

KAIST Machine Learning Tutorial

Metrics and Kernels A few recent topics

Marco Cuturi - Kyoto University

Data deluge



ACM, Dec'08



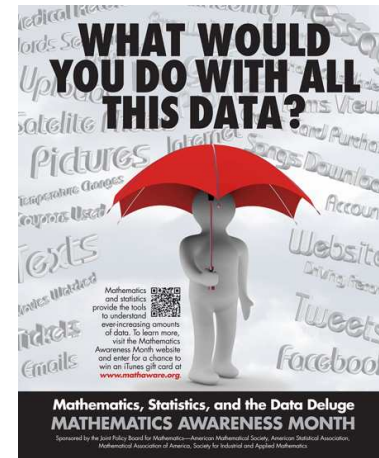
Economist Feb'10



IEEE Spectrum, Feb'11



Popular Science Nov. 2011



AMS/ASA Apr'12

Quantity and Complexity

Social lending



Download Loan Data

These files contain complete loan data, including the current loan status (Current, Late, F
We have removed all personally identifiable information to protect our members' privacy.



Download CSV

(44,533kb)

Quantity and Complexity

Social lending



448323,7000,7000,17.04%,36 months,10/3/09,10/17/09,10/15/09,E3,My Loan,debt_consolidation," I would like to know if I qualify for a 5 year term instead of a 3 year term. I would like to have my payments lower. 549920 added on 10/09/09 > My bills are currently almost \$3000 per month, This includes rent, car payment, utilities and other credit accounts. With this loan, I will be able to pay off my accounts much faster by consolidating 4 of my smaller debt accounts into one larger one. By giving me this loan I will be able to pay off this particular account in three years or less very easily. I just finished twelve years of military service, and was asked to return to my old job as a DOD employee. The doctor I worked with thought very highly of me to offer this job to me. My pay is salary based so my pay stays the same regardless of hours worked, with the exception of overtime. My job is very stable being in the medical field and taking care of Soldiers, Retirees, Veterans and their families. Please consider me for this loan as it will help me with my monthly payments. Thank you so very much. 549920 added on 10/10/09 > The first sentence was for a different amount than the \$7,000 currently being offered. With the new amount of \$7,000 I will have no problem with paying the monthly payment that will be required. With my current budget adding up to almost \$3,000 per month with this consolidation will put me about \$2,500 per month which will allow me to pay extra on this debt and start a savings account. 549920 added on 10/11/09 > I am trying to raise my credit score by deleting the amount of open credit accounts I have. This consolidation will allow me to do so. 549920 added on 10/13/09 > Would like to also add that by granting me this loan the balance of my monies after all bills are paid will be almost \$2,000 a month. This amount will allow me to double up payments and pay off this and my other debts alot earlier than anticipated.",249.72,Current,7000,23.39%,0,8376.96,0,8376.96,Member_549920,EL PASO,TX,RENT,6029,660-678,10/31/98,7,34,4908,34.60%,0,0,0,2,10,0,,Hawaii Pacific University,< 1 year,72177

Quantity and Complexity

Social lending

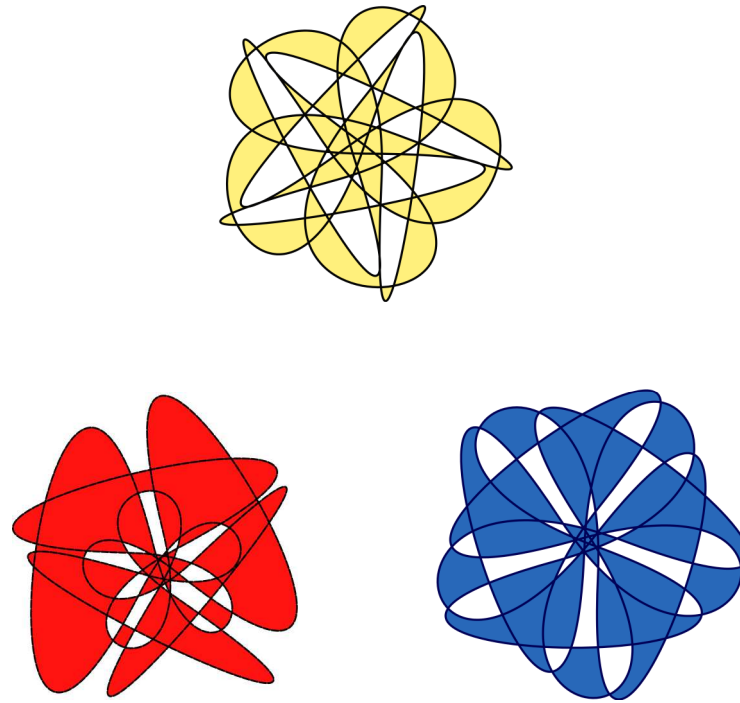
 LendingClub



AQUSH
アクシュ

448323,7000,7000,17.04%,36 months,10/3/09,10/17/09,10/15/09,E3,My Loan,debt_consolidation," I would like to know if I qualify for a 5 year term instead of a 3 year term. I would like to have my payments lower. **549920 added on 10/09/09** > My bills are currently almost \$3000 per month, This includes rent, car payment, utilities and other credit accounts. With this loan, I will be able to pay off my accounts much faster by consolidating 4 of my smaller debt accounts into one larger one. By giving me this loan I will be able to pay off this particular account in three years or less very easily. I just finished twelve years of military service, and was asked to return to my old job as a DOD employee. The doctor I worked with thought very highly of me to offer this job to me. My pay is salary based so my pay stays the same regardless of hours worked, with the exception of overtime. My job is very stable being in the medical field and taking care of Soldiers, Retirees, Veterans and their families. Please consider me for this loan as it will help me with my monthly payments. Thank you so very much. **549920 added on 10/10/09** > The first sentence was for a different amount than the \$7,000 currently being offered. With the new amount of \$7,000 I will have no problem with paying the monthly payment that will be required. With my current budget adding up to almost \$3,000 per month with this consolidation will put me about \$2,500 per month which will allow me to pay extra on this debt and start a savings account. **549920 added on 10/11/09** > I am trying to raise my credit score by deleting the amount of open credit accounts I have. This consolidation will allow me to do so. **549920 added on 10/13/09** > Would like to also add that by granting me this loan the balance of my monies after all bills are paid will be almost \$2,000 a month. This amount will allow me to double up payments and pay off this and my other debts alot earlier than anticipated.",249.72,Current,7000,23.39%,0,8376.96,0,8376.96,Member_549920,EL PASO,TX,RENT,6029,660-678,10/31/98,7,34,4908,34.60%,0,0,0,2,10,0,,Hawaii Pacific University,< 1 year,72177

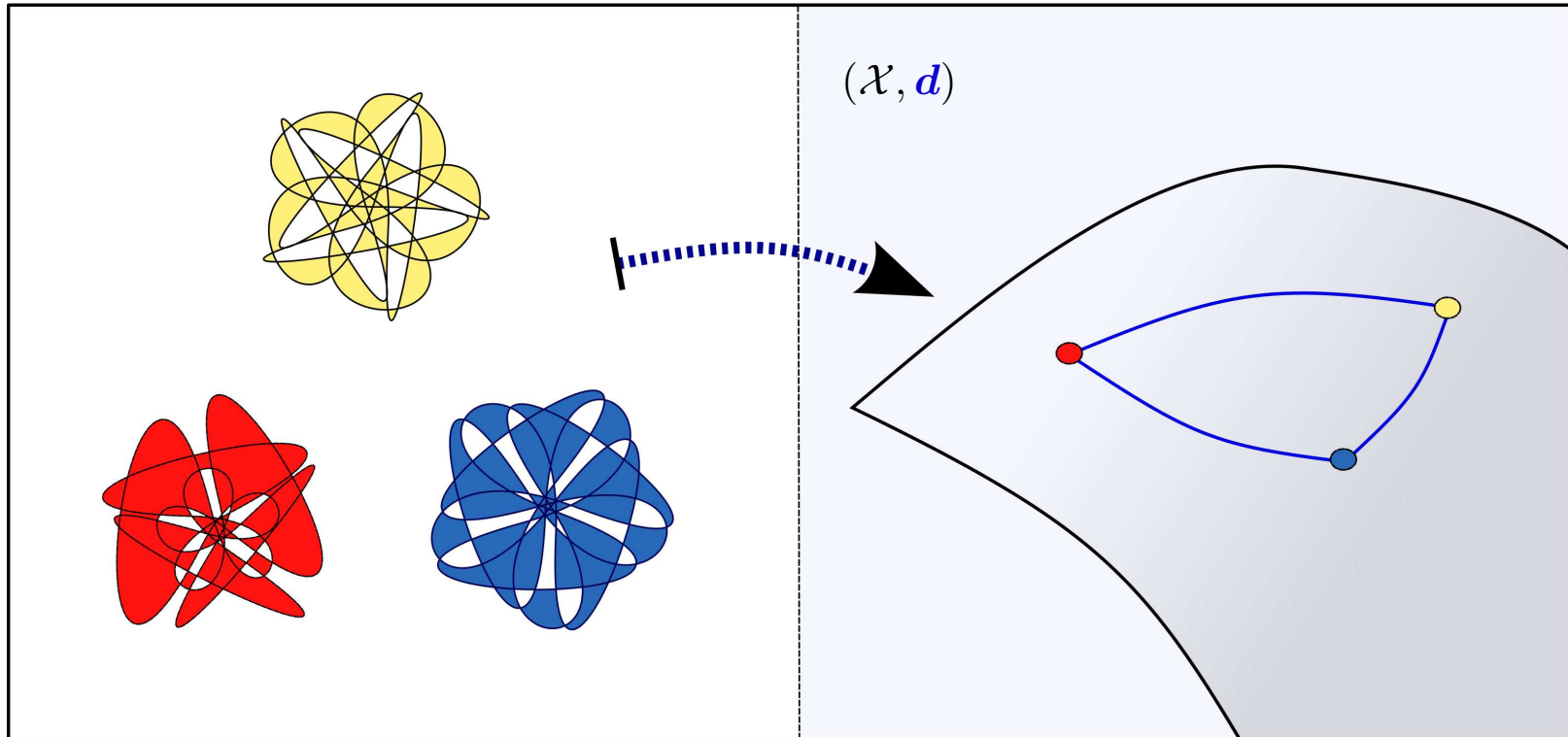
Address Complexity in Machine Learning



- Embed structures in metric spaces
- Embed structures in Hilbert spaces
- Create feature vectors: embed in \mathbb{R}^d .

↪ metric Space

$$d(x, y)$$

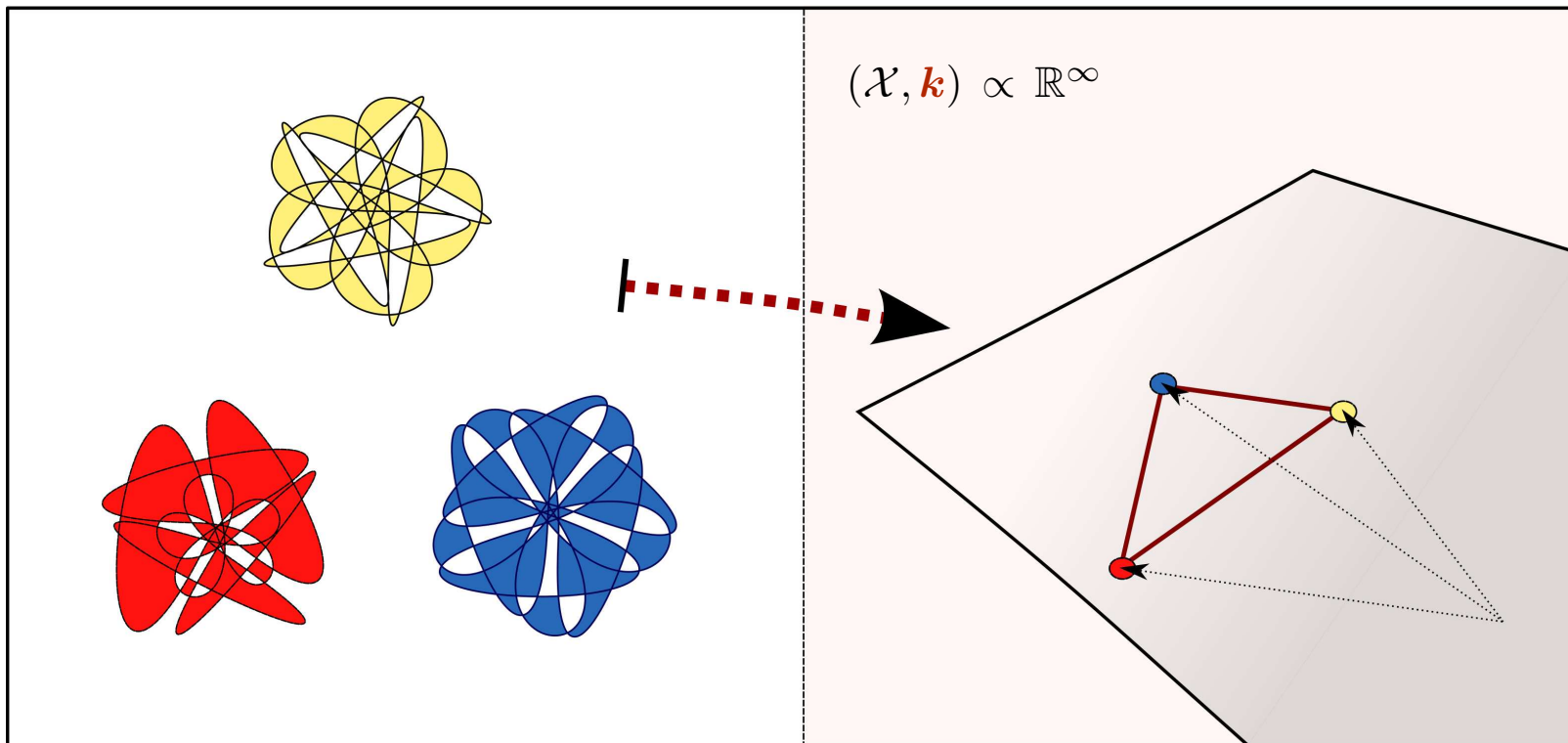


Retrieval & Clustering - Nearest Neighbor Methods

+ **implicit map, flexibility** / -- **limited operations**

↪ Hilbert Space

$$\mathbf{k}(x, y) = \langle \Phi(x), \Phi(y) \rangle_{\mathcal{H}}, \quad \mathbf{d}(x, y) = \|\Phi(x) - \Phi(y)\|_{\mathcal{H}}$$

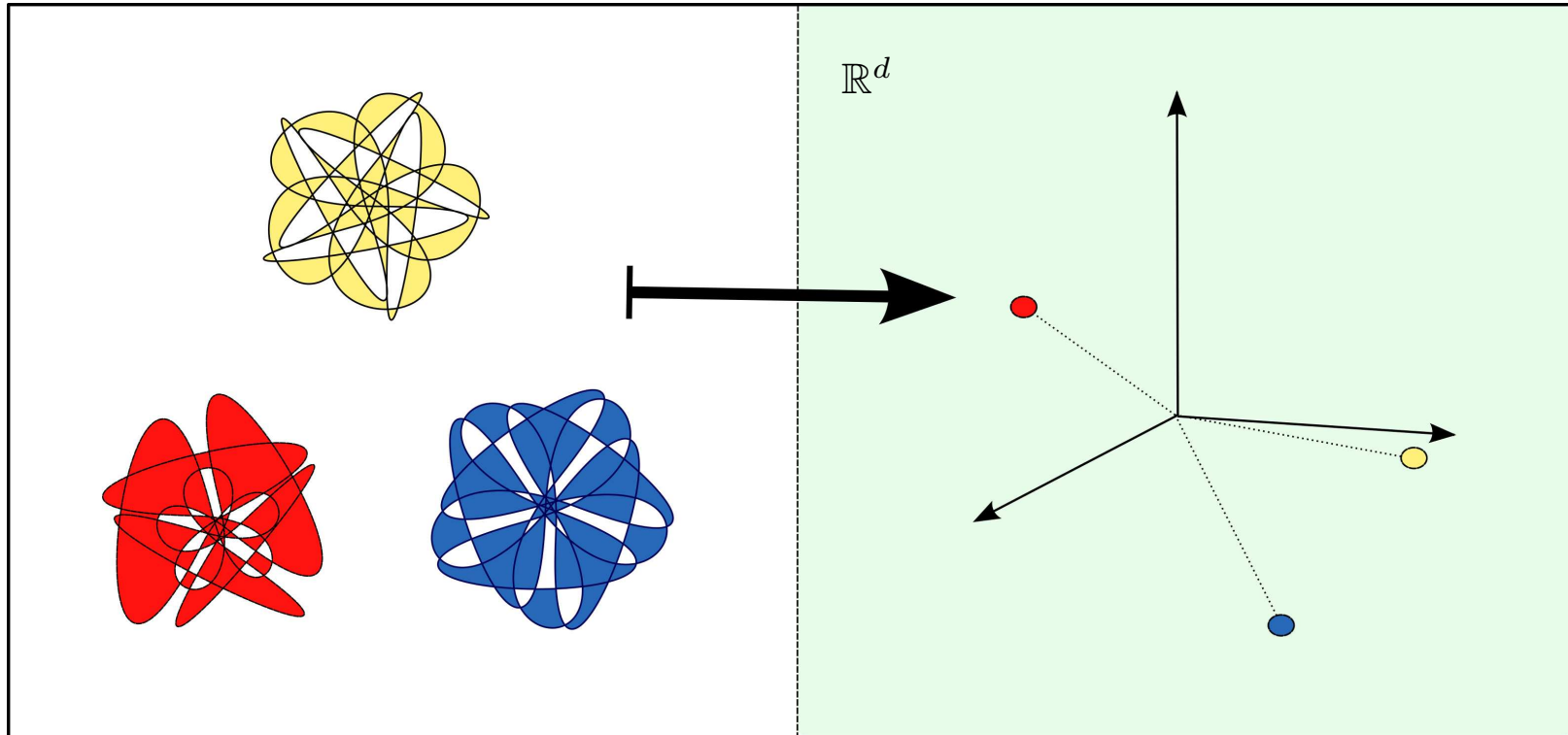


Kernel Methods, Support Vector Machines, Gaussian Processes ...

++ **kernel trick, convexity** / – **limited operations**

↪ Feature Vectors

$$\mathbf{x} = [x_1 \ x_2 \ x_3 \ \cdots \ x_d]^T$$



Linear Models - High-dimensional statistics - Deep Learning ...

++ **fast, feature selection** / - **explicit map**

Outline of this Tutorial

Talk about **distance** and **kernel** based methods

- Start with **distances**;
 - **why distances** in machine learning? example of k -nearest neighbors;
 - On **learning** a **good** distance for vectors with **Mahalanobis distances**

Outline of this Tutorial

- Continue with **kernels**
 - define **positive definite kernels**
 - introduce **support vector machines** to illustrate

- Conclude by mentioning the interface between **kernels** and **distances**

Distances

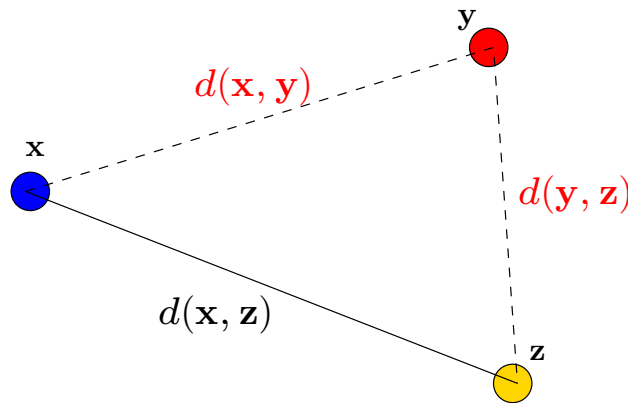
Distances

Definition: A **distance** defined on a set \mathcal{X} is a function

$$\begin{aligned} d: \mathcal{X} \times \mathcal{X} &\rightarrow \mathbb{R}_+ \\ (\mathbf{x}, \mathbf{y}) &\mapsto d(\mathbf{x}, \mathbf{y}) \end{aligned}$$

such that $\forall \mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathcal{X}$,

- $d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x})$, **symmetry**
- $d(\mathbf{x}, \mathbf{y}) = 0 \Leftrightarrow \mathbf{x} = \mathbf{y}$, **definiteness**
- $d(\mathbf{x}, \mathbf{z}) \leq d(\mathbf{x}, \mathbf{y}) + d(\mathbf{y}, \mathbf{z})$, **triangle inequality**



Why Distances in Machine Learning?

- In machine learning problems, we consider **training sets** of points

$$T = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$$

and use these points to “learn” a way to handle **new** observations \mathbf{x} .

For distance and kernel based algorithms,
learning from a training set T
actually means use T to check how similar \mathbf{x} is
to the instances \mathbf{x}_i in T
and make a decision based on that.

similar? for distance based tools \rightarrow how small $d(\mathbf{x}, \mathbf{x}_i)$ is.

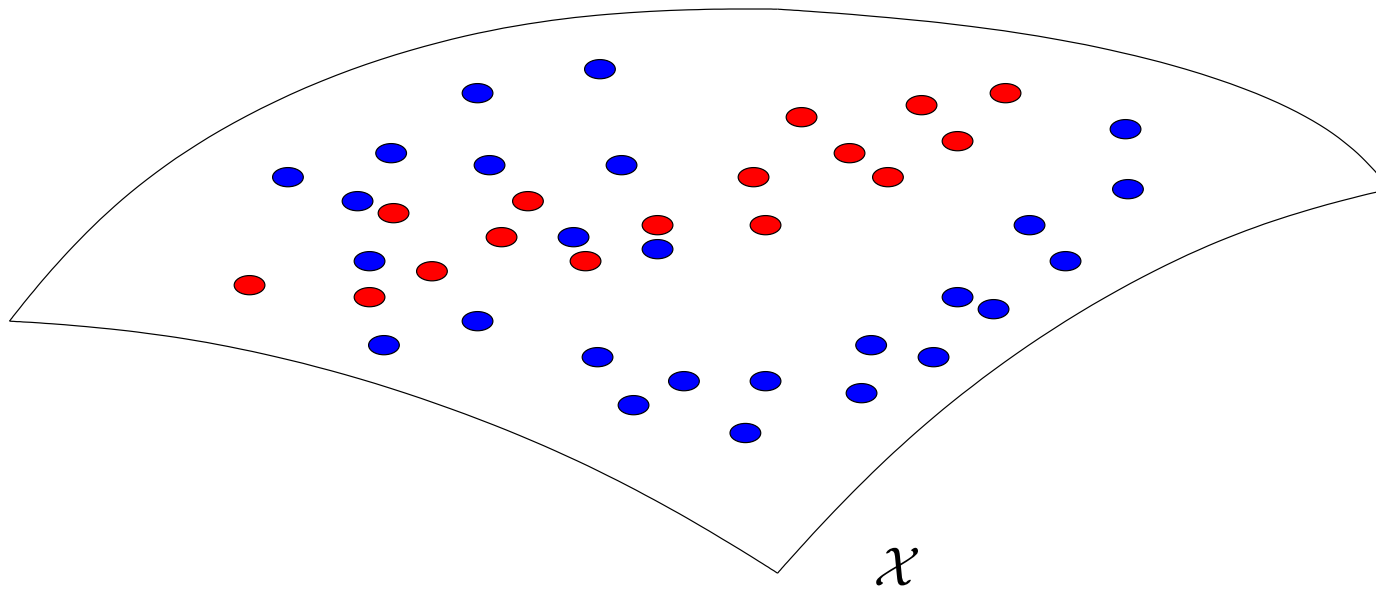
Example: k -nearest neighbors

A very elementary prediction tool

- Used in supervised tasks: given a point \mathbf{x} , guess its label y .
- All it takes to use **k -nearest neighbors** is
 - A data sample of labeled points $\{(\mathbf{x}_i, y_i), i = 1, \dots, N\}$.
 - A **distance function** for two points $d(\mathbf{x}, \mathbf{x}')$.
 - a parameter k that defines the size of the neighborhood.

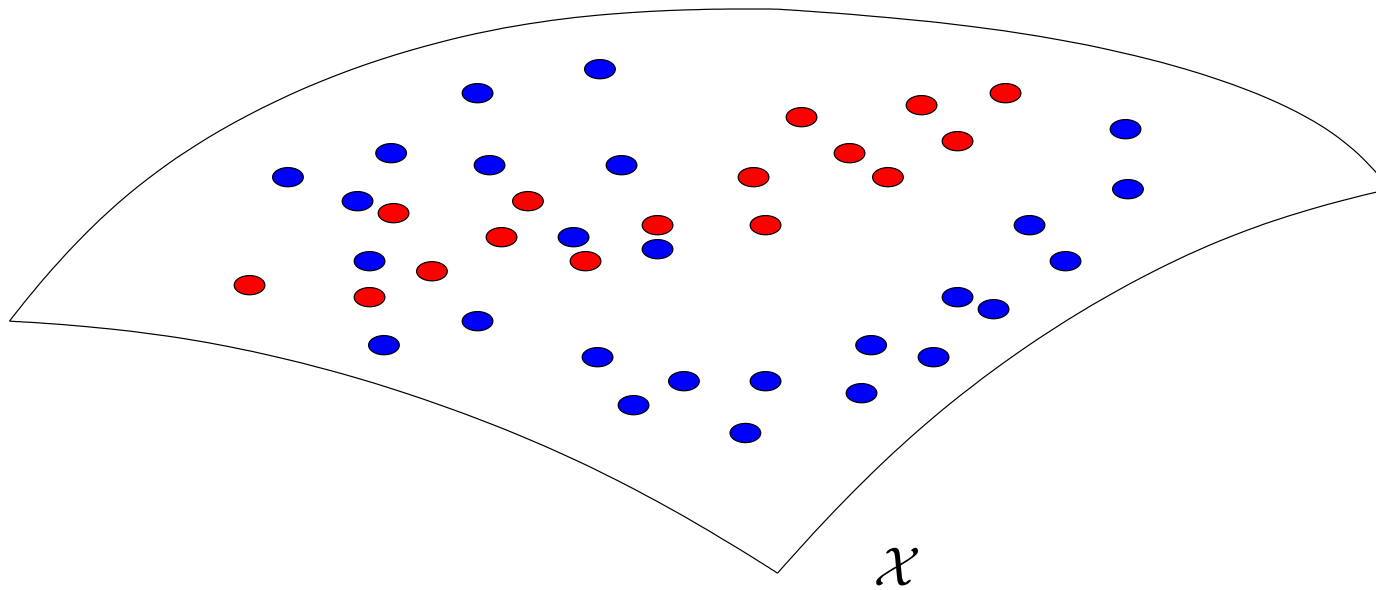
How does it work?

Example: k -nearest neighbors



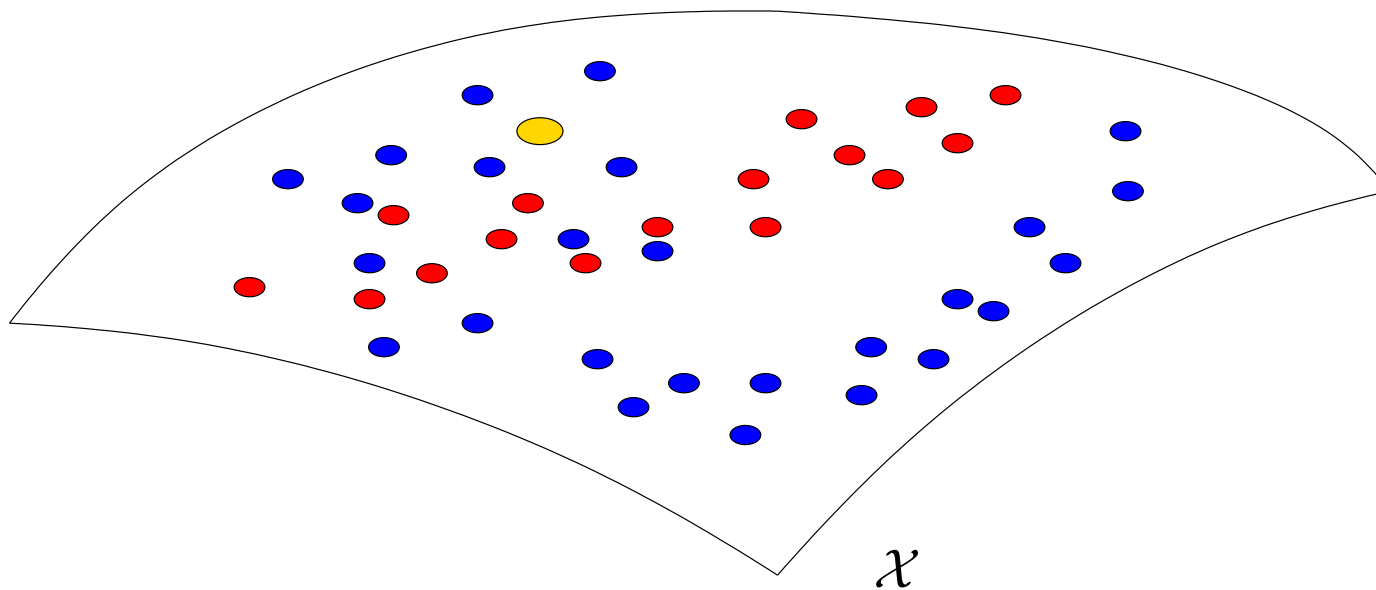
A database represented as a bunch of colored points

Example: k -nearest neighbors



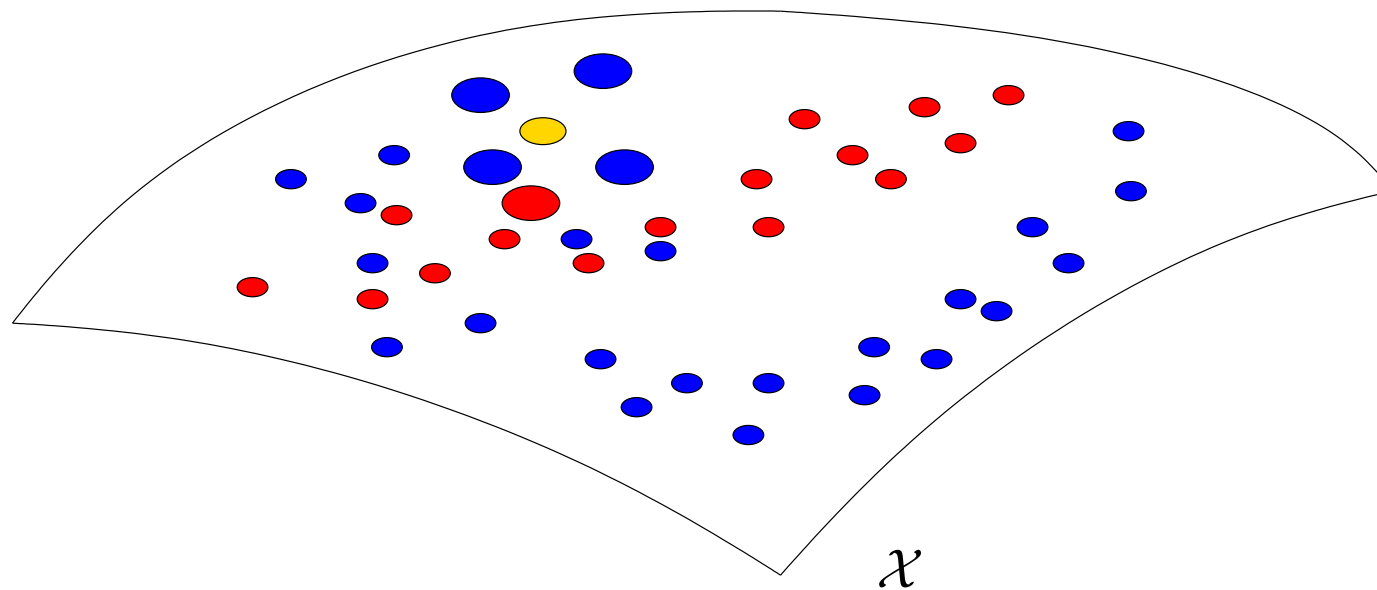
Points are split between **blue** and **red** points.

