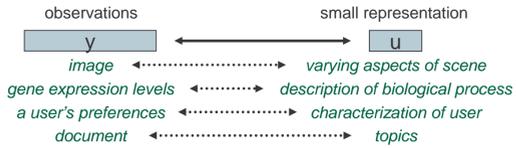


# Linear Dependent Dimensionality Reduction

Factor models are often natural in the analysis of multi-dimensional data. The underlying premise of such models is that **important aspects of the data can be captured via a low-dimensional representation.**



In many situations, including collaborative filtering and structure exploration, the "important" aspects are the dependencies between different attributes. Accordingly, we seek to identify a low-dimensional space that captures the **dependent** aspects of the data, and separate them from the **independent** variations. Our goal is to relax restrictions on the form of each these components, such as Gaussianity, additivity and linearity.

In this work we:

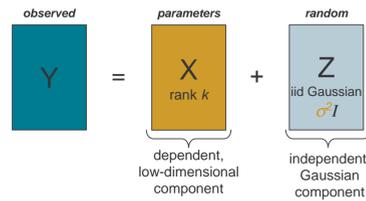
- Present a general framework for the problem: **Dependent Dimensionality Reduction**

Focusing on linear dependencies, we:

- Show that the **standard approach (PCA)** is consistent for **additive i.i.d. noise**, even if it is not Gaussian
- Show that a **variance-ignoring estimator** is appropriate for non-additive noise models
- Present a method for maximum likelihood estimation in the presence of **Gaussian mixture additive noise**
- Study the **consistency of maximum likelihood estimation** in this context, and show that the **ML estimator is not always consistent** (for example for Exponential-PCA).

## Dependent Dimensionality Reduction

Starting point: Identifying linear dependencies in the presence of i.i.d. Gaussian noise

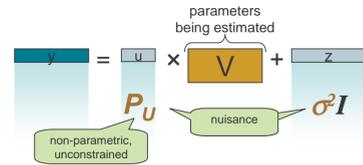


Log-likelihood(X)  $\propto$  sum-squared (Frobenius) distance to Y  
ML estimator is rank-k matrix minimizing  $\|X-Y\|_{Fro}$  given by leading components of SVD of Y

This formulation valid, but displeasing:

- Entire matrix X are parameters, estimated with a single observation Y
- Number of parameters linear in data
- Even with more data, cannot estimate beyond a fixed precision

What we *can* estimate with more data rows is the rank-k subspace of X:



Standard parametric analysis: imposes some (parametric) distribution  $P_u$ .

We do not make any assumptions about the distribution of  $u$ :

- Model class is non-parametric
- Can estimate a parametric aspect of the model – the subspace V.

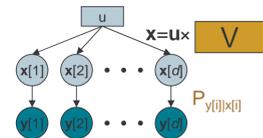
We do impose a strict form on the conditional distributions  $y_i|u$ .

Goal: relax this, and concentrate on the structural constraint – **u captures all the dependencies in y**

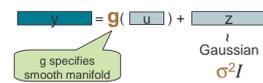
Relaxing Gaussianity leads to **Linear Additive** models:



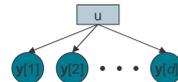
Relaxing additivity is appropriate, e.g., when the noise has a multiplicative component, or when the features in  $y$  are not real numbers, yielding **Linear Dependent Dimensionality Reduction**:



Going beyond linear models, fitting a non-linear manifold by minimizing sum-squared distance can be seen as a ML estimator for:



Combining these ideas leads us to discuss  $y_i|u$  directly:



**Dependent Dimensionality Reduction:**  
Low-dimensional representation  $u$  such that coordinates of  $y$  are **independent given u**

## Second Moment Methods

L2 approach to low-rank approximation: minimize sum-squared distance  $\|X-Y\|_{Fro}$ .

Subspace spanned by the leading  $k$  eigenvectors of empirical covariance of  $y$ .

