

机器学习中的范数正则化

L1范数和L0范数可以实现稀疏，L1因具有比L0更好的优化求解特性而被广泛应用

L0范数

- 定义：向量中非0元素的个数
- 使用L0范数正则化参数矩阵W，即希望W的大部分元素都为0，即W是稀疏的
- L0范数很难优化求解（NP难问题）

L1范数

- 定义：向量中各个元素绝对值之和
- “稀疏正则化”（Lasso regularization）
- L1是L0范数的最优凸近似
- 参数稀疏的原因
 - 特征选择 (Feature Selection)
 - 可解释性 (Interpretability)

L2范数

- 定义：||W||2，向量中各元素的平方和再开方
- L2范数正则项||W||2最小，可以使得W的每个元素都很小，都接近于0，越小的参数说明模型越简单
- L2范数的好处
 - 学习理论角度
 - 优化计算角度
- L1范数与L2范数的区别
- 模型空间的限制

核范数

- 定义：核范数||W||*是指矩阵奇异值的和 (Nuclear Norm)
- 如果矩阵表达的是结构性信息，例如图像、用户-推荐信息等，那么这个矩阵各行之间存在一定的相关性，则矩阵是低秩的
- 低秩矩阵
 - 低秩矩阵包含大量冗余信息，这些冗余信息可以用于对缺失数据进行恢复，也可以对数据进行特征提取
 - 低秩矩阵是非凸的，需要寻找凸近似求解
- 约束Low-Rank（低秩）
 - 矩阵填充 (Matrix Completion)
 - 鲁棒PCA
 - 鲁棒主成分分析 (Robust PCA)
- 约束低秩的实例
 - 背景建模
 - 变换不变低秩纹理 (TILT)
- 正则化参数的选择

监督学习可视化为最小化目标函数

- 第一项表示Loss，是L(y,f(x;w)) 衡量我们的模型（分类或者回归）对第i个样本的预测值f(x;w)和真实的标签y之间的误差
- 第二项表示正则化参数，对参数w的正则化函数Ω(w)去约束模型
- 在所有可能选择的模型中，我们应该选择能够很好地解释已知数据并且十分简单的模型
- 正则化项对应于模型的先验概率
- 在经验风险上加一个正则化项(regularizer)或惩罚项(penalty term)

带参模型

- Loss = Square loss : 最小二乘
- Loss = Hinge loss : SVM
- Loss = Exp-Loss : Boosting
- Loss = Log-Loss : Logistic Regression

L2范数为什么是使||W||2最小？可以理解为：假设权重只有两项，并且权重之和为1，则一定是两项越接近||W||2越小，权重越接近，说明不同特征对输出的影响几乎同等重要，所以模型就更简单了