

Kernel methods

MATH-412 - Statistical Machine Learning

Making models non-linear with a feature map

Idea : make non-linear transformation of the data first

- Quadratic map :

$$\phi(\mathbf{x}) = (x_1, \dots, x_p, x_1^2, \dots, x_p^2, x_1x_2, x_1x_3, \dots, x_{p-1}x_p)$$

- Fourier basis, spline basis, wavelet basis

Regularized empirical risk minimization with a mapping ϕ :

$$\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}^\top \phi(x_i), y_i) + \lambda \|\mathbf{w}\|^2.$$

Representer theorem (simple version with the feature map)

Theorem (Kimmeldorf and Wahba, 1971)

Consider the optimization problem

$$\min_{\mathbf{w} \in \mathbb{R}^d} L(\mathbf{w}^\top \phi(x_1), \dots, \mathbf{w}^\top \phi(x_n)) + \lambda \|\mathbf{w}\|^2$$

Then any local minimum is of the form $\mathbf{w} = \sum_{i=1}^n \alpha_i \phi(x_i),$

for some vector $\alpha \in \mathbb{R}^n$. Interpretation : $\mathbf{w} \in \text{span}(\phi(x_1), \dots, \phi(x_n))$.

So that

$$f_{\mathbf{w}}(x) = \mathbf{w}^\top \phi(x) = \sum_{i=1}^n \alpha_i \langle \phi(x_i), \phi(x) \rangle = \sum_{i=1}^n \alpha_i K(x_i, x).$$

Applying the representer theorem to the ERM problem

$$\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}^\top \boldsymbol{\phi}(x_i), y_i) + \lambda \|\mathbf{w}\|^2.$$

By the theorem of Kimmeldorf and Wahba, $\mathbf{w}^* = \sum_{j=1}^n \alpha_j^* \boldsymbol{\phi}(x_j)$.

So replacing in the previous expression, we get

$$\begin{aligned} \min_{\boldsymbol{\alpha}} \frac{1}{n} \sum_{i=1}^n \ell\left(\sum_{j=1}^n \alpha_j \langle \boldsymbol{\phi}_j(x_j), \boldsymbol{\phi}_i(x_i) \rangle, y_i\right) + \lambda \left\| \sum_{j=1}^n \alpha_j \boldsymbol{\phi}(x_j) \right\|^2. \\ \min_{\boldsymbol{\alpha}} \frac{1}{n} \sum_{i=1}^n \ell\left(\sum_{j=1}^n \alpha_j K_{ij}, y_i\right) + \lambda \sum_{1 \leq i, j \leq n} \alpha_i \alpha_j K_{ij}, \end{aligned}$$

with $K_{ij} = K(x_i, x_j) = \langle \boldsymbol{\phi}(x_i), \boldsymbol{\phi}(x_j) \rangle$ the values of a kernel function on pairs of input datapoints.

The ERM expressed with the kernel matrix

We rewrote $\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}^\top \phi(x_i), y_i) + \lambda \|\mathbf{w}\|^2$ as :

$$\min_{\boldsymbol{\alpha}} \frac{1}{n} \sum_{i=1}^n \ell\left(\sum_{j=1}^n \alpha_j K_{ij}, y_i\right) + \lambda \sum_{1 \leq i, j \leq n} \alpha_i \alpha_j K_{ij},$$

with $K_{ij} = K(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle$.

This can be rewritten in matrix vector form as

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^n} \frac{1}{n} \sum_{i=1}^n \ell\left(\mathbf{K}_{i \cdot} \boldsymbol{\alpha}, y_i\right) + \lambda \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha}.$$

Furthermore to make a prediction, our predictor is computed as

$$\hat{f}(x) = \mathbf{w}^{\star \top} \phi(x) = \sum_{j=1}^n \alpha_j^{\star} K(x_j, x).$$

The kernel matrix when $\phi(\mathbf{x}) = \mathbf{x}$.

