Optimization in Machine Learning

Optimization Problems Constrained problems





Learning goals

- Definition
- LP, QP, CP
- Ridge and Lasso
- Soft-margin SVM

CONSTRAINED OPTIMIZATION PROBLEM

 $\min_{\mathbf{x}\in\mathcal{S}} f(\mathbf{x}), \text{ with } f: \mathcal{S} \to \mathbb{R}.$

- **Constrained**, if domain S is restricted: $S \subsetneq \mathbb{R}^d$.
- **Convex** if *f* convex function and *S* convex set
- Typically ${\mathcal S}$ is defined via functions called constraints

 $\mathcal{S} := \{ \mathbf{x} \in \mathbb{R}^d \mid g_i(\mathbf{x}) \leq 0, h_j(\mathbf{x}) = 0 \ \forall \ i, j \}, \text{ where }$

- $g_i : \mathbb{R}^d \to \mathbb{R}, i = 1, ..., k$ are called inequality constraints,
- $h_j : \mathbb{R}^d \to \mathbb{R}, j = 1, ..., l$ are called equality constraints.

Equivalent formulation:

$$\begin{array}{ll} \min & f(\mathbf{x}) \\ \text{such that} & g_i(\mathbf{x}) \leq 0 \quad \text{ for } i = 1, \dots, k \\ & h_j(\mathbf{x}) = 0 \quad \text{ for } j = 1, \dots, l. \end{array}$$

× < 0 × × ×

LINEAR PROGRAM (LP)

• *f* linear s.t. linear constraints. Standard form:

 $\begin{array}{ll} \min_{\mathbf{x}\in\mathbb{R}^d} & \boldsymbol{c}^{\top}\mathbf{x} \\ \text{s.t.} & \boldsymbol{A}\mathbf{x}\geq\boldsymbol{b} \\ & \mathbf{x}\geq\mathbf{0} \end{array}$

for $\boldsymbol{c} \in \mathbb{R}^d$, $\boldsymbol{A} \in \mathbb{R}^{k \times d}$ and $\boldsymbol{b} \in \mathbb{R}^k$.



Visualization of constraints of 2D and 3D linear program (Source right figure: Wikipedia).



QUADRATIC PROGRAM (QP)

• f quadratic form s.t. linear constraints. Standard form:

$$\min_{\mathbf{x} \in \mathbb{R}^d} \quad \frac{1}{2} \mathbf{x}^\top \mathbf{A} \mathbf{x} + \mathbf{b}^\top \mathbf{x} + c$$
s.t. $\mathbf{E} \mathbf{x} \le \mathbf{f}$
 $\mathbf{G} \mathbf{x} = \mathbf{h}$

× × 0 × × ×

 $m{A} \in \mathbb{R}^{d imes d}, m{b} \in \mathbb{R}^{d}, m{c} \in \mathbb{R}, m{E} \in \mathbb{R}^{k imes d}, m{f} \in \mathbb{R}^{k}, m{G} \in \mathbb{R}^{l imes d}, m{h} \in \mathbb{R}^{l}.$



Visualization of quadratic objective (dashed) over linear constraints (grey). Source: Ma, Signal Processing Optimization Techniques, 2015.

CONVEX PROGRAM (CP)

• *f* convex, convex inequality constraints, linear equality constraints. Standard form:

$$\begin{array}{ll} \min\limits_{\mathbf{x}\in\mathbb{R}^d} & f(\mathbf{x})\\ \text{s.t.} & g_i(\mathbf{x}) \leq 0, i = 1, ..., k\\ & \mathbf{A}\mathbf{x} = \mathbf{b} \end{array}$$

for
$$\pmb{A} \in \mathbb{R}^{l imes d}$$
 and $\pmb{b} \in \mathbb{R}^{l}$.



Convex program (left) vs. nonconvex program (right). Source: Mathworks.

FURTHER TYPES



Quadratically constrained linear program (QCLP) and quadratically constrained quadratic program (QCQP). × × 0 × × ×

EXAMPLE 1: UNIT CIRCLE

min
$$f(x_1, x_2) = x_1 + x_2$$

s.t. $h(x_1, x_2) = x_1^2 + x_2^2 - 1 = 0$





f, h smooth. Problem **not convex** (S is not a convex set).

