

Office hours today: MOVED to 4:30-5:30 on Zoom only.

Convex optimization problems

Lecture 15

November 16, 2021

Beck, sections 8.1-8.3

Convex optimization problems

An optimization problem is convex if it's of the form

$$\min_{x \in C} f(x),$$

f is convex over C

C is a convex set (typically in \mathbb{R}^n)

C : "feasible" set

f : "objective"

$x \in \mathbb{R}^n$ is "feasible" if $x \in C$.

Note: $\max_{x \in C} g(x)$ is also a convex problem if g is concave.

The above formulation is "implicit". "Explicit" formulations typically (not always) take the form:

$$\min_{\underline{x} \in \mathbb{R}^n} f(\underline{x}) \quad \begin{array}{l} \text{"such that"} \\ \text{subject to} \end{array} \quad \begin{array}{l} g_i(\underline{x}) \leq 0, \quad i=1, \dots, m \\ h_j(\underline{x}) = 0, \quad j=1, \dots, n \end{array}$$

for some given convex functions $\{g_i\}_{i=1}^m$ and affine functions $\{h_j\}_{j=1}^n$.

$$\text{Feasible set: } \{ \underline{x} \in \mathbb{R}^n \mid \begin{array}{l} g_i(\underline{x}) \leq 0, \quad i=1, \dots, m \\ h_j(\underline{x}) = 0, \quad j=1, \dots, n \end{array} \} = C \text{ (convex)}$$

Why are convex optimization problems nice?

We have "nice" characterizations of existence/uniqueness to solutions.

Theorem: Suppose f is a convex function over a convex set $C \subset \mathbb{R}^n$. If $\underline{x} \in C$ is a local minimum of f over C , then it solves:

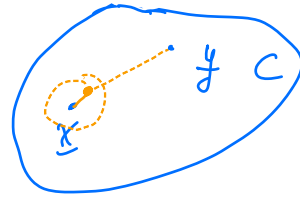
$$\underset{\underline{x} \in C}{\operatorname{argmin}} f(\underline{x}).$$

("Local minima are global minima")

Proof: Let $y \in C$. Since \underline{x} is a local minimum, then there's a small neighborhood around \underline{x} where f is

greater than or equal to $f(\underline{x})$.

I.e., $\exists \varepsilon > 0$ s.t. $\forall \underline{z} \in B(\underline{x}, \varepsilon) \cap C$
 $f(\underline{z}) \geq f(\underline{x})$.



Let $\delta = \varepsilon/2$ (assume $\delta \leq 1$).

$$\text{Note: } \underline{x} + \delta(y - \underline{x}) = \delta y + (1 - \delta)\underline{x} \in C$$

$$f(\underline{x}) \leq f(\underline{x} + \delta(y - \underline{x})) = f(\delta y + (1 - \delta)\underline{x})$$

\nearrow
 $\underline{x} + \delta(y - \underline{x}) \in C, B(\underline{x}, \varepsilon)$

$$\leq \delta f(y) + (1 - \delta)f(\underline{x})$$

$$\delta > 0 \implies f(\underline{x}) \leq f(y). \quad \blacksquare$$

Another property:

Theorem: Suppose f is convex over a convex set C .

Then:

$S = \underset{\underline{x} \in C}{\operatorname{argmin}} f(\underline{x})$ is a convex set.

(The set of all solutions is "nice".)

Proof: Assume $\underline{x}, y \in S$. Let $\lambda \in (0, 1)$

$$f(\lambda \underline{x} + (1 - \lambda)y) \leq \lambda f(\underline{x}) + (1 - \lambda)f(y)$$

$$= [\lambda + (1-\lambda)] \min_{z \in C} f(z)$$

$$= \min_{z \in C} f(z)$$

$$\Rightarrow \lambda x + (1-\lambda)y \in S. \quad \square$$

Corollary: If f is strictly convex over a convex set C , then a global minimizer is unique.

Rest of today = examples.

Linear programming

Applications: scheduling, allotment, cost.

Given $\underline{c} \in \mathbb{R}^n$, $\underline{A} \in \mathbb{R}^{m \times n}$, $\underline{b} \in \mathbb{R}^m$, $\underline{C} \in \mathbb{R}^{k \times n}$, $\underline{d} \in \mathbb{R}^k$

$\min_{\underline{x} \in \mathbb{R}^n} \underline{c}^T \underline{x}$ subject to $\underline{A} \underline{x} \leq \underline{b}$
 $\underline{C} \underline{x} = \underline{d}$) is convex.

