

Module 4: Supplementary Slides

(Optional, for Interested Students)



Obtaining Taylor Formula for Multivariable Functions

$$f(\mathbf{x} + h\mathbf{s}) = f(\mathbf{x}) + h\nabla f(\mathbf{x})^T \mathbf{s} + o(h)$$

$$f(\mathbf{x} + h\mathbf{s}) = f(\mathbf{x}) + h\nabla f(\mathbf{x})^T \mathbf{s} + \frac{1}{2}h^2 \mathbf{s}^T H(\mathbf{x}) \mathbf{s} + o(h^2)$$

As before let us define:

$$w(h) = f(\mathbf{x} + h\mathbf{s}).$$

From the lecture, we know that

$$w'(0) = \nabla f(\mathbf{x})^T \mathbf{s}$$

and

$$w''(0) = \mathbf{s}^T H(\mathbf{x}) \mathbf{s}.$$

Now, expanding the Taylor formula for the single variable function $w(h)$ around $h = 0$, we have:

$$f(\mathbf{x} + h\mathbf{s}) = w(h) = w(0) + hw'(0) + o(h) = f(\mathbf{x}) + h\nabla f(\mathbf{x})^T \mathbf{s} + o(h)$$

$$f(\mathbf{x} + h\mathbf{s}) = w(h) = w(0) + hw'(0) + \frac{1}{2}h^2 w''(0) + o(h^2) = f(\mathbf{x}) + h\nabla f(\mathbf{x})^T \mathbf{s} + \frac{1}{2}h^2 \mathbf{s}^T H(\mathbf{x}) \mathbf{s} + o(h^2)$$



Approximating a Multivariable Function



Approximating a Multivariable Function

$f(\mathbf{x} + \Delta\mathbf{x})$ can be approximated by using the Taylor expansion as:

$$f(\mathbf{x} + \Delta\mathbf{x}) \approx f(\mathbf{x}) + \nabla f(\mathbf{x}) \cdot \Delta\mathbf{x} = f(\mathbf{x}) + \sum_{k=1}^n \frac{\partial f}{\partial x_k} \Delta x_k$$

Example: If $f(x, y) = \ln(x) + \ln(y)$, then approximate $f(1.06, 1.02)$ using differentials.

Solution:

$$f(x + \Delta x, y + \Delta y) \approx f(x, y) + \frac{\partial f}{\partial x} \Delta x + \frac{\partial f}{\partial y} \Delta y$$

$$\ln(x + \Delta x) + \ln(y + \Delta y) \approx [\ln(x) + \ln(y)] + \frac{1}{x} \Delta x + \frac{1}{y} \Delta y$$

$$\ln(1.06) + \ln(1.02) \approx [\ln(1) + \ln(1)] + \frac{1}{1}(0.06) + \frac{1}{1}(0.02) = 0.08$$



Matrix Calculus



Matrix Calculus

Two Most Useful Results

Proposition 1: Let $\mathbf{x}, \mathbf{c} \in \mathbb{R}^{n \times 1}$. Then

$$\frac{\partial}{\partial \mathbf{x}} (\mathbf{c}^T \mathbf{x}) = \frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^T \mathbf{c}) = \mathbf{c}$$

Proof: To obtain $\nabla_{\mathbf{x}} \mathbf{c}^T \mathbf{x}$, we first expand $\mathbf{c}^T \mathbf{x}$:

$$\mathbf{c}^T \mathbf{x} = [c_1 x_1 + c_2 x_2 + \cdots + c_n x_n]$$

The partial derivative with respect to a single coordinate is:

$$\frac{\partial}{\partial x_i} \mathbf{c}^T \mathbf{x} = c_i$$

Thus, the gradient is:

$$\nabla_{\mathbf{x}} \mathbf{c}^T \mathbf{x} = \nabla_{\mathbf{x}} \mathbf{x}^T \mathbf{c} = \mathbf{c}$$



