

机器学习概论习题课

Yuanhao Pu & Jin Zhang

2024 年 1 月 11 日

目录

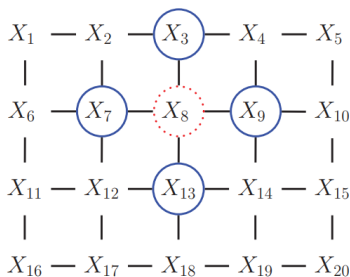
1 基于平均场理论的变分推断

- Example: Mean Field Theory for Ising Model
- Example: VB for linear regression
- Example: Boltzmann Machine

2 作业题

- 1.3
- 3.1
- 5.4
- 9.1
- 10.1

Ising Model



统计物理建模磁性材料中原子的行为。
 $X_i \in \{-1, +1\} \rightarrow$ 表示原子自旋方向；
 令 $s \sim t$ 表示 s 与 t 互为邻居，则（未归一化的）先验分布为

$$\begin{aligned} \log \tilde{p}(\mathbf{x}) &= - \sum_{s \sim t} w_{st} x_s x_t - \sum_s b_s x_s \\ &= -\frac{1}{2} \mathbf{x}^\top \mathbf{W} \mathbf{x} - \mathbf{b}^\top \mathbf{x} (\text{有时忽略 bias}) \end{aligned}$$

Image Denoising

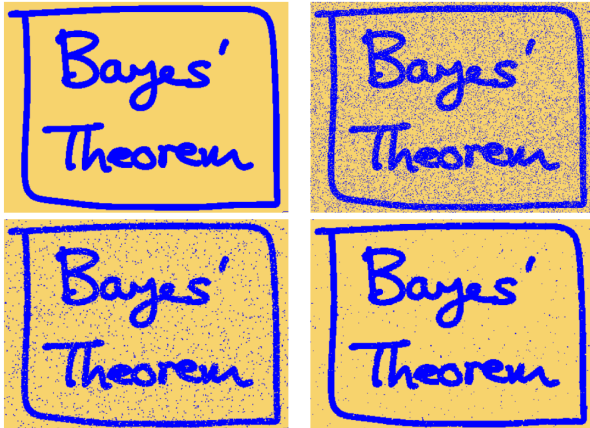
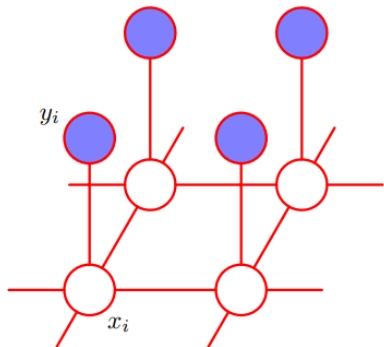


Figure 8.30 Illustration of image de-noising using a Markov random field. The top row shows the original binary image on the left and the corrupted image after randomly changing 10% of the pixels on the right. The bottom row shows the restored images obtained using iterated conditional models (ICM) on the left and using the graph-cut algorithm on the right. ICM produces an image where 96% of the pixels agree with the original image, whereas the corresponding number for graph-cut is 99%.

Ising Model for Image Denoising



$y_i \in \{-1, +1\}$ 含噪声图中的像素点
 $x_i \in \{-1, +1\}$ 无噪声图中的像素点
建模联合分布

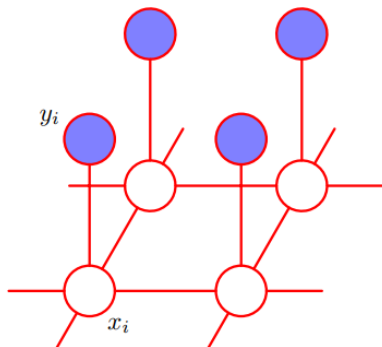
$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y}|\mathbf{x})$$

