

Sobolev Spaces and Finite-Dimensional Approximation

April 15, 2025

Please do exercises 3 and 4 by class-time on Thursday, 2/27. The rest are due the following Tuesday, 3/4, at midnight.

1 Introduction

In this homework, we'll look at a measure of the size of a Sobolev-type model. A model like this.¹

$$\mathcal{M} = \left\{ m(x) = \sum_j m_j \phi_j(x) : \sum_j \lambda_j m_j^2 \leq 1 \right\} \quad (1)$$

where $\langle \phi_j, \phi_k \rangle = \begin{cases} 1 & \text{if } j = k \\ 0 & \text{otherwise} \end{cases}$ and $\lambda_j \geq 0$.

We'll do this in abstract terms, leaving both the inner product $\langle \cdot, \cdot \rangle$ and the corresponding orthonormal basis ϕ_0, ϕ_1, \dots unspecified. So that we have a concrete example to think about, we'll think about what our results mean in the context of the first Sobolev regression model we talk about this semester, for which we use the inner product $\langle \cdot, \cdot \rangle_{L_2}$, the basis functions $\phi_j(x) = \sqrt{2} \cos(j\pi x)$, and $\lambda_j = (\pi j)^2$.

What do we mean by the *size* of a model? When we talk about linear models, we tend to think about *dimension*: the number of basis functions we need to span all the functions in the model. We need two basis functions for lines, four basis functions for cubic polynomials, etc. In this homework, we'll think about a way to generalize that idea so we can get a meaningful measure of the size of infinite-dimensional models. To do that, we'll think about finite-dimensional approximations. We'll think about how high-dimensional a vector space has to be to include an ϵ -approximation to every function m in the model. That is, about the smallest number N for which there exists a basis $b_0 \dots b_{N-1}$ that, for every function $m \in \mathcal{M}$, spans some function $m_\epsilon = m_0 b_0 + \dots + m_{N-1} b_{N-1}$ with $\|m - m_\epsilon\| \leq \epsilon$

Because it's a bit easier, we'll work with the 'inverse' of this measure. By that, I mean that we'll think about the optimal level ϵ of (uniform) approximation accuracy we can get using functions in *any* N -dimensional basis $b_0 \dots b_{N-1}$. This is called the *Kolmogorov N -width* of the model. Here's its definition.²

$$w_N(\mathcal{M}) = \min_{\substack{\text{basis functions} \\ b_0 \dots b_{N-1}}} \max_{m \in \mathcal{M}} \min_{\substack{\text{coefficients} \\ m_1 \dots m_n}} \left\| m - \sum_{j=0}^{N-1} m_j b_j \right\| \quad (2)$$

To make sense of this definition, let's go step by step from the inside out.

1. Start by thinking about some basis $b_0 \dots b_{N-1}$ and some function $m \in \mathcal{M}$ that we want to approximate. Imagine finding m_ϵ , the best approximation to m spanned by that basis. The inner min in (2) is the approximation error $\|m_\epsilon - m\|$.
2. Keeping that basis in mind, imagine finding the function m in our model \mathcal{M} that's hardest to approximate. The one for which our approximation error is largest. The max in (2) is our error in approximating that function.
3. Finally, imagine finding the best basis. The one for which this worst-case approximation error is smallest. The outer min in (2), i.e. the N -width of the model, is the best worst-case approximation error we can get using *any* N -dimensional basis.

In this homework, we'll calculate the Kolmogorov N -width of the abstract Sobolev-type model \mathcal{M} . We'll take it step by step, from the inside out. And we won't really 'minimize' over basis functions. Instead, we'll start out with a convenient guess, calculate the worst-case approximation error when we use that basis, then show that we can't do any better.

1.1 Warm Up

Before we get into the real stuff, let's take a minute to think a bit more about what the definition (2) means and why it is what it is.

1.1.1 Why Uniform Approximation?

Why are we considering the (best) worst-case approximation error for functions in a model \mathcal{M} ? Why not error for one particular function m ? That'd be the Kolmogorov N -width of the set $\{m\}$ containing only that function. And it's not very interesting. For $N \geq 1$, it's zero.

