

## Pisa part 2

### The Fenchel-Rockafellar duality

#### The Fenchel conjugate

$X$  Hilbert space

$f: X \rightarrow ]-\infty, +\infty]$  proper, convex l.s.c.

$$f^*: X \rightarrow ]-\infty, +\infty] \quad f^*(u) = \sup_x \langle u, x \rangle - f(x)$$

$f^*$  is proper, convex, l.s.c. and  $f^{**} = f$  (Fenchel-Moreau theorem).

#### Few facts.

1.  $f(x) + f^*(u) \geq \langle x, u \rangle$  (Young inequality)
2.  $f(x) + f^*(u) = \langle x, u \rangle \iff u \in \partial f(x)$
3.  $u \in \partial f(x) \iff x \in \partial f^*(u)$ .

$$\Gamma_0(X) = \{ f: X \rightarrow ]-\infty, +\infty] \mid f \text{ is proper, convex and l.s.c.} \}$$

Let  $f \in \Gamma_0(X)$ ,  $g \in \Gamma_0(Y)$  and  $A: X \rightarrow Y$  bounded linear operator. Consider the minimization problem:

$$(P) \quad \min_{x \in X} f(x) + g(Ax) =: \Phi(x)$$

The dual problem (in the sense of Fenchel-Rockafellar) is

$$(D) \quad \min_{u \in Y} g^*(u) + f^*(-A^*u) =: \Psi(u)$$

It is of the same form as (P),  $A^*: Y \rightarrow X$  is the transpose operator of  $A$ .

The dual problem can be obtained "formally" via Lagrange duality.

$$(P) \Leftrightarrow \inf_{\substack{(x, y) \\ Ax=y}} f(x) + g(y) \quad (\text{constrained problem})$$

The Lagrangian is  $L(x, y, u) = f(x) + g(y) + \langle Ax - y, u \rangle$

$$\text{We have } \sup_u L(x, y, u) = \begin{cases} f(x) + g(y) & \text{if } Ax = y \\ +\infty & \text{if } Ax \neq y. \end{cases}$$

Therefore (P)  $\Leftrightarrow \inf_{(x,y)} \sup_u L(x,y,u)$

Now, if we swap  $\inf$  and  $\sup$  we obtain

$$\begin{aligned} & \sup_u \inf_{(x,y)} L(x,y,u) \\ &= \sup_u \inf_{x,y} f(x) - \langle x, -A^*u \rangle + g(y) - \langle y, u \rangle \\ &= \sup_u -f^*(-A^*u) - g^*(u) \\ &= \sup_u -\Psi(u) = -\inf_u \Psi(u) \end{aligned}$$

So, in the end we obtain the dual problem (D).

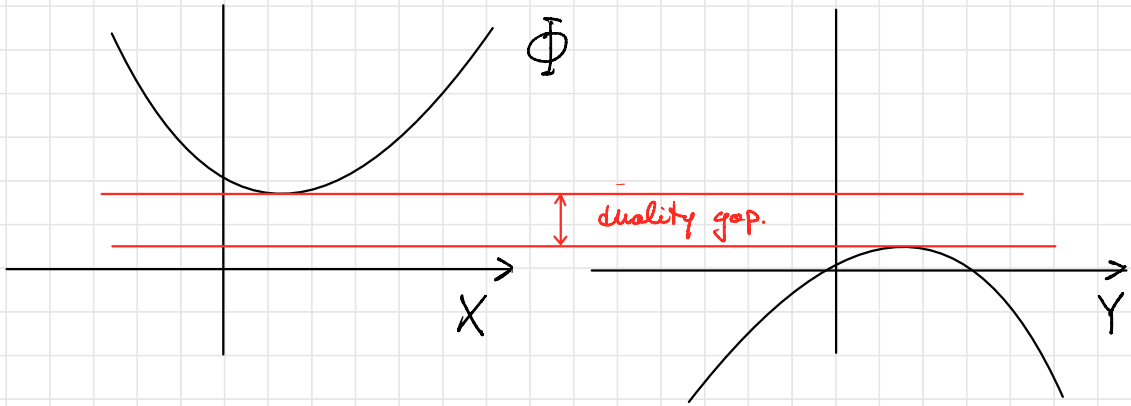
What is the relation between the primal and the dual problem?

