

Linear models

1 Projections

The projection of a vector x onto a subspace \mathcal{S} is the vector in \mathcal{S} that is closest to x . In order to define this rigorously, we start by introducing the concept of direct sum. If two subspaces are disjoint, i.e. their only common point is the origin, then a vector that can be written as a sum of a vector from each subspace is said to belong to their direct sum.

Definition 1.1 (Direct sum). *Let \mathcal{V} be a vector space. For any subspaces $\mathcal{S}_1, \mathcal{S}_2 \subseteq \mathcal{V}$ such that*

$$\mathcal{S}_1 \cap \mathcal{S}_2 = \{0\} \quad (1)$$

the direct sum is defined as

$$\mathcal{S}_1 \oplus \mathcal{S}_2 := \{x \mid x = s_1 + s_2 \quad s_1 \in \mathcal{S}_1, s_2 \in \mathcal{S}_2\}. \quad (2)$$

The representation of a vector in the direct sum of two subspaces is unique.

Lemma 1.2. *Any vector $x \in \mathcal{S}_1 \oplus \mathcal{S}_2$ has a **unique** representation*

$$x = s_1 + s_2 \quad s_1 \in \mathcal{S}_1, s_2 \in \mathcal{S}_2. \quad (3)$$

Proof. If $x \in \mathcal{S}_1 \oplus \mathcal{S}_2$ then by definition there exist $s_1 \in \mathcal{S}_1, s_2 \in \mathcal{S}_2$ such that $x = s_1 + s_2$. Assume $x = s'_1 + s'_2$, $s'_1 \in \mathcal{S}_1, s'_2 \in \mathcal{S}_2$, then $s_1 - s'_1 = s_2 - s'_2$. This implies that $s_1 - s'_1$ and $s_2 - s'_2$ are in \mathcal{S}_1 and also in \mathcal{S}_2 . However, $\mathcal{S}_1 \cap \mathcal{S}_2 = \{0\}$, so we conclude $s_1 = s'_1$ and $s_2 = s'_2$. \square

We can now define the projection of a vector x onto a subspace \mathcal{S} by separating the vector into a component that belongs to \mathcal{S} and another that belongs to its orthogonal complement.

Definition 1.3 (Orthogonal projection). *Let \mathcal{V} be a vector space. The orthogonal projection of a vector $x \in \mathcal{V}$ onto a subspace $\mathcal{S} \subseteq \mathcal{V}$ is a vector denoted by $\mathcal{P}_{\mathcal{S}} x$ such that $x - \mathcal{P}_{\mathcal{S}} x \in \mathcal{S}^{\perp}$.*

Theorem 1.4 (Properties of orthogonal projections). *Let \mathcal{V} be a vector space. Every vector $x \in \mathcal{V}$ has a **unique** orthogonal projection $\mathcal{P}_{\mathcal{S}} x$ onto any subspace $\mathcal{S} \subseteq \mathcal{V}$ of finite dimension. In particular x can be expressed as*

$$x = \mathcal{P}_{\mathcal{S}} x + \mathcal{P}_{\mathcal{S}^{\perp}} x. \quad (4)$$

For any vector $s \in \mathcal{S}$

$$\langle x, s \rangle = \langle \mathcal{P}_{\mathcal{S}} x, s \rangle. \quad (5)$$

For any orthonormal basis b_1, \dots, b_m of \mathcal{S} ,

$$\mathcal{P}_{\mathcal{S}} x = \sum_{i=1}^m \langle x, b_i \rangle b_i. \quad (6)$$

Proof. Since \mathcal{V} has finite dimension, so does \mathcal{S} , which consequently has an orthonormal basis with finite cardinality b'_1, \dots, b'_m by Theorem 3.7 in Lecture Notes 8. Consider the vector

$$p := \sum_{i=1}^m \langle x, b'_i \rangle b'_i. \quad (7)$$

It turns out that $x - p$ is orthogonal to every vector in the basis. For $1 \leq j \leq m$,

$$\langle x - p, b'_j \rangle = \left\langle x - \sum_{i=1}^m \langle x, b'_i \rangle b'_i, b'_j \right\rangle \quad (8)$$

$$= \langle x, b'_j \rangle - \sum_{i=1}^m \langle x, b'_i \rangle \langle b'_i, b'_j \rangle \quad (9)$$

$$= \langle x, b'_j \rangle - \langle x, b'_j \rangle = 0, \quad (10)$$

so by Lemma 3.3 in Lecture Notes 8 $x - p \in \mathcal{S}^\perp$ and p is an orthogonal projection. Since $\mathcal{S} \cap \mathcal{S}^\perp = \{0\}$ ¹ there cannot be two other vectors $x_1 \in \mathcal{S}, x_2 \in \mathcal{S}^\perp$ such that $x = x_1 + x_2$ so the orthogonal projection is unique.

Notice that $o := x - p$ is a vector in \mathcal{S}^\perp such that $x - o = p$ is in \mathcal{S} and therefore in $(\mathcal{S}^\perp)^\perp$. This implies that o is the orthogonal projection of x onto \mathcal{S}^\perp and establishes (4).

Equation (5) follows immediately from the orthogonality of any vector $s \in \mathcal{S}$ and $\mathcal{P}_{\mathcal{S}} x$.

Equation (6) follows from (5) and Lemma 3.5 in Lecture Notes 8. \square

Computing the norm of the projection of a vector onto a subspace is easy if we have access to an orthonormal basis (as long as the norm is induced by the inner product).

Lemma 1.5 (Norm of the projection). *The norm of the projection of an arbitrary vector $x \in \mathcal{V}$ onto a subspace $\mathcal{S} \subseteq \mathcal{V}$ of dimension d can be written as*

$$\|\mathcal{P}_{\mathcal{S}} x\|_{\langle \cdot, \cdot \rangle} = \sqrt{\sum_{i=1}^d \langle b_i, x \rangle^2} \quad (11)$$

for any orthonormal basis b_1, \dots, b_d of \mathcal{S} .

¹For any vector v that belongs to both \mathcal{S} and \mathcal{S}^\perp $\langle v, v \rangle = \|v\|_2^2 = 0$, which implies $v = 0$.

Proof. By (6)

$$\|\mathcal{P}_{\mathcal{S}} x\|_{\langle \cdot, \cdot \rangle}^2 = \langle \mathcal{P}_{\mathcal{S}} x, \mathcal{P}_{\mathcal{S}} x \rangle \quad (12)$$

$$= \left\langle \sum_i^d \langle b_i, x \rangle b_i, \sum_j^d \langle b_j, x \rangle b_j \right\rangle \quad (13)$$

$$= \sum_i^d \sum_j^d \langle b_i, x \rangle \langle b_j, x \rangle \langle b_i, b_j \rangle \quad (14)$$

$$= \sum_i^d \langle b_i, x \rangle^2. \quad (15)$$

□

Finally, we prove indeed that of a vector x onto a subspace \mathcal{S} is the vector in \mathcal{S} that is closest to x in the distance induced by the inner-product norm.

