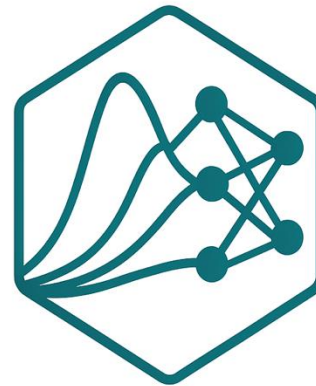


Lecture 1: Introduction

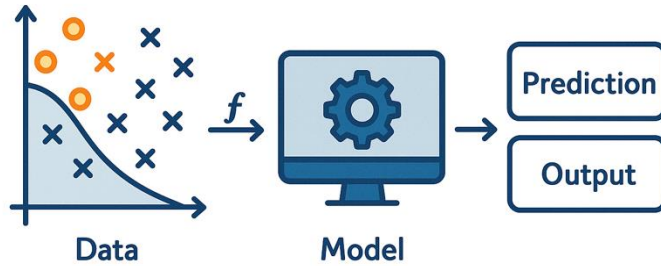
Krikamol Muandet
17 October 2025



IPML

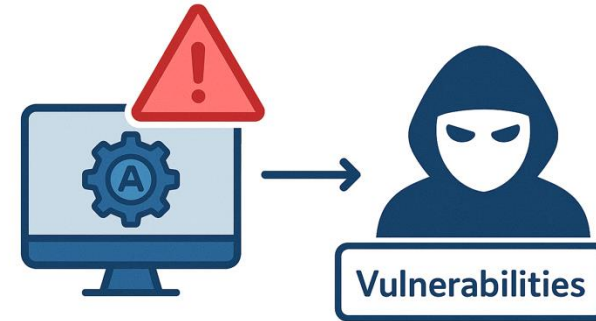
IMPRECISE
PROBABILISTIC
MACHINE LEARNING

Outline



- **Statistical Machine Learning**

- Classification problems
- Regression problems
- Unsupervised learning
- Generative modeling
- Semi-supervised learning
- Reinforcement learning



- **Pitfalls of Machine Learning**

- Overconfident prediction
- OOD generalisation
- Algorithmic biases
- Adversarial robustness
- Trustworthiness
- AI safety and misalignment



Statistical Machine Learning

The warrior of modern-day computation

Statistical Machine Learning

- A **hypothesis space** \mathcal{H} , a data set \mathcal{D} of finite **observations**, and a **loss function** ℓ .
- A **statistical learning** aims to find the “best” hypothesis in \mathcal{H} based on the observations in \mathcal{D} :

$$\hat{h} \in \arg \min_{h \in \mathcal{H}} \ell_h(\mathcal{D})$$

- How does \hat{h} compare to the **optimal hypothesis** h^* as the number of observations approaches infinity?

Statistical Machine Learning

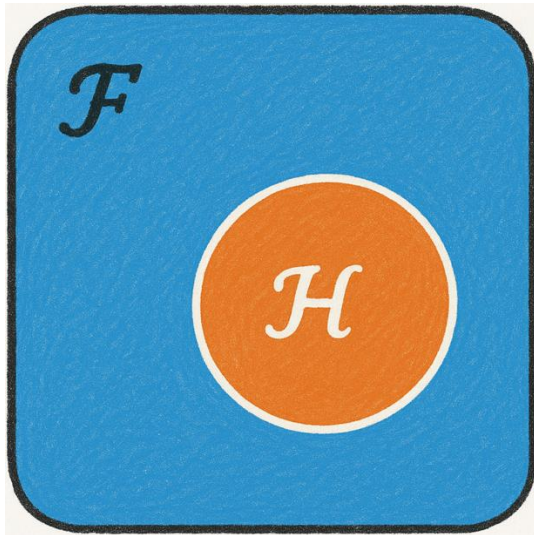
- An **input space** \mathcal{X} and **target space** \mathcal{Y}
- Consider a finite hypothesis class $\mathcal{H}_m = \{h_1, h_2, \dots, h_m\}$, $h_i: \mathcal{X} \rightarrow \mathcal{Y}$
- A data set $\mathcal{D} = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$, $(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}$
- A real-valued loss function $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow R_+$

$$\hat{h} \in \arg \min_{h \in \mathcal{H}_m} \hat{\mathbf{R}}(h) =: \frac{1}{n} \sum_{i=1}^n \ell(y_i, h(x_i))$$
$$h^* = \arg \min_{h \in \mathcal{H}_m} \mathbf{R}(h) =: \int \ell(Y, h(X)) dP(X, Y)$$

Learning here amounts to enumerating the hypothesis class and selecting the hypothesis with the lowest average loss.

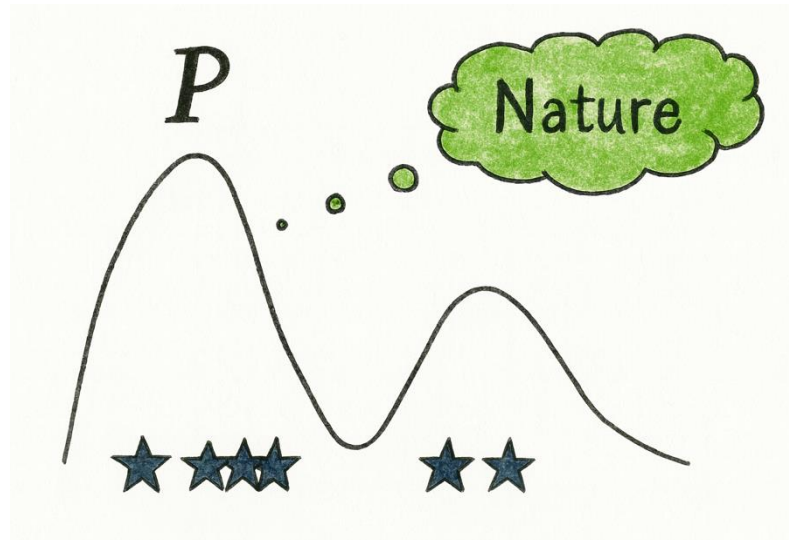
Statistical Machine Learning

Hypothesis Space



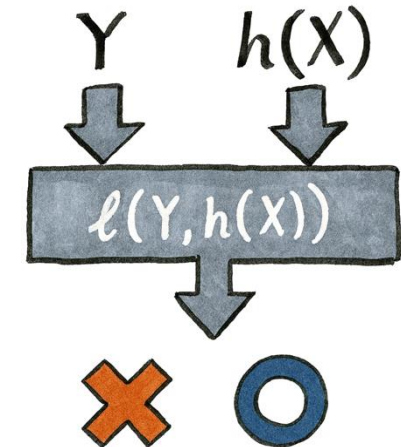
A **feasible** hypothesis space \mathcal{H} is a subset of all **conceivable** hypotheses \mathcal{F}

