

# Topic 6: THE KERNEL TRICK

---

STAT 37710/CMSC 25400 Machine Learning  
Risi Kondor, The University of Chicago

# Enriching linear classifiers

- The hypothesis space of linear classifiers is limited.
- Could be enriched by a **feature mapping**  $\phi: \mathbf{x} \rightarrow \bar{\mathbf{x}}$  that throws in various nonlinear features, e.g.,

$$\bar{\mathbf{x}} = (x_1, x_2, \dots, x_n, x_1x_2, x_3x_5x_9, \sin(x_1), e^{-(x_1-4)^2}, \dots).$$

- However, this is ad-hoc and potentially very expensive.

# The kernel trick

Observations:

1. For both the perceptron and the SVM,  $\mathbf{w}$  is of the form  $\mathbf{w} = \sum_{j=1}^m \gamma_j \mathbf{x}_j$ .
2. The algorithm only uses  $\mathbf{w}$  through inner products such as

$$\mathbf{w} \cdot \mathbf{x}_i = \sum_{j=1}^m \gamma_j (\mathbf{x}_j \cdot \mathbf{x}_i)$$

So all that we need are dot products  $\mathbf{x}_i \cdot \mathbf{x}_j$ , or, in the feature space,  $\bar{\mathbf{x}}_i \cdot \bar{\mathbf{x}}_j$ .

**IDEA:** Define the **kernel**  $k(\mathbf{x}, \mathbf{x}') = \bar{\mathbf{x}} \cdot \bar{\mathbf{x}}'$ , write everything in terms of  $k$ , and don't ever bother computing  $\bar{\mathbf{x}} = \phi(\mathbf{x})$  or  $\bar{\mathbf{x}}' = \phi(\mathbf{x}')$  !



Beam me up, Scotty!

# Inner product

In infinite dimensional spaces (e.g., function spaces) the notion of dot product  $\mathbf{x} \cdot \mathbf{x}' = \mathbf{x}^\top \mathbf{x}'$  does not make sense. Instead, we have inner products.

## Definition

An inner product on a vector space  $V$  (over  $\mathbb{R}$ ) is a function  $\langle \cdot, \cdot \rangle : V \rightarrow \mathbb{R}$  satisfying

1.  $\langle u, v \rangle = \langle v, u \rangle$  (symmetry)
  2.  $\langle \alpha u, v \rangle = \alpha \langle u, v \rangle$  (linearity I)
  3.  $\langle u, v + w \rangle = \langle u, v \rangle + \langle u, w \rangle$  (linearity II)
  4.  $\langle u, u \rangle \geq 0$  with equality only if  $u = 0$  (positivity)
- for all  $u, v, w \in V$  and  $\alpha \in \mathbb{R}$ .

Two vectors  $u, v \in V$  are said to be **orthogonal** if  $\langle u, v \rangle = 0$ .

# Hilbert spaces

## Definition

A vector space  $\mathcal{H}$  is said to be a Hilbert space if

1.  $\mathcal{H}$  has an inner product  $\langle \cdot, \cdot \rangle : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$  ;
2.  $\mathcal{H}$  is complete with respect to the norm

$$\|x\| = \sqrt{\langle x, x \rangle}$$

induced by  $\langle \cdot, \cdot \rangle$  .

Hilbert spaces sound scary but are really just the natural generalization of  $\mathbb{R}^n$  to possibly infinite dimensions. They are inner product in which the inner product works “as expected” and consequently most of linear algebra carries over to them with no problems.

# The kernel trick (more explicitly)

Algorithms like the Perceptron and the SVM can work in any Hilbert space  $\mathcal{H}$ .

- To invoke the algorithm in  $\mathcal{H}$ , use a feature mapping  $\phi: \mathcal{X} \rightarrow \mathcal{H}$ .
- However, never compute  $\phi(\mathbf{x})$  explicitly. Instead, use the kernel (pull-back of the inner product)

$$k(\mathbf{x}, \mathbf{x}') := \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle.$$

- Since  $\mathcal{H}$  is typically much higher dimensional than  $\mathcal{X}$ , the decision surface in  $\mathcal{X}$  corresponding to a linear classifier in  $\mathcal{H}$  can be much more complex.

Learning algorithms that can exploit this trick are called Hilbert space methods or kernel methods.