其中先验 $p(\mathbf{x})$ 有

$$p(\mathbf{x}) = \frac{1}{Z_0} \exp(-E_0(\mathbf{x}))$$

$$E_0(\mathbf{x}) = -\frac{1}{2} \mathbf{x}^\top \mathbf{W} \mathbf{x} = -\sum_{i=1}^D \sum_{j \in \text{nbr}_i} w_{ij} x_i x_j$$

Ising Model for Image Denoising



设 L_i 为 $x \sim y$ 间的能量函数,

$$p(\mathbf{y}|\mathbf{x}) = \prod_i p(y_i|x_i) \propto \exp(-\sum_i L_i(x_i))$$

则联合分布

$$p(\mathbf{x}, \mathbf{y}) = \frac{1}{Z} \exp(-E(\mathbf{x}))$$

$$E(\mathbf{x}) = E_0(\mathbf{x}) - \sum_i L_i(x_i)$$

Mean field

由平均场假设，考虑

$$q(\mathbf{x}) = \prod_i q(x_i, \mu_i)$$

为了得到 μ_i 的更新策略，展开 $\log \tilde{p}(\mathbf{x})$ 并提取出与 x_i 相关的项：

$$\log \tilde{p}(\mathbf{x}) = x_i \sum_{j \in \text{nbr}_i} w_{ij} x_j + L_i(x_i) + \text{const}$$

基于课上的结论，对所有的 $j \neq i$ 求期望，

$$q_i(x_i) \propto \exp(x_i \sum_{j \in \text{nbr}_i} w_{ij} \mu_j + L_i(x_i))$$

进一步的，令

$$m_i = \sum_{j \in \text{nbr}_i} w_{ij} \mu_j, \quad L_i^+ = L_i(+1), \quad L_i^- = L_i(-1)$$

则边缘后验可以直接求得

$$\begin{aligned} q_i(x_i = 1) &= \frac{\exp(1 * m_i + L_i^+)}{\exp(1 * m_i + L_i^+) + \exp(-1 * m_i + L_i^-)} \\ &= \frac{1}{1 + \exp(-2m_i + L_i^- - L_i^+)} = \sigma(2a_i) \end{aligned}$$

其中 $a_i = m_i + 0.5(L_i^+ - L_i^-)$ ，则参数 μ_i 有

$$\mu_i = \mathbb{E}_{q_i}[x_i] = q_i(x_i = +1) - q_i(x_i = -1) = \tanh(a_i)$$

$$\text{Update: } \mu_i^t = \tanh\left(\sum_{j \in \text{nbr}_i} w_{ij} \mu_j^{t-1} + 0.5(L_i^+ - L_i^-)\right)$$

Bayesian Linear Regression

对线性模型 $y = \mathbf{w}^\top \phi(\mathbf{x}) + \varepsilon, \varepsilon \sim \mathcal{N}(0, \lambda^{-1})$

写作分布形式 $p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \lambda) = \prod_{n=1}^N \mathcal{N}(y_n | \mathbf{w}^\top \phi(\mathbf{x}_n), \lambda^{-1})$

Bayesian: 引入关于 \mathbf{w} 的先验分布 $p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0)$

则由 $p(\mathbf{w}|\mathbf{y}) \propto p(\mathbf{y}|\mathbf{w})p(\mathbf{w})$, 得到后验

$$p(\mathbf{w}|\mathbf{t}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N)$$

$$\mathbf{m}_N = \mathbf{S}_N(\mathbf{S}_0^{-1} \mathbf{m}_0 + \lambda \Phi^\top \mathbf{y})$$

$$\mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \lambda \Phi^\top \Phi$$

(过程详见 PRML)

Bayesian Linear Regression

实际操作时考虑简化问题, $p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I})$, 则

$$\begin{aligned} \mathbf{m}_N &= \lambda S_N \Phi^\top \mathbf{y} \\ \mathbf{S}_N^{-1} &= \alpha \mathbf{I} + \lambda \Phi^\top \Phi \end{aligned}$$

