

Spectral Pollution

Alex H. Room

December 15, 2023

Abstract Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

CONTENTS

1	Introduction: the spectrum and its approximation	1
1.1	Spectra	2
1.2	Approximating spectra; Ritz and Galerkin methods	3
1.3	Spectral pollution	3
2	Spectral pollution in a multiplication operator	4
2.1	The spectrum of a multiplication operator	5
2.1.1	Estimating spectra computationally	6
2.1.2	The structure of approximations	6
2.2	Toeplitz operators	7
2.2.1	Toeplitz operators, Toeplitz matrices, and their equivalence	7
2.2.2	The spectrum of a Toeplitz operator with real symbol	8
2.2.3	Multiplication operators, Toeplitz operators, and spectral pollution	9
3	Bounding spectral pollution	10
3.1	The essential spectrum of an operator	11
3.2	Essential spectrum and compact perturbation	13
3.3	Numerical range and essential numerical range	14
3.4	Essential numerical range	15
4	Detecting spectral pollution	16
4.1	Dissipative barrier methods	17

The computation of spectra can boldly be considered the 'fundamental problem of operator theory' [1]. The spectrum of an operator holds the same key to the entire structure as eigenvalues do for linear algebra, however (as often occurs when passing from the finite- to infinite-dimensional case) we lose the ease of representation that allows the creation of such algorithms and formulae as those used for matrices.

1.1 Spectra

We must first define our quantity of interest: the spectrum of an operator.

Definition. (*Resolvent and spectrum*) (Adapted from [2]) Let T be a linear operator on a Banach space.

- The resolvent of T is the set $\rho(T) := \{\eta \in \mathbb{C} : (T - \eta I) \text{ has a bounded inverse}\}$, where I is the identity operator. If it exists, we call its inverse the 'resolvent operator' for η and T .
- The spectrum of T , denoted $\text{Spec}(T)$, is $\mathbb{C} \setminus \rho(T)$, i.e. the set of all complex numbers λ such that the operator $(T - \lambda I)$ does not have a bounded inverse.

In the remainder of the text, the identity operator I will be implicit in the operator - i.e. we will simply write $(T - \lambda I)$ as $(T - \lambda)$.

One can see that this concept generalises the eigenvalues of a matrix to any Banach space, and that for a finite-dimensional Banach space (of which \mathbb{R}^n and \mathbb{C}^n are major examples) we can recover the eigenvector equation. Indeed, the term 'resolvent' originates from its relation to the solution of certain types of differential equation; if for some differential operator L we have the equation $Lu = \lambda u + g$, a solution u would be related to g as $(L - \lambda)^{-1}g = u$.

For an operator on an infinite-dimensional space, values which satisfy the eigenvector equation *are* in the spectrum of the operator, but they do not comprise the entire spectrum - nor do we necessarily have a spectrum made up of a discrete set of values.

Example 1. Let S be the 'right-shift' operator on the sequence space $\ell^2(\mathbb{N})$, which has the following action: $S(x_1, x_2, x_3, \dots) = (0, x_1, x_2, x_3, \dots)$; that is, for a sequence u we map u_n to u_{n+1} , and pad the vector with a zero in u_1 's place. S has no eigenvalues, and its spectrum is the open unit disc $D = \{z \in \mathbb{C} : |z| < 1\}$.

Proof. (sketch) For the lack of eigenvalues; we see that if $Su = \lambda u$,

$$(0, u_1, u_2, u_3, \dots) = (\lambda u_1, \lambda u_2, \lambda u_3, \dots)$$

so in particular, we would require $\lambda u_1 = 0$. Then either $\lambda = 0$ and thus $u_n = 0$ for all n , or as $\lambda \neq 0$, $u_1 = 0$. Then $\lambda u_2 = u_1 = 0$, and continuing this we can see that $u_n = 0$ for all n , i.e. u must be the zero vector. Thus there is no non-zero vector satisfying the eigenvector equation.

To calculate its spectrum, we make use of a natural relation between the spectrum of an operator T and the spectrum of its adjoint:

$$\begin{aligned} \lambda \in \text{Spec}(T) &\text{ iff } (T - \lambda) \text{ is not invertible} \\ &\text{ iff } (T - \lambda)^* \text{ is not invertible (this is true for any operator)} \\ &= (T^* - \bar{\lambda}) \text{ not invertible, i.e. } \bar{\lambda} \in \text{Spec}(T^*). \end{aligned}$$

in this case, the adjoint is the left-shift operator $S^*u = (u_2, u_3, u_4, \dots)$; we can see that for any value λ we have for the vector $(\lambda, \lambda^2, \lambda^3, \dots)$

$$S^*(\lambda, \lambda^2, \lambda^3, \dots) = (\lambda^2, \lambda^3, \lambda^4, \dots) = \lambda(\lambda, \lambda^2, \lambda^3, \dots)$$

which is in $\ell^2(\mathbb{N})$ and is thus an eigenvector iff $|\lambda| < 1$; hence the unit disc is in $\text{Spec}(S^*)$ and hence its complex conjugate (which is also the unit disc) is in $\text{Spec}(S)$. \square

As we have alluded to, there is no universal algorithm for the calculation of operator spectra in the way that there is the QR algorithm [3] for matrices. To devise a formula for the spectrum of even a specific subset of a class of operators is a mathematical feat, and varieties of operators important to fields such as quantum physics ([4]), hydrodynamics ([5]), and crystallography ([6]) still withhold the structure of their spectra from decades-long attempts at discovery. To this end, we must employ numerical methods. We shall find that even approximating spectra computationally is not so easy.

1.2 Approximating spectra; Ritz and Galerkin methods

Definition. (Compressions, truncations, and Ritz matrices) (Adapted from [7]) Let T be an operator on a Hilbert space H , $\mathcal{L} \subseteq \text{Dom}(T)$ a closed linear subspace, and $P_{\mathcal{L}}$ the orthogonal projection of H onto \mathcal{L} .

- The **compression** of the operator T , which we will often denote $T_{\mathcal{L}}$, is defined

$$T_{\mathcal{L}} := P_{\mathcal{L}}T|_{\mathcal{L}}$$

where $|_{\mathcal{L}}$ denotes domain restriction to \mathcal{L} .

- If $\{\phi_n\}_{n \in \mathbb{N}}$ is an orthonormal basis for H , the n 'th **truncation** of T is the compression of T to $\text{Span}\{\phi_1, \phi_2, \dots, \phi_n\}$.
- We will call the matrix representation of T truncated to $\text{Span}\{\phi_1, \phi_2, \dots, \phi_n\}$ the **Ritz matrix** of T , and denote it T_n when the context is obvious. T_n is an $n \times n$ matrix with entries

$$(T_n)_{i,j} := (T\phi_i, \phi_j) \quad \forall i, j \leq n.$$

We have a natural yet unanswered question; do the eigenvalues of the matrices T_n converge to the spectrum of the operator T , as n increases? This question was investigated by Walther H. W. Ritz for whom we name our matrices, as well as by Boris Galerkin; the method of approximating the spectrum of an operator by the eigenvalues of its truncations is often called the Ritz method, Galerkin method, or indeed the Ritz-Galerkin method. We will formulate answers to this question in due course, but for our introduction it suffices to say that the answer is "not quite". What can go wrong?

1.3 Spectral pollution

Definition. (Spectral pollution) (Adapted from [7]) Let $(T_n)_{n \in \mathbb{N}}$ be an increasing sequence of truncations of an operator T . A value $\lambda \in \mathbb{C}$ is said to be a point of **spectral pollution** if there is a sequence $\lambda_n \in \text{Spec}(T_n)$ such that $\lambda_n \rightarrow \lambda$ but $\lambda \notin \text{Spec}(T)$.

Points of spectral pollution are, intuitively, artefacts of the approximation which will never converge to a point in the actual spectrum. We will see that they exist, that they are relatively common, and that they *get worse* as the approximation goes to higher iterations. Unless we already know what the spectrum of the operator is, it can be incredibly hard for us to decide whether a point is actually in the spectrum or whether it is spurious. In applications of spectral theory, this difference can be beyond a simple 'noisy data' nuisance - rather, a confounding problem.

Example 2. Consider a one-dimensional Sturm-Liouville operator

$$Ly = \frac{d}{dx}\left(P(x)\frac{dy}{dx}\right) + Q(x)y$$

where P and Q are scalar-valued functions. Q is called the 'potential' of the operator. These operators occur frequently in partial differential equations; for example, both the heat and wave equations are Sturm-Liouville problems.

In particular, let us look at the following equation. This is the Schrödinger operator governing the movement of an electron in a hydrogen atom:

TODO

The spectrum of this operator, fascinatingly, corresponds to the energy levels of certain ‘stable’ states of the atom in quantum mechanics. The mathematical analysis of this for larger atoms is an open problem; for example, the spectrum for the equivalent equation in the helium atom has only been calculated numerically [7]. The hydrogen atom’s spectrum can be found concretely, and is equal to

Let us see how well a Galerkin method can approximate this known spectrum.

