

Machine Learning Algorithms

Support Vector Machines and duality

Caio Corro

Université Paris-Saclay, CNRS, Laboratoire Interdisciplinaire des Sciences du Numérique,
91400, Orsay, France

Lagrangian relaxation and duality

Lagrangian relaxation 1/2

Let $U \subseteq \mathbb{R}^n$ be a convex set and $f : U \rightarrow \mathbb{R}$ be a function.

Primal mathematical program:

$$(P) \quad \min_{\mathbf{u} \in U} f(\mathbf{u})$$
$$\text{s.t.} \quad \mathbf{A}\mathbf{u} \leq \mathbf{b}$$

where $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$ defines a set of m linear inequalities.

Lagrangian relaxation 1/2

Let $U \subseteq \mathbb{R}^n$ be a convex set and $f : U \rightarrow \mathbb{R}$ be a function.

Primal mathematical program:

$$(P) \quad \min_{\mathbf{u} \in U} f(\mathbf{u}) \\ \text{s.t.} \quad \mathbf{A}\mathbf{u} \leq \mathbf{b}$$

where $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$ defines a set of m linear inequalities.

Lagrangian

The Lagrangian of (P) is the function $L : U \times \mathbb{R}_+^m \rightarrow \mathbb{R}$ defined as follows:

$$L(\mathbf{u}, \boldsymbol{\lambda}) = f(\mathbf{u}) + \boldsymbol{\lambda}^\top (\mathbf{A}\mathbf{u} - \mathbf{b})$$

- ▶ $\mathbf{u} \in U$: primal variables
- ▶ $\boldsymbol{\lambda} \in \mathbb{R}_+^m$, i.e. $\boldsymbol{\lambda} \geq 0$: dual variables / Lagrangian multipliers

Lagrangian relaxation 2/2

Relaxation of the primal problem

The relaxed problem $L : \mathbb{R}_+^m \rightarrow \mathbb{R}$ is defined as:

$$L(\boldsymbol{\lambda}) = \min_{\boldsymbol{u} \in U} L(\boldsymbol{u}, \boldsymbol{\lambda})$$

Lagrangian relaxation 2/2

Relaxation of the primal problem

The relaxed problem $L : \mathbb{R}_+^m \rightarrow \mathbb{R}$ is defined as:

$$L(\boldsymbol{\lambda}) = \min_{\mathbf{u} \in U} L(\mathbf{u}, \boldsymbol{\lambda})$$

Weak Lagrangian duality (lower bound)

Let $\hat{\mathbf{u}}$ be the optimal solution of the primal problem (P). Then:

$$\forall \boldsymbol{\lambda} \in \mathbb{R}_+^m : f(\hat{\mathbf{u}}) \geq L(\boldsymbol{\lambda})$$

Proof:

$$\begin{aligned} f(\hat{\mathbf{u}}) &\geq f(\hat{\mathbf{u}}) + \boldsymbol{\lambda}^\top (\mathbf{A}\hat{\mathbf{u}} - b) \quad (\text{because } \hat{\mathbf{u}} \text{ satisfies the constraints}) \\ &\geq \min_{\mathbf{u} \in U} f(\mathbf{u}) + \boldsymbol{\lambda}^\top (\mathbf{A}\mathbf{u} - b) \\ &= L(\boldsymbol{\lambda}) \end{aligned}$$

Lagrangian dual problem

Definition

- ▶ $L(\boldsymbol{\lambda})$ is a lower bound to the primal problem, $\forall \boldsymbol{\lambda} \in \mathbb{R}_+^m : f(\hat{\mathbf{u}}) \geq L(\boldsymbol{\lambda})$
- ▶ The dual problem search for the best lower bound

$$(D) \quad \max_{\boldsymbol{\lambda} \in \mathbb{R}_+^m} L(\boldsymbol{\lambda})$$

Lagrangian dual problem

Definition

- ▶ $L(\boldsymbol{\lambda})$ is a lower bound to the primal problem, $\forall \boldsymbol{\lambda} \in \mathbb{R}_+^m : f(\hat{\mathbf{u}}) \geq L(\boldsymbol{\lambda})$
- ▶ The dual problem search for the best lower bound

$$(D) \quad \max_{\boldsymbol{\lambda} \in \mathbb{R}_+^m} L(\boldsymbol{\lambda})$$

Concavity

The dual problem is concave (no matter if f is convex or not)

Strong Lagrangian duality

Let $\boldsymbol{\lambda} \in \mathbb{R}_+^m$ be a dual feasible solution and $\bar{\mathbf{u}} = \arg \max_{\mathbf{u} \in \mathcal{U}} L(\mathbf{u}, \boldsymbol{\lambda})$. If:

- ▶ $\mathbf{A}\bar{\mathbf{u}} \leq \mathbf{b}$ (primal feasibility condition)
- ▶ and $\boldsymbol{\lambda}^\top (\mathbf{A}\bar{\mathbf{u}} - \mathbf{b}) = 0$ (complementary slackness condition)

then $\bar{\mathbf{u}} = \hat{\mathbf{u}}$ is a primal optimal solution.

