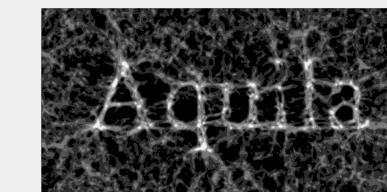


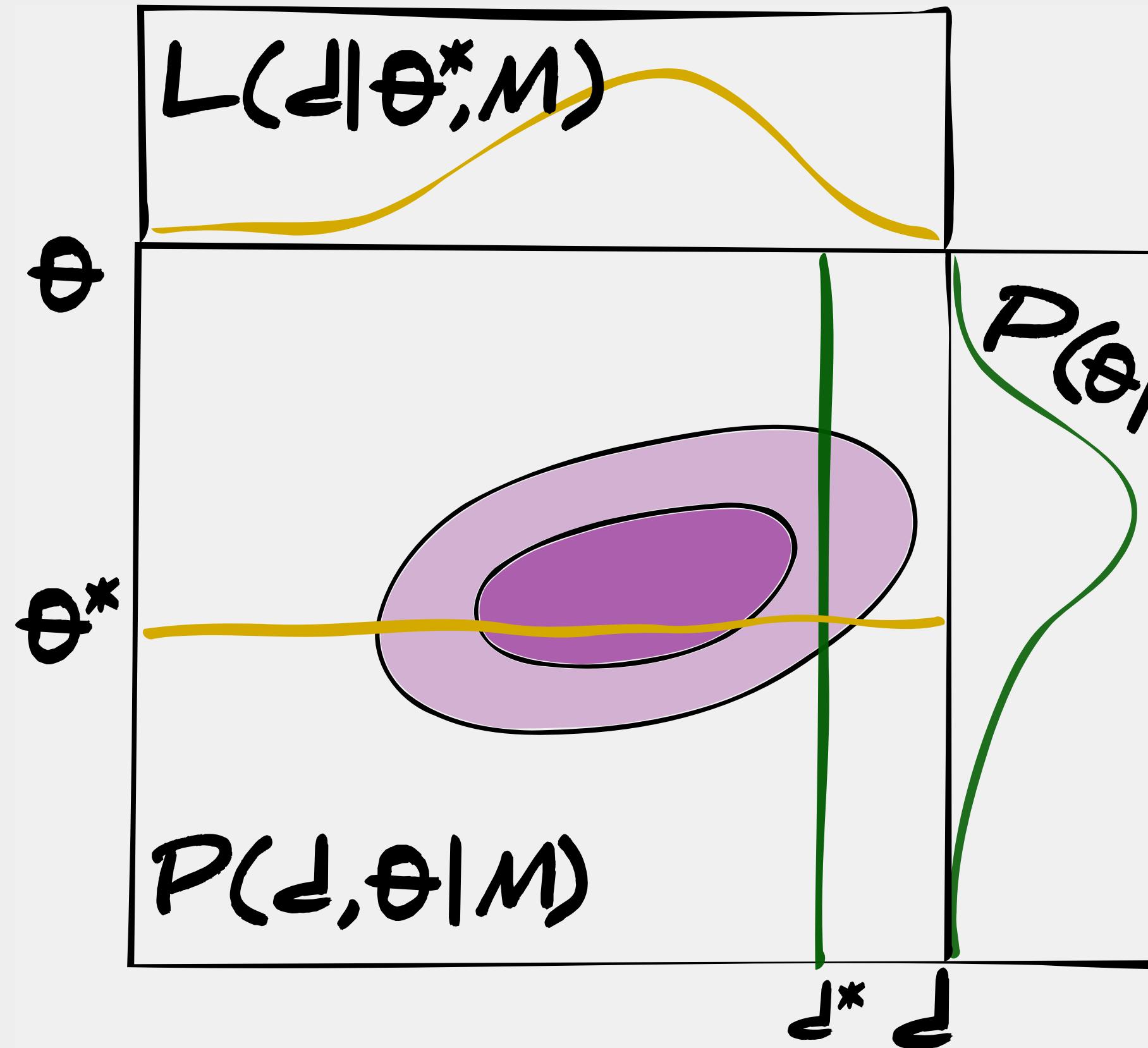
Parameter inference using neural networks

Tom Charnock

Institut d'Astrophysique de Paris



Joint distribution of parameters, θ , and possible data, d ,
 which can be generated by a model, $\mathcal{M} : \theta \rightarrow d$



$$P(\theta | d, M) = \frac{\mathcal{L}(d | \theta, M) p(\theta | M)}{p(d | M)}$$

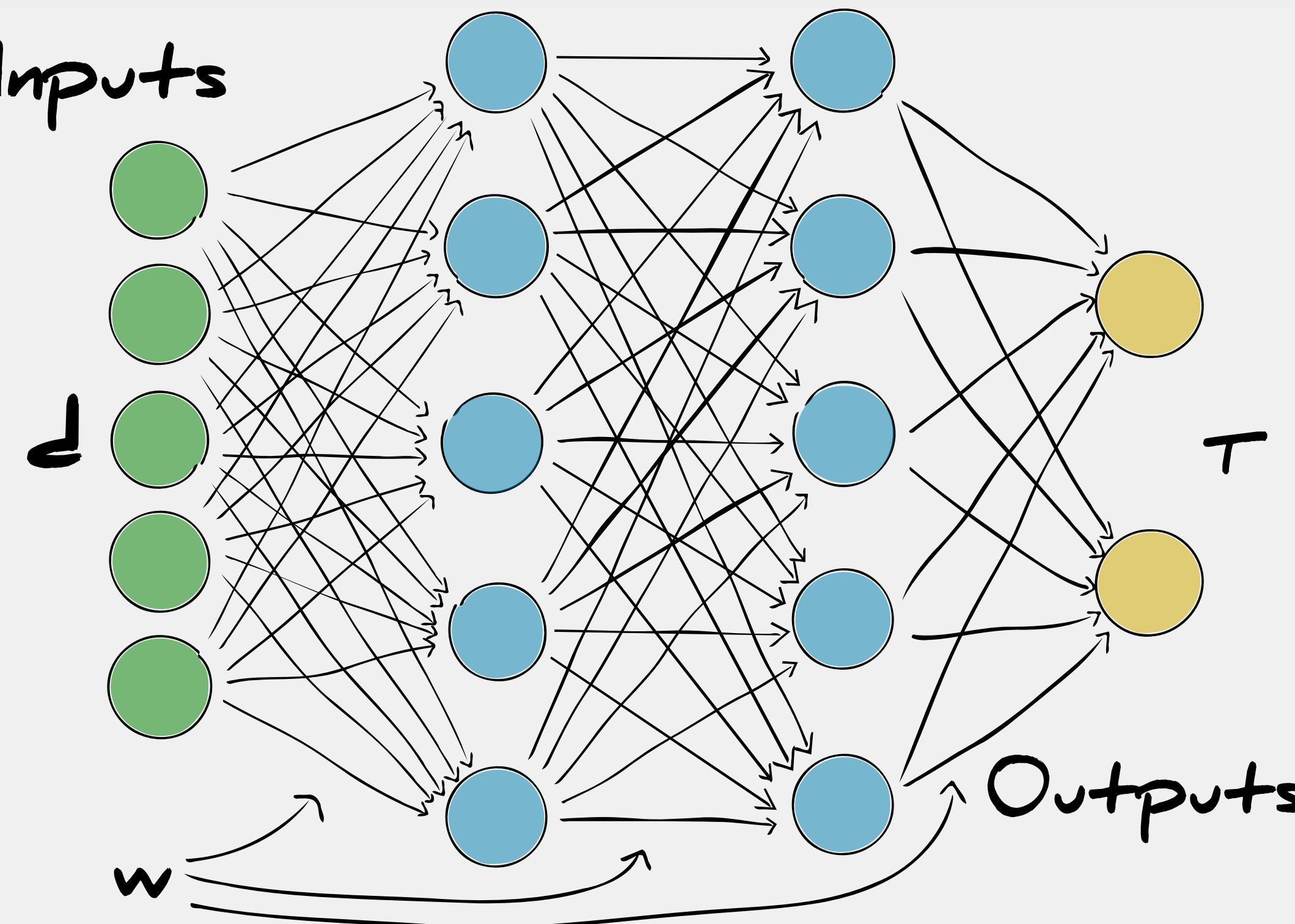
Obligatory Bayes
Theorem

Likelihood
 $\mathcal{L}(d | \theta^*, M)$

Posterior
 $P(\theta | d^*, M)$

Parameter inference using neural networks

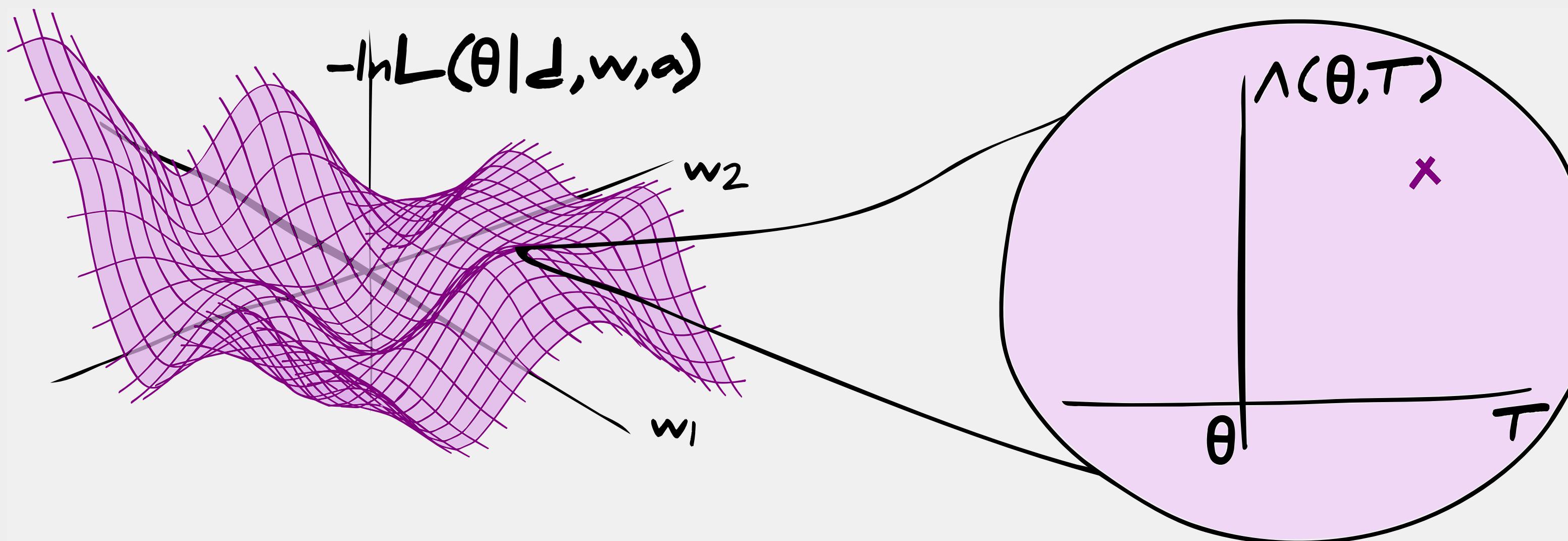
It's common to try and predict model parameters by regression



$$\text{NN}(\omega, \alpha) : \mathbf{d} \rightarrow \tau$$

An approximation to a model, $\mathcal{M} : \mathbf{d} \rightarrow \theta$

How likely are we to obtain any particular output from the network



$$\Lambda(\theta, \tau) \propto -\ln \mathcal{L}(\theta | \mathbf{d}, \omega, \alpha)$$

The posterior becomes the posterior predictive density

$$\mathcal{P}(\theta|\mathbf{d}) = \int d\omega d\alpha \mathcal{L}(\theta|\mathbf{d}, \omega, \alpha) \mathcal{P}(\omega, \alpha)$$

The posterior becomes the posterior predictive density

$$\mathcal{P}(\theta|\mathbf{d}) = \int d\omega d\alpha \mathcal{L}(\theta|\mathbf{d}, \omega, \alpha) \mathcal{P}(\omega, \alpha)$$

$\mathcal{P}(\theta|\mathbf{d})$ - Posterior predictive density

How likely are the true parameters given some data?

