# Ch13. Linear Factor Models

Jihyeong Jung

Department of Mathematics, KAIST

# Contents

# Intro.
## - Linear Factor Models

- About Linear Factor Model
  - The simplest probabilistic. models with latent variables
  - Defined by the use of stochastic linear decoder function; $x = Wh + b + \epsilon$
- Data generation process
  - Sample $h$ from dist. $h \sim p(h)$ s.t. $p(h) = \prod_i p(h_i)$
  - Sample $x$ based on given factors; $x = Wh + b + \epsilon$
    - Noise $\epsilon$ is typically Gaussian & diagonal; indep. across dimensions.
  - Illustration :



$x = Wh + b + \text{noise}$

- Probabilistic PCA & FA, etc. are special cases of prev. slide
  - Difference : Choosing dist. of noise $\epsilon$ and prior of $h$ before observing $x$
- FA : prior of latent variable $h$ is just; Gaussian $h \sim \mathcal{N}(h; 0, I)$
  - Given $h$, observed $x_i$s are assumed to be conditionally indep.
  - Noise $\epsilon$ is assumed to be; $\epsilon \sim \mathcal{N}(\epsilon; 0, \Psi)$, $s.t. \Psi = diag(\sigma^2)$ & $\sigma^2 = [\sigma_i^2]^T$
- The role of latent variables : *'capture the dependencies'* b/w the different observed $x$, it's easy to see that $x \sim \mathcal{N}(x; b, WW^T + \Psi)$

# PPCA & FA

- PCA in a 'Probabilistic' framework; modify to FA model
  - Set $\sigma_i^2 = \sigma^2 \, for \, \forall i = 1, 2, \ldots, n; \, so \, that \, x \sim \mathcal{N}(x; b, WW^T + \sigma^2 I)$
  - Equivalently, $x = Wh + b + \sigma\epsilon \, s.t. \, \epsilon \sim \mathcal{N}(\epsilon; 0, I)$
  - By iterative EM algorithm, we can estimate $W$ & $\sigma^2$

- Based on the observation that most variations in the data can be captured by $h$, up to some small residual reconstruction error
  - As $\sigma \to 0$, then $x = Wh + b + \sigma\epsilon \overset{p}{\to} x = Wh + b$ (PPCA $\overset{p}{\to}$ PCA)

# ICA
- About Independent Component Analysis

- Forming observed data by separating observed signal into many underlying signals that are scaled & added together
  - An approach to modeling linear factors
  - Signals are intended to be not only uncorrelated from each other, but also fully independent.

- There exists many variances referred to as ICA

- Introduced variant : training fully parametric model

  - Prior of the underlying factors, $p(\boldsymbol{h})$, must be fixed in advance.
  - Then the model generates $\boldsymbol{x} = \boldsymbol{W}\boldsymbol{h}$, determine $p(\boldsymbol{x})$ by transformation
  - And learning model using maximum likelihood

# ICA
- Motivation of 'The' Variant/some examples

- By choosing $p(\boldsymbol{h})$ to be independent, we can recover underlying factors as close as possible to be independent
  - Commonly used to recover mixed low-level signals
  - Ex) commonly used in neuroscience, for the electro-encephalography
    - To separate signals of the heart from signals of the brain
    - And to separate signals in different regions of the brain from each other

# ICA
- Non-Gaussian prior requirement

- All variants have common requirement: non-Gaussian prior $p(\boldsymbol{h})$
  - We cannot identify $\boldsymbol{W}$ if $p(\boldsymbol{h}) \sim Gaussian\ Normal$
  - In other words, we can obtain same $p(\boldsymbol{x})$ for many $\boldsymbol{W}$
- Quite different from other LFMs such as PPCA, FA which require Gaussian prior $p(\boldsymbol{h})$ to make operations have closed solutions
- For 'the' variant, we typically use $p(h_i) = \frac{d}{dh_i}\sigma(h_i) = \frac{d}{dh_i}(\frac{1}{1+e^{-h_i}})$
  - Have larger peaks near 0 than Gaussian
- So most ICA implementations aim at learning sparse features