For any i.i.d. additive noise:

$$\hat{\Sigma}_y \propto \Sigma_y = \Sigma_x + \sigma^2 I$$

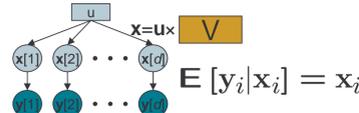
**L<sub>2</sub> estimation of the low-rank subspace (PCA) is consistent in the presence of any i.i.d. additive noise with finite variance**

Independent, non-identical additive noise:

$$\hat{\Sigma}_y \propto \Sigma_y = \Sigma_x + \Sigma_z$$

When the additive noise is independent, but not identically distributed, the L2 estimator is biased towards the high-variance coordinates. Instead, the **Variance Ignoring Estimator** seeks a rank- $k$  matrix approximating (minimizing the sum-squared distance to) the non-diagonal entries of the empirical covariance. This is a **Weighted Frobenius Low Rank Approximation (WLRA)** problem.

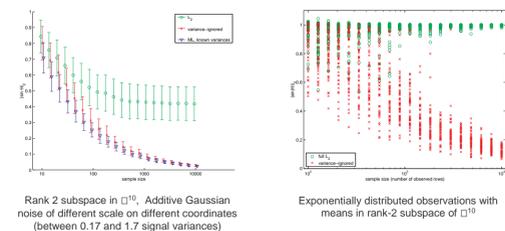
Unbiased Non-additive Noise:



$$\begin{aligned} \text{Cov}[y_i, y_j] &= E[y_i y_j] - E[y_i]E[y_j] \\ &= E[E[y_i y_j | x_i]] - E[E[y_i | x_i]]E[E[y_j | x_j]] \\ &= E[E[y_i | x_i]E[y_j | x_j]] - E[x_i]E[x_j] \\ &= E[x_i x_j] - E[x_i]E[x_j] \\ &= \text{Cov}[x_i, x_j] \end{aligned}$$

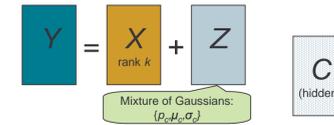
$$\hat{\Sigma}_y \propto \Sigma_y = \Sigma_x + E[\text{Var}[y_i | x_i]]$$

**Variance Ignoring Estimator is appropriate for any unbiased independent conditional model**



## Maximum Likelihood Estimation with Gaussian Mixture Noise

Model additive noise as a Gaussian mixture:



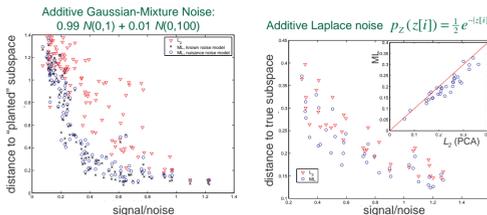
E step: calculate posteriors of C  
M step:

$$\begin{aligned} E_{\square|Y}[\log \Pr(\square = \square + \square | \square)] \\ &= -\sum_{i \in C} E_{\square_{i \in C} | Y_{i \in C}} \left[ \frac{1}{2} \log 2 \pi \sigma_{i \in C}^2 + \frac{(X_{i \in C} - Y_{i \in C}) - \square_{i \in C}}{2 \sigma_{i \in C}^2} \right]^2 \\ &= -\sum_{i \in C} \sum_{\square_{i \in C}} \frac{\Pr(\square_{i \in C} = \square_{i \in C})}{2 \sigma_{i \in C}^2} (\square_{i \in C} - (\square_{i \in C} + \square))^2 + \text{Const} \\ &= -\frac{1}{2} \sum_{i \in C} \square_{i \in C} (\square_{i \in C} - \square_{i \in C})^2 + \text{Const} \end{aligned}$$

Weighted Low-Rank Approximation with:

$$\square_{i \in C} = \Pr(\square_{i \in C} = \square_{i \in C}) \quad \square_{i \in C} = \square_{i \in C} + \Pr(\square_{i \in C} = \square_{i \in C}) \square_{i \in C}$$

and update mixture parameters.



