y_t|z_t) \cdot P(z_t|z_{t-1}) \cdot P(z_{1:t-1}|y_{1:t-1}) \end{aligned}$$

$$q(Z_{1:t} | X_{1:t}) = q(Z_t | Z_{1:t-1}, X_{1:t}) \cdot q(Z_{1:t-1} | X_{1:t})$$

假定有: $q(Z_{1:t} | X_{1:t}) = q(Z_{1:t-1} | X_{1:t-1}) \cdot q(Z_t | Z_{1:t-1}, X_{1:t})$, q 是我们指定的, 具有任意性, 独立性
根据条件独立, 若分布也符合Markov齐次, 则
此步 $\alpha = \frac{P(X_t | Z_t) \cdot P(Z_t | Z_{t-1})}{q(Z_t | Z_{1:t-1}, X_{1:t})}$

$$W_t^{(i)} \propto \frac{P(Z_{1:t} | X_{1:t})}{q(Z_{1:t} | X_{1:t})} \propto \frac{P(Y_t | Z_t) \cdot P(Z_t | Z_{t-1}) / P(Z_{1:t-1} | X_{1:t-1})}{q(Z_t | Z_{1:t-1}, X_{1:t})} = \frac{P(Y_t | Z_t) \cdot P(Z_t | Z_{t-1})}{q(Z_t | Z_{1:t-1}, X_{1:t})} \cdot W_{t-1}^{(i)} \rightarrow \text{权重计算}$$

假定初值时 $W_1^{(i)}$ 已算好.

则 $t=2$. $W_2^{(i)} = \alpha W_1^{(i)}$ 以此类推, 解决了 W 计算难的问题

权值退化, 随着 t 增大,
 $\lim_{t \rightarrow \infty} W_t = 0$, 因为每一个 W_t 会变小,
而 $\sum_i W_i = 1$, 就像铺满 $[0,1]$ 区间, 权值退化

15.3 Particle Filter \leftrightarrow Re-Sampling

Sequential Importance Sampling Filter (SIS)

$$W_t^{(i)} \propto \frac{P(Z_{1:t} | X_{1:t})}{q(Z_{1:t} | X_{1:t})} \propto \frac{P(Y_t | Z_t) P(Z_t | Z_{t-1})}{q(Z_t | Z_{1:t-1}, X_{1:t})} W_{t-1}^{(i)}$$

问题: 权值退化 $W_t^{(i)}$

比如有 $100 W_t^{(i)}$, 99 个都为 0.001, 1 个接近 1

solution: Re-Sampling (Basic Particle Filter)

选择一个合适的 proposal dist $q(z)$ (SIR Filter) Basic Particle Filter 衡量权重

Basic Particle Filter: SIS + Resampling 最基本的 resampling 方法 加上的 Resampling

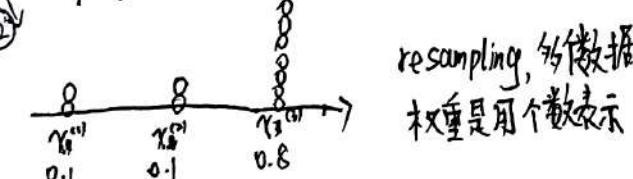
假设采样完成得到 X_1, X_2, X_3

①

第一次采样, 得到 3 个数据
以及各自权重

$$\begin{array}{c|ccc} & X_1^{(1)} & X_2^{(2)} & X_3^{(3)} \\ \hline & 0.1 & 0.1 & 0.8 \end{array}$$

resampling (以权重作为概率分布采样)



②

resampling, 多数据

权重是用个数表示

algorithm

前提 $t-1$ 时刻采样已完成 $\Rightarrow W_{t-1}^{(i)}$ 已知

t 时刻: for $i=1, \dots, N$

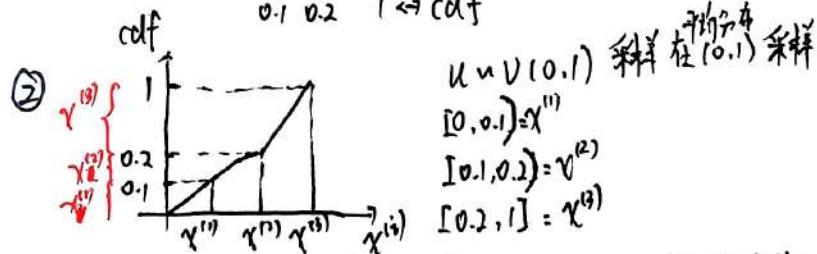
$$Z_t^{(i)} \sim q(Z_t | Z_{t-1}^{(i)}, X_{1:t})$$

$$W_t^{(i)} \propto W_{t-1}^{(i)} \cdot \frac{P(Y_t | Z_t^{(i)})}{q(Z_t^{(i)} | Z_{t-1}^{(i)}, X_{1:t})}$$

② Normalized $W_t^{(i)} \rightarrow 1$, $\sum_i W_t^{(i)} = 1$ $\Rightarrow W_t^{(i)}$

③ Resampling $Z_t^{(i)} \sim \hat{W}_t^{(i)} = \frac{1}{N}$, 所有粒子权重相同, 以做衡衡量权重

① 对于 3 个样本 $X^{(1)}, X^{(2)}, X^{(3)}$
0.1 0.1 0.8 \rightarrow pdf
0.1 0.2 1 \leftrightarrow cdf



总结: 这样就避免了权值退化,一开始可采样多个样本,之后停止采样,做 cdf 图,通过在 $[0,1]$ 平均分布中采样,落到哪个样本区间就采样哪个

15.4 SIR Filter 选择一个更好的提议分布

Sampling - Importance - Resampling: SIS + Resampling + $q(Z_t | Z_{1:t-1}, X_{1:t}) = P(Z_t | Z_{t-1})$

algorithm:

前提: $t-1$ 时刻采样已完成 $\Rightarrow W_{t-1}^{(i)}$ 已知

t 时刻

① Sampling: for $i=1, \dots, N$:

$$Z_t^{(i)} \sim P(Z_t | Z_{t-1}^{(i)})$$

$$W_t^{(i)} \sim P(Y_t | Z_t^{(i)}) \cdot W_{t-1}^{(i)}$$

end

② Normalized: $W_t^{(i)} \rightarrow \hat{W}_t^{(i)}$, $\sum_i \hat{W}_t^{(i)} = 1$

③ Resampling $\hat{W}_t^{(i)} = \frac{1}{N}$

"generate and test" 从自己身上找 $Z_t^{(i)}$

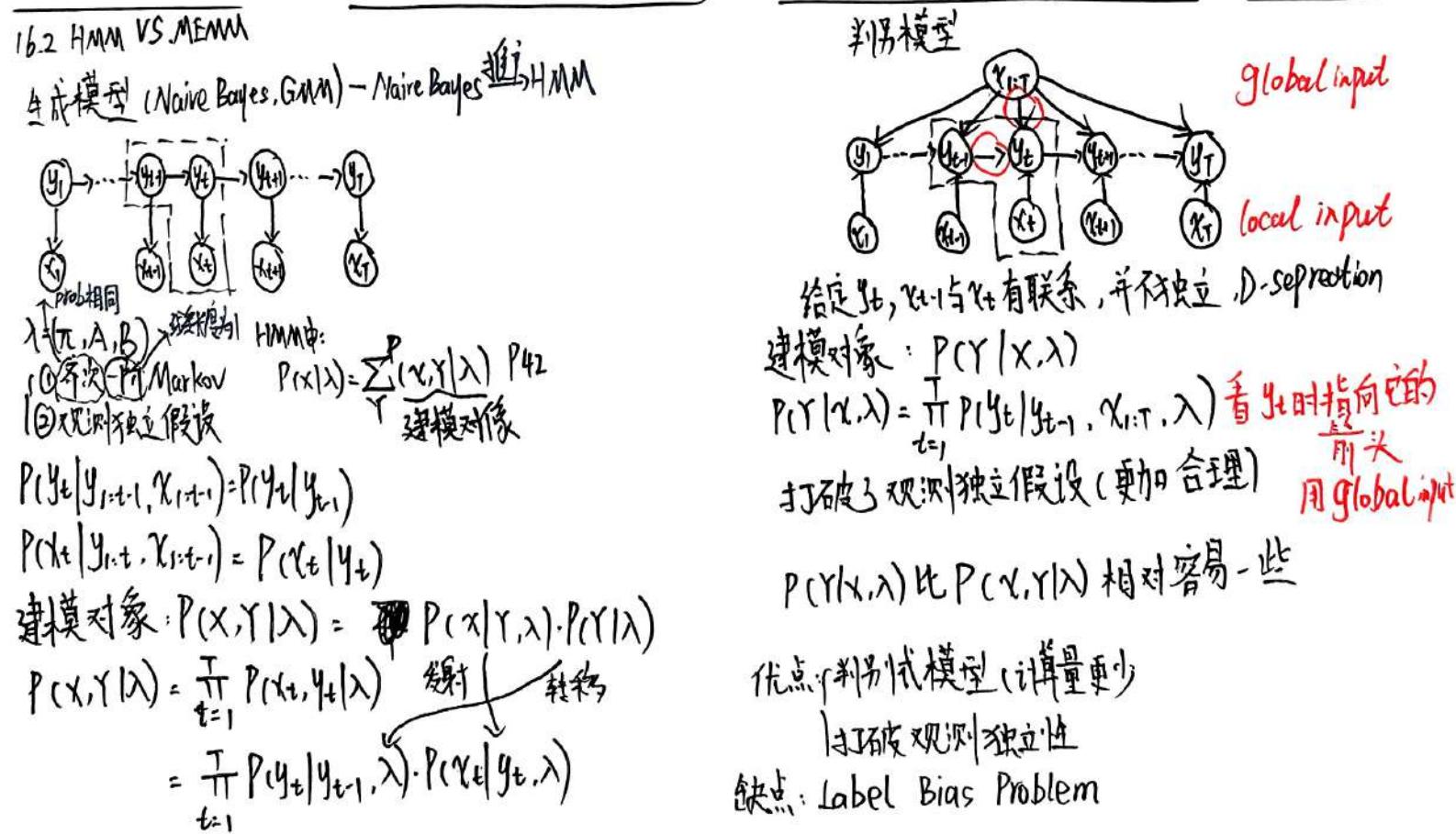
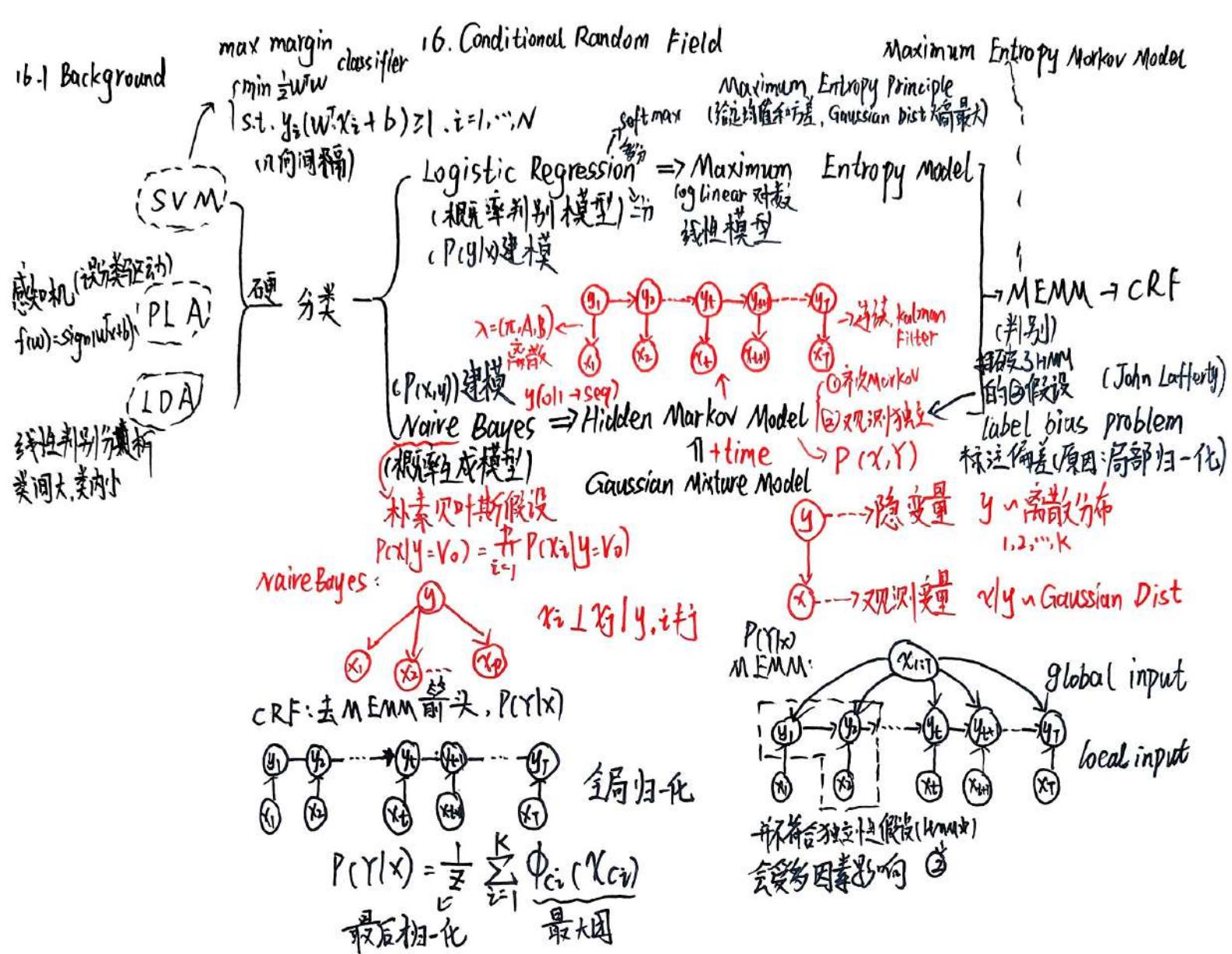
$$Z_t^{(i)} \sim P(Z_t | Z_{t-1}^{(i)})$$

$W_t^{(i)} \sim P(Y_t | Z_t^{(i)}) \cdot W_{t-1}^{(i)}$

新生成的

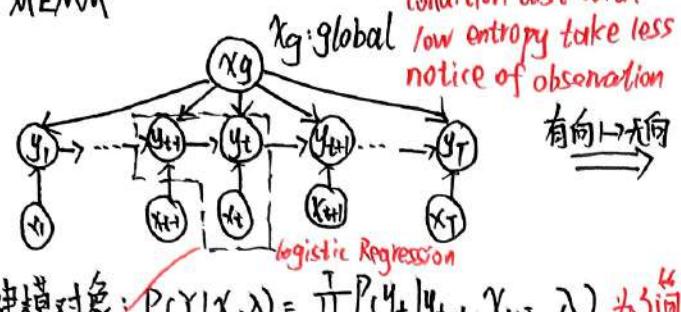
上一个权重

若 $W_t^{(i)}$ 越大, 则说明 $Z_t^{(i)}$ 取的越合适, 更能反映分布



16.3 MEMM VS. CRF

MEMM:

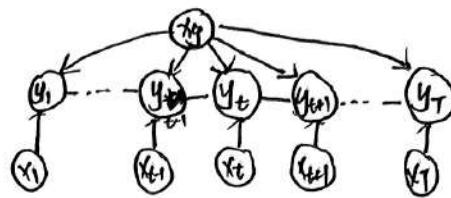


condition dist with
low entropy take less
notice of observation

有向 \$\rightarrow\$ 无向

建模对象: $P(Y|X, \lambda) = \prod_{t=1}^T P(y_t | y_{t-1}, x_{1:T}, \lambda)$ 为简化
省去 \$y_t\$ 的初态
状态的表述
\$y_t\$ 受两个 \$y_{t-1}\$ 与 \$x_t\$ 影响
mass score
function \$\geq 0\$
根本原因
由 \$y_{t-1}, y_t, x_t\$ 组成

chain-structured CRF

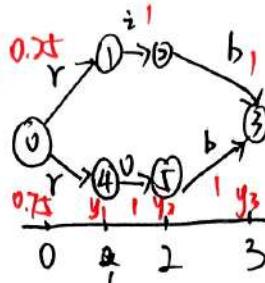


克服了 Label Bias problem

$$P(2|1, i) = 1 = P(2|1)$$

$$P(5|4, 0) = 1 = P(5|4)$$

可以看出从 \$0 \rightarrow 2, 4 \rightarrow 5\$, 根本没有关注 observation (\$i, o\$)



缺点: Label Bias Problem

Decoding: $Y = (y_1, y_2, y_3)$

$$Y = \arg \max_{y_1, y_2, y_3} P(y_1, y_2, y_3 | \text{rib})$$

$$= 0 \rightarrow 4 \rightarrow 5 \rightarrow 3$$

= rob. 显然不合理 (这就是归一化, 忽视 observation 的结果) 就是说, 图中 \$1 \rightarrow 2\$; \$2 \rightarrow 3, 4 \rightarrow 5, 5 \rightarrow 3\$, 上的 observation 没有任何作用, 只与 training 数据相关

eg. 样本 4 个, 3 个 rob, 1 个 rib, 则在图上标注概率

$$\text{training: } P(1|00, r) = 0.25$$

$$P(4|0, r) = 0.75$$

viterbi: Decoding.

K 个最大团
C_i: 最大团

16.4 pdf-parameter perspective

MRF 因子分解: Markov Random Field

$$P(x) = \frac{1}{Z} \prod_{i=1}^K \psi_i(x_{ci}) = \frac{1}{Z} \prod_{i=1}^K \exp [-E_i(x_{ci})] = \frac{1}{Z} \exp \sum_{i=1}^K F_i(x_{ci})$$



$$P(Y|X) = \frac{1}{Z} \exp \sum_{i=1}^K F_i(x_{ci})$$

$$= \frac{1}{Z} \exp \sum_{t=1}^T f_t(y_{t-1}, y_t, x_{1:T})$$

由于是线性链, \$f_t\$ 其实有相同

$$= \frac{1}{Z} \exp \sum_{t=1}^T F_t(y_{t-1}, y_t, x_{1:T})$$

$$= \frac{1}{Z} \exp \sum_{t=1}^T [\sum_{k=1}^K \lambda_k f_k(y_{t-1}, y_t, x_{1:T}) + \sum_{l=1}^L y_l g_l(y_t, x_{1:T})]$$

CRF-PDF

$$P(Y|X) = \frac{1}{Z} \exp \sum_{t=1}^T [\sum_{k=1}^K \lambda_k f_k(y_{t-1}, y_t, x_{1:T}) + \sum_{l=1}^L y_l g_l(y_t, x_{1:T})]$$

状态函数

转移函数

$$\Delta y_{t-1}, y_t, x_{1:T} = \Delta y_{t-1, n} + \Delta y_{t, m} + \Delta y_{t, v, adj_m}$$

$$\Delta y_{t-1, y_t, x_{1:T}} = \Delta y_{t-1, n} + \Delta y_{t, m}$$

$$y_{t-1}, y_t \in \{n, v, adj_m\}$$

f_k : 特征函数

$$\Delta y_{t-1, y_t, x_{1:T}} = \sum_{k=1}^K \lambda_k f_k(y_{t-1}, y_t, x_{1:T})$$

$$f_k = \begin{cases} 1 & y_{t-1} = n, y_t = v, x_{1:T} \\ 0 & \text{otherwise} \end{cases}$$

$$\Delta y_{t, x_{1:T}} = \sum_{l=1}^L y_l g_l(y_t, x_{1:T})$$

$$f_k = \begin{cases} 1 & y_{t-1} = n, y_t = v, x_{1:T} \\ 0 & \text{otherwise} \end{cases}$$

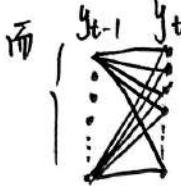
P49 f_k, g_l 是给定的特征函数, λ_k, y_l 参数

16.5 pdf-vector perspective

$$P(Y|X) = \frac{1}{Z} \exp \sum_{t=1}^T [\Delta_{\text{转移}} + \Delta_{\text{状态}}] = \frac{1}{Z} \exp \sum_{t=1}^T \left[\sum_{k=1}^K \lambda_k f_k(y_{t-1}, y_t, X_{1:T}) + \sum_{l=1}^L \gamma_l g_l(y_t, X_{1:T}) \right]$$

关于 K, L 的取值

假设 $y_t \in S = \{\text{动名副助, } \dots\}$



共 $|S|$ 个单

总共 $|S|^2$ 个组合

则 $K \leq |S|^2$, 且重要程度由 λ 决定

λ 的取值与 y_t 的集合中数相关

化简替换:

$$y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_T \end{pmatrix}; x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_T \end{pmatrix}; \lambda = \begin{pmatrix} \lambda_1 \\ \lambda_2 \\ \vdots \\ \lambda_K \end{pmatrix}; \gamma = \begin{pmatrix} \gamma_1 \\ \gamma_2 \\ \vdots \\ \gamma_L \end{pmatrix};$$

$$f = \begin{pmatrix} f_1 \\ f_2 \\ \vdots \\ f_K \end{pmatrix} = f(y_{t-1}, y_t, x); g = \begin{pmatrix} g_1 \\ g_2 \\ \vdots \\ g_L \end{pmatrix} = g(y_t, x)$$

$$\text{则有: } P(Y=y|X=x) = \frac{1}{Z(\chi, \lambda, \gamma)} \exp \sum_{t=1}^T [\lambda^T \cdot f(y_{t-1}, y_t, x) + \gamma^T \cdot g(y_t, x)]$$

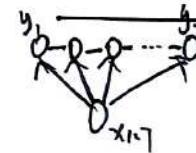
$$= \frac{1}{Z(\chi, \lambda, \gamma)} \exp \left[\lambda^T \cdot \sum_{t=1}^T f(y_{t-1}, y_t, x) + \gamma^T \cdot \sum_{t=1}^T g(y_t, x) \right]$$

$$\Theta = \begin{pmatrix} \lambda \\ \gamma \end{pmatrix}_{K+L} \quad H = \begin{pmatrix} \sum_{t=1}^T f \\ \sum_{t=1}^T g \end{pmatrix}_{K+L}$$

$$\therefore P(Y=y|X=x) = \frac{1}{Z(\Theta, H)} \exp [\Theta^T \cdot H(y_t, y_{t-1}, x)]$$

$$P(y|x) = \frac{1}{Z(\Theta, H)} \exp [\langle \Theta, H \rangle]$$

16.6. model problems Given training data: $\{(x^{(i)}, y^{(i)})\}_{i=1}^N$, x, y 均是 T 维



Learning: parameter estimation $\hat{\theta} = \arg \max \sum_{i=1}^N P(y^{(i)}|x^{(i)})$ ②

Inference: $\begin{cases} \text{marginal prob: } P(x_1), \dots, P(x_T) \\ \text{conditional prob: } P(y_t|x) \end{cases}$ 生成模型 ①

MAP Inference: decoding $\arg \max_y P(y|x)$ ③

$$\hat{y} = \arg \max_{y=y_1, \dots, y_T} P(y|x)$$

16.7 marginal probability

Given $P(Y=y|X=x)$, find $P(y_t=i|x)$

$$P(y|x) = \frac{1}{Z} \prod_{t=1}^T \psi_t(y_{t-1}, y_t, x)$$

状态
 $y_t \in S$. 考虑 $O(|S|^T, T)$ 为最大

$$P(y_t=i|x) = \sum_{y_1, y_2, \dots, y_{t-1}, y_t} P(y|x) = \sum_{y_1, y_2, \dots, y_{t-1}} \sum_{y_t} \frac{1}{Z} \prod_{t'=1}^T \psi_{t'}(y_{t'-1}, y_{t'}, x) = \frac{1}{Z} \sum_{y_1, y_2, \dots, y_{t-1}} \psi_t(y_{t-1}, y_t=i, x) \psi_1(y_0, y_1, x) \psi_2(y_1, y_2, x) \dots \psi_{t-1}(y_{t-2}, y_{t-1}, x)$$

$$\cdot \psi_t(y_{t-1}, y_t=i, x) \cdot \sum_{y_{t+1}, \dots, y_T} \psi_{t+1}(y_t=i, y_{t+1}, x) \dots \psi_T(y_{T-1}, y_T, x) = \frac{1}{Z} \sum_{y_{t-1}} \psi_t(y_{t-1}, y_t=i, x) \sum_{y_{t+1}, \dots, y_T} \psi_{t+1}(y_{t-1}, y_{t+1}, x) \dots \psi_T(y_{T-1}, y_T, x) \dots$$

$$\left[\sum_{y_1} \psi_t(y_1, y_2, x) \left(\sum_{y_0} \psi_1(y_0, y_1, x) \right) \right] \sum_{y_{t+1}} \psi_{t+1}(y_{t-1}, y_{t+1}, x) \dots \sum_{y_T} \psi_T(y_{T-1}, y_T, x) + \Delta_{\text{左}} + \Delta_{\text{右}} = \frac{1}{Z} (\Delta_{\text{左}} + \Delta_{\text{右}})$$

