

Kernel Based Support Vector Machines



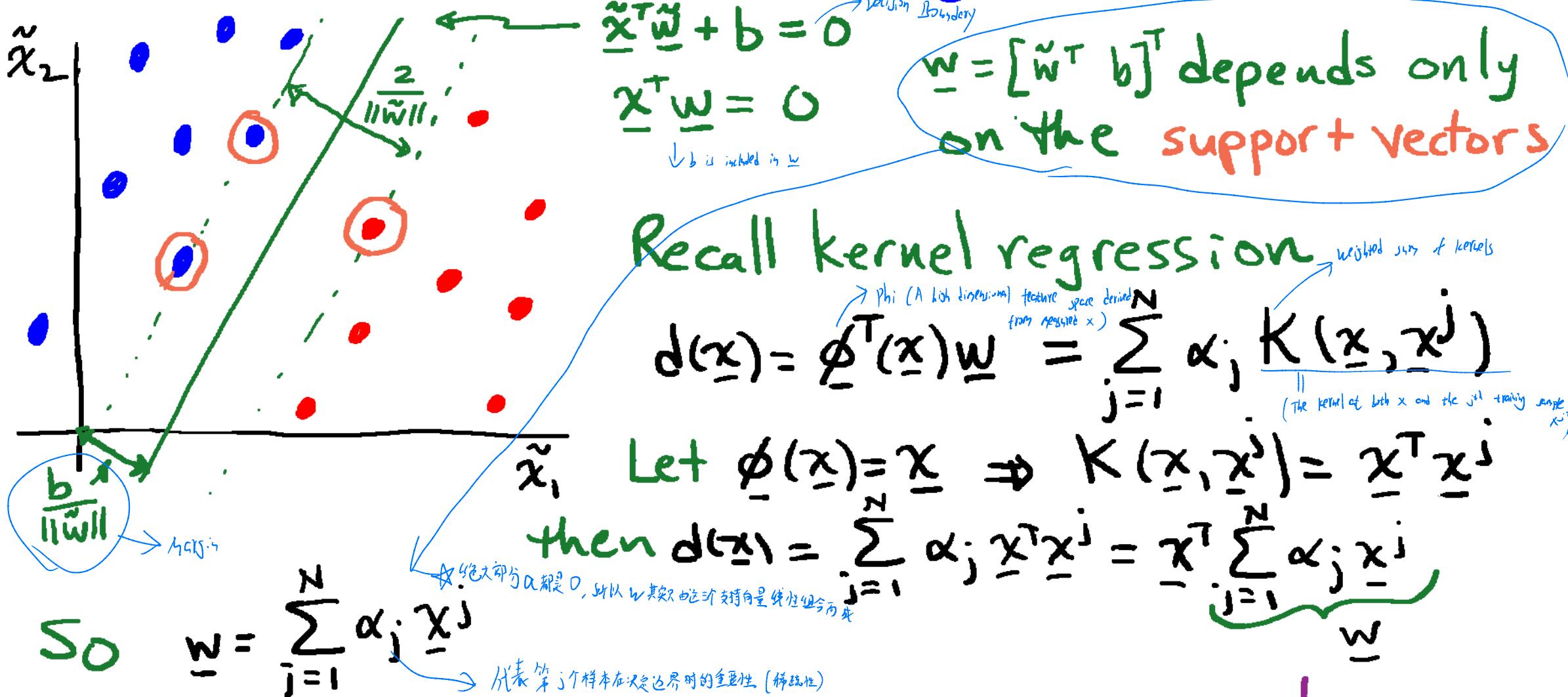
= kernel Classification

Objectives

1

- reformulate linear max margin classifier in terms of support vectors
- derive kernel version of hinge loss with ridge regression
- summarize features of support vector machines

Support vectors define max-margin classifier 2



All $\alpha_j = 0$ except Support vectors!

