

Last time :

□ uncertainty quantification I

↳ epistemic vs. aleatoric

↳ frequentist vs. Bayesian

Lecture 13

SP '26

Andrea Bajcsy

This Time:

□ uncertainty quantification II

↳ practical techniques! ensembles, conformal prediction

Summary & UQ Methods

	<u>Frequentism</u>	<u>Bayesianism</u>
✓	data-driven, easy comp. MLE $\hat{\theta} = \underset{\theta}{\operatorname{argmax}} P(y x; \theta)$	prior dist "jump starts" learning, posterior models uncertainty over θ params $P(\theta D) = \frac{P(\theta) \prod P(y x; \theta)}{P(D)}$
✗	misled by sampling noise, dataset size, etc.	computation usually very costly

frontiers: \Rightarrow beef this up!
ex. ② below

\Rightarrow approximate teis!
ex. ④ below

Practical Methods for UQ

Here are some keywords / techniques you may come across when looking for UQ methods for deep neural networks.

TODAY!

1) Deep Ensemble

\hookrightarrow [Lakshminarayanan et al. "Simple & Scalable UQ...". NeurIPS 2017]

2) Conformal Prediction (CP)

\hookrightarrow [Angelopoulos & Bates. "Gentle Intro to CP." arXiv 2021]

3) Gaussian Processes

\hookrightarrow [Rasmussen & Williams. "Gaussian Processes for ML." MIT Press 2006]

4) Laplace Approximation

\hookrightarrow [Daxberger et. al. "Laplace Redux". NeurIPS 2021]

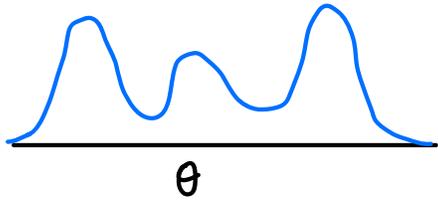
5) Monte Carlo Dropout

\hookrightarrow [Gal & Ghahramani. "Dropout as Bayesian Approx." ICML 2016]

DEEP ENSEMBLES

Recall how ideally, we would have uncertainty $P(\theta | \mathcal{D})$ over space of solutions given dataset. BUT, this is hard! (recall, full Bayesian approach)

Recall how we ideally had:

① dist. over all weights: $P(\theta | \mathcal{D})$ 

② predictions: $p(y | x, \mathcal{D}) = \int p(y | x, \theta) \cdot P(\theta | \mathcal{D}) d\theta$
 ↳ INTRACTABLE INTEGRAL ↳
 ex. $\theta \in \mathbb{R}^7 \text{ BILLION}$

Instead, deep ensembles approximate this empirically by replacing this complex integral w/ weighted sum of point estimates. Specifically:

① Train K independent neural networks via MLE.

This gives you $\theta_1^*, \theta_2^*, \dots, \theta_K^*$ different parameter values (b/c randomization + stochastic gradient descent).



Mathematically, you can think of it as mixture of Dirac's:

$$P(\theta | \mathcal{D}) \approx \frac{1}{K} \sum_{k=1}^K \delta_{\theta_k^*}(\theta)$$

↑ dirac delta function where all prob. mass is on θ_k^*

② To form predictions with uncertainty, we do:

$$p(y|x, \mathcal{D}) = \frac{1}{K} \sum_{k=1}^K p(y|x; \theta_k^*) \cdot \underbrace{p(\theta_k^* | \mathcal{D})}_{\uparrow}$$

since we just have K Dirac deltas, this is a uniformly weighted mixture model

⚠ In practice, we represent the ensemble's predictions as a Gaussian whose mean + variance are mean + variance of the ensemble:

$$p(y|x, \mathcal{D}) \approx \mathcal{N}(\mu_{\text{ens}}(x), \sigma_{\text{ens}}^2(x)) \quad \checkmark \text{ e.g. for regression problems}$$

where:

$$\mu_{\text{ens}}(x) = \frac{1}{K} \sum_{k=1}^K \mu_{\theta_k^*}(x) \quad \leftarrow \text{get mean predictions from all } K \text{ models}$$

$$\sigma_{\text{ens}}^2(x) = \frac{1}{K} \sum_{k=1}^K (\sigma_{\theta_k^*}^2 + \mu_{\theta_k^*}^2(x)) - \mu_{\text{ens}}^2(x)$$

