

Linear Algebra and Robot Modeling

Nathan Ratliff

Apr 8, 2014

Abstract

Linear algebra is fundamental to robot modeling, control, and optimization. This document reviews some of the basic kinematic equations and uses them to motivate an SVD-centric geometric perspective on linear algebra. This perspective illuminates the underlying structure and behavior of linear maps and simplifies analysis, especially for reduced rank matrices. We later review some properties of multidimensional quadratic forms that are central to optimal control. Most of this material should already be familiar to the reader; this document explores some specifics and offers a potentially unique intuition oriented perspective.

1 Basic kinematic equations

Let $\mathcal{C} \subset \mathbb{R}^d$ be a configuration space. For instance, the configuration space may be a space of valid joint angles for a robot manipulator. Consider as a running example the differentiable map mapping a point in the configuration space (a particular set of joint angles) to the three-dimensional location of the manipulator's fingertip. We denote this *forward kinematics* map abstractly as $\phi: \mathcal{C} \rightarrow \mathbb{R}^3$.

This map is no different from any other multidimensional differentiable mapping from calculus, so we can ask the typical calculus questions such as how does the output change with changes in the input, etc. Kinematic equations are none other than equations relating derivatives of the inputs to derivatives of the outputs.

Two common kinematic equations are

$$\dot{\mathbf{x}} = \mathbf{J}_\phi \dot{\mathbf{q}} \quad \text{and} \quad \ddot{\mathbf{x}} = \mathbf{J}_\phi \ddot{\mathbf{q}} + \dot{\mathbf{J}}_\phi \dot{\mathbf{q}}, \quad (1)$$

where $\mathbf{J}_\phi = \frac{\partial \phi}{\partial \mathbf{q}}$ is the Jacobian (total derivative) of ϕ . These equations are easy to write down, especially when referencing tables of derivative formulas in an Advanced Calculus text, but in robotics we need a strong intuitive understanding of what they mean. For instance, two questions that may pop to mind are

1. How does that first equation relate to the equation $\delta \mathbf{x} = \mathbf{J}_\phi \delta \mathbf{q}$?
2. What does $\dot{\mathbf{J}}_\phi$ actually mean, and how do we compute it?

1.1 Interpreting time derivatives

Whenever an equation uses time derivatives, such as $\dot{\mathbf{q}}$, there is an implicit assumption that \mathbf{q} is really a trajectory $\mathbf{q} : [0, T] \rightarrow \mathbb{R}^d$. The time derivatives give the position $\mathbf{q} = \mathbf{q}(t)$, velocity $\dot{\mathbf{q}} = \frac{d}{dt}\mathbf{q}(t)$, and acceleration $\ddot{\mathbf{q}} = \frac{d}{dt}\dot{\mathbf{q}}(t)$ of the trajectory at some time t .

Thus, the equation $\dot{\mathbf{x}} = \mathbf{J}_\phi \dot{\mathbf{q}}$ relates how the velocity $\dot{\mathbf{q}}$ in the system's configuration space relates to the velocity of the point on the end-effector in Cartesian space. Concretely, using our manipulator example, it tells us how joint velocities (the rate at which each joint is changing over time) relate to the finger tip's velocity (the rate at which the Cartesian dimensions of the fingertip are changing over time). The implicit assumption here is always that the system is moving along some trajectory $\mathbf{q}(t)$, and that that trajectory gives us the time variable that allows us to address questions of how the system evolves over time.

Now given that q refers implicitly to an underlying trajectory, we can now interpret what \mathbf{J}_ϕ means. As the system moves along its trajectory from $\mathbf{q}(t)$ to a point dt in the future $\mathbf{q}(t + dt)$, the associated Jacobian changes slightly as well since it's a non-constant function of t . $\dot{\mathbf{J}}_\phi$ gives explicitly the rate at which the Jacobian is changing with time:

$$\dot{\mathbf{J}}_\phi = \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} \left(\mathbf{J}_\phi(\mathbf{q}(t + \Delta t)) - \mathbf{J}_\phi(\mathbf{q}(t)) \right) \quad (2)$$

