

Convex Optimization

Prof. Nati Srebro

Lecture 8: From the Fenchel Conjugate to Matrix Inequalities

Reading: Boyd and Vandenberghe Sections 3.3, 4.6, 5.1.6, 5.2.4, 5.7

Recommended reading on specific applications: 6.2, 8.6

$$\begin{array}{ll} \min_{\mathbf{x} \in \mathbb{R}^n} & f_0(\mathbf{x}) \\ s.t. & f_i(\mathbf{x}) \leq 0 \quad i = 1..m \\ & h_i(\mathbf{x}) = 0 \quad j = 1..p \end{array}$$

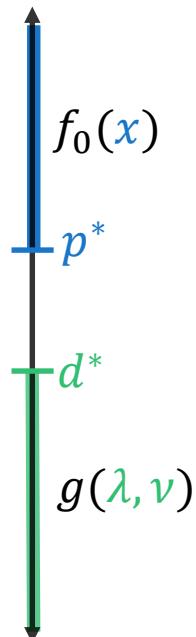
$$\begin{array}{ll} \max_{\boldsymbol{\lambda} \in \mathbb{R}^m, \boldsymbol{\nu} \in \mathbb{R}^p} & g(\boldsymbol{\lambda}, \boldsymbol{\nu}) \\ s.t. & \boldsymbol{\lambda}_i \geq 0 \end{array}$$

$$L(\mathbf{x}, (\boldsymbol{\lambda}, \boldsymbol{\nu})) = f_0(\mathbf{x}) + \sum_{i=1}^m \boldsymbol{\lambda}_i f_i(\mathbf{x}) + \sum_{j=1}^p \boldsymbol{\nu}_j h_j(\mathbf{x})$$

$$g(\boldsymbol{\lambda}, \boldsymbol{\nu}) = \inf_{\mathbf{x}} L(\mathbf{x}, (\boldsymbol{\lambda}, \boldsymbol{\nu}))$$

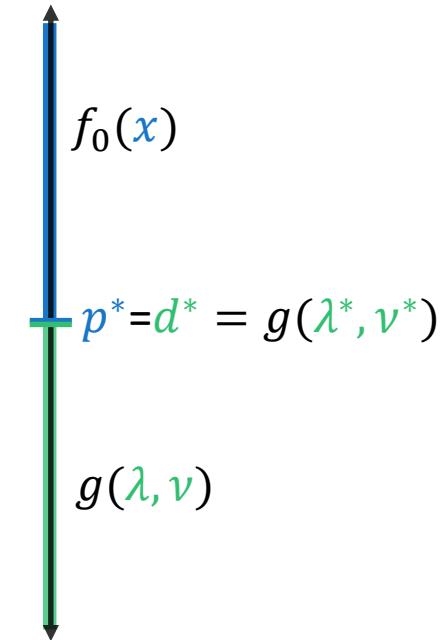
Weak Duality (always holds):

$$d^* \leq p^*$$



Convex + Slater \rightarrow Strong Duality

$$d^* = p^*$$



Dual of a Function

$$\begin{array}{ll} \min_{\mathbf{x} \in \mathbb{R}^n} & f(\mathbf{x}) \\ \text{s.t.} & \mathbf{x} = 0 \end{array}$$

