

Function Approximation Methods-V

CS771: Introduction to Machine Learning
Purushottam Kar



Outline of today's discussion

PART I

- Revisit Lagrangian duality
- Apply duality to SVMs
- Observe a beautiful connection b/w CD and SGD
- See how some state-of-the-art solvers for SVMs are built

PART II

- Learn how to model data (data features that is)
- Look at the naïve Bayes technique for supervised problems
- Get introduced to the GMM technique for unsupervised problems

Duality

How to take the problem you want to solve
... and convert it to a problem you can solve

Fenchel, Lagrange, Wolfe, Pontryagin

Duality

How to take the problem you want to solve
... and convert it to a problem you can solve

Fenchel, Lagrange, Wolfe, Pontryagin

Duality

How to take the problem you want to solve
... and convert it to a problem you can solve

The Lagrangian

Sept 6, 2017



The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$\hat{\mathbf{w}} = \arg \min \sum_{i=1}^n (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2$$
$$\text{s.t. } \|\mathbf{w}\|_2 \leq r$$

The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

Want $g(\mathbf{w}) \geq 0$?
 $-g(\mathbf{w}) \leq 0$

The Lagrangian

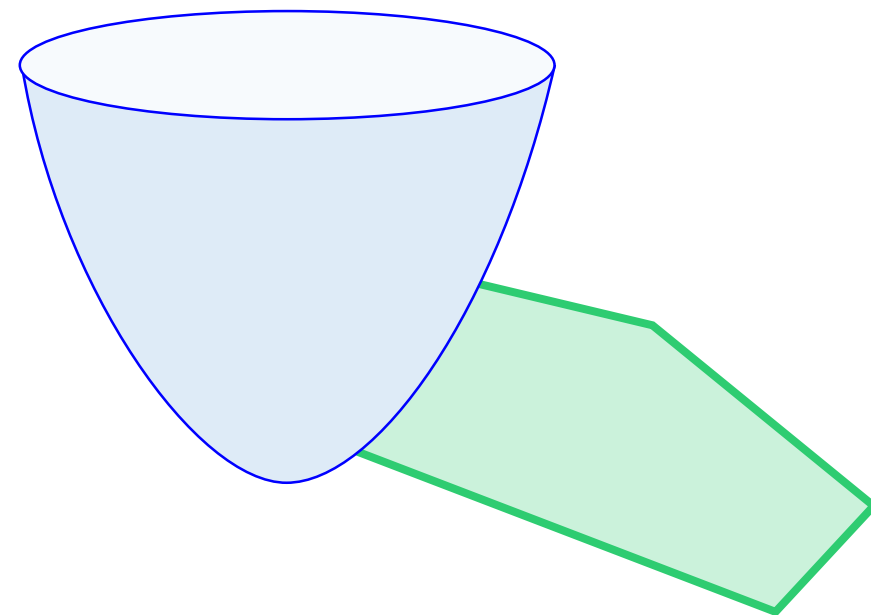
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

The Lagrangian

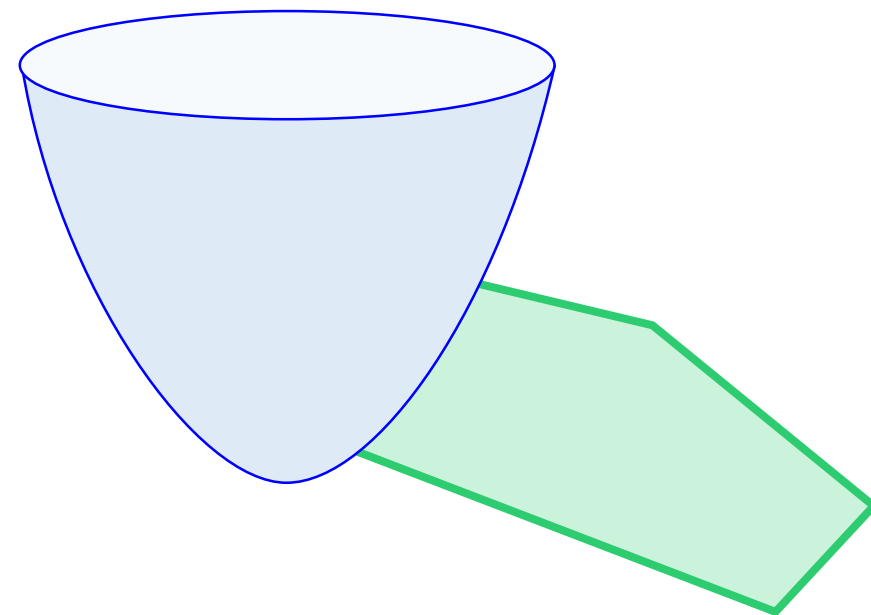
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$



The Lagrangian

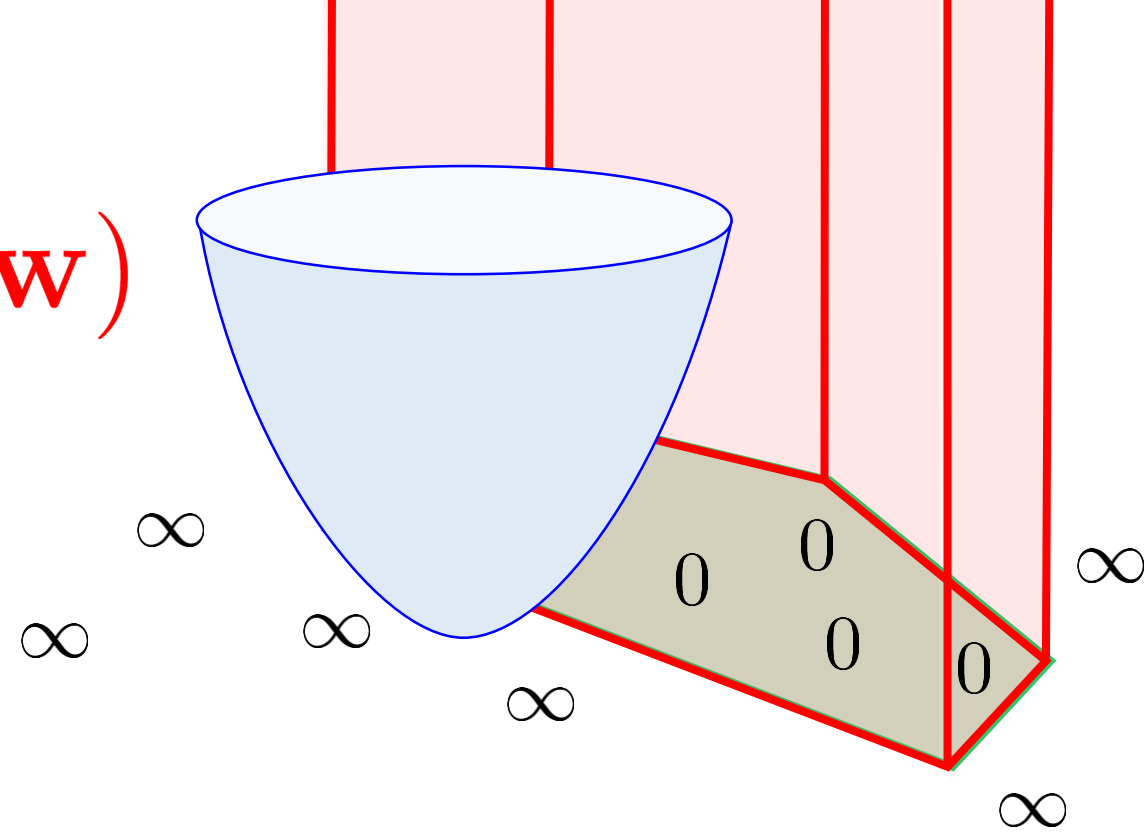
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + r(\mathbf{w})$$



The Lagrangian

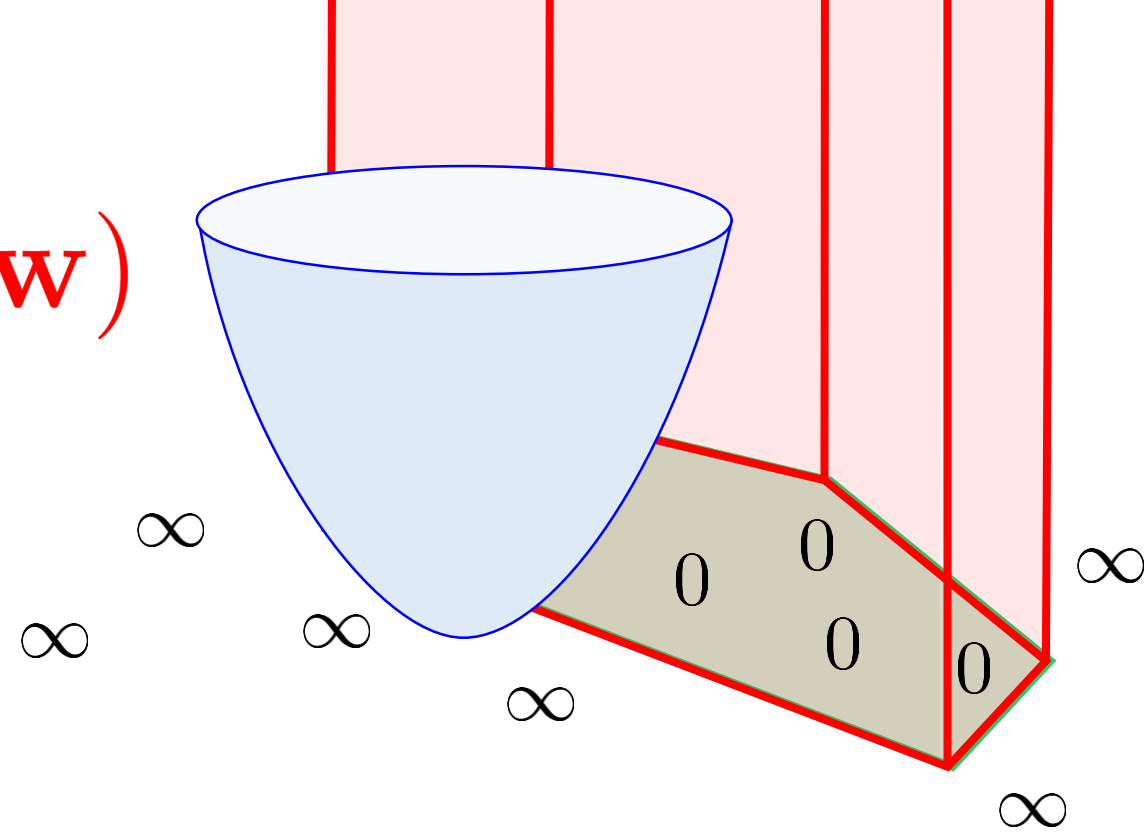
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + r(\mathbf{w})$$

Barrier Function



The Lagrangian

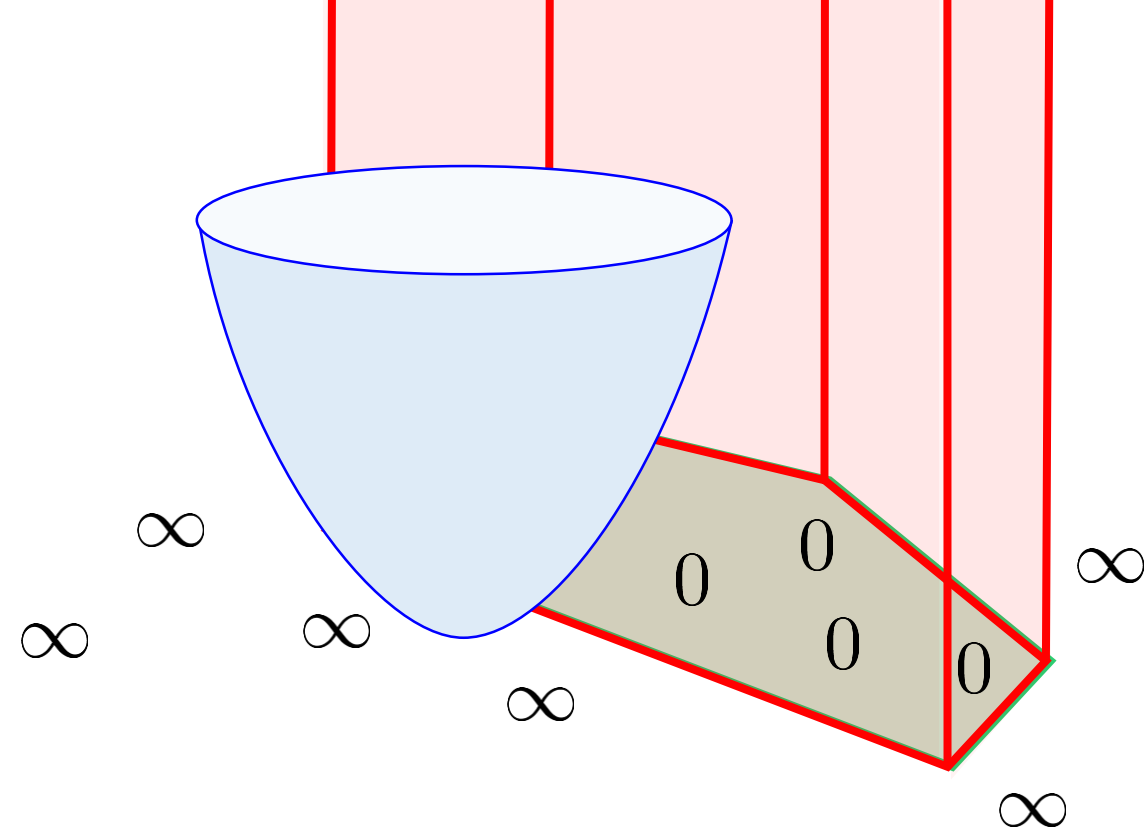
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + r(\mathbf{w})$$



The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

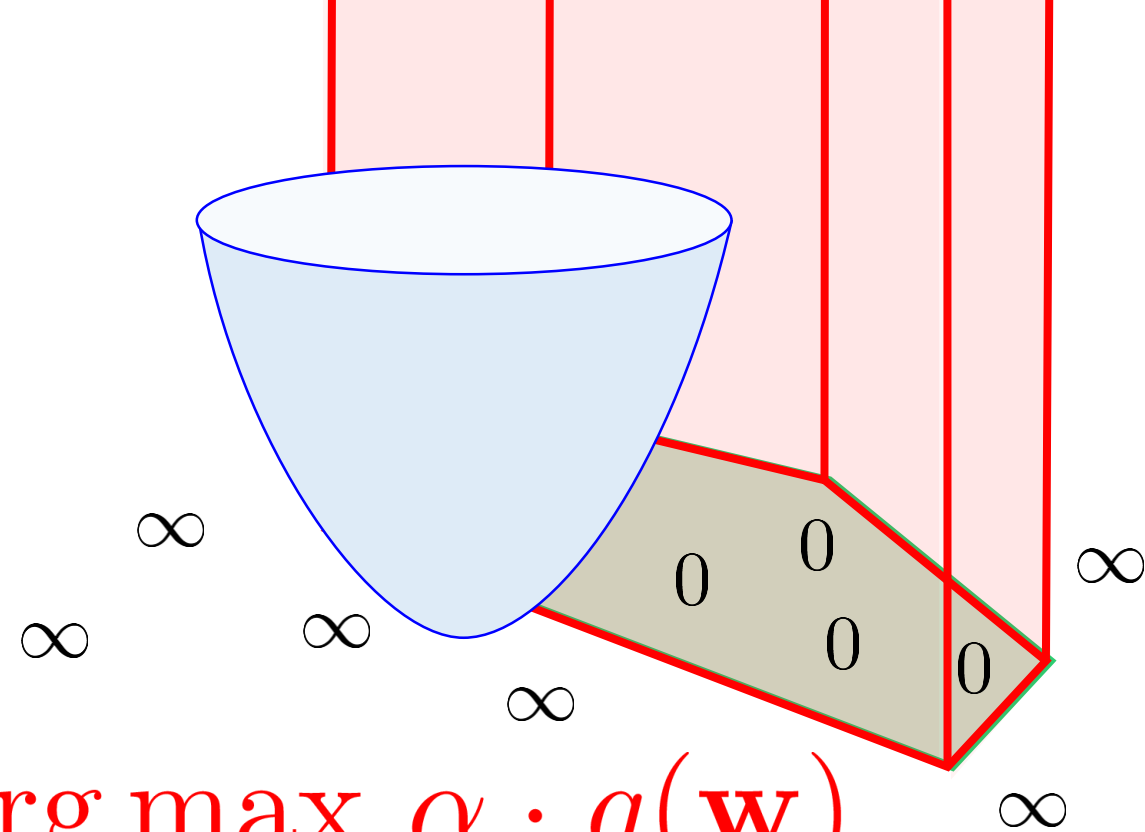


The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$\arg \max_{\alpha \geq 0} \alpha \cdot g(\mathbf{w})$$