Exercise 1 *Explain why, if model $\mathcal{M} = \{m\}$ contains only one function m , its Kolmogorov N -width $w_N(\mathcal{M})$ is 0 for $N \geq 1$. A sentence or two should be enough.*

Solution 1 *The model $\{m\}$ itself is one-dimensional. We get zero approximation error for $N \geq 1$ using any basis $b_0 \dots b_{N-1}$ with $b_0 = m$.*

This is an instance of a more general phenomenon. For example, the best estimator for a specific signal μ is the one that completely ignores the data and just spits out μ . It's not a useful one. If we want to talk about real estimators, we need to be thinking about at least two signals.

1.1.2 What We're Really Optimizing Over.

What we're optimizing over in (2) isn't really sets of N basis functions. It's N -dimensional subspaces. Prove it!

Exercise 2 *Prove that if $b_0 \dots b_{N-1}$ and $\tilde{b}_0 \dots \tilde{b}_{N-1}$ have the same span \mathcal{V} , then we get same worst-case approximation error for any set \mathcal{M} , i.e., that*

$$\max_{m \in \mathcal{M}} \min_{\substack{\text{coefficients} \\ m_1 \dots m_n}} \left\| m - \sum_{j=0}^{N-1} m_j b_j \right\| = \max_{m \in \mathcal{M}} \min_{\substack{\text{coefficients} \\ \tilde{m}_1 \dots \tilde{m}_n}} \left\| m - \sum_{j=0}^{N-1} \tilde{m}_j \tilde{b}_j \right\|$$

Tip. Recall that the span of a set of vectors $b_0 \dots b_{N-1}$ is the set of all linear combinations of them. $\text{span } b_0 \dots b_{N-1} = \mathcal{V}$ means that every linear combination $\sum_{j=0}^{N-1} a_j b_j$ of these vectors is some vector $v \in \mathcal{V}$ and that every vector $v \in \mathcal{V}$ can be written as a linear combination like that.

Solution 2 *The inner min is over functions spanned by the basis $b_0 \dots b_{N-1}$, so if $\tilde{b}_0 \dots \tilde{b}_{N-1}$ has the same span, the inner min will have the same value.*

One implication is that whenever we're trying to calculate the Kolmogorov N -width of a set \mathcal{M} , we can consider only orthonormal bases $b_0 \dots b_{N-1}$ (i.e. bases with $\langle b_i, b_j \rangle = 1$ if $i = j$ and 0 otherwise). That makes our calculations easier and it doesn't change the minimum we get, since every N -dimensional subspace \mathcal{V} has an orthonormal basis.

1.2 Conventions

1.3 The Norm

In the definition of the Kolmogorov N -width above, the norm $\|\cdot\|$ wasn't specified. It's conventional to use the two-norm $\|\cdot\|_{L_2}$, but the definition (2) makes sense for other norms. In this homework, we'll **use the norm induced by the inner product in our model's description**, i.e. the norm $\|m\| = \sqrt{\langle m, m \rangle}$ where $\langle \cdot, \cdot \rangle$ is the inner product for which $\langle \phi_j, \phi_k \rangle = 1$ if $j = k$ and 0 otherwise in (1). When we're thinking about the Sobolev regression model from class, this aligns with convention.

Keep in mind that the results we'll prove are specific to the case that the norms in (1) and (2) are the same. That's not as limiting as it sounds because if we have two norms that are related, e.g. norms $\|\cdot\|_a$ and $\|\cdot\|_b$ with $\|u\|_a \leq c\|u\|_b$ for some constant c , Kolmogorov N -widths based on these norms will be related too.

1.3.1 The Index Set

In the model description (1), the set of values that j takes on as we count out the terms in our sum wasn't specified. We call this our *index set* and j itself our *index*. People usually use the natural numbers $0, 1, 2, \dots$ when talking about these models in abstract terms. That's what we'll do in this homework. When we want to talk about a finite-dimensional model spanned by K basis functions $\phi_0, \dots, \phi_{K-1}$, we can think of our index set as being the numbers $0, 1, \dots, K-1$ for some K . Or, if we want to stick to the natural numbers, we can imagine we actually have infinitely many basis functions, but $\lambda_j = \infty$ for $j \geq K$ so we "aren't allowed to use" the rest of them.

