

机器学习中的各种范数与正则化

对于统计机器学习算法一般为缓解过拟合现象的发生需要在进行正则化操作,通过正则化以偏差的增加换取方差的减小。最常见的方法即为在损失函数中引入矩阵范数,以对模型的复杂程度做出惩罚,其目标函数一般如下式所示:

$$\min_{\theta} \tilde{J}(\theta; X, y) = \underbrace{J(\theta; X, y)}_{(1)} + \underbrace{\alpha \Omega(\theta)}_{(2)} \tag{1}$$

上式中,第一项即为经验风险,第二项即为正则化项。其中 $\alpha \geq 0$,为调整两者之间关系的系数。当 $\alpha = 0$ 时,则表示无正则化项, α 越大则表示对应正则化惩罚越大。

L^2 范数正则化

$L_2: \Omega(\theta) = \frac{1}{2} \|w\|^2$ 这里我们假设模型的偏置参数均为0,则参数 θ 即为 w ,因此其目标函数为: $\tilde{J}(w; X, y) = J(w; X, y) + \frac{\alpha}{2} w^T w$

对上式求其梯度有: $\nabla_w \tilde{J}(w; X, y) = \nabla_w J(w; X, y) + \alpha w$

使用梯度下降法更新权重 w ,则 w 将向负梯度方向收敛,如下: $w; \leftarrow w - \epsilon (\nabla_w J(w; X, y) + \alpha w) = (1 - \epsilon \alpha) w - \epsilon \nabla_w J(w; X, y)$

从中可以看出每次权值 w 更新时都将乘以 $(1 - \epsilon \alpha)$,该常数因子小于0,即权重将会逐渐收缩,趋近于0。

进一步地,这里令 $w^* = \arg \min_w J(w)$,将 $J(w; X, y)$ 简记为 $J(w)$,即 w^* 为损失函数 $J(w; X, y)$ 取得最小训练误差的权值。并在 w^* 的邻域对损失函数 $J(w)$ 做二次近似(泰勒展开),记为 $\hat{J}(w)$,如下:

$$\hat{J}(w) = J(w^*) + \frac{1}{2} (w - w^*)^T H (w - w^*) \tag{6}$$

上式中 H 为 J 在 w^* 处计算的Hessian矩阵,且该矩阵为半正定矩阵。由上述知, w^* 为损失函数的最优解,因此 $\hat{J}(w)$ 的梯度将为0,即式(6)对 w 求偏导为0,如下所示:

$$\nabla_w \hat{J}(w) = (w - w^*) H = 0 \tag{7}$$

记 \tilde{w} 为最优权值 w^* ,将式(7)代入式(4): $\nabla_w \tilde{J}(\tilde{w}) = \alpha \tilde{w} + (\tilde{w} - w^*) H = 0$

故:

$$(H + \alpha I) \tilde{w} = H w^* \quad \tilde{w} = (H + \alpha I)^{-1} H w^* \tag{9}$$

(由于Hessian矩阵为半正定矩阵,故其为实对称阵。因此有 $(\tilde{w} - w^*)^T H = H(\tilde{w} - w^*)$)

当 α 趋向于0时, \tilde{w} 将趋近于 w^* 。我们将实对称Hessian矩阵 H 分解为一个对角矩阵 Λ 和一组特征向量的标准正交基 Q ,因此有 $H = Q \Lambda Q^T$,代入式(9),可得

$$\begin{aligned} \tilde{w} &= (Q \wedge Q^T + \alpha I)^{-1} Q \wedge Q^T w = [Q(\wedge + \alpha I)Q^T]^{-1} Q \wedge \\ Q^T w &= Q^T \{I^{-1}\}(\wedge + \alpha I)^{-1} Q^T Q \wedge Q^T w = Q(\wedge + \alpha I)^{-1} \wedge \\ Q^T w &\end{aligned} \tag{10}$$

从上式中可以看出经过正则化后，权重 \tilde{w} 将会沿着由 H 特征向量所定义的轴缩放未经标准化的损失函数最优解 w^* 。具体来说，我们会根据 $\frac{\lambda_i}{\lambda_i + \alpha}$ 因子收缩与 H 第 i 个特征向量对齐的 w^* 的分量。如下图所示。

