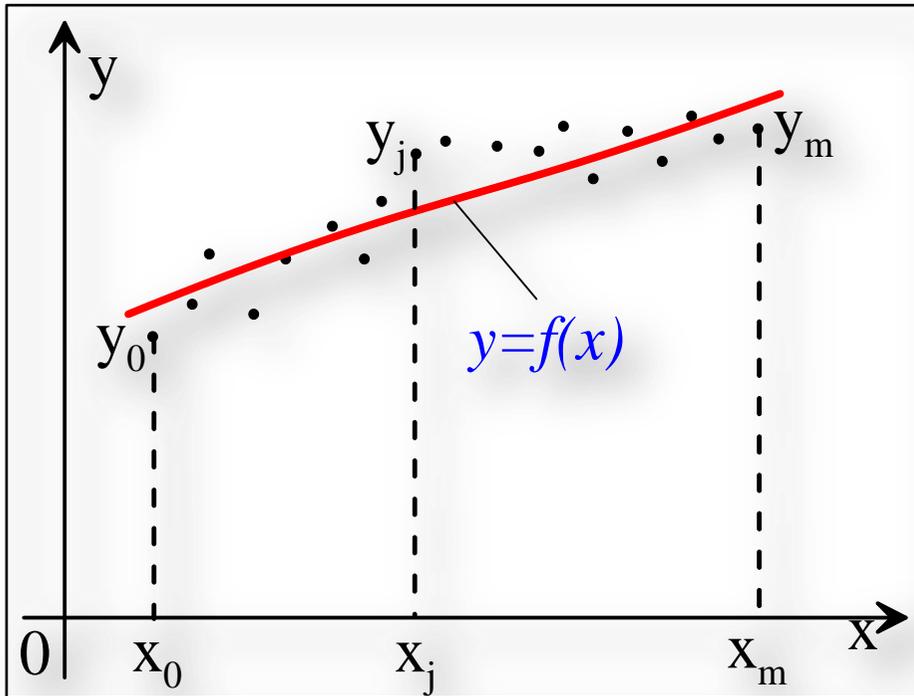


# LECTURE 2

# POLYNOMIAL APPROXIMATION



## FORMULATION OF THE LEAST-SQUARES APPROXIMATION PROBLEM



Consider the given set of the node

$$\{(x_k, y_k), k = 0, \dots, m\}$$

Typically, the number  $m$  is large, so the standard interpolation by the global polynomial is not sensible. Such polynomial would most likely exhibit rapid oscillations and it would be also extremely sensitive to small (and inevitable) inaccuracies in the input data.

**We will take a different approach. Our aim is to determine a function will capture “the general trend” (red line in the plot) presented by the given set, and such that its plot is – in a certain sense – close to all nodes.**

To this aim, we choose the prescribed finite set of the basic functions (e.g., polynomials, trigonometric functions, others)

$$\{\varphi_k(x), k = 0, \dots, n\}$$

and then look for the function in the form of the sum

$$f(x) = a_0\varphi_0(x) + a_1\varphi_1(x) + \dots + a_n\varphi_n(x) \equiv \sum_{j=0}^n a_j\varphi_j(x)$$

such that the following quantity is **minimal**

$$R = R(a_0, a_1, \dots, a_n) = \sum_{k=0}^m [f(x_k) - y_k]^2 = \sum_{k=0}^m \left[ \sum_{j=0}^n a_j\varphi_j(x_k) - y_k \right]^2 = \min$$

Such problem is called the **approximation in the sense of the least squares**.

**Necessary conditions for the minimum of the function  $R$  are**

$$\frac{\partial R}{\partial a_i} = 0 \quad , \quad i = 0, \dots, n$$

Then, the linear system for the unknown coefficients  $\{a_j, j = 0, \dots, n\}$  is obtained as follows

$$\frac{\partial R}{\partial a_i} = 2 \sum_{k=0}^m \varphi_i(x_k) \left[ \sum_{j=0}^n a_j \varphi_j(x_k) - y_k \right] = 0 \quad , \quad i = 0, 1, \dots, n$$

⇓

$$\sum_{j=0}^n \left[ \sum_{k=0}^m \varphi_i(x_k) \varphi_j(x_k) \right] a_j = \sum_{k=0}^m y_k \varphi_i(x_k)$$

This system can be written in the matrix/vector form as follows  $\mathbf{M}\mathbf{a} = \mathbf{z}$ , where

$$\mathbf{a} = [a_0, a_1, \dots, a_n]^T \quad , \quad M_{ij} = \sum_{k=0}^m \varphi_i(x_k) \varphi_j(x_k) = M_{ji} \quad , \quad i, j = 0, 1, \dots, n$$
$$z_i = \sum_{k=0}^m y_k \varphi_i(x_k) \quad , \quad i = 0, 1, \dots, n$$

