- So far, we only saw linear discriminative classifiers (SVM, LR, perceptron, LSC, ...)

- What if the data is separable, but not linearly separable?



$\Phi = \begin{bmatrix} x \\ x^2 \end{bmatrix}$

$\Phi = x_1 x_2$

$z = x_1 x_2$

$\Phi = $ high-dim xfm

- **IDEA** - transform the input space
  - mapping: $\Phi : X \to Z$
  - learn linear classifier in new space $(Z)$

- if the new dimension is large enough, the data should be separable by a hyperplane.
  $$\dim(Z) > \dim(X)$$

- in the limit, $\dim(Z) \to \infty$
  - we are mapping a vector $x_i \to$ into a function
  $$\Phi(x) = \begin{bmatrix} \phi_1(x) \\ \vdots \\ \phi_B(x) \end{bmatrix} \to \phi(x; \xi) \quad \Bigg| \quad x \to e^{-\frac{1}{2}(z-x)^2}$$

## Kernel SVM

- consider the SVM dual problem → replace $x_i$ with $\Phi(x_i)$

**Training**
$$\max_{\alpha} \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \underbrace{\Phi(x_i)^T \Phi(x_j)}_{K(x_i, x_j)}$$
$$\text{s.t.} \sum_i \alpha_i y_i = 0, \quad \alpha_i \geq 0 \; \forall i$$

**Decision boundary**
$$w^* = \sum_i \alpha_i y_i \Phi(x_i) \quad \leftarrow \; w \text{ has same dimension as } \Phi(\cdot)$$

$SV = \{x_i \mid \alpha_i > 0\}$
(support vector set)

**Bias term**
$$b^* = \frac{1}{|SV|} \sum_{i \in SV} (y_i - w^T \Phi(x_i))$$
$$= \frac{1}{|SV|} \sum_{i \in SV} \left( y_i - \sum_j \alpha_j y_j \underbrace{\Phi(x_j)^T \Phi(x_i)}_{K(x_i, x_j)} \right)$$

**Decision function**
$$y_* = \text{sign}(f(x_*))$$
$$f(x_*) = w^T \Phi(x_*) + b = \sum_i \alpha_i y_i \underbrace{\Phi(x_i)^T \Phi(x_*)}_{K(x_i, x_*)} + b$$

**Note:** The whole algorithm only depends on the data through $\Phi(x_i)^T \Phi(x_j)$ !!!

- Define a kernel function: (dot-product kernel)
$$K(x_i, x_j) = \Phi(x_i)^T \Phi(x_j) = \langle \Phi(x_i), \Phi(x_j) \rangle$$

- The dual SVM can be written w/ this kernel function → __nonlinear classifier__.

---

- Actually, we <u>only</u> need the kernel function, not the x form $\Phi(x)$.

- Just define the kernel ⇒ called the "kernel trick"!

__Why is it good?__

  - Save the calculation of $\Phi(x_i)$ (could be a large vector)

  - Need to calculate $O(n^2)$ terms of $K(x_i, x_j)$

  ⇒ kernel matrix:

$$K = \begin{bmatrix} K(x_1,x_1) & \cdots & K(x_1,x_n) \\ \vdots & \ddots & \vdots \\ K(x_n,x_1) & \cdots & K(x_n,x_n) \end{bmatrix}$$

  (Gram matrix)

Example: polynomial kernel

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} \in \mathbb{R}^d$$

$$K(x,x') = (x^Tx')^2 = \left(\sum_{i=1}^{d} x_i x_i'\right)^2 = \sum_{i=1}^{d}\sum_{j=1}^{d} x_i x_i' x_j x_j'$$

$$= \sum_i \sum_j (x_i x_j)(x_i' x_j')$$

$$= [x_1x_1 \quad x_1x_2 \quad \cdots\cdots \quad x_dx_d] \begin{bmatrix} x_1'x_1' \\ x_1'x_2' \\ \vdots \\ x_d'x_d' \end{bmatrix} \quad \Phi(x')$$

$\underbrace{\phantom{[x_1x_1 \quad x_1x_2]}}_{\Phi(x)}$

$$= \Phi(x)^T \Phi(x')$$

$$\Phi: \mathbb{R}^d \to \mathbb{R}^{d^2} \quad \to O(d^2)$$

$$K(x,x') = (x^Tx')^2 \quad \to O(d) \qquad \leftarrow \text{more efficient than explicitly computing } \Phi(x)$$

$\underbrace{\phantom{(x^Tx')}}_{\substack{d \\ d+1}}$

- What about the classifier?

$$f(x_*) = \sum_i \alpha_i y_i\, K(x_i, x_*) + b = \sum_i \alpha_i y_i (x_i^T x_*)^2 + b$$

$$= \sum_i \alpha_i y_i (x_*^T x_i x_i^T x_*) + b = x_*^T \left(\underbrace{\sum_i \alpha_i y_i x_i x_i^T}_{A}\right) x_* + b$$

$$= x_*^T A x_* + b \quad \leftarrow \text{quadratic function in } x_*.$$

✱ the kernel function specifies the <u>class</u> of nonlinear functions that are learned.

