

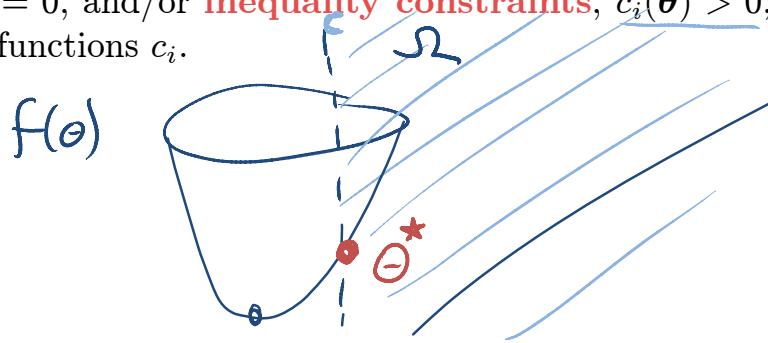
Constrained optimization

- Consider the following **constrained optimization problem**

$$\theta^* = \arg \min_{\theta \in \Omega} f(\theta)$$

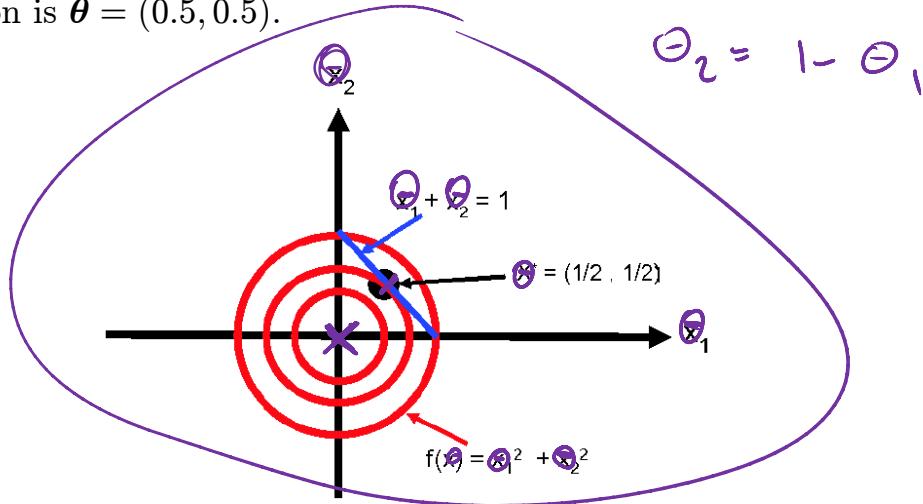
where Ω is some **feasible set**. If the parameters are real-valued, we typically assume $\Omega \subseteq \mathbb{R}^D$, but it could be a more abstract space, such as the set of positive definite matrices.

- The feasible set is then often defined in terms of a set of **equality constraints**, $c_i(\theta) = 0$, and/or **inequality constraints**, $c_i(\theta) > 0$, for certain constraint functions c_i .



Constrained optimization

- Suppose that we have a single equality constraint $c(\theta) = 0$.
- For example, we might have a quadratic objective, $f(\theta) = \theta_1^2 + \theta_2^2$, subject to a linear equality constraint, $c(\theta) = 1 - \theta_1 - \theta_2 = 0$.
- What we are trying to do is find the point θ^* that lives on the line, but which is closest to the origin. It is geometrically obvious that the optimal solution is $\theta = (0.5, 0.5)$.

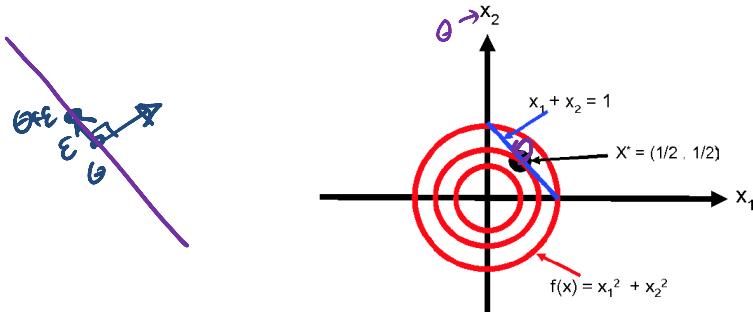


Constrained optimization

- The gradient of the constraint function $\nabla c(\theta)$ will be orthogonal to the constraint surface.
- To see why, consider a point θ on the constraint surface, and another point nearby, $\theta + \epsilon$, that also lies on the surface. If we make a Taylor expansion around θ we have

