

Computational Intelligence Lab

Cristian Perez Jensen

September 16, 2024

Note that these are not the official lecture notes of the course, but only notes written by a student of the course. As such, there might be mistakes. The source code can be found at github.com/cristianpjensen/eth-cs-notes. If you find a mistake, please create an issue or open a pull request.

Contents

1	<i>Preliminaries</i>	1
1.1	<i>Vector spaces</i>	1
1.2	<i>Norms</i>	1
1.3	<i>Matrices</i>	2
1.4	<i>Eigenvalues and eigenvectors</i>	3
1.5	<i>Convexity</i>	4
2	<i>Dimensionality reduction</i>	5
2.1	<i>Linear autoencoders</i>	5
2.2	<i>Projection</i>	6
2.3	<i>Principal component analysis</i>	10
2.4	<i>Learning algorithms</i>	11
3	<i>Matrix completion</i>	13
3.1	<i>Fully observed case</i>	14
3.2	<i>Incompletely observed case</i>	18
3.3	<i>Randomized methods for SVD</i>	23
4	<i>Latent variable models</i>	25
4.1	<i>Probabilistic clustering models</i>	25
4.2	<i>Topic models</i>	28
4.3	<i>Embeddings</i>	30
5	<i>Deep neural networks</i>	34
5.1	<i>Backpropagation</i>	35
5.2	<i>Gradient methods</i>	36
5.3	<i>Convolutional neural networks</i>	38
6	<i>Generative models</i>	40
6.1	<i>Autoregressive models</i>	40
6.2	<i>Variational autoencoders</i>	42
6.3	<i>Generative adversarial networks</i>	43
6.4	<i>Diffusion models</i>	44

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>I.e.</i> , $\{i, i+1, \dots, j\}$
$f : A \rightarrow B$	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)
$\boldsymbol{v} \in \mathbb{R}^n$	n -dimensional vector
$\boldsymbol{M} \in \mathbb{R}^{m \times n}$	$m \times n$ matrix
\boldsymbol{M}^\top	Transpose of matrix \boldsymbol{M}
\boldsymbol{M}^{-1}	Inverse of matrix \boldsymbol{M}
$\det(\boldsymbol{M})$	Determinant of \boldsymbol{M}
$\frac{\mathrm{d}}{\mathrm{d}x}f(x)$	Ordinary derivative of $f(x)$ w.r.t. x at point $x \in \mathbb{R}$
$\frac{\partial}{\partial \boldsymbol{x}}f(\boldsymbol{x})$	Partial derivative of $f(\boldsymbol{x})$ w.r.t. \boldsymbol{x} at point $\boldsymbol{x} \in \mathbb{R}^n$
$\nabla_{\boldsymbol{x}}f(\boldsymbol{x}) \in \mathbb{R}^n$	Gradient of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at point $\boldsymbol{x} \in \mathbb{R}^n$
$\nabla_{\boldsymbol{x}}^2f(\boldsymbol{x}) \in \mathbb{R}^{n \times n}$	Hessian of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at point $\boldsymbol{x} \in \mathbb{R}^n$

1 Preliminaries

1.1 Vector spaces

The *vector space* \mathbb{R}^m consists of all column vectors with m elements. For a set of vectors $\mathcal{C} = \{c_1, \dots, c_n \mid c_i \in \mathbb{R}^m\}$, we can define a subspace spanned by this set, denoted by $\text{span}(\mathcal{C})$. It is the set of all possible linear combinations of elements of \mathcal{C} . If a set of vectors that span a subspace are independent, they are called a *basis*, $b_1, \dots, b_k \in \mathbb{R}^m$. The number of basis vectors defines the *dimensionality* of the subspace.¹

Observation. The following facts hold about subspaces,

- Every subspace contains the zero vector $\mathbf{0}$;
- If x and y are in the subspace, then $x + y$ is also in the subspace;
- If x is in the subspace and $a \in \mathbb{R}$, then ax is also in the subspace.

¹ We know that the amount of basis vectors must be smaller than the amount of vectors that span the subspace, which must be smaller than the dimensionality of the space,

$$k \leq m \leq n.$$

Definition 1.1 (Orthogonal subspaces). Subspaces \mathcal{V} and \mathcal{W} are orthogonal when $v^\top w = 0$ for all $v \in \mathcal{V}, w \in \mathcal{W}$.

1.2 Norms

Definition 1.2 (Inner product). An inner product $\langle \cdot, \cdot \rangle : \mathcal{V} \times \mathcal{V} \rightarrow \mathbb{R}$ is an operation defined on a vector space \mathcal{V} that satisfies the following properties $\forall x, y, z, a, b \in \mathbb{R}$,

- Commutativity: $\langle x, y \rangle = \langle y, x \rangle$;
- Linearity: $\langle x, ay + bz \rangle = a\langle x, y \rangle + b\langle x, z \rangle$;
- Positive definiteness,

$$\begin{aligned} x \neq \mathbf{0} &\implies \langle x, x \rangle > 0 \\ x = \mathbf{0} &\iff \langle x, x \rangle = 0; \end{aligned}$$

- Bilinearity (follows from commutativity and linearity),

$$\begin{aligned} \langle x + y, z \rangle &= \langle x, z \rangle + \langle y, z \rangle \\ \langle x, y + z \rangle &= \langle x, y \rangle + \langle x, z \rangle. \end{aligned}$$

Corollary. $\langle x + y, x + y \rangle = \langle x, x \rangle + 2\langle x, y \rangle + \langle y, y \rangle$.

Corollary. $\langle Ax, y \rangle = \langle x, A^\top y \rangle$.

The vector space \mathcal{V} , along with an inner product, defines an inner vector space. During this course, we will assume that we always work with real vectors in \mathbb{R}^n . An example of an inner product is the dot product,²

² Usually, this operation is what is meant by the inner product.

$$\mathbf{x} \cdot \mathbf{y} = \mathbf{x}^\top \mathbf{y} = \sum_{i=1}^n x_i y_i.$$

Definition 1.3 (Norm). A norm $\|\cdot\| : \mathbb{R}^n \rightarrow \mathbb{R}$ is a function that can be thought of as a way of measuring the distance from the origin. Norms satisfy the following properties,

- Positive definiteness, $\mathbf{x} \neq \mathbf{0} \implies \|\mathbf{x}\| > 0$;
- Triangle inequality, $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$;
- Cauchy-Schwarz inequality, $|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\| \|\mathbf{y}\|$.

Corollary. For the Euclidean norm, the following holds,

$$\cos \theta = \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{x}\| \|\mathbf{y}\|},$$

where θ is the angle between \mathbf{x} and \mathbf{y} .

Corollary (Cosine theorem).

$$\|\mathbf{x} - \mathbf{y}\|^2 = \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 - 2\langle \mathbf{x}, \mathbf{y} \rangle.$$

Each inner product defines a canonical norm $\|\mathbf{x}\| \doteq \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}$. For example, the Euclidean norm is defined by the dot product,

$$\|\mathbf{x}\|_2 = \sqrt{\mathbf{x}^\top \mathbf{x}} = \sqrt{\sum_{i=1}^n x_i^2}.$$

The p -norm is a generalization of the Euclidean norm,

$$\|\mathbf{x}\|_p = \sqrt[p]{\sum_{i=1}^n |x_i|^p}.$$

1.3 Matrices

The rank r of a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ is the dimensionality of its *column space*. It is bounded by

$$r \leq \min\{m, n\}.$$

The matrix is full-rank if $r = \min\{m, n\}$.

A matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ defines 4 fundamental subspaces,

- Column space $\subseteq \mathbb{R}^m$ (r dimensional), $\{\mathbf{b} \mid \mathbf{A}\mathbf{x} = \mathbf{b}\}$;
- Null space $\subseteq \mathbb{R}^n$ ($n - r$ dimensional), $\{\mathbf{x} \mid \mathbf{A}\mathbf{x} = \mathbf{0}\}$;
- Row space $\subseteq \mathbb{R}^n$ (r dimensional), $\{\mathbf{b} \mid \mathbf{A}^\top \mathbf{x} = \mathbf{b}\}$;
- Left null space $\subseteq \mathbb{R}^m$ ($m - r$ dimensional), $\{\mathbf{x} \mid \mathbf{A}^\top \mathbf{x} = \mathbf{0}\}$.

The row space $\text{row}(A)$ is the *orthogonal complement* of the null space $\text{null}(A)$, thus $\text{row}(A) + \text{null}(A) = \mathbb{R}^n$. Similarly, $\text{col}(A) + \text{null}(A^\top) = \mathbb{R}^m$.

Definition 1.4 (Orthonormal matrix). An orthonormal matrix is an invertible matrix whose columns q_1, \dots, q_n are all orthogonal to each other and of unit length, i.e.,

$$\begin{aligned} q_i^\top q_j &= 0, \quad \forall i \neq j \in [n] \\ q_i^\top q_i &= 1, \quad \forall i \in [n]. \end{aligned}$$

Properties. Let $Q \in \mathbb{R}^{n \times n}$ be an orthogonal matrix,

$$\begin{aligned} Q^\top &= Q^{-1} \\ \langle x, y \rangle &= \langle Qx, Qy \rangle. \end{aligned}$$

Definition 1.5 (Trace). The trace of a square matrix $A \in \mathbb{R}^{n \times n}$ is the sum of its diagonal,

$$\text{tr}(A) = \sum_{i=1}^n a_{ii}.$$

Properties. Let $A, B, C \in \mathbb{R}^{n \times n}, x, y \in \mathbb{R}^n, c, d \in \mathbb{R}$, then

$$\begin{aligned} \text{tr}(cA + dB) &= c \cdot \text{tr}(A) + d \cdot \text{tr}(B) \\ \text{tr}(A) &= \text{tr}(A^\top) \\ \text{tr}(ABC) &= \text{tr}(CAB) = \text{tr}(BCA) \\ x^\top y &= \text{tr}(x^\top y) = \text{tr}(xy^\top). \end{aligned}$$

Furthermore, $\text{tr}(A)$ is equal to the sum of the eigenvalues of A .

1.4 Eigenvalues and eigenvectors

The eigenvector $v \in \mathbb{R}^n$ of a matrix $A \in \mathbb{R}^{n \times n}$ has its direction unchanged by its transformation,

$$Av = \lambda v,$$

which is equivalent to

$$(A - \lambda I)v = 0.$$

This matrix must be singular, thus we get the characteristic polynomial,

$$\det(A - \lambda I) = 0.$$

Any λ that satisfies the characteristic polynomial is an eigenvalue of A . Its corresponding eigenvector can then be found by solving the following linear system of equations for v ,

$$(A - \lambda I)v = 0.$$

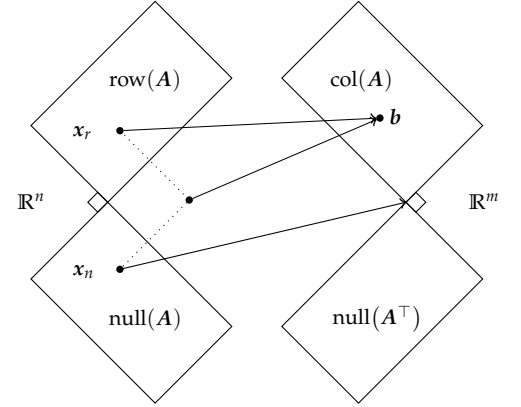


Figure 1.1. Illustration of the 4 spaces defined by a matrix A . It shows the perpendicular spaces. Furthermore, it shows that $Ax_r = b$ for some $x_r \in \text{col}(A)$. Also, if you add a vector from the null space to the row vector, it still maps to the same b , $A(x_r + x_n) = Ax_r + Ax_n = Ax_r = b$.

Linearity.

Cyclic property.

xy^\top is a rank-1 matrix.

The following lemmas are useful for computing the characteristic polynomial of different types of matrices.

Lemma 1.6 (Determinant of 2×2 matrix). Let $A \in \mathbb{R}^{2 \times 2}$ be the following matrix,

$$A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}.$$

Its determinant is computed by

$$\det(A) = ad - bc.$$

Lemma 1.7 (Determinant of triangular matrices). The determinant of a triangular matrix $A \in \mathbb{R}^{n \times n}$ is equal to the product of its diagonal,

$$A = \prod_{i=1}^n a_{ii}.$$

As a consequence the eigenvalues of a triangular matrix are its diagonal entries.

Lemma 1.8 (Determinant of matrix products). Let $A, B \in \mathbb{R}^{n \times n}$, then

$$\det(AB) = \det(A)\det(B).$$

1.5 Convexity

Definition 1.9 (Convexity). A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y), \quad x, y \in \mathbb{R}^n, \lambda \in [0, 1].$$

Lemma 1.10 (Jensen's inequality). Let f be convex, then

$$f(\mathbb{E}[x]) \leq \mathbb{E}[f(x)].$$

2 Dimensionality reduction

The motivation behind dimensionality reduction is to find a low-dimensional representation of high-dimensional data.³ Dimensionality reduction has two goals: (1) compressing the data, while preserving as much relevant information as possible, and (2) interpreting the data, which is easier in low-dimensional space.

Dimensionality reduction is often performed by an autoencoder, which typically has a bottleneck of low dimensionality and aims to predict its input; see Figure 2.1. Let the original data space be d -dimensional and the latent space be k -dimensional with $k \ll d$. An autoencoder consists of an encoder F and a decoder G ,

$$F : \mathbb{R}^d \rightarrow \mathbb{R}^k, \quad G : \mathbb{R}^k \rightarrow \mathbb{R}^d.$$

The idea is that $x \in \mathbb{R}^d$ is mapped to a latent vector $z \in \mathbb{R}^k$ by the encoder, which is mapped to a reconstruction $\hat{x} \in \mathbb{R}^d$ by the decoder. The idea is that the encoder must compress the information well for the decoder to be able to reconstruct its input. The reconstruction function is then the following function,

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

which aims to resemble the identity function $(G \circ F)(x) = x$. Generally, this is only possible if the data is intrinsically k -dimensional.

2.1 Linear autoencoders

In order to build a nice theory, we will only consider a single layer linear autoencoder.⁴ As a result, we have the following functions,

$$\begin{aligned} F(x) &= Wx, \quad W \in \mathbb{R}^{k \times d} \\ G(z) &= Vz, \quad V \in \mathbb{R}^{d \times k} \\ (G \circ F)(x) &= VWx. \end{aligned}$$

The objective to minimize is the following,

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

Intuitively, it is the difference between the identity function and $G \circ F$.

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

Proof. Let $F(x) = Wx + a$ and $G(z) = Vz + b$, then an affine model's function is computed by

$$(G \circ F)(x) = VWx + c, \quad c = Va + b.$$

³Often, the original raw representation is high-dimensional and redundant, e.g., images, audio, time series.



Figure 2.1. Diagram of a single layer linear autoencoder.

⁴Considering non-linear parametrizations will result in a much more powerful autoencoder.

Proof by contradiction. Assume there exists a $c \neq 0$ such that $\mathbb{E}\|x - (Px + c)\|^2 \leq \mathbb{E}\|x - Px\|^2$. Let $c \neq 0$, then

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

Cosine theorem.

Linearity of expectation.

Centered data: $\mathbb{E}[x] = 0$

$c \neq 0$.

This contradicts the assumption. ■

Thus, we will assume that the data is centered, which makes the analysis easier, since we do not need to consider the affine case. This is a reasonable assumption, because data can always be centered by subtracting $\mathbb{E}[x]$ from all data points.

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

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

with $A^{-1}W, VA$.

The weight matrices V, W are *non-identifiable*. As a result, we must not over-interpret the found representation.

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

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

P is constrained to be, at most, a rank- k matrix.

2.2 Projection

The rank constraint and linearity of P means that the image (column space) of P is a linear subspace $\mathcal{U} \subseteq \mathbb{R}^d$ of dimension at most k . 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.⁵

⁵ We do not search for the weight matrices W, V , since they are not unique, but P is unique.

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

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

Definition 2.1 (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. *Idempotency*: $P(P(x)) = P(x)$.
3. *Orthogonality*: $\text{null}(P) \perp \text{col}(P)$, which is equivalent to the column space and row space of P being equal. It is also equivalent to self-adjointness holding, $\langle P(x), y \rangle = \langle x, P(y) \rangle$.

