

Machine Intelligence:: Deep Learning

Week 6

Beate Sick, Oliver Dürr, Jonas Brändli

Probabilistic models with flexible CPDs

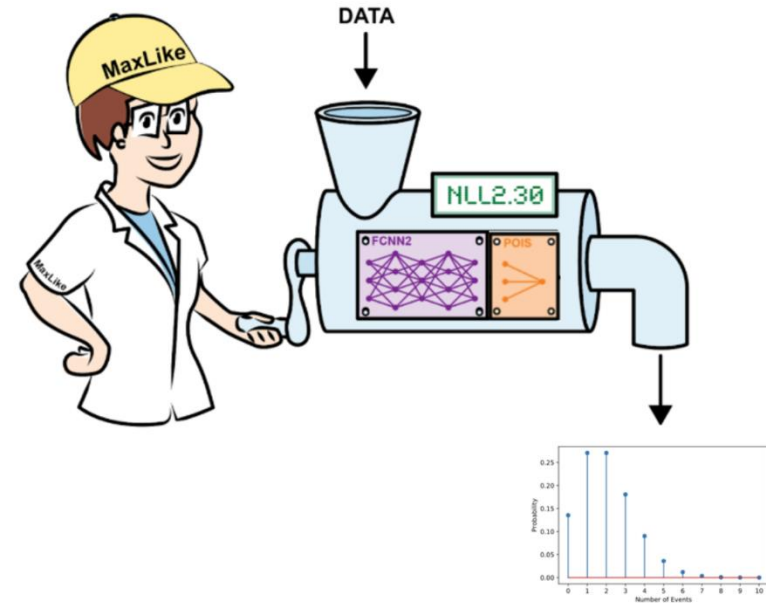
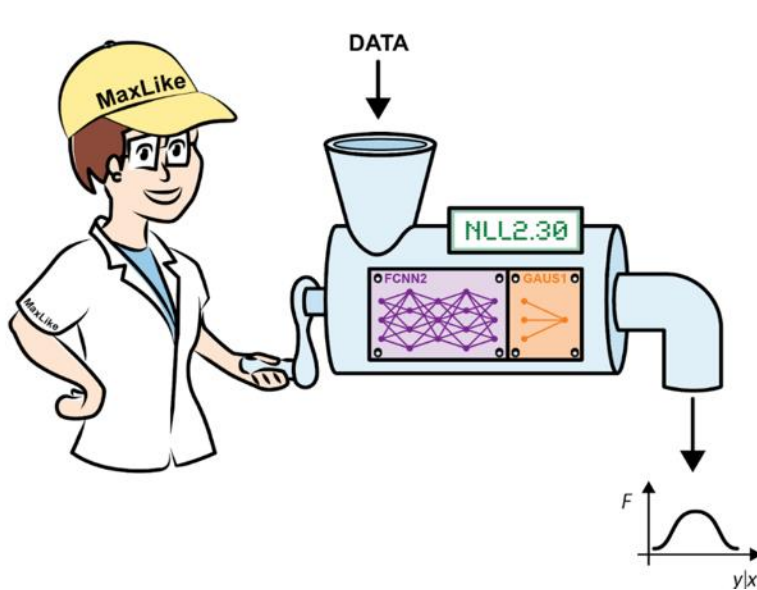
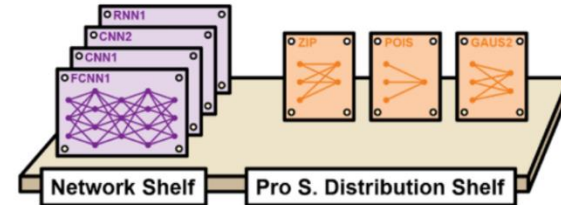
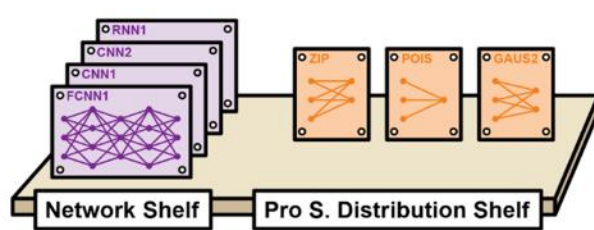
Outline of the DL Module (tentative)

- Day 1: Jumpstart to DL
 - What is DL
 - Basic Building Blocks
 - Keras
- Day 2: CNN I
 - ImageData
- Day 3: CNN II and RNN
 - Tips and Tricks
 - Modern Architectures
 - 1-D Sequential Data
- Day 4: Looking at details
 - Linear Regression
 - Backpropagation
 - Resnet
 - Likelihood principle
- Day 5: Probabilistic Aspects
 - Likelihood principle (cont'd)
 - TensorFlow Probability (TFP)
 - Negative Loss Likelihood NLL
 - Count Data
- Day 6: Probabilistic models in the wild
 - Complex conditional distributions
 - Grad CAM
 - Deep Ensembles
- Day 7: Uncertainty in DL
 - Bayesian Modeling
- Day 8: Uncertainty cont'd
 - Bayesian Neural Networks
 - Projects

Projects please register (see website)

<https://docs.google.com/spreadsheets/d/18VFrPbKq3YSOg8Ebc1q1wGgkfqaWI7IkCCEIGEDGj6Q/edit#gid=0>

We have a flexible tool where the choice of the architecture and the choice of the outcome distribution is independent



Modeling count data continued

Recall the camper example

N=250 groups visiting a national park

Y=count: number of fishes caught

X1=persons: number of persons in group

X2=child: number of children in the group

X3=bait: indicates if live bait was used

X4=camper: indicates if camper is brought



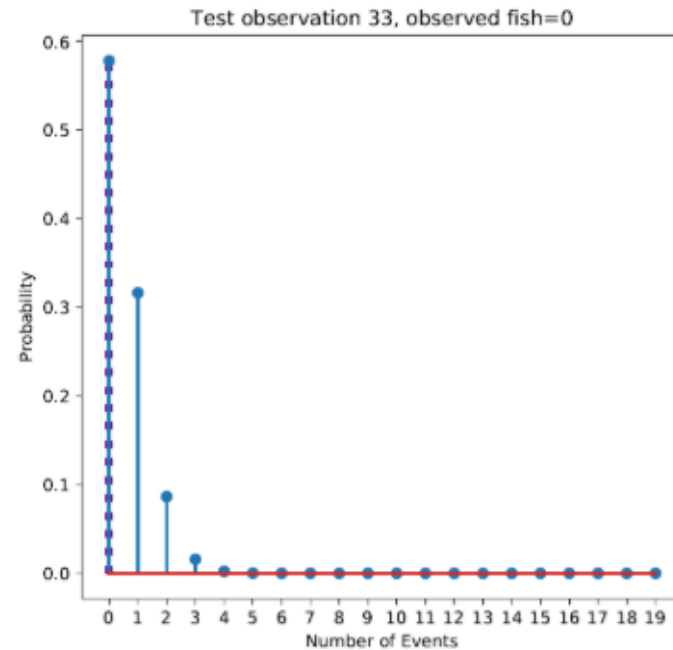
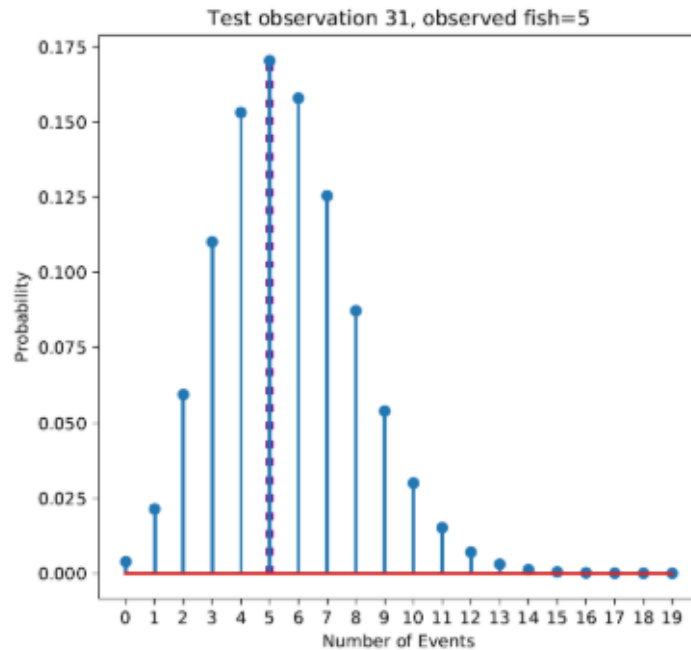
Data: <https://stats.idre.ucla.edu/r/dae/zip>

Model 2: Poisson regression, predicted CPDs for test observations

Predict CPD for outcome in test data:

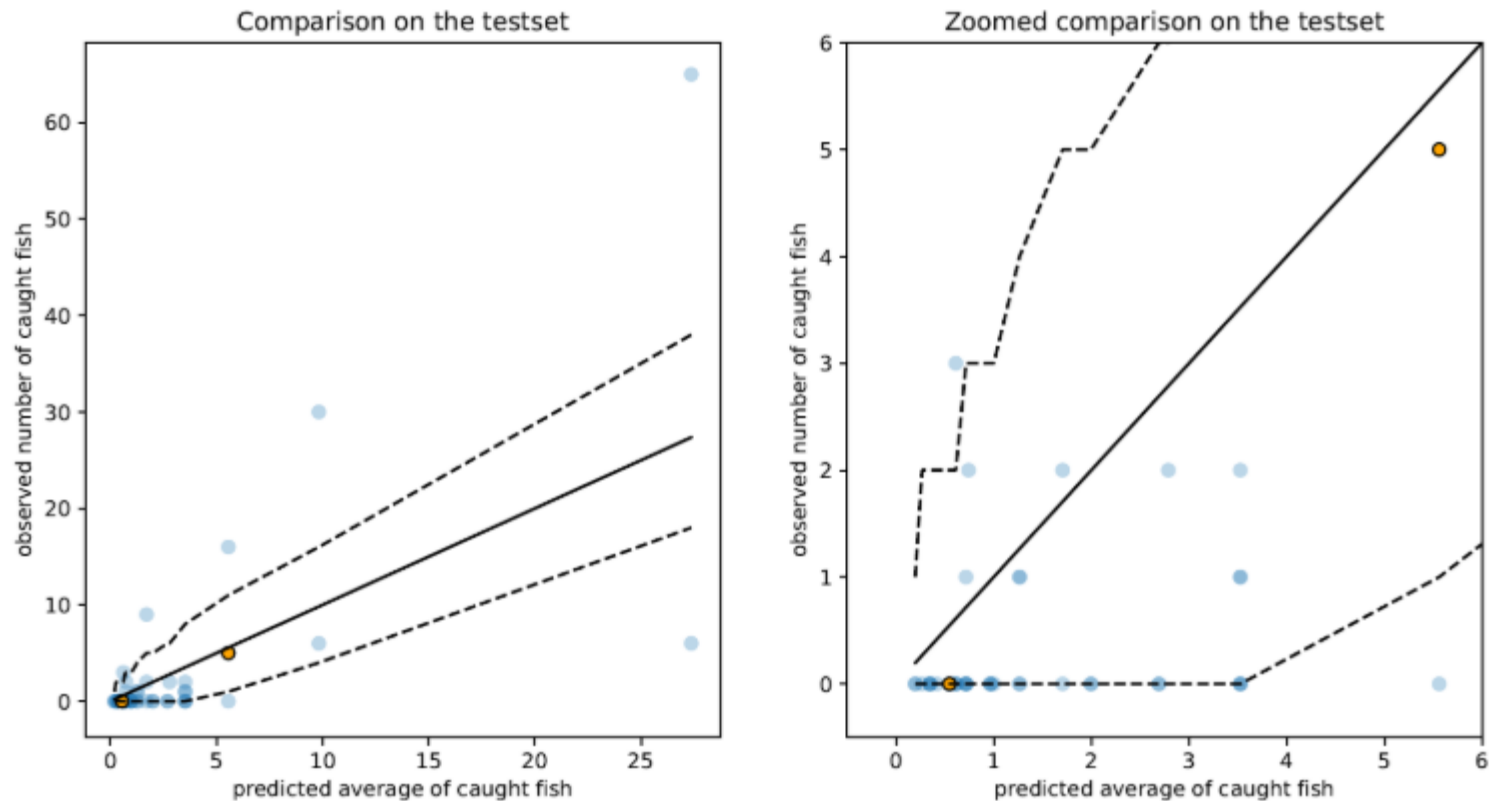
Group 31 used livebait, had a camper and were 4 persons with one child. $Y=5$ fish.

Group 33 used livebait, didn't have a camper and were 4 persons with two children. $Y=0$ fish.



What is the likelihood of the observed outcome in test obs 31 and 33?

Model 2: Poisson regression, visualize the CPDs by quantiles



The mean of the CPD is depicted by the solid lines.

The dashed lines represent the 0.025 and 0.975 quantiles, yielding the borders of a 95% prediction interval.

Note that different combinations of predictor values can yield the same parameters of the CPD.

