

# NUMERICAL ANALYSIS FOR MACHINE LEARNING

Lorenzo Bozzoni

October 27, 2023

# Contents

<b>1 Basic concepts of linear algebra</b>	<b>3</b>
<b>2 Matrix-vector multiplication</b>	<b>3</b>
2.1 Row-reduced echelon form . . . . .	4
<b>3 Matrix-matrix multiplication</b>	<b>4</b>
<b>4 Factorizations</b>	<b>4</b>
4.1 Orthogonal matrices . . . . .	5
4.1.1 Rotation . . . . .	5
4.1.2 Reflection . . . . .	5
<b>5 Null spaces</b>	<b>6</b>
<b>6 Null space cardinality</b>	<b>7</b>
<b>7 Eigenvalues and eigenvectors</b>	<b>8</b>
7.1 Eigenvectors of matrix power . . . . .	8
7.2 Power method . . . . .	8
7.3 Similar matrices . . . . .	8
7.4 QR factorization . . . . .	9
7.4.1 QR iteration . . . . .	9
7.5 Positive-definite symmetric matrices (SPD) . . . . .	10
<b>8 Singular Value Decomposition (SVD)</b>	<b>11</b>
8.1 Economy SVD . . . . .	12
8.2 Proof of the existence of SVD . . . . .	13
8.3 Geometrical interpretation of SVD . . . . .	15
8.3.1 Properties of SVD . . . . .	15
8.4 Snapshots method . . . . .	15
8.5 Matrix norms . . . . .	16
8.6 Eckart-Young theorem . . . . .	16
8.6.1 Proof considering $\ \cdot\ _F$ . . . . .	17
8.6.2 Proof considering $\ \cdot\ _2$ . . . . .	18
<b>9 PCA</b>	<b>19</b>
9.1 Choose rank of truncated SVD . . . . .	19
9.2 Randomize SVD . . . . .	20
<b>10 Least squares approximation</b>	<b>21</b>
10.1 Geometrical interpretation . . . . .	21
10.2 Optimization . . . . .	22
<b>11 Matrix completion</b>	<b>25</b>
11.1 Ideal estimator . . . . .	26
11.2 Practical estimator . . . . .	26
11.3 Ridge regression (regularization) . . . . .	27
<b>12 Page Rank</b>	<b>28</b>
<b>13 Kernel Methods</b>	<b>31</b>
<b>14 Computing derivatives</b>	<b>34</b>
14.1 Dual-numbers . . . . .	36

# 1 Basic concepts of linear algebra

The following are the main concepts of linear algebra we are going to face during the starting phase of the course:

1. Linear systems of equations:  $A\underline{x} = \underline{b}$
2. Eigenvalues and eigenvectors:  $A\underline{x} = \lambda\underline{x}$
3. Singular value decomposition (SVD):  $A\underline{v} = \sigma\underline{u}$
4. Minimization problem
5. Factorization:  $PA = LU$

## 2 Matrix-vector multiplication

$$\underline{c} = \underbrace{\begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix}}_{A_1} \underbrace{\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}}_{\underline{x}} = \begin{bmatrix} 1x_1 + 2x_2 \\ 3x_1 + 4x_2 \\ 5x_1 + 6x_2 \end{bmatrix} = \underbrace{\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}}_{\text{linear combination}} x_1 + \underbrace{\begin{bmatrix} 4 \\ 5 \\ 6 \end{bmatrix}}_{\text{linear combination}} x_2$$

We say that the vector  $\underline{c}$  belongs to the **column space** of  $A_1$ , i.e.  $\underline{c} \in \mathcal{C}(A_1)$ .

$$\underbrace{\begin{bmatrix} 1 & 4 & 7 \\ 2 & 5 & 8 \\ 3 & 6 & 9 \end{bmatrix}}_{\substack{A_2 \\ \underline{a_1} \quad \underline{a_2} \quad \underline{a_3}}} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

In this case we can easily notice that  $\underline{a_3} = \underline{a_1} + \underline{a_2}$ , which means that one column can be expressed as a linear combination of the other two (this means that the matrix  $A_2$  is singular). Because of this, we can say that  $\mathcal{C}(A_2) = \mathcal{C}(A_1)$ , i.e. the column space of  $A_2$  is the same as the column space of  $A_1$ .

Those columns spaces are a plane passing through the origin and spanned by the two vectors  $\underline{a_1}$  and  $\underline{a_2}$  (they define the slope of that plane).

Let's now consider these matrix:

$$\underbrace{\begin{bmatrix} 1 & 4 & 7 \\ 2 & 5 & 8 \\ 3 & 6 & 10 \end{bmatrix}}_{A_3} \qquad \underbrace{\begin{bmatrix} 1 & 2 & 3 \\ 2 & 4 & 6 \\ 3 & 6 & 9 \end{bmatrix}}_{A_4}$$

This left-hand matrix column space is  $\mathcal{C}(A_3) = \mathbb{R}^3$ , i.e. the entire real space of three dimensions. This is because the three vector columns of  $A_3$  are linearly independent so they span the entire space and not just a plane. While the column space of  $A_4$  is instead:  $\mathcal{C}(A_4) = [1 \ 2 \ 3]^T$  i.e. just a line since the three columns are linearly dependent and so they lie on the same line (they are parallel) just with different magnitude.

Another measure regarding matrices is the **dimension** or **rank**:

- $rank(A_1) = 2$
- $rank(A_2) = 2$
- $rank(A_3) = 3$
- $rank(A_4) = 1$

The rank is the number of linearly independent columns (or rows) of a matrix.

Let's consider again the matrix  $A_3$ :

$$\begin{bmatrix} 1 & 4 & 7 \\ 2 & 5 & 8 \\ 3 & 6 & 10 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \underbrace{\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} x_1 + \begin{bmatrix} 4 \\ 5 \\ 6 \end{bmatrix} x_2 + \begin{bmatrix} 7 \\ 8 \\ 10 \end{bmatrix} x_3}_{\in \mathcal{C}(A_3)} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}$$

so  $\underline{b}$  must be in  $\mathcal{C}(A_3)$  in order to have the system solvable. If  $\underline{b} \notin \mathcal{C}(A_3)$ , then the system is not solvable. In this particular case we have that the columns of  $A_3$  are linearly independent, so the system is solvable for any  $\underline{b}$  because

$\mathcal{C}(A_3) = \mathbb{R}^3$  and so  $\underline{b}$  is for sure inside that space.

Given  $A$ , find  $\mathcal{C}(A)$ . How can we solve this problem? Considering  $\underline{a}_1, \dots, \underline{a}_n$  as  $A$  columns, we can use the following iterative algorithm:

- put  $\underline{a}_1$  in  $\mathcal{C}(A)$
- if  $\underline{a}_2 = \alpha \underline{a}_1 \rightarrow \underline{a}_2 \notin \mathcal{C}(A)$ , otherwise put  $\underline{a}_2$  in  $\mathcal{C}(A)$

Until you reach the last column.

## 2.1 Row-reduced echelon form

Given the matrix  $A$ , defined as follow:

$$A = \begin{bmatrix} 1 & 4 & 7 \\ 2 & 5 & 8 \\ 3 & 6 & 9 \end{bmatrix}$$

we can obtain the **row-reduced echelon form** of  $A$  by applying the following operations:

$$A = CR = \begin{bmatrix} 1 & 4 \\ 2 & 5 \\ 3 & 6 \end{bmatrix} \begin{bmatrix} 1 & 0 & -1 \\ 0 & 1 & 2 \end{bmatrix}$$

where  $C$  is the matrix containing the columns of  $A$  that are linearly independent (i.e.  $\mathcal{C}(A)$ ) and  $R$  is the matrix of the coefficients of the linear combination of the columns of  $A$  that gives the columns of  $C$ .

Let's now consider the following matrix:

$$A_1 = \begin{bmatrix} 1 & 4 \\ 2 & 5 \\ 3 & 6 \end{bmatrix} \quad A_1^\top = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}$$

What we can say about  $A_1^\top$  column space? Is there any relationship with the column space of  $A_1$ ?

In order to compute its column space, we can start noticing that:  $\underline{a}_3 = 2\underline{a}_2 - \underline{a}_1$ . So, in general, we can say that:

$$\dim(\mathcal{C}(A)) = \dim(\mathcal{C}(A^\top)) = r \leq n \quad \text{where } n \text{ is the number of columns of } A$$

## 3 Matrix-matrix multiplication

$$C = AB = \begin{bmatrix} | & | & | \\ \underline{a}_1 & \dots & \underline{a}_n \\ | & | & | \end{bmatrix} \begin{bmatrix} - & \underline{b}_1 & - \\ - & \vdots & - \\ - & \underline{b}_n & - \end{bmatrix} = \overset{\text{col} \downarrow \text{row} \downarrow}{\underline{a}_1 \underline{b}_1} + \dots + \underline{a}_n \underline{b}_n$$

All the products that are summed at the end of the equation are matrices of rank 1.

Example:

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} 2 & 1 \\ 2 & 3 \end{bmatrix} = \underbrace{\begin{bmatrix} 2 & 1 \\ 6 & 3 \end{bmatrix}}_{\text{rank} = 1} + \underbrace{\begin{bmatrix} 4 & 6 \\ 8 & 12 \end{bmatrix}}_{\text{rank} = 1} = \begin{bmatrix} 6 & 7 \\ 14 & 15 \end{bmatrix}$$

## 4 Factorizations

1.  $A = LU$  or  $PA = LU$
2.  $A = QR$  where  $Q$  is orthogonal and  $R$  is upper triangular This is an improved version of the Row-reduced echelon form because that worked only for square matrices, while this works for any matrix.
3. Eigenvalues and eigenvectors decomposition: when  $S = S^\top$  (symmetric matrix) we can factorize it as  $S = Q\Lambda Q^\top$  where  $\Lambda$  is a diagonal matrix and  $Q$  is an orthogonal matrix (they are all squared matrices)
4. Generalization of the above:  $A = X\Lambda X^{-1}$  where  $X$  is a non-orthogonal matrix
5.  $A = U\Sigma V^\top$  where  $U$  and  $V$  are orthogonal matrices and  $\Sigma$  is a pseudo-diagonal matrix

A matrix is said to be pseudo-diagonal if it has the following form:

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & 0 & 0 \\ 0 & \sigma_2 & 0 & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & 0 & \sigma_n \\ 0 & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 0 \end{bmatrix} \quad m \text{ rows} \times n \text{ columns}$$

So it has diagonal elements for the first  $n$  rows then it has all zeros.

## 4.1 Orthogonal matrices

A matrix  $Q$  is orthogonal if  $Q^\top Q = I$  (i.e.  $Q^\top = Q^{-1}$ ). This means that the columns of  $Q$  are orthonormal, i.e. they are orthogonal and have unit norm.

The determinant of a orthogonal matrix is  $\pm 1$ .

Properties:

- $\|Q\underline{x}\| = \|\underline{x}\|$
- $\|Q\underline{x}\|^2 = (Q\underline{x})^\top Q\underline{x} = \underbrace{\underline{x}^\top Q^\top Q \underline{x}}_I = \|\underline{x}\|^2$

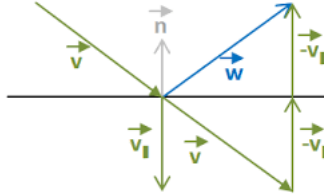
The first property is particularly easy to interpret since it means that when we multiply an orthogonal matrix to a vector, the norm of the vector doesn't change. As a proof of this, we can consider the following examples:

### 4.1.1 Rotation

A classical rotation matrix is:

$$\begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix}$$

### 4.1.2 Reflection



The horizontal line in the figure represent a plane  $\pi$  while  $\underline{n}$  is its normal vector of length 1. Given  $\underline{v}$  to obtain  $\underline{w}$  we can use the following formula:

$$\underline{w} = \underline{v} - 2(\underline{v}^\top \underline{n})\underline{n} = \underbrace{(I - 2\underline{n}\underline{n}^\top)}_{\text{reflection matrix}} \underline{v}$$

Moreover, the reflection matrix  $R$  is not only orthogonal, but also the inverse of itself, i.e.  $R^{-1} = R^\top$ . This makes sense because if we apply the reflection matrix twice, we obtain the original vector  $\underline{v}$ , i.e. the reflection of the reflection is the starting vector.

If we didn't have the 2 in the formula, we would obtain the projection of  $\underline{v}$  on the plane  $\pi$  which is called orthogonal projection and the matrix  $R$  would be singular.

Let's now dive a bit into the third point of the factorization list. We said that when  $S = S^\top$  (symmetric matrix) we can factorize it as  $S = Q\Lambda Q^\top$  where  $\Lambda$  is a diagonal matrix and  $Q$  is an orthogonal matrix.

$$S = S^\top = \underbrace{(Q\Lambda)}_{\tilde{Q}} Q^\top = \tilde{Q} Q^\top$$

$$\tilde{Q} = \underline{q}_1 \lambda_1 + \cdots + \underline{q}_n \lambda_n$$

Where the  $q$  vectors are columns and  $\lambda$  vectors are rows. So we can reformulate:

$$S = (\underline{q_1}\lambda_1 + \dots + \underline{q_n}\lambda_n)Q^\top = \underline{q_1}\lambda_1\underline{q_1}^\top + \dots + \underline{q_n}\lambda_n\underline{q_n}^\top$$

This is called **spectral decomposition** of matrix  $S$  and  $q_1, \dots, q_n$  are the eigenvectors of  $S$  while  $\lambda_1, \dots, \lambda_n$  are the eigenvalues of  $S$ .

$$S\underline{q_1} = \lambda_1\underline{q_1} = (\underline{q_1}\lambda_1\underline{q_1}^\top + \dots + \underline{q_n}\lambda_n\underline{q_n}^\top)\underline{q_1} = \lambda_1\underline{q_1}(\underline{q_1}^\top\underline{q_1})$$

All the other products are null since the vector  $\underline{q_1}$  is orthogonal to all the other vectors  $\underline{q_i}$  for  $i \neq 1$  (recall that they are eigenvectors).

## 5 Null spaces

Let's consider the starting problem for a linear system of equations:

$$A\underline{x} = \underline{b} \quad \text{with} \quad A \in \mathbb{R}^{m \times n}, \text{rank}(A) = r$$

We are going to introduce 2 more spaces other than the column ones. To do so we consider:

$$A\underline{x} = \underline{0} \quad \rightarrow \quad N(A) \equiv \ker(A) = \{\underline{x} \in \mathbb{R}^n : A\underline{x} = \underline{0}\}$$

$$A^\top \underline{x} = \underline{0} \quad \rightarrow \quad N(A^\top) \equiv \ker(A^\top) = \{\underline{x} \in \mathbb{R}^n : A^\top \underline{x} = \underline{0}\}$$

So now, adding the so called **null spaces** we have that:

1.  $\mathcal{C}(A) \subset \mathbb{R}^m$  and  $\dim(\mathcal{C}(A)) = r$
2.  $\mathcal{C}(A^\top) \subset \mathbb{R}^n$  and  $\dim(\mathcal{C}(A^\top)) = r$
3.  $N(A) \subset \mathbb{R}^n$  and  $\dim(N(A)) = ?$
4.  $N(A^\top) \subset \mathbb{R}^m$  and  $\dim(N(A^\top)) = ?$

We still do not know the dimensions of those spaces.

