

# 《机器学习》课程系列

判别函数的线性分类基础\*

武汉纺织大学数学与计算机学院

杜小勤

2019/02/26

## Contents

<b>1</b>	<b>概述</b>	<b>2</b>
<b>2</b>	<b>分类</b>	<b>4</b>
2.1	二分类基础 . . . . .	4
2.2	多分类基础 . . . . .	6
<b>3</b>	<b>基本判别函数</b>	<b>9</b>
3.1	最小平方分类方法 . . . . .	9
3.2	Fisher 二分类方法 . . . . .	11
3.3	最小平方分类方法与 Fisher 分类方法 . . . . .	15
3.4	Fisher 多分类方法 . . . . .	16
3.5	感知机 . . . . .	18
3.5.1	学习策略 . . . . .	19
3.5.2	学习的原始算法 . . . . .	22
3.5.3	算法的收敛性 . . . . .	24
3.5.4	学习的对偶算法 . . . . .	26

---

\*本系列文档属于讲义性质，仅用于学习目的。Last updated on: July 10, 2021。

4 附录	28
4.1 标量函数对矩阵的求导 . . . . .	28
4.2 公式证明一 . . . . .	39
4.3 公式证明二 . . . . .	41
4.4 感知机实验 . . . . .	44
5 参考文献	51

## 1 概述

分类问题指的是, 将输入变量  $\boldsymbol{x}$  划分到  $K$  个离散类别  $c_k$  中的某一个, 其中  $k = 1, \dots, K$ 。一般情况下, 这些类别互不相交, 每个输入变量  $\boldsymbol{x}$  被划分到唯一的一个类别中。

从分类的观点看, 在输入空间, 这些输入变量  $\boldsymbol{x}$  位于不同的决策区域 (Decision Region), 它们的边界被称为决策边界 (Decision Boundary) 或者决策面 (Decision Surface)。

在本章, 我们将讨论线性分类的判别函数 (Discriminant Function)。所谓线性分类的判别函数, 指的是利用判别函数直接将样本点  $\boldsymbol{x}$  映射到某个类别  $c_k (k = 1, \dots, K)$ , 判别函数也被称为决策函数。其决策面是关于输入变量  $\boldsymbol{x}$  的线性函数。例如, 在  $D$  维输入空间中, 决策面是  $D - 1$  维的超平面 (Hyperplane)<sup>1</sup>。分类判别函数有别于分类的概率生成式模型与概率判别式模型这 2 种方法, 稍后会作一个简单的对比介绍。

如果数据集可以被线性决策面精确地分类, 那么称该数据集是线性可分的 (Linearly Separable), 否则称数据集是线性不可分的 (Linearly Inseparable)。

在监督的分类问题中, 目标变量  $y$  是一个实值向量或标量, 用来表示类别标签 (Class Label)。例如, 在二分类问题中,  $y$  是一个标量, 其取值为  $\{c_1, c_2\}$ 。但是, 为了便于数学或算法表达, 一般将取值限定为  $\{0, 1\}$  或  $\{-1, 1\}$ 。而在多分类问题中, 为简化数学公式, 一般使用 2 种常用的表示方法: 标量形式,  $y$  的取值为  $\{c_1, c_2, \dots, c_K\}$ , 其中  $K$  表示类别个数; One-Hot 向量形式, 例如, 在手写体数字识别问题中,  $K = 10$ , 数字 3 和 4 被分别编码成  $y = (0, 0, 0, 1, 0, 0, 0, 0, 0, 0)$

<sup>1</sup>例如, 在具有 2 个特征分量的二维平面中, 线性决策面是一维直线或线段。

和  $y = (0, 0, 0, 0, 1, 0, 0, 0, 0, 0)$ ，其余数字的编码，依此类推。在后文，为了描述方便，在不引起歧义的情况下，公式中将混用这两种形式：使用“ $y = c_k$ ”描述“ $y$  属于类别  $c_k$ ”这样的事实，虽然  $y$  同时也是一个 One-Hot 向量。因此，在多分类问题中，不会使用粗体字表示 One-Hot 向量  $y$ 。

针对分类问题，有三种不同的解决方案，依照复杂度从高到低，它们分别是：

- 概率生成式方法

首先，对类的条件概率分布  $p(\mathbf{x}|c_k)$  及类的先验概率分布  $p(c_k)$  建模；然后使用贝叶斯定理计算后验概率分布：

$$p(c_k|\mathbf{x}) = \frac{p(\mathbf{x}|c_k)p(c_k)}{p(\mathbf{x})} \quad (1)$$

最后，使用训练集最大化求解该后验概率分布（即最大后验估计），得到最优的模型参数；

- 概率判别式方法

直接对后验概率分布  $p(c_k|\mathbf{x})$  建模，例如把该分布表示为参数模型，然后使用训练集最大化求解该参数模型，得到最优的模型参数；

- 判别函数或决策函数

该方法直接把样本点  $\mathbf{x}$  映射到具体的类别中；

在线性回归模型中，“线性”指的是，模型预测函数  $y(\mathbf{x}, \mathbf{w})$  是关于模型参数  $\mathbf{w}$  的线性函数<sup>2</sup>。在最简单的情况下，线性模型  $y(\mathbf{x}, \mathbf{w})$  也是关于输入变量  $\mathbf{x}$  的线性函数，即  $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$ （为了简化符号，省略了模型参数  $\mathbf{w}$ ）。

而在分类问题中，输入变量  $\mathbf{x}$  的预测值是离散类别标签的后验概率分布。此时，为了完成分类任务，可以考虑将最简单的线性模型  $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$  进行扩展——在线性变换的基础上，添加一个非线性变换  $f$ ：

$$y(\mathbf{x}) = f(\mathbf{w}^T \mathbf{x} + w_0) \quad (2)$$

其中，函数  $f$  在机器学习领域被称为激活函数（Activation Function），其反函数在统计学领域被称为连接函数（Link Function）。

<sup>2</sup>在多数情况下，线性模型  $y(\mathbf{x}, \mathbf{w})$  与  $\mathbf{x}$  之间的关系可能是非线性关系。

如果使用公式 (2) 进行分类决策, 那么  $y(\mathbf{x})$  的决策面是  $\mathbf{w}^T \mathbf{x} + w_0$ , 这意味着决策面仍然是线性决策面。因此, 不论  $f$  是否线性函数, 决策面始终是  $\mathbf{x}$  的线性函数, 公式 (2) 被称为广义线性模型 (Generalized Linear Model, GLM)<sup>3</sup>。

但是, 需要注意的是, 与线性回归模型中的情况不一样,  $y(\mathbf{x})$  不再是模型参数  $\mathbf{w}$  的线性模型, 因为  $f$  是一个非线性函数。这将导致线性分类模型的计算比线性回归模型的计算要复杂些。

与线性基函数的回归模型类似, 在后面的内容中, 我们也将引入基函数到分类模型 (公式2) 中, 即使用基函数向量  $\phi(\mathbf{x})$  替换输入向量  $\mathbf{x}$ :

$$y(\mathbf{x}) = f(\mathbf{w}^T \phi(\mathbf{x}) + w_0) \quad (3)$$

在上面的公式中, 先使用非线性基函数  $\phi$  对输入向量  $\mathbf{x}$  进行变换, 然后再执行常规线性分类。这意味着, 分类模型不再是  $\mathbf{x}$  的线性分类模型, 但其功能更加灵活。这是引入基函数之后带来的好处。需要注意的是, 下面讨论的算法同样适用于基于线性基函数的分类模型。

## 2 分类

### 2.1 二分类基础

首先, 考虑二分类模型, 然后将它扩展到  $K > 2$  的分类模型。

最简单的二分类判别函数具有如下形式:

$$y(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + w_0) \quad (4)$$

其中,  $\mathbf{w}$  被称为权值向量 (Weight Vector),  $w_0$  被称为偏置 (Bias),  $\text{sign}$  表示符号函数<sup>4</sup>:

$$\text{sign}(z) = \begin{cases} +1, & z \geq 0 \\ -1, & z < 0 \end{cases} \quad (5)$$

<sup>3</sup>关于广义线性模型更复杂的定义及其应用, 请阅读相关文献。

<sup>4</sup>可以把符号函数  $\text{sign}$  看作是一种阈值函数, 当输入变量值大于等于某个阈值 (例如, 本例中阈值为 0) 时, 函数的输出值为 +1; 否则, 函数的输出值为 -1。在一些情况下, 为处理方便, 也可以允许相应的函数输出值为 1 和 0。

因此，当  $\mathbf{w}^T \mathbf{x} + w_0 \geq 0$  时， $y(\mathbf{x}) = +1$ ，表示  $\mathbf{x}$  被划分到类别 +1 中；否则， $y(\mathbf{x}) = -1$ ，表示  $\mathbf{x}$  被划分到类别 -1 中。有时候，偏置  $w_0$  的相反数被称为阈值 (Threshold)，因为：当  $\mathbf{w}^T \mathbf{x} \geq -w_0$  时， $y(\mathbf{x}) = +1$ ；否则， $y(\mathbf{x}) = -1$ 。

在公式 (4) 中，(线性) 决策边界或决策面由  $\mathbf{w}^T \mathbf{x} + w_0 = 0$  确定，它是一个超平面。其中， $\mathbf{w}$  是该超平面的“法向量”，即向量  $\mathbf{w}$  与决策面上的任何向量都正交：假设  $\mathbf{x}_1$  和  $\mathbf{x}_2$  都位于决策面上，那么有  $\mathbf{w}^T \mathbf{x}_1 + w_0 = 0$  和  $\mathbf{w}^T \mathbf{x}_2 + w_0 = 0$ ，从而有  $\mathbf{w}^T (\mathbf{x}_1 - \mathbf{x}_2) = 0$ 。因此， $\mathbf{w}$  决定了决策面的方向。

在明确了向量  $\mathbf{w}$  的意义后，我们来看一下，如何计算原点到决策面的带符号垂直距离。假设  $\mathbf{x}$  是决策面上的点，那么有  $\mathbf{w}^T \mathbf{x} + w_0 = 0$ ，则：

$$\begin{aligned}\mathbf{w}^T \mathbf{x} &= -w_0 \\ \frac{\mathbf{w}^T \mathbf{x}}{\|\mathbf{w}\|} &= -\frac{w_0}{\|\mathbf{w}\|} \\ \frac{\mathbf{w}^T}{\|\mathbf{w}\|} \mathbf{x} &= -\frac{w_0}{\|\mathbf{w}\|}\end{aligned}\tag{6}$$

在上面的第 3 行公式中， $\frac{\mathbf{w}^T}{\|\mathbf{w}\|}$  表示  $\mathbf{w}$  的单位向量， $\frac{\mathbf{w}^T}{\|\mathbf{w}\|} \mathbf{x}$  表示  $\mathbf{w}$  的单位向量与向量  $\mathbf{x}$  的点积 (内积)，其含义就是向量  $\mathbf{x}$  在决策面的单位“法向量”方向上的带符号长度。由于向量  $\mathbf{x}$  与原点之间的特殊关系，则原点到决策面的带符号距离就是  $\frac{w_0}{\|\mathbf{w}\|}$ <sup>5</sup>。因此，偏置  $w_0$  决定了决策面的位置。

实际上，对于任意向量  $\mathbf{x}$  而言，公式  $\frac{\mathbf{w}^T \mathbf{x} + w_0}{\|\mathbf{w}\|}$  给出了该向量到决策面的带符号距离。为了说明这一点，考虑任意一点  $\mathbf{x}$  在决策面上的投影  $\mathbf{x}_\perp$ ，如图2-1所示，可得：

$$\begin{aligned}\mathbf{x} &= \mathbf{x}_\perp + r \frac{\mathbf{w}}{\|\mathbf{w}\|} \\ \mathbf{w}^T \mathbf{x} &= \mathbf{w}^T \mathbf{x}_\perp + r \frac{\mathbf{w}^T \mathbf{w}}{\|\mathbf{w}\|} \\ \mathbf{w}^T \mathbf{x} + w_0 &= \mathbf{w}^T \mathbf{x}_\perp + w_0 + r \frac{\mathbf{w}^T \mathbf{w}}{\|\mathbf{w}\|} \\ r &= \frac{\mathbf{w}^T \mathbf{x} + w_0}{\|\mathbf{w}\|}\end{aligned}\tag{7}$$

其中，上述第 3 行推导到第 4 行，利用了公式  $\mathbf{w}^T \mathbf{x}_\perp + w_0 = 0$  和  $\mathbf{w}^T \mathbf{w} = \|\mathbf{w}\|^2$ 。

<sup>5</sup>注意，向量  $\mathbf{x}$  在  $\mathbf{w}$  的单位“法向量”上的带符号长度，与原点到决策面的带符号距离，刚好是相反的。更一般地，设决策面为  $\mathbf{w}^T \mathbf{x} + w_0$ ，对该决策面按法向量  $\mathbf{w}$  进行归一化处理： $\frac{\mathbf{w}^T \mathbf{x}}{\|\mathbf{w}\|} + \frac{w_0}{\|\mathbf{w}\|}$ ，这就是任意点  $\mathbf{x}$  到该决策面的带符号距离。例如，将原点代入，得到  $\frac{w_0}{\|\mathbf{w}\|}$ ，这与前面的推导结果一致。

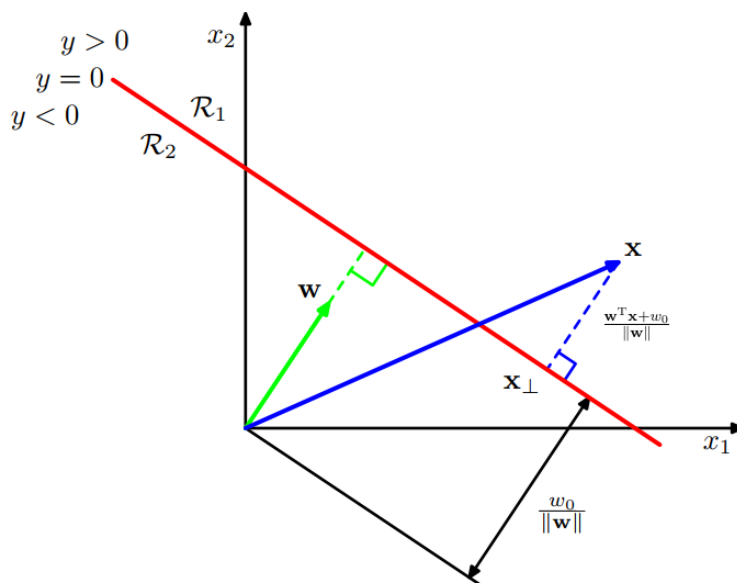


图 2-1: 二维线性判别函数中的一维决策面（一条线）

有时候，为了简化公式，将偏置  $w_0$  并入到权值向量  $w$  中，然后在输入向量  $x$  中引入一个对应的输入 1，仍然使用  $w$  和  $x$  分别表示扩展之后的权值向量与输入向量，可得：

$$y(x) = \text{sign}(w^T x) \quad (8)$$

本节中，将会混用这 2 种表示方法。

## 2.2 多分类基础

对于多分类问题，可以尝试组合前面的二分类判别函数，即使用多个二分类判别函数来构造一个  $K$  类判别函数。

首先，可以想到 2 种较为直接的从二分类判别函数构造  $K$  类判别函数的方法。第 1 种方法是，使用  $K - 1$  个二分类器，每个二分类器解决一个二分类问题：把属于类别  $c_k$  和不属于类别  $c_k$  的点分开，这种分类器被称为一对多 (One-Versus-The-Rest) 分类器。然而，这种方法会产生无法分类的区域<sup>6</sup>，如图 2-2 所示，绿色区域表示无法分类的区域。第 2 种方法是，使用  $\frac{K(K-1)}{2}$  个二分类器，即为每对类别设置一个二分类器，而每个点的类别将依据多数表决方式来确定：考虑所有判别函数的分类结果，选择表决最多的分类作为最后的分类结果，

<sup>6</sup>这种区域被称为歧义性区域 (Ambiguous Region)。

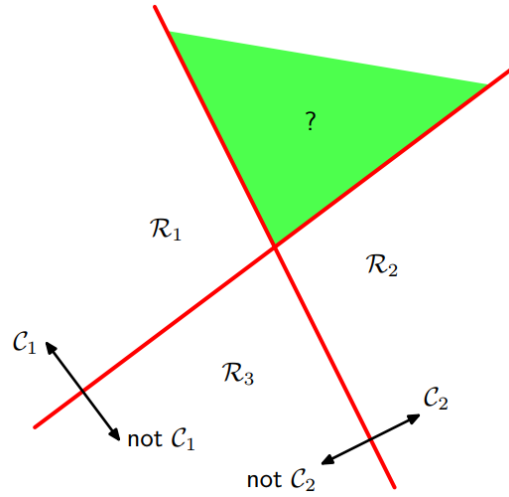


图 2-2: 一对多分类器存在的问题

这种分类器被称为一对一（One-Versus-One）分类器。但是，这也会造成无法分类的区域，如图2-3所示。

实际上，存在一种简单的解决方案：引入  $K$  个最简形式的决策面来解决  $K$  类分类问题，从而形成一个  $K$  类别判别函数<sup>7</sup>：

$$\begin{cases} y_1(\mathbf{x}) = \mathbf{w}_1^T \mathbf{x} + w_{10}, & c_1 \\ y_2(\mathbf{x}) = \mathbf{w}_2^T \mathbf{x} + w_{20}, & c_2 \\ \dots \\ y_K(\mathbf{x}) = \mathbf{w}_K^T \mathbf{x} + w_{K0}, & c_K \end{cases} \quad (9)$$

其分类规则如下：对每个待分类的点  $\mathbf{x}$ ，如果  $y_k(\mathbf{x}) > y_j(\mathbf{x})$  for all  $j \neq k$  成立，那么  $\mathbf{x}$  将被划分为  $c_k$  类别。

显然，任意 2 个类别  $c_k$  和  $c_j$  之间的决策面存在如下的关系：

$$\begin{aligned} y_k(\mathbf{x}) &= y_j(\mathbf{x}) = 0 \\ (\mathbf{w}_k - \mathbf{w}_j)^T \mathbf{x} + (w_{k0} - w_{j0}) &= 0 \end{aligned} \quad (10)$$

这种  $K$  类判别函数在输入空间上划分形成的决策区域具有如下特点：

- 决策区域之间的决策面是超平面；

<sup>7</sup>例如，用于分类的神经网络也采取类似的解决方案，它为每个类别设计一个输出。训练时，为每个输入  $\mathbf{x}$  配置一个理想输出  $y$ ，整体输出则采取 One-Hot 编码形式。分类预测时，神经网络在所有的  $K$  个输出中，选取输出值最高的输出作为整体的分类结果。

