

Support vector machines

CS 446 / ECE 449

2022-02-16 19:35:44 -0600 (0670b06)

Another algorithm for linear prediction. Why?!

“pytorch meta-algorithm”.

⋮

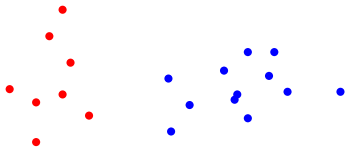
5. Tweak 1-4 until training error is small.
6. Tweak 1-5, possibly reducing model complexity, until testing error is small.

Support vector machines (SVMs) have three purposes for us.

1. Demonstrate **maximum margin predictors**, an example of “low complexity models”, which appear throughout machine learning (not just linear predictors).
2. Demonstrate **nonlinear kernels**, also pervasive.
3. Exercise convex optimization and duality.

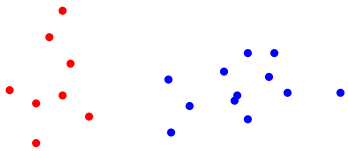
- ▶ Hard-margin SVM.
- ▶ Soft-margin SVM.
- ▶ SVM duality.
- ▶ Nonlinear SVM: kernels

Maximum margin linear separators

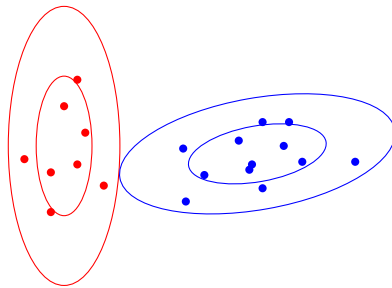


Which linear separator is best?

Maximum margin linear separators

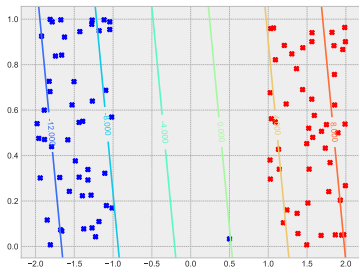
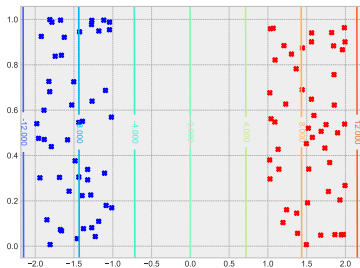


Which linear separator is best?



Which linear separator is best?

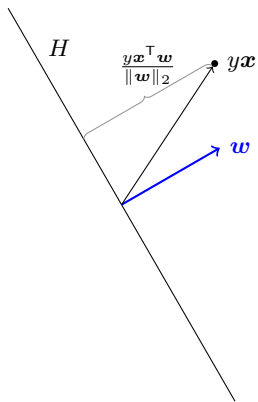
The max margin separator is **one** choice for a good predictor.
It is not always the best idea. Recall from lecture 3:



Even so, the maximum margin concept is pervasive in machine learning.

How to write “maximum margin classifier”?

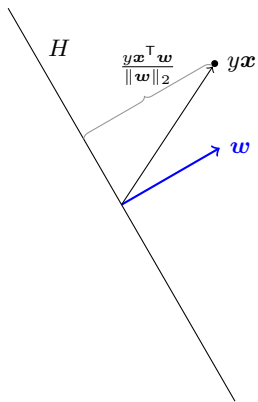
Input $\mathbf{x} \in \mathbb{R}^d$, label $y \in \{\pm 1\}$, predictor $\mathbf{w} \in \mathbb{R}^d$ with $H := \{\mathbf{z} \in \mathbb{R}^d : \mathbf{w}^\top \mathbf{z} = 0\}$.



How to write “maximum margin classifier”?

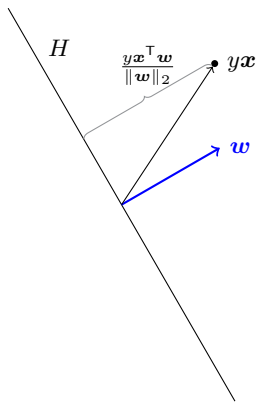
Input $\mathbf{x} \in \mathbb{R}^d$, label $y \in \{\pm 1\}$, predictor $\mathbf{w} \in \mathbb{R}^d$ with $H := \{\mathbf{z} \in \mathbb{R}^d : \mathbf{w}^\top \mathbf{z} = 0\}$.

Formulation #1: constrain the norm.



How to write “maximum margin classifier”?

Input $\mathbf{x} \in \mathbb{R}^d$, label $y \in \{\pm 1\}$, predictor $\mathbf{w} \in \mathbb{R}^d$ with $H := \{\mathbf{z} \in \mathbb{R}^d : \mathbf{w}^\top \mathbf{z} = 0\}$.



Formulation #1: constrain the norm.

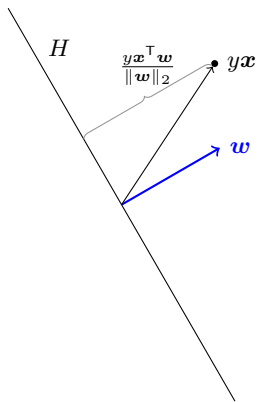
Single margin $\frac{y\mathbf{x}^\top \mathbf{w}}{\|\mathbf{w}\|}$.

Overall margin $\min_i \frac{y_i \mathbf{x}_i^\top \mathbf{w}}{\|\mathbf{w}\|}$.

Max margin $\max_{\|\mathbf{u}\|=1} \min_i y_i \mathbf{x}_i^\top \mathbf{u}$.

How to write “maximum margin classifier”?

Input $\mathbf{x} \in \mathbb{R}^d$, label $y \in \{\pm 1\}$, predictor $\mathbf{w} \in \mathbb{R}^d$ with $H := \{\mathbf{z} \in \mathbb{R}^d : \mathbf{w}^\top \mathbf{z} = 0\}$.



Formulation #1: constrain the norm.

Single margin $\frac{y\mathbf{x}^\top \mathbf{w}}{\|\mathbf{w}\|}$.

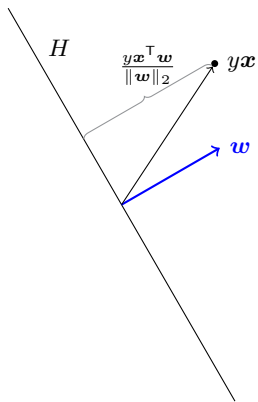
Overall margin $\min_i \frac{y_i \mathbf{x}_i^\top \mathbf{w}}{\|\mathbf{w}\|}$.

Max margin $\max_{\|\mathbf{u}\|=1} \min_i y_i \mathbf{x}_i^\top \mathbf{u}$.

Formulation #2: constrain the margins.

How to write “maximum margin classifier”?

Input $\mathbf{x} \in \mathbb{R}^d$, label $y \in \{\pm 1\}$, predictor $\mathbf{w} \in \mathbb{R}^d$ with $H := \{\mathbf{z} \in \mathbb{R}^d : \mathbf{w}^\top \mathbf{z} = 0\}$.



Formulation #1: constrain the norm.

Single margin $\frac{y\mathbf{x}^\top \mathbf{w}}{\|\mathbf{w}\|}$.

Overall margin $\min_i \frac{y_i \mathbf{x}_i^\top \mathbf{w}}{\|\mathbf{w}\|}$.

Max margin $\max_{\|\mathbf{w}\|=1} \min_i y_i \mathbf{x}_i^\top \mathbf{w}$.

Formulation #2: constrain the margins.

Consider any \mathbf{v} with $\min_i y_i \mathbf{x}_i^\top \mathbf{v} \geq 1$.

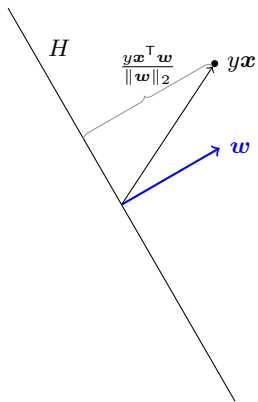
Can suppose $y_k \mathbf{x}_k^\top \mathbf{v} = 1$ for some k (why?).

Since margin scales with $\frac{1}{\|\mathbf{v}\|}$, choose

$$\begin{aligned} \min \quad & \frac{1}{2} \|\mathbf{v}\|^2 \\ \text{subject to} \quad & \mathbf{v} \in \mathbb{R}^d, \\ & y_i \mathbf{x}_i^\top \mathbf{v} \geq 1 \quad \forall i. \end{aligned}$$

How to write “maximum margin classifier”?

Input $\mathbf{x} \in \mathbb{R}^d$, label $y \in \{\pm 1\}$, predictor $\mathbf{w} \in \mathbb{R}^d$ with $H := \{\mathbf{z} \in \mathbb{R}^d : \mathbf{w}^\top \mathbf{z} = 0\}$.