Besprechung der Aufgabe



Use NN and tfp to fit a poisson model



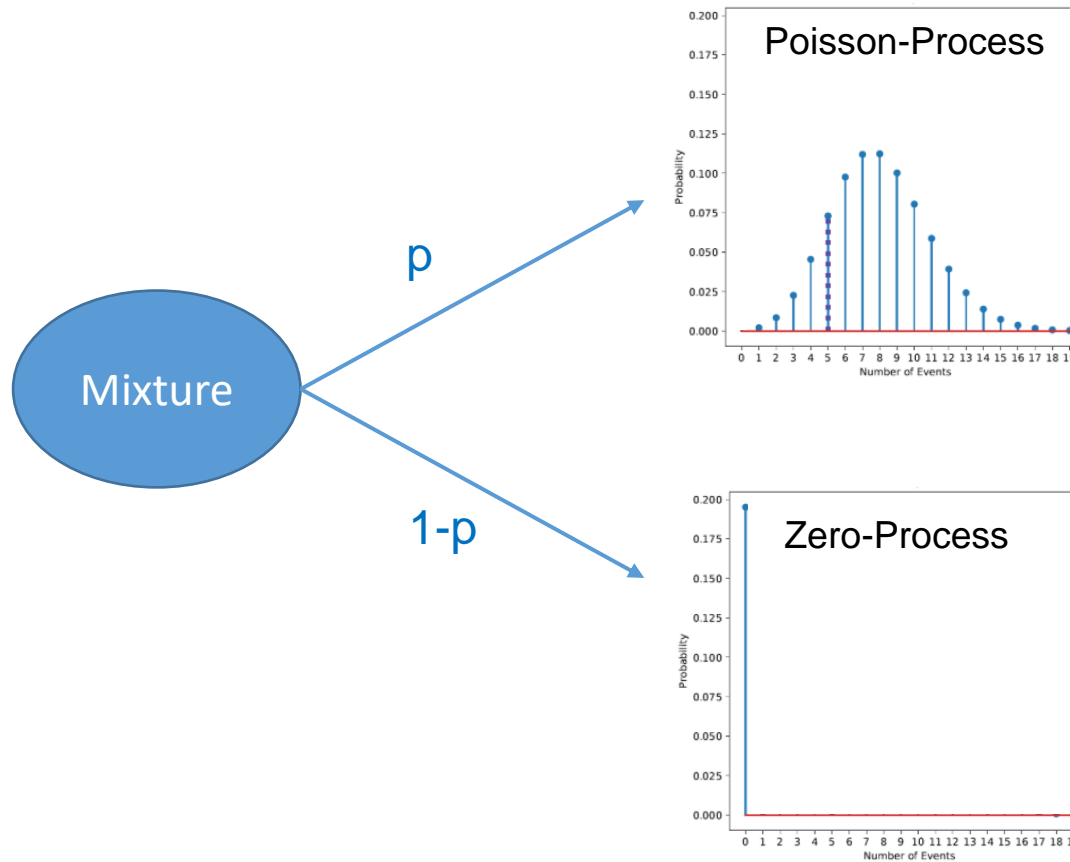
https://github.com/tensorchiefs/dl_course_2022/blob/master/notebooks/14_poisreg_with_tfp.ipynb

Modeling count data:

M3: ZIP regression

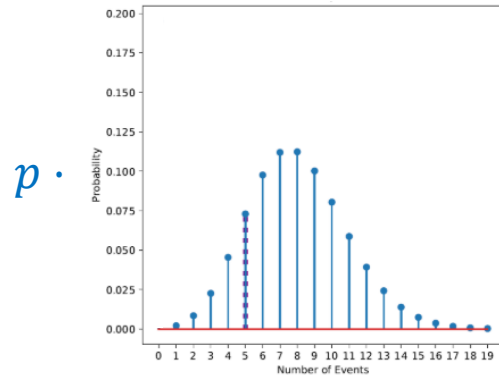
Zero-Inflated Poisson (ZIP) as Mixture Process

How many fish a group catches does not only depend on luck ;-)



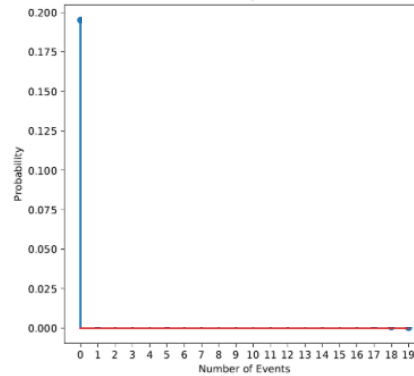
Zero-Inflated Poisson (ZIP) can be seen as Mixture Distribution

Poisson-Process



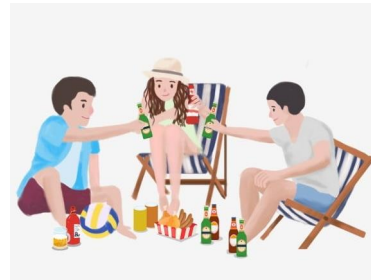
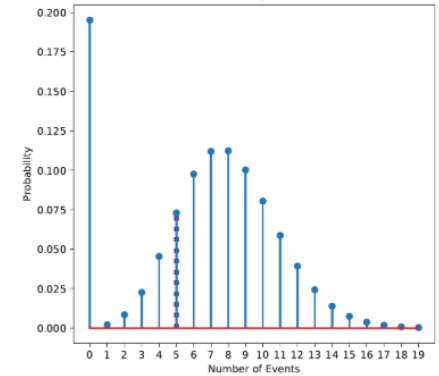
$+(1-p) \cdot$

Zero-Process

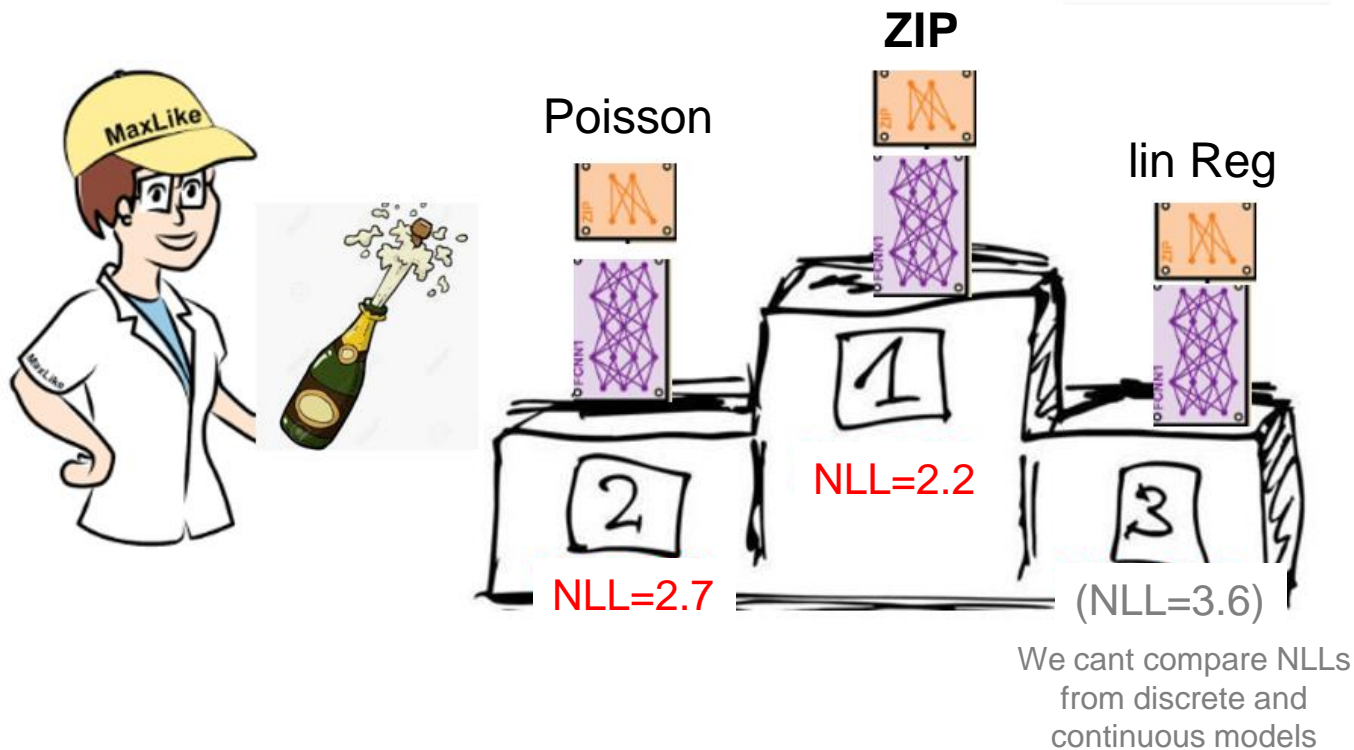


$=$

Zero-inflated Poisson

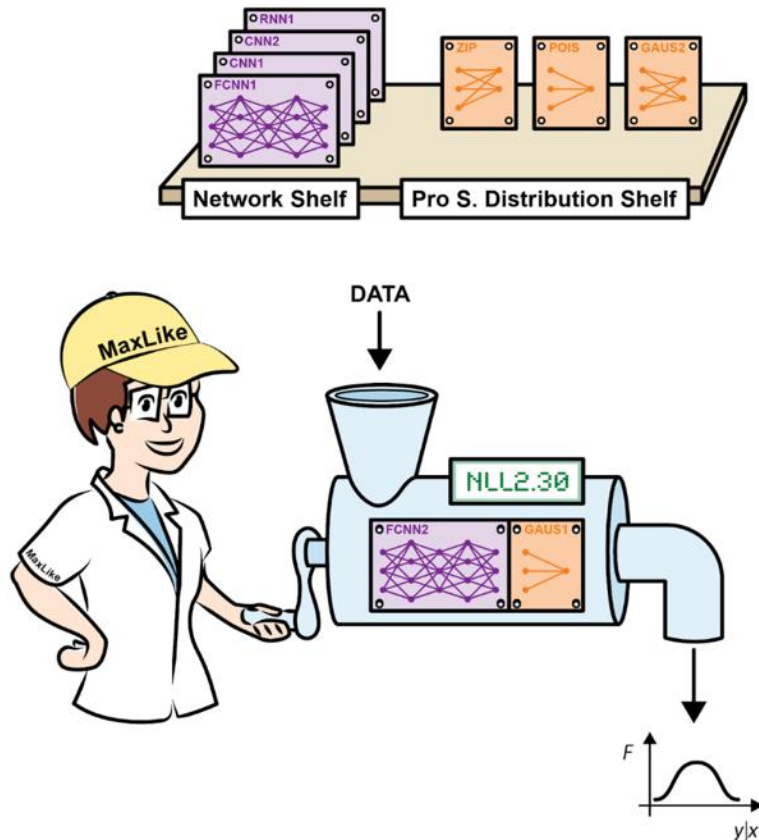


Validation NLL allows to rank different probabilistic models



Count Distributions in the real world

Building state of the art networks



Two Design choices

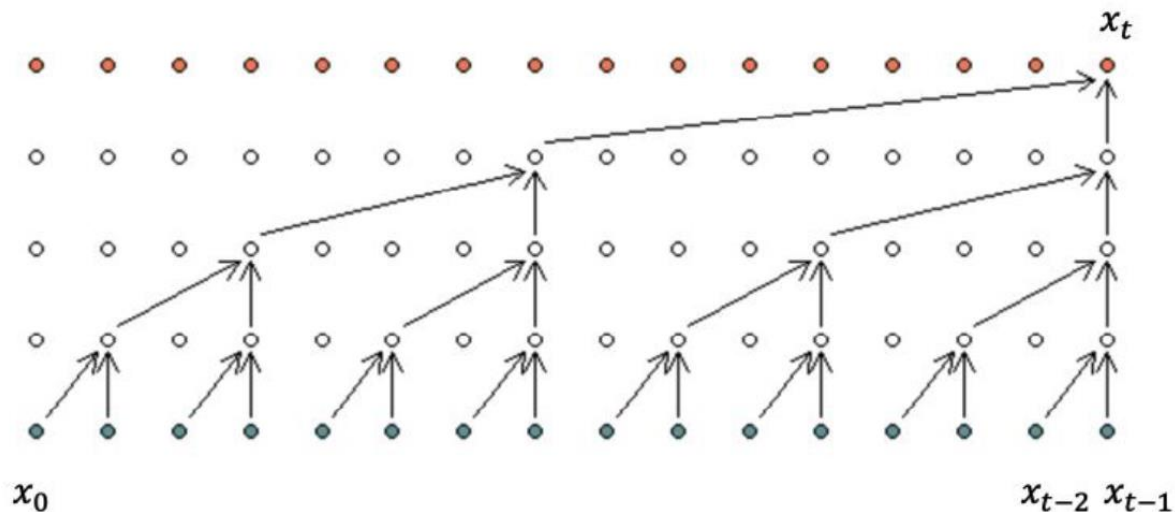
1. Choose the right network architecture to exploit structure of your problem ("inductive bias")
 - Images CNN
 - Sequences RNN or 1D convolution
2. Choose the right CPD for your data
 - Train with NLL
 - Evaluate performance on validation with NLL

Two state-of-the-art models PixelCNN and Wavenet have been improved by using *discretized logistic mixture* distributions as CPD.

