18-661 Introduction to Machine Learning

SVM - III

Spring 2020

ECE - Carnegie Mellon University

Midterm Information

Midterm will be on Wednesday, 2/26 in-class.

- Closed-book except for one double-sided letter-size handwritten page of notes that you can prepare as you wish.
- We will provide formulas for relevant probability distributions.
- You will not need a calculator. Only pen/pencil and scratch paper are allowed.

Will cover all topics presented through this Wednesday in class (SVM and before).

- (1) point estimation/MLE/MAP, (2) linear regression, (3) naive Bayes, (4) logistic regression, and (5) SVMs.
- This friday's recitation will go over practice exam questions.

1

Midterm: Concepts That You Should Know

This is a quick overview of the most important concepts/methods/models that you should expect to see on the midterm.

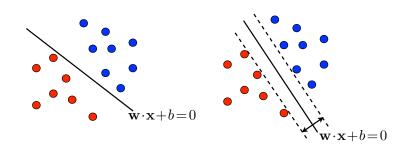
- MLE/MAP: how to find the likelihood of one or more observations given a system model, how to incorporate knowledge of a prior distribution, how to optimize the likelihood, loss functions
- Linear regression: how to formulate the linear regression optimization problem, how it relates to MLE/MAP, ridge regression, overfitting and regularization, gradient descent, bias-variance trade-off
- Naive Bayes: Bayes' rule, naive classification rule, why it is naive
- Logistic regression: how to formulate logistic regression, how it relates to MLE, comparison to naive Bayes, sigmoid function, softmax function, cross-entropy function
- SVMs: hinge loss formulation, max-margin formulation, dual of the SVM problem, kernel functions

Outline

- 1. Review of SVM Max Margin Formulation
- 2. A Dual View of SVMs (the short version)
- 3. Lagrange Duality and KKT conditions (optional)
- 4. Dual Derivation of SVMs (optional)
- 5. Kernel SVM

Review of SVM Max Margin Formulation

Intuition: Where to put the decision boundary?



Idea: Find a decision boundary in the 'middle' of the two classes that:

- Perfectly classifies the training data
- Is as far away from every training point as possible

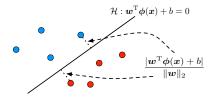
Let us apply this intuition to build a classifier that MAXIMIZES THE MARGIN between training points and the decision boundary

Defining the Margin

Margin

Smallest distance between the hyperplane and all training points

$$MARGIN(\boldsymbol{w}, b) = \min_{n} \frac{y_{n}[\boldsymbol{w}^{\top}\boldsymbol{x}_{n} + b]}{\|\boldsymbol{w}\|_{2}}$$



Rescaled Margin to Avoid Scaling Ambiguity

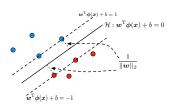
We can further constrain the problem by scaling (w, b) such that

$$\min_{n} y_{n}[\mathbf{w}^{\top} \mathbf{x}_{n} + b] = 1$$

We've fixed the numerator in the MARGIN(\boldsymbol{w}, b) equation, and we have:

MARGIN
$$(\boldsymbol{w}, b) = \frac{\min_{n} y_{n} [\boldsymbol{w}^{\top} \boldsymbol{x}_{n} + b]}{\|\boldsymbol{w}\|_{2}} = \frac{1}{\|\boldsymbol{w}\|_{2}}$$

Hence the points closest to the decision boundary are at distance $\frac{1}{\|\mathbf{w}\|_2}$!



SVM: max margin formulation for separable data

Assuming separable training data, we thus want to solve:

$$\max_{\boldsymbol{w},b} \frac{1}{\|\boldsymbol{w}\|_2} \quad \text{such that} \quad \underbrace{y_n[\boldsymbol{w}^\top \boldsymbol{x}_n + b] \geq 1, \ \forall \ n}_{\text{scaling of } \boldsymbol{w}, b}$$

This is equivalent to

$$\min_{\boldsymbol{w},b} \quad \frac{1}{2} \|\boldsymbol{w}\|_2^2$$
s.t. $y_n[\boldsymbol{w}^\top \boldsymbol{x}_n + b] \ge 1, \quad \forall \quad n$

Given our geometric intuition, SVM is called a **max margin** (or large margin) classifier. The constraints are called **large margin constraints**.

SVM for non-separable data

Constraints in separable setting

$$y_n[\mathbf{w}^{\top}\mathbf{x}_n + b] \ge 1, \quad \forall \quad n$$

Constraints in non-separable setting

Idea: modify our constraints to account for non-separability! Specifically, we introduce slack variables $\xi_n \geq 0$:

$$y_n[\mathbf{w}^{\top}\mathbf{x}_n + b] \ge 1 - \xi_n, \ \forall \ n$$

- For "hard" training points, we can increase ξ_n until the above inequalities are met
- What does it mean when ξ_n is very large?

8

Soft-margin SVM formulation

We do not want ξ_n to grow too large, and we can control their size by incorporating them into our optimization problem:

$$\min_{\boldsymbol{w},b,\boldsymbol{\xi}} \quad \frac{1}{2} \|\boldsymbol{w}\|_{2}^{2} + C \sum_{n} \xi_{n}$$
s.t. $y_{n} [\boldsymbol{w}^{\top} \boldsymbol{x}_{n} + b] \ge 1 - \xi_{n}, \quad \forall \quad n$

$$\xi_{n} \ge 0, \quad \forall \quad n$$

What is the role of C?

- User-defined hyperparameter
- Trades off between the two terms in our objective
- Same idea as the regularization term in ridge regression,

How to solve this problem?

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

- This is a convex quadratic program: the objective function is quadratic in w and linear in ξ and the constraints are linear (inequality) constraints in w, b and ξ_n.
- We can solve the optimization problem using general-purpose solvers, e.g., Matlab's quadprog() function.

A Dual View of SVMs (the short version)

What is duality?

Duality is a way of transforming a constrained optimization problem.

It tells us sometimes-useful information about the problem structure, and can sometimes make the problem easier to solve.

- Dual problem is always convex-easy to solve.
- Primal and dual problems are not always equivalent.
- Dual variables tell us "how bad" constraints are.

What is duality?

Duality is a way of transforming a constrained optimization problem.

It tells us sometimes-useful information about the problem structure, and can sometimes make the problem easier to solve.

- Dual problem is always convex—easy to solve.
- Primal and dual problems are not always equivalent.
- Dual variables tell us "how bad" constraints are.

