
PCA, Probabilistic PCA, Kernel PCA, ICA

Boyang Liu

Abstract

Summary of PCA, PPCA, KPCA, ICA and other related approaches. Mainly based on PRML chapter 12 and FOML. We also includes some disentangled representation objective in deep learning.

1 PCA

Denote our data matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$, we want to find the subspace $\mathbf{Z} \in \mathbb{R}^{n \times k}$, where $k \ll d$, to compress the data.

There are two approaches to understand PCA. The first one is maximum variance principle and the second one is minimum reconstruction loss principle.

1.1 Variance Maximization

We first consider the case when $k = 1$. Let $\mathbf{u}_1 \in \mathbb{R}^k$ be the projection vector, the variance in subspace can be written as:

$$\frac{1}{N} \sum_{n=1}^N \{\mathbf{u}_1^T \mathbf{x}_n - \mathbf{u}_1^T \bar{\mathbf{x}}\}^2 = \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1,$$
$$\text{where } \mathbf{S} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}}) (\mathbf{x}_n - \bar{\mathbf{x}})^T,$$
$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n.$$

The problem becomes to:

$$\max_{\mathbf{u}_1} \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$$

The above optimization problem is ill posed since we can easily set \mathbf{u}_1 to be infinitely large or small to make the objective not upper bounded. Notice we only care about the direction of \mathbf{u}_1 , we can constrain \mathbf{u}_1 to be unit length. Then, we have:

$$\max_{\mathbf{u}_1^T \mathbf{u}_1 = 1} \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$$

Introducing Lagrangian multiplier, and take the derivative w.r.t. \mathbf{u}_1 , we have:

$$\max_{\mathbf{u}_1, \lambda_1} \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^T \mathbf{u}_1)$$

$$\mathbf{S} \mathbf{u}_1 = \lambda_1 \mathbf{u}_1$$

$$\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 = \lambda_1$$

Noticing $\mathbf{S} \mathbf{u}_1 = \lambda_1 \mathbf{u}_1$ indicates that the largest value of variance is taken when λ_1 is the eigenvalue of \mathbf{S} while \mathbf{u}_1 is the corresponding eigenvector, and $\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$ indicates that λ_1 should be the largest eigenvalue.

For the case $k = 1$, note every time, we are finding the orthogonal subspace, which requires that \mathbf{u}_1 and \mathbf{u}_2 be orthogonal to each other, which has the form $\mathbf{u}_1^T \mathbf{u}_2 = 0$. Add this constrain and follow the similar rule, we can get that \mathbf{u}_2 is the the eigenvector of second largest eigenvalue of centered sample matrix. When $k > 2$, it can be easily proofed by induction.

1.2 Minimum Reconstruction Error

Let \mathbf{X} be a mean-centered data matrix, \mathcal{P}_k as the set of N-dimensional rank-k orthogonal projection matrix. PCA tries to seek the projection matrix, such that after projection and inverse-projection, the reconstruction error is minimum. Consider the following optimization problem (here $\mathbf{X} \in \mathbb{R}^{d \times n}$):

$$\min_{\mathbf{P} \in \mathcal{P}_k} \|\mathbf{P}\mathbf{X} - \mathbf{X}\|_F^2,$$

where $\mathcal{P}_k = \mathbf{U}_k \mathbf{U}_k^T$, and $\mathbf{U}_k \in \mathbb{R}^{d \times k}$. After we solve this optimization problem, the subspace representation is $\mathbf{Y} = \mathbf{U}_k^T \mathbf{X}$.