Dual concavity proof 1/2

$L(\boldsymbol{\lambda})$ is concave if and only if:

1. the domain of $L(\boldsymbol{\lambda})$ is convex (trivial)
2. $\forall \boldsymbol{\lambda}^{(1)}, \boldsymbol{\lambda}^{(2)} \in \boldsymbol{\lambda} \in \mathbb{R}_+^m, \epsilon \in [0, 1] : L(\epsilon \boldsymbol{\lambda}^{(1)} + (1 - \epsilon) \boldsymbol{\lambda}^{(2)}) \geq \epsilon L(\boldsymbol{\lambda}^{(1)}) + (1 - \epsilon) L(\boldsymbol{\lambda}^{(2)})$

Dual concavity proof 1/2

$L(\boldsymbol{\lambda})$ is concave if and only if:

1. the domain of $L(\boldsymbol{\lambda})$ is convex (trivial)
2. $\forall \boldsymbol{\lambda}^{(1)}, \boldsymbol{\lambda}^{(2)} \in \boldsymbol{\lambda} \in \mathbb{R}_+^m, \epsilon \in [0, 1] : L(\epsilon \boldsymbol{\lambda}^{(1)} + (1 - \epsilon) \boldsymbol{\lambda}^{(2)}) \geq \epsilon L(\boldsymbol{\lambda}^{(1)}) + (1 - \epsilon) L(\boldsymbol{\lambda}^{(2)})$

Let $\boldsymbol{\lambda} = \epsilon \boldsymbol{\lambda}^{(1)} + (1 - \epsilon) \boldsymbol{\lambda}^{(2)}$ and $\bar{\boldsymbol{u}} = \arg \min_{\boldsymbol{u} \in \mathcal{X}} L(\boldsymbol{u}, \boldsymbol{\lambda})$:

$$\left. \begin{array}{l} L(\bar{\boldsymbol{u}}, \boldsymbol{\lambda}^{(1)}) \geq L(\boldsymbol{\lambda}^{(1)}) \\ L(\bar{\boldsymbol{u}}, \boldsymbol{\lambda}^{(2)}) \geq L(\boldsymbol{\lambda}^{(2)}) \end{array} \right\} \Rightarrow \epsilon L(\bar{\boldsymbol{u}}, \boldsymbol{\lambda}^{(1)}) + (1 - \epsilon) L(\bar{\boldsymbol{u}}, \boldsymbol{\lambda}^{(2)}) \geq \epsilon L(\boldsymbol{\lambda}^{(1)}) + (1 - \epsilon) L(\boldsymbol{\lambda}^{(2)})$$

To understand the left-hand side, remember that:

$$L(\boldsymbol{\lambda}) = \min_{\boldsymbol{u} \in \mathcal{X}} L(\boldsymbol{u}, \boldsymbol{\lambda})$$

So:

- ▶ the right-hand side has the expected form
- ▶ the left-hand side is different, we need to fix this to finish the proof

Dual concavity proof 2/2

To simplify notations, we write $c(\mathbf{u}) = \mathbf{A}\mathbf{u} - \mathbf{b}$.

The left-hand side of the inequality can be rewritten as:

$$\begin{aligned}\epsilon L(\bar{\mathbf{u}}, \boldsymbol{\lambda}^{(1)}) + (1 - \epsilon)L(\bar{\mathbf{u}}, \boldsymbol{\lambda}^{(2)}) &= \epsilon \left(f(\bar{\mathbf{u}}) + \boldsymbol{\lambda}^{(1)\top} c(\bar{\mathbf{u}}) \right) + (1 - \epsilon) \left(f(\bar{\mathbf{u}}) + \boldsymbol{\lambda}^{(2)\top} c(\bar{\mathbf{u}}) \right) \\ &= \epsilon f(\bar{\mathbf{u}}) + \epsilon \boldsymbol{\lambda}^{(1)\top} c(\bar{\mathbf{u}}) + (1 - \epsilon)f(\bar{\mathbf{u}}) + (1 - \epsilon)\boldsymbol{\lambda}^{(2)\top} c(\bar{\mathbf{u}}) \\ &= f(\bar{\mathbf{u}}) + \left(\epsilon \boldsymbol{\lambda}^{(1)} + (1 - \epsilon)\boldsymbol{\lambda}^{(2)} \right)^\top c(\bar{\mathbf{u}}) \\ &= f(\bar{\mathbf{u}}) + \boldsymbol{\lambda}^\top c(\bar{\mathbf{u}}) \\ &= L(\bar{\mathbf{u}}, \boldsymbol{\lambda}) \\ &= L(\boldsymbol{\lambda}) \\ &= L(\epsilon \boldsymbol{\lambda}^{(1)} + (1 - \epsilon)\boldsymbol{\lambda}^{(2)})\end{aligned}$$

Hence, we obtain the inequality:

$$L(\epsilon \boldsymbol{\lambda}^{(1)} + (1 - \epsilon)\boldsymbol{\lambda}^{(2)}) \geq \epsilon L(\boldsymbol{\lambda}^{(1)}) + (1 - \epsilon)L(\boldsymbol{\lambda}^{(2)})$$

which proves that the Lagrangian dual objective is concave.

Strong Lagrangian duality proof

Let $\hat{\mathbf{u}}$ be a primal optimal solution.

By weak Lagrangian duality, we know that:

$$\begin{aligned} f(\hat{\mathbf{u}}) &\geq L(\boldsymbol{\lambda}) \\ &= L(\bar{\mathbf{u}}, \boldsymbol{\lambda}) \\ &= f(\bar{\mathbf{u}}) + \boldsymbol{\lambda}^\top (\mathbf{A}\bar{\mathbf{u}} - \mathbf{b}) \end{aligned}$$

From the prerequisites (complementary slackness) the second term is null:

$$f(\hat{\mathbf{u}}) \geq f(\bar{\mathbf{u}})$$

Moreover: $f(\hat{\mathbf{u}}) \leq f(\bar{\mathbf{u}})$ because $\hat{\mathbf{u}}$ is primal feasible,
Therefore $f(\bar{\mathbf{u}}) = f(\hat{\mathbf{u}})$.

