

Lecture 2

Linear Algebra (II)

明玉瑞 Yurui Ming

yrming@gmail.com

声明

Disclaimer

- ▶ 本讲义在准备过程中由于时间所限，所用材料来源并未规范标示引用来源。所引材料仅用于教学所用，作者无意侵犯原著者之知识产权，所引材料之知识产权均归原著者所有；若原著者介意之，请联系作者更正及删除。

The time limit during the preparation of these slides incurs the situation that not all the sources of the used materials (texts or images) are properly referenced or clearly manifested. However, all materials in these slides are solely for teaching and the author is with no intention to infringe the copyright bestowed on the original authors or manufacturers. All credits go to corresponding IP holders. Please address the author for any concern for remedy including deletion.

范数

Norm

- 有时我们需要计算向量的大小，特别是在机器学习中。向量的大小称为范数，形式上，向量的 p 范数 L^p 由下式给出，其中 $p \in \mathbb{R}$ 且 $p \geq 1$ ：

Sometimes we need to measure the size or magnitude of a vector, especially in machine learning. The magnitude of vectors are called norms. Formally, for $p \in \mathbb{R}$, $p \geq 1$, the L^p norm is given by:

$$\|\mathbf{x}\|_p = \left(\sum_i |x_i|^p \right)^{\frac{1}{p}}$$

- 范数是将向量映射到非负值的函数，直观的说，向量 \mathbf{x} 的范数衡量的是从原点到点 \mathbf{x} 的距离。更严格地说，范数是满足以下属性的任何函数 f ：

Norms are functions mapping vectors to non-negative values. On an intuitive level, the norm of a vector \mathbf{x} measures the distance from the origin to the point \mathbf{x} . More rigorously, a norm is any function f that satisfies the following properties:

- $f(\mathbf{x}) = 0 \implies \mathbf{x} = 0$
- $f(\mathbf{x} + \mathbf{y}) \leq f(\mathbf{x}) + f(\mathbf{y})$
- $\forall \alpha \in \mathbb{R}, f(\alpha \mathbf{x}) = |\alpha| f(\mathbf{x})$

范数

Norm

- $p = 2$ 的 L^2 范数又被称为欧几里得范数，它是从原点到由 \mathbf{x} 标识的点的欧几里得距离。 L^2 范数在机器学习中的使用如此频繁，以至于下标 2 通常被省略，简单地表示为 $\|\mathbf{x}\|$ 。使用 L^2 范数的平方计算向量的大小也很常见，因为它可以简单地表示为 $\mathbf{x}^T \mathbf{x}$ 。

The L^2 norm, with $p = 2$, is known as the Euclidean norm, which is simply the Euclidean distance from the origin to the point identified by \mathbf{x} . The L^2 norm is used so frequently in machine learning that it is often denoted simply as $\|\mathbf{x}\|$, with the subscript 2 omitted. It is also common to measure the size of a vector using the squared L^2 norm, which can be calculated simply as $\mathbf{x}^T \mathbf{x}$.

- 与 L^2 范数本身相比，平方 L^2 范数在数学和计算上更加便于使用。例如，平方 L^2 范数对 \mathbf{x} 的每个元素的导数仅依赖于 \mathbf{x} 的对应元素，而 L^2 范数的所有导数都依赖于整个向量。

The squared L^2 norm is more convenient to work with mathematically and computationally than the L^2 norm itself. For example, each derivative of the squared L^2 norm with respect to each element of \mathbf{x} depends only on the corresponding element of \mathbf{x} , while all the derivatives of the L^2 norm depend on the entire vector..

范数

Norm

- 在许多情况下，平方 L^2 范数可能是不太好用，因为它在原点附近增加得非常缓慢。在一些机器学习应用程序中，区分恰好为零的元素和小但非零的元素很重要。在这种情况下，我们引入另一个函数，该函数在所有位置都以相同的速率增长，但保留了数学上的简单性： L^1 范数。 L^1 范数可以简化为： $\|\mathbf{x}\|_1 = \sum_i |x_i|$

In many contexts, the squared L^2 norm may be undesirable because it increases very slowly near the origin. In several machine learning applications, it is important to discriminate between elements that are exactly zero and elements that are small but nonzero. In these cases, we turn to a function that grows at the same rate in all locations, but that retains mathematical simplicity: the L^1 norm. The L^1 norm may be simplified to $\|\mathbf{x}\|_1 = \sum_i |x_i|$

- L^1 范数常用于机器学习，特别是在零元素和非零元素之间的差异非常重要的场景。每当 \mathbf{x} 的一个元素从0移开 ϵ 时， L^1 范数就会增加 ϵ 。

The L^1 norm is commonly used in machine learning when the difference between zero and nonzero elements is very important. Every time an element of \mathbf{x} moves away from 0 by ϵ , the L^1 norm increases by ϵ .

范数

Norm

- 我们有时会通过计算非零元素的个数来计算向量的大小。一些作者将此函数称为“ L^0 范数”，但这是不正确的，向量中非零元素的个数不是范数，因为按 α 缩放向量不会改变非零元素的个数。 L^1 范数通常用作非零元素个数的替代。

We sometimes measure the size of the vector by counting its number of nonzero elements. Some authors refer to this function as the “ L^0 norm”, but this is incorrect terminology. The number of nonzero entries in a vector is not a norm, because scaling the vector by α does not change the number of nonzero entries. The L^1 norm is often used as a substitute for the number of nonzero entries.

- 机器学习中另一种常见的范数是 L^∞ 范数，也称为最大范数。该范数简化为向量中具有最大绝对值的元素，即 $\|\mathbf{x}\|_\infty = \max_i |x_i|$ 。

One other norm that commonly arises in machine learning is the L^∞ norm, also known as the max norm. This norm simplifies to the absolute value of the element with the largest magnitude in the vector, that is $\|\mathbf{x}\|_\infty = \max_i |x_i|$.

