

Computational Intelligence Lab

Cristian Perez Jensen

April 14, 2024

Contents

1	<i>Dimensionality reduction</i>	1
1.1	<i>Linear autoencoders</i>	1
1.2	<i>Projection</i>	2
1.3	<i>Principal component analysis</i>	6
1.4	<i>Learning algorithms</i>	6
1.5	<i>Non-linear autoencoders</i>	7

List of symbols

\doteq	Equality by definition
\approx	Approximate equality
\propto	Proportional to
\mathbb{N}	Set of natural numbers
\mathbb{R}	Set of real numbers
$i : j$	Set of natural numbers between i and j . I.e., $\{i, i+1, \dots, j\}$
$f : A \rightarrow B$	Function f that maps elements of set A to elements of set B
$\mathbb{1}\{\text{predicate}\}$	Indicator function (1 if predicate is true, otherwise 0)
$\mathbf{v} \in \mathbb{R}^n$	n -dimensional vector
$\mathbf{M} \in \mathbb{R}^{m \times n}$	$m \times n$ matrix
$\mathbf{T} \in \mathbb{R}^{d_1 \times \dots \times d_n}$	Tensor
\mathbf{M}^\top	Transpose of matrix \mathbf{M}
\mathbf{M}^{-1}	Inverse of matrix \mathbf{M}
$\det(\mathbf{M})$	Determinant of \mathbf{M}
$\frac{d}{dx}f(x)$	Ordinary derivative of $f(x)$ w.r.t. x at point $x \in \mathbb{R}$
$\frac{\partial}{\partial x}f(\mathbf{x})$	Partial derivative of $f(\mathbf{x})$ w.r.t. x at point $\mathbf{x} \in \mathbb{R}^n$
$\nabla_{\mathbf{x}}f(\mathbf{x}) \in \mathbb{R}^n$	Gradient of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at point $\mathbf{x} \in \mathbb{R}^n$
$\nabla_{\mathbf{x}}^2f(\mathbf{x}) \in \mathbb{R}^{n \times n}$	Hessian of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at point $\mathbf{x} \in \mathbb{R}^n$
$\boldsymbol{\theta} \in \Theta$	Parametrization of a model, where Θ is a compact subset of \mathbb{R}^K
\mathcal{X}	Input space
\mathcal{Y}	Output space
$\mathcal{D} \subseteq \mathcal{X} \times \mathcal{Y}$	Labeled training data

1 Dimensionality reduction

The motivation behind dimensionality reduction is to find a low-dimensional representation of high-dimensional data.¹ This can be split into two goals: (1) compressing the data, while preserving as much as possible of the relevant information, and (2) interpreting the data in low dimensionality is easier than high dimensionality.

Dimensionality reduction is often performed by an autoencoder, which typically has a bottleneck and aims to predict its input; see Figure 1. In general, we have an encoder F and a decoder G ,

$$F : \mathbb{R}^n \rightarrow \mathbb{R}^m, \quad G : \mathbb{R}^m \rightarrow \mathbb{R}^n. \quad m \ll n.$$

The reconstruction function is then the following function,

$$G \circ F : \mathbb{R}^n \rightarrow \mathbb{R}^n,$$

which is ideally the identity function.

1.1 Linear autoencoders

To build a nice theory, we will only consider a single layer linear autoencoder, which means that we have the following functions,

$$\begin{aligned} F : x \mapsto z &= Wx, \quad W \in \mathbb{R}^{m \times n}. \\ G : z \mapsto \hat{x} &= Vz, \quad V \in \mathbb{R}^{n \times m}. \end{aligned}$$

The objective to minimize of the linear encoder is the following,

$$\mathcal{R}(W, V) = \mathcal{R}(P \doteq VW) \doteq \mathbb{E} \left[\frac{1}{2} \|x - Px\|^2 \right].$$

Corollary. For centered data, i.e. $E[x] = 0$, optimal affine maps degenerate to linear ones.

