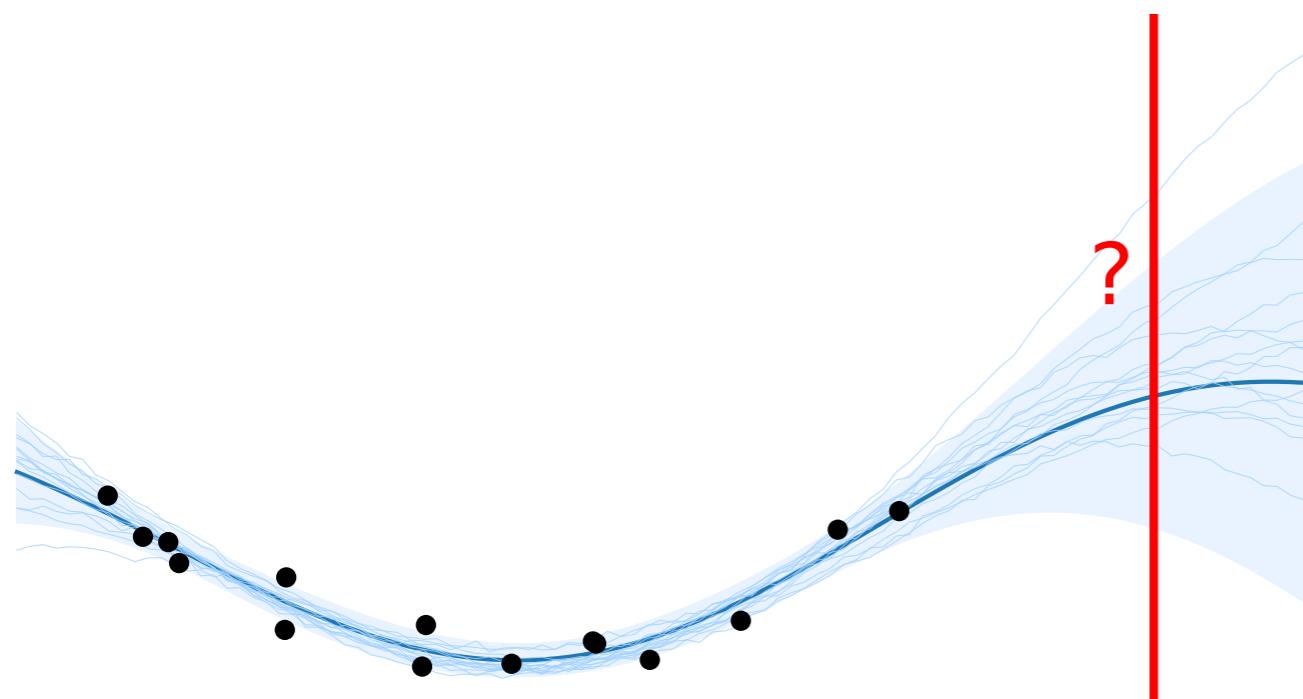


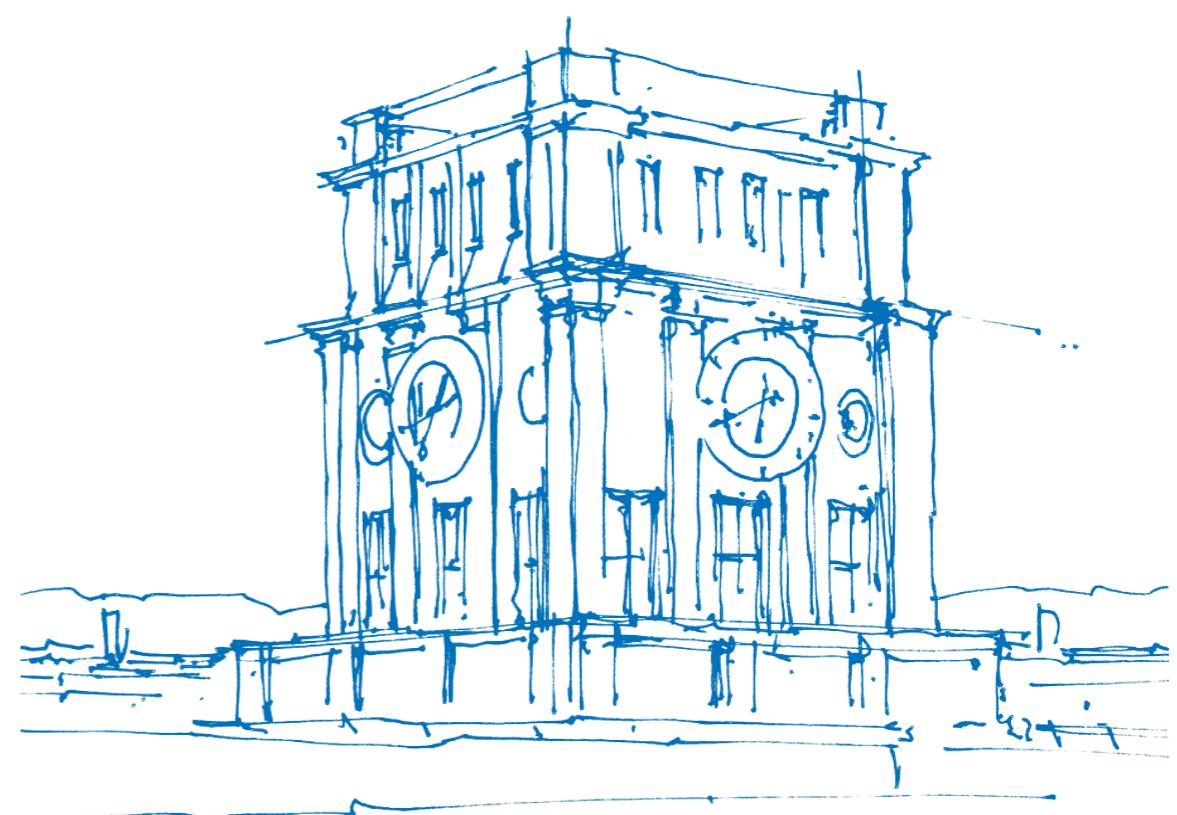
Probabilistic ML - Part II



Ai in Medicine Lecture Series

Chair for Ai in Medicine
Faculty of Informatics & Medicine
Technical University Munich

