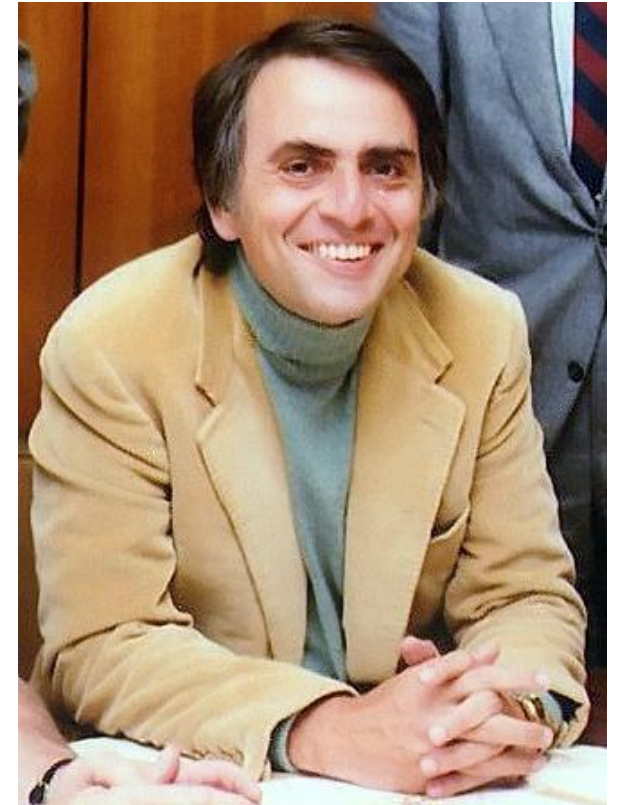


Lecture-6: Uncertainty in Neural Networks

Motivation

“Every time a scientific paper presents a bit of data, it's accompanied by an error bar – a quiet but insistent reminder that no knowledge is complete or perfect. It's a calibration of how much we trust what we think we know.”

— **Carl Sagan**,
The Demon-Haunted World: Science as a Candle in the Dark



https://commons.wikimedia.org/wiki/File:Carl_Sagan_Planetary_Society.JPG

Further Reading

Sagan, Carl. *The demon-haunted world: Science as a candle in the dark.*
Ballantine Books, 2011.

Types of Uncertainty

Classical Neural Network (or any parametric model)

$$y = f(x, \theta) ; D = \{x, t\}$$

Conventions for this lecture

- θ = neural network weights
- x = neural network input
- y = neural network output
- D = training data (input and targets)

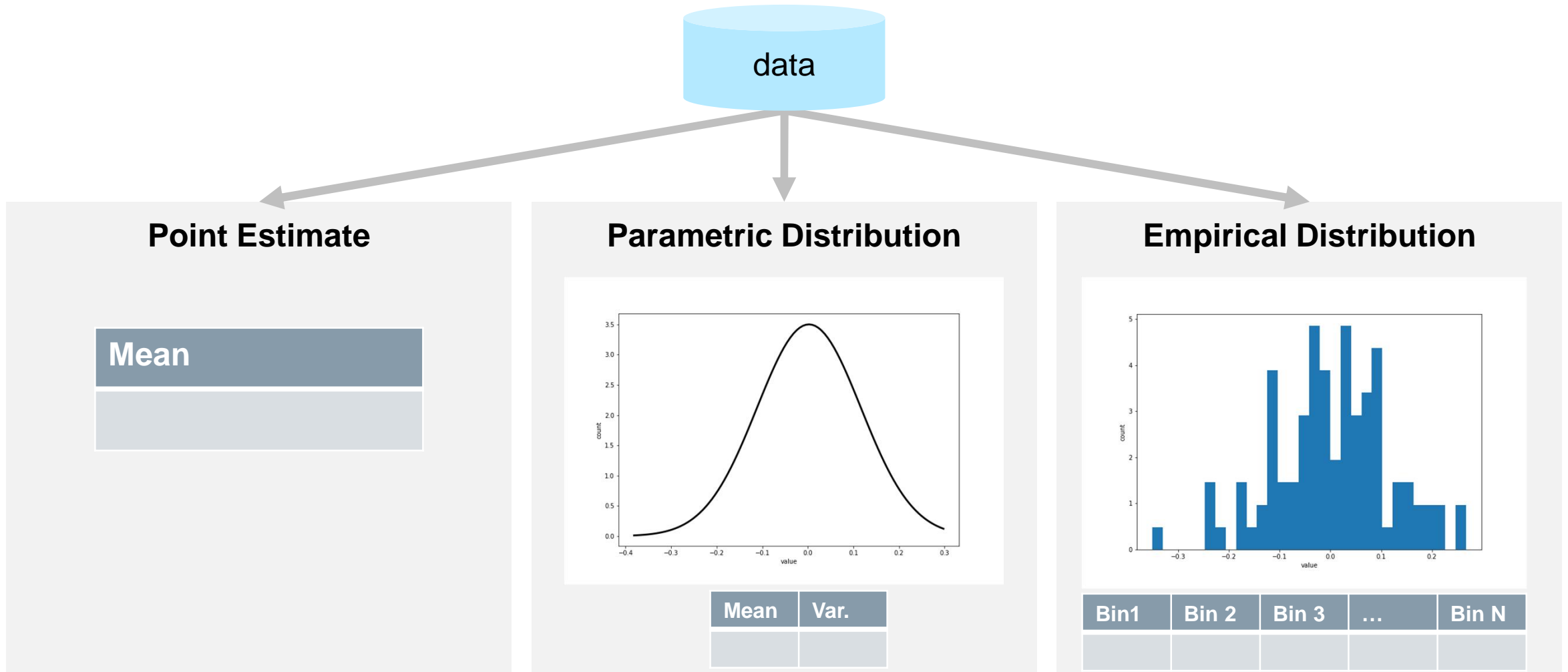
Output/Prediction Uncertainty

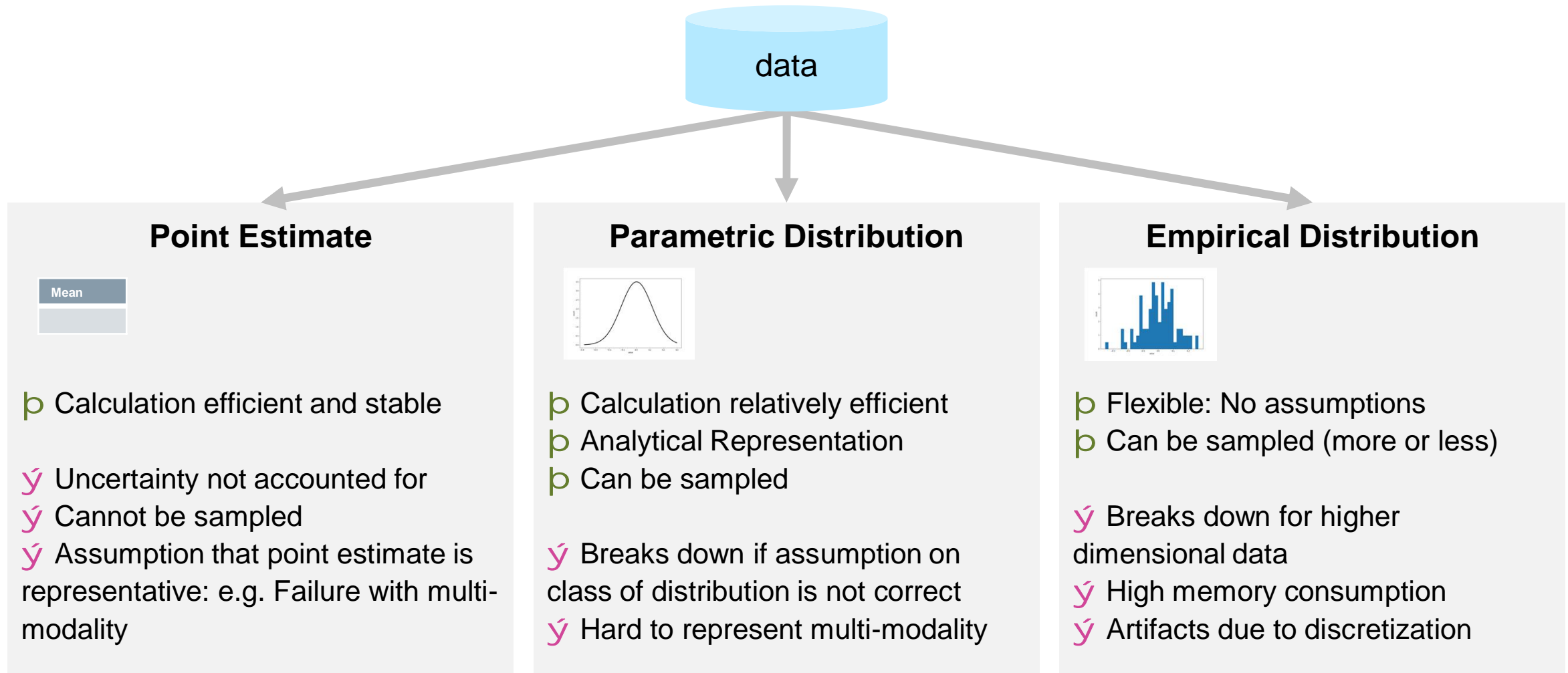
$$y \sim p(y|x, \theta)$$

Model (Weight) Uncertainty

$$\theta \sim p(\theta|D)$$

Three Different Ways to represent Distributions in Practice



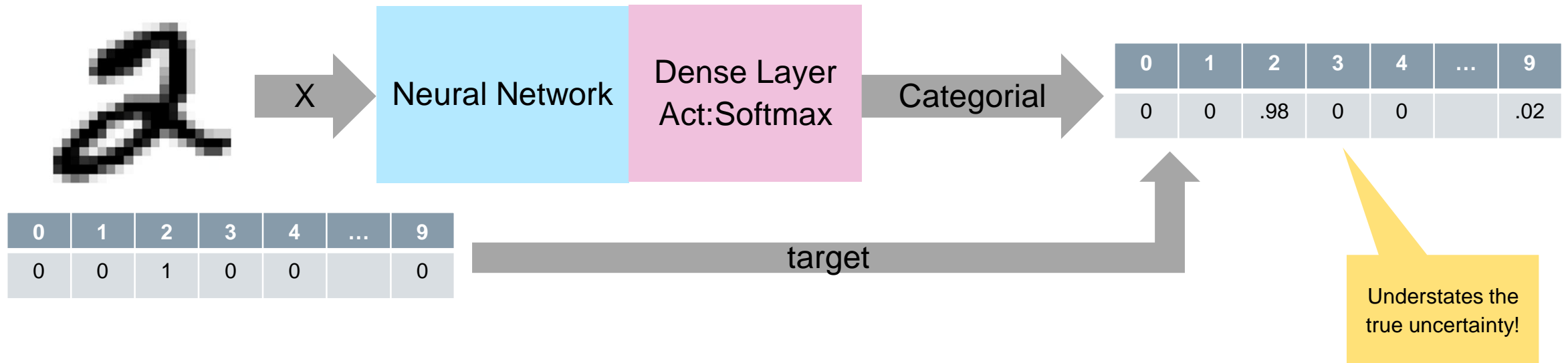


Output Uncertainty for Regression

Recap: Classification

The standard setup:

- Classification predicts a categorical distribution of the classes
- Targets are one-hot-encoded
- The loss function is the cross-entropy



We will focus on Regression henceforth ...