The above optimization problem can be formalized as following

$$\begin{aligned} \|\mathbf{P}\mathbf{X} - \mathbf{X}\|_F^2 &= \text{Tr}[(\mathbf{P}\mathbf{X} - \mathbf{X})^T(\mathbf{P}\mathbf{X} - \mathbf{X})] = \text{Tr}[\mathbf{X}^T \mathbf{P}^2 \mathbf{X} - 2\mathbf{X}^T \mathbf{P}\mathbf{X} + \mathbf{X}^T \mathbf{X}] \\ &= -\text{Tr}[\mathbf{X}^T \mathbf{P}\mathbf{X}] + \text{Tr}[\mathbf{X}^T \mathbf{X}] \end{aligned}$$

, The above derivation use the fact that \mathbf{P} is orthogonal matrix and projection matrix. ($\mathbf{P}^T \mathbf{P} = \mathbf{I}$, $\mathbf{P}\mathbf{P} = \mathbf{P}$).

Remove the constant part, we have:

$$\underset{\mathbf{P} \in \mathcal{P}_k}{\text{argmin}} \|\mathbf{P}\mathbf{X} - \mathbf{X}\|_F^2 = \underset{\mathbf{P} \in \mathcal{P}_k}{\text{argmax}} \text{Tr}[\mathbf{X}^T \mathbf{P}\mathbf{X}],$$

By the constrain $\mathcal{P}_k = \mathbf{U}_k \mathbf{U}_k^T$, we get:

$$\text{Tr}[\mathbf{X}^T \mathbf{P}\mathbf{X}] = \text{Tr}[\mathbf{U}^T \mathbf{X} \mathbf{X}^T \mathbf{U}] = \sum_{i=1}^k \mathbf{u}_i^T \mathbf{X} \mathbf{X}^T \mathbf{u}_i.$$

This form is the same as maximize the variance and thus we can get the same conclusion.

2 Probabilistic PCA

Probabilistic PCA (PPCA) provides the probabilistic view of PCA. It is more fast the PCA, and can be used when data contains missing value (EM algorithm).

Assumption of PPCA:

$$\begin{aligned} p(\mathbf{z}) &= \mathcal{N}(\mathbf{z} | \mathbf{0}, \mathbf{I}) \\ \mathbf{x} &= \mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\epsilon} \\ p(\boldsymbol{\epsilon}) &= \mathcal{N}(\boldsymbol{\epsilon} | \mathbf{0}, \sigma^2 \mathbf{I}) \end{aligned}$$

2.1 MLE for PPCA

Derive the distribution of \mathbf{x} by the above assumption, we get

$$\begin{aligned} p(\mathbf{x}) &= \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \mathbf{C}) \\ \mathbf{C} &= \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I} \end{aligned}$$

Then, we want to maximize the likelihood function of \mathbf{x} w.r.t $\boldsymbol{\mu}$, \mathbf{W} , σ .

The log-likelihood of \mathbf{X} is

$$\begin{aligned} \ln p(\mathbf{X} | \boldsymbol{\mu}, \mathbf{W}, \sigma^2) &= \sum_{n=1}^N \ln p(\mathbf{x}_n | \mathbf{W}, \boldsymbol{\mu}, \sigma^2) \\ &= -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln |\mathbf{C}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \end{aligned}$$

μ has the closed-form solution as sample mean. Substitute $\mu = \frac{1}{n} \sum_i \mathbf{x}_i$ back to the above likelihood function, we have:

$$\ln p(\mathbf{X}|\mathbf{W}, \boldsymbol{\mu}, \sigma^2) = -\frac{N}{2} \{D \ln(2\pi) + \ln |\mathbf{C}| + \text{Tr}(\mathbf{C}^{-1}\mathbf{S})\}$$

Then, we take derivative w.r.t \mathbf{W} and σ , we can get the closed-form solution as follows:

$$\begin{aligned} \mathbf{W}_{\text{ML}} &= \mathbf{U}_M (\mathbf{L}_M - \sigma^2 \mathbf{I})^{1/2} \mathbf{R} \\ \sigma_{\text{ML}}^2 &= \frac{1}{D-M} \sum_{i=M+1}^D \lambda_i \end{aligned}$$

where \mathbf{U}_M is a $D \times M$ matrix whose columns are given by any subset (of size M) of the eigenvectors of the data covariance matrix \mathbf{S} , the $M \times M$ diagonal matrix \mathbf{L}_M has elements given by the corresponding eigenvalues λ_i , and \mathbf{R} is an arbitrary $M \times M$ orthogonal matrix. Furthermore, Tipping and Bishop (1999b) showed that the maximum of the likelihood function is obtained when the M eigenvectors are chosen to be those whose eigenvalues are the M largest (all other solutions being saddle points).