Based on the design matrix \mathbf{X} , two symmetric p.s.d. matrices are natural :

- the *empirical covariance matrix* (assuming \mathbf{X} is centered)

$$\widehat{\Sigma} = \frac{1}{n} \mathbf{X}^\top \mathbf{X}$$

$$\widehat{\Sigma}_{k\ell} = \widehat{\text{Cov}}(X^{(k)}, X^{(\ell)}) = \left\langle \frac{1}{\sqrt{n}} \mathbf{x}^k, \frac{1}{\sqrt{n}} \mathbf{x}^\ell \right\rangle$$

- the *kernel matrix* or *Gram matrix*

$$\mathbf{K} = \mathbf{X} \mathbf{X}^\top$$

$$K_{ij} = \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

\mathbf{K} is simply the matrix of all dot products. \mathbf{K} encodes information about the data vectors $\mathbf{x}_i = \mathbf{X}_{i\cdot}^\top$ while $\widehat{\Sigma}$ encodes information about the variables $\mathbf{x}^k = \mathbf{X}_{\cdot k}$

Properties of the kernel matrix when $\phi(\mathbf{x}) = \mathbf{x}$.

The kernel matrix contains a lot of information about the data :

- It contains the information about all the distances between all pairs of data points (and between each data points and the origin). Indeed,

$$\|\mathbf{x}_i - \mathbf{x}_j\|_2^2 = \mathbf{K}_{ii} - 2\mathbf{K}_{ij} + \mathbf{K}_{jj}.$$

- As a consequence, any factorization of the matrix \mathbf{K} of the form

$$\mathbf{K} = \mathbf{R}\mathbf{R}^\top,$$

retrieves a representation of the data up to an isometry. This can be obtained for example by the Cholesky decomposition.

Why is this useful?

Dot products in feature space

Let $\mathbf{x} = (x_1, x_2) \in \mathbb{R}^2$ and $\phi(\mathbf{x}) = (x_1, x_2, x_1^2, x_2^2, \sqrt{2}x_1x_2)^\top$.

$$\begin{aligned}\langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle &= x_1y_1 + x_2y_2 + x_1^2y_1^2 + x_2^2y_2^2 + 2x_1x_2y_1y_2 \\ &= x_1y_1 + x_2y_2 + (x_1y_1)^2 + (x_2y_2)^2 + 2(x_1y_1)(x_2y_2) \\ &= \langle \mathbf{x}, \mathbf{y} \rangle + \langle \mathbf{x}, \mathbf{y} \rangle^2\end{aligned}$$

For $\mathbf{w} = (0, 0, 1, 1, 0)^\top$, $\mathbf{w}^\top \phi(\mathbf{x}) - 1 \leq 0 \Leftrightarrow \|\mathbf{x}\|^2 \leq 1$.

Linear separators in \mathbb{R}^5 correspond to conic separators in \mathbb{R}^2 .

<https://www.youtube.com/watch?v=Q7vT0--5VII>

Let $\mathbf{x} = (x_1, \dots, x_p) \in \mathbb{R}^p$ and

$$\phi(\mathbf{x}) = (x_1, \dots, x_p, x_1^2, \dots, x_p^2, \sqrt{2}x_1x_2, \dots, \sqrt{2}x_i x_j, \dots, \sqrt{2}x_{p-1}x_p)^\top.$$

Still have

$$\langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle = \langle \mathbf{x}, \mathbf{y} \rangle + \langle \mathbf{x}, \mathbf{y} \rangle^2$$

But explicit mapping too expensive to compute : $\phi(\mathbf{x}) \in \mathbb{R}^{p+p(p+1)/2}$.

Which abstract space is a good predictor space ?