Formulation #1: constrain the norm.

Single margin $\frac{y\mathbf{x}^\top \mathbf{w}}{\|\mathbf{w}\|}$.

Overall margin $\min_i \frac{y_i \mathbf{x}_i^\top \mathbf{w}}{\|\mathbf{w}\|}$.

Max margin $\max_{\|\mathbf{u}\|=1} \min_i y_i \mathbf{x}_i^\top \mathbf{u}$.

Formulation #2: constrain the margins.

Consider any \mathbf{v} with $\min_i y_i \mathbf{x}_i^\top \mathbf{v} \geq 1$.

Can suppose $y_k \mathbf{x}_k^\top \mathbf{v} = 1$ for some k (why?).

Since margin scales with $\frac{1}{\|\mathbf{v}\|}$, choose

$$\begin{aligned} \min \quad & \frac{1}{2} \|\mathbf{v}\|^2 \\ \text{subject to} \quad & \mathbf{v} \in \mathbb{R}^d, \\ & y_i \mathbf{x}_i^\top \mathbf{v} \geq 1 \quad \forall i. \end{aligned}$$

These two are equivalent
(up to scaling, when solutions exist).

Hard-margin SVM.

Take the solution to either optimization problem:

$$\begin{aligned} & \max \quad \min_i y_i \mathbf{x}_i^\top \mathbf{u}, \\ \text{subject to} \quad & \mathbf{u} \in \mathbb{R}^d, \\ & \|\mathbf{u}\| = 1; \end{aligned}$$

$$\begin{aligned} & \min \quad \frac{1}{2} \|\mathbf{v}\|^2 \\ \text{subject to} \quad & \mathbf{v} \in \mathbb{R}^d, \\ & y_i \mathbf{x}_i^\top \mathbf{v} \geq 1 \quad \forall i. \end{aligned}$$

Hard-margin SVM.

Take the solution to either optimization problem:

$$\begin{array}{ll} \max & \min_i y_i \mathbf{x}_i^\top \mathbf{u}, \\ \text{subject to} & \mathbf{u} \in \mathbb{R}^d, \\ & \|\mathbf{u}\| = 1; \end{array} \qquad \begin{array}{ll} \min & \frac{1}{2} \|\mathbf{v}\|^2 \\ \text{subject to} & \mathbf{v} \in \mathbb{R}^d, \\ & y_i \mathbf{x}_i^\top \mathbf{v} \geq 1 \quad \forall i. \end{array}$$

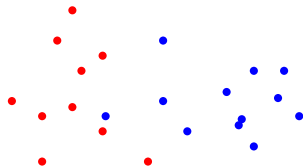
Remarks.

- ▶ Since the second is a **convex program**, many approaches exist. Homework will investigate a simple SGD-based strategy.
- ▶ What happens if the second formulation is **infeasible**? (That is, what if no vector \mathbf{v} satisfies $\min_i y_i \mathbf{x}_i^\top \mathbf{v} \geq 1$?)

Soft-margin SVM

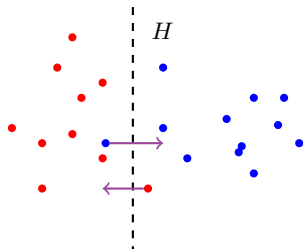
Soft-margin SVM

What is the max margin predictor for the following data?



Soft-margin SVM

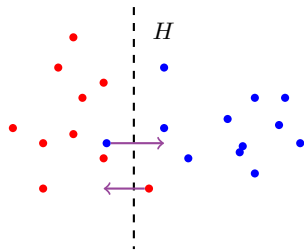
What is the max margin predictor for the following data?



Idea. pay a price for each $y_i \mathbf{x}_i^\top \mathbf{v} < 1$ with slack variables $(\xi_i)_{i=1}^n$:

Soft-margin SVM

What is the max margin predictor for the following data?



Idea. pay a price for each $y_i \mathbf{x}_i^\top \mathbf{v} < 1$ with slack variables $(\xi_i)_{i=1}^n$:

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^d, \xi_1, \dots, \xi_n \in \mathbb{R}} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_i \mathbf{x}_i^\top \mathbf{w} \geq 1 - \xi_i && \text{for all } i = 1, 2, \dots, n, \\ & \xi_i \geq 0 && \text{for all } i = 1, 2, \dots, n. \end{aligned}$$

Re-formulation as regularized ERM.

Formulation with slack variables $(\xi_i)_{i=1}^n$.

$$\min_{\mathbf{w} \in \mathbb{R}^d, \xi_1, \dots, \xi_n \in \mathbb{R}} \quad \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \xi_i$$

$$\text{s.t.} \quad y_i \mathbf{x}_i^\top \mathbf{w} \geq 1 - \xi_i$$

$$\xi_i \geq 0$$

for all $i = 1, 2, \dots, n$,

for all $i = 1, 2, \dots, n$.

Re-formulation as regularized ERM.

Formulation with slack variables $(\xi_i)_{i=1}^n$.

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^d, \xi_1, \dots, \xi_n \in \mathbb{R}} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_i \mathbf{x}_i^\top \mathbf{w} \geq 1 - \xi_i && \text{for all } i = 1, 2, \dots, n, \\ & \xi_i \geq 0 && \text{for all } i = 1, 2, \dots, n. \end{aligned}$$

Regularized ERM formulation.

Given any \mathbf{w} , choose $\xi_i := \max\{0, 1 - y_i \mathbf{x}_i^\top \mathbf{w}\}$, whereby

$$\min_{\mathbf{w} \in \mathbb{R}^d} \quad \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \max\{0, 1 - y_i \mathbf{x}_i^\top \mathbf{w}\},$$

where $\ell_{\text{hinge}}(y_i \mathbf{x}_i^\top \mathbf{w}) := \max\{0, 1 - y_i \mathbf{x}_i^\top \mathbf{w}\}$ is the **hinge loss**.

Re-formulation as regularized ERM.

Formulation with slack variables $(\xi_i)_{i=1}^n$.

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^d, \xi_1, \dots, \xi_n \in \mathbb{R}} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_i \mathbf{x}_i^\top \mathbf{w} \geq 1 - \xi_i && \text{for all } i = 1, 2, \dots, n, \\ & \xi_i \geq 0 && \text{for all } i = 1, 2, \dots, n. \end{aligned}$$

Regularized ERM formulation.

Given any \mathbf{w} , choose $\xi_i := \max\{0, 1 - y_i \mathbf{x}_i^\top \mathbf{w}\}$, whereby

$$\min_{\mathbf{w} \in \mathbb{R}^d} \quad \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \max\{0, 1 - y_i \mathbf{x}_i^\top \mathbf{w}\},$$

where $\ell_{\text{hinge}}(y_i \mathbf{x}_i^\top \mathbf{w}) := \max\{0, 1 - y_i \mathbf{x}_i^\top \mathbf{w}\}$ is the **hinge loss**.

Remarks.

- ▶ Normally we'd write $\frac{1}{n} \sum_{i=1}^n \ell_{\text{hinge}}(y_i \mathbf{x}_i^\top \mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|^2$.
- ▶ C (or λ) is a **hyper-parameter**; it has no good search procedure.

A convex program is an optimization problem (minimization or maximization) where a convex objective is minimized over a convex constraint (feasible) set.

A convex program is an optimization problem (minimization or maximization) where a convex objective is minimized over a convex constraint (feasible) set.

Every convex program has a corresponding **dual program**.

For the SVM, the dual has many nice properties:

- ▶ Clarifies the role of support vectors.
- ▶ Leads to a nice nonlinear approach via kernels.
- ▶ Gives another choice for optimization algorithms.

SVM hard-margin duality.

Define the two optimization problems

$$\min \left\{ \frac{1}{2} \|\mathbf{w}\|^2 : \mathbf{w} \in \mathbb{R}^d, \forall i. 1 - y_i \mathbf{x}_i^\top \mathbf{w} \leq 0 \right\} \quad (\text{primal}),$$

$$\max \left\{ \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j : \alpha \in \mathbb{R}^n, \alpha \geq 0 \right\} \quad (\text{dual}).$$

If the primal is feasible,

then the primal optimal value equals the dual optimal value.

Given a primal optimum $\bar{\mathbf{w}}$ and a dual optimum $\bar{\alpha}$, they satisfy

$$\bar{\mathbf{w}} = \sum_{i=1}^n \bar{\alpha}_i y_i \mathbf{x}_i.$$

SVM hard-margin duality.