2.2 EM for PPCA

Note \mathbf{Z} is latent variable, we can directly borrow EM to optimize it. The complete data log likelihood function is:

$$\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\mu}, \mathbf{W}, \sigma^2) = \sum_{n=1}^N \{\ln p(\mathbf{x}_n|\mathbf{z}_n) + \ln p(\mathbf{z}_n)\}$$

Initialize the \mathbf{W} , σ E step: (take expectation over the latent variable):

$$\begin{aligned} \mathbb{E}[\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\mu}, \mathbf{W}, \sigma^2)] &= -\sum_{n=1}^N \left\{ \frac{D}{2} \ln(2\pi\sigma^2) + \frac{1}{2} \text{Tr}(\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^T]) \right. \\ &\quad \left. + \frac{1}{2\sigma^2} \|\mathbf{x}_n - \boldsymbol{\mu}\|^2 - \frac{1}{\sigma^2} \mathbb{E}[\mathbf{z}_n]^T \mathbf{W}^T (\mathbf{x}_n - \boldsymbol{\mu}) \right. \\ &\quad \left. + \frac{1}{2\sigma^2} \text{Tr}(\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^T] \mathbf{W}^T \mathbf{W}) \right\} \end{aligned}$$

The expectation of \mathbf{z} given \mathbf{W} , σ are

$$\begin{aligned} \mathbb{E}[\mathbf{z}_n] &= \mathbf{M}^{-1} \mathbf{W}^T (\mathbf{x}_n - \bar{\mathbf{x}}) \\ \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^T] &= \sigma^2 \mathbf{M}^{-1} + \mathbb{E}[\mathbf{z}_n] \mathbb{E}[\mathbf{z}_n]^T \end{aligned} \tag{1}$$

M step: maximize the likelihood respect to \mathbf{W} and σ

$$\begin{aligned} \mathbf{W}_{\text{new}} &= \left[\sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}}) \mathbb{E}[\mathbf{z}_n]^T \right] \left[\sum_{n=1}^N \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^T] \right]^{-1} \\ \sigma_{\text{new}}^2 &= \frac{1}{ND} \sum_{n=1}^N \left\{ \|\mathbf{x}_n - \bar{\mathbf{x}}\|^2 - 2 \mathbb{E}[\mathbf{z}_n]^T \mathbf{W}_{\text{new}}^T (\mathbf{x}_n - \bar{\mathbf{x}}) \right. \\ &\quad \left. + \text{Tr}(\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^T] \mathbf{W}_{\text{new}}^T \mathbf{W}_{\text{new}}) \right\} \end{aligned} \tag{2}$$

3 Kernel PCA

PCA is for linear dimension reduction. Now, we use kernel to mapping data \mathbf{X} to $\phi(\mathbf{X})$

By the same derivative of maximum variance, we have:

$$\begin{aligned} \mathbf{C} &= \frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \\ \mathbf{C} \mathbf{v}_i &= \lambda_i \mathbf{v}_i \\ \frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \left\{ \phi(\mathbf{x}_n)^T \mathbf{v}_i \right\} &= \lambda_i \mathbf{v}_i \end{aligned}$$

