

Accounting for Probe-level Noise in Principal Component Analysis of Microarray Data

Guido Sanguinetti, Marta Milo, Magnus Rattray and Neil D. Lawrence

<http://www.bioinf.man.ac.uk/resources/puma/>



Machine Learning Group
Computer Science
University of Sheffield, U.K.



Artificial Intelligence Group
Computer Science
University of Manchester, U.K.

Overview

- Principal Component Analysis (PCA) is one of the most popular dimensionality reduction techniques for the analysis of high dimensional datasets.
- In its standard form, PCA assumes all data points to be i.i.d. and corrupted by i.i.d. noise.
- This is often an unreasonable assumption when dealing with microarray data, where different genes and different conditions have different levels of experimental and biological noise.
- We propose a new model-based approach to PCA that takes into account the variances associated with each gene in each experiment [8].
- The model provides significantly better results than standard PCA, and avoids arbitrary manipulations such as setting cut-offs on expression levels.
- We demonstrate how the model can be used to denoise a data set, leading to improved expression profiles and tighter clusterings.

Motivation

- PCA is one of the most popular techniques for extracting information from high dimensional datasets.
- It is very popular in microarray data analysis, where the principal components are interpreted as the (few) physiological processes driving the variability in the data set.
- PCA makes two implicit assumptions: the first is that the data is normally distributed, the second is that the uncertainty associated with each gene under each condition is constant. This second assumption is often very unrealistic in biological data.
- Traditionally, this problem has been avoided by introducing cut-offs at the preprocessing stage. This involves a large degree of arbitrariness in selecting the cut-off and potentially throws away useful information about low expressed genes.
- Recent techniques allow to estimate credibility intervals for each gene expression in each time point in a microarray experiment ([4],[5],[3] and [6]).
- We seek to avoid ad hoc manipulations and propagate the uncertainties through a probabilistic model as a principled way of avoiding the problems inherent with PCA.

Probabilistic PCA

- Our model is a generalisation of the Probabilistic PCA algorithm ([9]).
- This is a latent variable model where each d -dimensional data point y_n can be reconstructed from a q -dimensional latent point x_n via a linear transformation W and a corrupting noise vector ϵ_n :

$$y_n = Wx_n + \mu + \epsilon_n$$
- The noise vector is assumed to come from a spherical Gaussian distribution $\epsilon_n \sim N(0, \sigma^2 I)$, which implies a likelihood

$$y_n | x_n \sim N(Wx_n + \mu, \sigma^2 I)$$
- Integrating over the latent variable x_n one obtains the marginal likelihood

$$y_n \sim N(\mu, WW^T + \sigma^2 I) \quad (1)$$
- Maximising the likelihood one obtains an estimate of the matrix W s.t. its columns span the principal subspace of the data space.

Limits of Probabilistic PCA and Factor Analysis

- Probabilistic PCA assumes the data to be i.i.d. and explains all noise effects with an indiscriminate spherical Gaussian.
- A more flexible model is Factor Analysis [1], which allows each dimension to have a different noise variance

$$\epsilon_n \sim N(0, B^{-1})$$
- where B is a diagonal matrix containing the individual precisions.
- Factor Analysis however is still not flexible enough to handle common microarray situations, where errors can vary greatly between genes and between different conditions for the same gene. We therefore need a model which can handle non-i.i.d. data.

Propagating Uncertainty

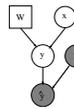
- We consider a modified model where the "true" expression levels are i.i.d. data but each measurement is corrupted by white noise with a different (known) variance. We assume that these variances have been measured at a preprocessing stage, using one of the various existing methods for obtaining credibility intervals from microarray experiments.
- Take y_n to be a d -dimensional vector which represents the true log expression level associated with the n -th gene under d different conditions. Rather than observing y_n directly we assume that we observe a corrupted form \tilde{y}_n

$$\tilde{y}_n = y_n + \nu_n \quad (2)$$

where ν_n is noise which is distributed as

$$\nu_n \sim N(0, B_n^{-1})$$

Here, B_n is a diagonal matrix whose i th diagonal element is given by β_{ni} which is the precision associated with the i -th experiment for the n -th gene. This precision can be obtained through one of the probabilistic analysis methods mentioned above. A graphical model representation of our model is given below.



Graphical representation of the noisy PCCA model.

