

Constructive Approximation of Functions

Benjamin Shih
Spring 2024

Abstract. These are notes from the Spring 2024 Brown University APMA DRP on the constructive approximation of functions, specifically polynomials and rational functions of one variable, under the direction of [Dr. Wenjun Zhao](#). We closely followed the text *Approximation Theory and Approximation Practice* by Trefethen [1] and additionally engaged in material on universal approximation [2] as it relates to neural networks out of interest.

Contents

1	Week 1: 2/22/2024	1
1.1	Chebyshev Series and Polynomials	1
1.2	Introduction to Linear Approximation with Neural Networks	2
2	Week 2: 2/29/2024	3
2.1	Interpolants and projections	3
2.2	Barycentric Interpolation Formula	4
2.3	Weierstrass Approximation Theorem	5
3	Week 3: 3/7/2024	5
3.1	Convergence for Differentiable Functions	5
3.2	Gibbs Phenomenon	6
3.3	Density of Two Layer Neural Networks	6
4	Week 4: 3/14/2024	8
4.1	Best Approximation	8
5	Skipped/To-do	8
6	Bibliography	9

1 Week 1: 2/22/2024

Chebyshev interpolants have a very strong approximation properties, as opposed to uniformly spaced points. They are the points that correspond to the real part of equispaced points on the unit circle in the complex plane. That is, the Chebyshev points are

$$x_j = \cos\left(\frac{j\pi}{n}\right)$$

and a key property is that they “collect” near the ends of the interval in higher density. Namely, this key property is that each point is, on average, the same distance away from every other point. For the most part, we deal with approximating functions on the interval $[-1, 1]$, which any function on any interval $[a, b]$ can be scaled to.

There are connections that can be drawn between the Chebyshev, Fourier, and Laurent settings, with each being used in numerical, complex, and real analysis heavily, respectively. In the Chebyshev settings, we approximate functions $f(x), x \in [-1, 1]$ with the form

$$f(x) \approx \sum_{k=0}^n a_k T_k(x)$$

while using $z \in S^1 \subset \mathbb{C}$ equispaced points on the complex plane gives us the Laurent setting with Laurent polynomials

$$F(z) = F(z^{-1}) = \frac{1}{2} \sum_{k=0}^n a_k (z^k + z^{-k})$$

and finally using the angle $\theta \in [-\pi, \pi]$ to define $\mathcal{F}(\theta) = F(e^{i\theta}) = f(\cos(\theta))$ gives us Fourier series as

$$\mathcal{F}(\theta) \approx \frac{1}{2} \sum_{k=0}^n a_k (e^{ik\theta} + e^{-ik\theta})$$

Their corresponding canonical grid systems are as follow:

Chebyshev points	$x_j = \cos\left(\frac{j\pi}{n}\right), \quad 0 \leq j \leq n$
Roots of unity (Laurent)	$z_j = e^{\frac{j\pi}{n}}, \quad -n+1 \leq j \leq n$
Equispaced points (Fourier)	$\theta_j = \frac{j\pi}{n}, \quad -n+1 \leq j \leq n$

1.1 Chebyshev Series and Polynomials

Definition 1.1 (k -th Chebyshev polynomial). The k -th Chebyshev polynomial is the real part of the function z^k on the unit circle; i.e.

$$T_k(x) = \Re(z^k) = \frac{1}{2}(z^k + z^{-k}) = \cos(k\theta)$$

Theorem 1.2 (Existence of Chebyshev Series). Suppose that f is Lipschitz on $[-1, 1]$, i.e. that there exists $C \in \mathbb{R}$ such that $|f(x) - f(y)| \leq C|x - y|$ for any $x, y \in \mathbb{R}$. Then f admits a unique

representation as a Chebyshev series

$$f(x) = \sum_{k=0}^{\infty} a_k T_k(x)$$

which is absolutely and uniformly convergent. The coefficients a_k are given by

$$a_k = \frac{2}{\pi} \int_{-1}^1 \frac{f(x) T_k(x)}{\sqrt{1-x^2}} dx$$

for $k \geq 1$, and for $k = 0$ by the same formula with a $\frac{1}{\pi}$ factor instead.