Lemma 2.2. $\Pi_{\mathcal{U}}$ defined as

$$\Pi_{\mathcal{U}}(x) = \underset{x' \in \mathcal{U}}{\operatorname{argmin}} \|x - x'\|^2$$

is an orthogonal projection of \mathbb{R}^n onto \mathcal{U} .



Figure 2.2. 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 2.1 (linearity, projection, orthogonality, and idempotency).

1. *Linearity*: A function f is linear if it satisfies homogeneity and additivity,

$$f(\alpha x) = \alpha f(x), \quad f(x + y) = f(x) + f(y).$$

This can easily be shown for $\Pi_{\mathcal{U}}$,

$$\begin{aligned} \Pi_{\mathcal{U}}(\alpha x) &= \underset{x'' \in \mathcal{U}}{\operatorname{argmin}} \|\alpha x - x''\|^2 \\ &= \underset{\alpha x' \in \mathcal{U}}{\operatorname{argmin}} \|\alpha x - \alpha x'\|^2 \\ &= \alpha \underset{x' \in \mathcal{U}}{\operatorname{argmin}} \|\alpha\| \|x - x'\|^2 \\ &= \alpha \underset{x' \in \mathcal{U}}{\operatorname{argmin}} \|x - x'\|^2 \\ &= \alpha \Pi_{\mathcal{U}}(x). \\ \Pi_{\mathcal{U}}(x + y) &= \underset{z \in \mathcal{U}}{\operatorname{argmin}} \|z - (x + y)\|^2 \\ &= \underset{\substack{\Pi_{\mathcal{U}}(x) + y' \\ y' \in \mathcal{U}}}{\operatorname{argmin}} \|\Pi_{\mathcal{U}}(x) + y' - x - y\|^2 \\ &= \Pi_{\mathcal{U}}(x) + \underset{y' \in \mathcal{U}}{\operatorname{argmin}} \|(y' - y) + (\Pi_{\mathcal{U}}(x) - x)\|^2 \\ &= \Pi_{\mathcal{U}}(x) + \underset{y' \in \mathcal{U}}{\operatorname{argmin}} \|y' - y\|^2 + \|\Pi_{\mathcal{U}}(x) - x\|^2 \\ &\quad - 2\langle \Pi_{\mathcal{U}}(x) - x, y' - y \rangle \\ &= \Pi_{\mathcal{U}}(x) + \underset{y' \in \mathcal{U}}{\operatorname{argmin}} \|y' - y\|^2 - 2\langle \Pi_{\mathcal{U}}(x) - x, y' \rangle \\ &= \Pi_{\mathcal{U}}(x) + \Pi_{\mathcal{U}}(y); \end{aligned}$$

$$\operatorname{argmin}_{\lambda x} f(x) = \lambda \operatorname{argmin}_x f(x).$$

$$z \doteq \Pi_{\mathcal{U}}(x) + y' \text{ for some } y' \in \mathcal{U}.$$

$$\operatorname{argmin}_{y+x} f(x) = y + \operatorname{argmin}_x f(x).$$

Cosine theorem.

$\|\Pi_{\mathcal{U}}(x) - x\|^2$ and $\langle \Pi_{\mathcal{U}}(x) - x, y \rangle$ do not depend on y' .

$\Pi_{\mathcal{U}}(x) - x \in \mathcal{U}^\perp$ and $y' \in \mathcal{U}$, so their inner product is 0.

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

$$\|\Pi_{\mathcal{U}}(\mathbf{x}) - \mathbf{x}\|^2 = \min_{\mathbf{u} \in \mathcal{U}} \|\mathbf{u} - \mathbf{x}\|^2 \leq \|\tilde{\mathbf{u}} - \mathbf{x}\|^2, \quad \forall \tilde{\mathbf{u}} \in \mathcal{U}. \quad (1)$$

Decompose $\Pi_{\mathcal{U}}(\mathbf{x}) - \mathbf{x} = \mathbf{u} + \mathbf{u}^{\perp}$, where $\mathbf{u} \in \mathcal{U}$ and $\mathbf{u}^{\perp} \in \mathcal{U}^{\perp}$. Then, we only need to show $\mathbf{u} = \mathbf{0}$, which we will prove by contradiction. Let $\mathbf{u} \neq \mathbf{0}$, then

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

Cosine theorem.

$\mathbf{u} \neq \mathbf{0}$ and $\langle \mathbf{u}, \mathbf{u}^{\perp} \rangle = 0$.

This contradicts with Equation (1), because $\Pi_{\mathcal{U}}(\mathbf{x}) - \mathbf{u} \in \mathcal{U}$.

■

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

Lemma 2.3. Given an orthonormal basis \mathbf{U} of \mathcal{U} ($\mathbf{u}_i \in \mathcal{U}, \|\mathbf{u}_i\| = 1, \langle \mathbf{u}_i, \mathbf{u}_j \rangle = 0, \forall i \neq j$), we can compute the optimal projection matrix,

$$\mathbf{P} = \mathbf{U}\mathbf{U}^{\top}.$$

Note that in this case \mathbf{W} and \mathbf{V} share parameters, and $\mathbf{U}\mathbf{U}^{\top}$ is the optimal weight matrix if we enforce parameter sharing via $\mathbf{V} = \mathbf{W}^{\top}$.

Proof. We need to show that \mathbf{P} is the orthogonal projection matrix onto \mathcal{U} :

1. *Projection*: $\mathbf{P}\mathbf{x} = \sum_{i=1}^k \mathbf{u}_i \mathbf{u}_i^{\top} \mathbf{x} = \sum_{i=1}^k \langle \mathbf{u}_i, \mathbf{x} \rangle \mathbf{u}_i \in \mathcal{U}$;
2. *Self-adjointness*: $\mathbf{P}^{\top} = (\mathbf{U}\mathbf{U}^{\top})^{\top} = \mathbf{U}\mathbf{U}^{\top} = \mathbf{P}$;
3. *Idempotency*: $\mathbf{P}\mathbf{P} = \mathbf{U}\mathbf{U}^{\top} \mathbf{U}\mathbf{U}^{\top} = \mathbf{U}\mathbf{U}^{\top} = \mathbf{P}$;
4. *Orthogonality*: $\langle \mathbf{P}\mathbf{x}, \mathbf{x} - \mathbf{P}\mathbf{x} \rangle = \langle \mathbf{x}, \mathbf{P}(\mathbf{x} - \mathbf{P}\mathbf{x}) \rangle = \langle \mathbf{x}, (\mathbf{P} - \mathbf{P}^2)\mathbf{x} \rangle = \langle \mathbf{x}, (\mathbf{P} - \mathbf{P})\mathbf{x} \rangle = 0$.

■

The problem is that we do not have an orthonormal basis of \mathcal{U} in general. The next lemma gives a more general result for any basis.

Lemma 2.4. For a non-orthonormal basis $V \in \mathbb{R}^{d \times k}$ of \mathcal{U} , we can recover the projection matrix,

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

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

Proof. We need to show that P is the orthogonal projection matrix of \mathcal{U} :

1. *Projection:* $PV = VV^+V = V(V^\top V)^{-1}V^\top V = V$. Together with the rank constraint, this yields $Pu^\perp = \mathbf{0}$ for all $u^\perp \in \mathcal{U}^\perp$;
2. *Self-adjointness:* $P^\top = (V(V^\top V)^{-1}V^\top)^\top = V(V^\top V)^{-1}V^\top = P$;
3. *Idempotency:* $PP = V(V^\top V)^{-1}V^\top V(V^\top V)^{-1}V^\top = VV^+ = P$.

■

Note that we have not made any reference to the data $x \sim \nu$ until now. This means that all of the above results are a priori by the structure of the model and choice of objective.

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

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

Cosine theorem.

Linearity of expectation and idempotency.

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

Because our data is centered, we know the following,

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

Assuming centered data,

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

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

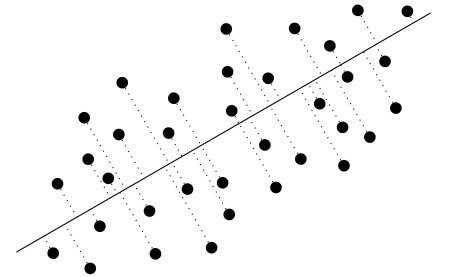


Figure 2.3. Intuitively, the projection onto the shown line preserves the most information. This projection from \mathbb{R}^2 onto \mathbb{R} maximizes variance.

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

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

$$\|\mathbf{P}\mathbf{x}\|^2 = \langle \mathbf{P}\mathbf{x}, \mathbf{P}\mathbf{x} \rangle = \langle \mathbf{x}, \mathbf{P}^2\mathbf{x} \rangle = \langle \mathbf{x}, \mathbf{P}\mathbf{x} \rangle.$$

Trace of scalar is equal to scalar.

Cyclic property of trace.

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

2.3 Principal component analysis

Intuitively, eigenvectors represent axes along which matrices have the largest variance. Thus, we want to project onto the k eigenvectors of $\mathbb{E}[\mathbf{x}\mathbf{x}^\top]$ that are associated with the k largest eigenvalues.

Theorem 2.5 (Spectral theorem). Any symmetric and positive semidefinite matrix Σ can be non-negatively diagonalized with an orthogonal matrix,

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

where $\lambda_1 \geq \dots \geq \lambda_d \geq 0$ and $\mathbf{Q} \in \mathbb{R}^d$ is orthogonal.

The spectral theorem gives us the eigendecomposition of $\mathbb{E}[\mathbf{x}\mathbf{x}^\top]$, because it is symmetric. The eigenvectors form an orthonormal basis, so the k principal eigenvectors are also orthonormal. Hence, we can use Lemma 2.3 to form the projection matrix, which is essentially what the PCA theorem tells us.

Theorem 2.6 (PCA theorem). The variance maximizing projection matrix \mathbf{P} for a covariance matrix $\mathbb{E}[\mathbf{x}\mathbf{x}^\top] = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top$ as in the spectral theorem is given by

$$\mathbf{P} = \mathbf{u}_k \mathbf{u}_k^\top, \quad \mathbf{u}_k = \mathbf{Q} \begin{bmatrix} \mathbf{I}_k \\ \mathbf{0} \end{bmatrix},$$

which means that \mathbf{u}_k consists of the k principal eigenvectors.

Proof.

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

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

Cyclic property.

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

In conclusion, given a dataset of n points $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, we perform dimensionality reduction by first centering the data,

$$\begin{aligned}\boldsymbol{\mu} &= \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \\ \bar{\mathbf{x}}_i &= \mathbf{x}_i - \boldsymbol{\mu}.\end{aligned}$$

Then, we compute the covariance matrix $\boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}$,

$$\boldsymbol{\Sigma} = \frac{1}{n} \sum_{i=1}^n \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^\top.$$

For that matrix, we compute the eigendecomposition (which exists because $\boldsymbol{\Sigma}$ is symmetric),

$$\boldsymbol{\Sigma} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top,$$

where $\mathbf{Q} \in \mathbb{R}^{d \times d}$ is orthogonal and $\mathbf{\Lambda}$ is diagonal. We then discard the $d - k$ last dimensions to obtain $\mathbf{U}_k \in \mathbb{R}^{d \times k}$,

$$\mathbf{U}_k = \mathbf{Q} \begin{bmatrix} \mathbf{I}_k \\ \mathbf{0} \end{bmatrix}.$$

The latent vectors are then computed by

$$\mathbf{z}_i = \mathbf{U}_k^\top \bar{\mathbf{x}}_i.$$

Their reconstructions are computed by

$$\hat{\mathbf{x}}_i = \mathbf{U}_k \mathbf{z}_i.$$

The squared reconstruction error is equal to the sum of the lower $d - k$ eigenvalues.

2.4 Learning algorithms

Eigenvalue decomposition of the (symmetric) sample covariance matrix has $\mathcal{O}(d^3)$ complexity. Furthermore, the complexity of computing $\mathbb{E}[\mathbf{x}\mathbf{x}^\top]$ is $\mathcal{O}(nd^2)$, where n is the number of data points.⁶ This is quite costly, thus we need to search for algorithms that have lower runtime complexity.

⁶ Typically, $n \gg d$.

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

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

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

Lemma 2.7. 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,

$$\mathbf{v}^{(k)} \propto \mathbf{A}^k \mathbf{v}^{(0)}$$

$$= \sum_{i=1}^d \alpha_i \lambda_i^k \mathbf{u}_i$$

$$\propto \alpha_1 \mathbf{u}_1 + \sum_{i=2}^d \alpha_i \left(\frac{\lambda_i}{\lambda_1} \right)^k \mathbf{u}_i.$$

$$\mathbf{A}^k \mathbf{v}^{(0)} = \sum_{i=1}^d \alpha_i \mathbf{A}^k \mathbf{u}_i = \sum_{i=1}^d \alpha_i \lambda_i^k \mathbf{u}_i.$$

Divide by λ_1^k .

$\lambda_i/\lambda_1 < 1$ 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 by $\mathbf{A}_2 = \mathbf{A} - \lambda_1 \mathbf{u}_1 \mathbf{u}_1^\top$ and then applying the algorithm again to \mathbf{A}_2 to recover \mathbf{u}_2 , and continue doing that until we have the k principal eigenvectors.

Thus, the total complexity of finding the k principal eigenvectors is $\mathcal{O}(Tkd^2)$ instead of $\mathcal{O}(d^3)$. However, this does not get rid of the $\mathcal{O}(nd^2)$ complexity for computing the sample covariance matrix, which is a bigger problem.

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,

$$\nabla_{\mathbf{V}} \mathcal{R}(\mathbf{W}, \mathbf{V}) = (\mathbf{V}\mathbf{W} - \mathbf{I})\mathbf{x}\mathbf{x}^\top \mathbf{W}^\top$$

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

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

3 Matrix completion

The goal of matrix completion is to fill in the missing entries of a sparse matrix. E.g., Netflix might want to use matrix completion to predict which users will give which ratings to which movies, given the ratings given by all users of Netflix. However, not all users have rated all movies, thus the matrix is sparse and must thus be filled in. Netflix can use this prediction to decide which movies they should recommend to their users.

If we assume that all entries are independent, then no information is carried by one entry about another. Because of this, we cannot reconstruct the matrix. Thus, we need to make an assumption about the dependency within the matrix.

A minimal assumption that we make is that entries within the same row or the same column are not independent, *i.e.*,

$$a_{ij} \perp\!\!\!\perp \{a_{kl} \mid k \neq i \wedge l \neq j\} \mid \{a_{il} \mid l \neq j\} \cup \{a_{kj} \mid k \neq i\}.$$

This states that a_{ij} is independent from all entries not on the same row or column, given that we know all entries on the same row or column. However, in reality, we do not have all values on the row and column for all entries. Thus, effectively, we have an indirect coupling between all entries, since a_{kl} influences a_{kj} and a_{il} (if they are unknown), which influence a_{ij} .

Formally, we have an underlying rating matrix $A \in \mathbb{R}^{n \times m}$ and an observation matrix $\Omega \in \{0, 1\}^{n \times m}$, where $\omega_{ij} = 1$ means that a_{ij} is observed. The observed matrix will be denoted by \tilde{A} to differentiate between the underlying matrix and its observations. The goal is to predict all values a_{kl} , where $\omega_{kl} = 0$.