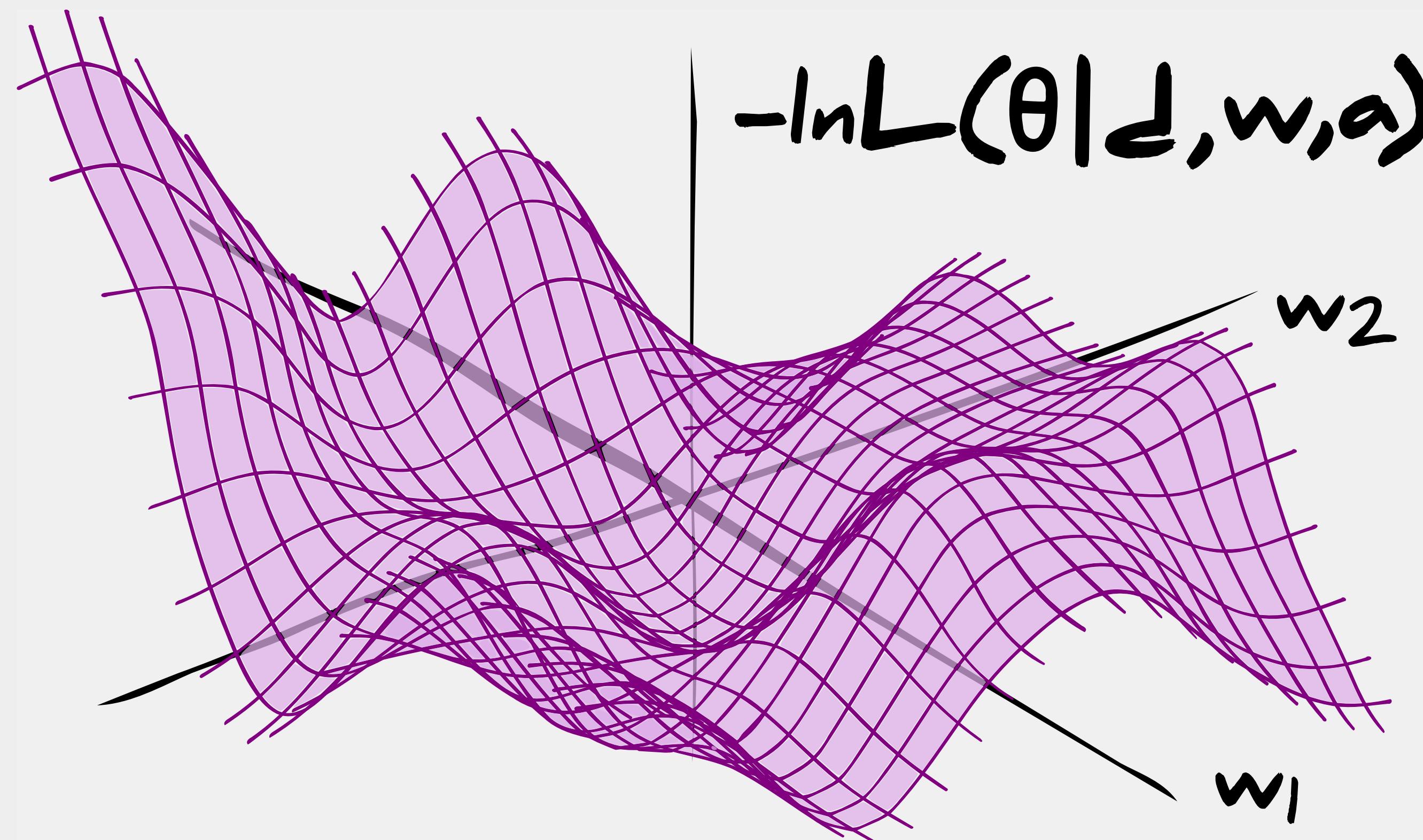
$\mathcal{L}(\theta|\mathbf{d}, \omega, \alpha)$ - Likelihood

How likely are the parameters to be generated by a particular network?

$\mathcal{P}(\omega, \alpha)$ - Probability density

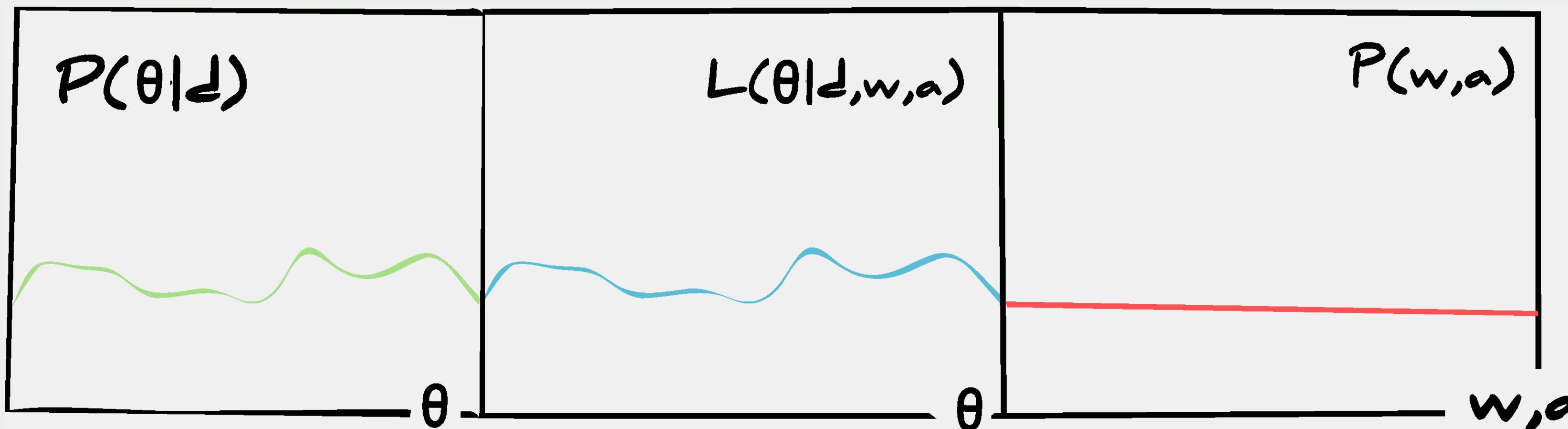
What is the probability of obtaining a particular network with particular parameter values?

The chance of getting any parameter value from any network given some data is *almost* equal



So all the information comes from prior knowledge of the weights and hyperparameters

$$\mathcal{P}(\theta|d) = \int d\omega d\alpha \mathcal{L}(\theta|d, \omega, \alpha) \mathcal{P}(\omega, \alpha)$$



Where does the information about the weights and hyperparameters come from?

Training and validation data

Training and validation data

$$\{\mathbf{d}, \theta\}^{\text{train}} \equiv \{\mathbf{d}_i^{\text{train}}, \theta_i^{\text{train}} | i \in [1, n_{\text{train}}]\}$$

Training and validation data

$$\{\mathbf{d}, \theta\}^{\text{train}} \equiv \{\mathbf{d}_i^{\text{train}}, \theta_i^{\text{train}} | i \in [1, n_{\text{train}}]\}$$

$$\{\mathbf{d}, \theta\}^{\text{val}} \equiv \{\mathbf{d}_i^{\text{val}}, \theta_i^{\text{val}} | i \in [1, n_{\text{val}}]\}$$

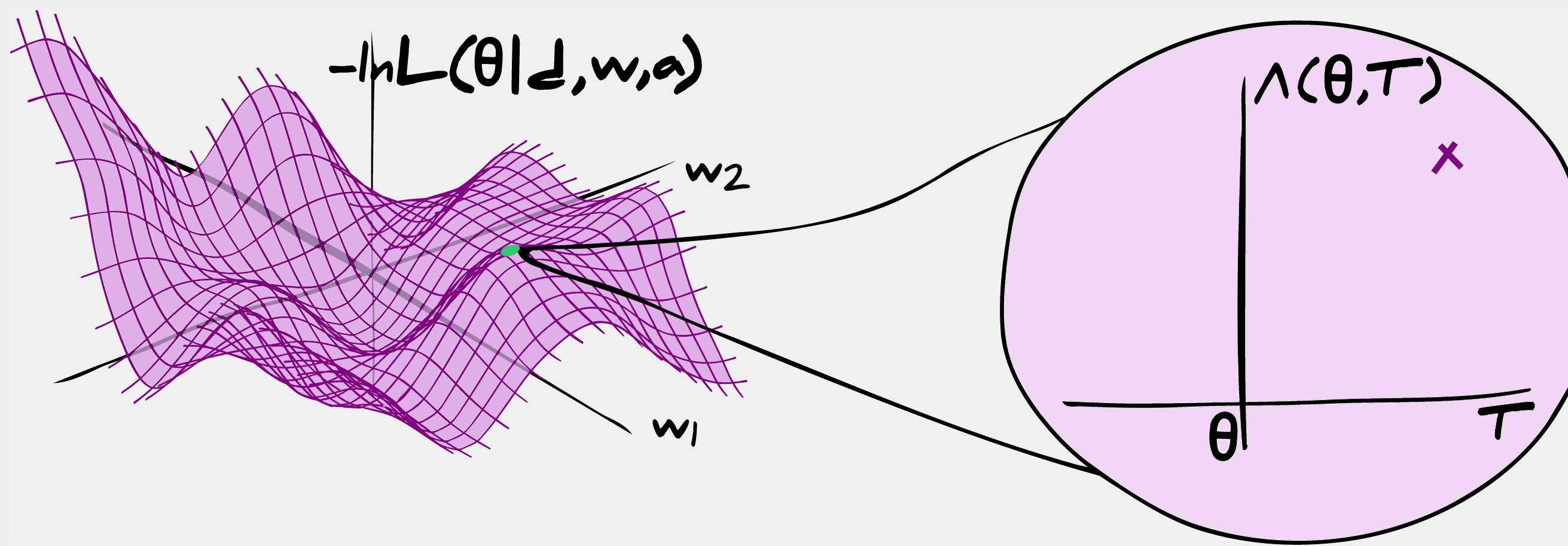
The prior distribution of weights and hyperparameters becomes the posterior

$$\begin{aligned}\mathcal{P}(\omega, \alpha | \{\mathbf{d}, \theta\}^{\text{train}}, \{\mathbf{d}, \theta\}^{\text{val}}) &\propto \\ \mathcal{L}(\omega, \alpha | \{\mathbf{d}, \theta\}^{\text{train}}, \{\mathbf{d}, \theta\}^{\text{val}}) p(\omega, \alpha)\end{aligned}$$