$$\max_{\mathbf{v} \in \mathbb{R}^n} -f^*(-\mathbf{v})$$

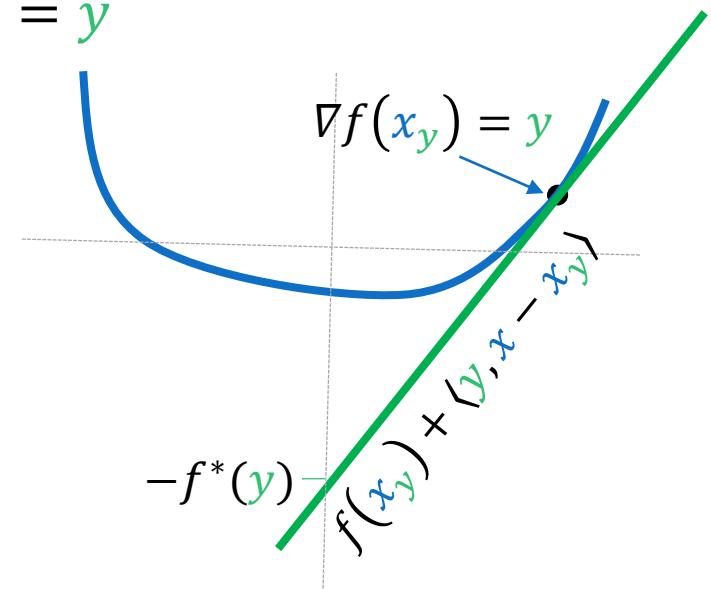
$$g(\mathbf{v}) = \inf_{\mathbf{x}} L(\mathbf{x}, \mathbf{v}) = \inf_{\mathbf{x}} f(\mathbf{x}) + \langle \mathbf{v}, \mathbf{x} \rangle$$

- Def: The Fenchel Conjugate of $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is a function $f^*: (\mathbb{R}^n)^* \rightarrow \mathbb{R}$

$$f^*(\mathbf{y}) = \sup_{\mathbf{x}} \langle \mathbf{y}, \mathbf{x} \rangle - f(\mathbf{x})$$

- Claim: $f^*(\mathbf{y}) = \langle \mathbf{y}, \mathbf{x}_y \rangle - f(\mathbf{x}_y)$ where $\nabla f(\mathbf{x}_y) = \mathbf{y}$

Proof: $0 = \nabla_x (\langle \mathbf{y}, \mathbf{x} \rangle - f(\mathbf{x})) = \mathbf{y} - \nabla f(\mathbf{x})$

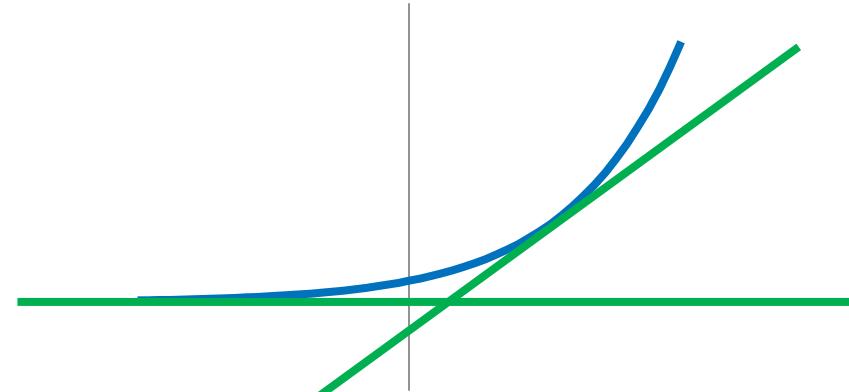


$f^*(\textcolor{teal}{y}) \stackrel{\text{def}}{=} \sup_x \langle \textcolor{teal}{y}, \textcolor{blue}{x} \rangle - f(\textcolor{blue}{x}) = \langle \textcolor{teal}{y}, \textcolor{blue}{x}_y \rangle - f(\textcolor{blue}{x}_y)$ where $\nabla f(\textcolor{blue}{x}_y) = \textcolor{teal}{y}$

- $f(x) = e^x$

- $y = f'(x) = e^x \rightarrow x = \log y$

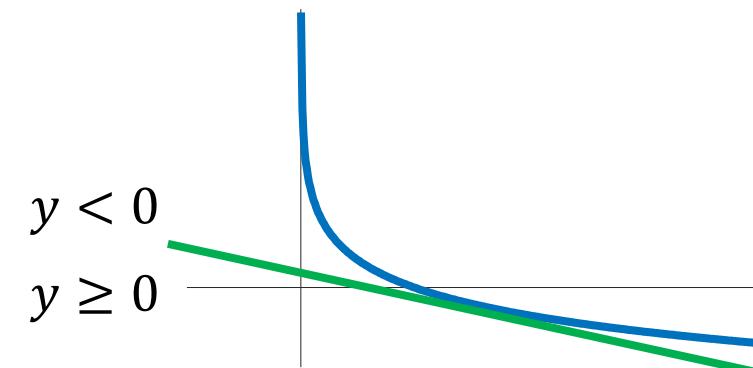
- $f^*(y) = \begin{cases} y \log y - y & y > 0 \\ 0 & y = 0 \\ \infty & y < 0 \end{cases}$



