


Constrained Optimisation

Yingzhen Li

Department of Computing
Imperial College London

 @liyzhen2
yingzhen.li@imperial.ac.uk

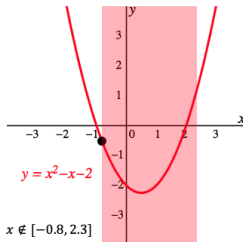
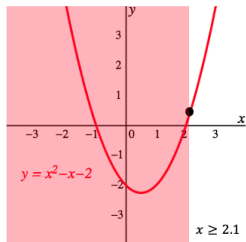
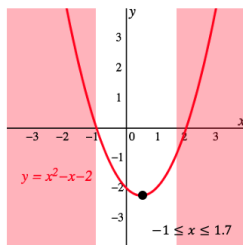
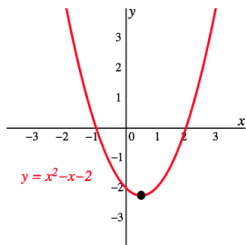
November 19, 2021

Reading for this week and next week

Read MML book: Section 7.2
Extra note will be uploaded to course materials.

From unconstrained to constrained optimisation

Find minimum for function $f(x) = x^2 - x - 2$:



Constrained optimisation: Set-up

A constrained optimisation problem typically has the following form:

$$\min_{\mathbf{x}} L(\mathbf{x})$$

subject to $g_i(\mathbf{x}) \leq 0, i = 1, \dots, N$ (inequality constraints)

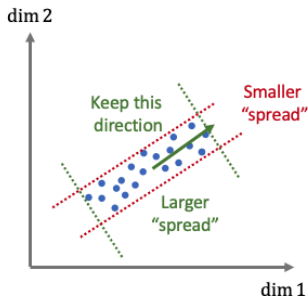
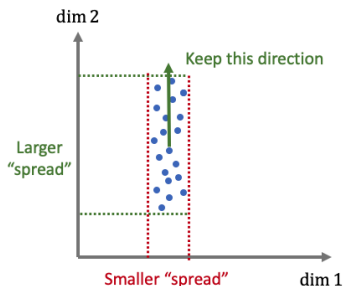
$h_j(\mathbf{x}) = 0, j = 1, \dots, M$ (equality constraints)

Strategies to solve a constrained problem:

- reparameterise \mathbf{x} to directly satisfy the constraints
- **This lecture:** Lagrange multiplier method

PCA: maximum variance perspective

PCA projects data onto directions where the datapoints “vary the most”



“Spread” is defined as the variance along a given direction

PCA: maximum variance perspective

Recall the problem set-up:

- Data: $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^{D \times 1}$ s.t. $\text{mean}(\mathbf{x}_n) = \mathbf{0}$
- Find projections in a **lower-dimensional** space:

$$\mathbf{z}_n := \mathbf{B}^\top \mathbf{x}_n \quad \Leftrightarrow \quad \mathbf{z}_{nj} := \mathbf{b}_j^\top \mathbf{x}_n$$

using an **orthonormal basis**

$$\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_M], \quad \mathbf{b}_m \in \mathbb{R}^{D \times 1}, \quad M < D$$

PCA: maximum variance perspective

Recall the problem set-up:

- ▶ Data: $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^{D \times 1}$ s.t. $\text{mean}(\mathbf{x}_n) = \mathbf{0}$
- ▶ Find projections in a **lower-dimensional** space:

$$\mathbf{z}_n := \mathbf{B}^\top \mathbf{x}_n \quad \Leftrightarrow \quad \mathbf{z}_{nj} := \mathbf{b}_j^\top \mathbf{x}_n$$

using an **orthonormal basis**

$$\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_M], \quad \mathbf{b}_m \in \mathbb{R}^{D \times 1}, \quad M < D$$

- ▶ Solve for \mathbf{B} such that

$$\sum_{m=1}^M \mathbb{V}[\mathbf{b}_m^\top \mathbf{x}_n] \quad \text{is maximised}$$

PCA: maximum variance perspective

Solve for the following constrained optimisation problem:

$$\max_{\{\mathbf{b}_1, \dots, \mathbf{b}_M\}} \sum_{m=1}^M \mathbb{V}[\mathbf{b}_m^\top \mathbf{x}_n]$$

subject to $\|\mathbf{b}_m\|_2^2 = 1, m = 1, \dots, M$
 $\mathbf{b}_i \perp \mathbf{b}_j, \quad \forall i \neq j$

