

• How can we do SVM for non-linearly boundaries?

$\begin{matrix} X & X & X & X \\ 0 & 0 & 0 & X \\ 0 & 0 & 0 & X \\ & & X & X \end{matrix}$ } No hyperplane does a very good job of separating the classes... if we could use a curved boundary, then maybe things would be possible.

• Recall the kernel trick from our discussion of kernel PCA: we want to

find a feature map $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^D$ that discriminates in the way we want, i.e. separates the classes.

ex: For example above, $\phi(x, y) = (x, y, x - y^2)$ would probably work, since the two classes lie on different parabolas.

• If D is big, i.e. we need to embed in a large number of dimensions, things may be slow computationally. To get around this, we never consider ϕ explicitly but rather by consider a kernel $K(x, y) = \phi(x)^T \phi(y)$, which allows to compute the distances according to ϕ as:

Exercise:
$$\|\phi(x) - \phi(y)\|_2 = \sqrt{K(x, x) - 2K(x, y) + K(y, y)}$$

The kernel $K(x,y) = \phi(x)^T \phi(y)$ is symmetric ($K(x,y) = K(y,x)$) and non-negative ($K(x,x) \geq 0$ for all x), and positive semidefinite ($\sum_{i,j=1}^n \alpha_i \alpha_j K(x_i, x_j) \geq 0$ for all $\{\alpha_i\}_{i=1}^n \subset \mathbb{R}, \{x_i\}_{i=1}^n \subset \mathbb{R}^d$)

Exercise: Positive semidefiniteness as above is related to matrix positive semidefiniteness as

K is psd \Leftrightarrow the matrix $(K(x_i, x_j))_{i,j=1, \dots, n}$ is psd as a matrix.

In fact, these properties completely characterize kernels!

Theorem (Moore-Aronszajn): Let ~~kernel~~ $K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ be symmetric and psd. Then $K(x,y) = \langle \phi(x), \phi(y) \rangle$ for some $\phi: \mathbb{R}^d \rightarrow \mathcal{F}$, where \mathcal{F} is a (possibly infinite dimensional) feature space.

Common kernel is the Gaussian kernel $K(x,y) = \exp(-\|x-y\|_2^2 / 2\sigma^2)$, some $\sigma > 0$.

Remark: The \mathcal{F} in the MAT may be infinite dimensional, so we don't want to think about writing out ϕ explicitly. Better to use the kernel trick and only ~~construct~~ K .

So, how to incorporate kernels into SVM? Let's replace x_i with $\phi(x_i)$, and see if we can ultimately bring the kernel $K(x_i, x_j) = \phi(x_i)^T \phi(x_j)$ into the picture.

Hard Margin SVM: original space

feature space

$$(w^*, b^*) = \arg \min_w \|w\|_2^2 \text{ s.t.}$$

$$(w^*, b^*) = \arg \min_w \|w\|_2^2 \text{ s.t.}$$

$$y_i (w^T x_i + b) \geq 1, i=1, \dots, n$$

$$y_i (w^T \phi(x_i) + b) \geq 1, i=1, \dots, n$$

Exercise: If (w^*, b^*) solve $(*)$, then $w^* = \sum_{i=1}^n \alpha_i \phi(x_i)$, some $\{\alpha_i\}_{i=1}^n \in \mathbb{R}$.

Hence, $(*)$ may be re-cast as

$$(\alpha_1^*, \dots, \alpha_n^*, b^*) = \arg \min_{(\alpha_1, \dots, \alpha_n, b)} \left\| \sum_{j=1}^n \alpha_j \phi(x_j) \right\|_2^2 \text{ s.t.}$$

$(**)$

$$y_i (w^T [\sum_{j=1}^n \alpha_j \phi(x_j)] \phi(x_i) + b) \geq 1, i=1, \dots, n.$$

Let $K = (K(x_i, x_j))_{i,j=1, \dots, n} \in \mathbb{R}^{n \times n}$ and $\alpha = (\alpha_1, \dots, \alpha_n) \in \mathbb{R}^{n \times 1}$. Then

we may write $(**)$ compactly as

$$(\alpha_1^*, \dots, \alpha_n^*, b^*) = \arg \min_{(\alpha, b)} \alpha^T K \alpha \text{ s.t.}$$

$$-(K\alpha)_i - b \leq -1, y_i = 1$$

$$-(K\alpha)_i + b \leq -1 \text{ if } y_i = -1$$

Once the $\{\alpha_i\}_{i=1}^n, b$ are learned, new points ^{call it x} are labeled as