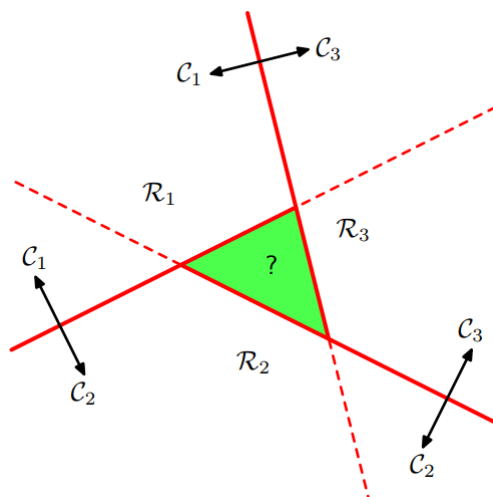
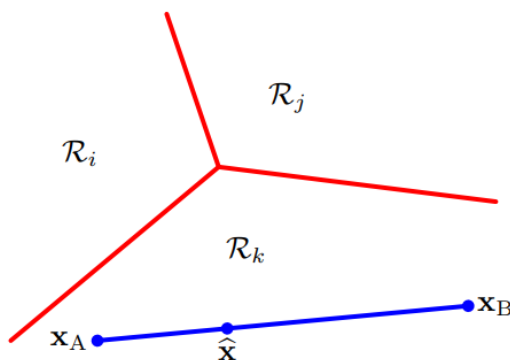


图 2-3: 一对一分类器存在的问题

图 2-4:  $K$  类判别函数形成的凸单连通决策区域

- 决策区域总是单连通的凸区域；

为了表明这个结论是正确的，假设在决策区域  $\mathcal{R}_k$  中有 2 个点  $\mathbf{x}_A$  和  $\mathbf{x}_B$ ，如图2-4所示。那么，任何位于  $\mathbf{x}_A$  和  $\mathbf{x}_B$  连线上的点都可以表示成下面的形式：

$$\hat{\mathbf{x}} = \lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B \quad (11)$$

其中  $0 \leq \lambda \leq 1$ 。根据  $K$  类判别函数的线性性质，下式成立：

$$y_k(\hat{\mathbf{x}}) = \lambda y_k(\mathbf{x}_A) + (1 - \lambda) y_k(\mathbf{x}_B) \quad (12)$$

由于  $\mathbf{x}_A$  和  $\mathbf{x}_B$  位于  $\mathcal{R}_k$  的内部，则  $y_k(\mathbf{x}_A) > y_j(\mathbf{x}_A)$  和  $y_k(\mathbf{x}_B) > y_j(\mathbf{x}_B)$ ，则可得  $y_k(\hat{\mathbf{x}}) > y_j(\hat{\mathbf{x}})$ 。同理，可得  $y_k(\hat{\mathbf{x}}) > y_i(\hat{\mathbf{x}})$ 。这表明  $\hat{\mathbf{x}}$  也在区域  $\mathcal{R}_k$  的内部，即  $\mathcal{R}_k$  是凸单连通区域。



有了判别函数，下面将具体讨论使用训练数据对参数进行调整的方法：最小平方方法、Fisher 线性判别函数与感知机算法。

### 3 基本判别函数

#### 3.1 最小平方分类方法

对于简单的二分类情形，既可以使用  $K$  类别判别函数（其中  $K = 2$ ），也可以使用第2.1节中介绍的单判别函数。因此，下面以通用的  $K$  类别判别函数进行讨论。

对于  $K$  类判别函数，每个类  $c_k$  有自己的线性模型：

$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0} \quad (13)$$

为了简化公式，可以将  $w_{k0}$  并入向量  $\mathbf{w}_k$  中，并在输入变量  $\mathbf{x}$  中增加一个虚输入  $x_0 = 1$ ，仍然使用原向量符号  $\mathbf{w}$  和  $\mathbf{x}$  来表示扩展后的向量，即：

$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} \quad (14)$$

其中， $\mathbf{w}_k = (w_{k0}, w_{k1}, w_{k2}, \dots)^T$  和  $\mathbf{x} = (1, x_1, x_2, \dots)^T$ 。

下面将  $K$  个线性判别函数使用矩阵来表示：

$$y(\mathbf{x}) = \mathbf{W}^T \mathbf{x} \quad (15)$$

其中， $\mathbf{W}$  是一个矩阵，第  $k$  列是向量  $\mathbf{w}_k$ 。在这种表示方式下，新输入  $\mathbf{x}$  将被划分到  $y_k = \mathbf{w}_k^T \mathbf{x}$  取得最大值的类别  $c_k$  中<sup>8</sup>。

考虑训练数据集  $\{\mathbf{x}_n, y_n\}$ ， $n = 1, \dots, N$ ，定义矩阵  $\mathbf{Y}$ ，它的第  $n$  行是  $y_n$  的 One-Hot 向量；定义矩阵  $\mathbf{X}$ ，它的第  $n$  行是向量  $\mathbf{x}_n^T$ 。因此，该训练数据集的平方误差函数可以表示成：

$$E_D(\mathbf{W}) = \frac{1}{2} \text{Tr} \{ (\mathbf{X}\mathbf{W} - \mathbf{Y})^T (\mathbf{X}\mathbf{W} - \mathbf{Y}) \} \quad (16)$$

然后，令平方误差函数关于  $\mathbf{W}$  的导数等于 0，可以得到  $\mathbf{W}$  的解析解：

$$\mathbf{W} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = \mathbf{X}^\dagger \mathbf{Y} \quad (17)$$

<sup>8</sup>注意，虽然  $y$  表示多类别，但是仍然使用非粗体字来表示。原因在于，为公式表达的简便性， $y = c_k$  及其 One-Hot 向量形式会混用，前文已经讨论过。

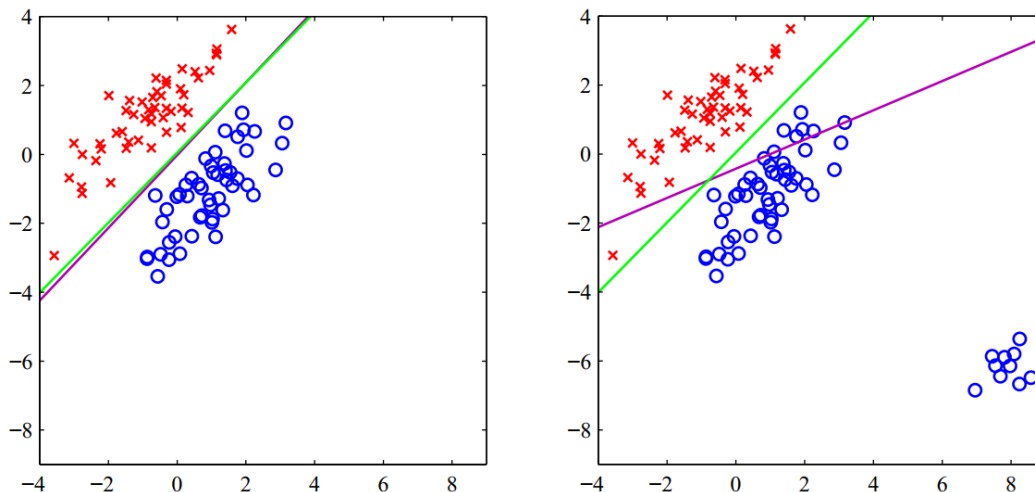


图 3-5: 最小平方方法对离群点很敏感

其中,  $\mathbf{X}^\dagger$  被称为  $\mathbf{X}$  的 Moore-Penrose 伪逆矩阵 (Pseudo-Inverse Matrix), 形式为:

$$\mathbf{X}^\dagger \equiv (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \quad (18)$$

因此, 得到了如下形式的  $K$  类判别函数:

$$y(\mathbf{x}) = \mathbf{W}^T \mathbf{x} = \mathbf{Y}^T (\mathbf{X}^\dagger)^T \mathbf{x} \quad (19)$$

值得注意的是, 如果训练集里每个目标向量  $y_n$  都满足某个线性限制:

$$\mathbf{a}^T y_n + b = 0 \quad (20)$$

其中  $\mathbf{a}$  和  $b$  分别为常向量和常量, 那么对于任何输入变量  $\mathbf{x}$ , 模型的预测值也满足同样的限制:

$$\mathbf{a}^T y(\mathbf{x}) + b = 0 \quad (21)$$

因此, 在目标向量  $y_n$  使用 “1-of-K” 或 One-Hot 形式时, 相应的模型预测值将具有如下性质: 对于任意的输入变量  $\mathbf{x}$ ,  $y(\mathbf{x})$  的所有分量之和等于 1。但是, 需要注意, 这并不意味着模型的预测输出表现为概率形式, 因为预测值并没有限制在  $0-1$  范围内。

虽然最小平方方法给出了  $K$  类判别函数的精确解析解, 但是, 同回归中的问题一样, 最小平方方法对于离群点缺少鲁棒性, 如图3-5所示。图中, 洋红色直线表示最小平方方法找到的决策边界, 绿色直线表示逻辑回归模型找到的决策边界。可以

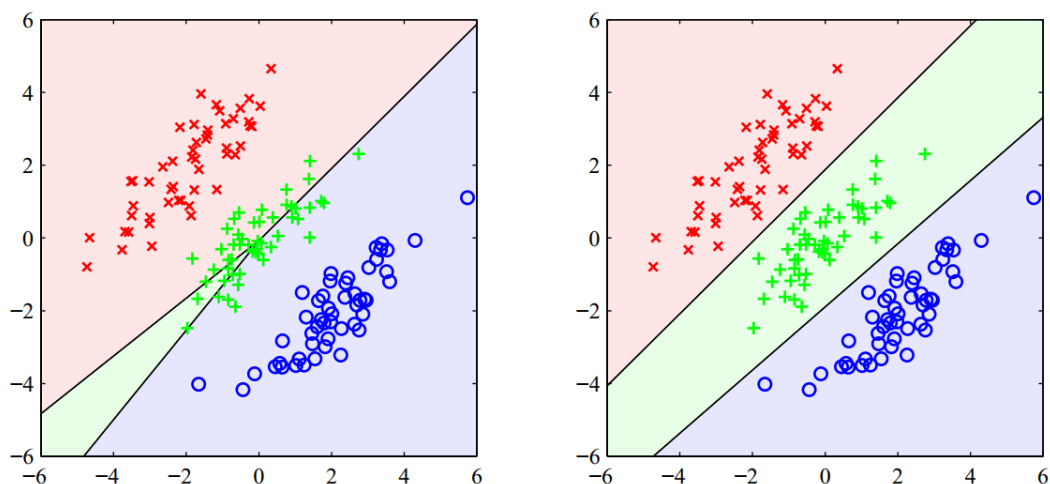


图 3-6: 最小平方法不能正确地进行多分类的例子

看出，最小平方法对离群点非常敏感，而逻辑回归模型却具有很好的鲁棒性。这是由最小平方法的性质造成的，它会惩罚那些远离决策边界（即过于正确）的预测。后面将会考虑另外几种误差函数，它们不会出现这种问题。

然而，最小平方法还存在更加严重的问题，如图3-6所示。该训练数据集由三个类别组成，本身是线性可分的。左图是最小平方法分类的结果，右图是逻辑回归分类的结果。显然，最小平方法没有产生正确的分类，而逻辑回归给出了正确的分类。

最小平方法无法产生正确的分类，其根本原因在于，这种方法对应于高斯条件分布假设下的最大似然法，而  $y_n$  中的每个分量是二值  $\{0, 1\}$  形式，它显然不服从高斯分布。因此，可以使用更加恰当的概率分布并转换为相应的损失函数来解决这个问题。

### 3.2 Fisher 二分类方法

在前面的线性分类判别方法中，无论是二分类问题，还是多分类问题，都需要定义线性变换函数  $\mathbf{w}^T \mathbf{x}$ ，它的作用是形成一个决策边界，从而对输入空间产生一个划分。

从维度降低的角度来看，它的作用是将一个  $D$  维输入向量  $\mathbf{x}$  经过权值向量  $\mathbf{w}$  投影到 1 维上。对于二分类问题，如果投影的结果（标量）大于等于 0（阈值），那么该输入变量  $\mathbf{x}$  将被划分到类别  $c_1$ ，否则将被划分到类别  $c_2$ 。对于多分类问

题，其分类规则为：对每个待划分的点  $\mathbf{x}$ ，如果  $y_k(\mathbf{x}) > y_j(\mathbf{x}), \forall j \neq k$  成立，那么  $\mathbf{x}$  将被划分到  $c_k$  类别，其效果相当于确定了多个阈值，并依据阈值对输入空间进行划分。

通常来说，向 1 维进行投影会造成相当多的信息损失，并使得原本能够在  $D$  维输入空间线性可分的训练数据集，在投影到 1 维空间后造成重叠，从而变得线性不可分。

但是，假设依然要将  $D$  维输入空间投影到 1 维空间上，我们可以选取类别之间分离效果最好的投影，即寻求最佳的投影权值向量，这就是 Fisher 线性分类器的基本思想。

为了推导出 Fisher 线性分类器，首先考虑二分类情形，然后推广到  $K$  分类情形。假设  $c_1$  类中存在  $N_1$  个点， $c_2$  类中存在  $N_2$  个点。那么，这 2 个类别的均值向量分别为：

$$\mathbf{m}_1 = \frac{1}{N_1} \sum_{n \in c_1} \mathbf{x}_n, \quad \mathbf{m}_2 = \frac{1}{N_2} \sum_{n \in c_2} \mathbf{x}_n \quad (22)$$

它们在原输入空间中的均值向量差是  $\mathbf{m}_2 - \mathbf{m}_1$ ，经过  $\mathbf{w}$  投影后，在 1 维空间中的类别距离是：

$$m_2 - m_1 = \mathbf{w}^T (\mathbf{m}_2 - \mathbf{m}_1) \quad (23)$$

其中：

$$m_k = \mathbf{w}^T \mathbf{m}_k \quad (24)$$

$m_k$  是来自类别  $c_k$  的投影均值。公式 (23) 表明，原输入空间中的均值向量差的投影，等于两个类别均值向量分别投影后的差。显然，这是一种非常简单的类别距离度量方式。为了获得最佳的类别间分离效果，可以关于权值向量  $\mathbf{w}$  最大化求解公式 (23)。

但是，需要注意的是，必须对权值向量  $\mathbf{w}$  施加限制，否则，权值向量  $\mathbf{w}$  可以无限制地增大。因此，可以将权值向量  $\mathbf{w}$  限制为单位向量，即  $\sum_i w_i^2 = 1$ ——问题转化为带约束的拉格朗日优化问题，可以使用拉格朗日乘数法来求解最优权值向量  $\mathbf{w}$ ，并得到结果： $\mathbf{w} \propto (\mathbf{m}_2 - \mathbf{m}_1)$ 。

这个结论很有趣，它表明，最佳权值向量  $\mathbf{w}$  与均值向量差  $\mathbf{m}_2 - \mathbf{m}_1$  的方向一致，且分离超平面垂直于均值向量差，分类情形如图3-7所示。可以看出，稍长的绿色直线表示分离超平面，稍短的绿色直线表示均值向量差，直方图形式可以直观地表示类别均值向量投影之后的距离  $m_2 - m_1$  的最大化——要最大化投影之

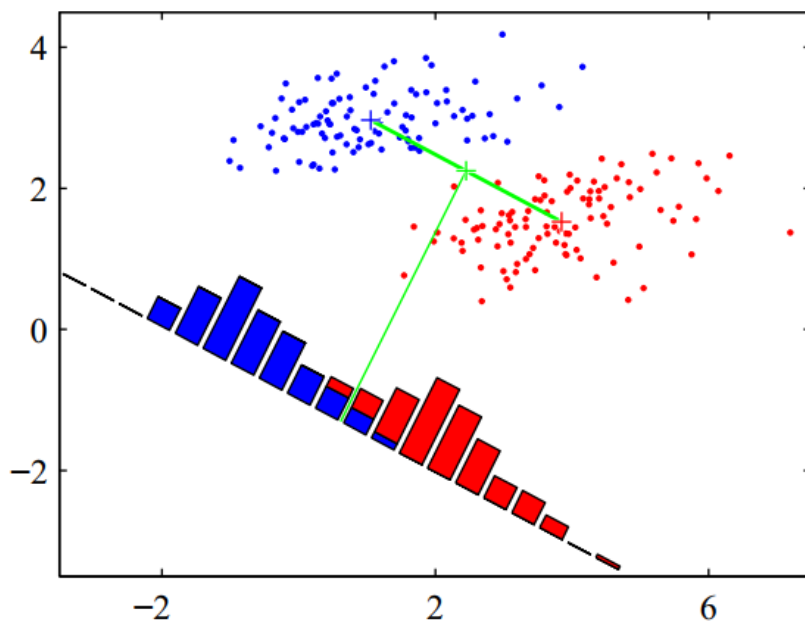


图 3-7: 类别均值投影的距离最大化方法存在的问题

后的均值距离，就要选择这样的分离超平面，它垂直于均值向量差，且经过 2 个均值的中点——这的确是一种非常简单的类别距离度量方式。

然而，该方法存在一个问题：原本 2 个类别在原始空间中是线性可分的，但是在投影之后，2 个类别之间就存在一定程度的重叠！显然，问题出在投影向量  $\mathbf{w}$  上，需要对它进行调整。

如何调整投影向量  $\mathbf{w}$  呢？注意到，如果类别分布的非对角协方差较大，那么这种问题就会出现。这给了我们一种改进算法的思路——寻找的投影向量  $\mathbf{w}$ ，既要最大化类别间的投影距离，又要最小化类内的方差，这样可以确保找到较为完美的分离超平面，它既能够最大限度地分离输入点，同时又使得类别重叠的部分最小化。这便是 Fisher 线性分类器的设计思路。

首先，定义  $c_k$  的类内方差 (Within-Class Variance)：

$$s_k^2 = \sum_{n \in c_k} (y_n - m_k)^2 \quad (25)$$

其中， $y_n = \mathbf{w}^T \mathbf{x}_n$ ， $m_k$  是来自类别  $c_k$  的投影均值。

然后，定义总的类内方差为  $s_1^2 + s_2^2$ ，定义 Fisher 准则为：

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2} \quad (26)$$

它表示类间方差 (Between-Class Variance) 与类内方差的比值。使用前面的相关公式替换其中的  $m_1$ 、 $m_2$ 、 $s_1$  和  $s_2$ ，可以得到关于  $w$  的公式：

$$J(w) = \frac{w^T S_B w}{w^T S_W w} \quad (27)$$

其中， $S_B$  表示类间协方差矩阵：

$$S_B = (m_2 - m_1)(m_2 - m_1)^T \quad (28)$$

$S_W$  表示类内协方差矩阵：

$$S_W = \sum_{n \in c_1} (x_n - m_1)(x_n - m_1)^T + \sum_{n \in c_2} (x_n - m_2)(x_n - m_2)^T \quad (29)$$

对公式 (27) 关于  $w$  最大化，得到<sup>9</sup>：

$$(w^T S_B w) S_W w = (w^T S_W w) S_B w \quad (30)$$

从公式 (28) 可以看出：

$$\begin{aligned} S_B w &= (m_2 - m_1)(m_2 - m_1)^T w \\ &= c(m_2 - m_1) \end{aligned} \quad (31)$$

其中  $c = (m_2 - m_1)^T w$  表示某个常量。上式表明， $S_B w$  总是在向量  $(m_2 - m_1)$  的方向上。此外，我们仅仅需要  $w$  的方向，而不关心它的长度，因此，可以将标量因子  $(w^T S_B w)$  和  $(w^T S_W w)$  从公式 (30) 中去掉，再使用公式 (31) 进行相应的替换（也去掉常量因子），最后公式 (30) 两端同时乘以  $S_W^{-1}$ ，从而得到：

$$w \propto S_W^{-1} (m_2 - m_1) \quad (32)$$

注意，改进方法比标准方法多了一个矩阵变换  $S_W^{-1}$ ，即类内协方差矩阵的逆变换，它的意义也非常明确——除了考虑类间距离最大化之外，还显式地考虑了类内协方差的最小化<sup>10</sup>，这是改进算法的优势。

