

KERNELS

MEASURE SIMILARITY BETWEEN OBJECTS **KERNEL FUNCTION**: $K(x, x') \in \mathbb{R}$, TYPICALLY SYMMETRIC AND NON-NEGATIVE, USED WHEN NO CLEAR HOW TO MAKE FEATURE VECTORS, IE STRINGS

GAUSSIAN / RBF KERNELS

$K = \exp\left(-\frac{1}{2}(x-x')^T \Sigma^{-1}(x-x')\right)$ **SQUARED EXPONENTIAL** • Σ^{-1} SCALE DIMENSION / BANDWIDTH

$K = \exp\left(-\frac{1}{2}(x-x')^T \Sigma^{-1}(x-x')\right)$ **ARD IF $\Sigma^{-1} \rightarrow \mathbb{I}$, Σ DIAGONAL**

$K = \exp\left(-\frac{\|x-x'\|^2}{2\sigma^2}\right)$ **ISOTROPIC / SPHERICAL** • **RBF** BECAUSE ONLY FUNCTION OF $\|x-x'\|$

DOCUMENT COMPARISON

$K = \frac{x^T x'}{\|x\|_2 \|x'\|_2}$ **COSINE SIMILARITY**, ANGLE BETWEEN VECTORS (0,1); SENSITIVE TO STOP-WORDS
 → USE WITH **TF-IDF TRANSFORM** $TF_{ij} = \log(1+x_{ij})$ **IDF** = $\frac{N}{1+\sum 1(x_{ij} > 0)}$ **TF-IDF** = $[TF(x_{ij}) \times IDF(j)]_i$
TERM FREQUENCY INVERSE DOCUMENT FREQUENCY

MERCER / POSITIVE DEFINITE KERNELS

GRAM MATRIX = $K = \begin{bmatrix} K(x_1, x_1) & \dots & K(x_1, x_M) \\ \vdots & & \vdots \\ K(x_N, x_1) & \dots & K(x_N, x_M) \end{bmatrix}$

- A KERNEL IS MERCER IFF GRAM MATRIX IS POSITIVE DEFINITE \forall INPUT
- IF KERNEL MERCER $\exists \phi: x \mapsto \mathbb{R}^D$ $K(x, x') = \phi(x)^T \phi(x')$, ϕ DEPENDS ON EIGENFUNCTIONS OF K
- **POLYNOMIAL KERNEL** $K = (\gamma x^T x' + R), R > 0$
- RBF AND COSIM ARE MERCER
- **SIGMOID** $K = \tanh(\gamma x^T x' + R)$ IS NOT MERCER
- IF K_1, K_2 MERCER $\rightarrow K_1 + K_2$ MERCER
- CAN GO KERNEL \rightarrow FEATURE VECTOR

LINEAR KERNELS

IF $\phi(x) = x$ $K = x^T x'$. GOOD WHEN HIGH DIM AND ORIGINAL FEATURES ARE INFORMATIVE, LINEAR DECISION BOUNDARY WORKS, NO NEED TO TRANSFORM SPACE

MATERN KERNEL

$R = \|x-x'\|, \nu > 0, R > 0, K_\nu$ BESSEL FUNCTION; $\nu \rightarrow \infty$ SE KERNEL
 $K = \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu} R}{l}\right)^\nu K_\nu\left(\frac{\sqrt{2\nu} R}{l}\right)$ • $\nu = 1/2 \rightarrow K = \exp(-R/l)$, IF USED IN GAUSSIAN PROCESS; ORNSTEIN-UHLENBECK PROCESS
 DESCRIBES VELOCITY OF BROWNIAN MOTION PARTICLE

STRING KERNELS

$W_S > 0$; A^* SET OF ALL POSSIBLE STRINGS IN ALPHABET (UNLESS STAR), $\phi_S(x)$ NO OF TIMES SUBSTRING s APPEARS IN x

• $K = \sum_{s \in A^*} W_s \phi_S(x) \phi_S(x')$ • IS MERCER • $O(|x| + |x'|)$

- $W_s = 0$ FOR $|s| > 1 \rightarrow$ BAG OF CHARACTERS
- s BAG OF W/WHITESPACE \rightarrow BAG OF WORDS
- FIXED LENGTH $k \rightarrow$ **k-SPECTRUM**
- CAN GENERALIZE TO PARSING / COMPARISON TREES