Define the two optimization problems

$$\min \left\{ \frac{1}{2} \|\mathbf{w}\|^2 : \mathbf{w} \in \mathbb{R}^d, \forall i. 1 - y_i \mathbf{x}_i^\top \mathbf{w} \leq 0 \right\} \quad (\text{primal}),$$

$$\max \left\{ \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j : \alpha \in \mathbb{R}^n, \alpha \geq 0 \right\} \quad (\text{dual}).$$

If the primal is feasible,

then the primal optimal value equals the dual optimal value.

Given a primal optimum $\bar{\mathbf{w}}$ and a dual optimum $\bar{\alpha}$, they satisfy

$$\bar{\mathbf{w}} = \sum_{i=1}^n \bar{\alpha}_i y_i \mathbf{x}_i.$$

- ▶ The **dual variables** α have dimension n , same as examples.
- ▶ We can write the primal optimum as a linear combination of examples.
- ▶ The dual objective is a **concave quadratic**.
- ▶ We will derive this duality using Lagrange multipliers.

Lagrange multipliers

Move constraints to objective using Lagrange multipliers.

$$\begin{aligned} \text{Original problem: } \quad & \min_{\mathbf{w} \in \mathbb{R}^d} \quad \frac{1}{2} \|\mathbf{w}\|_2^2 \\ & \text{s.t.} \quad 1 - y_i \mathbf{x}_i^\top \mathbf{w} \leq 0 \quad \text{for all } i = 1, \dots, n. \end{aligned}$$

- ▶ For each constraint $1 - y_i \mathbf{x}_i^\top \mathbf{w} \leq 0$, associate a dual variable (Lagrange multiplier) $\alpha_i \geq 0$.
- ▶ Move constraints to objective by adding $\sum_{i=1}^n \alpha_i (1 - y_i \mathbf{x}_i^\top \mathbf{w})$ and maximizing over $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_n) \in \mathbb{R}^n$ s.t. $\boldsymbol{\alpha} \geq \mathbf{0}$ (i.e., $\alpha_i \geq 0$ for all i).

Lagrange multipliers

Move constraints to objective using Lagrange multipliers.

$$\begin{aligned} \text{Original problem: } \quad & \min_{\mathbf{w} \in \mathbb{R}^d} \quad \frac{1}{2} \|\mathbf{w}\|_2^2 \\ & \text{s.t.} \quad 1 - y_i \mathbf{x}_i^\top \mathbf{w} \leq 0 \quad \text{for all } i = 1, \dots, n. \end{aligned}$$

- ▶ For each constraint $1 - y_i \mathbf{x}_i^\top \mathbf{w} \leq 0$, associate a dual variable (Lagrange multiplier) $\alpha_i \geq 0$.
- ▶ Move constraints to objective by adding $\sum_{i=1}^n \alpha_i (1 - y_i \mathbf{x}_i^\top \mathbf{w})$ and maximizing over $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_n) \in \mathbb{R}^n$ s.t. $\boldsymbol{\alpha} \geq \mathbf{0}$ (i.e., $\alpha_i \geq 0$ for all i).

Lagrangian $L(\mathbf{w}, \boldsymbol{\alpha})$:

$$L(\mathbf{w}, \boldsymbol{\alpha}) := \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y_i \mathbf{x}_i^\top \mathbf{w}).$$

Maximizing over $\boldsymbol{\alpha} \geq \mathbf{0}$ recovers primal problem: for any $\mathbf{w} \in \mathbb{R}^d$,

$$P(\mathbf{w}) := \sup_{\boldsymbol{\alpha} \geq \mathbf{0}} L(\mathbf{w}, \boldsymbol{\alpha}) = \begin{cases} \frac{1}{2} \|\mathbf{w}\|_2^2 & \text{if } \min_i y_i \mathbf{x}_i^\top \mathbf{w} \geq 1, \\ \infty & \text{otherwise.} \end{cases}$$

Lagrange multipliers

Move constraints to objective using Lagrange multipliers.

$$\begin{aligned} \text{Original problem: } \quad & \min_{\mathbf{w} \in \mathbb{R}^d} \quad \frac{1}{2} \|\mathbf{w}\|_2^2 \\ & \text{s.t.} \quad 1 - y_i \mathbf{x}_i^\top \mathbf{w} \leq 0 \quad \text{for all } i = 1, \dots, n. \end{aligned}$$

- ▶ For each constraint $1 - y_i \mathbf{x}_i^\top \mathbf{w} \leq 0$, associate a dual variable (Lagrange multiplier) $\alpha_i \geq 0$.
- ▶ Move constraints to objective by adding $\sum_{i=1}^n \alpha_i (1 - y_i \mathbf{x}_i^\top \mathbf{w})$ and maximizing over $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_n) \in \mathbb{R}^n$ s.t. $\boldsymbol{\alpha} \geq \mathbf{0}$ (i.e., $\alpha_i \geq 0$ for all i).

Lagrangian $L(\mathbf{w}, \boldsymbol{\alpha})$:

$$L(\mathbf{w}, \boldsymbol{\alpha}) := \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y_i \mathbf{x}_i^\top \mathbf{w}).$$

Maximizing over $\boldsymbol{\alpha} \geq 0$ recovers primal problem: for any $\mathbf{w} \in \mathbb{R}^d$,

$$P(\mathbf{w}) := \sup_{\boldsymbol{\alpha} \geq 0} L(\mathbf{w}, \boldsymbol{\alpha}) = \begin{cases} \frac{1}{2} \|\mathbf{w}\|_2^2 & \text{if } \min_i y_i \mathbf{x}_i^\top \mathbf{w} \geq 1, \\ \infty & \text{otherwise.} \end{cases}$$

What if we leave $\boldsymbol{\alpha}$ fixed, and minimize \mathbf{w} ?

Dual problem

Lagrangian

$$L(\mathbf{w}, \boldsymbol{\alpha}) := \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y_i \mathbf{x}_i^\top \mathbf{w}).$$

Primal hard-margin SVM

$$P(\mathbf{w}) = \sup_{\boldsymbol{\alpha} \geq \mathbf{0}} L(\mathbf{w}, \boldsymbol{\alpha}) = \sup_{\boldsymbol{\alpha} \geq \mathbf{0}} \left[\frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y_i \mathbf{x}_i^\top \mathbf{w}) \right].$$

Dual problem

Lagrangian

$$L(\mathbf{w}, \boldsymbol{\alpha}) := \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y_i \mathbf{x}_i^\top \mathbf{w}).$$

Primal hard-margin SVM

$$P(\mathbf{w}) = \sup_{\boldsymbol{\alpha} \geq \mathbf{0}} L(\mathbf{w}, \boldsymbol{\alpha}) = \sup_{\boldsymbol{\alpha} \geq \mathbf{0}} \left[\frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y_i \mathbf{x}_i^\top \mathbf{w}) \right].$$

Dual problem $D(\boldsymbol{\alpha}) = \min_{\mathbf{w}} L(\mathbf{w}, \boldsymbol{\alpha})$:

Dual problem

Lagrangian

$$L(\mathbf{w}, \boldsymbol{\alpha}) := \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \alpha_i (1 - \mathbf{y}_i \mathbf{x}_i^\top \mathbf{w}).$$

Primal hard-margin SVM

$$P(\mathbf{w}) = \sup_{\boldsymbol{\alpha} \geq \mathbf{0}} L(\mathbf{w}, \boldsymbol{\alpha}) = \sup_{\boldsymbol{\alpha} \geq \mathbf{0}} \left[\frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \alpha_i (1 - \mathbf{y}_i \mathbf{x}_i^\top \mathbf{w}) \right].$$

Dual problem $D(\boldsymbol{\alpha}) = \min_{\mathbf{w}} L(\mathbf{w}, \boldsymbol{\alpha})$: given $\boldsymbol{\alpha} \geq \mathbf{0}$, then $\mathbf{w} \mapsto L(\mathbf{w}, \boldsymbol{\alpha})$ is a convex quadratic with minimum $\mathbf{w} = \sum_{i=1}^n \alpha_i \mathbf{y}_i \mathbf{x}_i$, giving

$$\begin{aligned} D(\boldsymbol{\alpha}) &= \min_{\mathbf{w} \in \mathbb{R}^d} L(\mathbf{w}, \boldsymbol{\alpha}) = L\left(\sum_{i=1}^n \alpha_i \mathbf{y}_i \mathbf{x}_i, \boldsymbol{\alpha}\right) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \left\| \sum_{i=1}^n \alpha_i \mathbf{y}_i \mathbf{x}_i \right\|_2^2 \\ &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j \mathbf{y}_i \mathbf{y}_j \mathbf{x}_i^\top \mathbf{x}_j. \end{aligned}$$

Summarizing,

$$L(\mathbf{w}, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y_i \mathbf{x}_i^\top \mathbf{w})$$

$$P(\mathbf{w}) = \max_{\boldsymbol{\alpha} \geq 0} L(\mathbf{w}, \boldsymbol{\alpha})$$

$$D(\boldsymbol{\alpha}) = \min_{\mathbf{w}} L(\mathbf{w}, \boldsymbol{\alpha})$$

Lagrangian,

primal problem,

dual problem.

