

Optimization: Part 2

Abhinav Anand, IIMB

2018/06/21

Background

The problems seeking to maximize profits or minimize costs often feature nontrivial constraints which the optimal needs to satisfy. The solution must lie at the intersection of constraints (for equality constraints) or on one side of the constraint surface (for inequality constraints).

The problem in a general form is:

$$\max f(x) : x \in \mathbb{R}^n, x \geq 0$$

$$g_1(x_1, \dots, x_n) \leq b_1, \dots, g_k(x_1, \dots, x_n) \leq b_k$$

$$h_1(x_1, \dots, x_n) = c_1, \dots, h_m(x_1, \dots, x_n) = c_m$$

The objective function f is real valued, i.e., $f : \mathbb{R}^n \rightarrow \mathbb{R}$; $g(\cdot)$ are functional forms of the *inequality* constraints while $h(\cdot)$ are functional forms for the *equality* constraints.

Equality Constraints

Consider the case when $x = (x_1, x_2)$ and there is a single equality constraint $h_1(x) = c_1$.

$$\max f(x) : x \in \mathbb{R}^2, x \geq 0$$

$$h_1(x) = c_1$$

To make the illustration more concrete, consider $f(x) = x_1x_2$ and $h_1(x) = c_1 : x_1 + x_2 = 5$.

```

x_1 <- seq(0.1, 10, 0.05)
x_2 <- seq(0.1, 10, 0.05)

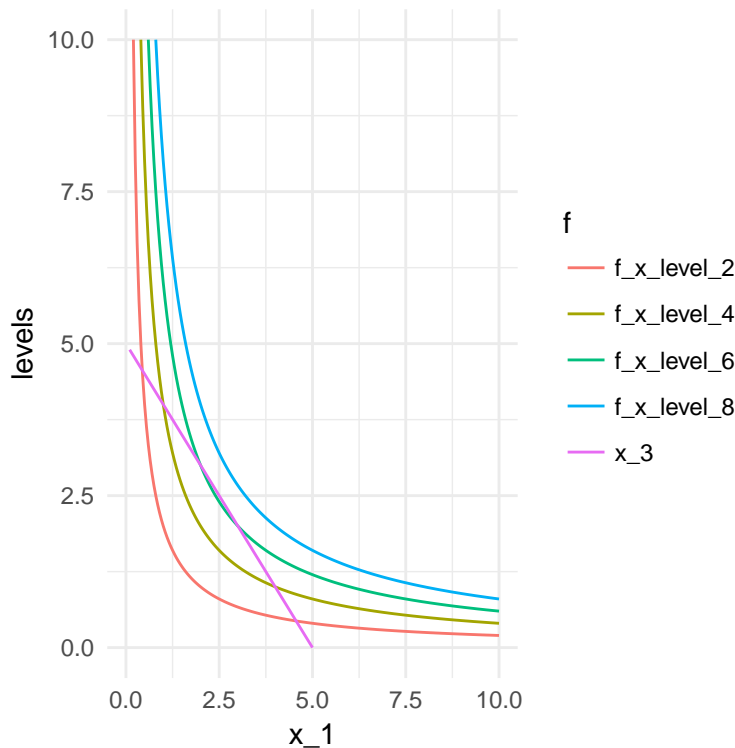
f_x_level_2 <- 2/x_2 #level sets
f_x_level_4 <- 4/x_2
f_x_level_6 <- 6/x_2
f_x_level_8 <- 8/x_2

x_3 <- 5 - x_2 #constraint set

data_obj_1 <- cbind(x_1,
                    f_x_level_2,
                    f_x_level_4,
                    f_x_level_6,
                    f_x_level_8,
                    x_3
                    ) %>%
  dplyr::as_tibble() %>%
  tidyr::gather(.,
                f_x_level_2:x_3,
                key = "f",
                value = "levels")

ggplot(data_obj_1, aes(x_1, levels, color = f)) +
  geom_line() +
  scale_y_continuous(limits = c(0, 10)) +
  scale_x_continuous(limits = c(0, 10)) +
  theme_minimal()

```



Geometrically we need to find the highest valued level set for $f(x) = x_1x_2$ that satisfies $x_1 + x_2 = 5, x_1, x_2 \geq 0$. The key observation is the following: at the optimal, the levels sets and the constraint set must be tangent—just touching (intersecting) each other at exactly one point. (Why must this be so? What happens if the plots cross over? Can we improve the objective function then?)