Data



Data are drawn **IID** from some unknown but fixed distribution $(X, Y) \sim P(X, Y)$

Loss Function



The loss function ℓ is **task-specific** and **quantifies the correctness** of the prediction

Probably Approximately Correct (PAC)

- Let $c: \mathcal{X} \rightarrow \{0,1\}$ be **target concept** to learn and \mathcal{C} a **concept class** comprising target concepts c .
- A **learning algorithm** \mathcal{A} receives sample S and selects a hypothesis h_S from \mathcal{H} approximating the target c .

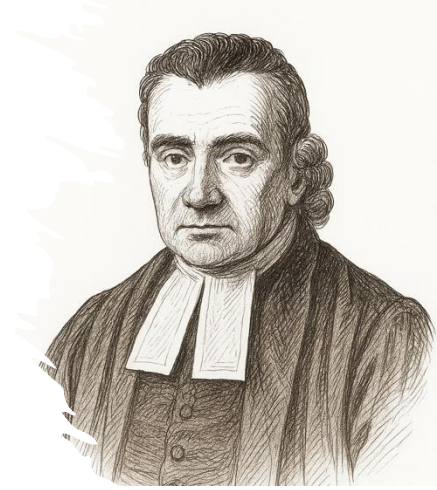
A hypothesis class \mathcal{H} is **PAC-learnable** if there exists a learning algorithm \mathcal{A} such that for all concepts $c \in \mathcal{C}$, $\epsilon > 0$, $\delta > 0$, and all distributions P ,

$$\Pr_{S \sim P^m} [R(h_S) \leq \epsilon] \geq 1 - \delta$$

for sample S of size $m = \text{poly}(1/\epsilon, 1/\delta)$ for a fixed polynomial.

Valiant (1984)

Learning as Belief Update

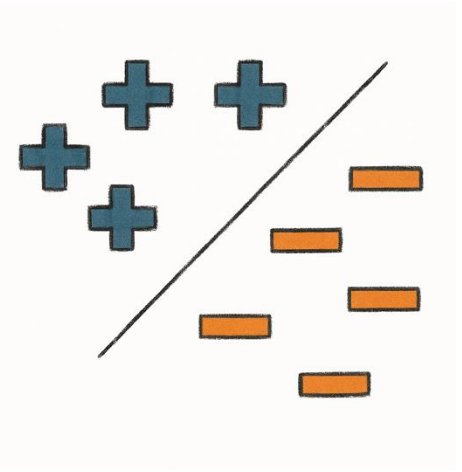


- Let Θ be a **parameter** space and \mathcal{D} an observed **evidence**
 - $P(\theta)$ – A **prior** probability over the parameter space Θ
 - $P(\mathcal{D} | \theta)$ – A **likelihood** function of \mathcal{D} given θ
 - $P(\theta | \mathcal{D})$ – A **posterior** probability of θ given \mathcal{D}
- The **Bayes** rule relates *posterior* to *likelihood* and *prior*:

$$P(\theta | \mathcal{D}) \propto P(\mathcal{D} | \theta)P(\theta)$$

- The most popular update rule in the precise setting, known for its **many desirable properties**, though not the only possible choice in the imprecise setting.

Classification Problems



- $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$, $\mathcal{X} \subset \mathbb{R}^d$ and $\mathcal{Y} = \{+1, -1\}$
- A **0-1** loss function $\ell(y, \hat{y}) = \mathbb{1}_{y \neq \hat{y}}$ (misclassification loss)

$$\hat{h} \in \arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{y_i \neq h(\mathbf{x}_i)}, \quad h_{\mathbf{w}, b}(\mathbf{x}) = \begin{cases} +1, & \langle \mathbf{w}, \mathbf{x} \rangle + b > 0 \\ -1, & \langle \mathbf{w}, \mathbf{x} \rangle + b \leq 0 \end{cases}$$

- For K classes, a **cross-entropy** loss $\ell(y, \hat{y}) = -\sum_{k=1}^K y^{(k)} \log \hat{y}^{(k)}$ where y is a **one-hot encoding** and \hat{y} is a **probability score**

$$y = (0, 0, \mathbf{1}, 0), \quad \hat{y} = (0.02, 0.1, \mathbf{0.7}, 0.18)$$

Kullback-Leibler (KL) Divergence

- The **Shannon entropy** $H(P) = -E_{x \sim P}[\log P(x)]$.
- For two probability distribution $P(x)$ and $Q(x)$:

$$D_{KL}(P \parallel Q) = E_{x \sim P} \left[\log \frac{P(x)}{Q(x)} \right] = E_{x \sim P} [\log P(x) - \log Q(x)]$$

