

INTRODUCTION TO MACHINE LEARNING COMPSCI 4ML3

LECTURE 5

HASSAN ASHTIANI

THE TRADE-OFF

- A POWERFUL/FLEXIBLE CURVE-FITTING METHOD
 - SMALL TRAINING ERROR
 - REQUIRES MORE TRAINING DATA TO GENERALIZE
 - OTHERWISE LARGE TEST ERROR
- A LESS FLEXIBLE CURVE-FITTING METHOD
 - LARGER TRAINING ERROR
 - REQUIRES LESS TRAINING DATA
 - SMALLER DIFFERENCE BETWEEN TRAINING AND TEST ERROR
- THE SO-CALLED “BIAS-VARIANCE” TRADE-OFF

THE CASE OF MULTIVARIATE POLYNOMIALS

- ASSUME $M \gg d$
- NUMBER OF TERMS (MONOMIALS): $\approx (\frac{M}{d})^d$
- #TRAINING SAMPLES \approx #PARAMETERS $\approx (\frac{M}{d})^d$
 - #TRAINING SAMPLES SHOULD INCREASE EXPONENTIALLY WITH d
 - SUSCEPTIBLE TO OVER-FITTING...
 - AN EXAMPLE OF **CURSE OF DIMENSIONALITY!**
- WE CAN SAY **SAMPLE COMPLEXITY** OF LEARNING MULTIVARIATE POLYNOMIALS IS EXPONENTIAL IN d
 - ORTHOGONAL TO COMPUTATIONAL COMPLEXITY

MODEL SELECTION: HOW TO AVOID OVERFITTING?

- SELECTING M (THE COMPLEXITY OF THE MODEL)
 - BASED ON d (DIMENSION) AND n (NUMBER OF SAMPLES)
- MORE PRACTICALLY, TRY SEVERAL OPTIONS FOR M
 - USE A **HOLDOUT (EVALUATION) SAMPLE**
 - NEVER USE TEST DATA TO TUNE PARAMETERS!

AVOID OVERFITTING WITH REGULARIZED LEAST SQUARES

$$\min_{W \in \mathcal{R}^d} \|XW - Y\|_2^2 + \lambda \|W\|_2^2$$

- ENCOURAGE A SOLUTION WITH A SMALLER NORM
- $W^{RLS} = (X^T X + \lambda I)^{-1} X^T Y$
 - EXERCISE: PROVE THAT THIS IS THE OPTIMAL SOLUTION
- DOES THE INVERSE ALWAYS EXIST?
 - YES! (EXERCISE: PROVE)
- HOW TO CHOOSE λ ?

POLYNOMIAL CURVE-FITTING REVISITED

- MAP THE INPUTS x^i TO A HIGHER DIMENSIONAL SPACE
 - A KIND OF “PRE-PROCESSING” THE DATA
- DO LINEAR REGRESSION ON THE HIGH-DIMENSIONAL SPACE
 - EQUIVALENT TO PERFORMING NON-LINEAR REGRESSION IN THE ORIGINAL SPACE
- MAP $\phi(x): \mathbb{R}^{d_1} \mapsto \mathbb{R}^{d_2}$ WHERE $d_2 \gg d_1$

- $\phi(x) = \begin{pmatrix} \phi_1(x) \\ \dots \\ \phi_{d_2}(x) \end{pmatrix}$ IS NONLINEAR, E.G., $x \in \mathbb{R}$ AND $\phi(x) = \begin{pmatrix} x \\ x^2 \\ x^3 \\ \dots \\ x^{d_2} \end{pmatrix}$
- WHAT IF d_2 IS MUCH LARGER THAN THE NUMBER OF SAMPLES?

CURVE-FITTING WITH BASIS FUNCTIONS

- FEATURE MAP: $\phi(x): \mathbb{R}^{d_1} \mapsto \mathbb{R}^{d_2} \quad d_2 \gg d_1$
- $\Phi_{n \times d_2} = [\phi(x^1) \quad \dots \quad \phi(x^n)]^T$
- TRAINING
 - $W^* = \min_W \|\Phi W - Y\|_2^2 + \lambda \|W\|_2^2$
 - $W^* = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T Y$
- PREDICTION
 - $\hat{y} = \langle W^*, \phi(x) \rangle = W^{*T} \phi(x)$

OTHER CHOICES OF $\phi(x)$

- PICK A FIXED (NONLINEAR) $\Phi(x)$
 - ENCODES YOUR PRIOR KNOWLEDGE ABOUT THE DATA
 - FEATURE ENGINEERING!
- POLYNOMIAL BASIS FUNCTIONS
- GAUSSIAN BASIS FUNCTIONS:
 - $\phi_i(x) = e^{-\frac{\|x-\mu_i\|_2^2}{2\sigma^2}}$
- DFT (FFT), WAVELET FOR TIME SERIES
- IS IT POSSIBLE TO LEARN THE MAPPING $\phi_i(x)$ ITSELF?
 - LATER, E.G., NEURAL NETWORKS

|

COMPUTATIONAL COMPLEXITY OF NAÏVE RLS

- TRAINING: CALCULATE $w^{RLS} = (\phi^T \phi + \lambda I)^{-1} \phi^T Y$
- BOTTLENECK: MATRIX INVERSION
 - HOW MANY OPERATIONS?
- PREDICTION: $\hat{y} = \langle \phi(x), w^{RLS} \rangle$
 - HOW MANY OPERATIONS?
- REGULARIZATION ALLOWS US TO GO INTO HIGH-DIMENSIONAL SPACE WITHOUT OVERRFITTING, BUT IT DOES NOT SOLVE THE COMPUTATIONAL PROBLEM

COMPUTATIONAL COMPLEXITY

- MATRIX MULTIPLICATION (N-BY-N MATRICES)
 - NATIVE METHOD: $O(N^3)$
 - STRASSEN'S ALGORITHM: $O(N^{2.8074})$
 - COPPERSMITH–WINOGRAD-LIKE ALGORITHMS [CURRENT BEST
 $O(N^{2.3728639})$]
- MATRIX INVERSION
 - GAUSSIAN ELIMINATION: $O(N^3)$
 - POSSIBLE TO REDUCE IT TO MULTIPLICATION

THE COMPUTATIONAL PROBLEM

- CAN WE SOLVE THE REGULARIZED LEAST SQUARES IN \mathbb{R}^{d_2} WITHOUT EXPLICITLY MAPPING THE DATA INTO \mathbb{R}^{d_2} ?
 - $W^* = \min_{W \in \mathbb{R}^{d_2}} \|\Phi W - Y\|_2^2 + \lambda \|W\|_2^2$
- SOMETHING LIKE MULTIPLICATION USING FFT
- IF SO, WE COULD EVEN MAP THE DATA TO AN INFINITE DIMENSIONAL SPACE!!

FFT AND MULTIPLICATION