### Example

$$A = \begin{bmatrix} 1 & 4 & 7 \\ 2 & 5 & 8 \\ 3 & 6 & 9 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \implies \begin{cases} x_1 + 4x_2 + 7x_3 = 0 \\ 2x_1 + 5x_2 + 8x_3 = 0 \\ 3x_1 + 6x_2 + 9x_3 = 0 \end{cases}$$

We compute the first equation

$$x_1 = -4x_2 - 7x_3 \implies \begin{cases} -3x_2 - 6x_3 = 0 \\ -6x_2 - 12x_3 = 0 \end{cases}$$

What is important to notice is that  $A$  has rank=2 so we have  $3 - 2 = 1$  **degrees of freedom**, i.e. we can choose one variable and the other two are automatically defined. This is visible in the last two equations of the system for example. In general, the degrees of freedom are given by  $n - r$  where  $n$  is the number of columns of  $A$  and  $r$  is the rank of  $A$ .

If we had 10 instead of 9 in  $A$  we would have had  $3 - 3 = 0$  degrees of freedom. This would translate in having the matrix  $A$  full rank and  $N(A) = \{\underline{0}\}$  so the only solution would be the null vector.

## 6 Null space cardinality

In the first lecture, we defined 4 spaces:  $N(A)$ ,  $N(A^\top)$ ,  $\mathcal{C}(A)$ ,  $\mathcal{C}(A^\top)$ . For the last two we defined also their cardinality whilst for the first ones we weren't able to tell yet. In this lecture we are going to find those values and prove them. In order to do so, we start from few useful properties:

1.  $\underline{x} = \underline{0} \in N(A)$  for any matrix  $A$
2. if  $\underline{x}, \underline{y} \in N(A) \implies A(\underline{x} + \underline{y}) = \underline{0}$
3. if  $\underline{x} \in N(A) \implies \alpha \underline{x}$  with  $\alpha \in \mathbb{R} \implies A(\alpha \underline{x}) = \underline{0}$

Consider, once again, the matrix  $A \in \mathbb{R}^{m \times n}$ ,  $\text{rank}(A) = r \leq n$ . We have seen the decomposition  $A = CR$ , where  $C$  contains the linearly independent columns of  $A$  and  $R$  contains the coefficients that allow to recover the columns of  $A$  starting from its independent columns. So, the matrix  $A$  can be rewritten as:

$$A = [A_1 \quad A_2] \quad A_1 \in \mathbb{R}^{m \times r} \quad A_2 \in \mathbb{R}^{m \times (n-r)}$$

Where  $A_1$  contains the independent columns of  $A$  and  $A_2$  the dependent ones. Example:

$$A = \begin{bmatrix} 1 & 4 & 7 \\ 2 & 5 & 8 \\ 3 & 6 & 9 \end{bmatrix} = \underbrace{\begin{bmatrix} 1 & 4 \\ 2 & 5 \\ 3 & 6 \end{bmatrix}}_{A_1} \underbrace{\begin{bmatrix} 1 & 0 & -1 \\ 0 & 1 & 2 \end{bmatrix}}_B$$

Since we have the last column of  $A$ , that is linearly dependent so it belongs to  $A_2$ , we can reformulate it in this way:

$$A = [A_1 \quad A_2] = [A_1 \quad A_1 B]$$

We build a new matrix  $K$  defined as follows:

$$K = \begin{bmatrix} -B \\ I_{n-r} \end{bmatrix} \quad K \in \mathbb{R}^{n \times (n-r)} \quad B \in \mathbb{R}^{r \times (n-r)}$$

$$AK = [A_1 \quad A_1 B] \begin{bmatrix} -B \\ I_{n-r} \end{bmatrix} = A_1(-B) + A_1 B = 0$$

Where the last 0 is actually a matrix of zeros of dimension  $m \times (n-r)$  because  $A$  has size  $m \times n$  and  $K$  has size  $n \times (n-r)$ . We have that:

$$AK = 0 \implies A \underline{k}_i = 0 \quad \forall i \in \{1, \dots, n-r\}$$

Where  $\underline{k}_i$  is the  $i$ -th column of  $K$ . This means that:  $\underline{k}_i \in N(A) \quad \forall i$ .

Now, we want to demonstrate that:  $K\underline{u} = 0 \implies \underline{u} = \underline{0}$ . To do so, we start from expanding  $K$  from its definition:

$$K = \begin{bmatrix} -B \\ I \end{bmatrix} \underline{u} = 0 \implies \begin{bmatrix} -B\underline{u} \\ \underline{u} \end{bmatrix} = \begin{bmatrix} \underline{0} \\ \underline{0} \end{bmatrix}$$

Where the two zero vectors have dimension  $r$  and  $n-r$  respectively! Considering the second row of the matrix we get:  $\underline{u} = \underline{0}$  so all columns of  $K$  are linearly independent.

If we consider the problem  $(\star) A\underline{x} = \underline{0}$ , we want to prove that each  $\underline{x}$  that satisfy  $(\star)$  must be a linear combination of the columns of  $K$ .

$$A_1 \underline{x} = \underline{0} \in \mathbb{R}^m \implies \underline{x} = \underline{0} \in \mathbb{R}^r$$

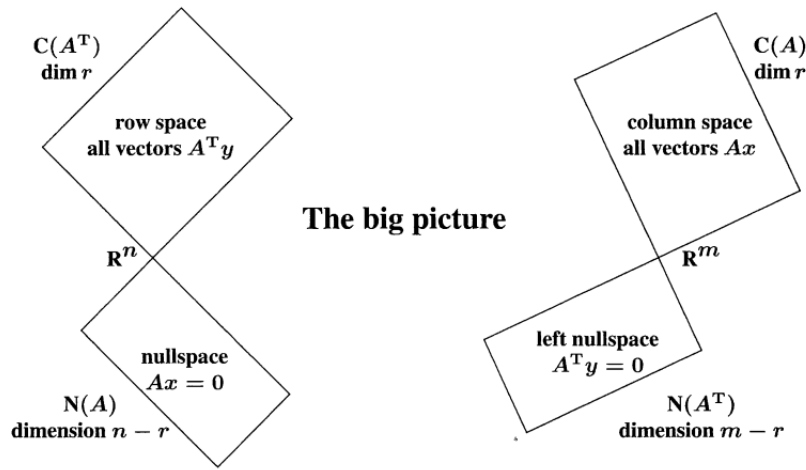
Because  $A_1$  has linearly independent columns, i.e. has full rank.

$$A\underline{u} = \underline{0} \in \mathbb{R}^m \implies [A_1 \quad A_1 B] \begin{bmatrix} \underline{u}_1 \\ \underline{u}_2 \end{bmatrix} = [A_1 \underline{u}_1 + A_1 B \underline{u}_2] = A_1 [\underline{u}_1 + B \underline{u}_2] = \underline{0}$$

We can notice that the last formulation obtained in the equation has the same form as the one from where we started the prove, so we can say that:

$$\underline{u}_1 + B \underline{u}_2 = \underline{0} \implies \underline{u}_1 = -B \underline{u}_2$$

$$\underline{u} = \begin{bmatrix} -B \underline{u}_2 \\ \underline{u}_2 \end{bmatrix} = \underbrace{\begin{bmatrix} -B \\ I \end{bmatrix}}_K \underline{u}_2 = K \underline{u}_2 \implies \dim(N(A)) = n - r$$



## 7 Eigenvalues and eigenvectors

Start considering a generic square matrix  $n \times n$ . We are going later to discuss even the symmetry and positive definite properties. Here below are the vectorial and the matrix form of the eigenvalue problem:

$$A\underline{x}_i = \lambda_i \underline{x}_i \quad i = 1, \dots, n \quad X^{-1}AX = \Lambda$$

Where in the right-hand side there is a diagonal matrix  $\Lambda$  with the eigenvalues of  $A$  on the diagonal while the matrix  $X$  with the eigenvectors of  $A$  as columns.

### 7.1 Eigenvectors of matrix power

What can we say about the eigenvectors and eigenvalues of  $A^2$ ?

$$A^2 \underline{x}_i = A(A\underline{x}_i) = A(\lambda_i \underline{x}_i) = \lambda_i (A\underline{x}_i) = \lambda_i^2 \underline{x}_i$$

So the eigenvalues of  $A^2$  are the eigenvalues of  $A$  squared. This is valid for any power of  $A$  since this method can be applied recursively and it is very useful when there are problems in which a matrix is iteratively multiplied many times.

#### Important:

Given  $A \in \mathbb{R}^{n \times n}$  full rank, then any vector  $\underline{v} \in \mathbb{R}^n$  can be written as a linear combination of the eigenvectors ( $\underline{x}_i$ ) of  $A$ .

### 7.2 Power method

In mathematics, power iteration (also known as the power method) is an eigenvalue algorithm: given a diagonalizable matrix  $A$ , the algorithm will produce a number  $\lambda$ , which is the greatest (in absolute value) eigenvalue of  $A$ , and a nonzero vector  $\underline{v}$ , which is a corresponding eigenvector of  $\lambda$ , that is,  $A\underline{v} = \lambda \underline{v}$ . The algorithm is also known as the Von Mises iteration.

There is also an **inverse PM** which is applied to  $A^{-1}$  to find the minimum eigenvalue or also **PM with a shift** applied to  $(A - \alpha I)^{-1}$ ,  $\alpha \in \mathbb{R}$  to find the closest eigenvalue to  $\alpha$ .

Can even be used in "deflation method" iteratively:

$$\begin{bmatrix} \lambda_1 & b_1^T \\ 0 & A_1 \end{bmatrix}$$

the original matrix could be reduced in that form and at every iteration the procedure is applied to the  $A_1$  matrix. This works only if we have different eigenvectors (or values) [check].

### 7.3 Similar matrices

Given two matrices  $A, B \in \mathbb{R}^n$ , they are said to be similar if  $B = M^{-1}AM$ , with  $M$  invertible.

$$\underbrace{M^{-1}AM}_{B} \underline{y} = \lambda \underline{y} \implies A \underbrace{M\underline{y}}_{\underline{w}} = \lambda \underbrace{M\underline{y}}_{\underline{w}} = A\underline{w} = \lambda \underline{w}$$



Where  $\lambda, \underline{y}$  contain respectively the eigenvalues and the eigenvectors of  $B$ . What we get from this equation is that **similar matrices share the same eigenvectors with scaled eigenvalues**.

## 7.4 QR factorization

Here is introduced in the context of eigenvalues. Let's consider a matrix  $A \in \mathbb{R}^{m \times n}$  where  $m \geq n$  and  $\text{rank}(A) = n$  (it has all independent columns). We can factorize  $A$  in this way:

$$A = QR \quad Q \in \mathbb{R}^{m \times n} \quad R \in \mathbb{R}^{n \times n}$$

Where  $Q$  is an orthogonal matrix and  $R$  is an upper triangular matrix. Since we are dealing with eigenvalues and eigenvectors, we are now going to consider the matrix  $A$  squared with the dimension  $n \times n$ .

### 7.4.1 QR iteration

$$A = A^{(0)} = Q^{(0)} R^{(0)}$$

$$A^{(1)} = Q^{(0)\top} A^{(0)} Q^{(0)} = Q^{(1)} R^{(1)}$$

So, iterating this procedure we get:

$$A^{(2)}, \dots, A^{(S)} = \text{is upper triangular}$$

After  $S$  iterations you obtain an upper triangular matrix. The matrices  $A, A^{(0)}, A^{(1)}, \dots, A^{(S)}$  are similar, so they share the same eigenvalues.

But, how can i compute  $Q$ ?

With the **Gram-Schmidt** procedure. It works also for non-square matrices.

Let's start from a generic matrix  $A$ :

$$A = \begin{bmatrix} | & | & | \\ \underline{a_1} & \dots & \underline{a_n} \\ | & | & | \end{bmatrix}$$

The algorithm is iterative and it is applied to the columns of  $A$  in such way:

$$\underline{q_1} = \frac{\underline{a_1}}{\|\underline{a_1}\|}$$

The vector  $\underline{q_1}$  is obtained by normalizing the first column of  $A$ , in such manner the new obtained vector will have norm 1.

$$\underline{q_2} = \underline{a_2} - \underline{q_1}(\underline{q_1}^\top \underline{a_2}) \implies \underline{q_2} = \frac{\underline{q_2}}{\|\underline{q_2}\|}$$

The second vector is obtained by subtracting from the second column of  $A$  the projection of  $\underline{a_2}$  on  $\underline{q_1}$ , in such manner the new vector will be orthogonal to  $\underline{q_1}$  and will have norm 1.

$$\underline{q_3} = \underline{a_3} - \underline{q_1}(\underline{q_1}^\top \underline{a_3}) - \underline{q_2}(\underline{q_2}^\top \underline{a_3}) \implies \underline{q_3} = \frac{\underline{q_3}}{\|\underline{q_3}\|}$$

And so on... Recall that the orthogonality is needed since we want to obtain an orthogonal matrix  $Q$  useful for the factorization. With Gram-Schmidt the resulting matrix not only will be orthogonal but also orthonormal, this means that its columns will have norm unitary.

Let's now continue with the factorization journey. We have said in the introduction of types of factorizations that, given  $A \in \mathbb{R}^{n \times n}$ , we have:

$$A = X \Lambda X^{-1}$$

Where  $X$  has as columns the eigenvectors of  $A$ , while  $\Lambda$  is a diagonal matrix with the eigenvalues of  $A$  on the diagonal. Now, let's consider the case where the matrix  $S$  is symmetric.

$$S \in \mathbb{R}^{n \times n} \quad S = S^\top$$

We can factorize  $S$  as follows:

$$S = Q \Lambda Q^\top$$

Where  $Q$  is orthogonal (this is true only because  $S$  is symmetric) and  $\Lambda$  is diagonal. We can prove that  $Q$  is orthogonal by:

1. Consider the two vectors  $\underline{x}, \underline{y}$  such as:  $S\underline{x} = \lambda\underline{x}$  and  $S\underline{y} = 0\underline{y}$ . So, we are saying that both vectors are eigenvectors.

$$\left. \begin{array}{l} \underline{y} \in N(S) \\ \underline{x} \in \mathcal{C}(S) = \mathcal{C}(S^\top) \end{array} \right\} \implies \underline{x} \perp \underline{y}$$

This is confirmed also by the scheme done during lecture with the 4 blocks. Notice that we have not specified or made any assumption on the value of  $\lambda$ .

2. Similar to point 1, we consider the two vectors  $\underline{x}, \underline{y}$  such as:  $S\underline{x} = \lambda\underline{x}$  and  $S\underline{y} = \alpha\underline{y}$ . Now, consider the matrix  $(S - \alpha I)$ , we can write:

$$\begin{aligned} (S - \alpha I)\underline{y} = 0\underline{y} &\implies \underline{y} \in N(S - \alpha I) \\ (S - \alpha I)\underline{x} = (\lambda - \alpha)\underline{x} &\implies \underline{x} \in \mathcal{C}(S - \alpha I) = \mathcal{C}((S - \alpha I)^\top) \end{aligned}$$

So, again we obtain:  $\underline{x} \perp \underline{y}$ .