PCA: maximum variance perspective

Solve for the following constrained optimisation problem:

$$\max_{\{\mathbf{b}_1, \dots, \mathbf{b}_M\}} \sum_{m=1}^M \mathbb{V}[\mathbf{b}_m^\top \mathbf{x}_n]$$

$$\text{subject to } \|\mathbf{b}_m\|_2^2 = 1, \quad m = 1, \dots, M$$
$$\mathbf{b}_i \perp \mathbf{b}_j, \quad \forall i \neq j$$

- ▶ Expanding the objective (see previous lecture on PCA)

$$\mathbb{V}[\mathbf{b}_m^\top \mathbf{x}_n] = \mathbf{b}_m^\top \mathbf{S} \mathbf{b}_m, \quad \mathbf{S} := \frac{1}{N} \mathbf{X} \mathbf{X}^\top$$

- ▶ Rewrite the maximisation objective with a **minimisation** objective

PCA: maximum variance perspective

Solve for the following constrained optimisation problem:

$$\min_{\{\mathbf{b}_1, \dots, \mathbf{b}_M\}} - \sum_{m=1}^M \mathbf{b}_m^\top \mathbf{S} \mathbf{b}_m$$

subject to $\|\mathbf{b}_m\|_2^2 = 1, m = 1, \dots, M$
 $\mathbf{b}_i \perp \mathbf{b}_j, \quad \forall i \neq j$

PCA: maximum variance perspective

Solve for the following constrained optimisation problem:

$$\min_{\{\mathbf{b}_1, \dots, \mathbf{b}_M\}} - \sum_{m=1}^M \mathbf{b}_m^\top \mathbf{S} \mathbf{b}_m$$

subject to $\|\mathbf{b}_m\|_2^2 = 1, m = 1, \dots, M$
 $\mathbf{b}_i \perp \mathbf{b}_j, \quad \forall i \neq j$

Let's solve it by the Lagrange multiplier method.

Lagrangian with Lagrange multipliers $\{\lambda_m\}, \{\gamma_{ij}\}$:

$$L(\{\mathbf{b}_m\}, \{\lambda_m\}, \{\gamma_{ij}\}) = - \sum_{m=1}^M \mathbf{b}_m^\top \mathbf{S} \mathbf{b}_m + \sum_{m=1}^M \lambda_m (\|\mathbf{b}_m\|_2^2 - 1) + \sum_{i < j} \gamma_{ij} \mathbf{b}_i^\top \mathbf{b}_j$$

PCA: maximum variance perspective

Lagrangian:

$$L(\{\mathbf{b}_m\}, \{\lambda_m\}, \{\gamma_{ij}\}) = - \sum_{m=1}^M \mathbf{b}_m^T \mathbf{S} \mathbf{b}_m + \sum_{m=1}^M \lambda_m (\|\mathbf{b}_m\|_2^2 - 1) + \sum_{i < j} \gamma_{ij} \mathbf{b}_i^T \mathbf{b}_j$$

Solutions of $\{\mathbf{b}_1, \dots, \mathbf{b}_M\}$ and the Lagrange multipliers need to satisfy:

- ▶ **stationarity:** $\nabla_{\mathbf{b}_m} L(\{\mathbf{b}_m\}, \{\lambda_m\}, \{\gamma_{ij}\}) = \mathbf{0}, m = 1, \dots, M$
- ▶ **primal feasibility:** $\|\mathbf{b}_m\|_2^2 = 1, \mathbf{b}_i \perp \mathbf{b}_j$
(it also means $\nabla_{\lambda_m} L = 0$ and $\nabla_{\gamma_{ij}} L = 0$)

PCA: maximum variance perspective

Lagrangian:

$$L(\{\mathbf{b}_m\}, \{\lambda_m\}, \{\gamma_{ij}\}) = - \sum_{m=1}^M \mathbf{b}_m^\top \mathbf{S} \mathbf{b}_m + \sum_{m=1}^M \lambda_m (\|\mathbf{b}_m\|_2^2 - 1) + \sum_{i < j} \gamma_{ij} \mathbf{b}_i^\top \mathbf{b}_j$$