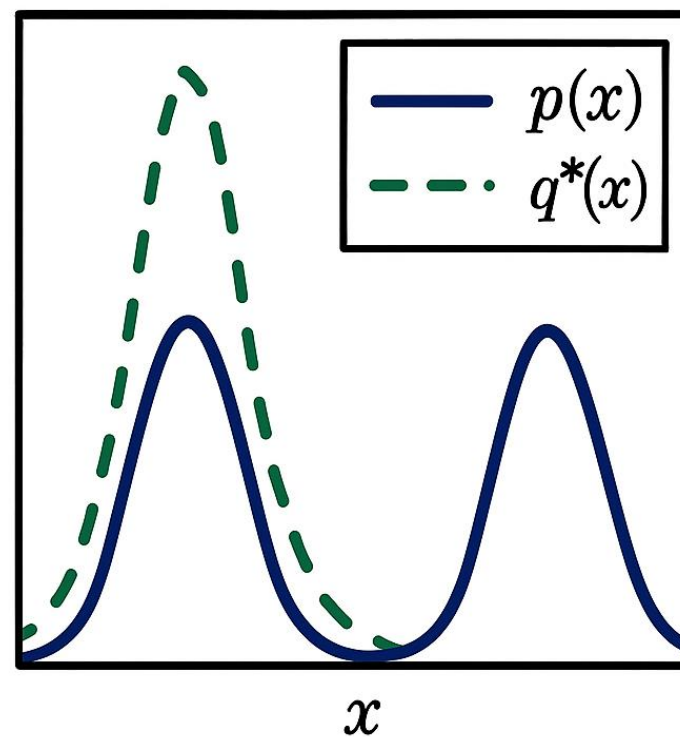
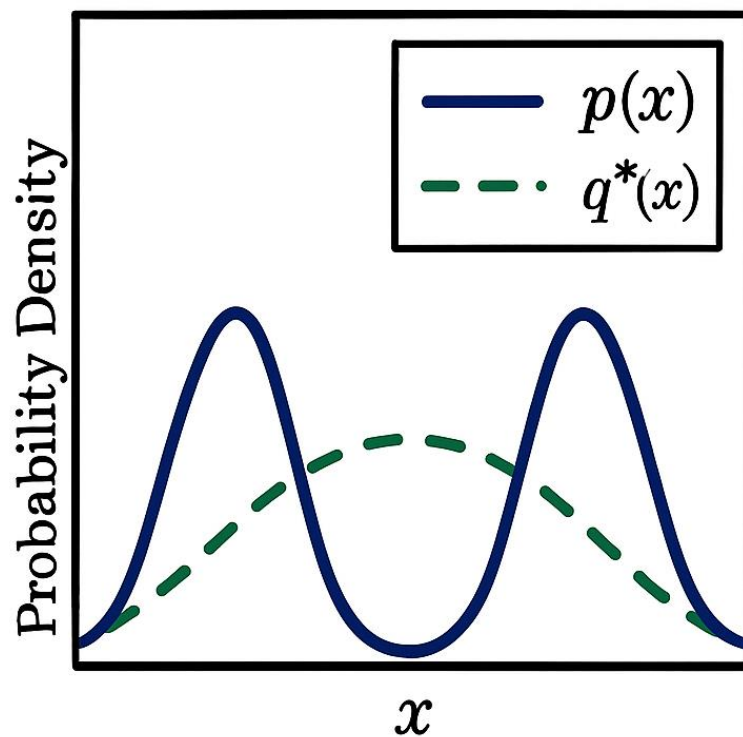
- Important properties: (1) $D_{KL}(P \parallel Q) \geq 0$ (2) $D_{KL}(P \parallel P) = 0$
(3) $D_{KL}(P \parallel Q) \neq D_{KL}(Q \parallel P)$
- The cross-entropy can be expressed as

$$H(P, Q) = H(P) + D_{KL}(P \parallel Q) = -E_{x \sim P}[\log Q(x)]$$

Kullback-Leibler (KL) Divergence

$$q^* = \arg \min_q D_{KL}(p \parallel q)$$

$$q^* = \arg \min_q D_{KL}(q \parallel p)$$



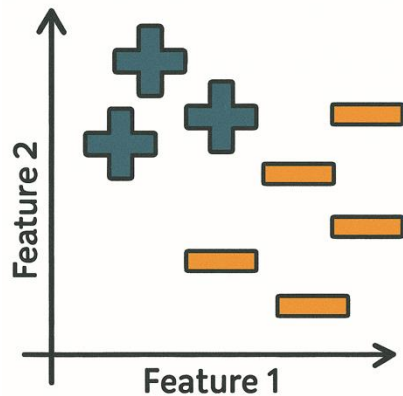
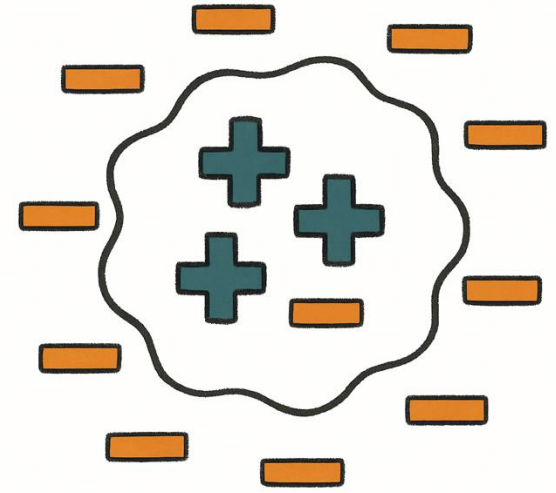
mode-covering vs. mode-seeking

Classification Problems

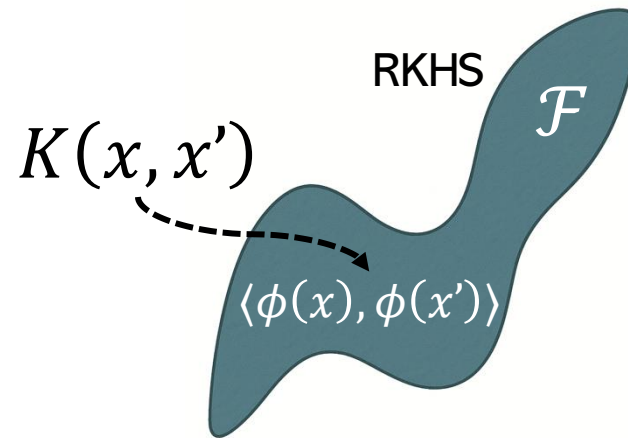
- $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}, \mathcal{X} \subset \mathbb{R}^d$ and $\mathcal{Y} = \{+1, -1\}$
- A hypothesis class

$$\mathcal{H} = \{h(\mathbf{x}) : h(\mathbf{x}) = \langle \mathbf{w}, \phi(\mathbf{x}) \rangle_{\mathcal{F}} + b\}$$

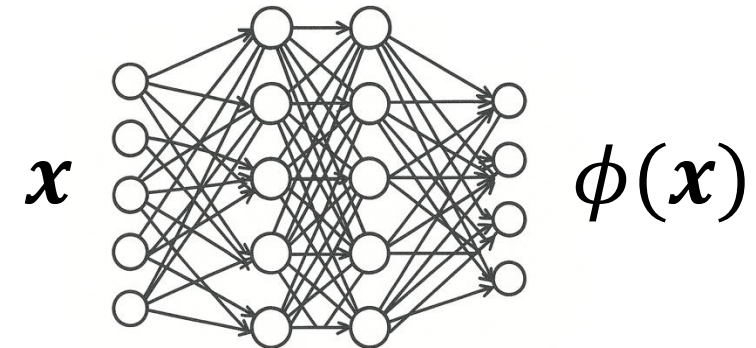
for some **non-linear feature map** $\phi: \mathcal{X} \rightarrow \mathcal{F}$.