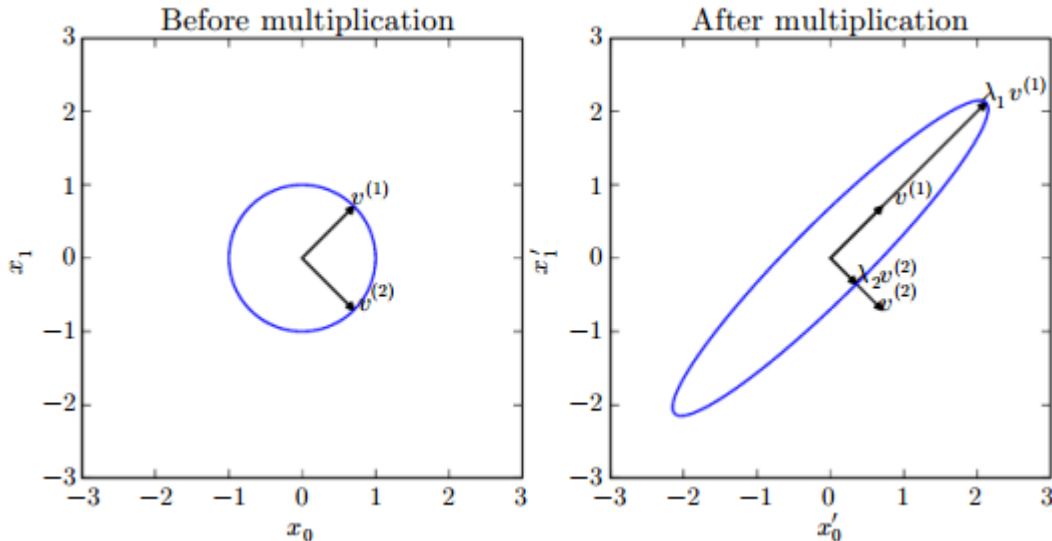


图1. 特征向量作用效果图

上图特征向量的作用效果图，这里矩阵有两个标准正交特征向量，对应的特征值分别为 $v^{(1)}$ 和 $v^{(2)}$ 。其中左图所有单位向量 $\mu \in \mathbb{R}^2$ 集合构成的单位圆。右图特征值的拉伸结果。

由于特征向量的缩放因子为 $\frac{\lambda_i}{\lambda_i + \alpha}$ ，故沿着特征值较大的方向（如 $\lambda_i \gg \alpha$ ）正则化影响较小。而 $\lambda_i \ll \alpha$ 将会收缩至几乎为零。因此 L^2 范数将使模型的参数趋近于0附近。

通过 L^2 正则化，在显著减小目标函数方向上的参数会被相对完整的保留，而对于无助于目标函数减小的方向（对应Hessian矩阵较小的特征值），其参数的改变将不会显著参加梯度，因此其在训练过程中会因正则化项而衰减至0。

此外，在线性回归的平方误差损失函数中引入二范数，即在原来逆矩阵的基础上加入对角阵，使得矩阵求逆可行，同时缓解过拟合的问题。而由于加入的对角矩阵，其就像一条“山岭”一样，因此， L^2 在统计学中也被称为岭回归或Tikhonov正则。

线性回归目标函数一般为：

$$L(w) = (Xw - y)^T (Xw - y) \tag{a} \quad \nabla_w L(w) = X^T (Xw - y) = 0 \tag{b}$$

固有：

$$X^T X w = X^T y \tag{c}$$

即：

$$w = (X^T X)^{-1} X^T y \tag{d}$$

当引入 L^2 正则项后，目标函数变为：

$$L(w) = (Xw - y)^T(Xw - y) + \frac{1}{2}\alpha w^T w \tag{e}$$

$$\nabla_w L(w) = X^T(Xw - y) + \alpha w = 0 \tag{f}$$

$$w = (X^T X + \alpha I)^{-1} X^T y \tag{g}$$

此即为岭回归。

上式中第（1）项 $X^T X$ 即为线性回归标准最小二乘项，第（2）项 αI 即为正则化项，为一对角阵。

另外从另一个角度理解，当 $X^T X$ 非正定时（某些参数线性相关），其无法求逆。此时解决的方法一般包括两种：

- （a）通过PCA构建线性无关组，进行降维处理，删除无关特征，求逆；
- （b）二范数正则化，即通过增加取值较小的对角阵使得矩阵求逆可行。由于PCA方法处理后其剔除了一些贡献程度较小的特征，而二范数只是将部分无关特征权值缩放置0附近，因此二范数也被称为Soft-PCA。 L^1 范数正则化

