# Optimization
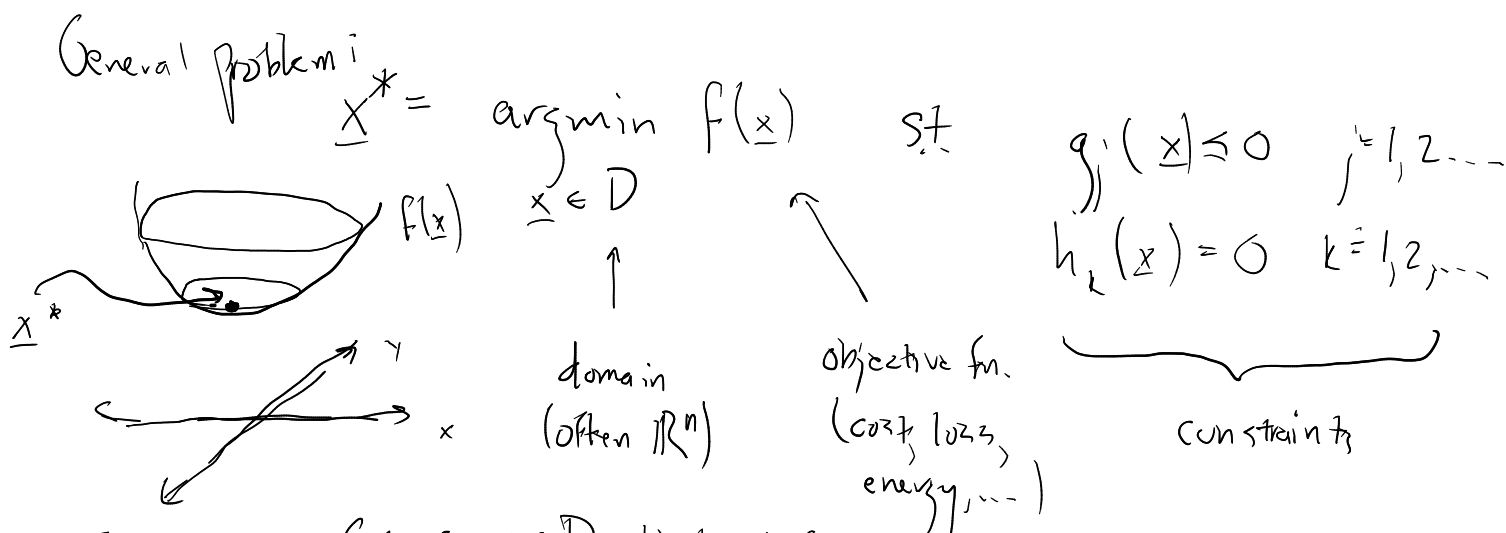
Topics: 1) Problem definition, types of problems
2) Convex problems
3) Solution methods
4) SVMs

Boyd & Vandenberghe

<u>Convex Optimization</u>

<span style="color:green">Green notes are extra examples and extensions not covered in class.</span>

General problem:



$$\underline{x}^* = \arg\min_{\underline{x} \in D} f(\underline{x}) \quad s.t. \quad g_j(\underline{x}) \leq 0 \quad j=1,2,\dots$$
$$h_k(\underline{x}) = 0 \quad k=1,2,\dots$$

$f(\underline{x})$

domain (often $\mathbb{R}^n$)

objective fn. (cost, loss, energy,...)

constraints

Feasible set: Set of $\underline{x} \in D$ that satisfy constraints.

Brute force (grid search): Divide $D$ into grid w/ width $\delta$. # points grows as $\left(1/\delta\right)^n$

| Type | Domain | Objective | Constraints | | Solution |
|------|--------|-----------|-------------|---|----------|
| Linear | $\mathbb{R}^n$ | $\underline{c}^T\underline{x}$ | $A\underline{x}\leq\underline{b},$ | $\underline{x}\geq 0$ | Easy (simplex method) |
| Integer | $\mathbb{N}^n$ | '' '' | '' | '' | NP-hard in general |
| Constraint sat. | $\{0,1\}^n$ | Constant | Boolean | | NP hard in general |
| Convex | $\mathbb{R}$ | Convex fn. | Convex set | | Easy! (interior-point method) |
| Quadratic | $\mathbb{R}^n$ | $\underline{x}^T Q\underline{x} + \underline{c}^T\underline{x}$ | $A\underline{x}\leq b$ | | Easy if $Q$ positive-definite |

...

