

### 算法 7.1 (最小间隔法)

输入: 线性可分训练数据集  $T = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ , 其中  $x_i \in \mathbb{R}^n$ ,  $y_i \in \{-1, +1\}$ ,  $i=1, 2, \dots, n$ .

输出: 最大间隔的分离超平面和分类决策函数.

1) 构造并求解约束最优化问题:

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

$$\text{s.t. } y_i(w \cdot x_i + b) - 1 \geq 0, i=1, 2, \dots, N.$$

求得最优解  $w^*, b^*$ .

2) 由最优解构造分离超平面:

$$w^* \cdot x + b^* = 0.$$

$$\text{分类决策函数 } f(x) = \text{sign}(w^* \cdot x + b^*).$$

### 算法 7.4 (非线性支持向量机学习算法)

输出: 分类决策函数.

1) 选取适当的核函数  $K(x, y)$  和适当的参数  $C$  构造最优化问题:

$$\min_{\alpha} \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j K(x_i, x_j) - \frac{N}{2} C.$$

$$\text{s.t. } \sum_{i=1}^N \alpha_i y_i = 0.$$

$$0 \leq \alpha_i \leq C, i=1, 2, \dots, N.$$

求得最优解:  $\alpha^* = (\alpha_1^*, \alpha_2^*, \dots, \alpha_N^*)^T$ .

2) 选择  $x^*$  的  $n$  个分量  $0 < \alpha_i < C$  计算:

$$b = y_i - \sum_{j=1}^n \alpha_j y_j K(x_j, x_i).$$

3) 构造决策函数:

$$f(x) = \text{sign} \left( \sum_{i=1}^n \alpha_i y_i K(x, x_i) + b \right).$$

### 算法 7.2 (线性可分支持向量机学习算法)

1) 构造并求解约束最优化问题:

$$\min_{\alpha} \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j (x_i - x_j) \cdot (x_i - x_j)$$

$$\text{s.t. } \sum_{i=1}^N \alpha_i y_i = 0.$$

$$0 \leq \alpha_i \leq 1, i=1, 2, \dots, N.$$

求得最优解  $\alpha^* = (\alpha_1^*, \alpha_2^*, \dots, \alpha_N^*)^T$ .

2) 计算:

$$w^* = \sum_{i=1}^n \alpha_i y_i x_i.$$

并选择  $x^*$  的  $n$  个分量  $0 < \alpha_i < 1$  计算:

$$b = y_i - \sum_{j=1}^n \alpha_j y_j (x_i - x_j) \cdot (x_i - x_j).$$

3) 得到分离超平面:

$$w^* \cdot x + b = 0.$$

$$\text{分类决策函数 } f(x) = \text{sign}(w^* \cdot x + b).$$

### 算法 7.3 (取巧解法)

输入: 分离超平面和分类决策函数.

1) 选择适当的参数  $C > 0$ .

构造并求解二次规划问题:

$$\min_{\alpha} \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j (x_i - x_j) \cdot (x_i - x_j)$$

$$\text{s.t. } \sum_{i=1}^N \alpha_i y_i = 0.$$

$$0 \leq \alpha_i \leq C, i=1, 2, \dots, N.$$

求得最优解  $\alpha^* = (\alpha_1^*, \alpha_2^*, \dots, \alpha_N^*)^T$ .

2) 计算:

$$w^* = \sum_{i=1}^n \alpha_i y_i x_i.$$

选择  $x^*$  的一个分量  $\alpha_i$  对应该条件:

$$0 < \alpha_i < C, \text{ 计算:}$$

$$b = y_i - \sum_{j=1}^n \alpha_j y_j (x_i - x_j) \cdot (x_i - x_j).$$

3) ...

1. 只有第一个算法的最优化问题是“弱凸的”，其他两个都是“凸的”。
2. 只有最后一个算法是“支持向量机”，其他两个不是。
3. 第一个算法，其参数  $C$  第三、四个角是  $C$  的约束。
4. 第二个算法，只有  $C$  乘于的角，这也是为什么叫“支持向量机”。

### 感知机学习算法的基本形式:

输入: 训练数据集  $T = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ , 其中  $x_i \in \mathbb{R}^n$ ,  $y_i \in \{-1, +1\}$ ,  $i=1, 2, \dots, n$ ; 学习率  $\eta$ ,  $0 < \eta \leq 1$ .