Moritz Knolle
15.12.2022

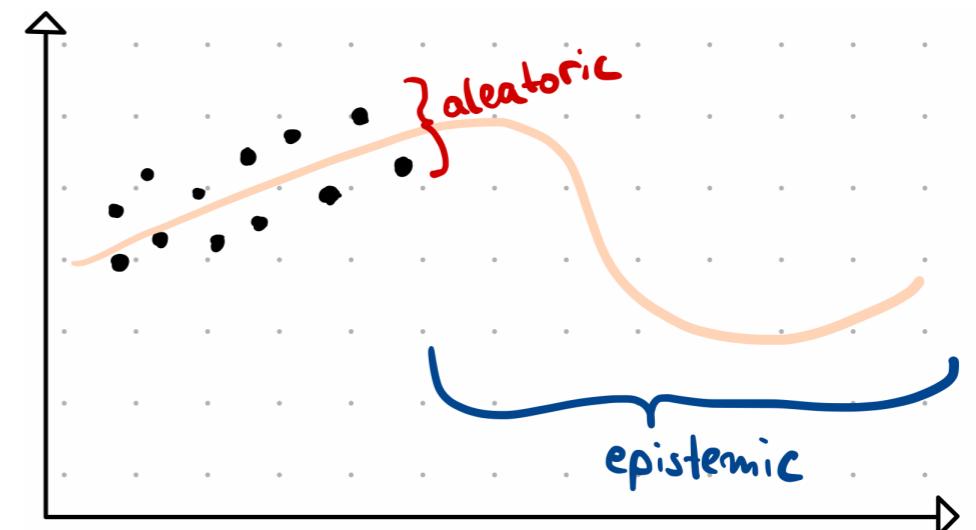


Last time:

- Pdf, cdf, random variables
- Joint & conditional probability
- Probabilistic process > Bayes Theorem

$$p(H|D) = \frac{p(D|H) p(H)}{p(D)}$$

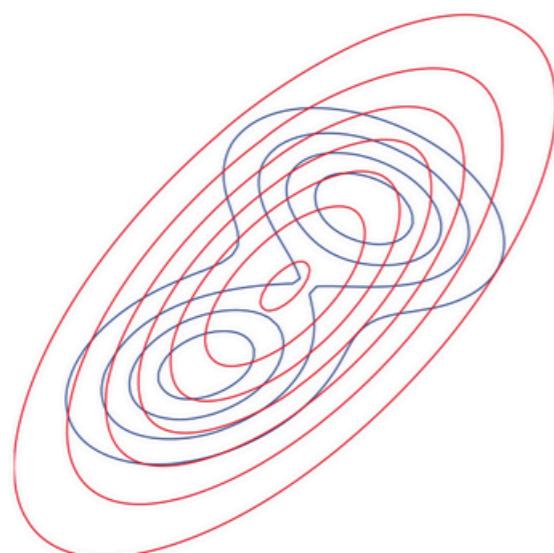
- Epistemic vs. Aleatoric uncertainty
- Sampling-based inference



Outline

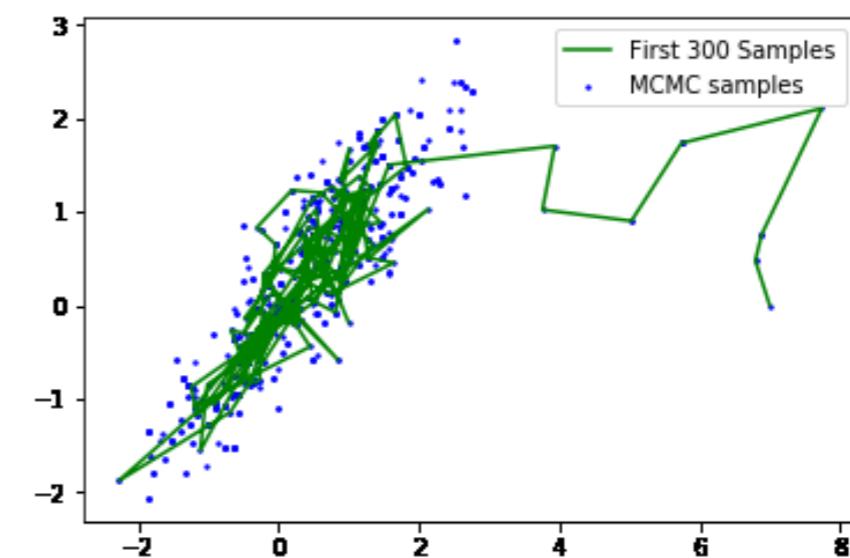
- I. Overview approximate inference
- II. From rejection sampling to Markov Chain Monte Carlo
- III. Variational methods
- IV. Measuring the quality of probabilistic predictions
- V. Practical methods for uncertainty quantification in deep learning

Today: Approximate Inference



Variational Methods

Approximation as Optimisation



Monte Carlo Methods

Sampling based

Recap: Sampling-based approximate inference

$$p(\theta|D) = \frac{p(D|\theta)p(\theta)}{\underbrace{p(D)}_{\text{expensive to evaluate}}}$$

cheap to evaluate

$$\int p(D|\theta)p(\theta)d\theta$$

Monte Carlo sampling:

$$\int p(D|\theta)p(\theta)d\theta = \frac{1}{N} \sum_i^N p(D|\theta_i)$$

$\theta_i \sim p(\theta)$

Rejection sampling in disguise

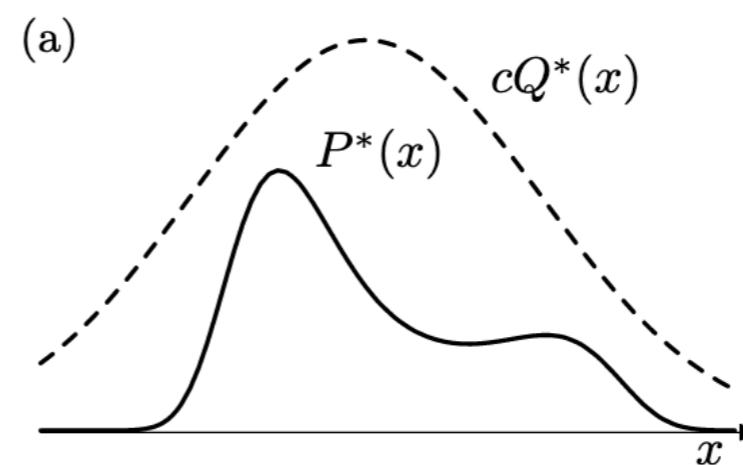
$$p(\theta|D) = \frac{p(D|\theta) p(\theta)}{p(D)}$$