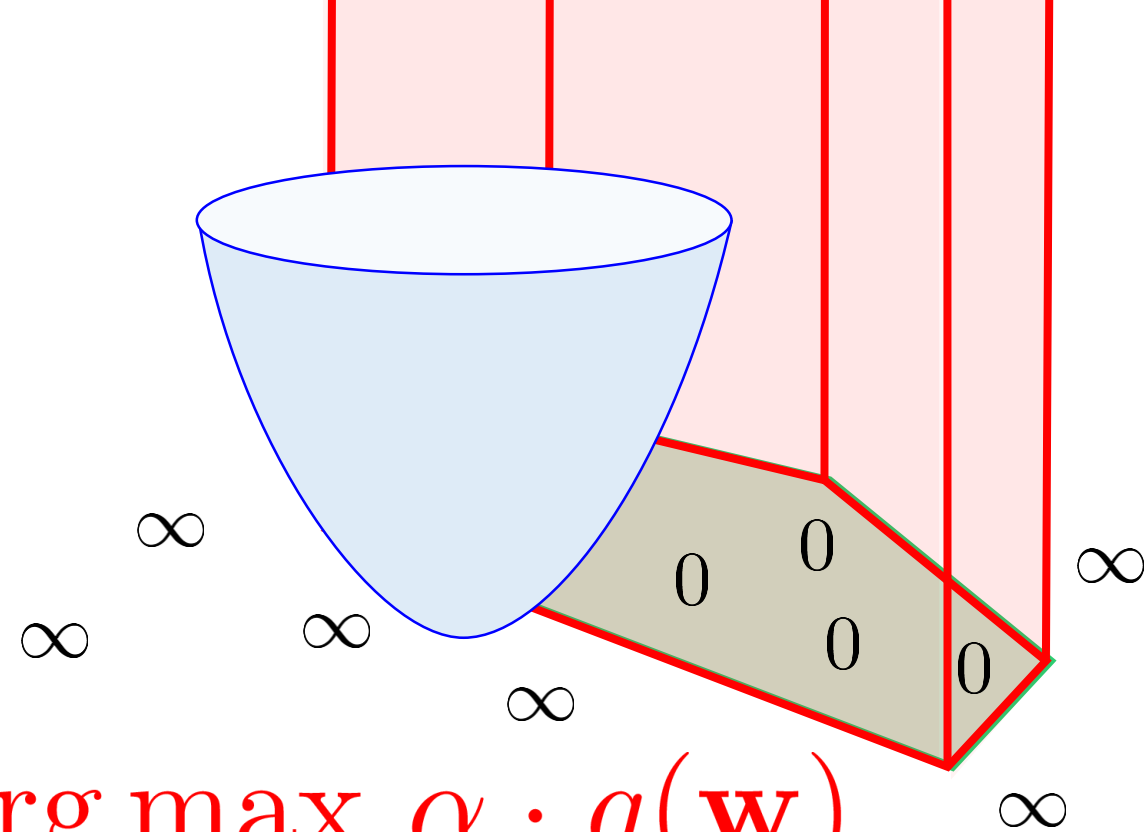
The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$\arg \max_{\alpha \geq 0} \alpha \cdot g(\mathbf{w})$$

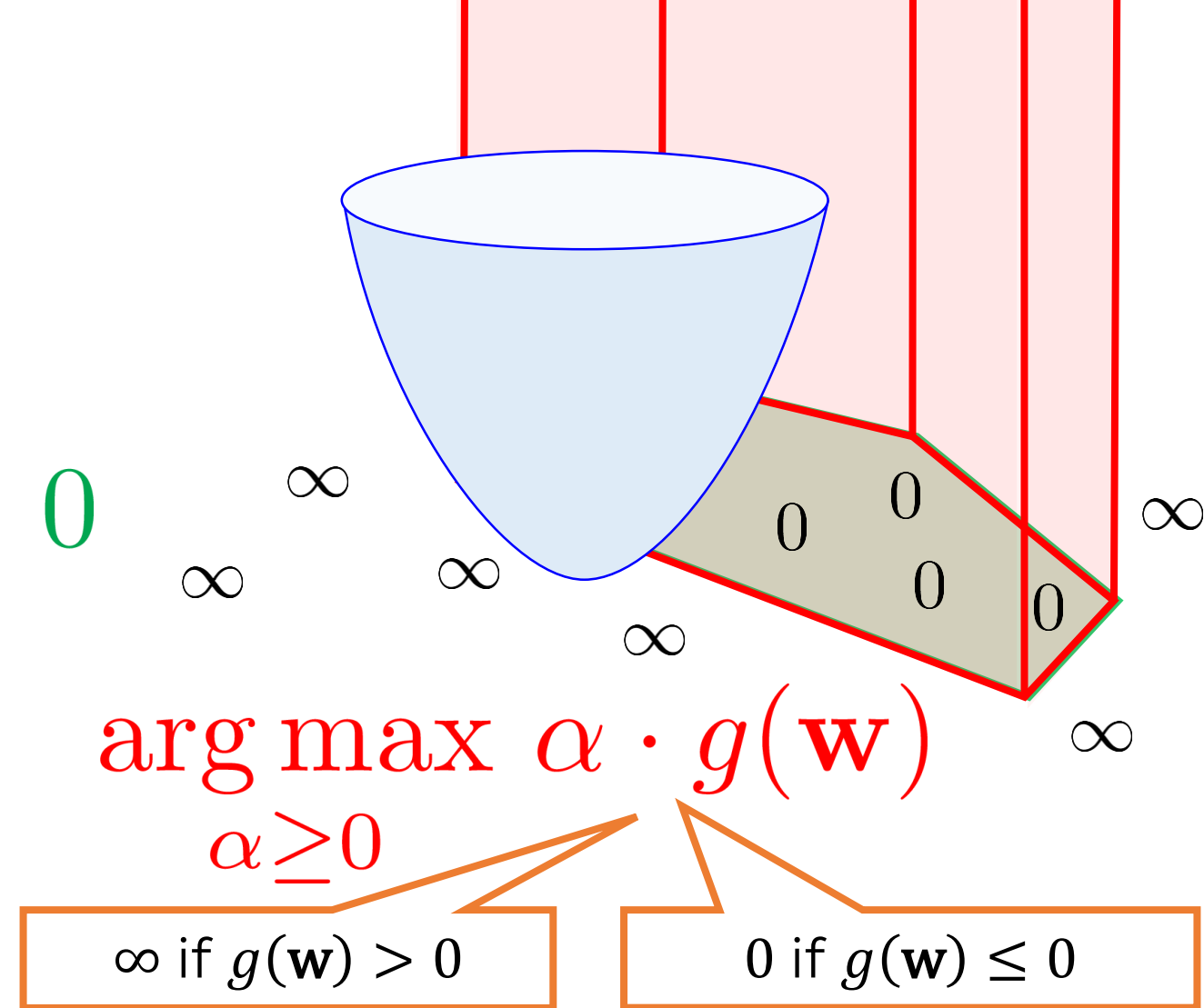
∞ if $g(\mathbf{w}) > 0$



The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$



The Lagrangian

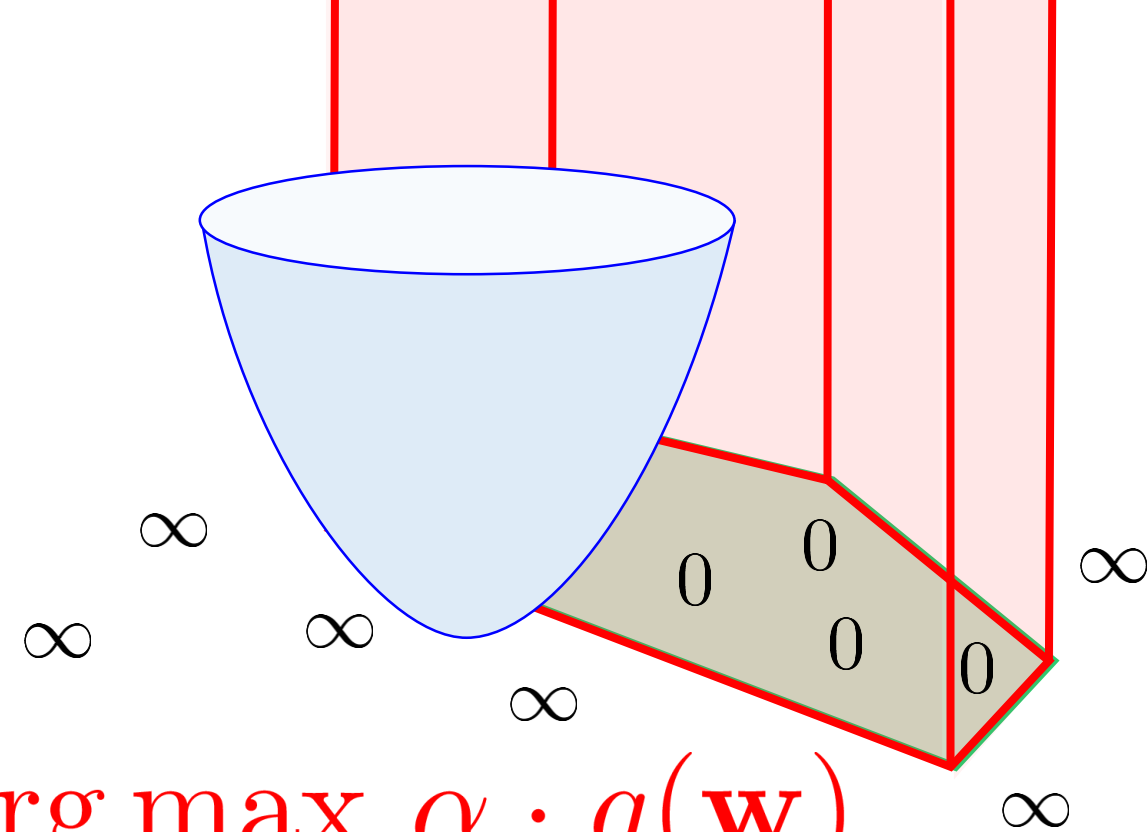
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$f(\mathbf{w}) + \arg \max_{\alpha \geq 0} \alpha \cdot g(\mathbf{w})$$

∞ if $g(\mathbf{w}) > 0$

0 if $g(\mathbf{w}) \leq 0$



The Lagrangian

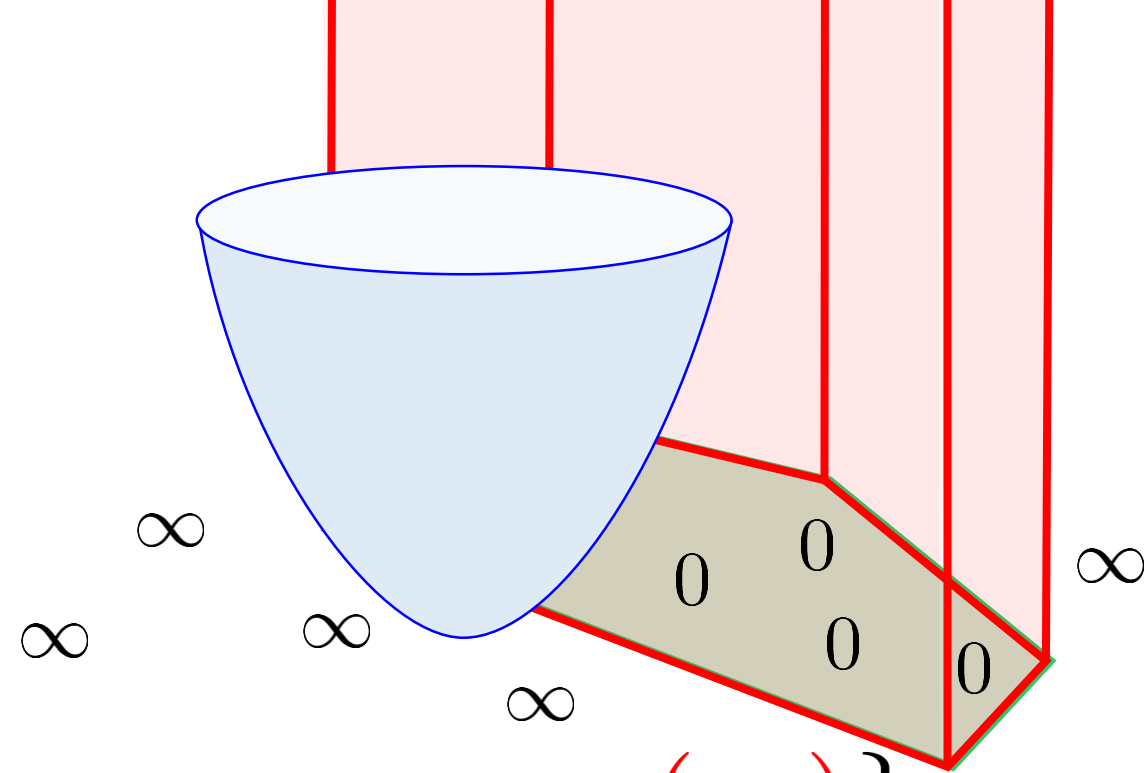
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$\arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ f(\mathbf{w}) + \arg \max_{\alpha \geq 0} \alpha \cdot g(\mathbf{w}) \right\}$$

∞ if $g(\mathbf{w}) > 0$

0 if $g(\mathbf{w}) \leq 0$



The Lagrangian

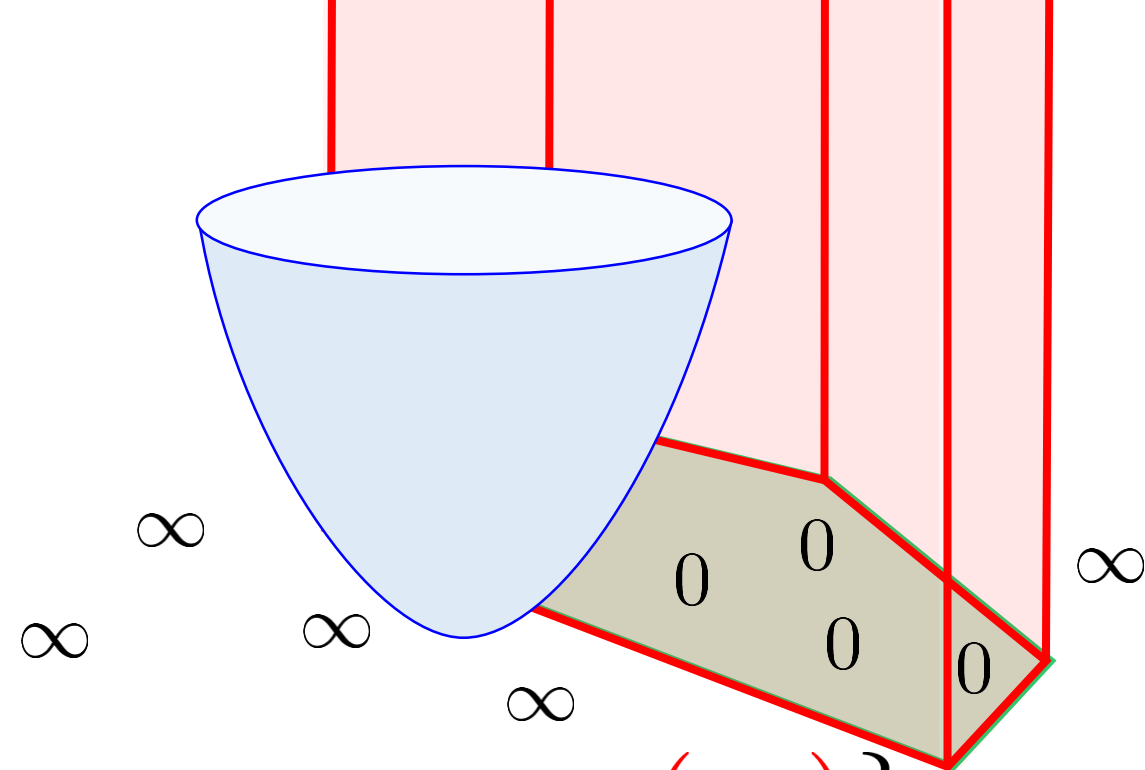
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ f(\mathbf{w}) + \arg \max_{\alpha \geq 0} \alpha \cdot g(\mathbf{w}) \right\}$$

