

CSCI-567: Machine Learning (Fall 2019)

Prof. Victor Adamchik

U of Southern California

Oct. 22, 2019

October 22, 2019 1 / 55

Outline

- 1 Support vector machines (primal formulation)
- 2 A detour: Linear Programming
- 3 A detour: Lagrangian duality
- 4 Support vector machines (dual formulation)

October 22, 2019 3 / 55

Outline

- 1 Support vector machines (primal formulation)
- 2 A detour: Linear Programming
- 3 A detour: Lagrangian duality
- 4 Support vector machines (dual formulation)

October 22, 2019 2 / 55

Support vector machines (SVM)

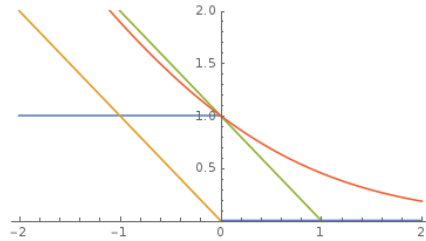
- One of the most commonly used classification algorithms
- Works well with the kernel trick
- Strong theoretical guarantees

We focus on **binary classification** here.

October 22, 2019 4 / 55

Primal formulation

In one sentence: linear model with L2 regularized hinge loss. Recall



- **perceptron loss** $\ell_{\text{perceptron}}(z) = \max\{0, -z\} \rightarrow$ Perceptron
- **logistic loss** $\ell_{\text{logistic}}(z) = \log(1 + \exp(-z)) \rightarrow$ logistic regression
- **hinge loss** $\ell_{\text{hinge}}(z) = \max\{0, 1 - z\} \rightarrow$ **SVM**

October 22, 2019 5 / 55

Primal formulation

For a linear model (\mathbf{w}, b) , this means

$$\min_{\mathbf{w}, b} \sum_n \max\{0, 1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b)\} + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

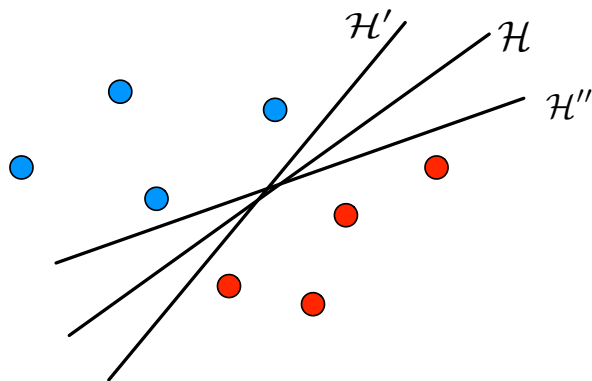
- recall $y_n \in \{-1, +1\}$
- a nonlinear mapping ϕ is applied
- the bias/intercept term b is used explicitly (think about why after this lecture)

So why L2 regularized hinge loss?

October 22, 2019 6 / 55

Geometric motivation: separable case

When data is **linearly separable**, there are *infinitely many hyperplanes with zero training error*.

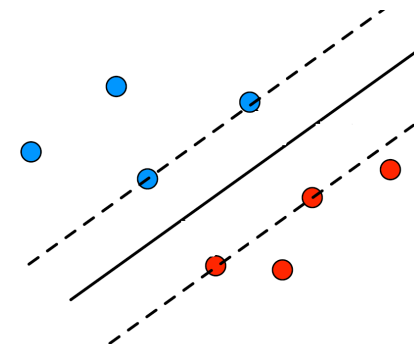


So which one should we choose?

October 22, 2019 7 / 55

Intuition

The further away from data points the better.



How to formalize this intuition?

October 22, 2019 8 / 55

Distance to hyperplane

What is the **distance** from a point x to a hyperplane $H : w^T x + b = 0$?

w is a normal vector perpendicular to H .

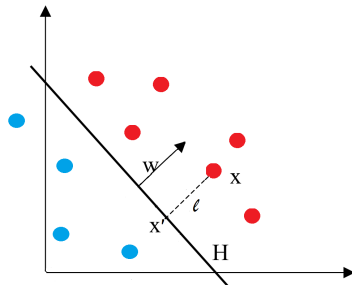
$x' \in H$ is the **projection** of x .

