

# Machine Intelligence:: Deep Learning

## Week 7

*Beate Sick, Jonas Brändli, Oliver Dürr*

Ensembling approaches for 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

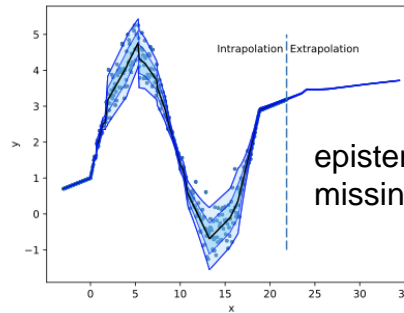
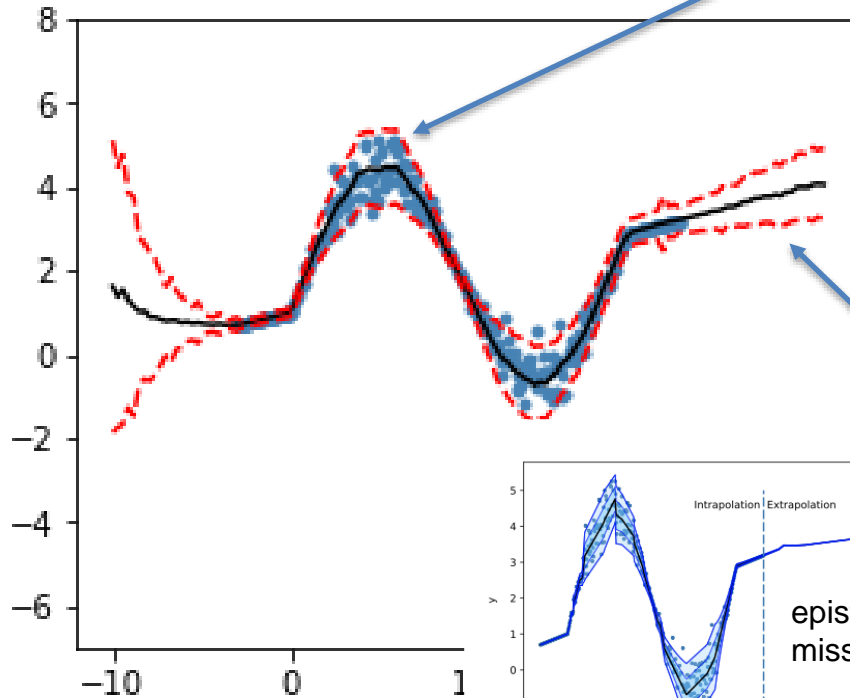
# Aleatoric vs. Epistemic Uncertainty

Much spread in train data → **aleatoric 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.



epistemic uncertainty is missing in classical NN

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



General [VIEW DOCS](#)

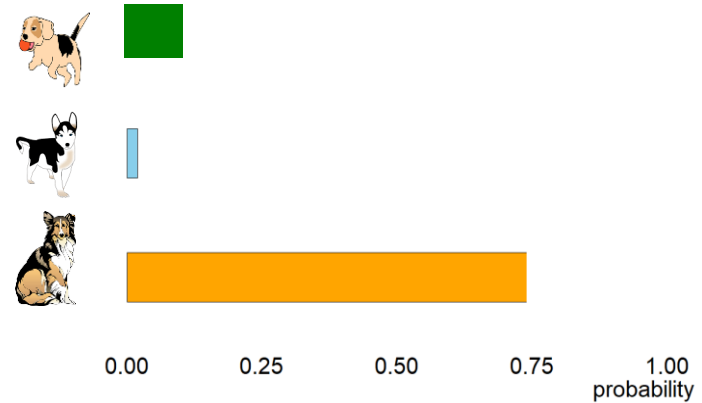
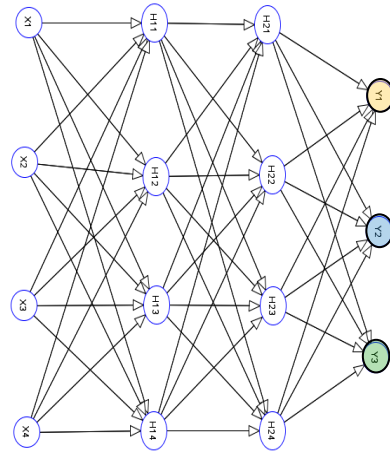
PREDICTED CONCEPT	PROBABILITY
group	0.979
adult	0.977
people	0.976
furniture	0.960
room	0.957
business	0.903
indoors	0.901
man	0.896
seat	0.895

# 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)

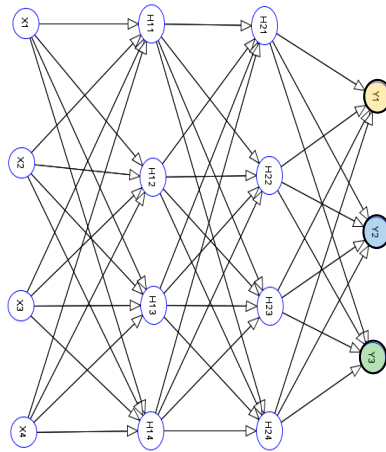
# CNNs have high performance on in-distribution examples



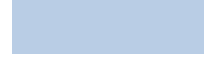
**CNNs yield high accuracy and calibrated probabilities, but...**

# 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?



You might expect:



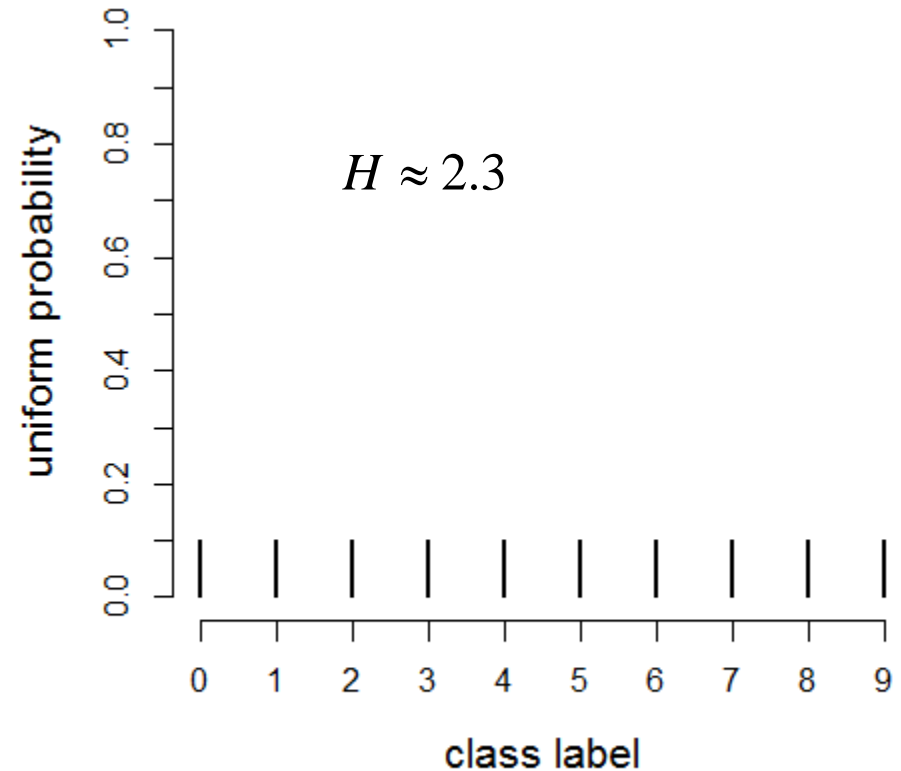
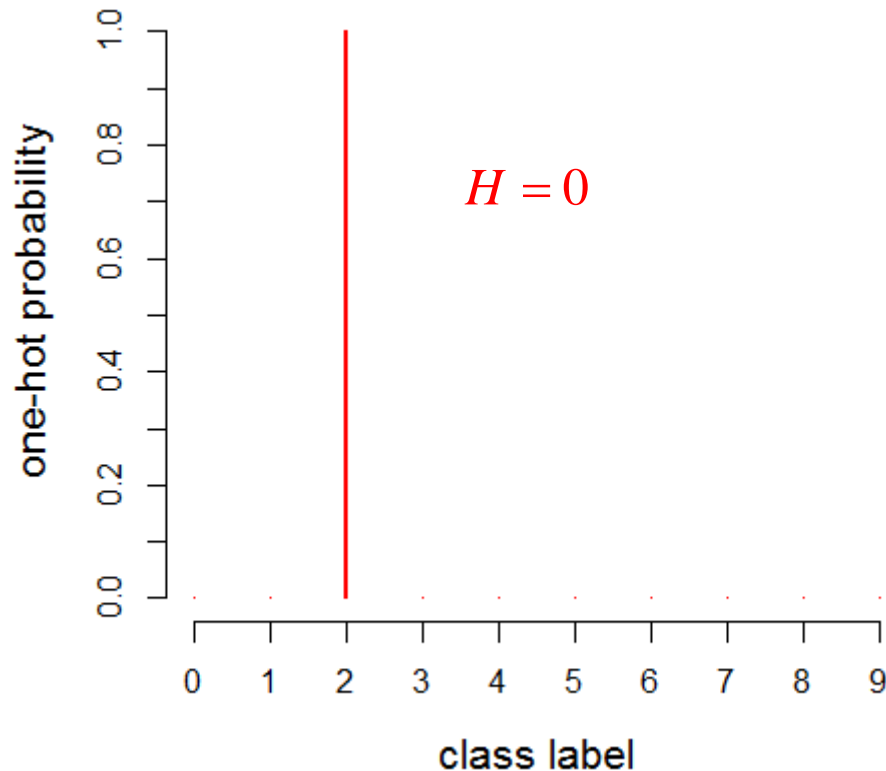
# Recall: Entropy as measure for uncertainty

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

Entropy is a measure for “untidyness” or uncertainty.