**Ex (least squares)**     $y = \underline{x} \cdot \beta$.  Given $\{\underline{x}_i, y_i\}$, $i = 1...P$, find optimal $\beta$

$$\beta^* = \underset{\beta}{\text{argmin}} \; \| X\beta - y \|^2$$

$\underset{P \times N}{\uparrow} \quad \underset{N \times 1}{\uparrow} \quad \underset{P \times 1}{\uparrow}$

$$f(\beta) = (X\beta - y)^T (X\beta - y)$$

$$= \beta^T X^T X \beta - 2 y^T X \beta + \underbrace{y^T y}_{\text{constant, Ignore}}$$

$$\beta^* = \underset{\beta}{\text{argmin}} \; \frac{1}{2} \beta^T Q \beta + \underline{c}^T \beta, \quad Q = X^T X$$

$$\underline{c} = -2 y^T X$$

Quadratic problem (convex)

Regularization:

$$f(\beta) = \underbrace{(X\beta - y)^T (X\beta - y)}_{\text{fitting to data}} + \underbrace{\lambda \beta^T \beta}_{\text{penalty on large } \beta_i^2}$$
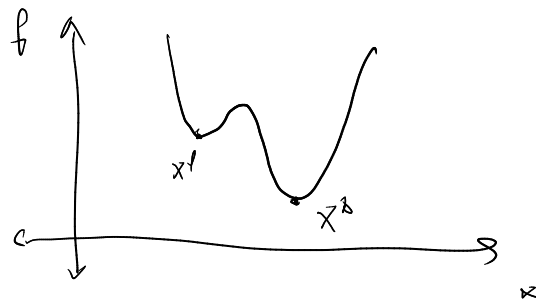
$$= \lambda \beta^T I \beta$$

Same as above with

$$Q \leftarrow Q + \lambda I$$
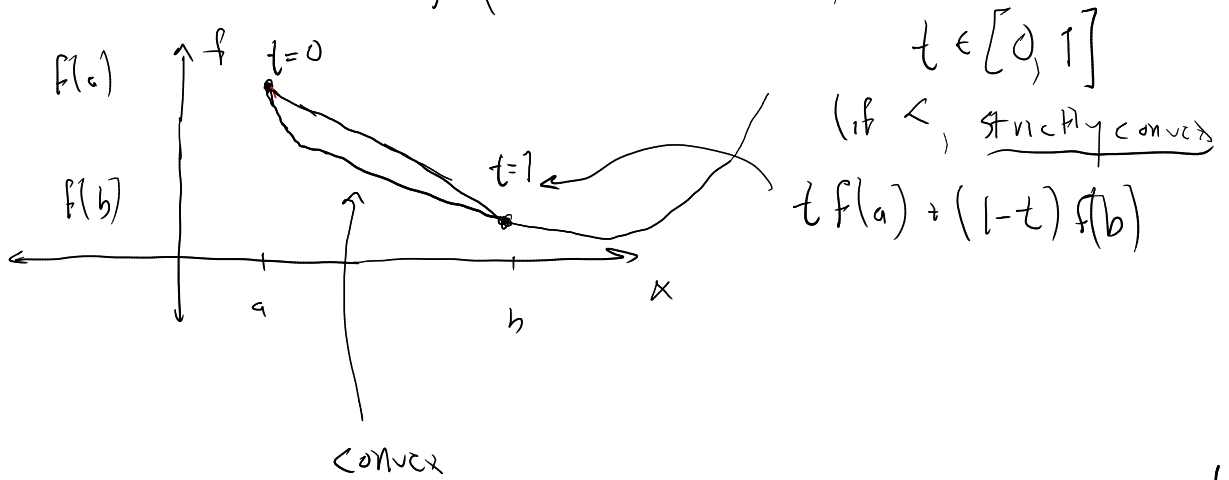
$\underline{x}^*$ is global optimum. May be local optima $\underline{x}^l$ s.t.