We usually calculate this quantity numerically by taking a finite-difference between the Jacobian now and the Jacobian we find a little bit in the direction of the current velocity. That approximation can be derived from the above using a first-order Taylor expansion and a finite time step:

$$\dot{\mathbf{J}}_\phi \approx \frac{1}{\Delta t} \left(\mathbf{J}_\phi(\mathbf{q}(t + \Delta t)) - \mathbf{J}_\phi(\mathbf{q}(t)) \right) \quad (3)$$

$$= \frac{1}{\Delta t} \left(\mathbf{J}_\phi(\mathbf{q}(t) + \Delta t \dot{\mathbf{q}}) - \mathbf{J}_\phi(\mathbf{q}(t)) \right). \quad (4)$$

1.2 Tangent vectors and co-tangent vectors

How does the expression $\delta \mathbf{x} = \mathbf{J}_\phi \delta \mathbf{q}$ differ from the expression above relating velocities $\dot{\mathbf{x}} = \mathbf{J}_\phi \dot{\mathbf{q}}$? Notationally, $\delta \mathbf{q}$ denotes a small displacement of \mathbf{q} which is qualitatively similar to a velocity vector, but depending on whether you're an engineer with a tool belt full of math tricks or a mathematical purist the two expressions can range from being "slightly different but heuristically the same" to being "fundamentally different". Clearly, it's superficially similar because derivatives work in very structured ways, but this equation doesn't assume that there's any underlying trajectory. The quantity $\delta \mathbf{q}$ should be thought of as a *small perturbation* in \mathbf{q} , and the equation tells us how small perturbations in \mathbf{q} result in small perturbations of \mathbf{x} .

There are technical details that we won't worry about in this class, but rigorous treatments of these ideas can be found in textbooks on differential geometry. There, $\dot{\mathbf{q}}$ is known as a *tangent* vector, and these small perturbations

$\delta\mathbf{q}$ are known as *co-tangent* vectors. The former generalizes tangents to curves across a manifold, and the latter generalizes gradients of scalar functions defined on the manifold.

Within robotics, rigor is often secondary to physical intuition, so definitions and manipulations of these ideas are less stringent. The bottom line is that $\dot{\mathbf{q}}$ always refers to a time derivative of an underlying trajectory, whereas $\delta\mathbf{q}$ will simply denote a perturbation (small movement of) the point \mathbf{q} in space with no notion of *time*.

We'll often use $\delta\mathbf{q}$ to also represent a finite-sized perturbation away from some approximation point \mathbf{q} , which lets us write the first-order Taylor expansion of a scalar function $c : \mathcal{C} \rightarrow \mathbb{R}$ as

$$c(\mathbf{q} + \delta\mathbf{q}) \approx c(\mathbf{q}) + \nabla c(\mathbf{q})^T \delta\mathbf{q}. \quad (5)$$

In this notation, the expression $\delta\mathbf{x} = (\mathbf{x} + \delta\mathbf{x}) - \mathbf{x} \approx \mathbf{J}_\phi \delta\mathbf{q}$ can be viewed as a first-order Taylor expansion of the map ϕ . Following this line of thought, we can also expand the approximation to a second order and write¹

$$\delta\mathbf{x} \approx \mathbf{J}_\phi \delta\mathbf{q} + \frac{1}{2} [\delta\mathbf{q}^T \nabla^2 \phi_k(\mathbf{q}) \delta\mathbf{q}]_k \quad (6)$$

This notion of second-order expansion doesn't make sense for $\dot{\mathbf{x}} = \mathbf{J}_\phi \dot{\mathbf{q}}$ since we always treat this latter equation as an exact differential expression relating time-rates of change.

2 A geometric perspective on linear algebra

Linear algebra is fundamental to control and optimization. This section reviews some of the most important properties of linear maps from the perspective of the Singular Value Decomposition (SVD). If we start from the notion that every matrix has an SVD, the structure and behavior of the matrix becomes immediately clear and algebraic manipulations of even reduced rank matrices are straightforward. We assume familiarity with the basic linear algebraic concepts and focus here on understanding the underlying geometry of linear maps.

2.1 Linear maps are bijections between fundamental spaces

The *row space* and *column space* are fundamental to a linear map. We'll see in this section that a linear map forms a bijection² between between these two spaces, and that all components in the orthogonal compliment to these spaces are simply removed when necessary. These orthogonal compliments are known as the left and right *null spaces*, respectively.