输出:  $w, b$ ; 感知机模型  $f(x) = \text{sign}(w \cdot x + b)$ .

1) 选取初值  $w, b$ .

2) 在训练集中选取数据  $(x_i, y_i)$ .

3) 如果  $y_i(w \cdot x_i + b) \leq 0$ ,

$$w \leftarrow w + \eta y_i x_i$$

$$b \leftarrow b + \eta y_i$$

4) 转至 2), 直至训练集中没有误分类点.

### 对偶形式:

输出:  $\alpha, b$ ; 感知机模型  $f(x) = \text{sign}(\sum_{i=1}^N \alpha_i y_i x_i \cdot x + b)$ .

其中  $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_N)^T$ .

1)  $\alpha \leftarrow 0, b \leftarrow 0$ .

2) 在训练集中选取数据  $(x_i, y_i)$ .

3) 如果  $y_i(\sum_{j=1}^N \alpha_j y_j x_j \cdot x_i + b) \leq 0$ ,

$$\alpha_i \leftarrow \alpha_i + \eta$$

$$b \leftarrow b + \eta y_i.$$

4) 转至 2), 直至没有误分类数据.

感知机模型:  $f(x) = \text{sign}(w \cdot x + b)$ .

学习策略: 极大化损失函数 —  $\min_{w, b} L(w, b) = \sum_{i \in M} y_i(w \cdot x_i + b)$ .

梯度下降法不断极小化目标函数.

### 第四章 朴素贝叶斯法

4.1 基本方法:

先验概率分布  $P(Y = c_k) = p_k, k=1, 2, \dots, K$ .

条件概率分布  $P(X = x_1 | Y = c_k) = P(X^{(1)} = x^{(1)}, \dots, X^{(n)} = x^{(n)} | Y = c_k)$ .

参数个数为  $K \prod_{i=1}^n S_i$ .

条件独立性假设:  $P(X = x | Y = c_k) = P(X^{(1)} = x^{(1)}, \dots, X^{(n)} = x^{(n)} | Y = c_k)$ .

$$= \prod_{i=1}^n P(X^{(i)} = x^{(i)} | Y = c_k).$$

朴素贝叶斯分类器: 良的分类函数.

后验概率分布:  $P(Y = c_k | X = x) = \frac{P(X = x | Y = c_k) P(Y = c_k)}{\sum_{j=1}^K P(X = x | Y = c_j) P(Y = c_j)}$ .

朴素贝叶斯分类器:  $y = f(x) = \arg \max_{c_k} \frac{P(Y = c_k | X = x) P(Y = c_k)}{\sum_{j=1}^K P(Y = c_j | X = x) P(Y = c_j)}$ .

$$y = f(x) = \arg \max_{c_k} P(Y = c_k | X = x) P(Y = c_k)$$

4.1.2 后验概率最大化的含义:

选择 0-1 损失函数:  $L(Y, f(x)) = \begin{cases} 1, & Y \neq f(x) \\ 0, & Y = f(x) \end{cases}$ .

为使期望损失最小, 只需  $X = x$  使  $f(x)$  最小, 由此得到:

$$f(x) = \arg \min_{c_k} \sum_{i=1}^n L(c_k, y_i) P(c_k | x_i = x).$$

$$= \arg \min_{c_k} \sum_{i=1}^n (1 - P(c_k | x_i = x)) P(c_k | x_i = x).$$

$$= \arg \max_{c_k} P(Y = c_k | X = x).$$

4.2.1 极大似然估计:

先验概率极大似然估计:  $P(Y = c_k) = \frac{N_k}{N}$ .

条件概率:  $P(X = x | Y = c_k) = \frac{\sum_{i=1}^N I(x_i = x, y_i = c_k)}{\sum_{i=1}^N I(y_i = c_k)}$ .

## 算法4.1

输入: 训练数据  $T = \{x^{(1)}, y^{(1)}, x^{(2)}, y^{(2)}, \dots, x^{(N)}, y^{(N)}\}$  其中  $x_i = (x_i^1, x_i^2, \dots, x_i^m)^T$ ,  $y_i \in \{y_1, y_2, \dots, y_K\}$ .  $y_i$  是第  $i$  个特征可能取的第  $i$  个值,  $j=1, 2, \dots, K$ ,  $L=1, 2, \dots, Q$ ,  $y_i \in \{y_1, y_2, \dots, y_K\}$ .

