

# Machine Intelligence:: Deep Learning

## Week 7

*Beate Sick, Jonas Brändli, Oliver Dürr*

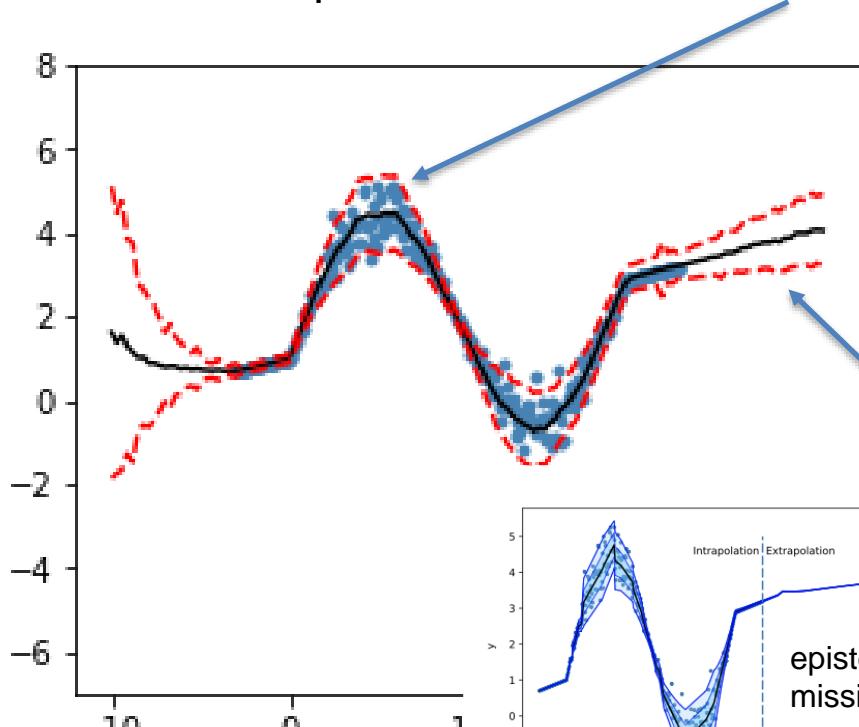
Improving the performance and uncertainty estimates of NN models by taking into account the epistemic uncertainty.

## Outline:

- Issues with current DL approach
  - No uncertainty for the fitted weights  
→ epistemic uncertainty is ignored causing different problems:
    - No increased uncertainty in case of extrapolation
    - Deficits in prediction performance
- Approaches to take epistemic uncertainty into account:
  - Deep Ensembling
  - MC Dropout
  - Bayesian Modelling
    - Introduction to Bayesian Statistics
    - A simple example (Bayes the Hacker's way)
    - Variational Approximation

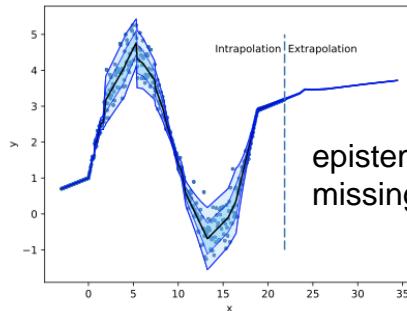
# Aleatory vs. Epistemic Uncertainty

Much spread in train data → aleatory uncertainty (from latin “[Alea Acta est](#)”)



Caesar: Die Würfel sind gefallen!

We understand a dice perfectly (no epistemic uncertainty), but still there is a aleatoric uncertainty about the the number that will show up next when rolling the dice.



No (or few) train data  
→ no (or few) “knowledge”  
(from [Ancient Greek ἐπιστήμη](#) (*epistémē*) 'knowledge')  
→ **epistemic uncertainty**



- *Aleatoric uncertainty* is due to the uncertainty, that is inherent in the data.
- The uncertainty when leaving the ‘known ground’ is called *epistemic uncertainty*.

# The elephant in the room

A high performant NN  
(trained on imageNet data)  
does not see the elephant!

GENERAL FACE NSFW COLOR MORE MODELS ▾



VIEW DOCS

PREDICTED CONCEPT	PROBABILITY
group	0.979
adult	0.977
people	0.976
furniture	0.960
room	0.957
business	0.903
indoors	0.901
man	0.896
seat	0.895

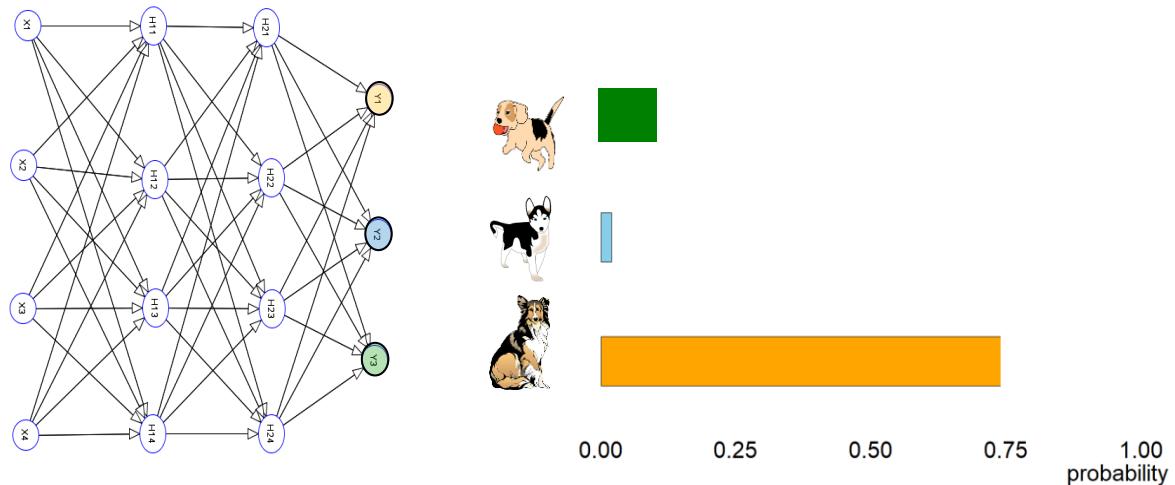
# Elephant in the room



- Aufgabe:

[https://github.com/tensorchiefs/dl\\_course\\_2022/blob/master/notebooks/18\\_elephant\\_in\\_the\\_room.ipynb](https://github.com/tensorchiefs/dl_course_2022/blob/master/notebooks/18_elephant_in_the_room.ipynb)

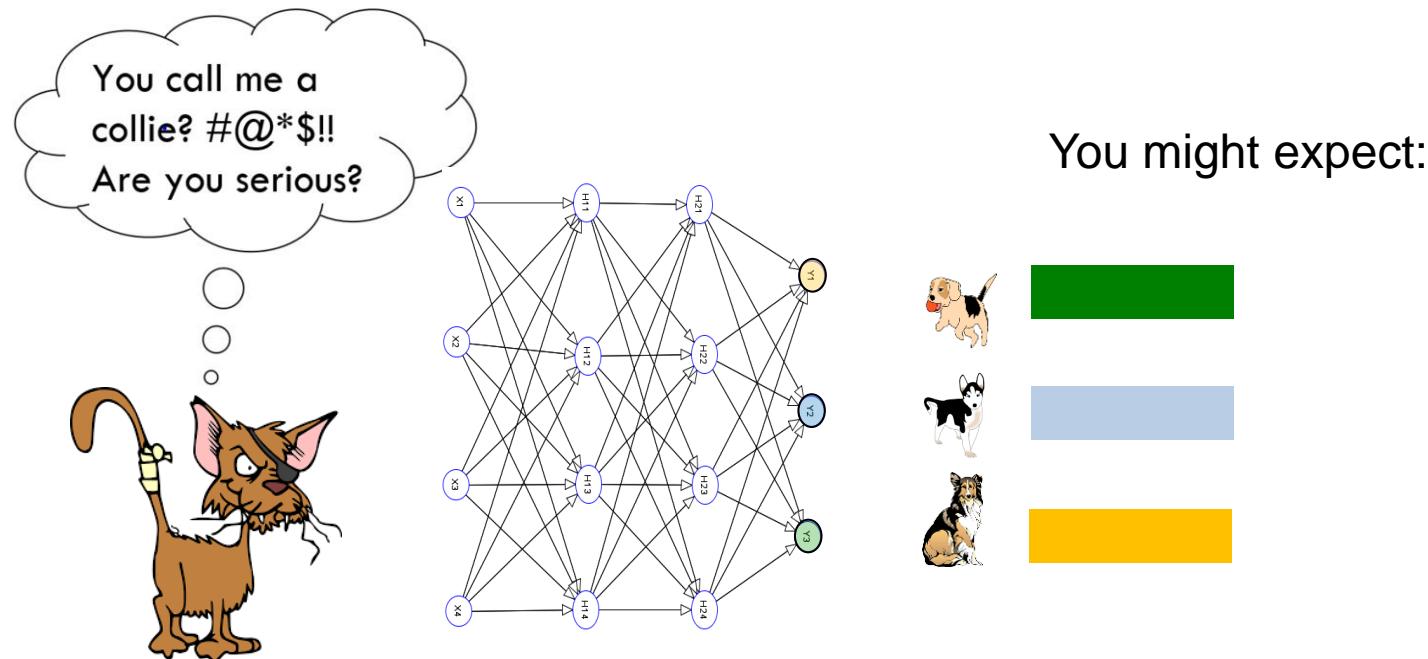
# CNNs have high performance on in-distribution examples



**CNNs yield high accuracy and calibrated probabilities, but...**

# A classical NN cannot ring the alarm in case of out-of-distribution (OOD) examples

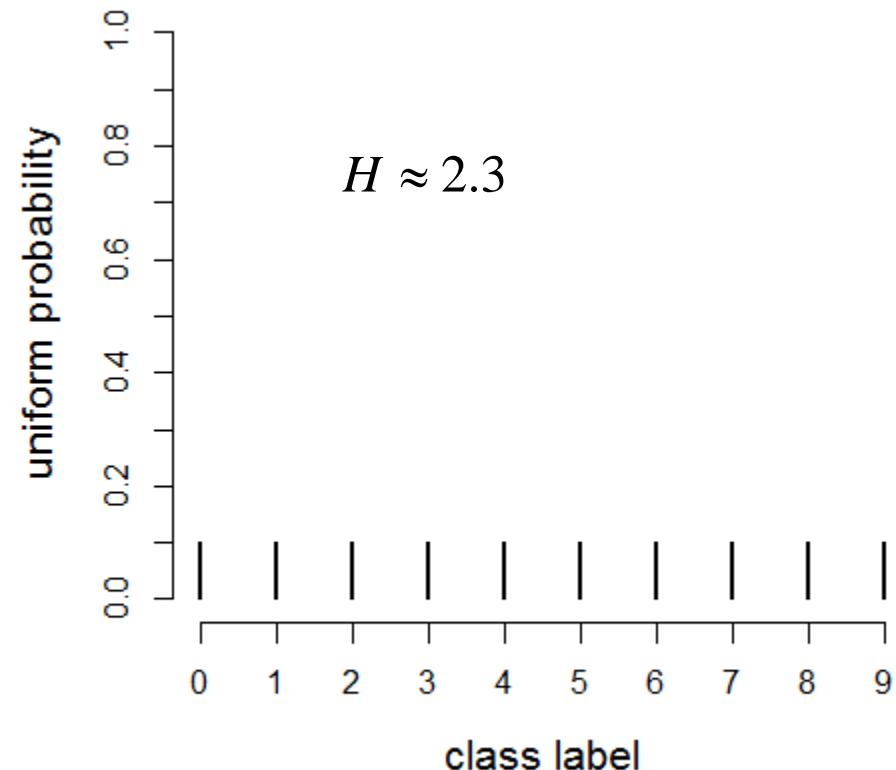
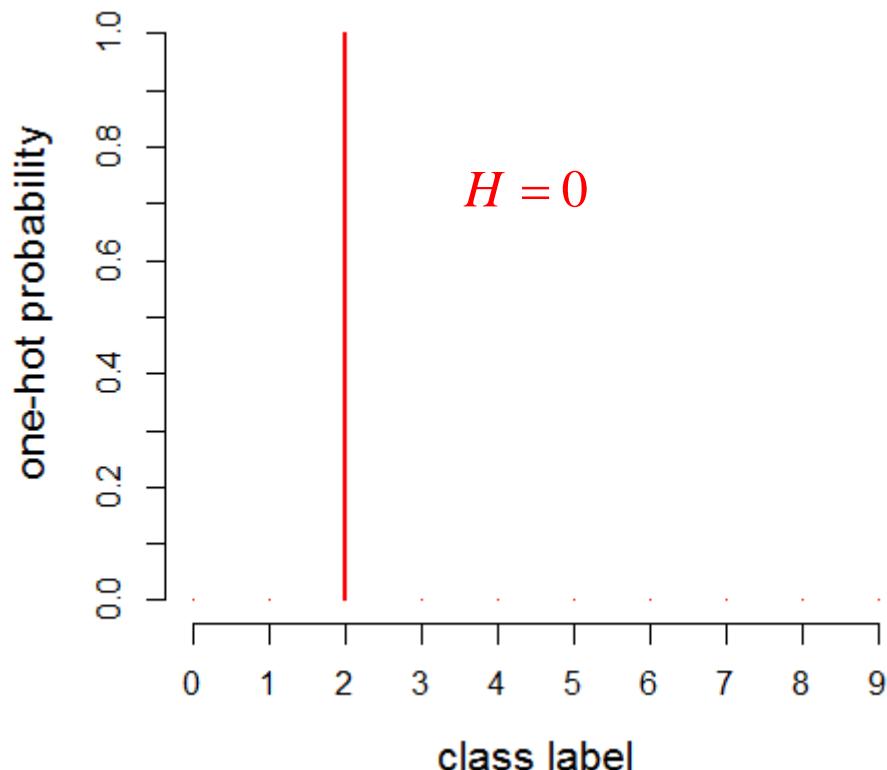
What happens if we present a novel class to the CNN?



## Recall: Entropy as measure for uncertainty

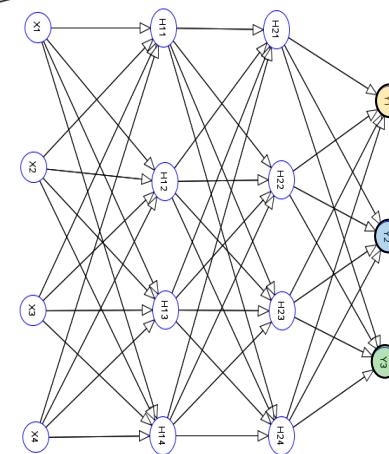
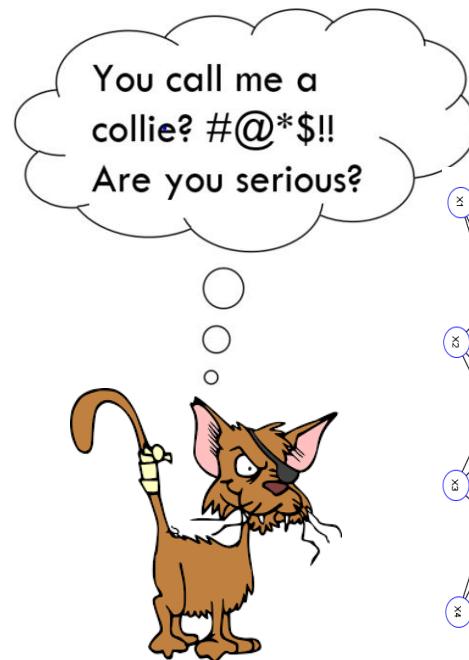
$$H(P) = - \sum_i p_i \cdot \log(p_i)$$

Entropy is a measure for “untidiness” or uncertainty.  
If a distribution has only one peak, it is tidy and  $H=0$   
If all outcomes are equally probable it is maximal untidy

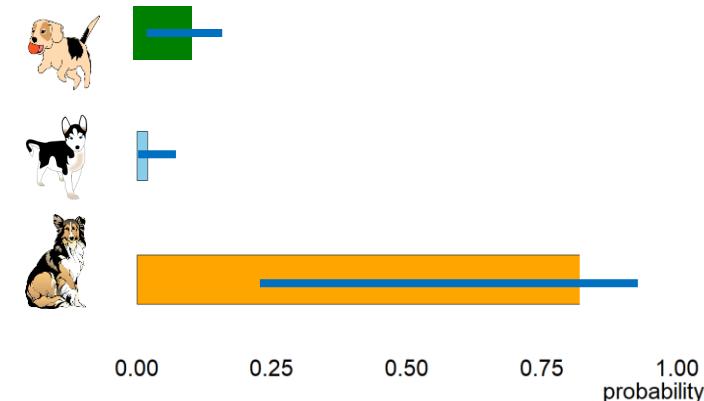


# A classical NN cannot ring the alarm in case of out-of-distribution (OOD) examples

What happens if we present a novel class to the CNN?



But you probably get:  
**Plain wrong !**



We need some error bars!  
We need epistemic uncertainty!

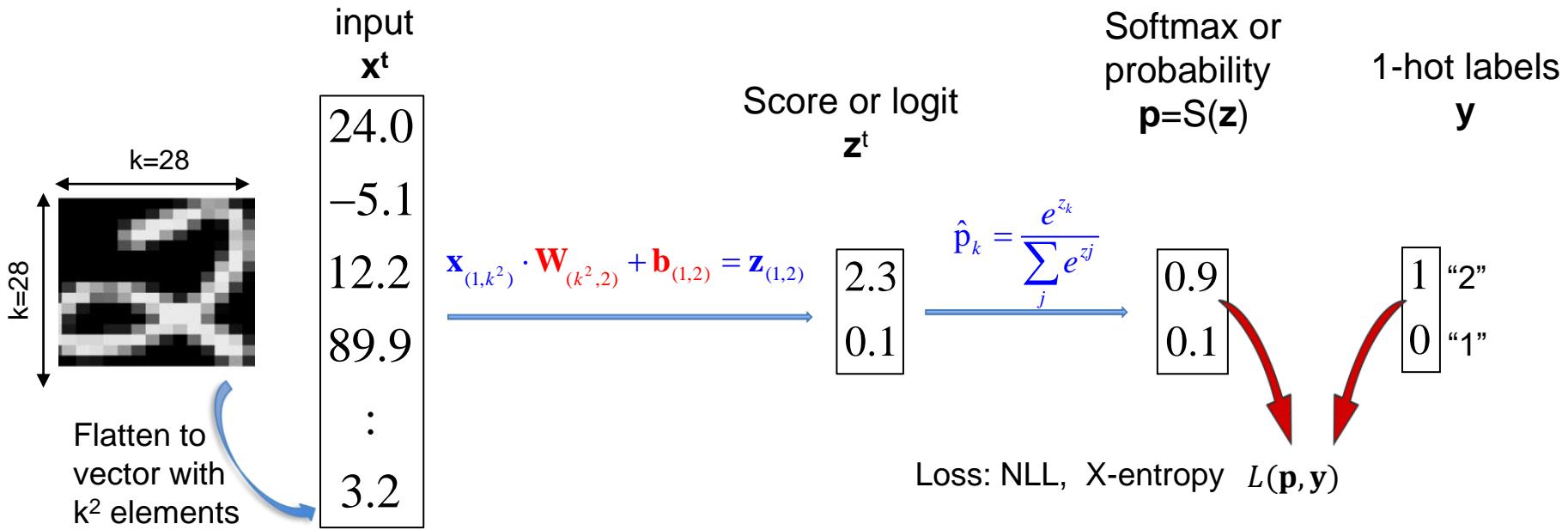
# Importance to detect OOD (out of distribution)



- Current DL Systems bad in out of distribution OOD situations
- Application need at least to detect OOD situations

# Deep ensembles

## Recall the training of NN models via SGD



Take step in direction of descent gradient:  
(the gradient is oriented orthogonal to contour lines)

$$w_i^{(t)} = w_i^{(t-1)} - \varepsilon^{(t)} \left. \frac{\partial L(\mathbf{w})}{\partial w_i} \right|_{w_i=w_i^{(t-1)}}$$

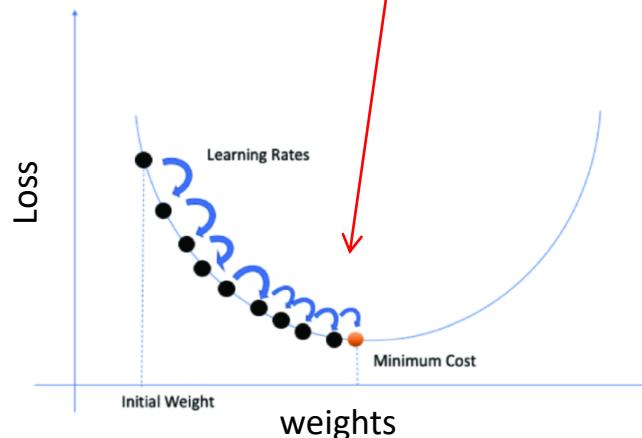
$$\text{II} \\ - \sum_{k=1}^2 y_k \cdot \log(p_k)$$

