

HW 机器学习概论

PB19151769 马宇骁

目录

1	HW1	2
1.1	2
1.2	3
1.3	3
1.4	4
1.5	4
2	HW2	5
2.1	5
2.2	6
2.3	6
2.4	7
3	HW3	8
3.1	8
3.2	9
3.3	9
3.4	9
3.5	10
4	HW4	11
4.1	11
4.2	11
4.3	12
4.4	12
4.4.1	12
4.4.2	12
4.4.3	13
4.4.4	14
4.4.5	14

5	HW5	15
5.1	15
5.2	16
5.3	16
5.4	16
6	HW6	18
6.1	18
6.2	19
6.3	19
6.4	19
7	HW7	21
7.1	21
7.2	21
7.3	22
7.4	22
8	HW8	23
8.1	23
8.2	23
9	HW9	24
9.1	24
9.2	24
9.3	25
10	HW10	25
10.1	26
10.2	26
10.3	27
11	HW11	28
11.1	28
11.2	29
11.3	29
11.3.1	29
11.3.2	29
12	HW12	30
13	HW13	30

14 HW14	30
14.1	31
14.1.1	31
14.1.2	31
14.1.3	31
14.2	32

1 HW1

作业



34

- 1. 计算 $\frac{\partial \ln \det(A)}{\partial x}$
- 2. 书习题1.2
- 3. 已知随机变量 $\mathbf{x} = [\mathbf{x}_1, \mathbf{x}_2] \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, 计算 $P(\mathbf{x}_1), P(\mathbf{x}_1|\mathbf{x}_2)$
- 4. 证明范数 $\|\mathbf{x}\|_p$ 是凸函数
- 5. 证明判定凸函数的0阶和1阶条件相互等价

$$\forall \mathbf{x}, \mathbf{y} \in \text{dom}(f), \forall t \in [0, 1], f(t\mathbf{x} + (1-t)\mathbf{y}) \leq tf(\mathbf{x}) + (1-t)f(\mathbf{y})$$



$$\forall \mathbf{x}, \mathbf{y} \in \text{dom}(f), f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})$$

1.1

先由:

$$\frac{\partial \det(A)}{\partial A_{ij}} = \frac{\partial \sum_k A_{ik} \text{adj}^T(A)_{ik}}{\partial A_{ij}} = \sum_k \frac{\partial A_{ik} \text{adj}^T(A)_{ik}}{\partial A_{ij}} = \sum_k \frac{\partial A_{ik}}{\partial A_{ij}} \text{adj}^T(A)_{ik} + \sum_k A_{ik} \frac{\partial \text{adj}^T(A)_{ik}}{\partial A_{ij}}$$

其中, $\frac{\partial \text{adj}^T(A)_{ik}}{\partial A_{ij}} = 0$, 有:

$$\frac{\partial \det(A)}{\partial A_{ij}} = \sum_k \frac{\partial A_{ik}}{\partial A_{ij}} \text{adj}^T(A)_{ik} = \text{adj}^T(A)_{ij}$$

因此,

$$\begin{aligned} \frac{\partial \ln \det(A)}{\partial x} &= \frac{1}{\det(A)} \sum_{i,j} \frac{\partial \det(A)}{\partial A_{ij}} \frac{\partial A_{ij}}{\partial x} \\ &= \frac{1}{\det(A)} \sum_{i,j} \text{adj}^T(A)_{ij} \frac{\partial A_{ij}}{\partial x} \\ &= \frac{1}{\det(A)} \sum_{i,j} \text{adj}(A)_{ji} \frac{\partial A_{ij}}{\partial x} \\ &= \frac{1}{\det(A)} \sum_{j=1}^n \left(A^* \frac{\partial A}{\partial x} \right)_{jj} \\ &= \text{tr} \left(A^{-1} \frac{\partial A}{\partial x} \right) \end{aligned}$$

1.2

题中的表 1.1 如图1.

编号	色泽	根蒂	敲声	好瓜
1	青绿	蜷缩	浊响	是
2	乌黑	蜷缩	沉闷	是
3	青绿	硬挺	清脆	否
4	乌黑	稍蜷	沉闷	否

图 1: 西瓜数据集

由题可知:

仅从表中发现有 3 个因素, 每个因素有 2, 3, 3 种取值, 假设空间总共有 $3 \times 4 \times 4 + 1 = 49$ 种假设, 其中有 1 种是空集, 不考虑空时:

全部不泛化 $2 \times 3 \times 3 = 18$ 种假设

一个属性泛化: $2 \times 3 + 3 \times 3 + 2 \times 3 = 21$ 种假设

两个属性泛化: $2 + 3 + 3 = 8$ 种假设

三属性泛化: 1 种假设

用这 48 种假设的排列组合来组成析合范式, 展开序列为:

$$1, C_{48}^1, C_{48}^2, \dots, C_{48}^{48}$$

共 49 个数, 左边的 1 代表 ‘空’, 一个都不选, 右边的 1 代表全部选。

如果 $k=48$, $2^{48} - 1$ 种 (排除一种都不选的情况).

如果 $0 < k < 48$, $\sum_{i=0}^k C_{48}^i - 1 = \sum_{i=1}^k C_{48}^i$ 种.

(以上结果均需去重)

若 k 足够大, 并考虑重复, 则演变为在不泛化的 18 种中选择, 共有 $2^{18}-1$ 种假设。

1.3

令

$$\Sigma = \begin{pmatrix} \sigma_1^2 & \text{Cov}(x_1, x_2) \\ \text{Cov}(x_2, x_1) & \sigma_2^2 \end{pmatrix}$$

其中, $\rho = \frac{\text{Cov}(x_1, x_2)}{\sigma_1 \sigma_2}$

由分布可知:

$$f(x_1, x_2) = \frac{1}{2\pi(1-\rho^2)^{\frac{1}{2}}\sigma_1\sigma_2} e^{-\frac{1}{2(1-\rho^2)}\left[\left(\frac{x_1-\mu_1}{\sigma_1}\right)^2 - 2\rho\left(\frac{x_1-\mu_1}{\sigma_1}\right)\left(\frac{x_2-\mu_2}{\sigma_2}\right) + \left(\frac{x_2-\mu_2}{\sigma_2}\right)^2\right]}$$

所以,

$$\begin{aligned}
P(x_1) &= \int_{-\infty}^{x_1} \int_{-\infty}^{+\infty} f(x_1, x_2) \, dx_2 \, dx_1 \\
&= \int_{-\infty}^{x_1} \int_{-\infty}^{+\infty} \frac{1}{2\pi(1-\rho^2)^{\frac{1}{2}}\sigma_1\sigma_2} e^{-\frac{1}{2(1-\rho^2)}\left[\left(\frac{x_1-\mu_1}{\sigma_1}\right)^2 - 2\rho\left(\frac{x_1-\mu_1}{\sigma_1}\right)\left(\frac{x_2-\mu_2}{\sigma_2}\right) + \left(\frac{x_2-\mu_2}{\sigma_2}\right)^2\right]} \, dx_2 \, dx_1 \\
&= \int_{-\infty}^{x_1} f_{x_1}(x_1) \, dx_1 \\
&= \int_{-\infty}^{x_1} \frac{1}{\sigma_1\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu_1}{\sigma_1}\right)^2} \, dx_1
\end{aligned}$$

$$\begin{aligned}
P(x_1|x_2) &= \int_{-\infty}^{x_1} \frac{f(x_1, x_2)}{f_{x_2}(x_2)} \, dx_1 \\
&= \int_{-\infty}^{x_1} \frac{1}{\sqrt{2\pi}\sigma_1\sqrt{1-\rho^2}} e^{-\frac{(w-\rho v)^2}{2(1-\rho^2)}} \, dx_1
\end{aligned}$$

其中,

$$w = \frac{(x_1 - \mu_1)^2}{\sigma_1^2}, \quad v = \frac{(x_2 - \mu_2)^2}{\sigma_2^2}.$$

1.4

p-范数 ($p \geq 1$):

$$f(X) = \|X\|_p = \left(\sum_{i=1}^n X_i^p\right)^{\frac{1}{p}}$$

要证明 f 为凸函数,

及证明 $f(\lambda \mathbf{x} + (1-\lambda)\mathbf{y}) \leq \lambda f(\mathbf{x}) + (1-\lambda)f(\mathbf{y})$ 对任意的 $0 \leq \lambda \leq 1$ 成立

由 Minkowski 不等式可知:

设 S 是一个度量空间, $1 \leq p \leq \infty$, $f, g \in L^p(S)$, 那么 $f+g \in L^p(S)$, 我们有:

$$\|f+g\|_p \leq \|f\|_p + \|g\|_p$$

即,

$$\left(\sum_{i=1}^n (x_i + y_i)^p\right)^{\frac{1}{p}} \leq \left(\sum_{i=1}^n x_i^p\right)^{\frac{1}{p}} + \left(\sum_{i=1}^n y_i^p\right)^{\frac{1}{p}}$$

如果 $1 < p < \infty$, 等号成立当且仅当 $k \geq 0$, $f = kg$ 或 $g = kf$

故,

$$\begin{aligned}
\|\lambda \mathbf{x} + (1-\lambda)\mathbf{y}\|_p &\leq \|\lambda \mathbf{x}\|_p + \|(1-\lambda)\mathbf{y}\|_p \\
&= \lambda \|\mathbf{x}\|_p + (1-\lambda) \|\mathbf{y}\|_p
\end{aligned}$$

凸函数得证.