Relaxing equality constraints

$$\begin{array}{ll} \min_{\mathbf{u}} & f(\mathbf{u}) \\ \text{s.t.} & \mathbf{A}\mathbf{u} = \mathbf{b} \end{array} \quad \Leftrightarrow \quad \begin{array}{ll} \min_{\mathbf{u}} & f(\mathbf{u}) \\ \text{s.t.} & \mathbf{A}\mathbf{u} \leq \mathbf{b} \\ & -\mathbf{A}\mathbf{u} \leq -\mathbf{b} \end{array}$$

Relaxing equality constraints

$$\begin{array}{ll} \min_{\mathbf{u}} & f(\mathbf{u}) \\ \text{s.t.} & \mathbf{A}\mathbf{u} = \mathbf{b} \end{array} \quad \Leftrightarrow \quad \begin{array}{ll} \min_{\mathbf{u}} & f(\mathbf{u}) \\ \text{s.t.} & \mathbf{A}\mathbf{u} \leq \mathbf{b} \\ & -\mathbf{A}\mathbf{u} \leq -\mathbf{b} \end{array}$$

Lagrangian:
$$\begin{aligned} L(\mathbf{u}, \boldsymbol{\lambda} \geq 0, \boldsymbol{\lambda}' \geq 0) &= f(\mathbf{u}) + \boldsymbol{\lambda}^\top (\mathbf{A}\mathbf{u} - \mathbf{b}) + \boldsymbol{\lambda}'^\top (-\mathbf{A}\mathbf{u} + \mathbf{b}) \\ &= f(\mathbf{u}) + \boldsymbol{\lambda}^\top (\mathbf{A}\mathbf{u} - \mathbf{b}) - \boldsymbol{\lambda}'^\top (\mathbf{A}\mathbf{u} - \mathbf{b}) \\ &= f(\mathbf{u}) + (\boldsymbol{\lambda} - \boldsymbol{\lambda}')^\top (\mathbf{A}\mathbf{u} - \mathbf{b}) \end{aligned}$$

Let $\boldsymbol{\lambda}'' = \boldsymbol{\lambda} - \boldsymbol{\lambda}'$:

$$L(\mathbf{u}, \boldsymbol{\lambda}'') = f(\mathbf{u}) + \boldsymbol{\lambda}''^\top (\mathbf{A}\mathbf{u} - \mathbf{b})$$

In this formulation, the dual variables $\boldsymbol{\lambda}'' \in \mathbb{R}^m$ is **unconstrained**.

Benefits of relaxing equalities

Unconstrained optimization problem + Simpler strong duality condition

Optimality conditions
Constrained optimization problems

KKT conditions for a minimization problems 1/3

Primal problem (MINIMIZATION)

$$\begin{aligned} \min_{\mathbf{u} \in \mathbb{R}^k} \quad & f(\mathbf{u}) \\ \text{s.t.} \quad & g^{(i)}(\mathbf{u}) \leq 0 && \forall 1 \leq i \leq m \\ & h^{(i)}(\mathbf{u}) = 0 && \forall 1 \leq i \leq n \end{aligned}$$

- ▶ $g^{(i)} \leq 0$: a set of m inequality constraints
- ▶ $h^{(i)}(\mathbf{u}) = 0$: a set of n equality constraints

Lagrangian

$$L(\mathbf{u}, \boldsymbol{\mu}, \boldsymbol{\lambda}) = f(\mathbf{u}) + \sum_j \mu_j g^{(j)}(\mathbf{u}) + \sum_j \lambda_j h^{(j)}(\mathbf{u})$$

- ▶ $\boldsymbol{\mu} \in \mathbb{R}_+^m$: dual variables associated with primal inequalities
- ▶ $\boldsymbol{\lambda} \in \mathbb{R}^n$: dual variables associated with primal equalities

KKT conditions for a minimization problem 2/3

Lagrangian

$$L(\mathbf{u}, \boldsymbol{\mu}, \boldsymbol{\lambda}) = f(\boldsymbol{\mu}) + \sum_j \mu_j g^{(j)}(\mathbf{u}) + \sum_j \lambda_j h^{(j)}(\mathbf{u})$$

Necessary optimal condition

Any optimal primal/dual triplet $(\hat{\mathbf{u}}, \hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\lambda}})$ satisfies the following conditions:

(stationarity) $\forall i : \partial_{\hat{\mathbf{u}}} \left(f(\hat{\mathbf{u}}) + \sum_j \hat{\mu}_j g^{(j)}(\hat{\mathbf{u}}) + \sum_j \hat{\lambda}_j h^{(j)}(\hat{\mathbf{u}}) \right) \ni \mathbf{0}$

(primal feasibility) $\forall i : g^{(i)}(\hat{\mathbf{u}}) \leq 0$

$$\forall i : h^{(i)}(\hat{\mathbf{u}}) = 0$$

(dual feasibility) $\forall i : \hat{\mu}_i \geq 0$

(complementary slackness) $\sum_i \hat{\mu}_i g^{(i)}(\hat{\mathbf{u}}) = 0$

KKT conditions for a minimization problem 3/3

Warning

KKT conditions are necessary constraints:

there may exist points $\mathbf{u}, \boldsymbol{\mu}, \boldsymbol{\lambda}$ that satisfy these constraints that are not optimal solutions

Necessary and sufficient conditions

If:

- ▶ f is a concave function,
- ▶ all $g^{(i)}$ are convex functions,
- ▶ and all $h^{(i)}$ are affine functions,

then the KKT are sufficient conditions.

