

Probabilistic Principal Component Analysis and the E-M algorithm

The Minh Luong

CS 3750

October 23, 2007

Outline

- Probabilistic Principal Component Analysis
 - Latent variable models
 - Probabilistic PCA
 - Formulation of PCA model
 - Maximum likelihood estimation
 - Closed form solution
 - EM algorithm
 - » EM Algorithms for regular PCA
 - » Sensible PCA (E-M algorithm for probabilistic PCA)
 - Mixtures of Probabilistic Principal Component Analysers

Review of PCA

- Primary uses:
 - Analyze data and extract variables with similar concepts (principal components)
 - Project the data onto a lower dimensional space
 - Principal components which explain a greater amount of the variance are considered to be more important
- Accomplishes this by:
 - Maximizing variance of the projected data \mathbf{x}
 - Represent matrix \mathbf{x} in a different (q -dimensional) space using a set of orthonormal vectors \mathbf{W}
 - Weight matrix \mathbf{W} is a $d \times q$ matrix that represents a re-mapping of original data \mathbf{y} into its “ideal” principal subspace, represented by \mathbf{x}
 - Each of q orthonormal columns of the weight matrix \mathbf{W} , \mathbf{w}_i , represents a separate principal component
 - Likelihood of a point in \mathbf{y} is the distance² between it and its reconstruction, $\mathbf{W}\mathbf{x}$

Limitations of PCA

- Non-parametric
 - no probabilistic model for observed data
- The variance-covariance matrix needs to be calculated
 - Can be very computation-intensive for large datasets with a high # of dimensions
- Does not deal properly with missing data
 - Incomplete data must either be discarded or imputed using ad-hoc methods
- Outlying data observations can unduly affect the analysis

Motivation behind probabilistic PCA

- Addresses limitations of regular PCA
- PCA can be used as a general Gaussian density model in addition to reducing dimensions
- Maximum-likelihood estimates can be computed for elements associated with principal components
- Captures dominant correlations with few parameters
- Multiple PCA models can be combined as a probabilistic mixture
- Can be used as a base for Bayesian PCA

Latent variable models

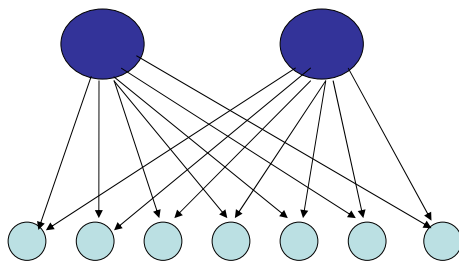
- Latent variable(s): unobserved variable (s)
 - Offer a lower dimensional representation of the data and their dependencies
- Latent variable model:
 - \mathbf{y} : observed variables (d -dimensions)
 - \mathbf{x} : latent variables (q -dimensions)
 - $q < d$
- Less dimensions results in more parsimonious models

Probabilistic PCA (PPCA)

- Latent variable model with linear relationship (factor analysis)
 - $y \sim Wx + \mu + \varepsilon$
 - Latent variables: $x \sim N(0, I)$
 - Error (or noise): $\varepsilon \sim N(0, \Psi)$
 - Location term (mean): μ
- Probabilistic PCA: Noise variances constrained to be equal ($\Psi_i = \sigma^2$)
 - Error: $\varepsilon \sim N(0, \sigma^2 I)$ (isotropic noise model)
 - $y|x \sim N(Wx + \mu, \sigma^2 I)$
 - $y \sim N(\mu, C_y)$, where $C_y = WW^T + \sigma^2 I$ (where C_y is the covariance matrix for the observed data y)
- Normal PCA is a limiting case of probabilistic PCA, taken as the limit as the covariance of the noise becomes infinitesimally small ($\Psi = \lim_{\sigma^2 \rightarrow 0} \sigma^2 I$)

Illustration of probabilistic PCA

Latent variables (x) $q = 2$
(hidden variables, underlying concepts)