# The vanilla perceptron

```
 $\mathbf{w} \leftarrow 0 ;$   
 $t \leftarrow 1 ;$   
while(true){  
  if  $\mathbf{w} \cdot \mathbf{x}_t \geq 0$  predict  $\hat{y}_t = 1$  ; else predict  $\hat{y}_t = -1$  ;  
  if  $((\hat{y}_t = -1) \text{ and } (y_t = 1))$  let  $\mathbf{w} \leftarrow \mathbf{w} + \mathbf{x}_t$  ;  
  if  $((\hat{y}_t = 1) \text{ and } (y_t = -1))$  let  $\mathbf{w} \leftarrow \mathbf{w} - \mathbf{x}_t$  ;  
   $t \leftarrow t + 1$  ;  
}
```

At any time  $t$  , the weight vector is of the form

$$\mathbf{w} = \sum_{i=1}^{t-1} c_i \mathbf{x}_i \quad \text{where} \quad c_i \in \{-1, 0, +1\} .$$



# The kernel perceptron

```
t ← 1 ;  
while(1){  
  if  $\sum_{i=1}^{t-1} c_i k(\mathbf{x}_i, \mathbf{x}_t) \geq 0$  predict  $\hat{y}_t = 1$  ; else predict  $\hat{y}_t = -1$  ;  
   $c_t \leftarrow 0$  ;  
  if  $((\hat{y}_t = -1) \text{ and } (y_t = 1))$  let  $c_t = 1$  ;  
  if  $((\hat{y}_t = 1) \text{ and } (y_t = -1))$  let  $c_t = -1$  ;  
   $t \leftarrow t + 1$  ;  
}
```

# The kernel SVM

Primal (forget about  $b$ )

$$\underset{w \in \mathcal{H}, \xi_1, \dots, \xi_m}{\text{minimize}} \quad \frac{1}{2} \|w\|^2 + \frac{C}{m} \sum_i \xi_i \quad \text{s.t.} \quad y_i \langle w, \phi(\mathbf{x}_i) \rangle \geq 1 - \xi_i \quad \xi_i \geq 0$$

Dual

$$\begin{aligned} \underset{\alpha_1, \dots, \alpha_m}{\text{maximize}} \quad L(\boldsymbol{\alpha}) &= \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j) \\ \text{subject to} \quad \sum_i y_i \alpha_i &= 0 \quad \text{and} \quad 0 \leq \alpha_i \leq \frac{C}{m} \quad \forall i \end{aligned}$$

Predict according to

$$\hat{y} = h(\mathbf{x}) = \text{sgn}(\langle w, \phi(\mathbf{x}) \rangle) = \text{sgn}\left(\sum_i \underbrace{y_i \alpha_i}_{\gamma_i} k(\mathbf{x}, \mathbf{x}_i)\right).$$

# Kernels

---

# So what does $k$ need to satisfy?

1. Symmetry:  $\langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle = \langle \phi(\mathbf{x}'), \phi(\mathbf{x}) \rangle \implies k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}', \mathbf{x})$ .
2. For any  $\mathbf{x}_1, \dots, \mathbf{x}_\ell$  and  $c_1, \dots, c_\ell \in \mathbb{R}$ , letting  $\xi = \sum_{i=1}^{\ell} c_i \phi(\mathbf{x}_i)$ , must have  $\langle \xi, \xi \rangle \geq 0 \implies$

$$\sum_{i=1}^{\ell} \sum_{j=1}^{\ell} c_i c_j \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle = \sum_{i=1}^{\ell} \sum_{j=1}^{\ell} c_i c_j k(\mathbf{x}_i, \mathbf{x}_j) \geq 0.$$

In general, any function  $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  satisfying 1. and 2. is called a symmetric, positive semi-definite (SPSD) kernel.

It turns out that this is all. Any SPD kernel  $k$  has a corresponding  $\mathcal{H}$  and  $\phi: \mathcal{X} \rightarrow \mathcal{H}$  such that  $\langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle = k(\mathbf{x}, \mathbf{x}')$  for any  $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$  (will see this below).