Then, $x' = x - \ell \frac{w}{\|w\|_2}$, we go ℓ units parallel to w .

Since x' belongs to a hyperplane, then

$$0 = w^T \left(x - \ell \frac{w}{\|w\|_2} \right) + b = w^T x - \ell \|w\| + b$$

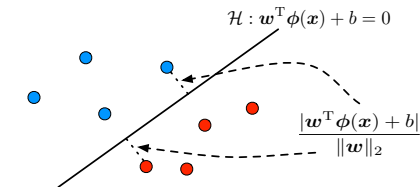
From this we find the distance $\ell = \frac{|w^T x + b|}{\|w\|_2}$.



Maximizing margin

Margin: the **smallest** distance from all training points to the hyperplane

$$\text{MARGIN OF } (w, b) = \min_n \frac{|w^T \phi(x_n) + b|}{\|w\|_2}$$



The intuition “**the further away the better**” translates to solving

$$\max_{w, b} \min_n \frac{|w^T \phi(x_n) + b|}{\|w\|_2}$$

October 22, 2019 9 / 55

October 22, 2019 10 / 55

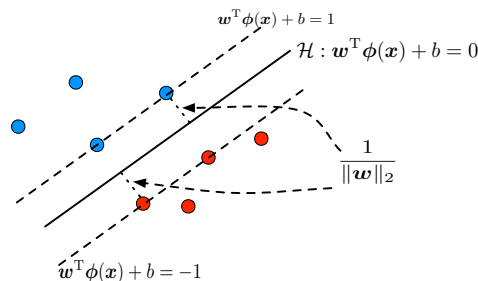
Rescaling

Note: rescaling (w, b) does not change the hyperplane at all.

We can thus always scale (w, b) s.t. $\min_n |w^T \phi(x_n) + b| = 1$

The margin then becomes

$$\begin{aligned} \text{MARGIN OF } (w, b) &= \min_n \frac{|w^T \phi(x_n) + b|}{\|w\|_2} \\ &= \frac{1}{\|w\|_2} \end{aligned}$$



Summary for separable data

Observe that $|w^T \phi(x_n) + b| = y_n(w^T \phi(x_n) + b)$.

Therefore, for a separable training set, we aim to solve the following optimization problem

$$\max_{w, b} \frac{1}{\|w\|_2} \quad \text{s.t.} \quad \min_n y_n(w^T \phi(x_n) + b) = 1$$

This is equivalent to

$$\min_{w, b} \frac{1}{2} \|w\|_2^2$$

subject to

$$y_n(w^T \phi(x_n) + b) \geq 1, \quad \forall n$$

SVM is thus also called **max-margin** classifier. The constraints above are called **hard-margin** constraints.

October 22, 2019 11 / 55

October 22, 2019 12 / 55

General non-separable case

If data is not linearly separable, the constraints

$$y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1, \quad \forall n$$

are obviously *not feasible*.

To deal with this issue, we relax them to **soft-margin** constraints:

$$y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1 - \xi_n, \quad \forall n$$

where we introduce **slack variables** $\xi_n \geq 0$.

We want ξ_n to be as small as possible.

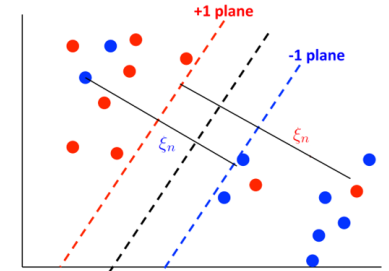
Meaning of slack variables ξ_n

The goal is to minimize the training errors (the number of misclassified points).

Instead we will minimize the distance between misclassified points and their correct hyperplane.

$0 < \xi_n \leq 1$ - data point falls within the margin on the correct side of the separating hyperplane; $\xi_n > 1$ - on the wrong side of the separating hyperplane.

We will introduce a hyperparameter C that represents a penalty for misclassifying points.



SVM Primal formulation

The objective function becomes

$$\min_{\mathbf{w}, b, \{\xi_n\}} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n$$

subject to

$$\begin{aligned} y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) &\geq 1 - \xi_n, \quad \forall n \\ \xi_n &\geq 0, \quad \forall n \end{aligned}$$

where C is a new hyperparameter.