There is another property:  $\lambda_i \in \mathbb{R}$ , so the eigenvalues on the diagonal of  $\Lambda$  are real. Proof:

$$S\underline{x} = \lambda\underline{x} \implies \overline{\underline{x}}^\top S\underline{x} = \lambda \overline{\underline{x}}^\top \underline{x}$$

The  $\overline{\underline{x}}$  represent the conjugate of the vector  $\underline{x}$ . If that vector has complex components, those elements are conjugated, otherwise, i.e. they are all real, they remain unaltered. In particular, once a complex number is conjugated, the result is a real number, as shown here:

$$(a + ib)(a - ib) = (a^2 + b^2) \in \mathbb{R}$$

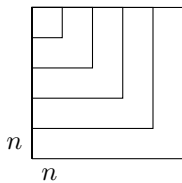
From the previous equation, we obtain:

$$\lambda = \frac{\overline{\underline{x}}^\top S\underline{x}}{\overline{\underline{x}}^\top \underline{x}} \in \mathbb{R}$$

## 7.5 Positive-definite symmetric matrices (SPD)

Characterizations:

- i  $\lambda_i > 0 \quad \forall i = 1, \dots, n$
- ii  $\underline{v}^\top S \underline{v} \geq 0 \quad \forall \underline{v} \in \mathbb{R}^n$ , with equality if and only if  $\underline{v} = 0$
- iii Leading determinants are positive.



This means that the determinant of the matrix obtained by taking the first  $k$  rows and columns of  $S$  is positive,  $\forall k = 1, \dots, n$ .

- iv Cholesky decomposition:  $S = B^\top B$ , with  $B$  upper triangular
- v All pivot elements are positive in the Gaussian elimination process

Let's consider  $\lambda > 0$  being a certain eigenvalue.

$$S\underline{x} = \lambda\underline{x}$$

We multiply both sides by  $\underline{x}^\top$ :

$$\underline{x}^\top S \underline{x} = \lambda \underline{x}^\top \underline{x} = \lambda \|\underline{x}\|^2 \geq 0$$

Recall that  $\underline{x}$  is an eigenvector while the before considered vector  $\underline{v}$  is a generic vector. With  $\underline{v}$ , instead, we have:

$$\underline{v} = (c_1 \underline{x}_1 + c_2 \underline{x}_2 + \dots + c_n \underline{x}_n)$$

So we are expressing  $\underline{v}$  as a linear combination of the eigenvectors of  $S$ .

$$\begin{aligned} & (c_1 \underline{x}_1 + c_2 \underline{x}_2 + \dots + c_n \underline{x}_n)^\top S (c_1 \underline{x}_1 + c_2 \underline{x}_2 + \dots + c_n \underline{x}_n) \\ & \left. \begin{aligned} c_1^2 \underline{x}_1^\top S \underline{x}_1 &= c_1^2 \lambda_1 \underline{x}_1^\top \underline{x}_1 = c_1^2 \lambda_1 \|\underline{x}_1\|^2 \\ c_1 c_2 \underline{x}_1^\top S \underline{x}_2 &= c_1 c_2 \lambda_2 \underline{x}_1^\top \underline{x}_2 = 0 \end{aligned} \right\} \text{there are two types of components} \end{aligned}$$

The first components is given by the eigenvectors with the same direction, while the second no so their scalar product is null (they are orthogonal).

From  $iv$ ):

$$S = B^\top B \implies \underline{v}^\top (B^\top B) \underline{v} = (\underline{v}^\top B^\top)(B \underline{v}) = (B \underline{v})^\top (B \underline{v}) = \|B \underline{v}\|^2 \geq 0$$

## 8 Singular Value Decomposition (SVD)

We are going to use it for:

- Least-squares approximation by introducing the pseudo-inverse of a matrix (Moore-Penrose inverse)
- Low-rank approximation with the Eckart-Young theorem

We start from:

$$A \in \mathbb{R}^{m \times n} \quad \begin{cases} m = \# \text{ of samples} \\ n = \# \text{ of features} \end{cases}$$

We can write:

$$A = U \Sigma V^\top$$

With:

- $U$  with dimensions  $m \times m$  and orthogonal
- $V^\top$  with dimensions  $n \times n$  and orthogonal
- $\Sigma$  with dimensions  $m \times n$  *almost* diagonal

If  $m > n$ , we can represent the matrices like this:

$$\underbrace{\left[ \begin{array}{c} \phantom{0} \\ \phantom{0} \\ \phantom{0} \end{array} \right]}_{m \times m} \underbrace{\left[ \begin{array}{c} \phantom{0} \\ \phantom{0} \\ \phantom{0} \end{array} \right]}_{m \times n} \underbrace{\left[ \begin{array}{c} \phantom{0} \\ \phantom{0} \end{array} \right]}_{n \times n}$$

What is the idea of SVD? Try to change features so variances are maximized and covariances are minimized. We don't want columns to be correlated.

In general:  $\text{rank}(A) = r < n$ .

$$AV = U\Sigma \iff V^\top V = I \iff V \text{ is orthogonal}$$

The component wise notation is:

$$A \underline{v}_i = \sigma_i \underline{u}_i$$

Given that the rank of  $A$  is  $r$ :

$$\begin{bmatrix} & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \end{bmatrix} \begin{bmatrix} \sigma_1 & & & & \\ & \sigma_2 & & & \\ & & \ddots & & \\ & & & \sigma_r & \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \end{bmatrix} \begin{cases} \sigma_1, \dots, \sigma_r > 0 \\ \sigma_{r+1}, \dots, \sigma_n = 0 \end{cases}$$

Typically:  $\sigma_1 > \sigma_2 > \dots > \sigma_r > \sigma_{r+1} = 0$ . We have:

$$\left\{ \begin{array}{l} \underline{Av}_1 = \sigma_1 \underline{u}_1 \\ \vdots \\ \underline{Av}_r = \sigma_r \underline{u}_r \end{array} \right\}^r \quad \left\{ \begin{array}{l} \underline{Av}_{r+1} = \sigma_{r+1} \underline{u}_{r+1} \\ \vdots \\ \underline{Av}_n = \sigma_n \underline{u}_n \end{array} \right\}^{n-r}$$

So the first  $r$  vectors span the column space of  $A$  while for the last  $n - r$  means that  $\underline{v}_i \in N(A)$  for  $i = r + 1, \dots, n$ . If we have  $A^\top$ , the decomposition is  $A^\top = (U\Sigma V^\top)^\top = V\Sigma^\top U^\top$ .

## 8.1 Economy SVD

What we've seen so far is the full SVD, but it can be optimized. Here is following the compact (reduced) representation, where once again we consider  $m > n$ :

$$\underbrace{\begin{bmatrix} & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \end{bmatrix}}_{m \times n} \underbrace{\begin{bmatrix} & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \end{bmatrix}}_{n \times n} \underbrace{\begin{bmatrix} & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \end{bmatrix}}_{n \times n}$$

This is caused by the fact that the last  $m - n$  rows in the central matrix are all 0 so multiply them for the last  $m - n$  columns of the left matrix is useless. This can be furthermore optimized by having matrix dimensions:  $(m \times r)(r \times r)(r \times n)$  because not all  $\sigma$  might be different than 0 (i.e. the rank of  $A$  is  $r$ ), so, in that case is useless even to multiply the last  $m - r$  rows of the central matrix.

**The SVD works for any matrix  $A$ .**

Let's suppose  $A$  is full rank  $n \times n$ :

$$A = U\Sigma V^\top = \sum_{i=1}^n \sigma_i \underbrace{\underline{u}_i \underline{v}_i^\top}_{\text{rank}=1}$$

View matrix-matrix multiplication for the rank 1 concept. If  $A$  is not full rank but instead has  $\text{rank}(A)=r$ , the same sum is no more computed until  $n$ , but instead  $r$ .

$$A = \sum_{i=1}^r \sigma_i \underline{u}_i \underline{v}_i^\top$$

What happens now if we pick a certain value  $\tilde{r} < r$ ?

$$A = U\Sigma V^\top \cong \sum_{i=1}^{\tilde{r}} \sigma_i \underline{u}_i \underline{v}_i^\top$$

We obtain a **rank  $\tilde{r}$  approximation of the matrix  $A$** . The rank of the matrix is known because is the sum of  $\tilde{r}$  matrices of rank 1. Moreover, that one, is the best approximation of rank  $\tilde{r}$  possible, i.e.:

$$\|A - \tilde{A}\| \leq \|A - B\| \quad \forall B \text{ of rank } = \tilde{r}$$

## 8.2 Proof of the existence of SVD

Once again, we start from matrix  $A \in \mathbb{R}^{m \times n}$  with rank  $= r$ . We consider the new matrix  $A^\top A$  which is:

- symmetric:  $(A^\top A)^\top = A^\top A$
- positive definite:  $\underline{x}^\top (A^\top A) \underline{x} = (\underline{x}^\top A^\top)(A \underline{x}) = (A \underline{x})^\top (A \underline{x}) = \|A \underline{x}\|^2 \geq 0$

We can use the following decomposition:

$$A^\top A = V \Lambda V^\top = \sum_{i=1}^n \lambda_i \underline{v}_i \underline{v}_i^\top$$

Recall that  $V$  contains the eigenvectors while  $\Lambda$  contains the eigenvalues. We rename  $\lambda_i = \sigma_i^2$ . The rank of  $A^\top A$  is  $r$ .

We want to prove that if  $\underline{x} \in N(A)$  then  $\underline{x} \in N(A^\top A)$ , to do so we proceed in both directions:

1. If we have  $A \underline{x} = 0 \implies \underline{x} \in N(A)$ . Is it possible to multiply both terms:

$$A^\top (A \underline{x}) = A^\top \underline{0} = \underline{0} \quad \text{so} \quad \underline{x} \in N(A) \implies \underline{x} \in N(A^\top A)$$

2. We start from  $(A^\top A) \underline{x} = 0 \implies \underline{x} \in N(A^\top A)$ . Again, we multiply:

$$\underline{x}^\top A^\top A \underline{x} = \|A \underline{x}\|^2 = 0 \quad \text{so} \quad \underline{x} \in N(A^\top A) \implies \underline{x} \in N(A)$$

Let's consider the couple of (eigenvalues, eigenvectors)  $= (\sigma_i^2, \underline{v}_i)$ :

$$A^\top A \underline{v}_i = \sigma_i^2 \underline{v}_i \xrightarrow{\text{component-wise}} A^\top A \underline{v}_i = \sigma_i^2 \underline{v}_i \quad (\dagger)$$

We introduce the quantity  $\underline{u}_i = \frac{A \underline{v}_i}{\sigma_i}$  which has some characteristics:

- i  $\underline{u}_i$  are unitary vectors:

$$\underline{u}_i^\top \underline{u}_i = \left( \frac{A \underline{v}_i}{\sigma_i} \right)^\top \left( \frac{A \underline{v}_i}{\sigma_i} \right) = \frac{\underline{v}_i^\top A^\top A \underline{v}_i}{\sigma_i^2} \stackrel{\dagger}{=} \frac{\sigma_i^2 \underline{v}_i^\top \underline{v}_i}{\sigma_i^2} = 1$$

The last passage of the equation is true because  $\underline{v}_i$  vectors are orthonormal.

- ii  $\underline{u}_i \perp \underline{u}_j$ :

$$\underline{u}_i^\top \underline{u}_j = \left( \frac{A \underline{v}_i}{\sigma_i} \right)^\top \left( \frac{A \underline{v}_j}{\sigma_j} \right) = \frac{\underline{v}_i^\top A^\top A \underline{v}_j}{\sigma_i \sigma_j} \stackrel{\dagger}{=} \frac{\sigma_j^2 \underline{v}_i^\top \underline{v}_j}{\sigma_i \sigma_j} = 0$$

- iii  $\underline{u}_i$  are eigenvectors of  $AA^\top$  with eigenvalues  $\sigma_i^2$ :

$$(AA^\top \underline{u}_i) = AA^\top \left( \frac{A \underline{v}_i}{\sigma_i} \right) = A \frac{A^\top A \underline{v}_i}{\sigma_i} \stackrel{\dagger}{=} A \frac{\sigma_i^2 \underline{v}_i}{\sigma_i} = \sigma_i^2 \left( \frac{A \underline{v}_i}{\sigma_i} \right) = \sigma_i^2 \underline{u}_i$$

We have demonstrated that  $A \underline{u}_i = \sigma_i \underline{u}_i$  and  $\underline{u}_i$  are orthonormal as well.

We have seen that  $\underline{u}_i = \frac{Av_i}{\sigma_i}$  but, what happen if  $\sigma_i = 0$ ?  
 Until now we have assumed that could not happen. Let's 2 examples, in the first we have the standard case while in the second is explained how to manage the case  $\sigma_i = 0$ .

**Example 1**

$$A = \begin{bmatrix} 4 & 4 \\ -3 & 3 \end{bmatrix} \quad \text{rank}(A) = 2$$

We build two new matrices  $X$  and  $Y$ :

$$X = A^T A = \begin{bmatrix} 25 & 7 \\ 7 & 25 \end{bmatrix} \quad \text{eigs}(X) = \begin{cases} \lambda_1 = 18 & \underline{v}_1 = \begin{bmatrix} -\frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} \end{bmatrix} \\ \lambda_2 = 32 & \underline{v}_2 = \begin{bmatrix} \frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} \end{bmatrix} \end{cases}$$

$$Y = AA^T = \begin{bmatrix} 32 & 0 \\ 0 & 18 \end{bmatrix} \quad \text{eigs}(Y) = \begin{cases} \lambda_1 = 18 & \underline{u}_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \lambda_2 = 32 & \underline{u}_2 = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \end{cases}$$

So we have all elements for constructing the SVD:

$$U = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} \sqrt{18} & 0 \\ 0 & \sqrt{32} \end{bmatrix} \quad V = \begin{bmatrix} -\frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \end{bmatrix}$$

It is easy to verify that  $A = U\Sigma V^T$ .

**Example 2**

$$A = \begin{bmatrix} 4 & 3 \\ 8 & 6 \end{bmatrix} \quad \text{rank}(A) = 1$$

As in the example before, we start building our matrices  $X$  and  $Y$ :

$$X = A^T A = \begin{bmatrix} 80 & 60 \\ 60 & 45 \end{bmatrix} \quad \text{eigs}(X) = \begin{cases} \lambda_1 = 125 & \underline{v}_1 = \begin{bmatrix} -4 \\ 5 \\ -3 \\ 5 \end{bmatrix} \\ \lambda_2 = 0 & \underline{v}_2 = ? \end{cases}$$

How can we compute the value of  $\underline{v}_2$ ? The idea is to build a vector  $\underline{v}_2$  such that it is orthogonal to  $\underline{v}_1$ . A possible solution is:  $\underline{v}_2 = \begin{bmatrix} 3 \\ 5 \\ -4 \\ 5 \end{bmatrix}$ .