代入对数似然, 并提取出相关项

$$\log p(\mathbf{w}|\mathbf{y}) = -\frac{\lambda}{2} \sum_{n=1}^N (y_n - \mathbf{w}^\top \phi(\mathbf{x}_n))^2 - \frac{\alpha}{2} \mathbf{w}^\top \mathbf{w} + \text{const}$$

能证明, 最大化梯度似然与最小化岭回归损失等价 (正则项系数 $\tau = \frac{\alpha}{\lambda}$)

Variational Bayesian for Linear Regression

若假定 $\eta = (\lambda, \alpha)$ 均未知, 选取超参数使边缘似然达到最大的方式称为 evidence procedure;
利用变分 Bayes 估计参数, 取先验

$$p(\mathbf{W}, \lambda, \alpha) = \mathcal{N}(\mathbf{w} | \mathbf{0}, (\lambda \alpha)^{-1} \mathbf{I}) Ga(\lambda | a_N^\lambda, b_N^\lambda) Ga(\alpha | a_N^\alpha, b_N^\alpha)$$

假设变分分布满足

$$q(\mathbf{w}, \lambda, \alpha) = q(\mathbf{w}, \lambda) q(\alpha)$$

Variational Bayesian for Linear Regression

$$q(\mathbf{w}, \alpha, \lambda) = \mathcal{N}(\mathbf{w} | \mathbf{w}_N, \lambda^{-1} \mathbf{V}_N) \text{Ga}(\lambda | a_N^\lambda, b_N^\lambda) \text{Ga}(\alpha | a_N^\alpha, b_N^\alpha)$$

where

$$\mathbf{V}_N^{-1} = \overline{\mathbf{A}} + \mathbf{X}^T \mathbf{X}$$

$$\mathbf{w}_N = \mathbf{V}_N \mathbf{X}^T \mathbf{y}$$

$$a_N^\lambda = a_0^\lambda + \frac{N}{2}$$

$$b_N^\lambda = b_0^\lambda + \frac{1}{2} (\|\mathbf{y} - \mathbf{X} \mathbf{w}\|^2 + \mathbf{w}_N^T \overline{\mathbf{A}} \mathbf{w}_N)$$

$$a_N^\alpha = a_0^\alpha + \frac{D}{2}$$

$$b_N^\alpha = b_0^\alpha + \frac{1}{2} \left(\frac{a_N^\lambda}{b_N^\lambda} \mathbf{w}_N^T \mathbf{w}_N + \text{tr}(\mathbf{V}_N) \right)$$

$$\overline{\mathbf{A}} = \langle \alpha \rangle \mathbf{I} = \frac{a_N^\alpha}{b_N^\alpha} \mathbf{I}$$

Variational Bayesian for Linear Regression

Update: 交替更新 $q(\mathbf{w}, \lambda)$ 和 $q(\alpha)$

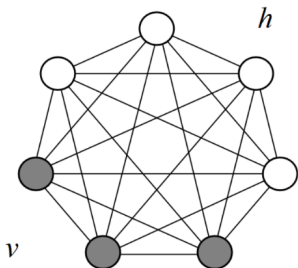
$$p(\mathcal{D}) = \int \int \int p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \lambda) p(\mathbf{w}|\alpha) p(\lambda) d\mathbf{w} d\alpha d\lambda$$

We can compute a lower bound on $\log p(\mathcal{D})$ as follows:

$$\begin{aligned} L(q) = & -\frac{N}{2} \log(2\pi) - \frac{1}{2} \sum_{i=1}^N \left(\frac{a_N^\lambda}{b_N^\lambda} (y_i - \mathbf{w}_N^T \mathbf{x}_i)^2 + \mathbf{x}_i^T \mathbf{V}_N \mathbf{x}_i \right) \\ & + \frac{1}{2} \log |\mathbf{V}_N| + \frac{D}{2} \\ & - \log \Gamma(a_0^\lambda) + a_0^\lambda \log b_0^\lambda - b_0^\lambda \frac{a_N^\lambda}{b_N^\lambda} + \log \Gamma(a_N^\lambda) - a_N^\lambda \log b_N^\lambda + a_N^\lambda \\ & - \log \Gamma(a_0^\alpha) + a_0^\alpha \log b_0^\alpha + \log \Gamma(a_N^\alpha) - a_N^\alpha \log b_N^\alpha \end{aligned}$$

