

Logistics 回归

上一页笔记提到两种模型：生成模型与鉴别模型

生成模型：假设每个类别有一个专属语言模型，将文本 d 放入每个模型中跑，找到生成该文本概率最大的模型，输出它的类别 C

核心动作：计算 C, d 联合概率 $P(C, d)$

核心公式：
$$C_{\max} = \underset{C \in C}{\operatorname{argmax}} P(d|C) P(C)$$

鉴别模型：直接计算 $P(C|d)$

二元 Logistics 回归：典型的鉴别模型

给定输入输出训练集 $(x^{(i)}, y^{(i)})$ 。 i 表示第 i 个样本， $x^{(i)}$ 表示这个样本的输入特征（如特征向量）， $y^{(i)}$ 表示其真实类别，在二元逻辑回归里 y 只有两个值：0（负例），1（正例）

对于每个 $x^{(i)}$ ，模型会先将其处理为特征向量 $[x_1, \dots, x_n]$ 然后计算输出预测类别 $\hat{y}^{(i)}$

具体计算：逻辑回归是一个线性模型，会先计算一个线性

$$\begin{aligned} \text{得分: } z^{(i)} &= w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b \\ &= W \cdot x^{(i)} + b \end{aligned}$$

$W = [w_1, \dots, w_n]$ 是给每个特征分配一个权重，表示该特征的重要程度

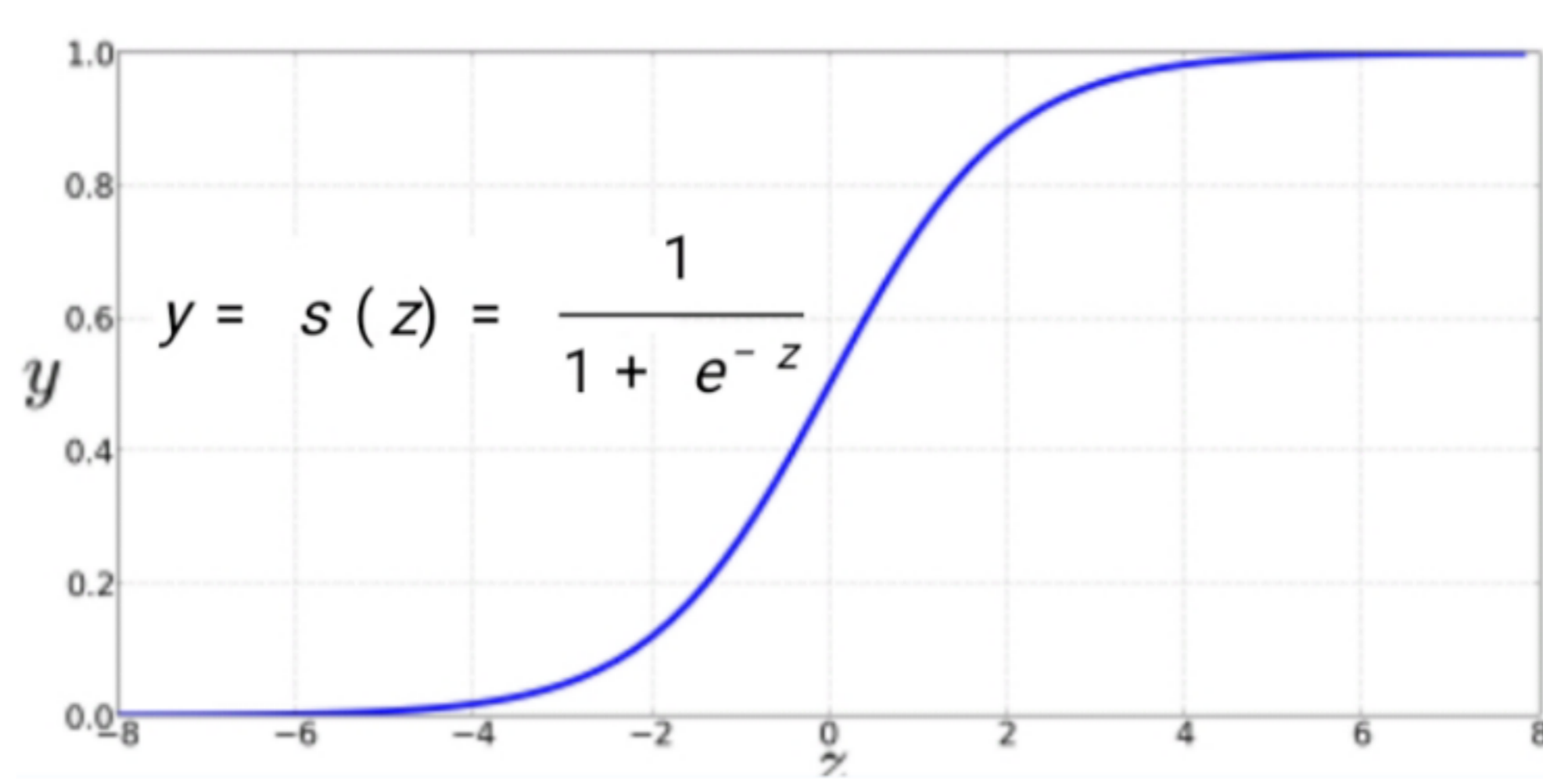
b 是偏置（截距）， z 是线性加权

如果 z 大于规定值，则 $y=1$ ，否则为 0

归一化：Sigmoid 激活函数，将 z 转化为概率，更利于将其用于做决策

$$y = S(z) = \frac{1}{1 + e^{-z}} = \frac{1}{1 + \exp(-z)}$$

$y \rightarrow 1$ 表示模型极度确信是正类



