

MLCC Laboratory 3: Dimensionality reduction, feature selection

In this laboratory we will address the problem of data analysis and dimensionality reduction with a reference to a classification problem.

Think hard before you call the instructors or you look at the solution file!

1 Warm up - data generation

You will generate a training and a test set of D -dimensional points (N points for each class), with $N = 100$ and $D = 30$. Only two of those dimensions will be meaningful, the other one will be a variable we will modify.

- **1.A** For each point, the first two variables will be generated by `MixGauss`, extracted from two gaussian distributions with centroids $(1, 1)$ and $(-1, -1)$ and standard deviation 0.7 (the first one with labels 1 , the second with label -1):

```
1 [Xtr, Ytr] = MixGauss(...);
2 Ytr(Ytr==2) = -1;
3 [Xts, Yts] = MixGauss(...);
4 Yts(Yts==2) = -1;
```

- **1.B.** You may want to plot the relevant variables of the data:

```
1 scatter(Xtr(:,1), Xtr(:,2), 50, Ytr, 'filled');
2 hold on; scatter(Xts(:,1), Xts(:,2), 50, Yts);
```

- **1.C** The remaining variables will be generated as gaussian noise:

```
1 sigma_noise = 0.01;
2 Xtr_noise = sigma_noise*randn(2*N, D-2);
3 Xts_noise = sigma_noise*randn(2*N, D-2);
```

To compose the final data matrix, run:

```
1 Xtr = [Xtr, Xtr_noise];
2 Xts = [Xts, Xts_noise];
```

2 Principal Component Analysis (PCA)

- **2.A** Compute the data principal components (see "help PCA").
- **2.B** Plot the first two components of `X_proj` using the following line:

```
1 scatter(X_proj(:, 1), X_proj(:, 2), 50, Ytr, 'filled');
```

- **2.C** Try now with the first 3 components, by using:

```
1 scatter3(X_proj(:, 1), X_proj(:, 2), X_proj(:, 3), 50, Ytr, 'filled');
```

Reason on the meaning of the results you are obtaining.

- **2.D** Display the `sqrt` of the first 10 eigenvalues and plot the coefficients (eigenvector) associated with the largest eigenvalue:

```
1 disp(sqrt(d(1:10)));
2 scatter(1:D, abs(V(:,1)));
```

- **2.E** Repeat the above steps with dataset generated using different `sigma_noise = 0, 0.01, 0.1, 0.5, 0.7, 1, 1.2, 1.4, 1.6, 2`. To what extent data visualization by PCA is affected by the noise?

3 Variable selection

- **3.A** Use the data generated in part 1. Standardize the data matrix, so that each column has mean 0 and standard deviation 1:

```
1 m = mean(Xtr); %(see "help mean")
2 s = std(Xtr); %(see "help std")
3 for i = 1:2*N
4     Xtr(i,:) = Xtr(i,:) - m;
5     Xtr(i,:) = Xtr(i,:) ./ s;
6 end
```

Do the same for `Xts`, by using `m` and `s` computed on `Xtr`.

- **3.B** Use the orthogonal matching pursuit algorithm (type `"help OMatchingPursuit"`).
- **3.C** You may want to check the predicted labels on the training set:

```
1 Ypred = sign(Xts * w);
2 err = calcErr(Yts, Ypred);
```

and plot the coefficients `w` with `scatter(1:D, abs(w))`. How the error changes with the number of iterations of the method?

- **3.D** By using the method `holdoutCVOMP` find the best number of iterations with `intIter = 2, ..., D` (and, for instance, `perc = 0.75, nrip = 20`).

Moreover, plot the training and validation error with the following lines:

```
1 figure; plot(intIter, Tm, 'r');
2 hold on; plot(intIter, Vm, 'b');
3 xlabel('Number of iterations for OMP'); ylabel('error');
4 legend('Test', 'Validation');
```

What is the behavior of the training and the validation errors with respect to the number of iterations?

- **3.E** Try to increase the number of relevant variables $D = 3, 5, \dots$ (and the corresponding standard deviation of the Gaussians) around the centroids:

```
1 ones(D, 1); % vector of all 1s
2 % and
3 -ones(D, 1); % vector of all -1s
```

and see how this change is reflected in the cross-validation.

4 If you have time - more experiments

- **4.A** Analyse the results you obtain on sections 2 and 3 once you choose:

- $N \gg D$
- $N \approx D$
- $N \ll D$,

and evaluate the benefits of the two different analyses.

- **4.B** Dimensionality reduction is often used as a pre-processing step to a learning (classification) algorithm. The idea is to perform classification in a lower dimensional space and therefore save computational time. You have the following task:
 - Generate a new training and test datasets as in Laboratory 1.
 - Perform dimensionality reduction on the training set.
 - Using the projection you just found, project the test set.
 - Perform kNN in the lower dimensional (projected) space. Compare the result (both accuracy and running time) with the one in Laboratory 1.