

KERNEL METHODS

MOTIVATING EXAMPLE: BINARY CLASSIFICATION: SUPPORT VECTOR MACHINES (IZENMAN Ch. 11)

SETTING: GIVEN DATA $\{(\vec{x}_i, y_i)\}_{i=1}^N$, $\vec{x}_i \in \mathbb{R}^d$, $y_i \in \{-1, 1\}$

GOAL: LEARN A FUNCTION $f: \mathbb{R}^d \rightarrow \mathbb{R}$ FROM THE GIVEN DATA SO THAT

$$c(\vec{x}) = \text{sgn}(f(\vec{x})) = \begin{cases} +1 & f(\vec{x}) \geq 0 \\ -1 & f(\vec{x}) < 0 \end{cases}$$

IS A CLASSIFIER, I.E. GIVEN NEW \vec{x}_{N+1} WE PREDICT $y_{N+1} = c(\vec{x}_{N+1})$.

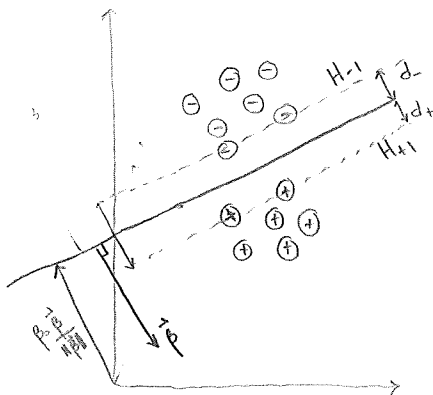
LINEARLY SEPARABLE CASE (SIMPLEST CASE)

ASSUME POSITIVE ($y_i = 1$) AND NEGATIVE ($y_i = -1$) POINTS

CAN BE PERFECTLY SEPARATED BY A HYPERPLANE

$$\{\vec{x} : \beta_0 + \vec{\beta}^T \vec{x} = 0\}$$

↑ CALLED A SEPARATING HYPERPLANE
(DECISION BOUNDARY)



* THE EXISTENCE OF A SEPARATING HYPERPLANE
IMPLIES EXISTENCE OF $\beta_0, \vec{\beta}$ S.T.

$$\beta_0 + \vec{\beta}^T \vec{x}_i \geq +1 \quad (y_i = +1); \quad \beta_0 + \vec{\beta}^T \vec{x}_i \leq -1 \quad (y_i = -1)$$

* IF d_+ (d_-) IS THE SMALLEST DISTANCE FROM
HYPERPLANE TO ^{NEAREST} NEGATIVE (POSITIVE) POINT

$$\text{MARGIN} = d_- + d_+$$

$$= \frac{2}{\|\vec{\beta}\|} \quad \leftarrow \text{MAXIMUM MARGIN GIVES OPTIMAL SEPARATING HYPERPLANE}$$

$$\text{MINIMIZE} \quad \frac{1}{2} \|\vec{\beta}\|^2$$

$$\text{SUBJECT TO} \quad y_i (\beta_0 + \vec{\beta}^T \vec{x}_i) \geq 1 \quad i=1, 2, \dots, N$$

PRIMAL DUAL OPTIMIZATION (KARUSH-KUHN-TUCKER CONDITIONS)

IT CAN BE SHOWN THAT THE OPTIMAL $\hat{\beta}_0, \hat{\beta}$ ARE GIVEN BY

$$\hat{\beta} = \sum_{i=1}^N \hat{\alpha}_i y_i \vec{x}_i$$

$$\hat{\beta}_0 = -\frac{1}{2} \left(\vec{x}_{-1}^T \hat{\beta} + \vec{x}_{+1}^T \hat{\beta} \right) \quad x_{-1} \text{ or } H_{-1}, x_{+1} \text{ or } H_{+1}$$

WHERE $\vec{\alpha} = (\alpha_1, \dots, \alpha_N)$ MAXIMIZES

$$\vec{1}_N^T \vec{\alpha} - \frac{1}{2} \vec{\alpha}^T H \vec{\alpha}$$

SUBJECT TO $\alpha_1, \dots, \alpha_N \geq 0$, $\sum \alpha_i y_i = 0$, $H_{ij} = y_i y_j \vec{x}_i^T \vec{x}_j$,