Prop.  $\forall x \in X \quad \forall u \in Y \quad \Phi(x) \geq -\Psi(u)$

Proof:  $\Phi(x) + \Psi(u) = f(x) + f^*(-A^*u) + g(Ax) + g^*(u)$

$$\geq \langle x, -A^*u \rangle + \langle Ax, u \rangle = 0$$

Therefore the situation is as follows.



The images of  $\Phi$  &  $\Psi$  are separated and hence

$$\inf \Phi \geq \sup \Psi = -\inf \Psi$$

When  $\inf \Phi = \sup \Psi$  (duality gap is zero)

we say that strong duality holds

Strong duality is important because in that case the optimal values of the primal and dual problems are the same and hence solving the dual problem becomes an alternative to solving the primal problem.

In order to ensure strong duality, qualification conditions are needed

The primal objective function is

$$\Phi(x) = f(x) + g(Ax)$$

The minimum requirement for  $\Phi$  is that it is a proper function.

$$\text{dom } \Phi \neq \emptyset \iff \exists x \in X \text{ s.t. } f(x) < +\infty \text{ and } g(Ax) < +\infty$$

$$\iff \exists x \in \text{dom } f \text{ s.t. } \underbrace{Ax}_{u} \in \text{dom } g$$

$$\iff 0 \in \text{dom } g - A(\text{dom } f)$$

Qualification conditions ask more. One of those is

$$\text{QC: } \underline{0 \in \text{int}(\text{dom } g - A(\text{dom } f))}$$

This is satisfied for instance if

$g$  is continuous at some  $Ax$  with  $x \in \text{dom } f$ .

Thm. If QC holds, then strong duality holds and the dual problem admits solutions

(this is usually written as  $\inf \Phi = - \min_u \Psi$ )

## Remark:

Since the dual of the dual problem is the primal problem, one can ensure strong duality also by requiring  $0 \in \text{int}(\text{dom} f^* + A^*(\text{dom} g^*))$ .

Let's address the characterization of primal and dual solutions.

Theorem Let  $S = \text{argmin} \Phi$ ,  $S^* = \text{argmin} \Psi$  and  $\bar{x} \in X$  and  $\bar{u} \in Y$ . Then the following are equivalent

a)  $\bar{x} \in S$ ,  $\bar{u} \in S^*$  and  $\inf \Phi = -\inf \Psi$ .

b)  $\Phi(\bar{x}) + \Psi(\bar{u}) = 0$ .

c)  $-A^*\bar{u} \in \partial f(\bar{x})$  and  $\bar{u} \in \partial g(A\bar{x})$

d)  $\bar{x} \in \partial f^*(-A^*\bar{u})$  and  $A\bar{x} \in \partial g^*(\bar{u})$

} KKT conditions.

Proof:

$$\Phi(\bar{x}) + \Psi(\bar{u}) = \underbrace{f(\bar{x}) + f^*(-A^*\bar{u})}_{\geq 0} - \langle \bar{x}, -A^*\bar{u} \rangle$$

$$+ \underbrace{g(A\bar{x}) + g^*(\bar{u})}_{\geq 0} - \langle A\bar{x}, \bar{u} \rangle$$

$$\text{Thus, b) } \Leftrightarrow f(\bar{x}) + f^*(-A^*\bar{u}) = \langle \bar{x}, -A^*\bar{u} \rangle$$

$$\text{and } g(A\bar{x}) + g^*(\bar{u}) = \langle A\bar{x}, \bar{u} \rangle$$

$$\Leftrightarrow -A^*\bar{u} \in \partial f(\bar{x}) \text{ and } \bar{u} \in \partial g(A\bar{x})$$

Remark: As a consequence of the previous theorem, we have that once strong duality is ensured, KKT conditions characterize primal and dual solutions simultaneously

---

Now we address the computational aspects of duality.

