

November 20, 2024



Inter-University Research Institute Corporation /
Research Organization of Information and Systems

National Institute of Informatics

Machine Learning for Graph Structured Data

Introduction to Big Data Science (ビッグデータ概論)

Mahito Sugiyama (杉山磨人)

Example of Learning from Data

(from mlss.tuebingen.mpg.de/2013/schoelkopf_whatismL_slides.pdf)

- 1, 2, 4, 7,... → What are succeeding numbers?

Example of Learning from Data

(from mlss.tuebingen.mpg.de/2013/schoelkopf_whatismL_slides.pdf)

- 1, 2, 4, 7, ... → What are succeeding numbers?

$$1, 2, 4, 7, 11, 16, \dots \quad (a_n = a_{n-1} + n - 1)$$

Example of Learning from Data

(from mlss.tuebingen.mpg.de/2013/schoelkopf_whatismL_slides.pdf)

- 1, 2, 4, 7, ... → What are succeeding numbers?

$$1, 2, 4, 7, 11, 16, \dots \quad (a_n = a_{n-1} + n - 1)$$

$$1, 2, 4, 7, 12, 20, \dots \quad (a_n = a_{n-1} + a_{n-2} + 1)$$

Example of Learning from Data

(from mlss.tuebingen.mpg.de/2013/schoelkopf_whatismL_slides.pdf)

- 1, 2, 4, 7, ... → What are succeeding numbers?

$$1, 2, 4, 7, 11, 16, \dots \quad (a_n = a_{n-1} + n - 1)$$

$$1, 2, 4, 7, 12, 20, \dots \quad (a_n = a_{n-1} + a_{n-2} + 1)$$

$$1, 2, 4, 7, 13, 24, \dots \quad (a_n = a_{n-1} + a_{n-2} + a_{n-3})$$

Example of Learning from Data

(from mlss.tuebingen.mpg.de/2013/schoelkopf_whatismL_slides.pdf)

- 1, 2, 4, 7, ... → What are succeeding numbers?

$$1, 2, 4, 7, 11, 16, \dots \quad (a_n = a_{n-1} + n - 1)$$

$$1, 2, 4, 7, 12, 20, \dots \quad (a_n = a_{n-1} + a_{n-2} + 1)$$

$$1, 2, 4, 7, 13, 24, \dots \quad (a_n = a_{n-1} + a_{n-2} + a_{n-3})$$

$$1, 2, 4, 7, 14, 28 \quad (\text{divisors of } 28)$$

Example of Learning from Data

(from mlss.tuebingen.mpg.de/2013/schoelkopf_whatismL_slides.pdf)

- 1, 2, 4, 7, ... → What are succeeding numbers?

$$1, 2, 4, 7, 11, 16, \dots \quad (a_n = a_{n-1} + n - 1)$$

$$1, 2, 4, 7, 12, 20, \dots \quad (a_n = a_{n-1} + a_{n-2} + 1)$$

$$1, 2, 4, 7, 13, 24, \dots \quad (a_n = a_{n-1} + a_{n-2} + a_{n-3})$$

$$1, 2, 4, 7, 14, 28 \quad (\text{divisors of } 28)$$

$$1, 2, 4, 7, 1, 1, 5, \dots \quad (\pi = 3.1415 \dots \text{ and } e = 2.718 \dots)$$

Example of Learning from Data

(from mlss.tuebingen.mpg.de/2013/schoelkopf_whatismL_slides.pdf)

- 1, 2, 4, 7, ... → What are succeeding numbers?
 - 1, 2, 4, 7, 11, 16, ... ($a_n = a_{n-1} + n - 1$)
 - 1, 2, 4, 7, 12, 20, ... ($a_n = a_{n-1} + a_{n-2} + 1$)
 - 1, 2, 4, 7, 13, 24, ... ($a_n = a_{n-1} + a_{n-2} + a_{n-3}$)
 - 1, 2, 4, 7, 14, 28 (divisors of 28)
 - 1, 2, 4, 7, 1, 1, 5, ... ($\pi = 3.1415 \dots$ and $e = 2.718 \dots$)
- 1344 results (!) in the online encyclopedia (<https://oeis.org/>)

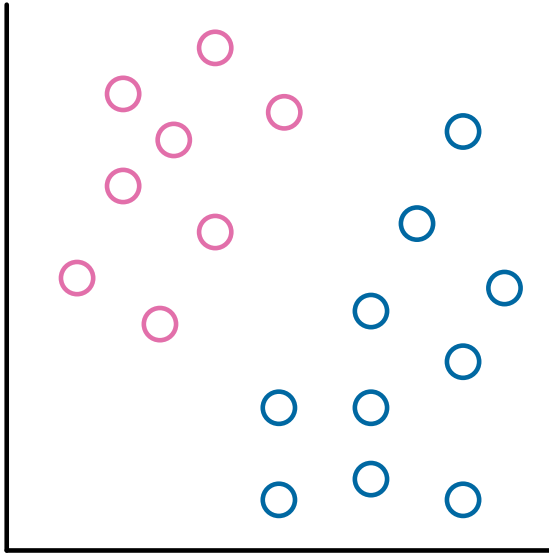
Learning as Scientific Problem

- Which is the correct answer (or **generalization**) for succeeding numbers of 1, 2, 4, 7, ... ?
 - Any answer is possible!