There are, of course, other methods for approximating operator spectra, such as the popular finite difference or ‘shooting’ methods [3], or specialised algebraic methods [8], which are not subject to pollution. So why do we care about truncation methods? The motivation for studying the Ritz-Galerkin method is that it makes *almost no* assumptions about the operator itself or the location of its spectrum. If the spectral pollution for a sequence of truncations can be discarded, detected or otherwise dealt with, it would provide a unified and powerful approach to numerical spectral theory for almost any operator.

II | SPECTRAL POLLUTION IN A MULTIPLICATION OPERATOR

2.1 The spectrum of a multiplication operator

Definition. (Multiplication operator) Let Ω be a σ -finite measure space. For a given function a , the multiplication operator M_a on $L^2(\Omega)$ is defined by the action $M_a f(x) = a(x)f(x)$; that is, it acts by pointwise multiplication with a . We call a the ‘symbol’ of the multiplication operator.

The spectrum of a multiplication operator is easy to calculate. We first need to define the ‘essential range’ of a function. Intuitively, this is similar to the standard range of a function, but ignoring values taken by the function on a set of measure zero - two functions which are equal almost everywhere will have the same essential range.

Definition. (Essential range and essential supremum) The essential range of a real-valued function f is the set:

$$\{k \in \mathbb{R} : \forall \varepsilon > 0, \mu\{x : |f(x) - k| < \varepsilon\} > 0\}$$

where μ is the Lebesgue measure.

The essential supremum of f , denoted $\text{esssup} f$, is the supremum of the essential range of f . We define $\text{essinf} f$ mutatis mutandis for the infimum.

The normed vector space L^∞ is the space of all bounded measurable functions, and its norm is

$$\|f\|_\infty = \text{esssup}(|f|).$$

Lemma 2.1. (Adapted from [1]) A multiplication operator is bounded if and only if its symbol is in L^∞ .

Proof. Let M_a be a multiplication operator on a Hilbert space H with symbol $a \in L^\infty$. We see that $|f(z)| \leq \|a\|_\infty$ for almost all z , so $|ag| = |a||g| \leq \|a\|_\infty |g|$ pointwise almost everywhere; thus

$$\|ag\|_2 = \sqrt{\int |ag|^2} \leq \sqrt{\int \|a\|_\infty^2 |g|^2} = \|a\|_\infty \sqrt{\int |g|^2} = \|a\|_\infty \|g\|_2. \quad (\star)$$

This implies $\|M_a\| = \sup_{g \in H} \|ag\|_2 \leq \|a\|_\infty$, so M_a is bounded.

Conversely, assume M_a is bounded, $a \neq 0$ and let $c < \|a\|_\infty$ (of course, we do not assume $\|a\|_\infty$ is finite). Then we know that the set $\{z : |a(x)| > c\}$ must have positive measure, and by σ -finiteness we have a subset $E \subseteq \{z : |a(x)| > c\}$, and its characteristic function $\mathbf{1}_E$ is in L^2 . Then $|a\mathbf{1}_E| \geq c\mathbf{1}_E$, and by a similar calculation to (\star) and taking the supremum, $\|M_a \mathbf{1}_E\|_2 \geq c\|\mathbf{1}_E\|_2$, so $\|M_a\| \geq c$. Taking the supremum over all c gives $\|M_a\| \geq \|a\|_\infty$, so $\|a\|_\infty$ is finite and so $a \in L^\infty$. \square

Note from this lemma we also get that $\|M_a\| = \|a\|_\infty$.

Theorem 2.2. The spectrum of the multiplication operator M_a is the essential range of its symbol a .

Proof. (Adapted from [9]) Let $(M_a - \lambda)f = g$. Then $f(z) = \frac{g}{a(z) - \lambda}$, and we see that $(M_a - \lambda)$ is invertible if and only if the operator $M_{(a-\lambda)^{-1}}$ is bounded, which is the case if and only if $(a - \lambda)^{-1} \in L^\infty$ by Lemma 2.1. This is the case exactly when $(a - \lambda) \geq \epsilon$ almost everywhere - by definition, this means that $(M_a - \lambda)$ is invertible if and only if λ is not in the essential range of a . \square

We will now take advantage of the ability to create simply-defined operators with easy-to-calculate spectra to observe the existence of spectral pollution.

2.1.1 Estimating spectra computationally

Example 3. Let M_f be the multiplication operator on $L^2(0,1)$ with symbol

$$f : x \mapsto \begin{cases} x & x < 1/2 \\ x + 1/2 & \text{otherwise.} \end{cases}$$

By Theorem 2.2, the spectrum of M_f is the set $[0, 1/2] \cup [1, 3/2]$. If our Ritz approximation works as expected, we should see the approximation create two dense clusters of eigenvalues which approach these intervals.

Figure 1 shows a plot of the approximate spectrum for various Ritz matrix sizes. This approximation was done for the sequence of truncations T_{25n} for $n \in \{2, 3, \dots, 20\}$ over the orthonormal basis $\phi_n(x) = \exp(2i\pi nx)$ for $n \in \mathbb{Z}$; each truncation was over the space $\text{span}\{\exp(2i\pi kx), k \in \mathbb{Z}, |k| < n/2\}$.

We do successfully and rather quickly get eigenvalues corresponding to our actual spectrum, but we also get a lot of other eigenvalues; some of them are converging to points in the spectrum, but as the approximation improves, the pollution doesn't fully dissipate.

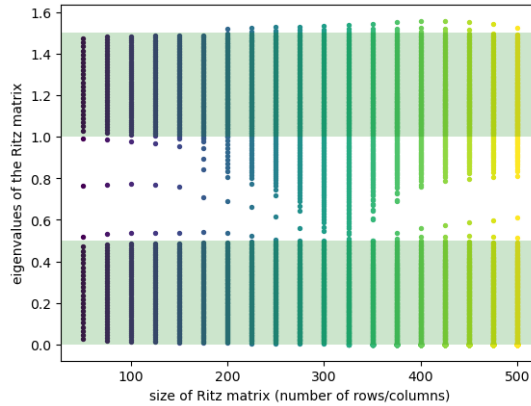


Figure 1: The approximate spectrum of the multiplication operator M_f for Ritz matrices of increasing size. The shaded green areas correspond to the actual spectrum of the operator.

Some of these extra values seem to eventually converge into the correct part of the spectrum, but others stay around; in particular, for all of these approximations we have an eigenvalue of roughly 0.95 which does not exist in the operator's spectrum.

Of course, this is not particularly rigorous - we do not yet know anything about the asymptotic behaviour of the eigenvalues of these approximations (it may well converge far beyond our biggest estimate here of a 500×500 matrix) - but this example is certainly motivating. If one had no knowledge of the actual spectrum of M_f , they would be forgiven for considering this to be strong empirical evidence that the spectrum contains the range $[0.8, 1.0]$ or at least some subset of it.

We have also raised a variety of other questions about the nature of this pollution - why are we only getting it in the gap between the intervals, and not far outside? Is it particular to our choice of sequence for our orthogonal projections (in particular, is there some choice which avoids pollution entirely)?

We will first discuss the nature of spectral pollution for a multiplication operator, taking advantage of an underlying structure to its Ritz matrices which make it possible to concretely identify the existence and location of spectral pollution. Following chapters will then devise technology that allows us to answer these questions in greater generality.

2.1.2 The structure of approximations

If we take a closer look at the structure of our approximating Ritz matrices we uncover deeper structure, which will lead to our next topic.

Example 4. Let M_f be a multiplication operator on $L^2(0,1)$. Now we create the Ritz matrix, $A_{j,k} = (M_f \phi_j, \phi_k)$, choosing the orthonormal basis $\phi_n(x) = \exp(2\pi i n x)$.

We now note that there is a structure to our matrix:

$$\begin{aligned} (M_f \phi_j, \phi_k) &= \int_0^1 f(x) \exp(2\pi i j x) \exp(-2\pi i k x) dx \\ &= \int_0^1 f(x) \exp(2\pi i (j - k)x) dx \\ &= c_{j-k} \end{aligned}$$

where c_n is the n 'th Fourier coefficient. Thus our Ritz matrix depends only on the value of $j-k$; in particular, it is constant along each diagonal. This is a special type of matrix known as a Toeplitz matrix.

As a result, the approximation of our operator M_f is equal to the approximation of a infinite matrix $(T_f)_{j,k} = c_{j-k}$, $j, k \in \mathbb{N}_0$ by its truncations $(T_{f,n})_{j,k} = c_{j-k}$ for $j, k \leq n$. And nowhere in this derivation did we use particular properties of f ; we can repeat this reasoning with any function capable of being represented by a Fourier series. Let us systematise what we have seen.

2.2 Toeplitz operators

2.2.1 Toeplitz operators, Toeplitz matrices, and their equivalence

The approximation of spectra provides a natural gateway to the study of Toeplitz operators, a type of operator with an elegant cluster of representations that will provide intuitive and concrete insight into spectral pollution. A full account of Toeplitz theory is the subject of whole monographs (such as [10]) and is an entire subfield in itself; here we will stick to exploring their spectra, and how these relate to the spectra of their multiplication operator neighbours.