Ideally one wants to go from a dual solution to a primal solution, so that the primal problem can be solved via the dual problem.

In order to do so, we need an additional assumption.

Assumption  $f$  is  $\mu$ -strongly convex

As a consequence of this assumption we have

1) the primal problem has a unique solution

2)  $\text{dom } f^* = X$  and  $f^*$  is  $\frac{1}{\mu}$ -Lipschitz smooth.

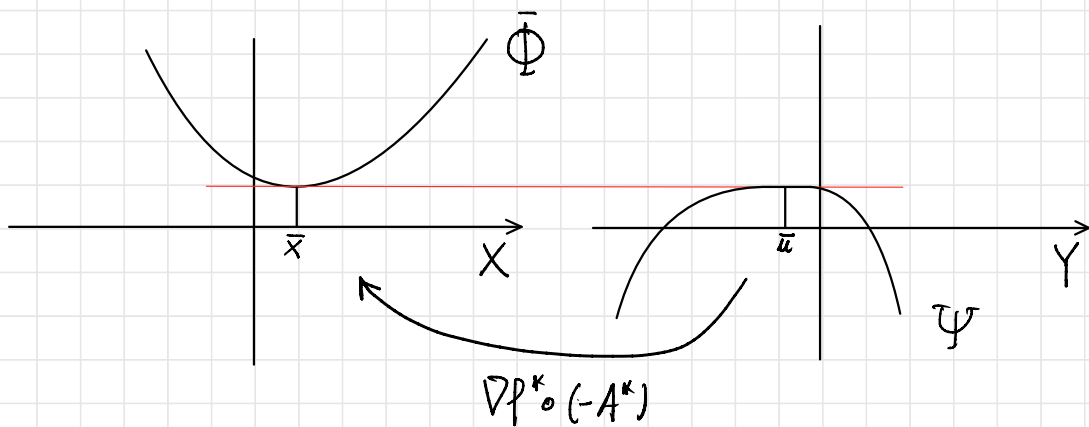
Then the KKT conditions become

$$\bar{x} = \nabla p^*(-A^* \bar{u}) \quad \text{and} \quad A\bar{x} \in \partial q^*(\bar{u})$$

The first of the KKT conditions determine uniquely the primal solution from a dual solution. So, the map.

$$u \mapsto \nabla p^*(-A^* u)$$

provides a way to go from the dual variable to the primal variable.



### Proposition

Let  $\bar{x} = \nabla p^*(-A^* \bar{u})$  and  $x = \nabla p^*(-A^* u)$ .

Then  $\frac{\mu}{2} \|x - \bar{x}\|^2 \leq \psi(u) - \psi(\bar{u})$



Proof.

$$f(\bar{x}) + f^*(-A^* \bar{u}) = \langle \bar{x}, -A^* \bar{u} \rangle$$

$$f(x) + f^*(-A^* u) = \langle x, -A^* u \rangle$$

$$\begin{aligned} \Psi(u) - \Psi(\bar{u}) &= g^*(u) - g^*(\bar{u}) + f^*(-A^* u) - f^*(-A^* \bar{u}) \\ &= g^*(u) - g^*(\bar{u}) + \langle x, -A^* u \rangle - f(x) + f(\bar{x}) + \langle \bar{x}, A^* \bar{u} \rangle \end{aligned}$$

$$\begin{aligned} &\geq \underbrace{\langle u - \bar{u}, A \bar{x} \rangle}_{A \bar{x} \in \partial g^*(\bar{u})} - \langle A x, u \rangle + \langle A \bar{x}, \bar{u} \rangle + f(\bar{x}) - f(x) \end{aligned}$$

$$= \langle x - \bar{x}, -A^* u \rangle + f(\bar{x}) - f(x)$$

since  $-A^* u \in \partial f(x)$  and  $f$  is  $\mu$ -strongly convex

$$f(\bar{x}) - f(x) = \langle \bar{x} - x, -A^* u \rangle \geq \frac{\mu}{2} \|x - \bar{x}\|^2. \quad \square$$

Remark The previous result shows that

if  $\Psi(u_n) \rightarrow \min \Psi$  then  $x_n = \nabla f^*(-A^* u_n) \rightarrow \bar{x}$ .