The main point you should understand is that we will solve the dual SVM problem in lieu of the max margin (primal) formulation

Derivation of the dual

Here is a skeleton of how to derive the dual problem.

Recipe

- 1. Formulate the generalized Lagrangian function that incorporates the constraints and introduces dual variables
- 2. Minimize the Lagrangian function over the primal variables
- 3. Substitute the primal variables for dual variables in the Lagrangian
- 4. Maximize the Lagrangian with respect to dual variables
- Recover the solution (for the primal variables) from the dual variables

Primal SVM

$$\min_{\boldsymbol{w},b,\boldsymbol{\xi}} \quad \frac{1}{2} \|\boldsymbol{w}\|_{2}^{2} + C \sum_{n} \xi_{n}$$
s.t. $y_{n} [\boldsymbol{w}^{\top} \boldsymbol{x}_{n} + b] \ge 1 - \xi_{n}, \quad \forall \quad n$

$$\xi_{n} \ge 0, \quad \forall \quad n$$

The constraints are equivalent to the following canonical forms:

$$-\xi_n \leq 0$$
 and $1 - y_n[\boldsymbol{w}^{\top} \boldsymbol{x}_n + b] - \xi_n \leq 0$

Lagrangian

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

Lagrangian

$$L(\mathbf{w}, b, \{\xi_n\}, \{\alpha_n\}, \{\lambda_n\}) = C \sum_{n} \xi_n + \frac{1}{2} ||\mathbf{w}||_2^2 - \sum_{n} \lambda_n \xi_n + \sum_{n} \alpha_n \{1 - y_n [\mathbf{w}^{\top} \mathbf{x}_n + b] - \xi_n \}$$

under the constraints that $\alpha_n \geq 0$ and $\lambda_n \geq 0$.

• Primal variables: \mathbf{w} , $\{\xi_n\}$, b; dual variables $\{\lambda_n\}$, $\{\alpha_n\}$

Lagrangian

$$L(\mathbf{w}, b, \{\xi_n\}, \{\alpha_n\}, \{\lambda_n\}) = C \sum_{n} \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 - \sum_{n} \lambda_n \xi_n + \sum_{n} \alpha_n \{1 - y_n [\mathbf{w}^{\top} \mathbf{x}_n + b] - \xi_n\}$$

- Primal variables: \mathbf{w} , $\{\xi_n\}$, b; dual variables $\{\lambda_n\}$, $\{\alpha_n\}$
- Minimize the Lagrangian function over the primal variables by setting $\frac{\partial L}{\partial \mathbf{w}} = 0$, $\frac{\partial L}{\partial b} = 0$, and $\frac{\partial L}{\partial \xi_n} = 0$.

Lagrangian

$$L(\mathbf{w}, b, \{\xi_n\}, \{\alpha_n\}, \{\lambda_n\}) = C \sum_{n} \xi_n + \frac{1}{2} ||\mathbf{w}||_2^2 - \sum_{n} \lambda_n \xi_n + \sum_{n} \alpha_n \{1 - y_n [\mathbf{w}^{\top} \mathbf{x}_n + b] - \xi_n\}$$

- Primal variables: \mathbf{w} , $\{\xi_n\}$, b; dual variables $\{\lambda_n\}$, $\{\alpha_n\}$
- Minimize the Lagrangian function over the primal variables by setting $\frac{\partial L}{\partial \mathbf{w}} = 0$, $\frac{\partial L}{\partial b} = 0$, and $\frac{\partial L}{\partial \xi_n} = 0$.
- Substitute the solutions to primal variables for dual variables in the Lagrangian

Lagrangian

$$L(\mathbf{w}, b, \{\xi_n\}, \{\alpha_n\}, \{\lambda_n\}) = C \sum_{n} \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 - \sum_{n} \lambda_n \xi_n + \sum_{n} \alpha_n \{1 - y_n [\mathbf{w}^{\top} \mathbf{x}_n + b] - \xi_n\}$$

- Primal variables: \mathbf{w} , $\{\xi_n\}$, b; dual variables $\{\lambda_n\}$, $\{\alpha_n\}$
- Minimize the Lagrangian function over the primal variables by setting $\frac{\partial L}{\partial \mathbf{w}} = 0$, $\frac{\partial L}{\partial b} = 0$, and $\frac{\partial L}{\partial \xi_n} = 0$.
- Substitute the solutions to primal variables for dual variables in the Lagrangian
- Maximize the Lagrangian with respect to dual variables

Lagrangian

$$L(\mathbf{w}, b, \{\xi_n\}, \{\alpha_n\}, \{\lambda_n\}) = C \sum_{n} \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 - \sum_{n} \lambda_n \xi_n + \sum_{n} \alpha_n \{1 - y_n [\mathbf{w}^{\top} \mathbf{x}_n + b] - \xi_n\}$$

- Primal variables: \mathbf{w} , $\{\xi_n\}$, b; dual variables $\{\lambda_n\}$, $\{\alpha_n\}$
- Minimize the Lagrangian function over the primal variables by setting $\frac{\partial L}{\partial \mathbf{w}} = 0$, $\frac{\partial L}{\partial b} = 0$, and $\frac{\partial L}{\partial \xi_n} = 0$.
- Substitute the solutions to primal variables for dual variables in the Lagrangian
- Maximize the Lagrangian with respect to dual variables
- After some further maths and simplifications, we have...

Dual formulation of SVM

Dual is also a convex quadratic program

$$\begin{aligned} \max_{\alpha} \quad & \sum_{n} \alpha_{n} - \frac{1}{2} \sum_{m,n} y_{m} y_{n} \alpha_{m} \alpha_{n} \mathbf{x}_{m}^{\top} \mathbf{x}_{n} \\ \text{s.t.} \quad & 0 \leq \alpha_{n} \leq C, \quad \forall \ n \\ & \sum_{n} \alpha_{n} y_{n} = 0 \end{aligned}$$

- There are N dual variables α_n , one for each data point
- Independent of the size d of x: SVM scales better for high-dimensional feature.
- May seem like a lot of optimization variables when N is large, but many of the α_n 's become zero. α_n is non-zero only if the n^{th} point is a support vector

Why do many α_n 's become zero?

$$\max_{\alpha} \sum_{n} \alpha_{n} - \frac{1}{2} \sum_{m,n} y_{m} y_{n} \alpha_{m} \alpha_{n} \mathbf{x}_{m}^{\top} \mathbf{x}_{n}$$
s.t. $0 \le \alpha_{n} \le C$, $\forall n$

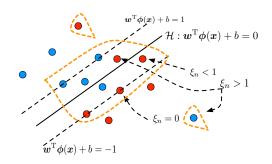
$$\sum_{n} \alpha_{n} y_{n} = 0$$

 By strong duality and KKT complementary slackness conditions, it tells us:

$$\alpha_n \{1 - \xi_n - y_n [\mathbf{w}^\top \mathbf{x}_n + b]\} = 0 \quad \forall n$$

- This tells us that $\alpha_n > 0$ iff $1 \xi_n = y_n[\mathbf{w}^\top \mathbf{x}_n + b]$
 - If $\xi_n = 0$, then support vector is on the margin
 - Otherwise, $\xi_n > 0$ means that the point is an outlier

Visualizing the support vectors



Support vectors $(\alpha_n > 0)$ are highlighted by the dotted orange lines.