# Kernel Functions

"Kernel trick" depends on $K(x,x')$ being a dot-product kernel.

**Define:** mapping $K: \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is a <u>dot-product</u> kernel iff

$$K(x,x') = \langle \Phi(x), \Phi(x') \rangle$$

where $\Phi: \mathcal{X} \to \mathcal{H}$ (vector space)

$\langle \cdot, \cdot \rangle$ is a dot product in $\mathcal{H}$.

How to check if $K(x,x')$ is a dot product kernel w/o knowing $\Phi$, $\langle \cdot, \cdot \rangle$?

**Define:** $K(x,x')$ is a <u>positive-definite</u> kernel on $\mathcal{X} \times \mathcal{X}$ iff

$\forall n$ and $\forall \{x_1, \dots, x_n\}$, $x_i \in \mathcal{X}$,

the kernel matrix $K = [k(x_i, x_j)]_{ij}$ is a <u>positive</u> definite matrix. <span style="color:red">(for all possible datasets, the kernel matrix is posdef.)</span>

$K$ is posdef <u>iff</u>
1) $y^T K y > 0 \quad \forall y$
2) eigenvalues of $K$ are positive
3) $K = A A^T$, $A$ has independent columns.

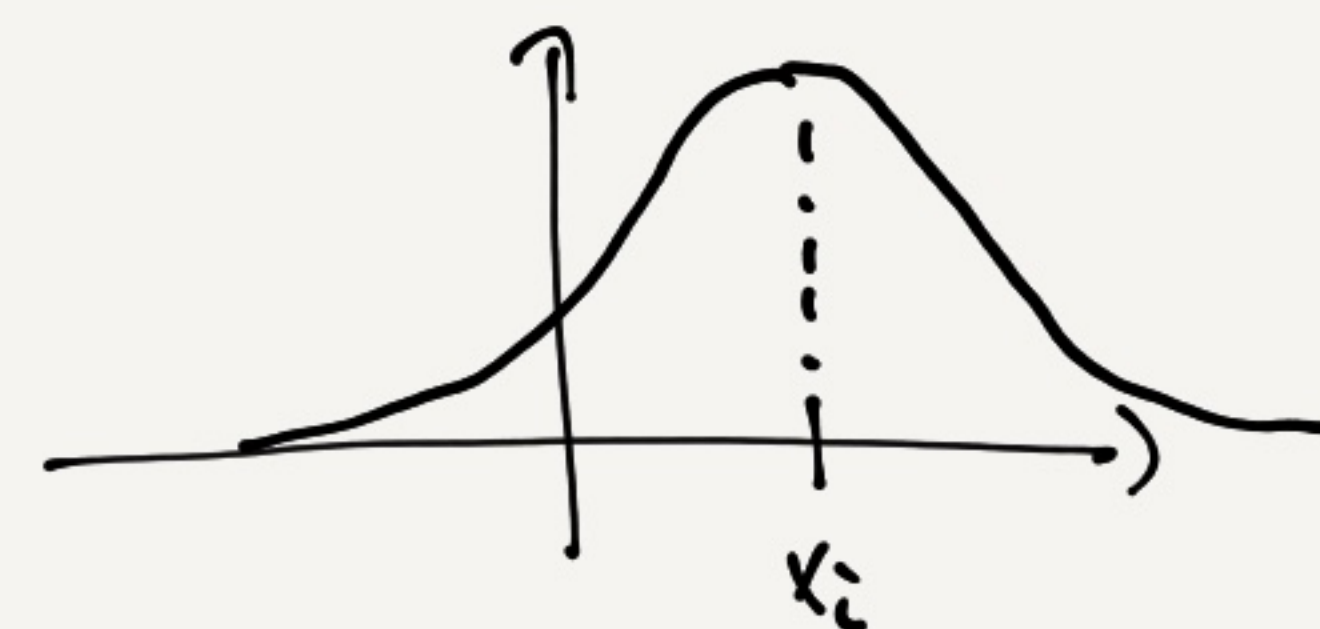☒☒☒ $\boxed{K(x,x') \text{ is a dotproduct kernel } \underline{\text{iff}} \text{ it is a posdef kernel.}}$

---

Given a posdef/dot product kernel, what is the high-dim xformation $\Phi$?