$$c(\theta + \epsilon) \approx c(\theta) + \epsilon^T \nabla c(\theta)$$

Since both θ and $\theta + \epsilon$ are on the constraint surface, we must have $c(\theta) = c(\theta + \epsilon)$ and hence $\epsilon^T \nabla c(\theta) \approx 0$. Since ϵ is parallel to the constraint surface, we see that the vector ∇c is normal to the surface.

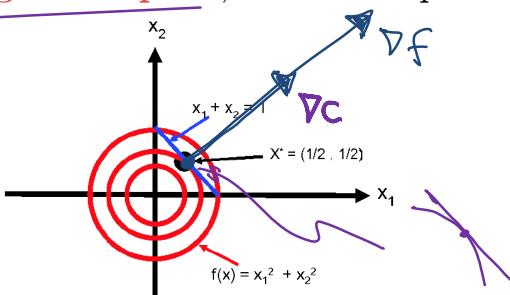


Constrained optimization

- We seek a point θ^* on the constraint surface such that $f(\theta)$ is minimized. Such a point must have the property that $\nabla f(\theta)$ is also orthogonal to the constraint surface, as otherwise we could decrease $f(\theta)$ by moving a short distance along the constraint surface.
- Since both $\nabla c(\theta)$ and $\nabla f(\theta)$ are orthogonal to the constraint surface at θ^* , they must be parallel (or anti-parallel) to each other. Hence there must exist a constant $\lambda^* \neq 0$ such that

$$\nabla f(\theta^*) = \lambda^* \nabla c(\theta^*)$$

λ^* is called a **Lagrange multiplier**, and can be positive or negative, but not zero.



Lagrangian

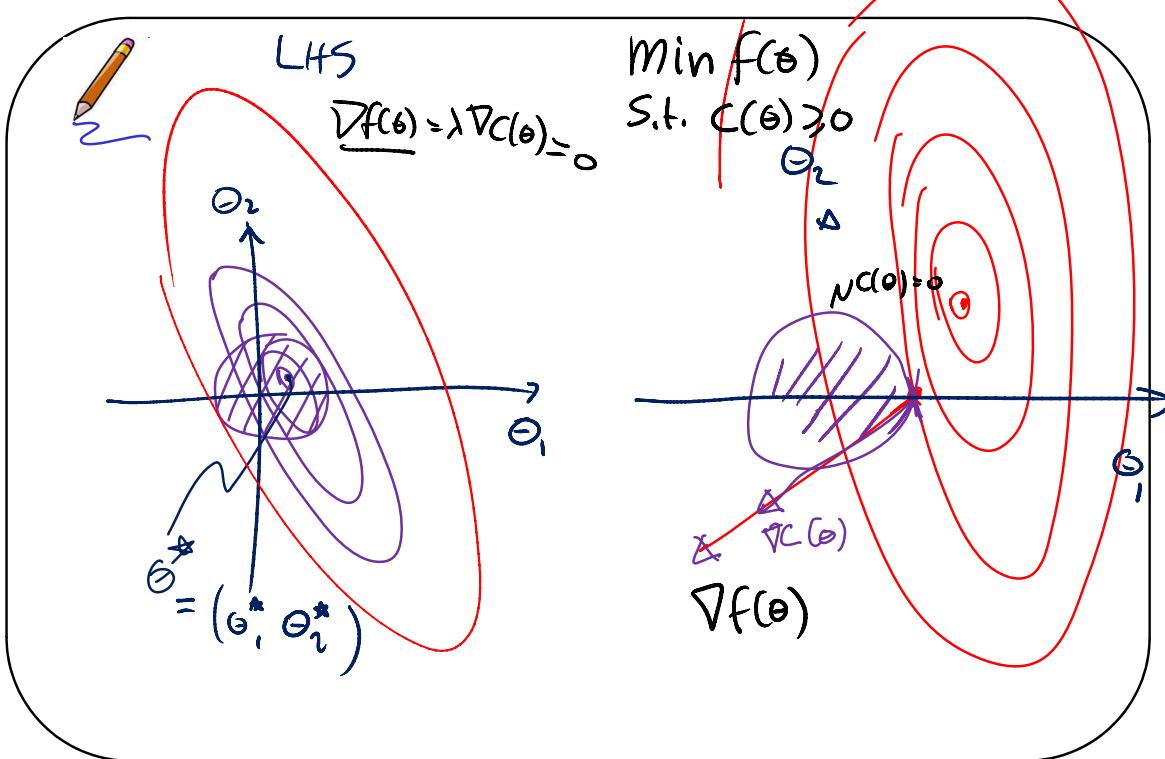
- We can now convert our constrained optimization problem into an unconstrained one by defining a new function called the **Lagrangian**:

$$L(\theta, \lambda) := f(\theta) - \lambda c(\theta)$$

We now have $D + 1$ equations in $D + 1$ unknowns, which we can solve for θ^* and λ . Why? Since we are only interested in θ^* , we can “throw away” the value λ ; hence it is sometimes called an **undetermined multiplier**.

$\nabla_{\theta} L(\theta, \lambda) = \nabla_{\theta} f(\theta) - \lambda \nabla_{\theta} c(\theta) = 0$
 $\nabla_{\theta} f(\theta) = \lambda \nabla_{\theta} c(\theta)$
 $\nabla_{\lambda} L(\theta, \lambda) = -c(\theta) = 0$
 $c(\theta) = 0$

Inequality constraints



Inequality constraints

- Now consider the case where we have a single **inequality constraint** $c(\theta) \geq 0$.
- If the solution lies in the region where $c(\theta) > 0$, the constraint is **inactive**, so we have the usual stationarity condition $\nabla f(\theta^*) = 0$. Our equations still hold, provided we set $\lambda^* = 0$. LHS
- If the solution lies on the boundary where $c(\theta) = 0$, the constraint is **active**, so $\nabla c(\theta)$ and $\nabla f(\theta)$ must be parallel, as for the equality constraint case. RHS
- However, this time we require that $\lambda^* > 0$, so the gradients point in the *same* direction. Since the gradients of c and f point in the same direction, we will follow c to its minimum, where $c(\theta^*) = 0$.
- We can summarize these two cases by writing $\lambda^* c(\theta^*) = 0$; either $\lambda^* = 0$ or $c(\theta^*) = 0$ (or both). This is called the **complementarity condition**.

Inequality constraints

- Putting it all together, the problem of minimizing $f(\boldsymbol{\theta})$ subject to $c(\boldsymbol{\theta}) \geq 0$ can be obtained by optimizing the Lagrangian subject to the following constraints:

$$\begin{cases} c(\boldsymbol{\theta}) \geq 0 \\ \lambda^* \geq 0 \\ \lambda^* c(\boldsymbol{\theta}^*) = 0 \end{cases}$$

Many constraints

- In general, if we have multiple equality constraints, $c_i(\boldsymbol{\theta}) = 0$ for $i \in \mathcal{E}$, and multiple inequality constraints, $c_i(\boldsymbol{\theta}) \geq 0$ for $i \in \mathcal{I}$, we can define the feasible set as

$$\Omega = \{\boldsymbol{\theta} \in \mathbb{R}^D : \underbrace{c_i(\boldsymbol{\theta}) = 0}_{\text{eq.}}, i \in \mathcal{E}, \underbrace{c_i(\boldsymbol{\theta}) \geq 0}_{\text{ineq.}}, i \in \mathcal{I}\}$$

and the Lagrangian as

$$L(\boldsymbol{\theta}, \boldsymbol{\lambda}) = f(\boldsymbol{\theta}) - \sum_{i \in \mathcal{E} \cup \mathcal{I}} \lambda_i c_i(\boldsymbol{\theta})$$

