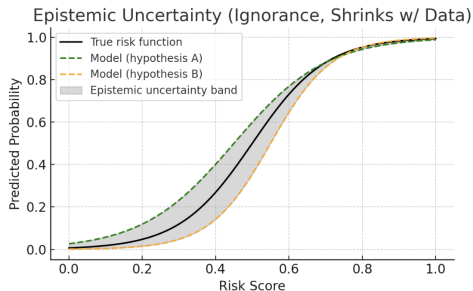# Rethinking Aleatoric and Epistemic Uncertainty
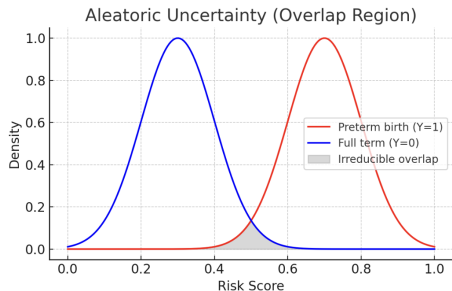
Freddie Bickford Smith, Jannik Kossen, Eleanor Trollope, Mark van der Wilk, Adam Foster, Tom Rainforth
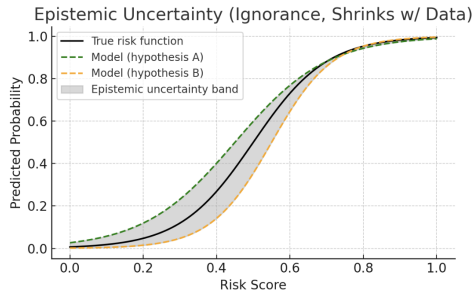
August 29, 2025

Presented by Mian Wei

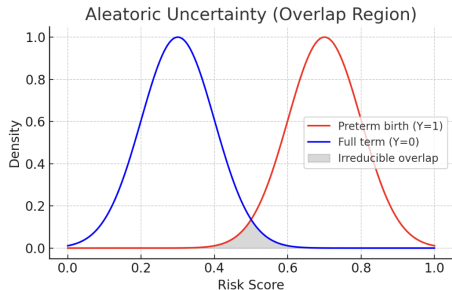# Aleatoric vs. Epistemic Uncertainty

We have a model to predict preterm birth, but there is predictive uncertainty. Should we use it cautiously, or should we collect more data and/or change the model?



- Aleatoric: Even with perfect features, some patients with identical profiles will deliver early, while others won't.
- Epistemic: Your dataset might be missing key variables or need better modeling.

# Aleatoric vs. Epistemic Uncertainty



- **Aleatoric**: unavoidable risk → best to hedge.
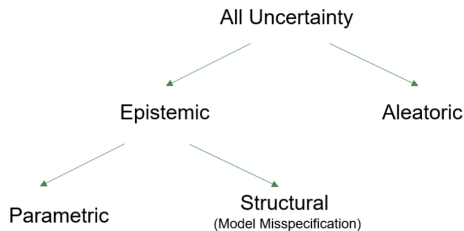- **Epistemic**: knowledge gap → best to reduce it.

# Uncertainty Decomposition



Figure: Adapted from [Che22]

- **Parametric**: uncertainty related to model parameter estimations under current model specification.
- **Structural**: the discrepancy between the assumed model specification and the true, unknown data-generating process

# An area of ongoing debate

- Model-based prediction vs. the true data-generating process
- Uncertainty on unseen data
- Uncertainty vs. prediction accuracy
- Different approaches to epistemic uncertainty: density-based, information-based, variance-based, ...

**Contribution of this paper:**

- A critique of the popular uncertainty decomposition view
- A new decision-based framework for uncertainty

# Notations

- an action, $a \in \mathcal{A}$
- a ground-truth variable, $z \in \mathcal{Z}$
- a loss function, $\ell : \mathcal{A} \times \mathcal{Z} \to \mathbb{R}$
- a policy $\pi \in \Pi$ that controls data generation
- some training data $y_{1:n} \sim p_{train}(y_{1:n}|\pi)$
- predictive model $p_n(z) = p(z; y_{1:n})$ or the predictive distribution $p_n(z) = \mathbb{E}_{p_n(\theta)}\left[p_n(z|\theta)\right]$, where $\theta \sim p_n(\theta) = p(\theta; y_{1:n})$ is a set of stochastic model parameters
- data generating process, $p_{train}(y_i|\pi(y_{<i}, y_{<i})$
- a ground-truth realization of $z$ or a reference distribution, $p_{eval}(z)$