Training a network

Training a network

What are the maximum likelihood estimates of the weights?



$$\omega^{\text{MLE}} = \underset{\omega}{\operatorname{argmax}} \left[\mathcal{L}(\{\theta\}^{\text{train}} | \{d\}^{\text{train}}, \omega, \alpha^*) \right]$$

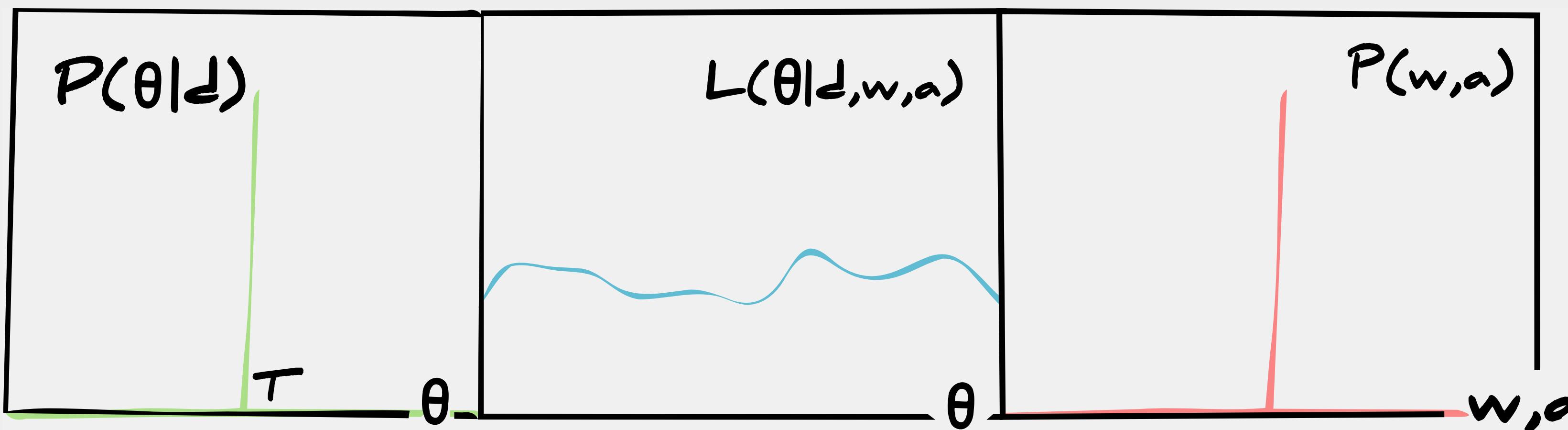
We degenerate the posterior

$$\mathcal{P}(\omega, \alpha | \{\mathbf{d}, \theta\}^{\text{train}}) \propto \mathcal{L}(\omega, \alpha | \{\mathbf{d}, \theta\}^{\text{train}}) p(\omega, \alpha)$$

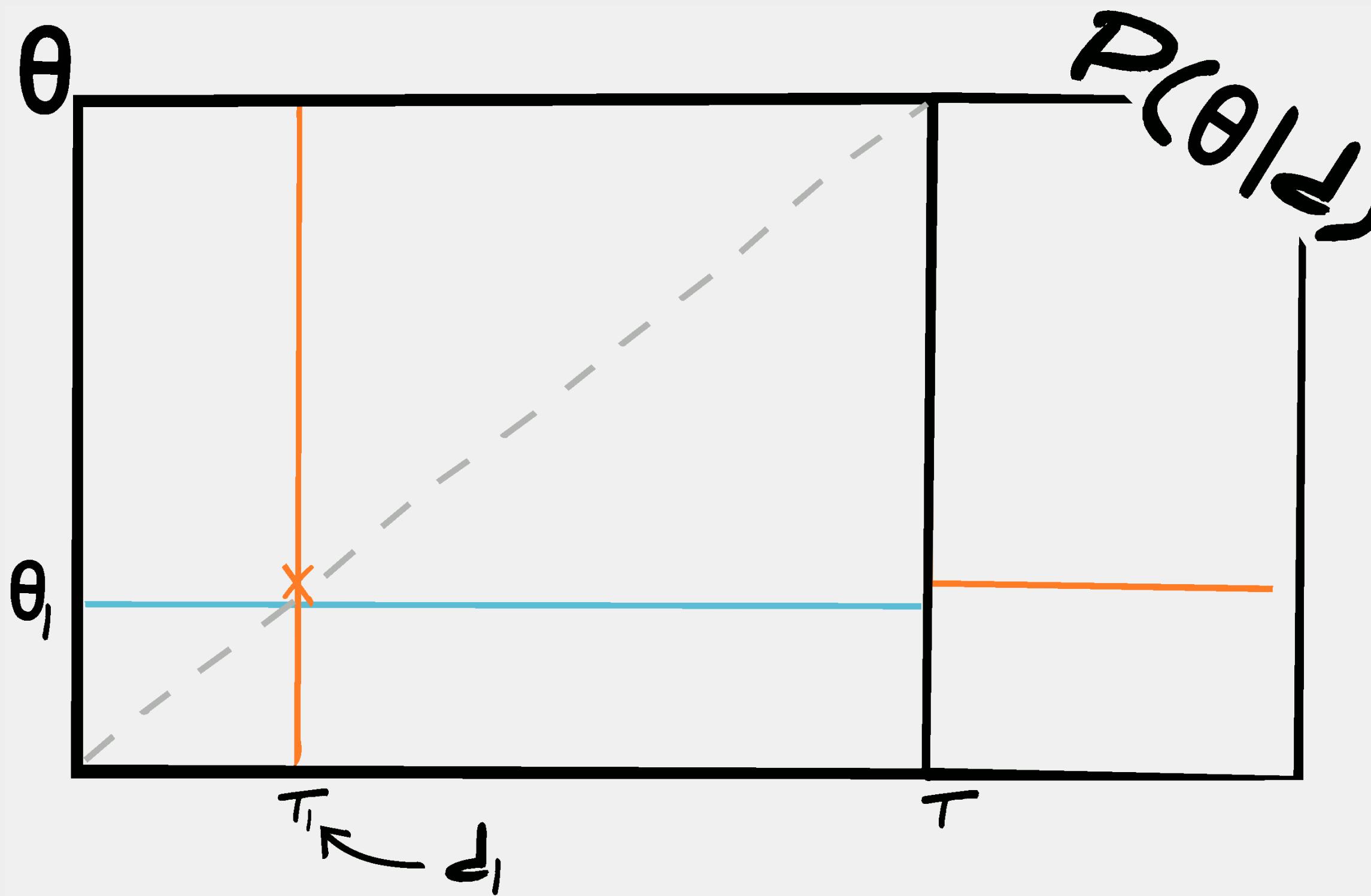
$$\rightarrow \delta(\omega - \omega^{\text{MLE}}, \alpha - \alpha^*)$$

$$\mathcal{P}(\theta | \mathbf{d}) \propto \int d\omega d\alpha \mathcal{L}(\theta | \mathbf{d}, \omega, \alpha) \delta(\omega - \omega^{\text{MLE}}, \alpha - \alpha^*)$$

