

# Introduction to Machine Learning

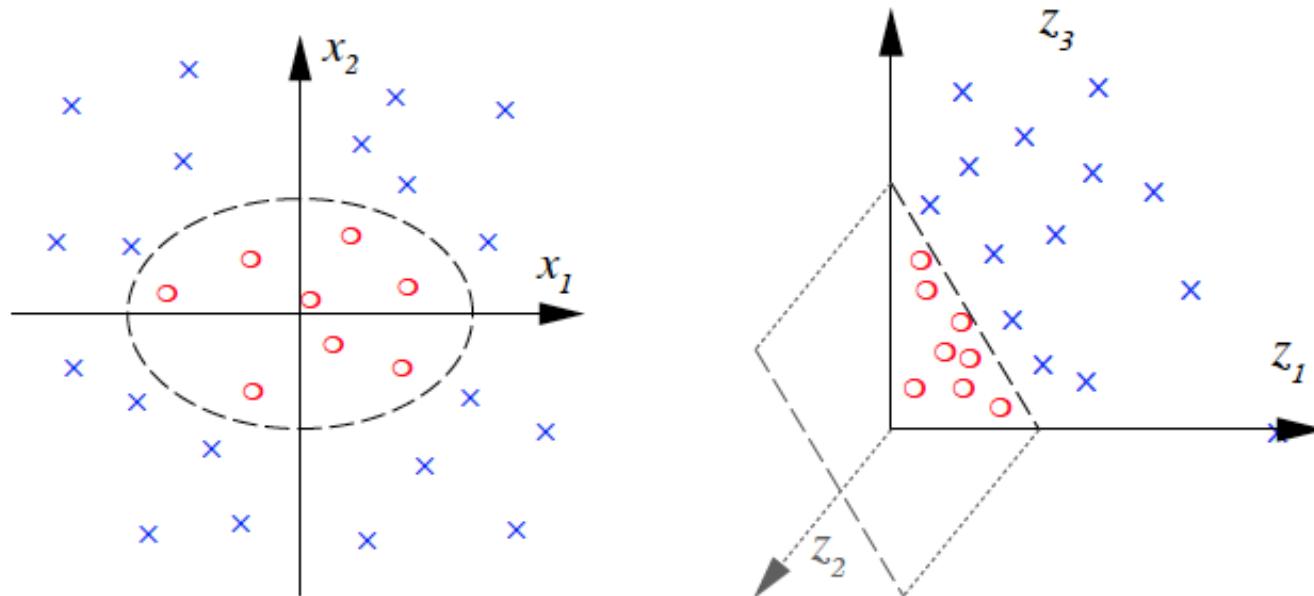
Non-linear prediction with kernels (continued)

Prof. Andreas Krause  
Learning and Adaptive Systems ([las.ethz.ch](http://las.ethz.ch))

# Solving non-linear classification tasks

- How can we find nonlinear classification boundaries?
- Similar as in regression, can use **non-linear transformations** of the feature vectors, followed by linear classification

$$\begin{aligned}\Phi : \mathbb{R}^2 &\rightarrow \mathbb{R}^3 \\ (x_1, x_2) &\mapsto (z_1, z_2, z_3) := (x_1^2, \sqrt{2}x_1x_2, x_2^2)\end{aligned}$$



# Avoiding the feature explosion

---

- Need  $O(d^k)$  dimensions to represent (multivariate) polynomials of degree  $k$  on  $d$  features
- **Example:**  $d=10000, k=2 \rightarrow$  Need  $\sim 100M$  dimensions
- In the following, we can see how we can efficiently **implicitly** operate in such high-dimensional feature spaces (i.e., without ever explicitly computing the transformation)

# The „Kernel Trick“

- Express problem s.t. it only depends on inner products
- Replace inner products by kernels

$$\mathbf{x}_i^T \mathbf{x}_j \quad \Rightarrow \quad k(\mathbf{x}_i, \mathbf{x}_j)$$

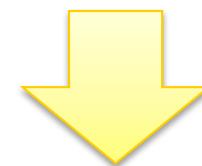
- This „trick“ is very widely applicable!