- $f(x) = -\log x$

- $y = f'(x) = -\frac{1}{x} \rightarrow x^* = -\frac{1}{y}$ for $y < 0$

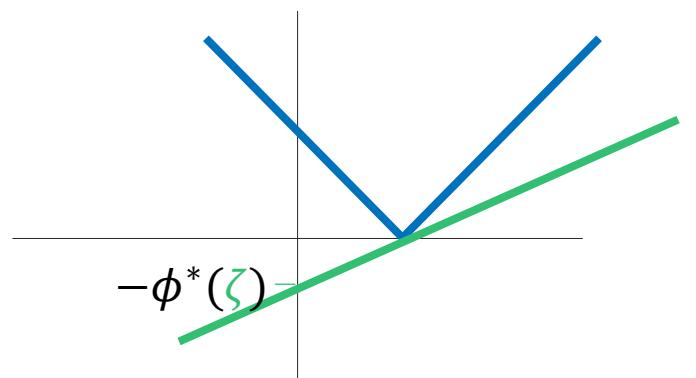
- $f^*(y) = \begin{cases} y \left(\frac{-1}{y} \right) + \log \frac{-1}{y} = -1 - \log(-y) & y < 0 \\ \infty & y \geq 0 \end{cases}$



- $f(x) = |x - 1|$

- $y = f'(x) \rightarrow x = 1$ for $-1 \leq y \leq 1$

- $f^*(y) = \begin{cases} y \cdot 1 - |1 - 1| = y & -1 \leq y \leq 1 \\ \infty & \text{otherwise} \end{cases}$



Optimization with Linear Constraints

$$\begin{array}{ll} \min_{\mathbf{x} \in \mathbb{R}^n} & f(\mathbf{x}) \\ \text{s.t.} & A\mathbf{x} = b \\ & G\mathbf{x} \leq h \end{array}$$

$$\begin{array}{ll} \max & -\langle \lambda, h \rangle - \langle v, b \rangle - f^*(-G^\top \lambda - A^\top v) \\ \text{s.t.} & \lambda \geq 0 \end{array}$$

$$\begin{aligned} g(\lambda, v) &= \inf_x f(x) + \langle \lambda, Gx - h \rangle + \langle v, Ax - b \rangle \\ &= -\langle \lambda, h \rangle - \langle v, b \rangle - \sup_x (\langle -G^\top \lambda - A^\top v, x \rangle - f(x)) \\ &= -\langle \lambda, h \rangle - \langle v, b \rangle - f^*(-G^\top \lambda - A^\top v) \end{aligned}$$

Example: Linear Programming

- $f(\mathbf{x}) = \langle c, \mathbf{x} \rangle$
- $f^*(y) = \begin{cases} 0 & y = c \\ \infty & \text{otherwise} \end{cases}$

$$\begin{array}{ll} \max & -\langle \lambda, h \rangle - \langle v, b \rangle \\ \text{s.t.} & -G^\top \lambda - A^\top v = c \\ & \lambda \geq 0 \end{array}$$

Example: Min Norm Solution

$$\begin{array}{ll} \min_{\mathbf{x} \in \mathbb{R}^n} & \|\mathbf{x}\|_2 \\ \text{s.t.} & A\mathbf{x} = b \end{array}$$

$$\begin{array}{ll} \max_{\mathbf{v}} & \langle -b, \mathbf{v} \rangle \\ \text{s.t.} & \|A^\top \mathbf{v}\|_2 \leq 1 \end{array}$$

- $f(x) = \|\mathbf{x}\|_2$
- $f^*(y) = \sup_x (\langle y, x \rangle - \|\mathbf{x}\|_2) = \begin{cases} 0 & \|\mathbf{y}\|_2 \leq 1 \\ \infty & \|\mathbf{y}\|_2 > 1 \end{cases}$