Require that

- (1) the space should be a Hilbert space $(\mathcal{H}, \|\cdot\|_{\mathcal{H}})$
- (2) $\forall x \in \mathcal{X}$, the *evaluation functional* $f \mapsto f(x)$ is *continuous* from $(\mathcal{H}, \|\cdot\|_{\mathcal{H}})$ to \mathbb{R} .
 - This is equivalent to requiring that for a given $x \in \mathcal{X}$:
if $\|f - g\|_{\mathcal{H}}$ is small then $|f(x) - g(x)|$ should be small.
 - The motivation is that we would like that
$$(\|\hat{f}_n - f^*\|_{\mathcal{H}} \rightarrow 0) \Rightarrow (\hat{f}_n(x) \rightarrow f^*(x))$$

Riesz Representation Theorem

Let \mathcal{H} be a Hilbert space, and $\psi : \mathcal{H} \rightarrow \mathbb{R}$ be a *continuous* linear form, then there exists $h_{\psi} \in \mathcal{H}$ such that

$$\forall f \in \mathcal{H}, \psi(f) = \langle h_{\psi}, f \rangle_{\mathcal{H}}.$$

Under (1) and (2) by this theorem, there must exist an element $h_x \in \mathcal{H}$ such that

$$\forall f \in \mathcal{H}, f(x) = \langle h_x, f \rangle_{\mathcal{H}}.$$

Reproducing Kernel Hilbert Space

So if \mathcal{H} is a Hilbert space of functions in which the *evaluation functionals* are continuous, then by the Riesz representation theorem, there must exist an element $h_x \in \mathcal{H}$ such that

$$\forall f \in \mathcal{H}, \quad f(x) = \langle h_x, f \rangle_{\mathcal{H}}.$$

But then by definition $h_y(x) = \langle h_x, h_y \rangle_{\mathcal{H}} = h_x(y)$.

Define the *reproducing kernel* as the function

$$\begin{aligned} K : \mathcal{X} \times \mathcal{X} &\rightarrow \mathbb{R} \\ (x, y) &\mapsto \langle h_x, h_y \rangle_{\mathcal{H}}. \end{aligned}$$

By definition $h_x(\cdot) = K(x, \cdot)$ so that

$$f(x) = \langle K(x, \cdot), f \rangle_{\mathcal{H}} \quad \text{and} \quad \langle K(x, \cdot), K(y, \cdot) \rangle_{\mathcal{H}} = K(x, y).$$

A space with these properties is called a *reproducing kernel Hilbert space* (RKHS).

Positive definite functions

$$(x, y) \mapsto K(x, y)$$

is a *positive definite function* if the matrix constructed as

$$\mathbf{K} = \begin{bmatrix} K(x_1, x_1) & \dots & \dots & K(x_1, x_n) \\ K(x_2, x_1) & \dots & \dots & K(x_2, x_n) \\ \vdots & & & \vdots \\ K(x_n, x_1) & \dots & \dots & K(x_n, x_n) \end{bmatrix}$$

is a positive semi-definite matrix

$$\text{i.e., } \forall \boldsymbol{\alpha} \in \mathbb{R}^n, \quad \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha} \geq 0,$$

for any choice of x_1, \dots, x_n and any value of n .

A reproducing kernel is a positive definite function

Proposition

A reproducing kernel is a positive definite function.

Proof of the claim The reproducing kernel is necessarily a *symmetric positive definite function* since for all $x_1, \dots, x_n \in \mathcal{X}$, we have $\langle K(x_i, \cdot), K(x_j, \cdot) \rangle_{\mathcal{H}} = K(x_i, x_j)$, and thus for all $\alpha_1, \dots, \alpha_n \in \mathbb{R}$.

$$0 \leq \left\langle \sum_i \alpha_i K(x_i, \cdot), \sum_j \alpha_j K(x_j, \cdot) \right\rangle_{\mathcal{H}} = \sum_{i,j} \alpha_i \alpha_j K(x_i, x_j),$$

with equality if and only if $\alpha_i = 0$ for all i .

Converse ?

Yes, any symmetric positive definite function is the reproducing kernel of a RKHS (Aronszajn, 1950).

Moore-Aronszajn theorem

Theorem

A symmetric function K on \mathcal{X} is positive definite if and only if there exists a Hilbert space \mathcal{H} and a mapping

$$\begin{aligned}\phi : \mathcal{X} &\rightarrow \mathcal{H} \\ x &\mapsto \phi(x)\end{aligned}$$

such that $K(x, y) = \langle \phi(x), \phi(y) \rangle_{\mathcal{H}}$.