We'll order our terms so that λ_j is **increasing**, i.e. $\lambda_0 \leq \lambda_1 \leq \dots$. When you think about the Sobolev model from class, this means our first term corresponds to the lowest frequency 'cosine' $\phi_0(x) = \cos(0\pi x) = 1$; our second to the next lowest frequency one $\phi_1(x) = \cos(\pi x)$, etc. Our model includes smaller and smaller multiples of these basis functions. The largest multiple of ϕ_j it includes is $\lambda_j^{-1/2} \phi_j$, which has norm $\|\lambda_j^{-1/2} \phi_j\| = \lambda_j^{-1/2} \|\phi_j\| = \lambda_j^{-1/2}$.

Often, when we're talking about a concrete model, it's convenient to sum over other index sets. For example, when we're talking about regression models with covariates $x \in \mathbb{R}^2$, we might want to think of j as pair (j_1, j_2) and our sum as being over all pairs of natural numbers (j_1, j_2) . This is equivalent to summing over the natural numbers because there are 'as many' natural numbers as there are pairs of natural numbers. Because we can order these pairs and count them out as the 'first', 'second', 'third', etc. pair, we could use natural numbers instead of pairs to 'name' the terms in our sum, but there's no reason to do that. It'd just be a more awkward way of naming the same things in that concrete case.

Don't worry about the sums $\sum_j m_j \phi_j$ converging. When we're using a version of this model where $\lambda_j \rightarrow \infty$ as $j \rightarrow \infty$, it does. It's easy to show that $\|\sum_{j=K}^{\infty} m_j \phi_j\|^2 \rightarrow 0$ as $K \rightarrow \infty$. If you're interested, try it! You may want to do Exercise 4 first.

If it makes you more comfortable, focus on the finite-dimensional case, i.e. act as if every sum you see is over the numbers $0, 1, \dots, K-1$ for some K . That way, there's no need to worry about convergence and you can use matrix/vector notation without it feeling like you're doing something weird.

2 Approximation using the Basis ϕ_0, ϕ_1, \dots

In this section, we'll talk about how well we can approximate functions in using the basis ϕ_0, ϕ_1, \dots . This isn't just a step towards calculating the Kolmogorov N -width of these models. It's something that comes up when we want to actually use these models in practice. For example, when we go to solve a least squares problem like this.

$$\begin{aligned} \hat{\mu} = \operatorname{argmin}_{m \in \mathcal{M}} \sum_{i=1}^n \{Y_i - m(X_i)\}^2 \\ \text{for } \mathcal{M} = \left\{ m(x) = \sum_{j=0}^{\infty} m_j \phi_j(x) : \sum_{j=0}^{\infty} \lambda_j m_j^2 \leq 1 \right\}. \end{aligned} \quad (3)$$

If this were a finite sum, we could just type it up in notation **CVXR** understands and ask for the solution. And infinite sums, we can do the same thing, but using truncated-series approximations to the functions in \mathcal{M} . Like this.

$$\begin{aligned} \hat{\mu}^N = \operatorname{argmin}_{m \in \mathcal{M}^N} \sum_{i=1}^n \{Y_i - m(X_i)\}^2 \\ \text{for } \mathcal{M}^N = \left\{ m(x) = \sum_{j=0}^{N-1} m_j \phi_j(x) : \sum_{j=0}^{N-1} \lambda_j m_j^2 \leq 1 \right\}. \end{aligned} \quad (4)$$

Intuitively, if \mathcal{M}^N contains a good approximation to $\hat{\mu}$, then $\hat{\mu}^N$ will be a good approximation to $\hat{\mu}$, and that'll be the case if \mathcal{M}^N contains a good approximation to every function in \mathcal{M} . It's not quite true that if \mathcal{M}^N contains an ϵ -approximation to every function in \mathcal{M} , then $\hat{\mu}^N$ will be an ϵ -approximation to $\hat{\mu}$, but that's usually a pretty good approximation. We'll act as if it were true for now and worry about the 'pretty good approximation' later in the semester when we have the right tools.

Before we go ahead and write the code in lab, there are a couple things we should work out. First, we should check that using truncated-series approximations like this is actually a reasonable thing to do. Second, we should figure out how many terms we need to include to get a good approximation.

Exercise 3 Let $m(x) = \sum_{j=0}^{\infty} m_j \phi_j(x)$ be an arbitrary function in the span of our basis functions ϕ_0, ϕ_1, \dots . Prove that truncating this series after N terms gives the best approximation to m in the span of the first N basis functions, i.e. that

$$m_0 \dots m_{N-1} = \operatorname{argmin}_{a_0 \dots a_{N-1}} \left\| m - \sum_{j=0}^{N-1} a_j \phi_j \right\|^2.$$

What is the error $\|m - \sum_{j=0}^{N-1} m_j \phi_j\|$ of this approximation?