∞ if $g(\mathbf{w}) > 0$

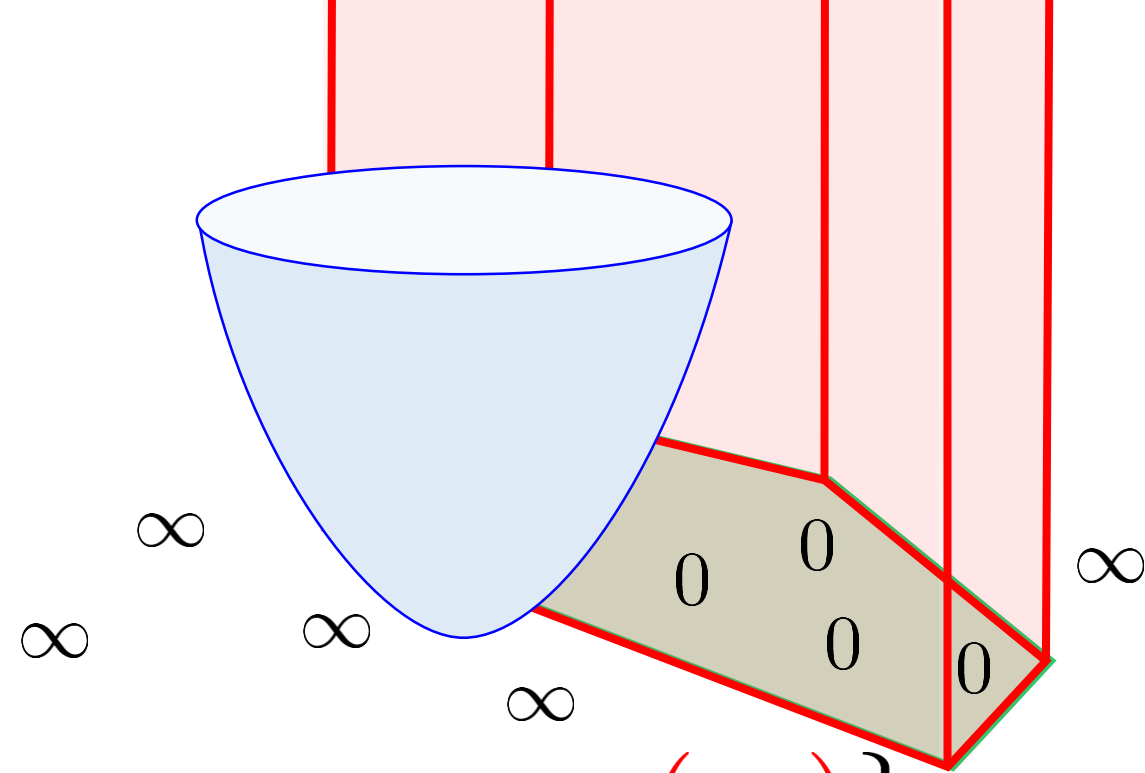
0 if $g(\mathbf{w}) \leq 0$



The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$



$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ f(\mathbf{w}) + \arg \max_{\alpha \geq 0} \alpha \cdot g(\mathbf{w}) \right\}$$

∞ if $g(\mathbf{w}) > 0$

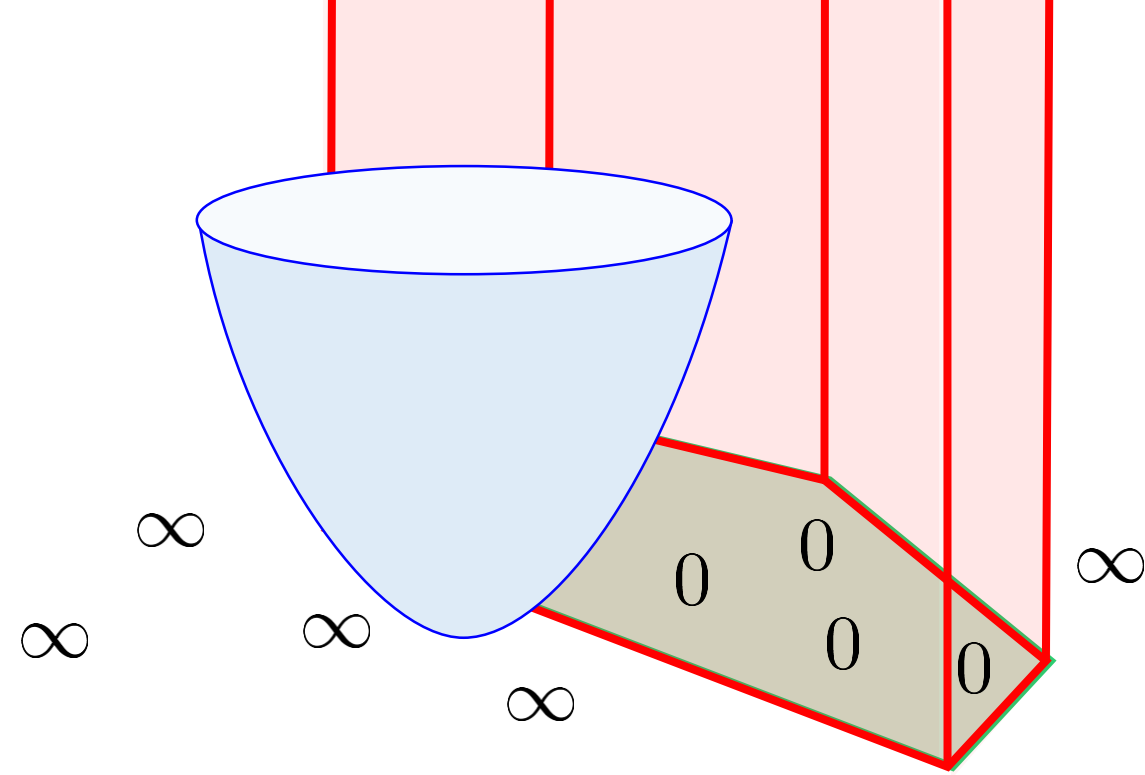
0 if $g(\mathbf{w}) \leq 0$

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\alpha \geq 0} \left\{ f(\mathbf{w}) + \alpha \cdot g(\mathbf{w}) \right\} \right\}$$

The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

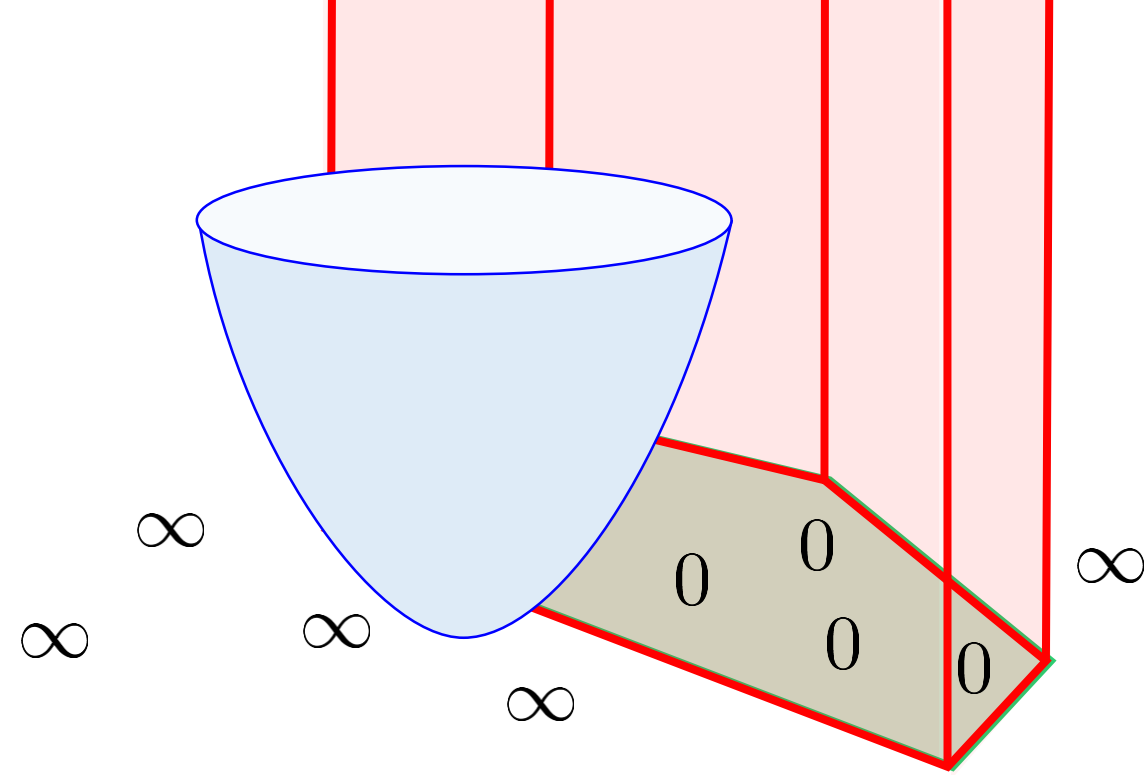


$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\alpha \geq 0} \{f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})\} \right\}$$

The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$



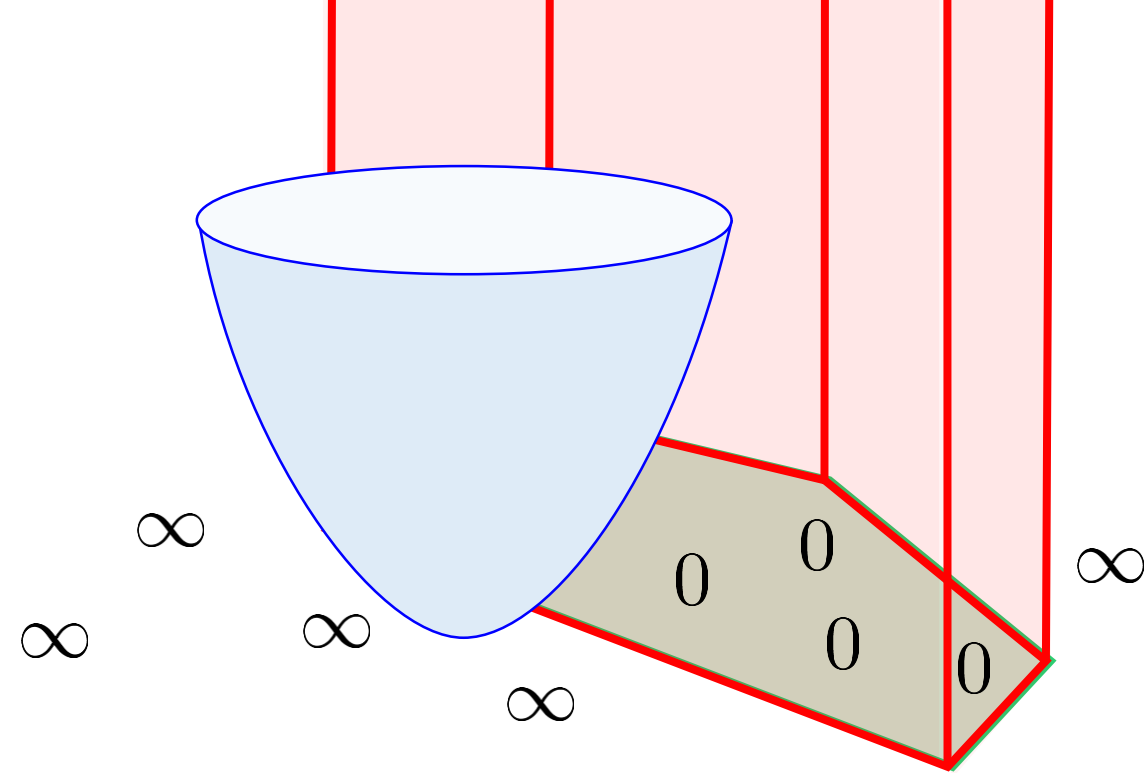
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\alpha \geq 0} \{f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})\} \right\}$$

$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$



$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\alpha \geq 0} \{f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})\} \right\}$$

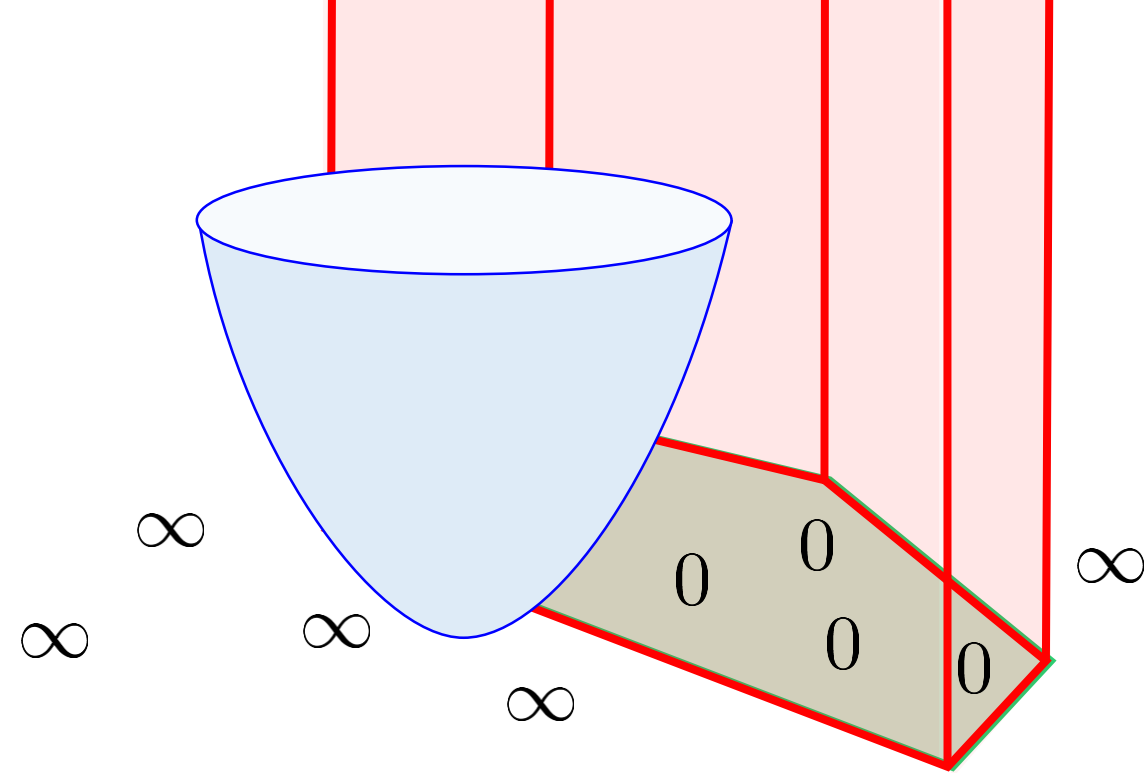
$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

Lagrangian

The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$



$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\alpha \geq 0} \{f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})\} \right\}$$

Lagrange multiplier

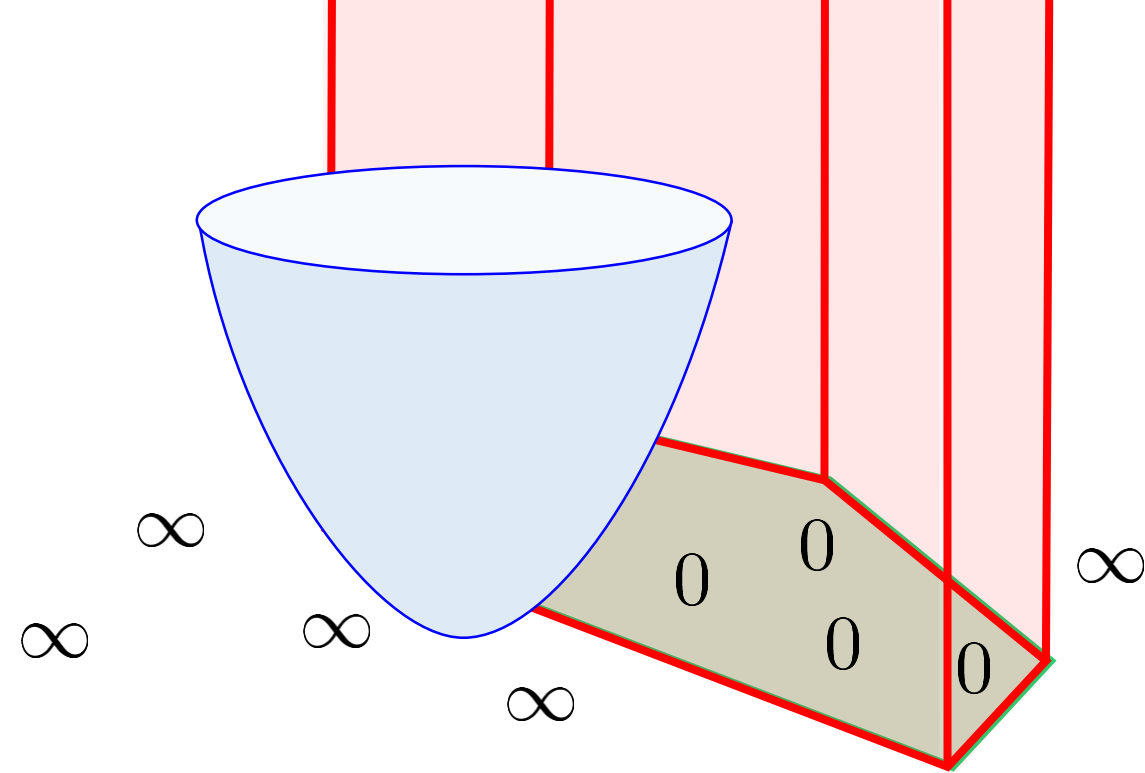
$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

Lagrangian

The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$



$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\alpha \geq 0} \left\{ \mathcal{L}(\mathbf{w}, \alpha) \right\} \right\}$$