Example 1.6 (Projection onto a one-dimensional subspace). To compute the projection of a vector x onto a one-dimensional subspace spanned by a vector v , we use the fact that $\left\{ v / \|v\|_{\langle \cdot, \cdot \rangle} \right\}$ is a basis for $\text{span}(v)$ (it is a set containing a unit vector that spans the subspace) and apply (6) to obtain

$$\mathcal{P}_{\text{span}(v)} x = \frac{\langle v, x \rangle}{\|v\|_{\langle \cdot, \cdot \rangle}^2} v. \quad (16)$$

Theorem 1.7 (The orthogonal projection is closest). *The orthogonal projection of a vector x onto a subspace \mathcal{S} belonging to the same inner-product space is the closest vector to x that belongs to \mathcal{S} in terms of the norm induced by the inner product. More formally, $\mathcal{P}_{\mathcal{S}} x$ is the solution to the optimization problem*

$$\underset{u}{\text{minimize}} \quad \|x - u\|_{\langle \cdot, \cdot \rangle} \quad (17)$$

$$\text{subject to} \quad u \in \mathcal{S}. \quad (18)$$

Proof. Take any point $s \in \mathcal{S}$ such that $s \neq \mathcal{P}_{\mathcal{S}} x$

$$\|x - s\|_{\langle \cdot, \cdot \rangle}^2 = \|\mathcal{P}_{\mathcal{S}^\perp} x + \mathcal{P}_{\mathcal{S}} x - s\|_{\langle \cdot, \cdot \rangle}^2 \quad (19)$$

$$= \|\mathcal{P}_{\mathcal{S}^\perp} x\|_{\langle \cdot, \cdot \rangle}^2 + \|\mathcal{P}_{\mathcal{S}} x - s\|_{\langle \cdot, \cdot \rangle}^2 \quad (20)$$

$$> 0 \quad \text{because } s \neq \mathcal{P}_{\mathcal{S}} x, \quad (21)$$

where (20) follows from the Pythagorean theorem since because $\mathcal{P}_{\mathcal{S}^\perp} x$ belongs to \mathcal{S}^\perp and $\mathcal{P}_{\mathcal{S}} x - s$ to \mathcal{S} . □

2 Linear minimum-MSE estimation

We are interested in estimating the sample of a continuous random variable X from the sample y of a random variable Y . If we know the joint distribution of X and Y then the optimal estimator in terms of MSE is the conditional mean $E(X|Y = y)$. However often it is very challenging to completely characterize a joint distribution between two quantities, but it is more tractable to obtain an estimate of their first and second order moments. In this case it turns out that we can obtain the optimal *linear* estimate of X given Y by using our knowledge of linear algebra.

Theorem 2.1 (Best linear estimator). *Assume that we know the means μ_X, μ_Y and variances σ_X^2, σ_Y^2 of two random variables X and Y and their correlation coefficient ρ_{XY} . The best linear estimate of the form $aY + b$ of X given Y in terms of mean-square error is*

$$g_{LMMSE}(y) = \frac{\rho_{XY} \sigma_X (y - \mu_Y)}{\sigma_Y} + \mu_X. \quad (22)$$

Proof. First we determine the value of b . The cost function as a function of b is

$$h(b) = E((X - aY - b)^2) = E((X - aY)^2) + b^2 - 2bE(X - aY) \quad (23)$$

$$= E((X - aY)^2) + b^2 - 2b(\mu_X - a\mu_Y). \quad (24)$$

The first and second derivative with respect to b are

$$h'(b) = 2b - 2(\mu_X - a\mu_Y), \quad (25)$$

$$h''(b) = 2. \quad (26)$$

Since h'' is positive the function is convex so the minimum is obtained by setting $h'(b)$ to zero, which yields

$$b = \mu_X - a\mu_Y. \quad (27)$$

Consider the centered random variables $\tilde{X} := X - \mu_X$ and $\tilde{Y} := Y - \mu_Y$. It turns out that to find a we just need to find the best estimate of \tilde{X} of the form $a\tilde{Y}$

$$E((X - aY - b)^2) = E((X - \mu_X - a(Y - \mu_Y) + \mu_X - a\mu_Y - b))^2 \quad (28)$$

$$= E\left(\left(\tilde{X} - a\tilde{Y}\right)^2\right), \quad (29)$$

clearly any a that minimizes the left-hand side also minimizes the right-hand side and vice versa.

Consider the vector space of zero-mean random variables. \tilde{X} and \tilde{Y} belong to this vector space. In fact,

$$\langle \tilde{X}, \tilde{Y} \rangle = E(\tilde{X}\tilde{Y}) \quad (30)$$

$$= \text{Cov}(X, Y) \quad (31)$$

$$= \sigma_X \sigma_Y \rho_{XY}, \quad (32)$$

$$\|\tilde{X}\|_{\langle \cdot, \cdot \rangle}^2 = E(\tilde{X}^2) \quad (33)$$

$$= \sigma_X^2, \quad (34)$$

$$\|\tilde{Y}\|_{\langle \cdot, \cdot \rangle}^2 = E(\tilde{Y}^2) \quad (35)$$

$$= \sigma_Y^2. \quad (36)$$

Any random variable of the form $a\tilde{Y}$ belongs to the subspace spanned by \tilde{Y} . Since the distance in this vector space is induced by the mean-square norm, by Theorem 1.7 the vector of the form $a\tilde{Y}$ that approximates \tilde{X} better is just the projection of \tilde{X} onto the subspace spanned by \tilde{Y} , which we will denote by $\mathcal{P}_{\tilde{Y}}\tilde{X}$. This subspace has dimension 1, so $\{\tilde{Y}/\sigma_Y\}$ is a basis for the subspace. The projection is consequently equal to

$$\mathcal{P}_{\tilde{Y}}\tilde{X} = \left\langle \tilde{X}, \frac{\tilde{Y}}{\sigma_Y} \right\rangle \frac{\tilde{Y}}{\sigma_Y} \quad (37)$$

$$= \langle \tilde{X}, \tilde{Y} \rangle \frac{\tilde{Y}}{\sigma_Y^2} \quad (38)$$

$$= \frac{\sigma_X \rho_{XY} \tilde{Y}}{\sigma_Y}. \quad (39)$$

So $a = \sigma_X \rho_{XY} / \sigma_Y$, which concludes the proof. \square

In words, the linear estimator of X given Y is obtained by

1. centering Y by removing its mean,
2. normalizing \tilde{Y} by dividing by its standard deviation,
3. scaling the result using the correlation between Y and X ,
4. scaling again using the standard deviation of X ,
5. recentering by adding the mean of X .