输出: 实例  $x$  的分类.

1) 计算先验概率及条件概率

$$P(Y=j) = \frac{\sum_{i=1}^N I(y_i=j)}{N}, \quad j=1, 2, \dots, K.$$

$$P(X^k=l | Y=j) = \frac{\sum_{i=1}^N I(x_i^k=l, y_i=j)}{\sum_{i=1}^N I(y_i=j)}, \quad k=1, 2, \dots, m; \quad l=1, 2, \dots, Q.$$

$$j=1, 2, \dots, K; \quad l=1, 2, \dots, Q; \quad k=1, 2, \dots, m.$$

2) 对给定的实例  $x = (x^{(1)}, x^{(2)}, \dots, x^{(m)})^T$ , 计算

$$P(Y=j) \prod_{k=1}^m P(X^k=l_k | Y=j) | Y=j, \quad j=1, 2, \dots, K.$$

3) 确定实例  $x$  的分类

$$\hat{y} = \arg \max_{j=1, 2, \dots, K} P(Y=j) \prod_{k=1}^m P(X^k=l_k | Y=j) | Y=j.$$

学习朴素贝叶斯分类器并确定类标记  $y$ .  $(x, y, x, y, x, y)$

4.2.1. 贝叶斯估计:

$$\text{条件概率: } P(X^k=l | Y=j) = \frac{\sum_{i=1}^N I(x_i^k=l, y_i=j)}{\sum_{i=1}^N I(y_i=j) + \lambda}$$

$\lambda=1 \rightarrow$  拉普拉斯平滑

$$P(X^k=l | Y=j) = \frac{\sum_{i=1}^N I(x_i^k=l, y_i=j) + 1}{\sum_{i=1}^N I(y_i=j) + Q}$$

$$\text{得: } b^* = b^* = -\frac{1}{2} [(w^* \cdot (x_1 - x_2) + w^* \cdot (x_1 - x_3))]$$

$$\text{又 } w^* \cdot x_1 + b^* \geq 1 = w^* \cdot x_1 + b^*.$$

$$w^* \cdot x_1 + b^* \geq 1 = w^* \cdot x_2 + b^*$$

$$\therefore w^* \cdot (x_1 - x_2) = 0. \text{ 同理得 } w^* \cdot (x_1 - x_3) = 0. \Rightarrow b^* = b^*, \text{ 得证}$$

拉格朗日函数:

$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^N \alpha_i (w \cdot x_i + b) + \sum_{i=1}^N \alpha_i.$$

$$\min_{w, b} L(w, b, \alpha) = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j (x_i \cdot x_j) + \sum_{i=1}^N \alpha_i.$$

$$\max_{\alpha} -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j (x_i \cdot x_j) + \sum_{i=1}^N \alpha_i.$$

$$\text{对偶问题: } \min_{\alpha} \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j (x_i \cdot x_j) - \sum_{i=1}^N \alpha_i.$$

$$\left\{ \begin{aligned} w^* &= \sum_{i=1}^N \alpha_i^* x_i \\ b^* &= y_i - \sum_{i=1}^N \alpha_i^* (x_i \cdot x_i) \end{aligned} \right. \quad \text{--- 原始问题的解}$$

$$\text{分离超平面可以写成 } \sum_{i=1}^N \alpha_i y_i (x \cdot x_i) + b^* = 0.$$

$$\text{分类决策函数 } f(x) = \text{sign}(\sum_{i=1}^N \alpha_i y_i (x \cdot x_i) + b^*).$$

核函数:  $X$  是输入空间,  $H$  是特征空间, 存在映射  $\phi(x): X \rightarrow H$

$$\text{核函数 } K(x, z) = \langle \phi(x), \phi(z) \rangle.$$

$$\text{那么称 } K(x, z) \text{ 为核函数, } \phi(x) \text{ 为映射函数}$$

$$W(x) = \sum_{i=1}^N \alpha_i \phi(x) \cdot \phi(x_i) = \sum_{i=1}^N \alpha_i K(x, x_i).$$

$$\text{分类决策函数: } f(x) = \text{sign}(\sum_{i=1}^N \alpha_i y_i K(x, x_i) + b^*).$$

常用核函数:  $\left\{ \begin{aligned} &\text{多项式核函数: } K(x, z) = (x \cdot z + 1)^p. \\ &\text{高斯核函数: } K(x, z) = \exp(-\frac{\|x - z\|^2}{2\sigma^2}). \end{aligned} \right.$

最大间隔分离超平面的存在性.