Use kernels for nonlinear decision boundaries³

High-dimensional feature space: $\underline{x} \rightarrow \underline{\phi}(\underline{x})$

e.g., $\underline{\phi}(\underline{x}) = [x_1^2 \ x_2^2 \ \dots \ x_2 x_4 \ \dots \ x_{n-1} \ x_n \ 1]$

$$\hat{d}(\underline{x}) = \text{sign}(\underline{\phi}^T(\underline{x}) \underline{w})$$

Hinge loss with ridge regression

$$\min_{\underline{w}} \sum_{i=1}^N (1 - d^i \underline{\phi}^T(\underline{x}^i) \underline{w})_+ + \lambda \|\underline{w}\|_2^2$$

Cost function of SVM



Claim:

$$\underline{w} = \sum_{j=1}^N \underline{\phi}(\underline{x}^j) \alpha_j$$

可以替换成 solve for α (Representer Theorem)
(proof in notes)

Kernel "trick" replaces $\underline{\phi}^T(\underline{x}^i) \underline{\phi}(\underline{x}^i)$ with $K(\underline{x}^i, \underline{x}^i)$ 4

Restate: $\underline{w} = \sum_{j=1}^N \alpha_j \underline{\phi}(\underline{x}^j)$

Hinge loss with ridge regression

$$\min_{\underline{\alpha}} \sum_{i=1}^N \left(1 - d^i \underline{\phi}^T(\underline{x}^i) \sum_{j=1}^N \alpha_j \underline{\phi}(\underline{x}^j) \right)_+ + \lambda \sum_{i=1}^N \alpha_i \underline{\phi}^T(\underline{x}^i) \sum_{j=1}^N \alpha_j \underline{\phi}(\underline{x}^j)$$

w^T w

$$\min_{\underline{\alpha}} \sum_{i=1}^N \left(1 - d^i \sum_{j=1}^N \alpha_j \underline{\phi}^T(\underline{x}^i) \underline{\phi}(\underline{x}^j) \right)_+ + \lambda \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j \underline{\phi}^T(\underline{x}^i) \underline{\phi}(\underline{x}^j)$$

$K(\underline{x}^i, \underline{x}^j)$ $K(\underline{x}^i, \underline{x}^i)$

Kernel "trick"

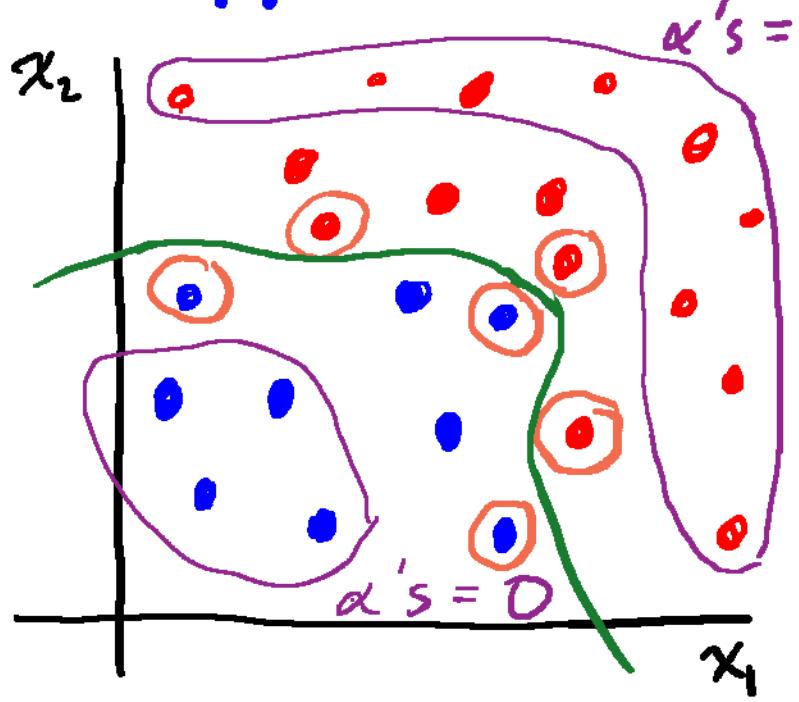
SVM

$$\min_{\underline{\alpha}} \sum_{i=1}^N \left(1 - d^i \sum_{j=1}^N \alpha_j K(\underline{x}^i, \underline{x}^j) \right)_+ + \lambda \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j K(\underline{x}^i, \underline{x}^j)$$