Example: k -nearest neighbors



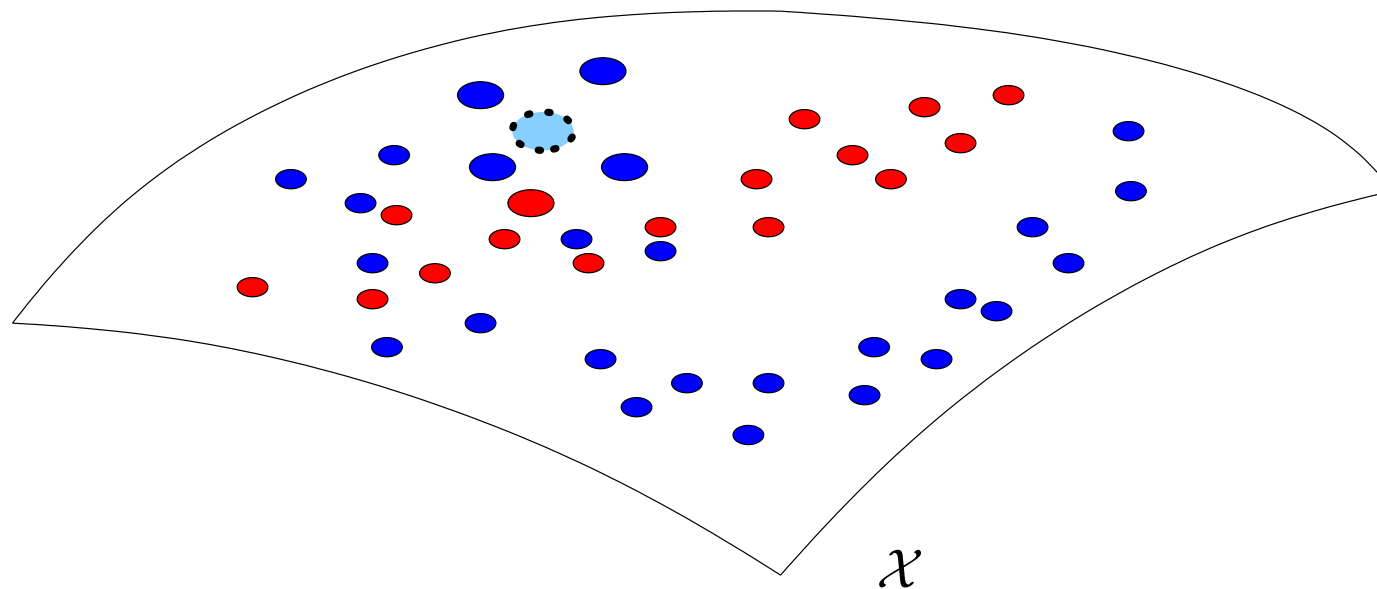
Here we have a new point to label... **blue** or **red**?

Example: k -nearest neighbors



Look at **5** closest neighbors... that is set $k = 5$.

Example: k -nearest neighbors



4 **blue** and 1 **red**: majority vote \rightarrow **5-nearest neighbors'** guess is **blue**.

k nearest neighbors

- A computer program cannot “see” the dataset the way we just did.
- Neither do we when $d > 3$...
- A computer would only rely on
 - the dataset $\{(\mathbf{x}_i, y_i), i = 1, \dots, N\}$.
 - the distance function d .
 - the parameter k .
- How would this work in practice?
- Suppose $N = 100$ and $k = 3$, labels are **blue** or **red**.

k nearest neighbors

- How would this work in practice?
 - get a new point x ... **we want to guess its label**

k nearest neighbors

- How would this work in practice?
 - get a new point \mathbf{x} .

- Compute a vector of distances $\begin{bmatrix} d(\mathbf{x}, \mathbf{x}_1) \\ d(\mathbf{x}, \mathbf{x}_2) \\ d(\mathbf{x}, \mathbf{x}_3) \\ \vdots \\ d(\mathbf{x}, \mathbf{x}_{100}) \end{bmatrix}$.

k nearest neighbors

- How would this work in practice?
 - get a new point \mathbf{x} .

- Compute a vector of distances
$$\begin{bmatrix} d(\mathbf{x}, \mathbf{x}_1) \\ d(\mathbf{x}, \mathbf{x}_2) \\ d(\mathbf{x}, \mathbf{x}_3) \\ \vdots \\ d(\mathbf{x}, \mathbf{x}_{100}) \end{bmatrix} = \begin{bmatrix} 0.32 \\ 4.56 \\ 0.112 \\ \vdots \\ 2.892 \end{bmatrix} .$$

k nearest neighbors

- How would this work in practice?

- get a new point \mathbf{x} .

- Compute a vector of distances
$$\begin{bmatrix} d(\mathbf{x}, \mathbf{x}_1) \\ d(\mathbf{x}, \mathbf{x}_2) \\ d(\mathbf{x}, \mathbf{x}_3) \\ \vdots \\ d(\mathbf{x}, \mathbf{x}_{100}) \end{bmatrix} = \begin{bmatrix} 0.32 \\ 4.56 \\ 0.112 \\ \vdots \\ 2.892 \end{bmatrix} .$$

- **Sort (at the top)** this vector of distances
$$\begin{bmatrix} d(\mathbf{x}, \mathbf{x}_3) \\ d(\mathbf{x}, \mathbf{x}_{86}) \\ d(\mathbf{x}, \mathbf{x}_{13}) \\ \vdots \\ d(\mathbf{x}, \mathbf{x}_2) \end{bmatrix} = \begin{bmatrix} 0.112 \\ 0.132 \\ 0.133 \\ \vdots \\ 4.56 \end{bmatrix} .$$

k nearest neighbors

- How would this work in practice?

- get a new point \mathbf{x} .

- Compute a vector of distances
$$\begin{bmatrix} d(\mathbf{x}, \mathbf{x}_1) \\ d(\mathbf{x}, \mathbf{x}_2) \\ d(\mathbf{x}, \mathbf{x}_3) \\ \vdots \\ d(\mathbf{x}, \mathbf{x}_{100}) \end{bmatrix} = \begin{bmatrix} 0.32 \\ 4.56 \\ 0.112 \\ \vdots \\ 2.892 \end{bmatrix} .$$

- **Sort (at the top)** this vector of distances
$$\begin{bmatrix} d(\mathbf{x}, \mathbf{x}_3) \\ d(\mathbf{x}, \mathbf{x}_{86}) \\ d(\mathbf{x}, \mathbf{x}_{13}) \\ \vdots \\ d(\mathbf{x}, \mathbf{x}_2) \end{bmatrix} = \begin{bmatrix} 0.112 \\ 0.132 \\ 0.133 \\ \vdots \\ 4.56 \end{bmatrix} .$$

- **Select** the three closest neighbors $\mathbf{x}_3, \mathbf{x}_{86}, \mathbf{x}_{13}$.

k nearest neighbors

- How would this work in practice?

- get a new point \mathbf{x} .

- Compute a vector of distances
$$\begin{bmatrix} d(\mathbf{x}, \mathbf{x}_1) \\ d(\mathbf{x}, \mathbf{x}_2) \\ d(\mathbf{x}, \mathbf{x}_3) \\ \vdots \\ d(\mathbf{x}, \mathbf{x}_{100}) \end{bmatrix} = \begin{bmatrix} 0.32 \\ 4.56 \\ 0.112 \\ \vdots \\ 2.892 \end{bmatrix} .$$

- **Sort (at the top)** this vector of distances
$$\begin{bmatrix} d(\mathbf{x}, \mathbf{x}_3) \\ d(\mathbf{x}, \mathbf{x}_{86}) \\ d(\mathbf{x}, \mathbf{x}_{13}) \\ \vdots \\ d(\mathbf{x}, \mathbf{x}_2) \end{bmatrix} = \begin{bmatrix} 0.112 \\ 0.132 \\ 0.133 \\ \vdots \\ 4.56 \end{bmatrix} .$$

- **Select** the three closest neighbors \mathbf{x}_3 , \mathbf{x}_{86} , \mathbf{x}_{13} with their labels.

k nearest neighbors

- How would this work in practice?

- get a new point \mathbf{x} .

- Compute a vector of distances
$$\begin{bmatrix} d(\mathbf{x}, \mathbf{x}_1) \\ d(\mathbf{x}, \mathbf{x}_2) \\ d(\mathbf{x}, \mathbf{x}_3) \\ \vdots \\ d(\mathbf{x}, \mathbf{x}_{100}) \end{bmatrix} = \begin{bmatrix} 0.32 \\ 4.56 \\ 0.112 \\ \vdots \\ 2.892 \end{bmatrix} .$$

- **Sort (at the top)** this vector of distances
$$\begin{bmatrix} d(\mathbf{x}, \mathbf{x}_3) \\ d(\mathbf{x}, \mathbf{x}_{86}) \\ d(\mathbf{x}, \mathbf{x}_{13}) \\ \vdots \\ d(\mathbf{x}, \mathbf{x}_2) \end{bmatrix} = \begin{bmatrix} 0.112 \\ 0.132 \\ 0.133 \\ \vdots \\ 4.56 \end{bmatrix} .$$

- **Select** the three closest neighbors \mathbf{x}_3 , \mathbf{x}_{86} , \mathbf{x}_{13} with their labels.

The 3-nearest neighbor algorithm outputs a **red** label for \mathbf{x} .

k nearest neighbors

3 blocks of k -nearest neighbor:

- database of colored points $\{ \mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_5, \mathbf{x}_6, \mathbf{x}_7, \mathbf{x}_8, \mathbf{x}_9, \dots, \mathbf{x}_{13}, \dots, \mathbf{x}_{100} \}$
- neighborhood size $k = 3$
- a distance function d .

the dataset is fixed. k is a simple parameter.
the **distance** is the most challenging to define.

Can we **fine tune** this distance so that we get improved results?

Choosing a Distance for k -nearest neighbors

Choosing a Distance

For many applications,
both training points $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ and observations \mathbf{x}
are vectors of **features**,

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} \in \mathbb{R}^d.$$

What distances can we use for vectors?

Choosing a Distance

- Euclidean/ l_2 distance in \mathbb{R}^d

$$d_2(\mathbf{x}, \mathbf{x}') = \sqrt{\sum_{i=1}^d (x_i - x'_i)^2}.$$

Choosing a Distance

- Euclidean/ l_2 distance in \mathbb{R}^d

$$d_2(\mathbf{x}, \mathbf{x}') = \sqrt{\sum_{i=1}^d (x_i - x'_i)^2} = \sqrt{(\mathbf{x} - \mathbf{x}')^T (\mathbf{x} - \mathbf{x}')}$$

- Manhattan/ l_1 distance in \mathbb{R}^d

$$d_1(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^d |x_i - x'_i|.$$

Choosing a Distance

- Euclidean/ l_2 distance in \mathbb{R}^d

$$d_2(\mathbf{x}, \mathbf{x}') = \sqrt{\sum_{i=1}^d (x_i - x'_i)^2}.$$

- Manhattan/ l_1 distance in \mathbb{R}^d

$$d_1(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^d |x_i - x'_i|.$$

- Infinite/ l_∞ distance in \mathbb{R}^d

$$d_\infty(\mathbf{x}, \mathbf{x}') = \max_{i=1..d} |x_i - x'_i|.$$

Choosing a Distance

- Euclidean/ l_2 distance in \mathbb{R}^d

$$d_2(\mathbf{x}, \mathbf{x}') = \sqrt{\sum_{i=1}^d (x_i - x'_i)^2}.$$

- Manhattan/ l_1 distance in \mathbb{R}^d

$$d_1(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^d |x_i - x'_i|.$$

- Infinite/ l_∞ distance in \mathbb{R}^d

$$d_\infty(\mathbf{x}, \mathbf{x}') = \max_{i=1..d} |x_i - x'_i|.$$

- the list of possibilities is **very long** [*Encyclopedia of Distances*, Deza & Deza, 2009]

Learning Distances

Rather than **choose**, **learn** using *data*.

Consider a
parameterized and **rich** enough **family of distances**
rather than a
long and predefined list of candidates

- always preferable if easy to implement and optimize...
- this is the goal of **metric learning**

Learning Distances: Starting Point

- The Euclidean distance,

$$d_2(\mathbf{x}, \mathbf{x}') = \sqrt{\sum_{i=1}^d (x_i - x'_i)^2} = \sqrt{(\mathbf{x} - \mathbf{x}')^T (\mathbf{x} - \mathbf{x}')}$$

can be seen as a particular case of **Mahalanobis distances** (40's),

$$d^\Omega(\mathbf{x}, \mathbf{x}') = \sqrt{(\mathbf{x} - \mathbf{x}')^T \Omega (\mathbf{x} - \mathbf{x}')}$$

a family of distances parameterized by a **positive definite matrix** Ω in $\mathbb{R}^{d \times d}$.

- Indeed, $d_2 = d^I$.

Positive Definiteness

$$d^{\Omega}(\mathbf{x}, \mathbf{x}') = \sqrt{(\mathbf{x} - \mathbf{x}')^T \Omega (\mathbf{x} - \mathbf{x}')}$$

a family of distances parameterized by a **positive definite matrix** Ω in $\mathbb{R}^{d \times d}$.

- Why require that Ω is symmetric positive definite?
 - if Ω is not positive definite $\rightarrow \exists \mathbf{x} \in \mathbb{R}^d \mid \mathbf{x}^T \Omega \mathbf{x} < 0$, then for any \mathbf{y} ,

$$d^{\Omega}(\mathbf{y} + \mathbf{x}, \mathbf{y}) = \mathbf{x}^T \Omega \mathbf{x} < 0 !$$

- If Ω is positive **semidefinite** and $\exists \mathbf{x} \neq 0 \mid \mathbf{x}^T \Omega \mathbf{x} = 0$, d^{Ω} is a **pseudo-metric**.
- This is usually ok for most algorithms.

Positive Definiteness

$$d^{\Omega}(\mathbf{x}, \mathbf{x}') = \sqrt{(\mathbf{x} - \mathbf{x}')^T \Omega (\mathbf{x} - \mathbf{x}')}$$

a family of distances parameterized by a **positive definite matrix** Ω in $\mathbb{R}^{d \times d}$.

- Note also that if Ω is positive definite, then $\exists L$ such that

$$\Omega = L^T L \text{ (Cholesky),}$$

then $d^{\Omega}(\mathbf{x}, \mathbf{x}') = d_2(L\mathbf{x}, L\mathbf{x}')$.

- Choosing a Mahalanobis distance \Leftrightarrow defining an arbitrary linear map L .

Learning Mahalanobis Distances

- Imagine you would like to use Mahalanobis distances,

$$d^{\Omega}(\mathbf{x}, \mathbf{x}') = \sqrt{(\mathbf{x} - \mathbf{x}')^T \Omega (\mathbf{x} - \mathbf{x}')}$$

question : **How to set Ω ?**

Learning Mahalanobis Distances

- Imagine you would like to use Mahalanobis distances,

$$d^{\Omega}(\mathbf{x}, \mathbf{x}') = \sqrt{(\mathbf{x} - \mathbf{x}')^T \Omega (\mathbf{x} - \mathbf{x}')}$$

question : **How to set Ω ?**

- Up to 10 years ago, basically set $\Omega = \hat{\Sigma}^{-1}$, inverse of empirical variance.
- Computed using a dataset $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ and

$$\hat{\Sigma} = \frac{1}{N-1} \sum_{i=1}^N \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^T,$$

where $\tilde{\mathbf{x}}_i \stackrel{\text{def}}{=} \mathbf{x}_i - \bar{\mathbf{x}}$ and $\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$

- also known as *whitening* the data.

Learning Mahalanobis Distances

- Imagine you would like to use Mahalanobis distances,

$$d^{\Omega}(\mathbf{x}, \mathbf{x}') = \sqrt{(\mathbf{x} - \mathbf{x}')^T \Omega (\mathbf{x} - \mathbf{x}')}$$

question : **How to set Ω ?**

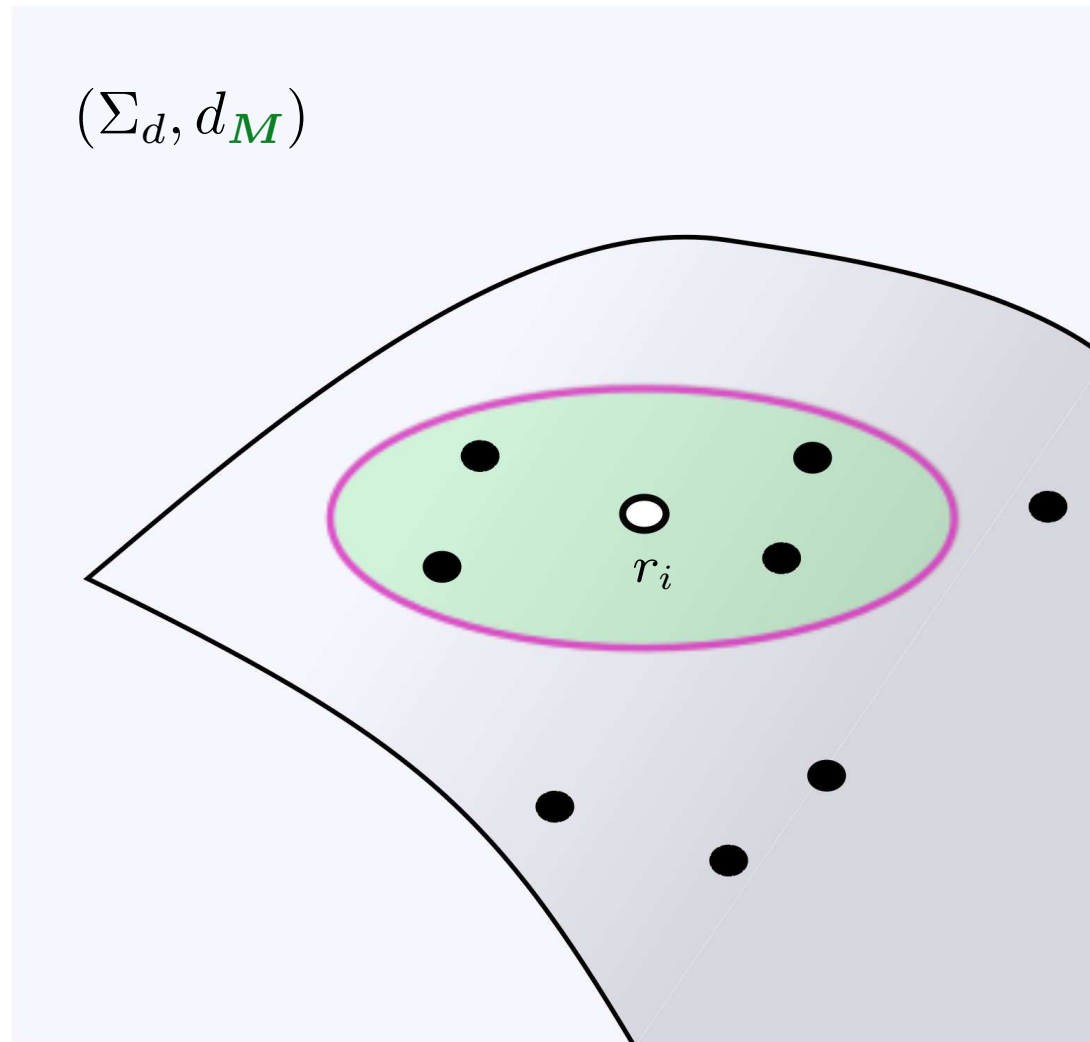
- **New answer in recent machine learning research** (>2003),
 - Distance Learning (Xing et al. 2003)
 - Pseudo-Metric Online Learning Algorithm (Shalev-Schwartz et al., 2004)
 - Large Margin Nearest Neighbor (Weinberger & Saul, 2006),
 - Information Theoretic Metric Learning (Kulis et al., 2007),
 - *etc.*

Works well for *many* datasets and vectors...

Basic Idea Behind Metric Learning

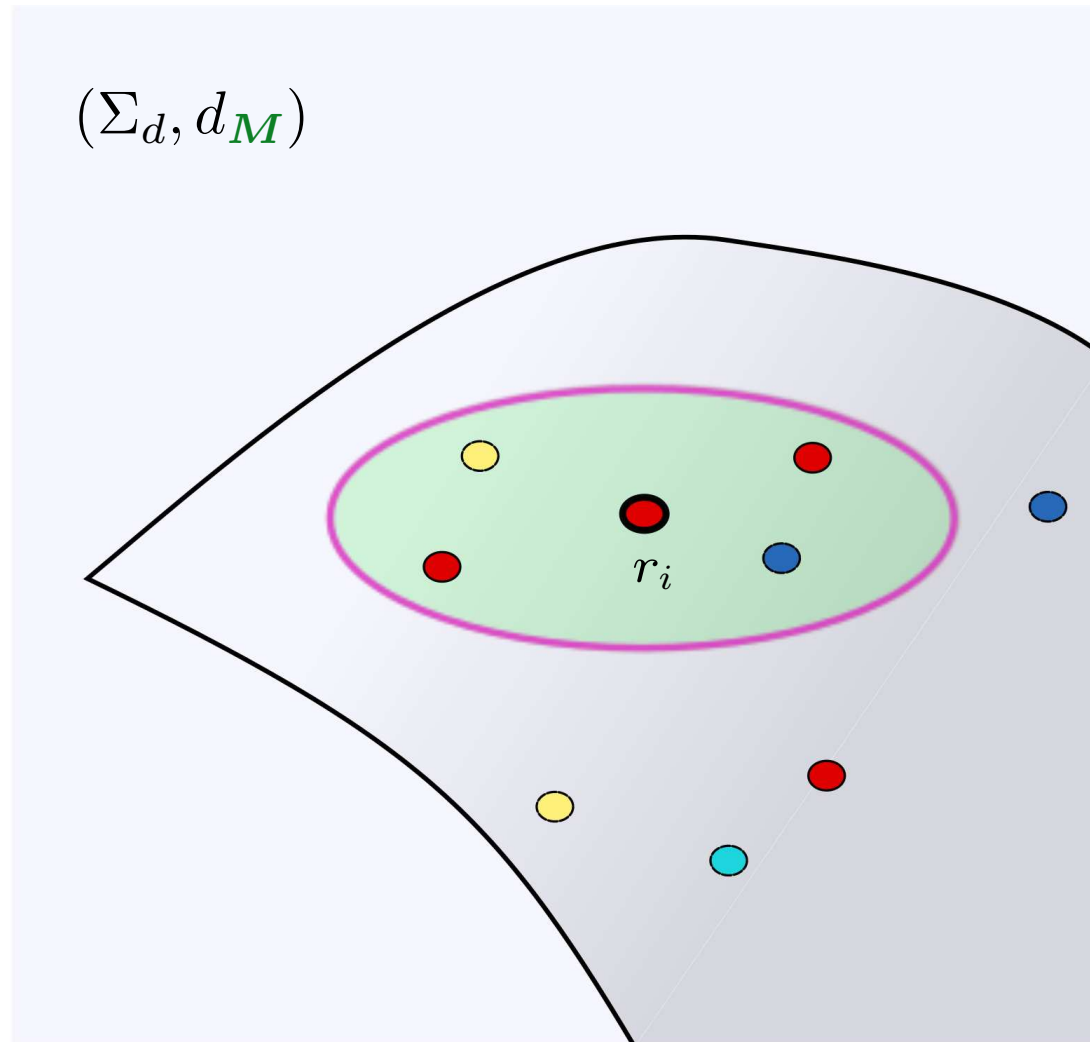
Learning a parameter M

Consider neighborhood of the i^{th} datum induced by d_M



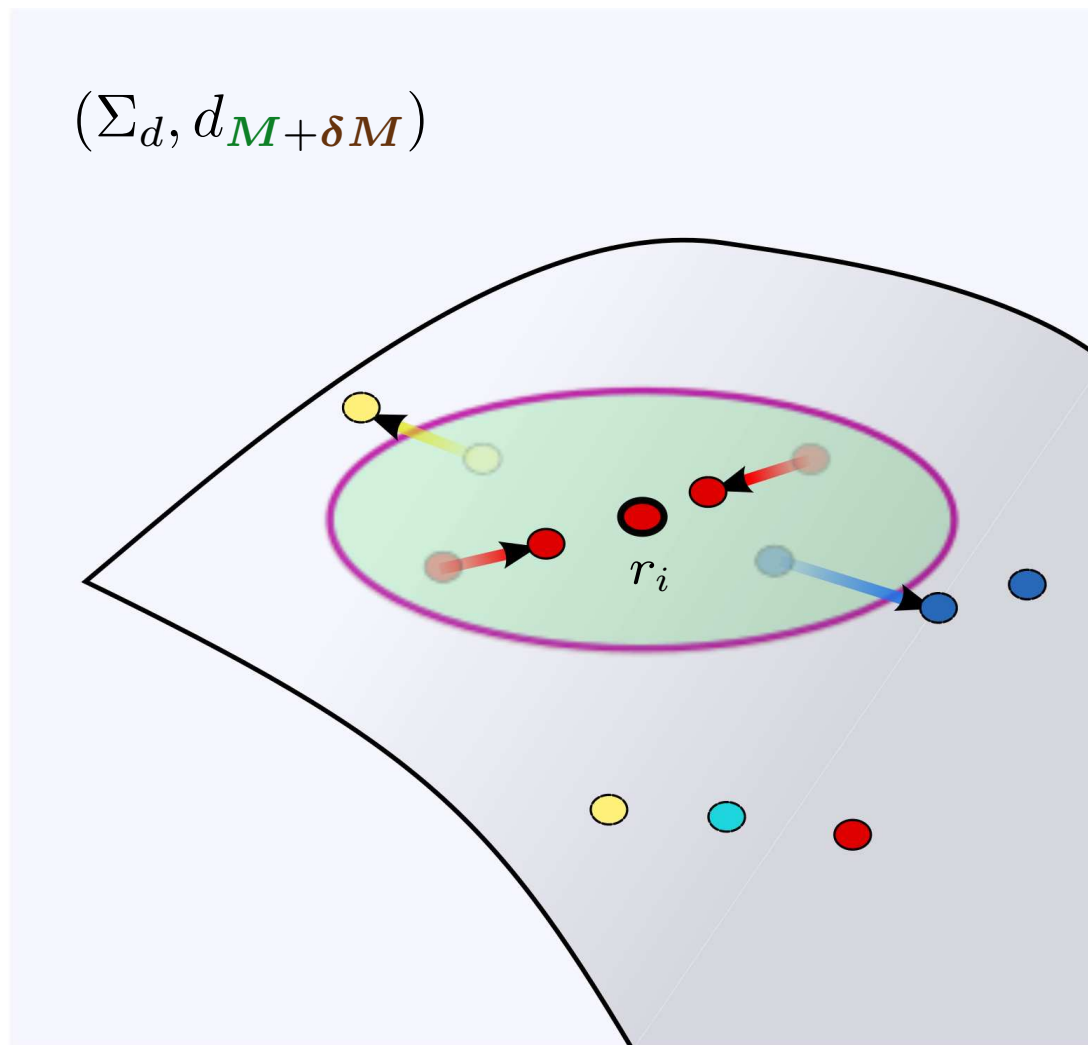
Learning a parameter M

Assume we have side-information, typically labels



Learning a parameter M

change M to $M + \delta M$ to improve local geometry



trick proposed in the metric learning literature: [Xing et al. '03, Weinberger Saul '06]

Mathematical Tools: Semidefinite Programming

Linear Programs

- $x \in \mathbb{R}^d, A \in \mathbb{R}^{d \times m}, \mathbf{b} \in \mathbb{R}^m$.
- Consider the **positive orthant** \mathbb{R}_+^d

$$\begin{array}{ll} \text{minimize} & \mathbf{c}^T \mathbf{x} \\ \text{subject to} & A\mathbf{x} \{=, \leq, \geq\} \mathbf{b} \\ & \mathbf{x} \in \mathbb{R}_+^d \end{array}$$

- History: simplex method (Danzig, '47), interior point methods (80's)

Conic Linear Program

- $x \in \mathbb{R}^d, A \in \mathbb{R}^{d \times m}, \mathbf{b} \in \mathbb{R}^m$.
- Consider an arbitrary **cone** \mathbf{K} .

$$\begin{array}{ll} \text{minimize} & \mathbf{c}^T \mathbf{x} \\ \text{subject to} & A\mathbf{x} \{=, \leq, \geq\} \mathbf{b} \\ & \mathbf{x} \in \mathbf{K} \end{array}$$

- History: 90's. Interior point methods.
- When \mathbf{K} is the cone of positive semidefinite matrices, **semidefinite program**.