# A Popular Decomposition View

## Usual formula

$$\underbrace{\text{EIG}_\theta}_{\text{"epistemic"}} = \underbrace{\text{H}\big[p_n(z)\big]}_{\text{"total"}} - \underbrace{\mathbb{E}_{p_n(\theta)}\big[\text{H}(p_n\big[z \mid \theta\big])\big]}_{\text{"aleatoric"}}$$

- H: Shannon entropy
- $\text{EIG}_\theta$: the expected information gain about $\theta$ from observing $z$

For finite $n$, epistemic and aleatoric uncertainty are only estimators of the true quantities and can be highly inaccurate.

# A Popular Decomposition View

**Aleatoric uncertainty**

| "captures noise inherent in the observations" $\mathrm{H}[p_{\mathrm{train}}(y_{1:n}|\pi)]$ or $\mathrm{H}[p_{\mathrm{eval}}(z)]$ | $\neq$ | "cannot be reduced even if more data were to be collected" $\mathrm{H}[p_\infty(z)]$ | $\neq$ | Expected parameter-conditional predictive entropy $\mathbb{E}_{p_n(\theta)}[\mathrm{H}[p_n(z|\theta)]]$ |

**Epistemic uncertainty**

| "uncertainty in the model parameters" $\mathrm{H}[p_n(\theta)]$ | $\neq$ | "can be explained away given enough data" $\mathrm{H}[p_n(z)] - \mathrm{H}[p_\infty(z)]$ | $\neq$ | Expected information gain in the model parameters $\mathrm{H}[p_n(z)] - \mathbb{E}_{p_n(\theta)}[\mathrm{H}[p_n(z|\theta)]]$ |

<span style="color:red">Model world $\neq$ Real world</span>

- Bayes-optimal action:

$$a_n^* = \arg \min_a \mathbb{E}_{p_n(z)}[\ell(a, z)]$$

- Loss-grounded uncertainty measure:

$$h[p_n(z)] = \mathbb{E}_{p_n(z)}[\ell(a_n^*, z)]$$

**Takeaway:** The choice of loss (and thus the uncertainty measure) depends on the decision problem at hand.

# Expected Uncertainty Reduction (EUR)

**Definition:**

$$UR_z(y_{1:m}^+) = h[p_n(z)] - h[p_{n+m}(z)].$$

$$EUR_z^{\text{true}}(\pi, m) = \mathbb{E}_{p_{\text{train}}(y_{1:m}^+|\pi)}[\, UR_z(y_{1:m}^+)\,].$$

# Expected Uncertainty Reduction (EUR)

**Definition:**

$$UR_z(y_{1:m}^+) = h[p_n(z)] - h[p_{n+m}(z)].$$

$$EUR_z^{\text{true}}(\pi, m) = \mathbb{E}_{p_{\text{train}}(y_{1:m}^+|\pi)}[\, UR_z(y_{1:m}^+)\,].$$

**As** $m \to \infty$, the decomposition:

$$h[p_n(z)] = \underbrace{EUR_z^{\text{true}}(\pi, \infty)}_{\text{Reducible}} + \underbrace{\mathbb{E}_{p_{\text{train}}(y_{1:m}^+|\pi)}[h[p_\infty(z)]]}_{\text{Irreducible}}.$$

# Expected Uncertainty Reduction (EUR)

**Definition:**

$$UR_z(y_{1:m}^+) = h[p_n(z)] - h[p_{n+m}(z)].$$

$$EUR_z^{\text{true}}(\pi, m) = \mathbb{E}_{p_{\text{train}}(y_{1:m}^+|\pi)}[\, UR_z(y_{1:m}^+) \,].$$

**As** $m \to \infty$, the decomposition:

$$h[p_n(z)] = \underbrace{EUR_z^{\text{true}}(\pi, \infty)}_{\text{Reducible}} + \underbrace{\mathbb{E}_{p_{\text{train}}(y_{1:m}^+|\pi)}[h[p_\infty(z)]]}_{\text{Irreducible}}.$$