Summarizing,

$$L(\mathbf{w}, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y_i \mathbf{x}_i^\top \mathbf{w})$$

Lagrangian,

$$P(\mathbf{w}) = \max_{\boldsymbol{\alpha} \geq 0} L(\mathbf{w}, \boldsymbol{\alpha})$$

primal problem,

$$D(\boldsymbol{\alpha}) = \min_{\mathbf{w}} L(\mathbf{w}, \boldsymbol{\alpha})$$

dual problem.

- For general Lagrangians, have weak duality

$$P(\mathbf{w}) \geq D(\boldsymbol{\alpha}),$$

since $P(\mathbf{w}) = \max_{\boldsymbol{\alpha}' \geq 0} L(\mathbf{w}, \boldsymbol{\alpha}') \geq L(\mathbf{w}, \boldsymbol{\alpha}) \geq \min_{\mathbf{w}'} L(\mathbf{w}', \boldsymbol{\alpha}) = D(\boldsymbol{\alpha})$.

Summarizing,

$$L(\mathbf{w}, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y_i \mathbf{x}_i^\top \mathbf{w})$$

Lagrangian,

$$P(\mathbf{w}) = \max_{\boldsymbol{\alpha} \geq 0} L(\mathbf{w}, \boldsymbol{\alpha})$$

primal problem,

$$D(\boldsymbol{\alpha}) = \min_{\mathbf{w}} L(\mathbf{w}, \boldsymbol{\alpha})$$

dual problem.

- ▶ For general Lagrangians, have weak duality

$$P(\mathbf{w}) \geq D(\boldsymbol{\alpha}),$$

since $P(\mathbf{w}) = \max_{\boldsymbol{\alpha}' \geq 0} L(\mathbf{w}, \boldsymbol{\alpha}') \geq L(\mathbf{w}, \boldsymbol{\alpha}) \geq \min_{\mathbf{w}'} L(\mathbf{w}', \boldsymbol{\alpha}) = D(\boldsymbol{\alpha})$.

- ▶ By **convexity**, have **strong duality** $\min_{\mathbf{w}} P(\mathbf{w}) = \max_{\boldsymbol{\alpha} \geq 0} D(\boldsymbol{\alpha})$,
and an optimum $\bar{\boldsymbol{\alpha}}$ for D gives an optimum $\bar{\mathbf{w}}$ for P via

$$\bar{\mathbf{w}} = \sum_{i=1}^n \bar{\alpha}_i y_i \mathbf{x}_i = \arg \min_{\mathbf{w}} L(\mathbf{w}, \bar{\boldsymbol{\alpha}}).$$

Support vectors

Optimal solutions $\bar{\mathbf{w}}$ and $\bar{\boldsymbol{\alpha}} = (\bar{\alpha}_1, \dots, \bar{\alpha}_n)$ satisfy

$$\blacktriangleright \bar{\mathbf{w}} = \sum_{i=1}^n \bar{\alpha}_i y_i \mathbf{x}_i = \sum_{i:\bar{\alpha}_i > 0} \bar{\alpha}_i y_i \mathbf{x}_i,$$

$\blacktriangleright \bar{\alpha}_i > 0 \Rightarrow y_i \mathbf{x}_i^T \bar{\mathbf{w}} = 1$ for all $i = 1, \dots, n$ (complementary slackness).

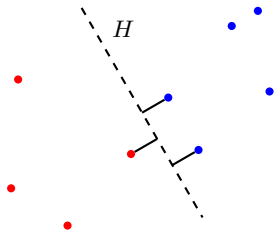
Support vectors

Optimal solutions $\bar{\mathbf{w}}$ and $\bar{\alpha} = (\bar{\alpha}_1, \dots, \bar{\alpha}_n)$ satisfy

$$\blacktriangleright \bar{\mathbf{w}} = \sum_{i=1}^n \bar{\alpha}_i y_i \mathbf{x}_i = \sum_{i:\bar{\alpha}_i > 0} \bar{\alpha}_i y_i \mathbf{x}_i,$$

$\blacktriangleright \bar{\alpha}_i > 0 \Rightarrow y_i \mathbf{x}_i^T \bar{\mathbf{w}} = 1$ for all $i = 1, \dots, n$ (complementary slackness).

The $y_i \mathbf{x}_i$ where $\bar{\alpha}_i > 0$ are called support vectors.



Support vectors

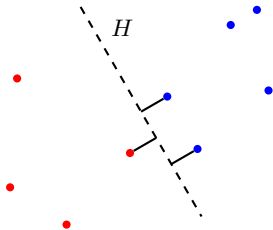
Optimal solutions $\bar{\mathbf{w}}$ and $\bar{\alpha} = (\bar{\alpha}_1, \dots, \bar{\alpha}_n)$ satisfy

$$\blacktriangleright \bar{\mathbf{w}} = \sum_{i=1}^n \bar{\alpha}_i y_i \mathbf{x}_i = \sum_{i:\bar{\alpha}_i > 0} \bar{\alpha}_i y_i \mathbf{x}_i,$$

$\blacktriangleright \bar{\alpha}_i > 0 \Rightarrow y_i \mathbf{x}_i^T \bar{\mathbf{w}} = 1$ for all $i = 1, \dots, n$ (complementary slackness).

The $y_i \mathbf{x}_i$ where $\bar{\alpha}_i > 0$ are called support vectors.

\blacktriangleright Support vector examples satisfy “margin” constraints with equality.



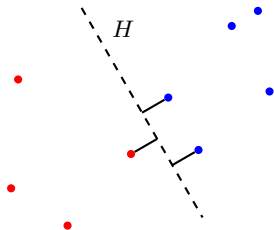
Support vectors

Optimal solutions $\bar{\mathbf{w}}$ and $\bar{\alpha} = (\bar{\alpha}_1, \dots, \bar{\alpha}_n)$ satisfy

$$\blacktriangleright \bar{\mathbf{w}} = \sum_{i=1}^n \bar{\alpha}_i y_i \mathbf{x}_i = \sum_{i:\bar{\alpha}_i > 0} \bar{\alpha}_i y_i \mathbf{x}_i,$$

$\blacktriangleright \bar{\alpha}_i > 0 \Rightarrow y_i \mathbf{x}_i^T \bar{\mathbf{w}} = 1$ for all $i = 1, \dots, n$ (complementary slackness).

The $y_i \mathbf{x}_i$ where $\bar{\alpha}_i > 0$ are called support vectors.



- \blacktriangleright Support vector examples satisfy “margin” constraints with equality.
- \blacktriangleright Get same solution if non-support vectors deleted.

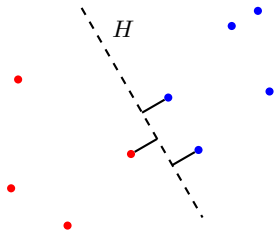
Support vectors

Optimal solutions $\bar{\mathbf{w}}$ and $\bar{\boldsymbol{\alpha}} = (\bar{\alpha}_1, \dots, \bar{\alpha}_n)$ satisfy

$$\blacktriangleright \bar{\mathbf{w}} = \sum_{i=1}^n \bar{\alpha}_i y_i \mathbf{x}_i = \sum_{i:\bar{\alpha}_i > 0} \bar{\alpha}_i y_i \mathbf{x}_i,$$

$\blacktriangleright \bar{\alpha}_i > 0 \Rightarrow y_i \mathbf{x}_i^T \bar{\mathbf{w}} = 1$ for all $i = 1, \dots, n$ (complementary slackness).

The $y_i \mathbf{x}_i$ where $\bar{\alpha}_i > 0$ are called support vectors.



- \blacktriangleright Support vector examples satisfy “margin” constraints with equality.
- \blacktriangleright Get same solution if non-support vectors deleted.
- \blacktriangleright Primal optimum is a linear combination of support vectors.

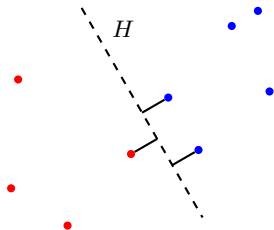
Support vectors

Optimal solutions $\bar{\mathbf{w}}$ and $\bar{\boldsymbol{\alpha}} = (\bar{\alpha}_1, \dots, \bar{\alpha}_n)$ satisfy

$$\blacktriangleright \bar{\mathbf{w}} = \sum_{i=1}^n \bar{\alpha}_i y_i \mathbf{x}_i = \sum_{i:\bar{\alpha}_i > 0} \bar{\alpha}_i y_i \mathbf{x}_i,$$

$\blacktriangleright \bar{\alpha}_i > 0 \Rightarrow y_i \mathbf{x}_i^T \bar{\mathbf{w}} = 1$ for all $i = 1, \dots, n$ (complementary slackness).