# ICA
- Another Variants

- There's a variant that adds some noise in the generation of $\boldsymbol{x}$ rather than deterministic decoder
    - It aims to make the elements of $\boldsymbol{h} = \boldsymbol{W}^{-1}\boldsymbol{x}$ independent from each other
- Some variants constrains $\boldsymbol{W}$ to be orthogonal to avoid possibly problematic operation determinant
- Some variants is not a 'generative model'
    - Many variants only know transformation b/w $\boldsymbol{x}$ & $\boldsymbol{h}$, but do not have any way of representing p($\boldsymbol{h}$) so it does not impose distribution over p($\boldsymbol{x}$)
    - Ex) many variants aim to increase sample kurtosis of $\boldsymbol{h} = \boldsymbol{W}^{-1}\boldsymbol{x}$
        - d

# ICA

- NICE(Non-linear Independent Components Estimation)

- ISA(Independent Subspace Analysis)

- Topographic ICA

# SFA

- LFM that uses information from time signals to learn invariant features

- Based on the Slowness principle

  - Idea : important characteristics of scenes change very slowly

  - Can be applied to any differentiable model trained with GD generally

  - By adding a term to cost function of the form $\lambda \sum_t L(f(x^{(t+1)}), f(x^{(t)}))$ to apply the Slowness principle

- Efficient application of the slowness principle; because it is applied to a linear feature extractor

# SFA

- Defining $f(x; \theta)$ as a linear transformation
    - And solve; $\min_{\theta} \mathbb{E}_t \left[ \left( f(x^{(t+1)})_i - f(x^{(t)})_i \right)^2 \right]$,
    - With constraints; $\mathbb{E}_t f(x^{(t)})_i = 0$ & $\mathbb{E}_t \left[ \left( f(x^{(t)})_i \right)^2 \right] = 1$
        - 1st constraint : to obtain a unique solution
        - 2nd constraint : to prevent solution where all features collapse to 0
- SFA features are ordered like PCA
    - The first feature is the slowest one

# SFA
## - SFA Algorithm/Generalizations/Advantages

- Additional constraint to learn multiple features;
  - $\mathbb{E}_t \left[ f(\boldsymbol{x}^{(t)})_i f(\boldsymbol{x}^{(t)})_j \right] = 0 \ \ for \ any \ i < j$
  - So the learned features must be linearly uncorrelated from each other
- Typically used to learn non-linear features by applying nonlinear expansion to input signals before running it
- Major advantage : possible to theoretically predict which features SFA will learn
  - Even in the deep nonlinear setting

# Sparse Coding
## - Training Sparse Coding model

- Encoder of the SC model is kind of optimization algorithm;
  - Solving $h^* = \text{argmax}_h\, p(h|x) = \text{argmin}_h\, \lambda||h||_1 + \beta||x - Wh||_2^2$
  - This yields a sparse $h^*$ because of the L1 norm imposition on $h$
- To train model, we alternate b/w minimization wrt $h$ and $W$
  - We treat $\beta$ as a hyper-parameter; we typically set it to 1 because it's role is shared with $\lambda$ in opt. problem
  - If we want to learn $\beta$, we must consider discarded terms

# Sparse Coding

- Sparse Coding +  non-parametric encoder
  - Can minimize the (reconstruction error + log-prior) better than any specific parametric encoder in principle
  - There's no generalization error; so applying sparse coding to feature extractor for a classifier makes better generalization than applying it as parametric func. to predict the code

# Sparse Coding

- non-parametric encoder

  - The primary disadvantage : large time cost of calculating $(h|x)$ because of iterative algorithm of non-parametric approach

  - Not easy to use back-prop. through the non-parametric encoder

- Often produces poor samples like other LFMs

  - Motivated the development of improved models