PixelCNN and Wavenet In/Output

- Both are autoregressive models $p(x_t | x_{t-1}, \dots, x_0)$
 - Wavenet
 - Predicts the next audio sample x_t from all previous x_{t-1}, \dots, x_0
 - x_t can take a number between 0 and $2^{16}-1 = 65,535$
 - PixelCNN
 - Predicts the next pixel x_t from all previous x_{t-1}, \dots, x_0
 - x_t can take a value between 0 and $2^8-1 = 255$

Wavenet architecture (see lecture 06_sequences)



PixelCNN and Wavenet In/Output

Modelling $p(x_t | x_{t-1}, \dots x_0)$

- Naive approach output probabilities for the different classes
 - This has been done in first versions
 - Does not take nature of the output data into account
- Better find a *flexible* distribution to model the discrete data between 0 and say 65536 (2^{16})
 - Both architectures have been improved by mixtures of discrete logistics

We now have a look at the mixtures of discrete logistics.

The code for the figures is in:

https://github.com/tensorchiefs/dl_book/blob/master/chapter_06/nb_ch06_01.ipynb

Making sense of discretized logistic mixture

- The base distribution: the logistic with 2 parameters
 - loc kind of center
 - scale spread

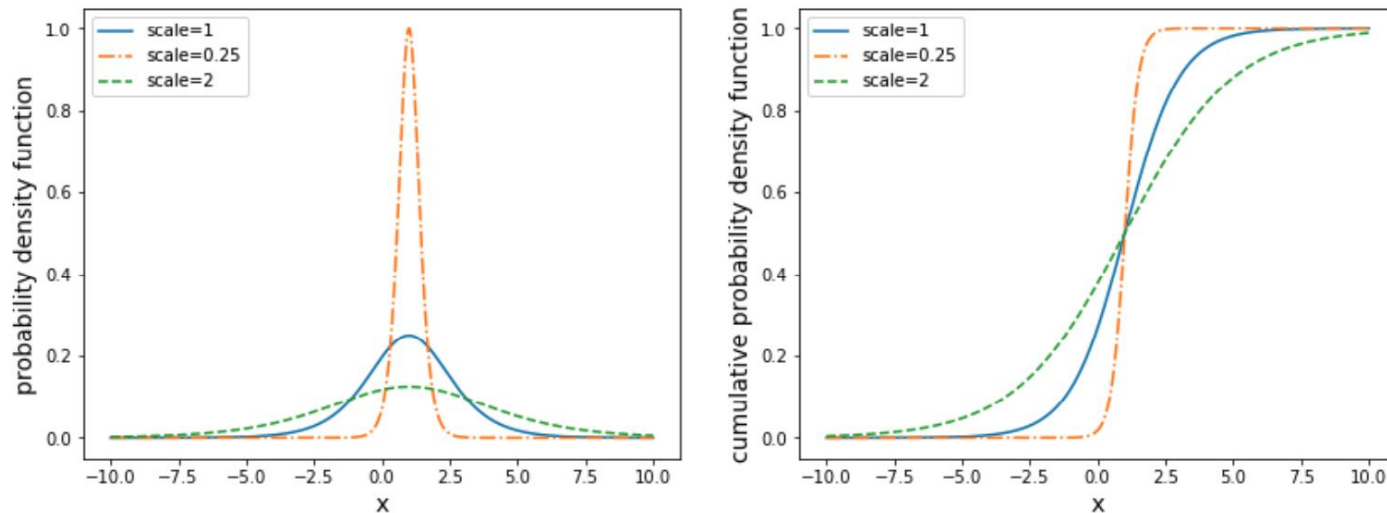
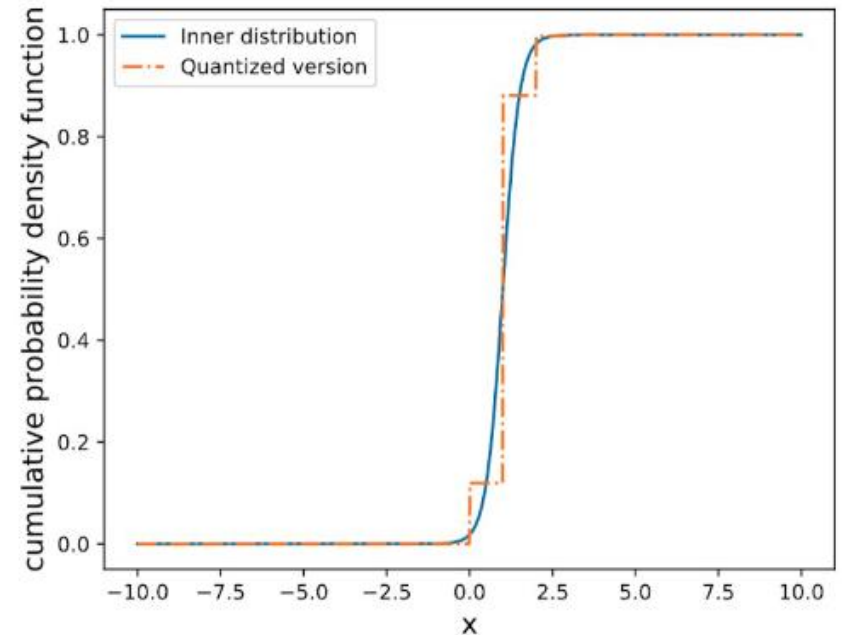
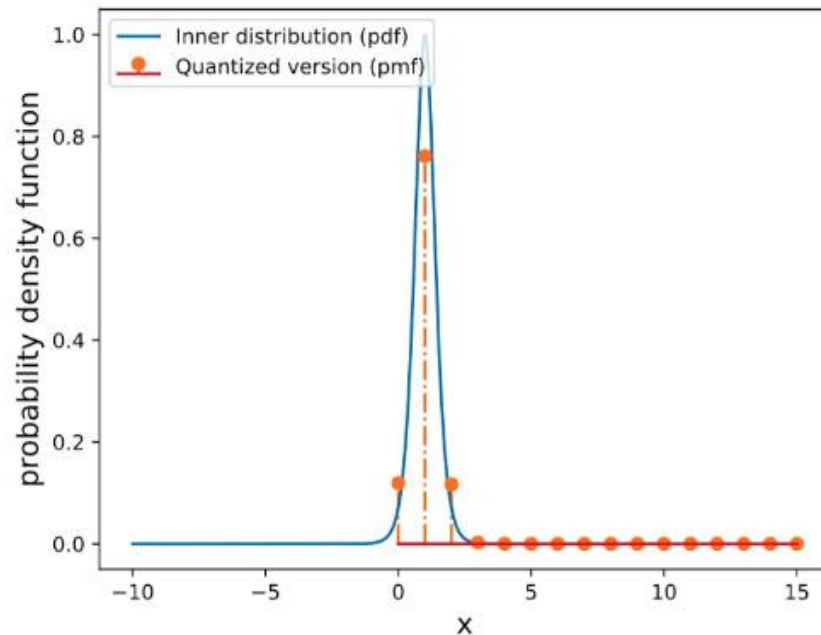


Figure 6.3 Three logistic functions created using `tfd.Logistic(loc=1, scale=scale)` with values of 0.5, 1.0, and 2.0 for the scale parameter. On the left is the probability density function (PDF) and on the right, the cumulative probability density function (CDF).

- The CDF has the same functional form as the sigmoid activation
- Still not discrete

Making sense of discretized logistic mixture

- The quantized (discretized) base distribution the logistic



Making sense of discretized logistic mixture

- Mixing several (here 2) discretized logistics

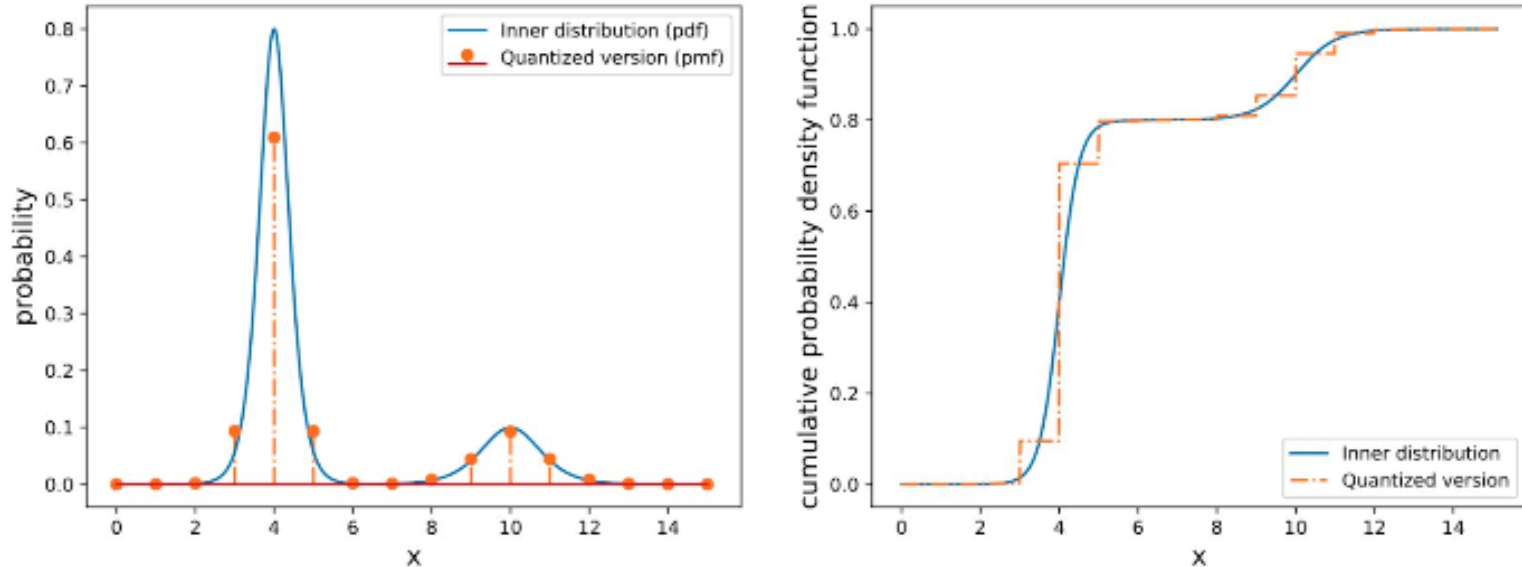


Figure 6.5 The resulting discrete distribution when mixing two logistic distributions (see listing 6.2 for the code that produces these plots)

Number of components, which need to be determined by network?

- 2 (or 1) for mixing
- 2 for scale
- 2 for location

Case Study for discretized logistic mixture

The original PixelCNN++ and Wavenet take too long to train. The fish data is too small for deep learning models. Let's use a large count data set.

Goal: Probabilistic model for deer activity conditioned on the time (in day and year).



wild	year	time	daytime	weekday
0	2002.0	0.000000	night.am	Sunday
0	2002.0	0.020833	night.am	Sunday
...
1	2002.0	0.208333	night.am	Sunday
0	2002.0	0.229167	pre.sunrise.am	Sunday
0	2002.0	0.270833	pre.sunrise.am	Sunday

The columns have the following meaning:

- wild: the number deers killed in a road accident in Bavaria
- year: the year (from 2002 to 2009 in the training set from 2010 to 2011 for the test set)
- time: the number of days to the first event. These numbers are measured in fractions of a day.

Data on deer related car accidents in the years 2002 until 2011 in Bavaria, Germany.

Target variable (wild): use number of deers killed during 30 minute period as surrogate

Modeling count data

Goal:

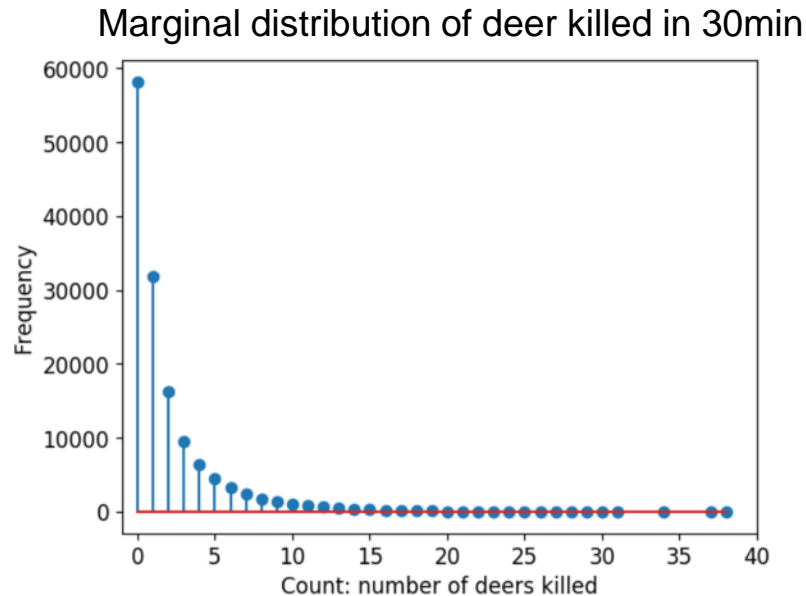
Predict CPD for $y=\text{\#deers-killed-in-30min}$, given x (time and derived variables).

