

Vector Spaces (Lab 1)

BST 235: Advanced Regression and Statistical Learning

Alex Levis, Fall 2019

1 Some motivation

In the first lecture, we laid out the basic regression framework that we will study in this course. At hand we have a collection of inputs $\mathbf{X}_1, \dots, \mathbf{X}_n \in \mathcal{X} \subseteq \mathbb{R}^d$, a corresponding collection of outputs $Y_1, \dots, Y_n \in \mathcal{Y}$ (for now we consider $\mathcal{Y} \subseteq \mathbb{R}$), and we are tasked vaguely with using this data to construct a function $g : \mathcal{X} \rightarrow \mathcal{Y}$ that can predict new outputs given new inputs, in some ‘good’ way. We introduced two possible probabilistic frameworks for describing how the data arise,

- Random design: $(Y_1, \mathbf{X}_1), \dots, (Y_n, \mathbf{X}_n) \stackrel{\text{iid}}{\sim} P$
- Fixed design: $(\mathbf{X}_1, \dots, \mathbf{X}_n) \sim P^*$, and given $(\mathbf{X}_1, \dots, \mathbf{X}_n)$, the outputs Y_1, \dots, Y_n are independently drawn from $P_{Y|\mathbf{X}}$, induced by the joint distribution P . In particular, we assume that for $i \in \{1, \dots, n\}$,

$$Y_i \perp ((Y_1, \mathbf{X}_1), \dots, (Y_{i-1}, \mathbf{X}_{i-1}), (Y_{i+1}, \mathbf{X}_{i+1}), \dots, (Y_n, \mathbf{X}_n)) \mid \mathbf{X}_i.$$

Using the notions of *loss* and *risk*, we decided that in a certain sense, the conditional mean $g_P(\mathbf{X}) = \mathbb{E}_P[Y \mid \mathbf{X}]$, called the *regression function*, was the optimal choice for g (how precisely was this decided?). Since we cannot possibly know this function exactly, a good idea is to use the data at hand to construct a *regression estimator* \hat{g}_n to approximate g_P , and subsequently *empirical risk minimization* (ERM) was proposed as a reasonable way to choose \hat{g}_n . Unfortunately, even in the simple setting where \mathbf{X} is a one-dimensional continuous variable, this approach fails unless we make further assumptions on the form of $\mathbb{E}_P[Y \mid \mathbf{X}]$.

Indeed, we will spend much of this course assuming a *linear model* for $\mathbb{E}_P[Y \mid \mathbf{X}]$. Specifically, we will postulate the existence of $\beta(P) \in \mathbb{R}^d$ such that

$$\mathbb{E}_P[Y \mid \mathbf{X}] = \mathbf{X}^T \beta(P) = \sum_{j=1}^d \beta_j(P) X_j.$$

Strictly speaking, with this modeling restriction and the mechanism of ERM, we could go through the brute force calculus to derive the least squares estimator $\hat{\beta}$ and end the story with the regression estimator $\hat{g}_n(\mathbf{X}) = \mathbf{X}^T \hat{\beta}$. However, doing only this would be foolish, as it would ignore the deep and beautiful structure of both the statistical model, and the optimal estimators that fall out. The mathematical language best suited for understanding linear models is, quite naturally, linear algebra, and this is what we will explore in the coming weeks.

2 Fundamental vector space concepts

Our study of linear algebra in BST 235 will cover the following broad topics:

- Vector spaces, subspaces, basis and dimension
- Inner products and norms, orthogonal projection
- Linear maps and matrices, rank, inverses
- Spectral / singular value decomposition of matrices, generalized inverses

Today, we will take as our goal to understand the ideas in the first bullet point.

Definition 1. Suppose we have a non-empty set V , a field \mathbb{F} and operations $\oplus : V \times V \rightarrow V$, $\odot : \mathbb{F} \times V \rightarrow V$. The triple (V, \oplus, \odot) , or just V if clear from context, is called a *vector space* if the following axioms are satisfied.

(1) **Vector addition:**

- (a) associativity: $(v_1 \oplus v_2) \oplus v_3 = v_1 \oplus (v_2 \oplus v_3)$, for all $v_1, v_2, v_3 \in V$,
- (b) identity element: $\exists 0_V \in V$ such that $v \oplus 0_V = v$, for all $v \in V$,
- (c) commutativity: $v_1 \oplus v_2 = v_2 \oplus v_1$, for all $v_1, v_2 \in V$,
- (d) inverse element: $\forall v \in V, \exists -v \in V$ such that $-v \oplus v = 0_V$.