This formulation is called the **soft-margin SVM**.

Optimization

$$\min_{\mathbf{w}, b, \{\xi_n\}} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n$$

subject to

$$\begin{aligned} y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) &\geq 1 - \xi_n, \quad \forall n \\ \xi_n &\geq 0, \quad \forall n \end{aligned}$$

- It is a convex (**quadratic** in fact) problem
- we can apply any convex optimization algorithms, e.g. SGD
- there are **more specialized and efficient** algorithms
- but usually we apply *kernel trick*, which requires solving the *dual problem*

Hinge Loss

How does this formulation

$$\min_{\mathbf{w}, b, \{\xi_n\}} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n$$

subject to

$$y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1 - \xi_n, \quad \forall n$$
$$\xi_n \geq 0, \quad \forall n$$

is related to L2 regularized hinge loss?

October 22, 2019 17 / 55

Equivalent form

Formulation

$$\min_{\mathbf{w}, b, \{\xi_n\}} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n$$

subject to

$$\xi_n \geq 1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b), \quad \forall n$$
$$\xi_n \geq 0, \quad \forall n$$

is equivalent to

$$\min_{\mathbf{w}, b, \{\xi_n\}} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n$$

subject to

$$\xi_n = \max \{0, 1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b)\}, \quad \forall n$$

October 22, 2019 18 / 55

Equivalent form

Formulation

$$\min_{\mathbf{w}, b, \{\xi_n\}} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n$$

subject to

$$\xi_n = \max \{0, 1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b)\}, \quad \forall n$$

is equivalent to

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \max \{0, 1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b)\}$$

This is exactly minimizing L2 regularized hinge loss!

October 22, 2019 19 / 55

Outline

- 1 Support vector machines (primal formulation)
- 2 A detour: Linear Programming
- 3 A detour: Lagrangian duality
- 4 Support vector machines (dual formulation)

October 22, 2019 20 / 55

Example Optimization Problem

Web server company wants to buy new servers.

Standard Model

- \$400
- 300W power
- Two shelves of rack
- Handles 1000 hits/min

Cutting-edge model

- \$1600
- 500W power
- One shelf
- 2000 hits/min

Budget:

- \$36,800
- 44 shelves of space
- 12,200W power

Goal: maximize the number of hits we serve per minute.

October 22, 2019 21 / 55

The approach: linear programming

- Introduce variables x_1 and x_2
(the number of servers of each model we buy)
- The number of hits per minute we get is:

$$1000x_1 + 2000x_2$$

- The budget places three limitations on us:
 - ▶ The financial budget:

$$400x_1 + 1600x_2 \leq 36800$$

- ▶ The number of shelves available:

$$2x_1 + x_2 \leq 44$$

- ▶ Power used collectively

$$300x_1 + 500x_2 \leq 12200$$

October 22, 2019 22 / 55

Summarize the optimization problem

$$\max_{x_1, x_2} 1000x_1 + 2000x_2$$

subject to:

$$400x_1 + 1600x_2 \leq 36800$$

$$2x_1 + x_2 \leq 44$$

$$300x_1 + 500x_2 \leq 12200$$

$$x_1, x_2 \geq 0$$

Various algorithms exist to solve the problem

October 22, 2019 23 / 55

Applications

Maximum Flow as a Linear Program

- Given a flow network with source, sink, edge capacities
- Flow through an edge must be at most capacity of edge.
- Flow into a vertex must equal flow out
(Exceptions: source, sink)

maximize: $\sum_{e \in \text{out}(s)} f_e$ where s is the source.

subject to: $0 \leq f_e \leq c_e$ for all edges e

$\sum_{e \in \text{in}(v)} f_e = \sum_{e \in \text{out}(v)} f_e$ for all vertices v
except the source and sink.