- $\xi_n = 0$ and $0 < \alpha_n < C$ when $y_n[\mathbf{w}^{\top} \mathbf{x}_n + b] = 1$.
- $\xi_n > 0$ and $\alpha_n = 0$ if $y_n[\mathbf{w}^\top \mathbf{x}_n + b] < 1$.

Once we solve for α_n 's, how to get w and b?

Recovering w

$$\frac{\partial L}{\partial \mathbf{w}} = 0 \to \mathbf{w} = \sum_{n} \alpha_{n} y_{n} \mathbf{x}_{n}$$

Only depends on support vectors, i.e., points with $\alpha_n > 0!$

Once we solve for α_n 's, how to get w and b?

Recovering w

$$\frac{\partial L}{\partial \mathbf{w}} = 0 \to \mathbf{w} = \sum_{n} \alpha_{n} \mathbf{y}_{n} \mathbf{x}_{n}$$

Only depends on support vectors, i.e., points with $\alpha_n > 0!$

Recovering b

If $0 < \alpha_n < C$ and $y_n \in \{-1, 1\}$:

$$y_n[\mathbf{w}^{\top} \mathbf{x}_n + b] = 1$$

$$b = y_n - \mathbf{w}^{\top} \mathbf{x}_n$$

$$b = y_n - \sum_m \alpha_m y_m \mathbf{x}_m^{\top} \mathbf{x}_n$$

Outline

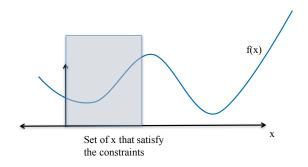
- 1. Review of SVM Max Margin Formulation
- 2. A Dual View of SVMs (the short version)
- 3. Lagrange Duality and KKT conditions (optional)
- 4. Dual Derivation of SVMs (optional)
- 5. Kernel SVM

Lagrange Duality and KKT

conditions (optional)

$$\begin{cases} \min_{\mathbf{x}} & f(\mathbf{x}) \\ \text{s.t.} & g_i(\mathbf{x}) \leq 0, \quad \forall i \\ & h_i(\mathbf{x}) = 0, \quad \forall j \end{cases}$$

This is the 'primal' problem.



$$\begin{cases} \min_{\mathbf{x}} & f(\mathbf{x}) \\ \text{s.t.} & g_i(\mathbf{x}) \leq 0, \quad \forall i \\ & h_i(\mathbf{x}) = 0, \quad \forall j \end{cases}$$

This is the 'primal' problem. The generalized Lagrangian is defined as follows:

$$\begin{cases} \min_{\mathbf{x}} & f(\mathbf{x}) \\ \text{s.t.} & g_i(\mathbf{x}) \leq 0, \quad \forall i \\ & h_i(\mathbf{x}) = 0, \quad \forall j \end{cases}$$

This is the 'primal' problem. The generalized Lagrangian is defined as follows:

$$L(\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = f(\mathbf{x}) + \sum_{i} \alpha_{i} g_{i}(\mathbf{x}) + \sum_{j} \beta_{j} h_{j}(\mathbf{x})$$

$$\begin{cases} \min_{\mathbf{x}} & f(\mathbf{x}) \\ \text{s.t.} & g_i(\mathbf{x}) \leq 0, \quad \forall i \\ & h_i(\mathbf{x}) = 0, \quad \forall j \end{cases}$$

This is the 'primal' problem. The generalized Lagrangian is defined as follows:

$$L(\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = f(\mathbf{x}) + \sum_{i} \alpha_{i} g_{i}(\mathbf{x}) + \sum_{j} \beta_{j} h_{j}(\mathbf{x})$$

Consider the following function:

$$\theta_P(\mathbf{x}) = \max_{\alpha, \beta, \alpha_i \geq 0} L(\mathbf{x}, \alpha, \beta)$$

$$\begin{cases} \min_{\mathbf{x}} & f(\mathbf{x}) \\ \text{s.t.} & g_i(\mathbf{x}) \leq 0, \quad \forall i \\ & h_i(\mathbf{x}) = 0, \quad \forall j \end{cases}$$

This is the 'primal' problem. The generalized Lagrangian is defined as follows:

$$L(\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = f(\mathbf{x}) + \sum_{i} \alpha_{i} g_{i}(\mathbf{x}) + \sum_{j} \beta_{j} h_{j}(\mathbf{x})$$

Consider the following function:

$$\theta_P(\mathbf{x}) = \max_{\boldsymbol{\alpha}, \boldsymbol{\beta}, \alpha_i \geq 0} L(\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta})$$

• If x violates a primal constraint, $\theta_P(x) = \infty$; otherwise $\theta_P(x) = f(x)$

$$\begin{cases} \min_{\mathbf{x}} & f(\mathbf{x}) \\ \text{s.t.} & g_i(\mathbf{x}) \leq 0, \quad \forall i \\ & h_i(\mathbf{x}) = 0, \quad \forall j \end{cases}$$

This is the 'primal' problem. The generalized Lagrangian is defined as follows:

$$L(\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = f(\mathbf{x}) + \sum_{i} \alpha_{i} g_{i}(\mathbf{x}) + \sum_{j} \beta_{j} h_{j}(\mathbf{x})$$

Consider the following function:

$$\theta_P(\mathbf{x}) = \max_{\boldsymbol{\alpha}, \boldsymbol{\beta}, \alpha_i \geq 0} L(\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta})$$