L^1 范数形式如下：

$$\Omega(\theta) = \|w\|_1 = \sum_i |w_i| \tag{11}$$

如上式所示， L^1 范数为各参数的绝对值之和。（ L^1 范数求导、优化困难，因此较 L^2 范数相比使用较少）对于， L^1 范数其目标函数如下所示：

$$\tilde{J}(w; X, y) = J(w; X, y) + \alpha \|w\|_1 \tag{12}$$

其对应的梯度如下：

$$\nabla_w \tilde{J}(w; X, y) = \nabla_w J(w; X, y) + \alpha \text{sign}(w) \tag{13}$$

上式中， $\text{sign}(w)$ 为符号函数，其取值结果只与个权值 w 的正负号有关。

同理，这里令 $w^* = \arg \min_w J(w)$ ，我们可以将 L^1 正则化目标函数的二次近似解分解为关于各参数求和的形式：

$$\hat{J}(w) = J(w^*) + \sum_i \left[\frac{1}{2} H_{\{i, i\}} (w_i - w_i^*)^2 + \alpha |w_i| \right] \tag{14}$$

对每一维 w_i 求梯度，以最小化式（14）。由于 w^* 为 $J(w)$ 的最优解，因此有 $\nabla_w J(w) = 0$ ，故：

$$\nabla_{w_i} \hat{J}(w) = [H_{\{i, i\}} (w_i - w_i^*) + \alpha \text{sign}(w_i^*)] = 0 \tag{15}$$

即：

$$w_i = w_i^* - \frac{\alpha}{H_{\{i, i\}}} \text{sign}(w_i^*)$$

对 w_i^* 的正负号分类讨论，则上式将等价于：

$$w_i = \text{sign}(w_i^*) \max\left\{w_i^* - \frac{\alpha}{H_{i,i}}, 0\right\} \tag{16}$$

当 $w_i > 0$ 时，会有两种结果：

- $w_i^* \leq \frac{\alpha}{H_{i,i}}$ 。此时式（16）的最优值为 $w_i = 0$ 。
- $w_i^* > \frac{\alpha}{H_{i,i}}$ 。此时则有 $w_i = w_i^* - \frac{\alpha}{H_{i,i}}$ ，即 w_i 在此方向上向0收缩 $\frac{\alpha}{H_{i,i}}$ 个距离。

同理，当 $w_i < 0$ 时， w_i 也将等于0或向0收缩。

与 L^2 范数相比， L^1 范数正则化使得部分参数为0。因此， L^1 范数会产生更稀疏的解，该性质已被大量应用于特征的选择机制。著名的LASSO回归即将 L^1 范数引入至线性模型中，并使用最小二乘代价函数。通过 L^1 正则化使得部分权值为零，而忽略相应的特征。如图2所示。

