Math for Machine Learning Week 4.1: Optimization and the Lagrangian Method

By: Samuel Deng

Logistics & Announcements

Lesson Overview

Optimization. Minimize an <u>objective function</u> $f : \mathbb{R}^d \to \mathbb{R}$ with the possible requirement that the minimizer \mathbf{x}^* belongs to a constraint set $\mathscr{C} \subseteq \mathbb{R}^d$.

Lagrangian. For optimization problems with \mathscr{C} defined by equalities/inequalities, the <u>Lagrangian</u> is a function $L : \mathbb{R}^d \times \mathbb{R}^m \times \mathbb{R}^r \to \mathbb{R}$ that "unconstrains" the problem.

Unconstrained local optima. With no constraints, the standard tools of calculus give conditions for a point x^* to be optimal, at least to all points close to it.

Constrained local optima (Lagrangian and KKT). When \mathscr{C} is represented by inequalities and equalities, we can use the method of Lagrange multipliers and the KKT Theorem to "unconstrain" the problem.

Ridge regression and minimum norm solutions. By constraining the norm of $\mathbf{w}^* \in \mathbb{R}^d$ of least squares (i.e. $\|\mathbf{w}^*\|$), we obtain more "stable" solutions.

Lesson Overview Big Picture: Least Squares





unconstrained min.
 constrained min.
 C

Lesson Overview

Big Picture: Gradient Descent





Optimization Problems Definition and examples

Motivation **Optimization in calculus**

In much of machine learning, we design algorithms for well-defined optimization problems. In an optimization problem, we want to minimize an <u>objective function</u> $f: \mathbb{R}^d \to \mathbb{R}$ with respect

to a set of constraints $\mathscr{C} \subseteq \mathbb{R}^d$:

minimize $\mathbf{x} \in \mathbb{R}^d$

 $f(\mathbf{X})$

subject to $x \in \mathscr{C}$

Motivation

Components of an optimization problem

 $\mathbf{x} \in \mathbb{R}^d$

 $f: \mathbb{R}^d \to \mathbb{R}$ is the <u>objective function</u>. $\mathscr{C} \subseteq \mathbb{R}^d$ is the <u>constraint/feasible set</u>.

x* is an <u>optimal solution (global minimum)</u> if

The <u>optimal value</u> is $f(\mathbf{x}^*)$. Our goal is to find \mathbf{x}^* and $f(\mathbf{x}^*)$.

- minimize $f(\mathbf{x})$
- subject to $x \in \mathscr{C}$

- $\mathbf{x}^* \in \mathscr{C}$ and $f(\mathbf{x}^*) \leq f(\mathbf{x})$, for all $\mathbf{x} \in \mathscr{C}$.
- Note: to maximize $f(\mathbf{x})$, just minimize $-f(\mathbf{x})$. So we'll only focus on minimization problems.

Motivation

Optimization in single-variable calculus

Ultimate goal: Find the global minimum of functions.

Intermediary goal: Find the local minima.

Now we will focus on constraints!





 \boldsymbol{x}

global mi

Motivation Example: Linear Programming

Let $\mathbf{c} \in \mathbb{R}^d$, $\mathbf{A} \in \mathbb{R}^{n \times d}$, $\mathbf{b} \in \mathbb{R}^n$ be fixed.

Let $\mathbf{x} \in \mathbb{R}^d$ be the <u>decision/free variables</u>.

minimize $\mathbf{c}^{\mathsf{T}}\mathbf{x}$ $\mathbf{x} \in \mathbb{R}^d$ subject to $Ax \leq b$

 \leq is element-wise inequality: $\mathbf{a}_i^{\mathsf{T}}\mathbf{x} \leq b_i$ for all $i \in [n]$.

We're cooking some NYC classics again. Suppose we have:

Bacon egg and cheese (BEC) requires 1 bacon, 1 egg, 1 cheese, and 1 roll.

Cost (including labor): \$3

Egg and cheese (EC) requires 0 bacon, 2 egg, 1 cheese, and 1 roll.

Cost (including labor): \$2

Bacon egg omelette (BEO) requires 1 bacon, 3 egg, 1/2 cheese, and 0 roll.

Cost (including labor): \$1

- 100 bacon, 120 egg, 150 cheese, and 300 (sandwich) rolls.

We're cooking some NYC classics again. Suppose v 100 bacon, 120 egg, 150 cheese, and 300Bacon egg and cheese (BEC) requires 1 bacon, 1 eg Cost (including labor): \$3 Egg and cheese (EC) requires 0 bacon, 2 egg, 1 che Cost (including labor): \$2 Bacon egg omelette (BEO) requires 1 bacon, 3 egg

Cost (including labor): \$1

we have:	Decision variables?
0 (sandwich) rolls.	$\mathbf{x} = (x_1, x_2, x_3) \in \mathbb{R}^3$
	$x_1 = $ number of BEC,
gg, 1 cheese, and 1 roll.	$x_2 = $ number of EC,
	$x_3 = $ number of BEO
	Constraints?
eese, and 1 roll.	Bacon: $\mathbf{a}_1 = (1,0,1)$, $b_1 = 100$
	Egg: $\mathbf{a}_2 = (1,2,3), b_2 = 120$
	Cheese: $\mathbf{a}_3 = (1,1,1/2)$, $b_3 = 150$
g, 1/2 cheese, and 0 roll.	Roll: $\mathbf{a}_4 = (1,1,0), b_4 = 300$
	Objective?

$$\mathbf{c}^{\mathsf{T}}\mathbf{x} = 3x_1 + 2x_2 + x_3$$

Decision variables?

Linear program: $\mathbf{x} = (x_1, x_2, x_3) \in \mathbb{R}^3$ x_1 = number of BEC, minimize x_2 = number of EC, subject to $x_1 + x_3 \le 100$ x_3 = number of BEO **Constraints?** Bacon: $\mathbf{a}_1 = (1,0,1), b_1 = 100$ Egg: $\mathbf{a}_2 = (1,2,3), b_2 = 120$ Cheese: $\mathbf{a}_3 = (1, 1, 1/2), b_3 = 150$

Roll:
$$\mathbf{a}_4 = (1,1,0), b_4 = 300$$

Objective?

$$\mathbf{c}^{\mathsf{T}}\mathbf{x} = 3x_1 + 2x_2 + x_3$$

 $3x_1 + 2x_2 + x_3$ $x_1 + 2x_2 + 3x_3 \le 120$ $x_1 + x_2 + 0.5x_3 \le 150$ $x_1 + x_2 \le 300$ $x_1 \ge 0$ $x_2 \ge 0$

$x_3 \ge 0$

minimize

 $3x_1 + 2x_2 + x_3$ subject to $x_1 + x_3 \le 100$ $x_1 + 2x_2 + 3x_3 \le 120$ $x_1 + x_2 + 0.5x_3 \le 150$ $x_1 + x_2 \le 300$ $x_1 \ge 0$ $x_2 \ge 0$ $x_3 \ge 0$



Regression Setup (Example View)

<u>**Observed:**</u> Matrix of training samples $\mathbf{X} \in \mathbb{R}^{n \times d}$ and vector of training labels $\mathbf{y} \in \mathbb{R}^{n}$.

$$\mathbf{X} = \begin{bmatrix} \leftarrow \mathbf{x}_1^\top \rightarrow \\ \vdots \\ \leftarrow \mathbf{x}_n^\top \rightarrow \end{bmatrix} \mathbf{y}$$

<u>**Unknown:**</u> Weight vector $\mathbf{w} \in \mathbb{R}^d$ with weights w_1, \ldots, w_d .

<u>Goal</u>: For each $i \in [n]$, we predict: $\hat{y}_i = \mathbf{w}^\top \mathbf{x}_i = w_1 x_{i1} + \ldots + w_d x_{id} \in \mathbb{R}$.

Choose a weight vector that "fits the training data": $\mathbf{w} \in \mathbb{R}^d$ such that $y_i \approx \hat{y}_i$ for $i \in [n]$, or:

$$= \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \text{ where } \mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d.$$

 $\mathbf{X}\mathbf{w} = \hat{\mathbf{y}} \approx \mathbf{y}$.

Regression Setup (Feature View)

<u>**Observed:**</u> Matrix of training samples $\mathbf{X} \in \mathbb{R}^{n \times d}$ and vector of training labels $\mathbf{y} \in \mathbb{R}^{n}$.

$$\mathbf{X} = \begin{bmatrix} \uparrow & & \uparrow \\ \mathbf{x}_1 & \dots & \mathbf{x}_d \\ \downarrow & & \downarrow \end{bmatrix} \mathbf{y} = \mathbf{y}$$

<u>**Unknown:**</u> Weight vector $\mathbf{w} \in \mathbb{R}^d$ with weights w_1, \ldots, w_d .

Choose a weight vector that "fits the training data": $\mathbf{w} \in \mathbb{R}^d$ such that $y_i \approx \hat{y}_i$ for $i \in [n]$, or:

$$= \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \text{ where } \mathbf{x}_1, \dots, \mathbf{x}_d \in \mathbb{R}^n.$$

 $\mathbf{X}\mathbf{w} = \hat{\mathbf{y}} \approx \mathbf{y}$.

Least Squares OLS Theorem

<u>Theorem (Ordinary Least Squares).</u> Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ and $\mathbf{y} \in \mathbb{R}^{n}$. Let $\hat{\mathbf{w}} \in \mathbb{R}^{d}$ be the least squares minimizer:

$$\hat{\mathbf{w}} = \arg\min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$$

If $n \ge d$ and $rank(\mathbf{X}) = d$, then:

$$\hat{\mathbf{w}} = (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}.$$

To get predictions $\hat{\mathbf{y}} \in \mathbb{R}^n$:

 $\hat{\mathbf{y}} = \mathbf{X}\hat{\mathbf{w}} = \mathbf{X}(\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}.$



Least Squares OLS Theorem

Proof (Calculus proof of OLS).

 $f(\mathbf{w}) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 \iff f(\mathbf{w}) = \mathbf{w}^\top \mathbf{X}^\top \mathbf{X}\mathbf{w} - 2\mathbf{w}^\top \mathbf{X}^\top \mathbf{y} + \mathbf{y}^\top \mathbf{y}$ "First derivative test." $\nabla_{\mathbf{w}} f(\mathbf{w}) = 2(\mathbf{X}^\top \mathbf{X})\mathbf{w} - 2\mathbf{X}^\top \mathbf{y}$. $2(\mathbf{X}^\top \mathbf{X})\mathbf{w} - 2\mathbf{X}^\top \mathbf{y} = \mathbf{0} \implies \mathbf{X}^\top \mathbf{X}\mathbf{w} = \mathbf{X}^\top \mathbf{y}$ rank $(\mathbf{X}) = d \Longrightarrow$ rank $(\mathbf{X}^\top \mathbf{X}) = d \Longrightarrow \mathbf{X}^\top \mathbf{X}$ is invertible:

$$\hat{\mathbf{w}} = (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}.$$

"Second derivative test." $\nabla^2_{\mathbf{w}} f(\mathbf{w}) = 2\mathbf{X}^{\mathsf{T}}\mathbf{X}$.

 $\operatorname{rank}(\mathbf{X}) = d \implies \operatorname{rank}(\mathbf{X}^{\mathsf{T}}\mathbf{X}) = d \implies \lambda_1, \dots, \lambda_d > 0$

 $\implies \mathbf{X}^{\mathsf{T}}\mathbf{X}$ is positive definite!



Local and global minima Definition of "locality" and different minima

Motivation

Optimization in single-variable calculus

Ultimate goal: Find the global minimum of functions.

Intermediary goal: Find the local minima.