类似这样拆分

从后往前求和得

$$\Delta_{\text{左}} = \sum_{y_{t-1}} \psi_t(y_{t-1}, y_{t-i}, x) \sum_{y_{t-2}} \psi_{t-1}(y_{t-2}, y_{t-1}, x) \dots \sum_{y_1} \psi_2(y_1, y_i, x) \sum_{y_0} \psi_1(y_0, y_i, x)$$

求递推:

$$\Delta_{\text{左}} \Leftarrow d_{t-1}(i) : y_0, y_1, y_2, \dots, y_{t-1}, \underbrace{y_t = i}_{\text{所有势函数的乘积}} \quad \begin{array}{c} \psi_t \\ \text{右半部势函数} \end{array} \quad \begin{array}{c} \psi_{t-1} \\ \text{左半部势函数} \end{array}$$

$$d_{t-1}(j) : \underbrace{y_0, y_1, y_2, \dots, y_{t-1}}_{\text{所有势函数}}, \underbrace{y_t = j}_{\text{右半部势函数}}$$

$$\psi_t(y_{t-1}=j, y_t=i, x) \cdot d_{t-1}(j)$$

$$\sum_{y_{t-1}=j} \psi_t(y_{t-1}=j, y_t=i, x) \cdot d_{t-1}(j)$$

$$d_{t-1}(j) = \sum_{j \in S} \psi_t(y_{t-1}=j, y_t=i, x) \cdot d_{t-1}(j) = d_{t-1}(i) \rightarrow \text{递推}$$

$$\therefore \Delta_{\text{左}} = d_{t-1}(i)$$

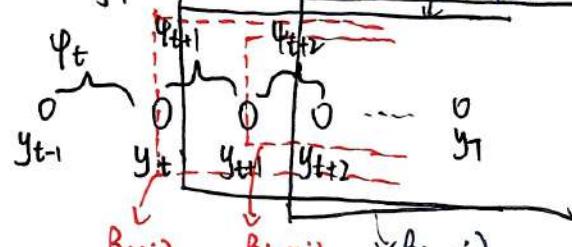
$$\Delta_{\text{右}} = \sum_{t+1, T} \psi_{t+1}(y_t=i, y_{t+1}, x) \cdot \psi_{t+2}(y_{t+1}, y_{t+2}, x) \dots \psi_{T-1}(y_{T-2}, y_{T-1}, x) \cdot \psi_T(y_{T-1}, y_T, x)$$

$$= \sum_{y_{t+1}} \psi_{t+1}(y_t=i, y_{t+1}, x) \cdot \sum_{y_{t+2}} \psi_{t+2}(y_{t+1}, y_{t+2}, x) \dots \sum_{y_{T-1}} \psi_{T-1}(y_{T-2}, y_{T-1}, x) \sum_{y_T} \psi_T(y_{T-1}, y_T, x) \quad \beta_{t-1}(i)$$

求递推:

$$\beta_{t-1}(i) = y_{t+1} y_{t+2} \dots y_T \quad \begin{array}{c} y_t \text{ 的右半部} \\ \text{所有势函数} \end{array}$$

$$\text{即 } \beta_{t-1}(i) = y_t = i \quad y_{t+1} y_{t+2} \dots y_T \quad \text{从 } t \text{ 到 } T \quad \begin{array}{c} y_t \\ \text{右半部} \\ \text{所有势函数} \end{array}$$



red more clear

$$\beta_{t+1}(j) = y_{t+1} = j \quad y_{t+2} y_{t+3} \dots y_T \quad \text{从 } t+1 \text{ 到 } T \quad \begin{array}{c} y_{t+1} \\ \text{右半部} \\ \text{所有势函数} \end{array}$$

$$\psi_{t+1}(y_t=i, y_{t+1}=j, x) \cdot \beta_{t+1}(j)$$

$$\beta_{t-1}(i) = \sum_{j \in S} \psi_{t+1}(y_t=i, y_{t+1}=j, x) \beta_{t+1}(j)$$

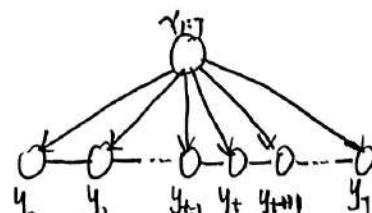
$$\therefore \Delta_{\text{右}} = \beta_{t-1}(i)$$

$$\text{有: } \beta_{t-1}(i) = y_t \text{ 的右半部} + \text{所有 } (y_{t+1} - y_T)$$

$$\beta_{t+1}(j) = y_{t+1} \text{ 的右半部} + \text{所有 } (y_{t+2} - y_T)$$

$$\beta_{t-1}(i) = \beta_{t+1}(j) + y_{t+1} \text{ 的左半部}$$

可以用 Inference: 变量消除法: sum product / Belief propagation
Forward-Backward 也可以理解上述过程



$$\therefore P(y_t=i|x) = \frac{1}{Z} d_{t-1}(i) \beta_{t-1}(i)$$

16.8 Learning

$$\hat{\theta} = \operatorname{argmax}_{\theta} \prod_{i=1}^N P(y^{(i)} | x^{(i)}) \quad N: \text{size of training data}$$

$$\lambda, \eta = \operatorname{argmax}_{\lambda, \eta} \prod_{i=1}^N P(y^{(i)} | x^{(i)})$$

$$P(y|x) = \frac{1}{Z(x, \lambda, \eta)} \exp \sum_{t=1}^T [\lambda^T \cdot f(y_{t-1}, y_t, x) + \eta^T \cdot g(y_t, x)]$$

$$= \operatorname{argmax}_{\lambda, \eta} \log \prod_{i=1}^N P(y^{(i)} | x^{(i)})$$

$$= \operatorname{argmax}_{\lambda, \eta} \sum_{i=1}^N \log P(y^{(i)} | x^{(i)})$$

$$= \operatorname{argmax}_{\lambda, \eta} \sum_{i=1}^N \left\{ -\log Z(x^{(i)}, \lambda, \eta) + \sum_{t=1}^T [\lambda^T \cdot f(y_{t-1}^{(i)}, y_t^{(i)}, x) + \eta^T \cdot g(y_t^{(i)}, x)] \right\}$$

$Z(x)$ 充分统计量

$$\triangleq \operatorname{argmax}_{\lambda, \eta} L(\lambda, \eta, x^{(i)})$$

方法：梯度上升. $\nabla_\lambda L$, $\nabla_\eta L$

$$\nabla_\lambda L = \sum_{i=1}^N \left[\sum_{t=1}^T f(y_{t-1}^{(i)}, y_t^{(i)}, x^{(i)}) - \underbrace{\frac{\nabla_\lambda \log Z(x^{(i)}, \lambda, \eta)}{\log \text{Partition Function}}}_{\text{log-Partition Function}} \right]$$

$$\begin{aligned} P(x|h) &= h(x) \exp \{ \eta^T \phi(x) - A(h) \}, \\ A(h) &= E_{P(x|h)}[\phi(x)] \end{aligned}$$

P23. 7.3.

$$\begin{aligned} &\rightarrow E \left[\sum_{t=1}^T f(y_{t-1}^{(i)}, y_t^{(i)}, x^{(i)}) \right] \\ &= \sum_y P(y|x^{(i)}) \cdot \sum_{t=1}^T f(y_{t-1}^{(i)}, y_t^{(i)}, x^{(i)}) \\ &= \sum_{t=1}^T \left[\sum_y P(y|x^{(i)}) \cdot f(y_{t-1}^{(i)}, y_t^{(i)}, x^{(i)}) \right] \\ &= \sum_{t=1}^T \sum_{y_{t+1, t+2}} \sum_{y_t} \sum_{y_{t+1, t+2}} P(y|x^{(i)}) \cdot f(y_{t-1}^{(i)}, y_t^{(i)}, x^{(i)}) \\ &= \sum_{t=1}^T \sum_{y_{t-1}} \sum_{y_t} \left(\sum_{y_{t+1, t+2}} \sum_{y_{t+1, t+2}} P(y|x^{(i)}) \cdot f(y_{t-1}^{(i)}, y_t^{(i)}, x^{(i)}) \right) \\ &= \sum_{t=1}^T \sum_{y_{t-1}} \sum_{y_t} \underbrace{P(y_{t-1}, y_t | x^{(i)}) f(y_{t-1}, y_t, x^{(i)})}_{\substack{\text{与PSO计算 } P(y|x) \text{ 的方法雷同} \\ \text{假设已求出, } P(y_{t-1}, y_t | x^{(i)}) = A(y_{t-1}, y_t)}} \end{aligned}$$

$$\nabla_\eta L = \sum_{i=1}^N \left[\sum_{t=1}^T [f(y_{t-1}^{(i)}, y_t^{(i)}, x^{(i)}) - \sum_{y_{t-1}} \sum_{y_t} A(y_{t-1}, y_t) \cdot f(y_{t-1}, y_t, x^{(i)})] \right]$$

梯度上升算法:

$$\lambda^{(t+1)} = \lambda^{(t)} + \text{step} \cdot \nabla_\lambda L(\lambda^{(t)}, \eta^{(t)})$$

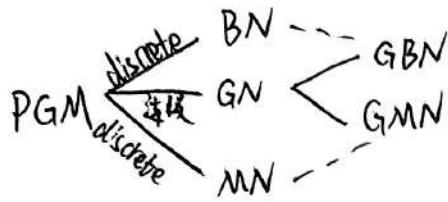
$$\eta^{(t+1)} = \eta^{(t)} + \text{step} \cdot \nabla_\eta L(\lambda^{(t)}, \eta^{(t)})$$

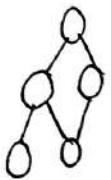
HMM: decoding \rightarrow Viterbi \rightarrow 动态规划

CRF: decoding

17. Gaussian Network

17.1 Background



 $X_i \sim N(\mu_i, \Sigma_i)$
 $X = (X_1, X_2, \dots, X_p)^T$

$$P(x) = \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}} \cdot \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

$$\Sigma = (6_{ij}) = \begin{pmatrix} 6_{11} & 6_{12} & \cdots & 6_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ 6_{p1} & 6_{p2} & \cdots & 6_{pp} \end{pmatrix}_{p \times p}$$

$$X_i \perp X_j \Leftrightarrow \sigma_{ij} = 0$$

条件独立性: $X_A \perp X_B \mid X_C$

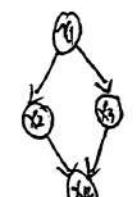
$$\Lambda = \Sigma^{-1} = \begin{pmatrix} \lambda_{11} & \lambda_{12} & \cdots & \lambda_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_{p1} & \lambda_{p2} & \cdots & \lambda_{pp} \end{pmatrix}_{p \times p}$$

precision matrix \Rightarrow information matrix

$$X_i \perp X_j \mid \underbrace{\{X_i, X_j\}}_{\text{除 } X_i, X_j \text{ 外}} \Leftrightarrow \lambda_{ij} = 0$$

17.2 Gaussian Bayesian Network

连续型的PGM. 有向GBN



$$P(x) = \prod_{i=1}^p P(X_i \mid \underbrace{X_{\text{pa}(i)}}_{\text{-一个集合(父亲)}}) \text{ BN 的因子分解 ... ①}$$

GBN is based on linear Gaussian Model

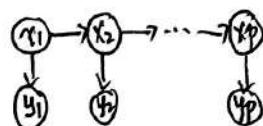
属于GBN (chain structure)
Kalman Filter (HMM). $\lambda = (\pi, A, B)$

$$P(x_t | y_{t-1})$$

$$P(y_t | v_t)$$

$$x_{\text{pa}(i)} = \{x_2, x_3\}$$

$$\text{global model}$$



Review.
kalman filter.

$$P(x_t | y_{t-1}) = N(x_t | A x_{t-1} + B, Q)$$

$$P(y_t | x_t) = N(y_t | C x_t + D, R)$$

$$\begin{cases} x_t = Ax_{t-1} + B + \varepsilon, \quad \varepsilon \sim N(0, Q) \\ y_t = Cx_t + D + S, \quad S \sim N(0, R) \end{cases}$$

$$\begin{cases} x_t = Ax_{t-1} + B + \varepsilon, \quad \varepsilon \sim N(0, Q) \\ y_t = Cx_t + D + S, \quad S \sim N(0, R) \end{cases}$$

$$\left. \begin{array}{l} x \sim N(\mu, \Sigma) \\ x = (x_1, x_2, \dots, x_p)^T \\ \mu = (\mu_1, \mu_2, \dots, \mu_p)^T \\ \varepsilon = (\varepsilon_1, \varepsilon_2, \dots, \varepsilon_p)^T \\ \Sigma = \text{diag}(\sigma_i^2) \\ W = [w_{ij}] \rightarrow \text{无 parents 处理 0} \end{array} \right\} *$$

$$\left. \begin{array}{l} P(x) = N(x \mid \mu_x, \Sigma_x) \\ P(y|x) = N(y \mid Ax + b, \Sigma_y) \\ P(x) = \prod_{i=1}^p P(X_i \mid X_{\text{pa}(i)}) \\ X_{\text{pa}(i)} = (x_1, x_2, \dots, x_k)^T \end{array} \right.$$

$$P(X_i \mid X_{\text{pa}(i)}) = N(X_i \mid \mu_i + W_i^T \cdot X_{\text{pa}(i)}, \sigma_i^2)$$

X_i 是 - 随机的

$$X_i = \mu_i + \sum_{j \in \text{pa}(i)} w_{ij} \cdot (X_j - \mu_j) + \varepsilon_i, \quad \varepsilon_i \text{ is r.v., } \varepsilon_i \sim N(0, 1)$$

$$X_i - \mu_i = \sum_{j \in \text{pa}(i)} w_{ij} (X_j - \mu_j) + \varepsilon_i, \quad X - \mu = W \cdot (X - \mu) + S \cdot \varepsilon$$

$$(I - W) \cdot (X - \mu) = S \cdot \varepsilon$$

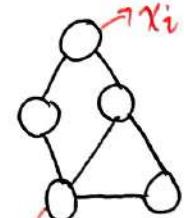
$$X - \mu = (I - W)^{-1} \cdot S \cdot \varepsilon$$

$$\left. \begin{array}{l} \Sigma = \text{cov}(x) = \text{cov}(X - \mu) \\ = \text{cov}[(I - W)^{-1} \cdot S \cdot \varepsilon] \end{array} \right.$$

17.3 Gaussian MRF learning: 结构参数

$$P(x) = \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}} \cdot \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\} \rightarrow \text{pdf (Gaussian dist)}$$

$$P(x) = \frac{1}{Z} \prod_{i=1}^p \underbrace{\psi_i(x_i)}_{\text{node potential}} \prod_{i,j \in \mathcal{E}} \underbrace{\psi_{i,j}(x_i, x_j)}_{\text{edge potential}}$$



x_i, x_j 边
 $\lambda_{ij} = 0$

$x_i = \begin{pmatrix} x_i \\ x_{ip} \end{pmatrix}; \Delta = (\lambda_{ij})_{pp}$

相等: $x_{ii} - \frac{1}{2} x_i \cdot \lambda_{ii} + h_i x_i$

不同: $x_i, x_j : -\frac{1}{2} (\lambda_{ij} x_i x_j + \lambda_{ji} x_j x_i)$, Δ 对称矩阵

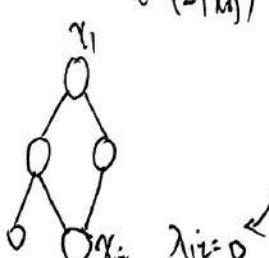
① $x_i \perp x_j, (\Sigma = (\delta_{ij})) \Leftrightarrow \delta_{ij} = 0 \rightarrow \text{marginal independence}$

② $x_i \perp x_j |_{-\{x_i, x_j\}} (\Delta = \Sigma^{-1} = (\lambda_{ij})) \Leftrightarrow \lambda_{ij} = 0 \rightarrow \text{条件独立}$

③ $\forall x_i, x_i |_{-\{x_i\}} \sim N \left(\sum_{j \neq i} \frac{\lambda_{ij}}{\lambda_{ii}} x_j, \lambda_{ii}^{-1} \right)$

解概率分布

$$x_i = \begin{pmatrix} x_i \\ x_{ip} \end{pmatrix} = \begin{pmatrix} x_a \\ x_b \end{pmatrix}$$



Δ : precision matrix

$\Delta \mu$: potentiation vector

$$\Delta \mu = \begin{pmatrix} h_1 \\ h_2 \\ \vdots \\ h_p \end{pmatrix}$$

$$\begin{aligned} & (2\pi)^{\frac{p}{2}}, |\Sigma|^{\frac{1}{2}} \text{ 都是 parameters} \quad \Delta = \Sigma^{-1} \\ & \therefore P(x) \propto \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\} = \exp \left\{ -\frac{1}{2} (x^T \Delta - \mu^T \Delta) (x - \mu) \right\} \\ & = \exp \left\{ -\frac{1}{2} (x^T \Delta x - x^T \Delta \mu - \mu^T \Delta x + \mu^T \Delta \mu) \right\} \\ & = \exp \left\{ -\frac{1}{2} (x^T \Delta x - 2 \mu^T \Delta x + \mu^T \Delta \mu) \right\} \\ & \propto \exp \left\{ -\frac{1}{2} x^T \Delta x + \mu^T \Delta x \right\} = \exp \left\{ -\frac{1}{2} x^T \Delta x + (\Delta \mu)^T x \right\} \end{aligned}$$

把 ψ_i 看作 node potential
把 x_i, x_j 看作 edge potential

此处把 Gaussian Dist 与 极端关联

全局/绝对独立

$x_{ij} \neq 0$, 固中有在 (x_i, x_j)

$x_{ij} = 0$, 固中不存在 (x_i, x_j)

由 P6 全局 Markov Property.

若两个节点 i 和 j 被集合 S 分离 (即该两节点间所有路径都经过 S) 则

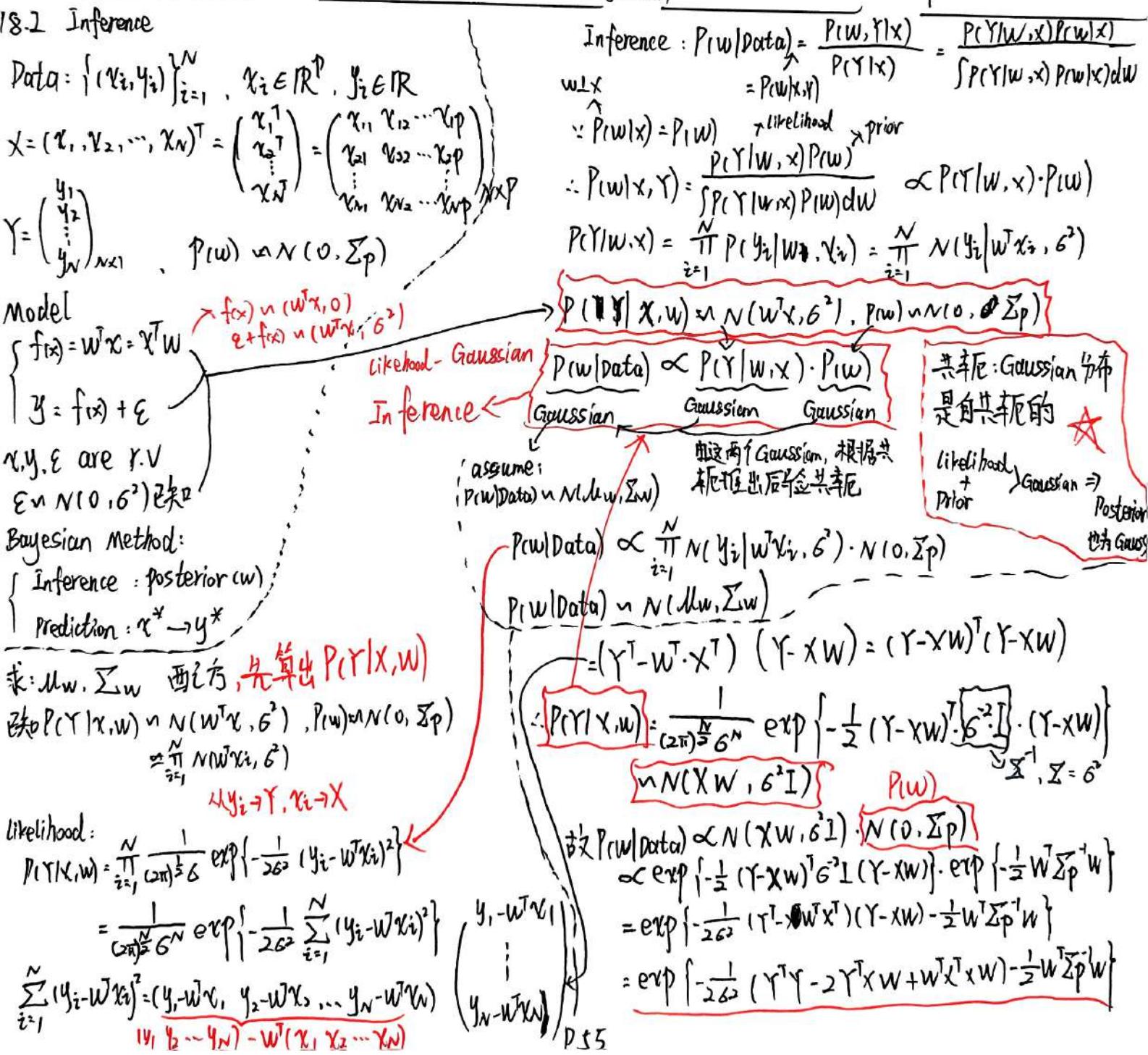
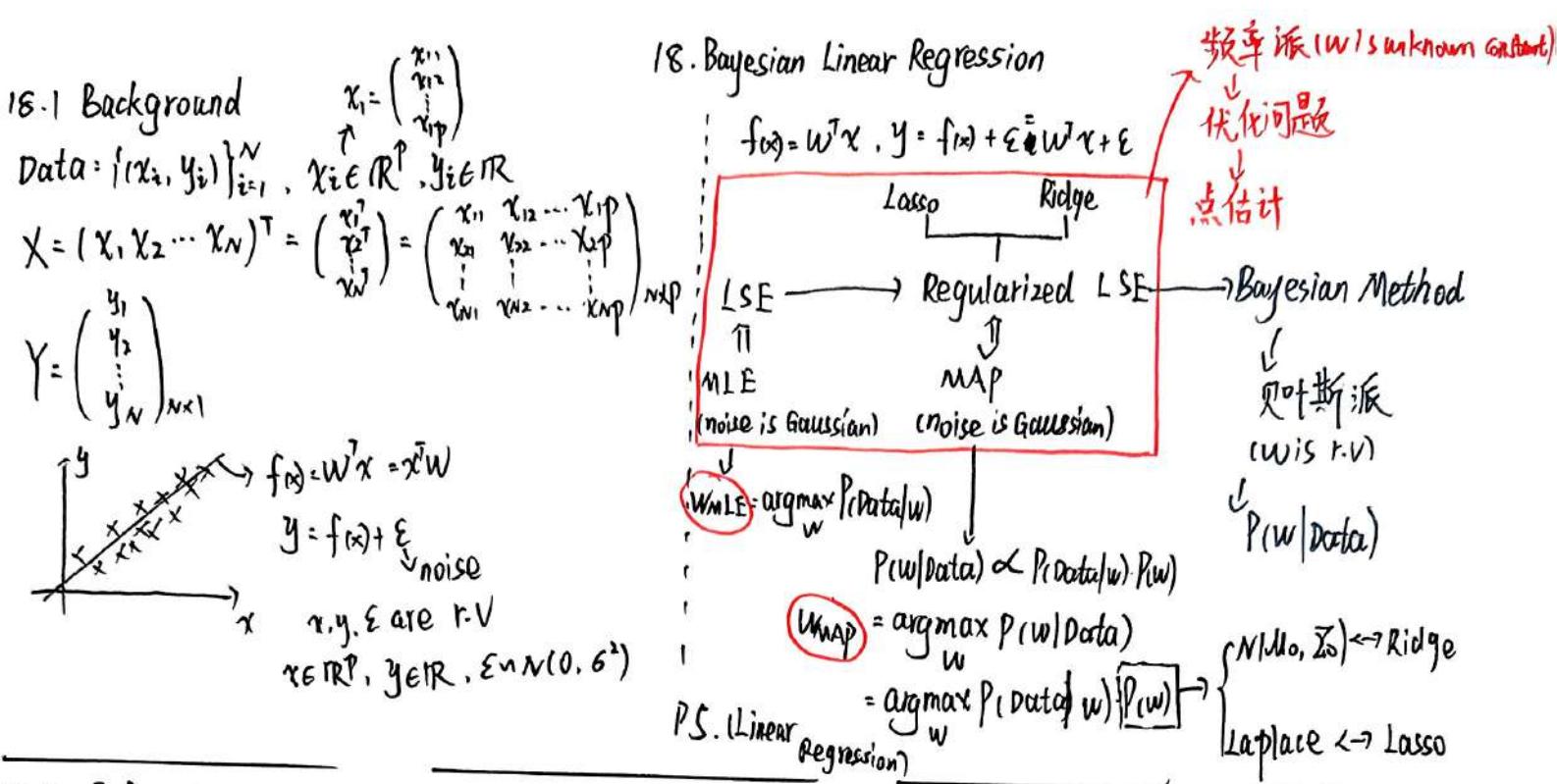
$x_i \perp x_j |_{-\{x_i, x_j\}}$