Solution 3 See the beginning of the Lab on Implementing Sobolev Regression.

Now that we've got a formula for the error of the best approximation to a specific function m , we can maximize it over the functions m in our model \mathcal{M} to get the worst-case approximation error. To make sure you stay on track, I'll tell you what it is. It's $\lambda_N^{-1/2}$. Prove it!

Exercise 4 Prove that $\lambda_N^{-1/2} = \max_{m \in \mathcal{M}} \min_{a_0 \dots a_{N-1}} \|m - \sum_{j=0}^{N-1} a_j \phi_j\|$.

Then ‘invert’ this to work out how many you how many terms you need to include so the set \mathcal{M}^N contains an ϵ -approximation to every function in \mathcal{M} , i.e., find the smallest N for which $\max_{m \in \mathcal{M}} \min_{a_0 \dots a_{N-1}} \|m - \sum_{j=0}^{N-1} a_j \phi_j\| \leq \epsilon$. Report this in abstract terms, as something we can calculate using the sequence $\lambda_0 \leq \lambda_1 \leq \dots$, and in concrete terms for the Sobolev model from class, where $\lambda_j = (\pi j)^2$.

Tip. Keep in mind that $\sum_j m_j^2 = \sum_j m_j^2 \lambda_j / \lambda_j$. Look over the inner product spaces homework if you’re stuck.

Solution 4 See the beginning of the Lab on Implementing Sobolev Regression.

Now, if we ignore the ‘later in the semester problem’ mentioned above, we’ve got everything we need to implement Sobolev regression. If we can decide how accurately we want to calculate $\hat{\mu}$, we know how to approximate our model well enough to get that accuracy. One remaining question is whether we’re doing this the best way. That is, if we’re going to throw some finite-dimensional approximation to our model into **CVXR**, is \mathcal{M}^N the best one to use? In the next section, we’ll show that it is.

This isn’t the kind of thing we always need to know in practice. It’s ok to do things sub-optimally and long as we’re ok with how well we’re doing them. But knowing that we’re using the best possible N -dimensional basis means, for example, that if we’re at the limit of what our computer can handle and still not happy with our accuracy, we shouldn’t waste time trying to do better *using this approach based on finite-dimensional uniform approximation*. There are other ways to do computation in models like this. One approach, based on something called the *kernel trick*, isn’t that much harder to implement and can be a lot faster, but it requires that we do a bit more pen-and-paper math. Or, if we’re lucky, that somebody else has already done it for the version of the model (1) we’re using. We’re not going to have time to cover the kernel trick this semester, but I have covered it in previous versions of this class that emphasized different topics, so I can give you slides and exercises on the topic that use familiar-enough terms if you’re interested.

3 The Optimality of the Basis $\phi_0 \dots \phi_{N-1}$.

In this section, we’ll show that we can’t improve on the uniform error bound we got in Exercise 4 by using any other N -dimensional approximation to our model. Throughout, when we write $b_0 \dots b_{N-1}$, we will think of these as the first N vectors in an orthonormal basis b_0, b_1, \dots that spans the same vector space as ϕ_0, ϕ_1, \dots . This is ‘free’, or in more standard math language, without loss of generality. We lose nothing by assuming that $b_0 \dots b_{N-1}$ orthonormal, as we saw in Exercise 2, the uniform approximation error we get is the same for any two bases that span the same vector space. We lose nothing by thinking of $b_0 \dots b_{N-1}$ as the first N vectors in an orthonormal basis that *contains* ϕ_0, ϕ_1, \dots in its span, as we can *extend our basis* by adding vectors b_N, b_{N+1}, \dots until we

get (an orthonormal) one that contains ϕ_0, ϕ_1, \dots . And we'll see momentarily in Exercise 6 that, given any such basis b_0, b_1, \dots , we can find another one $\tilde{b}_0, \tilde{b}_1, \dots$ that's also *contained in* the span of ϕ_0, ϕ_1, \dots and gives us N -dimensional approximations that are at least as good, so we lose nothing by focusing on bases like $\tilde{b}_0, \tilde{b}_1, \dots$ in the first place.

