

GAUSSIAN PROCESSES

FROM DISTRIBUTION OVER FCN

ARE SMOOTH

OPTIMAL SUPERVISED LEARNING: INFER DISTRIBUTION OF FUNCTIONS OVER DATA, USE IT TO PREDICT, $P(f|D)$ VS $P(D|D)$

$$P(y|x, X, y) = \int P(y|f, x) P(f|x, X) df$$

• MAIN GP DRAWBACK: THEY TAKE $O(N^3)$ DUE TO N INVERSION/DECOMPOSITION

GP: DEFINES PRIOR OVER FCN \rightarrow ~~FCN~~ INTO POSTERIOR AFTER SEEING SOME DATA. DEFINE DISTRIBUTION OVER FCN VALUES AT FINITE SET OF POINTS $x_1 \dots x_n$. $f(x_1) \dots f(x_n)$ IS JOINTLY GAUSSIAN $\mu(x)$, $\Sigma(x) = \sum_{i,j} k(x_i, x_j)$ KERNEL FUNCTION. KERNEL SIMILAR \rightarrow FCN OUTPUT SIMILAR.

IN REGRESSION CLOSED FORM IN $O(N^3)$, IN CLASSIFICATION APPROXIMATIONS. k IS COV. BAYESIAN ALTERNATIVE TO KERNEL METHODS. GRAPHICALLY f VALUES ARE HIDDEN NODES

GP FOR REGRESSION AND NOISING

$$f(x) \sim GP(\mu(x), k(x, x'))$$

MEAN FCN COV FCN / KERNEL JOINT GAUSSIAN

$$\mu(x) = E[f(x)] \quad k(x, x') = E[(f(x) - \mu(x))(f(x') - \mu(x'))^T] \quad \rightarrow \quad P(f|x) = N(f|\mu, k)$$

NOISE-FREE

WE WANT SAME VALUES FOR ALREADY SEEN x_s , INTERPOLATION. $x_s \rightarrow f_s$

• JOINT: $\begin{pmatrix} f \\ f_s \end{pmatrix} \sim N\left(\begin{pmatrix} \mu \\ \mu_s \end{pmatrix}, \begin{pmatrix} k & k_s \\ k_s^T & k_{ss} \end{pmatrix}\right)$

• POSTERIOR $P(f_s | x_s, X, f) = N(f_s | \mu_s, \Sigma_s)$

• $\mu_s = \mu(x_s) + k_s^T k^{-1} (f - \mu(x))$

• $\Sigma_s = k_{ss} - k_s^T k^{-1} k_s$

- IS SAMPLING FROM MVN OF DIM OF TRAINING SET
- MAXIMIZE MARGINAL LIKELIHOOD

SQUARE EXPONENTIAL KERNEL:

$$k(x, x') = \sigma^2 \exp\left(-\frac{1}{2\ell^2} (x - x')^2\right)$$

ℓ = HORIZONTAL SCALE
 σ^2 = VERTICAL SCALE

- SUM OF SET OF BASIS FCN POINTING ON TRAINING DATA x_s ROW VECT
- EM HAS 29 COV FCN

NOISY OBSERVATIONS

WE OBSERVE $y = f(x) + \epsilon$, $\epsilon \sim N(0, \sigma_y^2)$. $\text{COV}[y|x] = k + \sigma_y^2 I_N = k_y$. DIAGONAL BECAUSE WE ASSUME INDEPENDENT NOISE.

POSTERIOR

$$\mu_s = K_s^T K_y^{-1} y, \quad \Sigma_s = k_{ss} - K_s^T K_y^{-1} K_s$$

• POSTERIOR MEAN: $\bar{f}_s = K_s^T K^{-1} y = \sum \alpha_i k(x_i, x_s)$, $\alpha = K_y^{-1} y$

• PERFORMANCE DEPENDS EXCLUSIVELY ON KERNEL AND ITS PARAMS. IN RBF ℓ NOISY ℓ, σ_y^2 ARE H-V SCALE AND σ_y^2 NOISE VARIANCE. LARGE SCALE \rightarrow IRRELEVANT DIMENSION

IN MULTI-D SE IS: $k(x_1, x_2) = \sigma_f^2 \exp\left(-\frac{1}{2} (x_1 - x_2)^T M (x_1 - x_2)\right) + \sigma_y^2 \delta_{pq}$. M IS ISOTROPIC $\ell^{-2} I$, $\text{DIAG}(\ell^{-2})$

KERNEL PARAMETER ESTIMATION

CAN GO AND SEARCH BUT SLOW. EMPIRICAL BAYES, MAXIMIZE MARGINAL LIKELIHOOD; USE STANDARD GRADIENT OPTIMIZERS

$$P(y|x) = \int P(y|f, x) P(f|x) df \rightarrow \log P(y|x) = \log N(y|0, K_y) = -\frac{1}{2} y^T K_y^{-1} y - \frac{1}{2} \log |K_y| - \frac{N}{2} \log(2\pi)$$