# Some typical kernels

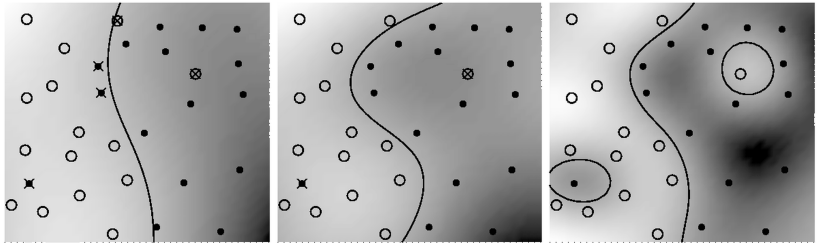
- Linear:  $k(\mathbf{x}, \mathbf{x}') = \mathbf{x} \cdot \mathbf{x}'$  (boring!)
- Polynomial:  $k(\mathbf{x}, \mathbf{x}') = (1 + \mathbf{x} \cdot \mathbf{x}')^p$
- Gaussian:

$$k(\mathbf{x}, \mathbf{x}') = \exp \left( -\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2} \right)$$

- Laplacian:  $k(x, x') = e^{-|x-x'|/\lambda}$
- String, graph, etc...

Since the hypothesis is  $h(\mathbf{x}) = \sum_i \alpha_i k(\mathbf{x}_i, \mathbf{x})$ , the kernel is like a similarity measure.

# SVM solutions with Gaussian kernel



# Closure properties of kernels

If  $k_1, k_2: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  are SPSPD kernels, then

- $k_+(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}')$
- $k_\times(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') \cdot k_2(\mathbf{x}, \mathbf{x}')$

are SPSPD kernels.

If  $k_1, k_2, \dots: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is a pointwise convergent sequence of SPSPD kernels, then  $\lim_{i \rightarrow \infty} k_i(\mathbf{x}, \mathbf{x}')$  is a SPSPD kernel.

If  $k_1: \mathcal{X}_1 \times \mathcal{X}_1 \rightarrow \mathbb{R}$  is an SPSPD kernel on  $\mathcal{X}_1$  and  $k_2: \mathcal{X}_2 \times \mathcal{X}_2 \rightarrow \mathbb{R}$  is a SPSPD kernel on  $\mathcal{X}_2$ , then

$$k((x_1, x_2), (x'_1, x'_2)) = k_1(x_1, x'_1) \cdot k_2(x_2, x'_2)$$

is a PSD kernel on  $\mathcal{X}_1 \times \mathcal{X}_2$ .

# Support vector machine example

Performance on classifying full MNIST dataset:

Classifier	Test Error
linear	8.4%
3-nearest-neighbor	2.4%
RBF-SVM	1.4 %
Tangent distance	1.1 %
LeNet	1.1 %
Boosted LeNet	0.7 %
Translation invariant SVM	0.56 %



# Kernel mystique

- But what does the kernel really capture about the data?
- What is this magical Hilbert space  $\mathcal{H}$  ?
- If  $\mathcal{H}$  is so high dimensional, what stops the SVM from overfitting like crazy?
- How does the kernel SVM control the complexity of the returned hypothesis?

## Reproducing kernel Hilbert spaces

---

# Reproducing kernel Hilbert space

Let  $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  be a positive definite kernel.

1. Define  $k_x: \mathcal{X} \rightarrow \mathbb{R}$  as the function  $k_x(x') = k(x, x')$ .
2. Define the function space  $\mathcal{H}_k^{\text{pre}}$  as the space of linear combinations

$$\xi(x') = \sum_{i=1}^m \alpha_i k_{x_i}(x')$$

for any  $m \in \mathbb{N}$ , any  $x_1, \dots, x_m \in \mathcal{X}$  and  $\alpha_1, \dots, \alpha_m \in \mathbb{R}$ .

3. Define  $\langle k_x, k_{x'} \rangle = k(x, x')$  and extend it to the rest of  $\mathcal{H}_k^{\text{pre}}$  by linearity.
4. Add to  $\mathcal{H}_k^{\text{pre}}$  the limit points of all Cauchy sequences.

The resulting space  $\mathcal{H}_k$  is the **Reproducing Kernel Hilbert Space (RKHS)** induced by  $k$ . Clearly, if  $\phi: x \mapsto k_x$ , then  $k(x, x') = \langle \phi(x), \phi(x') \rangle$  proving that for any PSD kernel there is a  $\mathcal{H}$  and an  $\phi$  which realize it in this way.