范数

Norm

- 有时我们也可能希望计算矩阵的大小。在深度学习的背景下，最常见的方式是使用类似于向量的 L^2 范数的稍微晦涩一点儿的Frobenius范数： $\|\mathbf{A}\|_F = \sqrt{\sum_{i,j} A_{i,j}^2}$

Sometimes we may also wish to measure the size of a matrix. In the context of deep learning, the most common way to do this is with the otherwise obscure Frobenius norm

which is analogous to the L^2 norm of a vector: $\|\mathbf{A}\|_F = \sqrt{\sum_{i,j} A_{i,j}^2}$

- 两个向量的点积可以用范数重写。具体来说， $\mathbf{x}^T \mathbf{y} = \|\mathbf{x}\|_2 \|\mathbf{y}\|_2 \cos \theta$ ，其中 θ 是 \mathbf{x} 和 \mathbf{y} 之间的夹角。但由于实际很难知道在高维空间中向量的夹角，因此读者知道其几何意义即可。

The dot product of two vectors can be rewritten in terms of norms. Specifically, $\mathbf{x}^T \mathbf{y} = \|\mathbf{x}\|_2 \|\mathbf{y}\|_2 \cos \theta$, where θ is the angle between \mathbf{x} and \mathbf{y} . But in practice it is hard to know the angle between two vectors especially in high dimensional space, the reader only needs to know its geometric meaning here.

矩阵

Matrices

- ▶ 下面代码展示了TensorFlow对向量与矩阵范数的支持：

The following code demonstrates the support of TensorFlow for norms of vectors and matrices:

```
import numpy as np; import tensorflow as tf

a = tf.Variable([1, 2, 3], dtype=tf.float32); a_l1 = tf.linalg. normalize(a, 1)
a_l2 = tf.linalg.normalize(a, 2); a_linf = tf.linalg. normalize(a, np.inf)
print(f"a_l1 = {a_l1}, a_l2 = {a_l2}, a_linf = {a_linf}")

# a_l1 = 6.0, a_l2 = 3.7416574954986572, a_linf = 3.0

A = tf.Variable([[1, 2, 3], [4, 5, 6]], dtype=tf.float32)
A_l2 = tf.linalg. normalize(A, 2); A_lfro = tf.linalg.normalize(A, "fro", axis=(0,1))
A_l2_r = tf.linalg.normalize(A, 2, axis=0); A_l2_c = tf.linalg.norm(A, 2, axis=1)
print(f"A_l2 = {A_l2}, A_lfro = {A_lfro}, A_l2_r = {A_l2_r}, A_l2_c = {A_l2_c}")
```


特殊矩阵

Special Matrices

- ▶ 对角矩阵是沿主对角线元素非零、其它位置元素均为零的矩阵。形式上，一个矩阵 \mathbf{D} 是对角矩阵当且仅当 $D_{i,j} = 0$ 对所有 $i \neq j$ 成立。对角元素都是1的单位矩阵是对角矩阵。我们用 $\text{diag}(\mathbf{v})$ 来表示一个方对角矩阵，其对角元素由向量 \mathbf{v} 的元素给出。对角矩阵在计算上比较高效，要计算 $\text{diag}(\mathbf{v}) \mathbf{x}$ ，我们只需要将每个元素 x_i 缩放 v_i 倍，也就是说， $\text{diag}(\mathbf{v}) \mathbf{x} = \mathbf{v} \odot \mathbf{x}$ 。求方对角矩阵的逆也很简单，但只有当每个对角线元素都非零时，逆矩阵才存在，此时， $\text{diag}(\mathbf{v})^{-1} = \text{diag}([1/v_1, 1/v_2, \dots, 1/v_n]^T)$ 。一般地，我们可以根据任意矩阵推导出一些通用的机器学习算法，但通过将某些矩阵限制为对角矩阵时描述和推导都更为简单。

Diagonal matrices consist of nonzero entries only along the main diagonal and zeros elsewhere. Formally, a matrix \mathbf{D} is diagonal if and only if $D_{i,j} = 0$ for all $i \neq j$. The identity matrix with all the diagonal entries 1 is a diagonal matrix. We write $\text{diag}(\mathbf{v})$ to denote a square diagonal matrix whose diagonal entries are given by the entries of the vector \mathbf{v} . Multiplying by a diagonal matrix is computationally efficient. To compute $\text{diag}(\mathbf{v}) \mathbf{x}$, we only need to scale each element x_i by v_i . In other words, $\text{diag}(\mathbf{v}) \mathbf{x} = \mathbf{v} \odot \mathbf{x}$. Inverting a square diagonal matrix is also efficient. The inverse exists only if every diagonal entry is non-zero, and in that case, $\text{diag}(\mathbf{v})^{-1} = \text{diag}([1/v_1, 1/v_2, \dots, 1/v_n]^T)$. In general, we may derive some general machine learning algorithm in terms of arbitrary matrices but obtain a less expensive and less descriptive algorithm by restricting some matrices to be diagonal.

特殊矩阵

Special Matrices

- 并非所有对角矩阵都需要是方阵，但非方对角矩阵没有逆矩阵。对于非方对角矩阵 \mathbf{D} ，乘积 $\mathbf{D}\mathbf{x}$ 将涉及对 \mathbf{x} 的每个元素进行缩放，如果 \mathbf{D} 的高度大于它的宽度，则将一些零连接到结果，或者如果 \mathbf{D} 的宽度大于它的高度，则丢弃向量的最后一些元素。

Not all diagonal matrices need be square, however, non-square diagonal matrices do not have inverses. For a non-square diagonal matrix \mathbf{D} , the product $\mathbf{D}\mathbf{x}$ will involve scaling each element of \mathbf{x} and either concatenating some zeros to the result, if \mathbf{D} is taller than it is wide, or discarding some of the last elements of the vector, if \mathbf{D} is wider than it is tall.