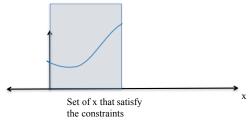
- If x violates a primal constraint, $\theta_P(x) = \infty$; otherwise $\theta_P(x) = f(x)$
- Thus $\min_{\mathbf{x}} \theta_P(\mathbf{x}) = \min_{\mathbf{x}} \max_{\alpha, \beta, \alpha_i \geq 0} L(\mathbf{x}, \alpha, \beta)$ has same solution as the primal problem, which we denote as p^*

$$\begin{cases} \min_{\mathbf{x}} & f(\mathbf{x}) \\ \text{s.t.} & g_i(\mathbf{x}) \leq 0, \quad \forall i \\ & h_i(\mathbf{x}) = 0, \quad \forall j \end{cases}$$

This is the 'primal' problem.

$$\theta_P(\mathbf{x}) = \max_{\alpha, \beta, \alpha_i \geq 0} L(\mathbf{x}, \alpha, \beta)$$

 $\max_{\alpha>0,\beta}L(x,\alpha,\beta) \ \ \text{is equal to f(x) for the feasible} \\ \text{x and infinity everywhere else}$



Constrained Optimization – Inequality Constraints

Primal Problem

$$p^* = \min_{\mathbf{x}} \theta_P(\mathbf{x}) = \min_{\mathbf{x}} \max_{\alpha, \beta, \alpha_i > 0} L(\mathbf{x}, \alpha, \beta)$$

Dual Problem

Consider the function:

$$\theta_D(\alpha,\beta) = \min_{\mathbf{x}} L(\mathbf{x},\alpha,\beta)$$

$$d^* = \max_{\boldsymbol{\alpha}, \boldsymbol{\beta}, \alpha_i \geq 0} \theta_D(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \max_{\boldsymbol{\alpha}, \boldsymbol{\beta}, \alpha_i \geq 0} \min_{\boldsymbol{x}} L(\boldsymbol{x}, \boldsymbol{\alpha}, \boldsymbol{\beta})$$

Relationship between primal and dual?

- $p^* \ge d^*$ (weak duality)
- 'min max' of any function is always greater than the 'max min'
- https://en.wikipedia.org/wiki/Max%E2%80%93min_inequality

How to find the solution $p^* = d^*$? Use KKT Conditions

Strong duality implies that there exist x^*, α^*, β^* such that:

- x^* is the solution to the primal and α^*, β^* is the solution to the dual
- $p^* = d^* = L(x^*, \alpha^*, \beta^*)$
- x^*, α^*, β^* satisfy the KKT conditions (which in fact are necessary and sufficient for optimality)

The Karush-Kuhn-Tucker (KKT) conditions are:

- Stationarity: $\frac{\partial L(\mathbf{x}, \mathbf{\alpha}^*, \boldsymbol{\beta}^*)}{\partial \mathbf{x}}|_{\mathbf{x}^*} = 0$. \mathbf{x}^* is a local extremum of the Lagrangian L for fixed $\mathbf{\alpha}^*, \boldsymbol{\beta}^*$.
- Feasibility: $g_i(\mathbf{x}^*) \leq 0$ and $h_i(\mathbf{x}^*) = 0$ (primal) and $\alpha_i^* \geq 0$ (dual) for all i. All primal and dual constraints are satisfied.
- Complementary slackness: $\alpha_i^* g_i(\mathbf{x}^*) = 0$ for all i. Either the Lagrange multiplier α_i^* is 0, or the corresponding constraint $g_i(\mathbf{x}^*) \leq 0$ is tight (i.e., $g_i(\mathbf{x}^*) = 0$).

Nutrients	Food 1 (\$2)	Food 2 (\$5)	Food 3 (\$15)
Carbs	20	1	1
Protein	1	30	40
Vitamins	1	10	5

To satisfy daily nutritional requirements we need at least

- 200 units of carbs
- 50 units of protein
- 40 units of vitamins

Primal problem: How do we minimize the cost of satisfying these requirements?

Nutrients	Food 1 (\$2)	Food 2 (\$5)	Food 3 (\$15)
Carbs	20	1	1
Protein	1	30	40
Vitamins	1	10	5

$$\min_{\substack{x_1, x_2, x_3 \\ \text{s.t.}}} 2x_1 + 5x_2 + 15x_3$$

$$-20x_1 - x_2 - x_3 \le -200$$

$$-x_1 - 30x_2 - 40x_3 \le -50$$

$$-x_1 - 10x_2 - 5x_3 \le -40$$

Nutrients	Food 1 (\$2)	Food 2 (\$5)	Food 3 (\$15)
Carbs	20	1	1
Protein	1	30	40
Vitamins	1	10	5

$$\min_{\substack{x_1, x_2, x_3 \\ \text{s.t.}}} 2x_1 + 5x_2 + 15x_3$$

$$-20x_1 - x_2 - x_3 \le -200$$

$$-x_1 - 30x_2 - 40x_3 \le -50$$

$$-x_1 - 10x_2 - 5x_3 \le -40$$

Primal Solution: $x_1 \approx 9.84$, $x_2 \approx 3$, $x_3 = 0$ Dual Solution: $\alpha_1 = 0.07$, $\alpha_2 = 0$, $\alpha_3 = 0.5$

Nutrients	Food 1 (\$2)	Food 2 (\$5)	Food 3 (\$15)
Carbs	20	1	1
Protein	1	30	40
Vitamins	1	10	5

$$\min_{\substack{x_1, x_2, x_3 \\ \text{s.t.}}} 2x_1 + 5x_2 + 15x_3$$

$$-20x_1 - x_2 - x_3 \le -200$$

$$-x_1 - 30x_2 - 40x_3 \le -50$$

$$-x_1 - 10x_2 - 5x_3 \le -40$$

Primal Solution: $x_1 \approx 9.84$, $x_2 \approx 3$, $x_3 = 0$ Dual Solution: $\alpha_1 = 0.07$, $\alpha_2 = 0$, $\alpha_3 = 0.5$

 $\alpha_2=0$ means that protein requirement is easy to satisfy

The Diet problem: Dual

Nutrients	Food 1 (\$2)	Food 2 (\$5)	Food 3 (\$15)
Carbs	20	1	1
Protein	1	30	40
Vitamins	1	10	5

A pharmacist wants to create a diet pill to satisfy the requirements and maximize profit. α_1 , α_2 , α_3 are the shadow prices of the nutrients.