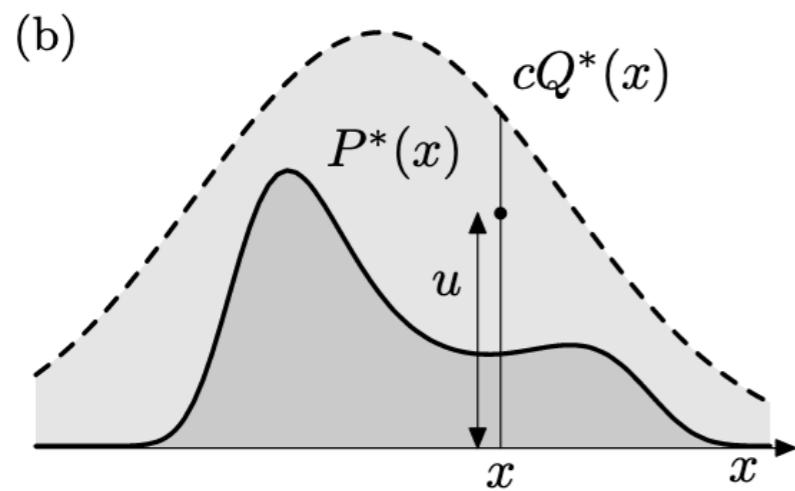
cheap to evaluate
expensive to evaluate

$$\int p(D|\theta) p(\theta) d\theta$$

- want to evaluate $P(x)$
- **BUT** can only evaluate an unnormalised version

$$P^*(x) = \frac{P^*(x)}{Z}$$



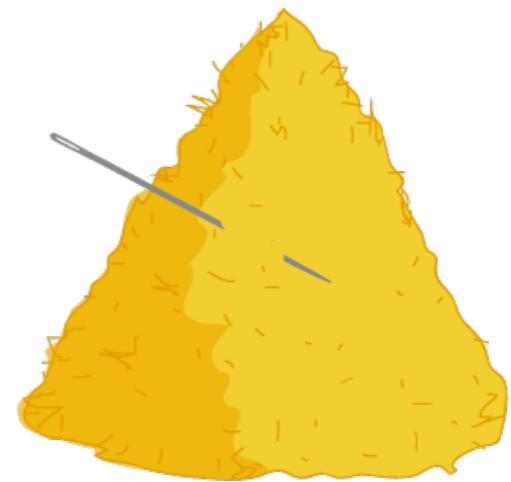
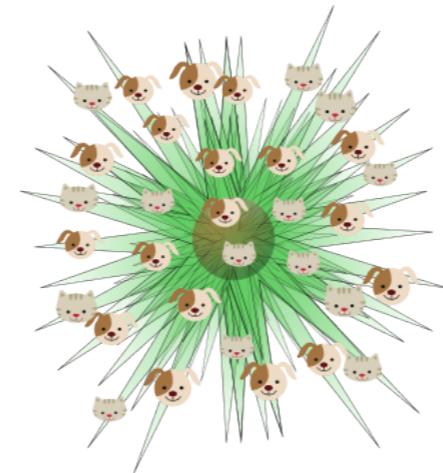
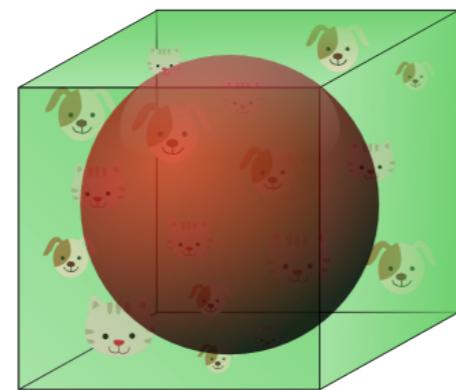
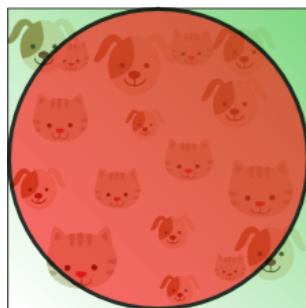


Rejection Sampling:

1. Construct proposal density $Q(x)$ (un-normalised version $Q^*(x)$)
2. Generate a sample x from $Q(x)$
3. Sample uniformly from
 $u \sim [0, cQ^*(x)]$
4. If $u > P^*(x)$ discard x otherwise add x to list of samples $[\dots, x^{(r)}]$

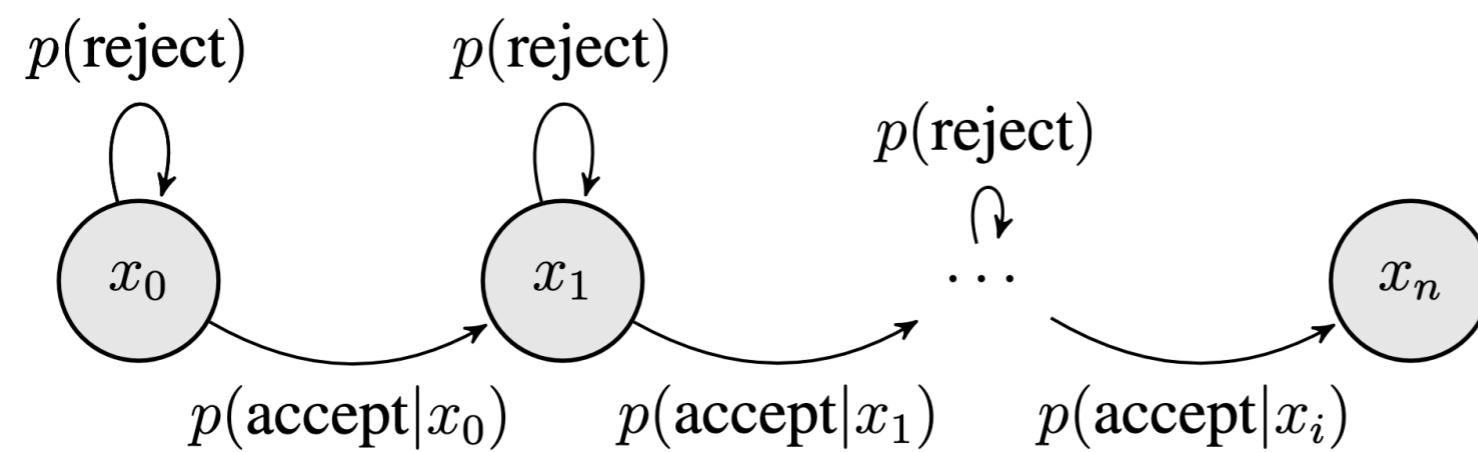
Problems: Rejection sampling

- $Q(x)$ needs to be similar to $P(x)$
- Rejection rate r is a function of dimensionality d :
 - As $d \rightarrow \infty$; $r \rightarrow 0$
- Samples are uncorrelated
 - repeatedly finding needle in haystack



Can we improve rejection sampling?

