

ST790 — Fall 2022

Imprecise-Probabilistic Foundations of Statistics

Ryan Martin

North Carolina State University

www4.stat.ncsu.edu/~rmartin

Week 12b

- Recap some of general ML details
- More classification (with imprecise probability)
- In particular:
 - Denoeux's *evidential neural network classifier*
 - conformal prediction and IMs
-

- Ingredients:
 - data, e.g., features X_i and labels Y_i
 - class \mathcal{F} of functions, hopefully $y \approx f^*(x)$ for some $f^* \in \mathcal{F}$
 - loss function, ℓ_f , to rate quality of f
- Note the absence of a statistical model...
- Training step boils down to “estimating” f via empirical risk minimization,¹ i.e., $\hat{f}_n = \arg \min_{f \in \mathcal{F}} n^{-1} \sum_{i=1}^n \ell_f(X_i, Y_i)$
- e.g., in classification,

$$\hat{f}_n(x) = \arg \max_{y \in \mathbb{Y}} \underbrace{\hat{\mathbb{P}}(Y = y \mid X = x)}_{\text{estimated predictive prob}}$$

- Use the trained \hat{f}_n to predict/classify new examples

¹Stochastic gradient descent is commonly used

- Huge \mathcal{F} and fancy algorithms/technology won't eliminate uncertainty, so UQ will always be relevant
- Two dominant statistical schools of thought?
 - *frequentist*
 - estimation is relatively easy
 - UQ isn't at all automatic
 - if it can be done, then likely inefficient ("model agnostic")
 - *Bayesian*
 - difficult to do (if one's being "honest")
 - UQ is an immediate by-product
 - meaningfulness of UQ wrt a single posterior dist?
- Imprecise-prob methods are a promising middle-ground...?

- **Last time:** *naive credal classifier*
- Extension/imprecise version of naive Bayes classifier
- Key features:
 - weaker prior assumptions (re: Manski)
 - able to classify examples to multiple labels
 - computationally tractable (thanks to IDM connection)
- **Today:** belief function/Dempster–Shafer approaches
 - *evidential neural net classifier*²
 - deep version, based on convolutional neural nets³

²Denoeux (*IEEE SMC* 2000)

³Tong, Xu, and Denoeux (*Neurocomputing* 2021)

Evidential classifier

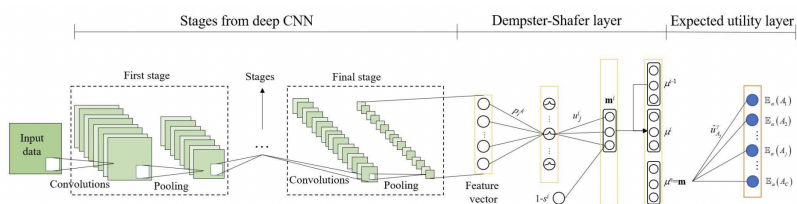
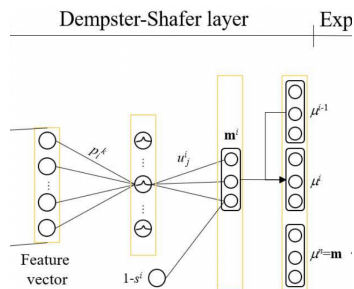


Figure 1: Architecture of an evidential deep-learning classifier.

- Multiple layers/stages:⁴
 - input gets processed through neural nets
 - neural net output gets converted into a mass/belief function
 - “expected utility” calculation for decision-making
- I'll focus exclusively here on the *DS* layer, which itself consists of several steps

⁴Screenshot from Tong, Xu, and Denoeux (2021)



- DS-layer consists of three steps:
 - distance-based support between input and references
 - mass function constructed for each reference
 - reference-specific mass functions combined via Dempster's rule
- Depends on parameters to-be-learned from training set

- Data consists of (X, Y) pairs
 - Y 's are labels
 - X 's represent images, chunks of text, etc
- Processing: $X \rightsquigarrow Z = Z(X) \in \mathbb{R}^q$
 - “ \rightsquigarrow ” designed to extract important characteristics
 - depends on the form of the input
 - depends on lots of to-be-learned parameters
- For our purposes, it suffices to proceed as if (Y, Z) is the available data, ignoring the processing
- Focus on mapping Z to a belief/mass function for Y

- Fix a set of *prototypes* p^1, \dots, p^R in \mathbb{R}^q
- Assign weight vectors α^r to each prototype:
 - $\beta_y^r := p^r$'s degree of membership to class y
 - with constraint $\sum_y \beta_y^r = 1$ for each r
 - these are to-be-learned parameters
- For a generic $z \in \mathbb{R}^q$, calculate the distance to prototypes

$$d^r = d^r(z) = \|z - p^r\|, \quad r = 1, \dots, R$$

- Factors influencing association between input z and label y
 - distance of z from prototypes
 - prototype membership degree with label y

- Given z , for each prototype $r = 1, \dots, R$, define a random set with mass function $m^r(\cdot)$,

$$m^r(\{y\}) = \alpha^r \beta_y^r \exp\{-\gamma^r (d^r)^2\}, \quad y \in \mathbb{Y}$$
$$m^r(\mathbb{Y}) = 1 - \alpha^r \exp\{-\gamma^r (d^r)^2\}$$