$$= \delta(\tau)$$

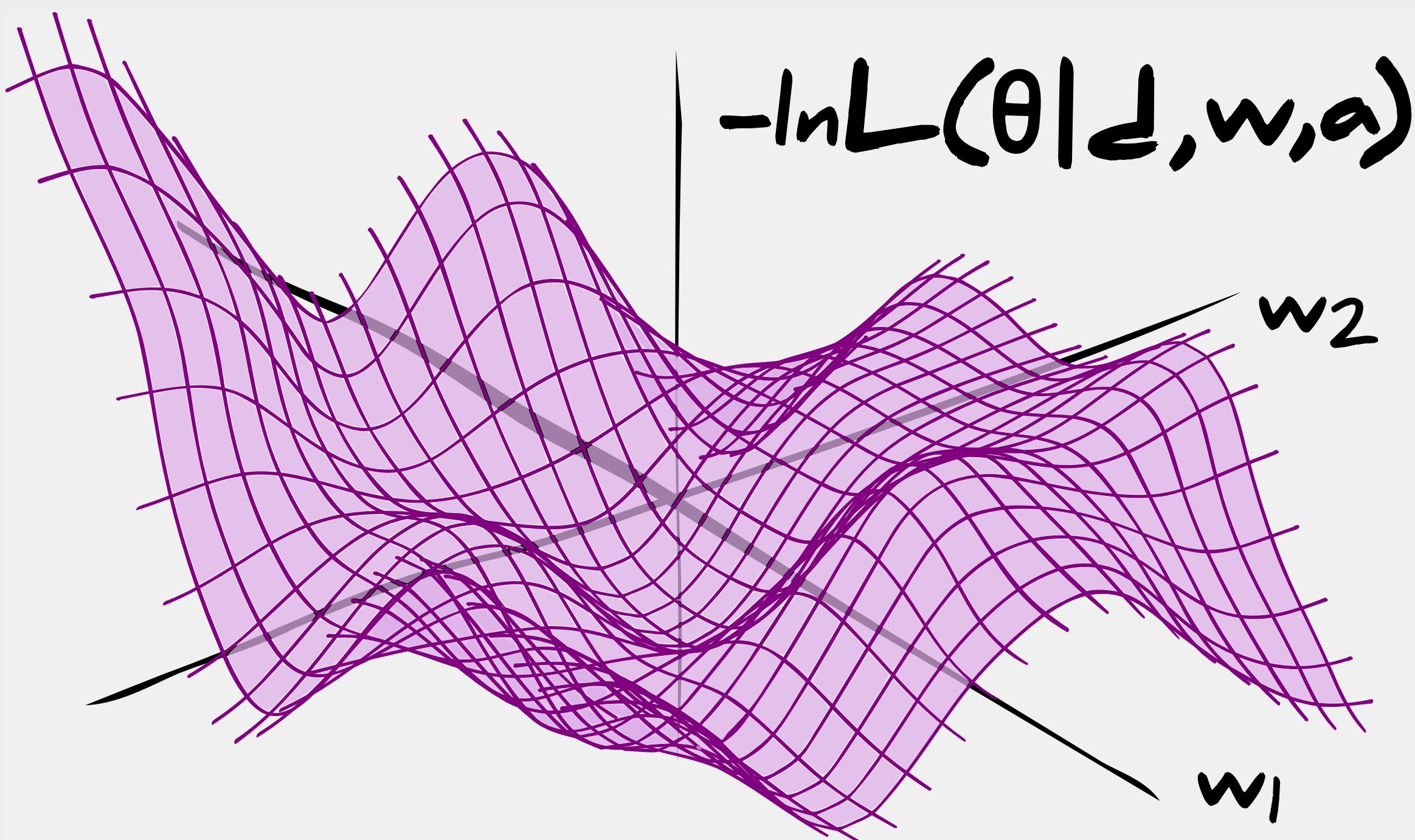


All predictions are (unknowably incorrect) estimates

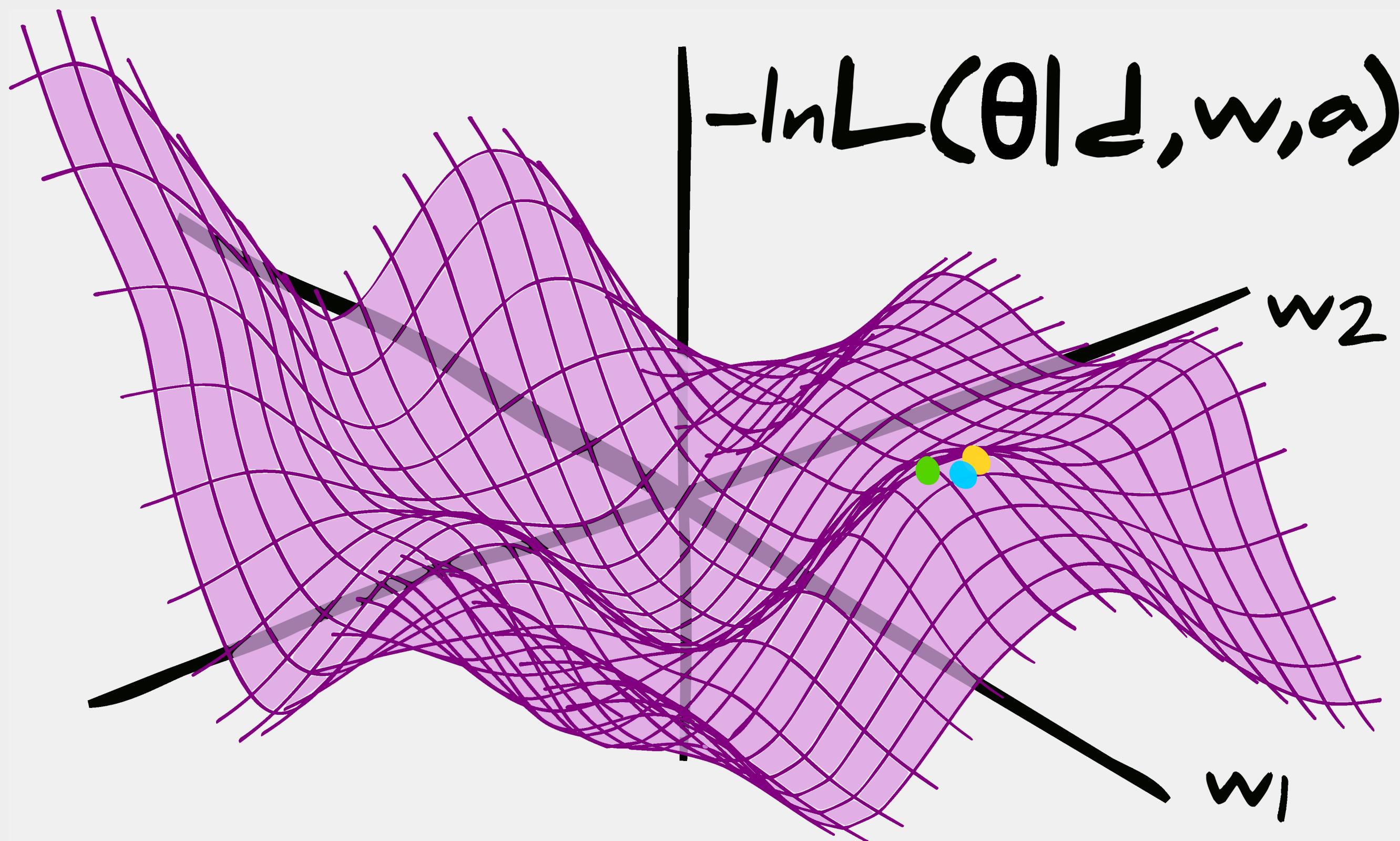


$$P(\theta|d) = \delta(\tau)$$

There is no way to interpret how close τ is to θ ...

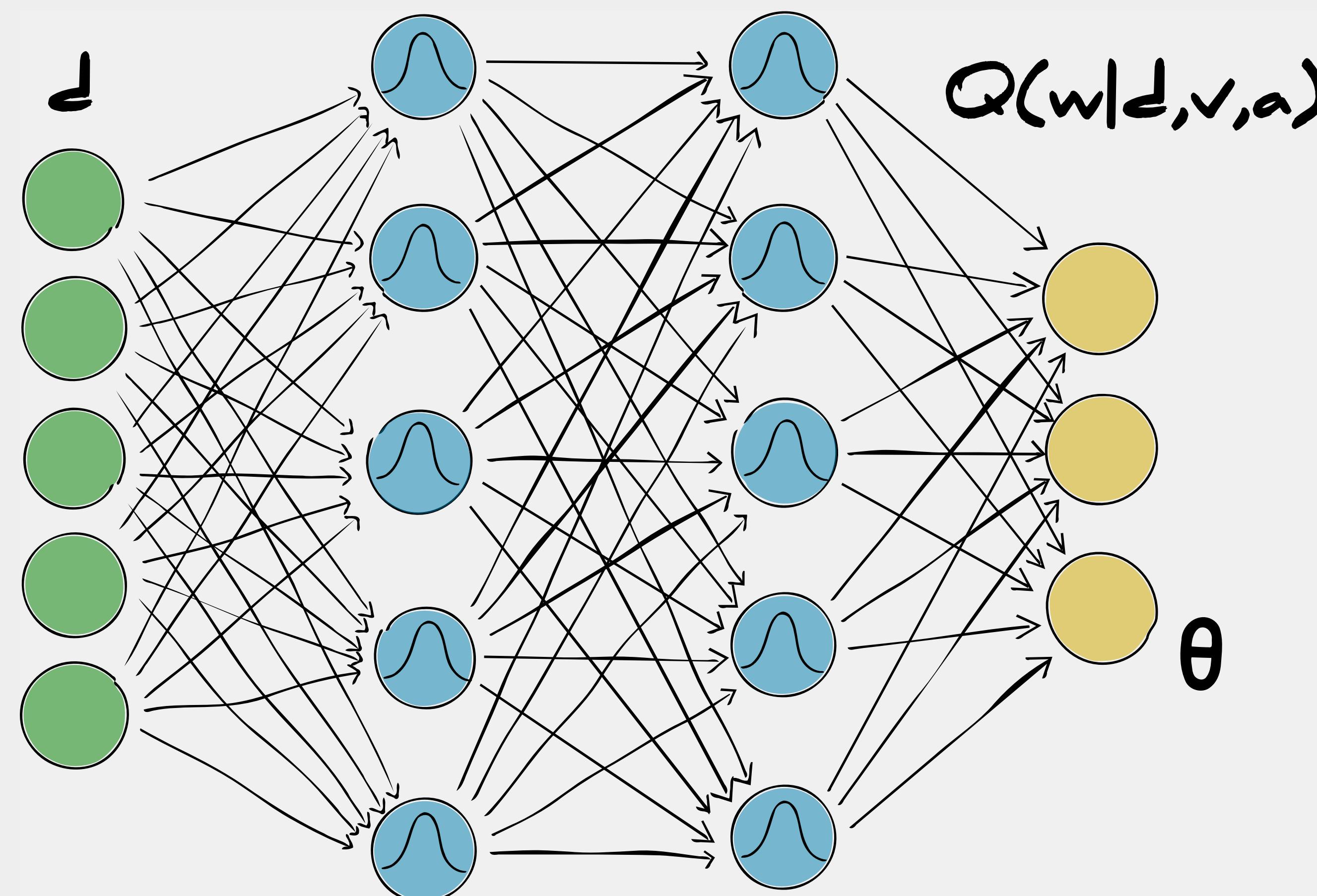


Local maximum likelihood estimates

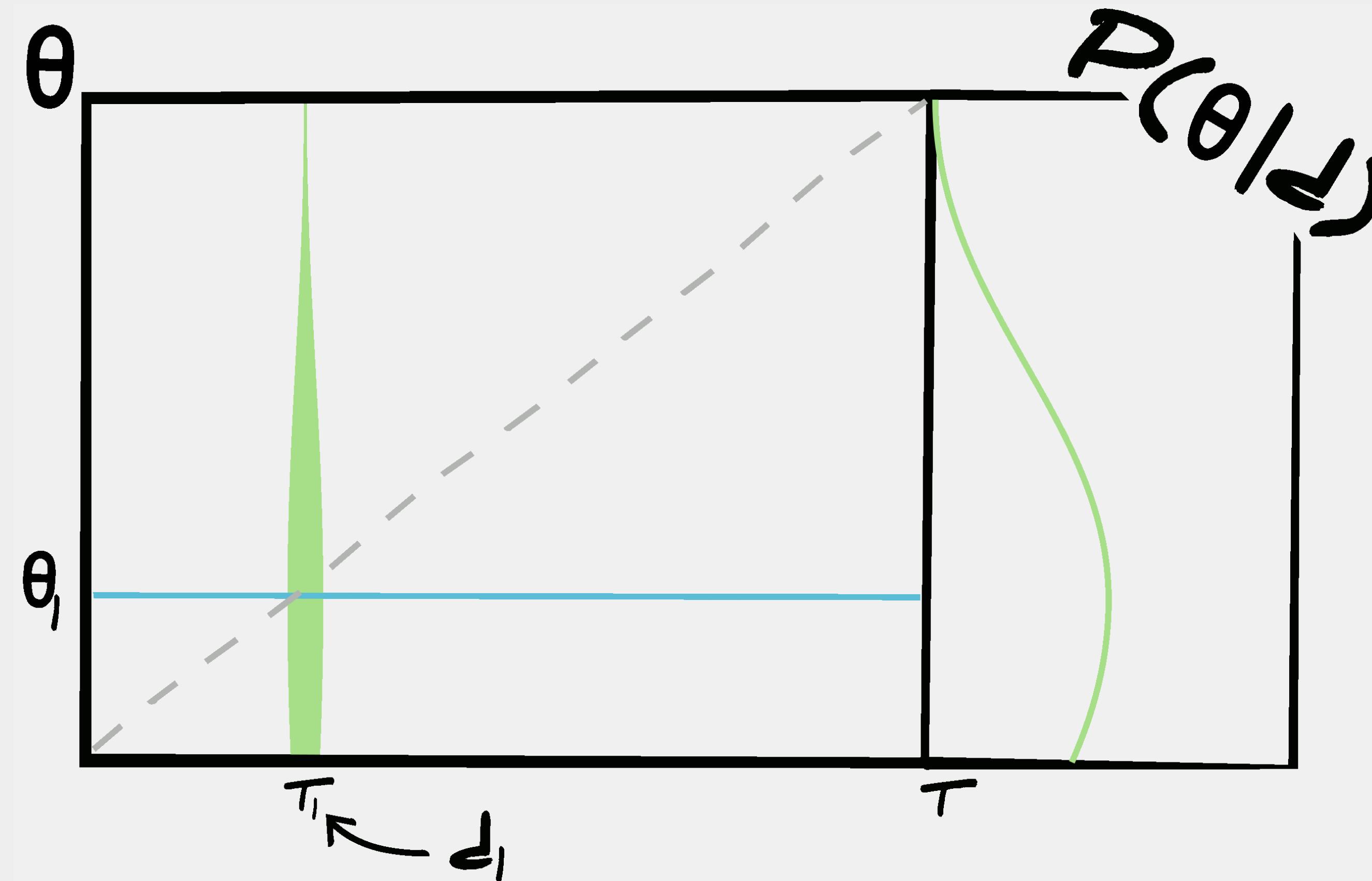


Are there better methods?

