

Abstract (摘要)

幂娜桑，空你吉哇，我是生而为弟。

在学习NLP的过程中，我对于公式的推导不是很擅长完全不会，因此我决定从头学习机器学习的理论推导。

本期主要分为两个部分：

- 第一部分将从矩阵，几何，概率三个视角，对线性回归-最小二乘法的闭式解进行推导，并提供参考代码。
- 第二部分将从矩阵和概率两个视角，对带正则化的最小二乘法的闭式解进行推导，并构造一个完整的线性回归类，同时实现闭式解求法与梯度下降解法。

Introduction (介绍)

线性回归模型是利用线性函数对一个或多个自变量和因变量 (y) 之间关系进行拟合的模型。

目标变量 (y) 为连续数值型，如：房价，人数，降雨量 回归模型是寻找一个输入变量到输出变量之间的映射函数。

回归问题的学习等价于函数拟合：使用一条函数曲线使其很好的拟合已知数据且能够预测未知数据。

回归问题分为模型的学习和预测两个过程。基于给定的训练数据集构建一个模型，根据新的输入数据预测相应的输出。

1.Algorithm (算法)

不带正则化

A.矩阵视角

注：一般情况下，我们讨论的向量都是列向量，因此推导过程中为保证矩阵的形状，会大量使用转置符

已知数据集 $D = \{(x_1, y_1), (x_2, y_2) \dots (x_n, y_n)\}$

其中 $x_i \in R^p$, $y_i \in R$, $i = 1, 2, \dots, n$

$$X = (x_1, x_2, \dots, x_n)^T = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ x_{n1} & x_{n2} & \dots & x_{np} \end{pmatrix}_{n \times p} \quad (33)$$

$$Y = (y_1, y_2, \dots, y_n)^T_{n \times 1} \quad (34)$$

这是我们建立的模型： $f(w) = w^T * x + w_0 * x_0$ 。

一般令 $x_0 = 1$, 而 $b = w_0 * x_0$, b 是偏置(bias), w 为权重weight, 下面为了推导的方便，我们将 w_0 并入 w 中， x_0 并入 X 中

因此模型更为 $f(w) = w^T * x$

最小二乘法的损失函数为：

$$L(w) = \sum_{i=1}^{i=n} \|y_i - w^T * x_i\|_2^2 \quad (35)$$

$$= (y_1 - w^T * x_1 \quad y_2 - w^T * x_2 \quad \dots \quad y_n - w^T * x_n) \begin{pmatrix} y_1 - w^T * x_1 \\ y_2 - w^T * x_2 \\ \vdots \\ y_n - w^T * x_n \end{pmatrix} \quad (36)$$

$$= (Y^T - w^T * X^T)(Y^T - w^T * X^T)^T \quad (37)$$

$$= (Y^T - w^T * X^T)(Y - X * w) \quad (38)$$

$$= Y^T * Y - w^T * X^T * Y - Y^T * X * w + w^T * X^T * X * w \quad (39)$$

仔细观察发现第二三项是互相转置的，而观察它的矩阵形状: $(1, p) * (p, n) * (n, 1) = (1, 1)$

得知这两项为标量，而标量的转置还是本身，因此可将两项合并，得

$$L(w) = Y^T * Y - 2 * w^T * X^T * Y + w^T * X^T * X * w \quad (40)$$

因此 $\hat{w} = \operatorname{argmin}(L(w))$

下面要求出 $L(w)$ 的最小值，对 $L(w)$ 求导

可以看到式子共三项，第一项与 w 无关，可以去掉。那么剩余两项就要涉及到矩阵求导了

关于矩阵求导，笔者推荐一位博主的三篇文章(比教科书还详细，严谨，每个公式都有证明)

[矩阵求导的本质与分子布局、分母布局的本质（矩阵求导——本质篇） - lterator的文章 - 知乎](#)

[矩阵求导公式的数学推导（矩阵求导——基础篇） - lterator的文章 - 知乎](#)

[矩阵求导公式的数学推导（矩阵求导——进阶篇） - lterator的文章 - 知乎](#)

下面为上述两项的导数求解过程：

因为 X, Y 为常数矩阵，因此可直接求出导数，但因为是对 w 求导，因此要对结果进行转置

$$\frac{d(2 * w^T * X^T * Y)}{dw} = 2 * X^T * Y \quad (41)$$

下面求解第三项

$$d(w^T * X^T * X * w) = \operatorname{tr}(d(w^T * X^T * X * w)) = \operatorname{tr}(X^T * X * d(w^T * w)) \quad (42)$$

$$= \operatorname{tr}(X^T * X * (d(w^T) * w + w^T * d(w))) = \operatorname{tr}(X^T * X * w * (dw)^T) + \operatorname{tr}(X^T * X * w^T * dw) \quad (43)$$

$$= \operatorname{tr}(w^T * X^T * X * dw) + \operatorname{tr}(X^T * X * w^T * dw) = \operatorname{tr}(2 * X^T * X * w^T * dw) \quad (44)$$

