

机器学习 - 周志华 - 线性模型

分子布局

如果我们设 $\mathbf{y} = (y_1, y_2, \dots, y_m)^T$, 那么对标量 x 求导的符号如下:

$$\frac{\partial \mathbf{y}}{\partial x} = \left(\frac{\partial y_1}{\partial x}, \dots, \frac{\partial y_m}{\partial x} \right)^T$$

这种符号也称为分子布局符号. 也就是对每个位置进行标量求导. 相反的即称为分母布局符号.

反过来就是相当于求梯度: 对 $\mathbf{x} = (x_1, \dots, x_n)^T$,

$$\frac{\partial y}{\partial \mathbf{x}} = \left(\frac{\partial y}{\partial x_1}, \dots, \frac{\partial y}{\partial x_n} \right)$$

注意这个求导得到的是列向量, 一些经典的标量求导:

❓ 例: 范数求导

对 $\|\mathbf{x}\|^2$ 求导.

使用欧氏范数有

$$\|\mathbf{x}\|^2 = \sum_{k=1}^n x_k^2$$

求导有

$$\frac{\partial \|\mathbf{x}\|^2}{\partial \mathbf{x}} = (2x_k) = 2\mathbf{x}^T$$

□

❓ 例: 内积求导

对 $\langle \mathbf{u}, \mathbf{v} \rangle$ 求导, \mathbf{u}, \mathbf{v} 均为 \mathbf{x} 的函数.

首先 $\langle \mathbf{u}, \mathbf{v} \rangle = \mathbf{u}^T \mathbf{v} = \sum_{k=1}^n u_k v_k$. 那么我们求导有

$$\begin{aligned}
\frac{\partial \mathbf{u}^T \mathbf{v}}{\partial \mathbf{x}} &= \left(\sum_{k=1}^n \left[v_k \frac{\partial u_k}{\partial x_i} + u_k \frac{\partial v_k}{\partial x_i} \right] \right)_i \\
&= \left(\sum_{k=1}^n v_k \frac{\partial u_k}{\partial x_i} \right)_i + \left(\sum_{k=1}^n u_k \frac{\partial v_k}{\partial x_i} \right)_i \\
&= \mathbf{v}^T \frac{\partial \mathbf{u}}{\partial \mathbf{x}} + \mathbf{u}^T \frac{\partial \mathbf{v}}{\partial \mathbf{x}}
\end{aligned}$$

□

向量对向量求导

如果我们记 $\mathbf{x} = (x_1, x_2, \dots, x_n)^T$, $\mathbf{y} = (y_1, y_2, \dots, y_m)^T$, 那么此时有向量对向量求导:

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \dots & \frac{\partial y_1}{\partial x_n} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \dots & \frac{\partial y_2}{\partial x_n} \\ \vdots & \vdots & & \vdots \\ \frac{\partial y_m}{\partial x_1} & \frac{\partial y_m}{\partial x_2} & \dots & \frac{\partial y_m}{\partial x_n} \end{pmatrix}$$

相当于就是说: 向量对向量求导将会拉伸一个维度到矩阵的维度, 我们看接下来的例子就可以大概明白其原理.

下面的 \mathbf{a} 和 \mathbf{A} 不是 \mathbf{x} 的函数, 均为常量向量、矩阵. 于是有

- $\frac{\partial}{\partial \mathbf{x}}(\mathbf{a} \mathbf{u}) = \mathbf{a} \frac{\partial \mathbf{u}}{\partial \mathbf{x}}.$
- $\frac{\partial \mathbf{x}}{\partial \mathbf{x}} = \mathbf{I}.$
- $\frac{\partial}{\partial \mathbf{x}}(\mathbf{A} \mathbf{x}) = \mathbf{A}.$
- $\frac{\partial}{\partial \mathbf{x}}(\mathbf{x}^T \mathbf{A}) = \mathbf{A}^T.$

我们可以拓展上述的内容, 矩阵对向量求导将会得到一个三维的**张量**, 以此类推有更高维的内容.

链式法则和自动求导

向量、矩阵求导也符合链式法则, 我们也和正常求导一样拆开即可.