The equation at the beginning of this page (the one regarding  $\underline{u}_i$ ) can still be computed for the vector  $\underline{v}_1$  in this case:

$$\underline{u}_1 = \frac{Av_1}{\sigma_1} = \frac{1}{\sqrt{125}} \begin{bmatrix} 4 & 3 \\ 8 & 6 \end{bmatrix} \begin{bmatrix} -4 \\ 5 \\ -3 \\ 5 \end{bmatrix} = \frac{1}{\sqrt{125}} \begin{bmatrix} -5 \\ -10 \end{bmatrix}$$

Now, the problem rises again because  $\sigma_2 = 0$  so we cannot compute  $\underline{u}_2$ . In reality, we do not need to compute  $\underline{u}_2$  because  $\underline{u}_1$  and  $\underline{v}_1$  are sufficient to recover the original matrix through (reduced) SVD, indeed:

$$\underbrace{\frac{1}{\sqrt{125}} \begin{bmatrix} -5 \\ -10 \end{bmatrix}}_U \underbrace{[\sqrt{125}]}_\Sigma \underbrace{\begin{bmatrix} -4 \\ 5 \\ -3 \\ 5 \end{bmatrix}}_V = \begin{bmatrix} 4 & 3 \\ 8 & 6 \end{bmatrix}$$

We could have used also the full SVD with the following matrices:

$$U = \frac{1}{\sqrt{125}} \begin{bmatrix} -5 & -10 \\ -10 & 5 \end{bmatrix} \quad \Sigma = \begin{bmatrix} \sqrt{125} & 0 \\ 0 & 0 \end{bmatrix} \quad V = \begin{bmatrix} -4 \\ 5 \\ -3 \\ 5 \end{bmatrix}$$

Where  $\underline{u}_2$  is the second column of the matrix  $U$  and is a vector orthogonal to  $\underline{u}_1$ . All elements written here and not in the previous formulation are a waste of memory. A concise recap of all versions of SVD:

1. full SVD:  $A = \underset{(m \times n)}{U} \underset{(m \times m)(m \times n)(n \times n)}{\Sigma} V^T$

2. economy SVD:  $A \underset{(m \times n)}{=} \underset{(m \times n)}{U} \underset{(n \times n)}{\Sigma} \underset{(n \times n)}{V^T}$

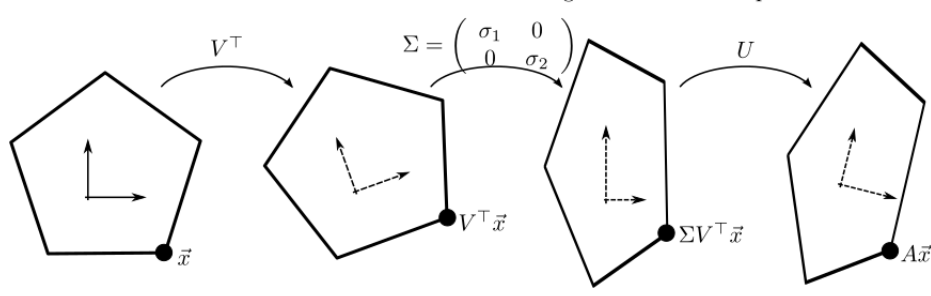
3. reduced SVD:  $A \underset{(m \times n)}{=} \underset{(m \times r)}{U} \underset{(r \times r)}{\Sigma} \underset{(r \times n)}{V^T}$

4. truncated SVD approximation:  $\sum_{i=1}^{\tilde{r}} \sigma_i \underline{u}_i \underline{v}_i^T$

For the first two cases the rank of  $A$  is  $n$ , for the third is  $r$  while in the first the approximation rank is decided.

### 8.3 Geometrical interpretation of SVD

Matrices  $U, V^T$  are orthonormal so, as mentioned in the section above, they represent a rotation (or reflection) while  $\Sigma$ , being a diagonal matrix, correspond to a scaling transformation.



Since for SVD no assumptions are made on the starting matrix, this means that any matrix can be obtained from 2 rotations and 1 scaling.

When  $A$  is squared and symmetric ( $A \in \mathbb{R}^{n \times n}, A = A^T$ )?

We can apply **polar decomposition**:

$$A = QR$$

Where  $Q$  is orthogonal and  $S$  is symmetric positive semi-definite. Why? From SVD.

$$A = U\Sigma V^T = \underbrace{(UV^T)}_Q \underbrace{(V\Sigma V^T)}_S$$

In particular, the product of the two orthogonal matrices in the first parenthesis is always another orthogonal matrix. In this case of decomposition, matrix  $A$  is obtained just by one rotation and one scaling (missing one more rotation with respect to classic SVD). Indeed, the second parenthesis results just in a stretch because  $V$  and  $V^{intercal}$  are two rotations identical and opposite, i.e. they cancel out.

#### 8.3.1 Properties of SVD

i If  $A$  is orthogonal, then  $\sigma_i = 1$  because if  $A$  orthogonal then  $A^T A = I$

ii All eigenvalues of a square matrix are  $\leq \sigma_1$ . Proof:

$$\|A\underline{x}\| = \|U\Sigma V^T \underline{x}\| = \|\Sigma V^T \underline{x}\| \leq \sigma_1 \|V^T \underline{x}\| = \sigma_1 \|\underline{x}\| \implies \|A\underline{x}_i\| \leq \sigma_1 \|\underline{x}_i\|$$

The matrix  $U$  disappear because an orthogonal matrix, when multiplied, does not change the magnitude of a vector. If we consider

$$\begin{aligned} \|A\underline{x}_i\| &= \|\lambda_i \underline{x}_i\| = |\lambda_i| \cdot \|\underline{x}_i\| \\ |\lambda_i| \cdot \|\underline{x}_i\| &\leq \sigma_1 \|\underline{x}_i\| \implies |\lambda_i| \leq \sigma_1 \quad \forall i \end{aligned}$$

### 8.4 Snapshots method

During real case scenarios, we will have a certain matrix  $A \in \mathbb{R}^{m \times n}$  and it might happen that  $m \gg n$  i.e. the number of samples is much greater than the number of features. In these cases, we can use the following trick to be more efficient.

For the SVD we need to compute both  $AA^T$  and  $A^T A$ . Which one is better to start with? And why? We would have  $AA^T : (m \times m)$  and  $A^T A : (n \times n)$ . Given  $m \gg n$  it is clear that the second one is better to start with because, being much smaller, it will be easier to compute its eigenvalues and eigenvectors.

## 8.5 Matrix norms

This concept is an extension of the vector norm. Given a matrix  $A \in \mathbb{R}^{m \times n}$ , we define the matrix norm as:

- **Frobenius norm:**  $\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2} = \sqrt{\text{Tr}(A^\top A)} = \sqrt{\text{Tr}(AA^\top)}$

Recall that  $\text{Tr}(AB) = \text{Tr}(BA)$ .

Example:

$$A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix} \quad \|A\|_F = \sqrt{1^2 + 2^2 + 3^2 + 4^2 + 5^2 + 6^2} = \sqrt{91}$$

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix} \begin{bmatrix} 1 & 3 & 5 \\ 2 & 4 & 6 \end{bmatrix} = \begin{bmatrix} 1^2 + 2^2 & \dots & \dots \\ \dots & 3^2 + 4^2 & \dots \\ \dots & \dots & 5^2 + 6^2 \end{bmatrix} \Rightarrow \sqrt{\text{Tr}(AA^\top)} = \sqrt{91}$$

If  $U$  is orthogonal, what happen to  $\|AU\|_F^2$ ?

$$\|AU\|_F^2 = \text{Tr}((AU)^\top AU) = \text{Tr}(U^\top A^\top AU) = \text{Tr}(A^\top A \underbrace{UU^\top}_I) = \text{Tr}(A^\top A) = \|A\|_F^2$$

This means that the Frobenius norm is invariant with respect to orthogonal transformations (†).

The Frobenius norm is also equal to  $\left(\sqrt{\sum_{i=1}^r \sigma_i^2}\right)$  where  $r$  is the rank of  $A$ . Proof:

$$\|A\|_F = \|U\Sigma V^\top\|_F \stackrel{\dagger}{=} \|\Sigma\|_F \triangleq \text{Tr} \sqrt{(\Sigma\Sigma^\top)} = \sqrt{\sum_{i=1}^r \sigma_i^2}$$

- **P-norms:** Recall the p-norm for vectors:

$$\underline{p} \in \mathbb{R}^n \Rightarrow \|\underline{p}\|_p = \left(\sum_{i=1}^n |p_i|^p\right)^{\frac{1}{p}}$$

Given  $A \in \mathbb{R}^{m \times n}$  we can define the p-norm for that matrix as:

$$\|A\|_p = \sup_{\underline{x} \in \mathbb{R}^n} \frac{\|A\underline{x}\|_p}{\|\underline{x}\|_p} = \sup_{\underline{x} \in \mathbb{R}^n, \|\underline{x}\|_p=1} \|A\underline{x}\|_p$$

If you choose  $p = 2$  you get the operator norm and  $\|A\|_2 = \sigma_1$ .

## 8.6 Eckart-Young theorem

Having defined the norms, we can now proceed with the proof of the **Eckart-Young theorem** which states:

For either  $\|\cdot\|_F$ ,  $\|\cdot\|_2$  we have:

$$\|A - A_k\| \leq \|A - B\| \quad \forall B \text{ of rank } k$$

where  $A_k = \sum_{i=1}^k \sigma_i \underline{u}_i \underline{v}_i^\top$  i.e. is the SVD approximation of rank  $k$  of the matrix  $A$ . Depending on the chosen norm you get:

$$\|A - A_k\| = \begin{cases} \sigma_{k+1} & \text{if } \|\cdot\|_2 \text{ considered} \\ \left(\sum_{i=1}^r \sigma_i^2\right)^{\frac{1}{2}} & \text{if } \|\cdot\|_F \text{ considered} \end{cases}$$

There will be 2 proofs, one for each type of norm.



### 8.6.1 Proof considering $\|\cdot\|_F$

We start from the **Weyl inequality**:

$$\sigma_{i+j-1}(X + Y) = \sigma_i(X) + \sigma_j(Y)$$

Where a generic  $\sigma_k(E)$  is the  $k$ -th singular value of the matrix  $E$ . We define

$$X = A_k - B \quad Y = B$$

So we have

$$\sigma_{i+k}(A) \leq \sigma_i(A - B) + \underbrace{\sigma_{k+1}(B)}_0$$

The last component has value of 0 because the matrix  $B$  has rank equal to  $k$ .

$$\|A - A_k\|_F^2 = \left( \sum_{i=k+1}^r \sigma_i^2(A) \right) \stackrel{\text{shift}}{=} \left( \sum_{i=1}^{r-k} \sigma_{i+k}^2(A) \right) \leq \left( \sum_{i=1}^{r-k} \sigma_{i+k}^2(A - B) \right) \leq \left( \sum_{i=1}^{\min(m,n)} \sigma_{i+k}^2(A - B) \right) = \|A - B\|_F^2$$

In particular, recall that  $A = \sum_{i=1}^r \sigma_i \underline{u}_i \underline{v}_i^\top$  and  $A_k = \sum_{i=1}^k \sigma_i \underline{u}_i \underline{v}_i^\top$ . The theorem is proved just by picking the first and last components of the inequality written above.

### 8.6.2 Proof considering $\|\cdot\|_2$

We are now going to prove the theorem with the last norm we did not use yet. Let's consider the matrix  $B$  with dimensions  $n \times d$  and  $\text{rank}(B)=k$ . This means that:

$$N(B) \subset \mathbb{R}^d \quad \dim(N(B)) = d - k$$

If we consider matrix  $V$  of the SVD decomposition of  $A$ :

$$V_{k+1} = \underbrace{\begin{bmatrix} | & | & | \\ v_1 & \dots & v_{k+1} \\ | & | & | \end{bmatrix}}_{k+1 \text{ cols}}$$

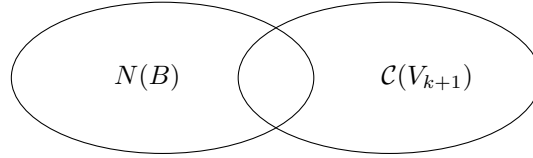
Those are the  $k + 1$  columns of  $V$ .

$$\mathcal{C}(V_{k+1}) \subset \mathbb{R}^d \quad \dim(\mathcal{C}(V_{k+1})) = k + 1$$

By adding the previous information, we have:

$$\dim(N(B)) + \dim(\mathcal{C}(V_{k+1})) = d - k + k + 1 = d + 1$$

Since both spaces are subset of  $\mathbb{R}^d$  and their summed dimensions are  $d + 1$  this means that the intersection of those spaces is not empty.



Consider  $\underline{w} \in N(B) \cap \mathcal{C}(V_{k+1})$ , we suppose for easyness that  $\|\underline{w}\|_2 = 1$ .

$$\underline{w} = \sum_{i=1}^{k+1} c_i \underline{v}_i \stackrel{(*)}{=} V_{k+1} \underline{c} \quad \sum_{i=1}^{k+1} c_i^2 = 1$$

We want to measure the following quantity:

$$\|A - B\|_2^2 = \underbrace{\sup_{\|\underline{w}\|} \|(A - B)\underline{w}\|_2^2}_{\text{particular } \|\underline{w}\|} \geq \underbrace{\|(A - B)\underline{w}\|_2^2}_{\text{generic } \|\underline{w}\|} \implies$$

Recall that  $\underline{w} \in N(B) \implies B\underline{w} = 0$ .

$$\implies \|A - B\|_2^2 \geq \|A\underline{w}\|_2^2 = \underline{w}^\top A^\top A \underline{w} \stackrel{\text{SVD}}{=} \underline{w}^\top V \Sigma^\top \underbrace{U^\top U}_I \Sigma V^\top \underline{w} =$$

$$\underline{w}^\top V \Sigma^\top \Sigma V^\top \underline{w} \stackrel{(*)}{=} \underline{c}^\top \underbrace{V_{k+1}^\top V_{k+1}}_I \Sigma^\top \Sigma \underbrace{V_{k+1}^\top V_{k+1}}_I \underline{c} = \underline{c}^\top \Sigma^\top \Sigma \underline{c} = \sum_{i=1}^{k+1} c_i^2 \sigma_i^2 \geq$$

Since singular values are ordered

$$\geq \sigma_{k+1}^2 \underbrace{\sum_{i=1}^{k+1} c_i^2}_{=1} = \sigma_{k+1}^2$$

So

$$\|A - B\|_2^2 \geq \sigma_{k+1}^2 = \|A - A_k\|_2^2$$

And, erasing the squares:

$$\|A - B\|_2 \geq \sigma_{k+1} = \|A - A_k\|_2$$

Where  $A_k$  is the rank  $k$  truncated SVD approximation, therefore

$$\|A - A_k\|_2 \leq \|A - B\|_2 \quad \forall B \text{ of rank } k$$

## 9 PCA

Here is a first real-world application of the SVD. PCA has the same aim as SVD i.e. find a way of projecting the dataset in a new space where variances are maximized and covariances are minimized.

We start from  $A \in \mathbb{R}^{n \times d}$  and we follow these points:

### i Center the matrix $A$

$\bar{A}$  is the mean centered with respect to columns while  $H$  is called the centering matrix and is obtained as follows:

$$H = I_n - \frac{1}{n} \underline{1}_n \underline{1}_n^\top$$

Where  $\underline{1}_n$  is the vector of dimension  $n$  containing all ones. The centered matrix is obtained:

$$\bar{A} = HA$$

### ii Build the covariance matrix

$$S = \frac{\bar{A}^\top \bar{A}}{n-1}$$

Where the denominator is  $n-1$  is because we want an unbiased estimator and it's not  $n$  because we have already taken 1 degree of freedom by centering the matrix. The covariance matrix is semidefinite positive so we can use eigenvalues and eigenvectors decomposition.

$$SV = VD \implies VDV^\top, D = V^\top SV$$

If you order the eigenvalues in decreasing order the corresponding eigenvectors are called principal components. To notice the relationship between SVD and PCA we can write:

$$S = \frac{1}{n-1} \bar{A}^\top \bar{A} = \frac{1}{n-1} V \Sigma^\top U^\top U \Sigma V^\top = \frac{1}{n-1} V \Sigma^2 V^\top$$

$$D = \frac{1}{n-1} \Sigma^2 \implies \lambda_k \frac{\sigma_k^2}{n-1}$$

PCA is SVD applied to a particular matrix.