To approximate the underlying matrix $A \in \mathbb{R}^{n \times m}$, we will make the assumption that it is a rank- k matrix. This is enforced by approximating A as the product of two matrices $U \in \mathbb{R}^{n \times k}$ and $V \in \mathbb{R}^{m \times k}$,

$$A \approx UV^\top.$$

Then, we get the following loss function that we wish to minimize,

$$U, V \in \operatorname{argmin}_{U, V} \ell(U, V) \doteq \frac{1}{2} \left\| \Pi_\Omega(\tilde{A} - UV^\top) \right\|_F^2,$$

where $\Pi_\Omega(M) = M \odot \Omega$ and $\|\cdot\|_F$ is the Frobenius norm. Rewriting as sums, we get the following loss function,

$$\ell(U, V) = \sum_{i=1}^n \sum_{j=1}^m \omega_{ij} \left(\tilde{a}_{ij} - u_i^\top v_j \right)^2.$$

A possible issue with the values is that some users might generally give higher ratings than others. Thus, we want to account for this by variance normalization. We can do this per row (user) or per column (item). (While it makes more sense intuitively that we should normalize ratings per user, it has been shown to be more effective to normalize over item ratings.) Normalization yields a vector with zero mean and unit variance,

$$Z = \frac{X - \mu}{\sigma}.$$

We view $\mathbf{u}_i \in \mathbb{R}^k$ as representing each user and $\mathbf{v}_j \in \mathbb{R}^k$ as representing each item. We estimate a_{ij} to be their inner product,

$$a_{ij} \approx \mathbf{u}_i^\top \mathbf{v}_j.$$

Furthermore, we have that \mathbf{U} and \mathbf{V} are fully determined by the observed values \tilde{A} , but can be used to extrapolate to a full matrix A , achieving our objective. Moreover, we have that all rows and columns are coupled, which satisfies our dependence assumption.

3.1 Fully observed case

In this section, we consider the fully observed case, where $\Omega = \mathbf{1}_{n \times m}$. Thus, we optimize w.r.t. the following loss function,

$$\ell(\mathbf{U}, \mathbf{V}) = \frac{1}{2} \|\mathbf{A} - \mathbf{UV}^\top\|_F^2.$$

Scalar case. To observe the properties of the gradients, we will first consider the 1-dimensional case, where $a \approx uv$, with $u, v \in \mathbb{R}$. We have the following loss function,

$$\ell(u, v) = \frac{1}{2}(a - uv)^2$$

with the following derivatives,

$$\frac{\partial}{\partial u} \ell(u, v) = v(uv - a), \quad \frac{\partial}{\partial v} \ell(u, v) = u(uv - a),$$

which induces the negative gradient field shown in Figure 3.1. Furthermore, we have the following Hessian,

$$\nabla^2 \ell(u, v) = \begin{bmatrix} v^2 & 2uv - a \\ 2uv - a & u^2 \end{bmatrix}.$$

Lemma 3.1 (Second-order characterization of convexity). If $f : \mathcal{X} \rightarrow \mathbb{R}$ is twice differentiable, then f is convex if and only if

$$\nabla^2 f(\mathbf{x}) \succeq \mathbf{0}, \quad \forall \mathbf{x} \in \text{Int}(\mathcal{X}).$$

At the origin, we have the following Hessian,

$$\nabla^2 \ell(0, 0) = \begin{bmatrix} 0 & -a \\ -a & 0 \end{bmatrix}.$$

In the scalar case, the objective function is non-convex for $a \neq 0$. The same result can be generalized to any $n, m \geq 1$.

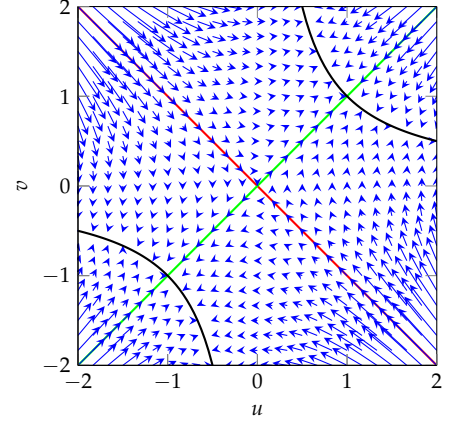


Figure 3.1. Negative gradient field for $a = 1$ in the scalar case with minima indicated by black. As can be seen, $[0, 0]$ is a 2-way saddle point and any vector $[-z, z]$ for $z \in \mathbb{R}$ moves toward it. We can start at any other point and use gradient-based optimization to converge to the minimum.

As a result, we do not have any general convergence guarantees using gradient-based methods. However, we can make guarantees for this specific problem if we study the gradient flow of the loss function with ordinary differential equations (ODE). Consider a balanced initialization $u_0 = v_0$. We see that u and v will evolve identically, because their partial derivatives are equal in this case. We have the following update rules,

$$\begin{aligned} u_{t+1} &= u_t - \eta \frac{\partial}{\partial u} \ell(u, v) = u_t - \eta v(uv - a) \\ v_{t+1} &= v_t - \eta \frac{\partial}{\partial v} \ell(u, v) = v_t - \eta u(uv - a), \end{aligned}$$

where η is an arbitrarily small stepsize. We then have the following ODEs of u and v w.r.t. the step t ,

$$\frac{du}{dt} = \frac{u_{t+1} - u_t}{\eta} = -v(uv - a), \quad \frac{dv}{dt} = \frac{v_{t+1} - v_t}{\eta} = -u(uv - a).$$

Consider $x = uv = u^2 = v^2$ (because of balanced initialization and identical evolution), then we have the following ODE,

$$\frac{dx}{dt} = u \frac{dv}{dt} + v \frac{du}{dt} = -u^2(uv - a) - v^2(uv - a) = -2x(x - a).$$

Solving this ODE yields the following solution for $x(t)$,

$$x(t) = a + \frac{ac - a^2}{ce^{2at} + a - c} \xrightarrow{t \rightarrow \infty} a.$$

Starting from a balanced initialization, gradient descent with a small enough $\eta > 0$ will converge uv to a .

Rank-1 model. Now consider the fully observed rank-1 model,

$$\ell(\mathbf{u}, \mathbf{v}) = \frac{1}{2} \|\mathbf{A} - \mathbf{u}\mathbf{v}^\top\|_F^2, \quad \mathbf{A} \in \mathbb{R}^{n \times m}, \mathbf{u} \in \mathbb{R}^n, \mathbf{v} \in \mathbb{R}^m.$$

We can rewrite this to a more useful form,

$$\begin{aligned} &= \frac{1}{2} \text{tr} \left((\mathbf{A} - \mathbf{u}\mathbf{v}^\top)^\top (\mathbf{A} - \mathbf{u}\mathbf{v}^\top) \right) & \|\mathbf{M}\|_F^2 &= \text{tr}(\mathbf{M}^\top \mathbf{M}). \\ &= \frac{1}{2} \text{tr} \left(\mathbf{A}^\top \mathbf{A} - \mathbf{v}\mathbf{u}^\top \mathbf{A} - \mathbf{A}^\top \mathbf{u}\mathbf{v}^\top + \mathbf{v}\mathbf{u}^\top \mathbf{u}\mathbf{v}^\top \right) \\ &= \frac{1}{2} \left(\text{tr}(\mathbf{A}^\top \mathbf{A}) - \text{tr}(\mathbf{v}\mathbf{u}^\top \mathbf{A}) - \text{tr}(\mathbf{A}^\top \mathbf{u}\mathbf{v}^\top) + \text{tr}(\mathbf{v}\mathbf{u}^\top \mathbf{u}\mathbf{v}^\top) \right) & \text{Linearity of trace.} \\ &= \frac{1}{2} \left(\text{tr}(\mathbf{A}^\top \mathbf{A}) - \text{tr}(\mathbf{u}^\top \mathbf{A}\mathbf{v}) - \text{tr}(\mathbf{v}^\top \mathbf{A}^\top \mathbf{u}) + \text{tr}(\mathbf{v}^\top \mathbf{v}\mathbf{u}^\top \mathbf{u}) \right) & \text{Cyclic property of trace.} \\ &= \frac{1}{2} \left(\text{tr}(\mathbf{A}^\top \mathbf{A}) - 2\text{tr}(\mathbf{u}^\top \mathbf{A}\mathbf{v}) + \|\mathbf{u}\|^2 \|\mathbf{v}\|^2 \right) & \text{tr}(\mathbf{A}) = \text{tr}(\mathbf{A}^\top). \\ &\propto \frac{1}{2} \|\mathbf{u}\|^2 \|\mathbf{v}\|^2 - \mathbf{u}^\top \mathbf{A}\mathbf{v}. \end{aligned}$$

We observe that the direction of \mathbf{u} and \mathbf{v} is fully decided by the second term. Thus, if we first solve for that, we can then easily solve for the norm of the two vectors by also considering the first term.

Specifically, let $\mathbf{u} = c_1 \tilde{\mathbf{u}}$ and $\mathbf{v} = c_2 \tilde{\mathbf{v}}$, such that $\tilde{\mathbf{u}}$ and $\tilde{\mathbf{v}}$ are unit vectors, and $c \doteq c_1 \cdot c_2$. Then, we have the following loss function with an additional parameter c ,

$$\ell(\tilde{\mathbf{u}}, \tilde{\mathbf{v}}, c) = \frac{1}{2}c^2 - c\tilde{\mathbf{u}}^\top \mathbf{A} \tilde{\mathbf{v}}, \quad \|\tilde{\mathbf{u}}\| = \|\tilde{\mathbf{v}}\| = 1.$$

Hence, we first maximize $\mathbf{u}^\top \mathbf{A} \mathbf{v}$ under the constraint that \mathbf{u} and \mathbf{v} are unit vectors, and then decide which c minimizes the loss function. Thus, first we have the following problem

$$\mathbf{u}, \mathbf{v} \in \operatorname{argmax}_{\|\mathbf{u}\|=\|\mathbf{v}\|=1} \mathbf{u}^\top \mathbf{A} \mathbf{v}.$$

We can solve this problem with its constraints by using Lagrangian multipliers,

$$\mathcal{L}(\mathbf{u}, \mathbf{v}, \alpha, \beta) = \mathbf{u}^\top \mathbf{A} \mathbf{v} - \alpha (\|\mathbf{u}\|^2 - 1) - \beta (\|\mathbf{v}\|^2 - 1).$$

Using the first-order optimality condition, we get the following solutions,

$$\begin{aligned} \nabla_{\mathbf{u}} \mathcal{L}(\mathbf{u}, \mathbf{v}, \alpha, \beta) &= \mathbf{A} \mathbf{v} - 2\alpha \mathbf{u} \stackrel{!}{=} \mathbf{0} \\ \nabla_{\mathbf{v}} \mathcal{L}(\mathbf{u}, \mathbf{v}, \alpha, \beta) &= \mathbf{A}^\top \mathbf{u} - 2\beta \mathbf{v} \stackrel{!}{=} \mathbf{0}. \end{aligned}$$

Choosing α and β to satisfy the unit vector constraints, we get the following solutions,

$$\mathbf{u} = \frac{\mathbf{A} \mathbf{v}}{\|\mathbf{A} \mathbf{v}\|}, \quad \mathbf{v} = \frac{\mathbf{A}^\top \mathbf{u}}{\|\mathbf{A}^\top \mathbf{u}\|}.$$

Thus, the solutions must satisfy

$$\mathbf{u} \propto (\mathbf{A} \mathbf{A}^\top) \mathbf{u}, \quad \mathbf{v} \propto (\mathbf{A}^\top \mathbf{A}) \mathbf{v}.$$

As a result, we know that \mathbf{u} is an eigenvector of $\mathbf{A} \mathbf{A}^\top$ and \mathbf{v} is an eigenvector of $\mathbf{A}^\top \mathbf{A}$. Using this fact, we can rewrite the intermediate objective as

$$\begin{aligned} \mathbf{u}^\top \mathbf{A} \mathbf{v} &= \frac{\mathbf{u}^\top \mathbf{A} \mathbf{A}^\top \mathbf{u}}{\|\mathbf{A}^\top \mathbf{u}\|} = \frac{\mathbf{u}^\top \lambda \mathbf{u}}{\sqrt{\mathbf{u}^\top \mathbf{A} \mathbf{A}^\top \mathbf{u}}} \\ &= \frac{\mathbf{u}^\top \lambda \mathbf{u}}{\sqrt{\mathbf{u}^\top \lambda \mathbf{u}}} = \frac{\lambda \|\mathbf{u}\|^2}{\sqrt{\lambda} \|\mathbf{u}\|} = \sqrt{\lambda}. \end{aligned}$$

In order to maximize this term, we must thus select \mathbf{u} and \mathbf{v} to be the principal eigenvectors of their respective matrices. Furthermore, $\ell(\tilde{\mathbf{u}}, \tilde{\mathbf{v}}, c)$ with this selection of $\tilde{\mathbf{u}}$ and $\tilde{\mathbf{v}}$ is minimized by selecting $c = \sqrt{\lambda_1}$, where λ_1 is the principal eigenvalue of $\mathbf{A} \mathbf{A}^\top$ and $\mathbf{A}^\top \mathbf{A}$.

The loss for a rank-1 model is minimized by selecting

$$A = \sqrt{\lambda_1} \mathbf{u} \mathbf{v}^\top,$$

where \mathbf{u} and \mathbf{v} are the principal eigenvectors of AA^\top and $A^\top A$, respectively, and λ_1 is the principal eigenvalue of AA^\top and $A^\top A$. We can compute these vectors by the power iteration algorithm.

Rank-k model. We can generalize the previous result by making use of the singular value decomposition (SVD).

Theorem 3.2 (SVD theorem). For each matrix $A \in \mathbb{R}^{n \times m}$, there exists a diagonal matrix $\Sigma \in \mathbb{R}^{n \times m}$ with ordered entries $\sigma_i \geq \sigma_{i+1} \geq 0, \forall i \in \min\{n, m\}$ and orthogonal matrices $\mathbf{U} \in \mathbb{R}^{n \times n}, \mathbf{V} \in \mathbb{R}^{m \times m}$ such that A can be expressed as

$$A = \mathbf{U} \Sigma \mathbf{V}^\top.$$

The SVD of A is intimately related to the eigendecomposition of AA^\top and $A^\top A$,

$$\begin{aligned} AA^\top &= \mathbf{U} \Sigma \Sigma^\top \mathbf{U}^\top \\ A^\top A &= \mathbf{V}^\top \Sigma^\top \Sigma \mathbf{V}. \end{aligned}$$

Thus, \mathbf{U} and \mathbf{V}^\top are the eigenvectors of AA^\top and $A^\top A$, respectively. Furthermore, σ_i^2 is the i -th eigenvalue of both.

Lemma 3.3. Let the SVD of $A \in \mathbb{R}^{n \times m}$ be given by $A = \mathbf{U} \Sigma \mathbf{V}^\top$, then

$$\|A\|_F^2 = \sum_{i=1}^{\min\{n, m\}} \sigma_i^2.$$

Proof. This can be shown by the properties of the trace operator,

$$\begin{aligned} \|A\|_F^2 &= \text{tr}(A^\top A) = \text{tr}(\mathbf{V} \Sigma^\top \mathbf{U}^\top \mathbf{U} \Sigma \mathbf{V}^\top) = \text{tr}(\mathbf{V} \Sigma^\top \Sigma \mathbf{V}^\top) \\ &= \text{tr}(\mathbf{V}^\top \mathbf{V} \Sigma^\top \Sigma) = \text{tr}(\Sigma^\top \Sigma) = \sum_{i=1}^{\min\{n, m\}} \sigma_i^2. \end{aligned}$$

■

Lemma 3.4. Let the SVD of $A \in \mathbb{R}^{n \times m}$ be given by $A = \mathbf{U} \Sigma \mathbf{V}^\top$, then

$$\|A\|_2 \doteq \sup_{\|x\|=1} \|Ax\| = \sigma_1.$$

Proof. We use the fact that orthogonal matrices preserve the Euclidean norm,

$$\begin{aligned}\sup_{\|x\|=1} \|Ax\| &= \sup_{\|x\|=1} \|\mathbf{U}\Sigma\mathbf{V}^\top x\| = \sup_{\|x\|=1} \|\Sigma\mathbf{V}^\top x\| \\ &= \sup_{\|z\|=1} \|\Sigma z\| = \|\Sigma\|_2 = \sigma_1.\end{aligned}$$

■

The Eckart-Young theorem is widely used for linear matrix approximation. It states that if we prune the singular values below σ_k in the SVD representation, we get an optimal rank- k approximation of a matrix.