DATA FIT MODEL COMPLEXITY CONSTANT

TRADEOFF

- VIA MARGINAL LIKELIHOOD AND GRADIENT DESCENT

$$\frac{\partial}{\partial \theta_i} \log P(y|x) = \frac{1}{2} \text{TR}\left[(\alpha \alpha^T - K_y^{-1}) \frac{\partial K_y}{\partial \theta_i}\right], \quad \alpha = K_y^{-1} y$$

$O(N^3)$ FOR K_y^{-1} AND THEN $O(N^2)$ PER HYPERPARAM

ALTERNATIVE: BAYESIAN INFERENCE

COMPUTE POSTERIOR AND NOT POINT ESTIMATES $P(f|D) \propto \sum P(f|D, \theta_s) P(\theta_s|D) \delta_s$

CENTRAL COMPOSITE DESIGN

PUT GRID POINT AT MODE AND AT $\pm 1\sigma$ ON EACH DIMENSION

PUT CURSE OF DIMENSIONS IF GRID AND HIGH DIM \rightarrow MONTECARLO

ALTERNATIVE: MULTIPLE KERNEL LEARNING

$k(x, x') = \sum w_i k_i(x, x')$, OPTIMIZE w_i S \rightarrow DATA FUSION

- COMPUTATIONAL ISSUES: \bar{f}_x , STUPID TO DIRECTLY INVERT K_y . USE CHOLESKY DECOMPOSITION $K_y = LL^T$ OR SOLVE $K_y \alpha = y$ WITH CONJUGATE GRADIENTS

• 1 $O(N^3)$ CHO, $O(N^2)$ FOR d , $O(N)$ MEAN, $O(N^2)$ VARIANCE

• 2 $O(N^2)$ FOR N ITERATIONS. $O(N^3)$ FOR FULL.

SEMI-PARAMETRIC GP

ADAPTS LINEAR MODEL FOR PROCESS MEAN $f(x) = \beta^T \phi(x) + R(x)$. $R(x) \sim GP(0, k(x, x'))$ MODELS THE RESIDUALS

IF $\beta \sim N(b, B) \rightarrow f(x) \sim GP(\phi(x)^T b, k(x, x') + \phi(x)^T B \phi(x))$ • INTEGRATE β OUT FOR PREDICTIVE DISTRIBUTION FOR X

INTEGRATE OUT PARAMS

• PREDICTIVE MEAN IS LINEAR MODEL OUTPUT + CORRECTION DUE TO GP
+ 1 COVARIANCE IS GP COV + UNCERTAINTY IN β

GP + GLM FOR CLASSIFICATION

MAIN PROBLEM: GAUSSIAN PRIOR NOT CONJUGATE TO BERNULLI/MULTINOUM LIKELIHOOD. APPROXIMATIONS: EXPECT-BADAPPROXIMATION, VARIATIONAL, MCMC, GAUSSIAN

• BINARY CLASSIFICATION: $P(y_i | x_i) = \sigma(y_i, f(x_i))$, $\sigma(z) = \text{SIGM}(z)$ LOGISTIC, OR $\Phi(z)$ PROBIT. $f \sim GP(0, k)$

- POSTERIOR: $P(f | x, y) \sim N(f, (K^{-1} + W)^{-1})$ - POSTERIOR PREDICTIVE: $\pi_x = P(y=1 | x, x', y) = \int \sigma(f_x) P(f_x | x, x', y) df_x$

- MARGINAL LIKELIHOOD: $\log p(y | x) = \log \int p(y | f, x) p(f) df = \log \int \prod_i \sigma(f_i) p(f_i) df = \log \int \prod_i \sigma(f_i) \prod_i p(f_i) df$
• DO NOT DIRECTLY INVERT K, L . USE CHOLESKY AND OTHER TRICKS
FITTING IS $O(N^3)$, PREDICTION IS $O(N^2 N_x)$ ALGO 15.2 P 624 MURPHY

• MULTI-CLASS: $P(y_i | x_i) = \text{CAT}(y_i | S(f_i))$, $f_i = (f_{i1}, f_{i2}, \dots, f_{iC})$, $f_{ic} = GP(0, k_c)$ ONE LATENT FCN PER CLASS. ARGOSY INDEPENDENT, MAY DIFFERENT KERNELS

- POSTERIOR! $\odot \odot$ - POST PREDICTIVE: $P(y | x, x', y) \approx \int \text{CAT}(y | S(f_x)) N(f_x | E[f_x], \text{COV}[f_x]) df_x$
 $O(CN^2)$

- MARGINAL LIKELIHOOD: $\odot \odot$ - COMPLEXITY: FIT IN $O(TCN^3)$ TIME AND $O(CN^2)$ SPACE

• FINDING THE PRIOR OVER $f(x)$. (LATENT FCN) \rightarrow PUT IT THROUGH LOGISTIC FUNCTION
TO FIND PRIOR ON PREDICTED CLASS
PREDICTION IN $O(CN^3 + CN^2 N_x)$