October 22, 2019 24 / 55

Standard form

A linear program is in **standard** form if it is in the following form:

$$\max_{x_n} \sum_n c_n x_n$$

subject to

$$\sum_n a_{mn} x_n \leq b_m, \quad \forall m$$
$$x_n \geq 0, \quad \forall n$$

We can write the standard form more compactly:

$$\max_x \mathbf{c}^T \mathbf{x}$$

subject to

$$\mathbf{A} \mathbf{x} \leq \mathbf{b} \text{ and } \mathbf{x} \geq 0$$

October 22, 2019 25 / 55

Duality

Primal (in x):

$$\begin{array}{ll} \text{maximize:} & \mathbf{c}^T \mathbf{x} \\ \text{subject to:} & \mathbf{A} \mathbf{x} \leq \mathbf{b} \\ & \mathbf{x} \geq \mathbf{0} \end{array}$$

Dual (in y):

$$\begin{array}{ll} \text{minimize:} & \mathbf{b}^T \mathbf{y} \\ \text{subject to:} & \mathbf{A}^T \mathbf{y} \geq \mathbf{c} \\ & \mathbf{y} \geq \mathbf{0} \end{array}$$

October 22, 2019 26 / 55

Weak Duality

Weak Duality: Let x be any feasible solution to the primal and y be any feasible solution for the dual. Then, $\mathbf{c}^T \mathbf{x} \leq \mathbf{b}^T \mathbf{y}$.

Proof.

$$\mathbf{c}^T \mathbf{x} = \mathbf{x}^T \mathbf{c} \leq \mathbf{x}^T (\mathbf{A}^T \mathbf{y}) = (\mathbf{A} \mathbf{x})^T \mathbf{y} \leq \mathbf{b}^T \mathbf{y}$$

The first inequality follows from the fact that y is feasible solution, the second inequality follows since x is feasible solution.

Recall a flow network: for any flow and any cut, $|f| \leq \text{cap}(A, B)$

October 22, 2019 27 / 55

Strong Duality

Strong Duality: Let x be any feasible solution to the primal and y be any feasible solution for the dual. Then, $\mathbf{c}^T \mathbf{x} = \mathbf{b}^T \mathbf{y}$.

Proof.

The proof of this theorem is beyond the scope of this lecture.

Recall the max-flow min-cut theorem from 570.

October 22, 2019 28 / 55

- 1 Support vector machines (primal formulation)
- 2 A detour: Linear Programming
- 3 A detour: Lagrangian duality
- 4 Support vector machines (dual formulation)

Primal problem

Suppose we want to solve

$$\min_{\mathbf{w}} F(\mathbf{w}) \quad \text{s.t.} \quad h_j(\mathbf{w}) \leq 0 \quad \forall j \in [J]$$

where functions h_1, \dots, h_J define J **constraints**.

SVM primal formulation is clearly of this form with $J = 2N$ constraints:

$$\begin{aligned} F(\mathbf{w}, b, \{\xi_n\}) &= C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 \\ h_n(\mathbf{w}, b, \{\xi_n\}) &= 1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) - \xi_n \quad \forall n \in [N] \\ h_{N+n}(\mathbf{w}, b, \{\xi_n\}) &= -\xi_n \quad \forall n \in [N] \end{aligned}$$

Lagrangian duality

Extremely important and powerful tool in analyzing optimizations

We will introduce basic concepts and derive the **KKT conditions**

Applying it to SVM reveals an important aspect of the algorithm

Lagrangian

Let us define the **Lagrangian** of the previous problem as:

$$L(\mathbf{w}, \{\lambda_j\}) = F(\mathbf{w}) + \sum_{j=1}^J \lambda_j h_j(\mathbf{w})$$

where $\lambda_1, \dots, \lambda_J \geq 0$ are new variables (called **Lagrangian multipliers**).

Note that

$$\max_{\{\lambda_j\} \geq 0} L(\mathbf{w}, \{\lambda_j\}) = \begin{cases} F(\mathbf{w}) & \text{if } h_j(\mathbf{w}) \leq 0 \quad \forall j \in [J] \\ +\infty & \text{else} \end{cases}$$

and thus,

$$\min_{\mathbf{w}} \max_{\{\lambda_j\} \geq 0} L(\mathbf{w}, \{\lambda_j\}) \iff \min_{\mathbf{w}} F(\mathbf{w}) \quad \text{s.t.} \quad h_j(\mathbf{w}) \leq 0 \quad \forall j \in [J]$$

Duality

We call this the **primal problem**

$$\min_{\mathbf{w}} \max_{\{\lambda_j\} \geq 0} L(\mathbf{w}, \{\lambda_j\})$$

We define the **dual problem** by swapping the min and max:

$$\max_{\{\lambda_j\} \geq 0} \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j\})$$

How are the primal and dual connected?