Lagrange multiplier

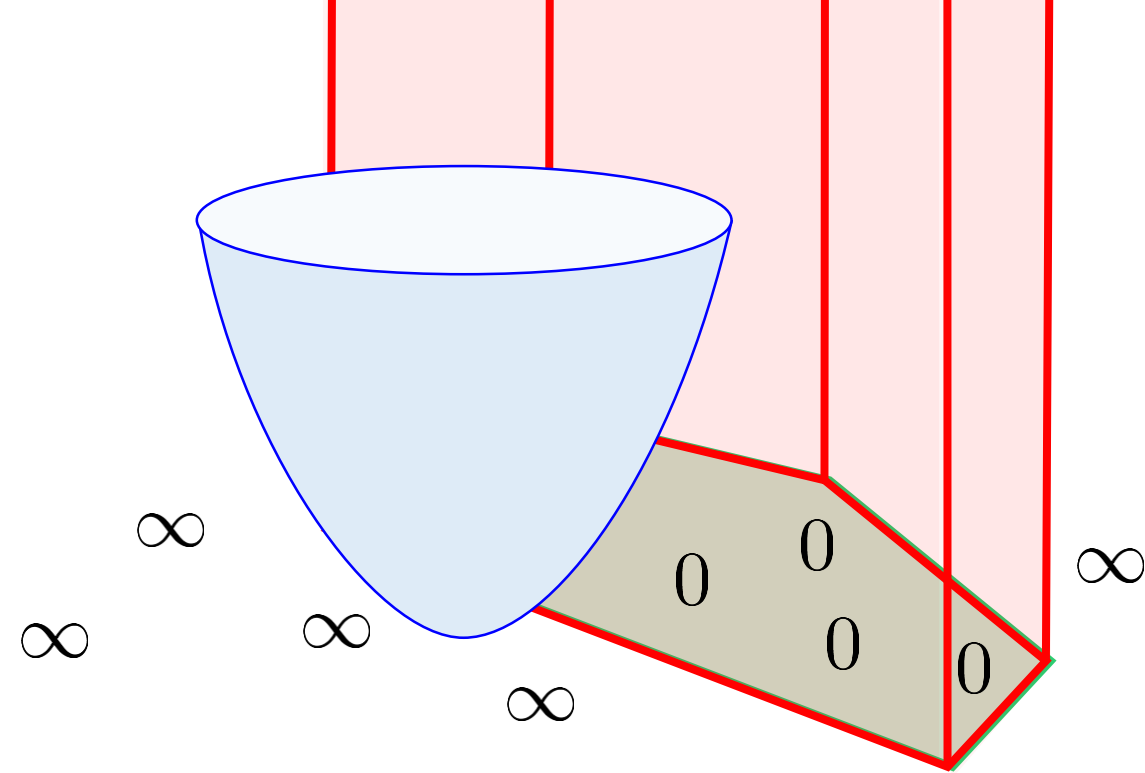
$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

Lagrangian

The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$



$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\alpha \geq 0} \left\{ \mathcal{L}(\mathbf{w}, \alpha) \right\} \right\}$$

Lagrange multiplier

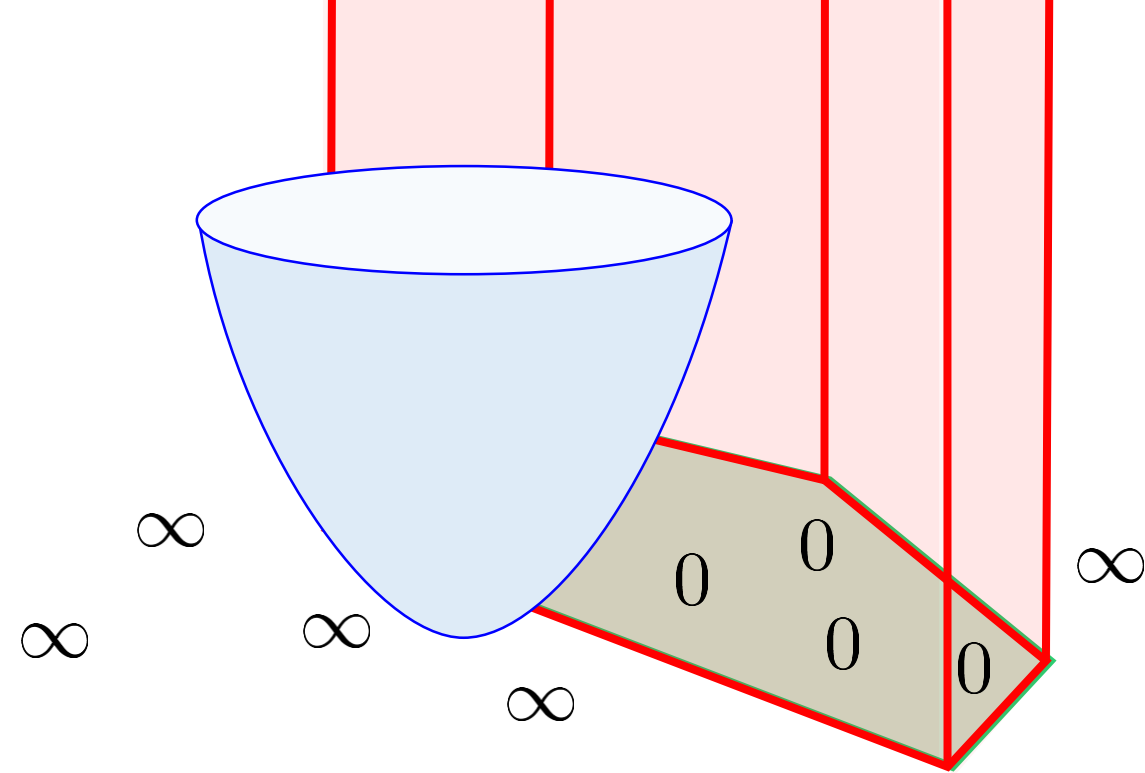
$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

Lagrangian

The Lagrangian

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$



$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\alpha \geq 0} \left\{ \mathcal{L}(\mathbf{w}, \alpha) \right\} \right\}$$

Lagrange multiplier

Primal problem

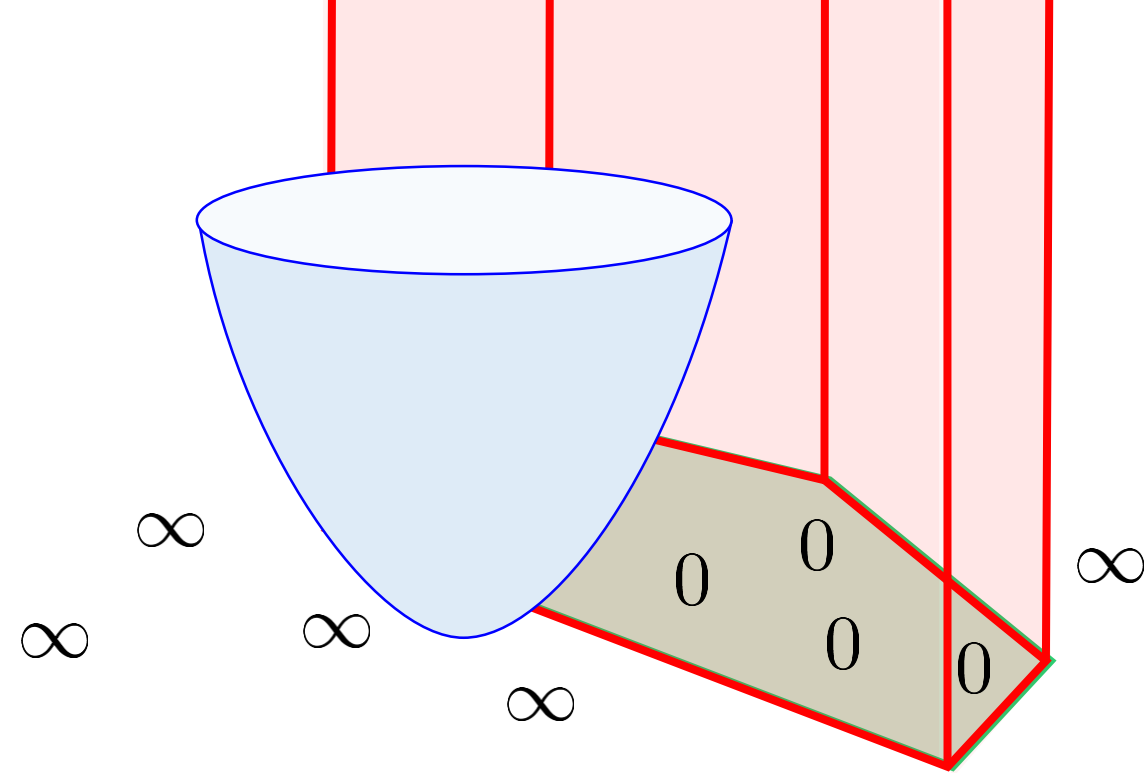
$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

Lagrangian

The Lagrangian

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$



$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\alpha \geq 0} \left\{ \mathcal{L}(\mathbf{w}, \alpha) \right\} \right\}$$

Lagrange multiplier

Primal problem

$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

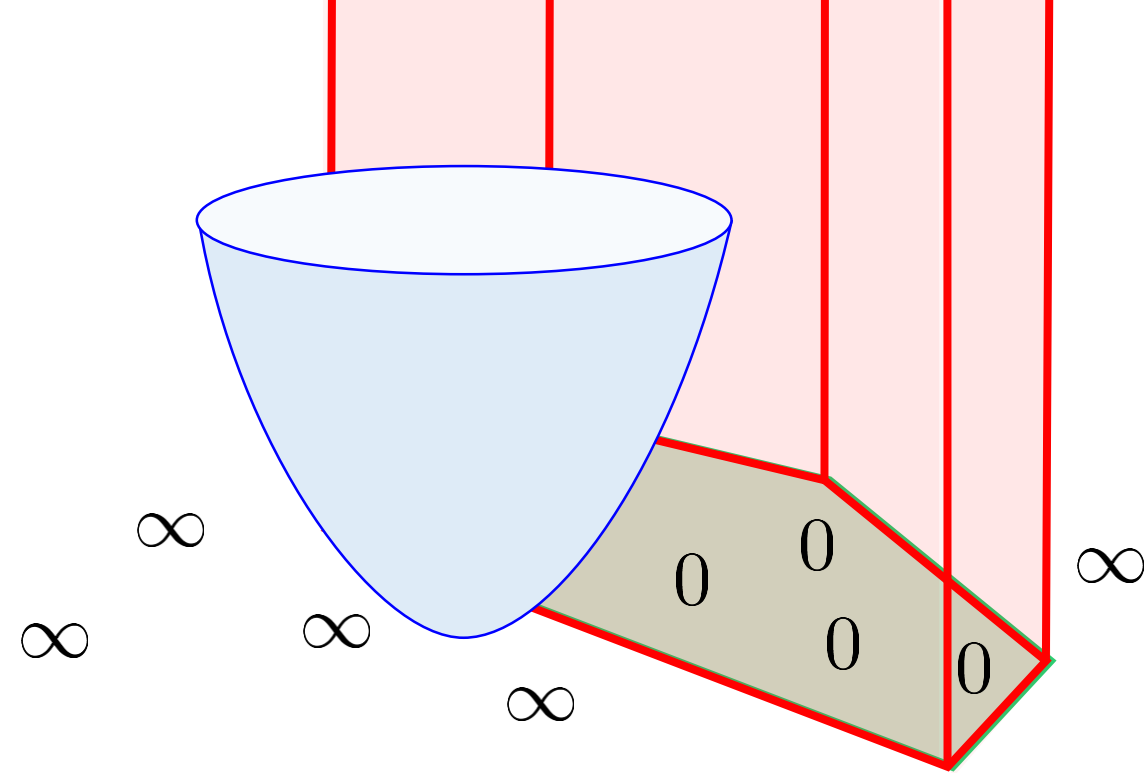
Lagrangian

The Lagrangian

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

Primal problem

$$\text{s.t. } g(\mathbf{w}) \leq 0$$



$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\alpha \geq 0} \left\{ \mathcal{L}(\mathbf{w}, \alpha) \right\} \right\}$$

Lagrange multiplier

Primal problem

$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

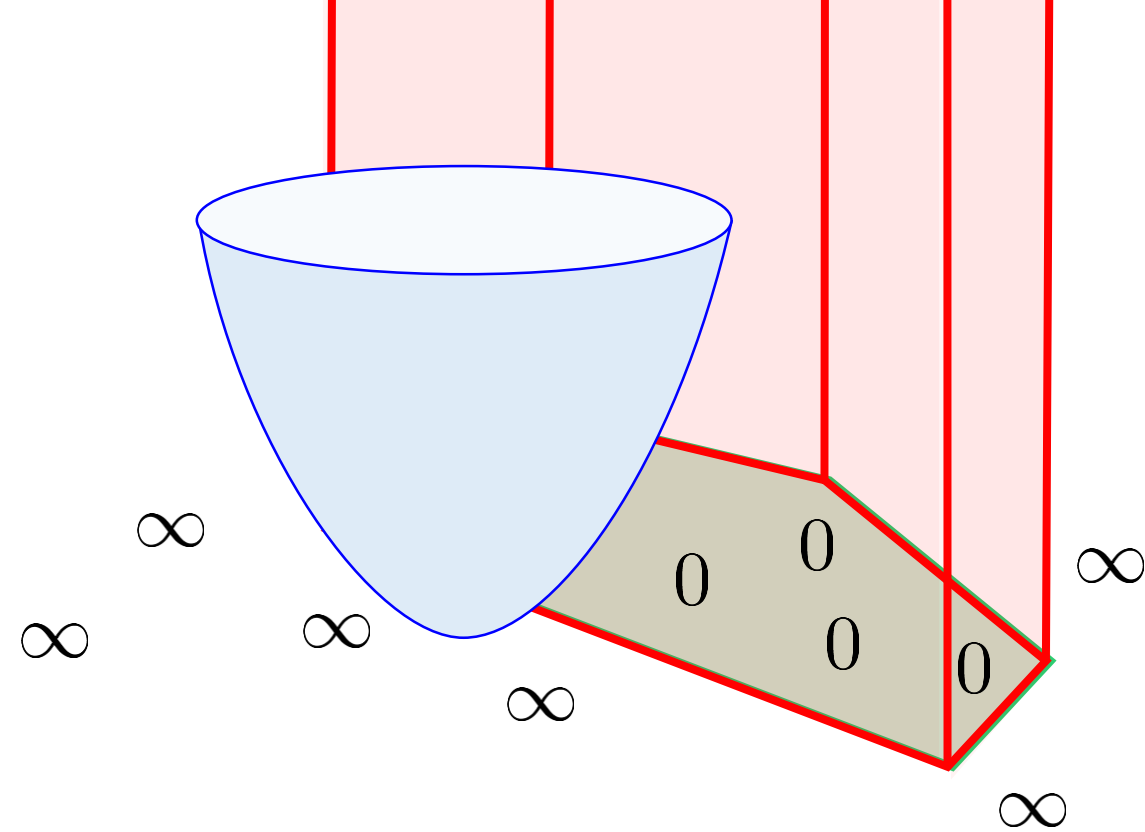
Lagrangian

The Lagrangian

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

Primal problem

$$\text{s.t. } g(\mathbf{w}) \leq 0$$



$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\alpha \geq 0} \{ \mathcal{L}(\mathbf{w}, \alpha) \} \right\}$$

Lagrange multiplier

Primal problem

$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

Lagrangian

The Lagrangian

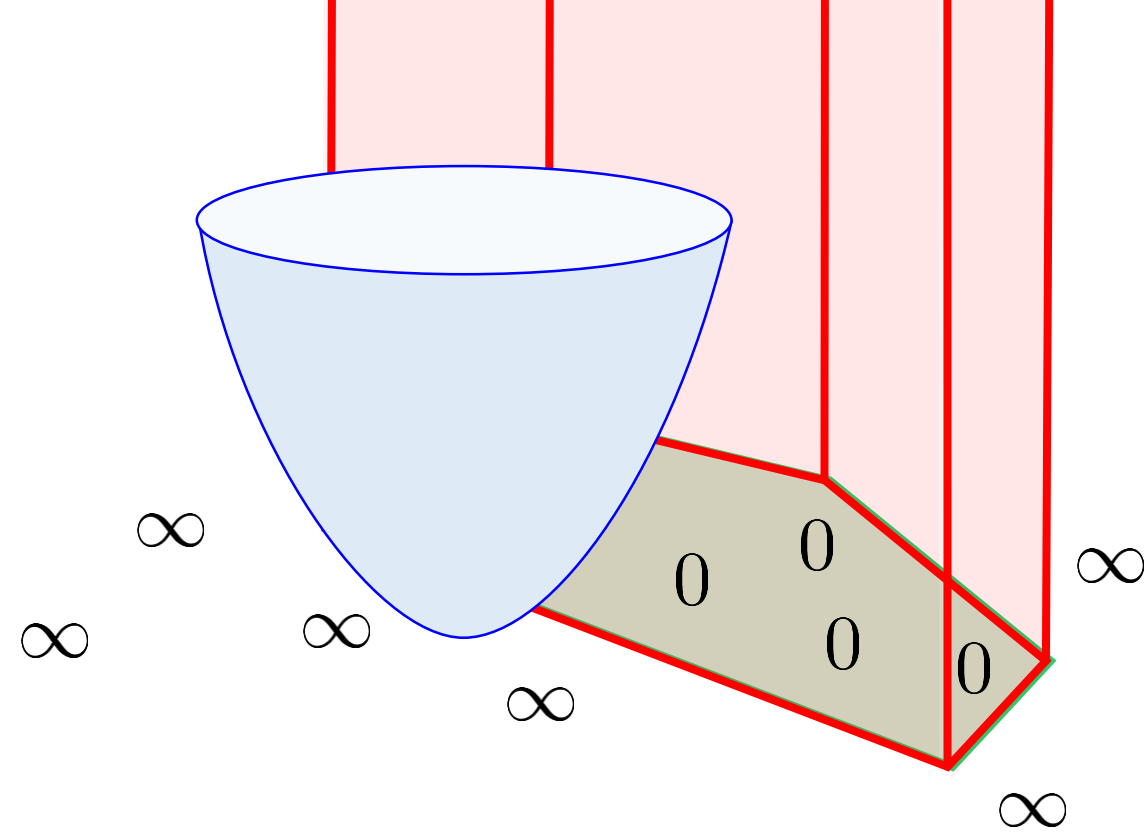
$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