If $\|\mathbf{y}\|_2 \leq 1$, then $\langle y, x \rangle \leq \|\mathbf{x}\|_2$
If $\|\mathbf{y}\|_2 > 1$, take $x = t\mathbf{y}, t \rightarrow \infty$

$$\begin{array}{ll} \min_{\mathbf{x} \in \mathbb{R}^n} & f(\mathbf{x}) \\ \text{s.t.} & A\mathbf{x} = b \\ & G\mathbf{x} \leq h \end{array}$$

$$\begin{array}{ll} \max & -\langle \lambda, h \rangle - \langle \nu, b \rangle - f^*(-G^\top \lambda - A^\top \nu) \\ \text{s.t.} & \lambda \geq 0 \end{array}$$

Example: Min Norm Solution

$$\begin{array}{ll} \min_{\mathbf{x} \in \mathbb{R}^n} & \|\mathbf{x}\|_2 \\ \text{s.t.} & A\mathbf{x} = b \end{array}$$

$$\begin{array}{ll} \max_{\mathbf{v}} & \langle -b, \mathbf{v} \rangle \\ \text{s.t.} & \|A^\top \mathbf{v}\|_2 \leq 1 \end{array}$$

- $f(x) = \|\mathbf{x}\|_2$
- $f^*(y) = \sup_x (\langle y, x \rangle - \|\mathbf{x}\|_2) = \begin{cases} 0 & \|\mathbf{y}\|_2 \leq 1 \\ \infty & \|\mathbf{y}\|_2 > 1 \end{cases}$

If $\|\mathbf{y}\|_2 \leq 1$, then $\langle y, x \rangle \leq \|\mathbf{x}\|_2$
 If $\|\mathbf{y}\|_2 > 1$, take $x = t\mathbf{y}, t \rightarrow \infty$
- For a general norm, $f(x) = \|\mathbf{x}\|$: $f^*(x) = \begin{cases} 0 & \|\mathbf{y}\|_* \leq 1 \\ \infty & \|\mathbf{y}\|_* > 1 \end{cases}$
- Dual norm: $\|\mathbf{y}\|_* = \sup_{\|\mathbf{x}\| \leq 1} \langle y, x \rangle$

$$\begin{array}{ll} \min_{\mathbf{x} \in \mathbb{R}^n} & \|\mathbf{x}\| \\ \text{s.t.} & A\mathbf{x} = b \end{array}$$

$$\begin{array}{ll} \max_{\mathbf{v}} & \langle -b, \mathbf{v} \rangle \\ \text{s.t.} & \|A^\top \mathbf{v}\|_* \leq 1 \end{array}$$

Combining Fenchel Conjugates

$$f(x_1, x_2, \dots, x_n) = \sum_{i=1}^n f_i(x_i)$$

$$\begin{aligned} f^*(y) &= \sup_x (\langle y, x \rangle - f(x)) = \sup_x \sum_{i=1}^n (y_i x_i - f_i(x_i)) \\ &= \sum_{i=1}^n \sup_{x_i} (y_i x_i - f_i(x_i)) = \sum_{i=1}^n f_i^*(y_i) \end{aligned}$$

Example: Linear Classification

$$\min_{\mathbf{w} \in \mathbb{R}^n} \quad \sum_{i=1}^m \phi(y_i \langle \mathbf{w}, \mathbf{x}_i \rangle)$$

$$\max g$$

$$g = \inf_{\mathbf{w}} \sum_{i=1}^m \phi(y_i \langle \mathbf{w}, \mathbf{x}_i \rangle)$$

Example: Linear Classification

$$\min_{\mathbf{w} \in \mathbb{R}^n} \sum_{i=1}^m \phi(y_i \langle \mathbf{w}, \mathbf{x}_i \rangle)$$