Loss = NLL averaged over all images in mini-batch

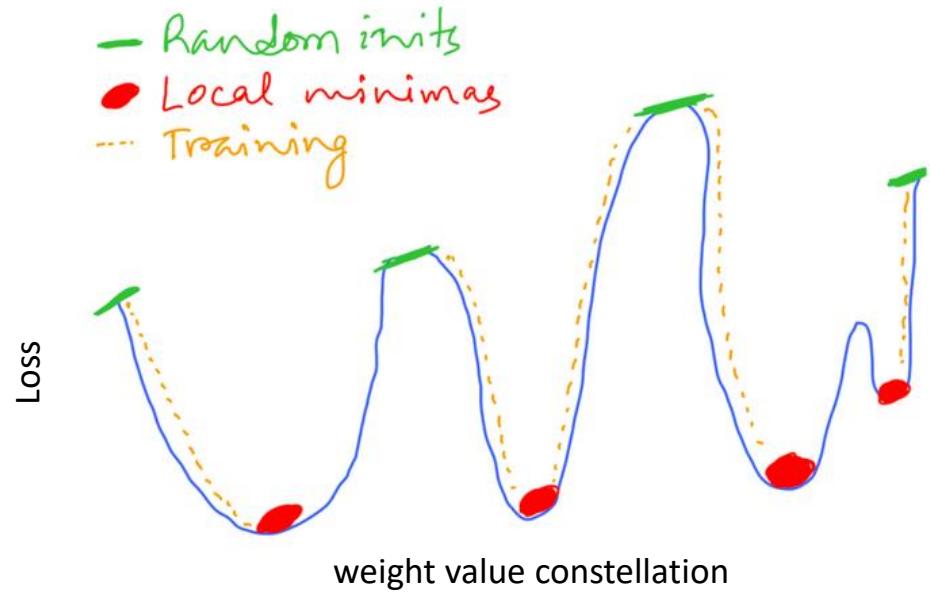
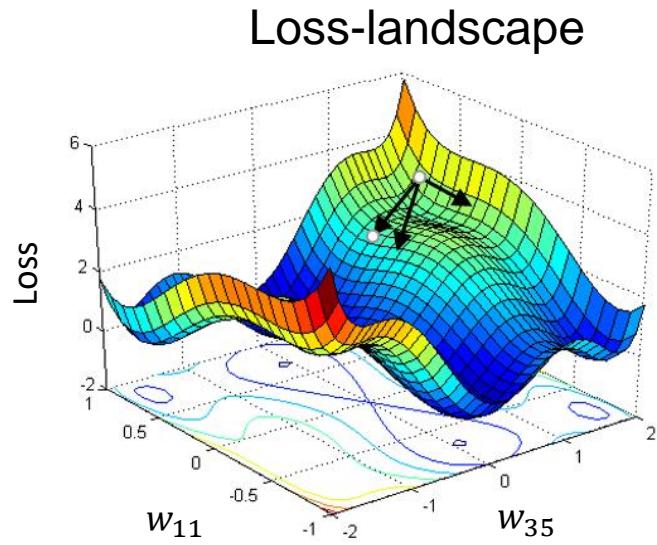
$$\text{NLL} = \frac{1}{N} \sum_i L(\mathbf{p}_i, \mathbf{y}_i)$$

$$\text{Loss} \\ C(w_1, w_2)$$

[Image credits](#)



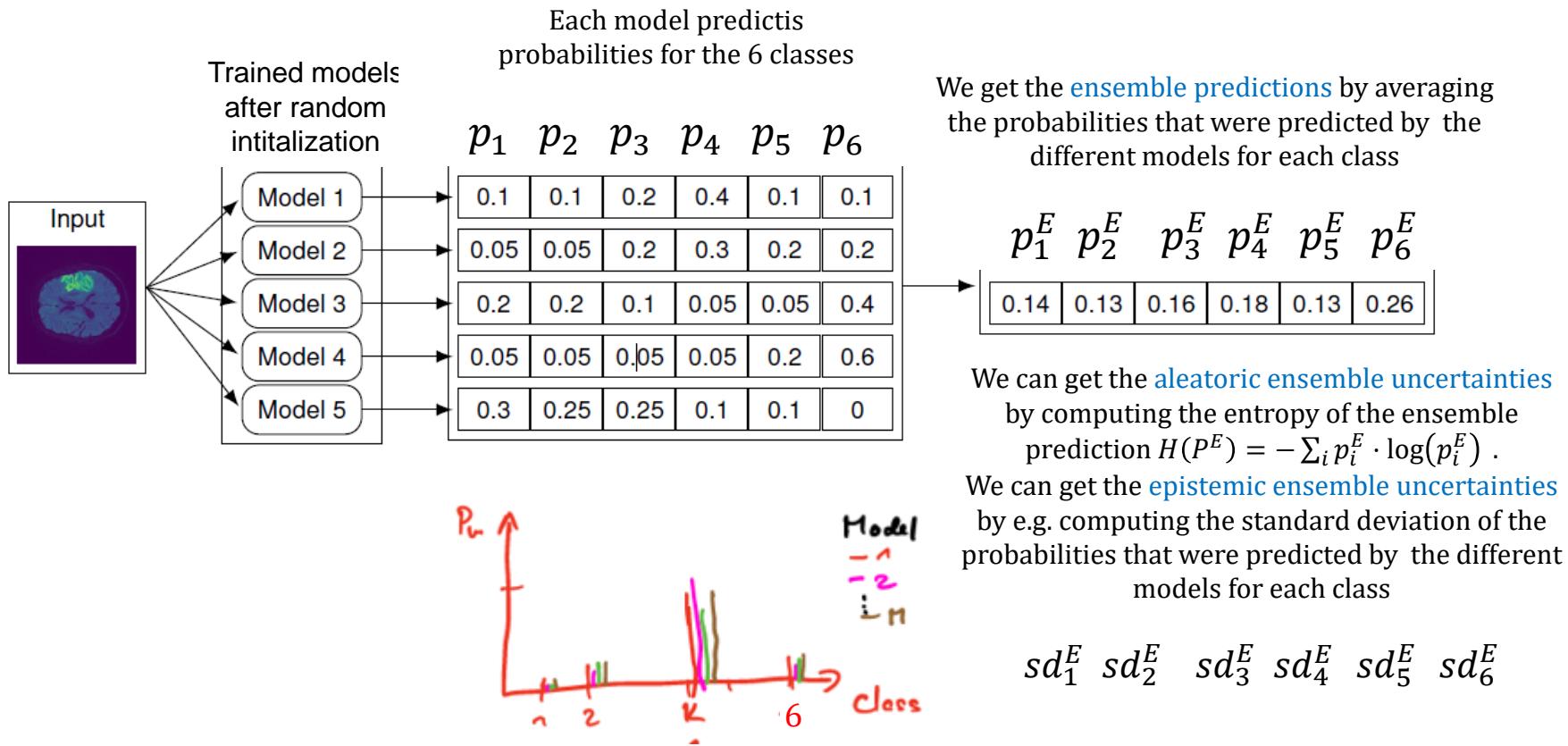
# The loss-landscape in DL is usually not convex



The loss-landscape of DL models has many local minima with similar depth.

Training is started with a random weight value initialization → training the NN with the same data several times is usually ending in different local minima.

# Deep ensembling: Train several NN models and average their predictions



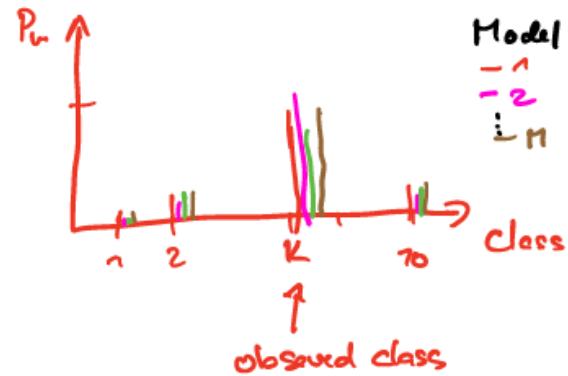
Nice:

For the convex NLL loss, it is guaranteed, that the NLL of the ensemble prediction is better (smaller or equal) than the average NLL of the individual models.

# Ensembling improves the NLL performance

Ensemble prediction for an observation  
with observed class =  $K_1$ , based on  $M$  models:

$$P_K^E = \frac{1}{M} \sum_{m=1}^M P_K^m \quad K = 1 \dots 10 \text{ for 10 classes}$$

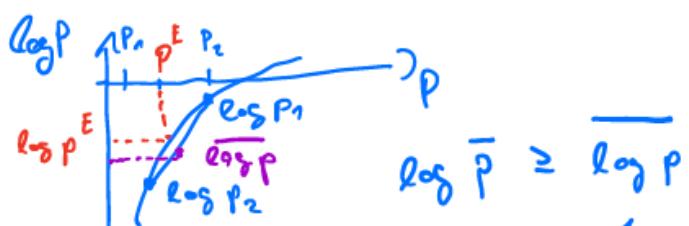


associated NLL contribution  $\ell$ :

$$\text{Model } m: \ell_m = -\log P_{km} ; \text{ Ensemble: } \ell^E = -\log P_K^E$$

to show  $\bar{\ell} \geq \ell^E$

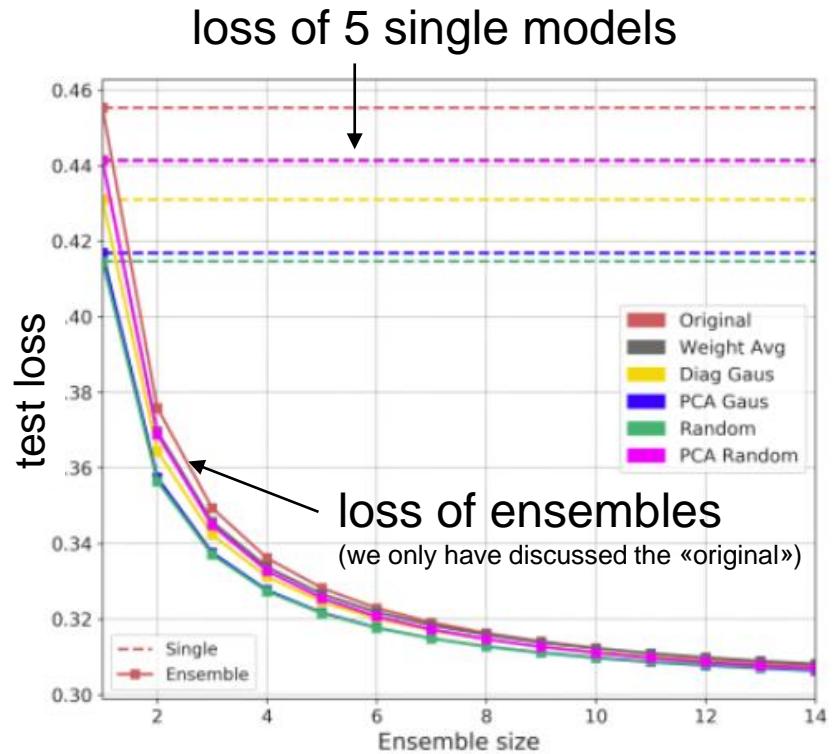
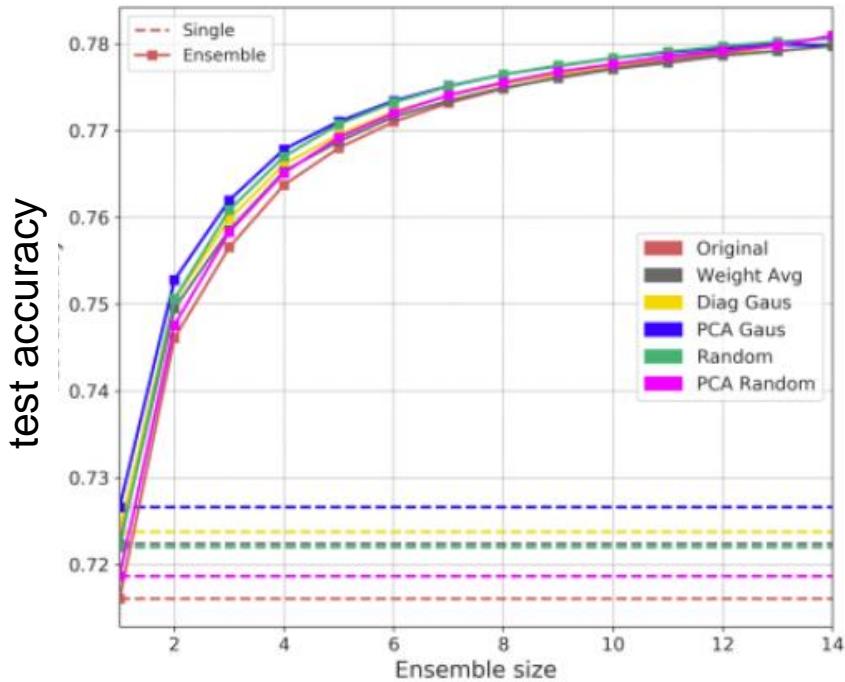
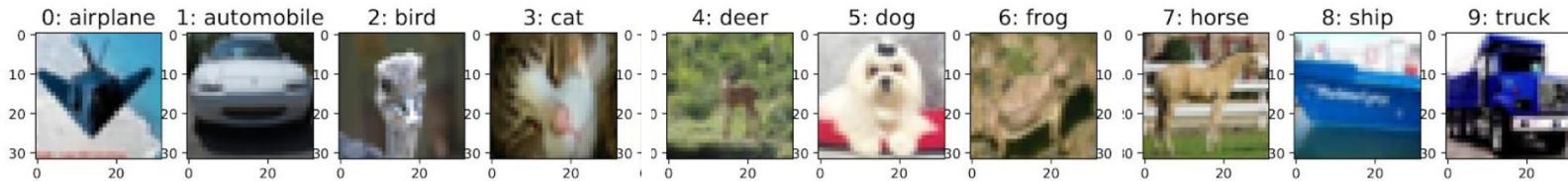
use Jensen inequality



$$\begin{aligned} \text{with } \bar{\ell} &= \frac{1}{M} \sum_{m=1}^M \ell_m = \frac{1}{M} \sum_m -\log P_{km} \\ &= -\underbrace{\frac{1}{M} \sum_m \log P_{km}}_{\leq \log \frac{1}{M} \sum_m P_{km}} \\ &\geq -\log P_K^E = \ell^E \quad \square \end{aligned}$$

$$\log \frac{1}{M} \sum_m P_m \geq \frac{1}{M} \sum_m \log P_m$$

# Deep ensembling improves prediction power

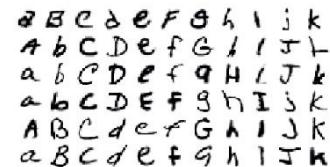
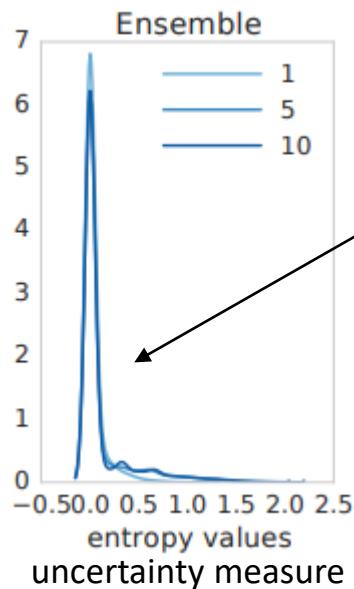


Ensembles with as few as 3 or 5 members are typically enough to achieve a performance gain.

# Deep ensembles improve uncertainty measures

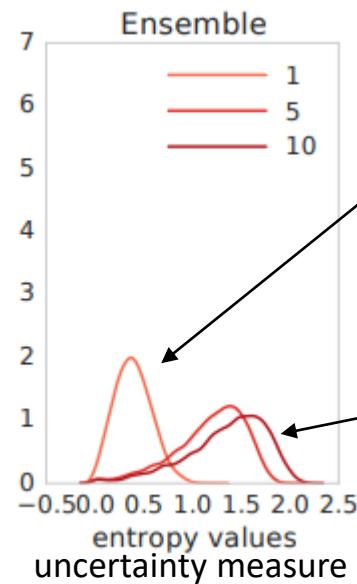
We want that a model, that is trained on normal MNIST letter data, should provide large uncertainties when applied on novel (NOT-MNIST) letter images.

Uncertainties for prediction of normal MNIST letter images



Ensemble uncertainty for known class stay small for all sizes of the ensemble

Uncertainties for prediction of NOT-MNIST images

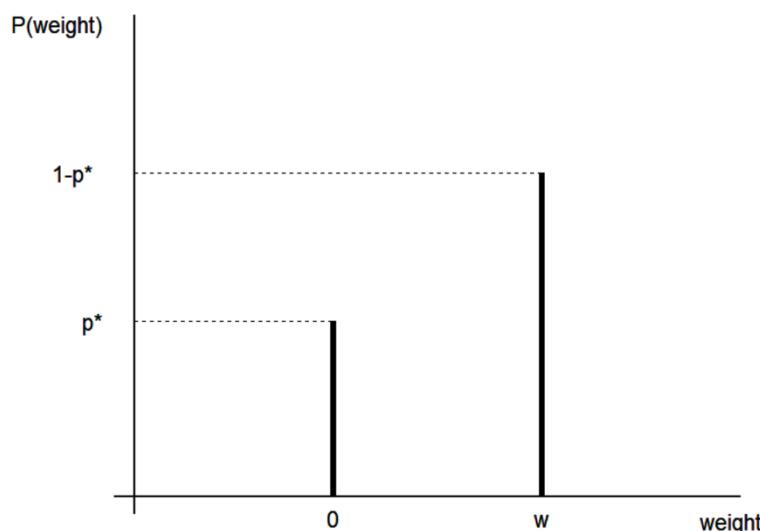
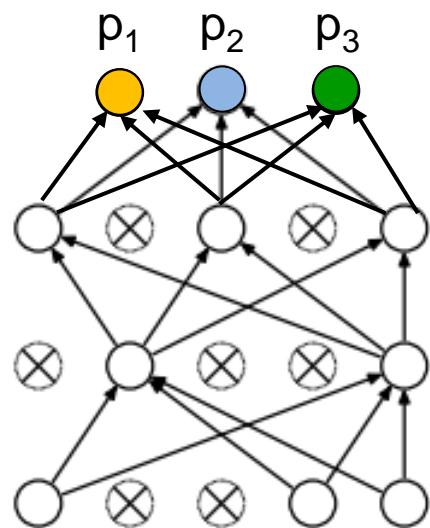