The above can be rewritten as:

$$\mathbf{v}_i = \sum_{n=1}^N a_{in} \phi(\mathbf{x}_n)$$

We know the projection vector is a linear combination of $\phi(\mathbf{x}_n)$. Left multiply $\phi(\mathbf{X})^T$ on both hand side of equation, we have

$$\frac{1}{N} \sum_{n=1}^N k(\mathbf{x}_l, \mathbf{x}_n) \sum_{m=1}^m a_{im} k(\mathbf{x}_n, \mathbf{x}_m) = \lambda_i \sum_{n=1}^N a_{in} k(\mathbf{x}_l, \mathbf{x}_n)$$

In matrix form, the above equation is:

$$\mathbf{K}^2 \mathbf{a}_i = \lambda_i N \mathbf{K} \mathbf{a}_i \quad (3)$$

Since \mathbf{K} is PSD kernel, then we can get the optimum \mathbf{a} by solving the eigenvalue problem:

$$\mathbf{K} \mathbf{a}_i = \lambda_i N \mathbf{a}_i \quad (4)$$

After we solve the eigenvalue problem, we substitute the linear span expression of projection vector, we have

$$y_i(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{v}_i = \sum_{n=1}^N a_{in} \phi(\mathbf{x})^T \phi(\mathbf{x}_n) = \sum_{n=1}^N a_{in} k(\mathbf{x}, \mathbf{x}_n) \quad (5)$$

However, the above derivation assumes the high dimensional projected data has zero mean. We need to calibrate this. Suppose the mean centered high dimensional data is

$$\tilde{\phi}(\mathbf{x}_n) = \phi(\mathbf{x}_n) - \frac{1}{N} \sum_{l=1}^N \phi(\mathbf{x}_l) \quad (6)$$

Then, we can compute the kernel of centered kernel feature:

$$\begin{aligned} \tilde{K}_{nm} &= \tilde{\phi}(\mathbf{x}_n)^T \tilde{\phi}(\mathbf{x}_m) \\ &= \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) - \frac{1}{N} \sum_{l=1}^N \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_l) \\ &\quad - \frac{1}{N} \sum_{l=1}^N \phi(\mathbf{x}_l)^T \phi(\mathbf{x}_m) + \frac{1}{N^2} \sum_{j=1}^N \sum_{l=1}^N \phi(\mathbf{x}_j)^T \phi(\mathbf{x}_l) \\ &= k(\mathbf{x}_n, \mathbf{x}_m) - \frac{1}{N} \sum_{l=1}^N k(\mathbf{x}_l, \mathbf{x}_m) - \frac{1}{N} \sum_{l=1}^N k(\mathbf{x}_n, \mathbf{x}_l) + \frac{1}{N^2} \sum_{j=1}^N \sum_{l=1}^N k(\mathbf{x}_j, \mathbf{x}_l) \end{aligned} \quad (7)$$

Matrix form of mean centered kernel is:

$$\tilde{\mathbf{K}} = \mathbf{K} - \mathbf{1}_N \mathbf{K} - \mathbf{K} \mathbf{1}_N + \mathbf{1}_N \mathbf{K} \mathbf{1}_N \quad (8)$$

4 Independent Component Analysis

- ICA is a linear method.
- ICA tries to recover the independent components rather than orthogonal components (minimize the mutual information).
- ICA will fail when more than one subspace is Gaussian

4.1 Minimize the Mutual Information

$$\begin{aligned} &l(\{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_d\}) \\ &= \sum_{i=1}^d H(\mathbf{y}_i) - H(\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_d) \\ &H(\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_d) = H(\mathbf{X}\mathbf{W}) = H(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_d) + \log \det(\mathbf{W}) \end{aligned} \quad (9)$$

The MI can be rewritten as:

$$l(\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_d) = \sum_{i=1}^n H(y_i) - H(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_d) - \log \det(W) \quad (10)$$

Assume \mathbf{y} to be unit variance

$$\begin{aligned} E(\mathbf{y}^T \mathbf{y}) &= \mathbf{W} E\{\mathbf{x} \mathbf{x}^T\} \mathbf{W}^T = \mathbf{I} \\ \det(\mathbf{I}) &= 1 = \det(\mathbf{W}) (\det E\{\mathbf{x} \mathbf{x}^T\}) \det(\mathbf{W}^T) \end{aligned} \quad (11)$$

This indicates $\det(W)$ is a constant. Now the MI can be written as

$$l(\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_d) = \sum_{i=1}^n H(y_i) - H(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_d) \quad (12)$$

Maximize the non-gaussianity can achieve ICA.