b/c $\text{var}(x) := E[x^2] - E[x]^2$

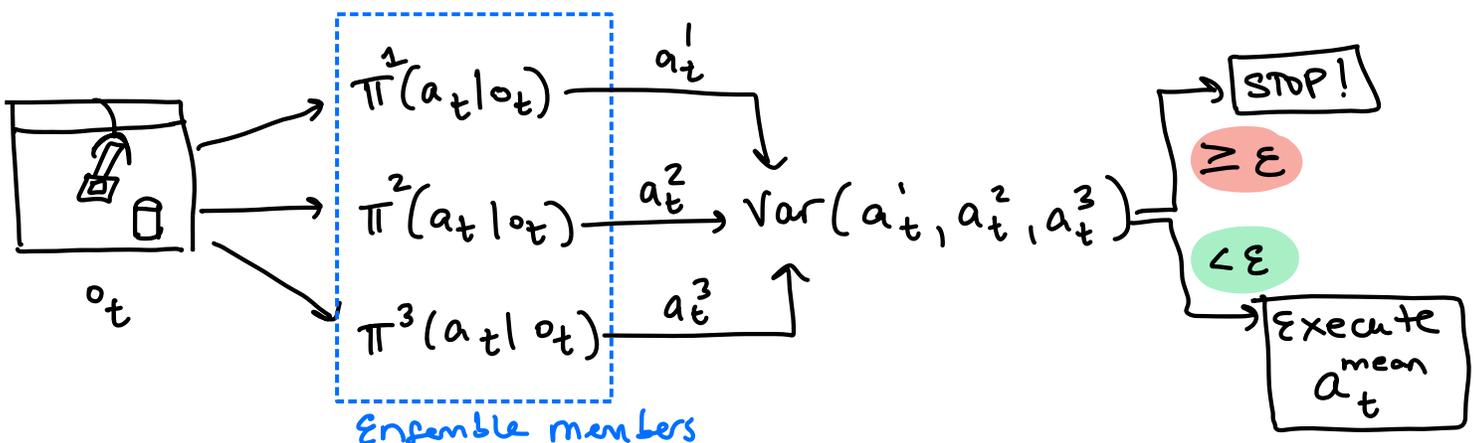
↑ variance of 1 ensemble member
 ↑ squared mean of 1 ensemble member
 ↑ mean of ensemble

What can we do with the ensemble?

↳ ex. if disagreement (i.e. variance $\sigma_{\text{ens}}^2(x)$) is too high

then this means the model is not confident

(i; maybe robot should stop, ask for help, etc.)