If a distribution has only one peak, it is tidy and  $H=0$

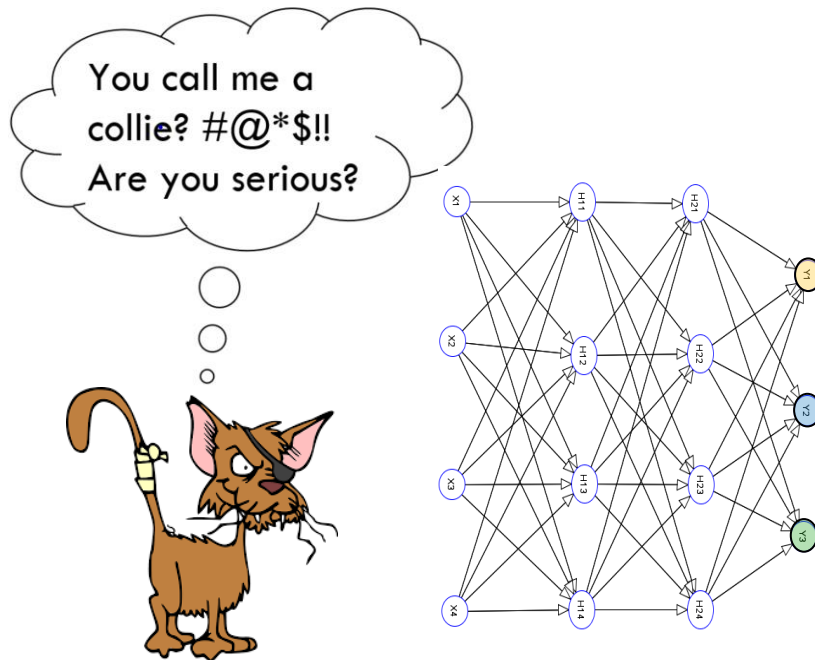
If all outcomes are equally probable it is maximal untidy





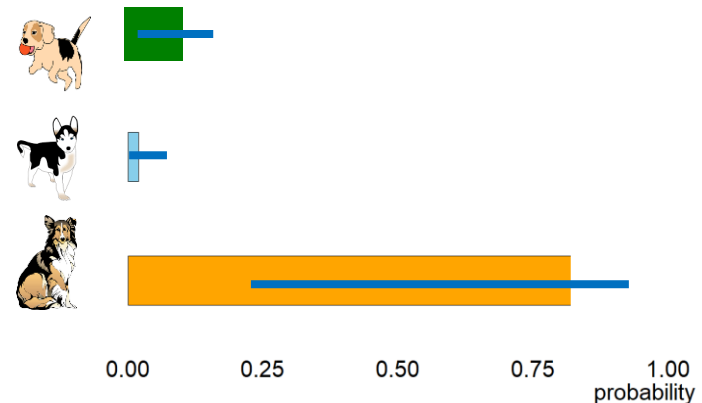
# A classical NN cannot ring the alarm in case of out-of-distribution (OOD) examples

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)

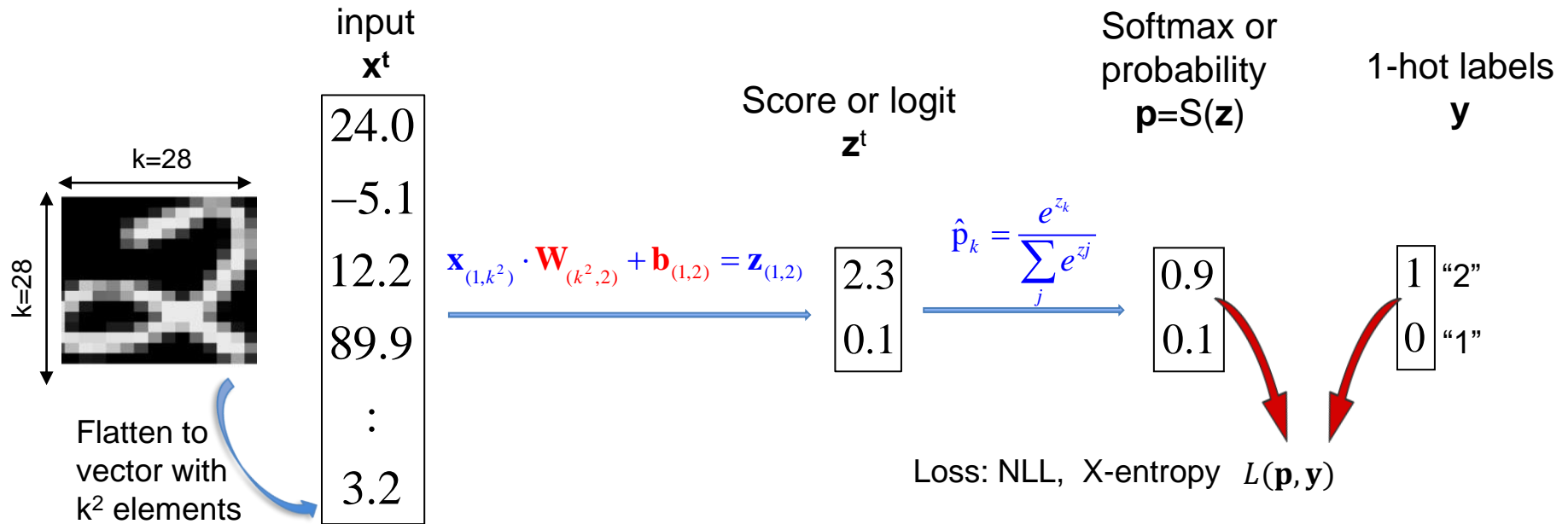


- Current DL Systems bad in out of distribution OOD situations
- Application need at least to detect OOD situations

# Deep ensembles

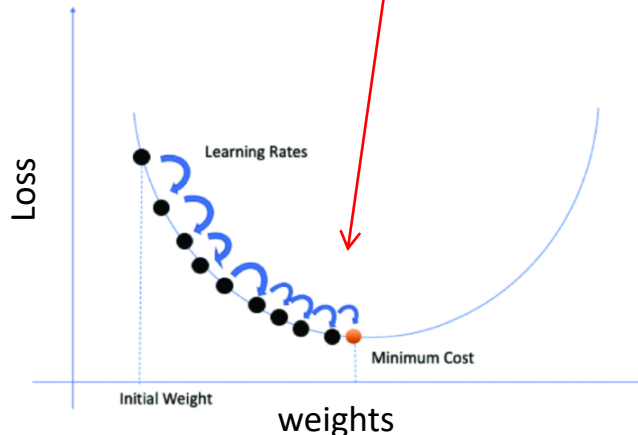
<https://www.tutorialspoint.com/how-can-keras-be-used-to-implement-ensembling-in-python>

# Recall the training of NN models via SGD



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

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



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

Loss = NLL averaged over all images in mini-batch

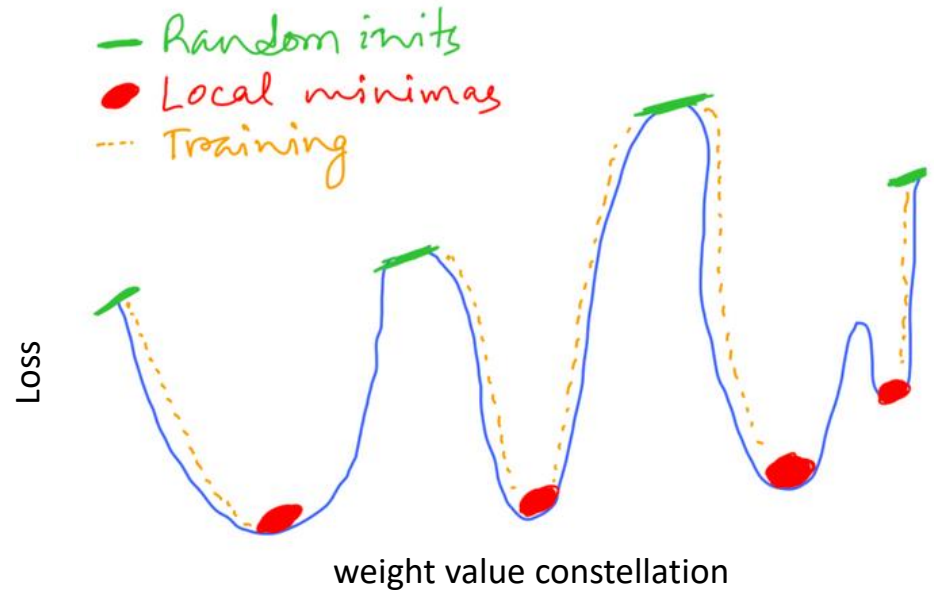
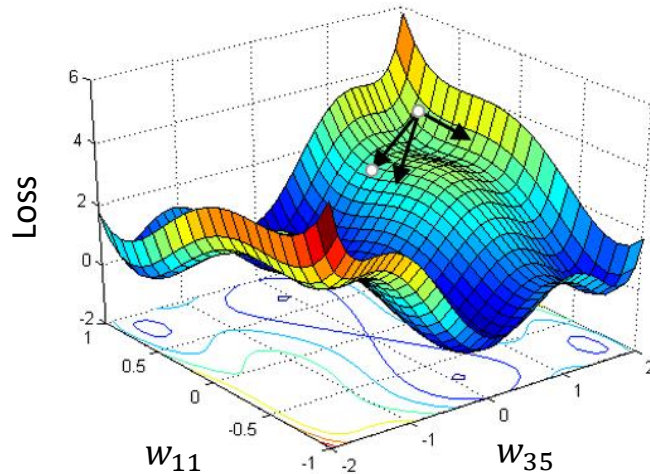
$$\text{NLL} = \frac{1}{N} \sum_i L(\mathbf{p}_i, \mathbf{y}_i)$$