The $y_i \mathbf{x}_i$ where $\bar{\alpha}_i > 0$ are called support vectors.



- \blacktriangleright Support vector examples satisfy “margin” constraints with equality.
- \blacktriangleright Get same solution if non-support vectors deleted.
- \blacktriangleright Primal optimum is a linear combination of support vectors.
- \blacktriangleright Dual solution and support vectors not necessarily unique! (Why not?)

Proof of complementary slackness

For the optimal (feasible) solutions $\bar{\mathbf{w}}$ and $\bar{\boldsymbol{\alpha}}$, we have

$$P(\bar{\mathbf{w}}) = D(\bar{\boldsymbol{\alpha}}) = \min_{\mathbf{w} \in \mathbb{R}^d} L(\mathbf{w}, \bar{\boldsymbol{\alpha}}) \quad (\text{by strong duality})$$

Proof of complementary slackness

For the optimal (feasible) solutions $\bar{\mathbf{w}}$ and $\bar{\boldsymbol{\alpha}}$, we have

$$\begin{aligned} P(\bar{\mathbf{w}}) = D(\bar{\boldsymbol{\alpha}}) &= \min_{\mathbf{w} \in \mathbb{R}^d} L(\mathbf{w}, \bar{\boldsymbol{\alpha}}) && \text{(by strong duality)} \\ &\leq L(\bar{\mathbf{w}}, \bar{\boldsymbol{\alpha}}) \end{aligned}$$

Proof of complementary slackness

For the optimal (feasible) solutions $\bar{\mathbf{w}}$ and $\bar{\alpha}$, we have

$$\begin{aligned} P(\bar{\mathbf{w}}) = D(\bar{\alpha}) &= \min_{\mathbf{w} \in \mathbb{R}^d} L(\mathbf{w}, \bar{\alpha}) \quad (\text{by strong duality}) \\ &\leq L(\bar{\mathbf{w}}, \bar{\alpha}) \\ &= \frac{1}{2} \|\bar{\mathbf{w}}\|_2^2 + \sum_{i=1}^n \bar{\alpha}_i (1 - y_i \mathbf{x}_i^\top \bar{\mathbf{w}}) \end{aligned}$$

Proof of complementary slackness

For the optimal (feasible) solutions $\bar{\mathbf{w}}$ and $\bar{\alpha}$, we have

$$\begin{aligned} P(\bar{\mathbf{w}}) = D(\bar{\alpha}) &= \min_{\mathbf{w} \in \mathbb{R}^d} L(\mathbf{w}, \bar{\alpha}) \quad (\text{by strong duality}) \\ &\leq L(\bar{\mathbf{w}}, \bar{\alpha}) \\ &= \frac{1}{2} \|\bar{\mathbf{w}}\|_2^2 + \sum_{i=1}^n \bar{\alpha}_i (1 - y_i \mathbf{x}_i^\top \bar{\mathbf{w}}) \\ &\leq \frac{1}{2} \|\bar{\mathbf{w}}\|_2^2 \quad (\text{constraints are satisfied}) \\ &= P(\bar{\mathbf{w}}). \end{aligned}$$

Proof of complementary slackness

For the optimal (feasible) solutions $\bar{\mathbf{w}}$ and $\bar{\alpha}$, we have

$$\begin{aligned} P(\bar{\mathbf{w}}) = D(\bar{\alpha}) &= \min_{\mathbf{w} \in \mathbb{R}^d} L(\mathbf{w}, \bar{\alpha}) \quad (\text{by strong duality}) \\ &\leq L(\bar{\mathbf{w}}, \bar{\alpha}) \\ &= \frac{1}{2} \|\bar{\mathbf{w}}\|_2^2 + \sum_{i=1}^n \bar{\alpha}_i (1 - y_i \mathbf{x}_i^\top \bar{\mathbf{w}}) \\ &\leq \frac{1}{2} \|\bar{\mathbf{w}}\|_2^2 \quad (\text{constraints are satisfied}) \\ &= P(\bar{\mathbf{w}}). \end{aligned}$$

Therefore, every term in sum $\sum_{i=1}^n \bar{\alpha}_i (1 - y_i \mathbf{x}_i^\top \bar{\mathbf{w}})$ must be zero:

$$\bar{\alpha}_i (1 - y_i \mathbf{x}_i^\top \bar{\mathbf{w}}) = 0 \quad \text{for all } i = 1, \dots, n.$$

Proof of complementary slackness

For the optimal (feasible) solutions $\bar{\mathbf{w}}$ and $\bar{\alpha}$, we have

$$\begin{aligned} P(\bar{\mathbf{w}}) = D(\bar{\alpha}) &= \min_{\mathbf{w} \in \mathbb{R}^d} L(\mathbf{w}, \bar{\alpha}) \quad (\text{by strong duality}) \\ &\leq L(\bar{\mathbf{w}}, \bar{\alpha}) \\ &= \frac{1}{2} \|\bar{\mathbf{w}}\|_2^2 + \sum_{i=1}^n \bar{\alpha}_i (1 - y_i \mathbf{x}_i^\top \bar{\mathbf{w}}) \\ &\leq \frac{1}{2} \|\bar{\mathbf{w}}\|_2^2 \quad (\text{constraints are satisfied}) \\ &= P(\bar{\mathbf{w}}). \end{aligned}$$

Therefore, every term in sum $\sum_{i=1}^n \bar{\alpha}_i (1 - y_i \mathbf{x}_i^\top \bar{\mathbf{w}})$ must be zero:

$$\bar{\alpha}_i (1 - y_i \mathbf{x}_i^\top \bar{\mathbf{w}}) = 0 \quad \text{for all } i = 1, \dots, n.$$

If $\bar{\alpha}_i > 0$, then must have $1 - y_i \mathbf{x}_i^\top \bar{\mathbf{w}} = 0$. (Not iff!)

SVM (hard-margin) duality summary

Lagrangian

$$L(\mathbf{w}, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y_i \mathbf{x}_i^\top \mathbf{w}).$$

Primal maximum margin problem was

$$P(\mathbf{w}) = \sup_{\boldsymbol{\alpha} \geq \mathbf{0}} L(\mathbf{w}, \boldsymbol{\alpha}) = \sup_{\boldsymbol{\alpha} \geq \mathbf{0}} \left[\frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y_i \mathbf{x}_i^\top \mathbf{w}) \right].$$

Dual problem

$$D(\boldsymbol{\alpha}) = \min_{\mathbf{w} \in \mathbb{R}^d} L(\mathbf{w}, \boldsymbol{\alpha}) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \left\| \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \right\|_2^2.$$

Given dual optimum $\bar{\boldsymbol{\alpha}}$,

- ▶ Corresponding primal optimum $\bar{\mathbf{w}} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i$;
- ▶ Strong duality $P(\bar{\mathbf{w}}) = D(\bar{\boldsymbol{\alpha}})$;
- ▶ $\bar{\alpha}_i > 0$ implies $y_i \mathbf{x}_i^\top \bar{\mathbf{w}} = 1$,
and these $y_i \mathbf{x}_i$ are support vectors.