所以

$$\frac{d(w^T * X^T * X * w)}{dw} = 2 * w * X^T * X \quad (45)$$

所以 $\frac{dL(w)}{dw} = 2 * X^T * X * w - 2 * X^T * Y$

令导数等于0，得出最小二乘的闭式解：

$$\hat{w} = (X^T * X)^{-1} * X^T * Y \quad (46)$$

B.几何视角

$$X = (x_1, x_2, \dots, x_n)^T = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ x_{n1} & x_{n2} & \dots & x_{np} \end{pmatrix}_{n \times p} \quad (47)$$

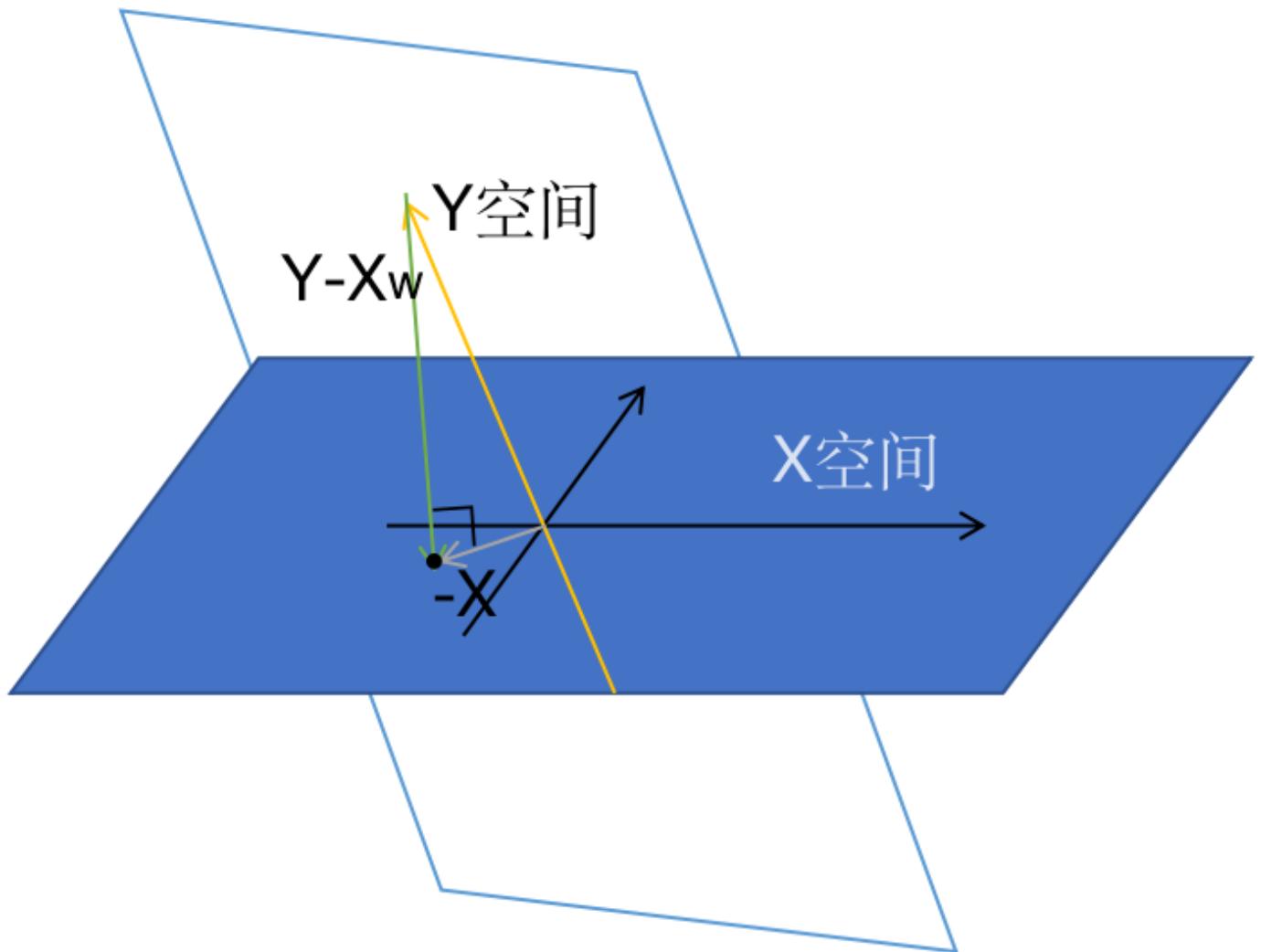
$$Y = (y_1, y_2, \dots, y_n)^T_{n \times 1} \quad (48)$$