CONFORMAL PREDICTION

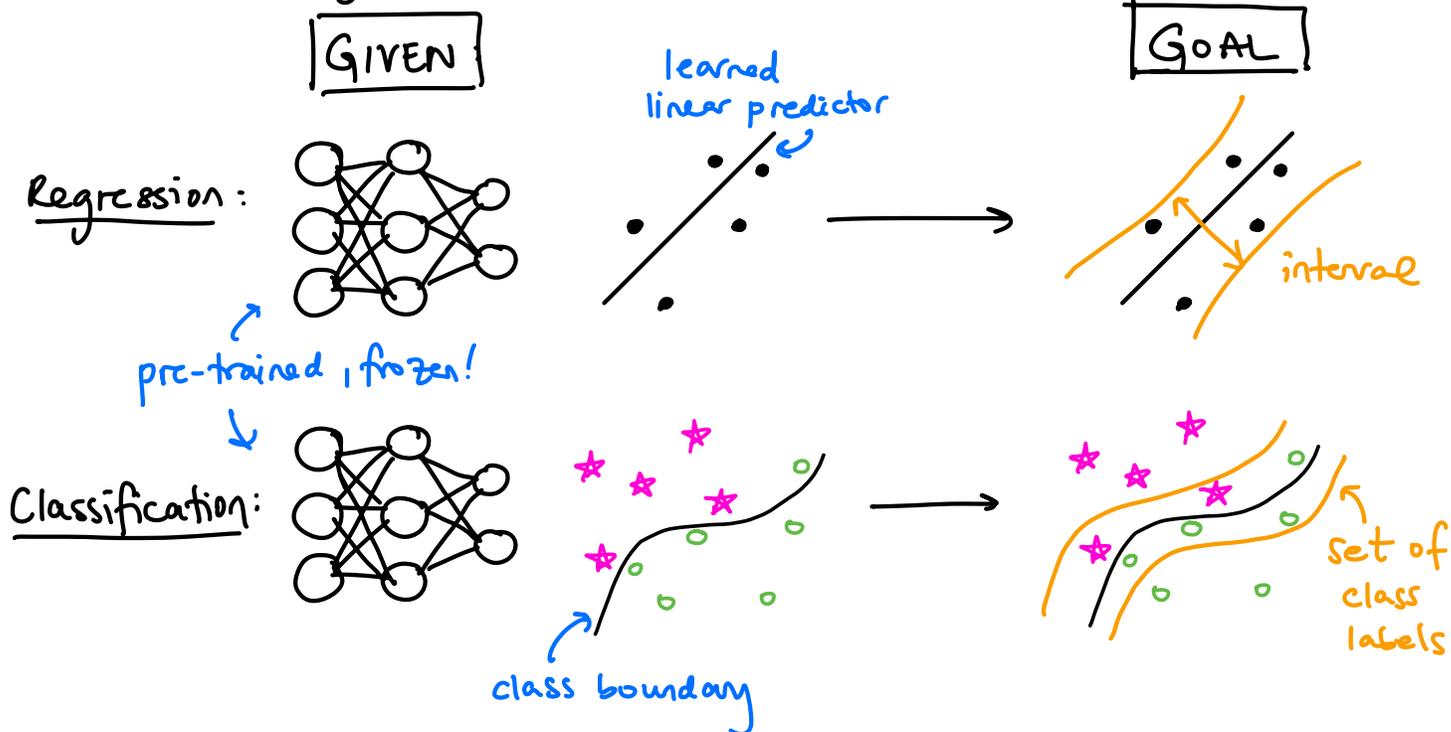
Ensembling has been very successful empirically. BUT, it requires you to train K separate "versions" of your predictive model. This is not possible when:

- ① you download pre-trained model
- ② you are training a foundation model

↳ ex. it took ≈ 3 months to train GPT-4....
now imagine having to do that K times!

So, what kind of uncertainty quantification technique could we use in such a situation?

conformal Prediction is a way to generate **prediction sets** for any (pre-trained) model.



Let's understand it via a image classification problem.

GOAL

Given : *fresh, unseen by model!* $\text{ex. image} \in \mathbb{R}^d$ $\text{ex. class} \in \mathcal{Y} = \{1, \dots, K\}$

- a calibration dataset $\{(X_i, Y_i)\}_{i=1}^N \sim \mathbb{P}$ i.i.d
- a model $\hat{f}_y(x)$ which estimates $\mathbb{P}(Y|X)$
 - ↑ inputs are x
 - ↑ outputs are over \mathcal{Y}
 - ↑ true distribution!
- a new input X_{N+1} *ex. new image*

Predict a SET $C(X_{N+1}) \subseteq \mathcal{Y}$ ← returns subset of label space
 function C depends on the model \hat{f} (ex. NN)

which contains the TRUE CLASS Y_{N+1} with high probability:

COVERAGE GUARANTEED:

$$\mathbb{P}(Y_{N+1} \in C(X_{N+1})) \geq 1 - \alpha$$

← user-defined error rate
 ↑ ex. $\alpha = 0.1$, then prob. that label is in set is 90%.
 new image w/ UNKNOWN label!



Figure 1: Prediction set examples on Imagenet. We show three progressively more difficult examples of the class fox squirrel and the prediction sets (i.e., $C(X_{\text{test}})$) generated by conformal prediction.

Desired Properties of Algo:

① Exact coverage

② small set size $|C(x)|$

③ "Adaptive" sets: **small sets** when **easy**; **big sets** when **hard**!