(2) **Scalar multiplication:**

- (a) associativity: $(a_1 \cdot a_2) \odot v = a_1 \odot (a_2 \odot v)$, for all $a_1, a_2 \in \mathbb{F}, v \in V$,
- (b) identity element: $1_{\mathbb{F}} \odot v = v$, for all $v \in V$,
- (c) distributivity wrt vector addition: $a \odot (v_1 \oplus v_2) = (a \odot v_1) \oplus (a \odot v_2)$, $\forall a \in \mathbb{F}, v_1, v_2 \in V$.
- (d) distributivity wrt to field addition: $(a_1 + a_2) \odot v = (a_1 \odot v) \oplus (a_2 \odot v)$, $\forall a_1, a_2 \in \mathbb{F}, v \in V$.

Remark 1. In this course, we will exclusively consider $\mathbb{F} = \mathbb{R}$, and we refer to V as a *real vector space*. In this case we can unambiguously write 0 and 1 for $0_{\mathbb{F}}$ and $1_{\mathbb{F}}$, respectively. Note that these axioms were historically landed upon due to their efficiency — they are minimal, but imply all the properties we would expect. As an exercise, you might show the following consequences of the above definition:

- If $u \oplus v = u \oplus w$ then $v = w$ (cancellation).
- The zero vector 0_V is unique, as are additive inverses.
- For any $v \in V$, $0 \odot v = 0_V$.
- For any $a \in \mathbb{F}$, $a \odot 0_V = 0_V$.
- $(-1) \odot v = -v$, for all $v \in V$.

Finally, note that we will often use $+$ instead of \oplus , av for $a \odot v$, and infer whether we mean vector or field operations based on the context.

Definition 2. Suppose (V, \oplus, \odot) is a vector space over the field \mathbb{F} . We say that $U \subseteq V$ is a *linear subspace*, or *subspace*, of V if (U, \oplus, \odot) is a vector space over \mathbb{F} , where \oplus and \odot are restricted to U .

Lemma 1. Given vector space (V, \oplus, \odot) over \mathbb{F} , the set $U \subseteq V$ is a subspace of V if and only if

- (i) $U \neq \emptyset$ (equivalently, check $0_V \equiv 0_U \in U$),
- (ii) $u, v \in U \implies u \oplus v \in U$ (closure under addition), and
- (iii) $a \in \mathbb{F}, v \in U \implies av \in U$ (closure under scalar multiplication).

Now that we have defined vector spaces as sets of elements for which addition and scalar multiplication are well-behaved, we can now introduce two dual concepts that characterize collections of vectors. One concept is *span*, which describes the overall expressiveness of a set of vectors. The other concept is *linear independence*, which pertains to the non-redundancy in such a collection. Formally, these are defined as follows:

Definition 3. Suppose (V, \oplus, \odot) is a vector space over \mathbb{F} . The *linear span* of the collection of vectors $\{v_1, \dots, v_k\} \subseteq V$ is the set of all linear combinations of these vectors,

$$(a_1 \odot v_1) \oplus \dots \oplus (a_k \odot v_k) =: \sum_{\ell=1}^k a_\ell v_\ell.$$

We write

$$\mathcal{L}(v_1, \dots, v_k) = \left\{ \sum_{\ell=1}^k a_\ell v_\ell \mid a_1, \dots, a_k \in \mathbb{F} \right\}.$$

We say that the vectors v_1, \dots, v_k *span* a subspace $U \subseteq V$ if $\mathcal{L}(v_1, \dots, v_k) = U$.

The vectors v_1, \dots, v_k are called *linearly independent* if

$$0_V = \sum_{\ell=1}^k a_\ell v_\ell \implies 0 = a_1 = \dots = a_k,$$

i.e., there is no non-trivial way to combine the vectors to yield 0_V . If there exist, $a_1, \dots, a_k \in \mathbb{F}$ not all zero such that $0_V = \sum_{\ell=1}^k a_\ell v_\ell$, then v_1, \dots, v_k are called *linearly dependent*.

Definition 4. Let (V, \oplus, \odot) be a vector space over \mathbb{F} . A collection $\{v_1, \dots, v_k\} \subseteq V$ is called a *basis* for V if v_1, \dots, v_k span V and are linearly independent.

Remark 2. A basis for a vector space essentially comprises a coordinate system — a minimal (cf. linearly independent) set of points that are sufficiently expressive (cf. spanning) to describe the whole space through linear combinations.

Two issues arise immediately: existence and uniqueness. With a slight expanding of the above definitions of span and linear independence to account for infinite sets, it **can be shown** that bases always exist. On the other hand, coordinate systems are not unique, and therefore neither are bases. Nevertheless, the number of elements in a basis is a constant for a given vector space, as we will demonstrate via a key lemma in class. This will allow us to define dimension, $\dim(V)$, as the cardinality of an arbitrary basis for V .

3 Exercises

Exercise 1. Convince yourself that the following are real vector spaces, by specifying natural choices for \oplus , \odot , 0_V , and additive inverses:

(a) $\mathbb{R}^d = \{[x_1 \cdots x_d]^T \mid x_1, \dots, x_d \in \mathbb{R}\}$, for $d \in \mathbb{N}$.

