# Machine Intelligence:: Deep Learning
# Week 8

*Beate Sick, Jonas Brändli, Oliver Dürr*

Bayesian approaches for improving the performance and uncertainty estimates of NN models by taking into account the epistemic uncertainty.
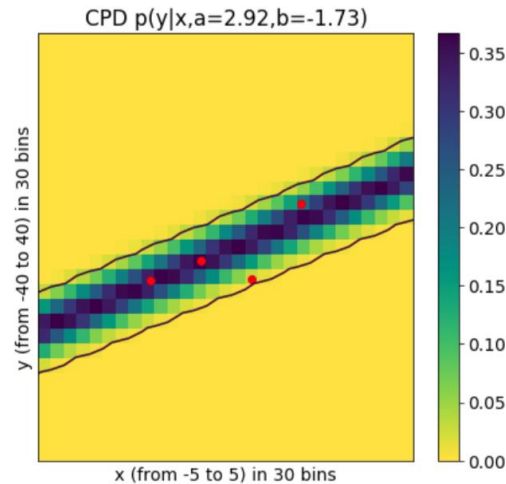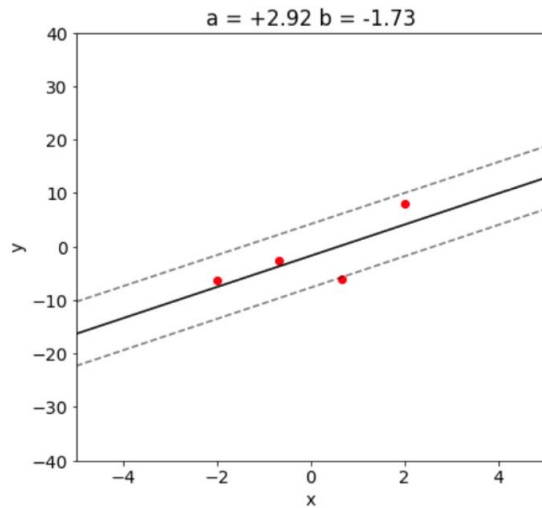
# Outline:

- Bayesian Modelling
    - Introduction to Bayesian Statistics
    - A simple example (Bayes the Hacker's way)
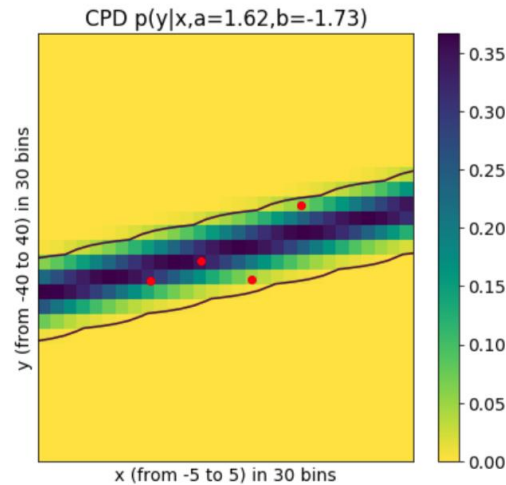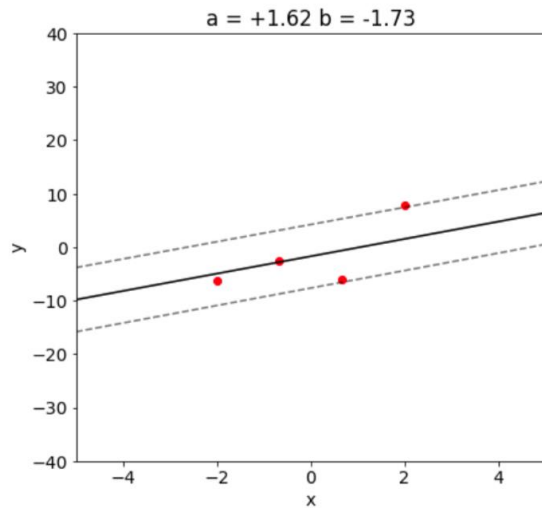    - Variational Approximation

# Bayesian Deep Learning

→ Improving prediction performance

→ Providing epistemic uncertainty measures

# Bayes the hackers' way

Let's look at good old linear regression to understand the gist of the Bayes idea. Assume $\sigma = 3$ to be known.
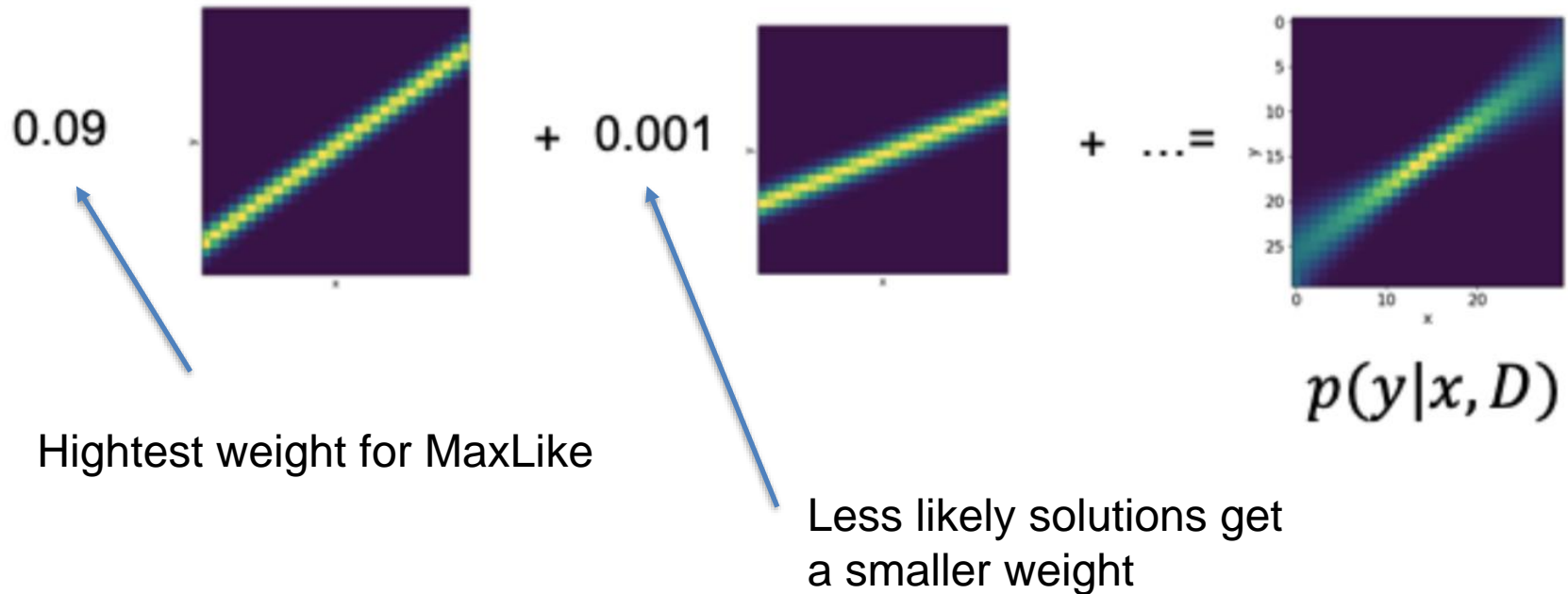