- 对称矩阵是转置等于其自身的矩阵，即 $\mathbf{A} = \mathbf{A}^T$ 。当元素由不依赖于参数顺序的具有两个参数的某些函数生成时，通常会出现对称矩阵。例如，如果 \mathbf{A} 是距离测量矩阵，其中 $A_{i,j}$ 给出点 i 到点 j 的距离，则 $A_{i,j} = A_{j,i}$ ，因为距离函数是对称的。

A symmetric matrix is any matrix that is equal to its own transpose, aka., $\mathbf{A} = \mathbf{A}^T$. Symmetric matrices often arise when the entries are generated by some function of two arguments that does not depend on the order of the arguments. For example, if \mathbf{A} is a matrix of distance measurements, with $A_{i,j}$ giving the distance from point i to point j , then $A_{i,j} = A_{j,i}$, because distance functions are symmetric.

特殊矩阵

Special Matrices

- ▶ 单位向量是欧几里得范数为1或具有单位范数的向量，即 $\|\mathbf{x}\|^2 = 1$ 。如果 $\mathbf{x}^T \mathbf{y} = 0$ ，则称向量 \mathbf{x} 和向量 \mathbf{y} 正交。如果两个向量都具有非零范数，这意味着它们彼此成90度角。在 \mathbb{R}^n 中，至多 n 个向量可能是非零范数相互正交的。如果多个向量不仅正交而且欧几里得范数均为1，我们称它们是规范正交的或标准正交的。

A unit vector is a vector with Euclidean norm equal to 1, or with unit norm, namely, $\|\mathbf{x}\|^2 = 1$. A vector \mathbf{x} and a vector \mathbf{y} are orthogonal to each other if $\mathbf{x}^T \mathbf{y} = 0$. If both vectors have nonzero norm, this means that they are at a 90 degree angle to each other. In \mathbb{R}^n , at most n vectors may be mutually orthogonal with nonzero norm. If the vectors not only are orthogonal but also have unit norm, we call them orthonormal.

- ▶ 正交矩阵是行相互正交且列相互正交的方阵，即 $\mathbf{A}^T \mathbf{A} = \mathbf{A} \mathbf{A}^T = \mathbf{I}$ ，这意味着 $\mathbf{A}^{-1} = \mathbf{A}^T$ ，因此计算正交矩阵的逆矩阵计算成本非常低。请注意正交矩阵的定义，它们的行或列之间不仅是正交的，而且是标准正交的。行或列正交但非标准正交的矩阵没有特殊名称定义。

An orthogonal matrix is a square matrix whose rows are mutually orthonormal and whose columns are mutually orthonormal as well, that is $\mathbf{A}^T \mathbf{A} = \mathbf{A} \mathbf{A}^T = \mathbf{I}$. This implies that $\mathbf{A}^{-1} = \mathbf{A}^T$, means their inverse are very cheap to compute. Pay careful attention to the definition of orthogonal matrices, their rows or columns are not merely orthogonal but fully orthonormal. There is no special term for a matrix whose rows or columns are orthogonal but not orthonormal.

本征分解

Eigendecomposition

- ▶ 许多数学对象可以通过将它们分解成基本的组成部分，或者找到它们的一些不依赖所选择的表示方式的一些普适性质，从而达到更好地理解。例如，整数可以分解为质因数。数字12的呈现方式会当我们选择基数是2以二进制形式表示时而改变，但 $12=2 \times 2 \times 3$ 始终为真。从这个表示中我们可以得出有用的属性，例如，12不能被5整除，并且12的任何整数倍都可以被3整除。就像我们可以通过将整数分解为质因数来发现一些关于整数的本质一样，我们还可以分解矩阵，以呈现有关其性质的信息，这些信息在矩阵表示为元素的数组时并不明显。

Many mathematical objects can be understood better by breaking them into constituent parts, or finding some properties of them that are universal, not caused by the way we choose to represent them. For example, integers can be decomposed into prime factors. The way we represent the number 12 will change depending on whether we write it in base ten or in binary, but it will always be true that $12 = 2 \times 2 \times 3$. From this representation we can conclude useful properties, for example, that 12 is not divisible by 5, and that any integer multiple of 12 will be divisible by 3. Much as we can discover something about the true nature of an integer by decomposing it into prime factors, we can also decompose matrices in ways that show us information about their functional properties that is not obvious from the representation of the matrix as an array of elements.

本征分解

Eigendecomposition

- 一种广为使用的矩阵分解称为特征分解，即将矩阵分解为一组特征向量和特征值。

One of the most widely used kinds of matrix decomposition is called eigen-decomposition, in which we decompose a matrix into a set of eigenvectors and eigenvalues.

- 方阵 \mathbf{A} 的特征向量是一个非零向量 \mathbf{v} ，乘以 \mathbf{A} 只会改变 \mathbf{v} 的尺度： $\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$ 。标量 λ 被称为对应于该特征向量的特征值。顺便说一下，也可以找到一个左特征向量使得 $\mathbf{v}^T\mathbf{A} = \lambda\mathbf{v}^T$ ，但我们通常关心的是右特征向量。如果 \mathbf{v} 是 \mathbf{A} 的特征向量，那么对于 $s \in \mathbb{R}$ 且 $s \neq 0$ ，任何重新缩放的向量 $s\mathbf{v}$ 也是特征向量。此外， $s\mathbf{v}$ 的特征值只需进行相反的缩放： λ/s 。出于这个原因，我们通常只求解单位特征向量。

An eigenvector of a square matrix \mathbf{A} is a nonzero vector \mathbf{v} such that multiplication by \mathbf{A} alters only the scale of \mathbf{v} : $\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$. The scalar λ is known as the eigenvalue corresponding to this eigenvector. BTW, One can also find a left eigenvector such that $\mathbf{v}^T\mathbf{A} = \lambda\mathbf{v}^T$, but we are usually concerned with right eigenvectors. If \mathbf{v} is an eigenvector of \mathbf{A} , then so is any rescaled vectors \mathbf{v} for $s \in \mathbb{R}$, $s \neq 0$. Moreover, $s\mathbf{v}$ has the reciprocal rescaled eigenvalue λ/s . For this reason, we usually look only for unit eigenvectors.

