

Machine Learning

Lecture 11: Constrained Optimization

Prof. Dr. Aleksandar Bojchevski

29.05.24

Constrained problems

Projected gradient descent

Lagrangian and duality

Non-negative least squares

In many important use cases, negative weights would not be sensible, e.g., when fitting physical properties like resistance or density.

Hence, we need an additional constraint:

$$\begin{aligned} \text{minimize}_{\mathbf{w}} \quad & \mathcal{L}_{\text{LS}}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i - y_i)^2 \\ \text{subject to} \quad & w_i \geq 0 \quad \text{for } i = 1, \dots, d. \end{aligned}$$

How do we solve this?

Constrained optimization problem: Given $f_0 : \mathbb{R}^d \rightarrow \mathbb{R}$ and $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$,

$$\text{minimize}_{\boldsymbol{\theta}} \quad f_0(\boldsymbol{\theta})$$

$$\text{subject to} \quad f_i(\boldsymbol{\theta}) \leq 0 \quad \text{for } i = 1, \dots, M.$$

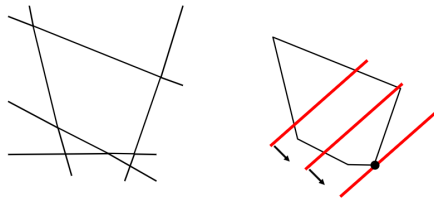
Feasibility: A point $\boldsymbol{\theta} \in \mathbb{R}^d$ is called **feasible** if and only if it satisfies the constraints of the optimization problem, i.e., $f_i(\boldsymbol{\theta}) \leq 0$ for all $i \in \{1, \dots, M\}$.

Minimum and minimizer: We call the optimal value the **minimum** p^* , and the point where the minimum is obtained the **minimizer** $\boldsymbol{\theta}^*$. Thus $p^* = f_0(\boldsymbol{\theta}^*)$.

Standard problems with inequality constraints

Linear Programming (LP)

$$\begin{array}{ll}\text{minimize}_{\theta} & \mathbf{c}^T \theta \\ \text{subject to} & \mathbf{A}\theta - \mathbf{b} \leq 0\end{array}$$



The feasible set defines a convex polytope – a convex set resulting from the intersection of half spaces.

The optimum is always at a vertex of the polytope (assuming a unique solution).¹

The [simplex algorithm](#) solves LPs by moving from vertex to vertex along edges.²

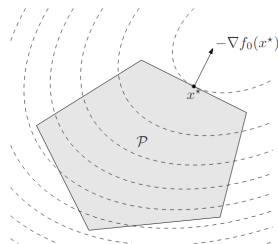
¹If there are multiple solutions, \mathbf{c} will be parallel to a face of the polytope.

²It can take exponential time in the dimension, but it is very efficient in practice. There are polynomial-time algorithms, e.g., interior point methods, which are often slower in practice.

Standard problems with inequality constraints

Quadratic Programming (QP)

$$\begin{array}{ll}\text{minimize}_{\theta} & \frac{1}{2}\theta^T Q \theta + c^T \theta \\ \text{subject to} & A\theta - b \leq 0\end{array}$$



If Q is positive semidefinite \rightarrow convex.

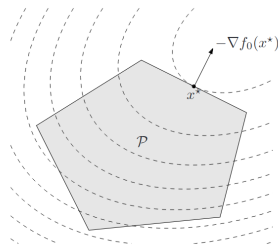
Non-negative least squares is an example of QP, with

$$\theta = w, \quad Q = XX^T, \quad c = Xy, \quad A = -I_D, \quad b = 0$$

Standard problems with inequality constraints

Quadratic Programming (QP)

$$\begin{array}{ll}\text{minimize}_{\theta} & \frac{1}{2}\theta^T Q \theta + c^T \theta \\ \text{subject to} & A\theta - b \leq 0\end{array}$$



If Q is positive semidefinite \rightarrow convex.

Non-negative least squares is an example of QP, with

$$\theta = w, \quad Q = XX^T, \quad c = Xy, \quad A = -I_D, \quad b = 0$$

We also have: [Second-order cone](#) programming, [Geometric](#) programming, ...