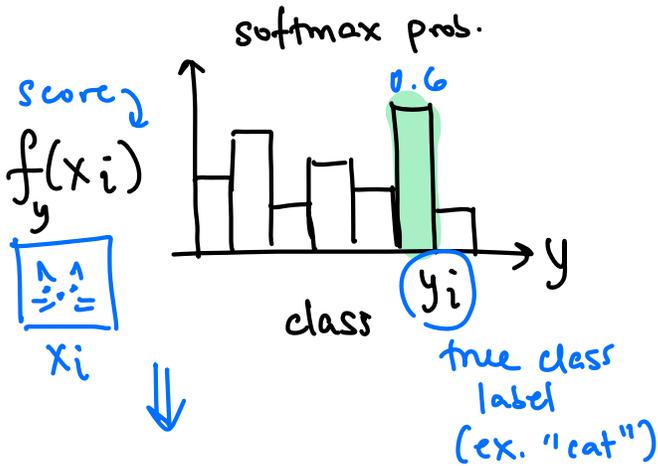
Ex.: naive way to get coverage:

$$C(X_{N+1}) = \begin{cases} \mathcal{Y} & \text{w.p. } 1 - \alpha \text{ (all labels)} \\ \emptyset & \text{w.p. } \alpha \text{ (no labels)} \end{cases}$$

⇒ get coverage but trivial/useless!

CONFORMAL PREDICTION ALGORITHM (Vovk et al., 2005)

① SCORE how poorly the model predicted the true class for all data in calibration set $\{(x_i, y_i)\}_{i=1}^N = \mathcal{D}_{\text{calib}}$



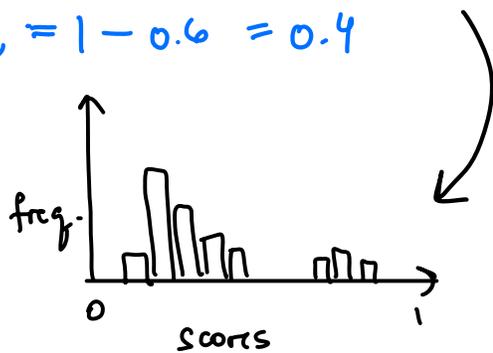
NON-CONFORMITY SCORE:

scalar $s_i := 1 - f_{y_i}(x_i)$

↳ high score \Rightarrow model is MORE WRONG

↳ low score \Rightarrow model is MORE RIGHT

$s_i = 1 - 0.6 = 0.4$



compute s_i for all $(x_i, y_i) \in \mathcal{D}_{\text{calib}}$
 You get this empirical distribution over the non-conformity scores

② compute the $1-\alpha$ empirical quantile of non-conformity scores

Recall: 0.5-quantile is just the median — it's the value of a random variable such that 50% of the probability mass is "below" or to the "left" of this value.

0.9-quantile is value where 90% of probability mass is "below" or "to the left" of this value

Defⁿ: p -quantile of distribution \mathbb{P} is the value x such that:

$$\mathbb{P}(X \leq x) \geq p.$$

future non-conformity score

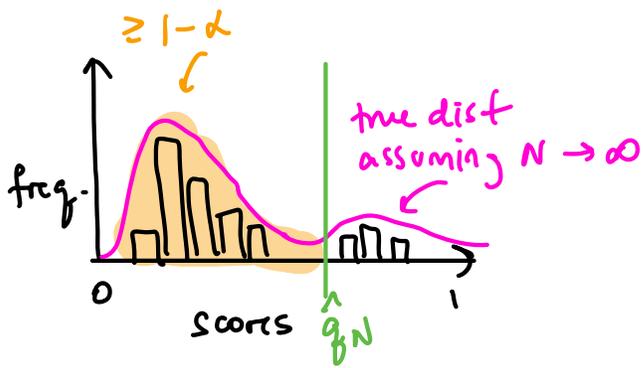
Theoretically: find $\hat{q}_{1-\alpha}$ s.t. $\mathbb{P}(s_{N+1} \leq \hat{q}_{1-\alpha}) \geq 1-\alpha$

$\alpha=0.1 \Rightarrow$

"find non-conformity score value large enough that only 10% of future values will ever exceed it (i.e. our error rate)"

! CHALLENGE: theoretical vs. empirical quantile. If we had

infinite data in our calibration set, then



\hat{q}_N would be perfect. But, we have finite data (ex. 100 pts.) - so, if we just took $(1-\alpha)$ -quantile, we would pick 90th smallest value out

of $N=100$ pts, we would not account for the $N+1=101$ th pt, and our coverage would only be $\frac{90}{101} \approx 89.1\% \neq 90\%$

Answer: do a finite sample adjustment, and instead of taking $(1-\alpha)$ -quantile, we take $\lceil (1-\alpha)(n+1) \rceil$ -quantile.