This is a prototypical linear programming problem.

Note: $\max_{x \in \mathbb{R}^n} \underline{c}^T \underline{x}$ s.t. $\underline{A} \underline{x} \leq \underline{b}$
 $\underline{C} \underline{x} = \underline{d}$

is also convex since $f(x) = \underline{c}^T \underline{x}$ is both
convex and concave.

Quadratic programming

L15-S04

$$\min_{\underline{x} \in \mathbb{R}^n} \underline{x}^T \underline{A} \underline{x} + 2 \underline{b}^T \underline{x} + c, \quad \underline{A} \in \mathbb{R}^{n \times n}, \quad \underline{A} \succeq \underline{0}, \quad \underline{b} \in \mathbb{R}^n, \quad c \in \mathbb{R}$$

This problem is convex.

A related problem:

$$\min_{\underline{x} \in \mathbb{R}^n} \underline{x}^T \underline{A} \underline{x} + 2 \underline{b}^T \underline{x} + c \quad \text{subject to}$$

$$\underline{x}^T \underline{A}_i \underline{x} + 2 \underline{b}_i^T \underline{x} + c_i \leq 0, \quad i = 1, \dots, m$$

$$\underline{x}^T \underline{C}_i \underline{x} + 2 \underline{d}_i^T \underline{x} + e_i = 0, \quad i = 1, \dots, k$$

QCQP's and convexity

The above problem is a quadratically constrained quadratic program (QCQP).

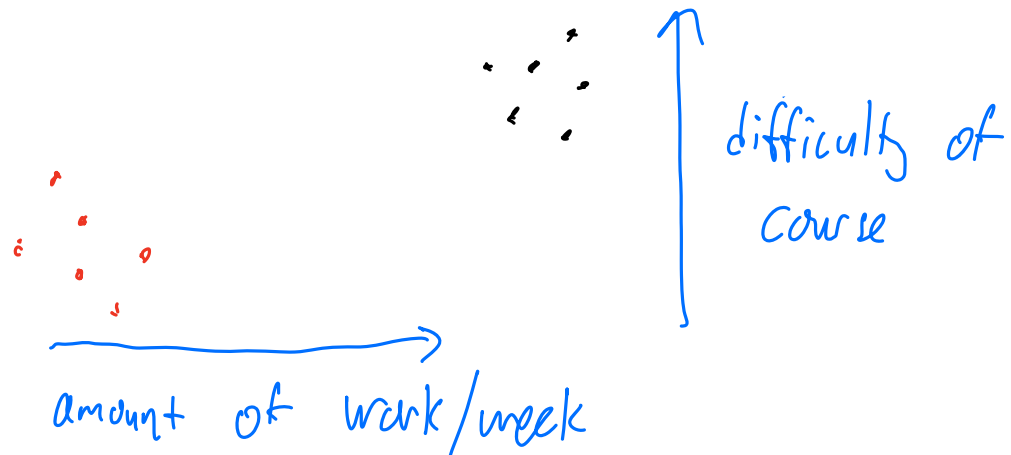
If there are no equality constraints ($k=0$) and $\underline{A}, \underline{A}_i \succeq \underline{0}$, then this is a convex problem.

Example: Chebyshev centers

L15-S06

Goal: generate a region that "best" describes a collection of points.

Courses @ U.



Given all red data: try to fit the smallest circle possible around all the red dots. \rightarrow "Chebyshev center" problem.

Goal: compute the radius and the center of the smallest enclosing ball.

Optimization:

$$\min_{\underline{c}, r} r \quad \text{s.t.} \quad \|\underline{c} - \underline{x}_i\|_2 \leq r, \quad i=1 \dots m$$

↖
of
data points.

This is a convex problem. (QCQP)

Hidden convexity

Hidden convexity: some problems are nonconvex, but under an equivalent "remapping", become convex.

Ex. $\min_{x \in [1, 2]} \log x$ is not convex.

$x = e^y$

$\min_{y \in [0, \log 2]} y$. (this is convex)

"Real" example: QCQP

$$\min_{\underline{x} \in \mathbb{R}^n} \underline{x}^T \underline{A} \underline{x} + 2\underline{b}^T \underline{x} + c \quad \text{subject to}$$
$$\|\underline{x}\|_2 \leq 1.$$

$$\underline{A} \in \mathbb{R}^{n \times n}, \text{ indefinite}$$

This ^{is} _{not} convex.