A single models shows less uncertainty for a novel class than an ensemble model

Ensemble uncertainty measures for novel classes improve with size of ensemble

# Dropout

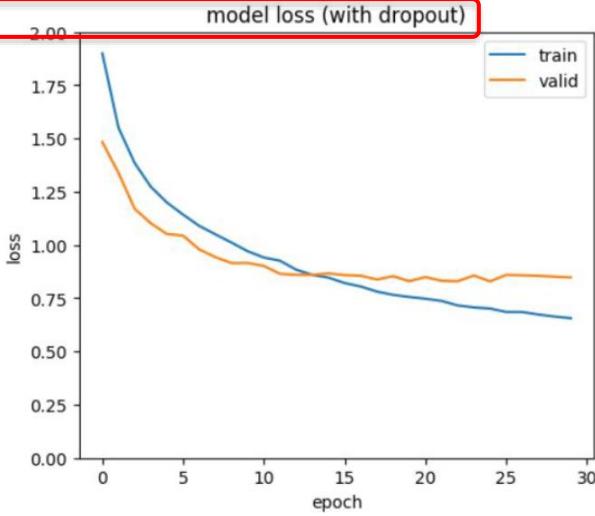
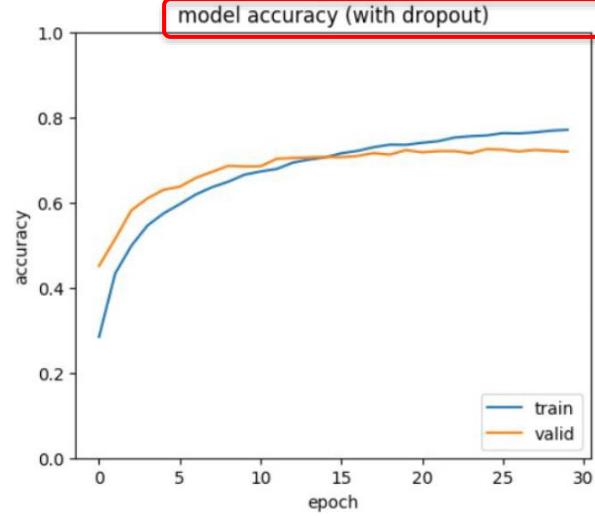
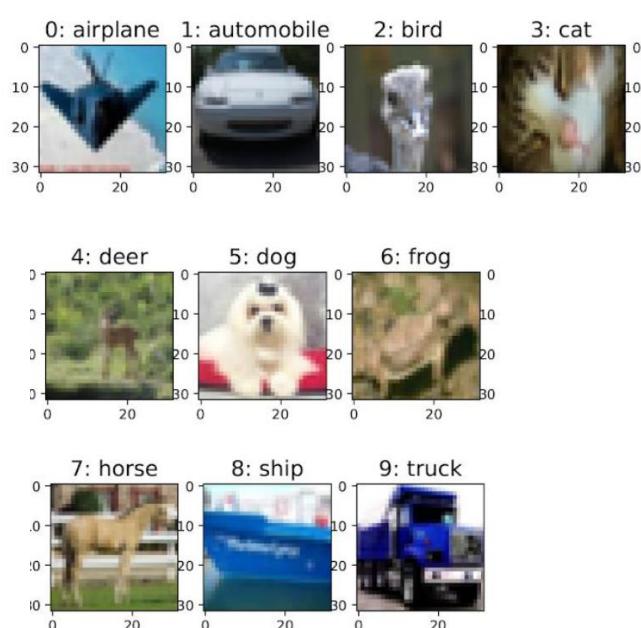
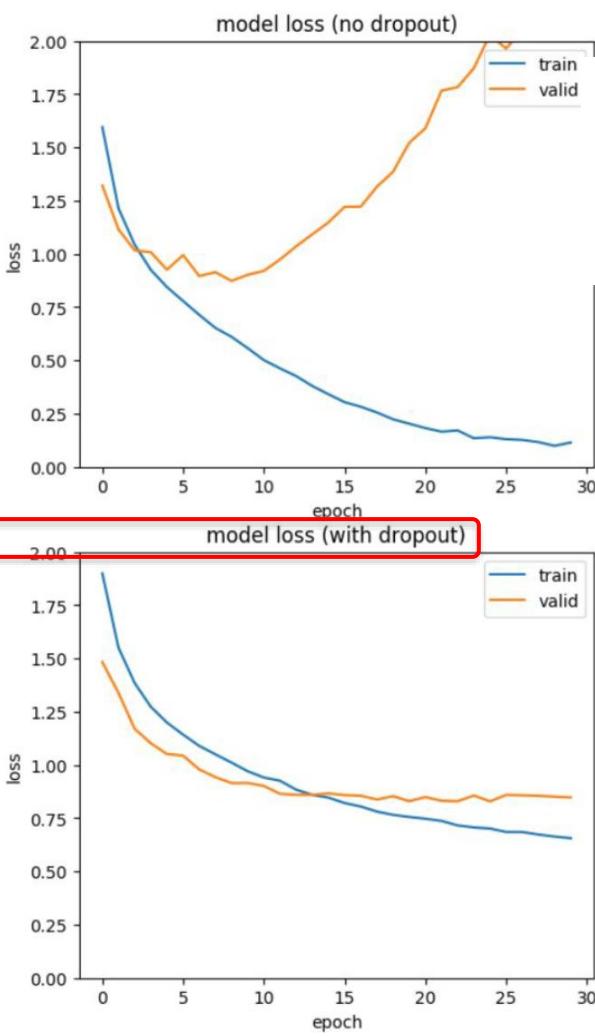
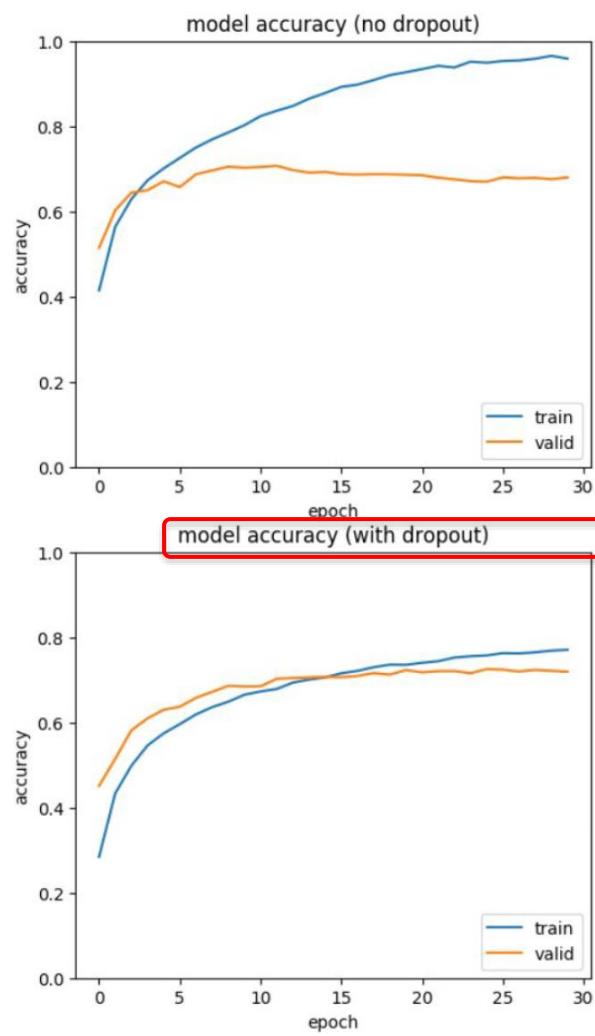
## Recall: Classical Dropout only during training



Using **dropout during training** implies:

- In each training step only weights to not-dropped units are updated → we train a sparse sub-model NN
- For non-Bayesian NN we freeze the weights after training to a value  $w \cdot p^*$

# Recall: Dropout fights overfitting in a CIFAR10 CNN



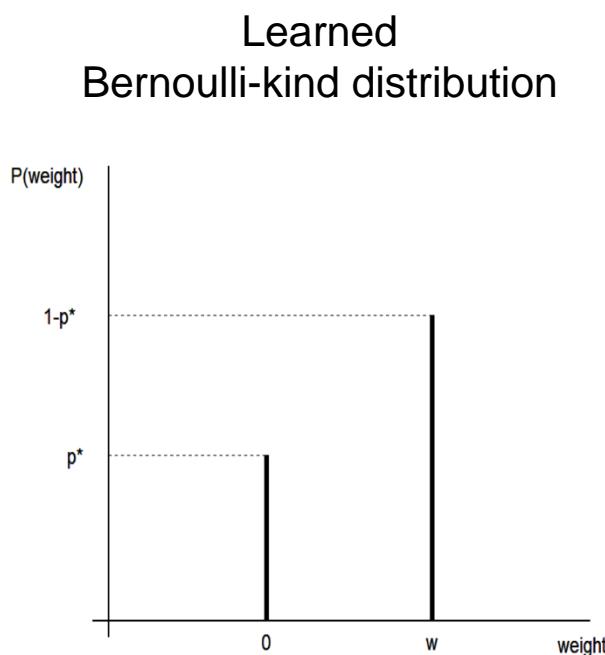
From Dropout during training  
to MC Dropout during test time

# Bayesian NN via MC Dropout

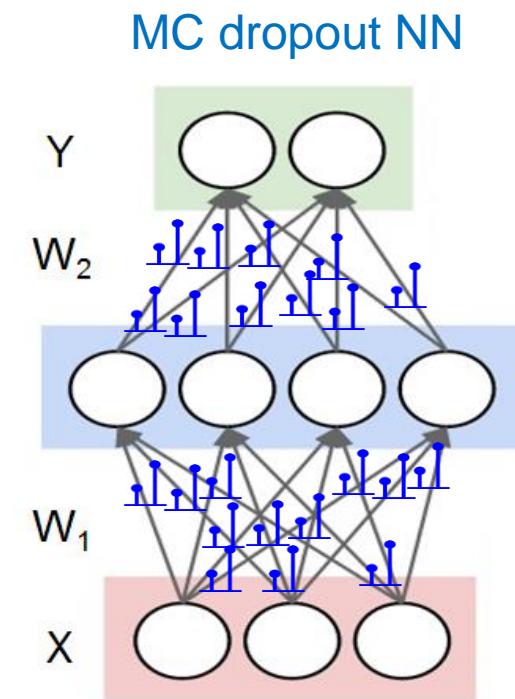
Yarin Gal et al. (2015):

Via Dropout training we learned a whole weight distribution for each connection.  
We can sample from this Bernoulli-kind weight distribution by performing  
dropout during test time and use the dropout-trained NN as Bayesian NN.

Gal showed that doing dropout approximates VI with a Bernoulli-kind variational distribution  $q_\theta$  (instead of a Gaussian).



Dropout in test time



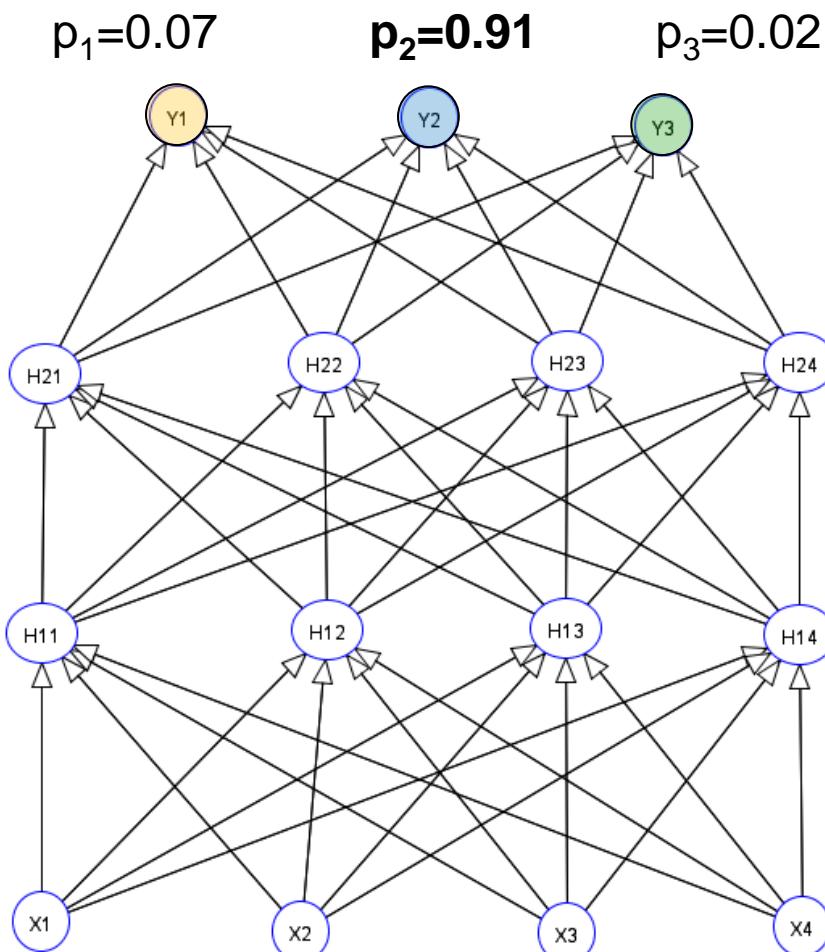
Weights have  
Bernoulli-kind distribution

Which parameter has this  $q_\theta$ ?

The value  $w$ .

# When using Dropout only during training

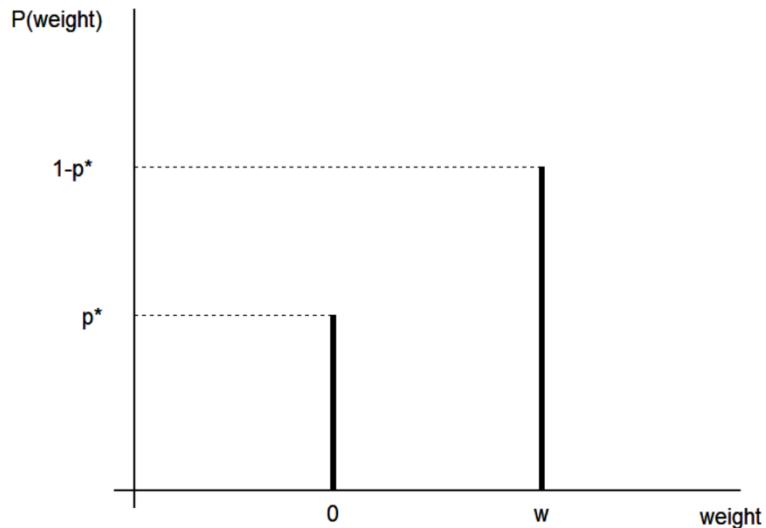
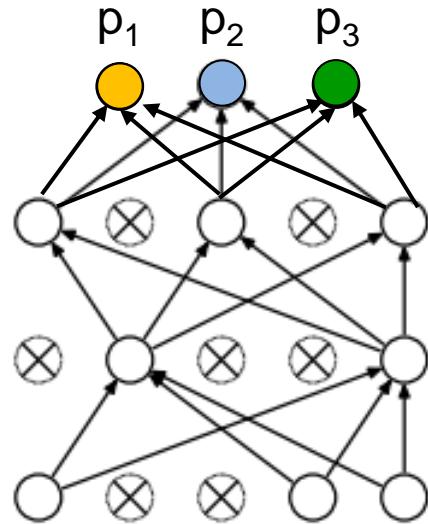
For non-Bayesian NN we freeze the weights after training to a value  $w \cdot p^*$  and use then the trained NN for prediction:



Probability of predicted class:  $p_{\max}$

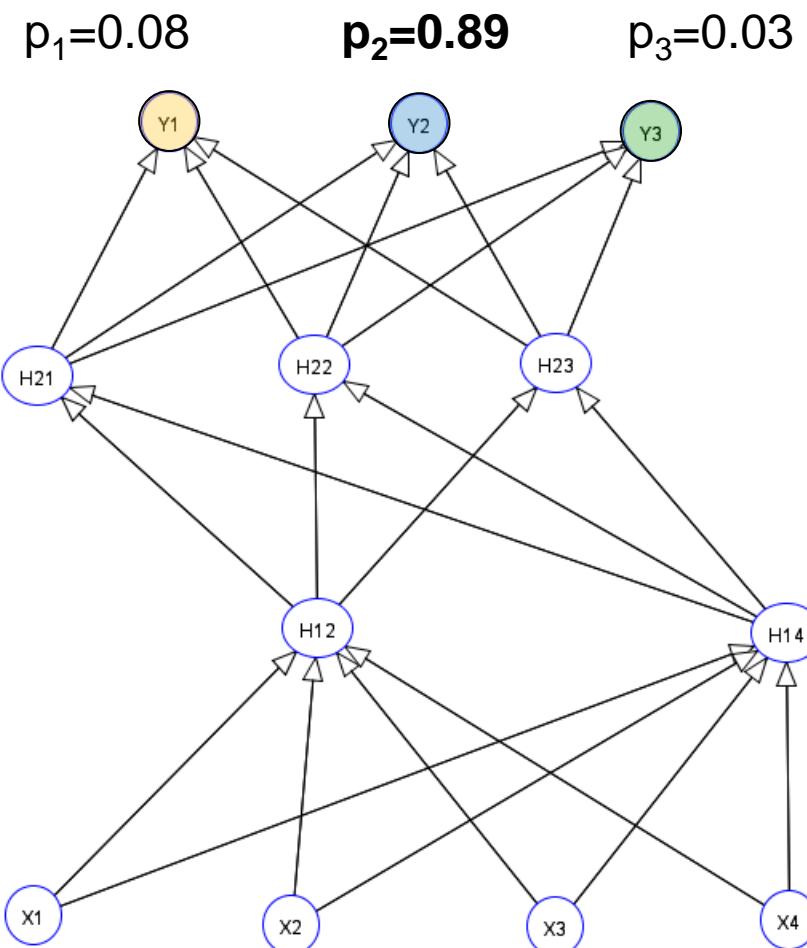
Input: image pixel values

# MC Dropout: we also perform dropout during test time



In each prediction instance we dropout a random subset of nodes, which corresponds to setting all weights starting from these nodes to zero.

# MC Dropout during test time: Run 1

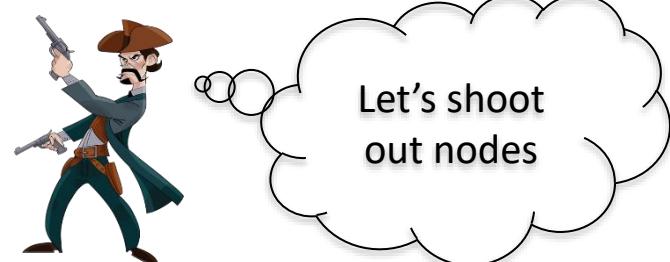
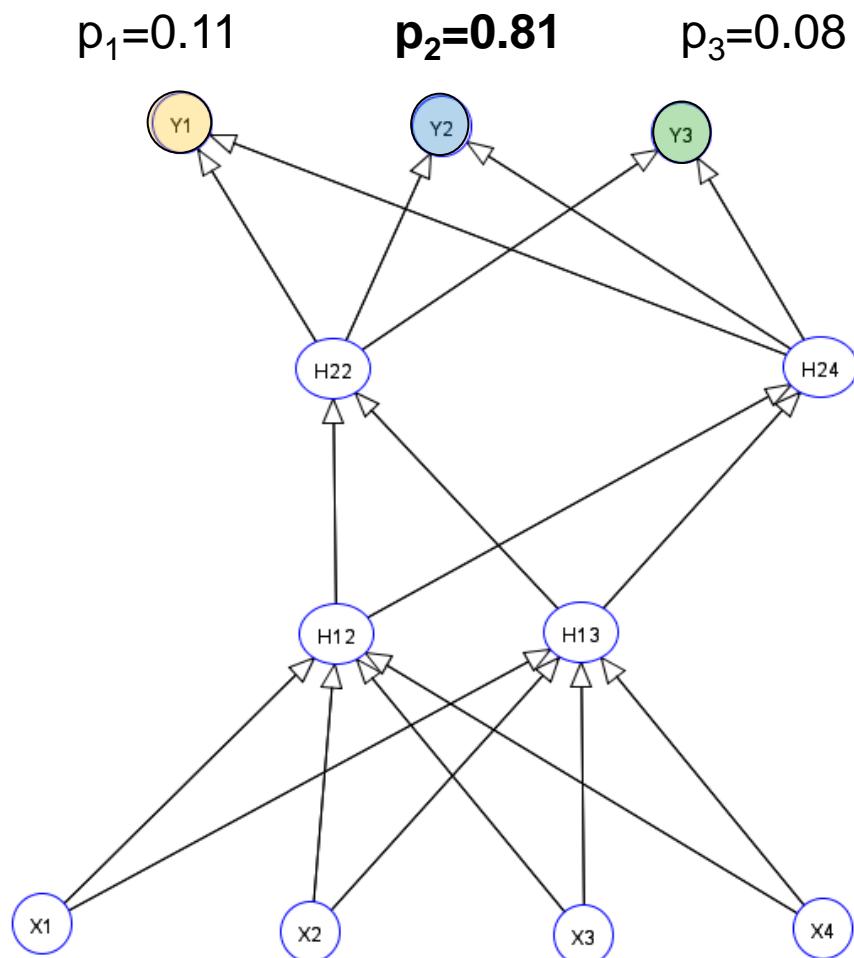


Let's shoot  
out nodes

Stochastic dropout of units

Same input image

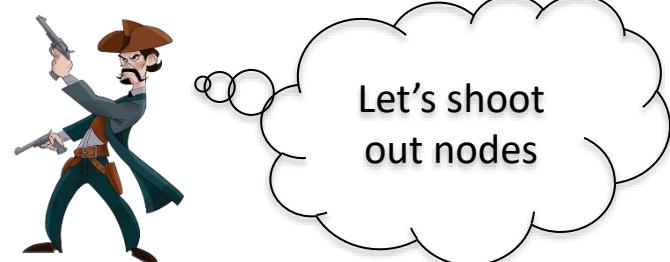
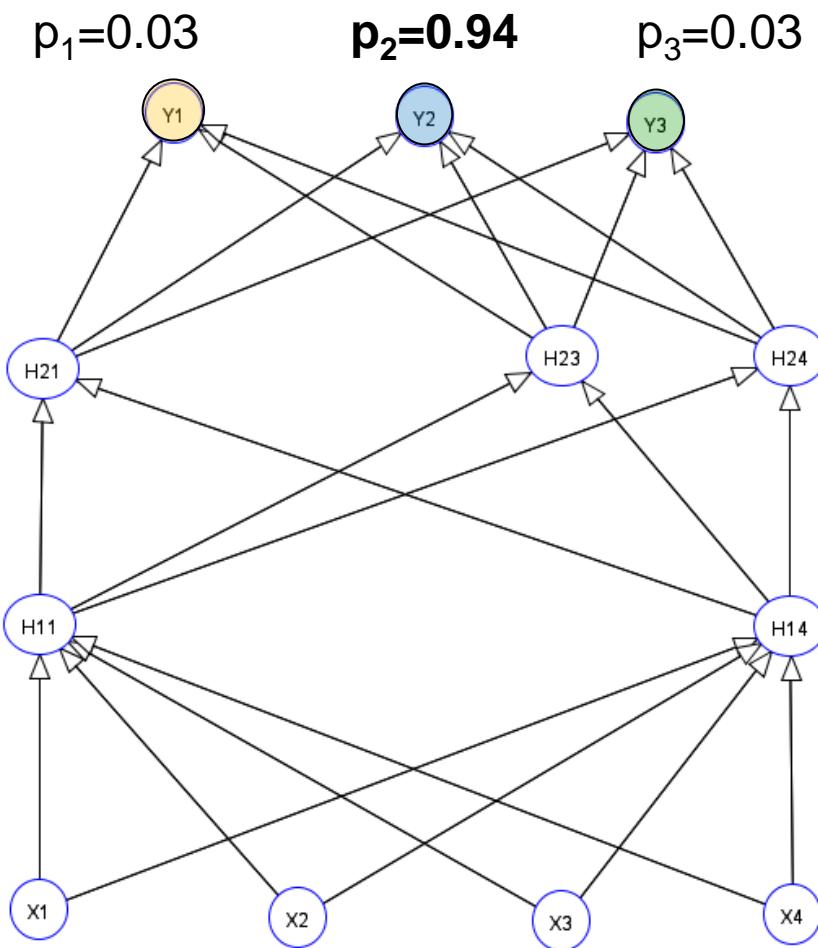
# MC Dropout during test time: Run 2