- We will assume a probabilistic PCA model as the marginal distribution for the true expression level y_n , as given in (1), and obtain,

$$\tilde{y}_n | x_n \sim N(Wx_n + \mu, \sigma^2 I + B_n^{-1}) \quad (3)$$

We denote collectively $A_n = \sigma^2 I + B_n^{-1}$ and using $x_n \sim N(0, I)$, we have the following marginalised likelihood,

$$\tilde{y}_n \sim N(\mu, WW^T + A_n) \quad (4)$$

- The corrupted data is Gaussian distributed with mean μ and covariance $C_n \doteq WW^T + A_n$. Notice however that, as the data is not i.i.d., the maximum likelihood estimator of the mean vector μ does no longer coincide with the empirical mean, but must be learnt alongside the other parameters.

Efficient Likelihood Optimisation

- Given the marginal likelihood of equation (4), we can optimise the parameters through a non-linear optimisation such as scaled conjugate gradients.
- This is generally computationally unfeasible for large data sets. A more efficient algorithm can be obtained through an expectation maximisation (EM) approach [2].
- Generally EM algorithms lead to a simplified optimisation problem (the M-step) by incorporating an additional step (the E-step).
- For our corrupted data PCA model this additional step is the computation of the posterior distribution for the latent space. This posterior is obtained through Bayes' theorem

$$x_n | \tilde{y}_n \sim N(M_n W^T A_n^{-1} (\tilde{y}_n - \mu), M_n) \quad (5)$$
- where we have defined

$$M_n = [W^T A_n^{-1} W + I]^{-1}$$
- The EM algorithm then iteratively updates the posterior distribution (E-step) and maximises a lower bound on the likelihood with respect to the model parameters (M-step). The algorithm provably converges to a maximum of the likelihood.

E-step

- The lower bound on the likelihood is given by

$$\mathcal{L}_c = -\frac{1}{2} \sum_{n=1}^N \log |A_n| + \sum_{n=1}^N \text{tr} \left(\left(x_n x_n^T \right) \right) + \sum_{n=1}^N \text{tr} \left((y_n - \mu)^T A_n^{-1} (y_n - \mu) \right) - 2 \sum_{n=1}^N (x_n)^T W^T A_n^{-1} (y_n - \mu) + \sum_{n=1}^N \text{tr} \left(W^T A_n^{-1} W \left(x_n x_n^T \right) \right), \quad (6)$$
- where the notation $\langle \cdot \rangle$ denotes the expectation under the posterior distribution over x_n (5).
- The E-step evaluates these expectations from the sufficient statistics of the posterior distribution over x_n .

M-step

- Taking the gradients of eq. (6), one obtains 6 fixed point equations for W and μ which give an approximate M-step

$$\mu = \left(\sum_{n=1}^N A_n^{-1} \right)^{-1} \sum_{n=1}^N A_n^{-1} (y_n - W x_n),$$

$$W_j = \sum_{n=1}^N H_{nj} (L_n)^{-1}, \quad (7)$$

where A_{ij} is the j -th diagonal element of A_n , and we have introduced the two matrices

$$H = \sum_{n=1}^N A_n^{-1} (y_n - \mu) (x_n)^T,$$

$$L_j = \sum_{n=1}^N A_n^{-1} (x_n x_n^T).$$

- These update equations can be iterated to find the maximum likelihood solution with respect to W and μ .
- A fixed point equation for σ^2 (which accounts for any residual variance) cannot be obtained as the gradient with respect to σ^2 is not linear. An efficient update for σ^2 can be obtained by using Newton's method.

Number of Principal Components

- The usual approach when implementing PCA for microarray data is to retain a reduced number of principal directions, q , and project the log expression levels along these directions before further processing.
- In general, the number of principal components retained is pre-determined according to the specific problem under consideration.
- In our model, however, the estimate of the noise allows to evaluate the statistical significance of a direction.
- Therefore, we automatically obtain the maximum number of principal components that can be retained.

Data sets

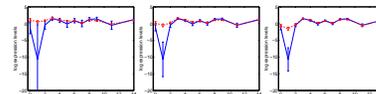
- Data set consisted of a temporal assay of Affymetrix GeneChip arrays that measured the gene expression profiles of a conditionally immortal cell line, UB/OC-1, from mouse cochlear epithelial cells at embryonic day 13.5 (E13.5), across 14 days of differentiation [7].
- Of particular interest in this study is the identification of targets regulated by the transcription factor gata-3, which is essential for normal inner ear development.
- In vivo the expression values of gata-3 are low before day 4 when they start to rise. They peak at day 8-9 and after a couple of days the expression level decreases again to then stabilize around a constant value.
- The raw data was processed using a modified version of the gMOS algorithm [6][5].