# The „Kernel Trick“

- Express problem s.t. it only depends on inner products
- Replace inner products by kernels

- Example: Perceptron

$$\hat{\alpha} = \arg \min_{\alpha_{1:n}} \frac{1}{n} \sum_{i=1}^n \max \left\{ 0, - \sum_{j=1}^n \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \right\}$$



$$\hat{\alpha} = \arg \min_{\alpha_{1:n}} \frac{1}{n} \sum_{i=1}^n \max \left\{ 0, - \sum_{j=1}^n \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j) \right\}$$

- Will see further examples later

# Kernelized Perceptron

Training

- Initialize  $\alpha_1 = \dots = \alpha_n = 0$
- For  $t=1,2,\dots$ 
  - Pick data point  $(\mathbf{x}_i, y_i)$  uniformly at random
  - Predict
$$\hat{y} = \text{sign}\left(\sum_{j=1}^n \alpha_j y_j k(\mathbf{x}_j, \mathbf{x}_i)\right)$$
  - If  $\hat{y} \neq y_i$  set  $\alpha_i \leftarrow \alpha_i + \eta_t$

Prediction

- For new point  $\mathbf{x}$ , predict

$$\hat{y} = \text{sign}\left(\sum_{j=1}^n \alpha_j y_j k(\mathbf{x}_j, \mathbf{x})\right)$$

# Properties of kernel functions

- Data space  $X$
- A kernel is a function  $k : X \times X \rightarrow \mathbb{R}$
- Can we use any function?

$$\exists \phi : X \rightarrow \mathbb{R}^D \text{ s.t. } k(x, x') = \phi(x)^T \phi(x')$$

- $k$  must be an **inner product** in a suitable space  
→  $k$  must be **symmetric**!

$$\forall x, x' \in X : k(x, x') = \phi(x)^T \phi(x') = \phi(x')^T \phi(x) = k(x', x)$$

→ Are there other properties that it must satisfy?

# Positive semi-definite matrices

Symmetric matrix  $M \in \mathbb{R}^{n \times n}$  is positive semidefinite iff

- $\equiv$
- (i)  $\forall x \in \mathbb{R}^n : x^T M x \geq 0$
  - (ii) All eigenvalues of  $M \geq 0$

---

(i)  $\Rightarrow$  (ii) :  $M$  is symmetric  $\Rightarrow M = U D U^T$  for  $D = \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{pmatrix}$ ,  $U^T U = I = U U^T$   
 $U = (u_1 | \dots | u_n)$  s.t.  $M u_i = \lambda_i u_i$

WtP:  $\lambda_i \geq 0$  &

$$u_i^T M u_i = u_i^T (\lambda_i u_i) = \lambda_i u_i^T u_i = \lambda_i \stackrel{\text{by (i)}}{\geq} 0$$

D

# Kernels → semi-definite matrices

- Data space  $X$  (possibly infinite)
- Kernel function  $k : X \times X \rightarrow \mathbb{R}$
- Take any finite subset of data  $S = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq X$
- Then the **kernel (gram) matrix**

$$\mathbf{K} = \begin{pmatrix} k(\underline{\mathbf{x}_1}, \mathbf{x}_1) & \dots & k(\mathbf{x}_1, \mathbf{x}_n) \\ \vdots & & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & \dots & k(\mathbf{x}_n, \mathbf{x}_n) \end{pmatrix} = \begin{pmatrix} \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_1) & \dots & \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_n) \\ \vdots & & \vdots \\ \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_1) & \dots & \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_n) \end{pmatrix}$$

is positive semidefinite

$$k = \phi^T \phi, \text{ where } \phi = \begin{pmatrix} \phi(x_1) \\ \vdots \\ \phi(x_n) \end{pmatrix} \in \mathbb{R}^{D \times n}$$

$$\text{Sps. } x \in \mathbb{R}^n : \text{ Then } x^T k x = \underbrace{x^T \phi^T}_{v^T} \underbrace{\phi x}_{v} = v^T v \geq 0$$

□

# Semi-definite matrices $\rightarrow$ kernels

- Suppose the data space  $X = \{1, \dots, n\}$  is finite, and we are given a positive semidefinite matrix  $K \in \mathbb{R}^{n \times n}$
- Then we can always construct a feature map

$$\phi : X \rightarrow \mathbb{R}^n$$

such that  $K_{i,j} = \phi(i)^T \phi(j)$

$K$  is s.p.d.  $\Rightarrow K = UDU^T$ , where  $D = \begin{pmatrix} \lambda_1 & 0 \\ 0 & \ddots & 0 \\ & & \lambda_n \end{pmatrix}$

and  $\lambda_i \geq 0 \forall i$

$$D = D^{\frac{1}{2}} D^{\frac{1}{2}T}, \text{ where } D^{\frac{1}{2}} = \begin{pmatrix} \sqrt{\lambda_1} & 0 \\ 0 & \ddots & 0 \\ & & \sqrt{\lambda_n} \end{pmatrix}$$

