

A1: Introduction to Optimization theory

very short review

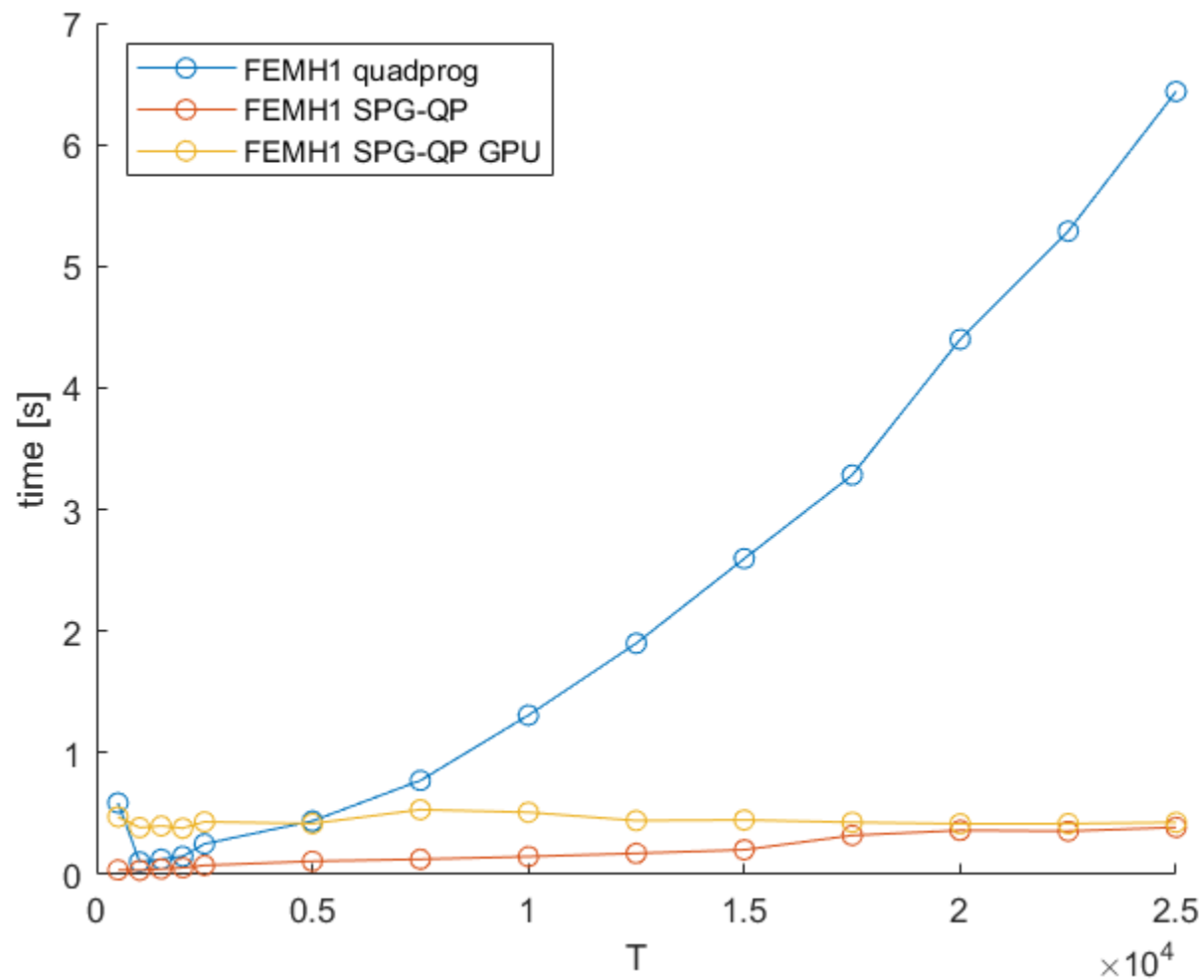
CECAM workshop, Mainz, 2019

Algorithm matters

solve QP problem

$$\Gamma^* = \arg \min_{\gamma} \gamma^T H \gamma + g^T \gamma$$

subject to $B\gamma = c$ and $\gamma \geq 0$

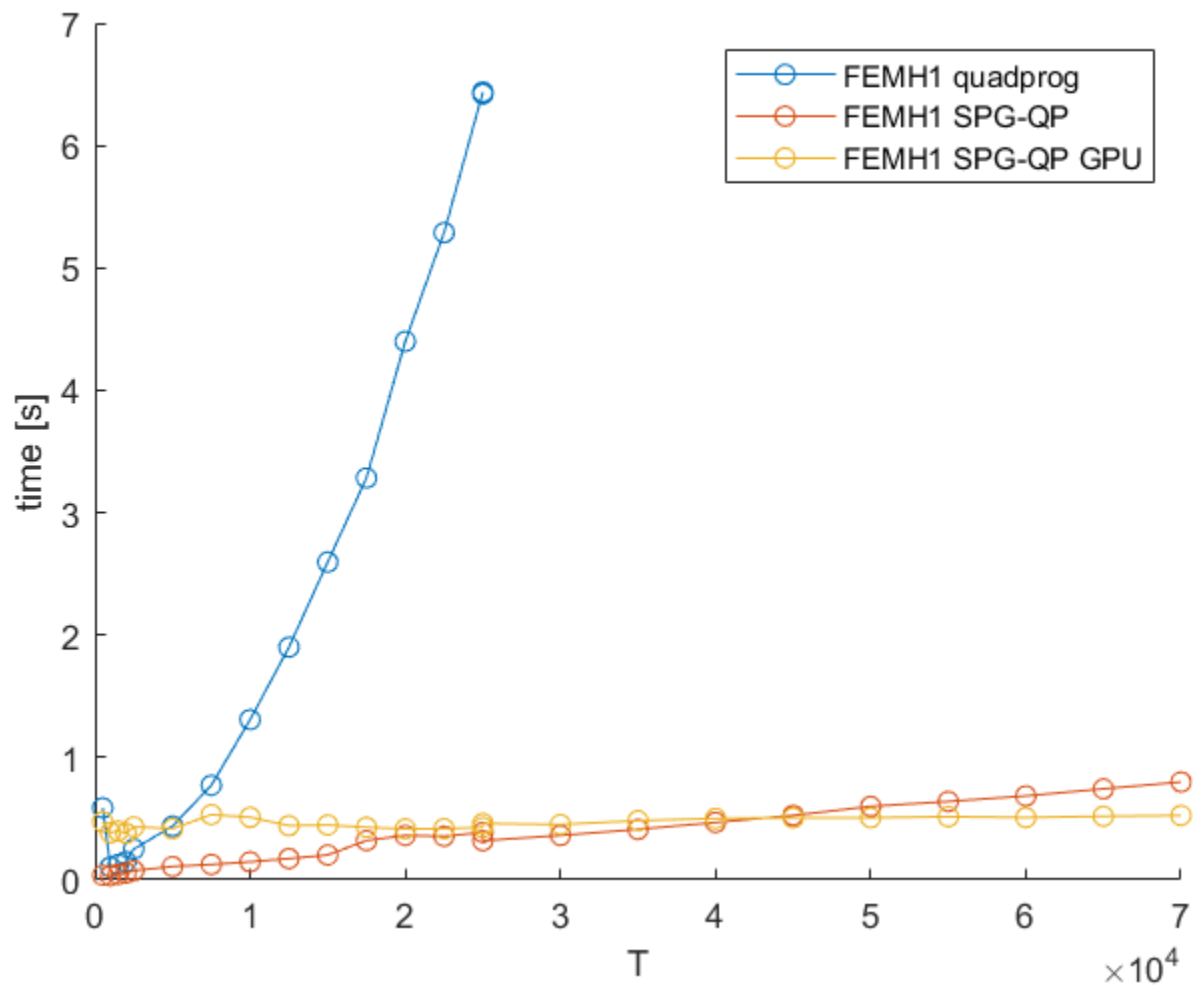


Algorithm matters

solve QP problem

$$\Gamma^* = \arg \min_{\gamma} \gamma^T H \gamma + g^T \gamma$$

subject to $B\gamma = c$ and $\gamma \geq 0$

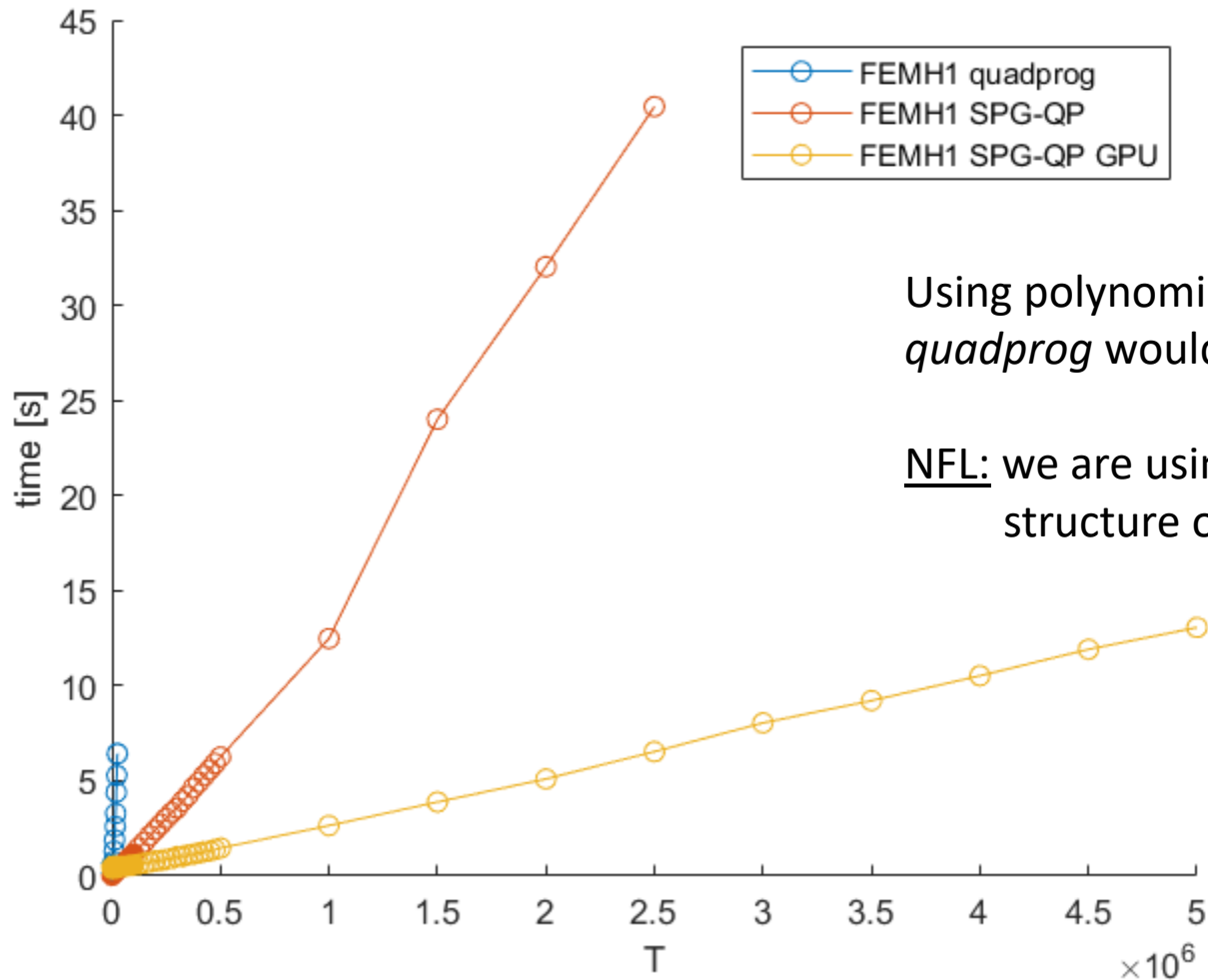


Algorithm matters

solve QP problem

$$\Gamma^* = \arg \min_{\gamma} \gamma^T H \gamma + g^T \gamma$$

subject to $B\gamma = c$ and $\gamma \geq 0$



Using polynomial regression (assuming $O(T^2)$)
quadprog would solve $T=5 \cdot 10^6$ in 2,5 days

NFL: we are using the specific
structure of the problem

Outline

1.) Unconstrained optimization problem

(some theory and definitions)

2.) Equality constrained optimization problem

(Lagrange function, Lagrange multipliers)

3.) Inequality constrained optimization problem

(Karush-Kuhn-Tucker)

1.) Unconstrained optimization problems

Optimization problem

What is maximum on Ω ? (of $f : \Omega \rightarrow \mathbb{R}$) $\Omega \subset \mathcal{V}$

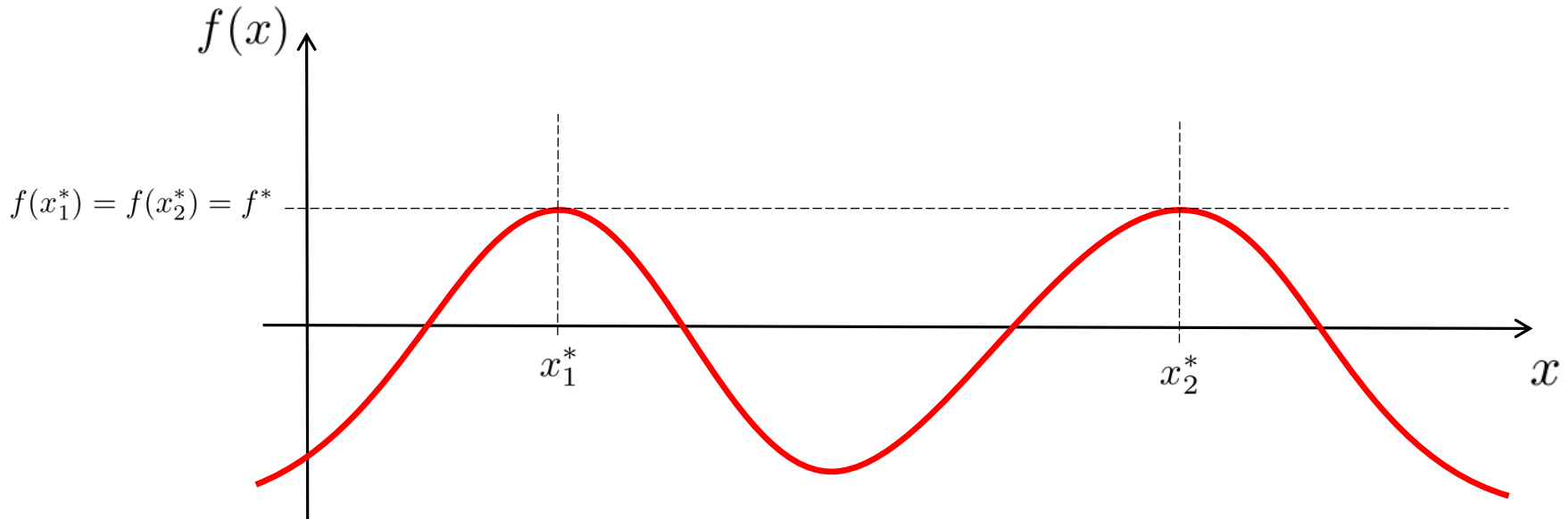
$$f^* = \max_{x \in \Omega} f(x)$$

Find $f^* \in \{f(x) | x \in \Omega\}$ such that $\forall x \in \Omega : f(x) \leq f^*$

$$x^* = \arg \max_{x \in \Omega} f(x)$$

Find $x^* \in \Omega$ such that $\forall x \in \Omega : f(x) \leq f(x^*)$

Example:



What is minimum?

What is maximum on Ω ? (of $f : \Omega \rightarrow \mathbb{R}$) $\Omega \subset \mathcal{V}$

$$f^* = \max_{x \in \Omega} f(x)$$

Find $f^* \in \{f(x) | x \in \Omega\}$ such that $\forall x \in \Omega : f(x) \leq f^*$

$$x^* = \arg \max_{x \in \Omega} f(x)$$

Find $x^* \in \Omega$ such that $\forall x \in \Omega : f(x) \leq f(x^*)$



What is minimum on Ω ? (of $f : \Omega \rightarrow \mathbb{R}$) $\Omega \subset \mathcal{V}$

$$f^* = \min_{x \in \Omega} f(x)$$

Find $f^* \in \{f(x) | x \in \Omega\}$ such that $\forall x \in \Omega : f(x) \geq f^*$

$$x^* = \arg \min_{x \in \Omega} f(x)$$

Find $x^* \in \Omega$ such that $\forall x \in \Omega : f(x) \geq f(x^*)$

What is minimum?