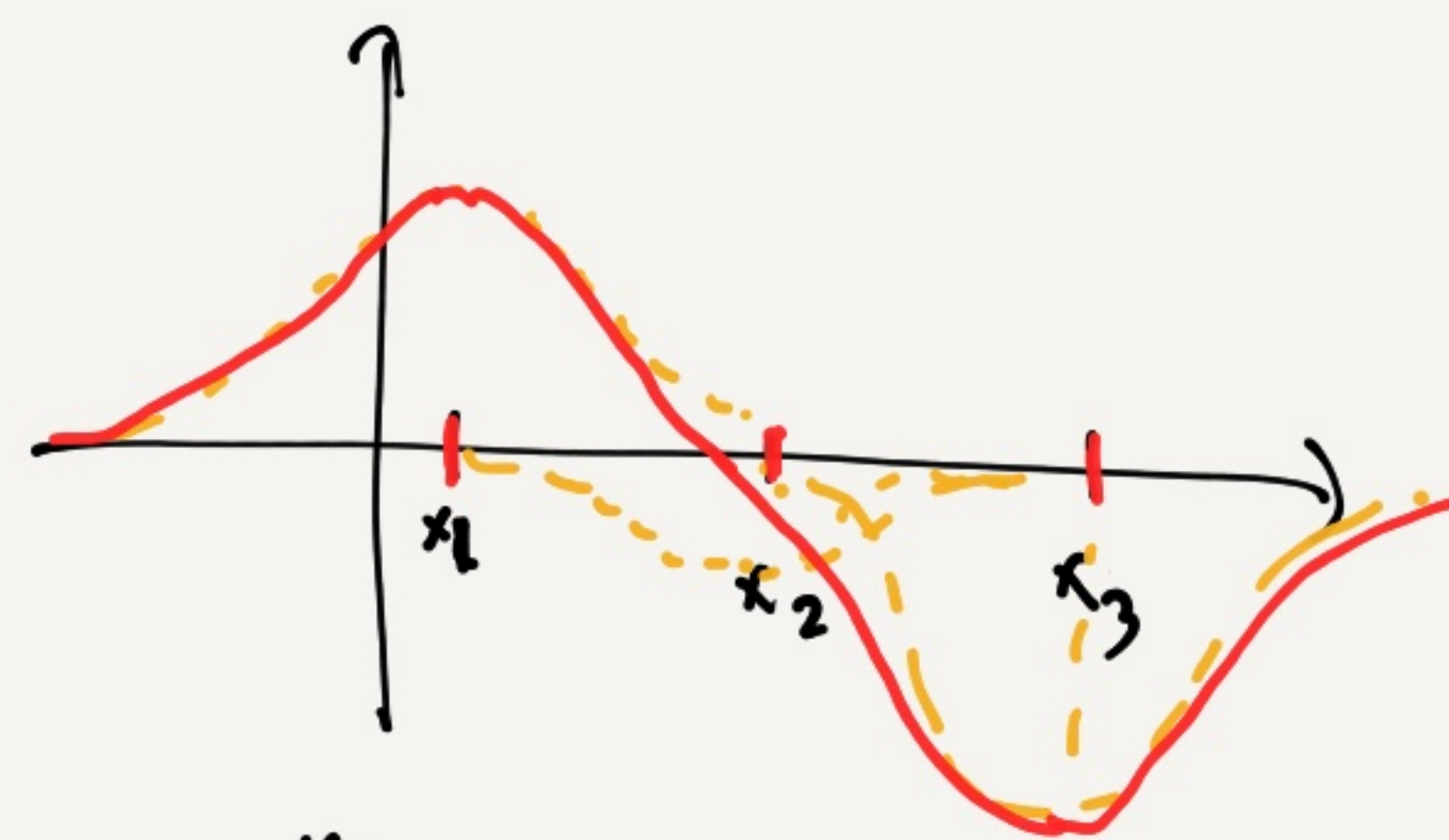let $\mathcal{H}$ = space of all linear combinations of function $K(\cdot, x_i)$.

<span style="color:red">function of $\cdot$    $x_i$ is fixed</span>

$$\mathcal{H} = \left\{ f(\cdot) \mid f(\cdot) = \sum_{i=1}^{m} \alpha_i K(\cdot, x_i), \ \forall m, \forall x_i \in \mathcal{X} \right\}$$

e.g. use a Gaussian kernel:
$$K(\cdot, x_i) = e^{-\frac{1}{\sigma^2} \| \cdot - x_i \|^2}$$

$f(\cdot)$ = combinations of Gaussians

• Let $f(\cdot) = \sum_{i=1}^{m} \alpha_i K(\cdot, x_i)$, $\quad g(\cdot) = \sum_j \beta_j K(\cdot, x_j)$