Compared with popular split:

$$H[p_n(z)] = \underbrace{\mathbb{E}_{p_n(\theta)}[H(p_n(z \mid \theta))]}_{\text{aleatoric}} + \underbrace{H[p_n(z)] - \mathbb{E}_{p_n(\theta)}[H(p_n(z \mid \theta))]}_{\text{BALD/epistemic}}$$

# Expected Uncertainty Reduction (EUR)

**Definition:**

$$UR_z(y^+_{1:m}) = h[p_n(z)] - h[p_{n+m}(z)].$$

$$EUR_z^{\text{true}}(\pi, m) = \mathbb{E}_{p_{\text{train}}(y^+_{1:m}|\pi)}[\, UR_z(y^+_{1:m}) \,].$$

**As** $m \to \infty$, the decomposition:

$$h[p_n(z)] = \underbrace{EUR_z^{\text{true}}(\pi, \infty)}_{\text{Reducible}} + \underbrace{\mathbb{E}_{p_{\text{train}}(y^+_{1:m}|\pi)}[h[p_\infty(z)]]}_{\text{Irreducible}}.$$

Compared with popular split:

$$H[p_n(z)] = \underbrace{\mathbb{E}_{p_n(\theta)}[H(p_n(z \mid \theta))]}_{\text{aleatoric}} + \underbrace{H[p_n(z)] - \mathbb{E}_{p_n(\theta)}[H(p_n(z \mid \theta))]}_{\text{BALD/epistemic}}$$

- Paper's: decision/loss grounded, any learner, depends on data process.
- BALD: an *estimator*, not a universal decomposition.

# EUR in Practice

**Problem:** We cannot access the true data-generating process or infinite data.

So we approximate with:

- Use model-based simulator $p_n(y_{1:m}^+ \mid \pi')$ instead of true $p_{\text{train}}$.
- Use approximate update $q_{n+m}(z)$ instead of true $p_{n+m}(z)$.

**Problem:** We cannot access the true data-generating process or infinite data.

So we approximate with:

- Use model-based simulator $p_n(y_{1:m}^+ \mid \pi')$ instead of true $p_{\text{train}}$.
- Use approximate update $q_{n+m}(z)$ instead of true $p_{n+m}(z)$.

**Estimator:**

$$EUR_z^{\text{est}}(\pi', m) = h[p_n(z)] - \mathbb{E}_{p_n(y_{1:m}^+ \mid \pi')}\big[h[q_{n+m}(z)]\big].$$

# EUR in Practice

**Problem:** We cannot access the true data-generating process or infinite data.

So we approximate with:

- Use model-based simulator $p_n(y_{1:m}^+ \mid \pi')$ instead of true $p_{\text{train}}$.
- Use approximate update $q_{n+m}(z)$ instead of true $p_{n+m}(z)$.

**Estimator:**

$$EUR_z^{\text{est}}(\pi', m) = h[p_n(z)] - \mathbb{E}_{p_n(y_{1:m}^+|\pi')}\big[h[q_{n+m}(z)]\big].$$

**Sources of error:**

1. Simulator mismatch ($p_n$ vs $p_{\text{train}}$).
2. Update approximation ($q_{n+m}$ vs $p_{n+m}$).

**Predictive Uncertainty**

- $h[p_n(z)] = \mathbb{E}_{p_n(z)}[\ell(a_n^*, z)]$
- How uncertain *my model* thinks the future is
- Subjective, depends on $p_n(z)$

**Predictive Uncertainty**

- $h[p_n(z)] = \mathbb{E}_{p_n(z)}[\ell(a_n^*, z)]$
- How uncertain *my model* thinks the future is
- Subjective, depends on $p_n(z)$

**Predictive Performance**

- $\text{Perf}(p_n) = \mathbb{E}_{p_{\text{eval}}}[\ell(a_n^*, z)]$
- How good the predictions are compared with reality
- Requires $p_{\text{eval}}(z)$