 \boldsymbol{x}

local mi global mi

"Local" to a Point

Definition of an open ball/neighborhood

Let $\mathbf{x} \in \mathbb{R}^d$ be a point. For some real value $\delta > 0$, the <u>open ball</u> or <u>neighborhood of radius</u> δ around **x** is the set of all points:

 $B_{\delta}(\mathbf{x}) := \{ \mathbf{a} \in \mathbb{R}^d : \|\mathbf{x} - \mathbf{a}\| < \delta \}.$

"Local" to a Point

Definition of an open ball/neighborhood

Example. Consider $\mathbf{x} = (1,1) \in \mathbb{R}^2$. What is the open ball of radius $\delta = 1$ around \mathbf{x} ?

"Local" to a Point

Definition of the interior of a set

Let $S \subseteq \mathbb{R}^d$ be a set. A point $\mathbf{x} \in S$ is an <u>interior point</u> if there exists a neighborhood $B_{\delta}(\mathbf{x})$ around **x** such that $B_{\delta}(\mathbf{x}) \subset S$ (where \subset is proper subset).

The interior of the set int(S) is the set of all interior points of S, i.e.

 $int(S) := \{ \mathbf{x} \in S : N_{\delta}(\mathbf{x}) \subset S \}.$

$B_{\delta}(\mathbf{x}) := \{ \mathbf{a} \in \mathbb{R}^d : \|\mathbf{x} - \mathbf{a}\| < \delta \}$

Types of Minima Local and global minima



Types of Minima Local and global minima

minimize $f(\mathbf{x})$ subject to $x \in \mathscr{C}$

 $\hat{\mathbf{x}} \in \mathscr{C}$ is a <u>(constrained) local minimum</u> if there is a neighborhood $B_{\delta}(\hat{\mathbf{x}})$ around $\hat{\mathbf{x}}$ such that

 $f(\hat{\mathbf{x}}) \leq f(\mathbf{x})$ for all $\mathbf{x} \in \mathscr{C} \cap B_{\delta}(\hat{\mathbf{x}})$.

 $\mathbf{x}^* \in \mathscr{C}$ is a <u>global minimum</u> if

 $f(\mathbf{x}^*) \leq f(\mathbf{x})$ for all $\mathbf{x} \in \mathscr{C}$.





Types of Minima Local and global minima

$\begin{array}{ll} \text{minimize} & f(\mathbf{x}) \\ \text{subject to} & \mathbf{x} \in \mathscr{C} \end{array}$

 $\hat{\mathbf{x}} \in \mathscr{C}$ is an <u>unconstrained local minimum</u> if there is a neighborhood $B_{\delta}(\hat{\mathbf{x}}) \subset \mathscr{C}$ around $\hat{\mathbf{x}}$ such that

 $f(\hat{\mathbf{x}}) \leq f(\mathbf{x})$ for all $\mathbf{x} \in B_{\delta}(\hat{\mathbf{x}})$.

Unconstrained local minima are in $int(\mathscr{C})$.

Constrained local minima can be on the "edge" of the constraint set.



Types of Minima Which type of minima are each of these points?

f(x)

 $\begin{array}{ll} \text{minimize} & f(\mathbf{x}) \\ \text{subject to} & \mathbf{x} \in \mathscr{C} \end{array}$

constrained local:

 $f(\hat{\mathbf{x}}) \leq f(\mathbf{x})$ for all $\mathbf{x} \in \mathscr{C} \cap B_{\delta}(\hat{\mathbf{x}})$

<u>unconstrained local:</u>

 $f(\hat{\mathbf{x}}) \leq f(\mathbf{x})$ for all $\mathbf{x} \in B_{\delta}(\hat{\mathbf{x}})$ and $B_{\delta}(\hat{\mathbf{x}}) \subset \mathscr{C}$. global:

 $f(\mathbf{x}^*) \leq f(\mathbf{x})$ for all $\mathbf{x} \in \mathscr{C}$.





Types of Minima Big picture

We want to find **global minima**.

Global minima could be either unconstrained local minima or constrained local minima.

Without *C*, global minima are just an *unconstrained local minima*.

f(x)

With *C*, global minima may lie on the boundary of the constraint set.

Find local minima, then test!





Finding local minima Big Picture

Necessary and sufficient conditions Review

Q is necessary for P. P is sufficient for Q.
sufficiency: If you assume this, you get your property.
A sufficient (not necessary) condition to get an A in this class is to get 100 on every assignment.
necessity: Your property cannot hold unless you assume this.
A necessary (not sufficient) condition to get an A in this class is to turn in every assignment.

$P \implies Q$

Unconstrained Minima How do we find unconstrained minima?

- $\hat{\mathbf{x}} \in \mathscr{C}$ is an <u>unconstrained local minimum</u> if there is a neighborhood $B_{\delta}(\hat{\mathbf{x}}) \subset \mathscr{C}$ around $\hat{\mathbf{x}}$ s.t. $f(\hat{\mathbf{x}}) \leq f(\mathbf{x})$ for all $\mathbf{x} \in B_{\delta}(\hat{\mathbf{x}})$.
- From single-variable calculus, this is true if:
 - $f'(x) = 0 \text{ and } f''(x) \ge 0.$

Unconstrained Minima Intuition from Taylor series

Let $\delta \in \mathbb{R}$ be a scalar increment.

At $x_0 \in \mathbb{R}$, the second-order Taylor approximation tells us all we need to know:

 $f(x_0 + \delta) \approx f(x_0)$

$$f'(x_0)\delta + \frac{1}{2}f''(x_0)\delta^2.$$

$$f'(x) = 0 \qquad f''(x) \ge 0$$

$$f''(x) > 0$$

Single-variable example

$$f(x) = e^{x/2}$$

Second-order Taylor expansion at $x_0 = 1$:

$$T^{2}(x) = e^{1/2} + \frac{e^{1/2}(x-1)}{2} + \frac{e^{1/2}(x-1)^{2}}{8}$$



Unconstrained Minima Intuition from Taylor series

Let $\delta \in \mathbb{R}$ be a scalar increment.

At $x_0 \in \mathbb{R}$, the second-order Taylor approximation tells us all we need to know:

What are the *necessary* conditions for *x* to be a minimum?

What are the *sufficient* conditions for *x* to be a minimum?

$f'(x) = 0 \qquad 1^{f''(x) \ge 0}$ $f(x_0 + \delta) \approx f(x_0) + f'(x_0)\delta + \frac{1}{2}f''(x_0)\delta^2.$ $f'(x) = 0 \qquad \qquad f''(x) \ge 0$

Unconstrained Minima Sufficient conditions met

$$f(x_0 + \delta) \approx f(x_0) + f'(x_0)\delta + \frac{1}{2}f''(x_0)\delta^2$$

Necessary conditions: $f'(x_0) = 0, f''(x_0) \ge 0.$

Sufficient conditions: $f'(x_0) = 0, f''(x_0) > 0.$



Unconstrained Minima Necessary, not sufficient

$$f(x_0 + \delta) \approx f(x_0) + f'(x_0)\delta + \frac{1}{2}f''(x_0)\delta^2$$

Necessary conditions: $f'(x_0) = 0, f''(x_0) \ge 0.$
Sufficient conditions: $f'(x_0) = 0, f''(x_0) > 0.$



$$f(x) = (x - 1)^3 + 1$$
Taylor's Theorem Intuition

How much do we lose by approximating f with a Taylor approximation? **Remainder**: how much more Taylor series is left after "chopping it off" at order n.

First-order approximation:

 $f(\mathbf{x}) \approx f(\mathbf{x}_0)$

The remainder is:

 $f(\mathbf{x}) - (f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^{\mathsf{T}}(\mathbf{x} - \mathbf{x}_0))$

$$+ \nabla f(\mathbf{x}_0)^{\mathsf{T}}(\mathbf{x} - \mathbf{x}_0)$$

Taylor's Theorem Intuition

How much do we lose by approximating *f* with a Taylor approximation? **Remainder**: how much more Taylor series is left after "chopping it off" at order *n*.

Second-order approximation:

$$f(\mathbf{x}) \approx f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^{\mathsf{T}} (\mathbf{x} - \mathbf{x}_0) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_0)^{\mathsf{T}} \nabla^2 f(\mathbf{x}_0) (\mathbf{x} - \mathbf{x}_0).$$

The remainder is:

$$f(\mathbf{x}) - \left(f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^{\mathsf{T}}(\mathbf{x} - \mathbf{x}_0) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^{\mathsf{T}} \nabla^2 f(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0)\right).$$

Remainder of Taylor Polynomial Definition

The <u>remainder</u> of a function and its Taylor polynomial at \mathbf{x}_0 is the function:

- $R^{n}(\mathbf{x}) := f(\mathbf{x}) T^{n}_{\mathbf{x}_{0}}(\mathbf{x})$
- What behavior would we like?
- Ideally, $R^n(\mathbf{x}) \to 0$ as $\mathbf{x} \to \mathbf{x}_0$ (the approximation gets better as we approach \mathbf{x}_0).

Remainder of Taylor Polynomial Definition

The <u>remainder</u> of a function and its Taylor polynomial at \mathbf{x}_0 is the function:

$$R^n(\mathbf{x}) := f(\mathbf{x}) - T^n_{\mathbf{x}_0}(\mathbf{x})$$

What behavior would we like?

Ideally, $R^n(\mathbf{x}) \to 0$ as $\mathbf{x} \to \mathbf{x}_0$ (the approximation gets better as we approach \mathbf{x}_0).



Taylor's Theorem Peano's Form

 \mathbf{x}_0 and let $\mathbf{d} \in \mathbb{R}^d$. For every $\epsilon > 0$, there exists a neighborhood $B_{\delta}(\mathbf{0})$ such that

$$\left| f(\mathbf{x}_0 + \mathbf{d}) - \left(f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^{\mathsf{T}} \mathbf{d} + \frac{1}{2} \mathbf{d}^{\mathsf{T}} \nabla^2 f(\mathbf{x}_0) \mathbf{d} \right) \right| \le \epsilon \|\mathbf{d}\|^2$$

for all $\mathbf{d} \in B_{\delta}(\mathbf{0})$.

However small you want the remainder (ϵ), as long as you are δ -close to \mathbf{x}_0 , the remainder can get $\epsilon \|\mathbf{d}\|^2$ small.

Theorem (2nd Order Taylor's Theorem: Peano's Form). Let $f : \mathbb{R}^d \to \mathbb{R}$ be twice differentiable at

Unconstrained local minima Necessary conditions

Least Squares OLS Theorem

Proof (Calculus proof of OLS).

 $f(\mathbf{w}) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 \iff f(\mathbf{w}) = \mathbf{w}^{\mathsf{T}}\mathbf{X}^{\mathsf{T}}\mathbf{X}\mathbf{w} - 2\mathbf{w}^{\mathsf{T}}\mathbf{X}^{\mathsf{T}}\mathbf{y} + \mathbf{y}^{\mathsf{T}}\mathbf{y}$ "First derivative test." $\nabla_{\mathbf{w}}f(\mathbf{w}) = 2(\mathbf{X}^{\mathsf{T}}\mathbf{X})\mathbf{w} - 2\mathbf{X}^{\mathsf{T}}\mathbf{y}$. $2(\mathbf{X}^{\mathsf{T}}\mathbf{X})\mathbf{w} - 2\mathbf{X}^{\mathsf{T}}\mathbf{y} = \mathbf{0} \implies \mathbf{X}^{\mathsf{T}}\mathbf{X}\mathbf{w} = \mathbf{X}^{\mathsf{T}}\mathbf{y}$

 $\operatorname{rank}(\mathbf{X}) = d \Longrightarrow \operatorname{rank}(\mathbf{X}^{\mathsf{T}}\mathbf{X}) = d \Longrightarrow \mathbf{X}^{\mathsf{T}}\mathbf{X}$ is invertible:

$$\hat{\mathbf{w}} = (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}.$$

"Second derivative test." $\nabla^2_{\mathbf{w}} f(\mathbf{w}) = 2\mathbf{X}^{\mathsf{T}} \mathbf{X}$.