~~Sp.~~  $A \subset \{1, \dots, n\}$

$K^A$  kernel matrix for set  $A$ . Then

$$K^A = K_{AA}$$

$$\begin{aligned} \phi : X &\rightarrow \mathbb{R}^n \\ \phi : i &\mapsto \phi_i \end{aligned}$$

$$\Rightarrow K = \underbrace{U}_{\Phi^T} \underbrace{D^{\frac{1}{2}}} \underbrace{D^{\frac{1}{2}T}}_{\Phi} U^T, \text{ where } \Phi = [\phi_1 | \dots | \phi_n]$$

Now it holds that  $k(i, j) = K_{i,j} = \phi_i^T \phi_j$



# Outlook: Mercer's Theorem

Let  $X$  be a compact subset of  $\mathbb{R}^d$  and  $k : X \times X \rightarrow \mathbb{R}$  a **kernel function**

Then one can expand  $k$  in a uniformly convergent series of bounded functions  $\phi_i$  s.t.

$$k(x, x') = \sum_{i=1}^{\infty} \lambda_i \phi_i(x) \phi_i(x')$$

Can be generalized even further

# Definition: kernel functions

- Data space  $X$
- A **kernel** is a function  $k : X \times X \rightarrow \mathbb{R}$  satisfying
- 1) **Symmetry**: For any  $\mathbf{x}, \mathbf{x}' \in X$  it must hold that

$$k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}', \mathbf{x})$$

- 2) **Positive semi-definiteness**: For any  $n$ , any set  $S = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq X$ , the kernel (Gram) matrix

$$\mathbf{K} = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \dots & k(\mathbf{x}_1, \mathbf{x}_n) \\ \vdots & & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & \dots & k(\mathbf{x}_n, \mathbf{x}_n) \end{pmatrix}$$

must be positive semi-definite

# Examples of kernels on $\mathbb{R}^d$

- Linear kernel:

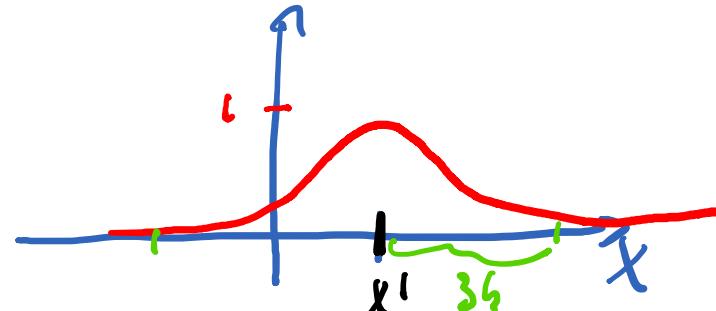
$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$$

- Polynomial kernel:

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + 1)^d$$

- Gaussian (RBF,  
squared exp. kernel):

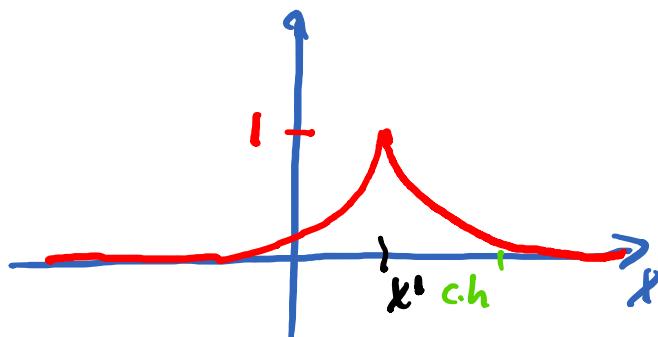
$$k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|_2^2 / h^2)$$



Bandwidth/  
Lengthscale  
parameter

- Laplacian kernel:

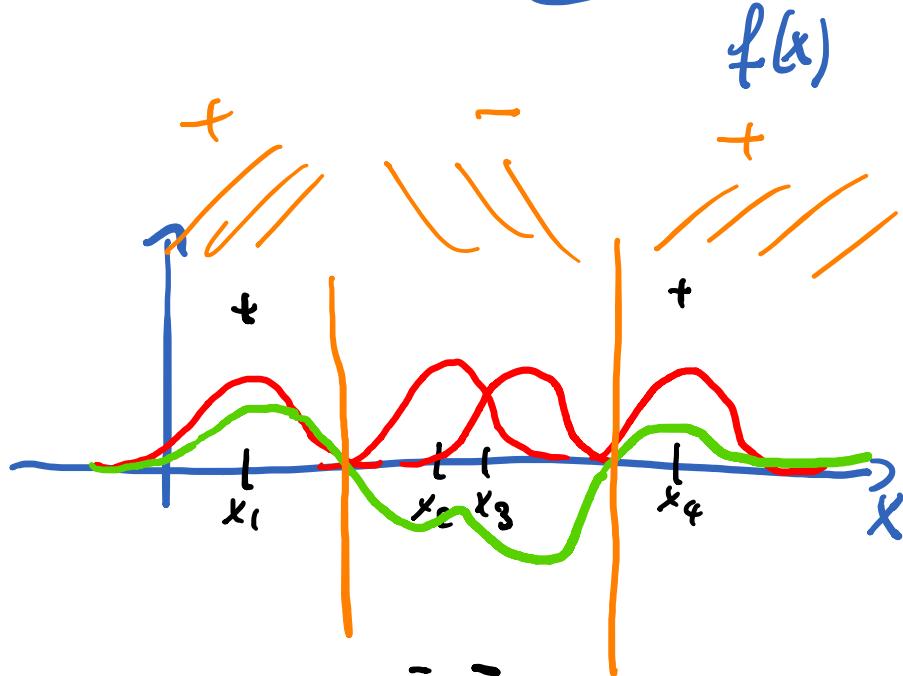
$$k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|_1 / h)$$



# Effect of kernel on function class

- Given kernel  $k$ , predictors (for kernelized classification) have the form

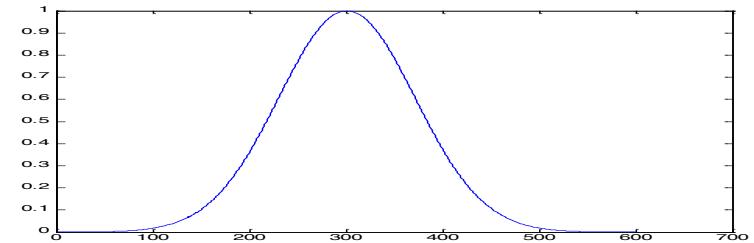
$$\hat{y} = \text{sign}\left(\sum_{j=1}^n \alpha_j y_j \overbrace{k(\mathbf{x}_j, \mathbf{x})}^{\text{f(x)}}\right)$$



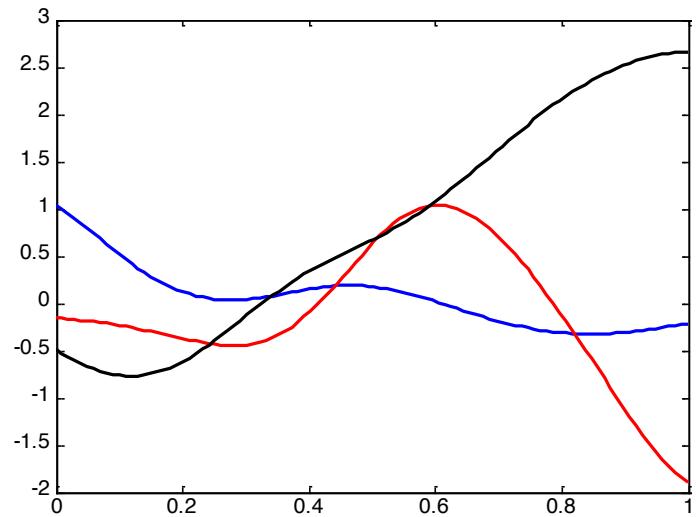
$$D = \{(x_1, +), (x_2, -), (x_3, -), (x_4, +)\}$$

# Example: Gaussian kernel

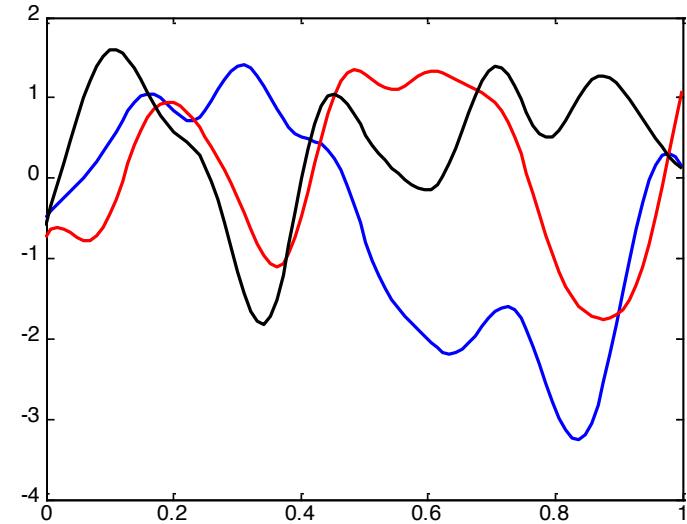
$$k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|_2^2/h^2)$$



