

# A Geometric Characterization of Solutions to The Algebraic Riccati Equation

## 1 Introduction

The Algebraic Riccati Equation (ARE) is the following matrix equation for the  $n \times n$  self-adjoint matrix  $P$

$$A^*P + PA + PDP + C = 0, \tag{1}$$

where  $A$ ,  $C$ , and  $D$  are  $n \times n$  matrices (in  $\mathbb{R}$  or  $\mathbb{C}$ ), with  $C = C^*$  and  $D = D^*$ . Generally, the question of solution existence and uniqueness is nontrivial. For instance, in the case of  $1 \times 1$  matrices, (1) becomes the familiar quadratic equation  $dp^2 + 2ap + c = 0$ , which may have distinct solutions, a unique solution, or no solution (when  $d$ ,  $a$ ,  $c$  and  $p$  are real). Although our interest is in the self-adjoint solutions of (1), non self-adjoint solutions to (1) may also exist, for instance

$$P = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \text{ solves (1) with } A = \begin{bmatrix} a & 0 \\ 0 & -a \end{bmatrix}, D = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}, \text{ and } C = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}. \tag{2}$$

In the case of non self-adjoint solutions, we note that if  $P$  solves (1) then so does its adjoint  $P^*$ .

Research on the ARE and its variants has been extremely active for over half a century, see [7] for references. Our work begins in Section 2 with a description of the Linear Quadratic Regulation (LQR) problem from control theory, which is one of the principle motivations for the ARE. In this context there are well-known results on the existence and uniqueness of certain ARE solutions. Essentially, the optimal control for any reasonable physical system corresponds to the solution of an ARE which exists, is unique, and can be computed easily using a variety of methods. Equation (1) however is much deeper than this particular application. Our main goal in this project is to describe solutions of (1) using the geometry of Lagrangian subspaces and symplectic forms. The connection between ARE solutions and Lagrangian subspaces is well known (see [5] and references given there). The geometric characterization of the ARE solutions that we present is related to the classical work on the symplectic classification of quadratic forms- see [6, §21.5] and references given there. We undertake our main task in Section 3, where we characterize the real symmetric solutions of (1) in the case where  $A$ ,  $C$ , and  $D$  are real. As part of this characterization, we obtain sufficient conditions for the existence of these solutions. Core objects in Section 3 are established using the natural and aesthetic language of dual spaces. This material is unfamiliar to many engineers, and so in Section 4 we establish these core objects using familiar but less elegant tools like inner products. In low enough dimensions, this development can be visualized, as we show in Section 4.1. After showing in Section 5 that our geometric characterization from Section 3 leads to a method of computing ARE solutions in Matlab, we end in Section 6 with concluding remarks and a discussion of ways in which the current work can be extended.

A substantial portion of our work is contained in the Appendices. In Appendix A, we derive conditions that must be satisfied by a solution to the LQR problem. In Appendix B we define terms used in Sections 2 and 3, and we discuss several uses of the word *adjoint*. In Appendix C, we prove the many assertions from Section 3, and in Appendix D, we discuss real objects in vector spaces over  $\mathbb{C}$ , as well as the complexification of vector spaces over  $\mathbb{R}$ .

## 2 Linear Quadratic Regulation

In this section we describe the Linear Quadratic Regulation (LQR) problem and its connection to the ARE. The LQR problem is to find a control function  $u(t)$  for the system

$$\dot{x} = Ax + Bu, \quad x(0) = x_0 \tag{3}$$

which minimizes the cost function

$$J = \int_0^T (x^T Q x + u^T R u) dt + \hat{x}^T \hat{Q} \hat{x}, \quad (4)$$

where  $Q = Q^T \geq 0$ ,  $\hat{Q} = \hat{Q}^T \geq 0$ , and  $R = R^T > 0$ , (and where  $\hat{x}$  denotes  $x(T)$ ). We refer to a  $u(t)$  that minimizes  $J$  as *optimal*. The matrices  $Q$  and  $R$  penalize nonzero state values  $x(t)$  and control values  $u(t)$  respectively for  $t \in [0, T]$ , while the matrix  $\hat{Q}$  penalizes nonzero state values  $\hat{x}$  at the final time  $T$ .

In Appendix A, we show that if a smooth  $u(t)$  minimizes  $J$ , then it satisfies

$$u(t) = -R^{-1} B^T P(t) x(t), \quad (5)$$

where  $P(t) = P^T(t) \geq 0$  satisfies what we will call the *Riccati initial value problem*, consisting of the differential equation

$$-\dot{P} = A^T P + P A - P B R^{-1} B^T P + Q, \quad (6)$$

and the starting condition  $P(T) = \hat{Q}$ , with integration *backwards* in time, from  $T$  to  $t < T$ . The LQR problem therefore is associated with the ARE given by

$$0 = A^T P + P A - P B R^{-1} B^T P + Q. \quad (7)$$

We note that solutions of (7) are fixed points in the flow induced by (6). Collecting results from the excellent book [7] by Lancaster and Rodman, we obtain the following

**Theorem 2.1** *The Riccati initial value problem with  $\hat{Q} = 0$  has a well-defined solution for every  $t < T$ . If  $(A, B)$  is stabilizable<sup>1</sup>, then  $P(t) \rightarrow P$  as  $T \rightarrow \infty$ , where  $P$  is a positive semi-definite solution of (7). If in addition  $(Q, A)$  is observable, then this  $P$  is positive definite, and is the only matrix which satisfies (7) among all positive semi-definite matrices.*

The positive definite  $P$  in Theorem 2.1 is particularly important because it plays a role in the *infinite horizon* LQR problem, in which the cost function is given by

$$J_\infty = \int_0^\infty (x^T Q x + u^T R u) dt. \quad (8)$$

If all the hypotheses of Theorem 2.1 are satisfied, then the infinite horizon problem and its solution can be obtained from the finite time problem by taking the limit as  $T \rightarrow \infty$ . In particular, because  $P(t)$  at any fixed time  $t$  approaches the constant positive definite  $P$  from Theorem 2.1, the optimal  $u$  at this fixed time must satisfy

$$u(t) = -R^{-1} B^T P x(t). \quad (9)$$

The corresponding minimum value of  $J_\infty$  is given by  $x_0^T P x_0$ . Also, we note that using (9) as a feedback rule for  $u(t)$  in  $\dot{x} = Ax + Bu$  causes the closed loop system  $\dot{x} = (A - BR^{-1}B^T P)x$  to be stable, (i.e., the eigenvalues of  $A - BR^{-1}B^T P$  are in the open left half plane).

In addition to an existence and uniqueness result for solutions to the ARE given in (7), the LQR problem suggests a practical way of computing solutions; simply start at 0 and integrate (6) until  $P(t)$  stops changing. With this in mind, we note that with  $Y = PX$ , the nonlinear (6) has the following *linear* equivalent

$$\frac{d}{dt} \begin{bmatrix} X \\ Y \end{bmatrix} = \begin{bmatrix} A & -BR^{-1}B^T \\ -Q & -A^T \end{bmatrix} \begin{bmatrix} X \\ Y \end{bmatrix} \quad (10)$$

## 2.1 The One Dimensional Setting

Here we illustrate the LQR solution and some of its possible degeneracies in the easy to understand one dimensional setting. To emphasize that that all objects are now scalars, we write (3) as

$$\dot{x} = ax + bu, \quad x(0) = x_0, \quad (11)$$

---

<sup>1</sup>The term stabilizable is defined along with its relatives in Appendix B.1.

and we write  $J$  as

$$J = \int_0^T (qx^2 + ru^2)dt + \hat{q}\hat{x}^2, \quad (12)$$

where  $\hat{x} = x(T)$ . From Appendix A we know that if a smooth optimal  $u(t)$  exists, then it must satisfy

$$u(t) = -\frac{b}{r}p(t)x(t), \quad (13)$$

where  $p(t)$  satisfies the Riccati initial value problem given by  $p(T) = \hat{q}$ , and the differential equation

$$-\dot{p} = q + 2ap - \frac{b^2}{r}p^2 =: f(p). \quad (14)$$

Our ultimate interest is in the fixed points of the flow induced by (14), as these are the solutions of the ARE given by  $f(p) = 0$ . This flow can be visualized by plotting  $f(p)$  as in Figure 1. We note that when there are no fixed points (i.e., when  $ra^2 + 4qb^2 < 0$ ), every trajectory that satisfies (14) experiences a finite time singularity. If the trajectory through  $p(T) = \hat{q}$  experiences its singularity over the interval on which an optimal control is desired, then of course this optimal control is an impossibility.

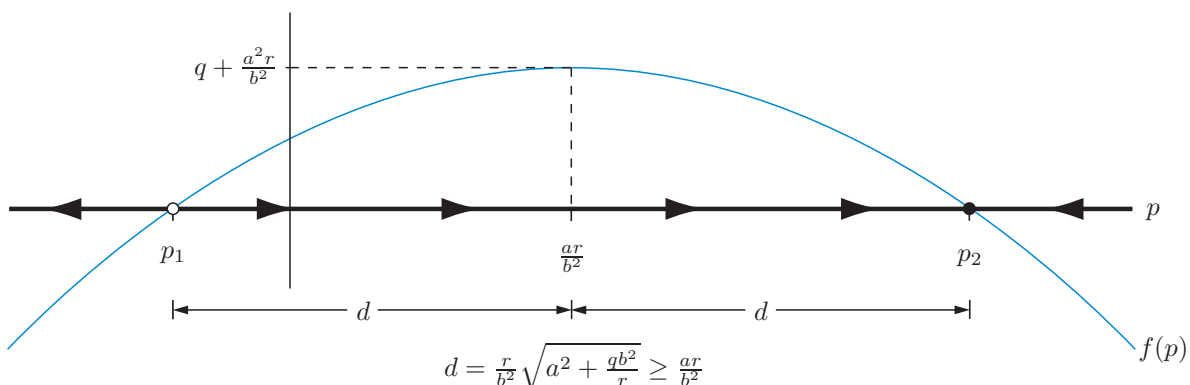


Figure 1: Here we graph the function  $f(p)$  which governs the evolution of  $p(t)$  according to  $-\dot{p} = f(p)$ . As  $t$  decreases, points  $p$  on the  $x$ -axis move as indicated by the large black arrows. We consider here the case  $b \neq 0$ ,  $r > 0$ , and  $q \geq 0$ .

When the conditions  $q \geq 0$  and  $r > 0$  given with the LQR problem are satisfied, the first segment of Theorem 2.1 assures us that the solution  $p(t)$  to the Riccati initial value problem with  $\hat{q} = 0$  is well-defined for all  $t < T$ . The second segment of Theorem 2.1 guarantees that the finite limit of  $p(t)$  is nonnegative if the pair  $(a, b)$  is stabilizable. In one dimension, the pair  $(a, b)$  is stabilizable if and only if either  $b \neq 0$  (i.e., we can control  $x$ ), or  $a < 0$  (i.e.,  $x$  goes to 0 on its own). By examining plots of  $f(p)$ , it is easy to verify that these conditions cause (14) to have a stable nonnegative fixed point. The third segment of Theorem 2.1 requires the pair  $(q, a)$  to be observable; in one dimension this occurs if and only if  $q \neq 0$ .

In the LQR problem,  $q$  and  $\hat{q}$  are nonnegative and  $r$  is positive. As an example of what can happen without these restrictions, note that when  $q = 0$ ,  $\hat{q} = 0$ , and  $r < 0$ , the resulting set of possible  $J$  values is not bounded below, (consider for instance the family of control functions  $\{u_k(t)\}_{k=1}^{\infty}$  where  $u_k(t) = k$  for all  $t$ ). Difficulties can also arise when  $r = 0$ . For instance with  $q = 1$ ,  $\hat{q} = 0$ ,  $r = 0$ , and  $x_0 = 1$ , the resulting set of possible  $J$  values is given by  $(0, \infty)$ . This set is bounded below, but it has no minimum and so no optimal  $u(t)$  exists.

### 3 A Geometric Characterization

In this section, we consider Algebraic Riccati Equations (1) in which the matrices  $A$ ,  $C$ , and  $D$  are real,

$$A^T P + PA + PCP + D = 0, \quad A, C, D \in \mathbb{R}^{n \times n}, \quad C = C^T, \quad D = D^T. \quad (15)$$

The solution matrices  $P$  can be either real or complex; our interest here is with matrices  $P = P^T$  that are real. Our first step is to move from matrices to mappings on abstract vector spaces. Let  $V$  be an

$n$ -dimensional vector space over  $\mathbb{C}$ , and let  $V'$  be its dual<sup>2</sup>. With respect to a basis  $\{e_i\}$  on  $V$  and its dual basis  $\{e'_i\}$  on  $V'$ , the matrices in (15) define the following mappings

$$A : V \longrightarrow V, \quad C : V' \longrightarrow V, \quad D : V \longrightarrow V', \quad \text{and} \quad P : V \longrightarrow V'. \quad (16)$$

We use  $\{e_i\}$  and  $\{e'_i\}$  to define real vectors on  $V$  and  $V'$  respectively (see Appendix D), and so the mappings in (16) are real. The same symbol denotes a matrix and its associated mapping, with the meaning of a symbol always clear from its context. As we discuss in Appendix B.3, the matrix corresponding to the natural adjoint  $P'$  of the mapping  $P : V \longrightarrow V'$  is the transpose (no conjugate) of the matrix corresponding to  $P$ , even when the matrix entries are complex. The natural adjoints of the mappings in (16) are given by

$$\begin{aligned} A' : V' \longrightarrow V & \quad \text{such that } \langle A'\xi, x \rangle_V = \langle \xi, Ax \rangle_V & \quad \text{for all } x \in V, \xi \in V', \\ C' : V' \longrightarrow V & \quad \text{such that } \langle \xi, C'\eta \rangle_V = \langle \eta, C\xi \rangle_V & \quad \text{for all } \xi, \eta \in V', \\ D' : V \longrightarrow V' & \quad \text{such that } \langle D'x, y \rangle_V = \langle Dy, x \rangle_V & \quad \text{for all } x, y \in V, \\ P' : V \longrightarrow V' & \quad \text{such that } \langle P'x, y \rangle_V = \langle Py, x \rangle_V & \quad \text{for all } x, y \in V, \end{aligned} \quad (17)$$

where  $\langle \xi, x \rangle_V$  denotes  $\xi(x)$ , with  $\xi \in V'$  and  $x \in V$ . We define a new vector space  $W = V \oplus V'$ , with elements written as  $\begin{bmatrix} x \\ \xi \end{bmatrix}^T$ , where  $x \in V$  and  $\xi \in V'$ , and we use the union of real bases on  $V$  and  $V'$  to define real vectors on  $W$ . We turn  $W$  into a *symplectic vector space* by pairing it with a *symplectic bilinear form*  $\sigma$  on  $W$ , given by

$$\sigma\left(\begin{bmatrix} x \\ \xi \end{bmatrix}, \begin{bmatrix} y \\ \eta \end{bmatrix}\right) = \langle \xi, y \rangle_V - \langle \eta, x \rangle_V. \quad (18)$$

The basic properties of  $\sigma$  are

$$\sigma(X, Y) = -\sigma(Y, X), \quad \text{and} \quad \sigma(X, Y) = 0 \quad \forall Y \in W \implies X = 0, \quad (19)$$

and in fact we show in C.1 that any symplectic bilinear form on an even dimensional vector space can be expressed as (18). Let  $Q(X, Y)$  be the following *symmetric bilinear form* on  $W$

$$Q\left(\begin{bmatrix} x \\ \xi \end{bmatrix}, \begin{bmatrix} y \\ \eta \end{bmatrix}\right) = \langle Dx, y \rangle_V + \langle \eta, Ax \rangle_V + \langle \xi, Ay \rangle_V + \langle \eta, C\xi \rangle_V. \quad (20)$$

We now show that matrix solutions  $P = P^T$  of (15) correspond to Lagrangian graph subspaces of  $W$  over which  $Q = 0$ .