- What if we could create samples that are correlated?



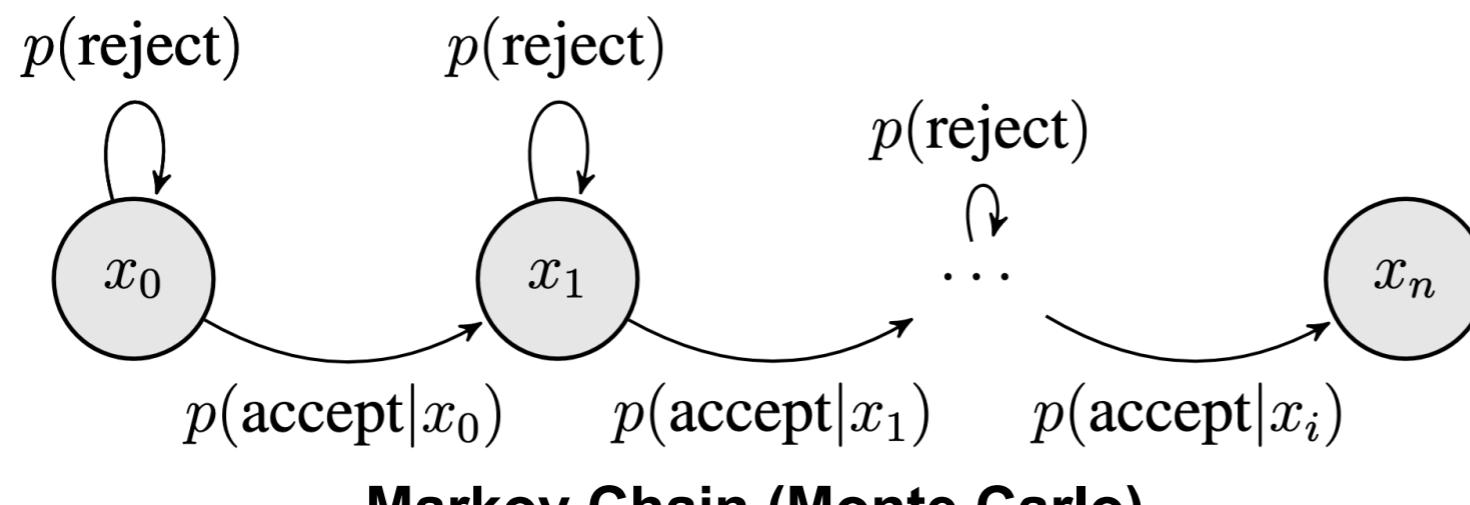
Can we improve rejection sampling?

- For every x_i , we need to model the transition probability $p(x_{i+1} | x_i)$
- next state only depends on current state!

Transition operator:

$$\bullet \quad a = \frac{P^*(x')/Z}{P^*(x^{(t)})/Z} = \frac{P^*(x')}{P^*(x^{(t)})};$$

- If $a > 1$ likely to accept x' , otherwise likely to reject x'

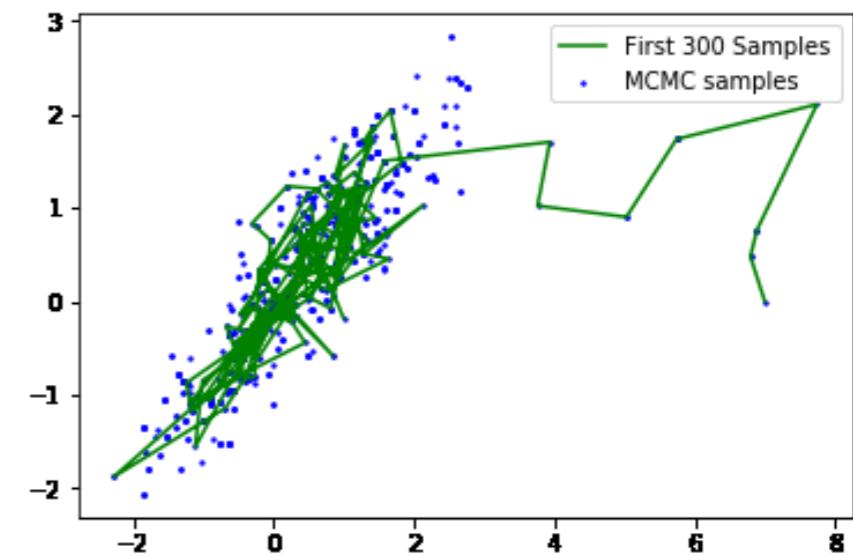


Metropolis-Hastings Algorithm

- A Markov Chain Monte Carlo (MCMC) method

Metropolis Hastings:

1. Sample $x' \sim Q(x' | x^{(t)})$
2. $a = \frac{P^*(x')}{P^*(x^{(t)})}, u \sim U(0,1)$
3. If $u \leq a$: **accept**
then $x_{t+1} = x'$
4. Else: **reject**
 $x_{t+1} = x_t$



Cool Demo

Can we come up with a similarity measure to see how close we are to the exact solution?