AND THE OPTIMAL HYPERTHANE IS

$$f(\vec{x}) = \hat{\beta}_0 + \sum_{i=1}^N \hat{\alpha}_i y_i (\vec{x}_i^T \vec{x})$$

NOTE: • TYPICALLY, $\vec{\alpha}$ IS SPARSE (MANY ZEROES) SO $\sum_{i=1}^N \dots$ CAN BE REDUCED TO SUPPORT VECTORS FOR WHICH $\hat{\alpha}_i > 0$.

• CLASSIFICATION OF \vec{x}_{N+1} REQUIRES INNER PRODUCTS $\vec{x}_i^T \vec{x}_N$

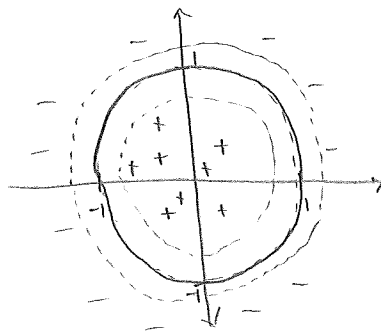
$$C(\vec{x}_{N+1}) = \text{SIGN} \left(f(\vec{x}_{N+1}) \right)$$

SO WE CAN DERIVE A CLASSIFICATION RULE IF WE HAVE ONLY THE INNER PRODUCTS OF ANY PAIR OF POINTS BUT NOT THE POINTS THEMSELVES (SIMILAR TO INNER PRODUCT MATRIX; CLASSICAL SCALING MDS)!

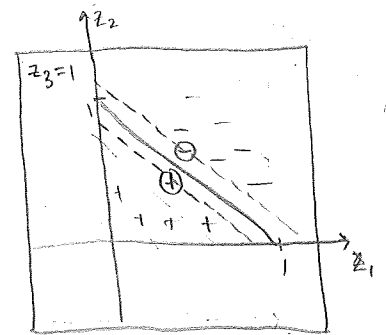
• LINEARLY SEPARABLE (HARD MARGIN) CASE CAN BE RELATED TO SOFT-MARGIN (STRICT SEPARATION NOT REQUIRED) w/ CONNECTIONS TO LOGISTIC REGRESSION.

Nonlinear SVM

Ex: Non-Plane Decision Boundary



$$\varphi \left((x_1, x_2) \right) = \begin{pmatrix} x_1^2 \\ x_2^2 \\ 1 \end{pmatrix} = \begin{pmatrix} z_1 \\ z_2 \\ z_3 \end{pmatrix}$$



* A CURVED DECISION BOUNDARY CAN BE REPRESENTED IMPLICITLY AS A FUNCTION (POLYNOMIAL IN SOME CASES) OF COMPONENTS OF \vec{x} , i.e. $x_1^2 + x_2^2 - 1 = 0$ IN CASE ABOVE

* CAN FIND A NONLINEAR MAPPING FROM ORIGINAL SPACE (\mathbb{R}^2 HERE) TO A HIGHER DIMENSIONAL "FEATURE" SPACE (HERE \mathbb{R}^3) WHERE DATA ARE LINEARLY SEPARABLE!

$$\varphi: \mathbb{R}^d \xrightarrow{\text{WHERE OUR DATA LIE}} \mathcal{H} \subset \mathbb{R}^p$$

↑ ↑
WHERE OUR DATA LIE FEATURE SPACE

* LINEAR SVM IN \mathcal{H}

• DECISION BOUNDARY IS $f(\vec{x}) = \beta_0 + \beta^T \varphi(x)$

$$= \beta_0 + \sum_{i=1}^N \alpha_i y_i \underbrace{\varphi(x_i)^T \varphi(\vec{x})}_{\text{REQUIRES AN INNER PRODUCT ON } \mathcal{H}!!}$$