推导过于复杂，重点掌握使用变分推断的流程，并理解 ELBO 为边缘似然的下界（参考 github 上 HW11&HW12.pdf 往年作业题）

Boltzmann Machine



可观测变量 $v_i \in \{0, 1\}$

隐变量 $h_i \in \{0, 1\}$

能量函数

$$\begin{aligned} E(\mathbf{v}, \mathbf{h}) &= -\left(\sum_{i < j} w_{ij} x_i x_j + \sum_i b_i x_i\right) \\ &= -\left(\mathbf{v}^\top W \mathbf{h} + \frac{1}{2} \mathbf{v}^\top L \mathbf{v} + \frac{1}{2} \mathbf{h}^\top J \mathbf{h} + \mathbf{b}_v^\top \mathbf{v} + \mathbf{b}_h^\top \mathbf{h}\right) \\ P(\mathbf{v}, \mathbf{h}) &= \frac{1}{Z} \exp(-E(\mathbf{v}, \mathbf{h})) \end{aligned}$$

图 1: Boltzmann machine 模型的概率图

VI for BM

变分推断通过 ELBO 优化变分分布 $Q_\phi(\mathbf{h}|\mathbf{v})$ 推断 $P(\mathbf{h}|\mathbf{v})$

$$\begin{aligned} ELBO = \mathcal{L} &= \sum_{\mathbf{h}} Q_\phi(\mathbf{h}|\mathbf{v}) \log P_\theta(\mathbf{v}, \mathbf{h}) + H(Q_\phi) \\ &= \sum_{\mathbf{h}} Q_\phi(\mathbf{h}|\mathbf{v}) [-\log Z + \mathbf{v}^\top W \mathbf{h} + \frac{1}{2} \mathbf{v}^\top L \mathbf{v} + \frac{1}{2} \mathbf{h}^\top J \mathbf{h} + \mathbf{b}_v^\top \mathbf{v} + \mathbf{b}_h^\top \mathbf{h}] + H(Q_\phi) \end{aligned}$$

基于平均场理论, 拆分 $Q_\phi(\mathbf{h}|\mathbf{v}) = \prod_{j=1}^P Q_\phi(h_j|\mathbf{v})$, h_j 是二值的, 令 $Q_\phi(h_j = 1|\mathbf{v}) = \phi_j$, 只考虑 \mathcal{L} 中与 ϕ_j (即 h_j) 相关的项

$$\begin{aligned} \mathcal{L} &= \sum_{\mathbf{h}} Q_\phi(\mathbf{h}|\mathbf{v}) [\mathbf{v}^\top W \mathbf{h} + \frac{1}{2} \mathbf{h}^\top J \mathbf{h} + \mathbf{b}_h^\top \mathbf{h}] + H(Q_\phi) \\ &= \sum_{\mathbf{h}} Q_\phi(\mathbf{h}|\mathbf{v}) \mathbf{v}^\top W \mathbf{h} + \sum_{\mathbf{h}} Q_\phi(\mathbf{h}|\mathbf{v}) \frac{1}{2} \mathbf{h}^\top J \mathbf{h} + \sum_{\mathbf{h}} Q_\phi(\mathbf{h}|\mathbf{v}) \mathbf{b}_h^\top \mathbf{h} + H(Q_\phi) \end{aligned}$$