值得注意的是，如果类内协方差矩阵  $S_W$  是各向同性的 (Isotropic)，那么它与单位矩阵成正比，此时  $w$  就只决定于类间均值向量差  $(m_2 - m_1)$ ，这与标准方法得到的结果一样。

<sup>9</sup>关于公式 (30) 的证明，请阅读附录第4.2节。

<sup>10</sup>类内协方差逆的最大化，相当于类内协方差的最小化。

公式 (32) 被称为 Fisher 线性判别式。但是, 严格来说, 它并不是一个判别式, 而是样本数据点投影方向的一个特定选择。然而, 将数据向  $\mathbf{w}$  投影之后, 确实可以构造出一个判别式用来将输入  $\mathbf{x}$  进行分类: 如果  $y(\mathbf{x}) \geq y_0$ , 那么  $\mathbf{x}$  属于类别  $c_1$ ; 否则,  $\mathbf{x}$  属于类别  $c_2$ 。

如何确定最优阈值  $y_0$  呢? 我们可以假设类条件概率密度  $p(y|c_k)$  服从高斯分布<sup>11</sup>, 然后通过最大似然法求解出高斯分布的参数。在对每个类别求解了高斯分布参数之后, 可以通过最小化误分类率的方式获得最优阈值<sup>12</sup>。

### 3.3 最小平方分类方法与 Fisher 分类方法

最小平方方法的训练准则是, 尽可能地让模型预测与理想值接近。而 Fisher 方法的训练准则是, 尽可能地拉开类别间的距离, 而缩小类别内的距离。

可以看出, 它们之间存在相似性——对于最小平方方法而言, 它需要确定出最佳的决策面, 以确保误分类率最小。而从分类的投影解释来看, 选取最佳决策面, 也是为了尽可能地拉开类别间的距离。因此, 直觉上来看, 这 2 种方法之间必然存在着某种联系。

下面, 我们将表明, 对于二分类问题, 如果选取适当的目标向量编码方式, 那么这 2 种方法实际上是等价的。

在前面的讨论中, 我们已经知道, 对于二分类目标值, 可以使用  $\{c_1, c_2\}$  这样的编码方式, 例如  $\{0, 1\}$  或  $\{-1, 1\}$ ; 而对于多分类目标值, 可以使用 1-of-K 或 One-Hot 编码方式。

现在, 我们使用另一种编码方式。令类别  $c_1 = \frac{N}{N_1}$ , 类别  $c_2 = -\frac{N}{N_2}$ , 其中  $N_1$  和  $N_2$  分别表示训练数据集中类别  $c_1$  和  $c_2$  的个数。

最小平方方法使用如下的误差函数:

$$E = \frac{1}{2} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - y_n)^2 \quad (33)$$

然后, 求解误差函数  $E$  关于参数  $w_0$  和  $\mathbf{w}$  的偏导数, 并令其分别等于 0 和  $\mathbf{0}$ , 可

<sup>11</sup>注意到,  $y = \mathbf{w}^T \mathbf{x}$  是一组随机变量的和, 根据中心极限定理, 我们可以做出高斯分布的假设。

<sup>12</sup>详细方法, 请参考文献 [2], P39, 1.5.1 节。

得：

$$\begin{aligned}\frac{\partial E}{\partial w_0} &= \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - y_n) = 0 \\ \frac{\partial E}{\partial \mathbf{w}} &= \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - y_n) \mathbf{x}_n = \mathbf{0}\end{aligned}\quad (34)$$

求解  $w_0$ ，得到：

$$w_0 = -\mathbf{w}^T \mathbf{m} \quad (35)$$

其中，我们使用了如下公式：

$$\sum_{n=1}^N y_n = N_1 \frac{N}{N_1} - N_2 \frac{N}{N_2} = 0 \quad (36)$$

$\mathbf{m}$  表示整个数据集的均值向量：

$$\mathbf{m} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n = \frac{1}{N} (N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2) \quad (37)$$

求解  $\mathbf{w}$ ，可得<sup>13</sup>：

$$\left( \mathbf{S}_W + \frac{N_1 N_2}{N} \mathbf{S}_B \right) \mathbf{w} = N (\mathbf{m}_1 - \mathbf{m}_2) \quad (38)$$

其中， $\mathbf{S}_W$  和  $\mathbf{S}_B$  分别表示类内协方差矩阵和类间协方差矩阵，与 Fisher 方法中的定义一样。注意到， $\mathbf{S}_B \mathbf{w}$  总是在向量  $(\mathbf{m}_2 - \mathbf{m}_1)$  的方向上，于是：

$$\mathbf{w} \propto \mathbf{S}_W^{-1} (\mathbf{m}_2 - \mathbf{m}_1) \quad (39)$$

该公式与 Fisher 方法中得到的结果一致。

考虑到  $w_0 = -\mathbf{w}^T \mathbf{m}$ ，因此，对于待预测输入向量  $\mathbf{x}$ ，当  $y(\mathbf{x}) = \mathbf{w}^T (\mathbf{x} - \mathbf{m}) > 0$  时， $\mathbf{x}$  属于类别  $c_1$ ，否则属于类别  $c_2$ 。

### 3.4 Fisher 多分类方法

下面将 Fisher 二分类线性判别函数泛化到多分类线性判别函数，此处的多分类指的是类别个数  $K > 2$ 。设输入向量  $\mathbf{x}$  的维度为  $D$ ，并引入输出向量  $\mathbf{y}$ <sup>14</sup>，其维度为  $D'$ ，它的每个元素  $y_k = \mathbf{w}_k^T \mathbf{x}$ ， $k = 1, \dots, D'$ 。为了方便起见，将使用矩阵形式表示：

$$\mathbf{y} = \mathbf{W}^T \mathbf{x} \quad (40)$$

<sup>13</sup>关于公式 (38) 的证明，请阅读附录第4.3节。

<sup>14</sup>由于前述原因，不使用粗体字表示它。



其中，权值矩阵  $\mathbf{W}$  的列是权值向量  $\mathbf{w}_k$ 。

为了得到 Fisher 多分类判别函数，需要仿照二分类情形，定义相应的投影前后的协方差矩阵：输入空间中的类内协方差矩阵  $\mathbf{S}_W$  和类间协方差矩阵  $\mathbf{S}_B$ 、输出空间（投影空间）中的类内协方差矩阵  $\mathbf{s}_W$  和类间协方差矩阵  $\mathbf{s}_B$ ，并基于上述协方差矩阵，定义多分类情形的 Fisher 优化准则  $J(\mathbf{W})$ ——类间方差与类内方差的比值，它是一个标量。 $J(\mathbf{W})$  的最大化，既确保了最大化类间方差，又确保了最小化类内方差，因而确保了能够找到较为完美的分离超平面。

首先，定义输入空间中的  $K$  类别类内协方差矩阵  $\mathbf{S}_W$ ，定义相当直接：

$$\mathbf{S}_W = \sum_{k=1}^K \mathbf{S}_k \quad (41)$$

其中， $\mathbf{S}_k$  定义为：

$$\begin{aligned} \mathbf{S}_k &= \sum_{n \in c_k} (\mathbf{x}_n - \mathbf{m}_k)(\mathbf{x}_n - \mathbf{m}_k)^T \\ \mathbf{m}_k &= \frac{1}{N_k} \sum_{n \in c_k} \mathbf{x}_n \end{aligned} \quad (42)$$

其中， $N_k$  是类别  $c_k$  的样本点个数。

输入空间中的  $K$  类别类间协方差矩阵  $\mathbf{S}_B$  的定义不是很直接，需要借助于全协方差矩阵：

$$\mathbf{S}_T = \sum_{n=1}^N (\mathbf{x}_n - \mathbf{m})(\mathbf{x}_n - \mathbf{m})^T \quad (43)$$

其中， $\mathbf{m}$  是整个数据集的均值向量：

$$\mathbf{m} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n = \frac{1}{N} \sum_{k=1}^K N_k \mathbf{m}_k \quad (44)$$

其中， $N = \sum_k N_k$  是整个数据集的样本点个数。有了全协方差矩阵和类内协方差矩阵，就可以按如下方式定义类间协方差矩阵：

$$\mathbf{S}_T = \mathbf{S}_W + \mathbf{S}_B \quad (45)$$

从而得到：

$$\mathbf{S}_B = \sum_{k=1}^K N_k (\mathbf{m}_k - \mathbf{m})(\mathbf{m}_k - \mathbf{m})^T \quad (46)$$

下面，在投影空间中定义相应的协方差矩阵。类似地，类内协方差矩阵定义为：

$$\mathbf{s}_W = \sum_{k=1}^K \sum_{n \in c_k} (y_n - \boldsymbol{\mu}_k)(y_n - \boldsymbol{\mu}_k)^T \quad (47)$$

类间协方差矩阵定义为：

$$\mathbf{s}_B = \sum_{k=1}^K N_k (\boldsymbol{\mu}_k - \boldsymbol{\mu})(\boldsymbol{\mu}_k - \boldsymbol{\mu})^T \quad (48)$$

其中：

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n \in c_k} y_n, \quad \boldsymbol{\mu} = \frac{1}{N} \sum_{k=1}^K N_k \boldsymbol{\mu}_k \quad (49)$$

有了  $K$  类别投影空间中的类间和类内协方差矩阵，仿照二分类情形，定义如下的 Fisher 最优化准则：

$$J(\mathbf{W}) = \text{Tr} \{ \mathbf{s}_W^{-1} \mathbf{s}_B \} \quad (50)$$

最大化该公式将确保同时最大化类间方差与最小化类内方差。与二分类情形类似，使用上述的权值矩阵  $\mathbf{W}$ 、输入空间中的类内协方差矩阵  $\mathbf{S}_W$  和类间协方差矩阵  $\mathbf{S}_B$  进行替换，可以得到如下形式的公式：

$$J(\mathbf{W}) = \text{Tr} \{ (\mathbf{W} \mathbf{S}_W \mathbf{W}^T)^{-1} (\mathbf{W} \mathbf{S}_B \mathbf{W}^T) \} \quad (51)$$

最大化  $J(\mathbf{W})$  将获得最优权值矩阵  $\mathbf{W}$ ，它是  $\mathbf{S}_W^{-1} \mathbf{S}_B$  的前  $D'$  个（最大特征值对应的）特征向量。

值得注意的是，观察公式 (46)，可以发现， $\mathbf{S}_B$  是由  $K$  个矩阵叠加而成的，而每个矩阵的秩为 1，<sup>15</sup>再由公式 (44) 可知， $\mathbf{S}_B$  的  $K$  个外积（矩阵）只有  $K-1$  个是相互独立的，故矩阵  $\mathbf{S}_B$  的秩至多是  $K-1$ ，其非零特征值至多只有  $K-1$  个。这表明，向  $\mathbf{S}_B$  的特征向量张成的  $K-1$  维子空间进行投影，并不会改变  $J(\mathbf{W})$  的值，因而投影空间的维度  $D'$  不会超过  $K-1$ 。

### 3.5 感知机

感知机 (Perceptron) 于 1957 年由 Rosenblatt 提出，虽然其形式与学习算法都非常简单，也不是第 1 个神经元模型<sup>16</sup>，但是它在模式识别与机器学习的发展历史上占有重要的地位。

<sup>15</sup>注意到，每个矩阵都是由 2 个向量的外积得到的，这种矩阵的秩为 1。原因是，无论是列向量，还是行向量，都可以看作是特殊的矩阵，它们的秩为 1，设列向量为  $A$ ，行向量为  $B$ ，则外积  $AB$  是一个矩阵，其秩满足下面的关系： $R(AB) \leq R(A)R(B) = 1$ 。

<sup>16</sup>第 1 个神经元模型是 McCulloch-Pitts 模型 (M-P 模型，1943 年提出)，一般被认为是神经网络的开端之作，它使用 Sigmoid 或 Tanh 作为激活函数。实际上，感知机是一种基于 M-P 结构的模型，但它使用了不连续函数 Sign 作为激活函数。虽然感知机的思想并不新颖，但是 Rosenblatt 使用计算机程序展示了感知机的功能，并通过简单的实验验证了感知机模型具有一定的泛化能力。

感知机的出现，标志着研究者们开始真正地对学习过程进行数学理论上的研究。1962 年，Novikoff 证明了关于感知机的第 1 个定理，这标志着学习理论研究的开始。该定理在创建学习理论的过程中起了十分重要的作用。从某种意义上说，它建立了学习算法的泛化能力与训练集误差的最小化之间的联系。

感知机的决策函数为：

$$f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + b) \quad (52)$$

其中， $\mathbf{w}$  被称为权值向量 (Weight Vector)； $b$  被称为偏置 (Bias)； $\text{sign}$  为符号函数：

$$\text{sign}(z) = \begin{cases} +1 & z \geq 0 \\ -1 & z < 0 \end{cases} \quad (53)$$

可以看出，感知机是一种二分类判别模型，其输入为样本点的特征向量  $\mathbf{x}$ ；输出  $y$  分为 2 类：为算法实现上的便利，取值分别为  $+1$  和  $-1$ 。通过第 2.1 部分的讨论，可知，权值向量  $\mathbf{w}$  用于控制决策面  $\mathbf{w}^T \mathbf{x} + b = 0$  的方向，而偏置  $b$  用于控制决策面的位置。感知机模型的假设空间是输入特征空间  $\mathbf{x}$  上的所有线性分类模型，即  $\{f | f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b\}$ 。

与 Fisher 分类方法不同，感知机直接以决策面  $\mathbf{w}^T \mathbf{x} + b = 0$  为界，将输入空间划分为 2 部分： $\mathbf{w}$  所指向的一边，类别为  $+1$ ，称为正类；反之，类别为  $-1$ ，称为负类。而 Fisher 分类方法首先需要将数据样本点进行降维，并在确定最佳投影方向之后，还需要找到最佳分界点，最后完成对样本点的分类。另外，无论从历史发展的角度看，还是从技术视角上看，神经网络和 SVM 都是感知机模型的一种延伸。

### 3.5.1 学习策略

首先，给出数据集的线性可分性定义。

定义 3.1 (数据集的线性可分性) 给定一个数据集

$$T = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\},$$

其中  $\mathbf{x}_i \in \mathbf{R}^n$ ,  $y_i = \{+1, -1\}$ ,  $i = 1, 2, \dots, N$ ，如果存在某个决策面  $\mathbf{w}^T \mathbf{x} + b = 0$ ，能够将数据集的正负样本点完全正确地划分到决策面的两侧，即对所有的  $y_i = +1$  样本点，有  $\mathbf{w}^T \mathbf{x}_i + b > 0$ ，对所有的  $y_j = -1$  样本点，有  $\mathbf{w}^T \mathbf{x}_j + b < 0$ ，则称数

数据集  $T$  是线性可分的 (*Linearly Separable*); 否则, 称数据集  $T$  是线性不可分的 (*Linearly Inseparable*)。

如何从数据集  $T$  中学习权值向量  $\mathbf{w}$  和偏置  $b$ ? 对于学习算法而言, 一般要定义一个目标函数。常用的目标函数, 包括效用函数 (Agent 或智能体)、值函数 (强化学习)、损失函数或负对数似然函数 (监督学习) 等。对于监督学习, 常常选取损失函数或负对数似然函数作为目标函数, 并执行最优化过程, 用于从数据集中学习权值参数  $\mathbf{w}$  和偏置  $b$ 。

对于损失函数的选取, 需要考虑以下几个重要因素:

- 能够准确地反映误分类及其代价;
- 便于最优化地求取参数: 凸函数; 连续可导函数;

一种很自然的想法是, 对于单个样本点而言, 正确分类时, 损失为 0; 错误分类时, 损失为 1, 即得到如下的 (单样本点)0-1 损失函数:

$$l(y_i, f(\mathbf{x}_i)) = \begin{cases} 1 & y_i \neq f(\mathbf{x}_i) \\ 0 & y_i = f(\mathbf{x}_i) \end{cases} \quad (54)$$

其中,  $f(\mathbf{x}_i) = \text{sign}(\mathbf{w}^T \mathbf{x}_i + b)$ 。利用  $y$  和  $f(\mathbf{x}_i)$  取值为  $\{+1, -1\}$  的特点, 将上述损失函数改写为:

$$l(y_i, f(\mathbf{x}_i)) = \begin{cases} 1 & y_i f(\mathbf{x}_i) < 0 \\ 0 & y_i f(\mathbf{x}_i) > 0 \end{cases} \quad (55)$$

进一步改写为:

$$l(y_i, f(\mathbf{x}_i)) = \begin{cases} -y_i f(\mathbf{x}_i) (\equiv 1) & y_i f(\mathbf{x}_i) < 0 \\ 0 & y_i f(\mathbf{x}_i) > 0 \end{cases} \quad (56)$$

然而, 0-1 损失函数是一个非凸函数, 且不是参数  $\mathbf{w}$  和  $b$  的可导函数, 不易进行优化, 如图3-8中红线所示。其中的一个关键因素在于, 使用  $\text{sign}$  对  $\mathbf{w}^T \mathbf{x}_i + b$  进行了符号钳制 (即二值化为  $+1$  或  $-1$ ), 使得  $-y_i \text{sign}(\mathbf{w}^T \mathbf{x}_i + b)$  在  $(0, 0)$  处, 出现了阶跃。

因此, 可以考虑将  $\text{sign}$  从损失函数中去掉<sup>17</sup>, 于是便得到了感知机的 (单样本

<sup>17</sup>当然, 仍然需要保留正确分类时函数值为 0 的情形。

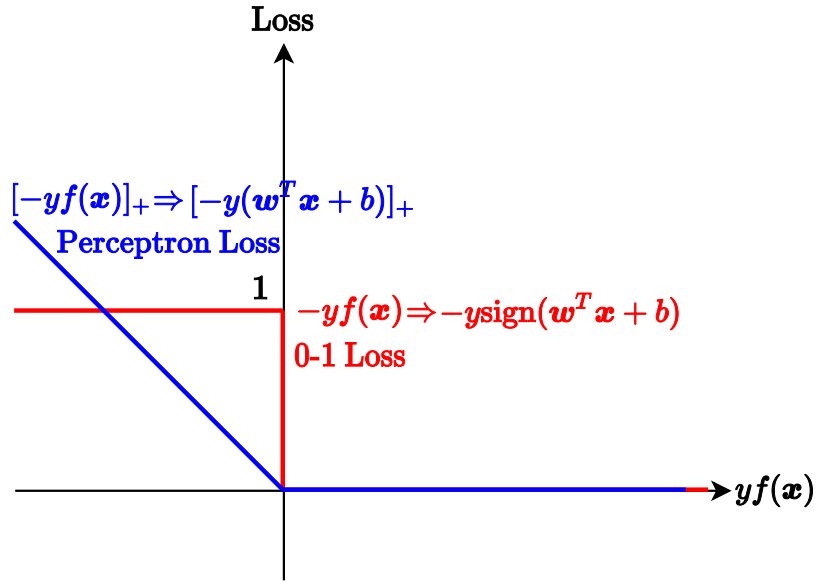


图 3-8: 0-1 损失函数与感知机损失函数

点) 损失函数  $-y_i(\mathbf{w}^T \mathbf{x}_i + b)$ , 如图3-8中蓝线所示<sup>18</sup>:

$$l(y_i, f(\mathbf{x}_i)) = \begin{cases} -y_i f(\mathbf{x}_i) (\equiv -y_i(\mathbf{w}^T \mathbf{x}_i + b)) & y_i f(\mathbf{x}_i) < 0 \\ 0 & y_i f(\mathbf{x}_i) \geq 0 \end{cases} \quad (58)$$