$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \mathbf{x})$$



Bandwidth  $h=.3$



Bandwidth  $h=.1$

# Examples of (non)-kernels

$$k(x, x') = \sin(x) \cos(x')$$

Not symmetric!  $k(x, x') \neq k(x', x)$

e.g. for  $x=0, x' = \frac{\pi}{2}$

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T M \mathbf{x}'$$

$k$  symmetric  $\Leftrightarrow M$  is symmetric positive

$k$  positive (semi-)definite  $\Leftrightarrow M$  is (semi-)definite

(D): Sps.  $M = -I$ . Then  $k(x, x') = -x x'$ , hence,  $k(x, x) = -x^2 \Rightarrow k$  not pos. def.

Sps.  $M$  is pos. def. Then  $M = U D^{\frac{1}{2}} D^{\frac{1}{2}T} U^T = V^T V$ . Then

$$k(x, x') = \mathbf{x}^T M \mathbf{x}' = \mathbf{x}^T V^T V \mathbf{x}' = (V\mathbf{x})^T (V\mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}') \text{ for } \phi(\mathbf{x}) = V\mathbf{x}$$

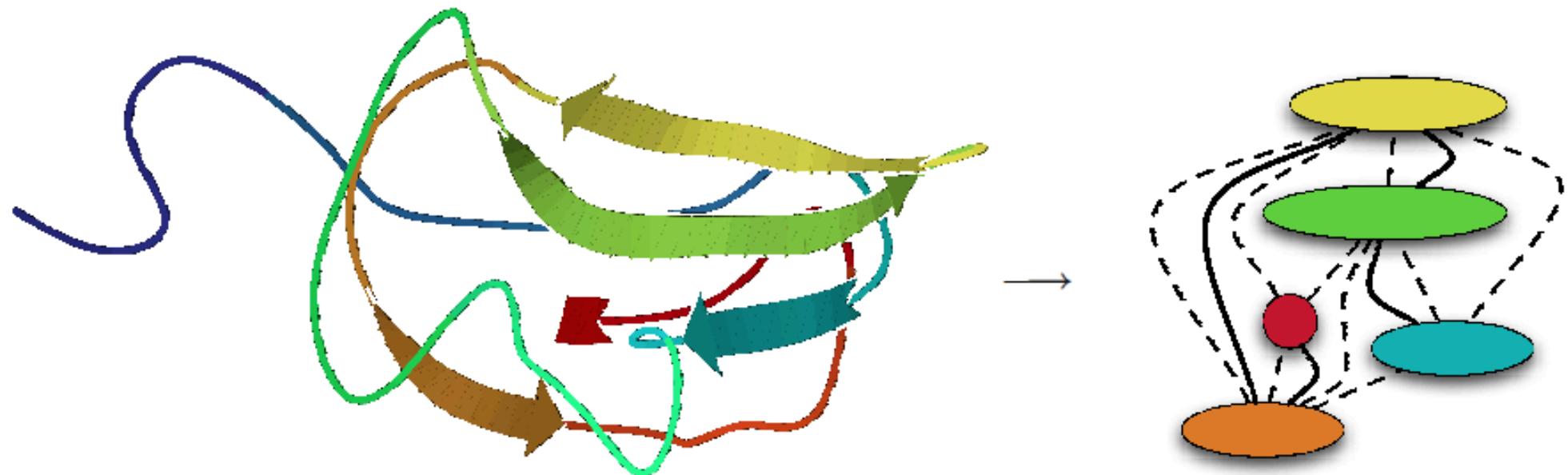
□

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

---

- Can define kernels on a variety of objects:
- Sequence kernels
- Graph kernels
- Diffusion kernels
- Kernels on probability distributions
- ...