Motivation for Output Uncertainty: Mean Average Precision @ 7

Challenge:

Give at most 7 recommendations from a large product catalog. You are evaluated based on the precision on these at most 7 items.

Idea:

Select at most 7 recommendations where the model is certain about the recommendation.

In General:

Situations where a response/action can be rejected.

Further Reading

Brando, Axel, et al. "Uncertainty Modelling in Deep Networks: Forecasting Short and Noisy Series." arXiv preprint arXiv:1807.09011 (2018).

Related Kaggle Challenge

<https://www.kaggle.com/c/santander-product-recommendation>

Neural Network with Output Uncertainty

$$y \sim p(y|x, \theta)$$

Let's commit to a parametric distribution:

$$y \sim \mathcal{N}(y|\mu, \sigma)$$

Your choice!
Also popular: The
Laplace Distribution

We will model μ as a Neural Network: $\mu(x, \theta)$

We either model σ as a scalar parameter under the assumption of homoskedastic uncertainty or as a Neural Network: $\sigma(x, \theta)$ for heteroskedastic uncertainty

$$\mathcal{N}(y|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma(x, \theta)^2}} \exp\left(-\frac{(y - \mu(x, \theta))^2}{2\sigma(x, \theta)^2}\right)$$

Neural Network with Output Uncertainty

$$\mathcal{N}(y|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma(x, \theta)^2}} \exp\left(-\frac{(y - \mu(x, \theta))^2}{2\sigma(x, \theta)^2}\right)$$

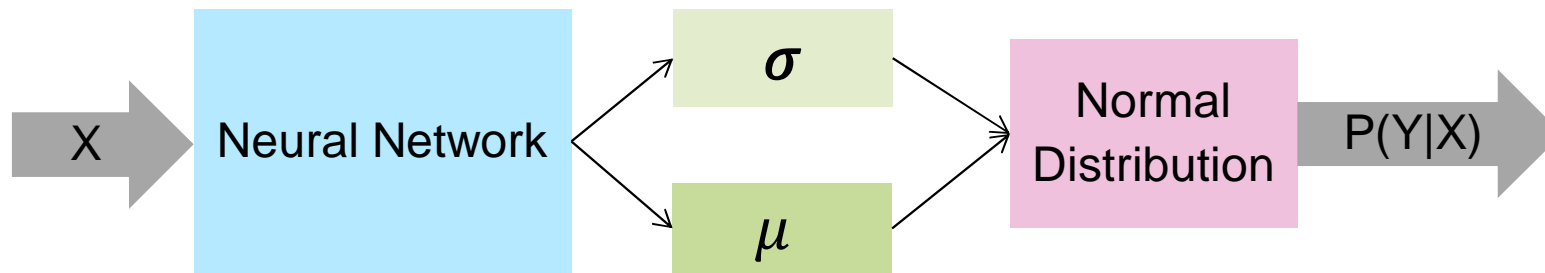
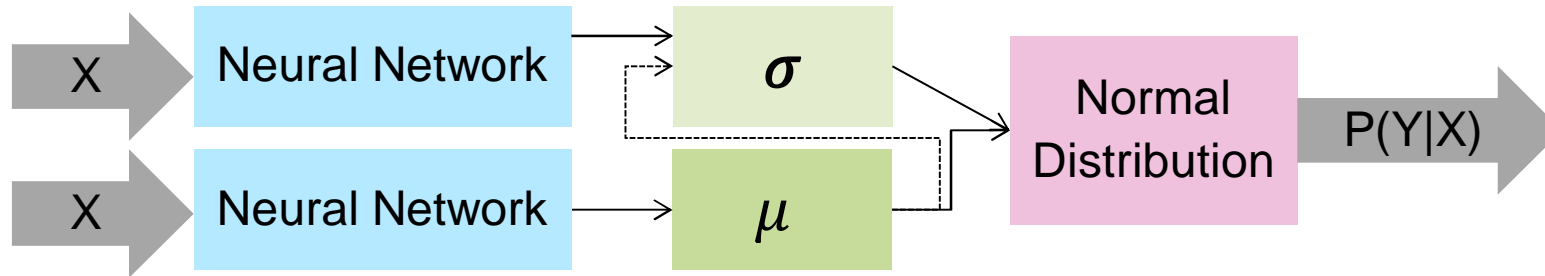
$$D = \{x, t\}$$

We will optimize the Log-Likelihood with respect to the weights θ

$$\mathcal{L}(\theta, D) = - \sum_{i=0}^{|D|} \left[-\log(\sigma(x_i, \theta)) - \frac{(t_i - \mu(x_i, \theta))^2}{2\sigma(x_i, \theta)^2} + C \right]$$

Use a variant of Stochastic Gradient Descent to train.

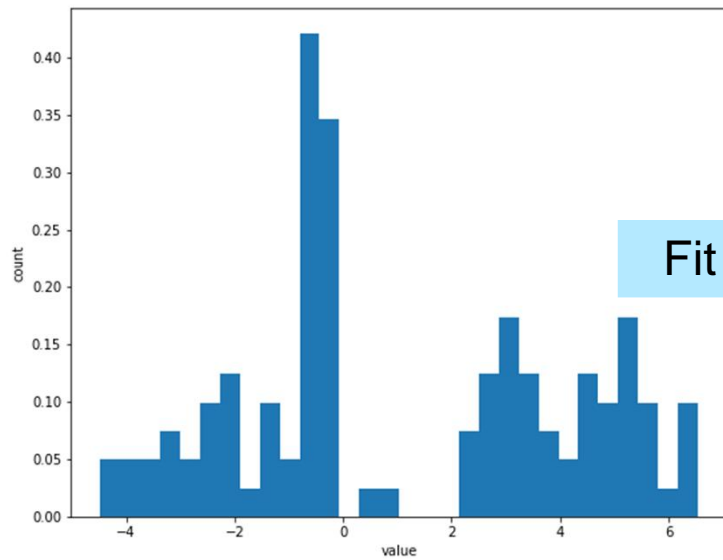
Network with Output Uncertainty: Architecture Variants



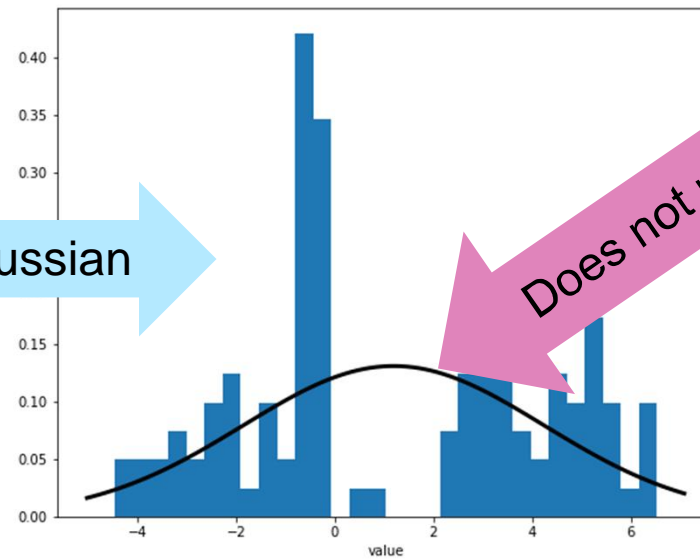
Loss: $-\log P(\mathbf{Y}|\mathbf{X})$

Understates the true uncertainty!

Multimodality



Fit Gaussian



Does not make sense

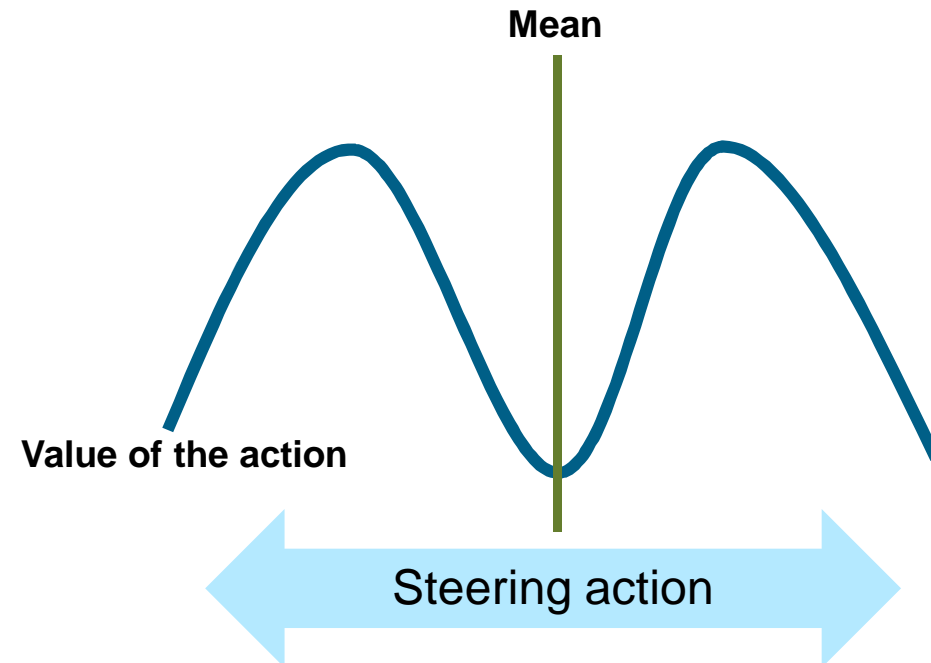
Example of Multimodality

Learning the value of an action
(see also *Reinforcement Learning*)

Further Reading
Depeweg, Stefan, et al. "Learning and policy search in stochastic dynamical systems with bayesian neural networks." arXiv preprint arXiv:1605.07127 (2016).