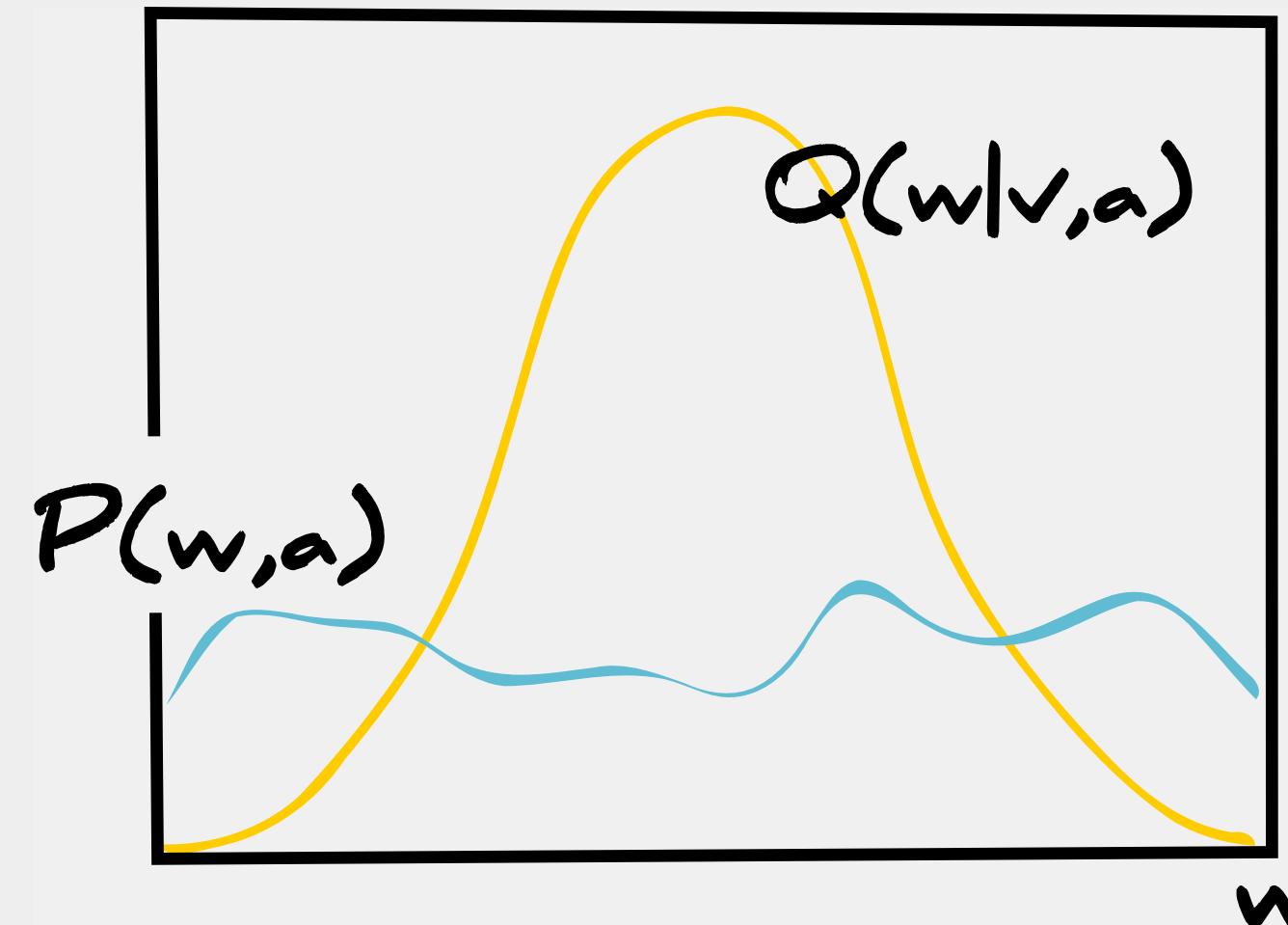
Variational inference



$$\mathcal{P}(\theta|\mathbf{d}) = \int d\omega d\nu d\alpha \mathcal{L}(\theta|\mathbf{d}, \omega, \alpha) Q(\omega|\nu, \alpha, \{\mathbf{d}, \theta\}^{\text{train}}) p(\nu, \alpha)$$



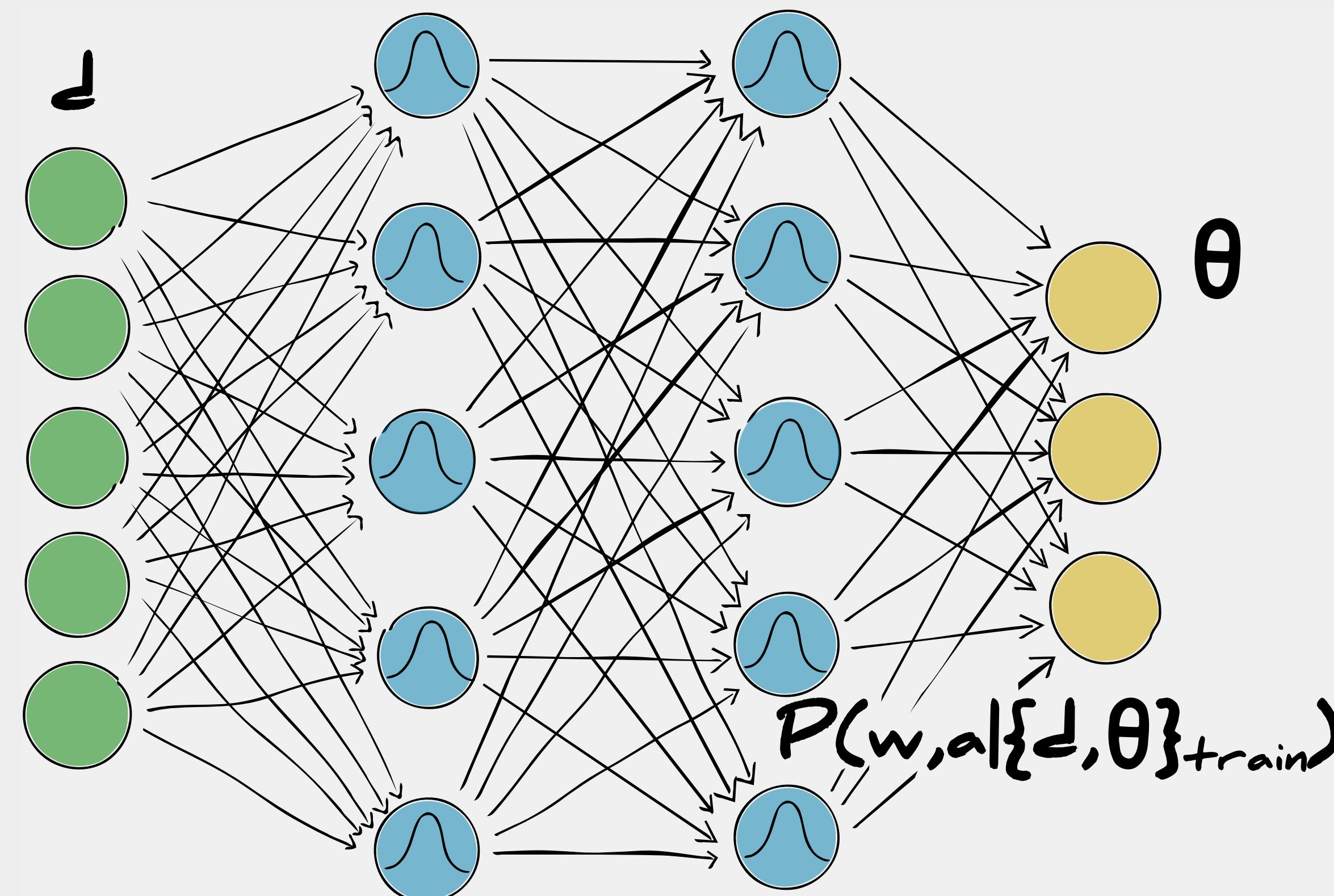
Still depends on fixed weights in the complex likelihood surface
and choice of variational distribution



$$\begin{aligned}\mathcal{P}(\theta|\mathbf{d}) &= \int d\omega d\nu d\alpha \mathcal{L}(\theta|\mathbf{d}, \omega, \alpha) Q(\omega|\nu, \alpha, \{\mathbf{d}, \theta\}^{\text{train}}) \\ &\quad \times \delta(\nu - \nu^{\text{MLE}}, \alpha - \alpha^*) \\ &= \int d\omega \mathcal{L}(\theta|\mathbf{d}, \omega, \alpha^*) Q(\omega|\nu^{\text{MLE}}, \alpha^*, \{\mathbf{d}, \theta\}^{\text{train}}).\end{aligned}$$

Bayesian neural networks

Bayesian neural networks



Approximate the posterior distribution of weights and hyperparameters

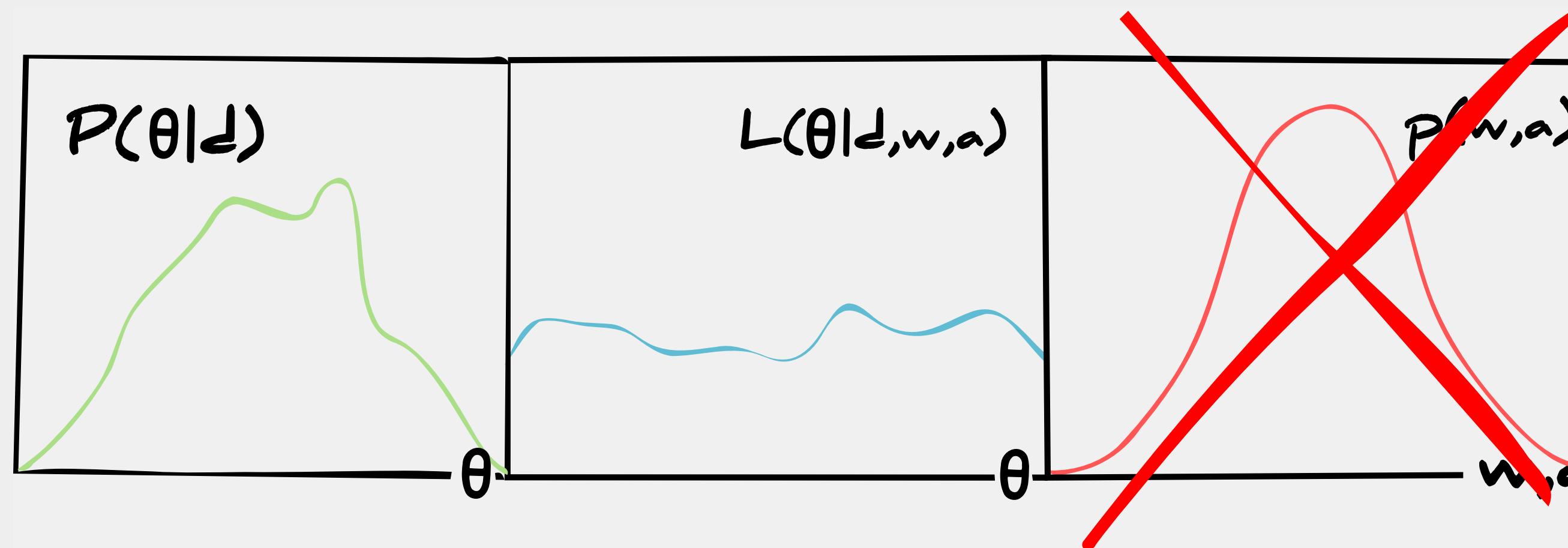
$$\begin{aligned}\mathcal{P}(\theta|\mathbf{d}) &= \int d\omega d\alpha \mathcal{L}(\theta|\mathbf{d}, \omega, \alpha) \mathcal{P}(\omega, \alpha|\{\mathbf{d}, \theta\}^{\text{train}}) \\ &\propto \int d\omega d\alpha \mathcal{L}(\theta|\mathbf{d}, \omega, \alpha) \\ &\quad \times \prod_i^{n_{\text{train}}} \mathcal{L}(\theta_i^{\text{train}}|\mathbf{d}_i^{\text{train}}, \omega, \alpha) p(\omega, \alpha).\end{aligned}$$

Everything is dependent on the training data!