Theorem 3.5 (Eckart-Young theorem). Let the SVD of $A \in \mathbb{R}^{n \times m}$ be given by $A = \mathbf{U}\Sigma\mathbf{V}^\top$. Then, for all $1 \leq k \leq \min\{n, m\}$, we have

$$A_k \doteq \mathbf{U} \text{diag}([\sigma_1, \dots, \sigma_k]) \mathbf{V}^\top \in \underset{\text{rank}(B) \leq k}{\text{argmin}} \|A - B\|_F.$$

From this, we can easily find the optimal rank- k approximation of A by

$$A \approx A_k = \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^\top.$$

Furthermore, we can directly compute the squared error of low rank approximations as

$$\|A - A_k\|_F^2 = \sum_{i=k+1}^{\text{rank}(A)} \sigma_i^2.$$

The SVD is computable in $\mathcal{O}(\min\{nm^2, mn^2\})$, which is fast for a non-convex problem. However, we cannot use this technique in the case of incomplete observations.

3.2 Incompletely observed case

In practice, the matrix is often incompletely observed, meaning that we can only compute the loss function w.r.t. the observed entries. This task is called “Matrix completion”, and we cannot use SVD in this case. Even worse, this problem is NP-hard, thus we have to resort to approximation algorithms.

Alternating least squares. As we saw, we parametrize the approximation of A as a factorization of two matrices $\mathbf{U} \in \mathbb{R}^{n \times k}$ and $\mathbf{V} \in \mathbb{R}^{m \times k}$,

$$A \approx \mathbf{U}\mathbf{V}^\top.$$

Since we cannot solve the problem exactly, we want to add regularization to the loss function to increase numerical stability of the solution,⁷

⁷ This will also be very important later.

$$\ell(\mathbf{U}, \mathbf{V}) = \frac{1}{2} \|\Pi_\Omega(\tilde{A} - \mathbf{U}\mathbf{V}^\top)\|_F^2 + \frac{\lambda}{2} (\|\mathbf{U}\|_F^2 + \|\mathbf{V}\|_F^2), \quad \lambda > 0.$$

If we fully expand the norms, we will find that this objective is a 4-th degree polynomial in the parameters (entries of \mathbf{U} and \mathbf{V}) with the following monomials,

$$\begin{aligned}\omega_{ij}u_{ir}v_{jr}u_{is}v_{js}, & \quad 1 \leq r, s \leq k \\ \omega_{ij}u_{ir}v_{jr}, & \quad 1 \leq r \leq k \\ u_{ir}^2, & \quad 1 \leq r \leq k \\ v_{jr}^2, & \quad 1 \leq r \leq k.\end{aligned}$$

Importantly, every 4-th order term involves exactly one row index i of \mathbf{U} and one row index j of \mathbf{V} . In other words, the objective w.r.t. u_i depends further only on \mathbf{V} and not any other row of \mathbf{U} .⁸ Using this information, we can separate out part of the objective function depending on a row v_j , where we treat \mathbf{U} as fixed,

⁸ The parameter dependencies form a bipartite graph between the rows of \mathbf{U} and \mathbf{V} .

$$\begin{aligned}\ell_{\mathbf{U}}(v_j) &= \frac{1}{2} \left\| \Pi_{\Omega}(\tilde{\mathbf{A}} - \mathbf{U}\mathbf{V}^{\top}) \right\|_F^2 + \frac{\lambda}{2} (\|\mathbf{U}\|_F^2 + \|\mathbf{V}\|_F^2) \\ &\propto \frac{1}{2} \sum_{i=1}^n \sum_{j'=1}^m \omega_{ij'} (a_{ij'} - \mathbf{u}_i^{\top} \mathbf{v}_{j'})^2 + \frac{\lambda}{2} \sum_{j'=1}^m \sum_{r=1}^k v_{j'r}^2 \\ &\propto \frac{1}{2} \sum_{i=1}^n \omega_{ij} (a_{ij} - \mathbf{u}_i^{\top} \mathbf{v}_j)^2 + \frac{\lambda}{2} \sum_{r=1}^k v_{jr}^2 \\ &= \frac{1}{2} \sum_{i=1}^n \omega_{ij} a_{ij}^2 + \omega_{ij} \mathbf{u}_i^{\top} \mathbf{v}_j \mathbf{u}_i^{\top} \mathbf{v}_j - 2\omega_{ij} a_{ij} \mathbf{u}_i^{\top} \mathbf{v}_j + \frac{\lambda}{2} \mathbf{v}_j^{\top} \mathbf{v}_j \\ &\propto \frac{1}{2} \sum_{i=1}^n \omega_{ij} \mathbf{v}_j^{\top} \mathbf{u}_i \mathbf{u}_i^{\top} \mathbf{v}_j + \frac{\lambda}{2} \mathbf{v}_j^{\top} \mathbf{v}_j - \sum_{i=1}^n \omega_{ij} a_{ij} \mathbf{u}_i^{\top} \mathbf{v}_j \\ &= \frac{1}{2} \left(\mathbf{v}_j^{\top} \left(\sum_{i=1}^n \omega_{ij} \mathbf{u}_i \mathbf{u}_i^{\top} \right) \mathbf{v}_j + \lambda \mathbf{v}_j^{\top} \mathbf{v}_j \right) - \left(\sum_{i=1}^n \omega_{ij} a_{ij} \mathbf{u}_i^{\top} \right) \mathbf{v}_j \\ &= \frac{1}{2} \mathbf{v}_j^{\top} \left(\sum_{i=1}^n \omega_{ij} \mathbf{u}_i \mathbf{u}_i^{\top} + \lambda \mathbf{I}_k \right) \mathbf{v}_j - \left(\sum_{i=1}^n \omega_{ij} a_{ij} \mathbf{u}_i^{\top} \right) \mathbf{v}_j.\end{aligned}$$

Expand norms and remove $\|\mathbf{U}\|_F^2$, because it does not depend on v_j .

Remove all terms that do not contain entries of v_j .

$$\mathbf{u}_i^{\top} \mathbf{v}_j = \mathbf{v}_j^{\top} \mathbf{u}_i.$$

Analogously, we can construct the following objective w.r.t. u_i ,

$$\ell_{\mathbf{V}}(u_i) = \frac{1}{2} \mathbf{u}_i^{\top} \left(\sum_{j=1}^m \omega_{ij} \mathbf{v}_j \mathbf{v}_j^{\top} + \lambda \mathbf{I}_k \right) \mathbf{u}_i - \left(\sum_{j=1}^m \omega_{ij} a_{ij} \mathbf{v}_j^{\top} \right) \mathbf{u}_i.$$

Since the matrix in the first term is rank- k , we can invert it.⁹ We have the following first-order optimal solution to the above loss functions,

⁹ This is because of the identity matrix, which is a result of the regularization.

$$\begin{aligned}\mathbf{v}_j^* &= \left(\sum_{i=1}^n \omega_{ij} \mathbf{u}_i \mathbf{u}_i^{\top} + \lambda \mathbf{I}_k \right)^{-1} \left(\sum_{i=1}^n \omega_{ij} a_{ij} \mathbf{u}_i \right) \\ \mathbf{u}_i^* &= \left(\sum_{j=1}^m \omega_{ij} \mathbf{v}_j \mathbf{v}_j^{\top} + \lambda \mathbf{I}_k \right)^{-1} \left(\sum_{j=1}^m \omega_{ij} a_{ij} \mathbf{v}_j \right).\end{aligned}$$

This involves taking the inverse of a $k \times k$ matrix, thus computing \mathbf{u}_i^* and \mathbf{v}_j^* has $\mathcal{O}(k^3)$ complexity, which is generally quite fast, because k is small.

Alternating least squares (ALS) makes use of these subproblems by alternating between optimizing \mathbf{U} given \mathbf{V} and \mathbf{V} given \mathbf{U} ,

$$\begin{aligned}\mathbf{V}^{(t+1)} &= \underset{\mathbf{V}}{\operatorname{argmin}} \ell_{\mathbf{U}^{(t)}}(\mathbf{V}) \\ \mathbf{U}^{(t+1)} &= \underset{\mathbf{U}}{\operatorname{argmin}} \ell_{\mathbf{V}^{(t+1)}}(\mathbf{U}).\end{aligned}$$

An iteration of ALS has complexity $\mathcal{O}((n+m) \cdot k^3)$. However, since each row of both matrices can be optimized in parallel, this can be done much faster in practice. Furthermore, since we are optimizing the objective fully over one half of the parameter space, the objective monotonically decreases. Hence, it converges to a first-order optimal fixed point, which may be a saddle point.

Moreover, we have the advantage that we can easily augment the model by adding additional dimensions and minimizing w.r.t. the existing dimensions, which may be useful when new users enter the system or new items become available.

Projection methods. We will now consider optimizing the objective by gradient descent, where the gradient is

$$\begin{aligned}\frac{\partial}{\partial \mathbf{B}} \frac{1}{2} \|\Pi_{\Omega}(\tilde{\mathbf{A}} - \mathbf{B})\|_F^2 &= \Pi_{\Omega}(\tilde{\mathbf{A}} - \mathbf{B}) \odot \frac{\partial}{\partial \mathbf{B}} \Pi_{\Omega}(\tilde{\mathbf{A}} - \mathbf{B}) \\ &= \Pi_{\Omega}(\tilde{\mathbf{A}} - \mathbf{B}) \odot \Pi_{\Omega}(-\mathbf{1}_{n \times m}) \\ &= -\Pi_{\Omega}(\tilde{\mathbf{A}} - \mathbf{B}).\end{aligned}$$

$$\frac{\partial}{\partial \mathbf{X}} \|\mathbf{X}\|_F^2 = 2\mathbf{X}.$$

Projection is linear.

However, gradient descent does not constrain its iterates to be rank- k matrices. Projected gradient descent solves this by projecting to the constrained space in between every gradient step. In general, it is hard to project to a space of rank- k matrices, but SVD makes it possible by making use of the Eckart-Young theorem. We will denote $[\mathbf{M}]_k$ as the projection to the space of rank- k matrices, which is computed by pruning all singular values below σ_k , as shown by the Eckart-Young theorem. Hence, the update rule is the following,

$$\mathbf{A}^{(t+1)} = \left[\mathbf{A}^{(t)} + \eta \Pi_{\Omega}(\tilde{\mathbf{A}} - \mathbf{A}^{(t)}) \right]_k.$$

This converges to a first-order optimal solution if η is small enough as a general result of projected gradient descent.

The problem with this approach is that the space of rank- k matrices is non-convex and is thus not guaranteed to converge to a global optimum. Thus, the next idea is to find the tightest convex relaxation of this space. *I.e.*, we want to find a convex space that contains all rank- k matrices, but not too many more. For this, we use the nuclear norm,

$$\|\mathbf{M}\|_* \doteq \sum_{i=1}^{\operatorname{rank}(\mathbf{M})} \sigma_i,$$



Figure 3.2. Illustration of projected gradient descent, where the dotted lines indicate projection steps.

where σ_i is the i -th singular value of M .

Definition 3.6 (Convex envelope). The convex envelope of a function $f : \mathcal{X} \rightarrow \mathbb{R}$ is the largest convex function g such that

$$g(x) \leq f(x), \quad \forall x \in \mathcal{X}.$$

Theorem 3.7. The convex envelope of $\text{rank}(\cdot)$ on $\{M \mid \|M\|_2 \leq 1\}$ is the nuclear norm $\|\cdot\|_*$.

Proof. Let $M \in \mathbb{R}^{n \times m}$ with $\|M\|_2 \leq 1$. By Lemma 3.4, we thus have $\sigma_1 \leq 1$. Thus,

$$\|M\|_* = \sum_{i=1}^{\text{rank}(M)} \sigma_i \leq \sum_{i=1}^{\text{rank}(M)} 1 = \text{rank}(M).$$

Furthermore, any valid norm is a convex function. Hence, by Definition 3.6, this concludes the proof. ■

We will use this to approximate the objective by a convex function. We have the following objective,

$$\ell(B) = \frac{1}{2} \|\Pi_{\Omega}(\tilde{A} - B)\|_F^2, \quad \text{rank}(B) \leq k, \|B\|_2 \leq c.$$

This can be rewritten using Lagrange multipliers as a non-convex objective,

$$\begin{aligned} \ell(B) &= \frac{1}{2} \|\Pi_{\Omega}(\tilde{A} - B)\|_F^2 + r \cdot \text{rank}(B) + \mu \|B\|_2 \\ &\approx \frac{1}{2} \|\Pi_{\Omega}(\tilde{A} - B)\|_F^2 + \tau \|B\|_* + \gamma \|B\|_F. \end{aligned}$$

This approximation of this objective is convex, so, using gradient descent and correct η , we can optimize it to a global optimum.

Definition 3.8 (Shrinkage operator). Let $X \in \mathbb{R}^{n \times m}$ with SVD $X = U \text{diag}(\sigma) V^T$ be given, then

$$\text{shrink}_{\tau}(X) \doteq U \text{diag}(\sigma - \tau)_+ V^T,$$

where M_+ clips all entries of M at zero.

Note that the rank of $\text{shrink}_{\tau}(C)$ decreases monotonically with τ , because more singular values are zeroed out as τ increases.

Theorem 3.9. The shrinkage operator minimizes the following optimization problem,

$$\text{shrink}_{\tau}(X) \in \arg\min_B \frac{1}{2} \|X - B\|_F^2 + \tau \|B\|_*.$$

Let $\sigma(M)$ be the vector of singular values of M , then, as shown by Lemma 3.3, the Frobenius norm can be computed by

$$\|M\|_F = \|\sigma(M)\|_2.$$

Similarly, the nuclear norm can be computed by

$$\|M\|_* = \|\sigma(M)\|_1.$$

Furthermore, the spectral norm can be computed by

$$\|M\|_2 = \|\sigma(M)\|_{\infty}.$$

The second constraint adds regularization and is necessary due to the condition of the space of Theorem 3.7.

The nuclear norm is the convex envelope of the rank function and the Frobenius norm upper bounds the spectral norm.

Using the shrink operator, we can define a new update rule, starting from $\mathbf{Y}^{(0)} = \mathbf{0}$,

$$\begin{aligned}\mathbf{X}^{(t+1)} &= \text{shrink}_{\tau}(\mathbf{Y}^{(t)}) \\ \mathbf{Y}^{(t+1)} &= \mathbf{Y}^{(t)} + \eta_t \Pi_{\Omega}(\tilde{\mathbf{A}} - \mathbf{X}^{(t+1)}).\end{aligned}$$

With a suitable schedule $\eta_t > 0$, the sequence $\{\mathbf{X}^{(t)}\}$ will converge to the solution of the following optimization problem,

$$\lim_{t \rightarrow \infty} \mathbf{X}^{(t)} \in \underset{\mathbf{B}}{\operatorname{argmin}} \tau \|\mathbf{B}\|_* + \frac{1}{2} \|\mathbf{B}\|_F^2, \quad \Pi_{\Omega}(\tilde{\mathbf{A}}) = \Pi_{\Omega}(\mathbf{B}).$$

Hence, the result of this algorithm will be exactly reproducing the observed entries, which the non-relaxed problem cannot guarantee, as there may be no rank- k algorithm with a projected residual that is zero. However, the convex relaxation approach does not define a low-rank sequence of matrices. Rather, the shrinkage operator is used to implicitly encourage low-rank approximations.

The singular value projection approach maintains a rank- k matrix, while the convex relaxation approach maintains a sparse iterate sequence and converges to reproduce the observed entries exactly. Both are beneficial relative to a dense representation of \mathbf{A} .

Exact recovery. So far we have only been interested in how to find the “best” completed matrix, given an incomplete matrix $\tilde{\mathbf{A}}$. However, now we are interested in the conditions under which we can exactly recover the matrix, assuming that the underlying matrix \mathbf{A} is of rank k .

Assume $\mathbf{A} \in \mathbb{R}^{n \times n}$, then the degrees of freedom are the k singular values and k left and right singular vectors. The i -th singular vector has $n - i$ degrees of freedom due to the constraints of unit length and pairwise orthogonality, thus the SVD has the following degrees of freedom,

$$k + 2 \sum_{i=1}^k n - i = k + 2 \left(nk - \sum_{i=1}^k i \right) = 2nk - k^2.$$

Thus, a necessary condition to exactly recover $\mathbf{A} \in \mathbb{R}^{n \times n}$ is that we have at least that number of observations. This condition is not sufficient, which we will prove by an example. Let \mathbf{A} have the SVD $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^{\top}$, then

$$\mathbf{A} = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^{\top} + \sum_{i=2}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^{\top}.$$

Let $\mathbf{u}_1 = \mathbf{e}_i$ and $\mathbf{v}_1 = \mathbf{e}_j$, then in order to recover σ_1 , we need to have sampled a_{ij} .