SVM soft-margin dual

SVM soft-margin dual

Similarly,

$$L(\mathbf{w}, \boldsymbol{\xi}, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \xi_i + \sum_{i=1}^n \alpha_i (1 - \xi_i - y_i \mathbf{x}_i^T \mathbf{w}) \quad (\text{Lagrangian}),$$

$$P(\mathbf{w}, \boldsymbol{\xi}) = \sup_{\boldsymbol{\alpha} \geq 0} L(\mathbf{w}, \boldsymbol{\xi}, \boldsymbol{\alpha}) \quad (\text{Primal}),$$

$$= \begin{cases} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \xi_i & \forall i, 1 - \xi_i - y_i \mathbf{x}_i^T \mathbf{w} \leq 0, \\ \infty & \text{otherwise,} \end{cases}$$

$$D(\boldsymbol{\alpha}) = \min_{\mathbf{w} \in \mathbb{R}^d, \boldsymbol{\xi} \in \mathbb{R}_{\geq 0}^n} L(\mathbf{w}, \boldsymbol{\xi}, \boldsymbol{\alpha}) \quad (\text{Dual}),$$

$$= \max_{\substack{\boldsymbol{\alpha} \in \mathbb{R}^n \\ 0 \leq \alpha_i \leq C}} \left[\sum_{i=1}^n \alpha_i - \frac{1}{2} \left\| \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \right\|^2 \right].$$

SVM soft-margin dual

Similarly,

$$L(\mathbf{w}, \boldsymbol{\xi}, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \xi_i + \sum_{i=1}^n \alpha_i (1 - \xi_i - y_i \mathbf{x}_i^\top \mathbf{w}) \quad (\text{Lagrangian}),$$

$$P(\mathbf{w}, \boldsymbol{\xi}) = \sup_{\boldsymbol{\alpha} \geq 0} L(\mathbf{w}, \boldsymbol{\xi}, \boldsymbol{\alpha}) \quad (\text{Primal}),$$

$$= \begin{cases} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \xi_i & \forall i, 1 - \xi_i - y_i \mathbf{x}_i^\top \mathbf{w} \leq 0, \\ \infty & \text{otherwise,} \end{cases}$$

$$D(\boldsymbol{\alpha}) = \min_{\mathbf{w} \in \mathbb{R}^d, \boldsymbol{\xi} \in \mathbb{R}_{\geq 0}^n} L(\mathbf{w}, \boldsymbol{\xi}, \boldsymbol{\alpha}) \quad (\text{Dual}),$$

$$= \max_{\substack{\boldsymbol{\alpha} \in \mathbb{R}^n \\ 0 \leq \alpha_i \leq C}} \left[\sum_{i=1}^n \alpha_i - \frac{1}{2} \left\| \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \right\|^2 \right].$$

Remarks.

- ▶ Dual solution $\bar{\boldsymbol{\alpha}}$ still gives primal solution $\bar{\mathbf{w}} = \sum_{i=1}^n \bar{\alpha}_i y_i \mathbf{x}_i$.
- ▶ Can take $C \rightarrow \infty$ to recover hard-margin case.
- ▶ Dual is still a constrained concave quadratic (used in many solvers).
- ▶ Some presentations include bias in primal ($\mathbf{x}_i^\top \mathbf{w} + b$); this introduces a constraint $\sum_{i=1}^n y_i \alpha_i = 0$ in dual.

Nonlinear SVM: feature mapping annoying in the primal?

Nonlinear SVM: feature mapping annoying in the primal?

SVM hard-margin primal, with a feature mapping $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^p$:

$$\min \left\{ \frac{1}{2} \|\mathbf{w}\|^2 : \mathbf{w} \in \mathbb{R}^p, \forall i. \phi(\mathbf{x}_i)^\top \mathbf{w} \geq 1 \right\}.$$

Now the search space has p dimensions; potentially $p \gg d$.

Can we do better?

Feature mapping in the dual

Feature mapping in the dual

SVM hard-margin dual, with a feature mapping $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^p$:

$$\max_{\alpha_1, \alpha_2, \dots, \alpha_n \geq 0} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j).$$

Given dual optimum $\bar{\alpha}$, since $\bar{\mathbf{w}} = \sum_{i=1}^n \bar{\alpha}_i y_i \phi(\mathbf{x}_i)$, we can predict on future \mathbf{x} with

$$\mathbf{x} \mapsto \phi(\mathbf{x})^\top \bar{\mathbf{w}} = \sum_{i=1}^n \bar{\alpha}_i y_i \phi(\mathbf{x})^\top \phi(\mathbf{x}_i).$$

Feature mapping in the dual

SVM hard-margin dual, with a feature mapping $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^p$:

$$\max_{\alpha_1, \alpha_2, \dots, \alpha_n \geq 0} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j).$$

Given dual optimum $\bar{\alpha}$, since $\bar{\mathbf{w}} = \sum_{i=1}^n \bar{\alpha}_i y_i \phi(\mathbf{x}_i)$, we can predict on future \mathbf{x} with

$$\mathbf{x} \mapsto \phi(\mathbf{x})^\top \bar{\mathbf{w}} = \sum_{i=1}^n \bar{\alpha}_i y_i \phi(\mathbf{x})^\top \phi(\mathbf{x}_i).$$

- ▶ Dual form never needs $\phi(\mathbf{x}) \in \mathbb{R}^p$, only $\phi(\mathbf{x})^\top \phi(\mathbf{x}_i) \in \mathbb{R}$.
- ▶ **Kernel trick:** replace every $\phi(\mathbf{x})^\top \phi(\mathbf{x}')$ with **kernel evaluation** $k(\mathbf{x}, \mathbf{x}')$. Sometimes $k(\cdot, \cdot)$ is much cheaper than $\phi(\mathbf{x})^\top \phi(\mathbf{x}')$.
- ▶ This idea started with SVM, but appears in many other places.
- ▶ **Downside:** implementations usually store **Gram matrix** $G \in \mathbb{R}^{n \times n}$ where $G_{ij} := k(\mathbf{x}_i, \mathbf{x}_j)$.

Kernel example: affine features

Affine features: $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^{1+d}$, where

$$\phi(\mathbf{x}) = (1, x_1, \dots, x_d).$$

Kernel form:

$$\phi(\mathbf{x})^\top \phi(\mathbf{x}') = 1 + \mathbf{x}^\top \mathbf{x}'.$$

Kernel example: quadratic features

Consider re-normalized quadratic features

$\phi: \mathbb{R}^d \rightarrow \mathbb{R}^{1+2d+\binom{d}{2}}$, where

$$\phi(\mathbf{x}) = \left(1, \sqrt{2}x_1, \dots, \sqrt{2}x_d, x_1^2, \dots, x_d^2, \right. \\ \left. \sqrt{2}x_1x_2, \dots, \sqrt{2}x_1x_d, \dots, \sqrt{2}x_{d-1}x_d \right).$$

Just writing this down takes time $\mathcal{O}(d^2)$.

Meanwhile,

$$\phi(\mathbf{x})^\top \phi(\mathbf{x}') = (1 + \mathbf{x}^\top \mathbf{x}')^2,$$

time $\mathcal{O}(d)$.

Kernel example: quadratic features

Consider re-normalized quadratic features

$\phi: \mathbb{R}^d \rightarrow \mathbb{R}^{1+2d+\binom{d}{2}}$, where

$$\phi(\mathbf{x}) = \left(1, \sqrt{2}x_1, \dots, \sqrt{2}x_d, x_1^2, \dots, x_d^2, \right. \\ \left. \sqrt{2}x_1x_2, \dots, \sqrt{2}x_1x_d, \dots, \sqrt{2}x_{d-1}x_d \right).$$

Just writing this down takes time $\mathcal{O}(d^2)$.

Meanwhile,

$$\phi(\mathbf{x})^\top \phi(\mathbf{x}') = (1 + \mathbf{x}^\top \mathbf{x}')^2,$$

time $\mathcal{O}(d)$.

Tweaks:

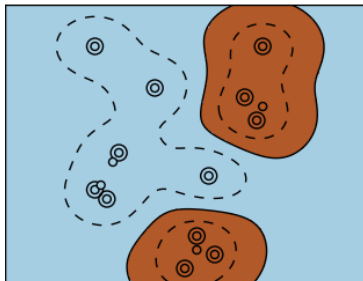
- ▶ What if we change exponent “2”?
- ▶ What if we replace additive “1” with 0?

RBF kernel

For any $\sigma > 0$, there is an infinite feature expansion $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^\infty$ such that

$$\phi(\mathbf{x})^\top \phi(\mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\sigma^2}\right),$$

which can be computed in $O(d)$ time.



This is called a Gaussian kernel or RBF kernel.

It has some similarities to nearest neighbor methods (later lecture).

ϕ maps to an infinite-dimensional space, but there's no reason to know that.

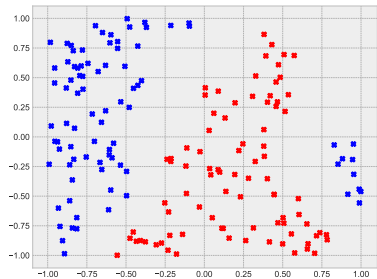
Defining kernels without ϕ

A (positive definite) **kernel function** $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is a symmetric function so that for any n and any data examples $(\mathbf{x}_i)_{i=1}^n$, the corresponding Gram matrix $\mathbf{G} \in \mathbb{R}^{n \times n}$ with $G_{ij} := k(\mathbf{x}_i, \mathbf{x}_j)$ is positive semi-definite.