1.5

- 充分性

令 $a = 1 - t, a \in [0, 1]$,

$$\begin{aligned} f(tx + (1-t)y) &\leq tf(x) + (1-t)f(y) \\ \implies af(y) &\geq f[(1-a)x + ay] - (1-a)f(x) = f[x + a(y-x)] + af(x) - f(x) \\ \implies f(y) &\geq f(x) + \lim_{a \rightarrow 0} \frac{f[x + a(y-x)] - f(x)}{a(y-x)}(y-x) \\ \implies f(y) &\geq f(x) + \nabla f(x)^T(y-x) \end{aligned}$$

• 必要性

令 $z = tx + (1-t)y$

$$f(y) \geq f(z) + \nabla f(z)^T(y-z) \quad (1)$$

$$f(x) \geq f(z) + \nabla f(z)^T(x-z) \quad (2)$$

$(1-t) * (1) + t * (2)$ 可得:

$$\begin{aligned} tf(x) + (1-t)f(y) &\geq f(z) + \nabla f(z)^T(t(1-t)(x-y) + t(1-t)(y-x)) = f(z) \\ \implies f(tx + (1-t)y) &\leq tf(x) + (1-t)f(y) \end{aligned}$$

2 HW2

作业



48

• 习题2.2

• 习题2.4

• 习题2.5

• 习题2.9

2.1

• 10 折交叉验证

分成 10 份每次取 9 份, 其中正反各一半即 45 个样本, 训练结果随机预测则错误率的期望是 0.5.

• 留一法

留出样本为正则训练 50 反 49 正预测反, 反之预测正, 结果一定错, 错误率期望 1.

2.2

- 真正例率 (TPR)

$$TPR = \frac{\text{预测为正例且真实为正例的数量}}{\text{真实为正例的数量}} = \frac{TP}{TP + FN}$$

- 假正例率 (FPR)

$$FPR = \frac{\text{预测为正例且真实为反例的数量}}{\text{真实为反例的数量}} = \frac{FP}{TN + FP}$$

- 查准率 (P)

$$P = \frac{\text{预测为正例且真实为正例的数量}}{\text{预测为正例的数量}} = \frac{TP}{TP + FP}$$

- 查全率 (R)

$$R = \frac{\text{预测为正例且真实为正例的数量}}{\text{真实为正例的数量}} = \frac{TP}{TP + FN}$$

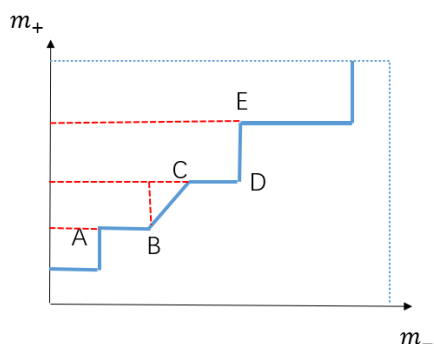
2.3

证明

$$\text{AUC} = 1 - \ell_{\text{rank}}$$

$$\ell_{\text{rank}} = \frac{1}{m^+m^-} \sum_{x^+ \in D^+} \sum_{x^- \in D^-} \left(\mathbb{I}(f(x^+) < f(x^-)) + \frac{1}{2} \mathbb{I}(f(x^+) = f(x^-)) \right)$$

可以将 ROC 曲线分为三部分：平行于 y 轴的；平行于 x 轴的；斜线段。



令 S_i 为 ROC 与 y 轴围成的面积块。则证明题目等式即证明 ℓ_{rank} 是 ROC 曲线与轴围成的全部面积。

- 平行于 y 轴的

此时, $\sum_{x^+ \in D^+} \sum_{x^- \in D^-} \left(\frac{1}{2} \mathbb{I}(f(x^+) = f(x^-)) \right) = 0$,

且 $\sum_{x^+ \in D^+} \sum_{x^- \in D^-} \frac{1}{m^+ m^-} \mathbb{I}(f(x^+) < f(x^-)) = \sum_{i, \text{平行 y 轴}} x_i n_i = \sum_{i, \text{平行 y 轴}} S_i$,
故, 不妨令 ℓ_{rank} 前半部分求和是 A, 后半部分求和是 B (方便打字), 则:

$$\sum_{i, \text{平行 y 轴}} S_i = A + B \text{ 得证;}$$

- 平行于 x 轴的

此时, $A + B = 0$, $S = 0$ 显然, 因此,

$$\sum_{i, \text{平行 x 轴}} S_i = A + B \text{ 得证;}$$

- 斜线段

此时, 面积可以看成开始的点做 y 轴平行的线向上至结束的点的高度的与 y 轴围成的矩形加上斜线段与刚刚做的辅助线围成的三角形的面积之和。

其中, A 为前者矩形的面积, B 是三角形的面积, 故

$$\sum_{i, \text{斜线段}} S_i = A + B \text{ 得证.}$$

综上, $\text{AUC} = 1 - \ell_{\text{rank}}$ 得证。

2.4

χ^2 检验过程:

对于我们要研究的指标, 总体被分为 k 类: $A_i, i = 1, 2, \dots, k$ 。理论上我们知道或者猜想这 k 类数据占总体的比例为:

$H_0: A_i$ 的占比为 $p_i, i = 1, 2, \dots, k$

检验这个原假设, 根据收集的每个分类 A_i 的实际频数 n_i , 结合理论频数, 有:

$$\sum_{i=1}^k \frac{(n_i - np_i)^2}{np_i} \sim \chi^2(k-1)$$

再通过查表比较大小, $\chi^2 < \chi_{\alpha, k-1}^2$ 则没有充分理由拒绝原假设, 认为实际符合理论。

(如果是连续数据则分段然后同上述离散的计算方式)

3 HW3

作业



- 3.2
- 3.7
- 在LDA多分类情形下，试计算类间散度矩阵 S_b 的秩并证明
- 给出公式3.45的推导过程
- 证明 $\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$ 是投影矩阵，并对线性回归模型从投影角度解释

3.1

- (3.18)

$$y = \frac{1}{1 + e^{-(w^T x + b)}}$$

一阶导数:

$$\frac{dy}{dw} = \frac{x e^{-(w^T x + b)}}{(1 + e^{-(w^T x + b)})^2} = x(y - y^2)$$

二阶导数:

$$\frac{d}{dw^T} \left(\frac{dy}{dw} \right) = x(1 - 2y) \left(\frac{dy}{dw} \right)^T = x x^T y(y - 1)(1 - 2y)$$

对于 XPX^T ，有任意向量 Z ，使得:

$$Z^T X P X^T Z = \sum_i P_{ii} V_i^2 \geq 0$$

故， XPX^T 是半正定的。再由 $y \in (0.5, 1)$ 时， $y(y - 1)(1 - 2y) < 0$ ，此时，(3.18) 的 Hessian 矩阵不是半正定的。故，(3.18) 是非凸的。

- (3.27)

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

二阶导数:

$$\frac{d}{d\beta^T} \left(\frac{d\ell}{d\beta} \right) = x x^T p_1(x; \beta)(1 - p_1(x; \beta))$$

其中， $p \in (0, 1)$ ，故 $p_1(x; \beta)(1 - p_1(x; \beta)) \geq 0$ ，(3.27) 的 Hessian 矩阵是半正定的。由此，(3.27) 是凸的。

3.2

当 $3 \leq k \leq 7$ 时, 需要长度 $2^{k-1} - 1$ 的 ECOC 编码。在类别为 4 时, 其可行的编码有 7 种,

	f0	f1	f2	f3	f4	f5	f6
c1	1	1	1	1	1	1	1
c2	0	0	0	0	1	1	1
c3	0	0	1	1	0	0	1
c4	0	1	0	1	0	1	0

当码长为 9 时, 那么之后加任意两个编码, 即为最优编码, 因为此时再加任意的编码都是先有编码的反码, 此时, 类别之间最小的海明距离都为 4, 不会再增加。

3.3

在 LDA 多分类情形下, 试计算类间散度矩阵 $S_b = \sum_{i=1}^N m_i (\mu_i - \mu) (\mu_i - \mu)^T$ 的秩并证明。

$$S_b = [(\mu_1 - \mu), (\mu_2 - \mu) \dots (\mu_N - \mu)] \begin{pmatrix} m_1 & & \\ & \ddots & \\ & & m_N \end{pmatrix} \begin{pmatrix} (\mu_1 - \mu)^T \\ (\mu_2 - \mu)^T \\ \vdots \\ (\mu_N - \mu)^T \end{pmatrix}$$

令 $M = \text{diag}(m_1, m_2, \dots, m_N)$, $A = (\mu_1 - \mu, \mu_2 - \mu, \dots, \mu_N - \mu)^T$, 则:

$$\begin{aligned} \text{rank}(S_b) &= \text{rank}(A^T M A) \\ &= \text{rank}(A^T M^{\frac{1}{2}} M^{\frac{1}{2}} A) \\ &= \text{rank}\left(A^T M^{\frac{1}{2}}\right) \left(A^T M^{\frac{1}{2}}\right)^T \\ &= \text{rank}\left(A^T M^{\frac{1}{2}}\right) \\ &= \text{rank}(A^T) \\ &= \text{rank}(\mu_1 - \mu, \mu_2 - \mu, \dots, \mu_N - \mu) \end{aligned}$$

再由

$$\begin{aligned} \sum_{i=1}^N m_i \mu_i &= \left(\sum_{i=1}^N m_i\right) \mu \\ \sum_{i=1}^N m_i (\mu_i - \mu) &= \mathbf{0} \end{aligned}$$

可知: $\text{rank}(S_b) \leq N - 1$.