Stochastic dropout of units

Same input image

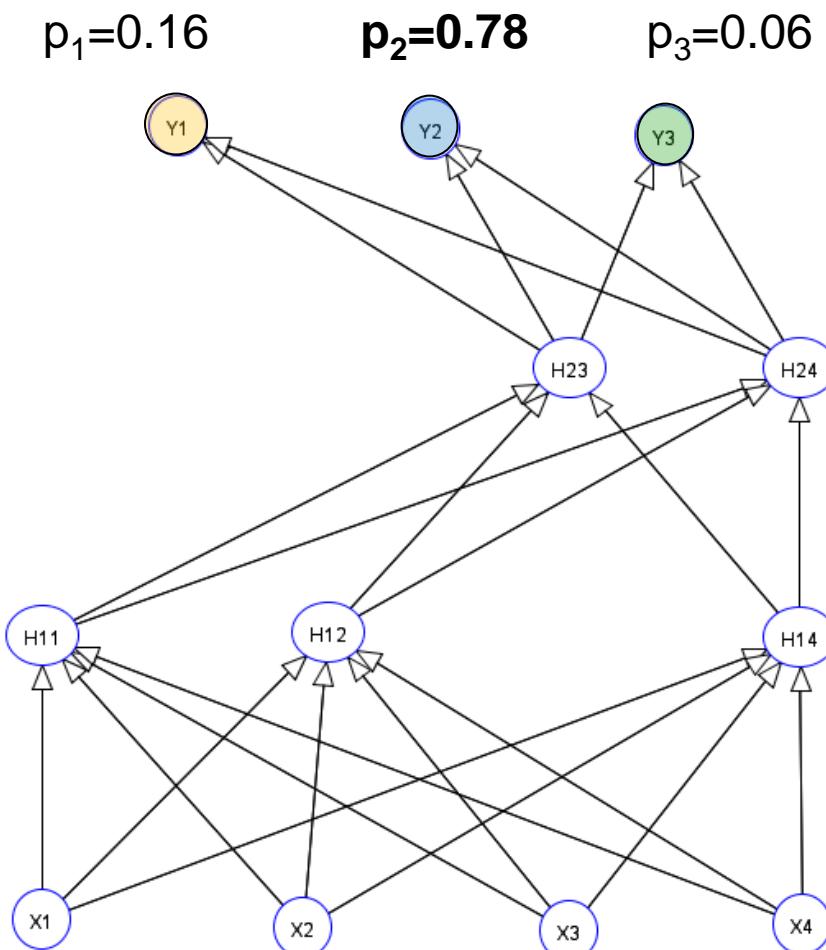
# MC Dropout during test time: Run 3



Stochastic dropout of units

Same input image

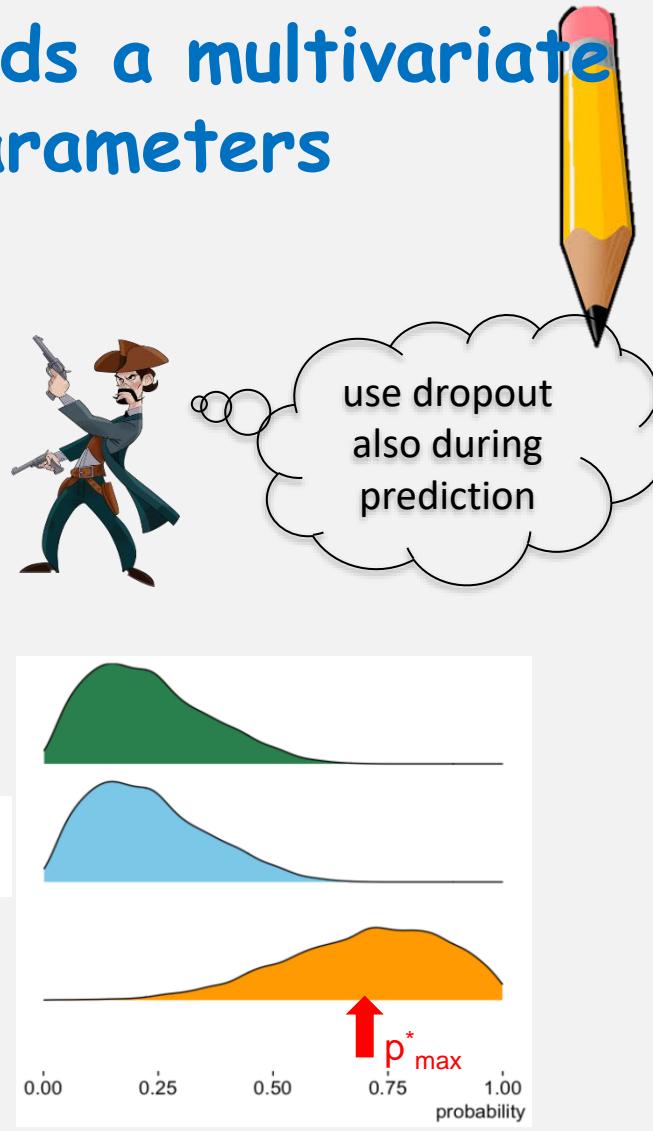
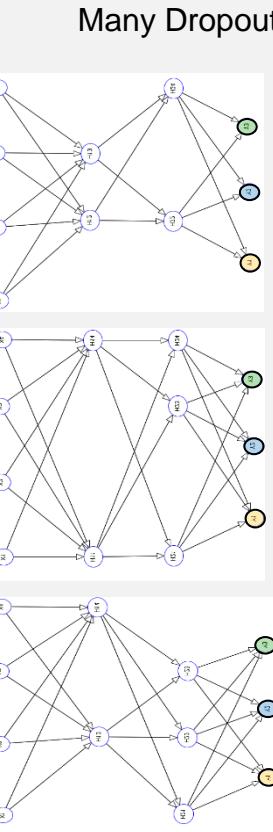
# MC Dropout during test time: Run 4



Stochastic dropout of units

Same input image

# MC Dropout during test time yields a multivariate predictive distribution for the parameters

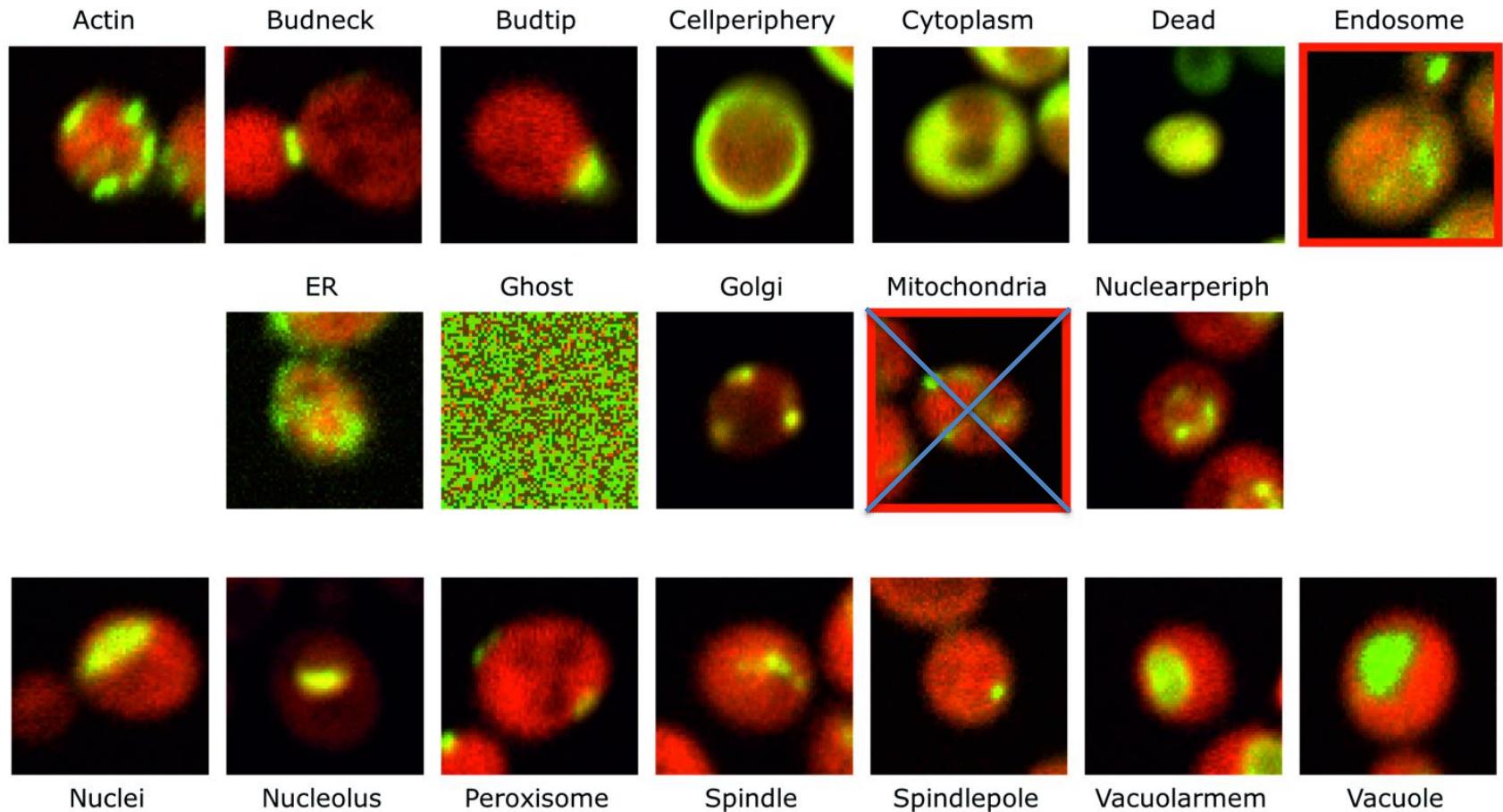


CNN predicts class “collie”  
but with high uncertainty

...

Remark: Mean of marginal give components of mean in multivariate distribution.

# Experiment with unknown phenotype

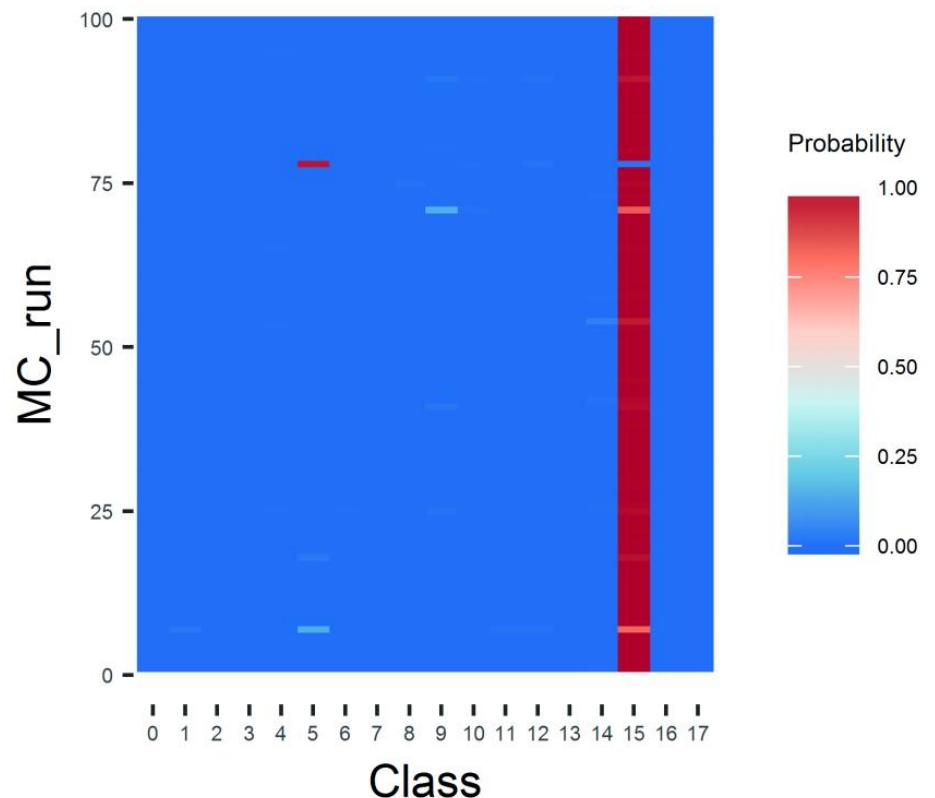


Dürr O, Murina E, Siegismund D, Tolkachev V, Steigrale S, Sick B. Know when you don't know, Assay Drug Dev Technol. 2018

# Probability distribution from MC dropout runs

## Image with known class 15

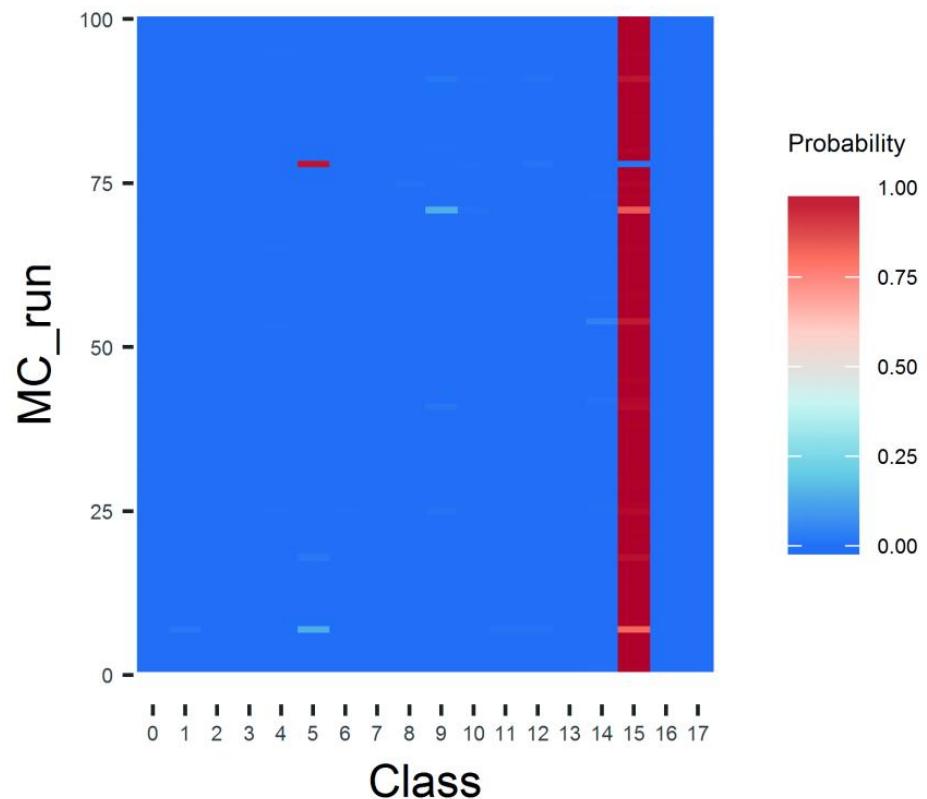
100 MC predictions for an image with known phenotype 15



# Probability distribution from MC dropout runs

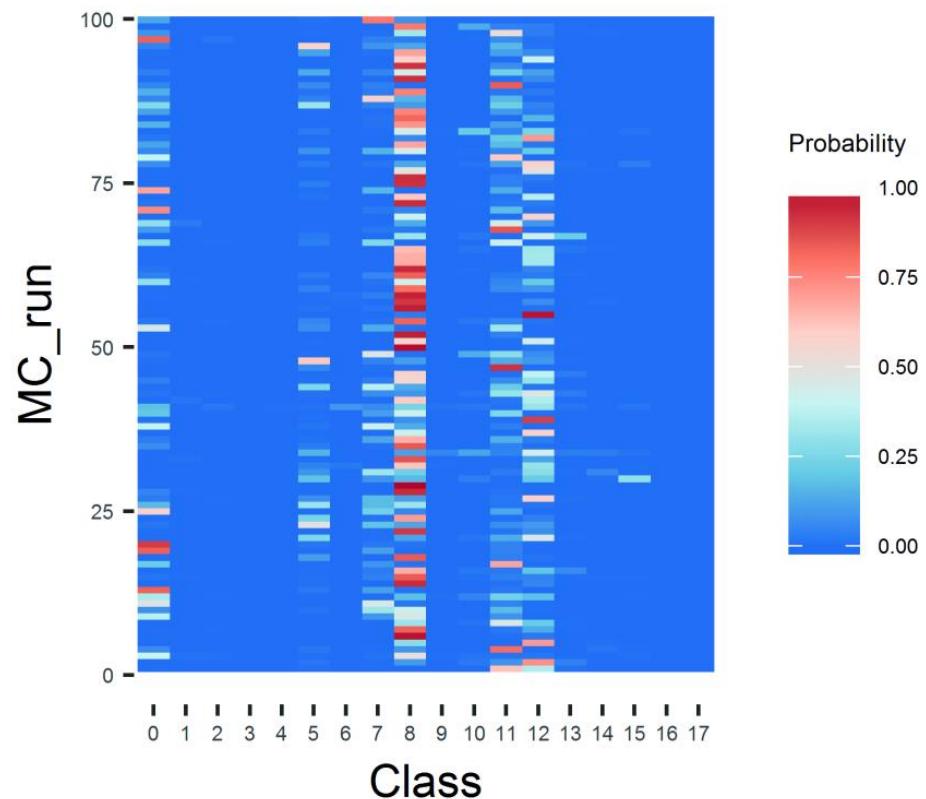
**Image with known class 15**

100 MC predictions for an image with known phenotype 15



**Image with unknown class**

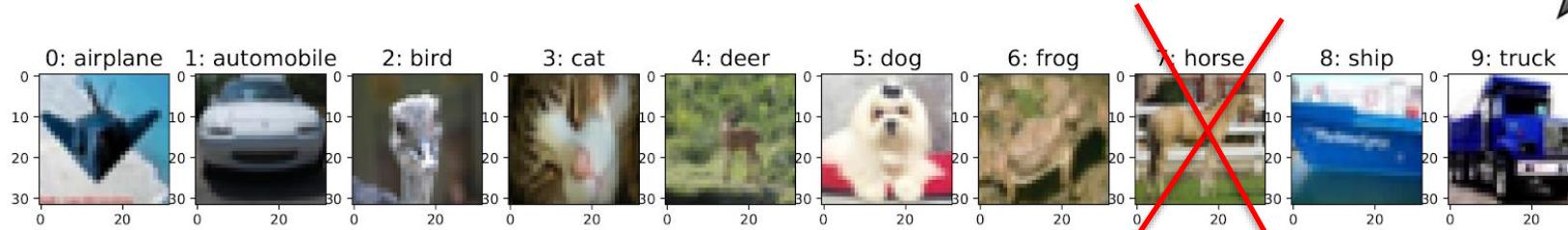
100 MC predictions for an image with an unknown phenotype



# Hands-on Time



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.