Primal
problem

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\alpha \geq 0} \{ \mathcal{L}(\mathbf{w}, \alpha) \} \right\}$$

Primal problem

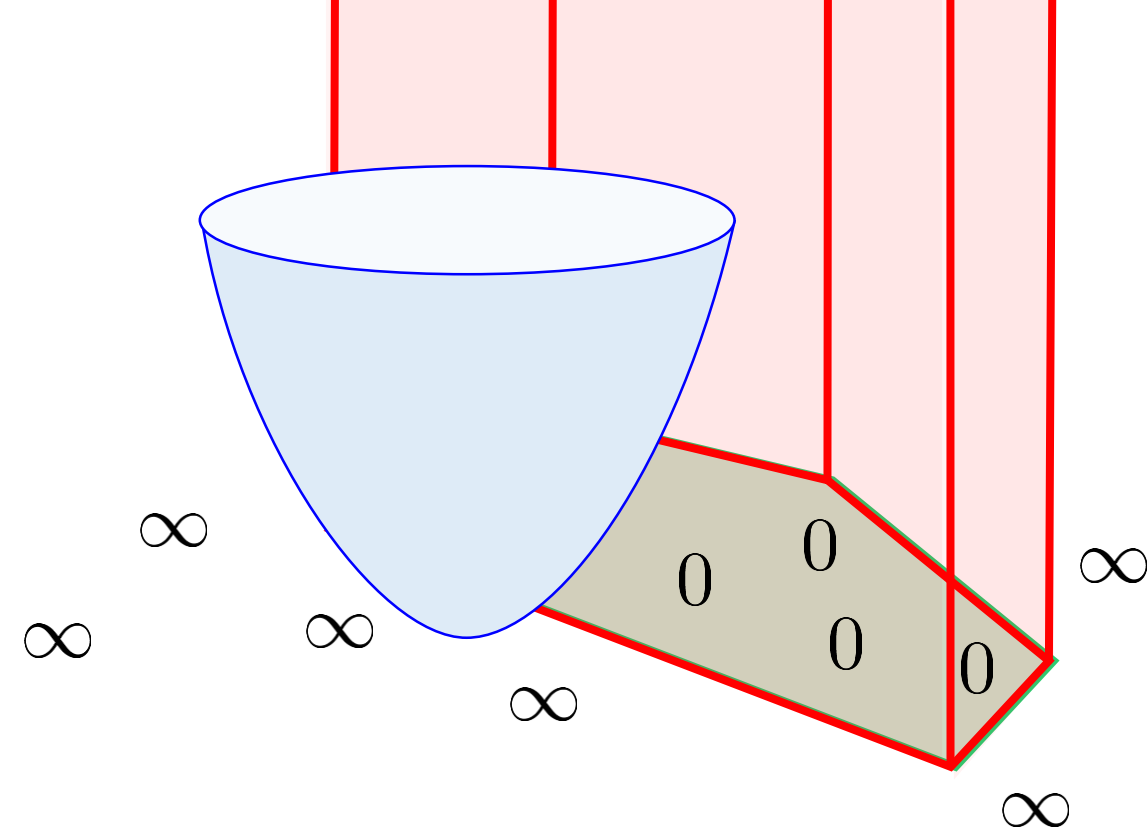


The Lagrangian

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

Primal
problem



$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\alpha \geq 0} \{ \mathcal{L}(\mathbf{w}, \alpha) \} \right\}$$

Primal problem

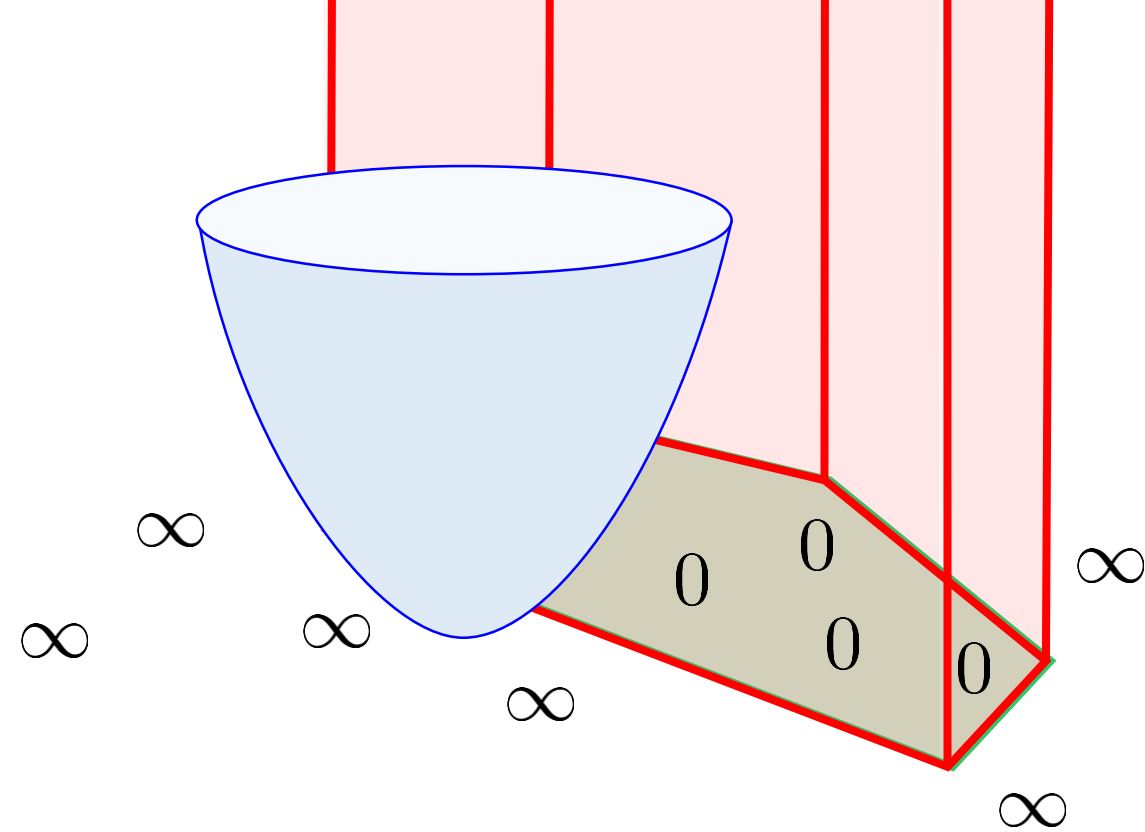
$$\hat{\mathbf{w}}_D = \arg \max_{\alpha \geq 0} \left\{ \arg \min_{\mathbf{w} \in \mathbb{R}^d} \{ \mathcal{L}(\mathbf{w}, \alpha) \} \right\}$$

The Lagrangian

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

Primal
problem

$$\text{s.t. } g(\mathbf{w}) \leq 0$$



$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\alpha \geq 0} \{ \mathcal{L}(\mathbf{w}, \alpha) \} \right\}$$

Primal problem

$$\hat{\mathbf{w}}_D = \arg \max_{\alpha \geq 0} \left\{ \arg \min_{\mathbf{w} \in \mathbb{R}^d} \{ \mathcal{L}(\mathbf{w}, \alpha) \} \right\}$$

Dual problem

The Lagrange Dual Problem

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

The Lagrange Dual Problem

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

The Lagrange Dual Problem

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

$$\hat{\mathbf{w}}_D = \arg \max_{\alpha \geq 0} \left\{ \arg \min_{\mathbf{w} \in \mathbb{R}^d} \{ \mathcal{L}(\mathbf{w}, \alpha) \} \right\}$$

The Lagrange Dual Problem

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

$$\hat{\mathbf{w}}_D = \arg \max_{\alpha \geq 0} \left\{ \arg \min_{\mathbf{w} \in \mathbb{R}^d} \{ \mathcal{L}(\mathbf{w}, \alpha) \} \right\}$$

$$\hat{\mathbf{w}}_P \stackrel{?}{=} \hat{\mathbf{w}}_D$$

The Lagrange Dual Problem

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$\hat{\mathbf{w}}_P \stackrel{?}{=} \hat{\mathbf{w}}_D$$

$\hat{\mathbf{w}}_p = \hat{\mathbf{w}}_d$ for nice problems*

$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

$$\hat{\mathbf{w}}_D = \arg \max_{\alpha \geq 0} \left\{ \arg \min_{\mathbf{w} \in \mathbb{R}^d} \{ \mathcal{L}(\mathbf{w}, \alpha) \} \right\}$$

The Lagrange Dual Problem

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

$$\hat{\mathbf{w}}_D = \arg \max_{\alpha \geq 0} \left\{ \arg \min_{\mathbf{w} \in \mathbb{R}^d} \{ \mathcal{L}(\mathbf{w}, \alpha) \} \right\}$$

$$\hat{\mathbf{w}}_P \stackrel{?}{=} \hat{\mathbf{w}}_D$$

$\hat{\mathbf{w}}_p = \hat{\mathbf{w}}_d$ for nice problems*

E.g. if $f(\cdot)$ is convex and $\{\mathbf{w} : g(\mathbf{w}) \leq 0\}$ is convex

The Lagrange Dual Problem

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

$$\hat{\mathbf{w}}_D = \arg \max_{\alpha \geq 0} \left\{ \arg \min_{\mathbf{w} \in \mathbb{R}^d} \{ \mathcal{L}(\mathbf{w}, \alpha) \} \right\}$$

$$\alpha_D \cdot g(\hat{\mathbf{w}}_D) = 0$$

$$\hat{\mathbf{w}}_P \stackrel{?}{=} \hat{\mathbf{w}}_D$$

$\hat{\mathbf{w}}_p = \hat{\mathbf{w}}_d$ for nice problems*

E.g. if $f(\cdot)$ is convex and $\{\mathbf{w} : g(\mathbf{w}) \leq 0\}$ is convex

The Lagrange Dual Problem

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

$$\hat{\mathbf{w}}_D = \arg \max_{\alpha \geq 0} \left\{ \arg \min_{\mathbf{w} \in \mathbb{R}^d} \{ \mathcal{L}(\mathbf{w}, \alpha) \} \right\}$$

$$\alpha_D \cdot g(\hat{\mathbf{w}}_D) = 0$$

$$\hat{\mathbf{w}}_P \stackrel{?}{=} \hat{\mathbf{w}}_D$$

$\hat{\mathbf{w}}_p = \hat{\mathbf{w}}_d$ for nice problems*

E.g. if $f(\cdot)$ is convex and $\{\mathbf{w} : g(\mathbf{w}) \leq 0\}$ is convex

The Lagrange Dual Problem

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

$$\hat{\mathbf{w}}_D = \arg \max_{\alpha \geq 0} \left\{ \arg \min_{\mathbf{w} \in \mathbb{R}^d} \{ \mathcal{L}(\mathbf{w}, \alpha) \} \right\}$$

The
maximizer

$$\alpha_D \cdot g(\hat{\mathbf{w}}_D) = 0$$

$$\hat{\mathbf{w}}_P \stackrel{?}{=} \hat{\mathbf{w}}_D$$

$\hat{\mathbf{w}}_p = \hat{\mathbf{w}}_d$ for nice
problems*

E.g. if $f(\cdot)$ is convex and
 $\{\mathbf{w} : g(\mathbf{w}) \leq 0\}$ is convex

The Lagrange Dual Problem

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g(\mathbf{w}) \leq 0$$

$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha \cdot g(\mathbf{w})$$

$$\hat{\mathbf{w}}_D = \arg \max_{\alpha \geq 0} \left\{ \arg \min_{\mathbf{w} \in \mathbb{R}^d} \{ \mathcal{L}(\mathbf{w}, \alpha) \} \right\}$$

The
maximizer

$$\alpha_D \cdot g(\hat{\mathbf{w}}_D) = 0$$

$$\hat{\mathbf{w}}_P \stackrel{?}{=} \hat{\mathbf{w}}_D$$

$\hat{\mathbf{w}}_p = \hat{\mathbf{w}}_d$ for nice
problems*

E.g. if $f(\cdot)$ is convex and
 $\{\mathbf{w} : g(\mathbf{w}) \leq 0\}$ is convex

Complimentary slackness
KKT Condition

Multiple Constraints!

Multiple Constraints!

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g_1(\mathbf{w}) \leq 0, g_2(\mathbf{w}) \leq 0 \\ h(\mathbf{w}) = 0$$

Multiple Constraints!

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

$$\text{s.t. } g_1(\mathbf{w}) \leq 0, g_2(\mathbf{w}) \leq 0$$

$$h(\mathbf{w}) = 0$$

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\substack{\alpha, \beta \geq 0 \\ \gamma \in \mathbb{R}}} \{f(\mathbf{w}) + \alpha \cdot g_1(\mathbf{w}) + \beta \cdot g_2(\mathbf{w}) + \gamma \cdot h(\mathbf{w})\} \right\}$$

Multiple Constraints!

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

One Lagrange multiplier for each constraint

$$\text{s.t. } g_1(\mathbf{w}) \leq 0, g_2(\mathbf{w}) \leq 0 \\ h(\mathbf{w}) = 0$$

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\substack{\alpha, \beta \geq 0 \\ \gamma \in \mathbb{R}}} \{f(\mathbf{w}) + \alpha \cdot g_1(\mathbf{w}) + \beta \cdot g_2(\mathbf{w}) + \gamma \cdot h(\mathbf{w})\} \right\}$$

Multiple Constraints!

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

One Lagrange multiplier for each constraint

$$\text{s.t. } g_1(\mathbf{w}) \leq 0, g_2(\mathbf{w}) \leq 0$$
$$h(\mathbf{w}) = 0$$

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\substack{\alpha, \beta \geq 0 \\ \gamma \in \mathbb{R}}} \{f(\mathbf{w}) + \alpha \cdot g_1(\mathbf{w}) + \beta \cdot g_2(\mathbf{w}) + \gamma \cdot h(\mathbf{w})\} \right\}$$

Doesn't allow
 $g_1(\mathbf{w}) > 0$

Multiple Constraints!

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

One Lagrange multiplier for each constraint

$$\text{s.t. } g_1(\mathbf{w}) \leq 0, g_2(\mathbf{w}) \leq 0$$
$$h(\mathbf{w}) = 0$$

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\substack{\alpha, \beta \geq 0 \\ \gamma \in \mathbb{R}}} \{ f(\mathbf{w}) + \alpha \cdot g_1(\mathbf{w}) + \beta \cdot g_2(\mathbf{w}) + \gamma \cdot h(\mathbf{w}) \} \right\}$$

Doesn't allow
 $g_1(\mathbf{w}) > 0$

Doesn't allow
 $g_2(\mathbf{w}) > 0$

Multiple Constraints!

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

One Lagrange multiplier for each constraint

$$\text{s.t. } g_1(\mathbf{w}) \leq 0, g_2(\mathbf{w}) \leq 0 \\ h(\mathbf{w}) = 0$$

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\substack{\alpha, \beta \geq 0 \\ \gamma \in \mathbb{R}}} \{f(\mathbf{w}) + \alpha \cdot g_1(\mathbf{w}) + \beta \cdot g_2(\mathbf{w}) + \gamma \cdot h(\mathbf{w})\} \right\}$$

Doesn't allow
 $g_1(\mathbf{w}) > 0$

Doesn't allow
 $g_2(\mathbf{w}) > 0$

Doesn't allow
 $h(\mathbf{w}) \neq 0$

Multiple Constraints!