3.4

由题得:

若 \mathbf{W} 为 (3.45) 的解。则对于任意常数与其的积也为解。不失一般性，令 $\text{tr}(\mathbf{W}^T \mathbf{S}_w \mathbf{W}) = 1$ ，则 (3.45) 等价于：

$$\begin{aligned} \min_{\mathbf{W}} & -\mathbf{W}^T \mathbf{S}_b \mathbf{W} \\ \text{s.t. } & \mathbf{W}^T \mathbf{S}_w \mathbf{W} = 1. \end{aligned}$$

拉格朗日函数为：

$$L(\mathbf{W}, \lambda) = -\text{tr}(\mathbf{W}^T \mathbf{S}_b \mathbf{W}) + \lambda (\text{tr}(\mathbf{W}^T \mathbf{S}_w \mathbf{W}) - 1)$$

于是有：

$$\begin{aligned} dL &= -\text{tr}((d\mathbf{W})^T \mathbf{S}_b \mathbf{W} + \mathbf{W}^T \mathbf{S}_b d\mathbf{W}) + \lambda \text{tr}((d\mathbf{W})^T \mathbf{S}_w \mathbf{W} + \mathbf{W}^T \mathbf{S}_w d\mathbf{W}) \\ &= -\text{tr}((\mathbf{S}_b \mathbf{W})^T d\mathbf{W}) - \text{tr}(\mathbf{W}^T \mathbf{S}_b d\mathbf{W}) + \lambda (\text{tr}((\mathbf{S}_w \mathbf{W})^T d\mathbf{W}) + \text{tr}(\mathbf{W}^T \mathbf{S}_w d\mathbf{W})) \\ &= \text{tr}([\lambda \mathbf{W}^T (\mathbf{S}_w^T + \mathbf{S}_w) - \mathbf{W}^T (\mathbf{S}_b^T + \mathbf{S}_b)] d\mathbf{W}) \\ &= \text{tr}([\lambda (\mathbf{S}_w^T + \mathbf{S}_w) - (\mathbf{S}_b^T + \mathbf{S}_b)] \mathbf{W}]^T d\mathbf{W}) \end{aligned}$$

故可得到：

$$\frac{\partial L}{\partial \mathbf{W}} = (\lambda (\mathbf{S}_w^T + \mathbf{S}_w) - (\mathbf{S}_b^T + \mathbf{S}_b)) \mathbf{W}$$

偏导数为 0 时，有：

$$\mathbf{S}_b \mathbf{W} = \lambda \mathbf{S}_w \mathbf{W}$$

证毕。

3.5

证明 $X(X^T X)^{-1} X^T$ 是投影矩阵，并对线性回归模型从投影角度解释。

令 $P = X(X^T X)^{-1} X^T$ ，

$$\begin{aligned} P^2 &= (X(X^T X)^{-1} X^T)(X(X^T X)^{-1} X^T) \\ &= X I (X^T X)^{-1} X^T \\ &= X(X^T X)^{-1} X^T \\ &= P \end{aligned}$$

且，

$$P^T = (X(X^T X)^{-1} X^T)^T = X(X^T X)^{-1} X^T = P$$

对 X ，

$$PX = X(X^T X)^{-1} X^T X = X I = X$$

故 P 是投影矩阵。

- 解释

对于 $Ax = b$ 来说，并不是任何时候都有解，实际上大多数这种类型的方程都无解。 A 的列空间的含义是方程组有解时 b 的取值空间，当 b 不在 A 的列空间时，方程无解。具体来说，当 A 是行数大于列数的长方矩阵时，方程组中的方程大于未知数个数，肯定无解。

虽然方程无解，这就需要转而寻找最接近可解问题的解。对于无解方程 $Ax = b$ ， Ax 总是在列空间里（因为列空间本来就是由 Ax 确定的，和 b 无关），而 b 就不一定了，所以需要微调 b ，将 b 变成列空间中最接近它的一个， $Ax = b$ 变成了： $A\hat{x} = p$

p 就是 A 的列空间的投影， \hat{x} 表示此时的 x 已经不是原来 $Ax = b$ 中的 x ， $A\hat{x} \neq b$ ，因为方程无解，所以不可能有 $Ax = b$ 。此时问题转换为寻找最好的 A ，让它与原方程的差距最小：

假设 A 的秩是 2，此时列空间的维度也是 2，列空间将是一个平面，平面上的向量有无数个，其中最接近 b 的当然是 b 在平面上的投影，只有 $b - p$ 才能产生最小的误差向量。

所以说，要求解不等方程 $Ax = b$ ，需要将 b 微调成它在 A 的列空间上的投影（列空间上的向量很多， b 在列空间上的投影是唯一的）。如果 b 在列空间中， b 的投影就是它自己；如果 b 正交于列空间，它的投影是 0。

4 HW4

习题

• 4.1

• 4.9

• 假设离散随机变量 $X \in \{1, \dots, K\}$ ，其取值为 k 的概率 $P(X = k) = p_k$ ，其熵为 $H(p) = -\sum_k p_k \log_2 p_k$ ，试用朗格朗日乘法证明熵最大的分布为均匀分布

习题

• 下表表示的二分类数据集，具有三个属性 A, B, C，样本标记为两类“+”，“-”。请运用你学过的知识完成如下问题：

索引	A	B	C	类别
1	T	T	1.0	+
2	T	T	6.0	+
3	T	F	5.0	-
4	F	F	4.0	+
5	F	T	7.0	-
6	F	T	3.0	-
7	F	F	8.0	-
8	T	F	7.0	+
9	F	T	5.0	-
10	F	F	2.0	+

(a) 整个训练样本关于类属性的熵是多少
(b) 数据集中 A, B 两个属性的信息增益各是多少
(c) 对于属性 C，计算所有可能划分的信息增益
(d) 根据 Gini 指数，A 和 B 两个属性哪个是最优划分
(e) 采用算法 C4.5，构造决策树

4.1

因为决策树是通过属性来划分，相同属性的样本最终肯定会进入相同的叶节点。一个叶节点只有一个分类，如果样本属性相同而分类不同，必然产生训练误差。反之，决策树只会在当前样本集合是同一类或者所有属性相同时才会停止划分，最终得到训练误差为 0 的决策树。

4.2

数据集 D 的纯度 Gini:

$$\text{Gini}(D) = 1 - \sum_{k=1}^{|y|} p_k^2 = \sum_{k=1}^{|y|} p_k(1 - p_k)$$

属性 a 的纯度 Gini-index:

$$\text{Gini_index}(D, a) = \sum_{v=1}^v \frac{|D^v|}{|D|} \text{Gini}(D^v)$$

记 \tilde{D} 为 D 中在属性 a 上没有缺失值的样本子集, \tilde{D}^v 为 \tilde{D} 在属性 a 上取值为 a^v 的样本子集。属性 a 的基尼指数可推广为:

$$\text{Gini_index}(D, a) = p \times \text{Gini_index}(\tilde{D}, a) = p \times \sum_{v=1}^V \tilde{v} \text{Gini}(D^v)$$

4.3

拉格朗日函数为:

$$L(\mathbf{p}, \lambda) = - \sum_{i=1}^K p_i \log_2 p_i + \lambda \left(\sum_{i=1}^K p_i - 1 \right)$$

则:

$$\frac{\partial L(\mathbf{p}, \lambda)}{\partial p_i} = -\log_2 p_i - \frac{1}{\ln 2} + \lambda = 0$$

$$\frac{\partial L(\mathbf{p}, \lambda)}{\partial \lambda} = \sum_{i=1}^K p_i - 1 = 0$$

且 $p_1 = p_2 = \cdots = p_K = 2^{\lambda - \frac{1}{\ln 2}}$, 由 $\sum_{i=1}^K p_i = 1$ 得:

$$p_1 = p_2 = \cdots = p_K = \frac{1}{K}$$

拉朗格由上述方程组解出 p_i 及 λ , 如此求得的 p_i , 就是函数 $L(\mathbf{p}, \lambda)$ 在附加条件 $\sum_{i=1}^K p_i = 1$ 下的可能的极值点。若这样的点只有一个, 由实际问题可直接确定此即所求的点, 因此熵的最大分布为均匀分布。

4.4

4.4.1

10 个样本中, 正例 5 个, 负例 5 个, 故:

$$\text{Ent}(D) = -(0.5 \log_2 0.5) \times 2 = 1$$

4.4.2

A 中'T' 4 个 (里面正例 3 个), 'F' 6 个 (里面正例 2 个); B 中'T' 5 个 (里面正例 2 个), 'F' 5 个 (里面正例 3 个), 故:

$$\begin{aligned} \text{Gain}(D, A) &= \text{Ent}(D) - \sum_{v=1}^V \frac{|D^v|}{|D|} \text{Ent}(D^v) \\ &= 1 + 0.4 \times (0.75 \log_2 0.75 + 0.25 \log_2 0.25) + 0.6 \times \left(\frac{2}{6} \log_2 \frac{2}{6} + \frac{4}{6} \log_2 \frac{4}{6} \right) \\ &= 0.1245 \end{aligned}$$

$$\begin{aligned}
\text{Gain}(D, B) &= \text{Ent}(D) - \sum_{v=1}^V \frac{|D^v|}{|D|} \text{Ent}(D^v) \\
&= 1 + 0.5 \times (0.6 \log_2 0.6 + 0.4 \log_2 0.4) + 0.5 \times (0.6 \log_2 0.6 + 0.4 \log_2 0.4) \\
&= 0.0290
\end{aligned}$$

4.4.3

属性 C 有 8 个取值 1,2,3,4,5,6,7,8. 其中 5（2 负）和 7（一正一负）有两个样本其余一个。7 种二分阈值的信息增益分别为：

1. 1, 2-8:

$$\begin{aligned}
\text{Gain}(D, C) &= \text{Ent}(D) - \sum_{v=1}^V \frac{|D^v|}{|D|} \text{Ent}(D^v) \\
&= 1 - \frac{1}{10} \times 0 - \frac{9}{10} \times 0.9911 \\
&= 0.1080
\end{aligned}$$