The Lagrangian

If the curves are tangents to each other at the optimal point then it must be so that the tangents to the curves at the optimal are in the same direction.

The slope of the level set of f at x^* is:

$$-\frac{\frac{\partial f}{\partial x_1}(x^*)}{\frac{\partial f}{\partial x_2}(x^*)}$$

and that of the equality constraint is:

$$-\frac{\frac{\partial h}{\partial x_1}(x^*)}{\frac{\partial h}{\partial x_2}(x^*)}$$

Since they're equal

$$\frac{\frac{\partial f}{\partial x_1}(x^*)}{\frac{\partial f}{\partial x_2}(x^*)} = \frac{\frac{\partial h}{\partial x_1}(x^*)}{\frac{\partial h}{\partial x_2}(x^*)} = \lambda$$

This can be rearranged as:

$$\frac{\frac{\partial f}{\partial x_1}(x^*)}{\frac{\partial h}{\partial x_1}(x^*)} = \frac{\frac{\partial f}{\partial x_2}(x^*)}{\frac{\partial h}{\partial x_2}(x^*)} = \lambda$$

or,

$$\begin{aligned}\frac{\partial f}{\partial x_1}(x^*) - \lambda \frac{\partial h}{\partial x_1}(x^*) &= 0 \\ \frac{\partial f}{\partial x_2}(x^*) - \lambda \frac{\partial h}{\partial x_2}(x^*) &= 0\end{aligned}$$

There are three unknowns: (x_1^*, x_2^*, λ) . There are two equations above, and there is a third equation—the constraint equation $h(x_1, x_2) = c_1$. Together, we can find x^* and λ .

The following function is formally referred to as the *Lagrangian*, and λ as the Lagrange multiplier:

$$\mathcal{L}(x_1, x_2, \lambda) := f(x_1, x_2) - \lambda \cdot (h(x_1, x_2) - c)$$

We consider the critical points of the Lagrangian: $\frac{\partial \mathcal{L}}{\partial x_1}(x^*), \frac{\partial \mathcal{L}}{\partial x_2}(x^*), \frac{\partial \mathcal{L}}{\partial \lambda}(x^*) = 0$.

Essentially by forming the Lagrangian, we are transforming a constrained optimization program featuring the objective function $f(\cdot)$ into an *unconstrained* optimization program featuring the Lagrangian $\mathcal{L}(\cdot)$. However, there is an extra variable λ , the Lagrange multiplier, that is introduced in the new program.

Constraint Qualification

In order for the slopes to be well-defined, $\frac{\partial h}{\partial x_1}(x^*) \neq 0$ and $\frac{\partial h}{\partial x_2}(x^*) \neq 0$. Since this is a restriction on the constraint set, it's called a *constraint qualification*.

Theorem: For the function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $f \in C^1$, if $x^* \in \mathbb{R}^n$ is a solution of

$$\max f(x) : x \geq 0, h(x) = c_1$$

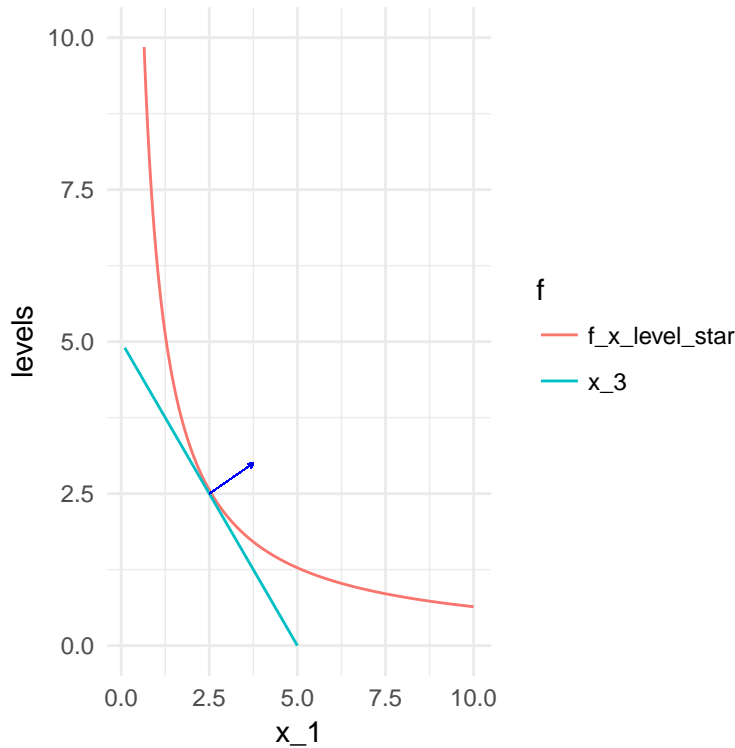
and x^* is *not* a critical point of h ; then, there is a real number λ^* such that (x^*, λ^*) is a critical point of $\mathcal{L} = f(x) - \lambda \cdot (h(x) - c)$.