本征分解

Eigendecomposition

- 假设一个矩阵 A 有 n 个线性独立的特征向量 $\{v^{(1)}, v^{(2)}, \dots, v^{(n)}\}$, 具有相应的特征值 $\{\lambda_1, \lambda_2, \dots, \lambda_n\}$ 。我们可以将所有特征向量连接起来形成一个矩阵 V , $V = [v^{(1)}, v^{(2)}, \dots, v^{(n)}]$ 。同样, 我们可以连接特征值以形成向量 $\lambda = [\lambda_1, \lambda_2, \dots, \lambda_n]^T$ 。 A 的特征分解为 $A = V \text{diag}(\lambda) V^{-1}$ 。

Suppose that a matrix A has n linearly independent eigenvectors $\{v^{(1)}, v^{(2)}, \dots, v^{(n)}\}$ with corresponding eigenvalues $\{\lambda_1, \lambda_2, \dots, \lambda_n\}$. We may concatenate all the eigenvectors to form a matrix V with one eigenvector per column: $V = [v^{(1)}, v^{(2)}, \dots, v^{(n)}]$. Likewise, we can concatenate the eigenvalues to form a vector $\lambda = [\lambda_1, \lambda_2, \dots, \lambda_n]^T$. The eigendecomposition of A is then given by $A = V \text{diag}(\lambda) V^{-1}$.

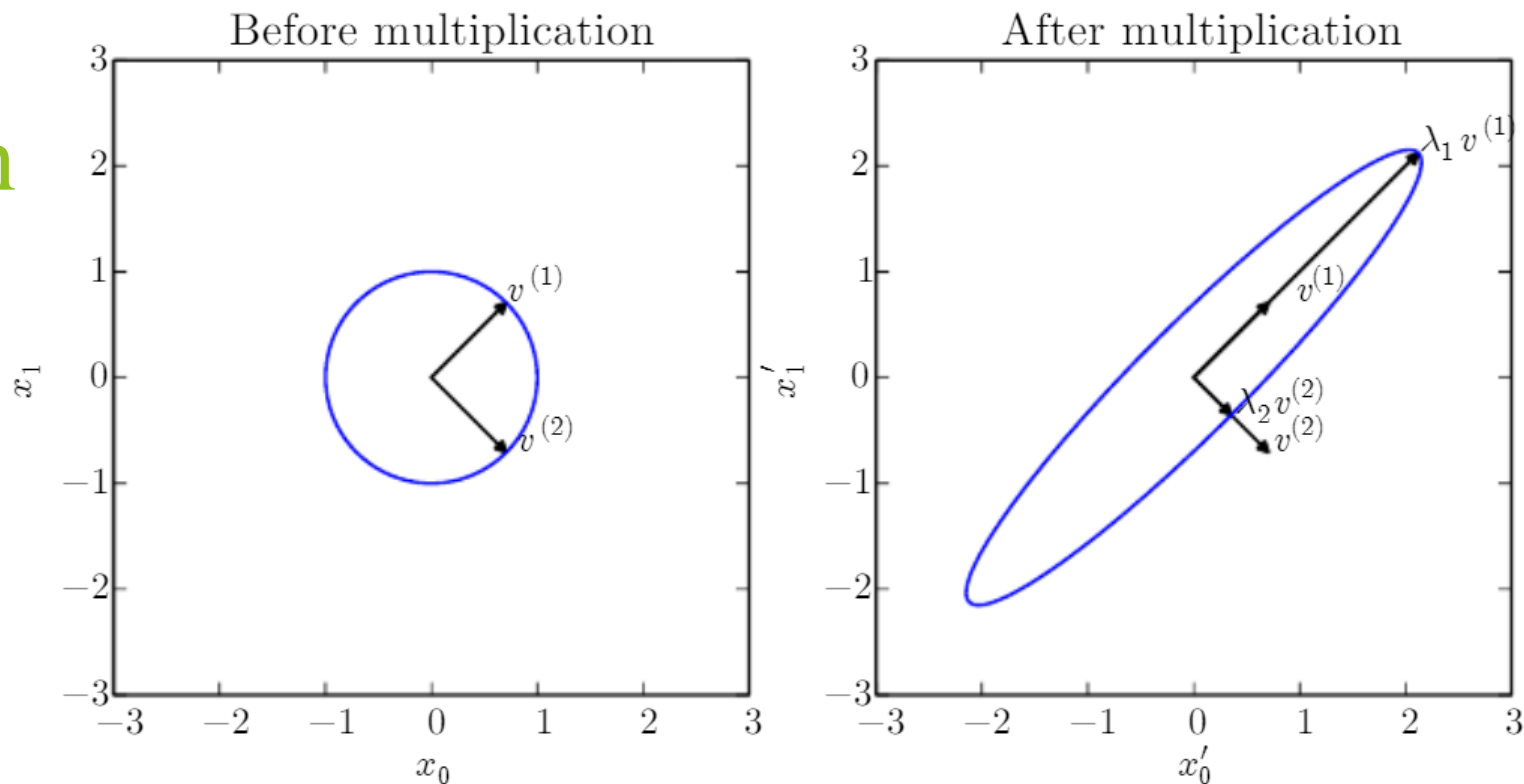
- 我们已经看到构造具有特定特征值和特征向量的矩阵使我们能够在期望的方向上拉伸空间。然而, 我们经常想将矩阵分解为它们的特征值和特征向量。这样做可以帮助我们分析矩阵的某些属性, 就像将一个整数分解成它的质因数可以帮助我们理解该整数的行为一样。

We have seen that constructing matrices with specific eigenvalues and eigenvectors enables us to stretch space in desired directions. Yet we often want to decompose matrices into their eigenvalues and eigenvectors. Doing so can help us analyze certain properties of the matrix, much as decomposing an integer into its prime factors can help us understand the behavior of that integer.