Matrix Calculus

Two Most Useful Results

Proposition 2: Let $\mathbf{x} \in \mathbb{R}^{n \times 1}$ and $\mathbf{A} \in \mathbb{R}^{n \times n}$. Then

$$\frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = (\mathbf{A} + \mathbf{A}^T) \mathbf{x}$$

Proof: To obtain $\nabla_{\mathbf{x}} \mathbf{x}^T \mathbf{A} \mathbf{x}$ for a square matrix \mathbf{A} , we first expand $\mathbf{x}^T \mathbf{A} \mathbf{x}$:

$$\mathbf{x}^T \mathbf{A} \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}^T \begin{bmatrix} a_{11} & \cdots & a_{1n} \\ a_{21} & \cdots & a_{2n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \cdots & a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}^T \begin{bmatrix} x_1 a_{11} + x_2 a_{12} + \cdots + x_n a_{1n} \\ x_1 a_{21} + x_2 a_{22} + \cdots + x_n a_{2n} \\ \vdots \\ x_1 a_{n1} + x_2 a_{n2} + \cdots + x_n a_{nn} \end{bmatrix} = \begin{bmatrix} x_1^2 a_{11} + x_2^2 a_{12} + \cdots + x_n^2 a_{1n} \\ x_1^2 a_{21} + x_2^2 a_{22} + \cdots + x_n^2 a_{2n} \\ \vdots \\ x_1^2 a_{n1} + x_2^2 a_{n2} + \cdots + x_n^2 a_{nn} \end{bmatrix}$$

The partial derivative with respect to the i th component is:

$$\frac{\partial}{\partial x_i} \mathbf{x}^T \mathbf{A} \mathbf{x} = \sum_{j=1}^n x_j (a_{ij} + a_{ji})$$



Matrix Calculus

Proof (continued): Thus, the gradient is:

$$\nabla_x x^T A x = \begin{bmatrix} \sum_{j=1}^n x_j (a_{1j} + a_{j1}) \\ \sum_{j=1}^n x_j (a_{2j} + a_{j2}) \\ \vdots \\ \sum_{j=1}^n x_j (a_{nj} + a_{jn}) \end{bmatrix} = \begin{bmatrix} a_{11} + a_{11} & a_{21} + a_{12} & \cdots & a_{1n} + a_{n1} \\ a_{12} + a_{21} & a_{22} + a_{22} & \cdots & a_{n2} + a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{1n} + a_{n1} & a_{2n} + a_{n2} & \cdots & a_{nn} + a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = (A + A^T)x$$



Matrix Calculus

Other Helpful and General Results

Let $a \in \mathbb{R}$, $\mathbf{x}, \mathbf{c} \in \mathbb{R}^n$, $\mathbf{U}(\mathbf{x}), \mathbf{V}(\mathbf{x}), \mathbf{A} \in \mathbb{R}^{n \times n}$. Then

1. $\frac{\partial}{\partial \mathbf{x}} (\mathbf{c}^T \mathbf{x}) = \frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^T \mathbf{c}) = \mathbf{c}$ (Proposition 1)

2. $\frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = (\mathbf{A} + \mathbf{A}^T) \mathbf{x}$ (Proposition 2)

3. $\frac{\partial \mathbf{x}^T \mathbf{x}}{\partial \mathbf{x}} = 2\mathbf{x}$ (Special Case of Proposition 2)

4. $\frac{\partial^2 \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{x} \partial \mathbf{x}^T} = \mathbf{A} + \mathbf{A}^T$ (Resulting from using Prop 1 on Prop 2)

5. $\frac{\partial \mathbf{c}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = \mathbf{A}^T \mathbf{c}$ (Another representation of Proposition 1)

6. $\frac{\partial a}{\partial \mathbf{x}} = 0$ (Trivial result)

7. $\frac{\partial}{\partial \mathbf{x}} (a \mathbf{U}(\mathbf{x})) = a \frac{\partial \mathbf{U}(\mathbf{x})}{\partial \mathbf{x}}$ (Trivial result)