¹For those familiar with tensor notation and the Einstein summation convention, this second term can be more compactly written as $\frac{1}{2} (\partial_{ij} \phi_k) \delta q^i \delta q^j$.

²A bijection is formally a *one-to-one* and *onto* function. One may view it as a full pairing between points in the domain and range. Every domain point \mathbf{x} has an corresponding range point \mathbf{y} and vice versa under the bijection.

Usually, the row and column spaces are defined as the space spanning the rows and the space spanning the columns. Those definitions, though true, aren't very insightful. This section shows that the SVD provides a nice geometric view of what these spaces are and how the linear map operates on them. This decomposition gives geometric insight into the fundamental nature of linear algebra.

Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ be a matrix representing a linear map from \mathbb{R}^n to \mathbb{R}^m . We know from what is sometimes referred to as the *fundamental theorem of linear algebra* that \mathbf{A} has a Singular Value Decomposition (SVD) of the form $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^T$, where \mathbf{U} and \mathbf{V} are orthogonal matrices and $\mathbf{S} \in \mathbb{R}^{m \times n}$ is a diagonal matrix of singular values $\sigma_1, \dots, \sigma_k$. We don't assume that \mathbf{A} is full rank, so k may be less than both m and n . Since \mathbf{S} is a non-square diagonal matrix with only $k \leq m, n$ nonzero entries, we can better reveal its structure by writing it as

$$\mathbf{S} = \begin{pmatrix} \mathbf{\Sigma} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} = \begin{pmatrix} \mathbf{I} \\ \mathbf{0} \end{pmatrix} \mathbf{\Sigma} \begin{pmatrix} \mathbf{I} & \mathbf{0} \end{pmatrix} \quad (7)$$

where the $\mathbf{0}$ matrices are appropriately sized and $\mathbf{\Sigma} = \mathbf{diag}(\sigma_1, \dots, \sigma_k)$.

Since \mathbf{S} as shown in Equation 7 decomposes as a square diagonal matrix $\mathbf{\Sigma}$ with zero buffers on either side, we can further decompose both \mathbf{U} and \mathbf{V} into

$$\mathbf{U} = \begin{pmatrix} \mathbf{U}_{//} & \mathbf{U}_{\perp} \end{pmatrix} \quad \text{and} \quad \mathbf{V} = \begin{pmatrix} \mathbf{V}_{//} & \mathbf{V}_{\perp} \end{pmatrix}, \quad (8)$$

where the columns with the subscript \perp are the ones annihilated by the zeros and the columns with the subscript $//$ are those that remain. This decomposition allows us to rewrite the SVD as

$$\mathbf{A} = \begin{bmatrix} \mathbf{U}_{//} & \mathbf{U}_{\perp} \end{bmatrix} \begin{pmatrix} \mathbf{\Sigma} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \begin{bmatrix} \mathbf{V}_{//}^T \\ \mathbf{V}_{\perp}^T \end{bmatrix} \quad (9)$$

$$= \left(\begin{bmatrix} \mathbf{U}_{//} & \mathbf{U}_{\perp} \end{bmatrix} \begin{bmatrix} \mathbf{I} \\ \mathbf{0} \end{bmatrix} \right) \mathbf{\Sigma} \left(\begin{bmatrix} \mathbf{V}_{//} & \mathbf{V}_{\perp} \end{bmatrix} \begin{bmatrix} \mathbf{I} \\ \mathbf{0} \end{bmatrix} \right)^T \quad (10)$$

$$= \mathbf{U}_{//} \mathbf{\Sigma} \mathbf{V}_{//}^T. \quad (11)$$

It's fairly easy to show that, in terms of this decomposition, $\mathbf{span}(\mathbf{V}_{//})$ is the column space, $\mathbf{span}(\mathbf{V}_{\perp})$ is the right null space, $\mathbf{span}(\mathbf{U}_{//})$ is the row space, and $\mathbf{span}(\mathbf{U}_{\perp})$ is the left null space (which we'll see below is the null space of the natural generalized inverse).

The last expression in Equation 11, known as the *thin* SVD, reveals the fundamental structure of any (possibly reduced rank) matrix. Some, depending on background, find it more insightful to view that expression as

$$\mathbf{A} = \mathbf{U}_{//} \mathbf{\Sigma} \mathbf{V}_{//}^T = \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^T, \quad (12)$$

where \mathbf{u}_i and \mathbf{v}_i are the k columns of $\mathbf{U}_{//}$ and $\mathbf{V}_{//}$, respectively.