- A subspace  $\Lambda \subset W$  is called *Lagrangian* if it is equal to its *symplectic complement* (given by the set of vectors in  $W$  that are  $\sigma$ -orthogonal to  $\Lambda$ ).
- A subspace  $\Lambda \subset W$  is called a *graph space* if it can be written as

$$\Lambda = \left\{ \begin{bmatrix} x \\ Px \end{bmatrix} : x \in V \right\}, \quad (21)$$

where  $P \in L(V', V)$ . With  $\pi : W \longrightarrow V$  the natural projection defined by  $\pi(\begin{bmatrix} x \\ y \end{bmatrix}^T) = x$ , we show in C.2 that  $\Lambda \subset W$  can be written as (21) if and only if  $\pi|_{\Lambda}$  is a bijection.

In C.3, we show that the graph space  $\Lambda \subset W$  is Lagrangian if and only if  $P$  from the graph space representation is equal to its natural adjoint. Now suppose that  $Q$  vanishes identically on a Lagrangian graph subspace of  $W$ . Then, using (20),

$$\begin{aligned} 0 &= Q\left(\begin{bmatrix} x \\ Px \end{bmatrix}, \begin{bmatrix} y \\ Py \end{bmatrix}\right) = \langle Dx, y \rangle_V + \langle Py, Ax \rangle_V + \langle Px, Ay \rangle_V + \langle Py, CPx \rangle_V \\ &= \langle (PA + A'P + D + PCP)x, y \rangle_V, \end{aligned} \quad (22)$$

for all  $x, y \in V$ , and so (15) holds<sup>3</sup>. Conversely, if  $P = P^T$  is a solution to (15), then it corresponds to a Lagrangian graph subspace of  $W$  over which  $Q$  vanishes. These findings are illustrated in Figure 2.

<sup>2</sup>As discussed in Section 4, the following development works just as well (but with less elegance) if  $V'$  is replaced by  $V$ , and if the natural pairing between vectors and functionals is replaced by an inner product on  $V$ .

<sup>3</sup>In detail,  $x$  in (22)<sub>2</sub> is acted on by a linear *map* from  $V$  to  $V'$ , which we have argued is equal to zero. The corresponding equation in *matrices* is (15), and so the matrix corresponding to the map  $P$  satisfies the ARE. Because  $P = P'$ , the matrix is equal to its transpose.

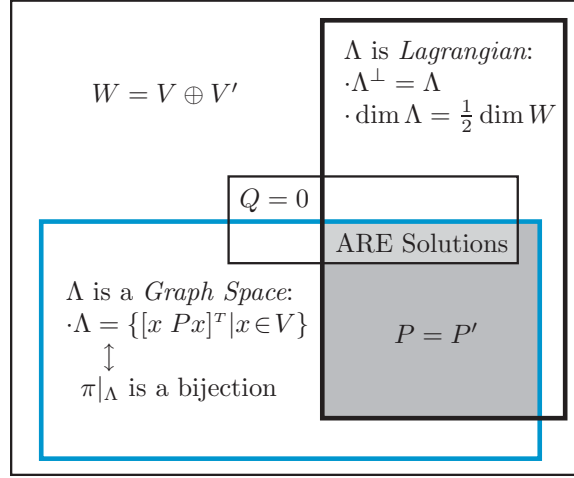


Figure 2: This diagram summarizes the subspaces of  $W$  that we use to characterize solutions to (15). Graph subspaces are Lagrangian if and only if  $P = P'$ . Lagrangian graph subspaces on which  $Q = 0$  correspond to the desired solutions.

At this point we have characterized the symmetric matrix solutions of (15). To characterize the *real* symmetric matrix solutions of (15), we must consider the structure of  $W$  in further detail. Let  $q(X) = Q(X, X)$  be the *quadratic form* associated with  $Q$ ,

$$q\left(\begin{bmatrix} x \\ \xi \end{bmatrix}\right) = \langle Dx, x \rangle_V + \langle \xi, Ax \rangle_V + \langle \xi, Ax \rangle_V + \langle \xi, C\xi \rangle_V. \quad (23)$$

The *Hamiltonian*<sup>4</sup>  $F$  of  $q$  is a real operator on  $W$  defined by

$$F = \frac{1}{2} \begin{bmatrix} p''_{\xi x} & p''_{\xi \xi} \\ p''_{xx} & p''_{x\xi} \end{bmatrix} = \begin{bmatrix} A & C \\ -D & -A' \end{bmatrix}. \quad (24)$$

In C.4 we show that  $F$  can also be defined by the following (coordinate invariant) equation,

$$\sigma(X, FY) = Q(X, Y), \quad X, Y \in W. \quad (25)$$

It follows (see C.7) that a Lagrangian subspace on which  $Q$  vanishes is equivalent to one which is invariant under  $F$ , and that solutions to (15) correspond to  $F$ -invariant Lagrangian graph subspaces of  $W$ . We let  $V_\lambda \subset W$  denote the space of generalized eigenvectors of  $F$  belonging to the eigenvalue  $\lambda \in \text{Spec}(F)$ . In C.8, we establish that

$$\lambda_1 + \lambda_2 \neq 0 \implies \sigma(V_{\lambda_1}, V_{\lambda_2}) = 0, \quad (26)$$

from which it follows (see C.9) that

$$\lambda \in \text{Spec}(F) \implies -\lambda, \bar{\lambda}, -\bar{\lambda} \in \text{Spec}(F), \quad (27)$$

and also (see C.10) that for  $\lambda \neq 0$ ,  $V_\lambda$  is the dual space of  $V_{-\lambda}$  according to

$$V_\lambda \stackrel{\cong}{\simeq} V'_{-\lambda}, \quad \alpha(X)(Y) = \sigma(X, Y), \quad X \in V_\lambda, \quad Y \in V_{-\lambda}. \quad (28)$$

It follows (see C.11) that  $V_\lambda \oplus V_{-\lambda}$  is a symplectic vector space, and (see C.13) that  $V_0$  is symplectic as well.

In C.15, we show that a Lagrangian subspace  $\Lambda \subset W$  is invariant under  $F$  if and only if

$$\Lambda = \left( \bigoplus_{\lambda \in \Sigma} \Lambda_\lambda \right) \oplus \Lambda_0, \quad (29)$$

<sup>4</sup>An alternate terminology is to call  $q$  the Hamiltonian of the vector field  $F$  on  $W$ .

where each  $\Lambda_\lambda$  is an  $F$ -invariant Lagrangian subspace of  $V_\lambda \oplus V_{-\lambda}$ , where  $\Lambda_0$  is an  $F$ -invariant Lagrangian subspace of  $V_0$ , and where  $\Sigma$  is some set satisfying  $\Sigma \cup (-\Sigma) = \text{Spec}(F) \setminus \{0\}$  and  $\Sigma \cap (-\Sigma) = \emptyset$ . The set  $\Sigma$  causes (29) to include an  $F$ -invariant Lagrangian subspace of each  $V_\lambda \oplus V_{-\lambda}$ , (we note that  $V_\lambda \oplus V_{-\lambda}$  is the same as  $V_{-\lambda} \oplus V_\lambda$ ). It is the specification of these subspaces that affects the resulting  $\Lambda$ .

We now consider the special case in which

$$\text{Spec}(F) \cap i\mathbb{R} = \emptyset. \quad (30)$$

Let  $\lambda_1, \dots, \lambda_m$  be the real eigenvalues of  $F$  with positive real parts, and let  $\mu_1, \dots, \mu_n$  be the eigenvalues of  $F$  in the first quadrant of  $\mathbb{C}$ , (i.e., with positive real and imaginary parts). Note that  $W$  can be decomposed as the direct sum of  $W_1, W_2$ , and  $W_3$ , where

$$W_1 = \bigoplus_{i=1}^m (V_{\lambda_i} \oplus V_{-\lambda_i}), \quad W_2 = \bigoplus_{i=1}^n (V_{\mu_i} \oplus V_{-\mu_i}), \quad W_3 = \bigoplus_{i=1}^n (V_{\bar{\mu}_i} \oplus V_{-\bar{\mu}_i}). \quad (31)$$

To construct  $\Lambda$  as in (29), we need to choose  $F$ -invariant Lagrangian subspaces of the terms in parenthesis. Two obvious such subspaces of  $V_\alpha \oplus V_{-\alpha}$  (with  $\alpha$  equal to  $\lambda_i, \mu_i$ , or  $\bar{\mu}_i$ ) are  $V_\alpha$  and  $V_{-\alpha}$  (see C.16 for details). The possible combinations of these subspaces give  $2^{m+2n}$  different  $F$ -invariant Lagrangian subspaces  $\Lambda$  of  $W$ . Several of these lead to the desired *real* solutions of (15). In particular, if we choose  $V_{\mu_i}$  and  $V_{\bar{\mu}_i}$  (or  $V_{-\mu_i}$  and  $V_{-\bar{\mu}_i}$ ) as the subspaces of  $V_{\mu_i} \oplus V_{-\mu_i}$  and  $V_{\bar{\mu}_i} \oplus V_{-\bar{\mu}_i}$  respectively, then when these are combined in the direct sum (29), the resulting space is spanned by real vectors in  $W$  (see C.17). Of course,  $V_{\pm\lambda_i}$  is spanned by real vectors as well. It follows that  $\Lambda$  is a *real*  $F$ -invariant Lagrangian subspace of  $W$ .

In order for the  $F$ -invariant Lagrangian  $\Lambda$  to correspond to a symmetric matrix solution of (15),  $\Lambda$  must also be a graph space. We now claim that  $\pi|_\Lambda$  is a bijection (recall that this implies  $\Lambda$  is a graph space) if  $\langle \xi, C\xi \rangle_V$  is non-degenerate, (i.e., if  $\langle \xi, C\xi \rangle_V$  only if  $\xi = 0$ ). The domain  $\Lambda$  and target space  $V$  of  $\pi|_\Lambda$  have the same dimension, and so bijectivity is equivalent to injectivity. We establish injectivity by showing that  $\ker(\pi|_\Lambda) = \{[0 \ 0]^T\}$ . If  $[x \ \xi]^T \in \Lambda$  gets mapped to 0 by  $\pi|_\Lambda$ , then clearly  $x = 0$ . Next, because the  $F$ -invariance of  $\Lambda$  is equivalent to the vanishing of  $Q$  (and hence  $q$ ) on  $\Lambda$ , we have

$$q\left(\begin{bmatrix} 0 \\ \xi \end{bmatrix}\right) = \langle \xi, C\xi \rangle_V = 0. \quad (32)$$

But  $\langle \xi, C\xi \rangle_V$  is non-degenerate and so  $\xi = 0$  as desired. Thus the  $F$ -invariant Lagrangian  $\Lambda$  is a graph space as desired, and can be written as  $\Lambda = \{[v \ Pv]^T | v \in V\}$ . The mapping  $P$  in this representation satisfies  $P = P'$ , and the associated symmetric matrix solves (15).

We now argue that the map  $P : V \rightarrow V'$  found above is real. It then follows that the matrix of  $P$  with respect to the bases  $\{e_i\}$  and  $\{e'_i\}$  (introduced at the beginning of this section) consists of real numbers. Real vectors in  $W = V \oplus V'$  are combinations over  $\mathbb{R}$  of the following real basis on  $W$

$$\left\{ \begin{bmatrix} e_1 \\ 0 \end{bmatrix}, \begin{bmatrix} e_2 \\ 0 \end{bmatrix}, \dots, \begin{bmatrix} e_n \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ e'_1 \end{bmatrix}, \begin{bmatrix} 0 \\ e'_2 \end{bmatrix}, \dots, \begin{bmatrix} 0 \\ e'_n \end{bmatrix} \right\}. \quad (33)$$

We showed that  $\Lambda$  is spanned over  $\mathbb{C}$  by real vectors in  $W$ , and so it follows that  $\Lambda$  has a real basis  $\{[b_i \ c_i]^T\}$ , each element of which is an  $\mathbb{R}$  combination of the elements in (33). We also showed that

$$\Lambda = \left\{ \begin{bmatrix} v \\ Pv \end{bmatrix} \mid v \in V \right\}, \quad (34)$$

and so it must be that the  $\mathbb{C}$ -span of the  $b_i$ 's is  $V$ , (in (34) we can pick  $v$  to be any element of  $V$ ). Pick any *real*  $u \in V$ , that is, pick  $u = \sum \gamma_i b_i$  where every  $\gamma_i$  is real. Note that  $\sum \gamma_i [b_i \ c_i]^T$  is in  $\Lambda$ . We know from (34) that every element of  $\Lambda$  is given by  $[v \ Pv]^T$  where  $v \in V$ , and so it must be that  $Pu = \sum \gamma_i c_i$ , which is a *real* vector in  $V'$ . Thus  $P$  is a real map as desired.

## 4 An Alternative Development

Here we consider an alternative method of establishing the core objects  $W, \sigma$ , and  $Q$  from Section 3. Rather than have  $V'$  be the dual of  $V$ , we set  $V'$  equal to  $V$ , so that  $W$  becomes  $V \oplus V$ . In this context,  $\langle \bullet, \bullet \rangle_V$  now

stands for an inner product on  $V$ , instead of a pairing between functionals and vectors. This inner product is an addition to the discourse beyond what is needed using the natural development from Section 3, however this aesthetic disadvantage is balanced by the benefit of working on familiar ground (engineers are generally more at home with inner products than with dual spaces). Because the matrices in (15) are real, the formal development from the first part of Section 3 is unaffected by this change. We now undertake this approach in one dimension, where the resulting constructions are easy to visualize. Spatial intuition in this simple case is a valuable guide for similar constructions in higher dimensions.