In particular, the basic functions can be chosen as the monomials

$$\varphi_0(x) \equiv 1, \varphi_1(x) = x, \dots, \varphi_j(x) = x^j, \dots, \varphi_n(x) = x^n$$

Then

$$M_{ij} = \sum_{k=0}^m x_k^{i+j}, \quad i, j = 0, 1, \dots, n, \quad z_i = \sum_{k=0}^m y_k x_k^i, \quad i = 0, 1, \dots, n$$

Note that the matrix  $\mathbf{M}$  can be expressed as

$$\mathbf{M} = \mathbf{W}^T \mathbf{W}$$

where the elements of the matrix  $\mathbf{W}$  are

$$W_{kj} = x_k^j, \quad k = 0, \dots, m, \quad j = 0, \dots, m$$

Indeed, we have

$$(\mathbf{W}^T \mathbf{W})_{ij} = \sum_{k=0}^m (\mathbf{W}^T)_{ik} (\mathbf{W})_{kj} = \sum_{k=0}^m (\mathbf{W})_{ki} (\mathbf{W})_{kj} = \sum_{k=0}^m x_k^i x_k^j = (\mathbf{M})_{ij}$$

Also, the right-hand side vector  $\mathbf{z}$  can be expressed as

$$\mathbf{z} = \mathbf{W}^T \mathbf{y} \quad , \quad \mathbf{y} = [y_0, y_1, \dots, y_m]^T$$

This form of the obtained linear system is not accidental. In fact, the polynomial approximations problem can be viewed as the **over determined interpolation problem**.

We will take an algebraic rather than the analytical approach. Consider an over determined linear system in the form of

$$\mathbf{A}\mathbf{x} = \mathbf{b} \quad , \quad \dim(\mathbf{A}) = (m, n) \quad , \quad m > n$$

This means that the number of equations  $m$  is larger than the number of the unknowns  $n$  (the matrix  $\mathbf{A}$  is rectangular). Usually, so the above system of equations is – there is no vector  $\mathbf{x}$  in  $R^n$  such that all equations are **simultaneously** satisfied. Yet, there is a possibility to re-define the problem of determination of the vector  $\mathbf{x}$  in such a way, that there **exists** a solution. Moreover, if the rank of the matrix  $\mathbf{A}$  is equal  $n$  (i.e., it is the maximal possible) then the solution is **unique**!

The idea is simple: we will find the vector  $\mathbf{x}$  which makes the residual vector

$$\mathbf{r} = \mathbf{b} - \mathbf{A}\mathbf{x}$$

as small (short) as possible. In other words, the Euclidean norm of the vector  $\mathbf{r}$

$$\|\mathbf{r}\|_2 := \sqrt{r_1^2 + \dots + r_m^2}$$

must be **minimal**.

To this aim, we will refer to the concept of the **range of the matrix  $\mathbf{A}$** . The range of  $\mathbf{A}$  is defined as the following set of vectors (the linear subspace of  $R^m$ )

$$\text{range}(\mathbf{A}) := \{ \mathbf{y} \in R^m : \mathbf{y} = \mathbf{A}\mathbf{x}, \mathbf{x} \in R^n \}$$

Next, we also define the **kernel** of the transpose matrix  $\mathbf{A}^T$ . The kernel of  $\mathbf{A}^T$  is the following set of vectors (also the linear subspace of  $R^m$ )

$$\ker(\mathbf{A}^T) = \{ \mathbf{y} \in R^m : \mathbf{A}^T \mathbf{y} = \mathbf{0} \in R^n \}$$

Now, we claim that these two subspaces in  $R^m$  are **orthogonal**, i.e.

$$\text{range}(\mathbf{A}) \perp \ker(\mathbf{A}^T) \Leftrightarrow \forall_{\mathbf{v} \in \text{range}(\mathbf{A})} \forall_{\mathbf{w} \in \ker(\mathbf{A}^T)} (\mathbf{v}, \mathbf{w}) \equiv \sum_{j=1}^m v_j w_j = 0$$