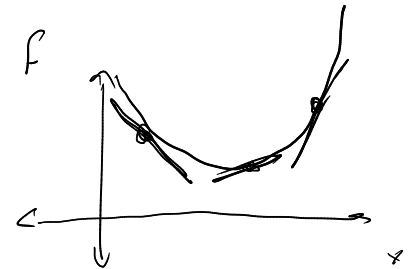$$f(\underline{x}_l) < f(\underline{x}) \quad \text{for} \quad \|\underline{x} - \underline{x}^l\| < \varepsilon.$$



**Def:** $f(\underline{x})$ is convex if $f\left(t\underline{a} + (1-t)\underline{b}\right) \leq t f(\underline{a}) + (1-t)f(\underline{b})$

$$t \in [0, 1]$$

(if $<$, strictly convex)

$$t f(a) + (1-t) f(b)$$



convex

$\Rightarrow$ If $\nabla f(a)$ is slope at $a$,

$$f(a) + \nabla f(a)\left[x - a\right] < f(x)$$

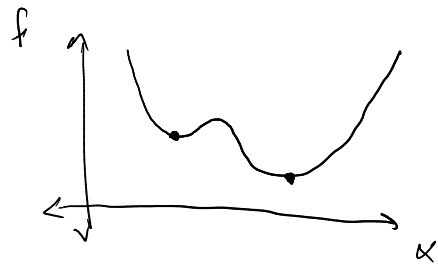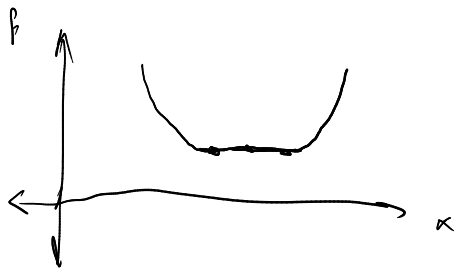If $f$ convex, all local min. are global min.

Strictly convex, one global min.

**Ex:** $f(x) = x^2$. $\quad f(ta + (1-t)b) = t^2 a^2 + (1-t)^2 b^2 + 2t(1-t)ab$   (1)

$$t f(a) + (1-t)f(b) = ta^2 + (1-t)b^2 \tag{2}$$

$$(1) - (2) = (t^2 - t)a^2 + ((1-t)^2 - (1-t))b^2 + 2t(1-t)ab$$

$$= t(t-1)a^2 + t(t-1)b^2 - 2t(t-1)ab = t(t-1)(a-b)^2 < 0$$

**Higher d:**

Gradient of $f(\underline{x})$: $\quad \nabla f(\underline{x}) = \begin{pmatrix} \partial f / \partial x_1 \\ \partial f / \partial x_2 \\ \vdots \\ \partial f / \partial x_n \end{pmatrix}$

Hessian:

$$H(\underline{x}) = \begin{pmatrix} \partial^2 f / \partial x_1^2 & \partial^2 f / \partial x_1 \partial x_2 \\ \partial^2 f / \partial x_2 \partial x_2 & \partial^2 f / \partial x_2^2 \\ & & \ddots \end{pmatrix}$$

1d: $\quad f(x) \approx f(a) + f'(a)(x-a) + \frac{1}{2} f''(a)(x-a)^2$

Minimum: $f' = 0, \quad f'' > 0$.

Strictly convex if $f'' > 0$ everywhere.

Higher-d: $\quad f(\underline{x}) \approx f(\underline{a}) + \nabla f(\underline{a})(\underline{x}-\underline{a}) + \frac{1}{2}(\underline{x}-\underline{a})^T H(\underline{a})(\underline{x}-\underline{a})$

Minimum: $\nabla f = 0$, $H$ positive definite $\left(\underline{v}^T H \underline{v} > 0 \; \forall \underline{v}\right)$

Strictly convex if $H$ positive definite everywhere

Analytical approaches:

Unconstrained problem: look for $\underline{x}^*$ w/ $\nabla f(\underline{x}^*) = 0$

Equality constraints $\longrightarrow$ method of Lagrange multipliers.