8. $\frac{\partial}{\partial \mathbf{x}} (\mathbf{U}(\mathbf{x}) + \mathbf{V}(\mathbf{x})) = \frac{\partial \mathbf{U}(\mathbf{x})}{\partial \mathbf{x}} + \frac{\partial \mathbf{V}(\mathbf{x})}{\partial \mathbf{x}}$

9. $\frac{\partial}{\partial \mathbf{x}} (\mathbf{U}(\mathbf{x})^T \mathbf{V}(\mathbf{x})) = \frac{\partial \mathbf{U}(\mathbf{x})}{\partial \mathbf{x}}^T \mathbf{V}(\mathbf{x}) + \frac{\partial \mathbf{V}(\mathbf{x})}{\partial \mathbf{x}}^T \mathbf{U}(\mathbf{x})$ (Product rule)

10. $\frac{\partial}{\partial \mathbf{x}} [\mathbf{U}(\mathbf{x})^T \mathbf{A} \mathbf{V}(\mathbf{x})] = \frac{\partial \mathbf{U}(\mathbf{x})}{\partial \mathbf{x}}^T \mathbf{A} \mathbf{V}(\mathbf{x}) + \frac{\partial \mathbf{V}(\mathbf{x})}{\partial \mathbf{x}}^T \mathbf{A}^T \mathbf{U}(\mathbf{x})$ (Product rule)

11. $\frac{\partial}{\partial \mathbf{x}} (\mathbf{c}^T \mathbf{U}(\mathbf{x})) = \frac{\partial \mathbf{U}(\mathbf{x})}{\partial \mathbf{x}}^T \mathbf{c}$ (Generalization to Proposition 1)

12. $\frac{\partial}{\partial \mathbf{x}} (\mathbf{c}^T \mathbf{x}^T \mathbf{c} \mathbf{x}) = \frac{\partial}{\partial \mathbf{x}} ((\mathbf{x} \mathbf{c})^T \mathbf{c} \mathbf{x}) = 2 \mathbf{c} \mathbf{c}^T \mathbf{x}$



Application of Matrix Calculus: Least Square Linear Regression



Linear Regression

In a multiple regression, we examine the linear relationship between one dependent variable and several independent variables. Let m observations of one dependent variable be denoted with $\mathbf{y} = [y_i]_{m \times 1}$ and the corresponding m observations of independent variables be denoted by $\mathbf{X} = [\mathbf{1} \ x_1 \ x_2 \ \dots \ x_k] = [x_{ij}]_{m \times (k+1)}$, $\mathbf{x}_j \in \mathbb{R}^m$, $\mathbf{1} = [1]_{m \times 1}$ where x_{ij} is the i^{th} observation of the independent variable \mathbf{x}_j . We wish to estimate the value of \mathbf{Y} by considering a linear function in the form of $\hat{y}_i = \beta_0(1) + \beta_1 x_{i1} + \dots + \beta_k x_{ik}$ where \hat{y}_i is the predicted value for observation i . In a concise form, we can write $\hat{\mathbf{y}} = \mathbf{X}\boldsymbol{\beta} \equiv [\hat{y}_i]_{m \times 1}$, and

$$\hat{\mathbf{y}}_{m \times 1} = \mathbf{X}\boldsymbol{\beta}, \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix}_{(k+1) \times 1}, \quad \mathbf{X} = \begin{bmatrix} 1 & x_{11} & \cdots & x_{1k} \\ 1 & x_{21} & \cdots & x_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{m1} & \cdots & x_{mk} \end{bmatrix}_{m \times (k+1)}$$

The value $\mathbf{h}_i = y_i - \hat{y}_i$ is prediction error corresponding to the i^{th} observations $[1 \ x_{i1} \ x_{i2} \ \dots \ x_{ik}]$ and y_i .