$s(z)$ 也写作 $\sigma(w \cdot x + b)$

Sigmoid 函数计算概率

$$\text{正类 } y=1 \text{ 的概率: } P(y=1|x) = \frac{1}{1 + \exp(-(w \cdot x + b))}$$

$$\begin{aligned} \text{负类 } y=0 \text{ 的概率: } P(y=0|x) &= 1 - P(y=1|x) \\ &= \frac{1}{1 + \exp(w \cdot x + b)} \end{aligned}$$

$$\text{决策: } \hat{y} = \begin{cases} 1 & \text{if } P(y=1|x) > 0.5 \\ 0 & \text{otherwise} \end{cases}$$

0.5 被称为决策边界

Logistics 回归学习

模型: 观测特征向量: x

类别集合: C

输出: $\hat{y} = \sigma(w \cdot x + b)$

真实类别: y

参数学习: 设置参数 w 和 b 最小化 $\hat{y}^{(i)}$ 与 $y^{(i)}$ 的距离 (损失, $L(\hat{y}, y)$)

似然概率: $P(y|x) = \hat{y}^y \cdot (1 - \hat{y})^{1-y}$

这是单个样本的似然函数, 衡量的是预测结果与

真实类别的匹配程度, 值越接近 1 说明模型越准

当 $y=1$ $P(y|x) = \hat{y} \rightarrow$ 模型预测为正类的概率

当 $y=0$ $P(y|x) = 1 - \hat{y} \rightarrow$ 模型预测为负类的概率

我们的目的: 最大化 $\log P(y|x)$

$$\begin{aligned}\log P(y|x) &= \log [\hat{y}^y (1-\hat{y})^{1-y}] \\ &= y \log \hat{y} + (1-y) \log (1-\hat{y})\end{aligned}$$