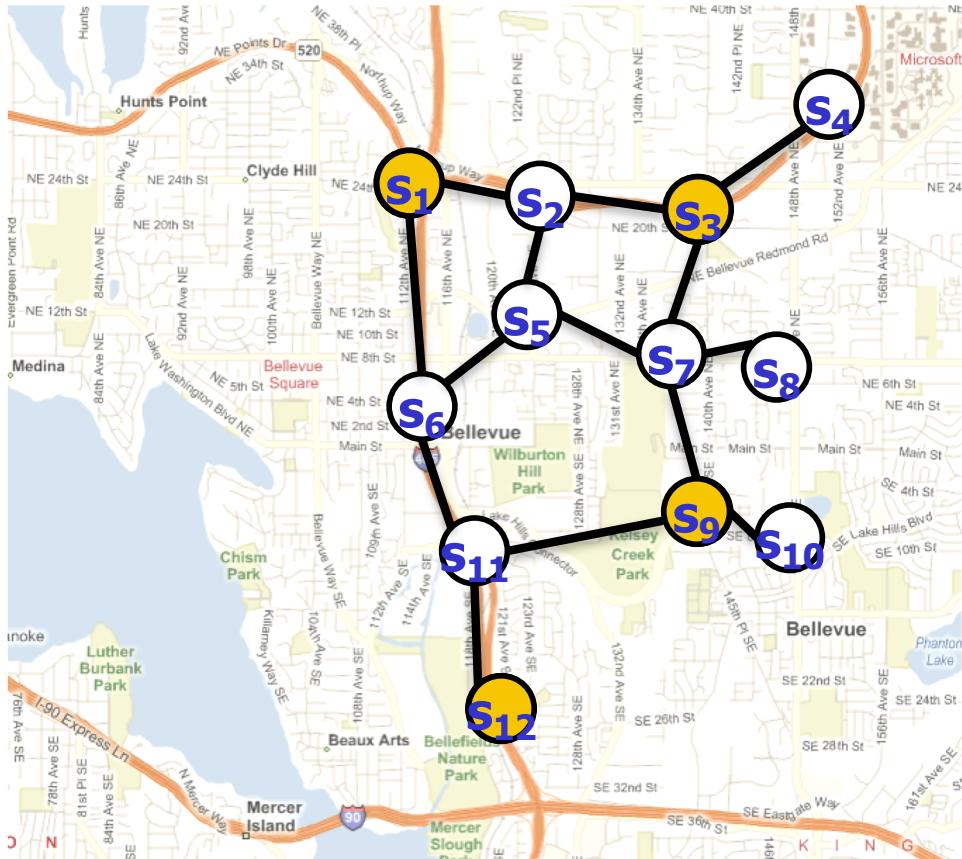
# Example: Graph kernels



[Borgwardt et al.]

- Can define a kernel for measuring similarity between graphs by comparing random walks on both graphs (not further defined here)

# Example: Diffusion kernels on graphs



$$K = \exp(-\beta L)$$

- Can measure similarity among nodes in a graph via diffusion kernels (not defined here)

# Kernel engineering (composition rules)

- Suppose we have two kernels

$$k_1 : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$$

$$k_2 : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$$

defined on data space  $X$

- Then the following functions are valid kernels:

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}')$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') k_2(\mathbf{x}, \mathbf{x}')$$

$$k(\mathbf{x}, \mathbf{x}') = c k_1(\mathbf{x}, \mathbf{x}') \text{ for } c > 0$$

$$k(\mathbf{x}, \mathbf{x}') = f(k_1(\mathbf{x}, \mathbf{x}'))$$

Sps.  $\mathcal{V} : \mathcal{Z} \rightarrow \mathcal{X}$

$k(\mathbf{z}, \mathbf{z}') = k_1(\mathcal{V}(\mathbf{z}), \mathcal{V}(\mathbf{z}'))$   
is valid kernel.

Pf:  $k_1(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$

$k(\mathbf{z}, \mathbf{z}') = \phi(\mathcal{V}(\mathbf{z}))^T \phi(\mathcal{V}(\mathbf{z}'))$

where  $f$  is a polynomial with positive coefficients or  
the exponential function

# Example: ANOVA kernel

$$k(x, x') = \sum_{j=1}^d k_j(x_j, x'_j)$$

where  $x \in \mathbb{R}^d$ ,  $k: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$   
 $k_j: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$