For exact recovery, it does not only matter how many entries we have sampled, but also where we have sampled them. Intuitively, we need the information to be “sufficiently distributed”.

The following gives us conditions for which we can exactly reconstruct the underlying matrix with high probability, under the assumption that the entries have been uniformly sampled.

Definition 3.10 (Incoherence). Let A be a rank- k matrix with SVD $A = U\Sigma V^\top$. Furthermore, define the following matrices,

$$P \doteq \sum_{i=1}^k u_i u_i^\top, \quad Q \doteq \sum_{i=1}^k v_i v_i^\top, \quad E \doteq \sum_{i=1}^k u_i v_i^\top.$$

Then, A is incoherent with parameter μ if and only if the following conditions are satisfied,

$$\begin{aligned} |p_{ij}|, |q_{ij}| &\leq \frac{\mu\sqrt{k}}{n}, \quad i \neq j \\ \left| p_{ii} - \frac{k}{n} \right|, \left| q_{ii} - \frac{k}{n} \right| &\leq \frac{\mu\sqrt{k}}{n} \\ |e_{ij}| &\leq \frac{\mu\sqrt{k}}{n}. \end{aligned}$$

Theorem 3.11. Let $A \in \mathbb{R}^{n \times n}$ be a rank- k matrix that is incoherent with $\mu \geq 1$ and for which S samples have been observed at random. Then, there is a universal constant C such that if $S \geq C\mu^2 nk(\log n)^6$, then with probability at least $1 - n^{-3}$, A fulfills

$$A = \underset{B}{\operatorname{argmin}} \|B\|_*, \quad \Pi_\Omega(B) = \Pi_\Omega(A).$$

This says that if the matrix is recoverable from its sampled entries, then it can be recovered via nuclear norm minimization.

3.3 Randomized methods for SVD

Considering that most algorithms introduced in this section make use of the singular value decomposition as a fundamental operator, we will now consider algorithms for computing it approximately in a low time complexity.

Let $A \in \mathbb{R}^{n \times m}$. The idea is to find an orthogonal matrix $Q \in \mathbb{R}^{n \times 2k}$ such that $A \approx QQ^\top A$. Then, we perform SVD of the smaller matrix $B \doteq Q^\top A = \tilde{U}\tilde{\Sigma}\tilde{V}^\top$, where $B \in \mathbb{R}^{2k \times m}$, and extend this SVD to A , $A \approx (Q\tilde{U})\tilde{\Sigma}\tilde{V}^\top$. The same can be done for the columns of B to further reduce the time complexity.

However, the key question is how to choose Q appropriately. The following is a scheme for finding a suitable Q ,

1. Generate a random matrix with i.i.d. Gaussian entries, $R \in \mathbb{R}^{m \times 2k}$;
2. Compute $Y = (AA^\top)^q AR$, where $Y \in \mathbb{R}^{n \times 2k}$, by repeated matrix



Figure 3.3. Schematic view of how to use random projections to compute the matrix decompositions more efficiently.

multiplications;

3. Construct an orthonormal basis Q for the image of Y (e.g., via Gram-Schmidt).

4 Latent variable models

The philosophy behind latent variable models is that we have *observables* X , which are augmented by *latent variables* Z , which happens by specifying a *complete data model* $p(X, Z)$, which implies a *marginal model*,

$$p(X = x) = \sum_{z \in \mathcal{Z}} p(X = x, Z = z).$$

The marginal model can be specified by a conditional $p(X | Z)$ and prior $p(Z)$,

$$p(X = x) = \sum_{z \in \mathcal{Z}} p(X = x | Z = z) p(Z = z).$$

4.1 Probabilistic clustering models

Assume we are given a dataset of s observables $\{x_t | t \in [s]\}$. The conceptually simplest family of latent variable models assigns a k -class categorical random variable Z_t to each observable. The latent information tags an observable as a member of that class.

Specifically, we have the following prior categorical distribution,

$$Z_t \sim \text{Categorical}(\pi_1, \dots, \pi_k), \quad P(Z_t = z) = \pi_z,$$

where $\pi \in \Delta^{k-1}$ is an unknown parameter.¹⁰ To fully parametrize the latent variable model, we need a class-conditional distribution for each class, $p(x | z)$. By parametrizing the class-conditional distribution of z by θ_z , we have the following parametrized distribution over X ,

$$p(x; \theta) = \sum_{z=1}^k p(z) p(x | z) = \sum_{z=1}^k \pi_z p(x; \theta_z).$$

We interpret this as a mixture distribution – convex combinations of class-specific distributions. Moreover, this model is fully parametrized by the prior and class-conditional distribution parameters,

$$\theta = [\pi, \theta_1, \dots, \theta_k].$$

For given parameters θ , we can use Bayes' rule to compute the latent posteriors,

$$p(z | x; \theta) = \frac{\pi_z p(x; \theta_z)}{\sum_{\zeta=1}^k \pi_{\zeta} p(x; \theta_{\zeta})}.$$

We can interpret these probabilities as observable-specific probabilistic cluster memberships.

Learning the parameters. A common approach to learn the best parameters is to maximize the likelihood of the data, which is called maximum likelihood estimation (MLE). In other words, we choose the model parameters that maximize the probability of the observed data,

$$\ell(\theta) = \log \prod_{t=1}^s p(x_t; \theta) = \sum_{t=1}^s \log p(x_t; \theta) = \sum_{t=1}^s \log \sum_{z=1}^k \pi_z p(x_t; \theta_z).$$

¹⁰ Δ^{k-1} is the k -dimensional probability simplex, which is the space of vectors where all elements are non-negative and sum to 1.

This gives us the following optimization problem,

$$\theta^* \in \operatorname{argmax}_{\theta} \ell(\theta),$$

which we optimize by the expectation-maximization (EM) algorithm. The EM algorithm is a general tool for learning in latent variable models, but we will introduce it in the context of mixture models.

Let $q_t \in \Delta^{k-1}, \forall t \in [s]$ be the variational parameters. Using this we can derive the evidence lower bound (ELBO),

$$\begin{aligned} \ell(\theta) &= \sum_{t=1}^s \log \sum_{z=1}^k q_{tz} \frac{\pi_z}{q_{tz}} p(\mathbf{x}_t; \theta_z) \\ &\geq \sum_{t=1}^s \sum_{z=1}^k q_{tz} \log \left(\frac{\pi_z}{q_{tz}} p(\mathbf{x}_t; \theta_z) \right) && \text{Jensen's inequality, log is concave.} \\ &= \sum_{t=1}^s \sum_{z=1}^k q_{tz} (\log \pi_z + \log p(\mathbf{x}_t; \theta_z) - \log q_{tz}) \\ &\doteq \ell(\theta; \{q_t\}_{t=1}^s). \end{aligned}$$

Since this lower bounds the MLE, we can use the ELBO as an objective to implicitly maximize the log-likelihood. Maximizing w.r.t. q increases the tightness of the bound of ELBO on the MLE,¹¹ while maximizing w.r.t. θ improves the model fit.

¹¹ q does not affect the likelihood, so an increase can only tighten the bound.

Firstly, we will solve for the optimal choice of q_t ,

$$\ell(\theta; \{q_t\}_{t=1}^s), \quad q_t \in \Delta^{k-1}, \forall t \in [s].$$

We can solve for every q_t independently, because the ELBO is separable w.r.t. q_t . Furthermore, we enforce the normalization constraint on q_t by introducing a Lagrange multiplier λ . This results in the following objective to maximize per q_t ,

$$q_t^* \in \operatorname{argmax}_{q_t} \sum_{z=1}^k q_{tz} (\log \pi_z + \log p(\mathbf{x}_t; \theta_z) - \log q_{tz}) - \lambda \left(\sum_{z=1}^k q_{tz} - 1 \right).$$

The derivative w.r.t. q_{tz} of this function is

$$\begin{aligned} \frac{\partial \ell(\theta; q_t)}{\partial q_{tz}} &= \frac{\partial}{\partial q_{tz}} \sum_{z'=1}^k q_{tz'} (\log p(\mathbf{x}_t; \theta_{z'}) + \log \pi_{z'} - \log q_{tz'}) \\ &\quad - \frac{\partial}{\partial q_{tz}} \lambda \left(\sum_{z'=1}^k q_{tz'} - 1 \right) \\ &= \log p(\mathbf{x}_t; \theta_z) + \log \pi_z - \frac{\partial}{\partial q_{tz}} q_{tz} \log q_{tz} - \lambda \\ &= \log p(\mathbf{x}_t; \theta_z) + \log \pi_z - \log q_{tz} - \frac{q_{tz}}{q_{tz}} - \lambda \\ &= \log p(\mathbf{x}_t; \theta_z) + \log \pi_z - \log q_{tz} - \lambda - 1. \end{aligned}$$

Thus, we have the following first-order optimality condition,

$$\log p(\mathbf{x}_t; \theta_z) + \log \pi_z - \log q_{tz} \stackrel{!}{=} \lambda + 1.$$

Exponentiating both sides yields

$$q_{tz} \stackrel{!}{=} \frac{\pi_z p(\mathbf{x}_t; \boldsymbol{\theta}_z)}{e^{\lambda+1}}.$$

Enforcing the constraint of $q_t \in \Delta^{k-1}$, we get

$$q_{tz} = \frac{\pi_z p(\mathbf{x}_t; \boldsymbol{\theta}_z)}{\sum_{\zeta=1}^k \pi_{\zeta} p(\mathbf{x}_t; \boldsymbol{\theta}_{\zeta})}.$$

As we saw before, this is the posterior of the latent class variable $p(z \mid \mathbf{x}_t; \boldsymbol{\theta})$. Note that the optimal choice of the variational parameters q_t depends on the parameters $\boldsymbol{\theta}$. Hence, it is only a partial step, called the expectation (E) step.

Now we also need a step to maximize the model parameters $\boldsymbol{\theta}$. We can easily solve for π by

$$\pi_z^* = p(z) = \sum_{t=1}^s p(\mathbf{x}_t, z) = \sum_{t=1}^s p(z \mid \mathbf{x}_t) p(\mathbf{x}_t) = \frac{1}{s} \sum_{t=1}^s q_{tz}.$$

This assumes that the data is distributed uniformly.

This solution can also be derived by the first-order optimality condition with the Lagrangian.

Moreover, the solution for $\boldsymbol{\theta}_z$ depends on the choice of the model distribution, but we can generally get to separable problems,

$$\boldsymbol{\theta}_z^* \in \operatorname{argmax}_{\boldsymbol{\theta}_z} \sum_{t=1}^s q_{tz} \log p(\mathbf{x}_t; \boldsymbol{\theta}_z).$$

This means that the parameters for different classes z are decoupled given the variational parameters q . Thus, for each component, we only have to solve a weighted MLE problem, which is often possible to do analytically. This partial step is called the maximization (M) step.

Gaussian mixture model. We will consider a common special case of the above framework, where we specify the component models $p(\mathbf{X} \mid Z)$ by Gaussians with unit variance,

$$p(\mathbf{x}; \pi, \{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k\}) = \sum_{z=1}^k \pi_z \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_z, \mathbf{I}_d),$$

where

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-d/2} \det(\boldsymbol{\Sigma})^{-1/2} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right).$$

The EM algorithm then consists of the following alternating equations,

$$q_{tz} \doteq p(z \mid \mathbf{x}_t, \boldsymbol{\theta}) = \frac{\pi_z \exp\left(-\frac{1}{2}\|\mathbf{z}_t - \boldsymbol{\mu}_z\|^2\right)}{\sum_{\zeta=1}^k \pi_{\zeta} \exp\left(-\frac{1}{2}\|\mathbf{z}_t - \boldsymbol{\mu}_{\zeta}\|^2\right)} \quad \text{E-step.}$$

$$\boldsymbol{\mu}_z = \frac{\sum_{t=1}^s q_{tz} \mathbf{x}_t}{\sum_{t=1}^s q_{tz}}, \quad \pi_z = \frac{1}{s} \sum_{t=1}^s q_{tz}. \quad \text{M-step.}$$

Intuitively, the E-step (soft) clusters the data points and the M-step computes the weighted centroids of each component.

4.2 Topic models

We will now consider a class of latent variable models known as topic models. These are used to analyze document collections and to discover the topical content of documents. Informally, topical content is the information that the words of a document carry about what the document is about.

Let Σ be the word vocabulary with $|\Sigma| = m$ words. A document d_i of length s_i is part of a collection of n documents and is a field of random variables X_{it} , whose realizations are words $x_{it} \in \Sigma, \forall t \in [s_i]$. We will assume that topical content is invariant to word order.¹² As a consequence, the sufficient statistics of a document are the frequencies of word occurrences within it. Hence, we can reduce the data to a bag-of-words representation of occurrence counts,

$$N_{ij} = |\{x_{it} = w_j \mid t \in [s_i]\}|.$$

In words, N_{ij} denotes how often word w_j occurred in document d_i . We can thus summarize the full document corpus in an occurrence matrix,

$$\mathbf{N} = [N_{ij}] \in \mathbb{N}^{n \times m}.$$

Moreover, we have to conceptualize what we mean by “topic”. Generally, topics refer to things that people are interested to talk or write about. Hence, we can represent topics by latent variables $Z \in [k]$, akin to what we saw in mixture models. We will then associate these topics with word occurrences, which induce topics of entire documents.¹³ Note that topics are not mutually exclusive, since documents can concern multiple.

We will assume that the documents in the corpus were created according to the following generative process.¹⁴ For each word, we sample a topic from $p(z \mid d)$, and then sample the words from $p(w \mid z)$. Thus, in order to define the latent variable data model, we need a document-conditional distribution over latent variables, $p(z \mid d)$, and a topic-conditional distribution over words, $p(w \mid z)$. We can then define a document-conditional word distribution,

$$p(w \mid d) = \sum_{z=1}^k p(w \mid z)p(z \mid d).$$

From this, we define the log-likelihood objective as

$$\ell(\boldsymbol{\theta}; \mathbf{N}) = \log p(\mathbf{N}; \boldsymbol{\theta}) = \log \prod_{i=1}^n \prod_{j=1}^m p(w_j \mid d_i)^{N_{ij}} = \sum_{i=1}^n \sum_{j=1}^m N_{ij} \log p(w_j \mid d_i).$$

Learning. We can equivalently write the log-likelihood objective in terms of the raw data,

$$\ell(\boldsymbol{\theta}) = \sum_{i=1}^n \sum_{t=1}^{s_i} \log \sum_{z=1}^k p(x_{it} \mid z)p(z \mid d_i).$$

¹² This is known as *exchangeability*, which says that the distribution of a sequence of random variables does not change under any permutation of their order.

¹³ If we were to associate the topic variables with documents, we would effectively be performing document clustering. However, in this case, we want the topics of documents to be induced by the words that they contain.

¹⁴ This is important so we can infer/reverse engineer a model from it. Also, it is important to make assumptions about the data explicit.

Similarly to the mixture models, we define an ELBO,

$$\ell(\theta) \geq \sum_{i=1}^n \sum_{t=1}^{s_i} \sum_{z=1}^k q_{itz} (\log p(x_{it} | z) + \log p(z | d_i) - \log q_{itz}).$$

Moreover, following the same steps, we can derive the EM equations,

$$q_{itz} = \frac{p(x_{it} | z)p(z | d_i)}{\sum_{\zeta=1}^k p(x_{it} | \zeta)p(\zeta | d_i)} \quad \text{E-step.}$$

$$p(w_j | z) = \frac{\sum_{i=1}^n \sum_{t=1}^{s_i} q_{itz} \cdot \mathbb{1}\{x_{it} = w_j\}}{\sum_{i=1}^n \sum_{t=1}^{s_i} q_{itz}}, \quad p(z | d_i) = \frac{1}{s_i} \sum_{t=1}^{s_i} q_{itz}. \quad \text{M-step.}$$

Intuitively, the E-step computes the posterior probabilities, $q_{itz} = p(z | x_{it}, d_i)$, where $p(z | d_i)$ acts as a prior. This yields a probabilistic clustering of word occurrences. The M-step computes the MLE for a q -weighted multinomial sample. This algorithm will converge, but not necessarily to the global maximizer.