Linear Regression

Ordinary Least Squares Regression (OLS): The OLS regression line is the linear function $\hat{\mathbf{y}} = \mathbf{X}\boldsymbol{\beta}$ that minimizes the sum of squared error values, i.e., $\min_{\boldsymbol{\beta} \in \mathbb{R}^k} (\mathbf{e}^T \mathbf{e})$, $\mathbf{e} = \mathbf{y} - \mathbf{X}\boldsymbol{\beta}$. The OLS regression line coefficient vector is given by:

$$\boldsymbol{\beta}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}.$$

Proof:

Expanding $\mathbf{e}^T \mathbf{e}$ gives:

$$\mathbf{e}^T \mathbf{e} = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{y}^T - \boldsymbol{\beta}^T \mathbf{X}^T) (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X}\boldsymbol{\beta} - \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{y} + \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X}\boldsymbol{\beta}.$$

From [Proposition 1](#), we note that

$$\frac{\partial}{\partial \boldsymbol{\beta}} (\mathbf{y}^T \mathbf{X}\boldsymbol{\beta}) = \frac{\partial}{\partial \boldsymbol{\beta}} (\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{y}) = \mathbf{X}^T \mathbf{y}.$$

Also, by [Proposition 2](#),

$$\frac{\partial}{\partial \boldsymbol{\beta}} (\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X}\boldsymbol{\beta}) = \mathbf{X}^T \mathbf{X} + (\mathbf{X}^T \mathbf{X})^T.$$



Linear Regression

Using the [first derivative test](#) to find the minimum values, using matrix calculus properties, we have:

$$\frac{\partial}{\partial \boldsymbol{\beta}} \left((\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right) = 0 - \mathbf{X}^T \mathbf{y} - \mathbf{X}^T \mathbf{y} + \left(\mathbf{X}^T \mathbf{X} + (\mathbf{X}^T \mathbf{X})^T \right) \boldsymbol{\beta} = -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \boldsymbol{\beta} = -2\mathbf{X}^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}).$$

Solving for the optimal coefficients vector gives:

$$-2\mathbf{X}^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^*) = \mathbf{0} \Rightarrow \mathbf{X}^T \mathbf{X} \boldsymbol{\beta}^* = \mathbf{X}^T \mathbf{y} \Rightarrow \boldsymbol{\beta}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}.$$

Now, checking the [second derivative test](#) to verify that $\boldsymbol{\beta}^*$ is a minimizer, we have

$$\frac{\partial}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T} \left((\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right) = \frac{\partial}{\partial \boldsymbol{\beta}} (-2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \boldsymbol{\beta}) = 2\mathbf{X}^T \mathbf{X},$$

in which we used [Proposition 1](#), to calculate $\frac{\partial}{\partial \boldsymbol{\beta}} (\mathbf{X}^T \mathbf{X} \boldsymbol{\beta}) = (\mathbf{X}^T \mathbf{X})^T = \mathbf{X}^T \mathbf{X}$. Also, any constant comes out of derivative.

Note that $\mathbf{X}^T \mathbf{X}$ is [positive semidefinite](#) because $\forall \mathbf{w} \in \mathbb{R}^m$:

$$\mathbf{w}^T (\mathbf{X}^T \mathbf{X}) \mathbf{w} = \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} = (\mathbf{X} \mathbf{w})^T (\mathbf{X} \mathbf{w}) = \|\mathbf{X} \mathbf{w}\|^2 \geq 0.$$

Hence, $\boldsymbol{\beta}^*$ is a minimizer.



Linear Regression

Example

Find the linear regression line for the points $(1,1), (2,2), (3,4), (4,4), (5,6)$.

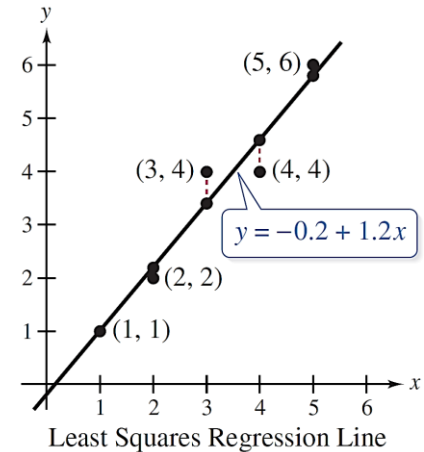
Solution.

