

Linear Algebra Review

Jing Xiang

March 18, 2014

1 PROPERTIES OF MATRICES

Below are a few basic properties of matrices:

- Matrix Multiplication is associative: $(\mathbf{AB})\mathbf{C} = \mathbf{A}(\mathbf{BC})$
- Matrix Multiplication is distributive: $\mathbf{A}(\mathbf{B} + \mathbf{C}) = \mathbf{AB} + \mathbf{AC}$
- Matrix Multiplication is NOT commutative in general, that is $\mathbf{AB} \neq \mathbf{BA}$. For example, if $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{B} \in \mathbb{R}^{n \times q}$, the matrix product \mathbf{BA} does not exist.

2 TRANSPOSE

The transpose of a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, is written as $\mathbf{A}^\top \in \mathbb{R}^{n \times m}$ where the entries of the matrix are given by:

$$(\mathbf{A}^\top)_{ij} = \mathbf{A}_{ji} \tag{2.1}$$

Properties:

- Transpose of a scalar is a scalar $a^\top = a$
- $(\mathbf{A}^\top)^\top = \mathbf{A}$
- $(\mathbf{AB})^\top = \mathbf{B}^\top \mathbf{A}^\top$
- $(\mathbf{A} + \mathbf{B})^\top = \mathbf{A}^\top + \mathbf{B}^\top$

3 TRACE

The trace of a square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is written as $\text{Tr}(\mathbf{A})$ and is just the sum of the diagonal elements:

$$\text{Tr}(\mathbf{A}) = \sum_{i=1}^n \mathbf{A}_{ii} \quad (3.1)$$

The trace of a product can be written as the sum of entry-wise products of elements.

$$\text{Tr}(\mathbf{A}^\top \mathbf{B}) = \text{Tr}(\mathbf{A}\mathbf{B}^\top) = \text{Tr}(\mathbf{B}^\top \mathbf{A}) = \text{Tr}(\mathbf{B}\mathbf{A}^\top) \quad (3.2)$$

$$= \sum_{i,j}^n \mathbf{A}_{i,j} \mathbf{B}_{i,j} \quad (3.3)$$

$$(3.4)$$

Properties:

- Trace of a scalar is a scalar $\text{Tr}(a) = a$
- $\mathbf{A} \in \mathbb{R}^{n \times n}$, $\text{Tr}(\mathbf{A}) = \text{Tr}(\mathbf{A}^\top)$
- $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$, $\text{Tr}(\mathbf{A} + \mathbf{B}) = \text{Tr}(\mathbf{A}) + \text{Tr}(\mathbf{B})$
- $\mathbf{A} \in \mathbb{R}^{n \times n}$, $c \in \mathbb{R}$, $\text{Tr}(c\mathbf{A}) = c \text{Tr}(\mathbf{A})$
- \mathbf{A}, \mathbf{B} such that $\mathbf{A}\mathbf{B}$ is square, $\text{Tr}(\mathbf{A}\mathbf{B}) = \text{Tr}(\mathbf{B}\mathbf{A})$
- $\mathbf{A}, \mathbf{B}, \mathbf{C}$ such that $\mathbf{A}\mathbf{B}\mathbf{C}$ is square, $\text{Tr}(\mathbf{A}\mathbf{B}\mathbf{C}) = \text{Tr}(\mathbf{B}\mathbf{C}\mathbf{A}) = \text{Tr}(\mathbf{C}\mathbf{A}\mathbf{B})$, this is called **trace rotation**.

4 VECTOR NORMS

\mathbf{A} norm of a vector $\|\mathbf{x}\|$ is a measure of it's "length" or "magnitude". The most common is the Euclidean or ℓ_2 norm.