Loss

$C(w_1, w_2)$

# The loss-landscape in DL is usually not convex

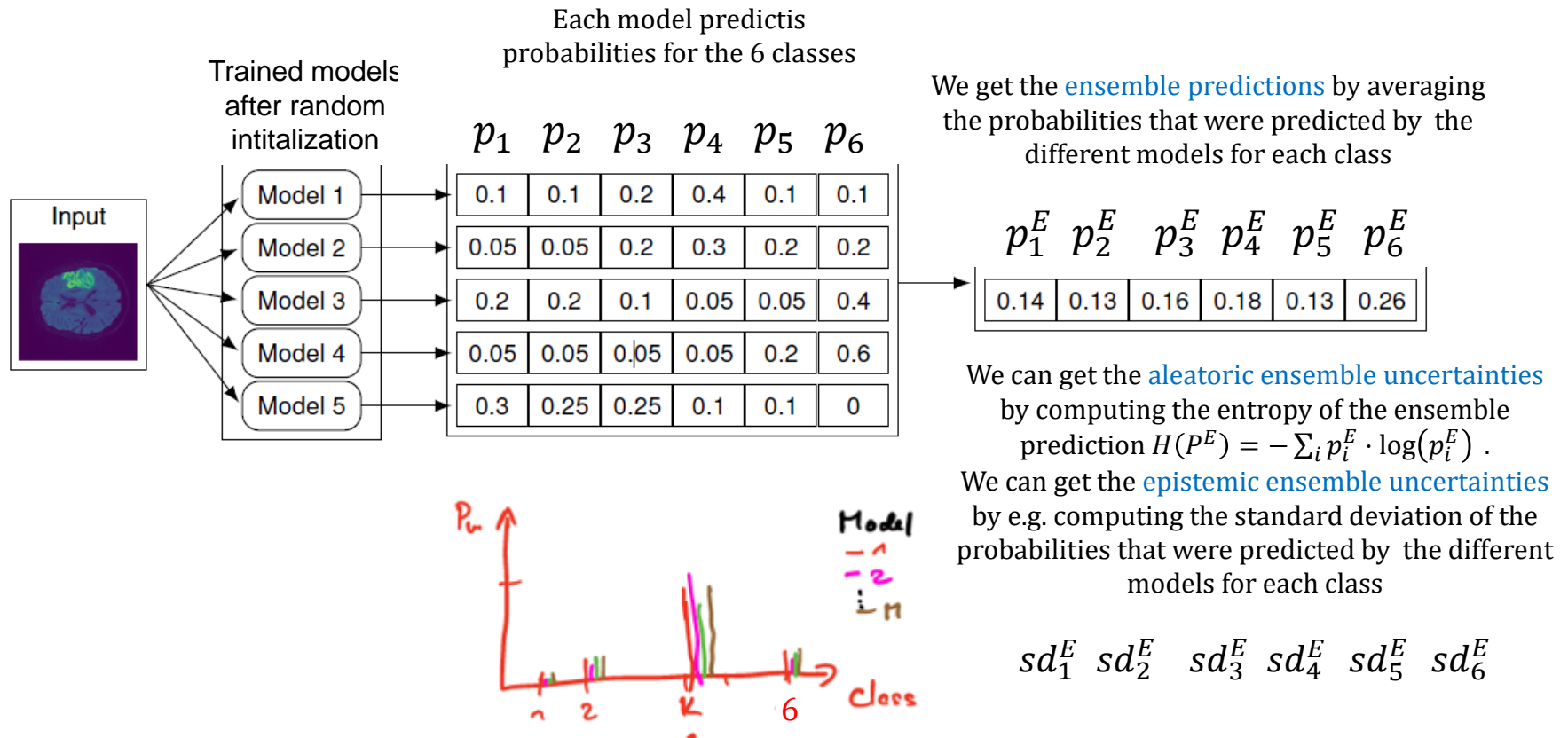
Loss-landscape



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 =  $k$ , based on  $M$  models:

$$P_k^E = \frac{1}{M} \sum_{m=1}^M P_{km} \quad k = 1 \dots 10 \text{ for 10 classes}$$

associated NLL contribution  $l$ :

Model  $m$ :  $l_m = -\log P_{km}$  / Ensemble:  $l^E = -\log P_k^E$

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

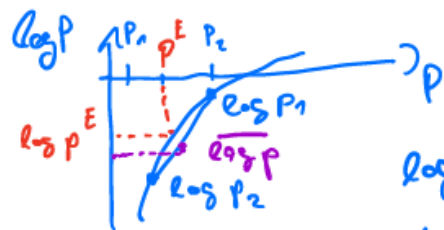
$$\text{with } \bar{l} = \frac{1}{M} \sum_{m=1}^M l_m = \frac{1}{M} \sum_{m=1}^M -\log P_{km}$$

$$= - \frac{1}{M} \sum_{m=1}^M \log P_{km}$$

$$\leq \log \frac{1}{M} \sum_{m=1}^M P_{km}$$

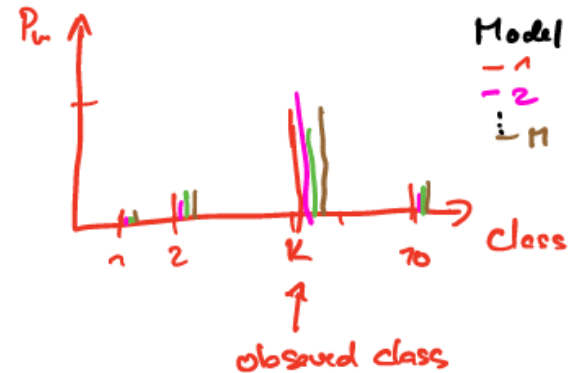
$$\geq -\log P_k^E = l^E \quad \square$$

use Jensen inequality



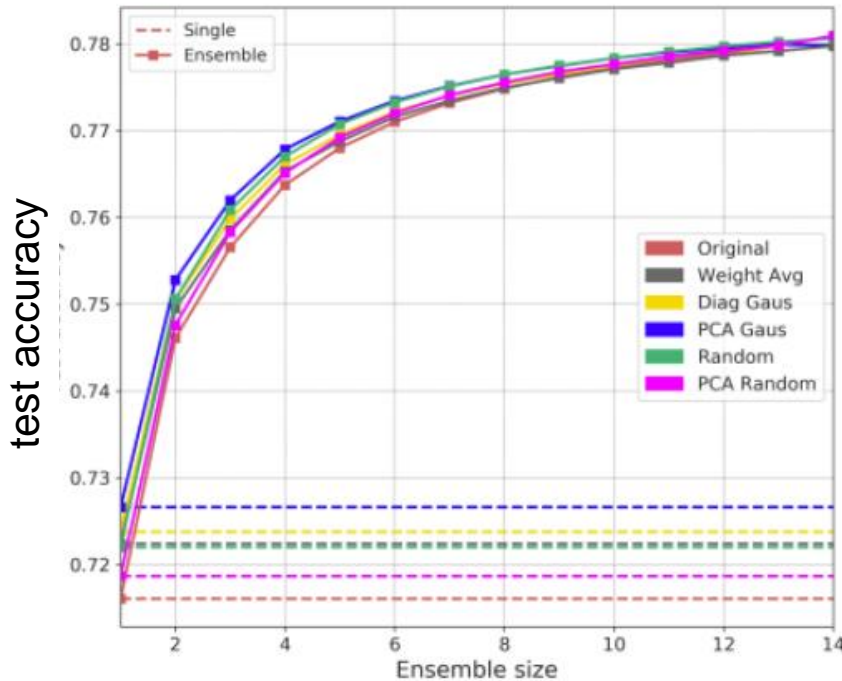
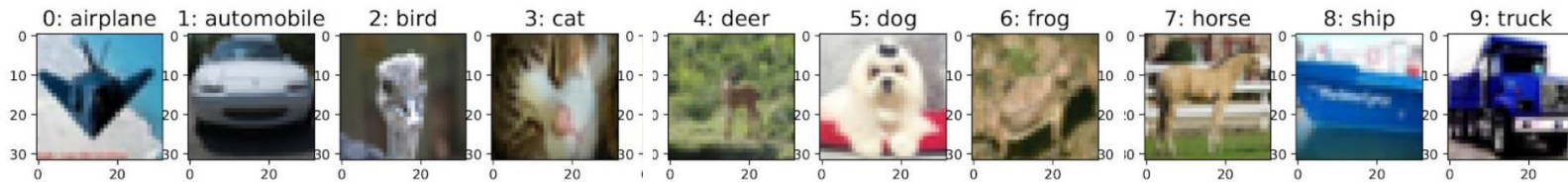
$$\log \bar{p} \geq \overline{\log p}$$

$$\log \frac{1}{M} \sum_{m=1}^M P_m \geq \frac{1}{M} \sum_{m=1}^M \log P_m$$

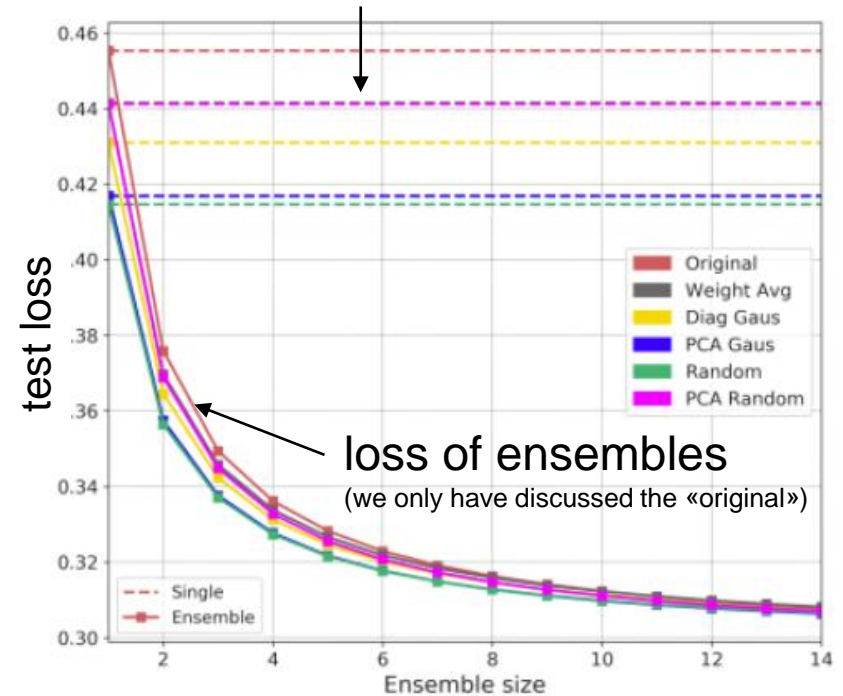




# Deep ensembling improves prediction power



loss of 5 single models



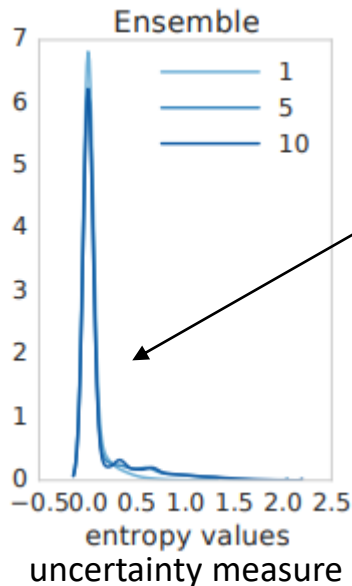
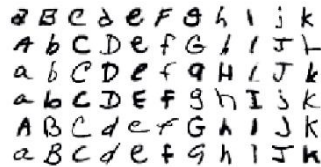
Ensembles with as few as 3 or 5 members are typically enough to achieve a performance gain.



# Deep ensembles improve uncertainty measures

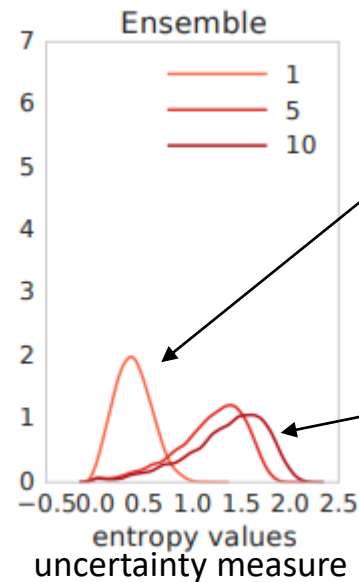
We want that a model, that is trained on normal MNIST letter data, should provide large uncertainties when applied on novel (NOT-MNIST) letter images.

