

Optimization

Topics: 1) Problem definition, types of problems

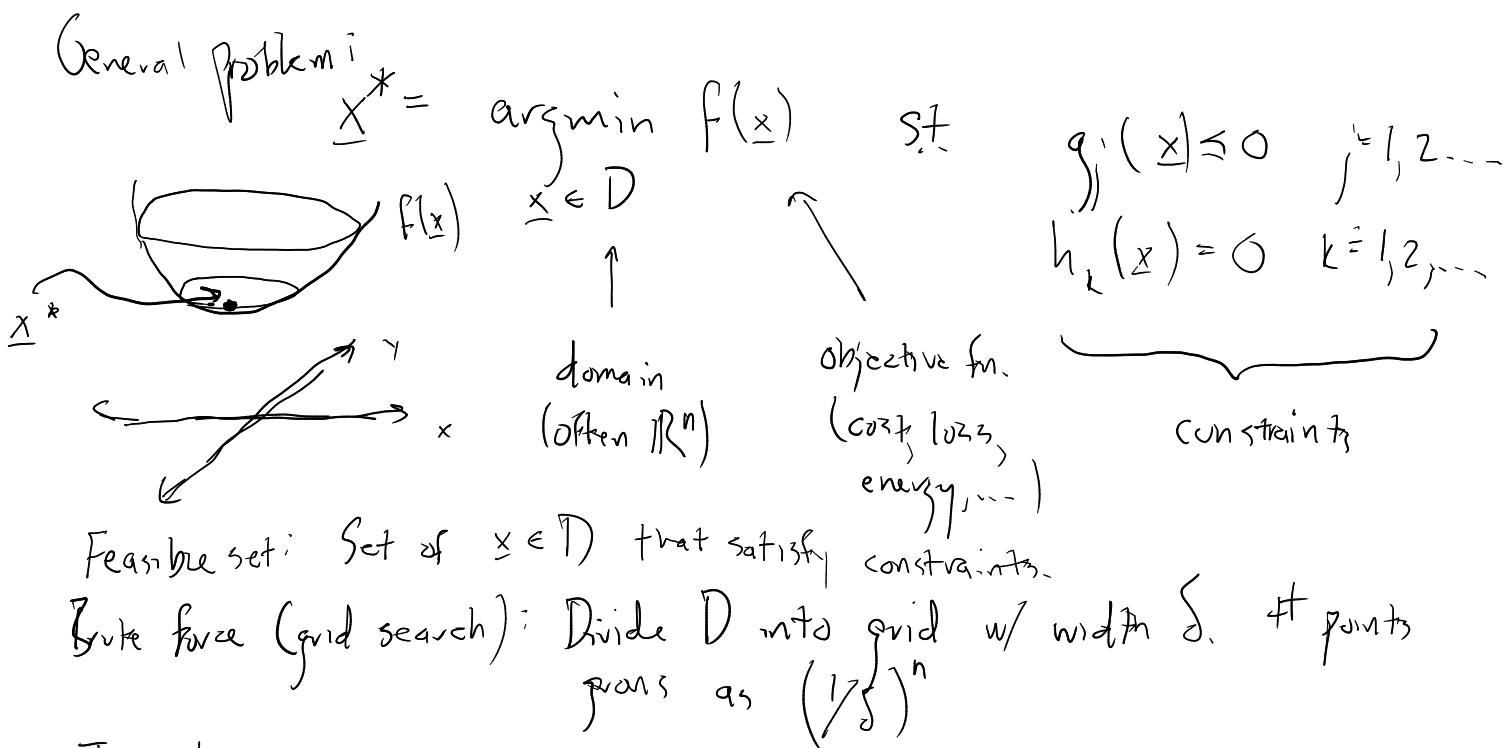
2) Convex problems

3) Solution methods,

4) SVMs

Boyd & Vandenberghe

Convex Optimization



Type	Domain	Objective	Constraints	Solution
Linear	\mathbb{R}^n	$\underline{c}^T \underline{x}$	$A \underline{x} \leq b, \underline{x} \geq 0$	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) $y = X\beta$. Given $\{x_i, y_i\}$, $i=1..P$, find optimal β