本征分解

Eigendecomposition

- 一个关于特征向量和特征值的影响示例如右上。这里，我们有一个具有两个正交特征向量 $\mathbf{v}^{(1)}$ 和 $\mathbf{v}^{(2)}$ 的矩阵 \mathbf{A} 。 $\mathbf{v}^{(1)}$ 对应于特征值 λ_1 ， $\mathbf{v}^{(2)}$ 对应于特征值 λ_2 。（左）我们将所有单位向量 $\mathbf{u} \in \mathbb{R}^2$ 的集合绘制为一个单位圆。（右）我们绘制了所有点 $\mathbf{A}\mathbf{u}$ 的集合。通过观察 \mathbf{A} 扭曲单位圆的方式，我们可以看到它在方向 $\mathbf{v}^{(i)}$ 上按 λ_i 缩放空间。



- An example of the effect of eigenvectors and eigenvalues is as top-right. Here, we have a matrix \mathbf{A} with two orthonormal eigenvectors, $\mathbf{v}^{(1)}$ with eigenvalue λ_1 and $\mathbf{v}^{(2)}$ with eigenvalue λ_2 . (Left) We plot the set of all unit vectors $\mathbf{u} \in \mathbb{R}^2$ as a unit circle. (Right) We plot the set of all points $\mathbf{A}\mathbf{u}$. By observing the way that \mathbf{A} distorts the unit circle, we can see that it scales space in direction $\mathbf{v}^{(i)}$ by λ_i .

本征分解

Eigendecomposition

- 并非每个矩阵都可以分解为特征值和特征向量。在某些情况下，分解存在但涉及复数而不是实数。在本书中，我们只分解具有简单分解形式的实对称矩阵，它们都可以分解为仅使用实值特征向量和特征值的表达式： $\mathbf{A} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T$ ，其中 \mathbf{Q} 是由 \mathbf{A} 的特征向量组成的正交矩阵， $\mathbf{\Lambda}$ 是对角矩阵。特征值 $\Lambda_{i,i}$ 与 \mathbf{Q} 的第 i 列的特征向量相关联，记为 $\mathbf{Q}_{:,i}$ 。因为 \mathbf{Q} 是一个正交矩阵，我们可以认为 \mathbf{A} 在方向 \mathbf{v}_i 上以 λ_i 倍缩放空间。注意，虽然任何实对称矩阵都具有特征分解，但特征分解可能不是唯一的。如果任何两个或多个特征向量具有相同的特征值，则位于其张成的空间的任何一组正交向量，也是具有该特征值的特征向量。

Not every matrix can be decomposed into eigenvalues and eigenvectors. In some cases, the decomposition exists but involves complex rather than real numbers. In this book, we usually just decompose real symmetric matrix which has a simple decomposition, aka., they can be decomposed into an expression using only real-valued eigenvectors and eigenvalues: $\mathbf{A} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T$, where \mathbf{Q} is an orthogonal matrix composed of eigenvectors of \mathbf{A} , and $\mathbf{\Lambda}$ is a diagonal matrix. The eigenvalue $\Lambda_{i,i}$ is associated with the eigenvector in column i of \mathbf{Q} , denoted as $\mathbf{Q}_{:,i}$. Because \mathbf{Q} is an orthogonal matrix, we can think of \mathbf{A} as scaling space by λ_i in direction \mathbf{v}_i . Notably, While any real symmetric matrix is guaranteed to have an eigendecomposition, the eigendecomposition may not be unique. If any two or more eigenvectors share the same eigenvalue, then any set of orthogonal vectors lying in their span are also eigenvectors with that eigenvalue.

本征分解

Eigendecomposition

- 按照惯例，我们通常按降序对 Λ 的条目进行排序。此时，若所有的特征值都是唯一的，那么特征分解唯一。

By convention, we usually sort the entries of Λ in descending order. Under this convention, the eigendecomposition is unique only if all the eigenvalues are unique.

- 矩阵的特征分解告诉我们关于矩阵的许多有用的事实。例如，当且仅当有特征值为零时，矩阵是奇异的。特征值均为正的矩阵称为正定矩阵，特征值全部为正值或零值的矩阵称为半正定矩阵。同样，如果所有特征值都是负数，则矩阵是负定的，如果所有特征值都是负数或零值，则矩阵是半负定的。半正定矩阵很有趣，因为它们保证 $\forall \mathbf{x}$, $\mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0$ 。若此时 \mathbf{x} 为单位向量，令 $f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x}$ ，则 $f_{max} = \max_i \{\lambda_i\}$ 。

The eigendecomposition of a matrix tells us many useful facts about the matrix. The matrix is singular if and only if any of the eigenvalues are zero. A matrix whose eigenvalues are all positive is called positive definite. A matrix whose eigenvalues are all positive or zero valued is called positive semi-definite. Likewise, if all eigenvalues are negative, the matrix is negative definite, and if all eigenvalues are negative or zero valued, it is negative semidefinite. Positive semidefinite matrices are interesting because they guarantee that $\forall \mathbf{x}$, $\mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0$. If \mathbf{x} are unit norm vectors, let $f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x}$, then $f_{max} = \max_i \{\lambda_i\}$.

奇异值分解

Singular Value Decomposition

- 奇异值分解(SVD)提供了另一种将矩阵因式分解为奇异向量和奇异值的方法。SVD使我们能够发现一些与特征分解所揭示的相同类型的信息。每个实矩阵未必有特征值分解，但都有奇异值分解，因此SVD更具有普遍的适用性。例如，如果矩阵不是方阵，则特征分解没有定义，我们必须使用奇异值分解来代替。

The singular value decomposition (SVD) provides another way to factorize a matrix, into singular vectors and singular values. The SVD enables us to discover some of the same kind of information as the eigendecomposition reveals. Every real matrix might not have has an eigenvalue decomposition, but it always has a singular value decomposition, which means the SVD is more generally applicable. For example, if a matrix is not square, the eigendecomposition is not defined, and we must use a singular value decomposition instead.