PYRAMID MATCH KERNELS

USED FOR IMAGES, **SIFT**, FEATURES VECOR EXTRACTED FROM IMAGE POIS, VECTORIZED, BAG OF SYMBOLS, FEATURES MAPPED TO HISTOGRAM, MULTI RESOLUTION COMPARISON W/ INTERSECTION. APPROXIMATES OPTIMAL MATCH OF SIMILARITIES BETWEEN MATCHING POINTS. IS MERCER

PROBABILISTIC KERNELS

PROBABILISTIC GENERATIVE MODEL OF FEATURE VECTORS $p(x|\theta)$

- PROBABILITY PRODUCT KERNEL

$$k = \int p(x|x_1)^{\theta_1} p(x|x_2)^{\theta_2} dx, \theta > 0, p(x|x_1) \rightarrow p(x|\hat{\theta}(x_1)), \hat{\theta} \text{ PARAM COMPUTED FROM SINGLE VECTOR}$$

- FITS MODEL TO SINGLE DATAPoint \rightarrow USES IT TO MEASURE SIMILARITY
- IF $p(x|\theta) = N(\mu, \sigma^2 I), \sigma \text{ FIXED}, \theta = 1, \rightarrow$ OBTAIN PDF

- FISHER KERNELS

$$g(x) = \nabla_{\theta} \log p(x|\theta)|_{\hat{\theta}} \quad F = \nabla \nabla^T \log p(x|\theta)|_{\hat{\theta}} \quad \hat{\theta} \text{ ALL DATA.}$$

\rightarrow LL GRADIENT, SCORE

\rightarrow FISHER INFO MATRIX, HESSIAN

- $g(x)$ DIRECTIONAL GRADIENT TO MAXIMIZE LIKELIHOOD

- KERNEL MEASURES SIMILARITY IF $g(x)$'s ARE SIMILAR WRT GEOMETRY OF LIKELIHOOD FCN CURVATURE

$$K = g(x)^T F^{-1} g(x')$$

KERNEL MACHINES

GLM, WHERE INPUT FEATURE VECTOR IS $\phi(x) = [n(x, \mu_1) \dots n(x, \mu_n)] \quad \mu_n = \text{CENTROIDS}$

\rightarrow KERNELIZED FEATURE VECTOR

- THEN USE KFM FOR LOGISTIC $p(y|x, \theta) = \text{BER}(w^T \phi(x))$ OR LINEAR $p(y|x, \theta) = N(w^T \phi(x))$ REGRESSION.

- OK FOR NON LINEARLY SEPARABLE BOUNDARIES
- BANDWIDTH AFFECTS FREQUENCY, RANGE

• HOW TO PICK μ_n ?

IF FEW DIM \rightarrow TILE THE SPACE UNIFORMLY

IF DIM HIGH \rightarrow NUMERICAL OPTIMIZATION / MCMC

IDEA: FIND CLUSTERS IN DATA AND ASSIGN ONE PROTOTYPE PER CLUSTER

BUT MOST DENSITY OF POINTS \neq MOST USEFUL

BUT I STILL HAVE TO PICK NO. OF CLUSTERS!

IDEA #2: MAKE EACH EXAMPLE A PROTOTYPE!

$$\phi(x) = [n(x, x_1), \dots, n(x, x_N)] \rightarrow D = N \rightarrow$$

- SPARSITY PRIOR ON $w \rightarrow$ SPARSITY VECTOR MACHINE, I.E. ℓ_1 VM, OR GROUP LASSO FOR MULTICLASS

- ℓ_2 VM \rightarrow OFC NOT SPARSE

- ARD/SBL \rightarrow RELEVANT VECTOR MACHINES

- SUPPORT VECTOR MACHINE \rightarrow MODIFIES LIKELIHOOD TERM! NO SPARSITY PRIOR

TEH KERNEL TRICK!

NOT DEFINE KERNELIZED FEATURE VECTOR \rightarrow WORK WITH ORIGINAL BUT REPLACE INNER PRODUCTS $\langle x, x' \rangle$ WITH $k(x, x')$. ONLY IF K IS MERCEL

- KERNELIZED NN CLASSIFICATION $\|x - x'\|_2^2 = \langle x, x \rangle + \langle x', x' \rangle - 2 \langle x, x' \rangle$

$$\langle x, y \rangle = x^T y = \sum_i x_i y_i$$

- KERNELIZED K-MEANS CLUSTERING: EACH CENTROID IS ONE OF DATA VECTORS, NOT ARBITRARY POINTS

WHEN UPDATE MEASURES, DISTANCE OF ALL CLUSTER OBJECTS TO ALL OTHER IN CLUSTER; PICK ONE WITH LOWEST $M_k = \text{ARG-MIN} \sum d(i, i')$