(Mixed) Integer Linear Programming

Integer Linear Programming (ILP)

$$\begin{aligned} & \text{minimize}_{\boldsymbol{\theta}} && \boldsymbol{c}^T \boldsymbol{\theta} \\ & \text{subject to} && \boldsymbol{A}\boldsymbol{\theta} - \boldsymbol{b} \leq \mathbf{0} \\ & && \boldsymbol{\theta} \in \mathbb{Z}^D \end{aligned}$$

We now restrict the parameters $\boldsymbol{\theta}$ to be integer (\mathbb{Z}), or often just binary.

When a subset of the variables is real-valued, it is called a mixed ILP (or MILP).

MILPS have many applications but are NP-hard to solve in general (unlike LPs).

Constrained problems

Projected gradient descent

Lagrangian and duality

Motivation for projected gradient descent

Most of the time we use gradient descent to solve optimization problems.

Can we adapt gradient descent to solve constrained problems?

Constrained optimization with gradient descent

Why can't we just directly apply gradient descent for constrained optimization?

Assume parameters are restricted to convex set \mathcal{X} (the domain). After some gradient descent steps, θ^{t+1} might be outside of region \mathcal{X} .

Even stronger: if current point θ^t is at the “border” of \mathcal{X} , each move along the gradient might lead to an infeasible solution:

$$\theta^t \in \mathcal{X}, \text{ but } \theta^{t+1} = \theta^t - \tau \nabla f(\theta^t) \notin \mathcal{X} \text{ for all } \tau > 0$$

Constrained optimization with gradient descent

Why can't we just directly apply gradient descent for constrained optimization?

Assume parameters are restricted to convex set \mathcal{X} (the domain). After some gradient descent steps, θ^{t+1} might be outside of region \mathcal{X} .

Even stronger: if current point θ^t is at the “border” of \mathcal{X} , each move along the gradient might lead to an infeasible solution:

$$\theta^t \in \mathcal{X}, \text{ but } \theta^{t+1} = \theta^t - \tau \nabla f(\theta^t) \notin \mathcal{X} \text{ for all } \tau > 0$$

Idea: project a new point back to (the closest point in) the convex set \mathcal{X}

$$\theta^{t+1} \leftarrow \pi_{\mathcal{X}}(\theta^t - \tau \nabla f(\theta^t))$$

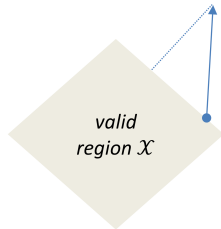
where $\pi_{\mathcal{X}}(\mathbf{p}) = \arg \min_{\theta \in \mathcal{X}} \|\theta - \mathbf{p}\|_2^2$ is the projection.

Projected gradient descent

Idea: Project a new point back on the convex set \mathcal{X}

$$\boldsymbol{\theta}^{t+1} \leftarrow \pi_{\mathcal{X}}(\boldsymbol{\theta}^t - \tau \nabla f(\boldsymbol{\theta}^t))$$

where $\pi_{\mathcal{X}}(\mathbf{p}) = \arg \min_{\boldsymbol{\theta} \in \mathcal{X}} \|\boldsymbol{\theta} - \mathbf{p}\|_2^2$.



But projection itself is a convex optimization problem!

If all constraints are linear and L_2 norm, we have a quadratic program.

- Possible to use standard solvers
- Still costly since projection has to be done after each gradient descent step

Efficient projections

Goal: $\pi_{\mathcal{X}}(\mathbf{p}) = \arg \min_{\boldsymbol{\theta} \in \mathcal{X}} \|\boldsymbol{\theta} - \mathbf{p}\|_2^2$.

Some (frequently observed) cases can be solved efficiently.

Projection onto box $\mathcal{X} = \{\boldsymbol{\theta} \in \mathbb{R}^d : l_i \leq \theta_i \leq u_i \text{ for all } i = 1, \dots, d\}$

$$(\pi_{\mathcal{X}}(\mathbf{p}))_i = \min(\max(l_i, p_i), u_i)$$

Projection onto L_2 -ball $\mathcal{X} = \{\boldsymbol{\theta} \in \mathbb{R}^d : \|\boldsymbol{\theta}\|_2 \leq c\}$

$$\pi_{\mathcal{X}}(\mathbf{p}) = \begin{cases} \mathbf{p} & \text{if } \|\mathbf{p}\|_2 \leq c \\ \frac{c}{\|\mathbf{p}\|_2} \mathbf{p} & \text{otherwise} \end{cases}$$

Goal: $\pi_{\mathcal{X}}(\mathbf{p}) = \arg \min_{\boldsymbol{\theta} \in \mathcal{X}} \|\boldsymbol{\theta} - \mathbf{p}\|_2^2$.

