

Recitation 2: Matrix Derivatives

2.1 Minimizing/Maximizing Functions

In machine learning, many problems we'll want to solve can be cast as minimizing (or maximizing) a function $J : \mathbb{R}^d \rightarrow \mathbb{R}$.

Examples:

- the sum of square errors in least squares: $J(\theta) = \|\mathbf{y} - \mathbf{X}\theta\|_2^2$
- the sum of square distances to cluster centers in KMeans: $J_{avg^2} = \sum_{j=1}^k \sum_{x \in C_j} d(\mathbf{x}, \mathbf{m}_j)^2$
- the Rayleigh Quotient: $\frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{x}}$

Recall from lecture, we can often use gradient descent to minimize such a function J

```

 $\theta \leftarrow \mathbf{0}$ 
until {convergence} {
   $\theta \leftarrow \theta - \alpha \nabla_{\theta} J(\theta)$ 
}

```

So we just need to be able to compute $\nabla_{\theta} J$. Note: the subscripted θ in $\nabla_{\theta} J$ means that we are taking the gradient or the derivative with respect to θ . When the context is clear, we may drop the θ subscript in $\nabla_{\theta} J$ for convenience. Suppose $\theta \in \mathbb{R}^d$. This gradient ∇J is given by:

$$\nabla J = \begin{bmatrix} \frac{\partial J}{\partial \theta_1} \\ \frac{\partial J}{\partial \theta_2} \\ \vdots \\ \frac{\partial J}{\partial \theta_d} \end{bmatrix}$$

2.1.1 Useful Matrix Derivatives

Most of the kinds of functions we'll end up computing the gradient of can be expressed in terms of matrix-vector and vector-vector operations. So first we'll take a look at some useful matrix derivatives that show up over and over again. In the following identities, let $\mathbf{w} \in \mathbb{R}^d$, $\mathbf{x} \in \mathbb{R}^d$, $\mathbf{A} \in \mathbb{R}^{d \times d}$.

$$\nabla_{\mathbf{x}}(\mathbf{w}^T \mathbf{x}) = \mathbf{w} \quad (2.1)$$

$$\nabla_{\mathbf{x}}(\mathbf{x}^T \mathbf{w}) = \mathbf{w} \quad (2.2)$$

$$\nabla_{\mathbf{x}}(\mathbf{x}^T \mathbf{x}) = 2\mathbf{x} \quad (2.3)$$

$$\nabla_{\mathbf{x}}(\mathbf{x}^T \mathbf{A} \mathbf{x}) = (\mathbf{A} + \mathbf{A}^T)\mathbf{x} \quad (2.4)$$

Proof: $\nabla_{\mathbf{x}}(\mathbf{w}^T \mathbf{x}) = \mathbf{w}$

Consider the i th index of $\nabla_{\mathbf{x}}(\mathbf{w}^T \mathbf{x})$:

$$\frac{\partial}{\partial x_i}(\mathbf{w}^T \mathbf{x}) = \frac{\partial}{\partial x_i} \left(\sum_j w_j x_j \right) = w_i$$

Repeating this for all indices of vector \mathbf{x} and plugging this back into the definition of the gradient, we see that:

$$\nabla_{\mathbf{x}}(\mathbf{w}^T \mathbf{x}) = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix} = \mathbf{w}$$

The derivation for $\nabla_{\mathbf{x}}(\mathbf{x}^T \mathbf{w}) = \mathbf{w}$ and $\nabla_{\mathbf{x}}(\mathbf{x}^T \mathbf{x}) = 2\mathbf{x}$ are identical. ■

Proof: $\nabla_{\mathbf{x}}(\mathbf{x}^T \mathbf{A} \mathbf{x}) = (\mathbf{A} + \mathbf{A}^T)\mathbf{x}$

$$\begin{aligned} \mathbf{x}^T \mathbf{A} \mathbf{x} &= \sum_i \sum_j A_{i,j} x_i x_j \\ &= A_{k,k} x_k^2 + \sum_{j \neq k} A_{k,j} x_k x_j + \sum_{i \neq k} A_{i,k} x_i x_k + \sum_{i \neq k} \sum_{j \neq k} A_{i,j} x_i x_j \\ \frac{\partial}{\partial x_k}(\mathbf{x}^T \mathbf{A} \mathbf{x}) &= 2A_{k,k} x_k + \sum_{j \neq k} A_{k,j} x_j + \sum_{i \neq k} A_{i,k} x_i + 0 \\ \frac{\partial}{\partial x_k}(\mathbf{x}^T \mathbf{A} \mathbf{x}) &= \sum_j A_{k,j} x_j + \sum_i A_{i,k} x_i \\ &= [\mathbf{A} \mathbf{x}]_k + [\mathbf{A}^T \mathbf{x}]_k \end{aligned}$$

where $[\mathbf{A} \mathbf{x}]_k$ denotes index k of the vector $\mathbf{A} \mathbf{x}$. Repeating this for all other indices and plugging them into the definition of the gradient, we see that: $\nabla_{\mathbf{x}}(\mathbf{x}^T \mathbf{A} \mathbf{x}) = (\mathbf{A} + \mathbf{A}^T)\mathbf{x}$ ■