\Rightarrow we can solve this set of equations to find the optimal $\hat{\mathbf{u}}$, i.e. the optimization problem (may) have a closed form solution

KKT conditions for linear constraints 1/2

Primal problem

$$\begin{aligned} \min_{\mathbf{u} \in \mathbb{R}^k} \quad & f(\mathbf{u}) \\ \text{s.t.} \quad & \mathbf{A}\mathbf{u} \leq \mathbf{b} \\ & \mathbf{C}\mathbf{u} = \mathbf{d} \end{aligned}$$

Lagrangian

$$L(\mathbf{u}, \boldsymbol{\mu}, \boldsymbol{\lambda}) = f(\mathbf{u}) + \langle \boldsymbol{\mu}, \mathbf{A}\mathbf{u} - \mathbf{b} \rangle + \langle \boldsymbol{\lambda}, \mathbf{C}\mathbf{u} - \mathbf{d} \rangle$$

- ▶ $\boldsymbol{\mu} \in \mathbb{R}_+^m$: dual variables associated with primal inequalities
- ▶ $\boldsymbol{\lambda} \in \mathbb{R}^n$: dual variables associated with primal equalities

KKT conditions for linear constraints 2/2

Lagrangian

$$L(\mathbf{u}, \boldsymbol{\mu}, \boldsymbol{\lambda}) = f(\mathbf{u}) + \langle \boldsymbol{\mu}, \mathbf{A}\mathbf{u} - \mathbf{b} \rangle + \langle \boldsymbol{\lambda}, \mathbf{C}\mathbf{u} - \mathbf{d} \rangle$$

Necessary and sufficient optimality condition

An optimal primal/dual triplet $(\hat{\mathbf{u}}, \hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\lambda}})$ satisfies the following constraints:

(stationarity) $\partial_{\hat{\mathbf{u}}} \left(f(\hat{\mathbf{u}}) + \langle \hat{\boldsymbol{\mu}}, \mathbf{A}\hat{\mathbf{u}} - \mathbf{b} \rangle + \langle \hat{\boldsymbol{\lambda}}, \mathbf{C}\hat{\mathbf{u}} - \mathbf{d} \rangle \right) \ni \mathbf{0}$

(primal feasibility) $\mathbf{A}\hat{\mathbf{u}} \leq \mathbf{b}$

$$\mathbf{C}\hat{\mathbf{u}} = \mathbf{d}$$

(dual feasibility) $\hat{\boldsymbol{\mu}} \geq \mathbf{0}$

(complementary slackness) $\langle \hat{\boldsymbol{\mu}}, \mathbf{A}\hat{\mathbf{u}} - \mathbf{b} \rangle = 0$

KKT conditions for a maximization problem 1/3

Primal problem (MAXIMIZATION)

$$\begin{aligned} \max_{\mathbf{u} \in \mathbb{R}^k} \quad & f(\mathbf{u}) \\ \text{s.t.} \quad & g^{(i)}(\mathbf{u}) \geq 0 && \forall 1 \leq i \leq m \\ & h^{(i)}(\mathbf{u}) = 0 && \forall 1 \leq i \leq n \end{aligned}$$

- ▶ $g^{(i)} \geq 0$: a set of m inequality constraints
- ▶ $h^{(i)}(\mathbf{u}) = 0$: a set of n equality constraints

Lagrangian

$$L(\mathbf{u}, \boldsymbol{\mu}, \boldsymbol{\lambda}) = f(\mathbf{u}) + \sum_j \mu_j g^{(j)}(\mathbf{u}) + \lambda_j \sum_j h^{(j)}(\mathbf{u})$$

- ▶ $\boldsymbol{\mu} \in \mathbb{R}_+^m$: dual variables associated with primal inequalities
- ▶ $\boldsymbol{\lambda} \in \mathbb{R}^n$: dual variables associated with primal equalities

Fenchel duality

Let $f : \mathbb{R}^m \rightarrow \mathbb{R}$, $h : \mathbb{R}^n \rightarrow \mathbb{R}$, $\mathbf{U} \in \mathbb{R}^{m \times n}$ and assume the following problem:

$$\min_{\mathbf{v} \in \mathbb{R}^n} f(\mathbf{U}\mathbf{v}) + h(\mathbf{v})$$

Let $f : \mathbb{R}^m \rightarrow \mathbb{R}$, $h : \mathbb{R}^n \rightarrow \mathbb{R}$, $\mathbf{U} \in \mathbb{R}^{m \times n}$ and assume the following problem:

$$\min_{\mathbf{v} \in \mathbb{R}^n} f(\mathbf{U}\mathbf{v}) + h(\mathbf{v})$$

We can rewrite the problem by introducing a term $\mathbf{t} \in \mathbb{R}^m$:

$$= \min_{\mathbf{v} \in \mathbb{R}^n, \mathbf{t} \in \mathbb{R}^m} f(\mathbf{t}) + h(\mathbf{v}) \quad \text{s.t.} \quad \mathbf{t} = \mathbf{U}\mathbf{v}$$

Let $f : \mathbb{R}^m \rightarrow \mathbb{R}$, $h : \mathbb{R}^n \rightarrow \mathbb{R}$, $\mathbf{U} \in \mathbb{R}^{m \times n}$ and assume the following problem:

$$\min_{\mathbf{v} \in \mathbb{R}^n} f(\mathbf{U}\mathbf{v}) + h(\mathbf{v})$$

We can rewrite the problem by introducing a term $\mathbf{t} \in \mathbb{R}^m$:

$$= \min_{\mathbf{v} \in \mathbb{R}^n, \mathbf{t} \in \mathbb{R}^m} f(\mathbf{t}) + h(\mathbf{v}) \quad \text{s.t.} \quad \mathbf{t} = \mathbf{U}\mathbf{v}$$