可以看出, 感知机损失函数在误分类 (即  $y_i(\mathbf{w}^T \mathbf{x}_i + b) < 0$ ) 时, 能够提供梯度信息。而 0-1 损失函数在误分类时, 无法提供梯度信息, 不能使学习算法进行有效的优化。

对于上述的感知机损失函数, 也可以从几何的角度来理解:

- 对于分类正确的样本点, 其损失或代价为 0;
- 对于分类错误的样本点, 令  $\|\mathbf{w}\| = 1$  (单位向量), 则其损失或代价为该样本点到决策面  $\mathbf{w}^T \mathbf{x} + b = 0$  的 (无符号) 距离:

$$-y_i \frac{(\mathbf{w}^T \mathbf{x}_i + b)}{\|\mathbf{w}\|} \equiv -y_i(\mathbf{w}^T \mathbf{x}_i + b) \quad (59)$$

<sup>18</sup>实际上, 感知机损失函数的核心函数可以表示为:

$$[z]_+ = \begin{cases} z & z > 0 \\ 0 & z \leq 0 \end{cases} \quad (57)$$

其中, 下标 + 表示当  $z > 0$  时, 函数值取  $z$ ; 否则, 函数值为 0。

其中,  $\frac{(\mathbf{w}^T \mathbf{x}_i + b)}{\|\mathbf{w}\|}$  表示样本点  $\mathbf{x}_i$  到决策面  $\mathbf{w}^T \mathbf{x} + b = 0$  的有符号距离; 对于误分类样本点  $\mathbf{x}_i$ ,  $-y_i(\mathbf{w}^T \mathbf{x}_i + b) > 0$  成立, 它表示 (无符号) 距离 (因  $\|\mathbf{w}\| = 1$ )。

于是, 对于给定的数据集  $T$ , 假设误分类点的集合为  $M$ , 那么所有误分类点到决策面  $\mathbf{w}^T \mathbf{x} + b$  的总距离为:

$$L(\mathbf{w}, b) = - \sum_{\mathbf{x}_i \in M} y_i(\mathbf{w}^T \mathbf{x}_i + b) \quad (60)$$

其中  $\|\mathbf{w}\| = 1$ 。上式就是感知机模型的损失函数。

感知机的学习策略是从假设空间  $\{f|f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b\}$  中选取使损失函数 (公式 (60)) 最小的模型参数  $\mathbf{w}$  和  $b$ 。于是, 将感知机模型的学习问题变成了一个最优化问题。

### 3.5.2 学习的原始算法

在损失函数表示为公式 (60) 时, 感知机模型的学习问题可以表示为如下的最优化问题:

$$\min_{\mathbf{w}, b} L(\mathbf{w}, b) = \min_{\mathbf{w}, b} - \sum_{\mathbf{x}_i \in M} y_i(\mathbf{w}^T \mathbf{x}_i + b) \quad (61)$$

其中,  $M$  为误分类样本点的集合。

对于误分类样本点的集合  $M$ , 损失函数  $L(\mathbf{w}, b)$  关于权值参数  $\mathbf{w}$  和偏置参数  $b$  的梯度可分别表示为:

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{w}} &= - \sum_{\mathbf{x}_i \in M} y_i \mathbf{x}_i \\ \frac{\partial L}{\partial b} &= - \sum_{\mathbf{x}_i \in M} y_i \end{aligned} \quad (62)$$

对于这些误分类点, 可以应用梯度下降方法, 对初始权值向量  $\mathbf{w}_0$  和偏置  $b_0$  进行迭代更新, 直至其收敛为止, 从而得到最优权值向量  $\mathbf{w}^*$  和偏置  $b^*$ 。稍后将证明, 在数据集线性可分的情况下, 感知机学习算法将确保收敛。易知, 感知机模型的学习是由误分类样本点驱动的。

在实际应用中, 一般不使用误分类点集合  $M$  执行批梯度下降<sup>19</sup>, 而是使用随机梯度下降对权值向量与偏置进行更新。具体而言, 权值向量  $\mathbf{w}$  和偏置  $b$  的更新

<sup>19</sup>在批梯度下降法中, 每次迭代时, 需要执行一次遍历, 以找出误分类点集合  $M$ 。因此, 使用随机梯度法替代。

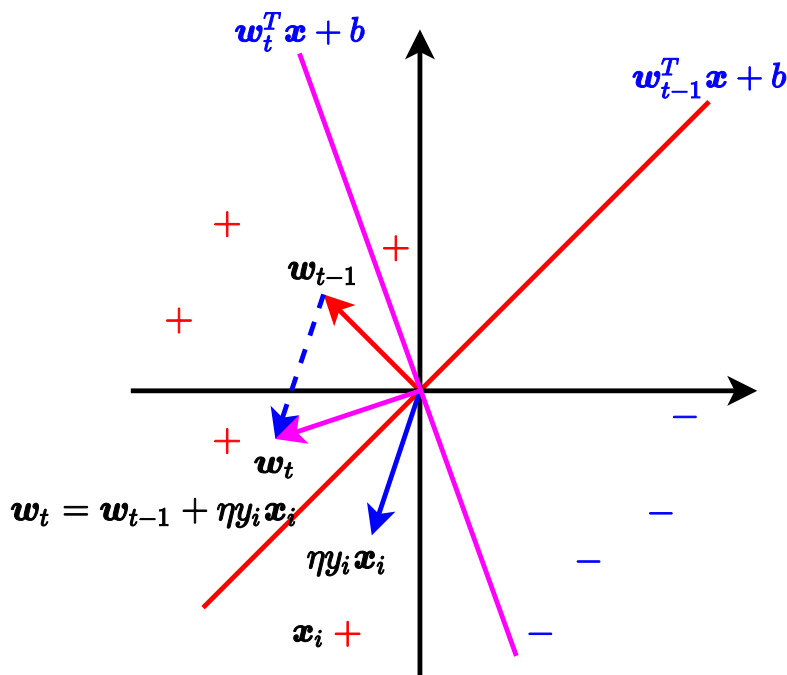


图 3-9: 权值向量更新公式的几何解释

公式如下：

$$\begin{aligned} w_t &= w_{t-1} + \eta y_i x_i \\ b_t &= b_{t-1} + \eta y_i \end{aligned} \quad (63)$$

其中， $\eta$  为学习率， $(x_i, y_i)$  为误分类样本点对。上述权值向量与偏置的更新公式，具有很好的几何直观性与解释性，如图3-9所示。图中，红色直线表示  $t-1$  时刻的决策面，紫色直线表示  $t$  时刻的决策面。当一个实例点  $x_i$  被误分类时，学习算法将调整  $w$  和  $b$  (图中只展示了权值向量  $w$  的更新)，使决策面向误分类点  $x_i$  的一侧旋转和移动，以减少该误分类点与决策面的距离，或使得新决策面能够正确地识别误分类点<sup>20</sup>，从而降低了总体损失。上述模型参数的更新过程，迭代进行，一直进行到没有误分类点出现为止，算法停止。需要注意的是，感知机模型的最优参数  $w^*$  和  $b^*$  不具有唯一性，它们与初始值有很大的关系。

感知机学习算法有 2 种形式：原始形式与对偶形式。下面，首先给出感知机学习算法的原始形式。

### 算法 3.1 (感知机学习算法的原始形式)

<sup>20</sup>当然，在调整后，可能会使得原本正确分类的样本点被误分类。但是，整体损失或代价在不断地降低。

---

```

1 Input:
2     Dataset:  $T$ 
3     Learning Rate:  $\eta$ 
4 Output:
5      $w, b$ 
6 def train( $T, \eta$ ):
7      $w, b = 0, 0$ 
8     finish = False
9     while not finish:
10        finish = True
11        for ( $x_i, y_i$ ) in  $T$ :
12            if  $y_i * (w * x_i + b) \leq 0$ :
13                 $w = w + \eta * y_i * x_i$ 
14                 $b = b + \eta * y_i$ 
15                finish = False
16        if finish:
17            return  $w, b$ 

```

---

### 3.5.3 算法的收敛性

下面将证明，在数据集线性可分的情况下，感知机的学习算法经过有限次迭代后，参数收敛，算法停止，且能够正确地将数据集进行分类。

为了便于表达，将决策面  $w^T x + b$  写成如下形式：

$$\hat{w}^T \hat{x} \equiv w^T x + b \quad (64)$$

即令  $\hat{w} = (w^T, b)^T$ ,  $\hat{x} = (x^T, 1)^T$ 。

**定理 3.1** (Novikoff收敛定理) 设数据集  $T$ :

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

是线性可分的，其中  $x_i \in \mathbf{R}^n$ ,  $y_i = \{+1, -1\}$ ,  $i = 1, 2, \dots, N$ ，则：

1. 一定存在决策面  $\hat{w}_*^T \hat{x} = 0$  (其中  $\|\hat{w}_*\| = 1$ )，它能够完全正确地将数据集分开，且存在  $\gamma > 0$ ，下式成立：

$$y_i (\hat{w}_*^T \hat{x}_i) \geq \gamma \quad \forall i, i = 1, 2, \dots, N \quad (65)$$

2. 令  $R = \max_{1 \leq i \leq N} \|\hat{x}_i\|$ ，则感知机学习算法在训练数据集  $T$  上的误分类次数  $k$  满足不等式：

$$k \leq \left( \frac{R}{\gamma} \right)^2 \quad (66)$$

证明：



1. 由于数据集  $T$  是线性可分的, 根据定义3.1可知, 存在某个决策面 (令其为  $\hat{\mathbf{w}}_*^T \hat{\mathbf{x}} = 0$ , 且  $\|\hat{\mathbf{w}}_*\| = 1$ ), 能够完全正确地将数据集分类, 且对于  $\forall (\hat{\mathbf{x}}_i, y_i) (i = 1, 2, \dots, N)$ ,  $y_i \hat{\mathbf{w}}_*^T \hat{\mathbf{x}}_i > 0$  成立, 于是令  $\gamma = \min_{(\hat{\mathbf{x}}_i, y_i)} y_i \hat{\mathbf{w}}_*^T \hat{\mathbf{x}}_i$ , 则  $y_i \hat{\mathbf{w}}_*^T \hat{\mathbf{x}}_i \geq \gamma (\forall i, i = 1, 2, \dots, N)$  成立。
2. 设误分类样本点为  $(\hat{\mathbf{x}}_i, y_i)$ , 根据感知机学习算法的原始形式 (算法3.1) 可知, 此时的权值需要更新, 其更新公式如下:

$$\begin{aligned} \mathbf{w}_k &= \mathbf{w}_{k-1} + \eta y_i \mathbf{x}_i \\ b_k &= b_{k-1} + \eta y_i \end{aligned} \Rightarrow \hat{\mathbf{w}}_k = \hat{\mathbf{w}}_{k-1} + \eta y_i \hat{\mathbf{x}}_i \quad (67)$$

其中, 误分类点  $(\hat{\mathbf{x}}_i, y_i)$  满足条件  $y_i \hat{\mathbf{w}}_{k-1}^T \hat{\mathbf{x}}_i \leq 0$ ; 且当  $k = 0$  时, 有  $\hat{\mathbf{w}}_0 = \mathbf{0}$ 。  
对公式 (67) 两端乘以  $\hat{\mathbf{w}}_*^T$ , 可得:

$$\hat{\mathbf{w}}_*^T \hat{\mathbf{w}}_k = \hat{\mathbf{w}}_*^T \hat{\mathbf{w}}_{k-1} + \eta y_i \hat{\mathbf{w}}_*^T \hat{\mathbf{x}}_i \quad (68)$$

将  $y_i \hat{\mathbf{w}}_*^T \hat{\mathbf{x}}_i \geq \gamma$  代入上式, 得到不等式:

$$\hat{\mathbf{w}}_*^T \hat{\mathbf{w}}_k \geq \hat{\mathbf{w}}_*^T \hat{\mathbf{w}}_{k-1} + \eta \gamma \quad (69)$$

于是, 对上述不等式中的  $\hat{\mathbf{w}}_*^T \hat{\mathbf{w}}_{k-1}$  继续应用类似的过程, 并一直持续到  $\hat{\mathbf{w}}_0 = \mathbf{0}$  为止, 可得:

$$\hat{\mathbf{w}}_*^T \hat{\mathbf{w}}_k \geq \hat{\mathbf{w}}_*^T \hat{\mathbf{w}}_{k-1} + \eta \gamma \geq \hat{\mathbf{w}}_*^T \hat{\mathbf{w}}_{k-2} + 2\eta \gamma \geq \dots \geq k\eta \gamma \quad (70)$$

由此, 得到不等式:

$$k\eta \gamma \leq \hat{\mathbf{w}}_*^T \hat{\mathbf{w}}_k \quad (71)$$

利用  $\|\hat{\mathbf{w}}_*\| = 1$ , 可得:

$$k\eta \gamma \leq \hat{\mathbf{w}}_*^T \hat{\mathbf{w}}_k = \|\hat{\mathbf{w}}_*\| \|\hat{\mathbf{w}}_k\| \cos \theta \leq \|\hat{\mathbf{w}}_*\| \|\hat{\mathbf{w}}_k\| = \|\hat{\mathbf{w}}_k\| \quad (72)$$

其中,  $\theta$  为  $\hat{\mathbf{w}}_*$  和  $\hat{\mathbf{w}}_k$  的夹角。继续对上式变形, 可得:

$$k\eta \gamma \leq \|\hat{\mathbf{w}}_k\| = \sqrt{\hat{\mathbf{w}}_k^T \hat{\mathbf{w}}_k} \quad (73)$$

将  $\hat{\mathbf{w}}_k = \hat{\mathbf{w}}_{k-1} + \eta y_i \hat{\mathbf{x}}_i$  代入, 可得:

$$\begin{aligned} k\eta \gamma &\leq \sqrt{(\hat{\mathbf{w}}_{k-1} + \eta y_i \hat{\mathbf{x}}_i)^T (\hat{\mathbf{w}}_{k-1} + \eta y_i \hat{\mathbf{x}}_i)} \\ &= \sqrt{\hat{\mathbf{w}}_{k-1}^T \hat{\mathbf{w}}_{k-1} + 2\eta y_i \hat{\mathbf{w}}_{k-1}^T \hat{\mathbf{x}}_i + \eta^2 y_i^2 \hat{\mathbf{x}}_i^T \hat{\mathbf{x}}_i} \end{aligned} \quad (74)$$

利用  $y_i^2 = 1$ , 继续变形, 可得:

$$\begin{aligned}
 k\eta\gamma &\leq \sqrt{\hat{\mathbf{w}}_{k-1}^T \hat{\mathbf{w}}_{k-1} + 2\eta y_i \hat{\mathbf{w}}_{k-1}^T \hat{\mathbf{x}}_i + \eta^2 \hat{\mathbf{x}}_i^T \hat{\mathbf{x}}_i} \\
 &= \sqrt{\hat{\mathbf{w}}_{k-1}^T \hat{\mathbf{w}}_{k-1} + 2\eta y_i \hat{\mathbf{w}}_{k-1}^T \hat{\mathbf{x}}_i + \eta^2 \|\hat{\mathbf{x}}_i\|^2} \\
 &\leq \sqrt{\hat{\mathbf{w}}_{k-1}^T \hat{\mathbf{w}}_{k-1} + \eta^2 \|\hat{\mathbf{x}}_i\|^2} \\
 &\leq \sqrt{\hat{\mathbf{w}}_{k-1}^T \hat{\mathbf{w}}_{k-1} + \eta^2 R^2}
 \end{aligned} \tag{75}$$

其中,  $R = \max_{1 \leq i \leq N} \|\hat{\mathbf{x}}_i\|$ 。于是, 对上述不等式中的  $\hat{\mathbf{w}}_{k-1}^T \hat{\mathbf{w}}_{k-1}$  继续应用类似的过程, 并一直持续到  $\hat{\mathbf{w}}_0 = \mathbf{0}$  为止, 可得:

$$k\eta\gamma \leq \sqrt{\hat{\mathbf{w}}_{k-1}^T \hat{\mathbf{w}}_{k-1} + \eta^2 R^2} \leq \sqrt{\hat{\mathbf{w}}_{k-2}^T \hat{\mathbf{w}}_{k-2} + 2\eta^2 R^2} \leq \dots \leq \sqrt{k\eta^2 R^2} = \sqrt{k}\eta R \tag{76}$$

最后, 得到:

$$k \leq \left(\frac{R}{\gamma}\right)^2 \tag{77}$$

可以看出, 对于数据集线性可分情形, 误分类的次数  $k$  (即权值向量  $\hat{\mathbf{w}}$  的调整次数) 是有上界的, 感知机学习算法的原始形式 (算法3.1) 必定收敛。

□

### 3.5.4 学习的对偶算法

感知机的学习算法, 除了原始形式之外, 还存在另一种形式——对偶算法。在原始算法中, 每次迭代直接更新权值  $\mathbf{w}$  和偏置  $b$ 。而在对偶算法中, 首先为每个样本点配置一个标量  $\alpha_i$ , 用来记录每个样本点被误分的“次数”<sup>21</sup>; 然后, 无论是  $\mathbf{w}$ , 还是  $b$ , 都可以使用  $\alpha$  进行 (间接) 更新。其中, 向量  $\alpha$  的元素由  $\alpha_i$  组成 ( $i = 1, 2, \dots, N$ )。

对偶算法的优点在于:

- 将  $\mathbf{w}$  和  $b$  表示成样本点对  $(\mathbf{x}_i, y_i)$  的线性组合形式, 为核函数的引入提供了基础;
- 在每次迭代的过程中, 数据样本点之间仅以向量内积  $\mathbf{x}_i^T \mathbf{x}_j$  的形式出现。因此, 可以预先计算它们之间的内积, 从而形成了 Gram 矩阵:

$$\mathbf{G} = [\mathbf{x}_i^T \mathbf{x}_j]_{N \times N} \tag{78}$$

<sup>21</sup>注意, “次数”与学习率  $\eta$  的取值有关: 当  $\eta = 1$  时,  $\alpha_i$  就是第  $i$  个样本点的误分次数; 否则, 如果  $0 < \eta < 1$ , 则  $\frac{\alpha_i}{\eta}$  才是误分次数。综合两种情况,  $\frac{\alpha_i}{\eta}$  就表示误分次数。

下面将讨论，如何从原始算法得到等价的对偶算法。在原始算法3.1中，权值向量与偏置的更新公式分别如下：

$$\begin{aligned} \mathbf{w} &\leftarrow \mathbf{w} + \eta y_i \mathbf{x}_i \\ b &\leftarrow b + \eta y_i \end{aligned} \quad (79)$$