#### 4.1 A One Dimensional Example

Here we consider the simple but nontrivial one dimensional ARE, which is the following scalar quadratic equation in  $x$

$$cx^2 + 2ax + d = 0. \quad (35)$$

The low dimensionality of this equation makes it possible to visualize the associated geometric constructions from Section 3, for instance  $V$  is simply  $\mathbb{R}$  and  $W = V \oplus V$  is simply  $\mathbb{R}^2$ . With an inner product on  $V = \mathbb{R}$  given by scalar multiplication, the symplectic bilinear form  $\sigma$  on  $W$  defined by (18) becomes

$$\sigma\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}\right) = x_2y_1 - y_2x_1 = \det\left(\begin{bmatrix} y_1 & x_1 \\ y_2 & x_2 \end{bmatrix}\right). \quad (36)$$

Therefore  $\sigma(u, v)$  returns the (signed) area of the parallelogram associated with  $u$  and  $v$ .

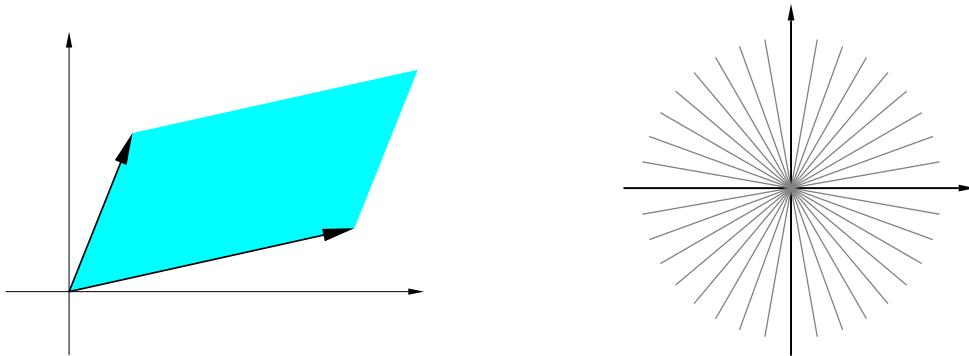


Figure 3: The symplectic bilinear form  $\sigma$  returns the (signed) area of the parallelogram associated with its input vectors, and the Lagrangian subspaces of  $W = \mathbb{R}^2$  are lines in  $\mathbb{R}^2$  through the origin.

A Lagrangian subspace  $\Lambda$  of  $W = \mathbb{R}^2$  must satisfy  $\dim \Lambda = \frac{1}{2} \dim W = 1$  (which requires  $\Lambda$  to be a line through the origin). The additional requirement that every vector in  $\Lambda$  be  $\sigma$ -orthogonal to every other vector in  $\Lambda$  is trivially satisfied in this case; if  $u$  and  $v$  correspond to points on a line through the origin, then  $\sigma(u, v) = 0$  because the parallelogram associated with  $u$  and  $v$  has zero area.

Graph subspaces of  $W = \mathbb{R}^2$  can be written as  $\{[\alpha \ x\alpha]^T | \alpha \in \mathbb{R}\}$ , and therefore consist of non-vertical lines through the origin. We have already established that lines through the origin correspond to the Lagrangian subspaces of  $W$ , and so the Lagrangian graph subspaces of  $W$  are simply the graph subspaces of  $W$ .

The major geometric characterization from Section 3 is that ARE solutions correspond to Lagrangian graph subspaces of  $W$  on which the symmetric bilinear form  $Q$  vanishes. In our one dimensional example,  $Q$  is given by

$$Q\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}\right) = \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} d & a \\ a & c \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \quad (37)$$

Lagrangian graph subspaces in our problem are non-vertical lines through the origin, and so our interest is in  $Q$  restricted to these lines. In particular, we are interested in finding vectors  $u = [x_1 \ x_2]^T$  with  $x_1 \neq 0$  such that  $Q(\alpha u, \beta u) = \alpha\beta Q(u, u) = 0$  for all  $\alpha, \beta \in \mathbb{R}$ . Note that  $Q(u, u)$  defines the quadratic form  $q(u)$ , which can be visualized (see Figure 4) as a surface above  $W = \mathbb{R}^2$ . The claim in Section 3 is that solutions

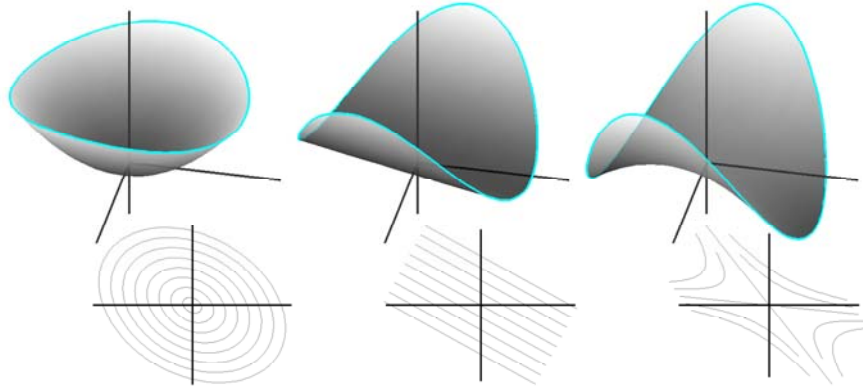


Figure 4: Here we illustrate three possible types of quadratic form  $p$  as surfaces above  $W$ . Each surface is accompanied by its corresponding level set. Non-vertical lines in the level set which includes the origin correspond to real solutions of (35). In the left most image there are no real solutions, in the middle image there is one real solution, and in the image at right there are two real solutions.

to the quadratic equation correspond to non-vertical lines through the origin along which  $q$  equals zero. The set in  $W$  over which  $q = 0$  is simply the level set of  $q$  which includes the origin, (see Figure 4 for an example of these level sets and their corresponding surfaces). This correspondence is easy to verify analytically; any non-vertical line in  $W$  can be written as the span of a vector of the form  $[1 \ x]^T$ , and

$$q\left(\begin{bmatrix} 1 \\ x \end{bmatrix}\right) = 0 \iff cx^2 + 2ax + d = 0. \quad (38)$$

## 5 Computing Solutions

Here we demonstrate the development in Section 3 by constructing solutions to an ARE in Matlab. Suppose our interest is with solutions to (15), where  $A$ ,  $C$ , and  $D$  are given by

$$A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}, \quad C = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}, \quad D = \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}. \quad (39)$$

We start by entering these arrays into Matlab, and by concatenating them to create  $F$

```
A=[1 2;3 4];
C=[2 0;0 1];
D=[1 2;2 1];
F=[A C;-D -A'];
```

Next we find the eigenvalues and eigenvectors of  $F$

```
[evecs,evals]=eig(F);
evals=diag(evals);
```

There are four distinct real eigenvalues and corresponding eigenvectors. Following the development from the last paragraph in Section 3, the span of any two of the eigenvectors is  $F$ -invariant Lagrangian subspace  $\Lambda$  (and thus a prospective ARE solution). In Matlab, these two eigenvectors are  $4 \times 1$  arrays, and stacking them next to one another gives a  $4 \times 2$  array  $[X;Y]$ , where  $X$  and  $Y$  are  $2 \times 2$  arrays. The ARE solution is the array  $P$  for which  $PX=Y$ . Matlab commands for constructing  $P$  in this way, and for testing that it satisfies the ARE are

```
L=evecs(:, [1 4]);
P=L(3:4,:)/L(1:2,:);
Error=A'*P+P'*A+P'*C*P+D
P =
```



```

    0.61431613579618  -0.54712432051509
    -0.54712432051509  0.07305870788308
Error =
    1.0e-14 *
    -0.04440892098501  -0.13322676295502
    -0.13322676295502  0.11102230246252

```

It follows that the computed solution  $P$  satisfies the ARE to numerical precision. Trying the other three possible eigenvector combinations gives three additional solutions, all of which also satisfy the ARE to numerical precision:

```

P =
    -0.70796721149692  -0.09939673897232
    -0.09939673897232  -0.07854267746637
P =
    -2.46072866126682  -3.90060326102768
    -3.90060326102768  -3.58406557700616
P =
    -2.53652935394154  -3.45287567948491
    -3.45287567948491  -6.22863227159236

```

Note that only the first of our four solutions is positive definite.

## 6 Additional Questions

We conclude our geometric characterization of the ARE solutions set with a host of interesting questions that require further investigation. The ARE (1) is a special case of the Algebraic Riccati Inequality,

$$A^*P + PA + PDP + C \geq 0, \quad (40)$$

and so it would be interesting to use the geometry of Lagrangian subspaces and symplectic forms to consider the matrices  $P$  that satisfy (40). To this end, Andrew Packard has suggested consideration of the papers by Gohberg, Lancaster, and Rodman.

Alan Weinstein has suggested that the non-graph  $F$ -invariant Lagrangian subspaces of  $W$  may correspond to something interesting in the context of AREs. Doubtless this is so, however we leave an investigation of this idea for another thesis.

Although in Section 3 we did obtain a sufficient condition for the existence of real symmetric solutions to (15), (i.e.,  $\text{Spec}(F) \cap i\mathbb{R} = \emptyset$  and  $\langle \xi, C\xi \rangle_V = 0 \implies \xi = 0$ ), we did not prove the existence and uniqueness results from LQR theory. Such a proof would begin with the translation of conditions such as controllability and observability into the language of symplectic forms used in Section 3.

A final tantalizing challenge is given by the characterization of basin boundaries for the flow associated with the Riccati differential equation. In the one dimensional case pictured in Figure 1, the positive ARE solution attracts initial conditions on the interval  $(p_1, \infty)$ , where  $p_1$  is a negative ARE solution. In higher dimensions, the basin boundaries are doubtless beautiful and intriguing.

## A LQR Optimality Conditions

Given  $\dot{x} = Ax + Bu$  with  $x \in \mathbb{R}^m$  and  $u \in \mathbb{R}^n$ , and an arbitrary initial condition  $x(0) = z$ , we wish to find a control signal  $u(t)$  which minimizes the cost function

$$J = \int_0^T (x^T Qx + u^T Ru)dt + \hat{x}^T \hat{Q} \hat{x}, \quad (41)$$

where  $\hat{x} = x(T)$ . Let  $\{t_0, t_1, t_2, \dots, t_N\}$  be a collection of  $N + 1$  evenly spaced points on  $[0, T]$ , (with  $h = T/N$  and  $t_k = kh$ ). If  $u : [0, t] \rightarrow \mathbb{R}^m$  is a smooth<sup>5</sup> control signal, then the collection of state vectors  $\{x_0, x_1, \dots, x_N\}$  established by

$$x_{k+1} = x_k + h(Ax_k + Bu(t_k)) \quad (42)$$

with  $x_0 = z$  satisfies  $\max |x_k - x(t_k)| \rightarrow 0$  as  $N \rightarrow \infty$ , where  $x(t)$  is the solution to the continuous problem. Equivalently, these state vectors can be defined by the matrix equation

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_N \end{bmatrix} = \begin{bmatrix} hB & & & & \\ AhB & hB & & & \\ A^2hB & AhB & hB & & \\ \vdots & \vdots & \vdots & \ddots & \\ A^{N-1}hB & A^{N-2}hB & A^{N-3}hB & \dots & hB \end{bmatrix} \begin{bmatrix} u(t_0) \\ u(t_1) \\ u(t_2) \\ \vdots \\ u(t_{N-1}) \end{bmatrix} + \begin{bmatrix} z \\ Az \\ A^2z \\ \vdots \\ A^{N-1}z \end{bmatrix} \quad (43)$$

where  $A = (I + hA)$ . We abbreviate this equation as  $\mathbf{X} = \mathbf{AU} + \mathbf{Z}$ . The cost function  $J$  defined in (41) can be approximated arbitrarily well as  $N$  gets big by  $\tilde{J} = \mathbf{X}^T \mathbf{QX} + \mathbf{U}^T \mathbf{RU}$ , where

$$\mathbf{Q} = \begin{bmatrix} hQ & & & & \\ & hQ & & & \\ & & \ddots & & \\ & & & hQ & \\ & & & & hQ + \hat{Q} \end{bmatrix} \quad \text{and} \quad \mathbf{R} = \begin{bmatrix} hR & & & & \\ & hR & & & \\ & & \ddots & & \\ & & & hR & \\ & & & & hR \end{bmatrix}. \quad (44)$$

Having described the effect of a continuous control signal on the given system in the discrete setting, we now consider minimizing  $J$ . Our strategy is to find a  $\mathbf{U}$  which minimizes  $\tilde{J}$  (for each  $N$ ), and then to take the limit as  $N \rightarrow \infty$ . Substituting  $\mathbf{X} = \mathbf{AU} + \mathbf{Z}$  into the expression for  $\tilde{J}$ , and differentiating with respect to the components of  $\mathbf{U}$ , we find that

$$\frac{\partial \tilde{J}}{\partial \mathbf{U}} = \mathbf{0} \iff \mathbf{U} = -\mathbf{R}^{-1} \mathbf{A}^T \mathbf{QX}. \quad (45)$$

In matrix form, this necessary condition for the minimization of  $\tilde{J}$  can be written as

$$\begin{bmatrix} u(t_0) \\ u(t_1) \\ u(t_2) \\ \vdots \\ u(t_{N-3}) \\ u(t_{N-2}) \\ u(t_{N-1}) \end{bmatrix} = -R^{-1}B^T \begin{bmatrix} hQ & A^T hQ & (A^T)^2 hQ & \dots & (A^T)^{N-2} hQ & (A^T)^{N-1} (hQ + \hat{Q}) \\ & hQ & A^T hQ & \dots & (A^T)^{N-3} hQ & (A^T)^{N-2} (hQ + \hat{Q}) \\ & & hQ & \dots & (A^T)^{N-4} hQ & (A^T)^{N-3} (hQ + \hat{Q}) \\ & & & \ddots & \vdots & \vdots \\ & & & & A^T hQ & (A^T)^2 (hQ + \hat{Q}) \\ & & & & hQ & A^T (hQ + \hat{Q}) \\ & & & & & hQ + \hat{Q} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_{N-2} \\ x_{N-1} \\ x_N \end{bmatrix}. \quad (46)$$