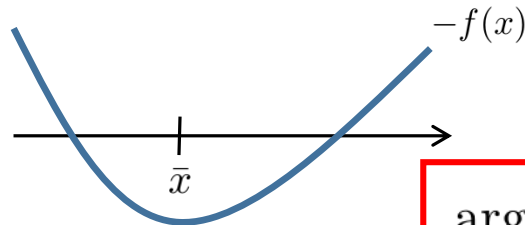
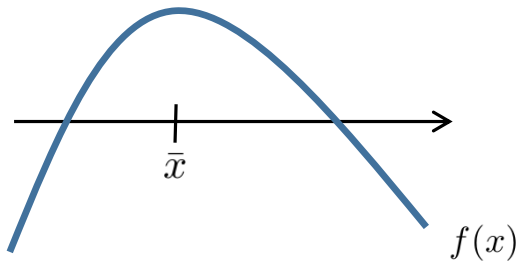
What is maximum on Ω ? (of $f : \Omega \rightarrow \mathbb{R}$) $\Omega \subset \mathcal{V}$

$$f^* = \max_{x \in \Omega} f(x)$$

Find $f^* \in \{f(x) | x \in \Omega\}$ such that $\forall x \in \Omega : f(x) \leq f^*$

$$x^* = \arg \max_{x \in \Omega} f(x)$$

Find $x^* \in \Omega$ such that $\forall x \in \Omega : f(x) \leq f(x^*)$



$$\arg \max f(x) = \arg \min -f(x)$$

$$\max f(x) = -\min(-f(x))$$

What is minimum on Ω ? (of $f : \Omega \rightarrow \mathbb{R}$) $\Omega \subset \mathcal{V}$

$$f^* = \min_{x \in \Omega} f(x)$$

Find $f^* \in \{f(x) | x \in \Omega\}$ such that $\forall x \in \Omega : f(x) \geq f^*$

$$x^* = \arg \min_{x \in \Omega} f(x)$$

Find $x^* \in \Omega$ such that $\forall x \in \Omega : f(x) \geq f(x^*)$

Using the definition

Exercise:

using the definition prove that:

$$\forall \alpha > 0 \forall \beta \in \mathbb{R} : \arg \min_{x \in \Omega} f(x) = \arg \min_{x \in \Omega} (\alpha f(x) + \beta)$$

$$\text{if } \forall x \in \Omega : f(x) \geq 0 \text{ then } \arg \min_{x \in \Omega} f(x) = \arg \min_{x \in \Omega} f^2(x)$$

$$\text{if } \forall x \in \Omega : f(x) > 0 \text{ then } \arg \max_{x \in \Omega} f(x) = \arg \max_{x \in \Omega} \log(f(x))$$

Reminder:

$$x^* = \arg \max_{x \in \Omega} f(x)$$

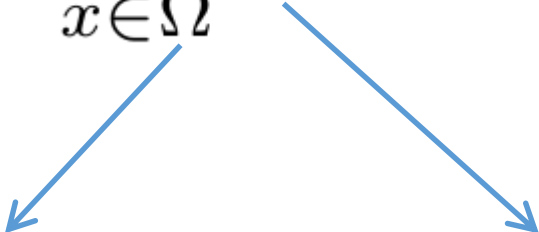
Find $x^* \in \Omega$ such that $\forall x \in \Omega : f(x) \leq f(x^*)$

$$x^* = \arg \min_{x \in \Omega} f(x)$$

Find $x^* \in \Omega$ such that $\forall x \in \Omega : f(x) \geq f(x^*)$

Function f is called *monotonically increasing* on Ω if $\forall x, y \in \Omega : x < y \Leftrightarrow f(x) < f(y)$

Solvability issues

$$\min_{x \in \Omega} f(x)$$


The diagram shows two blue arrows originating from the expression $\min_{x \in \Omega} f(x)$. One arrow points from the domain $x \in \Omega$ to the 'feasible set' section, and the other points from the function $f(x)$ to the 'objective function' section.

feasible set (defined by constraints)

- non-empty?
- convex?
- closed?
- bounded?
- integer?

...

objective function ("cost" function)

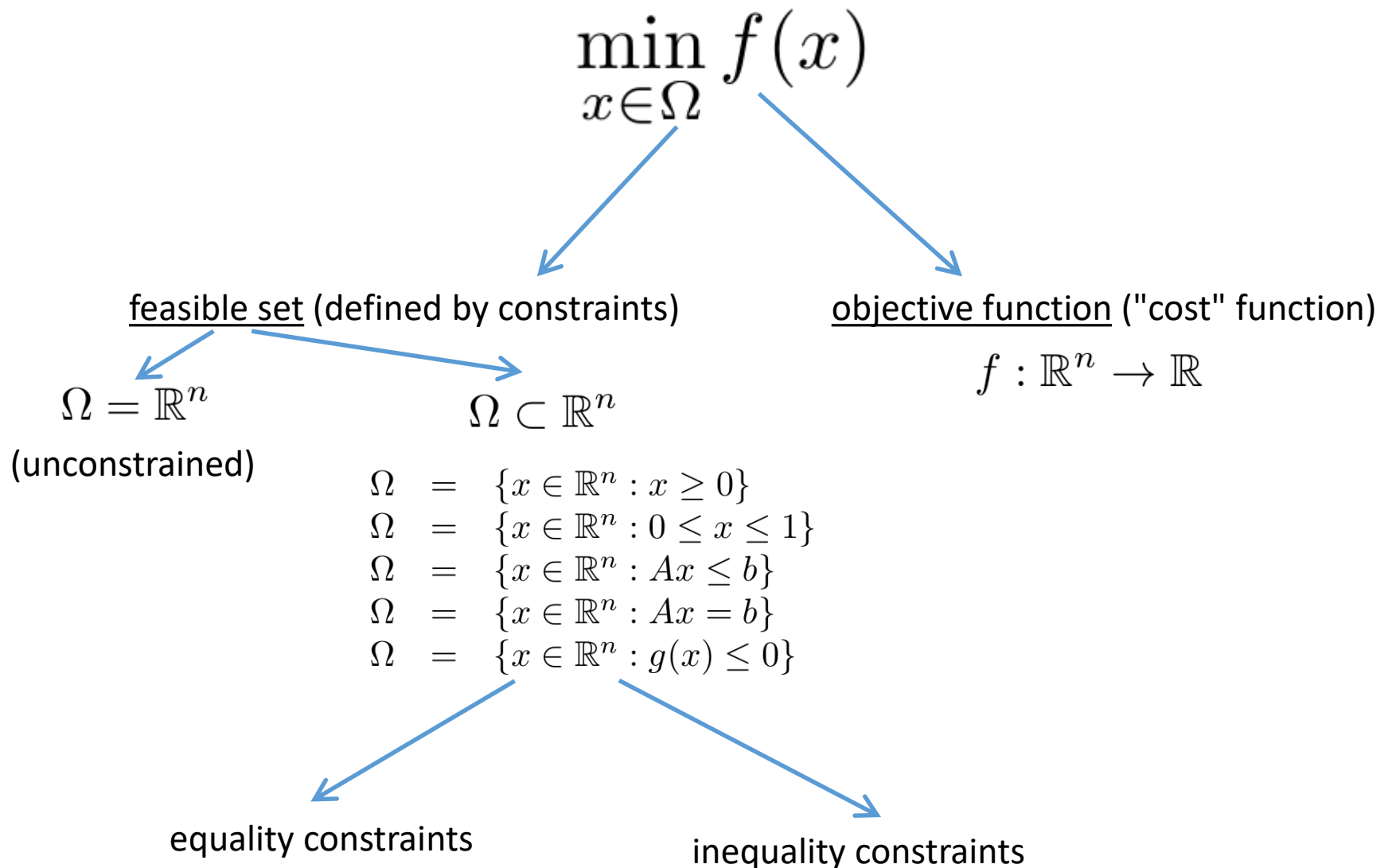
- continuous?
- defined everywhere in feasible set?
- convex? strictly/quasi convex?
- bounded from below?
- differentiable?

...

Properties define the solvability of the optimization problem.

... and many problems are still open!

What is optimization problem?



Basic theorem of solvability

$$\min_{x \in \Omega} f(x)$$

(WEIERSTRASS EXTREME VALUE THEOREM.)

If f is a real-valued continuous function on a non-empty compact (i.e. bounded and closed) domain Ω , then there exists $x \in \Omega$ such that $f(x) \geq f(y)$ for all $y \in \Omega$.

(BOUNDED SET.)

The set Ω is bounded, if there exists a real constant $M > 0$ such that

$$\forall x \in \Omega : \|x\| \leq M .$$

(CLOSED SET.)

The set Ω is closed if for any sequence of points $\{x^k\}$ in Ω , all limit points of this sequence belong to Ω .

$$x^* = \arg \min_{x \in \mathbb{R}^n} f(x) \quad (P)$$

First order necessary condition for unconstrained problem

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be differentiable at point $x^* \in \mathbb{R}^n$.
If x^* is a solution of (P) , then $\nabla f(x^*) = 0$.

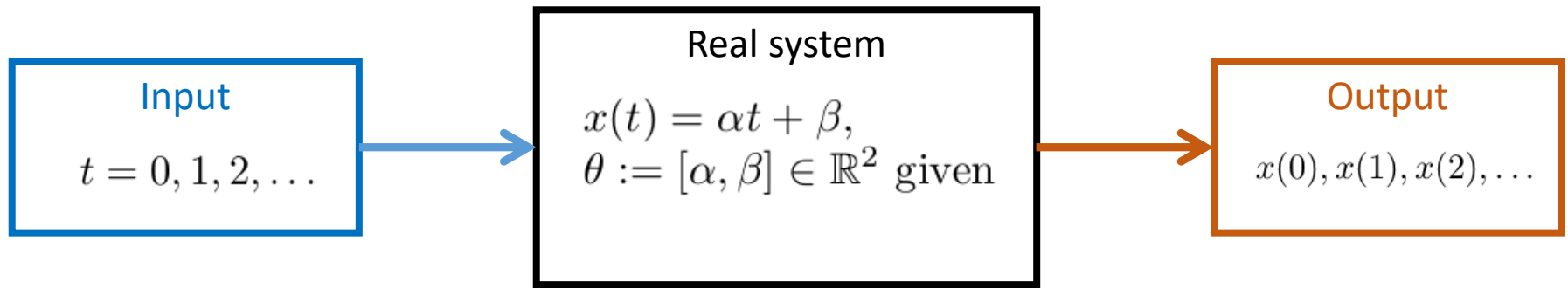
Second order necessary condition for unconstrained problem

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be twice differentiable.
If x^* is a solution of (P) , then $\nabla f(x^*) = 0$ and $\nabla^2 f(x^*) \succ 0$.

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be twice differentiable, $\nabla f(x^*) = 0$, and $\nabla^2 f(x^*) \succ 0$.
Then x^* is *local* solution of (P) .

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is strictly convex function.
Then if (P) has solution, then this solution is unique.

Example: (Linear) Regression



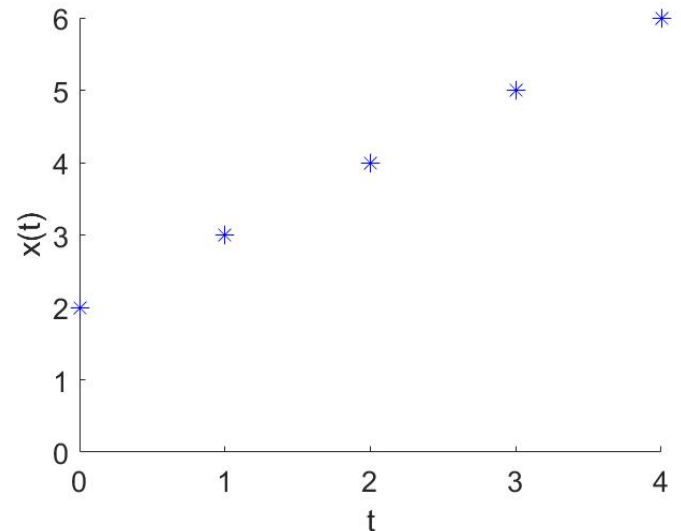
Example:

Uniform motion describes an object that is moving in a specific direction at a constant speed.

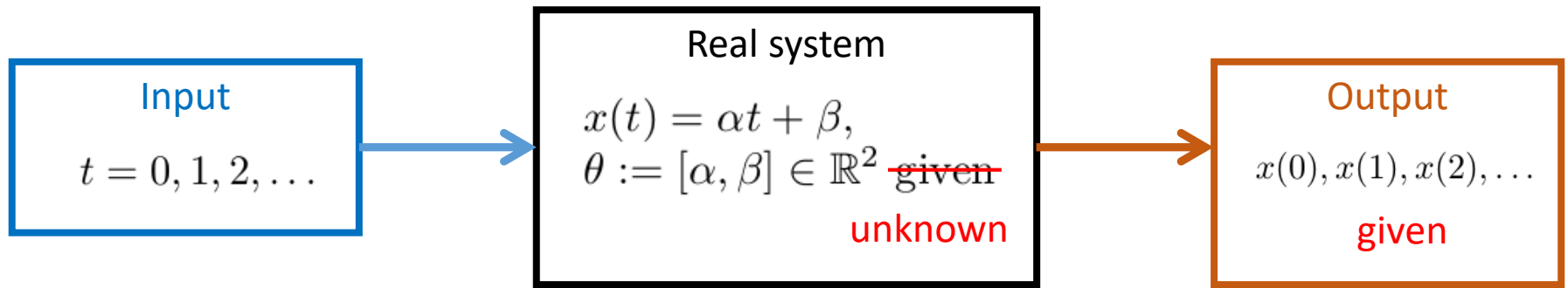
$$s = v \cdot t + s_0$$

$$v = 1, s_0 = 2$$

$$x(0) = 2, x(1) = 3, x(2) = 4, x(3) = 5, x(4) = 6$$



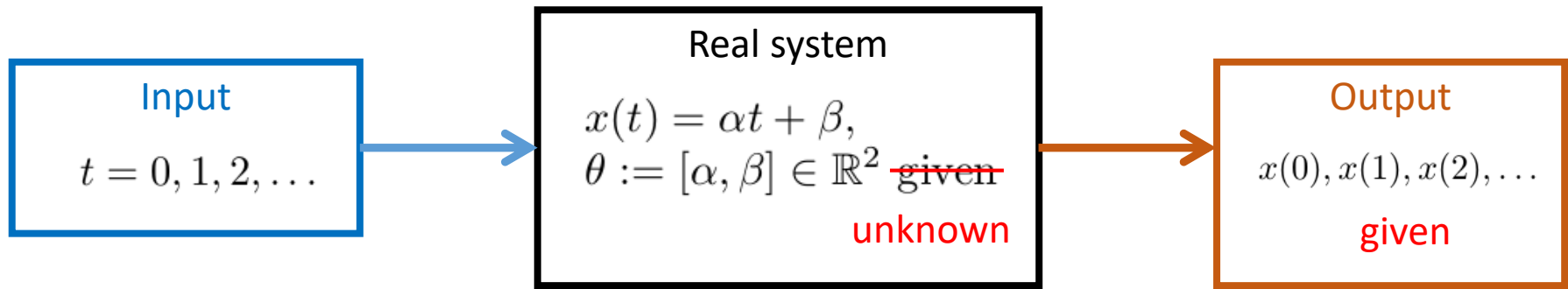
Example: (Linear) Regression



Example:

$$x(0) = 2, x(1) = 3, x(2) = 4, x(3) = 5, x(4) = 6$$

Example: (Linear) Regression

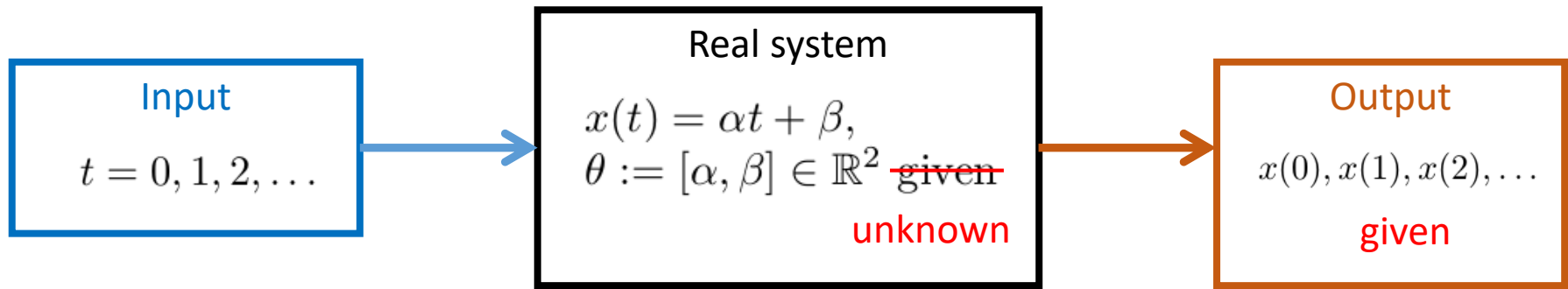


Example:

$$x(0) = 2, x(1) = 3, x(2) = 4, x(3) = 5, x(4) = 6$$

$$\begin{array}{lclcl} x(0) & = & \alpha \cdot 0 + \beta & \Rightarrow & 2 = \beta \\ x(1) & = & \alpha \cdot 1 + \beta & \Rightarrow & 3 = 1\alpha + \beta \\ x(2) & = & \alpha \cdot 2 + \beta & \Rightarrow & 4 = 2\alpha + \beta \\ x(3) & = & \alpha \cdot 3 + \beta & \Rightarrow & 5 = 3\alpha + \beta \\ x(4) & = & \alpha \cdot 4 + \beta & \Rightarrow & 6 = 4\alpha + \beta \end{array} \quad \left. \vphantom{\begin{array}{lclcl} x(0) \\ x(1) \\ x(2) \\ x(3) \\ x(4) \end{array}} \right\} \begin{array}{l} B\theta = c \\ \text{(system of linear equations)} \end{array}$$