1.2 Introduction to Linear Approximation with Neural Networks

Generally, we will work with compact and convex domains $X \subset \mathbb{R}^d$ and target functions $f : X \rightarrow \mathbb{R}$ from some function space $\mathcal{F}(X)$. We denote by $f_\theta : \mathbb{R}^d \rightarrow \mathbb{R}$ a neural network with n_θ parameters. The study of the subject of neural network approximation is mainly concerned with the following three problems:

1. **Density:** When $n_\theta \rightarrow \infty$, does there exist f_θ that approximates f well?
2. **Covergence:** If $n_\theta < \infty$ is fixed, how close can f_θ be to f ?
3. **Complexity:** If we want $\|f_\theta - f\| < \epsilon$, how large is n_θ ?

Before discussing neural networks directly, we motivate with some background of density in polynomial approximation.

Definition 1.3 (Uniform Convergence). A sequence of functions $\{f_n\}_{n=1}^{\infty}$ is said to converge uniformly to a limiting function f on a set X if for all $\epsilon > 0$, there exist $N \in \mathbb{N}$ such that $|f(x) - f_m(x)| < \epsilon$ for all $m \geq N$ and all $x \in X$ and denote this by

$$f_m \rightarrow f \quad \text{uniformly}$$

We can also equivalently define uniform convergence in terms of the supremum or infinity norm, where

$$\|f\|_\infty = \sup_{x \in X} |f(x)|$$

by the condition $\|f - f_m\|_\infty \rightarrow 0$. Weierstrass gave the following result, which says that any continuous function f on a closed subinterval of the real line can be approximated arbitrarily well by an algebraic polynomial.

Theorem 1.4 (Weierstrass). Given a function $f \in C([a, b])$ and $\epsilon > 0$, there exists an algebraic polynomial p such that $|f(x) - p(x)| < \epsilon$ for all $x \in [a, b]$, or equivalently, $\|f - p\|_\infty < \epsilon$.

It should be noted that a similar result exists for 2π -periodic continuous functions and trigonometric polynomials; i.e., that trigonometric polynomials are dense in the class of 2π -periodic continuous functions.

2 Week 2: 2/29/2024

2.1 Interpolants and projections

Take $f(x)$ to be Lipschitz continuous on $[-1, 1]$ with Chebyshev coefficients $\{a_k\}$ so that

$$f(x) = \sum_{k=0}^{\infty} a_k T_k(x)$$

Then we can approximate f in the space of n -degree polynomials by *interpolation*; that is, as

$$p_n(x) = \sum_{k=0}^n c_k T_k(x)$$

Another approximation of f is the *truncation* or *projection* to degree n , where the coefficients through degree n match those of f :

$$f_n(x) = \sum_{k=0}^n a_k T_k(x)$$

Theorem 2.1 (Aliasing of Chebyshev polynomials). For any $n \geq 1$ and $0 \leq m \leq n$, the following Chebyshev polynomials take the same values on the $(n+1)$ -point Chebyshev grid:

$$T_m, T_{2n-m}, T_{2n+m}, T_{4n-m}, T_{4n+m}, T_{6n-m}, \dots$$

Equivalently, for any $k \geq 0$, T_k takes the same value on the grid as T_m with

$$m = |(k + n - 1) \pmod{2n} - (n - 1)|$$

which is a number in the range $0 \leq m \leq n$.

This leads to the connection between the coefficients of the polynomial approximant $\{a_k\}$ and the Chebyshev coefficients $\{c_k\}$.