Because the  $x_i$ 's are related by (42), we now see that condition (46) can be written as the difference equation

$$u(t_{k-1}) = -R^{-1}B^T P_k x_k. \quad (47)$$

A sequence of matrices  $\{P_k\}$  which makes (47) equivalent to (46) as desired can be obtained by substituting (47) into (46). From the last row of (46), we obtain  $P_N = \hat{Q} + hQ$ , and from the other rows, in combination with (42), we obtain

$$\begin{aligned} P_k &= (I + hA^T)P_{k+1}(I + hBR^{-1}B^T P_{k+1})^{-1}(I + hA) + hQ \\ &= P_{k+1} + h(Q + A^T P_{k+1} + P_{k+1}A - P_{k+1}BR^{-1}B^T P_{k+1}) + O(h^2). \end{aligned} \quad (48)$$

for  $k = 1, 2, \dots, N-1$ . From the first equality, it follows that if  $R = R^T$ , then  $P_{k+1} = P_{k+1}^T$  implies  $P_k = P_k^T$ , and so if  $Q$  and  $\hat{Q}$  are symmetric, then  $P_k = P_k^T$  for all  $k$ . Assuming this symmetry in  $R$ ,  $Q$ , and  $\hat{Q}$ , we note

<sup>5</sup>We require  $u$  to be smooth because this causes the LQR problem to give rise to the ARE, which is our main interest in this work. We note however that it is entirely possible for  $u$  to be chosen from a less restricted class of functions. Many optimal control problems are solved by ‘‘bang-bang’’ control strategies that consist of impulses, (see for instance the problem of moving a space craft from one orbit to another with a minimum of fuel).

that in the limit as  $N \rightarrow \infty$ , (48) becomes the following *Riccati differential equation* for the symmetric matrix  $P(t)$ ,

$$-\dot{P} = Q + A^T P + P A - P B R^{-1} B^T P, \quad (49)$$

and  $P_N = \hat{Q} + hQ$  becomes the final condition  $P(T) = \hat{Q}$ ;  $P(t)$  is found by integrating (49) backwards in time from  $T$  to  $t < T$ . Also as  $N \rightarrow \infty$ , (47) becomes the following condition on  $u(t)$ ,

$$u(t) = -R^{-1} B^T P(t) x(t) \quad (50)$$

Although we derived (50) as a necessary condition on a control signal  $u(t)$  that minimizes  $J$ , this condition obviously (and amazingly) comprises a feedback control law that can be used to compute  $u(t)$  as the system evolves. We note that this law is independent of the initial state  $z$ . Also,  $P(t)$  in (50) is independent of  $x(t)$ , and so it can be computed ahead of time, so as to provide via feedback the optimal control signal  $u(t)$  at any time, and for any system state  $x(t)$ .

## A.1 Causality in the Discrete Setting

The transition from a necessary condition to a feedback control law is more involved in the discrete case. Because of (42),  $u(t_k)$  can only affect  $x_l$  with  $l > k$ , (this is apparent in (46) and (47)). A control law however needs to return  $u_k$  given nothing further advanced in time than  $x_k$ . We can get around this causal barrier by combining (42) with (47) to obtain

$$(I + hR^{-1} B^T P_{k+1} B) u_k = -R^{-1} B^T P_{k+1} (I + hA) x_k \quad (51)$$

Thus as in the continuous case, we now have a feedback control law for computing  $u_k$ . We note that (51) approaches (50) as  $N$  gets big.

## B Definitions

In view of the varied conventions used in different texts, in this appendix, we provide a short list of our definitions.

### B.1 Definitions from Linear Systems Theory

The definitions here are needed for Theorem 2.1 in Section 2. Although these definitions hold just as well for matrices over  $\mathbb{C}$ , our LQR problem is over  $\mathbb{R}$ , and so we let  $A$ ,  $B$ , and  $Q$  be real  $n \times n$ ,  $n \times m$ , and  $m \times n$  matrices respectively.

*controllable subspace*

The *controllable subspace*  $C_{A,B}$  of the matrix pair  $(A, B)$  is defined as the range of the  $n \times mn$  matrix  $[B \ AB \ A^2 B \ \dots \ A^{n-1} B]$ .

*controllable*

The matrix pair  $(A, B)$  is called *controllable* if its controllable subspace  $C_{A,B}$  has dimension  $n$ . We note that  $(A, B)$  is controllable if and only if the eigenvalues of  $A + BK$  can take on arbitrary prescribed values by appropriately choosing  $K$ .

*stabilizable*

The matrix pair  $(A, B)$  is called *stabilizable* if there exists a  $K$  such that the eigenvalues of  $A + BK$  are in the open left half plane. Obviously, if  $(A, B)$  is controllable, then it is stabilizable.

*observable and detectable*

The matrix pair  $(Q, A)$  is called *observable (detectable)* if  $(A^T, Q^T)$  is controllable (stabilizable).

These terms derive from the behavior of the linear system  $\dot{x} = Ax + Bu$ ,  $y = Qx$  associated with the matrices  $A$ ,  $B$ , and  $Q$ . For instance, if  $(A, B)$  is controllable, then for an arbitrary initial state  $x_1$  at time  $t_1$ , and an arbitrary target state  $x_2$ , there exists a finite time  $t_2 > t_1$  and a control signal  $u(t)$  over  $[t_1, t_2]$  which moves the system from  $x_1$  at time  $t_1$  to  $x_2$  at time  $t_2$ .

## B.2 Definitions from Section 3

### *direct sum*

If  $V$  and  $W$  are vector spaces over  $\mathbb{F}$ , then we define their *direct sum*  $V \oplus W$  to be a vector space over  $\mathbb{F}$  given by the Cartesian product  $V \times W$  endowed with the obvious addition and scalar multiplication:

$$(v_1, w_1) + (v_2, w_2) \mapsto (v_1 + v_2, w_1 + w_2), \quad \text{and} \quad \alpha(v, w) \mapsto (\alpha v, \alpha w). \quad (52)$$

### *bilinear operator*

Given three vector spaces  $V$ ,  $W$ , and  $X$  over the same base field  $\mathbb{F}$ , a *bilinear operator* is a function  $B : V \times W \rightarrow X$  such that for any  $w \in W$ , the map  $v \mapsto B(v, w)$  is a linear operator from  $V$  to  $X$ , and for any  $v \in V$ , the map  $w \mapsto B(v, w)$  is a linear operator from  $W$  to  $X$ .

- If  $V = W$  and  $B(v, w) = B(w, v)$  for all  $v, w \in V$ , then we say that  $B$  is *symmetric*.
- If  $V = W$  and  $B(v, w) = \overline{B(w, v)}$  for all  $v, w \in V$ , we say  $B$  is *conjugate symmetric*.
- When  $X = \mathbb{F}$ , we call  $B$  a *bilinear form*.

### *quadratic form*

Let  $V$  be a vector space over a field  $\mathbb{F}$ . A map  $Q : V \rightarrow \mathbb{F}$  is called a *quadratic form* on  $V$  if

- $Q(\alpha v) = \alpha^2 Q(v)$  for all  $\alpha \in \mathbb{F}$  and  $v \in V$ .
- $B(u, v) = Q(u + v) - Q(u) - Q(v)$  is a bilinear form on  $V$ .

### *non-degeneracy*

A bilinear form  $B : U \times V \rightarrow \mathbb{F}$  is called *non-degenerate* when  $B(u, v) = 0 \forall u \implies v = 0$ , and  $B(u, v) = 0 \forall v \implies u = 0$ . Of course for this to happen, we need  $\dim(U) = \dim(V)$ . If  $U = V$ , then  $B(u, u) = 0$  only if  $u = 0$ , and if  $B$  maps to  $\mathbb{R}$  when both its arguments are the same, then by continuity, either  $B(x, x) > 0$  or  $B(x, x) < 0$  for all  $x \neq 0$ . In the first case  $B$  is called *positive*.

### *inner product*

An inner product on a vector space  $V$  over  $\mathbb{C}$  is a conjugate symmetric, positive, non-degenerate bilinear form on  $V$ , (with the conjugate requirement vanishing in the case that  $V$  is a vector space over  $\mathbb{R}$ ).

### *positive definite*

An operator  $C$  on an inner product space  $(V, \langle \bullet, \bullet \rangle)$  is said to be *positive definite* if  $\langle Cx, x \rangle \geq 0 \forall x \in V$ , with  $\langle Cx, x \rangle = 0$  only if  $x = 0$ .

### *symplectic*

A bilinear form  $B$  is called *symplectic* if it is non-degenerate and if  $B(u, v) = -B(v, u)$ . As for etymological origins, the adjective *symplectic* derives from the Greek *symplektikos* which means intertwining. The associated Greek verb *symplekein* means to plait together or intertwine [8].

### *symplectic vector space*

A vector space endowed with a symplectic bilinear form is called a *symplectic vector space*.

### *symplectic complement*

Let  $W$  be a symplectic vector space with symplectic form  $\sigma$ . If  $\Lambda$  is a subspace of  $W$ , then we define its *symplectic complement*  $\Lambda^\perp$  by

$$\Lambda^\perp = \{v \in W \mid \sigma(v, w) = 0 \text{ for all } w \in \Lambda\} \quad (53)$$

We note that  $(\Lambda^\perp)^\perp = \Lambda$  and that  $\dim \Lambda + \dim \Lambda^\perp = \dim W$ . Also,  $\Lambda \subset \Lambda^\perp$  is equivalent to

$$X, Y \in \Lambda \implies \sigma(X, Y) = 0, \quad (54)$$

and  $\Lambda^\perp \subset \Lambda$  is equivalent to

$$\sigma(X, Y) = 0 \forall X \in \Lambda \implies Y \in \Lambda. \quad (55)$$

There are many different ways in which  $\Lambda$  can relate to  $\Lambda^\perp$ . The two of interest to us are as follows:

- if  $\Lambda \cap \Lambda^\perp = \{0\}$ , then  $\Lambda$  is called *symplectic*.  $\Lambda$  is symplectic if and only if  $\sigma$  restricts to a non-degenerate form on  $\Lambda$ . A symplectic subspace with the restricted form is a symplectic vector space in its own right.

- if  $\Lambda = \Lambda^\perp$  then  $\Lambda$  is called *Lagrangian*. If  $\Lambda$  is a Lagrangian subspace of  $W$ , then  $\dim(\Lambda) = \frac{1}{2} \dim(W)$ .

*dual space*

The *dual space* of a finite dimensional vector space  $V$  is given by  $L(V, \mathbb{F})$  and is denoted  $V'$ . The mappings comprising  $V'$  are referred to as *functionals*. The dual space  $V'$  has the same dimension as  $V$ , as can be seen by noting that  $V'$  is isomorphic to the space of  $1 \times \dim(V)$  matrices.

### B.3 Adjoints

In this section, we discuss the meanings of the word *adjoint* used in this report. We begin with a definition of the adjoint of a map between inner product spaces.

**A** *adjoint*:

If  $A : U \rightarrow V$ , then the *adjoint*  $A^*$  of  $A$  is given by  $A^* : V \rightarrow U$  such that  $\langle A^*v, u \rangle_U = \langle v, Au \rangle_V$  for all  $u \in U$  and for all  $v \in V$ , where  $\langle \bullet, \bullet \rangle_U$  and  $\langle \bullet, \bullet \rangle_V$  are inner products on  $U$  and  $V$  respectively. As a diagram, this definition becomes

$$\begin{array}{ccc} A : U & \longrightarrow & V & \forall u \in U \\ \langle A^*v, u \rangle_U = \langle v, Au \rangle_V & & \forall v \in V \\ U & \longleftarrow & V : A^* \end{array}$$

We let  $M(A, \{u_i\}, \{v_i\})$  denote the *matrix* of a map  $A : U \rightarrow V$  with respect to a basis  $\{u_i\}_{i=1}^m$  on  $U$  and a basis  $\{v_i\}_{i=1}^n$  on  $V$ . If  $a_{ij}$  is a generic element at row  $i$  and column  $j$  of this matrix, then  $Au_j = a_{1j}v_1 + a_{2j}v_2 + \dots + a_{nj}v_n$ . It follows that if the bases  $\{u_i\}$  and  $\{v_i\}$  are orthonormal with respect to  $\langle \bullet, \bullet \rangle_U$  and  $\langle \bullet, \bullet \rangle_V$  respectively, then  $M(A^*, \{v_i\}, \{u_i\})$  is the conjugate transpose of  $M(A, \{u_i\}, \{v_i\})$ . With this in mind, we define the adjoint of a matrix in the obvious way,

**B** *matrix adjoint*:

The *matrix adjoint* of a matrix is its conjugate transpose.

The next definition makes use of the dual space  $V'$  naturally associated with a vector space  $V$ . If  $x \in V$  and  $\xi \in V'$ , then we let  $\langle \xi, x \rangle_V$  denote  $\xi(x)$ , where the subscript  $V$  indicates that elements from  $V$  are considered vectors, and elements from the dual space  $V'$  are considered functionals. This pairing of elements from  $V$  and  $V'$  is highly symmetrical; because the dual of  $V'$  is simply  $V$  again, we find it pleasing to think of  $V$  and  $V'$  not as one deriving from the other but as two sides of the same coin. In particular, we note that  $\langle \xi, x \rangle_V = \langle x, \xi \rangle_{V'}$ , where on the right,  $\xi \in V'$  is considered a vector and  $x \in V$  is considered a functional.

**C** *natural adjoint*:

If  $A : U \rightarrow V$ , then we define the *natural adjoint*  $A'$  of  $A$  by  $A' : V' \rightarrow U'$  such that  $\langle A'\omega, v \rangle_{V'} = \langle \omega, Av \rangle_U$  for all  $\omega \in U'$  and for all  $v \in V$ , where  $\langle \bullet, \bullet \rangle_U$  and  $\langle \bullet, \bullet \rangle_{V'}$  are as discussed immediately above. As a diagram, this definition becomes

$$\begin{array}{ccc} A : U & \longrightarrow & V & \forall v \in V \\ \langle A'\omega, v \rangle_{V'} = \langle \omega, Av \rangle_U & & \forall \omega \in V' \\ U' & \longleftarrow & V' : A' \end{array}$$

The natural adjoint **C** requires fewer ingredients than the adjoint **A** (in particular  $U$  and  $V$  don't have to have inner products), and so it charms us with the aesthetics of minimalism. Unfortunately however, **A** and **C** are irreconcilable in the complex setting, due to the necessary *conjugate* symmetry of any inner product on a complex vector space, (i.e., inner products on complex vector spaces are hermitian). In detail, we make the following

**Claim:** There are no isomorphisms  $I_U : U' \rightarrow U$  and  $I_V : V' \rightarrow V$ , and no inner products  $\langle \bullet, \bullet \rangle_U$  and  $\langle \bullet, \bullet \rangle_{V'}$  for which the mappings  $A^* : V \rightarrow U$  and  $I_U \circ A' \circ I_V^{-1} : V \rightarrow U$  are the same.