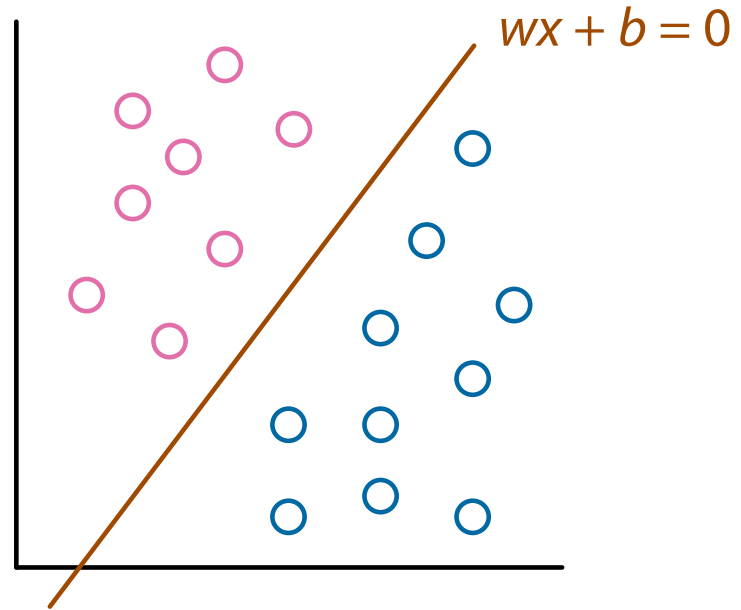
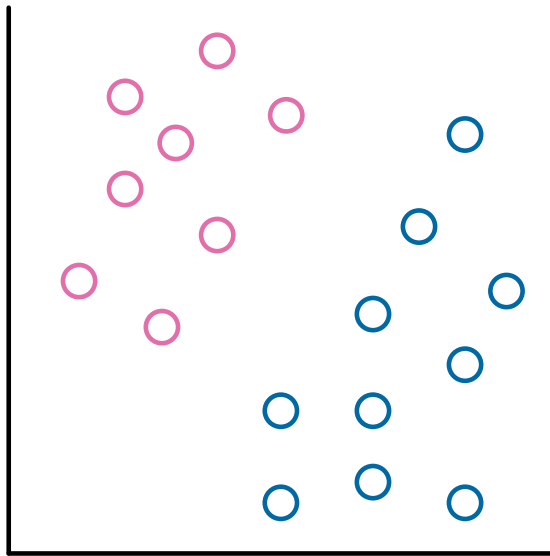
Learning as Scientific Problem

- Which is the correct answer (or **generalization**) for succeeding numbers of 1, 2, 4, 7, ... ?
 - Any answer is possible!
- We should take two points into consideration:
 - (i) We need to formalize the problem of “learning”
 - There are **two agents** (**teacher** and **learner**) in learning, which are different from “computation”
 - (ii) Learning is an **infinite process**

Learning of Binary Classifier



Learning of Binary Classifier



Example: Perceptron (by F. Rosenblatt, 1958)

- **Learning target:** two subsets $F, G \subseteq \mathbb{R}^d$ s.t. $F \cap G = \emptyset$
 - Assumption: F and G are **linearly separable**:
There exists a function (classifier) $f_*(\mathbf{x}) = \langle \mathbf{w}_*, \mathbf{x} \rangle + b$ s.t.
 $f_*(\mathbf{x}) > 0 \quad \forall \mathbf{x} \in F, \quad f_*(\mathbf{x}) < 0 \quad \forall \mathbf{x} \in G$

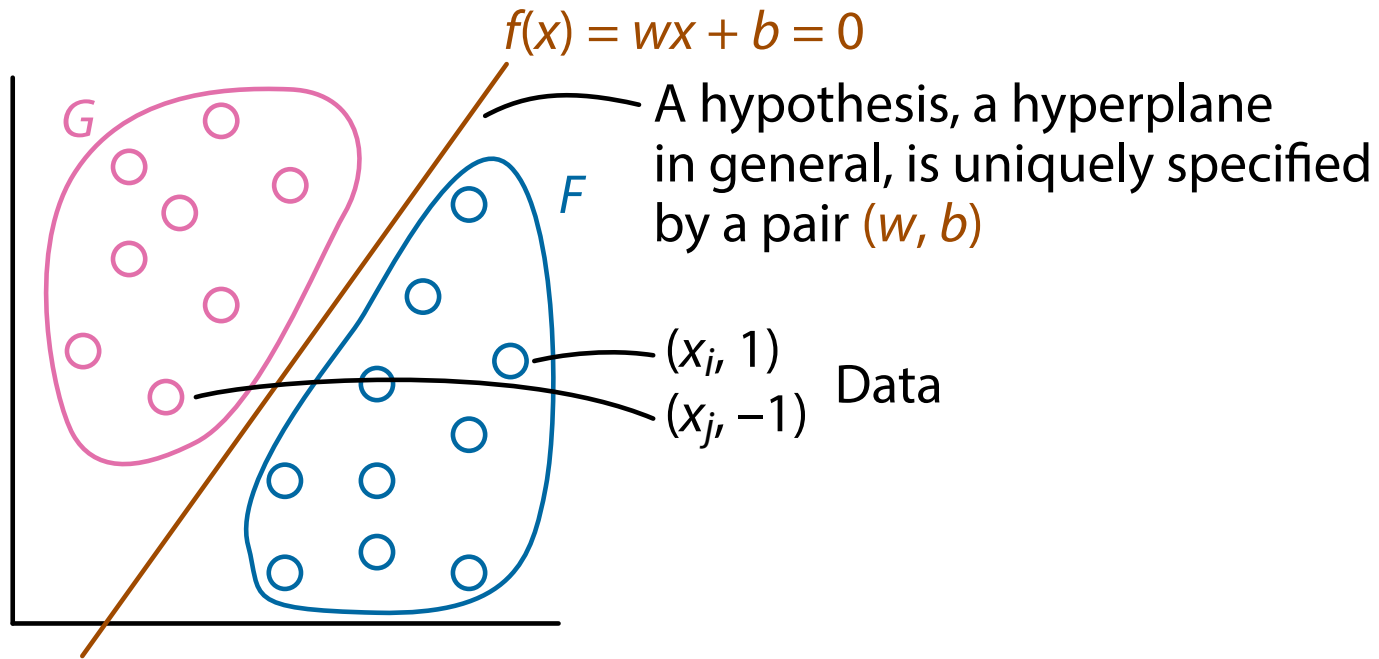
Example: Perceptron (by F. Rosenblatt, 1958)

- **Learning target:** two subsets $F, G \subseteq \mathbb{R}^d$ s.t. $F \cap G = \emptyset$
 - Assumption: F and G are **linearly separable**:
There exists a function (classifier) $f_*(\mathbf{x}) = \langle \mathbf{w}_*, \mathbf{x} \rangle + b$ s.t.
 $f_*(\mathbf{x}) > 0 \quad \forall \mathbf{x} \in F, \quad f_*(\mathbf{x}) < 0 \quad \forall \mathbf{x} \in G$
- **Hypotheses:** hyperplanes on \mathbb{R}^d
 - If we consider a linear equation $f(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + b$, each line can be uniquely specified by a pair of two parameters (\mathbf{w}, b) (**hypothesis**)

Example: Perceptron (by F. Rosenblatt, 1958)

