

Lecture 20:
Microarray Data Analysis: classification

Lubica Benuskova

<http://www.cs.otago.ac.nz/cosc348/>

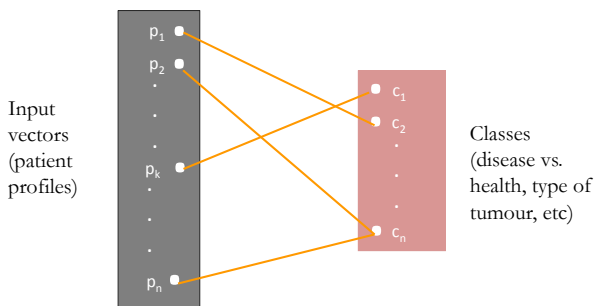
1

Steps of microarray analysis

- Gene Selection (after filtering and normalisation)
 - find genes, which would be the best predictors
- Clustering
 - K-means clustering
 - Hierarchical clustering
- Classification
 - K nearest neighbour
 - Support vector machine (SVM)
 - Artificial neural networks (ANN)
- Gene ontology:
 - record/compare/evaluate functions of genes

2

Classification: objects are assigned to classes



- Popular methods: K nearest neighbours, SVM and ANN.

3

Basic concepts for classification models

- **Training set** = a set of input vectors with known class labels
- **Test set** = a set of input vectors with *unknown* class labels
- **Stochastic training** = vectors are chosen randomly from the training set
- **Batch training** = all training vectors are presented to the model before the model parameters are updated
- **On-line protocol** = each pattern is presented once and model parameters are adjusted
- **Epoch** = a single presentation of all patterns in the training set.

4

Generalisation

- Generalisation:
 - to correctly classify data outside the training set, *e.g.* data in the test set.
 - *i.e.* the method must successfully classify (interpolate) to vectors it has not seen before.
- How can we test generalisation?
 - **Cross Validation (Leave one out)**: the class of each patient in the available data set is predicted while the rest of the data is regarded as the training set.

5

Classification results

- **Percentage of correctly classified objects,**
- **Confusion matrix** is a table with two rows and two columns that reports the number of True Negatives, **False Positives (type I error)**, **False Negatives (type II error)**, and True Positives.

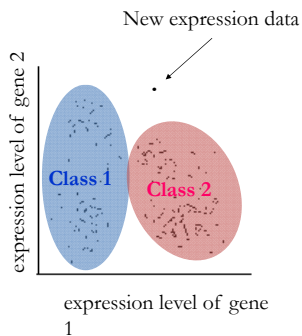
| | Predicted Negative | Predicted Positive |
|----------------|--------------------|--------------------|
| Negative Cases | True Negatives | False Positives |
| Positive Cases | False Negatives | True Positives |

| | Predicted Disease | Predicted Healthy |
|---------------|-------------------|-------------------|
| Disease Cases | True Disease | False Healthy |
| Healthy Cases | False Disease | True Healthy |

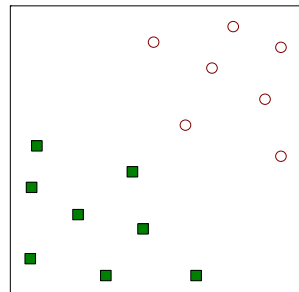
6

Classification based on K nearest neighbours

- Does the new expression data belong to class 1 or 2?
- K -NN algorithm computes the distance of a test sample to each of the training set samples, each of which has an associated class label, and then predicts the class of the test sample to be that of the majority of the K -closest samples.
- Optimal K is chosen based on cross-validation.



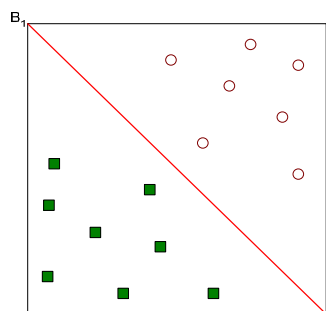
Linear classification



- Find a linear hyperplane (decision boundary) that will separate the data

8

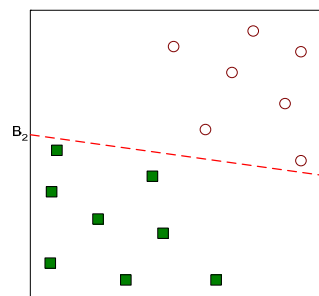
Linear classification



- One Possible Solution

9

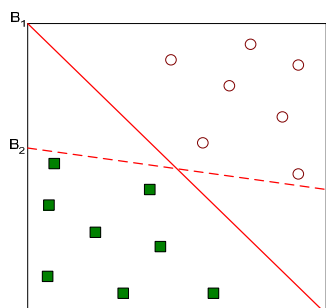
Linear classification



- Another possible solution

10

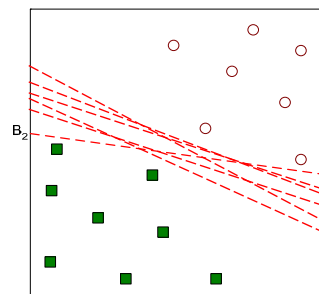
Linear classification



- Which one is better? B1 or B2?
- How do you define better?