Geometrically, this expression says that all rank k matrices are simultaneously a correspondence between a k -dimensional orthogonal basis in the domain space $\mathcal{B}_V = \{\mathbf{v}_i\}_{i=1}^k$ and a k -dimensional orthogonal basis in the range space $\mathcal{B}_U = \{\mathbf{u}_i\}_{i=1}^k$ with a rule for how vectors should be stretched or compressed (or reflected when $\sigma_i < 0$) along those dimensions. (We refer to these bases below as simply the matrices \mathbf{V}_\parallel and \mathbf{U}_\parallel , respectively.)

It's useful to think about each term in the expansion in relation to its corresponding *endomorphism*³ $\sum_i \sigma_i \mathbf{u}_i \mathbf{u}_i^T$. Each term of this endomorphism acts on a vector \mathbf{x} to find its component $(\mathbf{u}_i^T \mathbf{x}) \mathbf{u}_i$ along the basis element \mathbf{u}_i and stretch it by σ_i . Since each basis element is orthogonal to all others, this endomorphism simply stretches, compresses, and/or reflects the vector independently along the given basis directions.

Equation 12, therefore, shows that all matrices have this fundamental behavior. Except in this case, when the domain and range spaces differ, we need to both decompose the vector \mathbf{x} in terms of the domain basis and *simultaneously map* that decomposition onto the corresponding basis of the range space. Once that connection between the (sub)spaces is established the matrix can apply its underlying operation defined by its singular values σ_i .

The following equivalent expressions illustrate the underlying behavior of the linear map from multiple perspectives:

$$\mathbf{A} = \mathbf{U}_\parallel \boldsymbol{\Sigma} \mathbf{V}_\parallel^T = \underbrace{(\mathbf{U}_\parallel \mathbf{V}_\parallel^T)}_{\text{Stretch then map}} \underbrace{(\mathbf{V}_\parallel \boldsymbol{\Sigma} \mathbf{V}_\parallel^T)}_{\text{Map then stretch}} = \underbrace{(\mathbf{U}_\parallel \boldsymbol{\Sigma} \mathbf{U}_\parallel^T)}_{\text{Map then stretch}} \underbrace{(\mathbf{U}_\parallel \mathbf{V}_\parallel^T)}_{\text{Stretch then map}}, \quad (13)$$

Each gives a subtly different way of thinking about how the matrix \mathbf{A} operates. The first says that we can think of \mathbf{A} physically as a stretching/squishing of the domain space followed by an incompressible mapping between orthogonal bases (simply mapping each domain basis element to its corresponding range basis element), whereas the second says we can equivalently view \mathbf{A} as first a mapping between basis elements followed by a stretching of the space. It's natural then to think of \mathbf{A} 's action holistically as simply an association between the two k -dimensional subspaces defined by $\mathbf{U}_\parallel \mathbf{V}_\parallel^T$ and a corresponding stretching/squishing/reflection of that unified space.

2.2 A natural generalized inverse

Note that the basis-correspondence mappings of the form $\mathbf{U}_\parallel \mathbf{V}_\parallel^T = \sum_{i=1}^k \mathbf{u}_i \mathbf{v}_i^T$ are all between k -dimensional subspaces, not the full space. What happens to the dimensions of \mathbf{x} orthogonal to that space in the domain in the expression $\mathbf{A}\mathbf{x} = (\mathbf{U}_\parallel \boldsymbol{\Sigma} \mathbf{V}_\parallel^T) \mathbf{x}$? They vanish!⁴ There aren't enough singular values to represent those dimensions, so simply by the basic structure of a linear map, dimensions orthogonal to the fundamental k dimensional subspaces are removed.

³An endomorphism is a mapping from a space back onto the same space.