4.2 Non-Gaussianity Estimation

We will introduce a few metric for evaluating the non-gaussianity.

4.2.1 Central Limit Theorem

- Distribution of a sum of independent random variables tends toward a Gaussian distribution
- Sum of several independent random variables usually has a distribution that is closer to gaussian than any of the original random variables

4.2.2 Kurtosis

Kurtosis is defined by:

$$\text{kurt}(y) = E\{y^4\} - 3(E\{y^2\})^2 \quad (13)$$

When we constrain the y has zero mean and unit variance, the kurtosis has the form $E\{y^4\} - 3$. Kurtosis measures the degree of peakdness of a distribution and it is zero only for Gaussian distribution. However, it is sensitive to outliers due to the fourth moment.

4.2.3 Negentropy

The entropy of a random variable is defined as

$$H(y) = - \int_{-\infty}^{\infty} p(y) \log p(y) dy = -E\{\log p_i(y)\} \quad (14)$$

Denote y_G as the entropy of gaussian variable, the negentropy is defined as

$$J(y) = H(y_G) - H(y) \geq 0$$

It is hard to get $J(y)$ since we do not know the density distribution $p(y)$.

4.2.4 Approximation to Negentropy

Some approximation:

$$\begin{aligned} J(y) &\approx \frac{1}{12} E\{y^3\}^2 + \frac{1}{48} \text{kurt}(y)^2 && \text{not robust} \\ J(y) &\approx \sum_{i=1}^p k_i [E\{G_i(y)\} - E\{G_i(g)\}]^2 && \text{better approximation,} \end{aligned}$$

where k_i are some positive constants, y is assumed to have zero mean and unit variance, and g is a Gaussian variable with zero mean and unit variance. G_i are some non-quadratic functions such as

$$G_1(y) = \frac{1}{a} \log \cosh(ay), \quad G_2(y) = -\exp(-y^2/2), \quad (15)$$

where $1 \leq a \leq 2$ is some suitable constant. Although this approximation may no be accurate, it is always greater than zero except when x is Gaussian.

4.2.5 ICA as Projection Pursuit (FastICA)

We want to minimize:

$$\sum_i E \{G(y_i)\} = \sum_i E \{G(\mathbf{w}_i^T \mathbf{x})\}, \quad (16)$$

where w is a rotation matrix. Using Lagrangian Multiplier, we have:

$$O(\mathbf{w}) = E \{G(\mathbf{w}^T \mathbf{x})\} - \beta (\mathbf{w}^T \mathbf{w} - 1) / 2 \quad (17)$$

Take derivative respect to w :

$$F(\mathbf{w}) = \frac{\partial O(\mathbf{w})}{\partial \mathbf{w}} = E \{\mathbf{x}g(\mathbf{w}^T \mathbf{x})\} - \beta \mathbf{w} = 0 \quad (18)$$

We can use newton method to optimize it. The Hessian is

$$J_F(\mathbf{w}) = \frac{\partial F}{\partial \mathbf{w}} = E \{\mathbf{x}\mathbf{x}^T g'(\mathbf{w}^T \mathbf{x})\} - \beta \mathbf{I} \quad (19)$$

The first term can be approximated as

$$E \{\mathbf{x}\mathbf{x}^T g'(\mathbf{w}^T \mathbf{x})\} \approx E \{\mathbf{x}\mathbf{x}^T\} E \{g'(\mathbf{w}^T \mathbf{x})\} = E \{g'(\mathbf{w}^T \mathbf{x})\} \mathbf{I} \quad (20)$$