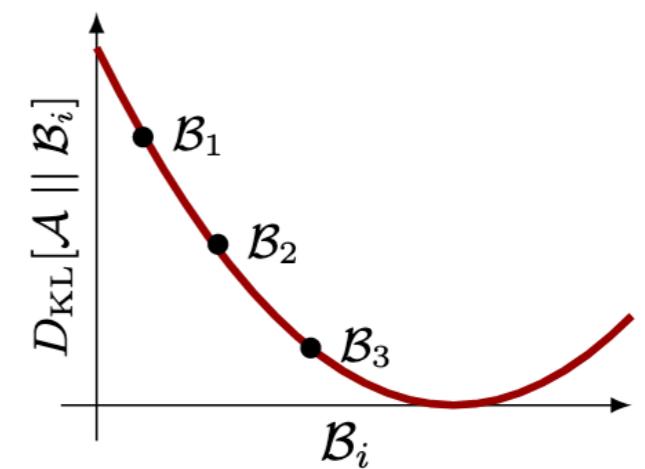
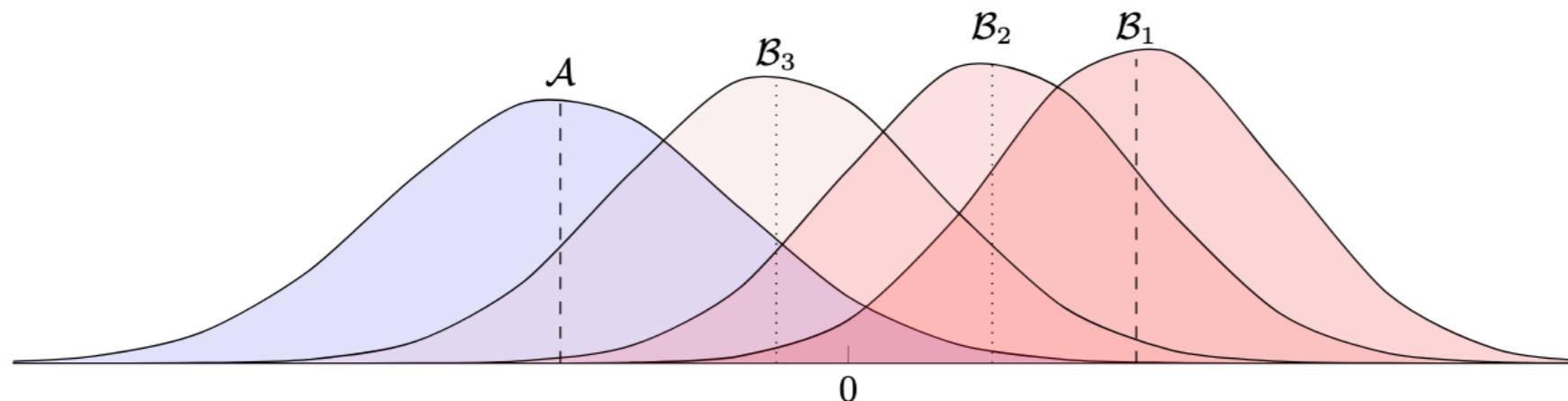
(To build an optimisation based inference approach)

The KL divergence - a measure of distributional similarity

- similarity between two probability distributions A and B

$$D_{\text{KL}}[A \parallel B] = \mathbb{E}_{x \sim A}[\log \frac{A(x)}{B(x)}]$$

- $D_{\text{KL}}[A \parallel B] = 0$ only if A and B are identical (under x)
- Nice error-coding interpretation



Variational Inference - Approximation through Optimisation:

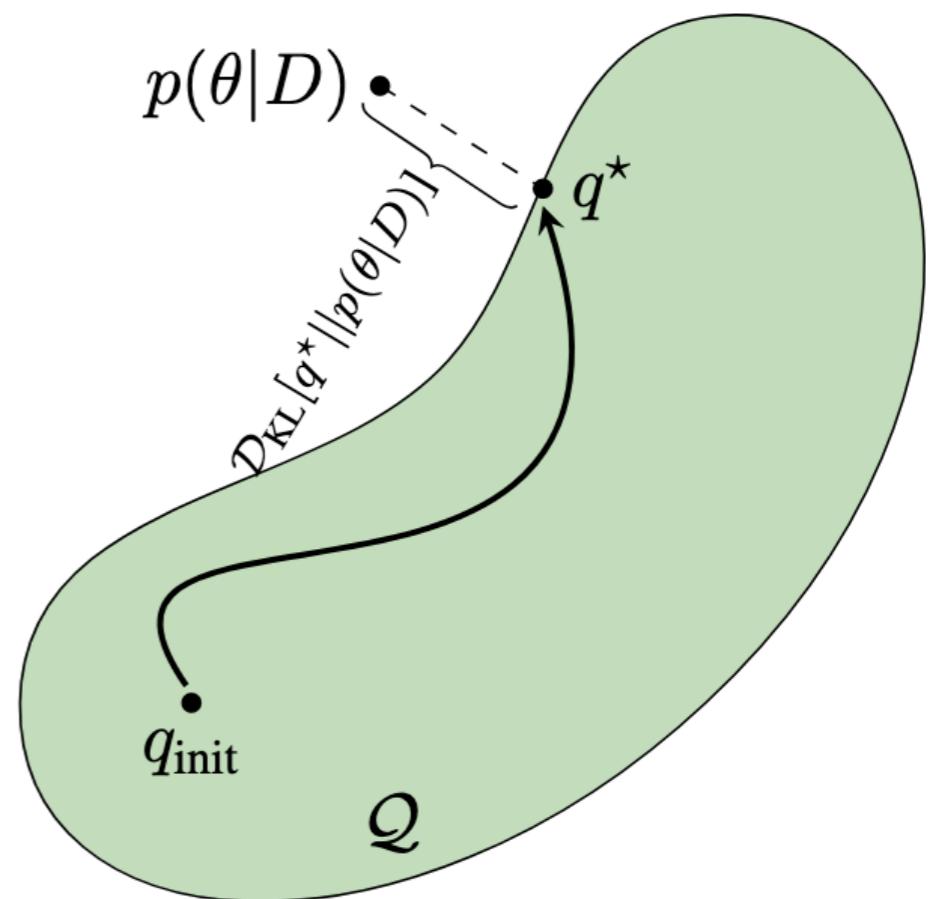
1. Define variational distribution family \mathcal{Q} , parameterised by ν
2. Measure quality of approximation through:

$$M(q_\nu^{(i)}) = D_{\text{KL}}[q_\nu^{(i)} \parallel p(\theta | D)]$$

3. Update ν

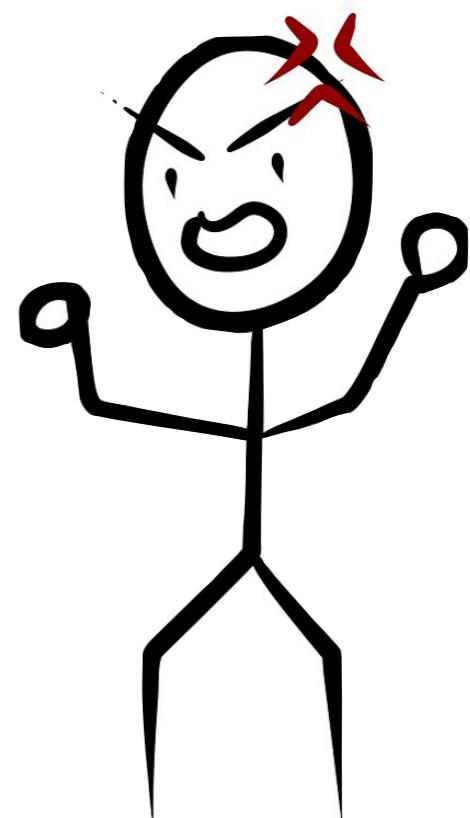
$$\nu^{(i+1)} = \nu^{(i)} - \nabla_\nu M(q_\nu^{(i)})$$

4. Repeat 2. & 3. until convergence



Can you spot the problem with this approach?