 $\operatorname{rank}(\mathbf{X}) = d \implies \operatorname{rank}(\mathbf{X}^{\mathsf{T}}\mathbf{X}) = d \implies \lambda_1, \dots, \lambda_d > 0$

 $\implies \mathbf{X}^{\mathsf{T}}\mathbf{X}$ is positive definite!



Necessary Conditions Comparison to single variable

 $f(x_0 + \delta) \approx f(x_0) + f'(x_0)\delta + \frac{1}{2}f''(x_0)\delta^2$

when δ is small enough.

Necessary conditions:

$$f'(x_0) = 0, f''(x_0) \ge 0.$$

$$f(\mathbf{x}_0 + \mathbf{d}) \approx f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^{\mathsf{T}} \mathbf{d} + \frac{1}{2} \mathbf{d}^{\mathsf{T}} \nabla^2 f(\mathbf{x}_0) \mathbf{d}^{\mathsf{T}} \mathbf{d} + \frac{1}{2} \mathbf{d}^{\mathsf{T}$$

when ||**d**|| is small enough.

Necessary conditions:

$$\nabla f(\mathbf{x}_0) = \mathbf{0}, \ \nabla^2 f(\mathbf{x}_0) \text{ is PSD.}$$



Differential Calculus Review: Derivative

at the point where we're taking derivative...

If $f : \mathbb{R}^d \to \mathbb{R}$ is differentiable at $\mathbf{x}_0 \in \mathbb{R}^d$...

 $\lim \frac{f(\mathbf{x}) - (f(\mathbf{x}_0))}{\|}$ $\mathbf{X} \rightarrow \mathbf{X}_0$

Throughout this section, $\mathbf{d} = \mathbf{x} - \mathbf{x}_0$.

$$\frac{\text{linear approximation}}{|\mathbf{x} - \mathbf{x}_0||} = 0$$

as x gets closer to x_0 the function is closer and closer to its linear approximation!

Unconstrained Minima Necessary conditions

Theorem (Necessary Conditions for Unconstrained Local Minimum).

Suppose $\mathbf{x}^* \in int(\mathscr{C})$ is an <u>unconstrained local minimum</u>. Then,

First-order condition. If f is differentiable at \mathbf{x}^* , then $\nabla f(\mathbf{x}^*) = \mathbf{0}$.

Second-order condition. If f is twice-differentiable at \mathbf{x}^* , then $\nabla^2 f(\mathbf{x}^*)$ is positive semidefinite, i.e. $\mathbf{v}^{\top} \nabla^2 f(\mathbf{x}^*) \mathbf{v} \ge 0$ for all $\mathbf{v} \in \mathbb{R}^d$.

- minimize $f(\mathbf{x})$
- subject to $\mathbf{x} \in \mathscr{C}$

Proof of first order necessary condition Step 1: Use definition of gradient for $\alpha \mathbf{d}$

Choose an arbitrary direction $\alpha \mathbf{d} \in \mathbb{R}^d$, where $\|\mathbf{d}\| = 1$ is a unit vector and $\alpha > 0$ is a scalar.

f is differentiable, so...

 $\lim \frac{f(\mathbf{x}^* + \alpha \mathbf{d}) - \mathbf{d}}{\alpha \mathbf{d}}$ $\alpha \rightarrow 0$

which is the same as stating:

$$\lim_{\alpha \to 0} \frac{f(\mathbf{x}^* + \alpha \mathbf{d}) - f(\mathbf{x}^*)}{\alpha} = \nabla f(\mathbf{x}^*)^{\mathsf{T}} \mathbf{d}.$$

- First-order condition. If f is differentiable at \mathbf{x}^* , then $\nabla f(\mathbf{x}^*) = \mathbf{0}$.

$$\frac{f(\mathbf{x}^*) - \alpha \nabla f(\mathbf{x}^*)^{\mathsf{T}} \mathbf{d}}{\alpha \|\mathbf{d}\|} = 0$$

Proof of first order necessary condition Step 2: Use local optimality on difference $f(\mathbf{x}^* + \alpha \mathbf{d}) - f(\mathbf{x}^*)$

From Step 1,



 \mathbf{x}^* is an <u>unconstrained local minimum</u>, so there exists a neighborhood $B_{\delta}(\mathbf{x}^*)$ such that $f(\mathbf{x}) \ge f(\mathbf{x}^*)$ for all $\mathbf{x} \in B_{\delta}(\mathbf{x}^*)$. So if $\alpha < \delta$ (sufficiently small),

 $f(\mathbf{x}^* + \alpha \mathbf{d}) \ge f(\mathbf{x}^*)$

First-order condition. If f is differentiable at \mathbf{x}^* , then $\nabla f(\mathbf{x}^*) = \mathbf{0}$.

$$\frac{f(\mathbf{x}^*)}{\mathbf{x}^*} = \nabla f(\mathbf{x}^*)^{\mathsf{T}} \mathbf{d}.$$

$$\mathbf{x}^*) \Longrightarrow \nabla f(\mathbf{x}^*)^\top \mathbf{d} \ge 0.$$

Proof of first order necessary condition Step 3: $\mathbf{d} \in \mathbb{R}^n$ was an arbitrary direction.

- $\mathbf{d} = \mathbf{e}_1 \implies \nabla f(\mathbf{x}^*)_1 \ge 0 \text{ and } \mathbf{d} = -\mathbf{e}_1 \implies \nabla f(\mathbf{x}^*)_1 < 0$
- $\mathbf{d} = \mathbf{e}_2 \implies \nabla f(\mathbf{x}^*)_2 \ge 0 \text{ and } \mathbf{d} = -\mathbf{e}_2 \implies \nabla f(\mathbf{x}^*)_2 < 0$
- •
- $\mathbf{d} = \mathbf{e}_d \implies \nabla f(\mathbf{x}^*)_d \ge 0 \text{ and } \mathbf{d} = -\mathbf{e}_d \implies \nabla f(\mathbf{x}^*)_d < 0$ Therefore, $\nabla f(\mathbf{x}^*) = \mathbf{0}$.

- First-order condition. If f is differentiable at \mathbf{x}^* , then $\nabla f(\mathbf{x}^*) = \mathbf{0}$.
- From Step 2, if $\alpha < \delta$ (sufficiently small), $\nabla f(\mathbf{x}^*)^{\mathsf{T}} \mathbf{d} \ge 0$. But $\mathbf{d} \in \mathbb{R}^d$ was an arbitrary direction with $\|\mathbf{d}\| = 1$.

Unconstrained Minima Necessary conditions

Theorem (Necessary Conditions for Unconstrained Local Minimum).

Suppose $\mathbf{x}^* \in int(\mathscr{C})$ is an <u>unconstrained local minimum</u>. Then,

First-order condition. If f is differentiable at \mathbf{x}^* , then $\nabla f(\mathbf{x}^*) = \mathbf{0}$.

Second-order condition. If f is twice-differentiable at \mathbf{x}^* , then $\nabla^2 f(\mathbf{x}^*)$ is positive semidefinite, i.e. $\mathbf{v}^{\top} \nabla^2 f(\mathbf{x}^*) \mathbf{v} \ge 0$ for all $\mathbf{v} \in \mathbb{R}^d$.

- minimize $f(\mathbf{x})$
- subject to $\mathbf{x} \in \mathscr{C}$

Proof of second order necessary condition Step 1: Use second-order Taylor approximation

Choose an arbitrary direction $\alpha \mathbf{d} \in \mathbb{R}^d$ where $\alpha > 0$ is a scalar. By Taylor's Theorem (Peano's form) there exists $\delta > 0$ such that for all $\mathbf{d} \in B_{\delta}(\mathbf{0})$:

$$f(\mathbf{x}^* + \alpha \mathbf{d}) - \left(f(\mathbf{x}^*) + \alpha \nabla f(\mathbf{x}^*)^{\mathsf{T}} \mathbf{d} + \frac{\alpha^2}{2} \mathbf{d}^{\mathsf{T}} \nabla^2 f(\mathbf{x}^*) \mathbf{d} \right) \le \alpha \|\mathbf{d}\|^2.$$

Second-order condition. If f is twice-differentiable at \mathbf{x}^* , then $\nabla^2 f(\mathbf{x}^*)$ is PSD.

Proof of second order necessary condition Step 2: Use first-order condition so $\alpha \nabla f(\mathbf{x}^*)^{\mathsf{T}} \mathbf{d} = 0$

$$f(\mathbf{x}^* + \alpha \mathbf{d}) - \left(f(\mathbf{x}^*) + \alpha \nabla f(\mathbf{x}^*)^{\mathsf{T}} \mathbf{d} + \frac{\alpha^2}{2} \mathbf{d}^{\mathsf{T}} \nabla^2 f(\mathbf{x}^*) \mathbf{d} \right) \le \alpha \|\mathbf{d}\|^2$$

x^{*} is an *unconstrained local minimum,* so by first-order condition (just proved):

$$f(\mathbf{x}^* + \alpha \mathbf{d}) - f(\mathbf{x}^*) \le \frac{\alpha^2}{2} \mathbf{d}^\top \nabla^2 f(\mathbf{x}^*) \mathbf{d} + \alpha \|\mathbf{d}\|^2$$

Second-order condition. If f is twice-differentiable at \mathbf{x}^* , then $\nabla^2 f(\mathbf{x}^*)$ is PSD.

Proof of second order necessary condition Step 3: Divide by $\|\mathbf{d}\|^2$ and use local optimality: $f(\mathbf{x}^* + \alpha \mathbf{d}) - f(\mathbf{x}^*) \ge 0$

Second-order condition. If f is twice-differentiable at \mathbf{x}^* , then $\nabla^2 f(\mathbf{x}^*)$ is PSD.

 $f(\mathbf{x}^* + \alpha \mathbf{d}) - f(\mathbf{x}^*) \leq$

Divide by $\|\mathbf{d}\|^2$ everywhere and take the limit as

$$\lim_{\alpha \to 0} \frac{f(\mathbf{x}^* + \alpha \mathbf{d}) - f(\mathbf{x}^*)}{\|\mathbf{d}\|^2} - \frac{\alpha^2}{2\|\mathbf{d}\|^2} \mathbf{d}^\top \nabla^2 f(\mathbf{x}^*) \mathbf{d} = 0$$

By local optimality of \mathbf{x}^* and arbitrary $\mathbf{d} \in \mathbb{R}^d$:

$$0 \leq \frac{f(\mathbf{x}^* + \alpha \mathbf{d}) - f(\mathbf{x}^*)}{\|\mathbf{d}\|^2}, \text{ so } 0 \leq \frac{1}{2} (\alpha \mathbf{d} / \|\mathbf{d}\|)^\top \nabla^2$$

$$\leq \frac{\alpha^2}{2} \mathbf{d}^{\mathsf{T}} \nabla^2 f(\mathbf{x}^*) \mathbf{d} + \alpha \|\mathbf{d}\|^2.$$

s
$$\alpha \to 0$$
:

 ${}^{2}f(\mathbf{x}^{*})(\alpha \mathbf{d}/\|\mathbf{d}\|) \Longrightarrow \nabla^{2}f(\mathbf{x}^{*})$ is PSD (definition of PSD).