VI for BM

$$\mathcal{L} = \textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4}$$

$$\textcircled{1} = \sum_i \sum_j \phi_j v_i w_{ij} \quad \textcircled{2} = \sum_j \sum_{m \neq j} \phi_j \phi_m J_{jm} \quad \textcircled{3} = \sum_j \phi_j b_{hj}$$

$$\textcircled{4} = - \sum_j [\phi_j \log \phi_j + (1 - \phi_j) \log(1 - \phi_j)]$$

$$\frac{\partial \mathcal{L}}{\partial \phi_j} = 0 \implies \phi_j = \sigma\left(\sum_i v_i w_{ij} + \sum_{m \neq j} \phi_m J_{jm} + b_{hj}\right)$$

迭代求解所有的 ϕ_j 得到 $Q_\phi(\mathbf{h}|\mathbf{v}) \simeq P(\mathbf{h}|\mathbf{v})$

1.3

已知随机变量 $\mathbf{x} = [\mathbf{x}_1, \mathbf{x}_2] \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, 计算 $P(\mathbf{x}_1), P(\mathbf{x}_1 | \mathbf{x}_2)$
Solve.

$$\mathbf{x}_1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11}).$$

待定 λ 使得 $\mathbf{x}_1 - \lambda \mathbf{x}_2$ 与 \mathbf{x}_2 相互独立, 等价于

$$\text{Cov}(\mathbf{x}_1 - \lambda \mathbf{x}_2, \mathbf{x}_2) = 0 \Rightarrow \lambda = \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{12}$$

$\mathbf{x}_1 | \mathbf{x}_2 = (\mathbf{x}_1 - \lambda \mathbf{x}_2) | \mathbf{x}_2 + \lambda \mathbf{x}_2 | \mathbf{x}_2 = (\mathbf{x}_1 - \lambda \mathbf{x}_2) + \lambda \mathbf{x}_2 | \mathbf{x}_2$ 是一正态分布, 求出均值和标准差即可

$$\mathbb{E}[\mathbf{x}_1 | \mathbf{x}_2] = \boldsymbol{\mu}_1 - \lambda \boldsymbol{\mu}_2 + \lambda \mathbf{x}_2$$

$$\text{Cov}(\mathbf{x}_1 - \lambda \mathbf{x}_2, \mathbf{x}_1 - \lambda \mathbf{x}_2) = \text{Cov}(\mathbf{x}_1 - \lambda \mathbf{x}_2, \mathbf{x}_1) = \boldsymbol{\Sigma}_{11} - \lambda \boldsymbol{\Sigma}_{21}$$

$$\mathbf{x}_1 | \mathbf{x}_2 \sim \mathcal{N}(\boldsymbol{\mu}_1 - \lambda \boldsymbol{\mu}_2 + \lambda \mathbf{x}_2, \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{21})$$

3.1

[课本习题 3.2] 试证明, 对于参数 w , 对率回归的目标函数 (3.18) 是非凸的, 但其对数似然函数 (3.27) 是凸的。

Sol.

$$y = \frac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x} + b}}$$

$$\ell(\beta) = \sum_{i=1}^m \left(-y_i \beta^\top \hat{\mathbf{x}}_i + \ln \left(1 + e^{\beta^\top \hat{\mathbf{x}}_i} \right) \right)$$

$$\frac{\partial y}{\partial \mathbf{w}} = \frac{\mathbf{x} e^{-(\mathbf{w}^\top \mathbf{x} + b)}}{\left(1 + e^{-(\mathbf{w}^\top \mathbf{x} + b)} \right)^2} = \mathbf{x} y (1 - y)$$

$$\frac{\partial^2 y}{\partial \mathbf{w}^\top \partial \mathbf{w}} = \frac{\partial y}{\partial \mathbf{w}^\top} \mathbf{x} (1 - y) + \frac{\partial (1 - y)}{\partial \mathbf{w}^\top} \mathbf{x} y = \mathbf{x} \mathbf{x}^\top y (1 - 2y) (1 - y)$$