Hand-crafted Features



Implicit Features



Explicit Features

Regression Problems

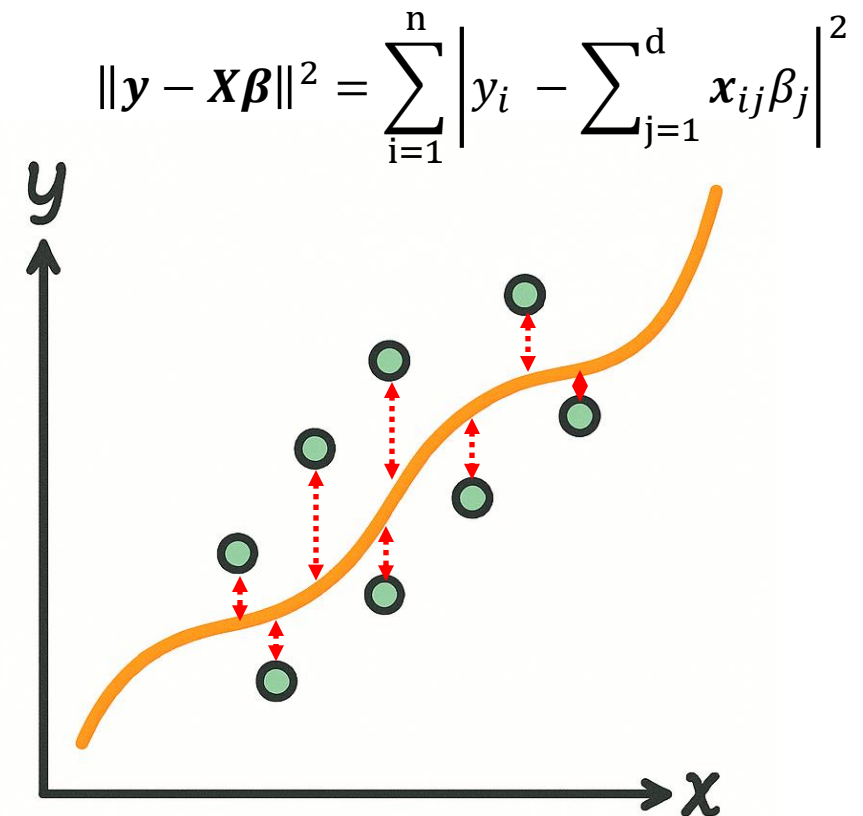
- $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$, $\mathcal{X} \subset \mathbb{R}^d, \mathcal{Y} \subset \mathbb{R}$
- A **square loss** function $\ell(y, \hat{y}) = (y - \hat{y})^2$
- A hypothesis class

$$\mathcal{H} = \{h(x) : h(x) = \langle w, \phi(x) \rangle_{\mathcal{F}} + b\}$$

- **Least Square (LS)**: $h_{\beta}(x) = \beta^T x$

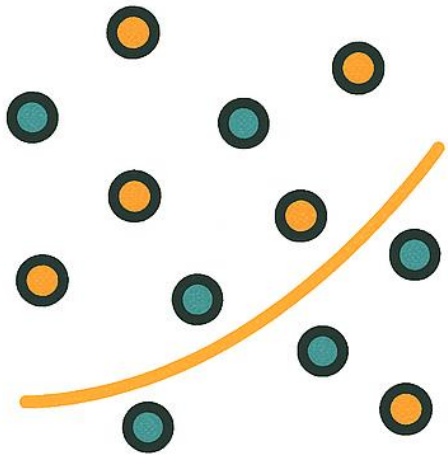
$$\hat{\beta} = \operatorname{argmin}_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|^2 = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

where $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T$, $\mathbf{y} = (y_1, \dots, y_n)^T$



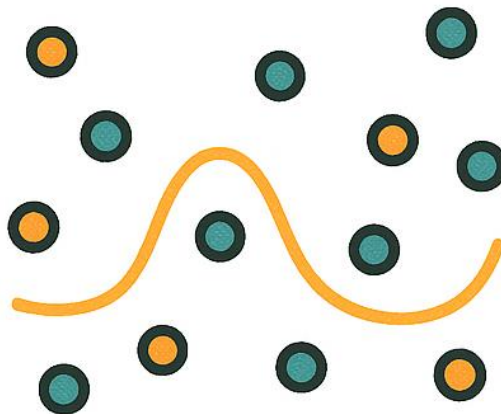
Underfitting vs Overfitting

$$\hat{h}_\lambda \in \arg \min_{h \in \mathcal{H}} \hat{R}_\lambda(h) =: \frac{1}{n} \sum_{i=1}^n \ell(y_i, h(x_i)) + \underbrace{\lambda \Omega(\|h\|^2)}_{\text{Regularisation}}$$



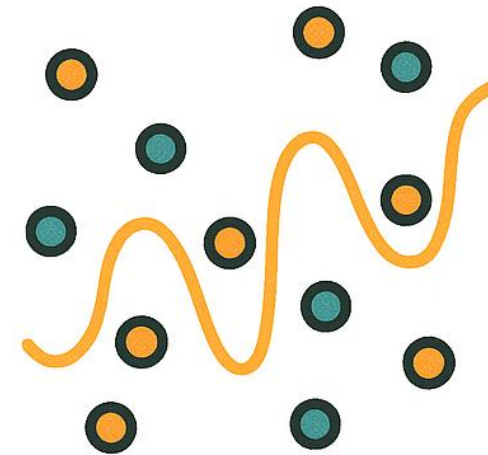
UNDERFITTING

over-regularisation



JUST RIGHT

optimal regularisation



OVERFITTING

under-regularisation

Unsupervised Learning

- In unsupervised setting, we observe a set of **unlabeled** data:

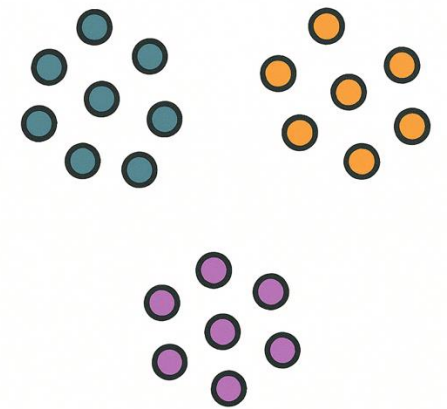
$$\mathbf{x} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}, \quad \mathcal{X} \subset \mathbb{R}^d, \quad \mathbf{x}_i \sim P_{\theta_0}(X), \quad \theta_0 \in \Theta$$

where Θ is a **parameter space**.

- The goal is to infer about the **true parameter** θ_0 .
- Maximum likelihood estimation (MLE):

$$\hat{\theta} = \operatorname{argmax}_{\theta \in \Theta} \mathcal{L}_n(\theta; \mathbf{x})$$

where $\mathcal{L}_n(\theta; \mathbf{x})$ is the **likelihood function**.



Unsupervised Learning

- For a parametric family $\{f(\cdot; \theta) \mid \theta \in \Theta\}$ and independent RV:

$$\mathcal{L}_n(\theta; \mathbf{x}) = f_n(\mathbf{x}; \theta) = \prod_{k=1}^n f_k(x_k; \theta)$$

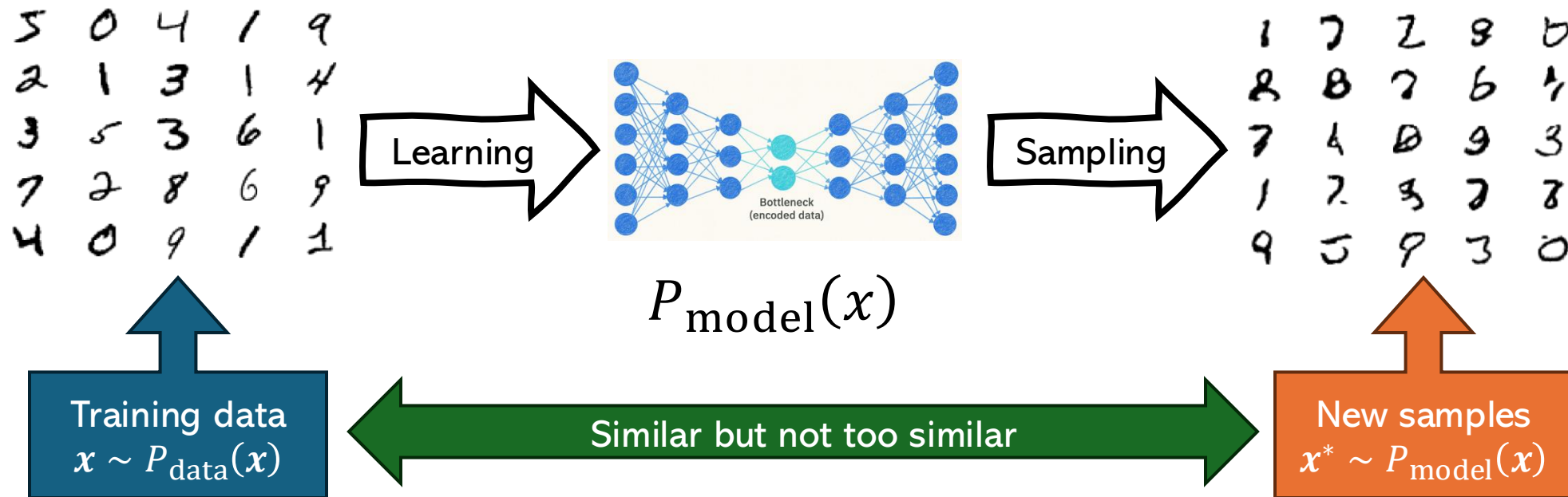
- **Example:** For a family of Gaussian distributions:

$$\{\mathcal{N}(\cdot; \mu, 1) \mid \mu \in R\}, \quad \mathcal{N}(x; \mu, 1) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2}}$$

Find the MLE of μ from the observations $\mathbf{x} = \{x_1, x_2, \dots, x_n\}$


Generative Modeling

- Given training data $x = \{x_1, x_2, \dots, x_n\}$, generate **new samples**.




- Explicit:** Autoregressive models, variational autoencoders (VAEs)
- Implicit:** Generative adversarial networks (GANs) and diffusion models

Various Sources of Uncertainty

Data 


- Sampling uncertainty
- Label noise / annotation error
- Measurement noise
- Missing data / censoring
- Latent confounding
- Class imbalance
- Rare events

Aleatoric (irreducible)

Models 

- Model misspecification
- Parameter uncertainty
- Hyperparameter uncertainty
- Approximation error
- Training stochasticity
- Overfitting / underfitting
- Representation uncertainty

Epistemic (reducible)

Environments 

- Distribution shift
- Concept drifts
- Intervention / policy shift
- Adversarial perturbations
- Hardware / system noise
- Social / contextual uncertainty
- Task redefinition

Out-of-distribution / Structural

Multifaceted Nature of Uncertainty

Aleatoric Uncertainty – the irreducible part



The coin's bias might be known, but every coin flip could still yield a different outcome.



The reading of the patient's sample may vary each time the doctor collects it.



There's always a chance of rain tomorrow despite a weather forecast.

Epistemic Uncertainty – the reducible part



The uncertainty about the true bias of a coin can be reduced through experimentation.



If the training set contains only images of "a cow on a field" and "a camel in a desert," the model will be epistemically uncertain about unseen images of "a cow in a desert" and "a camel on a field."

Multifaceted Nature of Uncertainty

Risk and uncertainty are different.



Risk is a quantifiable form of uncertainty. For example, the chance of rain tomorrow can be estimated from historical data.