若训练数据集  $T$  线性可分, 则可训练数据集  $T$  中的样本点完全由分离的超平面唯一确定.

证明: 1) 线性性.

2) 训练数据集线性可分

$$\text{最优化的问题 } \min_w \|w\|^2 \quad \text{--- 反证法可行解}$$

$$\text{s.t. } y_i(w \cdot x_i + b) - 1 \geq 0.$$

又: 目标函数有界.

3) 训练数据集线性可分, 记作  $(w^*, b^*)$ .

4) 训练数据集  $T$  中有线性点, 又有非线性点.

5)  $(w, b) = (0, 0)$  不是可行解

6)  $(w, b)$  以  $w^*$  满足  $w^* \neq 0$ , 故分离超平面存在.

7) 唯一性.

8) 首先证明  $w^*$  的唯一性. 假设没有两个最优解  $(w^*, b^*)$  和  $(w^*, b^*)$ .

$$\text{显然 } \|w^*\| = \|w^*\| = 0. \text{ 令 } w = \frac{w^* + w^*}{2}, \quad b = \frac{b^* + b^*}{2}.$$

易知  $(w, b)$  是可行解, 从而有:

$$0 \leq \|w\| \leq \frac{1}{2} (\|w^*\| + \|w^*\|) = 0.$$

由线性性可得:  $\|w\| = \frac{1}{2} (\|w^*\| + \|w^*\|) = 0$ , 从而  $w^* = w^* = 0$ .

1)  $\|x\| = 1$  若  $x = -1$ ,  $w = 0$ ,  $(w, b)$  不是解, 则  $1 \neq w \cdot x$ .

2) 设  $x_1$  和  $x_2$  是集合  $\{x_i | y_i = +1\}$  中分别对应于  $(w, b^*)$  和  $(w, b^*)$  使得  $y_i(w \cdot x_i + b) = 1$  或  $y_i(w \cdot x_i + b) = -1$  中分别对应于  $(w, b^*)$  和  $(w, b^*)$  使得  $\dots$  成立的点.

$$7.4. \quad g(w) = \sum_{i=1}^N y_i k(x_i, x) + b.$$

$$E = g(x) - y_i = \sum_{i=1}^N y_i k(x_i, x) + b - y_i.$$

SMD: 违反 KKT 条件最严重的样本点.

1) 有足够大的变化.

7.5. SMO 算法.

输入: 精度  $\epsilon$ .

输出: 近似解.

1) 取初值  $\alpha^{(0)} = 0$ ,  $k=0$ .

2) 选取优化变量  $\alpha^{(k)}, \alpha^{(k+1)}$ , 解析求解优化问题, 得到  $\alpha^{(k+1)}, \alpha^{(k+1)}$ .

更新  $\alpha^{(k+1)}$ .

3) 若在精度  $\epsilon$  范围内满足条件和

$$\sum_{i=1}^N \alpha_i y_i = 0.$$

$$0 \leq \alpha_i \leq C, \quad i=1, 2, \dots, N.$$

$$y_i g(x_i) = \begin{cases} \geq 1 & \text{if } x_i \cdot x_i = 0 \\ = 1 & \text{if } x_i \cdot x_i = C \\ \leq 1 & \text{if } x_i \cdot x_i = C \end{cases}$$

$$\text{其中, } g(x) = \sum_{i=1}^N y_i y_j (x_i \cdot x_j) + b.$$

则轻 (1): 否则  $k=k+1$ , 轻 (2);

(4) 取  $\alpha = \alpha^{(k+1)}$ .

例: 最小化二阶范数正则化的合页函数:  $\sum_{i=1}^N [1 - y_i (w \cdot x_i + b)] + \lambda \|w\|^2$ .