3 Matrices

A **matrix** is a rectangular array of numbers. We denote the vector space of $m \times n$ matrices by $\mathbb{R}^{m \times n}$. We denote the i th row of a matrix A by $A_{i,:}$, the j th column by $A_{:,j}$ and the (i, j) entry by A_{ij} . The transpose of a matrix is obtained by switching its rows and columns.

Definition 3.1 (Transpose). *The **transpose** A^T of a matrix $A \in \mathbb{R}^{m \times n}$ is a matrix in $\mathbb{R}^{n \times m}$*

$$(A^T)_{ij} = A_{ji}. \quad (40)$$

A **symmetric** matrix is a matrix that is equal to its transpose.

Matrices map vectors to other vectors through a linear operation called matrix-vector product.

Definition 3.2 (Matrix-vector product). *The product of a matrix $A \in \mathbb{R}^{m \times n}$ and a vector $x \in \mathbb{R}^n$ is a vector in \mathbb{R}^m , such that*

$$(Ax)_i = \sum_{j=1}^n A_{ij} x[j] \quad (41)$$

$$= \langle A_{i,:}, x \rangle, \quad (42)$$

i.e. the i th entry of Ax is the dot product between the i th row of A and x .

Equivalently,

$$Ax = \sum_{j=1}^n A_{:,j} x[j], \quad (43)$$

i.e. Ax is a linear combination of the columns of A weighted by the entries in x .

One can easily check that the transpose of the product of two matrices A and B is equal to the the transposes multiplied in the inverse order,

$$(AB)^T = B^T A^T. \quad (44)$$

We can now we can express the dot product between two vectors x and y as

$$\langle x, y \rangle = x^T y = y^T x. \quad (45)$$

The identity matrix is a matrix that maps any vector to itself.

Definition 3.3 (Identity matrix). *The identity matrix in $\mathbb{R}^{n \times n}$ is*

$$I = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ & & \ddots & \\ 0 & 0 & \cdots & 1 \end{bmatrix}. \quad (46)$$

Clearly, for any $x \in \mathbb{R}^n$ $Ix = x$.

Definition 3.4 (Matrix multiplication). *The product of two matrices $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{n \times p}$ is a matrix $AB \in \mathbb{R}^{m \times p}$, such that*

$$(AB)_{ij} = \sum_{k=1}^n A_{ik} B_{kj} = \langle A_{i,:}, B_{:,j} \rangle, \quad (47)$$

i.e. the (i, j) entry of AB is the dot product between the i th row of A and the j th column of B .

Equivalently, the j th column of AB is the result of multiplying A and the j th column of B

$$AB = \sum_{k=1}^n A_{ik} B_{kj} = \langle A_{i,:}, B_{:,j} \rangle, \quad (48)$$

and i th row of AB is the result of multiplying the i th row of A and B .

Square matrices may have an inverse. If they do, the inverse is a matrix that reverses the effect of the matrix of any vector.

Definition 3.5 (Matrix inverse). *The inverse of a square matrix $A \in \mathbb{R}^{n \times n}$ is a matrix $A^{-1} \in \mathbb{R}^{n \times n}$ such that*

$$AA^{-1} = A^{-1}A = I. \quad (49)$$

Lemma 3.6. *The inverse of a matrix is unique.*

Proof. Let us assume there is another matrix M such that $AM = I$, then

$$M = A^{-1}AM \quad \text{by (49)} \quad (50)$$

$$= A^{-1}. \quad (51)$$

□

An important class of matrices are **orthogonal matrices**.

Definition 3.7 (Orthogonal matrix). *An orthogonal matrix is a square matrix such that its inverse is equal to its transpose,*

$$U^T U = U U^T = I \quad (52)$$

By definition, the columns $U_{:1}, U_{:2}, \dots, U_{:n}$ of any orthogonal matrix have unit norm and orthogonal to each other, so they form an orthonormal basis (it's somewhat confusing that orthogonal matrices are not called orthonormal matrices instead). We can interpret applying U^T to a vector x as computing the coefficients of its representation in the basis formed by the columns of U . Applying U to $U^T x$ recovers x by scaling each basis vector with the corresponding coefficient:

$$x = U U^T x = \sum_{i=1}^n \langle U_{:i}, x \rangle U_{:i}. \quad (53)$$

Applying an orthogonal matrix to a vector does not affect its norm, it just rotates the vector.

Lemma 3.8 (Orthogonal matrices preserve the norm). *For any orthogonal matrix $U \in \mathbb{R}^{n \times n}$ and any vector $x \in \mathbb{R}^n$,*

$$\|Ux\|_2 = \|x\|_2. \quad (54)$$

Proof. By the definition of an orthogonal matrix

$$\|Ux\|_2^2 = x^T U^T U x \quad (55)$$

$$= x^T x \quad (56)$$

$$= \|x\|_2^2. \quad (57)$$

□

4 Eigendecomposition

An **eigenvector** v of A satisfies

$$Av = \lambda v \quad (58)$$

for a real number λ which is the corresponding **eigenvalue**. Even if A is real, its eigenvectors and eigenvalues can be complex.

Lemma 4.1 (Eigendecomposition). *If a square matrix $A \in \mathbb{R}^{n \times n}$ has n linearly independent eigenvectors v_1, \dots, v_n with eigenvalues $\lambda_1, \dots, \lambda_n$ it can be expressed in terms of a matrix Q , whose columns are the eigenvectors, and a diagonal matrix containing the eigenvalues,*

$$A = \begin{bmatrix} v_1 & v_2 & \cdots & v_n \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ & & \ddots & \\ 0 & 0 & \cdots & \lambda_n \end{bmatrix} \begin{bmatrix} v_1 & v_2 & \cdots & v_n \end{bmatrix}^{-1} \quad (59)$$

$$= Q\Lambda Q^{-1} \quad (60)$$

Proof.

$$AQ = \begin{bmatrix} Av_1 & Av_2 & \cdots & Av_n \end{bmatrix} \quad (61)$$

$$= \begin{bmatrix} \lambda_1 v_1 & \lambda_2 v_2 & \cdots & \lambda_n v_n \end{bmatrix} \quad (62)$$

$$= Q\Lambda. \quad (63)$$

As we will establish later on, if the columns of a square matrix are all linearly independent, then the matrix has an inverse, so multiplying the expression by Q^{-1} on both sides completes the proof. \square