$\rightarrow O(n_k^2)$ PER CLUSTER VS $O(n_k D)$ OF K-MEANS

- CAN TURN INTO CLASSIFIER

- CAN KERNELIZE

$$d(i, i') = \|x - x'\|_2^2$$

- KERNELIZED RIDGE REGRESSION

$$\text{PRIMAL PROBLEM: } w = (\sum x_i x_i^T + \lambda I_D)^{-1} \sum x_i y_i \xrightarrow{\text{MATRIX INVERSION LEMMA}} w = X^T (X X^T + \lambda I_D)^{-1} y$$

$$\text{DUAL PROBLEM: } \alpha = (K + \lambda I_N)^{-1} y \rightarrow w = X^T \alpha = \sum \alpha_i x_i$$

\downarrow $K X X^T = \text{GRAM MATRIX}$

- SOLUTION IS LINEAR COMBINATION/SUM OF TRAINING VECTORS

$$\text{PREDICTION: } \hat{f}(x) = w^T x = \sum \alpha_i x x^T = \sum \alpha_i k(x, x_i)$$

- COMPUTING α IS $O(N^3) \rightarrow$ w IS $O(D^2)$ SO USEFUL IN HIGH DIM JUST BECAUSE SPEED; PREDICTION WITH α IS $O(ND)$, WITH w IS $O(D)$

\rightarrow SPEEDUP WITH SPARSITY

KERNEL PCA

PCA BY FINDING EIGENVECTORS OF XX^T AND KERNEL TRICK.

- U, λ EIGENVECTOR, VALUES OF $XX^T \rightarrow$ EIGENVECTORS OF $X^T X = V = X^T U$ $\xrightarrow{\text{NORMALIZE}}$ $V_{PCA} = X^T U \lambda^{-1/2}$
 - USEFUL EVEN W/O KERNEL BECAUSE IS $D \times D$ VS $N \times N$
- $K = XX^T$, Φ = NATURAL DESIGN MATRIX IN FEATURE SPACE $\rightarrow V_{KPCA} = \Phi^T U \lambda^{-1/2}$ (CANNOT COMPUTE! POSSIBLY ∞ DIM)
 - $V = \sum_i \alpha_i \phi(x_i)$ α_i EVES NORMAL PROBLEM
 - α_i EVES KERNEL PROBLEM
- CAN PROJECT x_* IN FEATURE SPACE
 - $\Phi_*^T V_{KPCA} = \Phi_*^T \Phi^T U \lambda^{-1/2} = K_*^T U \lambda^{-1/2}$
 - $K_* = [k(x_*, x_1) \dots k(x_*, x_N)]$
- PROJECTED DATA NOT CENTERED, CANNOT SIMPLY SUBTRACT MEAN $\rightarrow \tilde{\Phi}_i = \phi(x_i) - \frac{1}{N} \sum \phi(x_i)$
 - CENTERED FEATURE VECTOR
 - $\tilde{K} = H K H$
- $H = I - \frac{1}{N} \mathbf{1} \mathbf{1}^T$ CENTERING MATRIX
- WITH $k(x, x')$ WE IMPLICITLY REPLACE x WITH $\Phi(x)$
- KPCA CAN DO UP TO N COMPONENTS
- USEFUL FOR CLASSIFICATION, NOT SO MUCH FOR VISUALIZATION

ALGO: PRODUCE $U = U_{(1)}(k(1,1))$, NORMALIZE + CENTER
 COMPUTE EVES EVALS OF U
 RETAIN TOP EVES; NORMALIZE BY SQRT OF EVALS
 USE THEN TO PROJECT DATA

SUPPORT VECTOR MACHINES

$J(w, \lambda) = \sum L(y_i, \hat{y}_i) + \lambda \|w\|^2$ IS OF EMPIRICAL RISK

HMMS: PROBABILISTICALLY UNNATURAL, SPARSITY IN LOSS AND NOT IN PRED. ENCODE KERNEL VIA ALGORITHMIC TRICK DO NOT RESULT IN PROBABILISTIC OUTPUTS

IDEA: WE CAN REDUCE LOSS TO SOMETHING THAT WILL ENFORCE SPARSITY SO THAT PREDICTIONS ONLY DEPEND ON SUBSET OF TRAINING DATA