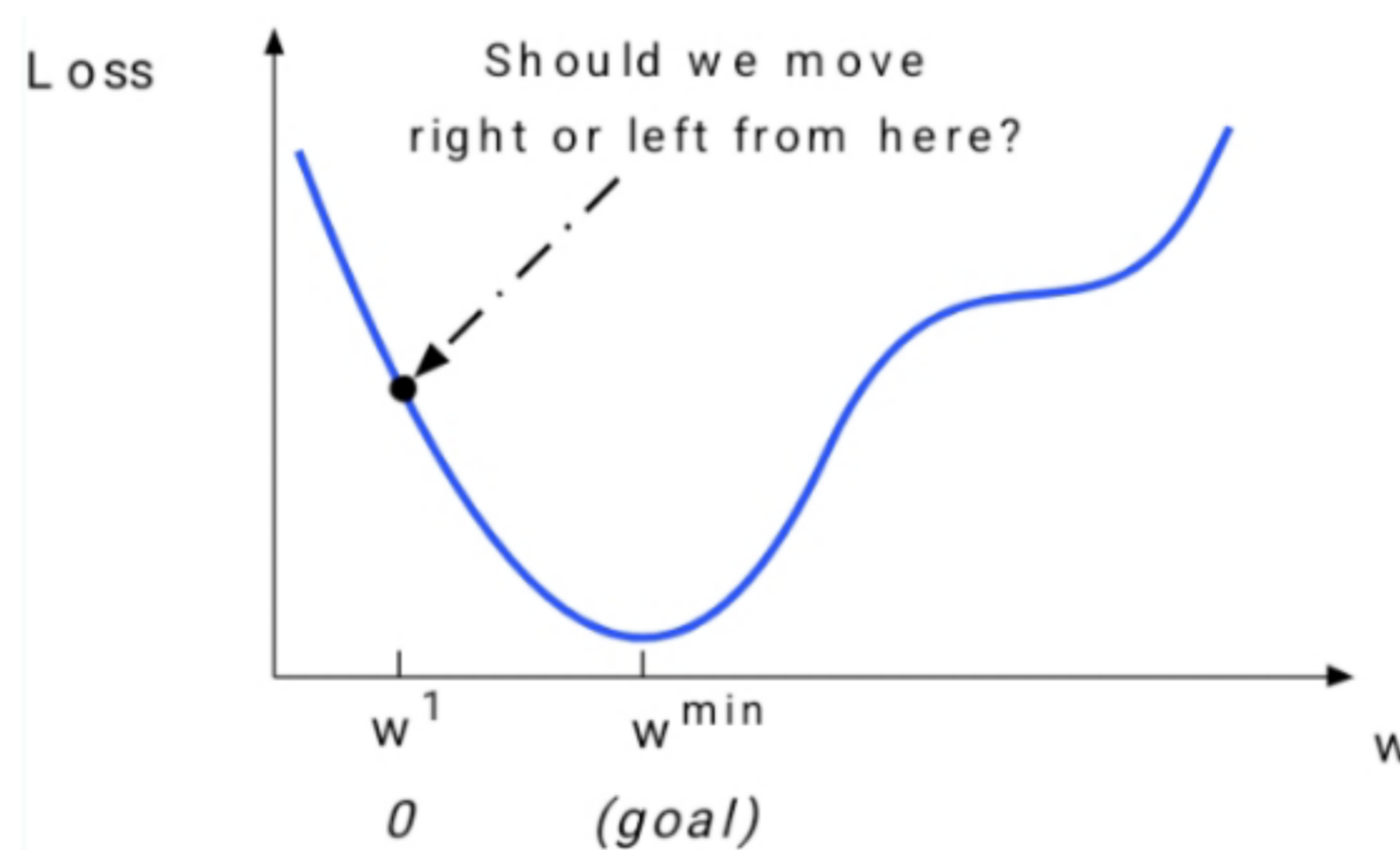
交叉熵损失: $L_{CE}(\hat{y}, y) = -\log P(y|x) = -[y \log \hat{y} + (1-y) \log (1-\hat{y})]$

最大化对数似然 \rightarrow 最小化交叉熵损失

梯度下降: 参数集合表示为 $\theta = (w, b)$, $\hat{y} = f(x; \theta)$

最小化参数使所有样本平均损失最小化:

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^m L_{CE}(f(x^{(i)}; \theta), y^{(i)})$$



$$w^{t+1} = w^t - \eta \frac{d}{dw} \hat{\theta}$$

$$\text{梯度: } \frac{d}{dw} L(f(x; w), y)$$

步长: η

随机梯度下降 (SGD): 每次迭代随机选 1 个样本计算梯度而非用全量样本 (用全量样本叫批量梯度下降, 用部分样本的叫小批量梯度下降)