We can use the class-conditional word distribution, $p(w | z)$, to find similar words that are connected by a common topic.

Latent Dirichlet allocation. The problem with the above topic model is that it assumes a fixed set of documents and selects the parameters to maximize the predictability of words within these. The natural next step is to extend the above in a way that accounts for modeling unseen documents. The generative process of this model is that we first choose a distribution over topics, $p(z | \alpha)$. Then, for each word, sample a topic from $p(z | \alpha)$ and sample the word from $p(w | z)$.

Effectively, the latent Dirichlet allocation (LDA) model takes a “Bayesian step up” from latent variable modeling. Its main idea is that it models a distribution over priors, which are modeled by mixture vectors,

$$v \in \Delta^{k-1}.$$

This parameter is distributed according to a Dirichlet distribution, and has hyperparameter α ,

$$p(v; \alpha) \propto \prod_{z=1}^k v_z^{\alpha_z - 1}, \quad \alpha_z > 0, \forall z \in [k].$$

The Dirichlet distribution is chosen as the prior because it is a conjugate prior of the categorical distribution. Typically, we set $\alpha_k = \alpha$ and optimize α on held-out validation data.

Let $\mathbf{U} = [p(w_j | z)] \in [0, 1]^{m \times k}$, then we have the following distribution,

$$p(X = [x_1, \dots, x_s] | \mathbf{U}) = \int \prod_{t=1}^s p(x_t | \mathbf{U}, v) p(v; \alpha) dv$$

$$p(w_j | \mathbf{U}, v) = \sum_{z=1}^k u_{jz} v_z.$$

As can be seen, the probabilities are not conditioned on the observed documents. As a result, when a new document is introduced, we can use the same distributions, whose parameters were learned from an existing document collection. In this sense, the LDA model is more robust than the topic model that was initially introduced.

Probabilistic matrix decomposition. The topic model introduced above is intimately related to matrix decomposition, where we see that the topic variable $z \in [k]$ plays the role of a rank constraint. More specifically, we define the following matrices,

$$\mathbf{U} \doteq [p(w_j | z)] \in [0, 1]^{m \times k}, \quad \mathbf{V} \doteq [p(z | d_i)] \in [0, 1]^{n \times k}.$$

We then form $\hat{\mathbf{N}} = \mathbf{UV}^\top$, which is a rank- k matrix. We interpret the entries of this new matrix $\hat{\mathbf{N}}$ as

$$\hat{N}_{ji} = \mathbf{u}_j^\top \mathbf{v}_i = \sum_{z=1}^k p(w_j | z) p(z | d_i) = p(w_j | d_i).$$

Note that $\hat{\mathbf{N}}$ is a relative version of \mathbf{N} , where $N_{ij} \approx s_i \hat{N}_{ij}$. Thus, we have a matrix decomposition with additional constraints,

$$u_{jz} \geq 0, \forall j \in [m], z \in [k], \quad v_{zi} \geq 0, \forall i \in [n], z \in [k].$$

This is known as non-negative matrix factorization (NMF).¹⁵ The objective follows from the maximum likelihood principle,

$$\ell(\mathbf{U}, \mathbf{V}; \mathbf{N}) = \sum_{i=1}^n \sum_{j=1}^m N_{ij} \log \hat{N}_{ij}, \quad \hat{\mathbf{N}} = \mathbf{UV}^\top.$$

Further, we have normalization constraints on the row and column level, namely

$$\sum_{j=1}^m u_{jz} = 1, \forall z \in [k], \quad \sum_{z=1}^k v_{zi} = 1, \forall i \in [n].$$

Thus, we have a special case of non-negative matrix decomposition with a log-likelihood objective.

4.3 Embeddings

In natural language, the atomic units of meaning are symbols, such as words. Generally, these symbols do not carry their meaning with them. Rather, their meaning come from their use within its language. The idea of embeddings is to learn representations that capture semantics by embedding symbols in vector space. Specifically, we want to learn a mapping from words to vectors such that the vectors represent word semantics.

The latent variable approach to this problem is to learn latent representations that are predictive of observations. Thus, we need to design a task such that the latent representations of the words must have some



Figure 4.1. Factors identified by non-negative matrix factorization in a face reconstruction task.



Figure 4.2. Factors identified by principal component analysis in a face reconstruction task.

¹⁵ For NMF, we can use the ALS algorithm, where we add projection steps to enforce non-negativity,

$$u_{jz} = \max\{0, u_{jz}\}, \quad v_{zi} = \max\{0, v_{zi}\}.$$

Figures 4.1 and 4.2 show identified factors of NMF and PCA on a face reconstruction task. As can be seen, NMF tends to identify part-based representations and its features are sparse, because there is no way of removing added features, due to the non-negativity constraint. In other words, NMF models must be careful in what they add.

form of semantics to be able to perform the task well. In the skip-gram model, we use the task of predicting whether a word is in the context of another.¹⁶ Effectively, we are treating the embedding of each word as a latent variable that predicts the co-occurring of context words. After the model has converged, we do not care about performing the task well, but rather the parameters (embeddings) of the model that lead to the task being performed well.

This task has the following likelihood function that we wish to maximize,

$$\ell(\theta) = \log \prod_{t=1}^T \prod_{\delta \in \mathcal{I}} p(x_{t+\delta} \mid x_t; \theta) = \sum_{t=1}^T \sum_{\delta \in \mathcal{I}} \log p(x_{t+\delta} \mid x_t; \theta),$$

where θ contains the embeddings as parameters of this model and \mathcal{I} is a set of displacements, e.g., $\mathcal{I} = \{-R, \dots, -1, 1, \dots, R\}$. The probabilities are computed by normalized inner products,

$$p(v \mid w; \theta) = \frac{\exp\langle \mathbf{z}_w, \mathbf{z}_v \rangle}{\sum_{u \in \Sigma} \exp\langle \mathbf{z}_w, \mathbf{z}_u \rangle}.$$

We can further refine this by introducing biases $b_v \in \mathbb{R}$ to explicitly control the marginal probability and using different embeddings for conditioned and predicted words,¹⁷ leading to

$$p(v \mid w; \theta) = \frac{\exp(\langle \zeta_w, \mathbf{z}_v \rangle + b_v)}{\sum_{u \in \Sigma} \exp(\langle \zeta_w, \mathbf{z}_u \rangle + b_u)}$$

with parameters

$$\theta_w = (\mathbf{z}_w, \zeta_w, b_w) \in \mathbb{R}^{2m+1}.$$

Sufficient statistics for this model is a co-occurrence matrix, where

$$N_{vw} = |\{t \mid x_t = w, x_{t+\delta} = v, \delta \in \mathcal{I}\}|.$$

The log-likelihood can then be computed by

$$\ell(\theta) = \sum_{v \in \Sigma} \sum_{w \in \Sigma} N_{vw} \left(\langle \zeta_w, \mathbf{z}_v \rangle + b_v - \log \sum_{u \in \Sigma} \exp(\langle \zeta_w, \mathbf{z}_u \rangle + b_u) \right).$$

The problem with this approach is that computing the normalization constant is expensive. We can discard normalization by reformulating the prediction problem as a classification problem. For this, we also need positive samples,

$$\mathcal{S}^+ = [(x_t, x_{t+\delta}) \mid t \in [T], \delta \in \mathcal{I}].$$

And, additionally we need negative samples, which we randomly sample,

$$\mathcal{S}^- = [(x_t, v_{tj}) \mid t \in [T], v_{tj} \stackrel{\text{iid}}{\sim} q, j \in [r]],$$

where q is a distribution over Σ and r is the sampling factor. Generally, q is chosen to satisfy

$$q(w) \propto p(w)^\alpha,$$

¹⁶ Rupert Firth: “You shall know a word by the company it keeps.”

¹⁷ This is necessary, because words are often not within their own context window, but will have a high inner product, because their vectors are equal.

The $[]$ notation indicates a multiset.

where typically $\alpha = 3/4$. The intuition behind this is that what matters most in learning semantic representations is not the very frequent words, which carry little meaning, and also not the infrequent words, but the in-between range. This choice of α makes them more likely to be sampled, as shown in Figure 4.3.

We can then define a logistic log-likelihood function that we wish to minimize,

$$\ell(\theta) = \sum_{(w,v) \in \mathcal{S}^+} \log \sigma(\langle \zeta_w, z_v \rangle + b_v) + \sum_{(w,u) \in \mathcal{S}^-} \log(1 - \sigma(\langle \zeta_w, z_u \rangle + b_u)).$$

This does not require computing a normalization constant.

Pointwise mutual information. Let $p(v, w)$ denote the true distribution of co-occurring words and by $q(v, w) = p(w)p(v)$ the distribution used for negative sampling. The optimal Bayesian classifier would then be

$$\mathbb{P}((v, w) = \text{true}) = \frac{\pi p(v, w)}{\pi p(v, w) + (1 - \pi)q(v, w)},$$

where π is the prior probability of a true pair. Considering the pre-image (logit) of the logistic function, we get

$$\begin{aligned} h_{vw}^* &= \sigma^{-1}\left(\frac{\pi p(v, w)}{\pi p(v, w) + (1 - \pi)q(v, w)}\right) \\ &= \log\left(\frac{\pi p(v, w)}{\pi p(v, w) + (1 - \pi)q(v, w)} \cdot \frac{\pi p(v, w) + (1 - \pi)q(v, w)}{(1 - \pi)q(v, w)}\right) \\ &= \log\left(\frac{\pi p(v, w)}{(1 - \pi)q(v, w)}\right) \\ &= \log \frac{p(v, w)}{q(v, w)} + \log \frac{\pi}{1 - \pi}. \end{aligned}$$

Hence, in the case of balanced classes $\pi = 1/2$ (equivalent to $r = 1$) and $\alpha = 1$, we get

$$h_{vw}^* = \log \frac{p(v, w)}{p(w)p(v)},$$

which is the pointwise mutual information.

GloVe. Global word vectors (GloVe) considers word embedding models in terms of matrix factorization. It maximizes a different objective,

$$\ell(\theta, N) = \sum_{\substack{v, w \in \Sigma \\ N_{vw} > 0}} f(N_{vw}) (\log N_{vw} - \log \hat{N}_{vw})^2, \quad \hat{N}_{vw} = p(v, w; \theta),$$

which is a weighted square loss on the log-scale. In practice, the following weighting function is used,

$$f(N) = \min\left\{1, \left(\frac{N}{N_{\max}}\right)^\alpha\right\},$$

where often $\alpha = 3/4$.



Figure 4.3. Plot of $p(w)^\alpha$ for $\alpha = 3/4$.

$$\sigma^{-1}(p) = \log \frac{p}{1-p}, \quad p \in (0, 1).$$

The idea of this objective is that we can work with an unnormalized conditional probability distribution and simply choose

$$\log p(v, w) = \langle \zeta_w, z_v \rangle.$$

The reason for this is that the squared objective is two-sided, whereas a likelihood objective will always increase if we increase probabilities and the balancing effect comes purely from the normalization.

GloVe can be interpreted as a low-rank matrix factorization with

$$\mathbf{U} \doteq [\zeta_1, \dots, \zeta_n]^\top, \quad \mathbf{V} \doteq [z_1, \dots, z_n]^\top.$$

Then, we have the following matrix of unnormalized probabilities,

$$\log \hat{N} = \mathbf{U}\mathbf{V}^\top.$$

The GloVe objective is a weighted Frobenius norm of the approximation residual between the observed log-count matrix and a low-rank factorization of embedding matrices. As a special case, consider

$$f(N) = \min\{1, N\},$$

which results in a matrix completion problem,

$$\begin{aligned} \mathbf{U}, \mathbf{V} &\in \underset{\mathbf{U}, \mathbf{V}}{\operatorname{argmin}} \sum_{i=1}^n \sum_{j=1}^n \mathbb{1}\{N_{ij} > 0\} \left(\log N_{ij} - (\mathbf{U}\mathbf{V}^\top)_{ij} \right) \\ &= \left\| \Pi_{\mathbb{1}\{N>0\}} (\log \mathbf{N} - \mathbf{U}\mathbf{V}^\top) \right\|_F^2. \end{aligned}$$

This is a matrix completion problem with $\mathbf{A} = \log \mathbf{N}$ and $\omega_{ij} = \mathbb{1}\{N_{ij} > 0\}$.

We can optimize \mathbf{U} and \mathbf{V} by stochastic gradient descent,

$$\begin{aligned} \zeta_w &\leftarrow \zeta_w + 2\eta f(N_{vw})(\log N_{vw} - \langle \zeta_w, z_v \rangle) z_v \\ z_v &\leftarrow z_v + 2\eta f(N_{vw})(\log N_{vw} - \langle \zeta_w, z_v \rangle) \zeta_w, \end{aligned}$$

where we sample (v, w) at random.

5 Deep neural networks

Generally, deep learning models consist of a function $H : \mathbb{R}^n \rightarrow \mathbb{R}^p$ that extracts p -dimensional features from n -dimensional data and a linear map g that makes the final prediction. The final learned function is then formalized by $\psi = g \circ H$. Machine learning philosophies differ in the way that they extract features from the data. There are three main philosophies,

- Feature engineering: H is hand-crafted to extract intuitive features that have good predictive power of the label. The problem with this approach is that it requires domain expertise;
- Feature expansion: H maps to a high-dimensional feature space using kernels and implicit models;
- Compositionality: The feature extraction function H is learned through the composition of L simple building blocks,

$$H = H_L \circ H_{L-1} \circ \dots \circ H_2 \circ H_1, \quad H_l : \mathbb{R}^{n_{l-1}} \rightarrow \mathbb{R}^{n_l}.$$

In compositional models, the partial maps $H_{1:l} \doteq H_l \circ \dots \circ H_1$ produce intermediate representations. These satisfy the Markov property and, as such, need to preserve task-relevant information. Once information is lost, it cannot be recovered. The idea of the layers is to make relevant information more accessible and explicit with increasing depth, such that g can easily make a prediction.¹⁸

The key question is how to define the basic building blocks such that their composition is more powerful than any one block can be. Consider two linear layers,

$$F(x; W_2) = W_2 x, \quad G(x; W_1) = W_1 x.$$

Their composition is again a linear layer,

$$(F \circ G)(x) = W_2 W_1 x = W x, \quad W = W_2 W_1.$$

So, linear layers are not appropriate building blocks.

The key idea is to combine a linear map with a non-linearity,

$$H(x; W) = \Phi(Wx),$$

where Φ is a non-linear element-wise activation function,

$$\Phi(z) = [\phi(z_1), \dots, \phi(z_m)], \quad \phi : \mathbb{R} \rightarrow \mathbb{R}.$$

Theoretically, a neural network with one hidden layer and a non-polynomial activation function is a universal function approximator. This means that any function can be represented using a single hidden layer and a non-linear activation function.¹⁹ However, in practice, this single hidden layer

¹⁸ For example, the early layers of CNN typically learn low-level features, whereas the later layers will learn higher level features that make use of the low-level features.

¹⁹ The following are commonly used activation functions with their derivatives,

- Sigmoid, $\sigma : \mathbb{R} \rightarrow (0, 1)$,

$$\begin{aligned} \sigma(z) &\doteq \frac{1}{1 + e^{-z}} \\ \sigma'(z) &= \sigma(z)(1 - \sigma(z)) \\ &= \sigma(z)\sigma(-z); \end{aligned}$$

- Hyperbolic tangent, $\tanh : \mathbb{R} \rightarrow (-1, 1)$,

$$\begin{aligned} \tanh(z) &\doteq \frac{e^z - e^{-z}}{e^z + e^{-z}} \\ \tanh'(z) &= 1 - \tanh^2(z); \end{aligned}$$

- Rectified linear unit, $\text{ReLU} : \mathbb{R} \rightarrow \mathbb{R}_{\geq 0}$,

$$\begin{aligned} \text{ReLU}(z) &\doteq \max\{0, z\} \\ \text{ReLU}'(z) &= \mathbb{1}\{z > 0\}. \end{aligned}$$

may need to be infinitely large. A single hidden layer multi-layer perceptron (MLP) is formalized by the following function,

$$\psi(x; \beta, W) = \beta^\top \sigma(Wx) = \sum_{j=1}^m \beta_j \sigma(w_j^\top x) = \sum_{j=1}^m \frac{\beta_j}{1 + \exp(-w_j^\top x)}.$$

In this case, $g(y; \beta) = \beta^\top y$ and $H(x; W) = Wx$.