- 回想一下，特征分解涉及分析矩阵 \mathbf{A} 以求出特征向量的矩阵 \mathbf{V} 和特征值向量 $\boldsymbol{\lambda}$ ，这样我们就可以将 \mathbf{A} 重写为 $\mathbf{A} = \mathbf{V} \text{diag}(\boldsymbol{\lambda}) \mathbf{V}^{-1}$ 。奇异值分解类似，只是这次我们将 \mathbf{A} 写成三个矩阵的乘积： $\mathbf{A} = \mathbf{U} \mathbf{D} \mathbf{V}^T$ 。

Recall that the eigendecomposition involves analyzing a matrix \mathbf{A} to discover a matrix \mathbf{V} of eigenvectors and a vector $\boldsymbol{\lambda}$ of eigenvalues such that we can rewrite \mathbf{A} as $\mathbf{A} = \mathbf{V} \text{diag}(\boldsymbol{\lambda}) \mathbf{V}^{-1}$. The singular value decomposition is similar, except this time we will write \mathbf{A} as a product of three matrices: $\mathbf{A} = \mathbf{U} \mathbf{D} \mathbf{V}^T$.

奇异值分解

Singular Value Decomposition

- 在矩阵 \mathbf{A} 的奇异值分解 $\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{V}^T$ 的表示中, 假设 \mathbf{A} 是一个 $m \times n$ 矩阵, 则 \mathbf{U} 被定义为一个 $m \times m$ 矩阵, \mathbf{D} 为一个 $m \times n$ 矩阵, \mathbf{V} 为一个 $n \times n$ 矩阵。这些矩阵中的每一个都被定义为具有特殊结构的矩阵。矩阵 \mathbf{U} 和 \mathbf{V} 都被定义为正交矩阵, 矩阵 \mathbf{D} 定义为对角矩阵, 注意 \mathbf{D} 不一定是方阵。沿 \mathbf{D} 对角线的元素称为矩阵 \mathbf{A} 的奇异值。 \mathbf{U} 的列称为左奇异向量, \mathbf{V} 的列被称为右奇异向量。我们实际上可以用关于 \mathbf{A} 的函数的特征分解来解释 \mathbf{A} 的奇异值分解。 \mathbf{A} 的左奇异向量是 $\mathbf{A}\mathbf{A}^T$ 的特征向量。 \mathbf{A} 的右奇异向量是 $\mathbf{A}^T\mathbf{A}$ 的特征向量。 \mathbf{A} 的非零奇异值是 $\mathbf{A}^T\mathbf{A}$ 的特征值的平方根。SVD 最有用的特征是我们可以用它来将矩阵求逆部分推广到非方矩阵。

In the expression of The singular value decomposition of matrix \mathbf{A} , suppose that \mathbf{A} is an $m \times n$ matrix, then \mathbf{U} is defined to be an $m \times m$ matrix, \mathbf{D} to be an $m \times n$ matrix, and \mathbf{V} to be an $n \times n$ matrix. Each of these matrices is defined to have a special structure. The matrices \mathbf{U} and \mathbf{V} are both defined to be orthogonal matrices. The matrix \mathbf{D} is defined to be a diagonal matrix. Note that \mathbf{D} is not necessarily square. The elements along the diagonal of \mathbf{D} are known as the singular values of the matrix \mathbf{A} . The columns of \mathbf{U} are known as the left-singular vectors. The columns of \mathbf{V} are known as the right-singular vectors. We can actually interpret the singular value decomposition of \mathbf{A} in terms of the eigendecomposition of functions of \mathbf{A} . The left-singular vectors of \mathbf{A} are the eigenvectors of $\mathbf{A}\mathbf{A}^T$. The right-singular vectors of \mathbf{A} are the eigenvectors of $\mathbf{A}^T\mathbf{A}$. The nonzero singular values of \mathbf{A} are the square roots of the eigenvalues of $\mathbf{A}^T\mathbf{A}$. Perhaps the most useful feature of the SVD is that we can use it to partially generalize matrix inversion to non square matrices.

奇异值分解

Singular Value Decomposition

- ▶ 矩阵不是方阵时，逆矩阵没定义。但实际上，假设矩阵 A 存在一个左逆 A_l^{-1} ，我们就可以通过每边左乘 A_l^{-1} ，求解线性方程 $Ax = b$ ，得到 $x = A_l^{-1}b$ 。当然，一般情况下，左逆并不唯一，并且，方程组可能无解，也可能有多组解。

Matrix inversion is not defined for matrices that are not square. Suppose we want to make a left-inverse of a matrix A , like A_l^{-1} , then we can solve a linear equation $Ax = b$ by left-multiplying each side to obtain $x = A_l^{-1}b$. In general, the left inverse of an matrix is not necessary unique, meantime, it is possible for this equation to have no solution, or multiple possible solutions.

- ▶ 实际上，计算这个左逆 A_l^{-1} 基于的正是奇异值分解： $A_l^{-1} = VD^{-1}U^T$ ，其中 U 、 D 和 V 是对 A 进行奇异值分解时的相应的矩阵，对角矩阵 D^{-1} 称为 D 的伪逆，通过取 D 非零元素的倒数，然后进行转置得到。验证如右：

In fact, the computation of left inverse is based on SVD: $A_l^{-1} = VD^{-1}U^T$, where U , D and V are the singular value decomposition of A , and D^{-1} is a diagonal matrix, called the pseudoinverse of D , obtained by taking the reciprocal of nonzero elements of D then taking the transpose of the resulting matrix. To prove it, we have deduction on the right:

$$\begin{aligned} A_l^{-1}A &= VD^{-1}U^T A \\ \Rightarrow A_l^{-1}A &= VD^{-1}U^T UDV^T \\ \Rightarrow A_l^{-1}A &= VD^{-1}I_m DV^T \\ \Rightarrow A_l^{-1}A &= VD^{-1}DV^T \\ \Rightarrow A_l^{-1}A &= VE_n V^T \\ \Rightarrow A_l^{-1}A &= E_n VV^T \\ \Rightarrow A_l^{-1}A &= E_n I_n \\ \Rightarrow A_l^{-1}A &= E_n \end{aligned}$$