- The **active set** is defined as the constraints that are active at a point:

$$\mathcal{A}(\boldsymbol{\theta}) = \mathcal{E} \cup \{i \in \mathcal{I} : c_i(\boldsymbol{\theta}) = 0\}$$

Karush-Kuhn-Tucker conditions

- We have the following necessary first-order conditions for being at a local minimum:

$$\begin{aligned}\nabla_{\theta} L(\theta, \lambda) &= 0 \\ c_i(\theta^*) &= 0 \quad \forall i \in \mathcal{E} \checkmark \\ c_i(\theta^*) &\geq 0 \quad \forall i \in \mathcal{I} \checkmark \\ \lambda_i^* &\geq 0 \quad \forall i \in \mathcal{I} \checkmark \\ \lambda_i^* c_i(\theta^*) &= 0 \quad \forall i \in \mathcal{I} \cup \mathcal{E} \checkmark\end{aligned}$$

- These are called the **Karush-Kuhn-Tucker** or **KKT** conditions.
- If f and the c_i are convex, the KKT conditions are sufficient for a minimum as well.

Example

Maximize $f(\theta) = 1 - \theta_1^2 - \theta_2^2$ subject to the constraint that $\theta_1 + \theta_2 = 1$.

(i) $L(\theta_1, \theta_2, \lambda) = f(\theta_1, \theta_2) - \lambda c(\theta_1, \theta_2)$

$$= 1 - \theta_1^2 - \theta_2^2 - \lambda [1 - \theta_1 - \theta_2]$$

(ii) $\nabla_{\theta_1} L(\theta_1, \theta_2, \lambda) = 0 \Rightarrow$
 $\nabla_{\theta_2} L(\theta_1, \theta_2, \lambda) = 0 \Rightarrow$
 $\nabla_{\lambda} L(\theta_1, \theta_2, \lambda) = 0 \Rightarrow$

Example



$$\theta^* = \left(\frac{1}{2}, \frac{1}{2} \right)$$

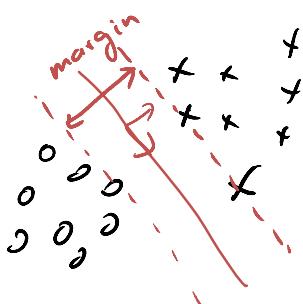
Quadratic programs

- A generic **quadratic program** or **QP** has the form

$$\min_{\theta} \frac{1}{2} \theta^T H \theta + d^T \theta \quad \text{s.t. } A\theta \leq b, \quad A_{eq}\theta = b_{eq}, \quad b_l \leq \theta \leq b_u$$

The constraints $b_l \leq \theta \leq b_u$ are known as **box constraints**, and can always be rewritten as linear inequality constraints.

- QPs arise in several areas of machine learning, including **support vector machines** and **lasso**.



- Assume we want to minimize:

$$f(\theta) = (\theta_1 - \frac{3}{2})^2 + (\theta_2 - \frac{1}{8})^2 = \frac{1}{2} \theta^T \mathbf{H} \theta + \mathbf{d}^T \theta + \text{const}$$

where $\mathbf{H} = 2\mathbf{I}$ and $\mathbf{d} = -(3, 1/4)$, subject to

\hookrightarrow Norm:

$$\|\theta\|_1 = |\theta_1| + |\theta_2|$$

$$|\theta_1| + |\theta_2| \leq 1$$

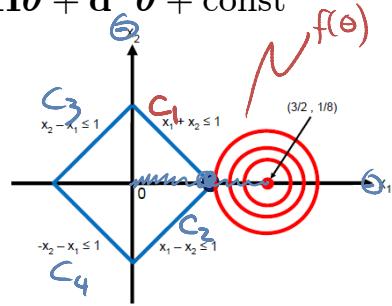
We can rewrite the constraints as

$$\begin{array}{l} \theta_1 + \theta_2 \leq 1, \quad \theta_1 - \theta_2 \leq 1, \quad -\theta_1 + \theta_2 \leq 1, \quad -\theta_1 - \theta_2 \leq 1 \\ \downarrow \quad \downarrow \quad \downarrow \quad \downarrow \\ \theta_1 + \theta_2 \geq 0, \quad \theta_1 - \theta_2 \geq 0, \quad -\theta_1 + \theta_2 \geq 0, \quad -\theta_1 - \theta_2 \geq 0 \end{array} \quad (C)$$