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

$\begin{matrix} \uparrow & \uparrow \\ P \times N & N \times 1 \end{matrix}$

$$\begin{aligned} f(\beta) &= (X\beta - y)^T (X\beta - y) \\ &= \beta^T X^T X \beta - 2y^T X \beta + \underbrace{y^T y}_{\text{constant, ignore}} \end{aligned}$$

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

$$c = -2y^T X$$

Quadratic problem (convex)

Regularization:

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

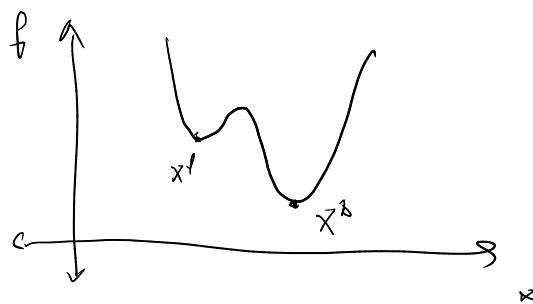
Same as above with

$$= \lambda \beta^T I \beta$$

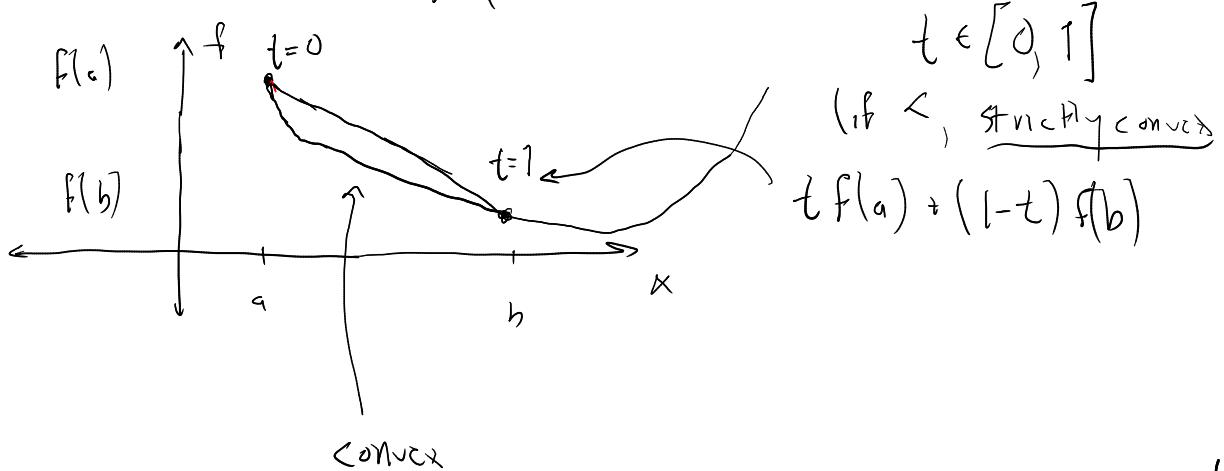
$$Q \leftarrow Q + \lambda I$$

\underline{x}^* is global optimum. May be local optima \underline{x}^l s.t.

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

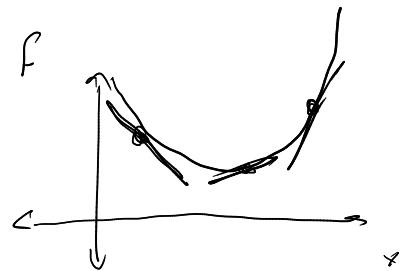


Def: $f(\underline{x})$ is convex if $f(t\underline{a} + (1-t)\underline{b}) \leq t f(\underline{a}) + (1-t)f(\underline{b})$



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

$$f(\underline{a}) + \nabla f(\underline{a})[\underline{x} - \underline{a}] \leq f(\underline{x})$$