In practice

Inputs: $\{ d^{\Omega}(\mathbf{x}_i, \mathbf{x}_j), \mathbf{x}_i, \mathbf{x}_j \in T \}$

In practice

Inputs: $\{ d^{\Omega}(\mathbf{x}_i, \mathbf{x}_j), \mathbf{x}_i, \mathbf{x}_j \in T \}$

- As a function of Ω , $d^{\Omega}(\mathbf{x}, \mathbf{x}')$ is **concave** w.r.t. Ω

$$d^{\Omega}(\mathbf{x}, \mathbf{x}') = \sqrt{\langle \Omega, (\mathbf{x} - \mathbf{x}')(\mathbf{x} - \mathbf{x}')^T \rangle}.$$

In practice

Inputs: $\{ d^{\Omega}(\mathbf{x}_i, \mathbf{x}_j), \mathbf{x}_i, \mathbf{x}_j \in T \}$

- As a function of Ω , $d^{\Omega}(\mathbf{x}, \mathbf{x}')$ is **concave** w.r.t. Ω

$$d^{\Omega}(\mathbf{x}, \mathbf{x}') = \sqrt{\langle \Omega, (\mathbf{x} - \mathbf{x}')(\mathbf{x} - \mathbf{x}')^T \rangle}.$$

- **squared** distance (**not the distance itself**) is linear w.r.t. Ω

$$d^{\Omega}(\mathbf{x}, \mathbf{x}')^2 = \langle \Omega, (\mathbf{x} - \mathbf{x}')(\mathbf{x} - \mathbf{x}')^T \rangle.$$

In practice

Inputs: $\{ d^{\Omega}(\mathbf{x}_i, \mathbf{x}_j), \mathbf{x}_i, \mathbf{x}_j \in T \}$

- **squared** distance (**not the distance itself**) is linear w.r.t. Ω

$$d^{\Omega}(\mathbf{x}, \mathbf{x}')^2 = \langle \Omega, (\mathbf{x} - \mathbf{x}')(\mathbf{x} - \mathbf{x}')^T \rangle.$$

- for two points \mathbf{x}, \mathbf{y} , constraints of the kind

$$d^{\Omega}(\mathbf{x}, \mathbf{y})^2 \geq \varepsilon \Leftrightarrow \langle \Omega, (\mathbf{x} - \mathbf{y})(\mathbf{x} - \mathbf{y})^T \rangle \geq \varepsilon$$

while for three points $\mathbf{x}, \mathbf{y}, \mathbf{z}$,

$$d^{\Omega}(\mathbf{x}, \mathbf{z}) \leq d^{\Omega}(\mathbf{x}, \mathbf{y}) \Leftrightarrow \langle \Omega, (\mathbf{x} - \mathbf{y})(\mathbf{x} - \mathbf{y})^T - (\mathbf{x} - \mathbf{z})(\mathbf{x} - \mathbf{z})^T \rangle \geq 0$$

In practice

Inputs: $\{ d^\Omega(\mathbf{x}_i, \mathbf{x}_j), \mathbf{x}_i, \mathbf{x}_j \in T \}$

- **squared** distance (**not the distance itself**) is linear w.r.t. Ω

$$d^\Omega(\mathbf{x}, \mathbf{x}')^2 = \langle \Omega, (\mathbf{x} - \mathbf{x}')(\mathbf{x} - \mathbf{x}')^T \rangle.$$

- for two points \mathbf{x}, \mathbf{y} , constraints of the kind

$$d^\Omega(\mathbf{x}, \mathbf{y})^2 \geq \varepsilon \Leftrightarrow \langle \Omega, (\mathbf{x} - \mathbf{y})(\mathbf{x} - \mathbf{y})^T \rangle \geq \varepsilon$$

while for three points $\mathbf{x}, \mathbf{y}, \mathbf{z}$,

$$d^\Omega(\mathbf{x}, \mathbf{z}) \leq d^\Omega(\mathbf{x}, \mathbf{y}) \Leftrightarrow \langle \Omega, (\mathbf{x} - \mathbf{y})(\mathbf{x} - \mathbf{y})^T - (\mathbf{x} - \mathbf{z})(\mathbf{x} - \mathbf{z})^T \rangle \geq 0$$

- Falls into a SDP setting... as long as **squared** distances are always considered...
- Some implementations proposed in the literature use dedicated solvers.

Examples

Metric Learning (Xing, Ng, Jordan, Russel 2003)

- Elementary idea:
 - pull all pairs of similar points **together** (\mathcal{S}),
 - push all pairs of dissimilar points **apart** (\mathcal{D})

$$\max_{\Omega \succeq 0} \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{S}} d^{\Omega}(\mathbf{x}_i, \mathbf{x}_j)^2$$

$$\text{such that } \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{D}} d^{\Omega}(\mathbf{x}_i, \mathbf{x}_j) \geq 1$$

- d^{Ω} is concave, so the constraint is convex
- The objective is linear in Ω .
- If use $d^{\Omega}(\cdot, \cdot)^2$ in the constraint, problem is degenerate.

Pseudo-Metric Online Learning Algorithm (POLA, 2004)

- Considers an online setting where pairs $\mathbf{x}_t, \mathbf{x}'_t$ that are similar ($y_t = 1$) or dissimilar ($y_t = -1$) are received sequentially at time t .
- A new Ω is selected at each iteration t depending on previous observations.
- Objective: minimize, up to T ,

$$L = \sum_{t=1}^T l_t(\Omega_t, \mathbf{b}_t)$$

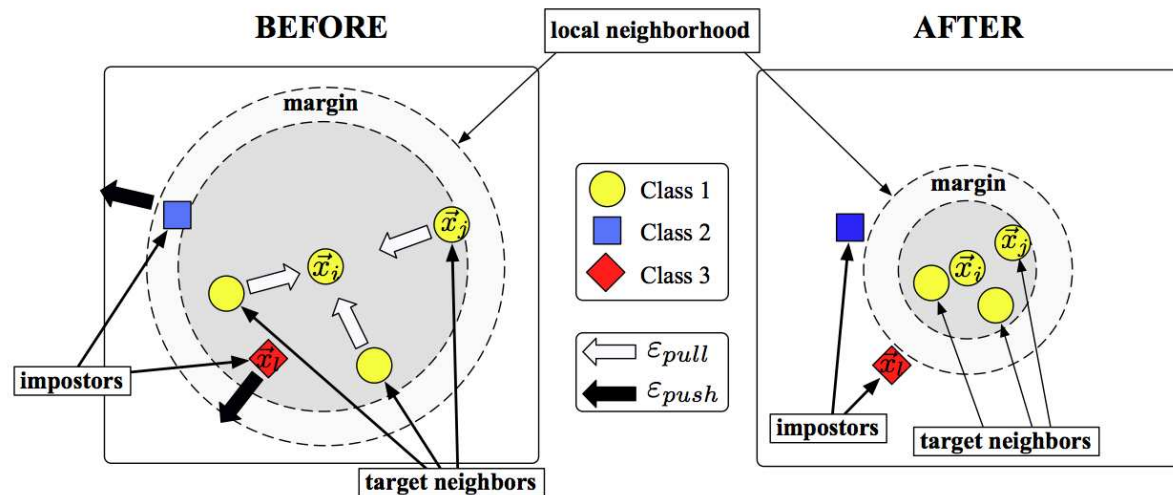
where

$$l_t(\Omega, \mathbf{b}) \stackrel{\text{def}}{=} \max(0, y_t(d^\Omega(\mathbf{x}_t, \mathbf{x}'_t))^2 - \mathbf{b}) + 1).$$

- To do so, use a simple updating step of Ω_t :
 - Let $\Omega_t \leftarrow \Omega_{t-1} - y_t \alpha_t (\mathbf{x}_t - \mathbf{x}'_t)(\mathbf{x}_t - \mathbf{x}'_t)^T$
 - Project on cone of positive definite matrices by thresholding

Large Margin Nearest Neighbor (LMNN, 2006)

- Original idea



- LMNN solves the following SDP:

$$\min_{\Omega} (1 - \mu) \sum_{i, j \rightsquigarrow i} (\mathbf{x}_i - \mathbf{x}_j)^T \Omega (\mathbf{x}_i - \mathbf{x}_j) + \mu \sum_{i, j \rightsquigarrow i, l} (1 - y_{il}) \xi_{ijl}$$

$$\text{such that } (\mathbf{x}_i - \mathbf{x}_j)^T \Omega (\mathbf{x}_i - \mathbf{x}_j) - (\mathbf{x}_i - \mathbf{x}_l)^T \Omega (\mathbf{x}_i - \mathbf{x}_l) \geq 1 - \xi_{ijl}$$

$$\xi_{ijl} \geq 0$$

Information Theoretic Metric Learning (ITML, 2007)

- Suppose we have a **prior candidate** for Ω (e.g. $\hat{\Sigma}^{-1}$)
- Call this prior candidate Ω_0 .
- ITML proposes to find an Ω that is not too far from Ω_0 while still splitting similar and dissimilar points apart:

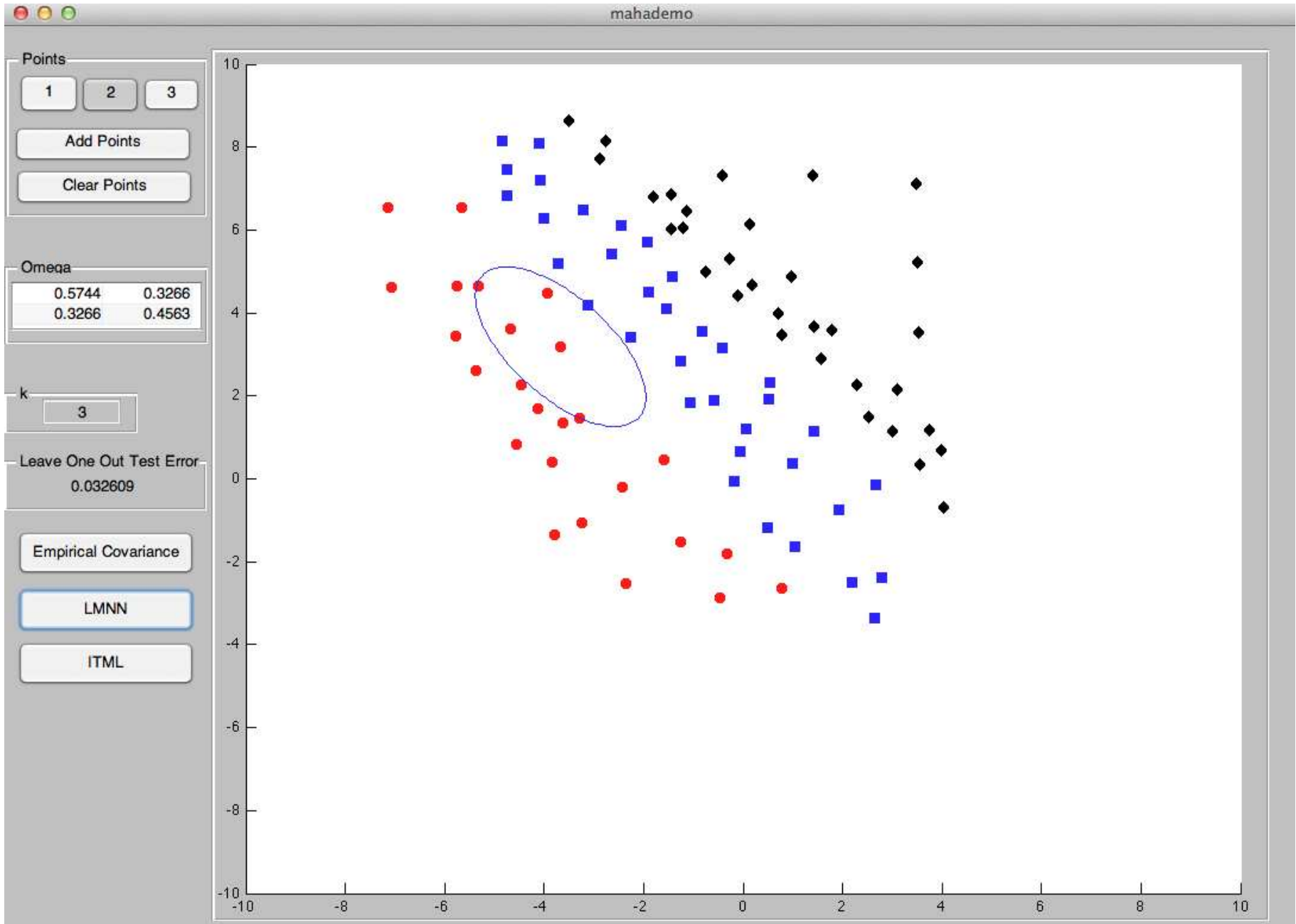
$$\min_{\Omega} \mathcal{L}(\Omega, \Omega_0)$$

such that for $(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{S}$, $d^{\Omega}(\mathbf{x}_i, \mathbf{x}_j)^2 \leq u$

such that for $(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{D}$, $d^{\Omega}(\mathbf{x}_i, \mathbf{x}_j) \geq l$

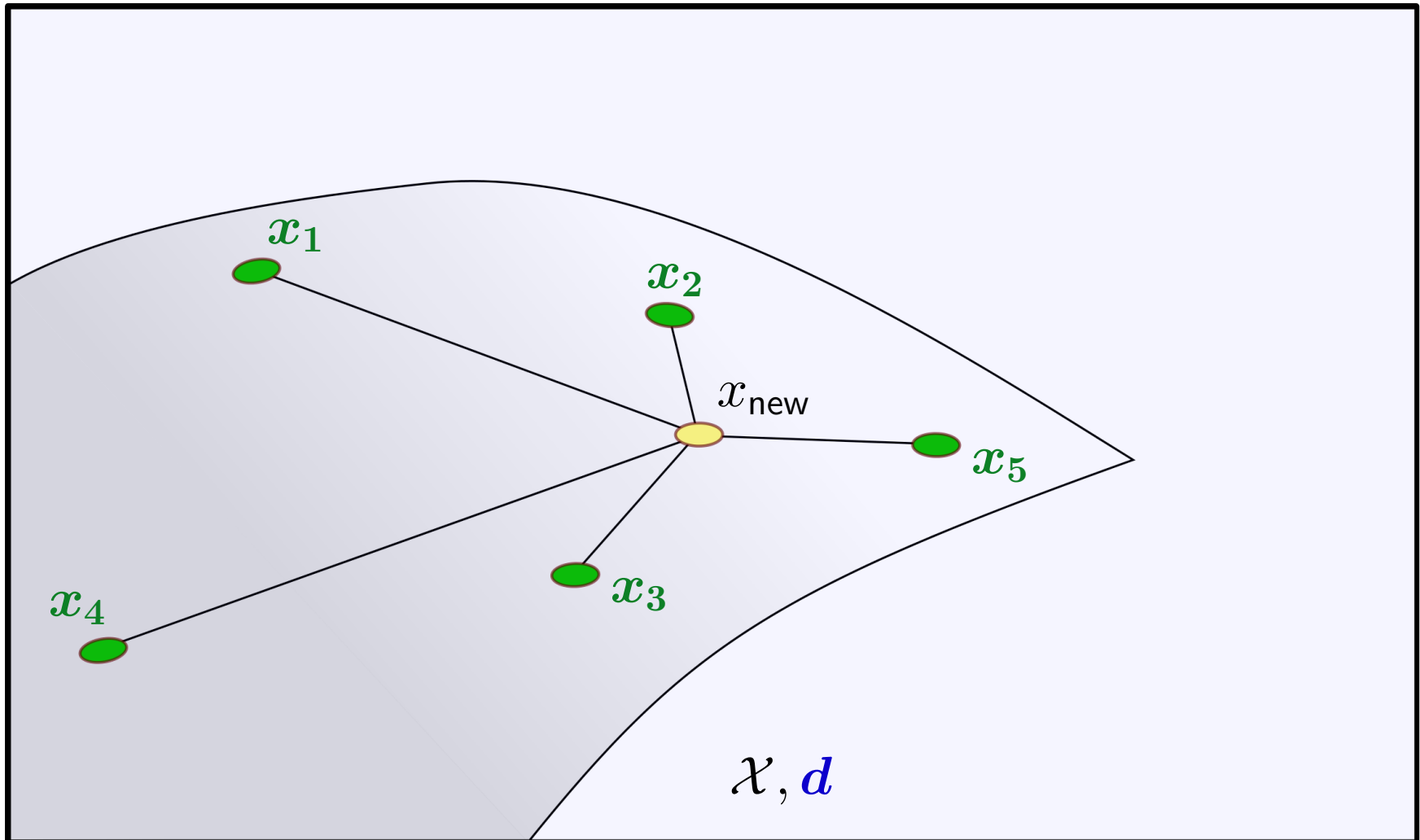
where $\mathcal{L}(A, B) \stackrel{\text{def}}{=} \text{tr}(AB^{-1}) - \log \det(AB^{-1})$ is a matrix-divergence

Illustration



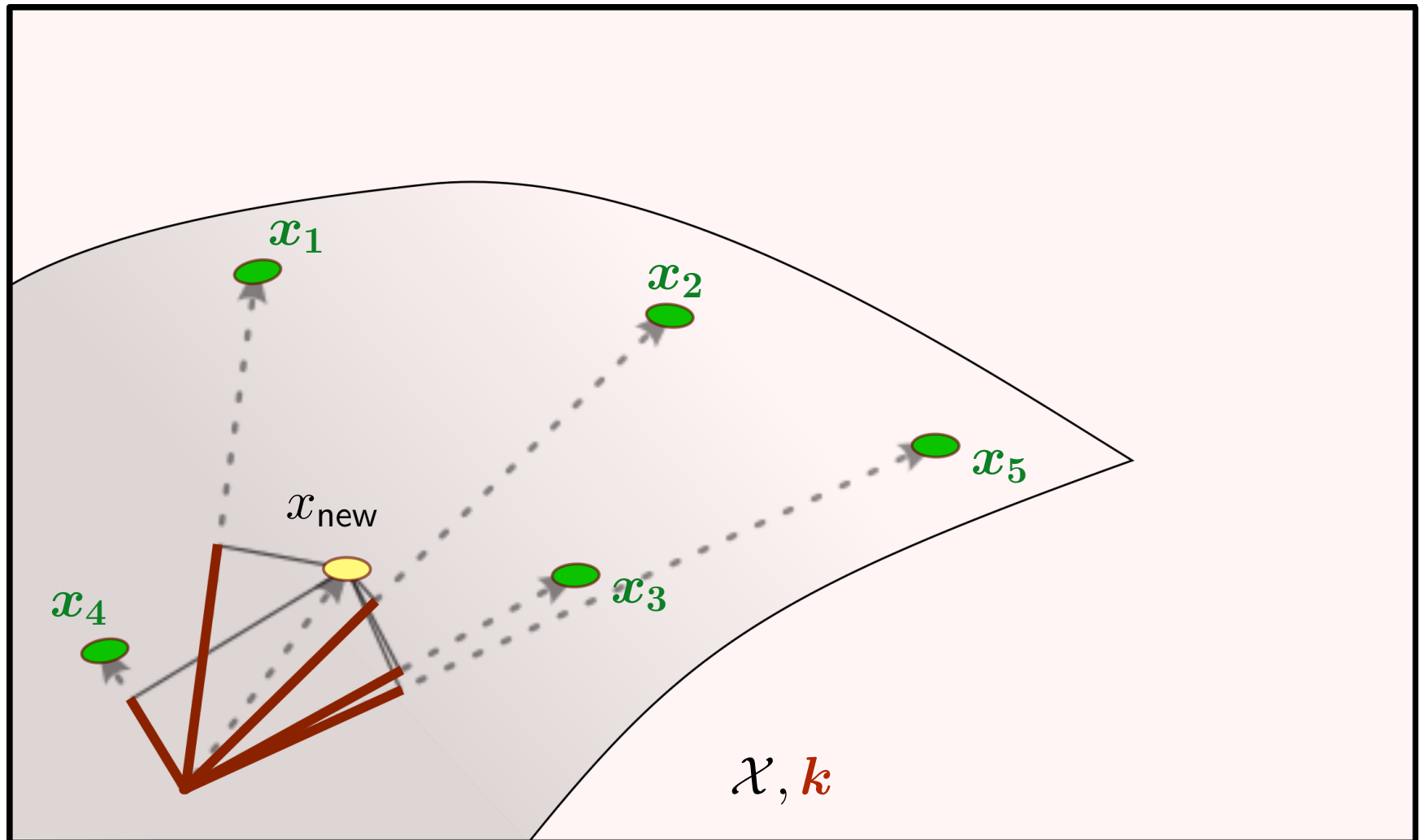
Kernels

While Distance Based Algorithms...



... rely **exclusively** on all distances $\{d(x_{\text{new}}, x_i), i = 1, \dots, 5\}$

Kernel Based Algorithms...



... rely exclusively on **dot-products** $k(x_{\text{new}}, x_i)$, $i = 1, \dots, 5$

Positive Definite Kernels

A bivariate function defined on a set \mathcal{X}

$$\begin{aligned} \mathbf{k} : \mathcal{X} \times \mathcal{X} &\rightarrow \mathbb{R}_+ \\ (\mathbf{x}, \mathbf{y}) &\mapsto \mathbf{k}(\mathbf{x}, \mathbf{y}) \end{aligned}$$

is a **positive definite kernel** if $\forall \mathbf{x}, \mathbf{y} \in \mathcal{X}$,

- $\mathbf{k}(\mathbf{x}, \mathbf{y}) = \mathbf{k}(\mathbf{y}, \mathbf{x})$, \rightarrow *symmetry*
- and $\forall n \in \mathbb{N}, \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \in \mathcal{X}^n, c \in \mathbb{R}^n$,

$$\sum_{i=1}^n c_i c_j \mathbf{k}(\mathbf{x}_i, \mathbf{x}_j) \geq 0, \rightarrow \text{positive definiteness}$$

$$\text{equivalently } K = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \cdots & k(\mathbf{x}_1, \mathbf{x}_i) & \cdots & k(\mathbf{x}_1, \mathbf{x}_n) \\ \vdots & \ddots & \vdots & \vdots & \vdots \\ k(\mathbf{x}_i, \mathbf{x}_1) & \cdots & k(\mathbf{x}_i, \mathbf{x}_i) & \cdots & k(\mathbf{x}_2, \mathbf{x}_n) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & \cdots & k(\mathbf{x}_n, \mathbf{x}_i) & \cdots & k(\mathbf{x}_n, \mathbf{x}_n) \end{bmatrix} \succeq 0$$

is positive semidefinite (has a nonnegative spectrum).

In the context of these lectures...

- A kernel k is a function

$$\begin{aligned} k : \mathcal{X} \times \mathcal{X} &\longmapsto \mathbb{R} \\ (\mathbf{x}, \mathbf{y}) &\longrightarrow k(\mathbf{x}, \mathbf{y}) \end{aligned}$$

- which compares two objects of a space \mathcal{X} , *e.g.*...

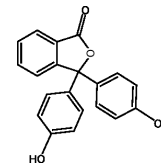
- strings, texts and sequences,



- images, audio and video feeds,



- graphs, interaction networks and 3D structures



- whatever actually... time-series of graphs of images? graphs of texts?...

What can we do with a kernel?

The setting

- Pretty simple setting: a set of objects $\mathbf{x}_1, \dots, \mathbf{x}_n$ of \mathcal{X}
- **Sometimes** additional information on these objects
 - labels $\mathbf{y}_i \in \{-1, 1\}$ or $\{1, \dots, \#(\text{classes})\}$,
 - scalar values $\mathbf{y}_i \in \mathbb{R}$,
 - associated object $\mathbf{y}_i \in \mathcal{Y}$
- A kernel $k : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$.

A few intuitions on the possibilities of kernel methods

Important concepts and perspectives

- The functional perspective: represent **points as functions**.
- **Nonlinearity** : linear combination of kernel evaluations.
- Summary of a sample through its **kernel matrix**.

Represent any point in \mathcal{X} as a function

For every \mathbf{x} , the map
 $\mathbf{x} \longrightarrow k(\mathbf{x}, \cdot)$
associates to \mathbf{x} a function $k(\mathbf{x}, \cdot)$ from \mathcal{X} to \mathbb{R} .

- Suppose we have a kernel k on bird images



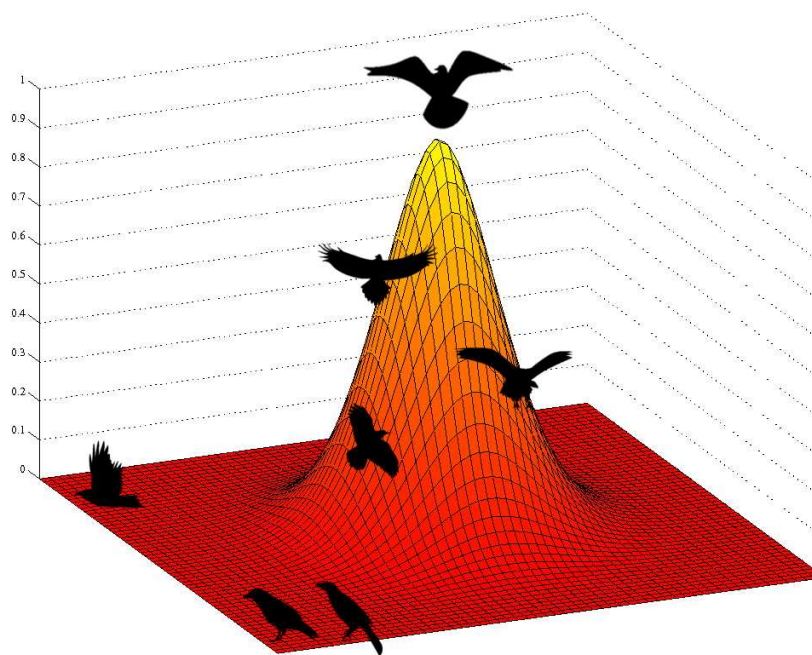
- Suppose for instance

$$k \left(\text{bird in flight}, \text{bird on ground} \right) = .32$$

Represent any point in \mathcal{X} as a function



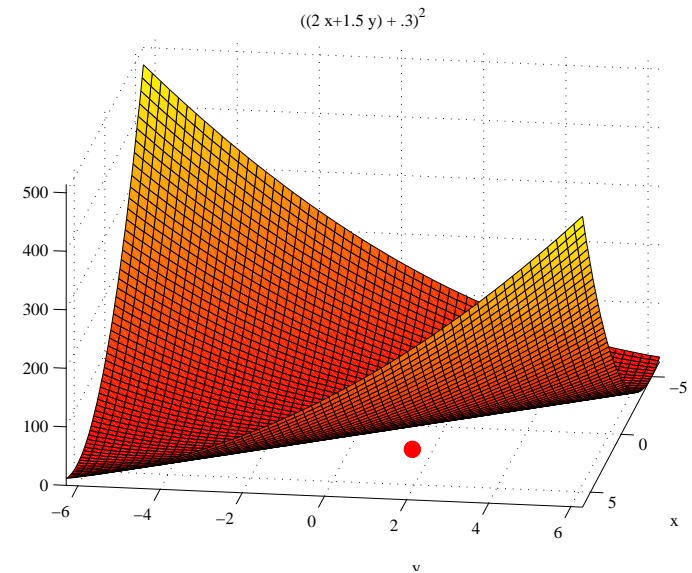
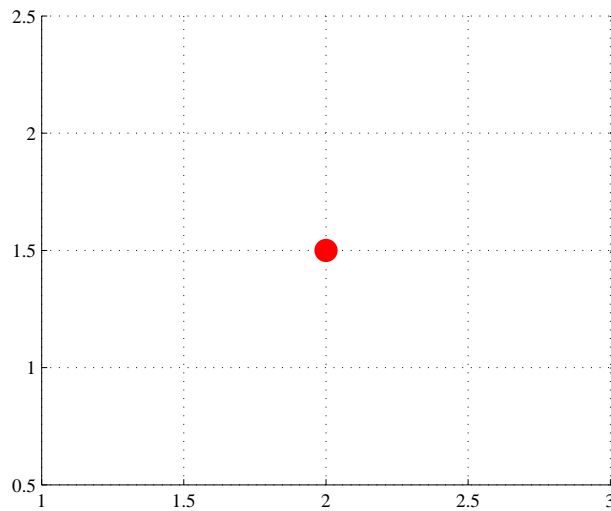
- We examine one image in particular:
- With kernels, we get a **representation** of that bird as a real-valued function, defined on the space of birds, represented here as \mathbb{R}^2 for simplicity.



schematic plot of $k(\text{bird}, \cdot)$.

Represent any point in \mathcal{X} as a function

- If the bird example was confusing...
- $k\left(\begin{bmatrix} x \\ y \end{bmatrix}, \begin{bmatrix} x' \\ y' \end{bmatrix}\right) = \left(\begin{bmatrix} x & y \end{bmatrix} \begin{bmatrix} x' \\ y' \end{bmatrix} + .3\right)^2$
- From a point in \mathbb{R}^2 to a function defined over \mathbb{R}^2 .



- We assume implicitly that the **functional representation** will be more useful than the **original representation**.

Decision functions as linear combination of kernel evaluations

- Linear decision functions are a major tool in statistics, that is functions

$$f(\mathbf{x}) = \beta^T \mathbf{x} + \beta_0.$$

- Implicitly, a point \mathbf{x} is processed depending on its characteristics x_i ,

$$f(\mathbf{x}) = \sum_{i=1}^d \beta_i x_i + \beta_0.$$

the free parameters are scalars $\beta_0, \beta_1, \dots, \beta_d$.

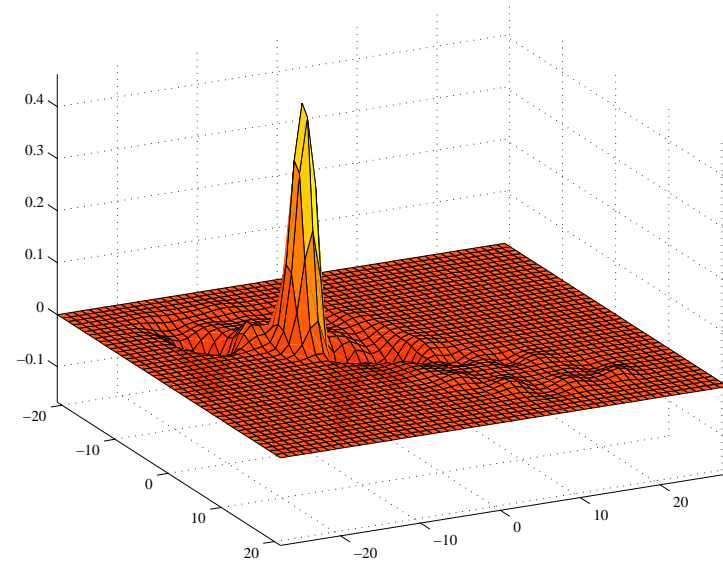
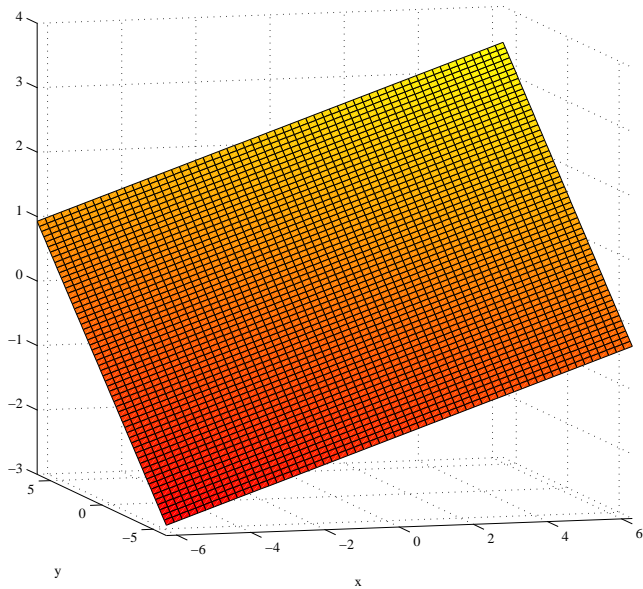
- Kernel methods yield candidate decision functions

$$f(\mathbf{x}) = \sum_{j=1}^n \alpha_j k(\mathbf{x}_j, \mathbf{x}) + \alpha_0.$$

the free parameters are scalars $\alpha_0, \alpha_1, \dots, \alpha_n$.