Under an appropriate mapping, this is a convex problem.

Projections onto convex sets

Idea: find closest point in a convex set to a point outside that set.

$$C \subset \mathbb{R}^n, \text{ convex}$$

$$\underline{x} \in \mathbb{R}^n, \text{ assume } \underline{x} \notin C$$



$$\text{Goal: compute } \underset{y \in C}{\operatorname{argmin}} \|\underline{x} - y\|$$

(if $\underline{x} \in C$, then the argmin is \underline{x})

Theorem: Assume $\|\cdot\|$ is strictly convex (e.g. $\|\cdot\| = \|\cdot\|_2$).

Then given $\underline{x} \in \mathbb{R}^n$, $C \subset \mathbb{R}^n$ that is convex, the problem

$$\operatorname{argmin}_{y \in C} \|\underline{x} - y\|$$

has a unique solution.

and non-empty, and closed

Proof (sketch): The problem is

$$\begin{aligned} \min_{y \in \mathbb{R}^n} \|\underline{x} - y\| \\ \text{s.t. } y \in C \end{aligned}$$

We know: $f(y) = \|\underline{x} - y\|$ is strictly convex.

($\|\cdot\|$ is strictly convex, and $\phi(y) = \underline{I}y - \underline{x}$ is affine)

So: this is optimization of a strictly convex function over a convex set.

\Rightarrow local minima are global minima.

(can show that \exists a single local minimum). \square

Since the above minimization problem is well posed (\exists a unique solution), we can call the solution to this problem as a projection.

Def: Given $C \subset \mathbb{R}^n$ that is closed, convex, then

$P_C: \mathbb{R}^n \rightarrow C$ is the projection operator of the convex set C , defined by solving

$$\operatorname{argmin}_{y \in C} \|x - y\|$$

for any $\underline{x} \in \mathbb{R}^n$. Then $P_C(\underline{x})$ is called the projection (of \underline{x}) onto C .

Although we have nice theory for P_C , actually computing it is pretty difficult, except in special cases.

Projection examples

L15-S09

Ex 1 (Projection onto hypercube)

$$C = [a_1, b_1] \times [a_2, b_2] \times \dots \times [a_n, b_n]$$

$$a_i \leq b_i, \quad i=1, \dots, n.$$

a_i could be $-\infty$

b_i could be $+\infty$



C is closed and convex. ($C = \mathbb{R}_+^n$ is one example)

In ℓ^2 : The minimization problem is

$$\min_{y \in C} \sum_{i=1}^n (y_i - x_i)^2 \quad \text{s.t.} \quad a_i \leq \cancel{x_i} \leq b_i \quad \forall i=1, \dots, n.$$

This problem is "separable": each dimension can be treated independently.

I.e., can consider $\min_{y_j \in \mathbb{R}} (y_j - x_j)^2$ s.t. $a_j \leq y_j \leq b_j$

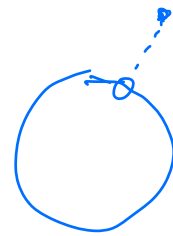
which is solved by:

$$x_i^* = \begin{cases} x_i, & x_i \in [a_i, b_i] \\ a_i, & x_i < a_i \\ b_i, & x_i > b_i \end{cases} \quad (\text{"clamping"})$$

The full solution is $\underline{x}^* = (x_1^*, \dots, x_n^*) = P_C(\underline{x})$

Ex. 2: $\|\cdot\| = \|\cdot\|_2$

$$C = B[0, r], \quad r > 0$$



$$\operatorname{argmin}_{y \in \mathbb{R}^n} \|\underline{x} - y\|_2 \quad \text{s.t. } y \in C$$



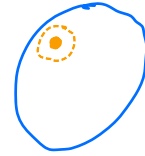
$$\operatorname{argmin}_{y \in C} \|\underline{x} - y\|_2^2$$

($t \mapsto t^2$ is strictly monotone on $[0, \infty)$)

Note: $f(y) = \|x - y\|_2^2$ is quadratic.

Suppose that $P_C(x)$ lies in the interior of C

Since f is strictly convex, then $P_C(x)$ must be a local minimum in some neighborhood.



$$\Rightarrow \nabla f(P_C(x)) = \underline{0}$$

$$\nabla f(y) = \underline{0} \Rightarrow y = x$$

If $x \notin C$, this is not possible.
(because $P_C(x) \notin C$)

contradiction
 \implies

$P_C(x)$ cannot lie in the interior of C .