Solution of the stationarity condition: make

$$\nabla_{\mathbf{b}_m} L(\{\mathbf{b}_m\}, \{\lambda_m\}, \{\gamma_{ij}\}) = \mathbf{0}, \quad m = 1, \dots, M$$

$$\nabla_{\mathbf{b}_m} L(\{\mathbf{b}_m\}, \{\lambda_m\}, \{\gamma_{ij}\}) = -2\mathbf{S}\mathbf{b}_m + 2\lambda_m \mathbf{b}_m + \sum_{i < m} \gamma_{im} \mathbf{b}_i + \sum_{j > m} \gamma_{mj} \mathbf{b}_j = \mathbf{0}$$

PCA: maximum variance perspective

Lagrangian:

$$L(\{\mathbf{b}_m\}, \{\lambda_m\}, \{\gamma_{ij}\}) = - \sum_{m=1}^M \mathbf{b}_m^\top \mathbf{S} \mathbf{b}_m + \sum_{m=1}^M \lambda_m (\|\mathbf{b}_m\|_2^2 - 1) + \sum_{i < j} \gamma_{ij} \mathbf{b}_i^\top \mathbf{b}_j$$

Solution of the stationarity condition: make

$$\nabla_{\mathbf{b}_m} L(\{\mathbf{b}_m\}, \{\lambda_m\}, \{\gamma_{ij}\}) = \mathbf{0}, \quad m = 1, \dots, M$$

$$\begin{aligned} \nabla_{\mathbf{b}_m} L(\{\mathbf{b}_m\}, \{\lambda_m\}, \{\gamma_{ij}\}) &= -2\mathbf{S}\mathbf{b}_m + 2\lambda_m \mathbf{b}_m + \sum_{i < m} \gamma_{im} \mathbf{b}_i + \sum_{j > m} \gamma_{mj} \mathbf{b}_j = \mathbf{0} \\ \Rightarrow \quad \mathbf{b}_m &= (-2\mathbf{S} + 2\lambda_m \mathbf{I})^{-1} \left(\sum_{i < m} \gamma_{im} \mathbf{b}_i + \sum_{j > m} \gamma_{mj} \mathbf{b}_j \right) \end{aligned}$$

Looks difficult to solve...

PCA: maximum variance perspective

Gradient of the Lagrangian w.r.t. \mathbf{b}_m :

$$\nabla_{\mathbf{b}_m} L(\{\mathbf{b}_m\}, \{\lambda_m\}, \{\gamma_{ij}\}) = -2\mathbf{S}\mathbf{b}_m + 2\lambda_m\mathbf{b}_m + \sum_{i < m} \gamma_{im}\mathbf{b}_i + \sum_{j > m} \gamma_{mj}\mathbf{b}_j = \mathbf{0}$$

Notice the **primal feasibility** condition:

$$\mathbf{b}_m^\top \mathbf{b}_j = 0, \forall j \neq m$$

PCA: maximum variance perspective

Gradient of the Lagrangian w.r.t. \mathbf{b}_m :

$$\nabla_{\mathbf{b}_m} L(\{\mathbf{b}_m\}, \{\lambda_m\}, \{\gamma_{ij}\}) = -2\mathbf{S}\mathbf{b}_m + 2\lambda_m\mathbf{b}_m + \sum_{i<m} \gamma_{im}\mathbf{b}_i + \sum_{j>m} \gamma_{mj}\mathbf{b}_j = \mathbf{0}$$

Notice the **primal feasibility** condition:

$$\mathbf{b}_m^\top \mathbf{b}_j = 0, \forall j \neq m$$

Left-multiply the gradient $\nabla_{\mathbf{b}_m} L(\{\mathbf{b}_m\}, \{\lambda_m\}, \{\gamma_{ij}\})$ with \mathbf{b}_m^\top :

$$\begin{aligned} \mathbf{b}_m^\top \nabla_{\mathbf{b}_m} L(\{\mathbf{b}_m\}, \{\lambda_m\}, \{\gamma_{ij}\}) &= -2\mathbf{b}_m^\top \mathbf{S}\mathbf{b}_m + 2\lambda_m \mathbf{b}_m^\top \mathbf{b}_m \\ &\quad + \underbrace{\sum_{i<m} \gamma_{im} \mathbf{b}_m^\top \mathbf{b}_i + \sum_{j>m} \gamma_{mj} \mathbf{b}_m^\top \mathbf{b}_j}_{=0 \text{ (required by primal feasibility)}} \end{aligned}$$