Note: If the constraint is replaced by $g(x_1, x_2) = x_1^2 + x_2^2 - 1 \le 0$, the problem is a convex program, even a quadratically constrained linear program (QCLP).

EXAMPLE 2: MAXIMUM LIKELIHOOD

Experiment: Draw *m* balls from a bag with balls of *k* different colors. Color *j* has a probability of p_j of being drawn.

The probability to realize the outcome $\mathbf{x} = (x_1, ..., x_k)$, x_j being the number of balls drawn in color *j*, is:

$$f(\mathbf{x}, m, \boldsymbol{p}) = \begin{cases} \frac{m!}{x_1! \cdots x_k!} \cdot \boldsymbol{p}_1^{x_1} \cdots \boldsymbol{p}_k^{x_k} & \text{if } \sum_{i=1}^k x_i = m \\ 0 & \text{otherwise} \end{cases}$$

The parameters p_j are subject to the following constraints:

$$0 \le p_j \le 1$$
 for all i
 $\sum_{j=1}^m p_j = 1.$



EXAMPLE 2: MAXIMUM LIKELIHOOD / 2

For a fixed *m* and a sample $\mathcal{D} = (\mathbf{x}^{(1)}, ..., \mathbf{x}^{(n)})$, where $\sum_{j=1}^{k} \mathbf{x}_{j}^{(i)} = m$ for all i = 1, ..., n, the negative log-likelihood is:

$$-\ell(\boldsymbol{p}) = -\log\left(\prod_{i=1}^{n} \frac{m!}{\mathbf{x}_{1}^{(i)}!\cdots\mathbf{x}_{k}^{(i)}!} \cdot p_{1}^{\mathbf{x}_{1}^{(i)}}\cdots p_{k}^{\mathbf{x}_{k}^{(i)}}\right)$$
$$= \sum_{i=1}^{n} \left[-\log(m!) + \sum_{j=1}^{k} \log(\mathbf{x}_{j}^{(i)}!) - \sum_{j=1}^{k} \mathbf{x}_{j}^{(i)} \log(p_{j})\right]$$
$$\propto -\sum_{i=1}^{n} \sum_{j=1}^{k} \mathbf{x}_{j}^{(i)} \log(p_{j})$$

× × ×

f, g, h are smooth.

Convex program: convex^(*) objective + box/linear constraints).

(*): log is concave, $-\log$ is convex, and the sum of convex functions is convex.

EXAMPLE 3: RIDGE REGRESSION

Ridge regression can be formulated as regularized ERM:

$$\hat{\theta}_{\mathsf{Ridge}} = \arg\min_{\boldsymbol{\theta}} \left\{ \sum_{i=1}^{n} \left(y^{(i)} - \boldsymbol{\theta}^{\top} \mathbf{x} \right)^2 + \lambda ||\boldsymbol{\theta}||_2^2 \right\}$$

Equivalently it can be written as constrained optimization problem:



f, g smooth. Convex program (convex objective, quadratic constraint).



β,

× 0 0 × 0 × ×

EXAMPLE 4: LASSO REGRESSION

Lasso regression can be formulated as regularized ERM:

$$\hat{\theta}_{\text{Lasso}} = \arg\min_{\boldsymbol{\theta}} \left\{ \sum_{i=1}^{n} \left(y^{(i)} - \boldsymbol{\theta}^{\top} \mathbf{x} \right)^{2} + \lambda ||\boldsymbol{\theta}||_{1} \right\}$$

Equivalently it can be written as constrained optimization problem:

$$\min_{\boldsymbol{\theta}} \sum_{i=1}^{n} \left(\boldsymbol{\theta}^{\top} \mathbf{x}^{(i)} - y^{(i)} \right)^{2}$$

s.t. $\|\boldsymbol{\theta}\|_{1} \leq t$

f smooth, g not smooth. Still convex program.