1. ℓ_2 norm : $\|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$

For example, this is used in ridge regression: $\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda\|\boldsymbol{\beta}\|_2^2$

2. ℓ_1 norm : $\|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|$

For example, this is used in ℓ_1 penalized regression: $\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda\|\boldsymbol{\beta}\|_1$

3. ℓ_∞ norm : $\|\mathbf{x}\|_\infty = \max_i |x_i|$

4. The above are all examples of the family of ℓ_p norms : $\|\mathbf{x}\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{\frac{1}{p}}$

5 RANK

A set of vectors $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \subset \mathbb{R}^m$ is said to be linearly independent if no vector can be represented as a linear combination of the remaining vectors. The rank of a matrix is size of the largest subset of columns of \mathbf{A} that constitute a linearly independent set. This is often referred to as the number of linearly independent columns of \mathbf{A} . Note the amazing fact that $\text{rank}(\mathbf{A}) = \text{rank}(\mathbf{A}^\top)$. This means that column rank = row rank. For $\mathbf{A} \in \mathbb{R}^{m \times n}$ $\text{rank}(\mathbf{A}) \leq \min(m, n)$. If $\text{rank}(\mathbf{A}) = \min(m, n)$, then \mathbf{A} is full rank.

6 INVERSE

The inverse of a symmetric matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is written as \mathbf{A}^{-1} and is defined such that:

$$\mathbf{A}\mathbf{A}^{-1} = \mathbf{A}^{-1}\mathbf{A} = \mathbf{I}$$

If \mathbf{A}^{-1} exists, the matrix is said to be **nonsingular**, otherwise it is **singular**. For a square matrix to be invertible, it must be full rank. Non-square matrices are not invertible.

Properties:

- $(\mathbf{A}^{-1})^{-1} = \mathbf{A}$
- $(\mathbf{AB})^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1}$
- $(\mathbf{A}^{-1})^\top = (\mathbf{A}^\top)^{-1}$

Sherman-Morrison-Woodbury Matrix Inversion Lemma

$$(\mathbf{A} + \mathbf{XBX}^\top)^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{X}(\mathbf{B}^{-1} + \mathbf{X}^\top\mathbf{A}^{-1}\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{A}^{-1}$$

This comes up and can often make a hard inverse into an easy inverse. \mathbf{A} and \mathbf{B} are square and invertible but they don't need to be the same dimension.

7 ORTHOGONAL MATRICES

- Two vectors are orthogonal if $\mathbf{u}^\top \mathbf{v} = 0$. A vector is normalized if $\|\mathbf{x}\| = 1$.
- A square matrix is orthogonal if all its columns are orthogonal to each other and are normalized (columns are orthonormal).
- If \mathbf{U} is an orthogonal matrix $\mathbf{U}^\top = \mathbf{U}^{-1}$, then $\mathbf{U}^\top \mathbf{U} = \mathbf{I} = \mathbf{U}\mathbf{U}^\top$.
- Note if \mathbf{U} is not square, but the columns are orthonormal, then $\mathbf{U}^\top \mathbf{U} = \mathbf{I}$ but $\mathbf{U}\mathbf{U}^\top \neq \mathbf{I}$. Orthogonal usually refers to the first case.

8 MATRIX CALCULUS

Gradient

Given $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$ is a function that takes as input a matrix and returns a real values. Then the gradient of f with respect to \mathbf{A} is the matrix of partial derivatives, that is the $m \times n$ matrix defined below.

$$(\nabla_{\mathbf{A}} f(\mathbf{A}))_{ij} = \frac{\partial f(\mathbf{A})}{\partial A_{ij}}$$

Note that the size of $\nabla_{\mathbf{A}} f(\mathbf{A})_{ij}$ is the same as the size of \mathbf{A} .

The gradient of a vector $\mathbf{x} \in \mathbb{R}^n$ is the following:

$$(\nabla_{\mathbf{x}} f(\mathbf{x}))_i = \frac{\partial f(\mathbf{x})}{\partial x_i}$$

The gradient of a function is only defined if that function is real-valued, that is it returns a real scalar value.

Hessian

Given $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a function that takes a vector and returns a real number. Then the Hessian of f with respect to \mathbf{x} is a $n \times n$ matrix of partial derivatives as defined below.

$$(\nabla_{\mathbf{x}}^2 f(\mathbf{x}))_{ij} = \frac{\partial^2 f(\mathbf{x})}{\partial x_i \partial x_j}$$

Just like the gradient, the Hessian is only defined when the function is real-valued. For the purposes of this class, we will only be taking the Hessian of a vector.