We relax the constraint using dual vars $\boldsymbol{\lambda} \in \mathbb{R}^m$ and build the Lagrangian dual:

$$\geq \max_{\boldsymbol{\lambda} \in \mathbb{R}^m} \min_{\mathbf{v} \in \mathbb{R}^n, \mathbf{t} \in \mathbb{R}^m} f(\mathbf{t}) + h(\mathbf{v}) + \boldsymbol{\lambda}^\top (\mathbf{U}\mathbf{v} - \mathbf{t})$$

Let $f : \mathbb{R}^m \rightarrow \mathbb{R}$, $h : \mathbb{R}^n \rightarrow \mathbb{R}$, $\mathbf{U} \in \mathbb{R}^{m \times n}$ and assume the following problem:

$$\min_{\mathbf{v} \in \mathbb{R}^n} f(\mathbf{U}\mathbf{v}) + h(\mathbf{v})$$

We can rewrite the problem by introducing a term $\mathbf{t} \in \mathbb{R}^m$:

$$= \min_{\mathbf{v} \in \mathbb{R}^n, \mathbf{t} \in \mathbb{R}^m} f(\mathbf{t}) + h(\mathbf{v}) \quad \text{s.t.} \quad \mathbf{t} = \mathbf{U}\mathbf{v}$$

We relax the constraint using dual vars $\boldsymbol{\lambda} \in \mathbb{R}^m$ and build the Lagrangian dual:

$$\geq \max_{\boldsymbol{\lambda} \in \mathbb{R}^m} \min_{\mathbf{v} \in \mathbb{R}^n, \mathbf{t} \in \mathbb{R}^m} f(\mathbf{t}) + h(\mathbf{v}) + \boldsymbol{\lambda}^\top (\mathbf{U}\mathbf{v} - \mathbf{t})$$

We re-arrange terms as follows:

$$\begin{aligned} &= \max_{\boldsymbol{\lambda} \in \mathbb{R}^m} \left(\min_{\mathbf{t} \in \mathbb{R}^m} f(\mathbf{t}) - \boldsymbol{\lambda}^\top \mathbf{t} \right) + \left(\min_{\mathbf{v} \in \mathbb{R}^n} h(\mathbf{v}) + \boldsymbol{\lambda}^\top \mathbf{U}\mathbf{v} \right) \\ &= \max_{\boldsymbol{\lambda} \in \mathbb{R}^m} - \left(\max_{\mathbf{t} \in \mathbb{R}^m} \boldsymbol{\lambda}^\top \mathbf{t} - f(\mathbf{t}) \right) - \left(\max_{\mathbf{v} \in \mathbb{R}^n} -\boldsymbol{\lambda}^\top \mathbf{U}\mathbf{v} - h(\mathbf{v}) \right) \end{aligned}$$

Fenchel duality

Let $f : \mathbb{R}^m \rightarrow \mathbb{R}$, $h : \mathbb{R}^n \rightarrow \mathbb{R}$, $\mathbf{U} \in \mathbb{R}^{m \times n}$ and assume the following problem:

$$\min_{\mathbf{v} \in \mathbb{R}^n} f(\mathbf{U}\mathbf{v}) + h(\mathbf{v})$$

Then the Fenchel dual problem is defined as follows:

$$\max_{\boldsymbol{\lambda} \in \mathbb{R}^m} -f^*(\boldsymbol{\lambda}) - h^*(-\boldsymbol{\lambda}^\top \mathbf{U})$$

- ▶ Under certain conditions, the two problem have the same solution
- ▶ Using stationarity condition from the Lagrangian dual, we can retrieve primal-dual relationship (example later)

Support vector machines for binary classification

Binary classification

Task

Given a vector of feature values $\mathbf{x} \in \mathbb{R}^d$ we want to predict $y \in \{-1, 1\}$, i.e. either class -1 or 1 .

$$y = \begin{cases} 1 & \text{if } \langle \mathbf{a}, \mathbf{x} \rangle \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

\Rightarrow predict the class which has the same sign as $\langle \mathbf{a}, \mathbf{x} \rangle$ (the choice of setting $\langle \mathbf{a}, \mathbf{x} \rangle = 0$ to class 1 is arbitrary).

In our framework

$$s_{\mathbf{a}}(\mathbf{u}) = \langle \mathbf{a}, \mathbf{x} \rangle$$

$$\hat{y}(w) = \begin{cases} 1 & \text{if } w \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

Support Vector Machine (SVM)

SVM training problem

Given a training set D of n examples, $(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(n)}, y^{(n)})$
 \Rightarrow find an hyperplane that separates the two classes

$$\begin{aligned} \min_{\mathbf{a}} \quad & \alpha r(\mathbf{a}) \\ \text{s.t.} \quad & \langle \mathbf{a}, \mathbf{x}^{(i)} \rangle y^{(i)} \geq m \quad \forall (\mathbf{x}^{(i)}, y^{(i)}) \in D \end{aligned}$$

- ▶ m : minimum required margin, usually $m = 1$
- ▶ $r(\mathbf{a}) = 0 \Rightarrow$ any hyperplane (no objective function!)
- ▶ $r(\mathbf{a}) = \frac{1}{2} \|\mathbf{a}\|_2^2 \Rightarrow$ hyperplane that maximizes the margin
[Boser et al., 1992, Section 2.1]

SVM reformulation 1/3

Downside the previous SVM formulation

- ▶ constrained optimization is difficult
- ▶ what about non separable data?