MaxLike Solution

A bit off the MaxLike Solution

# Combining different fits

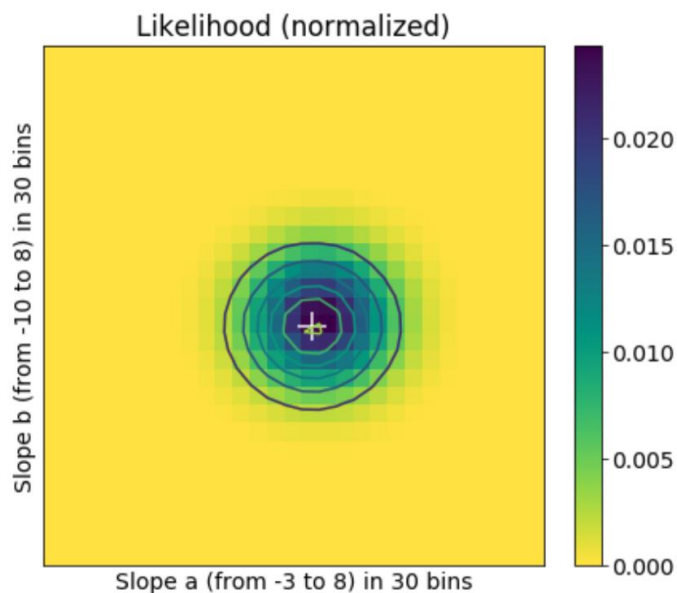Also take the other fits with different parameters into account and weight them



$0.09$

$+ \ 0.001$

$+ \ ... =$

$p(y|x, D)$

Hightest weight for MaxLike

Less likely solutions get
a smaller weight

Question: How to get the weight?

Idea: use the (normalized) likelihood $p_{\mathrm{norm}}(D|(a, b))$ !

# Don't put all egg's in one Basket

- Also take other solutions for parameters $a$, $b$ into account



$$p_{\text{norm}}(D|(a,b)) = \frac{p(D|(a,b))}{\sum_w p(D|(a,b))}$$
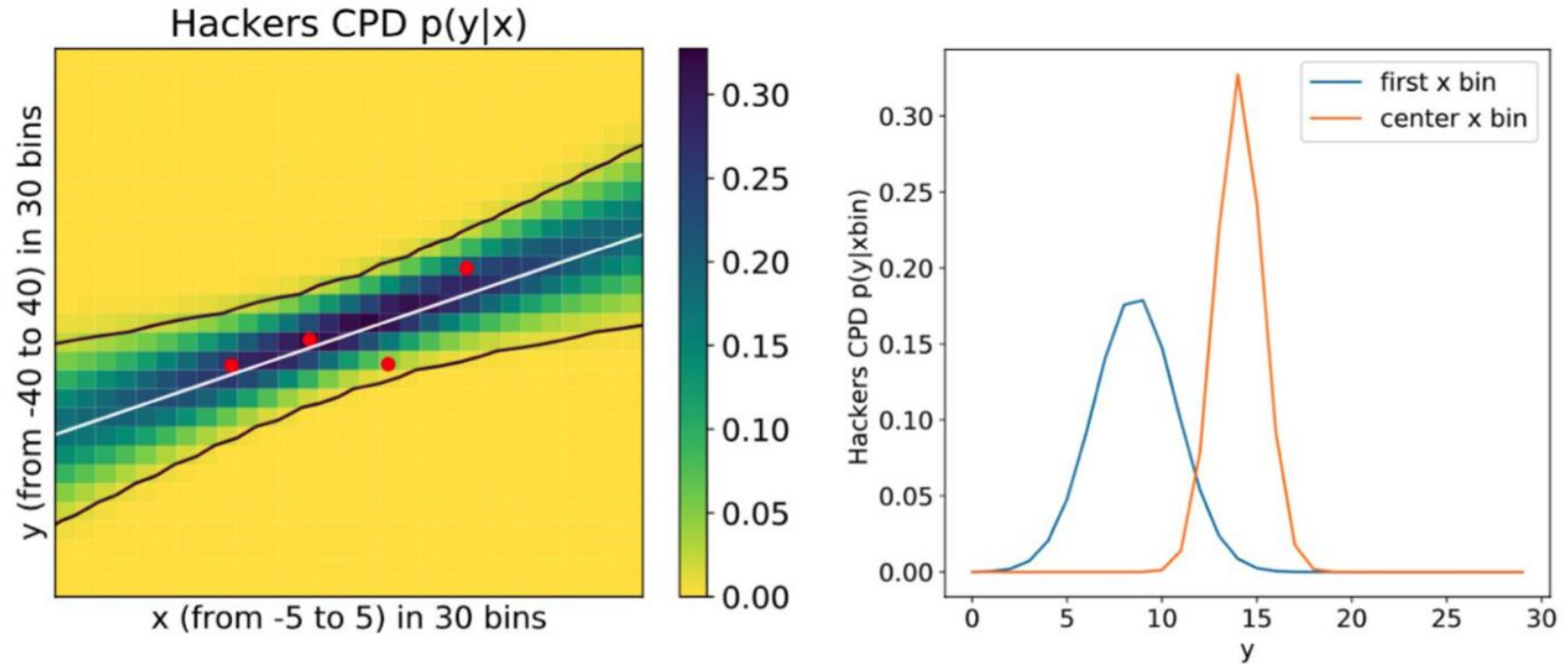
$$\text{p}(y|x,D) = \sum_a \sum_b p(y|x,(a,b)) \cdot p_{\text{norm}}(D|(a,b))$$