### 9.1 Choose rank of truncated SVD

When using truncated SVD, how to choose the rank  $k$ ? One possibility is to use  $k$  such that a predefined percentage of the variance is retained. Another idea starts from:

$$A = A_{true} + \gamma A_{noise}$$

Where:

- $A$ : is our dataset
- $A_{true}$ : is the underlying low-rank representation of our data
- $\gamma$ : magnitude of the noise
- $A_{noise}$ : is a gaussian noise with 0 mean and unitary variance

By defining  $\tau$  as threshold we have that if  $\sigma_i > \tau$  we are picking  $A_{true}$ . There are two cases:

- $\gamma$  is known, i.e. we know the magnitude of the noise:
  - if  $A \in \mathbb{R}^{n \times n}$  (square) then  $\tau = \frac{4}{\sqrt{3}} \gamma \sqrt{n}$
  - if  $A \in \mathbb{R}^{m \times n}$  we have two more cases:
    - \* if  $n \ll m \implies \tau = \lambda(\beta)$  where  $\beta = \frac{n}{m}$  and  $\lambda$  is the following function:

$$\lambda(\beta) = \sqrt{2(\beta+1) + \frac{8\beta}{(\beta+1) + \sqrt{(\beta^2 + 14\beta + 1)}}}$$

\* if  $m \ll n \implies \tau = \lambda(\beta)$  where  $\beta = \frac{m}{n}$  so it's equal as before but in this case the numerator and denominator of  $\beta$  are swapped.

- $\gamma$  is unknown. We define  $\tau$  as follows:

$$\tau = \omega(\beta)\sigma_{med} \quad \omega(\beta) = \frac{\lambda(\beta)}{\mu_\beta}$$

In which  $\sigma_{med}$  is the median of the singular values and  $\mu_\beta$  is the median of the Marcenko-Pastur distribution.  $\lambda(\beta)$  is the same function as before. In particular:

$$\mu_\beta = \int_{(1-\beta)^2}^{\mu_\beta} \frac{\sqrt{((1-\sqrt{\beta})^2 - t)(t - (1-\sqrt{\beta})^2)}}{2\pi t} dt = \frac{1}{2}$$

## 9.2 Randomize SVD

## 10 Least squares approximation

Consider, with  $n > p$ :

$$\begin{aligned} X \in \mathbb{R}^{n \times p} \quad & n = \# \text{ samples}, p = \# \text{ features}, \text{rank}(X) = p \\ \underline{y} \in \mathbb{R}^n \quad & \text{labels of each sample} \end{aligned}$$

We have:

$$\underbrace{\begin{bmatrix} - & x_1^\top & - \\ - & x_2^\top & - \\ & \vdots & \\ - & x_n^\top & - \end{bmatrix}}_p \quad \underbrace{x_j}_{\text{j-th column of } X}$$

If we provide a new sample  $\tilde{x} \rightarrow \tilde{y}$  we want to predict  $\tilde{y}$ , i.e. the label, using the information we have from the training set.

We can use a linear model:

$$\tilde{y} = \tilde{x}^\top \underline{w} \quad \underline{w} \in \mathbb{R}^p$$

Typically  $\underline{y} \neq X\underline{w}$  so  $\underline{y}$  won't be precisely obtained with that multiplication. We will have to the approximation:  $\hat{\underline{y}} = X\underline{w}$ . So we can say that  $\hat{\underline{y}} \in \mathcal{C}(X)$  while in general  $\underline{y} \notin \mathcal{C}(X)$ . The error of the prediction is given by:

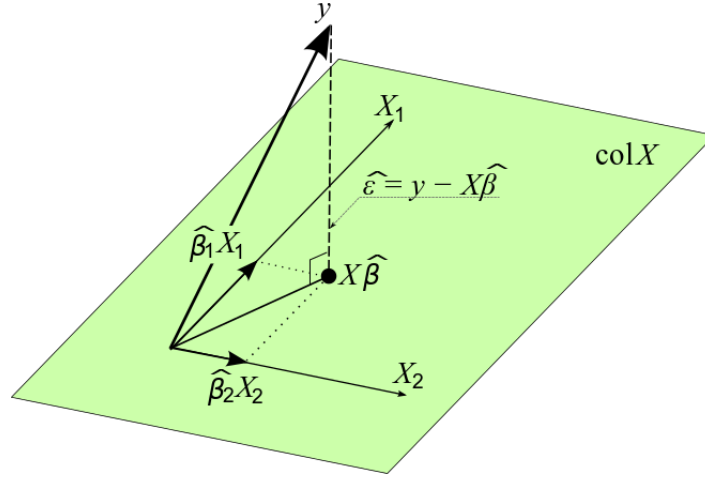
$$r_i(\underline{w}) = y_i - \hat{y}_i$$

We define the **residual vector** as  $\underline{r}(\underline{w})$  and our goal is to minimize its l2-norm squared:  $\|\underline{r}(\underline{w})\|_2^2$ .

Mathematically we can describe the problem as finding:  $\hat{\underline{w}} = \arg \min_{\underline{w}} \|\underline{r}(\underline{w})\|_2^2$ . There are two approaches possible:

1. Geometrical interpretation
2. By means of optimization procedures

### 10.1 Geometrical interpretation



So, a brief description of the figure above:

- the plane is the space of the predictions  $\mathcal{C}(X)$  indeed it is spanned by the column vectors of  $X$
- the vector  $\underline{y}$  is the vector of the label to be predicted
- the vector  $X\hat{\beta}$  in the figure is our prediction  $\hat{\underline{y}}$  and correspond to the projection of  $\underline{y}$  on the plane  $\mathcal{C}(X)$
- $\epsilon$  is the residual vector  $\underline{r}(\underline{w})$  and represent the error between the predicted value and the actual label

If we pick  $\underline{\bar{y}} \in \mathcal{C}(X)$  different from  $\hat{\underline{y}}$  and their respective residuals:  $\underline{\bar{r}} = \underline{y} - \underline{\bar{y}}$  and  $\underline{\hat{r}} = \underline{y} - \hat{\underline{y}}$  we have that:

$$\|\underline{\hat{r}}\|_2 \leq \|\underline{\bar{r}}\|_2$$

Indeed,  $\hat{r}$  is the orthogonal projection of  $\underline{y}$  on  $\mathcal{C}(X)$  and the orthogonal projection is the closest point to the vector  $\underline{y}$  in the subspace  $\mathcal{C}(X)$ .

Because of this, we can also say that:

$$\underline{x}_j^\top \hat{r} = 0 \quad j = 1, \dots, p$$

in matrix form:

$$X^\top \hat{r} = \underline{0} \implies X^\top (\underline{y} - \hat{\underline{y}}) = \underline{0} \implies X^\top (\underline{y} - X\hat{\underline{w}}) = \underline{0} \implies X^\top \underline{y} = X^\top X \hat{\underline{w}}$$

In general the rank of the matrix  $X^\top X$  is the same as the rank of  $X$  so if  $X$  is full rank then also  $X^\top X$  is full rank and this imply that is invertible and we can find  $\hat{\underline{w}}$  like this:

$$\hat{\underline{w}} = (X^\top X)^{-1} X^\top \underline{y}$$

If we are using this linear model for a binary classification we have to apply to the predicted value  $\hat{\underline{y}}$  a sign function in order to have the output restrained to  $\{-1, 1\}$ .

## 10.2 Optimization

Here we exploit mathematics for solving the initial problem:

$$\begin{aligned} \hat{\underline{w}} &= \arg \min_{\underline{w}} \|\underline{r}(\underline{w})\|_2^2 = \arg \min_{\underline{w}} \|\underline{y} - X\underline{w}\|_2^2 = \arg \min_{\underline{w}} (\underline{y} - X\underline{w})^\top (\underline{y} - X\underline{w}) = \\ &= \arg \min_{\underline{w}} \left[ \underline{y}^\top \underline{y} - (X\underline{w})^\top \underline{y} - \underline{y}^\top (X\underline{w}) + (X\underline{w})^\top X\underline{w} \right] = \arg \min_{\underline{w}} \underbrace{\left[ \underline{y}^\top \underline{y} - 2\underline{y}^\top X\underline{w} + \underline{w}^\top X^\top X \underline{w} \right]}_{F(\underline{w})} \end{aligned}$$

$F(\underline{w})$  is a **quadratic functional**.  $X$  is full rank,  $X^\top X$  is positive definite so this means that the functional is strictly convex and has a unique minimum. So we can compute the gradient and set it to zero:

$$\nabla_{\underline{w}} F(\underline{w}) = -2X^\top \underline{y} + 2X^\top X \underline{w} = \underline{0}$$

### Example 1 - Derivation

$$F(\underline{w}) = \underline{w}^\top \underline{c} = w_1 c_1 + w_2 c_2 + w_3 c_3 = \underline{c}^\top \underline{w} \quad \underline{w}, \underline{c} \in \mathbb{R}^3$$

The gradient is:

$$\nabla F = \begin{bmatrix} \frac{\partial F}{\partial w_1} \\ \frac{\partial F}{\partial w_2} \\ \frac{\partial F}{\partial w_3} \end{bmatrix} = \begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix} = \underline{c}$$

### Example 2 - Derivation

$$F(\underline{w}) = \underline{w}^\top A \underline{w} = \sum_{i=1}^n \sum_{j=1}^n w_j a_{ij} w_i \quad A \in \mathbb{R}^{n \times n}$$

$$\frac{\partial}{\partial w_k} (w_j a_{ij} w_i) = \begin{cases} a_{ij} w_i & k = j \neq i \\ w_j a_{ij} & k = i \neq j \\ 0 & k \neq i, k \neq j \\ 2w_k a_{ij} & k = i = j \end{cases}$$

In the last lecture we have introduced the least squares method. In particular we have mentioned the linear model for which:

$$\underline{\hat{y}} = X\underline{\hat{w}} \quad X \in \mathbb{R}^{n \times p} \quad \underline{\hat{y}} \in \mathbb{R}^n \quad n \geq p \quad X \text{ full rank}$$

We have obtained:

$$\underline{\hat{w}} = (X^\top X)^{-1} X^\top \underline{y} \implies \underline{\hat{y}} = \underbrace{X(X^\top X)^{-1} X^\top}_{P_x} \underline{y}$$

The matrix  $P_x$  dimension is given by the product of:  $(n \times p)(p \times p)(p \times n) = (n \times n)$  and has this properties:

- $P_x = P_x^2$
- $P_x$  is a projection matrix

Let's consider  $U$  an orthogonal ( $U^\top U = I$ ) matrix that contains the basis for  $\mathcal{C}(X)$  this means that  $\mathcal{C}(X) = \mathcal{C}(U)$ . We can write:

$$\underline{\hat{y}} = X\underline{\hat{w}} = U\underline{\tilde{w}}$$

So this basically means that the predicted value of  $y$  still a projection on a plane but this time the plane is spanned by the columns of  $U$  and not by the columns of  $X$ . By substituting last equation in the minimization method for least squares we have:

$$\underline{\tilde{w}} = \arg \min_{\underline{w}} \|\underline{y} - U\underline{w}\|_2^2 \implies \underline{\hat{y}} = U\underline{\tilde{w}} = U(U^\top U)^{-1} U^\top \underline{y} = U U^\top \underline{y}$$

This formulation is possible because this time in the parenthesis we have an orthogonal matrix and this means that  $(U^\top U)^{-1} = U^\top U = I$ . In general  $U U^\top \neq I$  because it might be rectangular (while  $U^\top U$  is always square).

**Example of usage of U** We start from  $X$ :

$$X = \begin{bmatrix} 1 & 1 \\ 1 & 3 \\ 0 & 0 \end{bmatrix} \quad \underline{x}_1 = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} \quad \underline{x}_2 = \begin{bmatrix} 1 \\ 3 \\ 0 \end{bmatrix}$$

How do we build the orthogonal matrix  $U$ ? We can use the Gram-Schmidt procedure:

$$\underline{u}_1 = \frac{\underline{x}_1}{\|\underline{x}_1\|} = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \\ 0 \end{bmatrix}$$

$$\underline{x}'_2 = \underline{x}_2 - (\underline{x}_2^\top \underline{u}_1) \underline{u}_1 = \underline{x}_2 - (\underline{u}_1^\top \underline{x}_2) \underline{u}_1 \implies \underline{u}_2 = \frac{\underline{x}'_2}{\|\underline{x}'_2\|} = \begin{bmatrix} -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \\ 0 \end{bmatrix}$$

So, the overall matrix  $U$  is:

$$U = \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ 0 & 0 \end{bmatrix} \quad U^\top U = I \text{ and } U U^\top \neq I$$

A drawback of Gram Schmidt is that, depending on the order chosen for the columns of  $X$ , the matrix  $U$  can be different. Moreover, the order of vector columns of  $U$  is meaningless.

Now, we want to exploit the SVD for computing the orthogonal matrix  $U$ . We start from the SVD of  $X$ :

$$X = U \Sigma V^\top$$

So

$$\begin{aligned}
\hat{\underline{w}} &= (X^\top X)^{-1} X^\top \underline{y} \\
&= (V \Sigma^\top U^\top U \Sigma V^\top)^{-1} V \Sigma^\top U^\top \underline{y} \\
&= (V \Sigma^\top \Sigma V^\top)^{-1} V \Sigma^\top U^\top \underline{y} \\
&= V (V \Sigma^\top \Sigma V^\top)^{-1} V \Sigma^\top U^\top \underline{y} \\
&= V (\Sigma^\top \Sigma)^{-1} \underbrace{V^\top V}_I \Sigma^\top U^\top \underline{y} \\
&= V \underbrace{(\Sigma^\top \Sigma)^{-1} \Sigma^\top}_{\Sigma^+} U^\top \underline{y} \\
&= V \Sigma^+ U^\top \underline{y}
\end{aligned}$$

Recall that:

$$\begin{aligned}
(AB)^{-1} &= B^{-1} A^{-1} \\
(V^\top)^{-1} &= V
\end{aligned}$$

Because  $V$  is orthogonal.

$\Sigma^+$  is called the pseudo-inverse of  $\Sigma$ .

Eventually, we have:

$$\begin{aligned}
\hat{\underline{y}} &= X (X^\top X)^{-1} X^\top \underline{y} = X X^+ \underline{y} \\
&= U (U^\top U)^{-1} U^\top \underline{y} = U U^+ \underline{y}
\end{aligned}$$

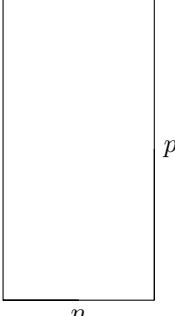
Recalling that we are considering the case in which  $n \geq p$ , the matrices have these dimensions:  $U = (n \times n)$ ,  $\Sigma = (n \times p)$ ,  $V^\top = (p \times p)$  and:

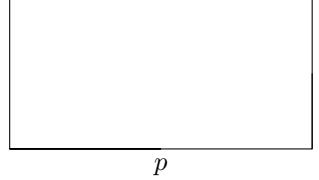
$$\begin{aligned}
\Sigma &= \begin{bmatrix} \sigma_1 & 0 & \dots & 0 \\ 0 & \sigma_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_p \\ 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix} \\
\Sigma^\top \Sigma &= \begin{bmatrix} \sigma_1 & 0 & \dots & 0 & 0 & \dots & 0 \\ 0 & \sigma_2 & \dots & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_p & 0 & \dots & 0 \end{bmatrix} \begin{bmatrix} \sigma_1 & 0 & \dots & 0 \\ 0 & \sigma_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_p \\ 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix} = \begin{bmatrix} \sigma_1^2 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_p^2 \end{bmatrix} \\
(\Sigma^\top \Sigma)^{-1} &= \begin{bmatrix} \frac{1}{\sigma_1^2} & 0 & \dots & 0 \\ 0 & \frac{1}{\sigma_2^2} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \frac{1}{\sigma_p^2} \end{bmatrix} \\
\Sigma^+ = (\Sigma^\top \Sigma)^{-1} \Sigma^\top &= \begin{bmatrix} \frac{1}{\sigma_1} & 0 & \dots & 0 & 0 & \dots & 0 \\ 0 & \frac{1}{\sigma_2} & \dots & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \frac{1}{\sigma_p} & 0 & \dots & 0 \end{bmatrix}
\end{aligned}$$

Let's consider now the case in which  $p \geq n$  and  $X$  has  $n$  linearly independent rows (before we had  $p$  linearly independent columns). This means that we have more unknowns than equations and we would find infinite solutions for  $\hat{\underline{w}}$  such that  $\hat{\underline{y}} = X \hat{\underline{w}}$ .