Definition. (Toeplitz matrix) A matrix A (finite or infinite) is Toeplitz if it is constant along its diagonals; that is, $A_{i,j} = A_{i+1,j+1}$ for any i, j (where $i, j < N \in \mathbb{N}$ if A is finite, or $i, j \in \mathbb{Z}_+$ when A is infinite).

Note that an infinite Toeplitz matrix induces an operator on $\ell^2(\mathbb{Z}_+)$. We see from our example that if $f \in L^\infty$ is represented by the Fourier series

$$f(z) = \sum_{k=-\infty}^{\infty} c_k e^{ik\theta},$$

then it induces a Toeplitz matrix T_f with $(T_f)_{j,k} = c_{j-k}$. A natural question is to then ask whether a Toeplitz matrix induces a function, and the answer is affirmative. Firstly, we must define a relevant setting for our matrices; indeed, an important part of our example was that $M_f \exp(2\pi i j x)$ was well-defined. The following definitions are adapted from [1].

Definition. (Hardy space) Let ζ be the monomial function on $L^2(\mathbb{T})$, $\zeta(z) = z$, where \mathbb{T} is the unit circle. Then the Hardy space H^2 is the span of all non-negative exponents of ζ ; $\text{span}\{1, \zeta, \zeta^2, \dots\}$.

As we are on the unit circle, z is more familiarly $e^{i\theta}$ for some θ ; then the basis of Hardy space becomes $\{e^{in\theta}\}_{n \in \mathbb{Z}_+}$. Then we can identify any element $f \in H^2$ as any square-integrable function defined on the unit circle with the Fourier series

$$f(e^{i\theta}) \sim \sum_{k=0}^{\infty} c_k e^{ik\theta},$$

that is, with all negative Fourier coefficients equal to zero. This can be identified with the operator on $\ell^2(\mathbb{Z}_+)$ via the isomorphism $\sum_{k=0}^{\infty} c_k e^{ik\theta} \mapsto (c_k)_{k \in \mathbb{Z}_+}$ [10].

A generalised definition can be made for H^p via any $L^p(\mathbb{T})$; we will almost entirely use H^2 (with the exception of needing H^1 later on), which is often defined as ‘the’ Hardy space.

Definition. (Toeplitz operator) Let $\phi \in L^\infty(\mathbb{T})$ be bounded and measurable on the unit circle. The Toeplitz operator T_ϕ is the compression of M_ϕ to the Hardy space: $T_\phi = P_{H^2} M_\phi|_{H^2}$. We call ϕ the ‘symbol’ of T_ϕ .

One may notice what appears like a clash of notation between the induced Toeplitz matrix that was just discussed with the Toeplitz operator. This is not so; there is an elegant relation between Toeplitz operators and Toeplitz matrices.

Theorem 2.3. Let A be a bounded operator on H^2 such that $(A\zeta^j, \zeta^k) = a_{j-k}$ for some sequence $(a_n)_{n \in \mathbb{Z}}$. Then there is some function $\phi \in L^\infty$ such that $A = T_\phi$ and a_n are the Fourier coefficients of ϕ .

Proof. □

Many properties of Toeplitz operators are hard to see via infinite matrices. Being able to represent them as both Fourier series and equally as the compressions of multiplication operators puts us on the firmer ground of functional analysis, rather than asymptotic linear algebra. From this, we are now in the position to exactly calculate the spectrum of a Toeplitz operator with real-valued symbol.

2.2.2 The spectrum of a Toeplitz operator with real symbol

To begin, we require a pair of properties regarding functions in $L^1(\mathbb{T})$ ’s Hardy space, H^1 .

Lemma 2.4. (Properties of H^1 functions) Let H^1 be the space of all functions $f \in L^1(\mathbb{T})$ where f has the Fourier series

$$f(e^{i\theta}) \sim \sum_{n=0}^{\infty} a_n e^{in\theta}$$

i.e. has no negative Fourier coefficients. Then the following properties hold:

1. If $f, g \in H^2$, then $fg \in H^1$;
2. If $f \in H^1$ is real-valued, then f is constant.

Proof. (1.) Let f, g be in H^2 . Then in particular, $f, g \in L^2$, and so their product fg is in L^1 (this can be seen directly by the Hölder inequality). Now if f has Fourier coefficients a_k and g has Fourier coefficients b_k , the k ’th Fourier coefficient of fg is given by $\sum_{n \in \mathbb{Z}} a_n b_{k-n}$. For negative k , we now have that either $n < 0$ so $a_n = 0$, or $n \geq 0$ so $b_{k-n} = 0$ as $k-n < 0$; this means that the Fourier series of fg satisfies the Hardy space property, but we must justify that the Fourier series of a product is equal to the product of the Fourier series - indeed, they converge to f in L^2 as

(2.) For any real-valued function, we have the relation $\overline{c_{-n}} = c_n$ for the Fourier coefficients of f :

$$\overline{c_{-n}} = \overline{\int f(x) e^{-inx} dx} = \int \overline{f(x)} e^{inx} dx = \int f(x) e^{inx} dx = c_n.$$

But if f is in Hardy space, $c_{-n} = 0$ for all $n \in \mathbb{N}$, and so $c_n = \overline{c_{-n}} = 0$ for all $n \in \mathbb{N}$. Thus the only coefficient remaining is c_0 , and a function with a constant Fourier series is a constant function. (This proof is valid for any $H^p, p \in [1, \infty)$.) □

We are now in the position to calculate the form of the spectrum for any Toeplitz operator. This spectrum was first found by Toeplitz himself under a stronger set of regularity assumptions, and then weakened by Wiener via his Tauberian theorem [11]. Our proof will consist of a series of claims about the resolvent set of T_ϕ , following a proof outlined in an exercise of Arveson ([1], Chapter 4.6, exercises 2-5) which is based on a proof by Hartman and Wintner.

Theorem 2.5. (Hartman-Wintner) Let $\phi \in L^\infty$ be real-valued. Then $\text{Spec}(T_\phi) = [m, M]$, where m and M are the infimum and supremum of the essential range of ϕ (Definition 2.1) respectively.

Proof. Note we already know that as ϕ is real-valued, T_ϕ is self-adjoint, and so its spectrum is on the real line. Furthermore, we will assume that ϕ is non-constant, as if ϕ is constant then T_ϕ is some constant multiple of the identity operator and its spectrum is simply the set containing that constant, and the result holds.

- I. Let $\lambda \in \mathbb{R}$ be such that $T_\phi - \lambda$ is invertible. Then there is a non-zero function $f \in H^2$ such that $(T_\phi - \lambda)f(z) = k$ for any z . By the definition of invertibility this is true for $f = (T_\phi - \lambda)^{-1}\kappa$, where κ is the constant function in H^2 ; $\kappa : z \mapsto k$ for some $k \in \mathbb{C}$.
- II. Now we claim that $(\phi - \lambda)|f|^2$ is constant almost everywhere. To do this, we first show that it is in H^1 . Indeed, we have $(\phi - \lambda)|f|^2 = ((\phi - \lambda)\bar{f})f$. Then by the previous part,

$$(\phi - \lambda)\bar{f} = \overline{(\phi - \lambda)f} = \bar{\kappa}$$

where $\bar{\kappa}$ maps z to \bar{k} . Then it is still a constant function so is still in H^2 , and by Lemma 2.4.1, $((\phi - \lambda)\bar{f})f$ is in H^1 as a product of two H^2 functions. Now we use Lemma 2.4.2: because ϕ and λ are real-valued, $(\phi - \lambda)|f|^2$ is also real-valued, so must be constant. Let this constant be called c .

- III. Our next claim is that $\phi - \lambda$ crosses the x -axis almost nowhere. This is almost immediate once we invoke a theorem of F. and M. Riesz¹, which states

Let f be a non-zero function in H^2 . Then the set $\{z \in \mathbb{T} : f(z) = 0\}$ has Lebesgue measure zero.

This means that $(\phi(z) - \lambda) = \frac{c}{|f(z)|^2}$ is well-defined on L^∞ . Then because $|f|^2 > 0$ a.e., $(\phi - \lambda)$ is positive almost everywhere if $c > 0$, and negative almost everywhere if $c < 0$. Note that if $c = 0$, then $(\phi(z) - \lambda) = 0$, so ϕ is the constant function with value λ and the result holds by our remark at the beginning of this proof.

- IV. Finally, we show $\text{Spec}(T_\phi) = [m, M]$. By our previous claim, if $T_\phi - \lambda$ is invertible, then either:

- $\phi(z) - \lambda < 0$ a.e., so $\phi(z) < \lambda$ a.e., so $\lambda > M$, or
- $\phi(z) - \lambda > 0$ a.e., so by the same argument $\lambda < m$.

Thus $T_\phi - \lambda$ is invertible only outside of the set $[m, M]$, as required. □

2.2.3 Multiplication operators, Toeplitz operators, and spectral pollution

Let us now bring the discussion back to spectral pollution, and to our discovery in Example 4. The Ritz matrices corresponding to the multiplication operator M_f is *identical* to the Ritz matrices corresponding to the Toeplitz operator T_f (which are just truncations of the corresponding Toeplitz matrix), but these operators have different spectra. No wonder we are seeing extra eigenvalues in the gap; the approximation ‘cannot tell’ the difference between an operator which has the essential range as its spectrum from one which has the same maximum and minimum but with all gaps filled in. Heuristically, one may even expect that with a large enough approximation, the *entire gap* could fill up with spurious eigenvalues².