Likelihood at 30x30 different positions of a and b. Normalized to be one.
https://github.com/tensorchiefs/ dl_book/blob/master/chapter_07 /nb_ch07_02.ipynb

# Result

$$p(y|x, D) = \sum_a \sum_b p(y|x, (a, b)) \cdot p_{\text{norm}}(D|(a, b))$$



**Figure 7.6 The predictive distribution for the Bayesian linear regression model, trained with the four data points shown on the left side by the color-code and on the right as conditional distribution at two different *x* positions. You can clearly see that the uncertainty gets larger when leaving the *x*-regions where there's data.**

# Bayesian statistics

- The Bayesian Theorem

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

- Applied to $A = \theta$ and $B = D$

  – Parameters $\theta$ of a model e.g. weights *w* of NN

  – Training data *D*

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

# Revisiting the Hackers' example

- Hacker's way

  - $p(y|x,D) = \sum p(y|x,w) \cdot p_{\text{norm}}(D|w)$

    - $p_{\text{norm}}(D|w) = \dfrac{p(D|w)}{\sum_w p(D|w)}$

- Bayes

  - $p(y|x,D) = \sum_w p(y|x,w) \cdot p(w|D)$

  - $p(w|D) = \dfrac{p(D|w)\boldsymbol{p(w)}}{p(D)} = \dfrac{p(D|w)p(w)}{\sum_w p(D|w)p(w)}$

→ The Hacker's way is Bayes, if we assume a uniform prior $p(w) = \text{const}$

# Bayesian modeling has less problems with complex models

**Frequentist's strategy:**
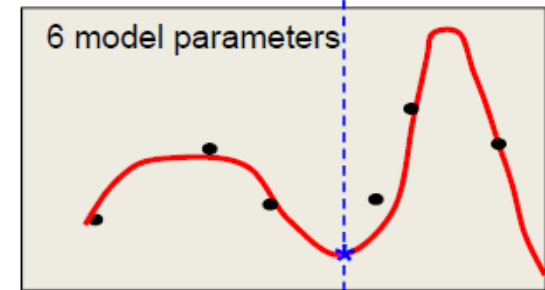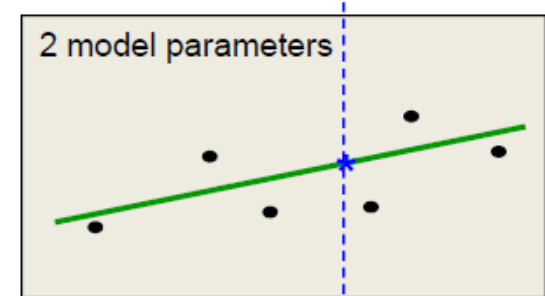You can only use a complex model if you have enough data!

**Bayesian's strategy:**
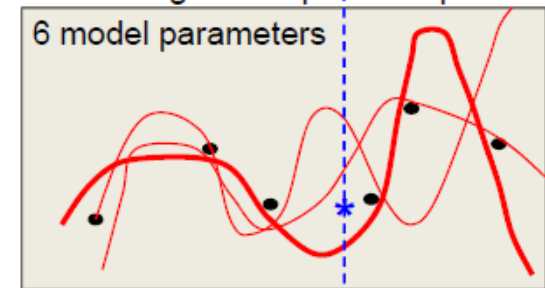Use the model complexity you believe in.

Do not just use the best fitting model.

Do use the full posterior distribution over parameter settings leading to vague predictions since many different parameter settings have significant posterior probability.

$$\mathrm{p}(y|x^*, \text{traindata}) = \int p(y|x^*, \theta) \cdot p_{\mathrm{norm}}(\theta|\text{traindata}) \; d\theta$$

Image credits: Hinton coursera course



2 model parameters



6 model parameters

Models with significant posterior probability



6 model parameters

x*: new input

# Bayesian terminology

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

- $p(\theta|D)$     posterior
- $p(D|\theta)$     likelihood
- $p(\theta)$     prior
- $p(D)$     evidence     just a normalization constant
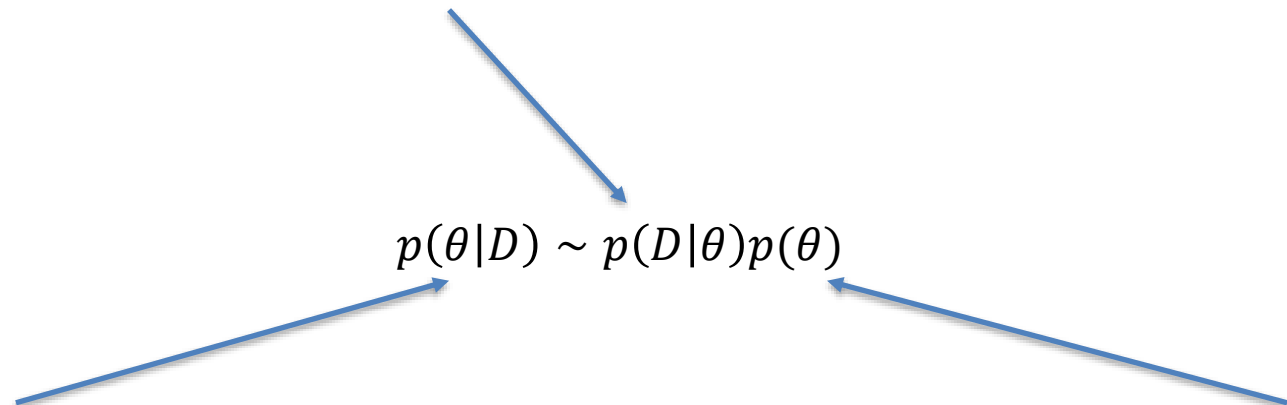
# The Bayesian Mantra (say it loud)

$$p(\boldsymbol{\theta}|\boldsymbol{D}) \sim p(\boldsymbol{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})$$

"The posterior is proportional to the likelihood, times the prior"

# Interpretation updating the degree of belief

- Parameters $\theta$ are random variables, following a distribution
  - This distribution reflects our belief about which values are probable for $\theta$

- The Bayes formula is seen as an update of our belief in the light of data

likelihood updates degree of belief

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

posterior (after seeing data)

prior belief (before seeing data)

# Example Coin Toss

Head or tail?

