

KDG Background Reading

- What are Polytopes?

a geometric object with flat sides.

- Kernel Density Estimation (KDE)

- a non-parametric way to estimate the PDF of a random variable.

- Let (x_1, x_2, \dots, x_n) be i.i.d. sampled drawn from a known distribution f at a given point x . We are interested in estimating the shape of the function f .

$$\hat{f}_h(x) = \frac{1}{n} \sum_{i=1}^n K_h(x-x_i) = \frac{1}{nh} \sum_{i=1}^n k\left(\frac{x-x_i}{h}\right)$$

where $K_h \rightarrow$ Kernel (a non-negative f_h)
 $h \rightarrow$ Bandwidth (a smoothing parameter)

- Performance Metrics

- ① Expected Calibration Error (ECE)

Measures the expected difference between accuracy and confidence by grouping all samples (size N) into K bins & calculating

$$ECE = \sum_{i=1}^k \frac{|B_i|}{N} |acc_i - conf_i|$$

acc_i and $conf_i$ are accuracy & average confidence in the i -th bin & $|B_i|$ is the number of samples in bin B_i

- * The pseudo-probabilities are class probabilities we get from the final layer of a NN.
The pseudo-probability of the predicted class generally over-estimates the actual probability of getting a correct answer.
- * If this over-estimation can be measured, it can be used to calibrate the NN such that its pseudo-probabilities would match the actual probability of the classes.
confidence

② Cohen's Kappa

$$\kappa = \frac{P_o - P_e}{1 - P_e}$$

$P_o \rightarrow$ relative observed agreement among raters

$P_e \rightarrow$ hypothetical probability of chance agreement.

- * Measures the agreement between 2 raters who each classify N times in to C mutually exclusive categories.
- * It is a quantitative measure of reliability for 2 raters that are rating the same thing, corrected for how often that the raters may agree by chance.

		rater 2	
		correct	incorrect
rater 1	correct	A	B
	incorrect	C	D

In A & D, the two raters are in agreement.

$$P_0 = P(\text{agreement}) = \frac{A+D}{A+B+C+D} = \frac{\text{agreements}}{\text{Total}}$$

Expected probability that both would say correct \Rightarrow

$$P(\text{correct}) = \frac{A+B}{A+B+C+D} \times \frac{A+C}{A+B+C+D}$$

Expected probability that both would say incorrect \Rightarrow

$$P(\text{incorrect}) = \frac{C+D}{A+B+C+D} \times \frac{B+D}{A+B+C+D}$$

P_e = overall random agreement probability that they agreed on either yes or no.

$$P_e = P(\text{correct}) + P(\text{incorrect})$$

Then $K = \frac{P_o - P_e}{1 - P_e}$

Understanding KDN

* Understanding polytopes formed by MLPs

- Partition & vote
- representation space learnt by DNNs is a partitioning of feature space into a union of convex polytopes.

classical statistical formulation of classification problem \Rightarrow

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n) \stackrel{iid}{\sim} f_{xy}$$

$$T_n = \{(x_i, y_i) | i \in \{1, \dots, n\}\} : \text{training data}$$

(x, y) : to-be classified test observation & its true label

$$x \in \mathbb{R}^d, y \in \{0, 1\}$$

Learn a classification rule $g_n = g(\cdot; \mathcal{T}_n)$ that maps feature vectors to class labels s.t. probability of misclassification $L(g_n) = P[g(x; \mathcal{T}_n) \neq Y | \mathcal{T}_n]$ is small.

Stone's Th^m for universally consistent classification:

- a successful classifier can be constructed by
- 1) partitioning the input space into cells depending on n , such that the # training points in each cell goes to infinity but slowly in n
 - 2) Estimating the posterior $\eta(x) = P[Y=1 | x=x]$ locally by voting based on the training class labels associated with the training feature vectors in cell $C(x) \subset \mathbb{R}^d$ in which the test observation falls.

Then $L(g_n) \rightarrow L^*$ almost surely for any F_{xy} where L^* is the Bayes optimal probability of misclassification.

Polytopes in Decision Forests

* given training data \mathcal{T}_n each tree t in RF constructs a partition $P_{n,t}$ by splitting the

input space.