⁴More explicitly, the operator $\mathbf{V}_\parallel \mathbf{V}_\parallel^T$ projects a vector onto the space $\text{span}(\mathbf{V}_\parallel)$, and $\mathbf{V}_\perp \mathbf{V}_\perp^T$ projects onto $\text{span}(\mathbf{V}_\perp)$, so we can always decompose a domain vector \mathbf{x} as $\mathbf{x} = (\mathbf{V}_\parallel \mathbf{V}_\parallel^T) \mathbf{x} +$

In other words, the mapping between the subspaces is bijective, and any components orthogonal to those spaces cannot be represented by the linear map. Thus, if $\mathbf{A} = \mathbf{U}_{\parallel}\mathbf{\Sigma}\mathbf{V}_{\parallel}^T$ is the forward map implementing a forward mapping between fundamental spaces and removing any domain element orthogonal to the column space then it reasons to say that the opposite mapping procedure which implements the inverse of that bijection between fundamental space and removes any component orthogonal to the row space is a natural generalized inverse for this map. This generalized inverse is given explicitly by the map

$$\mathbf{A}^{\dagger} = \mathbf{V}_{\parallel}\mathbf{\Sigma}^{-1}\mathbf{U}_{\parallel}^T. \quad (14)$$

It's straightforward to show that \mathbf{A}^{\dagger} gives the exact inverse on the bijection between the fundamental k -dimensional subspaces, but any component orthogonal to \mathbf{V}_{\parallel} is removed. This expression is exactly the Moore-Penrose pseudoinverse:

$$\mathbf{A}^{\dagger} = \mathbf{A}^T(\mathbf{A}\mathbf{A}^T)^{-1}, \quad (15)$$

as can be shown simply by expanding each $\mathbf{A} = \mathbf{U}_{\parallel}\mathbf{\Sigma}\mathbf{V}_{\parallel}^T$. In other words, by construction, we've shown that the Moore-Penrose pseudoinverse is a natural generalization of the notion of inverse. Prescriptively, it says that between the domain and range spaces the right thing to do is perform either the forward or inverse bijection between the fundamental k -dimensional subspaces, and simply remove any components orthogonal to those spaces. For the forward map, we remove any component orthogonal to the column space and perform the forward bijection. And for the inverse map, we remove any component orthogonal to the row space and perform the backward bijection.

2.3 Using the SVD to solve reduced rank linear equations

Let $\mathbf{A} = \mathbf{U}_{\parallel}\mathbf{\Sigma}\mathbf{V}_{\parallel}^T \in \mathbb{R}^{m \times n}$ be any rank k matrix as above. \mathbf{A} may be reduced rank, making k smaller than either m or n . Solving the equation

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

fully for the entire space of valid solutions is straightforward using this decomposition as we show here:

$$\begin{aligned} \mathbf{A}\mathbf{x} &= \mathbf{b} \\ \Rightarrow \mathbf{U}_{\parallel}\mathbf{\Sigma}\mathbf{V}_{\parallel}^T\mathbf{x} &= \mathbf{b} \\ \Rightarrow \mathbf{V}_{\parallel}^T\mathbf{x} &= \mathbf{\Sigma}^{-1}\mathbf{U}_{\parallel}^T\mathbf{b}. \end{aligned} \quad (16)$$

$(\mathbf{V}_{\perp}\mathbf{V}_{\perp}^T)\mathbf{x}$. Using that decomposition,

$$\begin{aligned} \mathbf{A}\mathbf{x} &= \mathbf{U}_{\parallel}\mathbf{\Sigma}\mathbf{V}_{\parallel}^T \left[(\mathbf{V}_{\parallel}\mathbf{V}_{\parallel}^T)\mathbf{x} + (\mathbf{V}_{\perp}\mathbf{V}_{\perp}^T)\mathbf{x} \right] \\ &= \mathbf{U}_{\parallel}\mathbf{\Sigma}\mathbf{V}_{\parallel}^T\mathbf{x}, \end{aligned}$$

since $\mathbf{V}_{\parallel}^T\mathbf{V}_{\perp} = \mathbf{0}$ by construction.