- In fact, this mapping is $\phi(x) = h_x$
- Such symmetric p.d. functions are often called *Mercer kernels*.
- We will not show this theorem in this course.

Common RKHSes for $\mathcal{X} = \mathbb{R}^p$

Linear kernel

- $K(\mathbf{x}, \mathbf{y}) = \mathbf{x}^\top \mathbf{y}$
- $\mathcal{H} = \{f_{\mathbf{w}} : \mathbf{x} \mapsto \mathbf{w}^\top \mathbf{x} \mid \mathbf{w} \in \mathbb{R}^p\}$
- $\|f_{\mathbf{w}}\|_{\mathcal{H}} = \|\mathbf{w}\|_2$

Polynomial kernel

- $K_h(\mathbf{x}, \mathbf{y}) = (\gamma + \mathbf{x}^\top \mathbf{y})^d$
- \mathcal{H}

Radial Basis Function kernel (RBF)

- $K_h(\mathbf{x}, \mathbf{y}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{y}\|_2^2}{2h}\right)$
- $\mathcal{H} = \text{Gaussian RKHS}$

Representer theorem

Theorem (Kimmeldorf and Wahba, 1971)

Consider the optimization problem

$$\min_{f \in \mathcal{H}} L(f(x_1), \dots, f(x_n)) + \lambda \|f\|_{\mathcal{H}}^2$$

Then any local minimum is of the form $f = \sum_{i=1}^n \alpha_i K(x_i, \cdot),$

where K is the reproducing kernel associated with the RKHS \mathcal{H} and α is a vector in \mathbb{R}^n .

Proof Indeed, let f be a local minimum and consider the subspace

$$\mathcal{S} = \{g \mid g = \sum_{i=1}^n \alpha_i K(x_i, \cdot), \quad \alpha \in \mathbb{R}^n\}.$$

Representer theorem

We can decompose $f = f_{\parallel} + f_{\perp}$ with $f_{\parallel} = \text{Proj}_{\mathcal{S}}(f)$. We then have

$$f_{\perp}(x_i) = \langle f_{\perp}, K(x_i, \cdot) \rangle_{\mathcal{H}} = 0 \quad \text{and} \quad \langle f_{\perp}, f_{\parallel} \rangle_{\mathcal{H}} = 0.$$

Thus

$$\begin{aligned} & L(f(x_1), \dots, f(x_n)) + \lambda \|f\|_{\mathcal{H}}^2 \\ = & L(f_{\parallel}(x_1), \dots, f_{\parallel}(x_n)) + \lambda (\|f_{\parallel}\|_{\mathcal{H}}^2 + 2\langle f_{\perp}, f_{\parallel} \rangle_{\mathcal{H}} + \|f_{\perp}\|_{\mathcal{H}}^2) \\ = & L(f_{\parallel}(x_1), \dots, f_{\parallel}(x_n)) + \lambda \|f_{\parallel}\|_{\mathcal{H}}^2 + \lambda \|f_{\perp}\|_{\mathcal{H}}^2 \end{aligned}$$

So that we must have $f_{\perp} = 0$.

Regularized ERM for f in a RKHS

$$\min_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i) + \lambda \|f\|_{\mathcal{H}}^2 \quad (\text{P})$$

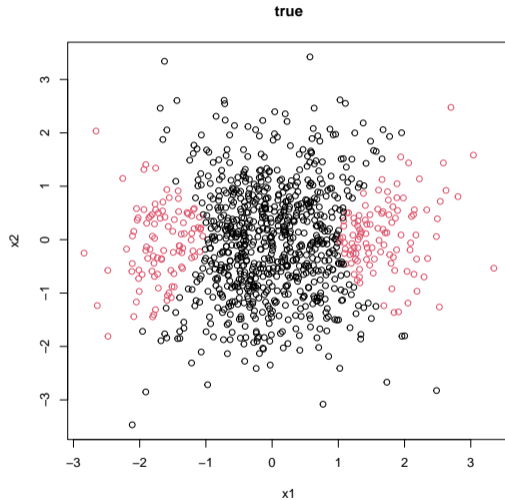
By the representer theorem, the solution of the regularized empirical risk minimization problem lies in the subspace of \mathcal{H} generated by the point x_i , i.e.,

$$f^* = \sum_{j=1}^n \alpha_j K(x_j, \cdot) \quad \text{for some } \alpha_i \in \mathbb{R}. \quad (\text{R})$$