Uncertainties for prediction of normal MNIST letter images



Ensemble uncertainty for a known class stays small for all sizes of the ensemble

Uncertainties for prediction of NOT-MNIST images

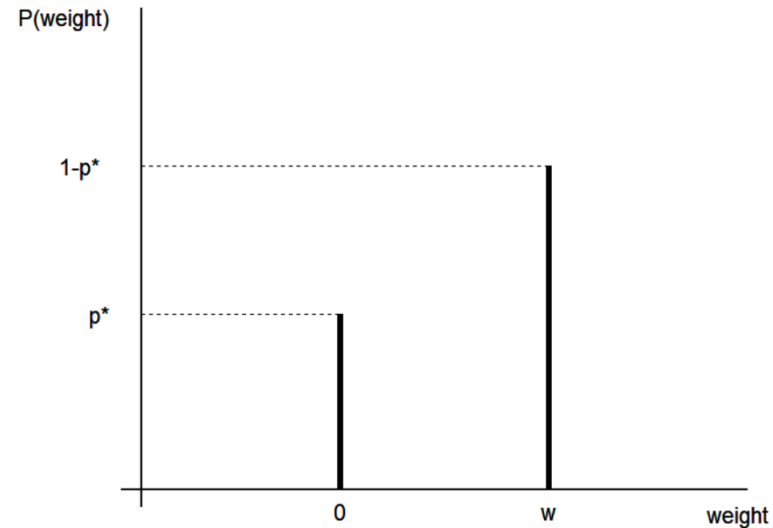
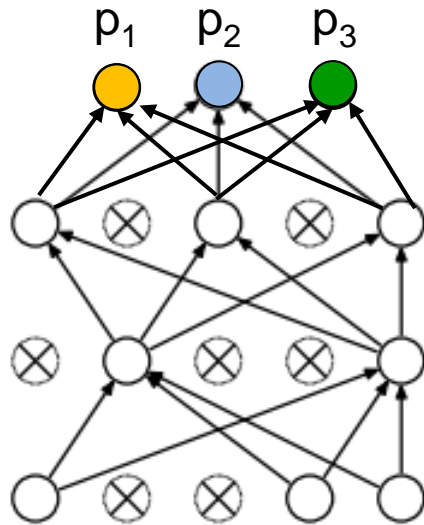


A single models shows less uncertainty for a novel class than an ensemble model

Ensemble uncertainty measures for novel classes improve with size of ensemble

# Dropout

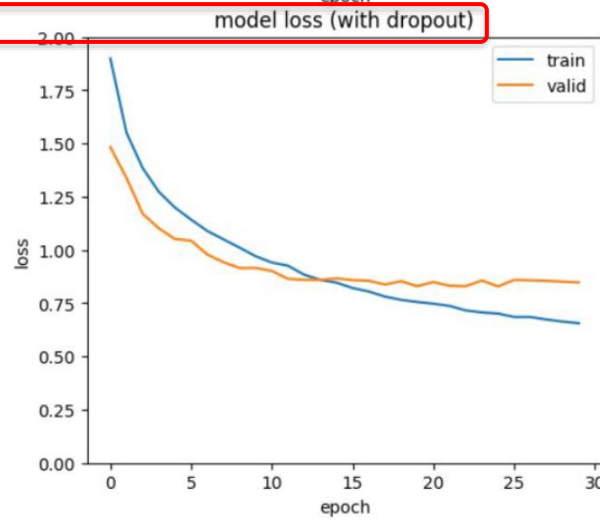
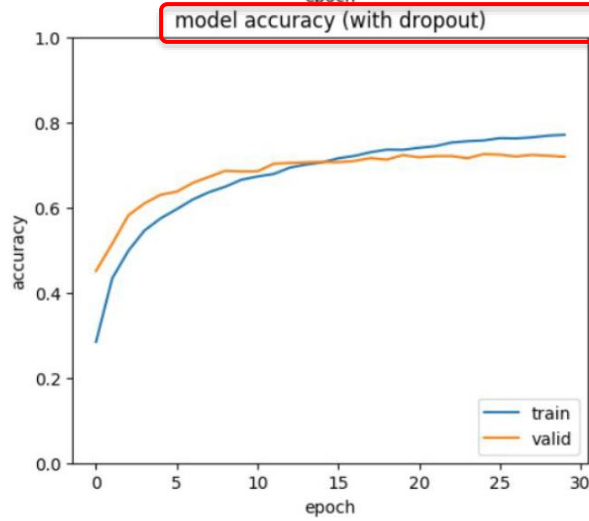
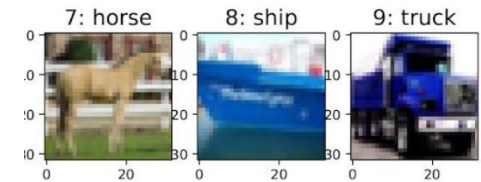
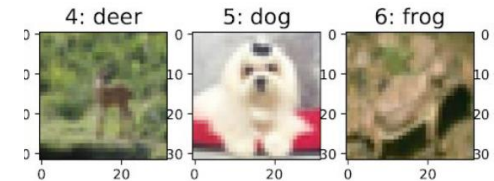
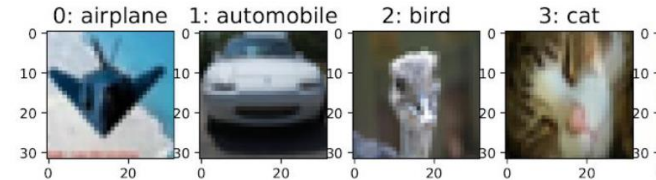
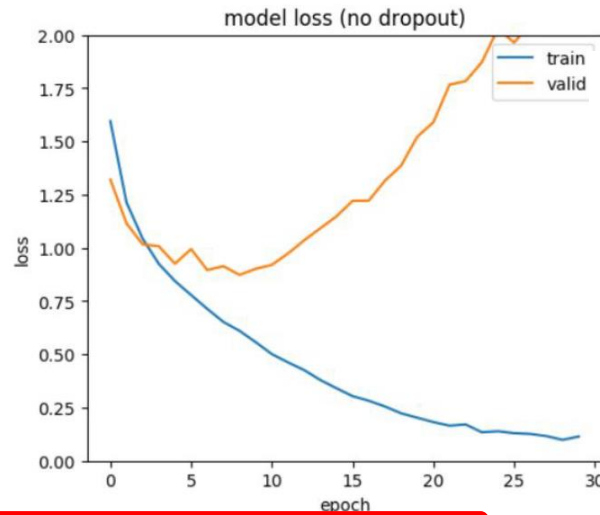
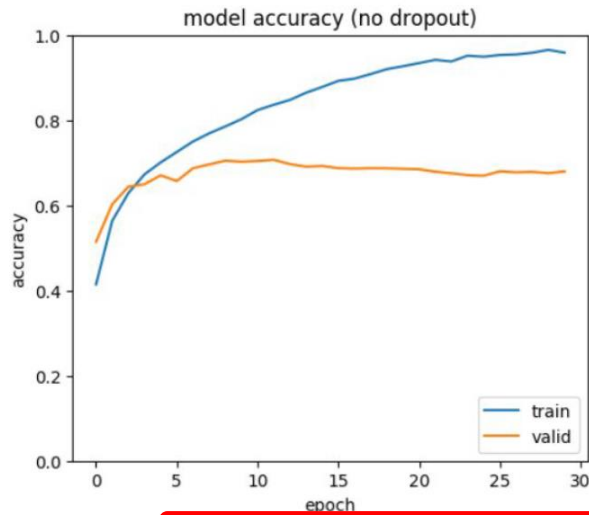
## Recall: Classical Dropout only during training



Using dropout during training implies:

- In each training step only weights to not-dropped units are updated  $\rightarrow$  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

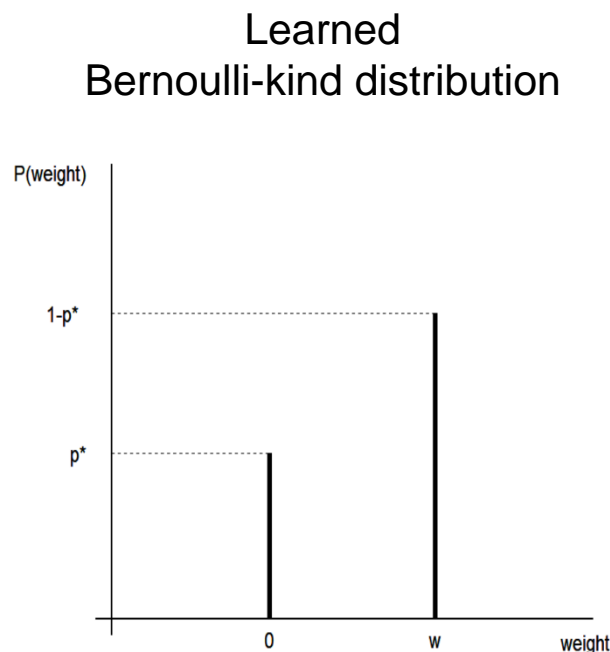
# Bayesian NN via MC Dropout

Yarin Gal et al. (2015):

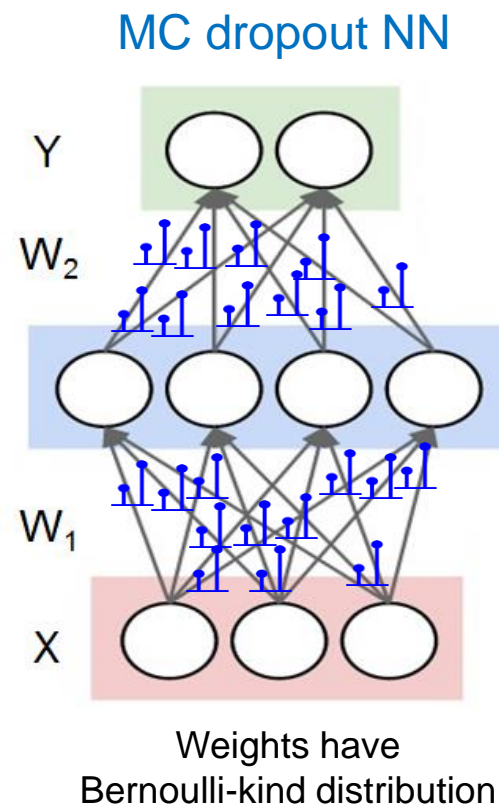
Via Dropout training we learned a whole weight distribution for each connection.