$Y_i \in \{0, 1\}$, here 0 stands for tail and 1 for head
$Y \sim \mathrm{Ber}(\theta)$, with 1 parameter $\theta = P(Y = 1)$

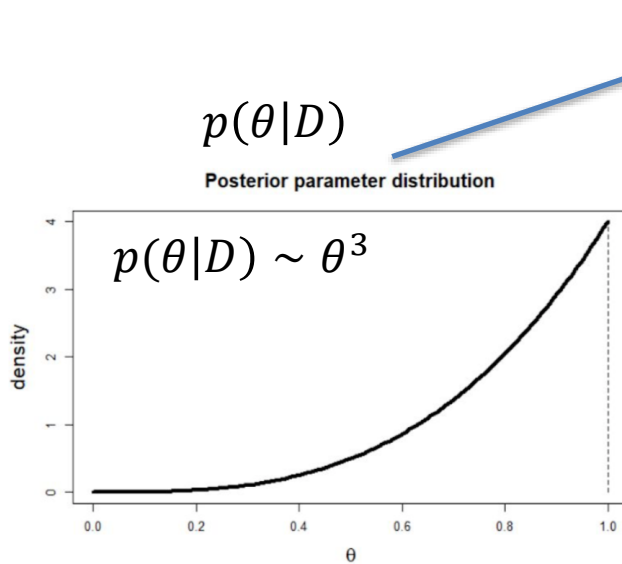For a fair coin $\theta = 0.5$
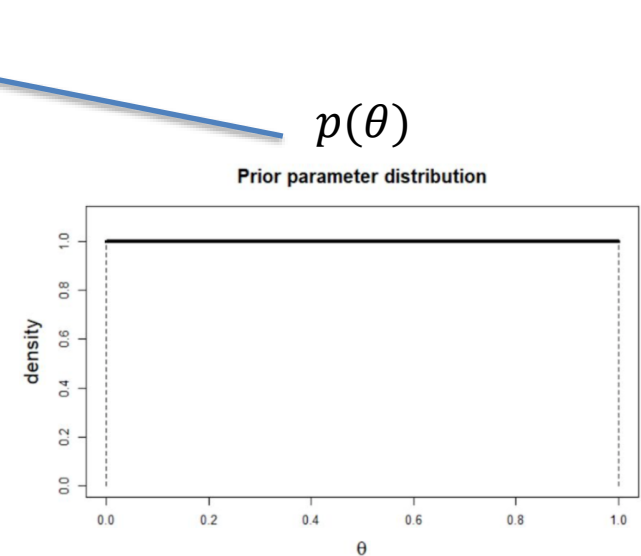


prior predictive distribution

But how to know if a coin is fair?

# Analyzing a Coin Toss Experiment

- We do an experiment and observe 3 times head → D='3 heads'
- $\theta$ parameter for the Bernoulli-distribution (probability of head)
- Before the experiment we assume all value of $\theta$ are equally likely $p(\theta) = \text{const}$
- Calculate likelihood $p(D|\theta) = p(y = 1) \cdot p(y = 1) \cdot p(y = 1) = \theta \cdot \theta \cdot \theta = \theta^3$
- Posterior $p(\theta|D) \sim p(D|\theta)p(\theta) = \theta^3$ beliefs more in head

**Posterior (after)**                                        **Prior (before)**

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

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

$\theta^3$

$p(\theta|D) \sim \theta^3$



Posterior parameter distribution



Prior parameter distribution

$p(\theta|D) = 4 \cdot \theta^3$ (the factor 4 is needed for normalization so that the posterior integrates to 1)

# Posterior Predictive Distribution

Posterior predictive distribution

$$\mathrm{p}(y|x, D) = \int_\theta p(y|x, \theta) \cdot p(\theta|D) \ d\theta$$

The coin example is unconditional (there is no predictor x)

$$\mathrm{p}(y|D) = \int_\theta p(y|\theta) \cdot p(\theta|D) \ d\theta$$

To compute the posterior predictive probability for head in the coin example, we need:

- $\mathrm{p}(y = 1|D) = \theta$
- $p(\theta|D) = 4 \cdot \theta^3$

Posterior predictive distribution

$$\mathrm{p}(y = 1|D) = \int_\theta \theta \ \cdot 4 \cdot \theta^3 d\theta = 0.8$$

$$P(Y = 1|D) = \int_0^1 \theta \cdot 4 \cdot \theta^3 d\theta = \frac{4}{5} \cdot \theta^5 \ \Big|_0^1$$