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^n, \mathbf{z} \in \mathbb{R}^m} \quad & \sum_{i=1}^m \phi(\mathbf{z}_i) \\ \text{s.t.} \quad & \mathbf{z}_i = y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \end{aligned}$$

$$\begin{aligned} \max_{\mathbf{v} \in \mathbb{R}^m} \quad & -\sum_{i=1}^m \phi^*(\mathbf{v}_i) \\ \text{s.t.} \quad & \sum_{i=1}^m y_i \mathbf{v}_i \mathbf{x}_i = 0 \end{aligned}$$

- $f(\mathbf{z}, \mathbf{w}) = \sum_{i=1}^m \phi(\mathbf{z}_i) + \langle \mathbf{0}, \mathbf{w} \rangle$
- $f^*(\boldsymbol{\zeta}, \boldsymbol{\omega}) = \sum_{i=1}^m \phi^*(\boldsymbol{\zeta}_i) + \begin{cases} 0 & \boldsymbol{\omega} = 0 \\ \infty & \boldsymbol{\omega} \neq 0 \end{cases}$
- $\tilde{x}_i = y_i x_i$, i.e. $\tilde{X} = X \cdot \text{diag}(y) \in \mathbb{R}^{n \times m}$
- Constraint: $[-I, \tilde{X}^\top] [\mathbf{z}; \mathbf{w}] = 0$
- Dual objective: $-\langle \mathbf{v}, \mathbf{0} \rangle - f^*([I\mathbf{v}, -\tilde{X}\mathbf{v}]) = -\sum_{i=1}^m \phi^*(\mathbf{v}_i)$ s.t. $\sum_{i=1}^m y_i \mathbf{v}_i \mathbf{x}_i = 0$

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^n} \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & A\mathbf{x} = b, G\mathbf{x} \leq h \end{aligned}$$

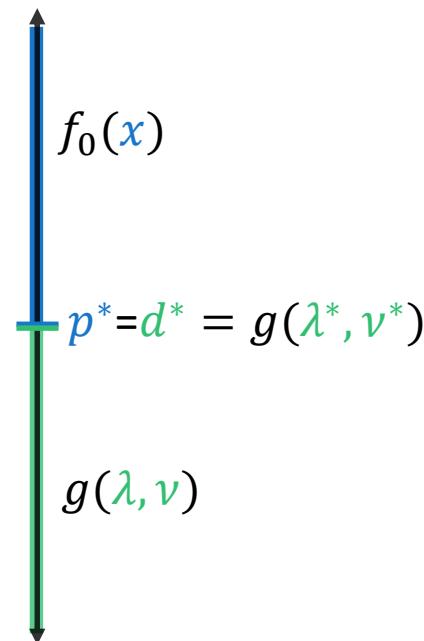
$$\begin{aligned} \max \quad & -\langle \boldsymbol{\lambda}, h \rangle - \langle \mathbf{v}, b \rangle - f^*(-G^\top \boldsymbol{\lambda} - A^\top \mathbf{v}) \\ \text{s.t.} \quad & \boldsymbol{\lambda} \geq 0 \end{aligned}$$

Example: Linear Classification

$$\min_{\mathbf{w} \in \mathbb{R}^n} \quad \sum_{i=1}^m \phi(y_i \langle \mathbf{w}, \mathbf{x}_i \rangle)$$

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^n, \mathbf{z} \in \mathbb{R}^m} \quad & \sum_{i=1}^m \phi(z_i) \\ \text{s.t.} \quad & z_i = y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \end{aligned}$$

$$\begin{aligned} \max_{\mathbf{v} \in \mathbb{R}^m} \quad & - \sum_{i=1}^m \phi^*(v_i) \\ \text{s.t.} \quad & \sum_{i=1}^m y_i v_i x_i = 0 \end{aligned}$$



Example: Linear Classification

$$\min_{\mathbf{w} \in \mathbb{R}^n} \sum_{i=1}^m \phi(y_i \langle \mathbf{w}, \mathbf{x}_i \rangle)$$

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^n, \mathbf{z} \in \mathbb{R}^m} & \sum_{i=1}^m \phi(z_i) \\ \text{s.t.} & z_i = y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \end{aligned}$$