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

One Lagrange multiplier for each constraint

$$\text{s.t. } g_1(\mathbf{w}) \leq 0, g_2(\mathbf{w}) \leq 0 \\ h(\mathbf{w}) = 0$$

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\substack{\alpha, \beta \geq 0 \\ \gamma \in \mathbb{R}}} \{f(\mathbf{w}) + \alpha \cdot g_1(\mathbf{w}) + \beta \cdot g_2(\mathbf{w}) + \gamma \cdot h(\mathbf{w})\} \right\}$$

Multiple Constraints!

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

One Lagrange multiplier for each constraint

$$\text{s.t. } g_1(\mathbf{w}) \leq 0, g_2(\mathbf{w}) \leq 0 \\ h(\mathbf{w}) = 0$$

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\substack{\alpha, \beta \geq 0 \\ \gamma \in \mathbb{R}}} \{f(\mathbf{w}) + \alpha \cdot g_1(\mathbf{w}) + \beta \cdot g_2(\mathbf{w}) + \gamma \cdot h(\mathbf{w})\} \right\}$$

$$\hat{\mathbf{w}}_D = \arg \max_{\substack{\alpha, \beta \geq 0 \\ \gamma \in \mathbb{R}}} \left\{ \arg \min_{\mathbf{w} \in \mathbb{R}^d} \{f(\mathbf{w}) + \alpha \cdot g_1(\mathbf{w}) + \beta \cdot g_2(\mathbf{w}) + \gamma \cdot h(\mathbf{w})\} \right\}$$

Multiple Constraints!

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

Complimentary Slackness

$$\alpha_D \cdot g_1(\hat{\mathbf{w}}_D) = 0$$

$$\beta_D \cdot g_2(\hat{\mathbf{w}}_D) = 0$$

$$\gamma_D \cdot h(\hat{\mathbf{w}}_D) = 0$$

One Lagrange multiplier for each constraint

$$\text{s.t. } g_1(\mathbf{w}) \leq 0, g_2(\mathbf{w}) \leq 0$$
$$h(\mathbf{w}) = 0$$

$$\hat{\mathbf{w}}_P = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ \arg \max_{\substack{\alpha, \beta \geq 0 \\ \gamma \in \mathbb{R}}} \{f(\mathbf{w}) + \alpha \cdot g_1(\mathbf{w}) + \beta \cdot g_2(\mathbf{w}) + \gamma \cdot h(\mathbf{w})\} \right\}$$

$$\hat{\mathbf{w}}_D = \arg \max_{\substack{\alpha, \beta \geq 0 \\ \gamma \in \mathbb{R}}} \left\{ \arg \min_{\mathbf{w} \in \mathbb{R}^d} \{f(\mathbf{w}) + \alpha \cdot g_1(\mathbf{w}) + \beta \cdot g_2(\mathbf{w}) + \gamma \cdot h(\mathbf{w})\} \right\}$$

SVM Revisited

SVM Revisited

$$\begin{aligned}\hat{\mathbf{w}} &= \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t. } y^i \langle \mathbf{w}, \mathbf{x}^i \rangle &\geq 1\end{aligned}$$

SVM Revisited

$$\begin{aligned}\hat{\mathbf{w}} &= \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t. } 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle &\leq 0\end{aligned}$$

SVM Revisited

$$\begin{aligned}\hat{\mathbf{w}} &= \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t. } 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle &\leq 0\end{aligned}$$

n constraints, so n
Lagrange multipliers

SVM Revisited

n constraints, so n
Lagrange multipliers

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\text{s.t. } 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle \leq 0$$

$$\mathcal{L}(\mathbf{w}, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle)$$

SVM Revisited

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\text{s.t. } 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle \leq 0$$

$$\mathcal{L}(\mathbf{w}, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle)$$

n constraints, so n
Lagrange multipliers

$$\boldsymbol{\alpha} \in \mathbb{R}^n$$

SVM Revisited

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\text{s.t. } 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle \leq 0$$

$$\mathcal{L}(\mathbf{w}, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle)$$

$$\hat{\mathbf{w}} = \hat{\mathbf{w}}_D = \max_{\boldsymbol{\alpha} \geq 0} \left\{ \min_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \boldsymbol{\alpha}) \right\}$$

n constraints, so n
Lagrange multipliers

$$\boldsymbol{\alpha} \in \mathbb{R}^n$$

SVM Revisited

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\text{s.t. } 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle \leq 0$$

$$\mathcal{L}(\mathbf{w}, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle)$$

$$\hat{\mathbf{w}} = \hat{\mathbf{w}}_D = \max_{\boldsymbol{\alpha} \geq 0} \left\{ \min_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \boldsymbol{\alpha}) \right\}$$

n constraints, so n
Lagrange multipliers

$$\boldsymbol{\alpha} \in \mathbb{R}^n$$

$\boldsymbol{\alpha} \geq 0$ is notation
for $\alpha_i \geq 0$ for all i

SVM Revisited

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\text{s.t. } 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle \leq 0$$

$$\mathcal{L}(\mathbf{w}, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle)$$

$$\hat{\mathbf{w}} = \hat{\mathbf{w}}_D = \max_{\boldsymbol{\alpha} \geq 0} \left\{ \min_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \boldsymbol{\alpha}) \right\}$$

n constraints, so n
Lagrange multipliers

$$\boldsymbol{\alpha} \in \mathbb{R}^n$$

Inner problem is an
Unconstrained problem
so Gradients must
vanish at optimum

$\boldsymbol{\alpha} \geq 0$ is notation
for $\alpha_i \geq 0$ for all i

SVM Revisited

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\text{s.t. } 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle \leq 0$$

$$\mathcal{L}(\mathbf{w}, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle)$$

$$\hat{\mathbf{w}} = \hat{\mathbf{w}}_D = \max_{\boldsymbol{\alpha} \geq 0} \left\{ \min_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \boldsymbol{\alpha}) \right\}$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \mathbf{0}$$

n constraints, so n
Lagrange multipliers

$$\boldsymbol{\alpha} \in \mathbb{R}^n$$

Inner problem is an
Unconstrained problem
so Gradients must
vanish at optimum

$\boldsymbol{\alpha} \geq 0$ is notation
for $\alpha_i \geq 0$ for all i

SVM Revisited

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\text{s.t. } 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle \leq 0$$

$$\mathcal{L}(\mathbf{w}, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle)$$

$$\hat{\mathbf{w}} = \hat{\mathbf{w}}_D = \max_{\boldsymbol{\alpha} \geq 0} \left\{ \min_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \boldsymbol{\alpha}) \right\}$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \mathbf{0} \Rightarrow \hat{\mathbf{w}}_D = \sum_{i=1}^n \alpha_i y^i \mathbf{x}^i$$

n constraints, so n
Lagrange multipliers

$$\boldsymbol{\alpha} \in \mathbb{R}^n$$

Inner problem is an
Unconstrained problem
so Gradients must
vanish at optimum

$\boldsymbol{\alpha} \geq 0$ is notation
for $\alpha_i \geq 0$ for all i

SVM Revisited

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\text{s.t. } 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle \leq 0$$

$$\mathcal{L}(\mathbf{w}, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle)$$

$$\hat{\mathbf{w}} = \hat{\mathbf{w}}_D = \max_{\boldsymbol{\alpha} \geq 0} \left\{ \min_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \boldsymbol{\alpha}) \right\}$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \mathbf{0} \Rightarrow \hat{\mathbf{w}}_D = \sum_{i=1}^n \alpha_i y^i \mathbf{x}^i$$

Optimal model $\hat{\mathbf{w}}$ is a weighted sum of training points!

n constraints, so n Lagrange multipliers

$$\boldsymbol{\alpha} \in \mathbb{R}^n$$

Inner problem is an Unconstrained problem so Gradients must vanish at optimum

$\boldsymbol{\alpha} \geq 0$ is notation for $\alpha_i \geq 0$ for all i



SVM Revisited

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\text{s.t. } 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle \leq 0$$

$$\mathcal{L}(\mathbf{w}, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle)$$

$$\hat{\mathbf{w}} = \hat{\mathbf{w}}_D = \max_{\boldsymbol{\alpha} \geq 0} \left\{ \min_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \boldsymbol{\alpha}) \right\}$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \mathbf{0} \Rightarrow \hat{\mathbf{w}}_D = \sum_{i=1}^n \alpha_i y^i \mathbf{x}^i$$

Optimal model $\hat{\mathbf{w}}$ is a weighted sum of training points!

Points with $\alpha_i \neq 0$ support vectors

n constraints, so n Lagrange multipliers

$$\boldsymbol{\alpha} \in \mathbb{R}^n$$

Inner problem is an Unconstrained problem so Gradients must vanish at optimum

$\boldsymbol{\alpha} \geq 0$ is notation for $\alpha_i \geq 0$ for all i



SVM Revisited

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\text{s.t. } 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle \leq 0$$

$$\mathcal{L}(\mathbf{w}, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle)$$

$$\hat{\mathbf{w}} = \hat{\mathbf{w}}_D = \max_{\boldsymbol{\alpha} \geq 0} \left\{ \min_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \boldsymbol{\alpha}) \right\}$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \mathbf{0} \Rightarrow \hat{\mathbf{w}}_D = \sum_{i=1}^n \alpha_i y^i \mathbf{x}^i$$

Optimal model $\hat{\mathbf{w}}$ is a weighted sum of training points!

Points with $\alpha_i \neq 0$ support vectors

n constraints, so n Lagrange multipliers

$$\boldsymbol{\alpha} \in \mathbb{R}^n$$

Inner problem is an Unconstrained problem so Gradients must vanish at optimum

$\boldsymbol{\alpha} \geq 0$ is notation for $\alpha_i \geq 0$ for all i

Lets substitute \mathbf{w} in \mathcal{L} and eliminate it!



SVM Revisited

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\text{s.t. } 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle \leq 0$$

$$\mathcal{L}(\boldsymbol{\alpha}) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle$$

$$\hat{\mathbf{w}} = \hat{\mathbf{w}}_D = \max_{\boldsymbol{\alpha} \geq 0} \left\{ \min_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \boldsymbol{\alpha}) \right\}$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \mathbf{0} \Rightarrow \hat{\mathbf{w}}_D = \sum_{i=1}^n \alpha_i y^i \mathbf{x}^i$$

Optimal model $\hat{\mathbf{w}}$ is a weighted sum of training points!

Points with $\alpha_i \neq 0$ support vectors

n constraints, so n Lagrange multipliers

$$\boldsymbol{\alpha} \in \mathbb{R}^n$$

Inner problem is an Unconstrained problem so Gradients must vanish at optimum

$\boldsymbol{\alpha} \geq 0$ is notation for $\alpha_i \geq 0$ for all i

Lets substitute \mathbf{w} in \mathcal{L} and eliminate it!



SVM Revisited

$$\begin{aligned}\hat{\mathbf{w}} &= \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t. } &1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle \leq 0\end{aligned}$$

SVM Revisited

$$\begin{aligned}\hat{\mathbf{w}} &= \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t. } &1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle \leq 0\end{aligned}$$

$$\hat{\mathbf{w}} = \hat{\mathbf{w}}_{\text{D}} = \sum_{i=1}^n \hat{\alpha}_i y^i \mathbf{x}^i$$

SVM Revisited

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\text{s.t. } 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle \leq 0$$

$$\hat{\mathbf{w}} = \hat{\mathbf{w}}_D = \sum_{i=1}^n \hat{\alpha}_i y^i \mathbf{x}^i$$

$$\hat{\alpha} = \arg \max_{\alpha \in \mathbb{R}^n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle$$

$$\text{s.t. } \alpha_i \geq 0$$

SVM Revisited

$$\begin{aligned}\hat{\mathbf{w}} &= \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t. } & 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle \leq 0\end{aligned}$$

$$\hat{\mathbf{w}} = \hat{\mathbf{w}}_D = \sum_{i=1}^n \hat{\alpha}_i y^i \mathbf{x}^i$$

$$\begin{aligned}\hat{\alpha} &= \arg \max_{\alpha \in \mathbb{R}^n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle \\ \text{s.t. } & \alpha_i \geq 0\end{aligned}$$

$$\hat{\alpha}_i (1 - y^i \langle \hat{\mathbf{w}}_D, \mathbf{x}^i \rangle) = 0$$

SVM Revisited

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2$$
$$\text{s.t. } 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle \leq 0$$

$$\hat{\mathbf{w}} = \hat{\mathbf{w}}_D = \sum_{i=1}^n \hat{\alpha}_i y^i \mathbf{x}^i$$

$$\hat{\alpha} = \arg \max_{\alpha \in \mathbb{R}^n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle$$

$$\text{s.t. } \alpha_i \geq 0$$

$$\hat{\alpha}_i (1 - y^i \langle \hat{\mathbf{w}}_D, \mathbf{x}^i \rangle) = 0$$

Complimentary
slackness

SVM Revisited

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2$$
$$\text{s.t. } 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle \leq 0$$

$$\hat{\mathbf{w}} = \hat{\mathbf{w}}_D = \sum_{i=1}^n \hat{\alpha}_i y^i \mathbf{x}^i$$

$$\hat{\alpha} = \arg \max_{\alpha \in \mathbb{R}^n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle$$

$$\text{s.t. } \alpha_i \geq 0$$

$$\hat{\alpha}_i (1 - y^i \langle \hat{\mathbf{w}}_D, \mathbf{x}^i \rangle) = 0$$

Complimentary
slackness

If $\hat{\alpha}_i \neq 0$, then
 $y^i \langle \hat{\mathbf{w}}_D, \mathbf{x}^i \rangle = 1$

SVM Revisited

$$\begin{aligned}\hat{\mathbf{w}} &= \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t. } 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle &\leq 0\end{aligned}$$

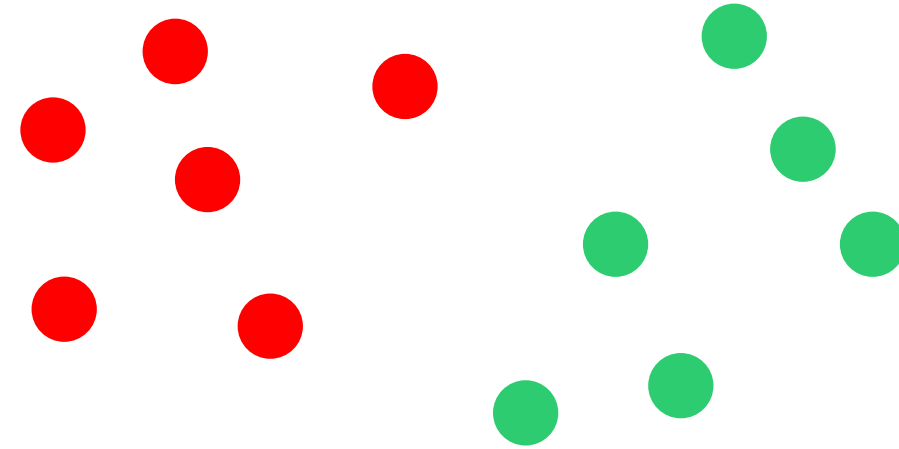
$$\hat{\mathbf{w}} = \hat{\mathbf{w}}_D = \sum_{i=1}^n \hat{\alpha}_i y^i \mathbf{x}^i$$

$$\begin{aligned}\hat{\alpha} &= \arg \max_{\alpha \in \mathbb{R}^n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle \\ \text{s.t. } \alpha_i &\geq 0\end{aligned}$$

$$\hat{\alpha}_i (1 - y^i \langle \hat{\mathbf{w}}_D, \mathbf{x}^i \rangle) = 0$$

Complimentary
slackness

If $\hat{\alpha}_i \neq 0$, then
 $y^i \langle \hat{\mathbf{w}}_D, \mathbf{x}^i \rangle = 1$



SVM Revisited

$$\begin{aligned}\hat{\mathbf{w}} &= \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t. } 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle &\leq 0\end{aligned}$$