2. 1-2, 3-8:

$$\begin{aligned}
\text{Gain}(D, C) &= \text{Ent}(D) - \sum_{v=1}^V \frac{|D^v|}{|D|} \text{Ent}(D^v) \\
&= 1 - \frac{2}{10} \times 0 - \frac{8}{10} \times 0.9544 \\
&= 0.2364
\end{aligned}$$

3. 1-3, 4-8:

$$\begin{aligned}
\text{Gain}(D, C) &= \text{Ent}(D) - \sum_{v=1}^V \frac{|D^v|}{|D|} \text{Ent}(D^v) \\
&= 1 - \frac{3}{10} \times 0.9183 - \frac{7}{10} \times 0.9852 \\
&= 0.0349
\end{aligned}$$

4. 1-4, 5-8:

$$\begin{aligned}
\text{Gain}(D, C) &= \text{Ent}(D) - \sum_{v=1}^V \frac{|D^v|}{|D|} \text{Ent}(D^v) \\
&= 1 - \frac{4}{10} \times 0.8113 - \frac{6}{10} \times 0.9183 \\
&= 0.1245
\end{aligned}$$

5. 1-5, 6-8:

$$\begin{aligned}
\text{Gain}(D, C) &= \text{Ent}(D) - \sum_{v=1}^V \frac{|D^v|}{|D|} \text{Ent}(D^v) \\
&= 1 - \frac{6}{10} \times 1 - \frac{4}{10} \times 1 \\
&= 0
\end{aligned}$$

6. 1-6, 7-8:

$$\begin{aligned}\text{Gain}(D, C) &= \text{Ent}(D) - \sum_{v=1}^V \frac{|D^v|}{|D|} \text{Ent}(D^v) \\ &= 1 - \frac{7}{10} \times 0.9852 - \frac{3}{10} \times 0.9183 \\ &= 0.0148\end{aligned}$$

7. 1-7, 8:

$$\begin{aligned}\text{Gain}(D, C) &= \text{Ent}(D) - \sum_{v=1}^V \frac{|D^v|}{|D|} \text{Ent}(D^v) \\ &= 1 - \frac{9}{10} \times 0.9911 - \frac{1}{10} \times 0 \\ &= 0.1080\end{aligned}$$

4.4.4

由基尼值的定义：

$$\begin{aligned}\text{Gini}(D) &= \sum_{k=1}^{|\mathcal{Y}|} \sum_{k' \neq k} p_k p_{k'} \\ &= 1 - \sum_{k=1}^{|\mathcal{Y}|} p_k^2\end{aligned}$$

对于 A, B 的 Gini 指数：

$$\begin{aligned}\text{Gini index}(D, A) &= \sum_{v=1}^V \frac{|D^v|}{|D|} \text{Gini}(D^v) \\ &= 0.4 \times (1 - 0.75^2 - 0.25^2) + 0.6 \times (1 - (\frac{2}{6})^2 - (\frac{4}{6})^2) \\ &= 0.4167\end{aligned}$$

$$\begin{aligned}\text{Gini index}(D, B) &= \sum_{v=1}^V \frac{|D^v|}{|D|} \text{Gini}(D^v) \\ &= 0.5 \times (1 - 0.4^2 - 0.6^2) \times 2 \\ &= 0.48\end{aligned}$$

所以 A 的 Gini 指数小，划分最优。

4.4.5

由前几小题可知：C 的信息增益选择第 2 种划分方式增益最大，由此选第二种作为 C 的处理方式。此时，

$$\begin{aligned}\text{GainRatio}(D, A) &= \frac{\text{Gain}(D, A)}{\text{IV}(A)} = \frac{0.1245}{0.9710} = 0.1232 \\ \text{GainRatio}(D, B) &= \frac{\text{Gain}(D, B)}{\text{IV}(B)} = \frac{0.0290}{1} = 0.0290 \\ \text{GainRatio}(D, C) &= \frac{\text{Gain}(D, C)}{\text{IV}(C)} = \frac{0.2364}{0.7219} = 0.3275\end{aligned}$$

所以第一层选择 C。此时，对于 C 的 1-2 子树全为正例，没必要继续计算，该子树直接判断为正例；对于 3-7 子树有：

$$\begin{aligned}\text{GainRatio}(D', A) &= \frac{\text{Gain}(D', A)}{\text{IV}'(A)} = \frac{0.9544 - 0.9512}{0.9544} = 0.003388 \\ \text{GainRatio}(D', B) &= \frac{\text{Gain}(D', B)}{\text{IV}'(B)} = \frac{0.9544 - 0.9056}{1} = 0.04876\end{aligned}$$

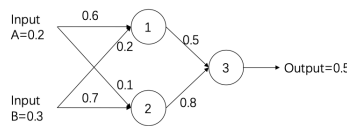
所以该子树选择 B 作为第二层。对于第三层，由表格数据可知：B 的 ‘T’ 类：A 为 ‘T’ 则判断为正类，A 为 ‘F’ 则判断为负类；B 的 ‘F’ 类：A 为 ‘T’ 由于正负都各出现一次，因此都可以，A 为 ‘F’ 则判断为正类。

5 HW5

作业



- 5.1
- 讨论 $\frac{\exp(x_i)}{\sum_{j=1}^C \exp(x_j)}$ 和 $\log \sum_{j=1}^C \exp(x_j)$ 的数值溢出问题
- 计算 $\frac{\exp(x_i)}{\sum_{j=1}^C \exp(x_j)}$ 和 $\log \sum_{j=1}^C \exp(x_j)$ 关于向量 $\mathbf{x} = [x_1, \dots, x_C]$ 的梯度
- 考虑如下简单网络，假设激活函数为 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（给出详细计算过程），并计算给定输入值 $\mathbf{x} = (0.2, 0.3)$ 时初始时和更新后的输出值，检查参数更新是否降低了平方损失值。



5.1

理想中的激活函数是阶跃函数，但是阶跃函数非连续，在 0 处不可导；线性激活函数没办法完全的拟合阶跃函数。线性函数在定义域内变换情况相同。使用线性函数作为激活函数时，无论是在隐藏层还是在输出层，不管传递多少层，其单元值都还是输入值和 x 的线性组合，若输出层也使用线性函数作为激活函数，那么就等价于线性回归，达不到激活与筛选的目的。

5.2

- $\frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}$

如果 x_i 很大, 那么指数操作时可能会大于数据类型容许的最大数字, 造成上溢, 这将使分母或分子变为 ∞ , 这就使得结果是 $0, \infty$ 还是 nan 难以得到准确结论。但可以在计算之前对所有的 x_i 减去 $\max(x_i)$ 避免上溢。

- $\log \sum_{j=1}^n e^{x_j}$

一样的原因, e^{x_i} 计算可能溢出, 也可以做同样处理, 取 $c = \max(x_i)$,

$$\begin{aligned} \log \left(\sum_{i=1}^n e^{x_i} \right) &= \log \left(\sum_{i=1}^n e^{x_i - c} e^c \right) \\ &= \log \left(e^c \sum_{i=1}^n e^{x_i - c} \right) \\ &= \log \left(\sum_{i=1}^n e^{x_i - c} \right) + \log(e^c) \\ &= \log \left(\sum_{i=1}^n e^{x_i - c} \right) + c \end{aligned}$$

5.3

记 $f(x) = \frac{e^{x_i}}{\sum_{j=1}^C e^{x_j}}$

$$\frac{\partial f(x)}{\partial x_k} = \begin{cases} -\frac{\exp(x_i + x_k)}{\sum_{j=1}^C \exp(x_j)^2}, & k \neq i \\ \frac{\exp(x_k) \sum_{j=1, j \neq k}^C \exp(x_j)}{\sum_{j=1}^C \exp(x_j)^2}, & k = i \end{cases}$$

记 $g(x) = \log \frac{e^{x_i}}{\sum_{j=1}^C e^{x_j}}$

$$\frac{\partial g(x)}{\partial x_k} = \begin{cases} -\frac{\exp(x_k)}{\sum_{j=1}^C \exp(x_j)}, & k \neq i \\ 1 - \frac{\exp(x_k)}{\sum_{j=1}^C \exp(x_j)}, & k = i \end{cases}$$

5.4

由题得:

$$h_1 = 0.6 \times 0.2 + 0.2 \times 0.3 = 0.18$$

$$h_2 = 0.1 \times 0.2 + 0.7 \times 0.3 = 0.23$$

经过激活函数 ReLU: $\max(0, x)$ 变换后值为:

$$\text{out}_{h_1} = \max(0, h_1) = 0.18$$

$$\text{out}_{h_2} = \max(0, h_2) = 0.23$$

所以输出值:

$$\text{Output} = 0.18 \times 0.5 + 0.23 \times 0.8 = 0.274$$

此时误差为：

$$E = \frac{1}{2}(0.5 - 0.274)^2 = 0.025538$$

BP(标准)：记 $out_{h_1} : o_1, out_{h_2} : o_2$ 由于：

$$E = \frac{1}{2}(y - w_5 o_1 - w_6 o_2)^2$$

故可得到：

$$\frac{\partial E}{\partial w_5} = -o_1(y - w_5 o_1 - w_6 o_2)$$

一次更新：

$$\begin{aligned} w_{5new} &= w_{5old} - \eta \frac{\partial E}{\partial w_{5old}} \\ &= 0.5 - 1 \times (-0.18(0.5 - 0.274)) \\ &= 0.54068 \end{aligned}$$