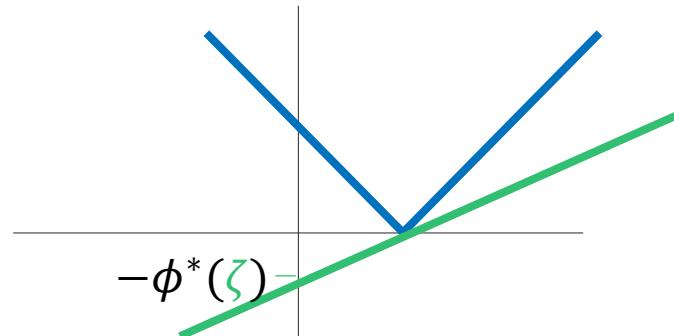
$$\begin{aligned} \max_{\mathbf{v} \in \mathbb{R}^m} & -\sum_{i=1}^m \phi^*(v_i) \\ \text{s.t.} & \sum_{i=1}^m y_i v_i x_i = 0 \end{aligned}$$

Absolute Error Loss:

$$\phi(z) = |z - 1| = |y_i \langle \mathbf{w}, \mathbf{x}_i \rangle - 1| = |\langle \mathbf{w}, \mathbf{x}_i \rangle - y_i|$$

$$\phi^*(\zeta) = \begin{cases} \zeta & -1 \leq \zeta \leq 1 \\ \infty & \text{otherwise} \end{cases}$$

$$\begin{aligned} \max_{\mathbf{v} \in \mathbb{R}^m} & -\sum_{i=1}^m v_i \\ \text{s.t.} & -1 \leq v_i \leq 1 \\ & \sum_{i=1}^m y_i v_i x_i = 0 \end{aligned}$$



Example: Linear Classification

$$\min_{\mathbf{w} \in \mathbb{R}^n} \quad \sum_{i=1}^m \phi(y_i \langle \mathbf{w}, \mathbf{x}_i \rangle)$$

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^n, \mathbf{z} \in \mathbb{R}^m} \quad & \sum_{i=1}^m \phi(z_i) \\ \text{s.t.} \quad & z_i = y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \end{aligned}$$

Absolute Error Loss:

$$\phi(z) = |z - 1| = |y_i \langle \mathbf{w}, \mathbf{x}_i \rangle - 1| = |\langle \mathbf{w}, \mathbf{x}_i \rangle - y_i|$$

$$\phi^*(\zeta) = \begin{cases} \zeta & -1 \leq \zeta \leq 1 \\ \infty & \text{otherwise} \end{cases}$$

Squared Loss:

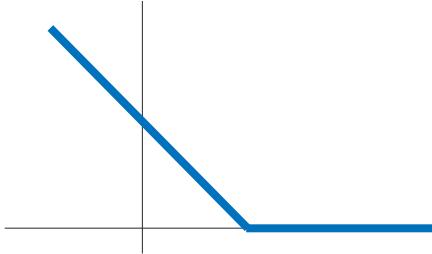
$$\phi(z - 1) = \frac{1}{2}(z - 1)^2 = (\langle \mathbf{w}, \mathbf{x}_i \rangle - y_i)^2$$

$$\phi^*(\zeta) = \frac{1}{2}\zeta^2 + \zeta$$

$$\begin{aligned} \max_{\mathbf{v} \in \mathbb{R}^m} \quad & -\sum_{i=1}^m \phi^*(v_i) \\ \text{s.t.} \quad & \sum_{i=1}^m y_i v_i x_i = 0 \end{aligned}$$

$$\begin{aligned} \max_{\mathbf{v} \in \mathbb{R}^m} \quad & -\sum_{i=1}^m v_i \\ \text{s.t.} \quad & -1 \leq v_i \leq 1 \\ & \sum_{i=1}^m y_i v_i x_i = 0 \end{aligned}$$

$$\begin{aligned} \max_{\mathbf{v} \in \mathbb{R}^m} \quad & -\frac{1}{2} \sum_{i=1}^m (v_i + v_i^2) \\ \text{s.t.} \quad & \sum_{i=1}^m y_i v_i x_i = 0 \end{aligned}$$