The matrices \mathbf{X} and \mathbf{y} are $\mathbf{X} = \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \\ 1 & 5 \end{bmatrix}, \mathbf{y} = \begin{bmatrix} 1 \\ 2 \\ 4 \\ 4 \\ 6 \end{bmatrix}$

$$\mathbf{X}^T \mathbf{X} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 & 5 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \\ 1 & 5 \end{bmatrix} = \begin{bmatrix} 5 & 15 \\ 15 & 55 \end{bmatrix}, \quad \mathbf{X}^T \mathbf{y} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 & 5 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 4 \\ 4 \\ 6 \end{bmatrix} = \begin{bmatrix} 17 \\ 63 \end{bmatrix}$$

$$\boldsymbol{\beta}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \frac{1}{50} \begin{bmatrix} 55 & -15 \\ -15 & 5 \end{bmatrix} \begin{bmatrix} 17 \\ 63 \end{bmatrix} = \begin{bmatrix} -0.2 \\ 1.2 \end{bmatrix}$$

Hence, the least squares regression line is $y = -0.2 + 1.2x$.



Constrained Optimization

KKT Conditions for Multiple Constraints



Equality Constraint

Proof for One Constraint Using Gradients

Consider an optimization problem with a single equality constraint

$$\begin{aligned} \min_x & f(x) \\ \text{subject to } & g(x) = 0 \end{aligned}$$

where f and g have continuous partial derivatives. Then $f(x)$ has an interior local minimum at x^* if, for any feasible direction s ,

$$\nabla f(x^*) \cdot s = 0 \Rightarrow \nabla f(x^*) \perp s$$

Feasible directions are those that satisfy:

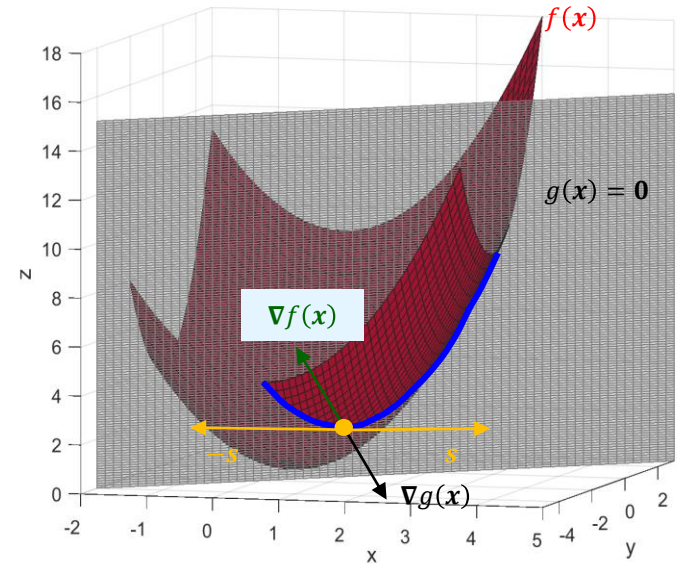
$$g(x + hs) = g(x) = 0 \Rightarrow Dg_s(x) = 0 \Rightarrow \nabla g(x) \cdot s = 0 \Rightarrow \nabla g(x) \perp s$$

So, at the optimal point x^* :

$$\nabla f(x^*) \parallel \nabla g(x^*) \Rightarrow \nabla f(x^*) = \lambda \nabla g(x^*)$$

λ is called a **Lagrange multiplier** and $L(x, \lambda) = f(x) - \lambda g(x)$ is called **Lagrangian function**. At x^* :

$$\nabla L(x^*) = \nabla f(x^*) - \lambda \nabla g(x^*) = 0$$



Inequality Constraints

Consider an optimization problem with a single inequality constraint

$$\begin{aligned} \min_{\mathbf{x}} f(\mathbf{x}) \\ \text{subject to } g(\mathbf{x}) \leq 0 \end{aligned}$$