在几何视角下，我们将 X 看作是一个 p 维的向量

X 的第一维是 $(x_{11}, x_{21}, \dots, x_{n1})$ ， X 的第 p 维是 $(x_{1p}, x_{2p}, \dots, x_{np})$

而这里的 Y 被看作是一个一维的向量

现在我们假设 $p = 2$ ，因为比较好画。示意图如下（俺真的画了好久，观众老爷们给波三连吧～～）



将模型改为 $f(w) = Xw$ ，意为对 X 向量施以 w 权重的放缩

而最小二乘的几何意义就是找到一个 w ，使得 $Y - Xw$ 这个向量到 X 空间的距离最小，那最小的情况当然就是与 X 空间垂直

所以我们有式子 $X^T * (Y - Xw) = 0$

从而求解 w :

$$X^T * X * w = X^T * Y \quad (49)$$

$$\hat{w} = (X^T * X)^{-1} * X^T * Y \quad (50)$$

可以看到求出的 w 与矩阵视角的结果相同。

C. 概率视角

首先明确，现实中是很难用一条直线去拟合分布的。真实的数据必然存在一定的随机性，也就是噪声。

因此我们假设噪声 $\epsilon \sim N(0, \sigma^2)$

所以 $y = f(w) + \epsilon = w^T * x + \epsilon$

所以 $y|x; w \sim N(w^T * x, \sigma^2)$

带入高斯分布的概率密度函数:

$$p(y|x; w) = \frac{1}{\sqrt{2\pi} * \sigma} * e^{-\frac{(y-w^T * x)^2}{2\sigma^2}} \quad (51)$$

下面使用 MLE (极大似然估计):

注: 所谓极大似然估计, 即通过大量的采样得到相对频率, 去逼近概率

我们设一个函数 $\zeta(w) = \log p(Y|X; w)$

因为 n 个数据之间是独立的, 因此可以将概率改为连乘的形式。

$$\zeta(w) = \log \prod_{i=1}^n p(y_i|x_i; w) = \sum_{i=1}^n \log p(y_i|x_i; w)$$

将高斯分布的概率密度函数带入式子:

$$\zeta(w) = \sum_{i=1}^n \left(\log \frac{1}{\sqrt{2\pi} * \sigma} - \frac{(y - w^T * x)^2}{2\sigma^2} \right)$$

因为前一项与 w 无关, 所以可以忽略

所以:

$$\hat{w} = \operatorname{argmax} \zeta(w) \quad (52)$$

$$= \operatorname{argmax} \sum_{i=1}^n - \frac{(y - w^T * x)^2}{2\sigma^2} \quad (53)$$

$$= \operatorname{argmin} \sum_{i=1}^n (y - w^T * x)^2 \quad (54)$$

而使用极大似然估计得到的结论正是最小二乘法的定义。

这也恰好说明, 最小二乘法隐藏着一个噪声为高斯分布的假设。

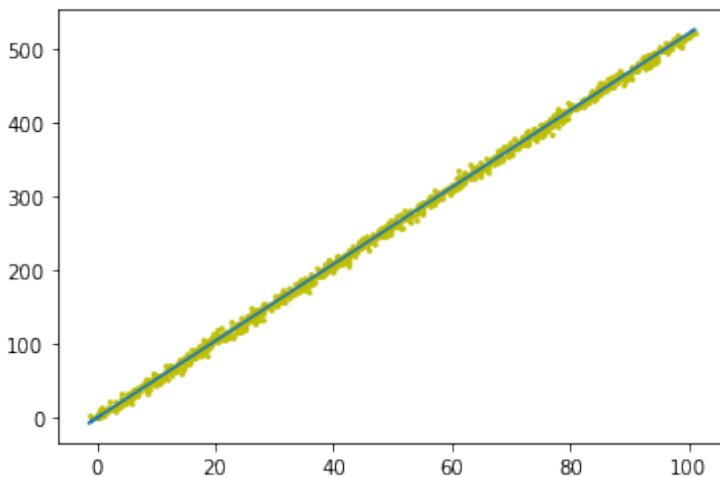
Implement (实作)

```
%matplotlib inline
import numpy as np
import matplotlib.pyplot as plt

# 样本数
n = 1000
```

```
# 噪声
epsilon = 1
X = np.expand_dims(np.linspace(0,100,1000), axis=-1)
w = np.asarray([5.2])
Y = X * w
# 增加噪声扰动
X += np.random.normal(scale=epsilon, size=(X.shape))
X_T = X.transpose()
w_hat = np.matmul(np.linalg.pinv((np.matmul(X_T, X))), np.matmul(X_T, Y))
print(w_hat)
plt.scatter(X, Y, s=3, c="y")
Y_hat = X * w_hat
plt.plot(X, Y_hat)
plt.show()
```

```
[[5.1992704]]
```