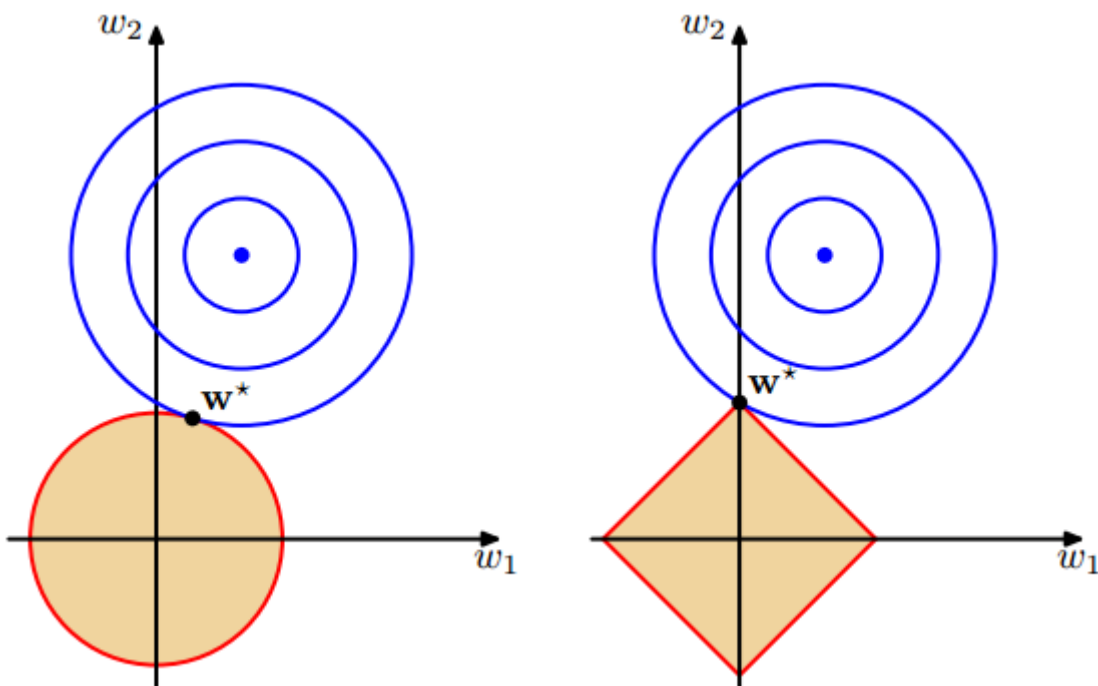


图2. 一范数与二范数示意图

上图中，蓝色的圆圈表示原问题可能的解范围，橘色的表示正则项可能的解范围。而整个目标函数（原问题（损失函数）+正则项）有解当且仅当两个解范围相切。从上图可以很容易地看出，由于 L^2 范数解范围是圆，所以相切的点有很大可能不在坐标轴上，而由于 L^1 范数是菱形，其相切的点更可能在坐标轴上。因此其只有一个坐标分量不为零，其它坐标分量为零，即 L^1 的解是稀疏的。

L^2 范数正则化，模型权重服从高斯分布， L^1 范数正则化，模型参数服从各向同性的拉普拉斯分布。即 L^1 正则化项 $\alpha \Omega(w) = \alpha \sum_i |w_i|$ 与通过MAP贝叶斯推断最大化对数先验项等价。此外，相较于 L^2 而言，采用 L^1 模型的鲁棒性较差（每次至零的权值均不相同）。

$$\begin{aligned} \log(p(w)) &= \sum_i \log \text{Laplace}(w_i; 0, \frac{1}{\alpha}) = \sum_i \log\left(\frac{1}{\alpha} \exp\left(-\frac{|w_i|}{\alpha}\right)\right) \\ &= \sum_i \log\left(\frac{1}{\alpha} \cdot \frac{1}{2} \exp(-\alpha |w_i|)\right) = \sum_i \log\left(\frac{1}{2\alpha}\right) + \sum_i \log(-\alpha |w_i|) \\ &= -n \log 2 - n \log \alpha - \sum_i \log |w_i| \end{aligned} \tag{17}$$

(拉普拉斯分布: $\text{Laplace}(x; \mu, \gamma) = \frac{1}{2\gamma} \exp\left(-\frac{|x - \mu|}{\gamma}\right)$)

由于目标函数是关于 w 的最大化进行学习，因此可以忽略 $n\log\alpha - n\log 2$ 。

L^0 范数

L^0 范数如下所示：

$$\|w\|_0 = \#\{i; w_i \neq 0\} \tag{18}$$

L^0 范数即为模型参数中不为0的数目。在实际问题中，模型一般较为复杂，参数数目较多，因此求解 L^0 范数为NP难问题，故一般不使用。在过去的前几年做压缩感知即稀疏表达时一般会使用 L^0 范数，但由于其优化困难因此会用 L^1 范替代。

Frobenius范数

Frobenius范数如下所示：

$$\|w\|_F = (\text{tr}(w^T w))^{\frac{1}{2}} = (\sum_i \sum_j |w_{ij}|^2)^{\frac{1}{2}} \tag{19}$$

从上式可以明显看出，矩阵的Frobenius范数就是将矩阵张成向量后的 L^2 范数。（在此就不做推导了）

各类范数一般作为各类机器学习模型权值的约束项（惩罚项）出现在目标函数中，训练过程以使得结构风险最小化，防止过拟合的发生。