Unconstrained Minima Necessary conditions

Theorem (Necessary Conditions for Unconstrained Local Minimum).

Suppose $\mathbf{x}^* \in int(\mathscr{C})$ is an <u>unconstrained local minimum</u>. Then,

First-order condition. If f is differentiable at \mathbf{x}^* , then $\nabla f(\mathbf{x}^*) = \mathbf{0}$.

Second-order condition. If f is twice-differentiable at \mathbf{x}^* , then $\nabla^2 f(\mathbf{x}^*)$ is positive semidefinite, i.e. $\mathbf{v}^{\mathsf{T}} \nabla^2 f(\mathbf{x}^*) \mathbf{v} \ge 0$ for all $\mathbf{v} \in \mathbb{R}^d$.

- minimize $f(\mathbf{x})$
- subject to $\mathbf{x} \in \mathscr{C}$

Unconstrained local minima Sufficient conditions

Least Squares OLS Theorem

Proof (Calculus proof of OLS).

 $f(\mathbf{w}) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 \iff f(\mathbf{w}) = \mathbf{w}^{\mathsf{T}}\mathbf{X}^{\mathsf{T}}\mathbf{X}\mathbf{w} - 2\mathbf{w}^{\mathsf{T}}\mathbf{X}^{\mathsf{T}}\mathbf{y} + \mathbf{y}^{\mathsf{T}}\mathbf{y}$ "First derivative test." $\nabla_{\mathbf{w}} f(\mathbf{w}) = 2(\mathbf{X}^{\mathsf{T}}\mathbf{X})\mathbf{w} - 2\mathbf{X}^{\mathsf{T}}\mathbf{y}$. $2(\mathbf{X}^{\mathsf{T}}\mathbf{X})\mathbf{w} - 2\mathbf{X}^{\mathsf{T}}\mathbf{y} = \mathbf{0} \implies \mathbf{X}^{\mathsf{T}}\mathbf{X}\mathbf{w} = \mathbf{X}^{\mathsf{T}}\mathbf{y}$ rank $(\mathbf{X}) = d \Longrightarrow \operatorname{rank}(\mathbf{X}^{\mathsf{T}}\mathbf{X}) = d \Longrightarrow \mathbf{X}^{\mathsf{T}}\mathbf{X}$ is invertible:

$$\hat{\mathbf{w}} = (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}.$$

"Second derivative test." $\nabla_{\mathbf{w}}^2 f(\mathbf{w}) = 2\mathbf{X}^{\mathsf{T}}\mathbf{X}$.

 $\operatorname{rank}(\mathbf{X}) = d \implies \operatorname{rank}(\mathbf{X}^{\mathsf{T}}\mathbf{X}) = d \implies \lambda_1, \dots, \lambda_d > 0$ $\implies \mathbf{X}^{\mathsf{T}}\mathbf{X} \text{ is positive definite!}$



Sufficient Conditions Comparison to single variable

 $f(x_0 + \delta) \approx f(x_0) + f'(x_0)\delta + \frac{1}{2}f''(x_0)\delta^2$

when δ is small enough.

Necessary conditions:

$$f'(x_0) = 0, f''(x_0) > 0.$$

$$f(\mathbf{x}_0 + \mathbf{d}) \approx f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^{\mathsf{T}} \mathbf{d} + \frac{1}{2} \mathbf{d}^{\mathsf{T}} \nabla^2 f(\mathbf{x}_0) \mathbf{d}^{\mathsf{T}} \mathbf{d}^{\mathsf{T}} \mathbf{d}^{\mathsf{T}} \mathbf{d}^{\mathsf{T}} \nabla^2 f(\mathbf{x}_0) \mathbf{d}^{\mathsf{T}} \mathbf{d}^$$

when ||**d**|| is small enough.

Necessary conditions:

$$\nabla f(\mathbf{x}_0) = \mathbf{0}, \ \nabla^2 f(\mathbf{x}_0) \text{ is PD.}$$



Unconstrained Minima Sufficient conditions

Theorem (Sufficient Conditions for Unconstrained Local Minimum).

Let $\mathbf{x}^* \in int(\mathscr{C})$. If $f \in \mathscr{C}^2$ and

then **x**^{*} is a *strict* unconstrained local minimum.

- minimize $f(\mathbf{x})$
- subject to $\mathbf{x} \in \mathscr{C}$

- $\nabla f(\mathbf{x}^*) = \mathbf{0}$ and $\nabla^2 f(\mathbf{x}^*)$ is positive definite,

Proof of second order sufficient condition Step 1: Use second-order Taylor approximation

Choose an arbitrary direction $\alpha \mathbf{d} \in \mathbb{R}^d$ where $\alpha > 0$ is a scalar. By Taylor's Theorem (Peano's form) there exists $\delta > 0$ such that for all $\mathbf{d} \in B_{\delta}(\mathbf{0})$:

$$f(\mathbf{x}^* + \alpha \mathbf{d}) - \left(f(\mathbf{x}^*) + \alpha \nabla f(\mathbf{x}^*)^{\mathsf{T}} \mathbf{d} + \frac{\alpha^2}{2} \mathbf{d}^{\mathsf{T}} \nabla^2 f(\mathbf{x}^*) \mathbf{d} \right) \ge -\alpha \|\mathbf{d}\|^2.$$

Note: Used the negative direction of the statement (which is absolute value).

Second-order condition. If $\nabla^2 f(\mathbf{x}^*)$ is PD, then \mathbf{x}^* is an unconstrained local minimum.

Proof of second order sufficient condition Step 2: Eigenvalues of PD matrix are positive

From Step 1, for any $\mathbf{d} \in \mathbb{R}^d$ with $\|\mathbf{d}\| = 1$ and $\alpha > 0$,

$$f(\mathbf{x}^* + \alpha \mathbf{d}) - \left(f(\mathbf{x}^*) + \alpha \nabla f(\mathbf{x}^*)^{\mathsf{T}} \mathbf{d} + \frac{\alpha^2}{2} \mathbf{d}^{\mathsf{T}} \nabla^2 f(\mathbf{x}^*) \mathbf{d} \right) \ge -\alpha \|\mathbf{d}\|^2.$$

Let the eigenvalues of $\nabla^2 f(\mathbf{x}^*)$ be $\lambda_1 \ge \dots \ge \lambda_d > 0$, and consider the smallest eigenvalue, $\lambda_d > 0$ with unit eigenvector \mathbf{v}_d with $\|\mathbf{v}_d\| = 1$.

$$\implies \frac{\alpha^2}{2} \mathbf{d}^{\mathsf{T}} \nabla^2 f(\mathbf{x}^*) \mathbf{d} \ge \frac{\alpha^2}{2} \mathbf{v}_d^{\mathsf{T}} \nabla f(\mathbf{x}^*) \mathbf{v}_d = \frac{\lambda_d \alpha^2}{2}.$$

Second-order condition. If $\nabla^2 f(\mathbf{x}^*)$ is PD, then \mathbf{x}^* is an unconstrained local minimum.

Proof of second order sufficient condition Step 3: $\alpha \nabla f(\mathbf{x}^*)^{\mathsf{T}} \mathbf{d} = 0$ from first-order condition

Cancel out the first-order term $\alpha \nabla f(\mathbf{x}^*)^{\mathsf{T}} \mathbf{d} = \mathbf{0}$ and plugin the eigenvalue lower bound

 $f(\mathbf{x}^* + \alpha \mathbf{d}) - f(\mathbf{x}^*) \ge \alpha \nabla f(\mathbf{x}^*)$

so this simplifies to...

$$f(\mathbf{x}^* + \alpha \mathbf{d}) - f(\mathbf{x}^*) \ge \frac{\lambda_d \alpha^2}{2} - \alpha \|\mathbf{d}\|^2 = \left(\frac{\lambda_d}{2} - \frac{\|\mathbf{d}\|^2}{\alpha}\right) \alpha^2.$$

- Second-order condition. If $\nabla^2 f(\mathbf{x}^*)$ is PD, then \mathbf{x}^* is an unconstrained local minimum.

$$(\mathbf{x}^*)^{\mathsf{T}}\mathbf{d} + \frac{\alpha^2}{2} \mathbf{d}^{\mathsf{T}} \nabla^2 f(\mathbf{x}^*) \mathbf{d} - \alpha \|\mathbf{d}\|^2$$

Proof of second order sufficient condition Step 4: Divide by $\|d\|^2$ and consider small enough $d \rightarrow 0$

Second-order condition. If $\nabla^2 f(\mathbf{x}^*)$ is PD, then \mathbf{x}^* is an unconstrained local minimum.

Take our inequality

$$f(\mathbf{x}^* + \alpha \mathbf{d}) - f(\mathbf{x}^*) \ge \frac{\lambda_d \alpha^2}{2} - \alpha \|\mathbf{d}\|^2 = \left(\frac{\lambda_d}{2} - \frac{\|\mathbf{d}\|^2}{\alpha}\right) \alpha^2.$$

and divide by $\|\mathbf{d}\|^2$ to get:

$$\frac{f(\mathbf{x}^* + \alpha \mathbf{d}) - f(\mathbf{x}^*)}{\|\mathbf{d}\|^2} \ge \left(\frac{\lambda_d}{2\|\mathbf{d}\|^2} - \frac{1}{\alpha}\right) \alpha^2, \text{ an}$$

In d sufficiently small $\mathbf{d}
ightarrow \mathbf{0}$ makes the RHS positive.

Least Squares OLS Theorem

Proof (Calculus proof of OLS).

 $f(\mathbf{w}) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 \iff f(\mathbf{w}) = \mathbf{w}^{\mathsf{T}}\mathbf{X}^{\mathsf{T}}\mathbf{X}\mathbf{w} - 2\mathbf{w}^{\mathsf{T}}\mathbf{X}^{\mathsf{T}}\mathbf{y} + \mathbf{y}^{\mathsf{T}}\mathbf{y}$ "First derivative test." $\nabla_{\mathbf{w}} f(\mathbf{w}) = 2(\mathbf{X}^{\mathsf{T}}\mathbf{X})\mathbf{w} - 2\mathbf{X}^{\mathsf{T}}\mathbf{y}$. $2(\mathbf{X}^{\mathsf{T}}\mathbf{X})\mathbf{w} - 2\mathbf{X}^{\mathsf{T}}\mathbf{y} = \mathbf{0} \implies \mathbf{X}^{\mathsf{T}}\mathbf{X}\mathbf{w} = \mathbf{X}^{\mathsf{T}}\mathbf{y}$ rank $(\mathbf{X}) = d \Longrightarrow \operatorname{rank}(\mathbf{X}^{\mathsf{T}}\mathbf{X}) = d \Longrightarrow \mathbf{X}^{\mathsf{T}}\mathbf{X}$ is invertible:

$$\hat{\mathbf{w}} = (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}.$$

"Second derivative test." $\nabla^2_{\mathbf{w}} f(\mathbf{w}) = 2\mathbf{X}^{\mathsf{T}} \mathbf{X}$.

 $\operatorname{rank}(\mathbf{X}) = d \implies \operatorname{rank}(\mathbf{X}^{\mathsf{T}}\mathbf{X}) = d \implies \lambda_1, \dots, \lambda_d > 0$

 $\implies \mathbf{X}^{\mathsf{T}}\mathbf{X}$ is positive definite!



Finding global minima Introducing constraint sets

Types of Minima Big picture

We want to find **global minima**.

Global minima could be either unconstrained local minima or constrained local minima.

Without *C*, global minima are just an *unconstrained local minima*.

f(x)

With *C*, global minima may lie on the boundary of the constraint set.