More rigorously, it is possible to show that for any point in a gap of the multiplication operator’s spectrum, some subsequence of truncations will have an approximate eigenvalue converging to that point.

Theorem 2.6. (Schmidt-Spitzer [11]) *Let T_n be the $n \times n$ matrix created by taking the first n rows and columns of a Toeplitz operator T_f . Consider the set*

$$B = \{\lambda \in \mathbb{C} : \lambda = \lim_{m \rightarrow \infty} \lambda_{i_m}, \lambda_{i_m} \in \text{Spec}(T_{i_m}), i_m \rightarrow \infty\}.$$

where i_m is a subsequence of \mathbb{N} . Then if f is real-valued, $B = \text{Spec}(T_f)$.

¹Which can be proven as a corollary of a famous theorem of Buerling; see [1], chapter 4.5.

²We also have the inverse question; if we are approximating the Toeplitz operator T_f , why does it take so much longer to approximate the spectrum in the interval $\text{Spec}(T_f) \setminus \text{Spec}(M_f)$?

Proof. If f is real-valued, then the corresponding Toeplitz matrix is Hermitian; its Fourier coefficients have the property $\overline{c_{-n}} = c_n$ (see the proof of Lemma 2.4.2), so if A is the corresponding Toeplitz matrix, $A_{i,j} = c_{i-j} = \overline{c_{j-i}} = \overline{A_{j,i}}$. Label the eigenvalues of T_n as $\lambda_{n_1}, \dots, \lambda_{n_{n+1}}$.

We now invoke a theorem of Szegő [12] regarding the distribution of eigenvalues: on any real interval $[a, b]$,

$$\lim_{n \rightarrow \infty} \frac{1}{n+1} \sum_{i=1}^{n+1} \mathbf{1}_{[a,b]}(\lambda_{n_i}) = \frac{1}{2\pi} \mu\{\theta : a \leq f(e^{i\theta}) \leq b\} \quad (1)$$

where μ is the Lebesgue measure, and $\mathbf{1}$ the characteristic function of $[a, b]$ (i.e. $\mathbf{1}_{[a,b]}(\lambda_{n_i})$ is 1 if $\lambda_{n_i} \in [a, b]$ and zero otherwise). Note that this makes sense, as T_n is Hermitian and therefore every eigenvalue is real-valued. Note also that by Theorem 2.5 and the definition of the unit circle \mathbb{T} , the spectrum of T_f is $\{f(e^{i\theta}) : \theta \in [0, 2\pi]\}$. The claim of equation (1) is therefore “as the size of the Toeplitz matrix increases, the proportion of eigenvalues of T_n in $[a, b]$ converges to the proportion of the spectrum in $[a, b]$ ”. Of course, this proportion will never be 0 for any interval $[a, b] \subseteq [\text{essinf} f, \text{esssup} f]$ of positive measure, so we can guarantee that a sequence of truncations has an eigenvalue which converges in $[a, b]$. The result follows. \square

This is a striking result - for any point in the spectral gap $\text{Spec}(T_f) \setminus \text{Spec}(M_f)$ (or indeed in $\text{Spec}(M_f)$, but they aren't really polluting anything), we can choose a sequence of truncations such that there is guaranteed to be spectral pollution at that point!

III | BOUNDING SPECTRAL POLLUTION

3.1 The essential spectrum of an operator

The essential spectrum has several definitions, the most popular usually denoted $\text{Spec}_{e,i}$ for $i \in \{1, 2, 3, 4, 5\}$ in order of size. For most well-behaved operators the definitions are equivalent. This particular definition is known as Weyl's criterion, $\text{Spec}_{e,2}$. [2]

Definition. (Essential spectrum) *The essential spectrum of an operator T on a Hilbert space H is defined as the set of all λ such that a **Weyl sequence** u_n exists for T and λ , i.e. a sequence with the properties:*

- $\|u_n\| = 1 \quad \forall n \in \mathbb{N}$;
- $u_n \rightharpoonup 0$ (where \rightharpoonup denotes weak convergence: $u_n \rightharpoonup u \Leftrightarrow (u_n, g) \rightarrow (u, g) \quad \forall g \in H$);
- $\lim_{n \rightarrow \infty} \|(T - \lambda)u_n\| \rightarrow 0$.

Proposition 3.1. *If λ is in the essential spectrum of T , then it is in the spectrum of T .*

Proof. Let λ satisfy the criterion in the definition of essential spectrum. Now assume for contradiction that the resolvent $(T - \lambda)^{-1}$ exists. Then:

$$\begin{aligned} 0 \leq \lim_{n \rightarrow \infty} \|u_n\| &= \lim_{n \rightarrow \infty} \|(T - \lambda)^{-1}(T - \lambda)u_n\| \\ &\leq \|(T - \lambda)^{-1}\| \lim_{n \rightarrow \infty} \|(T - \lambda)u_n\| \quad (\text{as } (T - \lambda)^{-1} \text{ is bounded}) \\ &= 0 \end{aligned}$$

and so $\|u_n\| \rightarrow 0$. But $\|u_n\| = 1$ for every n , so it cannot converge to zero! Thus this bounded inverse does not exist. \square

Corollary. *Note that the previous proof did not require the weak convergence of u_n ; indeed, if there is a sequence with $\|u_n\| = 1$ and $\lim_{n \rightarrow \infty} \|(T - \lambda)u_n\| \rightarrow 0$, then λ is in the overall spectrum of T (but not necessarily the essential spectrum!)*

We can loosen the definition of weak convergence to just require convergence in a dense subspace of H :

Lemma 3.2. *A bounded sequence u_n , $\|u_n\| \leq C$, is weakly convergent to u in H if and only if it is weakly convergent to u in L where L is a dense subspace of H .*

Proof. Weak convergence in H implying the same in L is obvious by the definition. Conversely, take $g \in H$. For any $\varepsilon > 0$, we have $\|g - \varphi\| < \varepsilon$ for $\varphi \in L$; furthermore by the weak convergence of u_n in L we have $N \in \mathbb{N}$ such that $(u_n - u, \varphi) < \varepsilon$ for $n \geq N$. Then:

$$(u_n - u, g) = (u_n - u, g - \varphi + \varphi) = (u_n - u, g - \varphi) + (u_n - u, \varphi) < \|u_n - u\| \|g - \varphi\| + \varepsilon < \varepsilon(C + 1) \rightarrow 0.$$

\square

We will now construct a Weyl sequence for the essential spectrum of the 'Laplacian' or 'free Schrödinger' operator $T = -\Delta$ on $L^2(\mathbb{R})$, where on \mathbb{R}^1 , Δ is the operator $\Delta f(x) = \frac{d^2}{dx^2} f(x)$

Proposition 3.3. *The essential spectrum of the operator $T = -\Delta$ is the closed half-axis $[0, +\infty)$.*

Proof. First, note that for the exponential function we have

$$T(\exp)(i\omega x) = \omega^2 \exp(i\omega x). \quad (2)$$

This gives much of the intuition for this proof; the function $\exp_{i\omega} : x \mapsto \exp(i\omega x)$ is not an eigenvector as it is not in $L^2(\mathbb{R})$, but it satisfies the eigenvalue equation for T and so any number $\lambda = \omega^2$ - and thus any $\lambda \in [0, +\infty)$ - is 'almost' an eigenvalue for T .

We take advantage of this by choosing some smooth bump function $\rho \in C_c^\infty(\mathbb{R})$ with $\|\rho\|_2 = 1$. We then define $\rho_n = \frac{1}{\sqrt{n}}\rho(x/n)$. ρ_n has some nice properties: by a substitution of variables and direct calculation we have $\|\rho_n\|_2 = \|\rho\|_2$, and furthermore any k 'th derivative $\rho_n^{(k)}$ of ρ_n converges to 0 in L^2 . Indeed:

$$\|\rho_n^{(k)}\|_2 = \frac{1}{n^k} \left\| \frac{1}{\sqrt{n}} \rho^{(k)}(x/n) \right\|_2 = \frac{\|\rho^{(k)}\|_2}{n^k} \rightarrow 0 \quad (3)$$

where one can see $\|\frac{1}{\sqrt{n}}\rho^{(k)}(x/n)\|_2 = \|\rho^{(k)}\|_2$ by the same calculation as $\|\rho_n\|_2 = \|\rho\|_2$.