3.1

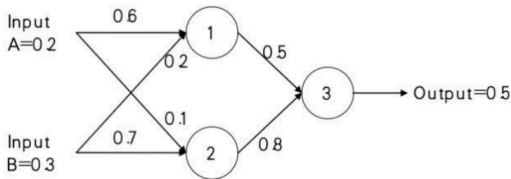
$\mathbf{x}\mathbf{x}^\top \geq 0$ 恒成立, 当 $0.5 < y < 1$ 时, $y(1-2y)(1-y) < 0$, 此时 $\frac{\partial^2 \ell}{\partial \mathbf{w}^\top \partial \mathbf{w}} < 0$, 因此函数 (3.18) 非凸。

$$\begin{aligned}\frac{\partial \ell}{\partial \boldsymbol{\beta}} &= \sum_{i=1}^m \left(-y_i \hat{\mathbf{x}}_i + \frac{1}{1 + \exp \boldsymbol{\beta}^\top \hat{\mathbf{x}}_i} \hat{\mathbf{x}}_i \exp \boldsymbol{\beta}^\top \hat{\mathbf{x}}_i \right) \\ \frac{\partial^2 \ell}{\partial \boldsymbol{\beta}^\top \partial \boldsymbol{\beta}} &= \frac{\partial}{\partial \boldsymbol{\beta}^\top} \sum_{i=1}^m \left(-y_i \hat{\mathbf{x}}_i + \frac{1}{1 + \exp \boldsymbol{\beta}^\top \hat{\mathbf{x}}_i} \hat{\mathbf{x}}_i \exp \left(\boldsymbol{\beta}^\top \hat{\mathbf{x}}_i \right) \right) \\ &= \sum_{i=1}^m \frac{\exp \left(\boldsymbol{\beta}^\top \hat{\mathbf{x}}_i \right)}{\left(1 + \exp \left(\boldsymbol{\beta}^\top \hat{\mathbf{x}}_i \right) \right)^2} \hat{\mathbf{x}}_i \hat{\mathbf{x}}_i^\top\end{aligned}$$

由于 $\hat{\mathbf{x}}_i \hat{\mathbf{x}}_i^\top \geq 0$ 且 $\frac{e^{\boldsymbol{\beta}^\top \hat{\mathbf{x}}_i}}{(1+e^{\boldsymbol{\beta}^\top \hat{\mathbf{x}}_i})^2} \geq 0$, 因此函数 (3.27) 为凸函数。

5.4 神经网络参数更新

激活函数为 ReLU，用平方损失 $\frac{1}{2}(y - \hat{y})^2$ 计算误差，请用 BP 算法更新一次所有参数 (学习率为 1)，给出更新后的参数值给定输入值 $x = (0.2, 0.3)$ 时初始时和更新后的输出值



Sol:

已知 $v_{11} = 0.6, v_{12} = 0.1, v_{21} = 0.2, v_{22} = 0.7, w_1 = 0.5, w_2 = 0.8$
 , 令 $\alpha_1, \alpha_2, \gamma$ 为结点 1, 2, 3 的输入, $\{\beta_1, \beta_2, \hat{y}\}$ 为对应输出,

5.4

正向传播:

$$\alpha_1 = v_{11}A + v_{21}B = 0.18 = \beta_1$$

$$\alpha_2 = v_{12}A + v_{22}B = 0.23 = \beta_2$$

$$\gamma = w_1\beta_1 + w_2\beta_2 = 0.274 = \hat{y}$$

$$E = \frac{1}{2}(y - \hat{y})^2 = 0.025538$$

反向传播 (链式法则):

$$\frac{\partial E}{\partial v_{11}} = \frac{\partial E}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \gamma} \frac{\partial \gamma}{\partial \beta_1} \frac{\partial \beta_1}{\partial \alpha_1} \frac{\partial \alpha_1}{\partial v_{11}} = -0.0226$$