PCA: maximum variance perspective

Gradient of the Lagrangian w.r.t. \mathbf{b}_m :

$$\nabla_{\mathbf{b}_m} L(\{\mathbf{b}_m\}, \{\lambda_m\}, \{\gamma_{ij}\}) = -2\mathbf{S}\mathbf{b}_m + 2\lambda_m\mathbf{b}_m + \sum_{i < m} \gamma_{im}\mathbf{b}_i + \sum_{j > m} \gamma_{mj}\mathbf{b}_j = \mathbf{0}$$

Notice the **primal feasibility** condition:

$$\mathbf{b}_m^\top \mathbf{b}_j = 0, \forall j \neq m$$

Left-multiply the gradient $\nabla_{\mathbf{b}_m} L(\{\mathbf{b}_m\}, \{\lambda_m\}, \{\gamma_{ij}\})$ with \mathbf{b}_m^\top :

$$\begin{aligned} \mathbf{b}_m^\top \nabla_{\mathbf{b}_m} L(\{\mathbf{b}_m\}, \{\lambda_m\}, \{\gamma_{ij}\}) &= -2\mathbf{b}_m^\top \mathbf{S}\mathbf{b}_m + 2\lambda_m \mathbf{b}_m^\top \mathbf{b}_m \\ &= \underbrace{-2\mathbf{b}_m^\top (\mathbf{S}\mathbf{b}_m - \lambda_m \mathbf{b}_m)}_{=0 \text{ (stationarity condition)}} \end{aligned}$$

PCA: maximum variance perspective

Gradient of the Lagrangian w.r.t. \mathbf{b}_m :

$$\nabla_{\mathbf{b}_m} L(\{\mathbf{b}_m\}, \{\lambda_m\}, \{\gamma_{ij}\}) = -2\mathbf{S}\mathbf{b}_m + 2\lambda_m\mathbf{b}_m + \sum_{i < m} \gamma_{im}\mathbf{b}_i + \sum_{j > m} \gamma_{mj}\mathbf{b}_j = \mathbf{0}$$

Notice the **primal feasibility** condition:

$$\mathbf{b}_m^\top \mathbf{b}_j = 0, \forall j \neq m$$

Left-multiply the gradient $\nabla_{\mathbf{b}_m} L(\{\mathbf{b}_m\}, \{\lambda_m\}, \{\gamma_{ij}\})$ with \mathbf{b}_m^\top :

$$\begin{aligned} \mathbf{b}_m^\top \nabla_{\mathbf{b}_m} L(\{\mathbf{b}_m\}, \{\lambda_m\}, \{\gamma_{ij}\}) &= -2\mathbf{b}_m^\top \mathbf{S}\mathbf{b}_m + 2\lambda_m \mathbf{b}_m^\top \mathbf{b}_m \\ &= \underbrace{-2\mathbf{b}_m^\top (\mathbf{S}\mathbf{b}_m - \lambda_m \mathbf{b}_m)}_{=0 \text{ (stationarity condition)}} \end{aligned}$$

- ▶ $\|\mathbf{b}_m\|_2^2 = 1$ (primal feasibility) \Rightarrow make $\mathbf{S}\mathbf{b}_m - \lambda_m \mathbf{b}_m = \mathbf{0}$
 $\Rightarrow \mathbf{b}_m$ is an eigenvector of \mathbf{S}

PCA: maximum variance perspective

Denote $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top$ with $\mathbf{Q} = [\mathbf{q}_1, \dots, \mathbf{q}_D]$.

Notice that for

$$\{(\mathbf{b}_1, \lambda_1), \dots, (\mathbf{b}_M, \lambda_M)\} \subset \{(\mathbf{q}_1, \Lambda_{11}), \dots, (\mathbf{q}_D, \Lambda_{DD})\}, \gamma_{ij} = 0$$

both conditions are satisfied:

- ▶ **stationarity**: $\nabla_{\mathbf{b}_m} L(\{\mathbf{b}_m\}, \{\lambda_m\}, \{\gamma_{ij}\}) = -2(\mathbf{S}\mathbf{b}_m - \lambda_m\mathbf{b}_m) = \mathbf{0}$
- ▶ **primal feasibility**: $\|\mathbf{b}_m\|_2^2 = 1, \mathbf{b}_i \perp \mathbf{b}_j$

PCA: maximum variance perspective

Denote $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top$ with $\mathbf{Q} = [\mathbf{q}_1, \dots, \mathbf{q}_D]$.