Now, let our candidate Weyl sequence be $u_n : x \mapsto \rho_n(x) \exp(i\omega x)$, which truncates $\exp(i\omega x)$ to $\text{supp}\rho_n$; this means u_n is in $L^2(\mathbb{R})$. $\|u_n\| = \|\rho_n\|_2 = \|\rho\|_2 = 1$ by direct calculation, and $u_n \rightharpoonup 0$: we can bound u_n by $\frac{1}{\sqrt{n}}M\mathbf{1}_{(\text{supp}u_n)}$, where $\mathbf{1}_A$ is the characteristic function of the set A and M is the maximum value of ρ . Then by Lemma 3.2, we can simply show weak convergence for any $\varphi \in C_0^\infty$, which is dense in L^2 :

$$\begin{aligned} (u_n, \varphi) &= \int_{\mathbb{R}} u_n \varphi \\ &\leq \int_{\mathbb{R}} \frac{1}{\sqrt{n}} M \mathbf{1}_{(\text{supp}u_n)} \varphi \\ &\leq \int_{\text{supp}\varphi} \frac{1}{\sqrt{n}} M \varphi \\ &= \frac{M}{\sqrt{n}} \int_{\text{supp}\varphi} \varphi \rightarrow 0, \end{aligned} \quad \text{as the integral of } \varphi \text{ is finite and independent of } n.$$

Finally, we show that $\lim_{n \rightarrow \infty} \|(T - \lambda)u_n\|_2 \rightarrow 0$ for $\lambda = \omega^2$:

$$\begin{aligned} \|(T - \lambda)u_n\|_2 &= \|(T(\exp_{i\omega} \rho_n) - \omega^2(\exp_{i\omega} \rho_n))\|_2 \\ &= \|(T(\exp_{i\omega} \rho_n) - T(\exp_{i\omega})\rho_n)\|_2 && \text{(by equation (2))} \\ &= \|\exp_{i\omega} T\rho_n - 2\omega \exp_{i\omega} \frac{d}{dx} \rho_n\|_2 && \text{(by the product rule)} \\ &= \|T\rho_n - 2\omega \frac{d}{dx} \rho_n\|_2 && \text{(see } \|\exp_{i\omega} \phi\|_2 = \|\phi\|_2 \text{ for any } \phi \in L^2) \\ &\leq \left\| -\frac{d^2}{dx^2} \rho_n \right\|_2 + 2\omega \left\| \frac{d}{dx} \rho_n \right\|_2 \rightarrow 0, \end{aligned}$$

converging by equation (3). Thus u_n forms a Weyl sequence for T and $\lambda \in [0, +\infty)$, as required. \square

We can use a similar idea for another example to find the essential spectrum of the multiplication operator:

Proposition 3.4. *The essential spectrum of the operator M_f on $L^2(0, 1)$, where $M_f u(x) = f(x)u(x)$, is the range of f .*

Proof. Similar to before, our initial idea comes from an 'almost-eigenvector'. In this case, if λ is in the range of f with $f(x_0) = \lambda$, we see that $M_f \delta_{x_0} = \lambda \delta_{x_0}$, where δ_{x_0} is the Dirac delta centred at x_0 . Again, δ_{x_0} is not an eigenfunction of M_f as it is not in the correct domain - this time, it isn't even strictly a function (it is a distribution).

Now consider a Friedrichs mollifier ρ . This is a function in $C_0^\infty(\mathbb{R})$ with the property that $\sqrt{n}\rho(ny) \rightarrow \delta_0$ as $n \rightarrow \infty$; we renormalise it such that $\|\rho\|_2 = 1$, and take the sequence

$$u_n : x \mapsto \begin{cases} \sqrt{n}\rho(n(x - x_0)) & x \in (0, 1) \\ 0 & \text{otherwise} \end{cases}$$

thus u_n "converges to δ_{x_0} " in the sense of distributions. Note that $\|u_n\|_2 = \|\rho\|_2 = 1$ for all n , and this

sequence converges weakly to 0:

$$\begin{aligned}
|(u_n, g)| &= \int_{\text{supp} u_n} \sqrt{n} \rho(n(x - x_0)) g(x) && (\text{for any } g \in L^2(0, 1)) \\
&\leq \|u_n\|_2 \sqrt{\int_{\text{supp} u_n} |g(x)|^2} && (\text{by Hölder's inequality}) \\
&= \sqrt{\int_{\text{supp} u_n} |g(x)|^2} \rightarrow 0, \quad \text{as } \text{supp}(u_n) \text{ decreases to } 0.
\end{aligned}$$

Then we see $\|(M_f - \lambda)u_n\|_2$ converges to zero by similar reasoning:

$$\begin{aligned}
\|(M_f - \lambda)u_n\|_2^2 &= \int_{\text{supp} \rho_n} |(f(x) - f(x_0))\sqrt{n} \rho(n(x - x_0))|^2 && (\text{using that } \lambda = f(x_0)) \\
&= \|(f(x) - f(x_0))^2\|_{L^\infty(\text{supp} \rho_n)} \|\rho_n^2\|_1 && (\text{by Hölder's inequality}) \\
&= \sup_{x \in \text{supp} \rho_n} \|(f(x) - f(x_0))^2\| \rightarrow 0 && (\text{note } \|\rho_n^2\|_{L^1} = \|\rho_n\|_2 = 1)
\end{aligned}$$

converging to zero as $\text{supp} \rho_n$ shrinks around x_0 by the continuity of f . \square

Compare this example to Theorem 2.2, and see that the *entire spectrum* of a multiplication operator is essential spectrum; the operator has no eigenvalues.

3.2 Essential spectrum and compact perturbation

One interesting property of the essential spectrum that is not easily visible from our earlier definition is that it is invariant under compact perturbations. This is not the case for eigenvalues!

Definition. (Compact operators and rank) An operator T on a normed vector space X is compact if for every bounded sequence $(x_n)_{n \in \mathbb{N}}$ in X , the sequence $(Tx_n)_{n \in \mathbb{N}}$ has a convergent subsequence.

The rank of an operator T , denoted $\text{rank} T$, is the dimension of its range.

Note in particular that if a bounded operator T has finite rank, then T is compact; as its image is finite-dimensional and bounded, the Bolzano-Weierstrass theorem holds for $(Tx_n)_{n \in \mathbb{N}}$.

Note also that any compact operator is necessarily bounded, as otherwise we could choose a bounded sequence $(x_n)_{n \in \mathbb{N}}$ in H such that $\|Tx_n\| \rightarrow \infty$, and then it would not be possible for $(Tx_n)_{n \rightarrow \infty}$ to have a bounded subsequence.

Theorem 3.5. Let λ be in the essential spectrum of an operator T on a Hilbert space H . Then

$$\lambda \in \bigcap_{K \in \mathcal{K}(H)} \text{Spec}_e(T + K) \quad (4)$$

where $\mathcal{K}(H)$ is the space of all compact linear operators on H .

Proof. First, let $\lambda \in \text{Spec}_e(T)$ have the Weyl sequence x_n for the operator T . Then

$$\|(T + K - \lambda)x_n\| = \|(T - \lambda)x_n + Kx_n\| \leq \|(T - \lambda)x_n\| + \|Kx_n\|$$

And as K is compact, Kx_n has a convergent subsequence Kx_{n_k} , where $x_{n_k} \rightharpoonup 0$ because $x_n \rightharpoonup 0$. Then Kx_{n_k} also weakly converges to 0 by two applications of the Riesz representation theorem and the boundedness of K : for any $\phi \in H$,

$$\begin{aligned}
(Kx_{n_k}, \phi) &= f(Kx_{n_k}) && \text{for some bounded linear functional } f \in H^* \\
&= (f \circ K)(x_{n_k}) && \text{which is also a bounded linear functional in } H^* \\
&= (x_{n_k}, \psi) && \text{for some } \psi \in H \\
&\rightarrow 0 && \text{by weak convergence of } x_{n_k} \text{ to } 0
\end{aligned}$$

and as it is (strongly) convergent to some value K , it must also be weakly convergent to the same value; thus $Kx_{n_k} \rightarrow 0$.

Thus $\|(T + K - \lambda)x_{n_k}\| \leq \|(T - \lambda)x_{n_k}\| + \|Kx_{n_k}\| \rightarrow 0$, so x_{n_k} is a Weyl sequence for λ and $T + K$ for any compact operator K , so $\lambda \in \bigcap_{K \in \mathcal{K}(H)} \text{Spec}_e(T + K)$. \square

Remark. As mentioned before, eigenvalues do not have this property: let λ be an eigenvalue of T with eigenvector u , and P the orthogonal projection onto the space $\text{span}\{u\}$ (so $\text{rank } P = 1$). Then:

$$(T + P)u = Tu + u = \lambda u + u = (\lambda + 1)u$$

so λ is not an eigenvalue of $(T + P)$. This will become a very useful property in discussing a method of detecting spectral pollution known as dissipative barrier methods, which we will explore in section 4.1.

3.3 Numerical range and essential numerical range

Definition. (Numerical range of an operator) Let T be an operator on a Hilbert space H . The numerical range $W(T)$ is defined

$$W(T) = \{(Tu, u) : u \in \text{Dom}(T), \|u\| = 1\}$$

where $\text{Dom}(T)$ is the domain of T .

The numerical range has a variety of interesting properties which make them useful for roughly approximating the location of spectra.

Proposition 3.6. The numerical range $W(T)$ of an operator T has the following properties:

1. $W(T) \in \mathbb{R}$ if and only if T is self-adjoint;
2. $W(T_{\mathcal{L}}) \subseteq W(T)$, where $T_{\mathcal{L}}$ is the compression of T to the closed subspace \mathcal{L} ;
3. (Toeplitz-Hausdorff theorem) $W(T)$ is a convex set;
4. $\text{Spec}(T) \subseteq \overline{W(T)}$, where $\overline{W(T)}$ is the closure of the numerical range of T .