注意， E_n 是对角矩阵，但对角线上元素可能为零。

迹、行列式

Trace, Determinant

- 矩阵的迹定义为矩阵所有对角线元素的总和： $\text{Tr}(\mathbf{A}) = \sum_i A_{i,i}$ ，标量的迹是它本身： $a = \text{Tr}(a)$ 。迹有许多有用的性质，如矩阵与其转置矩阵具有相同的迹： $\text{Tr}(\mathbf{A}) = \text{Tr}(\mathbf{A}^T)$ 。其满足交换律，即对 $\mathbf{A} \in \mathbb{R}^{m \times n}$ ， $\mathbf{B} \in \mathbb{R}^{n \times m}$ ， $\text{Tr}(\mathbf{AB}) = \text{Tr}(\mathbf{BA})$ 。矩阵 \mathbf{A} 的Frobenius范数可写作 $\|\mathbf{A}\|_F = \sqrt{\text{Tr}(\mathbf{AA}^T)}$ 。

The trace of a matrix is the sum of all the diagonal entries: $\text{Tr}(\mathbf{A}) = \sum_i A_{i,i}$. A scalar is its own trace: $a = \text{Tr}(a)$. The trace operator is endowed with many useful properties, for example, it is invariant to the transpose: $\text{Tr}(\mathbf{A}) = \text{Tr}(\mathbf{A}^T)$. It is invariant to permutation, that is for $\mathbf{A} \in \mathbb{R}^{m \times n}$ ， $\mathbf{B} \in \mathbb{R}^{n \times m}$ ， $\text{Tr}(\mathbf{AB}) = \text{Tr}(\mathbf{BA})$. An alternative way of writing the Frobenius norm of a matrix \mathbf{A} is $\|\mathbf{A}\|_F = \sqrt{\text{Tr}(\mathbf{AA}^T)}$.

- 方阵的行列式是将矩阵映射到实数的函数，表示为 $\det(\mathbf{A})$ 。行列式等于矩阵所有特征值的乘积，其绝对值可以被认为是矩阵扩展或收缩空间的量度。如果行列式为0，则空间至少沿一个维度完全塌缩，导致它失去高维体积；如果行列式为1，则转换保留体积。

The determinant of a square matrix, denoted $\det(\mathbf{A})$, is a function that maps matrices to real scalars. The determinant is equal to the product of all the eigenvalues of the matrix. The absolute value of the determinant can be thought of as a measure of how much multiplication by the matrix expands or contracts space. If the determinant is 0, then space is contracted completely along at least one dimension, causing it to lose all its volume. If the determinant is 1, then the transformation preserves volume.

例子：主成分分析

Example: Principal Components Analysis

- ▶ 一个可以仅使用基本线性代数知识推导出来的简单的机器学习算法是主成分分析 (PCA)。假设有 \mathbb{R}^n 中的 m 个点的集合 $\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)}\}$ ，我们想对这些点应用有损压缩，这意味着用较少内存存储这些点但以损失一些精度为代价。当然，我们希望损失的精度尽可能小。

One simple machine learning algorithm, principal components analysis (PCA), can be derived using only knowledge of basic linear algebra. Suppose we have a collection of m points $\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)}\}$ in \mathbb{R}^n and we want to apply lossy compression to these points. Lossy compression means storing the points in a way that requires less memory but may lose some precision. We want to lose as little precision as possible.

- ▶ 我们可以选择的一种方法是在低维空间对这些点进行编码。对于任一 $\mathbf{x}^{(i)} \in \mathbb{R}^n$ ，我们将找到一个对应的编码向量 $\mathbf{c}^{(i)} \in \mathbb{R}^l$ 。如果 l 比 n 小，则存储新编码值将比存储原始数据占用更少的内存。我们希望要找到一个编码函数来为输入生成编码： $\mathbf{c} = f(\mathbf{x})$ ，并找到一个解码函数来生成给定其代码的重构输入： $\mathbf{x} \approx g(\mathbf{c})$ 。PCA即由选择的解码函数定义。

One way we can encode these points is to represent a lower-dimensional version of them. For each point $\mathbf{x}^{(i)} \in \mathbb{R}^n$ we will find a corresponding code vector $\mathbf{c}^{(i)} \in \mathbb{R}^l$. If l is smaller than n , storing the code points will take less memory than storing the original data. We will want to find some encoding function that produces the code for an input: $\mathbf{c} = f(\mathbf{x})$, and a decoding function that produces the reconstructed input given its code: $\mathbf{x} \approx g(\mathbf{c})$. PCA is thus defined by our choice of the decoding function.

例子：主成分分析

Example: Principal Components Analysis

- 为了使解码器简单，我们选择使用矩阵乘法将编码映射回 \mathbb{R}^n 。令 $g(\mathbf{c}) = \mathbf{D}\mathbf{c}$ ，其中 $\mathbf{D} \in \mathbb{R}^{n \times l}$ 是定义解码的矩阵。为了简化编码问题，PCA将 \mathbf{D} 的列约束为彼此正交。为了给问题一个唯一解，我们进一步将 \mathbf{D} 的所有列都限制为具有单位范数。

To make the decoder very simple, we choose to use matrix multiplication to map the code back into \mathbb{R}^n . Let $g(\mathbf{c}) = \mathbf{D}\mathbf{c}$, where $\mathbf{D} \in \mathbb{R}^{n \times l}$ is the matrix defining the decoding. To keep the encoding problem easy, PCA constrains the columns of \mathbf{D} to be orthogonal to each other. To give the problem a unique solution, we constrain all the columns of \mathbf{D} to have unit norm..

- 我们需要做的第一件事就是确定每个输入点 \mathbf{x} 的最佳编码值 \mathbf{c}^* 。一种方法是最小化输入 \mathbf{x} 与其重构 $g(\mathbf{c}^*)$ 之间的距离。在主成分算法中，我们使用 L^2 范数来衡量这个距离。由于平方 L^2 范数与 L^2 范数在最小化 \mathbf{c} 值上等价，我们实际上使用平方 L^2 范数：