Lemma 4.2. *Not all matrices have an eigendecomposition*

Proof. Consider for example the matrix

$$\begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}. \quad (64)$$

Assume λ has a nonzero eigenvalue corresponding to an eigenvector with entries v_1 and v_2 , then

$$\begin{bmatrix} v_2 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} \lambda v_1 \\ \lambda v_2 \end{bmatrix}, \quad (65)$$

which implies that $v_2 = 0$ and hence $v_1 = 0$, since we have assumed that $\lambda \neq 0$. This implies that the matrix does not have nonzero eigenvalues associated to nonzero eigenvectors. \square

An interesting use of the eigendecomposition is computing successive matrix products very fast. Assume that we want to compute

$$AA \cdots Ax = A^k x, \quad (66)$$

i.e. we want to apply A to x k times. A^k *cannot* be computed by taking the power of its entries (try out a simple example to convince yourself). However, if A has an eigendecomposition,

$$A^k = Q\Lambda Q^{-1}Q\Lambda Q^{-1}\dots Q\Lambda Q^{-1} \quad (67)$$

$$= Q\Lambda^k Q^{-1} \quad (68)$$

$$= Q \begin{bmatrix} \lambda_1^k & 0 & \dots & 0 \\ 0 & \lambda_2^k & \dots & 0 \\ & & \dots & \\ 0 & 0 & \dots & \lambda_n^k \end{bmatrix} Q^{-1}, \quad (69)$$

using the fact that for diagonal matrices applying the matrix repeatedly is equivalent to taking the power of the diagonal entries. This allows to compute the k matrix products using just 3 matrix products and taking the power of n numbers.

From high-school or undergraduate algebra you probably remember how to compute eigenvectors using determinants. In practice, this is usually not a viable option due to stability issues. A popular technique to compute eigenvectors is based on the following insight. Let $A \in \mathbb{R}^{n \times n}$ be a matrix with eigendecomposition $Q\Lambda Q^{-1}$ and let x be an arbitrary vector in \mathbb{R}^n . Since the columns of Q are linearly independent, they form a basis for \mathbb{R}^n , so we can represent x as

$$x = \sum_{i=1}^n \alpha_i Q_{:i}, \quad \alpha_i \in \mathbb{R}, \quad 1 \leq i \leq n. \quad (70)$$

Now let us apply A to x k times,

$$A^k x = \sum_{i=1}^n \alpha_i A^k Q_{:i} \quad (71)$$

$$= \sum_{i=1}^n \alpha_i \lambda_i^k Q_{:i}. \quad (72)$$

If we assume that the eigenvectors are ordered according to their magnitudes and that the magnitude of one of them is larger than the rest, $|\lambda_1| > |\lambda_2| \geq \dots$, and that $\alpha_1 \neq 0$ (which happens with high probability if we draw a random x) then as k grows larger the term $\alpha_1 \lambda_1^k Q_{:1}$ dominates. The term will blow up or tend to zero unless we normalize every time before applying A . Adding the normalization step to this procedure results in the **power method** or power iteration, an algorithm of great importance in numerical linear algebra.

Algorithm 4.3 (Power method).

Input: A matrix A .

Output: An estimate of the eigenvector of A corresponding to the largest eigenvalue.

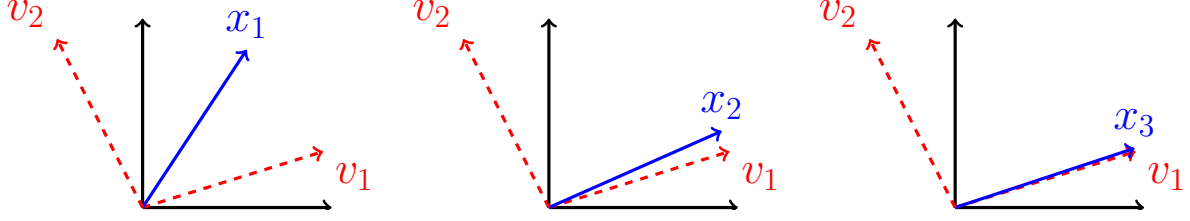


Figure 1: Illustration of the first three iterations of the power method for a matrix with eigenvectors v_1 and v_2 , whose corresponding eigenvalues are $\lambda_1 = 1.05$ and $\lambda_2 = 0.1661$.

Initialization: Set $x_1 := x / \|x\|_2$, where x contains random entries.
For $i = 1, \dots, k$, compute

$$x_i := \frac{Ax_{i-1}}{\|Ax_{i-1}\|_2}. \quad (73)$$

Figure 1 illustrates the power method on a simple example, where the matrix— which was just drawn at random— is equal to

$$A = \begin{bmatrix} 0.930 & 0.388 \\ 0.237 & 0.286 \end{bmatrix}. \quad (74)$$

The convergence to the eigenvector corresponding to the eigenvalue with the largest magnitude is very fast.

5 Time-homogeneous Markov chains

A **Markov chain** is a sequence of discrete random variables X_0, X_1, \dots such that

$$p_{X_{k+1}|X_0, X_1, \dots, X_k}(x_{k+1}|x_0, x_1, \dots, x_k) = p_{X_{k+1}|X_k}(x_{k+1}|x_k). \quad (75)$$

In words, X_{k+1} is conditionally independent of X_j for $j \leq k-1$ conditioned on X_k . If the value of the random variables is restricted to a finite set $\{\alpha_1, \dots, \alpha_n\}$ with probability one, the Markov chain is said to be **time homogeneous**. More formally,

$$P_{ij} := p_{X_{k+1}|X_k}(\alpha_i|\alpha_j) \quad (76)$$

only depends on i and j , **not** on k for all $1 \leq i, j \leq n, k \geq 0$.

If a Markov chain is time homogeneous we can group the transition probabilities P_{ij} in a **transition matrix** P . We express the pmf of X_k restricted to the set where it can be

nonzero as a vector,

$$\pi_k = \begin{bmatrix} p_{X_k}(\alpha_1) \\ p_{X_k}(\alpha_2) \\ \dots \\ p_{X_k}(\alpha_n) \end{bmatrix}. \quad (77)$$

By the Chain Rule, the pmf of X_k can be computed from the pmf of X_0 using the transition matrix,

$$\pi_k = P P \cdots P \pi_0 = P^k \pi_0. \quad (78)$$

In some cases, no matter how we initialize the Markov chain, the Markov Chain *forgets* its initial state and converges to a **stationary distribution**. This is exploited in Markov-Chain Monte Carlo methods that allow to sample from arbitrary distributions by building the corresponding Markov chain. These methods are very useful in Bayesian statistics.

A Markov Chain that converges to a stationary distribution π_∞ is said to be **ergodic**. Note that necessarily $P\pi_\infty = \pi_\infty$, so that π_∞ is an eigenvector of the transition matrix with a corresponding eigenvalue equal to one.