# The reproducing property

For any  $f \in \mathcal{H}_k$  expressible as a finite linear combination

$$f(x) = \sum_{i=1}^m \alpha_i k_{x_i}(x),$$

function evaluation and the inner product are linked by the remarkable property

$$f(x) = \sum_{i=1}^m \alpha_i k(x_i, x) = \sum_{i=1}^m \alpha_i \langle k_{x_i}, k_x \rangle = \left\langle \sum_{i=1}^m \alpha_i k_{x_i}, k_x \right\rangle = \langle f, k_x \rangle.$$

# Regularized Risk Minimization

---

# RRM in RKHS

Recall that at the abstract level many ML algorithms just perform some form of regularized risk minimization:

$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{F}} \left[ \underbrace{\frac{1}{m} \sum_{i=1}^m \ell(f(x_i), y_i)}_{\text{training error}} + \underbrace{\Omega(f)}_{\text{regularizer}} \right].$$

In general, searching an infinite dimensional space of functions on a computer is pretty challenging. However, if  $\mathcal{F} = \mathcal{H}_k$  and  $\Omega(f) = \|f\|^2$ , then

$$f(x_i) = \langle f, k_{x_i} \rangle, \quad \Omega(f) = \langle f, f \rangle,$$

so it reduces to a problem in linear algebra!

# The Representer Theorem

## Theorem (Wahba)

Let  $k: \mathcal{X} \rightarrow \mathcal{X} \rightarrow \mathbb{R}$  be a PSD kernel and  $\mathcal{H}_k$  be the corresponding RKHS. Then for any loss function  $\ell$  and any monotonically increasing  $\chi: \mathbb{R} \rightarrow \mathbb{R}$ , the solution to

$$\underset{f \in \mathcal{H}_k}{\text{minimize}} \frac{1}{m} \sum_{i=1}^m \ell(f(x_i), y_i) + \chi(\|f\|_{\mathcal{H}_k})$$

is of the form

$$f(x) = \sum_{i=1}^m \alpha_i k(x_i, x).$$

This is the key to making kernel machines implementable on computers, since we only need to optimize for the  $m$  real numbers  $\alpha_1, \dots, \alpha_m$ .

# RRM form of the SVM

$$\underset{w \in \mathcal{H}_k, \xi_1, \dots, \xi_m}{\text{minimize}} \quad \left[ \frac{1}{2} \|w\|^2 + \frac{C}{m} \sum_i \xi_i \right] \quad \text{s.t.} \quad y_i \langle w, \phi(\mathbf{x}_i) \rangle \geq 1 - \xi_i \quad \xi_i \geq 0$$

$$\Updownarrow$$

$$\underset{f \in \mathcal{H}_k, \xi_1, \dots, \xi_m}{\text{minimize}} \quad \left[ \frac{1}{2} \|f\|^2 + \frac{C}{m} \sum_i \xi_i \right] \quad \text{s.t.} \quad y_i f(\mathbf{x}_i) \geq 1 - \xi_i \quad \xi_i \geq 0.$$

$$\Updownarrow$$

$$\underset{f \in \mathcal{H}_k, \xi_1, \dots, \xi_m}{\text{minimize}} \quad \left[ \frac{1}{2} \|f\|^2 + \frac{C}{m} \sum_i (1 - y_i f(\mathbf{x}_i))_{\geq 0} \right]$$

So the soft margin SVM is a special case of RRM with  $\mathcal{F} = \mathcal{H}_k$ ,

$$\ell(\hat{y}, y) = (1 - y\hat{y})_{\geq 0}, \quad \Omega(f) = \frac{1}{2C} \|f\|^2, \quad h(x) = \text{sgn}(f(x)).$$

Question: But what does the regularizer  $\|f\|^2$  express?



# The RKHS of the Gaussian kernel

Theorem (Girosi et al., 1995)

For the Gaussian kernel  $k(x, x') = \exp(-\|x - x'\|^2 / (2\sigma^2))$  on  $\mathbb{R}^n$ ,

$$\|f\|_{\mathcal{H}_k}^2 \propto \int e^{\sigma^2 \omega^2 / 2} |\tilde{f}(\omega)|^2 d\omega,$$

where

$$\tilde{f}(\omega) = \int f(x) e^{-2\pi i \omega \cdot x} dx$$

is the Fourier transform of  $f$ .