The solution found before  $\hat{\underline{w}} = V \Sigma^+ U^\top \underline{y}$  it's still valid but now  $\Sigma^+ = \Sigma^\top (\Sigma \Sigma^\top)^{-1}$  so it has the same shape as before but transposed. Summary:



$n > p$   
  
 $\hat{\underline{w}} = V\Sigma^+U^\top \underline{y}$   
 $\Sigma^+ = (\Sigma^\top \Sigma)^{-1} \Sigma^\top$

$p > n$   
  
 $\hat{\underline{w}} = V\Sigma^+U^\top \underline{y}$   
 $\Sigma^+ = \Sigma^\top (\Sigma \Sigma^\top)^{-1}$

Let's consider two vectors:  $\underline{w}$  and  $\hat{\underline{w}}$ :

$$\underline{y} = x\underline{w} = X\hat{\underline{w}}$$

So they both are solutions of the system of equations.

$$\begin{aligned}
\|\underline{w}\|_2^2 &= \|(\underline{w} - \hat{\underline{w}}) + \hat{\underline{w}}\|_2^2 \\
&= \|\underline{w} - \hat{\underline{w}}\|_2^2 - 2 \underbrace{(\underline{w} - \hat{\underline{w}})^\top \hat{\underline{w}}}_\star + \|\hat{\underline{w}}\|_2^2 \\
&= \underbrace{\|\underline{w} - \hat{\underline{w}}\|_2^2}_{\geq 0} + \|\hat{\underline{w}}\|_2^2
\end{aligned}$$

$$\begin{aligned}
\star &= (\underline{w} - \hat{\underline{w}})^\top X^\top (X X^\top)^{-1} \underline{y} \\
&= (X(\underline{w} - \hat{\underline{w}}))^\top (X X^\top)^{-1} \underline{y} \\
&= (X\underline{w} - X\hat{\underline{w}})^\top (X X^\top)^{-1} \underline{y} \\
&= (0)^\top (X X^\top)^{-1} \underline{y} \\
&= 0
\end{aligned}$$

So

$$\implies \|\underline{w}\|_2^2 \geq \|\hat{\underline{w}}\|_2^2$$

This means that  $\hat{\underline{w}}$  has minimum norm solution, i.e. among all possible vectors  $w$  that satisfy the initial equation  $\underline{y} = x\underline{w} = X\hat{\underline{w}}$ ,  $\hat{\underline{w}}$  is the one with the minimum norm. In order to understand why that property is important, consider the following example.

Suppose to have:

$$X = \begin{bmatrix} 1 & 0 & 0.01 \\ 0 & 1 & 0 \end{bmatrix} \quad \underline{y} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

We could choose:

$$\underline{w}_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad \underline{w}_2 = \begin{bmatrix} 0 \\ 0 \\ 100 \end{bmatrix}$$

It is easily understandable that, if would choose the second vector as a solution instead of the first one, we would have an important amplification of the error related to the measure of the value of the column feature. The same cannot be said for the first vector.

Notice that, it is not that improbable to have more features than samples, consider for example images, for which each pixel is a feature.

## 11 Matrix completion

We start from

$$X \in \mathbb{R}^{n \times p} \quad \text{rank}(X) = r \ll \min(n, p)$$

We know only partially  $X$ , we know  $X_{i,j}$  for  $(i, j) \in \Omega$ .

Matrix completion is a method for filling missing values. If we had full rank matrix we would have very independent columns so we would not be able to retrieve/obtain missing values. Being low rank, in this sense, helps us accomplishing this goal.

## 11.1 Ideal estimator

$$\begin{cases} \hat{X} = \arg \min_{z \in \mathbb{R}^{n \times p}} [\text{rank}(z)] \\ \text{subject to } \hat{X}_{ij} = X_{ij}(i, j) \in \Omega \end{cases}$$

This formulation is computationally unfeasible and the object function is non-convex. What is in practical terms saying is that, among all possible solutions, we want the one with minimal rank.

How can we solve it?

## 11.2 Practical estimator

We are going to use the **Nuclear norm**:  $\|\cdot\|_* = \sum_{i=1}^{\min(n,p)} \sigma_i$ .

The idea is that, instead of minimizing the rank, we minimize the norm since, when performing the SVD, the number of non-zero singular values is equal to the rank of the matrix.

$$\begin{cases} \hat{X} = \arg \min_{z \in \mathbb{R}^{n \times p}} \|z\|_* \\ \text{subject to } \hat{X}_{ij} = X_{ij}(i, j) \in \Omega \end{cases}$$

This new formulation is convex-optimal. How are we going to solve this problem? With the **SVT (Singular Value Threshold)**. Here is the algorithm:

- Initialize  $\hat{X} = \text{zeros}(n, p)$
- Set  $\hat{X}_{ij} = X_{ij}$  for  $(i, j) \in \Omega$
- For  $k = 1, 2, \dots, N$ 
  - $\hat{X}_{old} = \hat{X}$
  - $[U, \Sigma, V^\top] = \text{svd}(\hat{X})$
  - $\Sigma \rightarrow \hat{\Sigma} \begin{cases} \hat{\sigma}_i = \sigma_i & \text{if } \sigma_i > \tau \\ \hat{\sigma}_i = 0 & \text{otherwise} \end{cases}$
  - $\hat{X} = U \hat{\Sigma} V^\top$
  - $\hat{X}_{ij} = X_{ij}$  for  $(i, j) \in \Omega$
- $\|\hat{X} - X_{old}\| < \epsilon$

The constant  $\tau$  is imposing some sort of thresholding on the acceptance of singular values. This is an example of non-monothone algorithm because after the reduction of rank with the condition on sigmas, we override some value of the matrix to the one contained in  $\Omega$ . Is it proved that for a large enough index  $k$  the algorithm converge to the solution.

Let's recap what we have seen so far regarding the Least Squares method.

$$(x_i, y_i) \quad i = 1, \dots, n \quad \underline{x}_i \in \mathbb{R}^p \text{ and } y_i \in \mathbb{R}$$

In matrix form:

$$X \in \mathbb{R}^{n \times p} \text{ and } \underline{y} \in \mathbb{R}^n$$

We have made the hypothesis, until now, of having a full rank matrix  $X$ . The linear model is defined as:

$$\hat{\underline{y}} = X\hat{\underline{w}}_{LS} \quad \text{with} \quad \hat{\underline{w}}_{LS} = (X^\top X)^{-1} X^\top \underline{y}$$

Moreover, we said that the actual value of  $y$  is not exactly given by this linear model but it will be an approximation:

$$\underline{y} \approx X\underline{w}_{LS}$$

And this can be written by expliciting an error:

$$\underline{y} = X\underline{w}^* + \underline{\epsilon}$$

where  $\underline{\epsilon}$  is the error or noise vector. If we replace this last equation in the one with  $\hat{\underline{w}}_{LS}$  we get:

$$\hat{\underline{w}}_{LS} = (X^\top X)^{-1} X^\top (X\underline{w}^* + \underline{\epsilon}) = \underline{w}^* + (X^\top X)^{-1} X^\top \underline{\epsilon} =$$

Now we can apply the SVD to  $X$  so, if we call  $X = U\Sigma V^\top$  we get, from previous lectures, that:  $(X^\top X)^{-1} X^\top = V\Sigma^+ U^\top$ . We have said before that  $X$  is full rank so we will have all singular components different than 0 and  $\sigma_1 > \sigma_2 > \dots > \sigma_p > 0$ . The matrix  $\Sigma^+$  is a psuedo-diagonal matrix with the inverse of the singular values on the diagonal (check few pages above).

$$= \underline{w}^* + V\Sigma^+ U^\top \underline{\epsilon}$$

- $U^\top \underline{\epsilon}$ : we multiply a vector with an orthogonal matrix so its norm will not change.
- $V\Sigma^+ U^\top \underline{\epsilon}$ : we multiply a vector with a diagonal matrix so we will have a scaling of the vector. But **if a singular value is very very small, its inverse in the matrix SigmaPlus will be huge and will scale the error vector a lot! So the error will be amplified and the original model vector  $\underline{w}^*$  will be negligible.**

### 11.3 Ridge regression (regularization)

This method will help us in preventing the problem mentioned just before. We start from the definition of the weight vector for linear model explicited for the optimization method of the Least Squares:

$$\hat{\underline{w}}_{LS} = \arg \min_{\underline{w}} \underbrace{\|\underline{y} - X\underline{w}\|_2^2 + \lambda \|\underline{w}\|_2^2}_{f(\underline{w})}$$

In particular we have added a term.

$$f(\underline{w}) = \underline{y}^\top \underline{y} - 2\underline{w}^\top X\underline{y} + \underline{w}^\top X^\top X \underline{w} + \lambda \underline{w}^\top \underline{w}$$

We can now compute the gradient of this function:

$$\nabla_{\underline{w}}(f(\underline{w})) = -2X^\top \underline{y} + 2X^\top X \underline{w} + 2\lambda \underline{w} = 0$$

$$X^\top \underline{y} = (X^\top X + \lambda I) \underline{w}$$

$$\hat{\underline{w}}_R = (X^\top X + \lambda I)^{-1} X^\top \underline{y}$$

It's easy to notice that if  $\lambda = 0$  we get the Least Squares solution. If  $\lambda > 0$  we will have a different solution. We can now compute the SVD of  $X$ :

$$\begin{aligned} \hat{\underline{w}}_R &= (V\Sigma^\top \underbrace{U^\top U}_I \Sigma V^\top + \lambda \underbrace{V^\top V}_I)^{-1} V\Sigma^\top U^\top \underline{y} \\ &= [V(\Sigma^\top \Sigma + \lambda I)V^\top]^{-1} V\Sigma^\top U^\top \underline{y} \\ &= V \underbrace{(\Sigma^\top \Sigma + \lambda I)\Sigma^\top}_M U^\top \underline{y} \end{aligned}$$

Where

$$M = \begin{bmatrix} \frac{\sigma_1}{\sigma_1^2 + \lambda} & 0 & \dots & 0 & 0 & \dots & 0 \\ 0 & \frac{\sigma_2}{\sigma_2^2 + \lambda} & \dots & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \frac{\sigma_p}{\sigma_p^2 + \lambda} & 0 & \dots & 0 \end{bmatrix}$$

- if  $\sigma_i$  is big compared to  $\lambda$  then  $\frac{\sigma_i}{\sigma_i^2 + \lambda} \approx \frac{1}{\sigma_i}$
- if  $\sigma_i$  is close to 0 then  $\frac{\sigma_i}{\sigma_i^2 + \lambda} \approx 0$

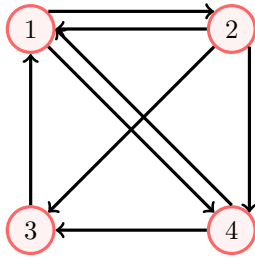
If we now consider again the problem of having a singular value close to 0 we can see that it is now solved because we would reduct to the second case just above and the pseudo inverse would be almost 0. If  $\lambda$  is very small the matrix of  $\Sigma$ 's is almost equal to  $\Sigma^+$ .

$$\begin{aligned}\hat{w}_R &= (X^\top X + \lambda I)^{-1} X^\top \underline{y} \\ &= (X^\top X + \lambda I)^{-1} X^\top (X \underline{w}^* + \underline{\epsilon}) \\ &= (X^\top X + \lambda I)^{-1} X^\top X \underline{w}^* + (X^\top X + \lambda I)^{-1} X^\top \underline{\epsilon}\end{aligned}$$

If  $\underline{\epsilon} = \underline{0}$  then there is only the first term and, since we are not finding the perfect projection on the plane, we make an higher error with respect to Least Squares.

## 12 Page Rank

Consider 4 websites, the arrows represents the link from one site to another.



The idea is to surf the web randomly (random walks on the graph) and if you do that long enough you will reach a **Steady state** where  $\pi_i$  will be the probability of being on the  $i$ -th website. In this case the vector  $\underline{\pi} \in \mathbb{R}^4$  because there are 4 websites.

We are assuming there are not separated sites. How can we represent the matrix? It an **Adjacency matrix**

$$\tilde{A} = \begin{bmatrix} 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 \end{bmatrix} \rightarrow \tilde{A}_{ij} = \begin{cases} 1 & \text{if there is a link from } i \text{ to } j \\ 0 & \text{otherwise} \end{cases} \quad \xrightarrow{\text{normalized version}} \quad A = \begin{bmatrix} 0 & 1 & 1/3 & 1/2 \\ 0 & 0 & 1/3 & 1/2 \\ 1/2 & 0 & 0 & 0 \\ 1/2 & 0 & 1/3 & 0 \end{bmatrix}$$

In the normalized version all columns sum up to 1. What happen if we multiply the matrix  $A$  times a canonical basis vector  $\underline{e}_3 = [0 \ 0 \ 1 \ 0]^\top$ ?

$$A \underline{e}_3 = \begin{bmatrix} 1/3 \\ 1/3 \\ 0 \\ 1/3 \end{bmatrix}$$

As it was trivial to figure, we obtain, in this case, the third column of the adjacency matrix. In a probability perspective, the vector  $\underline{e}_3$  represents the probability of starting from the third website. The vector we obtain is the probability of reaching the other websites starting from the third one. This means that in this case we are certain that we start from the third website and for sure we won't be in that site in the following iteration.

