

Problem 1: Gradients and Hessians

1. Using that \mathbf{A} is symmetric, we obtain:

$$\begin{aligned} \frac{1}{2} \mathbf{x}^\top \mathbf{A} \mathbf{x} &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n A_{ij} x_j x_i \\ \implies \frac{1}{2} \left(\nabla_x \mathbf{x}^\top \mathbf{A} \mathbf{x} \right)_k &= \frac{1}{2} \frac{\partial \mathbf{x}^\top \mathbf{A} \mathbf{x}}{\partial x_k} = \frac{1}{2} \left(\sum_{i=1}^n A_{ik} x_i + \sum_{j=1}^n A_{kj} x_j \right) = \sum_{i=1}^n A_{ki} x_i = (\mathbf{A} \mathbf{x})_k \\ \implies \frac{1}{2} \nabla_x \mathbf{x}^\top \mathbf{A} \mathbf{x} &= \mathbf{A} \mathbf{x}. \end{aligned}$$

Also:

$$\begin{aligned} \mathbf{b}^\top \mathbf{x} &= \sum_{i=1}^n b_i x_i \\ \implies \left(\nabla_x \mathbf{b}^\top \mathbf{x} \right)_k &= \frac{\partial \mathbf{b}^\top \mathbf{x}}{\partial x_k} = b_k \\ \implies \nabla_x \mathbf{b}^\top \mathbf{x} &= \mathbf{b}. \end{aligned}$$

Therefore, $\nabla_x f(x) = \mathbf{A} \mathbf{x} + \mathbf{b}$ as desired.

2. Using the chain rule for derivative, we can calculate each element of the gradients as follows:

$$[\nabla_x f(x)]_k = \frac{\partial f(x)}{\partial x_k} = \frac{\partial g(h(x))}{\partial x_k} = \frac{\partial g(y)}{\partial y} \frac{\partial h(x)}{\partial x_k}.$$

This implies: $\nabla_x f(x) = g'(h(x)) \nabla_x h(\mathbf{x})$.

3. Using the result in point 1, we can calculate each component of the Hessian as follows:

$$[\nabla_x^2 f(x)]_{ij} = \frac{\partial [\nabla_x f(x)]_j}{\partial x_i} = \frac{\partial (\sum_{k=1}^n A_{jk} x_k + b_j)}{\partial x_i} = A_{ji} = A_{ij}.$$

Therefore, $\nabla_x^2 f(x) = \mathbf{A}$.

4. Using the results in point 2 and then the result in point 1 we obtain:

$$\nabla_{\mathbf{x}} g(\mathbf{a}^\top \mathbf{x}) = g'(\mathbf{a}^\top \mathbf{x}) \mathbf{a}.$$

Taking the derivatives again, we get:

$$\nabla_{\mathbf{x}}^2 g(\mathbf{a}^\top \mathbf{x}) = g''(\mathbf{a}^\top \mathbf{x}) \mathbf{a} \mathbf{a}^\top.$$

Problem 2: Positive definite matrices

1. First check that $\mathbf{z}\mathbf{z}^\top$ is symmetric:

$$\left(\mathbf{z}\mathbf{z}^\top\right)^\top = \left(\mathbf{z}^\top\right)^\top \mathbf{z}^\top = \mathbf{z}\mathbf{z}^\top.$$

Now, let's check that it is positive semidefinite. Take $\mathbf{x} \in \mathbb{R}^n$, then

$$\mathbf{x}^\top \mathbf{z}\mathbf{z}^\top \mathbf{x} = \left(\mathbf{z}^\top \mathbf{x}\right)^\top \left(\mathbf{z}^\top \mathbf{x}\right) = \left\|\mathbf{z}^\top \mathbf{x}\right\|_2^2 \geq 0.$$

2. By definition, $\mathbf{x} \in \mathcal{N}(\mathbf{A})$ iff $\mathbf{z}\mathbf{z}^\top \mathbf{x} = \mathbf{0}$. This is possible iff $\mathbf{z}^\top \mathbf{x} = 0$. Therefore,

$$\mathcal{N}(\mathbf{A}) = \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{z}^\top \mathbf{x} = 0\}.$$

These are all the vectors in \mathbb{R}^n that are perpendicular to \mathbf{z} .

All columns of $\mathbf{z}\mathbf{z}^\top$ are collinear and are just a multiple of \mathbf{z} . Therefore, $\text{rank}(\mathbf{A}) = 1$.

3. First, check the symmetry: $(\mathbf{B}\mathbf{A}\mathbf{B}^\top)^\top = \mathbf{B}(\mathbf{B}\mathbf{A})^\top = \mathbf{B}\mathbf{A}^\top \mathbf{B}^\top = \mathbf{B}\mathbf{A}\mathbf{B}^\top$. Next, let $\mathbf{x} \in \mathbb{R}^m$ and $\mathbf{y} = \mathbf{B}^\top \mathbf{x} \in \mathbb{R}^n$. Then, because \mathbf{A} is PSD, $\mathbf{x}^\top \mathbf{B}\mathbf{A}\mathbf{B}^\top \mathbf{x} = \mathbf{y}^\top \mathbf{A}\mathbf{y} \geq 0$. Therefore, $\mathbf{B}\mathbf{A}\mathbf{B}^\top$ is PSD.