Classical network : $\mathcal{P}(\omega, \alpha | \{d, \theta\}^{\text{train}}) \rightarrow \delta(\omega - \omega^{\text{MLE}}, \alpha - \alpha^*)$

Variational inference : $\mathcal{P}(\omega, \alpha | \{d, \theta\}^{\text{train}}) = \mathcal{Q}(\omega | \nu^{\text{MLE}}, \alpha^*, \{d, \theta\}^{\text{train}})$

Bayesian networks : $\mathcal{P}(\omega, \alpha | \{d, \theta\}^{\text{train}}) = \prod_i^{n_{\text{train}}} \mathcal{L}(t_i^{\text{train}} | \theta_i^{\text{train}}, \omega, \alpha) p(\omega, \alpha)$



If the training data, network optimisation, or
(pretty much) anything else is not perfect,
the inference will be biased

Parameter inference with summaries

Consider when the likelihood of some data, $\mathcal{L}(\mathbf{d}|\theta)$ is unknown

Consider when the likelihood of some data, $\mathcal{L}(\mathbf{d}|\theta)$ is unknown

Summarise the data, $f : \mathbf{d} \rightarrow \mathbf{t}$, to have a nicer likelihood for the model $\mathcal{M} \otimes f : \theta \rightarrow \mathbf{t}$

$$P(\theta|\mathbf{t}, \mathcal{M} \otimes f) = \frac{\mathcal{L}(\mathbf{t}|\theta, \mathcal{M} \otimes f)p(\theta|\mathcal{M})}{p(\mathbf{t}|\mathcal{M} \otimes f)}$$

The data no longer enters the inference, only the summaries

We can use a neural network as this summarising function

We can use a neural network as this summarising function

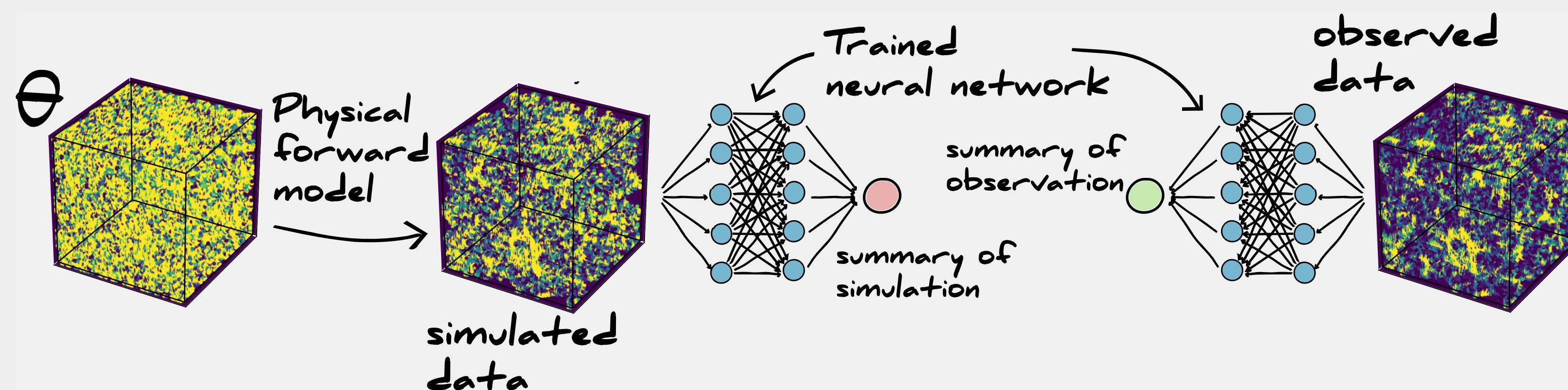
But the likelihood of the summaries is still difficult in this situation

Forward modelling

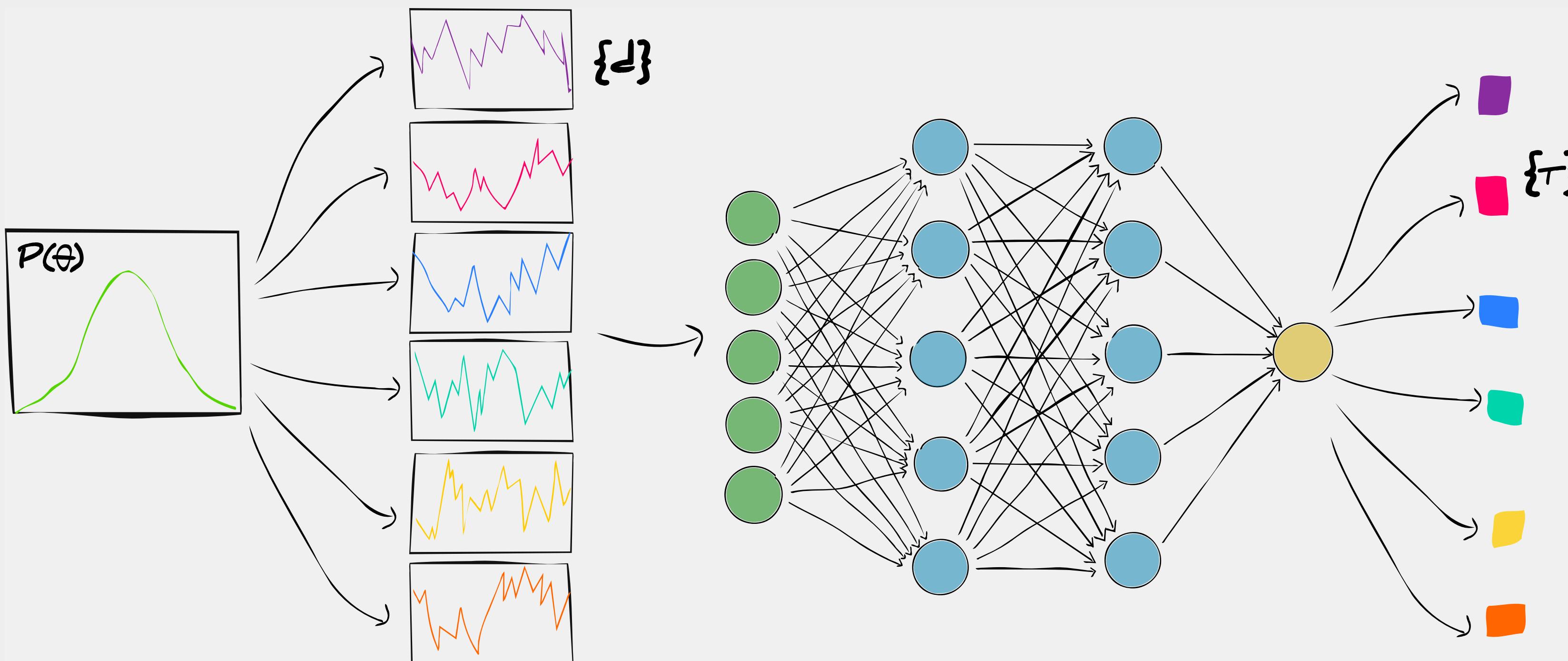
Method 1 :

Likelihood-free inference

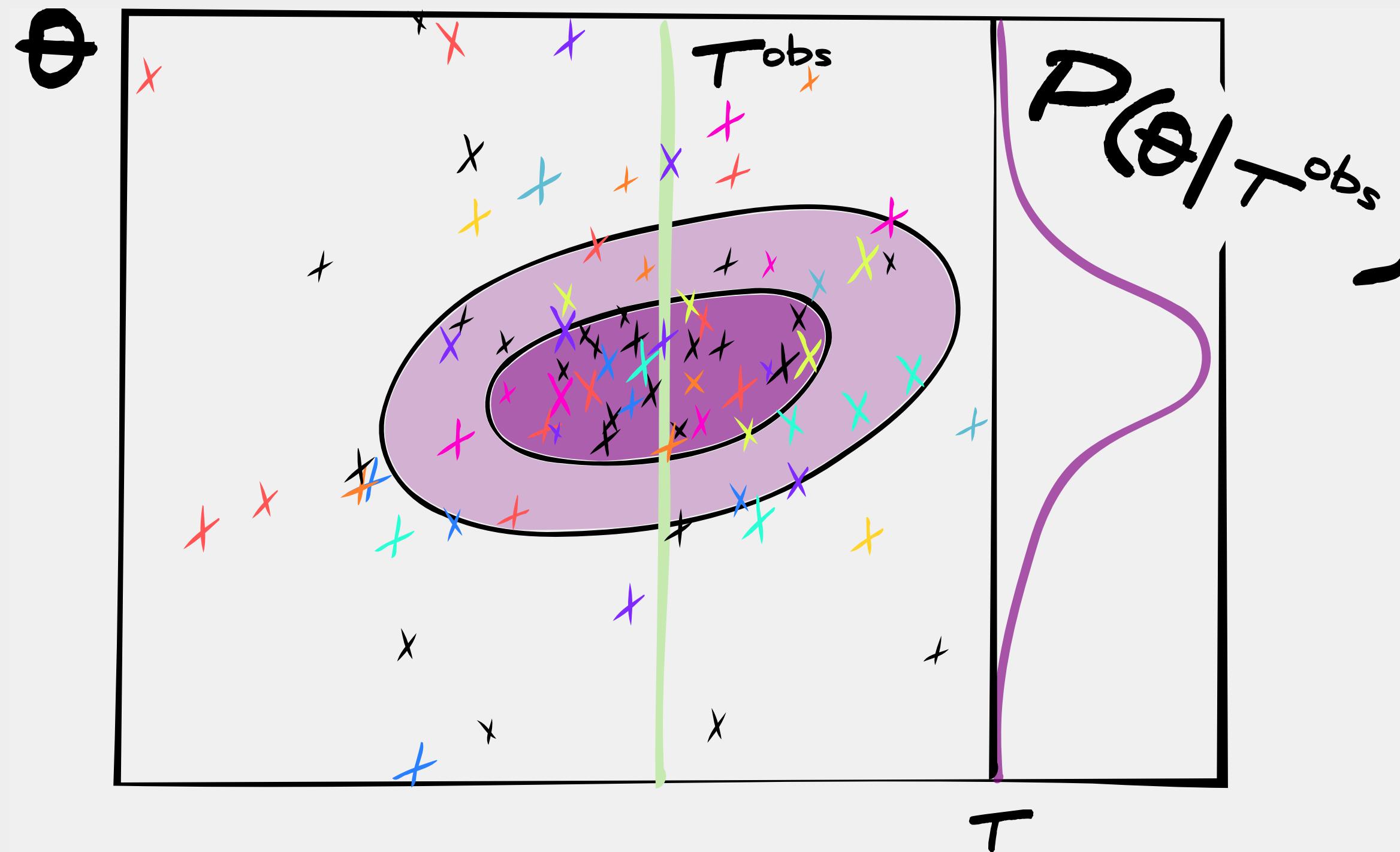
Use the same neural network to summarise observations and simulations