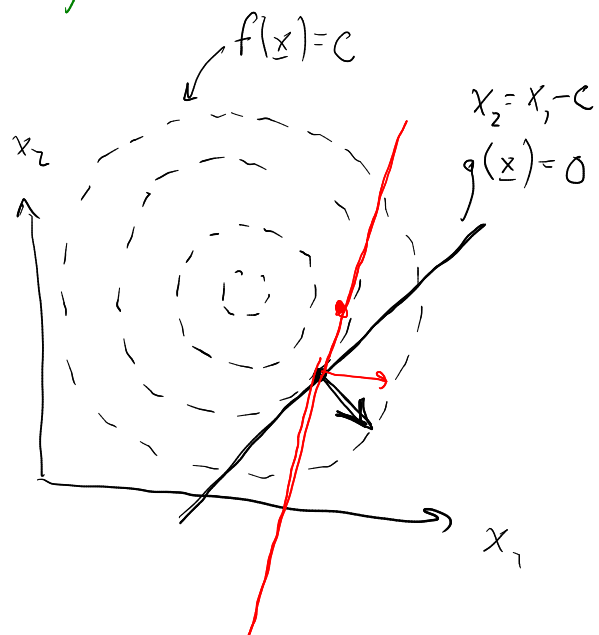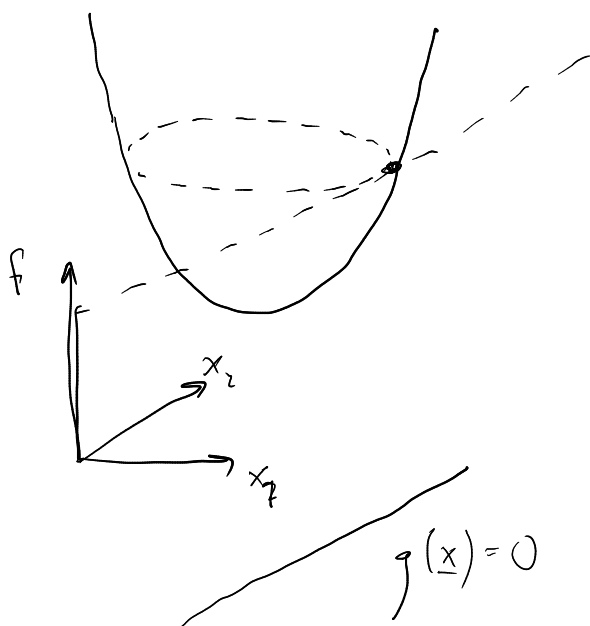$\min_{\underline{x}} f(\underline{x})$    s.t.    $g_j(\underline{x}) = 0$     $j = 1, 2, \cdots$

Let $\mathcal{L}(\underline{x}, \underline{\lambda}) = f(\underline{x}) - \sum_j \lambda_j \, g_j(\underline{x})$

$\lambda_j$ are Lagrange multipliers.

Look for $\underline{x}^*, \underline{\lambda}$   s.t.   $\nabla \mathcal{L}(\underline{x}, \underline{\lambda}) = 0$

Note $\dfrac{\partial}{\partial \lambda_j} \mathcal{L} = - g_j(\underline{x}) = 0$

$\nabla_{\underline{x}} f(\underline{x}) - \sum_j \lambda_j \nabla_{\underline{x}} g_j(\underline{x}) = 0$   "level sets"



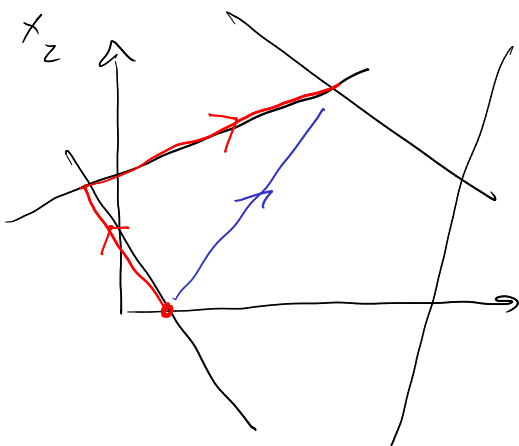$f(\underline{x}) = c$

$x_2 = x_1 - c$
$g(\underline{x}) = 0$

$g(\underline{x}) = 0$

With inequality constraints $g_j(\underline{x}) \leq 0$, $h_k(\underline{x}) = 0$.