$$P(Y = 1|D) = \frac{4}{5} \cdot 1^5 - \frac{4}{5} \cdot 0^5 = 0.8$$
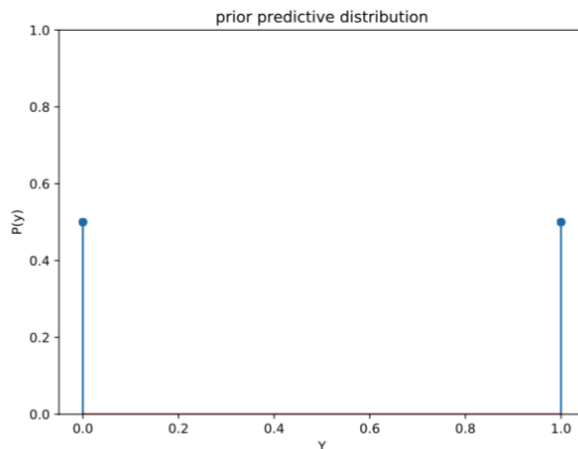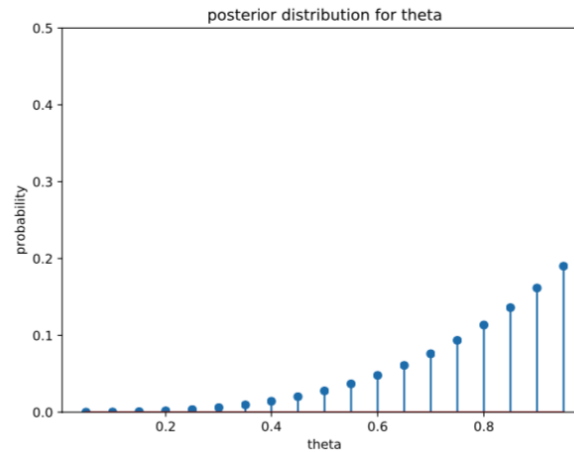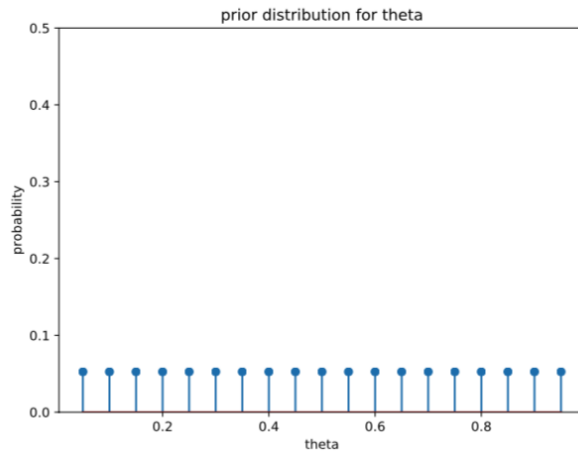
$$P(Y = 0) = 1 - P(Y = 1) = 0.2$$

# Prior and Posterior Predictive distribution

# Coin example «the hacker's way»

- Work through the notebook
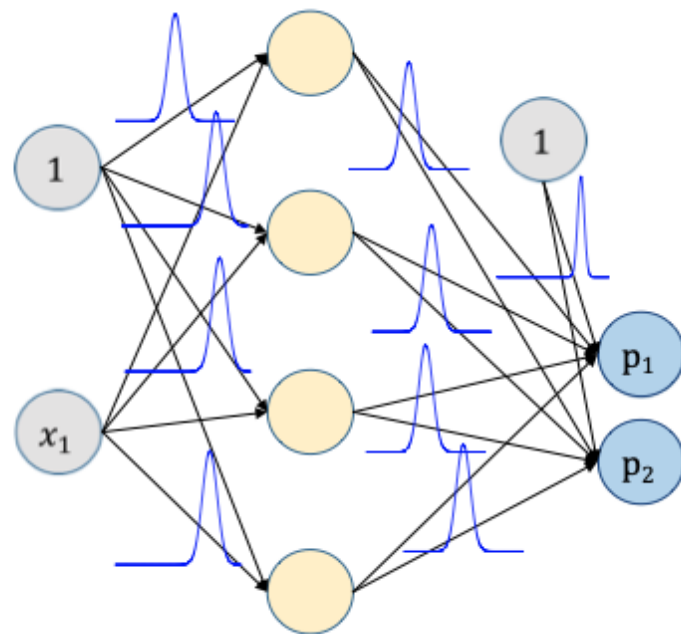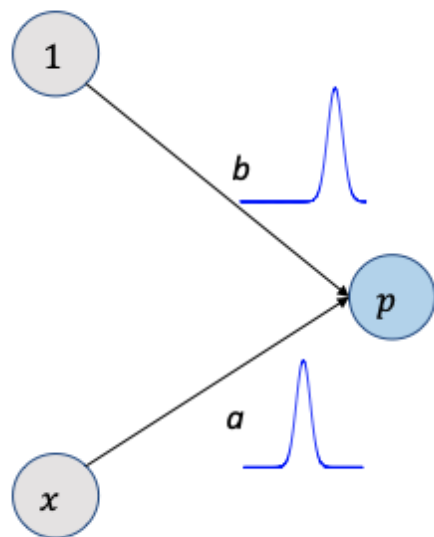- **https://github.com/tensorchiefs/dl_book/blob/master/chapter_07/nb_ch07_03.ipynb**

# Bayesian Neural Networks

# Bayesian Neural* Networks (BNN)

- Linear Regression with Gaussian Prior and fixed Sigma can be solved analytically



- Bayesian Neural Network cannot be solved analytically

*Don't confuse Bayesian Networks (DAGs) with Bayesian Neural Networks

# Approximations to BNN

- A BNN would require to calculate

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

- Usually no analytical solution exists (only for simple problems)