由此有 3 种情况可用
 $x_i \perp x_j |_{-\{x_i, x_j\}}$

$$\begin{cases} \Sigma^{-1} = \left(\begin{array}{cc} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{array} \right)^{-1} & \text{分块矩阵的求逆} \\ (A - BCD)^{-1} \rightarrow \text{Woodbury formula} \end{cases}$$

Schur complementary

P 3 1.4



由一个 $P(x) = N(\mu, \Sigma)$, 变量为 x
指数部分: $\exp\left\{-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right\}$

这是讨论 $P(x) = N(\mu, \Sigma)$ 的一般情况

$$= \exp\left\{-\frac{1}{2}(x^T \Sigma^{-1} - \mu^T \Sigma^{-1})(x - \mu)\right\}$$

$$\star = \exp\left\{-\frac{1}{2}\left[\underbrace{x^T \Sigma^{-1} x}_{\text{二次项}} - 2\mu^T \Sigma^{-1} x + \mu^T \Sigma^{-1} \mu\right]\right\}, \text{一次项系数}-2与外部}-\frac{1}{2}\text{可约掉: } \mu^T \Sigma^{-1} x$$

此题:

$$P(w|\text{Data}) \text{ 的二次项: } -\frac{1}{2\sigma^2} w^T X^T X w - \frac{1}{2} w^T \Sigma_p^{-1} w = -\frac{1}{2} [w^T (\underbrace{\sigma^{-2} X^T X + \Sigma_p^{-1}}_{\Sigma_w^{-1}}) w], \text{变量为 } w$$

$$P(w|\text{Data}) \text{ 的一次项: } -\frac{1}{2\sigma^2} \cdot (-2 Y^T \times w) = \underbrace{\sigma^{-2} Y^T \times w}_{\mu_w^T A} \rightarrow \text{precision matrix}$$

$$\mu_w^T A = \mu_w^T \Sigma_w^{-1} = \sigma^{-2} Y^T X$$

$$A^T = A \text{ (precision matrix)}$$

$$\therefore \mu_w^T A = \sigma^{-2} Y^T X$$

$$A \mu_w = \sigma^{-2} X^T Y$$

$$\mu_w = \sigma^{-2} A^{-1} X^T Y$$

故 $P(w|\text{Data}) \sim N(\sigma^{-2} A^{-1} X^T Y, A^{-1})$ 且 $A = \sigma^{-2} X^T X + \Sigma_p^{-1}$

18.3 Prediction.

Given χ^* , get y^* .

$$\text{Model: } \begin{cases} f(x) = w^T x = x^T w \\ y = f(x) + \varepsilon \quad \varepsilon \sim N(0, \sigma^2) \end{cases}$$

① $f(\chi^*)$: noise-free

$$\begin{aligned} f(x) &= x^T w \\ f(\chi^*) &= \chi^{*T} w \xrightarrow{\text{已知 Data } w, \text{做 prediction}} P(w|\text{Data}) = N(\mu_w, \Sigma_w) \end{aligned}$$

$$w \sim N(\mu_w, \Sigma_w)$$

$$\chi^{*T} w \sim N(\chi^{*T} \mu_w, \chi^{*T} \Sigma_w \chi^*)$$

$$\therefore P(f(\chi^*)|\text{Data}, \chi^*) = N(\chi^{*T} \mu_w, \chi^{*T} \Sigma_w \chi^*)$$

② $y^*, y^* = f(\chi^*) + \varepsilon$, noise

$$\therefore P(y^*|\text{Data}, \chi^*) = N(\chi^{*T} \mu_w, \chi^{*T} \Sigma_w \chi^* + \sigma^2)$$

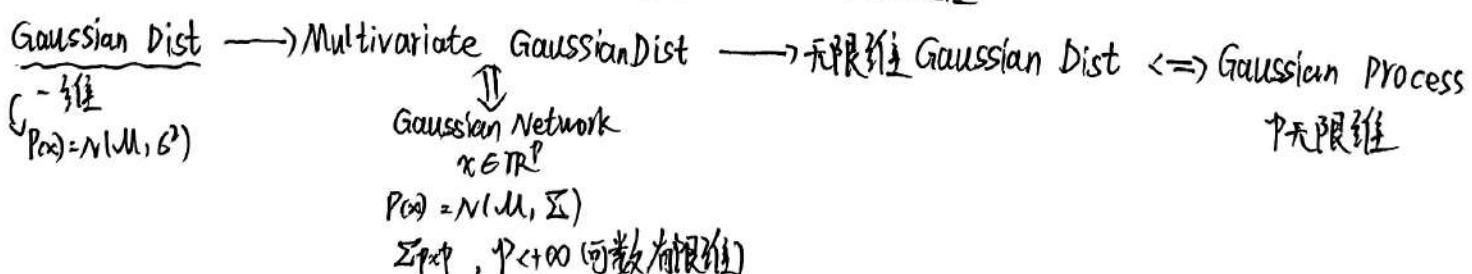
$$\begin{cases} \mu_w \\ \Sigma_w \end{cases} \xrightarrow{\text{参见上节}} \begin{cases} \mu_w = \sigma^{-2} A^{-1} X^T Y \\ \Sigma_w = A^{-1}, A = \sigma^{-2} X^T X + \Sigma_p^{-1} \end{cases}$$

总结: 先做 Inference 求出 Σ_w 与 μ_w , 再利用它们做 Prediction

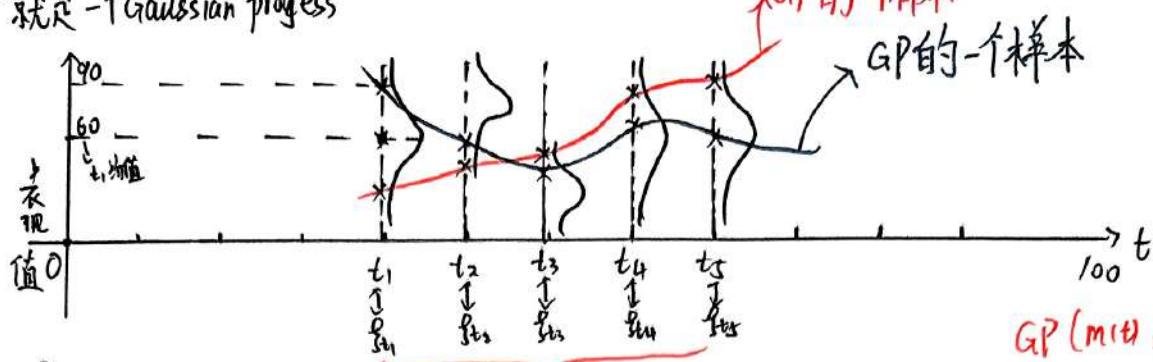
19.1 Background

19. Gaussian process

distribution \rightarrow 随机过程



$\{f_t\}_{t \in T}$, T 连续域, if $\forall n \in N^+, t_1, t_2, \dots, t_n \in T$. s.t. $\{f_{t_1}, f_{t_2}, \dots, f_{t_n}\} \triangleq f_{t_1-t_n} \sim N(\mu_{t_1-t_n}, \Sigma_{t_1-t_n})$
那么 $\{f_t\}_{t \in T}$ 就是一个 Gaussian process



人的一生: $[0, 100]$

$$t \in [0, 100], f_t \sim N(\mu_t, \sigma_t^2)$$

这代表在 t 时刻他的表现 $(0-100)$

$t=0$, 他的一生基本已定, $t>0$, 每一时刻 μ_t, σ_t^2 是已经确定

$$GP(m(t), k(s, t))$$

$$m(t) = E[f_t]$$

$$k(s, t) = E[(f_s - m(s))(f_t - m(t))^T]$$

$m(t) \rightarrow$ mean function

$k(s, t) \rightarrow$ covariance function

kernel function

19.2 weight space perspective 非线性 Bayesian Regression / Gauss process Regression

Recall Bayesian Linear Regression

$$\text{Model: } \begin{cases} f(x) = w^T x = x^T w \\ y = f(x) + \epsilon \quad [\epsilon \sim N(0, \sigma^2)] \end{cases}$$

$$\textcircled{1} P(w | \text{Data}) = N(w | \mu_w, \Sigma_w)$$

$$\begin{cases} \mu_w = G^{-2} A^{-1} x^T Y \\ \Sigma_w = A^{-1} \quad (A = G^{-2} x^T x + \Sigma_p^{-1}) \end{cases}$$

接下来介绍 kernel Trick:

若 $y = f(x) + \epsilon$, $f(x)$ is not linear function

对于 Non-linear 问题: \textcircled{1} Non-linear Transformation

\textcircled{2} Bayesian LR

\textcircled{2} Given x^* , prediction

1) noise-free

$$P(f(x^*) | \text{Data}, x^*) = N(x^{*T} \mu_w, x^{*T} \Sigma_w x^*)$$

2) noise

$$P(y^* | \text{Data}, x^*) = N(x^{*T} \mu_w + \sigma^2, x^{*T} \Sigma_w x^* + \sigma^2)$$

Noise-free:

$$x = (x_1, x_2, \dots, x_N)^T, Y = (y_1, y_2, \dots, y_N)^T$$

$$f(x^*) | x, Y, \mu^* \sim N(x^{*T} (G^{-2} A^{-1} x^T Y), x^{*T} A^{-1} x^*)$$

$$\text{If } \phi: x \mapsto z, x \in \mathbb{R}^p, z = \phi(x) \in \mathbb{R}^q, q > p, \quad P(w) \sim N(0, \Sigma_p), \phi(w) = w^T \sim N(0, \Sigma_q), \text{ noise-free}$$

$$\text{Define: } \Phi = \phi(x) = (\phi(x_1) \ \phi(x_2) \ \dots \ \phi(x_N))^T \underset{N \times q}{\underset{\uparrow}{w^T}} \underset{q \times q}{\underset{\uparrow}{\Sigma_q}}$$

$$\text{Then: Model: } f(x) = \phi(x)^T w^T, w \text{ 变为 } q \text{ 维}$$

$$\text{故有: } f(x^*) | x, Y, \mu^* \sim N(G^{-2} \phi(x^*)^T A^{-1} \Phi^T Y, \phi(x^*)^T A^{-1} \phi(x^*))$$

$$\text{How to compute } A^{-1} ? \rightarrow \text{Woodbury formula: } (A + UCV)^{-1} = A^{-1} - A^{-1} V (C^{-1} + V A^{-1} V^T)^{-1} V A^{-1}$$

后面 Proof this is kernel function, 用到 kernel trick

PS: 可以 kernel Trick, 避免

↓ 求出 $\phi(x)$, 直接定义 kernel function = $k(x, x')$

\textcircled{2} Kernel \rightarrow kernel BLR

x, w, ϵ 均升到 q 维

$$P(w) \sim N(0, \Sigma_q)$$

$$\epsilon \sim N(0, \sigma^2), \text{ 此处 } \epsilon \text{ 为 } q \text{ 维}$$

$q \times q$

\uparrow

但此 ϵ 为 $q \times q$
之前为 $p \times p$

$q \times q$

\uparrow

此处 ϵ 为 q 维

$q \times q$

在不计算 A^{-1} 的情况下，求均值

$$A = G^{-2} \bar{\Sigma} + \Sigma_q^{-1}$$

$$A \Sigma_q = G^{-2} \bar{\Sigma} \bar{\Sigma} \Sigma_q + I$$

$$A \Sigma_q \bar{\Sigma}^T = G^{-2} \bar{\Sigma}^T \bar{\Sigma} \Sigma_q \bar{\Sigma}^T + \bar{\Sigma}^T$$

$$= G^{-2} \bar{\Sigma}^T (k + \sigma^2 I), k = \bar{\Sigma} \Sigma_q \bar{\Sigma}^T$$

$$\Sigma_q \bar{\Sigma}^T = G^{-2} A^{-1} \bar{\Sigma}^T (k + \sigma^2 I)$$

$$G^{-2} A^{-1} \bar{\Sigma}^T = \Sigma_q \bar{\Sigma}^T (k + \sigma^2 I)^{-1}$$

$$\underbrace{G^{-2} \phi(x^*)^T A^{-1} \bar{\Sigma}^T}_f Y = \phi(x^*) \Sigma_q \bar{\Sigma}^T (k + \sigma^2 I)^{-1} Y \rightarrow \text{通过 } A^{-1}$$

$f(x^*) | X, Y, x^*$'s expectation

均值

Likewise: $f(x^*) | X, Y, x^*$'s covariance: $\phi(x^*)^T \Sigma_q \phi(x^*) - \phi(x^*)^T \Sigma_q \bar{\Sigma}^T (k + \sigma^2 I)^{-1} \bar{\Sigma} \Sigma_q \phi(x^*)$

$\therefore f(x^*) | X, Y, x^* \sim N\left(\underbrace{\phi(x^*)^T \Sigma_q \bar{\Sigma}^T (k + \sigma^2 I)^{-1} Y}_{A}, \underbrace{\phi(x^*)^T \Sigma_q \phi(x^*) - \phi(x^*)^T \Sigma_q \bar{\Sigma}^T (k + \sigma^2 I)^{-1} \bar{\Sigma} \Sigma_q \phi(x^*)}\right) \star$

$k = \bar{\Sigma} \Sigma_q \bar{\Sigma}^T, \phi(x^*)^T \Sigma_q \bar{\Sigma}^T, \phi(x^*)^T \Sigma_q \phi(x^*), \phi(x^*)^T \Sigma_q \bar{\Sigma}^T, \bar{\Sigma} \Sigma_q \phi(x^*)$ 形式类似

且 $\bar{\Sigma} = \phi(x) = (\phi(x_1), \phi(x_2) \dots \phi(x_n))^T$ 展开后相乘仍为 $\phi^T(x) \Sigma_q \phi(x)$ 这种形式

$X \mapsto P$

判断: $k(x, x') = \phi^T(x) \Sigma_q \phi(x')$ $\xleftrightarrow{?}$ kernel function

参见 P22. Definition

$\because \Sigma_q$: positive definite, 并且对称, $\Sigma_q = (\Sigma_q^{1/2})^2$

$$\therefore k(x, x') = \phi^T(x) \Sigma_q^{1/2} \Sigma_q^{1/2} \phi(x') = (\Sigma_q^{1/2} \phi(x))^T \cdot \Sigma_q^{1/2} \phi(x') = \langle \psi(x), \psi(x') \rangle$$

$$\psi(x) = \Sigma_q^{1/2} \phi(x)$$

$\therefore k(x, x') = \phi^T(x) \Sigma_q \phi(x')$ is kernel function

即说明不用找到具体的 $\psi(x)$, 运用 kernel Trick 定义 kernel function (P22)

Gaussian process Regression $\hat{=}$ kernel trick \leftrightarrow (Non-linear Transformation, inner product) + Bayesian Linear Regression

Two views
① weight-space view
② function-space view
 $f(x) = \phi^T(x) \psi$ weight
 $y = f(x) + \epsilon$
equal result

$f(x)$ is r.v., $f(x) \sim GP(m(x), k(x, x'))$

Gauss Process Regression is the extension of
Bayesian Linear Regression with kernel Trick.

19.3 from weight space to function space

Recall Gaussian Process:

$\{f_t\}_{t \in T}$, T : continuous time/space. $\forall n \in N(nz)$, $t_1, t_2, \dots, t_n \rightarrow \underbrace{\beta_1, \beta_2, \dots, \beta_n}_{r.v}$. 全 $\beta_{1:n} = (\beta_1, \beta_2, \dots, \beta_n)^T$,
 If $\beta_{1:n} \sim N(\mu_{1:n}, \Sigma_{1:n})$, Then $\{f_t\}_{t \in T}$ is Gaussian process. $f_t \sim GP(m_t, k_{(t,s)})$

mean function covariance function
kernel function

回到 weight-space view (关注对象为 w)

$$f(x) = \phi(x)^T w$$

$$y = f(x) + \epsilon, \epsilon \sim N(0, \sigma^2)$$

Bayesian Method:

给定 prior: $w \sim N(0, \Sigma_q)$ 此处 w 已做从 $p \rightarrow q$ 升维

$$\therefore f(x) = \phi(x)^T w$$

$$\therefore E[f(x)] = E[\phi(x)^T w] = \phi(x)^T E[w] = 0$$

$$\forall x, x' \in \mathbb{R}^q$$

$$\text{cov}(f(x), f(x')) = E[(f(x) - E[f(x)])(f(x') - E[f(x')])] = E[f(x)f(x')] = E[\phi(x)^T w \cdot \underbrace{\phi(x')^T w}_{\text{升维}}] = E[\phi(x)^T w w^T \phi(x')]$$

$$= \phi(x)^T E[w \cdot w^T] \phi(x') = \phi(x)^T E[(w-0)(w-0)] \phi(x')$$

$$\therefore w \sim N(0, \Sigma_q)$$

$$\therefore \text{cov}(f(x), f(x')) = \boxed{\phi(x)^T \Sigma_q \phi(x')} = \langle \phi(x), \phi(x') \rangle = k(x, x')$$

$$\text{Kernel function: } \psi(x) = \sum_q \frac{1}{q} \phi(x)$$

启发: $f(x)$ 是否可看作一个高斯过程

$\{f(x)\}_{x \in \mathbb{R}^q}$ 一族 就是说, 这里推出来的 $E[f(x)]$ 与 $\text{cov}(f(x), f(x'))$ 与 P51 对 Gaussian Process 的定义有相同之处, 同为 kernel function

GPR:

① weight-space view: 关注的是 w ② function-space view: 关注的是 $f(x)$

$$x^* \rightarrow y^*$$

$$P(y^* | \text{Data}, x^*) = \int_w P(y^* | w, x^*) \cdot P(w) dw \quad P(y^* | \text{Data}, x^*) = \int_f P(y^* | f, x^*) \cdot P(f) df$$

\uparrow process
Gaussian function:
 $f(x) \sim GP(m(x), k(x, x'))$
mean function covariance function

① $f(x)$ is function

② $f(x)$ is \mathbb{R}^q Gaussian Dist

$t \rightarrow \beta_t$, $\{\beta_t\}_{t \in T} \sim GP$

$x \rightarrow f(x)$, $\{f(x)\}_{x \in \mathbb{R}^q} \sim GP$

$T = \mathbb{R}^q$ 紧密的

$t \leftrightarrow x$

$\beta_t \leftrightarrow f(x)$

19.4 function space perspective

$\{f(x)\}_{x \in \mathbb{R}^q} \sim GP(m(x), k(x, x'))$ 已知的 GP, $m(x), k(x, x')$ 均已知

Regression:

Data: $\{(x_i, y_i)\}_{i=1}^N$, $y = f(x) + \epsilon, \epsilon \sim N(0, \sigma^2)$

$$X = (x_1, x_2, \dots, x_N)^T \in \mathbb{R}^{q \times N}$$

$$Y = (y_1, y_2, \dots, y_N)^T \in \mathbb{R}^{N \times 1}$$

$$f(x) \sim N(\mu(x), k(x, x))$$

$$Y = f(x) + \epsilon \sim N(\mu(x), k(x, x) + \sigma^2 I)$$

$f(x)$ is normal Dist

Prediction: Given $x^* = (x_1^*, x_2^*, \dots, x_m^*)$

$$Y^* = f(x^*) + \epsilon$$

$$(Y^*, f(x^*)) \sim N\left(\begin{pmatrix} \mu(x^*) \\ \mu(x^*) \end{pmatrix}, \begin{pmatrix} k(x^*, x^*) & k(x^*, x^*) \\ k(x^*, x^*) & k(x^*, x^*) \end{pmatrix}\right)$$

分四个块

$$\begin{aligned} \text{cov}(Y, Y) &= \text{cov}(Y, f(x^*)) \\ \text{cov}(f(x^*), Y) &= \text{cov}(f(x^*), f(x^*)) \\ \text{cov}(Y, f(x^*)) &= \text{cov}(f(x^*), f(x^*)) \\ &= k(x, x^*) \end{aligned}$$

已知联合高斯分布，求条件概率

$$(Y, f(x^*)) \sim N\left(\begin{pmatrix} \mu_{(Y)} \\ \mu_{f(x^*)} \end{pmatrix}, \begin{pmatrix} k(x, x) + \sigma^2 I & k(x, x^*) \\ k(x^*, x) & k(x^*, x^*) \end{pmatrix}\right)$$

求 $P(f(x^*) | Y, x, x^*) \rightarrow$ 条件概率 $\sim N(\mu^*, \Sigma^*)$

舍在 $f(x^*)$ 中

公式: $X \sim N(\mu, \Sigma)$
 $X = \begin{pmatrix} x_a \\ x_b \end{pmatrix}, \mu = \begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix}; \Sigma = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}$
 $x_b | x_a \sim N(\mu_{b|a}, \Sigma_{b|a})$
 $\mu_{b|a} = \Sigma_{ba} \Sigma_{aa}^{-1} (\mu_a - \mu_b) + \mu_b$
 $\Sigma_{b|a} = \Sigma_{bb} - \Sigma_{ba} \Sigma_{aa}^{-1} \Sigma_{ab}$

$$\mu^* = k(x^*, x) \cdot [k(x, x) + \sigma^2 I]^{-1} (Y - \mu_{(x)}) + \mu_{(x^*)}$$

$$\Sigma^* = k(x^*, x^*) - k(x^*, x) [k(x, x) + \sigma^2 I]^{-1} \cdot k(x, x^*)$$

① noise free:

$$P(f(x^*) | Y, x, x^*) = N(\mu^*, \Sigma^*) \quad P(f(x^*) | Y, x, x^*) = N(\mu^*, \Sigma^* + \sigma^2 I)$$

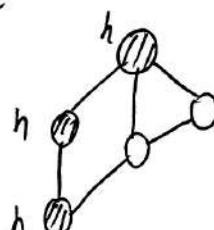
function space & weight space 一样, 但更加简单
一开始就认为 $f(x) \sim GP(m(x), k(x, x))$ \rightarrow 求出 μ^*, Σ^* 才证明 $\sim GP(m(x), k(x, x))$

20. Restricted Boltzmann Machine

20.1 Background

Boltzmann Machine: Markov Random Field with hidden nodes

Nodes \rightarrow V, h
 observed variable: v
 hidden variable: h



因子分解: (Hammersley-Clifford Theorem) 基于最大团 (C 最大团, $\psi_i(x_{ci})$: 势函数 potential function, Z : 带积因子, partition function)

$$P(x) = \frac{1}{Z} \prod_{i=1}^K \psi_i(x_{ci}) \quad \text{s.t. } \psi_i \geq 0, Z = \sum_{i=1}^K \prod_{j=1}^K \psi_j(x_{cj}) = \sum_{x_1} \sum_{x_2} \dots \sum_{x_p} \prod_{i=1}^K \psi_i(x_{ci})$$

$\star \psi_i(x_{ci}) = \exp\{-E(x_{ci})\}$ E : energy function

$$P(x) = \frac{1}{Z} \prod_{i=1}^K \psi_i(x_{ci}) = \underbrace{\frac{1}{Z} \exp\left\{-\sum_{i=1}^K E(x_{ci})\right\}}_{\text{指数型分布族}} \Rightarrow \text{简化 } P(x) = \underbrace{\frac{1}{Z} \exp\{-E(x)\}}_{\text{Boltzmann Dist (Gibbs Dist)}} \rightarrow \text{pdf}$$

history: Boltzmann Dist. 热力学物理: 一个物理系统 \rightarrow



物理解释: 粒子: 原子/分子
 每个粒子的 V 受到其它
 粒子的影响 (碰撞)
 和 E (energy function) 与
 粒子本身相关.