其中， $(\mathbf{x}_i, y_i)$  为误分类点对。仔细观察，可以发现，每次  $(\mathbf{x}_i, y_i)$  被误分类时，权值  $\mathbf{w}$  和偏置  $b$  就分别执行一次更新  $\eta y_i \mathbf{x}_i$  和  $\eta y_i$ 。除去样本点对  $(\mathbf{x}_i, y_i)$  之外，更新量只与  $\eta$  有关。因此，可令  $n_i$  表示样本点对  $(\mathbf{x}_i, y_i)$  的更新次数<sup>22</sup>，并令  $\alpha_i = n_i \eta$ ，于是得到公式 (79) 的等价形式：

$$\begin{aligned} \mathbf{w} &= \sum_{i=1}^N n_i \eta y_i \mathbf{x}_i & \Rightarrow & \quad \mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i \\ b &= \sum_{i=1}^N n_i \eta y_i & & \quad b = \sum_{i=1}^N \alpha_i y_i \end{aligned} \quad (80)$$

于是，在迭代过程中，对偶算法只需要对  $\alpha_i$  的值进行更新：

$$\alpha_i \leftarrow \alpha_i + \eta \quad (81)$$

前提条件是，当样本点对  $(\mathbf{x}_i, y_i)$  是误分类点，即满足如下条件时：

$$\begin{aligned} y_i (\mathbf{w}^T \mathbf{x}_i + b) &\leq 0 \quad \Rightarrow \\ y_i (\mathbf{x}_i^T \mathbf{w} + b) &\leq 0 \quad \Rightarrow \\ y_i \left( \mathbf{x}_i^T \sum_{j=1}^N \alpha_j y_j \mathbf{x}_j + \sum_{j=1}^N \alpha_j y_j \right) &\leq 0 \quad \Rightarrow \\ y_i \left( \sum_{j=1}^N \alpha_j y_j \mathbf{x}_j^T \mathbf{x}_i + \sum_{j=1}^N \alpha_j y_j \right) &\leq 0 \end{aligned} \quad (82)$$

$\alpha_i$  的意义也是非常明确的：其值越大，表明样本点对  $(\mathbf{x}_i, y_i)$  被误分类的次数越多，它离决策面也越近，对决策面的形成有较大的影响<sup>23</sup>；相反，那些  $\alpha_i = 0$  的点，离决策面较远，对决策面的形成没有任何影响。

下面给出感知机学习的对偶算法。

### 算法 3.2 (感知机学习算法的对偶形式)

<sup>22</sup>注意，在算法3.1中， $\mathbf{w}_0 = \mathbf{0}$  和  $b_0 = 0$ 。

<sup>23</sup>在支持向量机中，这样的样本点被称为支持向量。

---

```

1 Input:
2     Dataset:  $T=(x, y)$ , Vector:  $x$ , Vector:  $y$ 
3     Learning Rate:  $\eta$ 
4 Output:
5     Vector:  $\alpha$ 
6 def train(T, eta):
7      $\alpha = 0$ 
8      $finish = False$ 
9     while not finish:
10         $finish = True$ 
11        for (xi, yi) in T:
12             $wxi = \sum(\alpha[j] * y[j] * \text{dot}(x[j], xi))$  for all j
13             $b = \sum(\alpha[j] * y[j])$  for all j
14            if yi * (wxi + b) <= 0:
15                 $\alpha[i] = \alpha[i] + \eta$ 
16             $finish = False$ 
17        if finish:
18            return alpha

```

---

## 4 附录

### 4.1 标量函数对矩阵的求导

标量函数  $f(\mathbf{X})$  对矩阵  $\mathbf{X}$  的导数，定义为：

$$\frac{\partial f}{\partial \mathbf{X}} = \left[ \frac{\partial f}{\partial X_{ij}} \right] \quad (83)$$

即  $f(\mathbf{X})$  对  $\mathbf{X}$  逐元素求导，且排列成与矩阵  $\mathbf{X}$  各维度相同的矩阵。在实际应用中，可以利用一些基本的运算规则，进行有效的推导，下面将逐一进行介绍。

为此，首先需要定义一些基本的操作与运算规则。在推导的过程中，需要使用微分与导数或梯度的关系。在一元微积分中，标量函数  $f(x)$  对标量变量  $x$  的导数为  $f'(x)$ ，其与微分  $df$  的关系如下：

$$df = f'(x)dx \quad (84)$$

在多元微积分中，梯度<sup>24</sup>与微分的关系如下：

$$df = \sum_i \frac{\partial f}{\partial x_i} dx_i = \frac{\partial f}{\partial \mathbf{x}}^T d\mathbf{x} \quad (85)$$

其中， $df$  表示全微分， $\frac{\partial f}{\partial \mathbf{x}}$  为偏导数向量， $d\mathbf{x}$  为微分向量， $df = \frac{\partial f}{\partial \mathbf{x}}^T d\mathbf{x}$  表示全微分等于两个向量的内积。

---

<sup>24</sup>即标量函数对向量的导数或标量函数对向量各分量的偏导数。

类似地，也可以在全微分  $df$  与标量函数  $f(x)$  对矩阵  $\mathbf{X}$  的导数之间建立联系：

$$df = \sum_{i=1}^m \sum_{j=1}^n \frac{\partial f}{\partial X_{ij}} dX_{ij} = \text{tr} \left( \frac{\partial f}{\partial \mathbf{X}}^T d\mathbf{X} \right) \quad (86)$$

其中， $\text{tr}$  表示迹运算<sup>25</sup>。全微分  $df$  等于  $\frac{\partial f}{\partial \mathbf{X}}$  和  $d\mathbf{X}$  这 2 个矩阵的标准内积，上式表明，它可以方便地使用迹来表示。

下面是关于全微分的一些基本运算：

- $d(\mathbf{X} \pm \mathbf{Y}) = d\mathbf{X} \pm d\mathbf{Y}$ ;
- $d(\mathbf{X}\mathbf{Y}) = d(\mathbf{X})\mathbf{Y} + \mathbf{X}d\mathbf{Y}$ ;
- $d(\mathbf{X}^T) = (d\mathbf{X})^T$ ;
- $d\text{tr}(\mathbf{X}) = \text{tr}(d\mathbf{X})$ ;
- $d\mathbf{X}^{-1} = -\mathbf{X}^{-1}d\mathbf{X}\mathbf{X}^{-1}$ ; <sup>26</sup>
- $d|\mathbf{X}| = \text{tr}(\mathbf{X}^*d\mathbf{X})$ ，其中  $|\mathbf{X}|$  表示  $n$  阶矩阵  $\mathbf{X}$  的行列式； $\mathbf{X}^*$  表示  $\mathbf{X}$  的伴随矩阵，当  $\mathbf{X}$  可逆时， $\mathbf{X}^* = |\mathbf{X}|\mathbf{X}^{-1}$ ，则  $d|\mathbf{X}| = |\mathbf{X}|\text{tr}(\mathbf{X}^{-1}d\mathbf{X})$ ; <sup>27</sup>
- $d(\mathbf{X} \odot \mathbf{Y}) = d\mathbf{X} \odot \mathbf{Y} + \mathbf{X} \odot d\mathbf{Y}$ ， $\odot$  表示 Hadamard 积 (Hadamard Product)，即两个矩阵对应元素的积，这两个矩阵的维度必须相同，结果矩阵的维度与原矩阵保持一致；
- $d\sigma(\mathbf{X}) = \sigma'(\mathbf{X}) \odot d\mathbf{X}$ ，其中  $\sigma(\mathbf{X}) = [\sigma(X_{ij})]$  是对矩阵  $\mathbf{X}$  的每个元素进行逐元素运算的标量函数， $\sigma'(\mathbf{X}) = [\sigma'(X_{ij})]$  是与  $\sigma(\mathbf{X})$  对应的逐元素求导的

<sup>25</sup>迹运算表示方阵主对角线上的元素之和，这是一种应用广泛的运算。其应用之一是，定义矩阵  $\mathbf{X}, \mathbf{Y} \in \mathbf{R}^{m \times n}$  的标准内积为：

$$\langle \mathbf{X}, \mathbf{Y} \rangle = \text{tr}(\mathbf{X}^T \mathbf{Y}) = \sum_{i=1}^m \sum_{j=1}^n X_{ij} Y_{ij} \quad (87)$$

注意，如果  $\mathbf{A}$  与  $\mathbf{B}$  符合矩阵乘法规则，则  $\text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{BA})$  成立。如果  $\mathbf{A}$ 、 $\mathbf{B}$  与  $\mathbf{C}$  均为方阵，则  $\text{tr}(\mathbf{ABC}) = \text{tr}(\mathbf{BCA}) = \text{tr}(\mathbf{CAB})$  成立。请参考《机器学习》课程系列之《基础知识》，Chapter1-CN.pdf。

<sup>26</sup> $\mathbf{X}\mathbf{X}^{-1} = \mathbf{I} \Rightarrow d(\mathbf{X}\mathbf{X}^{-1}) = d\mathbf{I} \Rightarrow (d\mathbf{X})\mathbf{X}^{-1} + \mathbf{X}d\mathbf{X}^{-1} = \mathbf{0} \Rightarrow d\mathbf{X}^{-1} = -\mathbf{X}^{-1}d\mathbf{X}\mathbf{X}^{-1}$ 。

<sup>27</sup>可以使用 Laplace 展开证明，详见张贤达《矩阵分析与应用》第 279 页。

标量函数。例如：

$$\mathbf{X} = \begin{bmatrix} X_{11} & X_{12} \\ X_{21} & X_{22} \end{bmatrix}, d \sin(\mathbf{X}) = \begin{bmatrix} \cos X_{11} dX_{11} & \cos X_{12} dX_{12} \\ \cos X_{21} dX_{21} & \cos X_{22} dX_{22} \end{bmatrix} = \cos(\mathbf{X}) \odot d\mathbf{X}$$

现在，可以利用全微分的迹运算公式 (86) 来获取标量函数  $f(\mathbf{X})$  对矩阵  $\mathbf{X}$  的导数：通过迹运算，在导数  $\frac{\partial f}{\partial \mathbf{X}}$ 、微分  $d\mathbf{X}$  及全微分  $df$  之间建立联系，将全微分变换成该公式的右侧形式，这样就直接获得了导数  $\frac{\partial f}{\partial \mathbf{X}}$  的整体形式，而不需要依据其定义 (公式 (83)) 进行低效率的逐元素求解。

为了应用上述策略，需要使用下面的迹技巧 (Trace Trick)：

- $a = \text{tr}(a)$ ;
- $\text{tr}(\mathbf{X}^T) = \text{tr}(\mathbf{X})$ ;
- $\text{tr}(\mathbf{X} \pm \mathbf{Y}) = \text{tr}(\mathbf{X}) \pm \text{tr}(\mathbf{Y})$ ;
- $\text{tr}(\mathbf{X}\mathbf{Y}) = \text{tr}(\mathbf{Y}\mathbf{X})$ ，其中  $\mathbf{X}$  与  $\mathbf{Y}$  符合矩阵乘法规则，运算的结果是  $\sum_{i,j} X_{ij}Y_{ji}$ ;
- $\text{tr}(\mathbf{X}^T(\mathbf{Y} \odot \mathbf{Z})) = \text{tr}((\mathbf{X} \odot \mathbf{Y})^T \mathbf{Z})$ ，其中， $\mathbf{X}$ 、 $\mathbf{Y}$ 、 $\mathbf{Z}$  的维度都相同，运算的结果是  $\sum_{i,j} X_{ij}Y_{ij}Z_{ij}$ ;

综上所述，可以直接利用公式 (85) 和公式 (86)，通过适当的变换，就可以直接求解出标量函数关于向量的偏导数以及标量函数关于矩阵的偏导数。这些适当的变换包括上述关于微分与迹的加减乘法、逆、行列式、逐元素运算函数等，也包括将迹运算应用于全微分  $df$  并利用迹中矩阵乘积的可交换性，将  $d\mathbf{X}$  右侧项交换至左侧等技巧，最后比照全微分公式，就可以直接获得所需的偏导数。

需要注意的是，在求解偏导数的过程中，经常遇到的一种求导问题是所谓的复合求导。例如，设  $f(\mathbf{Y})$  是矩阵  $\mathbf{Y}$  的函数，而矩阵  $\mathbf{Y}$  又是矩阵  $\mathbf{X}$  的函数。假设已经求解出  $\frac{\partial f}{\partial \mathbf{Y}}$ ，那么如何求解出  $\frac{\partial f}{\partial \mathbf{X}}$ ？

求解  $\frac{\partial f}{\partial \mathbf{X}}$ ，不能简单地套用标量函数关于标量变量求导的链式法则。但是，通过函数嵌套的方式来使用全微分公式 ((85) 和 (86))，经过适当的变换与变形，最终也能够得到偏导数。例如，在公式  $df = \text{tr}\left(\frac{\partial f}{\partial \mathbf{Y}}^T d\mathbf{Y}\right)$  中，对其中的  $d\mathbf{Y}$  再次应用全微分公式  $d\mathbf{Y} = \text{tr}\left(\frac{\partial \mathbf{Y}}{\partial \mathbf{X}}^T d\mathbf{X}\right)$ ，并将后者代入前者，最后使用迹技巧将  $d\mathbf{X}$  右侧所有项交换到左侧，即可得到  $\frac{\partial f}{\partial \mathbf{X}}$ 。

例如, 已知  $\frac{\partial f}{\partial \mathbf{Y}}$ ,  $\mathbf{Y} = \mathbf{A}\mathbf{X}\mathbf{B}$ , 求解  $\frac{\partial f}{\partial \mathbf{X}}$  的过程如下:

$$\begin{aligned} df &= \text{tr} \left( \frac{\partial f}{\partial \mathbf{Y}}^T d\mathbf{Y} \right) = \text{tr} \left( \frac{\partial f}{\partial \mathbf{Y}}^T \mathbf{A} d\mathbf{X} \mathbf{B} \right) = \text{tr} \left( \mathbf{B} \frac{\partial f}{\partial \mathbf{Y}}^T \mathbf{A} d\mathbf{X} \right) \\ &= \text{tr} \left( \left( \mathbf{A}^T \frac{\partial f}{\partial \mathbf{Y}} \mathbf{B}^T \right)^T d\mathbf{X} \right) \end{aligned} \quad (88)$$

在上面的推导过程中, 使用了如下推导:

$$d\mathbf{Y} = (d\mathbf{A})\mathbf{X}\mathbf{B} + \mathbf{A}d\mathbf{X}\mathbf{B} + \mathbf{A}\mathbf{X}d\mathbf{B} = \mathbf{A}d\mathbf{X}\mathbf{B} \quad (89)$$

其中,  $\mathbf{A}$  和  $\mathbf{B}$  均为常量。最后, 得到  $\frac{\partial f}{\partial \mathbf{X}} = \mathbf{A}^T \frac{\partial f}{\partial \mathbf{Y}} \mathbf{B}^T$ 。

下面, 将列举若干实例, 应用上述原则求解所需偏导数:

- 已知  $f(\mathbf{X}) = \mathbf{a}^T \mathbf{X} \mathbf{b}$ , 其中  $\mathbf{a}$  和  $\mathbf{b}$  分别为  $m$  维和  $n$  维列向量,  $\mathbf{X}$  为  $m \times n$  维矩阵,  $f$  为标量函数, 求解  $\frac{\partial f}{\partial \mathbf{X}}$ :

$$\begin{aligned} f(\mathbf{X}) &= \mathbf{a}^T \mathbf{X} \mathbf{b} \Rightarrow df = d(\mathbf{a}^T \mathbf{X} \mathbf{b}) \Rightarrow \\ df &= d\mathbf{a}^T \mathbf{X} \mathbf{b} + \mathbf{a}^T d\mathbf{X} \mathbf{b} + \mathbf{a}^T \mathbf{X} d\mathbf{b} = \mathbf{a}^T d\mathbf{X} \mathbf{b} \Rightarrow \\ df &= \text{tr}(df) = \text{tr}(\mathbf{a}^T d\mathbf{X} \mathbf{b}) = \text{tr}(\mathbf{b} \mathbf{a}^T d\mathbf{X}) = \text{tr} \left( (\mathbf{a} \mathbf{b}^T)^T d\mathbf{X} \right) \end{aligned} \quad (90)$$

可得  $\frac{\partial f}{\partial \mathbf{X}} = \mathbf{a} \mathbf{b}^T$ 。

- 已知  $f(\mathbf{X}) = \mathbf{a}^T e^{\mathbf{X} \mathbf{b}}$ , 其中  $\mathbf{a}$  和  $\mathbf{b}$  分别为  $m$  维和  $n$  维列向量,  $\mathbf{X}$  为  $m \times n$  维矩阵,  $f$  为标量函数, 此处的指数函数  $e$  执行逐元素运算, 求解  $\frac{\partial f}{\partial \mathbf{X}}$ :

$$\begin{aligned} f(\mathbf{X}) &= \mathbf{a}^T e^{\mathbf{X} \mathbf{b}} \Rightarrow df = d(\mathbf{a}^T e^{\mathbf{X} \mathbf{b}}) \Rightarrow \\ df &= d\mathbf{a}^T e^{\mathbf{X} \mathbf{b}} + \mathbf{a}^T d(e^{\mathbf{X} \mathbf{b}}) = \mathbf{a}^T (e^{\mathbf{X} \mathbf{b}} \odot d(\mathbf{X} \mathbf{b})) \Rightarrow \\ df &= \text{tr}(df) = \text{tr}(\mathbf{a}^T (e^{\mathbf{X} \mathbf{b}} \odot d(\mathbf{X} \mathbf{b}))) = \text{tr} \left( (\mathbf{a} \odot e^{\mathbf{X} \mathbf{b}})^T d(\mathbf{X} \mathbf{b}) \right) \\ &= \text{tr} \left( \mathbf{b} (\mathbf{a} \odot e^{\mathbf{X} \mathbf{b}})^T d\mathbf{X} \right) = \text{tr} \left( ((\mathbf{a} \odot e^{\mathbf{X} \mathbf{b}}) \mathbf{b}^T)^T d\mathbf{X} \right) \end{aligned} \quad (91)$$

可得  $\frac{\partial f}{\partial \mathbf{X}} = (\mathbf{a} \odot e^{\mathbf{X} \mathbf{b}}) \mathbf{b}^T$

- 已知  $f(\mathbf{Y}) = \text{tr}(\mathbf{Y}^T \mathbf{A} \mathbf{Y})$ ,  $\mathbf{Y} = \sigma(\mathbf{B} \mathbf{X})$ , 其中  $\mathbf{Y}$  为  $l \times n$  矩阵,  $\mathbf{A}$  为  $l \times l$  对称矩阵,  $\mathbf{B}$  为  $l \times m$  矩阵,  $\mathbf{X}$  为  $m \times n$  矩阵,  $\sigma$  为逐元素运算函数, 求