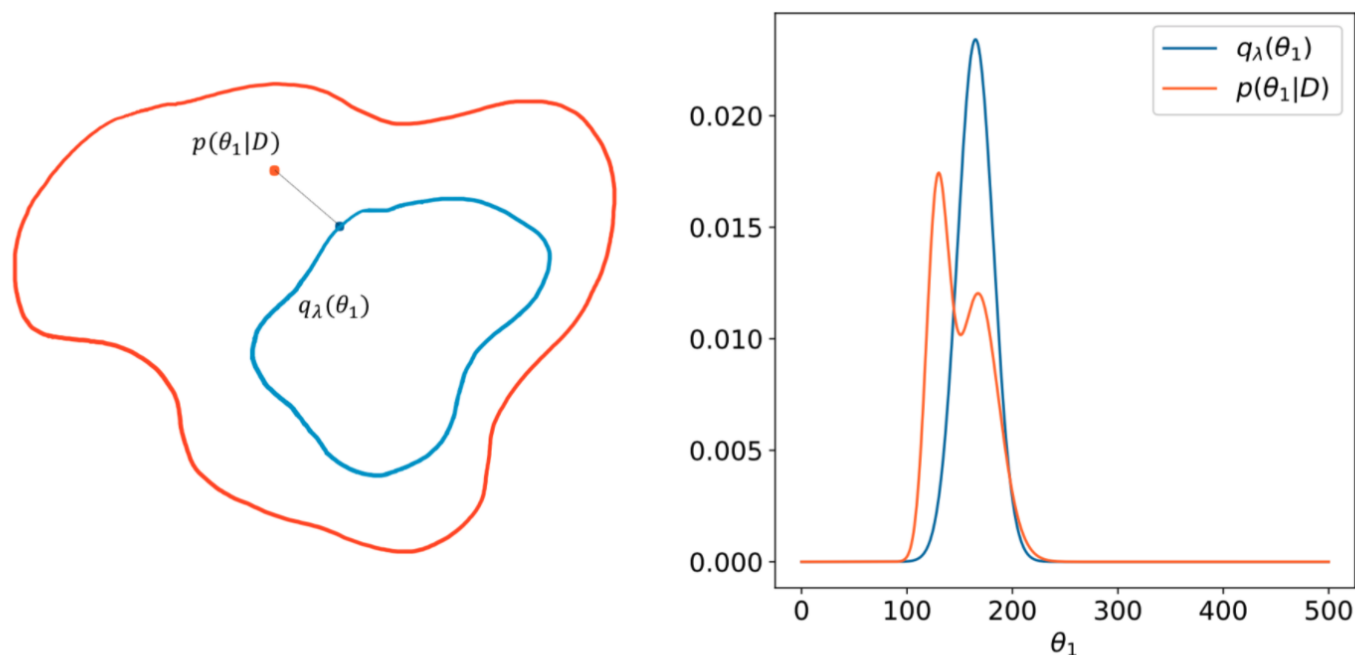- Computing $\sum_\theta p(D|\theta)p(\theta)$ is impossible for high-dimension $\theta$

**Approximations**
- MCMC (only for very small NN feasible)
    - Sample from $p(\theta|D)$ with knowledge of $p(D|\theta')p(\theta')$ / $p(D|\theta'')p(\theta'')$
- Gaussian variational Inference VI
    - Approximate $p(w|D)$ by a Gaussian $N(\mu, \sigma)$ and tune $\mu, \sigma$
- MC-Dropout
    - MC-Dropout during predictions (magically) samples from a variational approximation

# Variational Inference

# The principle of VI

- Replace Posterior $p(\theta|D)$ with variation distribution $q_\lambda(\theta)$
- Typically independent Gaussian for each weight $\lambda = (\mu, \sigma)$
  - $p(\theta|D) = q_{\mu,\sigma}(\theta)$



**Figure 8.3 The principle idea of variational inference (VI). The larger region on the left depicts the space of all possible distributions, and the dot in the upper left represents the posterior** $p(\theta_1|D)$ **(corresponding to the dotted density on the right panel). The inner region depicts the space of possible variational distributions** $q_\lambda(\theta_1)$**. The optimized variational distribution** $q_\lambda(\theta_1)$ **(illustrated by the point in the inner loop in the left panel, corresponding to the solid density on the right panel) has the smallest distance to the posterior (shown by the dotted line on the right).**

# Particularies

- Layers for VI:
  - `DenseReparameterization`
  - `Convolution{1D,2D,3D}Reparameterization`
  - `Further a method called Flipout to speed up training`

From documentation (Convolution2DFlipout)

When doing minibatch stochastic optimization, <u>make sure to scale this loss such that it is applied just once per epoch</u> (e.g. if kl is the sum of losses for each element of the batch, you should pass kl / num_examples_per_epoch to your optimizer)
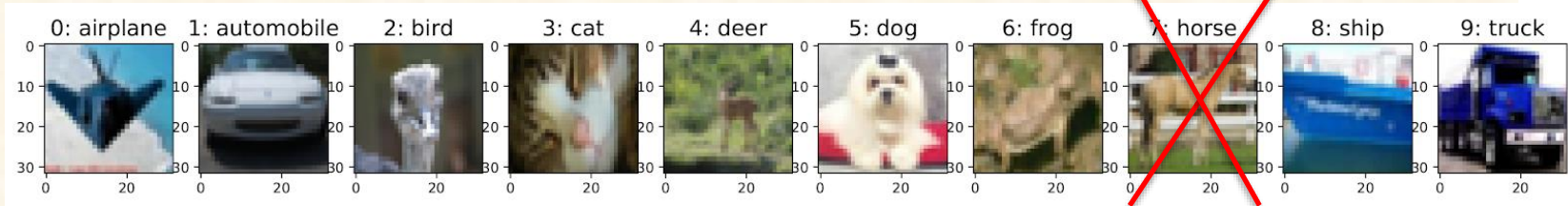
```
kl = tfp.distributions.kl_divergence
divergence_fn=lambda q, p, _: kl(q, p) / (num * 1.0)

DenseReparameterization(1,kernel_divergence_fn=divergence_fn)
```

# Hands-on Time cntd.: Fit the VI Bayesian NN



Train a CNN with only 9 of the 10 classes and investigate if the uncertainties are different when predicting images from known or unknown classes.

https://github.com/tensorchiefs/dl_course_2021/blob/master/notebooks/20_cifar10_classification_mc_and_vi.ipynb
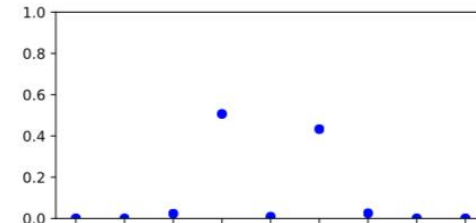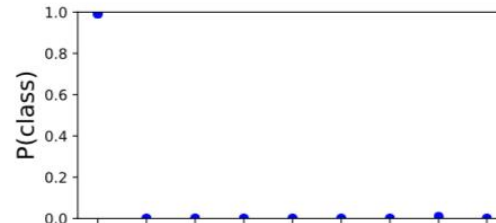
# Comparing non-Bayesian with Bayesian NN
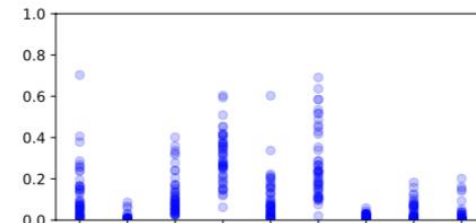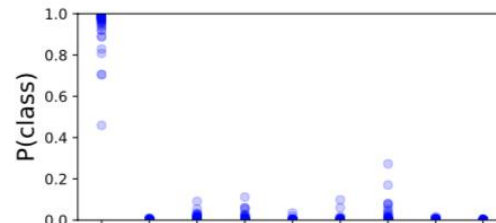
# Looking at the predictive distribution!