Cifar10 classification case study with novel class

Imports

Loading and preparation of the dataset

Non-Bayesian CNN

Variational Inference

MC Dropout

Getting mc dropout predictions

Accuracy on the the known lables in the train set for all three models

Non-Bayesian prediction

Bayesian VI prediction

Bayesian MC prediction

Predicted classes for the unknown class

Compare the predictions for a known and unknown image

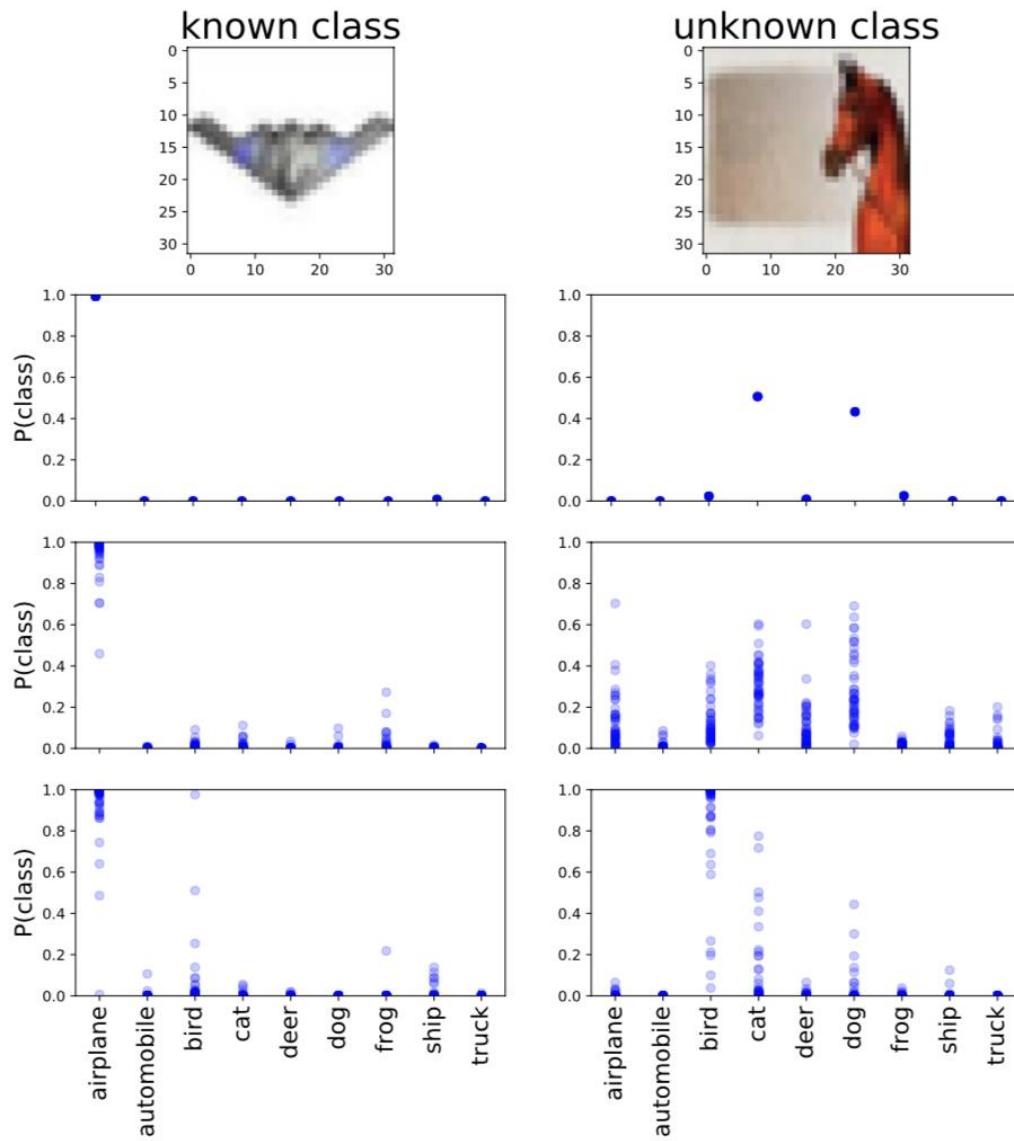
Please focus on MC Dropout

[https://github.com/tensorchiefs/dl\\_course\\_2022/blob/master/notebooks/20\\_cifar10\\_classification\\_mc\\_and\\_vi.ipynb](https://github.com/tensorchiefs/dl_course_2022/blob/master/notebooks/20_cifar10_classification_mc_and_vi.ipynb)

Run the Variational Inference without trying to understand it – we will care later!

# Looking at the predictive distribution!

Classical CNN (no Bayes)

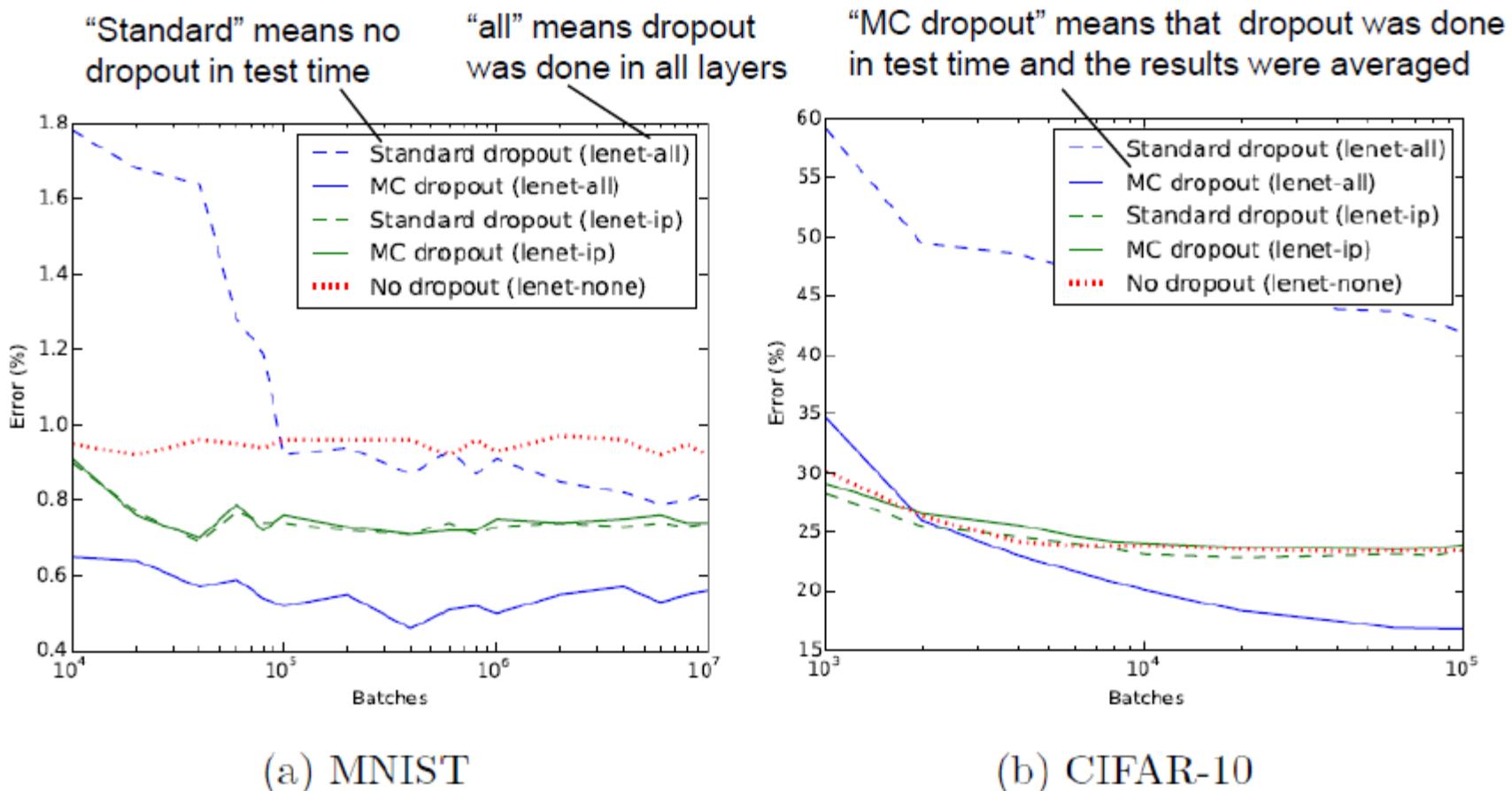


Gaussian VI

MC dropout

# MC Dropout increases prediction performance

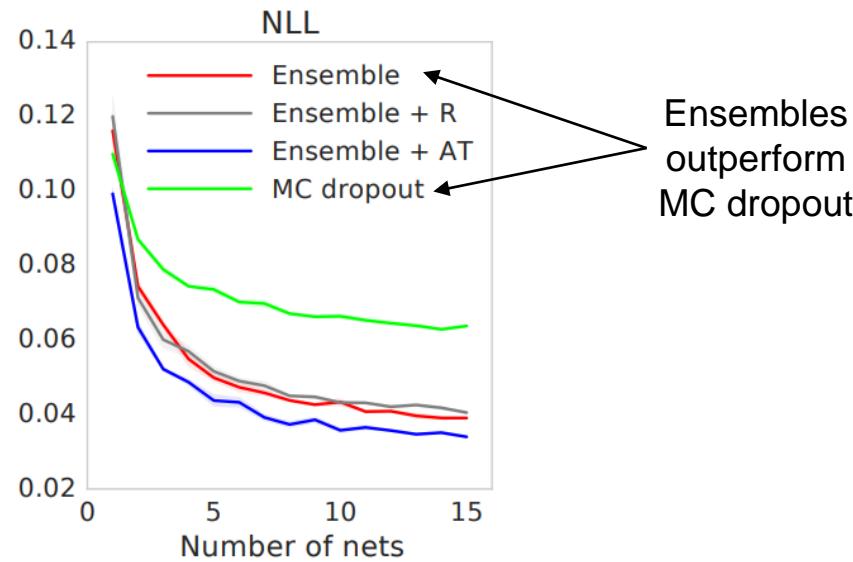
MC dropout is equivalent to performing *several* stochastic forward passes through the network and averaging the results. By doing so Yarin Gal was able to **outperform state of the art error rates** on MNIST and CIFAR-10 without changing the architecture of the used CNNs or anything else beside dropout.



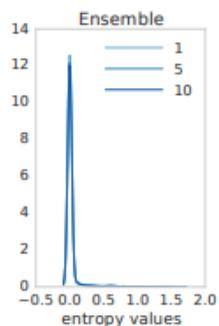
# How does MC dropout compare with deep ensembles?

- For MC dropout we only need to compare one NN with as many parameters as a classical NN. We then average different MC dropout predictions
- For deep ensembles we need to train several NNs (typical 3 to 5) with different random initialization. We then average the predictions of these NNs
- Deep ensembles are computationally more costly but provide typically better prediction performance (and also better uncertainty measures) than MC dropout

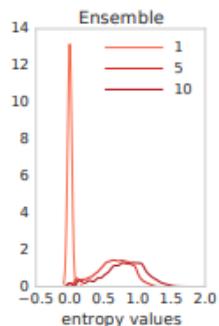
MNIST classification



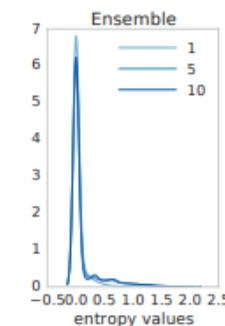
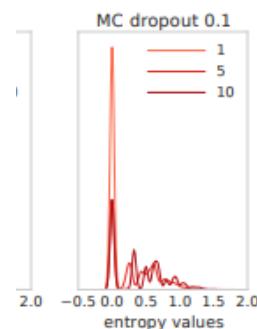
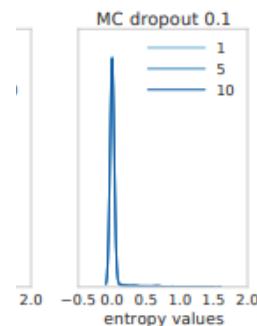
# Compare MC dropout and deep ensembles uncertainties



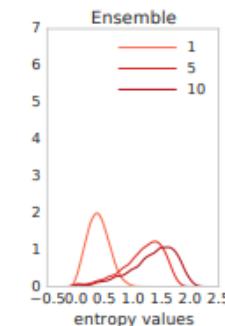
The novel class uncertainty is better in ensemble models than in MC Dropout models



(a) MNIST-NotMNIST



The novel class uncertainty is better in ensemble models than in MC Dropout models



(b) SVHN-CIFAR10

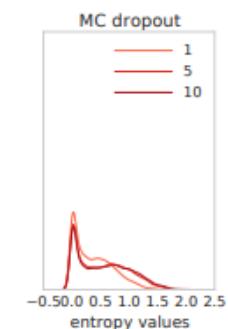
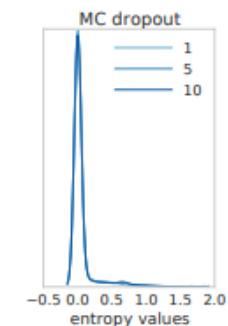
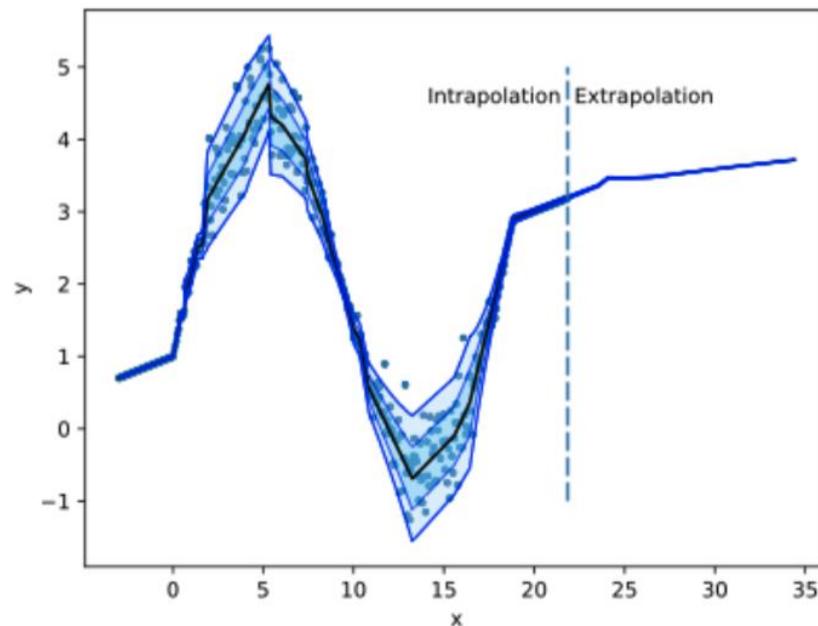
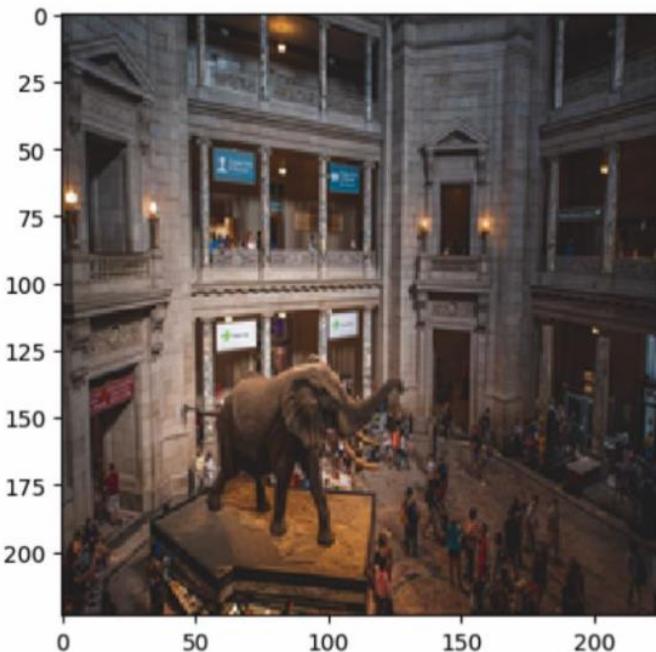


Figure 3: Histogram of the predictive entropy on test examples from known classes (top row) and unknown classes (bottom row), as we vary ensemble size  $M$ .

# Bayesian Deep Learning

# Extrapolation

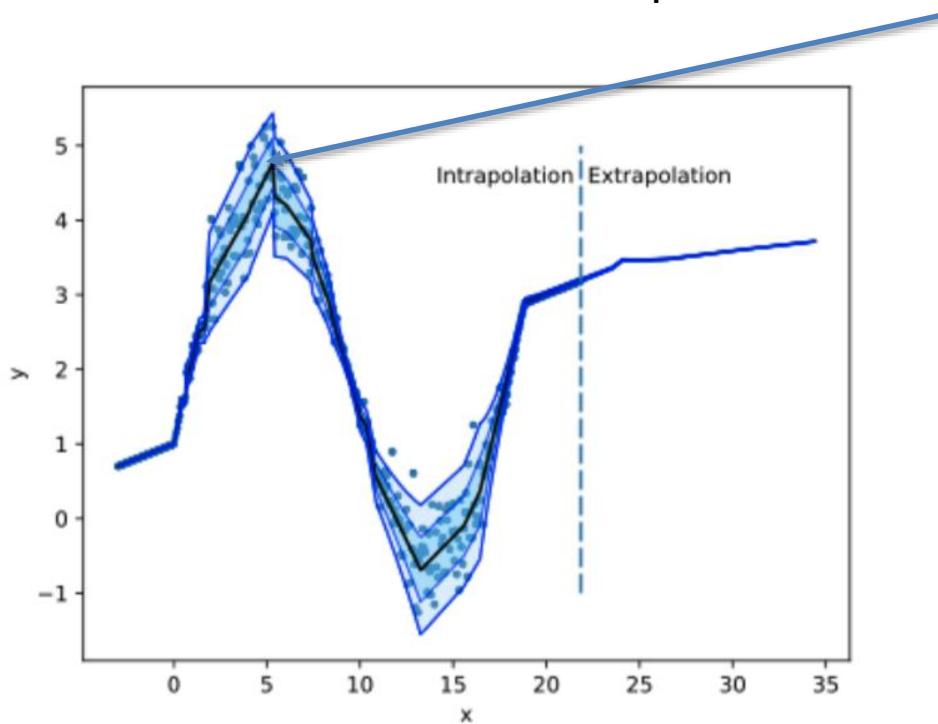


**Figure 7.2 Bad case of DL.** The high performant VGG16-CNN trained on ImageNet data fails to see the elephant in the room. The five highest-ranked class predictions of the objects in the image are horse\_cart, shopping\_cart, palace, streetcar, gondola; the elephant is not found! This image is an extrapolation of the training set. In the regression problem on the right side of the dashed vertical line, there's zero uncertainty in the regions where there's no data (extrapolation).

# Aleatory vs. Epistemic Uncertainty

Much spread in data, aleatoryic (from “[Alea Acta est](#)”)

Die Würfel sind gefallen!



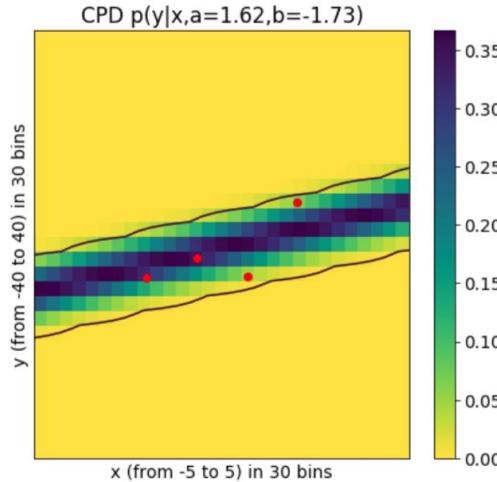
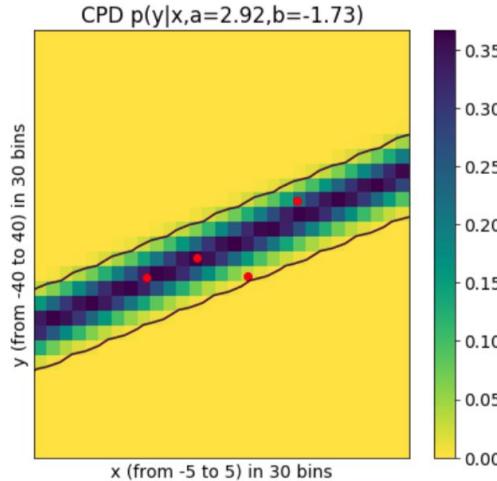
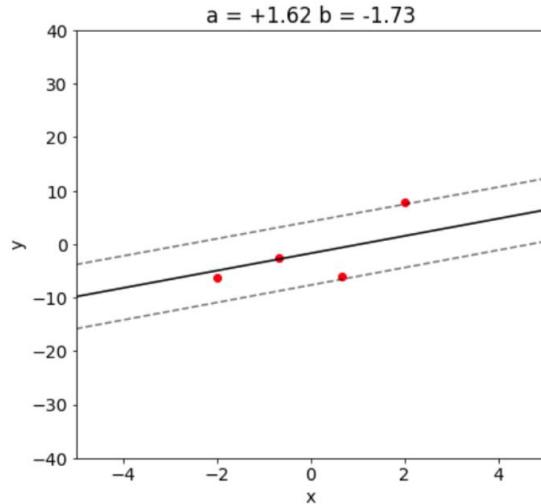
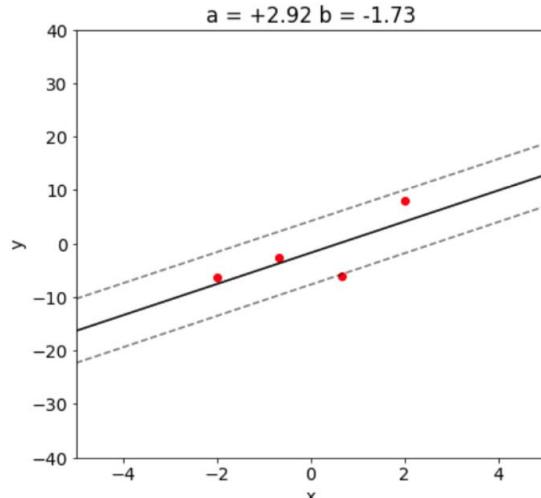
We understand a dice perfectly (no epistemic uncertainty), but still there is a aleatoric uncertainty about the the number that will show up next when rolling the dice.

- *Aleatoric* uncertainty is due to the uncertainty, that is inherent in the data.
- The uncertainty when leaving the ‘known ground’ is called *epistemic* uncertainty.

We use Bayesian methods to model the epistemic model uncertainty, that is due to the uncertainty of the estimated model parameters, called weights in NN models.