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

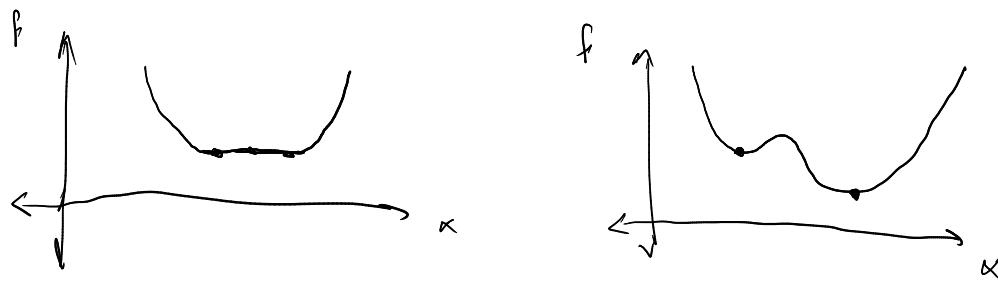
strictly convex, one global min.

$$\text{Ex: } f(\underline{x}) = \underline{x}^2. \quad f(t\underline{a} + (1-t)\underline{b}) = t^2 \underline{a}^2 + (1-t)^2 \underline{b}^2 + 2t(1-t)\underline{a}\underline{b} \quad (1)$$

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

$$(1)-(2) = (t^2 - t)\underline{a}^2 + ((1-t)^2 - (1-t))\underline{b}^2 + 2t(1-t)\underline{a}\underline{b}$$

$$= t(t-1)\underline{a}^2 + t(t-1)\underline{b}^2 - 2t(t-1)\underline{a}\underline{b} = t(t-1)(\underline{a} - \underline{b})^2 \leq 0$$



Higher d:

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

Hessian:

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

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

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

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

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

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

Strictly convex if H positive definite everywhere

Analytical approaches:

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

Equality constraints \rightarrow method of Lagrange multipliers.

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

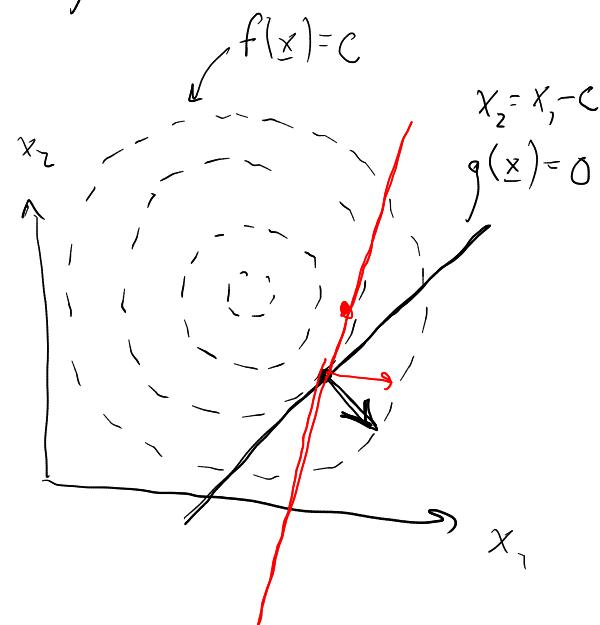
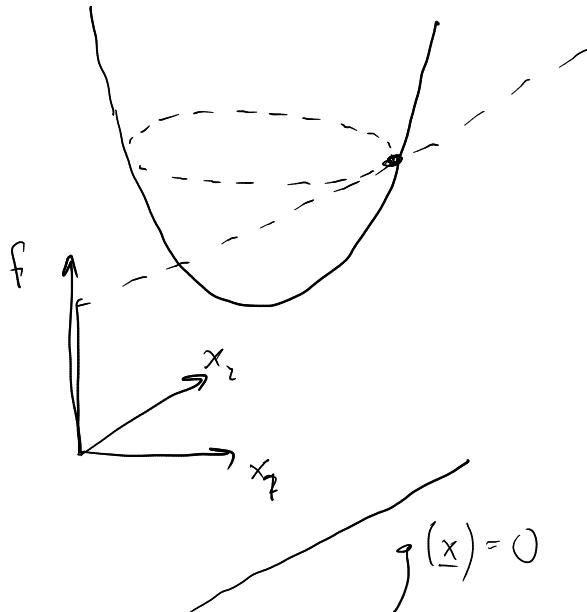
$$\text{let } \mathcal{L}(\underline{x}, \lambda) = f(\underline{x}) - \sum_j \lambda_j g_j(\underline{x})$$