Generate simulations drawn from prior and summarise

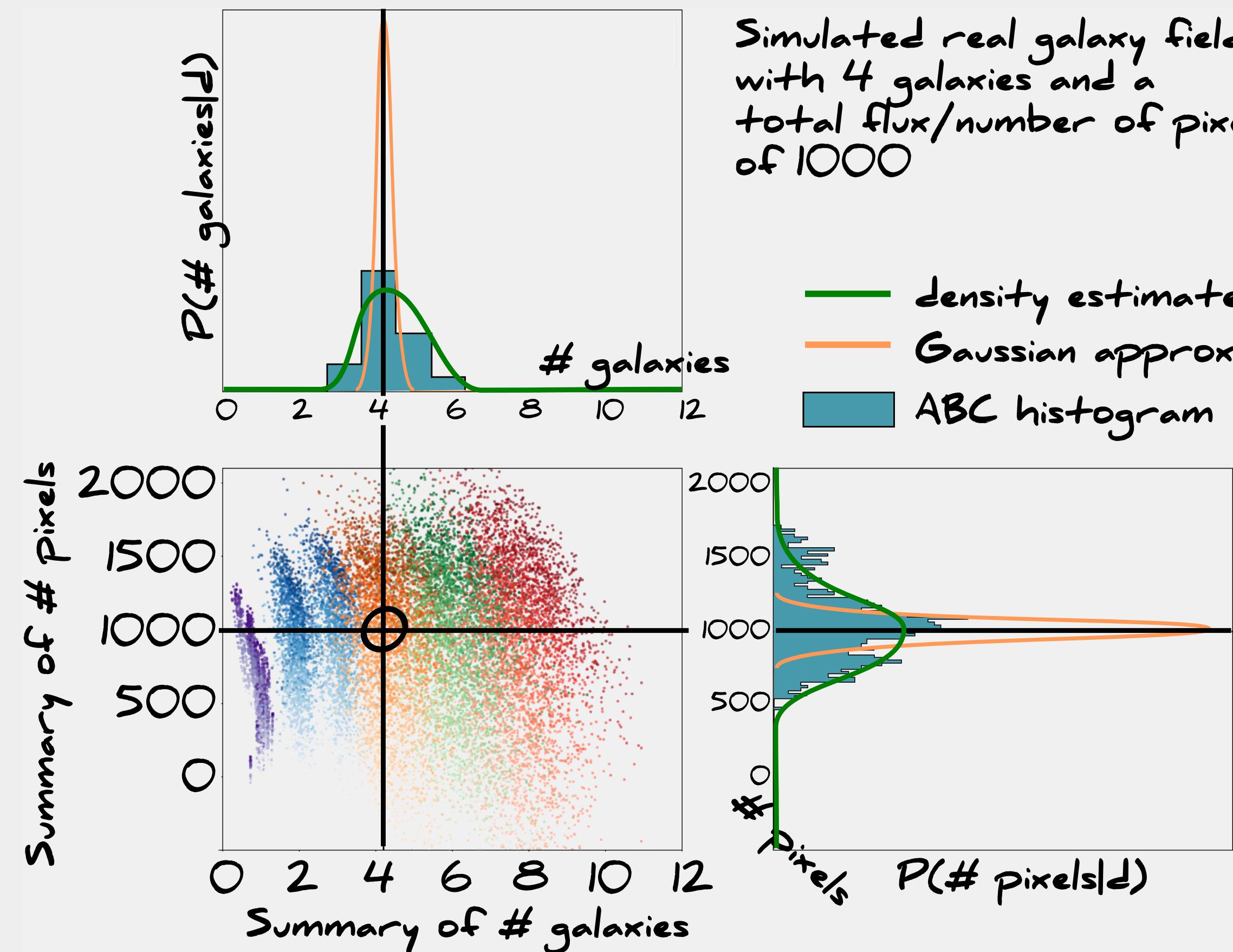


Make a kernel density estimate (or use ABC) on the simulations



A slice through at the summarised observation is the approximate posterior

Posterior distribution of galaxy counts and fluxes in fields



We're still biased by the form of the density estimator we use, and by the number of simulations we use for approximating the posterior

Method 2 :