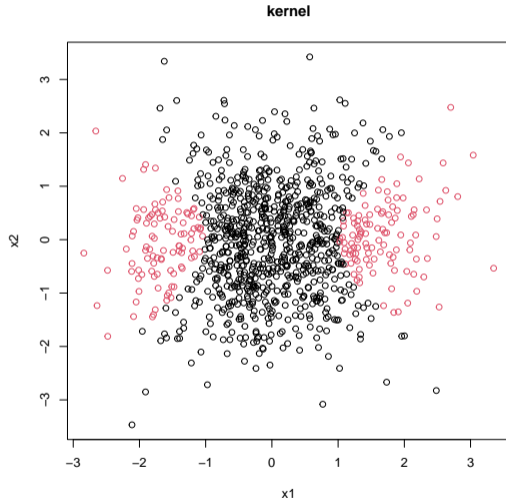
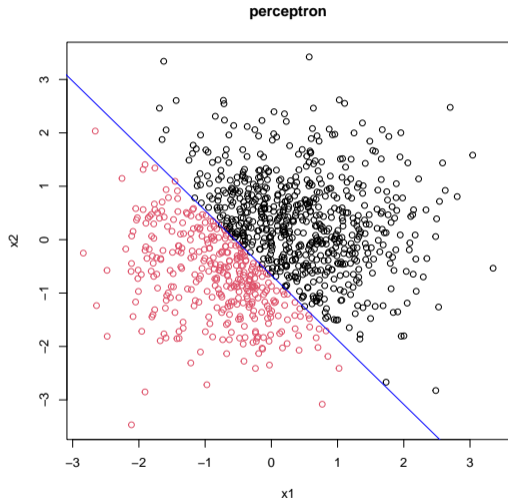
The solution of (P) is therefore of the form (R) with $\alpha \in \mathbb{R}^n$ the solution of

$$\min_{\alpha \in \mathbb{R}^n} \frac{1}{n} \sum_{i=1}^n \ell\left(\sum_{j=1}^n \alpha_j K(x_j, x_i), y_i\right) + \lambda \sum_{1 \leq i, j \leq n} \alpha_i \alpha_j K(x_i, x_j).$$