解  $\frac{\partial f}{\partial \mathbf{X}}$ :

$$\begin{aligned}
 f(\mathbf{Y}) &= \text{tr}(\mathbf{Y}^T \mathbf{A} \mathbf{Y}) \Rightarrow df = d(\text{tr}(\mathbf{Y}^T \mathbf{A} \mathbf{Y})) \Rightarrow \\
 df &= \text{tr}(d(\mathbf{Y}^T \mathbf{A} \mathbf{Y})) = \text{tr}(d(\mathbf{Y}^T) \mathbf{A} \mathbf{Y} + \mathbf{Y}^T d(\mathbf{A} \mathbf{Y})) \\
 &= \text{tr}((d\mathbf{Y})^T \mathbf{A} \mathbf{Y} + \mathbf{Y}^T \mathbf{A} d\mathbf{Y}) = \text{tr}((d\mathbf{Y})^T \mathbf{A} \mathbf{Y}) + \text{tr}(\mathbf{Y}^T \mathbf{A} d\mathbf{Y}) \\
 &= \text{tr}((\mathbf{A} \mathbf{Y})^T d\mathbf{Y}) + \text{tr}(\mathbf{Y}^T \mathbf{A} d\mathbf{Y}) = \text{tr}(\mathbf{Y}^T \mathbf{A}^T d\mathbf{Y}) + \text{tr}(\mathbf{Y}^T \mathbf{A} d\mathbf{Y}) \\
 &= 2\text{tr}(\mathbf{Y}^T \mathbf{A} d\mathbf{Y}) = 2\text{tr}((\sigma(\mathbf{B}\mathbf{X}))^T \mathbf{A} d(\sigma(\mathbf{B}\mathbf{X}))) \\
 &= 2\text{tr}((\sigma(\mathbf{B}\mathbf{X}))^T \mathbf{A} (\sigma'(\mathbf{B}\mathbf{X}) \odot \mathbf{B} d\mathbf{X})) \\
 &= 2\text{tr}((\mathbf{A}^T \sigma(\mathbf{B}\mathbf{X}) \odot \sigma'(\mathbf{B}\mathbf{X}))^T \mathbf{B} d\mathbf{X}) \\
 &= 2\text{tr}((\mathbf{B}^T \mathbf{A}^T \sigma(\mathbf{B}\mathbf{X}) \odot \sigma'(\mathbf{B}\mathbf{X}))^T d\mathbf{X})
 \end{aligned} \tag{92}$$

可得  $\frac{\partial f}{\partial \mathbf{X}} = 2\mathbf{B}^T \mathbf{A}^T \sigma(\mathbf{B}\mathbf{X}) \odot \sigma'(\mathbf{B}\mathbf{X}) = 2\mathbf{B}^T \mathbf{A} \sigma(\mathbf{B}\mathbf{X}) \odot \sigma'(\mathbf{B}\mathbf{X})$ 。

- 线性回归模型的损失函数为  $l = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$ ，其中  $\mathbf{y}$  为  $m$  维列向量， $\mathbf{X}$  为  $m \times n$  矩阵， $\mathbf{w}$  为  $n$  维列向量，模型参数  $\mathbf{w}$  的解析解：

$$\begin{aligned}
 l &= \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 \Rightarrow dl = d(\|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2) \Rightarrow \\
 dl &= d((\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})) \Rightarrow \\
 dl &= (d\mathbf{X}\mathbf{w})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) + (\mathbf{X}\mathbf{w} - \mathbf{y})^T d\mathbf{X}\mathbf{w} \Rightarrow \\
 dl &= 2(\mathbf{X}\mathbf{w} - \mathbf{y})^T d\mathbf{X}\mathbf{w} \Rightarrow dl = 2(\mathbf{X}\mathbf{w} - \mathbf{y})^T \mathbf{X} d\mathbf{w}
 \end{aligned} \tag{93}$$

其中，倒数第 2 步的推导利用了向量内积等式： $\mathbf{u}^T \mathbf{v} = \mathbf{v}^T \mathbf{u}$ 。

根据公式 (85)，可得  $\frac{\partial f}{\partial \mathbf{x}} = 2\mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{y})$ ，令  $\frac{\partial f}{\partial \mathbf{x}} = 0$ ，则：

$$2\mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{y}) = 0 \Rightarrow \mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y} \Rightarrow \mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \tag{94}$$

其中，令  $\mathbf{X}^\dagger = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ ，它被称为矩阵  $\mathbf{X}$  的 Moore-Penrose 伪逆矩阵 (Pseudo-Inverse Matrix)。实际上，它可以看作逆矩阵概念的泛化——将逆矩阵推广到非方阵上。如果  $\mathbf{X}$  是方阵且可逆，那么 Moore-Penrose 伪逆矩阵就退化为逆矩阵  $\mathbf{X}^\dagger \equiv \mathbf{X}^{-1}$ 。

- 多元高斯分布的定义为：

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\boldsymbol{\Sigma}|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right\} \tag{95}$$



其中  $D$  维向量  $\boldsymbol{\mu}$  被称为均值向量,  $D \times D$  矩阵  $\boldsymbol{\Sigma}$  被称为协方差矩阵,  $|\boldsymbol{\Sigma}|$  表示  $\boldsymbol{\Sigma}$  的行列式。设样本  $T = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$ , 其中  $\mathbf{x}_i \sim \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}, \boldsymbol{\Sigma})$ , 试利用最大似然方法估计均值  $\boldsymbol{\mu}$  和协方差矩阵  $\boldsymbol{\Sigma}$ 。

样本数据集的对数似然函数为:

$$\begin{aligned} \mathcal{L}(\boldsymbol{\mu}, \boldsymbol{\Sigma}; \mathbf{X}) &= \log \prod_{i=1}^N \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\boldsymbol{\Sigma}|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right\} \\ &= -\frac{DN}{2} \log 2\pi - \frac{N}{2} \log |\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \end{aligned} \quad (96)$$

损失函数表示为负对数似然函数:

$$\begin{aligned} L(\boldsymbol{\mu}, \boldsymbol{\Sigma}; \mathbf{X}) &= -\mathcal{L}(\boldsymbol{\mu}, \boldsymbol{\Sigma}; \mathbf{X}) \\ &= \frac{DN}{2} \log 2\pi + \frac{N}{2} \log |\boldsymbol{\Sigma}| + \frac{1}{2} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \end{aligned} \quad (97)$$

观察上式, 可知,  $\boldsymbol{\mu}$  与  $\boldsymbol{\Sigma}$  相互独立, 可以单独优化。首先, 优化  $\boldsymbol{\mu}$ , 固定  $\boldsymbol{\Sigma}$ :

$$\begin{aligned} dL &= d \left( \frac{DN}{2} \log 2\pi + \frac{N}{2} \log |\boldsymbol{\Sigma}| + \frac{1}{2} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right) \Rightarrow \\ dL &= \frac{1}{2} \sum_{i=1}^N [d(\mathbf{x}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) + (\mathbf{x}_i - \boldsymbol{\mu})^T d(\boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}))] \Rightarrow \\ dL &= \frac{1}{2} \sum_{i=1}^N [-d\boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) - (\mathbf{x}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} d\boldsymbol{\mu}] \Rightarrow \\ dL &= - \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} d\boldsymbol{\mu} \end{aligned} \quad (98)$$

其中, 最后一步的推导, 利用了向量内积等式  $\mathbf{u}^T \mathbf{v} = \mathbf{v}^T \mathbf{u}$  以及  $\boldsymbol{\Sigma}$  为对称矩阵的事实。

根据公式 (85), 得到:  $\frac{\partial f}{\partial \boldsymbol{\mu}} = - \sum_{i=1}^N (\boldsymbol{\Sigma}^{-1})^T (\mathbf{x}_i - \boldsymbol{\mu})$ 。令  $\frac{\partial f}{\partial \boldsymbol{\mu}} = \mathbf{0}$ , 可得:

$$\begin{aligned} - \sum_{i=1}^N (\boldsymbol{\Sigma}^{-1})^T (\mathbf{x}_i - \boldsymbol{\mu}) &= \mathbf{0} \Rightarrow \sum_{i=1}^N (\boldsymbol{\Sigma}^{-1})^T \mathbf{x}_i = N (\boldsymbol{\Sigma}^{-1})^T \boldsymbol{\mu} \Rightarrow \\ \boldsymbol{\mu} &= \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \end{aligned} \quad (99)$$

然后，优化  $\Sigma$ ，固定  $\mu$ ：

$$\begin{aligned}
 dL &= d \left( \frac{DN}{2} \log 2\pi + \frac{N}{2} \log |\Sigma| + \frac{1}{2} \sum_{i=1}^N (\mathbf{x}_i - \mu)^T \Sigma^{-1} (\mathbf{x}_i - \mu) \right) \Rightarrow \\
 dL &= \frac{N}{2} \frac{1}{|\Sigma|} d|\Sigma| + \frac{1}{2} \sum_{i=1}^N (\mathbf{x}_i - \mu)^T d\Sigma^{-1} (\mathbf{x}_i - \mu) \Rightarrow \\
 dL &= \frac{N}{2} \frac{1}{|\Sigma|} |\Sigma| \text{tr}(\Sigma^{-1} d\Sigma) - \frac{1}{2} \sum_{i=1}^N (\mathbf{x}_i - \mu)^T \Sigma^{-1} d\Sigma \Sigma^{-1} (\mathbf{x}_i - \mu) \Rightarrow \\
 dL &= \text{tr}(dL) = \frac{N}{2} \text{tr}(\Sigma^{-1} d\Sigma) - \frac{1}{2} \sum_{i=1}^N \text{tr}((\mathbf{x}_i - \mu)^T \Sigma^{-1} d\Sigma \Sigma^{-1} (\mathbf{x}_i - \mu)) \Rightarrow \\
 dL &= \frac{N}{2} \text{tr}(\Sigma^{-1} d\Sigma) - \frac{1}{2} \sum_{i=1}^N \text{tr}(\Sigma^{-1} (\mathbf{x}_i - \mu) (\mathbf{x}_i - \mu)^T \Sigma^{-1} d\Sigma) \Rightarrow \\
 dL &= \text{tr} \left( \frac{N}{2} \Sigma^{-1} d\Sigma - \frac{1}{2} \sum_{i=1}^N \Sigma^{-1} (\mathbf{x}_i - \mu) (\mathbf{x}_i - \mu)^T \Sigma^{-1} d\Sigma \right)
 \end{aligned} \tag{100}$$

继续对上面的公式变形：

$$dL = \text{tr} \left( \left[ \frac{N}{2} \Sigma^{-1} - \frac{1}{2} \sum_{i=1}^N \Sigma^{-1} (\mathbf{x}_i - \mu) (\mathbf{x}_i - \mu)^T \Sigma^{-1} \right] d\Sigma \right) \tag{101}$$

于是，根据公式 (86)，可得  $\frac{\partial L}{\partial \Sigma} = \frac{N}{2} \Sigma^{-1} - \frac{1}{2} \sum_{i=1}^N \Sigma^{-1} (\mathbf{x}_i - \mu) (\mathbf{x}_i - \mu)^T \Sigma^{-1}$ ，其中利用了  $\Sigma$  为对称矩阵的事实。令  $\frac{\partial L}{\partial \Sigma} = \mathbf{0}$ ，则：

$$\begin{aligned}
 \frac{N}{2} \Sigma^{-1} &= \frac{1}{2} \sum_{i=1}^N \Sigma^{-1} (\mathbf{x}_i - \mu) (\mathbf{x}_i - \mu)^T \Sigma^{-1} \\
 \frac{N}{2} \mathbf{I} &= \frac{1}{2} \sum_{i=1}^N \Sigma^{-1} (\mathbf{x}_i - \mu) (\mathbf{x}_i - \mu)^T \\
 \Sigma &= \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \mu) (\mathbf{x}_i - \mu)^T
 \end{aligned} \tag{102}$$

其中  $\mathbf{I}$  表示单位矩阵。

- 在多分类逻辑回归模型中，使用负对数似然函数表示交叉熵损失函数<sup>28</sup>：

$$L(\mathbf{W}) = -\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K y_i^{(k)} \log \sigma(\mathbf{w}_k^T \mathbf{x}_i) \tag{103}$$

<sup>28</sup>在《机器学习》课程系列之《逻辑回归》中，已经做了详细的介绍，并使用类似一元微积分的方法求解了梯度。

其中  $K$  表示类别个数,  $N$  表示样本点个数,  $\sigma(\cdot)$  表示 Softmax 函数。下面, 将使用本节介绍的方法求解其梯度。为此, 首先将多分类逻辑回归模型的损失函数表示为矩阵向量方式:

$$L(\mathbf{W}) = -\frac{1}{N} \sum_{i=1}^N \mathbf{y}_i^T \log \sigma(\mathbf{W} \mathbf{x}_i) \quad (104)$$

其中,  $\mathbf{W}$  表示  $K \times M$  权值矩阵, 每行表示每个分类的权值向量;  $\mathbf{x}_i$  表示  $M$  维列向量 (第  $i$  个样本点);  $\mathbf{y}_i$  表示对应的  $K$  维 One-Hot 列向量 (第  $i$  个样本点);  $\sigma(\mathbf{W} \mathbf{x}_i) = \frac{\exp(\mathbf{W} \mathbf{x}_i)}{\mathbf{1}^T \exp(\mathbf{W} \mathbf{x}_i)}$  表示 Softmax 函数逐元素运算后得到的列向量, 其中  $\mathbf{1}$  表示全 1 向量,  $\exp(\cdot)$  为逐元素运算指数函数。于是:

$$\begin{aligned} L(\mathbf{W}) &= -\frac{1}{N} \sum_{i=1}^N \mathbf{y}_i^T \log \sigma(\mathbf{W} \mathbf{x}_i) \Rightarrow \\ L(\mathbf{W}) &= -\frac{1}{N} \sum_{i=1}^N \mathbf{y}_i^T (\mathbf{W} \mathbf{x}_i - \mathbf{1} \log (\mathbf{1}^T \exp(\mathbf{W} \mathbf{x}_i))) \end{aligned} \quad (105)$$

对  $L(\mathbf{W})$  两端应用全微分算子, 可得:

$$\begin{aligned} dL &= d \left( -\frac{1}{N} \sum_{i=1}^N \mathbf{y}_i^T (\mathbf{W} \mathbf{x}_i - \mathbf{1} \log (\mathbf{1}^T \exp(\mathbf{W} \mathbf{x}_i))) \right) \Rightarrow \\ dL &= -\frac{1}{N} \sum_{i=1}^N \mathbf{y}_i^T (d(\mathbf{W} \mathbf{x}_i) - d(\mathbf{1} \log (\mathbf{1}^T \exp(\mathbf{W} \mathbf{x}_i)))) \Rightarrow \\ dL &= -\frac{1}{N} \sum_{i=1}^N \mathbf{y}_i^T \left( d\mathbf{W} \mathbf{x}_i - \mathbf{1} \frac{1}{\mathbf{1}^T \exp(\mathbf{W} \mathbf{x}_i)} d(\mathbf{1}^T \exp(\mathbf{W} \mathbf{x}_i)) \right) \Rightarrow \\ dL &= -\frac{1}{N} \sum_{i=1}^N \mathbf{y}_i^T \left( d\mathbf{W} \mathbf{x}_i - \mathbf{1} \frac{1}{\mathbf{1}^T \exp(\mathbf{W} \mathbf{x}_i)} \mathbf{1}^T (\exp(\mathbf{W} \mathbf{x}_i) \odot d\mathbf{W} \mathbf{x}_i) \right) \end{aligned} \quad (106)$$

利用公式  $\mathbf{y}_i^T \mathbf{1} = 1$  和  $\mathbf{1}^T (\mathbf{u} \odot \mathbf{v}) = \mathbf{u}^T \mathbf{v}$ , 继续变形:

$$\begin{aligned}
 dL &= -\frac{1}{N} \sum_{i=1}^N \left( \mathbf{y}_i^T d\mathbf{W} \mathbf{x}_i - \frac{1}{\mathbf{1}^T \exp(\mathbf{W} \mathbf{x}_i)} \exp(\mathbf{W} \mathbf{x}_i)^T d\mathbf{W} \mathbf{x}_i \right) \Rightarrow \\
 dL &= -\frac{1}{N} \sum_{i=1}^N (\mathbf{y}_i^T d\mathbf{W} \mathbf{x}_i - \sigma^T(\mathbf{W} \mathbf{x}_i) d\mathbf{W} \mathbf{x}_i) \Rightarrow \\
 dL &= \text{tr}(dL) = -\frac{1}{N} \sum_{i=1}^N [\text{tr}(\mathbf{y}_i^T d\mathbf{W} \mathbf{x}_i) - \text{tr}(\sigma^T(\mathbf{W} \mathbf{x}_i) d\mathbf{W} \mathbf{x}_i)] \Rightarrow \\
 dL &= -\frac{1}{N} \sum_{i=1}^N [\text{tr}(\mathbf{x}_i \mathbf{y}_i^T d\mathbf{W}) - \text{tr}(\mathbf{x}_i \sigma^T(\mathbf{W} \mathbf{x}_i) d\mathbf{W})] \Rightarrow \\
 dL &= \text{tr} \left( -\frac{1}{N} \sum_{i=1}^N \mathbf{x}_i (\mathbf{y}_i^T - \sigma^T(\mathbf{W} \mathbf{x}_i)) d\mathbf{W} \right)
 \end{aligned} \tag{107}$$

于是, 根据公式 (86), 可得:

$$\frac{\partial L}{\partial \mathbf{W}} = -\frac{1}{N} \sum_{i=1}^N [\mathbf{x}_i (\mathbf{y}_i^T - \sigma^T(\mathbf{W} \mathbf{x}_i))]^T = -\frac{1}{N} \sum_{i=1}^N (\mathbf{y}_i - \sigma(\mathbf{W} \mathbf{x}_i)) \mathbf{x}_i^T \tag{108}$$

- 神经网络: 函数的多层嵌套视角

神经网络对输入样本  $\mathbf{x}$  连续地应用非线性函数——第  $i-1$  层的输出是第  $i$  层的输入, 最后在终端获得相应的输出, 这是一种典型的数据流水线。从函数的角度看, 神经网络相当于多层函数嵌套。

现假设有一个 2 层神经网络, 输出层使用 Softmax 函数 ( $\sigma_2(\cdot)$ ) 作为激活函数, 它的前一级网络使用 Sigmoid 函数 ( $\sigma_1(\cdot)$ ) 作为激活函数。对于给定的数据集, 该神经网络的损失函数可表示为<sup>29</sup>:

$$L(\mathbf{W}_2, \mathbf{W}_1, \mathbf{b}_2, \mathbf{b}_1) = -\frac{1}{N} \sum_{i=1}^N \mathbf{y}_i^T \log \sigma_2(\mathbf{W}_2 \sigma_1(\mathbf{W}_1 \mathbf{x}_i + \mathbf{b}_1) + \mathbf{b}_2) \tag{109}$$

其中,  $\mathbf{x}_i$  表示  $S$  维列向量 (第  $i$  个样本点);  $\mathbf{W}_1$  表示  $M \times S$  权值矩阵;  $\mathbf{b}_1$  为  $M$  维列向量;  $\sigma_1(\mathbf{z}_1) = \frac{1}{1 + \exp(-\mathbf{z}_1)}$  为逐元素运算函数,  $\mathbf{z}_1 = \mathbf{W}_1 \mathbf{x} + \mathbf{b}_1$ ;  $\mathbf{W}_2$  表示  $K \times M$  权值矩阵, 每行表示每个分类的权值向量;  $\mathbf{b}_2$  为  $K$  维列向量;  $\sigma_2(\mathbf{z}_2) = \frac{\exp(\mathbf{z}_2)}{\mathbf{1}^T \exp(\mathbf{z}_2)}$  表示 Softmax 函数逐元素运算后得到的列向量, 其中  $\mathbf{1}$  表示全 1 向量,  $\exp(\cdot)$  为逐元素运算指数函数,  $\mathbf{z}_2 = \mathbf{W}_2 \sigma_1(\cdot) + \mathbf{b}_2$ ;  $\mathbf{y}_i$