- **Data:** a sequence of pairs $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots$
 - (\mathbf{x}_i, y_i) : (a real-valued vector in \mathbb{R}^d , a label)
 - $\mathbf{x}_i \in F \cup G, y_i \in \{1, -1\}$,
 $y_i = 1$ ($y_i = -1$) if $\mathbf{x}_i \in F$ ($\mathbf{x}_i \in G$)

Learning Model for Perceptron



Learning Procedure of Perceptron

1. $\mathbf{w} \leftarrow 0, b \leftarrow 0$ (or a small random value) // initialization
2. for $i = 1, 2, 3, \dots$ do
3. Receive i -th pair (\mathbf{x}_i, y_i)
4. Compute $a = \sum_{j=1}^d w^j x_i^j + b$
5. if $y_i \cdot a < 0$ then // \mathbf{x}_i is misclassified
6. $\mathbf{w} \leftarrow \mathbf{w} + y_i \mathbf{x}_i$ // update the weight
7. $b \leftarrow b + y_i$ // update the bias
8. end if
9. end for

Correctness of Perceptron

- It is guaranteed that a perceptron always converges to a correct classifier
 - A correct classifier is a function f s.t.
 - $f(\mathbf{x}) > 0 \quad \forall \mathbf{x} \in F,$
 - $f(\mathbf{x}) < 0 \quad \forall \mathbf{x} \in G$
 - The convergence theorem
- Note: there are (infinitely) many functions that correctly classify F and G
 - A perceptron converges to one of them

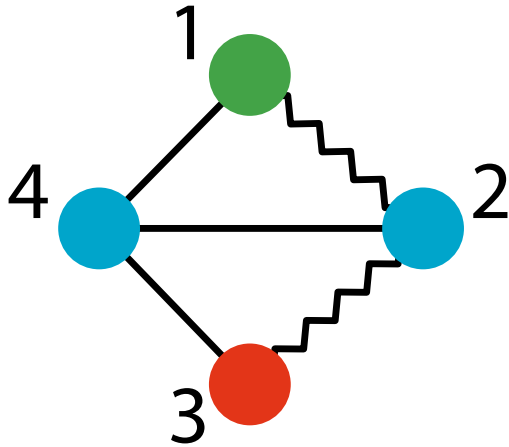
Summary: Perceptron

Target	Two disjoint subsets of \mathbb{R}^d
Representation	Two parameters (\mathbf{w}, b) of linear equation $f(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + b$
Data	Real vectors from target subsets
Algorithm	Perceptron
Correctness	Convergence theorem

What is Graph?

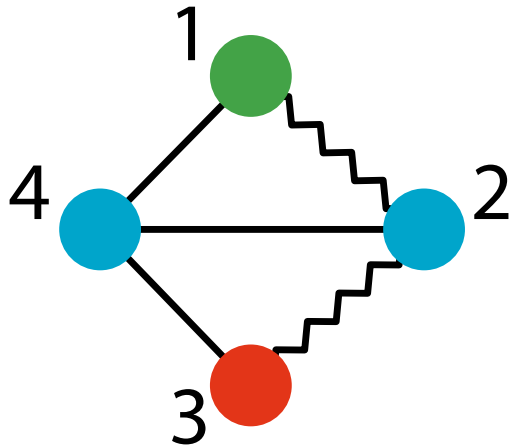
- A graph is an object consisting of **vertices** (nodes) connected with **edges**
 - Many examples in real-world, e.g., chemical compounds
- A graph is **directed** if the edges are directed, otherwise it is **undirected**
- A graph is written as $G = (V, E)$, where V is a vertex set and E is an edge set
- **Labels** can be associated with vertices and/or edges

Example of Graph



- A graph $G = (V, E, \phi)$
 - $V = \{1, 2, 3, 4\}$
 - $E = \{\{1, 2\}, \{1, 4\}, \{2, 3\}, \{2, 4\}, \{3, 4\}\}$
 - $\phi(1) = \text{green}, \phi(2) = \text{blue},$
 $\phi(3) = \text{red}, \phi(4) = \text{blue}$
 - $\phi(\{1, 2\}) = \text{zigzag}, \phi(\{1, 4\}) = \text{straight},$
 $\phi(\{2, 3\}) = \text{zigzag}, \phi(\{2, 4\}) = \text{straight},$
 $\phi(\{3, 4\}) = \text{straight}$

Example of Graph



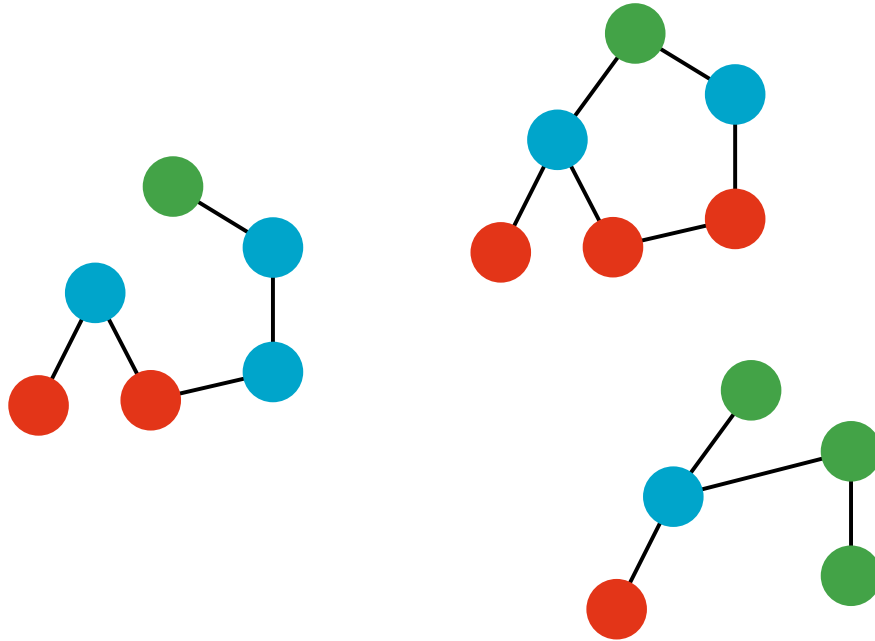
- The adjacency matrix

$$A = \begin{bmatrix} 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{bmatrix}$$