Common forms of Derivatives

$$\begin{aligned} \frac{\partial(\mathbf{a}^\top \mathbf{x})}{\partial \mathbf{x}} &= \frac{\partial(\mathbf{x}^\top \mathbf{a})}{\partial \mathbf{x}} = \mathbf{a} \\ \frac{\partial(\mathbf{x}^\top \mathbf{A} \mathbf{x})}{\partial \mathbf{x}} &= (\mathbf{A} + \mathbf{A}^\top) \mathbf{x} \\ \frac{\partial(\mathbf{a}^\top \mathbf{X} \mathbf{b})}{\partial \mathbf{X}} &= \mathbf{a} \mathbf{b}^\top \\ \frac{\partial(\mathbf{a}^\top \mathbf{X}^\top \mathbf{b})}{\partial \mathbf{X}} &= \mathbf{b} \mathbf{a}^\top \\ \frac{\partial(\mathbf{a}^\top \mathbf{X} \mathbf{a})}{\partial \mathbf{X}} &= \frac{\partial(\mathbf{a}^\top \mathbf{X}^\top \mathbf{a})}{\partial \mathbf{X}} = \mathbf{a} \mathbf{a}^\top \end{aligned}$$

$$\begin{aligned}
\frac{\partial(\mathbf{x}^\top \mathbf{A})}{\partial \mathbf{x}} &= \mathbf{A} \\
\frac{\partial(\mathbf{x}^\top)}{\partial \mathbf{x}} &= \mathbf{I} \\
\frac{\partial(\mathbf{A}\mathbf{x})}{\partial z} &= \mathbf{A} \frac{\partial \mathbf{x}}{\partial z} \\
\frac{\partial(\mathbf{X}\mathbf{Y})}{\partial z} &= \mathbf{X} \frac{\partial \mathbf{Y}}{\partial z} + \frac{\partial \mathbf{X}}{\partial z} \mathbf{Y} \\
\frac{\partial(\mathbf{X}^{-1})}{\partial z} &= -\mathbf{X}^{-1} \frac{\partial \mathbf{X}}{\partial z} \mathbf{X}^{-1} \\
\frac{\partial \ln |\mathbf{X}|}{\partial \mathbf{X}} &= (\mathbf{X}^{-1})^\top = (\mathbf{X}^\top)^{-1}
\end{aligned}$$

9 LINEAR REGRESSION

To begin, the likelihood can be derived from a multivariate normal distribution. The likelihood for linear regression is given by:

$$\begin{aligned}
P(\mathcal{D}|\boldsymbol{\beta}, \sigma^2) &= P(\mathbf{y}|\mathbf{X}, \boldsymbol{\beta}, \sigma^2) = \prod_{i=1}^n \mathcal{N}(y_i|\mathbf{x}_i, \boldsymbol{\beta}, \sigma^2) \\
&= (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left(-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right)
\end{aligned}$$

By taking the log and throwing away constants, we get the negative log-likelihood below.

$$-\log P(\mathcal{D}|\boldsymbol{\beta}, \sigma^2) = \frac{n}{2} \log(\sigma^2) + \frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

We can now define the residual sum of squares or least squares.

$$\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\| = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

Maximizing the likelihood is equivalent to minimizing the negative log likelihood and also equivalent to minimizing the residual sum of squares. You will also hear this being called finding the least squares solution. We can rewrite the expression as follows.

$$\begin{aligned}
\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\| &= (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \\
&= \mathbf{y}^\top \mathbf{y} - 2(\mathbf{X}^\top \mathbf{y})^\top \boldsymbol{\beta} + \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta}
\end{aligned}$$

To find the minimum, we first have to take the derivative. Note, we need two matrix derivative identities $\frac{\partial \mathbf{x}^\top \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = (\mathbf{A} + \mathbf{A}^\top) \mathbf{x}$ and $\frac{\partial \mathbf{a}^\top \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a}$. Also, note that $\mathbf{X}^\top \mathbf{X}$ is symmetric.