Find local minima, then test!





Unconstrained Minima Necessary conditions

Theorem (Necessary Conditions for Unconstrained Local Minimum).

Suppose $\mathbf{x}^* \in int(\mathscr{C})$ is an <u>unconstrained local minimum</u>. Then,

First-order condition. If f is differentiable at \mathbf{x}^* , then $\nabla f(\mathbf{x}^*) = \mathbf{0}$.

Second-order condition. If f is twice-differentiable at \mathbf{x}^* , then $\nabla^2 f(\mathbf{x}^*)$ is positive semidefinite, i.e. $\mathbf{v}^{\top} \nabla^2 f(\mathbf{x}^*) \mathbf{v} \ge 0$ for all $\mathbf{v} \in \mathbb{R}^d$.

- minimize $f(\mathbf{x})$
- subject to $\mathbf{x} \in \mathscr{C}$

Finding global minima Using necessary conditions with constraints

Necessary conditions for unconstrained local minima:

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

How do we find the *global* minimum from this?

- 1. Find *unconstrained local minima* from first-order condition $M := \{ \mathbf{x}^* \in \operatorname{int}(\mathscr{C}) : \nabla f(\mathbf{x}^*) = \mathbf{0} \}.$
- 2. Find the set of "boundary" points $B := \mathscr{C} \setminus int(\mathscr{C}) = \{ \mathbf{x} \in \mathscr{C} : \mathbf{x} \notin int(\mathscr{C}) \}$.
- 3. The global minimum must be in the set $M \cup B$, so evaluate f on all $\mathbf{x} \in M \cup B$.

- and $\nabla^2 f(\mathbf{x}^*) \ge 0$.

Finding global minima Using necessary conditions with constraints

Necessary conditions for unconstrained local minima:

 $\nabla f(\mathbf{x}^*) = \mathbf{0}$ and $\nabla^2 f(\mathbf{x}^*) \ge 0$.

How do we find the *global* minimum from this?

- 1. Find *unconstrained local minima* from first-order condition $M := \{\mathbf{x}^* \in int(\mathscr{C}) : \nabla f(\mathbf{x}^*) = \mathbf{0}\}.$
- 2. Find the set of "boundary" points $B := \mathscr{C} \setminus \operatorname{int}(\mathscr{C}) = \{ \mathbf{x} \in \mathscr{C} : \mathbf{x} \notin \operatorname{int}(\mathscr{C}) \}.$
- 3. The global minimum must be in the set $M \cup B$, so evaluate f on all $\mathbf{x} \in M \cup B$.



local min global min

Finding global minima Using necessary conditions without constraints

Necessary conditions for unconstrained local minima:

How do we find the global minimum from this when $\mathscr{C} = \mathbb{R}^d$?

- 1. Find unconstrained local minima from first-order condition $M := \{\mathbf{x}^* \in \mathbb{R}^d : \nabla f(\mathbf{x}^*) = \mathbf{0}\}$.
- 2. There are no boundary points! ($B := \mathscr{C} \setminus int(\mathscr{C}) = \{ \mathbf{x} \in \mathscr{C} : \mathbf{x} \notin int(\mathscr{C}) \} = \emptyset \}$
- 3. The global minimum must be in the set M, so evaluate f on all $\mathbf{x} \in M$.

- $\nabla f(\mathbf{x}^*) = \mathbf{0}$ and $\nabla^2 f(\mathbf{x}^*) \ge \mathbf{0}$.

Finding global minima Using necessary conditions without constraints

Necessary conditions for unconstrained local minima:

$$\nabla f(\mathbf{x}^*) = \mathbf{0}$$
 and $\nabla^2 f(\mathbf{x}^*) \ge 0$.

How do we find the *global* minimum from this when $\mathscr{C} = \mathbb{R}^d$?

- 1. Find *unconstrained local minima* from first-order condition $M := \{\mathbf{x}^* \in \mathbb{R}^d : \nabla f(\mathbf{x}^*) = \mathbf{0}\}.$
- 2. There are no boundary points! $(B := \mathscr{C} \setminus int(\mathscr{C}) = \{ \mathbf{x} \in \mathscr{C} : \mathbf{x} \notin int(\mathscr{C}) \} = \emptyset)$
- 3. The global minimum must be in the set M, so evaluate f on all $\mathbf{x} \in M$.



















10

Unconstrained Minima Example

When $f : \mathbb{R} \to \mathbb{R}$ is one-dimensional on $\mathscr{C} = [a, b]$ and differentiable on $int(\mathscr{C}) := (a, b)$.

minimize x^2 subject to $x \in [1,3]$

Unconstrained Minima Example

minimize x^2 subject to $x \in [1,3]$

When $f : \mathbb{R} \to \mathbb{R}$ is one-dimensional on $\mathscr{C} = [a, b]$ and differentiable on $int(\mathscr{C}) := (a, b)$.



constrained min unconstrained min
Unconstrained Minima Example: Why haven't we solved optimization?

How do we deal with the possible constrained local minima induced by \mathscr{C} ?

- minimize $f(x_1, x_2)$
- subject to $x_1^2 + x_2^2 \le 1$

Need to evaluate f on the infinite number of points on the boundary of the circle, $\mathscr{C}(\mathsf{int}(\mathscr{C})) := \{ \mathbf{x} \in \mathbb{R}^2 : x_1^2 + x_2^2 = 1 \}$

Unconstrained Minima Example: Why haven't we solved optimization?

minimize $f(x_1, x_2)$ subject to $x_1^2 + x_2^2 \le 1$

Need to evaluate *f* on the infinite number of points on the boundary of the circle, $\mathscr{C}\setminus\operatorname{int}(\mathscr{C}) := \{\mathbf{x} \in \mathbb{R}^2 : x_1^2 + x_2^2 = 1\}!$

How do we deal with the possible constrained local minima induced by C?



Unconstrained Minima Example: Why haven't we solved optimization?

minimize $f(x_1, x_2)$ subject to $x_1^2 + x_2^2 \le 1$

Need to evaluate *f* on the infinite number of points on the boundary of the circle, $\mathscr{C}\setminus\operatorname{int}(\mathscr{C}) := \{\mathbf{x} \in \mathbb{R}^2 : x_1^2 + x_2^2 = 1\}!$

How do we deal with the possible constrained local minima induced by C?





Constrained Minima Equality Constraints and the Lagrangian

Types of Minima Which type of minima are each of these points?

f(x)

 $\begin{array}{ll} \text{minimize} & f(\mathbf{x}) \\ \text{subject to} & \mathbf{x} \in \mathscr{C} \end{array}$

constrained local:

 $f(\hat{\mathbf{x}}) \leq f(\mathbf{x})$ for all $\mathbf{x} \in \mathscr{C} \cap B_{\delta}(\hat{\mathbf{x}})$

<u>unconstrained local:</u>

 $f(\hat{\mathbf{x}}) \leq f(\mathbf{x})$ for all $\mathbf{x} \in B_{\delta}(\hat{\mathbf{x}})$ and $B_{\delta}(\hat{\mathbf{x}}) \subset \mathscr{C}$. global:

 $f(\mathbf{x}^*) \leq f(\mathbf{x})$ for all $\mathbf{x} \in \mathscr{C}$.





Constrained Local Minima Minimum values on the "edge of the constraint set"



constrained min unconstrained min



Constrained Minima Equality constrained optimization

Objective function $f : \mathbb{R}^d \to \mathbb{R}$ like before.

 h_1, \ldots, h_m are \mathscr{C}^1 functions $h_i : \mathbb{R}^d \to \mathbb{R}$ that form \mathscr{C} , the constraint set.



Constrained Minima Equality constrained optimization

minimiz subject

The = 0 constraint is without loss of generality:

If we want $h_j(\mathbf{x}) = c$ then we can always consider $h'_j(\mathbf{x}) = h_j(\mathbf{x}) - c = 0$ instead.

ze
$$f(\mathbf{x})$$

to $h_1(\mathbf{x}) = 0$
:
 $h_m(\mathbf{x}) = 0$

Constrained Minima: Equality Constraints Example: Maximum Volume Box

- minimize $x_1 x_2 x_3$
- subject to $x_1x_2 + x_2x_3 + x_1x_3 c/2 = 0$

Objective function: $f(\mathbf{x}) = x_1 x_2 x_3$

Single equality constraint: $h : \mathbb{R}^3 \to \mathbb{R}$, defined as $h(\mathbf{x}) = x_1x_2 + x_2x_3 + x_1x_3 - c/2$.

Constrained Minima: Equality Constraints Idea

Convert constrained optimization problem into an *unconstrained* optimization problem.

 $j \in [m]$), represented by a vector $\lambda \in \mathbb{R}^m$ (the <u>Lagrange multipliers</u>).

- Then deal with unconstrained problem as we did before:
 - $\nabla f(\mathbf{x}) = \mathbf{0}$ and $\nabla^2 f(\mathbf{x}) \ge 0$.

The unconstrained optimization problem will have m more variables (for each constraint h_i for

Constrained Minima: Equality Constraints Definition of the Lagrangian

minimiz subject

The associated Lagrangian function $L : \mathbb{R}^d \times \mathbb{R}^m \to \mathbb{R}$ is

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

ze
$$f(\mathbf{x})$$

to $h_1(\mathbf{x}) = 0$
 \vdots
 $h_m(\mathbf{x}) = 0$

$$f(\mathbf{x}) + \sum_{i=1}^{m} \lambda_i h_i(\mathbf{x}).$$

Constrained Minima: Equality Constraints Regularity Conditions

minimize $f(\mathbf{x})$ subject to h_1

A point $\mathbf{x} \in \mathbb{R}^d$ is a <u>regular point</u> if:

1. **x** is feasible, i.e. $h_1(\mathbf{x}) = 0, ..., h_m(\mathbf{x}) = 0$.

2. The gradients $\nabla h_1(\mathbf{x}), \ldots, \nabla h_m(\mathbf{x})$ are linearly independent.

Constraints are "non-redundant." This is a property of how we write down our problem.

$$\mathbf{x}$$
$$(\mathbf{x}) = 0, \dots, h_m(\mathbf{x}) = 0$$

Constrained Minima: Equality Constraints Lagrange Multiplier Theorem: Necessary Conditions

Theorem (Lagrange Multiplier Theorem - Necessary). Let $\mathbf{x}^* \in \mathbb{R}^d$ be a local minimum that is a regular point. Then, there exists a unique vector $\lambda \in \mathbb{R}^m$ called a Lagrange multiplier such that

 $\nabla f(\mathbf{x}^*) +$

$$\sum_{i=1}^{m} \lambda_i \nabla h_i(\mathbf{x}^*) = 0$$

Constrained Minima: Equality Constraints Lagrange Multiplier Theorem: Necessary Conditions

 $\nabla f(\mathbf{x}^*) +$

If, in addition, f and h_1, \ldots, h_m are twice continuously differentiable,

$$\mathbf{d}^{\mathsf{T}}\left(\nabla^2 f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i \nabla^2 h_i(\mathbf{x}^*)\right) \mathbf{d} \ge 0$$

for all $\mathbf{d} \in \mathbb{R}^d$ such that $\nabla h_i(\mathbf{x}^*)^\top \mathbf{d} = 0$ for all $j \in [m]$.