Theorem 2.2 (Aliasing formula for Chebyshev coefficients). Let f be Lipschitz continuous on $[-1, 1]$, and p_n be its Chebyshev interpolant in \mathcal{P}_n . Let $\{a_k\}$ and $\{c_k\}$ be the Chebyshev coefficients of f and p_n . Then

$$c_0 = a_0 + a_{2n} + a_{4n} + \dots$$

$$c_n = a_n + a_{3n} + a_{5n} + \dots$$

and for $1 \leq k \leq n-1$,

$$c_k = a_k + (a_{k+2n} + a_{k+4n} + \dots) + (a_{-k+2n} + a_{-k+4n} + \dots)$$

Essentially, this says that any f is indistinguishable from a polynomial interpolant of degree n on the $(n+1)$ point grid obtained by reassigning all Chebyshev coefficients $\{a_k\}$ to their aliases up to degree

n . The errors between the two different approximations are

$$\begin{aligned} f(x) - f_n(x) &= \sum_{k=n+1}^{\infty} a_k T_k(x) \\ f(x) - p_n(x) &= \sum_{k=n+1}^{\infty} a_k (T_k(x) - T_m(x)) \end{aligned}$$

where $m = |(k + n - 1) \pmod{2n} - (n - 1)|$ which are absolutely convergent. We note that f_n often leads to a approximation on the basis of relative error to f compared to p_n .

2.2 Barycentric Interpolation Formula

We're now interested in evaluating Chebyshev interpolants; multiple approaches exist and range from $O(n \log n)$ to $O(n)$ work. We focus on the latter, called the *barycentric interpolation formula*, which is direct and numerically stable. The formula takes the form

$$p(x) = \sum_{j=0}^n f_j \ell_j(x) \quad (1)$$

which is in alternative Lagrange form and is the linear combination of unique *Lagrange* or *cardinal* polynomials $\ell_j \in \mathcal{P}_n$ defined as

$$\ell_j(x_k) = \begin{cases} 1 & \text{if } k = j \\ 0 & \text{if } k \neq j \end{cases}$$

In fact, we have an exact expression for ℓ_j :

$$\ell_j(x) = \frac{\prod_{k \neq j} (x - x_k)}{\prod_{k \neq j} (x_j - x_k)} \quad (2)$$

However, computational complexity wise, (2) is not so good, but we can remedy this.

Definition 2.3 (Node Polynomial). The *node polynomial* $\ell \in \mathcal{P}_{n+1}$ for a given grid is

$$\ell(x) = \prod_{k=0}^n (x - x_k)$$

Using this, (2) then becomes

$$\ell_j(x) = \frac{\ell(x)}{\ell'(x_j)(x - x_j)}$$

Making the substitutions

$$\lambda_j = \frac{1}{\prod_{k \neq j} (x_j - x_k)}, \quad \ell_j(x) = \ell(x) \frac{\lambda_j}{x - x_j}$$

(1) becomes the “type 1 barycentric formula”:

$$p(x) = \ell(x) \sum_{j=0}^n \frac{\lambda_j}{x - x_j} f_j$$

2.3 Weierstrass Approximation Theorem

Chapter 6 in the text goes over the Weierstrass Approximation Theorem, which we covered in the neural network survey paper as Theorem 1.4. We note that the result has been generalized by theorem due to Runge and Mergelyan, which state that a function f defined on a compact set $K \subset \mathbb{C}$ with connected complement which is continuous on K and analytic throughout K (resp. interior of K), then f can be approximated on K by polynomials.

3 Week 3: 3/7/2024

3.1 Convergence for Differentiable Functions

Generally, we accept the following to be a motif in approximation:

The smoother a function, the faster its approximants converge as $n \rightarrow \infty$.

Generally, if a function $f \in C^k([-1, 1])$, then convergence of its Chebyshev interpolants occurs on the rate of $O(n^{-k})$ with respect to the supremum norm. Note that if f is not *continuously* differentiable, then may not be true. However, if f has *bounded variation* (f has bounded variation if its total variation—the 1-norm of the derivative—is finite), then it enjoys this convergence.

Definition 3.1 (Absolute continuity (heuristic)). A function f is absolutely continuous if it is equal to the integral of its derivative, which exists almost everywhere and is Lebesgue integrable.

With this definition, we can establish a bound on the Chebyshev coefficients for differentiable functions.