λ_j are Lagrange multipliers.

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

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

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

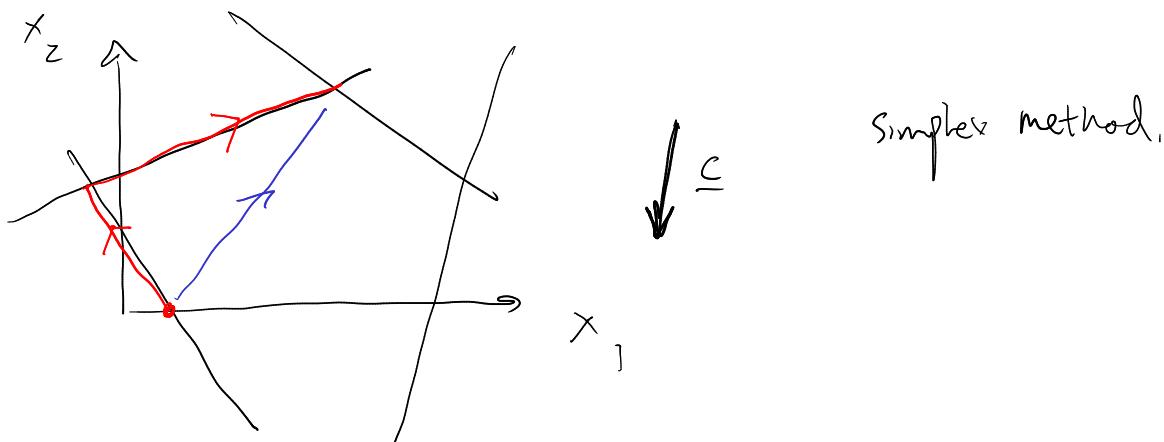


With inequality constraints $g_j(x) \leq 0, h_k(x) = 0$,
 Karush-Kuhn-Tucker (KKT) conditions

Numerical approaches:

Linear problems:

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



Simplex method.

Convex problems: Interior point methods.

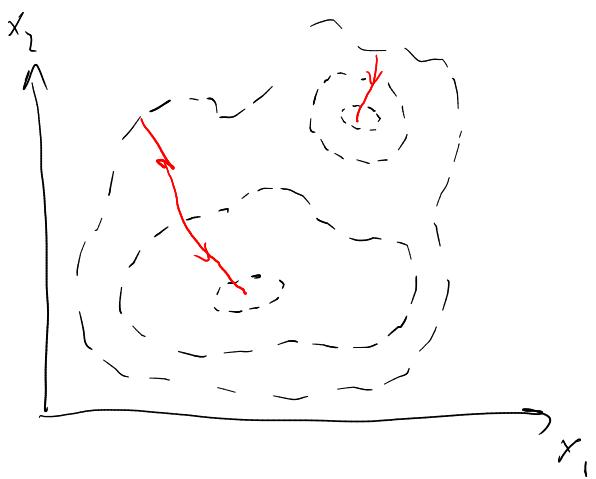
Gradient descent:

Given initial value \underline{x}_0 ,

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

↑
step size

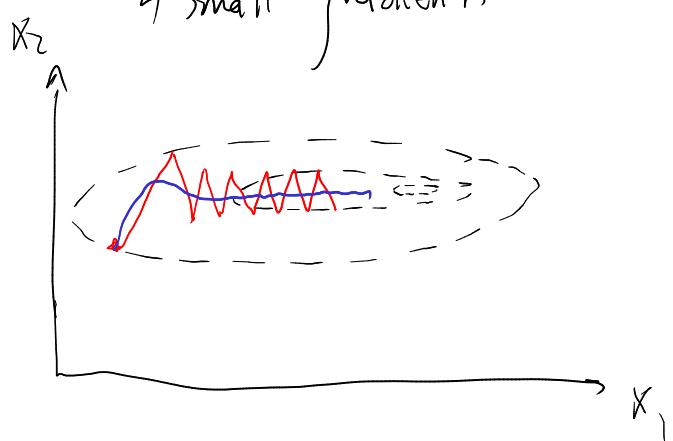
"learning rate"



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

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

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



Momentum:

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

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

Ex
Last time: Lyapunov fn. for Hopfield net.

$$\mathcal{L}(\underline{x}) = \sum_i \int_{x_i}^{x_i} dz_z F'(z) - F(x_i) I - \frac{1}{2} F(x_i) \sum_j W_{ij} F(x_j)$$

$$\frac{dx_i}{dt} = -\frac{\partial \mathcal{L}}{\partial x_i} F'(x_i) \Rightarrow \frac{d\underline{x}}{dt} = -(\nabla \mathcal{L}) \cdot (F'(\underline{x}))$$

If F linear, $F' = 1$, and $\underline{x}(t)$ follows gradient of \mathcal{L} .

If $F' > 0$, $\underline{x}(t)$ moves in direction that matches sign of gradient.

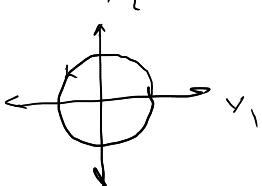
Can we always write ODE as gradient system?

$$\dot{\underline{x}} = -\nabla f(\underline{x})$$

No, Necessary & sufficient:

$$\begin{aligned}\dot{x}_i &= -\frac{\partial}{\partial x_i} f \Rightarrow \frac{\partial}{\partial x_j} \dot{x}_i = \frac{\partial}{\partial x_i} \dot{x}_j \quad \forall i, j \\ \dot{x}_j &= -\frac{\partial}{\partial x_j} f\end{aligned}$$

Counterexample:



$$\dot{\underline{x}} = \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix} \underline{x} \quad \begin{aligned}\dot{x}_1 &= x_2 & \frac{\partial}{\partial x_2} \dot{x}_1 &= 1 \\ \dot{x}_2 &= -x_1 & \frac{\partial}{\partial x_1} \dot{x}_2 &= -1\end{aligned}$$

Can write Lyapunov function $\mathcal{L}(t) = \frac{1}{2} \|\nabla f\|^2 \Rightarrow \frac{d\mathcal{L}}{dt} = \sum_i \frac{\partial \mathcal{L}}{\partial x_i} \frac{\partial x_i}{\partial t}$

$$= -\frac{1}{2} \sum_i \left(\frac{\partial f}{\partial x_i} \right)^2 = -\sum_i (\dot{x}_i)^2$$

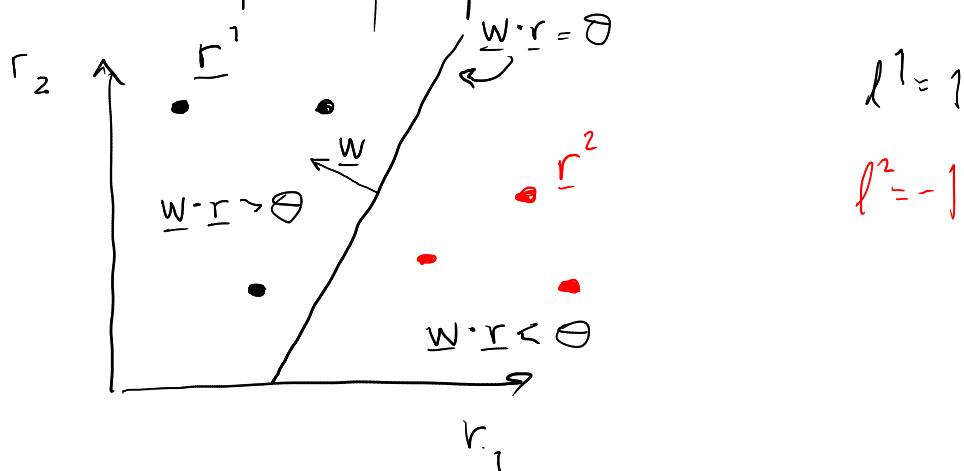
Gradient system cannot have periodic orbits:

If $\underline{x}(0) = \underline{x}(T)$, $\mathcal{L}(0) = \mathcal{L}(T)$. But

$$\mathcal{L}(T) = \int_0^T \frac{d\mathcal{L}}{dt} dt = \int_0^T -\sum_i (\dot{x}_i)^2 dt \Rightarrow \dot{x}_i = 0$$

- Convex optimization example: Support vector machines (SVM)

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

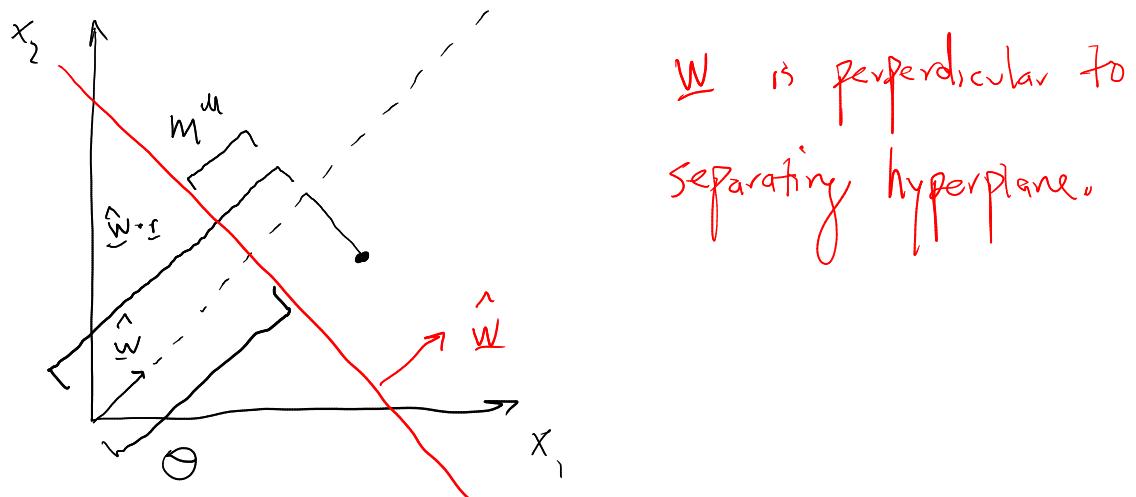


Classifier: $\ell = \text{Sign} \left(\frac{\underline{w} \cdot \underline{r} - \theta}{\uparrow \text{weights} \quad \uparrow \text{threshold}} \right)$

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

Idea: Maximize Margin m^u (^{smallest} distance from \underline{r}^u to boundary). (draw).

If $\|\underline{w}\|=1$, then $m^u = |\hat{w} \cdot \underline{r}^u - \theta|$



w is perpendicular to
separating hyperplane.

Optimization problem:

$$\underset{\underline{w}}{\text{maximize}} \quad \min_{\underline{w}} \left| \underline{w} \cdot \underline{r}^u - \theta \right| \quad \text{st. } \|\underline{w}\| = 1,$$

$$\text{Sign}(\underline{w} \cdot \underline{r}^u - \theta) = \underline{l}^u$$

Redefine constraints: $(\underline{w} \cdot \underline{r}^u - \theta) \cdot \underline{l}^u > 0$.