Theorem (Lagrange Multiplier Theorem - Necessary). Let $\mathbf{x}^* \in \mathbb{R}^d$ be a local minimum that is a regular point. Then, there exists a unique vector $\lambda \in \mathbb{R}^m$ called a <u>Lagrange multiplier</u> such that

$$\sum_{i=1}^{m} \lambda_i \nabla h_i(\mathbf{x}^*) = 0$$

Constrained Minima: Equality Constraints How to remember the Lagrange multiplier theorem

 $\nabla f(\mathbf{x}) + \sum_{i=1}^{n} f(\mathbf{x})$

Remember the necessary conditions for unconstrained local minima: $\nabla f(\mathbf{x}) = \mathbf{0}$ and $\nabla^2 f(\mathbf{x}) \ge 0$.

Applying first-order necessary conditions for Lagrangian, so local minimum $(\mathbf{x}^*, \lambda^*)$ must satisfy

 $\nabla_{\mathbf{x}} L(\mathbf{x}^*, \lambda^*) = \mathbf{0}$ and $\nabla_{\lambda} L(\mathbf{x}^*, \lambda^*) = \mathbf{0}$.

Notice that $\nabla_{\lambda} L(\mathbf{x}^*, \lambda^*) = \mathbf{0}$ is the same as requiring feasibility: $h_i(\mathbf{x}^*) = 0$ for all $j \in [m]$.

$$\sum_{i=1}^{m} \lambda_i \nabla h_i(\mathbf{x}) = 0$$

Constrained Minima: Equality Constraints Lagrange Multiplier Theorem: Sufficient Conditions

such that $\mathbf{x}^* \in \mathbb{R}^d$ and $\lambda^* \in \mathbb{R}^m$ satisfy

Then, \mathbf{x}^* is a local minimum.

- Theorem (Lagrange Multiplier Theorem Sufficient Conditions). Let f and h be \mathscr{C}^2 functions,
 - $\nabla_{\mathbf{x}} L(\mathbf{x}^*, \lambda^*) = 0$ and $\nabla_{\lambda} L(\mathbf{x}^*, \lambda^*) = 0$
 - $\mathbf{d}^{\top} \nabla^2_{\mathbf{x},\mathbf{x}} L(\mathbf{x}^*, \lambda^*) \mathbf{d} > 0$, for all $\mathbf{d} \in \mathbb{R}^d$ such that $\nabla h_i(\mathbf{x}^*)^{\top} \mathbf{d} = 0$ for all $j \in [m]$.

Constrained Minima: Equality Constraints How do we use the Lagrangian?

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

Assuming a global minimum exists, to find it...

1. Find the set $(\mathbf{x}^*, \lambda^*)$ of <u>regular points</u> satisfying the first-order necessary conditions:

 $\nabla_{\mathbf{x}} L(\mathbf{x}^*, \lambda^*) = 0$ and $\nabla_{\lambda} L(\mathbf{x}^*, \lambda^*) = 0$.

- 2. Find the set of all non-regular points.
- 3. The global minima must be among the points in (1) or (2).

$$\hat{\mathbf{x}}(\mathbf{x}) + \sum_{i=1}^{m} \lambda_i h_i(\mathbf{x}).$$

Constrained Minima: Equality Constraints Example: Maximum Volume Box

- minimize $x_1 x_2 x_3$

subject to $x_1x_2 + x_2x_3 + x_1x_3 - c/2 = 0$

Constrained Minima Inequality Constraints and the KKT Theorem

Constrained Minima Inequality constrained optimization



Objective function $f : \mathbb{R}^d \to \mathbb{R}$ like before. h_1, \ldots, h_m are \mathscr{C}^1 functions $h_i : \mathbb{R}^d \to \mathbb{R}$ that form \mathscr{C} , the constraint set. g_1, \ldots, g_r are \mathscr{C}^1 functions $g_i : \mathbb{R}^d \to \mathbb{R}$ that form \mathscr{C} , the constraint set.



Constrained Minima Inequality constrained optimization

- minimize $f(\mathbf{x})$

To solve: Reduce to equality constrained optimization.

The only difference is that each *inequality constraint* can either be <u>active</u> or not.



A constraint $j \in [r]$ is <u>active</u> if $g_i(\mathbf{x}) = 0$.

Constrained Minima: Inequality Constraints Definition of active constraints

For feasible $\mathbf{x} \in \mathbb{R}^d$ the set of <u>active inequality constraints</u> is

A point $\mathbf{x} \in \mathbb{R}^d$ is a <u>regular point</u> if it is feasible and the gradients $\{\nabla h_1(\mathbf{x}), \dots, \nabla h_m(\mathbf{x})\} \cup \{\nabla g_j(\mathbf{x}) : j \in \mathscr{A}(\mathbf{x})\}$

are linearly independent.

 $\mathscr{A}(\mathbf{x}) := \{j : g_j(\mathbf{x}) = 0\} \subseteq [r].$

Constrained Minima: Inequality Constraints Lagrangian in Inequality Constrained Optimization

- minimize $f(\mathbf{x})$
- subject to h_1
 - g_1
- The Lagrangian function $L : \mathbb{R}^d \times \mathbb{R}^m \times \mathbb{R}^r \to \mathbb{R}$ is the function
 - $L(\mathbf{x}, \lambda, \mu) := f(\mathbf{x}) +$

x)
(**x**) = 0,...,
$$h_m(\mathbf{x}) = 0$$

(**x**) $\leq 0,..., g_r(\mathbf{x}) \leq 0$

+
$$\sum_{i=1}^{m} \lambda_i h_i(\mathbf{x}) + \sum_{j=1}^{r} \mu_j g_j(\mathbf{x}).$$

Constrained Minima: Inequality Constraints Karush-Kuhn-Tucker (KKT) Theorem

Theorem (KKT Theorem - Necessary Conditions). Let $\mathbf{x}^* \in \mathbb{R}^d$ be a local minimum that is a <u>regular point</u>. Then, there exists unique vectors $\lambda \in \mathbb{R}^m$ and $\mu \in \mathbb{R}^r$ called <u>Lagrange multipliers</u> such that

$$\nabla f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i \nabla h_i(\mathbf{x}^*) + \sum_{j=1}^r \mu_j^* \nabla g_j(\mathbf{x}^*) = 0,$$

where $\mu_i^* \ge 0$ for all $j \in [r]$ and $\mu_i^* = 0$ for all non-active constraints $j \notin \mathscr{A}(\mathbf{x}^*)$ (complementary slackness).



Constrained Minima: Inequality Constraints Karush-Kuhn-Tucker (KKT) Theorem

Theorem (KKT Theorem - Necessary Conditions). Let $\mathbf{x}^* \in \mathbb{R}^d$ be a local minimum that is a <u>regular point</u>. Then, there exists unique vectors $\lambda \in \mathbb{R}^m$ and $\mu \in \mathbb{R}^r$ called <u>Lagrange multipliers</u> such that

$$\nabla f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i \nabla h_i(\mathbf{x}^*) + \sum_{j=1}^r \mu_j^* \nabla g_j(\mathbf{x}^*) = 0,$$

If, in addition, f and the h_i are all twice continuously differentiable,

$$\mathbf{d}^{\mathsf{T}}\left(\nabla^2 f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i \nabla^2 h_i(\mathbf{x}^*)\right) \mathbf{d} \ge 0$$

for all $\mathbf{d} \in \mathbb{R}^d$ such that $\nabla h_i(\mathbf{x}^*)^\top \mathbf{d} = 0$ for all $j \in [m]$.

where $\mu_i^* \ge 0$ for all $j \in [r]$ and $\mu_i^* = 0$ for all non-active constraints $j \notin \mathscr{A}(\mathbf{x}^*)$ (complementary slackness).

Constrained Minima: Inequality Constraints Karush-Kuhn-Tucker (KKT) Theorem

 $L(\mathbf{x}, \lambda, \mu) := f(\mathbf{x}) + f(\mathbf{x})$

Write the previous necessary conditions at the local optimum $(\mathbf{x}^*, \lambda^*, \mu^*)$ as:

 $\nabla_{\mathbf{x}} L(\mathbf{x}^*, \lambda^*, \mu^*) = 0, \ \mathbf{h}(\mathbf{x}^*) = 0, \ \mathbf{g}(\mathbf{x}^*) \le 0$

where we also require the <u>complementary slackness</u> conditions:

 $\mu^* \geq 0$ and μ_j^*

$$\sum_{i=1}^{m} \lambda_i h_i(\mathbf{x}) + \sum_{j=1}^{r} \mu_j g_j(\mathbf{x}),$$

$$*g_j(\mathbf{x}^*) = 0, \,\forall j \in [r].$$

Constrained Minima: Inequality Constraints Karush-Kuhn-Tucker (KKT) Theorem: Sufficient Conditions

Theorem (KKT Theorem - Sufficient Conditions). Let f, h, and g be \mathscr{C}^2 functions, such that $\mathbf{x}^* \in \mathbb{R}^d, \lambda \in \mathbb{R}^m, \mu^* \in \mathbb{R}^r$ satisfy

$$\nabla_{\mathbf{x}} L(\mathbf{x}^*, \lambda^*, \mu^*) =$$

 $\mu^* \geq 0$ and μ_i^*

 $\mathbf{d}^{\top} \nabla_{\mathbf{x} \mathbf{x}}^{2} L(\mathbf{x}^{*}, \lambda^{*}, \mu^{*}) \mathbf{d} > 0,$

for all **d** such that $\nabla h_i(\mathbf{x}^*)^{\mathsf{T}}\mathbf{d} = 0$ for all $i \in [m]$ and $\nabla g_i(\mathbf{x}^*)^{\mathsf{T}}\mathbf{d} = 0$, $\forall j \in \mathscr{A}(\mathbf{x}^*)$.

Then, \mathbf{x}^* is a local minimum.

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

$$*g_j(\mathbf{x}^*) = 0, \forall j \in [r]$$

Constrained Minima: Inequality Constraints How do we use the Lagrangian?

Assuming a global minimum exists, to find a global minimum...

Find the set $(\mathbf{x}^*, \lambda^*, \mu^*)$ satisfying the necessary conditions: 1.

$$\nabla_{\mathbf{x}} L(\mathbf{x}^*, \lambda^*, \mu^*) = 0, \ \mathbf{h}(\mathbf{x}^*) =$$

$$\mu^* \ge 0$$
 and $\mu_j^* g_j(\mathbf{x}^*) = 0, \forall_j$

- 2. Find the set of all non-regular points.
- 3. The global minima must be among the points in (1) or (2).

 $L(\mathbf{x}, \lambda, \mu) = f(\mathbf{x}) + \sum_{i=1}^{m} \lambda_i h_i(\mathbf{x}) + \sum_{j=1}^{r} \mu_j g_j(\mathbf{x})$

- $= 0, \mathbf{g}(\mathbf{x}^*) \leq 0$ (first-order conditions)
- $j \in [r]$ (<u>complementary slackness</u>)

Constrained Minima: Inequality Constraints Example: Smallest point in a halfspace

minimize $\frac{1}{2} \|\mathbf{x}\|_2^2$ subject to $x_1 + x_2 + x_3 \le -3$

Least Squares Regression Regularization and Ridge Regression

Regression Setup (Example View)

<u>**Observed:**</u> Matrix of training samples $\mathbf{X} \in \mathbb{R}^{n \times d}$ and vector of training labels $\mathbf{y} \in \mathbb{R}^{n}$.

$$\mathbf{X} = \begin{bmatrix} \leftarrow \mathbf{x}_1^\top \rightarrow \\ \vdots \\ \leftarrow \mathbf{x}_n^\top \rightarrow \end{bmatrix} \mathbf{y}$$

<u>**Unknown:**</u> Weight vector $\mathbf{w} \in \mathbb{R}^d$ with weights w_1, \ldots, w_d .