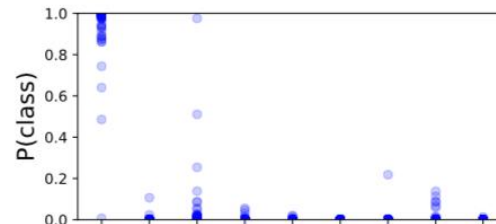# Non-Bayesian and Bayesian NNs



Non-Bayesian NN

MC dropout Bayesian NN

VI Bayesian NN
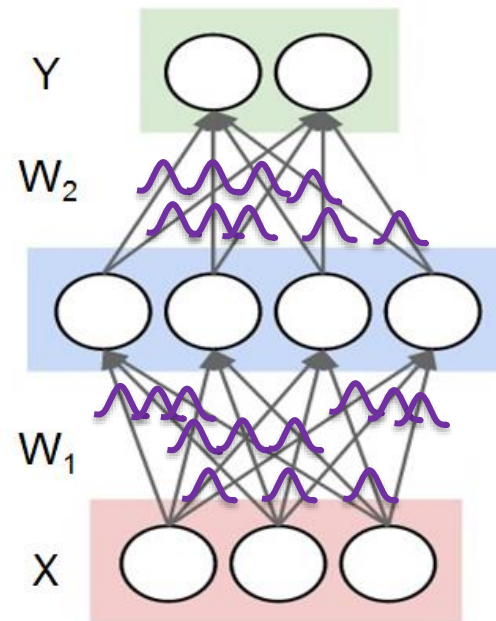
**Weights are fixed**
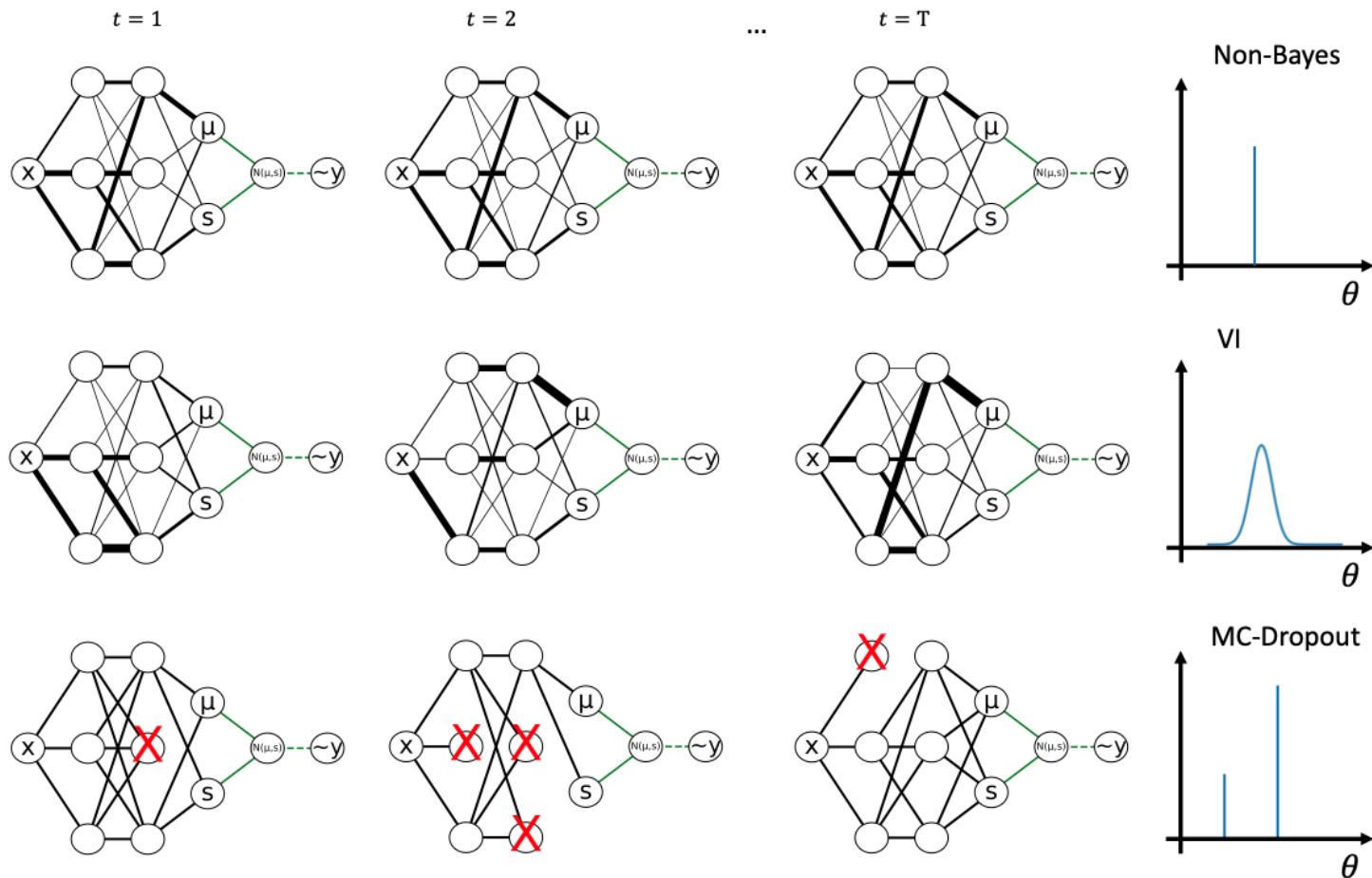
**Weights have Bernoulli-kind distribution**