Conversely, let a transition matrix P of a Markov chain have a valid eigendecomposition with n linearly independent eigenvectors v_1, v_2, \dots and corresponding eigenvalues $\lambda_1 > \lambda_2 \geq \lambda_3 \dots$. If the eigenvector corresponding to the largest eigenvalue has non-negative entries then

$$\pi_\infty := \frac{v_1}{\sum_{i=1}^n v_1(i)} \quad (79)$$

is a valid pmf and

$$P\pi_\infty = \lambda_1 \pi_\infty \quad (80)$$

is also a valid pmf, which is only possible if the largest eigenvalue λ_1 equals one. Now, if we represent any possible initial pmf π_0 in the basis formed by the eigenvectors of P we have

$$\pi_0 = \sum_{i=1}^n \alpha_i v_i, \quad \alpha_i \in \mathbb{R}, \quad 1 \leq i \leq n, \quad (81)$$

and

$$\pi_k = P^k \pi_0 \quad (82)$$

$$= \sum_{i=1}^n \alpha_i P^k v_i \quad (83)$$

$$= \alpha_1 \lambda_1 v_1 + \sum_{i=2}^n \alpha_i \lambda_i^k v_i. \quad (84)$$

Since the rest of eigenvalues are strictly smaller than one,

$$\lim_{k \rightarrow \infty} \pi_k = \alpha_1 \lambda_1 v_1 = \pi_\infty \quad (85)$$

where the last equality follows from the fact that the sequence of π_k all belong to the closed set $\{\pi \mid \sum_i^n \pi(i) = 1\}$ so the limit also belongs to the set and hence is a valid pmf. We refer the interested reader to more advanced texts treating Markov chains for further details.

6 Eigendecomposition of symmetric matrices

The following lemma shows that eigenvectors of a symmetric matrix corresponding to different nonzero eigenvalues are necessarily orthogonal.

Lemma 6.1. *If $A \in \mathbb{R}^{n \times n}$ is symmetric, then if u_i and u_j are eigenvectors of A corresponding to different nonzero eigenvalues $\lambda_i \neq \lambda_j \neq 0$*

$$u_i^T u_j = 0. \quad (86)$$

Proof. Since $A = A^T$

$$u_i^T u_j = \frac{1}{\lambda_i} (Au_i)^T u_j \quad (87)$$

$$= \frac{1}{\lambda_i} u_i^T A^T u_j \quad (88)$$

$$= \frac{1}{\lambda_i} u_i^T A u_j \quad (89)$$

$$= \frac{\lambda_j}{\lambda_i} u_i^T u_j. \quad (90)$$

This is only possible if $u_i^T u_j = 0$. □

It turns out that every $n \times n$ symmetric matrix has n linearly independent vectors. The proof of this is beyond the scope of these notes. An important consequence is that **all** symmetric matrices have an eigendecomposition of the form

$$A = U D U^T \quad (91)$$

where $U = [u_1 \ u_2 \ \cdots \ u_n]$ is an orthogonal matrix.

The eigenvalues of a symmetric matrix $\lambda_1, \lambda_2, \dots, \lambda_n$ can be positive, negative or zero. They determine the value of the **quadratic form**:

$$q(x) := x^T A x = \sum_{i=1}^n \lambda_i (x^T u_i)^2 \quad (92)$$

If we order the eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$ then the first eigenvalue is the maximum value attained by the quadratic if its input has unit ℓ_2 norm, the second eigenvalue is the maximum value attained by the quadratic form if we restrict its argument to be normalized and orthogonal to the first eigenvector, and so on.

Theorem 6.2. *For any symmetric matrix $A \in \mathbb{R}^n$ with normalized eigenvectors u_1, u_2, \dots, u_n with corresponding eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_n$*

$$\lambda_1 = \max_{\|u\|_2=1} u^T A u, \quad (93)$$

$$u_1 = \arg \max_{\|u\|_2=1} u^T A u, \quad (94)$$

$$\lambda_k = \max_{\|u\|_2=1, u \perp u_1, \dots, u_{k-1}} u^T A u, \quad (95)$$

$$u_k = \arg \max_{\|u\|_2=1, u \perp u_1, \dots, u_{k-1}} u^T A u. \quad (96)$$

The theorem is proved in Section A.1 of the appendix.

7 The singular-value decomposition

If we consider the columns and rows of a matrix as sets of vectors then we can study their respective spans.

Definition 7.1 (Row and column space). *The **row space** $\text{row}(A)$ of a matrix A is the span of its rows. The **column space** $\text{col}(A)$ is the span of its columns.*

It turns out that the row space and the column space of any matrix have the same dimension. We name this quantity the **rank** of the matrix.

Theorem 7.2. *The rank is well defined;*

$$\dim(\text{col}(A)) = \dim(\text{row}(A)). \quad (97)$$

Section A.2 of the appendix contains the proof.

The following theorem states that we can decompose any real matrix into the product of orthogonal matrices containing bases for its row and column space and a diagonal matrix with a positive diagonal. It is a fundamental result in linear algebra, but its proof is beyond the scope of these notes.

Theorem 7.3. *Without loss of generality let $m \leq n$. Every rank r real matrix $A \in \mathbb{R}^{m \times n}$ has a unique singular-value decomposition of the form (SVD)*

$$A = \begin{bmatrix} u_1 & u_2 & \cdots & u_m \end{bmatrix} \begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 \\ & & \ddots & \\ 0 & 0 & \cdots & \sigma_m \end{bmatrix} \begin{bmatrix} v_1^T \\ v_2^T \\ \vdots \\ v_m^T \end{bmatrix} \quad (98)$$

$$= USV^T, \quad (99)$$

where the **singular values** $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_m \geq 0$ are nonnegative real numbers, the matrix $U \in \mathbb{R}^{m \times m}$ containing the **left singular vectors** is orthogonal, and the matrix $V \in \mathbb{R}^{m \times n}$ containing the **right singular vectors** is a submatrix of an orthogonal matrix (i.e. its columns form an orthonormal set).

Note that we can write the matrix as a sum of rank-1 matrices

$$A = \sum_{i=1}^m \sigma_i u_i v_i^T \quad (100)$$

$$= \sum_{i=1}^r \sigma_i u_i v_i^T, \quad (101)$$

where r is the number of nonzero singular values. The first r left singular vectors $u_1, u_2, \dots, u_r \in \mathbb{R}^m$ form an orthonormal basis of the column space of A and the first r right singular vectors $v_1, v_2, \dots, v_r \in \mathbb{R}^n$ form an orthonormal basis of the row space of A . Therefore the rank of the matrix is equal to r .

8 Principal component analysis

The goal of **dimensionality-reduction** methods is to project high-dimensional data onto a lower-dimensional space while preserving as much information as possible. These methods are a basic tool in data analysis; some applications include visualization (especially if we project onto \mathbb{R}^2 or \mathbb{R}^3), denoising and increasing computational efficiency. **Principal component analysis** (PCA) is a *linear* dimensionality-reduction technique based on the SVD.