Proof. Proof by contradiction. Let $a \neq 0$, then

$$\begin{aligned} \mathbb{E} \left[\|x - (Px + a)\|^2 \right] &= \mathbb{E} [\langle x - Px - a, x - Px - a \rangle] \\ &= \mathbb{E} [\langle x - Px, x - Px \rangle + 2\langle x - Px, -a \rangle + \langle -a, -a \rangle] \\ &= \mathbb{E} \left[\|x - Px\|^2 \right] - 2 \underbrace{\langle \mathbb{E}[x] - P\mathbb{E}[x], a \rangle}_{=0} + \underbrace{\|a\|^2}_{>0} \\ &> \mathbb{E} \left[\|x - Px\|^2 \right]. \end{aligned}$$

Thus, the risk is strictly worse if $a \neq 0$. ■

Thus, we will assume that the data is centered, which makes the analysis easier, since we do not need to consider the affine case.

¹ Often, the original raw representation is high-dimensional and redundant, e.g., images, audio, time series.

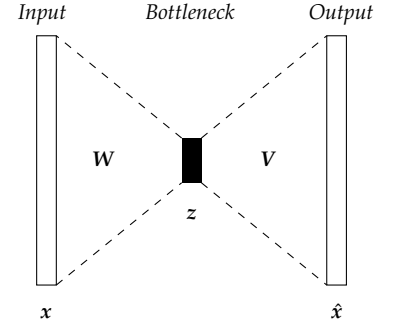


Figure 1. Diagram of a single layer linear autoencoder.

Note that while the optimal linear reconstruction map P is unique, its parametrization W, V is not unique, since for any invertible matrix $A \in \mathbb{R}^{m \times m}$, we can construct an optimal parametrization,

$$VW = VIW = V(AA^{-1})W = (VA)(A^{-1}W),$$

with $A^{-1}W, VA$. The weight matrices are non-identifiable.

Since P cannot be any $n \times n$ matrix, we want to know how the composition of $W \in \mathbb{R}^{m \times n}$ and $V \in \mathbb{R}^{n \times m}$ characterizes the matrix P and which constraints they impose. The answer to this is that the weight matrices impose a rank constraint on P ,

$$\text{rank}(P) = \min\{\text{rank}(W), \text{rank}(V)\} \leq \min\{m, n\} = m.$$

Thus, when optimizing for P , we are constrained to matrices with rank less or equal to m .

1.2 Projection

The rank constraint and linearity of P means that the image (column space) of P is a linear subspace $\mathcal{U} \subseteq \mathbb{R}^n$ of dimension at most m . We will break the solution to our problem into two parts: (1) finding the optimal subspace \mathcal{U} , and (2) finding the optimal mapping to that subspace.²

Finding the optimal mapping to a subspace. We will first focus on (2); given subspace \mathcal{U} , we need to determine the optimal linear map P^* , such that

$$P^* = \underset{P}{\operatorname{argmin}} \|x - Px\|^2, \quad \text{col}(P) = \mathcal{U}.$$

Definition 1 (Orthogonal projection). A linear transformation $P : \mathcal{V} \rightarrow \mathcal{V}$ is called an orthogonal projection onto \mathcal{U} if $\forall x \in \mathcal{V}$:

1. *Projection*: $P(x) \in \mathcal{U}$;
2. *Orthogonality*: $\text{null}(P) \perp \text{col}(P)$, which is equivalent to the following holding, $\langle P(x), y \rangle = \langle x, P(y) \rangle$ (self-adjointness);
3. *Idempotency*: $P(P(x)) = P(x)$.

Definition 2. The orthogonal projection to a linear subspace $\mathcal{U} \subseteq \mathbb{R}^n$ is defined as

$$\Pi_{\mathcal{U}} : \mathbb{R}^n \rightarrow \mathcal{U}, \quad \Pi_{\mathcal{U}}(x) = \underset{x' \in \mathcal{U}}{\operatorname{argmin}} \|x - x'\|.$$

² We do not search for the weight matrices W, V , since they are not unique, but P is unique.

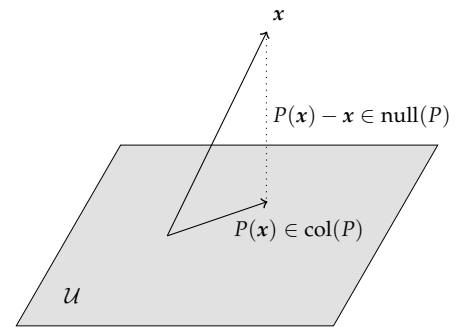


Figure 2. Orthogonal projection of x onto subspace plane \mathcal{U} .

Proof. We need to show that the definition of $\Pi_{\mathcal{U}}$ indeed is an orthogonal projection by showing that it adheres to the properties of Definition 1.