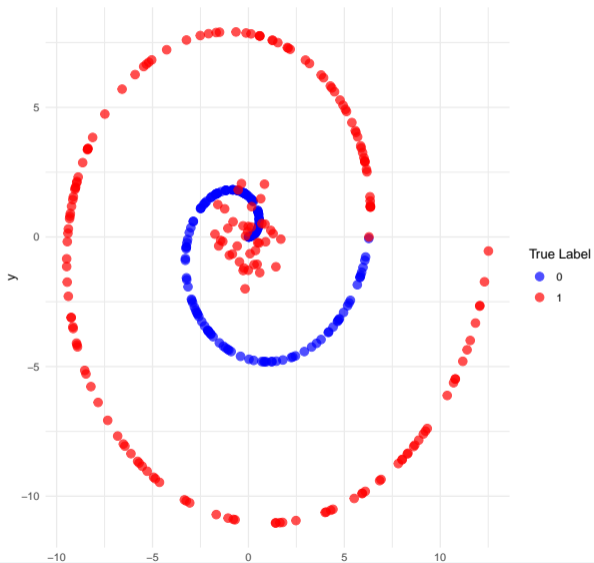
Hyperbola example



Hyperbola example

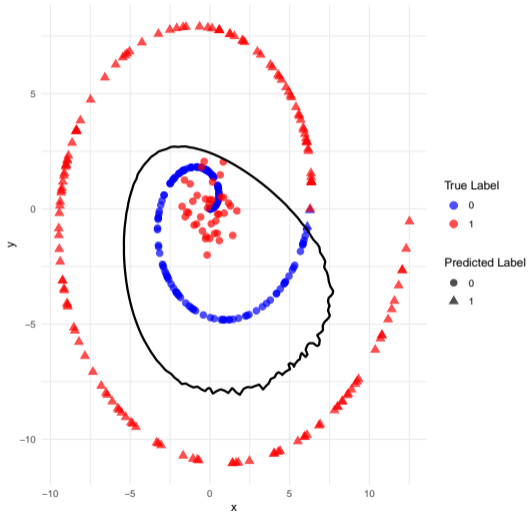


“Roll” example

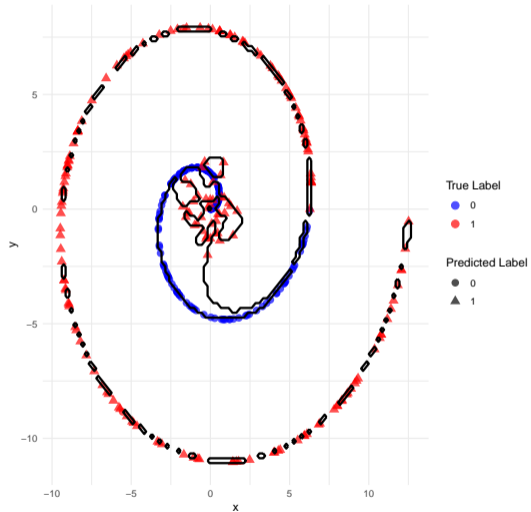


“Roll” example

Polynomial(3) Kernel



Gaussian(0.5) Kernel



$\|f\|_{\mathcal{H}}$ measures the smoothness of the function f

Indeed :

$$|f(x) - f(x')| = |\langle f, K(x, \cdot) - K(x', \cdot) \rangle_{\mathcal{H}}| \leq \|f\|_{\mathcal{H}} \|K(x, \cdot) - K(x', \cdot)\|_{\mathcal{H}}$$

→ f is Lipschitz with respect to the ℓ^2 distance induced by the RKHS

$$d(x, x') = \|K(x, \cdot) - K(x', \cdot)\|_{\mathcal{H}} = \sqrt{K(x, x) + K(x', x') - 2K(x, x')}$$

→ $\|f\|_{\mathcal{H}}$ is the Lipschitz constant

Kernel combinations

Assume K, K_1 and K_2 are positive definite functions,
then the following are still p.d. kernel functions :

Sum of kernels : For $\alpha_1, \alpha_2 > 0$, $\tilde{K}(x, y) = \alpha_1 K_1(x, y) + \alpha_2 K(x, y)$

Limits of kernels : $K(x, y) = \lim_{n \rightarrow \infty} K_n(x, y)$

Pointwise product : $\tilde{K}(x, y) = K_1(x, y) K_2(x, y)$

Pairwise kernel : $\tilde{K}(x, y) = \sum_{z \in \mathcal{Z}} K(x, z) K(z, y)$

Normalized kernel : $\tilde{K}(x, y) = \frac{K(x, y)}{\sqrt{K(x, x) K(y, y)}} = \cos \angle(\phi(x), \phi(y))$

In terms of kernel matrices

Pointwise product : $\tilde{K} = K_1 \odot K_2$ (Hadamard product)

Pairwise kernel : $\tilde{K} = K^2$ (Matrix product)

Scaling...

The kernelized form $\min_{\alpha} \frac{1}{n} \sum_{i=1}^n \ell(\alpha^\top \mathbf{k}_i, y_i) + \frac{\lambda}{2} \alpha^\top \mathbf{K} \alpha$

requires to compute $\mathbf{K} \in \mathbb{R}^{n \times n}$.