E 和 $P(\text{state})$ 成反比
 能量越大, 越能够挣脱束缚, 不稳定, 越可能
 从当前状态变为其它状态, 会向稳态发展

系统状态: state 表达 energy
 $P(\text{state}) \propto \exp\left\{-\frac{E}{kT}\right\}$ \downarrow temperature
 Boltzmann constant

假设一共有 M 个状态

$$\frac{S}{P_1 P_2 \dots P_i \dots P_M}$$

20.2 Representation

$\{X = \{V_1, V_2, \dots, V_p\} = \{h, v\}\} \rightarrow \text{Boltzmann Machine (latent与observed分开)}$

$h = \{h_1, h_2, \dots, h_m\}; V = \{V_1, V_2, \dots, V_n\} m+n=p$

Boltzmann Machine: 问题: Inference
精确 \rightarrow untractable
近似 \rightarrow 计算量太大
简化

Restricted Boltzmann Machine: (h, v) 之间有连接, h, v 内部无连接

$$\text{P}(x) = \frac{1}{Z} \exp\{-E(x)\} \Leftrightarrow \text{P}(v, h) = \frac{1}{Z} \exp\{-E(v, h)\}$$

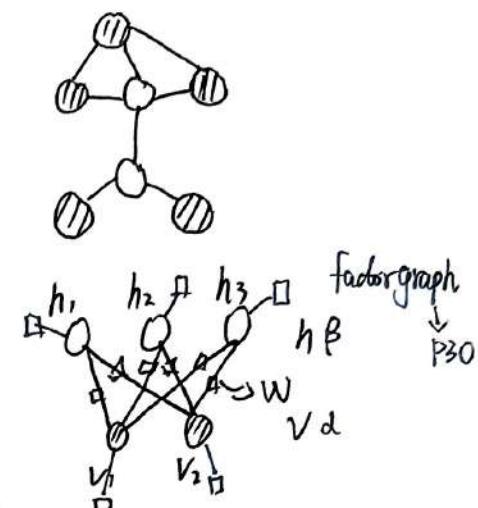
$$E(v, h) = -(\vec{h}^T W V + \vec{d}^T V + \vec{B}^T h)$$

$$\text{P}(v, h) = \frac{1}{Z} \exp\{\vec{h}^T W V + \vec{d}^T V + \vec{B}^T h\} = \frac{1}{Z} \underbrace{\exp\{\vec{h}^T W V\}}_{\text{factor graph view}} \cdot \underbrace{\exp\{\vec{d}^T V\}}_{\text{edge}} \cdot \underbrace{\exp\{\vec{B}^T h\}}_{\text{node}}$$

换为 Vector:

$$x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{pmatrix} = \begin{pmatrix} h \\ v \end{pmatrix}; h = \begin{pmatrix} h_1 \\ h_2 \\ \vdots \\ h_m \end{pmatrix}; V = \begin{pmatrix} V_1 \\ V_2 \\ \vdots \\ V_n \end{pmatrix} p=m+n$$

$$\therefore \text{RBM's pdf} \rightarrow \text{P}(x) = \text{P}(v, h) = \frac{1}{Z} \exp\{\vec{h}^T W V\} \cdot \exp\{\vec{d}^T V\} \cdot \exp\{\vec{B}^T h\} = \frac{1}{Z} \prod_{i=1}^m \prod_{j=1}^n \underbrace{\exp\{h_i w_{ij} v_j\}}_{\text{edge}} \prod_{j=1}^n \underbrace{\exp\{d_j \cdot v_j\}}_{\text{node } v} \cdot \prod_{i=1}^m \underbrace{\exp\{B_i \cdot h_i\}}_{\text{node } h}$$



$$\exp\{\vec{h}^T W V\} = \sum_{i=1}^m \sum_{j=1}^n \underbrace{h_i w_{ij} v_j}_{\text{edge}} \quad \underbrace{\exp\{h_i w_{ij} v_j\}}_{\text{node } v}$$

$$\prod_{i=1}^m \prod_{j=1}^n \exp\{h_i w_{ij} v_j\} = \prod_{j=1}^n \underbrace{\exp\{d_j \cdot v_j\}}_{\text{node } v} \cdot \prod_{i=1}^m \underbrace{\exp\{B_i \cdot h_i\}}_{\text{node } h}$$

20.3 Representation - review

Naive Bayes \leftarrow NB \rightarrow 独立贝叶斯假设

Gaussian Mixture Model \leftarrow GMM \rightarrow 引入隐变量

$y \rightarrow$ seq (线性连接) \rightarrow HMM

State Space Model \leftarrow SMM \rightarrow Kalman Filter

Logistic Regression \rightarrow MEM Model

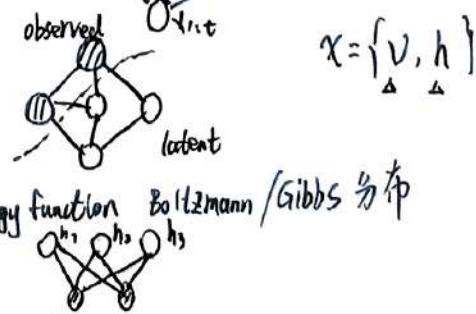
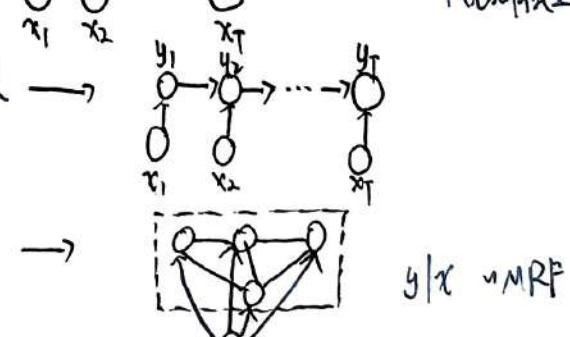
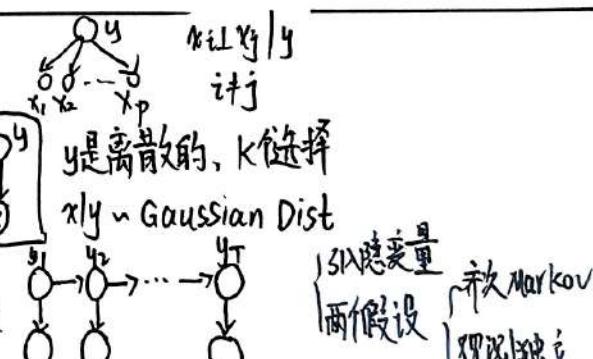
Maximum Entropy Markov Model \leftarrow MEMM \rightarrow 打破了观测独立假设

Conditional Random Field \leftarrow CRF \rightarrow 判别模型

Linear chain - CRF \leftarrow LC-CRF \rightarrow 特性相同

Boltzmann Machine \leftarrow BM \rightarrow 无向

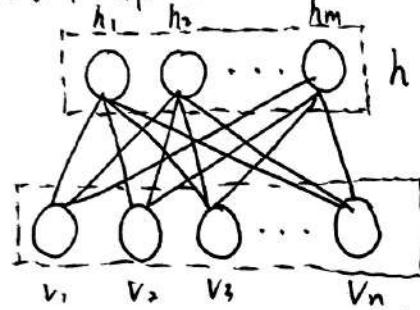
Restricted BM \leftarrow RBM \rightarrow 条件独立性



概率图模型

| | |
|--------------|------|
| 方向(有向/无向/混合) | → 边 |
| 离散/连续/混合 | → 点 |
| 条件独立性 | → 边 |
| 隐变量 | → 点 |
| 指数族分布 | → 结构 |

20.4 Inference - Posterior Distribution



local Markov. -

目的: Inference → posterior $\rightarrow P(h|v), P(v|h)$

求 $P(h|v)$

v : neighbour (h_t)

$$P(h|v) = \prod_{i=1}^m P(h_i|v) \quad h-t: \{h_i\}_{i \neq t}$$

这里假设 h_i : binary RBF h_i 只有 {0,1} 可选

$$P(h_i=1|v) = P(h_i=1|h_{-t}, v) = \frac{P(h_i=1, h_{-t}, v)}{P(h_{-t}, v)}$$

$$\textcircled{1} \quad \frac{P(h_i=1, h_{-t}, v)}{P(h_i=1, h_{-t}, v) + P(h_i=0, h_{-t}, v)} \quad P(h_i, h_{-t}, v) = P(h_{-t}, v)$$

ReLU $\rightarrow h_i \geq 0$

$$E(h, v) = -(\sum_{i=1}^m \sum_{j=1}^n h_i w_{ij} v_j + h_i \sum_{j=1}^n w_{ij} v_j + \sum_{j=1}^n d_j v_j + \sum_{i=1}^m \beta_i h_i)$$

$$\Delta_1 = \sum_{i \neq t} \beta_i h_i + \beta_0 h_0$$

$$\Delta_2 = \sum_{j=1}^n w_{ij} v_j$$

$$\Delta_3 = \sum_{j=1}^n d_j v_j$$

$$\Delta_4 = h_t \cdot H_L(v)$$

$$\Delta_5 = h_t \cdot H_U(v)$$

计算 $P(h_t, h_{-t}, v)$ 的通项, 把 $h_i=1, 0$ 代入即可

$$h_t=1$$

$$E(h, v) = h_t \cdot H_L(v) + \bar{H}_L(h_{-t}, v) = -E(h, v)$$

$$\therefore \text{分子} = P(h_t=1, h_{-t}, v) = \frac{1}{Z} \exp \left\{ h_t \cdot H_L(v) + \bar{H}_L(h_{-t}, v) \right\} = \frac{1}{Z} \exp \left\{ \cdot H_L(v) + \bar{H}_L(h_{-t}, v) \right\}$$

$$\text{分子} = \frac{1}{Z} \exp \left\{ \cdot H_L(v) + \bar{H}_L(h_{-t}, v) \right\} + \frac{1}{Z} \exp \left\{ \cdot \bar{H}_L(h_{-t}, v) \right\}$$

$$\therefore P(h_t=1|v) = \frac{1}{1 + \exp \{ H_L(v) \}} = \sigma(H_L(v)) = \sigma \left[\left(\sum_{j=1}^n w_{ij} v_j + \beta_i \right) \right] \star \quad \text{Sigmoid function}$$

$$\sigma \rightarrow \text{Sigmoid. function} \xrightarrow{\text{1}} \sigma(x) = \frac{1}{1+e^{-x}}$$

$$\text{Likewise: } P(h_t=0|v) = \frac{1}{1 + \exp \{ H_L(v) \}} = \sigma(-H_L(v)) = \sigma \left[- \left(\sum_{j=1}^n w_{ij} v_j + \beta_i \right) \right] \star$$

$$\text{分子} = P(h_t=0, h_{-t}, v) = \frac{1}{Z} \exp \{ h_t \cdot H_L(v) + \bar{H}_L(h_{-t}, v) \} = \frac{1}{Z} \exp \{ \bar{H}_L(h_{-t}, v) \}$$

$$\text{分子} = P(h_t=0, h_{-t}, v) + P(h_t=1, h_{-t}, v) = \frac{1}{Z} \exp \{ \bar{H}_L(h_{-t}, v) \} + \frac{1}{Z} \{ H_L(v) + \bar{H}_L(h_{-t}, v) \}$$

$$P(v|h) \text{ 与上述过程相同}$$

$$\text{可以求闭解 } P(h_t|v)$$

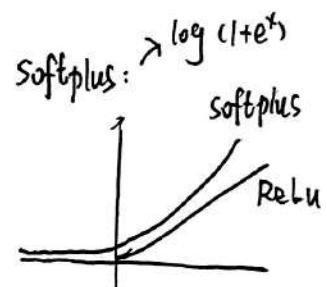
$$\text{RBM - 神经网络 sigmoid} \rightarrow \sigma(x) = \frac{1}{1+e^{-x}}$$

20.5. Inference - Marginal Distribution

目的: Inference \rightarrow marginal $\rightarrow P(v)$

$$\begin{aligned}
 P(v) &= \sum_h P(h, v) = \sum_h \frac{1}{Z} \exp\{-E(h, v)\} = \sum_h \frac{1}{Z} \exp\{\alpha^T v + h^T w v + \beta^T h\} \\
 &= \sum_{h_1} \cdots \sum_{h_m} \frac{1}{Z} \exp\left\{\frac{\alpha^T v + \underbrace{h_1^T w v}_{\Delta} + \underbrace{\beta^T h}_{\Delta}}{\Delta}\right\} = \frac{1}{Z} \exp\{\alpha^T v\} \cdot \sum_{h_1} \cdots \sum_{h_m} \exp\left\{\underbrace{h_1^T w v + \beta^T h}_{\sum_{i=1}^m (h_i w_i v + \beta_i h_i)}\right\} \\
 &= \frac{1}{Z} \exp\{\alpha^T v\} \cdot \sum_{h_1} \cdots \sum_{h_m} \exp\left\{\sum_{i=1}^m (h_i w_i v + \beta_i h_i)\right\} \\
 &= \frac{1}{Z} \exp\{\alpha^T v\} \cdot \sum_{h_1} \cdots \sum_{h_m} \exp\{h_1 w_1 v + \beta_1 h_1\} \cdot \exp\{h_2 w_2 v + \beta_2 h_2\} \cdots \cdot \exp\{h_m w_m v + \beta_m h_m\} \\
 &= \frac{1}{Z} \exp\{\alpha^T v\} \sum_{h_1} \exp\{h_1 w_1 v + \beta_1 h_1\} \cdots \sum_{h_m} \exp\{h_m w_m v + \beta_m h_m\} \quad \because h_i \in \{0, 1\} \\
 &= \frac{1}{Z} \exp\{\alpha^T v\} (1 + \exp\{w_1 v + \beta_1\}) \cdots (\exp\{w_m v + \beta_m\}) \\
 &= \frac{1}{Z} \exp\{\alpha^T v\} \cdot \exp\{\log(1 + \exp\{w_1 v + \beta_1\}) \cdots (1 + \exp\{w_m v + \beta_m\})\} \\
 &= \frac{1}{Z} \exp\{\alpha^T v + \sum_{i=1}^m [\log(1 + \exp\{w_i v + \beta_i\})]\}
 \end{aligned}$$

$$\begin{aligned}
 h &= \begin{pmatrix} h_1 \\ h_2 \\ \vdots \\ h_m \end{pmatrix} \quad h_i \in \{0, 1\} \\
 w &= [w_{ij}]_{m \times n} = \begin{pmatrix} -w_1 & \cdots & -w_n \\ \vdots & \ddots & \vdots \\ -w_m & \cdots & -w_1 \end{pmatrix} \\
 &\quad (h_1, h_2, \dots, h_m) \begin{pmatrix} -w_1 \\ -w_2 \\ \vdots \\ -w_m \end{pmatrix} \\
 &\quad (B_1, \dots, B_m) \begin{pmatrix} h_1 \\ h_2 \\ \vdots \\ h_m \end{pmatrix}
 \end{aligned}$$

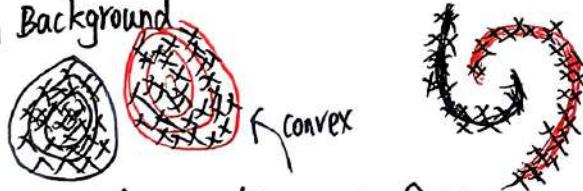


$$\therefore P(v) = \frac{1}{Z} \exp\{\alpha^T v + \sum_{i=1}^m [\text{softplus}(w_i v + \beta_i)]\}$$

Learning 放在后面. P68

21. Spectral Clustering

21.1 Background



GMM/ k-means 不好做

kernel + k-means 会映射到高维

| compactness: k-means, GMM

| connectivity: spectral clustering

21.2 Model Introduction

Graph-based (带权重的无向图)

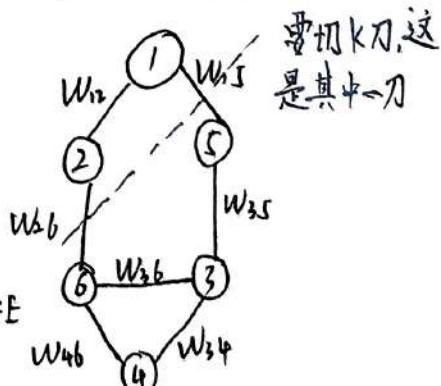
$$X = (x_1, x_2, \dots, x_N)^T = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{pmatrix}_{N \times p} \quad \left\{ \begin{array}{l} V = \bigcup_{k=1}^K A_k \\ A_i \cap A_j = \emptyset, \forall i, j \in \{1, 2, \dots, K\} \end{array} \right.$$

$$G = \{V, E\}; V = \{1, 2, \dots, N\} \Leftrightarrow X; E = W = [w_{ij}]_{N \times N}$$

w : similarity matrix (affinity matrix), 其 $w_{ij} = \begin{cases} K(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right), & (i, j) \in E \\ 0, & \text{otherwise} \Rightarrow (i, j) \notin E \end{cases}$

定义: $ACV, BCV, A \cap B = \emptyset$

$$W(A, B) = \sum_{i \in A} \sum_{j \in B} w_{ij}$$



$$\begin{aligned}
 \text{cut}(v) &= \text{cut}(A_1, \dots, A_K) \\
 &= \sum_{i=1}^K [W(A_i, v) - W(A_i, \bar{A}_i)] \\
 &= \sum_{i=1}^K [W(A_i, v) - W(A_i, A_i)]
 \end{aligned}$$

目标: $\min \text{Cut}(V)$. $\text{cut}(V) = \sum_{i=1}^k \frac{w(A_i, \bar{A}_i)}{\Delta}$, $\Delta = \text{degree}(A_k)$ 有向图: 出度. 入度
 $\{A_k\}_{k=1}^K$ 作平均
 \uparrow 无向图: 度
 $\therefore \Delta = \text{degree}(A_k) = \sum_{i \in A_k} d_i$, $d_i = \sum_{j=1}^N w_{ij}$
 $\text{如 } d_1 = \sum_{j=1}^N w_{1j}, d_2 = \sum_{j=1}^N w_{2j}$

Model
 $\therefore N\text{cut} = \sum_{k=1}^K \frac{w(A_k, \bar{A}_k)}{\sum_{i \in A_k} d_i}$, $d_i = \sum_{j=1}^N w_{ij}$

目标
 $\arg \min N\text{cut}$ - 优化问题
 $\{A_k\}_{k=1}^K$ A_k 是顶点的一个组合
 $\{\hat{A}_k\}_{k=1}^K = \arg \min N\text{cut}(V)$

2.1.3 Matrix - Indicator Vector / Diagonal Matrix

Indicator vector: $y_i \in \{0, 1\}^K$ 表示 $y_i = \begin{pmatrix} 1 \leq i \leq N \\ 1 \leq j \leq K \end{pmatrix}$ 样本个数 N 类别个数 K

$\sum_{j=1}^K y_{ij} = 1 \rightarrow \begin{pmatrix} y_{11} \\ y_{12} \\ \vdots \\ y_{1N} \end{pmatrix}$ 其中只有一个 1

$\{A_k\}_{k=1}^K = Y = (y_1, y_2, \dots, y_N)^T_{N \times K} = \begin{pmatrix} y_1^T \\ y_2^T \\ \vdots \\ y_N^T \end{pmatrix}$

$y_{ij} = 1 \Leftrightarrow$ 第 i 个样本属于第 j 种类别. 目标: $\boxed{Y = \arg \min N\text{cut}(V)}$

Diagonal Matrix

$N\text{cut}(V) = \sum_{k=1}^K \frac{w(A_k, \bar{A}_k)}{\sum_{i \in A_k} d_i} = \text{tr} \left(\frac{w(A, \bar{A})}{\sum_{i \in A_1} d_i} - \frac{w(A_2, \bar{A}_2)}{\sum_{i \in A_2} d_i} - \dots - \frac{w(A_K, \bar{A}_K)}{\sum_{i \in A_K} d_i} \right) = \text{tr} \left(\begin{array}{cccc} w(A, \bar{A}) & & & 0 \\ \vdots & \ddots & & \vdots \\ 0 & \dots & \ddots & 0 \end{array} \right) = \text{tr} \left(\begin{array}{cccc} \sum_{i \in A_1} d_i & & & 0 \\ \vdots & \ddots & & \vdots \\ 0 & \dots & \ddots & \sum_{i \in A_K} d_i \end{array} \right)^{-1}$

$O_{K \times K} \Leftrightarrow O^T \cdot P$

$\therefore N\text{cut}(V) = \text{tr} O \cdot P^T$

已知 W, Y , 求 O, P

$Y^T Y = (y_1, y_2, \dots, y_N) \begin{pmatrix} y_1^T \\ y_2^T \\ \vdots \\ y_N^T \end{pmatrix} = \sum_{i=1}^N y_i \cdot y_i^T = \begin{pmatrix} N_1 & & & \\ & N_2 & & \\ & & \ddots & \\ & & & N_K \end{pmatrix}_{K \times K}$ $\downarrow A_k$ 中的值

N_k : 在 N 个样本中, 属于类 k 的样本个数, $\sum_{k=1}^K N_k = N$, $N_k = |A_k| = \sum_{i \in A_k} 1$

$$Y^T Y = \left(\sum_{i \in A_1} 1 \quad \sum_{i \in A_2} 1 \quad \dots \quad \sum_{i \in A_K} 1 \right)_{K \times K}$$

$$\therefore P = Y^T \cdot D \cdot Y, D = \begin{pmatrix} d_1 & & & \\ & d_2 & & \\ & & \ddots & \\ & & & d_N \end{pmatrix}_{N \times N} \text{ 有 } N \text{ 个样本点}$$

$$= \text{diag}(w, 1_N)$$

$\rightarrow d_i \rightarrow 1 个 N 点, 参考 d_i = \sum_{j=1}^N w_{ij}$

$$w \cdot 1_N = \begin{pmatrix} w_{11} & w_{12} & \dots & w_{1N} \\ w_{21} & w_{22} & \dots & w_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ w_{NN} & w_{NN} & \dots & w_{NN} \end{pmatrix} \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}_{N \times 1} = \begin{pmatrix} d_1 \\ d_2 \\ \vdots \\ d_N \end{pmatrix}_{N \times 1}$$

$$\therefore P = Y^T \cdot \text{diag}(w, 1_N) \cdot Y$$

Laplacian Matrix:

$$O = \begin{pmatrix} W(A_1, \bar{A}_1) & & \\ & \ddots & \\ & & W(A_k, \bar{A}_k) \end{pmatrix}_{K \times K}$$

$$W(A_k, \bar{A}_k) = \underbrace{W(A_k, V)}_{\sum d_i} - \underbrace{W(A_k, A_k)}_{\sum_{i \in A_k} \sum_{j \in A_k} W_{ij}}$$

$$\therefore O = \left(\begin{array}{c|c|c|c} \sum d_i & \cdots & 0 & \\ \hline \vdots & \ddots & \vdots & \\ \hline 0 & \cdots & \sum d_i & \end{array} \right) - \left(\begin{array}{c|c|c|c} W(A_1, A_1) & \cdots & 0 & \\ \hline \vdots & \ddots & \vdots & \\ \hline 0 & \cdots & W(A_k, A_k) & \end{array} \right)$$

$$O = Y^T \cdot D \cdot Y - A, O' = Y^T D Y - Y^T W Y = \text{tr}(O P') = \text{tr}(O' P')$$

注: $\text{Ncut}(v) = \text{tr} O \cdot P^{-1} = \text{tr}(Y^T D \cdot Y - A) \cdot P^{-1} \Leftrightarrow \text{tr}(Y^T D \cdot Y - W) \cdot P^{-1}$