拉格朗日函数:  $L(w, b, \alpha, \mu)$

$$= \frac{1}{2} \|w\|^2 + 0 \sum_{i=1}^N \alpha_i - \sum_{i=1}^N \alpha_i [1 - y_i (w \cdot x_i + b) + \mu_i - 1] - \sum_{i=1}^N \mu_i$$

$$\text{解: } \min_w \max_{\alpha, \mu} \quad \text{求 } \max_{\alpha, \mu} \min_w$$

## 6.1 逻辑斯蒂分析

$$\text{几率函数 } F(x) = P(X=1) = \frac{1}{1 + e^{-w \cdot x}}$$

$$\begin{aligned} P(Y=1|x) &= \frac{\exp(w \cdot x + b)}{1 + \exp(w \cdot x + b)} \\ P(Y=0|x) &= \frac{1}{1 + \exp(w \cdot x + b)} \end{aligned}$$

将权值向量和输入向量并花

$$\begin{aligned} P(Y=1|x) &= \frac{\exp(w \cdot x)}{1 + \exp(w \cdot x)} \\ P(Y=0|x) &= \frac{1}{1 + \exp(w \cdot x)} = \frac{1}{1 + \exp(w \cdot x)} \end{aligned}$$

即事件发生与事件不发生概率之比

$$\text{对比模型而言 } \log \frac{P(Y=1|x)}{P(Y=0|x)} = w \cdot x$$

模型的参数估计：(应用极大似然估计法估计参数)

$$\text{设: } P(Y=1|x) = \pi(x), \quad P(Y=0|x) = 1 - \pi(x)$$

$$\text{似然函数: } f(x) = \prod_{i=1}^n [\pi(x_i)]^{y_i} [1 - \pi(x_i)]^{1-y_i}$$

$$\begin{aligned} \text{对数似然函数: } L(w) &= \sum_{i=1}^n [y_i \log \pi(x_i) + (1-y_i) \log (1-\pi(x_i))] \\ &= \sum_{i=1}^n [y_i \log \frac{\pi(x_i)}{1-\pi(x_i)} + \log (1-\pi(x_i))] \\ &= \sum_{i=1}^n [y_i \log (w \cdot x_i) - \log (1 + \exp(w \cdot x_i))] \end{aligned}$$

对  $L(w)$  求极大值得到

$$\text{则回归模型为 } P(Y=1|x) = \frac{\exp(w \cdot x)}{1 + \exp(w \cdot x)}$$

$$P(Y=0|x) = \frac{1}{1 + \exp(w \cdot x)}$$

$$\text{多类逻辑斯蒂回归模型: } \begin{cases} P(Y=1|x) = \frac{\exp(w_1 \cdot x)}{1 + \sum_{j=2}^K \exp(w_j \cdot x)} \\ P(Y=K|x) = \frac{1}{1 + \sum_{j=2}^K \exp(w_j \cdot x)} \end{cases}$$

## 6.2 最大熵模型

原理：在满足约束条件的模型集中选取熵最大的模型

X服从均匀分布时，熵最大

最大熵模型的意义：

首先根据  $P(X,Y)$  的经验分布和条件分布  $P(X|Y)$ 、 $P(Y|X)$

$$P(X=x, Y=y) = \frac{N_{xy}}{N}, \quad P(X=x) = \frac{N_{x+}}{N}$$

用特征函数  $\phi(x,y)$  描述输入  $x$  和输出  $y$  之间的某个事实

$$\phi(x,y) = \begin{cases} 1, & x, y \text{ 满足某事实} \\ 0, & \text{否则} \end{cases}$$