Level-Set Gradients and Lagrangians

```
f_x_level_star <- 6.4/x_2 #set manually

data_plot_grad <- cbind(x_1,
                        f_x_level_star,
                        x_3
                        ) %>%
  dplyr::as_tibble() %>%
  tidyr::gather(.,
                c(f_x_level_star, x_3),
                key = 'f',
                value = 'levels'
                )

ggplot(data_plot_grad, aes(x_1, levels, color = f)) +
  geom_line() +
  scale_y_continuous(limits = c(0, 10)) +
  geom_segment(aes(x = 2.5,
                  y = 2.5,
                  xend = 3.75, #set manually
                  yend = 3
                  ),
              color = "blue",
              arrow = arrow(length = unit(0.015, "npc")),
              size = 0.2
              ) +
  theme_minimal()
```



The gradient of f and h are respectively $[\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}]$ and $[\frac{\partial h}{\partial x_1}, \frac{\partial h}{\partial x_2}]$. These point to the directions of maximum change and are orthogonal to the level sets of f, h .

Since the level sets and the constraints are tangent at the optimal implies that their respective gradients—orthogonal to the tangent—must again point in the same direction.

$$[\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}] = \lambda [\frac{\partial h}{\partial x_1}, \frac{\partial h}{\partial x_2}]$$

This yields exactly the Lagrangian function for the optimization program.

Several Equality Constraints

Consider the following variation on the maximization problem where there are several equality constraints now:

$$\max f(x) : x \in \mathbb{R}^n, x \geq 0$$

$$C_h = \{h_1(x) = c_1, \dots, h_m(x) = c_m\}$$

The generalization from one to many equality constraints is straightforward.

Constraint Qualification

In the case of one constraint, the qualification is:

$$\left[\frac{\partial h}{\partial x_1}(x^*), \dots, \frac{\partial h}{\partial x_n}(x^*)\right] \neq (0, \dots, 0)$$

Likewise, in the case of m equality constraints: $\{h_1(x) = c_1, \dots, h_m(x) = c_m\}$, their *Jacobian* (first derivative matrix) must be *invertible* at the critical point x^* .

$$Dh(x^*) = \begin{bmatrix} \frac{\partial h_1}{\partial x_1}(x^*), \dots, \frac{\partial h_1}{\partial x_n}(x^*) \\ \frac{\partial h_2}{\partial x_1}(x^*), \dots, \frac{\partial h_2}{\partial x_n}(x^*) \\ \vdots \\ \frac{\partial h_m}{\partial x_1}(x^*), \dots, \frac{\partial h_m}{\partial x_n}(x^*) \end{bmatrix}$$

For the constraint Jacobian at the critical point required to be invertible implies that the matrix above must have *full rank*, further equivalent to the condition that the determinant be non-zero at the critical point. More formally, it is said that (h_1, \dots, h_m) satisfy the *non-degenerate constraint qualification* (NCDQ) at x^* if matrix $Dh(x^*)$ is invertible at x^* (has full rank).

The Geometry of the Lagrangian

When there are m constraints: $h_1(x) = c_1, \dots, h_m(x) = c_m$, the gradient of the objective function $\nabla f(x^*)$ at the optimal must be a linear combination of the gradients of the constraints $\sum_{i=1}^m \lambda_i \nabla h_i(x^*)$.