Observed variables (y) $d = 7$
(data)

$$x \sim N(0, I)$$

Remapping: Wx
(Weight matrix: W)

+

μ (location parameter)

+

Random error (noise): ε
 $\varepsilon \sim N(0, \sigma^2 I)$

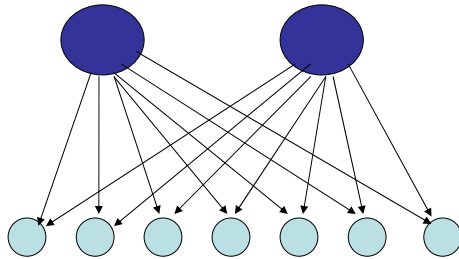
$$y = Wx + \mu + \varepsilon$$

$$y \sim N(\mu, WW^T + \sigma^2 I)$$

Parameters of interest: W (weight matrix), σ^2 (variance of noise)

Illustration of probabilistic PCA

Latent variables (\mathbf{x}) $q = 2$
(hidden variables, underlying concepts)



Note: Observed variables become independent of each other given latent factors

Observed variables (\mathbf{y}) $d = 7$
(data)

PPCA (Maximum likelihood PCA)

- Log-likelihood for Gaussian noise model:
 - $LL = -N/2 \{d \ln(2\pi) + \ln|\mathbf{C}_y| + \text{tr}(\mathbf{C}_y^{-1}\mathbf{S})\}$
 - $\mathbf{C}_y = \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}$
- Maximum likelihood estimates for above:
 - $\boldsymbol{\mu}$: mean of the data
 - \mathbf{S} (sample covariance matrix of the observations \mathbf{Y}):
 - $\mathbf{S} = (1/N) \sum_{n=1}^N (\mathbf{Y}_n - \boldsymbol{\mu})(\mathbf{Y}_n - \boldsymbol{\mu})^T$
- MLE's for \mathbf{W} and σ^2 can be solved in two ways:
 - closed form (Tipping and Bishop)
 - EM algorithm (Roweis)

$\text{Tr}(\mathbf{A})$ = sum of diagonal elements of \mathbf{A}

MLE's for probabilistic PCA (closed form)

- Likelihood of LL is maximized with respect to W and σ^2 , MLE's can be obtained in closed form:

$$- \sigma_{ML}^2 = \frac{1}{d-q} \sum_{j=q+1}^d \lambda_j$$

- Represents the variance lost in the projection, averaged over the # dim decreased

$$- W_{ML} = U_q (\Lambda_q - \sigma^2 I)^{1/2} R$$

- Represents the mapping of the latent space (containing X) to that of the principal subspace (containing Y)
- Columns of U_q ($d \times q$ matrix): principal eigenvectors of S
- Λ_q ($q \times q$ diagonal matrix): corresponding eigenvalues $\lambda_{1..q}$
- R : $q \times q$ arbitrary rotation matrix (can be set to $R=I$)

Derivation of MLEs

$$- LL = -N/2 \{d \ln(2\pi) + \ln|C_y| + \text{tr}(C^{-1}_y S)\}$$

The 1st derivative of LL w/ respect to W :

$$- dL/dW = N(C^{-1}SC^{-1}W - C^{-1}W), \text{ where } W = ULV^T = \sigma^2 I + WW^T$$

$$- \text{The stationary points are } SC^{-1}W = W.$$

$$- \text{Non-trivial case: } W \neq 0, C \neq S$$

$$- \text{SVD: } W = ULV^T, U: d \times q \text{ orthonormal vectors, } L: q \times q \text{ matrix of singular values, } V: q \times q \text{ orthogonal matrix,}$$

