

Data Description:

Let's start by examining the contents of the provided datasets. I'll first load the datasets and then provide a brief introduction to each.

Based on the initial inspection, here's a brief introduction to each dataset:

1. Copy of actual.csv (actual_df):

Columns:

`patient`: Represents an identifier for each patient.

`cancer`: Specifies the type of cancer for each patient.

Target Variable: `cancer` – This is the type of cancer we would want to predict.

Observation: This seems to be the ground truth data for the test set.

2. data_test.csv (data_test_df):

Columns:

`Gene Description`: Describes the gene or the probe.

`Gene Accession Number`: An identifier for the gene.

Multiple columns representing data for different patients. These columns alternate between the expression value and a "call" value (likely indicating the quality or reliability of the measurement).

Observation: This dataset appears to contain gene expression data for various patients. The columns named with numbers represent the patient identifier, matching the `patient` column in `actual_df`.

3. data_train.csv (data_train_df):

Columns: Similar structure to `data_test.csv`, with gene descriptions, accession numbers, and alternating columns for patient data and their associated "call" values.

Observation: This is the training set, containing gene expression data for another set of patients.

Support Vector Machines (SVM) are a class of supervised learning algorithms that can be used for classification or regression problems. In the context of cancer type prediction, SVMs are primarily used for classification.

In cancer type prediction, data often comes from genetic or molecular profiles of tumor samples. For instance, gene expression levels from microarrays or RNA sequencing could be used as features.

The core idea of SVM is to find a hyperplane that best separates the data points of different classes. In a two dimensional space, this hyperplane is a line. In higher dimensions, it becomes a plane or a hyperplane.

Real world data is often not linearly separable. SVM can handle such data by transforming the original feature space into a higher dimensional space using a mathematical function called a kernel. Popular kernels include linear, polynomial, radial basis function (RBF), and sigmoid.

- Cancer datasets, especially those derived from genomic, transcriptomic, or proteomic profiles, are high dimensional. Neural networks can handle vast amounts of high dimensional data effectively
- Pretrained neural network models on large datasets can be finetuned for specific cancer prediction tasks.

Grid search is a method used to perform hyperparameter tuning to determine the optimal values for a given model.

Define the Parameter Space:

1. Learning rate
2. Batch size
3. Number of layers
4. Number of neurons in each layer
5. Activation functions
6. Dropout rates
7. Regularization parameters

Cartesian Product of Parameters: The grid search algorithm will generate a grid of all possible hyperparameter combinations based on the specified ranges.

Train & Validate

Evaluation Metric: An evaluation metric is used to assess the model's performance for each hyperparameter combination.

Test Performance: Evaluate the model's performance on a separate test dataset to gauge its generalization capabilities.

Support Vector Machines (SVM):

Strengths:

1. Effective in High Dimensional Spaces: SVMs can handle high dimensional data, which is common in genomics and proteomics.
2. Maximal Margin Classifier: By focusing on the maximal margin, SVMs aim to achieve good generalization even if the training sample size is small.
3. Kernel Trick: SVMs can handle nonlinearly separable data by using the kernel trick, which implicitly maps data to a higher dimensional space.

Weaknesses:

1. Binary Classifier: SVMs are inherently binary classifiers. Multiclass classification requires methods like "onevsone" or "onevsall", which can be computationally intensive for many classes.
2. Computationally Intensive: For large datasets, SVMs can be slow and require a lot of memory.
3. Parameter Sensitivity: Performance can be highly sensitive to the choice of the kernel and regularization parameter. Proper choice often requires cross validation, which can be computationally intensive.

Neural Networks:

Strengths:

1. Automatic Feature Learning: Neural networks, especially deep networks, can automatically learn and extract important features from raw data.
2. Multiclass Classification: Neural networks can handle multiclass classification natively.
3. StateoftheArt Performance: For many tasks, deep learning models have achieved state of the art performance, surpassing traditional machine learning methods.

Weaknesses:

1. Require Large Datasets: To harness their full potential, neural networks often require large amounts of labeled data. In biomedical domains, acquiring such data can be challenging.
2. Computationally Intensive: Training deep neural networks requires significant computational resources, especially GPUs.
3. Risk of Overfitting: Without proper regularization or if given too many parameters, neural networks can easily overfit to training data.

Summary:

- Random Forest (RF) has an accuracy of 85.3%.
- Support Vector Machine with a linear kernel (SVC (linear)) achieved the highest accuracy of 97.05%.
- Support Vector Machine with an RBF kernel (SVC (rbf)) has the lowest accuracy at 61.8%.
- Support Vector Machine with a polynomial kernel (SVC (poly)) has an accuracy of 73.53%.
- Neural Network (NN) matches the performance of the linear SVC with an accuracy of 97.05%.