Profile Reconstruction

- The estimates of the parameters, together with the expectations of the hidden variables x_n , can be used in equation (3) to obtain estimates of the true gene expression levels and their covariances, given by

$$\tilde{y}_n = W(x_n + \mu) + \epsilon_n$$

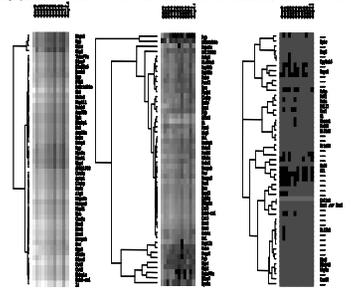
$$\Sigma_n = W \left[\langle x_n x_n^T \rangle - \langle x_n \rangle \langle x_n^T \rangle \right] W^T + \sigma^2 I$$
- We can then obtain an estimated profile for the "true" expression levels.
- To show the effect of the uncertainty in the data, we modelled the data three times, artificially reducing the variances by factors 1, 4 and 9. The corrected expression profiles are shown below.
- Note that as the uncertainty in the original profile is decreased the corrected profile tends to stay closer to its original course. As can be seen from the plots, any point with large associated uncertainty (such as the day 1 point for the gata-3 profile) can be significantly changed and this can lead to a large decrease in the associated uncertainty.



Corrected profile (thick dashed line) and original profile (thin solid line) for the gata-3 gene (a transcription factor) left: corrected profile based using the original uncertainties; middle: corrected profile with the uncertainty halved and right: corrected profile with a third of the original uncertainty.

Clustering

- Clustering is a widely used technique for summarising expression levels obtained from microarray data and as an exploratory technique for finding functional analogues.
- One suggested use of PCA in microarray analysis is as a preprocessing step before cluster analysis. The use of PCA before clustering can be justified by the fact that the larger principal components are expected to capture the structure in the data set.
- However, standard PCA does not always improve the clustering but often degrades it, since the dominant components, which contain most of the variation in the data, are highly influenced by the very noisy data points.
- By accounting for the variance in the log expression levels, our algorithm automatically downweights noisy values and ensures that the components we extract accurately reflect the structure of the data.
- The clustering is further improved when performed on the denoised reconstructed profiles, as these are the best estimates of the true profiles. This leads to much tighter and biologically plausible clusters in the data set under consideration, as shown below.



Hierarchical clustering of microarray data left: the top 50 genes in the second principal component obtained using our model (denoised profiles); middle: the top 50 genes in the second principal component obtained using our model (original profiles) and right: the top 50 genes in the second principal component obtained by standard PCA. Clustering was performed using the GeneCluster software from the Eisen Lab.

Acknowledgements

GS, MR and NL gratefully acknowledge support from a BBSRC award "Improved processing of microarray data with probabilistic models". MM is supported by an Advanced Tutorial Fellowship from the Wellcome Trust. The software used can be downloaded from <http://www.bioinf.man.ac.uk/~r08020208/puma/>

References

- [1] David J. Bartholomew. *Latent Variable Models and Factor Analysis*. Charles Griffin & Co. Ltd, London, 1987.
- [2] Arthur P. Dempster, Nan M. Laird, and Donald B. Rubin. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society, B*, 39(1):1-38, 1977.
- [3] Anne-Mette K. Hein, Sylvia Richardson, H. C. Causton, Graeme K. Ambler, and Peter J. Green. BGX: a fully Bayesian gene expression index for Affymetrix GeneChip data. *Bioinformatics*, 20(4):518-526, 2004.
- [4] Neil D. Lawrence, Marta Milo, Mahesan Niranjan, Penny Rabbas, and Stephan Soullier. Reducing the variability in cDNA microarray image processing by Bayesian inference. *Bioinformatics*, 20(4):518-526, 2004.
- [5] Xuejun Liu, Marta Milo, Neil D. Lawrence, and Magnus Rattray. A tractable probabilistic model for Affymetrix probe-level analysis across multiple chips. *Bioinformatics*, 21(18):3637-3644, 2005.
- [6] Marta Milo, Alireza Fazeli, Mahesan Niranjan, and Neil D. Lawrence. A probabilistic model for the extraction of expression levels from oligonucleotide arrays. *Biochemical Transactions*, 31(6):1510-1512, 2003.
- [7] Marcelo N. Rivolta, A. Halsall, C. Johnson, M. Tones, and Matthew C. Holley. Genetic profiling of functionally related groups of genes during conditional differentiation of a mammalian cochlear hair cell line. *Genome Research*, 12(7):1091-1099, 2002.
- [8] Guido Sanguinetti, Marta Milo, Magnus Rattray, and Neil D. Lawrence. Accounting for probe-level noise in principal component analysis of microarray data. *Bioinformatics*, 2005.
- [9] Michael E. Tipping and Christopher M. Bishop. Probabilistic principal component analysis. *Journal of the Royal Statistical Society, B*, 63(3):611-622, 1999.