SVM reformulation 1/3

Downside the previous SVM formulation

- ▶ constrained optimization is difficult
- ▶ what about non separable data?

SVM with slack variables

$$\begin{aligned} \min_{\mathbf{a}, \epsilon} \quad & \alpha r(\mathbf{a}) + \sum_{i=1}^n \epsilon_i \\ \text{s.t.} \quad & \langle \mathbf{a}, \mathbf{x}^{(i)} \rangle y^{(i)} \geq m - \epsilon_i \quad (\mathbf{x}^{(i)}, y^{(i)}) \in D \\ & \epsilon \geq 0 \end{aligned}$$

- ▶ $\epsilon \in \mathbb{R}^n$: the vector of slack variables, one slack variable per datapoint
- ▶ one slack variable per datapoint
- ▶ optimize both over \mathbf{a} and ϵ .

SVM reformulation 2/3

SVM with slack variables

$$\begin{aligned} \min_{\mathbf{a}, \epsilon} \quad & \alpha r(\mathbf{a}) + \sum_{i=1}^n \epsilon_i \\ \text{s.t.} \quad & \langle \mathbf{a}, \mathbf{x}^{(i)} \rangle y^{(i)} \geq m - \epsilon_i \quad (\mathbf{x}^{(i)}, y^{(i)}) \in D \\ & \epsilon \geq 0 \end{aligned}$$

Constraints reformulation

$$\begin{aligned} \langle \mathbf{a}, \mathbf{x}^{(i)} \rangle y^{(i)} &\geq m - \epsilon_i & \forall (\mathbf{x}^{(i)}, y^{(i)}) \in D \\ \langle \mathbf{a}, \mathbf{x}^{(i)} \rangle y^{(i)} - m &\geq -\epsilon_i & \forall (\mathbf{x}^{(i)}, y^{(i)}) \in D \\ m - \langle \mathbf{a}, \mathbf{x}^{(i)} \rangle y^{(i)} &\leq \epsilon_i & \forall (\mathbf{x}^{(i)}, y^{(i)}) \in D \\ \max(0, m - \langle \mathbf{a}, \mathbf{x}^{(i)} \rangle y^{(i)}) &= \epsilon_i & \forall (\mathbf{x}^{(i)}, y^{(i)}) \in D \end{aligned}$$

SVM reformulation 3/3

SVM with slack variables

$$\min_{\mathbf{a}, \epsilon} \alpha r(\mathbf{a}) + \sum_{i=1}^n \epsilon_i$$

$$\text{s.t. } \langle \mathbf{a}, \mathbf{x}^{(i)} \rangle y^{(i)} \geq m - \epsilon_i \quad \Leftrightarrow \quad \max(0, m - \langle \mathbf{a}, \mathbf{x}^{(i)} \rangle y^{(i)}) = \epsilon_i \quad \forall (\mathbf{x}^{(i)}, y^{(i)}) \in D$$
$$\epsilon \geq 0$$

Unconstrained SVM training problem

Remark

$$\min_{\mathbf{a}} \sum_{i=1}^n \max(0, m - \langle \mathbf{a}, \mathbf{x}^{(i)} \rangle y^{(i)}) + \alpha r(\mathbf{a})$$
$$= \min_{\mathbf{a}} \sum_{i=1}^n \ell_m(y^{(i)}, \langle \mathbf{a}, \mathbf{x}^{(i)} \rangle) + \alpha r(\mathbf{a})$$

- ▶ if $m > 0$: hinge loss (usually set $m = 1$)
- ▶ if $m = 0$: perceptron loss

where $\ell_m(y, w) = \max(0, m - yw)$ is the binary hinge loss function, parameterized by the margin parameter $m \geq 0$

Support vector machines for multiclass classification

Multiclass classification

Task

- ▶ Predicting a 1-in- k class given a set of feature values \mathbf{x}
- ▶ Usually $k > 2$, but works for $k = 2$

In our framework

Let d be the number of features and k the number of classes.

- ▶ Scoring function: $s(\mathbf{x}) = \mathbf{Ax} + \mathbf{b}$ where $\mathbf{A} \in \mathbb{R}^{k \times d}$, $\mathbf{b} \in \mathbb{R}^d$
- ▶ Non-probabilistic prediction function: $\hat{y}(\mathbf{w}) = \arg \max_{\mathbf{y} \in E(k)} \langle \mathbf{w}, \mathbf{y} \rangle$
- ▶ Probabilistic prediction function: $\hat{y}(\mathbf{w}) = \text{softmax}(\mathbf{w})$

Support Vector Machine (SVM)

Given a training set D of n datapoints $(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), (\mathbf{x}^{(2)}, \mathbf{y}^{(2)}), \dots, (\mathbf{x}^{(n)}, \mathbf{y}^{(n)})$ the SVM training problem is defined as follows:

$$\arg \min_{\mathbf{A}} \alpha r(\mathbf{A})$$

$$\text{s.t. } \langle \mathbf{y}', \mathbf{A}\mathbf{x}^{(i)} \rangle + m \leq \langle \mathbf{y}^{(i)}, \mathbf{A}\mathbf{x}^{(i)} \rangle \quad \forall (\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \in D, \mathbf{y}' \in \mathcal{Y}(k) \setminus \{\mathbf{y}^{(i)}\}$$

- ▶ $k - 1$ constraints per training point with this formulation
- ▶ if highest scoring class that is not the gold class satisfies the constraint, all other non-gold classes will also satisfy it