$$\begin{aligned} \max_{\alpha_1,\alpha_2,\alpha_3} & 200\alpha_1 + 50\alpha_2 + 40\alpha_3 \\ \text{s.t.} & 20\alpha_1 + \alpha_2 + \alpha_3 \leq 2 \\ & \alpha_1 + 30\alpha_2 + 10\alpha_3 \leq 5 \\ & \alpha_1 + 40\alpha_2 + 5\alpha_3 \leq 15 \\ & \alpha_1,\alpha_2,\alpha_3 \geq 0 \end{aligned}$$

The Diet problem: Dual

Nutrients	Food 1 (\$2)	Food 2 (\$5)	Food 3 (\$15)
Carbs	20	1	1
Protein	1	30	40
Vitamins	1	10	5

A pharmacist wants to create a diet pill to satisfy the requirements and maximize profit. α_1 , α_2 , α_3 are the shadow prices of the nutrients.

$$\begin{array}{ll} \max _{\alpha_{1},\alpha_{2},\alpha_{3}} & 200\alpha_{1}+50\alpha_{2}+40\alpha_{3} \\ \text{s.t.} & 20\alpha_{1}+\alpha_{2}+\alpha_{3}\leq 2 \\ & \alpha_{1}+30\alpha_{2}+10\alpha_{3}\leq 5 \\ & \alpha_{1}+40\alpha_{2}+5\alpha_{3}\leq 15 \\ & \alpha_{1},\alpha_{2},\alpha_{3}\geq 0 \end{array}$$

Primal Solution: $x_1 \approx 9.84$, $x_2 \approx 3$, $x_3 = 0$ Dual Solution: $\alpha_1 = 0.07$, $\alpha_2 = 0$, $\alpha_3 = 0.5$

Recap

- When working with constrained optimization problems with inequality constraints, we can write down primal and dual problems.
- The dual solution is always a lower bound on the primal solution (weak duality) and it is always convex
- The duality gap equals 0 under certain conditions (strong duality), and in such cases we can either solve the primal or dual problem.
- Strong duality (and thus the KKT conditions) hold for the SVM problem.
- See http://cs229.stanford.edu/notes/cs229-notes3.pdf for details

Dual Derivation of SVMs

(optional)

Derivation of the dual

We will next derive the dual formulation for SVMs.

Recipe

- 1. Formulate the generalized Lagrangian function that incorporates the constraints and introduces dual variables
- 2. Minimize the Lagrangian function over the primal variables
- 3. Substitute the primal variables for dual variables in the Lagrangian
- 4. Maximize the Lagrangian with respect to dual variables
- Recover the solution (for the primal variables) from the dual variables

Deriving the dual for SVM

Primal SVM

$$\min_{\boldsymbol{w},b,\boldsymbol{\xi}} \quad \frac{1}{2} \|\boldsymbol{w}\|_{2}^{2} + C \sum_{n} \xi_{n}$$
s.t. $y_{n} [\boldsymbol{w}^{\top} \boldsymbol{x}_{n} + b] \ge 1 - \xi_{n}, \quad \forall \quad n$

$$\xi_{n} \ge 0, \quad \forall \quad n$$

The constraints are equivalent to $-\xi_n \leq 0$ and $1 - v_n[\mathbf{w}^{\top} \mathbf{x}_n + b] - \xi_n \leq 0$.

Lagrangian

$$L(\boldsymbol{w}, b, \{\xi_n\}, \{\alpha_n\}, \{\lambda_n\}) = C \sum_n \xi_n + \frac{1}{2} \|\boldsymbol{w}\|_2^2 - \sum_n \lambda_n \xi_n$$
$$+ \sum_n \alpha_n \{1 - y_n [\boldsymbol{w}^\top \boldsymbol{x}_n + b] - \xi_n\}$$

under the constraints that $\alpha_n \geq 0$ and $\lambda_n \geq 0$.

Minimizing the Lagrangian

Taking derivatives with respect to the primal variables

$$\frac{\partial L}{\partial \mathbf{w}} = \frac{\partial}{\partial \mathbf{w}} \left(\frac{1}{2} \| \mathbf{w} \|_{2}^{2} - \sum_{n} \alpha_{n} y_{n} \mathbf{w}^{\top} \mathbf{x}_{n} \right) = \mathbf{w} - \sum_{n} y_{n} \alpha_{n} \mathbf{x}_{n} = 0$$

$$\frac{\partial L}{\partial b} = \frac{\partial}{\partial b} - \sum_{n} \alpha_{n} y_{n} b = -\sum_{n} \alpha_{n} y_{n} = 0$$

$$\frac{\partial L}{\partial \xi_{n}} = \frac{\partial}{\partial \xi_{n}} (C - \lambda_{n} - \alpha_{n}) \xi_{n} = C - \lambda_{n} - \alpha_{n} = 0$$

These equations link the primal variables and the dual variables and provide new constraints on the dual variables:

$$\mathbf{w} = \sum_{n} y_{n} \alpha_{n} \mathbf{x}_{n}$$
$$\sum_{n} \alpha_{n} y_{n} = 0$$
$$C - \lambda_{n} - \alpha_{n} = 0$$

Rearrange the Lagrangian

$$L(\cdot) = C \sum_{n} \xi_{n} + \frac{1}{2} \|\mathbf{w}\|_{2}^{2} - \sum_{n} \lambda_{n} \xi_{n} + \sum_{n} \alpha_{n} \{1 - y_{n} [\mathbf{w}^{\top} \mathbf{x}_{n} + b] - \xi_{n} \}$$

where $\alpha_n \geq 0$ and $\lambda_n \geq 0$. We now know that $\mathbf{w} = \sum_n y_n \alpha_n \mathbf{x}_n$.

$$g(\{\alpha_n\}, \{\lambda_n\}) = L(\mathbf{w}, b, \{\xi_n\}, \{\alpha_n\}, \{\lambda_n\})$$

$$= \underbrace{\sum_{n} (C - \alpha_n - \lambda_n) \xi_n}_{\text{gather terms with } \xi_n} + \frac{1}{2} \| \underbrace{\sum_{n} y_n \alpha_n \mathbf{x}_n}_{\text{substitute for } \mathbf{w}} \|_2^2 + \sum_{n} \alpha_n$$

$$- \underbrace{\sum_{n} \alpha_n y_n}_{\text{again substitute for } \mathbf{w}}_{\text{again substitute for } \mathbf{w}}$$

