# Lecture 10: Kernel Methods
# CS480/680 Intro to Machine Learning

2023-2-9

Pascal Poupart
David R. Cheriton School of Computer Science

UNIVERSITY OF
**WATERLOO**

# Non-linear Models Recap

- Generalized linear models:

  − fixed non-linear basis functions
  − limited hypothesis space
  − easy to optimize (usually convex)

- Neural networks:

  − adaptive non-linear basis functions
  − rich hypothesis space
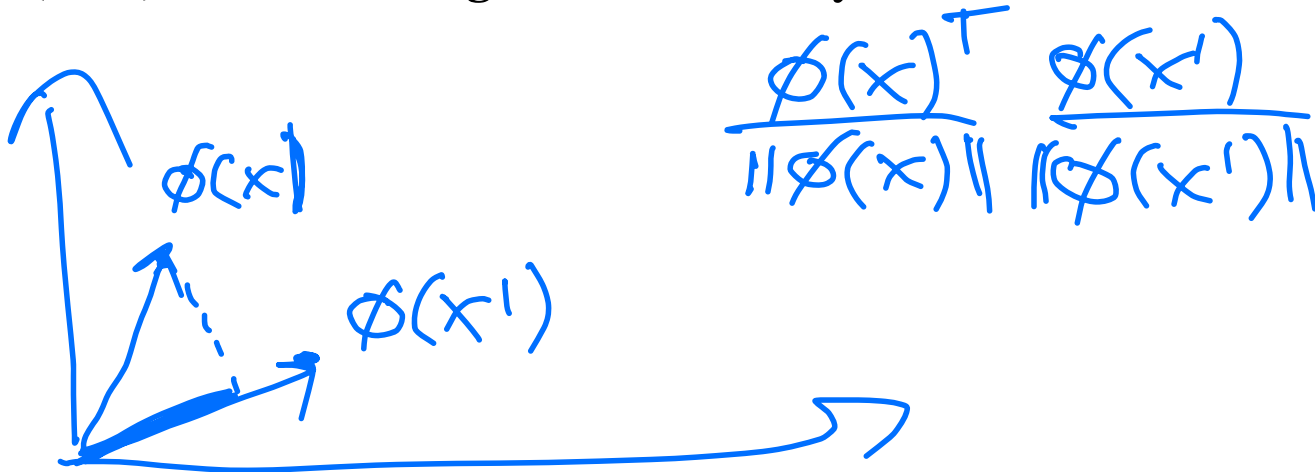  − hard to optimize (usually non-convex)

# Kernel Methods

- Idea: use large (possibly infinite) set of fixed non-linear basis functions
- Normally, complexity depends on number of basis functions, but by a "dual trick", **complexity depends on the amount of data**

- Examples:
  - **Gaussian Processes** (next class)
  - **Support Vector Machines** (next week)
  - Kernel perceptron
  - Kernel logistic regression

UNIVERSITY OF
**WATERLOO**

# Kernel Function

- Let $\phi(x)$ be a set of basis functions that map inputs $x$ to a feature space.

- In many algorithms, this feature space only appears in the dot product $\phi(x)^T \phi(x')$ of input pairs $x, x'$.

- Define the kernel function $k(x, x') = \phi(x)^T \phi(x')$ to be the dot product of any pair $x, x'$ in feature space.

    - **We only need to know $k(x, x')$**, not $\phi(x)$

UNIVERSITY OF
WATERLOO

# Illustration of Kernel Function

- $k(\boldsymbol{x}, \boldsymbol{x}') = \phi(\boldsymbol{x})^T \phi(\boldsymbol{x}')$

- Intuition: $k(\boldsymbol{x}, \boldsymbol{x}')$ measures degree of similarity

$$\frac{\phi(x)^T \; \phi(x')}{\|\phi(x)\| \; \|\phi(x')\|}$$

$\phi(x)$

$\phi(x')$

UNIVERSITY OF
WATERLOO

# Dual Representations

- Recall linear regression objective

$$E(\boldsymbol{w}) = \frac{1}{2}\sum_{n=1}^{N}[\boldsymbol{w}^T\phi(\boldsymbol{x}_n) - y_n]^2 + \frac{\lambda}{2}\boldsymbol{w}^T\boldsymbol{w}$$

- Solution: set gradient to 0

$$\nabla E(\boldsymbol{w}) = \sum_n(\boldsymbol{w}^T\phi(\boldsymbol{x}_n) - y_n)\phi(\boldsymbol{x}_n) + \lambda\boldsymbol{w} = 0$$

$$\boldsymbol{w} = -\frac{1}{\lambda}\sum_n(\boldsymbol{w}^T\phi(\boldsymbol{x}_n) - y_n)\phi(\boldsymbol{x}_n)$$

vector        scalar      vector

$\therefore$ **$w$ is a linear combination of inputs in feature space**

$$\{\phi(\boldsymbol{x}_n)\,|\,1 \le n \le N\}$$

UNIVERSITY OF
WATERLOO

# Dual Representations

- Substitute $\mathbf{w} = \mathbf{\Phi a}$

- Where $\mathbf{\Phi} = [\phi(\mathbf{x}_1)\ \phi(\mathbf{x}_2)\ \ldots\ \phi(\mathbf{x}_N)]$

$$\mathbf{a} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_N \end{bmatrix} \quad \text{and } a_n = -\frac{1}{\lambda}(\mathbf{w}^T\phi(\mathbf{x}_n) - y_n)$$