9.1

记 $\text{err}^*(\mathbf{x}) = 1 - \max_{c \in \mathcal{Y}} P(c | \mathbf{x})$, $\text{err}(\mathbf{x}) = 1 - \sum_c P(c | \mathbf{x})P(c | \mathbf{z})$,
其中 \mathbf{z} 为 \mathbf{x} 的最近邻, 试证明在样本无穷多时

$$\text{err}^*(\mathbf{x}) \leq \text{err}(\mathbf{x}) \leq \text{err}^*(\mathbf{x}) \left(2 - \frac{|\mathcal{Y}|}{|\mathcal{Y}| - 1} \times \text{err}^*(\mathbf{x}) \right)$$

证明. 先证明左边不等式:

$$\begin{aligned} \text{err}^*(\mathbf{x}) &= 1 - \max_{c \in \mathcal{Y}} P(c | \mathbf{x}) = 1 - \max_{c \in \mathcal{Y}} P(c | \mathbf{x}) \cdot \sum_c P(c | \mathbf{z}) \\ &= 1 - \sum_c \max_{c \in \mathcal{Y}} P(c | \mathbf{x}) \cdot P(c | \mathbf{z}) \\ &\leq 1 - \sum_c P(c | \mathbf{x}) \cdot P(c | \mathbf{z}) = \text{err}(\mathbf{x}) \end{aligned}$$

9.1

令 $c^* = \arg \max_c P(c | \mathbf{x})$, 再证明右边不等式:

$$\begin{aligned} \text{err}^*(\mathbf{x}) &= 1 - \sum_c P(c | \mathbf{x}) \cdot P(c | \mathbf{z}) \simeq 1 - \sum_c P(c | \mathbf{x})^2 \\ &\leq 1 - P(c^* | \mathbf{x})^2 - \sum_{c \neq c^*} P(c | \mathbf{x})^2 \\ &\leq 1 - P(c^* | \mathbf{x})^2 - \frac{1}{|\mathcal{Y}| - 1} \left(\sum_{c \neq c^*} P(c | \mathbf{x}) \right)^2 \\ &= 1 - P(c^* | \mathbf{x})^2 - \frac{1}{|\mathcal{Y}| - 1} (1 - P(c^* | \mathbf{x}))^2 \end{aligned}$$

9.1

$$\begin{aligned}
 &= (1 - P(c^* | \mathbf{x})) \cdot \left(1 + P(c^* | \mathbf{x}) - \frac{1}{|\mathcal{Y}| - 1} (1 - P(c^* | \mathbf{x})) \right) \\
 &= (1 - P(c^* | \mathbf{x})) \cdot \left(2 - \frac{|\mathcal{Y}|}{|\mathcal{Y}| - 1} (1 - P(c^* | \mathbf{x})) \right) \\
 &= \text{err}^*(\mathbf{x}) \cdot \left(2 - \frac{|\mathcal{Y}|}{|\mathcal{Y}| - 1} \cdot \text{err}^*(\mathbf{x}) \right)
 \end{aligned}$$

10.1

[课本习题 11.5] 结合图 11.2, 试举例说明 L_1 正则化在何种情形下不能产生稀疏解。

解. 如图 1 所示, 当平方误差项等值线的斜率较大的时候, 其与 L_1 范数等值线的交点就不再位于坐标轴上, 因此将无法产生稀疏解。

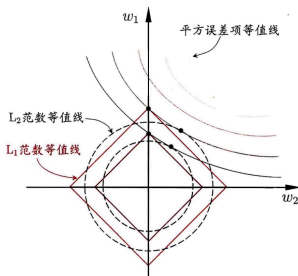


图 11.2 L_1 正则化比 L_2 正则化更易于得到稀疏解

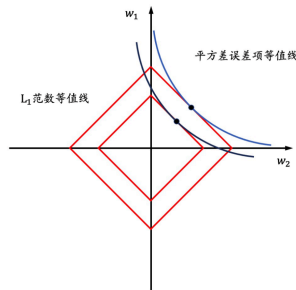


图 1: 情况演示图