$$\bullet C^{-1}W = W(\sigma^2 I + W^T W)^{-1} = UL(\sigma^2 I + L^2)^{-1}V^T$$

$$- \text{At the stationary points:}$$

$$\bullet SUL(\sigma^2 I + L^2)V^T = ULV^T$$

$$\bullet SUL = U(\sigma^2 I + L^2)L$$

$$- \text{Column vectors of } U, u_j, \text{ are eigenvectors of } S, \text{ with eigenvalue } \lambda_j, \text{ such that } \sigma^2 + l_j^2 = \lambda_j$$

$$\bullet l_j^2 = (\lambda_j - \sigma^2)^{1/2}$$

$$- (\text{substitute into SVD}) W = U_q (\Lambda_q - \sigma^2 I) R$$

$$\bullet U_q: d \times q \text{ with } q \text{ column eigenvectors } u_j \text{ of } S$$

$$\bullet \Lambda_q: \lambda_1 \dots \lambda_q, (q \text{ eigenvalues of } u_j), \text{ or } \sigma^2 \text{ (corresponding } d-q \text{ "discarded" rows of } W)$$

$$\bullet R: \text{arbitrary orthogonal matrix, equivalent to a rotation in principal subspace (or a re-parametrization)}$$

Derivation of MLEs (cont)

- Substitute above results into original LL expression
- $LL = -N/2 \{ d \ln(2\pi) + \sum_{j=1}^q \ln(\lambda_j) + \sum_{j=q+1}^d \lambda_j + (d - q) \ln \sigma^2 + q \}$
 - $\lambda_1 \dots \lambda_q$ are q non-zero eigenvalues of \mathbf{u}_j and $\lambda_{q+1} \dots \lambda_d$ are zero
- Taking derivative of above with respect to σ^2 and solving for zero gives:

$$\sigma_{ML}^2 = \frac{1}{d - q} \sum_{j=q+1}^d \lambda_j$$

Differences between factor analysis and probabilistic PCA (PPCA)

- Covariance
 - PPCA (and standard PCA) is covariant under rotation of the original data axes
 - Factor analysis is covariant under component-wise rescaling
- Principal components (or factors)
 - In PPCA: different principal components (axes) can be found incrementally
 - Factor analysis: factors from a two-factor model may not correspond to those from a one-factor model

Dimensionality reduction and optimal reconstruction

- Using Bayes rule, we can obtain a posterior estimate of the latent variables
 - $x|y \sim N(M^{-1}W^T(y - \mu), \sigma^2 M^{-1})$,
 - where $M = W^T W + \sigma^2 I$, M is a $q \times q$ matrix
 - **Cond. latent mean:** $E[x|y] = \langle x_n | y_n \rangle = M^{-1}W^T(y_n - \mu)$
- Reconstruction of the observed data with respect to the new subspace:
 - The latent projection of regular PCA is skewed towards the origin (due to marginal distribution for x)
 - $y_n = W_{ML} \langle x_n | y_n \rangle + \mu$ is not orthogonal and thus not optimal
 - Optimal reconstruction of the observed data may still be obtained from conditional latent mean:
 - $y_n = W_{ML}(W_{ML}^T W_{ML})^{-1} M \langle x_n | y_n \rangle + \mu$

Motivation behind using E-M for PCA

- Naive PCA and MLE PCA computation-heavy for high dimensional data or large data sets
- PCA does not deal properly with missing data
 - E-M algorithm estimates ML values of missing data at each iteration
- Naïve PCA uses simplistic way (distance² from observed data) to access covariance
 - Sensible PCA (SPCA) defines a proper covariance structure whose parameters can be estimated through the E-M algorithm

E-M algorithm (review)

- Iterative process to estimate parameters consisting of two steps for each iteration
 - Expectation (data step): complete all hidden and missing variables Θ (or latent variables) from current set of parameters
 - Maximization (likelihood step): Update set of parameters Θ' , using MLE, from complete set of data from previous step
- Likelihood obtained from MLEs guaranteed to improve in successive iterations
- Continue iterations until negligible improvement is found in likelihood