One can solve the primal problem via the dual problem

# A dual approach to regularized empirical risk minimization

We first introduce the hypothesis space

Let  $H$  be a Hilbert space and  $\Lambda: \mathcal{X} \rightarrow H$  be a map (called feature map).

Then this map serves to define real functions defined in  $\mathcal{X}$  and parametrized in  $H$ .

$$f: \mathcal{X} \rightarrow \mathbb{R} \quad f(x) = \langle w, \Lambda(x) \rangle, \quad w \in H.$$

The set of all these functions define a Hilbert space of functions

$$\mathcal{H} = \{ f: \mathcal{X} \rightarrow \mathbb{R} \mid f = \langle w, \Lambda(\cdot) \rangle \quad w \in H \}$$

$$\|f\|_{\mathcal{H}} = \inf \{ \|w\| \mid w \in H \text{ and } f = \langle w, \Lambda(\cdot) \rangle \}$$

If  $\overline{\text{span}\{\Lambda(x) \mid x \in \mathcal{X}\}} = H$ , then there is

a unique  $W \in H$  s.t.  $f = \langle w, \Lambda(\cdot) \rangle$  and  $\|f\|_{\mathcal{H}} = \|W\|$

One can prove that  $\mathcal{H}$  is a RKHS with kernel

$$\underline{K(x, x') = \langle \Lambda(x), \Lambda(x') \rangle.}$$

Given

$(x_1, y_1) \dots (x_n, y_n) \in \mathcal{X} \times \mathcal{Y}$  training points,

$l: \mathbb{R} \times \mathcal{Y} \rightarrow \mathbb{R}$  loss function (convex wrt the first v.),

the regularized empirical risk minimization

problem is

$$\min_{w \in H} \frac{1}{n\lambda} \sum_{i=1}^n l(\langle w, \Lambda(x_i) \rangle, y_i) + \frac{1}{2} \|W\|^2 \quad \text{R-ERM}$$

We will put the problem in the framework of the Fenchel-Rockafellar duality.

$$1. \quad A: H \rightarrow \mathbb{R}^n \quad Aw = \begin{bmatrix} \langle w, \Lambda(x_1) \rangle \\ \vdots \\ \langle w, \Lambda(x_n) \rangle \end{bmatrix}$$

$$2. \quad f(w) = \frac{1}{2} \|W\|^2 \quad (\text{which is } 1\text{-strongly convex})$$

$$3. g: \mathbb{R}^n \rightarrow \mathbb{R} \quad g(z) = \frac{1}{n\lambda} \sum_{i=1}^n \ell(-z_i, y_i)$$

Then problem R-ERM is written as

$$\min_{w \in H} g(-Aw) + f(w).$$

Thus, the dual problem is

$$\min_{\alpha \in \mathbb{R}^n} g^*(\alpha) + f^*(A^*\alpha).$$

Now we note that

$$1. A^*\alpha = \sum_{i=1}^n \alpha_i \Lambda(x_i)$$

$$2. f^*(v) = \frac{1}{2} \|v\|^2$$

$$\text{Therefore } f^*(A^*\alpha) = \frac{1}{2} \left\| \sum_{i=1}^n \alpha_i \Lambda(x_i) \right\|^2$$

$$= \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j \langle \Lambda(x_i), \Lambda(x_j) \rangle = \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j K(x_i, x_j)$$

The Gram matrix  $K_{ij} = K(x_i, x_j)$ ,

$$f^*(A^*\alpha) = \frac{1}{2} \langle K\alpha, \alpha \rangle = \frac{1}{2} \alpha^T K \alpha$$

3.  $g$  can be written as

$$g(z) = \sum_{i=1}^n g_i(z)$$

where  $g_i(t) = \frac{1}{n\lambda} \ell(-t, y_i)$   $g_i: \mathbb{R} \rightarrow \mathbb{R}$ .

