

INTRODUCTION TO
MACHINE LEARNING
COMPSCI 4ML3

LECTURE 7

HASSAN ASHTIANI

CALCULATING OLS WITH FEATURE MAPS

- FEATURE MAP: $\phi(x): \mathbb{R}^{d_1} \mapsto \mathbb{R}^{d_2}$

$$X_{n \times d_1} \rightarrow \Phi_{n \times d_2}$$

- TO CALCULATE: $W^* = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T Y$

- NEED TO INVERT A $d_2 \times d_2$ MATRIX

- d_2 CAN BE VERY LARGE, AND EVEN INFINITE!

- KERNEL TRICK: COMPUTE THE HIGH-DIMENSIONAL INNER PRODUCT EFFICIENTLY

- $K(x^i, x^j) = \langle \phi(x^i), \phi(x^j) \rangle$

- USE THIS AS A BUILDING-BLOCK FOR PERFORMING OTHER OPERATIONS

- REWRITE THE LEAST SQUARES SOLUTION SO THAT IT ONLY USES INNER PRODUCTS IN THE FEATURE MAP!

ROADMAP

- OPTIMAL OLS $W^* = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T Y$
- ASSUME d_2 IS VERY LARGE, EVEN $d_2 \gg n$
- INSTEAD OF FINDING w , TRY TO INTRODUCE NEW PARAMETER a WHOSE SIZE IS n RATHER THAN d_2
- NOW WE HAVE n PARAMETERS
- FIND OPTIMAL a AS A FUNCTION OF K

KERNELIZED LEAST SQUARES

- $W^* = \min_W \|\Phi W - Y\|_2^2 + \lambda \|W\|_2^2$

- STEP 1: SHOW THERE EXISTS $a \in \mathbb{R}^n$, SUCH THAT $W^* = \Phi^T a$

- IN OTHER WORDS, $W^* = \sum a_i \phi(x^i)$

- NUMBER OF PARAMETERS?

- n INSTEAD OF $d_2 \dots$

- PROOF?

$\phi(x^i) \in \mathbb{R}^{d_2}$
 $W^* \in \mathbb{R}^{d_2}$

$$0 = \frac{\partial}{\partial w} \left(\underbrace{\|\Phi w - Y\|_2^2}_{\text{---}} + \underbrace{\lambda \|w\|_2^2}_{\text{~~~~~}} \right)$$

$$= 2w^T (\Phi^T \Phi) - 2Y^T \Phi + 2\lambda w^T = 0$$

$$\Rightarrow \lambda w = \Phi^T Y - \Phi^T \Phi w$$

$$\Rightarrow \underbrace{w}_{\text{~~~~~}} = \Phi^T \left(\underbrace{\frac{Y - \Phi w}{\lambda}}_a \right) = \sum a_i \phi(x_i)$$

$$(AB)^T = B^T A^T, \quad (A+B)^T = A^T + B^T, \quad \dots$$

KERNEL FUNCTION NOTATIONS

→ • $k(x^i, x^j) = \langle \phi(x^i), \phi(x^j) \rangle$

- KERNEL OR GRAM MATRIX OF A DATA SET:

→ • $K_{n \times n} = [k(x^i, x^j)] = \Phi \Phi^T$



• $k(x) = \Phi \phi(x) = [k(x, x^1) \quad k(x, x^2) \quad \dots \quad k(x, x^n)]^T$