We know that if the solution lies at the constraint boundary, the Lagrange multiplier holds: $\nabla f(\mathbf{x}^*) - \lambda \nabla g(\mathbf{x}^*) = \mathbf{0}$, otherwise, we have $\nabla f(\mathbf{x}^*) = \mathbf{0}$. We could optimize the problem by introducing an infinite step penalty for infeasible points:

$$f_{\infty\text{-step}}(\mathbf{x}) = \begin{cases} f(\mathbf{x}) & \text{if } g(\mathbf{x}) \leq 0 \\ \infty & \text{if } g(\mathbf{x}) > 0 \end{cases} \Rightarrow f_{\infty\text{-step}}(\mathbf{x}) = f(\mathbf{x}) + \infty(g(\mathbf{x}) > 0)$$

Unfortunately, $f_{\infty\text{-step}}(\mathbf{x})$ is inconvenient to optimize because it is discontinuous and nondifferentiable. We can instead use a linear penalty $\lambda g(\mathbf{x})$, which forms a lower bound on $\infty(g(\mathbf{x}) > 0)$ and penalizes the objective as long as $\lambda > 0$. We can use this linear penalty to construct a Lagrangian function:

$$L(\mathbf{x}, \mu) = f(\mathbf{x}) - \lambda g(\mathbf{x})$$

We can recover $f_{\infty\text{-step}}(\mathbf{x})$ by maximizing with respect to μ .

$$f_{\infty\text{-step}}(\mathbf{x}) = \max_{\mu \geq 0} L(\mathbf{x}, \mu)$$

For any infeasible \mathbf{x} , we get infinity and for any feasible \mathbf{x} , we get $f(\mathbf{x})$.



Inequality Constraints

KKT Conditions

The new optimization problem is thus

$$\min_x \max_{\lambda \geq 0} L(\mathbf{x}, \lambda)$$

This formulation is known as the **primal problem**. Optimizing the primal problem will require finding \mathbf{x}^* such that:

1. **(Prime) Feasibility: $g(\mathbf{x}^*) \leq 0$**
 - the point is feasible
2. **Dual Feasibility: $\lambda \geq 0$**
 - The penalty must point in the right direction. μ is also called a dual variable.
3. **Complementary Slackness: $\lambda g(\mathbf{x}^*) = 0$**
 - A feasible point on the boundary will have $g(\mathbf{x}^*) = 0$ whereas a feasible point with $g(\mathbf{x}^*) < 0$ must have $\lambda = 0$ to recover $f(\mathbf{x}^*)$ from Lagrangian
4. **Stationarity: $\nabla f(\mathbf{x}^*) - \lambda \nabla g(\mathbf{x}^*) = \mathbf{0}$**
 - When the constraint is active, we require the Lagrange multiplier. When it is inactive, we require $\nabla f(\mathbf{x}^*) = \mathbf{0}$ and $\lambda = 0$.

These conditions are known as **Karush-Kuhn-Tucker (KKT) conditions**. (Mostly used for verification of optimality of a solution. It is hard to use it directly to find an optimal value)



Constrained Optimization

KKT Conditions For Multiple Constraints

KKT conditions: Consider the following problem:

$$\max_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$$

subject to :

$$c_k(\mathbf{x}) \leq 0, \quad k = 1, \dots, K$$

The first derivative necessary conditions for $\mathbf{x}^* \in \mathbb{R}^m$ to be a local maximizer is that $\exists \boldsymbol{\lambda}^* \in \mathbb{R}^K$ such that :

1. (Primal) feasibility: $c_k(\mathbf{x}) \leq 0, \quad k = 1, \dots, K$
2. Stationarity: $\frac{\partial f(\mathbf{x})}{\partial x_i} - \sum_{k=1}^K \lambda_k \frac{\partial c_k(\mathbf{x})}{\partial x_i} = 0, \quad i = 1, \dots, n$
3. Complementary Slackness: $\lambda_k^* c_k(\mathbf{x}) = 0, \quad k = 1, \dots, K,$
4. Positivity (/Dual feasibility): $\lambda_k^* \geq 0, \quad k = 1, \dots, K$