Therefore  $g^*(\alpha) = \sum_{i=1}^n g_i^*(\alpha)$

Then we need to compute  $g_i^*$

By definition

$$\begin{aligned} g_i^*(s) &= \sup_t st - g_i(t) = \sup_t st - \frac{1}{n\lambda} \ell(-t, y_i) \\ &= \frac{1}{n\lambda} \sup_t (-s n\lambda)(-t) - \ell(-t, y_i) \\ &= \frac{1}{n\lambda} \ell^*(-s n\lambda, y_i) \end{aligned}$$

Putting everything together we obtain

$$\min_{\alpha \in \mathbb{R}^n} \frac{1}{2} \langle K \alpha, \alpha \rangle + \frac{1}{n\lambda} \sum_{i=1}^n \ell^*(-\alpha_i n\lambda, y_i)$$

The KKT conditions in general are

$$\bar{w} = \underbrace{\nabla f^*}_{I} (A^* \bar{\alpha}) \text{ and } -A \bar{w} \in \partial g^*(\bar{\alpha})$$

and in this case they become

$$\bar{w} = \sum_{i=1}^n \bar{\alpha}_i \Lambda(x_i) \quad - \langle \bar{w}, \Lambda(x_i) \rangle \in \partial \varphi_i^*(\bar{\alpha}_i)$$

1

$$\bar{\alpha}_i \in \partial \varphi_i(-\langle \bar{w}, \Lambda(x_i) \rangle)$$

$$- \bar{\alpha}_i \in \frac{1}{n\lambda} \partial l(\langle \bar{w}, \Lambda(x_i) \rangle, y_i)$$

2

We note that the first of the KKT conditions is nothing more than the representer theorem

The corresponding function is

$$\begin{aligned} f_{\bar{w}}(x) &= \langle \bar{w}, \Lambda(x) \rangle = \sum_{i=1}^n \bar{\alpha}_i \langle \Lambda(x_i), \Lambda(x) \rangle \\ &= \sum_{i=1}^n \bar{\alpha}_i K(x_i, x) \end{aligned}$$

which shows that the estimator can be expressed in terms of the kernel function and the dual solution.

This aspect is what makes this approach attractive and is known as the kernel trick

This framework defines the so called kernel methods which has the following nice features:

1. The dual problem is a finite dimensional convex problem
2. The dual objective function is expressed in terms of the kernel functions, the data and the loss function. Therefore it is completely known.
3. The final estimator (function) is expressed in terms of the kernel function, the data and the dual solution.

Now, concerning the algorithmic approach to R-ERM,

We note that the mapping linking the dual variable to the primal variable is

$$\alpha \mapsto \sum_{i=1}^n \alpha_i \Lambda(x_i)$$

and for every  $\alpha$  if we define  $w = \sum_{i=1}^n \alpha_i \Lambda(x_i)$

we have

$$\frac{1}{2} \|w - \bar{w}\|^2 \leq \Psi(\alpha) - \min \Psi$$

---

This inequality is the basis of approaching the R-EAM problem via dual algorithms

Indeed if  $f = \langle w, \Lambda(\cdot) \rangle$  and  $\bar{f} = \langle \bar{w}, \Lambda(\cdot) \rangle$ , we have

$$\begin{aligned} |f(x) - \bar{f}(x)| &= |\langle w - \bar{w}, \Lambda(x) \rangle| \leq \|w - \bar{w}\| \|\Lambda(x)\| \\ &\leq \sqrt{2(\Psi(\alpha) - \min \Psi)} \sqrt{K(x, x)}. \end{aligned}$$

If  $(\alpha^{(k)})_{k \in \mathbb{N}}$  is such that  $\Psi(\alpha^{(k)}) \rightarrow \inf \Psi$

and if we define  $w_k = \sum_{i=1}^n \alpha^{(k)} \Lambda(x_i)$  and

the corresponding function  $f_k = \langle w_k, \Lambda(\cdot) \rangle \in \mathcal{H}$ ,

we have that at a new point  $x \in \mathcal{X}$ .

$$|f_k(x) - \bar{f}(x)| \leq \underbrace{\sqrt{2(\Psi(\alpha^{(k)}) - \min \Psi)}}_0 \sqrt{K(x, x)} \rightarrow 0$$

In case  $\sup_x K(x, x) < +\infty$ , then  $\|f_k - \bar{f}\|_\infty \rightarrow 0$ .