I.e., $P_C(x)$ lies on the boundary, $\|P_C(x)\|_2 = r$.

$$\text{I.e., } \operatorname{argmin}_{y \in C} \|x - y\|_2^2 = \operatorname{argmin}_{\|y\|_2 = r} \|x - y\|_2^2$$

$$= \operatorname{argmin}_{\|y\|_2 = r} \underbrace{\|x\|_2^2 + \|y\|_2^2 - 2\langle x, y \rangle}_{\text{constant}}$$

$$= \operatorname{argmin}_{\|y\|_2 = r} -2\langle x, y \rangle$$

$$= \operatorname{argmax}_{\|y\|_2 = r} \langle x, y \rangle$$

↑
this is maximized when y is parallel to x
because of Cauchy-Schwarz.

I.e., the max is achieved when $y = \frac{x}{\|x\|_2} r$

$$\implies P_C(x) = \begin{cases} x, & \|x\|_2 \leq r \\ \frac{r}{\|x\|_2} x, & \|x\|_2 > r \end{cases}$$

Ex. 3: Let $\underline{A} \in \mathbb{R}^{n \times n}$ be symmetric.

$$S_+^n = \{ \underline{B} \in \mathbb{R}^{n \times n} \mid \underline{B} \text{ symmetric, } \underline{B} \succeq \underline{0} \}. \quad (\text{set of positive semi-definite matrices.})$$

S_+^n is closed and convex.

$$P_{S_+^n}(\underline{A}) = ? \quad (\text{use matrix 2-norm } \|\cdot\|_2)$$

Let $\underline{A} = \underline{U} \underline{\Lambda} \underline{U}^T$ (eigenvalue decomp.)

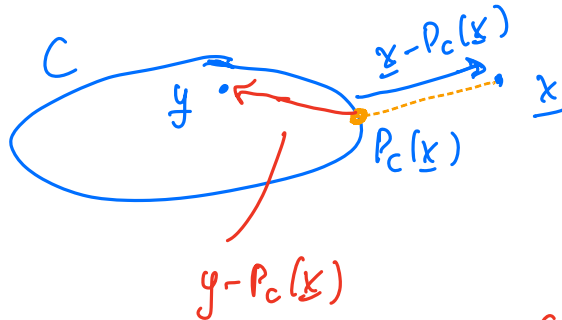
$$\underline{\Sigma} : n \times n \text{ diagonal matrix, with } \sigma_{i,i} = \begin{cases} \lambda_{i,i} & \text{if } \lambda_{i,i} \geq 0 \\ 0 & \text{if } \lambda_{i,i} < 0 \end{cases}$$

$$\text{Then } P_{S_+^n}(\underline{A}) = \underline{U} \underline{\Sigma} \underline{U}^T.$$

Orthogonal projection properties

L15-S10

There's a nice geometric property of this projection operator.



This picture suggests that

$$\langle \underline{x} - P_C(\underline{x}), y - P_C(\underline{x}) \rangle \leq 0 \quad \forall y \in C.$$

Theorem: Let C be ^aconvex, ^{set}closed, let $\underline{x} \notin C$. Then $\underline{x}^* \in C$ satisfies
 $\underline{x}^* = P_C(\underline{x})$ iff $\forall y \in C, \langle \underline{x} - \underline{x}^*, y - \underline{x}^* \rangle \leq 0$

Theorem (P_C is norm non-expansive)

Let C be closed, convex in \mathbb{R}^n . Then $\forall \underline{x}, \underline{y} \in \mathbb{R}^n$, we have

$$\|P_C(x) - P_C(y)\| \leq \|x - y\|.$$

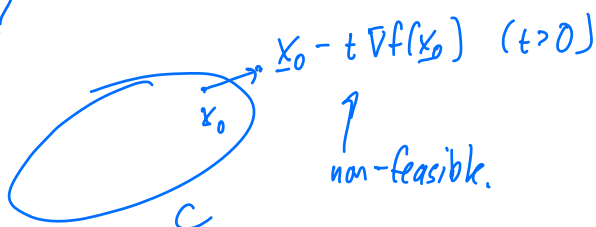
Why do we care about projection onto convex sets?

It's a subproblem in "projected gradient" / "gradient projection" methods.