We will establish “**weak duality**” and “**strong duality**” for a non-linear optimization.

Strong duality

When F, h_1, \dots, h_m are convex, under some conditions (KKT conditions):

$$\min_{\mathbf{w}} \max_{\{\lambda_j\} \geq 0} L(\mathbf{w}, \{\lambda_j\}) = \max_{\{\lambda_j\} \geq 0} \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j\})$$

This is called “**strong duality**”.

We will derive those conditions in the next slides.

Weak Duality

Let \mathbf{w}^* and $\{\lambda_j^*\}$ be the primal and dual solutions respectively, then

$$\begin{aligned} \max_{\{\lambda_j\} \geq 0} \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j\}) &= \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j^*\}) \\ &\leq L(\mathbf{w}^*, \{\lambda_j^*\}) \\ &\leq \max_{\{\lambda_j\} \geq 0} L(\mathbf{w}^*, \{\lambda_j\}) \\ &= \min_{\mathbf{w}} \max_{\{\lambda_j\} \geq 0} L(\mathbf{w}, \{\lambda_j\}) \end{aligned}$$

This is called “**weak duality**”.

Deriving the Karush-Kuhn-Tucker (KKT) conditions

Observe that if strong duality holds:

$$\begin{aligned} F(\mathbf{w}^*) &= \min_{\mathbf{w}} \max_{\{\lambda_j\} \geq 0} L(\mathbf{w}, \{\lambda_j\}) = \max_{\{\lambda_j\} \geq 0} \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j\}) = \\ &= \min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j^*\}) \leq L(\mathbf{w}^*, \{\lambda_j^*\}) = F(\mathbf{w}^*) + \sum_{j=1}^J \lambda_j^* h_j(\mathbf{w}^*) \leq \\ &\leq F(\mathbf{w}^*) \end{aligned}$$

Implications:

- *all inequalities above have to be equalities!*
- last equality implies $\lambda_j^* h_j(\mathbf{w}^*) = 0$ for all $j \in [J]$
- equality $\min_{\mathbf{w}} L(\mathbf{w}, \{\lambda_j^*\}) = L(\mathbf{w}^*, \{\lambda_j^*\})$ implies \mathbf{w}^* is a **minimizer** of $L(\mathbf{w}, \{\lambda_j^*\})$ and thus has **zero gradient**.

The Karush-Kuhn-Tucker (KKT) conditions

If \mathbf{w}^* and $\{\lambda_j^*\}$ are the primal and dual solution respectively, then:

Stationarity:

$$\nabla_{\mathbf{w}} L(\mathbf{w}^*, \{\lambda_j^*\}) = \nabla F(\mathbf{w}^*) + \sum_{j=1}^J \lambda_j^* \nabla h_j(\mathbf{w}^*) = \mathbf{0}$$

Complementary slackness:

$$\lambda_j^* h_j(\mathbf{w}^*) = 0 \quad \text{for all } j \in [J]$$

Feasibility:

$$h_j(\mathbf{w}^*) \leq 0 \quad \text{and} \quad \lambda_j^* \geq 0 \quad \text{for all } j \in [J]$$

These are *necessary conditions*. They are also *sufficient* when F is convex and h_1, \dots, h_J are continuously differentiable convex functions.

October 22, 2019 37 / 55

Outline

- 1 Support vector machines (primal formulation)
- 2 A detour: Linear Programming
- 3 A detour: Lagrangian duality
- 4 Support vector machines (dual formulation)

October 22, 2019 38 / 55

Writing down the Lagrangian

Recall the primal formulation

$$\min_{\mathbf{w}, b, \{\xi_n\}} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n$$

subject to

$$\begin{aligned} 1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) - \xi_n &\leq 0, \quad \forall n \\ -\xi_n &\leq 0, \quad \forall n \end{aligned}$$

Lagrangian is

$$\begin{aligned} L(\mathbf{w}, b, \{\xi_n\}, \{\alpha_n\}, \{\lambda_n\}) &= \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n - \sum_n \lambda_n \xi_n \\ &\quad + \sum_n \alpha_n (1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) - \xi_n) \end{aligned}$$

where $\alpha_1, \dots, \alpha_N \geq 0$ and $\lambda_1, \dots, \lambda_N \geq 0$ are Lagrangian multipliers.