- \rightarrow SUPPORT VECTORS: POINTS FOR WHICH ERROR IS OUTSIDE TUBE
- KERNEL TRICK + MODIFIED LOSS
- MINIMIZE CLASSIF. ERROR, MAXIMIZE GEOMETRIC MARGIN

FOR REGRESSION:

- EPSILON INSENSITIVE LOSS FUNCTION: $L_\epsilon(y, \hat{y}) = \begin{cases} 0 & \text{IF } |y - \hat{y}| \leq \epsilon \\ |y - \hat{y}| - \epsilon & \text{ELSE} \end{cases}$
 - ANY POINT OUTSIDE ϵ -TUBE IS PENALIZED L ON PREDICTION
- NOT DIFFERENTIABLE: USE SLACK VARIABLES TO REPRESENT DEGREE OF POINT OUTSIDE TUBE. ALLOW MISCLASSIFICATION SOFT MARGIN
- OPTIMAL SOLUTION $\hat{w} = \sum \alpha_i x_i$, α IS SPARSE BECAUSE DON'T CARE ABOUT ERRORS $\leq \epsilon$
- x FOR WHICH $\alpha_i > 0 \rightarrow$ SUPPORT VECTORS
- PREDICTIONS: $\hat{y}(x) = \hat{w}_0 + \sum \alpha_i k(x_i, x)$
 - $J = C \sum (\xi^+ + \xi^-) + \frac{1}{2} \|w\|^2$
 - LOSS / OBJECTIVE
 - C IS INVERSE REG. COEFFICIENT

FOR CLASSIFICATION:

- HINGE LOSS: $L_{\text{HINGE}}(y, \eta) = \max(0, 1 - y\eta) = (1 - y\eta)^+$ $\eta = f(x)$ = 'CONFIDENCE' OF PREDICTION. NO PROBABILISTIC. QUADRATIC PROBLEM/ALGO
- OBJECTIVE: $\min \frac{1}{2} \|w\|^2 + C \sum \xi_i, \xi_i \geq 0$ SOLVED IN $O(N^3)$, OPTIM TO $O(N^2)$, C REGULARIZERS
- x INCORRECTLY CLASSIFIED \rightarrow SUPPORT VECTORS
- PREDICTION: $\hat{y}(x) = \text{SGN}(\hat{w}_0 + \sum \alpha_i k(x_i, x))$ $O(ND)$ TO COMPUTE

MARGIN PRINCIPLE: SVM MAXIMIZES PERPENDICULAR DISTANCE TO CLOSEST POINT, MARGIN. WIDE MARGIN CLASSIFIER. SOFT MARGIN CONSTRAINTS WITH ξ $\sum \xi_i$ IS UPPER BOUND ON MISCLASSIFICATIONS AT TRAINING TIME. $C = 1/N$ CONTROLS THIS FRACTION $\rightarrow \nu$ -SVM

PROBABILISTIC OUTPUT: NO BUILT-IN SPARSE, CAN INTERPRET VIA LOG-ODDS RATIO. BUT RESULTS NOT WELL CALIBRATED

MULTI-CLASS CLASSIFICATION: OUTPUT NOT CALIBRATED \rightarrow CANNOT INTO SOFTMAX. ONE-VS-REST: C CLASSIFIERS, EACH ON $C^i, (1-C^i)$: PROMPTOMATIC CLASS IMBALANCE. ONE VS ONE: $C(C-1)/2$ CLASSIFIERS, ON EVERY PAIR, CLASSIFY ON HIGHEST VOTE COUNT.

PICKING C : RECOMMENDED PICK IT WITH CV OVER 2d GRID VALUES. ALTERNATIVELY \pm LASSO/LARS ALGO. START WITH LARGE λ

PROBABILISTIC INTERPRETATION: WITH A LOT OF RELAXATIONS WE CAN INTERPRET HINGE-LOSS AS GAUSSIAN SCALE MIXTURE \rightarrow CAN USE BAYESIAN TO SET HYPERTUNING

PERFORMANCE: ALL KERNEL METHODS SIMILAR ACCURACY OVER RANGE OF PROBLEMS AND SOME KERNEL. GP LLVM ARE $O(N^3)$ SVM IS $O(N)$ WHEN UNFOLNED. SVM USUALLY SLOWER THAN LLVM BUT GREEDY TRAINING IS FASTER. IF SPEED MATTERS: USE RVM

IF CALIBRATED PROBABILITY MATTERS: USE GP (GAUSSIAN PROCESS)