11

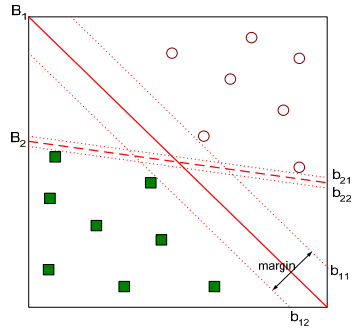
Linear classification



- There are many possible solutions – which one is the best?

12

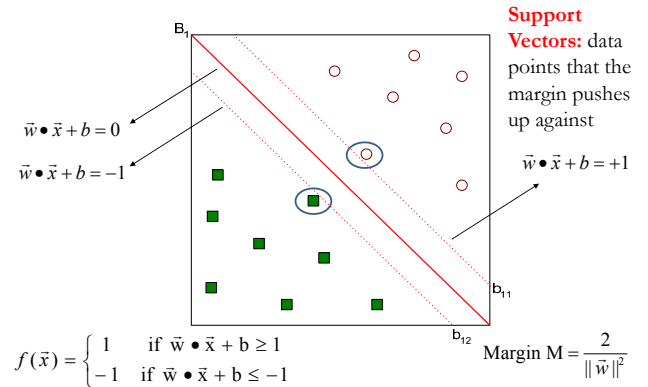
Linear Support Vector Machine (SVM)



- The hyperplane that maximizes the margin B1 is better than B2

13

SVM: find support vectors to maximize margin



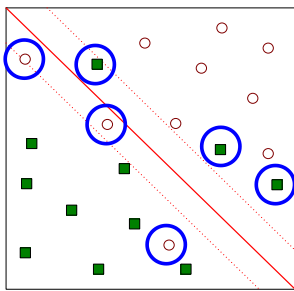
$$f(\bar{x}) = \begin{cases} 1 & \text{if } \bar{w} \cdot \bar{x} + b \geq 1 \\ -1 & \text{if } \bar{w} \cdot \bar{x} + b \leq -1 \end{cases}$$

$$\text{Margin } M = \frac{2}{\|\bar{w}\|^2}$$

14

Nonlinear problem

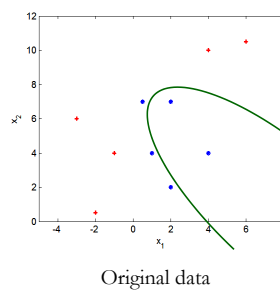
- What if the problem is not linearly separable?



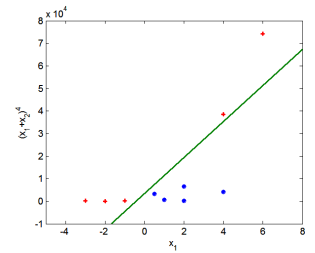
15

“Nonlinear” SVM

- Transform data into higher dimensional space
- Reduce the problem to a linear one



Original data



Transformed data

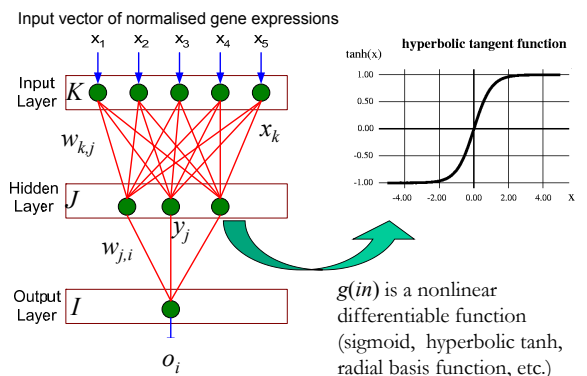
16

Properties of SVM

- The optimised margin function has a single minimum, this avoids the problem of local minima.
- The maximum margin requirements are appropriate conditions when the underlying statistical properties of the data are not known (which is often so).
- SVM works well when few observations are available (when the data space is very sparsely populated), which is often the case of medical problems.
- If there is a “high” nonlinearity in data, ANN may be more appropriate, as they are inherently nonlinear model.

17

MLP (Multilayer Perceptron)



$g(in)$ is a nonlinear differentiable function (sigmoid, hyperbolic tanh, radial basis function, etc.)

18

MLP representational power

- In principle any function can be approximated with arbitrarily small error by a two-layer feed-forward MLP.
- Sigmoid functions in the hidden layer act as a set of basis functions for composing more complex functions (like sine waves in Fourier analysis).
- The fit is adjusted thanks to a set of parameters of the MLP, which is the set of all weights.
- There are training algorithms which are able to find an optimal set of weights for each problem.