October 22, 2019 39 / 55

Applying the stationarity condition

$$L = \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n - \sum_n \lambda_n \xi_n + \sum_n \alpha_n (1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) - \xi_n)$$

\exists primal and dual variables $\mathbf{w}, b, \{\xi_n\}, \{\alpha_n\}, \{\lambda_n\}$ s.t. $\nabla_{\mathbf{w}, b, \{\xi_n\}} L = \mathbf{0}$, which means

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} - \sum_n \alpha_n y_n \phi(\mathbf{x}_n) = \mathbf{0}$$

$$\frac{\partial L}{\partial b} = - \sum_n \alpha_n y_n = 0$$

$$\frac{\partial L}{\partial \xi_n} = C - \lambda_n - \alpha_n = 0, \quad \forall n$$

October 22, 2019 40 / 55

Rewrite the Lagrangian in terms of dual variables

Replacing \mathbf{w} by $\sum_n y_n \alpha_n \phi(\mathbf{x}_n)$, after some simplification, we have

$$\begin{aligned} L &= \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_n C \xi_n - \sum_n \lambda_n \xi_n + \sum_n \alpha_n (1 - y_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) - \xi_n) \\ &= \frac{1}{2} \left\| \sum_n y_n \alpha_n \phi(\mathbf{x}_n) \right\|_2^2 + \sum_n C \xi_n - \sum_n \lambda_n \xi_n + \sum_n \alpha_n - \sum_n \alpha_n \xi_n - \\ &\quad \sum_n \alpha_n y_n \left(\left(\sum_m y_m \alpha_m \phi(\mathbf{x}_m) \right)^T \phi(\mathbf{x}_n) + b \right) \end{aligned}$$

Rewrite the Lagrangian in terms of dual variables

Since $C = \lambda_n + \alpha_n$ (see slide 40), we get

$$\begin{aligned} L &= \frac{1}{2} \left\| \sum_n y_n \alpha_n \phi(\mathbf{x}_n) \right\|_2^2 + \sum_n \alpha_n \\ &\quad - \sum_n \alpha_n y_n \left(\left(\sum_m y_m \alpha_m \phi(\mathbf{x}_m) \right)^T \phi(\mathbf{x}_n) + b \right) \end{aligned}$$

Rewrite the Lagrangian in terms of dual variables

Since $\sum_n \alpha_n y_n = 0$ (see slide 40), we have

$$\begin{aligned} L &= \frac{1}{2} \left\| \sum_n y_n \alpha_n \phi(\mathbf{x}_n) \right\|_2^2 + \sum_n \alpha_n \\ &\quad - \sum_n \alpha_n y_n \left(\sum_m y_m \alpha_m \phi(\mathbf{x}_m) \right)^T \phi(\mathbf{x}_n) \end{aligned}$$

which could be further simplified

$$\begin{aligned} L &= \sum_n \alpha_n + \frac{1}{2} \left\| \sum_n y_n \alpha_n \phi(\mathbf{x}_n) \right\|_2^2 - \sum_{m,n} \alpha_n \alpha_m y_m y_n \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n) \\ &= \sum_n \alpha_n - \frac{1}{2} \sum_{m,n} \alpha_n \alpha_m y_m y_n \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n) \end{aligned}$$

The dual formulation

So the **dual formulation of SVM** is:

$$\max_{\{\alpha_n\}, \{\lambda_n\}} \sum_n \alpha_n - \frac{1}{2} \sum_{m,n} y_m y_n \alpha_m \alpha_n \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n)$$

subject to (see slide 40)

$$\begin{aligned} \sum_n \alpha_n y_n &= 0, \\ C - \lambda_n - \alpha_n &= 0, \\ \alpha_n &\geq 0, \quad \forall n \\ \lambda_n &\geq 0, \quad \forall n \end{aligned}$$

Now it is clear that with a **kernel function** for the mapping ϕ , we can kernelize SVM. That is the reason why we need the dual SVM.