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

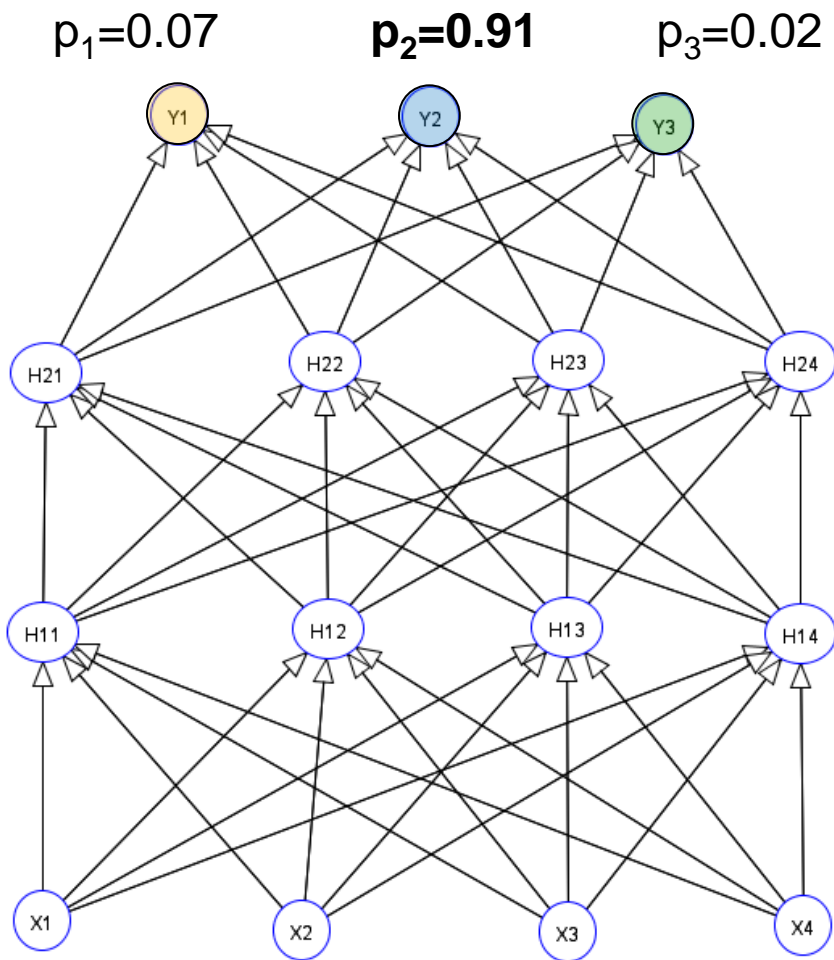


Which parameter has this  $q_\theta$ ?

The value  $w$ .

# When using Dropout only during training

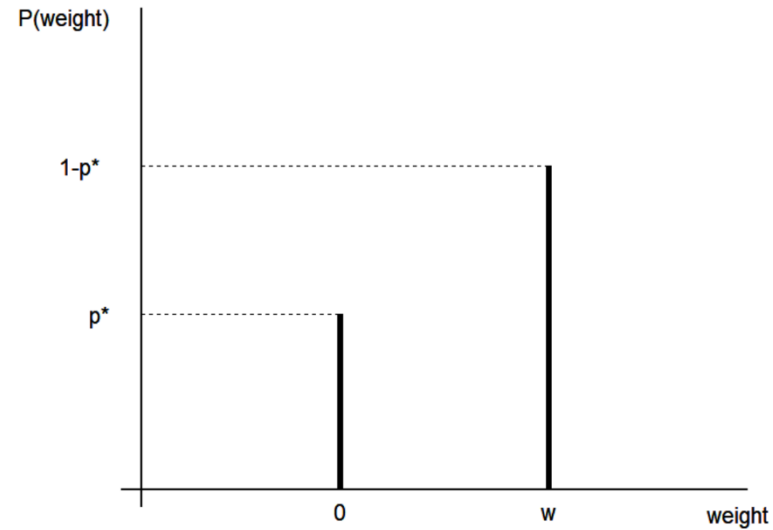
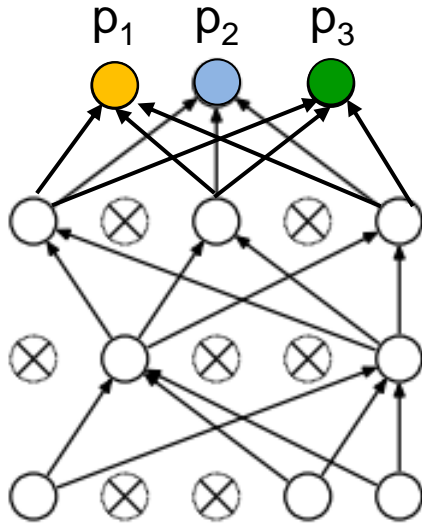
For non-Bayesian NN we freeze the weights after training to a value  $w \cdot p^*$  and use then the trained NN for prediction:



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

Input: image pixel values

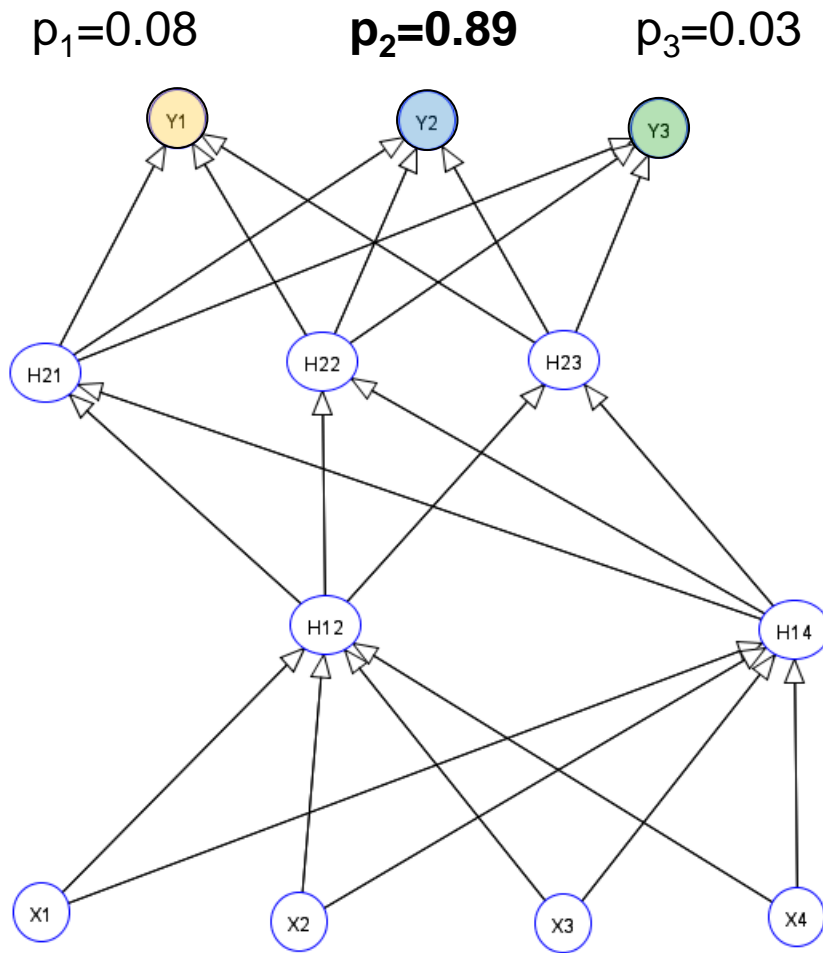
# MC Dropout: we also perform dropout during test time



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



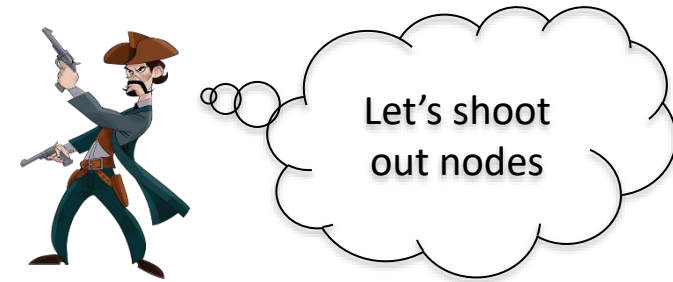
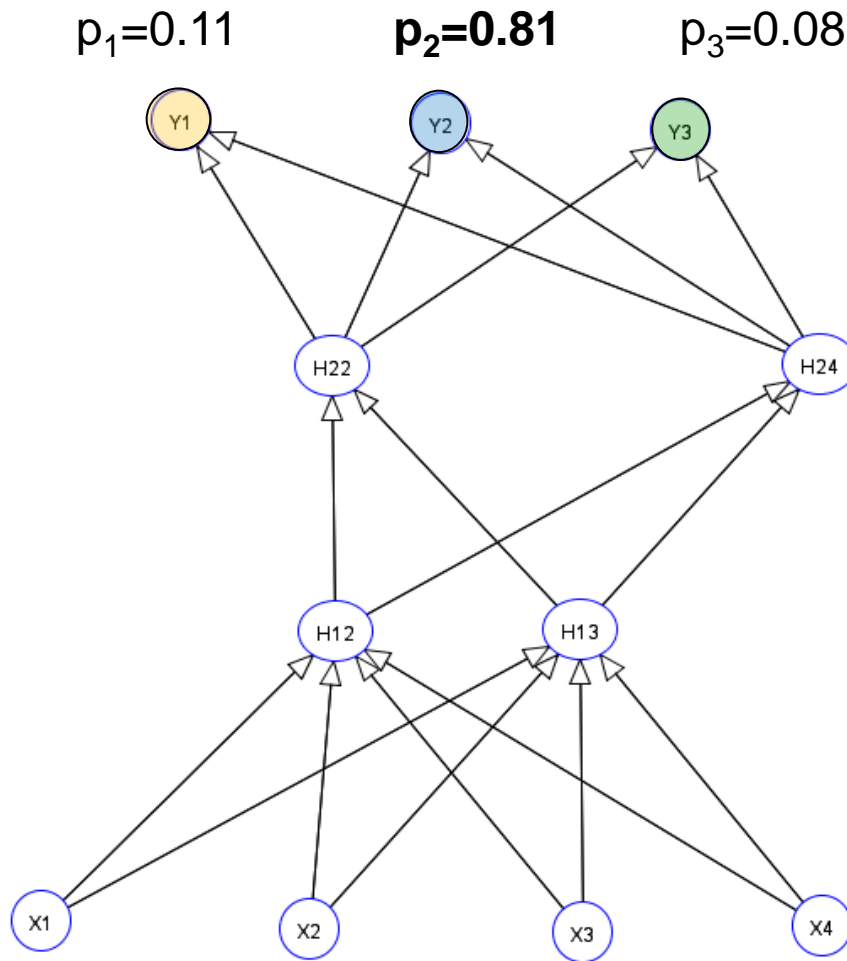
# MC Dropout during test time: Run 1