Suppose we want to numerically compute a solution to

$$\min_{x \in C} f(x), \quad C \text{ is convex, } f \text{ is } \cancel{\text{convex}}.$$

Maybe gradient descent?



Projected gradient method: initialize at x_0

For $i = 1, 2, \dots$

compute $\nabla f(x_{i-1})$

compute stepsize t_i

$$x_i = P_C(x_{i-1} - t_i \nabla f(x_{i-1}))$$

end

This allows us to solve (convex)-constrained optimization problems.

Support vector machines

Goal: classification (binary)

Given: $\{(\underline{x}_i, y_i)\}_{i=1}^M$ data pairs

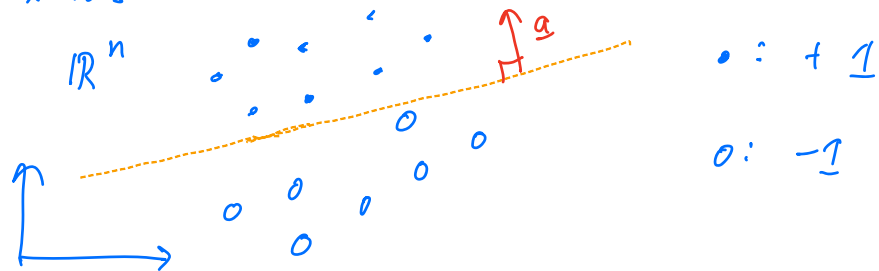
\underline{x}_i : "features" of an item

y_i : ± 1 , a known classification.
"label"

Output: predictor $f(\underline{x}) \rightarrow \pm 1$, $f: \mathbb{R}^n \rightarrow \{+1, -1\}$

Support vector machines is a family of classification techniques for building f .

Hypothesis: feature data (\underline{x}_i) are "linearly" separated in \mathbb{R}^n based on their labels

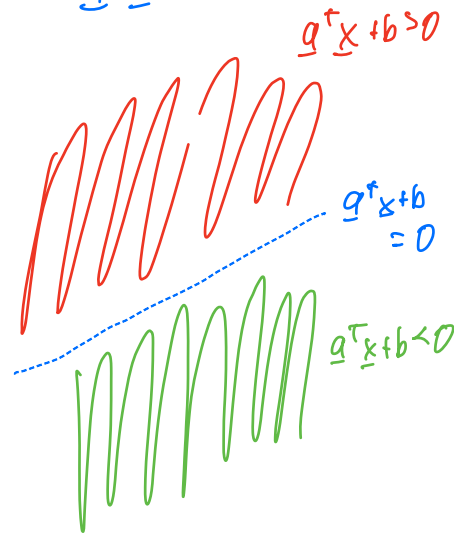


"Linearly separated": \exists a hyper-plane in \mathbb{R}^n that divides \mathbb{R}^n into two half-spaces, each of which contains only data of a single label

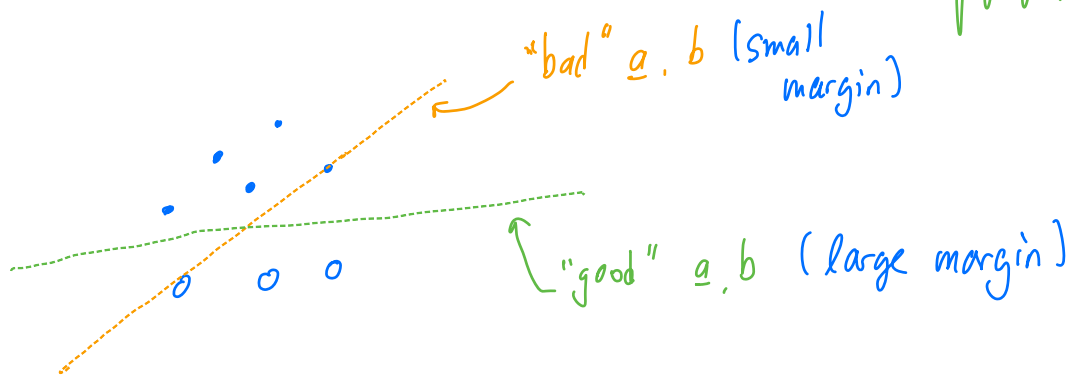
Equations of hyperplanes take the form $\underline{a}^T \underline{x} + b = 0$ for some $\underline{a} \in \mathbb{R}^n \setminus \{0\}$ and $b \in \mathbb{R}$.