Karush-Kuhn Tucker (KKT) conditions

Numerical approaches:

Linear problems:

$$\min \underline{c}^T \underline{x} \quad \text{s.t.} \quad A\underline{x} \leq \underline{b}, \quad \underline{x} > 0$$
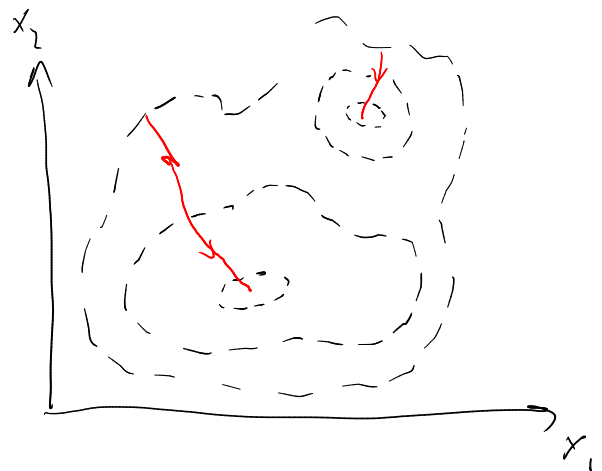
Simplex method.

Convex problems: Interior point methods.

Gradient descent:

Given initial value $\underline{x}_0$,

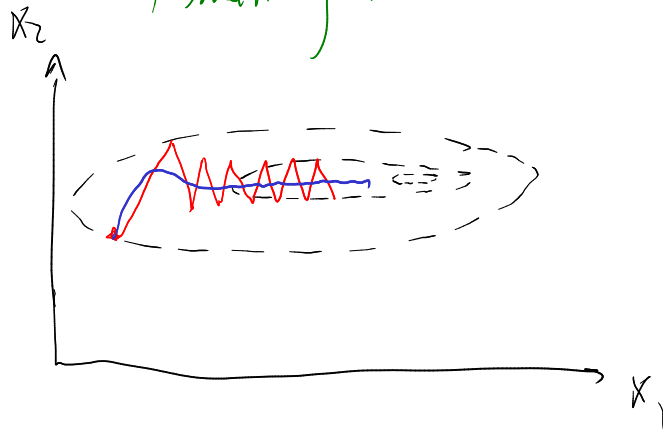$$\underline{x}_{n+1} = \underline{x}_n - \eta_n \nabla f(\underline{x}_n)$$

step size

"Learning rate"

Often $\eta_{n+1} < \eta_n$.

"line search": Choose $\eta_n$ to minimize $f(\underline{x}_{n+1})$.

Problems:   1) multiple minima   (multiple initial conditions, noise)
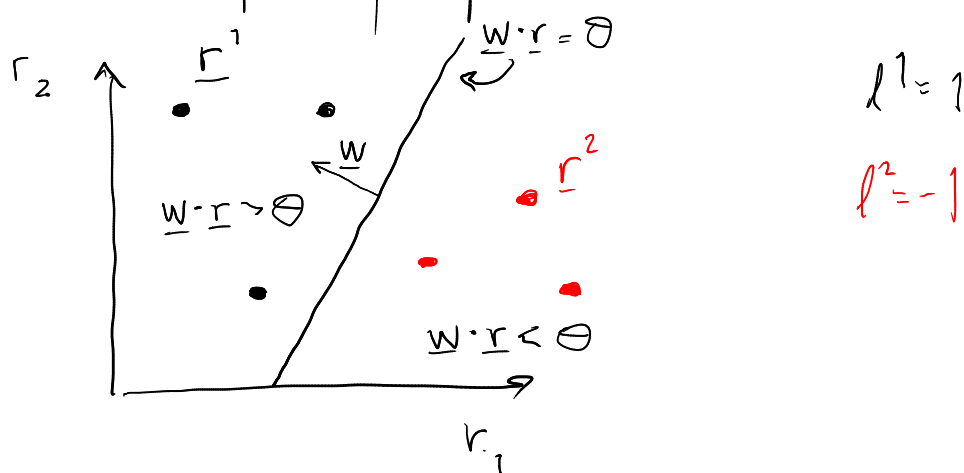
   2) small gradients



Momentum:   $\underline{z}_{n+1} = \beta \underline{z}_n + \nabla f(\underline{x}_n)$        $\beta = 0.99$

   $\underline{x}_{n+1} = \underline{x}_n - \eta_n \underline{z}_{n+1}$

Convex optimization example: Support vector machines (SVM)