Stochastic dropout of units

Same input image

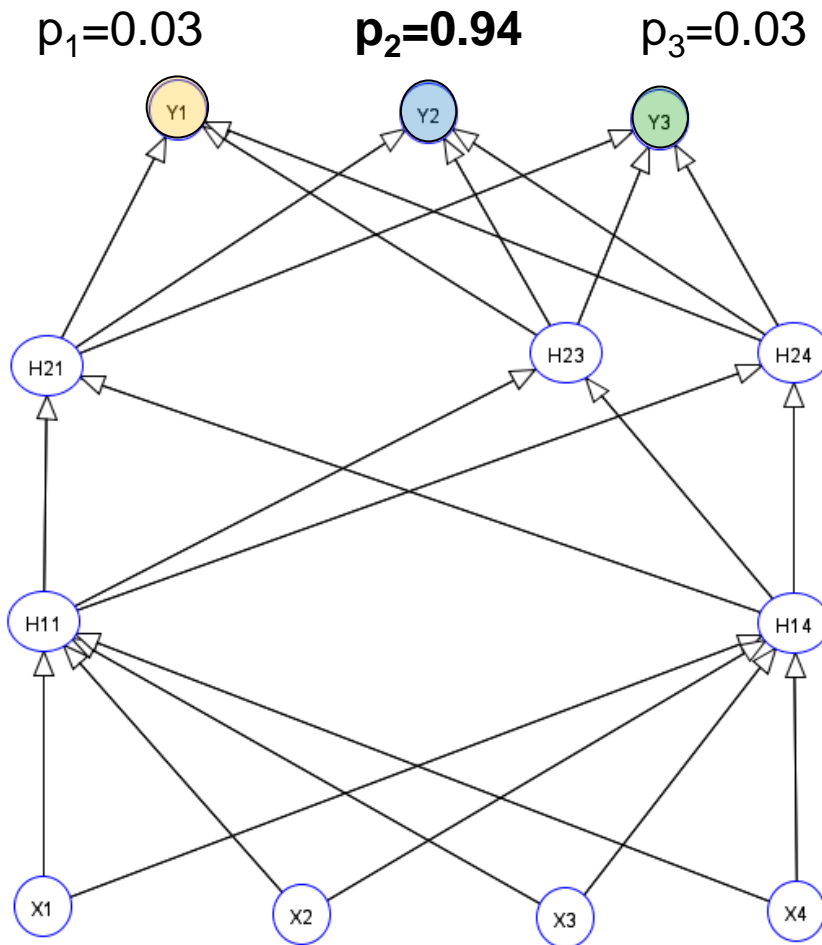
# MC Dropout during test time: Run 2



Stochastic dropout of units

Same input image

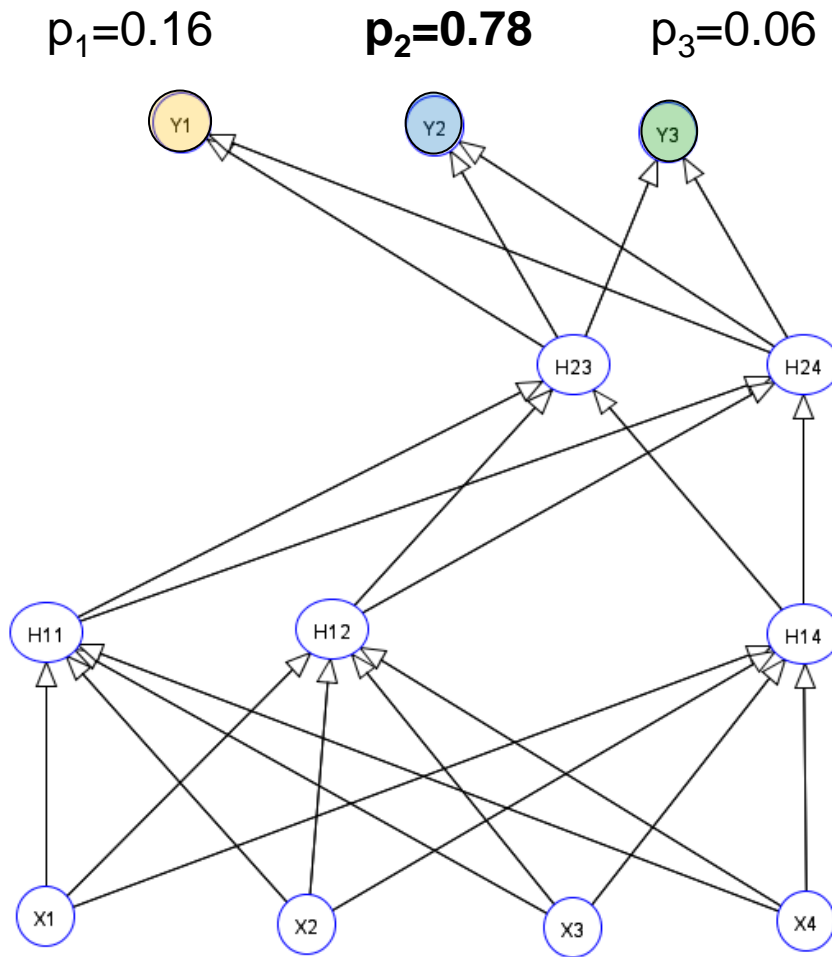
# MC Dropout during test time: Run 3



Stochastic dropout of units

Same input image

# MC Dropout during test time: Run 4



Stochastic dropout of units

Same input image

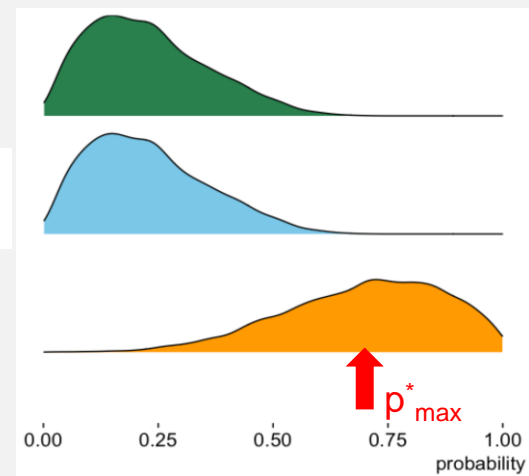
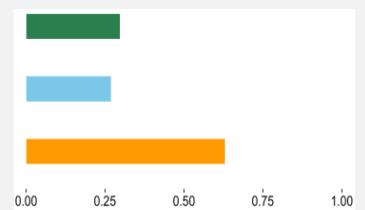
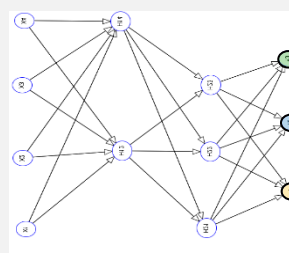
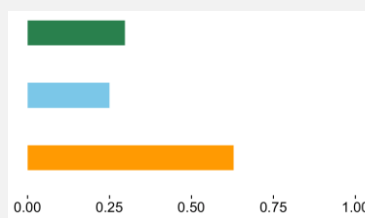
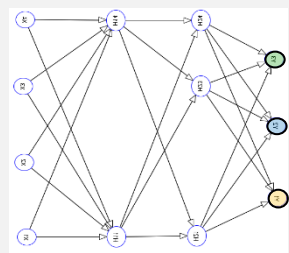
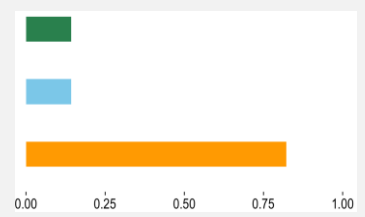
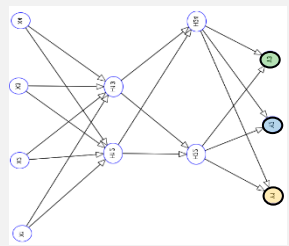
# MC Dropout during test time yields a multivariate predictive distribution for the parameters



use dropout  
also during  
prediction



Many Dropout Runs in forward pass

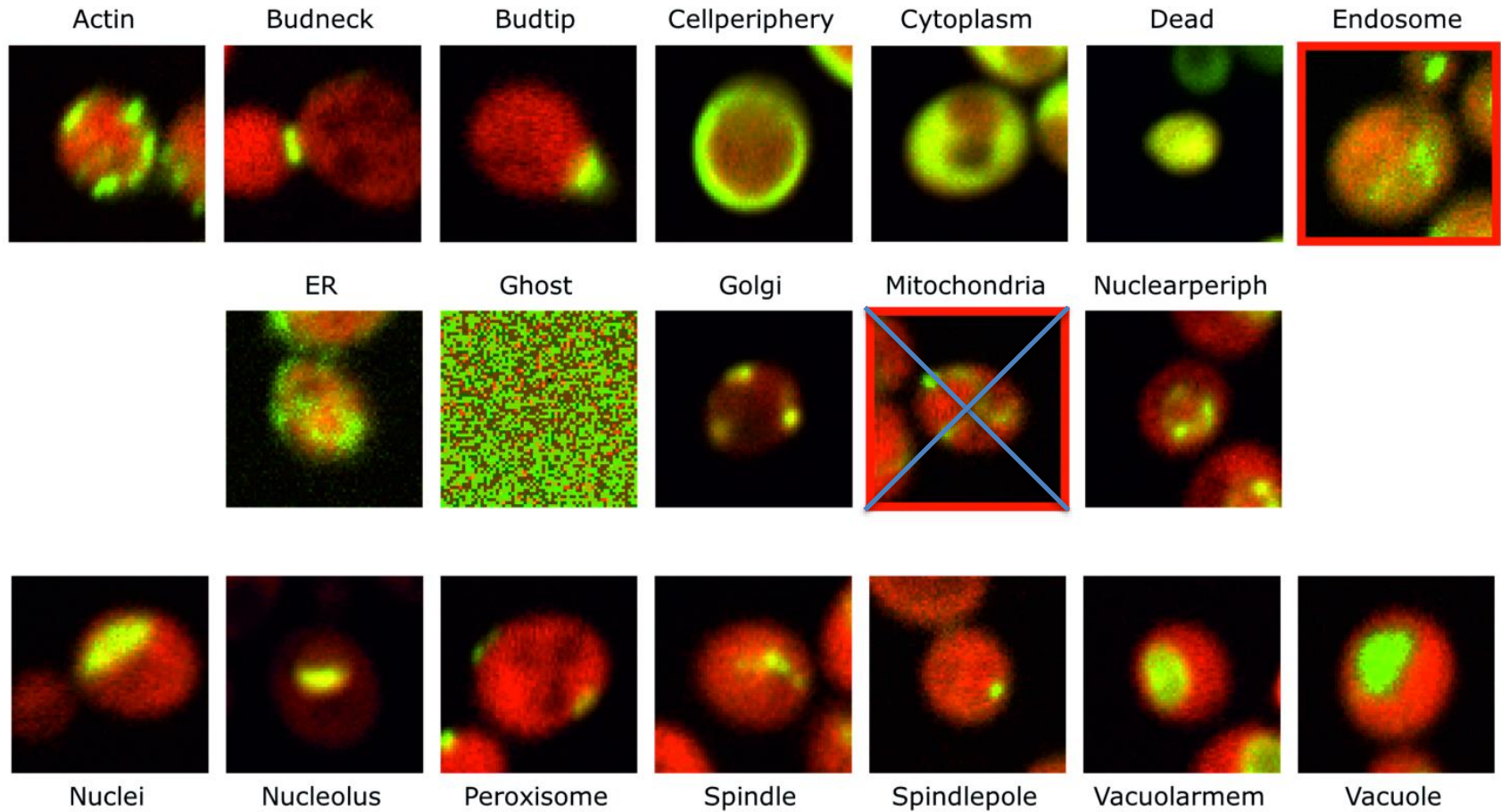


CNN predicts class “collie”  
but with high uncertainty

...