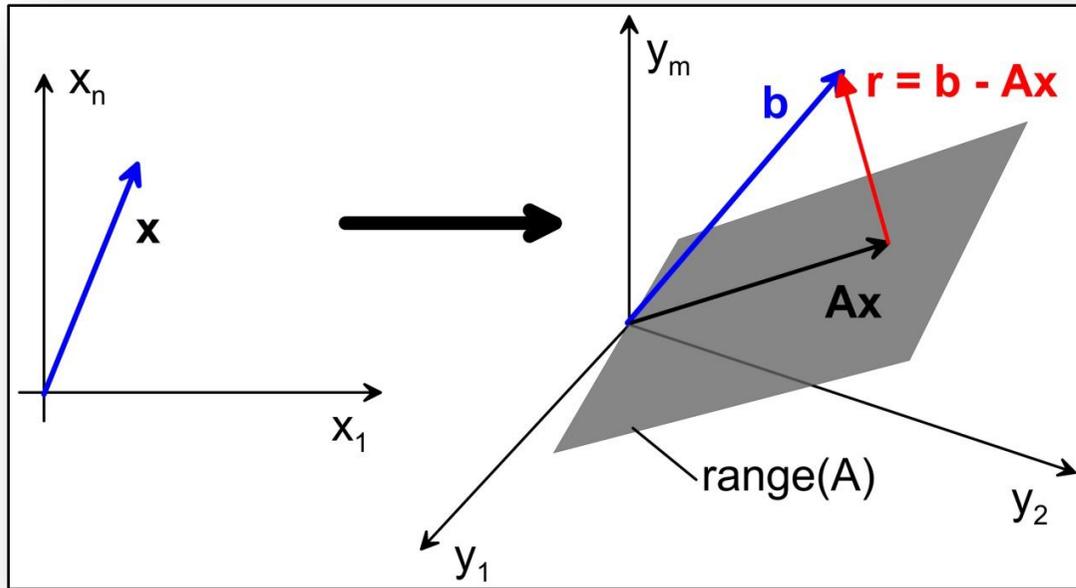
The **proof** of this crucial fact goes as follows. First, we note that

$$\mathbf{v} \in \text{range}(\mathbf{A}) \Leftrightarrow \exists_{\mathbf{p} \in R^n} \mathbf{v} = \mathbf{A}\mathbf{p}$$

Then the scalar product of the vectors  $\mathbf{v}$  and  $\mathbf{w}$  can be calculated

$$\begin{aligned} (\mathbf{v}, \mathbf{w}) &= (\mathbf{A}\mathbf{p}, \mathbf{w}) = \sum_{j=1}^m (\mathbf{A}\mathbf{p})_j w_j = \sum_{j=1}^m \left( \sum_{k=1}^n a_{jk} p_k \right) w_j = \sum_{k=1}^n \left( \sum_{j=1}^m a_{jk} w_j \right) p_k = \\ &= \sum_{k=1}^n \left( \sum_{j=1}^m a_{jk} w_j \right) p_k = \sum_{k=1}^n (\mathbf{A}^T \mathbf{w})_k p_k = (\mathbf{p}, \mathbf{A}^T \mathbf{w}) = (\mathbf{p}, \mathbf{0}) = 0 \end{aligned}$$

The conclusion follows from the fact that the choice of  $\mathbf{v}$  and  $\mathbf{w}$  is arbitrary.



Note that (see the picture) that the length (we say – the norm) of the residual vector  $\mathbf{r}$  is the smallest if this vector is perpendicular (we say – orthogonal) to the  $\text{range}(\mathbf{A})$ . But it means that the minimal residual vector  $\mathbf{r}$  belongs to the  $\text{ker}(\mathbf{A}^T)$ .

Thus, we have

$$\|\mathbf{r}\|_2 \equiv \sqrt{\sum_{k=1}^m r_k^2} = \min \Leftrightarrow \mathbf{r} \perp \text{range}(\mathbf{A}) \Leftrightarrow \mathbf{r} \in \text{ker}(\mathbf{A}^T)$$

⇓

$$\mathbf{A}^T (\mathbf{b} - \mathbf{A}\mathbf{x}) = 0 \Rightarrow \mathbf{A}^T \mathbf{A}\mathbf{x} = \mathbf{A}^T \mathbf{b}$$

In the context of the **polynomial approximation problem**, we can see that  $A = W$ . On the other hand, it is clear that the **over determined** linear system  $Wa = y$  corresponds exactly to the “impossible” interpolation conditions

$$y_j = f(x_j) = \sum_{k=0}^n x_j^k a_k = \sum_{k=0}^n W_{jk} a_k \quad , \quad j = 0, 1, \dots, m$$

The above mentioned method to solve the polynomial approximation problem is called the method of the **normal equations**.

In the remaining part of these notes, we present an alternative and computationally more effective method of the **orthogonal polynomials**.