$$\begin{aligned} & \frac{\partial (\mathbf{y}^\top \mathbf{y} - 2(\mathbf{X}^\top \mathbf{y})^\top \boldsymbol{\beta} + \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \\ &= -2(\mathbf{X}^\top \mathbf{y}) + (\mathbf{X}^\top \mathbf{X} + (\mathbf{X}^\top \mathbf{X})^\top) \boldsymbol{\beta} \\ &= -2(\mathbf{X}^\top \mathbf{y}) + 2\mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} \end{aligned}$$

After setting the derivation equal to zero and solving for $\boldsymbol{\beta}$, we get the following.

$$\begin{aligned} 0 &= -2(\mathbf{X}^\top \mathbf{y}) + 2\mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} \\ \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} &= \mathbf{X}^\top \mathbf{y} \\ \boldsymbol{\beta} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \end{aligned}$$

These are called the normal equations. To solve this in Octave/Matlab, you can implement the equations explicitly using the inverse. However, doing `beta = X \ y`; is a more stable way of solving the normal equations. It does a QR decomposition.

You can check that this solution is the global minimum and not just a stationary point. To do this, you need to evaluate the Hessian, or the second derivative. You should find that the result is a positive definite matrix. And since the Hessian is positive definite, the function is convex and thus the only stationary point is also the global minimum.

10 RIDGE REGRESSION

Now, we're going to derive ridge regression in a similar way. Recall that for linear regression, we found the MLE from forming the likelihood $P(\mathbf{y}|\boldsymbol{\beta})$. Here, we can derive the MAP estimate from the posterior which is constructed from the likelihood and the prior. Let $\boldsymbol{\beta} \sim \mathcal{N}(0, \frac{1}{\lambda} \mathbf{I}_p)$ be a prior on the parameter vector $\boldsymbol{\beta}$ where \mathbf{I}_p is an identity matrix of size p . The form of the posterior is given below.

$$\begin{aligned} P(\boldsymbol{\beta}|\mathbf{y}) &\propto P(\mathbf{y}|\boldsymbol{\beta})P(\boldsymbol{\beta}) \\ &\propto \mathcal{N}(\mathbf{y}|\mathbf{X}, \boldsymbol{\beta}, \sigma^2) \mathcal{N}(0, \frac{1}{\lambda} \mathbf{I}_p) \end{aligned}$$

Given that $\sigma^2 = 1$, we first want to derive the posterior for $\boldsymbol{\beta}$.

$$\begin{aligned}
P(\boldsymbol{\beta}|\mathbf{y}) &\propto P(\mathbf{y}|\boldsymbol{\beta})P(\boldsymbol{\beta}) \\
&\propto \mathcal{N}(\mathbf{y}|\mathbf{X}\boldsymbol{\beta}, \sigma^2)\mathcal{N}\left(0, \frac{1}{\lambda}\mathbf{I}_p\right) \\
&\propto (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left(-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right) \cdot (2\pi)^{-\frac{p}{2}} \left|\frac{1}{\lambda}\mathbf{I}_p\right|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}\boldsymbol{\beta}^\top\left(\frac{1}{\lambda}\mathbf{I}_p\right)^{-1}\boldsymbol{\beta}\right) \\
&\propto (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left(-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right) \cdot (2\pi)^{-\frac{p}{2}}\left(\frac{1}{\lambda}\right)^{-\frac{p}{2}} \exp\left(-\frac{\lambda}{2}\boldsymbol{\beta}^\top\mathbf{I}_p\boldsymbol{\beta}\right) \\
&\propto (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left(-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right) \cdot (2\pi\lambda^{-1})^{-\frac{p}{2}} \exp\left(-\frac{\lambda}{2}\boldsymbol{\beta}^\top\boldsymbol{\beta}\right) \\
&\propto (2\pi\sigma^2)^{-\frac{n}{2}} (2\pi\lambda^{-1})^{-\frac{p}{2}} \exp\left(-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) - \frac{\lambda}{2}\boldsymbol{\beta}^\top\boldsymbol{\beta}\right)
\end{aligned}$$