Example: (Linear) Regression



Example:

$$x(0) = 2, x(1) = 3, x(2) = 4, x(3) = 5, x(4) = 6$$

$$\begin{aligned} x(0) &= \alpha \cdot 0 + \beta &\Rightarrow \\ x(1) &= \alpha \cdot 1 + \beta &\Rightarrow \\ x(2) &= \alpha \cdot 2 + \beta &\Rightarrow \\ x(3) &= \alpha \cdot 3 + \beta &\Rightarrow \\ x(4) &= \alpha \cdot 4 + \beta &\Rightarrow \end{aligned}$$

$$\begin{aligned} 2 &= \beta \\ 3 &= 1\alpha + \beta \\ 4 &= 2\alpha + \beta \\ 5 &= 3\alpha + \beta \\ 6 &= 4\alpha + \beta \end{aligned}$$

redundant

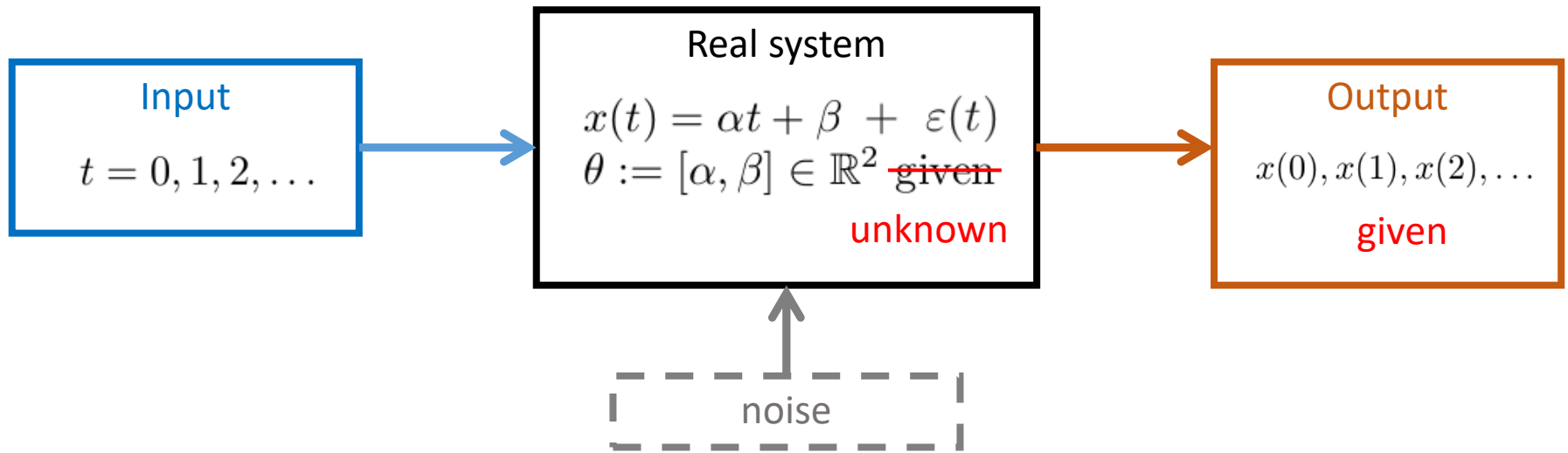
$$\alpha = 1, \beta = 2$$

satisfies all equations

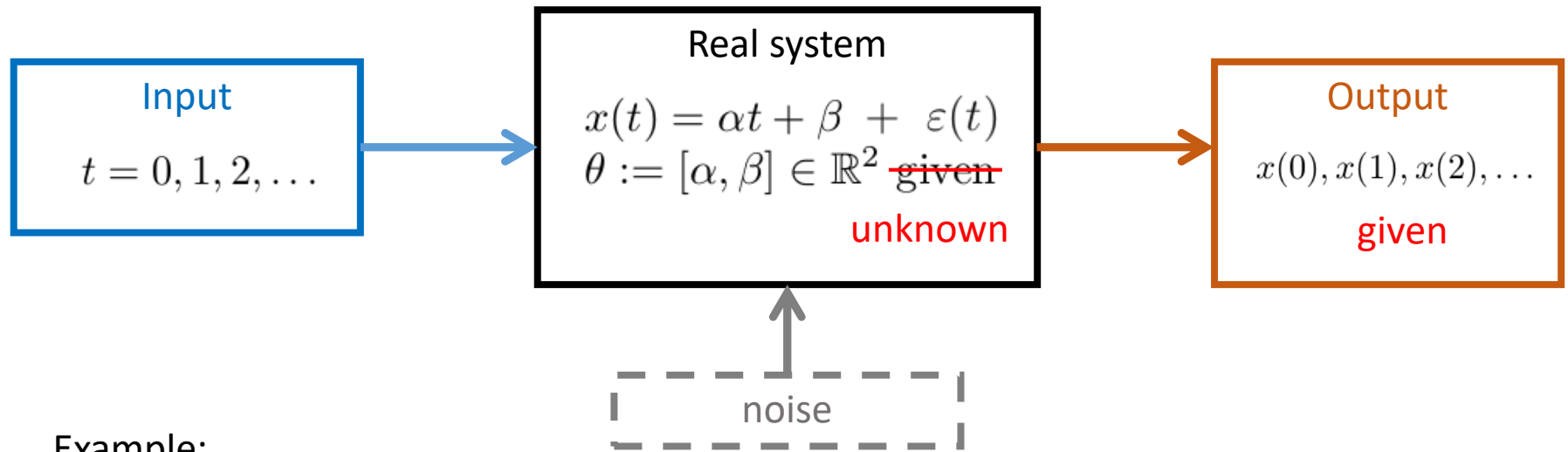
$$B\theta = c$$

(system of linear equations)

(Linear) Regression



(Linear) Regression



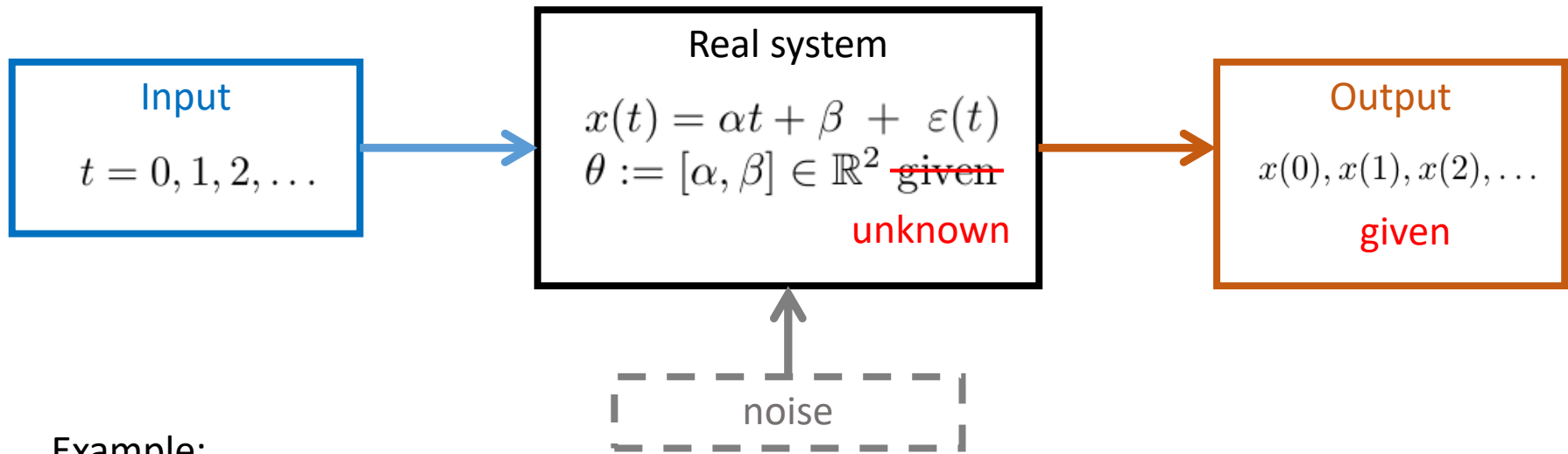
Example:

~~$x(0) = 2, x(1) = 3, x(2) = 4, x(3) = 5, x(4) = 6$~~

$x(0) = 1.9865, x(1) = 3.0303, x(2) = 4.0073, x(3) = 4.9994, x(4) = 6.0071$

$x(0)$	$=$	$\alpha \cdot 0 + \beta$	\Rightarrow	1.9865	$=$	β
$x(1)$	$=$	$\alpha \cdot 1 + \beta$	\Rightarrow	3.0303	$=$	$1\alpha + \beta$
$x(2)$	$=$	$\alpha \cdot 2 + \beta$	\Rightarrow	4.0073	$=$	$2\alpha + \beta$
$x(3)$	$=$	$\alpha \cdot 3 + \beta$	\Rightarrow	4.9994	$=$	$3\alpha + \beta$
$x(4)$	$=$	$\alpha \cdot 4 + \beta$	\Rightarrow	6.0071	$=$	$4\alpha + \beta$

(Linear) Regression



Example:

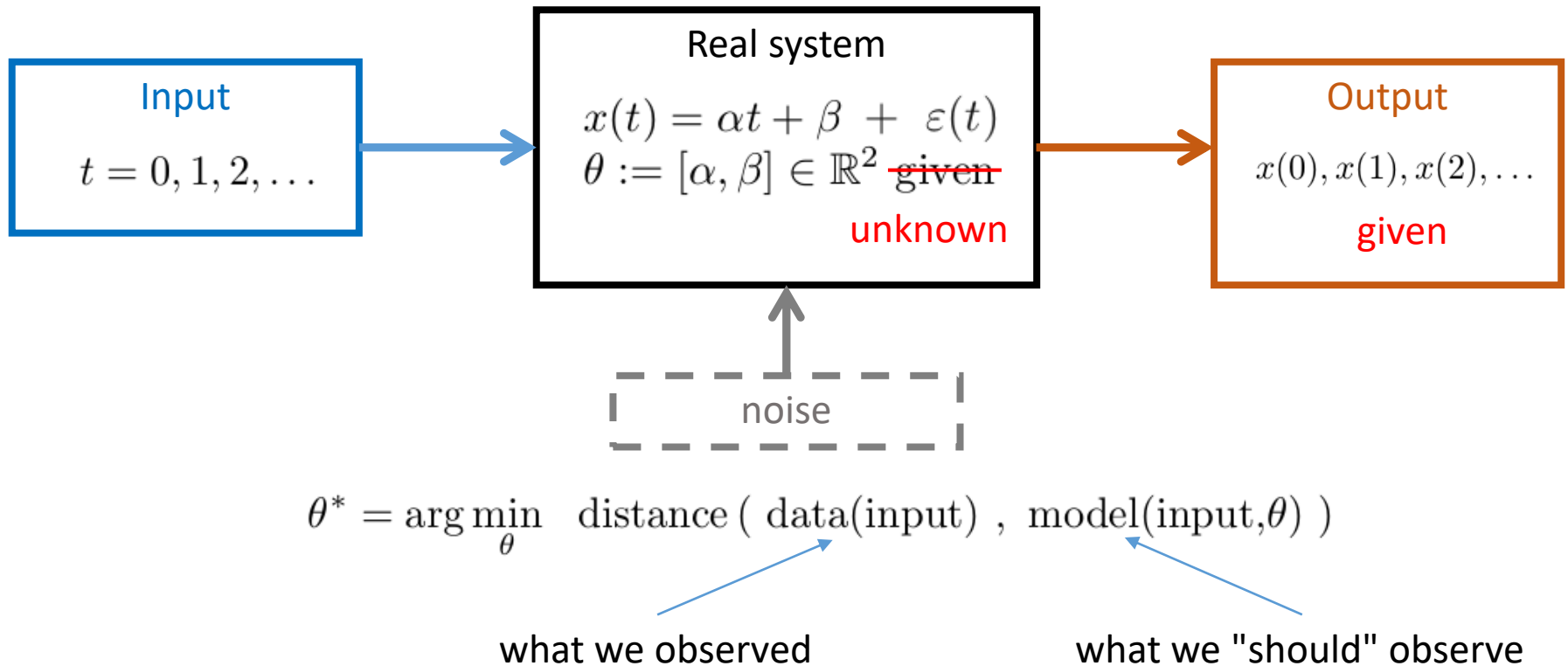
~~$x(0) = 2, x(1) = 3, x(2) = 4, x(3) = 5, x(4) = 6$~~

$x(0) = 1.9865, x(1) = 3.0303, x(2) = 4.0073, x(3) = 4.9994, x(4) = 6.0071$

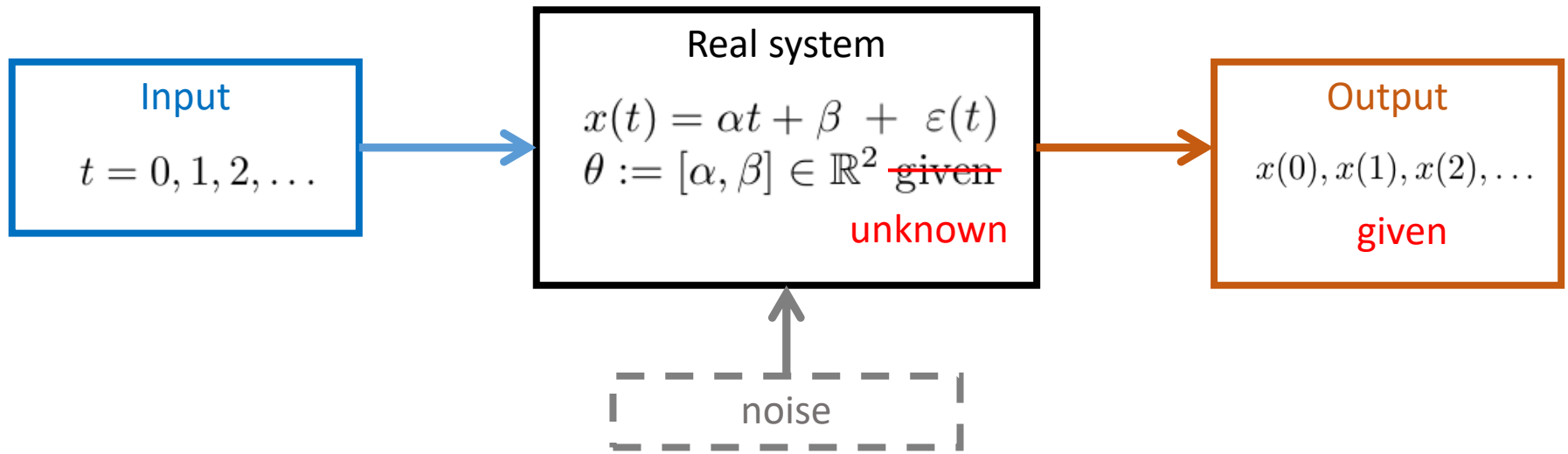
$$\begin{array}{llll} x(0) & = & \alpha \cdot 0 + \beta & \Rightarrow 1.9865 = \beta \\ x(1) & = & \alpha \cdot 1 + \beta & \Rightarrow 3.0303 = 1\alpha + \beta \\ x(2) & = & \alpha \cdot 2 + \beta & \Rightarrow 4.0073 = 2\alpha + \beta \\ x(3) & = & \alpha \cdot 3 + \beta & \Rightarrow 4.9994 = 3\alpha + \beta \\ x(4) & = & \alpha \cdot 4 + \beta & \Rightarrow 6.0071 = 4\alpha + \beta \end{array}$$

does not have solution :(

(Linear) Regression



(Linear) Regression



$$\theta^* = \arg \min_{\theta} \text{distance} (\text{data}(\text{input}) , \text{model}(\text{input}, \theta))$$