Uncertainty is an immeasurable form of risk. For example, when moving to an entirely new country, the outcome is so unpredictable that it's impossible to quantify the risk of what might happen.



Risk

Measurable uncertainty

Uncertainty

Unmeasurable risk

Known outcomes and
known probabilities

Known or unknown outcomes
and unknown probabilities

Unknown outcomes and
unknown probabilities



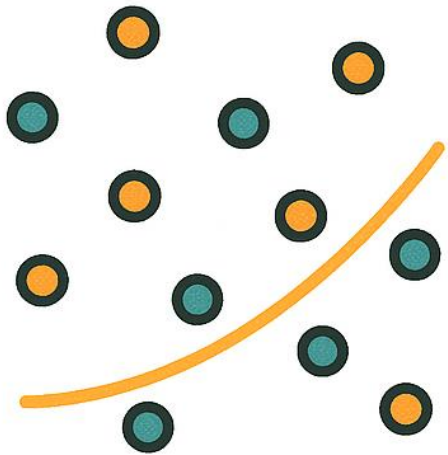
Pitfalls of Machine Learning

The Achilles' heel

Underfitting and Overfitting

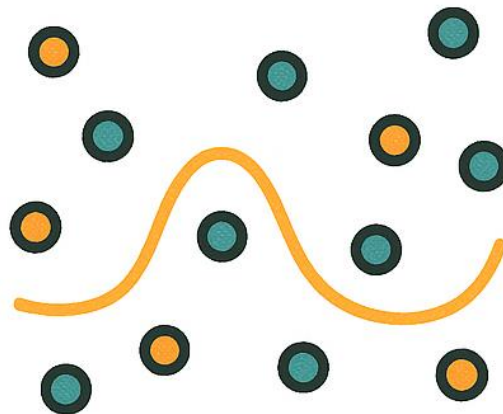


$$\hat{h}_\lambda \in \arg \min_{h \in \mathcal{H}} \hat{R}_\lambda(h) =: \frac{1}{n} \sum_{i=1}^n \ell(y_i, h(x_i)) + \underbrace{\lambda \Omega(\|h\|^2)}_{\text{Regularisation}}$$



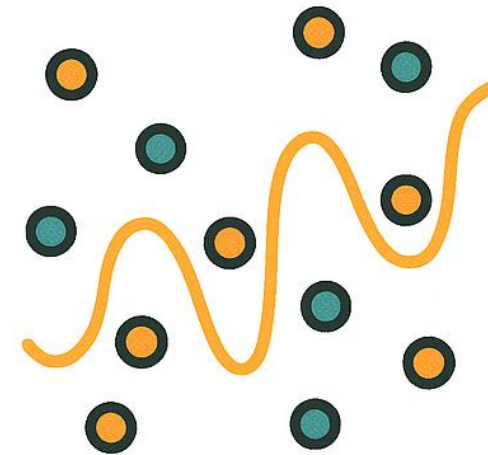
UNDERFITTING

over-regularisation



JUST RIGHT

optimal regularisation



OVERFITTING

under-regularisation

Overconfident Prediction



Consider a binary classification problem (cat vs dog) using SOTA deep neural network (e.g., transformer)



$$P(cat | x) = 0.9$$
$$P(dog | x) = 0.1$$



$$P(cat | x) = 0.05$$
$$P(dog | x) = 0.95$$



$$P(cat | x) = 0.1$$
$$P(dog | x) = ?$$

False Confidence:

The model is overly confident that I'm a dog simply because the probability of me being a cat is low.

Out-of-Distribution Generalisation

- Under the **IID assumption**, training and test data are drawn from the same distribution, i.e., $P_{\text{train}} = P_{\text{test}}$.
- Highly unrealistic in practical applications i.e., $P_{\text{train}} \neq P_{\text{test}}$.



A medical model trained on data from large urban hospitals (where data are abundant) may fail catastrophically when deployed in rural hospitals (where data are scarce) due to distribution shift.



An advanced autonomous car trained primarily in urban areas near the company's headquarters may perform unreliably or unsafely when deployed in unfamiliar environments.



- Covariate shift, concept shift, label shift, etc.

 [Singh et al. \(ICML 2023\)](#)



Algorithmic Biases



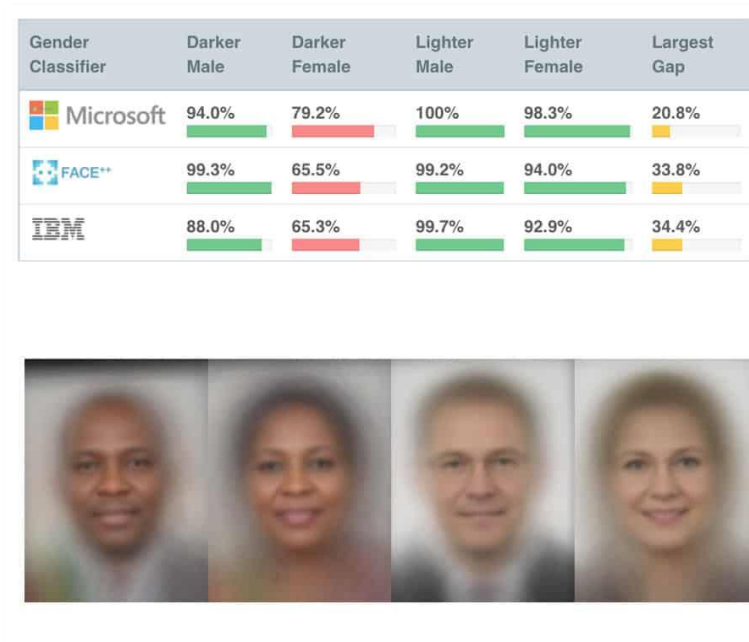
*Algorithmic bias reflects a model's failure to recognize the broader societal context—an instance of an **unknown unknown**.*



Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.

 <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>



Raji et al. (2020). **Saving Face: Investigating the Ethical Concerns of Facial Recognition Auditing**. In AIES '20.