Then, the Hessian becomes diagonal, and the inverse is easy to calculate.

$$J_F(\mathbf{w}) = [E \{g'(\mathbf{w}^T \mathbf{x})\} - \beta] \mathbf{I} \quad (21)$$

The Newton iteration becomes:

$$\mathbf{w} \leftarrow \mathbf{w} - \frac{1}{E \{g'(\mathbf{w}^T \mathbf{x})\} - \beta} [E \{\mathbf{x}g(\mathbf{w}^T \mathbf{x})\} - \beta \mathbf{w}] \quad (22)$$

Multiplying both sides by the scalar $\beta - E \{g'(\mathbf{w}^T \mathbf{x})\}$, we have

$$\mathbf{w} \leftarrow E \{\mathbf{x}g(\mathbf{w}^T \mathbf{x})\} - E \{g'(\mathbf{w}^T \mathbf{x})\} \mathbf{w} \quad (23)$$

We are just scale the parameters, we can renormalizes it. The projection pursuit for ICA can be summarized as follows:

- Choose initial random guess for w
- Iterate:

$$\mathbf{w} \leftarrow E \{\mathbf{x}g(\mathbf{w}^T \mathbf{x})\} - E \{g'(\mathbf{w}^T \mathbf{x})\} \mathbf{w}$$

- Normalize

$$\mathbf{w} \leftarrow \mathbf{w} / \|\mathbf{w}\|$$

- If not converged, go back to step 2

Estimate the 2^{nd} ICA component similarity, just add the orthogonal constraint.

4.2.6 MLE for ICA

Another view for deriving ICA is from MLE. The following equation always holds.

$$p_x(x) = p_s(Wx) \cdot |W| \quad (24)$$

Since we assume the source are independent:

$$p(s) = \prod_{i=1}^n p_s(s_i) \quad (25)$$

So the likelihood for our observation is

$$p(x) = \prod_{i=1}^n p_s(w_i^T x) \cdot |W| \quad (26)$$

We need to assume the probability distribution of the sources, and remember we cannot choose Gaussian. A reasonable CDF is the sigmoid function $g(s) = 1 / (1 + e^{-s})$, and the pdf is $p_s(s) = g'(s)$ The log likelihood is:

$$\ell(W) = \sum_{i=1}^m \left(\sum_{j=1}^n \log g'(w_j^T x^{(i)}) + \log |W| \right) \quad (27)$$

, Applying gradient descent or SGD, we can get the update rule. $(\nabla_W |W| = |W| (W^{-1})^T)$

5 Non-linear ICA

- It is an open problem
- It is largely related to Disentangling Representation
- Beta VAE, Factor VAE

5.1 VAE

$$\begin{aligned}
KL(q(z)||p(z|x)) &= \int q(z) \log \frac{q(z)}{p(z|x)} dz \\
&= \int q(z) [\log q(z) - \log p(z|x)] dz \\
&= \int q(z) [\log q(z) - \log p(X|z) - \log p(z)] dz + \log p(X) \\
&= \int q(z) \log p(X|z) dz - KL(q(z)||p(z)) + \log p(X) \\
&\leq \int q(z) \log p(X|z) dz - KL(q(z)||p(z))
\end{aligned} \tag{28}$$

5.2 Beta-VAE

$$\begin{aligned}
&\max_{\theta} \int p(x|z) dz \\
&\text{s.t. } KL(q_{\phi}(z|x)||p(z)) < \epsilon
\end{aligned} \tag{29}$$

Objective of Beta-VAE

$$\mathcal{F}(\theta, \phi, \beta; \mathbf{x}, \mathbf{z}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \beta (D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) - \epsilon) \tag{30}$$

5.3 Factor-VAE

Objective:

$$\mathbb{E}_{p(\mathbf{x})} [\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))] - \gamma D_{KL}\left(q(\mathbf{z})||\prod_{j=1}^d q(\mathbf{z}_j)\right) \tag{31}$$