Possible CPD models:

$$(y | x) \sim \text{Pois}(\lambda_x)$$

$$(y | x) \sim \text{ZIP}({}^z p_x, \lambda_x)$$

$$(y | x) \sim \text{discretizedLogisticMix}({}^1 p_x, {}^2 p_x, {}^3 p_x, {}^1 \mu_x, {}^2 \mu_x, {}^3 \mu_x, {}^1 \sigma_x, {}^2 \sigma_x, {}^3 \sigma_x)$$



Code

```
1 def quant_mixture_logistic(out, bits=8, num=3):
2     loc, un_scale, logits = tf.split(out, #A
3                                     num_or_size_splits=num,
4                                     axis=-1)
5     scale = tf.nn.softplus(un_scale) #B
6     discretized_logistic_dist = tfd.QuantizedDistribution(
7         distribution=tfd.TransformedDistribution( #C
8             distribution=tfd.Logistic(loc=loc, scale=scale),
9             bijector=tfb.AffineScalar(shift=-0.5)),
10        low=0.,
11        high=2**bits - 1.
12    )
13    mixture_dist = tfd.MixtureSameFamily(#D
14        mixture_distribution=tfd.Categorical(logits=logits),
15        components_distribution=discretized_logistic_dist)
16    return mixture_dist
17
```

```
19 inputs = tf.keras.layers.Input(shape=(100,))
20 h1 = Dense(10, activation='tanh')(inputs)
21 out = Dense(6)(h1) #E
22 p_y = tfp.layers.DistributionLambda(quant_mixture_logistic)(out)
```

Architectures

Layer (type)	Output Shape	Param #
=====	=====	=====
input_5 (InputLayer)	[(None, 16)]	0
dense_10 (Dense)	(None, 100)	1700
dense_11 (Dense)	(None, 100)	10100
dense_12 (Dense)	(None, 10)	1010
dense_13 (Dense)	(None, 9)	99
distribution_lambda_4 (Distr ((None,), (None,)))		0
=====	=====	=====
Total params: 12,909		
Trainable params: 12,909		
Non-trainable params: 0		

This depends CPD

Poisson

(None, 1)

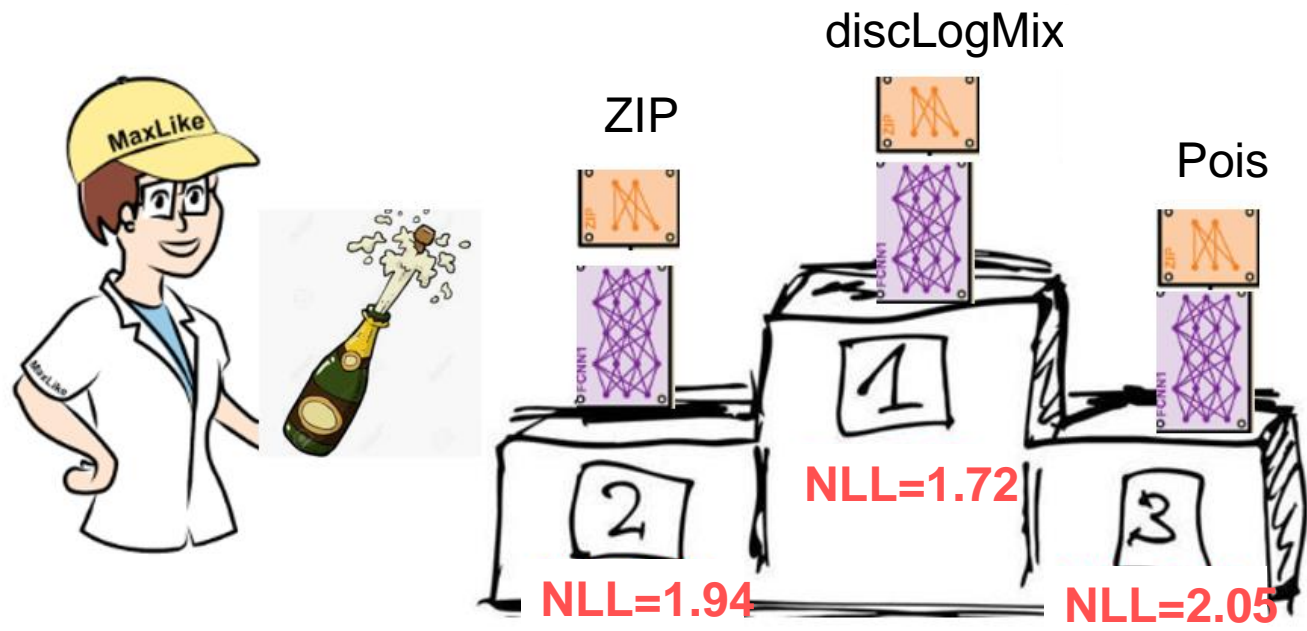
Zero Inflated?

(None, 2)

Discretized logistic mixture?

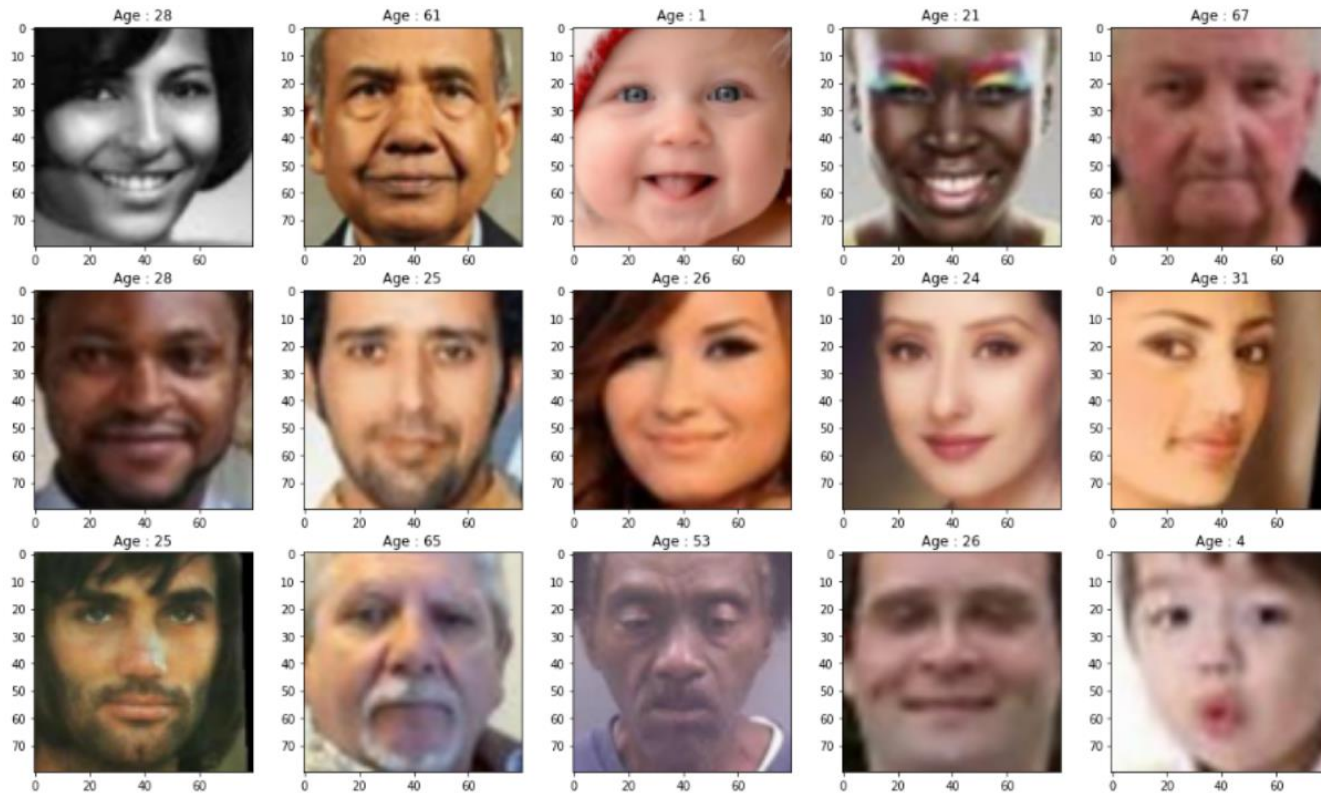
(None, 9)

Validation NLL allows to rank different probabilistic models



Probabilistic models with complex input data

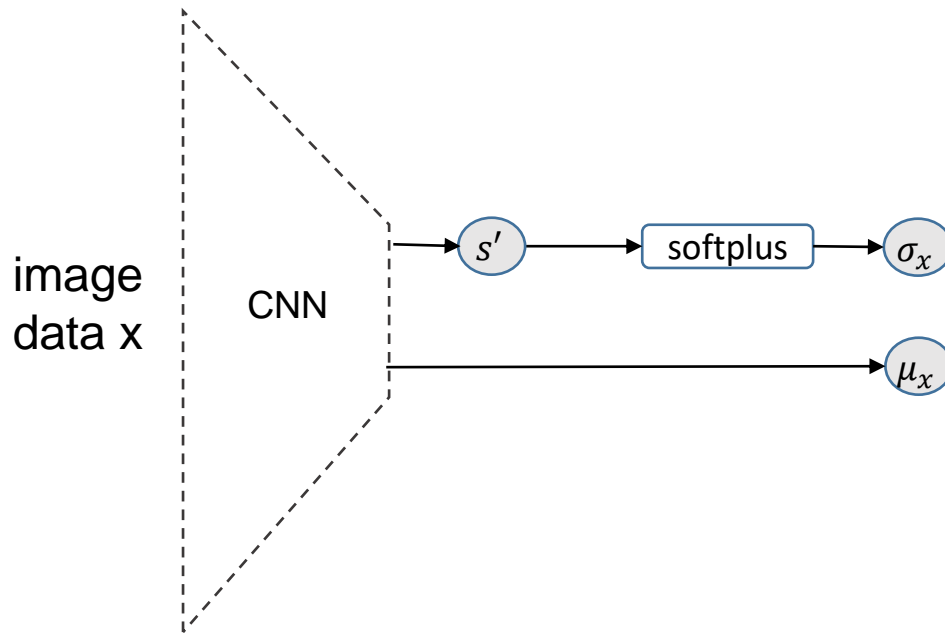
The UTK face data - face image data with known age



Data: <https://stats.idre.ucla.edu/r/dae/zip>

UTKFace data set containing N= 23'708 images of cropped faces of humans with known age ranging from 1 to 116 years.

Modeling a Gaussian CPD with flexible mean & variance



Minimize the negative log-likelihood (NLL):

$$\hat{\mathbf{w}}_{\text{ML}} = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{i=1}^n -\log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) + \frac{(y_i - \mu_{x_i})^2}{2\sigma^2}$$

CNNs for modeling Gaussian CPDs

```
def NLL(y, distr):
    return -distr.log_prob(y)

def my_dist(params):
    return tfd.Normal(loc=params[:,0:1], scale=1e-3 + tf.math.softplus(0.05 * params[:,1:2]))# both parameters are learnable

inputs = Input(shape=(80,80,3))
x = Convolution2D(16,kernel_size,padding='same',activation="relu")(inputs)
x = Convolution2D(16,kernel_size,padding='same',activation="relu")(x)
x = MaxPooling2D(pool_size=pool_size)(x)

x = Convolution2D(32,kernel_size,padding='same',activation="relu")(x)
x = Convolution2D(32,kernel_size,padding='same',activation="relu")(x)
x = MaxPooling2D(pool_size=pool_size)(x)

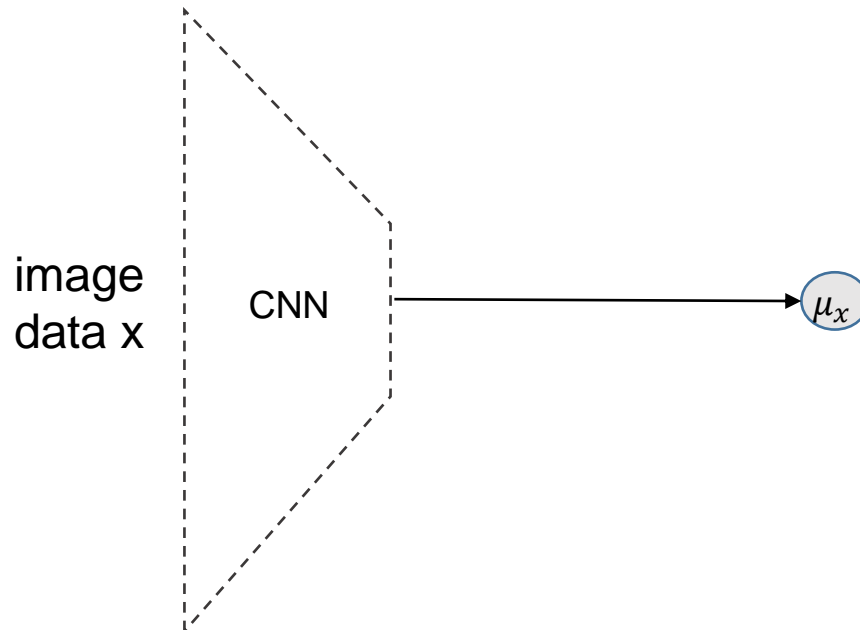
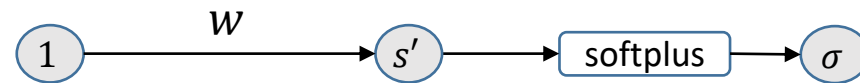
x = Convolution2D(32,kernel_size,padding='same',activation="relu")(x)
x = Convolution2D(32,kernel_size,padding='same',activation="relu")(x)
x = MaxPooling2D(pool_size=pool_size)(x)

x = Flatten()(x)
x = Dense(500,activation="relu")(x)
x = Dropout(0.3)(x)
x = Dense(50,activation="relu")(x)
x = Dropout(0.3)(x)
x = Dense(2)(x)
dist = tf.layers.DistributionLambda(my_dist)(x)

model_flex = Model(inputs=inputs, outputs=dist)
model_flex.compile(tf.keras.optimizers.Adam(), loss=NLL)
```