**Proof:** First note that if  $\{u_i\}$  and  $\{v_i\}$  are bases on  $U$  and  $V$ , and if we construct the corresponding bases  $\{\tilde{u}_i\}$  and  $\{\tilde{v}_i\}$  on  $U'$  and  $V'$  (so that  $\tilde{u}_i(u_j) = \delta_{ij}$  and  $\tilde{v}_i(v_j) = \delta_{ij}$ ), then  $M(A', \{\tilde{v}_i\}, \{\tilde{u}_i\})$  is the transpose (no conjugate) of  $M(A, \{u_i\}, \{v_i\})$ . Now let  $\{u_i\}$  and  $\{v_i\}$  be orthonormal with respect to inner products

on  $U$  and  $V$ . We established already that if  $M(A, \{u_i\}, \{v_i\}) = [a_{ij}]$ , then  $M(A^*, \{v_i\}, \{u_i\}) = [\bar{a}_{ji}]$ . The  $(i, j)$  element of  $M(I_U \circ A' \circ I_V^{-1}, \{v_i\}, \{u_i\})$  is given by

$$M(I_V^{-1}, \{v_i\}, \{\tilde{v}_i\})M(A', \{\tilde{v}_i\}, \{\tilde{u}_i\})M(I_U, \{\tilde{u}_i\}, \{u_i\}) = e_{ik}a_{lk}f_{lj}, \quad (56)$$

with summation on the repeated indices. There is no choice of  $[e_{ik}]$  and  $[f_{lj}]$  which gives  $\bar{a}_{ji}$  (the best we can obtain is  $a_{ji}$ ). $\square$

In the case of real mappings however (as with the matrices in Section 3), definitions **A** and **C** can be made to coincide.

## C Theorems from Section 3

Here we offer detail on claims made in Section 3. As noted in Section 4, this development works just as well if we replace  $V'$  with  $V$ , and if we replace  $\langle \bullet, \bullet \rangle_V$  with an inner product on  $V$ . Each proof in this section (and in fact in this entire thesis) is original, and not taken from any reference.

**Theorem C.1** *Any symplectic bilinear form  $\sigma(x, y)$  on an even dimensional vector space  $W$  can be expressed in terms of some non-degenerate bilinear form  $\langle \bullet, \bullet \rangle$  as  $\langle p_1(x), p_2(y) \rangle - \langle p_2(x), p_1(y) \rangle$ . Here  $p_1(x)$  is the projection of  $x$  onto a subspace  $W_1$  of  $W$ , and  $p_2(x)$  is the projection of  $x$  onto a subspace  $W_2$  of  $W$ , where  $W_1 \oplus W_2 = W$ , and  $\dim W_1 = \dim W_2 = \frac{1}{2} \dim W$ .*

**Proof:** We will build a basis  $(w_1, \dots, w_n, \tilde{w}_1, \dots, \tilde{w}_n)$  for the  $2n$  dimensional space  $W$  so that the given symplectic bilinear form has the desired representation, with  $W_1 = \text{span}(w_1, \dots, w_n)$ ,  $W_2 = \text{span}(\tilde{w}_1, \dots, \tilde{w}_n)$ , and

$$\begin{aligned} \langle x, y \rangle &:= x_1\tilde{y}_1 + x_2\tilde{y}_2 + \dots + x_n\tilde{y}_n, \\ &+ \tilde{x}_1y_1 + \tilde{x}_2y_2 + \dots + \tilde{x}_ny_n. \end{aligned} \quad (57)$$

where  $x_i$  and  $\tilde{x}_i$  are the components of  $x$  with respect to the constructed basis according to

$$x = x_1w_1 + \dots + x_nw_n + \tilde{x}_1\tilde{w}_1 + \dots + \tilde{x}_n\tilde{w}_n. \quad (58)$$

Our argument is inductive, and involves the construction of  $n$  subspaces  $V_1, V_2, \dots, V_n$  of  $W$ . Let  $V_1 := W$  and pick any nonzero  $w_1 \in V_1$ . Then pick  $\tilde{w}_1 \in V_1$  so that  $\sigma(w_1, \tilde{w}_1) = 1$ . This is possible to do because  $\sigma$  is symplectic and therefore non-degenerate on  $V_1$ . Note that  $\text{span}(\{w_1, \tilde{w}_1\})$  is a 2-dimensional subspace of  $V_1$ , ( $\tilde{w}_1 = \kappa w_1$  isn't possible because it would cause  $\sigma(w_1, \tilde{w}_1) = 0$ ), and that a dimension  $2(n-1)$  subspace  $V_2$  exists such that  $\sigma(\omega_1, v) = 0$  and  $\sigma(\tilde{\omega}_1, v) = 0$  for all  $v \in V_2$ . Note that  $V_2 \oplus \text{span}(w_1, \tilde{w}_1) = V_1$ , and that  $\sigma|_{V_2}$  is a symplectic form over  $V_2$ , and so by induction we obtain the desired basis elements. $\square$

**Theorem C.2**  $\Lambda \subset W$  can be written as  $\left\{ \begin{bmatrix} x \\ Px \end{bmatrix} : x \in V \right\}$  if and only if  $\pi|_\Lambda$  is a bijection.

**Proof:** Let  $\pi : W \rightarrow V$  be the natural projection defined by  $\pi([x \ y]^T) = x$ . If  $\Lambda \subset W$  is given by  $\{[x \ Px]^T | x \in V\}$ , then  $\tilde{\pi} := \pi|_\Lambda$  is clearly a bijection. For the converse, if  $\Lambda$  is a subspace of  $W$  for which  $\tilde{\pi} := \pi|_\Lambda$  is a bijection, then  $\tilde{\pi} \in L(\Lambda, V)$ ,  $\tilde{\pi}^{-1} \in L(V, \Lambda)$ , and  $x \in V$  gets sent by  $\tilde{\pi}^{-1}$  to  $[x \ y]^T \in \Lambda$ . Composing the linear map  $\tilde{\pi}^{-1}$  with the projection from  $[x \ y]^T \in \Lambda$  to  $y \in V$  gives  $P \in L(V)$  such that  $y = Px$ , allowing  $\Lambda$  to be expressed as  $\{[x \ Px]^T | x \in V\}$ . $\square$

**Theorem C.3** A subspace of  $W$  given by  $\Lambda = \{[x \ Px]^T | x \in V\}$  is Lagrangian if and only if  $P = P'$ .

**Proof:** If  $\Lambda$  is Lagrangian, then  $\sigma(X, Y) = 0$  for all  $X, Y \in \Lambda$ . Using (17)<sub>4</sub>, we note that

$$\sigma(X, Y) = \sigma\left(\begin{bmatrix} x \\ Px \end{bmatrix}, \begin{bmatrix} y \\ Py \end{bmatrix}\right) = \langle Px, y \rangle_V - \langle Py, x \rangle_V = \langle Px, y \rangle_V - \langle P'y, y \rangle_V = \langle (P - P')x, y \rangle_V \quad (59)$$

for all  $x, y \in V$ , and so we must have  $P = P'$ . For the converse, let  $\Lambda = \{[x \ Px]^T | x \in V\}$  with  $P = P'$ . The steps already taken show that  $\sigma(X, Y) = 0$  for all  $X, Y \in \Lambda$ . Now consider  $[x \ \eta]^T \in W$  not in  $\Lambda$ , that is,

with  $\eta \neq Px$ . We want to show that some  $u \in V$  exists so that  $[u \ Pu]^T \in \Lambda$  satisfies  $\sigma([x \ \eta]^T, [u \ Pu]^T) \neq 0$ . Note that

$$\sigma\left(\begin{bmatrix} x \\ \eta \end{bmatrix}, \begin{bmatrix} u \\ Pu \end{bmatrix}\right) \neq 0 \iff \langle \eta, u \rangle_V - (\langle Pu, x \rangle_V = \langle P'x, u \rangle_V = \langle Px, u \rangle_V) = \langle \eta - Px, u \rangle_V \neq 0 \quad (60)$$

where we've used  $(17)_4$  and the hypothesis that  $P = P'$ . If no such  $u$  exists, then  $\langle y - Px, u \rangle_V = 0$  for every  $u \in V$ , in particular for  $u = \eta - Px$ . This however implies that  $\eta - Px = 0$ , which is a contradiction, and so it must be that  $\sigma(X, Y) = 0 \forall X \in \Lambda \implies Y \in \Lambda$  as desired. Finally, note that  $\Lambda$  is isomorphic to  $V$ , and so  $\dim \Lambda = \dim V = \frac{1}{2} \dim W$ . Thus  $\Lambda$  is Lagrangian as desired.  $\square$

**Theorem C.4** *The Hamiltonian  $F$  defined by  $F = \begin{bmatrix} A & C \\ -D & -A' \end{bmatrix}$  can equivalently be defined by  $\sigma(X, FY) = Q(X, Y)$  for all  $X, Y \in W$ .*

**Proof:** Start with  $F : W \longrightarrow W$  written as

$$F = \begin{bmatrix} F_{11} & F_{12} \\ F_{21} & F_{22} \end{bmatrix}. \quad (61)$$

Then,  $\sigma(X, FY) = Q(X, Y)$  becomes

$$\langle x_2, F_{11}y_1 + F_{12}y_2 \rangle_V - \langle F_{21}y_1 + F_{22}y_2, x_1 \rangle_V = \langle Dx_1, y_1 \rangle_V + \langle y_2, Ax_1 \rangle_V + \langle x_2, Ay_1 \rangle_V + \langle y_2, Cx_2 \rangle_V \quad (62)$$

where  $x_1, x_2, y_1$ , and  $y_2$  are arbitrary vectors in  $V$ . Setting  $x_2$  and  $y_2$  to zero gives

$$\begin{aligned} -\langle F_{21}y_1, x_1 \rangle_V &= \langle Dx_1, y_1 \rangle_V, \\ \implies \langle -F_{21}y_1, x_1 \rangle_V &= \langle D'y_1, x_1 \rangle_V, \end{aligned} \quad (63)$$

from which it follows that  $F_{21} = -D' = -D$ . Similarly, we find that  $F_{22} = -A'$ ,  $F_{11} = A$ , and  $F_{12} = C$ .  $\square$

**Theorem C.5**  $\sigma(X, FY) + \sigma(FX, Y) = 0$

**Proof:**

$$\begin{aligned} \sigma(X, FY) &= Q(X, Y) && \text{from C.4} \\ &= Q(Y, X) && \text{from the symmetry of } Q \\ &= \sigma(Y, FX) && \text{from C.4} \\ &= -\sigma(FX, Y) && \text{from the skew-symmetry of } \sigma. \end{aligned}$$

**Theorem C.6** *Let  $B : V \times V \longrightarrow \mathbb{F}$  be a non-degenerate bilinear form on a finite dimensional vector space  $V$ . For every functional  $\varphi$  in the dual space  $V'$ , there exists a unique  $x \in V$  such that  $\varphi(y) = B(y, x)$  for all  $y \in V$ .*

**Proof:** Let  $\{e_i\}$  be a basis for  $V$ . We need to show that a unique  $x = \sum x_i e_i$  exists so that  $\varphi(e_i) = B(e_i, x)$  for each  $e_i$ . Finding the  $x_i$ 's is accomplished by solving

$$\begin{bmatrix} \varphi(e_1) \\ \varphi(e_2) \\ \vdots \\ \varphi(e_n) \end{bmatrix} = \begin{bmatrix} B(e_1, e_1) & B(e_1, e_2) & \cdots & B(e_1, e_n) \\ B(e_2, e_1) & B(e_2, e_2) & \cdots & B(e_2, e_n) \\ \vdots & \vdots & \ddots & \vdots \\ B(e_n, e_1) & B(e_n, e_2) & \cdots & B(e_n, e_n) \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad (64)$$

We'll be done if we can show that the matrix is invertible. If the matrix isn't invertible, then some nonzero  $[x_1, x_2, \dots, x_n]^T$  gets mapped to zero, corresponding to a nonzero  $x$  for which  $B(e_i, x) = 0$  for each  $e_i$ . This implies that  $B(x, x) = 0$  which contradicts the definition of  $B$  as a non-degenerate bilinear form.  $\square$

**Theorem C.7** *A Lagrangian subspace on which  $Q = 0$  is equivalent to one which is invariant under  $F$ .*

**Proof:** Let  $v \in \Lambda$  where  $\Lambda$  is a Lagrangian subspace. We need to show that  $Fv \in \Lambda$  if and only if  $Q$  vanishes on  $\Lambda$ . If  $Q$  vanishes on  $\Lambda$ , then  $Q(u, v) = 0$  for every  $u \in \Lambda$ . Using C.4, this becomes  $\sigma(u, Fv) = 0$  for all  $u \in \Lambda$ , which means that  $Fv \in \Lambda$ . For the converse, let  $Fv \in \Lambda$ , causing  $\sigma(u, Fv) = 0$  for every  $u \in \Lambda$ . Using

C.4 again we see that this implies  $Q(u, v) = 0$  for every  $v \in \Lambda$ , and so  $Q$  vanishes on  $\Lambda$  as desired, ( $v$  was chosen arbitrarily in  $\Lambda$  to begin with).  $\square$

**Definition of  $V_\Lambda$**

For any  $\lambda \in \mathbb{C}$ , we let  $V_\lambda$  denote the subspace of  $W$  consisting of vectors  $v$  for which  $(F - \lambda)^N v = 0$  for  $N$  sufficiently large. Note that  $V_\lambda$  is defined for *every*  $\lambda \in \mathbb{C}$ , whether or not  $\lambda \in \text{Spec}(F)$ . When  $\lambda$  is not an eigenvalue of  $F$ ,  $V_\lambda = \{0\}$ , otherwise  $V_\lambda$  is the nonzero space of generalized eigenvectors of  $F$  associated with  $\lambda$ .

**Theorem C.8**  $\lambda_1 + \lambda_2 \neq 0$  implies that  $\sigma(V_{\lambda_1}, V_{\lambda_2}) = 0$ .

**Proof:** Start with  $\lambda_1, \lambda_2 \in \mathbb{C}$  such that  $\lambda_1 + \lambda_2 \neq 0$ . We suppose that  $\lambda_1, \lambda_2 \in \text{Spec}(F)$  because if this is false, then the theorem conclusion holds trivially. We first establish that  $(F + \lambda_2)V_{\lambda_1} = V_{\lambda_1}$ , by showing that  $(F + \lambda_2)$  is a bijective linear operator on  $V_{\lambda_1}$ .