In order to tune this model, we need to be able to compute its gradients, which tell us how to locally optimize a loss function ℓ . In this case, we choose a squared loss, $\ell(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$. The gradients are computed by

$$\begin{aligned} \frac{\partial}{\partial \beta_j} \frac{1}{2} (\psi(x; \beta, W) - y)^2 &= (\psi(x; \beta, W) - y) \frac{\partial}{\partial \beta_j} \psi(x; \beta, W) \\ &= \frac{\psi(x; \beta, W) - y}{1 + \exp(-w_j^\top x)} \\ \frac{\partial}{\partial w_{ji}} \frac{1}{2} (\psi(x; \beta, W) - y)^2 &= (\psi(x; \beta, W) - y) \frac{\partial}{\partial w_{ji}} \psi(x; \beta, W) \\ &= (\psi(x; \beta, W) - y) \beta_j \frac{\partial}{\partial w_{ji}} \sigma(w_j^\top x) \\ &= (\psi(x; \beta, W) - y) \beta_j \sigma(w_j^\top x) \sigma(-w_j^\top x) \frac{\partial}{\partial w_{ji}} w_j^\top x \quad \sigma'(z) = \sigma(z)(1 - \sigma(z)) = \sigma(z)\sigma(-z). \\ &= \frac{\psi(x; \beta, W) - y}{1 + \exp(-w_j^\top x)} \frac{\beta_j x_i}{1 + \exp(w_j^\top x)}. \end{aligned}$$

Learning is typically carried out through stochastic gradient descent, which iteratively selects a random mini-batch $\mathcal{S}_t \subseteq \mathcal{S} = \{(x_i, y_i)\}_{i=1}^n$. The update rule is then

$$\theta_{t+1} = \theta_t - \eta \sum_{(x,y) \in \mathcal{S}_t} \nabla_{\theta} \ell(\psi(x; \theta), y).$$

5.1 Backpropagation

Calculating the gradient by hand for every model is very tedious and time consuming. Backpropagation is an algorithm that can compute the gradient of any function, which consists of building blocks with known gradients, in linear time. This algorithm makes use of dynamic programming, which means that it breaks the problem down into smaller subproblems and re-uses solutions to previously seen subproblems. In this case, the solution to the subproblems are the gradients, which can be re-used in later gradients by the chain rule and sum rule. In short, backpropagation exploits compositionality to efficiently compute the gradient.

Let $H_k : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be an intermediate layer of a compositional model. It's Jacobi map is defined as

$$[J_k]_{ij} \doteq \frac{\partial h_{ki}}{\partial z_j},$$

h_{ki} is the i -th output of layer k and z_j is its j -th input.

which is an implicit function of the input $\mathbf{z} \in \mathbb{R}^n$ to \mathbf{h}_k . Furthermore, define the error signal as

$$\delta_k \doteq \left[\frac{\partial \ell}{\partial \mathbf{h}_k} \right]^\top,$$

which quantifies how a change in loss is induced by a change in unit k 's activation. Using the chain rule, we can find a recurrence relation between error signals,

$$\delta_k \doteq \left[\frac{\partial \ell}{\partial \mathbf{h}_k} \right]^\top = \left[\frac{\partial \ell}{\partial \mathbf{h}_{k+1}} \frac{\partial \mathbf{h}_{k+1}}{\partial \mathbf{h}_k} \right]^\top \doteq \mathbf{J}_{k+1}^\top \delta_{k+1}.$$

Thus, we can use dynamic programming to compute all error signals efficiently. Lastly, in order to compute the gradient w.r.t. the parameters, we use the chain rule again,

$$\begin{aligned} \frac{\partial \ell}{\partial [\mathbf{W}_k]_{ij}} &= \frac{\partial \ell}{\partial \mathbf{h}_k} \frac{\partial \mathbf{h}_k}{\partial [\mathbf{W}_k]_{ij}} \\ &= \frac{\partial \ell}{\partial \mathbf{h}_k} \left[\frac{\partial \phi(\mathbf{y})}{\partial \mathbf{y}} \Big|_{\mathbf{y}=\mathbf{W}_k \mathbf{z}_{k-1}} \right] \frac{\partial}{\partial [\mathbf{W}_k]_{ij}} \mathbf{W}_k \mathbf{z}_{k-1} \\ &= \delta_k^\top \text{diag}(\phi'(\mathbf{W}_k \mathbf{z}_{k-1})) \text{vec}_i([\mathbf{z}_{k-1}]_j) \\ &= [\delta_k]_i \cdot \phi'([\mathbf{W}_k]_i^\top \mathbf{z}_{k-1}) \cdot [\mathbf{z}_{k-1}]_j. \end{aligned}$$

Intuitively, the local parameter gradient $\frac{\partial \ell}{\partial \mathbf{W}_k}$ is the product of an upstream vector \mathbf{z}_{k-1} , a downstream error signal δ_k , and the local sensitivity of the unit $\phi'(\mathbf{W}_k \mathbf{z}_{k-1})$. Note that we first need to perform a forward pass in order to compute these gradients in the backward pass.

Moreover, we have the following Jacobi maps for different activation functions,

- Linear activation, $\mathbf{J}_k = \mathbf{W}$;
- ReLU layer, $\mathbf{J}_k = \text{diag}(\mathbb{1}\{\mathbf{W}\mathbf{z} > 0\})\mathbf{W}$;
- Sigmoid layer, $\mathbf{J}_k = \text{diag}(\sigma(\mathbf{W}\mathbf{z}) \odot \sigma(-\mathbf{W}\mathbf{z}))\mathbf{W}$.

5.2 Gradient methods

In gradient descent, we iteratively update the parameters by

$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \eta \nabla \ell(\boldsymbol{\theta}_k), \quad \eta > 0.$$

A fundamental question is whether gradient descent will converge to an optimal solution. The key intuition is that gradient descent can only work if the gradient does not change too much relative to the step size. The gradient information must remain informative in a neighborhood around a point.

Definition 5.1 (Smoothness). A differentiable function $\ell : \mathbb{R}^d \rightarrow \mathbb{R}$ is L -smooth for some $L > 0$ if

$$\|\nabla \ell(\theta) - \nabla \ell(\theta')\| \leq L\|\theta - \theta'\|, \quad \forall \theta, \theta'.$$

This is equivalent to the gradient being L -Lipschitz continuous.

Smoothness implies a bound on the Hessian,

$$\|\nabla^2 \ell(\theta)\| \leq L, \quad \forall \theta.$$

As a consequence, this means that smoothness bounds the largest eigenvalue of the Hessian. A large L means that the gradient can change quickly, making it more unstable, thus we need to lower the stepsize. This is intuitively why the stepsize $\eta = 1/L$ works well.

Definition 5.2 (ϵ -critical point). Let ℓ be differentiable at θ , then θ is an ϵ -critical point if

$$\|\nabla \ell(\theta)\| \leq \epsilon.$$

Lemma 5.3. Gradient descent on an L -smooth, differentiable function $\ell : \mathbb{R}^d \rightarrow \mathbb{R}$ with step size $\eta = 1/L$ finds an ϵ -critical point in at most

$$k = \frac{2L(\ell(\theta_0) - \ell^*)}{\epsilon^2}$$

steps.

Thus, smoothness is sufficient to find local minima. The question is what properties of ℓ will ensure convergence to a global minimum.

Definition 5.4 (Polyak-Lojasiewicz condition). A differentiable function $\ell : \mathbb{R}^d \rightarrow \mathbb{R}$ satisfies the Polyak-Lojasiewicz (PL) condition with parameter $\mu > 0$ if

$$\frac{1}{2}\|\nabla \ell(\theta)\|^2 \geq \mu(\ell(\theta) - \ell^*), \quad \forall \theta.$$

Lemma 5.5. Let ℓ be differentiable, L -smooth, and μ -PL. Then, gradient descent with stepsize $\eta = 1/L$ converges at a geometric rate,

$$\ell(\theta_k) - \ell^* \leq \left(1 - \frac{\mu}{L}\right)^k (\ell(\theta_0) - \ell^*).$$

Definition 5.6 (Strong convexity). A differentiable function $\ell : \mathbb{R}^d \rightarrow \mathbb{R}$ is μ -strongly convex for some $\mu > 0$ if

$$\ell(\theta') \geq \ell(\theta) + \langle \nabla \ell(\theta), \theta' - \theta \rangle + \frac{\mu}{2}\|\theta' - \theta\|^2, \quad \forall \theta, \theta'.$$



Figure 5.1. Graph of the gradient of an L -smooth 1-dimensional function. Due to smoothness, the gradient may only stay within the grey area for every point. Higher L makes the area larger, which means that the current gradient is less informative about its surroundings. Thus, we have to be more careful when applying gradient descent.

Intuitively, strong convexity bounds the smallest eigenvalue of a locally quadratic approximation of ℓ .

Lemma 5.7 (Strong convexity \Rightarrow PL). Let ℓ be μ -strongly convex, then it fulfills the PL condition with the same μ .

Thus, strong convexity ensures convergence to a global optimum.

Momentum. The heavy ball method is an optimization algorithm with momentum,

$$\theta_{k+1} = \theta_k - \eta \nabla \ell(\theta_k) + \beta(\theta_k - \theta_{k-1}), \quad \beta \in (0, 1).$$

Intuitively, we are acting as if the iterates have mass. As a consequence, we will pass small gradient areas faster, and thus overcomes converging to local minima.

Adaptivity. AdaGrad uses the gradient history to adapt the stepsize per dimension,

$$[\theta_{k+1}]_i = [\theta_k]_i - \eta_i^k \nabla_i \ell(\theta_k), \quad \eta_i^k \doteq \frac{\eta}{\sqrt{\gamma_i^k + \delta}},$$

where

$$\gamma_i^k = \gamma_i^{k-1} + (\nabla_i \ell(\theta_k))^2.$$

Acceleration. Nesterov's accelerated gradient descent has the following update rule,

$$\begin{aligned} \theta'_{k+1} &= \theta_k + \beta(\theta_k - \theta_{k-1}) \\ \theta_{k+1} &= \theta'_{k+1} - \eta \nabla \ell(\theta'_{k+1}). \end{aligned}$$

While it is not intuitive why this works, it provides a faster convergence rate than vanilla gradient descent.

Adam. Adam combines momentum and adaptivity to increase convergence speed. It defines the following variables,

$$\begin{aligned} g_k &= \beta g_{k-1} + (1 - \beta) \nabla \ell(\theta_k), \quad \beta \in [0, 1] \\ h_k &= \alpha h_{k-1} + (1 - \alpha) \nabla \ell(\theta_k)^{\odot 2}, \quad \alpha \in [0, 1]. \end{aligned}$$

The update rule is then

$$\theta_{k+1} = \theta_k - \eta_k \odot g_k, \quad \eta_k = \eta \oslash (\sqrt{h_k} + \delta).$$

5.3 Convolutional neural networks

Images and audio have an extremely high dimensionality, which means that a fully connected layer would have an extremely high parameter count. However, we can exploit the locality, scale, and shift invariance of these types of data using a new operator; the convolution.

Definition 5.8 (Convolution). Given two functions f, h , their convolution is defined as

$$(f * h)(u) \doteq \int_{-\infty}^{\infty} h(u-t)f(t)dt = \int_{-\infty}^{\infty} f(u-t)h(t)dt.$$

The convolution operator is shift invariant, meaning that if we shift and then apply the operator, we get the same result as if we were to first apply the operator and then shift. Formally,

$$f_{\Delta} * h = (f * h)_{\Delta}$$

The converse is also true: any linear shift-invariant transformation can be written as a convolution.

In practice, we have discrete data, which means that we have to define a discrete convolution operator,

$$(f * h)[u] \doteq \sum_{t=-\infty}^{\infty} f[t]h[u-t].$$

This easily extends to two dimensions,

$$(f * h)[x, y] = \sum_{u=-\infty}^{\infty} \sum_{v=-\infty}^{\infty} f[u, v]h[x-u, y-v].$$

Typically, f and h are defined over a finite domain.

The cross-correlation operator is equal to a convolution with a flipped kernel,

$$(f \star h)[u] \doteq \sum_{t=-\infty}^{\infty} h[t]f[u+t].$$

Convolutional neural networks (CNN) learn the kernel of convolutional layers and stacks them in a compositional way to extract features from images. In this way, these layers exploit the shift invariance, locality, and scale of the data. Furthermore, it increases the statistical efficiency w.r.t. MLPs, because of the shared parameters.

As MLPs, CNNs alternate between (linear) convolutional layers and non-linear element-wise functions to increase model capacity. Moreover, it employs max-pooling layers to decrease the dimensionality of the input. It does so by taking only the maximum of every $k \times k$ patch. This has the effect that the input is subsampled k times. After many max-pooling layers, the data will no longer be location dependent, which allows us to throw away spatial information by flattening the data. From there, linear layers can be used to make the final prediction.

6 Generative models

Generally, a dataset has a true generative process $p(x)$ that we can sample from by querying the dataset. Considering that collecting data is often expensive, we wish to parametrize a generative model $p(x; \theta)$ that is indistinguishable from $p(x)$.

We want the models to not generate unnatural patterns, so we do not want the following,

$$p(x; \theta) \gg p(x) \approx 0, \quad \exists x.$$

Furthermore, we want the full domain of the distribution to be covered by our generative model. Thus, we also do not want the following to happen,

$$p(x) \gg p(x; \theta) \approx 0, \quad \exists x.$$

Moreover, we want to train generative models such that they give high probabilistic mass to samples from the dataset. Thus, a good objective function is to maximize the log-likelihood of the dataset (maximum likelihood estimation),

$$\ell(\theta) = \mathbb{E}_{x \sim p(x)} [\log p(x; \theta)].$$

However, we will see that some models cannot tractably compute the likelihood and must thus resort to lower bounds. We will also see that the likelihood might not be the best measure of success, which will lead to alternative loss functions.

6.1 Autoregressive models

Autoregressive models make use of a sequential ordering of variables and define the model directly on these observables. There are data types that have an inherent temporal ordering, such as text and audio, while other data types require defining an ordering. For example, in order to generate images pixel-by-pixel, a consistent ordering of the pixels must be defined. Empirically, it has been found that any ordering, even random, works well.

In order to compute the likelihood, autoregressive models make use of the chain rule,

$$p(x; \theta) = \sum_{t=1}^m p(x_t \mid x_1, \dots, x_{t-1}; \theta).$$

Typically, autoregressive models are modeled by modeling the conditional distribution $p(x_t \mid x_{1:t-1}; \theta)$ as a neural network. Specifically, the history is first mapped to a latent variable (embedding),

$$x_{1:t-1} \mapsto z \in \mathbb{R}^d,$$

which is used to predict the probability distribution over next observables x_t by applying the softmax operator,

$$p(x \mid x_{1:t-1}; \theta) = \frac{\exp\langle z, w_x \rangle}{\sum_{x' \in \mathcal{X}} \exp\langle z, w_{x'} \rangle}.$$

We sample from $p(x; \theta)$ by iteratively sampling $p(x_t \mid x_{1:t-1}; \theta)$ and appending it to the history.