同理：

$$\begin{aligned} w_{6new} &= w_{6old} - \eta \frac{\partial E}{\partial w_{6old}} \\ &= 0.8 - 1 \times (-0.23(0.5 - 0.274)) \\ &= 0.85198 \end{aligned}$$

此时有： $h_1 = o_1, h_2 = o_2$ ，故由于：

$$h_1 = w_1 A + w_3 B$$

有：

$$\begin{aligned} E &= \frac{1}{2}(y - w_5 o_1 - w_6 o_2)^2 \\ &= \frac{1}{2}(y - w_5 h_1 - w_6 h_2)^2 \\ &= \frac{1}{2}(y - w_5(w_1 A + w_3 B) - w_6 h_2)^2 \\ \frac{\partial E}{\partial w_1} &= -w_5 A(y - w_5(w_1 A + w_3 B) - w_6 h_2) \end{aligned}$$

所以：

$$\begin{aligned} w_{1new} &= w_{1old} - \eta \frac{\partial E}{\partial w_{1old}} \\ &= 0.6 - 1 \times (-0.5 \times 0.2(0.5 - 0.274)) \\ &= 0.6226 \end{aligned}$$

同理：

$$\begin{aligned} w_{3new} &= w_{3old} - \eta \frac{\partial E}{\partial w_{3old}} \\ &= 0.2 - 1 \times (-0.5 \times 0.3(0.5 - 0.274)) \\ &= 0.2339 \end{aligned}$$

$$\begin{aligned}
w_{2new} &= w_{2old} - \eta \frac{\partial E}{\partial w_{2old}} \\
&= 0.1 - 1 \times (-0.8 \times 0.2(0.5 - 0.274)) \\
&= 0.13616 \\
w_{4new} &= w_{4old} - \eta \frac{\partial E}{\partial w_{4old}} \\
&= 0.7 - 1 \times (-0.8 \times 0.3(0.5 - 0.274)) \\
&= 0.75424
\end{aligned}$$

此时，更新后的输出值为：

$$\begin{aligned}
Output_{new} &= \max((0.6226 \times 0.2 + 0.2339 \times 0.3), 0) \times 0.54068 \\
&\quad + \max((0.13616 \times 0.2 + 0.75424 \times 0.3), 0) \times 0.85198 \\
&= 0.32125
\end{aligned}$$

平方损失：

$$E_{new} = \frac{1}{2}(0.5 - 0.32125)^2 = 0.01598 < 0.025538$$

因此参数更新降低了平方损失值。

6 HW6

作业



42

- 6.4
- 6.6
- 6.9

支持向量回归的对偶问题如下，

$$\begin{aligned}
\max_{\alpha, \hat{\alpha}} g(\alpha, \hat{\alpha}) &= -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m (\alpha_i - \hat{\alpha}_i)(\alpha_j - \hat{\alpha}_j) \kappa(\mathbf{x}_i, \mathbf{x}_j) + \sum_{i=1}^m (y_i(\hat{\alpha}_i - \alpha_i) - \epsilon(\hat{\alpha}_i + \alpha_i)) \\
s.t. \quad C &\geq \alpha, \hat{\alpha} \geq 0 \text{ and } \sum_{i=1}^m (\alpha_i - \hat{\alpha}_i) = 0
\end{aligned}$$

请将该问题转化为类似于如下标准型的形式 ($\mathbf{u}, \mathbf{v}, \mathbf{K}$ 均已知)，

$$\begin{aligned}
\max_{\alpha} g(\alpha) &= \alpha^T \mathbf{v} - \frac{1}{2} \alpha^T \mathbf{K} \alpha \\
s.t. \quad C &\geq \alpha \geq 0 \text{ and } \alpha^T \mathbf{u} = 0
\end{aligned}$$

例如在软间隔SVM中 $\mathbf{v} = \mathbf{1}, \mathbf{u} = \mathbf{y}, \mathbf{K}[i, j] = y_i y_j \kappa(\mathbf{x}_i, \mathbf{x}_j)$.

若 $\kappa(x_i, x_j) = \phi(x_i)^T \phi(x_j) = (x_i^T x_j)^2$ ，求 $\phi(x_i)$ 表达式。

6.1

线性判别分析能够解决多分类问题，而 SVM 只能解决二分类问题。故当两类样本线性可分时，且处理二分类问题时等价。

6.2

SVM 的目的是求出与支持向量有最大化距离的直线, 以每个样本为圆心, 该距离为半径做圆, 可以近似认为圆内的点与该样本属于相同分类。如果出现了噪声, 那么这个噪声所带来的错误分类也将最大化, 所以 SVM 对噪声是很敏感的。

6.3

使用对率损失函数 ℓ_{\log} 来代替式 (6.29) 中的 0/1 损失函数, 则几乎就得到了对率回归模型 (3.27), 即:

由 SVM 模型:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{s.t.} \quad y_i (\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 \quad i = 1, \dots, m$$

软间隔的支持向量机:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \ell_{0/1} (y_i (\mathbf{w}^\top \phi(\mathbf{x}_i) + b) - 1)$$

令 $y_i (\mathbf{w}^\top \phi(\mathbf{x}_i) + b) = Z$, 对率损失:

$$\ell_{\log}(z) = \log(1 + e^{-z}) = \log\left(\frac{1 + e^z}{e^z}\right) = \log(1 + e^z) - z$$

即核对率回归:

$$\begin{aligned} \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (-z + \log(1 + e^z)) \\ \ell(\beta) = \sum_{i=1}^m \left(-y_i \beta^\top \hat{\mathbf{x}}_i + \log(1 + e^{\beta^\top \hat{\mathbf{x}}_i}) \right) \end{aligned}$$

6.4

令

$$\begin{aligned} \boldsymbol{\alpha} = \begin{pmatrix} \alpha \\ \hat{\alpha} \end{pmatrix} \in \mathbb{R}^{2m}, \mathbf{v} = \begin{pmatrix} -\mathbf{y} - \epsilon \\ \mathbf{y} - \epsilon \end{pmatrix} \in \mathbb{R}^{2m} \\ \mathbf{K} = \begin{pmatrix} K & -K \\ -K & K \end{pmatrix} \in \mathbb{R}^{2m \times 2m} \text{ 其中 } K_{ij} = \kappa(\mathbf{x}_i, \mathbf{x}_j) \end{aligned}$$

则可以得到结论.

$$\begin{aligned}
\kappa(x_i, x_j) &= (x_i^T x_j)^2 \\
&= (x_i^T x_j) (x_i^T x_j) \\
&= \left(\sum_{p=1}^d x_{pi} x_{pj} \right) \left(\sum_{q=1}^d x_{qi} x_{qj} \right) \\
&= \sum_{p=1}^d \sum_{q=1}^d x_{pi} x_{pj} x_{qi} x_{qj} \\
&= \begin{pmatrix} x_1 x_1 \\ x_1 x_2 \\ \vdots \\ x_1 x_d \\ \vdots \\ x_2 x_d \\ \vdots \\ x_d x_1 \\ \vdots \\ x_d x_d \end{pmatrix} \times \begin{pmatrix} x_1 x_1 \\ x_1 x_2 \\ \vdots \\ x_1 x_d \\ \vdots \\ x_2 x_d \\ \vdots \\ x_d x_1 \\ \vdots \\ x_d x_d \end{pmatrix}
\end{aligned}$$

则可以得到：

$$\phi : \mathbf{x} \in \mathbb{R}^d \rightarrow \begin{pmatrix} x_1 x_1 \\ x_1 x_2 \\ \vdots \\ x_1 x_d \\ \vdots \\ x_2 x_d \\ \vdots \\ x_d x_1 \\ \vdots \\ x_d x_d \end{pmatrix} \in \mathbb{R}^{d^2}$$

7 HW7

作业



61

- 习题7.4
- 习题7.5
- 证明EM算法的收敛性
- 在HMM中, 求解概率 $P(x_{n+1}|x_1, x_2, \dots, x_n)$

7.1

实践中使用式 (7.15) 决定分类类别时, 若数据的维数非常高, 则概率连乘 $\prod_{i=1}^d P(x_i|c)$ 的结果通常会非常接近于 0 从而导致下溢. 试述防止下溢的可能方案.

(7.15):

$$h_{nb}(\mathbf{x}) = \arg \max_{c \in \mathcal{Y}} P(c) \prod_{i=1}^d P(x_i | c)$$

若连乘的式子太多, 导致乘积接近 0。由于属性个数是已知的, 可以对每个乘式做适当次的开方处理, 可以保证结果不会为 0。另外也可以对各项取对数, 当累加太多时, 可能导致和接近负无穷。可以对每个加数除以属性的个数, 来防止溢出。

7.2

试证明: 二分类任务中两类数据满足高斯分布且方差相同时, 线性判别分析产生贝叶斯最优分类器.

假设 1 类样本均值为 u_1 , 2 类样本均值为 u_2 .