# 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.

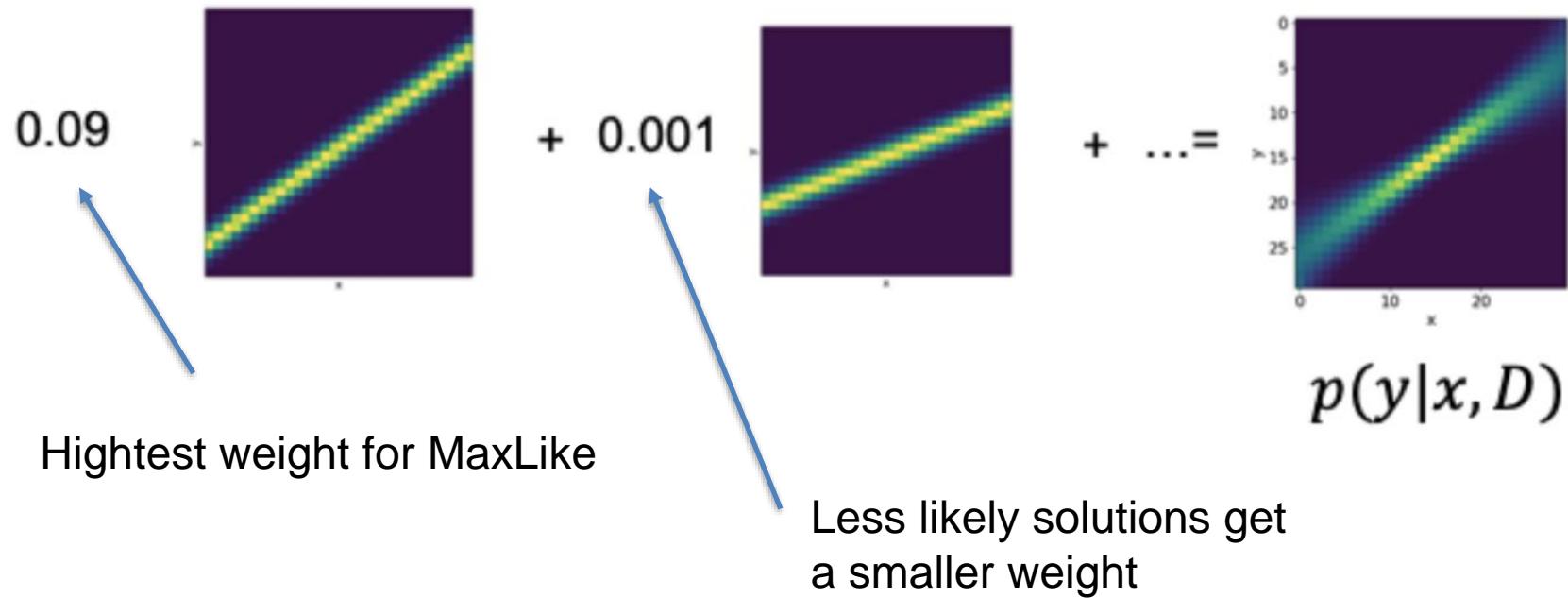


MaxLike Solution

A bit off the MaxLike Solution

# Combining different fits

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

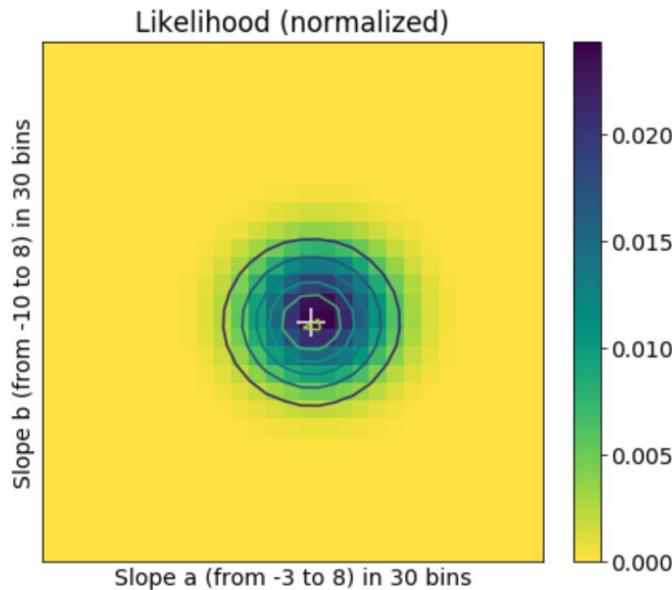


Question: How to get the weight?

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

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

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



$$p_{\text{norm}}(D|(a,b)) = \frac{p(D|(a,b))}{\sum_w p(D|(a,b))}$$

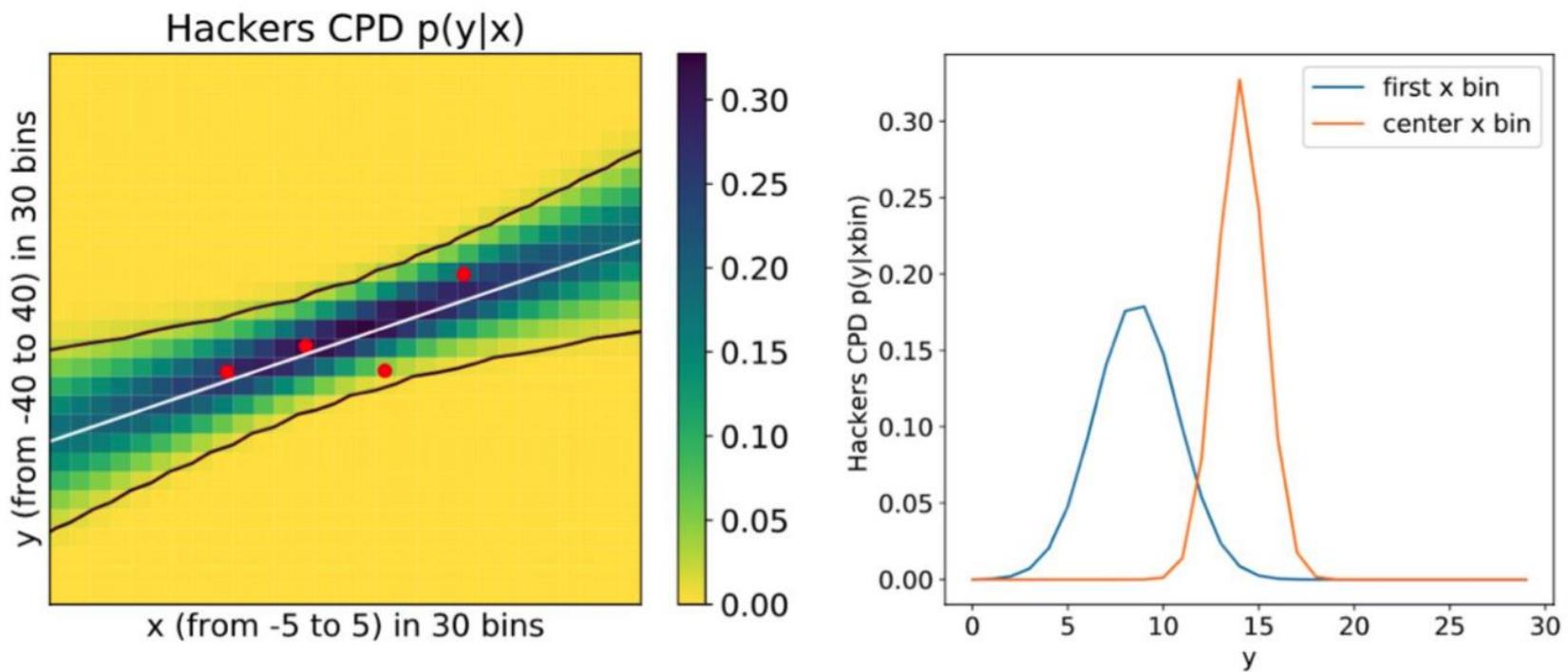
$$p(y|x, D) = \sum_a \sum_b p(y|x, (a, b)) \cdot p_{\text{norm}}(D|(a, b))$$

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](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.

# What have we done? (Comparison with MaxLike)

- Bayes the hacker's way (averaging over many)

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

- For normalization, we calculated likelihood

$$p_{\text{norm}}(D|(a, b)) = \frac{p(D|(a, b))}{\sum_{(a,b)} p(D|(a, b))}$$

- at nbins<sup>2</sup> positions (here 900)

- Maximum likelihood (just take the single best)

- Instead of just using the max likelihood method

$$\hat{a}, \hat{b} = \text{argmax}(p_{\text{norm}}(D|(a, b))) = \text{argmax}(p(D|(a, b)))$$

- The sum reduces to a single term

$$p(y|x, D) = p(y|x, (\hat{a}, \hat{b}))$$

# What have we done? (Abstraction with $(a,b) \rightarrow w$ )

- We calculated
  - $p(y|x, D) = \sum_a \sum_b p(y|x, (a, b)) \cdot p_{\text{norm}}(D|(a, b))$
- Or with the weights  $w$ 
  - $p(y|x, D) = \sum_w p(y|x, w) \cdot p_{\text{norm}}(D|w)$
  - $p(y|x, D) = \int_w p(y|x, w) \cdot p_{\text{norm}}(D|w) dw$  for continuous weights
- The weights now have a distribution, they are random variables as well.
- We choose
  - $p_{\text{norm}}(D|w) = \frac{p(D|w)}{\sum_w p(D|w)}$
  - We treated the weights as random variables
  - It turns out that we were not the first to do so...

# 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  $D$

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

$p(\theta D)$	posterior	see next slide
$p(D \theta)$	likelihood	our good old friend
$p(\theta)$	prior	see next slide
$p(D)$	evidence	just a normalization constant

# Revisiting the Hackers' example

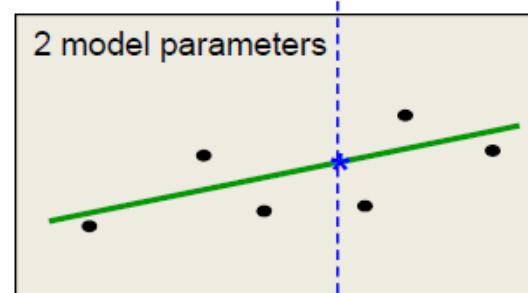
- Hacker's way
  - $p(y|x, D) = \sum p(y|x, w) \cdot p_{\text{norm}}(D|w)$ 
    - $p_{\text{norm}}(D|w) = \frac{p(D|w)}{\sum_w p(D|w)}$
- Bayes
  - $p(y|x, D) = \sum_w p(y|x, w) \cdot p_{\text{norm}}(w|D)$
  - $p(w|D) = \frac{p(D|w)p(w)}{p(D)} = \frac{p(D|w)p(w)}{\sum_w p(D|w)p(w)}$

→ 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!



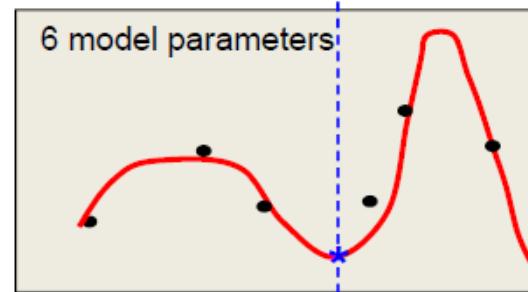
## Bayesian's strategy:

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**.

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



Models with significant posterior probability

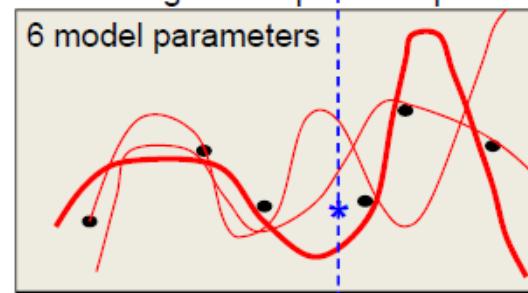


Image credits: Hinton coursera course

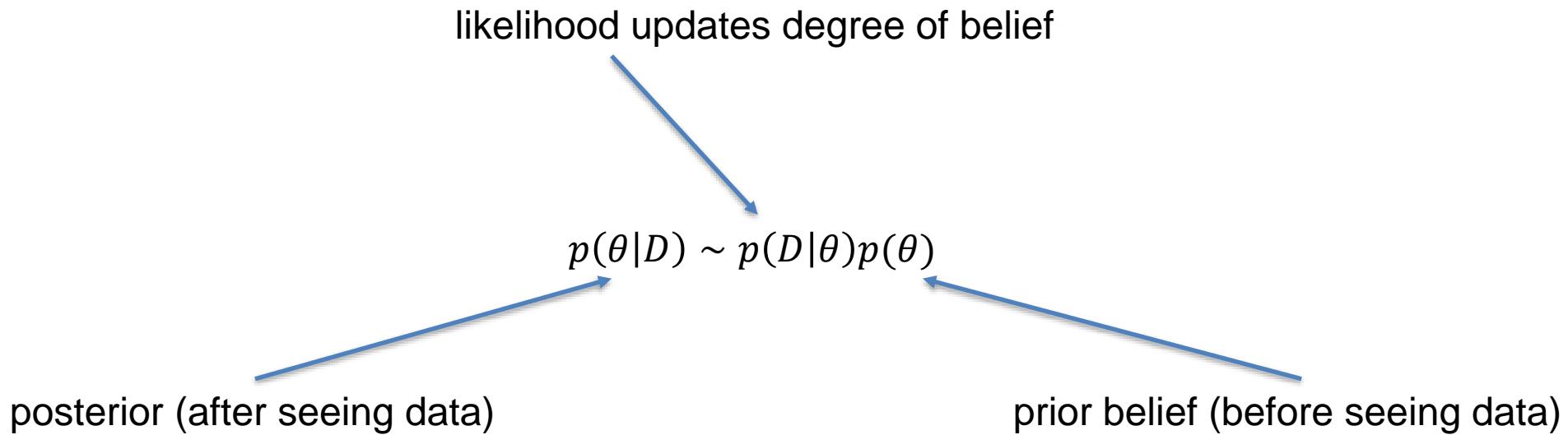
## The Bayesian Mantra (say it loud)

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

“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  $\theta$   
Geben Sie hier eine Formel ein.
- The Bayes formula is seen as an update of our belief in the light of data



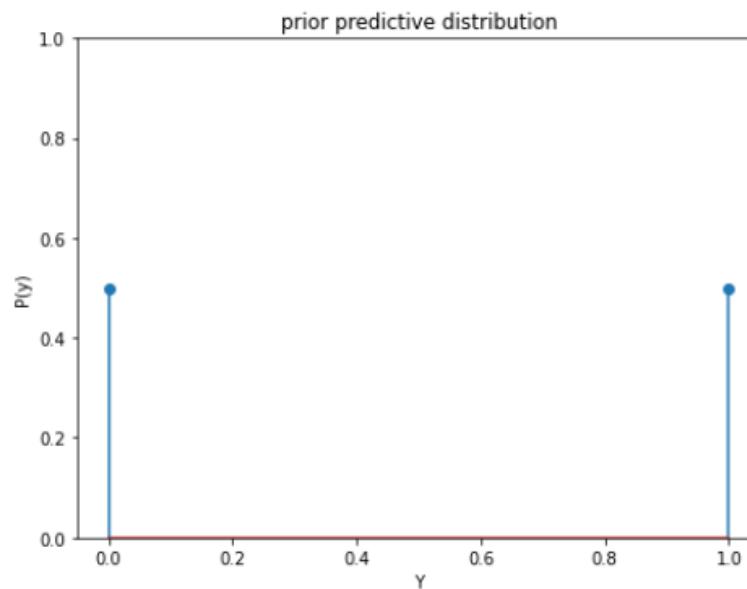
# Example Coin Toss

Head or tail?



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

For a fair coin  $\theta = 0.5$

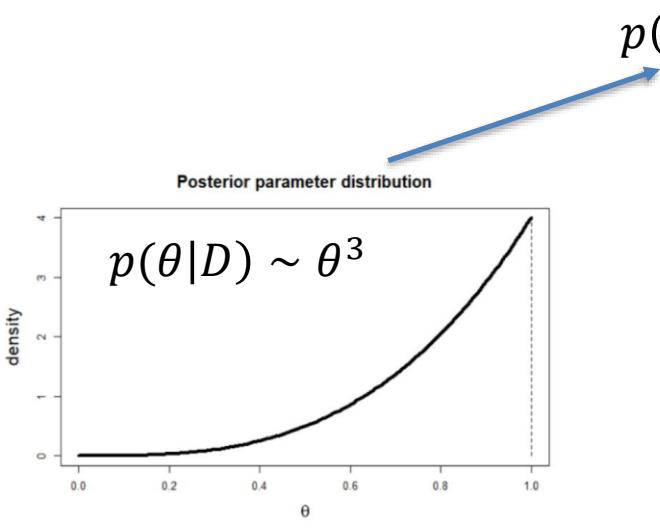


But how to know if a coin is fair?

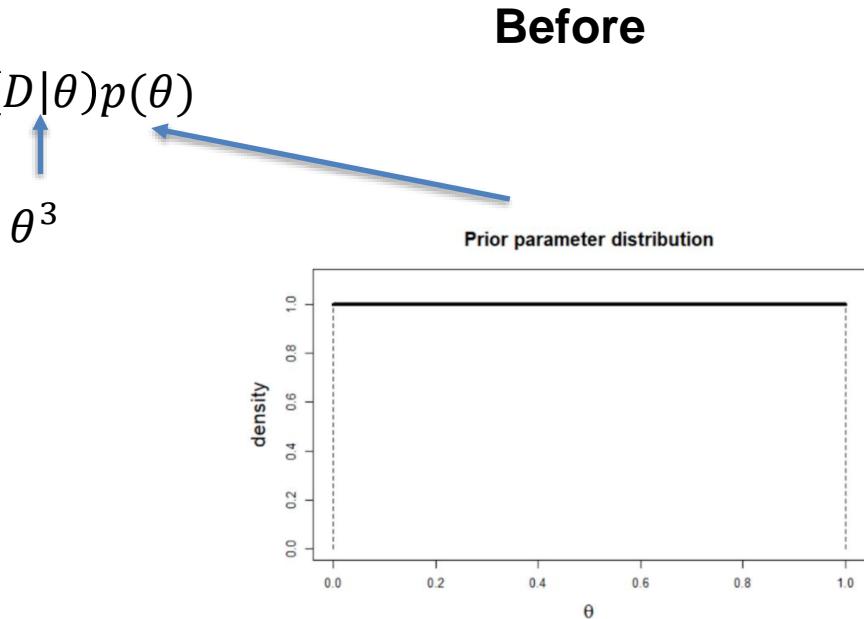
# Analyzing a Coin Toss Experiment

- We do an experiment and observe 3 times head  $\rightarrow D = \text{'3 heads'}$
- $\theta$  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 changes to head

After



Before



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

# Posterior Predictive Distribution

Posterior predictive distribution via marginalization over  $\theta$

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

For unconditional (no x)

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

Prob. of certain result given  $\theta$                           Posterior prob. of certain value  $\theta$

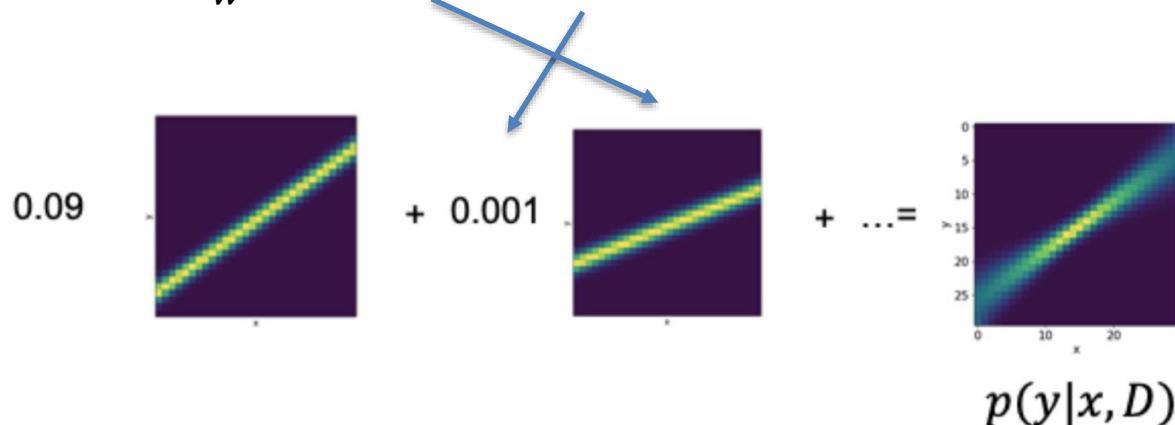
For coin interested in  $p(y = 1|D)$  prob. for head