## Predictive Uncertainty

- $h[p_n(z)] = \mathbb{E}_{p_n(z)}[\ell(a_n^*, z)]$
- How uncertain *my model* thinks the future is
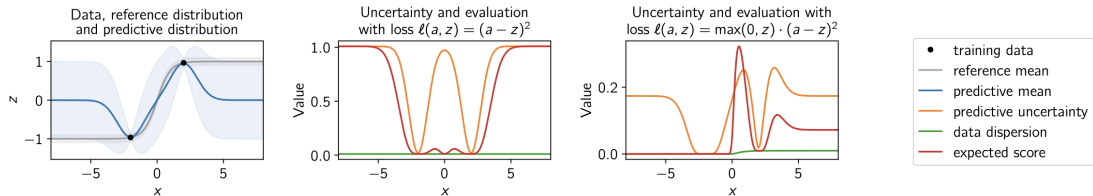- Subjective, depends on $p_n(z)$

## Predictive Performance

- $\text{Perf}(p_n) = \mathbb{E}_{p_{\text{eval}}}[\ell(a_n^*, z)]$
- How good the predictions are compared with reality
- Requires $p_{\text{eval}}(z)$

## Data Dispersion

- Dispersion $=$ entropy/variance of $p_{\text{eval}}(z)$
- How random the world really is, regardless of the model
- World-based, not model-based

Data, reference distribution and predictive distribution

Uncertainty and evaluation with loss $\ell(a, z) = (a - z)^2$

Uncertainty and evaluation with loss $\ell(a, z) = \max(0, z) \cdot (a - z)^2$

- training data
- reference mean
- predictive mean
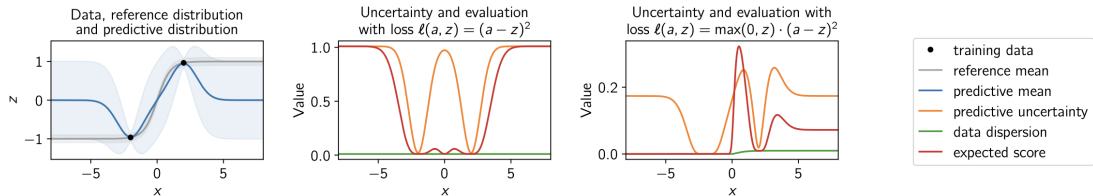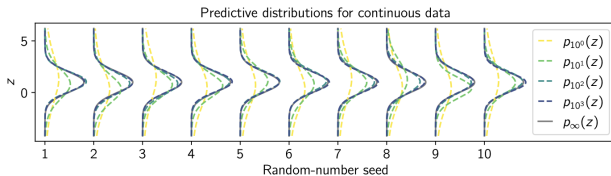- predictive uncertainty
- data dispersion
- expected score

**Prop 1:** Bayes estimator under quadratic loss = the posterior mean.

**Prop 2:** $h[p_n(z)]$ is the Bayes estimator of expected performance under $p_{\text{eval}}$.

**Prop 3:** $\mathbb{E}_{p_n(\theta)}[h[p_n(z \mid \theta)]]$ is the Bayes estimator of data dispersion $h[p_{\text{eval}}(z)]$.

# Some Concepts



Data, reference distribution and predictive distribution

Uncertainty and evaluation with loss $\ell(a, z) = (a - z)^2$

Uncertainty and evaluation with loss $\ell(a, z) = \max(0, z) \cdot (a - z)^2$

- training data
- reference mean
- predictive mean
- predictive uncertainty
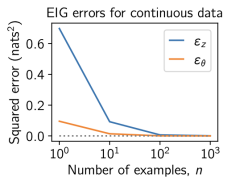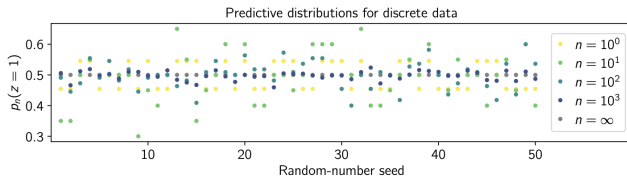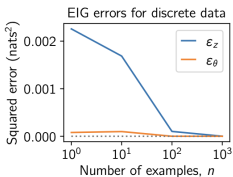- data dispersion
- expected score

**Prop 1:** Bayes estimator under quadratic loss = the posterior mean.

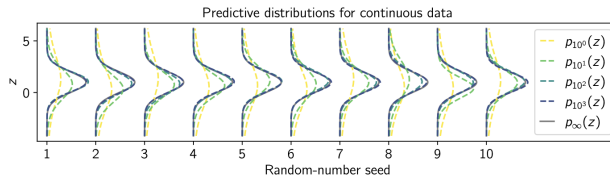**Prop 2:** $h[p_n(z)]$ is the Bayes estimator of expected performance under $p_{\text{eval}}$.
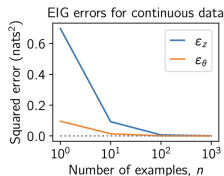
**Prop 3:** $\mathbb{E}_{p_n(\theta)}[h[p_n(z \mid \theta)]]$ is the Bayes estimator of data dispersion $h[p_{\text{eval}}(z)]$.