Decision functions as linear combination of kernel evaluations

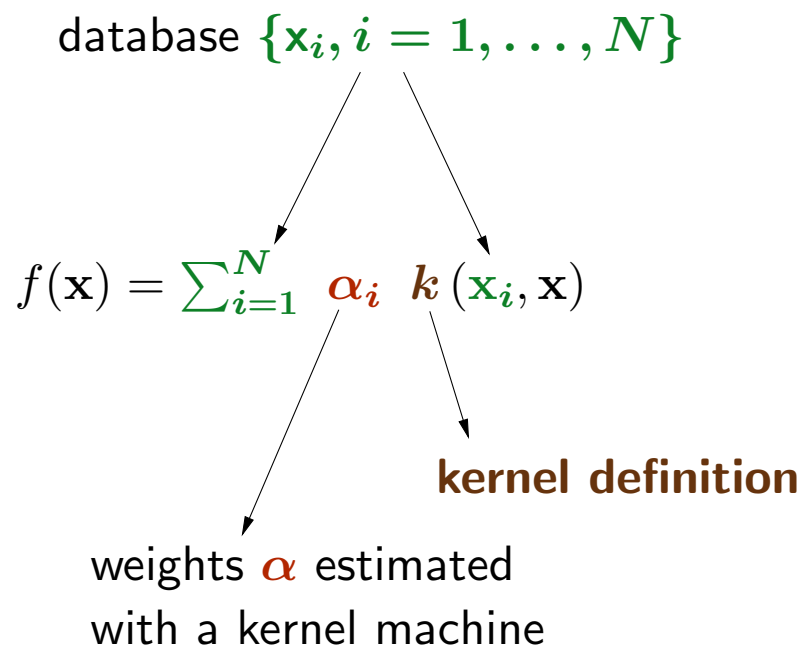
- linear decision surface / linear expansion of **kernel surfaces** (here $k_G(\mathbf{x}_i, \cdot)$)



- Kernel methods are considered **non-linear** tools.
- Yet not completely “nonlinear” → only one-layer of nonlinearity.

kernel methods use the data as a functional base to define decision functions

Decision functions as linear combination of kernel evaluations



- f is any predictive function of interest of a new point \mathbf{x} .
- Weights α are **optimized** with a kernel machine (*e.g.* support vector machine)

intuitively, kernel methods provide decisions based on how *similar* a point \mathbf{x} is to each instance of the training set

The Gram matrix perspective

- Imagine a little task: you have read 100 novels so far.



- You would like to know whether you will enjoy reading a **new** novel.
- A few options:
 - read the book...
 - have friends read it for you, read reviews.
 - try to guess, based on the novels you read, if you will like it

The Gram matrix perspective

Two distinct approaches

- Define what **features** can characterize a book.
 - Map each book in the library onto vectors



$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$$

typically the x_i 's can describe...

- ▷ # pages, language, year 1st published, country,
 - ▷ coordinates of the main action, keyword counts,
 - ▷ author's prizes, popularity, booksellers ranking
- Challenge: find a decision function using 100 ratings and features.

The Gram matrix perspective

- Define what makes **two novels similar**,
 - Define a kernel k which quantifies novel similarities.
 - Map the library onto a Gram matrix



$$\longrightarrow K = \begin{bmatrix} k(b_1, b_1) & k(b_1, b_2) & \cdots & k(b_1, b_{100}) \\ k(b_2, b_1) & k(b_2, b_2) & \cdots & k(b_2, b_{100}) \\ \vdots & \vdots & \cdots & \vdots \\ k(b_n, b_1) & k(b_n, b_2) & \cdots & k(b_{100}, b_{100}) \end{bmatrix}$$

- Challenge: find a decision function that takes this 100×100 matrix as an input.

The Gram matrix perspective

Given a new novel,

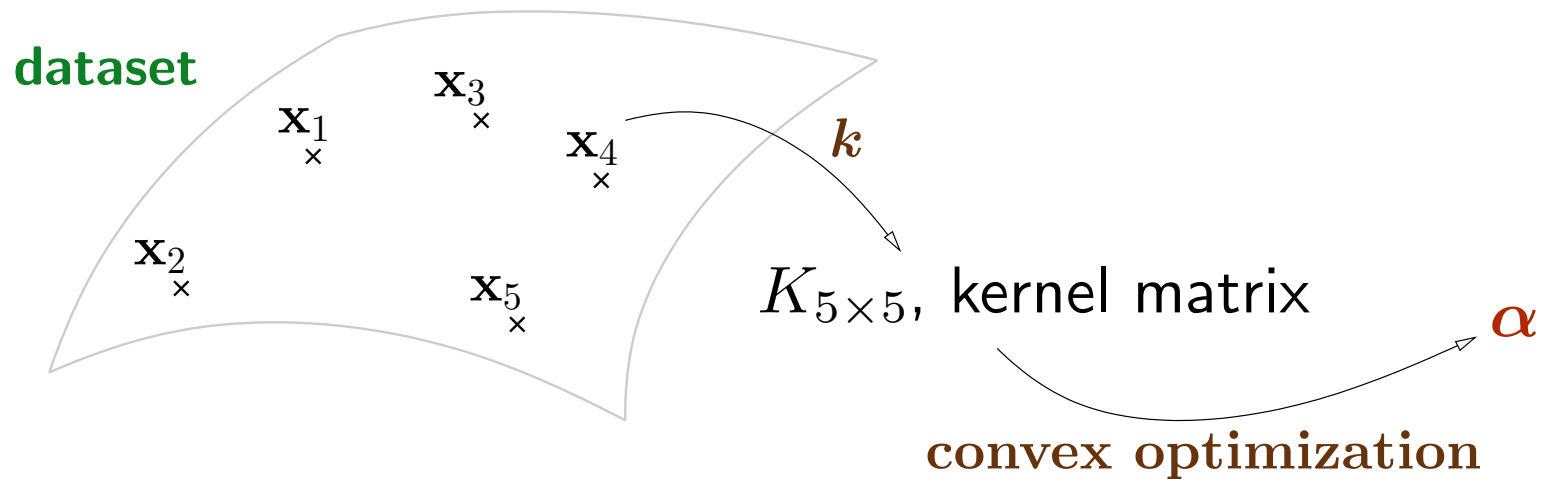
- with the **features approach**, the prediction can be rephrased as **what are the features of this new book?** what **features** have I found in the past that were good indicators of my taste?
- with the **kernel approach**, the prediction is rephrased as **which novels this book is similar or dissimilar to?** what **pool of books** did I find the most influentials to define my tastes accurately?

kernel methods **only use kernel similarities**, do not consider features.

Features can help define similarities, but **never considered elsewhere**.

The Gram matrix perspective

in kernel methods, clear separation between the kernel...



and **Convex optimization** (thanks to psdness of K , more later) to output the α 's.

Example of a Kernel Machine: Classification with the Support Vector Machine

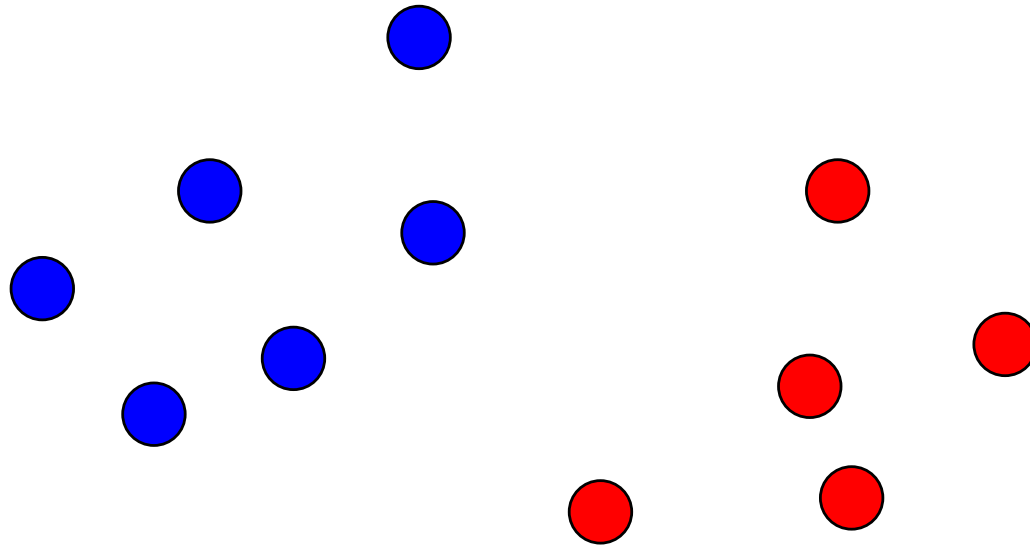
Data

- **Data** : vectors $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_N$.
- Ideally, to infer a “yes/no” rule, we need the **correct answer** for each vector.
- We consider thus a set of **pairs of (vector,bit)**

$$\text{“training set”} = \left\{ \left(\mathbf{x}_i = \begin{bmatrix} x_1^i \\ x_2^i \\ \vdots \\ x_d^i \end{bmatrix} \in \mathbb{R}^d, \mathbf{y}_i \in \{0, 1\} \right)_{i=1..N} \right\}$$

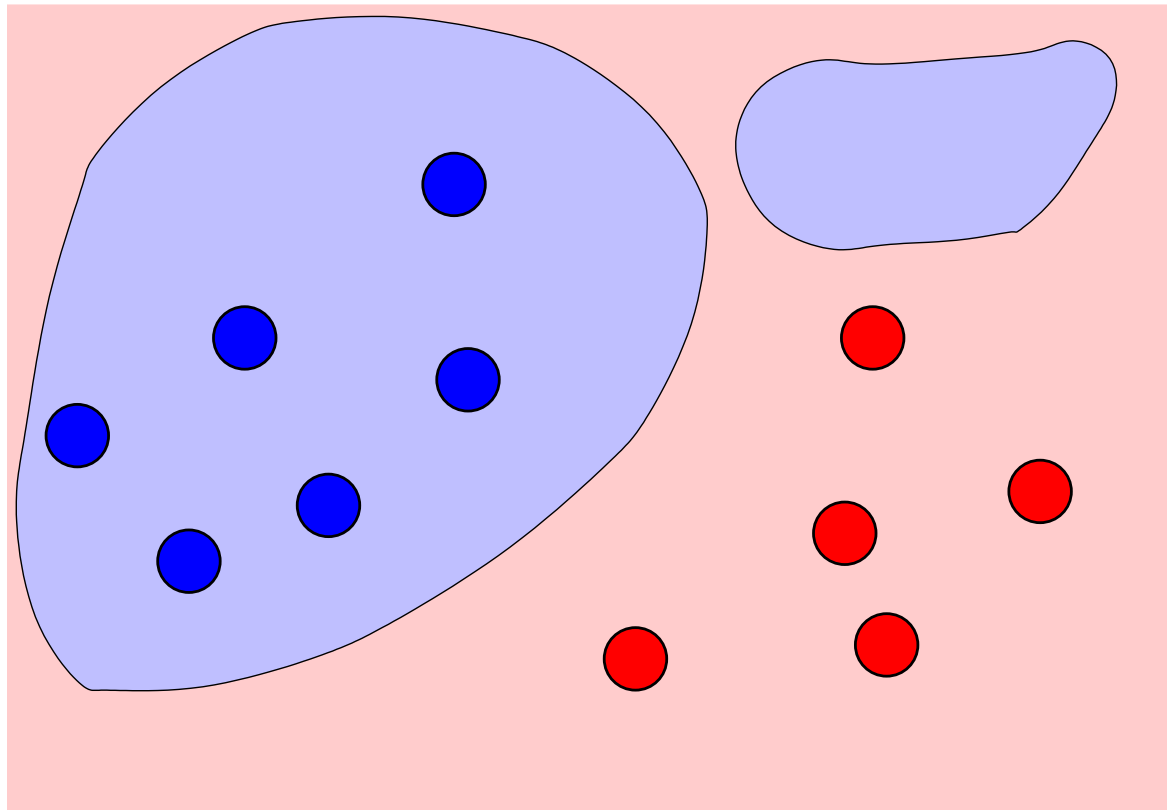
- For illustration purposes **only** we will consider **vectors in the plane**, $d = 2$.
- Points are easier to represent in 2 dimensions than in 20.000...
- The ideas for $d \gg 3$ are **exactly the same**.

Binary Classification Separation Surfaces for Vectors



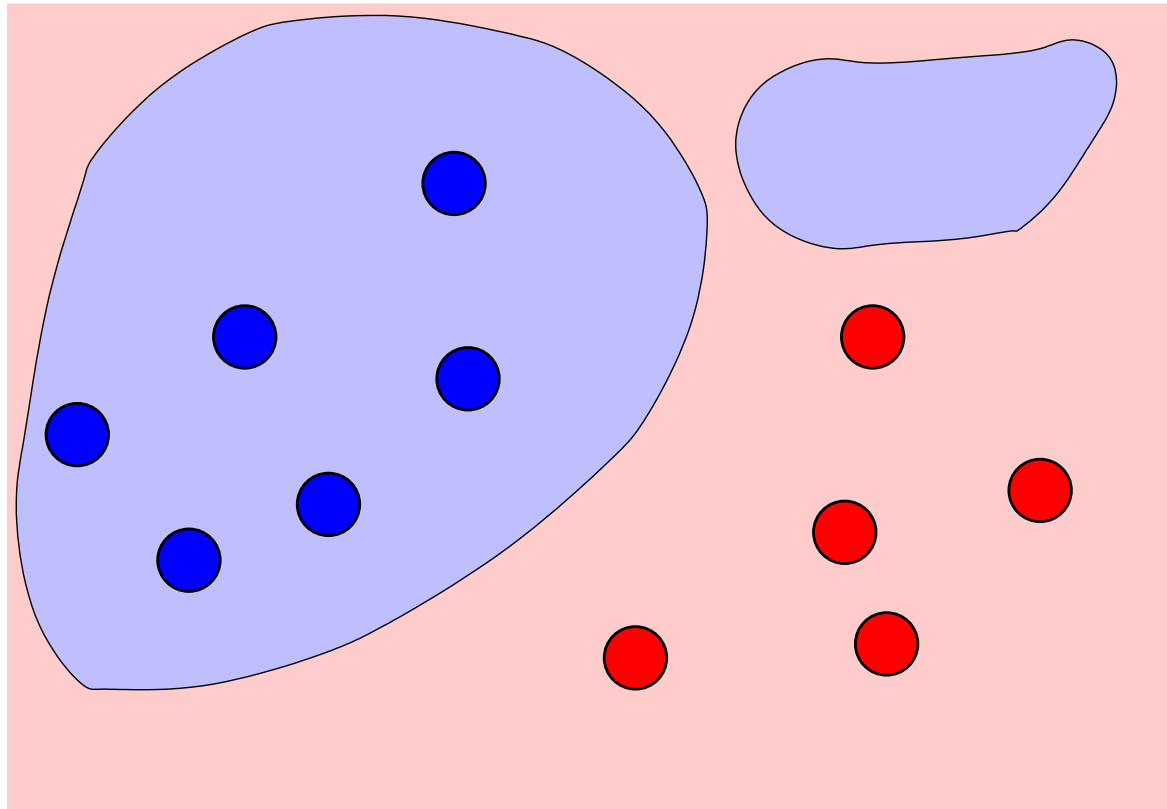
What is a classification rule?

Binary Classification Separation Surfaces for Vectors



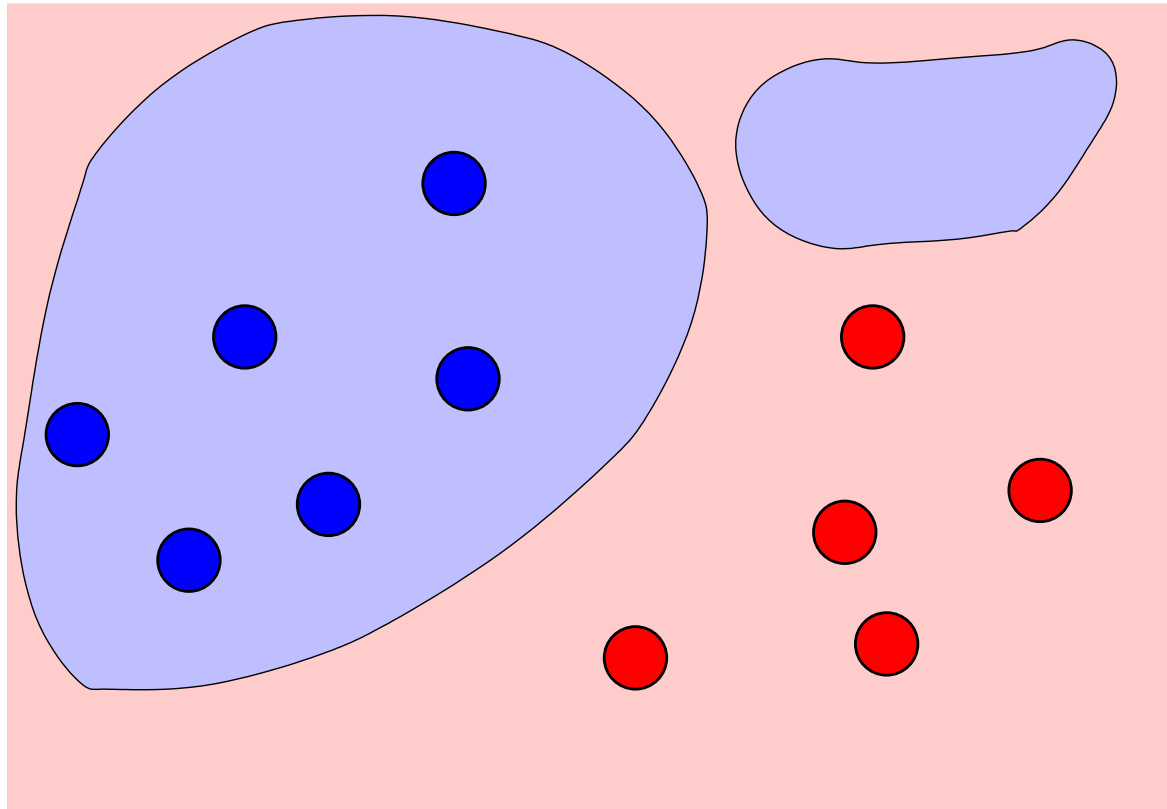
Classification rule = a partition of \mathbb{R}^d into two sets

Binary Classification Separation Surfaces for Vectors



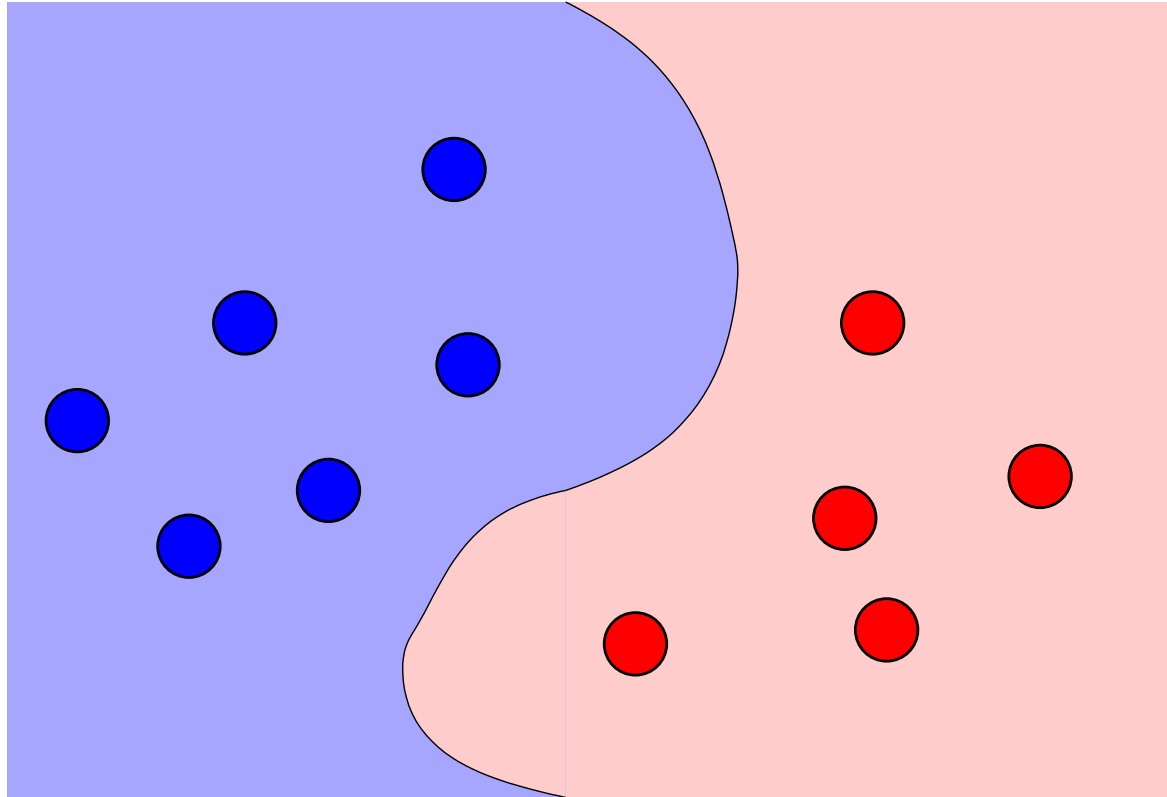
This partition is usually interpreted as the level set of function on \mathbb{R}^d

Binary Classification Separation Surfaces for Vectors



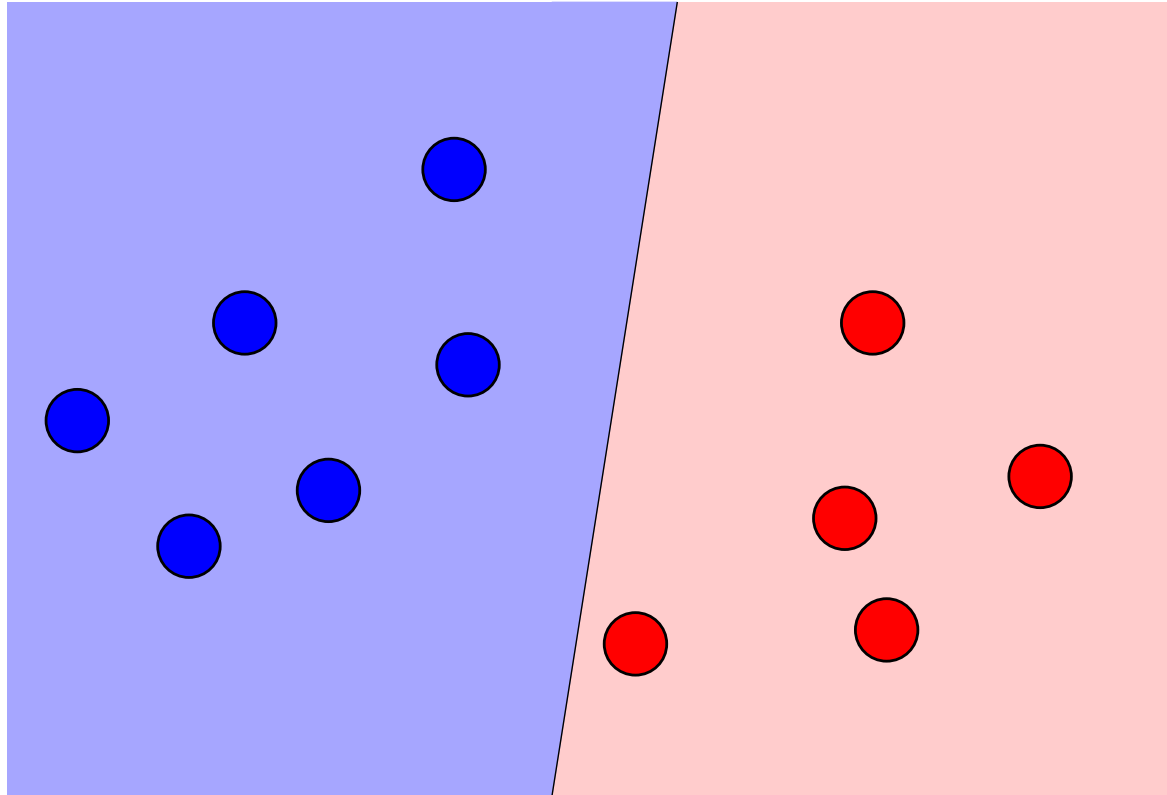
Typically, $\{\mathbf{x} \in \mathbb{R}^d \mid f(\mathbf{x}) > 0\}$ and $\{\mathbf{x} \in \mathbb{R}^d \mid f(\mathbf{x}) \leq 0\}$

Classification Separation Surfaces for Vectors



Can be defined by a single surface, *e.g.* a curved line

Classification Separation Surfaces for Vectors



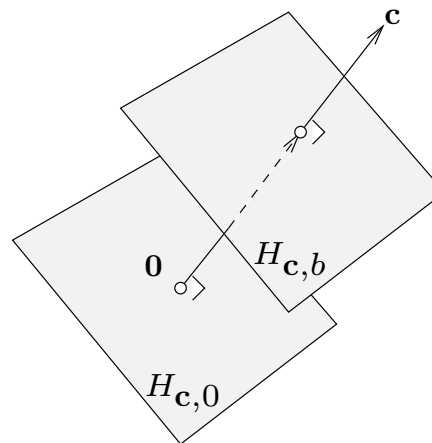
Even more **simple**: using **straight lines** and halfspaces.

Linear Classifiers

- **Straight lines** (hyperplanes when $d > 2$) are **the simplest type** of classifiers.
- A hyperplane $H_{\mathbf{c},b}$ is a set in \mathbb{R}^d defined by
 - a normal vector $\mathbf{c} \in \mathbb{R}^d$
 - a constant $b \in \mathbb{R}$. as

$$H_{\mathbf{c},b} = \{\mathbf{x} \in \mathbb{R}^d \mid \mathbf{c}^T \mathbf{x} = b\}$$

- Letting b vary we can “slide” the hyperplane across \mathbb{R}^d

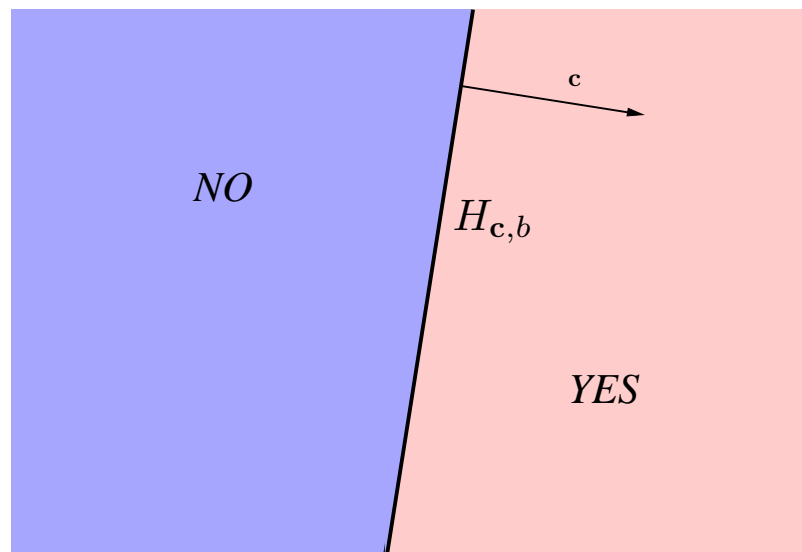


Linear Classifiers

- Exactly like lines in the plane, hypersurfaces **divide** \mathbb{R}^d into **two** halfspaces,

$$\{\mathbf{x} \in \mathbb{R}^d \mid \mathbf{c}^T \mathbf{x} < b\} \cup \{\mathbf{x} \in \mathbb{R}^d \mid \mathbf{c}^T \mathbf{x} \geq b\} = \mathbb{R}^d$$

- Linear classifiers attribute the “yes” and “no” answers given arbitrary \mathbf{c} and b .



- Assuming we only look at halfspaces for the decision surface...
...how to choose the **“best”** (\mathbf{c}^*, b^*) given a training sample?

Linear Classifiers

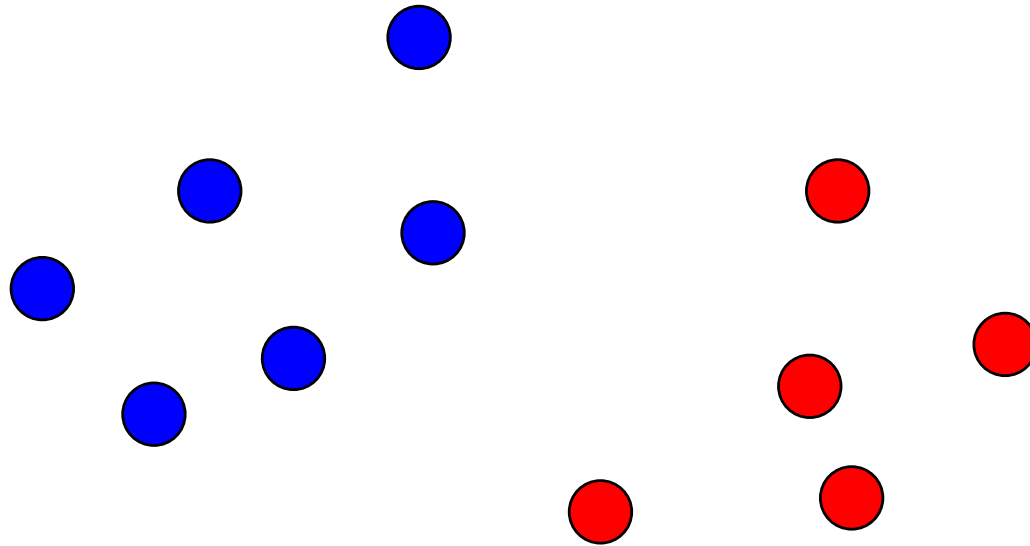
- This specific question,

“training set” $\{(\mathbf{x}_i \in \mathbb{R}^d, \mathbf{y}_i \in \{0, 1\})_{i=1..N}\} \xrightarrow{????}$ “best” \mathbf{c}^*, b^*

has different answers. Depends on the meaning of “best” ?:

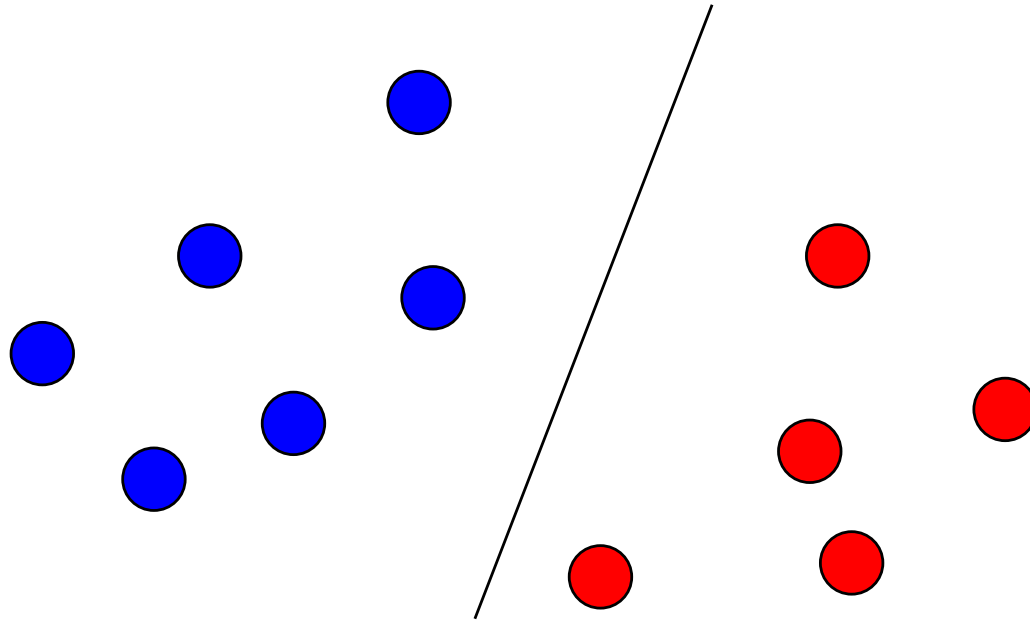
- **Linear Discriminant Analysis** (or Fisher’s Linear Discriminant);
- **Logistic regression** maximum likelihood estimation;
- **Perceptron**, a one-layer neural network;
- **Support Vector Machine**, the result of a convex program
- *etc.*

Classification Separation Surfaces for Vectors



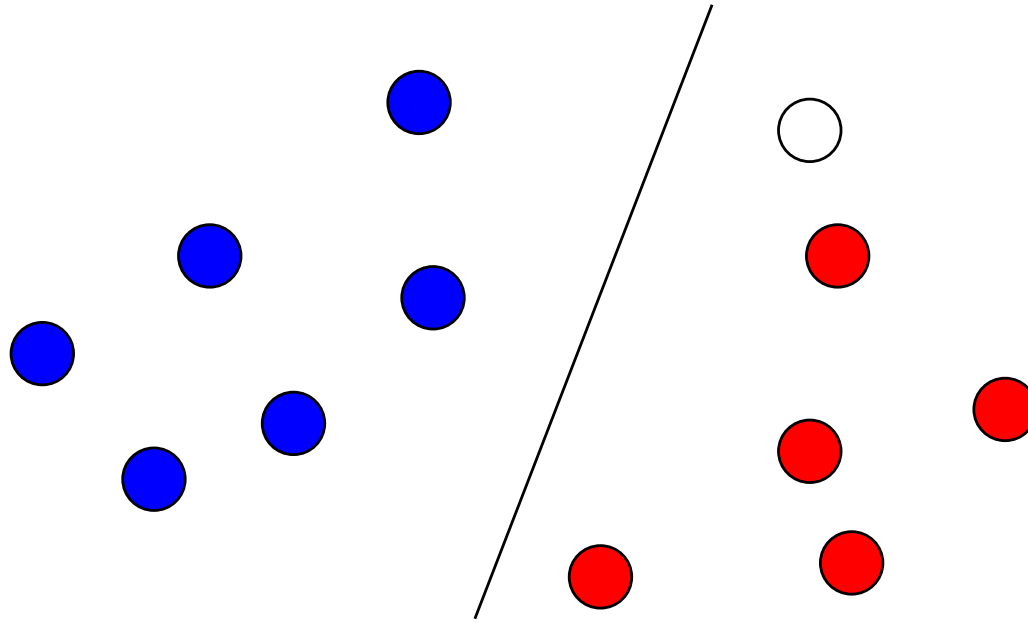
Given two sets of points...