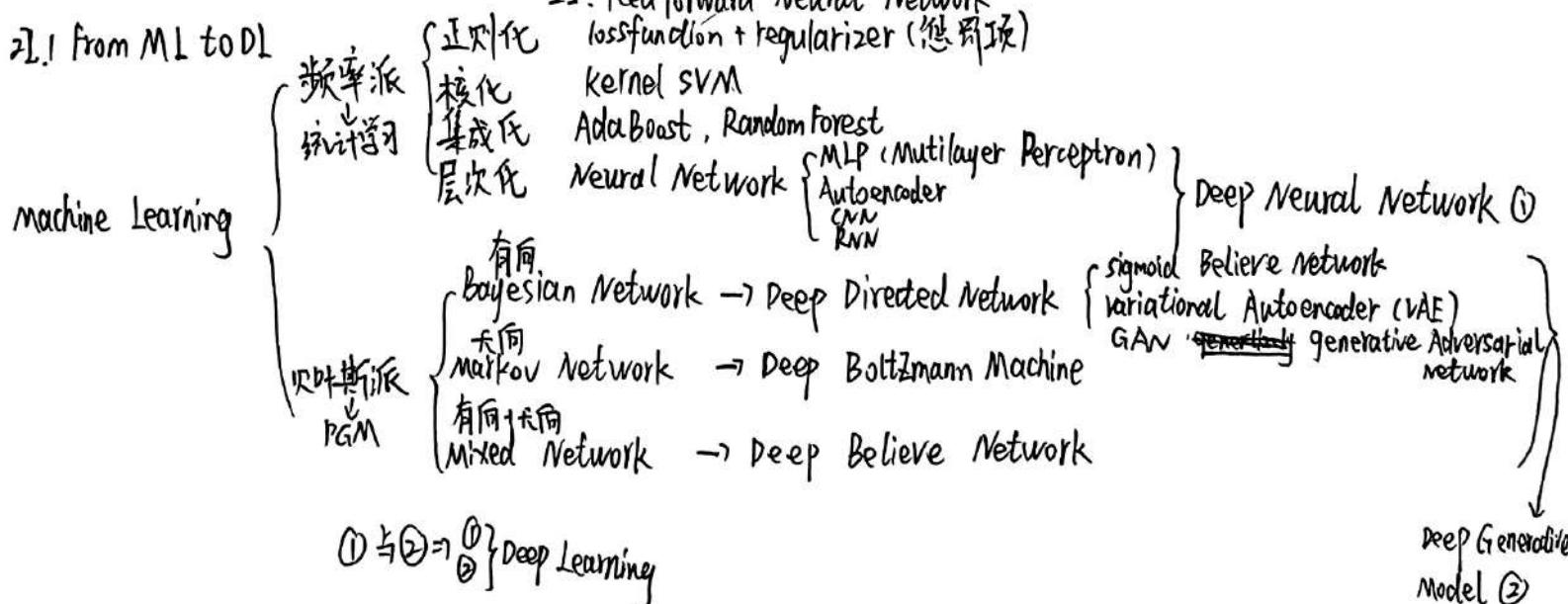
why?: 对角线, trace 只要对角线, 且 O, P, P^{-1} 为对角阵, 只看对角

$$\therefore = \text{tr}(Y^T D \cdot Y - Y^T W Y) \cdot P^{-1}$$

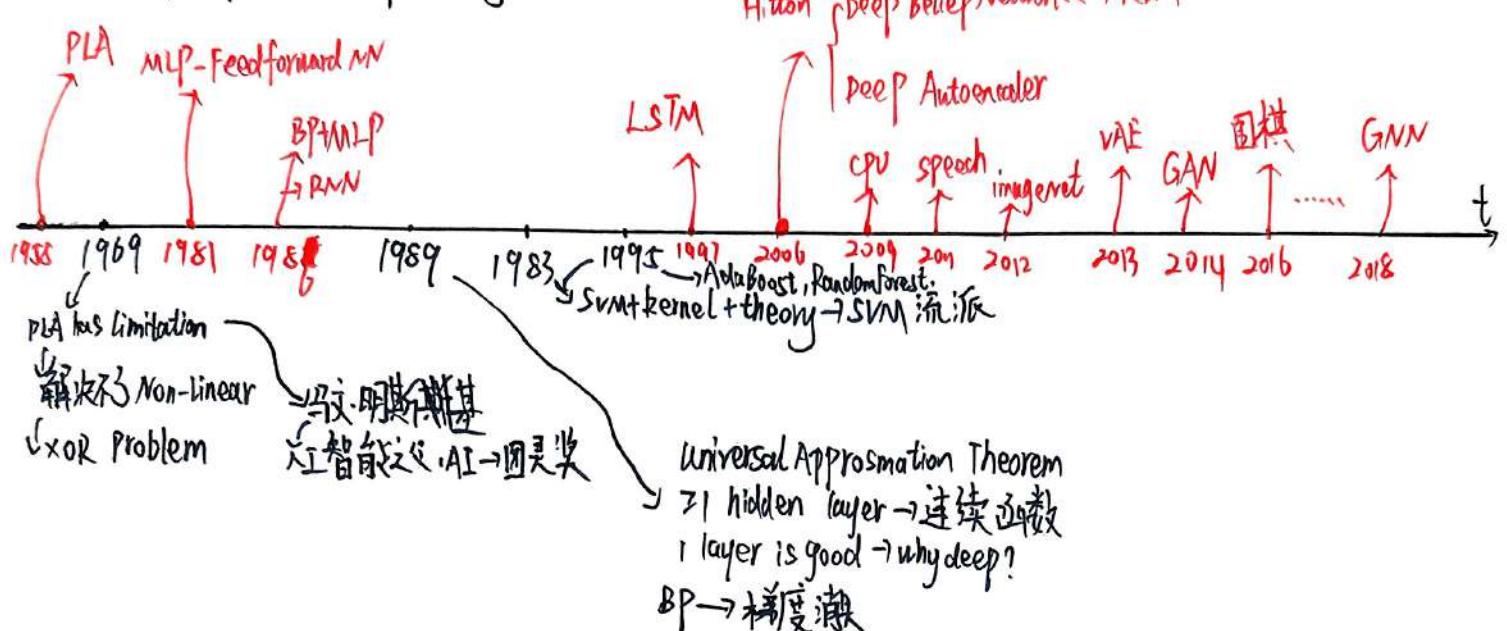
$$= \text{tr}(Y^T (D-W) Y \cdot (Y^T D \cdot Y)^{-1}), \text{且 } L = D-W, \text{为 Laplacian Matrix}$$

$$\therefore \hat{Y} = \underset{Y}{\operatorname{argmin}} \text{Ncut}(v) = \underset{Y}{\operatorname{argmin}} \text{tr} [Y^T (D-W) Y \cdot (Y^T D \cdot Y)^{-1}]$$

2]. From ML to DL



22.2 from perceptron to deep learning



22.3 non-linear transformation

solution:

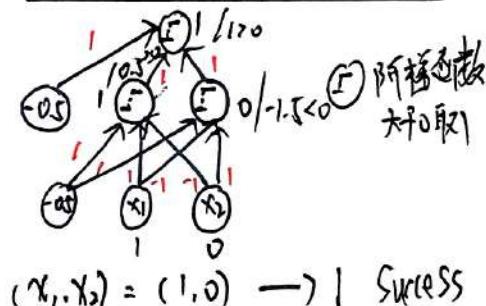
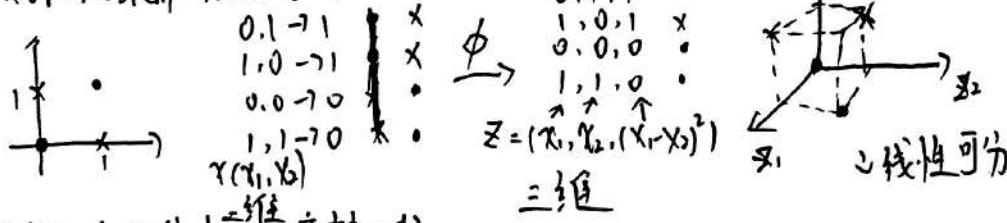
① Non-Transformation: 映射(ϕ)

$$\phi: \mathcal{X} \rightarrow \mathcal{Z}$$

input space feature space

Cover's Theorem: 高维比低维更易线性可分

XOR Problem: 相同为0, 不同为1



② Kernel Method → 映射(ϕ)

$$k(x, x') = \langle \phi(x), \phi(x') \rangle \text{ 隐藏了一个 } \phi$$

$$x, x' \in \mathcal{X}$$

③ 神经网络 (XOR prob) → 自转(ϕ)

XOR OR AND NOT

复合运算 基本运算

XOR → 相同为0, 不同为1

$$x_1 \oplus x_2 = (x_1 \wedge x_2) \vee (x_1 \wedge \neg x_2)$$

$$x_1 \oplus x_2 \rightarrow 1 = \begin{cases} 1, 0 \\ 0, 1 \end{cases} = \textcircled{1} \vee \textcircled{2} = \textcircled{1} \wedge \textcircled{2} = (x_1 \wedge x_2) \vee (x_1 \wedge \neg x_2)$$

复合运算 → 复合表达式

④ OR: $x_1 \text{ OR } x_2 \rightarrow 1$

$$1, 0 \rightarrow 1$$

$$0, 1 \rightarrow 1$$

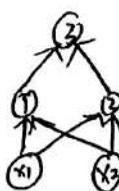
$$1, 1 \rightarrow 1$$

$$0 \rightarrow 1$$

$$1 \rightarrow 1$$

⑤ AND: $x_1 \text{ AND } x_2 \rightarrow 1$

$$1, 1 \rightarrow 1$$



有向无环图
神经网络

$$\begin{matrix} \textcircled{1} \wedge \\ \textcircled{2} \wedge \\ \textcircled{3} \vee \end{matrix}$$

复合运算 → 复合表达式 → 复合函数

23. Confronting Partition function

23.1 The log-Likelihood gradient

动机: Learning, evaluation

$$x \in \mathbb{R}^p, f_{(0,1)}^p$$

$$p(x|\theta) = \frac{1}{Z(\theta)} \hat{P}(x|\theta), Z(\theta) = \int \hat{P}(x|\theta) dx \quad \int \hat{P}(x|\theta) dx \neq 1, \hat{P}(x|\theta) 未归一化, P(x|\theta) 为归一化的概率分布$$

$$\text{MLE: } \theta \leftarrow \arg \max_{\theta} \sum_{i=1}^N \log \hat{P}(x_i|\theta)$$

$$\hat{\theta} = \arg \max_{\theta} p(x|\theta) = \arg \max_{\theta} \prod_{i=1}^N \hat{P}(x_i|\theta) = \arg \max_{\theta} \log \prod_{i=1}^N \hat{P}(x_i|\theta) = \arg \max_{\theta} \sum_{i=1}^N \log \hat{P}(x_i|\theta)$$

$$= \arg \max_{\theta} \sum_{i=1}^N [\log \hat{P}(x_i|\theta) - \log Z(\theta)] = \arg \max_{\theta} \sum_{i=1}^N \log \hat{P}(x_i|\theta) - \log Z(\theta) \cdot N \quad \text{代入 } P(x|\theta) = \frac{1}{Z(\theta)} \cdot \hat{P}(x|\theta)$$

$$= \arg \max_{\theta} \left[\frac{1}{N} \sum_{i=1}^N [\log \hat{P}(x_i|\theta)] \right] - [\log Z(\theta)] \quad \begin{array}{l} \text{正相} \\ \text{负相} \end{array}$$

positive phase negative phase

$$\nabla_{\theta} L(\theta) = \underbrace{\frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log \hat{P}(x_i|\theta)}_{①} - \underbrace{\nabla_{\theta} \log Z(\theta)}_{②}$$

相对后验积常数.

$$\begin{aligned} ② &= \frac{1}{Z(\theta)} \cdot \nabla_{\theta} \cdot Z(\theta) = \frac{P(x|\theta)}{\hat{P}(x|\theta)} \cdot \nabla_{\theta} \int \hat{P}(x|\theta) dx = \boxed{\frac{P(x|\theta)}{\hat{P}(x|\theta)}} \cdot \int \nabla_{\theta} \hat{P}(x|\theta) dx = \int \frac{P(x|\theta)}{\hat{P}(x|\theta)} \cdot \nabla_{\theta} \hat{P}(x|\theta) dx \\ &\because \nabla_{\theta} \log \hat{P}(x|\theta) = \frac{1}{\hat{P}(x|\theta)} \cdot \nabla_{\theta} \hat{P}(x|\theta) \quad \text{与 } \theta \text{ 相关, 不变.} \\ &\text{参见: } Z(\theta) = \int \hat{P}(x|\theta) dx, x \text{ 被积分了} \end{aligned}$$

$$\therefore ② = \int P(x|\theta) \cdot \nabla_{\theta} \log \hat{P}(x|\theta) dx = \underset{P(x|\theta)}{E} [\nabla_{\theta} \log \hat{P}(x|\theta)] \rightarrow \text{近似值 (MC)}$$

① 可以使用梯度上升法

$$\nabla_{\theta} L(\theta) = \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log \hat{P}(x_i|\theta) - \underset{P(x|\theta)}{E} [\nabla_{\theta} \log \hat{P}(x|\theta)]$$

23.2 Stochastic Maximum Likelihood

pdata distribution: P_{data} , 经验分布; $\hat{\theta} = \arg \max_{\theta} L(\theta)$

model distribution: $P_{\text{model}} \triangleq P(x|\theta)$; $L(\theta) = \frac{1}{N} \sum_{i=1}^N \log \hat{P}(x_i|\theta) - \log Z(\theta)$

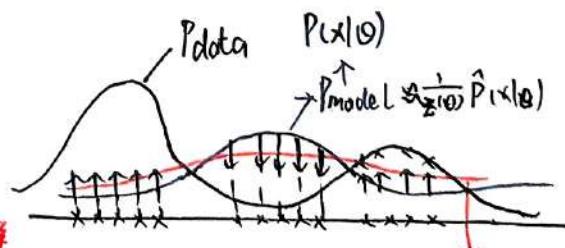
$$\therefore \nabla_{\theta} L(\theta) = \underset{\text{positive phase}}{E_{P_{\text{data}}} [\nabla_{\theta} \log \hat{P}(x_i|\theta)]} - \underset{\text{negative phase}}{E_{P_{\text{Model}}} [\nabla_{\theta} \log \hat{P}(x|\theta)]} \quad \text{Gradient Ascent based on MC MC}$$

Gradient Ascend:

$$\theta^{(t+1)} = \theta^{(t)} + \gamma \nabla_{\theta} L(\theta^{(t)}) \quad \text{相同用 } \nabla_{\theta} L(\theta) = 0, \text{ 停止 stop}$$

Gibbs Sampling from P_{model} : $i = 1, 2, \dots, m$, $\{x_i\}_{i=1}^m$, 从上一次迭代采样

$$\begin{aligned} \hat{x}_i &\sim P(x|\theta^{(t)}) \\ \vdots \\ \hat{x}_m &\sim P(x|\theta^{(t)}) \end{aligned} \quad \left\{ \text{fancy particles} \right.$$



$$Z(\theta) = \int \hat{P}(x|\theta) dx \quad \text{减小}$$

23.3 what is Contrastive Divergence

由上述两点总结：

$$x \in \mathbb{R}^d, \{x_i\}_{i=1}^N, P(x|\theta) = \frac{1}{Z(\theta)} \hat{P}(x|\theta), Z(\theta) = \int \hat{P}(x|\theta) dx$$

MLE: Given $X = \{x_i\}_{i=1}^N$, Estimate θ

$$\hat{\theta} = \operatorname{argmax}_{\theta} L(\theta) = \operatorname{argmin}_{\theta} KL(P_{\text{data}} || P_{\text{model}}) = \operatorname{argmin}_{\theta} KL(P^{(0)} || P^{(\infty)})$$

$$L(\theta) = \frac{1}{N} \sum_{i=1}^N \log \hat{P}(x_i|\theta) - \log Z(\theta)$$

$$\nabla_{\theta} L(\theta) = E_{P_{\text{data}}} [\nabla_{\theta} \log \hat{P}(x_i|\theta)] - E_{P_{\text{model}}} [\nabla_{\theta} \log \hat{P}(x|\theta)]$$

positive phase negative phase

Gradient Ascend:

$$\theta^{(t+1)} = \theta^{(t)} + \eta \nabla_{\theta} L(\theta)$$

23.4 The Name of Contrastive Divergence

$$\hat{\theta} = \operatorname{argmax}_{\theta} \frac{1}{N} \sum_{i=1}^N \log P(x_i|\theta) = \operatorname{argmax}_{\theta} E_{P_{\text{data}}} [\log P(x_i|\theta)] = \operatorname{argmax}_{\theta} \int P_{\text{data}} \log P_{\text{model}} dx$$

$$= \operatorname{argmax}_{\theta} \int P_{\text{data}} \cdot \log \frac{P_{\text{model}}}{P_{\text{data}}} dx = \operatorname{argmax}_{\theta} -KL(P_{\text{data}} || P_{\text{model}}) = \operatorname{argmin}_{\theta} KL(P_{\text{data}} || P_{\text{model}})$$

CD-K ($K=1, 2, \dots$)

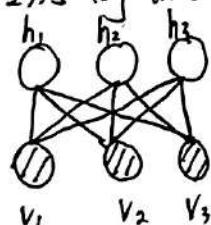
$$CD = KL(P^{(0)} || P^{(\infty)}) - KL(P^{(k)} || P^{(\infty)})$$

$$\hat{\theta} = \operatorname{argmin} [KL(P^{(0)} || P^{(\infty)}) - KL(P^{(k)} || P^{(\infty)})]$$

Contrastive Divergence

$$KL(P^{(0)} || P^{(\infty)}) - KL(P^{(k)} || P^{(\infty)}) \text{ 梯度 } \propto \sum_{i=1}^m \nabla_{\theta} \log \hat{P}(x_i|\theta^{(t)}) - \sum_{i=1}^m \nabla_{\theta} \log \hat{P}(x_i|\theta^{(t)})$$

23.5 log-likehood gradient of energy-based model / RBM Learning



$$\begin{cases} P(h, v) = \frac{1}{Z} \exp \{-E(h, v)\} \\ E(h, v) = -(\mathbf{h}^T \mathbf{W} \mathbf{v} + \mathbf{d}^T \mathbf{v} + \mathbf{b}^T \mathbf{h}) \end{cases}$$

$$\begin{aligned} \log P(v) &= \log \sum_h P(h, v) = \log \sum_h \frac{1}{Z} \exp \{-E(h, v)\} \\ &= \log \sum_h \exp \{-E(h, v)\} - \log Z \\ &= \log \sum_h \exp \{-E(h, v)\} - \log \sum_{h, v} \exp \{-E(h, v)\} \\ &= ① - ② \end{aligned}$$

log-Likelihood: training set $v \in S, |S| = N$

$$\frac{1}{N} \sum_{v \in S} \log P(v)$$

log-Likelihood gradient:

$$\frac{\partial}{\partial \theta} \frac{1}{N} \sum_{v \in S} \log P(v)$$

when $t+1$:

① Sampling for positive phase from P_{data}
 x_1, x_2, \dots, x_m themselves are training data

② Sampling for negative phase from CD Learning

$$P_{\text{model}} = P(x_i|\theta^{(t)}) \quad \hat{x}_i = x_i \text{, 指直接从训练数据中采样}$$

• Initialization mixing time $\hat{x}_1, \hat{x}_2, \dots, \hat{x}_m$ (any distribution)

0-step $\xrightarrow{\text{1-step}} \dots \xrightarrow{k\text{-step}} \dots \xrightarrow{\text{1-step}} \hat{x}_1$

Gibbs Sampling

$$\theta^{(t+1)} = \theta^{(t)} + \eta \cdot \left[\sum_{i=1}^m \nabla_{\theta} \log \hat{P}(x_i|\theta^{(t)}) - \sum_{i=1}^m \nabla_{\theta} \log \hat{P}(\hat{x}_i|\theta^{(t)}) \right]$$

CD-K: 在Markov链中在第K点停止不可达到平稳分布

在 23.3 中修改 Initialization 变为 CD-Learning

$$\begin{cases} P_{\text{data}} = P^{(0)} \\ P_{\text{model}} = P^{(\infty)} \end{cases}$$

$$\begin{aligned} \log P(v) &= \log \sum_h P(h, v) = \log \sum_h \frac{1}{Z} \exp \{-E(h, v)\} \\ &= \log \sum_h \exp \{-E(h, v)\} - \log Z \\ &= \log \sum_h \exp \{-E(h, v)\} - \log \sum_{h, v} \exp \{-E(h, v)\} \\ &= ① - ② \\ \therefore \frac{\partial}{\partial \theta} \log P(v) &= \frac{\partial}{\partial \theta} ① - \frac{\partial}{\partial \theta} ② \\ \frac{\partial}{\partial \theta} ① &= \frac{\partial}{\partial \theta} \log \sum_h \exp \{-E(h, v)\} = \frac{1}{\sum_h \exp \{-E(h, v)\}} \sum_h \exp \{-E(h, v)\} \cdot \frac{\partial}{\partial \theta} (-E(h, v)) \\ &= \frac{-1}{\sum_h \exp \{-E(h, v)\}} \sum_h \exp \{-E(h, v)\} \cdot \frac{\partial}{\partial \theta} E(h, v) \end{aligned}$$

$$\begin{aligned} \frac{\partial}{\partial \theta} \Theta &= \sum_h \cdot \frac{\frac{1}{Z} \exp\{-E(h, v)\}}{\sum_h \exp\{-E(h, v)\}} \cdot \frac{\partial}{\partial \theta} E(h, v) = (-1) \sum_h \cdot \frac{\frac{1}{Z} \exp\{-E(h, v)\}}{\sum_h \frac{1}{Z} \exp\{-E(h, v)\}} \cdot \frac{\partial E(h, v)}{\partial \theta} = - \sum_h \frac{P(h, v)}{\sum_h P(h, v)} \cdot \frac{\partial}{\partial \theta} E(h, v) \\ &= - \sum_h \frac{P(h, v)}{P(v)} \cdot \frac{\partial}{\partial \theta} E(h, v) = - \sum_h P(h|v) \cdot \frac{\partial}{\partial \theta} E(h, v) \\ \frac{\partial}{\partial \theta} \Theta &= \frac{\partial}{\partial \theta} \log \sum_{h, v} \exp\{-E(h, v)\} = \frac{1}{\sum_h \exp\{-E(h, v)\}} \sum_{h, v} \exp\{-E(h, v)\} \cdot (-1) \cdot \frac{\partial E(h, v)}{\partial \theta} \\ &= - \sum_{h, v} \underbrace{\frac{\exp\{-E(h, v)\}}{\sum_{h, v} \exp\{-E(h, v)\}}}_{Z} \cdot \frac{\partial E(h, v)}{\partial \theta} = - \sum_{h, v} P(h, v) \cdot \frac{\partial E(h, v)}{\partial \theta} \\ \therefore \frac{\partial}{\partial \theta} \log P(v) &= \sum_{h, v} P(h, v) \cdot \frac{\partial E(h, v)}{\partial \theta} - \sum_h P(h|v) \cdot \frac{\partial}{\partial \theta} E(h, v) \end{aligned}$$

基于 energy-based model, a log-likelihood

23.6 log-Likelihood gradient of RBM

Objective: Given training set: $v \in S$, $|S|=N$

log-likelihood gradient for RBM:

$$\frac{\partial}{\partial \theta} \cdot \frac{1}{N} \cdot \sum_{v \in S} \log P(v)$$

$$\theta = \{w, a, \beta\}, S: \text{training set}, \frac{\partial}{\partial \theta} \log P(v) = \sum_{h, v} P(h, v) - \frac{\partial E(h, v)}{\partial \theta} - \sum_h P(h|v) \cdot \frac{\partial}{\partial \theta} E(h, v)$$

$w = [w_{ij}]_{m \times n}$:

$$\frac{\partial}{\partial w_{ij}} \log P(v) = - \sum_h P(h|v) \cdot \frac{\partial E(h, v)}{\partial w_{ij}} + \sum_{h, v} P(h, v) \cdot \frac{\partial E(h, v)}{\partial w_{ij}}$$

$$\because E(h, v) = -(h^T w \cdot v + a^T v + \beta^T h) = -(h^T w v + \Delta) = -(\sum_{i=1}^m \sum_{j=1}^n h_i w_{ij} v_j + \Delta)$$

$$\therefore \frac{\partial}{\partial w_{ij}} \log P(v) = - \sum_h P(h|v) \cdot (-h_i v_j) + \sum_{h, v} P(h, v) \cdot (-h_i v_j)$$

$$= \underbrace{\sum_h P(h|v) \cdot h_i v_j}_{\text{①}} - \underbrace{\sum_{h, v} P(h, v) \cdot h_i v_j}_{\text{②}}$$

$h_i \neq 0, = 0$ 就约掉

$$\text{RBM: 值, } v_i \in \{0, 1\}, h_i \in \{0, 1\}$$

$$\text{①} = \sum_{h_1} \sum_{h_2} \dots \sum_{h_m} P(h_1, h_2, \dots, h_m | v) \cdot h_i v_j = \sum_{h_i} P(h_i | v) \cdot h_i v_j = \boxed{P(h_i=1 | v) \cdot v_j}$$

\downarrow P62 有求解公式

$$\text{②} = \sum_{h, v} P(v) \cdot P(h|v) \cdot h_i v_j = \sum_v P(v) \cdot \sum_h P(h|v) \cdot h_i v_j = \boxed{\sum_v P(v) \cdot P(h_i=1 | v) \cdot v_j}$$

同①

$v \in \{0, 1\}$, $\log(1, 2^{100})$, intractable
CD来解决.