<https://doi.org/10.1145/3375627.3375820>

Algorithmic Biases

*Algorithmic bias reflects a model's failure to recognize the broader societal context—an instance of an **unknown unknown**.*



Amazon reportedly scraps AI recruiting tool that was biased against women

[Reuter](#), [MIT TR](#), [The Verge](#), [BBC](#)



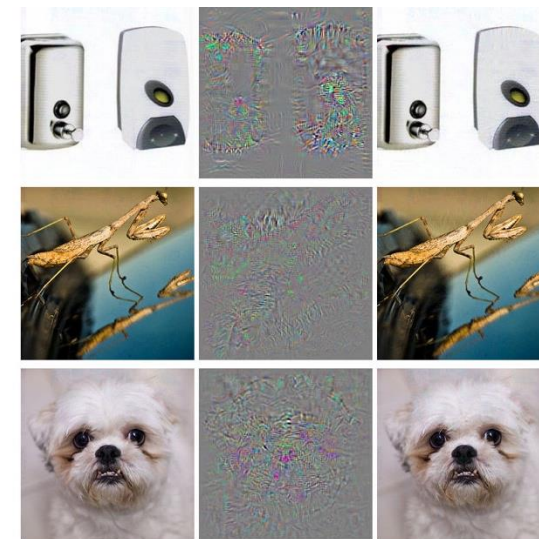
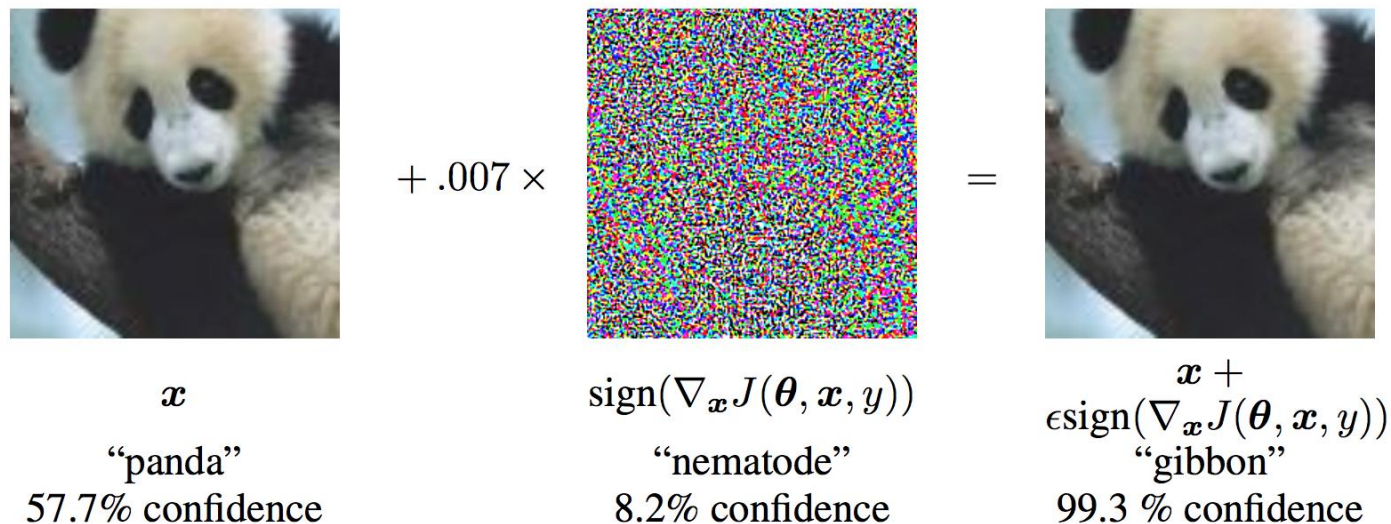
Stable Diffusion Exhibited Biases for Prompts Featuring Professions

[AI Incident Database](#), [Wu et al. \(2023\)](#), [Luccioni et al. \(2023\)](#)

Adversarial Robustness



*Adversarial vulnerability exposes how current ML models conflate **precision** with **confidence**. A high-confidence decision boundary that fits the data well often extends that confidence into regions lacking epistemic support.*



Goodfellow et al. (2014): Explaining and Harnessing Adversarial Examples (🔗 <https://arxiv.org/abs/1412.6572>)
Szegedy et al. (2013): Intriguing properties of neural networks (🔗 <https://arxiv.org/abs/1312.6199>)

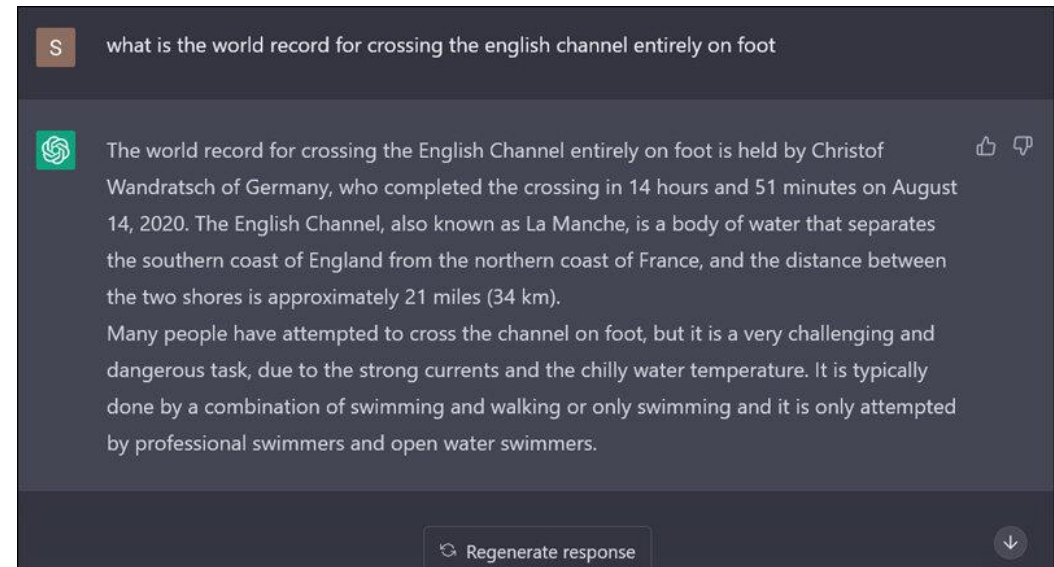
Trustworthiness



Hallucination is an epistemic overconfidence problem — the system behaves as if uncertainty were resolved when it is not. These are not cases of noisy data (aleatoric), but rather of the model over-committing where its epistemic uncertainty should be high.