1. *Projection:* This is true by definition of the values that the argmin are allowed to take on;
2. *Idempotency:* For all $\mathbf{u} \in \mathcal{U}$, $\Pi_{\mathcal{U}}(\mathbf{u}) = \operatorname{argmin}_{\mathbf{x}' \in \mathcal{U}} \|\mathbf{u} - \mathbf{x}'\| = \mathbf{u}$. Thus, $\Pi_{\mathcal{U}} = \Pi_{\mathcal{U}} \circ \Pi_{\mathcal{U}}$;
3. *Orthogonality:* We need to show that $\Pi_{\mathcal{U}}(\mathbf{x}) - \mathbf{x} \in \mathcal{U}^{\perp}$. Decompose it into $\mathbf{u} \in \mathcal{U}$ and $\mathbf{u}^{\perp} \in \mathcal{U}^{\perp}$ by $\Pi_{\mathcal{U}}(\mathbf{x}) - \mathbf{x} = \mathbf{u} + \mathbf{u}^{\perp}$. Then, we only need to show that $\mathbf{u} = \mathbf{0}$.

Proof by contradiction. Let $\mathbf{u} \neq \mathbf{0}$, then

$$\begin{aligned} \|\Pi_{\mathcal{U}}(\mathbf{x}) - \mathbf{x}\|^2 &= \underbrace{\|\mathbf{u}\|^2}_{>0} + \|\mathbf{u}^{\perp}\|^2 + 2 \underbrace{\langle \mathbf{u}, \mathbf{u}^{\perp} \rangle}_{=0} \\ &> \|\mathbf{u}^{\perp}\|^2 \\ &= \|\underbrace{\Pi_{\mathcal{U}}(\mathbf{x}) - \mathbf{u}}_{\in \mathcal{U}} - \mathbf{x}\|^2, \end{aligned}$$

which contradicts with

$$\|\Pi_{\mathcal{U}}(\mathbf{x}) - \mathbf{x}\| = \min_{\mathbf{x}' \in \mathcal{U}} \|\mathbf{x}' - \mathbf{x}\| \leq \|\tilde{\mathbf{u}} - \mathbf{x}\|, \quad \forall \tilde{\mathbf{u}} \in \mathcal{U}.$$

Hence, $\Pi_{\mathcal{U}}(\mathbf{x}) - \mathbf{x} \in \mathcal{U}^{\perp}$ and is unique;

4. *Linearity:* We need to show homogeneity,

$$\begin{aligned} \Pi_{\mathcal{U}}(\alpha \mathbf{x}) &= \operatorname{argmin}_{\mathbf{x}'' \in \mathcal{U}} \|\alpha \mathbf{x} - \mathbf{x}''\| \\ &= \operatorname{argmin}_{\alpha \mathbf{x}' \in \mathcal{U}} \|\alpha \mathbf{x} - \alpha \mathbf{x}'\| & \mathbf{x}' = \frac{1}{\alpha} \mathbf{x}'' \\ &= \alpha \operatorname{argmin}_{\mathbf{x}' \in \mathcal{U}} \|\mathbf{x} - \mathbf{x}'\| & \operatorname{argmin}_{\lambda \mathbf{x}} f(\mathbf{x}) = \lambda \operatorname{argmin}_{\mathbf{x}} f(\mathbf{x}). \\ &= \alpha \operatorname{argmin}_{\mathbf{x}' \in \mathcal{U}} \|\mathbf{x} - \mathbf{x}'\| \\ &= \alpha \Pi_{\mathcal{U}}(\mathbf{x}). \end{aligned}$$