**Takeaway:** Model-based uncertainty $\neq$ truth; only *estimators* of performance/dispersion.
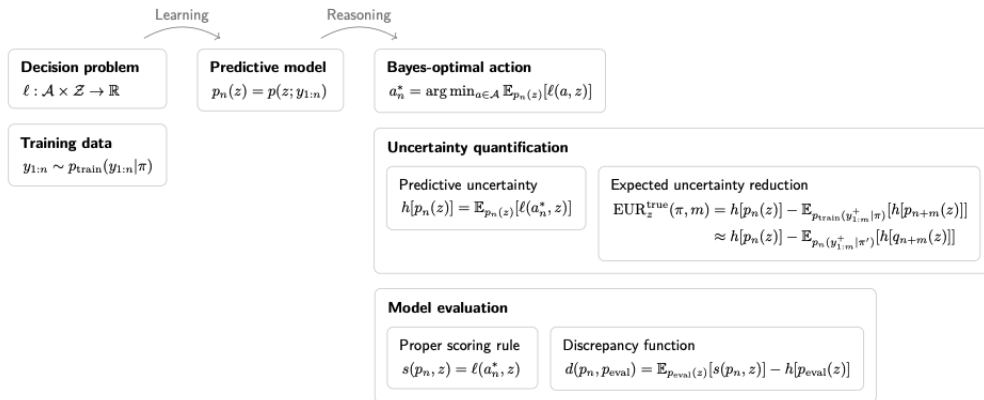
**Takeaway:**

- BALD is not the "epistemic uncertainty" truth. It's an estimator, and often poor for long-run reducibility.
- But good proxy for *short-run parameter IG* → explains success in active learning.

# Decision-Theoretic Framework



Learning     Reasoning

**Decision problem**
$\ell : \mathcal{A} \times \mathcal{Z} \to \mathbb{R}$

**Predictive model**
$p_n(z) = p(z; y_{1:n})$

**Bayes-optimal action**
$a_n^\star = \arg\min_{a \in \mathcal{A}} \mathbb{E}_{p_n(z)}[\ell(a, z)]$

**Training data**
$y_{1:n} \sim p_{\text{train}}(y_{1:n} | \pi)$

**Uncertainty quantification**

Predictive uncertainty
$h[p_n(z)] = \mathbb{E}_{p_n(z)}[\ell(a_n^\star, z)]$

Expected uncertainty reduction
$\text{EUR}_z^{\text{true}}(\pi, m) = h[p_n(z)] - \mathbb{E}_{p_{\text{train}}(y_{1:m}^+ | \pi)}[h[p_{n+m}(z)]]$
$\approx h[p_n(z)] - \mathbb{E}_{p_n(y_{1:m}^+ | \pi')}[h[q_{n+m}(z)]]$

**Model evaluation**

Proper scoring rule
$s(p_n, z) = \ell(a_n^\star, z)$

Discrepancy function
$d(p_n, p_{\text{eval}}) = \mathbb{E}_{p_{\text{eval}}(z)}[s(p_n, z)] - h[p_{\text{eval}}(z)]$

# Accommodation

**Is it worth reading?** Maybe.

1. Uncertainty is decision-specific, not one-size-fits-all.
2. Decomposition: reducible vs. irreducible (depends on DGS, not just model).
3. Model-based quantities are *estimators*, not ground truths.
4. BALD works in practice by estimating short-run parameter IG.

**Cons:**

1. The authors claim their decomposition is better, but the argument is unconvincing, as there are no experiments, no rigorous proof, and no empirical validation.
2. More like a conceptual critique + framework clarification paper.

Thank you! Any questions?

Shuo Chen.
Introduction and exemplars of uncertainty decomposition.
*arXiv preprint arXiv:2211.15475*, 2022.