We use component-wise addition and scalar multiplication. Let $a \in \mathbb{R}$, $\mathbf{x}_j = [x_{j1} \cdots x_{jd}]^T \in \mathbb{R}^d$, for $j = 1, 2$. Then define

$$\mathbf{x}_1 \oplus \mathbf{x}_2 = \begin{bmatrix} x_{11} \\ \vdots \\ x_{1d} \end{bmatrix} \oplus \begin{bmatrix} x_{21} \\ \vdots \\ x_{2d} \end{bmatrix} := \begin{bmatrix} x_{11} + x_{21} \\ \vdots \\ x_{1d} + x_{2d} \end{bmatrix}, \text{ and } a \odot \mathbf{x}_1 = a \odot \begin{bmatrix} x_{11} \\ \vdots \\ x_{1d} \end{bmatrix} := \begin{bmatrix} ax_{11} \\ \vdots \\ ax_{1d} \end{bmatrix}$$

Now that these are defined, we can clearly use $0_{\mathbb{R}^d} := \mathbf{0} = [0 \cdots 0]^T$ for the additive identity, and $-\mathbf{x}_j := (-1) \odot \mathbf{x}_j$ for an additive inverse. All the vector space axioms can be verified as direct consequences of corresponding properties of the real numbers.

(b) For a probability space (Ω, \mathcal{A}, P) , the space

$$L_2(P) = \{X : \Omega \rightarrow \mathbb{R} \mid X \text{ measurable, } \mathbb{E}_P(X^2) < \infty\}.$$

Now we use pointwise addition and scalar multiplication for functions in $L_2(P)$. Let $a \in \mathbb{R}$, $X_1, X_2 \in L_2(P)$, then define $X_1 \oplus X_2$ and $a \odot X_1$ via

$$(X_1 \oplus X_2)(\omega) := X_1(\omega) + X_2(\omega), \forall \omega \in \Omega, \text{ and } (a \odot X_1)(\omega) := aX_1(\omega), \forall \omega \in \Omega.$$

For the additive identity, we can use the constant zero random variable $0_{L_2(P)} := 0$ (in the sense of the function that is identically 0 over Ω), and also $-X_1$ is defined via $(-1) \odot X_1$. Again, the vector space axioms can be checked using properties of the real numbers.

One more technical detail: we need to ensure $L_2(P)$ is closed under linear combinations! Suppose $X_1, X_2 \in L_2(P)$. Scalar multiplication is easy: if $a \in \mathbb{R}$, $\mathbb{E}[(aX_1)^2] = a^2\mathbb{E}[X_1^2] < \infty$, so $aX_1 \in L_2(P)$, by linearity. Addition is a little trickier. Consider the following inequality for $x, y \in \mathbb{R}$:

$$(x + y)^2 = x^2 + 2xy + y^2 \leq x^2 + (x^2 + y^2) + y^2 = 2(x^2 + y^2),$$

since

$$0 \leq (x - y)^2 = x^2 - 2xy + y^2 \implies 2xy \leq x^2 + y^2.$$

It follows that

$$\mathbb{E}[(X_1 + X_2)^2] \leq \mathbb{E}[2(X_1^2 + X_2^2)] = 2\mathbb{E}[X_1^2] + 2\mathbb{E}[X_2^2] < \infty,$$

by linearity, so $X_1 + X_2 \in L_2(P)$. Alternatively, as Izzy explained in lab, we can use Cauchy-Schwarz! We will soon study inner products in lecture: for vector space V , $\langle \cdot, \cdot \rangle : V \times V \rightarrow \mathbb{R}$ is a (real) inner product on V if it satisfies certain axioms. We will see that Cauchy-Schwarz is a general result for any inner product space: if $(V, \langle \cdot, \cdot \rangle)$ is a real inner product space, and $v_1, v_2 \in V$, then $|\langle v_1, v_2 \rangle| \leq \|v_1\| \|v_2\|$, where $\|v_j\| := \sqrt{\langle v_j, v_j \rangle}$. Now, it turns out that $L_2(P)$ is an inner product space with $\langle X_1, X_2 \rangle := \mathbb{E}[X_1 X_2]$. Therefore we immediately know $\mathbb{E}[X_1 X_2] \leq \sqrt{\mathbb{E}[X_1^2] \mathbb{E}[X_2^2]} < \infty$ when $X_1, X_2 \in L_2(P)$, so $\mathbb{E}[(X_1 + X_2)^2] = \mathbb{E}[X_1^2] + 2\mathbb{E}[X_1 X_2] + \mathbb{E}[X_2^2] < \infty$.