And, we need to show additivity,

$$\begin{aligned} \Pi_{\mathcal{U}}(\mathbf{x} + \mathbf{y}) &= \operatorname{argmin}_{\mathbf{z} \in \mathcal{U}} \|\mathbf{z} - (\mathbf{x} + \mathbf{y})\| \\ &= \operatorname{argmin}_{\Pi_{\mathcal{U}}(\mathbf{x}) + \mathbf{y}' \in \mathcal{U}} \|\Pi_{\mathcal{U}}(\mathbf{x}) + \mathbf{y}' - \mathbf{x} - \mathbf{y}\| & \mathbf{z} \doteq \Pi_{\mathcal{U}}(\mathbf{x}) + \mathbf{y}' \text{ for some } \mathbf{y}' \in \mathcal{U}. \\ &= \Pi_{\mathcal{U}}(\mathbf{x}) + \operatorname{argmin}_{\mathbf{y}' \in \mathcal{U}} \|(\mathbf{y}' - \mathbf{y}) + (\Pi_{\mathcal{U}}(\mathbf{x}) - \mathbf{x})\|^2 & \operatorname{argmin}_{\mathbf{x} + \lambda} f(\mathbf{x}) = \lambda + \operatorname{argmin}_{\mathbf{x}} f(\mathbf{x}). \\ &= \Pi_{\mathcal{U}}(\mathbf{x}) + \operatorname{argmin}_{\mathbf{y}' \in \mathcal{U}} \|\mathbf{y}' - \mathbf{y}\|^2 + \|\Pi_{\mathcal{U}}(\mathbf{x}) - \mathbf{x}\|^2 + 2 \langle \Pi_{\mathcal{U}}(\mathbf{x}) - \mathbf{x}, \mathbf{y}' - \mathbf{y} \rangle \\ &= \Pi_{\mathcal{U}}(\mathbf{x}) + \operatorname{argmin}_{\mathbf{y}' \in \mathcal{U}} \|\mathbf{y}' - \mathbf{y}\|^2 + 2 \underbrace{\langle \Pi_{\mathcal{U}}(\mathbf{x}) - \mathbf{x}, \mathbf{y}' \rangle}_{\langle \mathcal{U}^{\perp}, \mathcal{U} \rangle = 0} & \langle \Pi_{\mathcal{U}}(\mathbf{x}) - \mathbf{x}, \mathbf{y}' - \mathbf{y} \rangle = \langle \Pi_{\mathcal{U}}(\mathbf{x}) - \mathbf{x}, \mathbf{y}' \rangle - \langle \Pi_{\mathcal{U}}(\mathbf{x}) - \mathbf{x}, \mathbf{y} \rangle. \\ &= \Pi_{\mathcal{U}}(\mathbf{x}) + \Pi_{\mathcal{U}}(\mathbf{y}). \end{aligned}$$

■

So, we know that $\Pi_{\mathcal{U}}$ is an orthogonal projection. Now, we want to find the matrix P representing that linear transformation.

Given an orthonormal basis U of \mathcal{U} , we can compute the optimal projection matrix,

$$P = UU^T.$$

Note that in this case W and V share parameters, and UU^T is the optimal weight matrix if we enforce parameter sharing via $V = W^T$.

Proof. The image of the projection matrix is \mathcal{U} ,

$$Px = \left(\sum_{i=1}^m u_i u_i^T \right) x = \sum_{i=1}^m u_i u_i^T x = \sum_{i=1}^m \langle u_i, x \rangle u_i.$$

Self-adjointness,

$$P^T = (UU^T)^T = UU^T = P.$$

Idempotency,

$$PP = U \underbrace{U^T U}_{I_m} U^T = UU^T = P.$$

Orthogonality, TODO

■

In general, we do not have an orthonormal basis for \mathcal{U} . In a non-orthonormal basis V for \mathcal{U} , we can recover the projection matrix,

$$P = VV^+, \quad V^+ \doteq (V^T V)^{-1} V^T.$$

Note that V^+ is the left Moore-Penrose pseudo-inverse of V .

Proof. P is the projection matrix of \mathcal{U} ,

$$PV = VV^+V = V(V^T V)^{-1} V^T V = V.$$

Together with the rank constraint, this yields $Pu^\perp = 0$ for all $u^\perp \in \mathcal{U}^\perp$.

Self-adjointness,

$$P^T = (V(V^T V)^{-1} V^T)^T = V(V^T V)^{-1} V^T = P.$$

Idempotency,

$$PP = V(V^T V)^{-1} V^T V(V^T V)^{-1} V^T = VV^+ = P.$$

■

Finding the optimal subspace. Now we need to find out which subspace of dimension m or less is optimal to project onto. First, we need to rewrite the objective function to find a new interpretation,