Example:

$t = 0, 1, \dots$ input

$x(0), x(1), x(2), \dots$ data

$m(t, \theta) := \alpha t + \beta$ model (linear)

$$\rho(x(\cdot), m(\cdot, \theta)) := \sum_{t=0}^{T-1} \|x(t) - m(t, \theta)\|_2^2 \quad (\text{least square error})$$

(Linear) Regression

$$\theta^* = \arg \min \sum_{t=0}^{T-1} \|x(t) - m(t, \theta)\|_2^2 = \arg \min \sum_{t=0}^{T-1} (x(t) - m(t, \theta))^2 = \arg \min \underbrace{\|B\theta - c\|_2^2}_{=: f(\theta)}$$

$x_t \in \mathbb{R}$ $m(t, \theta) := \alpha t + \beta$

Example:

$t = 0, 1, \dots$ input

$x(0), x(1), x(2), \dots$ data

$m(t, \theta) := \alpha t + \beta$ model (linear)

$$\rho(x(\cdot), m(\cdot, \theta)) := \sum_{t=0}^{T-1} \|x(t) - m(t, \theta)\|_2^2 \quad (\text{least square error})$$

(Linear) Regression

$$\theta^* = \arg \min \sum_{t=0}^{T-1} \|x(t) - m(t, \theta)\|_2^2 = \arg \min \sum_{t=0}^{T-1} (x(t) - m(t, \theta))^2 = \arg \min \|B\theta - c\|_2^2$$

$x_t \in \mathbb{R}$ $m(t, \theta) := \alpha t + \beta$ $=: f(\theta)$

Reminder:

$$\begin{aligned} x(0) &= \alpha \cdot 0 + \beta \\ x(1) &= \alpha \cdot 1 + \beta \\ x(2) &= \alpha \cdot 2 + \beta \\ x(3) &= \alpha \cdot 3 + \beta \\ x(4) &= \alpha \cdot 4 + \beta \end{aligned}$$

$$\begin{array}{cc} & t^1 \quad t^0 \\ \begin{bmatrix} 0 & 1 \\ 1 & 1 \\ 2 & 1 \\ 3 & 1 \\ 4 & 1 \end{bmatrix} & \begin{bmatrix} \alpha \\ \beta \end{bmatrix} = \begin{bmatrix} x(0) \\ x(1) \\ x(2) \\ x(3) \\ x(4) \end{bmatrix} \\ \underbrace{\hspace{1.5cm}} & \underbrace{\hspace{1.5cm}} & \underbrace{\hspace{1.5cm}} \\ B & \theta & c \end{array}$$

(Linear) Regression

$$\theta^* = \arg \min \sum_{t=0}^{T-1} \|x(t) - m(t, \theta)\|_2^2 = \arg \min \sum_{t=0}^{T-1} (x(t) - m(t, \theta))^2 = \arg \min \|B\theta - c\|_2^2$$

$x_t \in \mathbb{R}$ $m(t, \theta) := \alpha t + \beta$ $\underbrace{\hspace{10em}}_{=: f(\theta)}$

$$\begin{aligned} f(\theta) &= \|B\theta - c\|_2^2 = \langle B\theta - c, B\theta - c \rangle \\ &= \theta^T B^T B \theta - 2c^T B \theta + c^T c \end{aligned}$$

$$\begin{aligned} \nabla f(\theta) &= 2B^T B \theta - 2B^T c \\ \nabla^2 f(\theta) &= 2B^T B \end{aligned}$$

$$\nabla f(\theta) = 0 \quad \Leftrightarrow \quad (B^T B)\theta = B^T c$$

$$\forall y \in \mathbb{R}^2 : y^T B^T B y = \|By\|^2 \geq 0 \quad \Rightarrow \quad f(\theta) \text{ is convex.}$$

this system has always solution...

(Polynomial) Regression

degree of polynomial model

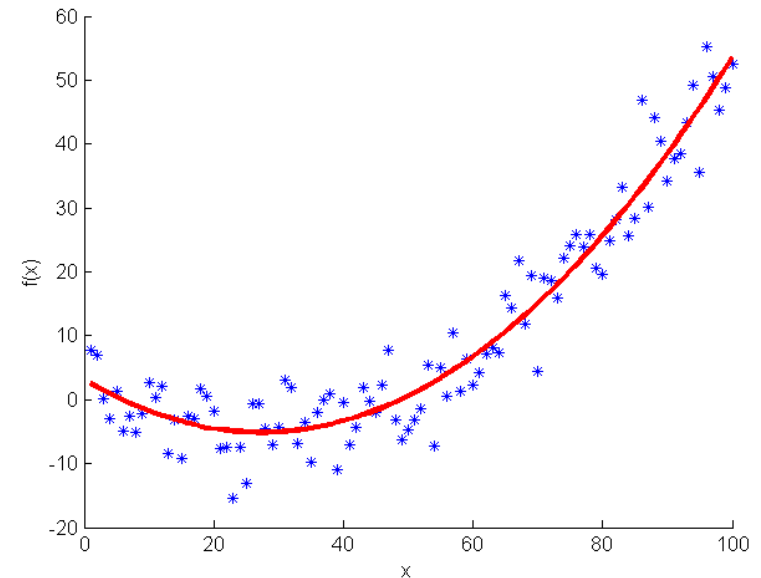
$$m(t, \theta) := \sum_{k=0}^P \theta_k t^k, \quad t = 0, 1, \dots, T-1$$

$$\theta := \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \\ \vdots \\ \theta_P \end{bmatrix} \in \mathbb{R}^{P+1} \quad c := \begin{bmatrix} x(0) \\ x(1) \\ x(2) \\ x(3) \\ \vdots \\ x(T-1) \end{bmatrix} \in \mathbb{R}^{T-1}$$

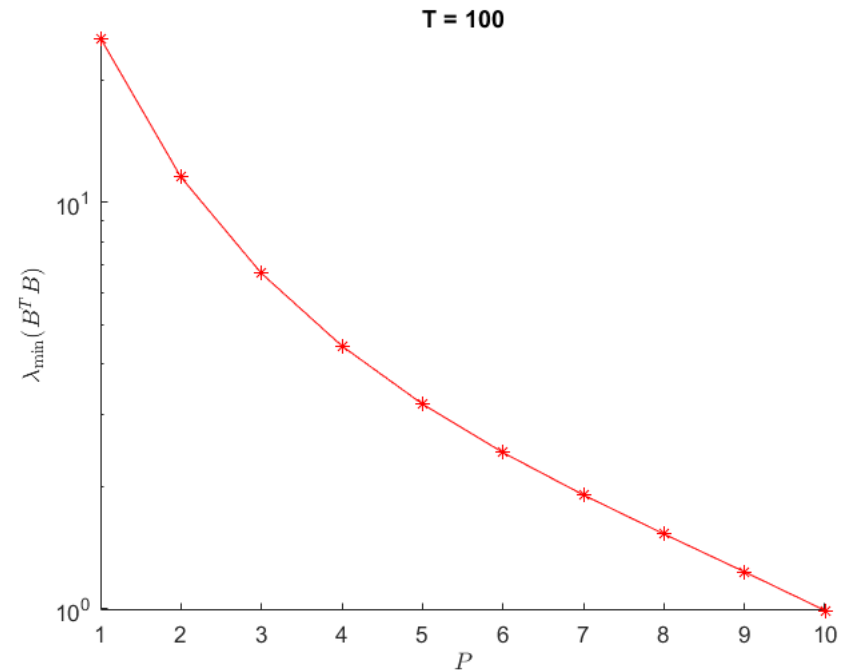
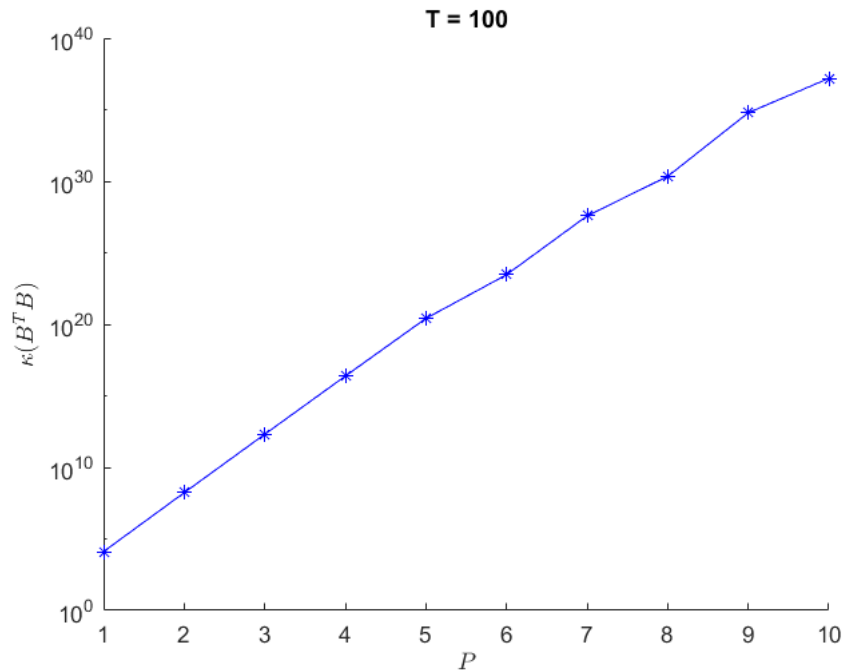
$$B := \begin{bmatrix} 1 & 0^1 & 0^2 & \dots & 0^P \\ 1 & 1^1 & 1^2 & \dots & 1^P \\ 1 & 2^1 & 2^2 & \dots & 2^P \\ 1 & 3^1 & 3^2 & \dots & 3^P \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & (T-1)^1 & (T-1)^2 & \dots & (T-1)^P \end{bmatrix} \in \mathbb{R}^{T-1, P+1}$$

solve:

$$(B^T B)\theta = B^T c$$



(Polynomial) Regression



$$B := \begin{bmatrix} 1 & 0^1 & 0^2 & \dots & 0^P \\ 1 & 1^1 & 1^2 & \dots & 1^P \\ 1 & 2^1 & 2^2 & \dots & 2^P \\ 1 & 3^1 & 3^2 & \dots & 3^P \\ \vdots & & & & \\ 1 & (T-1)^1 & (T-1)^2 & \dots & (T-1)^P \end{bmatrix} \in \mathbb{R}^{T-1, P+1}$$

solve:

$$(B^T B)\theta = B^T c$$

2.) Equality constrained optimization problem

Optimization with equality constraints

$$x^* = \arg \min_{x \in \Omega} f(x), \quad \Omega := \{x \in \mathbb{R}^n : h_i(x) = 0, i = 1, \dots, m\}$$

$f : \mathbb{R}^n \rightarrow \mathbb{R}$ is objective function,
 $\Omega \subset \mathbb{R}^n$ is feasible set,
 $h_i : \mathbb{R}^n \rightarrow \mathbb{R}$ is i -th equality constraint (m constraints).

Necessary optimality conditions:

Let x be an optimality point. Then there exist $\lambda \in \mathbb{R}^m$ such that

$$\begin{aligned} \nabla_x L(x, \lambda) &= \nabla f(x) + \sum_{i=1}^m \lambda_i \nabla h_i(x) = 0, \\ \nabla_{\lambda_i} L(x, \lambda) &= h_i(x) = 0, \quad i = 1, \dots, m, \end{aligned}$$

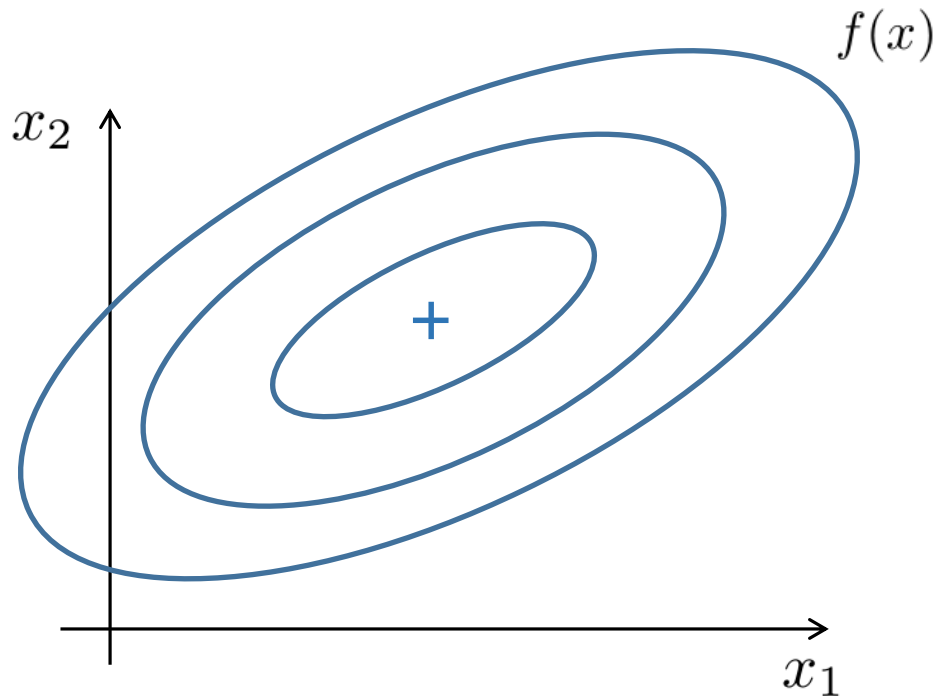
where $L : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$ is Lagrange function defined as

$$L(x, \lambda) := f(x) + \sum_{i=1}^m \lambda_i h_i(x).$$

Optimization with equality constraints

"proof"

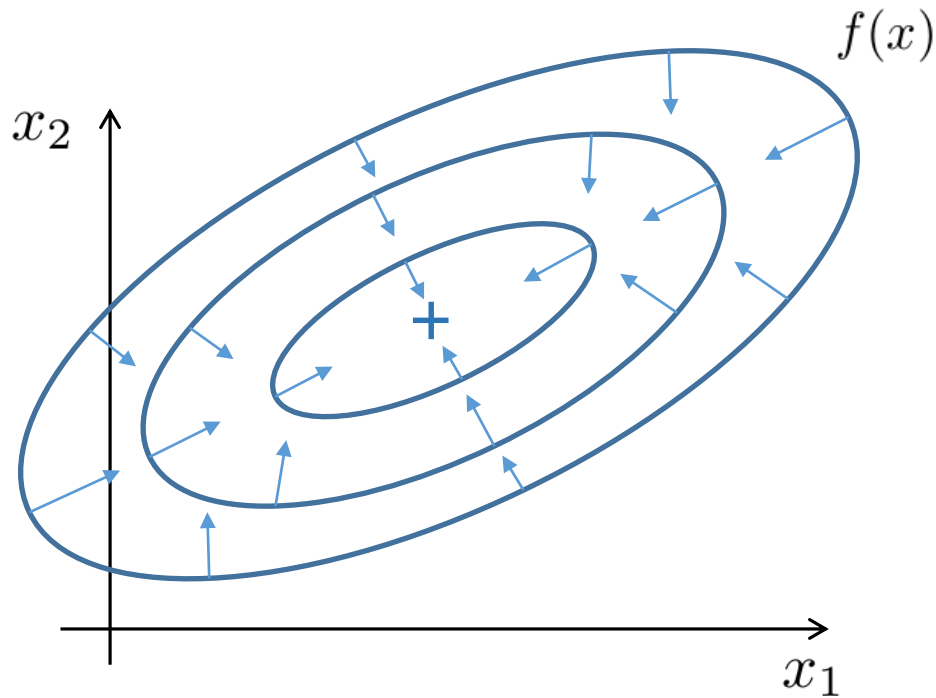
$$\begin{array}{rcl} \nabla f(x) + \lambda \nabla h(x) & = & 0 \\ h(x) & = & 0 \end{array}$$



