

MACHINE LEARNING

- **SUPERVISED** → CLASSIFICATION / REGRESSION | CLASS: Y IS CATEGORICAL, REG: Y IS REAL CONTINUOUS
- **UNSUPERVISED** → CLUSTERING, DENSITY ESTIMATION, DIMENSIONALITY REDUCTION, GRAPH STRUCTURE DISCOVERY, IMITATION, NOISE FILTERING
- **REINFORCEMENT** → REINFORCEMENT

KNN ALGORITHM

CLASSIFIED

$$P(Y=C | X, D, K) = \frac{1}{K} \sum_{i \in N_K(X, D)} I(Y_i = C)$$

I IS INDICATOR

- LOOKS AT K NEAREST POINTS (TRAINING) CLOSEST TO TEST INPUT.
- RETURNS EMPIRICAL FRACTION AS ESTIMATE

• IF $K=1$ → VORONOI TESSELLATION

- **CURSE OF DIMENSIONALITY** → DIMENSIONS GROW LINEARLY, SIZE VOLUME OF SAMPLE SPACE GROWS EXPONENTIALLY

→ **SOLUTION!** LET'S MAKE ASSUMPTIONS ABOUT THE MODEL, INDUCTIVE BIAS (PARAMETERS IN PARAMETRIC MODEL)

CROSS VALIDATION

- **K-FOLDS**: $K-1$ AS TRAINING, 1 AS TEST, AVG CRISF. ERROR ON FOLDS, • $K=1$: LEAVE-ONE-OUT CV

INFORMATION THEORY 4 DUMMIES

ENTROPY: $H(X) = - \sum_k P(X=x_k) \cdot \log_2 P(X=x_k)$

BINARY/BERNOULLI ENTROPY: $H(X) = -[P \log_2 P + (1-P) \log_2 (1-P)]$

- **KL DIVERGENCE**: MEASURES DISSIMILARITY OF PROBABILITY DISTRIBUTIONS, N° OF EXTRA BITS DUE TO USING Q INSTEAD OF P

$$KL(P||Q) = \sum_n P_n \log \frac{P_n}{Q_n} = \sum_n P_n \log P_n - \sum_n P_n \log Q_n = -H(P) + H(P, Q)$$

$H(P, Q) = \text{CROSS-ENTROPY} = - \sum_n P_n \log Q_n$ → N° BITS TO ENCODE MSG FROM P IF WE USE MODEL Q

$KL(P||Q) \geq 0$

- **UNIFORM** IS DISCRETE DISTRIBUTION WITH MAX ENTROPY → **LAPLACE'S PRINCIPLE OF INSUFFICIENT REASON** CHOOSE UNIFORM WHEN NO PRIORS

- **MUTUAL INFORMATION** HOW SIMILAR JOINT DISTRIBUTION IS TO FACTORED DISTRIBUTION

• $I(X, Y) = KL(P(X, Y) || P(X)P(Y)) = \sum_x \sum_y P(x, y) \cdot \log \frac{P(x, y)}{P(x)P(y)}$ • $I(X, Y) = 0$ IFF $P(X, Y) = P(X)P(Y)$ = INDEPENDENCE

• $I(X, Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)$, $H(X|Y)$ **CONDITIONAL ENTROPY** = $\sum_x P(x) H(Y|X=x)$

- **MI** IS REDUCTION IN UNCERTAINTY FOR X AFTER HAVING OBSERVED Y , OR VICEVERSA

• **POINTWISE MI (PMI)** = FOR EVENTS, NOT DISTRIBUTIONS = $\log \frac{P(X, Y)}{P(X)} = \log \frac{P(Y|X)}{P(Y)}$

- **MUTUAL INFORMATION COEFFICIENT (MIC)**: FOR CONTINUOUS DISTRIBUTIONS, DISCRETE 2D GRIDS: BIN THE DATA, TRY DIFF. BIN SIZES

$m(X, Y) = \frac{\max_{\mathcal{G}} I(X(\mathcal{G}), Y(\mathcal{G}))}{\sqrt{\min_{\mathcal{G}} I(X(\mathcal{G}), Y(\mathcal{G}))}}$

MIC: $\max [m(X, Y)]$ FOR GUN SIZES | 0 = NO RELATIONSHIP
1 = NOISE-FREE RELATIONSHIP, NOT JUST LINEAR