This is so because ∇h_i gives the directions of maximum increase of h_i . If the constraints $h_i = c_i$ are to be satisfied, we must move in the direction where h_i are constant (and equal to c_i) and hence in directions *orthogonal* (perpendicular) to ∇h_i and hence in directions *orthogonal* to $\sum_{i=1}^m \lambda_i \nabla h_i(x^*)$. In the same way, the contour lines for f are orthogonal to ∇f .

At the optimal, the objective function f must touch (be tangent to) the binding constraints and hence the gradient of f must lie in the *span* of the constraints' gradients: $\nabla f(x^*) = \sum_{i=1}^m \lambda_i \nabla h_i(x^*)$.

Theorem: Given $f, h_1, \dots, h_m : \mathbb{R}^n \rightarrow \mathbb{R} \in C^1$, where f is the objective function and $h_i = c_i$ are equality constraints which x must satisfy; and that h_i satisfy the NDCQ condition; then there are $\lambda^* = (\lambda_1^*, \dots, \lambda_m^*)$ such that (x^*, λ^*) is the critical point of the Lagrangian $\mathcal{L}(x, \lambda)$.

Hence

$$\mathcal{L}(x, \lambda) := f(x) - \lambda_1 \cdot (h_1(x) - c_1) + \dots + \lambda_m \cdot (h_m(x) - c_m)$$

and

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial x_1}(x^*, \lambda^*) &= 0, \dots, \frac{\partial \mathcal{L}}{\partial x_n}(x^*, \lambda^*) = 0 \\ \frac{\partial \mathcal{L}}{\partial \lambda_1}(x^*, \lambda^*) &= 0, \dots, \frac{\partial \mathcal{L}}{\partial \lambda_m}(x^*, \lambda^*) = 0 \end{aligned}$$

In total we get $n + m$ equations with $n + m$ unknowns: $(x_1, \dots, x_n, \lambda_1, \dots, \lambda_m)$.

Inequality Constraints

Consider now a maximization program in \mathbb{R}^2 with one inequality constraint:

$$\max f(x_1, x_2), (x_1, x_2) \geq 0, g(x_1, x_2) \leq b$$

Binding Constraint

To make the discussion more concrete, reconsider the maximization of the objective function $f = x_1 x_2$ with an inequality constraint $g(x) = x_1^2 + x_2^2 \leq 1$.

We can see that for such a maximization program, the optimal must occur at the boundary of the constraint set $g(x) = b$. The level sets of f, g are tangent to each other, or equivalently their gradients are in the same direction:

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

Additionally, however we note that $\lambda > 0$ i.e., the gradients point in the same direction and *not* in the opposite direction as is seen below.

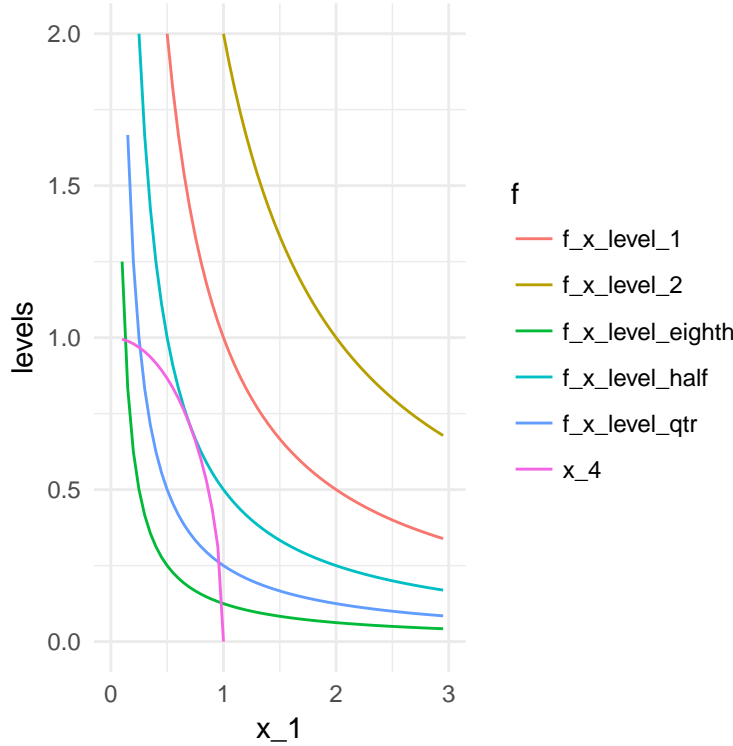
```
x_4 <- sqrt(1 - x_2^2) #inequality constraint

f_x_level_1 <- 1/x_2
f_x_level_half <- 1/(2*x_2)
f_x_level_qtr <- 1/(4*x_2)
f_x_level_eighth <- 1/(8*x_2)

data_plot_ineq <- cbind(x_1,
                        f_x_level_half,
                        f_x_level_1,
                        f_x_level_2,
                        f_x_level_eighth,
                        f_x_level_qtr,
                        x_4
                        ) %>%

dplyr::as_tibble() %>%
tidyr::gather(.,
              f_x_level_half:x_4,
              key = 'f',
              value = 'levels'
              )

ggplot(data_plot_ineq, aes(x_1, levels, color = f)) +
  geom_line() +
  scale_y_continuous(limits = c(0, 2)) +
  scale_x_continuous(limits = c(0, 3)) +
  theme_minimal()
```