https://commons.wikimedia.org/wiki/File:Newport_Whitepit_Lane_pot_hole.JPG



Result of Point Estimation



https://commons.wikimedia.org/wiki/File:Bus_in_hole.jpg

Inverse Problems

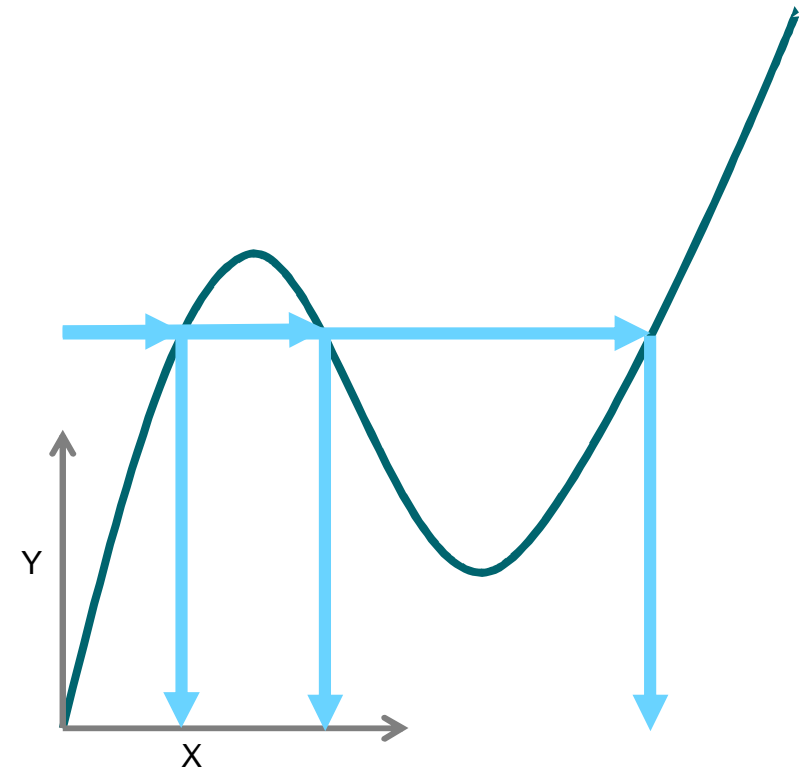
1. We are interested in X , but X can only be perceived indirectly via Y .
2. The relation between X and Y is given by a known, possibly stochastic function

$$y = g(x)$$

This means we would like to learn the inverse of g :

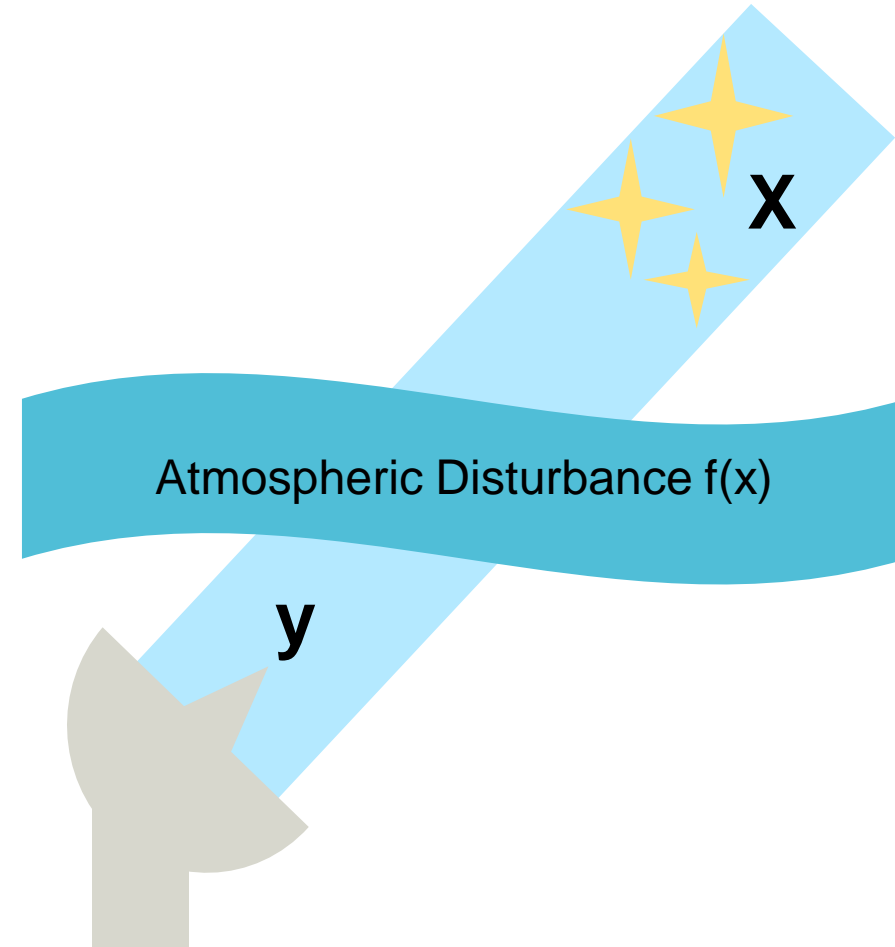
$$x = f(y)$$

The challenge is that g may not be invertible!
This means that f may be stochastic even if g is not.

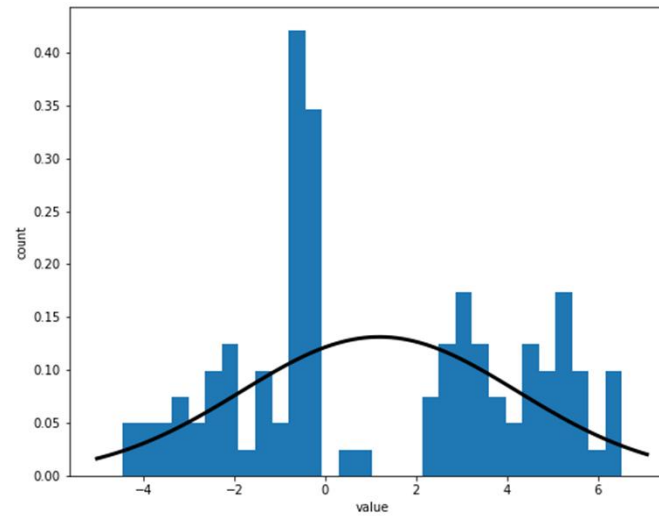
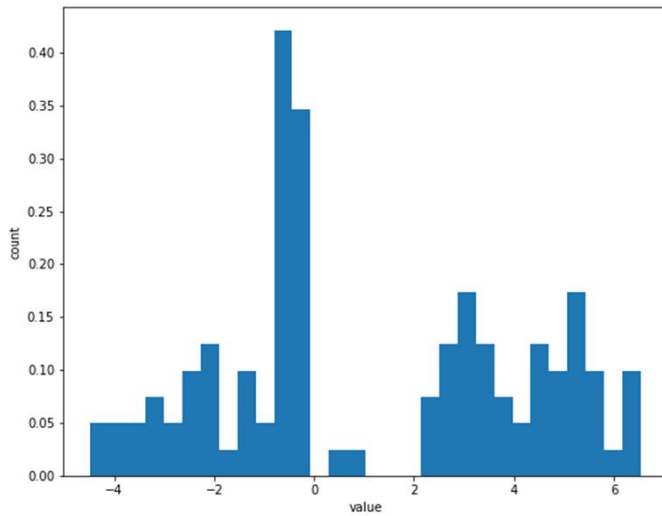


Inverse Problems: Examples

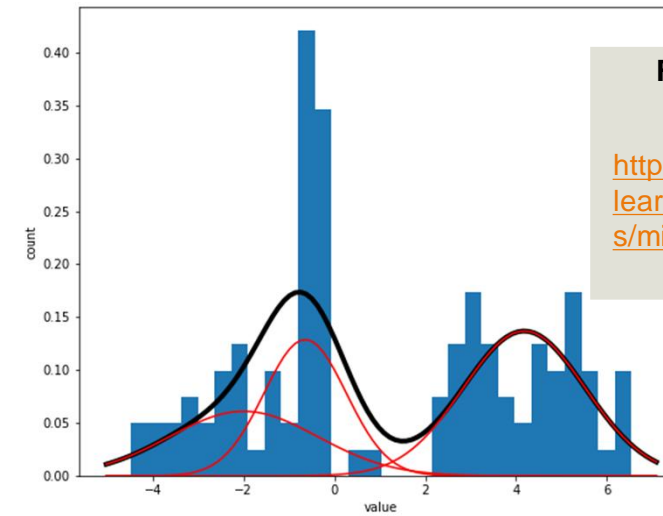
1. Image De-Noising
2. Astronomy Imaging
3. Deconvolution
4. Localization
5. ...



Gaussian Mixtures



Mean	Variance
1.2	3



	Mean	Var.	weight
Cluster 1	4.1	1.8	0.6
Cluster 2	-2.0	2.7	0.38
Cluster 3	-0.6	0.8	0.02

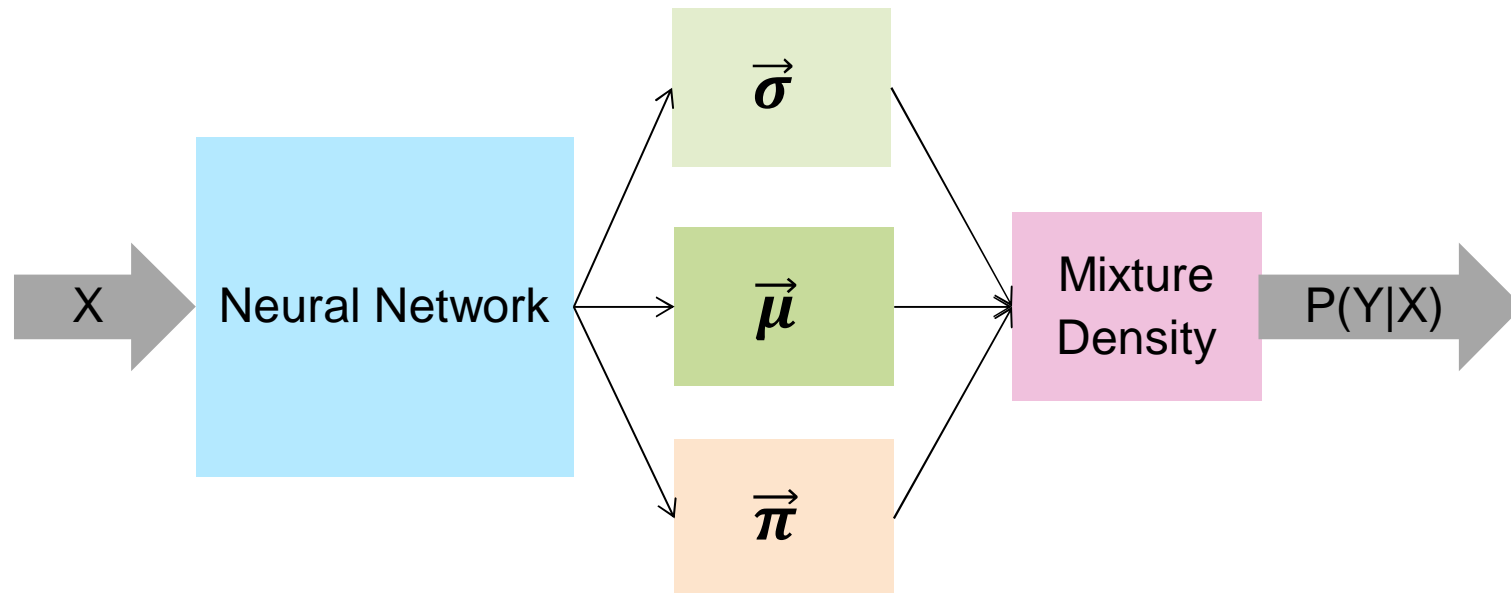
Further Practice with GMMs
<http://scikit-learn.org/stable/module/mixture.html#bgmm>

If you want to know how many clusters to assume:

Further Reading
Rasmussen, Carl Edward. "The infinite Gaussian mixture model." Advances in neural information processing systems. 2000.