Classification Separation Surfaces for Vectors



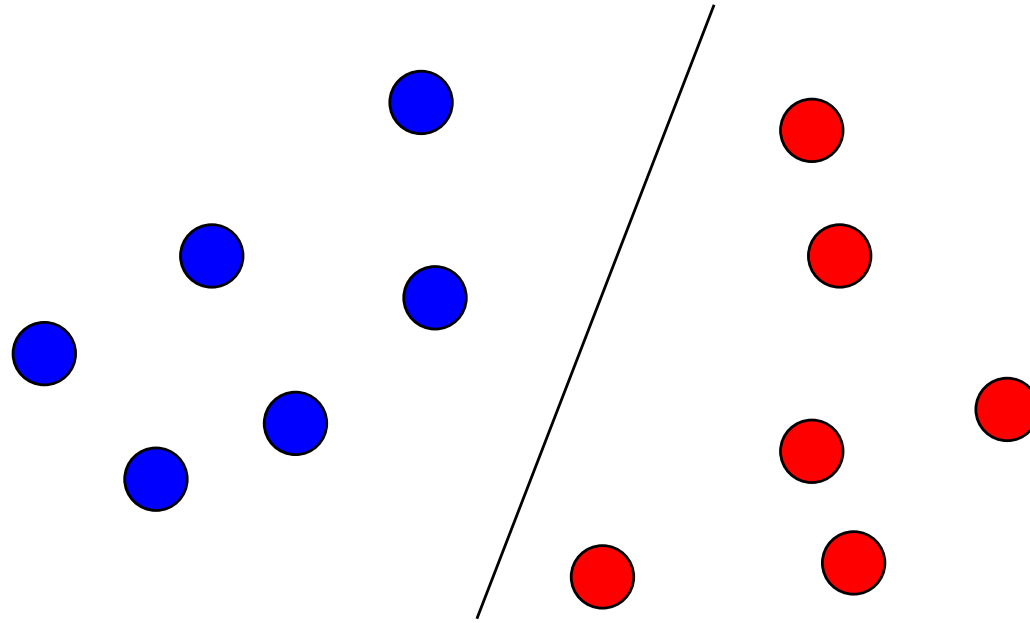
It is sometimes possible to separate them perfectly

Classification Separation Surfaces for Vectors

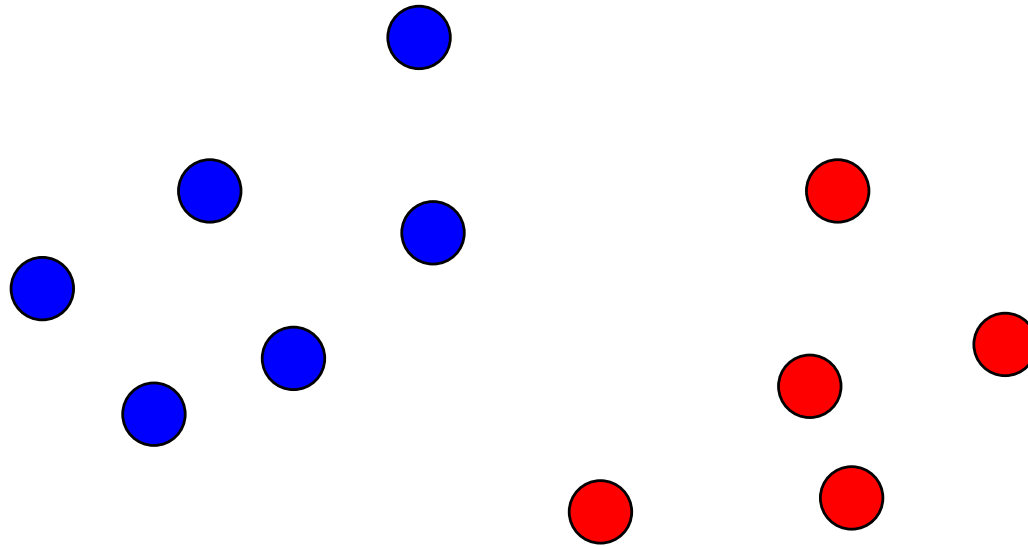


Each choice might look equivalently good on the training set,
but it will have obvious impact on new points

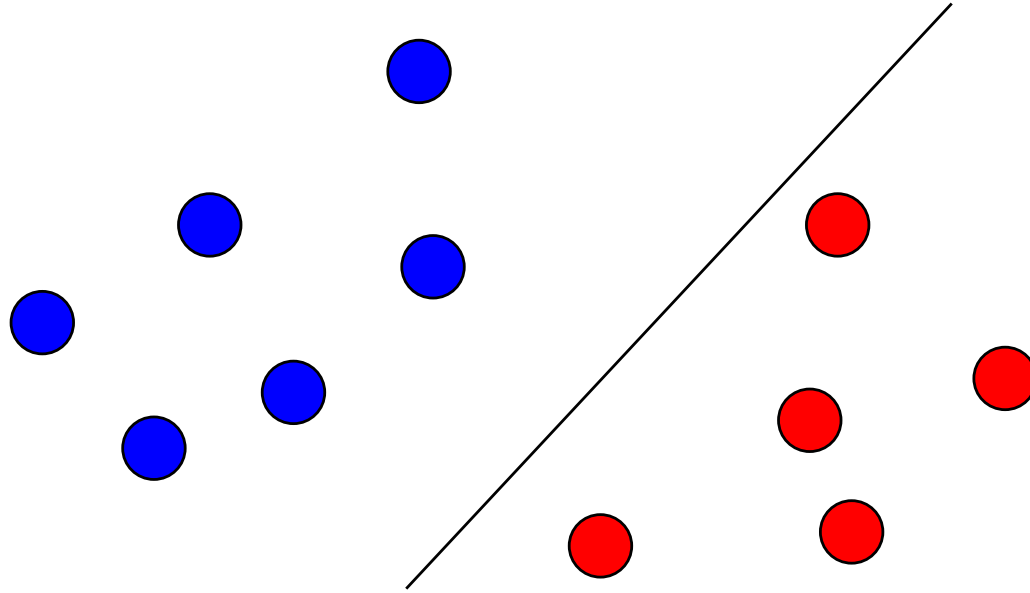
Classification Separation Surfaces for Vectors



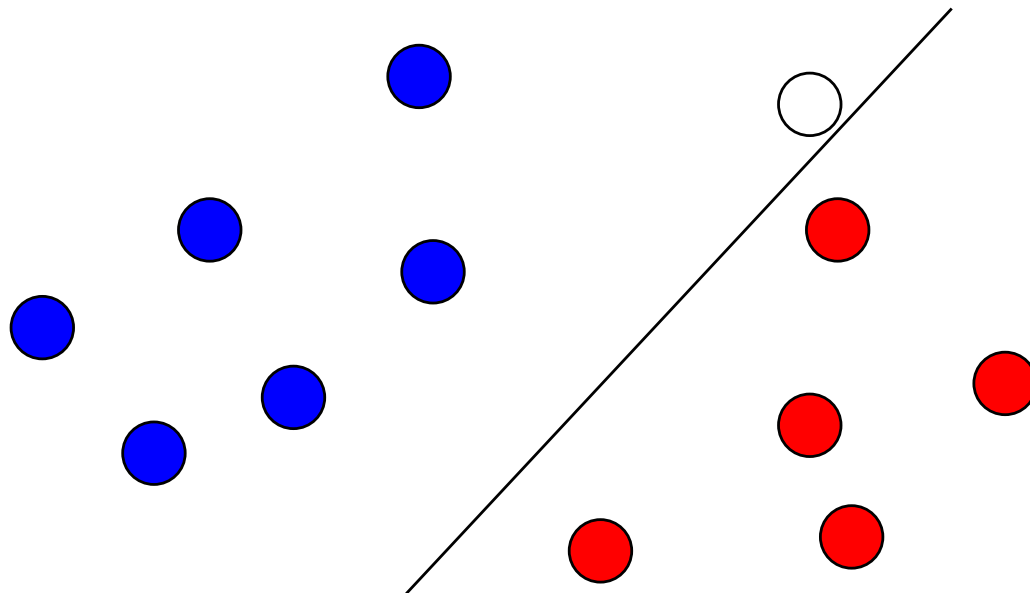
Linear classifier, some degrees of freedom



Linear classifier, some degrees of freedom

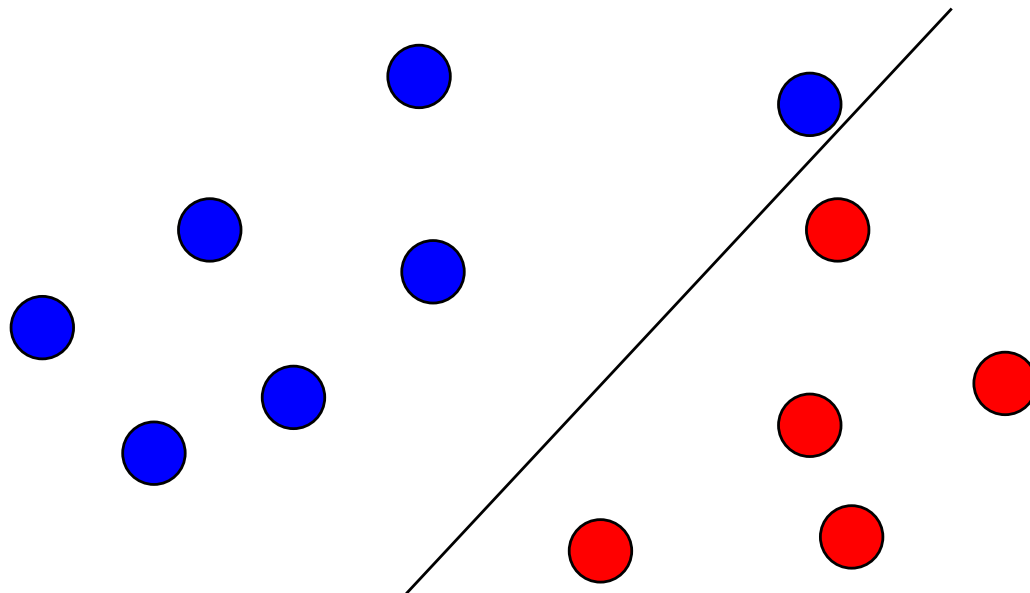


Linear classifier, some degrees of freedom

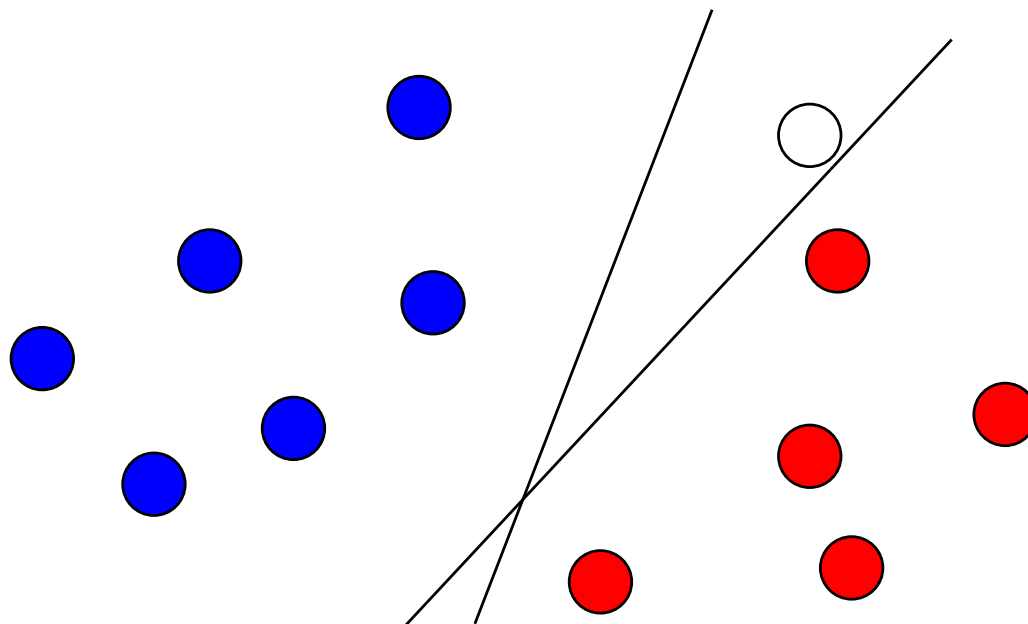


Specially close to the border of the classifier

Linear classifier, some degrees of freedom

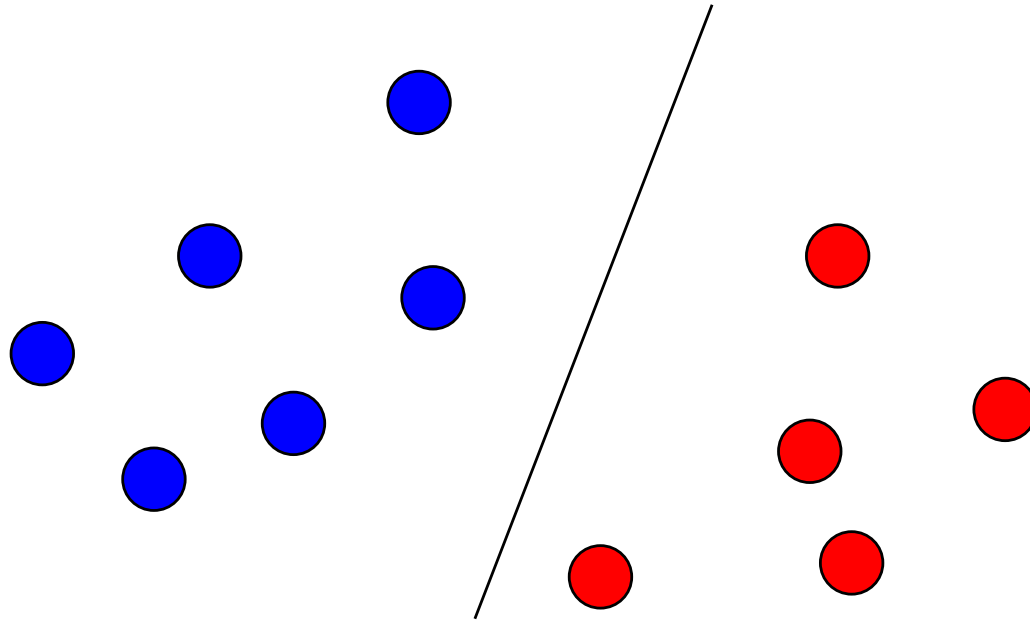


Linear classifier, some degrees of freedom

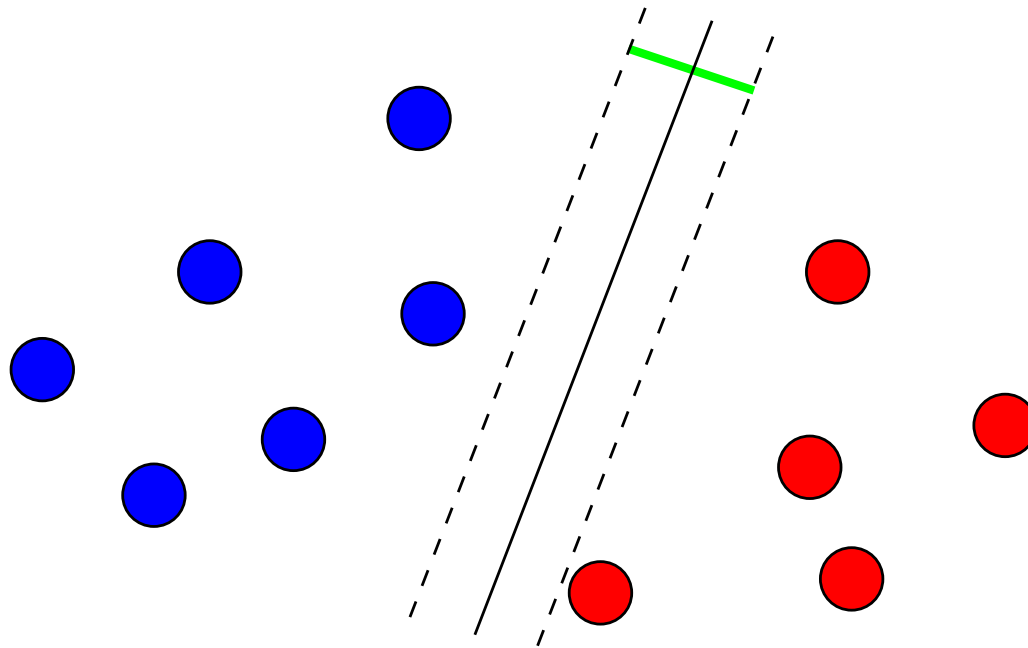


For each different technique, different results, different performance.

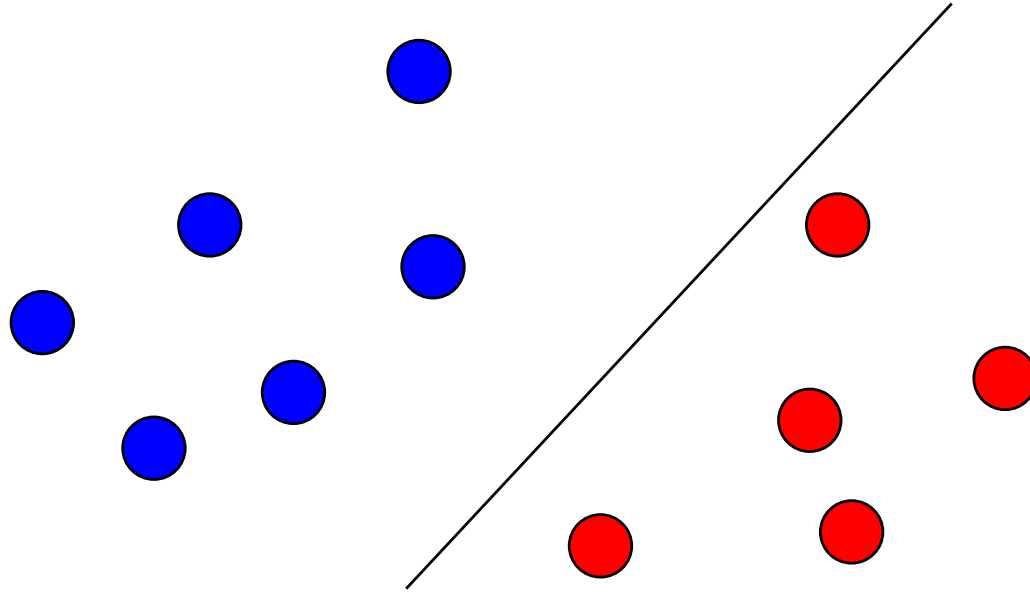
A criterion to select a linear classifier: the margin ?



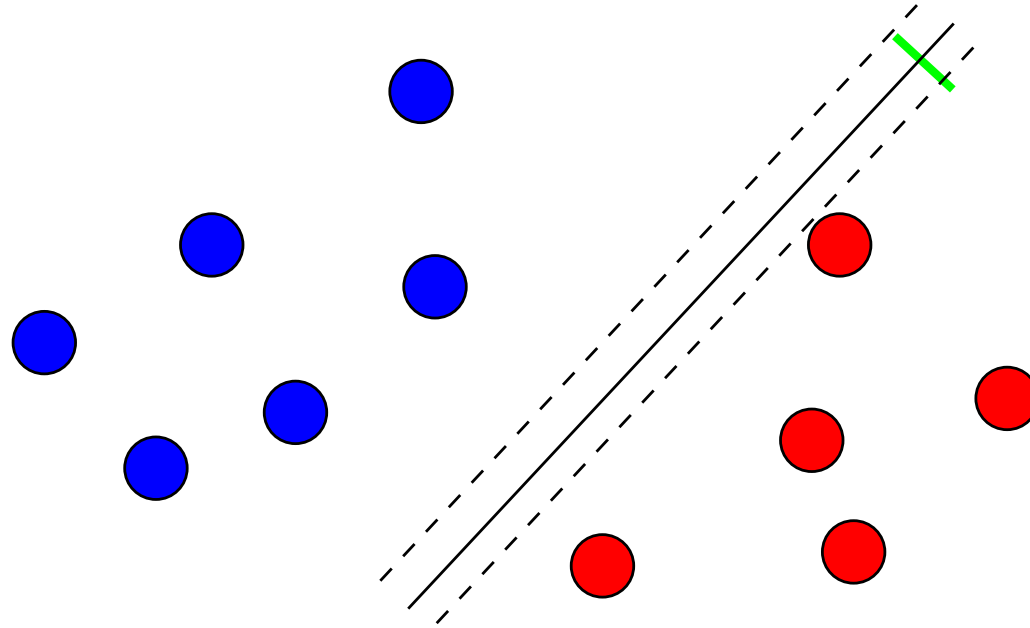
A criterion to select a linear classifier: the margin ?



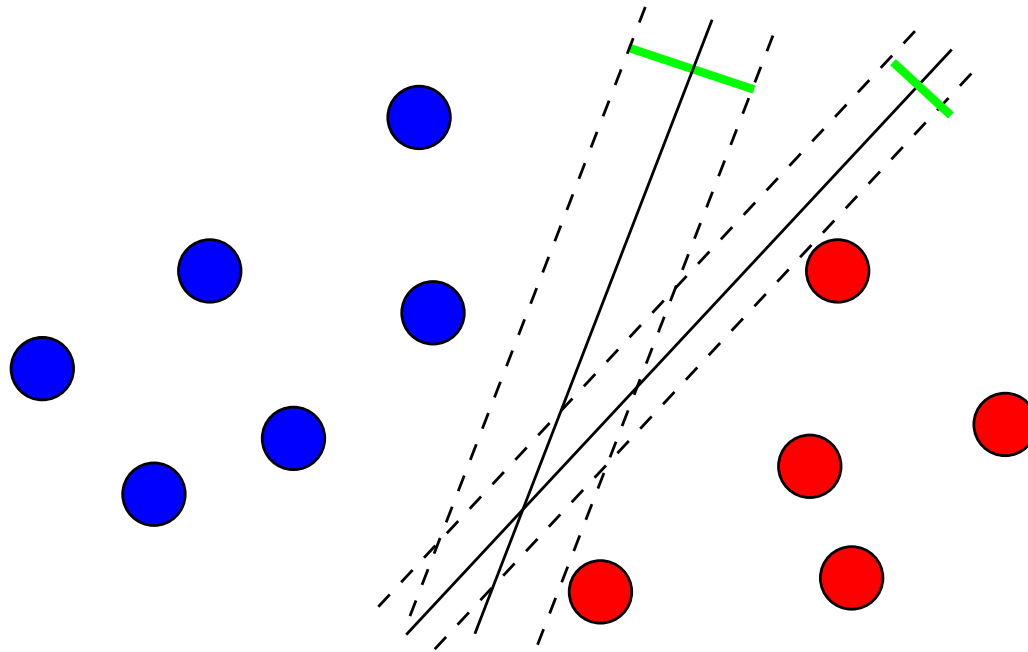
A criterion to select a linear classifier: the margin ?



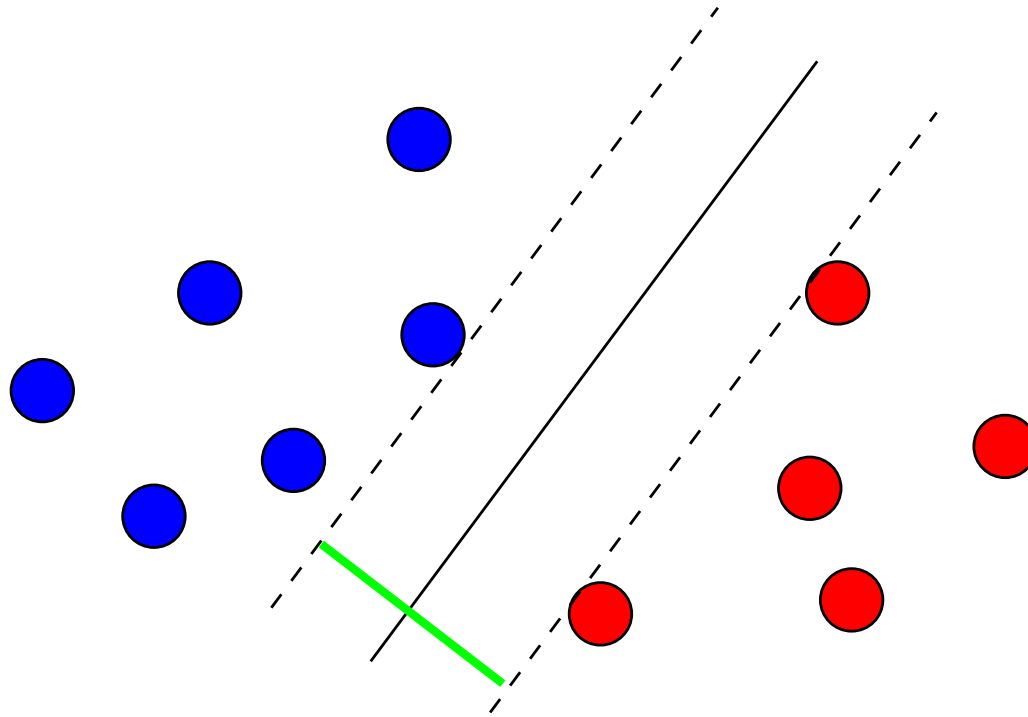
A criterion to select a linear classifier: the margin ?



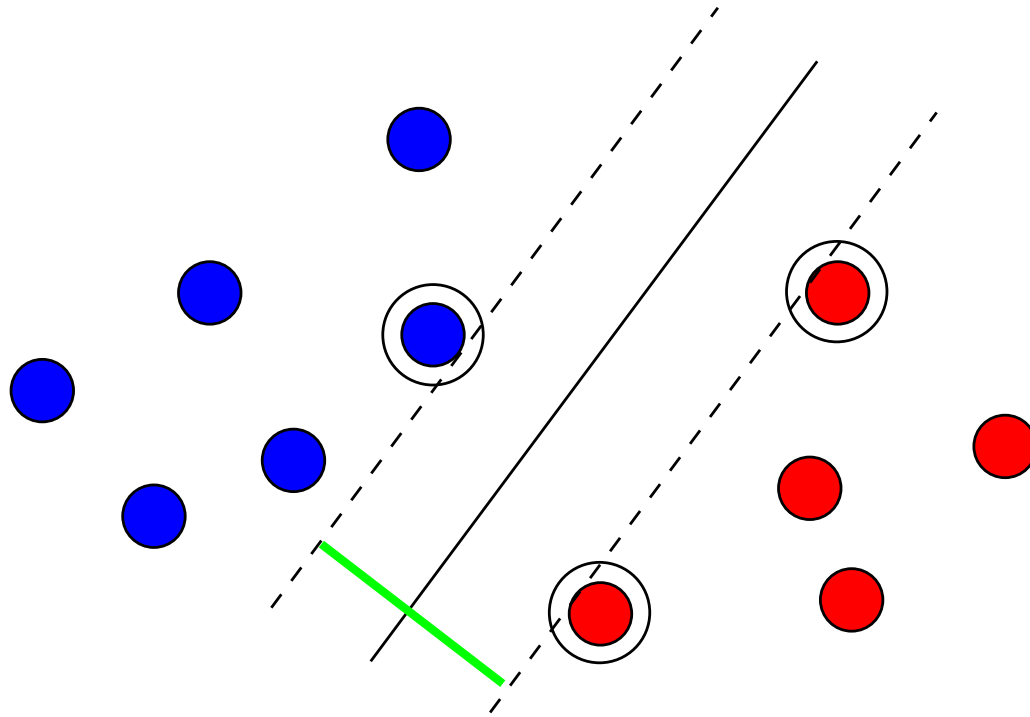
A criterion to select a linear classifier: the margin ?



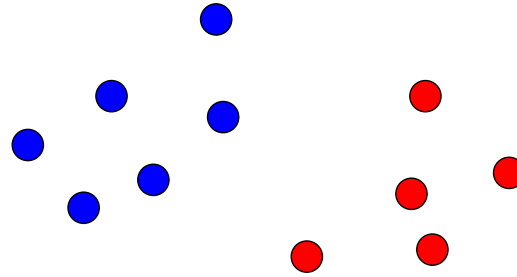
Largest Margin Linear Classifier ?



Support Vectors with Large Margin



In equations



- The **training set** is a finite set of n data/class pairs:

$$\mathcal{T} = \{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_N, \mathbf{y}_N)\},$$

where $\mathbf{x}_i \in \mathbb{R}^d$ and $\mathbf{y}_i \in \{-1, 1\}$.