由题意, 数据满足高斯分布且方差相同, 因此, 当样本足够大时, 认为等价于: 线性判别分析公式 $J = \frac{|w^T(u_1 - u_2)|^2}{w^T(\Sigma_1 + \Sigma_2)w}$ 求最大值; 对 $\frac{1}{J} = \frac{w^T(\Sigma_1 + \Sigma_2)w}{|w^T(u_1 - u_2)|^2} = \sum_i \frac{(1 - y_i) |w^T(x_i - u_1)|^2 + y_i |w^T(x_i - u_2)|^2}{|w^T(u_1 - u_2)|^2}$ 求最小值。

最优贝叶斯分类器使得每个训练样本的后验概率 $P(c|x)$ 最大, 对应线性判别分析中, 离对应分类的中心距离平方除以两个分类中心的距离平方越小。即求

$$\sum_i \frac{(1 - y_i) |w^T(x_i - u_1)|^2 + y_i |w^T(x_i - u_2)|^2}{|w^T(u_1 - u_2)|^2}$$

的最小值。

此时, 显然有两个式子相同。故, 线性判别分析产生最优贝叶斯分类器。

7.3

要证明 EM 算法收敛，即证明对数似然函数的值在迭代中一直增大。即

$$\sum_{i=1}^m \log P(x^{(i)}; \theta^{j+1}) \geq \sum_{i=1}^m \log P(x^{(i)}; \theta^j)$$

由于

$$L(\theta, \theta^j) = \sum_{i=1}^m \sum_{z^{(i)}} P(z^{(i)} | x^{(i)}; \theta^j) \log P(x^{(i)}, z^{(i)}; \theta)$$

令

$$H(\theta, \theta^j) = \sum_{i=1}^m \sum_{z^{(i)}} P(z^{(i)} | x^{(i)}; \theta^j) \log P(z^{(i)} | x^{(i)}; \theta)$$

相减得到

$$\sum_{i=1}^m \log P(x^{(i)}; \theta) = L(\theta, \theta^j) - H(\theta, \theta^j)$$

分别取 θ 为 θ^{j+1} 和 θ^j , 再相减得到

$$\sum_{i=1}^m \log P(x^{(i)}; \theta^{j+1}) - \sum_{i=1}^m \log P(x^{(i)}; \theta^j) = [L(\theta^{j+1}, \theta^j) - L(\theta^j, \theta^j)] - [H(\theta^{j+1}, \theta^j) - H(\theta^j, \theta^j)]$$

此时，要证明 EM 算法的收敛性，只需要证明上式的右边是非负的即可。

由于 θ^{j+1} 使得 $L(\theta, \theta^j)$ 极大，故

$$L(\theta^{j+1}, \theta^j) - L(\theta^j, \theta^j) \geq 0$$

再由于

$$\begin{aligned} H(\theta^{j+1}, \theta^j) - H(\theta^j, \theta^j) &= \sum_{i=1}^m \sum_{z^{(i)}} P(z^{(i)} | x^{(i)}; \theta^j) \log \frac{P(z^{(i)} | x^{(i)}; \theta^{j+1})}{P(z^{(i)} | x^{(i)}; \theta^j)} \\ &\leq \sum_{i=1}^m \log \left(\sum_{z^{(i)}} P(z^{(i)} | x^{(i)}; \theta^j) \frac{P(z^{(i)} | x^{(i)}; \theta^{j+1})}{P(z^{(i)} | x^{(i)}; \theta^j)} \right) \\ &= \sum_{i=1}^m \log \left(\sum_{z^{(i)}} P(z^{(i)} | x^{(i)}; \theta^{j+1}) \right) \\ &= 0 \end{aligned}$$

故，EM 算法的收敛性得证。

7.4

$$\begin{aligned} P(x_{n+1} | x_1, x_2, \dots, x_n) &= \frac{P(x_1, x_2, \dots, x_{n+1})}{P(x_1, x_2, \dots, x_n)} \\ &= \frac{P(y_1) P(x_1 | y_1) \prod_{k=2}^{n+1} P(y_k | y_{k-1}) P(x_k | y_k)}{P(y_1) P(x_1 | y_1) \prod_{k=2}^n P(y_k | y_{k-1}) P(x_k | y_k)} \\ &= P(y_{n+1} | y_n) P(x_{n+1} | y_{n+1}) \end{aligned}$$

8 HW8

HW8&9 如下:

- 8.2
- 8.8
- 给定任意的两个相同长度向量 \mathbf{x}, \mathbf{y} , 其余弦距离为 $1 - \frac{\mathbf{x}^\top \mathbf{y}}{|\mathbf{x}||\mathbf{y}|}$, 证明余弦距离不满足传递性, 而余弦夹角 $\arccos\left(\frac{\mathbf{x}^\top \mathbf{y}}{|\mathbf{x}||\mathbf{y}|}\right)$ 满足
- 证明k-means算法的收敛性
- 在k-means算法中替换欧式距离为其他任意的度量, 请问“聚类簇”中心如何计算?

8.1

式 (8.5)

$$\ell_{\text{exp}}(H \mid \mathcal{D}) = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [e^{-f(\mathbf{x})H(\mathbf{x})}]$$

由式 (8.4)

$$H(\mathbf{x}) = \sum_{t=1}^T \alpha_t h_t(\mathbf{x})$$

再由式 (8.11)

$$\alpha_t = \frac{1}{2} \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$$

可以看出分类器的权重只与分类器的错误率负相关。

对于式 (8.5) 的情况时, 由于 $f(x) \in \{+1, -1\}$, $H(x)$ 为实数, 当两者同号时, 有 $f(x)H(x) > 0$, 此时, $e^{-f(x)H(x)} = e^{-|H(x)|} < 1$, $|H(x)|$ 越大损失函数越小; 两者异号时, 有 $f(x)H(x) < 0$, 此时, $e^{-f(x)H(x)} = e^{|H(x)|} > 1$, $|H(x)|$ 越大损失函数越大。且由于指数函数的性质, 在两个区间的时候损失函数都是单调的。

由此, 当使用任意损失函数 $\ell(-f(x)H(x))$, 且对于 $H(x)$ 在区间 $[-\infty, \delta](\delta > 0)$ 单调递减的情况, 仍然满足损失函数是分类任务原本的 0 / 1 损失函数的一致的替代损失函数。

8.2

MultiBoosting 由于集合了 Bagging, Wagging, AdaBoost, 可以有效的降低误差和方差, 特别是误差, 但是训练成本和预测成本都会显著增加。

Iterative Bagging 相比 Bagging 会降低误差, 但是方差上升。由于 Bagging 本身就是一种降低方差的算法, 所以 Iterative Bagging 相当于 Bagging 与单分类器的折中。

9 HW9

9.1

- 余弦距离不满足传递性:

不妨令空间中三个点的坐标如下:

$$x = (1, 0), y = (0, 1), z = \left(\frac{\sqrt{2}}{2}, \frac{\sqrt{2}}{2}\right)$$

则, 余弦距离有:

$$d(x, y) = 1, d(x, z) = 1 - \frac{\sqrt{2}}{2}, d(z, y) = 1 - \frac{\sqrt{2}}{2}$$

此时, $d(x, z) + d(y, z) < d(x, y)$. 故余弦距离不满足传递性。

- 余弦夹角满足传递性:

证明:

$$\arccos\left(\frac{(x^\top y)}{|x||y|}\right) \leq \arccos\left(\frac{(x^\top z)}{|x||z|}\right) + \arccos\left(\frac{(z^\top y)}{|z||y|}\right)$$

两侧同时取 \cos , 由于在区间 $[0, \pi]$ 上单调递减, 于是在该区间有:

$$\begin{aligned} \frac{(x^\top y)}{|x||y|} &\geq \cos\left(\arccos\left(\frac{(x^\top z)}{|x||z|}\right) + \arccos\left(\frac{(z^\top y)}{|z||y|}\right)\right) \\ &= \frac{(x^\top z)}{|x||z|} * \frac{(z^\top y)}{|z||y|} - \sqrt{\left(1 - \left(\frac{(x^\top z)}{|x||z|}\right)^2\right) * \left(1 - \left(\frac{(z^\top y)}{|z||y|}\right)^2\right)} \\ &= \frac{x^\top (zz^\top) y}{|x||z||z||y|} - \sqrt{\left(1 - \left(\frac{(x^\top z)}{|x||z|}\right)^2\right) * \left(1 - \left(\frac{(z^\top y)}{|z||y|}\right)^2\right)} \\ &= \frac{(x^\top y)}{|x||y|} - \sqrt{\left(1 - \left(\frac{(x^\top z)}{|x||z|}\right)^2\right) * \left(1 - \left(\frac{(z^\top y)}{|z||y|}\right)^2\right)} \end{aligned}$$

由于: $\sqrt{\left(1 - \left(\frac{(x^\top z)}{|x||z|}\right)^2\right) * \left(1 - \left(\frac{(z^\top y)}{|z||y|}\right)^2\right)} \geq 0$, 故不等式得证。余弦夹角满足传递性。

9.2

证明 k-means 算法的收敛性

即证明:

$$\min J(u_1, u_2, \dots, u_k) = 1/2 \sum_{i=1}^m \sum_{j=1}^k (x_i - u_j)^2$$

1) 单调; 2) 有界。

1. 单调

- (a) 更新 m 个点的归属类别：将 m 个点归属到已有的 k 个中心，其中的判定就是离哪个中心点近，就归属到那一类，此时损失函数为 $J_0(u_1, u_2, \dots, u_k)$. 在中心确定的情况下，如果不归属到最近的中心，归属到其他中心，此刻损失函数假设为 $J_1(u_1, u_2, \dots, u_k)$ ，一定有 $J_0 \leq J_1$. (假如其中 $m-1$ 个点都归属到最近的中心了，某一个点没有归属到最近中心，那么如果把这个点放到最近的那个中心，损失函数就会降低)
- (b) 假设第一个类别有 N_1 个点，中心为 u_j ，这个类别的损失函数部分为 $1/2 \sum_{i=1}^{N_1} (x_i - u_1)^2$ ，且有： $1/2 \sum_{i=1}^{N_1} (x_i - u_1)^2 \geq 1/2 \sum_{i=1}^{N_1} (x_i - x_{ave})^2$ 其中 x_{ave} ， x_{ave} 为其均值中心。固定归属类别，通过把均值中心作为类别中心，降低了 J