Taking the negative log and dropping constants, we get:

$$\begin{aligned}
&\propto \frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \frac{\lambda}{2}\boldsymbol{\beta}^\top\boldsymbol{\beta} \\
&\propto (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda\boldsymbol{\beta}^\top\boldsymbol{\beta} \quad \text{Setting } \sigma^2 \text{ to 1 and dropping more constants.}
\end{aligned}$$

Now, since we wanted to maximize the posterior, we now need to minimize the negative log of the posterior. Note that minimizing the above expression is exactly the same as finding the ridge solution by minimizing the sum of squares plus the l_2 penalty (Eq. 10.1). These two expressions are equivalent, and thus minimizing them will yield identical solutions.

$$\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda\|\boldsymbol{\beta}\|_2^2 = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda\boldsymbol{\beta}^\top\boldsymbol{\beta} \quad (10.1)$$

Let's expand out and write the loss function in matrix form.

$$\begin{aligned}
&(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda\boldsymbol{\beta}^\top\boldsymbol{\beta} \\
&= \mathbf{y}^\top\mathbf{y} - 2(\mathbf{X}^\top\mathbf{y})^\top\boldsymbol{\beta} + \boldsymbol{\beta}^\top\mathbf{X}^\top\mathbf{X}\boldsymbol{\beta} + \lambda\boldsymbol{\beta}^\top\boldsymbol{\beta}
\end{aligned}$$

To find the value of $\boldsymbol{\beta}$ that minimizes the loss function, we first have to take the derivative. Note, we need two matrix derivative identities $\frac{\partial \mathbf{x}^\top \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = (\mathbf{A} + \mathbf{A}^\top)\mathbf{x}$ and $\frac{\partial \mathbf{a}^\top \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a}$. Also, note that $\mathbf{X}^\top \mathbf{X}$ is symmetric.

$$\begin{aligned}
& \frac{\partial}{\partial \boldsymbol{\beta}} (\mathbf{y}^\top \mathbf{y} - 2(\mathbf{X}^\top \mathbf{y})^\top \boldsymbol{\beta} + \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} + \lambda \boldsymbol{\beta}^\top \boldsymbol{\beta}) \\
&= -2(\mathbf{X}^\top \mathbf{y}) + (\mathbf{X}^\top \mathbf{X} + (\mathbf{X}^\top \mathbf{X})^\top) \boldsymbol{\beta} + 2\lambda \boldsymbol{\beta} \\
&= -2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} + 2\lambda \boldsymbol{\beta}
\end{aligned}$$

After setting the derivation equal to zero and solving for $\boldsymbol{\beta}$, we get the following.

$$\begin{aligned}
0 &= -2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} + 2\lambda \boldsymbol{\beta} \\
\mathbf{X}^\top \mathbf{y} &= \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} + \lambda \boldsymbol{\beta} \\
\mathbf{X}^\top \mathbf{y} &= (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}) \boldsymbol{\beta} \\
\boldsymbol{\beta} &= (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}
\end{aligned}$$

Just like linear regression, you can implement the equations explicitly in Matlab/Octave. In practice, you might have trouble calculating the inverse directly if the matrix is huge and λ is small. We can also derive a numerically stable way of computing $\boldsymbol{\beta}$ using the backslash operator. Define $\tilde{\mathbf{X}}$ and $\tilde{\mathbf{y}}$ such that $\boldsymbol{\beta}$ can be written as:

$$\boldsymbol{\beta} = (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}} \tilde{\mathbf{y}} \quad (10.2)$$

Then, you can use the backslash operator as shown below.

```
Xtil = [X; sqrt(lambda)*eye(p)];
ytil=[y; zeros(p,1)];
beta = Xtil\ytil;
```

11 QUADRATIC FORMS

For a square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ and a vector $\mathbf{x} \in \mathbb{R}^n$, the scalar value $\mathbf{x}^\top \mathbf{A} \mathbf{x}$ is referred to as quadratic form. We can write it explicitly as follows:

$$\mathbf{x}^\top \mathbf{A} \mathbf{x} = \sum_{i=1}^n x_i (\mathbf{A} \mathbf{x})_i = \sum_{i=1}^n x_i \left(\sum_{j=1}^n A_{ij} x_j \right) = \sum_{i=1}^n \sum_{j=1}^n A_{ij} x_i x_j$$

11.1 DEFINITIONS

Positive Definite (PD)

notation: $\mathbf{A} > 0$ or $\mathbf{A} \succ 0$ and the set of all positive definite matrices \mathbb{S}_{++}^n .