3.1 Approximating a Single Function

We'll start by tackling a generalization of Exercise 3.

Exercise 5 Let m be an arbitrary function and $b_0 \dots b_{N-1}$ be the first N functions in an orthonormal basis b_0, b_1, \dots with $m \in \text{span } b_0 \dots b_{N-1}$ (i.e. a sequence b_0, b_1, \dots with $\langle b_i, b_j \rangle = 1$ if $i = j$ and 0 otherwise for which $m = \sum_{j=0}^{\infty} a_j b_j$ for some coefficients a_0, a_1, \dots). Prove that the best approximation to m in the span of $b_0 \dots b_{N-1}$, in the sense that $\|m - m^N\|$ is as small as possible, is $m^N = \sum_{j=0}^{N-1} \langle m, b_j \rangle b_j$ and that the squared error of this approximation is $\|m^N - m\|^2 = \sum_{j \geq N} \langle m, b_j \rangle^2$.

Solution 5 This is the same argument as in Exercise 3 in different notation. In particular, with b_j replacing ϕ_j and $\langle m, b_j \rangle$ replacing m_j . Here's the argument.

To start, observe that $m = \sum_{j=0}^{\infty} \langle m, b_j \rangle b_j$. We know that because the coefficients $a_j = \langle m, b_j \rangle$ are the only ones for which, for every basis function b_k , we have $\langle b_k, \sum_{j=0}^{\infty} a_j b_j \rangle = a_k$ equal to $\langle m, b_k \rangle$. For any coefficients $a_0 \dots a_{N-1}$,

$$\begin{aligned} \left\| m - \sum_{j=0}^{N-1} a_j b_j \right\|^2 &= \left\| \sum_{j=0}^{\infty} \langle m, b_j \rangle b_j - \sum_{j=0}^{N-1} a_j b_j \right\|^2 \\ &= \left\| \sum_{j=0}^{N-1} (\langle m, b_j \rangle - a_j) b_j + \sum_{j=N}^{\infty} \langle m, b_j \rangle b_j \right\|^2 \end{aligned}$$

We can simplify this by observing that, because our basis functions are orthonormal, the squared norm of a linear combination of them is just the sum of the squared coefficients.

$$\begin{aligned} \left\| \sum_{j=0}^{\infty} a_j b_j \right\|^2 &= \left\langle \sum_{j=0}^{\infty} a_j b_j, \sum_{k=0}^{\infty} a_k b_k \right\rangle \\ &= \sum_{j=0}^{\infty} \sum_{k=0}^{\infty} a_j a_k \langle b_j, b_k \rangle \\ &= \sum_{j=0}^{\infty} a_j^2 \quad \text{because} \quad \langle b_j, b_k \rangle = \begin{cases} 1 & \text{if } j = k \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

Applying this to our approximation error, we get this.

$$\left\| m - \sum_{j=0}^{N-1} a_j b_j \right\|^2 = \sum_{j=0}^{N-1} (\langle m, b_j \rangle - a_j)^2 + \sum_{j=N}^{\infty} \langle m, b_j \rangle^2$$

This is minimized when the first sum is zero, i.e. when $a_j = \langle m, b_j \rangle$ for $j = 0 \dots N-1$.

Exercise 6 Consider a basis $b_0 \dots b_{N-1}$ like the one in Exercise 5 and let m be an arbitrary function that is in both (i) the span of ϕ_0, ϕ_1, \dots and (ii) the span of b_0, b_1, \dots . Let $b_j^\phi = \sum_k \langle b_j, \phi_k \rangle \phi_k$ be the projection of the basis function b_j onto the span of ϕ_0, ϕ_1, \dots and $\tilde{b}_j = b_j^\phi / \|b_j^\phi\|$ be a normalized version of b_j^ϕ or, if $b_j^\phi = 0$, let $\tilde{b}_j = 0$. Prove that, if \tilde{m}^N is the best approximation to m in the span of $\tilde{b}_0 \dots \tilde{b}_{N-1}$, then $\|\tilde{m}^N - m\| \leq \|m^N - m\|$.

Tip. $a^2 - b^2 = (a-b)^2 + 2b(a-b)$ for any a and b . When a and b are $\langle m, b_j \rangle$ and $\langle m, \tilde{b}_j \rangle$ respectively, what do you know about $a - b$?