<sup>29</sup>相当于在上例 (多分类逻辑回归模型) 的基础上添加了函数嵌套功能。

表示对应的  $K$  维 One-Hot 列向量 (第  $i$  个样本点)。试计算各权值矩阵与偏置向量的梯度<sup>30</sup>。

为便于表达, 令  $\mathbf{z}_i^{(1)} = \mathbf{W}_1 \mathbf{x}_i + \mathbf{b}_1$ ,  $\mathbf{z}_i^{(2)} = \mathbf{W}_2 \sigma_1(\mathbf{z}_i^{(1)}) + \mathbf{b}_2$ 。首先, 将损失函数 (公式 (109)) 表示为如下形式:

$$\begin{aligned} L &= -\frac{1}{N} \sum_{i=1}^N \mathbf{y}_i^T \log \sigma_2(\mathbf{z}_i^{(2)}) \Rightarrow \\ L &= -\frac{1}{N} \sum_{i=1}^N \mathbf{y}_i^T \left( \mathbf{z}_i^{(2)} - \mathbf{1} \log \left( \mathbf{1}^T \exp(\mathbf{z}_i^{(2)}) \right) \right) \Rightarrow \\ L &= -\frac{1}{N} \sum_{i=1}^N \left( \mathbf{y}_i^T \mathbf{z}_i^{(2)} - \log \left( \mathbf{1}^T \exp(\mathbf{z}_i^{(2)}) \right) \right) \end{aligned} \quad (110)$$

于是, 对上式两端应用全微分算子:

$$\begin{aligned} dL &= -\frac{1}{N} \sum_{i=1}^N \left( \mathbf{y}_i^T d\mathbf{z}_i^{(2)} - d \log \left( \mathbf{1}^T \exp(\mathbf{z}_i^{(2)}) \right) \right) \Rightarrow \\ dL &= -\frac{1}{N} \sum_{i=1}^N \left( \mathbf{y}_i^T d\mathbf{z}_i^{(2)} - \frac{1}{\mathbf{1}^T \exp(\mathbf{z}_i^{(2)})} \mathbf{1}^T d \exp(\mathbf{z}_i^{(2)}) \right) \Rightarrow \\ dL &= -\frac{1}{N} \sum_{i=1}^N \left( \mathbf{y}_i^T d\mathbf{z}_i^{(2)} - \frac{1}{\mathbf{1}^T \exp(\mathbf{z}_i^{(2)})} \mathbf{1}^T \left( \exp(\mathbf{z}_i^{(2)}) \odot d\mathbf{z}_i^{(2)} \right) \right) \Rightarrow \\ dL &= -\frac{1}{N} \sum_{i=1}^N \left( \mathbf{y}_i^T d\mathbf{z}_i^{(2)} - \frac{1}{\mathbf{1}^T \exp(\mathbf{z}_i^{(2)})} \exp^T(\mathbf{z}_i^{(2)}) d\mathbf{z}_i^{(2)} \right) \end{aligned} \quad (111)$$

继续变形, 得到:

$$\begin{aligned} dL &= -\frac{1}{N} \sum_{i=1}^N \left( \mathbf{y}_i^T d\mathbf{z}_i^{(2)} - \sigma_2^T(\mathbf{z}_i^{(2)}) d\mathbf{z}_i^{(2)} \right) \Rightarrow \\ dL &= -\frac{1}{N} \sum_{i=1}^N \left( \mathbf{y}_i^T - \sigma_2^T(\mathbf{z}_i^{(2)}) \right) d\mathbf{z}_i^{(2)} \end{aligned} \quad (112)$$

对上式两端应用迹运算, 得到:

$$\begin{aligned} dL &= \text{tr}(dL) = \text{tr} \left( -\frac{1}{N} \sum_{i=1}^N \left( \mathbf{y}_i^T - \sigma_2^T(\mathbf{z}_i^{(2)}) \right) d\mathbf{z}_i^{(2)} \right) \Rightarrow \\ dL &= -\frac{1}{N} \sum_{i=1}^N \text{tr} \left( \left( \mathbf{y}_i - \sigma_2(\mathbf{z}_i^{(2)}) \right)^T d\mathbf{z}_i^{(2)} \right) \end{aligned} \quad (113)$$

<sup>30</sup>实际上, 利用所谓的中间量概念, 可以建立层与层之间的梯度递推公式, 详见《人工智能》课程系列之《人工神经网络基础》, Chapter8-CN.pdf。在下文中, 也为矩阵与向量的梯度引入中间梯度量, 从而建立了相邻层间的梯度递推公式。

于是, 得到 (梯度) 中间量<sup>31</sup>:

$$\frac{\partial L}{\partial \mathbf{z}_i^{(2)}} = \mathbf{y}_i - \sigma_2 \left( \mathbf{z}_i^{(2)} \right) \quad (114)$$

可以将  $dL$  简化为:

$$dL = -\frac{1}{N} \sum_{i=1}^N \text{tr} \left( \left( \frac{\partial L}{\partial \mathbf{z}_i^{(2)}} \right)^T d\mathbf{z}_i^{(2)} \right) \quad (115)$$

将  $\mathbf{z}_i^{(2)} = \mathbf{W}_2 \sigma_1 \left( \mathbf{z}_i^{(1)} \right) + \mathbf{b}_2$  代入上式, 可得:

$$\begin{aligned} dL &= -\frac{1}{N} \sum_{i=1}^N \text{tr} \left( \left( \frac{\partial L}{\partial \mathbf{z}_i^{(2)}} \right)^T d \left( \mathbf{W}_2 \sigma_1 \left( \mathbf{z}_i^{(1)} \right) + \mathbf{b}_2 \right) \right) \Rightarrow \\ dL &= -\frac{1}{N} \sum_{i=1}^N \text{tr} \left( \left( \frac{\partial L}{\partial \mathbf{z}_i^{(2)}} \right)^T d\mathbf{W}_2 \sigma_1 \left( \mathbf{z}_i^{(1)} \right) \right) \\ &\quad - \frac{1}{N} \sum_{i=1}^N \text{tr} \left( \left( \frac{\partial L}{\partial \mathbf{z}_i^{(2)}} \right)^T \mathbf{W}_2 d\sigma_1 \left( \mathbf{z}_i^{(1)} \right) \right) \\ &\quad - \frac{1}{N} \sum_{i=1}^N \text{tr} \left( \left( \frac{\partial L}{\partial \mathbf{z}_i^{(2)}} \right)^T d\mathbf{b}_2 \right) \end{aligned} \quad (116)$$

从上式的第 1 项可以得到权值  $\mathbf{W}_2$  的梯度为:

$$\frac{\partial L}{\partial \mathbf{W}_2} = -\frac{1}{N} \sum_{i=1}^N \frac{\partial L}{\partial \mathbf{z}_i^{(2)}} \sigma_1^T \left( \mathbf{z}_i^{(1)} \right) \quad (117)$$

从第 3 项可以得到偏置  $\mathbf{b}_2$  的梯度为:

$$\frac{\partial L}{\partial \mathbf{b}_2} = -\frac{1}{N} \sum_{i=1}^N \frac{\partial L}{\partial \mathbf{z}_i^{(2)}} \quad (118)$$

对第 2 项继续变换, 得到:

$$\begin{aligned} & -\frac{1}{N} \sum_{i=1}^N \text{tr} \left( \left( \frac{\partial L}{\partial \mathbf{z}_i^{(2)}} \right)^T \mathbf{W}_2 d\sigma_1 \left( \mathbf{z}_i^{(1)} \right) \right) \Rightarrow \\ & -\frac{1}{N} \sum_{i=1}^N \text{tr} \left( \left( \frac{\partial L}{\partial \mathbf{z}_i^{(2)}} \right)^T \mathbf{W}_2 \left( \sigma_1' \left( \mathbf{z}_i^{(1)} \right) \odot d\mathbf{z}_i^{(1)} \right) \right) \Rightarrow \\ & -\frac{1}{N} \sum_{i=1}^N \text{tr} \left( \left( \mathbf{W}_2^T \frac{\partial L}{\partial \mathbf{z}_i^{(2)}} \odot \sigma_1' \left( \mathbf{z}_i^{(1)} \right) \right)^T d\mathbf{z}_i^{(1)} \right) \Rightarrow \end{aligned} \quad (119)$$

<sup>31</sup>需要注意的是, 在本小节中, 中间量的定义与常规有所不同, 它们之间相差一个“-”号。

于是, 得到 (梯度) 中间量:

$$\frac{\partial L}{\partial \mathbf{z}_i^{(1)}} = \mathbf{W}_2^T \frac{\partial L}{\partial \mathbf{z}_i^{(2)}} \odot \sigma'_1(\mathbf{z}_i^{(1)}) \quad (120)$$

对公式 (119) 进行简化, 可得:

$$-\frac{1}{N} \sum_{i=1}^N \text{tr} \left( \left( \frac{\partial L}{\partial \mathbf{z}_i^{(1)}} \right)^T d\mathbf{z}_i^{(1)} \right) \quad (121)$$

为得到权值  $\mathbf{W}_1$  与偏置  $\mathbf{b}_1$  的梯度, 将  $\mathbf{z}_i^{(1)} = \mathbf{W}_1 \mathbf{x}_i + \mathbf{b}_1$  代入上式, 可得:

$$\begin{aligned} & -\frac{1}{N} \sum_{i=1}^N \text{tr} \left( \left( \frac{\partial L}{\partial \mathbf{z}_i^{(1)}} \right)^T d(\mathbf{W}_1 \mathbf{x}_i + \mathbf{b}_1) \right) \Rightarrow \\ & -\frac{1}{N} \sum_{i=1}^N \text{tr} \left( \left( \frac{\partial L}{\partial \mathbf{z}_i^{(1)}} \right)^T d\mathbf{W}_1 \mathbf{x}_i \right) - \frac{1}{N} \sum_{i=1}^N \text{tr} \left( \left( \frac{\partial L}{\partial \mathbf{z}_i^{(1)}} \right)^T d\mathbf{b}_1 \right) \end{aligned} \quad (122)$$

由此, 得到权值  $\mathbf{W}_1$  和偏置  $\mathbf{b}_1$  的梯度分别为:

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{W}_1} &= -\frac{1}{N} \sum_{i=1}^N \frac{\partial L}{\partial \mathbf{z}_i^{(1)}} \mathbf{x}_i^T \\ \frac{\partial L}{\partial \mathbf{b}_1} &= -\frac{1}{N} \sum_{i=1}^N \frac{\partial L}{\partial \mathbf{z}_i^{(1)}} \end{aligned} \quad (123)$$

实际上, 使用  $i-1$  替换公式中的层索引 1, 使用  $i$  替换公式中的层索引 2, 上述公式均适用于任意相邻层。

## 4.2 公式 (30) 的证明<sup>32</sup>

证明: Fisher 准则可表示为:

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} \quad (124)$$

对上式两端应用全微分算子, 得到:

$$\begin{aligned} dJ &= d(\mathbf{w}^T \mathbf{S}_B \mathbf{w} (\mathbf{w}^T \mathbf{S}_W \mathbf{w})^{-1}) \\ &= d(\mathbf{w}^T) \mathbf{S}_B \mathbf{w} (\mathbf{w}^T \mathbf{S}_W \mathbf{w})^{-1} + \mathbf{w}^T \mathbf{S}_B d(\mathbf{w}) (\mathbf{w}^T \mathbf{S}_W \mathbf{w})^{-1} \\ &\quad - \mathbf{w}^T \mathbf{S}_B \mathbf{w} (\mathbf{w}^T \mathbf{S}_W \mathbf{w})^{-2} d(\mathbf{w}^T \mathbf{S}_W \mathbf{w}) \\ &= \frac{1}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} (d(\mathbf{w}^T) \mathbf{S}_B \mathbf{w} + \mathbf{w}^T \mathbf{S}_B d(\mathbf{w})) \\ &\quad - \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{(\mathbf{w}^T \mathbf{S}_W \mathbf{w})^2} d(\mathbf{w}^T \mathbf{S}_W \mathbf{w}) \end{aligned} \quad (125)$$

<sup>32</sup>本节的证明由周佳爽给出。

进一步地，应用迹技巧，可以得到：

$$\begin{aligned}
 dJ &= \text{tr}(dJ) \\
 &= \text{tr}\left(\frac{1}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} (d(\mathbf{w}^T) \mathbf{S}_B \mathbf{w} + \mathbf{w}^T \mathbf{S}_B d(\mathbf{w}))\right) \\
 &\quad - \text{tr}\left(\frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{(\mathbf{w}^T \mathbf{S}_W \mathbf{w})^2} d(\mathbf{w}^T \mathbf{S}_W \mathbf{w})\right)
 \end{aligned} \tag{126}$$

如果将上式直接展开，公式较长，不便于理解，这里将上式拆分为两部分。对于第一项：

$$\begin{aligned}
 &\text{tr}\left(\frac{1}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} (d(\mathbf{w}^T) \mathbf{S}_B \mathbf{w} + \mathbf{w}^T \mathbf{S}_B d(\mathbf{w}))\right) \\
 &= \frac{1}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} [\text{tr}(d(\mathbf{w}^T) \mathbf{S}_B \mathbf{w}) + \text{tr}(\mathbf{w}^T \mathbf{S}_B d(\mathbf{w}))] \\
 &= \frac{1}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} [\text{tr}((\mathbf{w}^T \mathbf{S}_B d(\mathbf{w}))^T) + \text{tr}(\mathbf{w}^T \mathbf{S}_B d(\mathbf{w}))] \\
 &= \frac{2}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} \text{tr}(\mathbf{w}^T \mathbf{S}_B d(\mathbf{w}))
 \end{aligned} \tag{127}$$

公式 (126) 的第二项化简过程如下：

$$\begin{aligned}
 &\text{tr}\left(\frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{(\mathbf{w}^T \mathbf{S}_W \mathbf{w})^2} d(\mathbf{w}^T \mathbf{S}_W \mathbf{w})\right) \\
 &= \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{(\mathbf{w}^T \mathbf{S}_W \mathbf{w})^2} \text{tr}(d(\mathbf{w}^T) \mathbf{S}_W \mathbf{w} + \mathbf{w}^T \mathbf{S}_W d(\mathbf{w})) \\
 &= \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{(\mathbf{w}^T \mathbf{S}_W \mathbf{w})^2} [\text{tr}(\mathbf{w}^T \mathbf{S}_W d(\mathbf{w}))^T + \text{tr}(\mathbf{w}^T \mathbf{S}_W d(\mathbf{w}))] \\
 &= 2 \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{(\mathbf{w}^T \mathbf{S}_W \mathbf{w})^2} \text{tr}(\mathbf{w}^T \mathbf{S}_W d(\mathbf{w}))
 \end{aligned} \tag{128}$$

于是得到：

$$\begin{aligned}
 dJ &= \text{tr}(dJ) \\
 &= \frac{2}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} \text{tr}(\mathbf{w}^T \mathbf{S}_B d(\mathbf{w})) - 2 \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{(\mathbf{w}^T \mathbf{S}_W \mathbf{w})^2} \text{tr}(\mathbf{w}^T \mathbf{S}_W d(\mathbf{w})) \\
 &= 2 \text{tr} \left[ \left( \frac{1}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} \mathbf{w}^T \mathbf{S}_B - \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{(\mathbf{w}^T \mathbf{S}_W \mathbf{w})^2} \mathbf{w}^T \mathbf{S}_W \right) d(\mathbf{w}) \right]
 \end{aligned} \tag{129}$$

由上式可得：

$$\begin{aligned}
 \frac{\partial J}{\partial \mathbf{w}} &= 2 \left( \frac{1}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} \mathbf{w}^T \mathbf{S}_B - \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{(\mathbf{w}^T \mathbf{S}_W \mathbf{w})^2} \mathbf{w}^T \mathbf{S}_W \right)^T \\
 &= 2 \left( \frac{1}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} \mathbf{S}_B \mathbf{w} - \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{(\mathbf{w}^T \mathbf{S}_W \mathbf{w})^2} \mathbf{S}_W \mathbf{w} \right)
 \end{aligned} \tag{130}$$



在  $J$  取极大值时, 有  $\frac{\partial J}{\partial \mathbf{w}} = \mathbf{0}$ 。于是, 得到:

$$\begin{aligned}\frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{(\mathbf{w}^T \mathbf{S}_W \mathbf{w})^2} \mathbf{S}_W \mathbf{w} &= \frac{1}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} \mathbf{S}_B \mathbf{w} \\ (\mathbf{w}^T \mathbf{S}_B \mathbf{w}) \mathbf{S}_W \mathbf{w} &= (\mathbf{w}^T \mathbf{S}_W \mathbf{w}) \mathbf{S}_B \mathbf{w}\end{aligned}\quad (131)$$

□

### 4.3 公式 (38) 的证明<sup>33</sup>

证明: 令类别  $c_1 = \frac{N}{N_1}$ , 类别  $c_2 = -\frac{N}{N_2}$ , 其中  $N_1$  和  $N_2$  分别表示训练数据集中类别  $c_1$  和  $c_2$  的个数。

最小平方法使用如下的误差函数:

$$E = \frac{1}{2} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - y_n)^2 \quad (132)$$

然后, 求解误差函数  $E$  关于参数  $w_0$  和  $\mathbf{w}$  的偏导数, 并令其分别等于 0 和  $\mathbf{0}$ , 可得:

$$\frac{\partial E}{\partial w_0} = \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - y_n) = 0 \quad (133)$$

$$\frac{\partial E}{\partial \mathbf{w}} = \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - y_n) \mathbf{x}_n = \mathbf{0} \quad (134)$$

首先计算  $w_0$ , 对公式 (133) 进行恒等变形, 计算得到:

$$\begin{aligned}Nw_0 &= -\sum_{n=1}^N \mathbf{w}^T \mathbf{x}_n + \sum_{n=1}^N y_n \\ &= -\mathbf{w}^T \sum_{n=1}^N \mathbf{x}_n + N_1 \frac{N}{N_1} - N_2 \frac{N}{N_2} \\ &= -\mathbf{w}^T \sum_{n=1}^N \mathbf{x}_n \Rightarrow \\ w_0 &= -\mathbf{w}^T \mathbf{m}\end{aligned}\quad (135)$$

对于公式 (134), 由于  $(\mathbf{w}^T \mathbf{x}_n + w_0 - y_n)$  是标量, 为计算方便, 可对  $\mathbf{x}_n$  进行转置处理, 得到:

$$\sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - y_n) \mathbf{x}_n^T = \mathbf{0} \quad (136)$$

<sup>33</sup>本节的证明由周佳爽给出。

将  $w_0$  带入到公式 (136) 的左端, 得到:

$$\begin{aligned}
 & \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n - \mathbf{w}^T \mathbf{m} - y_n) \mathbf{x}_n^T \\
 &= \sum_{n \in C1} \mathbf{w}^T (\mathbf{x}_n \mathbf{x}_n^T - \mathbf{m} \mathbf{x}_n^T) - \sum_{n \in C1} y_n \mathbf{x}_n^T + \sum_{n \in C2} \mathbf{w}^T (\mathbf{x}_n \mathbf{x}_n^T - \mathbf{m} \mathbf{x}_n^T) - \sum_{n \in C2} y_n \mathbf{x}_n^T \\
 &= \mathbf{w}^T \left( \sum_{n \in C1} \mathbf{x}_n \mathbf{x}_n^T - N_1 \mathbf{m} \mathbf{m}_1^T \right) - \frac{N}{N_1} N_1 \mathbf{m}_1^T + \mathbf{w}^T \left( \sum_{n \in C2} \mathbf{x}_n \mathbf{x}_n^T - N_2 \mathbf{m} \mathbf{m}_2^T \right) + \frac{N}{N_2} N_2 \mathbf{m}_2^T
 \end{aligned} \tag{137}$$