\underline{a} : normal vector to hyperplane.

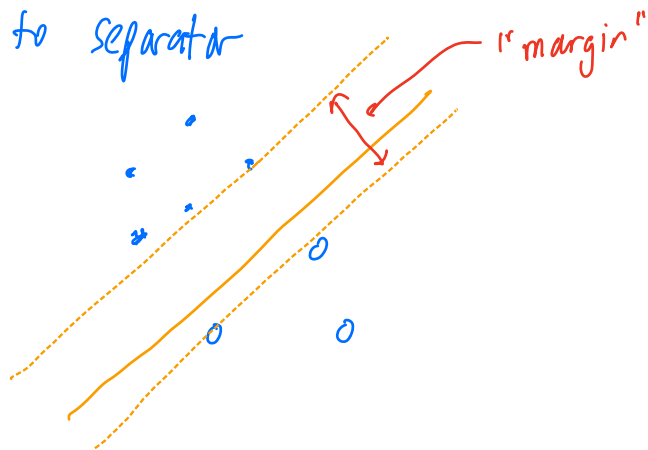
Predictor/classifier: $f(\underline{x}) = \text{sgn}(\underline{a}^T \underline{x} + b)$



This problem doesn't have a unique solution: find the "best" \underline{a}, b .



Margin: distance between hyperplanes defining data closest to separator



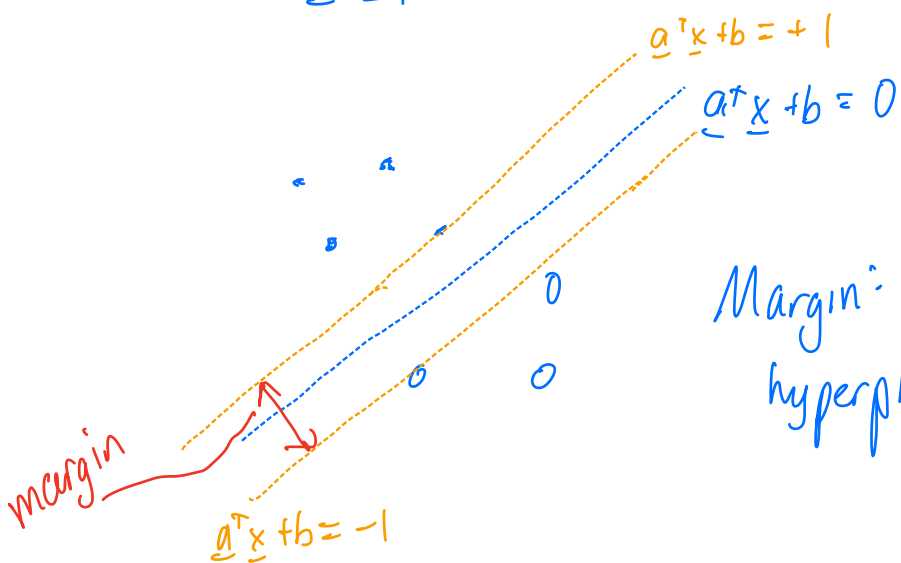
Simplify computation of the margin: can always choose

\underline{a}, b s.t. for all (-1) labelled data,

$\underline{a}^T \underline{x}_i + b = -1$ is the max value

s.t. for all $(+1)$ labelled data,

$\underline{a}^T \underline{x}_i + b = +1$ is the min value



Margin: distance between hyperplanes

The distance between hyperplanes defined by
 $\underline{a}^T \underline{x} + b = +1$ and $\underline{a}^T \underline{x} + b = -1$

is $\frac{2}{\|\underline{a}\|_2}$ (\underline{x} in $\underline{a}^T \underline{x} + b = -1$, replace
by $\underline{\tilde{x}} = \underline{x} + \frac{2\underline{a}}{\|\underline{a}\|_2^2}$, then
 $\underline{a}^T \underline{\tilde{x}} + b = +1$)

Margin 

Optimization for \underline{a}, b :

$\max_{\underline{a}, b}$ (Margin) subject to
 $f(\underline{x}_i) = y_i \quad \forall i = 1 \dots M.$

↓ simplification: $|\underline{a}^T \underline{x}_i + b| \geq 1$
↓ Margin is $2/\|\underline{a}\|_2$