- We first show that  $(F + \lambda_2)$  is a linear operator on  $V_{\lambda_1}$ .  
 $v \in V_{\lambda_1}$  implies that for some  $N$ ,  $(F - \lambda_1)^N v = 0$ , which implies that  $(F + \lambda_2)(F - \lambda_1)^N v = 0$ . But  $(F + \lambda_2)$  commutes with  $(F - \lambda_1)^N$ , and so  $(F - \lambda_1)^N (F + \lambda_2)v = 0$ , establishing that  $(F + \lambda_2)v \in V_{\lambda_1}$  as desired. (We had to make sure that  $(F + \lambda_2)$  didn't map from  $V_{\lambda_1}$  to outside of  $V_{\lambda_1}$ ).
- Next we show that  $(F + \lambda_2)$  is an injective linear operator on  $V_{\lambda_1}$ .

We need to show that 0 is the only vector in  $V_{\lambda_1}$  that is mapped to 0 by  $(F + \lambda_2)$ . If this is false, then some nonzero  $v \in V_{\lambda_1}$  satisfies  $Fv = -\lambda_2 v$ . Because  $v \in V_{\lambda_1}$ , some  $N$  exists for which  $(F - \lambda_1)^N v = 0$ . But

$$(F - \lambda_1)^N v = (F - \lambda_1)(F - \lambda_1) \cdots (F - \lambda_1)v = v(\lambda_1 + \lambda_2)^N (-1)^N \neq 0, \quad (65)$$

which is a contradiction.

An injective operator is also bijective, and so we are done. If  $(F + \lambda_2)V_{\lambda_1} = V_{\lambda_1}$ , then  $(F + \lambda_2)^N V_{\lambda_1} = V_{\lambda_1}$  for any  $N$ , and so

$$\sigma(V_{\lambda_1}, V_{\lambda_2}) = \sigma((F + \lambda_2)^N V_{\lambda_1}, V_{\lambda_2}) = \sigma(V_{\lambda_1}, (-F + \lambda_2)^N V_{\lambda_2}), \quad (66)$$

where the last equality is due to C.5 and the linearity of  $\sigma$ . Note that  $(-F + \lambda_2)^N V_{\lambda_2}$  is 0 for  $N$  sufficiently large, and so  $\sigma(V_{\lambda_1}, V_{\lambda_2}) = 0$  as desired.  $\square$

**Theorem C.9**  $\lambda \in \text{Spec}(F)$  implies that  $\bar{\lambda}, -\lambda \in \text{Spec}(F)$ .

**Proof:** Let  $\lambda \in \text{Spec}(F)$ . Then  $Fv = \lambda v$  implies  $F\bar{v} = \bar{\lambda}\bar{v}$  ( $F = \bar{F}$  because  $F$  is composed of maps that are real), and so  $\bar{\lambda} \in \text{Spec}(F)$  as desired. To show that  $-\lambda \in \text{Spec}(F)$ , let  $v$  be a nonzero vector in  $V_\lambda$ . Using C.8, we know that  $\sigma(v, V_\mu) = 0$  whenever  $\mu \neq -\lambda$ . If  $-\lambda$  was not an eigenvalue of  $F$ , then we would have  $V_{-\lambda} = \{0\}$ , causing  $\sigma(v, V_{-\lambda}) = 0$  as well. This means  $\sigma(v, w) = 0$  for every  $w \in W$ , which contradicts the fact that  $v \neq 0$  and that  $\sigma$  is non-degenerate.  $\square$

**Theorem C.10** When  $\lambda \neq 0$ ,  $V_\lambda$  is the dual space of  $V_{-\lambda}$ .

**Proof:** Let  $V_\lambda \subset W$  be given, and let  $V'_\lambda$  denote its dual. Noting that  $\sigma$  is a symplectic bilinear form on  $W$ , we know from C.6 that every  $\varphi \in V'_\lambda$  has a unique corresponding  $X \in W$  for which  $\varphi(Y) = \sigma(Y, X)$  for all  $Y \in V_\lambda$ . But in fact, from C.8, it must be that  $X$  only has components in  $V_{-\lambda}$ . With  $\{\varphi_i\}$  a basis for  $V'_\lambda$ , and with  $\{X_i\}$  a collection of corresponding vectors in  $V_{-\lambda}$ , we can establish a linear mapping from  $V'_\lambda$  into  $V_{-\lambda}$  by  $\varphi_i \mapsto X_i$ .

Suppose that this mapping is not injective. Then it sends some nonzero functional  $\varphi = \gamma_1 \varphi_1 + \cdots + \gamma_n \varphi_n$  to  $\gamma_1 X_1 + \cdots + \gamma_n X_n = 0$ . Because  $\varphi \in V'_\lambda$  is nonzero however, some nonzero  $Y \in V_\lambda$  exists for which  $\varphi(Y) \neq 0$ . Note that

$$\begin{aligned} \varphi(Y) &= \gamma_1 \varphi_1(Y) + \cdots + \gamma_n \varphi_n(Y) \\ &= \gamma_1 \sigma(Y, X_1) + \cdots + \gamma_n \sigma(Y, X_n) \\ &= \sigma(Y, \gamma_1 X_1 + \cdots + \gamma_n X_n) \\ &= \sigma(Y, 0) = 0 \end{aligned} \quad (67)$$

which is a contradiction, and so the mapping is injective.



To show surjectivity, suppose some nonzero  $X \in V_{-\lambda}$  is beyond the range of the mapping. Then, consider the functional  $\varphi(Y) := \sigma(Y, X)$ . From C.8,  $\sigma(Y, X) = 0$  for every  $Y$  not in  $V_\lambda$ . If this functional was nonzero for some  $Y \in V_\lambda$ ,  $X$  would be within the range of our mapping, contrary to our supposition. So it must be that  $\sigma(Y, X) = 0$  for every  $Y$  in  $V_\lambda$ , and thus for every  $Y$  in  $W$ . But  $X \neq 0$  and so this contradicts the non-degeneracy of  $\sigma$ .

Thus the proposed mapping is surjective and we have established a bijection between  $V'_\lambda$  and  $V_{-\lambda}$ . The two are isomorphic, and we say that  $V_\lambda$  is the dual space of  $V_{-\lambda}$ . In greater detail,  $V'_{-\lambda} = \{\alpha(X) \mid X \in V_\lambda\}$  where  $\alpha(X) : V_{-\lambda} \rightarrow \mathbb{C}$  according to  $\alpha(X)(Y) \mapsto \sigma(X, Y)$ .  $\square$

**Theorem C.11** *When  $\lambda \neq 0$ ,  $V_\lambda \oplus V_{-\lambda}$  is a symplectic vector space under  $\sigma$ .*

**Proof:** If  $\sigma$  is a symplectic bilinear form on  $W$  then clearly  $\sigma(X, Y) = -\sigma(Y, X)$  for  $X, Y$  from any subspace of  $W$ , however  $\sigma$  is not automatically non-degenerate on subspaces of  $W$ . For instance if  $\sigma$  is a symplectic bilinear form on  $\mathbb{C}^2$ , then  $\sigma$  is necessarily 0 on any one dimensional subspace, and so it is impossible to satisfy

$$X \neq 0 \implies \exists Y \text{ such that } \sigma(X, Y) \neq 0. \quad (68)$$

on this subspace. However this property does hold on  $V_\lambda \oplus V_{-\lambda}$  as we now show.

Pick any nonzero  $X \in V_\lambda \oplus V_{-\lambda}$ , and write  $X = X_1 + X_2$ , where  $X_1 \in V_\lambda$  and  $X_2 \in V_{-\lambda}$ . Either  $X_1$  or  $X_2$  (or both) are nonzero. We suppose  $X_1$  is nonzero (the proof for  $X_1 = 0$  is similar). Suppose that no  $Y \in V_{-\lambda}$  exists for which  $\sigma(X_1, Y) \neq 0$ . By C.8,  $\sigma(X_1, Y) = 0$  for  $Y$  taken from every other subspace of  $W$ , and so  $\sigma(X_1, Y) = 0$  for every  $Y$  in  $W$ . But  $X \neq 0$ , which contradicts the non-degeneracy of  $\sigma$  over  $W$ . Thus it must be that some  $Y \in V_{-\lambda}$  exists for which  $\sigma(X_1, Y) \neq 0$ , and so  $\sigma$  is non-degenerate on  $V_\lambda \oplus V_{-\lambda}$  as desired.  $\square$

**Theorem C.12** *If  $\alpha \neq \beta$  then  $V_\alpha \cap V_\beta = \{0\}$ .*

**Proof:** Let  $\alpha \neq \beta$ , and suppose that some nonzero  $v$  is in both  $V_\alpha$  and  $V_\beta$ . Then integers  $n_\alpha$  and  $n_\beta$  exist so that  $(F - \alpha)^{n_\alpha} v = 0$ ,  $(F - \beta)^{n_\beta} v = 0$ , and  $(F - \beta)^{n_\beta - 1} v := \tilde{v} \neq 0$ . Note that  $(F - \alpha)^{n_\alpha}$  commutes with  $(F - \beta)^{n_\beta - 1}$ , and so

$$(F - \alpha)^{n_\alpha} v = 0 \implies (F - \beta)^{n_\beta - 1} (F - \alpha)^{n_\alpha} v = 0 \implies (F - \alpha)^{n_\alpha} \tilde{v} = 0. \quad (69)$$

But  $F\tilde{v} = \beta\tilde{v}$ , and so the right-hand equality in (69) becomes  $(\beta - \alpha)^{n_\alpha} = 0$ , which is a contradiction because  $\alpha \neq \beta$ .  $\square$

**Theorem C.13**  *$V_0$  is a symplectic vector space under  $\sigma$ .*

**Proof:** If  $V_0 = \{0\}$  (that is, if  $0 \neq \text{Spec}(F)$ ), then  $\sigma$  is trivially a symplectic bilinear form over  $V_0$ , and  $V_0$  is a symplectic vector space under  $\sigma$ . Now suppose  $0 \in \text{Spec}(F)$ . In this case  $V_0$  contains more than the zero vector, and in fact has even dimension, because in the following decomposition of the even dimensional  $W$ ,

$$W = V_0 \oplus (V_{\lambda_1} \oplus V_{-\lambda_1}) \oplus (V_{\lambda_2} \oplus V_{-\lambda_2}) \oplus \cdots \oplus (V_{\lambda_n} \oplus V_{-\lambda_n}). \quad (70)$$

each term in parenthesis has even dimension, and none of these subspaces intersect (see C.12). Of course  $\sigma(X, Y) = -\sigma(Y, X)$  for all  $X$  and  $Y$  in  $V_0$ . To show the non-degeneracy of  $\sigma$  over  $V_0$ , note that if  $X$  is a nonzero vector in  $V_0$ , then from C.8,  $\sigma(X, Y) = 0$  for every  $Y$  not in  $V_0$ . As with C.11, it then follows that there must be some  $Y \in V_0$  for which  $\sigma(X, Y) \neq 0$ , or else we would have  $\sigma(X, Y) = 0$  for every  $Y \in W$  where  $X \neq 0$ , contradicting the non-degeneracy of  $\sigma$  over  $W$ .  $\square$

**Theorem C.14** *If  $U \subset W$  is  $F$ -invariant, then  $U$  is equal to a direct sum of subspaces  $U_i$  of the subspaces  $V_{\lambda_i}$  of generalized eigenvectors of  $F$ .*

**Example over  $\mathbb{R}$ :** Consider the independent vectors  $\mathbf{e}_1$  and  $\mathbf{e}_2$  in  $\mathbb{R}^2$ . Clearly  $\mathbb{R}^2 = V_1 \oplus V_2$  where  $V_1 = \text{span}(\{\mathbf{e}_1\})$  and  $V_2 = \text{span}(\{\mathbf{e}_2\})$ . The subspace  $U = \text{span}(\{\mathbf{e}_1 + \mathbf{e}_2\})$  however is *not* given by the direct sum of subspaces of the  $V_i$ 's.

**Proof:** Suppose the subspace  $U$  is not equal to the direct sum of any collection of subspaces  $U_i$  of  $V_{\lambda_i}$ , where  $V_{\lambda_i}$  is the space of generalized eigenvectors associated with  $\lambda_i \in \text{Spec}(F)$ . The direct sum of the  $V_{\lambda_i}$ 's is  $W$ , and so to every  $u \in U \subset W$  there corresponds a unique element  $u_i$  of each  $V_{\lambda_i}$  such that  $u = u_1 + \cdots + u_n$ . It is easy to verify that because  $U$  is a subspace, the collection of all possible such  $u_i$ 's forms a subspace  $U_i$

of  $V_{\lambda_i}$ . This establishes that  $U \subset U_1 \oplus \cdots \oplus U_n$ . Our hypothesis implies that this subset relation is strict, that is, that some vector  $v_1 + \cdots + v_n \in U_1 \oplus \cdots \oplus U_n$  exists that isn't in  $U$ , (keep in mind that  $v_i \in U_i \forall i$ ).

In the proof of C.8, we established that  $\lambda_1 + \lambda_2 \neq 0$  implies  $(F + \lambda_2)V_{\lambda_1} = V_{\lambda_1}$ , (equivalently that  $\lambda_1 \neq \lambda_2$  implies  $(F - \lambda_2)V_{\lambda_1} = V_{\lambda_1}$ ). Note that with  $d = \dim(W)$ ,  $(F - \lambda_i)^d$  annihilates  $V_{\lambda_i}$ . It follows that  $F_{\lambda_i} := (F - \lambda_2)^d (F - \lambda_3)^d \cdots (F - \lambda_n)^d$  is a bijective operator on  $V_{\lambda_1}$  that annihilates every other  $V_{\lambda_i}$ . Similarly, we define  $F_{\lambda_i}$  for each  $\lambda_i$ . Because  $F_{\lambda_i}$  is a bijective operator on  $V_{\lambda_i}$ , every  $v_i \in U_i \subset V_{\lambda_i}$  has a pre-image  $\hat{v}_i \in V_{\lambda_i}$ .

If one of the pre-images  $\hat{v}_i \in V_{\lambda_i}$  is not in  $U_i$ , then  $F_{\lambda_i}$  maps an entire subspace (e.g.,  $\alpha \hat{v}_i$ ) into  $U_i \subset V_{\lambda_i}$ . But  $F_{\lambda_i}$  is a bijective operator on  $V_{\lambda_i}$ , and so to make room for the image of  $\alpha \hat{v}_i$ , it must be that  $F_{\lambda_i}$  moves some vectors off of  $U_i$ . Let  $u_i$  be one of these vectors. The corresponding vector in  $U$  gets mapped by  $F_{\lambda_i}$  (a polynomial in  $F$ ) to some vector in  $V_{\lambda_i}$  that is not in  $U_i$ , that is, to a vector that is not in  $U$ . This establishes that  $U$  is not invariant under  $F$ .