objective function

Support vector machines have sparse $\underline{\alpha}$

5

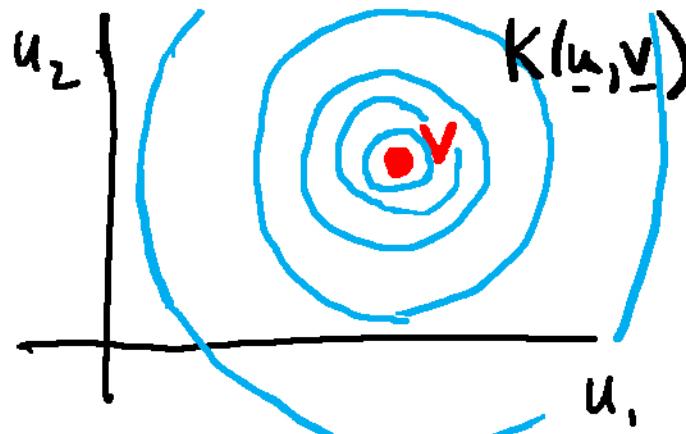


decision boundary

$$d(\underline{x}) = 0 = \underline{\phi}^T(\underline{x}) \underline{w} = \sum_{j=1}^N \alpha_j K(x, x^j)$$

Boundary (hinge loss) depends only on
the support vectors

Ex: Gaussian kernels



$$K(\underline{u}, \underline{v}) = \exp \left\{ - \frac{\|\underline{u} - \underline{v}\|_2^2}{2\sigma^2} \right\}$$

$K(\underline{u}, \underline{v})$ measures similarity/alignment
of $\underline{u}, \underline{v}$

Solve for $\underline{\alpha}$ using gradient
descent

1. 预测函数 $d(\underline{x})$

两者的底层数学形式其实是一样的，都是核函数的加权和，只是用法不同。

- 通用公式：

$$d(\underline{x}) = \sum_{i=1}^N \alpha_i K(\underline{x}, \underline{x}^i)$$

- 回归 (Regression)：直接使用 $d(\underline{x})$ 的值作为预测结果（例如预测房价）。
- 分类 (SVM)：使用 $\text{sign}(d(\underline{x}))$ 作为分类结果。
 - 如果 $d(\underline{x}) > 0 \rightarrow$ 预测为 +1 类。
 - 如果 $d(\underline{x}) < 0 \rightarrow$ 预测为 -1 类。

2. 损失函数 (Cost Function)

这是两者本质的区别，决定了模型是“钻牛角尖”还是“只求稳妥”。

- 回归用的：平方误差 (Squared Error)

$$(y_{true} - y_{pred})^2$$

- 笔记：它要求预测值必须无限接近真实值。
- SVM 用的：合页损失 (Hinge Loss)

$$(1 - y_{true} \cdot y_{pred})_+$$

- 笔记： $(z)_+$ 的意思是 $\max(0, z)$ 。
- 它的意思是：只要你分类正确且信心足够大（即 $y \cdot d(\underline{x}) > 1$ ），损失就是 0，我就不惩罚你了。这让 SVM 具有很强的鲁棒性。

3. 目标函数 (Objective Function)

这就是我们要让电脑去“最小化”的总公式。

- 核回归 (Kernel Ridge Regression)：

$$\min_{\underline{w}} \sum_{i=1}^N (d^i - \underline{\phi}^T(\underline{x}^i)\underline{w})^2 + \lambda \|\underline{w}\|^2$$

- 核分类 (Kernel SVM)：

$$\min_{\underline{w}} \sum_{i=1}^N (1 - d^i \underline{\phi}^T(\underline{x}^i)\underline{w})_+ + \lambda \|\underline{w}\|^2$$

4. 权重 \underline{w} 与稀疏性

根据表示定理， \underline{w} 总是数据的线性组合： $\underline{w} = \sum \alpha_j \phi(\underline{x}^j)$ 。

- 回归： α_j 通常都不为 0。因为很难有样本的预测误差恰好完全等于 0，所以每个样本都会对模型产生一点拉力。
- SVM： α_j 大部分为 0。只要样本被正确分类且在边界之外（安全区），它的 Hinge Loss 就是 0，对应的 α 也就是 0。只有那些在边界上的支持向量， α 才不为 0。

**Copyright 2019
Barry Van Veen**

Kernel Based Support Vector Machines

Proof: Weights Lie in Space Spanned by $\phi(\mathbf{x}^j)$

©Barry Van Veen 2019

Background: Classifier training features and labels are $\mathbf{x}^i, d^i, i = 1, 2, \dots, N$. Classification in a high-dimensional feature space is performed using the mapping $\phi(\mathbf{x})$ as $\hat{d}(\mathbf{x}) = \text{sign}\{\phi^T(\mathbf{x})\mathbf{w}\}$.