GAUSSIAN PROCESSES VS THE WORLD

ARM, NUMMS, PDF NETS, ARE ALL GPs!!
CRF

- **LINEAR MODELS** EQUIVALENT TO GP WITH $k(x, x') = x^T \Sigma x'$. DEGENERATE COV FCN BECAUSE AT MOST D NON ZERO EIGENVALUES
UNDERFITTING (NOT FLEXIBLE ENOUGH) OVERCONFIDENT (PRIOR TOO POOR \rightarrow POSTERIOR TOO CONCENTRATED)

- **LINEAR SMOOTHIES**
ARE STUFF OF THE FORM $f(x) = \sum w_i(x) y_i$. LINEAR FUNCTIONS OF TRAINING OUTPUTS. KERNEL REGRESSION, LOCALLY WEIGHTED REGRESSION
ARE GPs BECAUSE POSTERIOR PREDICTIVE MEAN CAN BE EXPRESSED AS SUCH. $w_i(x) = [k(x, x_i) / (k(x, x) + \sigma^2)]$
IN GPs KERNEL BANDWIDTH IS AUTO DECREASING WITH INCREASING N
DOF: $\text{TR}(k(x, x) + \sigma^2 I)^{-1} = \sum \frac{\lambda_i}{\lambda_i + \sigma^2}$ IF LL 1 CORRESPONDING BASIS FUNCTION HAS LITTLE INFLUENCE

- **SVM** CAN REWRITE SVM OBJECTIVE IN FORM ANALOGOUS TO MAP ESTIMATORS FOR GP. NO CAN FULLY CONVERT BECAUSE THERE IS NO LINEAR MODE MATCHING HINGE-LOSS.

- **L1VM, RVM** EQUIVALENT TO GP WITH $k(x, x') = \sum \frac{1}{\alpha_i} \phi_i(x) \phi_i(x')$ DEGENERATE AND DEPENDS ON TRAINING DATA \rightarrow OVERCONFIDENT

- **NEURAL NETS**: ARE A NONLINEAR GENERALIZATION OF GLM. IF NO. OF HIDDEN UNITS $\rightarrow \infty$ AND ACTIVATION FCN IS BOUNDED WE GET GP
CAPTURING TRUE KERNEL FORMS ACCORDING TO GP.

- **SMOOTHING SPLINES**: NONPARAMETRIC MODELS USED TO INTERPOLATE 1D OR 2D DATA
FIT f MINIMIZING DISCREPANCY OF DATA PLUS SMOOTHING PENALIZATION FOR WIGGLYNESS
 $J(f) = \sum (f(x_i) - y_i)^2 + \lambda \int \frac{d^m}{dx^m} f(x)^2 dx$ **PIECEWISE POLYNOMIAL**
CAN BE FIT WITH RIDGE REGRESSION **REGRESSION SPLINE**: PLACES POLYNOMIALS AT FIXED SET OF K LOCATIONS. UNITS.
(CHOOSING NO AND RULE OF UNITS IS LIKE PICKING SUPPORT VECTORS)
CUBIC SPLINE IS MAP OF $f(x) = \beta_0 + \beta_1 x + R(x)$ $R(x) \sim GP(0, \sigma_f^2 k(x, x'))$, $k(x, x') = \int_0^1 (x-u)(x'-u) du$ **GENERALIZES TO HIGHER DIMS**

- **RKHS**:
REPRODUCING-KERNEL HILBERT SPACES. GENERALIZES NOTION OF SMOOTHNESS. CPLEX FCNS W/ KERNEL \rightarrow LARGE NORMS BECAUSE MANY EIGENFUNCTIONS
- **PROBLEM**: $J(f) = \frac{1}{2\sigma_y^2} \sum (y_i - f(x_i))^2 + \frac{1}{2} \|f\|_H^2 \rightarrow$ FCN HAS TO BE $f(x) = \sum \alpha_i k(x, x_i) \rightarrow \hat{f}(x_i)$ IS SAME OF POSTERIOR MEAN OF GP PREDICTIVE DISTRIBUTION

GP LATENT VARIABLE MODEL

COMBINES KERNELS WITH PROBABILISTIC PCA \rightarrow GP-LVM • CAN FORMULATE PPCA DUALY. MAXIMIZE Z AM INTEGRATE W OUT.
• LIKELIHOOD COMES OUT DEPENDING ON $Y Y^T \rightarrow$ CAN SOLVE WITH EIGENVALUE METHOD. • IF LINEAR KERNEL \rightarrow PCA BUT ALSO $K_z = K + \sigma^2 I$
• NO MLE ANYMORE \rightarrow GRADIENT BASED OPTIMIZERS K IS GRAM MATRIX **GENERAL K**
• LEARNS MAPPING FROM LATENT SPACE TO OBSERVED SPACE
• PCA IS GP-LVM W/ LINEAR KERNEL
• PRIOR IS ON W
• LOCATIONS OF POINTS IN LATENT SPACE IS DONE BY MAXIMIZING LIKELIHOOD W/ X

SPLINES