A symmetric matrix $\mathbf{A} \in \mathbb{S}^n$ is positive definite if for all non-zero vectors $\mathbf{x} \in \mathbb{R}$, $\mathbf{x}^\top \mathbf{A} \mathbf{x} > 0$.

Positive Semidefinite (PSD) notation: $\mathbf{A} \geq 0$ or $\mathbf{A} \succeq 0$ and the set of all positive semidefinite matrices \mathbb{S}_+^n .

A symmetric matrix $\mathbf{A} \in \mathbb{S}^n$ is positive semidefinite if for all non-zero vectors $\mathbf{x} \in \mathbb{R}$, $\mathbf{x}^\top \mathbf{A} \mathbf{x} \geq 0$.

Negative Definite (ND) notation: $\mathbf{A} < 0$ or $\mathbf{A} \prec 0$.

Similarly, a symmetric matrix $\mathbf{A} \in \mathbb{S}^n$ is negative definite if for all non-zero vectors $\mathbf{x} \in \mathbb{R}$, $\mathbf{x}^\top \mathbf{A} \mathbf{x} < 0$.

Negative Semidefinite (NSD) notation: $\mathbf{A} \leq 0$ or $\mathbf{A} \preceq 0$.

Similarly, a symmetric matrix $\mathbf{A} \in \mathbb{S}^n$ is negative semidefinite if for all non-zero vectors $\mathbf{x} \in \mathbb{R}$, $\mathbf{x}^\top \mathbf{A} \mathbf{x} \leq 0$.

Indefinite Lastly, a symmetric matrix $\mathbf{A} \in \mathbb{S}^n$ is indefinite if it is neither positive semidefinite nor negative semidefinite, that is if there exists $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}$ such that $\mathbf{x}_1^\top \mathbf{A} \mathbf{x}_1 > 0$ and $\mathbf{x}_2^\top \mathbf{A} \mathbf{x}_2 < 0$.

If \mathbf{A} is positive definite, then $-\mathbf{A}$ is negative definite and vice versa. The same can be said about positive semidefinite and negative semidefinite. Also, positive definite and negative definite matrices are always full rank and invertible.

12 EIGENVALUES AND EIGENVECTORS

Given a square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$, $\lambda \in \mathbb{C}$ is an **eigenvalue** and $\mathbf{x} \in \mathbb{C}$ (complex set of numbers) the corresponding eigenvector if

$$\mathbf{A} \mathbf{x} = \lambda \mathbf{x}, \mathbf{x} \neq 0$$

This condition can be rewritten as:

$$(\mathbf{A} - \lambda \mathbf{I}) \mathbf{x} = 0$$

where \mathbf{I} is the identity matrix. Now for a non-zero vector to satisfy this equation, then $(\mathbf{A} - \lambda \mathbf{I})$ must not be invertible, which means that it is singular and the determinant is zero.

You can use the definition of the determinant to expand this expression into a polynomial in λ and then find the roots (real or complex) of the polynomial to find the n eigenvalues $\lambda_1, \dots, \lambda_n$. Once you have the eigenvalues λ_i , you can find the corresponding eigenvector by solving the system of equations $(\lambda_i \mathbf{I} - \mathbf{A})\mathbf{x} = 0$.

12.1 PROPERTIES

- The trace of a matrix \mathbf{A} is equal to the sum of its eigenvalues:

$$\text{Tr}(\mathbf{A}) = \sum_{i=1}^n \lambda_i$$

- The determinant of \mathbf{A} is equal to the product of its eigenvalues

$$|\mathbf{A}| = \prod_{i=1}^n \lambda_i$$

- The rank of \mathbf{A} is equal to the number of non-zero eigenvalues of \mathbf{A}
- The eigenvalues of a diagonal matrix $\mathbf{D} = \text{diag}(d_1, \dots, d_n)$ are just the diagonal entries d_1, \dots, d_n

12.2 DIAGONALIZATION

A square matrix \mathbf{A} is said to be diagonalizable if it is similar to a diagonal matrix. A diagonal matrix \mathbf{A} has the property that there exists an invertible matrix \mathbf{X} and a diagonal matrix Λ such that $\mathbf{A} = \mathbf{X}\Lambda\mathbf{X}^{-1}$.