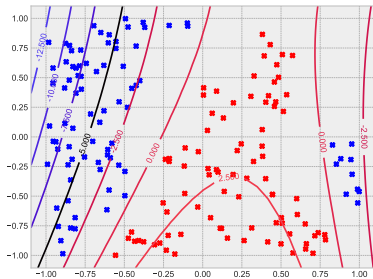
Defining kernels without ϕ

A (positive definite) **kernel function** $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is a symmetric function so that for any n and any data examples $(\mathbf{x}_i)_{i=1}^n$, the corresponding Gram matrix $\mathbf{G} \in \mathbb{R}^{n \times n}$ with $G_{ij} := k(\mathbf{x}_i, \mathbf{x}_j)$ is positive semi-definite.

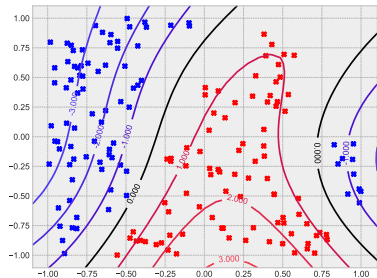
- ▶ There is a ton of theory about this formalism; e.g., keywords RKHS, representer theorem, Mercer's theorem.
- ▶ Given any such k , there always exists a corresponding ϕ .
- ▶ This definition ensures the SVM dual is still concave.



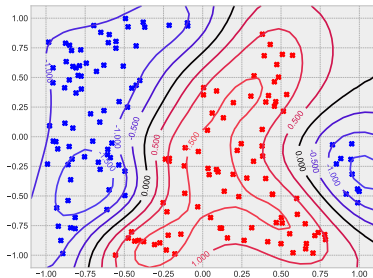
Source data.



Quadratic SVM.



RBF SVM ($\sigma = 1$).



RBF SVM ($\sigma = 0.1$).

Summary for SVM

- ▶ Hard-margin SVM.
- ▶ Soft-margin SVM.
- ▶ SVM duality.
- ▶ Nonlinear SVM: kernels

(Appendix.)

Equivalence of two hard-margin formulations.

Proof. First note that both have unique solutions (when feasible). For the first formulation, suppose we have two solutions \mathbf{u} and \mathbf{u}' , and define another vector $\mathbf{u}'' := (\mathbf{u} + \mathbf{u}')/2$. Then \mathbf{u}'' achieves the same margin value as \mathbf{u} and \mathbf{u}' , but if $\mathbf{u} \neq \mathbf{u}'$, then $\|\mathbf{u}''\| < 1$, which means $\mathbf{u}''/\|\mathbf{u}''\|$ achieves a larger margin value than the purported optima \mathbf{u} and \mathbf{u}' , a contradiction. For the second formulation, it suffices to note that the objective is strictly convex.

Now consider solution \mathbf{u} to the first, with margin γ . Then $\mathbf{v} := \mathbf{u}/\gamma$ is feasible for second, with optimal value $1/(2\gamma^2)$. So the optimal value is at most this; if it is exactly the optimal value, we are done, otherwise suppose the optimum $\bar{\mathbf{v}}$ has $\|\bar{\mathbf{v}}\|^2/2 = 1/(2\rho^2) < 1/(2\gamma^2)$. Then $\bar{\mathbf{u}} := \rho\bar{\mathbf{v}}$ is a unit vector, and moreover has $\min_i y_i \mathbf{x}_i^\top \bar{\mathbf{u}} = \rho > \gamma$, a contradiction since this is better than the supposed optimum \mathbf{u} .

□

Soft-margin dual derivation

Let's derive the final dual form:

$$L(\mathbf{w}, \boldsymbol{\xi}, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \xi_i + \sum_{i=1}^n \alpha_i (1 - \xi_i - \mathbf{y}_i \mathbf{x}_i^\top \mathbf{w})$$

$$D(\boldsymbol{\alpha}) = \min_{\mathbf{w} \in \mathbb{R}^d, \boldsymbol{\xi} \in \mathbb{R}_{\geq 0}^n} L(\mathbf{w}, \boldsymbol{\xi}, \boldsymbol{\alpha}) = \max_{\substack{\boldsymbol{\alpha} \in \mathbb{R}^n \\ 0 \leq \alpha_i \leq C}} \left[\sum_{i=1}^n \alpha_i - \frac{1}{2} \left\| \sum_{i=1}^n \alpha_i \mathbf{y}_i \mathbf{x}_i \right\|^2 \right].$$

Given $\boldsymbol{\alpha}$ and $\boldsymbol{\xi}$, the minimizing \mathbf{w} is still $\mathbf{w} = \sum_{i=1}^n \alpha_i \mathbf{y}_i \mathbf{x}_i$; plugging in,

$$D(\boldsymbol{\alpha}) = \min_{\boldsymbol{\xi} \in \mathbb{R}_{\geq 0}^n} L \left(\sum_{i=1}^n \alpha_i \mathbf{y}_i \mathbf{x}_i, \boldsymbol{\xi}, \boldsymbol{\alpha} \right) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \left\| \sum_{i=1}^n \alpha_i \mathbf{y}_i \mathbf{x}_i \right\|^2 + \sum_{i=1}^n \xi_i (C - \alpha_i).$$

The goal is to maximize D ; if $\alpha_i > C$, then $\xi_i \uparrow \infty$ gives $D(\boldsymbol{\alpha}) = -\infty$.

Otherwise, minimized at $\xi_i = 0$. Therefore the dual problem is

$$\max_{\substack{\boldsymbol{\alpha} \in \mathbb{R}^n \\ 0 \leq \alpha_i \leq C}} \left[\sum_{i=1}^n \alpha_i - \frac{1}{2} \left\| \sum_{i=1}^n \alpha_i \mathbf{y}_i \mathbf{x}_i \right\|^2 \right].$$

Gaussian kernel feature expansion

First consider $d = 1$, meaning $\phi: \mathbb{R} \rightarrow \mathbb{R}^\infty$.

What ϕ has $\phi(x)\phi(y) = e^{-(x-y)^2/(2\sigma^2)}$?

Gaussian kernel feature expansion

First consider $d = 1$, meaning $\phi: \mathbb{R} \rightarrow \mathbb{R}^\infty$.

What ϕ has $\phi(x)\phi(y) = e^{-(x-y)^2/(2\sigma^2)}$?

Reverse engineer using Taylor expansion:

$$e^{-(x-y)^2/(2\sigma^2)} = e^{-x^2/(2\sigma^2)} \cdot e^{-y^2/(2\sigma^2)} \cdot e^{xy/\sigma^2}$$

Gaussian kernel feature expansion

First consider $d = 1$, meaning $\phi: \mathbb{R} \rightarrow \mathbb{R}^\infty$.

What ϕ has $\phi(x)\phi(y) = e^{-(x-y)^2/(2\sigma^2)}$?

Reverse engineer using Taylor expansion:

$$\begin{aligned} e^{-(x-y)^2/(2\sigma^2)} &= e^{-x^2/(2\sigma^2)} \cdot e^{-y^2/(2\sigma^2)} \cdot e^{xy/\sigma^2} \\ &= e^{-x^2/(2\sigma^2)} \cdot e^{-y^2/(2\sigma^2)} \cdot \sum_{k=0}^{\infty} \frac{1}{k!} \left(\frac{xy}{\sigma^2}\right)^k \end{aligned}$$

Gaussian kernel feature expansion

First consider $d = 1$, meaning $\phi: \mathbb{R} \rightarrow \mathbb{R}^\infty$.

What ϕ has $\phi(x)\phi(y) = e^{-(x-y)^2/(2\sigma^2)}$?

Reverse engineer using Taylor expansion:

$$\begin{aligned} e^{-(x-y)^2/(2\sigma^2)} &= e^{-x^2/(2\sigma^2)} \cdot e^{-y^2/(2\sigma^2)} \cdot e^{xy/\sigma^2} \\ &= e^{-x^2/(2\sigma^2)} \cdot e^{-y^2/(2\sigma^2)} \cdot \sum_{k=0}^{\infty} \frac{1}{k!} \left(\frac{xy}{\sigma^2}\right)^k \end{aligned}$$

So let

$$\phi(x) := e^{-x^2/(2\sigma^2)} \left(1, \frac{x}{\sigma}, \frac{1}{\sqrt{2!}} \left(\frac{x}{\sigma}\right)^2, \frac{1}{\sqrt{3!}} \left(\frac{x}{\sigma}\right)^3, \dots \right).$$

Gaussian kernel feature expansion

First consider $d = 1$, meaning $\phi: \mathbb{R} \rightarrow \mathbb{R}^\infty$.

What ϕ has $\phi(x)\phi(y) = e^{-(x-y)^2/(2\sigma^2)}$?