## Example

1. The square loss  $l(t, y) = \frac{1}{2}(t - y)^2$

$$l^*(s, y) = \frac{1}{2}s^2 + y \cdot s. \quad \text{Hence}$$

$$\begin{aligned} \frac{1}{n\lambda} \sum_{i=1}^n l^*(-\alpha_i n\lambda, y_i) &= \frac{1}{n\lambda} \sum_{i=1}^n \left( \frac{1}{2} (\alpha_i n\lambda)^2 - y_i \alpha_i n\lambda \right) \\ &= \frac{n\lambda}{2} \|\alpha\|^2 - \langle y, \alpha \rangle \end{aligned}$$

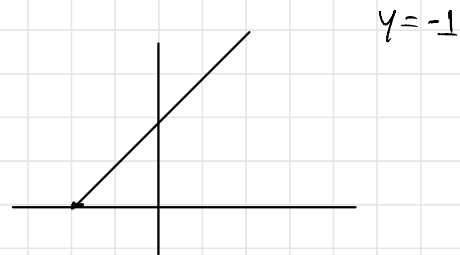
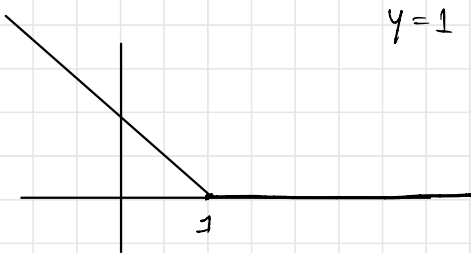
The dual problem is

$$\min_{\alpha \in \mathbb{R}^n} \frac{1}{2} \alpha^T K \alpha + \frac{n\lambda}{2} \|\alpha\|^2 - \langle y, \alpha \rangle$$

which has solution

$$\begin{aligned} K\alpha + n\lambda\alpha - y &= 0 \iff (K + n\lambda I)\alpha = y \\ \iff \alpha &= (K + n\lambda I)^{-1} y \end{aligned}$$

2. The hinge loss  $l(t, y) = (1 - ty)_+$  with  $y \in \mathcal{Y} = \{-1, 1\}$ .



$$l(t, y) = \chi(t, y) \quad \text{where} \quad \chi(r) = (1-r)_+$$

$$\begin{aligned} l^*(s, y) &= \sup_t st - \chi(t, y) = \sup_t (s, y)(t, y) - \chi(t, y) \\ &= \chi^*(s, y) \end{aligned}$$

Moreover the function  $\chi$  is the function  $(-r)_+$  translated by  $-1$ , therefore  $\chi(r) = [( -s )_+]_{s=r-1}$

$$\chi^*(q) = q + i_{[-1, 0]}(q) \quad [(-s)_+]^*(q) = i_{[-1, 0]}(q)$$

$$l^*(s, y) = s, y + i_{[-1, 0]}(s, y)$$

$$\frac{1}{n\lambda} \sum_{c=1}^n l^*(-\alpha c n \lambda, y_c) = \frac{1}{n\lambda} \sum_{c=1}^n -\alpha c n \lambda y_c + i_{[-1, 0]}(-\alpha c n \lambda y_c)$$

$$0 \leq \alpha c n \lambda y_c \leq 1 \quad \Leftrightarrow \quad 0 \leq \alpha y_c \leq \frac{1}{n\lambda}$$

$$\min_{\alpha \in \mathbb{R}^n} \frac{1}{2} \langle K \alpha, \alpha \rangle - \alpha^T y + i_{[0, \frac{1}{n\lambda}]^n}(\alpha \odot y)$$

This is a constrained quadratic minimization problem.