Some (frequently observed) cases can be solved efficiently.

Projection onto L_1 -ball $\mathcal{X} = \{\boldsymbol{\theta} \in \mathbb{R}^d : \|\boldsymbol{\theta}\|_1 \leq c\}$.³

Projection onto L_1 -ball with box-constraints

$\mathcal{X} = \{\boldsymbol{\theta} \in \mathbb{R}^d : \|\boldsymbol{\theta}\|_1 \leq c \text{ and } l_i \leq \theta_i \leq u_i \text{ for all } i = 1, \dots, d\}$.⁴

³Linear time algorithms: “Projection onto an L_1 -norm Ball with Application to Identification of Sparse Autoregressive Models” by Jitkomut Songsiri.

⁴Linear time algorithms: “ L_1 Projections with Box Constraints” by Mithun Das Gupta et al.

Projected gradient descent

Often used for solving (large-scale) constrained optimization problems.

Highly efficient if projection can be evaluated efficiently.

Each step leads to a feasible solution.⁵

Like GD, we need to choose the learning rate, etc.

Projected gradient descent is a special case of so called proximal methods.

⁵At each step, the solution remains inside the valid domain. There also exist methods that step outside the feasible region.

Constrained problems

Projected gradient descent

Lagrangian and duality

Constrained optimization problem

Given $f_0 : \mathbb{R}^d \rightarrow \mathbb{R}$ and $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$,

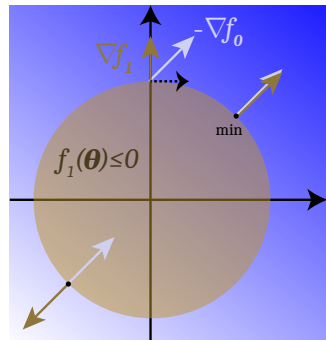
$$\text{minimize}_{\boldsymbol{\theta}} \quad f_0(\boldsymbol{\theta})$$

$$\text{subject to} \quad f_i(\boldsymbol{\theta}) \leq 0 \quad \text{for } i = 1, \dots, M.$$

Minimization with a single inequality constraint

$$\underset{\theta}{\text{minimize}} \quad \underbrace{-(\theta_1 + \theta_2)}_{=f_0(\theta)} \quad \text{s.t.} \quad \underbrace{\theta_1^2 + \theta_2^2 - 1}_{=f_1(\theta)} \leq 0$$

- a) If minimizer θ^* is in the interior of the circle ($f_1(\theta^*) < 0$) then $\nabla f_0(\theta^*) = 0$.
- b) If θ^* is on the circle ($f_1(\theta^*) = 0$) then $-\nabla f_0(\theta^*)$ and $\nabla f_1(\theta^*)$ are collinear.
- c) Otherwise, $-\nabla f_0(\theta^*) = \alpha \nabla f_1(\theta^*) + \beta \mathbf{v}$ for a tangent \mathbf{v} with $\mathbf{v}^T \nabla f_1(\theta^*) = 0$.



$f_0(\theta)$ is color coded.

If c) we can move⁶ in direction of \mathbf{v} to improve f_0 while still satisfying f_1 .

$$\begin{aligned} f_0(\theta^* + \epsilon \mathbf{v}) &\simeq f_0(\theta^*) + \epsilon \mathbf{v}^T \nabla f_0(\theta^*) = f_0(\theta^*) - \epsilon \beta \|\mathbf{v}\|_2^2 < f_0(\theta^*) \\ f_1(\theta^* + \epsilon \mathbf{v}) &\simeq f_1(\theta^*) + \epsilon \mathbf{v}^T \nabla f_1(\theta^*) = f_1(\theta^*) = 0 \rightarrow \text{feasible!} \end{aligned}$$

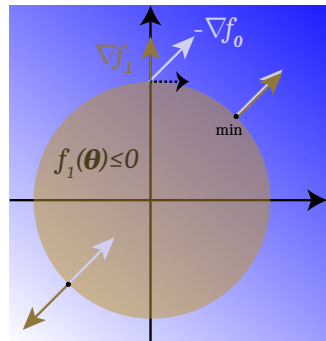
Minimization with a single inequality constraint

$$\underset{\theta}{\text{minimize}} \quad \underbrace{-(\theta_1 + \theta_2)}_{=f_0(\theta)} \quad \text{s.t.} \quad \underbrace{\theta_1^2 + \theta_2^2 - 1}_{=f_1(\theta)} \leq 0$$

Thus at the minimizer θ^* we have

$$\boxed{-\nabla f_0(\theta^*) = \alpha \nabla f_1(\theta^*)}$$

with $\alpha \geq 0$ to ensure that $-\nabla f_0$ and ∇f_1 do not point in opposite directions.