$\max_{\underline{a}, b} \frac{2}{\|\underline{a}\|_2}$ s.t. $y_i (\underline{a}^T \underline{x}_i + b) \geq 1$
 $\forall i = 1 \dots M$

Recall: x_i, y_i are given

$$\begin{array}{c} \updownarrow \\ \max t \leftrightarrow \min 1/t \quad (t \geq 0) \\ \downarrow \\ s \mapsto s^2 \text{ is monotonic on } [0, \infty) \end{array}$$

$$\min_{a, b} \|a\|_2^2 \quad \text{s.t.} \quad y_i (a^T x_i + b) \geq 1 \\ \forall i = 1, \dots, M$$

(Maximizing margin subject to data fidelity)

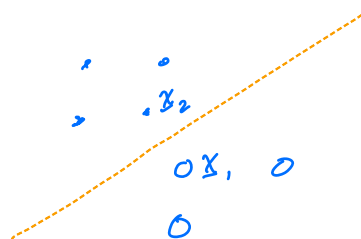
Note: objective $\|a\|_2^2$ is convex.

$y_i (a^T x_i + b) \geq 1$: level sets of convex (affine) functions

\Rightarrow feasible set is convex.

\Rightarrow Support vector machine opt. is a convex problem. (It's quadratic.)

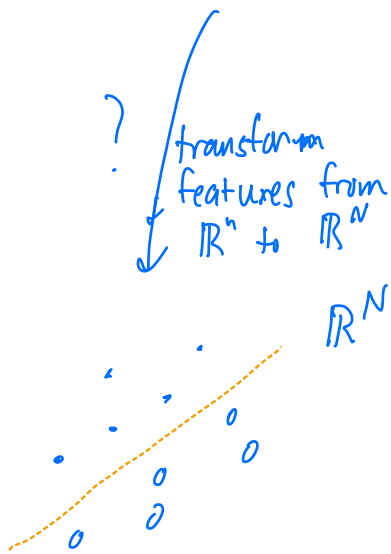
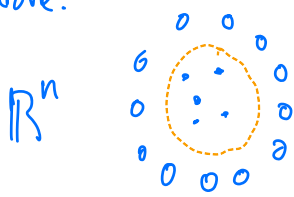
Why "support vector"?



Note: optimization result is unaffected by $\underline{x}_i, i \neq 1, 2$.

$\underline{x}_1, \underline{x}_2$ entirely define hyperplane $\rightarrow \underline{x}_i, \underline{x}_2$ called
"support vectors".

This procedure employs linear separability. This can't always be done.



"Modern" versions of SVM employ

"nonlinear" separators via

the "kernel trick": idea is to replace $\underline{a}^T \underline{x}$ by $k(\underline{a}, \underline{x})$.

(e.g. $k(\underline{x}, \underline{y}) = \underline{x}^T \underline{y} \rightarrow$ linear SVM

$$k(\underline{x}, \underline{y}) = 1 - \exp(\underline{x}^T \underline{y})$$

$$k(\underline{x}, \underline{y}) = 1 - \exp(-\|\underline{x} - \underline{y}\|_2^2)$$

Another practical augmentation: $y_i (\underline{a}^T \underline{x}_i + b) \geq 1 \quad i=1 \dots M$

is a "hard" cutoff condition.

Terminology: "hard margin"

This is unforgiving and the feasible set can be empty.

Instead: soften problem: hard requirement is $y_i (\underline{a}^T \underline{x}_i + b) \geq 1$

Alternative: $\min 1 - y_i(a^T x_i + b)$ when $1 - y_i(a^T x_i + b) \geq 0$.
measure of violation of $y_i(a^T x_i + b) \geq 1$.

Soft margin optimization:

$$\min_{a, b} \|a\|_2^2 + \lambda \sum_{i=1}^m \max\{0, 1 - y_i(a^T x_i + b)\}$$

λ : regularization parameter balancing margin maximization against soft violations.

$$(\lambda > 0), \lambda \sim \frac{1}{m}$$

This is still convex.

→ this soft margin SVM optimization is Tikhonov regularization in disguise.

$$\min: \sum_{i=1}^m \max\{0, 1 - y_i(a^T x_i + b)\} + \frac{1}{\lambda} \|a\|_2^2$$

Tikhonov regularization.

— in machine learning, minimizing $\sum_{i=1}^m \max\{0, 1 - y_i(a^T x_i + b)\}$ is called minimizing the empirical risk.