We can show the dot-product btwn $f, g$ is

$$\langle f, g \rangle = \sum_i \sum_j \alpha_i \beta_j K(x_i, x_j)$$

e.g. Gaussian kernel.
$$\langle f, g \rangle = \sum_i \sum_j \alpha_i \beta_j \, e^{-\frac{1}{\sigma^2} \| x_i - x_j \|^2}$$

## Special case

$$\left.\begin{array}{l} d_i = 1, \alpha_{x_i} = 0 \\ b_j = 1, \beta_{y_j} = 0 \end{array}\right\} \Rightarrow \langle \underbrace{k(\cdot, x_i)}_{\Phi(x_i)}, \underbrace{k(\cdot, x_j)}_{\Phi(x_j)} \rangle = K(x_i, x_j)$$

## Hence

$$K(x_i, x_j) = \langle \overline{\Phi(x_i)}, \overline{\Phi(x_j)} \rangle$$

where $\Phi: \mathcal{X} \to \mathcal{H}$

$$x_i \to \Phi(x_i) = k(\cdot, x_i)$$

*transformation from vector $x$ to a function (inf. dim space)*

## Final Note

given $f(\cdot) = \sum_i \alpha_i K(\cdot, x_i)$

then $\langle k(\cdot, x), f(\cdot) \rangle = \sum_i \alpha_i K(x_i, x) = f(x)$

"reproducing property" - dot product of $f$ w/ kernel gives back the $f$.
(similar to dirac delta & convolution)

$\mathcal{H}$ is called a "Reproducing Kernel Hilbert Space." (RKHS)
(Hilbert space = vector space & dot product & ....)

---

- RKHS uniquely specifies the kernel function & vice versa.

- Posdef kernels are also called _Mercer kernels_

## Representer Theorem

Empirical Risk: $R_{emp} = \sum_{i=1}^{n} L(y_i, f(x_i))$

Regularizer: $\Omega(\|f\|_p)$, $\Omega \geq 0$ & strictly monotonically increasing.

$\mathcal{H}, k(x, x')$

$$f^* = \underset{f}{\arg\min} \; R_{emp}(f) + \lambda \, \Omega(\|f\|_p)$$

then $f^*$ has the form:

$$f^* = \sum_i \alpha_i k(x, x_i), \quad f^* \in \mathcal{H} \text{ (RKHS)}$$

( Many ML algorithms fit this framework
→ similar form of $f$
→ they are kernelizable (make it nonlinear)
→ inf. dim space of $f$ → finite dim $\alpha$. )

# Kernel functions

## Kernels on $\mathbb{R}^d$

linear: $k(x,x') = x^T x'$

$k(x,x') = x^T A x'$ , $A$ posdef

poly: $k(x,x') = (x^T x' + c)^d$

$k(x,x') = (x^T x')^d$

Gaussian: $k(x,x') = e^{-\frac{1}{\gamma}\|x-x'\|^2}$ ← Radial Basis Function (RBF)

Expo: $k(x,x') = e^{-\frac{1}{\gamma}\|x-x'\|}$

- can combine kernels

$k(x,x') = \underbrace{x^T x'}_{\text{linear}} + \underbrace{e^{-\frac{1}{\gamma}\|x-x'\|^2}}_{\text{Gaussian (RBF)}}$

see PS for the rules.

## Kernel on histograms $x, x'$ are histograms

correlation kernel: $k(x,x') = x^T x' = \sum_{i=1}^{d} x_i x_i'$

Bhattacharyya kernel: $k(x,x') = \sum_{i=1}^{d} \sqrt{x_i} \sqrt{x_i'}$

$\chi$-squared kernel: $k(x,x') = e^{-\frac{1}{\sigma^2}\chi^2(x,x')}$
  RBF

$\chi^2(x,x') = \sum_i \frac{(x_i - x_i')^2}{\frac{1}{2}(x_i + x_i')}$

histogram intersection kernel: $k(x,x') = \sum_{i=1}^{d} \min(x_i, x_i')$



## Kernels on sets

$X = \{x_1, \ldots, x_n\}$

$X' = \{x_1', \ldots, x_m'\}$

# of common elements

$|X \cap X'|$

intersection kernel: $k(X,X') = 2$

distance kernel: $k(X,X') = e^{-\frac{1}{\gamma}\sum_i \sum_j d(x_i, x_j')}$

pyramid match kernel:

## Kernels on strings/trees/graphs

$k(x,x') = \sum_S w_S \phi_S(x) \phi_S(x')$

$w_S > 0$    # of times substring $S$ appears in $x'$.

※ **Kernels** allow us to learn classifiers directly on non-vector data.

$\Rightarrow$ need to define the appropriate PD kernel.