2.2 Least Squares

In the least squares problem, we are given data: $\mathbf{X} \in \mathbb{R}^{n \times d}$, $\mathbf{y} \in \mathbb{R}^n$. Our goal is to find some weight vector $\boldsymbol{\theta} \in \mathbb{R}^d$ that minimizes the sum of square errors:

$$J(\boldsymbol{\theta}) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|^2 = (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})$$

Before we compute the gradient, expand J :

$$\begin{aligned} J(\boldsymbol{\theta}) &= (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) \\ &= (\mathbf{y}^T - \boldsymbol{\theta}^T \mathbf{X}^T) (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) \\ &= \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X}\boldsymbol{\theta} - \boldsymbol{\theta}^T \mathbf{X}^T \mathbf{y} + \boldsymbol{\theta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\theta} \end{aligned}$$

Notice that every term is now in the form matrix/vector operations that we already saw in the matrix

identities from the previous section. Taking the gradients of each of the individual terms:

$$\begin{aligned}
 \nabla_{\theta}(\mathbf{y}^T \mathbf{y}) &= 0 \\
 \nabla_{\theta}(-\mathbf{y}^T \mathbf{X} \theta) &= -\mathbf{X}^T \mathbf{y} && \text{(Apply 2.1 with } \mathbf{w} = \mathbf{X}^T \mathbf{y}\text{)} \\
 \nabla_{\theta}(-\theta^T \mathbf{X}^T \mathbf{y}) &= -\mathbf{X}^T \mathbf{y} && \text{(Apply 2.2 with } \mathbf{w} = \mathbf{X}^T \mathbf{y}\text{)} \\
 \nabla_{\theta}(\theta^T \mathbf{X}^T \mathbf{X} \theta) &= (\mathbf{X}^T \mathbf{X} + (\mathbf{X}^T \mathbf{X})^T) \theta && \text{(Apply 2.4 with } \mathbf{A} = \mathbf{X}^T \mathbf{X}\text{)} \\
 &= 2\mathbf{X}^T \mathbf{X} \theta
 \end{aligned}$$

Putting this all together:

$$\nabla_{\theta} J = -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \theta$$

At a minimum, the gradient will be $\mathbf{0}$. So setting this to $\mathbf{0}$ and solving for θ gives us the normal equations:

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

Note: during recitation I had incorrectly written: $\nabla_{\theta}(-\mathbf{y}^T \mathbf{X} \theta) = -\mathbf{y}^T \mathbf{X}$. A student corrected me on this, but I didn't register this even after he explained it a few times. Apologies. This is now corrected in the derivation above.

2.3 Rayleigh Quotient using Lagrange Multipliers

In the first homework, we asked you to prove that given a real symmetric matrix A , the maximal eigenvector maximizes the Rayleigh Quotient: $\frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{x}}$.

Suppose we added the constraint that \mathbf{x} be unit length, $\mathbf{x}^T \mathbf{x} = 1$, allowing us to ignore the $\mathbf{x}^T \mathbf{x}$ in the denominator. Note that $\frac{\mathbf{x}}{\mathbf{x}^T \mathbf{x}}$ is already a unit vector in the direction of \mathbf{x} so we haven't really changed our problem by forcing \mathbf{x} to be a unit vector. Now this is precisely the sort of constrained optimization problem that we can use Lagrange multipliers to solve.

Let the function we'd like to minimize be: $f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x}$, and the constraint we must satisfy be $g(\mathbf{x}) = 0$, where $g(\mathbf{x}) = \mathbf{x}^T \mathbf{x} - 1$. Applying the method of [Lagrange multipliers](#) gives us:

$$\begin{aligned}
 \nabla f &= \lambda \nabla g \\
 \nabla(\mathbf{x}^T \mathbf{A} \mathbf{x}) &= \lambda \nabla(\mathbf{x}^T \mathbf{x} - 1) \\
 (\mathbf{A} + \mathbf{A}^T) \mathbf{x} &= 2\lambda \mathbf{x} \\
 2\mathbf{A} \mathbf{x} &= 2\lambda \mathbf{x} \\
 \mathbf{A} \mathbf{x} &= \lambda \mathbf{x}
 \end{aligned}$$

This tells us that the unit vector \mathbf{x} maximizing the Rayleigh quotient must be an eigenvector of \mathbf{A} . So the maximal eigenvector will maximize the Rayleigh Quotient.