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

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

# Summary

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

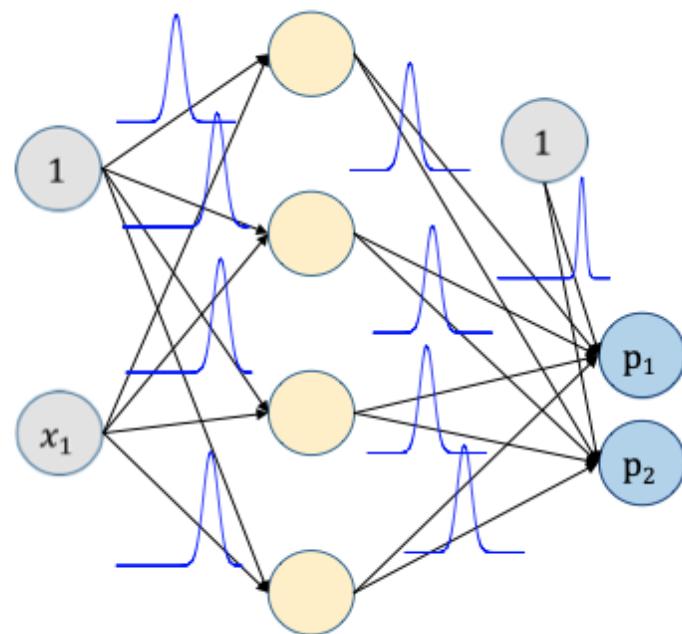
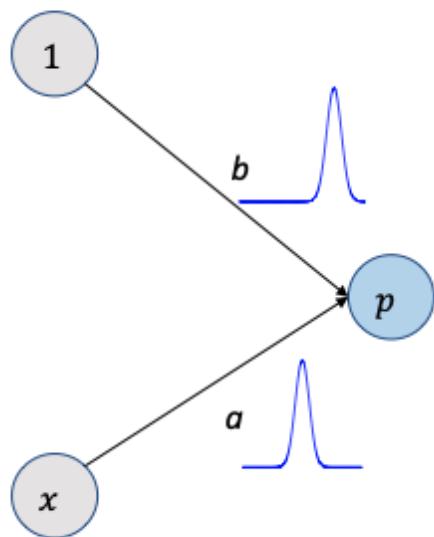


- Not just a single solution
  - “Marginalizing instead of optimization”
- Bayes usually achieves
  - Better uncertainty estimates
  - Better predictive performance

# 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)
- It's possible to calculate  $p(D|\theta)p(\theta)$  for a few values fast
- But calculate  $\sum_{\theta} p(D|\theta)p(\theta)$  is impossible for high-dimension  $\theta$ 
  - We need to calculate 10 values per dimension, then  $10^{\text{dimensions}}$  evals

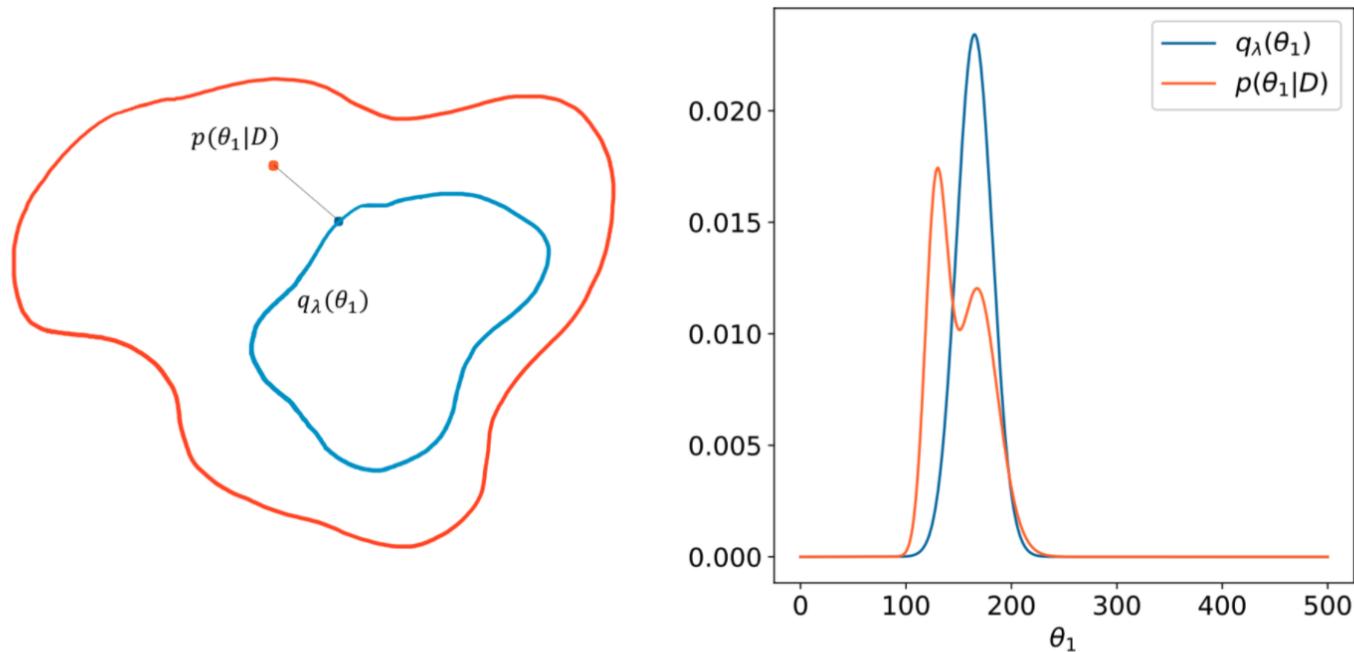
## Approximations

- MCMC (only for very small NN)
  - Sample from  $p(\theta|D)$  with knowledge of  $p(D|\theta')p(\theta')$  /  $p(D|\theta'')p(\theta'')$
- Variational Inference VI
  - Replace  $p(w|D)$  with an approximation e.g.  $N(\mu, \sigma)$  and determine  $\mu, \sigma$
- MC-Dropout
  - Dropout also during predictions (magically) samples from a posterior

# Variational Inference

# The principle of VI

- Replace  $p(\theta|D)$  with  $q_\lambda(\theta)$  (Variational Ansatz)
- 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

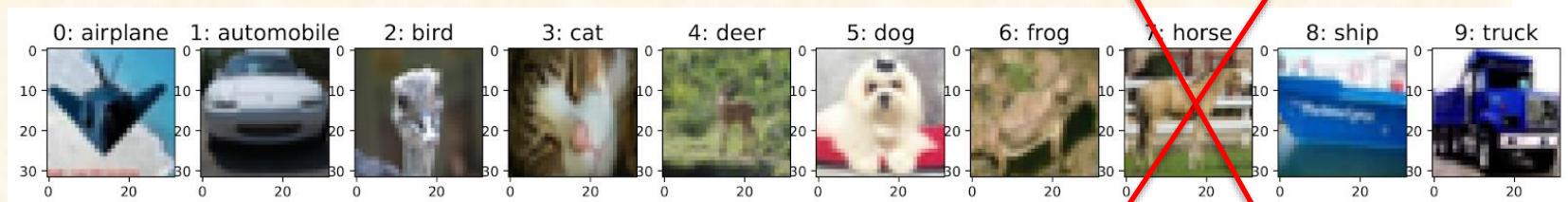
From documentation (Convolution2DFlipout)

When doing minibatch stochastic optimization, make sure to scale this loss such that it is applied just once per epoch (e.g. if  $\text{kl}$  is the sum of losses for each element of the batch, you should pass  $\text{kl} / \text{num\_examples\_per\_epoch}$  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



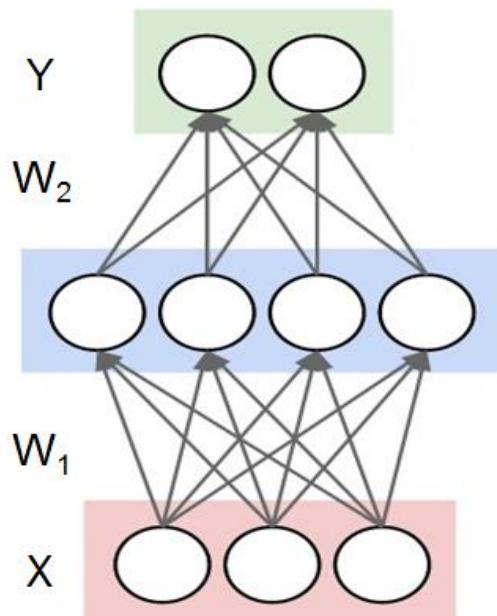
Train a CNN with only 9 of the 10 classes and investigate if the uncertainties are different when predicting images from known or unknown classes.

[https://github.com/tensorchiefs/dl\\_course\\_2021/blob/master/notebooks/20\\_cifar10\\_classification\\_mc\\_and\\_vi.ipynb](https://github.com/tensorchiefs/dl_course_2021/blob/master/notebooks/20_cifar10_classification_mc_and_vi.ipynb)

# Comparing non-Bayesian with Bayesian NN

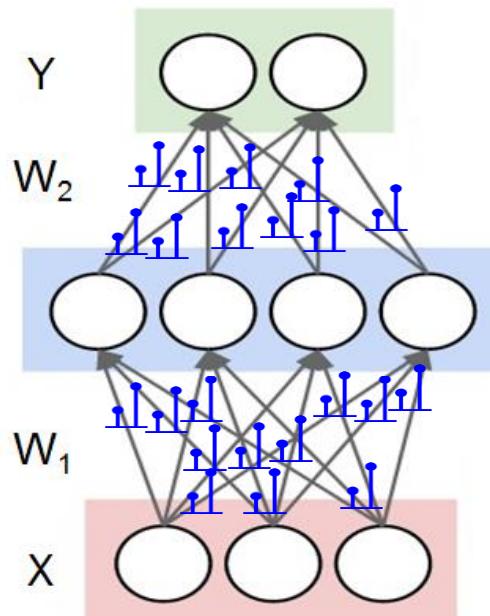
# Non-Bayesian and Bayesian NNs

Non-Bayesian  
NN



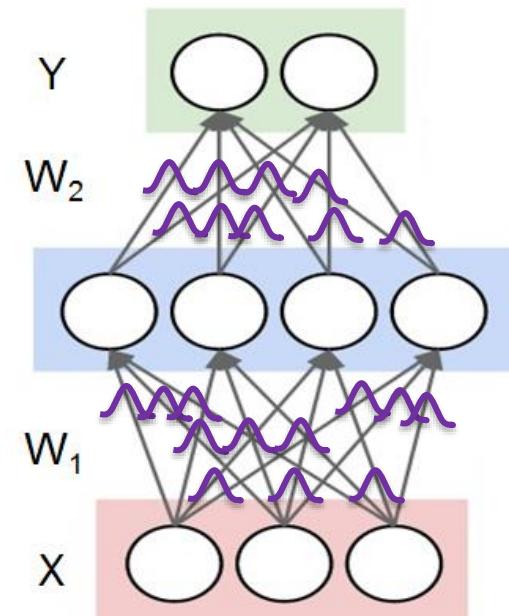
**Weights are fixed**

MC dropout  
Bayesian NN



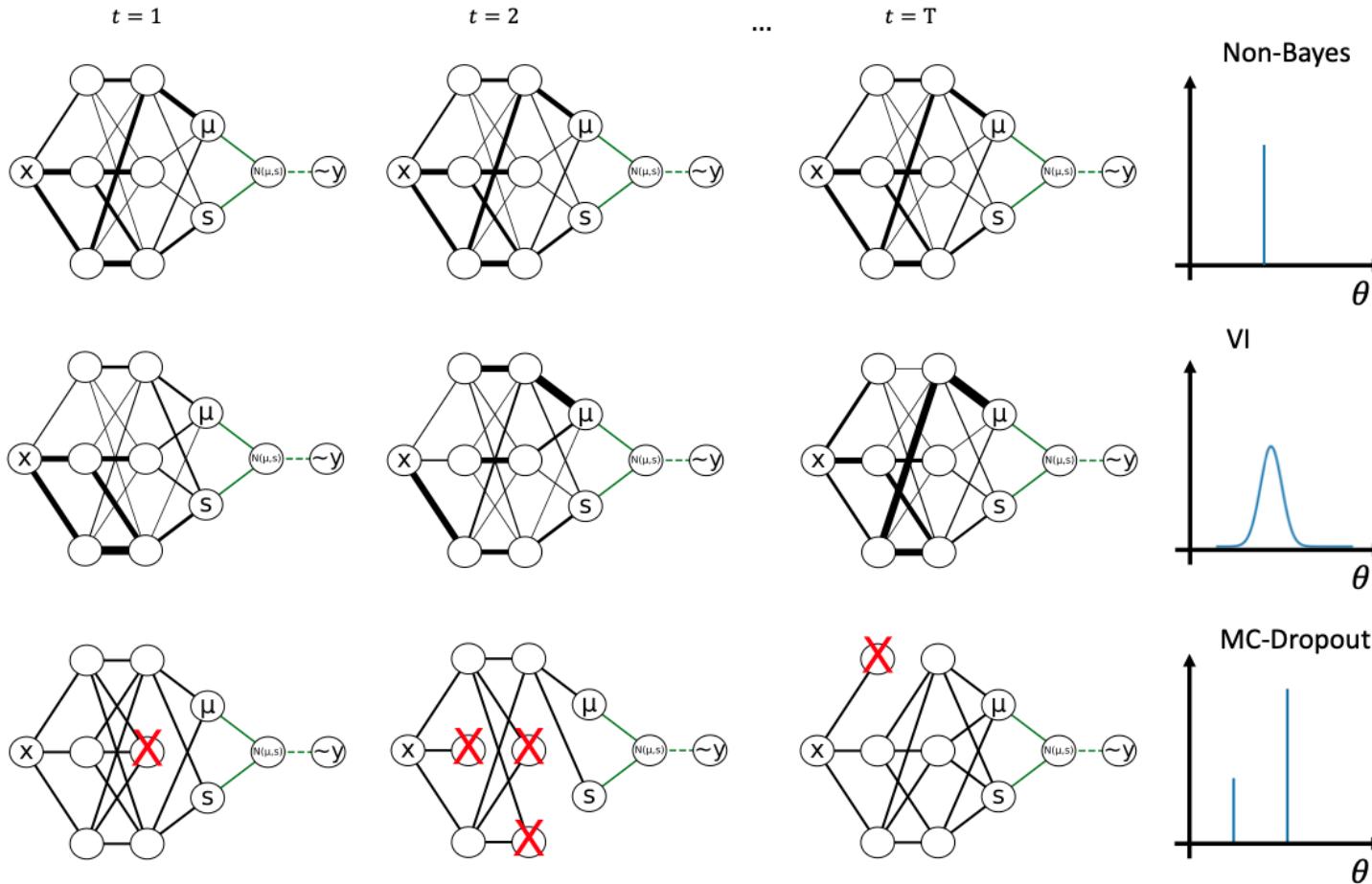
**Weights have  
Bernoulli-kind  
distribution**

VI  
Bayesian NN



**Weights have  
Gaussian  
distribution**

# Comparing different Network types



<https://www.youtube.com/watch?v=mQrUcUoT2k4>

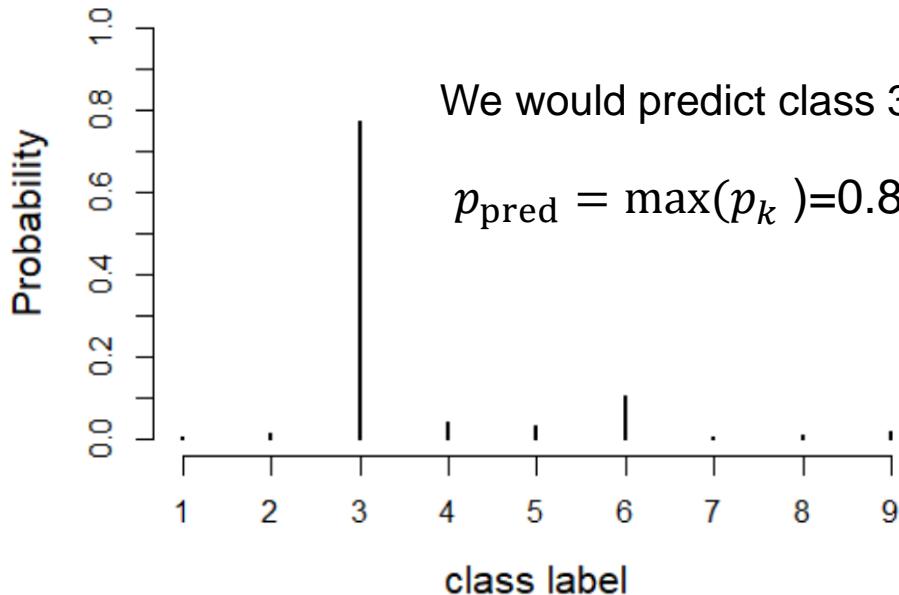
A Non-Bayesian NN learns one set of weights: the same input same output  
A Bayesian NN learns distribution of weights: same input different outputs

# Uncertainty measures in classification

# Uncertainty in non-Bayesian classification

Multinomial CPD

$$\text{MN}(p_1(x, w), p_2(x, w), \dots, p_9(x, w))$$



In a non-Bayesian NN we make for each input  $x$  ONE CPD:

<b>Image <math>x</math></b>
$\text{MN}(p_1(x, w), \dots, p_9(x, w))$

**Uncertainty measures capturing the aleatoric uncertainty :**

Negative log-Likelihood:  $NLL = -\log(p_{\text{pred}})$

Entropy:  $H = -\sum_{k=1}^9 p_k \cdot \log(p_k)$

# Uncertainty in Bayesian classification

In a Bayesian NN we sample T-times from the weight distributions and get each time a slightly different multinomial CPD

<b>predict_no</b>	<b>Image x</b>
1	MN( $p_1(x, w_1), \dots, p_9(x, w_1)$ )
2	MN( $p_1(x, w_2), \dots, p_9(x, w_2)$ )
...	
T	MN( $p_1(x, w_T), \dots, p_9(x, w_T)$ )

For each class  $k$  ( $k \in \{1, 2, \dots, 9\}$ ) we determine the mean probability:  $p_k^* = \frac{1}{T} \sum_{i=1}^T p_{ki}$

We predict the class with the highest mean probability:  $p_{pred}^* = \max(p_k^*)$

## Uncertainty measures including aleatoric and epistemic contributions:

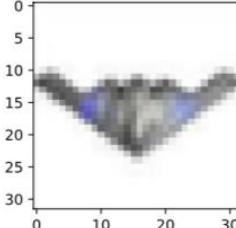
Entropy:  $H^* = - \sum_{k=1}^9 p_k^* \cdot \log(p_k^*)$

Total variance:  $V_{tot}^* = \sum_{k=1}^9 \text{var}(p_k) = \sum_{k=1}^9 \sum_{i=1}^T (p_{kt} - p_k^*)^2$

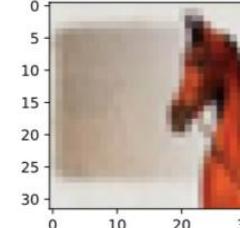
# Looking at the predictive distribution!