- The cost of working with kernels **quadratic in n** .
- ... unless \mathbf{K} is low rank, e.g. for the linear kernel
- This is a price to pay to work in very high/infinite dimensional spaces
- It is however possible to
 - compute low rank approximations to \mathbf{K} using *Nyström's method* (Williams and Seeger, 2001; Gittens and Mahoney, 2016) or
 - use greedy approximation schemes (Smola et al., 2000)
 - compute directly lower/finite dimensional approximation to the feature map using *random features expansions* (Rahimi and Recht, 2007; Bach, 2017; Yang et al., 2017)

Kernel ridge regression

$$\min_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda \|f\|_{\mathcal{H}}^2$$
$$\min_{f \in \mathcal{H}} \frac{1}{n} \|\mathbf{y} - \mathbf{f}\|_2^2 + \lambda \|\mathbf{f}\|_{\mathcal{H}}^2 \quad \text{with} \quad \mathbf{f} = (f(x_1), \dots, f(x_n)).$$

By the representer property $\hat{f}(x) = \sum_{i=1}^n \alpha_i K(x_i, x)$, so that $\frac{1}{2} \|\mathbf{f} - \mathbf{y}\|_2^2 = \frac{1}{2} \|\mathbf{K}\boldsymbol{\alpha} - \mathbf{y}\|_2^2$.

The regularized empirical risk is $\frac{1}{2n} \|\mathbf{K}\boldsymbol{\alpha} - \mathbf{y}\|_2^2 + \frac{\lambda}{2} \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha}$

and the minimizers are of the form $\boldsymbol{\alpha}^* + \mathbf{h}$ with $\boldsymbol{\alpha}^* = (\mathbf{K} + \lambda n \mathbf{I})^{-1} \mathbf{y}$, and $\mathbf{h} \in \text{Ker}(\mathbf{K})$.

Finally $\hat{f}(x) = \sum_{i=1}^n \alpha_i^* K(x_i, x)$ because $\sum_{i=1}^n h_i K(x_i, \cdot) = 0, \forall \mathbf{h} \in \text{Ker}(\mathbf{K})$.

Convolution vs Mercer kernels

In this course we encountered two type of kernels

Convolution kernels

Used for density estimation and by the Nadaraya-Watson estimator

$$K_{\delta}(x - y) = h\left(\frac{\|x - y\|}{\delta}\right)$$

- e.g. Epanechnikov, tricube or Gaussian kernel

Mercer kernels

... or simply positive definite kernel functions, which by Aronszajn's theorem provide the inner product of a RKHS

$$K(x, y)$$

- e.g. linear, polynomial, Laplace, Gaussian kernel and more

• Some are actually both

Summary

- Every positive definite function (Mercer kernel) is associated with a RKHS
- Regularized ERM can be kernelized (e.g. ridge regression)
- Many other algorithms in ML can be kernelized (e.g. kernel PCA)
- The representer theorem of Kimmeldorf and Wahba (1971) guarantees that a large class of optimization problems in RKHS can be reformulated as a finite-dimensional optimization problem.
- Using kernels directly has complexity n^2 but there are efficient approximation schemes (Nyström, random feature expansions)
- Mercer kernels should not be confused with convolution kernels used for density estimation and by Nadaraya-Watson estimators.

References I

- Aronszajn, N. (1950). Theory of reproducing kernels. *Transactions of the American mathematical society*, 68(3) :337–404.
- Bach, F. (2017). On the equivalence between kernel quadrature rules and random feature expansions. *The Journal of Machine Learning Research*, 18(1) :714–751.
- Gittens, A. and Mahoney, M. W. (2016). Revisiting the Nyström method for improved large-scale machine learning. *The Journal of Machine Learning Research*, 17(1) :3977–4041.
- Rahimi, A. and Recht, B. (2007). Random features for large-scale kernel machines. In *Proceedings of the 20th International Conference on Neural Information Processing Systems*, pages 1177–1184. Curran Associates Inc.
- Smola, A., Schölkopf, B., and Langley, P. (2000). Sparse greedy matrix approximation for machine learning. In *17th International Conference on Machine Learning, Stanford, 2000*, pages 911–911.
- Williams, C. K. and Seeger, M. (2001). Using the Nyström method to speed up kernel machines. In *Advances in Neural Information Processing Systems*, pages 682–688.
- Yang, Y., Pilanci, M., and Wainwright, M. (2017). Randomized sketches for kernels : Fast and optimal nonparametric regression. *The Annals of Statistics*, 45(3) :991–1023.