$f_0(\theta)$ is color coded.

Multiple inequality constraints

Given $f_0 : \mathbb{R}^d \rightarrow \mathbb{R}$ and $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$,

$$\begin{aligned} & \text{minimize}_{\boldsymbol{\theta}} && f_0(\boldsymbol{\theta}) \\ & \text{subject to} && f_i(\boldsymbol{\theta}) \leq 0 \quad \text{for } i = 1, \dots, M. \end{aligned}$$

For multiple constraints, we have reached the minimum when the negative gradient $-\nabla f_0$ has no component \boldsymbol{v} that is tangent to the boundary of the admissible set, that is $\boldsymbol{v}^T \nabla f_i(\boldsymbol{\theta}^*) = 0$ for $i = 1, \dots, M$.

This is the case when $-\nabla f_0$ is a linear combination of the ∇f_i 's with non-negative coefficients

$$\boxed{-\nabla f_0(\boldsymbol{\theta}^*) = \sum_{i=1}^M \alpha_i \nabla f_i(\boldsymbol{\theta}^*), \quad \alpha_i \geq 0} \tag{1}$$

$$\begin{array}{ll}\text{minimize}_{\boldsymbol{\theta}} & f_0(\boldsymbol{\theta}) \\ \text{subject to} & f_i(\boldsymbol{\theta}) \leq 0 \quad \text{for } i = 1, \dots, M\end{array}$$

We define the **Lagrangian** $L : \mathbb{R}^d \times \mathbb{R}^M \rightarrow \mathbb{R}$ associated with the above problem as

$$L(\boldsymbol{\theta}, \boldsymbol{\alpha}) = f_0(\boldsymbol{\theta}) + \sum_{i=1}^M \alpha_i f_i(\boldsymbol{\theta})$$

where $\alpha_i \geq 0$ is the **Lagrange multiplier** associated with the constraint $f_i(\boldsymbol{\theta}) \leq 0$.

Setting $\nabla_{\boldsymbol{\theta}^*} L(\boldsymbol{\theta}^*, \boldsymbol{\alpha}) = 0$ recovers the optimality criterion (Eq. 1) for $\boldsymbol{\theta}^*$,

$$\nabla_{\boldsymbol{\theta}^*} L(\boldsymbol{\theta}^*, \boldsymbol{\alpha}) = \nabla f_0(\boldsymbol{\theta}^*) + \sum_{i=1}^M \alpha_i \nabla f_i(\boldsymbol{\theta}^*) = 0.$$

Lagrange dual function

The **Lagrange dual function** $g : \mathbb{R}^M \rightarrow \mathbb{R}$ maps α to the minimum of the Lagrangian over θ (possibly $-\infty$ for some values of α),

$$g(\alpha) = \min_{\theta \in \mathbb{R}^d} L(\theta, \alpha) = \min_{\theta \in \mathbb{R}^d} \left(f_0(\theta) + \sum_{i=1}^M \alpha_i f_i(\theta) \right).$$

It is concave in α since it is the point-wise minimum of a family of affine functions of α .⁷

⁷See Boyd p. 216.

Interpretation of the Lagrangian

$$\begin{array}{ll}\text{minimize}_{\boldsymbol{\theta}} & f_0(\boldsymbol{\theta}) \\ \text{subject to} & f_i(\boldsymbol{\theta}) \leq 0, \quad i = 1, \dots, M\end{array}$$

$$L(\boldsymbol{\theta}, \boldsymbol{\alpha}) = f_0(\boldsymbol{\theta}) + \sum_{i=1}^M \alpha_i f_i(\boldsymbol{\theta})$$