SMOOTHING KERNELS

USED FOR NONPARAMETRIC DENSITY ESTIMATES. UNSUPERVISED DENSITY ESTIMATION OR GENERATIVE CLASSIF/REGRESSION MODELS

$$\bullet \int n(x) dx = 1, \int x n(x) dx = 0, \int x^2 n(x) dx > 0$$

GAUSSIAN/RBF: $k_h(x) = \frac{1}{h^D (2\pi)^{D/2}} \exp\left(-\frac{1}{2h^2} x^2\right)$ H IS BANDWIDTH, 'WIDTH' OF KERNEL, σ^2

• UNBOUNDED SUPPORT

EPANECHNIKOV: $k(x) = \frac{3}{4} (1-x^2) \mathbb{I}(|x| \leq 1)$ • COMPACT SUPPORT
• NONDIFF AT ITS BOUNDARY

TRI-CUBE: $k(x) = \frac{70}{81} (1-|x|^3) \mathbb{I}(|x| \leq 1)$

• COMPACT

• TWICE DIFFERENTIABLE AT ITS BOUNDARY

BOXCAR: UNIFORM DISTRIBUTION $k(x) = \mathbb{I}(|x| \leq 1)$

KERNEL DENSITY ESTIMATION (KDE)

ALTERNATIVE TO PARAMETRIC DENSITY ESTIMATION WITH GAMM \rightarrow NEED TO SPECIFY K, μ, σ

HERE WE JUST SET A CLUSTER CENTER PER DATAPoint $\mu_i = x_i$

KERNEL DENSITY/PARZEN WINDOW ESTIMATION: $\hat{P}(x) = \frac{1}{N} \sum k_h(x-x_i)$

- NO MODEL FITTING \rightarrow JUST SET H WITH CV, NO NEED TO PICK K (CV OR DIRECTED P-VALUE MIXTURE, BAYESIAN)
- LOT OF MEMORY TO STORE, USELESS FOR CLUSTERING TASKS
- CV MINIMIZE FREQUENTIST R_{JS}

EXAMPLES: BOXCAR \rightarrow HISTOGRAM COUNT
(1D) RBF \rightarrow SMOOTHER HISTO

KNN/CASSIFIER

WE GROW VOLUME AROUND x UNTIL K DATAPoints REGARDLESS OF CLASS ^{MODEL VOLUME}. $V(x), N_c(x)$ SAMPLES OF CLASS c

• CLASS-COND. DENSITY $P(x|y=c, D) = \frac{N_c(x)}{N_c V(x)}$, N_c TOTAL COUNT FOR c • CLASS PRIOR: $P(y=c|D) = \frac{N_c}{N}$

• CLASS POSTERIOR: $P(y=c|D) = \frac{N_c(x)}{K}$

KERNEL REGRESSION

COMPUTE CONDITIONAL EXPECTATION $f(x) = E(y|x) = \int y \cdot P(y|x) dy$, KDE APPROX $P(x,y) = \frac{1}{N} \sum k_h(x-x_i) k_h(y-y_i)$. PVA PVA. SUM TO ONE ZERO MEAN...

$$\bullet f(x) = \sum w_i(x) y_i$$

$$w_i(x) = \frac{k_h(x-x_i)}{\sum k_h(x-x_i)}$$

• PREDICTION IS WEIGHTED SUM OF OUTPUTS AT TRAINING POINTS, WEIGHT DEPEND ON SIMILARITY TO TRAINING POINTS

• ONLY 1 FREE PARAMETER

• OPTIMAL $H = \left(\frac{4}{3N}\right)^{1/5} \hat{\sigma}$

LOCALLY-WEIGHTED REGRESSION:

KERNEL REGRESSION FITS A CONSTANT FUNCTION LOCALLY. WE CAN IMPROVE BY FITTING A LINEAR REGRESSION MODEL FOR EACH POINT x
 $\min \sum w(x, x_i) [y_i - \beta(x)^T \phi(x_i)]^2$, $\phi(x) = [1, x]$, $\beta(x) = (\Phi^T D(x) \Phi)^{-1} \Phi^T D(x) y$, Φ DESIGN MATRIX, $D = \text{DIAG}(w(x, x_i))$

• $\hat{f}(x) = \phi(x)^T \hat{\beta}(x) = \sum w_i(x) y_i$ \leftarrow PREDICTION, $w_i(x)$ IS EQUIVALENT KERNEL!