which we can write more compactly as

C_1

$$\mathbf{b} - \mathbf{A}\theta \geq 0$$



where $\mathbf{b} = \mathbf{1}$ and

$$\left[\begin{array}{c} 1 \\ 1 \\ -1 \\ -1 \end{array} \right]$$

$$\mathbf{A} = \begin{pmatrix} 1 & 1 \\ 1 & -1 \\ -1 & 1 \\ -1 & -1 \end{pmatrix}$$

Quadratic programs

- The Lagrangian is

$$L(\theta, \lambda) = \frac{1}{2} \theta^T \mathbf{H} \theta + \mathbf{d}^T \theta + \lambda^T (\mathbf{A}\theta - \mathbf{b})$$

and the KKT conditions are

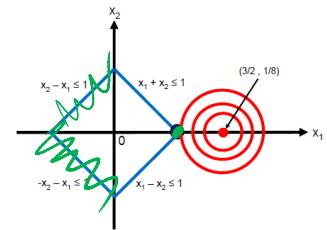
$$\mathbf{H}\theta + \mathbf{d} + \mathbf{A}^T \lambda = 0$$

$$\mathbf{b} - \mathbf{A}\theta \geq 0$$

If we treat the inequality as an equality, we can write

$$\begin{pmatrix} \mathbf{H} & \mathbf{A}^T \\ \mathbf{A} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \theta \\ \lambda \end{pmatrix} = \begin{pmatrix} -\mathbf{d} \\ \mathbf{b} \end{pmatrix}$$

Quadratic programs



- The KKT matrix on the LHS is singular. Note constraints c_3 and c_4 (corresponding to the two left faces of the diamond) are inactive, so $c_3(\boldsymbol{\theta}^*) > 0$ and $c_4(\boldsymbol{\theta}^*) > 0$ and hence, by complementarity, $\lambda_3^* = \lambda_4^* = 0$. We can therefore remove these inactive constraints to get the following:

$$\left(\begin{array}{cccc} 2 & 0 & 1 & 1 \\ 0 & 2 & 1 & -1 \\ 1 & 1 & 0 & 0 \\ 1 & -1 & 0 & 0 \end{array} \right) \begin{pmatrix} \theta_1 \\ \theta_2 \\ \lambda_1 \\ \lambda_2 \end{pmatrix} = \begin{pmatrix} 3 \\ 1/4 \\ 1 \\ 1 \end{pmatrix}$$

We see that the solution is

$$\boldsymbol{\theta}^* = (1, 0)^T, \boldsymbol{\lambda}^* = (0.875, 0.125, 0, 0)^T$$

Notice that the optimal value of $\boldsymbol{\theta}$ occurs at one of the vertices of the L1 simplex. Consequently the solution vector is **sparse**.

Lasso for feature selection

$\min_{\boldsymbol{\theta}} \| \mathbf{Y} - \mathbf{X}\boldsymbol{\theta} \|_2^2 + \lambda \| \boldsymbol{\theta} \|_1$

$\min_{\boldsymbol{\theta}: \| \boldsymbol{\theta} \|_1 = t} \| \mathbf{Y} - \mathbf{X}\boldsymbol{\theta} \|_2^2$

$t = \mathcal{F}(\lambda)$

Lasso for feature selection



Duality

- **Duality theory** provides an alternative way to express optimization problems that can often lead to faster algorithms, as well as new insights into a problem. It also relaxes some of the differentiation conditions.

