

DW_共读AI圣经：第四章 Part 1 线性回归

第一章中有一个核心思想：

机器学习就是从数据中学习概率

这里的概率需要通过函数（抑或是称为模型）来求解，那么“线性”就是最简单的函数关系。通过线性模型的数学原理推导，我们就可以看到机器学习从数学角度来说到底在做什么。虽然更复杂的回归模型我们已经无法人力计算，但是从思路上一致的。这就是学习本章的意义：以线性回归模型为例，探究机器学习的数学本质。

目录

- [细节推导和概念引入](#)
 - [问题引入](#)
 - [似然函数和误差函数的关系推导细节](#)
 - [似然函数的求解](#)
- [补充信息](#)
 - [其他常见的基函数](#)

细节推导和概念引入

问题引入

$$y(x, w) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(x)$$

这是一个简单的线性函数可以用来描述一个线性回归模型，其中含有自变量的部分叫做**基函数**。（常见基函数在后面介绍）

$$y(x, w) = \sum_{j=0}^{M-1} w_j \phi_j(x)$$

加入一个偏置参数就可以统一为以上函数，就可以看做是一个简单的神经网络。

那么核心问题就是如何确定该函数中的参数 w

似然函数和误差函数的关系推导细节

需要预测的目标变量 t 是由包含高斯噪声的函数 $y(x, w)$ 给出的

$$t = y(\mathbf{x}, \mathbf{w}) + \varepsilon$$

其中 ε 是方差 σ^2 的零均值高斯随机变量，那么预测值的概率求解就可以是

$$p(t \mid \mathbf{x}, \mathbf{w}, \sigma^2) = \mathcal{N}(t \mid y(\mathbf{x}, \mathbf{w}), \sigma^2)$$

当我们有多个观测数据时，就可以把 t 看成是一个列向量，由于数据点都是从分布式中独立产生的，那么就可以得到似然函数

$$p(\mathbf{t} \mid \mathbf{X}, \mathbf{w}, \sigma^2) = \prod_{n=1}^N \mathcal{N}(t_n \mid \mathbf{w}^T \phi(\mathbf{x}_n), \sigma^2)$$

以下证明最大化似然函数就是最小化误差函数，可以由此确定参数的值：

证明: 在 Gaussian 噪声分布下, $t = y(x, w) + \varepsilon$
 (不属于严格最大似然函数等价的等价证明)

假设观测值 t_n 服从高斯分布, 其均值为模型预测值 $w^T \phi(x_n)$, 方差为 σ^2 :

$$P(t_n | x_n, w, \sigma^2) = N(t_n | w^T \phi(x_n), \sigma^2).$$

① 假设 n 个样本独立同分布: (独立条件).

$$P(t | x_n, w, \sigma^2) = \prod_{n=1}^N N(t_n | w^T \phi(x_n), \sigma^2)$$

② 取对数 (运算技巧).

$$\ln p(t | x_n, w, \sigma^2) = \sum_{n=1}^N \ln [N(t_n | w^T \phi(x_n), \sigma^2)]$$

③ 展开高斯分布的概率密度函数.

$$\ln N(\cdot) = \ln \left[\frac{1}{\sqrt{2\pi\sigma^2}} \right] + \ln \left[\exp \left(-\frac{(t_n - w^T \phi(x_n))^2}{2\sigma^2} \right) \right]$$

$$\text{④ 简化对数: } \ln N(\cdot) = -\frac{1}{2} \ln(2\pi) - \frac{1}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2} \sum_{n=1}^N (t_n - w^T \phi(x_n))^2$$

(涉及对数的基本运算)

⑤ 代入原式.

$$\begin{aligned} \ln p(t | x_n, w, \sigma^2) &= -\frac{N}{2} \ln(2\pi) - \frac{N}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2} \sum_{n=1}^N (t_n - w^T \phi(x_n))^2 \\ &= -\frac{N}{2} \ln(2\pi) - \frac{N}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2} E_0(w) \end{aligned}$$

故最大化似然函数
 需要最小化误差函数.

前一个 task
 讲到的
 误差函数.

似然函数的求解

明确目标：求解似然函数的最大值，下面就是数学计算的问题了。

我们选择求梯度的方式来求解：

$$\nabla_w \ln p(t | \mathbf{X}, \mathbf{w}, \sigma^2) = \frac{1}{\sigma^2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\} \phi(\mathbf{x}_n)^T$$

用梯度的方法求解最大似然。(以求解 \mathbf{w} 为例)。

$$\begin{aligned} \nabla_w \ln p(t | \mathbf{X}, \mathbf{w}, \sigma^2) &= \nabla_w \left[-\frac{1}{2\sigma^2} E_D(\mathbf{w}) \right] \text{ 实际是求偏导的运算} \\ &= \nabla_w \left[-\frac{1}{2\sigma^2} \sum_{n=1}^N \frac{(t - \mathbf{w}^T \phi(\mathbf{x}_n))^2}{2} \right] \\ &= \nabla_w \left\{ -\frac{1}{2\sigma^2} \sum_{n=1}^N \left[\frac{t^2}{x} + (\mathbf{w}^T \phi(\mathbf{x}_n))^2 - 2t \mathbf{w}^T \phi(\mathbf{x}_n) \right] \right\} \\ &\text{求 } \mathbf{w} \text{ 的偏导} \quad \nabla_w (\mathbf{w}^T \mathbf{a}) = \mathbf{a}, \quad \nabla_w (\mathbf{w}^T \mathbf{A} \mathbf{w}) = (\mathbf{A} + \mathbf{A}^T) \mathbf{w} \\ &= -\frac{1}{2\sigma^2} \sum_{n=1}^N \left\{ 2\phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \mathbf{w} - 2t_n \phi(\mathbf{x}_n) \right\} \text{ (一些向量运算)} \\ &= -\frac{1}{2\sigma^2} \sum_{n=1}^N \left[t_n \phi(\mathbf{x}_n) - \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \mathbf{w} \right] \\ &= \frac{1}{2\sigma^2} \sum_{n=1}^N \left[t_n - \mathbf{w}^T \phi(\mathbf{x}_n) \right] \phi(\mathbf{x}_n) \quad \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \mathbf{w} = \phi(\mathbf{x}_n) (\phi(\mathbf{x}_n)^T \mathbf{w}) \text{ (同成分解)} \end{aligned}$$