<u>Goal</u>: For each $i \in [n]$, we predict: $\hat{y}_i = \mathbf{w}^\top \mathbf{x}_i = w_1 x_{i1} + \ldots + w_d x_{id} \in \mathbb{R}$.

Choose a weight vector that "fits the training data": $\mathbf{w} \in \mathbb{R}^d$ such that $y_i \approx \hat{y}_i$ for $i \in [n]$, or:

$$= \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \text{ where } \mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d.$$

 $\mathbf{X}\mathbf{w} = \hat{\mathbf{y}} \approx \mathbf{y}$.

Regression Setup (Feature View)

<u>**Observed:**</u> Matrix of training samples $\mathbf{X} \in \mathbb{R}^{n \times d}$ and vector of training labels $\mathbf{y} \in \mathbb{R}^{n}$.

$$\mathbf{X} = \begin{bmatrix} \uparrow & & \uparrow \\ \mathbf{x}_1 & \dots & \mathbf{x}_d \\ \downarrow & & \downarrow \end{bmatrix} \mathbf{y} = \mathbf{y}$$

<u>**Unknown:**</u> Weight vector $\mathbf{w} \in \mathbb{R}^d$ with weights w_1, \ldots, w_d .

Choose a weight vector that "fits the training data": $\mathbf{w} \in \mathbb{R}^d$ such that $y_i \approx \hat{y}_i$ for $i \in [n]$, or:

$$= \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \text{ where } \mathbf{x}_1, \dots, \mathbf{x}_d \in \mathbb{R}^n.$$

 $\mathbf{X}\mathbf{w} = \hat{\mathbf{y}} \approx \mathbf{y}$.

Least Squares OLS Theorem

<u>Theorem (Ordinary Least Squares).</u> Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ and $\mathbf{y} \in \mathbb{R}^{n}$. Let $\hat{\mathbf{w}} \in \mathbb{R}^{d}$ be the least squares minimizer:

$$\hat{\mathbf{w}} = \arg\min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$$

If $n \ge d$ and $rank(\mathbf{X}) = d$, then:

$$\hat{\mathbf{w}} = (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}.$$

To get predictions $\hat{\mathbf{y}} \in \mathbb{R}^n$:

 $\hat{\mathbf{y}} = \mathbf{X}\hat{\mathbf{w}} = \mathbf{X}(\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}.$



Least Squares Least norm exact solution

For $\mathbf{X} \in \mathbb{R}^{n \times d}$ with $rank(\mathbf{X}) = n$,

minimize $\mathbf{w} \in \mathbb{R}^d$

- $\|\mathbf{W}\|$
- subject to Xw = y
- We already know how to solve this use the pseudoinverse!

Least Squares Least norm exact solution

For $\mathbf{X} \in \mathbb{R}^{n \times d}$ with $\operatorname{rank}(\mathbf{X}) = n_{i}$

minimize $\|\mathbf{W}\|$ $\mathbf{w} \in \mathbb{R}^d$ subject to Xw = y

<u>Theorem (Minimum norm least squares solution).</u> Let $\mathbf{X} \in \mathbb{R}^{n \times d}$, let $d \ge n$, and let $rank(\mathbf{X}) = n$. Then, $\hat{\mathbf{w}} = \mathbf{X}^+ \mathbf{y} = \mathbf{V} \mathbf{\Sigma}^+ \mathbf{U}^\top \mathbf{y}$ is the exact solution $\mathbf{X} \hat{\mathbf{w}} = \mathbf{y}$ with smallest Euclidean norm:

 $\|\mathbf{w}\|_2^2 \ge \|\hat{\mathbf{w}}\|_2^2$ for all $\mathbf{w} \in \mathbb{R}^d$.





Least Squares Least norm exact solution

minimize $\mathbf{w} \in \mathbb{R}^d$

Alternate proof (through Lagrangian). For Lagrange multipliers $\lambda \in \mathbb{R}^n$,

$$L(\mathbf{w},\lambda) =$$

First-order conditions: $\nabla_{\mathbf{w}} L(\mathbf{w}, \lambda) = 2\mathbf{w} + \mathbf{X}^{\mathsf{T}} \lambda$ and $\nabla_{\lambda} L(\mathbf{w}, \lambda) = \mathbf{X}\mathbf{w} - \mathbf{y}$.

Setting equal to zero: $2\mathbf{w} + \mathbf{X}^{\mathsf{T}}\lambda = \mathbf{0}$ and $\mathbf{X}\mathbf{w} - \mathbf{y} = \mathbf{0} \Longrightarrow \mathbf{w} = -\frac{1}{2}\mathbf{X}^{\mathsf{T}}\lambda$ and $\mathbf{X}\mathbf{w} = \mathbf{y}$

Solve for
$$\lambda$$
: $\mathbf{X}\mathbf{w} = -\frac{1}{2}\mathbf{X}\mathbf{X}^{\mathsf{T}}\lambda \implies -\frac{1}{2}(\mathbf{X}\mathbf{X}^{\mathsf{T}})\lambda = \mathbf{y} \implies \mathbf{x}$

Plug λ back in to solve for \mathbf{w} : $\mathbf{w} = -\frac{1}{2}\mathbf{X}^{\mathsf{T}}\lambda = -\frac{1}{2}\mathbf{X}^{\mathsf{T}}\left(-2(\mathbf{X}\mathbf{X}^{\mathsf{T}})^{-1}\mathbf{y}\right) \implies \mathbf{w} = \mathbf{X}^{\mathsf{T}}(\mathbf{X}\mathbf{X}^{\mathsf{T}})^{-1}\mathbf{y} = \mathbf{X}^{+}\mathbf{y}$. The pseudoinverse!

 $||\mathbf{W}||$ subject to Xw = y $\|\mathbf{w}\| + \lambda^{\top} (\mathbf{X}\mathbf{w} - \mathbf{v})$ $\lambda = -2(\mathbf{X}\mathbf{X}^{\mathsf{T}})^{-1}\mathbf{y}.$
Least Squares Least norm exact solution

For $\mathbf{X} \in \mathbb{R}^{n \times d}$ with $\operatorname{rank}(\mathbf{X}) = n_{i}$

minimize $\mathbf{w} \in \mathbb{R}^d$

Then, $\hat{\mathbf{w}} = \mathbf{X}^+ \mathbf{y} = \mathbf{V} \mathbf{\Sigma}^+ \mathbf{U}^\top \mathbf{y}$ is the exact solution $\mathbf{X} \hat{\mathbf{w}} = \mathbf{y}$ with smallest Euclidean norm:

subject to Xw = y

- <u>Theorem (Minimum norm least squares solution).</u> Let $\mathbf{X} \in \mathbb{R}^{n \times d}$, let $d \ge n$, and let $rank(\mathbf{X}) = n$.
 - $\|\mathbf{w}\|_2^2 \ge \|\hat{\mathbf{w}}\|_2^2$ for all $\mathbf{w} \in \mathbb{R}^d$.

Our goal will now be to minimize two objectives:

Writing this as an optimization problem:

 $\mathbf{w} \in \mathbb{R}^d$

where $\gamma > 0$ is a fixed tuning parameter.

This optimization problem is known as <u>ridge/Tikhonov/ ℓ_2 -regularized regression</u>.

- $\|Xw v\|^2$ and $\|w\|^2$.

minimize $\|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \gamma \|\mathbf{w}\|^2$

Our goal will now be to minimize two objectives: $\|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$ and $\|\mathbf{w}\|^2$. Writing this as an optimization problem:

 $\begin{array}{ll} \text{minimize} & \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \gamma \|\mathbf{w}\|^2\\ \mathbf{w} \in \mathbb{R}^d \end{array}$

where $\gamma > 0$ is a fixed tuning parameter.

This optimization problem is known as <u>ridge/</u> <u>Tikhonov/ ℓ_2 -regularized regression.</u>



Our goal will now be to minimize two objectives: $\|Xw - y\|^2$ and $\|w\|^2$. Writing this as an optimization problem:

> minimize $\|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \gamma \|\mathbf{w}\|^2$ $\mathbf{w} \in \mathbb{R}^d$

where $\gamma > 0$ is a fixed tuning parameter.

This optimization problem is known as <u>ridge/</u> <u>Tikhonov/*l*₂-regularized regression.</u>





Ridge Regression Property: PSD to PD matrices

 $\mathbf{w} \in \mathbb{R}^d$

Property (Perturbing PSD matrices). Let $\mathbf{A} \in \mathbb{R}^{d \times d}$ be a positive semidefinite matrix. Then, for any $\gamma > 0$, the matrix $\mathbf{A} + \gamma \mathbf{I}$ is positive definite.

Proof. Let $\mathbf{v} \in \mathbb{R}^d$ be any vector. $\mathbf{v}^{\mathsf{T}}(\mathbf{A} + \gamma \mathbf{I})\mathbf{v}$

 $= \mathbf{v} \cdot \mathbf{A} \mathbf{v}$

 ≥ 0

minimize $\|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \gamma \|\mathbf{w}\|^2$

How do we solve this using the first and second order conditions?

$$= \mathbf{v}^{\mathsf{T}} (\mathbf{A}\mathbf{v} + \gamma \mathbf{v}) = \mathbf{v}^{\mathsf{T}} \mathbf{A}\mathbf{v} + \gamma \mathbf{v}^{\mathsf{T}} \mathbf{v}$$
$$+ \gamma ||\mathbf{v}||^{2}$$
$$\to 0 \text{ unless } \mathbf{v} = \mathbf{0}.$$

Ridge Regression First-order conditions

 $\mathbf{w} \in \mathbb{R}^d$

Take the gradient and set to **0**:

By property (perturbing PSD matrices), $\mathbf{X}^{\mathsf{T}}\mathbf{X} + \gamma \mathbf{I}$ is PD, so:

minimize $\|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \gamma \|\mathbf{w}\|^2$

$\nabla_{\mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \nabla_{\mathbf{w}} \|\mathbf{w}\|^2 = 2\mathbf{X}^{\mathsf{T}}\mathbf{X}\mathbf{w} - 2\mathbf{X}^{\mathsf{T}}\mathbf{y} + 2\gamma\mathbf{w}$ $2\mathbf{X}^{\mathsf{T}}\mathbf{X}\mathbf{w} - 2\mathbf{X}^{\mathsf{T}}\mathbf{y} + 2\gamma\mathbf{w} = \mathbf{0} \implies (\mathbf{X}^{\mathsf{T}}\mathbf{X} + \gamma\mathbf{I})\mathbf{w} = \mathbf{X}^{\mathsf{T}}\mathbf{y}$

$\mathbf{w}^* = (\mathbf{X}^\top \mathbf{X} + \gamma \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}.$

Least Squares Solving ridge regression

 $\mathbf{w} \in \mathbb{R}^d$

Candidate minimizer: $\mathbf{w}^* = (\mathbf{X}^\top \mathbf{X} + \gamma \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$. Gradient: $\nabla_{\mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \nabla_{\mathbf{w}} \|\mathbf{w}\|^2 = 2\mathbf{X}^{\mathsf{T}}\mathbf{X}\mathbf{w} - 2\mathbf{X}^{\mathsf{T}}\mathbf{y} + 2\gamma\mathbf{w}$ Taking the Hessian,

Sufficient condition for optimality applies!

minimize $\|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \gamma \|\mathbf{w}\|^2$

 $\nabla^2 f(\mathbf{w}) = \mathbf{X}^{\mathsf{T}} \mathbf{X} + \gamma \mathbf{I}$, which is positive definite.