↳ intuition: by multiplying by $(n+1)$ we are accounting for how s_{N+1} may be ranked relative to existing calibration pts.



In our 100 data pt. example, we would pick $\lceil (1-\alpha)(n+1) \rceil = \lceil 0.9 * 101 \rceil = 91$ st smallest value out of 100 pts and our coverage is $\frac{91}{101} \approx 90.09\% \geq 90\%$

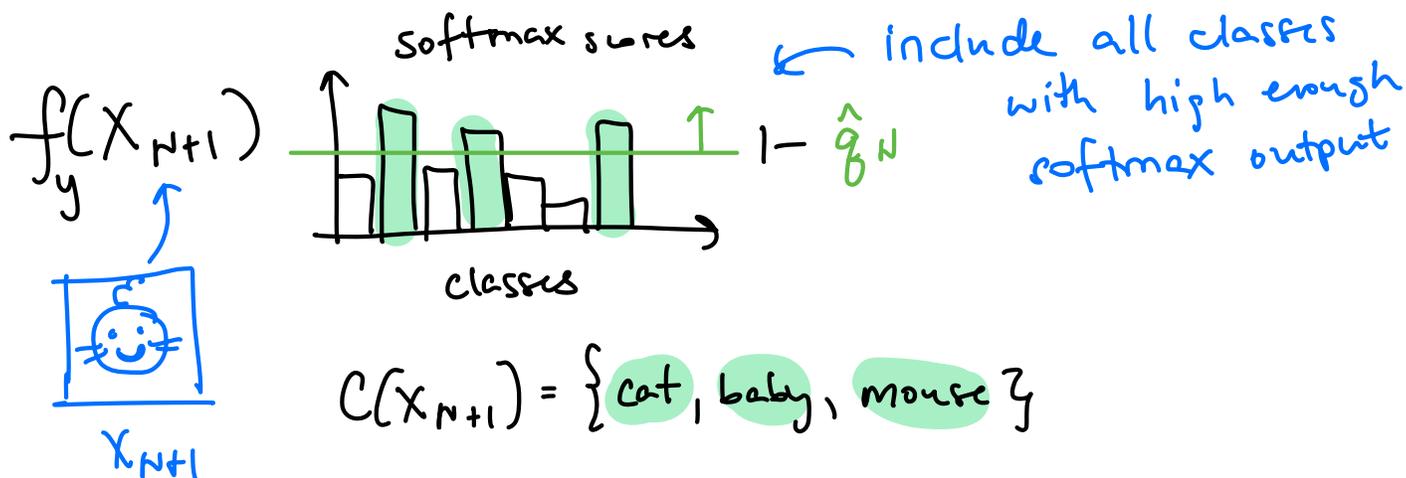
⇒ see proof for why this works in Appendix D of "A Gentle Intro to Conformal Prediction" and Prof. Ryan Tibshirani's notes (on class website)

! only works b/c data (x_i, y_i) is iid! (or, in general, exchangeable)

③ Form prediction sets using the \hat{q}_N we constructed before (with finite sample adjustment):

NOTE: you may see:
 $s(x_{N+1}, y) \leq \hat{q}_y$
 \downarrow
 $1 - \hat{f}(x_{N+1}) \leq \hat{q}_y$
 \downarrow
 $\hat{f}(x_{N+1}) \geq 1 - \hat{q}_y$

$$C(x_{N+1}) = \left\{ y : f_y(x_{N+1}) \geq 1 - \hat{q}_N \right\}$$



Theorem (Informal): For any model, dataset, α , N , you achieve the coverage guarantee:

$$1 - \alpha \leq \mathbb{P}(y_{N+1} \in C(x_{N+1})) \leq 1 - \alpha + \frac{1}{N+1}$$

ex. if $\alpha = 0.1$ and $N = 100$, then the size of the prediction set is not too conservative except for a factor of $\frac{1}{N+1}$. So $0.9\% \leq \mathbb{P}(\dots) \leq 0.91\%$.

Some Things to Consider:

- ① Exchangeable or i.i.d dataset is often NOT true in sequential decision-making domain!
- ② The design of non-conformity score REALLY matters! Key design question in CP
- ③ Think carefully about how to construct calibration data.