Theorem 3.2 (Chebyshev coefficients for differentiable functions). For an integer $\nu \geq 0$, let f and its derivatives through $f^{(\nu-1)}$ be absolutely continuous on $[-1, 1]$ and suppose the ν -th derivative $f^{(\nu)}$ is of bounded variation V . Then for $k \geq \nu + 1$, the Chebyshev coefficients of f satisfy

$$|a_k| \leq \frac{2V}{\pi k(k-1) \cdot (k-\nu)} \leq \frac{2V}{\pi(k-\nu)^{\nu+1}}$$

With this bound on the coefficients, we establish the following convergence result:

Theorem 3.3 (Convergence for differentiable functions). If f satisfies the conditions of Theorem 3.2, then with V denoting the total variation of $f^{(\nu)}$ for some $\nu \geq 1$, then for any $n > \nu$, its Chebyshev projections satisfy

$$\|f - f_n\| \leq \frac{2V}{\pi\nu(n-\nu)^\nu}$$

and its Chebyshev interpolants satisfy

$$\|f - p_n\| \leq \frac{4V}{\pi\nu(n-\nu)^\nu}$$

3.2 Gibbs Phenomenon

We now study the *Gibbs Phenomenon*, which describes how polynomial interpolants and projections oscillate and overshoot near discontinuities. Note that as $n \rightarrow \infty$, the overshoot gets narrower, but not shorter!

Theorem 3.4 (Gibbs Phenomenon for Chebyshev interpolants). Let p_n be the degree n Chebyshev interpolant of the function $f(x) = \text{sign}(x)$ on $[-1, 1]$. Then as $n \rightarrow \infty$,

$$\lim_{n \rightarrow \infty, n \text{ odd}} \|p_n\| = c_1 \approx 1.282$$

$$\lim_{n \rightarrow \infty, n \text{ even}} \|p_n\| = c_2 \approx 1.066$$

These values are the height or maximum attained by the Chebyshev interpolant on $f(x)$.

3.3 Density of Two Layer Neural Networks

We consider the family of two-layer feedforward networks consisting of:

- d input neurons
- r neurons in one hidden layer using the same activation $\sigma : \mathbb{R} \rightarrow \mathbb{R}$
- one output neuron with no activation nor bias.

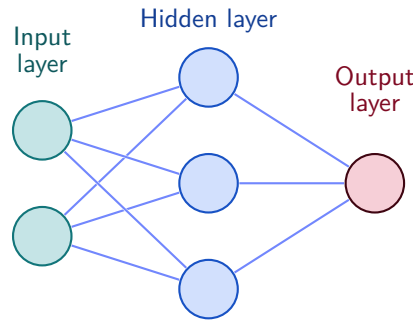


Figure 1: A feed-forward neural network with one hidden layer, where $d = 2, r = 3$.

If we let $x \in \mathbb{R}^d$ be the input to the network, then the output $y \in \mathbb{R}$ can be written

$$y = \sum_{i=1}^r W_{1,i}^2 \sigma \left(\sum_{j=1}^d W_{i,j}^1 x_j + b_i \right)$$

where $W^1 \in \mathbb{R}^{r \times d}$ and $W^2 \in \mathbb{R}^{1 \times r}$ are the weight matrices for the layer and $b \in \mathbb{R}^r$ is the bias in the hidden layer. Then the main density result for such networks is as follows.

Theorem 3.5 (Pinkus). Let

$$\mathcal{M}(\sigma) = \text{span}\{\sigma(w \cdot x + b) : w \in \mathbb{R}^d, b \in \mathbb{R}\}$$

where $\sigma \in C(\mathbb{R})$. Then $\mathcal{M}(\sigma)$ is dense in $C(\mathbb{R}^n)$ with respect to the supremum norm on compact sets, if and only if σ is not a polynomial.