We control both parameters (μ_x, σ_x) of a Gaussian CPD $N(\mu_x, \sigma_x)$ by a CNN
→ More flexible than in classical regression where $\sigma = \text{constant}$

Modeling Gaussian CPD with flexible mean & constant variance



Minimize the negative log-likelihood (NLL):

$$\hat{\mathbf{w}}_{\text{ML}} = \underset{w}{\operatorname{argmin}} \sum_{i=1}^n -\log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) + \frac{(y_i - \mu_{x_i})^2}{2\sigma^2}$$

CNNs for modeling Gaussian CPDs

```
def NLL(y, distr):
    return -distr.log_prob(y)

def my_dist(params):
    return tfd.Normal(loc=params[:,0:1], scale=1e-3 + tf.math.softplus(0.05 * params[:,1:2]))# both parameters are learnable

input1 = Input(shape=(80,80,3))
input2 = Input(shape=(1,))
x = Convolution2D(16,kernel_size,padding='same',activation="relu")(input1)
x = Convolution2D(16,kernel_size,padding='same',activation="relu")(x)
x = MaxPooling2D(pool_size=pool_size)(x)

x = Convolution2D(32,kernel_size,padding='same',activation="relu")(x)
x = Convolution2D(32,kernel_size,padding='same',activation="relu")(x)
x = MaxPooling2D(pool_size=pool_size)(x)

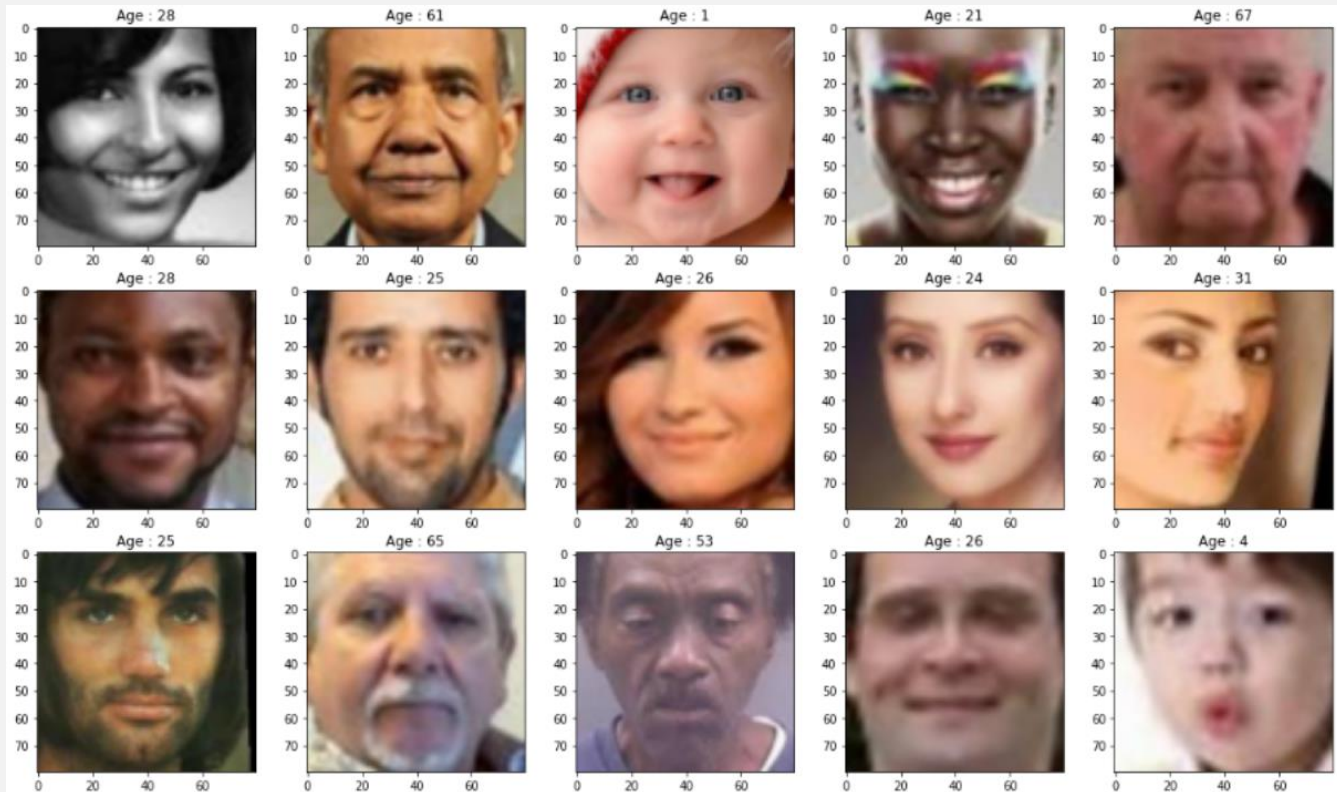
x = Convolution2D(32,kernel_size,padding='same',activation="relu")(x)
x = Convolution2D(32,kernel_size,padding='same',activation="relu")(x)
x = MaxPooling2D(pool_size=pool_size)(x)

x = Flatten()(x)
x = Dense(500,activation="relu")(x)
x = Dropout(0.3)(x)
x = Dense(50,activation="relu")(x)
x = Dropout(0.3)(x)
out1 = Dense(1)(x)
out2 = Dense(1)(input2)
params = Concatenate()([out1,out2])
dist = tfp.layers.DistributionLambda(my_dist)(params) #

model_const_sd = Model(inputs=[input1,input2], outputs=dist) ## use a trick with two inputs, input2 is just ones
model_const_sd.compile(tf.keras.optimizers.Adam(), loss=NLL)
```

We control both parameters (μ, σ) of a Gaussian CPD $N(\mu_x, \sigma) \rightarrow$ But assume a constant variance

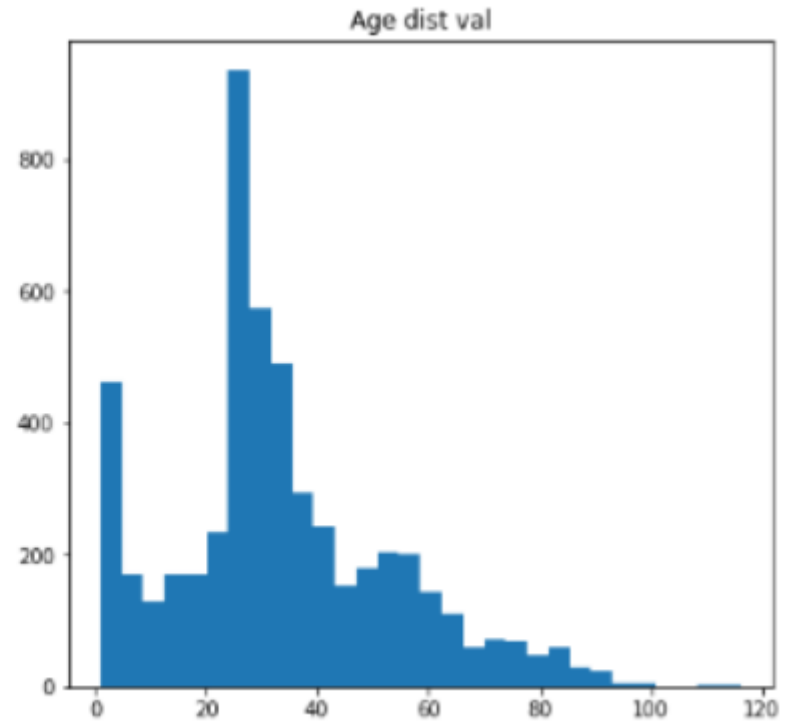
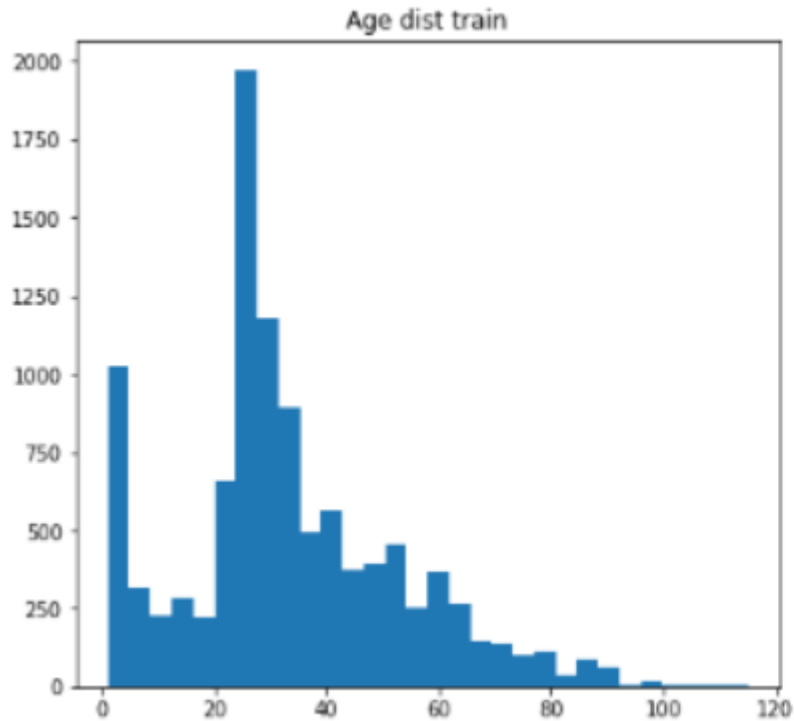
Exercise: bis um 16:00



Check out the models for age prediction with flexible and constant variance and try to understand the code and to answer the questions.

https://github.com/tensorchiefs/dl_course_2022/blob/master/notebooks/17_faces_regression.ipynb

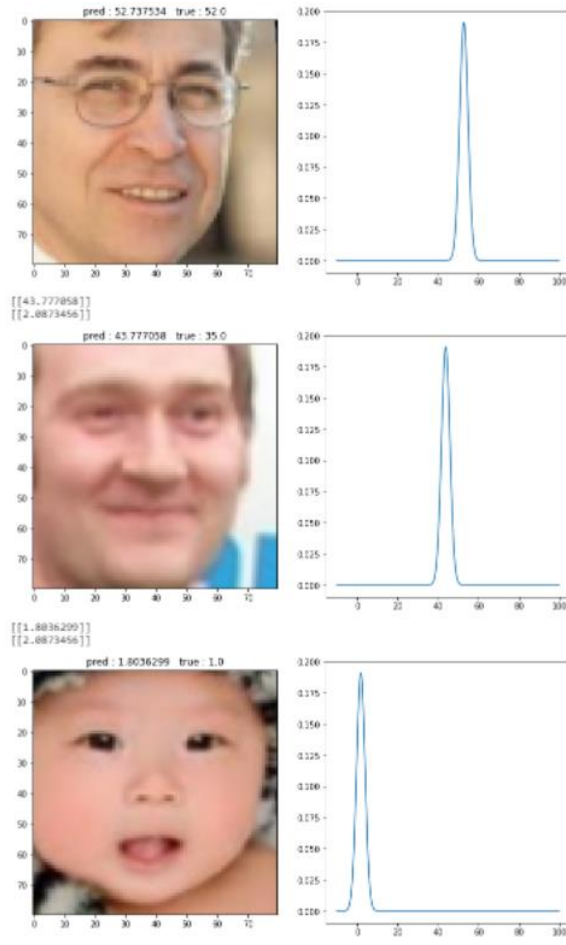
Age distribution



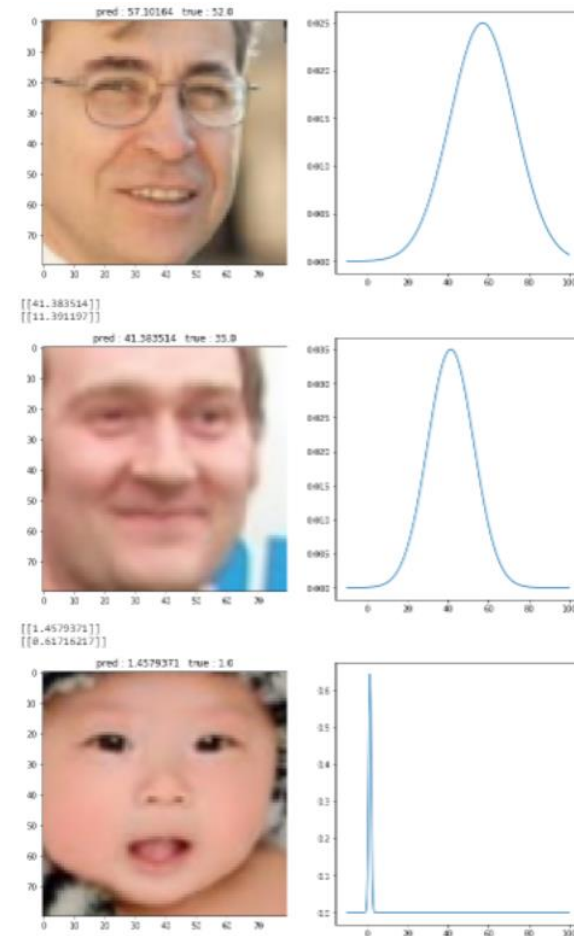
We have a lot of small children in the data set, for whom the age estimation is probably not so difficult.