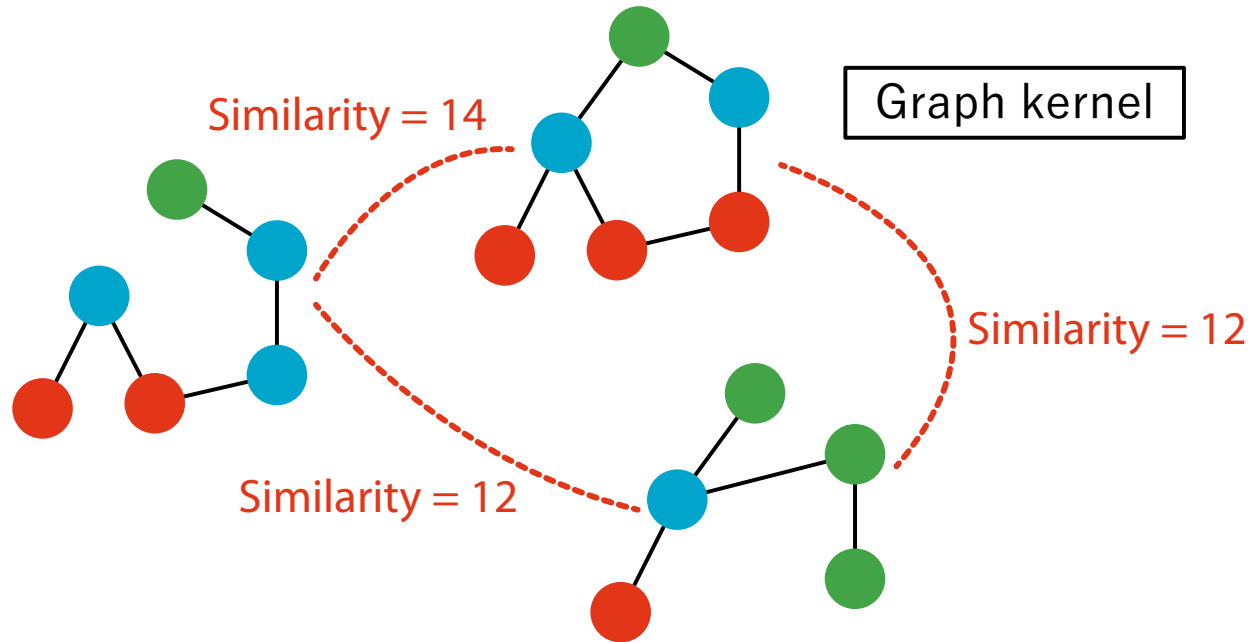
ML on Graphs

- How to perform machine learning on **graphs**?
 - Each object is a graph, so we have a collection of graphs
- Classification (or regression) on graphs is nontrivial problem
 - The difficulty comes from the fact that measuring the **similarity** (or distance) between graphs is nontrivial
- **Graph kernel** computes the similarity between graphs
- **Graph Neural Networks** are recently studied, while there is no significant difference between their performances
 - They share the core idea (message passing)

Similarity between Graphs



Similarity between Graphs



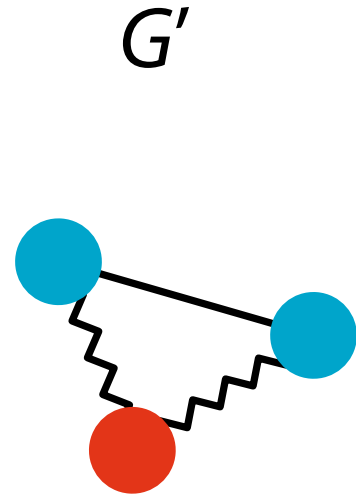
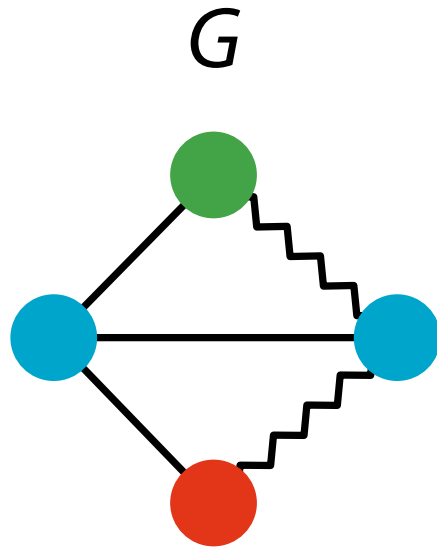
Kernels on Structured Data

- Given objects X and Y , **decompose** them into substructures S and T
- The **R-convolution kernel** K_R by Haussler (1999) is:

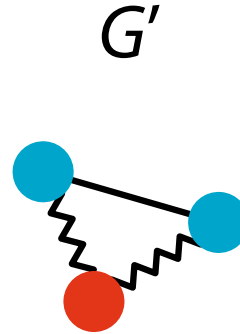
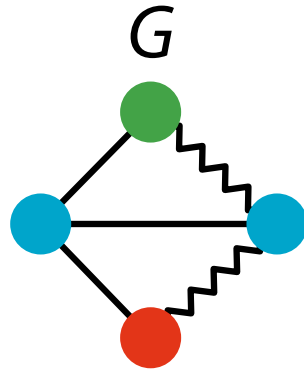
$$K_R(X, Y) = \sum_{s \in S, t \in T} K_{\text{base}}(s, t)$$




- e.g. X is a graph and S is the set of all subgraphs
- Since naively computing this kernel is expensive, many efficient graph kernels have been proposed

Example



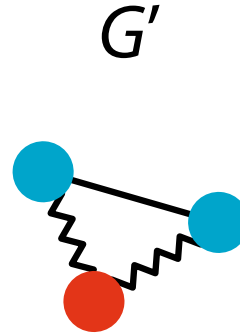
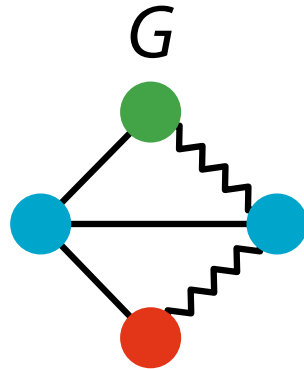
Vertex Label Histogram Kernel



