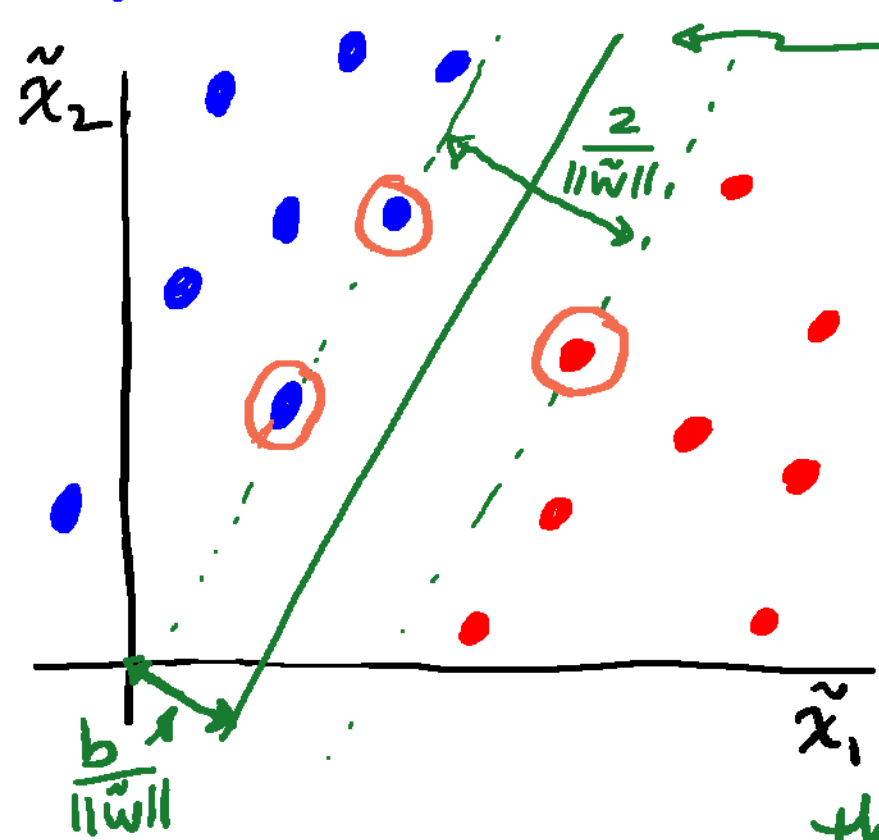


Kernel Based Support Vector Machines

Objectives

- 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



$$\tilde{\underline{x}}^T \tilde{\underline{w}} + b = 0$$

$$\underline{x}^T \underline{w} = 0$$

$\underline{w} = [\tilde{w}^T \ b]^T$ depends only on the support vectors

Recall kernel regression

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

Let $\phi(\underline{x}) = \underline{x} \Rightarrow K(\underline{x}, \underline{x}^j) = \underline{x}^T \underline{x}^j$

then $d(\underline{x}) = \sum_{j=1}^N \alpha_j \underline{x}^T \underline{x}^j = \underline{x}^T \underbrace{\sum_{j=1}^N \alpha_j \underline{x}^j}_{\underline{w}}$

So $\underline{w} = \sum_{j=1}^N \alpha_j \underline{x}^j$

All $\alpha_j = 0$ except support vectors!

Use kernels for nonlinear decision boundaries³

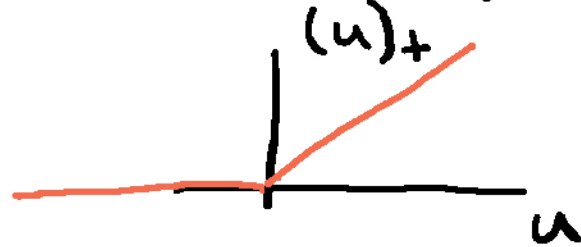
High-dimensional feature space: $\underline{x} \rightarrow \underline{\phi}(\underline{x})$

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

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

Hinge loss with ridge regression

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



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

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

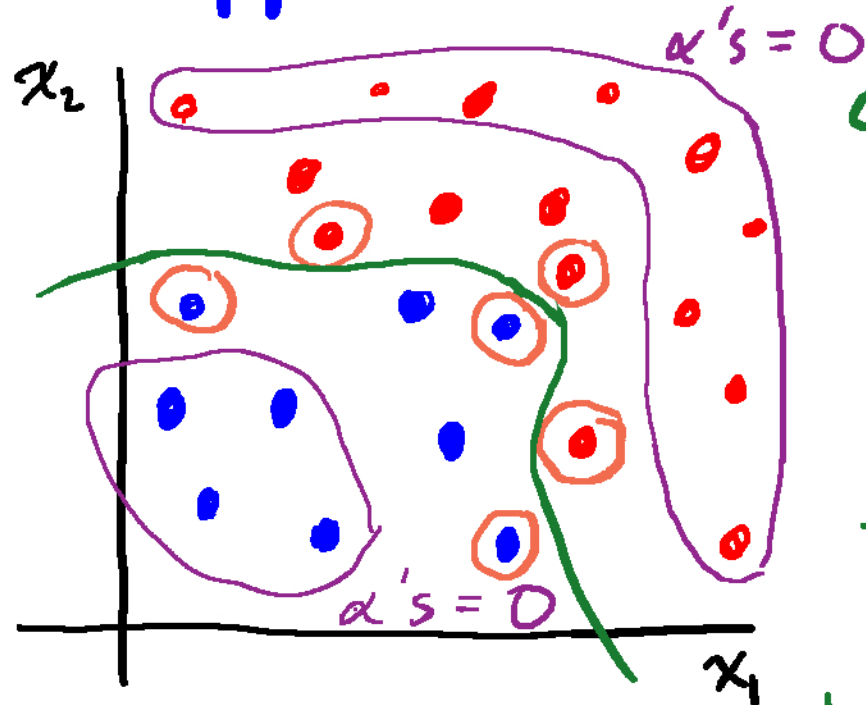
$$\min_{\underline{\alpha}} \sum_{i=1}^N \left(1 - d^i \sum_{j=1}^N \alpha_j \underbrace{\underline{\phi}^T(\underline{x}^i) \underline{\phi}(\underline{x}^j)}_{K(\underline{x}^i, \underline{x}^j)} \right)_+ + \lambda \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j \underbrace{\underline{\phi}^T(\underline{x}^i)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

Kernel "trick"

SVM

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

Support vector machines have sparse $\underline{\alpha}$ 5



decision boundary

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

Boundary (hinge loss) depends only on the support vectors

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

Ex: Gaussian kernels



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

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

Copyright 2019
Barry Van Veen