## METHOD OF THE ORTHOGONAL POLYNOMIALS

For the given sets of the nodes, one can define the inner (scalar) product of the functions as follows

$$\langle \varphi_i, \varphi_j \rangle := \sum_{k=0}^m \varphi_i(x_k) \varphi_j(x_k)$$

We say that the basic functions are (discretely) orthogonal on the given set of the nodes if and only if

$$\langle \varphi_i, \varphi_j \rangle = 0 \quad , \quad i \neq j$$

Using the orthogonal functions is advantageous because the corresponding matrix  $\mathbf{M}$  of the system of the normal equations is purely **diagonal**

$$M_{ij} = \begin{cases} \sum_{k=0}^m \varphi_i^2(x_k) & , \quad i = j \\ 0 & , \quad i \neq j \end{cases}$$

This means that the equations of this system are not coupled and they can be solved separately one by one. Indeed, one gets

$$\left( \sum_{k=0}^m \varphi_i^2(x_k) \right) a_i = \sum_{k=0}^m y_k \varphi_i(x_k) \quad , \quad i = 0, \dots, n$$

⇓

$$a_i = \frac{\sum_{k=0}^m y_k \varphi_i(x_k)}{\sum_{k=0}^m \varphi_i^2(x_k)}$$

The question remains how to construct efficiently orthogonal basic functions for a given set of the nodes. Since such options is particularly interesting, we will explain how to generate an **orthogonal set of polynomials**.

**The computational procedure is recursive.** It means that consecutive polynomials will be constructed by the use of those previously defined.

In order to initiate such procedure, two first polynomials must be prescribed. These polynomials are

$$q_0(x) \equiv 1 \quad , \quad q_1(x) = x - \alpha_1$$

The number  $\alpha_1$  is selected in order to assure that  $q_0$  and  $q_1$  are orthogonal:

$$\langle q_0, q_1 \rangle = 0 \Rightarrow \sum_{k=0}^m (x_k - \alpha_1) = 0 \Rightarrow \alpha_1 = \frac{1}{m+1} \sum_{k=0}^m x_k$$

Higher-order polynomials are obtained via two-point recurrence as follows

$$q_{j+1}(x) = xq_j(x) - \alpha_{j+1}q_j(x) - \beta_jq_{j-1}(x)$$

The coefficients  $\alpha_{j+1}$  and  $\beta_j$  are chosen to assure that the polynomial  $q_{j+1}$  is orthogonal to the previously defined polynomials  $q_j$  and  $q_{j-1}$ :

$$\langle q_{j+1}, q_j \rangle = 0 \quad , \quad \langle q_{j+1}, q_{j-1} \rangle = 0$$

These conditions lead to the following formulas:

$$\alpha_{j+1} = \frac{\langle xq_j, q_j \rangle}{\langle q_j, q_j \rangle} = \frac{\sum_{k=0}^m x_k q_j^2(x_k)}{\sum_{k=0}^m q_j^2(x_k)}, \quad \beta_j = \frac{\langle xq_j, q_{j-1} \rangle}{\langle q_{j-1}, q_{j-1} \rangle} = \frac{\sum_{k=0}^m x_k q_{j-1}(x_k) q_j(x_k)}{\sum_{k=0}^m q_{j-1}^2(x_k)}$$

**The crucial question arises: how come that  $\langle q_{j+1}, q_i \rangle = 0$  ,  $i = 0, 1, \dots, j-2$  ?!**

Surprisingly enough, this property holds “automatically”! To see this, we will calculate the inner product of the polynomial  $q_{j+1}$  and  $q_k$  for  $k < j-1$ :

$$\langle q_{j+1}, q_i \rangle = \langle xq_j, q_i \rangle - \underbrace{\alpha_{j+1} \langle q_j, q_i \rangle}_0 - \underbrace{\beta_j \langle q_{j-1}, q_i \rangle}_0 = \underbrace{\langle q_j, xq_i \rangle}_{xq_i(x) \text{ has order } i+1 < j} = \sum_{p=0}^{i+1 < j} \gamma_p \underbrace{\langle q_j, q_p \rangle}_{=0 \text{ for each } p < j} = 0$$

Final form of the approximating polynomial can be written as

$$f(x) = \sum_{j=0}^n \left[ \frac{\sum_{k=0}^m y_k q_j(x_k)}{\sum_{k=0}^m q_j^2(x_k)} \right] q_j(x)$$

For efficient and stable method of evaluation of the polynomial  $f(x)$  see the **Clenshaw recurrence formula** (see, for instance, Numerical Recipes in C++, 3<sup>rd</sup> Ed., p. 222)