Considering

$$\underline{\pi} = \begin{bmatrix} \pi_1 \\ \pi_2 \\ \pi_3 \\ \pi_4 \end{bmatrix}$$

The steady state is defined as:  $A \underline{\pi} = \underline{\pi}$ . This means that, the steady state is the eigenvector of the matrix  $A$  with eigenvalue 1. The matrix  $A$  is positive (not positive-definite) i.e. it has all non-negative coefficients. A positive matrix is denoted with  $A > 0$ . From Perron-Frobenius theory we know that  $\lambda_1 = 1$  is the largest eigenvalue.

$$\lambda_1 = 1 > \lambda_2 > \lambda_3 > \dots > \lambda_n \quad \lambda_i \neq 0$$

As mentioned in a previous lecture, we can use the power method in order to retrieve the largest eigenvalue. In particular, we start from:

$$\underline{\pi}^{(0)} \quad \text{with} \quad \|\underline{\pi}^{(0)}\| = 1$$

Then, for  $k = 1, 2, \dots$

$$\underline{\pi}^{(k)} = \frac{A\underline{\pi}^{(k-1)}}{\|A\underline{\pi}^{(k-1)}\|}$$

if  $\|\underline{\pi}^{(k)} - \underline{\pi}^{(k-1)}\| < \epsilon$  then stop

Obviously,  $\epsilon$  represents a tolerance value. As we can see, there is an recursive definition in a sense that, at each iteration, the same operation is made on the same variable. For example:

$$\underline{\pi}^{(1)} = \frac{A\underline{\pi}^{(0)}}{\|A\underline{\pi}^{(0)}\|} \quad \underline{\pi}^{(2)} = \frac{A\underline{\pi}^{(1)}}{\|A\underline{\pi}^{(1)}\|} = \frac{A^2\underline{\pi}^{(0)}}{\|A^2\underline{\pi}^{(0)}\|}$$

So, iterating  $k$  times:

$$\underline{\pi}^{(k)} = \frac{A^k \underline{\pi}^{(0)}}{\|A^k \underline{\pi}^{(0)}\|}$$

Since in  $A$  there are no columns that are completely equal to 0, the matrix can be diagonalized. This implies that there are some  $\underline{v}_i, i = 1, \dots, n$  that can be used as a basis for the space. In particular we can write:

$$\underline{\pi}^{(0)} = \sum_{i=1}^n \alpha_i \underline{v}_i$$

We want to plug this expression in the previous equation (power method). Before, notice that,  $A\underline{v}_i = \lambda_i \underline{v}_i$  because  $\underline{v}_i$  are the vectors which diagonalize  $A$  so the ones such that  $A = V\Lambda V^T$ . The numerator of  $\underline{\pi}^{(k)}$  can be written as:

$$A^k \underline{\pi}^{(0)} = \alpha_1 \lambda_1^k \left( \underline{v}_1 + \sum_{i=2}^n \frac{\alpha_i}{\alpha_1} \left( \frac{\lambda_i}{\lambda_1} \right)^k \underline{v}_i \right)$$

Proof:

$$\begin{aligned} A^k \underline{\pi}^{(0)} &= V \Lambda^k V^{-1} (\alpha_1 \underline{v}_1 + \dots + \alpha_n \underline{v}_n) \\ &\stackrel{\circledast}{=} V \Lambda^k (\alpha_1 \underline{e}_1 + \dots + \alpha_n \underline{e}_n) \\ &= V (\alpha_1 \lambda_1^k \underline{e}_1 + \dots + \alpha_n \lambda_n^k \underline{e}_n) \\ &\stackrel{\diamond}{=} \alpha_1 \lambda_1^k \underline{v}_1 + \dots + \alpha_n \lambda_n^k \underline{v}_n \\ &= \alpha_1 \lambda_1^k \left( \underline{v}_1 + \frac{\alpha_2}{\alpha_1} \left( \frac{\lambda_2}{\lambda_1} \right)^k \underline{v}_2 + \dots + \frac{\alpha_n}{\alpha_1} \left( \frac{\lambda_n}{\lambda_1} \right)^k \underline{v}_n \right) \end{aligned}$$

$\circledast \rightarrow V^{-1} \underline{v}_i = V^{-1} \underline{e}_i = \begin{bmatrix} - & \underline{v}_1^T & - \\ - & \underline{v}_2^T & - \\ - & \vdots & - \\ - & \underline{v}_n^T & - \end{bmatrix} \underline{v}_i = \underline{e}_i$

$\diamond \rightarrow A\underline{v} = \lambda \underline{v} \implies AA\underline{v} = \lambda A\underline{v} \implies A^2 \underline{v} = \lambda^2 \underline{v}$

If  $k \rightarrow \infty$  then  $\left( \frac{\lambda_i}{\lambda_1} \right)^k \rightarrow 0$  because  $\lambda_1 > \lambda_i$  for  $i = 2, \dots, n$ . So, the only remaining term is  $\underline{v}_1$  which is the eigenvector with the largest eigenvalue. So what we have just done is the proof of the convergence of the power method. The closer are the eigenvalues to  $\lambda_1$  the slower is the convergence.

So far we have considered:

$$\underline{y} = X\underline{w} \quad X \in \mathbb{R}^{n \times p} \text{ and } \left\{ \begin{array}{l} \text{Least Squares : } \hat{\underline{w}}_L^S \\ \text{Ridge Regression : } \hat{\underline{w}}_{RR} \end{array} \right\} \underline{w} \in \mathbb{R}^p \text{ dense vector, i.e. not many zeros}$$

We are going to see a method which obtain a vector  $\underline{w}$  with many zeros, as sparse as possible. We want to consider this model:

$$\underline{y} = X\underline{w} \quad X \in \mathbb{R}^{n \times p} \quad p > n$$

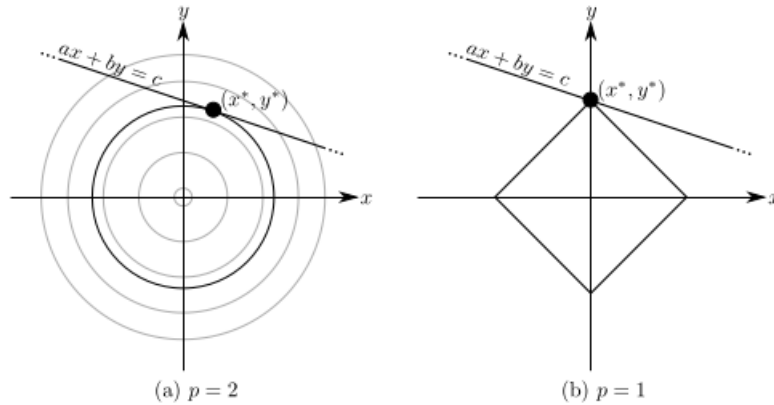
We have said that the system is undetermined since it has infinite solutions. Suppose to have 2 features and 1 sample.

$$X = \begin{bmatrix} 2 & 3 \end{bmatrix} \quad y = [1]$$

And so:

$$1 = \begin{bmatrix} 2 & 3 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \implies 1 = 2w_1 + 3w_2 \leftarrow \text{line}$$

As mentioned in a previous lecture, we want to find the minimum length solution so we can plot the line found before and the circles that represents the l2-norm distance.



On the right-hand side it is represented the equivalent plot but with the l1-norm. We can see that the solution on the right has one coordinate (are features!) that is equal to zero. This is true in general, i.e. we will obtain sparser solution using the l1-norm rather than the l2-norm.

With L1-norm it's still a convex optimization problem and

$$F(\underline{w}) = \|X\underline{w} - \underline{y}\|_2^2 + \lambda \|\underline{w}\|_1$$

This model is implemented by **Lasso** (Least Absolute Shrinkage and Selection Operator) and achieve both the shortest distance solution and the selection of some features.

This is important because by reducing the number of features, we increase the interpretability of the model.

There is also the **Elastic Net** which combines both Lasso and Ridge:

$$F(\underline{w}) = \|X\underline{w} - \underline{y}\|_2^2 + \lambda_1 \|\underline{w}\|_1 + \lambda_2 \|\underline{w}\|_2^2$$

Once again, we start considering:

$$\underline{y} \in \mathbb{R}^n \quad X \in \mathbb{R}^{n \times p} \text{ with } p \text{ lin. ind. cols } (\implies \sigma_i > 0, i = 1, \dots, p)$$

Now we write the formulation of the weights vector  $w$  we found for Ridge Regression.

$$\begin{aligned} \hat{\underline{w}}_R &= V \underbrace{(\Sigma^\top \Sigma + \lambda I)^{-1} \Sigma^\top U^\top}_{\Sigma^\top (\Sigma \Sigma^\top + \lambda I)^{-1}} \underline{y} \\ &= \underbrace{V \Sigma^\top U^\top}_{X^\top} \underbrace{U (\Sigma \Sigma^\top + \lambda I)^{-1} U^\top}_{\alpha \in \mathbb{R}^n} \underline{y} \\ &\stackrel{\square}{=} X^\top \underline{\alpha} \\ &= \sum_{i=1}^n \alpha_i \underline{x}_i \end{aligned}$$

In the second passage we have added the matrices  $U^T U$  because its the identity matrix. We obtain a weighted sum of the column vectors of  $X^T$ , where  $X$  is:

$$X = \begin{bmatrix} - & \underline{x_1^T} & - \\ - & \underline{x_2^T} & - \\ & \vdots & \\ - & \underline{x_n^T} & - \end{bmatrix}$$

So, until now, we have discussed about linear models. A generic form would be:

$$\hat{y}_i = w_1 x_{i1} + w_2 x_{i2} \quad \text{if} \quad \underline{x}_i = \begin{bmatrix} x_{i1} \\ x_{i2} \end{bmatrix}$$

Now a new model:

$$\hat{y}_i = w_1 x_{i1} + w_2 x_{i2} + w_3 x_{i1}^2 + w_4 x_{i2}^2 + w_5 x_{i1} x_{i2}$$

So, the original feature vector  $\underline{x}$  is transformed into a new feature vector by means of function called feature map  $\phi(x)$ :

$$\phi(x) = \begin{bmatrix} x_1 \\ x_2 \\ x_1^2 \\ x_2^2 \\ x_1 x_2 \end{bmatrix} \in \mathbb{R}^d \quad d > p \text{ typically}$$

And

$$\hat{y}_i = \phi(x_i)^T \underline{w}$$

In general  $d$  can be huge.

## 13 Kernel Methods

The aim of this methods is to avoid the necessity of computing huge vectors.

$$\Phi = \begin{bmatrix} - & \phi(\underline{x}_1)^T & - \\ - & \phi(\underline{x}_2)^T & - \\ & \vdots & \\ - & \phi(\underline{x}_n)^T & - \end{bmatrix} \in \mathbb{R}^{n \times d} \quad \underline{\hat{y}} = \Phi \underline{w}$$

The objective is still the same: i.e. finding  $\underline{w}$ . We are now going to consider ridge regression in order to achieve that. Instead of  $\underline{\hat{w}}_R = X^T \underline{\alpha}$  we can now write:

$$\underline{\hat{w}}_R = \Phi^T \underline{\alpha}$$

Where

$$\underline{\alpha} = U(\Sigma \Sigma^T + \lambda I)^{-1} U^T \underline{y} = (X X^T + \lambda I)^{-1} \underline{y}$$

Here is the proof:

$$(U \Sigma V^T V \Sigma^T U^T + \lambda U U^T)^{-1} \underline{y}$$

So, by similarity with the expression obtained in the previous page ( $\underline{\hat{w}}_R \stackrel{\square}{=} X^T \underline{\alpha}$ ) we can write:

$$\underline{\alpha} = (\Phi \Phi^T + \lambda I)^{-1} \underline{y}$$

The computation of  $\underline{\alpha}$  is practically impossible since  $\Phi$  is huge. How can we reduce the cost of that computation? Introducing the **Kernel Trick of Kernel Function**:

$$K(\underline{x}_i, \underline{x}_j) = \underline{x}_i^T \underline{x}_j \quad \leftarrow \text{scalar product between 2 samples}$$

Consider:

$$\underline{x}_i = \begin{bmatrix} x_{i1} \\ x_{i2} \end{bmatrix} \rightarrow \phi(\underline{x}_i) = \begin{bmatrix} x_{i1}^2 \\ \sqrt{2} x_{i1} x_{i2} \\ x_{i2}^2 \end{bmatrix}$$

$$\phi(\underline{x}_i)^T \phi(\underline{x}_j) = x_{i1}^2 x_{j1}^2 + 2 x_{i1} x_{i2} x_{j1} x_{j2} + x_{i2}^2 x_{j2}^2 = (\star)$$

In order to compute  $(\star)$  we have to apply the feature map twice and the scalar product between the two vectors. We can avoid this by using the kernel function:

$$K(\underline{x}_i, \underline{x}_j) = (\underline{x}_i^\top \underline{x}_j)^2 = (x_{i1}x_{j1} + x_{i2}x_{j2})^2 = (\star)$$

But here i'm computing this product immediately, without applying the feature map. We have computed the scalar product in a cheaper way.

Typical kernel functions family:

i Polynomials of degree  $q$ :  $K(\underline{x}_i, \underline{x}_j) = (\underline{x}_i^\top \underline{x}_j)^q$

ii Polynomials of degree less or equal to  $q$ :  $K(\underline{x}_i, \underline{x}_j) = (\underline{x}_i^\top \underline{x}_j + 1)^q$

iii Gaussian RBF:  $K(\underline{x}_i, \underline{x}_j) = \exp(-\frac{\|\underline{x}_i - \underline{x}_j\|_2^2}{2\sigma^2})$

Consider again the value of alpha:

$$\begin{aligned}\underline{\alpha} &= (\underbrace{\Phi\Phi^\top}_{K \in \mathbb{R}^{n \times n}} + \lambda I)^{-1} \underline{y} \\ &= (K + \lambda I)^{-1} \underline{y}\end{aligned}$$

$$\Phi \in \mathbb{R}^{n \times d}$$

$$K_{ij} = \phi(\underline{x}_i)^\top \phi(\underline{x}_j) = K(\underline{x}_i, \underline{x}_j)$$

How to predict a label from a new sample?

$$\begin{aligned}\hat{y} &= \phi(\underline{x})^\top \hat{\underline{w}}_R \\ &= \hat{\underline{w}}_R^\top \phi(\underline{x}) \\ &= (\Phi^\top \underline{\alpha}) \phi(\underline{x}) \\ &= \underline{\alpha}^\top \Phi \phi(\underline{x}) \\ &= \sum_{i=1}^n \alpha_i \phi(\underline{x}_i)^\top \phi(\underline{x}) \\ &= \sum_{i=1}^n \alpha_i K(\underline{x}_i, \underline{x})\end{aligned}$$

We are doing a scalar product between the already present samples in the dataset and the new one. All weighted by the value of  $\alpha$  which represent how much each sample contribute to the definition of the decision boundary. Recall that the scalar product is a similarity measure so there will be samples more or less similar to the one to predict.