It is important to state the usefulness of these properties with regards to spectral pollution. Not only does $W(T)$ bound the spectrum of T , it bounds the spectrum of $T_{\mathcal{L}}$ - effectively, bounding the region in which spectral pollution can occur to a convex set around $\text{Spec}(T)$.

Proof. (1.) T is self-adjoint if and only if $(Tu, u) = (u, Tu)$ for all u ; by conjugate symmetry of scalar products, $(u, Tu) = \overline{(Tu, u)}$; we combine these to find $(Tu, u) = \overline{(Tu, u)}$ and the result follows.

(2.) $W(T_{\mathcal{L}}) = \{(PTPu, u) : u \in \mathcal{L}, \|u\| = 1\}$. We then use the self-adjointness of P to see

$$(PTPu, u) = (TPu, Pu)$$

Then as $u \in \mathcal{L}$, $\|Pu\| = \|u\| = 1$, so $(TPu, Pu) \in W(T)$, and the result follows.

(3. [13]) Take $\lambda = (Tx, x)$, $\mu = (Ty, y) \in W(T)$. Define the line segment between them as $\nu = t\lambda + (1-t)\mu$ for $t \in [0, 1]$, and $T_{\mathcal{L}}$ the compression of T to the subspace $\mathcal{L} = \text{span}\{x, y\}$. Then we note that $(T_{\mathcal{L}}x, x) = (Tx, x)$ and $(T_{\mathcal{L}}y, y) = (Ty, y)$, so λ, μ are in $W(T_{\mathcal{L}})$. $T_{\mathcal{L}}$ is two-dimensional, so is a 2×2 matrix; it can be proven by direct calculation (see [13]) that the numerical range of a 2×2 matrix is an ellipse (with foci at either eigenvalue of the matrix!) and so ν is in $W(T_{\mathcal{L}})$. Then by property 2, $W(T_{\mathcal{L}}) \subseteq W(T)$, so ν is also in $W(T)$, as required.

(4.) $\overline{W(T)}$ is the set of all points λ such that there is a sequence of unit vectors u_n where

$$\lim_{n \rightarrow \infty} (Tu_n, u_n) = \lambda.$$

The approximate point spectrum $\sigma_{ap}(T) = \{\lambda \in \mathbb{C} : \lim_{n \rightarrow \infty} \|(T - \lambda)u_n\| = 0 \text{ for some sequence } (u_n)_{n \in \mathbb{N}}, \|u_n\| = 1\}$ can be shown to contain the boundary of the spectrum of T (e.g. in Halmos [14], problem 78). Then by the Cauchy-Schwarz inequality, $|(T - \lambda)u_n, u_n| \leq \|(T - \lambda)u_n\| \rightarrow 0$, and so

$$\begin{aligned} |((T - \lambda)u_n, u_n)| &= |(Tu_n, u_n) - (\lambda u_n, u_n)| \\ &= |(Tu_n, u_n) - \lambda \|u_n\|^2| \\ &= |(Tu_n, u_n) - \lambda| \rightarrow 0; \\ &\implies (Tu_n, u_n) \rightarrow \lambda. \end{aligned}$$

Thus the boundary of the spectrum is in $\overline{W(T)}$; as the numerical range is convex, this means the whole spectrum must be. □

3.4 Essential numerical range

A similar notion to that of the numerical range is the essential numerical range, $W_e(T)$. This set lowers its aim to simply estimating the essential spectrum, but in the process manages to do so much more accurately for some operators.

Definition. (*Essential numerical range*) (adapted from [15]) *The essential numerical range of an operator T is given by*³

$$W_e(T) := \left\{ \lim_{n \rightarrow \infty} (Tu_n, u_n) : (u_n)_{n \in \mathbb{N}} \text{ in } \text{Dom}(T), \|u_n\| = 1, u_n \rightharpoonup 0. \right\}$$

Note the parallels with our definition of the essential spectrum, $\text{Spec}_e(T)$. Indeed, these parallels are reflected in the properties of $W_e(T)$:

Proposition 3.7. *The essential numerical range $W_e(T)$ of an operator T has the following properties:*

1. $W_e(T)$ is convex;
2. $W_e(T) \subseteq \overline{W(T)}$;
3. $\text{conv}(\text{Spec}_e(T)) \subseteq W_e(T)$, with equality if T is self-adjoint and bounded.

Proof. (1. [16]) We prove this by applying the Toeplitz-Hausdorff theorem (Proposition 3.6.3) to a sequence. We take $\lambda = \lim_{n \rightarrow \infty} (Tx_n, x_n)$, $\mu = \lim_{n \rightarrow \infty} (Ty_n, y_n) \in W_e(T)$ and define a sequence

$$\nu_n = t(Tx_n, x_n) + (1 - t)(Ty_n, y_n) \quad t \in [0, 1],$$

which obviously converges to $\nu = t\lambda + (1 - t)\mu$ as $n \rightarrow \infty$. We then create a sequence of compressions T_n , where each compression is to $\text{span}\{x_n, y_n\}$. By the Toeplitz-Hausdorff theorem, we know that ν_n is in $\text{Spec}(T_n)$ and get a sequence $\nu_n = (Tz_n, z_n)$ converging to ν ; the elements z_n are unit vectors in $\text{span}\{x_n, y_n\}$, and they weakly converge to 0 because x_n and y_n both do:

$$(z_n, g) = (\alpha x_n + \beta y_n, g) = \alpha(x_n, g) + \beta(y_n, g) \rightarrow 0 \quad \forall g \in H.$$

This means that $\nu = \lim_{n \rightarrow \infty} (Tz_n, z_n)$ is in $\text{Spec}_e(T)$ as required.

(2.) This can be seen directly from looking at the two definitions. By definition, we have

$$\overline{W(T)} = \{\lim_{n \rightarrow \infty} (Tu_n, u_n) : (u_n)_{n \in \mathbb{N}} \text{ in } \text{Dom}(T), \|u_n\| = 1\},$$

³Much like the essential spectrum, there are multiple definitions of the essential numerical range. However (at least for bounded operators) there is much more equivalence between the definitions than we have for essential spectrum! [15] We choose the definition with the most natural relation to our choice of definition for essential spectrum.

and $W_e(T)$ is the subset of this with the extra condition that $u_n \rightarrow 0$.

(3.) The inclusion $\text{Spec}_e(T) \subseteq W_e(T)$ comes from an analogous argument to that of Proposition 3.6.4; then $\text{conv}(\text{Spec}_e(T)) \subseteq W_e(T)$ by this inclusion and that $W_e(T)$ is a convex set. \square

It remains to show that $\text{conv}(\text{Spec}_e(T)) = W_e(T)$ when T is self-adjoint and bounded. This will be proven after the following theorem, which describes another similarity with essential spectra; invariance under compact perturbation.

Theorem 3.8. *A value λ is in the essential numerical range of an operator T on a Hilbert space H if and only if*

$$\lambda \in \bigcap_{K \in \mathcal{K}(H)} \overline{W(T + K)}$$

where $\mathcal{K}(H)$ is the set of all compact linear operators on H .

Proof. Let $\lambda = \lim_{n \rightarrow \infty} (Tx_n, x_n)$ be in $W_e(T)$. Then $((T + K)x_n, x_n) = (Tx_n, x_n) + (Kx_n, x_n)$ which converges to λ if $(Kx_n, x_n) \rightarrow 0$, which is true by the weak convergence of x_n (using that K is bounded as it is compact). Thus for any compact operator K , $\lambda \in W_e(T + K) \subseteq \overline{W(T + K)}$.

Conversely, \square

We have seen that the essential numerical range estimates the bounds of the essential spectrum with quite astounding accuracy for some types of operator. But the essential numerical range far outdoes the regular numerical range on bounding spectral pollution; in fact, it provides an *exact* set on which it is possible for pollution to occur!

Theorem 3.9. *All spectral pollution in the Ritz approximation of $\text{Spec}(T)$ will be located inside of $W_e(T)$; within this set, it can occur anywhere in $W_e(T) \setminus \text{Spec}(T)$.*

For a bounded operator, this is the main result of a paper by Pokrzywa [17]. The main theorem of the paper has the corollary that for $\lambda \notin W_e(T)$, we have $\lambda \in \text{Spec}(T)$ iff $\text{dist}(\lambda, \text{Spec}(T_n)) \rightarrow 0$; that is, outside of the essential numerical range, every point in the approximate spectrum $\text{Spec}(T_n)$ converges to a point in the actual spectrum of T . This is followed by a lemma which claims that for any sequence $(\lambda_n)_{n \in \mathbb{N}}$ in the interior of $W_e(T)$, there is a sequence of orthogonal projections such that $\lambda_{n-1} \in \text{Spec}(T_n)$ - not only does all spectral pollution occur inside this range, but for *any point* in $W_e(T) \setminus \text{Spec}(T)$, spectral pollution occurs there in some approximation.

IV | DETECTING SPECTRAL POLLUTION

4.1 Dissipative barrier methods

Perhaps one of the simplest methods for separating eigenvalues from spectral pollution is the dissipative barrier method. This method leverages finite-rank perturbation of an operator.