Optimization with equality constraints

"proof"

$$\begin{array}{rcl} \nabla f(x) + \lambda \nabla h(x) & = & 0 \\ h(x) & = & 0 \end{array}$$

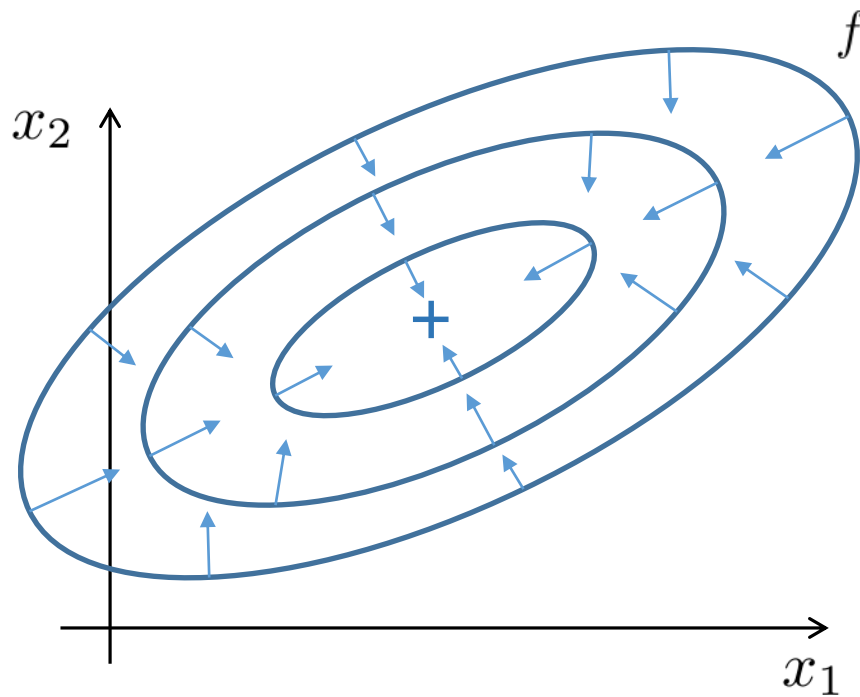


$-\nabla f(x)$ is vector of steepest descent

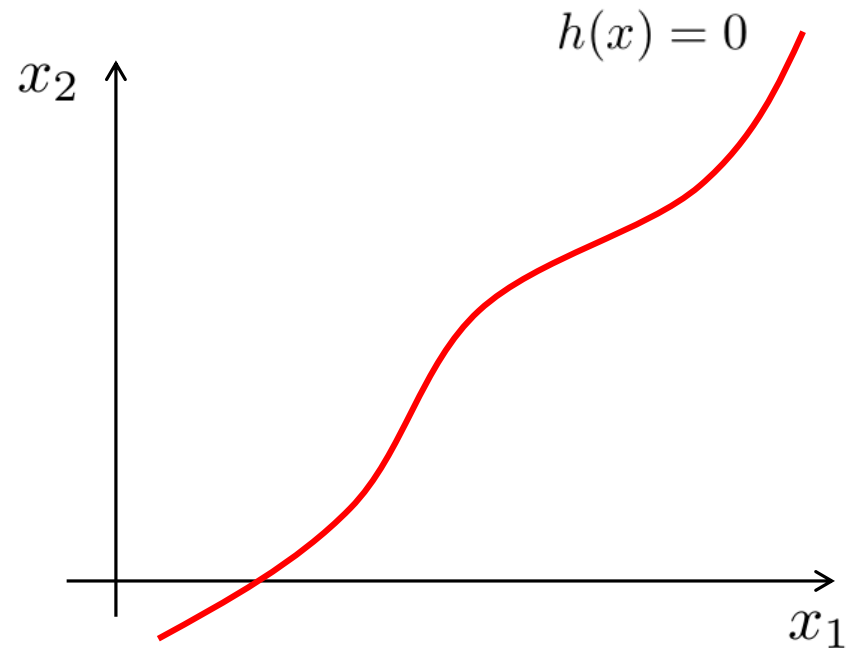
Optimization with equality constraints

"proof"

$$\begin{aligned}\nabla f(x) + \lambda \nabla h(x) &= 0 \\ h(x) &= 0\end{aligned}$$



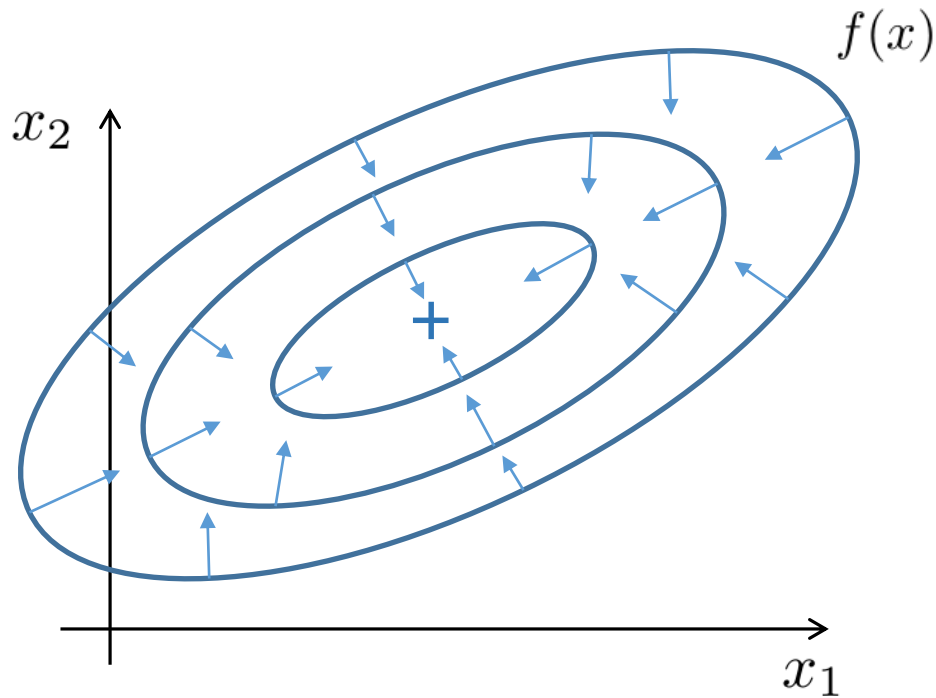
$-\nabla f(x)$ is vector of steepest descent



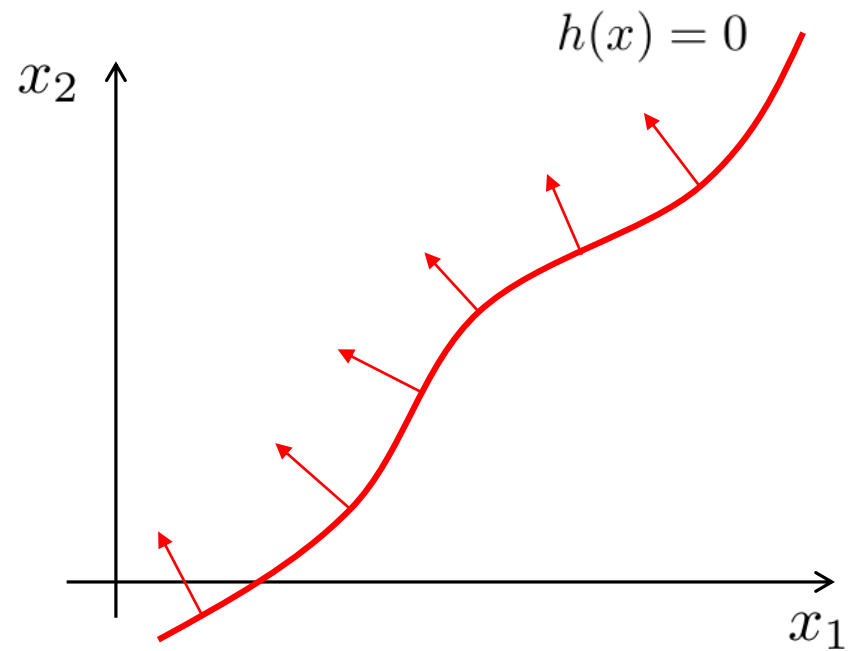
Optimization with equality constraints

"proof"

$$\begin{aligned}\nabla f(x) + \lambda \nabla h(x) &= 0 \\ h(x) &= 0\end{aligned}$$



$-\nabla f(x)$ is vector of steepest descent

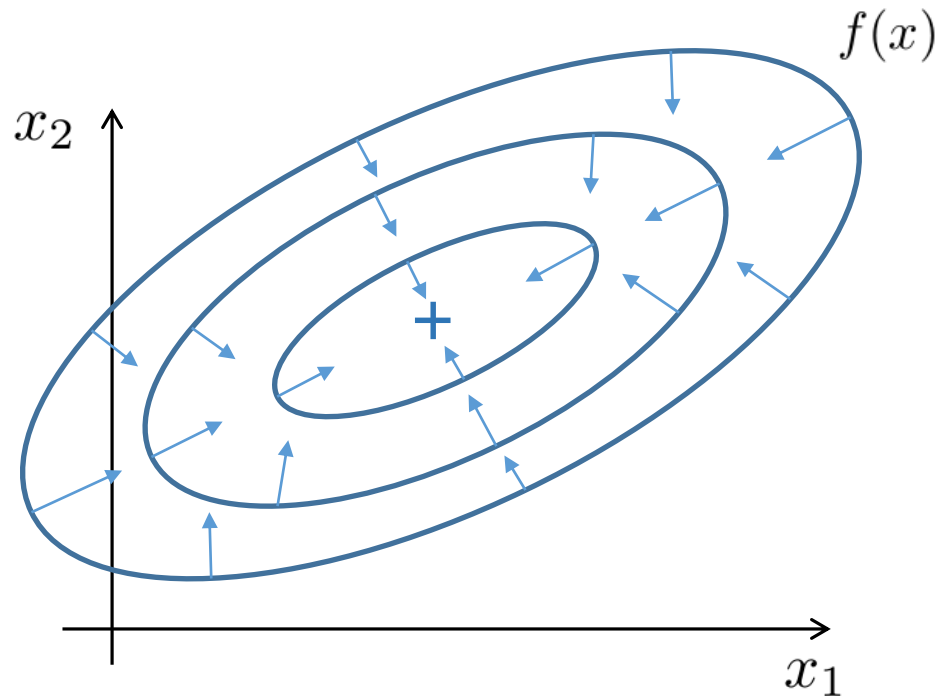


$\nabla h(x)$ is normal vector

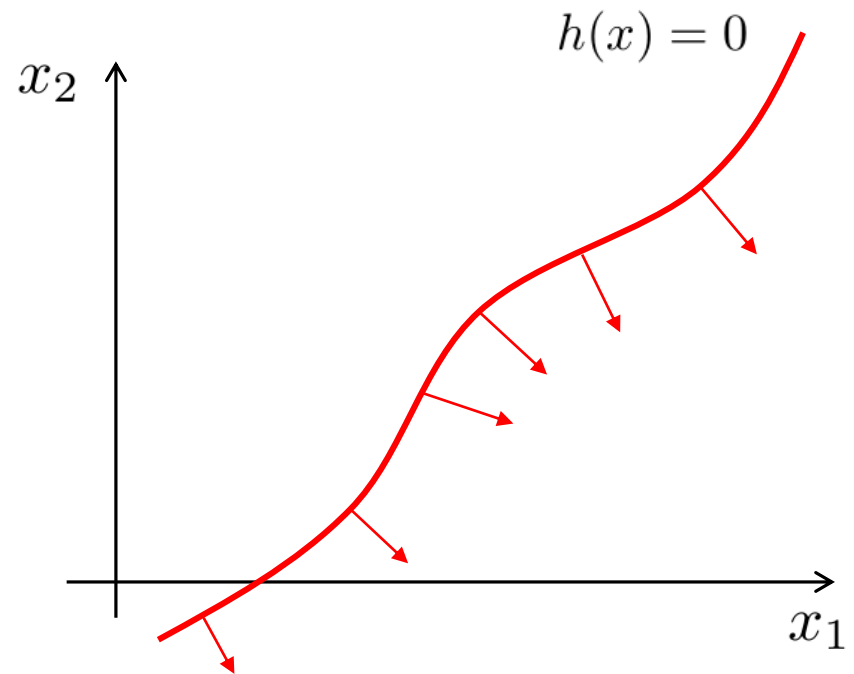
Optimization with equality constraints

"proof"

$$\begin{aligned}\nabla f(x) + \lambda \nabla h(x) &= 0 \\ h(x) &= 0\end{aligned}$$



$-\nabla f(x)$ is vector of steepest descent



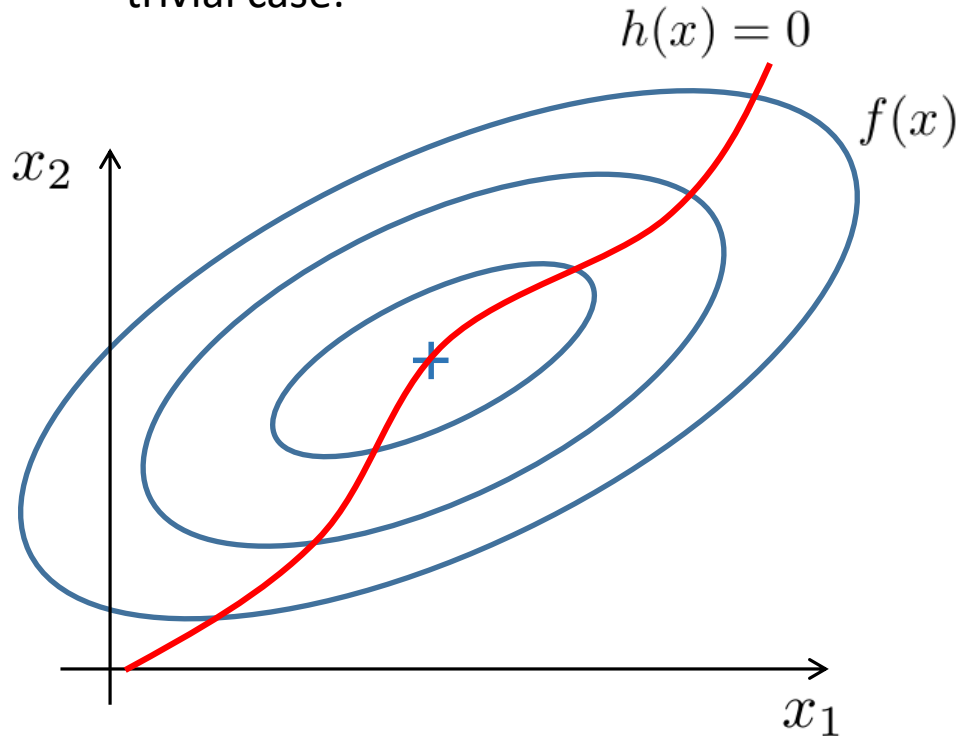
$\nabla h(x)$ is normal vector

$$\begin{aligned}h(x) = 0 &\Leftrightarrow -h(x) = 0 \\ \nabla h(x) &\quad \quad -\nabla h(x)\end{aligned}$$

Optimization with equality constraints

"proof"

- trivial case:



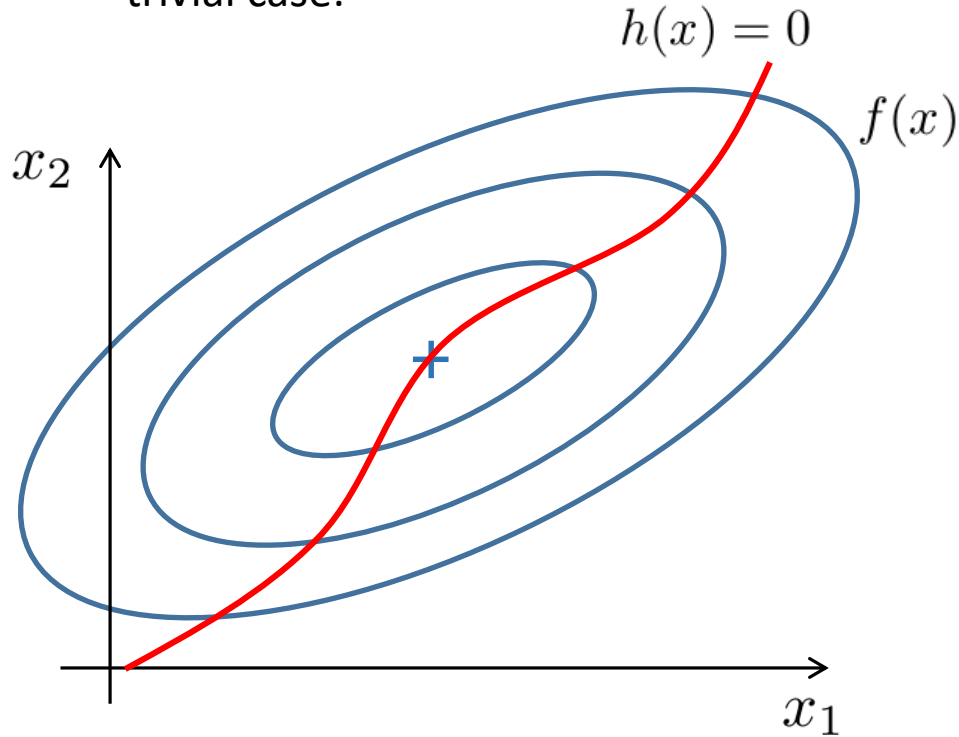
$$\begin{aligned}\nabla f(x) + \lambda \nabla h(x) &= 0 \\ h(x) &= 0\end{aligned}$$

$$\arg \min_{x \in \mathbb{R}^n} f(x) \in \Omega$$

Optimization with equality constraints

"proof"

- trivial case:



condition for unconstrained problem

$$\begin{aligned}\nabla f(x) + \cancel{\lambda \nabla h(x)} &= 0 \\ h(x) &= 0\end{aligned}$$

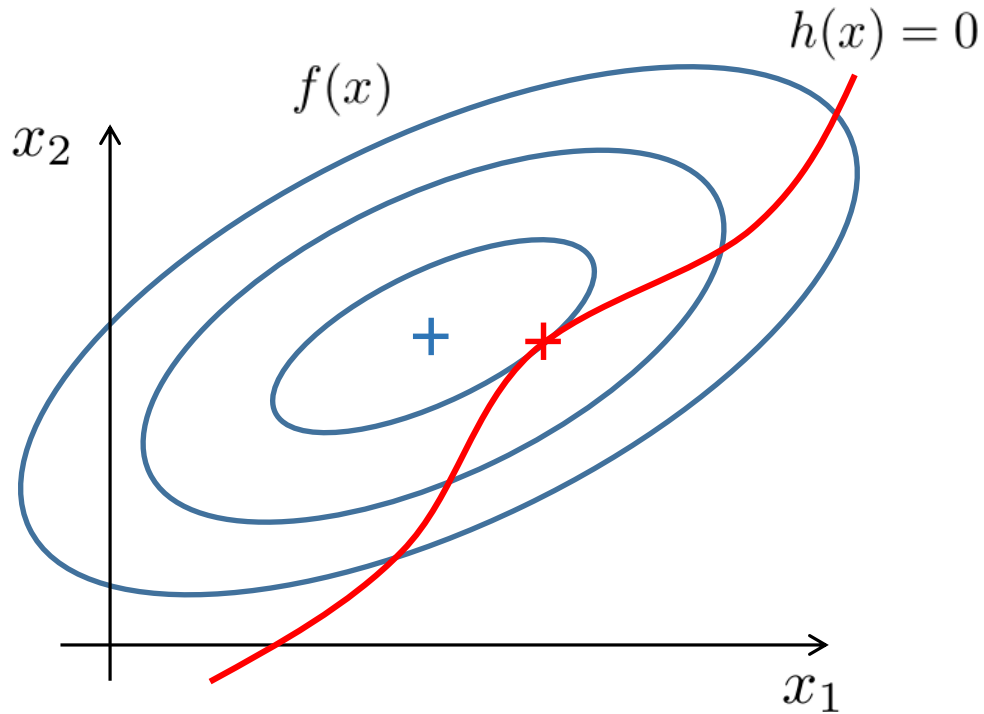
$$\arg \min_{x \in \mathbb{R}^n} f(x) \in \Omega$$

$$\lambda = 0$$

Optimization with equality constraints

"proof"

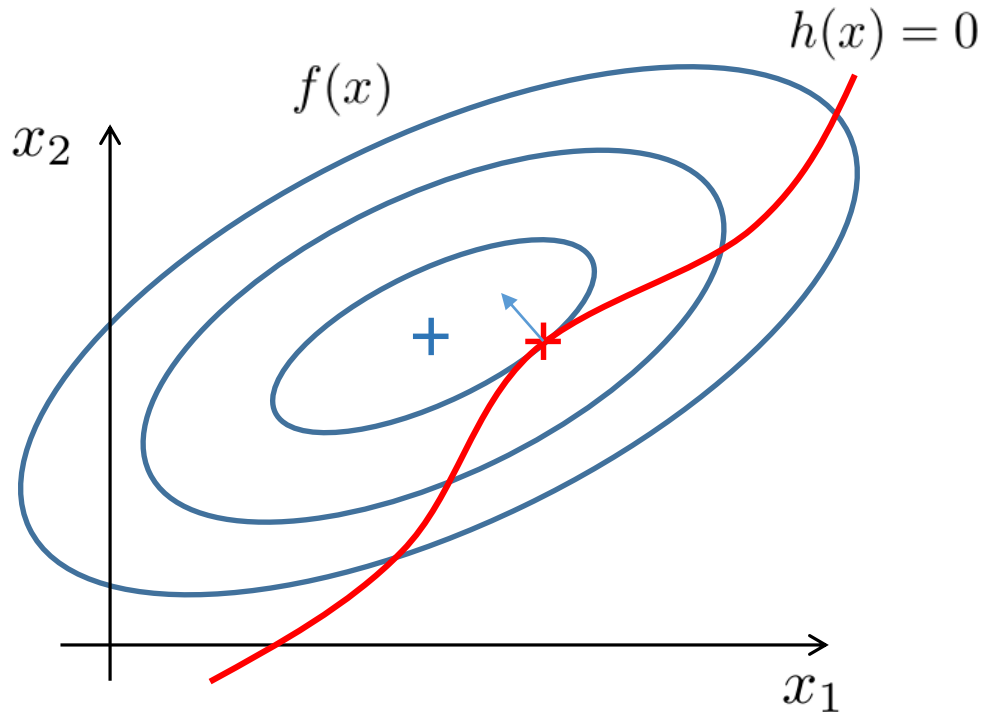
$$\begin{array}{rcl} \nabla f(x) + \lambda \nabla h(x) & = & 0 \\ h(x) & = & 0 \end{array}$$



Optimization with equality constraints

"proof"

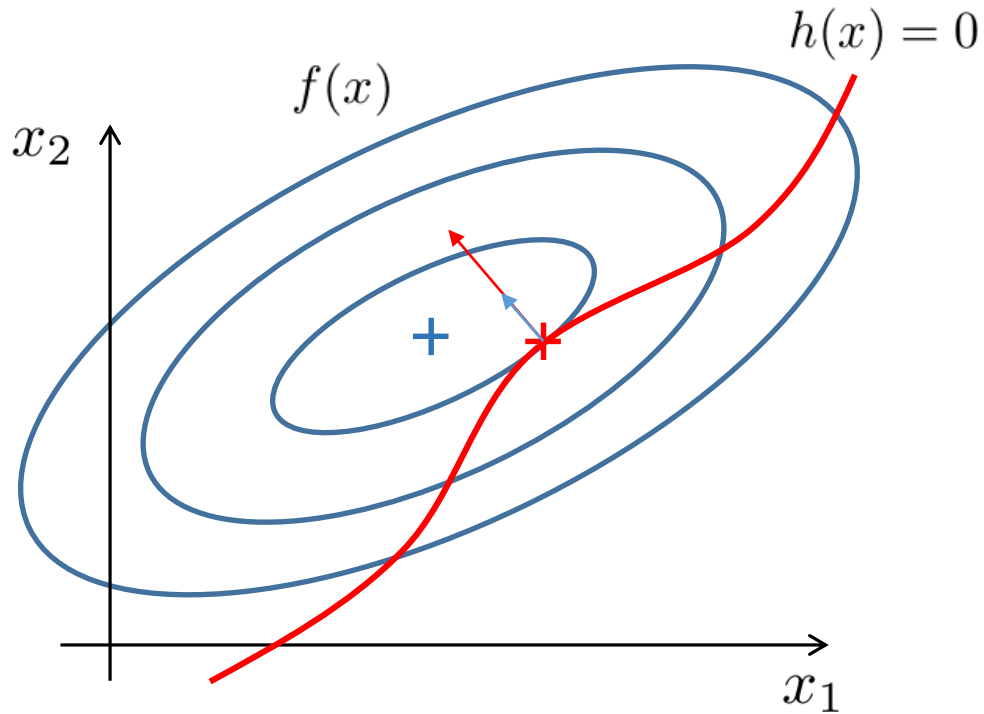
$$\begin{array}{rcl} \nabla f(x) + \lambda \nabla h(x) & = & 0 \\ h(x) & = & 0 \end{array}$$



Optimization with equality constraints

"proof"

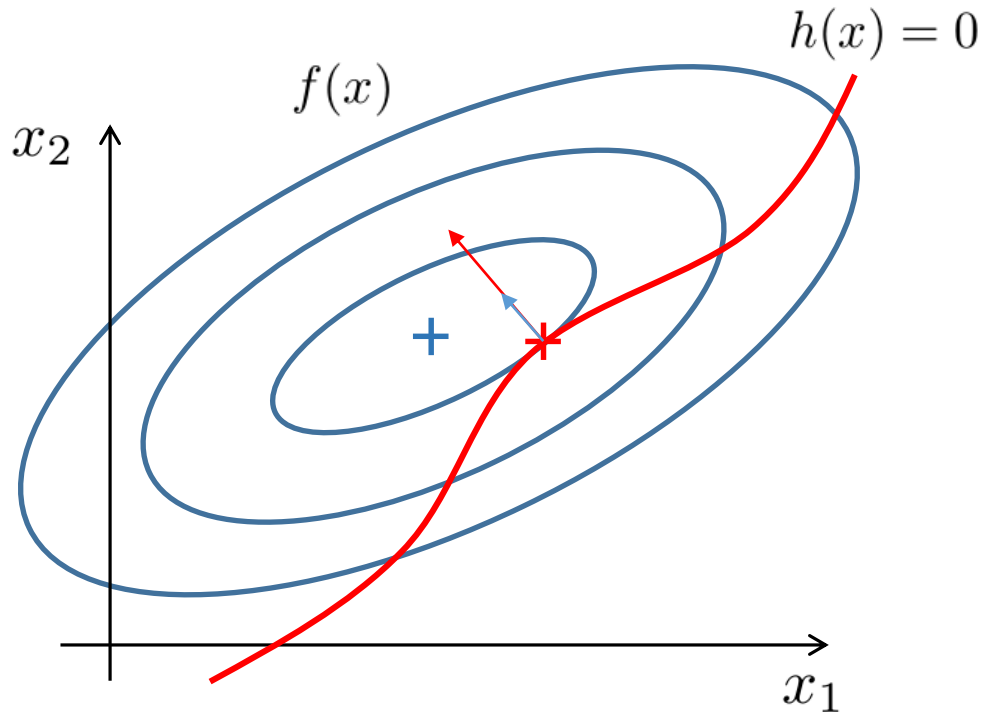
$$\begin{array}{rcl} \nabla f(x) + \lambda \nabla h(x) & = & 0 \\ h(x) & = & 0 \end{array}$$



Optimization with equality constraints

"proof"

$$\begin{array}{rcl} \nabla f(x) + \lambda \nabla h(x) & = & 0 \\ h(x) & = & 0 \end{array}$$

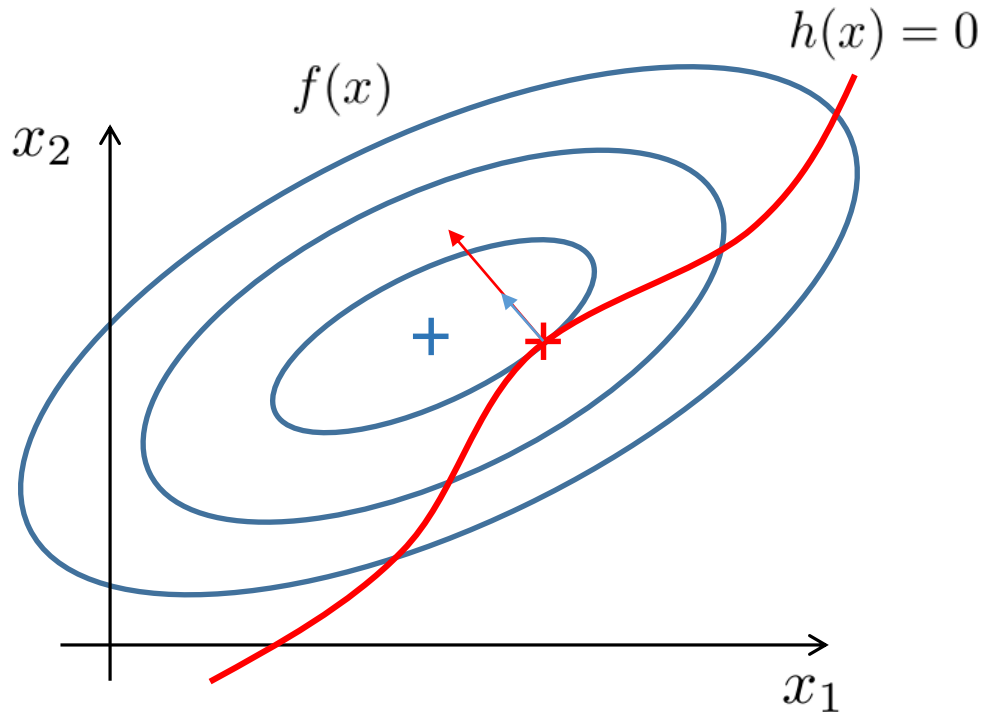


- $\Rightarrow \nabla f(x)$ and $\nabla h(x)$ have “same” direction (or opposite)
- $\Rightarrow \nabla f(x)$ and $\nabla h(x)$ are linearly dependent
- \Rightarrow there exists coefficient $\alpha \in \mathbb{R}$ such that $\nabla f(x) = \alpha \nabla h(x)$

Optimization with equality constraints

"proof"

$$\begin{array}{rcl} \nabla f(x) + \lambda \nabla h(x) & = & 0 \\ h(x) & = & 0 \end{array}$$



$\Rightarrow \nabla f(x)$ and $\nabla h(x)$ have “same” direction (or opposite)

$\Rightarrow \nabla f(x)$ and $\nabla h(x)$ are linearly dependent

\Rightarrow there exists coefficient $\cancel{\alpha} \in \mathbb{R}$ such that $\nabla f(x) = \cancel{\alpha} \nabla h(x)$

$-\lambda$

$-\lambda$

Equality constrained problems

Example:

$$\arg \min_{x_1+x_2=2} x_1^2 + x_2^2 + x_1x_2 + 1000$$

$$L(x_1, x_2, \lambda) = x_1^2 + x_2^2 + x_1x_2 + 1000 + \lambda(x_1 + x_2 - 2)$$

$$\frac{\partial L}{\partial x_1} = 2x_1 + x_2 + \lambda = 0$$

$$\frac{\partial L}{\partial x_2} = 2x_2 + x_1 + \lambda = 0$$

$$\frac{\partial L}{\partial \lambda} = x_1 + x_2 - 2 = 0$$

$$\bar{x}_1 = \bar{x}_2 = 1, \quad \lambda = -3$$

3.) Inequality constrained optimization problem

Inequality constrained problems

$$\arg \min_{x \in \Omega} f(x), \quad \Omega = \{x \in \mathbb{R}^n : h_i(x) \leq 0, i = 1, \dots, m\}$$

Inequality constrained problems

$$\arg \min_{x \in \Omega} f(x), \quad \Omega = \{x \in \mathbb{R}^n : h_i(x) \leq 0, i = 1, \dots, m\}$$

$$L(x, \lambda) = f(x) + \sum_{i=1}^m \lambda_i h_i(x) \quad \text{Lagrange function}$$

$$\begin{aligned} \nabla_x L(x, \lambda) &= \nabla f(x) + \sum_{i=1}^m \lambda_i \nabla_x h_i(x) &= 0 \\ \nabla_{\lambda_i} L(x, \lambda) &= h_i(x) &\leq 0, \quad i = 1, \dots, m \\ &\lambda_i &\geq 0 \\ &\lambda_i h_i(x) &= 0 \end{aligned}$$

Karush-Kuhn-Tucker optimality conditions

Inequality constrained problems

$$\arg \min_{x \in \Omega} f(x), \quad \Omega = \{x \in \mathbb{R}^n : h_i(x) \leq 0, i = 1, \dots, m\}$$

$$L(x, \lambda) = f(x) + \sum_{i=1}^m \lambda_i h_i(x) \quad \text{Lagrange function}$$

$$\begin{aligned} \nabla_x L(x, \lambda) &= \nabla f(x) + \sum_{i=1}^m \lambda_i \nabla_x h_i(x) = 0 \\ \nabla_{\lambda_i} L(x, \lambda) &= h_i(x) \leq 0, \quad i = 1, \dots, m \\ \lambda_i &\geq 0 \\ \lambda_i h_i(x) &= 0 \end{aligned}$$

!!!