Resulting age CPDs

Constant variance



flexible variance



In case of a flexible variance, a broad predicted Gaussian CPD does indicate high uncertainty about the age.

Where does a CNN look at?

Visualize patches yielding high values in activation maps

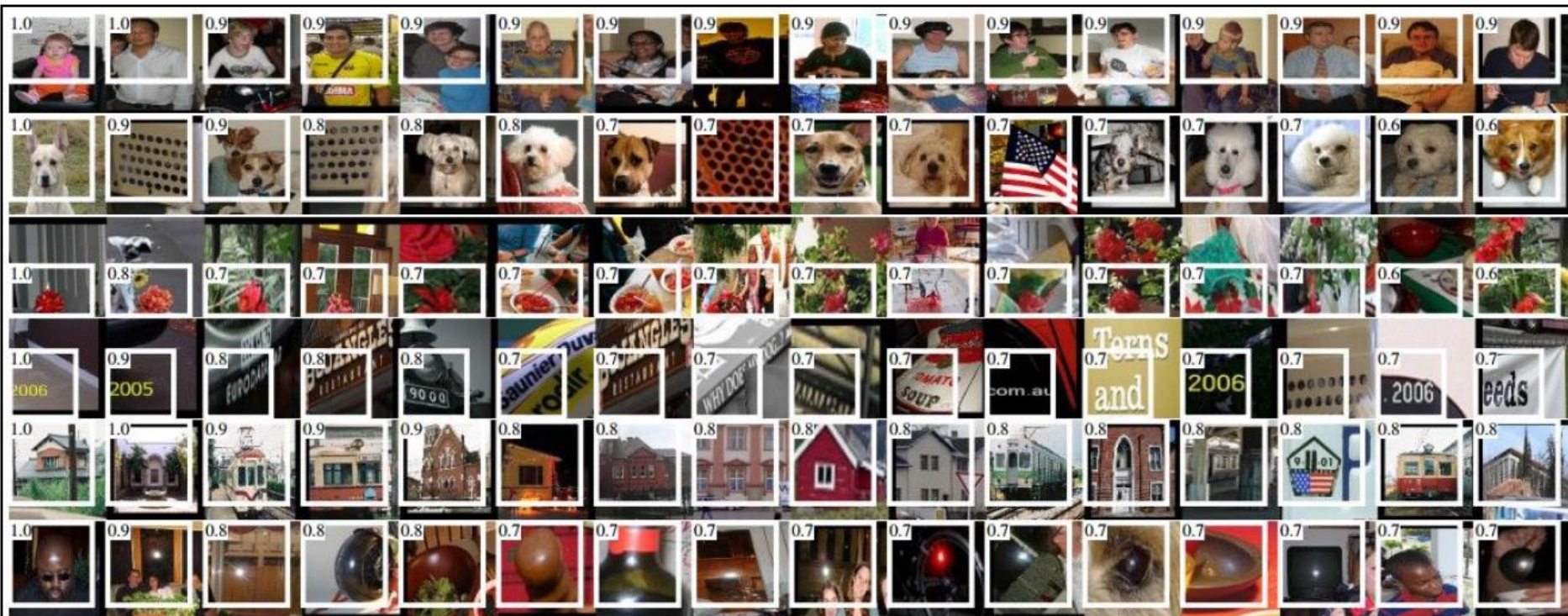
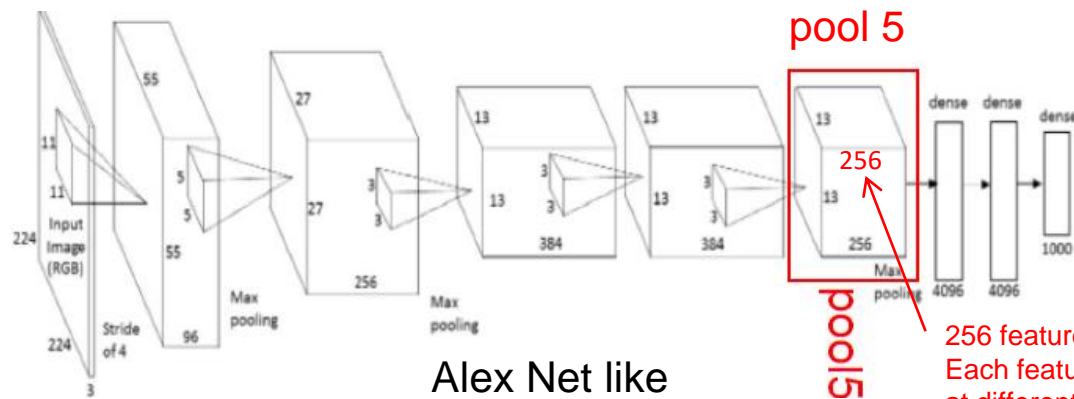


Figure 4: Top regions for six pool_5 units. Receptive fields and activation values are drawn in white. Some units are aligned to concepts, such as people (row 1) or text (4). Other units capture texture and material properties, such as dot arrays (2) and specular reflections (6).



<http://cs231n.github.io/understanding-cnn/>

Here we show image patches that activate maps in layer 5 most.

256 feature maps generated by 256 different 3x3 filters. Each feature map consists of equivalent neurons looking at different positions of the input volume.

What kind of image (patches) excites a certain neuron corresponding to a large activation in a feature map?

10 images from data set leading to high signals 6 feature maps of **conv6**



10 images from data set leading to high signals 6 feature maps of **conv9**

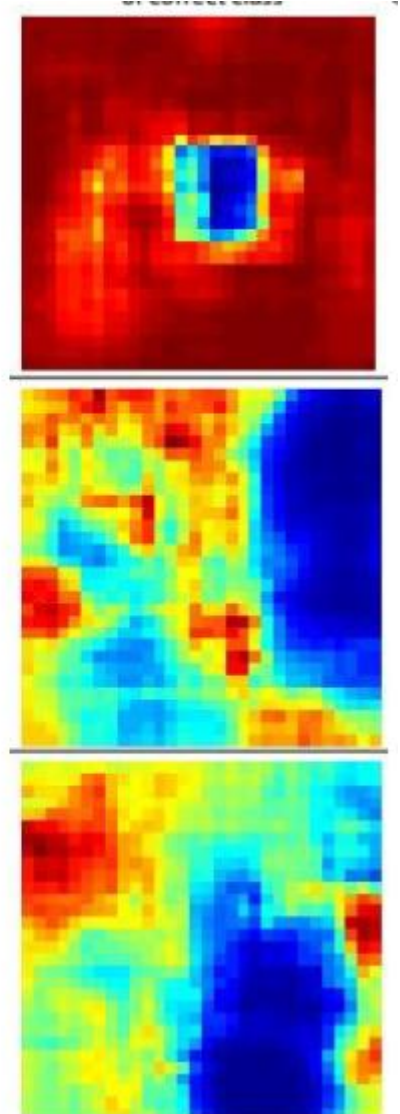


Which pixels are important for the classification?

Occlusion experiments

Occlude part of the image with a mask and check for each position of the mask how strongly the score for the correct class is changing.

Warning:
Usefulness depends on application...



Occlusion experiments [\[Zeiler & Fergus 2013\]](#)

Which pixels are important for the classification?

LIME: Local Interpretable Model-agnostic Explanations

Idea:

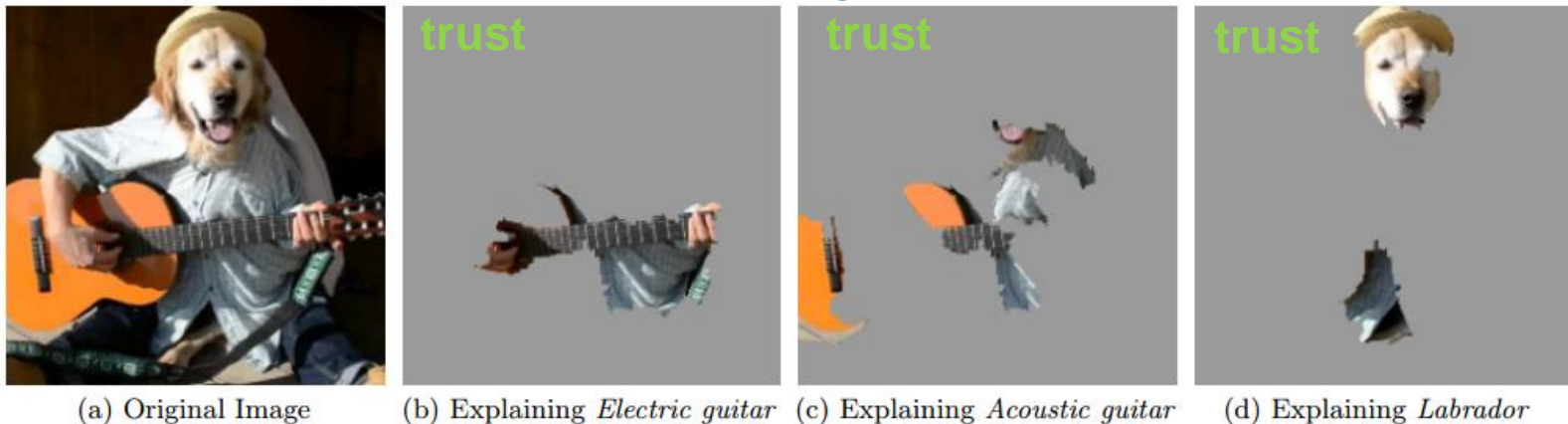
- 1) perturb interpretable features of the instance – e.g. randomly delete super-pixels in an image and track as perturbation vector such as $(0,1,1,0,\dots,1)=x$.
- 2) Classify perturbed instance by your model, here a CNN, and track the achieved classification-score= y
- 3) Identify for which features/super-pixels the presence in the perturbed input version are important to get a high classification score (use RF or lasso for $y \sim x$)