Now suppose that every pre-image  $\hat{v}_i \in V_{\lambda_i}$  is in its respective  $U_i$ . To each of these pre-images there corresponds a vector  $u_i \in U$  such that  $F_{\lambda_i} u_i = v_i$ . In words, the sum of polynomial functions of  $F$  applied to vectors  $u_i \in U$  is a vector  $v_1 + \cdots + v_n$  that is not in  $U$ . Thus  $U$  is not invariant under  $F$ .  $\square$

**Theorem C.15** *F-invariant Lagrangian subspaces of  $W$  can be represented as*

$$\Lambda = \left( \bigoplus_{\lambda \in \Sigma} \Lambda_\lambda \right) \oplus \Lambda_0, \quad (71)$$

where  $\Lambda_\lambda$  and  $\Lambda_0$  are  $F$ -invariant Lagrangian subspaces of the symplectic spaces  $V_\lambda \oplus V_{-\lambda}$  and  $V_0$  respectively, and  $\Sigma$  is a set satisfying  $\Sigma \cup (-\Sigma) = \text{Spec}(F) \setminus \{0\}$  and  $\Sigma \cap (-\Sigma) = \emptyset$ .

**Proof:** Let  $\Lambda$  be an  $F$ -invariant Lagrangian subspace of  $W$ . Note that  $W$  can be expressed in terms of the invariant subspaces of  $F$  as

$$W = (V_{\lambda_1} \oplus V_{-\lambda_1}) \oplus (V_{\lambda_2} \oplus V_{-\lambda_2}) \oplus \cdots \oplus (V_{\lambda_n} \oplus V_{-\lambda_n}) \oplus V_0 \quad (72)$$

where the  $V_0$  term is present only when  $0 \in \text{Spec}(F)$ . If  $A = B \oplus C$ , we let  $P_B(v)$  denote the unique vector in  $B$  corresponding to  $v \in A$ . With this in mind we define

$$\Lambda_\lambda = \{P_{V_\lambda \oplus V_{-\lambda}}(v) | v \in \Lambda\}. \quad (73)$$

Before proceeding, we caution that vector subspaces are generally not closed under a projection operation. For instance if  $A$  ( $B$ ) is the  $x_1$  ( $x_2$ ) axis, then  $A \oplus B$  is  $\mathbb{R}^2$ , and vectors in the subspace on which  $x_1 = 2x_2$  will project to vectors not in this subspace. To show that  $\Lambda_\lambda$  is a Lagrangian subspace of  $V_\lambda \oplus V_{-\lambda}$ , we note that

- $\Lambda_\lambda^\perp \subset \Lambda_\lambda \iff \sigma(v_\lambda, \hat{v}_\lambda) = 0 \forall v_\lambda \in \Lambda_\lambda$  implies  $\hat{v}_\lambda \in \Lambda_\lambda$   
Start with  $\hat{v}_\lambda \in V_\lambda \oplus V_{-\lambda}$  that is not in  $\Lambda_\lambda$ . Note that  $\hat{v}_\lambda \notin \Lambda$ , because if it was, we'd have  $\hat{v}_\lambda \in \Lambda_\lambda$  contrary to hypothesis (note  $P_{V_\lambda \oplus V_{-\lambda}}(\hat{v}_\lambda) = \hat{v}_\lambda$ ). Because  $\Lambda$  is Lagrangian, some  $v \in \Lambda$  exists for which  $\sigma(v, \hat{v}_\lambda) \neq 0$ . By the linearity of  $\sigma$  and C.8,  $\sigma(v_\lambda, \hat{v}_\lambda) \neq 0$ , where  $v_\lambda = P_{V_\lambda \oplus V_{-\lambda}}(v) \in \Lambda_\lambda \subset V_\lambda \oplus V_{-\lambda}$ . Thus we have shown that some  $v_\lambda \in V_\lambda \oplus V_{-\lambda}$  exists in  $\Lambda_\lambda$  for which  $\sigma(v_\lambda, \hat{v}_\lambda) \neq 0$ .
- $\Lambda_\lambda \subset \Lambda_\lambda^\perp \iff \sigma(v_\lambda, \hat{v}_\lambda) = 0 \forall v_\lambda, \hat{v}_\lambda \in \Lambda_\lambda$   
Given  $v_\lambda, \hat{v}_\lambda \in \Lambda_\lambda$ , we know from the definition (73) of  $\Lambda_\lambda$  that vectors  $v$  and  $\hat{v}$  exist in  $\Lambda$  that are projected by  $P_{V_\lambda \oplus V_{-\lambda}}(\bullet)$  to  $v_\lambda$  and  $\hat{v}_\lambda$  respectively. Because  $\Lambda$  is an  $F$ -invariant subspace of  $W$ , we know from C.14 that  $\Lambda$  is given by the direct sum of subspaces  $U_i$  of  $W$ , where each  $U_i$  is a subspace of  $V_{\lambda_i} \oplus V_{-\lambda_i}$ . Of course, because  $\Lambda$  is Lagrangian, these  $U_i$ 's will end up having additional properties, but these are of no concern to us right now. What matters is that any vector in the direct sum of the  $U_i$ 's is in  $\Lambda$ . In particular,  $v_\lambda$  and  $\hat{v}_\lambda$  are in  $\Lambda$ , and so  $\sigma(v_\lambda, \hat{v}_\lambda) = 0$ , which is what we wanted to show.
- $\dim(\Lambda_\lambda) = \frac{1}{2} \dim(V_\lambda \oplus V_{-\lambda})$   
Let  $\{b_i\}$  be a basis for  $\Lambda_\lambda$ , and define the operator  $T_\lambda$  on  $V_\lambda \oplus V_{-\lambda}$  by

$$T_\lambda(x) = b_1 \sigma(x, b_1) + b_2 \sigma(x, b_2) + \cdots + b_n \sigma(x, b_n) \quad (74)$$

Note that  $\dim(V_\lambda \oplus V_{-\lambda}) = \dim(\ker(T_\lambda)) + \dim(\text{range}(T_\lambda))$ , and so we'll be done if we can show that  $\ker(T_\lambda)$  and  $\text{range}(T_\lambda)$  both equal  $\Lambda_\lambda$ .

–  $\ker(T_\lambda) = \Lambda_\lambda$

We've already established that  $\Lambda_\lambda = \Lambda_\lambda^\perp$ , and so if  $x \in \Lambda_\lambda$ , then  $T_\lambda(x) = 0$ . If  $x \notin \Lambda_\lambda$ , then some  $y \in \Lambda_\lambda$  exists for which  $\sigma(x, y) \neq 0$ , and so  $T_\lambda(x) \neq 0$ .

–  $\text{range}(T_\lambda) = \Lambda_\lambda$

We need to show that for every  $b_i$ , there is some  $c_i \in V_\lambda \oplus V_\lambda$  for which  $\sigma(b_i, c_i)$  is nonzero for  $k = i$ , and zero for  $k \neq i$ . We do this by defining the operator  $T$  on  $W$  as the sum of all the  $T_\lambda$ 's. The same reasoning used to establish that  $\ker(T_\lambda) = \Lambda_\lambda$  (above) shows that  $\ker(T) = \Lambda$ . Of course  $\dim(W) = \dim(\ker(T)) + \dim(\text{range}(T))$ . Because  $\Lambda$  is Lagrangian, it follows that  $\dim(\text{range}(T)) = \dim(\Lambda)$ . Note that  $T$  has the same form as  $T_\lambda$ , but with more  $b_i$ 's. These  $b_i$ 's comprise a basis for  $\Lambda$ , and so  $\text{range}(T) \subset \Lambda$ . It follows immediately that  $\text{range}(T) = \Lambda$ , and so it must be that for every  $b_i$ , some  $c \in W$  causes  $\sigma(c, b_k)$  to be zero (nonzero) for  $k = i$  ( $k \neq i$ ). When  $b_i$  is one of the basis elements of  $\Lambda_\lambda$ , we write the corresponding  $c$  as  $c_i + \hat{c}$ , where  $c_i = P_{V_\lambda \oplus V_{-\lambda}}(c)$ . Then from C.8 it follows that  $0 \neq \sigma(c, b_i) = \sigma(c_i, b_i) + \sigma(\hat{c}, b_i) = \sigma(c_i, b_i)$ . Similarly, for any of the other  $b_k$ 's in  $\Lambda_\lambda$ ,  $0 = \sigma(c, b_k) = \sigma(c_i, b_k) + \sigma(\hat{c}, b_k) = \sigma(c_i, b_k)$ , and so we are done.

To show that  $\Lambda_\lambda$  is invariant under  $F$ , we note that  $v_\lambda \in \Lambda_\lambda$  implies the existence of some  $v \in \Lambda$  such that  $P_{V_\lambda \oplus V_{-\lambda}}(v) = v_\lambda$ . We know that  $\Lambda$  is  $F$ -invariant, and so we'll be done if we can show that  $F(P_{V_\lambda \oplus V_{-\lambda}}(v)) = P_{V_\lambda \oplus V_{-\lambda}}(Fv)$ . But this equality is obvious; the  $(V_{\lambda_i} \oplus V_{-\lambda_i})$ 's are  $F$ -invariant. Writing  $v = v_{\lambda_1} + v_{\lambda_2} + \cdots + v_{\lambda_n}$  where  $v_{\lambda_i} \in V_{\lambda_i} \oplus V_{-\lambda_i}$ , we have  $Fv = Fv_{\lambda_1} + Fv_{\lambda_2} + \cdots + Fv_{\lambda_n}$ , where  $Fv_{\lambda_i} \in V_{\lambda_i} \oplus V_{-\lambda_i}$ .

At this point we've constructed subspaces  $\Lambda_\lambda \subset V_\lambda \oplus V_{-\lambda}$  that are Lagrangian. Certainly  $\Lambda$  is a subset of the direct sum of these subspaces. Equality follows from the fact that the dimension of the direct sum equals the dimension of  $\Lambda$ . In detail,

$$\dim(\Lambda_{\lambda_1} \oplus \cdots \oplus \Lambda_{\lambda_n}) = \dim(\Lambda_{\lambda_1}) + \cdots + \dim(\Lambda_{\lambda_n}), \quad (75)$$

where  $\Lambda_{\lambda_i}$  has half the dimension of  $V_{\lambda_i} \oplus V_{-\lambda_i}$ . It follows that the numbers  $\dim(\Lambda_{\lambda_i})$  sum to half the dimension of  $W$ , which is the same as the dimension of  $\Lambda$ .

**Theorem C.16** *When  $\alpha \neq 0$ ,  $V_{\pm\alpha}$  is an  $F$ -invariant Lagrangian subspace of  $V_\alpha \oplus V_{-\alpha}$ .*

**Proof:** The  $F$ -invariance is obvious, as is the dimensionality requirement. To show that the orthogonal complement of  $V_{\pm\alpha}$  in  $V_\alpha \oplus V_{-\alpha}$  is exactly  $V_{\pm\alpha}$ , first note that from C.8,  $x, y \in V_{\pm\alpha}$  implies  $\sigma(x, y) = 0$ . Next, pick some  $x \in V_\alpha \oplus V_{-\alpha}$  that is not in  $V_{\pm\alpha}$ . We need to show that for some  $y \in V_{\pm\alpha}$ ,  $\sigma(x, y) \neq 0$ . If no such  $y$  exists, then  $\sigma$  vanishes identically on  $V_\alpha \oplus V_{-\alpha}$ . However this contradicts C.11, in which we establish that  $V_\alpha \oplus V_{-\alpha}$  is a symplectic vector space under  $\sigma$  (recall from (19) that if  $\sigma(x, y) = 0$  for all  $y \in \tilde{S}$ , then  $x = 0$ , where  $\tilde{S}$  is a symplectic vector space under  $\sigma$ ).  $\square$

**Theorem C.17** *If  $\mu \in \text{Spec}(F)$  has positive real and imaginary parts, then  $V_\mu \oplus V_{\bar{\mu}}$  is spanned by real vectors in  $W$ .*

**Proof:** If  $\{b_i\}$  is a basis for  $V_\mu$ , then every  $b_i$  satisfies  $(F - \mu)^N b_i = 0$  for some  $N$ , and so  $(F - \bar{\mu})^N \bar{b}_i = 0$  as well, showing that the vectors  $\{\bar{b}_i\}$  are all in  $V_{\bar{\mu}}$ . Every linear combination  $\gamma_i \bar{b}_i$  of the  $\bar{b}_i$ 's is nonzero, because if one wasn't, then taking the conjugate would give a nonzero linear combination of the  $b_i$ 's. It follows that there are  $\dim(V_\mu)$  independent vectors in  $V_{\bar{\mu}}$ . If  $\dim(V_{\bar{\mu}})$  was greater than  $\dim(V_\mu)$ , there would be a  $c \in V_{\bar{\mu}}$  independent of the  $\bar{b}_i$ 's. Reusing our previous arguments (this time going from  $V_{\bar{\mu}}$  to  $V_\mu$ ), we find that this  $c$  would be a vector in  $V_\mu$  independent of the  $b_i$ 's. But this is impossible because  $\{b_i\}$  is a basis for  $V_\mu$ , and so it must be that  $\{\bar{b}_i\}$  is a basis for  $V_{\bar{\mu}}$ . From C.12, we know that  $V_\mu \cap V_{\bar{\mu}} = \{0\}$  and so the  $b_i$ 's and  $\bar{b}_i$ 's comprise a basis  $B$  for  $V_\mu \oplus V_{\bar{\mu}}$ . Note that  $\text{span}(\text{Re}(b_i), \text{Im}(b_i)) = \text{span}(b_i, \bar{b}_i)$ . Replacing each conjugate pair  $b_i, \bar{b}_i$  in  $B$  with the real vectors  $\text{Re}(b_i)$  and  $\text{Im}(b_i)$  has no effect on the number of items in  $B$  or their span, and so the real vectors  $\text{Re}(b_i)$  and  $\text{Im}(b_i)$  comprise a basis for  $V_\mu \oplus V_{\bar{\mu}}$ .  $\square$

## D Field Switching

In this appendix we discuss moving back and forth between vector spaces over  $\mathbb{R}$  and  $\mathbb{C}$ .

## D.1 Real Objects in Vector Spaces over $\mathbb{C}$

Let  $V_{\mathbb{C}}$  be an  $n$ -dimensional vector space over  $\mathbb{C}$ . The following equivalent constructions establish the real part, imaginary part, and conjugate of a vector in  $V_{\mathbb{C}}$ . Either of these constructions will be said to endow  $V_{\mathbb{C}}$  with *real content*.