Karush-Kuhn-Tucker optimality conditions

Inequality constrained problems

$$\arg \min_{x \in \Omega} f(x), \quad \Omega = \{x \in \mathbb{R}^n : h_i(x) \leq 0, i = 1, \dots, m\}$$

$$L(x, \lambda) = f(x) + \sum_{i=1}^m \lambda_i h_i(x) \quad \text{Lagrange function}$$

$$\nabla_x L(x, \lambda) = \nabla f(x) + \sum_{i=1}^m \lambda_i \nabla_x h_i(x) = 0$$

$$\nabla_{\lambda_i} L(x, \lambda) = h_i(x) \leq 0, \quad i = 1, \dots, m$$

$$\lambda_i \geq 0$$

$$\boxed{\lambda_i h_i(x) = 0} \quad (\text{complementarity condition})$$

Karush-Kuhn-Tucker optimality conditions

"proof":

1.) $\lambda = 0$ and $h(x) \leq 0$

2.) $\lambda \geq 0$ and $h(x) = 0$

3.) $\lambda = 0$ and $h(x) = 0$

Inequality constrained problems

“proof”:

$$\begin{aligned}\nabla_x L(x, \lambda) &= \nabla f(x) + \sum_{i=1}^m \lambda_i \nabla_x h_i(x) &= 0 \\ \nabla_{\lambda_i} L(x, \lambda) &= h_i(x) &\leq 0, \quad i = 1, \dots, m \\ &\lambda_i &\geq 0 \\ &\lambda_i h_i(x) &= 0\end{aligned}$$

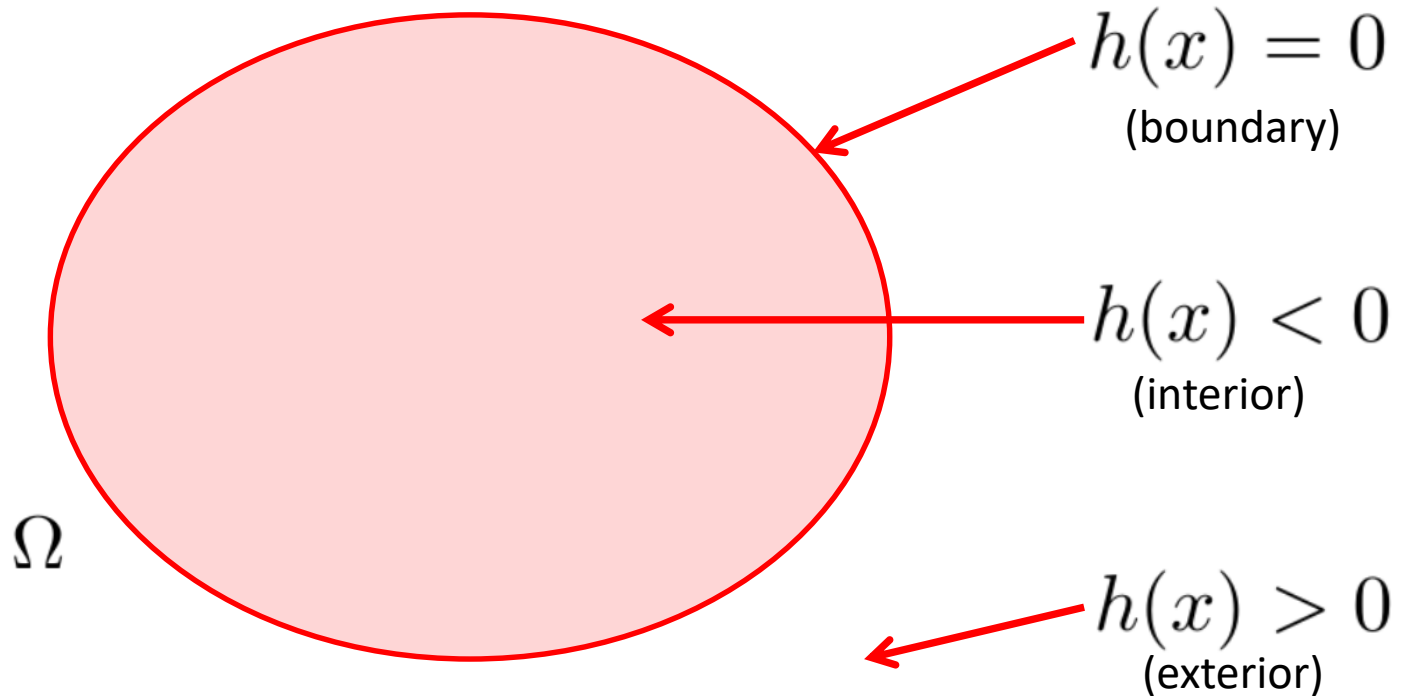
Ω



Inequality constrained problems

“proof”:

$$\begin{aligned}\nabla_x L(x, \lambda) &= \nabla f(x) + \sum_{i=1}^m \lambda_i \nabla_x h_i(x) = 0 \\ \nabla_{\lambda_i} L(x, \lambda) &= h_i(x) \leq 0, \quad i = 1, \dots, m \\ \lambda_i &\geq 0 \\ \lambda_i h_i(x) &= 0\end{aligned}$$

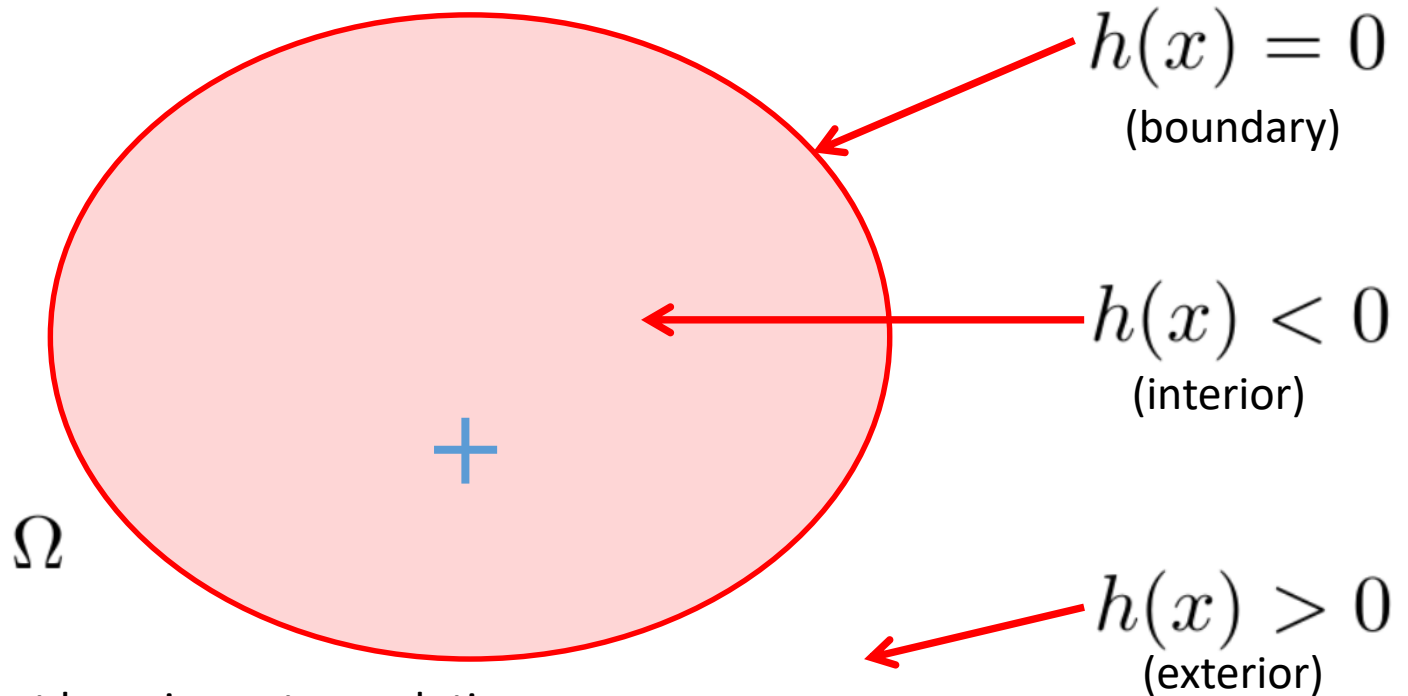


Inequality constrained problems

“proof”:

1.) $\lambda = 0$ and $h(x) \leq 0$

$$\begin{aligned}\nabla_x L(x, \lambda) &= \nabla f(x) + \sum_{i=1}^m \lambda_i \nabla_x h_i(x) = 0 \\ \nabla_{\lambda_i} L(x, \lambda) &= h_i(x) \leq 0, \quad i = 1, \dots, m \\ \lambda_i &\geq 0 \\ \lambda_i h_i(x) &= 0\end{aligned}$$



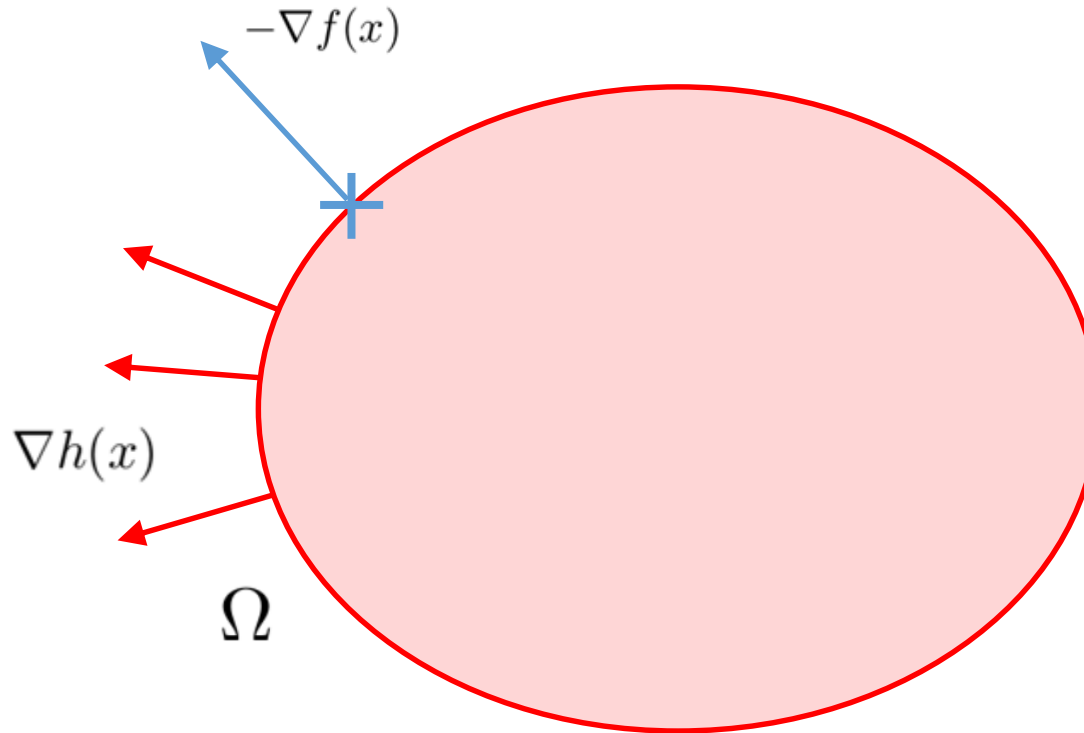
- Constraint does not have impact on solution
 - After substitution of lambda into KKT we get unconstrained condition
- Inequality constrained = Unconstrained

Inequality constrained problems

$$\begin{aligned}\nabla_x L(x, \lambda) &= \nabla f(x) + \sum_{i=1}^m \lambda_i \nabla_x h_i(x) = 0 \\ \nabla_{\lambda_i} L(x, \lambda) &= h_i(x) \leq 0, \quad i = 1, \dots, m \\ \lambda_i &\geq 0 \\ \lambda_i h_i(x) &= 0\end{aligned}$$

“proof”:

2.) $\lambda \geq 0$ and $h(x) = 0$



$h(x) = 0$
(boundary)

$h(x) < 0$
(interior)

$h(x) > 0$
(exterior)

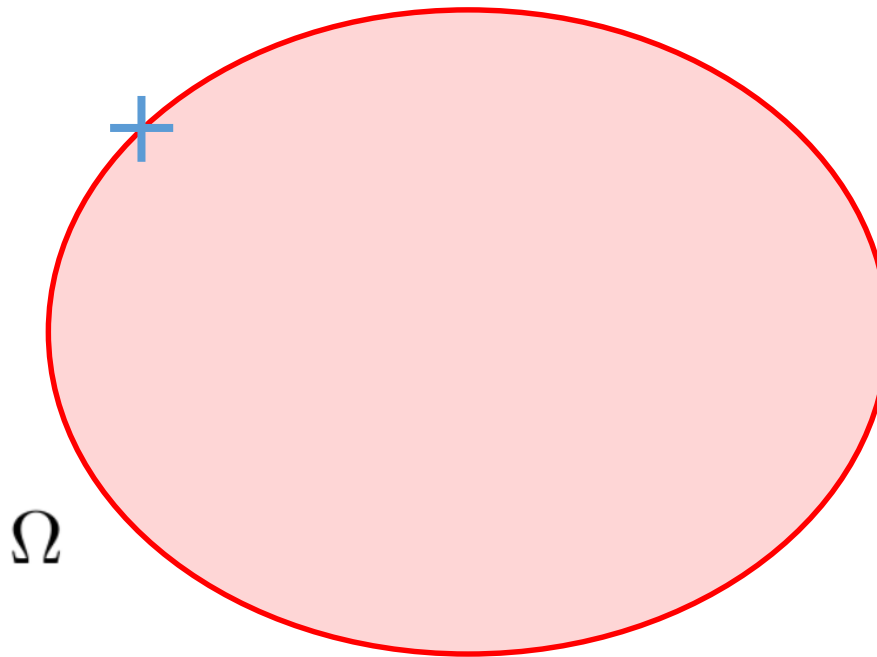
- Solution is on boundary (= equality constrained problem)
- Outer normal of feasible set and $-\text{gradient}$ have the same direction (“function is decreasing out of the feasible set”)

Inequality constrained problems

“proof”:

$$\begin{aligned}\nabla_x L(x, \lambda) &= \nabla f(x) + \sum_{i=1}^m \lambda_i \nabla_x h_i(x) = 0 \\ \nabla_{\lambda_i} L(x, \lambda) &= h_i(x) \leq 0, \quad i = 1, \dots, m \\ \lambda_i &\geq 0 \\ \lambda_i h_i(x) &= 0\end{aligned}$$

3.) $\lambda = 0$ and $h(x) = 0$



$$h(x) = 0$$

(boundary)

$$h(x) < 0$$

(interior)

$$h(x) > 0$$

(exterior)