$$\begin{aligned}
 \mathcal{R}(P) &= \frac{1}{2} \mathbb{E} [\|x - Px\|^2] \\
 &= \frac{1}{2} \mathbb{E} [\langle x, x \rangle + 2\langle x, -Px \rangle + \langle -Px, -Px \rangle] \\
 &= \frac{1}{2} \mathbb{E} [\|x\|^2] + \frac{1}{2} \mathbb{E} [\|Px\|^2] - \mathbb{E} [\langle x, Px \rangle] \\
 &= \frac{1}{2} \mathbb{E} [\|x\|^2] + \frac{1}{2} \mathbb{E} [\|Px\|^2] - \mathbb{E} [\langle x, P^2x \rangle] \\
 &= \frac{1}{2} \mathbb{E} [\|x\|^2] + \frac{1}{2} \mathbb{E} [\|Px\|^2] - \mathbb{E} [\|Px\|^2] \\
 &= \frac{1}{2} \mathbb{E} [\|x\|^2] - \frac{1}{2} \mathbb{E} [\|Px\|^2].
 \end{aligned}$$

Idempotency of projection matrices.

$$\langle x, P^2x \rangle = \langle Px, Px \rangle.$$

Because our data is centered, we know the following,

$$\begin{aligned}
 \text{Var}[x] &= \mathbb{E} [\|x\|^2] - \underbrace{\|\mathbb{E}[x]\|^2}_{=0} \\
 &= \mathbb{E} [\|x\|^2]. \\
 \text{Var}[Px] &= \mathbb{E} [\|Px\|^2] - \|\mathbb{E}[Px]\|^2 \\
 &= \mathbb{E} [\|Px\|^2] - \underbrace{\|P\mathbb{E}[x]\|^2}_{=0} \\
 &= \mathbb{E} [\|Px\|^2].
 \end{aligned}$$

$$\mathcal{R}(P) = \frac{1}{2} (\text{Var}[x] - \text{Var}[Px]) \propto -\frac{1}{2} \text{Var}[Px].$$

Hence, minimizing $\mathcal{R}(P)$ is equivalent to maximizing the variance $\text{Var}[Px]$.

We can further simplify this expression to find a sufficient statistic for the objective function,

$$\begin{aligned}
 -\frac{1}{2} \text{Var}[Px] &= -\frac{1}{2} \mathbb{E} [\|Px\|^2] \\
 &= -\frac{1}{2} \mathbb{E} [\langle x, Px \rangle] \\
 &= -\frac{1}{2} \mathbb{E} [\text{tr}(x^\top Px)] \\
 &= -\frac{1}{2} \text{tr}(\mathbb{E} [Pxx^\top]) \\
 &= -\frac{1}{2} \text{tr}(P\mathbb{E} [xx^\top]).
 \end{aligned}$$

$$\|Px\|^2 = \langle Px, Px \rangle = \langle x, P^2x \rangle = \langle x, Px \rangle.$$

Cyclic property of trace.

$\mathbb{E}[xx^\top]$ is a sufficient statistic for $\mathcal{R}(P)$.

The optimal projection is fully determined by the covariance matrix $\mathbb{E}[\mathbf{x}\mathbf{x}^\top]$, together with $\mathbb{E}[\mathbf{x}]$ for centering.

1.3 Principal component analysis

Theorem 3 (Spectral theorem). Any symmetric and positive semidefinite matrix Σ can be non-negatively diagonalized with an orthogonal matrix,

$$\Sigma = Q\Lambda Q^\top, \quad \Lambda = \text{diag}(\lambda_1, \dots, \lambda_n),$$

where $\lambda_1 \geq \dots \geq \lambda_n \geq 0$ and Q is orthogonal.

Remark. Q is composed of ordered eigenvectors of Σ , and Λ is composed of ordered eigenvalues of Σ .

Theorem 4 (PCA theorem). The variance maximizing projection matrix P for a covariance matrix $\mathbb{E}[\mathbf{x}\mathbf{x}^\top] = Q\Lambda Q^\top$ as in the spectral theorem is given by

$$P = UU^\top, \quad U = Q \begin{bmatrix} I_m \\ \mathbf{0} \end{bmatrix}.$$

Proof.

$$\begin{aligned} \text{Var}[P\mathbf{x}] &= \text{tr}(P\mathbb{E}[\mathbf{x}\mathbf{x}^\top]) \\ &= \text{tr}(UU^\top Q\Lambda Q^\top) \\ &= \text{tr}((Q^\top U)(Q^\top U)^\top \Lambda). \end{aligned}$$

$$P = UU^\top, \quad \mathbb{E}[\mathbf{x}\mathbf{x}^\top] = Q\Lambda Q^\top.$$

Cyclic property.