Notice that for

$$\{(\mathbf{b}_1, \lambda_1), \dots, (\mathbf{b}_M, \lambda_M)\} \subset \{(\mathbf{q}_1, \Lambda_{11}), \dots, (\mathbf{q}_D, \Lambda_{DD})\}, \gamma_{ij} = 0$$

both conditions are satisfied:

- ▶ **stationarity**: $\nabla_{\mathbf{b}_m} L(\{\mathbf{b}_m\}, \{\lambda_m\}, \{\gamma_{ij}\}) = -2(\mathbf{S}\mathbf{b}_m - \lambda_m\mathbf{b}_m) = \mathbf{0}$
- ▶ **primal feasibility**: $\|\mathbf{b}_m\|_2^2 = 1, \mathbf{b}_i \perp \mathbf{b}_j$

\Rightarrow Solutions from the Lagrangian provide **local minima** only.

PCA: maximum variance perspective

Denote $\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top$ with $\mathbf{Q} = [\mathbf{q}_1, \dots, \mathbf{q}_D]$.

Notice that for

$$\{(\mathbf{b}_1, \lambda_1), \dots, (\mathbf{b}_M, \lambda_M)\} \subset \{(\mathbf{q}_1, \Lambda_{11}), \dots, (\mathbf{q}_D, \Lambda_{DD})\}, \gamma_{ij} = 0$$

both conditions are satisfied:

- ▶ **stationarity**: $\nabla_{\mathbf{b}_m} L(\{\mathbf{b}_m\}, \{\lambda_m\}, \{\gamma_{ij}\}) = -2(\mathbf{S}\mathbf{b}_m - \lambda_m\mathbf{b}_m) = \mathbf{0}$
- ▶ **primal feasibility**: $\|\mathbf{b}_m\|_2^2 = 1, \mathbf{b}_i \perp \mathbf{b}_j$

\Rightarrow Solutions from the Lagrangian provide **local minima** only.

Get the **global minimum** by choosing the best from all local optima:

$$\min \sum_{m=1}^M -\mathbf{b}_m^\top \mathbf{S} \mathbf{b}_m \quad \text{for } \{\mathbf{b}_1, \dots, \mathbf{b}_M\} \subset \{\mathbf{q}_1, \dots, \mathbf{q}_D\}$$

$\Rightarrow \{\mathbf{b}_1, \dots, \mathbf{b}_M\}$ is the leading M eigenvectors of \mathbf{S}

Constrained optimisation: equality constraints

A constrained optimisation problem with equality constraints:

$$\min_{\mathbf{x}} L(\mathbf{x})$$

subject to $h_j(\mathbf{x}) = 0, j = 1, \dots, M$ (equality constraints)

The Lagrange multiplier method:

Constrained optimisation: equality constraints

A constrained optimisation problem with equality constraints:

$$\min_{\mathbf{x}} L(\mathbf{x})$$

subject to $h_j(\mathbf{x}) = 0, j = 1, \dots, M$ (equality constraints)

The Lagrange multiplier method:

- Construct the Lagrangian function

$$L(\mathbf{x}, \{\lambda_j\}) = L(\mathbf{x}) + \sum_{j=1}^M \lambda_j h_j(\mathbf{x})$$

Constrained optimisation: equality constraints

A constrained optimisation problem with equality constraints:

$$\min_{\mathbf{x}} L(\mathbf{x})$$

subject to $h_j(\mathbf{x}) = 0, j = 1, \dots, M$ (equality constraints)

The Lagrange multiplier method:

- ▶ Construct the Lagrangian function

$$L(\mathbf{x}, \{\lambda_j\}) = L(\mathbf{x}) + \sum_{j=1}^M \lambda_j h_j(\mathbf{x})$$

- ▶ Find stationary points for the Lagrangian: solve for

$$\underbrace{\nabla_{\mathbf{x}} L(\mathbf{x}, \{\lambda_j\}) = \mathbf{0}}_{\text{stationarity}} \text{ and } \underbrace{\nabla_{\lambda_j} L(\mathbf{x}, \{\lambda_j\}) = 0, j = 1, \dots, M}_{\text{primal feasibility}}$$

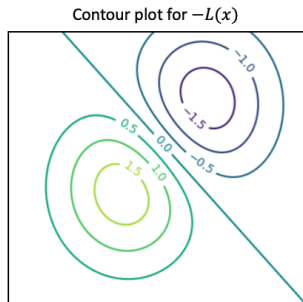
Intuition for the Lagrange multiplier method