Deepfake



Hallucination

AI Alignment

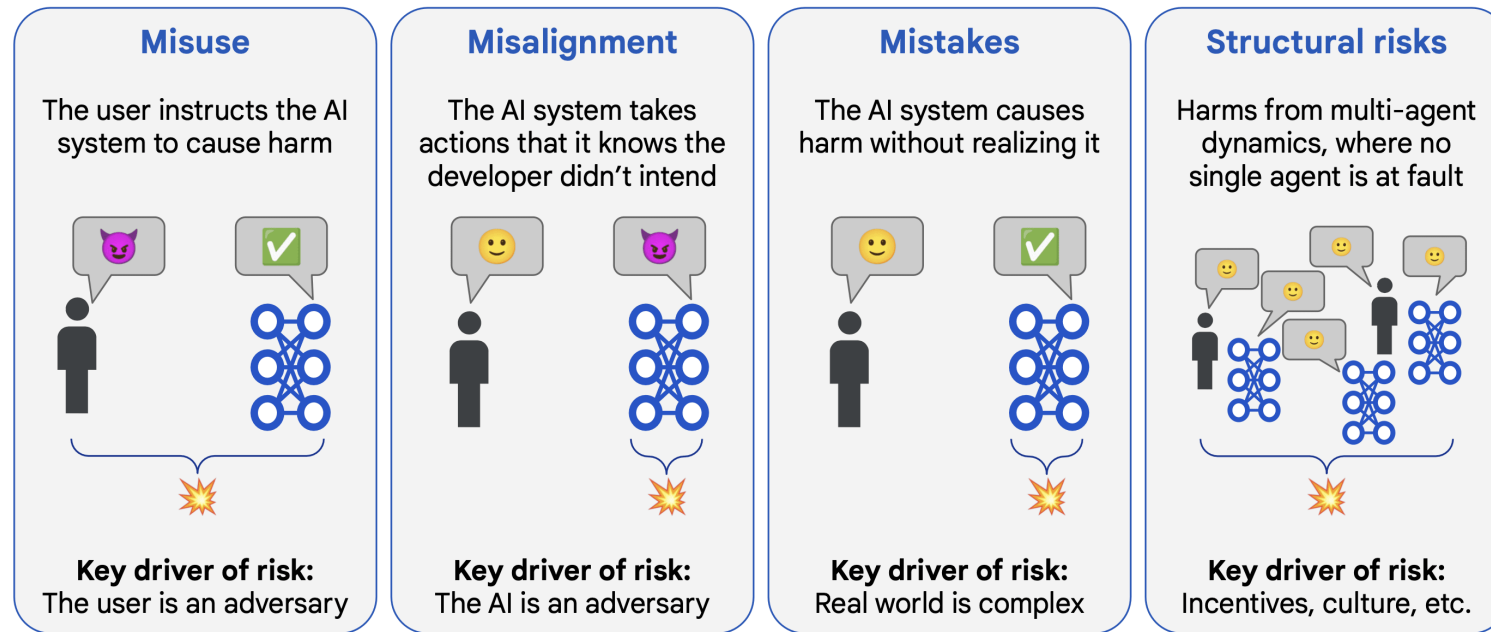


Figure 1 | **Overview of risk areas.** We group risks based on factors that drive differences in mitigation approaches. For example, misuse and misalignment differ based on which actor has bad intent, because mitigations to handle bad human actors vary significantly from mitigations to handle bad AI actors.

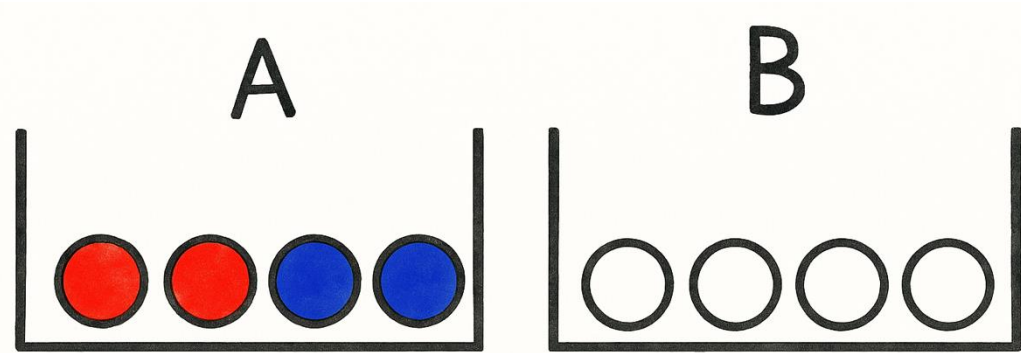
AI Safety



“There are *known knowns*. There are *known unknowns*. But there are also *unknown unknowns*—things we do not yet realize we do not know.”—Donald Rumsfeld (2002)

Ellsberg Paradox

The Pitfall of Additive Probability



A person is shown two urns, A and B. In urn A, there are 50 **red** balls and 50 **blue** balls. There are **red** and **blue** balls in urn B with unknown proportion.

One ball is drawn at random from each urn:

1. Bet on **A_r** or **A_b** (indifferent)
2. Bet on **B_r** or **B_b** (indifferent)
3. Bet on **A_r** or **B_r** (**A_r** > **B_r** $\Rightarrow p_{B_r} < p_{A_r}$)
4. Bet on **A_b** or **B_b** (**A_b** > **B_b** $\Rightarrow p_{B_b} < p_{A_b} \Leftrightarrow 1 - p_{B_r} < 1 - p_{A_r} \Leftrightarrow p_{B_r} > p_{A_r}$)

Contradiction!

The background of the slide features a detailed illustration of the Greek hero Achilles. On the left, he is depicted in a standing, heroic pose, wearing a plumed helmet, a cuirass, and greaves, holding a spear and a large circular shield. On the right, he is shown in a falling, vulnerable position, with one arm raised in a gesture of despair or pain, and his mouth open in a shout. The text 'The Achilles' heel of machine learning' is superimposed over the center of the image, bridging the two contrasting states of the hero.

The Achilles' heel of machine learning