			
G	2	1	1
G'	2	0	1

$$K_{\text{VH}}(G, G') = 2 \cdot 2 + 1 \cdot 0 + 1 \cdot 1 = 5$$

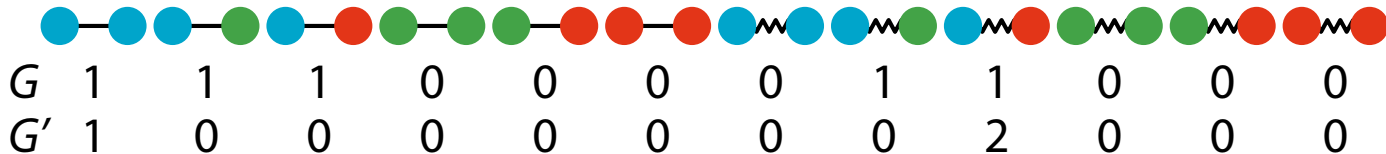
Edge Label Histogram Kernel



	—	~~~~
G	3	2
G'	1	2

$$K_{EH}(G, G') = 3 \cdot 1 + 2 \cdot 2 = 7$$

Vertex-Edge Label Histogram Kernel



G	1	1	1	0	0	0	0	0	1	1	0	0	0
G'	1	0	0	0	0	0	0	0	0	2	0	0	0

$$K_{\text{VEH}}(G, G') = 3$$

Product Graph

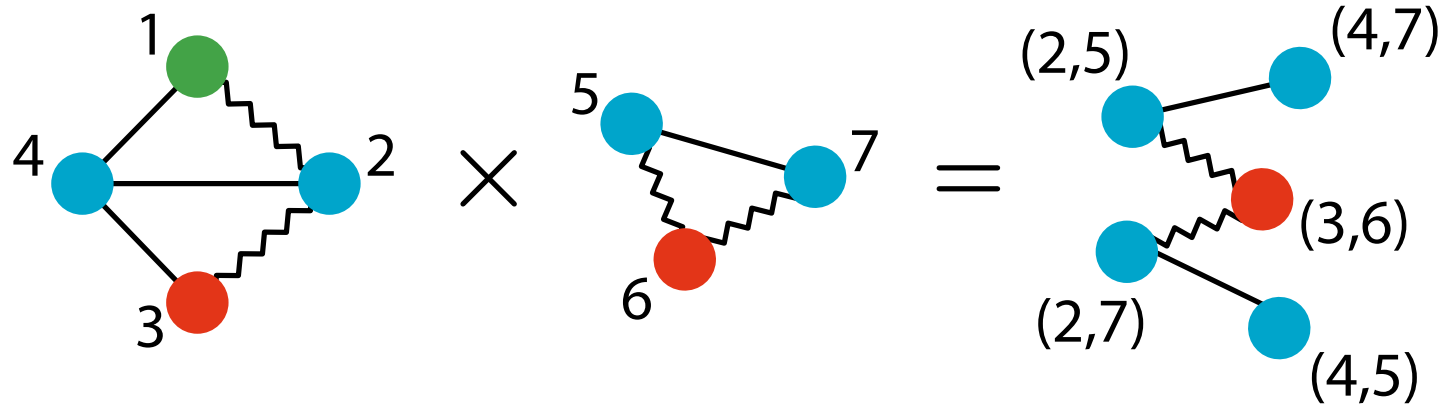
- The **direct product** $G_{\times} = (V_{\times}, E_{\times}, \phi_{\times})$ of $G = (V, E, \phi)$, $G' = (V', E', \phi')$:

$$V_{\times} = \{(v, v') \in V \times V' \mid \phi(v) = \phi'(v')\},$$

$$E_{\times} = \left\{ ((u, u'), (v, v')) \in V_{\times} \times V_{\times} \mid \begin{array}{l} (u, v) \in E, (u', v') \in E', \\ \phi(u, v) = \phi'(u', v') \end{array} \right\}$$

- All labels are inherited

Example of Product Graph



k -Step Random Walk Kernel

- The k -step (fixed-length- k) random walk kernel between G and G' :

$$K_{\times}^k(G, G') = \sum_{i,j=1}^{|V_{\times}|} \left[\lambda_0 A_{\times}^0 + \lambda_1 A_{\times}^1 + \lambda_2 A_{\times}^2 + \cdots + \lambda_k A_{\times}^k \right]_{ij}$$

$(\lambda_l > 0)$

- A_{\times} : The adjacency matrix of the product graph
- The ij entry of A_{\times}^n shows the number of paths from i to j

Geometric Random Walk Kernel (1/2)

- K_{\times}^{∞} can be directly computed if $\lambda_{\ell} = \lambda^{\ell}$ (geometric series), resulting in the geometric random walk kernel:

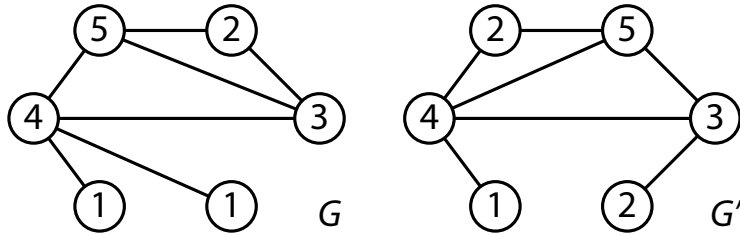
$$\begin{aligned} K_{\text{GR}}(G, G') &= \sum_{i,j=1}^{|V_{\times}|} \left[\lambda^0 A_{\times}^0 + \lambda^1 A_{\times}^1 + \lambda^2 A_{\times}^2 + \lambda^3 A_{\times}^3 + \dots \right]_{ij} \\ &= \sum_{i,j=1}^{|V_{\times}|} \left[\sum_{\ell=0}^{\infty} \lambda^{\ell} A_{\times}^{\ell} \right]_{ij} = \sum_{i,j=1}^{|V_{\times}|} \left[(\mathbf{I} - \lambda A_{\times})^{-1} \right]_{ij} \end{aligned}$$

Geometric Random Walk Kernel (2/2)

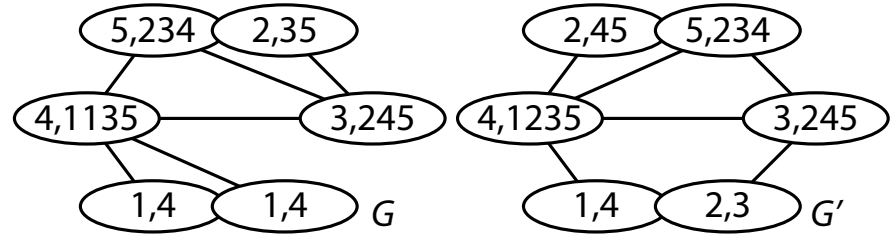
- Geometric random walk kernel is well-defined only if $\lambda < 1/\mu_{x,\max}$ ($\mu_{x,\max}$ is the max. eigenvalue of A_x)
- δ_x (min. degree) $\leq \bar{d}_x$ (average degree) $\leq \mu_{x,\max} \leq \Delta_x$ (max. degree)