E-M algorithm for normal PCA

- Amounts to an iterative procedure for finding subspace spanned by the q leading eigenvectors without computing covariance
- E-step: $X = (W^T W)^{-1} W^T Y$
 - Fix subspace and project data, y , into it to give values of hidden states x
 - Known: Y : d -dimensional observed data
 - Unknown (latent): X : q -dimensional unknown states
- M-step: $W_{\text{new}} = YX^T(XX^T)^{-1}$
 - Fix values of hidden states and choose subspace orientation that minimizes squared reconstruction errors

E-M algorithm and missing data

- Data with missing obs filled out: \mathbf{x} , Complete data (with blanks not filled out): \mathbf{y}

E-step (fill in missing variables):

- If data point \mathbf{y} is complete, then $\mathbf{y}^* = \mathbf{y}$ and \mathbf{x}^* is found as usual
- If the data point \mathbf{y} is not complete, \mathbf{x}^* and \mathbf{y}^* are the solution to the least squares problem. Compute \mathbf{x} by projecting the observed data \mathbf{y} into the current subspace.
 - For each (possibly incomplete) point \mathbf{y} , find the unique pair of points $(\mathbf{x}^*, \mathbf{y}^*)$ that minimize the norm $\|\mathbf{W}\mathbf{x}^* - \mathbf{y}^*\|$.
 - Constrain \mathbf{x}^* to be in the current principal subspace and \mathbf{y}^* in the subspace defined by known info about \mathbf{y}
 - If \mathbf{y} can be completely solved in system of equations, set corresponding column of \mathbf{X} to \mathbf{x}^* and the corresponding column of \mathbf{Y} to \mathbf{y}^*
 - Otherwise, QR factorization can be used on a particular constraint matrix to find least squares solution

E-M algorithm and missing data (E-step)

$$\mathbf{W} = \begin{pmatrix} 1 & 1 \\ 1 & 0.5 \\ 2 & 1 \end{pmatrix} \quad \mathbf{X} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \quad \mathbf{Y} = \begin{pmatrix} 3 \\ 1 \\ ? \end{pmatrix}$$

$$\begin{array}{l} \mathbf{W}\mathbf{x} = \mathbf{y} \\ x_1 + x_2 = 3 \\ x_1 + 0.5x_2 = 1 \\ 2x_1 + x_2 = y \end{array} \xrightarrow{\text{solve}} \mathbf{X}^* = \begin{pmatrix} -1 \\ 4 \end{pmatrix} \quad \mathbf{Y}^* = \begin{pmatrix} 3 \\ 1 \\ 2 \end{pmatrix}$$

Set $\mathbf{x} = (-1, 4)$, $\mathbf{y} = (3, 1, 2)$, proceed to M-step

If two elements are missing in \mathbf{Y} , then we use QR factorization to find the pair $(\mathbf{x}^*, \mathbf{y}^*)$ with the least squares of the norm $\|\mathbf{W}\mathbf{x}^* - \mathbf{y}^*\|$, according to the constraints specified in the set of equations $\mathbf{W}\mathbf{x} = \mathbf{y}$.

EM for probabilistic PCA (Sensible Principal Component Analysis)