\mathbf{x} can always be decomposed as a linear combination of the columns \mathbf{V}_{\parallel} and a linear combination of the columns \mathbf{V}_{\perp} as in $\mathbf{x} = \mathbf{V}_{\parallel}\alpha + \mathbf{V}_{\perp}\beta$, where α and β are the associated coefficients. Doing so shows that the left hand side of the expression in Equation 16 reduces to

$$\begin{aligned}\mathbf{V}_{\parallel}^T \mathbf{x} &= \mathbf{V}_{\parallel}^T (\mathbf{V}_{\parallel}\alpha + \mathbf{V}_{\perp}\beta) \\ &= (\mathbf{V}_{\parallel}^T \mathbf{V}_{\parallel}) \alpha + (\mathbf{V}_{\parallel}^T \mathbf{V}_{\perp}) \beta \\ &= \alpha.\end{aligned}$$

In other words, the equation $\mathbf{A}\mathbf{x} = \mathbf{b}$ constrains \mathbf{x} only through α , which is given by

$$\alpha = \Sigma^{-1} \mathbf{U}_{\parallel}^T \mathbf{b}. \quad (17)$$

The full solution is, therefore,

$$\mathbf{x} = \mathbf{V}_{\parallel}\alpha + \mathbf{V}_{\perp}\beta = \underbrace{(\mathbf{V}_{\parallel}\Sigma^{-1}\mathbf{U}_{\parallel}^T)}_{\mathbf{A}^{\dagger}} \mathbf{b} + \mathbf{V}_{\perp}\beta \quad (18)$$

for any $(n - k)$ -dimensional coefficient vector β . As indicated, the first term is just the well known pseudoinverse solution, and the second term is an element of the null space spanned by \mathbf{V}_{\perp} .

This notation is somewhat different than commonly used to describe these solutions (raw matrix operations such as in Equation 15), but even computationally, expressions derived through applications of SVDs can be easily and robustly implemented, too, since matrix libraries make SVDs computations readily available in practice. The trade-off is robustness, stability, and geometric insight (SVD) for computational speed (raw matrix operations).

2.4 Solving the system without algebra

Consider the equation $\mathbf{A}\mathbf{x} = \mathbf{b}$, where $\mathbf{A} \in \mathbb{R}^{m \times n}$ can be any possibly reduced rank matrix. Then, as discussed above, the decomposition $\mathbf{A} = \mathbf{U}_{\parallel}\Sigma\mathbf{V}_{\parallel}^T$ tells us that \mathbf{A} transports specifically only the portion of \mathbf{x} lying within $\mathbf{span}(\mathbf{V}_{\parallel})$ to a point in the space $\mathbf{span}(\mathbf{U}_{\parallel})$. That means that all vectors of the form

$$\mathbf{z}(\beta) = \mathbf{x}_0 + \mathbf{V}_{\perp}\beta \quad (19)$$

are transported to the same point under \mathbf{A} . In particular, when \mathbf{x}_0 lies on the column space $\mathbf{span}(\mathbf{V}_{\parallel})$, i.e. $\mathbf{x}_0 = \mathbf{V}_{\parallel}\alpha$, then we know where \mathbf{x}_0 ends up under the fundamental internal bijection of \mathbf{A} and that means we know where all points $\mathbf{z}(\beta)$ end up.

Thus, to solve the system, we just need to know which point in $\mathbf{span}(\mathbf{V}_{\parallel})$ maps to $\mathbf{b} \in \mathbf{span}(\mathbf{U}_{\parallel})$. Since the forward map takes the coordinates of the point \mathbf{x} in the domain basis \mathbf{V}_{\parallel} , inflates them by the factors σ_i , and applies them directly to the range basis \mathbf{U}_{\parallel} , the inverse operation must just do the opposite. It

should take the coefficients in range basis \mathbf{U}_{\parallel} , shrink them by inverse factors $\frac{1}{\sigma_i}$, and apply them to domain basis \mathbf{V}_{\parallel} . Thus (as we derived above), if the forward map is $\mathbf{A} = \mathbf{U}_{\parallel}\mathbf{\Sigma}\mathbf{V}_{\parallel}^T$, the inverse map between fundamental spaces must be $\mathbf{A}^\dagger = \mathbf{V}_{\parallel}\mathbf{\Sigma}^{-1}\mathbf{U}_{\parallel}^T$ (specifically, \mathbf{U}_{\parallel}^T finds the components in the k -dimensional orthogonal basis for $\mathbf{span}(\mathbf{U}_{\parallel})$, $\mathbf{\Sigma}^{-1}$ scales those components, and \mathbf{V}_{\parallel} applies the resulting components directly to the k orthogonal basis elements of $\mathbf{span}(\mathbf{V}_{\parallel})$).