Hinge Loss:

$$\phi(z) = [1 - z]_+$$

$$\phi^*(y) = \begin{cases} y & -1 \leq y \leq 0 \\ \infty & \text{otherwise} \end{cases}$$

Absolute Error Loss:

$$\phi(z) = |z - 1| = |y_i \langle \mathbf{w}, \mathbf{x}_i \rangle - 1| = |\langle \mathbf{w}, \mathbf{x}_i \rangle - y_i|$$

$$\phi^*(\zeta) = \begin{cases} \zeta & -1 \leq \zeta \leq 1 \\ \infty & \text{otherwise} \end{cases}$$

Squared Loss:

$$\phi(z) = \frac{1}{2}z^2 = (\langle \mathbf{w}, \mathbf{x}_i \rangle - y_i)^2$$

$$\phi^*(\zeta) = \frac{1}{2}\zeta^2$$

$$\max_{\mathbf{v} \in \mathbb{R}^m} -\sum_{i=1}^m v_i$$

$$\text{s.t.} \quad -1 \leq v_i \leq 0$$

$$\sum_{i=1}^m y_i v_i x_i = 0$$

$$\max_{\mathbf{v} \in \mathbb{R}^m} -\sum_{i=1}^m v_i$$

$$\text{s.t.} \quad -1 \leq v_i \leq 1$$

$$\sum_{i=1}^m y_i v_i x_i = 0$$

$$\max_{\mathbf{v} \in \mathbb{R}^m} -\frac{1}{2} \sum_{i=1}^m v_i^2$$

$$\text{s.t.} \quad \sum_{i=1}^m y_i v_i x_i = 0$$

Logistic Regression:

$$\phi(z) = \log(1 + e^{-z})$$

$$h(p) = -p \log p - (1-p) \log(1-p)$$

$$\phi^*(\zeta) = \begin{cases} -h(-\zeta) & -1 \leq \zeta \leq 0 \\ \infty & \text{otherwise} \end{cases}$$

Hinge Loss:

$$\phi(z) = [1 - z]_+$$

$$\phi^*(\zeta) = \begin{cases} \zeta & -1 \leq \zeta \leq 0 \\ \infty & \text{otherwise} \end{cases}$$

Absolute Error Loss:

$$\phi(z) = |z - 1| = |y_i \langle w, x_i \rangle - 1| = |\langle w, x_i \rangle - y_i|$$

$$\phi^*(\zeta) = \begin{cases} \zeta & -1 \leq \zeta \leq 1 \\ \infty & \text{otherwise} \end{cases}$$

Squared Loss:

$$\phi(z) = \frac{1}{2}z^2 = (\langle w, x_i \rangle - y_i)^2$$

$$\phi^*(\zeta) = \frac{1}{2}\zeta^2$$

$$\max_{v \in \mathbb{R}^m} \sum_{i=1}^m h(-v_i)$$

$$\text{s.t.} \quad -1 \leq v_i \leq 0$$

$$\sum_{i=1}^m y_i v_i x_i = 0$$

$$\max_{v \in \mathbb{R}^m} -\sum_{i=1}^m v_i$$

$$\text{s.t.} \quad -1 \leq v_i \leq 0$$

$$\sum_{i=1}^m y_i v_i x_i = 0$$

$$\max_{v \in \mathbb{R}^m} -\sum_{i=1}^m v_i$$

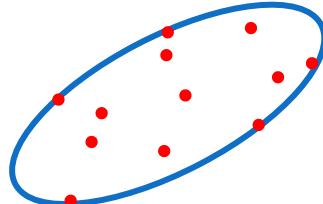
$$\text{s.t.} \quad -1 \leq v_i \leq 1$$

$$\sum_{i=1}^m y_i v_i x_i = 0$$

$$\max_{v \in \mathbb{R}^m} -\frac{1}{2} \sum_{i=1}^m v_i^2$$

$$\text{s.t.} \quad \sum_{i=1}^m y_i v_i x_i = 0$$

Minimum Volume Covering Ellipsoid



- Problem: Given points $\{a_i \in \mathbb{R}^n\}$ find minimum volume ellipsoid centered at origin containing all points