- Probabilistic PCA model:
 - $\mathbf{Y} \sim N(\boldsymbol{\mu}, \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I})$
- Similar to normal PCA model, the differences are:
 - We do not take the limit as σ^2 approaches 0
 - During E-M iterations, data can be directly generated from the SPCA model, and the likelihood estimated from the test data set
 - Likelihood much lower for data far away from the training set, even if they are near the principal subspace
- EM algorithm steps implemented as follows:
 - E: $\boldsymbol{\beta} = \mathbf{W}^T(\mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I})^{-1} \langle \mathbf{x}_n | \mathbf{y}_n \rangle = \boldsymbol{\beta}(\mathbf{Y} - \boldsymbol{\mu})$, $\boldsymbol{\Sigma}_x = n\mathbf{I} - n\boldsymbol{\beta}\mathbf{W} + \langle \mathbf{x}_n | \mathbf{y}_n \rangle \langle \mathbf{x}_n | \mathbf{y}_n \rangle^T$
 - Log-likelihood in terms of weight matrix \mathbf{W} , and a *centered* observed data matrix $\mathbf{Y} - \boldsymbol{\mu}$, noise covariance $\sigma^2\mathbf{I}$, and conditional latent mean $\langle \mathbf{x}_n | \mathbf{y}_n \rangle$
 - M: $\mathbf{W}^{new} = (\mathbf{Y} - \boldsymbol{\mu}) \langle \mathbf{x}_n | \mathbf{y}_n \rangle^T \boldsymbol{\Sigma}_x^{-1}$, $\sigma^{2\ new} = \text{trace}[\mathbf{X}\mathbf{X}^T - \mathbf{W} \langle \mathbf{x}_n | \mathbf{y}_n \rangle (\mathbf{Y} - \boldsymbol{\mu})^T] / n^2$
 - Differentiate LL in terms of \mathbf{W} and σ^2 and set to zero.

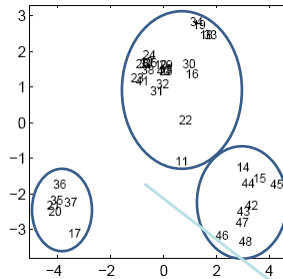
Advantages of using E-M algorithm in probabilistic PCA models

- Convergence:
 - Tipping and Bishop showed (1997) that the only stable local extremum is the *global maximum* at which the true principal subspace is found
- Complexity:
 - Methods that explicitly compute the sample covariance matrix have complexities $O(nd^2)$
 - E-M algorithm does not require computation of sample covariance matrix, $O(dnq)$
 - Huge advantage when $q \ll d$ (# of principal components is much smaller than original # of variables)

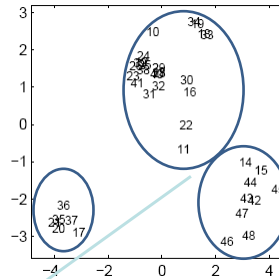
E-M algorithm for PPCA (illustration)

Example: 38 observations (with 18 data points each) from *Tobamovirus* data set (Ripley, 1996)

Standard PCA (on complete data)



Probabilistic PCA (using EM algorithm) with 20% (136) missing values



3 clusters

Other methods for PCA

- Power iteration methods
 - Iteratively update eigenvector estimates through repeated multiplication by matrix to be diagonalized
 - Extremely inefficient to calculate explicitly ($O(nq^2)$)
 - E-M algorithm provides efficient way to obtain sample covariance matrix, without explicitly calculating it
 - Iterative methods to compute SVD are closely related to the E-M algorithm
- Learning methods for the principal subspace
 - Sanger's and Oja's rule
 - Typically require more iterations and the learning parameter to be set by hand

Mixtures of probabilistic PCAs

- A combination of local probabilistic PCA models
- Multiple plots may reveal more complex data structures than a PCA projection alone
- Applications:
 - Image compression (Dony and Haykin 1995)
 - Visualization (Bishop and Tipping, 1998)
- Clustering mechanisms of mixture PPCA:
 - Local linear dimensionality reduction
 - Semi-parametric density estimation

Mixtures of probabilistic PCAs

- $LL = \sum_{n=1}^N \ln\{p(\mathbf{y}_n)\} = \sum_{n=1}^N \ln \left\{ \sum_{i=1}^M \pi_i p(\mathbf{y}_n|i) \right\}$
 - $p(\mathbf{y}|i)$ is a single PPCA model and π_i is the corresponding mixing proportion
 - Different mean vectors μ_i , weighting matrices \mathbf{W}_i , and noise error parameters σ_i^2 for each of M probabilistic PCA models
- An iterative E-M algorithm can be used to solve for parameters
- Guaranteed to find a *local* maximum of the log-likelihood