Further Reading
Blei, David M., and Michael I. Jordan. "Variational inference for Dirichlet process mixtures." Bayesian analysis 1.1 (2006): 121-143.

The Mixture Density Network: Architecture



- $\vec{\sigma}$: standard deviations
- $\vec{\mu}$: means
- $\vec{\pi}$: cluster weights
- **Loss:** $-\log P(\mathbf{Y}|\mathbf{X})$

Reference

Bishop, Christopher M. *Mixture density networks*. Technical Report NCRG/4288, Aston University, Birmingham, UK, 1994.

The Mixture Density Network: Implementation in TensorFlow Probability

```
import tensorflow as tf
import tensorflow_probability as tfp
tfd = tfp.distributions
```

```
X_ph = tf.placeholder(tf.float32, [None, D])
y_ph = tf.placeholder(tf.float32, [None])
```

```
K = 10
```

```
net = tf.layers.dense(X_ph, 15, activation=tf.nn.relu)
net = tf.layers.dense(net, 15, activation=tf.nn.relu)
locs = tf.layers.dense(net, K, activation=None)
scales = tf.layers.dense(net, K, activation=tf.exp)
logits = tf.layers.dense(net, K, activation=None)
```

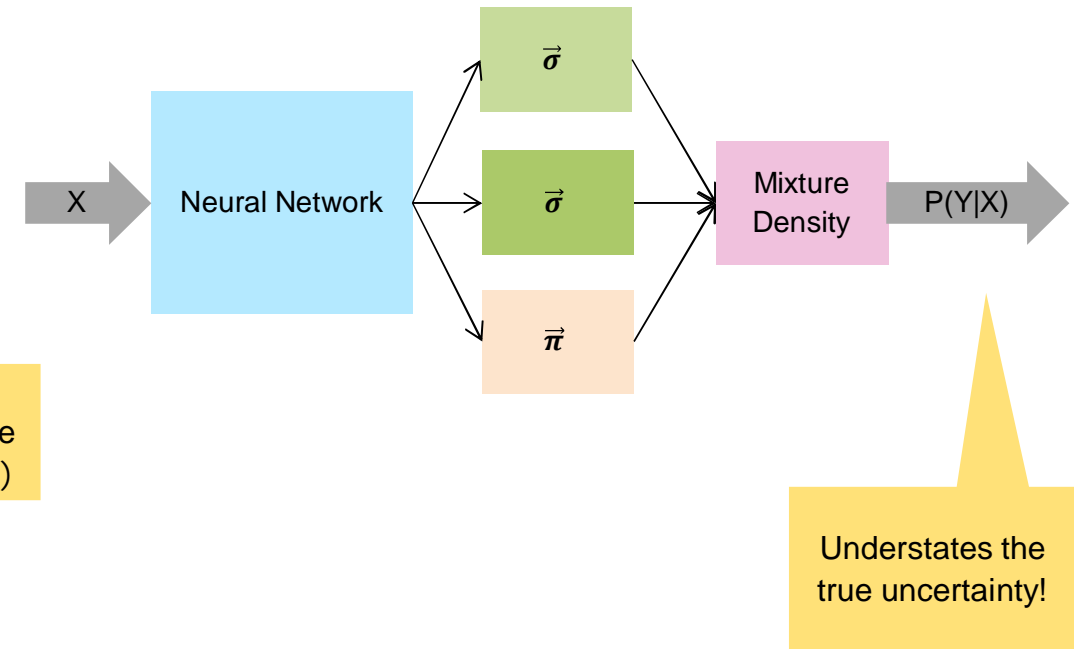
Note the activation! We want $\vec{\sigma} \in (0, \infty)$

```
components = [ tfd.Normal(loc=loc, scale=scale) for loc, scale
                in zip(tf.unstack(tf.transpose(locs)), tf.unstack(tf.transpose(scales)))]
```

```
cat = tfd.Categorical(logits=logits)
```

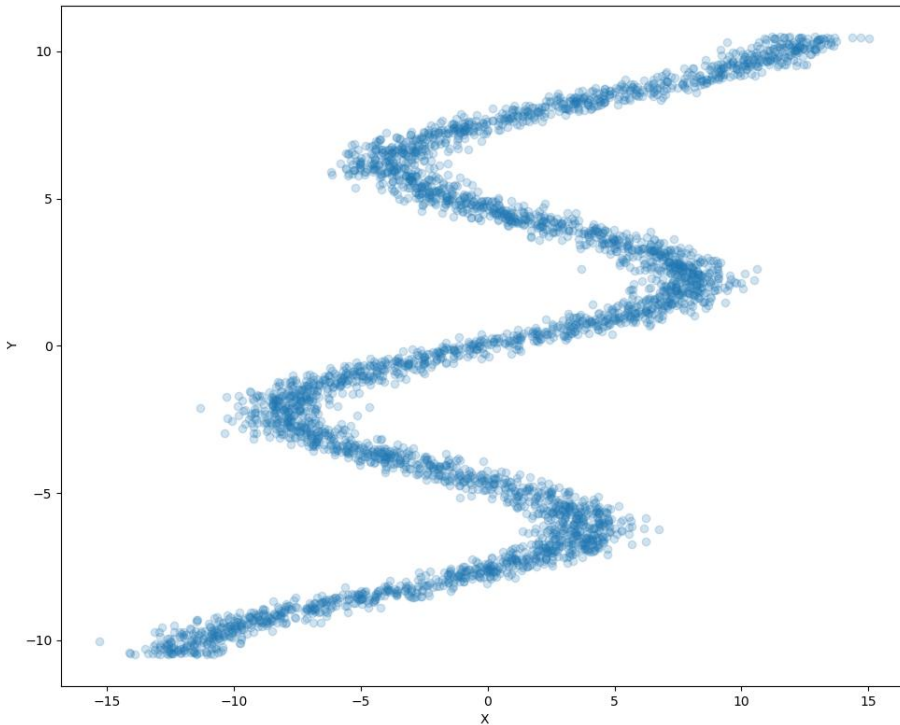
```
mdn = tfd.Mixture(cat=cat, components=components)
```

```
loss = -tf.reduce_mean(mdn.log_prob(y_ph))
```

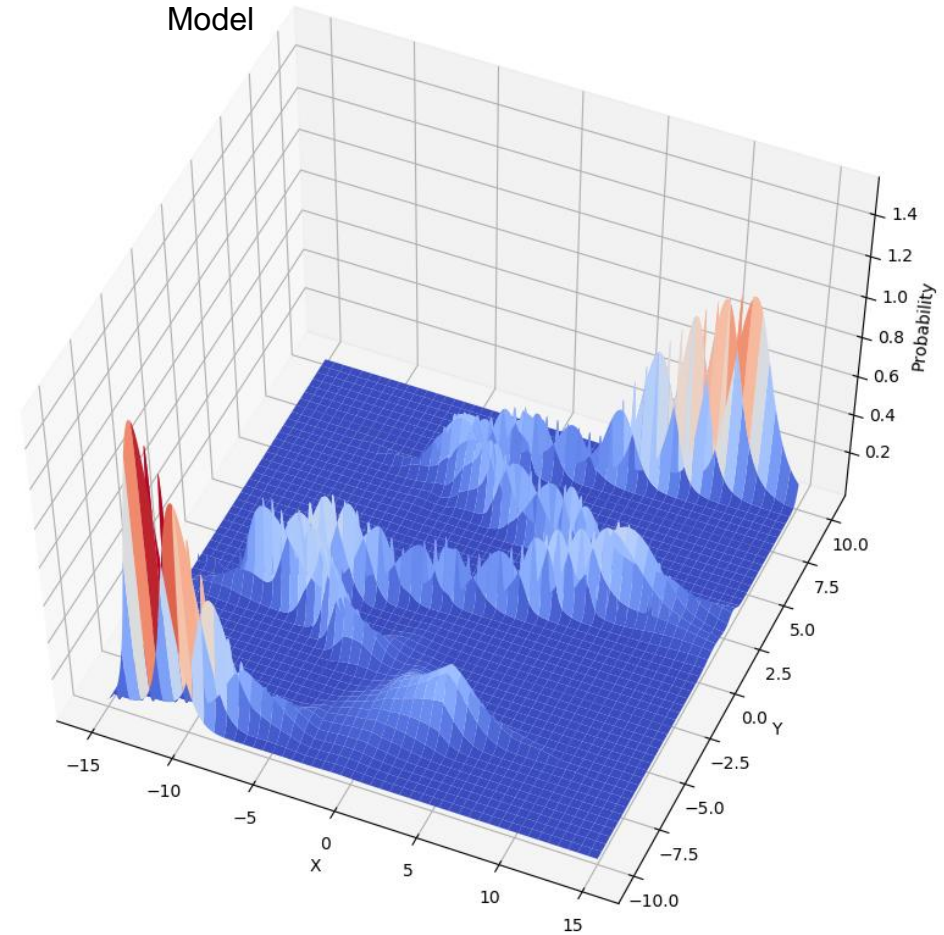


The Mixture Density Network: Result

Training Data



Mixture Density Model



Model Uncertainty

Model Uncertainty

$$y = f(x, \theta) ; D = \{x, t\}$$

Model (Weight) Uncertainty

$$\theta \sim p(\theta | D)$$

Implies an uncertainty in the output!

Is it qualitatively different than simple output uncertainty?

Example: Tossing a Coin

You have tossed a coin 1 000 000 times.

Head	Tails
501 071	498 929

What is your prediction for the outcome of the next flip?

How sure are you? In other words: How sure are you about your model?

Example: Predicting a Time Series

Data

t	0	1	2	3	4	5	6
y	0	1	?	?	?	?	?

What is your prediction for y?

How sure are you? In other words: How sure are you about your model?