If we interpret a set of data vectors as belonging to an ambient vector space, applying PCA allows to find directions in this space along which the data have a high variation. This is achieved by centering the data and then extracting the singular vectors corresponding to the largest singular values. The next two sections provide a geometric and a probabilistic justification.

Algorithm 8.1 (Principal component analysis).

Input: n data vectors $\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n \in \mathbb{R}^m$, a number $k \leq \min\{m, n\}$.

Output: The first k **principal components**, a set of orthonormal vectors of dimension m .

1. Center the data. Compute

$$x_i = \tilde{x}_i - \frac{1}{n} \sum_{i=1}^n \tilde{x}_i, \quad (102)$$

for $1 \leq i \leq n$.

2. Group the centered data in a data matrix X

$$X = \begin{bmatrix} x_1 & x_2 & \cdots & x_n \end{bmatrix}. \quad (103)$$

Compute the SVD of X and extract the left singular vectors corresponding to the k largest singular values. These are the first k principal components.

8.1 PCA: Geometric interpretation

Once the data are centered, the energy of the projection of the data points onto different directions in the ambient space reflects the variation of the dataset along those directions. PCA selects the directions that maximize the ℓ_2 norm of the projection and are mutually orthogonal. By Lemma 1.5, the sum of the squared ℓ_2 norms of the projection of the centered data x_1, x_2, \dots, x_n onto a 1D subspace spanned by a unit-norm vector u can be expressed as

$$\sum_{i=1}^n \|\mathcal{P}_{\text{span}(u)} x_i\|_2^2 = \sum_{i=1}^n u^T x_i x_i^T u \quad (104)$$

$$= u^T X X^T u \quad (105)$$

$$= \|X^T u\|_2^2. \quad (106)$$

If we want to maximize the energy of the projection onto a subspace of dimension k , an option is to choose orthogonal 1D projections sequentially. First we choose a unit vector u_1 that maximizes $\|X^T u\|_2^2$ and is consequently the 1D subspace that is better adapted to the data. Then, we choose a second unit vector u *orthogonal* to the first which maximizes $\|X^T u\|_2^2$ and hence is the 1D subspace that is better adapted to the data while being in the orthogonal complement of u_1 . We repeat this procedure until we have k orthogonal directions. This is exactly equivalent to performing PCA, as proved in Theorem 8.2 below. The k directions correspond to the first k principal components. Figure 2 provides an example in 2D. Note how the singular values are proportional to the energy that lies in the direction of the corresponding principal component.

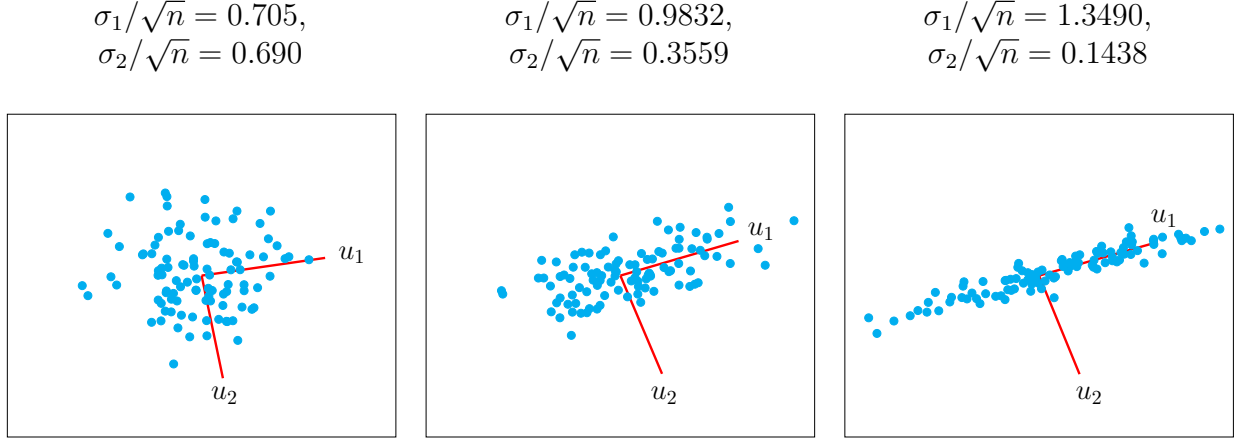


Figure 2: PCA of a dataset with $n = 100$ 2D vectors with different configurations. The two first singular values reflect how much energy is preserved by projecting onto the two first principal components.

Theorem 8.2. For any matrix $X \in \mathbb{R}^{m \times n}$, where $n > m$, with left singular vectors u_1, u_2, \dots, u_m corresponding to the nonzero singular values $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_m$,

$$\sigma_1 = \max_{\|u\|_2=1} \|X^T u\|_2, \quad (107)$$

$$u_1 = \arg \max_{\|u\|_2=1} \|X^T u\|_2, \quad (108)$$

$$\sigma_k = \max_{\substack{\|u\|_2=1 \\ u \perp u_1, \dots, u_{k-1}}} \|X^T u\|_2, \quad 2 \leq k \leq r, \quad (109)$$

$$u_k = \arg \max_{\substack{\|u\|_2=1 \\ u \perp u_1, \dots, u_{k-1}}} \|X^T u\|_2, \quad 2 \leq k \leq r. \quad (110)$$

Proof. If the SVD of X is USV^T then the eigendecomposition of XX^T is equal to

$$XX^T = USV^T V S U^T = U S^2 U^T, \quad (111)$$

where $V^T V = I$ because $n > m$ and the matrix has m nonzero singular values. S^2 is a diagonal matrix containing $\sigma_1^2 \geq \sigma_2^2 \geq \dots \geq \sigma_m^2$ in its diagonal.