PREDICTION, GIVEN a

- $W^* = \Phi^T a$
- PREDICTION ON TRAINING POINTS

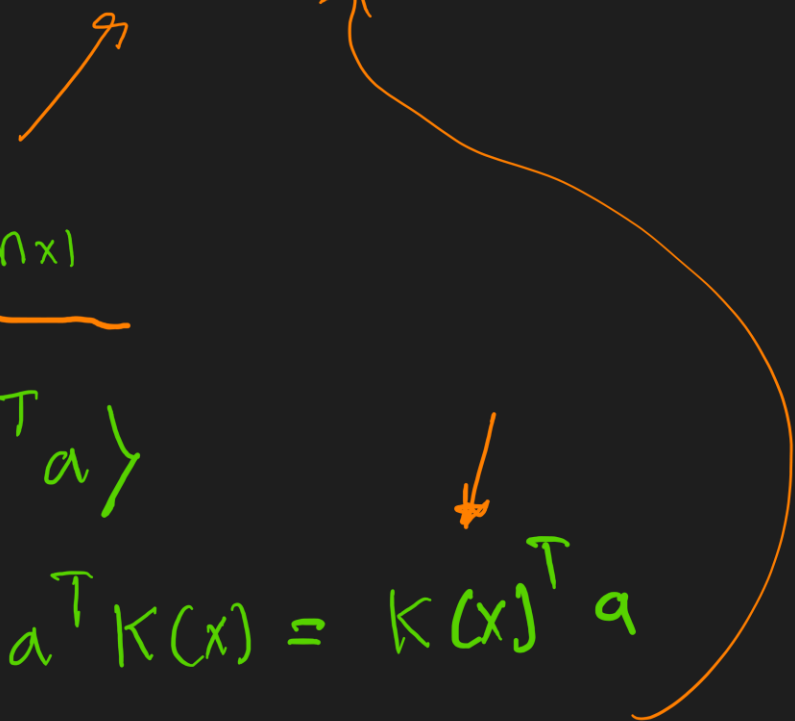
• $\hat{Y} = \Phi W^* = ?$ $\Phi \Phi^T a = K \cdot a$
 $n \times n$ $n \times 1$

- PREDICTION FOR NEW TEST POINT x :

• $\hat{y}(x) = \langle \phi(x), W^* \rangle = ?$ $\langle \phi(x), \Phi^T a \rangle$
 $= a^T \Phi \phi(x) = a^T K(x) = K(x)^T a$
 $= \langle a, K(x) \rangle$

$K(x) = [K(x, x^1), K(x, x^2), \dots]$

nice when $d_2 \gg n$



FINDING a USING DUAL FORM

- $W^* = \min_W \|\Phi W - Y\|_2^2 + \lambda \|W\|_2^2$
- STEP 2: USE $W^* = \Phi^T a$ TO REFORMULATE THE PROBLEM IN TERMS OF FINDING a (DUAL FORM)

- $\min_{\vec{a} \in \mathbb{R}^n} \|\Phi \Phi^T a - Y\|_2^2 + \lambda \|\Phi^T a\|_2^2$ OR...

$$\rightarrow \min_{\vec{a}} \|K a - Y\|_2^2 + \lambda a^T K a \leftarrow$$

- $a^* = (K + \lambda I)^{-1} Y$ (PROOF?)

- FASTER WHEN $d_2 \gg n$

CHOICE OF KERNEL

- KERNEL ENCODES SIMILARITY OF POINTS x^i AND x^j

- POLYNOMIAL: $k(x, z) = (1 + x^T z)^M$

- GAUSSIAN: $k(x, z) = e^{-\left(\frac{1}{2\sigma^2}\right)\|x-z\|_2^2} = e^{-\alpha\|x-z\|_2^2}$

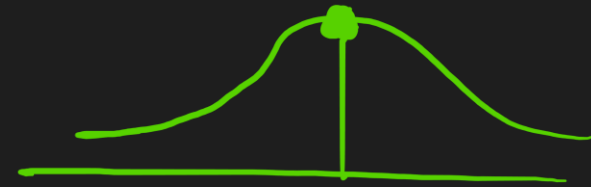
- $\phi(x)$ IS INFINITE DIMENSIONAL

- HOW TO CHOOSE A KERNEL?