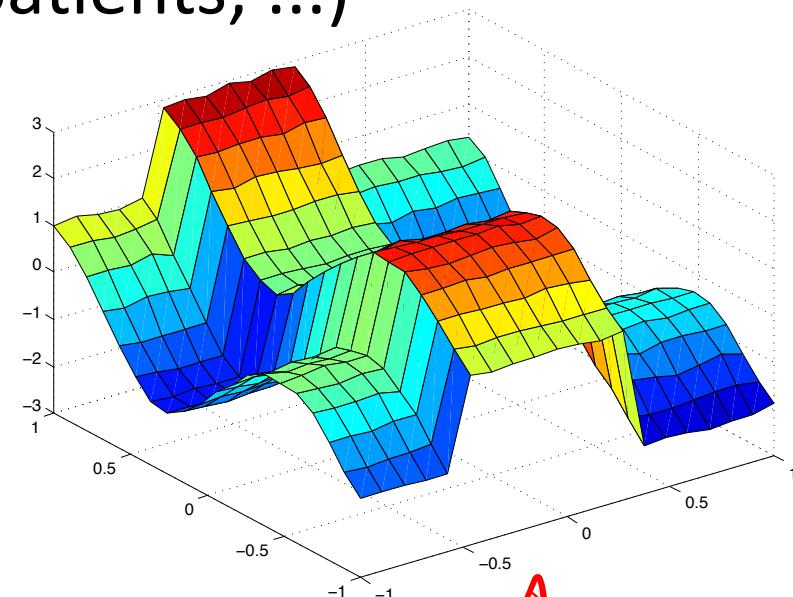
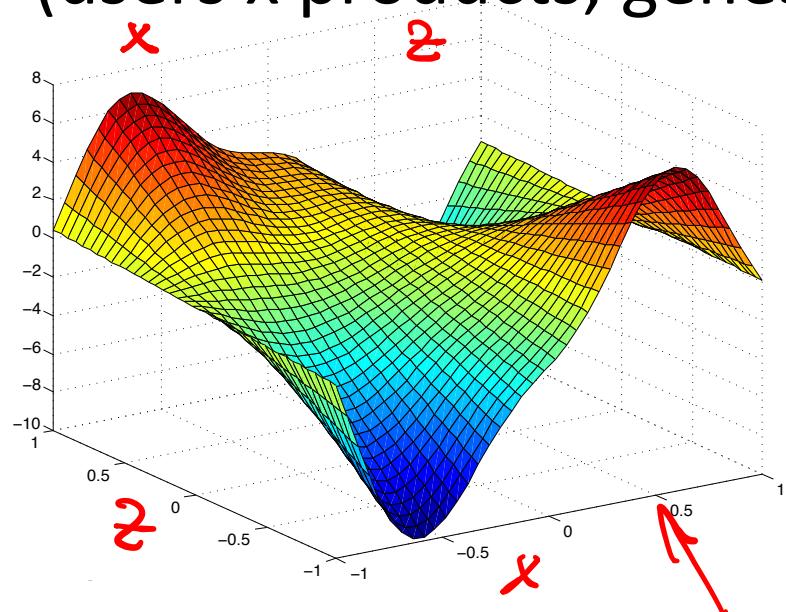
e.g.  $k_j(x_i, x'_i) = \exp\left(-\frac{(x_i - x'_i)^2}{h_j^2}\right)$

What functions does  $k$  model?  
↙  
ith data point

$$\begin{aligned} f(x) &= \sum_{i=1}^n \alpha_i k(x^{(i)}, \cdot) \\ &= \sum_{i=1}^n \alpha_i \sum_{j=1}^d k_j(x_j^{(i)}, x'^{(i)}_j) \\ &= \sum_{j=1}^d \underbrace{\sum_{i=1}^n \alpha_i k_j(x_j^{(i)}, x'^{(i)}_j)}_{f_j(x_j)} \end{aligned}$$

# Example: Modeling pairwise data

- May want to use kernels to model pairwise data  
(users x products; genes x patients; ...)



$$k_{\pi}((x, z), (x', z')) := k_x(x, x') \cdot k_z(z, z')$$

$$k_x(x, x') + k_z(z, z')$$

# Where are we?

---

- We've seen how to kernelize the perceptron
- Discussed properties of kernels, and seen examples
- Next questions:
  - What kind of predictors / decision boundaries do kernel methods entail?
  - Can we use the kernel trick beyond the perceptron?

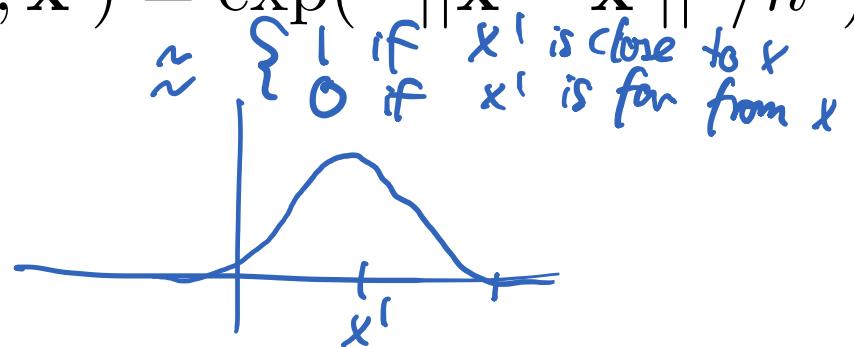
# Kernels as *similarity functions*

- Recall Perceptron (and SVM) classification rule:

$$y = \text{sign} \left( \sum_{i=1}^n \alpha_i y_i k(\mathbf{x}_i, \mathbf{x}) \right) \approx \text{sign} \left( \sum_{i=1}^n \alpha_i y_i [x \text{ "close" to } x_i] \right)$$