2. 有界

$$J(u_1, u_2, \dots, u_k) \geq 0$$

故，单调，有界，收敛。

9.3

在 k-means 算法中替换欧氏距离为其他任意的度量，请问聚类簇“中心”如何计算？

直接进行平均计算得到的均值是符合欧几里得距离定义下的均值的定义的，但是并不一定符合其它距离计算公式。不同距离的定义方式需要设计一套单独的“均值”的定义方式进行聚类簇“中心”的计算。若损失函数是凸的，则直接对其求梯度并令其等于 0，可得到聚类中心。同时，很可能完全定义不出来，如曼哈顿距离。

10 HW10

习题



42

- 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

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

提示：柯西-施瓦兹不等式 $(\sum_i a_i)^2 \leq n(\sum_i a_i^2)$

- 10.4

- 求解20页的优化问题

附加题：令 $\mathbf{M} = \mathbf{P}\mathbf{P}^T$ ，那么下列问题还是凸优化问题吗？试证明之。

$$\min_{\mathbf{P}} \sum_{(x_i, x_j) \in \mathcal{M}} \|x_i - x_j\|_{\mathbf{M}}^2 \quad \text{s. t.} \quad \sum_{(x_i, x_j) \in \mathcal{C}} \|x_i - x_j\|_{\mathbf{M}}^2 \geq 1$$

10.1

由题可知：

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

要证 $\text{err}^*(\mathbf{x}) \leq \text{err}(\mathbf{x})$

即证 $\max_{c \in Y} P(c | \mathbf{x}) \geq \sum_c P(c | \mathbf{x}) P(c | \mathbf{z})$

$$\begin{aligned}\sum_c P(c | \mathbf{x}) P(c | \mathbf{z}) &\approx \sum_c P(c | \mathbf{x})^2 \\ &\leq \max_{c \in Y} P(c | \mathbf{x}) \sum_c P(c | \mathbf{x}) \\ &= \max_{c \in Y} P(c | \mathbf{x})\end{aligned}$$

得证.

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

$$\begin{aligned}\text{err}(x) &= 1 - \sum_x P^2(c | x) \\ &= (2 - \text{err}^*(x)) \text{err}^*(x) - \sum_{c|c^*} P^2(c | x)\end{aligned}$$

由 Cauchy 不等式：

$$\sum_{c|c^*} P^2(c | x) \geq \frac{1}{|\mathcal{Y}|} (1 - P(c^* | x))^2 = \frac{1}{|\mathcal{Y}|} \text{err}^*(\mathbf{x})^2$$

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

10.2

在实践中,协方差矩阵 $\mathbf{X}\mathbf{X}^T$ 的特征值分解常由中心化后的样本矩阵 \mathbf{X} 的奇异值分解代替,试述其原因.

因为两者等价,下证明:

\mathbf{X} 的奇异值分解为:

$$\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T, \quad \mathbf{X}\mathbf{X}^T = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T (\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T)^T = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T \mathbf{V}\mathbf{\Sigma}^T \mathbf{U}^T$$

此时, $\mathbf{V}^T \mathbf{V} = \mathbf{I}, \mathbf{U}^T \mathbf{U} = \mathbf{I}$, 故,

$$\mathbf{X}\mathbf{X}^T = \mathbf{U}\mathbf{\Sigma}\mathbf{\Sigma}^T \mathbf{U}^T = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$$

\mathbf{X} 的特征值分解时, $\mathbf{X}^T \mathbf{X} = \mathbf{P}\mathbf{\Lambda}\mathbf{P}^T$, $\mathbf{P}=\mathbf{U}$ 则两者等价。

还有,用奇异值分解代替特征分解,节省了计算和存储的成本,且计算精度较高。

10.3

主成分分析—求解

$$\max_W \text{tr}(\mathbf{W}^T \mathbf{X} \mathbf{X}^T \mathbf{W}) \quad \text{s.t.} \quad \mathbf{W}^T \mathbf{W} = \mathbf{I}_{d'}$$

使用拉格朗日乘子法可得

$$\mathbf{X} \mathbf{X}^T \mathbf{W} = \Lambda \mathbf{W}$$

由题得：拉格朗日函数为：

$$L(w, \theta) = \text{tr}(\mathbf{W}^T \mathbf{X} \mathbf{X}^T \mathbf{W}) - \text{tr}(\theta^T (\mathbf{W}^T \mathbf{W} - \mathbf{I}))$$

其中， θ 乘子矩阵。考虑其中前 d' 个最大的特征值，令 $\theta = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_{d'})$ 则：

$$\begin{aligned} \frac{\partial L(W, \theta)}{\partial W} &= \frac{\partial (\text{tr}(\mathbf{W}^T \mathbf{X} \mathbf{X}^T \mathbf{W}))}{\partial W} - \frac{\partial (\text{tr}(\theta^T (\mathbf{W}^T \mathbf{W} - \mathbf{I})))}{\partial W} \\ &= 2\mathbf{X} \mathbf{X}^T \mathbf{W} - \mathbf{W} \theta - \mathbf{W} \theta^T \\ &= 2\mathbf{X} \mathbf{X}^T \mathbf{W} - 2\mathbf{W} \theta \end{aligned}$$

为 0 时，有

$$\mathbf{X} \mathbf{X}^T \mathbf{W} = \mathbf{W} \theta$$

代回，有：

$$\begin{aligned} &\max_W \text{tr}(\mathbf{W}^T \mathbf{X} \mathbf{X}^T \mathbf{W}) \\ &= \max_W \sum_{i=1}^{d'} \mathbf{W}_i^T \mathbf{X} \mathbf{X}^T \mathbf{W}_i \\ &= \max_{i=1}^I \sum_i \lambda_i \mathbf{W}_i^T \mathbf{W}_i \\ &= \max_x \sum_{i=1}^{d'} \lambda_i \end{aligned}$$

11 HW11

作业



45

- 11.5
- 11.7
- PPT 20页: 证明回归和对率回归的损失函数的梯度是否满足L-Lipschitz条件, 并求出L

11.1

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.

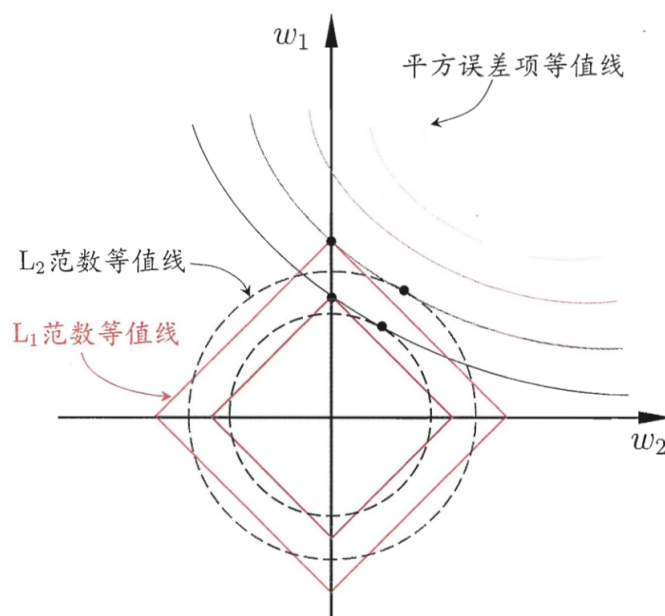


图 11.2 L_1 正则化比 L_2 正则化更易于得到稀疏解

L_1 正则化之所以可以产生稀疏解, 主要是因为平方误差项等值线与 L_1 等值线的第一个交点位于坐标轴上, 如书上图 11.2 所示, 当平方误差项等值线的曲率比较大时, 就会导致其与 L_1 等值线的第一个交点不再位于坐标轴上, 此时就无法产生稀疏解。

11.2

试述直接求解 L0 范数正则化会遇到的困难。

L0 范数是不连续的，而且是非凸函数，无法通过优化直接求解，必须采用遍历的方式，因此导致这个问题是个 NP 难问题。

11.3

11.3.1

回归：

$$f(w) = \|Xw - y\|^2$$

故，

$$\nabla f(w) = 2X^T(Xw - y)$$

任取 w_1, w_2 ,

$$\begin{aligned} \|\nabla f(w_1) - \nabla f(w_2)\|_2 &= 2\|X^T(Xw_1 - y) - X^T(Xw_2 - y)\|_2 \\ &= 2\|X^T X(w_1 - w_2)\|_2 \\ &\leq 2\|X^T X\|_2 \|w_1 - w_2\|_2 \end{aligned}$$

则，满足 L-Lipschitz 条件，且 $\mathcal{L} = 2\|X^T X\|_2 = 2\|X\|_2^2$

11.3.2

对率回归：

$$f(w) = \sum_{i=1}^m (-yw^T x_i + \ln(1 + e^{w^T x_i}))$$

故，

$$\nabla f(w) = -\sum_{i=1}^m (y_i - \frac{e^{w^T x_i}}{1 + e^{w^T x_i}}) x_i$$

任取 w_1, w_2 ,

$$\begin{aligned} \|\nabla f(w_1) - \nabla f(w_2)\|_2 &= \left\| \left(-\sum_{i=1}^m (y_i - \frac{e^{w_1^T x_i}}{1 + e^{w_1^T x_i}}) x_i \right) - \left(-\sum_{i=1}^m (y_i - \frac{e^{w_2^T x_i}}{1 + e^{w_2^T x_i}}) x_i \right) \right\|_2 \\ &= \left\| \sum_{i=1}^m \left(\frac{e^{w_1^T x_i}}{1 + e^{w_1^T x_i}} - \frac{e^{w_2^T x_i}}{1 + e^{w_2^T x_i}} \right) x_i \right\|_2 \end{aligned}$$

$$\text{令 } g(\alpha) = \frac{e^\alpha}{1 + e^\alpha}$$

$$\nabla g(\alpha) = \frac{e^\alpha}{(1 + e^\alpha)^2} \geq 0, \quad \nabla^2 g(\alpha) = -\frac{e^\alpha (e^\alpha - 1)}{(1 + e^\alpha)^3}$$