Noise modeled a bounded **Gaussian Scale Mixture** (mixture of zero-mean Gaussians with variance bounded away from zero) Captures many distributions, including heavy tailed (non log-concave)

Instead of Gaussian mixture modeling: **Newton's method on log-Likelihood** (Gordon NIPS02 [SJ, ICML 03])  
Distribution must be log-concave – not applicable to heavy tailed distributions

## Weighted Low Rank Approximation

Given  $W, A$ , find rank- $k$   $X$  minimizing **weighted** sum-square distance:

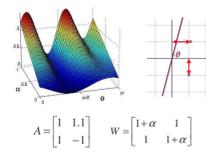
$$\sum_{ij} W_{ij} \left( A_{ij} - X_{ij} \right)^2$$

For fixed  $V$ , find optimal  $U$   
For fixed  $U$ , find optimal  $V$

$$\begin{aligned} J^*(V) &= \min_U J(UV) \\ \frac{\partial}{\partial V} J^*(V) &= 2U^*(U^*V - A) \otimes W \end{aligned}$$

$$X \leftarrow \text{LRA}(W \otimes A + (1 - W) \otimes X)$$

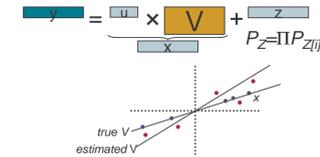
When weights introduced:  
• Not incremental with  $k$   
• Eigenmethods do not apply  
• Local minima:



## Consistency of Maximum Likelihood Estimation with a Known Noise model

General Conditions:

Assume a known additive noise model, and consider maximum likelihood estimation with respect to that model:



$$\Phi(V) = \max_u \log p_Z(y - uV)$$

ML estimator is consistent  $\iff E_{y=uV+z}(\Phi(V))$  is maximized on true V

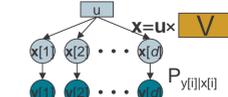
$$\Psi(V; x) = E_z \left[ \max_u \log p_Z((x+z) - uV) \right]$$

ML estimator is consistent for any  $P_u$   $\iff$  for all  $x$ ,  $V$  maximizes  $\Psi(V; x)$  iff  $V$  spans  $x$

ML estimator is consistent for any  $P_u$   $\implies \Psi(V; 0)$  is constant for all  $V$

**Maximum Likelihood Low-Rank estimation with non-Gaussian noise is not, in general, consistent**

Non-Additive Models:



These conditions can also be used to investigate the consistency of ML estimators with non-additive known conditionals  $y_i|x_i$ , where:

$$\Psi(V; x) = E_{y|x} \left[ \max_u \log p_{y|x}(y|uV) \right]$$

Of particular interest is "Exponential PCA", where the distribution  $y_i|x_i$  forms an exponential family with  $x_i$  the natural parameters [Collins Dasgupta Schapire, NIPS01].

**"Exponential-PCA" is not, in general, consistent**

When the *mean* parameters form a low-rank subspace, the variance-ignoring estimator is applicable, but when the *natural* parameters form a low-rank subspace no generally consistent estimator is known.

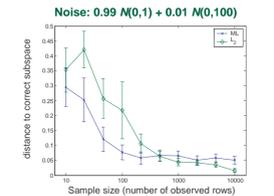
Challenge: Find a consistent estimator for the low-rank subspace of natural parameters

Gaussian Additive Noise:  
 $\Psi(V; x) = E[L_2 \text{ distance of } x+z \text{ from } V]$   
ML estimator is consistent

Laplace Additive Noise:

$$p_Z(z|V) = \frac{1}{2} e^{-|z|/V}$$

The ML estimator is **not** consistent!



Logistic Low-Rank Approximation:

$$P_{y_i|x_i}(+1|x) = \frac{1}{1 + e^{-x_i}}$$

The ML estimator is **not** consistent!