- α^r 's, β^r 's, and γ^r 's are to-be-learned parameters
- Easy to check that this is a genuine mass function

$$\sum_y m^r(\{y\}) + m^r(\mathbb{Y}) = 1$$

- Defines a belief/plausibility function on \mathbb{Y}
- This gives a prototype-specific quantification of uncertainty about which label y is associated with input z

- Goal is overall UQ, not a prototype-specific UQ
- Denoeux's idea:
 - since each prototype-specific UQ is a belief function
 - just combine m^1, \dots, m^R via Dempster's rule
- In symbols, $m = \bigoplus_{r=1}^R m^r$, Shafer's orthogonal sum
- Detailed formulas are messy⁵ and, hence, omitted
- Given this (z-dependent) mass function m , there are some options for carrying out classification:
 - naive strategy, $\arg \max_y m(\{y\})$
 - belief function yields a Choquet integral, so we can classify based on optimizing lower/upper expected utility⁶

⁵Denoeux uses some recursive relations...

⁶I'll cover general decision-theory details later

Evidential classifier: summary

- Process raw X through, say, a convolutional neural net
- Output Z and labels Y go into the DS-layer
- Returns a belief function on \mathbb{Y} for classification
- Parameters to be tuned in both the initial processing and the DS-layer, can be handled simultaneously via SGD
- For a new example, the feature X_{n+1} gets mapped to Z_{n+1} and then to a belief function on \mathbb{Y}
- Classification rule can be tailored so that set-valued classifications are made, more conservative, less error-prone

Conformal prediction (again)

- Roughly, Denoeux takes some existing machinery and uses the output to construct a belief function for UQ
- There are other ways to implement such a strategy
- *Conformal prediction*⁷ is a powerful method to leverage
- Recall:
 - set $Z_i = (X_i, Y_i)$ for $i = 1, \dots, n$
 - set $Z_{n+1} = (x, y)$ for generic (x, y)
 - define a non-conformity score $M(B, z)$
 - compute $\mu_i = M(\{Z_1, \dots, Z_{n+1}\} \setminus \{Z_i\}, Z_i)$, $i = 1, \dots, n + 1$
 - return $\pi_n(y | x) = (n + 1)^{-1} \sum_{i=1}^{n+1} \mathbf{1}\{\mu_i \geq \mu_{n+1}\}$
 - prediction region: $C_\alpha(Z^n; x) = \{y : \pi_n(y | x) > \alpha\}$

⁷Vovk et al's *Algorithmic Learning in a Random World*

- It turns out that conformal prediction can be related to (nested) random sets and belief functions⁸
- Conformal prediction's coverage reliability aligns with IM validity, so it's a special kind of belief function
- With finite \mathbb{Y} , the random set can be empty with non-zero probability; implies $\pi_n(y | x) < 1$ for all y
- This is bad — coherence & validity fail
- Two remedies:
 - condition on random set $\neq \emptyset$ (Dempster-style)
 - appropriately “stretch” random set⁹
- Both preserve validity, but latter is more efficient!

⁸Cella & M. (*IJAR* 2022), arXiv:2112.10234

⁹M. and Liu, *Inferential Models*, Ch. 5

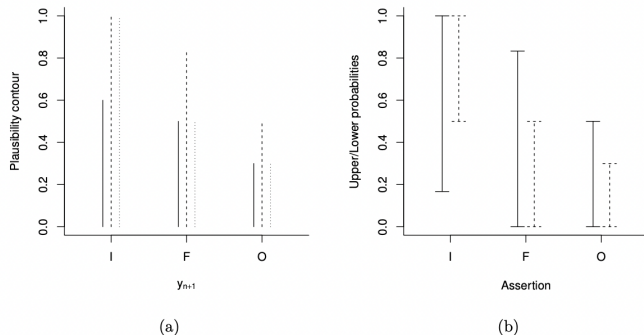


Figure 3: Panel (a): Plausibility contours in Equation (23), derived from an IM construction with no adjustment (solid lines), conditioning adjustment (dashed lines) and stretching adjustment (dotted lines). Panel (b): Upper and lower probabilities for the singleton assertions $\{I\}$, $\{F\}$ and $\{O\}$ derived from an IM construction with the conditioning adjustment (solid lines) and the stretching adjustment (dashed lines). These predictions are based on a new alligator of length $x_{n+1} = 2$ meters.

- I gave a high-level explanation of two imprecise-probability-based classification methods
 - evidential classifier: neural nets & Dempster–Shafer
 - IM classifier: conformal prediction & nested random sets
- Comparison:
 - conformal prediction can be used in conjunction with deep learning, but it's likely expensive¹⁰
 - evidential classifier (probably) doesn't have error rate controls
- Other methods...?

¹⁰ “Split” conformal prediction is faster, but validity is only approximate

- Prediction in regression
- i.e., supervised learning with continuous Y
- More IMs and conformal prediction
- Brand new stuff on *random fuzzy numbers*