We can write all the eigenvector equations simultaneously as $\mathbf{A}\mathbf{X} = \mathbf{X}\Lambda$ where the columns of $\mathbf{X} \in \mathbb{R}^{n \times n}$ are the eigenvectors of \mathbf{A} and Λ is a diagonal matrix whose entries are the eigenvalues of \mathbf{A} . If the eigenvectors of \mathbf{A} are linearly independent, then the matrix \mathbf{X} will be invertible, so $\mathbf{A} = \mathbf{X}\Lambda\mathbf{X}^{-1}$. This is known as the **eigenvalue decomposition** of the matrix.

Why is this useful? Because powers of diagonal matrices are easy to compute. Try computing \mathbf{A}^3 . Also, remember this form $\mathbf{A} = \mathbf{X}\Lambda\mathbf{X}^{-1} = \mathbf{X}\Lambda\mathbf{X}^\top = \sum_{i=1}^n \lambda_i \mathbf{x}_i \mathbf{x}_i^\top$. We will see this later when we cover SVMs with kernels: $\sum_{i=1}^n \lambda_i \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^\top$.

12.3 PROPERTIES OF EIGENVALUES/EIGENVECTORS FOR SYMMETRIC MATRICES

- For a symmetric matrix $\mathbf{A} \in \mathbb{S}^n$, all the eigenvalues are real.
- The eigenvectors of \mathbf{A} are orthonormal so that means the matrix \mathbf{X} is an orthogonal matrix (so we can denote the matrix of eigenvectors as \mathbf{U}).

We can then write

$$\mathbf{A} = \mathbf{X}\mathbf{\Lambda}\mathbf{X}^{-1}$$

$$\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^{-1} \text{ The inverse of an orthogonal matrix is just the inverse.}$$

$$\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^{\top}$$

This means that

$$\begin{aligned}\mathbf{x}^{\top}\mathbf{A}\mathbf{x} &= \mathbf{x}^{\top}\mathbf{U}\mathbf{\Lambda}\mathbf{U}^{\top}\mathbf{x} \\ &= \mathbf{y}^{\top}\mathbf{\Lambda}\mathbf{y} \\ &= \sum_{i=1}^n \lambda_i y_i^2\end{aligned}$$

Since y_i^2 is always positive, the sign of this expression depends entirely on the λ_i 's. If all $\lambda_i > 0$, then the matrix is positive definite; if all $\lambda_i \geq 0$, then \mathbf{A} is positive semidefinite. If $\lambda_i < 0$ and $\lambda_i \leq 0$, then the matrix is negative definite or negative semidefinite respectively. If \mathbf{A} has both positive and negative eigenvalues, then it is indefinite.

13 SINGULAR VALUE DECOMPOSITION

Any $n \times m$ matrix \mathbf{A} can be written as

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^{\top}$$

where

$$\mathbf{U} = \text{eigenvectors of } \mathbf{A}\mathbf{A}^{\top} \text{ (} n \times n \text{)}$$

$$\mathbf{\Sigma} = \sqrt{\text{diag}(\text{eig}(\mathbf{A}\mathbf{A}^{\top}))} \text{ (} n \times m \text{)}$$

$$\mathbf{V} = \text{eigenvectors of } \mathbf{A}^{\top}\mathbf{A} \text{ (} m \times m \text{)}$$

13.1 PROPERTIES

$$\mathbf{U}^{\top}\mathbf{U} = \mathbf{I}$$

$$\mathbf{U}\mathbf{U}^{\top} = \mathbf{I}$$

$$\mathbf{V}^{\top}\mathbf{V} = \mathbf{I}$$

$$\mathbf{V}\mathbf{V}^{\top} = \mathbf{I}$$

However, if you do the economy SVD, all the above properties are true except $\mathbf{U}\mathbf{U}^{\top} \neq \mathbf{I}$.