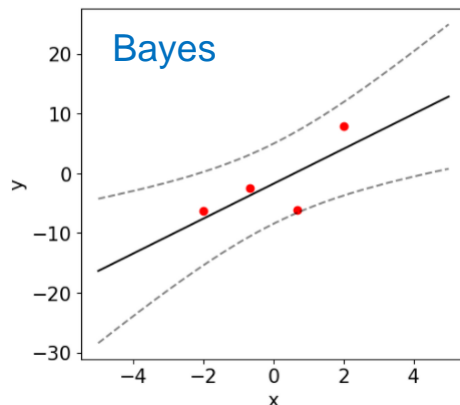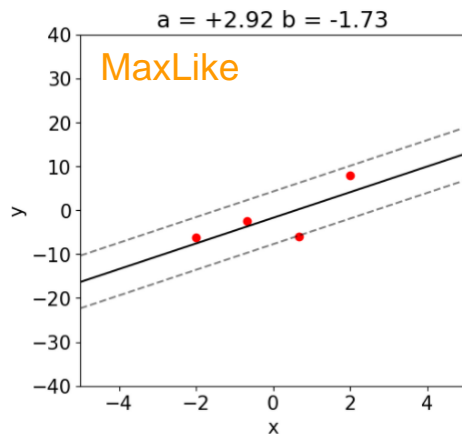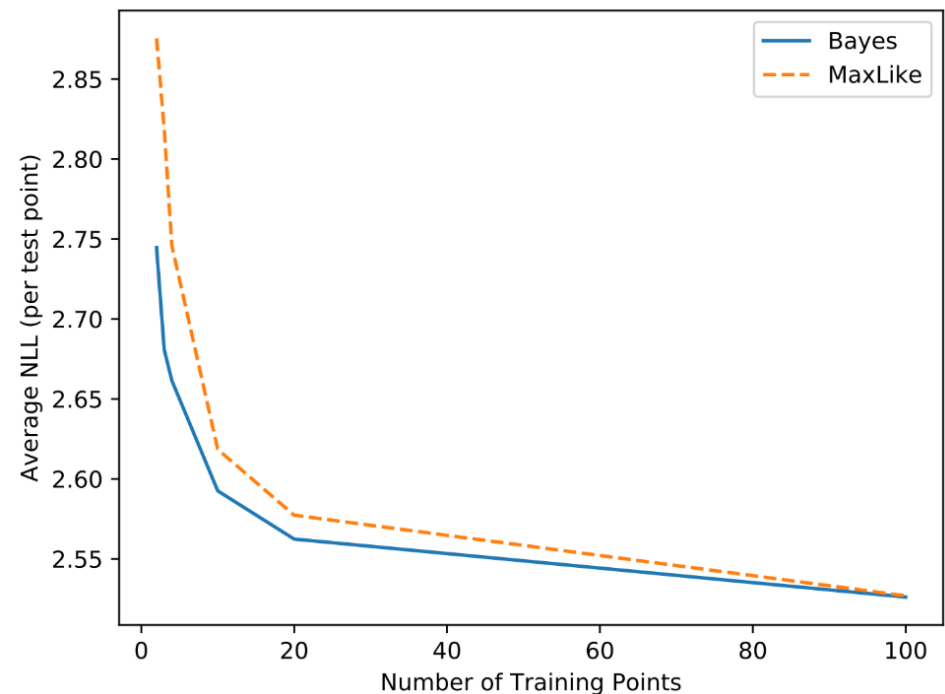**Weights have Gaussian distribution**

# Comparing different Network types

A Non-Baysian NN learns one set of weights: the same input same output
A Bayesian NN learns distribution of weights: same input different outputs

31

# With Bayes we get better prediction performance

A linear regression model via Maximum Likelihood and via Bayes, both trained on the same 4 data points
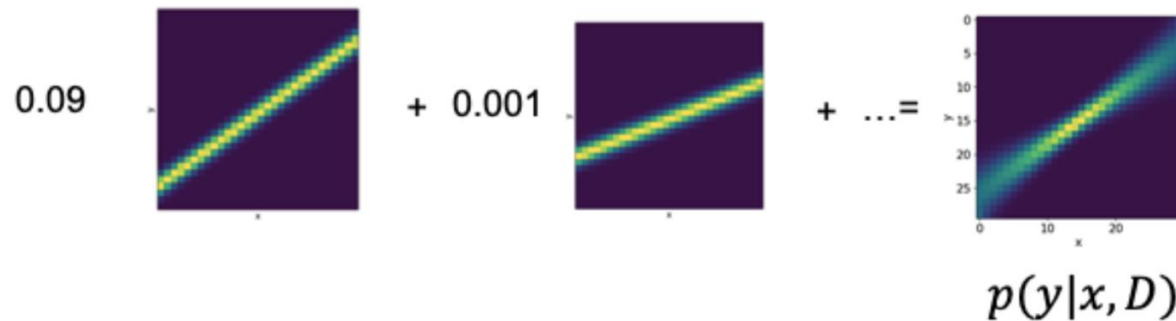
Evaluate the prediction performance of both models on 100 test data points sampled from the same distribution as the 4/10/20/40/70/100 train data points.







**https://github.com/tensorchiefs/dl_book/blob/master/chapter_07/nb_ch07_04.ipynb**

# Summary in 3 lines

$$p(y|x, D) = \int_w p(y|x, w) \cdot p(w|D) \ dw$$



$$p(y|x, D)$$

Not just a single solution

"Integration instead of optimization"

Bayes usually achieves

Better uncertainty estimates

Better predictive performance

# Conclusion

- Standard neural networks (NNs) fail to express their uncertainty (can't talk about the elephant in the room).

- Bayesian neural networks (BNNs) can express their uncertainty.

- BNNs often yield better performance than their non-Bayesian variant.

- Novel classes can be better identified with BNNs, which combine epistemic and aleatoric uncertainties compared to standard NNs.

- Variational inference (VI) and Monte Carlo dropout (MC dropout) are approximation methods that allow you to fit deep BNNs.

- TFP provides easy to use layers for fitting a BNN via VI.

- MC dropout can be used in Keras for fitting BNNs.