- i. Choose a basis  $\{e_i\}$  of  $V_{\mathbb{C}}$ . If  $v = \sum(a_i + ib_i)e_i$  is in  $V_{\mathbb{C}}$ , then we define the *real part* of  $v$  as  $\sum a_i e_i$ , the *imaginary part* of  $v$  as  $i \sum b_i e_i$ , and the *conjugate* of  $v$  as  $\sum(a_i - ib_i)e_i$ .
- ii. Choose a map  $K : V_{\mathbb{C}} \rightarrow V_{\mathbb{C}}$  with the properties  $K(\lambda v) = \bar{\lambda}K(v)$ ,  $K^2(v) = v$ , and  $K(u + v) = K(u) + K(v)$ . Such a map is referred to as a *real structure*, and serves as a conjugation operator on  $V_{\mathbb{C}}$ . If  $v \in V_{\mathbb{C}}$ , then we define the *real part* of  $v$  as  $\frac{1}{2}(v + K(v))$ , the *imaginary part* of  $v$  as  $\frac{1}{2}(v - K(v))$ , and the *conjugate* of  $v$  as  $K(v)$ .

We use  $\text{Re}(v)$ ,  $\text{Im}(v)$ , and  $\bar{v}$  to denote the real part, imaginary part, and conjugate of  $v$  respectively. A vector with no imaginary part is called *real*, and a vector with no real part is called *imaginary*.

### D.1.1 Real Maps

Let  $U$  and  $V$  be vector spaces over  $\mathbb{C}$  with real content. We call  $A \in L(U, V)$  a *real map* if it maps real vectors in  $U$  to real vectors in  $V$ . If  $A$  is real then the matrix of  $A$  with respect to real bases in  $U$  and  $V$  consists of real numbers.

### D.1.2 Equivalent Definitions

Different bases in i. and real structures in ii. can endow  $V_{\mathbb{C}}$  with the same real content. For instance,

- if  $\{\tilde{e}_i\}$  is a basis of  $V_{\mathbb{C}}$  such that the coefficients of each  $\tilde{e}_i$  with respect to the  $\{e_i\}$  basis from i. are real numbers, then the real content on  $V_{\mathbb{C}}$  established by using  $\{\tilde{e}_i\}$  in i. is the same as the real content on  $V_{\mathbb{C}}$  established by using  $\{e_i\}$ .
- if  $A$  is a real bijective linear operator on  $V_{\mathbb{C}}$  with respect to the real structure  $K$ , then the same real content established by the real structure  $K$  is also established by the real structure  $A^{-1} \circ K \circ A$ .

### D.1.3 Subspaces

In this section we discuss some of the many different conceptions of a *real subspace* of a complex vector space with real content. We use only the first of these in the thesis; the remaining concepts (like *totally real* and its relatives), are included for the sake of contrast. Let  $V_{\mathbb{C}}$  be an  $n$ -dimensional vector space over  $\mathbb{C}$  with real content.

- A subspace  $U \subset V_{\mathbb{C}}$  will be called *real* if it can be given as the  $\mathbb{C}$ -span of real vectors in  $V_{\mathbb{C}}$ , (which is equivalent to requiring that  $\bar{U} = U$ ). Obviously  $V_{\mathbb{C}}$  is a real subspace of itself. Also, we note that subspaces of  $V_{\mathbb{C}}$  exist which are not real. For instance, if  $\{e_i\}$  is a real basis of  $V_{\mathbb{C}}$ , then  $\text{span}_{\mathbb{C}}(\{e_1\})$  and  $\text{span}_{\mathbb{C}}(\{ie_2\})$  are real, but  $\text{span}_{\mathbb{C}}(\{e_1 + ie_2\})$  is not real. One peculiarity of this definition is that a real subspace of  $V_{\mathbb{C}}$  will always contain imaginary vectors, (for instance  $ie_1 \in \text{span}_{\mathbb{C}}(\{e_1\})$ ).

The elements in the  $n$ -dimensional complex vector space  $V_{\mathbb{C}}$  comprise a  $2n$ -dimensional vector space over  $\mathbb{R}$ , which we call  $V_{\mathbb{R}}$ . Unless otherwise noted, a subspace in this discussion inherits the field of its parent space.

- Let  $V_{\mathbb{C}}$  be a vector space over  $\mathbb{C}$  with real content. If  $U$  is a subspace of  $V_{\mathbb{C}}$ , then we define the *real projection*  $\text{Re}(U)$  of  $U$  by

$$\text{Re}(U) = \{\text{Re}(u) | u \in U\}. \quad (76)$$

We note that  $\text{Re}(U)$  is merely a subset (not a subspace) of  $V_{\mathbb{C}}$ . However  $\text{Re}(U)$  is a subspace of the  $2n$ -dimensional  $V_{\mathbb{R}}$ , and  $\text{Re}(U)$  is also a subspace of the  $n$ -dimensional  $\text{span}_{\mathbb{R}}(\{e_i\})$ , where  $\{e_i\}$  is a real basis of  $V_{\mathbb{C}}$ . The inequality

$$\dim(U) \leq \dim(\text{Re}(U)) \leq \min(2 \dim(U), n) \quad (77)$$

holds whether  $\text{Re}(U)$  is regarded as a subspace of  $V_{\mathbb{R}}$  or of  $\text{span}_{\mathbb{R}}(\{e_i\})$ . For instance if the subspace  $U$  of  $V_{\mathbb{C}}$  is given by the (1-dimensional)  $\mathbb{C}$ -span of  $e_1 + ie_2$ , then  $\text{Re}(U)$  is given by the (2-dimensional)  $\mathbb{R}$ -span of  $e_1$  and  $e_2$ .

- A subspace  $U$  of  $V_{\mathbb{R}}$  is called *totally real* if  $U \cap iU = \{0\}$ . We note that complex multiplication is well defined on  $V_{\mathbb{R}}$  even though  $V_{\mathbb{R}}$  is a vector space over  $\mathbb{R}$ . Also,  $V_{\mathbb{C}}$  need not have real content for this definition to make sense.
- We call a subspace of  $V_{\mathbb{R}}$  *entirely real* (*entirely imaginary*) if it can be given as the  $\mathbb{R}$ -span of real (imaginary) vectors in  $V_{\mathbb{C}}$ . We call a subspace of  $V_{\mathbb{R}}$  *entirely mixed* if it can be given as the  $\mathbb{R}$ -span of both real and imaginary vectors in  $V_{\mathbb{C}}$ .

The following examples illustrate the distinction between the above subspaces of  $V_{\mathbb{R}}$ . We also show these relationships in Figure 5.

- if  $U = \text{span}_{\mathbb{R}}(\{e_1\})$ , then  $U$  is entirely real and totally real.
- if  $U = \text{span}_{\mathbb{R}}(\{ie_2\})$ , then  $U$  is entirely imaginary and totally real.
- if  $U = \text{span}_{\mathbb{R}}(\{e_1, ie_2\})$ , then  $U$  is entirely mixed and totally real.
- if  $U = \text{span}_{\mathbb{R}}(\{e_1, ie_1\})$ , then  $U$  is entirely mixed but not totally real, (note that  $U = iU$ ).
- if  $U = \text{span}_{\mathbb{R}}(\{e_1 + ie_1\})$ , or if  $U = \text{span}_{\mathbb{R}}(\{e_1 + ie_2\})$ , then  $U$  is totally real, however  $U$  is neither entirely real, nor entirely imaginary, nor entirely mixed.

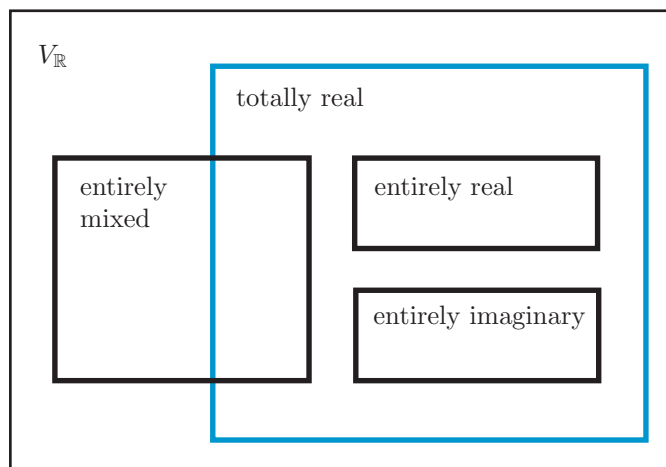


Figure 5: This diagram shows how the different subspaces of  $V_{\mathbb{R}}$  are related.

## D.2 Moving from $\mathbb{R}$ to $\mathbb{C}$

Here we use a vector space over  $\mathbb{R}$  to construct a vector space over  $\mathbb{C}$ . Although this material isn't used in the thesis, it is a natural complement to the previous section, and it can be used to understand some of the ideas presented there. If  $U$  is an  $n$ -dimensional vector space over  $\mathbb{R}$ , we define the *complexification*  $U_{\mathbb{C}}$  of  $U$  to be the vector space  $U \times U$  over  $\mathbb{C}$  with

- addition defined component-wise:  $(u_1, u_2) + (v_1, v_2) := (u_1 + v_1, u_2 + v_2)$ ,
- scalar multiplication defined by:  $(\alpha_r + i\alpha_i) \cdot (u_1, u_2) := (\alpha_r u_1 - \alpha_i u_2, \alpha_r u_2 + \alpha_i u_1)$ .

We call  $u$  the *real part* and  $v$  the *imaginary part* of  $(u, v) \in U_{\mathbb{C}}$ . Vectors in  $U_{\mathbb{C}}$  with zero real (imaginary) parts will be called *imaginary* (*real*), and we call  $(u, -v)$  the *conjugate* of  $(u, v)$ . If  $\{e_k\}$  is a basis for  $U$ , then  $\{(e_k, 0)\}$  is a basis for  $U_{\mathbb{C}}$ , and so  $\dim(U_{\mathbb{C}}) = \dim(U)$  as we might expect.

*Linear Operators:*

It is straightforward to show that  $F$  is a linear operator on  $U_{\mathbb{C}}$  if and only if  $F$  can be written as  $(F_R, F_I)$  where  $F_R$  and  $F_I$  are linear operators on  $U$ , and where the action of  $F$  on an element of  $U_{\mathbb{C}}$  is given by

$$(F_R, F_I)(u_r, u_i) = (F_R u_r - F_I u_i, F_R u_i + F_I u_r) \quad (78)$$

We call  $F_R$  ( $F_I$ ) the *real* (*imaginary*) parts of  $F$ , and we call  $(F_R, 0) \in L(U_{\mathbb{C}})$  the *complexification* of  $F_R \in L(U)$ .

*Inner Products:*

Starting with the inner product  $\langle \bullet, \bullet \rangle$  on  $U$ , we can build a hermitian (i.e., conjugate symmetric) inner product  $\langle \bullet, \bullet \rangle_{\mathbb{C}}$  on  $U_{\mathbb{C}}$  according to

$$\langle (u_r, u_i), (v_r, v_i) \rangle_{\mathbb{C}} := \langle u_r, v_r \rangle + \langle u_i, v_i \rangle + i(\langle u_i, v_r \rangle - \langle u_r, v_i \rangle). \quad (79)$$

*Symplectic Forms:*

Starting with the symplectic form  $\sigma : U \times U \rightarrow \mathbb{R}$ , we can build a symplectic form  $\sigma_{\mathbb{C}} : U_{\mathbb{C}} \times U_{\mathbb{C}} \rightarrow \mathbb{C}$  according to

$$\sigma_{\mathbb{C}}((u_r, u_i), (v_r, v_i)) = \sigma(u_r, v_r) - \sigma(u_i, v_i) + i(\sigma(u_r, v_i) + \sigma(u_i, v_r)). \quad (80)$$

*Symmetric Forms:*

Starting with the symmetric form  $Q : U \times U \rightarrow \mathbb{R}$ , we can build a symmetric form  $Q_{\mathbb{C}} : U_{\mathbb{C}} \times U_{\mathbb{C}} \rightarrow \mathbb{C}$  according to

$$Q_{\mathbb{C}}((u_r, u_i), (v_r, v_i)) = Q(u_r, v_r) - Q(u_i, v_i) + i(Q(u_r, v_i) + Q(u_i, v_r)) \quad (81)$$

## References

- [1] Brian D.O. Anderson and John B. Moore, *Optimal Control*, Prentice Hall, 1989. ISBN 81-203-0697-X.
- [2] Panos J. Antsaklis and Anthony N. Michel, *Linear Systems*, ©2006 Birkhäuser Boston, 2<sup>nd</sup> Corrected Printing, ISBN-10 0-8176-4434-2.
- [3] Sheldon Axler, *Linear Algebra Done Right*, Springer, 2<sup>nd</sup> edition, Corrected 3<sup>rd</sup> printing, 1999.
- [4] Stephen Boyd's, EE363 Notes, Stanford, Winter 2005-06, <http://www.stanford.edu/class/ee363/>
- [5] Gerhard Freiling, Volker Mehrmann, and Hongguo Xu, *Existence, Uniqueness, and Parametrization of Lagrangian Invariant Subspaces*, SIAM J. Matrix Anal. Appl. Vol. 23, No. 4, pp. 1045-1069, 2002.
- [6] L. Hörmander, *The Analysis of Linear Partial Differential Operators, vol. III, IV*, Springer Verlag, 1985.
- [7] Peter Lancaster and Leiba Rodman, *Algebraic Riccati Equations*, Oxford: Clarendon Press; New York: Oxford University Press, 1995.
- [8] Webster's Third New International Dictionary, Unabridged. Merriam-Webster, 2002. <http://unabridged.merriam-webster.com> (4 Aug. 2006).
- [9] Andrew Packard, *ME234 Notes*, <http://jagger.me.berkeley.edu/pack/me234/>
- [10] Annelise Raphael and Maciej Zworski, *Pseudospectral Effects and Basins of Attraction*, unpublished notes, 2005.

## E Acknowledgments

My heartfelt thanks to Maciej Zworski for his patient and kindly guidance throughout this project. I am especially grateful for his insistence that I employ the “natural and invariant” approach in my work. Trudging up this particular mountain seemed pointless until one day I looked around and found myself enjoying a spectacular view. Thanks to Alan Weinstein for taking my thesis seriously. I learned an incredible amount by fixing the mistakes he identified, especially those in Appendix D. I am grateful to Andrew Packard for always having time to sit and talk with me about my project, and for inspiring me by his own hard work and dedication. Finally, my many thanks to Oliver O'Reilly, for encouraging me to pursue the Masters degree in mathematics, and for being patient when it ended up taking a bit longer than expected.