Support Vector Machine (SVM)

Given a training set D of n datapoints $(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), (\mathbf{x}^{(2)}, \mathbf{y}^{(2)}), \dots, (\mathbf{x}^{(n)}, \mathbf{y}^{(n)})$ the SVM training problem is defined as follows:

$$\begin{aligned} \arg \min_{\mathbf{A}} \quad & \alpha r(\mathbf{A}) \\ \text{s.t.} \quad & \langle \mathbf{y}', \mathbf{A}\mathbf{x}^{(i)} \rangle + m \leq \langle \mathbf{y}^{(i)}, \mathbf{A}\mathbf{x}^{(i)} \rangle \quad \forall (\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \in D, \mathbf{y}' \in \mathcal{Y}(k) \setminus \{\mathbf{y}^{(i)}\} \end{aligned}$$

- ▶ $k - 1$ constraints per training point with this formulation
- ▶ if highest scoring class that is not the gold class satisfies the constraint, all other non-gold classes will also satisfy it

Alternative constraint set:

$$m + \max_{\mathbf{y}' \in \mathcal{Y}(k) \setminus \{\mathbf{y}^{(i)}\}} \langle \mathbf{y}', \mathbf{A}\mathbf{x}^{(i)} \rangle \leq \langle \mathbf{y}^{(i)}, \mathbf{A}\mathbf{x}^{(i)} \rangle \quad \forall (\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \in D$$

Slack variable reformulation

Constraints:

$$m + \max_{\mathbf{y}' \in \mathcal{Y}(k) \setminus \{\mathbf{y}^{(i)}\}} \langle \mathbf{y}', \mathbf{Ax}^{(i)} \rangle \leq \langle \mathbf{y}^{(i)}, \mathbf{Ax}^{(i)} \rangle \quad \forall (\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \in D$$

Add slack variables $\epsilon \in \mathbb{R}_+^n$ in case the problem is not separable:

$$\begin{aligned} m + \max_{\mathbf{y}' \in \mathcal{Y}(k) \setminus \{\mathbf{y}^{(i)}\}} \langle \mathbf{y}', \mathbf{Ax}^{(i)} \rangle &\leq \langle \mathbf{y}^{(i)}, \mathbf{Ax}^{(i)} \rangle + \epsilon_i & \forall (\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \in D \\ -\langle \mathbf{y}^{(i)}, \mathbf{Ax}^{(i)} \rangle + m + \max_{\mathbf{y}' \in \mathcal{Y}(k) \setminus \{\mathbf{y}^{(i)}\}} \langle \mathbf{y}', \mathbf{Ax}^{(i)} \rangle &\leq \epsilon_i & \forall (\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \in D \\ \max \left(0, -\langle \mathbf{y}^{(i)}, \mathbf{Ax}^{(i)} \rangle + m + \max_{\mathbf{y}' \in \mathcal{Y}(k) \setminus \{\mathbf{y}^{(i)}\}} \langle \mathbf{y}', \mathbf{Ax}^{(i)} \rangle \right) &= \epsilon_i & \forall (\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \in D \end{aligned}$$

The two SVM variants

Constrained optimization problem

$$\begin{aligned} \arg \min_{\mathbf{A}} \quad & \alpha r(\mathbf{A}) + \sum_{i=1}^n \epsilon_i \\ \text{s.t.} \quad & \langle \mathbf{y}', \mathbf{A} \mathbf{x}^{(i)} \rangle + m \leq \langle \mathbf{y}^{(i)}, \mathbf{A} \mathbf{x}^{(i)} \rangle + \epsilon \quad \forall (\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \in D, \mathbf{y}' \in \mathcal{Y}(k) \setminus \{\mathbf{y}^{(i)}\} \\ & \epsilon \geq 0 \end{aligned}$$

Unconstrained optimization problem

$$\arg \min_{\mathbf{A}} \quad \sum_{i=1}^n \max \left(0, -\langle \mathbf{y}^{(i)}, \mathbf{A} \mathbf{x}^{(i)} \rangle + m + \max_{\mathbf{y}' \in \mathcal{Y}(k) \setminus \{\mathbf{y}^{(i)}\}} \langle \mathbf{y}', \mathbf{A} \mathbf{x}^{(i)} \rangle \right) + \alpha r(\mathbf{A})$$

Dual SVM training problem

Binary SVM with slack variables

SVM with slack variables

$$\min_{\mathbf{a}, \epsilon} \alpha r(\mathbf{a}) + \sum_{i=1}^n \epsilon_i$$

$$\text{s.t. } \langle \mathbf{a}, \mathbf{x}^{(i)} \rangle y^{(i)} \geq m - \epsilon_i \quad \Leftrightarrow \quad \max(0, m - \langle \mathbf{a}, \mathbf{x}^{(i)} \rangle y^{(i)}) = \epsilon_i \quad \forall (\mathbf{x}^{(i)}, y^{(i)}) \in D$$
$$\epsilon \geq 0$$

Unconstrained SVM training problem

Remark

$$\min_{\mathbf{a}} \sum_{i=1}^n \max(0, m - \langle \mathbf{a}, \mathbf{x}^{(i)} \rangle y^{(i)}) + \alpha r(\mathbf{a})$$
$$= \min_{\mathbf{a}} \sum_{i=1}^n \ell_m(y^{(i)}, \langle \mathbf{a}, \mathbf{x}^{(i)} \rangle) + \alpha r(\mathbf{a})$$