As before, the constraint qualification condition holds—the optimal should not be a critical point of the inequality constraint.

Non-binding Constraint

When constraints are of the form $g(x) \leq b$, the optimal may as well lie strictly in the interior depending upon the objective function contour lines. In such cases, we can still form the Lagrangian provided we enforce the condition that $\lambda = 0$, which when combined with the slackness of the inequality: $g(x) - b < 0$ gives complementary slackness conditions: $\lambda \cdot (g(x) - b) = 0$.

Since we do not know which constraint is non-binding at the optimal, we cannot set $\frac{\partial \mathcal{L}}{\partial \lambda} = 0$ (which is equivalent to $g(x) - b = 0$). Hence such a condition is replaced by the complementary slackness condition, which implies that either the constraint is binding, in which case $\lambda > 0$ or the constraint is slack in which case $\lambda = 0$.

Theorem: Given functions $f, g : \mathbb{R}^2 \rightarrow \mathbb{R} \in C^1$ and that $g(x_1, x_2) \leq b$ (additionally, if $g(x^*) = b$, then $\nabla g(x^*) \neq 0$); then (x^*, λ^*) is the critical point of the Lagrangian

and the following are true:

1. $\nabla_{x^*, \lambda^*} \mathcal{L} = 0$
2. $\lambda^* \cdot (g(x^*) - b) = 0$
3. $\lambda^* \geq 0$

Remarks

1. $\nabla_x \mathcal{L} = 0$ for both equality and inequality constraints.
2. $\nabla_\lambda \mathcal{L} = 0$ is true *only* for equality constraints.
3. Constraint qualification needs to be checked *only* for binding inequality constraints.
4. $\lambda \geq 0$ *only* for inequality constraints.
5. Complementary slackness must hold for inequality constraints.

Several Inequality Constraints

The generalization from one inequality constraint to several is fairly straightforward:

Theorem: Given functions $f, g_1, \dots, g_k : \mathbb{R}^n \rightarrow \mathbb{R} \in C^1$ and $x^* \in \mathbb{R}^n$ is a local maximizer on the constraint set

$$C_g = \{g_1(x) \leq b_1, \dots, g_k(x) \leq b_k\}$$

Assume without loss of generality that the first k_0 constraints are binding at x^* and the rest $k - k_0$ are slack. Also assume the nondegeneracy constraint qualification is satisfied by stipulating that the following constraint derivative matrix is invertible at the optimal x^*

$$Dg(x^*) = \begin{bmatrix} \frac{\partial g_1}{\partial x_1}(x^*), \dots, \frac{\partial g_1}{\partial x_n}(x^*) \\ \vdots \\ \frac{\partial g_{k_0}}{\partial x_1}(x^*), \dots, \frac{\partial g_{k_0}}{\partial x_n}(x^*) \end{bmatrix}$$

That is, the constraint Jacobian at the critical point is invertible and has full rank.