Input image

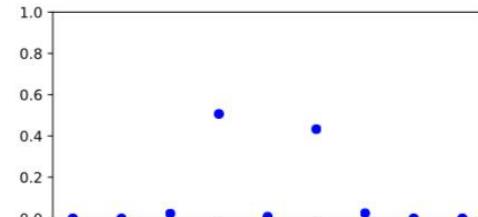
known class



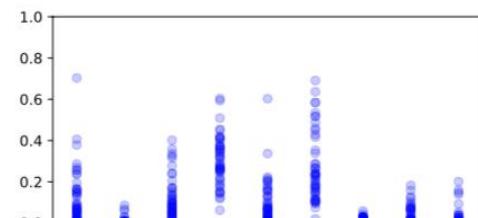
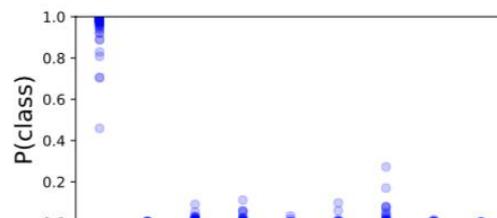
unknown class



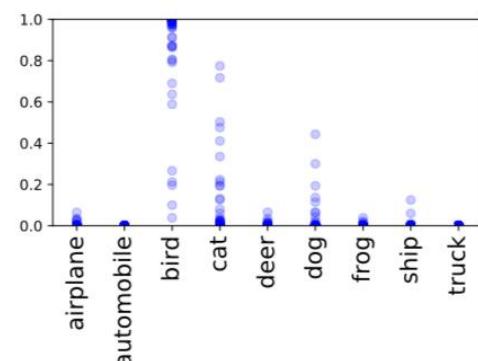
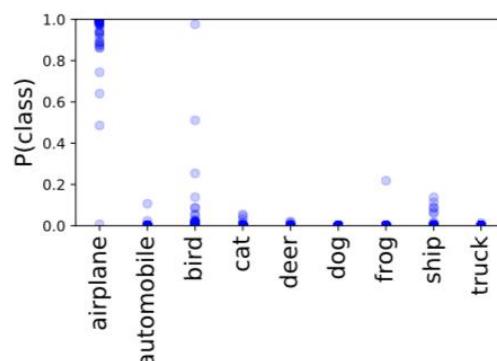
Non-Bayesian CNN



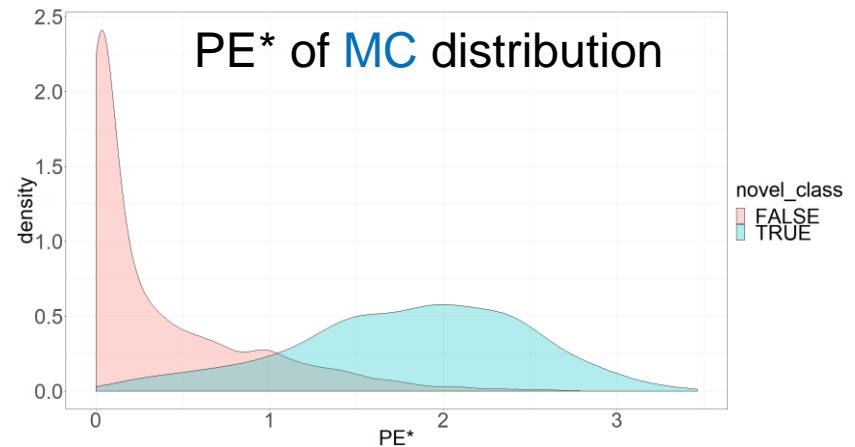
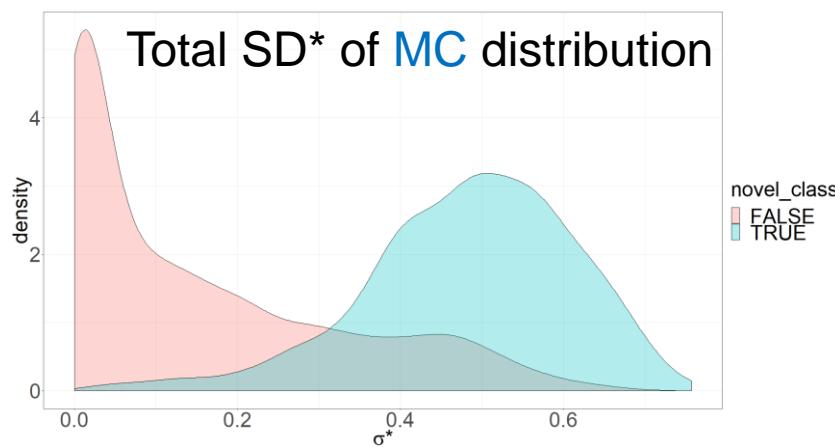
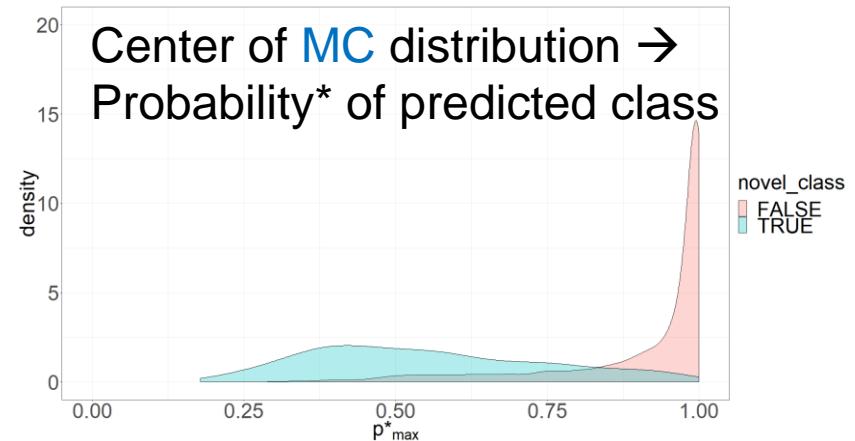
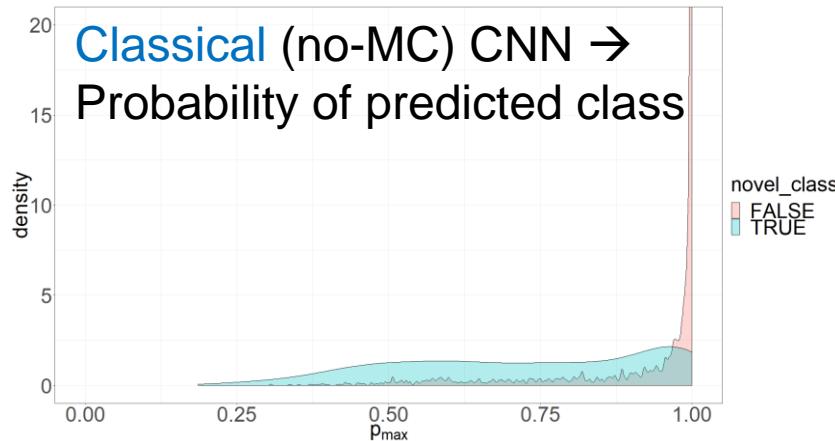
Bayesian CNN via VI



Bayesian CNN via dropout



# Do known/novel classes yield different values for probability and uncertainty measures?

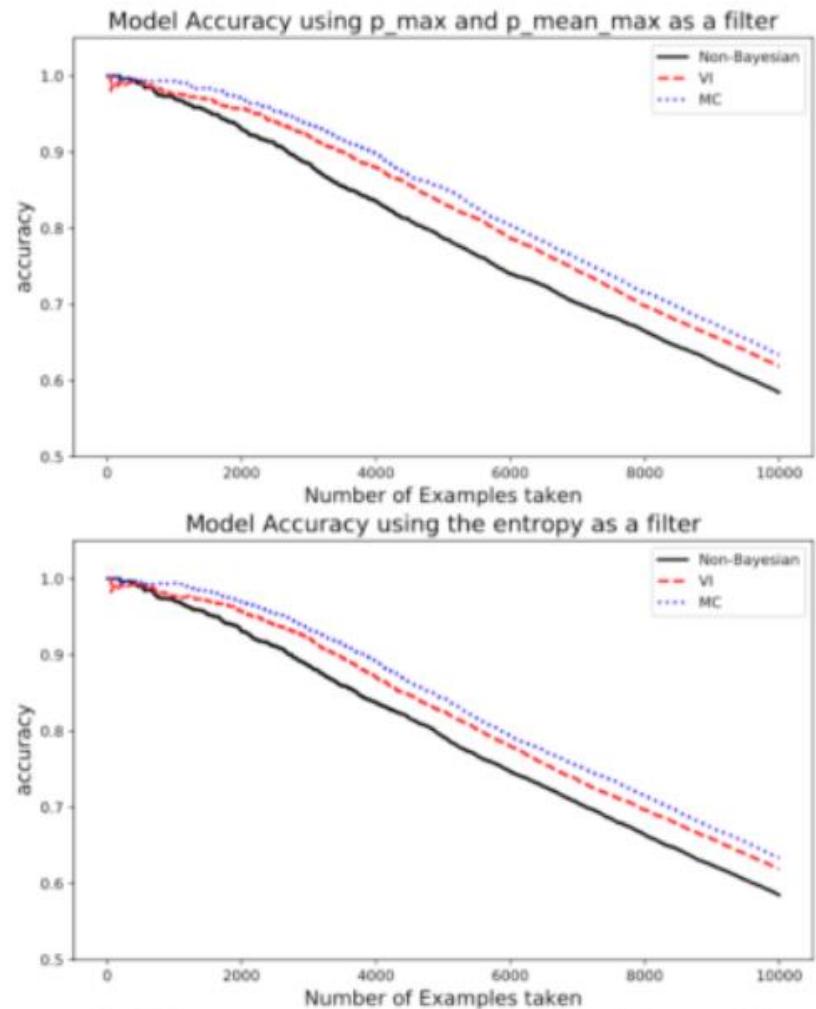
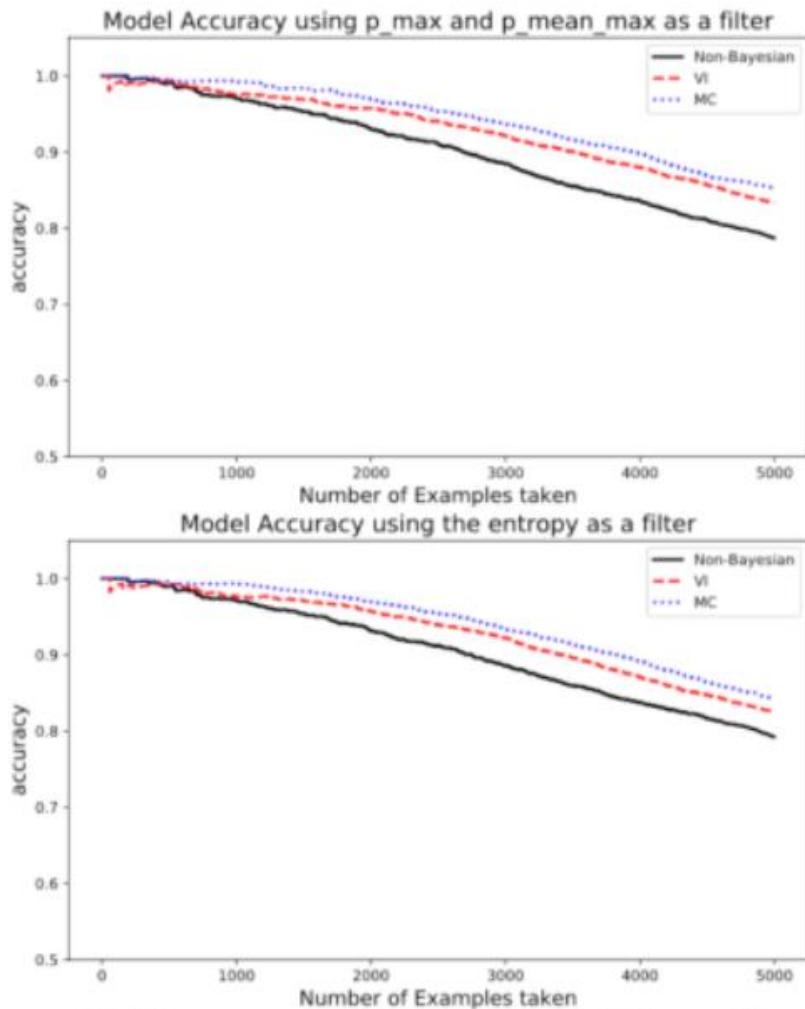


# Filtering experiment based on uncertainty

Goal: Get higher accuracy by filter only predictions which are quite certainly correct

- Each prediction has an attached uncertainty measure
- Sort predictions according to the uncertainty measures
- A set of predictions with very low uncertainties should achieve a high accuracy
- By successively adding predictions with increasing uncertainties should yield predictions sets with decreasing accuracies.

# Filtering experiment to compare uncertainty measures



Uncertainty from non-Bayesian NN is less good in filtering out wrong classifications than uncertainty measures from Bayesian variants of the NN.

# Uncertainty measures in regression

# Uncertainty in non-Bayesian NN

We do predictions for 400 x-values between -10 and 30 yielding for each x a Gaussian CPD.

x1= -10	x2= -9.9	...	x400= 30
$N(\mu_{x_{1,w}}, \sigma_{x_{1,w}})$	$N(\mu_{x_{2,w}}, \sigma_{x_{2,w}})$		$N(\mu_{x_{400,w}}, \sigma_{x_{400,w}})$

**Uncertainty** measures capturing the **aleatoric** uncertainty at  $x$ :

Standard deviation:  $\sigma_x$

$$95\% \text{ PI: } [q_{0.025}; q_{0.975}] = [\mu_x - 1.96 \cdot \sigma_x; \mu_x + 1.96 \cdot \sigma_x]$$

Remark:

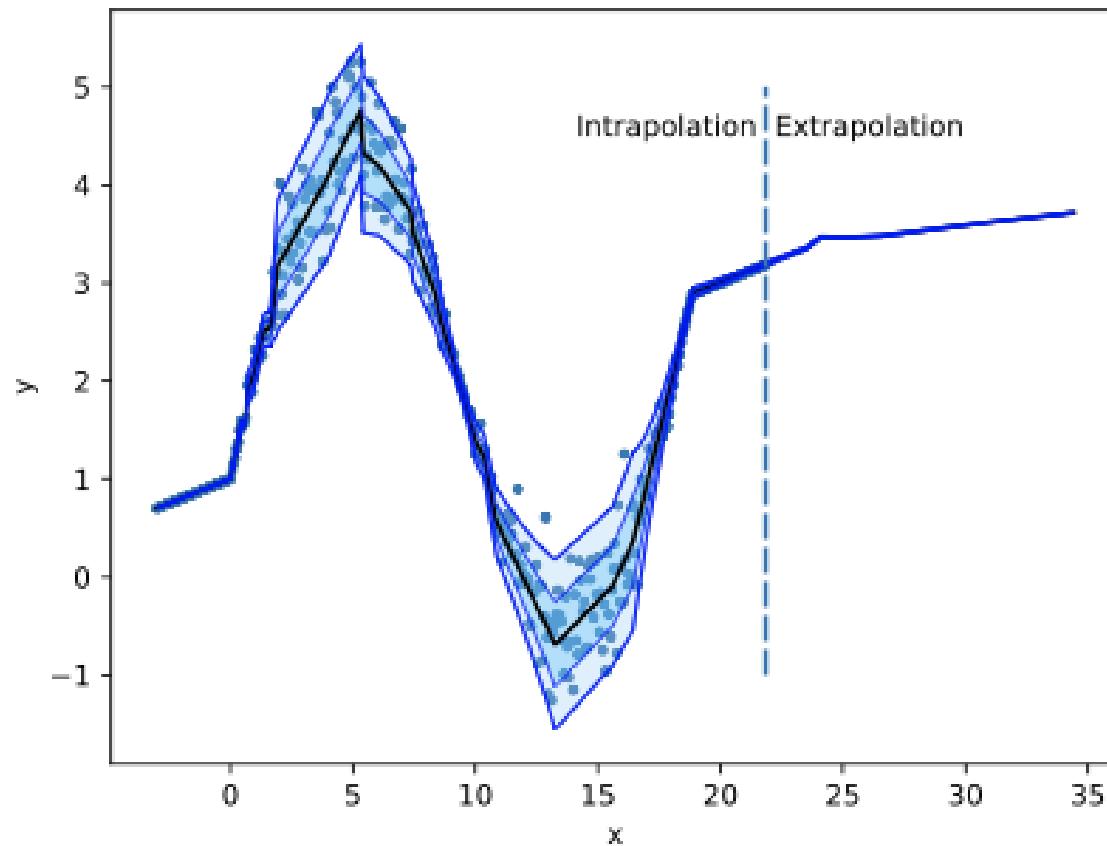
We could also estimate the 95% PI at position x by sampling several times from the CPD and determine the 0.025 and 0.975 quantiles, yielding :

$$95\% \text{ PI: } [q_{0.025}; q_{0.975}]$$

# The problem of non-Bayesian NN

Problem:

A non-Bayesian NN does extrapolation with very small uncertainty



# Uncertainty in Bayesian regression NN

In a Bayesian NN we sample T-times from the weight distributions and get each time a slightly different CPD. In regression the CPD is often Gaussian.

We do predictions for 400 x-values between -10 and 30 yielding in each of the T runs a different Gaussian CPD at each x-position.

<b>predict_no</b>	x1= -10	x2= -9.9	...	x400= 30
1	$N(x_1, w_1, x_1, w_1)$	$N(x_2, w_1, x_2, w_1)$		$N(x_{400}, w_1, x_{400}, w_1)$
2	$N(x_1, w_2, x_1, w_2)$	$N(x_2, w_2, x_2, w_2)$		$N(x_{400}, w_2, x_{400}, w_2)$
...				
T	$N(x_1, w_T, x_1, w_T)$	$N(x_2, w_T, x_2, w_T)$		$N(x_{400}, w_T, x_{400}, w_T)$

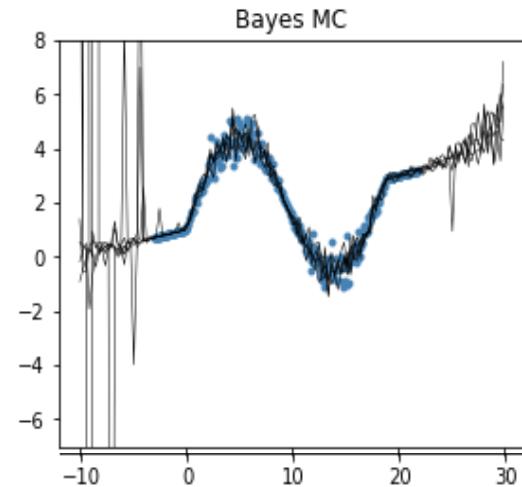
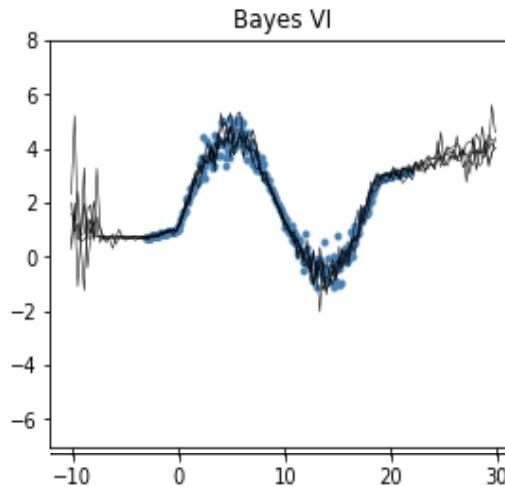
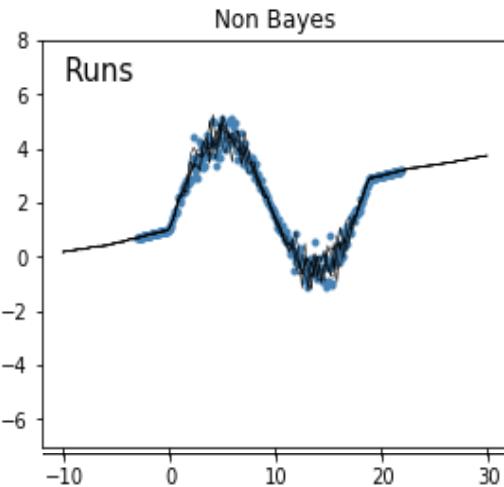
**Uncertainty measures including aleatoric and epistemic contributions:**

To estimate the 95% PI at position x we sample y-values from each of the T CPDs and determine from the samples the 0.025 and 0.975 quantiles, yielding :

95% PI:  $[q_{0.025}; q_{0.975}]$

# Can we see enhanced uncertainty in extrapolation

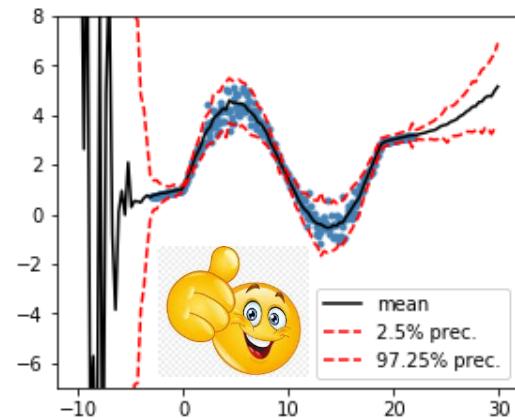
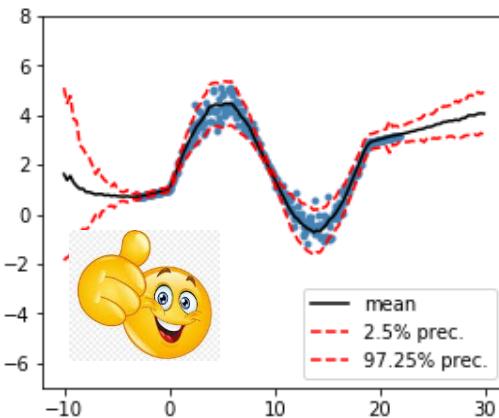
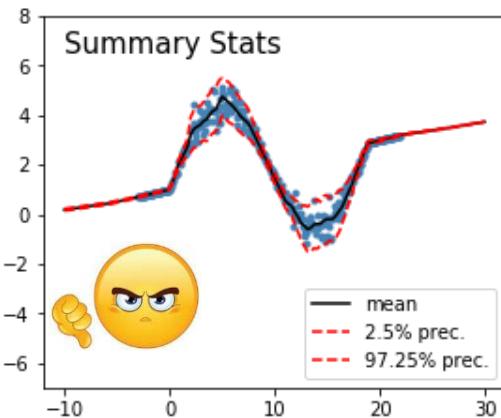
The solid lines show five predicted y-vectors corresponding to 5 CPDs at each x-position.



<https://youtu.be/FO5avm3XT4g>

<https://youtu.be/mQrUcUoT2k4>

<https://youtu.be/0-oyDeR9HrE>



# 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.

# VI in TensorFlow probability

```
model = tf.keras.Sequential([
    tfp.layers.DenseReparameterization(1,input_shape=(None,1)),
    tfp.layers.DenseReparameterization(2),
    tfp.layers.DenseReparameterization(3)
])
```

This builds a network, with hidden layers having:

- 1 Node
- 2 Node
- 3 Node

Question: How many parameters does this network have? Note that biases are included but don't have variational distributions. [3 minutes]

# Solution

```
model = tf.keras.Sequential([
    tfp.layers.DenseReparameterization(1, input_shape=(None,1)),
    tfp.layers.DenseReparameterization(2),
    tfp.layers.DenseReparameterization(3)
])
```

Model: "sequential\_1"

Layer (type)	Output Shape	Param #
dense_reparameterization_3 (None, None, 1)		3
dense_reparameterization_4 (None, None, 2)		6
dense_reparameterization_5 (None, None, 3)		15
<hr/>		
Total params:	24	
Trainable params:	24	
Non-trainable params:	0	