- IT SHOULD BE VALID (THERE MUST EXIST A ϕ)

- DOMAIN KNOWLEDGE

- KERNEL FUNCTION CAPTURES "SIMILARITY" BETWEEN POINTS



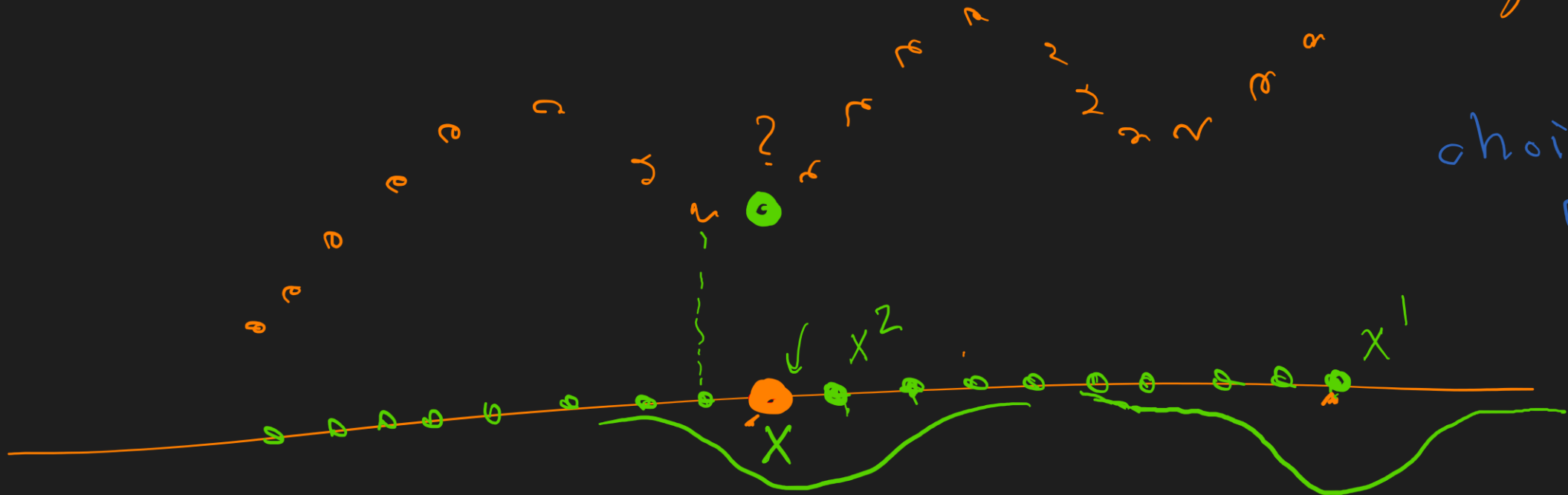
GAUSSIAN KERNEL: INTUITION

- JUPYTER NOTEBOOK

$x \in \mathbb{R}$

$y \in \mathbb{R}$

choice of σ



$\rightarrow \sum a_i \underbrace{K(x, x^i)}$, $K(x, x^1) \approx 0$
 $K(x, x^2) \gg 0$

COMPUTATIONAL COMPLEXITY

- MATRIX MULTIPLICATION (N-BY-N MATRICES)
 - NATIVE METHOD: $O(N^3)$
 - STRASSEN'S ALGORITHM: $O(N^{2.8074})$
 - CURRENTLY BEST KNOWN METHOD: COPPERSMITH–WINOGRAD ALGORITHM $O(N^{2.3755})$
- MATRIX INVERSION
 - GAUSSIAN ELIMINATION: $O(N^3)$
 - POSSIBLE TO REDUCE IT TO MULTIPLICATION (SO $O(N^{2.3755})$)

COMPUTATIONAL COMPLEXITY

• TRAINING COMPLEXITY (n TRAINING POINTS)

• REGULARIZED LEAST SQUARES

• $W = (X^T X + \lambda I)^{-1} X^T Y$ ✓

$d_2 \times d_2$ matrix \rightarrow

• KERNEL LEAST SQUARES

• $a = (K + \lambda I)^{-1} Y$ ✓

$n \times n$ matrix \rightarrow

• TEST COMPLEXITY (FOR A SINGLE TEST POINT)

• REGULARIZED LEAST SQUARES

• $x^T W$ ✓

$\rightarrow d_2$ computation

• KERNEL LEAST SQUARES

• $k(x)^T a$ ✓

$\rightarrow n$ computations