Exercise 2. Taking $V = \mathbb{R}^2$, give examples of a non-empty subset $U \subseteq V$ that is *not* a subspace, but that *does* satisfy

- (a) Closure under addition, and remaining vector space properties

Take $U = \mathbb{Q}^2 = \{(x_1, x_2) \mid x_1, x_2 \in \mathbb{Q}\}$. I'll leave it to you folks to think about why this example fails for the remaining property, scalar multiplication.

- (b) Closure under scalar multiplication, and remaining vector space properties

Any union of distinct lines through the origin should work. For instance, take

$$U = \{(x, 0) \mid x \in \mathbb{R}\} \cup \{(0, y) \mid y \in \mathbb{R}\}.$$

Same deal as above — you guys convince yourselves this doesn't satisfy closure under addition.

Exercise 3. Let V be a vector space, and $U, W \subseteq V$ be two linear subspaces. Show that

- (a) $U \cap W$ is a linear subspace of V .

Actually, this result holds in greater generality. Let $\mathcal{B} \subseteq \mathcal{P}(V)$ be a class of linear subspaces of V . We will show that $V^* = \bigcap_{W \in \mathcal{B}} W$ is a linear subspace. First, $0_V \in V^*$, since $0_V \in W$, for all $W \in \mathcal{B}$, as these are all subspaces. It remains to show V^* is closed under linear combinations. For any $\alpha, \beta \in \mathbb{F}$ and $v_1, v_2 \in V^*$, we must have $v_1, v_2 \in W$, for all $W \in \mathcal{B}$. Hence, $\alpha v_1 + \beta v_2 \in W$, for all $W \in \mathcal{B}$, since each W is a subspace. Thus, $\alpha v_1 + \beta v_2 \in \bigcap_{W \in \mathcal{B}} W = V^*$, so the intersection is closed under linear combinations, and therefore is a subspace.

- (b) $U \cup W$ is a subspace iff $U \subseteq W$ or $W \subseteq U$.

If $U \subseteq W$ (or $W \subseteq U$) then $U \cup W = W$ (or $U \cup W = U$), which is a subspace.

Conversely, assume that $U \cup W$ is a subspace, and suppose that one set is not contained in the other. Then there exist $v_1 \in U \setminus W$ and $v_2 \in W \setminus U$. Since $U \cup W$ is a subspace, it must hold that $v_1 + v_2 \in U \cup W$, because $v_1, v_2 \in U \cup W$. There are two cases: (1) if $v_1 + v_2 \in U$, then $v_2 = (v_1 + v_2) - v_1 \in U$, since U is a subspace; (2) $v_1 + v_2 \in W$, then $v_1 = (v_1 + v_2) - v_2 \in W$, since W is a subspace. Either case represents a contradiction of our choice of v_1, v_2 , so we conclude that one set must be contained in the other.

Exercise 4. Let V be a vector space, and suppose $\{v_1, \dots, v_k\}$ are linearly independent. Show that $0_V \notin \{v_1, \dots, v_k\}$.

We will show the contrapositive. Suppose $0_V \in \{v_1, \dots, v_k\}$, and let $j \in \{1, \dots, k\}$ be such that $v_j = 0_V$. Then we can take $a_j = 1$, and $a_\ell = 0$ for all $\ell \neq j$, so that

$$\sum_{\ell=1}^k a_\ell v_\ell = a_j v_j = v_j = 0_V.$$

By definition, we conclude that $\{v_1, \dots, v_k\}$ are linearly dependent.

Exercise 5. Let V be a vector space, and suppose $\{v_1, \dots, v_k\}$ are linearly dependent. Show that there exists $1 \leq j \leq k$, and scalars a_1, \dots, a_{j-1} such that

$$v_j = \sum_{\ell=1}^{j-1} a_\ell v_\ell,$$

where a sum from 1 to 0 is defined as 0_V .

By linear dependence, there exist scalars $a_1, \dots, a_k \in \mathbb{R}$, not all zero, such that

$$\sum_{\ell=1}^k a_\ell v_\ell = 0_V. \tag{1}$$

Consider the set $\mathcal{J} = \{\ell \in \{1, \dots, k\} \mid a_\ell \neq 0\}$, and choose $j = \max(\mathcal{J})$. If $j = 1$, then by (1) and the definition of j , we must have

$$a_1 v_1 = 0_V \implies v_1 = 0_V,$$

which satisfies the claim by our definition of a sum from $\ell = 1$ to 0. If $j \geq 2$, then again by (1) and our choice of j ,

$$\sum_{\ell=1}^j a_\ell v_\ell = 0_V \implies a_j v_j = - \sum_{\ell=1}^{j-1} a_\ell v_\ell.$$

Dividing through by $a_j \neq 0$, we obtain

$$v_j = \sum_{\ell=1}^{j-1} \left(-\frac{a_\ell}{a_j} \right) v_\ell,$$

proving the claim.