Given this intuition, we're now equipped to simply write down the full linear space solution to the system directly:

$$\mathbf{x}^* = \left(\mathbf{V}_{\parallel}\mathbf{\Sigma}^{-1}\mathbf{U}_{\parallel}^T\right)\mathbf{b} + \mathbf{V}_{\perp}\beta \quad (20)$$

for any $\beta \in \mathbb{R}^k$.

Note that this argument can also be used to solve, in the least squares sense, the system when \mathbf{b} doesn't actually lie within $\mathbf{span}(\mathbf{U}_{\parallel})$. In this case, we just need to find which element of $\mathbf{span}(\mathbf{V}_{\parallel})$ gets mapped onto the orthogonal projection of \mathbf{b} onto $\mathbf{span}(\mathbf{U}_{\parallel})$, which is just the point defined by the coordinates of \mathbf{b} in basis \mathbf{U}_{\parallel} . Considering only these coordinates, the above argument still unfolds in the same way and the solution doesn't change. In other words, we can use a geometric argument to show that that Equation 20 also solves the least squares problem.

3 Quadratic forms and their manipulation

Quadratic forms are important tools in optimization and control because of their close connection to linear systems. Above, for instance, we saw that the basic structure of a linear map encodes an implicit least squares problem, which itself may be viewed as a quadratic objective function.

This section reviews very briefly some of the basic rules governing quadratic forms that should ideally be understood at an intuitive level. These rules relate how quadratics combine with one another and how they behave under linear transformation. We'll state the rules here without proof. They can all be derived by calculation.

1. **Adding quadratics always gives another quadratic.** $Q_1(\mathbf{x}) + Q_2(\mathbf{x})$ is a quadratic over \mathbf{x} , and $Q_x(\mathbf{x}) + Q_q(\mathbf{q})$ is a quadratic function defined jointly over \mathbf{x} and \mathbf{q} . Generally, adding any two quadratic functions of any collection of variables gives a joint quadratic over the union of those variables. Calculating the coefficients can be tedious, but we always know that it'll still be a quadratic.
2. **Linearly transforming a quadratic gives another quadratic.** If $Q(\mathbf{x})$ is a quadratic and $\mathbf{x} = \mathbf{A}\mathbf{q} + \mathbf{B}\mathbf{u}$ is a linear transformation of \mathbf{q} and \mathbf{u} into \mathbf{x} , then the composition $Q(\mathbf{A}\mathbf{q} + \mathbf{B}\mathbf{u}) = \tilde{Q}(\mathbf{q}, \mathbf{u})$ is a quadratic function over \mathbf{q} and \mathbf{u} .

3. **Conditioning: Fixing one variable gives a quadratic over the others.** If $Q(\mathbf{x}, \mathbf{u})$ is a joint quadratic defined over both \mathbf{x} and \mathbf{u} , fixing one of the variables, say \mathbf{x} at a particular value \mathbf{x}_t , gives a quadratic function $Q(\mathbf{x}_t, \mathbf{u}) = \tilde{Q}_{\mathbf{x}_t}(\mathbf{u})$ over the remaining variables \mathbf{u} .
4. **Marginalization: Optimizing over a subset of variables gives a quadratic over the rest.** Let $Q(\mathbf{x}, \mathbf{u})$ be a quadratic over both \mathbf{x} and \mathbf{u} . We can solve analytically for the optimal \mathbf{u} as a function of \mathbf{x} to get an expression $\mathbf{u}^*(\mathbf{x})$ that tells us the optimal setting to the variable \mathbf{u} given the value \mathbf{x} . Plugging that expression back into the quadratic removes \mathbf{u} (since we're optimizing over it) and gives $\tilde{Q}(\mathbf{x}) = Q(\mathbf{x}, \mathbf{u}^*(\mathbf{x}))$. This new function $\tilde{Q}(\mathbf{x})$ is also a quadratic.

It's important to be able to perform the calculations to explicitly derive the above results, but it's more important to understand these results at an intuitive level since their calculation can be tedious and understanding them alone can give a lot of theoretical insight into optimization and control problems. In particular, the last two properties, denoted as *conditioning* and *marginalization* in reference to their corresponding Gaussian analogs (Gaussian inference and quadratic optimization are very closely related), are of central importance to optimal control.