Then the Lagrangian

$$\mathcal{L}(x, \lambda) := f(x) - \lambda_1 \cdot (g_1(x) - b_1) - \dots - \lambda_k (g_k(x) - b_k)$$

has the critical points (x^*, λ^*) such that:

1. $\nabla_x \mathcal{L}(x^*) = 0$
2. $\lambda_1^* \cdot (g_1(x) - b_1) = 0, \dots, \lambda_k^* \cdot (g_k(x) - b_k) = 0$
3. $\lambda_1^* \geq 0, \dots, \lambda_k^* \geq 0$

The General Case: Mixed Constraints

The Karush-Kuhn-Tucker (KKT) conditions are first order necessary conditions for a general optimization program including equality as well as inequality constraints. If there are no inequality constraints, the KKT conditions are equivalent to the method of Lagrange multipliers.

We consider the general case:

$$\max f(x), x \in \mathbb{R}^n, x \geq 0 :$$

$$C_g = \{g_1(x) \leq b_1, \dots, g_k(x) \leq b_k\}$$

$$C_h = \{h_1(x) = c_1, \dots, h_m(x) = c_m\}$$

We assume $f, \{g_i\}_{i=1}^k, \{h_j\}_{j=1}^m : \mathbb{R}^n \rightarrow \mathbb{R} \in C^1$.

For the inequality constraint set C_g , we assume without loss of generality that the first k_0 constraints are binding at x^* . We assume that the following derivative matrix is invertible at the optimal (NDCQ) (has full rank):

$$\begin{bmatrix} \frac{\partial g_1}{\partial x_1}(x^*), \dots, \frac{\partial g_1}{\partial x_n}(x^*) \\ \vdots \\ \frac{\partial g_{k_0}}{\partial x_1}(x^*), \dots, \frac{\partial g_{k_0}}{\partial x_n}(x^*) \\ \vdots \\ \frac{\partial h_1}{\partial x_1}(x^*), \dots, \frac{\partial h_1}{\partial x_n}(x^*) \\ \vdots \\ \frac{\partial h_m}{\partial x_1}(x^*), \dots, \frac{\partial h_m}{\partial x_n}(x^*) \end{bmatrix}$$

We form the Lagrangian as before:

$$\mathcal{L}(x, \lambda, \mu) := f(x) - \sum_{i=1}^k \lambda_i \cdot (g_i(x) - b_i) - \sum_{j=1}^m \mu_j \cdot (h_j(x) - c_j)$$

The critical points of the Lagrangian are found from the FOC:

$$\nabla_x \mathcal{L}(x^*, \lambda^*, \mu^*) := \nabla f(x^*) - \sum_{i=1}^k \lambda_i \nabla g_i(x^*) - \sum_{j=1}^m \mu_j \nabla h_j(x^*) = 0$$

The Lagrangian above has multipliers $\lambda^* \in \mathbb{R}^k, \mu^* \in \mathbb{R}^m$ such that

1. $\nabla_x \mathcal{L}(x^*, \lambda^*) = 0$
2. $\lambda_1 \cdot (g_1(x^*) - b_1) = 0, \dots, \lambda_k \cdot (g_k(x^*) - b_k) = 0$
3. $h_1(x^*) - c_1 = 0, \dots, h_m(x^*) - c_m = 0$
4. $\lambda_1^* \geq 0, \dots, \lambda_k^* \geq 0$
5. $g_1(x^*) - b_1 \leq 0, \dots, g_k(x^*) - b_k \leq 0$

Interpreting the Lagrange Multiplier

In plain language, the Lagrange multiplier measures the sensitivity of the optimal value of the objective function to changes in the constraints' right hand sides.

We illustrate this idea using a simple maximization program in two dimensions and with one equality constraint:

$$\max f(x), x \in \mathbb{R}^2, x \geq 0 :$$

$$h(x) = c$$

The central idea is to consider c as a parameter that varies from problem to problem. Given c , $(x_1^*(c), x_2^*(c))$ is the optimal as a function of the constraint parameter c ; $f(x_1^*(c), x_2^*(c))$ is the optimal value as a function of constraint parameter; and $\lambda^*(c)$ is the corresponding multiplier.

Theorem: Given $f, h : \mathbb{R}^2 \rightarrow \mathbb{R} \in C^1$; given $x_1^*(c), x_2^*(c), \lambda^*(c) \in C^1$ and NDCQ condition:

$$\lambda^*(c) = \frac{d}{dc} \cdot [f(x_1^*(c), x_2^*(c))]$$

The same idea can be generalized for the case with several mixed constraints.

Computation

Application: Regularized Regressions