Claim: The weights \mathbf{w} that satisfy

$$\min_{\mathbf{w}} f(\mathbf{w}) = \min_{\mathbf{w}} \sum_{i=1}^N (1 - d^i \phi^T(\mathbf{x}^i) \mathbf{w})_+ + \lambda \|\mathbf{w}\|_2^2 \quad (1)$$

are of the form

$$\mathbf{w} = \sum_{j=1}^N \phi(\mathbf{x}^j) \alpha_j$$

Proof: The proof proceeds by adding a component to \mathbf{w} that is orthogonal to the space spanned by the vectors $\{\phi(\mathbf{x}^j), j = 1, 2, \dots, N\}$ and then showing that component must be zero at the minimum of Eq. 1.

Suppose

$$\mathbf{w} = \sum_{j=1}^N \phi(\mathbf{x}^j) \alpha_j + \phi^\perp$$

where ϕ^\perp is orthogonal to the space spanned by $\{\phi(\mathbf{x}^j), j = 1, 2, \dots, N\}$, that is, $\phi^T(\mathbf{x}^j) \phi^\perp = 0, j = 1, 2, \dots, N\}$. Note that any vector \mathbf{w} can be expressed as a sum of a component in the space spanned by the $\phi(\mathbf{x}^j)$ and a component orthogonal to that same space.

The optimization problem Eq. 1 may be rewritten

$$\begin{aligned} \min_{\mathbf{w}} f(\mathbf{w}) &= \min_{\alpha, \phi^\perp} \sum_{i=1}^N \left(1 - d^i \phi^T(\mathbf{x}^i) \left(\sum_{j=1}^N \phi(\mathbf{x}^j) \alpha_j + \phi^\perp \right) \right)_+ + \lambda \left\| \sum_{j=1}^N \phi(\mathbf{x}^j) \alpha_j + \phi^\perp \right\|_2^2 \end{aligned} \quad (2)$$

$$\begin{aligned} &= \min_{\alpha, \phi^\perp} \sum_{i=1}^N \left(1 - d^i \left(\sum_{j=1}^N \phi^T(\mathbf{x}^i) \phi(\mathbf{x}^j) \alpha_j + \phi^T(\mathbf{x}^i) \phi^\perp \right) \right)_+ \\ &\quad + \lambda \left(\sum_{i=1}^N \sum_{j=1}^N \phi^T(\mathbf{x}^i) \phi(\mathbf{x}^j) \alpha_i \alpha_j + 2 \sum_{i=1}^N \alpha_i \phi^T(\mathbf{x}^i) \phi^\perp + \phi^{\perp T} \phi^\perp \right) \end{aligned} \quad (3)$$

where in the second line we have used the identity

$$\left\| \sum_{j=1}^N \phi(\mathbf{x}^j) \alpha_j + \phi^\perp \right\|_2^2 = \left(\sum_{i=1}^N \phi(\mathbf{x}^i) \alpha_i + \phi^\perp \right)^T \left(\sum_{j=1}^N \phi(\mathbf{x}^j) \alpha_j + \phi^\perp \right)$$

and multiplied out the terms in the product.

Now use the fact that $\phi^T(\mathbf{x}^i)\phi^\perp = 0$ to reexpress the optimization problem as

$$\min_{\mathbf{w}} f(\mathbf{w}) = \min_{\boldsymbol{\alpha}, \phi^\perp} \sum_{i=1}^N \left(1 - d^i \sum_{j=1}^N \phi^T(\mathbf{x}^i) \phi(\mathbf{x}^j) \alpha_j \right)_+ + \lambda \sum_{i=1}^N \sum_{j=1}^N \phi^T(\mathbf{x}^i) \phi(\mathbf{x}^j) \alpha_i \alpha_j + \lambda \phi^{\perp T} \phi^\perp \quad (4)$$

$$= \min_{\boldsymbol{\alpha}, \phi^\perp} \sum_{i=1}^N \left(1 - d^i \sum_{j=1}^N \phi^T(\mathbf{x}^i) \phi(\mathbf{x}^j) \alpha_j \right)_+ + \lambda \left\| \sum_{j=1}^N \phi(\mathbf{x}^j) \alpha_j \right\|_2^2 + \lambda \|\phi^\perp\|_2^2 \quad (5)$$

The only term containing ϕ^\perp is the last one, $\lambda \|\phi^\perp\|_2^2$, which is nonnegative since $\lambda > 0$. Consequently, we conclude the minimum is attained when $\phi^\perp = \mathbf{0}$ and thus

$$\mathbf{w} = \sum_{j=1}^N \phi(\mathbf{x}^j) \alpha_j$$