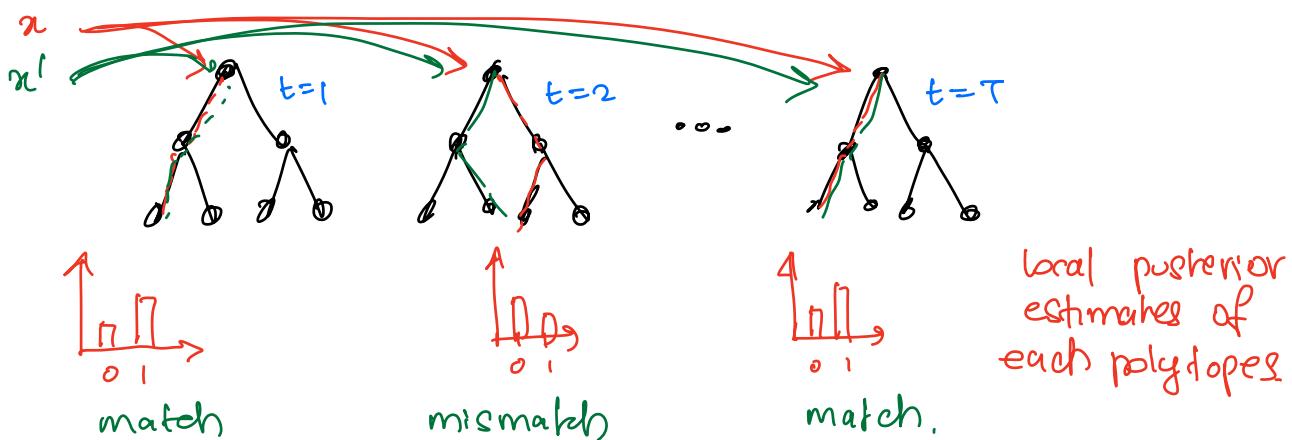
- * each tree results in a partition $P_{n,t} \in \mathbb{R}^d$
 (each leaf of each tree)

- * This partition admits a posterior estimate

$$\hat{n}_{n,t,j} = \left(\frac{1}{N_{n,t,j}}\right) \sum_{i: x_i \in C_{n,t,j}} I\{y_i = j\}$$

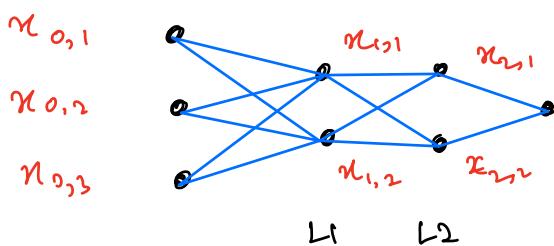
based on the class labels of training data

feature vectors in cell $C_{n,t,j}$



∴ The more the matched leaf ids are the more likely that x & x' belong in the same polytope

Polytopes in Deep Networks



$$x_{1,1} = f(w_{1,1} x_1 + w_{1,2,1} x_2 + w_{1,3,1} x_3 + b_{1,1})$$

$$x_{1,2} = f(\omega_{1,1,2} x_1 + \omega_{1,2,2} x_2 + \omega_{1,3,2} x_3 + b_{1,2})$$

$$\begin{bmatrix} x_{1,1} \\ x_{1,2} \end{bmatrix} = f \left\{ \begin{bmatrix} \downarrow & \omega_{1,1,1} & \omega_{1,2,1} & \omega_{1,3,1} \\ \omega_{1,1,2} & \omega_{1,2,2} & \omega_{1,3,2} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_{1,1} \\ b_{1,2} \end{bmatrix} \right\}$$

$$\therefore \begin{bmatrix} x_{2,1} \\ x_{2,2} \end{bmatrix} = f \left\{ \begin{bmatrix} \omega_{2,1,1} & \omega_{2,2,1} \\ \omega_{2,1,2} & \omega_{2,2,2} \end{bmatrix} \begin{bmatrix} x_{1,1} \\ x_{1,2} \end{bmatrix} + \begin{bmatrix} b_{2,1} \\ b_{2,2} \end{bmatrix} \right\}$$

$$x_{3,1} = \sigma (\omega_{3,1,1} x_{2,1} + \omega_{3,2,1} x_{2,2} + b_{3,1})$$

* all the samples that turn on/off the same ReLUs (nodes) ends up in the same polytope at the penultimate layer.