例：矩阵求导链式法则

设 $\mathbf{X} \in \mathbb{R}^{m \times n}$, $\mathbf{w} \in \mathbb{R}^n$, $\mathbf{y} \in \mathbb{R}^m$, 对如下函数求导:

$$z = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$$

我们依次分解:

$$\mathbf{a} = \mathbf{X}\mathbf{w}, \mathbf{b} = \mathbf{a} - \mathbf{y}, z = \|\mathbf{b}\|^2$$

求导有:

$$\begin{aligned} \frac{\partial z}{\partial \mathbf{w}} &= \frac{\partial z}{\partial \mathbf{b}} \frac{\partial \mathbf{b}}{\partial \mathbf{a}} \frac{\partial \mathbf{a}}{\partial \mathbf{w}} \\ &= \frac{\partial \|\mathbf{b}\|^2}{\partial \mathbf{b}} \frac{\partial (\mathbf{a} - \mathbf{y})}{\partial \mathbf{a}} \frac{\partial \mathbf{X}\mathbf{w}}{\partial \mathbf{w}} \\ &= 2\mathbf{b}^T \mathbf{I} \mathbf{X} = 2\mathbf{b}^T \mathbf{X} \\ &= 2(\mathbf{X}\mathbf{w} - \mathbf{y})^T \mathbf{X} \end{aligned}$$

□

回归任务

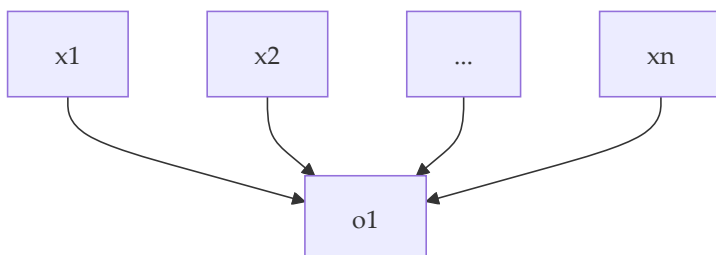
回归任务是通过输入去预测连续数值的任务, 线性回归就是通过如下的 n 维输入:

$$\mathbf{x} = (x_1, x_2, \dots, x_n)^T$$

计算出 n 维权重 $\mathbf{w} = (w_1, w_2, \dots, w_n)^T$ 和标量偏差 b , 从而可得输出为:

$$y = \langle \mathbf{w}, \mathbf{x} \rangle + b$$

线性模型可以看作是一个单层的神经网络:



上面的第一层为输入层, 第二层为输出层.

衡量回归质量

考虑 MSE :

$$\ell(y, \hat{y}) = \frac{1}{2}(y - \hat{y})^2$$

线性回归任务可以计算其解析解，我们先将训练损失计算为：

$$\ell(\mathbf{X}, \mathbf{y}, \mathbf{w}, b) = \frac{1}{2n} \sum_{i=1}^n (y_i - \langle \mathbf{x}_i, \mathbf{w} - b \rangle)^2 = \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\mathbf{w} - b\|^2$$

那么我们实际上就是要求出如下参数：

$$\mathbf{w}^*, \mathbf{b}^* = \arg \min_{\mathbf{w}, b} \ell(\mathbf{X}, \mathbf{y}, \mathbf{w}, b)$$

这实际上是一个优化类问题，其解法也比较简单，首先我们将偏差加入到权重当中，我们将偏差对应的项视为新增一个变量 1，权重为 b ，也就是说

$$\mathbf{X} \rightarrow (\mathbf{X}, 1); \mathbf{w} \rightarrow (\mathbf{w}, b)$$

因此有

$$\ell(\mathbf{X}, \mathbf{y}, \mathbf{w}) = \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2$$

求导后有

$$\frac{\partial}{\partial \mathbf{w}} \ell(\mathbf{X}, \mathbf{y}, \mathbf{w}) = \frac{1}{n} (\mathbf{y} - \mathbf{X}\mathbf{w})^T \mathbf{X}$$

损失函数是凸函数，上述导数取 0 即有最值，可得

$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X} \mathbf{y}$$

这是唯一解。注意这里我们默认了 $\mathbf{X}^T \mathbf{X}$ 是可逆的，

线性模型的变化

广义线性模型

若函数 $g(x)$ 为单调可微函数，则

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

则称该线性模型为**广义线性模型**。

线性模型的分类任务：对数几率回归

我们讨论的是线性回归模型，但是如果任务是分类的怎么办？考虑二分类任务，其输出标记 $y \in \{0, 1\}$ ，而线性回归模型产生的预测值 $z = \mathbf{w}^T \mathbf{x} + b$ 是实值，于是，我们需将实值转换为 0/1 值。最理想的是“单位阶跃函数”(unit-step function)：

$$y = \begin{cases} 0, & x < 0 \\ 0.5, & x = 0 \\ 1, & x > 0 \end{cases}$$

但是这种函数不可微，因此我们需要选取一种类似的可微函数。