二阶导数为 0 时, $\alpha = 0$, 一阶导数在此处取得极大值 $\frac{1}{4}$.

$$\begin{aligned}\|g(\alpha) - g(\beta)\|_2 &= \left\| \frac{e^\alpha}{1+e^\alpha} - \frac{e^\beta}{1+e^\beta} \right\|_2 \\ &= \|\nabla(\xi)(\alpha - \beta)\|_2 \quad \xi \in (\alpha, \beta) \\ &\leq \frac{1}{4}\|\alpha - \beta\|_2\end{aligned}$$

故,

$$\begin{aligned}\|\nabla f(w_1) - \nabla f(w_2)\|_2 &\leq \frac{1}{4} \left\| \sum (w_2^T x_i - w_1^T x_i) x_i \right\|_2 \\ &= \frac{1}{4} \|x^T x (w_1 - w_2)\|_2 \\ &\leq \frac{1}{4} \|x^T x\|_2 \|w_1 - w_2\|_2\end{aligned}$$

则, 满足 L-Lipschitz 条件, 且 $\mathcal{L} = \frac{1}{4} \|X^T X\|_2 = \frac{1}{4} \|X\|_2^2$

12 HW12

无

13 HW13

无

14 HW14

1. 假设数据集 $D = \{x_1, x_2, \dots, x_m\}$, 任意 x_i 是从均值为 μ 、方差为 λ^{-1} 的正态分布 $\mathcal{N}(\mu, \lambda^{-1})$ 中独立采样而得到。假设 μ 和 λ 的先验分布为 $p(\mu, \lambda) = \mathcal{N}(\mu|\mu_0, (\kappa_0\lambda)^{-1})\text{Gam}(\lambda|a_0, b_0)$ 其中 $\text{Gam}(\lambda|a_0, b_0) = \frac{1}{\Gamma(a_0)} b_0^{a_0} \lambda^{a_0-1} \exp(-b_0\lambda)$

(1) 请写出联合概率分布 $p(D, \mu, \lambda)$

(2) 请写出证据下界 (即变分推断的优化目标), 并证明其为观测数据边缘似然 $\sum_{i=1}^m \log p(x_i)$ 的下界

(3) (3) 请用变分推断法近似推断后验概率 $p(\mu, \lambda|D)$

2. PPT 46, 给出 CRF 的预测问题的解法

14.1

14.1.1

$$\begin{aligned}
p(D, \mu, \lambda) &= p(\mu, \lambda) \cdot p(D \mid \mu, \lambda) \\
&= p(\mu, \lambda) \cdot \prod_{i=1}^m p(x_i \mid \mu, \lambda) \\
&= \frac{1}{\sqrt{2\pi(\kappa_0\lambda)^{-1}}} \exp \left\{ -\frac{1}{2(\kappa_0\lambda)^{-1}} (\mu - \mu_0)^2 \right\} \\
&\quad \cdot \frac{1}{\Gamma(a_0)} b_0^{a_0} \cdot \lambda^{a_0-1} \exp \{-b_0\lambda\} \cdot \left(\frac{\lambda}{2\pi} \right)^{\frac{m}{2}} \exp \left\{ -\frac{\lambda}{2} \sum_{i=1}^m (x_i - \mu)^2 \right\}
\end{aligned}$$

14.1.2

由题可知，变分推断的目的是寻找一个参数的概率密度函数

$$q(z) \text{ s.t. } q^*(z) = \arg \min_{q(z)} KL(q(z) \parallel p(z|x))$$

其中，KL 散度为

$$KL(q(z) \parallel p(z \mid x)) = \mathbf{E}[\log q(z)] - \mathbf{E}[\log p(z, x)] + \log p(x)$$

再由于 KL 散度的非负性，故可得到：

$$\begin{aligned}
ELBO &= \mathbf{E}[\log p(z, x)] - \mathbf{E}[\log q(z)] \\
&= \mathbf{E}_q[\log p(\lambda)] + \mathbf{E}_q[\log p(\mu \mid \lambda)] + \mathbf{E}_q[\log p(x \mid \mu, \lambda)] - \mathbf{E}_q[\log q(\lambda)] - \mathbf{E}_q[\log q(\mu)]
\end{aligned}$$

14.1.3

由题可得，

$$\begin{aligned}
\frac{\partial L}{\partial q_u(u)} &= E_\lambda[\log p(\mu \mid \lambda)] + E_\lambda[\log p(D \mid \mu, \lambda)] - \log q(u) = 0 \\
\log q(u) &= \frac{E[\lambda]\kappa_0}{-2} (\mu - \mu_0)^2 - \frac{E[\lambda]}{2} \sum_{i=1}^m (x_i - \mu)^2 \\
&= -\frac{E[x]}{2} \left[(\kappa_0 + m) \mu^2 + \sum_{i=1}^m x_i^2 - 2\mu(\kappa_0\mu_0 + m\bar{x}) \right] \\
&= -\frac{E[\lambda]}{2} \left[(\kappa_0 + m) \left(\mu - \frac{\kappa_0\mu_0 + m\bar{x}}{\kappa_0 + m} \right)^2 + \sum_{i=1}^m x_i^2 - \frac{(\kappa_0\mu_0 + m\bar{x})^2}{\kappa_0 + m} \right]
\end{aligned}$$

令 $\mu_m = \frac{\kappa_0\mu_0 + m\bar{x}}{\kappa_0 + m}$, $\lambda_m = (\kappa_0 + m)E[\lambda]$, 故可以看出: $q_u^*(u) \sim \mathcal{N}(\mu \mid \mu_m, \lambda_m)$.

再由

$$\frac{\partial L}{\partial q_\lambda(\lambda)} = E_\mu[\log p(D \mid \lambda, \mu)] + E_\mu[\log(\mu \mid \lambda)] + E_\mu[\log p(\lambda)] - \log q(\lambda) = 0$$

故

$$\begin{aligned}\log q^*(\lambda) &= -\frac{\lambda}{2} E_{\mu} \left[\kappa_0 (\mu - \mu_0)^2 + \sum_{i=1}^m (x_i - \mu)^2 \right] + (a_0 - 1) \log \lambda - b_0 \lambda + \frac{m+1}{2} \log \lambda \\ &= \left(a_0 + \frac{m-1}{2} \right) \log \lambda - \left(b_0 + \frac{1}{2} E_{\mu} \left[\kappa_0 (\mu - \mu_0)^2 + \sum_{i=1}^m (x_i - \mu)^2 \right] \right) \lambda\end{aligned}$$

令 $a_m = a_0 + \frac{m+1}{2}$, $b_m = b_0 + \frac{1}{2} E_{\mu} \left[\kappa_0 (\mu - \mu_0)^2 + \sum_{i=1}^m (x_i - \mu)^2 \right]$,

可以看出: $q^*(\lambda) \sim \text{Gam}(\lambda \mid a_m, b_m)$.

故, 可以得到: $q^*(\mu, \lambda) \sim \mathcal{N}(\mu \mid N_m, \lambda_m^{-1}) \text{Gam}(\lambda \mid a_m, b_m)$

在无先验的情况下有:

$$\begin{cases} \mu_0 = a_0 = b_0 = \kappa_0 = 0 \\ E[\lambda] = \frac{a_m}{b_m} \\ E[\mu] = \mu_m = \bar{x} \\ E[\mu^2] = \bar{x}^2 + \frac{1}{mE[\lambda]}\end{cases}$$

可以解得:

$$\begin{aligned}E[\lambda] &= \frac{1}{E[x^2] - \bar{x}^2} = \frac{1}{\text{Var}(x)} \\ \lambda_m &= \frac{m}{\text{Var}(x)}\end{aligned}$$

$$b_m = b_0 + \frac{1}{2} E_{\mu} \left[\kappa_0 \mu^{-2} + \sum_{i=1}^m (x_i - \mu)^2 \right]$$

记带入后的 λ_m, b_m 为 λ', b' , 可以得到:

$$p(\mu, \lambda \mid D) \sim \mathcal{N}(\mu \mid \mu_m, \lambda') \text{Gam}(\lambda \mid am, b')$$

14.2

采用维特比算法, 首先求出位置 1 的各个标记 $j = 1, 2, \dots, m$ 的非规范化概率:

$$\delta_1(j) = \omega * F_1(y_0 = \text{start}, y_1 = j, x), \quad j = 1, 2, \dots, m$$

由递推公式, 求出位置 i 的各个标记 $l = 1, 2, \dots, m$ 的非规范化概率的最大值, 同时记录非规范化概率最大值的途径:

$$\delta_i(l) = \max_{1 \leq j \leq m} \{ \delta_{i-1}(j) + \omega * F_i(y_{i-1} = j, y_i = l, x) \}, \quad l = 1, 2, \dots, m$$

$$\Psi_i(l) = \arg \max_{1 \leq j \leq m} \{ \delta_{i-1}(j) + \omega * F_i(y_{i-1} = j, y_i = l, x) \}, \quad l = 1, 2, \dots, m$$

直到 $i = n$ 时终止, 此时求得非规范化概率的最大值为:

$$\max_y (\omega * F(y, x)) = \max_{1 \leq j \leq m} \delta_n(j)$$

以及最优路径的终点

$$y_n^* = \max_{1 \leq j \leq m} \delta_n(j)$$

由此最优路径终点返回

$$y_i^* = \Psi_{i+1}(y_{i+1}^*), \quad i = n-1, n-2, \dots, 1$$

最终求得最优路径.