2.Algorithm（算法）

带正则化

A.矩阵视角

首先给出带正则项的新损失函数：

$$\zeta(w) = \sum_{i=1}^n ||y_i - w^T * x_i||^2 + \lambda * ||w||^2 \quad (55)$$

然后引用不带正则化的矩阵视角的损失函数的推导形式：

$$\zeta(w) = Y^T * Y - 2 * w^T * X^T * Y + w^T * X^T * X * w + \lambda * ||w||^2 \quad (56)$$

所以 $\hat{w} = \operatorname{argmax}(\zeta(w))$

对 $\zeta(w)$ 求导，得到：

$$\frac{\partial \zeta(w)}{\partial w} = 2 * X^T * X * w - 2 * X^T * Y + 2\lambda * w \quad (57)$$

令导数为0，得到带正则化的最小二乘法的闭式解：

$$\hat{w} = (X^T * X + \lambda * I)^{-1} * X^T * Y \quad (58)$$

I 为单位矩阵

B.概率视角

假设噪声 $\epsilon \sim N(0, \sigma_1^2)$

$w \sim N(0, \sigma_2^2)$

因为 $y = w^T * x + \epsilon$

所以 $y|w \sim N(w^T * x, \sigma_1^2)$

下面我们使用 MAP （最大后验估计）：

由贝叶斯定理得：

$$P(w|Y) = \frac{P(Y|w) * P(w)}{P(Y)} \quad (59)$$

其中 $P(w)$ 为先验概率， $P(Y|w)$ 为似然概率， $P(Y)$ 为归一化概率，先验概率乘似然概率并归一化得到后验概率 $P(w|Y)$

其中 $P(Y)$ 实际上为常数，因此：

$$\hat{w} = \operatorname{argmax}(P(w|Y)) = \operatorname{argmax}(P(Y|w) * P(w)) = \operatorname{argmax}(\log(P(Y|w) * P(w))) \quad (60)$$

因为每个样本间是独立的，因此可以将概率连乘

$$= \operatorname{argmax}(\log(\prod_{i=1}^n P(y_i|w) * P(w))) = \operatorname{argmax}(\sum_{i=1}^n \log(P(y_i|w) + \log(P(w)))) \quad (61)$$

带入高斯分布的概率密度函数，得到：

$$\hat{w} = \operatorname{argmax}(\sum_{i=1}^n \log(\frac{1}{\sqrt{2\pi} * \sigma_1}) - \frac{(y_i - w^T * x_i)^2}{2\sigma_1^2} + \log(\frac{1}{\sqrt{2\pi} * \sigma_2}) - \frac{w^2}{2\sigma_2^2}) \quad (62)$$

因为 σ_1, σ_2 都为超参数，因此可以省略

所以：

$$\hat{w} = \operatorname{argmin}(\sum_{i=1}^n \frac{(y_i - w^T * x_i)^2}{2\sigma_1^2} + \frac{w^2}{2\sigma_2^2}) \quad (63)$$

$$= \operatorname{argmin}(\sum_{i=1}^n (y_i - w^T * x_i)^2 + \frac{\sigma_1^2}{\sigma_2^2} * w^2) \quad (64)$$

可以看到，使用 MAP 推导出的结果正是带正则项的最小二乘的定义

Implement（实现）

带正则项

```
import os
os.chdir("../")
import numpy as np
import matplotlib.pyplot as plt
from models.linear_models import LinearRegression
```

```
X_ = np.expand_dims(np.linspace(0, 10, 1000), axis=-1)
X = np.c_[X_, np.ones(1000)]
w = np.asarray([5.2, 1])
Y = X.dot(w)
X = np.r_[X, np.asarray([[11, 1], [12, 1], [13, 1]])]
Y = np.r_[Y, np.asarray([100, 110, 120])]

model = LinearRegression(l2_ratio=1e1, epoch_num=1000, lr=1e-2, batch_size=100,
if_standard=False)
model.fit(X[:, :-1], Y)
print(model.get_params())
model.draw(X[:, :-1], Y)
```

```
([[5.213677612426203]], [0.9234839073187315])
```