This sequence $\tilde{b}_0, \tilde{b}_1, \dots$ isn't necessarily a basis, but the only thing standing in the way is the possibility that some of them are zero, i.e. the possibility that $\langle b_j, \phi_k \rangle = 0$ for some j . Throw out the zero vectors and you've got a basis. And because the first N vectors in that basis include $\tilde{b}_0 \dots \tilde{b}_{N-1}$, they'll also give you approximation error that's as good or better than the first N vectors in the original basis b_0, b_1, \dots .

3.2 Uniform Approximation

Now that we've got a formula for the approximation error for a single function m , and we've established that we can restrict our attention to orthonormal bases b_0, b_1, \dots with exactly the same span as ϕ_0, ϕ_1, \dots , let's move on to the worst-case approximation error. What we're going to do is start with our formula for $\|m^N - m\|^2$ from Exercise 5, plug in the series expansion $m = \sum_j m_j \phi_j$, and interpret the maximization over these coefficients as an *eigenvalue problem*. To do this, we use the 'variational characterization' of the largest eigenvalue of a symmetric matrix A .

$$\max_{\|x\|_2 \leq 1} x^T A x \quad \text{is the largest eigenvalue of the symmetric matrix } A. \quad (5)$$

More generally, making the substitution $y = Sx$ for some invertible matrix S ,

$$\max_{\|Sx\|_2 \leq 1} x^T A x \quad \text{is the largest eigenvalue of the symmetric matrix } S^{-1} A S^{-1}. \quad (6)$$

Where does this symmetric matrix come from in our approximation problem? Let's plug in our series expansion $m = \sum_j m_j \phi_j$ into our approximation error formula and look.

$$\begin{aligned}
& \max_{m \in \mathcal{M}} \min_{a_0 \dots a_{N-1}} \left\| m - \sum_{j=0}^{N-1} a_j b_j \right\|^2 \\
&= \max_{m \in \mathcal{M}} \sum_{k \geq N} \langle m, b_k \rangle^2 \\
&= \max_{\substack{\text{sequences } m_0, m_1, \dots \\ \text{with } \sum_i \lambda_i m_i^2 \leq 1}} \sum_{k \geq N} \left\langle \sum_i m_i \phi_i, b_k \right\rangle \left\langle \sum_j m_j \phi_j, b_k \right\rangle \\
&= \max_{\substack{\text{sequences } m_0, m_1, \dots \\ \text{with } \sum_i \lambda_i m_i^2 \leq 1}} \sum_i \sum_j m_i m_j \sum_{k \geq N} B_{ik} B_{jk} \quad \text{for } B_{ik} = \langle \phi_i, b_k \rangle.
\end{aligned} \tag{7}$$

Take a look at the inner sum. If you're in the habit of writing out matrix products elementwise, i.e. writing $(BA)_{ij} = \sum_k B_{ik} A_{kj}$, it might jump out at you as something like the product $BB^T = \sum_k B_{ik} B_{jk}$. It's not quite that because we're summing over $k \geq N$, but we can include a diagonal 'selector matrix' S to throw out the terms we don't want.

$$\begin{aligned}
\max_{m \in \mathcal{M}} \min_{a_0 \dots a_{N-1}} \left\| m - \sum_{j=0}^{N-1} a_j b_j \right\|^2 &= \max_{\substack{\text{sequences } m_0, m_1, \dots \\ \text{with } \sum_i \lambda_i m_i^2 \leq 1}} \sum_i \sum_j m_i m_j \sum_k B_{ik} S_{kk} B_{jk} \\
&\quad \text{for } B_{ik} = \langle \phi_i, b_k \rangle \\
&\quad \text{and } S_{k\ell} = \begin{cases} 1 & \text{for } k = \ell \geq N \\ 0 & \text{otherwise} \end{cases}
\end{aligned} \tag{8}$$

At this point, let's translate things into matrix/vector notation. To avoid technicalities, we'll make a few simplifying assumptions from this point on.