Epistemic and Aleatoric Uncertainty

Epistemic

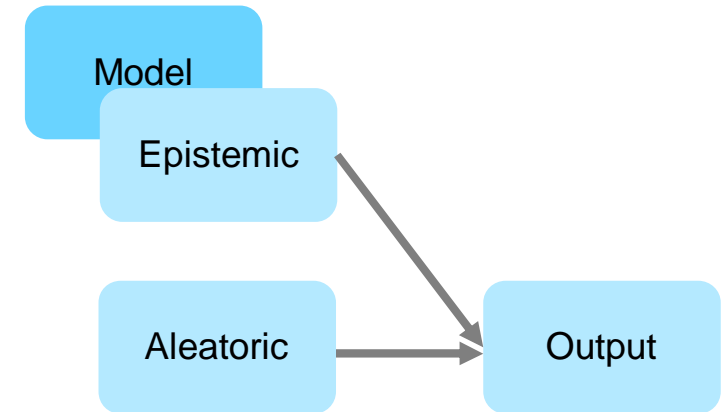
Caused by insufficient number of observations. In other words: The model is underdetermined.

More observations will reduce this type of uncertainty

Aleatoric

Caused by stochasticity or unobservability of an aspect of a system.

More observations will not reduce this type of uncertainty.
Different types of observations might.



Further Reading

Kendall, Alex, and Yarin Gal. "What uncertainties do we need in bayesian deep learning for computer vision?." Advances in neural information processing systems. 2017.

Model/Weight Uncertainty in Neural Networks

Bayes Theorem

$$P(\theta|D) = \frac{\begin{array}{|c|c|} \hline \text{Likelihood} & \text{Prior} \\ \hline \text{of Data} & \\ \hline \end{array} P(D|\theta)P(\theta)}{\begin{array}{|c|} \hline P(D) \\ \hline \end{array} \text{Normalization}}$$

Is this radically different than “normal” Neural Networks?

Preparation: Neural Networks in the Bayesian Framework

Parametric Model with Noise Term

$$t = f(x, \theta) + \varepsilon$$
$$\varepsilon \sim N(0, \sigma)$$
$$D = \{x, t\}$$

Inference on Weights

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

Preparation: Neural Networks in the Bayesian Framework

Likelihood

$$P(y|\theta) = N(y|f(x, \theta) - t, \sigma)$$

$$P(D|\theta) = \prod_{(x,t) \in D} N(y|f(x, \theta) - t, \sigma)$$

Prior

$$P(\theta) = \prod_{\theta} N(\theta|0,1)$$

Log Posterior

$$\mathcal{L}_{\theta}(\theta, D) = C_1 \sum_{(x,t) \in D} (f(x, \theta) - t)^2 + C_2 \sum_{\theta} \theta^2$$

MSE **L2 Weight Regularization**

$$P(\theta|D) = \frac{\overset{\text{Likelihood of Data}}{P(D|\theta)} \overset{\text{Prior}}{P(\theta)}}{\underset{\text{Normalization}}{P(D)}}$$

Preparation: Neural Networks in the Bayesian Framework

$$\mathcal{L}(\theta, D) = C_1 \sum_{(x,t) \in D} (f(x, \theta) - t)^2 + C_2 \sum_{\theta} \theta^2$$

MSE

L2 Weight Regularization

Optimization gives the **Maximum A Posterior (MAP)** Estimate

Further Reading

Barber, David, and Christopher M. Bishop. "Ensemble learning in Bayesian neural networks." NATO ASI SERIES F COMPUTER AND SYSTEMS SCIENCES 168 (1998): 215-238.

Why Bayes is Difficult in Practice

Two Main Challenges:

1. How to represent Distributions?
2. How to calculate $P(D) = \int P(D|\theta)P(\theta)d\theta$?

$$P(\theta|D) = \frac{\begin{array}{|c|c|} \hline \text{Likelihood} & \text{Prior} \\ \hline \text{of Data} & \\ \hline \end{array} \begin{array}{|c|} \hline P(D|\theta)P(\theta) \\ \hline \end{array}}{\begin{array}{|c|} \hline P(D) \\ \hline \end{array}}$$

Normalization

Some Recipes to deal with Bayes Theorem

Reduce it to a point estimate

→ use numerical optimization to find it

p *That's what we just did!*

Empirically estimate $P(\theta|D)$

→ use Markov Chain Monte Carlo Methods to sample $P(\theta|D)$. (MCMC can sample from un-normalized distributions)

Ü *Let's take a closer look!*

Assume conjugate distribution of the exponential family for $P(D|\theta)$ and $P(\theta)$

→ solve analytically

y *Not working for Neural Networks!*

Make simplifying assumptions such as factorization of $P(\theta|D)$

→ use Variational Methods to approximate $P(\theta|D)$

Ü *Let's take a closer look!*

Note: *There are more recipes for special types of models. E.g. Message passing for graphical models.*

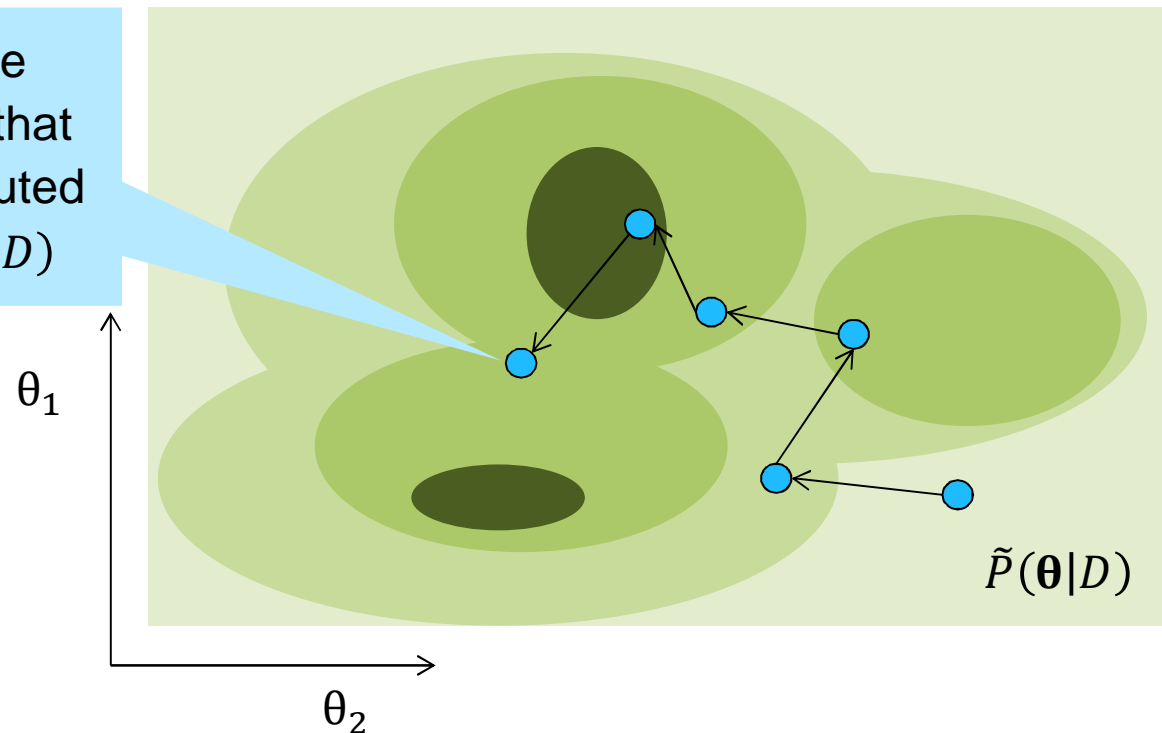
Empirical Estimate of Weights Distribution via Markov Chain Monte Carlo

Problem Statement:

1. We want to sample from the distribution of Neural Network weights given the Data: $P(\theta|D)$.
2. For a given θ we can evaluate the un-normalized $\tilde{P}(\theta|D) = P(D|\theta)P(\theta)$.

Idea:

Jump around in the weight space such that the points are distributed proportional to $\tilde{P}(\theta|D)$



How to Jump: The Metropolis-Hastings Algorithm

Split the Jump into two Phases

1. Propose a new point in weight space

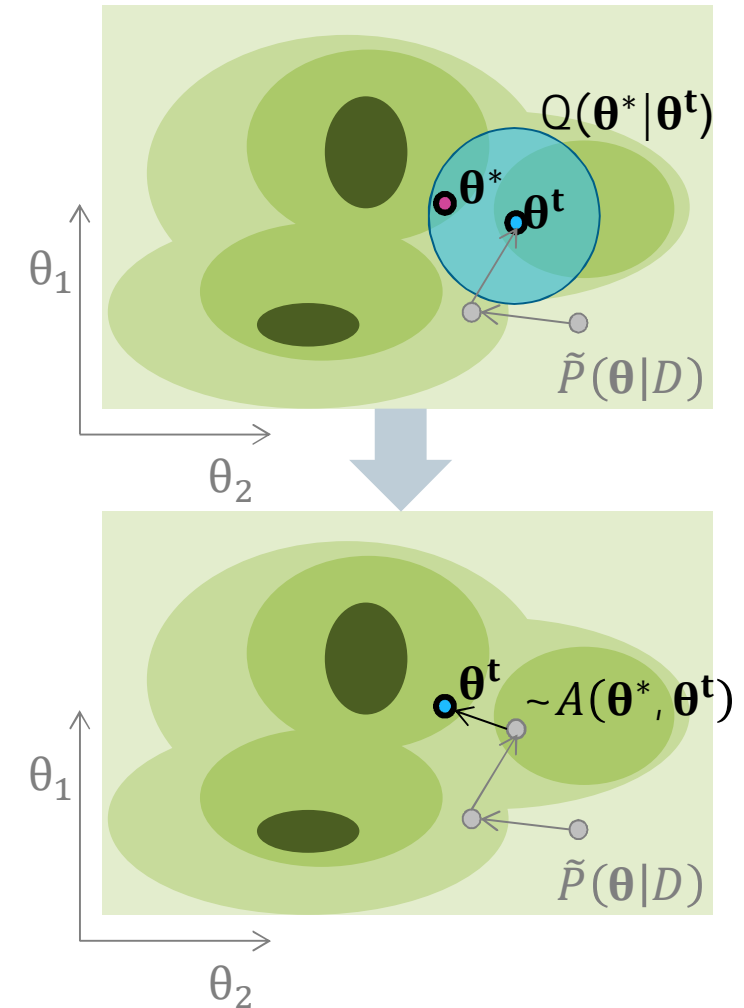
by sampling from a arbitrary but known proposal distribution

- We call the proposal distribution $Q(\theta^* | \theta^t)$
- We call the proposed weight vector θ^*
- We call the current weight vector θ^t