→ support vector machines with the Gaussian kernel totally make sense!

# Three faces of the kernel

1. Inner product in feature space — mostly magic.
2. Measure of similarity between data points — pragmatic.
3. Responsible for regularization — makes sense.

## Examples of Kernels

---

# Kernels

A kernel  $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  may be specified either

- Explicitly, like the Gaussian RBF kernel

$$k(x, x') = \exp(-\|x - x'\|^2 / (2\sigma^2)).$$

- Implicitly, via some algorithm that computes  $k(x, x')$  for given  $x, x'$ .

Criteria for a good kernel:

- Positive semi-definiteness (and, of course, symmetry).
- Good notion of similarity between data / good regularizer.
- Efficiently computable.

Note:  $\mathcal{X}$  can be almost anything, doesn't need to be  $\mathbb{R}^d$ .

# Kernels between distributions

Let  $\mathcal{X}$  be the space of distributions on some space  $\mathcal{S}$ . Bhattacharyya kernel:

$$k_{\text{Bhatta}}(p, p') = \int \sqrt{p(x)} \sqrt{p'(x)} dx$$

Question: Is this a valid kernel (i.e., PSD)? In general, difficult to compute, however if  $p = \mathcal{N}(\mu, \Sigma)$  and  $p' = \mathcal{N}(\mu', \Sigma')$ , then

$$k_{\text{Bhatta}}(p, p') = |\Sigma|^{-1/4} |\Sigma'|^{-1/4} |\Sigma^\dagger|^{1/2} \exp\left(-\frac{1}{4}\mu^\top \Sigma^{-1}\mu - \frac{1}{4}\mu'^\top \Sigma'^{-1}\mu' + \frac{1}{2}\mu^\dagger^\top \Sigma^\dagger \mu^\dagger\right) \quad (1)$$

where  $\Sigma^\dagger = \left(\frac{1}{2}\Sigma^{-1} + \frac{1}{2}\Sigma'^{-1}\right)^{-1}$  and  $\mu^\dagger = \frac{1}{2}\Sigma^{-1}\mu + \frac{1}{2}\Sigma'^{-1}\mu'$ .

[K & Jebara, 2003]

# Contiguous substring kernel

Let  $\Sigma$  be an alphabet and let  $n_u(a)$  denote the number of places that the string  $u \in \Sigma^\ell$  appears in the string  $a \in \Sigma^*$  as a contiguous substring.

Given two strings  $a, b \in \Sigma^*$  the contiguous substring kernel is

$$k(a, b) = \sum_{u \in \Sigma^\ell} n_u(a) n_u(b).$$

- Question: Is this a valid kernel? Yes, because it corresponds to the feature map  $\phi: \Sigma^* \rightarrow \mathbb{N}^{\Sigma^\ell}$  with  $[\phi(a)]_u = n_u(a)$ .
- Question: What is the complexity of computing it?  $O(\ell |a| |b|)$

# Gappy substring kernels

Given an index sequence  $\mathbf{i} = (i_1, \dots, i_\ell)$  with  $1 \leq i_1 < i_2 < \dots < i_\ell$  let  $\mathbf{a}_{\mathbf{i}}$  denote the gappy substring  $a_{i_1} a_{i_2} \dots a_{i_\ell}$ . Let  $n_u^{\text{gap}}(\mathbf{a}) = |\{ \mathbf{i} \mid \mathbf{a}_{\mathbf{i}} = u \}|$ . Consider the kernel

$$k_\ell(\mathbf{a}, \mathbf{b}) = \sum_{u \in \Sigma^\ell} n_u^{\text{gap}}(\mathbf{a}) n_u^{\text{gap}}(\mathbf{b}).$$

- Question: Is this a valid kernel? Yes.
- Question: What is the complexity of computing it?  $O(\ell |\mathbf{a}| |\mathbf{b}|)$  using the dynamic programming recursion

$$k_\ell(\mathbf{a}_{1:i+1}, \mathbf{b}_{1:j}) = k_\ell(\mathbf{a}_{1:i}, \mathbf{b}_{1:j}) + \sum_{p=\ell}^j \mathbb{I}(a_{i+1} = b_p) k_{\ell-1}(\mathbf{a}_{1:i}, \mathbf{b}_{1:p-1})$$