A simplified set-up: $M = 1$, $\lambda := \lambda_1$, $h(\mathbf{x}) := h_1(\mathbf{x})$, $\mathbf{x} \in \mathbb{R}^D$

$$L(\mathbf{x}, \lambda) = L(\mathbf{x}) + \lambda h(\mathbf{x})$$

Intuition for solving $\nabla_{\mathbf{x}, \lambda} L = \mathbf{0}$:

- ▶ A “surface” (i.e. low-dim manifold) in \mathbb{R}^D can be represented as $f(\mathbf{x}) = c$
 - ▶ **Contour plot:** plotting different manifolds $f(\mathbf{x}) = c_i$ with different c_i
- ▶ ...so we can plot the **contour of $-L(\mathbf{x})$**



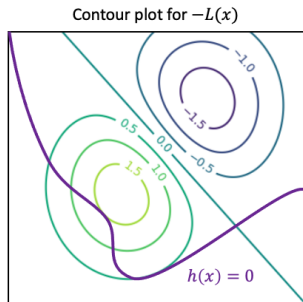
Intuition for the Lagrange multiplier method

A simplified set-up: $M = 1$, $\lambda := \lambda_1$, $h(\mathbf{x}) := h_1(\mathbf{x})$, $\mathbf{x} \in \mathbb{R}^D$

$$L(\mathbf{x}, \lambda) = L(\mathbf{x}) + \lambda h(\mathbf{x})$$

Intuition for solving $\nabla_{\mathbf{x}, \lambda} L = \mathbf{0}$:

- ▶ A “surface” (i.e. low-dim manifold) in \mathbb{R}^D can be represented as $f(\mathbf{x}) = c$
 - ▶ **Constraint:** $h(\mathbf{x}) = 0$ represents a low-dim manifold in \mathbb{R}^D
- ▶ ...so we can plot the “surface” of $h(\mathbf{x}) = 0$



Intuition for the Lagrange multiplier method

A simplified set-up: $M = 1$, $\lambda := \lambda_1$, $h(\mathbf{x}) := h_1(\mathbf{x})$, $\mathbf{x} \in \mathbb{R}^D$

$$L(\mathbf{x}, \lambda) = L(\mathbf{x}) + \lambda h(\mathbf{x})$$

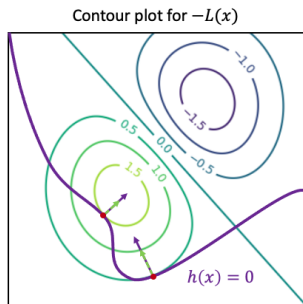
Intuition for solving $\nabla_{\mathbf{x}, \lambda} L = \mathbf{0}$:

- ▶ A “surface” (i.e. low-dim manifold) in \mathbb{R}^D can be represented as $f(\mathbf{x}) = c$
- ▶ Solving $\nabla_{\mathbf{x}, \lambda} L = \mathbf{0}$:

$$\Rightarrow -\nabla_{\mathbf{x}^*} L(\mathbf{x}^*) = \lambda^* \nabla_{\mathbf{x}^*} h(\mathbf{x}^*)$$

$$h(\mathbf{x}^*) = 0$$

\Rightarrow For some c^* , $-L(\mathbf{x}) = c^*$ and $h(\mathbf{x}) = 0$ are **tangent** at point \mathbf{x}^*



$$-\nabla_{\mathbf{x}^*} L(\mathbf{x}^*) = \lambda^* \nabla_{\mathbf{x}^*} h(\mathbf{x}^*)$$

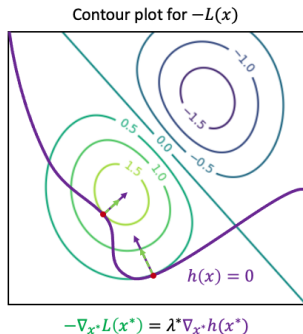
Intuition for the Lagrange multiplier method

A simplified set-up: $M = 1$, $\lambda := \lambda_1$, $h(\mathbf{x}) := h_1(\mathbf{x})$, $\mathbf{x} \in \mathbb{R}^D$

$$L(\mathbf{x}, \lambda) = L(\mathbf{x}) + \lambda h(\mathbf{x})$$

Intuition for solving $\nabla_{\mathbf{x}, \lambda} L = \mathbf{0}$:

- ▶ Why finding **“points of contact”**:
 - ▶ If can find $\delta \mathbf{x}$ s.t. $h(\mathbf{x} + \delta \mathbf{x}) = 0$,
 $\langle -\nabla_{\mathbf{x}} L(\mathbf{x}), \delta \mathbf{x} \rangle > 0$
 $\Rightarrow L(\mathbf{x}) > L(\mathbf{x} + \delta \mathbf{x})$



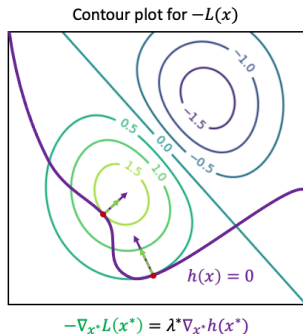
Intuition for the Lagrange multiplier method

A simplified set-up: $M = 1$, $\lambda := \lambda_1$, $h(\mathbf{x}) := h_1(\mathbf{x})$, $\mathbf{x} \in \mathbb{R}^D$

$$L(\mathbf{x}, \lambda) = L(\mathbf{x}) + \lambda h(\mathbf{x})$$

Intuition for solving $\nabla_{\mathbf{x}, \lambda} L = \mathbf{0}$:

- ▶ Why finding **“points of contact”**:
 - ▶ If can find $\delta \mathbf{x}$ s.t. $h(\mathbf{x} + \delta \mathbf{x}) = 0$,
 $\langle -\nabla_{\mathbf{x}} L(\mathbf{x}), \delta \mathbf{x} \rangle > 0$
 $\Rightarrow L(\mathbf{x}) > L(\mathbf{x} + \delta \mathbf{x})$
 - ▶ To make \mathbf{x}^* a local minimum of $L(\mathbf{x})$ on surface $h(\mathbf{x}) = 0$:
make sure $\forall \delta \mathbf{x}$ s.t. $h(\mathbf{x} + \delta \mathbf{x}) = 0$,
 $\langle -\nabla_{\mathbf{x}} L(\mathbf{x}), \delta \mathbf{x} \rangle \leq 0$
 $\Rightarrow \mathbf{x}^*$ is a **“point of contact”**



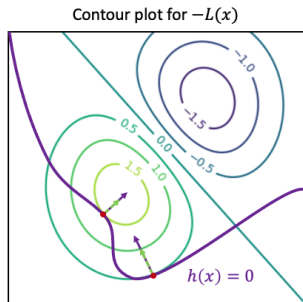
Intuition for the Lagrange multiplier method

A simplified set-up: $M = 1$, $\lambda := \lambda_1$, $h(\mathbf{x}) := h_1(\mathbf{x})$, $\mathbf{x} \in \mathbb{R}^D$

$$L(x, \lambda) = L(x) + \lambda h(x)$$

Intuition for solving $\nabla_{x,\lambda} L = \mathbf{0}$:

- ▶ Why finding **“points of contact”**:
 - ▶ $\nabla_x f(x)$ is a normal vector of surface $f(x) = c$ at x
 - ▶ Find **“points of contact”** between surfaces $-L(x) = c^*$ and $h(x) = 0$:
 \Leftrightarrow find x^* s.t. $-\nabla_{x^*} L(x^*)$ and $\nabla_{x^*} h(x^*)$ align on the same line



$$-\nabla_{x^*} L(x^*) = \lambda^* \nabla_{x^*} h(x^*)$$

Summary

Summary for today's lecture:

- Constrained optimisation with **equality** constraints
 - The Lagrange multiplier method
 - PCA (maximum variance perspective)
- Exercise for today:
 - formulate PCA (minimum reconstruction error perspective) as constrained optimisation
 - solve it with the Lagrange multiplier method

Summary

Summary for today's lecture:

- ▶ Constrained optimisation with **equality** constraints
 - ▶ The Lagrange multiplier method
 - ▶ PCA (maximum variance perspective)
- ▶ Exercise for today:
 - ▶ formulate PCA (minimum reconstruction error perspective) as constrained optimisation
 - ▶ solve it with the Lagrange multiplier method

Next lecture:

- ▶ Constrained optimisation with **inequality** constraints
 - ▶ The KKT conditions
 - ▶ Ridge regression: constrained optimisation perspective