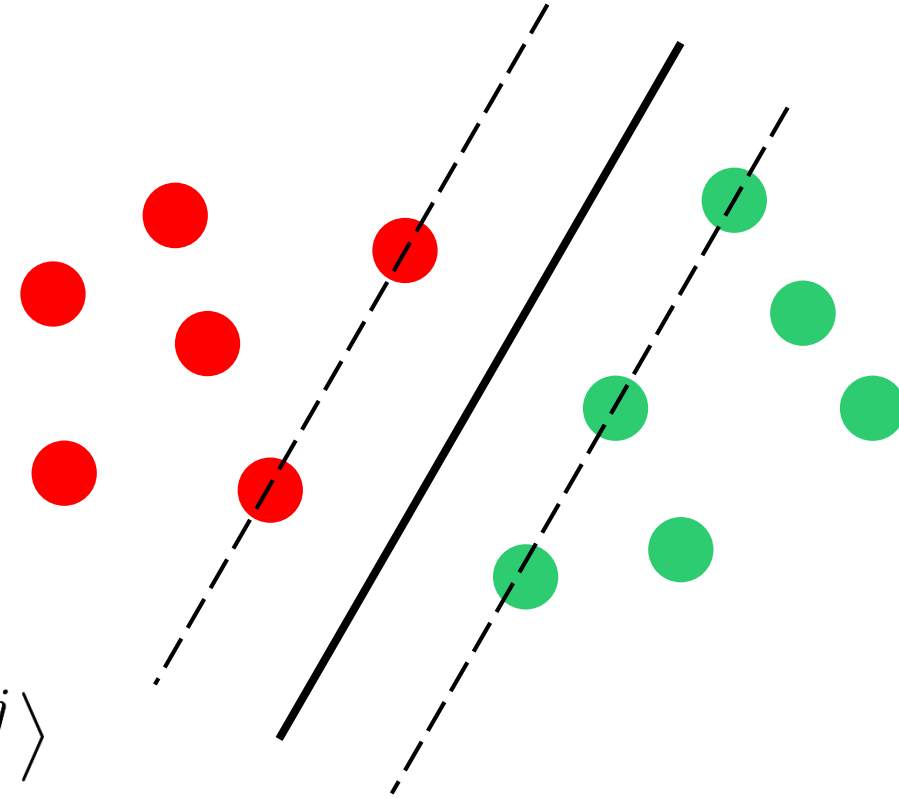
$$\hat{\mathbf{w}} = \hat{\mathbf{w}}_D = \sum_{i=1}^n \hat{\alpha}_i y^i \mathbf{x}^i$$

$$\begin{aligned}\hat{\alpha} &= \arg \max_{\alpha \in \mathbb{R}^n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle \\ \text{s.t. } \alpha_i &\geq 0\end{aligned}$$

$$\hat{\alpha}_i (1 - y^i \langle \hat{\mathbf{w}}_D, \mathbf{x}^i \rangle) = 0$$

Complimentary
slackness

If $\hat{\alpha}_i \neq 0$, then
 $y^i \langle \hat{\mathbf{w}}_D, \mathbf{x}^i \rangle = 1$



SVM Revisited

$$\begin{aligned}\hat{\mathbf{w}} &= \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t. } 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle &\leq 0\end{aligned}$$

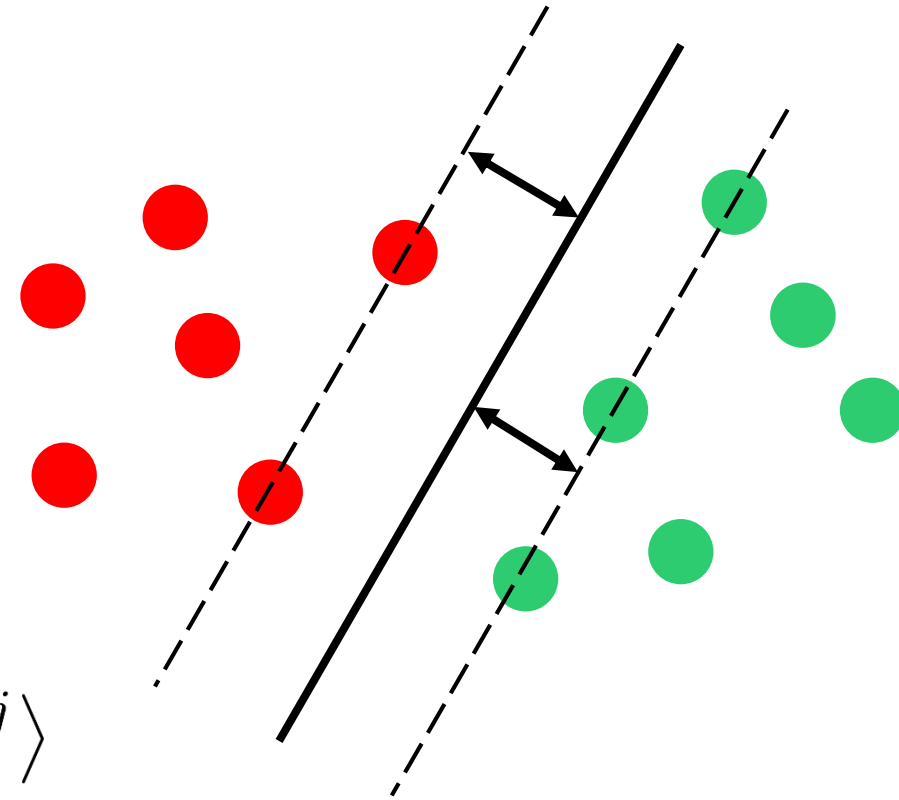
$$\hat{\mathbf{w}} = \hat{\mathbf{w}}_D = \sum_{i=1}^n \hat{\alpha}_i y^i \mathbf{x}^i$$

$$\begin{aligned}\hat{\alpha} &= \arg \max_{\alpha \in \mathbb{R}^n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle \\ \text{s.t. } \alpha_i &\geq 0\end{aligned}$$

$$\hat{\alpha}_i (1 - y^i \langle \hat{\mathbf{w}}_D, \mathbf{x}^i \rangle) = 0$$

Complimentary
slackness

If $\hat{\alpha}_i \neq 0$, then
 $y^i \langle \hat{\mathbf{w}}_D, \mathbf{x}^i \rangle = 1$



SVM Revisited

$$\begin{aligned}\hat{\mathbf{w}} &= \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t. } 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle &\leq 0\end{aligned}$$

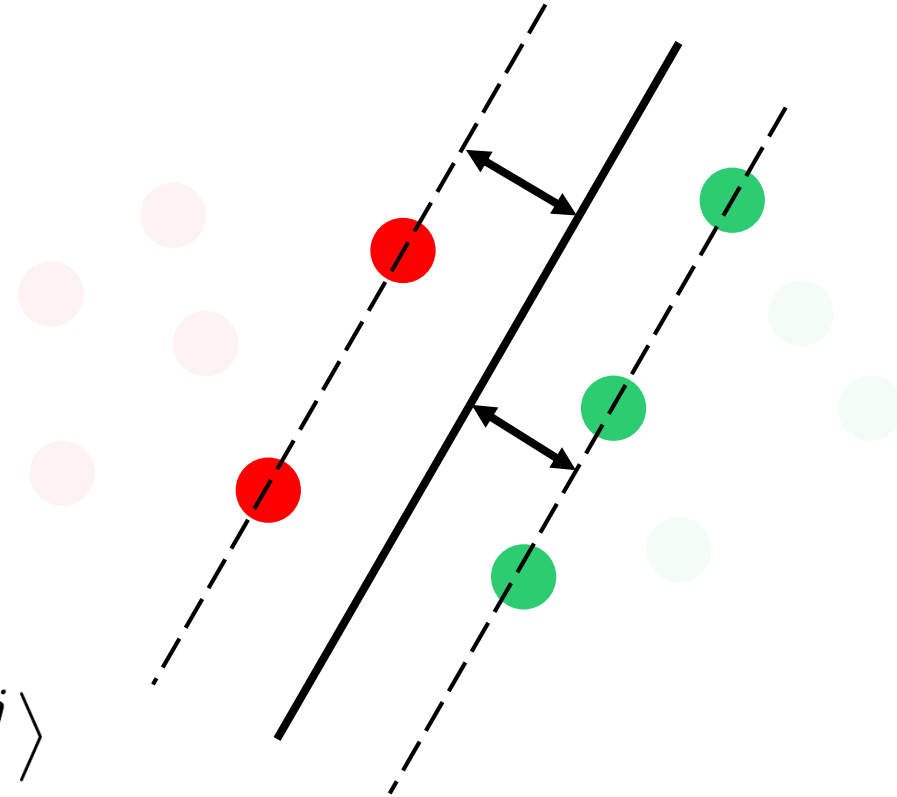
$$\hat{\mathbf{w}} = \hat{\mathbf{w}}_D = \sum_{i=1}^n \hat{\alpha}_i y^i \mathbf{x}^i$$

$$\begin{aligned}\hat{\alpha} &= \arg \max_{\alpha \in \mathbb{R}^n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle \\ \text{s.t. } \alpha_i &\geq 0\end{aligned}$$

$$\hat{\alpha}_i (1 - y^i \langle \hat{\mathbf{w}}_D, \mathbf{x}^i \rangle) = 0$$

Complimentary
slackness

If $\hat{\alpha}_i \neq 0$, then
 $y^i \langle \hat{\mathbf{w}}_D, \mathbf{x}^i \rangle = 1$



SVM Revisited

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\text{s.t. } 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle \leq 0$$

$$\hat{\mathbf{w}} = \hat{\mathbf{w}}_D = \sum_{i=1}^n \hat{\alpha}_i y^i \mathbf{x}^i$$

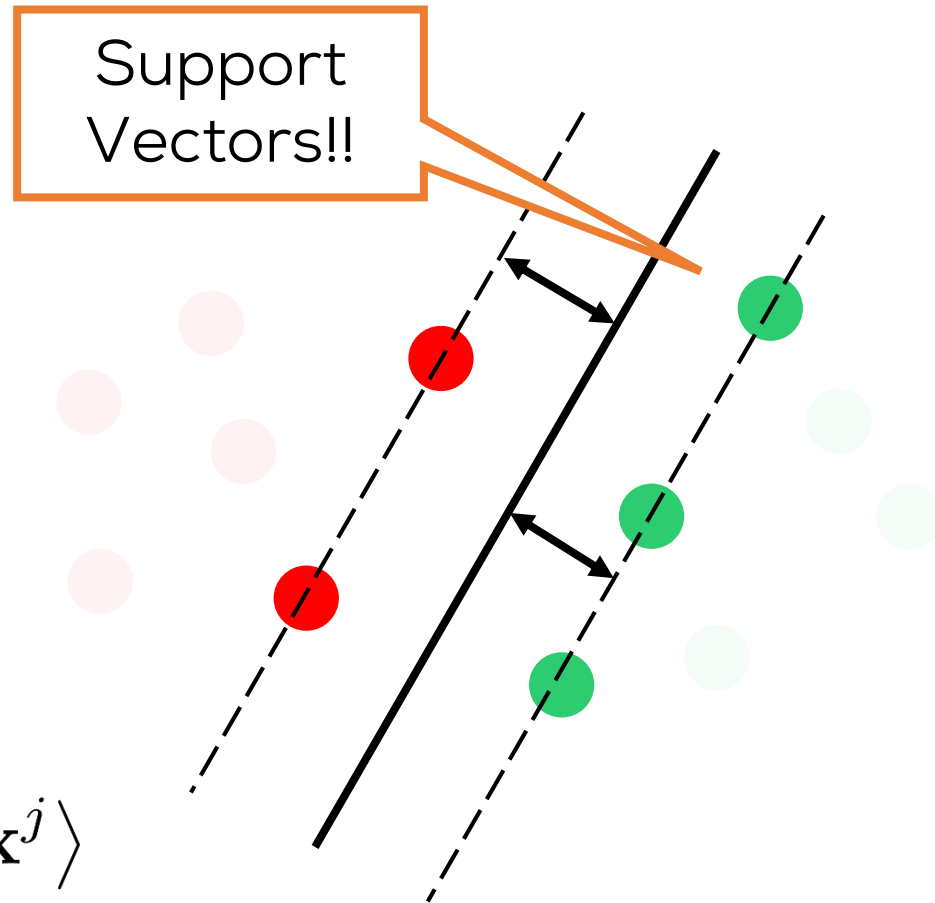
$$\hat{\alpha} = \arg \max_{\alpha \in \mathbb{R}^n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle$$

$$\text{s.t. } \alpha_i \geq 0$$

$$\hat{\alpha}_i (1 - y^i \langle \hat{\mathbf{w}}_D, \mathbf{x}^i \rangle) = 0$$

Complimentary
slackness

If $\hat{\alpha}_i \neq 0$, then
 $y^i \langle \hat{\mathbf{w}}_D, \mathbf{x}^i \rangle = 1$



SVM Revisited

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\text{s.t. } 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle \leq 0$$

$$\hat{\mathbf{w}} = \hat{\mathbf{w}}_D = \sum_{i=1}^n \hat{\alpha}_i y^i \mathbf{x}^i$$

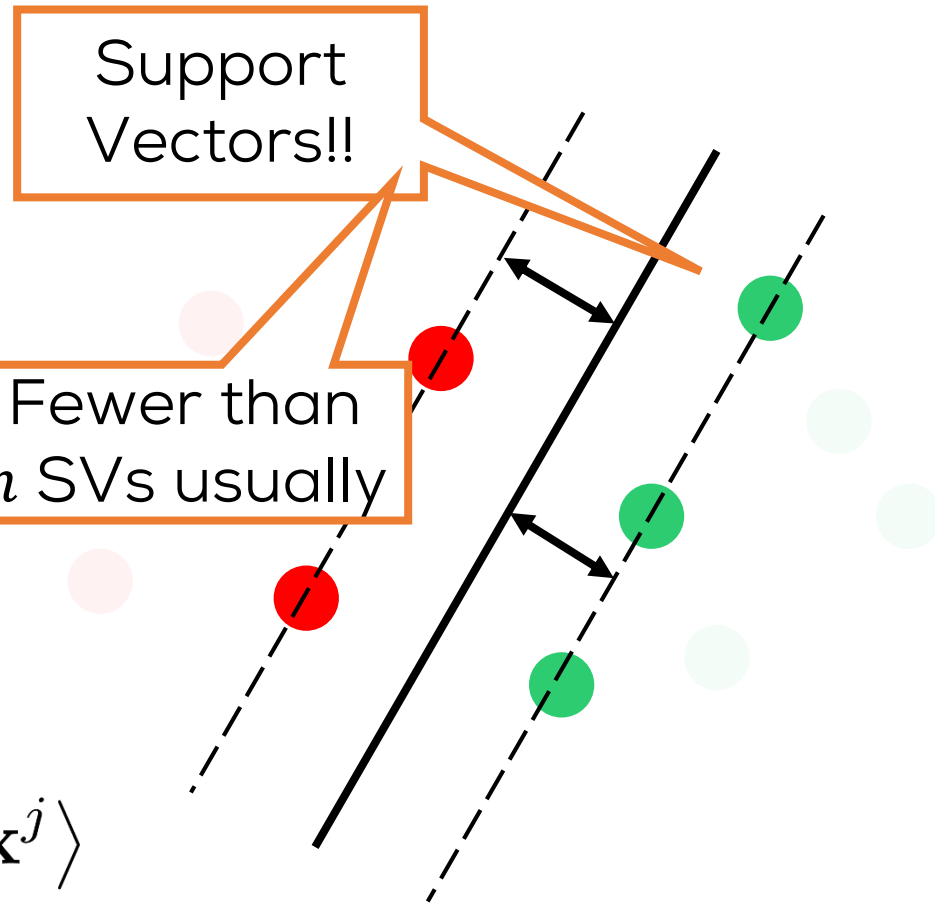
$$\hat{\alpha} = \arg \max_{\alpha \in \mathbb{R}^n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle$$

$$\text{s.t. } \alpha_i \geq 0$$

$$\hat{\alpha}_i (1 - y^i \langle \hat{\mathbf{w}}_D, \mathbf{x}^i \rangle) = 0$$

Complimentary slackness

If $\hat{\alpha}_i \neq 0$, then $y^i \langle \hat{\mathbf{w}}_D, \mathbf{x}^i \rangle = 1$



SVM Revisited

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\text{s.t. } 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle \leq 0$$

$$\hat{\mathbf{w}} = \hat{\mathbf{w}}_D = \sum_{i=1}^n \hat{\alpha}_i y^i \mathbf{x}^i$$

$$\hat{\alpha} = \arg \max_{\alpha \in \mathbb{R}^n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle$$

$$\text{s.t. } \alpha_i \geq 0$$

$$\hat{\alpha}_i (1 - y^i \langle \hat{\mathbf{w}}_D, \mathbf{x}^i \rangle) = 0$$

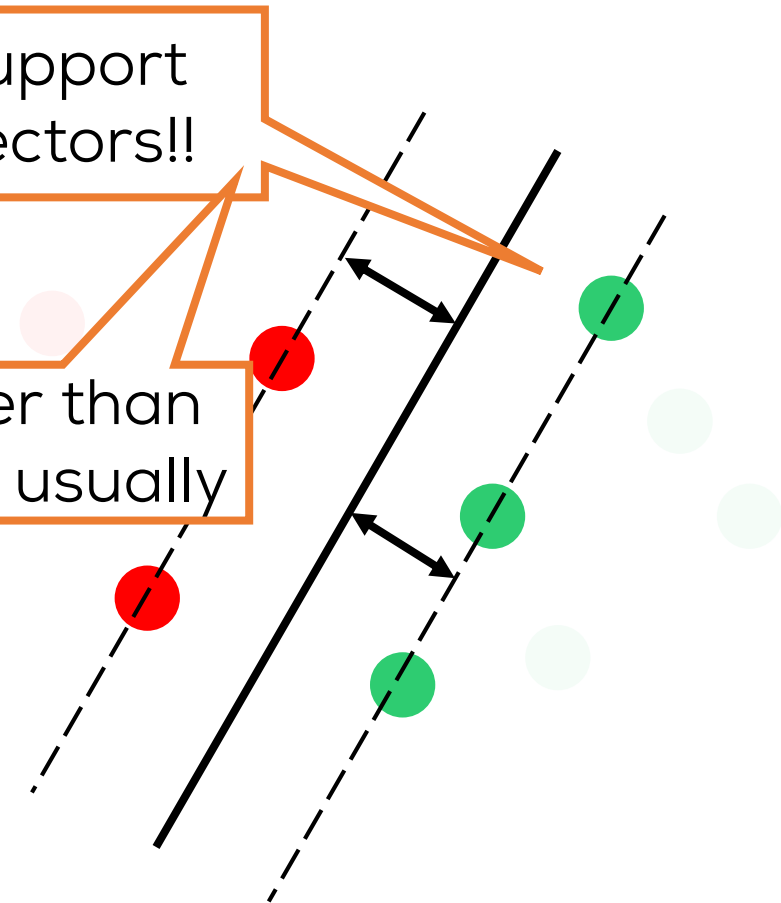
Support Vectors!!