Essentially, Pinkus' theorem states that the network space $\mathcal{M}(\sigma)$ is dense in the space of continuous functions with respect to the supremum norm as long as the class of activation function chosen is non-polynomial. This means that given a target function $f \in C(\mathbb{R}^d)$ and a compact subset $X \subset \mathbb{R}^d$, for any $\epsilon > 0$, there exists $g \in \mathcal{M}(\sigma)$ so that

$$\sup_{x \in X} |f(x) - g(x)| < \epsilon$$

Examples of such σ which would satisfy the necessary conditions of Pinkus' theorem would be $\sin(\cdot)$, $\cos(\cdot)$, etc. Let us consider the proof of the 1-dimensional case of the theorem; the extension to multiple dimensions can be found in [2].

Proof of the 1-dimensional case of Pinkus' theorem In this proof we consider the 1-dimensional case of $\mathcal{M}(\sigma)$, which we will denote

$$\mathcal{N}(\sigma) = \text{span}\{(wx + b) : w, b \in \mathbb{R}\}$$

We will show that for any non-polynomial σ , $\mathcal{N}(\sigma)$ is dense in $C(\mathbb{R})$. We employ the following lemma:

Lemma 3.6. Let $\sigma \in C^\infty((a, b))$ where $(a, b) \subset \mathbb{R}$ is an open interval. If for every point $x \in (a, b)$, there exists an integer $k = k(x)$ such that $\sigma^{(k)}$ vanishes at x , i.e. $\sigma^{(k)}(x) = 0$, then σ is a polynomial.

By our assumption, we have σ is smooth but nonpolynomial; hence, Lemma 3.6 implies that there exists a point $c \in \mathbb{R}$ at which $\sigma^{(k)} \neq 0$ for $k = \mathbb{Z}_+ \cup \{0\}$. Then

$$\frac{\sigma((w+h)x+c) - \sigma(wx+c)}{h} \in \mathcal{N}(\sigma), \quad \forall h \neq 0$$

where $w \in \mathbb{R}$. Taking the limit, we have

$$\left. \frac{d}{dw} \sigma(wx+c) \right|_{w=0} = x\sigma'(c) \in \overline{\mathcal{N}(\sigma)}$$

where $\overline{\mathcal{N}(\sigma)}$ is the closure of $\mathcal{N}(\sigma)$. In general, we have

$$\left. \frac{d^k}{dw^k} \sigma(wx+c) \right|_{w=0} = x^k \sigma^{(k)}(c) \in \overline{\mathcal{N}(\sigma)}$$

and since $\sigma^{(k)} \neq 0$ for $k \in \mathbb{N} \cup \{0\}$, $\overline{\mathcal{N}(\sigma)}$ is a set of monomials. Then we conclude by Theorem 1.4 that $\overline{\mathcal{N}(\sigma)}$ (and hence $\mathcal{N}(\sigma)$) is dense in $C(\mathbb{R})$.

4 Week 4: 3/14/2024

4.1 Best Approximation

Definition 4.1 (Best Approximant). The *best approximant* of a function f is a polynomial p^* with a specified degree n such that p^* minimizes the ∞ -norm on some interval with f .

This p^* is unique; however, Chebyshev interpolants are often as good or even better (that is, best approximation under the ∞ -norm may not be as useful as it sounds). Best approximations hold a property called the *equioscillation* property; that is, the error curve of the best approximant attains extreme magnitudes with alternating signs at a succession of values of x .

Theorem 4.2 (Equioscillation characterization of best approximants). For $f \in C([-1, 1])$, there is a unique best approximation $p^* \in \mathcal{P}_n$. If $f : [-1, 1] \rightarrow \mathbb{R}$, then p^* is real as well, and a polynomial $p \in \mathcal{P}_n$ is equal to p^* if and only if $f - p$ equioscillates in at least $n + 2$ extreme points.

5 Skipped/To-do

End of Chapter 4, Chapter 8, Chapter 11 (due to complex variables)

6 Bibliography

- [1] Lloyd N. Trefethen. *Approximation Theory and Approximation Practice (Other Titles in Applied Mathematics)*. Society for Industrial and Applied Mathematics, USA, 2012. ISBN 1611972396.
- [2] Mohammad Motamed. Approximation power of deep neural networks: an explanatory mathematical survey. preprint at <https://arxiv.org/pdf/2207.09511.pdf>, 2022.