Problem: Given $P$ patterns $\underline{r}^\mu$, $\mu = 1 \ldots P$, and labels $\ell^\mu = \pm 1$, find linear separating hyperplane that optimally separates $\ell = 1$ and $\ell = -1$.
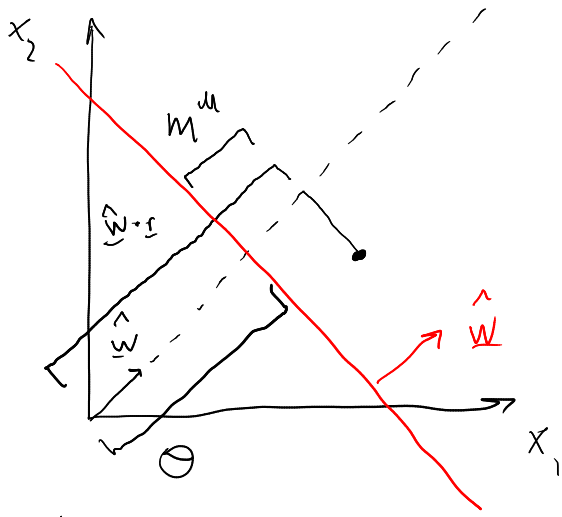


$\ell^1 = 1$

$\ell^2 = -1$

Classifier: $\ell = \text{sign}\left( \underline{w} \bullet \underline{r} - \Theta \right)$

weights      threshold

How to choose optimal $\underline{w}, \Theta$? May be multiple valid solutions (draw).

Idea: Maximize $\underline{\text{margin}}$ $m^\mu$ ($\overset{\text{smallest}}{\vee}$ distance from $\underline{r}^\mu$ to boundary). (draw).

If $\|\underline{w}\| = 1$, then $m^\mu = | \hat{\underline{w}} \bullet \underline{r}^\mu - \Theta |$

$\underline{W}$ is perpendicular to separating hyperplane.

Optimization problem:

$$\text{maximize}_{\underline{w}} \quad \min_{\mu} \left| \underline{w} \cdot \underline{r}^{\mu} - \theta \right| \quad \text{s.t.} \quad \|\underline{w}\| = 1,$$
$$\text{sign}\left( \underline{w} \cdot \underline{r}^{\mu} - \theta \right) = \ell^{\mu}$$

Redefine constraints: $\left( \underline{w} \cdot \underline{r}^{\mu} - \theta \right) \cdot \ell^{\mu} > 0.$

How to deal with $\|\underline{w}\|$?

If margin $= M$, $\left( \underline{w} \cdot \underline{r}^{\mu} - \theta \right) \ell^{\mu} \geq m$. Divide by $m$:

$$\left( \frac{\underline{w}}{m} \cdot \underline{r}^{\mu} - \frac{\theta}{m} \right) \ell^{\mu} \geq 1$$

Reparameterize: $\tilde{\underline{w}} \leftarrow \frac{\underline{w}}{M}$

$$\tilde{\theta} = \frac{\theta}{m}$$

$$\left( \tilde{\underline{w}} \cdot \underline{r}^{\mu} - \tilde{\theta} \right) \ell^{\mu} \geq 1.$$

Note: $\|\tilde{\underline{w}}\| = \left\|\frac{\underline{w}}{m}\right\| = \frac{1}{m}$

Maximize margin $\iff$ minimize $\|\tilde{\underline{w}}\|^2$ !

$$\underline{w}^*, \theta^* = \operatorname*{argmin}_{\underline{w}, \theta} \underline{w}^T \underline{w} \quad \text{s.t.} \left(\underline{w} \cdot \underline{r}^\mu - \theta\right) \ell^\mu \geq 1.$$

$$\implies \|\underline{w}^*\| = \frac{1}{m}.$$

Properties: 1) $\underline{w}$ determined only by closest points

(those on margin) — $\underline{\text{support vectors}}$.

2) Binary classification — linear boundary

3) Fully supervized ($\ell^\mu$ know $\forall \mu$).

4) Soln only exists if data $\underline{\text{linearly separable}}$

Extensions: 1) Multi-class

2) Kernel SVM (later)

3) "Soft margin"
$\downarrow$

Allow misclassifications:

Penalize w/ $c^\mu = \max\left(0, 1 - \ell^\mu\left(\underline{w} \cdot \underline{r}^\mu - \theta\right)\right)$

$\text{Min} \sum_{\mu=1}^{P} c^\mu + \lambda \|\underline{w}\|^2$

"hinge loss"

How to write as convex problem?