Pixel CNN. Pixel CNN is used to generate $n \times m$ images,

$$p(\mathbf{x}) = \prod_{t=1}^{n \times m} p(x_t \mid x_1, \dots, x_{t-1}).$$

However, it only conditions on the pixels within a predefined context window. Since it is only allowed to use information about pixels above and to the left of the current pixel, we have to use a masked kernel window that zeroes out all pixels below and to the left of the current pixel. This enables the model to learn much faster than other autoregressive models, because it does not have to wait for all previous pixels to be generated, because we already have access to them, and the latent representation depends only on the pixels within the context window. But, generation is still slow, because all pixels in the context window need to be generated first.

Transformers. Recently, transformers have been used to power state-of-the-art language models. They make use of multi-headed self-attention (MHSA), layer normalization, and an MLP. Let $\mathbf{X} \in \mathbb{R}^{T \times d}$ contain the T d -dimensional embeddings in the current context window. MHSA introduces a query matrix $\mathbf{W}_Q \in \mathbb{R}^{d \times k}$, a key matrix $\mathbf{W}_K \in \mathbb{R}^{d \times k}$, and a value matrix $\mathbf{W}_V \in \mathbb{R}^{d \times r}$. We multiply these with \mathbf{X} to compute query, key, and value representations for every timestep,

$$\mathbf{Q} = \mathbf{X}\mathbf{W}_Q \in \mathbb{R}^{T \times k}, \quad \mathbf{K} = \mathbf{X}\mathbf{W}_K \in \mathbb{R}^{T \times k}, \quad \mathbf{V} = \mathbf{X}\mathbf{W}_V \in \mathbb{R}^{T \times r}.$$

These interact like a soft mapping, where the queries and keys decide how much from the values should be taken,

$$\mathbf{\Xi} = \mathbf{A}\mathbf{V} \in \mathbb{R}^{T \times r}, \quad \mathbf{A} = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{k}}\right) \in \mathbb{R}^{T \times T},$$

where softmax is applied row-wise. Intuitively, the queries broadcast what information the timesteps want, the keys broadcast what information they have to be compatible with the queries, and the values contain the actual information. The outputs ξ_i are convex combinations of the values v_t , where the query q_i and keys k_t decide the attention weights a_{it} ,

$$\xi_i = \sum_{t=1}^T a_{it} v_t, \quad a_{it} \propto q_i^\top k_t.$$

Setting $d = k = r$, makes the output shape equal the input shape, which means that we can stack self-attention layers easily. Furthermore, we can add “heads” that are each smaller self-attention operations that are combined. Let h be the number of heads, then we first project X to h d/h -dimensional representations,

$$XW_i \in \mathbb{R}^{T \times d/h}, \quad W_i \in \mathbb{R}^{d \times d/h}, i \in [h].$$

Afterwards, we apply self-attention to each head. Finally, we concatenate them and project them to the output,

$$\Xi = \text{concat}(\text{head}_1, \dots, \text{head}_h)W_O \in \mathbb{R}^{T \times d}, \quad W_O \in \mathbb{R}^{h \cdot r \times d}.$$

These can efficiently be computed by parallelizing the self-attention operators of the heads.

In order to satisfy the autoregressive property, we need to mask out future timesteps. This is easily done by adding a mask matrix to the attention matrix, resulting in

$$A = \text{softmax}\left(\frac{QK^\top}{\sqrt{k}} + M\right), \quad M = \begin{bmatrix} -\infty & -\infty & \dots & -\infty \\ 0 & -\infty & \dots & -\infty \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & -\infty \end{bmatrix}.$$

The result is that $a_{it} = 0$ for $t \geq i$, so the i -th timestep is only a function of the timesteps before it.

Transformers also make use of layer normalization, which normalizes the data by

$$x' = \gamma \left(\frac{x - \mu}{\sigma} \right) + \beta, \quad \mu = \frac{1}{d} \sum_{i=1}^d x_i, \quad \sigma^2 = \frac{1}{d} \sum_{i=1}^d (x_i - \mu)^2,$$

where $\gamma, \beta \in \mathbb{R}$ are parameters. This normalization block is applied in order to balance activation magnitudes across tokens and dimensions.

6.2 Variational autoencoders

The variational autoencoder (VAE) is a latent variable model, where x is first mapped to an instantiation of an elementary distribution $p(z | x)$ by an encoder, from which we sample the latent variable z . Then, a decoder reconstructs the original input \hat{x} from this latent variable z . During training, we make sure that the output distributions of the encoder do not deviate too far from a prior distribution $p(z)$. After training, we can then generate new data points by sampling z from the prior and passing it through the decoder.

The log-likelihood is computed by

$$\log p(x; \theta) = \log \int p(z)p(x | z; \theta) dz,$$



Figure 6.1. Layer in the transformer architecture.

which is intractable, because we cannot integrate over all latent variables. Thus, we must derive an ELBO using a variational distribution q ,

$$\begin{aligned}
\log p(\mathbf{x}; \boldsymbol{\theta}) &= \log \int p(\mathbf{z}) p(\mathbf{x} | \mathbf{z}; \boldsymbol{\theta}) d\mathbf{z} \\
&= \log \int q(\mathbf{z}) \frac{p(\mathbf{z})}{q(\mathbf{z})} p(\mathbf{x} | \mathbf{z}; \boldsymbol{\theta}) d\mathbf{z} \\
&\geq \int q(\mathbf{z}) \log \frac{p(\mathbf{z})}{q(\mathbf{z})} p(\mathbf{x} | \mathbf{z}; \boldsymbol{\theta}) d\mathbf{z} \\
&= \int q(\mathbf{z}) \log p(\mathbf{x} | \mathbf{z}; \boldsymbol{\theta}) d\mathbf{z} - \int q(\mathbf{z}) \log \frac{q(\mathbf{z})}{p(\mathbf{z})} d\mathbf{z} \\
&= \mathbb{E}_{q(\mathbf{z})} [\log p(\mathbf{x} | \mathbf{z}; \boldsymbol{\theta})] - D_{\text{KL}}(q(\mathbf{z}) \| p(\mathbf{z}))
\end{aligned}$$

We parameterize q as the encoder with parameters $\boldsymbol{\psi}$ (instead of choosing q for every \mathbf{x} like in the case of mixture models), which results in the final loss function,

$$\ell(\boldsymbol{\theta}, \boldsymbol{\psi}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x}; \boldsymbol{\psi})} [\log p(\mathbf{x} | \mathbf{z}; \boldsymbol{\theta})] - D_{\text{KL}}(q(\mathbf{z} | \mathbf{x}; \boldsymbol{\psi}) \| p(\mathbf{z})).$$

Thus, the ELBO encourages good reconstruction by the first term and makes sure that the distributions produced by the encoder do not deviate too far from the prior by the second term. However, we can make a Monte Carlo approximation of the expectation by using the reparametrization trick. This involves sampling an instance from a fixed distribution and using this to compute an instance of the desired distribution with the same probability. In the case of q as a Gaussian, where the encoder outputs means $\boldsymbol{\mu}$ and variances $\boldsymbol{\sigma}$, we can obtain a sample from this distribution by only sampling from $\mathcal{N}(\mathbf{0}, \mathbf{I})$ as follows,

$$\begin{aligned}
\boldsymbol{\epsilon} &\sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \\
\mathbf{z} &= \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\epsilon}.
\end{aligned}$$

\mathbf{z} is a sample from $\mathcal{N}(\boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2))$. The second term of the ELBO is analytically computable if we further define the prior to be Gaussian $\mathcal{N}(\mathbf{0}, \mathbf{I})$.

6.3 Generative adversarial networks

In practice, the log-likelihood may not be the best to optimize for. The next idea views generative models in terms of a classification problem,

$$p(\mathbf{x}, Y; \boldsymbol{\theta}) = Y p(\mathbf{x}) + (1 - Y) p(\mathbf{x}; \boldsymbol{\theta}), \quad Y \in \{0, 1\},$$

where Y indicates from which distribution \mathbf{x} is sampled.

Assume that we have access to a classifier (or discriminator) that can distinguish between samples from the true data distribution and samples from the model's distribution,

$$\pi : \mathbb{R}^n \rightarrow [0, 1], \quad \pi(\mathbf{x}) \approx \mathbb{P}(Y = 1 | \mathbf{x}).$$

We then have the following logistic loss,

$$\begin{aligned}\ell(\theta; \pi) &= \mathbb{E}[Y \log \pi(x) + (1 - Y) \log(1 - \pi(x))] \\ &= \int p(x) \log \pi(x) dx + \int p(x; \theta) \log(1 - \pi(x)) dx.\end{aligned}$$

The Bayes-optimal classifier for this loss is

$$\begin{aligned}\pi^*(x) &= \mathbb{P}(Y = 1 \mid x) \\ &= \frac{p(x \mid Y = 1)p(Y = 1)}{p(x \mid Y = 1)p(Y = 1) + p(x \mid Y = 0)p(Y = 0)} && \text{Bayes' rule.} \\ &= \frac{p(x) \cdot 1/2}{p(x) \cdot 1/2 + p(x; \theta) \cdot 1/2} && \text{Set equal priors for both cases.} \\ &= \frac{p(x)}{p(x) + p(x; \theta)}.\end{aligned}$$

Furthermore, by using the Bayes-optimal classifier, we can rewrite the logistic loss function as optimizing the Jensen-Shannon divergence between $p(x)$ and $p(x; \theta)$,

$$\begin{aligned}\ell(\theta) &= \int p(x) \log \frac{p(x)}{p(x) + p(x; \theta)} dx + \int p(x; \theta) \log \frac{p(x; \theta)}{p(x) + p(x; \theta)} dx && \text{Insert Bayes-optimal classifier.} \\ &= \mathbb{E}_{p(x)} \left[\log \frac{p(x)}{p(x) + p(x; \theta)} \right] + \mathbb{E}_{p(x; \theta)} \left[\log \frac{p(x; \theta)}{p(x) + p(x; \theta)} \right] && \text{Rewrite as expectations.} \\ &= \mathbb{E}_{p(x)} \left[\log \frac{2 \cdot p(x)}{2 \cdot (p(x) + p(x; \theta))} \right] + \mathbb{E}_{p(x; \theta)} \left[\log \frac{2 \cdot p(x; \theta)}{2 \cdot (p(x) + p(x; \theta))} \right] && \text{Multiply by } 2/2 \text{ within the logarithms on both sides.} \\ &= \mathbb{E}_{p(x)} \left[\log \frac{2 \cdot p(x)}{p(x) + p(x; \theta)} \right] + \mathbb{E}_{p(x; \theta)} \left[\log \frac{2 \cdot p(x; \theta)}{p(x) + p(x; \theta)} \right] - \log 4 && \text{Take out the } -\log 2 \text{ of both expectations.} \\ &= D_{\text{KL}} \left(p(x) \left\| \frac{p(x) + p(x; \theta)}{2} \right\| \right) + D_{\text{KL}} \left(p(x; \theta) \left\| \frac{p(x) + p(x; \theta)}{2} \right\| \right) - \log 4 && D_{\text{KL}}(p \parallel q) = \mathbb{E}_p \left[\log \frac{p(x)}{q(x)} \right]. \\ &= 2D_{\text{JS}}(p(x) \parallel p(x; \theta)) - \log 4. && D_{\text{JS}}(p \parallel q) = \frac{1}{2} D_{\text{KL}}(p \parallel \frac{p+q}{2}) + \frac{1}{2} D_{\text{KL}}(q \parallel \frac{p+q}{2}).\end{aligned}$$

In practice, we do not have access to $p(x)$ and $p(x; \theta)$ to compute the optimal discriminator π^* . Despite this, it is not a hard job, so we can expect a neural network, parametrized by ϕ , to perform well on this task. This results in an adversarial loss function,

$$\theta^* = \underset{\theta}{\operatorname{argmin}} \max_{\phi} \ell(\theta, \phi), \quad \phi^* = \underset{\phi}{\operatorname{argmax}} \ell(\theta^*, \phi).$$

This can be optimized by gradient descent where every step involves optimizing both “players”. However, that is often non-converging. The extragradient method usually performs better.

6.4 Diffusion models

Diffusion models represent the new state-of-the-art in generative models. Instead of generating a full sample in one forward pass, it breaks it down into many small iterative steps. It does so by starting from pure noise,

$\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, and iteratively working toward a sample from the data distribution, $\mathbf{x}_0 \sim p(\mathbf{x})$.

The forward diffusion process forms a Markov chain,

$$\mathbf{x}_t = \alpha_t \mathbf{x}_{t-1} + \beta_t \boldsymbol{\epsilon}_t, \quad \boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}),$$

where $\alpha_t, \beta_t \in (0, 1)$. Typically, α_t and β_t are chosen to preserve signal variance, *i.e.*,

$$\alpha_t^2 + \beta_t^2 = 1, \quad \beta_t = \sqrt{1 - \alpha_t^2}.$$

Given \mathbf{x}_0 , we can directly compute \mathbf{x}_t for any t ,

$$\mathbf{x}_t = \bar{\alpha}_t \mathbf{x}_0 + \bar{\beta}_t \boldsymbol{\epsilon}_t, \quad \boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}),$$

where

$$\bar{\alpha}_t = \prod_{\tau=1}^t \alpha_\tau, \quad \bar{\beta}_t = \sqrt{1 - \bar{\alpha}_t^2}.$$

Because $\lim_{t \rightarrow \infty} \bar{\alpha}_t = 0$, we iteratively remove information in the diffusion process. In the limit, all information of \mathbf{x}_0 is lost,

$$\mathbf{x}_\infty \mid \mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I}).$$

The idea of diffusion models is to learn the time-reversed Markov chain,

$$p(\mathbf{x}_{t-1} \mid \mathbf{x}_t; \boldsymbol{\theta}).$$

Provided that $\lim_{t \rightarrow \infty} \beta_t = 0$, then we can approximate this distribution by a Gaussian,

$$\mathbf{x}_{t-1} \mid \mathbf{x}_t \sim \mathcal{N}(\boldsymbol{\mu}(\mathbf{x}_t, t; \boldsymbol{\theta}), \boldsymbol{\Sigma}(\mathbf{x}_t, t; \boldsymbol{\theta})),$$

which we can simplify to only learn the mean,

$$\mathbf{x}_{t-1} \mid \mathbf{x}_t \sim \mathcal{N}(\boldsymbol{\mu}(\mathbf{x}_t, t; \boldsymbol{\theta}), \sigma_t^2 \mathbf{I}), \quad \sigma_t^2 \in \{\beta_t, \bar{\beta}_t\}.$$

Optimizing the log-likelihood with this model reduces to optimizing the squared error criterion between \mathbf{x}_{t-1} and $\boldsymbol{\mu}$, which is easy to optimize. Because the next step \mathbf{x}_{t-1} can be computed from the noise $\boldsymbol{\epsilon}_t$, this is effectively the same as training a model to output the added noise and optimizing w.r.t. the squared error on the noise. The learning algorithm can be found in Algorithm 1 and the sampling algorithm can be found in Algorithm 2. For image generation, $\boldsymbol{\epsilon}(\mathbf{x}_t, t; \boldsymbol{\theta})$ is typically modeled as a U-net.

Require: $\{\beta_t\}_{t=1}^T$

```

1: while not converged do
2:    $x_0 \sim q$ 
3:    $t \sim \text{Unif}([T])$ 
4:    $\epsilon \sim \mathcal{N}(\mathbf{0}, I)$ 
5:    $\ell(\theta) = \frac{1}{2} \|\epsilon - \epsilon(\bar{\alpha}_t x_0 + \bar{\beta}_t \epsilon, t; \theta)\|^2$ 
6:    $\theta \leftarrow \theta - \eta \nabla \ell(\theta)$ 
7: end while

```

Algorithm 1. Diffusion model training algorithm.

Require: $\{\beta_t\}_{t=1}^T, \theta$

```

1:  $x_T \sim \mathcal{N}(\mathbf{0}, I)$ 
2: for  $t = T, \dots, 1$  do
3:    $z \sim \mathcal{N}(\mathbf{0}, I)$  if  $t > 1$  else  $z = \mathbf{0}$ 
4:    $x_{t-1} = \frac{1}{\alpha_t} \left( x_t - \frac{\beta_t^2}{\beta_t} \epsilon(x_t, t; \theta) \right) + \sigma_t z$ 
5: end for
6: return  $x_0$ 

```

Algorithm 2. Diffusion model sampling algorithm.