1. We'll assume we're working with a finite-dimensional model, so all the sums we're working with are over $j \in 0 \dots K$ for some K . We'll get bounds that don't depend on K , so addressing the infinite-dimensional case amounts to taking a $K \rightarrow \infty$ limit.
2. We'll assume our λ_j are strictly positive and strictly increasing, i.e. that $0 < \lambda_0 < \lambda_1 < \dots$. This basically saves us the trouble of special-casing zeros and ties in our calculations. Once we've got an argument that works for λ_j like this, it's easy to add those cases in so it works for any positive increasing sequence $0 \leq \lambda_0 \leq \lambda_1 \leq \dots$.

In terms of the 'change of basis' matrix B and 'selector matrix' S from (8) and a 'diagonal matrix' Λ with elements $\Lambda_{ii} = \lambda_i$, we can write our maximal

approximation error like this.³

$$\begin{aligned}
\max_{m \in \mathcal{M}} \min_{a_0 \dots a_{N-1}} \left\| m - \sum_{j=0}^{N-1} a_j b_j \right\|^2 &= \max_{\|\Lambda^{1/2} m\|_2 \leq 1} m^T B S^N B^T m \\
&= \text{the largest eigenvalue of } \Lambda^{-1/2} B S^N B^T \Lambda^{-1/2} \\
&= \text{the largest eigenvalue of } S^N B^T \Lambda^{-1} B S^N \\
&\text{where } \Lambda_{ij} = \begin{cases} \lambda_i & \text{for } i = j \\ 0 & \text{otherwise} \end{cases}.
\end{aligned} \tag{9}$$

The last step, in which we turn our matrix product ‘inside out’, is one we haven’t yet justified.

Exercise 7 Suppose the matrix $A^T A$ has an eigenvalue λ with corresponding eigenvector v . Show that λ is also an eigenvalue of AA^T . Then explain how this justifies the last step in our calculation of the maximal approximation error.

Tip. What’s $AA^T Av$? What’s $S^N S^N$?

Solution 7 If $A^T A$ has an eigenvalue λ with corresponding eigenvector v , then λ is also an eigenvalue of AA^T with corresponding eigenvector Av .

$$(AA^T)(Av) = A(A^T Av) = A(\lambda v) = \lambda Av.$$

Applying this mechanically to $A = \Lambda^{1/2} B S^N$ and observing that $A^T = (S^N)^T B^T (\Lambda^{-1/2})^T = S^N B^T \Lambda^{-1/2}$ due to symmetry of S^N and $\Lambda^{1/2}$, this says that these two matrices have the same eigenvalues:

$$\begin{aligned}
A^T A &= S^N B^T \Lambda^{-1/2} \Lambda^{-1/2} B S^N = S^N B^T \Lambda^{-1} B S^N \\
AA^T &= \Lambda^{-1/2} B S^N S^N B^T \Lambda^{-1/2} = \Lambda^{-1/2} B^T S^N B \Lambda^{-1/2} \quad (\text{check that } S^N S^N = S^N).
\end{aligned}$$

Since they have the same eigenvalues, they have the same largest eigenvalue.

So far, this has been a lot of reading. There were a lot of little steps that are about as mechanical as adding 2+2 if you’re used to this stuff, but can stop you in your tracks if you’re not. I’ve written them up for you to make sure you got a chance to tackle this last step, which calls for a little more thought.

Exercise 8 Prove that, if $b_0 \dots b_K$ is an orthonormal basis that spans the same vector space as $\phi_0 \dots \phi_K$, then $\max_{m \in \mathcal{M}} \min_{a_0 \dots a_{N-1}} \left\| m - \sum_{j=0}^{N-1} a_j b_j \right\|^2 \geq \lambda_N^{-1}$. Briefly explain why this, in combination with Exercise 4, implies that our model’s Kolmogorov N -width $w_N(\mathcal{M})$ is λ_N^{-1} .

Tip. If you can find some unit-length vector x for which $x^T S^N B^T \Lambda^{-1} B S^N x \geq \lambda_N^{-1}$, you’re done. Why?

Tip. What do you know about the change of basis matrix B ? If $B_{i \cdot}$ and $B_{\cdot j}$ are its i th row and j th column, what are $\langle B_{i \cdot}, B_{k \cdot} \rangle$ and $\langle B_{\cdot j}, B_{\cdot \ell} \rangle$?

³Here $\Lambda^{1/2}$ can be any symmetric matrix with $\Lambda^{1/2} \Lambda^{1/2} = \Lambda$, but we may as well make it the diagonal matrix with elements $\Lambda_{ii}^{1/2}$.