$\underbrace{\sum_{i=1}^n \alpha_i y_i k(\mathbf{x}_i, \mathbf{x})}_{f(x)}$

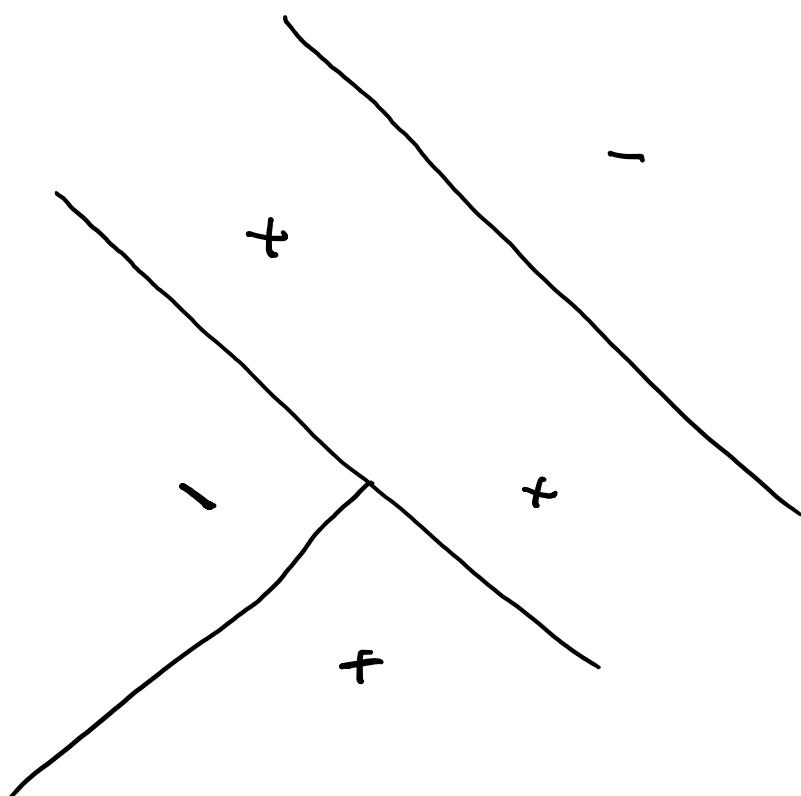
- Consider Gaussian kernel  $k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|^2/h^2)$



# Side note: Nearest-neighbor classifiers

- For data point  $\mathbf{x}$ , predict majority of labels of  $k$  nearest neighbors

$$y = \text{sign} \left( \sum_{i=1}^n y_i [\mathbf{x}_i \text{ among } k \text{ nearest neighbors of } \mathbf{x}] \right)$$



# Demo: k-NN

---

-

# Nearest-neighbor classifiers

- For data point  $\mathbf{x}$ , predict majority of labels of  $k$  nearest neighbors

$$y = \text{sign} \left( \sum_{i=1}^n y_i [\mathbf{x}_i \text{ among } k \text{ nearest neighbors of } \mathbf{x}] \right)$$

- How to choose  $k$ ?
  - Cross-validation! 😊

# K-NN vs. Kernel Perceptron

- k-Nearest Neighbor:

$$y = \text{sign} \left( \sum_{i=1}^n y_i [\mathbf{x}_i \text{ among } k \text{ nearest neighbors of } \mathbf{x}] \right)$$

- Kernel Perceptron:

$$y = \text{sign} \left( \sum_{i=1}^n y_i \alpha_i k(\mathbf{x}_i, \mathbf{x}) \right)$$

# Comparison: k-NN vs Kernelized Perceptron

<i>Method</i>	<i>k</i> -NN	<i>Kernelized Perceptron</i>
<b>Advantages</b>	No training necessary	Optimized weights can lead to improved performance Can capture „global trends“ with suitable kernels Depends on „wrongly classified“ examples only
<b>Disadvantages</b>	Depends on all data → inefficient	Training requires optimization