Rewrite constraints: $(\underline{w}\cdot\underline{r}^\mu - \Theta)\ell^\mu \geq 1 - c^\mu$

$$\underline{w}^*, \underline{c}^*, \Theta^* = \underset{\underline{w}, \underline{c}, \Theta}{\text{argmin}} \sum_{\mu=1}^{P} c^\mu + \lambda \underline{w}^T\underline{w} \quad \text{s.t.} \; (\underline{w}\cdot\underline{r}^\mu - \Theta)\ell^\mu \geq 1 - c^\mu$$

Interpretation of Lagrange mult:

$$\mathcal{L} = f(\underline{x}) + \lambda g(\underline{x}). \qquad \text{Let } g(\underline{x}) = \tilde{g}(\underline{x}) - c = 0$$

$$\Longrightarrow \frac{\partial \mathcal{L}}{\partial c} = \lambda. \quad \lambda \text{ is sensitivity of } \mathcal{L} \text{ to change in constraint}$$

For SVMs, $\lambda \neq 0$ only for SVs (on margin).

Duality:   "primal problem"

$$\underline{x}^* = \underset{\underline{x}}{\text{argmin}} \; f(\underline{x}), \qquad \begin{array}{l} g_j(\underline{x}) \leq 0 \\ h_k(\underline{x}) = 0 \end{array}$$

Write Lagrangian

$$\mathcal{L}(x, \underline{\lambda}, \underline{\nu}) = f(\underline{x}) + \sum_j \lambda_j \, g_j(\underline{x}) + \sum_k \nu_k \, h_k(\underline{x})$$

Dual function:

$$G(\underline{\lambda}, \underline{\nu}) = \underset{\underline{x} \in D}{\inf} \mathcal{L}(\underline{x}, \underline{\lambda}, \underline{\nu})$$

Note $G(\underline{\lambda}, \underline{v}) \leq f^*$ if $\lambda_j \geq 0 \ \forall j.$

Why? For any __feasible__ $\underline{x}$, $g_k(\underline{x}) \leq 0$, $h_k(\underline{x}) = 0$

$\implies \underbrace{\sum_j \lambda_j g_j(\underline{x})}_{\leq 0} + \underbrace{\sum_k v_k h_k(\underline{x})}_{= 0} \leq 0$

$\implies \mathcal{L}(\underline{x}, \underline{\lambda}, \underline{v}) \leq f(\underline{x})$

Can minimize $G$ "dual problem":

$\underline{\lambda}^*, \underline{v}^* = \arg\max_{\underline{\lambda}, \underline{v}} G(\underline{\lambda}, \underline{v})$ s.t. $\lambda_j \geq 0$

If primal problem convex, $G^* = f^*$

For SVM,

$$\mathcal{L} = \frac{1}{2} \underline{w}^T \underline{w} - \sum_\mu \lambda^\mu \left[ (\underline{w} \cdot \underline{r}^\mu - \theta) \ell^\mu - 1 \right]$$

$G(\underline{\lambda}) = \inf_{\underline{w}, \theta} \mathcal{L}$

$\dfrac{\partial \mathcal{L}}{\partial w_i} = 0 \implies \underline{w} - \sum_\mu \lambda_\mu \underline{r}^\mu \ell^\mu = 0$

$\dfrac{\partial \mathcal{L}}{\partial \theta} = 0 \implies \sum_\mu \lambda^\mu \ell^\mu = 0$

$\underline{w} = \sum_\mu \lambda_\mu \underline{r}^\mu \ell^\mu \quad \leftarrow \text{sum of SVs}$

Dual problem:

$$\max -\frac{1}{2} \sum_{\mu} \sum_{\nu} \lambda_\mu \lambda_\nu \ell_\mu \ell_\nu \left(\underline{r}^\mu\right)^T \underline{r}^\nu + \sum_{\mu} \lambda^\mu$$

$$s.t. \sum_{\mu} \lambda_\mu \ell_\mu = 0, \quad \lambda_\mu > 0.$$

Optimization over $P$ variables vs. $n$.