The dual formulation

The last three constraints can be simplified, therefore the **dual formulation of SVM** can be written as

$$\max_{\{\alpha_n\}} \sum_n \alpha_n - \frac{1}{2} \sum_{m,n} y_m y_n \alpha_m \alpha_n k(\mathbf{x}_m, \mathbf{x}_n)$$

subject to

$$\begin{aligned} \sum_n \alpha_n y_n &= 0, \\ 0 \leq \alpha_n &\leq C, \quad \forall n \end{aligned}$$

where $k(x, x')$ is a kernel.

October 22, 2019 45 / 55

Applying complementary slackness

Recall the SVM primal formulation

$$\min_{\mathbf{w}, b, \{\xi_n\}} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n$$

subject to

$$\begin{aligned} 1 - \xi_n - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) &\leq 0, \quad \forall n \\ -\xi_n &\leq 0, \quad \forall n \end{aligned}$$

Recall complementary slackness (slide 37):

$$\lambda_j^* h_j(\mathbf{w}^*) = 0 \quad \text{for all } j \in [J]$$

Therefore, for all n we have

$$\lambda_n^* \xi_n^* = 0, \quad \alpha_n^* (1 - \xi_n^* - y_n(\mathbf{w}^{*T} \phi(\mathbf{x}_n) + b^*)) = 0$$

October 22, 2019 47 / 55

Recover the primal solution

But how do we predict given the dual solution $\{\alpha_n^*\}$? Need to figure out the primal solution \mathbf{w}^* and b^* .

Based on previous observation (see slide 40,

$$\mathbf{w}^* = \sum_n \alpha_n^* y_n \phi(\mathbf{x}_n) = \sum_{n: \alpha_n^* > 0} \alpha_n^* y_n \phi(\mathbf{x}_n)$$

A point with $\alpha_n^* > 0$ is called a “**support vector**”. Hence the name SVM.

To identify b^* , we need to apply complementary slackness.

October 22, 2019 46 / 55

Applying complementary slackness

Complementary slackness:

$$\lambda_n^* \xi_n^* = 0, \quad \alpha_n^* (1 - \xi_n^* - y_n(\mathbf{w}^{*T} \phi(\mathbf{x}_n) + b^*)) = 0$$

For some support vector $\phi(\mathbf{x}_n)$ if we have $0 < \alpha_n^* < C$, then

$$\lambda_n^* = C - \alpha_n^* > 0$$

With the first condition we know $\xi_n^* = 0$.

With the second condition we know $1 = y_n(\mathbf{w}^{*T} \phi(\mathbf{x}_n) + b^*)$ and thus

$$b^* = y_n - \mathbf{w}^{*T} \phi(\mathbf{x}_n) = y_n - \sum_m y_m \alpha_m^* k(\mathbf{x}_m, \mathbf{x}_n)$$

Having both \mathbf{w}^* and b^* we can do prediction on a new point \mathbf{x} :

$$\text{SGN}(\mathbf{w}^{*T} \phi(\mathbf{x}) + b^*) = \text{SGN}\left(\sum_m y_m \alpha_m^* k(\mathbf{x}_m, \mathbf{x}) + b^*\right)$$

October 22, 2019 48 / 55

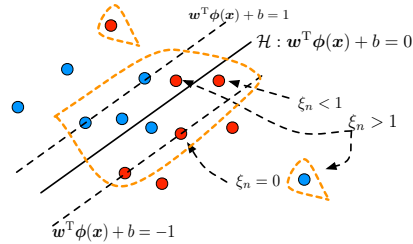
Geometric interpretation of support vectors

A support vector satisfies $\alpha_n^* \neq 0$ and

$$1 - \xi_n^* \leq y_n(\mathbf{w}^{*\top} \phi(\mathbf{x}_n) + b^*)$$

When

- $\xi_n^* = 0$, $y_n(\mathbf{w}^{*\top} \phi(\mathbf{x}_n) + b^*) = 1$ and thus the point is $1/\|\mathbf{w}^*\|_2$ away from the hyperplane.
- $\xi_n^* < 1$, the point is classified correctly but does not satisfy the large margin constraint.
- $\xi_n^* > 1$, the point is misclassified.