Take an operator T on a Hilbert space, and let P be an orthonormal projection onto a finite-dimensional space such that for an eigenvector u , $\|Pu - u\|$ is sufficiently small. Let λ be the eigenvalue for u . Then for the operator $T + iP$,

$$(T + iP)u = Tu + iP u \approx \lambda u + iu = (\lambda + i)u.$$

i.e. $(\lambda + i)$ is approximately an eigenvalue for $T + iP$.

Then, the eigenvalues of the perturbed operator will have imaginary part of approximately 1, and the set of eigenvalues produced by the approximation can be filtered to discard points without imaginary part close to 1 as being pollution. This is particularly useful for self-adjoint operators, where the entire spectrum and the essential numerical range are a subset of the real line; in that case, the perturbed eigenvalues will be the only points on the spectrum with significant imaginary part. Even better, if the operator is also bounded, then all pollution is in the essential numerical range (Theorem ??) which is invariant under compact perturbation (Theorem ??) so the pollution will converge to the real axis.

Let us see this concretely with an example.

Example 5. We return once again to our discontinuous multiplication operator M_f on $L^2(0, 1)$ with symbol

$$f : x \mapsto \begin{cases} x & x < 1/2 \\ x + 1/2 & \text{otherwise.} \end{cases}$$

This time, we add a ‘rank-one perturbation’ to get the operator \tilde{M} which has the action

$$\tilde{M}u = M_f u + (u, \varphi)\varphi,$$

where φ is a real-valued function on $L^2(0, 1)$. Then \tilde{M} has an eigenvector: we rearrange to get

$$\begin{aligned} f(x)u(x) + (u, \varphi)\varphi(x) &= \lambda u(x), \\ \text{therefore } (\lambda - f(x))u(x) &= (u, \varphi)\varphi(x), \\ u(x) &= c \frac{\varphi(x)}{\lambda - f(x)} \end{aligned}$$

for some constant c . We then normalise this to $\frac{\varphi(x)}{\lambda - f(x)}$, which requires

$$\begin{aligned} f(x)u + (u, \varphi)\varphi &= \lambda u \\ (f(x) - \lambda)c \frac{\varphi(x)}{\lambda - f(x)} + \left(\frac{c\varphi}{\lambda - f}, \varphi\right)\varphi(x) &= 0 \\ -1 + \left(\frac{\varphi}{\lambda - f}, \varphi\right) &= 0 \end{aligned}$$

so we can normalise provided that $\left(\frac{\varphi}{\lambda - f}, \varphi\right) = \int_0^1 \frac{|\varphi|^2}{\lambda - m} = 1$.

Thus for any value λ we can choose $\varphi \in L^2(0, 1)$ and scale it so that $\int_0^1 \frac{|\varphi|^2}{\lambda - m} = 1$ to get an operator \tilde{M} with spectrum $\text{Spec}(M) \cup \{\lambda \in \mathbb{C} : \int_0^1 \frac{|\varphi|^2}{\lambda - m} = 1\}$. In Figure 2 we can see the results of Ritz approximations with φ chosen such that the operator has an eigenvalue at 0.7.⁴ In Figure 3 we see the same approximation but with the dissipative barrier iP where P is the projection onto $\text{span}\{\phi_n, |n| \leq 25\}$. Note that the spectrum on the line $(x, y) : \text{Imag}(x) = 1$ converges to what we’d expect the spectrum to be!

⁴In particular, φ was chosen to be the constant $(\log(3) - 3\log(2) + \log(5) + \log(7) - \log(10))^{-1}$; this satisfies the normalisation condition at $\lambda = 0.7$ but also at $\lambda \approx 4.4$; the eigenvalue at $\lambda \approx 4.4$ has been cropped out of the figure to improve illustration of the idea.

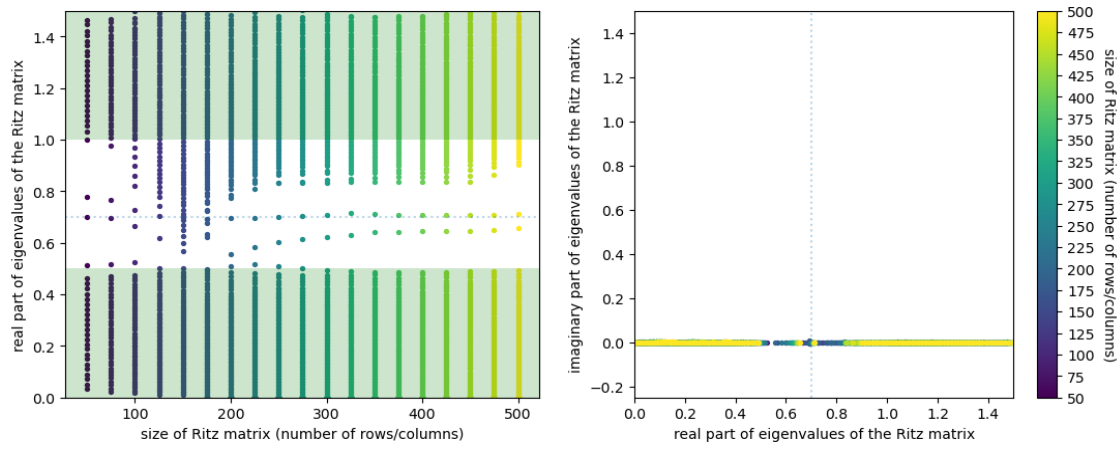


Figure 2: The real part of the approximate spectrum for \tilde{M} ; on the left, the real parts of the approximate spectrum as the size of the Ritz matrix increases; on the right, the complex approximate spectrum where colour is used to denote the size of the approximation. The green shaded regions correspond to the essential spectrum of \tilde{M} , and the dotted lines are at $\text{Re}(x) = 0.7$ to show where the added eigenvalue should be. (An intuitive way to view these figures is to see them as a three-dimensional plot, with the left figure ‘top-down’, and the right figure ‘from the east’)

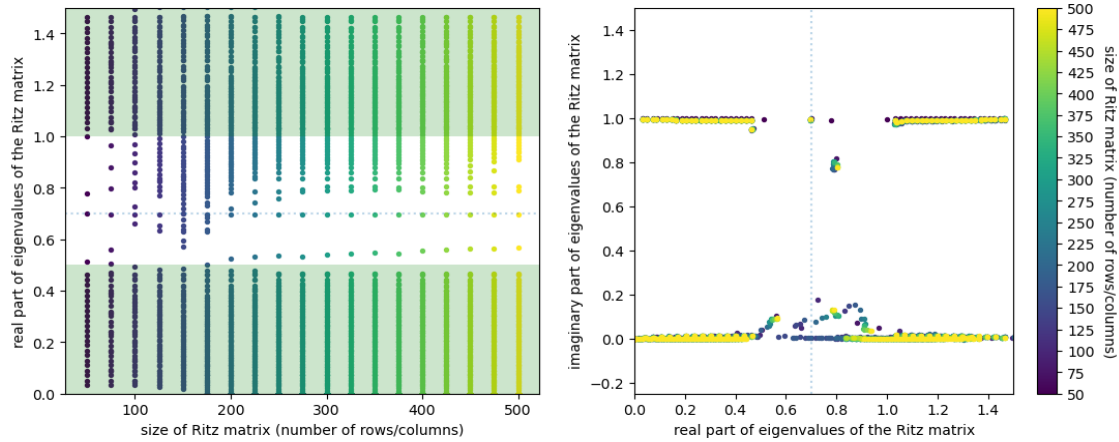


Figure 3: The real part of the approximate spectrum for $\tilde{M} + iP$; compare with Figure 2. See that the line at 1.0 on the imaginary axis converges to the actual spectrum of the operator, while the pollution remains below.

Remark. One may also note that there are bands corresponding to the essential spectrum with imaginary part 1. The reason why dissipative barriers ‘replicate’ the essential spectrum is an open problem; it has recently been investigated specifically for Schrödinger operators [stepanenkoTODO] but in general remains unknown.

As a second example, let us try to replicate some results from Aceto et al. (2006) [8]. In this paper, they use an algebraic method combined with a ‘shooting technique’ to find high-accuracy estimates of eigenvalues for a Sturm-Liouville operator; this highly-specialised algorithm is free from pollution but works only for a specific class of Sturm-Liouville operators. is a

Example 6. In particular, take the following eigenvalue problem on $L^2[0, \infty)$:

$$\begin{cases} -y'' + (\sin(x) - \frac{40}{1+x^2})y = \lambda y \\ y(0) \cos(\pi/8) + y'(0) \cos(\pi/8) = 0. \end{cases}$$

This operator has a ‘band-gap’ structure; it has intervals (bands) of essential spectrum, with eigenvalues dotted in the gaps between bands. In two of the spectral gaps $J_2 = (-0.34767, 0.59480)$ and $J_3 = (0.91806, 1.2932)$ (denoted in line with the paper and rounded to 5sf) the algebraic method finds the following eigenvalues:

J_2	J_3
0.33594	0.94963
0.53662	1.2447
0.58083	1.2919
0.59150	

Firstly, we will see whether we can reproduce this data. The algebraic method used first truncates the half-line to $[0, 70\pi]$. We will do the same and perform a Ritz approximation on the truncated interval, applying a dissipative barrier; this is a success, as one can see in Figure 4.