特征函数  $\phi(x,y)$  为经验分布  $P(X,Y)$  的期望值，用  $E\phi$  表示

$$E\phi(x,y) = \sum_{x,y} \phi(x,y) P(X,Y)$$

特征函数  $\phi(x,y)$  为模型  $P(X,Y)$  的期望值，用  $E_P\phi$  表示

$$E_P\phi(x,y) = \sum_{x,y} \phi(x,y) P(x,y)$$

假设两个期望值相等：  $E\phi = E_P\phi$  为模型学习的约束条件

设：  $C = \{P \in \mathcal{P} \mid E_P\phi_i = E\phi_i, i=1, \dots, n\}$

定义在条件概率分布  $P(X,Y)$  上的条件熵

$$H(P) = -\sum_{x,y} P(x,y) \log P(x,y)$$

则模型集合  $C$  中条件熵  $H(P)$  最大的模型即为最大熵模型

最大熵模型的学习：

等价于求解最优化问题：

$$\max_{P \in C} H(P) = -\sum_{x,y} P(x,y) \log P(x,y)$$

$$\text{s.t. } E_P\phi_i = E\phi_i, i=1, \dots, n$$

$$\sum_{x,y} P(x,y) = 1$$

改写成最小值问题：

$$\min_{P \in C} -H(P) = \sum_{x,y} P(x,y) \log P(x,y)$$

例：为约平果和化了的对偶问题：

① 设  $x$  拉格朗日函数  $L(P, w)$ ：

$$\begin{aligned} L(P, w) &= -H(P) + w_0 (1 - \sum_{x,y} P(x,y)) + \sum_{i=1}^n w_i (E_P\phi_i - E\phi_i) \\ &= -\sum_{x,y} P(x,y) \log P(x,y) + w_0 (1 - \sum_{x,y} P(x,y)) \\ &\quad + \sum_{i=1}^n w_i (\sum_{x,y} \phi_i(x,y) P(x,y) - \sum_{x,y} \phi_i(x,y) f_i(x,y)) \end{aligned}$$

② 原始问题转化为对偶问题：

$$\min_{P \in C} \max_w L(P, w) \rightarrow \max_w \min_{P \in C} L(P, w)$$

③ 解对偶问题：

$$a. \text{ 记 } \psi(w) = \min_{P \in C} L(P, w) = L(P_w, w)$$

$$b. \text{ 将 } \psi(w) \text{ 的解记为 } P_w = \arg \min_{P \in C} L(P, w) = P_w(y|x)$$

c. 对  $L(P_w, w)$  求  $P_w(y|x)$  的导数

$$\begin{aligned} \frac{\partial L(P_w, w)}{\partial P_w(y|x)} &= \sum_{x,y} P_w(x,y) (\log P_w(y|x) + 1) - \sum_{i=1}^n w_i \sum_{x,y} \phi_i(x,y) P_w(y|x) \\ &= \sum_{x,y} P_w(x,y) (\log P_w(y|x) + 1 - w_0 - \sum_{i=1}^n w_i \phi_i(x,y)) \\ &\stackrel{\text{令}=0}{=} \sum_{x,y} P_w(x,y) = \frac{\exp(\sum_{i=1}^n w_i \phi_i(x,y))}{\exp(1 - w_0)} \end{aligned}$$

$$\text{d. 由上化出得: } P_w(y|x) = \frac{1}{Z_w(x)} \exp(\sum_{i=1}^n w_i \phi_i(x,y))$$

$$\text{其中 } Z_w(x) = \sum_{y} \exp(\sum_{i=1}^n w_i \phi_i(x,y))$$

④ 解对偶问题等价于极大化问题：

$$\max_w \psi(w) \Rightarrow w^* = \arg \max_w \psi(w)$$

极大似然估计：

证明：对偶函数的极大化  $\Rightarrow$  最大熵模型的极大似然估计

条件概率分布  $P_w(y|x)$  的对数似然函数表示为：

$$L(P_w) = \log \prod_{x,y} P_w(y|x)^{N_{xy}} = \sum_{x,y} N_{xy} \log P_w(y|x)$$

当条件概率分布  $P_w(y|x)$  是  $P_w^*$  时，对数似然函数  $L(P_w)$  为

$$L(P_w) = \sum_{x,y} P_w^*(x,y) \log P_w^*(y|x)$$

又对偶函数  $\psi(w)$ ：

$$\begin{aligned} \psi(w) &= \sum_{x,y} P_w^*(x,y) \log P_w^*(y|x) + \sum_{i=1}^n w_i (\sum_{x,y} \phi_i(x,y) P_w^*(y|x) - \sum_{x,y} \phi_i(x,y) f_i(x,y)) \\ &= \sum_{x,y} P_w^*(x,y) \log P_w^*(y|x) - \sum_{i=1}^n w_i \log Z_w(x) \end{aligned}$$

$$\text{且 } \sum_{x,y} P_w^*(x,y) = 1, \text{ 则得: } \psi(w) = L(P_w)$$

最大熵模型：  $P_w(y|x) = \frac{1}{Z_w(x)} \exp(\sum_{i=1}^n w_i \phi_i(x,y))$

$$Z_w(x) = \sum_{y} \exp(\sum_{i=1}^n w_i \phi_i(x,y))$$