于是, 由公式 (136) 可进一步得到:

$$\mathbf{w}^T \left\{ \left( \sum_{n \in C1} \mathbf{x}_n \mathbf{x}_n^T + \sum_{n \in C2} \mathbf{x}_n \mathbf{x}_n^T \right) - \mathbf{m} (N_1 \mathbf{m}_1^T + N_2 \mathbf{m}_2^T) \right\} = N(\mathbf{m}_1^T - \mathbf{m}_2^T) \tag{138}$$

下面, 依据  $S_w$  定义, 对其进行公式变形:

$$\begin{aligned}
 S_W &= \sum_{n \in C1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^T + \sum_{n \in C2} (\mathbf{x}_n - \mathbf{m}_2)(\mathbf{x}_n - \mathbf{m}_2)^T \\
 &= \left( \sum_{n \in C1} \mathbf{x}_n \mathbf{x}_n^T + \sum_{n \in C2} \mathbf{x}_n \mathbf{x}_n^T \right) - \sum_{n \in C1} (\mathbf{m}_1 \mathbf{x}_n^T + \mathbf{x}_n \mathbf{m}_1^T) \\
 &\quad - \sum_{n \in C2} (\mathbf{m}_2 \mathbf{x}_n^T + \mathbf{x}_n \mathbf{m}_2^T) + \left( \sum_{n \in C1} \mathbf{m}_1 \mathbf{m}_1^T + \sum_{n \in C2} \mathbf{m}_2 \mathbf{m}_2^T \right) \\
 &= \left( \sum_{n \in C1} \mathbf{x}_n \mathbf{x}_n^T + \sum_{n \in C2} \mathbf{x}_n \mathbf{x}_n^T \right) - (N_1 \mathbf{m}_1 \mathbf{m}_1^T + N_1 \mathbf{m}_1 \mathbf{m}_1^T) \\
 &\quad - (N_2 \mathbf{m}_2 \mathbf{m}_2^T + N_2 \mathbf{m}_2 \mathbf{m}_2^T) + (N_1 \mathbf{m}_1 \mathbf{m}_1^T + N_2 \mathbf{m}_2 \mathbf{m}_2^T)
 \end{aligned} \tag{139}$$

移项后, 计算得到:

$$\left( \sum_{n \in C1} \mathbf{x}_n \mathbf{x}_n^T + \sum_{n \in C2} \mathbf{x}_n \mathbf{x}_n^T \right) = S_W + N_1 \mathbf{m}_1 \mathbf{m}_1^T + N_2 \mathbf{m}_2 \mathbf{m}_2^T \tag{140}$$

利用公式 (140), 对下列公式进行变形计算得到:

$$\begin{aligned}
 & \left( \sum_{n \in C1} \mathbf{x}_n \mathbf{x}_n^T + \sum_{n \in C2} \mathbf{x}_n \mathbf{x}_n^T \right) - \mathbf{m} (N_1 \mathbf{m}_1^T + N_2 \mathbf{m}_2^T) \\
 &= S_W + N_1 \mathbf{m}_1 \mathbf{m}_1^T + N_2 \mathbf{m}_2 \mathbf{m}_2^T - \mathbf{m} (N_1 \mathbf{m}_1^T + N_2 \mathbf{m}_2^T) \\
 &= S_W + N_1 \mathbf{m}_1 \mathbf{m}_1^T + N_2 \mathbf{m}_2 \mathbf{m}_2^T - \frac{1}{N} (N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2) (N_1 \mathbf{m}_1^T + N_2 \mathbf{m}_2^T) \\
 &= S_W + \frac{N_1 N_2}{N} \left[ \frac{N}{N_2} \mathbf{m}_1 \mathbf{m}_1^T + \frac{N}{N_1} \mathbf{m}_2 \mathbf{m}_2^T - \frac{N_1}{N_2} \mathbf{m}_1 \mathbf{m}_1^T - \frac{N_2}{N_1} \mathbf{m}_2 \mathbf{m}_2^T - \mathbf{m}_1 \mathbf{m}_2^T - \mathbf{m}_2 \mathbf{m}_1^T \right] \\
 &= S_W + \frac{N_1 N_2}{N} [\mathbf{m}_1 \mathbf{m}_1^T + \mathbf{m}_2 \mathbf{m}_2^T - \mathbf{m}_2 \mathbf{m}_1^T - \mathbf{m}_1 \mathbf{m}_2^T] \\
 &= S_W + \frac{N_1 N_2}{N} S_B
 \end{aligned} \tag{141}$$

将上式带入到公式 (138) 中, 可得:

$$\boldsymbol{w}^T \left\{ \boldsymbol{S}_W + \frac{N_1 N_2}{N} \boldsymbol{S}_B \right\} = N(\boldsymbol{m}_1^T - \boldsymbol{m}_2^T) \quad (142)$$

对等式两端进行转置, 可得到:

$$\left\{ \boldsymbol{S}_W + \frac{N_1 N_2}{N} \boldsymbol{S}_B \right\} \boldsymbol{w} = N(\boldsymbol{m}_1 - \boldsymbol{m}_2) \quad (143)$$

□

## 4.4 感知机实验

IRIS(鸢尾花) 数据集

由 Fisher 在 1936 年整理, 包含 4 个特征: Sepal.Length (花萼长度)、Sepal.Width (花萼宽度)、Petal.Length (花瓣长度)、Petal.Width (花瓣宽度), 它的特征值都为正浮点数, 单位为厘米。

目标值为鸢尾花的分类:

- Iris Setosa (山鸢尾): 0;
- Iris Versicolour (杂色鸢尾): 1;
- Iris Virginica (维吉尼亚鸢尾): 2;

```
In [1]: import pandas as pd
import numpy as np
from sklearn.datasets import load_iris
import matplotlib.pyplot as plt
%matplotlib inline
%config InlineBackend.figure_format = 'svg'

In [2]: from IPython.display import Image
Image(filename="iris.jpg",width=200, height=100)
```

Out[2]:



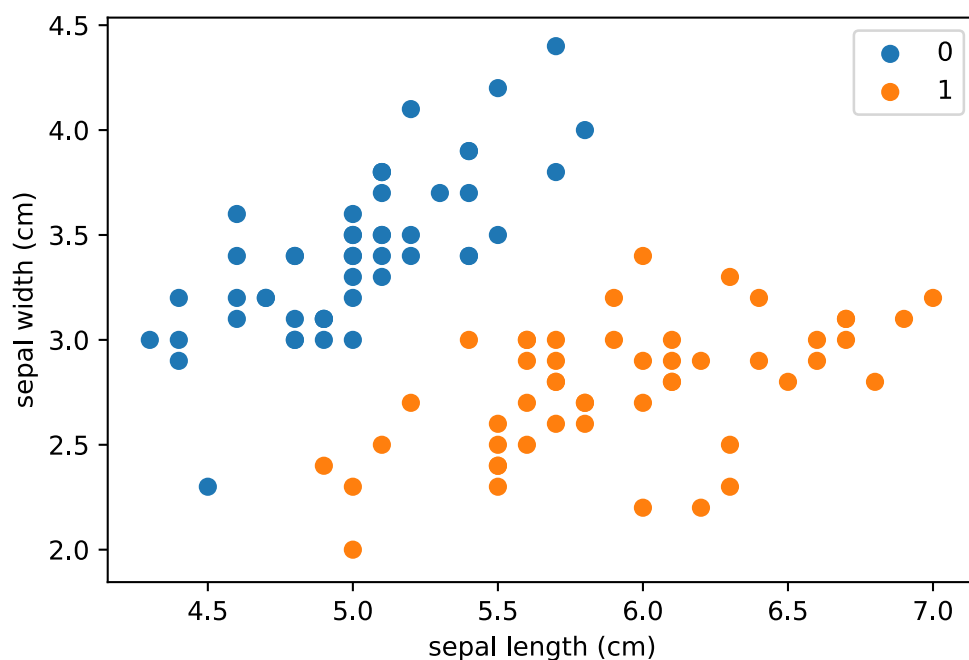
载入 sklearn 鸢尾花数据集

```
In [3]: iris = load_iris()
print(iris.keys()) # 数据集对象的属性
print(iris.data.shape) # 维度信息
print(iris.feature_names) # 特征名称
print(iris['target_names']) # 分类名称
print(iris['target']) # 分类值
```



显示 2 分类数据集（为便于可视化，只选取前 2 列数据，即花萼长度与宽度属性数据，并且只选取前 100 条记录，即 0 与 1 分类）

```
In [7]: plt.scatter(iris_df[:50]['sepal length (cm)'], iris_df[:50]
        ['sepal width (cm)'], label='0')
plt.scatter(iris_df[50:100]
        ['sepal length (cm)'], iris_df[50:100]['sepal width (cm)'], label='1')
plt.xlabel('sepal length (cm)')
plt.ylabel('sepal width (cm)')
plt.legend()
plt.savefig('PERCEPTRON_OUTPUT1.pdf', bbox_inches='tight')
```



提取训练数据

```
# 取出第 0 列、第 1 列、最后一列，前 100 条记录
In [8]: data = np.array(iris_df.iloc[:100, [0, 1, -1]])
        X, Y = data[:, :-1], data[:, -1]
        # 转换成适合感知机训练的分类: 1、-1
        Y = np.array([1 if y == 0 else -1 for y in Y])
        print(data.shape)
```

(100, 3)

定义感知机模型

```

In [9]: class Perceptron:
        def __init__(self):
            self.alpha = 0.05
            self.paras = []

        def w_x_b(self, x):
            return np.dot(self.W, x) + self.b

        # 随机梯度下降
        def fit(self, X, Y):
            self.W = np.ones(X.shape[1], dtype = np.float32)
            self.b = -5

            x0_mean = np.mean(X[0])
            x1_mean = np.mean(X[1])
            self.W[0] = x1_mean/x0_mean
            print(self.W, self.b)
            self.paras = [(self.W, self.b)]

            finish = False
            while not finish:
                finish = True
                for index in range(len(X)):
                    x = X[index]
                    y = Y[index]
                    if y * self.w_x_b(x) <= 0:
                        # 注意, 不能使用 += 运算符; 否则, 计算错误
                        self.W = self.W + self.alpha * y * x
                        self.b = self.b + self.alpha * y
                        finish = False

                self.paras.append((self.W, self.b))

            if finish:
                return 'Perceptron training completed!'

```

使用训练数据集训练感知机

```

In [10]: np.random.seed()
        model = Perceptron()
        model.fit(X, Y)

```

```

[ 0.91860467  1.          ] -5

```

```

Out[10]: 'Perceptron training completed!'

```

绘制训练动画

```

In [11]: import matplotlib.animation as animation
         from IPython.display import HTML
         # 动画播放, 如果是%matplotlib inline, 则会嵌入图片
         %matplotlib notebook

fig = plt.figure()
ax = plt.axes(xlim=(3, 7), ylim=(2, 5))

ax.plot(data[:50, 0], data[:50, 1], '+', label='0')
ax.plot(data[50:100, 0], data[50:100, 1], '.', label='1')
ax.set_xlabel('sepal length (cm)')
ax.set_ylabel('sepal width (cm)')
ax.set_title('Perceptron')
ax.legend(loc = 'upper left')

line, = ax.plot([], [], lw = 2)
def update(para):
    W = para[0]
    b = para[1]

    X_points = np.linspace(4, 7,10)
    Y_points = -(W[0] * X_points + b)/W[1]
    line.set_data(X_points, Y_points)
ani = animation.FuncAnimation(fig, update, model.paras)
# 视频播放
# 安装 FFMPEG 库
# 在 anaconda 环境下, 需要执行 “conda install -c conda-forge ffmpeg”
# 需要 FFMPEG.exe, 直接下载 BUILD, 解压到目录, 例如 E:\ffmpeg\bin
# 然后在环境变量 PATH 中添加 “E:\ffmpeg\bin”,
# 路径名称需要根据实际情况做相应的调整
HTML(ani.to_html5_video())

```

```

In [12]: ani.save('Fitting-duxiaoqin.mp4')

```

绘制训练结果

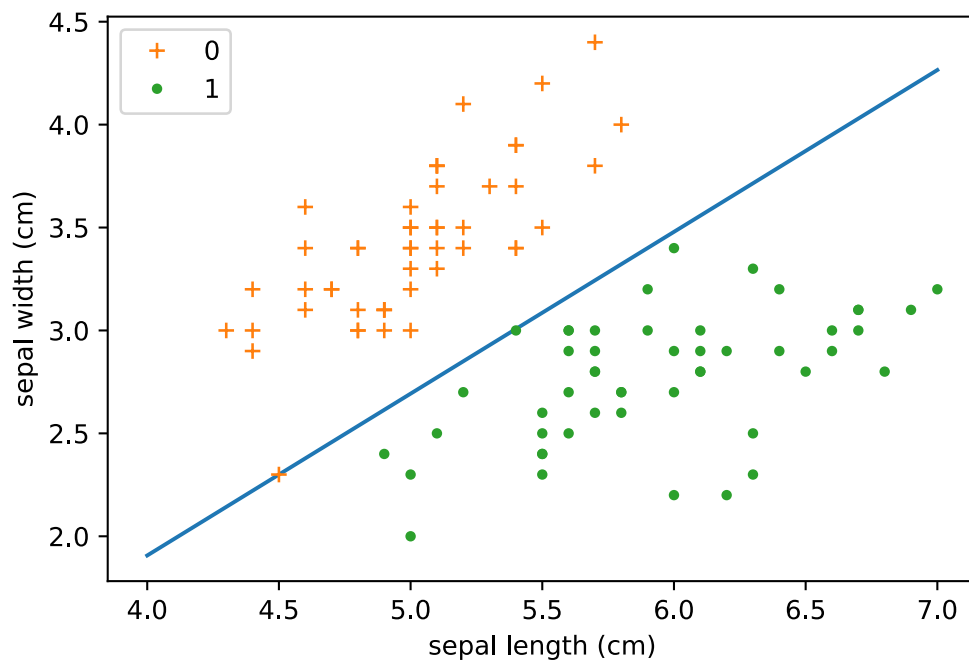
```

In [13]: %matplotlib inline
         X_points = np.linspace(4, 7,10)
         Y_points = -(model.W[0] * X_points + model.b)/model.W[1]
         plt.plot(X_points, Y_points)

         plt.plot(data[:50, 0], data[:50, 1], '+', label='0')
         plt.plot(data[50:100, 0], data[50:100, 1], '.', label='1')
         plt.xlabel('sepal length (cm)')
         plt.ylabel('sepal width (cm)')
         plt.legend()
         plt.savefig('PERCEPTRON_OUTPUT2.pdf', bbox_inches='tight')

```





### Perceptron (sklearn) 实验

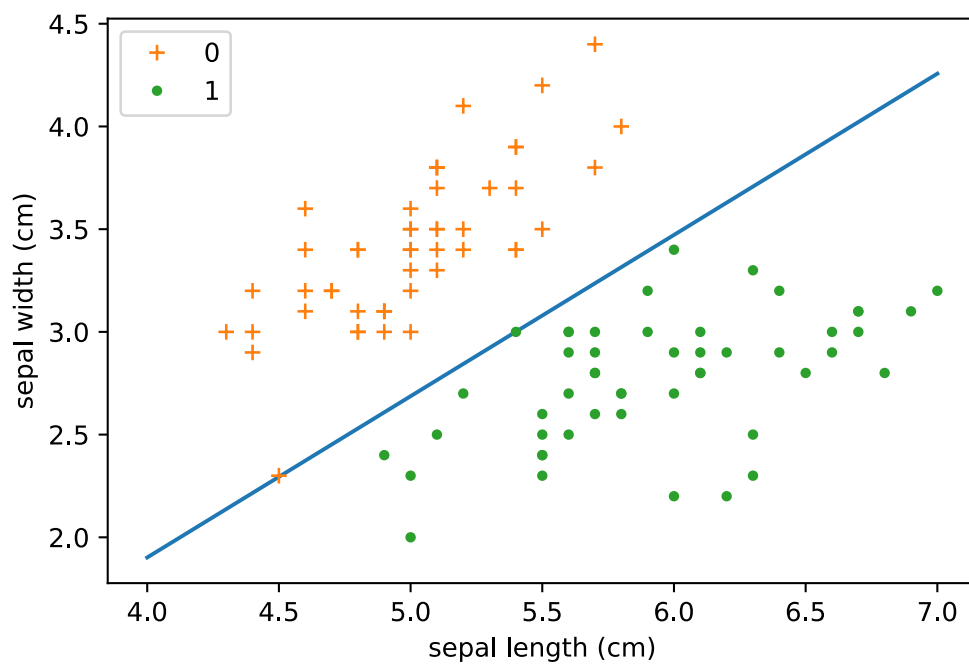
```
In [14]: from sklearn.linear_model import Perceptron
```

```
In [15]: model = Perceptron(max_iter = 1000)
          model.fit(X, Y)
```

```
Out[15]: Perceptron(alpha=0.0001, class_weight=None, eta0=1.0, fit_intercept=True,
                    max_iter=1000, n_iter=None, n_jobs=1, penalty=None, random_state=0,
                    shuffle=True, tol=None, verbose=0, warm_start=False)
```

```
In [16]: X_points = np.linspace(4, 7, 10)
          Y_points = -(model.coef_[0][0] * X_points + model.intercept_) /
                    model.coef_[0][1]
          plt.plot(X_points, Y_points)

          plt.plot(data[:50, 0], data[:50, 1], '+', label='0')
          plt.plot(data[50:100, 0], data[50:100, 1], '.', label='1')
          plt.xlabel('sepal length (cm)')
          plt.ylabel('sepal width (cm)')
          plt.legend()
          plt.savefig('PERCEPTRON_OUTPUT3.pdf', bbox_inches='tight')
```



## 5 参考文献

1. 周志华。《机器学习》，清华大学出版社，2016 年 1 月第 1 版。
2. Christopher M. Bishop. Pattern Recognition and Machine Learning. Springer, 2006.
3. Kevin P. Murphy. Machine Learning: A Probabilistic Perspective, The MIT Press, 2012.
4. 李航。《统计学习方法》，清华大学出版社，2019 年 5 月第 2 版。
5. Vladimir N. Vapnik。《统计学习理论的本质》，张学工译，2000 年，清华大学出版社。
6. Kaare Brandt Petersen, Michael Syskind Pedersen. The Matrix Cookbook, Nov. 15, 2012. <http://matrixcookbook.com>.
7. 矩阵求导 I. <https://zhuanlan.zhihu.com/p/24709748>.
8. 矩阵求导 II. <https://zhuanlan.zhihu.com/p/24863977>.
9. Fisher 线性判别与线性感知机。 <https://zhuanlan.zhihu.com/p/88657842>.
10. <https://github.com/wzyonggege/statistical-learning-method>;