KKT Conditions For Multiple Constraints

Example

Solve the following problem:

$$\max_{x_1, x_2} f(x_1, x_2) = (x_1 + 1)^2 + (x_2 + 1)^2$$

subject to:

$$\lambda_1 \quad x_1^2 + x_2^2 \leq 2,$$

$$\lambda_2 \quad x_2 \leq 1$$

Solution: The KKT conditions are

- 1) $2(x_1^* + 1) - \lambda_1(2x_1^*) = 0 \Rightarrow x_1^* + 1 = \lambda_1^* x_1^*$
- 2) $2(x_2^* + 1) - \lambda_1(2x_2^*) - \lambda_2(1) = 0 \Rightarrow x_2^* + 1 = \lambda_1^* x_2^* + \frac{1}{2}\lambda_2^*$
- 3) $\lambda_1^*(x_1^{*2} + x_2^{*2} - 2) = 0$
- 4) $\lambda_2^*(x_2^* - 1) = 0$
- 5) $\lambda_1^* \geq 0, \lambda_2^* \geq 0$



KKT Conditions For Multiple Constraints

Solution

We can examine the possibilities for the complementarity condition:

- $\lambda_1^* = 0, \lambda_2^* = 0$

$$x_1^* + 1 = 0, x_2^* + 1 = 0 \Rightarrow (x_1^*, x_2^*) = (-1, -1) \Rightarrow f(x_1^*, x_2^*) = 0 \quad \text{max}$$

- $\lambda_1^* = 0, \lambda_2^* > 0$

$$x_2 + 1 = \frac{1}{2}\lambda_2^*, \lambda_2^*(x_2^* - 1) = 0 \Rightarrow (x_1^*, x_2^*) = (-1, 1), \lambda_2^* = 4 \Rightarrow f(x_1^*, x_2^*) = -4 \quad \text{Saddle point}$$

- $\lambda_1^* > 0, \lambda_2^* = 0$

$$x_1^* + 1 = \lambda_1^* x_1^*, x_2^* + 1 = \lambda_1^* x_2^*, \lambda_1^*(x_1^{*2} + x_2^{*2} - 2) = 0 \Rightarrow (x_1^*, x_2^*) = (1, 1), \lambda_1^* = 2 \Rightarrow f(x_1^*, x_2^*) = -8 \quad \text{min}$$

- $\lambda_1^* > 0, \lambda_2^* > 0$

$$x_1^{*2} + x_2^{*2} = 2, x_2^* = 1 \Rightarrow (x_1^*, x_2^*) = (\pm 1, -1) \quad \text{But the answers coincide the previous items.}$$



KKT Conditions For Multiple Equality and Inequality Constraints

Consider the following problem:

$$\max_{x \in \mathbb{R}^n} f(x)$$

Subject to :

$$c_k(x) \leq 0, \quad k = 1, \dots, K$$

$$d_l(x) = 0, \quad l = 1, \dots, L$$

The first derivative necessary conditions for $x^* \in \mathbb{R}^n$ to be a local maximizer is that $\exists \lambda^* \in \mathbb{R}^K$ and $\exists \gamma^* \in \mathbb{R}^L$ such that :

1. (Primal) feasibility: $c_k(x) \leq 0, \quad k = 1, \dots, K$, and $d_l(x) = 0, \quad l = 1, \dots, L$
2. Stationarity: $\frac{\partial f(x)}{\partial x_i} - \sum_{k=1}^K \lambda_k \frac{\partial c_k(x)}{\partial x_i} - \sum_{l=1}^L \gamma_l \frac{\partial d_l(x)}{\partial x_i} = 0, \quad i = 1, \dots, n$
3. Complementary Slackness: $\lambda_k^* c_k(x) = 0, \quad k = 1, \dots, K$, and $\gamma_l^* d_l(x) = 0, \quad l = 1, \dots, L$
4. Positivity (/Dual feasibility): $\lambda_k^* \geq 0, \quad k = 1, \dots, K$ (only for inequality constraints)



Additional Examples in R



Plotting 3D in R

`persp (x,y,z)`

Creates a 3D surface graphic for the vectors (x,y,z) where z is a function of x and y defined previously by the "outer" code.