- Dual objective: minimize $E$ with respect to $\mathbf{a}$

$$E(\mathbf{a}) = \frac{1}{2}\mathbf{a}^T\mathbf{\Phi}^T\mathbf{\Phi}\mathbf{\Phi}^T\mathbf{\Phi}\mathbf{a} - \mathbf{a}^T\mathbf{\Phi}^T\mathbf{\Phi}\mathbf{y} + \frac{\mathbf{y}^T\mathbf{y}}{2} + \frac{\lambda}{2}\mathbf{a}^T\mathbf{\Phi}^T\mathbf{\Phi}\mathbf{a}$$

UNIVERSITY OF
WATERLOO

# Gram Matrix

- Let $K = \Phi^T\Phi$ be the Gram matrix

- Substitute in objective:

$$E(a) = \frac{1}{2}a^T K K a - a^T K y + \frac{y^T y}{2} + \frac{\lambda}{2}a^T K a$$

$k(x_1, x_1) = \phi(x_1)^T\phi(x_1)$

$k(x_1, x_2) = \phi(x_1)^T\phi(x_2)$

- Solution: set gradient to 0

$$\nabla E(a) = K K a - K y + \lambda K a = 0$$

$$K(K + \lambda I)a = Ky$$

$$a = (K + \lambda I)^{-1}y$$

- Prediction:

*row vector*

$$y_* = \phi(x_*)^T w = \phi(x_*)^T \Phi a = k(x_*, X)(K + \lambda I)^{-1}y$$

where $(X, y)$ is the training set and $(x_*, y_*)$ is a test instance

UNIVERSITY OF
WATERLOO

# Dual Linear Regression

- Prediction: $y_* = \phi(\boldsymbol{x}_*)^T \boldsymbol{\Phi} \boldsymbol{a}$

$$= k(\boldsymbol{x}_*, \boldsymbol{X})(\boldsymbol{K} + \lambda \boldsymbol{I})^{-1} \boldsymbol{y}$$

- Linear regression where we find dual solution $\boldsymbol{a}$ instead of primal solution **w**.

- Complexity:

  - Primal solution: depends on # of basis functions

  - Dual solution: depends on amount of data

    - Advantage: can use very large # of basis functions

    - Just need to know kernel $k$

UNIVERSITY OF
**WATERLOO**

# Constructing Kernels

- Two possibilities:
  - Find mapping $\boldsymbol{\phi}$ to feature space and let $\boldsymbol{K} = \boldsymbol{\phi}^T\boldsymbol{\phi}$
  - Directly specify $\boldsymbol{K}$

- Can any function that takes two arguments serve as a kernel?

- No, a valid kernel must be positive semi-definite
  - In other words, $k$ must factor into the product of a transposed matrix by itself (e.g., $\boldsymbol{K} = \boldsymbol{\phi}^T\boldsymbol{\phi}$)
  - Or all eigenvalues must be greater than or equal to 0.

UNIVERSITY OF
**WATERLOO**

# Example

- Let $k(x, z) = (x^T z)^2$

$$x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \qquad z = \begin{pmatrix} z_1 \\ z_2 \end{pmatrix}$$

$$= (x_1 z_1 + x_2 z_2)^2$$

$$= x_1^2 z_1^2 + 2 x_1 z_1 x_2 z_2 + x_2^2 z_2^2$$

$$= (x_1^2, \sqrt{2} x_1 x_2, x_2^2) \begin{pmatrix} z_1^2 \\ \sqrt{2} z_1 z_2 \\ z_2^2 \end{pmatrix}$$

$$= \phi(x)^T \phi(z)$$

$$\hookrightarrow \phi_1(x) = x_1^2$$
$$\phi_2(x) = \sqrt{2} x_1 x_2$$
$$\phi_3(x) = x_2^2$$

UNIVERSITY OF
WATERLOO

# Constructing Kernels

- Can we construct $k$ directly without knowing $\phi$?

- Yes, any positive semi-definite $k$ is fine since there is a corresponding implicit feature space. But positive semi-definiteness is not always easy to verify.