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

# Experiment with unknown phenotype

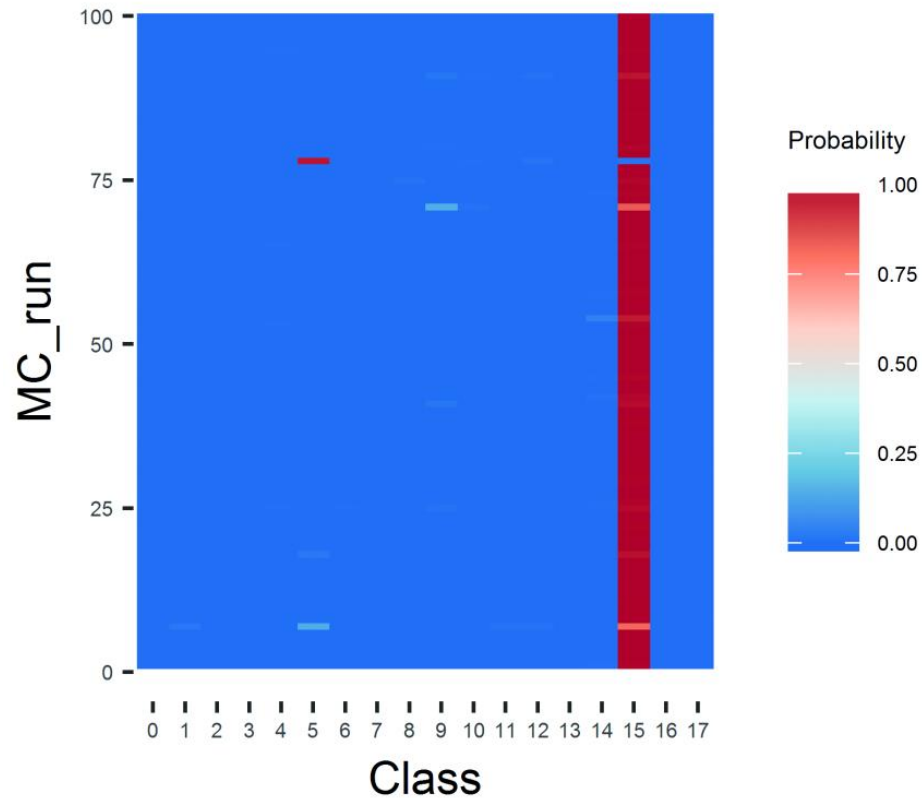


Dürr O, Murina E, Siegismund D, Tolkachev V, Steigle S, Sick B. Know when you don't know, Assay Drug Dev Technol. 2018

# Probability distribution from MC dropout runs

## Image with known class 15

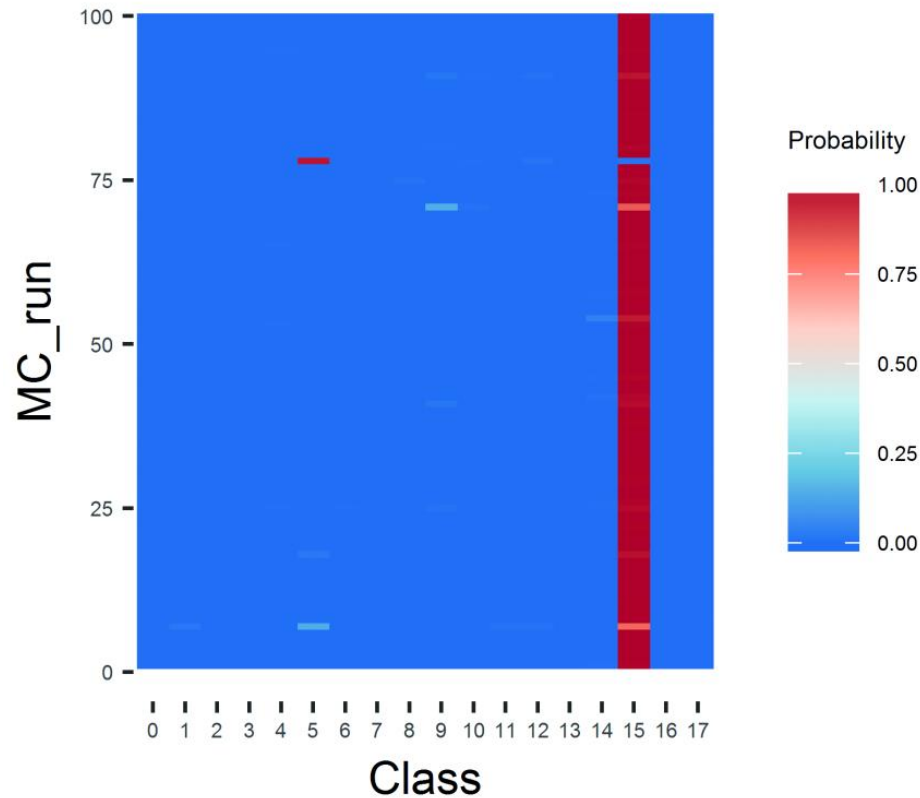
100 MC predictions for an image with known phenotype 15



# Probability distribution from MC dropout runs

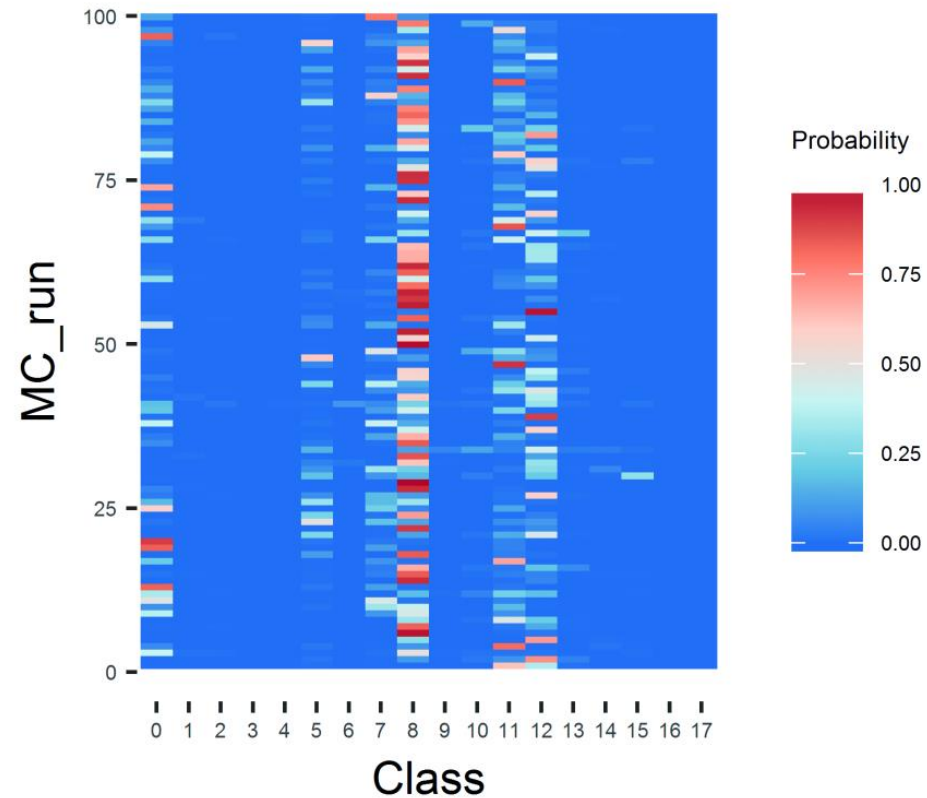
## Image with known class 15

100 MC predictions for an image with known phenotype 15



## Image with **unknown** class

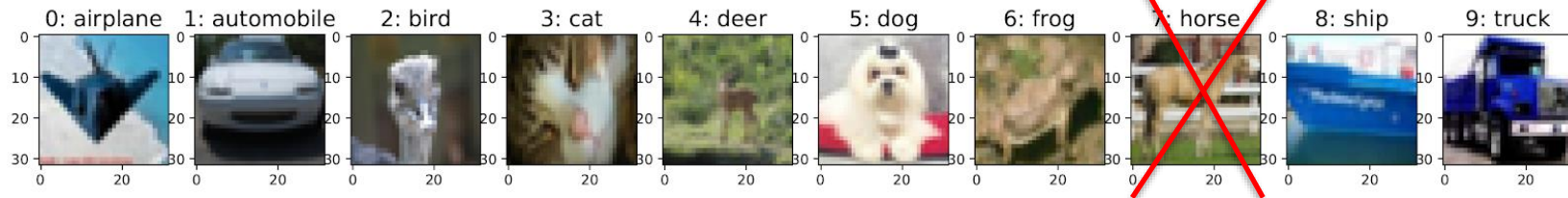
100 MC predictions for an image with an unknown phenotype