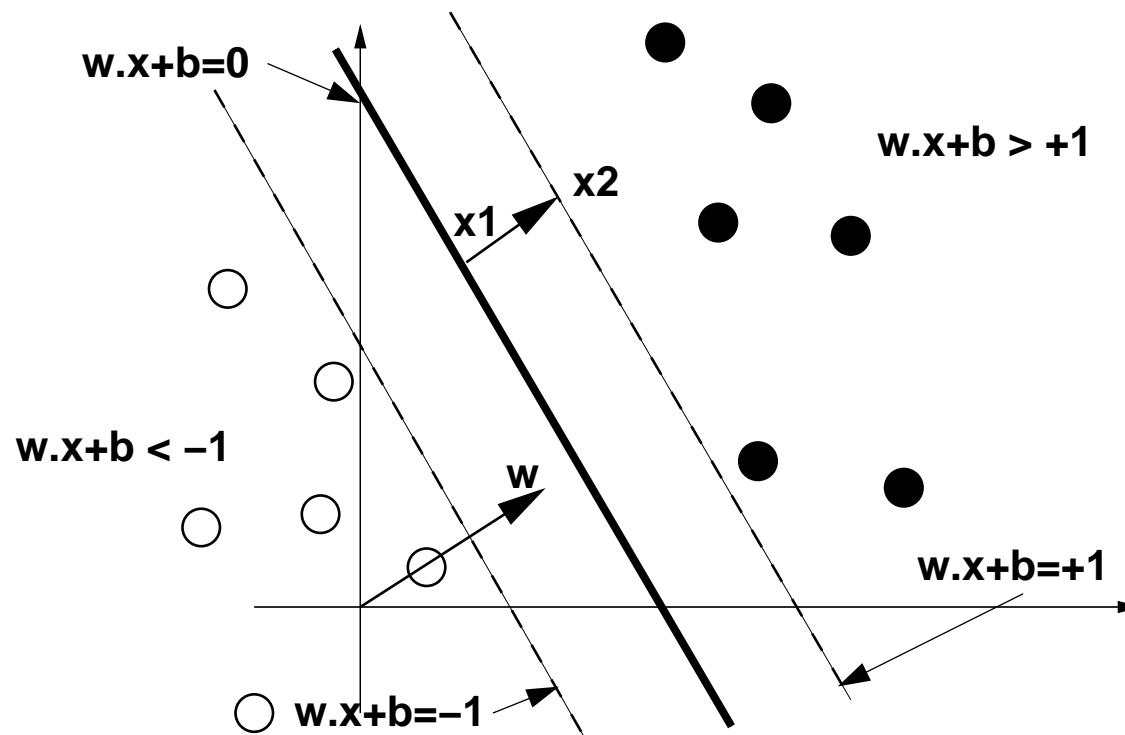
- We assume (for the moment) that the data are **linearly separable**, i.e., that there exists $(\mathbf{w}, b) \in \mathbb{R}^d \times \mathbb{R}$ such that:

$$\begin{cases} \mathbf{w}^T \mathbf{x}_i + b > 0 & \text{if } \mathbf{y}_i = 1, \\ \mathbf{w}^T \mathbf{x}_i + b < 0 & \text{if } \mathbf{y}_i = -1. \end{cases}$$

How to find the largest separating hyperplane?

For the linear classifier $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$ consider the *interstice* defined by the hyperplanes

- $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = +1$
- $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = -1$



The margin is $2/||\mathbf{w}||$

- Indeed, the points \mathbf{x}_1 and \mathbf{x}_2 satisfy:

$$\begin{cases} \mathbf{w}^T \mathbf{x}_1 + b = 0, \\ \mathbf{w}^T \mathbf{x}_2 + b = 1. \end{cases}$$

- By subtracting we get $\mathbf{w}^T (\mathbf{x}_2 - \mathbf{x}_1) = 1$, and therefore:

$$\gamma = 2||\mathbf{x}_2 - \mathbf{x}_1|| = \frac{2}{||\mathbf{w}||}.$$

where γ is the margin.

All training points should be on the appropriate side

- For positive examples ($y_i = 1$) this means:

$$\mathbf{w}^T \mathbf{x}_i + b \geq 1$$

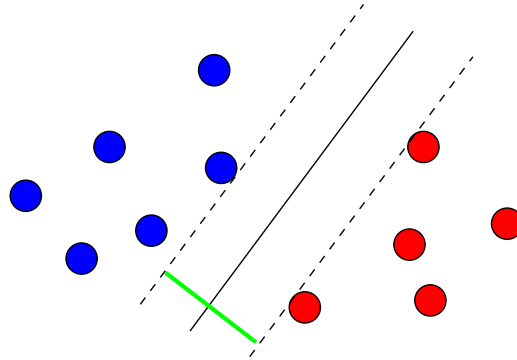
- For negative examples ($y_i = -1$) this means:

$$\mathbf{w}^T \mathbf{x}_i + b \leq -1$$

- in both cases:

$$\forall i = 1, \dots, n, \quad \mathbf{y}_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

Finding the optimal hyperplane



- Finding the optimal hyperplane is equivalent to finding (\mathbf{w}, b) which minimize:

$$\|\mathbf{w}\|^2$$

under the constraints:

$$\forall i = 1, \dots, n, \quad \mathbf{y}_i (\mathbf{w}^T \mathbf{x}_i + b) - 1 \geq 0.$$

This is a classical quadratic program on \mathbb{R}^{d+1}
linear constraints - **quadratic objective**

Lagrangian

- In order to minimize:

$$\frac{1}{2} \|\mathbf{w}\|^2$$

under the constraints:

$$\forall i = 1, \dots, n, \quad y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1 \geq 0.$$

- introduce **one dual variable** α_i for each constraint,
- one constraint for **each training point**.
- the **Lagrangian** is, for $\alpha \succeq 0$ (that is for each $\alpha_i \geq 0$)

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1).$$

The Lagrange dual function

$$g(\alpha) = \inf_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1) \right\}$$

the saddle point conditions give us that at the minimum in \mathbf{w} and b

$$\mathbf{w} = \sum_{i=1}^n \alpha_i \mathbf{y}_i \mathbf{x}_i, \quad (\text{derivating w.r.t } \mathbf{w}) \quad (*)$$

$$0 = \sum_{i=1}^n \alpha_i \mathbf{y}_i, \quad (\text{derivating w.r.t } b) \quad (**)$$

substituting (*) in g , and using (**) as a constraint, get the dual function $g(\alpha)$.

- To solve the dual problem, **maximize** g w.r.t. α .
- **Strong duality holds** : primal and dual problems have the **same optimum**.
- KKT gives us $\alpha_i (\mathbf{y}_i (\mathbf{w}^T \mathbf{x}_i + b) - 1) = 0$,
...hence, either **$\alpha_i = 0$** or **$\mathbf{y}_i (\mathbf{w}^T \mathbf{x}_i + b) = 1$** .
- $\alpha_i \neq 0$ **only** for points on the support hyperplanes $\{(\mathbf{x}, \mathbf{y}) \mid \mathbf{y}_i (\mathbf{w}^T \mathbf{x}_i + b) = 1\}$.

Dual optimum

The dual problem is thus

$$\begin{aligned} &\text{maximize} && g(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \\ &\text{such that} && \alpha \succeq 0, \sum_{i=1}^n \alpha_i \mathbf{y}_i = 0. \end{aligned}$$

This is a **quadratic program** in \mathbb{R}^n , with *box constraints*.
 α^* can be computed using optimization software
(*e.g.* built-in matlab function)

Recovering the optimal hyperplane

- With α^* , we recover (\mathbf{w}^T, b^*) corresponding to the **optimal hyperplane**.
- \mathbf{w}^T is given by $\mathbf{w}^T = \sum_{i=1}^n y_i \alpha_i \mathbf{x}_i^T$,
- b^* is given by the conditions on the support vectors $\alpha_i > 0$, $\mathbf{y}_i(\mathbf{w}^T \mathbf{x}_i + b) = 1$,

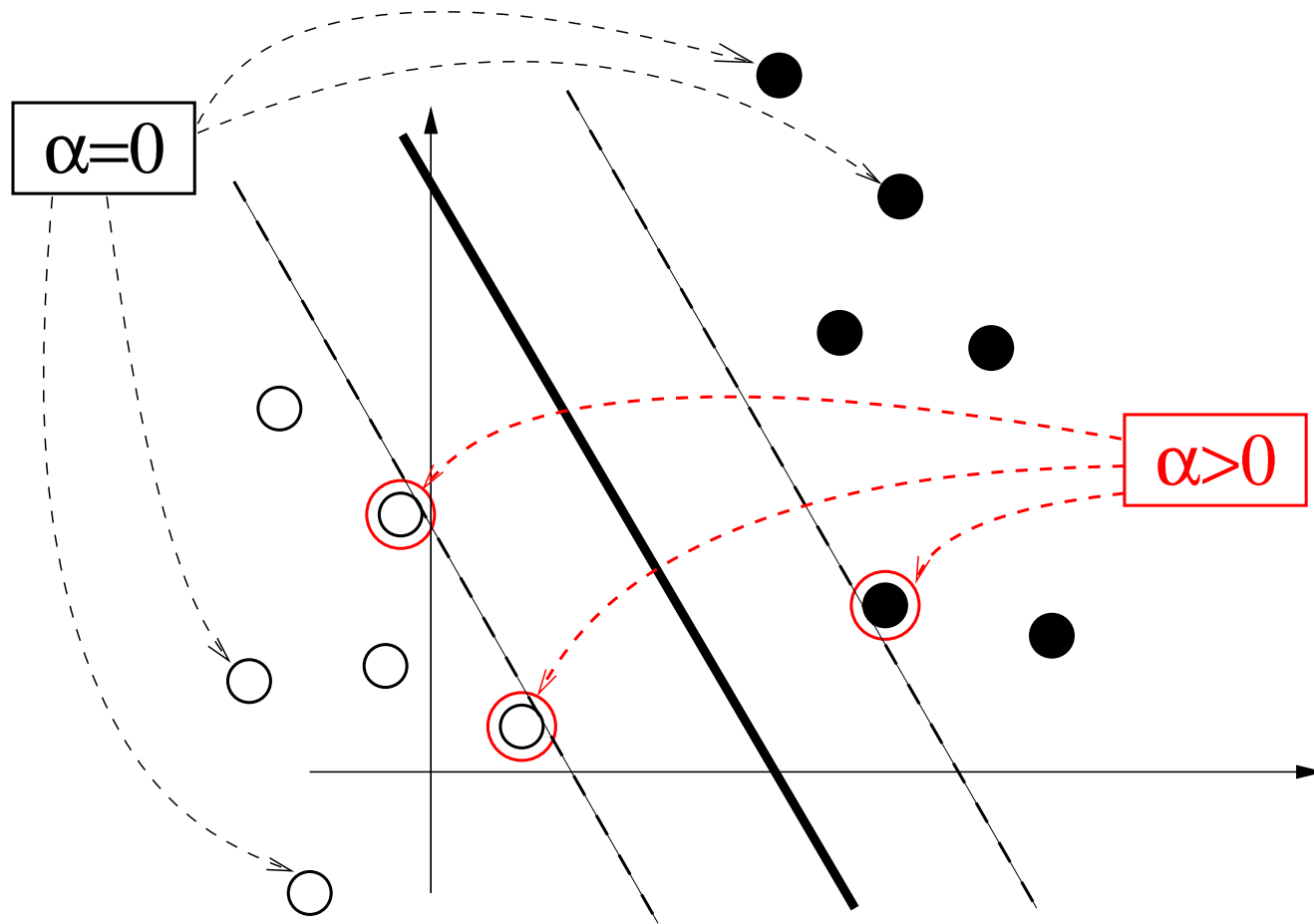
$$b^* = -\frac{1}{2} \left(\min_{\mathbf{y}_i=1, \alpha_i>0} (\mathbf{w}^T \mathbf{x}_i) + \max_{\mathbf{y}_i=-1, \alpha_i>0} (\mathbf{w}^T \mathbf{x}_i) \right)$$

- the **decision function** is therefore:

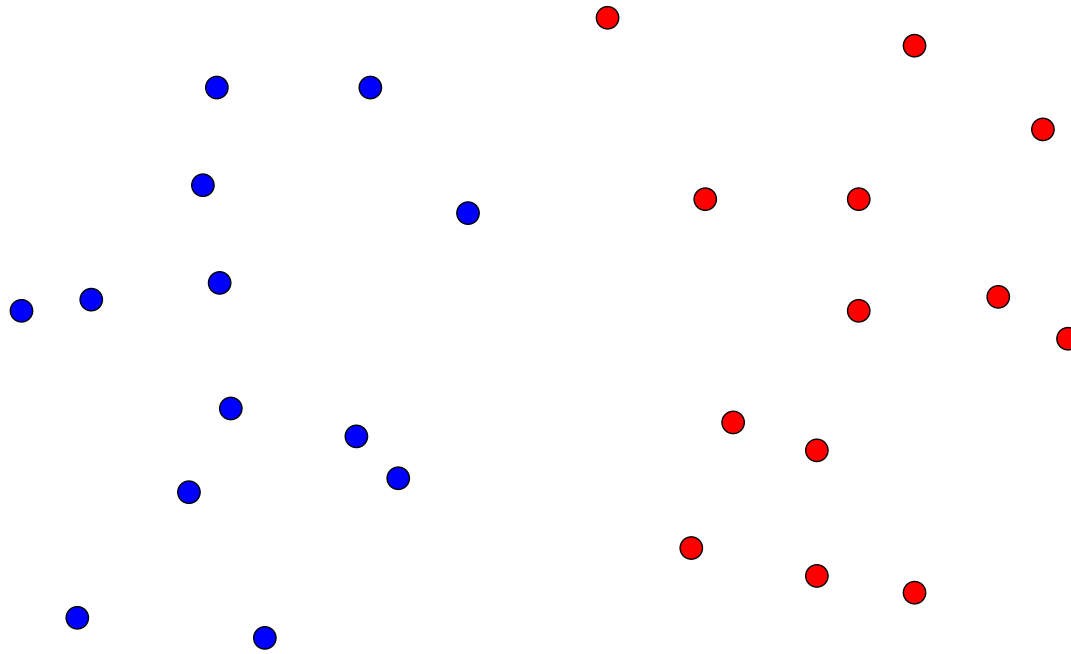
$$\begin{aligned} f^*(\mathbf{x}) &= \mathbf{w}^T \mathbf{x} + b^* \\ &= \left(\sum_{i=1}^n y_i \alpha_i \mathbf{x}_i^T \right) \mathbf{x} + b^*. \end{aligned}$$

- Here the **dual** solution gives us directly the **primal** solution.

Interpretation: support vectors

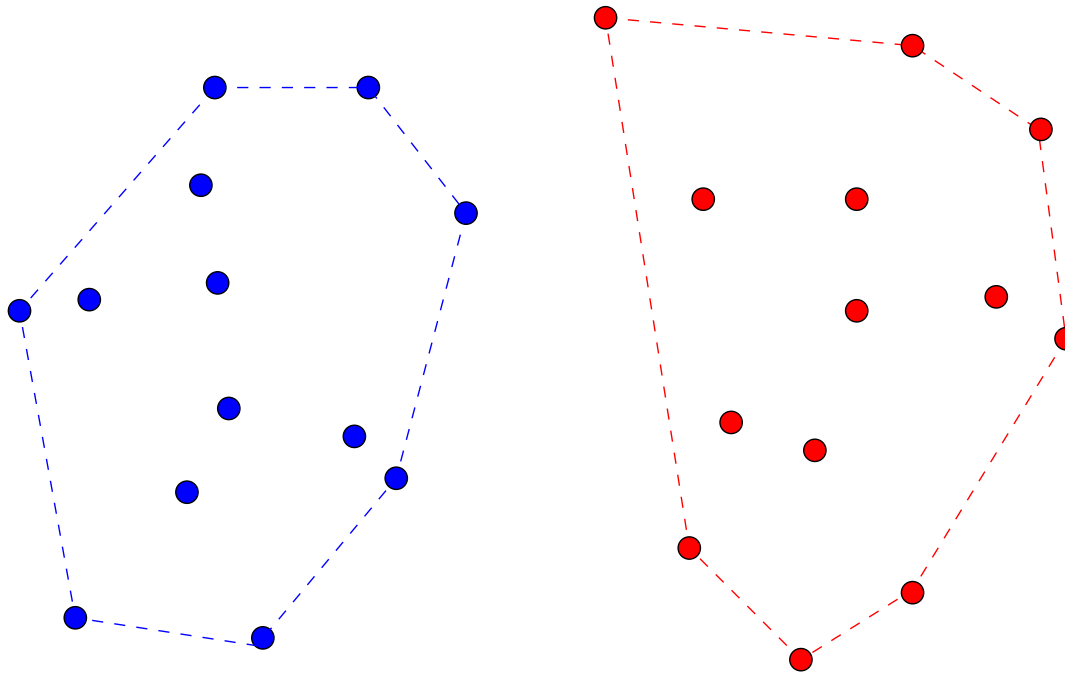


Another interpretation: Convex Hulls



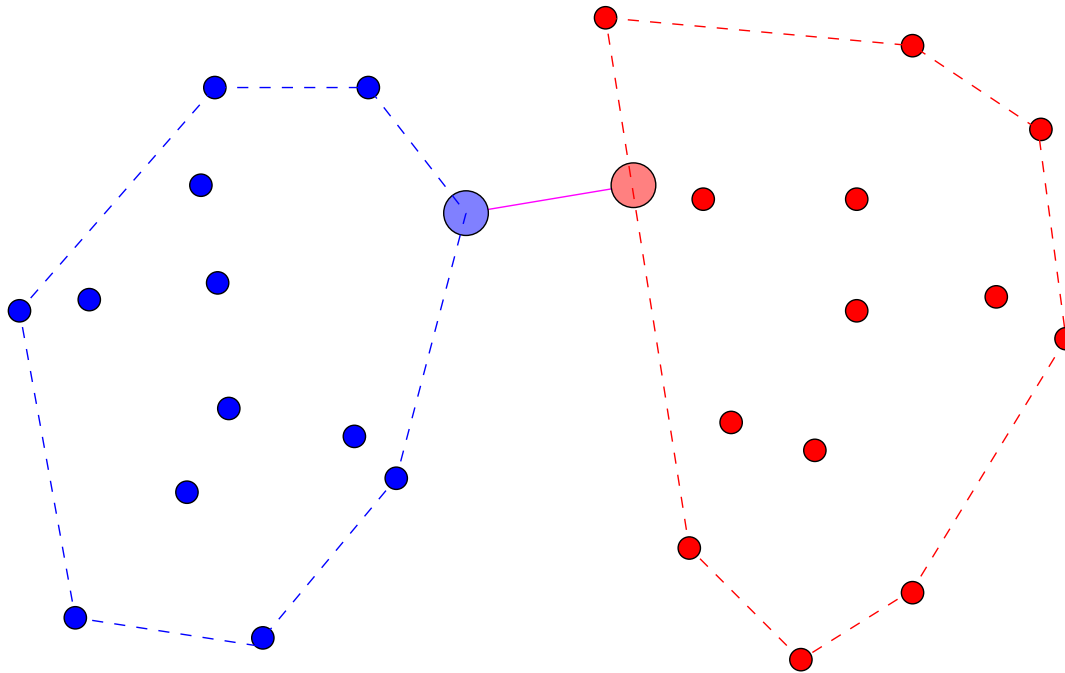
go back to 2 sets of points that are linearly separable

Another interpretation: Convex Hulls



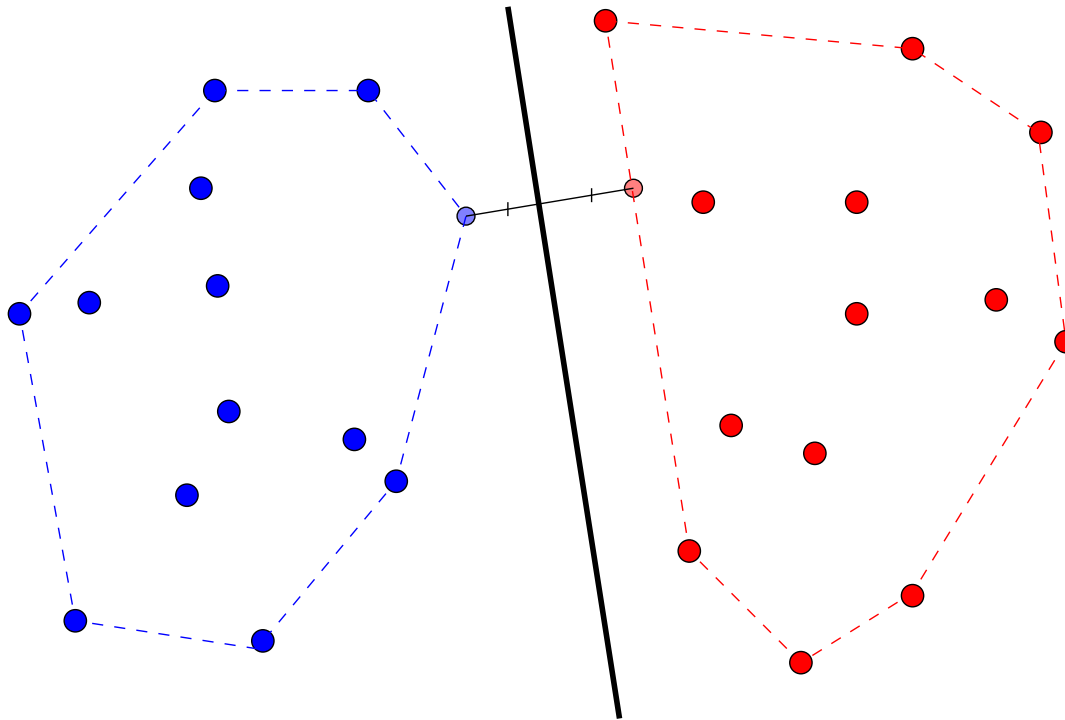
Linearly separable = convex hulls do not intersect

Another interpretation: Convex Hulls



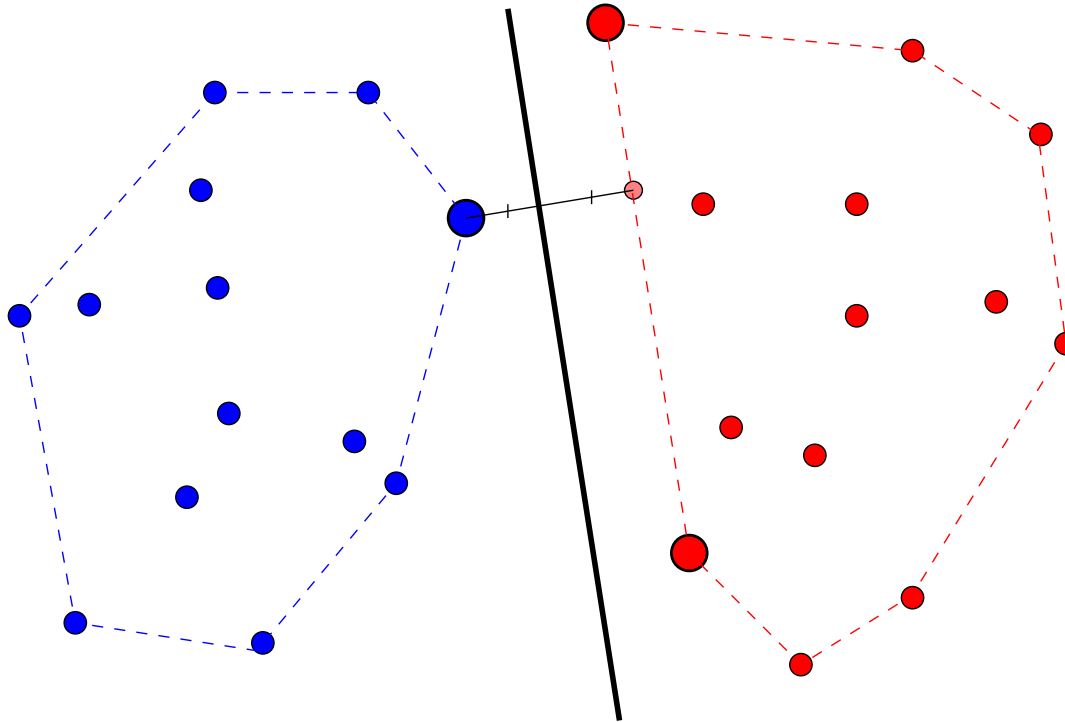
Find two closest points, one in each convex hull

Another interpretation: Convex Hulls



The SVM = bisection of that segment

Another interpretation: Convex Hulls

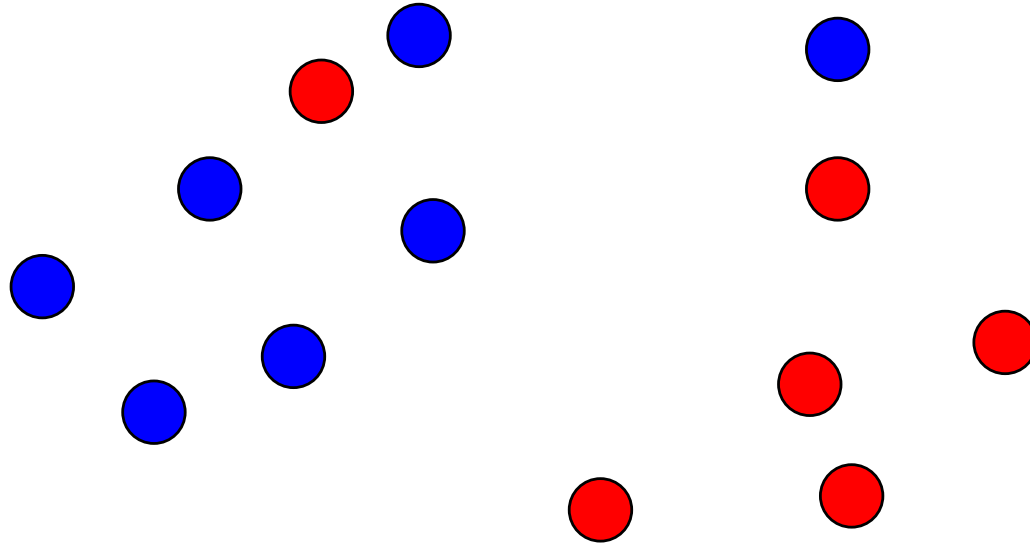


support vectors = extreme points of the faces on which the two points lie

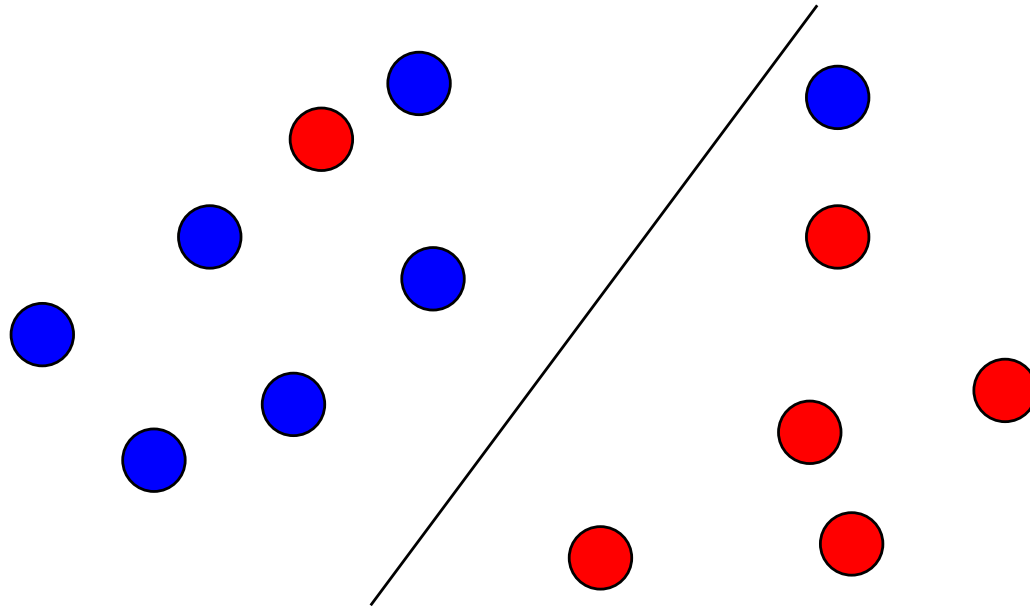
The non-linearly separable case

(when convex hulls intersect)

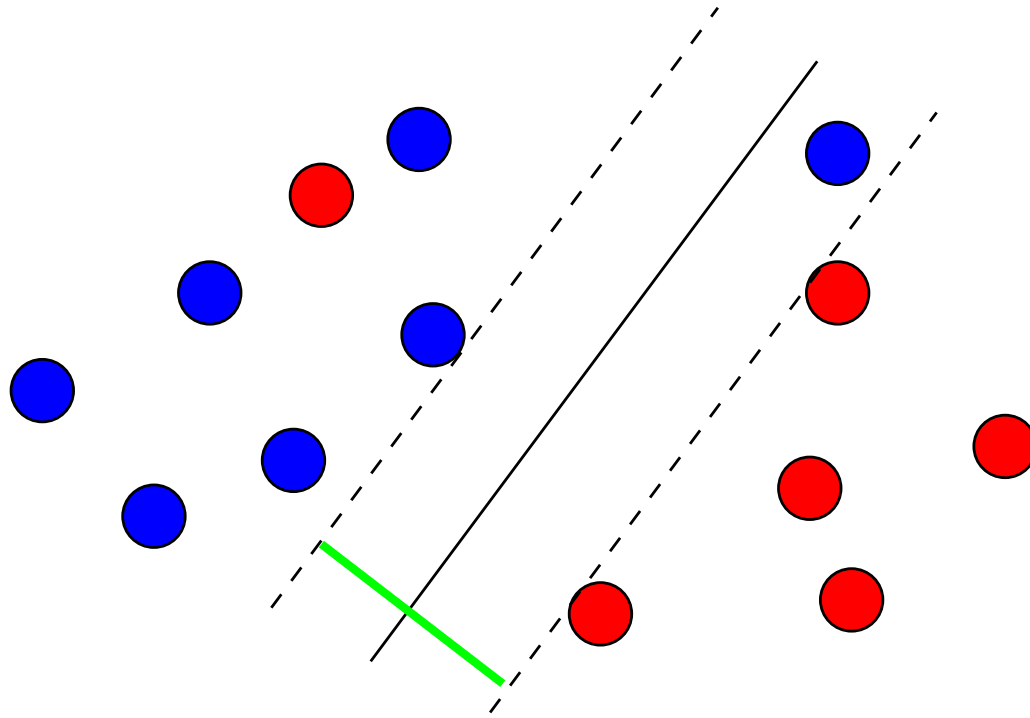
What happens when the data is not linearly separable?



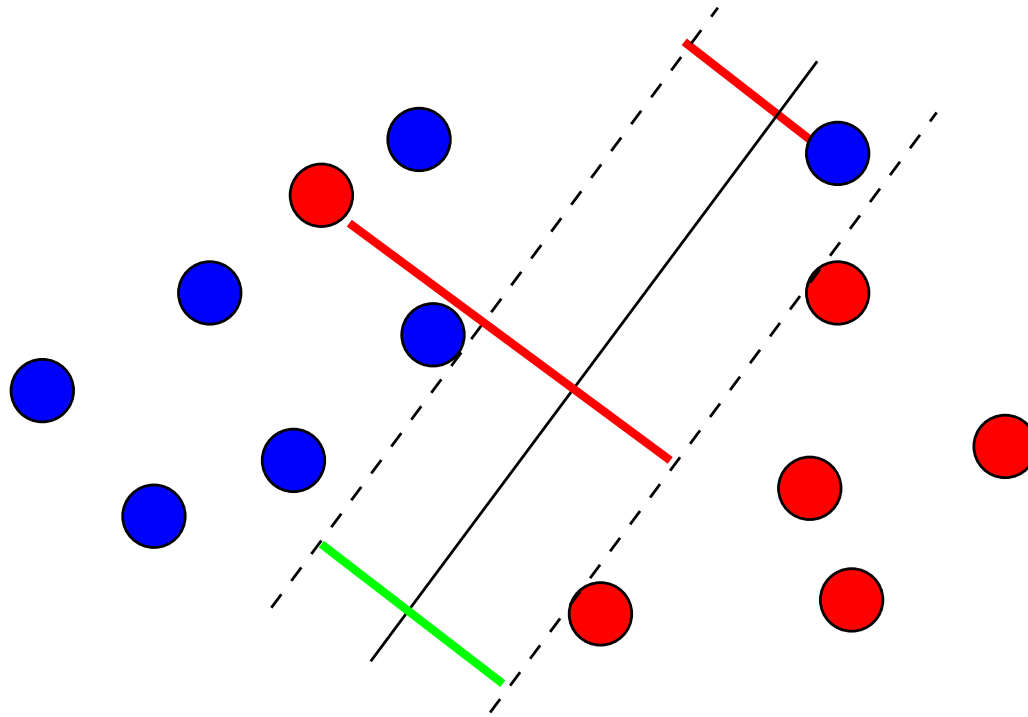
What happens when the data is not linearly separable?