Ridge Regression Theorem

<u>Theorem (Ridge Regression).</u> Let $\mathbf{X} \in \mathbb{R}^{n \times d}$, $\mathbf{y} \in \mathbb{R}^{n}$, and $\gamma > 0$. Then, $\mathbf{w} \in \mathbb{R}^d$

has the form:

To get predictions $\hat{\mathbf{y}} \in \mathbb{R}^n$:

$\hat{\mathbf{y}} = \mathbf{X}\hat{\mathbf{w}} = \mathbf{X}(\mathbf{X}^{\mathsf{T}}\mathbf{X} + \gamma\mathbf{I})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}.$

$\hat{\mathbf{w}} = (\mathbf{X}^{\mathsf{T}}\mathbf{X} + \gamma \mathbf{I})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}.$

- $\hat{\mathbf{w}} = \arg \min \|\mathbf{X}\mathbf{w} \mathbf{y}\|^2 + \gamma \|\mathbf{w}\|^2$

Least Squares Comparison with ridge solution

<u>Theorem (Ridge Regression).</u> Let $\mathbf{X} \in \mathbb{R}^{n \times d}$, $\mathbf{y} \in \mathbb{R}^n$, and $\gamma > 0$. Then, the ridge minimizer:

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \gamma \|\mathbf{w}\|^2$$

has the form:

$$\hat{\mathbf{w}} = (\mathbf{X}^{\mathsf{T}}\mathbf{X} + \gamma \mathbf{I})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}.$$

To get predictions $\hat{\mathbf{y}} \in \mathbb{R}^n$:

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\mathbf{w}} = \mathbf{X}(\mathbf{X}^{\mathsf{T}}\mathbf{X} + \gamma\mathbf{I})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}.$$

Theorem (Ordinary Least Squares). Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ and $\mathbf{y} \in \mathbb{R}^{n}$. Let $\hat{\mathbf{w}} \in \mathbb{R}^{d}$ be the least squares minimizer:

> $\hat{\mathbf{w}} = \arg \min \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$ $\mathbf{w} \in \mathbb{R}^d$

If $n \ge d$ and $rank(\mathbf{X}) = d$, then:

$$\hat{\mathbf{w}} = (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}.$$

To get predictions $\hat{\mathbf{y}} \in \mathbb{R}^n$:

 $\hat{\mathbf{y}} = \mathbf{X}\hat{\mathbf{w}} = \mathbf{X}(\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}.$



But $\hat{\mathbf{y}}$ might not be a perfect fit to \mathbf{y} !

Model this using a true weight vector $\mathbf{w}^* \in \mathbb{R}^d$ and an error term $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_n) \in \mathbb{R}^n$.

Choose a weight vector that "fits the training data": $\hat{\mathbf{w}} \in \mathbb{R}^d$ such that $y_i \approx \hat{y}_i$ for $i \in [n]$, or: $\mathbf{X}\hat{\mathbf{w}} = \hat{\mathbf{y}} \approx \mathbf{y}$.

- $y_i = \mathbf{x}_i^{\mathsf{T}} \mathbf{w}^* + \epsilon_i$ for all $i \in [n]$
 - $\mathbf{y} = \mathbf{X}\mathbf{w}^* + \epsilon$

True labels: $\mathbf{y} = \mathbf{X}\mathbf{w}^* + \epsilon$.

What happens when we use the OLS weights $\hat{\mathbf{w}} = (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}$?

 $\hat{\mathbf{w}} = (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}$ $= (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}(\mathbf{X}\mathbf{w}^{*} + \epsilon)$ $= \mathbf{w}^* + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \epsilon$

 $= (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{X}\mathbf{w}^{*} + (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\epsilon$

True labels: $\mathbf{y} = \mathbf{X}\mathbf{w}^* + \epsilon$.

What happens when we use the OLS weights $\hat{\mathbf{w}} = (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}$?

 $\hat{\mathbf{w}} = (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}$ $= (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}(\mathbf{X}\mathbf{w}^{*} + \epsilon)$ $= \mathbf{w}^* + (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\boldsymbol{\epsilon}$

When $\epsilon = 0$ (y is linearly related to X), this is perfect: $\hat{\mathbf{w}} = \mathbf{w}^*$!

 $= (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{X}\mathbf{w}^{*} + (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\boldsymbol{\epsilon}$

True labels: $\mathbf{y} = \mathbf{X}\mathbf{w}^* + \epsilon$.

What happens when we use the OLS weights $\hat{\mathbf{w}} = (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}$?

 $\hat{\mathbf{w}} = (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}$ $= (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}(\mathbf{X}\mathbf{w}^{*} + \epsilon)$ $= \mathbf{w}^* + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \epsilon$

When $\epsilon \neq 0$, we are off by $\hat{\mathbf{w}} - \mathbf{w}^* = (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\epsilon$.

- $= (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{X}\mathbf{w}^{*} + (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\boldsymbol{\epsilon}$

Error in (OLS) Regression **Eigendecomposition perspective**

Weight vector's error: $\hat{\mathbf{w}} - \mathbf{w}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \epsilon$. We know that $\mathbf{X}^{\mathsf{T}}\mathbf{X}$ (the covariance matrix) is PSD, so it is diagonalizable: $\mathbf{X}^{\mathsf{T}}\mathbf{X} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^{\mathsf{T}} =$

The inverse of the diagonal matrix Λ^{-1} :

$$\boldsymbol{\Lambda}^{-1} = \begin{bmatrix} 1/\lambda_1 & \dots & 0\\ \vdots & \ddots & \vdots\\ 0 & \dots & 1/\lambda_d \end{bmatrix}, \text{ so}$$

$$\Rightarrow (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1} = \mathbf{V}^{\mathsf{T}}\mathbf{\Lambda}^{-1}\mathbf{V}.$$

if λ_i is small, the entries of $\hat{\mathbf{w}}$ blow up!

Error in Regression Error using ridge regression

True labels: $\mathbf{y} = \mathbf{X}\mathbf{w}^* + \epsilon$.

What happens when we use the <u>ridge regression weights</u> $\hat{\mathbf{w}} = (\mathbf{X}^{\mathsf{T}}\mathbf{X} + \gamma \mathbf{I})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}$?

 $\hat{\mathbf{w}} = (\mathbf{X}^{\mathsf{T}}\mathbf{X} + \gamma \mathbf{I})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}$ $= (\mathbf{X}^{\mathsf{T}}\mathbf{X} + \gamma \mathbf{I})^{-1}\mathbf{X}^{\mathsf{T}}(\mathbf{X}\mathbf{w}^{*} + \epsilon)$ $= (\mathbf{X}^{\mathsf{T}}\mathbf{X} + \gamma \mathbf{I})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{X}\mathbf{w}^{*} + (\mathbf{X}^{\mathsf{T}}\mathbf{X} + \gamma \mathbf{I})^{-1}\mathbf{X}^{\mathsf{T}}\epsilon$

When $\epsilon = 0$ (y is linearly related to X), this is no longer perfect:

$\hat{\mathbf{w}} = (\mathbf{X}^{\mathsf{T}}\mathbf{X} + \gamma \mathbf{I})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{X}\mathbf{w}^{*}$, but...

Error in Regression Error using ridge regression

True labels: $\mathbf{y} = \mathbf{X}\mathbf{w}^* + \epsilon$.

What happens when we use the <u>ridge regression weights</u> $\hat{\mathbf{w}} = (\mathbf{X}^{\mathsf{T}}\mathbf{X} + \gamma \mathbf{I})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}$?

$$\hat{\mathbf{w}} = (\mathbf{X}^{\mathsf{T}}\mathbf{X} + \gamma \mathbf{I})^{-1}\mathbf{X}$$
$$= (\mathbf{X}^{\mathsf{T}}\mathbf{X} + \gamma \mathbf{I})^{-1}\mathbf{X}$$
$$= (\mathbf{X}^{\mathsf{T}}\mathbf{X} + \gamma \mathbf{I})^{-1}\mathbf{X}$$

When $\epsilon \neq 0$, we have more stable errors!

- **V** $\mathbf{X}^{\mathsf{T}}(\mathbf{X}\mathbf{w}^* + \epsilon)$
- $\mathbf{X}^{\mathsf{T}}\mathbf{X}\mathbf{w}^{*} + (\mathbf{X}^{\mathsf{T}}\mathbf{X} + \gamma\mathbf{I})^{-1}\mathbf{X}^{\mathsf{T}}\boldsymbol{\epsilon}$

Error in Ridge Regression Eigendecomposition perspective

Ridge weights: $\hat{\mathbf{w}} = (\mathbf{X}^{\mathsf{T}}\mathbf{X} + \gamma \mathbf{I})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}$.

We know that $\mathbf{X}^{\mathsf{T}}\mathbf{X}$ is positive semidefinite, so it is diagonalizable:

 $\mathbf{X}^{\mathsf{T}}\mathbf{X} + \gamma \mathbf{I} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^{\mathsf{T}} + \mathbf{V}(\gamma \mathbf{I})\mathbf{V}^{\mathsf{T}} =$

The inverse of the diagonal matrix $(\Lambda + \gamma I)^{-1}$:

$$\mathbf{\Lambda} + \gamma \mathbf{I})^{-1} = \begin{bmatrix} \frac{1}{\lambda_1 + \gamma} & \dots & 0\\ \vdots & \ddots & \vdots\\ 0 & \dots & \frac{1}{\lambda_d + \gamma} \end{bmatrix},$$

$$\implies (\mathbf{X}^{\mathsf{T}}\mathbf{X} + \gamma \mathbf{I})^{-1} = \mathbf{V}^{\mathsf{T}}(\mathbf{\Lambda} + \gamma \mathbf{I})^{-1}\mathbf{V}.$$

so
$$\frac{1}{\lambda_i + \gamma}$$
 entries are never bigger than $\frac{1}{\gamma}$

<u>Theorem (Ridge Regression).</u> Let $\mathbf{X} \in \mathbb{R}^{n \times d}$, $\mathbf{y} \in \mathbb{R}^{n}$, and $\gamma > 0$. Then,

 $\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \gamma \|\mathbf{w}\|^2$

has the form:

$$\hat{\mathbf{w}} = (\mathbf{X}^{\mathsf{T}}\mathbf{X} + \gamma \mathbf{I})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}.$$

To get predictions $\hat{\mathbf{y}} \in \mathbb{R}^n$:

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\mathbf{w}} = \mathbf{X}(\mathbf{X}^{\mathsf{T}}\mathbf{X} + \gamma\mathbf{I})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}.$$



For bigger γ, bigger "constraint" ball!



Recap

Lesson Overview

Optimization. Minimize an <u>objective function</u> $f : \mathbb{R}^d \to \mathbb{R}$ with the possible requirement that the minimizer \mathbf{x}^* belongs to a constraint set $\mathscr{C} \subseteq \mathbb{R}^d$.

Lagrangian. For optimization problems with \mathscr{C} defined by equalities/inequalities, the <u>Lagrangian</u> is a function $L : \mathbb{R}^d \times \mathbb{R}^m \times \mathbb{R}^r \to \mathbb{R}$ that "unconstrains" the problem.

Unconstrained local optima. With no constraints, the standard tools of calculus give conditions for a point x^* to be optimal, at least to all points close to it.

Constrained local optima (Lagrangian and KKT). When \mathscr{C} is represented by inequalities and equalities, we can use the method of Lagrange multipliers and the KKT Theorem to "unconstrain" the problem.

Ridge regression and minimum norm solutions. By constraining the norm of $\mathbf{w}^* \in \mathbb{R}^d$ of least squares (i.e. $\|\mathbf{w}^*\|$), we obtain more "stable" solutions.

Lesson Overview Big Picture: Least Squares





unconstrained min.
 constrained min.

Lesson Overview

Big Picture: Gradient Descent