• FINDING $\vec{\alpha}$ REQUIRES MAXIMIZING

$$\vec{1}^T \vec{\alpha} - \vec{\alpha}^T H \vec{\alpha}, \quad H_{ij} = y_i y_j \varphi(x_i)^T \varphi(x_j)$$

Ex: 2^{nd} degree polynomial in \mathbb{R}^d

- GIVEN $\vec{x} = (x_1 \dots x_d)^T$ WE WANT TO FIND THE DECISION BOUNDARY OF THE FORM

$$f(\vec{x}) = \beta_0 + \vec{p}^T \varphi(x), \quad \varphi(x) = \begin{pmatrix} x_1^2 \\ \vdots \\ x_d^2 \\ \sqrt{2} x_d x_{d-1} \\ \sqrt{2} x_d x_1 \\ \sqrt{2} x_{d-1} x_{d-2} \\ \vdots \\ \sqrt{2} x_d x_1 \\ \sqrt{2c} x_d \\ \vdots \\ \sqrt{2c} x_1 \end{pmatrix} \in \mathbb{R}^{\frac{(d+1)(d+2)}{2}}$$

- DIMENSION OF $\varphi(x)$ IS QUADRATIC IN d !
- EVEN WORSE IF WE INCLUDE CUBIC, QUARTIC, ... TERMS
- EXPONENTIALLY BAD AS WE INCREASE COMPLEXITY, FROM A COMPUTATIONAL PERSPECTIVE

• OBSERVE:

$$\begin{aligned} \varphi(x)^T \varphi(y) &= \sum_{i=1}^d (x_i^2)(y_i^2) + \sum_{i=2}^d \sum_{j=1}^{i-1} 2x_i x_j y_i y_j + \sum_{i=1}^d 2c x_i y_i + c \\ &= \sum_{i=1}^d (x_i y_i)^2 + 2 \sum_{i=1}^d \sum_{j=1}^{i-1} (x_i y_i)(x_j y_j) + 2c \sum_{i=1}^d x_i y_i + c^2 \\ &= (\vec{x}^T \vec{y} + c)^2 = k(\vec{x}, \vec{y}) \end{aligned}$$

* NO NEED TO COMPUTE $\varphi(x)$ OR $\varphi(y)$ IF WE ONLY NEED $\varphi(x)^T \varphi(y)$!

- COMPUTE INNER PRODUCTS IN FEATURE SPACE USING ONLY ORIGINAL VECTORS \vec{x}, \vec{y} !

"KERNEL TRICK": $k(\vec{x}, \vec{y}) = \varphi(\vec{x})^T \varphi(\vec{y})$

* HIGHER ORDER POLYNOMIALS USING $k(\vec{x}, \vec{y}) = (\vec{x}^T \vec{y} + c)^d$

↑
POLYNOMIAL KERNEL

- LINEAR METHODS REDUCING INNER PRODUCTS CAN BE "KERNELIZED"
TO GIVE A NONLINEAR METHOD w/o EVER COMPUTING $\phi(x_i)$!

Ex: RADIAL BASIS FUNCTION (RBF) KERNEL

$$k(\vec{x}, \vec{y}) = e^{-\frac{\|\vec{x} - \vec{y}\|^2}{2\sigma^2}} = \exp \left\{ -\frac{1}{2\sigma^2} (\|\vec{x}\|^2 + \|\vec{y}\|^2) + \frac{\vec{x}^T \vec{y}}{\sigma^2} \right\}$$

$$= \sum_{j=0}^{\infty} \frac{(\vec{x}^T \vec{y} / \sigma^2)^j}{j!} e^{-\frac{\|\vec{x}\|^2}{2\sigma^2}} e^{-\frac{\|\vec{y}\|^2}{2\sigma^2}} = \sum_{j=0}^{\infty} \sum_{\sum n_i = j} e^{-\frac{\|\vec{x}\|^2}{2\sigma^2}} \frac{n_1! \dots n_d!}{\sigma^j \sqrt{n_1! \dots n_d!}} e^{-\frac{\|\vec{y}\|^2}{2\sigma^2}} \frac{n_1! \dots n_d!}{\sigma^j \sqrt{n_1! \dots n_d!}}$$