Example: Plot $f(x,y) = x^3 + y^3 - xy$

```
x=y=seq(-1,1,length=20)
```

```
f=function(x,y)x^3+y^3-x*y
```

```
z=outer(x,y,f)
```

```
persp(x,y,z,
```

```
  main= "Perspective Plot", # (optional)
```

```
  zlab = "Height", # (optional)
```

```
  theta = -30,      # Rotation (vertical) (optional)
```

```
  phi = 0,         # Rotation (horizontal) (optional)
```

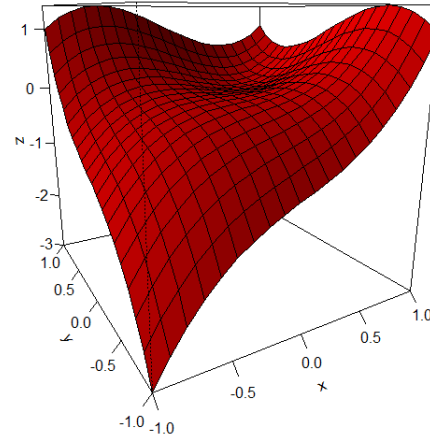
```
  expand = 1,      # Shrinking/growing factor (optional)
```

```
  col = "red",    # Shrinking/growing factor (optional)
```

```
  shade = 0.6,    # color shade (optional)
```

```
  ticktype = "detailed" # puts numbers on x,y, and z axes
```

```
)
```



Constrained Optimization

Example: Equality Constraint

Minimize $c(q_1, q_2) = 2q_1^2 + q_1q_2 + q_2^2 + 200$ subject to $q_1 + q_2 = 200$.

Solution:

```
f=function(x) (2*x[1]^2+x[1]*x[2]+x[2]^2+200)
equalities=function(x) {
    h=0
    h[1]=200-x[1]-x[2]
    return(h) }
p0=c(0,0)
y=constrOptim.nl(p0,f,heq = equalities);
print(y$par)
[1] 49.9983 150.0017
print(y$value)
[1] 35200
```



Constrained Optimization

Example: Equality Constraint

Maximize $f(x_1, x_2) = 4x_1^2 + 10x_2^2$ subject to $x_1^2 + x_2^2 = 4$.

Solution: Write the function as minimizing $-f(x_1, x_2)$.

```
f=function(x) -(4*x[1]^2+10*x[2]^2)
Equalities=function(x) {
    h=0
    h[1]=(x[1]^2+x[2]^2-4)
    return(h) }
p0=c(2,0)
y=constrOptim.nl(p0,f, heq=Equalities);
print(y$par)
[1] -4.99786e-05 -2.00000e+00  $(x_1^*, x_2^*) = (0, -2)$ 
print(y$value)
[1] -40  $f(x_1^*, x_2^*) = 40$ 
```



Constrained Optimization

Example: Inequality Constraint

Maximize $f(x_1, x_2) = 4x_1^2 + 10x_2^2$ subject to $x_1^2 + x_2^2 \leq 4$.

Solution: Write the function as minimizing $-f$.

```
f=function(x) -(4*x[1]^2+10*x[2]^2)
inequalities=function(x){
    h=0
    h[1]= 4-x[1]^2-x[2]^2 #note that h[1] must be positive, i.e., h[1]>=0
    return(h) }
p0=c(0,0)
y=constrOptim.nl(p0,f,hin=inequalities);
par: 0.02919113 1.999787
fval: -39.99489
```