Small recap considering both cases in which we consider both following cases: the features are on the horizontal axes and the vertical one.

$$X \in \mathbb{R}^{n \times p} \text{ centered and } \begin{cases} n & \# \text{ of samples} \\ p & \# \text{ of features} \end{cases}$$

$$C = \frac{1}{n-1} X^\top X \text{ (spd)}$$

can be factorized:

$$C = V \Lambda V^\top$$

The columns of  $V$  are the eigenvectors of  $C$  and  $\Lambda$  is a diagonal matrix with the eigenvalues of  $C$ . The eigenvectors are also called **principal directions**. The **principal components** instead, are obtained as:

$$XV$$

And they represent the coordinates of the original samples in the new basis. If we write  $X = U \Sigma V^\top$  then

$$\begin{aligned} C &= \frac{1}{n-1} V \Sigma^\top U^\top U \Sigma V^\top \\ &= \frac{1}{n-1} V \Sigma^2 V^\top \end{aligned}$$

And we have also that:

$$\lambda_i = \frac{1}{n-1} \sigma_i^2$$

Moreover, with SVD:

$$XV = U \Sigma V^\top V = U \Sigma$$

Are still principal components.

$$X \in \mathbb{R}^{p \times n} \text{ centered and } \begin{cases} n & \# \text{ of samples} \\ p & \# \text{ of features} \end{cases}$$

Now, in this case we simply need to interchange the roles for  $U$  and  $V$ .

$$C = \frac{1}{n-1} X X^\top = \frac{1}{n-1} U \Sigma^2 U^\top$$

The principal directions are in  $U$  (and indeed not in  $V$  this time). To get the principal components we have to:

$$X^\top U$$

In Ridge we have found that:

$$\hat{w}_R = \arg \min_{\underline{w}} \|\underline{y} - \Phi \underline{w}\|_2^2 + \lambda \|\underline{w}\|_2^2$$

Let's consider a binary classification problem so we have  $(\underline{x}_i, y_i)$  and  $y_i \in \{-1, 1\}$ ,  $\underline{x}_i \in \mathbb{R}^p$ . Now consider Least squares for this problem:

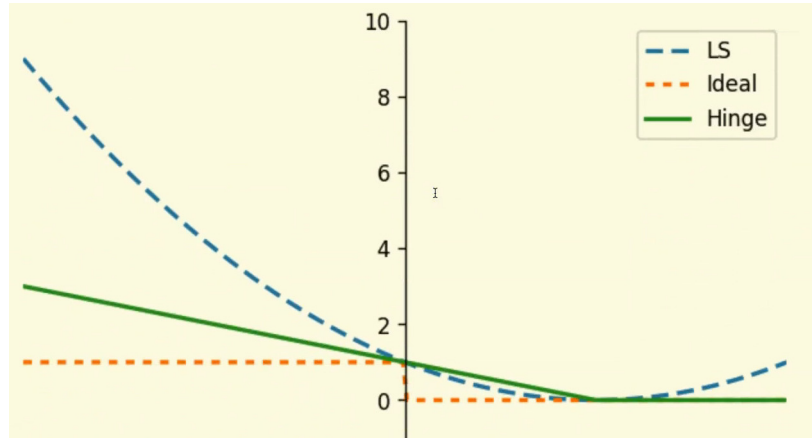
$$\hat{w}_{LS} = \arg \min_{\underline{w}} \sum_{i=1}^n (y_i - \underline{x}_i^\top \underline{w})^2 = \arg \min_{\underline{w}} \sum_{i=1}^n (1 - y_i \underline{x}_i^\top \underline{w})^2$$

You can easily check that the last equality holds in both values of  $y_i$  because the parenthesis is squared. There is an example on the notes that shows how the least squares is not a good idea for classification problems. Here are different loss functions we could use:

1. Square loss:  $L_i(y_i, \underline{x}_i^\top, \underline{w}) = (1 - y_i \underline{x}_i^\top \underline{w})^2$
2. Ideal loss:  $L_i(y_i, \underline{x}_i^\top, \underline{w}) = \begin{cases} 0 & y_i \\ 1 & y_i \underline{x}_i^\top \underline{w} < 0 \end{cases}$
3. Hinge loss:  $L_i(y_i, \underline{x}_i^\top, \underline{w}) = (1 - y_i \underline{x}_i^\top \underline{w})_+$

Where  $(a)_+ = \begin{cases} a & \text{if } a > 0 \\ 0 & \text{if } a \leq 0 \end{cases}$

Here is the plot of those functions:



On the x axis you have  $y_i \underline{X}_i^T \underline{w}$ . This picture can be divided in 3 points: if we are in the negative part (left part of the plot) this means that the predicted label is wrong and it has opposite sign with respect to the true label. In the middle part, the predicted label is correct but the margin is small. In the right part, the predicted label is correct and the margin is large. The square loss, differently from the other two, continue to grow even if the margin is large. The ideal loss is not differentiable and this is a problem for the optimization. The hinge loss is differentiable and it is the most used in practice. The ideal loss gives the same weight in the loss for a misclassified entry independently on our far is from the boundary. In the hinge instead the loss grows linearly with the distance and in the square grows squared.

## 14 Computing derivatives

Now we are going to consider different ways to compute the derivative of a function:

Method	PRO	CONS
Manual computation	Exact, good for the proofs	Prone to error, expensive for complex functions
Numerical differentiation	Easy to program	Floating point precision*, computational cost
Symbolic differentiation	Exact, good for proofs	Expression swelling**, not able to differentiate if conditions or loops
Automatic differentiation (AD)	Exact, fast	Implementation

\*: view jupyter notebook "Round\_off&Truncation.ipynb" for more details.

\*\* : big expressions of the derivative.

Consider this function:

$$y = f(x_1, x_2) = \sin((x_1 + x_2)x_2^2)$$

There are two possibilities for representing this function:

1. **Wengert list** for a generic function  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is defined as:

$$\begin{cases} \text{Input variables: } v_{i-n} = x_i & i = 1, \dots, n \\ \text{Intermediate variables: } v_i & i = 1, \dots, l \\ \text{Output variables: } v_{n-i} = v_{l-i} & i = m - 1, \dots, 0 \end{cases}$$

Now, let's go back to our example function, in this case its Wengert list formulation is the following:

$$\begin{array}{l} v_{-1} = x_1 = 1 \\ v_0 = x_2 = 2 \\ \hline v_1 = x_1 + x_2 = 3 \\ v_2 = x_2^2 = 4 \\ v_3 = v_1 v_2 = 12 \\ v_4 = \sin(v_3) = 0.536 \\ \hline y = v_4 \end{array}$$

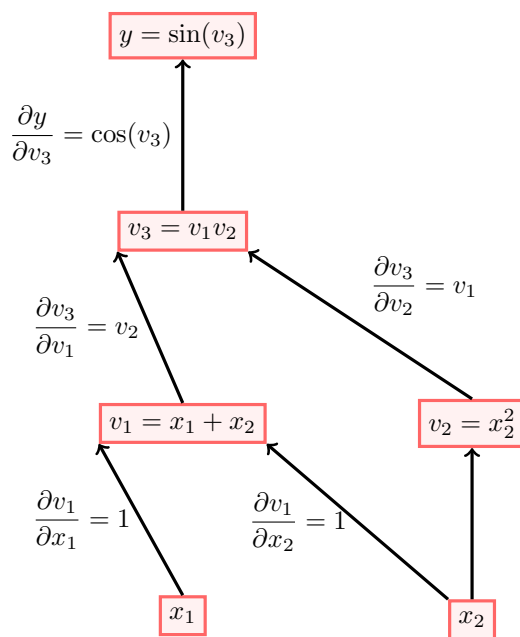
Notice that I've separated with the horizontal lines the different variables of the list. The input variables values have been chosen randomly just to compute the example with numbers.

How can we use the Wengert list for computing the derivative of a function? Let's consider the derivative of  $y$  with respect to  $x_1$ :

$$\begin{array}{l}
 \dot{v}_{-1} = \dot{x}_1 = 1 \\
 \dot{v}_0 = \dot{x}_2 = 0 \\
 \hline
 \dot{v}_1 = \dot{x}_1 + \dot{x}_2 = 1 \\
 \dot{v}_2 = 0 \\
 \dot{v}_3 = \dot{v}_1 v_2 + v_1 \dot{v}_2 = 4 \\
 \dot{v}_4 = \cos(v_3) \dot{v}_3 = 3.37 \\
 \hline
 \dot{y} = \dot{v}_4 = 3.37
 \end{array}$$

In this way we have splitted and computed only easy operations. This is called the **Forward Mode of Automatic Differentiation (AD)**. AD gives you the derivative computed for a certain point. If we change points the operations have to be done again.

## 2. Direct Acyclic Graph (DAG)



As you can see on the nodes there are the values of the variables and on the arrows there are the derivatives. The arrows are the partial derivatives of the function with respect to the variable of the node. **To compute the derivative of the function wrt a certain variable, you need to consider all paths that connect the top of the graph with the desired variable.** From  $y$  to  $x_1$  there is a single path so it's easy but if you consider  $y$  and  $x_2$  there are 2 possible paths, in this case you need to add them.

This is the **Reverse or Backward Mode of Automatic Differentiation (AD)**. It is more efficient than the forward mode because it computes the derivative of all the variables at the same time. It is used in neural networks because the number of variables is much higher than the number of samples.

Consider  $y = f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  and the two cases:  $\begin{cases} m \gg n \\ m \ll n \end{cases}$  In particular, the bottom case is the most interesting one to us since we might have many features in our model, especially for images and neural networks. For this case, the forward mode (FM) seen during last lecture is not very suitable because, considering NN, you need to compute a lot of derivatives wrt the weights.

Recall that, the Jacobian matrix is defined as:

$$\mathbf{J} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \cdots & \frac{\partial y_1}{\partial x_n} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \cdots & \frac{\partial y_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \frac{\partial y_m}{\partial x_2} & \cdots & \frac{\partial y_m}{\partial x_n} \end{bmatrix}$$

Suppose to have  $\underline{\dot{X}} = \underline{e}_i$  as initialization value for the Wengert list forward pass method. We would have:

$$\dot{y}_j = \frac{\partial y_j}{\partial x_i} \quad j = 1, \dots, n$$

You obtain a column of the jacobian. In general we are interested in the product of the jacobian times a vector which could be anything but often is the evaluation of the function at the previous iteration. It means that in practical computation, we are not interested in having the jacobian and then to perform the matrix-vector multiplication, but we want the result of it.

$$\mathbf{J} \underline{r}$$

**The FM of AD allows to have this result in a matrix free approach.** Any algorithm that is "matrix-free" means that in practise, even if the computation involves the presence of a matrix, you are able to avoid the construction of the complete matrix. You can resort to some tricks or algorithms that allow you to come up to the desired result without the need to explicitly computing the matrix. This is really good especially when dealing with big matrices.

Instead of  $\underline{\dot{X}} = \underline{e}_i$ , we can start from  $\underline{\dot{X}} = \underline{r}$  where  $\underline{r}$  is the vector we want to multiply the jacobian with.

## 14.1 Dual-numbers

They are similar to complex numbers.

$$a + b\epsilon \quad \text{with} \quad a, b \in \mathbb{R} \quad \text{and} \quad \epsilon \neq 0, \epsilon^2 = 0$$

Where  $a$  is the real part and  $b$  is the dual part. Basic operations with dual numbers:

$$(a + b\epsilon) + (c + d\epsilon) = (a + c) + (b + d)\epsilon$$

$$(a + b\epsilon)(c + d\epsilon) = ac + (ad + bc)\epsilon$$

Why this algebra is useful to our purposes? Consider the generic function  $f$ , we want to evaluate it with a dual number:

$$f(a + b\epsilon) = f(a) + f'(a)b\epsilon + \underbrace{\frac{f''(a)}{2}b^2\epsilon^2 + \dots}_0$$

We have used the Taylor expansion, the last terms are zero because of  $\epsilon^2$ . Notice that, if we have unitary dual part ( $b = 1$ ), we obtain the derivative of  $f$  in  $a$ . This is the key point of the dual numbers:

$$b = 1 \implies a + b\epsilon = a + \epsilon \implies f(a + \epsilon) = f(a) + f'(a)\epsilon$$

The last term is again a dual number in which the real part is the value of the function in  $a$  and the dual part is the derivative of the function in  $a$ .

Let's check if the derivative properties are actually respected, suppose to have  $f(x) = g(x)h(x)$ :

$$\begin{aligned} f(a + \epsilon) &= g(a + \epsilon)h(a + \epsilon) \\ &= [g(a) + g'(a)\epsilon][h(a) + h'(a)\epsilon] \\ &= g(a)h(a) + g(a)h'(a)\epsilon + g'(a)h(a)\epsilon + g'(a)h'(a)\epsilon^2 \\ &= g(a)h(a) + \underbrace{[g(a)h'(a) + g'(a)h(a)]\epsilon}_{\text{der. of the product}} \end{aligned}$$

Or  $f(x) = g(h(x))$ :

$$\begin{aligned} f(a + \epsilon) &= g(h(a + \epsilon)) \\ &= g(h(a) + h'(a)\epsilon) \\ &= \underbrace{g(h(a))}_{f(a)} + \underbrace{g'(h(a))h'(a)\epsilon}_{f'(a)\epsilon} \end{aligned}$$

Numerical example:

$$f(x) = \frac{1}{x} \implies f(x + \epsilon) = \frac{1}{x + \epsilon} = \frac{x - \epsilon}{(x + \epsilon)(x - \epsilon)} = \frac{x - \epsilon}{x^2 - \epsilon^2} = \frac{x - \epsilon}{x^2} = \underbrace{\frac{1}{x}}_{f(a)} - \underbrace{\frac{1}{x^2}\epsilon}_{f'(a)\epsilon}$$

Normally, we won't use the definition like this, but we will use real values, i.e.  $x = 2$  for example. There is formalism for which  $\epsilon$  are represented as matrices:

$$\epsilon = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \implies \epsilon^2 = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$$

But this is not so important because at the end we are interested in its coefficients.

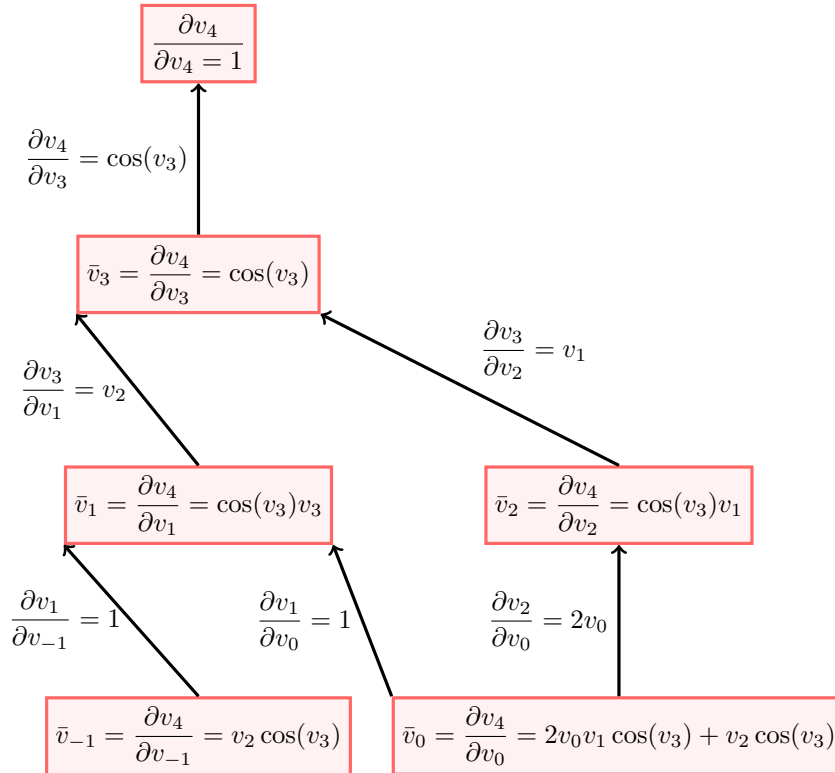
Let's now analyze the **Backward or reverse mode**.

It's related with the sensitivity of the output wrt the input. We introduce a new quantity:

$$\bar{v}_i = \frac{\partial y}{\partial v_i}$$

Where  $v_i$  is one of the variables we have written considering the forward mode and which might not be the input. This is what happen in backpropagation, a partial derivative of the output wrt to a certain weight of the NN.

We start from the output:



Where the bottom left box is  $\frac{\partial y}{\partial x_1}$  and the bottom right box is  $\frac{\partial y}{\partial x_2}$ . We have been able to compute all the intermediate sensitivities not just wrt of the input. If we have  $\mathbb{R}^n \rightarrow \mathbb{R}^m$  and  $n \gg m$  this method is very efficient because with just one sweep you compute the derivative of the output wrt all the input parameters. In the forward mode i would have had to perform  $n$  forward steps, and obviously this would be more costly.