Weisfeiler-Lehman Kernel

Given graphs



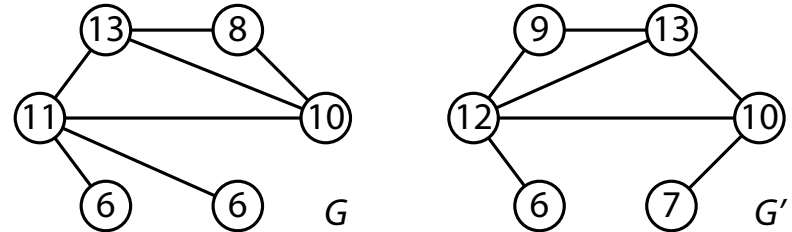
1st iteration



Re-labeling after 1st iteration

1,4 \rightarrow 6	3,245 \rightarrow 10
2,3 \rightarrow 7	4,1135 \rightarrow 11
2,35 \rightarrow 8	4,1235 \rightarrow 12
2,45 \rightarrow 9	5,234 \rightarrow 13

After 1st iteration



Weisfeiler-Lehman Kernel

- The kernel value becomes:

$$\begin{bmatrix} \text{label} \\ \phi(G)^{(1)} \\ \phi(G')^{(1)} \end{bmatrix} = \begin{bmatrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 \\ 2 & 1 & 1 & 1 & 1 & 2 & 0 & 1 & 0 & 1 & 1 & 0 & 1 \\ 1 & 2 & 1 & 1 & 1 & 1 & 1 & 0 & 1 & 1 & 0 & 1 & 1 \end{bmatrix},$$

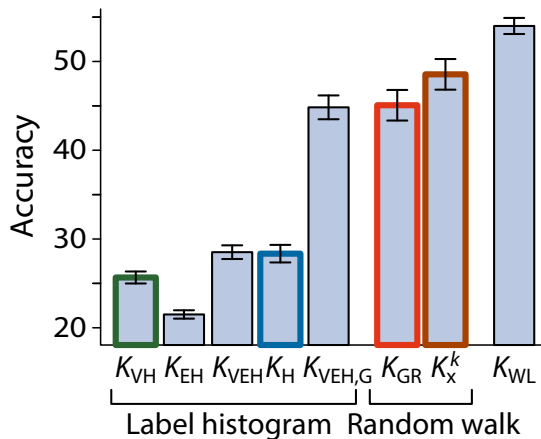
$$K_{\text{WL}}^1(G, G') = 11$$

- An important building block of GNNs

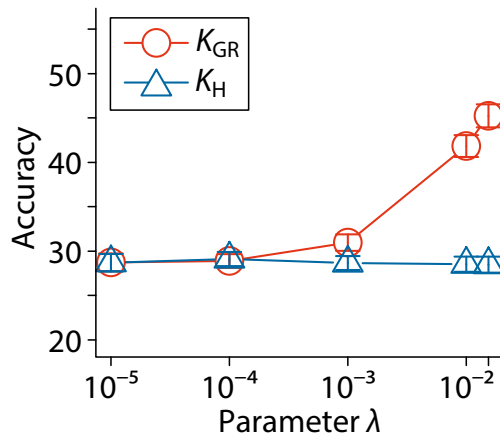
Performance Comparison

ENZYMES

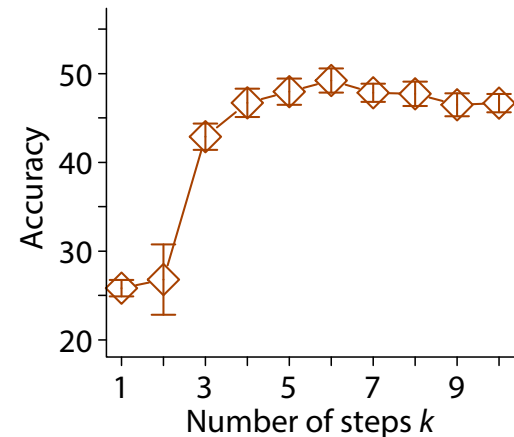
(i) Comparison of various graph kernels



(ii) Comparison of K_{GR} with K_H



(iii) k -step K_x^k



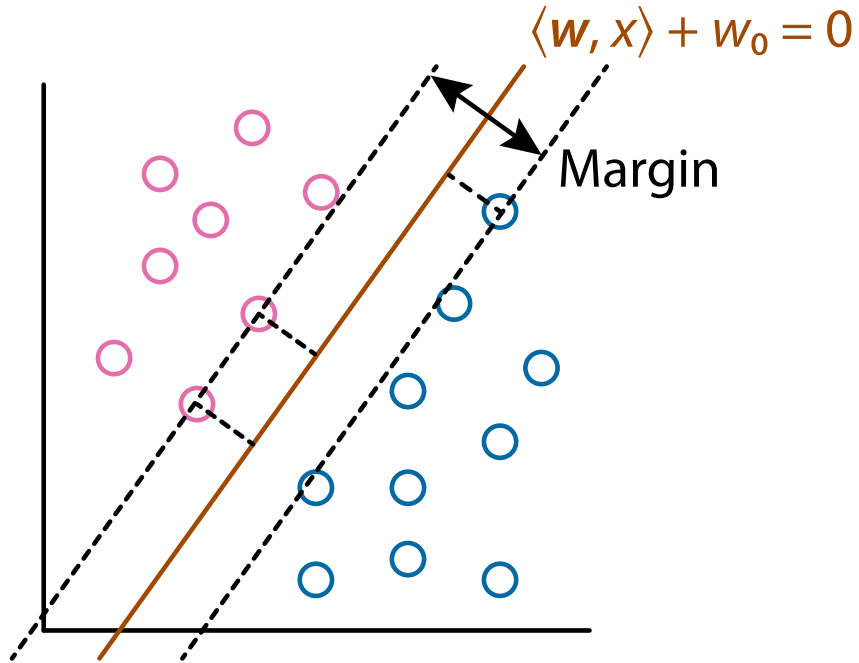
graphkernels Package

- A package for graph kernels available in R and Python
- R:
<https://CRAN.R-project.org/package=graphkernels>
- Python:
<https://pypi.org/project/graphkernels/>
- Paper:
<https://doi.org/10.1093/bioinformatics/btx602>