TRANSLATION/
INTEGRAL
(Depends only on $\|\vec{x} - \vec{y}\|$)

OTHER EXAMPLES

• LAPLACIAN $\exp(-\|\vec{x} - \vec{y}\| / \sigma)$

• THIN-PLATE SPINE $\left(\frac{\|\vec{x} - \vec{y}\|}{\sigma}\right)^2 \log\left(\frac{\|\vec{x} - \vec{y}\|}{\sigma}\right)$

• SIGMOID $\tanh(a(\vec{x}^T \vec{y}) + b)$ * ONLY TRANSITION SUPPORT FOR CERTAIN VALUES OF a, b

- INFINITE DIMENSIONAL FEATURE SPACE !!
- GIVEN ANY LABELED DATA, USING THE KERNEL GIVES A UNIVERSALLY SEPARABLE SET IN FEATURE SPACE
- BASICALLY UNLIMITED FLEXIBILITY (RISK OF OVERFITTING)

NOTE:

- KERNELS CAN BE COMBINED TOGETHER TO CAPTURE MORE COMPLEX RELATIONSHIPS

$$k(x, y) = k_1(x, y) + k_2(x, y)$$

$$k(x, y) = \alpha k(x, y) \quad \alpha > 0$$

$$k(x, y) = k_1(x, y) k_2(x, y)$$

$$k(x, y) = f(x) f(y) \quad f: \mathbb{R}^d \rightarrow \mathbb{R}$$

- THE CHOICE OF KERNEL IMPLICITLY DEFINES AN INNER PRODUCT (HILBERT) SPACE OF FEATURES AS MAP $\phi: \mathbb{R}^d \rightarrow \mathcal{H}$

- BEST CHOICE OF KERNEL IS DEPENDENT ON THE PROBLEM

- BAYES: CORRECT KNOWLEDGE OR SHAPE OF SOLUTION IS USEFUL

- GOOD NEWS: CAN EXTEND SOME TECHNIQUES (OPTIMAL PARAMETER TUNING) FOR A CHOSEN MODEL

THEORETICAL DETAILS:

- MOTIVATION: LINEAR SEPARABILITY EXISTS IN HIGH-DIMENSIONAL

THM. (COVER 1965)

"A COMPLEX PATTERN-CLASSIFICATION PROBLEM, CAST IN A HIGH-DIMENSIONAL SPACE NONLINEARLY, IS MORE LIKELY TO BE LINEARLY SEPARABLE THAN IN A LOW-DIMENSIONAL SPACE, PROVIDED THAT THE SPACE IS NOT DENSELY POPULATED."

PROOF:

$$x_i \mapsto \vec{e}_i \in \mathbb{R}^n$$



• MAP DATA TO VERTICES OF SIMPLEX

• ANY PARTITION INTO TWO SETS IS SEPARABLE BY A HYPERPLANE.

- REQUIREMENTS & PROPERTIES OF KERNELS

- CAN WE CHOOSE ANY FUNCTION $k: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ AS A KERNEL? EQUIVALENTLY, DOES A GIVEN k HAVE A CORRESPONDING HILBERT SPACE \mathcal{H} AND

$$\text{MAP } \varphi: \mathbb{R}^d \rightarrow \mathcal{H} \text{ S.T. } k(\vec{x}, \vec{y}) = \underbrace{\langle \varphi(\vec{x}), \varphi(\vec{y}) \rangle_{\mathcal{H}}}_{\text{INNER PRODUCT}} ?$$

INNER PRODUCT

OF \mathcal{H} (USUALLY THE STANDARD DOT PRODUCT)

- MERCER'S CONDITION

THERE EXISTS A HILBERT SPACE AND MAP $\varphi: \mathbb{R}^d \rightarrow \mathcal{H}$ S.T. $k(\vec{x}, \vec{y}) = \langle \varphi(\vec{x}), \varphi(\vec{y}) \rangle_{\mathcal{H}}$