Now we will aim for better; whether we can reproduce these eigenvalues using a Ritz approximation on $[0, \infty)$, with the orthonormal basis $\{\phi_n\}_{n \in \mathbb{N}}$, $\phi_n = \exp(-x/2)L_n$, where L_n is the n ’th Laguerre polynomial (see [18] for a proof that this is indeed an orthonormal basis)

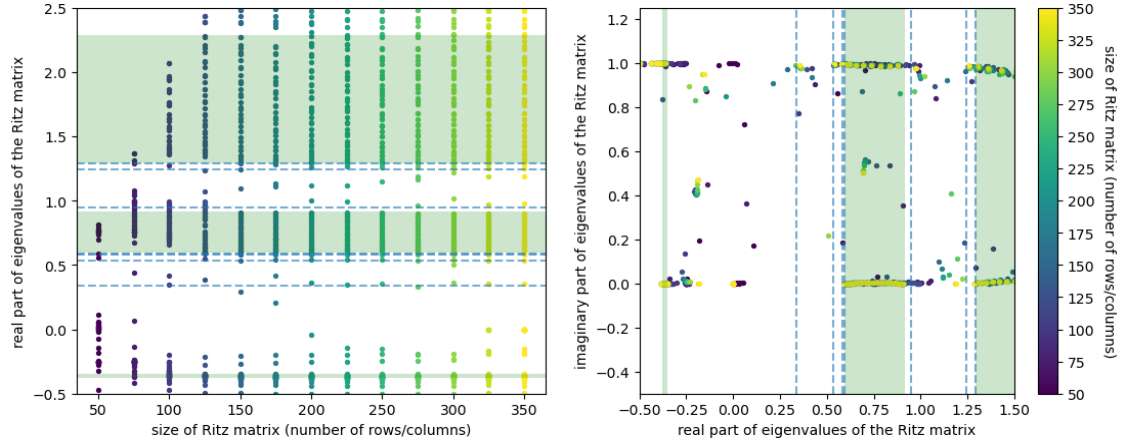


Figure 4: The results of truncating the domain and applying a dissipative barrier to the operator of Example 6. The green bands represent the essential spectrum of the operator, and the dotted blue lines indicate where the algebraic method found eigenvalues in two spectral gaps. Note that for larger approximations, the points in the spectral gaps with imaginary part 1.0 are very close to the algebraic method's approximation.

REFERENCES

- [1] William Arveson. *A Short Course on Spectral Theory*. Vol. 209. Graduate Texts in Mathematics. New York, NY: Springer, 2002. ISBN: 978-1-4419-2943-3 978-0-387-21518-1. DOI: [10.1007/b97227](https://doi.org/10.1007/b97227). (Visited on 11/15/2023).
- [2] David Edmunds and Des Evans. *Spectral Theory and Differential Operators*. Oxford University Press, May 2018. ISBN: 978-0-19-186113-0. DOI: [10.1093/oso/9780198812050.001.0001](https://doi.org/10.1093/oso/9780198812050.001.0001). (Visited on 11/23/2023).
- [3] Endre Süli and David F. Mayers. *An Introduction to Numerical Analysis*. Cambridge University Press, Aug. 2003. ISBN: 978-0-521-00794-8.
- [4] Mathieu Lewin and Eric Séré. “Spectral Pollution and How to Avoid It (With Applications to Dirac and Periodic Schrödinger Operators)”. In: *Proceedings of the London Mathematical Society* 100.3 (May 2010), pp. 864–900. ISSN: 00246115. DOI: [10.1112/plms/pdp046](https://doi.org/10.1112/plms/pdp046). arXiv: [0812.2153 \[math-ph\]](https://arxiv.org/abs/0812.2153). (Visited on 11/29/2023).
- [5] M. Lisa Manning, B. Bamieh, and J. M. Carlson. *Descriptor Approach for Eliminating Spurious Eigenvalues in Hydrodynamic Equations*. Dec. 2008. DOI: [10.48550/arXiv.0705.1542](https://doi.org/10.48550/arXiv.0705.1542). arXiv: [0705.1542 \[physics\]](https://arxiv.org/abs/0705.1542). (Visited on 11/29/2023).
- [6] Eric Cancès, Virginie Ehrlacher, and Yvon Maday. “Periodic Schrödinger Operators with Local Defects and Spectral Pollution”. In: *SIAM Journal on Numerical Analysis* 50.6 (Jan. 2012), pp. 3016–3035. ISSN: 0036-1429. DOI: [10.1137/110855545](https://doi.org/10.1137/110855545). (Visited on 11/29/2023).
- [7] E. Brian Davies. *Spectral Theory and Differential Operators*. Cambridge Studies in Advanced Mathematics. Cambridge: Cambridge University Press, 1995. ISBN: 978-0-521-58710-5. DOI: [10.1017/CB09780511623721](https://doi.org/10.1017/CB09780511623721). (Visited on 11/15/2023).
- [8] L Aceto, P Ghelardoni, and M Marletta. “Numerical Computation of Eigenvalues in Spectral Gaps of Sturm–Liouville Operators”. In: *Journal of Computational and Applied Mathematics* (2006).
- [9] Stephan Ramon Garcia, Javad Mashreghi, and William T. Ross. *Operator Theory by Example*. Oxford University Press, Jan. 2023. ISBN: 978-0-19-195456-6. DOI: [10.1093/oso/9780192863867.001.0001](https://doi.org/10.1093/oso/9780192863867.001.0001). (Visited on 11/16/2023).
- [10] Albrecht Böttcher and Bernd Silbermann. *Analysis of Toeplitz Operators*. Springer Monographs in Mathematics. Berlin, Heidelberg: Springer, 1990. ISBN: 978-3-662-02654-0 978-3-662-02652-6. DOI: [10.1007/978-3-662-02652-6](https://doi.org/10.1007/978-3-662-02652-6). (Visited on 11/15/2023).
- [11] Palle Schmidt and Frank Spitzer. “The Toeplitz Matrices of an Arbitrary Laurent Polynomial”. In: *Mathematica Scandinavica* 8.1 (1960), pp. 15–38. ISSN: 0025-5521. JSTOR: [24489115](https://www.jstor.org/stable/24489115). (Visited on 11/15/2023).
- [12] Ulf Grenander and Gabor Szegő. *Toeplitz Forms and Their Applications*. American Mathematical Society, 2001. ISBN: 978-0-8218-2844-1.
- [13] Karl E. Gustafson and Duggirala K. M. Rao. *Numerical Range*. Ed. by S. Axler, F. W. Gehring, and P. R. Halmos. Universitext. New York, NY: Springer, 1997. ISBN: 978-0-387-94835-5 978-1-4613-8498-4. DOI: [10.1007/978-1-4613-8498-4](https://doi.org/10.1007/978-1-4613-8498-4). (Visited on 11/15/2023).
- [14] Paul R. Halmos. *A Hilbert Space Problem Book*. Vol. 19. Graduate Texts in Mathematics. New York, NY: Springer, 1982. ISBN: 978-1-4684-9332-0 978-1-4684-9330-6. DOI: [10.1007/978-1-4684-9330-6](https://doi.org/10.1007/978-1-4684-9330-6). (Visited on 11/15/2023).
- [15] P. A. Fillmore, J. G. Stampfli, and James P. Williams. “On the Essential Numerical Range, the Essential Spectrum, and a Problem of Halmos”. In: *Acta Sci. Math. (Szeged)* 33.197 (1972), pp. 179–192.
- [16] Sabine Bögli, Marco Marletta, and Christiane Tretter. “The Essential Numerical Range for Unbounded Linear Operators”. In: *Journal of Functional Analysis* 279.1 (July 2020), p. 108509. ISSN: 00221236. DOI: [10.1016/j.jfa.2020.108509](https://doi.org/10.1016/j.jfa.2020.108509). (Visited on 11/15/2023).

- [17] Andrzej Pokrzywa. “Method of Orthogonal Projections and Approximation of the Spectrum of a Bounded Operator”. In: *Studia Mathematica* 65.1 (1979), pp. 21–29. ISSN: 0039-3223. (Visited on 11/15/2023).
- [18] Gábor Szegő. *Orthogonal Polynomials*. American Mathematical Society, 1975. ISBN: 978-0-8218-1023-1.

INDEX

- compression, 3
- Hardy space, 7
- matrix
 - Ritz, 3
 - Toeplitz, 7
- operator
 - compact, 13
 - multiplication, 5
 - Schrödinger, 3
 - Sturm-Liouville, 3
 - Toeplitz, 8
- range
 - essential, 5
 - essential numerical, 15
 - numerical, 14
- rank of an operator, 13
- resolvent, 2
- spectrum, 2
 - essential, 11
 - pollution of, 3
- theorem
 - Hartman-Wintner, 8
 - Schmidt-Spitzer, 9
 - Toeplitz-Hausdorff, 14
- truncation, 3
- Weyl sequence, 11
- Weyl's criterion, 11