Fewer than n SVs usually

$\|\alpha\|_1$ is a sparsity promoting reg.!!

Complimentary slackness

If $\hat{\alpha}_i \neq 0$, then $y^i \langle \hat{\mathbf{w}}_D, \mathbf{x}^i \rangle = 1$



SVM Revisited

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\text{s.t. } 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle \leq 0$$

$$\hat{\mathbf{w}} = \hat{\mathbf{w}}_D = \sum_{i=1}^n \hat{\alpha}_i y^i \mathbf{x}^i$$

$$\hat{\alpha} = \arg \max_{\alpha \in \mathbb{R}^n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle$$

$$\text{s.t. } \alpha_i \geq 0$$

$$\hat{\alpha}_i (1 - y^i \langle \hat{\mathbf{w}}_D, \mathbf{x}^i \rangle) = 0$$

What about the version with slack and all the notational agony?

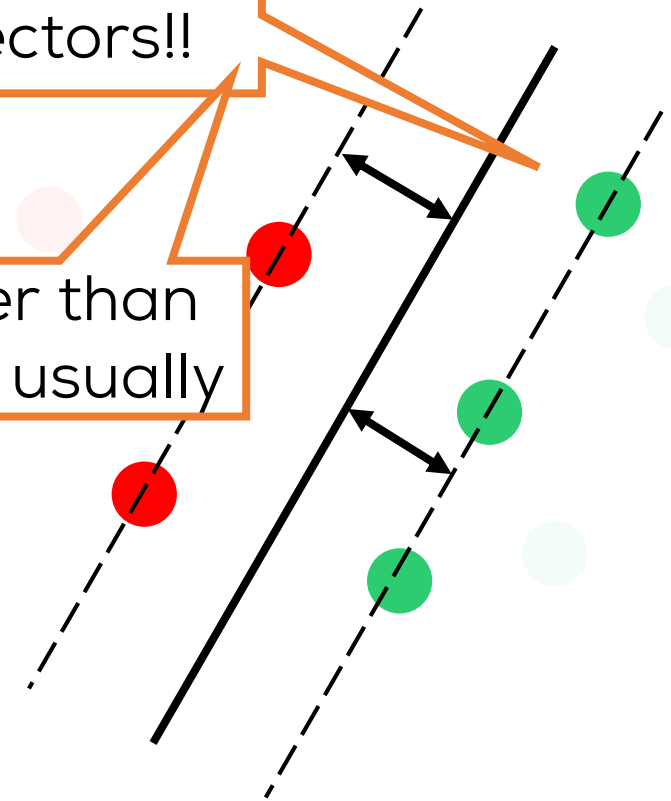
Support Vectors!!

Fewer than n SVs usually

$\|\alpha\|_1$ is a sparsity promoting reg.!!

Complimentary slackness

If $\hat{\alpha}_i \neq 0$, then $y^i \langle \hat{\mathbf{w}}_D, \mathbf{x}^i \rangle = 1$



SVM Revisited

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\text{s.t. } 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle \leq 0$$

$$\hat{\mathbf{w}} = \hat{\mathbf{w}}_D = \sum_{i=1}^n \hat{\alpha}_i y^i \mathbf{x}^i$$

$$\hat{\alpha} = \arg \max_{\alpha \in \mathbb{R}^n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle$$

$$\text{s.t. } \alpha_i \geq 0$$

$$\hat{\alpha}_i (1 - y^i \langle \hat{\mathbf{w}}_D, \mathbf{x}^i \rangle) = 0$$

What about the version with slack and all the notational agony?

Support Vectors!!

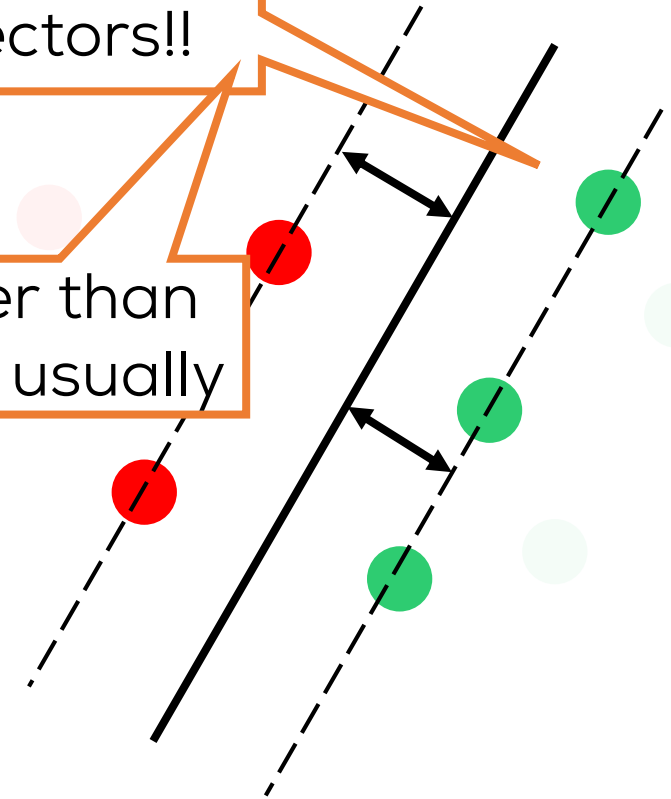
Slightly more tedious

Fewer than n SVs usually

$\|\alpha\|_1$ is a sparsity promoting reg.!!

Complimentary slackness

If $\hat{\alpha}_i \neq 0$, then $y^i \langle \hat{\mathbf{w}}_D, \mathbf{x}^i \rangle = 1$



SVM Revisited

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\text{s.t. } 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle \leq 0$$

$$\hat{\mathbf{w}} = \hat{\mathbf{w}}_D = \sum_{i=1}^n \hat{\alpha}_i y^i \mathbf{x}^i$$

$$\hat{\alpha} = \arg \max_{\alpha \in \mathbb{R}^n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle$$

$$\text{s.t. } \alpha_i \geq 0$$

$$\hat{\alpha}_i (1 - y^i \langle \hat{\mathbf{w}}_D, \mathbf{x}^i \rangle) = 0$$

What about the version with slack and all the notational agony?

Support Vectors!!

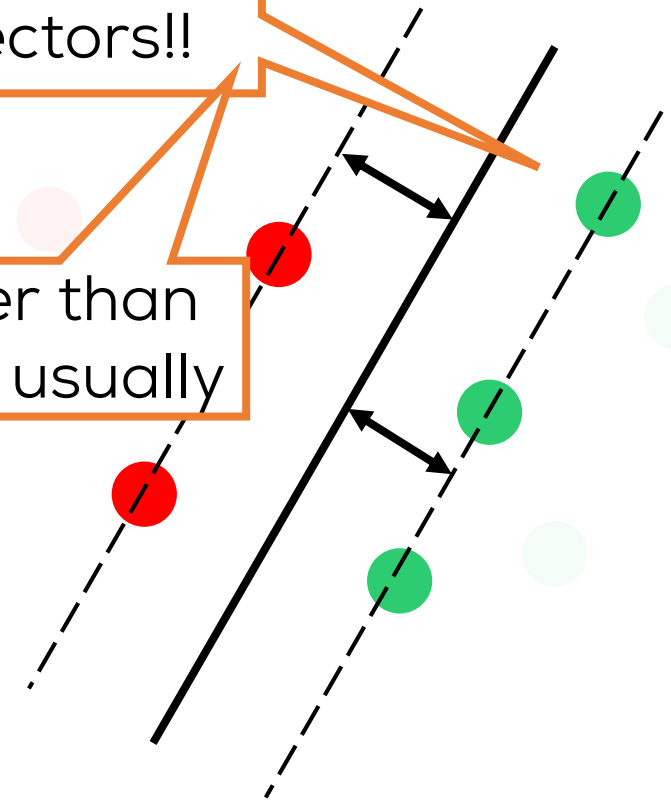
Slightly more tedious

Fewer than n SVs usually

$\|\alpha\|_1$ is a sparsity promoting reg.!!

Complimentary slackness

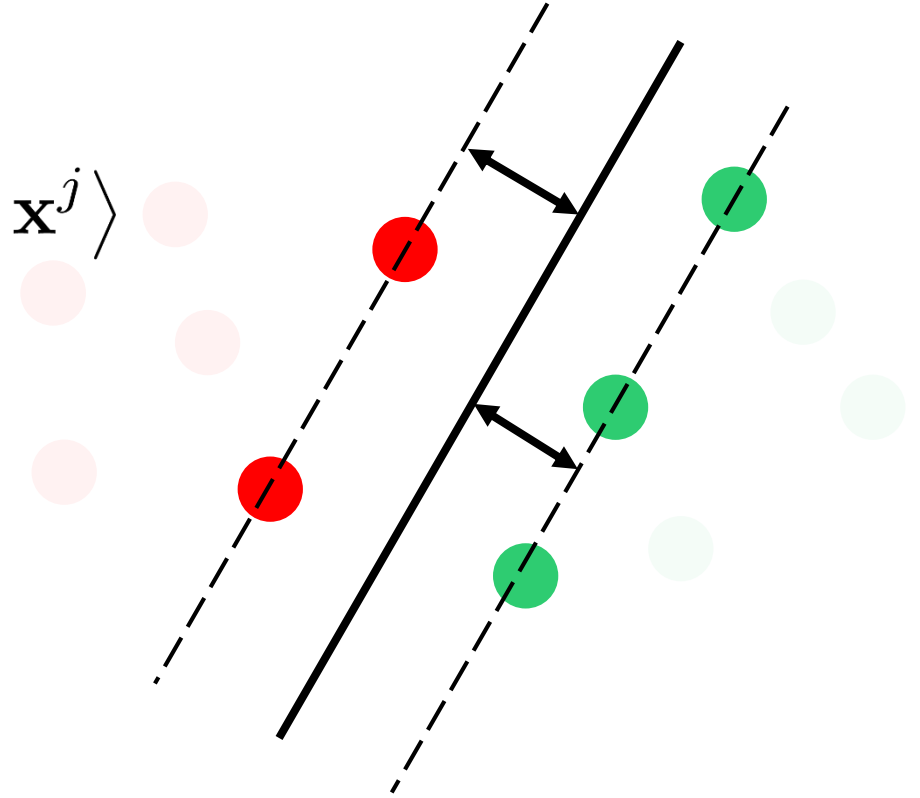
If $\hat{\alpha}_i \neq 0$, then $y^i \langle \hat{\mathbf{w}}_D, \mathbf{x}^i \rangle = 1$



App: SVMs via CD

$$\hat{\alpha} = \arg \max_{\alpha \in \mathbb{R}^n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle$$

s.t. $\alpha_i \geq 0$



Coordinate Descent

- Similar to gradient descent except update one variable at a time
- Notation: $\nabla_i f(\mathbf{w}) = [\nabla f(\mathbf{w})]_i$ i.e. the i -th directional derivative

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + r(\mathbf{w})$$

COORDINATE DESCENT

1. Initialize $\mathbf{w}^0 \in \mathbb{R}^d$
2. Select a coordinate $j_t \in [d]$
3. Let $\mathbf{g}_{j_t}^t \leftarrow \nabla_{j_t} f(\mathbf{w}^t) + \nabla_{j_t} r(\mathbf{w}^t) \in \mathbb{R}$
4. Update $\mathbf{w}_{j_t}^{t+1} \leftarrow \mathbf{w}_{j_t}^t - \eta_t \cdot \mathbf{g}_{j_t}^t$
5. Preserve other coord. $\mathbf{w}_k^{t+1} \leftarrow \mathbf{w}_k^t, k \neq j_t$
6. Repeat until convergence

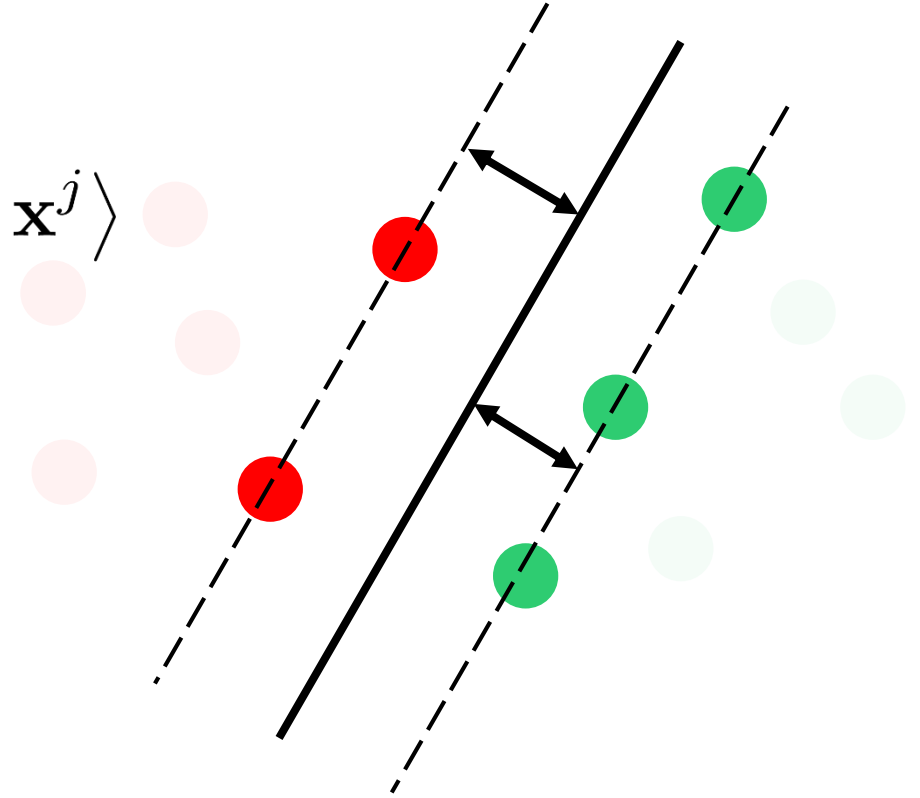
Random (Stochastic
CD), Cyclic
 $1, 2, \dots, d, 1, 2, 3 \dots d, \dots$

Only $O(n)$ time
per iter!

App: SVMs via CD

$$\hat{\alpha} = \arg \max_{\alpha \in \mathbb{R}^n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle$$

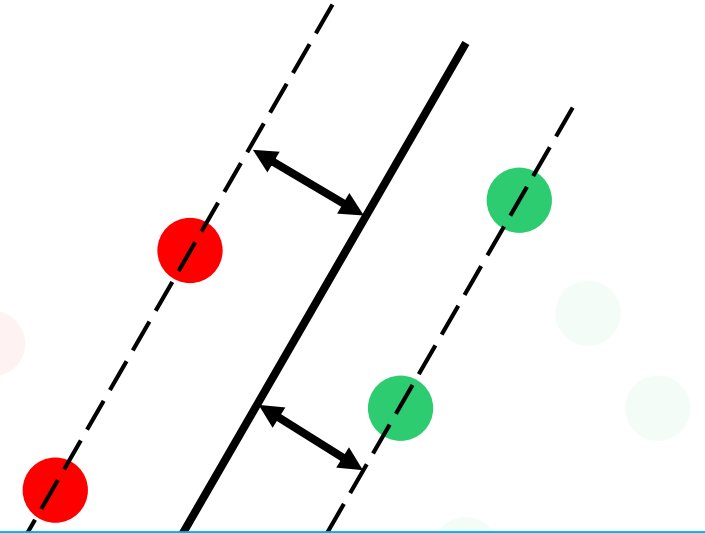
s.t. $\alpha_i \geq 0$



App: SVMs via CD

$$\hat{\alpha} = \arg \max_{\alpha \in \mathbb{R}^n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle$$

s.t. $\alpha_i \geq 0$



- Choose $i_t \in [n]$ at each time step (random, cyclic etc)
- Update only α_{i_t} , leave all other $\alpha_i, i \neq i_t$ untouched
 - Use projected gradient descent – constraint is pretty simple
- Note: α_{i_t} corresponds to the data point $(\mathbf{x}^{i_t}, y^{i_t})$
 - Only the i_t -th data point required to perform this update – exercise
 - Stochastic CD in dual looks like stochastic GD in primal!
- Current state of the art for SVMs – Liblinear, Scikit-learn

On to Data Modelling!

Hope you had fun with FA

Will revisit this method soon

For now, back to PML ...

Please give your Feedback

<http://tinyurl.com/ml17-18afb>