For every choice of $\boldsymbol{\alpha}$, the corresponding *unconstrained* $g(\boldsymbol{\alpha}) = \min_{\boldsymbol{\theta} \in \mathbb{R}^d} L(\boldsymbol{\theta}, \boldsymbol{\alpha})$ is a **lower bound** on the optimal value of the constrained problem:

$$\min_{\substack{\boldsymbol{\theta} \in \mathbb{R}^d \\ f_i(\boldsymbol{\theta}) \leq 0}} f_0(\boldsymbol{\theta}) = f_0(\boldsymbol{\theta}^*) \geq \underbrace{f_0(\boldsymbol{\theta}^*) + \sum_{i=1}^M \underbrace{\alpha_i}_{\geq 0} \underbrace{f_i(\boldsymbol{\theta}^*)}_{\leq 0}}_{\leq 0} = L(\boldsymbol{\theta}^*, \boldsymbol{\alpha}) \geq \underbrace{\min_{\boldsymbol{\theta} \in \mathbb{R}^d} L(\boldsymbol{\theta}, \boldsymbol{\alpha})}_{g(\boldsymbol{\alpha})}$$

Hence, $f_0(\boldsymbol{\theta}^*) \geq g(\boldsymbol{\alpha})$ for $\forall \boldsymbol{\alpha}$.

Lagrange dual problem

For each $\alpha \geq 0$, the Lagrange dual function $g(\alpha)$ gives us a **lower bound** on the optimal value p^* of the original optimization problem.

What is the best (highest) lower bound?

Lagrange dual problem

$$\begin{array}{ll}\text{maximize}_{\alpha} & g(\alpha) \\ \text{subject to} & \alpha_i \geq 0, \quad i = 1, \dots, m\end{array}$$

We call α **feasible** if and only if all $\alpha_i \geq 0$ and $L(\theta, \alpha)$ is bounded from below⁸ for $\theta \in \mathbb{R}^d$.

⁸Since we are maximizing g , we are not interested in dual multipliers α such that $g(\alpha) = -\infty$, so the condition $g(\alpha) \neq -\infty$ is usually added as an additional constraint to the dual problem.

Since for all $\alpha \geq 0$ it holds that $g(\alpha) \leq p^*$ we have **weak duality (always)**,

$$d^* \leq p^* .$$

The difference $p^* - d^* \geq 0$ between the solution of the original and the dual problem is called the **duality gap**.

Under certain conditions we have **strong duality** where

$$d^* = p^* ,$$

i.e. the maximum to the Lagrange dual problem is the minimum of the original (primal) constrained optimization problem ($f_0(\theta^*) = g(\alpha^*)$).

Dual solution

Let θ^* be a minimizer of the primal problem and α^* a maximizer of the dual problem. If strong duality holds and $L(\theta, \alpha)$ is convex in θ , then

$$\theta^* = \arg \min_{\theta} L(\theta, \alpha^*).$$

Proof: Recall the proof of $g(\alpha) \leq p^*$ from Slide 21, now set $\alpha = \alpha^*$ and use that $g(\alpha^*) = p^*$ is given (strong duality) \Rightarrow every “ \leq ” in the chain has to be a “ $=$ ”.

$$\underbrace{\min_{\theta \in \mathbb{R}^d} L(\theta, \alpha^*)}_{=g(\alpha^*)} \leq L(\theta^*, \alpha^*) = f_0(\theta^*) + \sum_{i=1}^M \alpha_i^* f_i(\theta^*) \leq f_0(\theta^*) = \underbrace{\min_{f_i(\theta) \leq 0} f_0(\theta)}_{=p^*}$$

Under strong duality the first inequality becomes an equality and yields the claim.

Constraint qualifications for convex problems

Consider the constrained optimization problem,

$$\begin{array}{ll}\text{minimize}_{\boldsymbol{\theta}} & f_0(\boldsymbol{\theta}) \\ \text{subject to} & f_i(\boldsymbol{\theta}) \leq 0 \quad \text{for } i = 1, \dots, M.\end{array}$$

Slater's constraint qualification: The duality gap is zero (i.e. strong duality holds) if f_0, f_1, \dots, f_M are **convex** and there exists a **feasible** $\boldsymbol{\theta} \in \mathbb{R}^d$ such that for every constraint $i = 1, \dots, M$ it holds, either

- the constraint is affine, that is $f_i(\boldsymbol{\theta}) = \mathbf{w}_i^T \boldsymbol{\theta} + b_i$, or
- the constraint is satisfied with “<”, that is $f_i(\boldsymbol{\theta}) \prec 0$
i.e., for the non-affine constraints we require a strict inequality.