23.7 CD-K Algorithm for RBM

log-Likelihood gradient for RBM:

$$\frac{\partial}{\partial w_{ij}} \log P(v) = P(h_i=1|v) v_j - \sum_v P(v) P(h_i=1|v) \cdot v_j = E_{P(v)} [P(h_i=1|v) v_j]$$

$$\frac{\partial}{\partial w_{ij}} \cdot \frac{1}{N} \cdot \sum_{v \in S} \log P(v) = \frac{1}{N} \sum_{v \in S} \frac{\partial}{\partial w_{ij}} \log P(v) \Leftrightarrow \frac{1}{N} \Delta w_{ij}$$

CD-K for RBM

For each $v \in S$

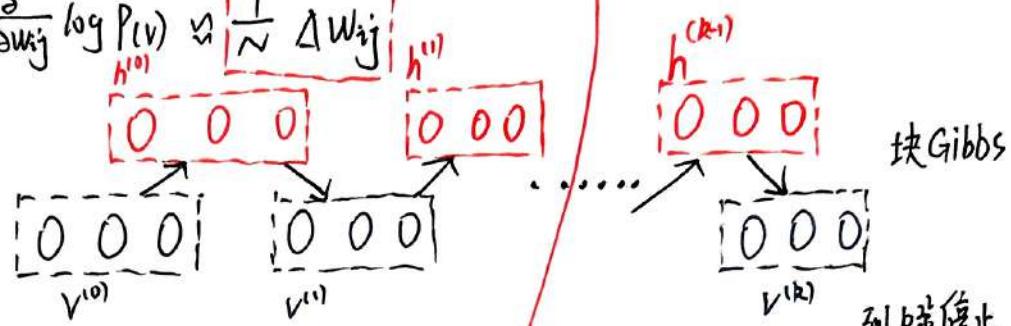
$$v^{(0)} \leftarrow v$$

For $k=0, 1, 2, \dots, k-1$:

For $i=1, 2, \dots, m$: Sample $h_i^{(k)} \sim P(h_i|v^{(k)})$
 For $j=1, 2, \dots, n$: Sample $v_j^{(k+1)} \sim P(v_j|h^{(k)})$

For $i=1, 2, \dots, m, j=1, 2, \dots, n$:

$$\Delta w_{ij} \leftarrow \Delta w_{ij} + \frac{\partial}{\partial w_{ij}} \log P(v) \approx P(h_i=1|v^{(0)}) \cdot v_j^{(0)} - P(h_i=1|v^{(k)}) \cdot v_j^{(k)}$$



24. Approximate Inference

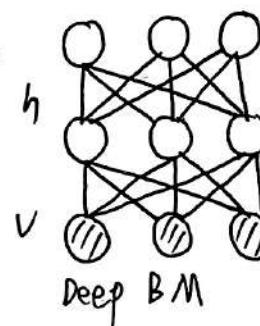
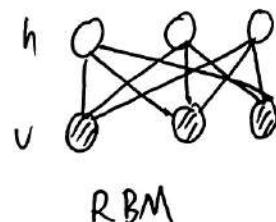
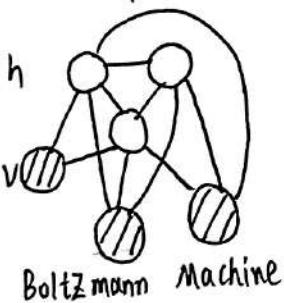
24.1. Introduction

① 推断的动机 { 推断本身 ($P(h|v)$)

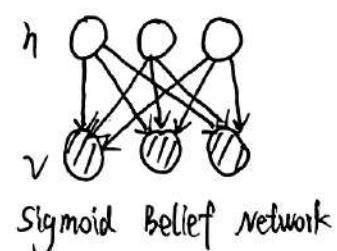
学习需要

② 推断是困难的 { 有向图: explain away

无向图: mutual interaction



explain away
head to head, 合成后, h 不是相互独立



24.2 Inference is optimization

$V = \{v\}$ 样本

log-likelihood: $\sum_{v \in V} \log P(v)$

$$\log P(v) = \log \frac{P(v, h)}{P(h|v)} = \log \frac{P(v, h)}{q(h|v)} \cdot \frac{q(h|v)}{P(h|v)} = \log \frac{P(v|h)}{q(h|v)} + \log \frac{q(h|v)}{P(h|v)}$$

$$\int (\log P(v)) \cdot q(h|v) dh = \int \log \frac{P(v, h)}{q(h|v)} \cdot q(h|v) dh + \int \log \frac{q(h|v)}{P(h|v)} \cdot q(h|v) dh$$

$$\log P(v) \int q(h|v) dh \rightarrow \log P(v) = E_{q(h|v)} [\log \frac{P(v, h)}{q(h|v)}] + E_{q(h|v)} [k L(q(h|v) || P(h|v))]$$

$$= E_{q(h|v)} [\log P(v, h) - \log q(h|v)] + k L(q(h|v) || P(h|v))$$

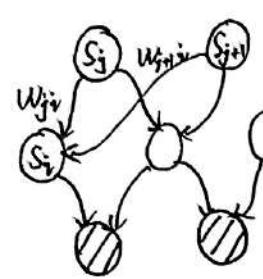
$$= E_{q(h|v)} [\log P(v, h)] - H[q(h|v)] + k L(q(h|v) || P(h|v))$$

$$\int q(h|v) \log q(h|v) = H[q(h|v)]$$

25. Sigmoid Belief Network

25.1 Introduction Neal → a student of Hinton

$$S = \{S_1, S_2, \dots, S_n\} = \{V, h\} = \{V, h^{(1)}, h^{(2)}\}, V, h \text{ 为 } \{0, 1\} \text{ 二值}$$



$$h^{(1)} \quad P(S_i=1) = \sigma(\sum_{j < i} w_{ji} S_j)$$

$$h^{(1)} \quad P(S_i=0) = 1 - P(S_i=1) = \sigma(-\sum_{j < i} w_{ji} S_j)$$

$$P(S_i | S_j : j < i) = \sigma(S_i^* \sum_{j < i} w_{ji} S_j)$$

$$\begin{cases} S_i^* = 2S_i - 1 \\ S_i = 1, S_i^* = 1 \\ S_i = 0, S_i^* = -1 \end{cases}$$

$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$

Neal → [1990]

$$1 - \sigma(x) = \frac{\exp(-x)}{1 + \exp(-x)} = \frac{1}{1 + \exp(x)} = \sigma(-x)$$

$$G(x) = G_m \cdot \sigma(-x)$$

-一个式子表达

25.2 gradient of log-likelihood

$$P(S_i | S_j : j < i) = \sigma(S_i^* \sum_{j < i} w_{ji} S_j) \quad P(S) = \prod_i P(S_i | S_j : j < i) = P(V, h)$$

$$\text{log-likelihood: } \sum_{v \in V} \log P(v)$$

$$\frac{\partial \log P(v)}{\partial w_{ji}} = \frac{1}{P(v)} \cdot \frac{\partial P(v)}{\partial w_{ji}} = \frac{1}{P(v)} \cdot \frac{\partial \sum_h P(h, v)}{\partial w_{ji}} = \sum_h \frac{1}{P(v)} \cdot \frac{\partial P(h, v)}{\partial w_{ji}} = \sum_h \frac{P(h|v)}{P(h, v)} \cdot \frac{\partial P(v, h)}{\partial w_{ji}}$$

$$= \sum_h P(h|v) \cdot \frac{1}{P(S)} \cdot \frac{\partial P(S)}{\partial w_{ji}}$$

$$\frac{1}{P(S)} \cdot \frac{\partial P(S)}{\partial w_{ji}} = \frac{1}{\prod_k P(S_k | S_j : j < k)} \cdot \frac{\Delta_b \cdot \partial P(S_i | S_j : j < i)}{\partial w_{ji}} = \frac{1}{P(S_i | S_j : j < i) \cdot \Delta_b} \cdot \frac{\Delta_b \cdot \partial P(S_i | S_j : j < i)}{\partial w_{ji}}$$

丁项中没有项与 w_{ji} 相关

$$= \frac{1}{P(S_i | S_j : j < i)} \cdot \frac{\partial \sigma(S_i^* \sum_{j < i} w_{ji} S_j)}{\partial w_{ji}}$$

$$= \frac{1}{\sigma(S_i^* \sum_{j < i} w_{ji} S_j)} \cdot \sigma(S_i^* \sum_{j < i} w_{ji} S_j) \cdot \sigma(-S_i^* \sum_{j < i} w_{ji} S_j) \cdot S_i \cdot S_j$$

$$= \sigma(-S_i^* \sum_{j < i} w_{ji} S_j) \cdot S_i^* \cdot S_j$$

$$\sigma(x) = \frac{1}{1 + \exp(-x)} = \frac{\exp(x)}{1 + \exp(x)}$$

$$\sigma'(x) = \frac{\exp(-x)}{(1 + \exp(-x))^2} = \frac{1}{1 + \exp(-x)} \cdot \frac{\exp(-x)}{1 + \exp(-x)} = \sigma(x) \cdot \sigma(-x)$$

$$\therefore \frac{\partial \log P(v)}{\partial w_{ji}} = \sum_h P(h|v) \cdot \sigma(-S_i^* \sum_{j < i} w_{ji} S_j) \cdot S_i^* \cdot S_j$$

$$\therefore \frac{\partial}{\partial w_{ji}} \sum_{v \in V} \log P(v) = \sum_{v \in V} \sum_h P(h|v) \cdot \sigma(-S_i^* \sum_{j < i} w_{ji} S_j) \cdot S_i^* \cdot S_j = \sum_{v \in V} \sum_S P(S|v) \cdot \sigma(-S_i^* \sum_{j < i} w_{ji} S_j) \cdot S_i^* \cdot S_j$$

$$= E_{(v, h) \sim P(s|v)} [\sigma(-S_i^* \sum_{k < i} w_{ki} S_k) \cdot S_i^* \cdot S_j]$$

$$\text{v.v} P_{\text{data}} = P_{\text{v}}$$

$$\text{豆包} = E_{v \sim P_{\text{data}}} [E_{s \sim P(s|v)} [\sigma(-S_i^* \sum_{k < i} w_{ki} S_k) \cdot S_i^* \cdot S_j]] \quad \text{全期期望表示}$$

$$E_{v, s} [f_{(v, s)}] = E_v [E_{s|v} [f_{(v, s)}|v]]$$

↓ MCMC 采样 (Neal 小规模)

$$P(h|v) = P(h, v|v) = P(s|v)$$

$P(h|v)$ 算不出 explain away

$$E_{v \sim P_{\text{data}}}[E_{s \sim P(s|v)}[f(v, s)|v]] \approx \frac{1}{M} \sum_{m=1}^M \left[\frac{1}{N} \sum_{n=1}^N f(v^{(m)}, s^{(m, n)}) \right]$$

$$f(v, s) = s \left(-S_i \sum_{j < i} w_{ji} s_j \right) \cdot S_i \cdot S_j \quad (\text{目标函数})$$

$v^{(m)}$ 是从 $P_{\text{data}}(v)$ 采样的观测样本

$s^{(m, n)}$ 是对每个 $v^{(m)}$, 从 $P(s|v^{(m)})$ 采样的隐变量样本

外层期望 $v \sim P_{\text{data}}$, 直接采

假设我们有训练数据集 $D = \{v^{(1)}, v^{(2)}, \dots, v^{(M)}\}$, 则 $P_{\text{data}}(v) \approx \frac{1}{|D|} \sum_{v \in D} \delta(v - \hat{v})$ (经验分布)

每次训练从 D 中随机抽取 M 个样本 $\{v^{(1)}, v^{(2)}, \dots, v^{(M)}\}$

内层期望 $s \sim P(s|v)$, Gibbs 高维

内层的 $P(s|v)$ 是高维依赖分布, 用 Gibbs 采样

25.3 Wake-sleep Algorithm - Introduction (Hinton)

$$\text{Gradient of log-likelihood: } \sum_{v \in V} \sum_S P(s|v) \cdot S_i^* \cdot S_j^* (-S_i \sum_{k < i} S_k w_{ki})$$



wake Phase:

① Bottom-up: 激活 neuron (获得各层样本)

② Learning Generative Connection (t \to w)

sleep Phase:

① Top-down: 激活 neuron (获得样本)

② Learning Recognition connection (t \to R)

$p(h|v)$ 求不出, 近似 - \{q(h|v), q(h|v)\} 的参数就是 R

25.4 Wake-sleep Algorithm - KL Divergence

Wake Phase: 固定 \uparrow R \downarrow w EM \Rightarrow M-step

$$E_{q_\phi(h|v)}[\log p_\theta(h, v)] \approx \bar{N}$$

$$\hat{\theta} = \arg \max_{\theta} E_{q_\phi(h|v)}[\log p_\theta(h, v)], \text{ with } \phi \text{ fixed}$$

H[q] 的参数是 \phi, \phi 是固定的, 与 \theta 无关

$$\hat{\theta} = \arg \max_{\theta} L(\theta) \quad \leftarrow [ELBO]_{\max} \leftarrow [KL]_{\min} \leftarrow \arg \min_{\theta} \text{KL}[q_\phi(h|v) || p_\theta(h|v)]$$

Sleep Phase: 固定 \downarrow w \uparrow R EM \Rightarrow E-step

$$\hat{\phi} = \arg \max_{\phi} E_{p_\theta(h|v)}[\log q_\phi(h|v)], \text{ with } \theta \text{ fixed}$$

$$= \arg \max_{\phi} \int p_\theta(h, v) \cdot \log q_\phi(h|v) dh$$

$$= \arg \max_{\phi} \int p_\theta(v) \cdot p_\theta(h|v) \cdot \log q_\phi(h|v) dh$$

只与 \theta 有关, 与 h, \phi 无关, 提出后省略

$$\log p(v) = ELBO + KL(q||p)$$

$$ELBO = L = \mathbb{E}_{q(h|v)} \left[\log \frac{p(v, h)}{q(h|v)} \right]$$

$$= \mathbb{E}_{q(h|v)} [\log p(v, h)] + H[q]$$

$$= \arg \max_{\phi} \int p_\theta(v) \int p_\theta(h|v) \cdot \log q_\phi(h|v) dh$$

$$= \arg \max_{\phi} \int p_\theta(h|v) \cdot \log q_\phi(h|v) dh$$

$$= \arg \max_{\phi} \int p_\theta(h|v) \cdot \log \left[\frac{q_\phi(h|v)}{p_\theta(h|v)} \cdot p_\theta(h|v) \right] dh$$

$$= \arg \max_{\phi} \int p_\theta(h|v) \cdot \log \frac{q_\phi(h|v)}{p_\theta(h|v)} + \underbrace{p_\theta(h|v) \cdot p_\theta(h|v)}_{\text{与 } \phi \text{ 无关}} dh$$

$$= \arg \max_{\phi} \int p_\theta(h|v) \cdot \log \frac{q_\phi(h|v)}{p_\theta(h|v)} dh$$

$$= \arg \max_{\phi} -KL[p_\theta(h|v) || q_\phi(h|v)]$$

$$= \arg \min_{\phi} \{KL[p_\theta(h|v) || q_\phi(h|v)]\}$$

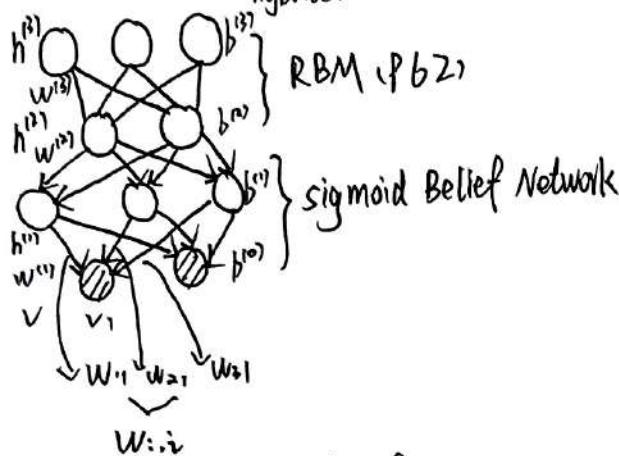
26. Deep Belief Network

有向图 Model

Belief Network \leftrightarrow Bayesian Network

26.1 Introduction

Hinton (2006) hybrid model



$w_{:,i}$: i-th column vector of w

$$\Theta = \{w^{(1)}, w^{(2)}, w^{(3)}, b^{(1)}, b^{(2)}, b^{(3)}\}$$

$$P(v, h^{(1)}, h^{(2)}, h^{(3)}) = P(v|h^{(1)}, h^{(2)}, h^{(3)}) \cdot P(h^{(1)}, h^{(2)}, h^{(3)})$$

$$= P(v|h^{(1)}) \cdot P(h^{(1)}, h^{(2)}, h^{(3)})$$

$$= P(v|h^{(1)}) \cdot P(h^{(1)}|h^{(2)}, h^{(3)}) \cdot P(h^{(2)}, h^{(3)})$$

$$= P(v|h^{(1)}) \cdot P(h^{(1)}|h^{(2)}) \cdot P(h^{(2)}, h^{(3)})$$

$$= \prod_{i=1}^3 P(v_i|h^{(1)}) \cdot \prod_{j=1}^2 P(h_j^{(1)}|h^{(2)}) \cdot P(h^{(2)}, h^{(3)})$$

$$P(v_i|h^{(1)}) = \text{sigmoid}(w_{:,i}^\top \cdot h^{(1)} + b_i)$$

到 v_i 该值的所有边 $\rightarrow w_{:,i}^\top$

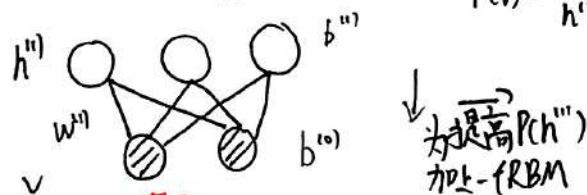
利用
Pb2

$$P(h_j^{(1)}|h^{(2)}) = \text{sigmoid}(w_{:,j}^\top \cdot h^{(2)} + b_j)$$

$$P(h^{(1)}, h^{(2)}) = \frac{1}{Z} \exp \left\{ h^{(1)^\top} \cdot w^{(1)} \cdot h^{(1)} + h^{(2)^\top} \cdot w^{(2)} \cdot h^{(2)} + h^{(1)^\top} \cdot b^{(1)} + h^{(2)^\top} \cdot b^{(2)} \right\}$$

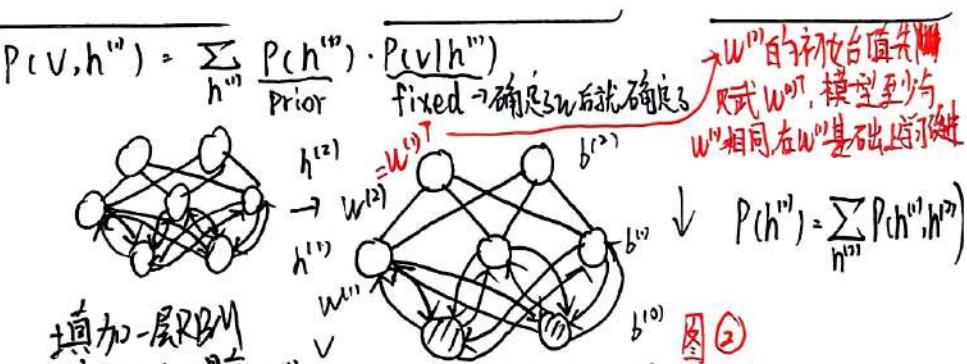
$$E(h, v) = -\text{dot product}$$

26.2 Stacking RBM



图①

将 RBM 由无向变为有向，需加上蓝色箭头，而假设中固定 $P(v|h^{(1)})$ ，则需擦去蓝色，留下黑色箭头，变成 DBN。固定 $P(v|h^{(1)})$ 即学习出了 $w^{(1)}$ ，为了改进 $w^{(1)}$ ，则为 $h^{(1)}$ 增加一层或多层。



图②

固定 $P(v|h^{(1)})$ 后，向上的蓝色箭头擦去，就变成了 DBN

26.3 Extra layers improve ELBO

关于上节图①

$$\begin{aligned} \log P(v) &= \log \sum_{h^{(1)}} P(v, h^{(1)}) \\ &= \log \sum_{h^{(1)}} q(h^{(1)}|v) \cdot \frac{P(v, h^{(1)})}{q(h^{(1)}|v)} \\ &= \log \sum_{h^{(1)}} q(h^{(1)}|v) \cdot \frac{P(v, h^{(1)})}{q(h^{(1)}|v)} \\ &= \log E_{q(h^{(1)}|v)} \left[\frac{P(v, h^{(1)})}{q(h^{(1)}|v)} \right] \\ &\geq E_{q(h^{(1)}|v)} \left[\log \frac{P(v, h^{(1)})}{q(h^{(1)}|v)} \right] \\ &= \sum_{h^{(1)}} q(h^{(1)}|v) \left[\log P(v, h^{(1)}) - \log q(h^{(1)}|v) \right] \end{aligned}$$

$$= \sum_{h^{(1)}} q(h^{(1)}|v) \left[\log P(h^{(1)}) + \log P(v|h^{(1)}) - \log q(h^{(1)}|v) \right]$$

improve $P(h^{(1)})$

都确定

关于图②，2nd layer RBM Learning \Leftrightarrow maximum log-likelihood over $P(h^{(1)})$

\Leftrightarrow maximum ELBO of $P(v)$

Learning 是求出最好的 $w^{(2)}$ 和 $b^{(2)}$ ，使 $P(h^{(1)})$ 的 likelihood 达到最大 ($P(h^{(1)})$)，从而在其余部分已知的基础上增大 ELBO，从而增大 $\log P(v)$

26.4 Pre-training

$$\log P(v) \geq E\text{LBO} = \sum_{h''} q(h''|v) \cdot \log P(v, h'') - \sum_{h''} q(h''|v) \cdot \log P(h''|v)$$

求不出，只能近似 (explain away)

当作 RBM 来算，避免用 explain away，用 RBM 近似，则 h'' 中两个 0 独立。

$$q(h''|v) = \prod_i q(h_i''|v) = \prod_i \text{sigmoid}(w_i^{(1)} \cdot v + b_i^{(1)})$$

向量 偏置

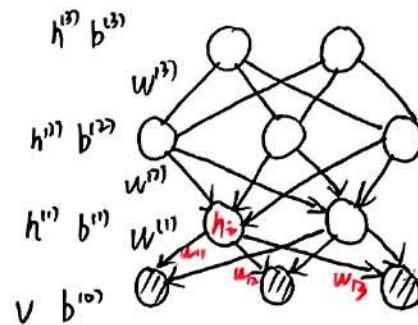
得到 w_{i1} 与 b_{i1}

从此层层递进，求 $w_2, b_2, w_3, b_3, \dots$

$DBN \rightarrow ELBO$ is relatively loose 比较松

~~此法忽略了 $w^{(1)}$ 对 $h^{(1)}$ 的影响~~

$P(h''|v)$ 选出 $w^{(1)}$ ，再抽样出 $h^{(1)}$ ，然后用 $h^{(1)}$ ， $P(h''|h^{(1)})$ 选出 $w^{(2)}$ ，此时忽略 $w^{(1)}$ 对 $h^{(2)}$ 影响，会用忽略 $w^{(1)}$ 用 $P(h^{(2)}; w^{(2)})$ 近似 $P(h^{(2)}; w^{(1)}, w^{(2)})$



Generative process
生成样本
容易
Gibbs Sampling

$q(h''|v)$ 是一个近似后验，因子分解后的一个近似分布
first, second layer 当作 RBM 求 $q(h''|v)$

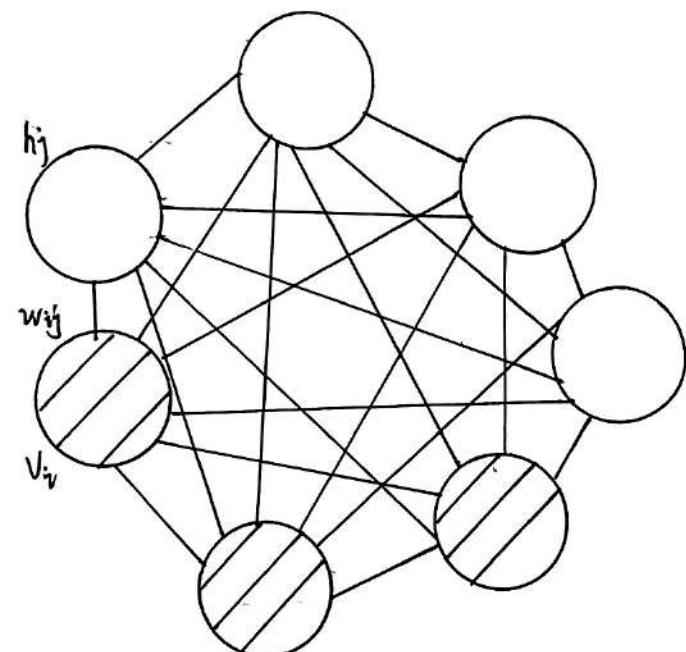
与真实的 $P(h''|v)$ 差距大， $P(h''|v)$ 并不是可因子分解的
故 ELBO 比较松 两层 与 P18 DBM 对应