# Random walk kernel on graphs

Let  $G = (V, E)$  be a graph with adjacency matrix  $A$ , and  $\text{path}_p(x, x')$  be the set of all path from  $x$  to  $x'$  in  $G$  of length  $p$ . Consider

$$k_{2\ell}(x, x') = |\text{path}_{2\ell}(x, x')|.$$

- Question: Is this a valid kernel? Yes.
- Question: What is the complexity of computing it?

$$k_{2\ell}(x, x') = [A^{2\ell}]_{x, x'}$$

Question: What is the problem with the random walk kernel?

[Gärtner]



# Diffusion kernel on graphs

Imagine a lazy random walker, which, when he is at a vertex of degree  $d$

- with probability  $dp$  moves to one of the neighbors (selected randomly)
- with probability  $1 - dp$  stays in place.

The transition matrix of this process is  $T = I + pL$ , where  $L$  is the **graph Laplacian**

$$[L]_{i,j} = \begin{cases} 1 & i \sim j \\ -d_i & i = j \\ 0 & \text{otherwise.} \end{cases}$$

After  $n$  timesteps the distribution of the random walker is given by  $T^n$ . Now take the continuous limit where  $n \rightarrow \infty$  and simultaneously  $p = 1/n$ .

# Diffusion kernel on graphs

After time  $t$  the distribution is given by

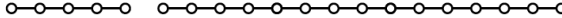
$$K_t = e^{tL} = \lim_{n \rightarrow \infty} \left( I + \frac{tL}{n} \right)^n = I + tL + \frac{t^2}{2}L^2 + \frac{t^3}{6}L^3 + \dots$$

The **diffusion kernel** on  $G$  with parameter  $t$  is  $k_t(x, x') = [K_t]_{x, x'}$ .

- Question: Is this a valid kernel? Yes, the exponential of a symmetric matrix is always PSD.
- Question: What is its computational complexity? Typically  $O(|V|^3)$ .

[K. & Lafferty, 2002]

# Eigenvectors of the graph Laplacian



(a) a linear unweighted graph with two segments



$$\lambda_1 = 0.00$$



$$\lambda_2 = 0.00$$



$$\lambda_3 = 0.04$$



$$\lambda_4 = 0.17$$



$$\lambda_5 = 0.38$$



$$\lambda_6 = 0.38$$



$$\lambda_7 = 0.66$$



$$\lambda_8 = 1.00$$



$$\lambda_9 = 1.38$$



$$\lambda_{10} = 1.38$$



$$\lambda_{11} = 1.79$$



$$\lambda_{12} = 2.21$$



$$\lambda_{13} = 2.62$$



$$\lambda_{14} = 2.62$$



$$\lambda_{15} = 3.00$$



$$\lambda_{16} = 3.34$$



$$\lambda_{17} = 3.62$$



$$\lambda_{18} = 3.62$$



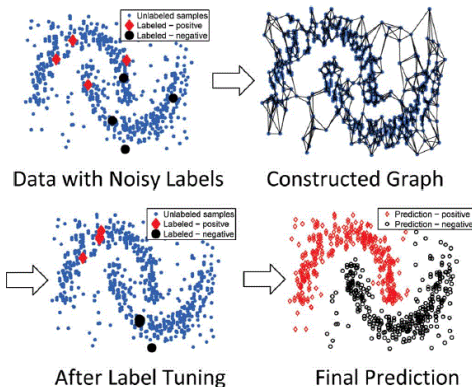
$$\lambda_{19} = 3.83$$



$$\lambda_{20} = 3.96$$

(b) the eigenvectors and eigenvalues of the Laplacian  $L$

# Graph kernels for semi-supervised ML



- Often labeled data is expensive but unlabeled data is abundant.
- Use the unlabeled data to construct a graph (mesh)  $\rightarrow$  graph kernel  $k$ .
- Supervised learning with  $k$  enforces that  $f$  be smooth wrt. this graph.

# FURTHER READING

- B. Schölkopf and A. J. Smola: **Learning with Kernels**
- N. Christianini and J-S Taylor: **An Introduction to Support Vector Machines...**
- I. Steinwart: **Support Vector Machines**