- Solution is on boundary without impact of the constraint, gradient is equal to zero

Example: Kullback-Leibler estimation of Markov transitional matrix

- consider categorical data:

$$x_t \in \{s_1, \dots, s_n\}, \quad t = 1, \dots, T$$

Example: Kullback-Leibler estimation of Markov transitional matrix

- consider categorical data:

$$x_t \in \{s_1, \dots, s_n\}, \quad t = 1, \dots, T$$

- discrete probability density vectors:

$$\pi_t \in [0, 1]^n, \quad \sum_{i=1}^n \{\pi_t\}_i = 1 \quad \pi_t = \begin{bmatrix} P(x_t = s_1) \\ \vdots \\ P(x_t = s_n) \end{bmatrix}$$

Example: Kullback-Leibler estimation of Markov transitional matrix

- consider categorical data:

$$x_t \in \{s_1, \dots, s_n\}, \quad t = 1, \dots, T$$

- discrete probability density vectors:

$$\pi_t \in [0, 1]^n, \quad \sum_{i=1}^n \{\pi_t\}_i = 1 \quad \pi_t = \begin{bmatrix} P(x_t = s_1) \\ \vdots \\ P(x_t = s_n) \end{bmatrix}$$

- assume Markov process:

$$\pi_{t+1} = \Lambda \pi_t, \quad \Lambda \in [0, 1]^{n,n}, \quad \forall j : \sum_{i=1}^n \{\Lambda\}_{i,j} = 1$$

$$\Lambda = \begin{bmatrix} P(x_{t+1} = s_1 \mid x_t = s_1) & \dots & P(x_{t+1} = s_1 \mid x_t = s_n) \\ \vdots & \ddots & \vdots \\ P(x_{t+1} = s_n \mid x_t = s_1) & \dots & P(x_{t+1} = s_n \mid x_t = s_n) \end{bmatrix}$$

Example: Kullback-Leibler estimation of Markov transitional matrix

- there is a “noise” in the data, therefore:

$$\pi_{t+1} \approx \Lambda \pi_t$$

In [mathematical statistics](#), the **Kullback–Leibler divergence** (also called **relative entropy**) is a measure of how one [probability distribution](#) is different from a second, reference probability distribution.^{[1][2]}

For [discrete probability distributions](#) P and Q defined on the same [probability space](#), the Kullback–Leibler divergence between P and Q is defined^[4] to be

$$D_{\text{KL}}(P \parallel Q) = - \sum_{x \in \mathcal{X}} P(x) \log \left(\frac{Q(x)}{P(x)} \right)$$

which is equivalent to

$$D_{\text{KL}}(P \parallel Q) = \sum_{x \in \mathcal{X}} P(x) \log \left(\frac{P(x)}{Q(x)} \right).$$

$$\begin{aligned} \Lambda^* &= \arg \min \sum_{t=1}^{T-1} D_{KL}(\pi_{t+1} \mid \Lambda \pi_t) \\ &= \arg \min \sum_{t=1}^{T-1} \sum_{i=1}^n \{\pi_{t+1}\}_i \log \frac{\{\pi_{t+1}\}_i}{\{\Lambda \pi_t\}_i} \quad \dots \end{aligned}$$

Example: Kullback-Leibler estimation of Markov transitional matrix

$$\begin{aligned}\Lambda^* &= \arg \min \sum_{t=1}^{T-1} D_{KL}(\pi_{t+1} \mid \Lambda \pi_t) \\&= \arg \min \sum_{t=1}^{T-1} \sum_{i=1}^n \{\pi_{t+1}\}_i \log \frac{\{\pi_{t+1}\}_i}{\{\Lambda \pi_t\}_i} \\&= \arg \min - \sum_{t=1}^{T-1} \sum_{i=1}^n \{\pi_{t+1}\}_i \log \{\Lambda \pi_t\}_i \\&= \arg \min - \sum_{t=1}^{T-1} \sum_{i=1}^n \{\pi_{t+1}\}_i \log \sum_{j=1}^n \{\Lambda\}_{i,j} \{\pi_t\}_j \\&\leq \arg \min - \underbrace{\sum_{t=1}^{T-1} \sum_{i=1}^n \sum_{j=1}^n \{\pi_{t+1}\}_i \{\pi_t\}_j \log \{\Lambda\}_{i,j}}_{:=f(\Lambda), \quad f: \mathbb{R}^{n,n} \rightarrow \mathbb{R}}\end{aligned}$$

Example: Kullback-Leibler estimation of Markov transitional matrix

$$\begin{aligned}
 \Lambda^* &= \arg \min \sum_{t=1}^{T-1} D_{KL}(\pi_{t+1} \mid \Lambda \pi_t) \\
 &= \arg \min \sum_{t=1}^{T-1} \sum_{i=1}^n \{\pi_{t+1}\}_i \log \frac{\{\pi_{t+1}\}_i}{\{\Lambda \pi_t\}_i} \\
 &= \arg \min - \sum_{t=1}^{T-1} \sum_{i=1}^n \{\pi_{t+1}\}_i \log \{\Lambda \pi_t\}_i \\
 &= \arg \min - \sum_{t=1}^{T-1} \sum_{i=1}^n \{\pi_{t+1}\}_i \log \sum_{j=1}^n \{\Lambda\}_{i,j} \{\pi_t\}_j \\
 &\leq \arg \min - \underbrace{\sum_{t=1}^{T-1} \sum_{i=1}^n \sum_{j=1}^n \{\pi_{t+1}\}_i \{\pi_t\}_j \log \{\Lambda\}_{i,j}}_{:=f(\Lambda), \quad f: \mathbb{R}^{n,n} \rightarrow \mathbb{R}}
 \end{aligned}$$

$$L(\Lambda, \lambda) := f(\Lambda) + \sum_{j=1}^n \lambda_j \left(\sum_{i=1}^n \{\Lambda\}_{i,j} - 1 \right)$$

$$\frac{\partial L}{\partial \{\Lambda\}_{\hat{i}, \hat{j}}} = - \frac{\sum_{t=1}^{T-1} \{\pi_{t+1}\}_{\hat{i}} \{\pi_t\}_{\hat{j}}}{\{\Lambda\}_{\hat{i}, \hat{j}}} + \lambda_{\hat{j}} = 0$$

$$\frac{\partial L}{\partial \{\lambda\}_{\hat{j}}} = \sum_{i=1}^n \{\Lambda\}_{i, \hat{j}} - 1 = 0$$

Example: Kullback-Leibler estimation of Markov transitional matrix

$$\begin{aligned}
 \Lambda^* &= \arg \min \sum_{t=1}^{T-1} D_{KL}(\pi_{t+1} \mid \Lambda \pi_t) \\
 &= \arg \min \sum_{t=1}^{T-1} \sum_{i=1}^n \{\pi_{t+1}\}_i \log \frac{\{\pi_{t+1}\}_i}{\{\Lambda \pi_t\}_i} \\
 &= \arg \min - \sum_{t=1}^{T-1} \sum_{i=1}^n \{\pi_{t+1}\}_i \log \{\Lambda \pi_t\}_i \\
 &= \arg \min - \sum_{t=1}^{T-1} \sum_{i=1}^n \{\pi_{t+1}\}_i \log \sum_{j=1}^n \{\Lambda\}_{i,j} \{\pi_t\}_j \\
 &\leq \arg \min - \underbrace{\sum_{t=1}^{T-1} \sum_{i=1}^n \sum_{j=1}^n \{\pi_{t+1}\}_i \{\pi_t\}_j \log \{\Lambda\}_{i,j}}_{:=f(\Lambda), \quad f: \mathbb{R}^{n,n} \rightarrow \mathbb{R}}
 \end{aligned}$$

$$L(\Lambda, \lambda) := f(\Lambda) + \sum_{j=1}^n \lambda_j \left(\sum_{i=1}^n \{\Lambda\}_{i,j} - 1 \right)$$

$$\begin{aligned}
 \frac{\partial L}{\partial \{\Lambda\}_{\hat{i}, \hat{j}}} &= - \frac{\boxed{\sum_{t=1}^{T-1} \{\pi_{t+1}\}_{\hat{i}} \{\pi_t\}_{\hat{j}}}}{\{\Lambda\}_{\hat{i}, \hat{j}}} + \lambda_{\hat{j}} = 0 \\
 \frac{\partial L}{\partial \{\lambda\}_{\hat{j}}} &= \sum_{i=1}^n \{\Lambda\}_{i, \hat{j}} - 1 = 0
 \end{aligned}$$

$\{\Lambda\}_{\hat{i}, \hat{j}} = \frac{c_{\hat{i}, \hat{j}}}{\lambda_{\hat{j}}}$

Example: Kullback-Leibler estimation of Markov transitional matrix

$$\begin{aligned}
 \Lambda^* &= \arg \min \sum_{t=1}^{T-1} D_{KL}(\pi_{t+1} \mid \Lambda \pi_t) \\
 &= \arg \min \sum_{t=1}^{T-1} \sum_{i=1}^n \{\pi_{t+1}\}_i \log \frac{\{\pi_{t+1}\}_i}{\{\Lambda \pi_t\}_i} \\
 &= \arg \min - \sum_{t=1}^{T-1} \sum_{i=1}^n \{\pi_{t+1}\}_i \log \{\Lambda \pi_t\}_i \\
 &= \arg \min - \sum_{t=1}^{T-1} \sum_{i=1}^n \{\pi_{t+1}\}_i \log \sum_{j=1}^n \{\Lambda\}_{i,j} \{\pi_t\}_j \\
 &\leq \arg \min - \underbrace{\sum_{t=1}^{T-1} \sum_{i=1}^n \sum_{j=1}^n \{\pi_{t+1}\}_i \{\pi_t\}_j \log \{\Lambda\}_{i,j}}_{:=f(\Lambda), \quad f: \mathbb{R}^{n,n} \rightarrow \mathbb{R}}
 \end{aligned}$$

$$L(\Lambda, \lambda) := f(\Lambda) + \sum_{j=1}^n \lambda_j \left(\sum_{i=1}^n \{\Lambda\}_{i,j} - 1 \right)$$

$$\frac{\partial L}{\partial \{\Lambda\}_{\hat{i}, \hat{j}}} = - \frac{\boxed{\sum_{t=1}^{T-1} \{\pi_{t+1}\}_{\hat{i}} \{\pi_t\}_{\hat{j}}}}{\{\Lambda\}_{\hat{i}, \hat{j}}} + \lambda_{\hat{j}} = 0$$

$$\frac{\partial L}{\partial \{\lambda\}_{\hat{j}}} = \sum_{i=1}^n \{\Lambda\}_{i, \hat{j}} - 1 = 0$$

$$\sum_{i=1}^n \frac{c_{i, \hat{j}}}{\lambda_{\hat{j}}} = 1$$

$$\{\Lambda\}_{\hat{i}, \hat{j}} = \frac{c_{\hat{i}, \hat{j}}}{\lambda_{\hat{j}}}$$

Example: Kullback-Leibler estimation of Markov transitional matrix

$$\Lambda^* = \arg \min \sum_{t=1}^{T-1} D_{KL}(\pi_{t+1} \mid \Lambda \pi_t)$$

$$= \arg \min \sum_{t=1}^{T-1} \sum_{i=1}^n \{\pi_{t+1}\}_i \log \frac{\{\pi_{t+1}\}_i}{\{\Lambda \pi_t\}_i}$$

$$= \arg \min - \sum_{t=1}^{T-1} \sum_{i=1}^n \{\pi_{t+1}\}_i \log \{\Lambda \pi_t\}_i$$

$$\lambda_{\hat{j}} = \sum_{i=1}^n c_{i,\hat{j}} = \sum_{t=1}^{T-1} \underbrace{\sum_{i=1}^n \{\pi_{t+1}\}_i \{\pi_t\}_{\hat{j}}}_{=1} = \sum_{t=1}^{T-1} \{\pi_t\}_{\hat{j}}$$

$$= \arg \min - \sum_{t=1}^{T-1} \sum_{i=1}^n \{\pi_{t+1}\}_i \log \sum_{j=1}^n \{\Lambda\}_{i,j} \{\pi_t\}_j$$

$$\leq \arg \min - \underbrace{\sum_{t=1}^{T-1} \sum_{i=1}^n \sum_{j=1}^n \{\pi_{t+1}\}_i \{\pi_t\}_j \log \{\Lambda\}_{i,j}}_{:=f(\Lambda), \quad f: \mathbb{R}^{n,n} \rightarrow \mathbb{R}}$$

$$\sum_{i=1}^n \frac{c_{i,\hat{j}}}{\lambda_{\hat{j}}} = 1$$

$$L(\Lambda, \lambda) := f(\Lambda) + \sum_{j=1}^n \lambda_j \left(\sum_{i=1}^n \{\Lambda\}_{i,j} - 1 \right)$$

$$\frac{\partial L}{\partial \{\Lambda\}_{\hat{i},\hat{j}}} = - \frac{\sum_{t=1}^{T-1} \{\pi_{t+1}\}_{\hat{i}} \{\pi_t\}_{\hat{j}}}{\{\Lambda\}_{\hat{i},\hat{j}}} + \lambda_{\hat{j}} = 0$$

$$\{\Lambda\}_{\hat{i},\hat{j}} = \frac{c_{\hat{i},\hat{j}}}{\lambda_{\hat{j}}}$$

$$\frac{\partial L}{\partial \{\lambda\}_{\hat{j}}} = \sum_{i=1}^n \{\Lambda\}_{i,\hat{j}} - 1 = 0$$

Example: Kullback-Leibler estimation of Markov transitional matrix

$$\Lambda^* = \arg \min \sum_{t=1}^{T-1} D_{KL}(\pi_{t+1} \mid \Lambda \pi_t)$$

$$= \arg \min \sum_{t=1}^{T-1} \sum_{i=1}^n \{\pi_{t+1}\}_i \log \frac{\{\pi_{t+1}\}_i}{\{\Lambda \pi_t\}_i}$$

$$= \arg \min - \sum_{t=1}^{T-1} \sum_{i=1}^n \{\pi_{t+1}\}_i \log \{\Lambda \pi_t\}_i$$

$$= \arg \min - \sum_{t=1}^{T-1} \sum_{i=1}^n \{\pi_{t+1}\}_i \log \sum_{j=1}^n \{\Lambda\}_{i,j} \{\pi_t\}_j$$

$$\leq \arg \min - \underbrace{\sum_{t=1}^{T-1} \sum_{i=1}^n \sum_{j=1}^n \{\pi_{t+1}\}_i \{\pi_t\}_j \log \{\Lambda\}_{i,j}}_{:=f(\Lambda), \quad f: \mathbb{R}^{n,n} \rightarrow \mathbb{R}}$$

$$L(\Lambda, \lambda) := f(\Lambda) + \sum_{j=1}^n \lambda_j \left(\sum_{i=1}^n \{\Lambda\}_{i,j} - 1 \right)$$

$$\frac{\partial L}{\partial \{\Lambda\}_{i,\hat{j}}} = - \frac{\sum_{t=1}^{T-1} \{\pi_{t+1}\}_{\hat{i}} \{\pi_t\}_{\hat{j}}}{\{\Lambda\}_{i,\hat{j}}} + \lambda_{\hat{j}} = 0$$

$$\frac{\partial L}{\partial \{\lambda\}_{\hat{j}}} = \sum_{i=1}^n \{\Lambda\}_{i,\hat{j}} - 1 = 0$$

$$\lambda_{\hat{j}} = \sum_{i=1}^n c_{i,\hat{j}} = \sum_{t=1}^{T-1} \underbrace{\sum_{i=1}^n \{\pi_{t+1}\}_i \{\pi_t\}_{\hat{j}}}_{=1} = \sum_{t=1}^{T-1} \{\pi_t\}_{\hat{j}}$$

$$\sum_{i=1}^n \frac{c_{i,\hat{j}}}{\lambda_{\hat{j}}} = 1$$

$$\{\Lambda\}_{i,\hat{j}} = \frac{c_{i,\hat{j}}}{\lambda_{\hat{j}}}$$

$$\{\Lambda\}_{i,\hat{j}} = \frac{\sum_{t=1}^{T-1} \{\pi_{t+1}\}_{\hat{i}} \{\pi_t\}_{\hat{j}}}{\sum_{t=1}^{T-1} \{\pi_t\}_{\hat{j}}}$$