2. Define an Acceptance Probability for the proposed weight vector:

$$\text{accept} \sim A(\theta^*, \theta^t)$$

- If accept: $\theta^{t+1} = \theta^*$
- else: $\theta^{t+1} = \theta^t$



How to Jump: The Metropolis-Hastings Algorithm

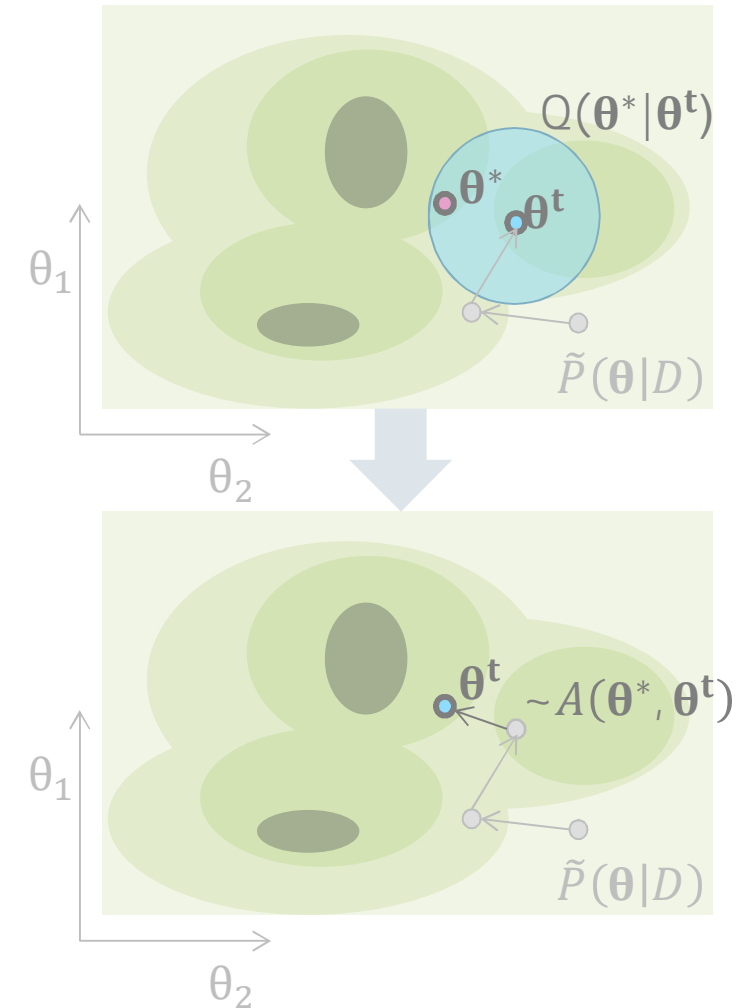
If the Acceptance Probability is chosen as

$$A(\theta^*, \theta^t) = \min\left(1, \frac{\tilde{P}(\theta^*|D)Q(\theta^t|\theta^*)}{\tilde{P}(\theta^t|D)Q(\theta^*|\theta^t)}\right)$$


It can be proven that distribution of the sampled θ^t , θ^{t+1} , θ^{t+2} ... converges to $P(\theta|D)$

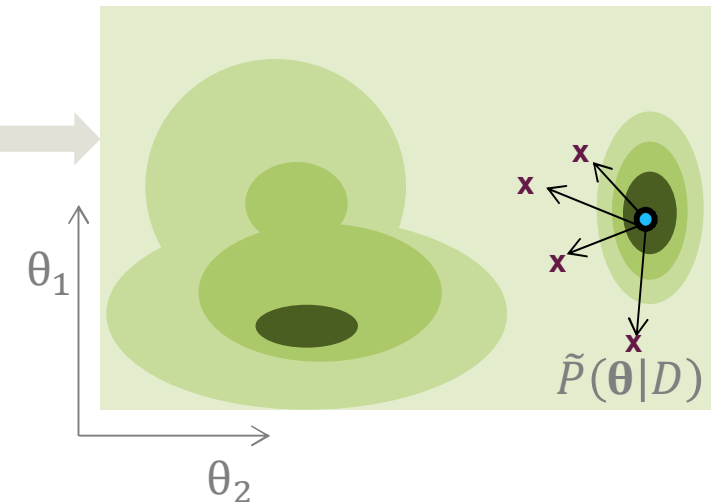
Further Reading

Chapter 11, Sampling Methods of
Christopher M. Bishop. 2006. Pattern
Recognition and Machine Learning
(Information Science and Statistics).
Springer-Verlag, Berlin.



Caveats of Markov Chain Monte Carlo

1. Computational Complexity: One sampling step requires computation of $P(D|\theta)P(\theta)$ which includes D .
2. Convergence to $P(\theta|D)$ can be slow: 
3. Trade-off between high acceptance and efficient traversal of weight-space
 - Long jumps make acceptance low
 - Short jumps makes movement in space slow



There is a large body of research addressing these points: Best use an existing framework or library! 

Python Libraries

- PyMC3: <https://docs.pymc.io/>
- Edward: <http://edwardlib.org/>
- TensorFlow Probability: <https://www.tensorflow.org/probability/>

Big Caveat of Markov Chain Monte Carlo for Deep Learning

What's the memory requirement?

- S samples needed for each weight to build histogram
- W weights in a Deep Neural Network

→ **$O(S*W)$**

For VGG19: 26578886 and 1000 samples for each weight: ~26 Billion!
One weight is float32, so ~100 GB for the sample traces.

Parametric Estimate via the Variational Method

So, $P(\theta|D)$ is intractable? Let's simplify and assume:

$$P(\theta|D) \approx Q_{\omega}(\theta)$$

$$Q_{\omega}(\theta) = \prod_i Q_{\omega_i}(\theta_i)$$

$$\forall_i: \int d\theta_i \cdot Q_{\omega_i}(\theta_i) = 1$$

Let's measure the discrepancy between $P(\theta|D)$ and $Q_{\omega}(\theta)$ with the Kullback–Leibler divergence:

$$D_{KL}(Q_{\omega}(\theta) || P(\theta|D)) = \int d\theta \cdot Q_{\omega}(\theta) \log \frac{Q_{\omega}(\theta)}{P(\theta|D)}$$

If we can measure it, maybe we can minimize it with respect to ω ...

Using D_{KL} is a design choice. There are good reasons to choose differently.

Further Reading

Hernández-Lobato, J. M., et al. "Black-Box α -divergence minimization." 33rd International Conference on Machine Learning, ICML 2016. Vol. 4. 2016

Further Reading

Ranganath, Rajesh, et al. "Operator variational inference." Advances in Neural Information Processing Systems. 2016.

Parametric Estimate via the Variational Method

$$D_{KL}(Q_{\omega}(\boldsymbol{\theta}) || P(\boldsymbol{\theta}|D)) = \int d\boldsymbol{\theta}. Q_{\omega}(\boldsymbol{\theta}) \log \frac{Q_{\omega}(\boldsymbol{\theta})}{P(\boldsymbol{\theta}|D)}$$

Still contains the intractable term $P(\boldsymbol{\theta}|D)$. Let's dissect it!

$$\int d\boldsymbol{\theta}. Q_{\omega}(\boldsymbol{\theta}) \log \frac{Q_{\omega}(\boldsymbol{\theta})}{P(\boldsymbol{\theta}|D)} = \int d\boldsymbol{\theta}. Q_{\omega}(\boldsymbol{\theta}) \log \frac{Q_{\omega}(\boldsymbol{\theta})}{P(\boldsymbol{\theta}, D)} + \log P(D)$$

Rearrange:

$$\log P(D) = D_{KL}(Q_{\omega}(\boldsymbol{\theta}) || P(\boldsymbol{\theta}|D)) + \int d\boldsymbol{\theta}. Q_{\omega}(\boldsymbol{\theta}) \log \frac{P(\boldsymbol{\theta}, D)}{Q_{\omega}(\boldsymbol{\theta})}$$

Const for given D

3 0

Lower bound on log P(D)!

The Evidence Lower Bound

Let's give it a name: **Evidence Lower Bound**

$$\text{ELBO}(D) := \int d\theta. Q_{\omega}(\theta) \log \frac{P(\theta, D)}{Q_{\omega}(\theta)}$$

Option A:

Use factorization and assume exponential family for Q

Coordinate Ascent Variational Inference (CAVI) / Mean Field Approximation

Option B:

Use sampling and stay in the gradient descent framework of neural networks.

Automatic Differentiation Variational Inference (ADVI) / Black Box Variational Inference

(for deep neural networks)

Further Reading

Chapter 10, Approximate Inference of Christopher M. Bishop. 2006. Pattern Recognition and Machine Learning (Information Science and Statistics). Springer-Verlag, Berlin.

The Evidence Lower Bound in Detail

$$\int d\boldsymbol{\theta}. Q_{\omega}(\boldsymbol{\theta}) \log \frac{P(\boldsymbol{\theta}, D)}{Q_{\omega}(\boldsymbol{\theta})} := \text{ELBO}(D)$$

Further deconstruction yields:

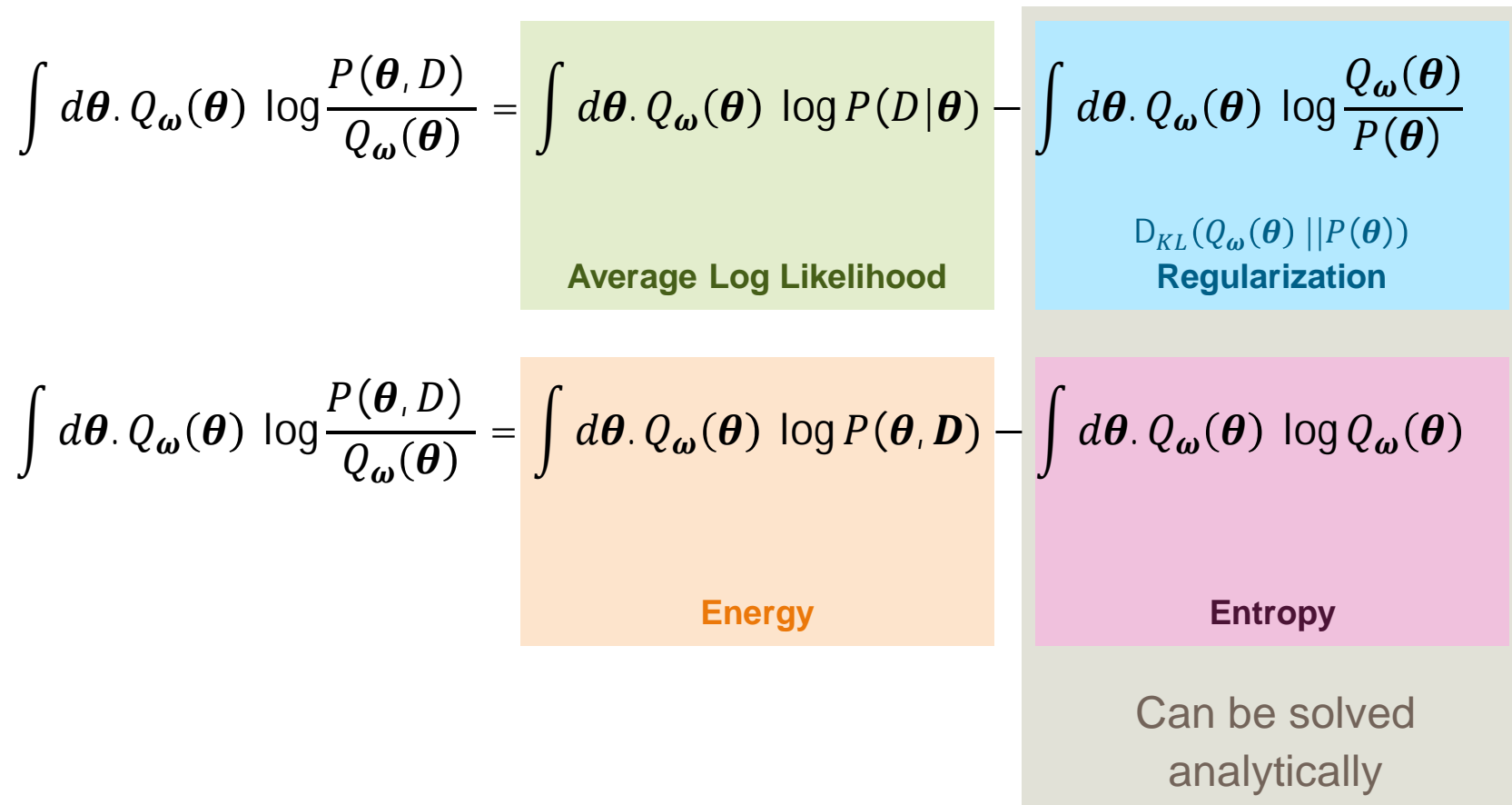
$$P(\boldsymbol{\theta}, D) = P(D|\boldsymbol{\theta})P(\boldsymbol{\theta})$$

$$\int d\boldsymbol{\theta}. Q_{\omega}(\boldsymbol{\theta}) \log \frac{P(\boldsymbol{\theta}, D)}{Q_{\omega}(\boldsymbol{\theta})} = \int d\boldsymbol{\theta}. Q_{\omega}(\boldsymbol{\theta}) \log P(D|\boldsymbol{\theta}) - \int d\boldsymbol{\theta}. Q_{\omega}(\boldsymbol{\theta}) \log \frac{Q_{\omega}(\boldsymbol{\theta})}{P(\boldsymbol{\theta})}$$

Average Log Likelihood

$D_{KL}(Q_{\omega}(\boldsymbol{\theta}) || P(\boldsymbol{\theta}))$
Regularization

Alternative Decomposition



Stochastic Gradient of the ELBO

$$\begin{aligned}\nabla_{\omega} \text{ELBO}(D) &= \nabla_{\omega} \int Q_{\omega}(\theta) \log P(D|\theta) - Q_{\omega}(\theta) \log \frac{Q_{\omega}(\theta)}{P(\theta)} d\theta \\ &= \nabla_{\omega} \mathbb{E}_{Q_{\omega}(\theta)} \left[\log P(d|\theta) - \log \frac{Q_{\omega}(\theta)}{P(\theta)} \right]\end{aligned}$$

If we could move the derivative into the Expectation, we could approximate it by sampling!

There are at least two options:

1. Log Derivative Trick

ŷ suffers from high variance

⊢ flexible with respect to the distribution

2. Reparametrization (a.k.a. Perturbation Analysis, Pathwise Derivatives)

⊢ easy to implement, ⊢ low gradient variance

ŷ does not work for every distribution

Viable and flexible option in combination with variance control techniques.
Also used in Reinforcement Learning

Further Reading

Williams, Ronald J. "Simple statistical gradient-following algorithms for connectionist reinforcement learning." Machine learning 8.3-4 (1992): 229-256.

Further Reading

Ranganath, Rajesh, Sean Gerrish, and David Blei. "Black box variational inference." Artificial Intelligence and Statistics. 2014.

Reparametrization Trick

$$\nabla_{\omega} \mathbb{E}_{p(z|\omega)}[f(z)] = \nabla_{\omega} \int p(z|\omega) f(z) dz$$

Find a way to reparametrize $p(z|\omega)$ such that $z \sim g(\varepsilon, \omega)$ and we just sample ε

$$\nabla_{\omega} \int p(z|\omega) f(z) dz = \nabla_{\omega} \int p(\varepsilon) f(g(\varepsilon, \omega)) d\varepsilon$$

We move the derivative inside the integral:

warranted by the
dominated
convergence
theorem

$$= \int p(\varepsilon) \nabla_{\omega} f(g(\varepsilon, \omega)) d\varepsilon = \mathbb{E}_{p(\varepsilon)}[\nabla_{\omega} f(g(\varepsilon, \omega))]$$

Example:

$$z \sim \mathcal{N}(\mu, \sigma) \rightarrow z \sim \mu + \sigma \varepsilon; \varepsilon \sim \mathcal{N}(0, 1)$$

ε is now independent of μ, σ and backpropagation works!

Further Reading

Kingma, Diederik P., and Max Welling. "Auto-encoding variational bayes." arXiv preprint arXiv:1312.6114 (2013).

Flipout

Challenge:

“because a network typically has many more weights than units, it is very expensive to compute and store separate weight perturbations for every example in a mini-batch. Therefore, stochastic weight methods are typically done with a single sample per mini-batch.”



Further Reading

Wen, Yeming, et al. "Flipout: Efficient Pseudo-Independent Weight Perturbations on Mini-Batches." arXiv preprint arXiv:1803.04386 (2018).

For VGG19: 26578886 weights and a batch size of 100 : ~2.6 Billion!

- **Same sample for the entire batch**
- **high variance in the gradient**
- **small learning rate needed**
- **slow convergence**

Idea:

Add independence to the weight samples in one batch with low computational cost.

Conditions:

1. weight distribution is symmetric
2. each weight is sampled independently

These conditions are met for our stochastic variable (perturbation) ε !

- Let ΔW be a matrix of independently sampled $\varepsilon \sim \mathcal{N}(0,1)$
- Let E be a matrix of independently sampled signs ± 1

$$\widehat{\Delta W} = \Delta W \circ E$$

$\widehat{\Delta W}$ and ΔW are equally distributed, but $\widehat{\Delta W}$ is less correlated!

$$\widehat{\Delta W} = \Delta W \circ E$$

***E* is as big as a fully sampled weights matrix for the batch. Nothing gained yet!**

Can we generate *E* cheaply?

Idea:

Use a rank one *E* induced by a product of two random sign vectors *r* and *s* :

$$E = r s^T$$

For a full analysis of the gradient variance reduction and computational cost see: 

Further Reading

Wen, Yeming, et al. "Flipout: Efficient Pseudo-Independent Weight Perturbations on Mini-Batches." arXiv preprint arXiv:1803.04386 (2018).

Implementation in TensorFlow Probability

The Bayes Version of the LeNet CNN Architecture:

Taken from the TensorFlow Probability Examples

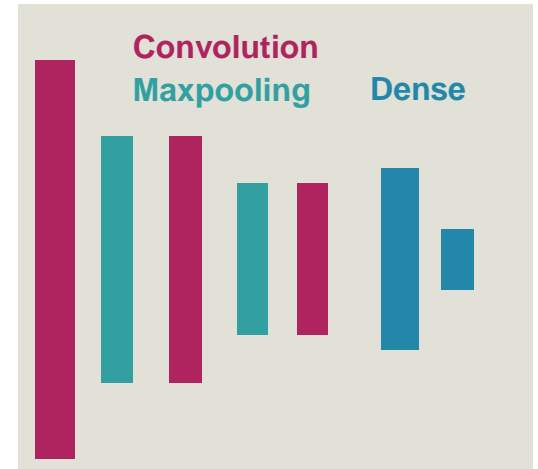
```
neural_net = tf.keras.Sequential([  
    tf.layers.Convolution2DFlipout(6, kernel_size=5, padding="SAME", activation=tf.nn.relu),  
    tf.keras.layers.MaxPooling2D(pool_size=[2, 2], strides=[2, 2], padding="SAME"),  
    tf.layers.Convolution2DFlipout(16, kernel_size=5, padding="SAME", activation=tf.nn.relu),  
    tf.keras.layers.MaxPooling2D(pool_size=[2, 2], strides=[2, 2], padding="SAME"),  
    tf.layers.Convolution2DFlipout(120, kernel_size=5, padding="SAME", activation=tf.nn.relu),  
    tf.keras.layers.Flatten(),  
    tf.layers.DenseFlipout(84, activation=tf.nn.relu),  
    tf.layers.DenseFlipout(10)  
])
```

```
labels_distribution = tfd.Categorical( logits = neural_net(images) )
```

```
neg_log_likelihood = -tf.reduce_mean(labels_distribution.log_prob(labels))
```

```
kl = sum(neural_net.losses) / mnist_data.train.num_examples
```