$D_{\text{KL}}[q_v^{(i)} \parallel p(\theta | D)]$ involves the expensive posterior
 $p(\theta | D)$ again



The ELBO Trick

- Maximise lower bound on the objective instead (ELBO, \mathcal{L})

$$\mathcal{L}(q_v) = \mathbb{E}_q[\log p(D | \theta)] - \underbrace{D_{KL}[q(\theta) || p(\theta)]}_{\text{KL-divergence between } q \text{ and prior on } \theta}$$

$$\arg \max_v \mathcal{L}(q_v)) = \arg \min_v D_{KL}[q_v || p(\theta | D)]$$

- Maximising \mathcal{L} , monotonically decreases our objective!

Quiz: MCMC vs. Variational Inference

| | Variational Inference | MCMC Methods |
|---------------------------------|-----------------------|--------------|
| • Approximation Quality Measure | ✓ | ✗ |
| • Recover exact posterior | (✓) | ✓ |
| • Fast? | ✓ | ✗ |

We covered how to approximate posteriors

But how to measure the quality
of a predictive distribution?



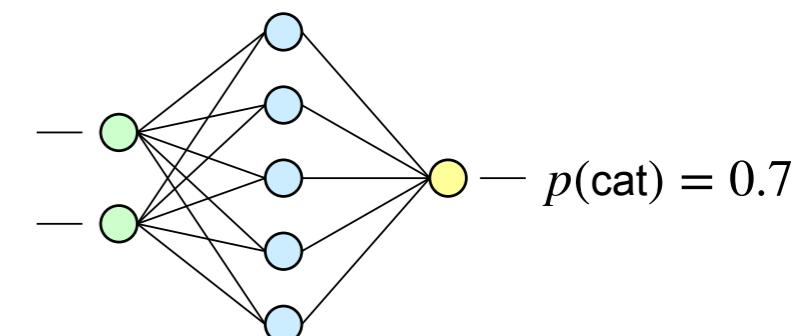
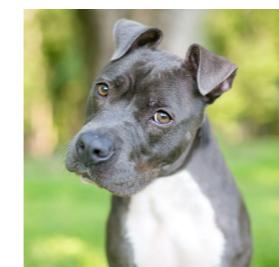
Assessing the quality of a probabilistic prediction?

Calibration

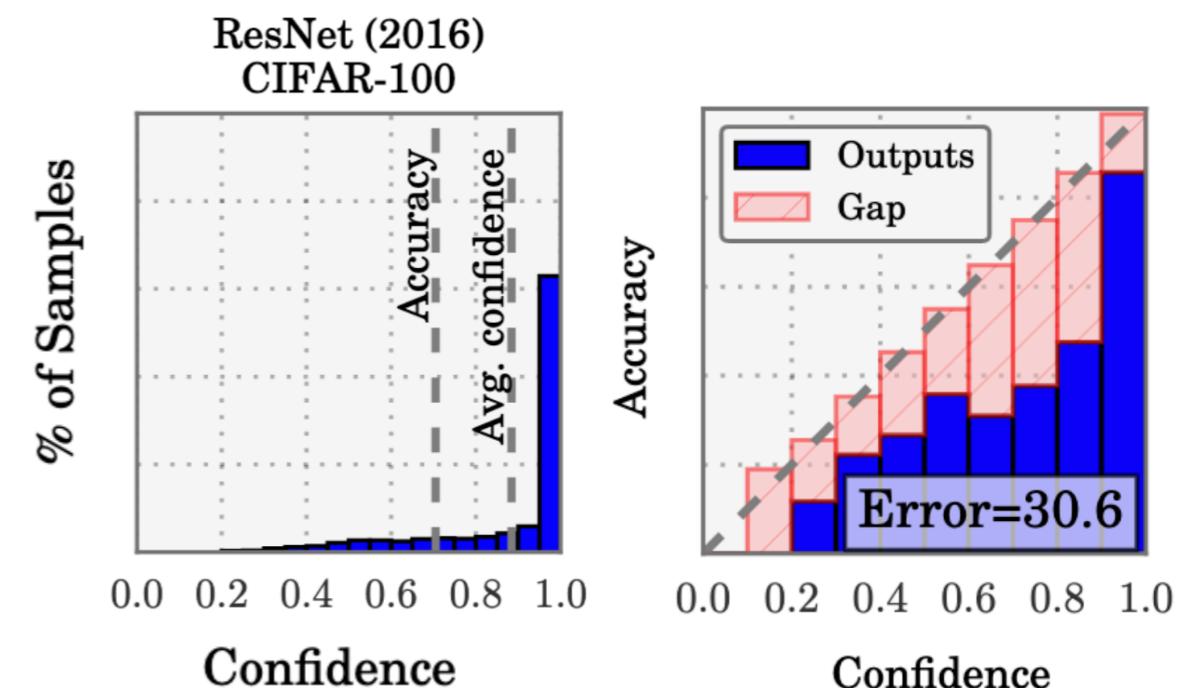
$$P(y|\hat{p}) = p \quad \forall p \in [0,1]$$

Intuitively (model confidence \simeq probability)

Prediction with confidence 20% should be correct 20% of the time.



- Softmax confidence for modern neural networks tend do be ***mis-calibrated***
- Expected calibration error (ECE)
- ***Calibrated* \neq *accurate*** prediction



Metrics for probabilistic predictions (scoring rules)

- Brier Score (mean squared error)

$$\frac{1}{N} \sum_{i=1}^N (p_i - y_i)^2$$

- Log likelihood:

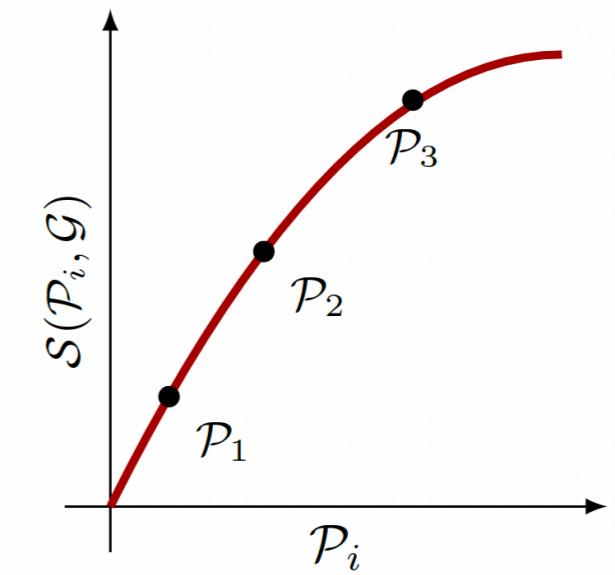
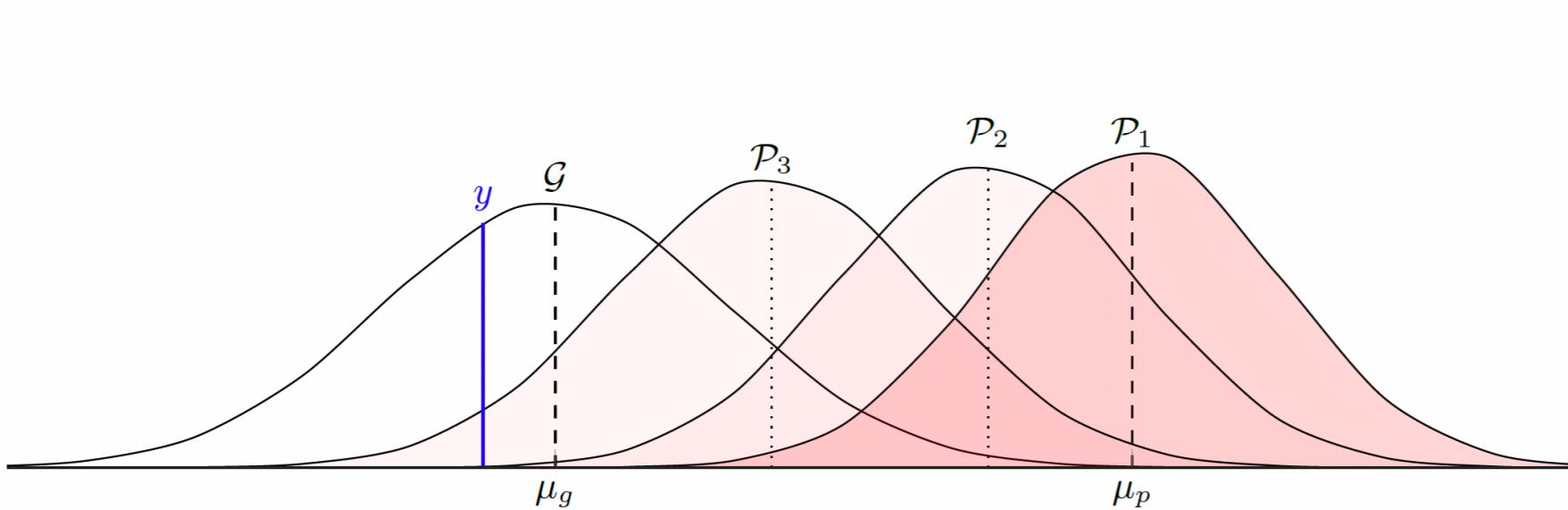
$$\frac{1}{N} \sum_{i=1}^N \log p_\theta(y_i)$$

- Also called **scoring rules**

Proper scoring rule

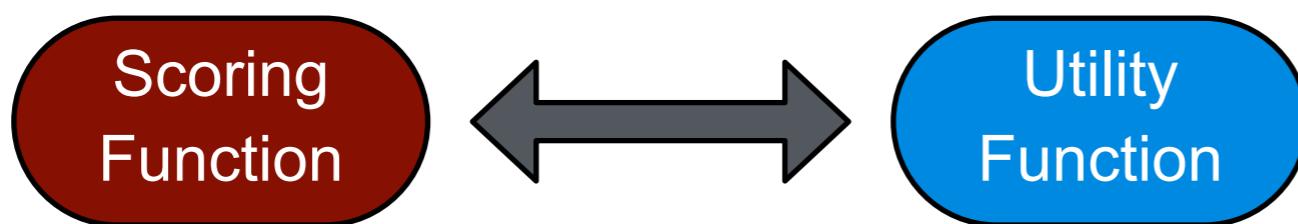
A **scoring rule** $\mathcal{S}(P, G)$ for a predictive distribution $P = (p_1, \dots, p_k)$ and underlying ground truth distribution G is called **proper** iff:

$$\mathcal{S}(P, P) \leq \mathcal{S}(P, G) \quad \forall P, G$$

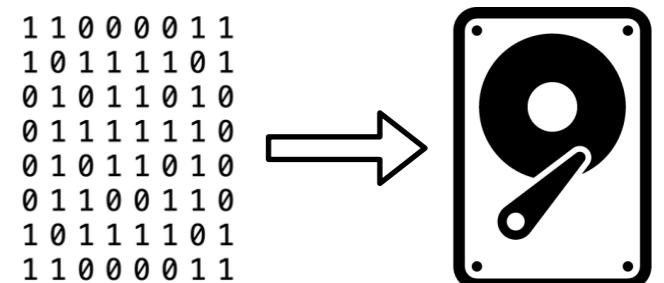


Scoring rules and decision making

- Close connection to *decision and information theory*



- Dual form (see Gneiting & Raftery):
 - **Primal:** scoring function maximisation
 - **Dual:** utility maximisation
- **Log likelihood**, direct interpretation as growth rate when betting on mutually exclusive outcomes (Cover & Thomas)