× o × ×

The SVM problem can be formulated in 3 equivalent ways: two primal, and one dual one (we will see later what "dual" means). Here, we only discuss the nature of the optimization problems. A more thorough statistical derivation of SVMs is given in "Supervised learning".

Formulation 1 (primal): ERM with Hinge loss



× 0 0 × 0 × ×

Formulation 2 (primal): Geometric formulation

- Find decision boundary which separates classes with **maximum** safety distance
- Distance to points closest to decision boundary ("safety margin γ ") should be **maximized**





Formulation 2 (primal): Geometric formulation

$$\min_{\boldsymbol{\theta},\boldsymbol{\theta}_0} \quad \frac{1}{2} \|\boldsymbol{\theta}\|^2$$
s.t. $y^{(i)} \left(\left\langle \boldsymbol{\theta}, \mathbf{x}^{(i)} \right\rangle + \theta_0 \right) \ge 1 \quad \forall i \in \{1, \dots, n\}$

× 0 0 × 0 × ×



Maximize safety margin γ . No point is allowed to violate safety margin constraint.

The problem is a **QP**: Quadratic objective with linear constraints.

Formulation 2 (primal): Geometric formulation (soft constraints)

$$\min_{\boldsymbol{\theta}, \boldsymbol{\theta}_0, \zeta^{(i)}} \quad \frac{1}{2} \|\boldsymbol{\theta}\|^2 + C \sum_{i=1}^n \zeta^{(i)}$$
s.t. $y^{(i)} \left(\left\langle \boldsymbol{\theta}, \mathbf{x}^{(i)} \right\rangle + \theta_0 \right) \ge 1 - \zeta^{(i)} \quad \forall i \in \{1, \dots, n\},$
and $\zeta^{(i)} \ge 0 \quad \forall i \in \{1, \dots, n\}.$



Maximize safety margin γ . Margin violations are allowed, but are minimized.

The problem is a **QP**: Quadratic objective with linear constraints.

Formulation 3 (dual): Dualizing the primal formulation

$$\begin{aligned} \max_{\boldsymbol{\alpha} \in \mathbb{R}^n} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j \boldsymbol{y}^{(i)} \boldsymbol{y}^{(j)} \left\langle \boldsymbol{x}^{(i)}, \boldsymbol{x}^{(j)} \right\rangle \\ \text{s.t.} \quad & 0 \le \alpha_i \le C \quad \forall \, i \in \{1, \dots, n\}, \quad \sum_{i=1}^n \alpha_i \boldsymbol{y}^{(i)} = 0 \end{aligned}$$

Matrix notation:

$$\begin{array}{ll} \max_{\boldsymbol{\alpha} \in \mathbb{R}^n} & \boldsymbol{\alpha}^\top \mathbf{1} - \frac{1}{2} \boldsymbol{\alpha}^\top \operatorname{diag}(\boldsymbol{y}) \, \mathbf{X}^\top \mathbf{X} \operatorname{diag}(\boldsymbol{y}) \, \boldsymbol{\alpha} \\ \text{s.t.} & \boldsymbol{0} \leq \alpha_i \leq \boldsymbol{C} \quad \forall \, i \in \{1, \dots, n\}, \quad \boldsymbol{\alpha}^\top \boldsymbol{y} = \boldsymbol{0} \end{array}$$

Kernelization: Replace dot product between **x**'s with $\mathbf{K}_{ij} = k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$, where $k(\cdot, \cdot)$ is a positive definite kernel function ($\Rightarrow \mathbf{K}$ positive semi-definite).

$$\max_{\boldsymbol{\alpha} \in \mathbb{R}^n} \quad \boldsymbol{\alpha}^\top \mathbf{1} - \frac{1}{2} \boldsymbol{\alpha} \operatorname{diag}(\boldsymbol{y}) \, \boldsymbol{K} \operatorname{diag}(\boldsymbol{y}) \, \boldsymbol{\alpha}$$

s.t. $0 \le \alpha_i \le C \quad \forall \, i \in \{1, \dots, n\}, \quad \boldsymbol{\alpha}^\top \boldsymbol{y} = 0$

This is QP with a single affine equality constraint and *n* box constraints.

× 0 0 × ×