The result now follows from applying Theorem 6.2 to the quadratic form

$$u X X^T u = \|X^T u\|_2^2. \quad (112)$$

□

This result shows that PCA is equivalent to choosing the *best* (in terms of ℓ_2 norm) k 1D subspaces following a *greedy* procedure, since at each step we choose the best 1D subspace

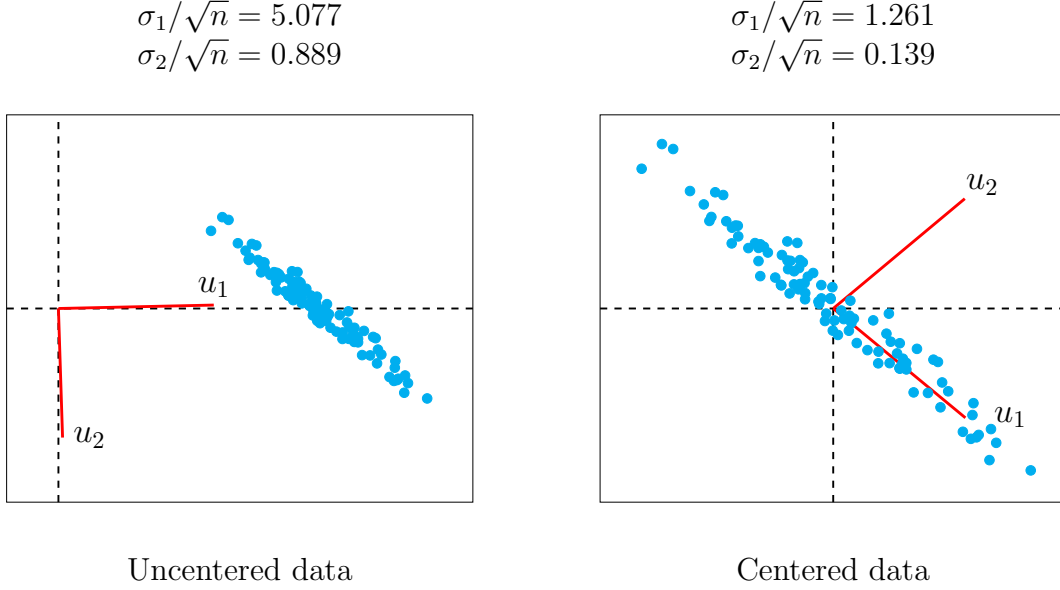


Figure 3: PCA applied to $n = 100$ 2D data points. On the left the data are not centered. As a result the dominant principal component u_1 lies in the direction of the mean of the data and PCA does not reflect the actual structure. Once we center, u_1 becomes aligned with the direction of maximal variation.

orthogonal to the previous ones. A natural question to ask is whether this method produces the best k -dimensional subspace. A priori this is not necessarily the case; many greedy algorithms produce suboptimal results. However, in this case the greedy procedure is indeed optimal: the subspace spanned by the first k principal components is the *best* subspace we can choose in terms of the ℓ_2 -norm of the projections. The theorem is proved in Section A.3 of the appendix.

Theorem 8.3. *For any matrix $X \in \mathbb{R}^{m \times n}$ with left singular vectors u_1, u_2, \dots, u_m corresponding to the nonzero singular values $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_m$,*

$$\sum_{i=1}^n \|\mathcal{P}_{\text{span}(u_1, u_2, \dots, u_k)} x_i\|_2^2 \geq \sum_{i=1}^n \|\mathcal{P}_{\mathcal{S}} x_i\|_2^2, \quad (113)$$

for any subspace \mathcal{S} of dimension $k \leq \min\{m, n\}$.

Figure 3 illustrates the importance of centering before applying PCA. Theorems 8.2 and 8.3 still hold if the data are not centered. However, the norm of the projection onto a certain direction no longer reflects the variation of the data. In fact, if the data are concentrated around a point that is far from the origin, the first principal component will tend to be aligned in that direction. This makes sense as projecting onto that direction captures more energy. As a result, the principal components do not capture the directions of maximum variation *within* the cloud of data.

8.2 PCA: Probabilistic interpretation

Let us interpret our data, x_1, x_2, \dots, x_n in \mathbb{R}^m , as samples of a random vector \mathbf{X} of dimension m . Recall that we are interested in determining the directions of maximum variation of the data in ambient space. In probabilistic terms, we want to find the directions in which the data have *higher variance*. The covariance matrix of the data provides this information. In fact, we can use it to determine the variance of the data in any direction.

Lemma 8.4. *Let u be a unit vector,*

$$\text{Var}(\mathbf{X}^T u) = u^T \Sigma_{\mathbf{X}} u. \quad (114)$$

Proof.

$$\text{Var}(\mathbf{X}^T u) = \mathbb{E}((\mathbf{X}^T u)^2) - \mathbb{E}^2(\mathbf{X}^T u) \quad (115)$$

$$= \mathbb{E}(u \mathbf{X} \mathbf{X}^T u) - \mathbb{E}(u^T \mathbf{X}) \mathbb{E}(\mathbf{X}^T u) \quad (116)$$

$$= u^T (\mathbb{E}(\mathbf{X} \mathbf{X}^T) - \mathbb{E}(\mathbf{X}) \mathbb{E}(\mathbf{X})^T) u \quad (117)$$

$$= u^T \Sigma_{\mathbf{X}} u. \quad (118)$$

□

Of course, if we only have access to samples of the random vector, we do not know the covariance matrix of the vector. However we can approximate it using the **empirical covariance matrix**.

Definition 8.5 (Empirical covariance matrix). *The empirical covariance of the vectors x_1, x_2, \dots, x_n in \mathbb{R}^m is equal to*

$$\bar{\Sigma}_n := \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)(x_i - \bar{x}_n)^T \quad (119)$$

$$= \frac{1}{n} X X^T, \quad (120)$$

where \bar{x}_n is the sample mean, as defined in Definition 1.3 of Lecture Notes 4, and X is the matrix containing the centered data as defined in (103).

If we assume that the mean of the data is zero (i.e. that the data have been centered using the true mean), then the empirical covariance is an unbiased estimator of the true covariance matrix:

$$\mathbb{E}\left(\frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^T\right) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(\mathbf{X}_i \mathbf{X}_i^T) \quad (121)$$

$$= \Sigma_{\mathbf{X}}. \quad (122)$$

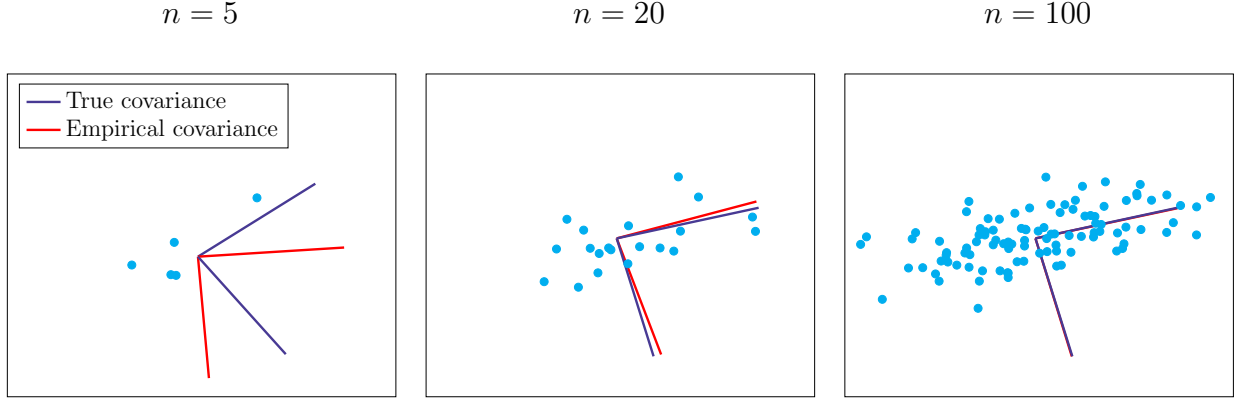


Figure 4: Principal components of n data vectors samples from a 2D Gaussian distribution. The eigenvectors of the covariance matrix of the distribution are also shown.