# 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)

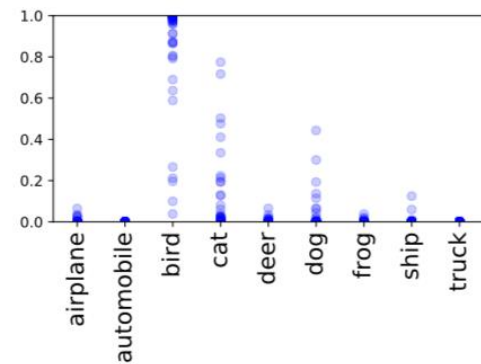
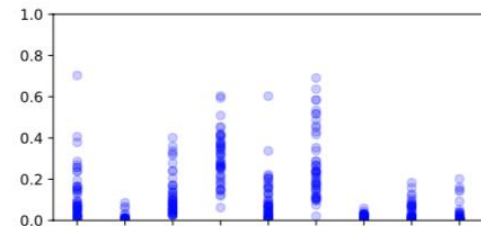
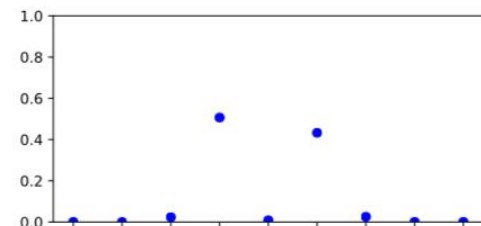
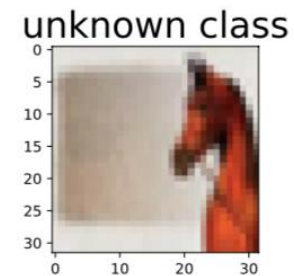
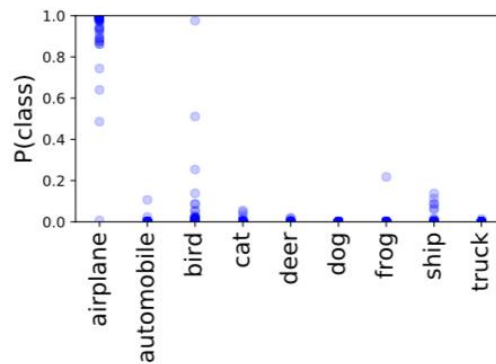
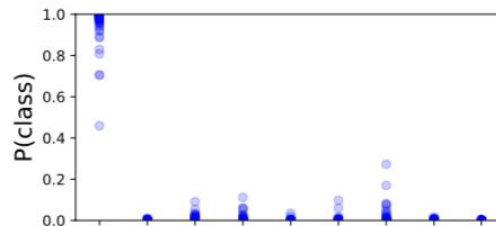
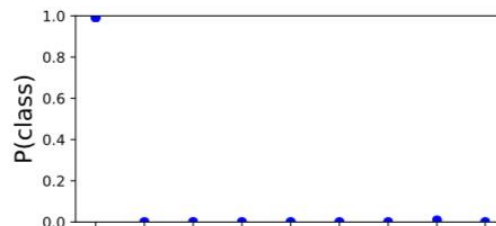
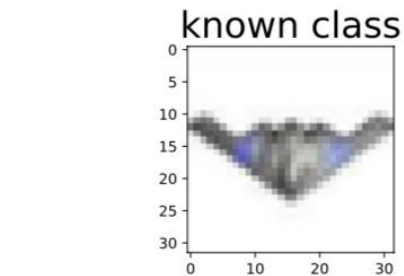
Run the Variational Inference without trying to understand it – we will care later!

# Looking at the predictive distribution!

Classical CNN (no Bayes)

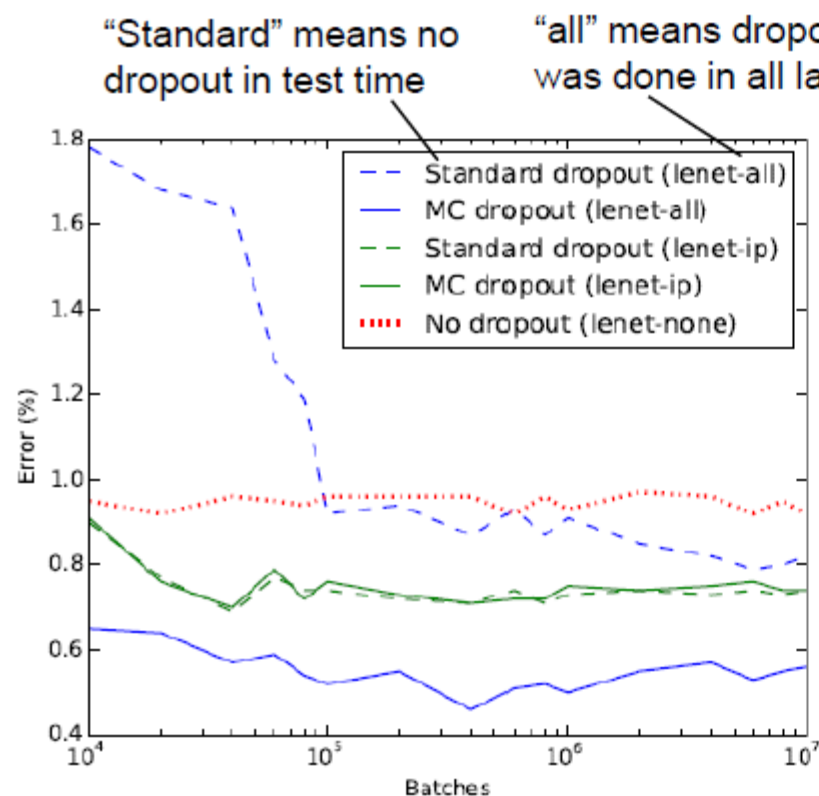
Gaussian VI

MC dropout

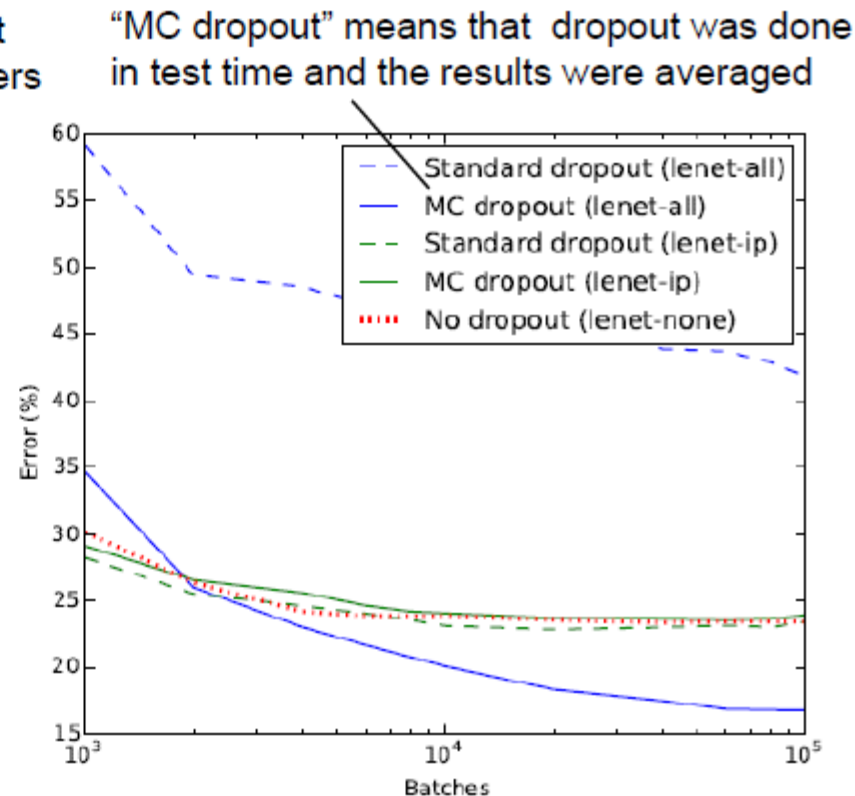


# MC Dropout increases prediction performance

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.



(a) MNIST

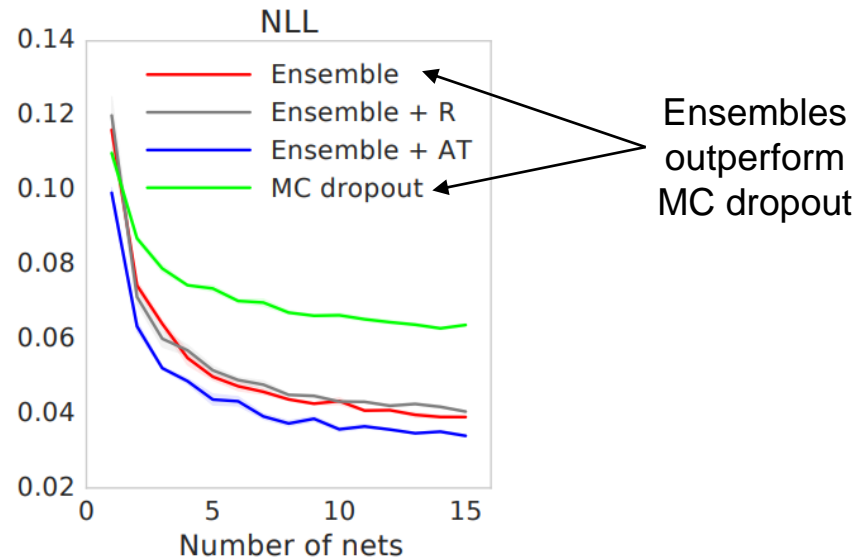


(b) CIFAR-10

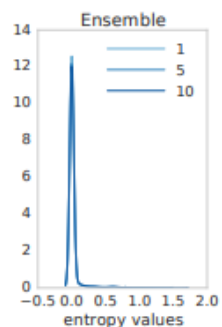
# How does MC dropout compare with deep ensembles?

- 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

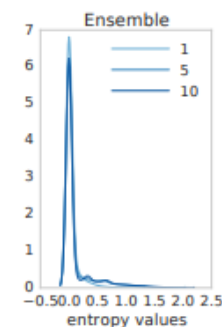
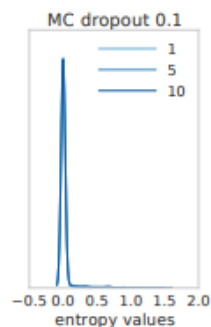
MNIST classification



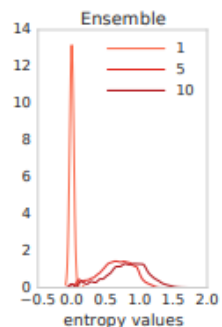
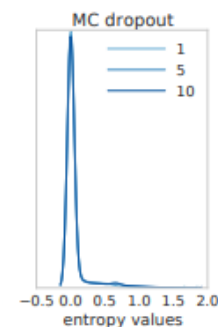
# Compare MC dropout and deep ensembles uncertainties



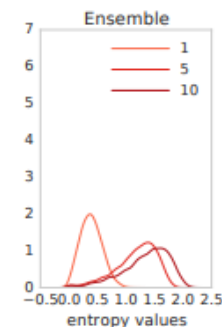
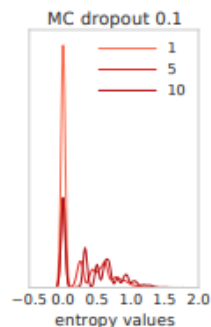
The novel class uncertainty is better in ensemble models than in MC Dropout models



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



(a) MNIST-NotMNIST



(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  $M$ .