Rearrange the Lagrangian

$$L(\cdot) = C \sum_{n} \xi_{n} + \frac{1}{2} \|\mathbf{w}\|_{2}^{2} - \sum_{n} \lambda_{n} \xi_{n} + \sum_{n} \alpha_{n} \{1 - y_{n} [\mathbf{w}^{\top} \mathbf{x}_{n} + b] - \xi_{n} \}$$

where $\alpha_n \geq 0$ and $\lambda_n \geq 0$. We now know that $\mathbf{w} = \sum_n y_n \alpha_n \mathbf{x}_n$.

$$g(\{\alpha_n\}, \{\lambda_n\}) = L(\mathbf{w}, b, \{\xi_n\}, \{\alpha_n\}, \{\lambda_n\})$$

$$= \underbrace{\sum_{n} (C - \alpha_n - \lambda_n) \xi_n}_{\text{gather terms with } \xi_n} + \frac{1}{2} \| \underbrace{\sum_{n} y_n \alpha_n \mathbf{x}_n}_{\text{substitute for } \mathbf{w}} \|_2^2 + \sum_{n} \alpha_n$$

$$- \underbrace{\sum_{n} \alpha_n y_n}_{\text{again substitute for } \mathbf{w}}_{\text{again substitute for } \mathbf{w}}$$

Then, set $\sum_{n} \alpha_{n} y_{n} = 0$ and $C - \lambda_{n} - \alpha_{n} = 0$ and simplify to get..

Incorporate the constraints

Constraints from partial derivatives: $\sum_{n} \alpha_{n} y_{n} = 0$ and $C - \lambda_{n} - \alpha_{n} = 0$.

$$g(\{\alpha_n\}, \{\lambda_n\}) = L(\mathbf{w}, b, \{\xi_n\}, \{\alpha_n\}, \{\lambda_n\})$$

$$= \sum_{n} \underbrace{(C - \alpha_n - \lambda_n)}_{\text{equal to 0!}} \xi_n + \frac{1}{2} \| \sum_{n} y_n \alpha_n \mathbf{x}_n \|_2^2 + \sum_{n} \alpha_n$$

$$- \underbrace{\left(\sum_{n} \alpha_n y_n\right)}_{\text{equal to 0!}} b - \sum_{n} \alpha_n y_n \left(\sum_{m} y_m \alpha_m \mathbf{x}_m\right)^{\top} \mathbf{x}_n$$

$$= \sum_{n} \alpha_n + \frac{1}{2} \| \sum_{n} y_n \alpha_n \mathbf{x}_n \|_2^2 - \sum_{m,n} \alpha_n \alpha_m y_m y_n \mathbf{x}_m^{\top} \mathbf{x}_n$$

$$= \sum_{n} \alpha_n - \frac{1}{2} \sum_{m,n} \alpha_n \alpha_m y_m y_n \mathbf{x}_m^{\top} \mathbf{x}_n$$

The dual problem

Maximizing the dual under the constraints

$$\max_{\alpha} g(\{\alpha_n\}, \{\lambda_n\}) = \sum_{n} \alpha_n - \frac{1}{2} \sum_{m,n} y_m y_n \alpha_m \alpha_n \mathbf{x}_m^{\top} \mathbf{x}_n$$
s.t. $\alpha_n \ge 0, \quad \forall n$

$$\sum_{n} \alpha_n y_n = 0$$

$$C - \lambda_n - \alpha_n = 0, \quad \forall n$$

$$\lambda_n \ge 0, \quad \forall n$$

We can simplify as the objective function does not depend on λ_n . Specifically, we can combine the constraints involving λ_n resulting in the following inequality constraint: $\alpha_n \leq C$:

$$C - \lambda_n - \alpha_n = 0, \ \lambda_n \ge 0 \iff \lambda_n = C - \alpha_n \ge 0$$

$$\iff \alpha_n \le C$$

Dual formulation of SVM

Dual is also a convex quadratic program

$$\begin{aligned} \max_{\alpha} \quad & \sum_{n} \alpha_{n} - \frac{1}{2} \sum_{m,n} y_{m} y_{n} \alpha_{m} \alpha_{n} \mathbf{x}_{m}^{\top} \mathbf{x}_{n} \\ \text{s.t.} \quad & 0 \leq \alpha_{n} \leq C, \quad \forall \ n \\ & \sum_{n} \alpha_{n} y_{n} = 0 \end{aligned}$$

- There are N dual variables α_n , one for each data point
- Independent of the size d of x: SVM scales better for high-dimensional feature.
- May seem like a lot of optimization variables when N is large, but many of the α_n 's become zero. α_n is non-zero only if the n^{th} point is a support vector

Advantages of SVM

We've seen that the geometric formulation of SVM is equivalent to minimizing the empirical hinge loss. This explains why SVM:

- 1. Is less sensitive to outliers.
- 2. Maximizes distance of training data from the boundary
- 3. Generalizes well to many nonlinear models.
- 4. Only requires a subset of the training points.
- 5. Scales better with high-dimensional data.

The last thing left to consider is non-linear decision boundaries, or kernel SVMs

Kernel SVM

Non-linear basis functions in SVM

- What if the data is not linearly separable?
- We can transform the feature vector x using non-linear basis functions. For example,

$$\phi(\mathbf{x}) = \left[egin{array}{c} 1 \ x_1 \ x_2 \ x_1x_2 \ x_2^2 \ x_2^2 \end{array}
ight]$$

ullet Replace old x by $\phi(old x)$ in both the primal and dual SVM formulations

Primal and Dual SVM Formulations: Kernel Versions

Primal

$$\min_{\boldsymbol{w},b,\xi} \quad \frac{1}{2} \|\boldsymbol{w}\|_{2}^{2} + C \sum_{n} \xi_{n}$$
s.t. $y_{n} [\boldsymbol{w}^{\top} \phi(\boldsymbol{x}_{n}) + b] \ge 1 - \xi_{n}, \quad \forall \quad n$

$$\xi_{n} \ge 0, \quad \forall \quad n$$

Primal and Dual SVM Formulations: Kernel Versions

Primal

$$\min_{\boldsymbol{w},b,\boldsymbol{\xi}} \quad \frac{1}{2} \|\boldsymbol{w}\|_{2}^{2} + C \sum_{n} \xi_{n}$$
s.t. $y_{n} [\boldsymbol{w}^{\top} \phi(\boldsymbol{x}_{n}) + b] \ge 1 - \xi_{n}, \quad \forall \quad n$