If the higher moments of the data $E(X_i^2 X_j^2)$ and $E(X_i^4)$ are finite, by Chebyshev's inequality the entries of the empirical covariance matrix converge to the entries of the true covariance matrix. This means that in the limit

$$\text{Var}(\mathbf{X}^T u) = u^T \Sigma_{\mathbf{X}} u \quad (123)$$

$$\approx \frac{1}{n} u^T X X^T u \quad (124)$$

$$= \frac{1}{n} \|X^T u\|_2^2 \quad (125)$$

for any unit-norm vector u . In the limit the principal components correspond to the directions of maximum variance of the underlying random vector. These directions also correspond to the eigenvectors of the true covariance matrix by Theorem 6.2. Figure 4 illustrates how the principal components converge to the eigenvectors of Σ .

A Proofs

A.1 Proof of Theorem 6.2

The eigenvectors are an orthonormal basis (they are mutually orthogonal and we assume that they have been normalized), so we can represent any unit-norm vector h_k that is orthogonal to u_1, \dots, u_{k-1} as

$$h_k = \sum_{i=k}^m \alpha_i u_i \quad (126)$$

where

$$\|h_k\|_2^2 = \sum_{i=k}^m \alpha_i^2 = 1, \quad (127)$$

by Lemma 1.5. Note that h_1 is just an arbitrary unit-norm vector.

Now we will show that the value of the quadratic form when the normalized input is restricted to be orthogonal to u_1, \dots, u_{k-1} cannot be larger than λ_k ,

$$h_k^T A h_k = \sum_{i=1}^n \lambda_i \left(\sum_{j=k}^m \alpha_j u_i^T u_j \right)^2 \quad \text{by (92) and (126)} \quad (128)$$

$$= \sum_{i=1}^n \lambda_i \alpha_i^2 \quad \text{because } u_1, \dots, u_m \text{ is an orthonormal basis} \quad (129)$$

$$\leq \lambda_k \sum_{i=k}^m \alpha_i^2 \quad \text{because } \lambda_k \geq \lambda_{k+1} \geq \dots \geq \lambda_m \quad (130)$$

$$= \lambda_k, \quad \text{by (127)}. \quad (131)$$

This establishes (93) and (95). To prove (108) and (110) we just need to show that u_k achieves the maximum

$$u_k^T A u_k = \sum_{i=1}^n \lambda_i (u_i^T u_k)^2 \quad (132)$$

$$= \lambda_k. \quad (133)$$

A.2 Proof of Theorem 7.2

It is sufficient to prove

$$\dim(\text{row}(A)) \leq \dim(\text{row}(A)) \quad (134)$$

for an arbitrary matrix A . We can apply the result to A^T to establish $\dim(\text{row}(A)) \geq \dim(\text{row}(A))$, since $\text{row}(A) = \text{row}(A)^T$ and $\text{row}(A) = \text{row}(A)^T$.

To prove (134) let $r := \dim(\text{row}(A))$ and let $x_1, \dots, x_r \in R^n$ be a basis for $\text{row}(A)$. Consider the vectors $Ax_1, \dots, Ax_r \in R^n$. They belong to $\text{row}(A)$ by (43), so if they are linearly independent then $\dim(\text{row}(A))$ must be at least r . We will prove that this is the case by contradiction.

Assume that Ax_1, \dots, Ax_r are linearly dependent. Then there exist coefficients $\alpha_1, \dots, \alpha_r \in \mathbb{R}$ such that

$$0 = \sum_{i=1}^r \alpha_i Ax_i = A \left(\sum_{i=1}^r \alpha_i x_i \right) \quad (\text{by linearity of the matrix product}), \quad (135)$$

This implies that $\sum_{i=1}^r \alpha_i x_i$ is orthogonal to every row of A and hence to every vector in $\text{row}(A)$. However it is in the span of a basis of $\text{row}(A)$ by construction! This is only possible if $\sum_{i=1}^r \alpha_i x_i = 0$, which is a contradiction because x_1, \dots, x_r are assumed to be linearly independent.

A.3 Proof of Theorem 8.3

We prove the result by induction. The *base case*, $k = 1$, follows immediately from (108). To complete the proof we need to show that if the result is true for $k - 1 \geq 1$ (this is the *induction hypothesis*) then it also holds for k .

Let \mathcal{S} be an arbitrary subspace of dimension k . We choose an orthonormal basis for the subspace b_1, b_2, \dots, b_k such that b_k is orthogonal to u_1, u_2, \dots, u_{k_1} . We can do this by using any vector that is linearly independent of u_1, u_2, \dots, u_{k_1} and subtracting its projection onto the span of u_1, u_2, \dots, u_{k_1} (if the result is always zero then \mathcal{S} is in the span and consequently cannot have dimension k).

By the induction hypothesis,

$$\sum_{i=1}^n \left\| \mathcal{P}_{\text{span}(u_1, u_2, \dots, u_{k_1})} x_i \right\|_2^2 = \sum_{i=1}^{k-1} \|X^T u_i\|_2^2 \quad \text{by (6)} \quad (136)$$

$$\leq \sum_{i=1}^n \left\| \mathcal{P}_{\text{span}(b_1, b_2, \dots, b_{k_1})} x_i \right\|_2^2 \quad (137)$$

$$= \sum_{i=1}^{k-1} \|X^T b_i\|_2^2 \quad \text{by (6)}. \quad (138)$$

By (110)

$$\sum_{i=1}^n \|\mathcal{P}_{\text{span}(u_k)} x_i\|_2^2 = \|X^T u_k\|_2^2 \quad (139)$$

$$\leq \sum_{i=1}^n \|\mathcal{P}_{\text{span}(b_k)} x_i\|_2^2 \quad (140)$$

$$= \|X^T b_k\|_2^2. \quad (141)$$

Combining (138) and (138) we conclude

$$\sum_{i=1}^n \|\mathcal{P}_{\text{span}(u_1, u_2, \dots, u_k)} x_i\|_2^2 = \sum_{i=1}^k \|X^T u_i\|_2^2 \quad (142)$$

$$\leq \sum_{i=1}^k \|X^T b_i\|_2^2 \quad (143)$$

$$\leq \sum_{i=1}^n \|\mathcal{P}_{\mathcal{S}} x_i\|_2^2. \quad (144)$$