IFF

$$\text{FOR ANY } g(\vec{x}) \text{ S.T. } \int_{\mathbb{R}^d} [g(\vec{x})]^2 d\vec{x} < \infty, \text{ THEN } \underbrace{\int_{\mathbb{R}^d} \int_{\mathbb{R}^d} g(\vec{x}) k(\vec{x}, \vec{y}) g(\vec{y}) d\vec{x} d\vec{y}}_{\text{NOTE SIMILARITY TO P.S.D. MATRICES!}} \geq 0.$$

* k IS SEMI-POSITIVE DEFINITE

NOTE SIMILARITY

TO P.S.D. MATRICES!

EQUIVALENTLY, GIVEN ANY FINITE SET OF POINTS $x_1, \dots, x_N \in \mathbb{R}^d$ AND $c_1, \dots, c_N \in \mathbb{R}$

$$\sum_{i=1}^N \sum_{j=1}^N c_i c_j k(\vec{x}_i, \vec{x}_j) = \vec{c}^T \mathbf{K} \vec{c}$$

$\mathbf{K}_{ij} = k(\vec{x}_i, \vec{x}_j)$ IS CALLED THE GRAM MATRIX

VIEWING KERNELS AS OPERATIONS ON FUNCTION SPACES

• LET $L^2(\mathbb{R}^d) = \left\{ g: \mathbb{R}^d \rightarrow \mathbb{R} \mid \int [g(x)]^2 dx < \infty \right\}$

• ASSOCIATED W/ A KERNEL $k: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ IS A LINEAR OPERATION

$$(T_k f)(x) = \int k(\vec{x}, \vec{y}) f(\vec{y}) d\vec{y} \quad (\text{ASSUMING } \iint |k(\vec{x}, \vec{y})|^2 d\vec{x} d\vec{y} < \infty)$$

$$T_k: L^2(\mathbb{R}^d) \longrightarrow L^2(\mathbb{R}^d)$$

• MERCER'S THM. (SPECTRAL THM.) $k(\vec{x}, \vec{y}) = k(\vec{y}, \vec{x})$

SUPPOSE k IS A CONTINUOUS, SYMMETRIC KERNEL SATISFYING MERCER'S CONDITIONS.
 THEN THERE EXISTS AN ORTHONORMAL BASIS $\{\phi_j\}_{j=1}^{\infty} \subset L^2(\mathbb{R}^d)$ OF EIGENFUNCTIONS
 OF T_k W/ A CORRESPONDING SEQUENCE OF NONNEGATIVE EIGENVALUES $\{\lambda_j\}_{j=1}^{\infty}$,
 AND

$$k(\vec{x}, \vec{y}) = \sum_{j=1}^{\infty} \lambda_j \phi_j(\vec{x}) \phi_j(\vec{y}) \quad (\text{DIAGONALIZATION})$$

• REPRODUCING PROPERTY

• A KERNEL k (WITH CORRESPONDING HILBERT SPACE \mathcal{H}) HAS THE REPRODUCING PROPERTY IF FOR ANY $f \in \mathcal{H}$

$$\langle f(\cdot), k(\cdot, \cdot) \rangle_{\mathcal{H}} = \langle f, \phi(\cdot) \rangle_{\mathcal{H}} = f(x)$$

AND WE SAY k IS A REPRODUCING KERNEL.

• FOR ANY POSITIVE SEMIDEFINITE (MERCEK) KERNEL, WE CAN CONSTRUCT A UNIQUE HILBERT SPACE \mathcal{H}_k FOR WHICH k IS ITS REPRODUCING KERNEL. \mathcal{H}_k IS CALLED A REPRODUCING KERNEL HILBERT SPACE (RKHS)

• RKHSs HAVE MANY APPLICATIONS IN STATISTICS (KERNEL PCA)
 FOR EXAMPLE

• MOORE-ARONZAJAN THEOREM PROVES THIS CLAIM IN A CONSTRUCTIVE MANNER.