This term is maximized by $Q^\top U = \begin{bmatrix} I_m & \mathbf{0} \end{bmatrix}^\top$. ■

1.4 Learning algorithms

Eigenvalue decomposition of the (symmetric) sample covariance matrix has $\mathcal{O}(n^3)$ complexity. Furthermore, the complexity of computing $\mathbb{E}[\mathbf{x}\mathbf{x}^\top]$ is $\mathcal{O}(Nn^2)$.³ This is quite costly, thus we need to search for algorithms that have lower runtime complexity.

³ Typically, $N \gg n$.

Power method. The power method is a recursive algorithm for computing principal eigenvectors. It initializes a vector at random $\mathbf{v}^{(0)} \sim \mathcal{N}(\mathbf{0}, I)$. Then, it iteratively improves this guess,

$$\mathbf{v}^{(t+1)} = \frac{A\mathbf{v}^{(t)}}{\|A\mathbf{v}^{(t)}\|}.$$

The computational complexity of this algorithm is $\mathcal{O}(Tn^2)$.

Lemma 5. Let \mathbf{u}_1 be the unique principal eigenvector of a diagonalizable matrix \mathbf{A} with eigenvalues $\lambda_1 > \lambda_2 \geq \dots \geq \lambda_n \geq 0$. If $\langle \mathbf{v}_0, \mathbf{u}_1 \rangle \neq 0$, then

$$\lim_{t \rightarrow \infty} \mathbf{v}^{(t)} = \mathbf{u}_1.$$

Proof. We can decompose vectors as a linear combination of eigenvectors, $\mathbf{v}^{(0)} = \sum_{i=1}^n \alpha_i \mathbf{u}_i$. Then,

$$\begin{aligned} \mathbf{v}^{(k)} &\propto \mathbf{A}^k \mathbf{v}^{(0)} \\ &= \sum_{i=1}^n \alpha_i \lambda_i^k \mathbf{u}_i & \mathbf{A}^k \mathbf{v}^{(0)} &= \sum_{i=1}^n \alpha_i \mathbf{A}^k \mathbf{u}_i. \\ &\propto \alpha_1 \mathbf{u}_1 + \sum_{i=2}^n \alpha_i \left(\frac{\lambda_i}{\lambda_1} \right)^k \mathbf{u}_i. & \text{Divide by } \lambda_1. \end{aligned}$$

$\lambda_i/\lambda_1 < 1$ for $i > 1$, thus the sum goes to 0 and $\mathbf{v}^{(k)} \rightarrow \mathbf{u}_1$. ■

We can use this algorithm to also compute the next principal eigenvectors by factoring out \mathbf{u}_1 and then doing the algorithm again to recover \mathbf{u}_2 , and continue doing that until we have the m principal eigenvectors.

Thus, the total complexity of finding the m principal eigenvectors is $\mathcal{O}(Tmn^2)$. However, this does not get rid of the $\mathcal{O}(Nn^2)$ complexity for computing the sample covariance matrix.

Gradient descent. By treating the autoencoder as a neural network, we can use deep learning techniques, such as gradient descent. Gradient descent iteratively updates the weights by

$$\mathbf{P}^{(t+1)} = \mathbf{P}^{(t)} - \eta \nabla_{\mathbf{P}} \mathcal{R}(\mathbf{P}).$$

The gradient is computed by $(\mathbf{P} - \mathbf{I})\mathbf{x}\mathbf{x}^\top$. The problem with this is that we cannot constrain \mathbf{P} to be a projection. Thus, we actually need to update \mathbf{V} and \mathbf{W} . Thus, by the chain rule for matrix derivatives,

$$\begin{aligned} \nabla_{\mathbf{W}} \mathcal{R}(\mathbf{W}, \mathbf{V}) &= (\mathbf{P} - \mathbf{I})\mathbf{x}\mathbf{x}^\top \mathbf{W}^\top \\ \nabla_{\mathbf{V}} \mathcal{R}(\mathbf{W}, \mathbf{V}) &= \mathbf{V}^\top (\mathbf{P} - \mathbf{I})\mathbf{x}\mathbf{x}^\top. \end{aligned}$$

The complexity for T iterations is then $\mathcal{O}(T(m+k)n^2)$, where k is the batch size.

1.5 Non-linear autoencoders

We can get much better performance by considering non-linearities, and we can easily use gradient descent to work with non-linear architectures.