19

Training MLP based on input vectors

- Initialize the weights $\mathbf{w} = (w_0, w_1, \dots, w_k)$ as small random numbers
- For each input vector calculate the MLP output and adjust the weights in such a way that the output of MLP is consistent with class labels of training examples c_i

– Let the error function be:

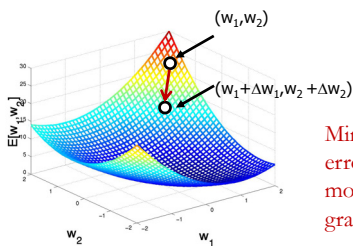
$$E = \sum_{i=1}^M (c_i - o_i)^2$$

– The task is to find the weights w_i 's that minimize the above error function:

- e.g, the error backpropagation algorithm:

20

Weights optimization by gradient descent



Minimization of the error function E – moving against the gradient of E

$$w_i \leftarrow w_i + \alpha \Delta w_i$$

$$0 < \alpha < 1$$

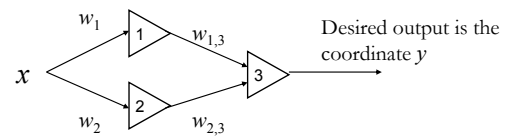
α is the learning speed

$$\Delta w_i = - \frac{\partial E}{\partial w_i}$$

21

Example of MLP

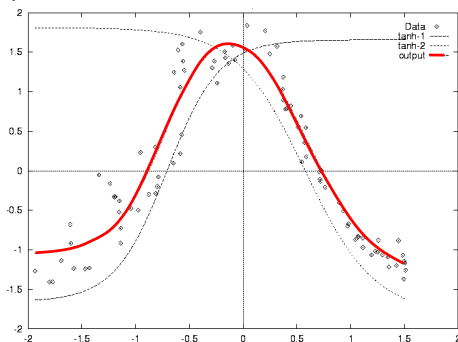
- Inputs are x coordinates of points in 2D space.
- 2 hidden units (1 and 2) have hyperbolic tangent activation function.
- One output unit (No. 3) sums linearly the outputs of 2 hidden units minus an adjustable bias and predicts the y coordinate.



22

Approximation of the nonlinear function $y = f(x)$

- Weights are adjusted based on input x data. MLP is then able to predict y values also for unseen x values.



23

MLP: advantages versus disadvantages

Advantages

- Guaranteed to converge to local minimum of error;
- Even when the training data contains high noise;
- Even when training data are nonlinearly separable – in fact MLP can approximate arbitrary nonlinear function.

Disadvantage

- Learning often gets stuck in a bad local minimum;
- No guarantee to converge to global minimum;
- Curse of dimensionality
 - For $N=2$ dimensional inputs, we need $100 = 10^2$ training data
 - For $N=3$ dimensional inputs, $10^3 = 1000$ training data are needed
 - In general, 10^N training data vectors are needed

24

NeuCom (c) Student v0.914

File Visualization Data Analysis Modeling Discovery Help

Statistical Methods
 Prediction
 Optimization
 Cross Validation

Evolutionary Methods
 Multi-Layer Perceptron (MLP)
 Evolving Connectionist Methods
 Radial Basis Function (RBF)

NeuCom

A Neuro-Computing Environment for Evolving Intelligence

Available Datasets List

Dataset_A1_46_complex54
 Dataset_C_4genesLog10_Normalised_134

Diagram components: New Inputs, Feature Selection, NNM, Higher Level Decision, Action Modules, Environment (Critique), Adaptation, Rule extraction.

Buttons: Save, View & Modify, Transpose, Rename, Extract, Split Plots, Delete, Delete All, Normalise, Join, Eigen Transforms, SPIR, 20 %

25

NeuCom - Multi-Layer Perceptron for Classification

Available Datasets: Dataset_C_4genesLog10_Normalised_2_141 #Samples: 60 #Variables: 5

Operation Parameters
 Mode: Train

Network Parameters
 Number of hidden units: 5
 Number of training cycles: 100
 Output Value Precision: 0
 Output Function Precision: 0
 Activation Function: linear
 Optimisation: jac

Visualisation Parameters
 Show: Confusion Table

Results
 Overall Accuracy: 76.667 %

Current Status
Training Finished

| | Actual Class 1 | Actual Class 2 |
|-------------------|----------------|----------------|
| Predicted Class 1 | 12 | 6 |
| Predicted Class 2 | 9 | 34 |

Buttons: Start, Reset

26

NeuCom - Multi-Layer Perceptron for Classification

Available Datasets: Dataset_C_4genesLog10_Normalised_134 #Samples: 60 #Variables: 5

Operation Parameters
 Mode: Train

Network Parameters
 Number of hidden units: 5
 Number of training cycles: 100
 Output Value Precision: 0
 Output Function Precision: 0
 Activation Function: linear
 Optimisation: jac

Visualisation Parameters
 Show: Accuracy of Each Class

Results
 Overall Accuracy: 76.667 %

Current Status
Training Finished

| Class | Accuracy % |
|-------|------------|
| 1 | ~70 |
| 2 | ~77 |

Buttons: Start, Reset

27