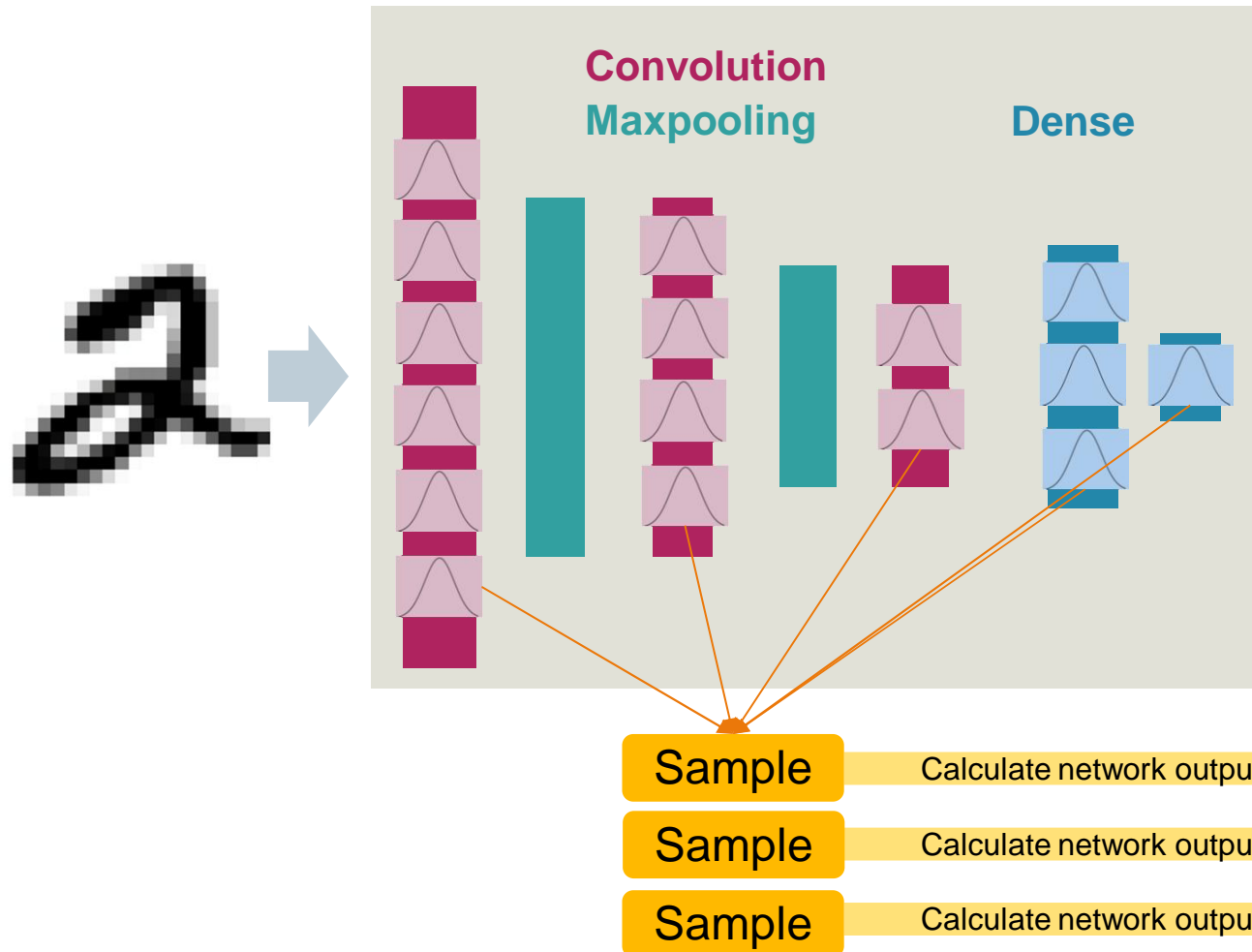
```
elbo_loss = neg_log_likelihood + kl
```

$$\approx - \int d\theta. Q_{\omega}(\theta) \log P(D|\theta) + \int d\theta. Q_{\omega}(\theta) \log \frac{Q_{\omega}(\theta)}{P(\theta)}$$


Further Reading on LeNet
LeCun, Yann, et al. "Gradient-based learning applied to document recognition." Proceedings of the IEEE 86.11 (1998): 2278-2324.

Further Practice
https://github.com/tensorflow/probability/tree/master/tensorflow_probability/examples

Sampling from a Bayesian Neural Network



Pros

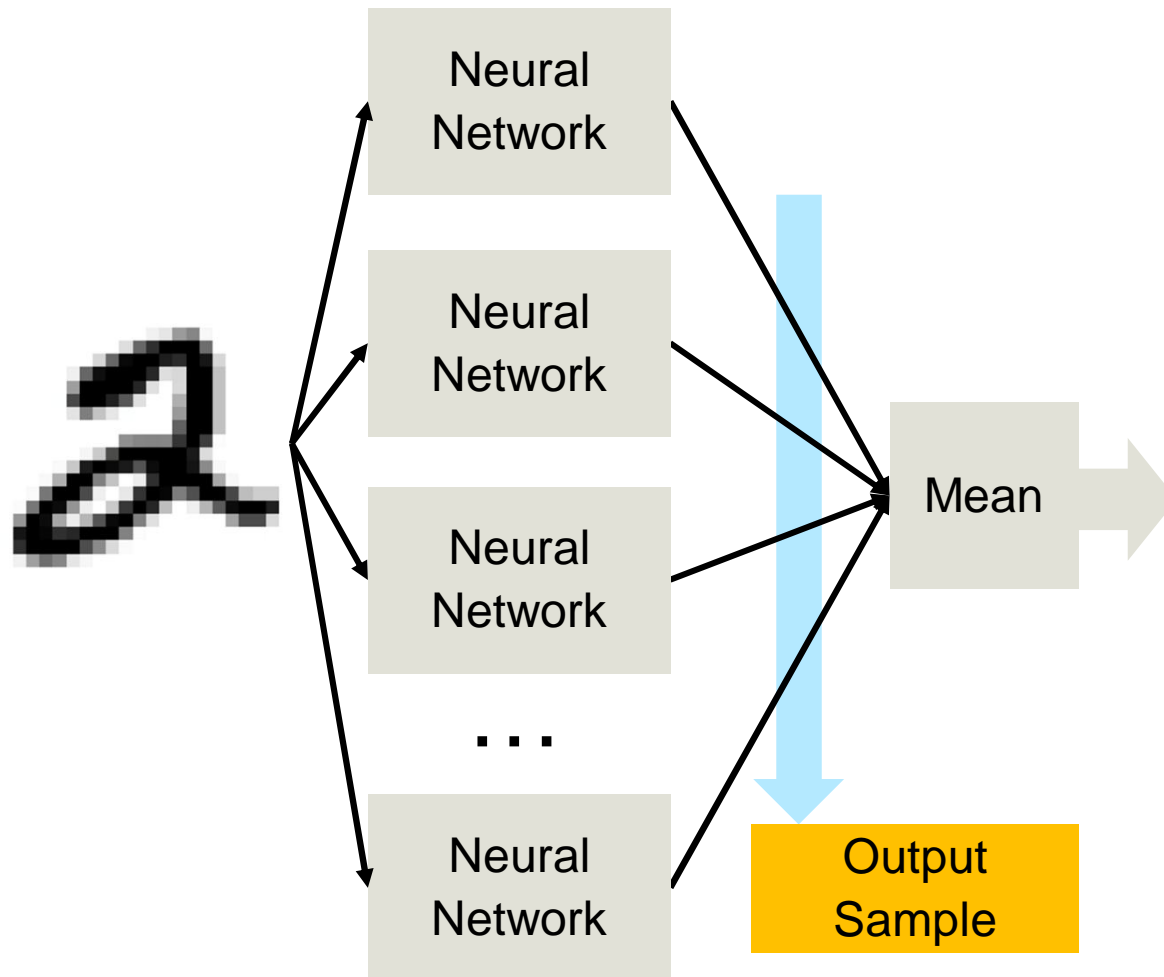
- **Estimate of Model Uncertainty**
(If your use-case benefits from it!)
- Implies Output Uncertainty
- Acts as Regularization

Con

- Computational / Memory Overhead

0	1	2	3	4	...	9
0	0	.98	0	0		.02
0	.01	.97	0	0		.01
0	0	.99	0	0		0

Approximating Bayesian Networks with Ensembles



Pros

- Simple
- Parallelizable
- Reduction of expected error

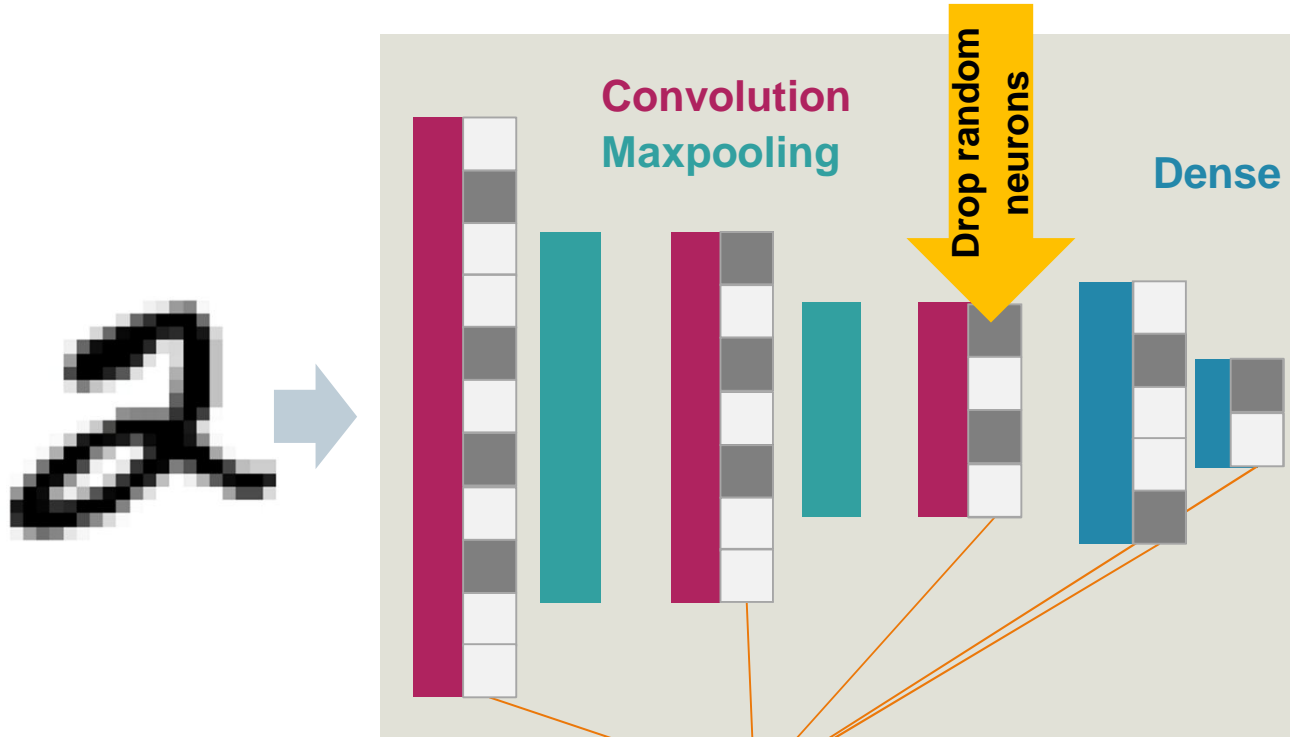
Con

- Excessive training time for modern Deep Neural Networks

Further Reading

Perrone, Michael P., and Leon N. Cooper. When networks disagree: Ensemble methods for hybrid neural networks. No. TR-61. BROWN UNIV PROVIDENCE RI INST FOR BRAIN AND NEURAL SYSTEMS, 1992.

Approximating Ensembles with Dropout



Pros

- Implies Output Uncertainty
- Acts as Regularization
- Computationally highly efficient

Con

- Only indirect estimate of Model Uncertainty

Further Reading

Gal, Yarin, and Zoubin Ghahramani. "Dropout as a Bayesian approximation: Insights and applications." Deep Learning Workshop, ICML. Vol. 1. 2015.

Sample

Calculate network output

Sample

Calculate network output

Sample

Calculate network output

0	1	2	3	4	...	9
0	0	.98	0	0		.02
0	.01	.97	0	0		.01
0	0	.99	0	0		0

In case you want to delve deeper into Bayesian Methods