The first thing we need to do is figure out how to generate the optimal code point \mathbf{c}^* for each input point \mathbf{x} . One way to do this is to minimize the distance between the input point \mathbf{x} and its reconstruction, $g(\mathbf{c}^*)$. In the principal components algorithm, we use the L^2 norm. However, since the squared L^2 and the L^2 norm are equivalent in minimizing the same value of \mathbf{c} , actually we adopt the squared L^2 norm:

$$\mathbf{c}^* = \arg \min_{\mathbf{c}} \|\mathbf{x} - g(\mathbf{c})\|_2^2$$

例子：主成分分析

Example: Principal Components Analysis

- ▶ 根据平方 L^2 范数的定义，上述最小化的等式等价于：

By the definition of the L^2 norm, The function being minimized simplifies to:

$$\begin{aligned} \mathbf{c}^* &= \arg \min_{\mathbf{c}} (\mathbf{x} - g(\mathbf{c}))^T (\mathbf{x} - g(\mathbf{c})) \\ &= \arg \min_{\mathbf{c}} (\mathbf{x}^T \mathbf{x} - 2\mathbf{x}^T g(\mathbf{c}) + g(\mathbf{c})^T g(\mathbf{c})) \end{aligned}$$

- ▶ 由于第一项为定值，我们可以去掉；我们将 \mathbf{D} 代入继续化简，由于 \mathbf{D} 的列向量正交且范数为1，我们最终有：

We can omit the first term since it is a fixed term. We can substitute $g(\mathbf{c})$ with \mathbf{D} for further simplification. By the orthogonality and unit norm constraints on \mathbf{D} , we further have:

$$\begin{aligned} \mathbf{c}^* &= \arg \min_{\mathbf{c}} (-2\mathbf{x}^T g(\mathbf{c}) + g(\mathbf{c})^T g(\mathbf{c})) \\ &= \arg \min_{\mathbf{c}} (-2\mathbf{x}^T \mathbf{D}\mathbf{c} + \mathbf{c}^T \mathbf{D}^T \mathbf{D}\mathbf{c}) \\ &= \arg \min_{\mathbf{c}} (-2\mathbf{x}^T \mathbf{D}\mathbf{c} + \mathbf{c}^T \mathbf{I}_l \mathbf{c}) = \arg \min_{\mathbf{c}} (-2\mathbf{x}^T \mathbf{D}\mathbf{c} + \mathbf{c}^T \mathbf{c}) \end{aligned}$$

- ▶ 我们不再证明， \mathbf{D} 的构造如下：

- 构造矩阵 \mathbf{X} ，其中 $\mathbf{X}_{i,:} = \mathbf{x}^{(1)}$
- 求解 $\mathbf{x}^T \mathbf{x}$ 的特征值 $\{\lambda_i\}$ 与对应的特征向量 $\{\mathbf{v}_i\}$ ，满足 $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$
- 构造矩阵 \mathbf{D} ，其中 $\mathbf{D}_{:,j} = \mathbf{v}_j, \quad j \leq l$ 。
- 最后 $\mathbf{c} = \mathbf{D}^T \mathbf{x}$, $\mathbf{x} = \mathbf{D}\mathbf{c}$

TensorFlow例子

TensorFlow Example

- ▶ 下面代码展示了TensorFlow对上面所讲内容的支持:

The following code demonstrates the support of TensorFlow for content lectured above:

```
A = tf.Variable([[1, 2, 3], [4, 5, 6], [7, 8, 9]], dtype=tf.float32); A_diag = tf.linalg.tensor_diag_part(A)
```

```
a = tf.Variable([1, 2, 3], dtype=tf.float32); diag_from_a = tf.linalg.diag(a)
```

```
S = A + tf.transpose(A)
```

```
w, v = tf.linalg.eig(S)
```

```
# w: <tf.Tensor: shape=(3,), dtype=complex64, numpy=
```

```
# array([ 3.2916481e+01+0.j, -2.9164727e+00+0.j,  1.0197117e-06+0.j],      dtype=complex64)>
```

```
# v: <tf.Tensor: shape=(3, 3), dtype=complex64, numpy=
```

```
# array([[ -0.35162535+0.j, -0.84243304+0.j,  0.4082484 +0.j],
```

```
        [ -0.5533561 +0.j, -0.16471203+0.j, -0.8164966 +0.j],
```

```
        [ -0.7550872 +0.j,  0.5130073 +0.j,  0.40824825+0.j]],      dtype=complex64)>
```

TensorFlow例子

TensorFlow Example

- ▶ 下面代码展示了TensorFlow对上面所讲内容的支持:

The following code demonstrates the support of TensorFlow for content lectured above:

```
A = tf.random.normal((3, 4)); S, U, V = tf.linalg.svd(A)
```

```
# U: [[ 0.09378654, 0.89235  , -0.44149217],  
      [-0.99545395, 0.09144013, -0.02664502],  
      [-0.01659341, -0.4419841 , -0.89686936]]
```

```
# S: [1.7239952 , 1.5212345 , 0.48055685]
```

```
# V: [[-0.8886927 , -0.21227351, -0.34561253],  
      [-0.17547224, -0.19447708, -0.0261912 ],  
      [ 0.3904554 , -0.79509723, -0.45282608],  
      [-0.16425498, -0.5337969 , 0.8214711 ]]
```

```
t = tf.linalg.trace(A); t1 = tf.linalg.trace(tf.linalg.matmul(A, tf.transpose(A)))
```

```
t2 = tf.linalg.trace(tf.linalg.matmul(tf.transpose(A), A)); assert tf.math.abs(t1 - t2) < 10e-6
```

TensorFlow例子

TensorFlow Example

- 下面代码展示了TensorFlow对上面所讲内容的支持:

The following code demonstrates the support of TensorFlow for content lectured above:

```
import numpy as np
import tensorflow as tf

def pca(X, l):
    n, m = X.get_shape().as_list() # Data matrix X, assumes 0-centered
    C = tf.linalg.matmul(tf.transpose(X), X) / (n-1) # Compute covariance matrix
    eigen_vals, eigen_vecs = tf.linalg.eig(C) # Eigen decomposition
    eigen_vecs = tf.math.real(tf.gather(eigen_vecs, [0, l], axis=1))
    X_pca = tf.linalg.matmul(X, eigen_vecs) # Project X onto PC space
    return X_pca

A = tf.random.normal((3, 4))
A = A - tf.math.reduce_mean(A, axis=0) # Normalize the data
A_pca = pca(A, 2); print(A_pca)
```