Duality

- Consider the **primal problem**

$$\min_{\theta} f(\theta) \text{ s.t. } \mathbf{c}(\theta) \geq \mathbf{0}$$

The Lagrangian is

$$L(\theta, \lambda) = f(\theta) - \lambda^T \mathbf{c}(\theta)$$

$$f(\theta) = \lambda^T \mathbf{c}(\theta) + L$$

We define the **dual** objective function as

$$g(\lambda) = \min_{\theta} L(\theta, \lambda) = \min_{\theta} f(\theta) - \lambda^T \mathbf{c}(\theta) = -f^*(\lambda)$$

where f^* is the **Fenchel conjugate** of f .

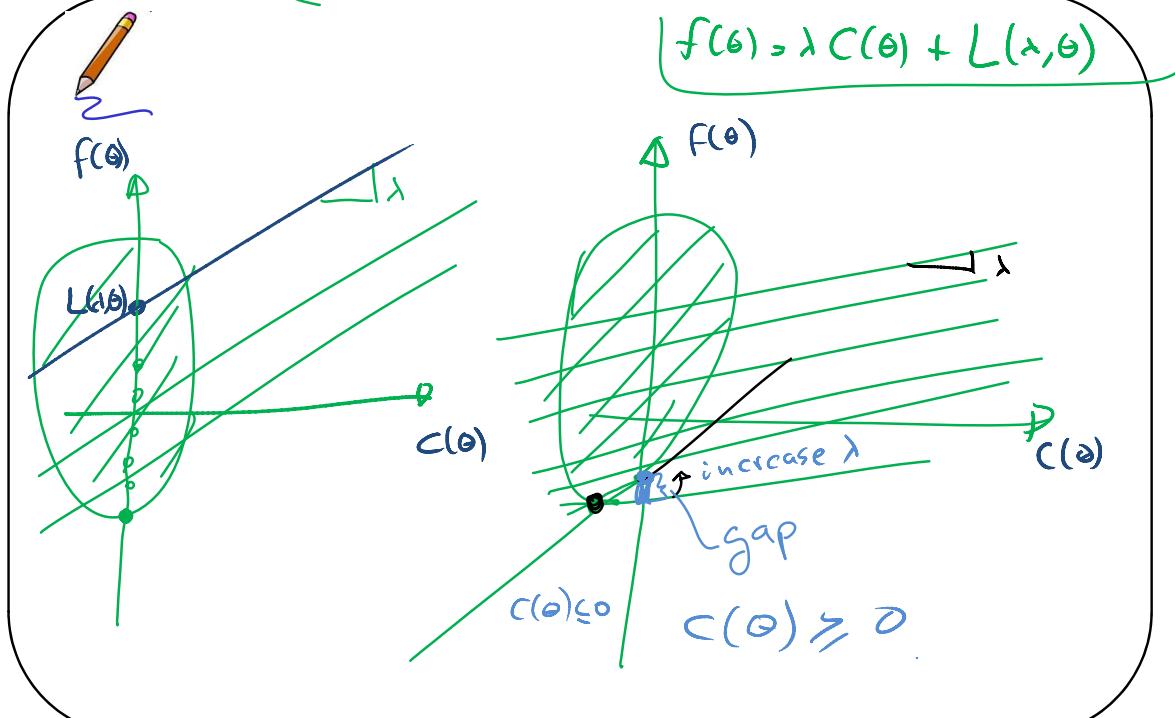
- We see that the dual objective g is a concave function, since it is a minimum over an affine function of λ . The corresponding **dual problem** is