What happens when the data is not linearly separable?



What happens when the data is not linearly separable?



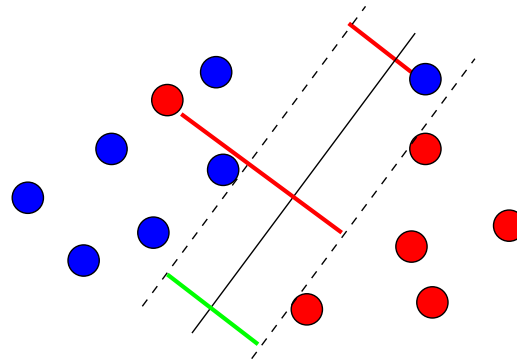
Soft-margin SVM ?

- Find a trade-off between **large margin** and **few errors**.

- Mathematically:

$$\min_f \left\{ \frac{1}{\text{margin}(f)} + C \times \text{errors}(f) \right\}$$

- C is a parameter



Soft-margin SVM formulation ?

- The **margin** of a labeled point (\mathbf{x}, \mathbf{y}) is

$$\text{margin}(\mathbf{x}, \mathbf{y}) = \mathbf{y} (\mathbf{w}^T \mathbf{x} + b)$$

- The **error** is
 - 0 if $\text{margin}(\mathbf{x}, \mathbf{y}) > 1$,
 - $1 - \text{margin}(\mathbf{x}, \mathbf{y})$ otherwise.

- The soft margin SVM solves:

$$\min_{\mathbf{w}, b} \{ \|\mathbf{w}\|^2 + C \sum_{i=1}^n \max\{0, 1 - \mathbf{y}_i (\mathbf{w}^T \mathbf{x}_i + b)\} \}$$

- $c(u, y) = \max\{0, 1 - yu\}$ is known as the **hinge loss**.
- $c(\mathbf{w}^T \mathbf{x}_i + b, \mathbf{y}_i)$ associates a mistake cost to the decision \mathbf{w}, b for example \mathbf{x}_i .

Dual formulation of soft-margin SVM

- The soft margin SVM program

$$\min_{\mathbf{w}, b} \left\{ \|\mathbf{w}\|^2 + C \sum_{i=1}^n \max\{0, 1 - \mathbf{y}_i (\mathbf{w}^T \mathbf{x}_i + b)\} \right\}$$

can be rewritten as

$$\begin{aligned} & \text{minimize} && \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \\ & \text{such that} && \mathbf{y}_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \end{aligned}$$

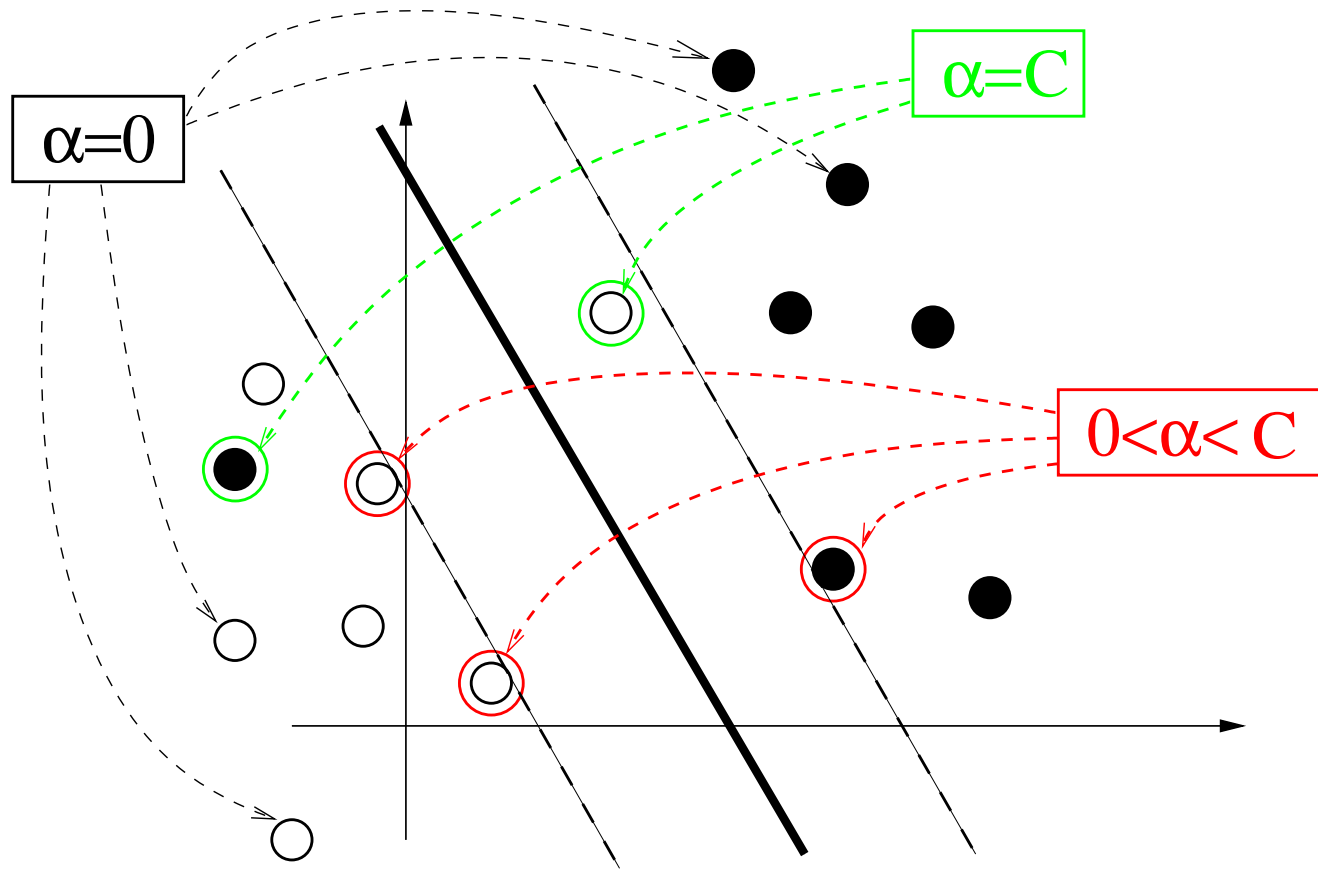
- In that case the dual function

$$g(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j \mathbf{y}_i \mathbf{y}_j \mathbf{x}_i^T \mathbf{x}_j,$$

which is finite under the constraints:

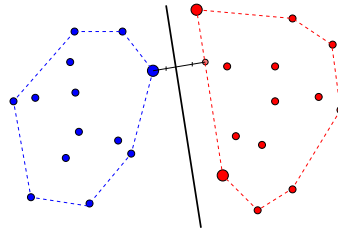
$$\begin{cases} 0 \leq \alpha_i \leq C, & \text{for } i = 1, \dots, n \\ \sum_{i=1}^n \alpha_i \mathbf{y}_i = 0. \end{cases}$$

Interpretation: bounded and unbounded support vectors

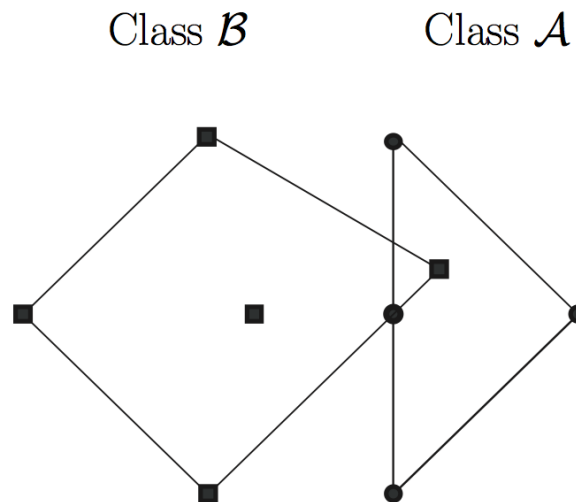


What about the convex hull analogy?

- Remember the separable case



- Here we consider the case where the two sets are not linearly separable, *i.e.* their convex hulls **intersect**.



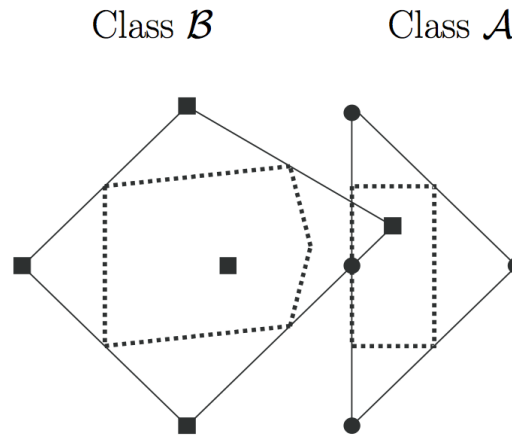
What about the convex hull analogy?

Definition 1. Given a set of n points \mathcal{A} , and $0 \leq C \leq 1$, the set of finite combinations

$$\sum_{i=1}^n \lambda_i \mathbf{x}_i, 1 \leq \lambda_i \leq C, \sum_{i=1}^n \lambda_i = 1,$$

is the (C) reduced convex hull of \mathcal{A}

- Using $C = 1/2$, the reduced convex hulls of \mathcal{A} and \mathcal{B} ,



- Soft-SVM with $C =$ closest two points of C -reduced convex hulls.

Kernels

Kernel trick for SVM's

- use a mapping ϕ from \mathcal{X} to a feature space,
- which corresponds to the **kernel** k :

$$\forall \mathbf{x}, \mathbf{x}' \in \mathcal{X}, \quad k(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$$

- Example: if $\phi(\mathbf{x}) = \phi \left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right) = \begin{bmatrix} x_1^2 \\ x_2^2 \end{bmatrix}$, then

$$k(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle = (x_1)^2(x_1')^2 + (x_2)^2(x_2')^2.$$

Training a SVM in the feature space

Replace each $\mathbf{x}^T \mathbf{x}'$ in the SVM algorithm by $\langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle = k(\mathbf{x}, \mathbf{x}')$

- **Reminder:** the dual problem is to maximize

$$g(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j),$$

under the constraints:

$$\begin{cases} 0 \leq \alpha_i \leq C, & \text{for } i = 1, \dots, n \\ \sum_{i=1}^n \alpha_i y_i = 0. \end{cases}$$

- The **decision function** becomes:

$$\begin{aligned} f(\mathbf{x}) &= \langle \mathbf{w}, \phi(x) \rangle + b^* \\ &= \sum_{i=1}^n y_i \alpha_i k(\mathbf{x}_i, \mathbf{x}) + b^*. \end{aligned} \tag{1}$$

The Kernel Trick ?

The explicit computation of $\phi(\mathbf{x})$ is not necessary.
The kernel $k(\mathbf{x}, \mathbf{x}')$ is enough.

- the SVM optimization for α works **implicitly** in the feature space.
- the SVM is a kernel algorithm: only need to input \mathbf{K} and \mathbf{y} :

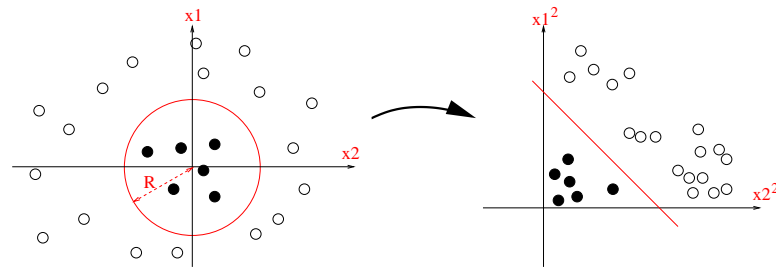
$$\begin{aligned} \text{maximize} \quad & g(\alpha) = \alpha^T \mathbf{1} - \frac{1}{2} \alpha^T (\mathbf{K} \odot \mathbf{y}\mathbf{y}^T) \alpha \\ \text{such that} \quad & 0 \leq \alpha_i \leq C, \quad \text{for } i = 1, \dots, n \\ & \sum_{i=1}^n \alpha_i \mathbf{y}_i = 0. \end{aligned}$$

- \mathbf{K} 's **positive definite** \Leftrightarrow **problem has an unique optimum**
- the decision function is $f(\cdot) = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \cdot) + b$.

Kernel example: polynomial kernel

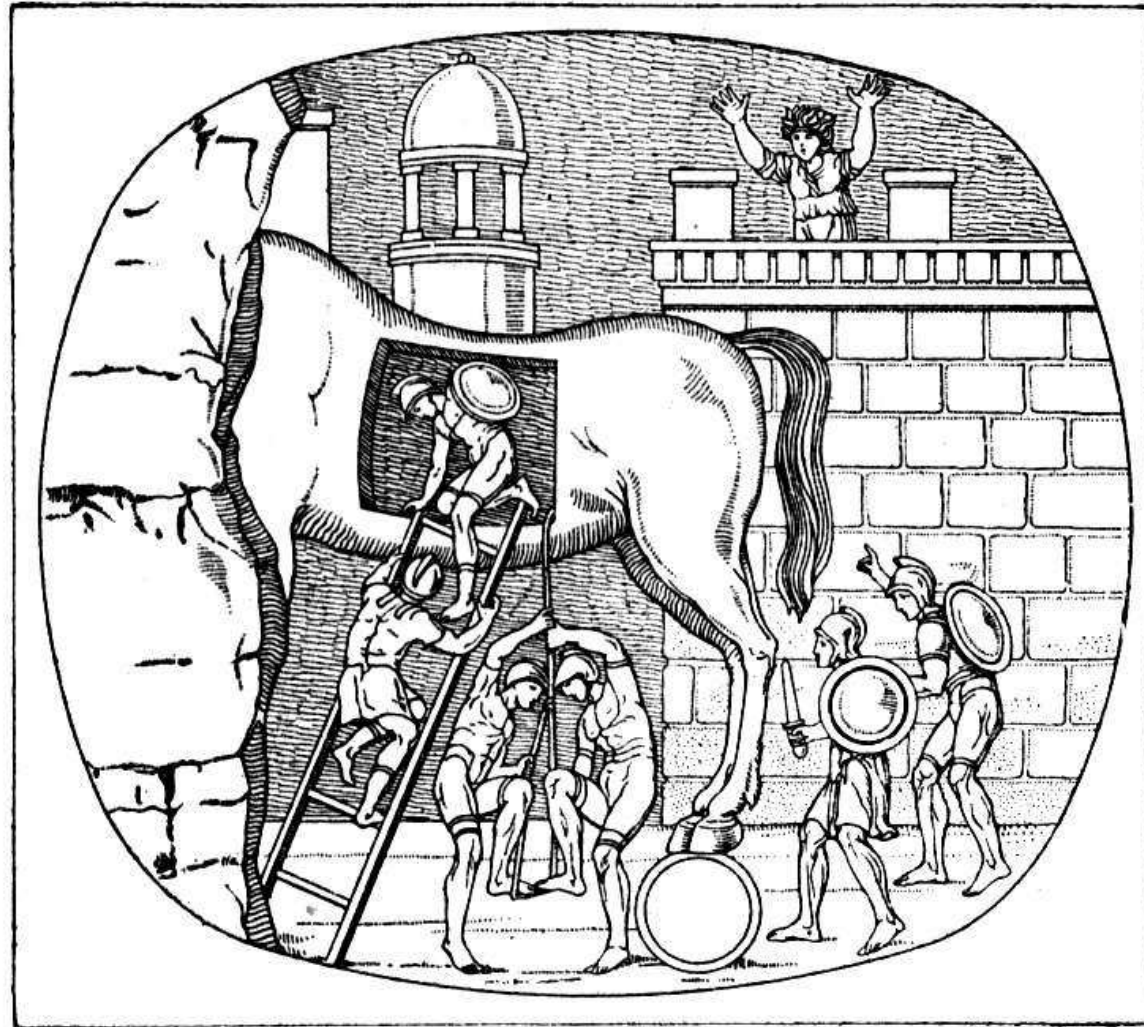
- For $\mathbf{x} = (x_1, x_2)^\top \in \mathbb{R}^2$, let $\phi(\mathbf{x}) = (x_1^2, \sqrt{2}x_1x_2, x_2^2) \in \mathbb{R}^3$:

$$\begin{aligned}K(\mathbf{x}, \mathbf{x}') &= x_1^2x_1'^2 + 2x_1x_2x_1'x_2' + x_2^2x_2'^2 \\ &= \{x_1x_1' + x_2x_2'\}^2 \\ &= \{\mathbf{x}^T \mathbf{x}'\}^2.\end{aligned}$$



Kernels are Trojan Horses onto Linear Models

- With kernels, complex structures can enter the realm of linear models



Designing Kernels

- As with distances, one can design kernels on any kind of object
 - time-series, strings, graphs, trees
 - images, video, audio, text
 - combination of the objects above!
- very large literature! too vast to discuss today.

Positive Definite Kernels & Combinatorial Distances for Structures

Structured Objects

- Objects in a countable set
 - variable length strings, trees, graphs, permutations
- Constrained vectors
 - Positive vectors, histograms
- Vectors of different sizes
 - variable length time series

Structured Objects

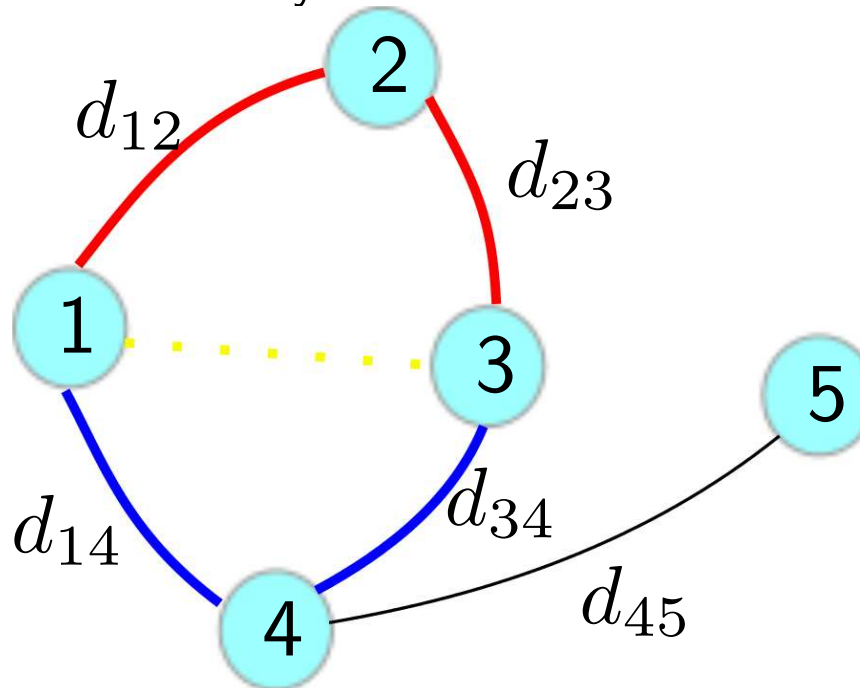
- Objects in a countable set
 - variable length strings, trees, graphs, permutations
- Constrained vectors
 - Positive vectors, histograms
- Vectors of different sizes
 - variable length time series

How can we define a **kernel** or a **distance** on such sets?

in most cases, applying standard distances on \mathbb{R}^n or even \mathbb{N}^n is meaningless

Back to fundamentals

- **Distances** are **optimal** by nature, and quantify **shortest length paths**.
 - Graph-metrics are defined that way

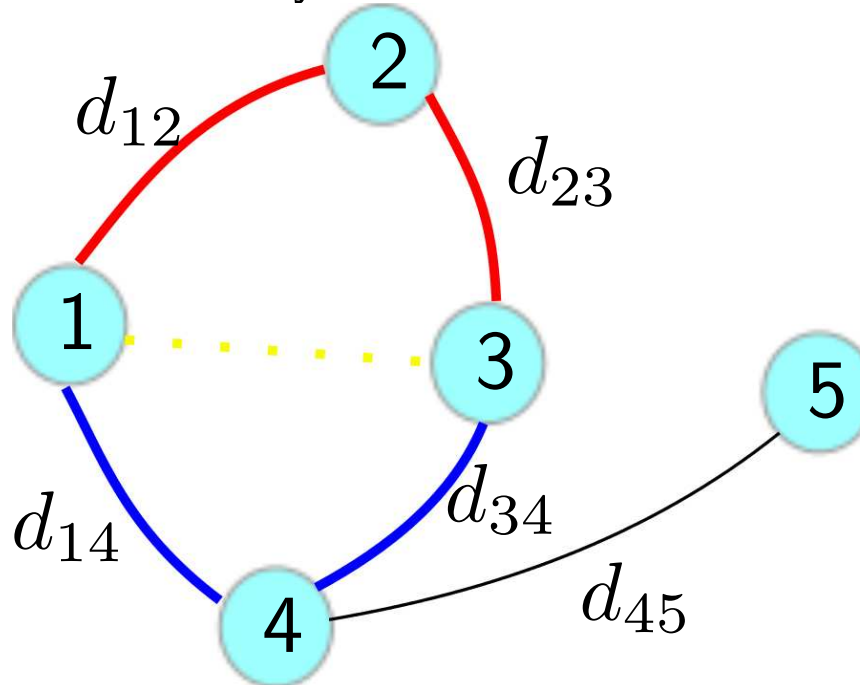


- **Triangle inequalities** are defined precisely to enforce this optimality

$$d(\mathbf{x}, \mathbf{y}) \leq d(\mathbf{x}, \mathbf{z}) + d(\mathbf{z}, \mathbf{y})$$

Back to fundamentals

- **Distances** are **optimal** by nature, and quantify **shortest length paths**.
 - Graph-metrics are defined that way



- **Triangle inequalities** are defined precisely to enforce this optimality

$$d(\mathbf{x}, \mathbf{y}) \leq d(\mathbf{x}, \mathbf{z}) + d(\mathbf{z}, \mathbf{y})$$

→ many **distances** on structured objects rely on **optimization**

Back to fundamentals

- **p.d. kernels** are additive by nature
 - k is positive definite $\Leftrightarrow \exists \varphi : \mathcal{X} \rightarrow \mathcal{H}$ such that

$$k(\mathbf{x}, \mathbf{y}) = \langle \varphi(x), \varphi(y) \rangle_{\mathcal{H}}.$$

- $X \in \mathcal{S}_n^+$ $\Leftrightarrow \exists L \in \mathbb{R}^{n \times n} | X = L^T L.$

Back to fundamentals

- **p.d. kernels** are additive by nature
 - k is positive definite $\Leftrightarrow \exists \varphi : \mathcal{X} \rightarrow \mathcal{H}$ such that

$$k(\mathbf{x}, \mathbf{y}) = \langle \varphi(x), \varphi(y) \rangle_{\mathcal{H}}.$$

- $X \in \mathcal{S}_n^+$ $\Leftrightarrow \exists L \in \mathbb{R}^{n \times n} | X = L^T L$.

→ many **kernels** on structured objects
rely on defining **explicitly** (possibly infinite) feature vectors

very large literature on this subject which we will not address here.

Combinatorial Distances

- To define a **distance**, an approach which has been repeatedly used is to,
 - Consider two inputs \mathbf{x}, \mathbf{y} ,
 - Define a **countable** set of **mappings** from \mathbf{x} to \mathbf{y} , $T(\mathbf{x}, \mathbf{y})$
 - Define a **cost** $c(\tau)$ for each element τ of $T(\mathbf{x}, \mathbf{y})$.
 - Define a distance between \mathbf{x}, \mathbf{y} as

$$d(\mathbf{x}, \mathbf{y}) = \min_{\tau \in T(\mathbf{x}, \mathbf{y})} c(\tau)$$

Combinatorial Distances

- To define a **distance**, an approach which has been repeatedly used is to,
 - Consider two inputs \mathbf{x}, \mathbf{y} ,
 - Define a **countable** set of **mappings** from \mathbf{x} to \mathbf{y} , $T(\mathbf{x}, \mathbf{y})$
 - Define a **cost** $c(\tau)$ for each element τ of $T(\mathbf{x}, \mathbf{y})$.
 - Define a distance between \mathbf{x}, \mathbf{y} as

$$d(\mathbf{x}, \mathbf{y}) = \min_{\tau \in T(\mathbf{x}, \mathbf{y})} c(\tau)$$

- **Symmetry, definiteness and triangle inequalities** depend on c and T .

Combinatorial Distances

- To define a **distance**, an approach which has been repeatedly used is to,
 - Consider two inputs \mathbf{x}, \mathbf{y} ,
 - Define a **countable** set of **mappings** from \mathbf{x} to \mathbf{y} , $T(\mathbf{x}, \mathbf{y})$
 - Define a **cost** $c(\tau)$ for each element τ of $T(\mathbf{x}, \mathbf{y})$.
 - Define a distance between \mathbf{x}, \mathbf{y} as

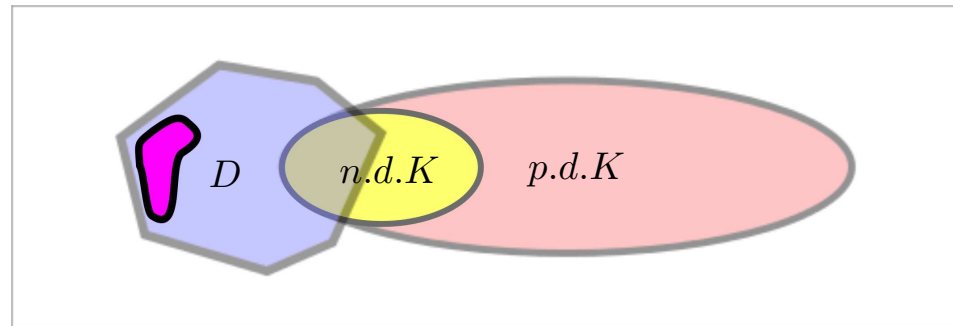
$$d(\mathbf{x}, \mathbf{y}) = \min_{\tau \in T(\mathbf{x}, \mathbf{y})} c(\tau)$$

- **Symmetry, definiteness and triangle inequalities** depend on c and T .
- In many cases, T is endowed with a dot product, $c(\tau) = \langle \tau, \theta \rangle$ for some θ .

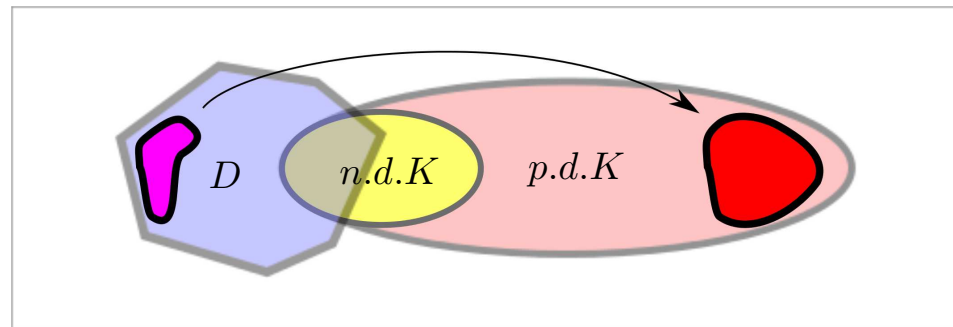
Combinatorial Distances are not Negative Definite

$$d(\mathbf{x}, \mathbf{y}) = \min_{\tau \in T(\mathbf{x}, \mathbf{y})} c(\tau)$$

- In most cases such distances are **not** negative definite



- Can we use them to define kernels?

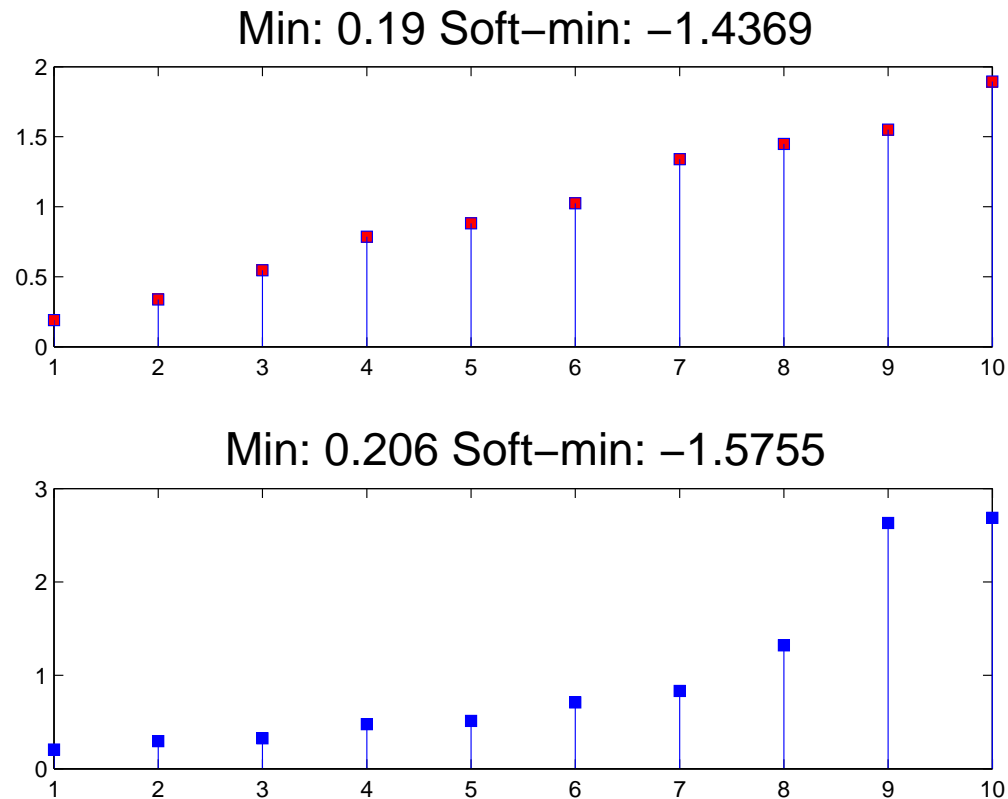


- **Yes** so far, using always the **same technique**.