Many other theorems with different conditions that guarantee zero duality gap.

Recipe for solving constrained optimization problems

The constrained optimization problem,

$$\text{minimize}_{\boldsymbol{\theta}} \quad f_0(\boldsymbol{\theta}) \quad \text{s.t.} \quad f_i(\boldsymbol{\theta}) \leq 0 \quad \text{for } i = 1, \dots, M,$$

with all f_0, f_1, \dots, f_M **convex** can be approached as follows:

Recipe for solving constrained optimization problems

The constrained optimization problem,

$$\text{minimize}_{\boldsymbol{\theta}} \quad f_0(\boldsymbol{\theta}) \quad \text{s.t.} \quad f_i(\boldsymbol{\theta}) \leq 0 \quad \text{for } i = 1, \dots, M,$$

with all f_0, f_1, \dots, f_M **convex** can be approached as follows:

1. Formulate the **Lagrangian** $L(\boldsymbol{\theta}, \boldsymbol{\alpha}) = f_0(\boldsymbol{\theta}) + \sum_{i=1}^M \alpha_i f_i(\boldsymbol{\theta})$.
2. For each $\boldsymbol{\alpha} \geq \mathbf{0}$ obtain the **dual function** $g(\boldsymbol{\alpha})$ by solving $g(\boldsymbol{\alpha}) = \min_{\boldsymbol{\theta}} L(\boldsymbol{\theta}, \boldsymbol{\alpha})$:
 - 2.1 Figure out for which $\boldsymbol{\alpha}$ the objective is unbounded
 - 2.2 For other compute $g(\boldsymbol{\alpha})$, e.g. solve $\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}, \boldsymbol{\alpha}) = \mathbf{0}$ to get $\boldsymbol{\theta}^*(\boldsymbol{\alpha})$, then $g(\boldsymbol{\alpha}) = L(\boldsymbol{\theta}^*(\boldsymbol{\alpha}), \boldsymbol{\alpha})$
3. Solve the **dual problem**, it's maximum is also the minimum of the original problem if **Slater's condition** is satisfied.

$$\text{maximize}_{\boldsymbol{\alpha}} \quad g(\boldsymbol{\alpha}) \quad \text{s.t.} \quad \alpha_i \geq 0 \quad \text{for } i = 1, \dots, M.$$

KKT conditions

Assume that we have the following optimization problem

$$\begin{array}{ll}\text{minimize}_{\boldsymbol{\theta}} & f_0(\boldsymbol{\theta}) \\ \text{subject to} & f_i(\boldsymbol{\theta}) \leq 0 \quad \text{for } i = 1, \dots, M,\end{array}$$

with f_0 **convex** and f_1, \dots, f_M **convex** such that strong duality holds.

KKT conditions: $\boldsymbol{\theta}^*$ and $\boldsymbol{\alpha}^*$ are the optimal solutions of the constrained optimization problem and the corresponding Lagrange dual problem if and only if they satisfy the **Karush-Kuhn-Tucker (KKT)** conditions for $i = 1, \dots, M$:

$$f_i(\boldsymbol{\theta}^*) \leq 0 \qquad \text{primal feasibility,}$$

$$\alpha_i^* \geq 0 \qquad \text{dual feasibility,}$$

$$\alpha_i^* f_i(\boldsymbol{\theta}^*) = 0 \qquad \text{complementary slackness,}$$

$$\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}^*, \boldsymbol{\alpha}^*) = 0 \qquad \boldsymbol{\theta}^* \text{ minimizes Lagrangian.}$$

In PGD we project the parameters to the feasible set in each step.

Lagrange formalism reformulates constrained optimization problem into an unconstrained dual problem.

For convex objective and constraints, solving the dual problem allows to solve the primal problem (duality gap is zero) under rather weak conditions.

KKT conditions characterize the solution and sometimes allow to obtain it without minimization.

Convexity is necessary to guarantee optimality.

Main reading

- “Probabilistic Machine Learning: An Introduction” by Murphy
[ch. 8.5.1 – 8.5.5, 8.6.1]

Additional reading

- “Convex Optimization” by Boyd
[ch. 5.1 – 5.3, 5.5.3]

Some slides are based on an older version by S. Günnemann. Some figures are from Bishop and Murphy.