<https://arxiv.org/abs/1602.04938>



-> presence of snow was used to distinguish wolf and husky

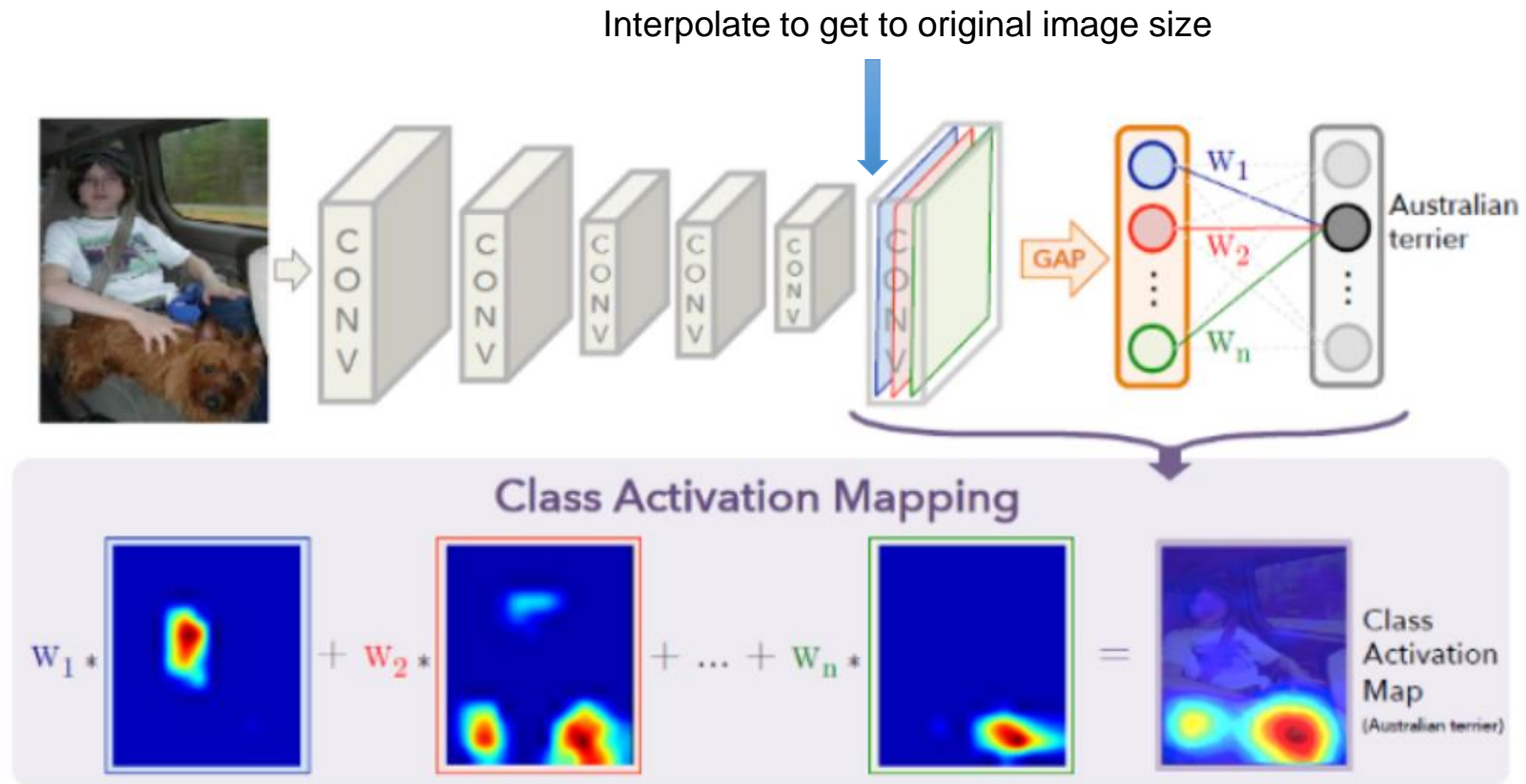
-> Explain the CNN classification by showing instance-specific important features
visualize important feature allows to judge the individual classification



Tf, python code: <https://github.com/marcotcr/lime> -> image tutorial

<https://github.com/marcotcr/lime/blob/master/doc/notebooks/Tutorial%20-%20Image%20Classification%20Keras.ipynb>

Class Activation Mapping (CAM)

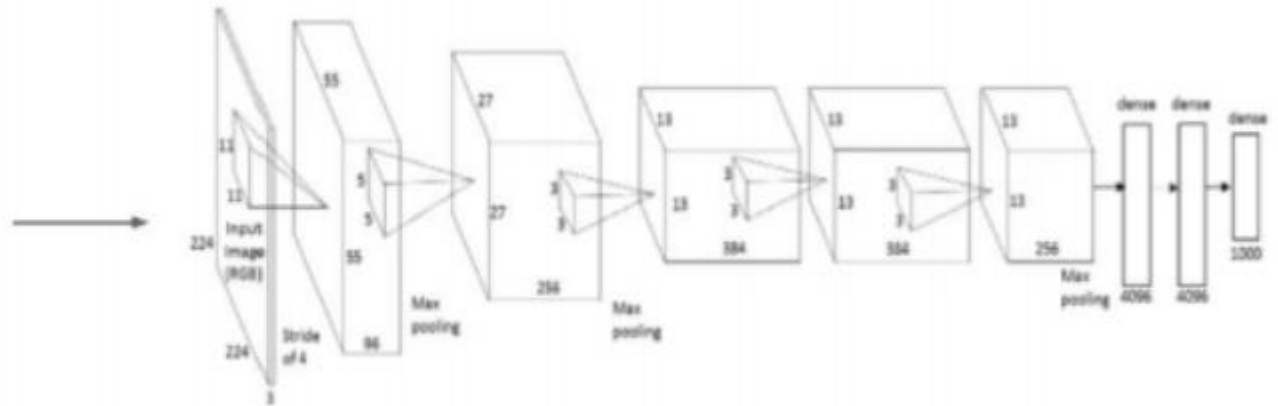


We compute the CAM by multiplying each feature map of the final convolution layer with the corresponding weight connected to the neuron of the winning class.

<https://arxiv.org/abs/1512.04150>

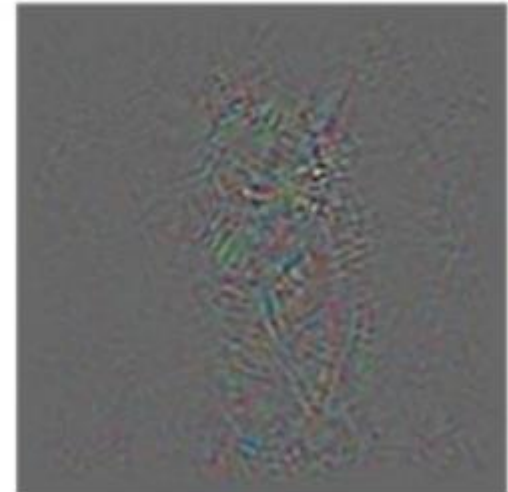
<https://towardsdatascience.com/class-activation-mapping-using-transfer-learning-of-resnet50-e8ca7cfd657e>

Gradient Backpropagation



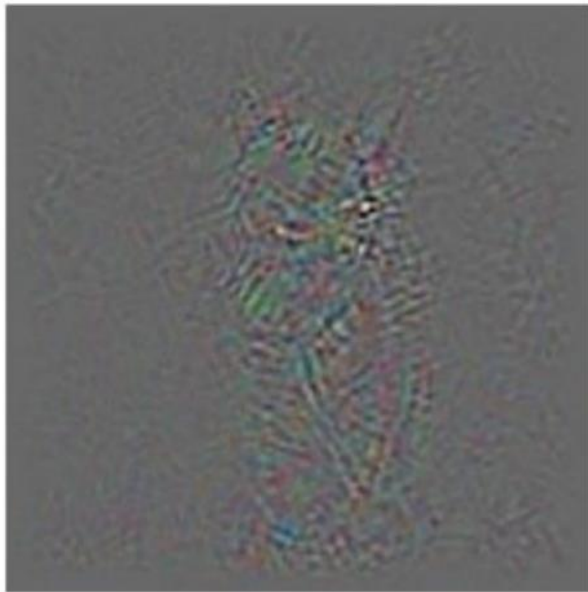
- 1) Do a forward pass with the image
- 2) Compute the gradient of the winning class neuron using backprop

Which pixels would make the class neuron not turn on if they had been different?
In other words, for which inputs is $\partial(\text{class} - \text{neuron})/\partial x_i$ large?



Guided Gradient Backpropagation

- We are **only interested to see which image features the neuron detects**, not in what kind of stuff it doesn't detect
- So **only propagating positive gradients** (set all the negative gradients to 0)



Backprop



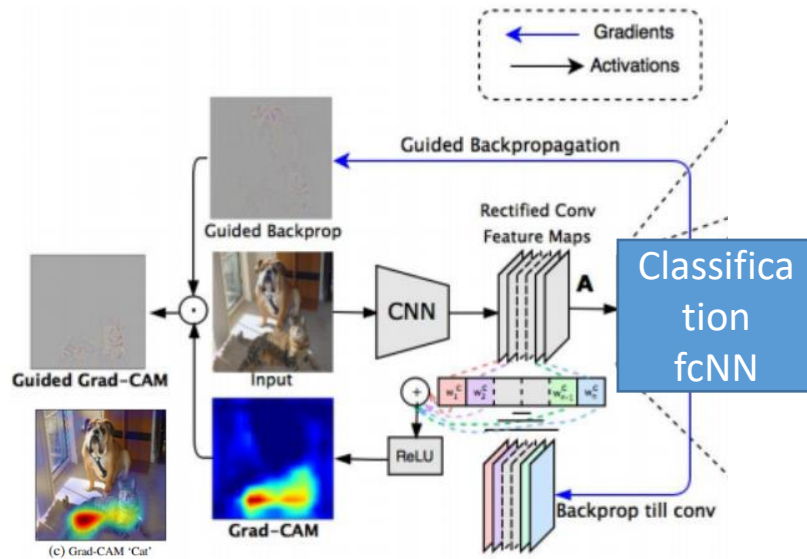
Guided Backprop

<https://www.cs.toronto.edu/~guerzhoy/321/lec/W07/HowConvNetsSee.pdf>

Code: see iNNvestigate colab NB https://github.com/avciidp/explainable_ai/blob/master/iNNvestigate_fashion_light.ipynb
paper: <https://arxiv.org/abs/1808.04260>

Grad-Cam: Gradient based Class Activation maps

Grad-Cam: Use Guided Backprop to scale contributions in CAMs



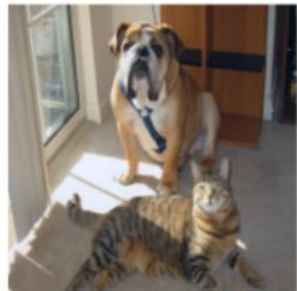
- Set gradient to zero for all classes except the predicted class
- Backpropagate this gradient to the last conv-layer
-> get coarse grad-CAM localization
- Multiply coarse grad-CAM with the guided back-prop signal

<https://arxiv.org/abs/1610.02391>

Code: https://keras.io/examples/vision/grad_cam/

<https://www.pyimagesearch.com/2020/03/09/grad-cam-visualize-class-activation-maps-with-keras-tensorflow-and-deep-learning/>

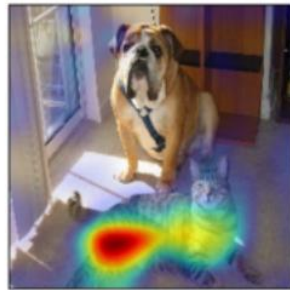
What in the image was important for the prediction?



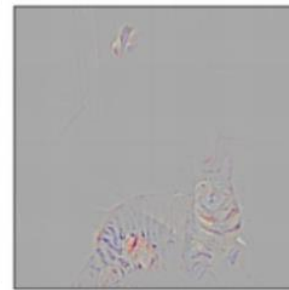
(a) Original Image



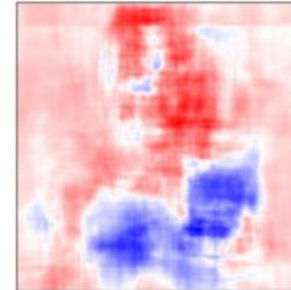
(b) Guided Backprop 'Cat'



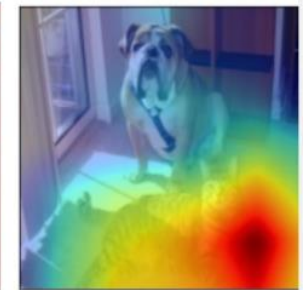
(c) Grad-CAM 'Cat'



(d) Guided Grad-CAM 'Cat'



(e) Occlusion map 'Cat'



(f) ResNet Grad-CAM 'Cat'



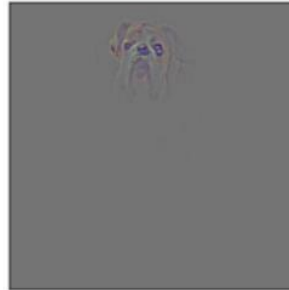
(g) Original Image



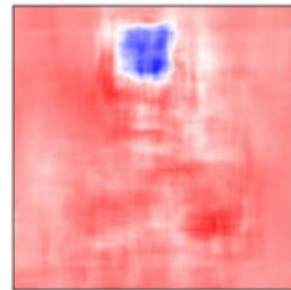
(h) Guided Backprop 'Dog'



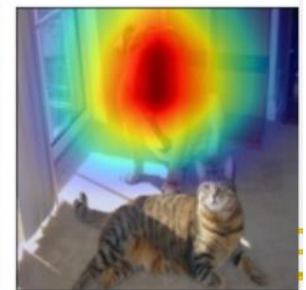
(i) Grad-CAM 'Dog'



(j) Guided Grad-CAM 'Dog'

