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 :

$$\phi(\mathbf{x}) = (x_1, \dots, x_p, x_1^2, \dots, x_p^2, x_1 x_2, x_1 x_3, \dots, x_{p-1} x_p)$$

• Fourier basis, spline basis, wavelet basis

Regularized empirical risk minimization with a mapping ϕ :

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

Math-412 Kernel methods 2/18

Representer theorem (simple version with the feature map)

Theorem (Kimmeldorf and Wahba, 1971)

Consider the optimization problem

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

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

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

So that
$$f_{\boldsymbol{w}}(x) = \boldsymbol{w}^{\top} \boldsymbol{\phi}(x) = \sum_{i=1}^{n} \alpha_{i} \langle \boldsymbol{\phi}(x_{i}), \boldsymbol{\phi}(x) \rangle = \sum_{i=1}^{n} \alpha_{i} K(x_{i}, x).$$

Kernel methods 3/18

Applying the representer theorem to the ERM problem

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

By the theorem of Kimmeldorf and Wahba, $\boldsymbol{w}^{\star} = \sum_{j=1}^{n} \alpha_{j}^{\star} \boldsymbol{\phi}(x_{j}).$

So replacing in the previous expression, we get

$$\min_{\alpha} \frac{1}{n} \sum_{i=1}^{n} \ell \left(\sum_{j=1}^{n} \alpha_{j} \langle \phi_{j}(x_{j}), \phi_{i}(x_{i}) \rangle, y_{i} \right) + \lambda \left\| \sum_{j=1}^{n} \alpha_{j} \phi(x_{j}) \right\|^{2}.$$

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

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

Math-412 Kernel methods 4/1

The ERM expressed with the kernel matrix

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

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

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

This can be rewritten in matrix vector form as

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

Furthermore to make a prediction, our predictor is computed as

$$\widehat{f}(x) = \boldsymbol{w}^{\star \top} \boldsymbol{\phi}(x) = \sum_{j=1}^{n} \alpha_{j}^{\star} K(x_{j}, x).$$

Math-412 Kernel methods 5/18

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

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

ullet the *empirical covariance matrix* (assuming $oldsymbol{X}$ is centered)

$$\widehat{oldsymbol{\Sigma}} = rac{1}{n} oldsymbol{X}^ op oldsymbol{X}$$

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

• the kernel matrix or Gram matrix

$$oldsymbol{K} = oldsymbol{X} oldsymbol{X}^{ op}$$

$$\boldsymbol{K}_{ij} = \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

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

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,

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

ullet As a consequence, any factorization of the matrix $oldsymbol{K}$ of the form

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

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 $\phi(\mathbf{x}) = (x_1, x_2, x_1^2, x_2^2, \sqrt{2}x_1x_2)^\top$.

$$\langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle = x_1 y_1 + x_2 y_2 + x_1^2 y_1^2 + x_2^2 y_2^2 + 2x_1 x_2 y_1 y_2$$

= $x_1 y_1 + x_2 y_2 + (x_1 y_1)^2 + (x_2 y_2)^2 + 2(x_1 y_1)(x_2 y_2)$
= $\langle \mathbf{x}, \mathbf{y} \rangle + \langle \mathbf{x}, \mathbf{y} \rangle^2$

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

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

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

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

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 : $\phi(\mathbf{x}) \in \mathbb{R}^{p+p(p+1)/2}$.

Math-412 Kernel methods 9/18

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

$$\left(\|\widehat{f}_n - f^*\|_{\mathcal{H}} \to 0\right) \Rightarrow \left(\widehat{f}_n(x) \to f^*(x)\right)$$

Riesz Representation Theorem

Let $\mathcal H$ be a Hilbert space, and $\psi:\mathcal H\to\mathbb R$ 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

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

Math-412 Kernel methods 10/18

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

$$K: \mathcal{X} \times \mathcal{X} \to \mathbb{R}$$

 $(x, y) \mapsto \langle h_x, h_y \rangle_{\mathcal{H}}.$

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

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

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

Math-412 Kernel methods 11/18

Positive definite functions

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

is a positive definite function if the matrix constructed as

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

is a positive semi-definite matrix

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

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

Math-412 Kernel methods 12/18

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 $x_1, \ldots, x_n \in \mathcal{X}$, we have $\langle K(x_i, \cdot), K(x_j, \cdot) \rangle_{\mathcal{H}} = K(x_i, x_j)$, and thus for all $\alpha_1, \ldots, \alpha_n \in \mathbb{R}$.

$$0 \le \left\langle \sum_{i} \alpha_{i} K(x_{i}, \cdot), \sum_{j} \alpha_{j} K(x_{j}, \cdot) \right\rangle_{\mathcal{H}} = \sum_{i, j} \alpha_{i} \alpha_{j} K(x_{i}, x_{j}),$$

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).

Math-412 Kernel methods 13/18

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

$$\phi: \mathcal{X} \to \mathcal{H}$$
$$x \mapsto \phi(x)$$

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.

14/18

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

Linear kernel

- $\bullet K(\mathbf{x}, \mathbf{y}) = \mathbf{x}^{\top} \mathbf{y}$
- $\bullet \ \mathcal{H} = \{ f_{\boldsymbol{w}} : \mathbf{x} \mapsto \boldsymbol{w}^{\top} \mathbf{x} \mid \boldsymbol{w} \in \mathbb{R}^p \}$
- $||f_{w}||_{\mathcal{H}} = ||w||_{2}$

Polynomial kernel

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

Radial Basis Function kernel (RBF)

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

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

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

Math-412 Kernel methods 16/18

Representer theorem

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

$$\mathit{f}_{\perp}(x_i) = \langle \mathit{f}_{\perp}, \mathit{K}(x_i, \cdot) \rangle_{\mathcal{H}} = 0 \quad \text{and} \quad \langle \mathit{f}_{\perp}, \mathit{f}_{/\!\!/} \rangle_{\mathcal{H}} = 0.$$

Thus

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

$$= L(f_{/\!/}(x_1), ..., f_{/\!/}(x_n)) + \lambda (||f_{/\!/}||_{\mathcal{H}}^2 + 2\langle f_{\perp}, f_{/\!/}\rangle_{\mathcal{H}} + ||f_{\perp}||_{\mathcal{H}}^2)$$

$$= L(f_{/\!/}(x_1), ..., f_{/\!/}(x_n)) + \lambda ||f_{/\!/}||_{\mathcal{H}}^2 + \lambda ||f_{\perp}||_{\mathcal{H}}^2$$

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

Math-412 Kernel methods 17,

Regularized ERM for f in a RKHS

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

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

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

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

Math-412 Kernel methods 18/18