$P(h''|v)$ 选出 $w^{(1)}$ ，再抽样出 $h^{(1)}$ ，然后用 $h^{(1)}$ ， $P(h''|h^{(1)})$ 选出 $w^{(2)}$ ，此时忽略 $w^{(1)}$ 对 $h^{(2)}$ 影响，会用忽略 $w^{(1)}$ 用 $P(h^{(2)}; w^{(2)})$ 近似 $P(h^{(2)}; w^{(1)}, w^{(2)})$

27. Boltzmann Machine

$$\Theta = \{W, L, J\}; V: 样本集合, |V|=N=D+P$$

27.1 Introduction.



$$V = \{0,1\}^D, h = \{0,1\}^P, 分量 0 或 1, 大 D 大 P 列向量$$

$$L = [L_{ij}]_{D \times D}, J = [J_{ij}]_{P \times P}, W = [W_{ij}]_{D \times P} \quad (h \text{ 与 } v \text{ 的权重})$$

$$P(v, h) = \frac{1}{Z} \exp \{-E(h, v)\}$$

$$E(h, v) = -(V^T \cdot W \cdot h + \frac{1}{2} V^T \cdot L \cdot V + \frac{1}{2} h^T \cdot J \cdot h)$$

$$Z = \int_h \int_v \exp \{-E(h, v)\} dv dh$$

L 与 J 对角线上为 0，自己不会相连

$$\begin{matrix} v_1 & v_2 & v_3 \\ v_1 & 0 & \Delta \\ v_2 & \Delta & 0 \\ v_3 & 0 & 0 \end{matrix} \quad \text{上三行下三列乘上对角}$$

$$\text{若无 } L, \sum_{i=1}^3 \sum_{j=1}^3 v_i w_{ij} v_j = 3 \times 3 = 9 \text{ 次}$$

$$\text{实际上只用 } (9-3) \times \frac{1}{2} = 3 \text{ 次。}$$

27.2 Gradient of log-Likelihood

$$P(v) = \sum_h P(v, h)$$

$$\frac{1}{N} \sum_{v \in V} \log P(v) \leftarrow \text{log-Likelihood}$$

$$\frac{\partial}{\partial \theta} \cdot \frac{1}{N} \sum_{v \in V} \log P(v) = \frac{1}{N} \sum_{v \in V} \frac{\partial \log P(v)}{\partial \theta} \leftarrow \text{gradient of log-Likelihood}$$

$$\frac{\partial \log P(v)}{\partial \theta} = \sum_v \sum_h P(v, h) \cdot \frac{\partial E(v, h)}{\partial \theta} - \sum_h P(h|v) \frac{\partial E(h, v)}{\partial \theta} \quad P68$$

$$\frac{\partial \log P(v)}{\partial w} = \sum_v \sum_h P(v, h) \cdot (-vh^T)_{\text{dep}} - \sum_h P(h|v) \cdot (-vh^T)$$

$$= \sum_h P(h|v) \cdot vh^T - \sum_v \sum_h P(v, h) \cdot vh^T$$

$$\therefore \frac{1}{N} \sum_{v \in V} \frac{\partial \log P(v)}{\partial w} = \frac{1}{N} \sum_{v \in V} \sum_h P(h|v) \cdot vh^T - \frac{1}{N} \sum_{v \in V} \sum_v \sum_h P(v, h) \cdot vh^T \quad \text{大V是样本集合}$$

看作是关于 v, h 的 function, v 和 h 已经被积掉, 与 $\sum_{v \in V}$ 来说
只是一个常数

L, J 求法与 w 类似

$$\frac{1}{N} \sum_{v \in V} = P_{\text{data}(w)} = \frac{1}{N} \sum_{v \in V} \sum_h P(h|v) \cdot vh^T - \sum_v \sum_h P(v, h) \cdot vh^T \rightarrow E_{v \sim P_{\text{data}}, h \sim P_{\text{model}}(h|v)} [vh^T] = E_{P_{\text{data}}} [vh^T]$$

从数据来
数据分布

$$= E_{P_{\text{data}}} [vh^T] - E_{P_{\text{model}}} [vh^T]$$

梯度只与 vh^T 有关, 有 MCMC sampling
比较干净

能直接得到 P_{data}

需要利用模型分布 P_{model}

这里 P_{data} 隐式包含
 $P_{\text{model}}(h|v)$, 因为的分布由
模型给定

$$\text{样本 } v \sim i.i.d. P_{\text{data}(w)}, P_{\text{data}} = P_{\text{data}}(v) \cdot P_{\text{model}}(h|v)$$

$$\frac{1}{N} \sum_{v \in V} \approx E_{v \sim P_{\text{data}}} \quad P_{\text{model}} = P_{\text{model}}(v, h)$$

27.3 Gradient ascend based on MCMC

$$\Delta w = d(E_{P_{\text{data}}} [vh^T] - E_{P_{\text{model}}} [vh^T])$$

$$\Delta L = d(E_{P_{\text{data}}} [vv^T] - E_{P_{\text{model}}} [vv^T])$$

$$\Delta J = d(E_{P_{\text{data}}} [hh^T] - E_{P_{\text{model}}} [hh^T])$$

$$P_{\text{data}} = P_{\text{data}}(v, h) = P_{\text{data}}(v) \cdot P_{\text{model}}(h|v)$$

$$P_{\text{model}} = P_{\text{model}}(v, h), d: 步长$$

上述 $\Delta w, \Delta L, \Delta J$ 都防矩阵

$$\Delta w_{ij} = d(E_{P_{\text{data}}} [v_i h_j] - E_{P_{\text{model}}} [v_i h_j])$$

positive phase negative phase

Both hard, intractable
MCMC (1983-1985 Hinton)

利用 Boltzmann Machine 的定义, P74 的 $P(v, h)$ 与上 $E(v, h)$

报告导出: 假设已证明

$$P(v_i=1|h, v_{-i}) = \sigma(\sum_{j=1}^p w_{ij} h_j + \sum_{k=1, k \neq i}^D J_{ik} v_k) \quad \sigma \rightarrow \text{Sigmoid} \rightarrow \text{Iter}$$

$$P(h_j=1|V, h_{-j}) = \sigma(\sum_{i=1}^D w_{ij} v_i + \sum_{m=1, m \neq j}^p J_{jm} h_m)$$

规模大效率特别差

这是 P76 导出的

而 sigmoid Belief network 中
的 Sigmoid Function 是人为定义的,

27.4 Conditional Probability

推导: $P(V_i=1 | h, v_{-i}) = \sigma(\sum_{j=1}^P W_{ij} h_j + \sum_{k=1, k \neq i}^D L_{ik} V_k)$

$$\sigma = \frac{e^x}{1+e^x} = \frac{1}{1+e^{-x}} \text{ sigmoid}$$

$$P(h_j=1 | v, h_{-j}) = \sigma(\sum_{i=1}^D W_{ij} v_i + \sum_{m=1, m \neq j}^P J_{jm} h_m)$$

$$P(V_i=1 | h, v_{-i}) = \frac{P(v, h)}{P(h, v_{-i})} = \frac{\frac{1}{Z} \exp\{-E(h, v)\}}{\sum_{v_i} \frac{1}{Z} \exp\{-E(h, v)\}} = \frac{\exp\{v^T w h + \frac{1}{2} v^T L v + \frac{1}{2} h^T J h\}}{\sum_{v_i} \exp\{v^T w h + \frac{1}{2} v^T L v + \frac{1}{2} h^T J h\}}$$

$$= \frac{\exp\{\frac{1}{2} h^T J h\} \cdot \exp\{v^T w h + \frac{1}{2} v^T L v\}}{\exp\{\frac{1}{2} h^T J h\} \cdot \sum_{v_i} \exp\{v^T w h + \frac{1}{2} v^T L v\}} = \frac{\exp\{v^T w h + \frac{1}{2} v^T L v\}|_{v_i=0}}{\exp\{v^T w h + \frac{1}{2} v^T L v\}|_{v_i=1}}$$

对 v 进行实例化令
 $v_{-i}, \forall i \neq i$ 为常数
 v_i 为变量
 \therefore 只有分子改变

$$P(v_{i=1} | h, v_{-i}) = \frac{\exp\{v^T w h + \frac{1}{2} v^T L v\}|_{v_i=1}}{\exp\{v^T w h + \frac{1}{2} v^T L v\}|_{v_i=0} + \exp\{v^T w h + \frac{1}{2} v^T L v\}|_{v_i=1}} = \frac{\Delta v_{i=1}}{\Delta v_{i=0} + \Delta v_{i=1}}$$

$\Delta v_i = \exp\{v^T w h + \frac{1}{2} v^T L v\} = \exp\{\sum_{j=1}^P \sum_{i=1}^D v_i^j W_{ij} h_j + \frac{1}{2} \sum_{j=1}^P \sum_{k=1, k \neq i}^D v_i^j L_{ik} v_k\}$

$$= \exp\{\sum_{i=1, i \neq i}^P \sum_{j=1}^D v_i^j W_{ij} h_j + \sum_{j=1}^P v_i^j W_{ij} h_j + \frac{1}{2} (\sum_{i=1, i \neq i}^P \sum_{k=1, k \neq i}^D v_i^j L_{ik} v_k + v_i^j L_{ii} v_i + \sum_{i=1, i \neq i}^P v_i^j L_{ik} v_k)$$

$L_{ii}=0, \forall i=0$ $k \neq i, i \neq i$ $k \neq i, i \neq i$

$\Delta v_{i=0} = \exp\{\sum_{i=1, i \neq i}^P \sum_{j=1}^D v_i^j W_{ij} h_j + \frac{1}{2} (\sum_{i=1, i \neq i}^P \sum_{k=1, k \neq i}^D v_i^j L_{ik} v_k + 2 \sum_{k=1, k \neq i}^D v_i^j L_{ik} v_k)\}$

$\Delta v_{i=1} = \exp\{\sum_{i=1, i \neq i}^P \sum_{j=1}^D v_i^j W_{ij} h_j + \frac{1}{2} \sum_{i=1, i \neq i}^P \sum_{k=1, k \neq i}^D v_i^j L_{ik} v_k\}$

$$\therefore P(v_{i=1} | h, v_{-i}) = \frac{\Delta v_{i=1}}{\Delta v_{i=0} + \Delta v_{i=1}} = \frac{\exp\{A+B + \sum_{j=1}^P W_{ij} h_j + \sum_{k=1, k \neq i}^D L_{ik} v_k\}}{\exp\{A+B\} + \exp\{A+B + \sum_{j=1}^P W_{ij} h_j + \sum_{k=1, k \neq i}^D L_{ik} v_k\}}$$

Sigmoid: $\frac{e^x}{1+e^x} = \frac{1}{1+e^{-x}}$

$$= \frac{\exp\{\sum_{j=1}^P W_{ij} h_j + \sum_{k=1, k \neq i}^D L_{ik} v_k\}}{1 + \exp\{\sum_{j=1}^P W_{ij} h_j + \sum_{k=1, k \neq i}^D L_{ik} v_k\}} = \sigma(\sum_{j=1}^P W_{ij} h_j + \sum_{k=1, k \neq i}^D L_{ik} v_k)$$

27.5 Variational Inference (Hinton)

根据P70类似: $\mathcal{L} = \text{ELBO} = \log P_\theta(v) - \text{KL}(q_\phi || P_\theta) = \sum_h q_\phi(h|v) \cdot \log P_\theta(v, h) + H[q]$

基于平移场

$$q_\phi(h|v) = \prod_{j=1}^P q_\phi(h_j|v), \quad q_\phi(h_j=1|v) = \phi_j, \quad q_\phi(h_j=0|v) = 1 - \phi_j, \quad \phi = \{\phi_j\}_{j=1}^P$$

$$\hat{\phi}_j = \arg \max_{\phi_j} \mathcal{L} = \arg \max_{\phi_j} \sum_h q_\phi(h|v) [-\log Z + v^T \cdot w \cdot h + \frac{1}{2} v^T \cdot L \cdot v + \frac{1}{2} h^T J \cdot h] + H[q]$$

$$= \arg \max_{\phi_j} \sum_h \underbrace{q_\phi(h|v)}_1 \cdot \underbrace{[-\log Z + \frac{1}{2} v^T \cdot L \cdot v]}_{与 \phi_j, h \text{ 无关}} + \sum_h q_\phi(h|v) [v^T \cdot w \cdot h + \frac{1}{2} h^T J \cdot h] + H[q]$$

$$= \arg \max_{\phi_j} \sum_h q_\phi(h|v) [v^T w h + \frac{1}{2} h^T J h] + H[q]$$

$$= \arg \max_{\phi_j} \sum_h q_\phi(h|v) \cdot v^T w h + \frac{1}{2} \sum_h q_\phi(h|v) \cdot h^T J h + H[q] = \arg \max_{\phi_j} \textcircled{1} + \textcircled{2} + \textcircled{3}$$

$$\textcircled{1} = \sum_h q_\phi(h|v) \cdot \sum_{i=1}^D \sum_{j=1}^P w_{ij} h_j = \sum_h \prod_{j=1}^P q_\phi(h_j|v) \cdot \sum_{i=1}^D \sum_{j=1}^P v_i w_{ij} h_j$$

$$\text{其中-项} := \sum_h \prod_{j=1}^P q_\phi(h_j|v) \cdot v_1 w_{12} h_2 = \sum_{h_2} q_\phi(h_2|v) \cdot v_1 w_{12} h_2 \cdot \underbrace{\sum_{h_1 h_2} \prod_{j=1/2}^P q_\phi(h_j|v)}_{\sum_{h_1} q_\phi(h_1|v) \cdot \sum_{h_3} q_\phi(h_3|v) \cdots \sum_{h_P} q_\phi(h_P|v)} = 1$$

$$= \sum_{h_2} q_\phi(h_2|v) \cdot v_1 w_{12} h_2$$

$$\because h_2 = \{0, 1\}$$

$$\therefore \textcircled{1} = q_\phi(h_2=1|v) \cdot v_1 w_{12} = \phi_2 \cdot v_1 w_{12} = \phi_2 \cdot v_1 w_{12}$$

$$\therefore \textcircled{1} = \sum_{i=1}^D \sum_{j=1}^P v_i w_{ij} \cancel{\phi_j} \sum_h$$

$$\text{关于} \textcircled{2} = \frac{1}{2} \sum_h q_\phi(h|v) h^T J h = \frac{1}{2} \prod_{j=1}^P q_\phi(h_j|v) \cdot \sum_{i=1}^D \sum_{j=1}^P h_i J_{ij} h_j$$

$$\text{其中-项}: \therefore \sum_h \prod_{j=1}^P q_\phi(h_j|v) \cdot h_1 J_{12} h_2 = \sum_{h_1} \sum_{h_2} q_\phi(h_1|v) q_\phi(h_2|v) \cdot h_1 J_{12} h_2 \underbrace{\sum_{h_1 h_2} \prod_{j=3}^P q_\phi(h_j|v)}_{\text{对称性} \rightarrow \text{常数} C \rightarrow \text{消除} \rightarrow \text{其中-项.}}$$

$$\therefore = \sum_{h_1} \sum_{h_2} q_\phi(h_1|v) q_\phi(h_2|v) \cdot h_1 J_{12} h_2$$

$$\therefore h_1, h_2 = \{0, 1\}$$

$$\therefore (0,0) \times ; (0,1) \times ; (1,0) \times ; (1,1) \times$$

$$\therefore = q_\phi(h_1=1|v) q_\phi(h_2=1|v) J_{12} = \phi_1 \phi_2 J_{12}$$

$$\therefore \textcircled{2} = \sum_{j=1}^P \sum_{m=1}^P \phi_j \phi_m J_{jm} \quad (\text{要扣去其中一个} h \text{的值})$$

$$\textcircled{3} = - \sum_j [\phi_j \log \phi_j + (1 - \phi_j) \log (1 - \phi_j)]$$

$$\frac{\partial \mathcal{G}}{\partial \phi_j} = \sum_{i=1}^D v_i w_{ij} \quad ; \quad \frac{\partial \mathcal{G}}{\partial \phi_j} = \sum_{m=1}^P \phi_m J_{jm} \quad ; \quad \frac{\partial \mathcal{G}}{\partial \phi_j} = \log \frac{\phi_j}{1 - \phi_j}$$

$$q_\phi(h|v) \approx p_{\text{model}}(h|v)$$

注意, 这里中 j 与其它 ϕ_m 有关, $\sum_{m=1}^P$, 故一个一个的求, 一耗一耗的找, 一直到 j

$G: \text{Sigmoid: } \frac{1}{1 + e^{-x}}$ (x, y_1, \dots, y_N)

$$\frac{\partial [\textcircled{1} + \textcircled{2} + \textcircled{3}]}{\partial \phi_j} \triangleq 0$$

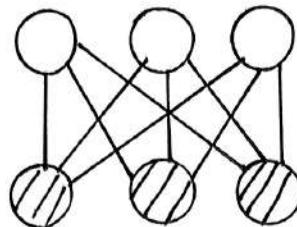
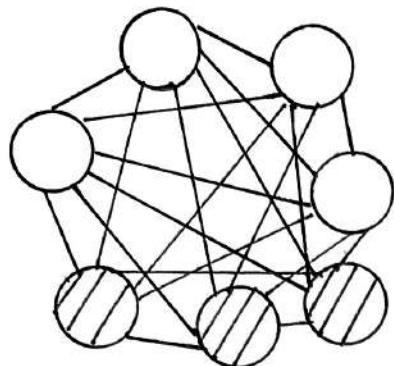
$$\phi_j = \sigma \left(\sum_{i=1}^D v_i w_{ij} + \sum_{m=1}^P \phi_m J_{jm} \right)$$

$$\hat{\phi} = \{\hat{\phi}_j\}_{j=1}^P$$

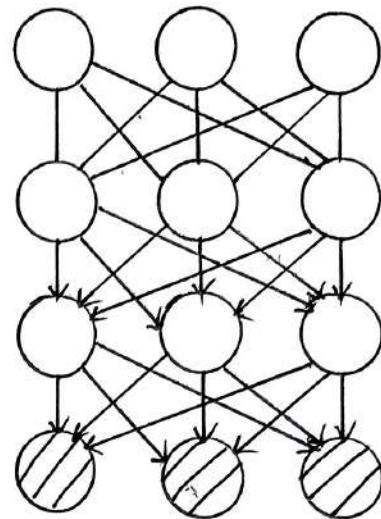
不动点方程 & 法 = 坐标上升

28. Deep Boltzmann Machine

28.1 Introduction



RBM 1980, 2002



DBN 2006

$\left\{ \begin{array}{l} \text{Pre-training (stacking RBM)} \\ \text{Fine-tuning (wake-sleep BP)} \end{array} \right.$

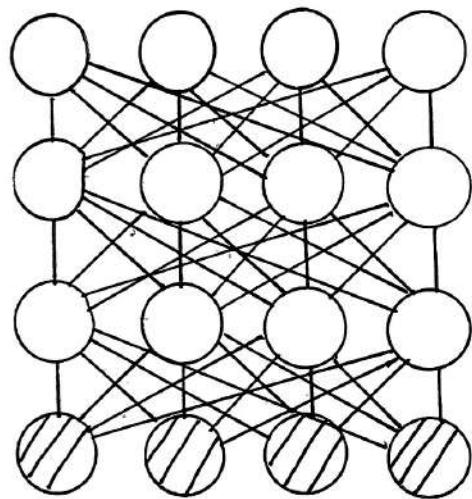
General Boltzmann Machine 1983

SGA = stochastic Gradient Ascend

$$\Delta W = \alpha (\underbrace{E_{P_{\text{data}}}[vh^\top]}_{\text{positive phase}} - \underbrace{E_{P_{\text{model}}}[vh^\top]}_{\text{negative phase}})$$

$$\left. \begin{array}{l} P_{\text{data}} = P_{\text{data}}(v, h) = P_{\text{data}}(v) \cdot P_{\text{data}, \text{model}}(h|v) \\ P_{\text{model}} = P_{\text{model}}(v, h) \end{array} \right\} \text{Variational Inference P77}$$

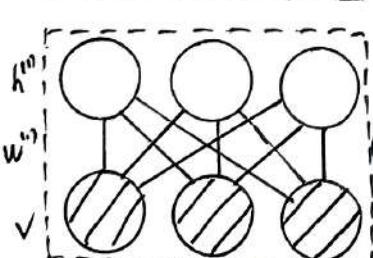
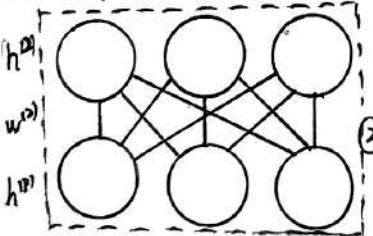
$$P_{\text{model}} = P_{\text{model}}(v, h) \rightarrow \text{PCD (逐层对比散度)}$$



$\left\{ \begin{array}{l} \text{Pre-training (stacking RBM)} \\ \text{SGA} \end{array} \right.$

DBM 2008

28.2 Pre-training



$$P(v) = \sum_{h^{(1)}} P(v, h^{(1)}) = \sum_{h^{(1)}} P(h^{(1)}) \cdot P(v|h^{(1)})$$

在 DBM 中，直接设 $w^{(1)}$ 给 $w^{(2)}$ 。在 $w^{(1)}$ 基础上训练，用 $P(h^{(1)}; w^{(1)})$

P74

与近似表达 $P(h^{(1)}; w^{(1)}, w^{(2)}) \rightarrow$ 用 $P(h^{(1)}; w^{(1)})$ 或 $P(h^{(1)}; w^{(2)})$ 表示 $P(h^{(1)}; w^{(1)}, w^{(2)})$
者不准确，下节介绍如何融合

$$\therefore P(v) = \sum_{h^{(1)}} P(v, h^{(1)}) \quad \text{在 DBM 中：} P(v) = \sum_{h^{(1)}} P(v, h^{(1)}) \quad P(h^{(1)}) = \sum_v P(v, h^{(1)})$$

$$\therefore P(v) = \sum_{h^{(1)}} P(h^{(1)}; w^{(1)}) \cdot P(v|h^{(1)}; w^{(1)})$$

对 $h^{(1)}$ Sampling:

$$P(h^{(1)}|v, w^{(1)}) = \prod_{i=1}^3 P(h_i^{(1)}|v; w^{(1)})$$

在 RBM 中训练解 P62
此时①中 $h^{(1)}$ 采样，采样后的数据在②中当作 v 重复上述①的步骤，叫做呼引 N 个 $h^{(1)}$

$$P(h^{(1)}; w^{(1)}) = \sum_{h^{(2)}} P(h^{(1)}, h^{(2)}; w^{(1)})$$

$$\therefore P(h^{(1)}; w^{(1)}) = \sum_{h^{(2)}, v} P(v, h^{(1)}, h^{(2)})$$

Intuition:
用 $P(h^{(1)}; w^{(1)})$ 和 $P(h^{(1)}; w^{(2)})$ 几何平均相似
 $P(h^{(1)}; w^{(1)}, w^{(2)}) \rightarrow h^{(1)} \text{ 与 } w^{(1)}, w^{(2)} \text{ 有关}$

P78

28.3 double counting Problem

真正的: $P(h^{(1)}; w^{(1)}, w^{(2)})$

直觉: 同时利用 $P(h^{(1)}; w^{(1)})$ 和 $P(h^{(1)}; w^{(2)})$ 去近似 $P(h^{(1)}; w^{(1)}, w^{(2)})$, $N \uparrow v \in V$

$$P(h^{(1)}; w^{(1)}) = \sum_v P(v, h^{(1)}; w^{(1)}) = \sum_v P(v) \cdot P(h^{(1)}|v; w^{(1)}) \approx \frac{1}{N} \sum_{v \in V} P(h^{(1)}|v; w^{(1)})$$

$$P(h^{(1)}; w^{(2)}) = \sum_{h^{(2)}} P(h^{(1)}, h^{(2)}; w^{(2)}) \stackrel{(1)}{=} \left(\sum_{h^{(2)}} P(h^{(2)}) \cdot P(h^{(1)}|h^{(2)}; w^{(2)}) \right) \stackrel{\text{Aggregated posterior}}{\approx} \frac{1}{N} \sum_{h^{(2)} \in H} P(h^{(1)}|h^{(2)}; w^{(2)}) \quad (2)$$

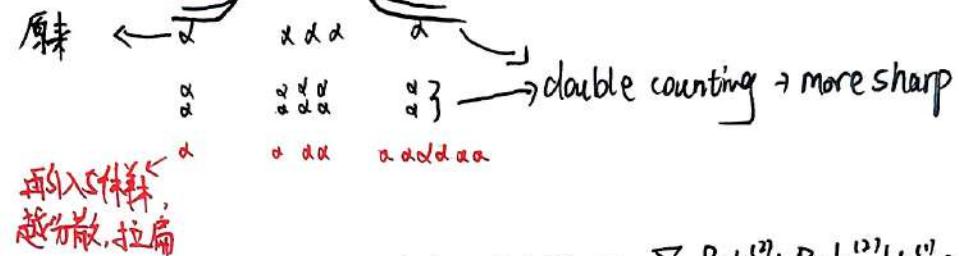
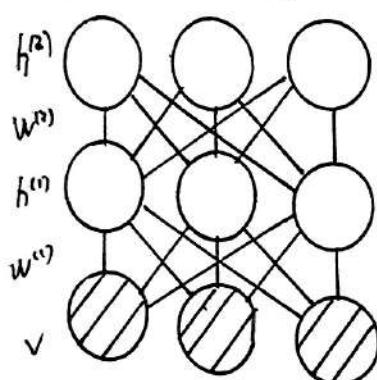
$$E_{P(v)} [P(h^{(1)}|v; w^{(1)})] \stackrel{\text{Mc sampling}}{\approx} \frac{1}{N} \sum_{v \in V} P(h^{(1)}|v; w^{(1)})$$

$\therefore P(h^{(1)}; w^{(1)}) \approx \frac{1}{N} \sum_{v \in V} P(h^{(1)}|v; w^{(1)})$; $P(h^{(1)}; w^{(2)}) \approx \frac{1}{N} \sum_{h^{(2)} \in H} P(h^{(1)}|h^{(2)}; w^{(2)})$, 虽简单相加

且 v 依赖于 V , $h^{(1)}$ 也依赖于 V , 由此可知, 若层数增加, $h^{(i)}$ 仍依赖于 V , 会导致重复计算

V 和 H 都从 V 来的, double counting, 不意图.