Problem 3: Eigenvectors, eigenvalues, and the spectral theorem

1. Let $\mathbf{e}_i \in \mathbb{R}^n$ denote the i th canonical basis vector of \mathbb{R}^n , $[\mathbf{e}_i]_j = \delta[i - j]$. Then,

$$\mathbf{A}\mathbf{t}_i = \mathbf{T}\mathbf{\Lambda}\mathbf{T}^{-1}\mathbf{t}_i = \mathbf{T}\mathbf{\Lambda}\mathbf{T}^{-1}\mathbf{T}\mathbf{e}_i = \mathbf{T}\mathbf{\Lambda}\mathbf{e}_i = \mathbf{T}\lambda_i\mathbf{e}_i = \lambda_i\mathbf{t}_i.$$

2. Follows directly from the problem before, as $\mathbf{U}^{-1} = \mathbf{U}^\top$.
3. Suppose there exists $i \in \{1 \dots n\}$ with $\lambda_i < 0$. Let $\mathbf{v}_i, \mathbf{v}_i \neq 0$ be a corresponding eigenvector. Then $\mathbf{v}_i^\top \mathbf{A}\mathbf{v}_i = \mathbf{v}_i^\top \lambda_i \mathbf{v}_i = \lambda_i \|\mathbf{v}_i\|_2^2 < 0$. But as \mathbf{A} is PSD $\mathbf{x}^\top \mathbf{A}\mathbf{x} \geq 0$ for all $\mathbf{x} \in \mathbb{R}^n$. This is a contradiction and therefore $\lambda_i \geq 0$, for all $i \in \{1 \dots n\}$.

Problem 5: Linear regression and gradient descent

See `ps1_sol_problem5.ipynb` file.

Problem 6: High-dimensional regression (exam practice)

1. The least squares solution is defined as

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$$

where

$$J(\boldsymbol{\theta}) = \frac{1}{2}(\mathbf{X}\boldsymbol{\theta} - \mathbf{y})^\top (\mathbf{X}\boldsymbol{\theta} - \mathbf{y}).$$

To minimize $J(\boldsymbol{\theta})$, let's calculate $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$:

$$\begin{aligned}\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) &= \nabla_{\boldsymbol{\theta}} \frac{1}{2}(\mathbf{X}\boldsymbol{\theta} - \mathbf{y})^{\top}(\mathbf{X}\boldsymbol{\theta} - \mathbf{y}) \\ &= \frac{1}{2} \nabla_{\boldsymbol{\theta}}(\boldsymbol{\theta}^{\top} \mathbf{X}^{\top} \mathbf{X} \boldsymbol{\theta} - \boldsymbol{\theta}^{\top} \mathbf{X}^{\top} \mathbf{y} - \mathbf{y}^{\top} \mathbf{X} \boldsymbol{\theta} + \mathbf{y}^{\top} \mathbf{y}) \\ &= \frac{1}{2} \nabla_{\boldsymbol{\theta}}(\boldsymbol{\theta}^{\top} \mathbf{X}^{\top} \mathbf{X} \boldsymbol{\theta} - 2\boldsymbol{\theta}^{\top} \mathbf{X}^{\top} \mathbf{y} + \mathbf{y}^{\top} \mathbf{y}) \\ &= \frac{1}{2}(2\mathbf{X}^{\top} \mathbf{X} \boldsymbol{\theta} - 2\mathbf{X}^{\top} \mathbf{y}).\end{aligned}$$

To minimize $J(\boldsymbol{\theta})$ we set $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbf{0}$ which yields:

$$\mathbf{X}^{\top}(\mathbf{X}\hat{\boldsymbol{\theta}} - \mathbf{y}) = \mathbf{0}.$$

Above, $\hat{\boldsymbol{\theta}}$ denotes the optimal value for $\boldsymbol{\theta}$, i.e. the one that minimizes $J(\boldsymbol{\theta})$.

Rearranging the terms we obtain the normal equations:

$$\mathbf{X}^{\top} \mathbf{X} \hat{\boldsymbol{\theta}} = \mathbf{X}^{\top} \mathbf{y}.$$

So the value of $\hat{\boldsymbol{\theta}}$ that minimizes $J(\boldsymbol{\theta})$ is given in closed form by the equation:

$$\hat{\boldsymbol{\theta}} = (\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{X}^{\top} \mathbf{y}.$$

2. The $p \times p$ matrix $\mathbf{X}^{\top} \mathbf{X}$ needs to be full-rank.

This is possible when all the columns of \mathbf{X} are linearly independent. The necessary condition is $n \geq p$.

If $n < p$, $\mathbf{X}^{\top} \mathbf{X}$ cannot be full-rank. Here is why:

$$\text{rank}(\mathbf{X}^{\top} \mathbf{X}) \leq \min(\text{rank}(\mathbf{X}^{\top}), \text{rank}(\mathbf{X})) \leq \min(n, p) = n < p. \quad (1)$$

3. First note,

$$\begin{aligned}\hat{\boldsymbol{\theta}} &= (\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{X}^{\top} \mathbf{y} \\ &= (\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{X}^{\top} (\mathbf{X}\boldsymbol{\theta} + \mathbf{m}) \\ &= \boldsymbol{\theta} + (\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{X}^{\top} \mathbf{m}.\end{aligned}$$

Therefore,

$$\mathbb{E}[\hat{\boldsymbol{\theta}}] = \boldsymbol{\theta} + \mathbb{E}[(\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{X}^{\top} \mathbf{m}] = \boldsymbol{\theta}.$$