+1 if $\sum_{i=1}^n \alpha_i K(x_i, x) + b > 0$

-1 if $\sum_{i=1}^n \alpha_i K(x_i, x) + b < 0$



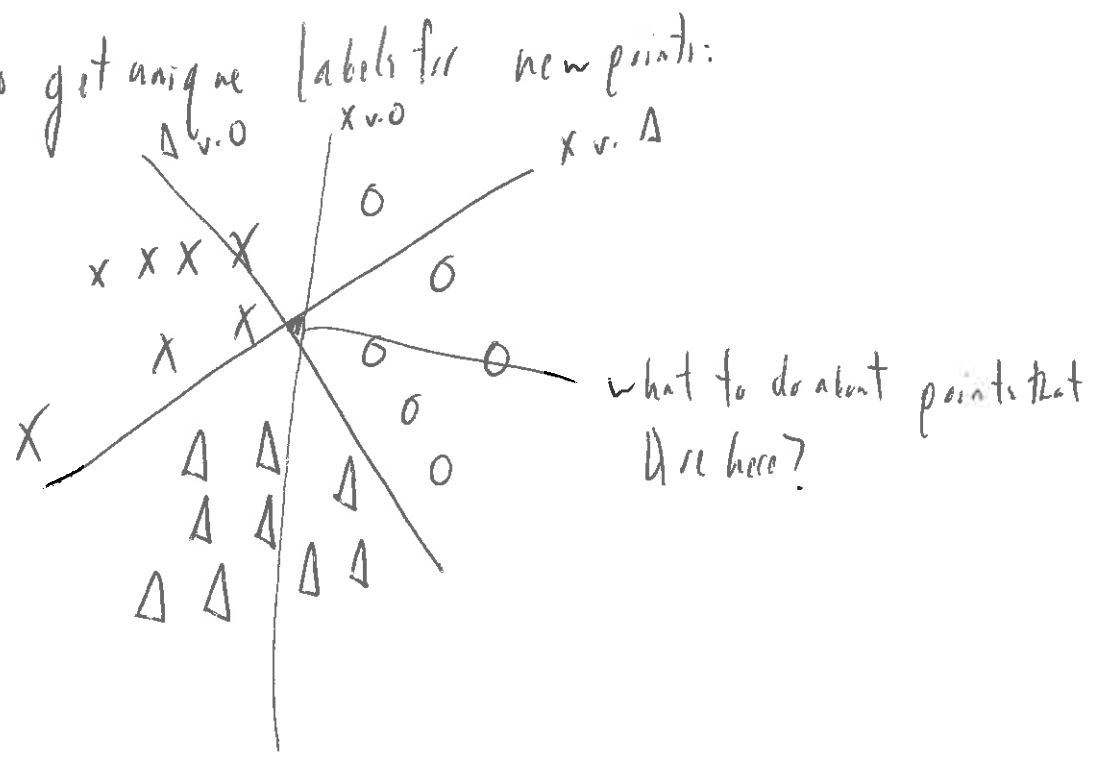
What if we have multiple classes we want to separate with SVM?

Natural approach: for each of the K classes, train a classifier against each other class. This is the so-called one-versus-one SVM.

More precisely, let $k_1, k_2 \in \{1, \dots, K\}$ be distinct labels. Build a classifier SVM by handling the two class problem with data $\{x_i \mid y_i = k_1 \text{ or } k_2\}$.

Problem: hard to get unique labels for new points:

ex (3 class)



- An alternative approach is to do a one-versus-all classification scheme for each class, then give a point the label that maximizes the margin among the K -one versus all classifiers.

- More precisely, for each $k \in \{1, \dots, K\}$, let $C_1 = \{x_i \mid y_i = k\}$
 $C_2 = \{x_i \mid y_i \neq k\}$.

Train an SVM to distinguish C_1 from C_2 , and let w_k, b_k be the corresponding hyperplane parameters. This yields a family of SVM with parameters $\{(w_k, b_k)\}_{k=1}^K$.

- One then chooses the class that maximizes margin.

$$y(x) = \arg \max_{k=1, \dots, K} w_k^T X + b_k$$

- Of course, this would all work for kernel SVM as well.

- So far, we have seen two ways to approach classification:
 - 1) K -NN algorithms focus on local similarities, but maybe poor if local similarities are not enough and global structure is needed.

2.7 SVM: fit a hyperplane or warped hyperplane between classes. (6)

- These methods are rather principled, easy to understand in terms of decisions made, and rely on a small number of parameters.
- We will end the course with neural networks / deep learning, which are harder to understand and have a lot of parameters, but can work amazingly well.