- Alternative, construct kernels from other kernels using rules that preserve positive semi-definiteness

UNIVERSITY OF
WATERLOO

# Rules to construct Kernels

- Let $k_1(\boldsymbol{x}, \boldsymbol{x}')$ and $k_2(\boldsymbol{x}, \boldsymbol{x}')$ be valid kernels
- The following kernels are also valid:
  1. $k(\boldsymbol{x}, \boldsymbol{x}') = c k_1(\boldsymbol{x}, \boldsymbol{x}') \quad \forall c > 0$
  2. $k(\boldsymbol{x}, \boldsymbol{x}') = f(\boldsymbol{x}) k_1(\boldsymbol{x}, \boldsymbol{x}') f(\boldsymbol{x}') \quad \forall f$
  3. $k(\boldsymbol{x}, \boldsymbol{x}') = q(k_1(\boldsymbol{x}, \boldsymbol{x}')) \quad q$ is polynomial with coeffs $\geq 0$
  4. $k(\boldsymbol{x}, \boldsymbol{x}') = \exp(k_1(\boldsymbol{x}, \boldsymbol{x}'))$
  5. $k(\boldsymbol{x}, \boldsymbol{x}') = k_1(\boldsymbol{x}, \boldsymbol{x}') + k_2(\boldsymbol{x}, \boldsymbol{x}')$
  6. $k(\boldsymbol{x}, \boldsymbol{x}') = k_1(\boldsymbol{x}, \boldsymbol{x}') k_2(\boldsymbol{x}, \boldsymbol{x}')$
  7. $k(\boldsymbol{x}, \boldsymbol{x}') = k_3(\phi(\boldsymbol{x}), \phi(\boldsymbol{x}'))$
  8. $k(\boldsymbol{x}, \boldsymbol{x}') = \boldsymbol{x}^T \boldsymbol{A} \boldsymbol{x}' \quad \boldsymbol{A}$ is symmetric positive semi-definite
  9. $k(\boldsymbol{x}, \boldsymbol{x}') = k_a(\boldsymbol{x_a}, \boldsymbol{x}_a') + k_b(\boldsymbol{x_b}, \boldsymbol{x}_b')$
  10. $k(\boldsymbol{x}, \boldsymbol{x}') = k_a(\boldsymbol{x_a}, \boldsymbol{x}_a') k_b(\boldsymbol{x_b}, \boldsymbol{x}_b')$

where $\boldsymbol{x} = \begin{pmatrix} \boldsymbol{x}_a \\ \boldsymbol{x}_b \end{pmatrix}$

UNIVERSITY OF
WATERLOO

# Common Kernels

- Polynomial kernel: $k(\boldsymbol{x}, \boldsymbol{x}') = (\boldsymbol{x}^T \boldsymbol{x}')^M$

  - $M$ is the degree

  - Feature space: all degree M products of entries in $\boldsymbol{x}$

  - Example: Let $\boldsymbol{x}$ and $\boldsymbol{x}'$ be two images, then feature space could be all products of M pixel intensities

- More general polynomial kernel:
$$k(\boldsymbol{x}, \boldsymbol{x}') = (\boldsymbol{x}^T \boldsymbol{x}' + c)^M \ \text{ with } c > 0$$

  - Feature space: all products of up to M entries in $\boldsymbol{x}$

UNIVERSITY OF
WATERLOO

# Example

$x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \qquad x' = \begin{pmatrix} x_1' \\ x_2' \end{pmatrix}$

- $k(x, x') = \left(x^T x' + c\right)^2$

$$= \left(x_1 x_1' + x_2 x_2' + c\right)^2$$

$$= x_1^2 x_1'^2 + 2x_1 x_1' x_2 x_2' + x_2^2 x_2'^2 + 2x_1 x_1' c + 2x_2 x_2' c$$
$$+ c^2$$

$$= \left(x_1^2, \sqrt{2} x_1 x_2, x_2^2, \sqrt{2c}\, x_1, \sqrt{2c}\, x_2, c\right)$$
$$\left(x_1'^2, \sqrt{2} x_1' x_2', x_2'^2, \sqrt{2c}\, x_1', \sqrt{2c}\, x_2', c\right)^T$$

# Common Kernels

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

- Valid Kernel because:

$$= e^{-\boldsymbol{x}^T\boldsymbol{x}/2\sigma^2} \; e^{\boldsymbol{x}^T\boldsymbol{x}'/\sigma^2} \; e^{-\boldsymbol{x}'^T\boldsymbol{x}'/2\sigma^2}$$

$\boldsymbol{x}^T\boldsymbol{x}'$ is a valid kernel by rule $8$ when $A = I$

$\boldsymbol{x}^T\boldsymbol{x}'/\sigma^2$   ''   ''     ''     ''     ''     $1$

$e^{\boldsymbol{x}^T\boldsymbol{x}'/\sigma^2}$   ''   ''     ''     ''     ''     $4$

$k(\boldsymbol{x}, \boldsymbol{x}')$   ''   ''     ''     ''     ''     $2$

- Implicit feature space is infinite!

UNIVERSITY OF
WATERLOO

# Non-vectorial Kernels

- Kernels can be defined with respect to other things than vectors such as sets, strings or graphs

- Example for strings: $k(d_1, d_2)$ = similarity between two documents (weighted sum of all non-contiguous strings that appear in both documents $d_1$ and $d_2$).

- Lodhi, Saunders, Shawe-Taylor, Christianini, Watkins, **Text Classification Using String Kernels**, JMLR, p. 419-444, 2002.

UNIVERSITY OF
**WATERLOO**