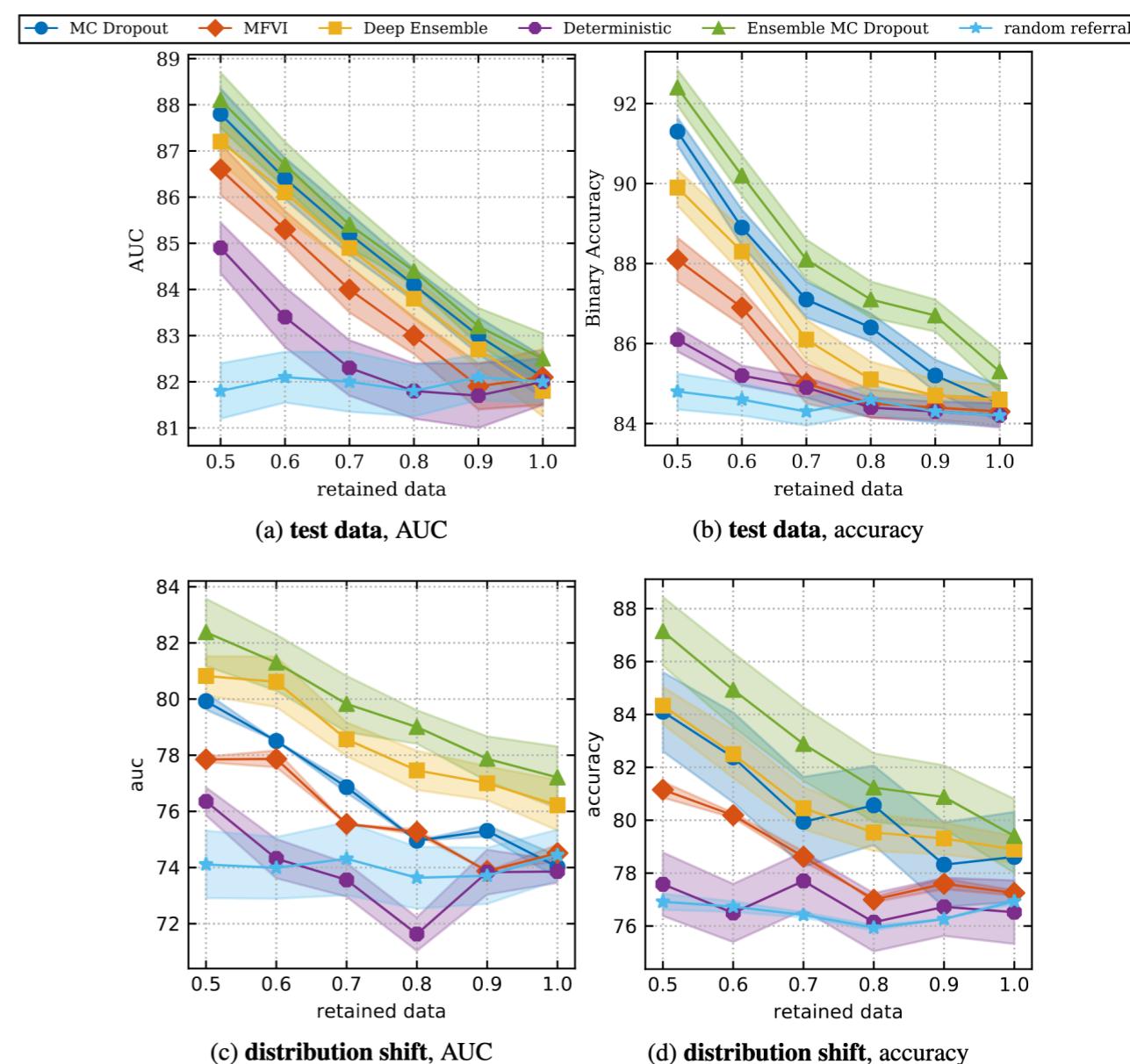
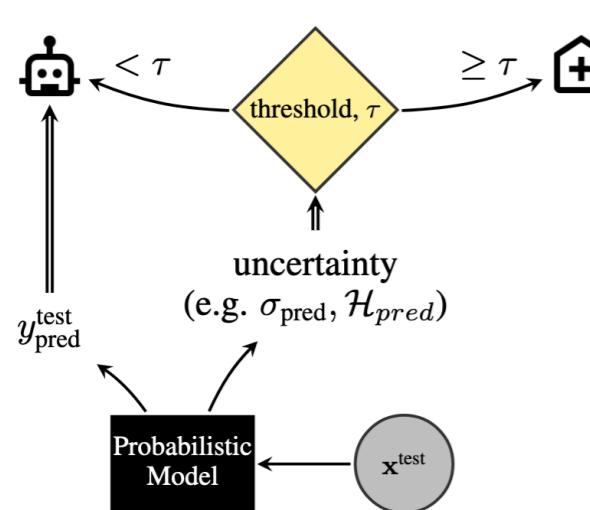


$$H(X) = - \sum_{x \in X} p(x) \log p(x)$$

Now onto some practical uncertainty methods for
Ai in Medicine

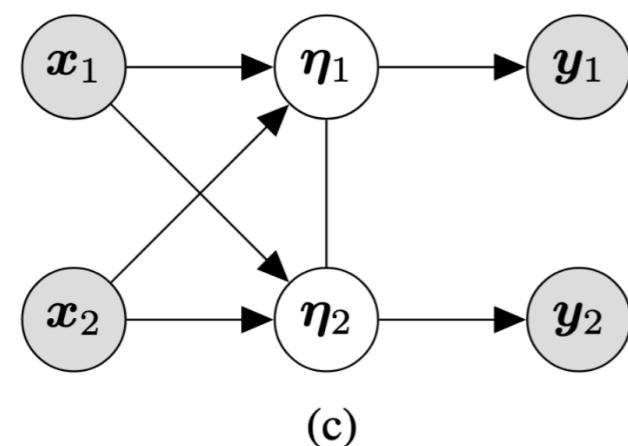
Do VI and MCMC scale to the tasks we care about?

- Sadly not! Computationally prohibitive with large NNs
 - Mean field VI gives sub-optimal posterior approximation
 - In practice: simply use ensembles of deep networks

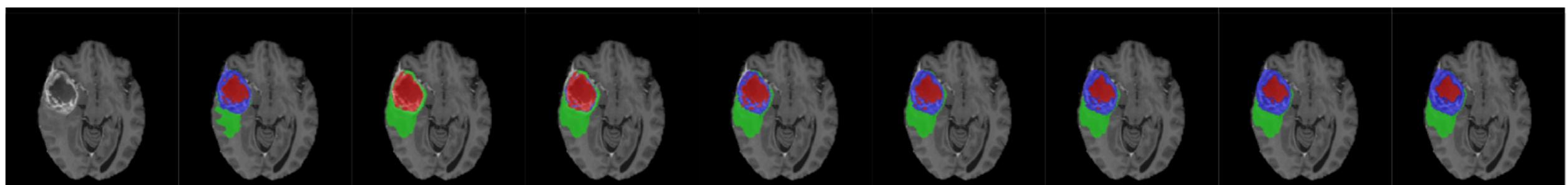


Stochastic Segmentation Networks (SSNs)

- Segmentation targets often ambiguous and can vary between ground truth raters
- Can we model aleatoric inter-rater segmentation uncertainty?
- Network output parameterises a low-rank Gaussian over image space



$$\boldsymbol{\eta}^{(m)} | \mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}(\mathbf{x}), \boldsymbol{\Sigma}(\mathbf{x})) ,$$



THE END

Questions?