An alternative definition of minimality

for a family of numbers $a_n, n \in \mathbb{N}$,

$$\text{soft-min} a_n = -\log \sum_n e^{-a_n}$$



Soft-min of costs - Generating Functions

$$d(\mathbf{x}, \mathbf{y}) = \min_{\tau \in T(\mathbf{x}, \mathbf{y})} c(\tau)$$

e^{-d} is **not** positive definite in the general case

Soft-min of costs - Generating Functions

$$d(\mathbf{x}, \mathbf{y}) = \min_{\tau \in T(\mathbf{x}, \mathbf{y})} c(\tau)$$

e^{-d} is **not** positive definite in the general case

$$\delta(\mathbf{x}, \mathbf{y}) = \text{soft-min}_{\tau \in T(\mathbf{x}, \mathbf{y})} c(\tau)$$

$e^{-\delta}$ has been proved to be **positive definite** in all known cases

Soft-min of costs - Generating Functions

$$d(\mathbf{x}, \mathbf{y}) = \min_{\tau \in T(\mathbf{x}, \mathbf{y})} c(\tau)$$

e^{-d} is **not** positive definite in the general case

$$\delta(\mathbf{x}, \mathbf{y}) = \text{soft-min}_{\tau \in T(\mathbf{x}, \mathbf{y})} c(\tau)$$

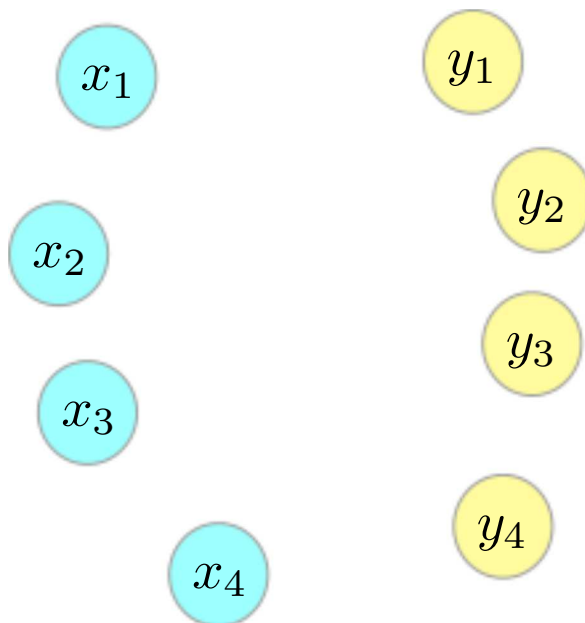
$e^{-\delta}$ has been proved to be **positive definite** in all known cases

$$e^{-\delta(\mathbf{x}, \mathbf{y})} = \sum_{\tau \in T(\mathbf{x}, \mathbf{y})} e^{-\langle \tau, \theta \rangle} = G_{T(\mathbf{x}, \mathbf{y})}(\theta)$$

$G_{T(\mathbf{x}, \mathbf{y})}$ is the **generating function** of the set of all mappings between \mathbf{x} and \mathbf{y} .

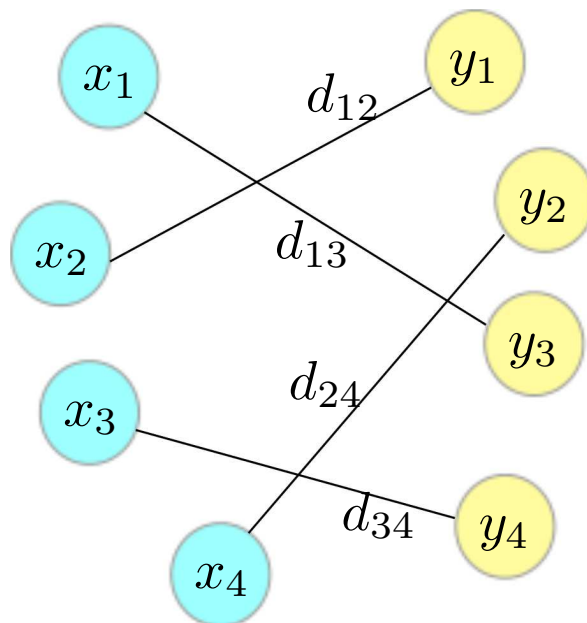
Example: Optimal assignment distance between two sets

- **Input:** $\mathbf{x} = \{x_1, \dots, x_n\}, \mathbf{y} = \{y_1, \dots, y_n\} \in \mathcal{X}^n$



Example: Optimal assignment distance between two sets

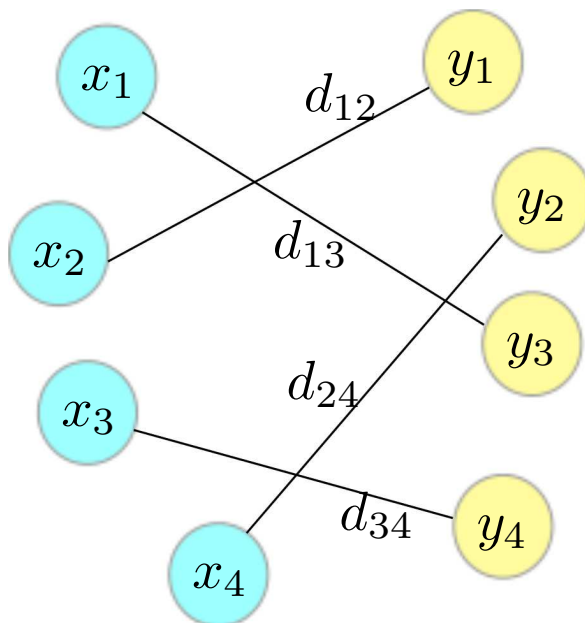
- **Input:** $\mathbf{x} = \{x_1, \dots, x_n\}, \mathbf{y} = \{y_1, \dots, y_n\} \in \mathcal{X}^n$



- **cost parameter:** distance d on \mathcal{X} . **mapping variable:** permutation σ in S_n
- **cost:** $\sum_{i=1}^n d(x_i, y_{\sigma(i)})$.

Example: Optimal assignment distance between two sets

- **Input:** $\mathbf{x} = \{x_1, \dots, x_n\}, \mathbf{y} = \{y_1, \dots, y_n\} \in \mathcal{X}^n$



- **cost parameter:** distance d on \mathcal{X} . **mapping variable:** permutation σ in S_n .
- **cost:** $\sum_{i=1}^n d(x_i, y_{\sigma(i)}) = \langle P_\sigma, D \rangle$ where $D = [d(x_i, y_j)]$

$$d_{\text{Assig.}}(\mathbf{x}, \mathbf{y}) = \min_{\sigma \in S_n} \sum_{i=1}^n d(x_i, y_{\sigma(i)}) = \min_{\sigma \in S_n} \langle P_\sigma, D \rangle$$

Example: Optimal assignment distance between two sets

$$d_{\text{Assig.}}(\mathbf{x}, \mathbf{y}) = \min_{\sigma \in S_n} \sum_{i=1}^n d(x_i, y_{\sigma(i)}) = \min_{\sigma \in S_n} \langle P_{\sigma}, D \rangle$$

define $k = e^{-d}$. If k is positive definite on \mathcal{X} then

$$k_{\text{Perm}}(\mathbf{x}, \mathbf{y}) = \sum_{\sigma \in S_n} e^{-\langle P_{\sigma}, D \rangle} = \text{Permanent}[k(x_i, y_j)]$$

is positive definite (C. 2007). $e^{-d_{\text{Assig.}}}$ is not (Frohlich et al. 2005, Vert 2008).

Example: Optimal alignment between two strings

- **Input:** $x = (x_1, \dots, x_n), y = (y_1, \dots, y_m) \in \mathcal{X}^n, \mathcal{X}$ finite

x = DOING, **y** =DONE

Example: Optimal alignment between two strings

- **Input:** $x = (x_1, \dots, x_n), y = (y_1, \dots, y_m) \in \mathcal{X}^n, \mathcal{X}$ finite

$\mathbf{x} = \text{DOING}, \mathbf{y} = \text{DONE}$

- **mapping variable:** alignment $\pi = \begin{pmatrix} \pi_1(1) & \cdots & \pi_1(q) \\ \pi_2(1) & \cdots & \pi_2(q) \end{pmatrix}$. (increasing path)

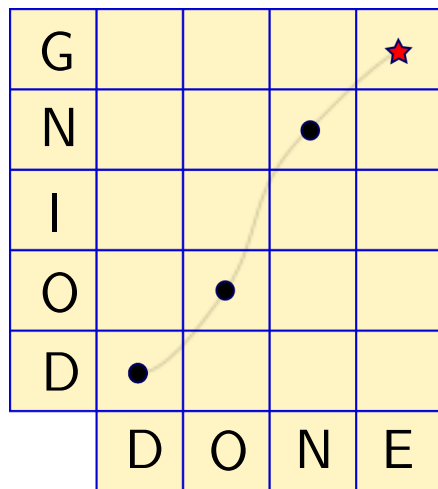


Example: Optimal alignment between two strings

- **Input:** $x = (x_1, \dots, x_n), y = (y_1, \dots, y_m) \in \mathcal{X}^n, \mathcal{X}$ finite

$\mathbf{x} = \text{DOING}, \mathbf{y} = \text{DONE}$

- **mapping variable:** alignment $\pi = \begin{pmatrix} \pi_1(1) & \cdots & \pi_1(q) \\ \pi_2(1) & \cdots & \pi_2(q) \end{pmatrix}$. (increasing path)



- **cost parameter:** distance d on \mathcal{X} + gap function $g : \mathbb{N} \rightarrow \mathbb{R}$.
- $c(\pi) = \sum_{i=1}^{|\pi|} d(x_{\pi_1(i)}, y_{\pi_2(i)}) + \sum_{i=1}^{|\pi|-1} g(\pi_1(i+1) - \pi_1(i)) + g(\pi_2(i+1) - \pi_2(i))$

Example: Optimal alignment between two strings

- **Input:** $x = (x_1, \dots, x_n), y = (y_1, \dots, y_m) \in \mathcal{X}^n, \mathcal{X}$ finite

$\mathbf{x} = \text{DOING}, \mathbf{y} = \text{DONE}$

- **mapping variable:** alignment $\pi = \begin{pmatrix} \pi_1(1) & \cdots & \pi_1(q) \\ \pi_2(1) & \cdots & \pi_2(q) \end{pmatrix}$. (increasing path)



- **cost parameter:** distance d on \mathcal{X} + gap function $g : \mathbb{N} \rightarrow \mathbb{R}$.

- $c(\pi) = \sum_{i=1}^{|\pi|} d(x_{\pi_1(i)}, y_{\pi_2(i)}) + \sum_{i=1}^{|\pi|-1} g(\pi_1(i+1) - \pi_1(i)) + g(\pi_2(i+1) - \pi_2(i))$

$$d_{\text{align}}(\mathbf{x}, \mathbf{y}) = \min_{\pi \in \text{Alignments}} c(\pi)$$

Example: Optimal alignment between two strings

$$d_{\text{align}}(\mathbf{x}, \mathbf{y}) = \min_{\pi \in \text{Alignments}} c(\pi)$$

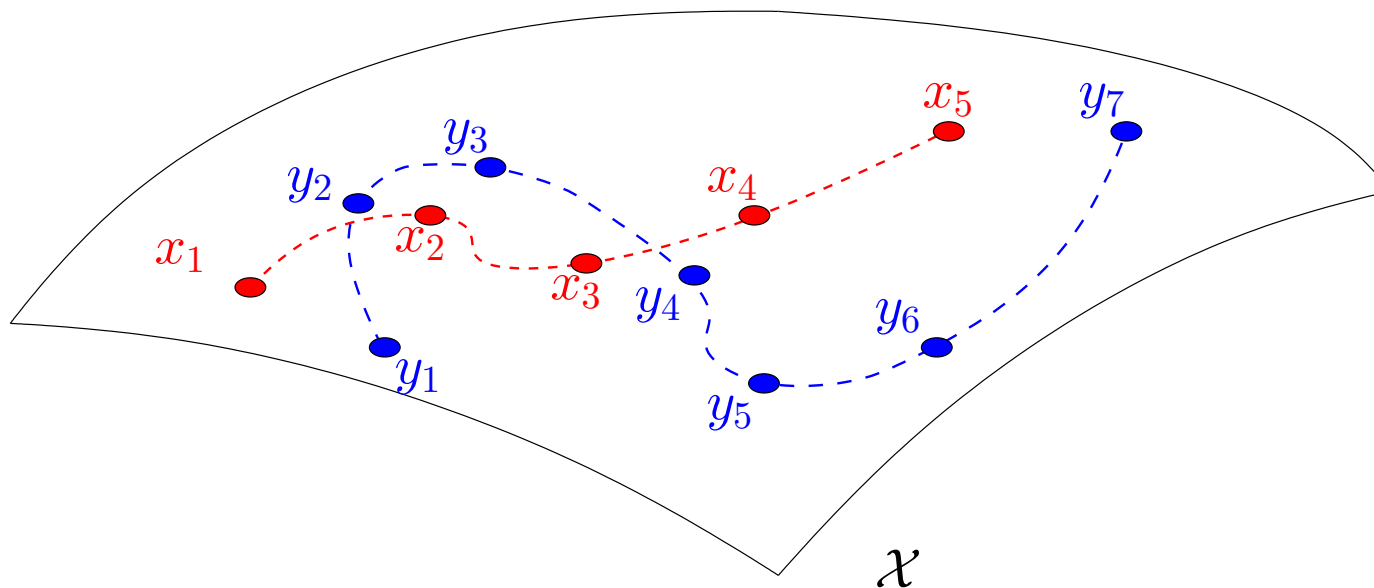
define $k = e^{-d}$. If k is positive definite on \mathcal{X} then

$$k_{\text{LA}}(\mathbf{x}, \mathbf{y}) = \sum_{\pi \in \text{Alignments}} e^{-c(\pi)}$$

is positive definite (Saigo et al. 2003).

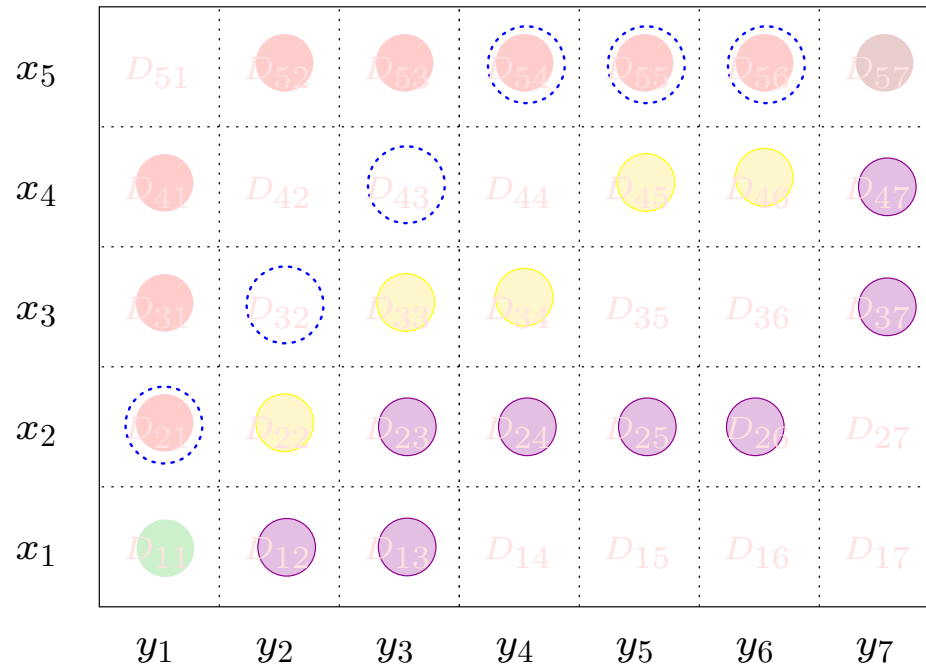
Example: Optimal time warping between two time series

- **Input:** $x = (x_1, \dots, x_n), y = (y_1, \dots, y_m) \in \mathbb{R}^n$



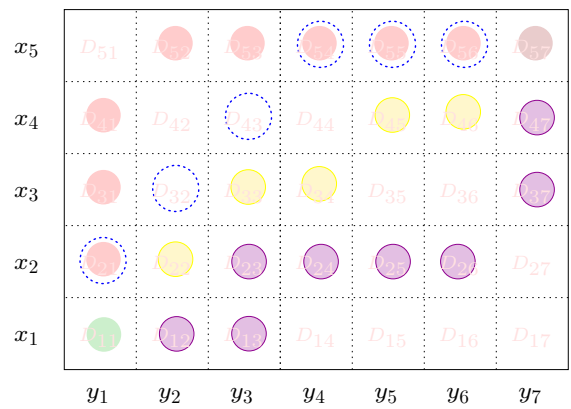
Example: Optimal time warping between two time series

- **Input:** $x = (x_1, \dots, x_n), y = (y_1, \dots, y_m) \in \mathbb{R}^n$
- **mapping variable:** $\pi = \begin{pmatrix} \pi_1(1) & \dots & \pi_1(q) \\ \pi_2(1) & \dots & \pi_2(q) \end{pmatrix}$. (increasing **contiguous** path)



Example: Optimal time warping between two time series

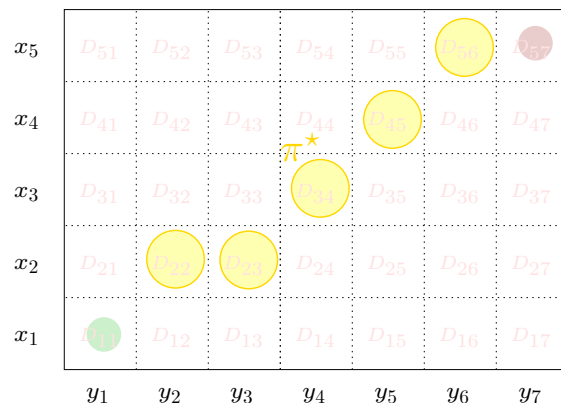
- **Input:** $x = (x_1, \dots, x_n), y = (y_1, \dots, y_m) \in \mathbb{R}^n$
- **mapping variable:** $\pi = \begin{pmatrix} \pi_1(1) & \dots & \pi_1(q) \\ \pi_2(1) & \dots & \pi_2(q) \end{pmatrix}$. (increasing **contiguous** path)



- **cost parameter:** distance d on \mathcal{X} . **cost:** $c(\pi) = \sum_{i=1}^{|\pi|} d(x_{\pi_1(i)}, y_{\pi_2(i)})$

Example: Optimal time warping between two time series

- **Input:** $x = (x_1, \dots, x_n), y = (y_1, \dots, y_m) \in \mathbb{R}^n$
- **mapping variable:** $\pi = \begin{pmatrix} \pi_1(1) & \dots & \pi_1(q) \\ \pi_2(1) & \dots & \pi_2(q) \end{pmatrix}$. (increasing **contiguous** path)



- **cost parameter:** distance d on \mathcal{X} . **cost:** $c(\pi) = \sum_{i=1}^{|\pi|} d(x_{\pi_1(i)}, y_{\pi_2(i)})$

$$d_{\text{DTW}}(\mathbf{x}, \mathbf{y}) = \min_{\pi \in \text{Alignments}} c(\pi)$$

Example: Optimal alignment between two strings

$$d_{\text{DTW}}(\mathbf{x}, \mathbf{y}) = \min_{\pi \in \text{Alignments}} c(\pi)$$

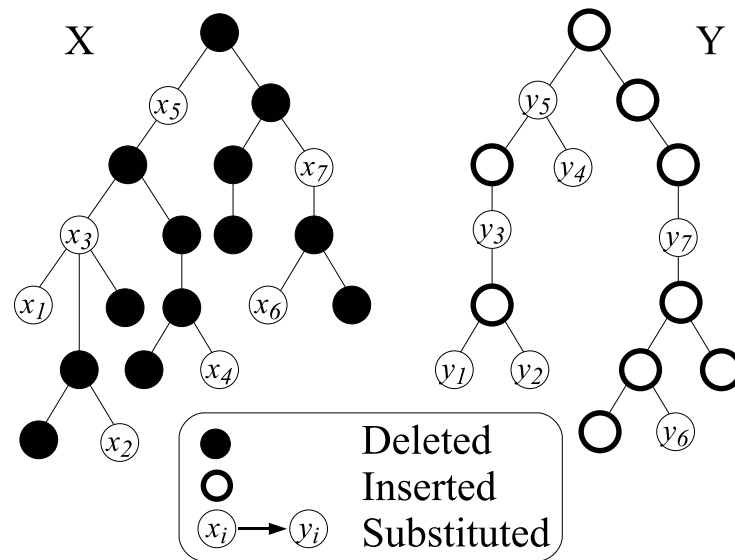
define $k = e^{-d}$. If k is positive definite and geometrically divisible on \mathcal{X} then

$$k_{\text{GA}}(\mathbf{x}, \mathbf{y}) = \sum_{\pi \in \text{Alignments}} e^{-c(\pi)}$$

is p.d. (C. et al. 2007, C. 2011). $e^{-d_{\text{DTW}}}$ is not (Shimodaira et al. 2001)

Example: Edit-distance between two trees

- **Input:** two labeled trees \mathbf{x}, \mathbf{y} .
- **mapping variable:** sequence of substitutions/deletions/insertions of vertices

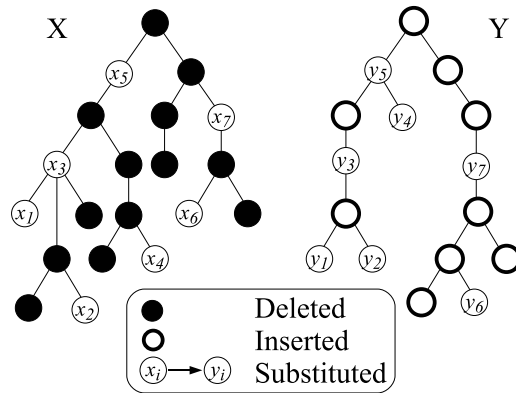


- **cost parameter:** γ distance between labels and cost for deletion/insertion

$$d_{\text{TreeEdit}}(\mathbf{x}, \mathbf{y}) = \min_{\sigma \in \text{EditScripts}(\mathbf{x}, \mathbf{y})} \sum \gamma(\sigma_i)$$

Example: Edit-distance between two trees

- **Input:** two labeled trees \mathbf{x}, \mathbf{y} .
- **mapping variable:** sequence of substitutions/deletions/insertions of vertices



- **cost parameter:** γ distance between labels and cost for deletion/insertion

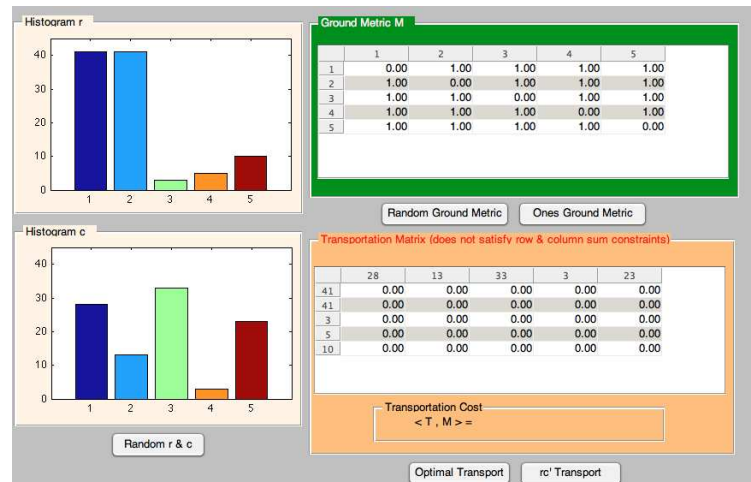
$$d_{\text{TreeEdit}}(\mathbf{x}, \mathbf{y}) = \min_{\sigma \in \text{EditScripts}(\mathbf{x}, \mathbf{y})} \sum \gamma(\sigma_i)$$

- if $e^{-\gamma}$ p.d. then p.definiteness of the generating function was proved by Shin & Kuboyama 2008; Shin, C., Kuboyama 2011.

Example: Transportation distance between discrete histograms

Monge-Kantorovich, Wasserstein, Earth Mover's, Mallow's etc...

- **Input:** two integer histograms $\mathbf{x}, \mathbf{y} \in \mathbb{N}^d$ such that $\sum_{i=1}^d x_i = \sum_{i=1}^d y_i = N$



- **mapping:** transportation matrices $U(r, c) = \{X \in \mathbb{N}^{d \times d} | X \mathbf{1}_d = \mathbf{x}, X^T \mathbf{1}_d = \mathbf{y}\}$
- **cost parameter:** M distance matrix in \mathcal{M}_d .

$$d_W(\mathbf{x}, \mathbf{y}) = \min_{X \in U(r, c)} \langle X, M \rangle$$

Example: Transportation distance between discrete histograms

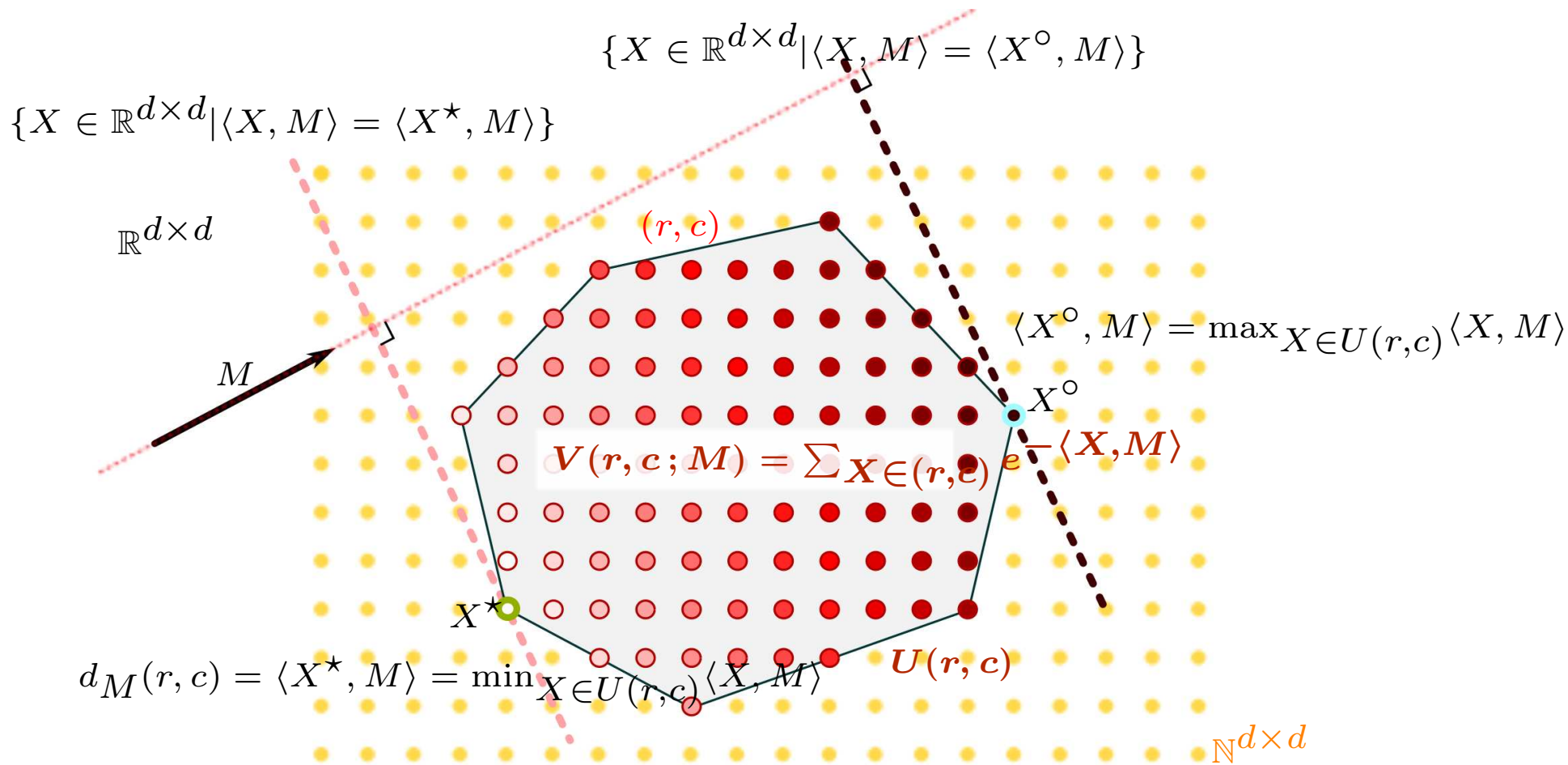
$$d_W(\mathbf{x}, \mathbf{y}) = \min_{X \in U(r,c)} \langle X, M \rangle$$

define $k_{ij} = e^{-m_{ij}}$. If $[k_{ij}]$ is positive definite on \mathcal{X} then

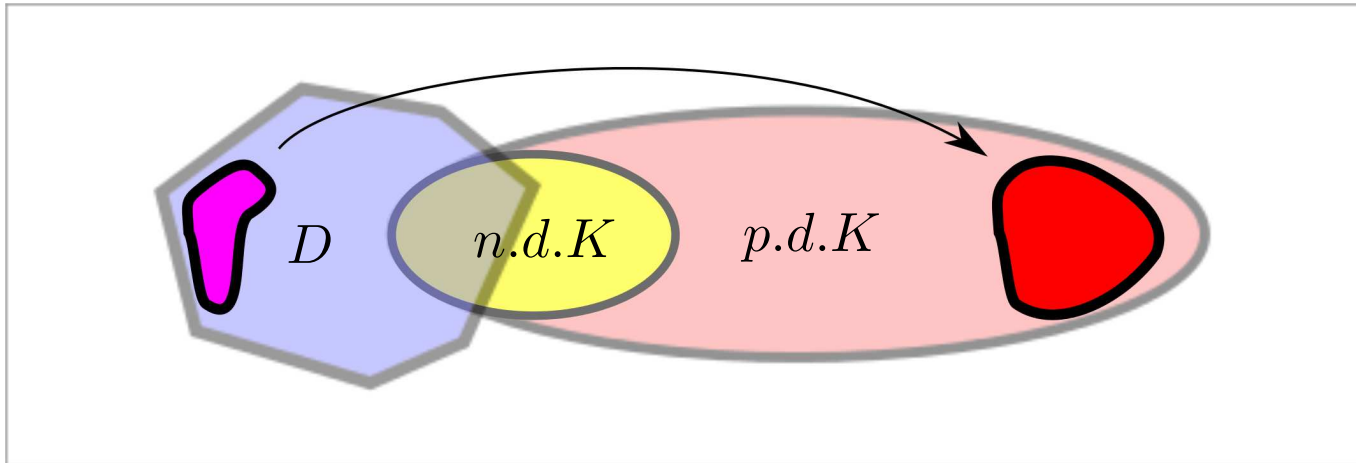
$$k_M(\mathbf{x}, \mathbf{y}) = \sum_{X \in U(r,c)} e^{-\langle X, M \rangle}$$

is positive definite (C., submitted).

Example: Transportation distance between discrete histograms



To wrap up



$$d(\mathbf{x}, \mathbf{y}) = \min_{\tau \in T(\mathbf{x}, \mathbf{y})} c(\tau), \quad \delta(\mathbf{x}, \mathbf{y}) = \text{soft-min}_{\tau \in T(\mathbf{x}, \mathbf{y})} c(\tau)$$

$e^{-\delta(\mathbf{x}, \mathbf{y})} = \sum_{\tau \in T(\mathbf{x}, \mathbf{y})} e^{-\langle \tau, \theta \rangle} = G_{T(\mathbf{x}, \mathbf{y})}(\theta)$ is positive definite in many (all) cases.