$$\max_{\lambda} g(\lambda) \text{ s.t. } \lambda \geq \mathbf{0}$$

Duality

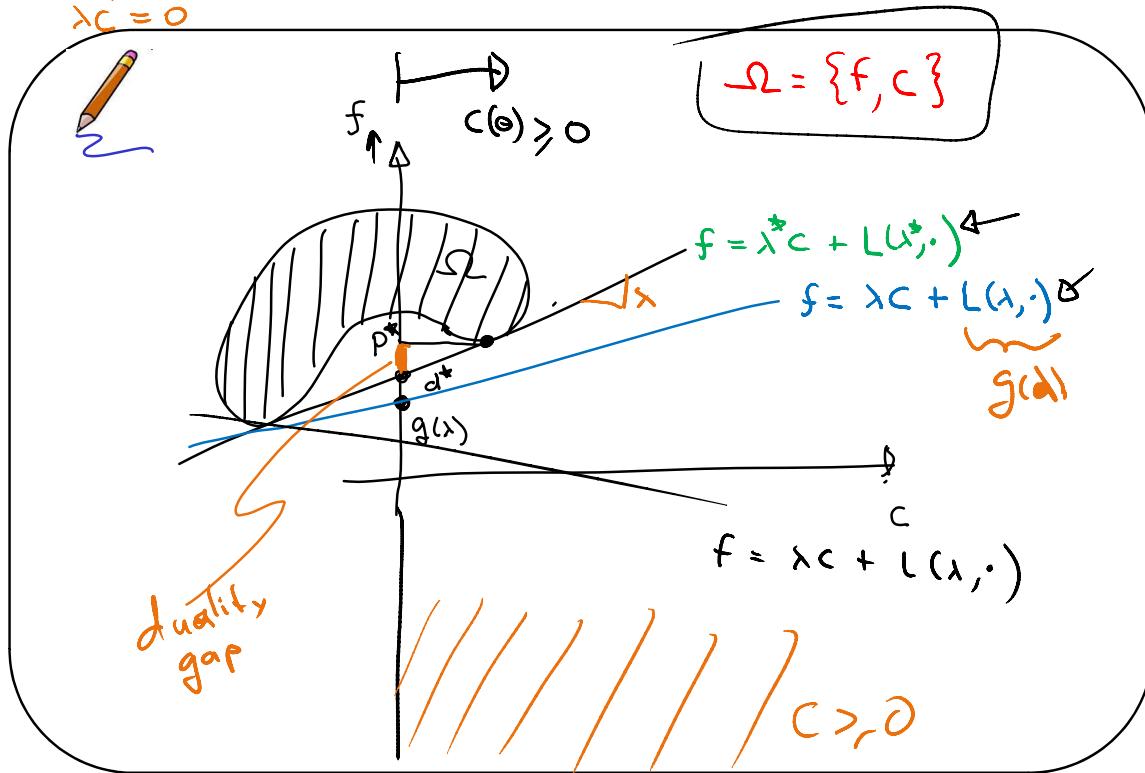
$$L(\theta, \lambda) = f(\theta) - \lambda^T \mathbf{c}(\theta)$$

$$f(\theta) > \lambda^T \mathbf{c}(\theta) + L(\lambda, \theta)$$



$$\begin{aligned} c > 0 \\ \lambda > 0 \\ \lambda c = 0 \end{aligned}$$

Duality



Duality

- Solving the dual has several advantages:
 1. It is always convex, even if the primal is not;
 2. The number of variables in the dual is equal to the number of constraints in the primal, which is often less than the number of variables in the primal
 3. It might enable us to deal with non-differentiable problems.

Duality

- The key question is, do the two methods give the same results? Let $p^* = f(\theta^*)$ be the optimal primal value, and $d^* = g(\lambda^*)$ be the optimal dual value. We have the following two important theorems:

- **Weak duality:** $d^* \leq p^*$. This always holds. To see this, note that for $\lambda \geq 0$, since $c(\theta) \geq 0$,

$$f(\theta) \geq L(\theta, \lambda) \geq \min_{\theta'} L(\theta', \lambda) = g(\lambda)$$

- **Strong duality:** $d^* = p^*$. This only holds for convex problems. The reason is that a convex function can be precisely represented either in primal or dual form.

Put another way, for any real function $L(\theta, \lambda)$, weak duality says we always have

$$\min_{\theta} \max_{\lambda} L(\theta, \lambda) \geq \max_{\lambda} \min_{\theta} L(\theta, \lambda)$$

If strong duality holds, the two terms are equal, so the **duality gap**, $p^* - d^*$, is zero. In this case, $L(\theta^*, \lambda^*)$ is a **saddle point**.

Further reading

- Please read the book section about linear programming as another example.
- Read on the algorithms
 1. Interior point methods
 2. Active set methods
 3. Projected gradient

Nocedal & Wright
Stephen Boyd
Bertsekas
Nonlinear programming

Mark Schmidt