$$\xi_{n} \ge 0, \quad \forall \quad n$$

Dual

$$\max_{\alpha} \sum_{n} \alpha_{n} - \frac{1}{2} \sum_{m,n} y_{m} y_{n} \alpha_{m} \alpha_{n} \phi(\mathbf{x}_{m})^{\top} \phi(\mathbf{x}_{n})$$
s.t. $0 \le \alpha_{n} \le C$, $\forall n$

$$\sum_{n} \alpha_{n} y_{n} = 0$$

Primal and Dual SVM Formulations: Kernel Versions

Primal

$$\min_{\boldsymbol{w},b,\boldsymbol{\xi}} \quad \frac{1}{2} \|\boldsymbol{w}\|_{2}^{2} + C \sum_{n} \xi_{n}$$
s.t. $y_{n} [\boldsymbol{w}^{\top} \phi(\boldsymbol{x}_{n}) + b] \ge 1 - \xi_{n}, \quad \forall \quad n$

$$\xi_{n} \ge 0, \quad \forall \quad n$$

Dual

$$\max_{\alpha} \sum_{n} \alpha_{n} - \frac{1}{2} \sum_{m,n} y_{m} y_{n} \alpha_{m} \alpha_{n} \phi(\mathbf{x}_{m})^{\top} \phi(\mathbf{x}_{n})$$
s.t. $0 \le \alpha_{n} \le C$, $\forall n$

$$\sum_{n} \alpha_{n} y_{n} = 0$$

IMPORTANT POINT: In the dual problem, we only need $\phi(x_m)^\top \phi(x_n)$.

Dual Kernel SVM

We replace the inner products $\phi(x_m)^{\top}\phi(x_n)$ with a kernel function

$$\max_{\alpha} \sum_{n} \alpha_{n} - \frac{1}{2} \sum_{m,n} y_{m} y_{n} \alpha_{m} \alpha_{n} k(\boldsymbol{x}_{m}, \boldsymbol{x}_{n})$$
s.t. $0 \le \alpha_{n} \le C, \quad \forall n$

$$\sum_{n} \alpha_{n} y_{n} = 0$$

Dual Kernel SVM

We replace the inner products $\phi(x_m)^{\top}\phi(x_n)$ with a kernel function

$$\max_{\alpha} \sum_{n} \alpha_{n} - \frac{1}{2} \sum_{m,n} y_{m} y_{n} \alpha_{m} \alpha_{n} k(\mathbf{x}_{m}, \mathbf{x}_{n})$$
s.t. $0 \le \alpha_{n} \le C$, $\forall n$

$$\sum_{n} \alpha_{n} y_{n} = 0$$

- $k(\mathbf{x}_m, \mathbf{x}_n)$ is a scalar and it is independent of the dimension of the feature vector $\phi(\mathbf{x})$.
- $k(\mathbf{x}_m, \mathbf{x}_n)$ roughly measures the similarity of \mathbf{x}_m and \mathbf{x}_n .
- $k(\mathbf{x}_m, \mathbf{x}_n)$ is a kernel function if it is symmetric and positive-definite $(k(\mathbf{x}, \mathbf{x}) > 0 \text{ for all } \mathbf{x} > 0)$.

We do not need to know the exact form of $\phi(x)$, which lets us define much more flexible nonlinearities.

• Dot product:

$$k(\mathbf{x}_m, \mathbf{x}_n) = \mathbf{x}_m^{\top} \mathbf{x}_n$$

We do not need to know the exact form of $\phi(x)$, which lets us define much more flexible nonlinearities.

• Dot product:

$$k(\mathbf{x}_m, \mathbf{x}_n) = \mathbf{x}_m^{\top} \mathbf{x}_n$$

• Dot product with PD matrix Q:

$$k(\mathbf{x}_m, \mathbf{x}_n) = \mathbf{x}_m^\top \mathbf{Q} \mathbf{x}_n$$

We do not need to know the exact form of $\phi(x)$, which lets us define much more flexible nonlinearities.

• Dot product:

$$k(\mathbf{x}_m, \mathbf{x}_n) = \mathbf{x}_m^\top \mathbf{x}_n$$

• Dot product with PD matrix Q:

$$k(\mathbf{x}_m, \mathbf{x}_n) = \mathbf{x}_m^\top \mathbf{Q} \mathbf{x}_n$$

· Polynomial kernels:

$$k(\mathbf{x}_m, \mathbf{x}_n) = (1 + \mathbf{x}_m^{\top} \mathbf{x}_n)^d, \quad d \in \mathbb{Z}^+$$

We do not need to know the exact form of $\phi(x)$, which lets us define much more flexible nonlinearities.

• Dot product:

$$k(\mathbf{x}_m, \mathbf{x}_n) = \mathbf{x}_m^\top \mathbf{x}_n$$

• Dot product with PD matrix Q:

$$k(\mathbf{x}_m, \mathbf{x}_n) = \mathbf{x}_m^\top \mathbf{Q} \mathbf{x}_n$$

• Polynomial kernels:

$$k(\mathbf{x}_m, \mathbf{x}_n) = (1 + \mathbf{x}_m^{\mathsf{T}} \mathbf{x}_n)^d, \quad d \in \mathbb{Z}^+$$

• Radial basis function:

$$k(\mathbf{x}_m, \mathbf{x}_n) = \exp\left(-\gamma \|\mathbf{x}_m - \mathbf{x}_n\|^2\right)$$
 for some $\gamma > 0$

and many more.

Test prediction

Learning w and b:

$$\mathbf{w} = \sum_{n} \alpha_{n} y_{n} \phi(\mathbf{x}_{n})$$
$$b = y_{n} - \mathbf{w}^{\top} \phi(\mathbf{x}_{n}) = y_{n} - \sum_{m} \alpha_{m} y_{m} k(\mathbf{x}_{m}, \mathbf{x}_{n})$$

But for test prediction on a new point \mathbf{x} , do we need the form of $\phi(\mathbf{x})$ in order to find the sign of $\mathbf{w}^{\top}\phi(\mathbf{x}) + b$?

Test prediction

Learning w and b:

$$\mathbf{w} = \sum_{n} \alpha_{n} y_{n} \phi(\mathbf{x}_{n})$$
$$b = y_{n} - \mathbf{w}^{\top} \phi(\mathbf{x}_{n}) = y_{n} - \sum_{m} \alpha_{m} y_{m} k(\mathbf{x}_{m}, \mathbf{x}_{n})$$

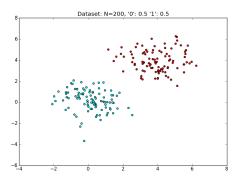
But for test prediction on a new point \mathbf{x} , do we need the form of $\phi(\mathbf{x})$ in order to find the sign of $\mathbf{w}^{\top}\phi(\mathbf{x}) + b$? Fortunately, no!