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

If margin = M, $(\underline{w} \cdot \underline{r}^u - \theta) \underline{l}^u \geq M$. Divide by M:

$$\left(\frac{\underline{w}}{M} \cdot \underline{r}^u - \frac{\theta}{M} \right) \underline{l}^u \geq 1$$

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

$$(\tilde{\underline{w}} \cdot \underline{r}^u - \tilde{\theta}) \underline{l}^u \geq 1.$$

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

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

$$\underline{w}^*, \theta^* = \underset{\underline{w}, \theta}{\operatorname{argmin}} \quad \underline{w}^T \underline{w} \quad \text{s.t. } (\underline{w} \cdot \underline{r}^u - \theta) l^u \geq 1.$$

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

Properties:

- 1) \underline{w} determined only by closest points (those on margin) — support vectors.

- 2) Binary classification — linear boundary

- 3) Fully supervised (l^u know $\forall u$)

- 4) Sol'n only exists if data linearly separable

Extensions:

- 1) Multi-class

- 2) Kernel SVM (later)

- 3) "Soft margin"

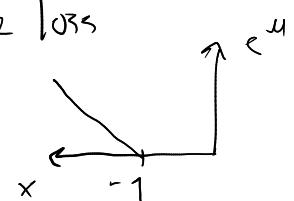


Allow misclassifications:

$$\text{Penalize w/ } c^u = \max(0, 1 - l^u(\underline{w} \cdot \underline{r}^u - \theta))$$

$$\min \sum_{u=1}^p c^u + \lambda \|\underline{w}\|^2$$

"Hinge loss"



How to write as convex problem?

Rewrite constraints: $(\underline{w} \cdot \underline{r}^u - \theta) l^u \geq 1 - c^u$

$$\begin{array}{ll} \underline{w}^*, \underline{l}, \theta = & \underset{\substack{\underline{w}, \underline{l}, \theta}}{\operatorname{argmin}} \quad \sum_{u=1}^P c^u + \lambda \underline{w}^T \underline{w} \\ & \text{st } (\underline{w} \cdot \underline{r}^u - \theta) l^u \\ & \geq 1 - c^u \end{array}$$

Interpretation of Lagrange mult:

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

$\Rightarrow \frac{\partial \mathcal{L}}{\partial c} = \lambda$. λ is sensitivity of \mathcal{L} to change in constraint.

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

Duality: "primal problem"

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

Write Lagrangian

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

Dual function: \underline{x}^{\min}

$$G(\underline{\lambda}, \underline{v}) = \inf_{\underline{x} \in D} \mathcal{L}(\underline{x}, \underline{\lambda}, \underline{v})$$

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

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

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

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

Can minimize G "dual problem":

$$\underline{\lambda}, \underline{v}^* = \underset{\underline{\lambda}, \underline{v}}{\operatorname{argmax}} G(\underline{\lambda}, \underline{v}) \text{ s.t. } \lambda_j \geq 0$$

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

For SVM,

$$L = \frac{1}{2} \underline{w}^T \underline{w} - \sum_m \lambda^m \left[(\underline{w} \cdot \underline{r}^m - \theta) l^m - 1 \right]$$

$$G(\underline{\lambda}) = \inf_{\underline{w}, \theta} L \quad \frac{\partial L}{\partial w_i} = 0 \Rightarrow \underline{w} - \sum_m \lambda^m \underline{r}^m l^m = 0$$

$$\frac{\partial L}{\partial \theta} = 0 \Rightarrow \sum_m \lambda^m l^m = 0$$

$$\underline{w} = \sum_m \lambda^m \underline{r}^m l^m \leftarrow \text{sum of SVs}$$

Dual problem:

$$\max \frac{1}{2} \sum_u \sum_v \lambda_u \lambda_v l_u l_v (\underline{r}^u)^T \underline{r}^v + \sum_u \lambda^u$$

s.t. $\sum_u \lambda_u l_u = 0, \quad \lambda_u > 0.$

Optimization over P variables vs. n .