$V, w^{(1)}, h^{(1)}, h^{(2)}, w^{(2)}$ 合并, P78 书上图类似:



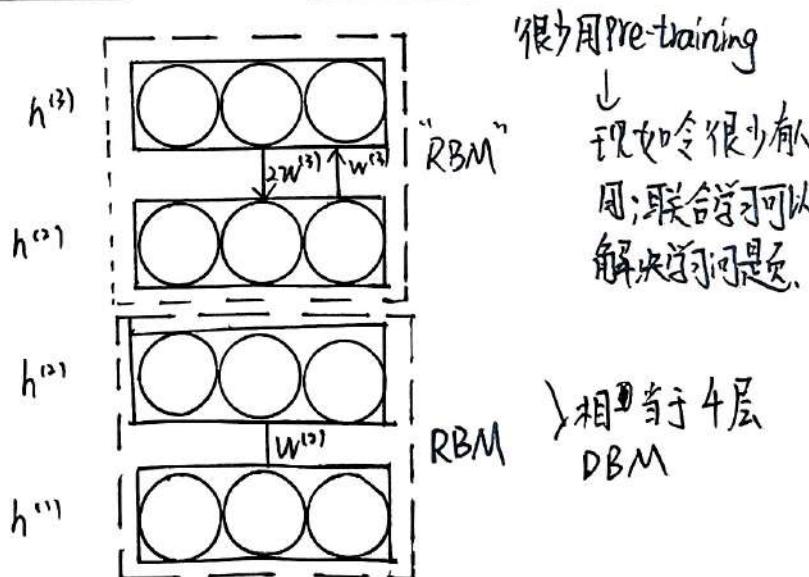
(?) $P(h^{(1)}; w^{(1)})$ 的表示存在问题 $\sum_{h^{(2)}} P(h^{(2)}) P(h^{(1)}|h^{(2)}; w^{(1)})$ 不合理. *, 应为 $\sum_{h^{(2)}} P(h^{(2)}) \cdot P(h^{(1)}|h^{(2)}; w^{(1)})$ 但 $P(h^{(1)}|h^{(2)}; w^{(1)})$ 并未明确说明如何计算

28.4. Pre-training Continue

True: $P(h^{(1)}; w^{(1)}, w^{(2)})$

$$P(h^{(1)}; w^{(1)}) = \sum_v P(v) \cdot P(h^{(1)}|v; w^{(1)}) \rightarrow \text{DBM}$$

$$P(h^{(1)}; w^{(2)}) = \sum_{h^{(2)}} P(h^{(2)}) \cdot P(h^{(1)}|h^{(2)}; w^{(2)}) \rightarrow \text{DBN}$$



29. Generative Model

29.1 Supervised learning vs unsupervised learning

target (二者对比) (分类, 回归, 标记, 降维, 聚类, 特征学习, 密度估计, 生成数据)

概率模型 判别模型 ($P(Y|X)$): LR, MEMM, CRF

监督 概率模型 生成模型

非概率模型: PLA, SVM, KNN, NN, Tree Model

非监督 概率模型 生成模型

非概率模型: PCA, LSA, k-means, Auto-encoder

深度生成模型 (加入Deep Learning)

Energy-based Model: Boltzmann Machine

VAE
GAN (GSN)

Auto-regressive Model

Flow-based Model

传统生成模型:

Naive Bayes

Mixture Model: GMM

Time-Series Model:

HMM | Kalman Filter | Particle Filter

non-parameteric:

Bayesian Model: Gauss Process
Dirichlet process

Mixed membership model: LDA

Factorial Model: Factorial Analysis
P-PCA, ICA,

生成模型: 关注样本分布本身

PCA → P-PCA → FA

k-means → GMM

Auto-encoder → VAE

LSA → PLSA → LDA

29.2 Presentation & Inference & Learning

presentation: 形神兼备: {
 形: Discrete vs Continuous
 Directed vs Undirected Model
 Latent Variable Model vs Fully-observed Model
 Shallow vs Deep (层)
 Sparse vs Dense (连接)
 HMM vs BM
 神: Parameteric Model vs Non-parameteric
Implicit Density vs Explicit Density GAN
 不基于 $p(A)$ 分布对称 基于 $p(A)$ 分布对称

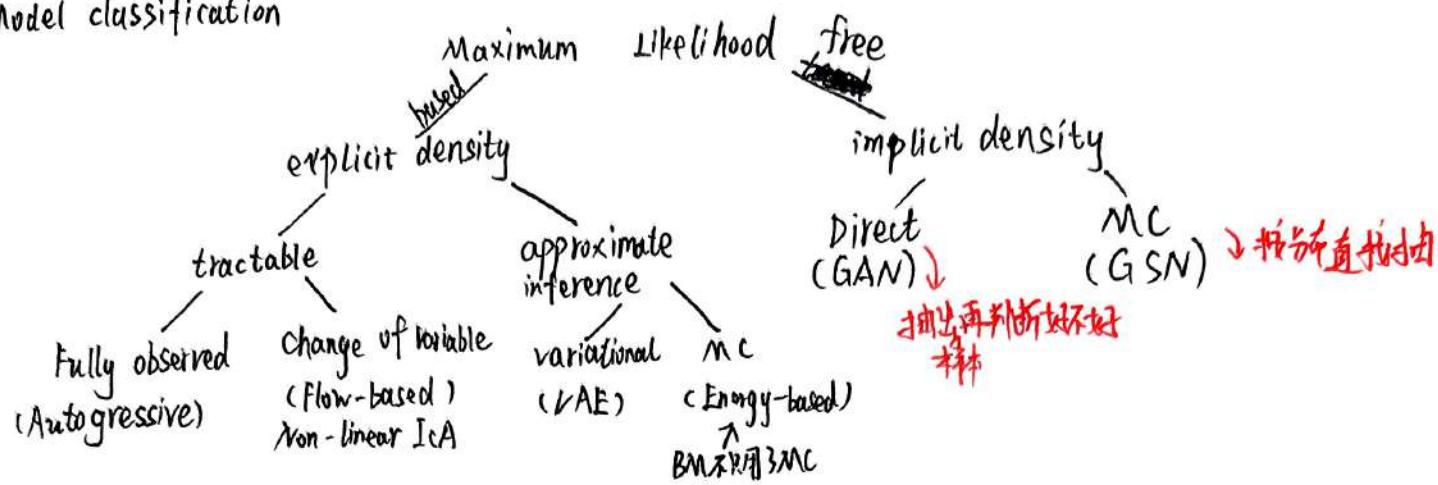
Inference:

tractable v.s. intractable

Learning:

Likelihood-based Model vs. Likelihood-free Model GAN

29.3 Model classification



29.4 probabilistic graph vs Neural network (Bayesian Network vs Neural Network)

probabilistic Graph $\Rightarrow P(x)$ 的表示
Neural Network \Rightarrow 函数逼近器

(以连接主义)

Boltzmann Machine
无向图模型
神经网络

神经网络
确定性
随机性

PG
森：结构化的。
稀疏，浅层。
条件独立假设
具备可解释性

NN
深层、稠密
计算图
可解释性未知

Inference: 精确近似(MC, 变分)
Learning: likelihood maximum estimate
适配: MLE
适合: high level reasoning

NN
推断很容易，但没意义(已知 w 后)
梯度下降 (BP)
表示学习, low level reasoning

链式求导法则
动态规划
递归+缓存

29.5 Reparametrization Trick

假设: $p_{\theta}(y) = N(\mu, \sigma^2)$ 实际上可能是一个很复杂的分布

$$z \sim N(0, 1), \quad \theta = \{ \mu, \sigma^2 \}$$

$$y = \mu + \sigma z$$

$z^{(i)} \sim N(0, 1)$ 具有随机性

$y^{(i)} = \mu + \sigma z^{(i)}$ 随机性转移到 $y^{(i)}$ 上

$$y = f(\mu, \sigma, z)$$

$$z \rightarrow \boxed{\text{NN}} \xrightarrow{\mu, \sigma^2} y$$

$J(y)$: 目标函数

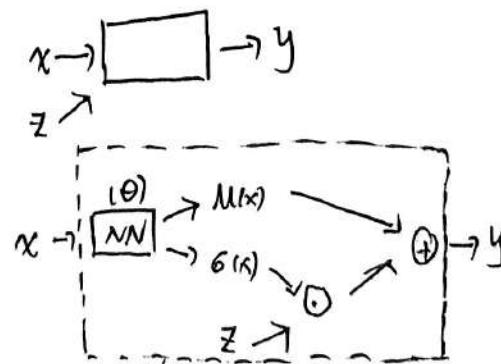
$$\frac{\nabla J(y)}{\nabla \theta} = \frac{\nabla J(y)}{\nabla y} \cdot \frac{\nabla y}{\nabla \theta}$$

② 条件概率分布

$$p(y|x) = N(x; \mu, \sigma^2)$$

$$z \sim N(0, 1)$$

$$y = \mu(x) + \sigma(x) \cdot z$$



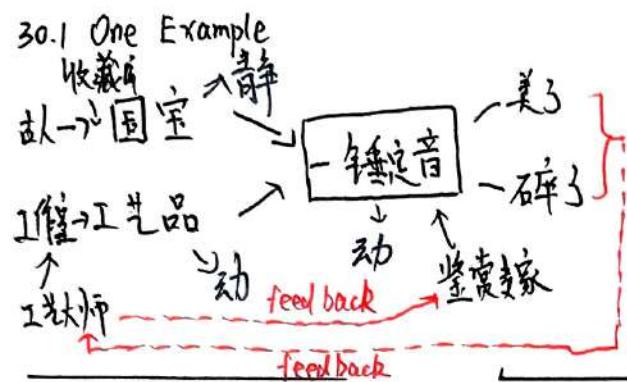
$$J_{\theta}(y) = \sum_{i=1}^N \|y - y^{(i)}\|^2$$

$$\frac{\nabla J_{\theta}(y)}{\nabla \theta} = \frac{\nabla J_{\theta}(y)}{\nabla y} \frac{\nabla y}{\nabla \mu} \frac{\nabla \mu}{\nabla \theta} + \frac{\nabla J_{\theta}(y)}{\nabla y} \frac{\nabla y}{\nabla \sigma} \frac{\nabla \sigma}{\nabla \theta}$$

30. Generative Adversarial Network → Implicit Density Model

目标：成为高水平，可以以假乱真的大师。

双赢 $\left\{ \begin{array}{l} \text{高水平的鉴赏专家} \rightarrow \text{手段} \\ \text{高水平的工艺品大师} \rightarrow \text{目标} \end{array} \right.$
 ↳ 有高专家是成为高大师的先决条件



30.2 Model representation

国宝： $\{x_i\}_{i=1}^N$: P_{data}

工艺品： $P_g(x; \theta_g)$: 假良 $\mathbb{E}^{P_z(z)}$

$\mathbb{E}^{P_z(z)}$ $\xrightarrow{\text{generator}}$ $\gamma = G(z; \theta_g)$ 找一个容易的分布
 $G(z; \theta_g)$ 遵循 $P_g(x; \theta_g)$, 用可微生成网络 $G(z; \theta_g)$ 去逼近 P_g , 从采样角度 \rightarrow Implicit Density

鉴赏家： $D(x; \theta_d)$: 代表 x 是真品(国宝)的概率

$P_{\text{data}} \xrightarrow{\text{NN}} \text{NN} \rightarrow Y/N$

$\log D(x) \uparrow$

高专家：
 1 if x is from P_{data} , then $D(x) \uparrow$ $\log(1 - D(G(z; \theta_g))) \uparrow$
 1 if x is from P_g , then $D(x) \downarrow$ $/ 1 - D(x) \uparrow = \log(1 - D(G(z; \theta_g))) \downarrow$
 $(z \text{ is from } P_z(z))$

$$\max_D \mathbb{E}_{x \sim P_{\text{data}}} [\log D(x)] + \mathbb{E}_{z \sim P_z} [\log (1 - D(G(z; \theta_g)))]$$

高大师：if x is from P_g , then $D(G(z; \theta_g)) \uparrow$ / $\log [1 - D(G(z; \theta_g))] \downarrow$

$$\min_G \mathbb{E}_{z \sim P_z} [\log (1 - D(G(z; \theta_g)))]$$

总目标： $\min_G \max_D \mathbb{E}_{x \sim P_{\text{data}}} [\log D(x)] + \mathbb{E}_{z \sim P_z} [\log (1 - D(G(z; \theta_g)))]$ $P_g \rightarrow P_{\text{data}}$

恒等？为什么

30.3 Global optimality



$$P_{\text{data}}(x) = \{x_i\}_{i=1}^N$$

$P_g(x; \theta_g)$: generator $(P_z(z) + G(z; \theta_g))$

y/x : discriminator $\frac{y/x}{P(D(x))} \frac{1}{1-D(x)}$

$$\text{if } V(D, G) = \mathbb{E}_{x \sim P_{\text{data}}} [\log D(x)] + \mathbb{E}_{z \sim P_z} [\log (1 - D(G(z; \theta_g)))]$$

For fixed G , fix $\max_D V(D, G)$

$$\max_D V(D, G) = \max_D \int P_{\text{data}} \cdot \log D dx + \int P_g \cdot \log (1 - D) dx$$

$$= \max_D \int [P_{\text{data}} \cdot \log D + P_g \cdot \log (1 - D)] dx$$

$$\frac{\partial}{\partial D} (\max_D V(D, G)) = \int \frac{\partial}{\partial D} [P_{\text{data}} \cdot \log D + P_g \cdot \log (1 - D)] dx$$

$$= \int [P_{\text{data}} \cdot \frac{1}{D} + P_g \cdot \frac{-1}{1-D}] dx \triangleq 0$$

$$\Rightarrow \frac{\partial}{\partial D} G = \frac{P_g}{P_{\text{data}} + P_g}$$

将 D_G^* 代入 $V(D, G)$, 则有:

$$\begin{aligned} \min_{G} \max_D V(D, G) &= \min_G V(D_G^*, G) = \min_G E_{\text{unp}_{\text{data}}} \left[\log \frac{P_d}{P_d + P_g} \right] + E_{\text{unp}_g} \left[\log \left(1 - \frac{P_d}{P_d + P_g} \right) \right] \\ &= \min_G E_{\text{unp}_{\text{data}}} \left[\log \frac{P_d}{P_d + P_g} \right] + E_{\text{unp}_g} \left[\log \frac{P_g}{P_d + P_g} \right] \rightarrow \text{与 KL 散度很像} \\ &= \min_G E_{\text{unp}_d} \left[\log \frac{P_d}{(P_d + P_g)/2} \cdot \frac{1}{2} \right] + E_{\text{unp}_g} \left[\log \frac{P_g}{(P_d + P_g)/2} \cdot \frac{1}{2} \right] \\ &= \min_G \text{KL}(P_d || \frac{P_d + P_g}{2}) + \text{KL}(P_g || \frac{P_d + P_g}{2}) - \log 4 \end{aligned}$$

G

$Z - \log 4$

当 $P_d = \frac{P_d + P_g}{2} = P_g$ 时, “=” 成立

故 $D_G^* = P_g$, $D_G^* = \frac{1}{2}$

$D_G^* = \frac{1}{2}$, 专家已经无法一锤定音了

KL 散度很像

$$KL(Q||P) = Q(x) \log \frac{Q(x)}{P(x)} = P(x) [\log \frac{Q(x)}{P(x)}]$$

但 $P_d + P_g$ 两个概率密度函数相加属于 $[0, 2]$ 了, 不是概率密度函数
稍加处理

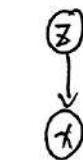
$$E(\log \frac{1}{2}) = -\log 2$$

KL 散度非负

31. Variational Autoencoder

31.1 representation

Latent Variable Model



GMM: k Gaussian Dist 混合: $Z \sim \text{Categorical Dist}$ 离散 $\frac{z}{p_1, p_2, \dots, p_k} \quad \sum_i p_i = 1 \quad x|z_i \sim N(\mu_i, \Sigma_i)$

VAE: infinite Gaussian Dist: $\begin{cases} Z \sim N(0, I) \\ x|z \sim N(\mu_\theta(z), \Sigma_\theta(z)) \end{cases}$

$$p_\theta(x) = \int_z p_\theta(x, z) dz = \int_z p_\theta(z) \cdot p_\theta(x|z) dz$$

intractable
 $p_\theta(z|x) = \frac{p_\theta(z)p_\theta(x|z)}{p_\theta(x)}$
intractable 无法直接求出 $p_\theta(z|x)$

31.2 Learning

Latent Variable Model



$$P_\theta(x|z) = N(\mu_\theta(z), \Sigma_\theta(z))$$

$P_\theta(z|x)$ is intractable 用 $q_\phi(z|x)$ 逼近

$$\log p_\theta(x) = \text{ELBO} + \text{KL}(q_\phi(z|x) || P_\theta(z|x))$$

EM: E-step: 当 $q = P_\theta(z|x)$ 时, $\text{KL} = 0$, $\arg\max_\theta P_\theta(x) = \arg\max_\theta \text{ELBO}$
expectation is ELBO

M-step: $\hat{\theta} = \arg\max_\theta \text{ELBO} = \arg\max_\theta E_{q_\phi(z|x)} [\log p_\theta(x, z)]$

$$\langle \hat{\theta}, \hat{\phi} \rangle = \arg\min \text{KL}(q_\phi(z|x) || P_\theta(z|x))$$

= argmax ELBO

$$= \arg\max_{q_\phi(z|x)} E_{q_\phi(z|x)} [\log P_\theta(x, z)] + H[q_\phi]$$

$$\langle \hat{\theta}, \hat{\phi} \rangle = \arg \max_{\theta, \phi} E_{q_{\phi}(z|x)} [\log P_{\theta}(x|z)] - KL(q_{\phi}(z|x) || P_{\theta}(z))$$

采用 SGVI / SGBV / SVI / Amortized Inference, P36

$z|x$

$$e \sim N(0, I)$$

$$z|x \sim N(\mu_{\phi}(x), \Sigma_{\phi}(x))$$

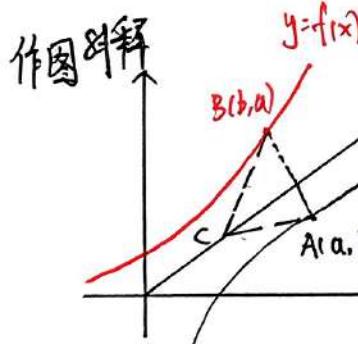
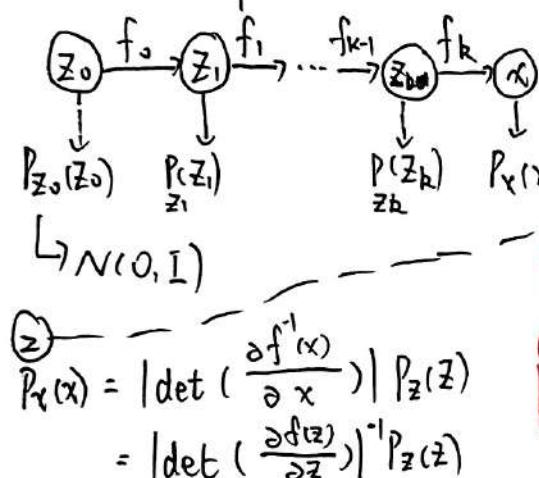
$$\downarrow z = \mu_{\phi}(x) + \sum_{i=1}^d \epsilon_i \cdot \Sigma_{\phi}(x)^{-1}$$

$$x \rightarrow \text{NN} \rightarrow \mu \quad x \rightarrow \text{NN} \rightarrow \Sigma \quad \mu + \Sigma \epsilon \rightarrow z$$

先从 $Z^{(1)}$ ~ $q_{\phi}(z|x)$ 采样
再圆过 $Z^{(1)}$, 计算 $P_{\theta}(x|z^{(1)})$ 再 decoder
 $(\mu, \Sigma) \rightarrow (\mu, \Sigma)$
 $\mu + \Sigma \epsilon \rightarrow \mu + \Sigma^{\frac{1}{2}} \epsilon$

32. Normalizing Flow

32.1 Model Representation



$A : (a, b)$ 过 A 点的导数:

$$B : (b, a) \quad f'(a) = \tan \theta = \frac{b-a}{a-c}$$

$C : (c, d)$ 过 B 点的导数:

$$f'(b) = \tan \theta' = \frac{d-c}{b-c}$$

用 MLE 学习, $P_x(x) = |\det(\frac{\partial f(x)}{\partial z})|^{-1} P_z(z)$

在本题: $P_x(x) = |\det(\frac{\partial f(z_k)}{\partial z_k})|^{-1} P_{z_k}(z_k)$

再用 change of variables Theorem 进行代换即可得一连串, 即每一个节点都利用 change of variables Theorem
代 $P_{z_k}(z_k)$ 与 $\frac{P_{z_{k-1}}(z_{k-1})}{P_{z_k}(z_k)}$

Change of Variables Theorem

① Assuming $x = f(z)$, $z, x \in \mathbb{R}^p$

$$z \sim P_z(z), x \sim P_x(x)$$

f is continuous, invertible 可逆

$$\therefore \int_z P_z(z) dz = 1 = \int_x P_x(x) dx$$

$$P_x(x), P_z(z) > 0$$

$$|\int_z P_z(z) dz| = |\int_x P_x(x) dx|$$

$$\therefore P_x(x) = \frac{|dz|}{|dx|} \cdot P_z(z)$$

$\because x = f(z)$, f is invertible

$$\text{①-②} \frac{dx}{dz} \cdot \frac{dz}{dx} = 1 \quad \therefore z = f^{-1}(x)$$

$$\therefore P_x(x) = \left| \frac{\partial f^{-1}(x)}{\partial x} \right| \cdot P_z(z) \quad -1 \text{ 位}$$

$$\therefore P_x(x) = \left| \frac{\nabla f^{-1}(x)}{\nabla x} \right| \cdot P_z(z) \quad -2 \text{ 位}$$

$\frac{\nabla f^{-1}(x)}{\nabla x}$: Jacobian Matrix

$$\det \left(\frac{\partial f^{-1}(x)}{\partial x} \right) : \text{Jacobian Determinant}$$

$$\therefore P_x(x) = \left| \det \left(\frac{\partial f^{-1}(x)}{\partial x} \right) \right| \cdot P_z(z)$$