Kernel-based Classification: SVM

- A dataset D is **separable** by $f \iff y_i f(\mathbf{x}_i) > 0, \forall i \in \{1, 2, \dots, n\}$
- The **margin** is the distance from the classification hyperplane to the closest data point
- Support vector machines (SVMs) tries to find a hyperplane that **maximizes** the margin
 - Can be viewed as an extension of perceptron

Margin



Formulation of SVMs

- The distance from a point \mathbf{x}_i to a hyperplane

$f(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + w_0 = 0$ is

$$\frac{|f(\mathbf{x}_i)|}{\|\mathbf{w}\|} = \frac{|\langle \mathbf{w}, \mathbf{x}_i \rangle + w_0|}{\|\mathbf{w}\|}$$

- The margin maximization problem can be written as

$$\max_{\mathbf{w}, w_0, M} \frac{M}{\|\mathbf{w}\|} \quad \text{subject to } y_i f(\mathbf{x}_i) \geq M, i \in \{1, 2, \dots, n\}$$

$$- M = \min_{i \in \{1, 2, \dots, n\}} |\langle \mathbf{w}, \mathbf{x}_i \rangle + w_0|$$

Hard Margin SVMs

- We can eliminate M and obtain

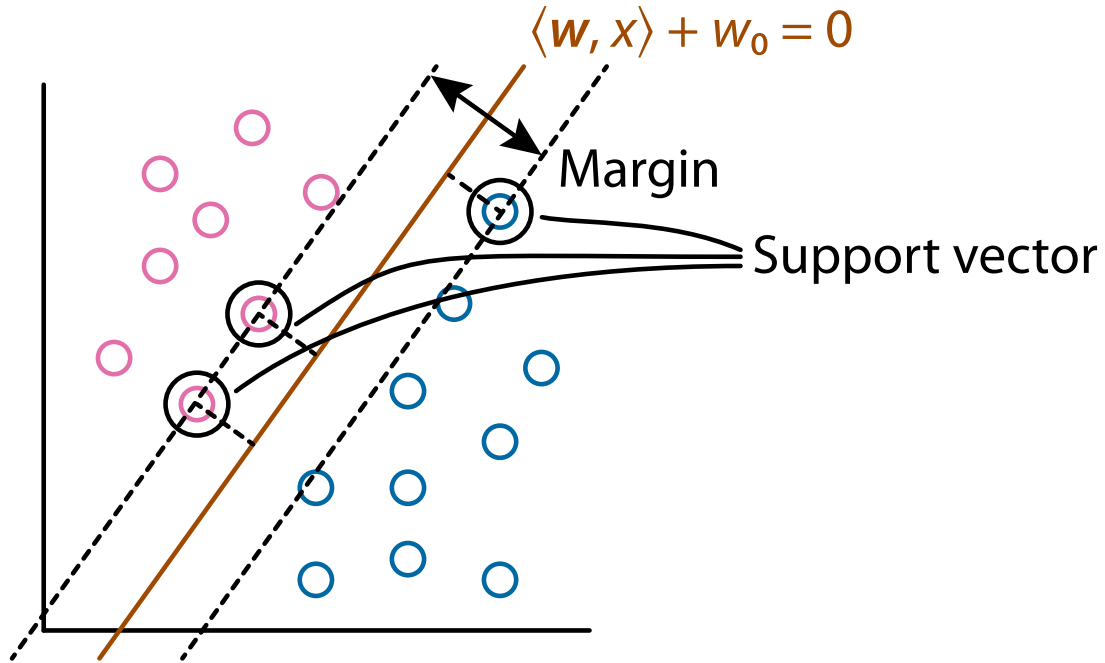
$$\max_{\mathbf{w}, w_0} \frac{1}{\|\mathbf{w}\|} \quad \text{subject to } y_i f(\mathbf{x}_i) \geq 1, i \in \{1, 2, \dots, n\}$$

- This is equivalent to

$$\min_{\mathbf{w}, w_0} \|\mathbf{w}\|^2 \quad \text{subject to } y_i f(\mathbf{x}_i) \geq 1, i \in \{1, 2, \dots, n\}$$

- The standard formulation of **hard margin SVMs**
- There are data points x_i satisfying $y_i f(\mathbf{x}_i) = 1$, called **support vectors**

Margin



Soft Margin (1/2)

- Datasets are not often separable
- Extend SV classification to **soft margin** by relaxing $\langle \mathbf{w}, \mathbf{x} \rangle + w_0 \geq 1$
- Change the constraint $y_i f(\mathbf{x}_i) \geq 1$ using the **slack variable** ξ_i to
$$y_i f(\mathbf{x}_i) = y_i (\langle \mathbf{w}, \mathbf{x} \rangle + w_0) \geq 1 - \xi_i, \quad i \in \{1, 2, \dots, n\}$$

Soft Margin (2/2)

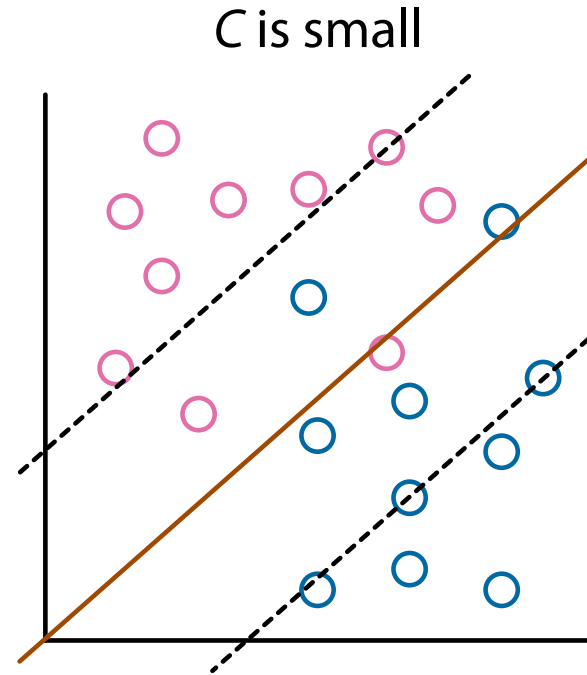
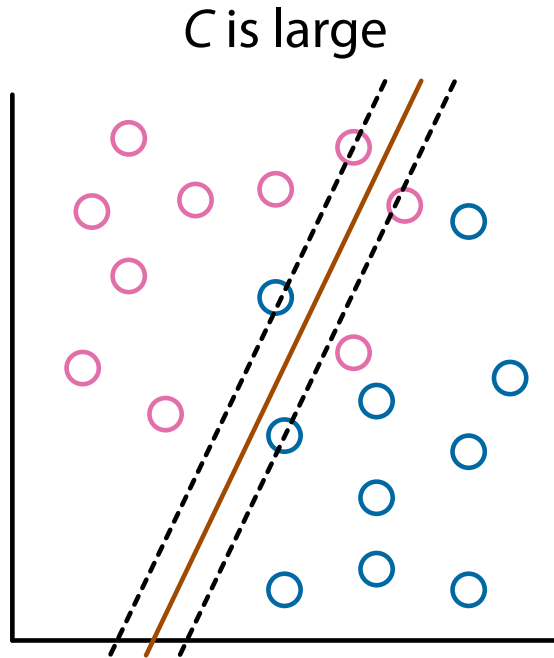
- The formulation of **soft margin SVM** (C-SVM) is