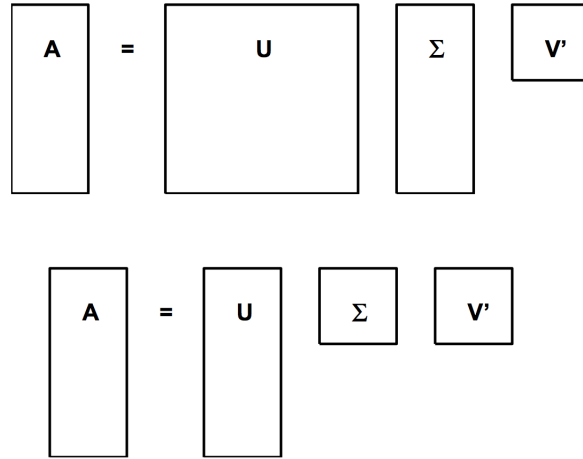


Figure 13.1: Taken from Matrix Cookbook.

13.2 RELATION TO EIGENVALUE DECOMPOSITION

$$\begin{aligned} \mathbf{A}^\top \mathbf{A} &= \mathbf{V} \boldsymbol{\Sigma}^\top \mathbf{U}^\top \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^\top = \mathbf{V} \boldsymbol{\Sigma}^2 \mathbf{V}^\top \\ \mathbf{A} \mathbf{A}^\top &= \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^\top \mathbf{V} \boldsymbol{\Sigma}^\top \mathbf{U}^\top = \mathbf{U} \boldsymbol{\Sigma}^2 \mathbf{U}^\top \end{aligned}$$

The columns of \mathbf{V} are the eigenvectors of $\mathbf{A}^\top \mathbf{A}$.

The columns of \mathbf{U} are the eigenvectors of $\mathbf{A} \mathbf{A}^\top$.

The values of $\boldsymbol{\Sigma}$, σ_i are the square roots of the eigenvalues of $\mathbf{A}^\top \mathbf{A}$ or $\mathbf{A} \mathbf{A}^\top$, so $\sigma_i = \sqrt{\lambda_i}$

14 PRINCIPAL COMPONENTS ANALYSIS

Often times when we have data in high-dimensional space, we can actually reduce the dimensions considerably while still capturing most of the variance of the data. This is called dimensionality reduction and one of the approaches is to use principal component analysis or PCA. PCA basically approximates some real $m \times n$ matrix \mathbf{A} with the sum of some simple matrices that are rank one outer products.

The SVD of matrix \mathbf{A} can be written:

$$\mathbf{A} = \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^\top$$

where

$$\mathbf{A} = \mathbf{E}_1 + \mathbf{E}_2 + \cdots + \mathbf{E}_p,$$

where $p = \min(m, n)$. The component matrices \mathbf{E}_i are rank one outer products:

$$\mathbf{E}_i = \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$$

The component matrices are orthogonal to each other so, the product is 0.

$$\mathbf{E}_j \mathbf{E}_k^\top = 0, \text{ where } j \neq k$$

The norm of each component matrix is the corresponding singular value.

$$\|\mathbf{E}\|_i = \sigma_i$$

So, the contribution that each component makes to reproducing A is determined by the size of the singular value. So, if you wanted to figure out how many components to include, you can plot the singular values and then cut it off where there is a significant drop in the value.

15 REFERENCES

The following are my sources for this tutorial and you should check them out for further reading.

Zico Kolter's Linear Algebra Review and Reference

<http://cs229.stanford.edu/section/cs229-linalg.pdf>

The Matrix Cookbook

<http://orion.uwaterloo.ca/~hwoikowi/matrixcookbook.pdf>

Matlab's Eigenvalues and Singular Values

<http://www.mathworks.com/moler/eigs.pdf>

Course Notes from Harvard on Eigenvalues and Eigenvectors

http://www.math.harvard.edu/archive/20_spring_05/handouts/ch05_notes.pdf

Machine Learning: A Probabilistic Perspective by Kevin Murphy

<http://www.cs.ubc.ca/~murphyk/MLbook/index.html>