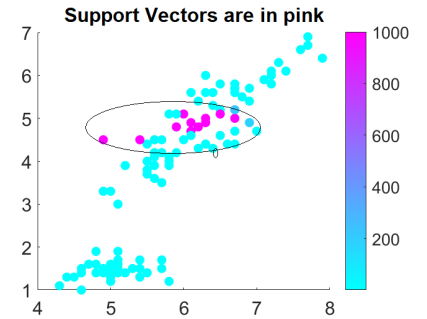
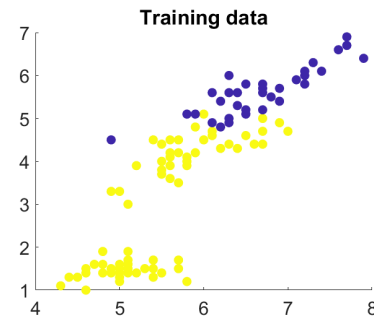


Support vectors (circled with the orange line) are *the only points that matter!*

An example

One drawback of kernel method: **non-parametric**, need to keep all training points potentially

However, for SVM, very often **#support vectors** $\ll N$



Summary

Interpretation: maximize the margin

For separable data

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & y_n[\mathbf{w}^\top \phi(\mathbf{x}_n) + b] \geq 1, \quad \forall n \end{aligned}$$

For non-separable data

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n \\ \text{s.t.} \quad & y_n[\mathbf{w}^\top \phi(\mathbf{x}_n) + b] \geq 1 - \xi_n, \quad \forall n \\ & \xi_n \geq 0, \quad \forall n \end{aligned}$$

where C is a hyperparameter and ξ_n are slack variables.

Summary

Interpretation: minimize loss

Minimize loss on all data

$$\min_{\mathbf{w}, b} \sum_n \max(0, 1 - y_n[\mathbf{w}^\top \phi(\mathbf{x}_n) + b]) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

equivalently

$$\begin{aligned} \min_{\mathbf{w}, b, \{\xi_n\}} \quad & C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & 1 - y_n[\mathbf{w}^\top \phi(\mathbf{x}_n) + b] \leq \xi_n, \quad \forall n \\ & \xi_n \geq 0, \quad \forall n \end{aligned}$$

Summary

SVM: **max-margin linear classifier**

Primal (equivalent to minimizing L2 regularized hinge loss):

$$\min_{\mathbf{w}, b, \{\xi_n\}} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n$$

subject to

$$\begin{aligned} \xi_n &\geq 1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b), \quad \forall n \\ \xi_n &\geq 0, \quad \forall n \end{aligned}$$

Dual (kernelizable, reveals what training points are support vectors):

$$\begin{aligned} \max_{\{\alpha_n\}} \quad & \sum_n \alpha_n - \frac{1}{2} \sum_{m,n} y_m y_n \alpha_m \alpha_n \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n) \\ \text{s.t.} \quad & \sum_n \alpha_n y_n = 0 \quad \text{and} \quad 0 \leq \alpha_n \leq C, \quad \forall n \end{aligned}$$

October 22, 2019 53 / 55

Summary

Typical steps of applying Lagrangian duality

- start with a primal problem
- write down the Lagrangian (one dual variable per constraint)
- apply KKT conditions to find the **connections between primal and dual solutions**
- **eliminate primal variables** and arrive at the dual formulation
- maximize the Lagrangian with respect to dual variables
- recover the primal solutions from the dual solutions

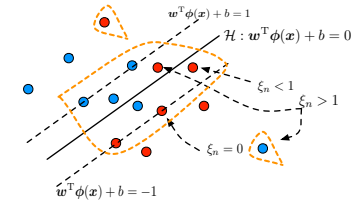
October 22, 2019 55 / 55

Geometric interpretation of support vectors

Nonzero α_n is called support vector

Some α_n will become zero

$$\begin{aligned} \min_{\alpha} \quad & \sum_n \alpha_n - \frac{1}{2} \sum_{m,n} y_m y_n \alpha_m \alpha_n k(\mathbf{x}_m, \mathbf{x}_n) \\ \text{s.t.} \quad & 0 \leq \alpha_n \leq C, \quad \forall n \\ & \sum_n \alpha_n y_n = 0 \end{aligned}$$



Support vectors are those being circled with the orange line. Removing them will change the solution.

October 22, 2019 54 / 55