**Problem 1.** In our proof of local optimality implying global optimality of convex functions, we define $\lambda = 1 - \frac{\varepsilon}{2||y-x||}$ and $z = \lambda x + (1 - \lambda)y$. Prove that $z$ is indeed in the $\varepsilon$-neighbor of $x$.

*Hint*: Calculate the distance between $x$ and $z$, and show it's less than $\varepsilon$.

**Answer:** (The answer essentially follows the lecture.) Note that
$$z - x = (\lambda x + (1 - \lambda)y) - x = (1 - \lambda)(x - y).$$
Taking the norm on both sides gives us the desired result:
$$||z - x|| = (1 - \lambda)||x - y|| = (1 - (1 - \frac{\varepsilon}{2||y-x||}))||x - y|| = \frac{\varepsilon}{2||y-x||}||y - x|| = \varepsilon/2.$$

**Problem 2.**
Suppose $f$ is a convex function and we have the gradient $\nabla f(x) = 0$ at the point $x$. Prove that $x$ is a global optimum of the function $f$.

**Answer:** Since $f(\cdot):\mathbb{R}\longrightarrow\mathbb{R}$ is a convex function, from the first-order condition (see the lecture on convexity), we know that
$$\forall y \in \text{dom } f, \qquad f(y) \geq f(x) + [\nabla f(x)]^T (y - x).$$
So if the gradient at $x$ is equal to zero, i.e. $\nabla f(x) = 0$, we get that
$$\forall y \in \text{dom } f, \qquad f(y) \geq f(x),$$
That is, $x$ is a global minimum of $f$.

## Problem 3.

In future lectures, we will use the $l_2$-penalized square error as the training objective, given by

$$J = \frac{1}{2M} \sum_{m=1}^{M} \left( \sum_{i=0}^{d} w_i x_i^{(m)} - t^{(m)} \right)^2 + \lambda \sum_{i=0}^{d} w_i^2$$

where $\lambda$ is a positive constant.

a) Express the loss in the vector form using $\mathbf{X} \in \mathbb{R}^{M \times (d+1)}$, $\mathbf{w} \in \mathbb{R}^{d+1}$ and $t \in \mathbb{R}^{M}$

b) Compute the gradient $\nabla J(\mathbf{w})$ and the Hessian $\nabla\nabla J(\mathbf{w})$ and show that $J$ is convex in $\mathbf{w}$.

c) Derive a closed-form solution to the problem.

*Hints:*

Suppose $S$ is a symmetric matrix.

$$\frac{\partial \mathbf{x}^\top S \mathbf{x}}{\partial \mathbf{x}} = 2S\mathbf{x} \qquad\qquad \frac{\partial A\mathbf{x}}{\partial \mathbf{x}} = A^\top$$
,

References: https://en.wikipedia.org/wiki/Matrix_calculus

Note: When matrix calculus is needed in exams, I will give the precise formulas that should be used for solving the problem.

**Answer:**

a) Let $x^{(m)} \in \mathbb{R}^{d+1}$ be the vector defined as $(x^{(m)})_i = x_i^{(m)}$. Note that

$$\sum_{i=0}^{d} w_i^2 = w^\top w = ||w||_2^2 \qquad \text{and} \qquad \sum_{i=0}^{d} w_i x_i^{(m)} = w^\top x^{(m)}.$$

Then the loss becomes

$$J = \frac{1}{2M} \sum_{m=1}^{M} (w^\top x^{(m)} - t^{(m)})^2 + \lambda ||w||_2^2.$$

Notice that the individual terms in the above summation is the $m$th component of the vector $Xw - t$, and the summation itself is the norm of this vector. We thus obtain our final solution:

$$J = \frac{1}{2M} ||Xw - t||_2^2 + \lambda ||w||_2^2.$$

b) (Traditionally, $\frac{\partial}{\partial w} J$ is assumed to be a row vector with its $i$th element defined as $(\frac{\partial}{\partial w} J)_i = \frac{\partial}{\partial w_i} J$, while the gradient $\nabla J = (\frac{\partial}{\partial w} J)^\top$ is a column vector.)

Using the hint given in the question, it is clear that the $\nabla||w||_2^2 = 2w$, and

$$\nabla||Xw - t||_2^2 = \frac{\partial(Xw-t)}{\partial w} \cdot \frac{\partial}{\partial(Xw-t)}||Xw - t||_2^2 = 2X^\top(Xw - t).$$

Combining these two equations, we obtain the gradient:

$$\nabla J = \frac{1}{M}X^\top(Xw - t) + 2\lambda w.$$

Computing the Hessian is now straightforward.

$$\nabla^2 J = \frac{\partial}{\partial w}\nabla J = \frac{1}{M}X^\top X + 2\lambda I.$$

c) The problem to solve is $\min_w J$. To solve this, we set the gradient of the objective to zero, and then solve for $w$:

$$\nabla J = 0 \quad \Rightarrow \quad (X^\top X + 2M\lambda I)w^* = X^\top t \quad \Rightarrow \quad w^* = (X^\top X + 2M\lambda I)^{-1}X^\top t.$$

Note that since the Hessian is PSD (why?), MSE is a convex function (also see the notes for Lecture 5). This means that the above solution $w^*$ is a global minimum of the MSE problem. (Note that, for all $M > 0$ and $\lambda > 0$, the matrix $X^\top X + 2M\lambda I$ is guaranteed to be non-singular (why?).

**END**