$$\mathcal{E} = \{z | \|Pz\|^2 \leq 1\} = \{z | z^\top X z \leq 1\}$$

$$X = P^\top P \geq 0$$

$$\text{Volume}(\mathcal{E}) \propto \sqrt{\det X^{-1}}$$

$S^n \subset \mathbb{R}^{n \times n}$
Symmetric

$$\begin{aligned} \min_{X \in S^n} \quad & -\log \det X \\ \text{s.t.} \quad & a_i^\top X a_i \leq 1 \quad i = 1..m \end{aligned}$$

$$\begin{aligned} a_i^\top X a_i &= \sum_{jk} (a_i[j] a_i[k]) X[j, k] \\ &= \langle a_i a_i^\top, X \rangle \end{aligned}$$

$$\log \det X = \log \prod_i \lambda_i(X) \stackrel{\text{def}}{=} \sum_i \log \lambda_i(X) = \begin{cases} \log (\det X) & \text{if } X \geq 0 \\ \infty & \text{otherwise} \end{cases}$$

Dual of MVCE

$$\begin{array}{ll} \min_{\mathbf{X} \in S^n} & -\log \det \mathbf{X} \\ \text{s.t.} & a_i^\top \mathbf{X} a_i \leq 1 \quad i = 1..m \\ & \mathbf{X} \succ 0 \end{array}$$

$$\begin{array}{ll} \max_{\lambda \in \mathbb{R}^m} & -\sum \lambda_i + n + \log \det(\sum_i \lambda_i a_i a_i^\top) \\ \text{s.t.} & \lambda_i \geq 0 \\ & \sum_i \lambda_i a_i a_i^\top \succ 0 \end{array}$$

- Fenchel conjugate of $f(X) = -\log \det X$:

$$Y = \nabla f(X) = -X^{-1} \rightarrow X = -Y^{-1} \text{ if } Y \prec 0$$

$$\langle A, B \rangle = \sum_{jk} A[j, k] B[j, k] = \text{tr}(A^\top B)$$

$$f^*(Y) = \langle Y, -Y^{-1} \rangle + \log \det(-Y^{-1}) = -\text{tr}(I) - \log \det(-Y)$$

And so:

$$f^*(Y) = \begin{cases} -n - \log \det(-Y) & \text{if } Y \prec 0 \\ \infty & \text{otherwise} \end{cases}$$

- Strong Duality?**

- Yes! Large enough ellipsoid always feasible, and $\log \det X$ finite.

$$\begin{array}{ll} \min_{\mathbf{x} \in \mathbb{R}^n} & f(\mathbf{x}) \\ \text{s.t.} & A\mathbf{x} = b, G\mathbf{x} \leq h \end{array}$$

$$\begin{array}{ll} \max & -\langle \boldsymbol{\lambda}, h \rangle - \langle \mathbf{v}, b \rangle - f^*(-G^\top \boldsymbol{\lambda} - A^\top \mathbf{v}) \\ \text{s.t.} & \boldsymbol{\lambda} \geq 0 \end{array}$$

Matrix Inequalities (Semi Definite Constraints)

$$\begin{array}{ll} \min_{\mathbf{x} \in \mathbb{R}^n} & f_0(\mathbf{x}) \\ s.t. & f_i(\mathbf{x}) \leq 0 \\ & h_i(\mathbf{x}) = 0 \end{array}$$

$f_0: \mathbb{R}^n \rightarrow \mathbb{R}$ convex

$f_i: \mathbb{R}^n \rightarrow S^{k_i} \subset \mathbb{R}^{k_i \times k_i}$ **matrix-convex:**

$$f_i(\theta x + (1 - \theta)x') \leq \theta f_i(x) + (1 - \theta)f_i(x')$$

$h_i: \mathbb{R}^n \rightarrow \mathbb{R}$ linear (or $\rightarrow S^k$, doesn't matter much)

Sometimes convenient to represent $x \in S^n$: doesn't matter much,
its just a vector space of dimensionality $n(n + 1)/2$

Semi-Definite Programming (SDP): f_0, f_i linear