Infer the data, physics and neural network

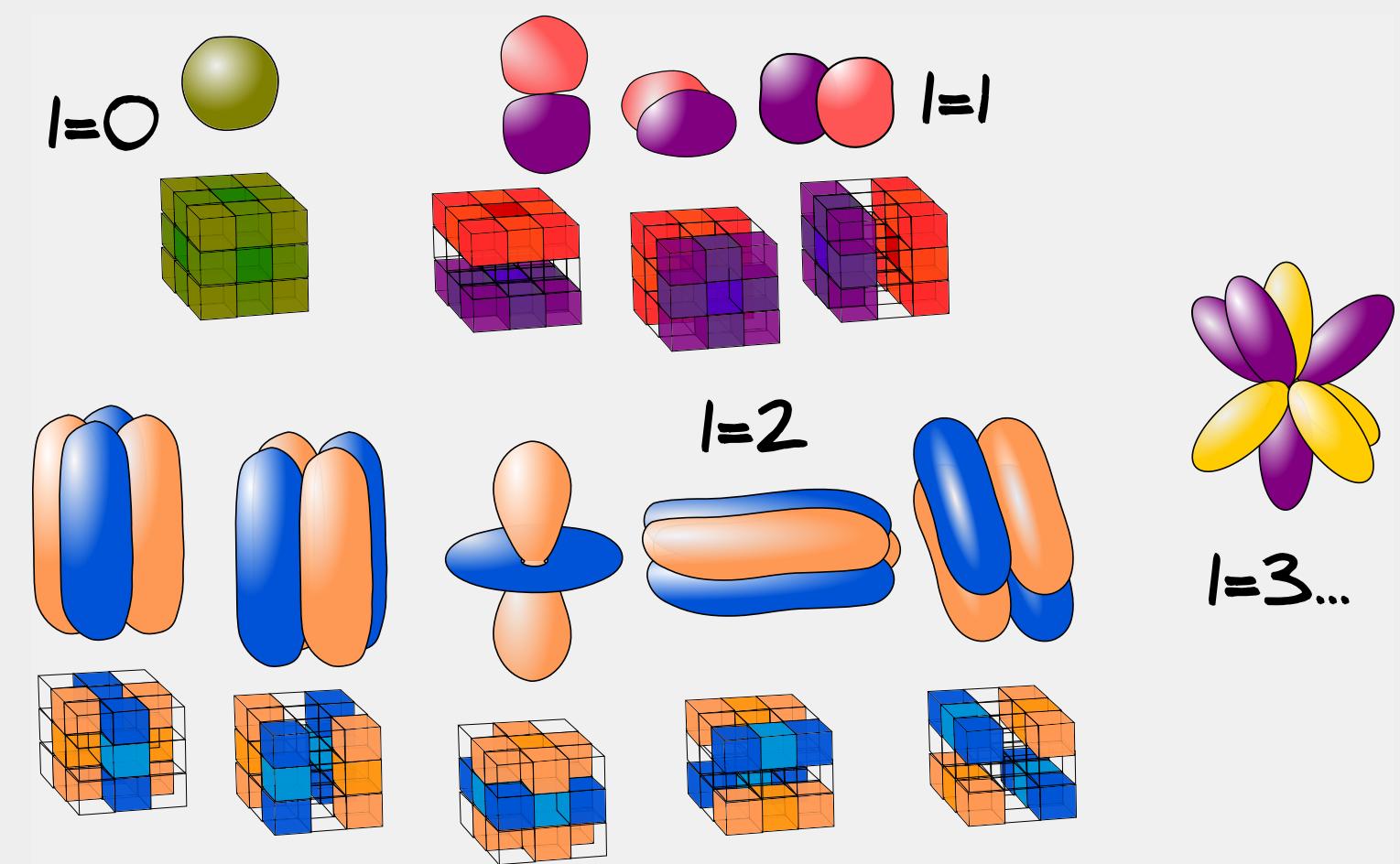
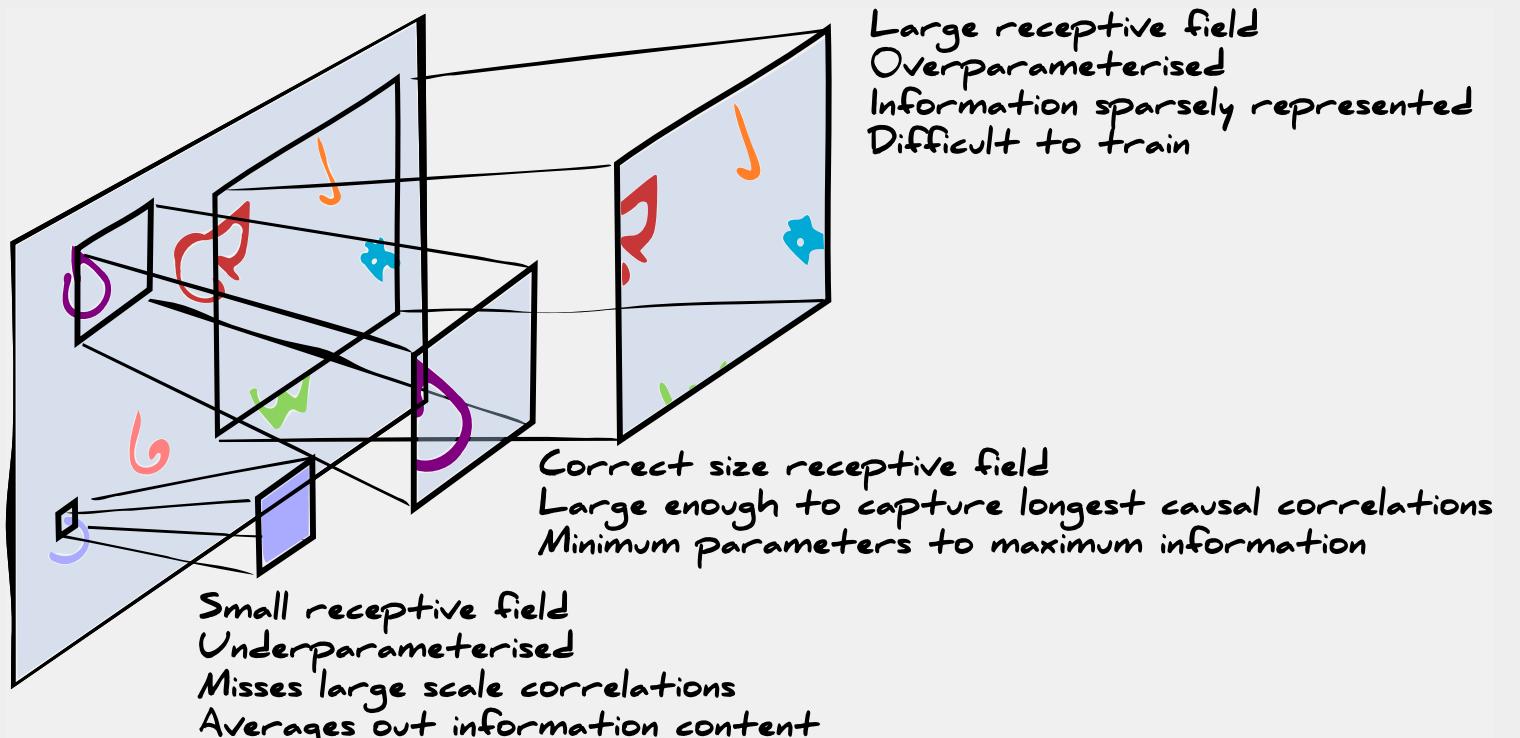
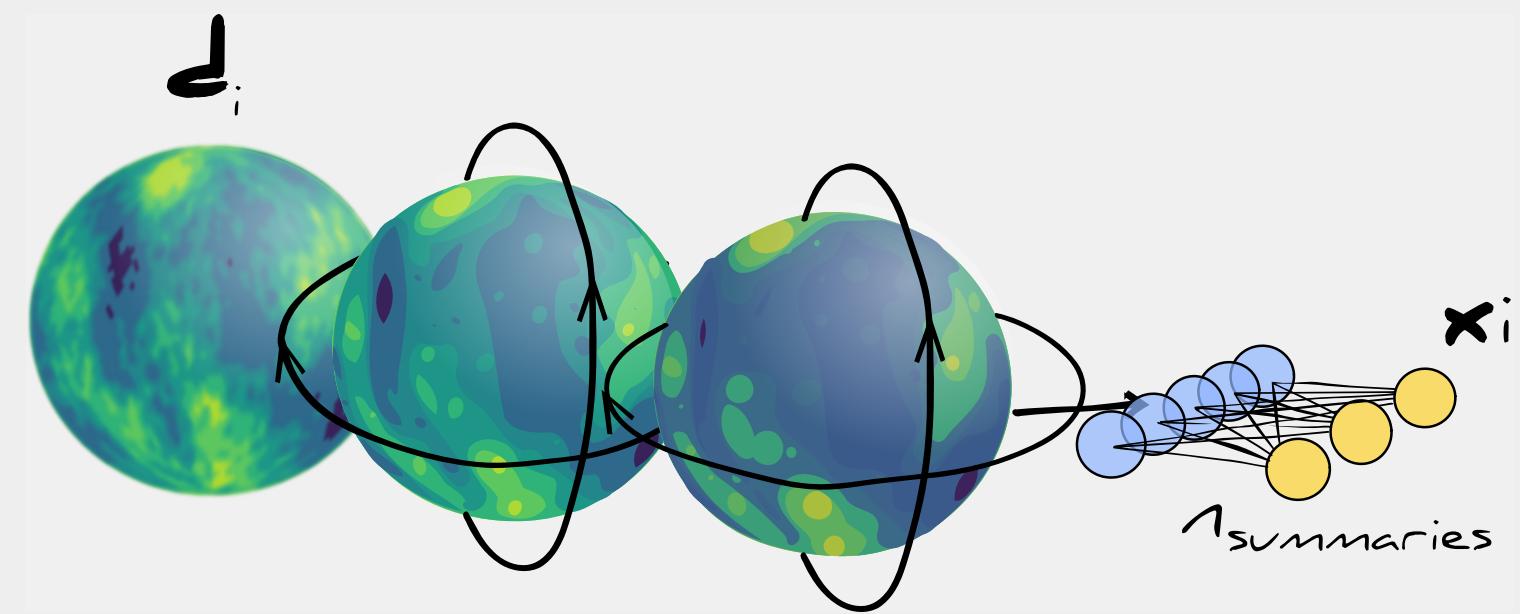
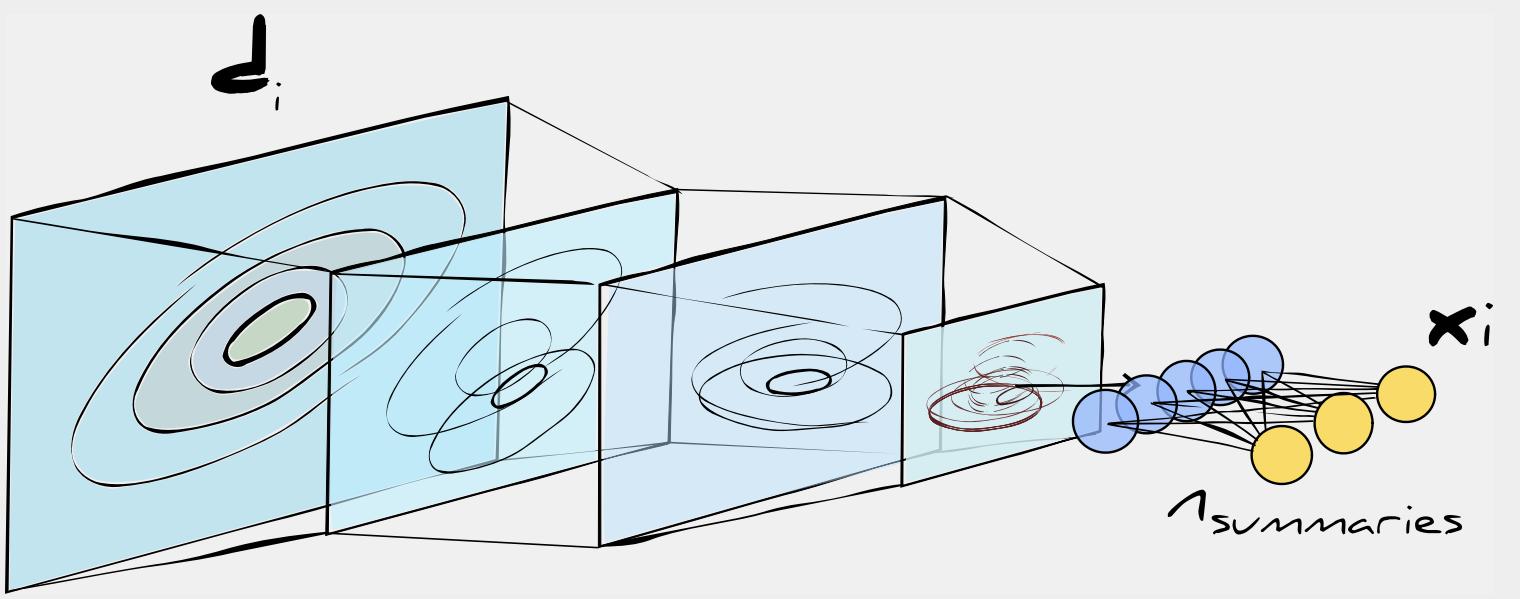
Need a good control on the likelihood of weights and hyperparameters

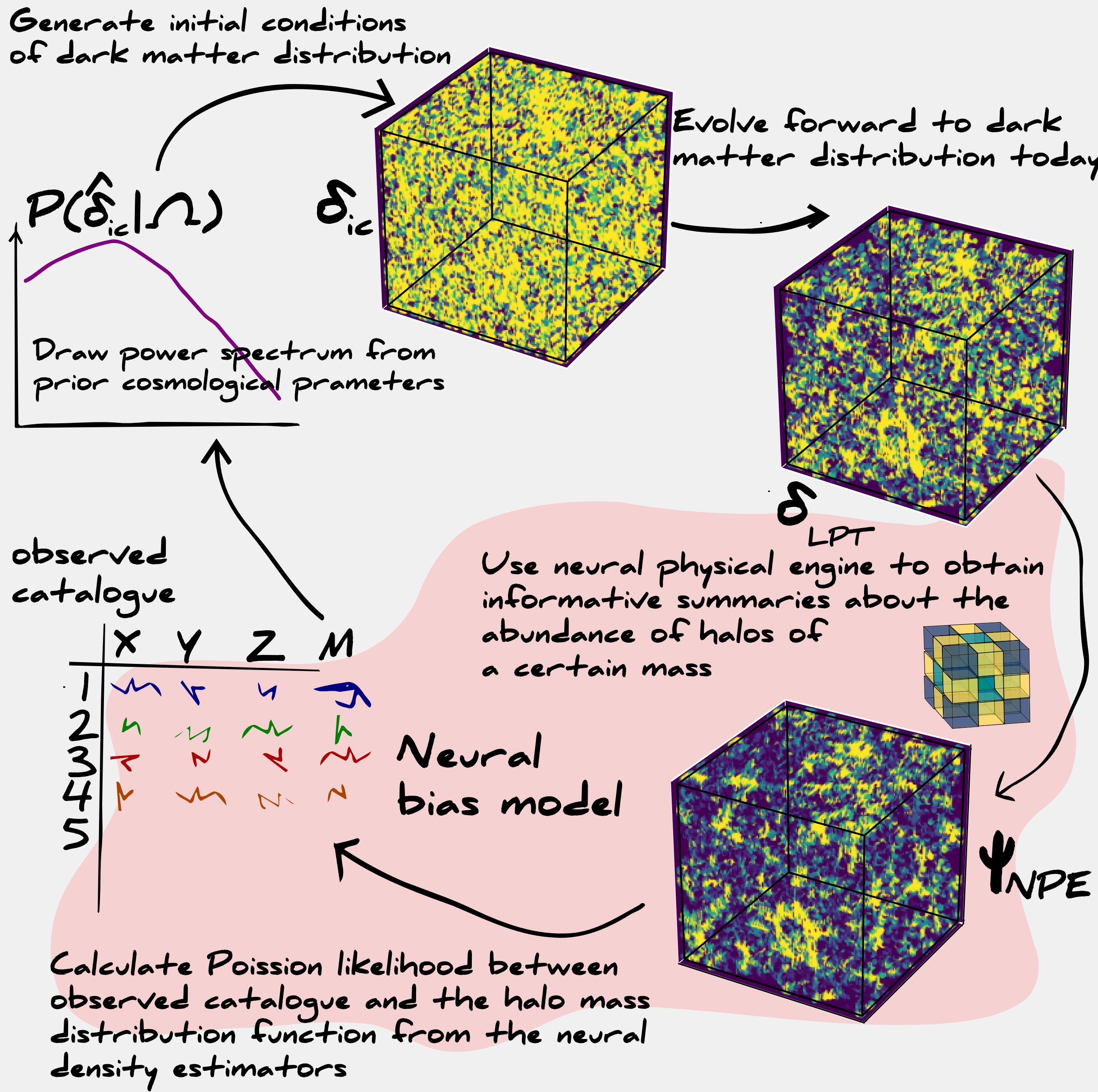
- remove degeneracies
- force to be convex
- etc.

Need a good control on the likelihood of weights and hyperparameters

- remove degeneracies
- force to be convex
- etc.

Use neural physical engines





Conclusions

Conclusions

Doing parameter estimation with neural networks is not as easy as predicting parameter values

We can still use results from neural networks as highly informative summaries

We can also infer all parameters of the forward model, including network to get unbiased results

Questions

When do we say we trust a neural network?

How approximate is too approximate?

Are using neural networks for parameter estimation worse than using other summary statistics?