Test Prediction:

$$h(\mathbf{x}) = \text{SIGN}(\sum_{n} y_{n} \alpha_{n} k(\mathbf{x}_{n}, \mathbf{x}) + b)$$

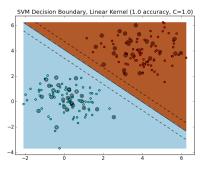
At test time it suffices to know the kernel function! So we really do not need to know ϕ .

Given a dataset $\{(x_n, y_n) \text{ for } n = 1, 2, ..., N\}$, how do you classify it using kernel SVM ?



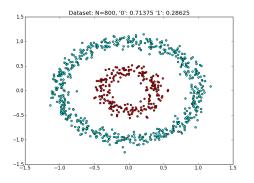
Given a dataset $\{(\mathbf{x}_n, y_n) \text{ for } n = 1, 2, ..., N\}$, how do you classify it using kernel SVM ?

Here is the decision boundary with linear soft-margin SVM



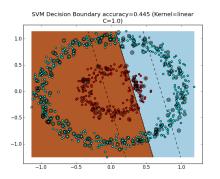
Given a dataset $\{(\mathbf{x}_n, y_n) \text{ for } n = 1, 2, ..., N\}$, how do you classify it using kernel SVM ?

What if the data is not linearly separable?



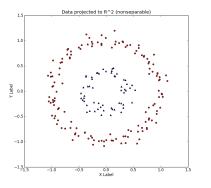
Given a dataset $\{(\mathbf{x}_n, y_n) \text{ for } n = 1, 2, ..., N\}$, how do you classify it using kernel SVM ?

The linear decision boundary is pretty bad



Given a dataset $\{(x_n, y_n) \text{ for } n = 1, 2, ..., N\}$, how do you classify it using kernel SVM ?

Use kernel $\phi(x) = [x_1, x_2, x_1^2 + x_2^2]$ to transform the data in a 3D space



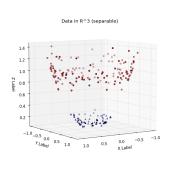


Image Source: https:

//www.eric-kim.net/eric-kim-net/posts/1/kernel_trick.html

Given a dataset $\{(x_n, y_n) \text{ for } n = 1, 2, ..., N\}$, how do you classify it using kernel SVM ?

Then find the decision boundary. How? Solve the Dual problem

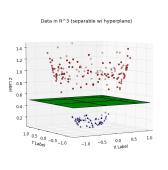
$$\max_{\alpha} \sum_{n} \alpha_{n} - \frac{1}{2} \sum_{m,n} y_{m} y_{n} \alpha_{m} \alpha_{n} \phi(\mathbf{x}_{m})^{\top} \phi(\mathbf{x}_{n})$$
s.t. $0 \le \alpha_{n} \le C$, $\forall n$

$$\sum_{n} \alpha_{n} y_{n} = 0$$

Then find **w** and *b*. Predict $y = \text{sign}(\mathbf{w}^T \phi(\mathbf{x}) + b)$.

Given a dataset $\{(x_n, y_n) \text{ for } n = 1, 2, ..., N\}$, how do you classify it using kernel SVM ?

Here is the resulting decision boundary



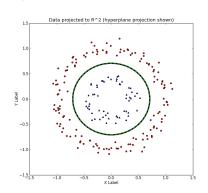


Image Source: https:

//www.eric-kim.net/eric-kim-net/posts/1/kernel_trick.html

Given a dataset $\{(x_n, y_n) \text{ for } n = 1, 2, ..., N\}$, how do you classify it using kernel SVM ?

In general, you don't need to concretely define $\phi(\mathbf{x})$. In the dual problem we can just use the kernel function $k(\mathbf{x}_m, \mathbf{x}_n)$. For cases where $\phi(\mathbf{x})$ is concretely defined, $k(\mathbf{x}_m, \mathbf{x}_n) = \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n)$.

$$\max_{\alpha} \sum_{n} \alpha_{n} - \frac{1}{2} \sum_{m,n} y_{m} y_{n} \alpha_{m} \alpha_{n} \phi(\mathbf{x}_{m})^{\top} \phi(\mathbf{x}_{n})$$
s.t. $0 \le \alpha_{n} \le C$, $\forall n$

$$\sum_{n} \alpha_{n} y_{n} = 0$$

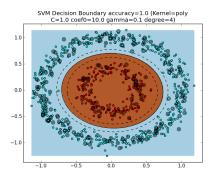
Given a dataset $\{(x_n, y_n) \text{ for } n = 1, 2, ..., N\}$, how do you classify it using kernel SVM ?

In general, you don't need to concretely define $\phi(\mathbf{x})$. In the dual problem we can just use the kernel function $k(\mathbf{x}_m, \mathbf{x}_n)$. For cases where $\phi(\mathbf{x})$ is concretely defined, $k(\mathbf{x}_m, \mathbf{x}_n) = \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n)$.

$$\begin{aligned} \max_{\alpha} \sum_{n} \alpha_{n} - \frac{1}{2} \sum_{m,n} y_{m} y_{n} \alpha_{m} \alpha_{n} k(\mathbf{x}_{m}, \mathbf{y}_{m}) \\ \text{s.t.} \quad 0 \leq \alpha_{n} \leq C, \quad \forall \ n \\ \sum_{n} \alpha_{n} y_{n} = 0 \end{aligned}$$

Given a dataset $\{(x_n, y_n) \text{ for } n = 1, 2, ..., N\}$, how do you classify it using kernel SVM ?

Effect of the choice of kernel: Polynomial kernel (degree 4)



Given a dataset $\{(x_n, y_n) \text{ for } n = 1, 2, ..., N\}$, how do you classify it using kernel SVM ?

Effect of the choice of kernel: Radial Basis Kernel

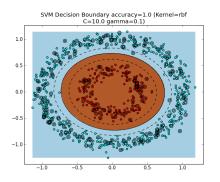


Image Source: https:

//www.eric-kim.net/eric-kim-net/posts/1/kernel_trick.html

Summary

You should know:

- Hinge loss function of SVM.
- How to derive the SVM dual.
- How to use the "kernel trick" in the dual SVM formulation to enable kernel SVM.
- How to compute an SVM prediction.