- ▶ if $m > 0$: hinge loss (usually set $m = 1$)
- ▶ if $m = 0$: perceptron loss

where $l_m(y, w) = \max(0, m - yw)$ is the binary hinge loss function, parameterized by the margin parameter $m \geq 0$

Motivations

SVM training problem

$$\min_{\mathbf{a}} \sum_{(\mathbf{x}, y) \in D} \max(0, m - \langle \mathbf{x}, \mathbf{x} \rangle y) + \frac{1}{2} \|\mathbf{a}\|_2^2$$

The hinge loss is not differentiable everywhere:

$$\ell_m(y, w) = \max(0, m - wy)$$

SVM dual formulation benefits

- ▶ Differentiable! \Rightarrow no need for a subgradient
- ▶ Can be trained with a hyper-parameter free algorithm! \Rightarrow no step-size

Notation change

Notation change

- ▶ $\mathbf{X} \in \mathbb{R}^{n \times d}$: matrix where each row consists of a training point, i.e. $X_{i,j} = x_j^{(i)}$
- ▶ $\mathbf{Y} \in \mathbb{R}^{n \times n}$: diagonal matrix containing labels, i.e. $Y_{i,i} = y^{(i)}$ and $\forall i \neq j: Y_{i,j} = 0$.
- ▶ $\ell(v_i) = \max(0, 1 - v_i)$

$$\min_{\mathbf{a}, \mathbf{v}} \sum_{i=1}^n \ell([\mathbf{YXa}]_i) + \frac{1}{2} \|\mathbf{a}\|_2^2$$

Fenchel conjugate of the Hinge loss 1/2

Conjugate of separable functions [Beck, Theorem 4.12]

Let $h : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ and $f_i : \mathbb{R} \rightarrow \mathbb{R} \cup \{\infty\}, i \in \{1 \dots n\}$ be functions defined as:

$$h(\mathbf{v}) = \sum_i f_i(v_i)$$

Then, the conjugate of h is defined as:

$$h^*(\mathbf{u}) = \sum_i f_i^*(u_i)$$

Hinge loss

The Hinge loss term in the SVM object is a separable function:

$$\sum_{i=1}^n \max(0, 1 - [\mathbf{YXa}]_i)$$

Fenchel conjugate of the Hinge loss 2/2

Hinge loss

The Hinge loss term in the SVM object is a separable function:

$$\sum_{i=1}^n \max(0, 1 - [\mathbf{YXa}]_i)$$

Hinge loss conjugate [Beck, Section 4.4.3]

Let $f(w) = \max(0, 1 - w)$, then $f^*(u) = u + \delta_{[-1,0]}(u)$.

Note: The indicator function act as constraint on the domain of the dual variables in the objective

Fenchel dual of the SVM training problem 1/2

Fenchel dual problem

Let $f : \mathbb{R}^m \rightarrow \mathbb{R}$, $h : \mathbb{R}^n \rightarrow \mathbb{R}$, $\mathbf{U} \in \mathbb{R}^{m \times n}$ and assume the following problem:

$$\begin{aligned} & \min_{\mathbf{v} \in \mathbb{R}^n} f(\mathbf{U}\mathbf{v}) + h(\mathbf{v}) \\ & \geq \max_{\boldsymbol{\lambda} \in \mathbb{R}^m} -f^*(\boldsymbol{\lambda}) - h^*(-\mathbf{U}^\top \boldsymbol{\lambda}) \end{aligned}$$

SVM case

- ▶ $\mathbf{U} = \mathbf{YX}$
- ▶ $\mathbf{v} = \mathbf{a}$

If d is the input dimension and n the number of datapoints:

$$\begin{aligned} & \min_{\mathbf{a} \in \mathbb{R}^d} f(\mathbf{YX}\mathbf{a}) + h(\mathbf{v}) \\ & \geq \max_{\boldsymbol{\lambda} \in \mathbb{R}^n} -f^*(\boldsymbol{\lambda}) - h^*(-(\mathbf{YX})^\top \boldsymbol{\lambda}) \end{aligned}$$

Fenchel dual of the SVM training problem 2/2

$$\max_{\lambda \in \mathbb{R}^n} -f^*(\lambda) - h^*(-(\mathbf{YX})^\top \lambda)$$

- ▶ $h^*(-(\mathbf{YX})^\top \lambda)$: this term is trivial because h is quadratic regularization!
- ▶ $f^*(\lambda)$: Fenchel conjugate of the Hinge loss

$$\begin{aligned} \max_{\lambda} \quad & - \sum_{i=1}^n \lambda_i - \frac{1}{2} \lambda^\top \mathbf{YX X}^\top \mathbf{Y} \lambda \\ \text{s.t.} \quad & -1 \leq \lambda_i \leq 0 \quad \forall 1 \leq i \leq n \end{aligned}$$

Recovering optimal primal variables from optimal dual variables

From the KKT's stationarity conditions, for optimal primal and dual variables we have:

$$\begin{aligned}\nabla_{\mathbf{a}} L(\mathbf{a}, \mathbf{v}, \boldsymbol{\lambda}) &= 0 \\ \nabla_{\mathbf{a}} \left(\sum_{i=1}^n \max(0, 1 - v_i) + \frac{1}{2} \|\mathbf{a}\|_2^2 + \boldsymbol{\lambda}^\top (\mathbf{YX}\mathbf{a} - \mathbf{v}) \right) &= 0 \\ \mathbf{a} + (\boldsymbol{\lambda}^\top \mathbf{YX})^\top &= 0 \\ \mathbf{a} &= -\mathbf{X}^\top \mathbf{Y} \boldsymbol{\lambda}\end{aligned}$$