$$\min_{\mathbf{w}, w_0, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i \in \{1, 2, \dots, n\}} \xi_i$$

$$\text{s.t. } y_i f(\mathbf{x}_i) \geq 1 - \xi_i, \xi_i \geq 0, i \in \{1, 2, \dots, n\}$$

- C is called the **regularization parameter**

Soft Margin



Data Point Location

- $y_i f(\mathbf{x}_i) > 1$: \mathbf{x}_i is outside margin
 - These points do not affect to the classification hyperplane
- $y_i f(\mathbf{x}_i) = 1$: \mathbf{x}_i is on margin
- $y_i f(\mathbf{x}_i) < 1$: \mathbf{x}_i is inside margin
 - These points do not exist in hard margin
- Points on margin and inside margin are support vectors

Dual Problem (1/2)

- The formulation of C-SVM

$$\min_{\mathbf{w}, w_0, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i \in \{1, 2, \dots, n\}} \xi_i$$

$$\text{s.t. } y_i f(\mathbf{x}_i) \geq 1 - \xi_i, \xi_i \geq 0, i \in \{1, 2, \dots, n\}$$

is called the **primal problem**

- This is usually solved via the **dual problem**

Dual Problem (2/2)

- The dual problem is given as

$$\max_{\alpha} -\frac{1}{2} \sum_{i,j \in [n]} \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle + \sum_{i \in [n]} \alpha_i$$

$$\text{s.t. } \sum_{i \in [n]} \alpha_i y_i = 0, 0 \leq \alpha_i \leq C, i \in [n]$$

Extension to Nonlinear Classification

- To achieve nonlinear classification, convert each data point \mathbf{x} to some point $\phi(\mathbf{x})$, and $f(\mathbf{x})$ becomes

$$f(\mathbf{x}) = \langle \mathbf{w}, \phi(\mathbf{x}) \rangle + w_0$$

- The dual problem becomes

$$\max_{\alpha} -\frac{1}{2} \sum_{i,j \in [n]} \alpha_i \alpha_j y_i y_j \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle + \sum_{i \in [n]} \alpha_i \quad \text{s.t.} \quad \sum_{i \in [n]} \alpha_i y_i = 0, \quad 0 \leq \alpha_i \leq C, \quad i \in [n]$$

- Only the dot product $\langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$ is used!
- We do not even need to know $\phi(\mathbf{x}_i)$ and $\phi(\mathbf{x}_j)$

C-SVM with Kernel Trick

- Using the **kernel function** K such that $K(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$, we have

$$\max_{\alpha} -\frac{1}{2} \sum_{i,j \in [n]} \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) + \sum_{i \in [n]} \alpha_i$$

$$\text{s.t. } \sum_{i \in [n]} \alpha_i y_i = 0, 0 \leq \alpha_i \leq C, i \in [n]$$

- The technique of using K is called **kernel trick**
- We can use graph kernels for K !

Positive Definite Kernel

- A kernel $K : \Omega \times \Omega \rightarrow \mathbb{R}$ is a **positive definite kernel** if

(i) $K(x, y) = K(y, x)$

- (ii) For x_1, x_2, \dots, x_n , the $n \times n$ matrix

$$(K_{ij}) = \begin{bmatrix} K(x_1, x_1) & \dots & K(x_n, x_1) \\ \dots & \dots & \dots \\ K(x_1, x_n) & \dots & K(x_n, x_n) \end{bmatrix}$$

is positive (semi-)definite, that is, $\sum_{i,j=1}^n c_i c_j K(x_i, x_j) \geq 0$
for any $c_1, c_2, \dots, c_n \in \mathbb{R}$

- $(K_{ij}) \in \mathbb{R}^{n \times n}$ is called the **Gram matrix**

Popular Positive Definite Kernels

- Linear Kernel

$$K(\mathbf{x}, \mathbf{y}) = \langle \mathbf{x}, \mathbf{y} \rangle$$

- Gaussian (RBF) kernel

$$K(\mathbf{x}, \mathbf{y}) = \exp\left(-\frac{1}{\sigma^2} \|\mathbf{x} - \mathbf{y}\|^2\right)$$

- Polynomial Kernel

$$K(\mathbf{x}, \mathbf{y}) = (\langle \mathbf{x}, \mathbf{y} \rangle + c)^d \quad c, d \in \mathbb{R}$$

Simple Kernels

- The all-ones kernel

$$K(\mathbf{x}, \mathbf{y}) = 1$$

- The delta (Dirac) kernel

$$K(\mathbf{x}, \mathbf{y}) = \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{y}, \\ 0 & \text{otherwise} \end{cases}$$

Closure Properties of Kernels

- For two kernels K_1 and K_2 , $K_1 + K_2$ is a kernel
- For two kernels K_1 and K_2 , the product $K_1 \cdot K_2$ is a kernel
- For a kernel K and a positive scalar $\lambda \in \mathbb{R}^+$, λK is a kernel
- For a kernel K on a set D , its zero-extension:

$$K_0(\mathbf{x}, \mathbf{y}) = \begin{cases} K(\mathbf{x}, \mathbf{y}) & \text{if } \mathbf{x}, \mathbf{y} \in D, \\ 0 & \text{otherwise} \end{cases}$$

is a kernel

Summary

- Introduction to ML and its basic concepts
- Graph kernels for graph structured data
- Kernel-based ML methods, such as SVM
 - There are many other options, e.g. kernel PCA, kernel k -means, kernel ridge regression, ...