

COMP6245(2021/22): Foundations of Machine Learning Lab Five

Issue	19 November 2021
Deadline	26 November 2021 (10:00 AM)

Objective

- To implement a Radial Basis Functions model on a regression problem and compare its performance with linear regression.
- To use ten-fold cross validation to quote uncertainty in empirical results when comparing the performances of two machine learning approaches.

Radial Basis Functions (RBF)

The RBF model is given by

$$g(\mathbf{x}) = \sum_{j=1}^M \lambda_j \phi(\|\mathbf{x} - \mathbf{m}_j\| / \sigma).$$

The model has a nonlinear part (with \mathbf{m}_j and σ as parameters within a basis function $\phi(\cdot)$) and a linear part with parameters λ_j . In problems where we can make sensible choices of the nonlinear part, the learning problem reduces to a linear problem in the λ_j s.

Constructing an $N \times M$ *design matrix* U with terms $u_{ij} = \phi(\|\mathbf{x}_i - \mathbf{x}_j\| / \sigma)$, and the targets in an $N \times 1$ vector \mathbf{f} , the least squares solution to estimate λ s is similar to linear regression: $\mathbf{l} = (U^t U)^{-1} U^t \mathbf{f}$, where \mathbf{l} is an $M \times 1$ vector containing the unknown λ s.

1. Study the skeleton implementation of a Gaussian RBF model given in the Appendix.
2. Make the following improvements to the given implementation:
 - Normalize each feature of the input data to have a mean of 0 and standard deviation of 1.
 - The width parameter of the basis functions σ is set to be the distance between two randomly chosen points. Could this sometimes cause an error? Change it to be the average of several pairwise distances.
 - The locations of the M basis functions \mathbf{m}_j are set at random points in the input space. Cluster the data using K-means clustering (with $K = M$) and set the basis function locations to the cluster centres.
 - Split the data into training and test sets, estimate the model on the training set and note the test set performance.
3. Implement ten-fold cross validation, where you split the data into ten parts, train on nine tenths of the data and test on the held out tenth set, repeating the process ten times.
4. Explore *overtraining* by changing the number of basis functions used.
5. Display the distributions of test set results for the RBF and Linear Regression Models as boxplots side by side.

Report

Describe the work you have done as a short report. Upload a *pdf* file **no longer than two pages** using the ECS handin system: <http://handin.ecs.soton.ac.uk>. Please use \LaTeX to typeset your report. Please make sure your name and email are included.

Appendix: Skeleton Implementation of RBF

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

from sklearn import datasets
from sklearn.linear_model import LinearRegression

diabetes = datasets.load_diabetes()
X = diabetes.data
y = diabetes.target

def gaussian(x, u, sigma):
    return(np.exp(-0.5 * np.linalg.norm(x-u) / sigma))

N, p = X.shape
print(N, p)

# Space for design matrix
#
M = 200
U = np.zeros((N,M))

# Basis function locations at random
#
C = np.random.randn(M,p)

# Basis function range as distance between two random data
#
x1 = X[np.floor(np.random.rand()*N).astype(int),:]
x2 = X[np.floor(np.random.rand()*N).astype(int),:]
sigma = np.linalg.norm(x1-x2)

# Construct the design matrix
#
for i in range(N):
    for j in range(M):
        U[i,j] = gaussian(X[i,:], C[j,:], sigma)

# Pseudo inverse solution for linear part
#
l = np.linalg.inv(U.T @ U) @ U.T @ y

# Predicted values on training data
#
yh = U @ l
fig, ax = plt.subplots(figsize=(3,3))
ax.scatter(y, yh, c='m', s=3)
ax.grid(True)
ax.set_title("Training Set", fontsize=14)
ax.set_xlabel("True Target", fontsize=12)
ax.set_ylabel("Prediction", fontsize=12)
```