(后续步骤书中已给出)

当梯度设为0的时候，就可以进行求解：

$$0 = \sum_{n=1}^N t_n \phi(\mathbf{x}_n)^T - \mathbf{w}^T \left(\sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \right)$$

$$\mathbf{w}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

同理也可以求得其他参数。

补充信息

其他常见的基函数

在深度学习中，**基函数**通常指神经网络中用于引入非线性的**激活函数**。它们是神经网络能够学习复杂模式的关键组件。以下是几种常用基函数：

1. Sigmoid (Logistic Function)

- 公式:
$$\sigma(z) = 1 / (1 + \exp(-z))$$
 - 输出范围: $(0, 1)$
 - 特点:
 - 将输入压缩到0-1之间，适合输出概率（二分类输出层）。
 - 历史重要，但现代深层网络较少使用。
 - 缺点:
 - **梯度消失**：当输入绝对值较大时，梯度接近0，导致深层网络训练困难。
 - **非零中心输出**：可能导致后续层输入分布偏移。
-

2. Tanh (Hyperbolic Tangent)

- 公式:
$$\tanh(z) = (\exp(z) - \exp(-z)) / (\exp(z) + \exp(-z))$$

(或 $2 * \text{sigmoid}(2z) - 1$)
 - 输出范围: $(-1, 1)$
 - 特点:
 - 类似Sigmoid，但输出以0为中心。
 - 梯度比Sigmoid稍强（最大梯度为1）。
 - 缺点:
 - 仍然存在**梯度消失**问题（尤其在绝对值大的区域）。
 - 常用于RNN、LSTM等循环网络。
-

3. ReLU (Rectified Linear Unit)

- 公式:
 $\text{ReLU}(z) = \max(0, z)$
 - 输出范围: $[0, \infty)$
 - 特点:
 - 计算高效: 仅需比较和阈值操作。
 - 缓解梯度消失: 正区间梯度恒为1。
 - 稀疏激活: 约50%神经元在训练中被激活。
 - 现代深度网络最常用。
 - 缺点:
 - **Dying ReLU问题**: 输入为负时梯度为0, 神经元可能永久"死亡"。
 - 输出非零中心。
-

4. Leaky ReLU

- 公式:
 $\text{LeakyReLU}(z) = \max(\alpha z, z)$
(α 是一个小的正数, 如0.01)
 - 输出范围: $(-\infty, \infty)$
 - 特点:
 - 解决ReLU的"死亡"问题: 负区间有微小梯度(α)。
 - 保留了ReLU在正区间的优点。
 - 变种:
 - **Parametric ReLU (PReLU)**: α 作为可学习参数。
-

5. ELU (Exponential Linear Unit)

- 公式:
 $\text{ELU}(z) = \{ z, \text{ if } z > 0; \alpha(\exp(z) - 1), \text{ if } z \leq 0 \}$
(α 通常设为1)
- 输出范围: $(-\alpha, \infty)$
- 特点:
 - 负区间平滑渐近到 $-\alpha$, 缓解Dying ReLU问题。
 - 输出接近零中心化。
 - 理论上可能比ReLU有更好的表现。
- 缺点:

- 计算量稍大（涉及指数运算）。

6. Swish

- 公式:
$$\text{Swish}(z) = z * \text{sigmoid}(\beta z)$$

(β 可以是常数或可学习参数，常取1)
- 输出范围: $(-\infty, \infty)$
- 特点:
 - 由Google提出，在部分任务上表现优于ReLU。
 - 平滑且非单调（负区间有微小"凸起"）。
 - 负值不会被完全抑制。

7. GELU (Gaussian Error Linear Unit)

- 公式 (近似):
$$\text{GELU}(z) \approx 0.5 * z * (1 + \tanh(\sqrt{2/\pi} * (z + 0.044715z^3)))$$

(或使用Sigmoid近似)
- 输出范围: $(-\infty, \infty)$
- 特点:
 - 受随机正则化思想启发（如Dropout）。
 - 在**Transformer模型**（如BERT, GPT）中广泛使用。
 - 表现常优于ReLU/ELU。

选择建议 & 总结

函数	适用场景	主要优势	主要劣势
Sigmoid	二分类输出层	输出概率直观	梯度消失严重
Tanh	RNN/LSTM隐藏层	零中心输出	梯度消失
ReLU	大多数CNN/MLP隐藏层 (默认首选)	计算高效，缓解梯度消失	Dying ReLU问题
Leaky ReLU/ELU	担心Dying ReLU时	解决负区间问题	计算稍复杂

函数	适用场景	主要优势	主要劣势
Swish/GELU	追求更高性能 (尤其 Transformer)	平滑，理论性质好	计算量最大

核心作用：引入非线性，使神经网络能够逼近任意复杂函数。没有它们，深度网络将退化为线性模型。