Reverse engineer using Taylor expansion:

$$\begin{aligned} e^{-(x-y)^2/(2\sigma^2)} &= e^{-x^2/(2\sigma^2)} \cdot e^{-y^2/(2\sigma^2)} \cdot e^{xy/\sigma^2} \\ &= e^{-x^2/(2\sigma^2)} \cdot e^{-y^2/(2\sigma^2)} \cdot \sum_{k=0}^{\infty} \frac{1}{k!} \left(\frac{xy}{\sigma^2}\right)^k \end{aligned}$$

So let

$$\phi(x) := e^{-x^2/(2\sigma^2)} \left(1, \frac{x}{\sigma}, \frac{1}{\sqrt{2!}} \left(\frac{x}{\sigma}\right)^2, \frac{1}{\sqrt{3!}} \left(\frac{x}{\sigma}\right)^3, \dots \right).$$

How to handle $d > 1$?

Gaussian kernel feature expansion

First consider $d = 1$, meaning $\phi: \mathbb{R} \rightarrow \mathbb{R}^\infty$.

What ϕ has $\phi(x)\phi(y) = e^{-(x-y)^2/(2\sigma^2)}$?

Reverse engineer using Taylor expansion:

$$\begin{aligned} e^{-(x-y)^2/(2\sigma^2)} &= e^{-x^2/(2\sigma^2)} \cdot e^{-y^2/(2\sigma^2)} \cdot e^{xy/\sigma^2} \\ &= e^{-x^2/(2\sigma^2)} \cdot e^{-y^2/(2\sigma^2)} \cdot \sum_{k=0}^{\infty} \frac{1}{k!} \left(\frac{xy}{\sigma^2}\right)^k \end{aligned}$$

So let

$$\phi(x) := e^{-x^2/(2\sigma^2)} \left(1, \frac{x}{\sigma}, \frac{1}{\sqrt{2!}} \left(\frac{x}{\sigma}\right)^2, \frac{1}{\sqrt{3!}} \left(\frac{x}{\sigma}\right)^3, \dots \right).$$

How to handle $d > 1$?

$$\begin{aligned} e^{-\|\mathbf{x}-\mathbf{y}\|^2/(2\sigma^2)} &= e^{-\|\mathbf{x}\|^2/(2\sigma^2)} \cdot e^{-\|\mathbf{y}\|^2/(2\sigma^2)} \cdot e^{\mathbf{x}^\top \mathbf{y}/\sigma^2} \\ &= e^{-\|\mathbf{x}\|^2/(2\sigma^2)} \cdot e^{-\|\mathbf{y}\|^2/(2\sigma^2)} \cdot \sum_{k=0}^{\infty} \frac{1}{k!} \left(\frac{\mathbf{x}^\top \mathbf{y}}{\sigma^2}\right)^k. \end{aligned}$$

Kernel example: products of all subsets of coordinates

Consider $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^{2^d}$, where

$$\phi(\mathbf{x}) = \left(\prod_{i \in S} x_i \right)_{S \subseteq \{1, 2, \dots, d\}}$$

Time $\mathcal{O}(2^d)$ just to write down.

Kernel evaluation takes time $\mathcal{O}(d)$:

$$\phi(\mathbf{x})^\top \phi(\mathbf{x}') = \prod_{i=1}^d (1 + x_i x'_i).$$

Suppose k_1 and k_2 are positive definite kernel functions.

Suppose k_1 and k_2 are positive definite kernel functions.

1. $k(x, y) := k_1(x, y) + k_2(x, y)$ define a positive definite kernel?

Suppose k_1 and k_2 are positive definite kernel functions.

1. $k(x, y) := k_1(x, y) + k_2(x, y)$ define a positive definite kernel?
2. $k(x, y) := c \cdot k_1(x, y)$ (for $c \geq 0$) define a positive definite kernel?

Suppose k_1 and k_2 are positive definite kernel functions.

1. $k(x, y) := k_1(x, y) + k_2(x, y)$ define a positive definite kernel?
2. $k(x, y) := c \cdot k_1(x, y)$ (for $c \geq 0$) define a positive definite kernel?
3. $k(x, y) := k_1(x, y) \cdot k_2(x, y)$ define a positive definite kernel?

Suppose k_1 and k_2 are positive definite kernel functions.

1. $k(x, y) := k_1(x, y) + k_2(x, y)$ define a positive definite kernel?
2. $k(x, y) := c \cdot k_1(x, y)$ (for $c \geq 0$) define a positive definite kernel?
3. $k(x, y) := k_1(x, y) \cdot k_2(x, y)$ define a positive definite kernel?

Another approach: random features $k(\mathbf{x}, \mathbf{x}') = \mathbb{E}_{\mathbf{w}} F(\mathbf{w}, \mathbf{x})^\top F(\mathbf{w}, \mathbf{x}')$ for some F ; we will revisit this with deep networks and the [neural tangent kernel \(NTK\)](#).

Kernel ridge regression

Kernel ridge regression:

$$\min_{\mathbf{w}} \frac{1}{2n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2.$$

Solution:

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X} + \lambda n \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}.$$

Linear algebra fact:

$$\mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top + \lambda n \mathbf{I})^{-1} \mathbf{y}.$$

Therefore predict with

$$\mathbf{x} \mapsto \mathbf{x}^\top \hat{\mathbf{w}} = (\mathbf{X}\mathbf{x})^\top (\mathbf{X}\mathbf{X}^\top + \lambda n \mathbf{I})^{-1} \mathbf{y} = \sum_{i=1}^n (\mathbf{x}_i^\top \mathbf{x})^\top \left[(\mathbf{X}\mathbf{X}^\top + \lambda n \mathbf{I})^{-1} \mathbf{y} \right]_i.$$

Kernel approach:

- ▶ Compute $\boldsymbol{\alpha} := (\mathbf{G} + \lambda n \mathbf{I})^{-1}$, where $\mathbf{G} \in \mathbb{R}^n$ is the Gram matrix: $G_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$.
- ▶ Predict with $\mathbf{x} \mapsto \sum_{i=1}^n \alpha_i \mathbf{x}_i^\top \mathbf{x}$.

There are a few ways;
one is to use one-against-all as in lecture.

Many researchers have proposed various multiclass SVM methods,
but some of them can be shown to fail in trivial cases.

New two-forms explanation idea **INCOMPLETE**

- ▶ All that matters is the direction; so for each direction which is a strict separator (meaning $\min_i y_i \mathbf{x}_i^T \mathbf{v} > 0$), we will pick one vector, and then compare the various directions. The final step is to write this as a convex (or concave) program.
- ▶ First approach: in each strict separation direction, normalize it so that $\|\mathbf{v}\| = 1$. If $\|\mathbf{v}\| = 1$, then the margin on the training set is $\min_i y_i \mathbf{x}_i^T \mathbf{v}$. We want the maximum amongst all these directions, which is $\max_{\|\mathbf{v}\|=1, \min_i y_i \mathbf{x}_i^T \mathbf{v} > 0} \min_i y_i \mathbf{x}_i^T \mathbf{v}$. Firstly we can drop the strict separation constraint since the objective enforces it. The first constraint however is not a convex set, so we don't have a concave program yet, but note that we have the same solutions if we instead write $\max_{\|\mathbf{v}\| \leq 1} \min_i y_i \mathbf{x}_i^T \mathbf{v}$. To see this, note that a vector \mathbf{u} with $\|\mathbf{u}\| < 1$ can not be a maximizer, since $\mathbf{u}/\|\mathbf{u}\|$ achieves a $1/\|\mathbf{u}\| > 1$ factor larger margin thus the solution to both problems lies on the boundary.
- ▶ Second approach: in each strict separation direction, normalize it so that $\min_i y_i \mathbf{x}_i^T \mathbf{v} = 1$. Since the distance to the separator is $y \mathbf{x}^T \mathbf{v} / \|\mathbf{v}\|$, then for a minimizing example k with $y_k \mathbf{x}_k^T \mathbf{v} = 1$, it follows that the margin is $y_k \mathbf{x}_k^T \mathbf{v} / \|\mathbf{v}\| = 1/\|\mathbf{v}\|$, and therefore the overall margin after this normalization is $1/\|\mathbf{v}\|$. We want to maximize the margin, so we could maximize $1/\|\mathbf{v}\|$, but this is a little awkward for optimization algorithms, so we instead minimize the reciprocal, $\|\mathbf{v}\|$. This too is a little awkward, so we square it, which preserves the optima, and minimize $\|\mathbf{v}\|^2/2$.
- ▶ Lastly, we can't write the program as $\min \{\|\mathbf{v}\|^2/2 : \forall i, y_i \mathbf{x}_i^T \mathbf{v} > 0\}$, because the infimal value is 0, which is not attained despite the set being bounded, so this is not a well-posed formulation for the maximum margin direction.

- ▶ Shalev-Shwartz/Ben-David: chapter 15.
- ▶ Murphy: chapter 14.