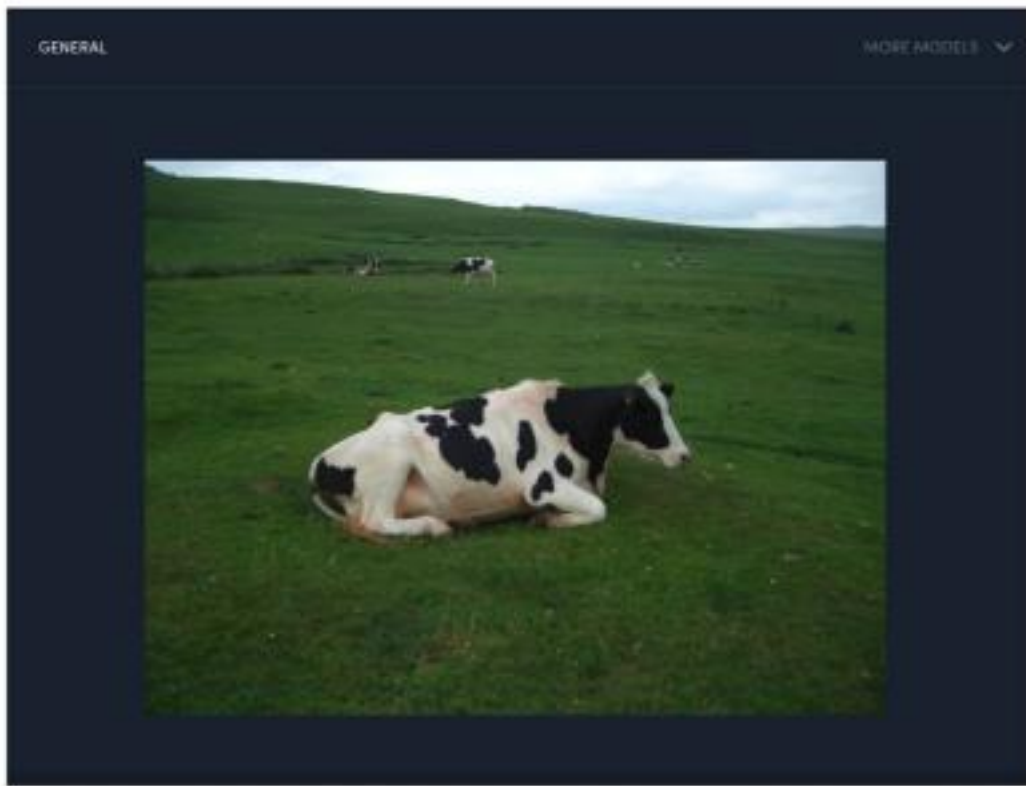


(k) Occlusion map 'Dog'



(l) ResNet Grad-CAM 'Dog'

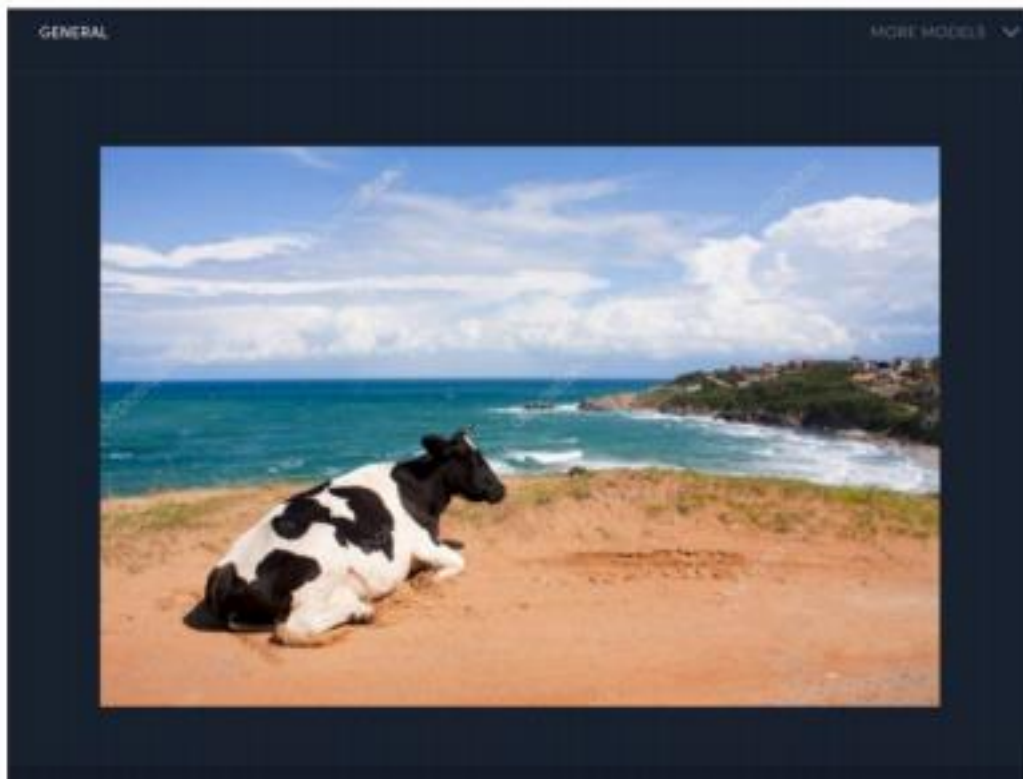
Typisches Beispiel einer Kuh



General	VIEW DOCS
cow	0.992
cattle	0.983
mammal	0.979
grass	0.978
livestock	0.966
farm	0.964
landscape	0.963
pasture	0.954
grassland	0.949
agriculture	0.948
no person	0.945

Untypisches Beispiel einer Kuh

Dieses Bild ist zu weit weg von den Trainingsdaten. Kuh nicht unabhängig von Hintergrund




Neuronales Netz erkennt keine Kuh!

General	VIEW DOCS
no person	0.991
beach	0.990
water	0.985
sand	0.981
sea	0.980
travel	0.978
seashore	0.972
summer	0.954
sky	0.946
outdoors	0.944
ocean	0.936

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

If time allows:

- Issues with current DL approach
 - No uncertainty for the fitted weights
 - if algorithmic and epistemic uncertainty is ignored :
 - No increased uncertainty in case of extrapolation
 - Deficits in prediction performance
- Approaches to take algorithmic and epistemic uncertainty into account:
 - Deep Ensembling
 - MC Dropout

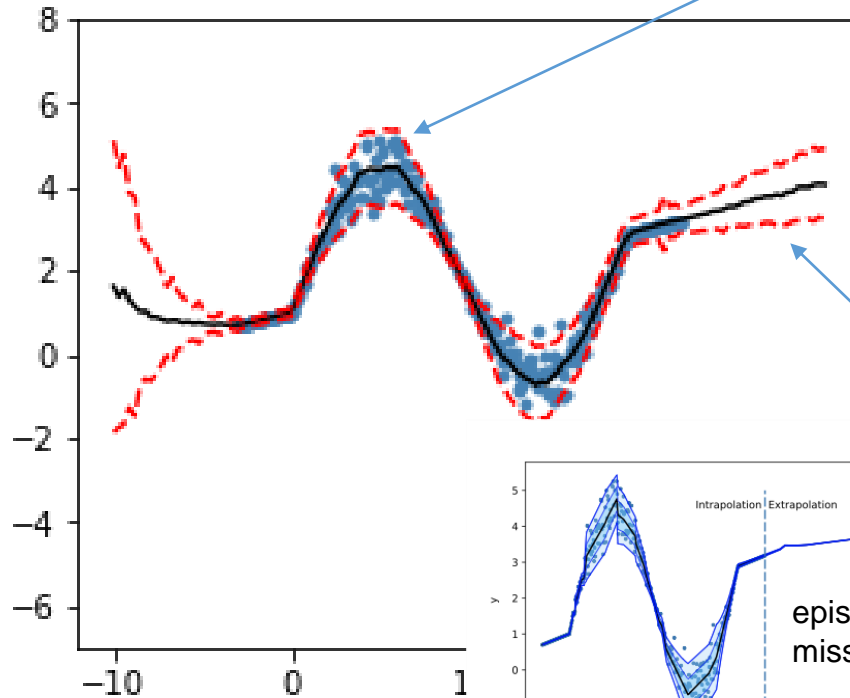
Aleatoric vs. Epistemic Uncertainty

Much spread in train data → **aleatoric uncertainty** (from latin "[Alea Acta est](#)")

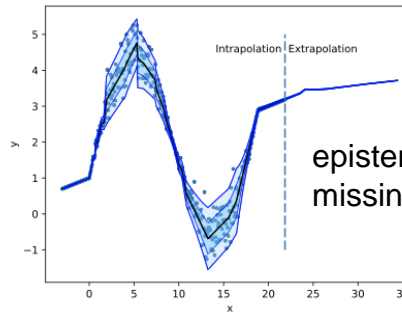
Caesar: Die Würfel sind gefallen!



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.



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

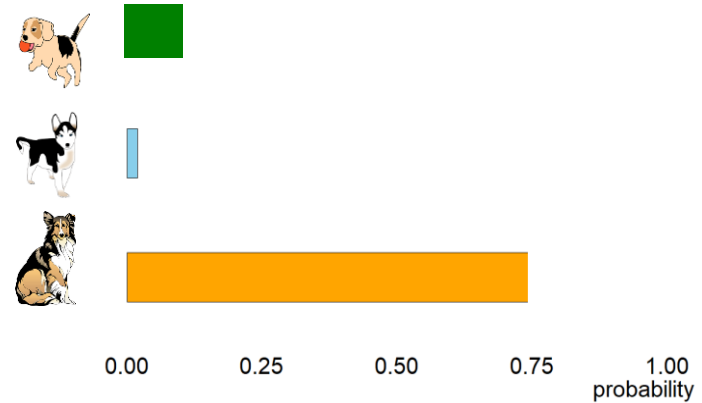
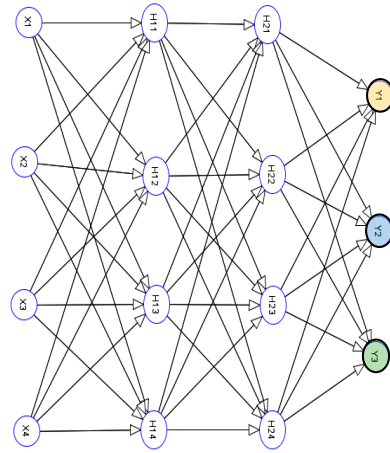


epistemic uncertainty is missing in classical NN



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

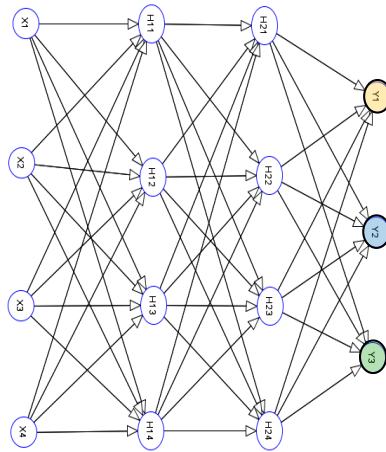
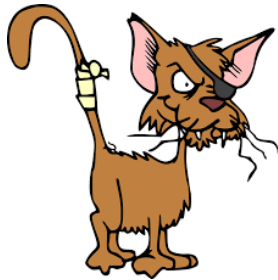
CNNs have high performance on in-distribution examples



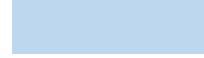
CNNs yield high accuracy and calibrated probabilities, but...

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?

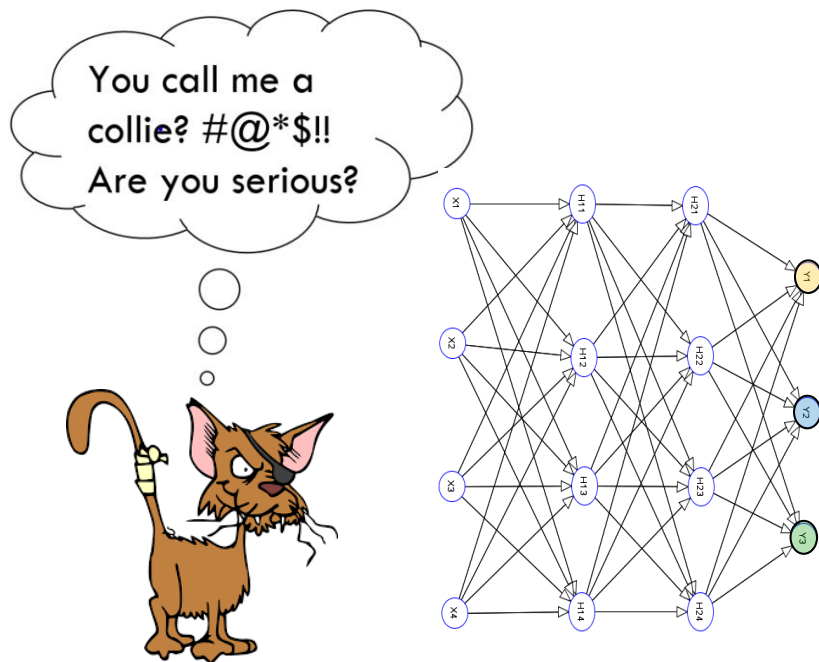


You might expect:



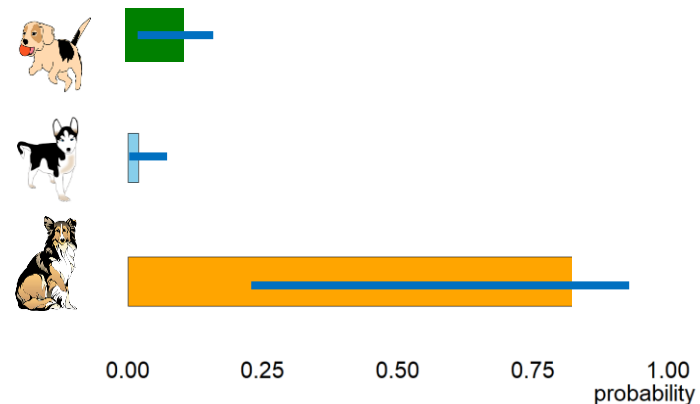
A classical NN cannot ring the alarm in case of out-of-distribution (OOD) examples

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)



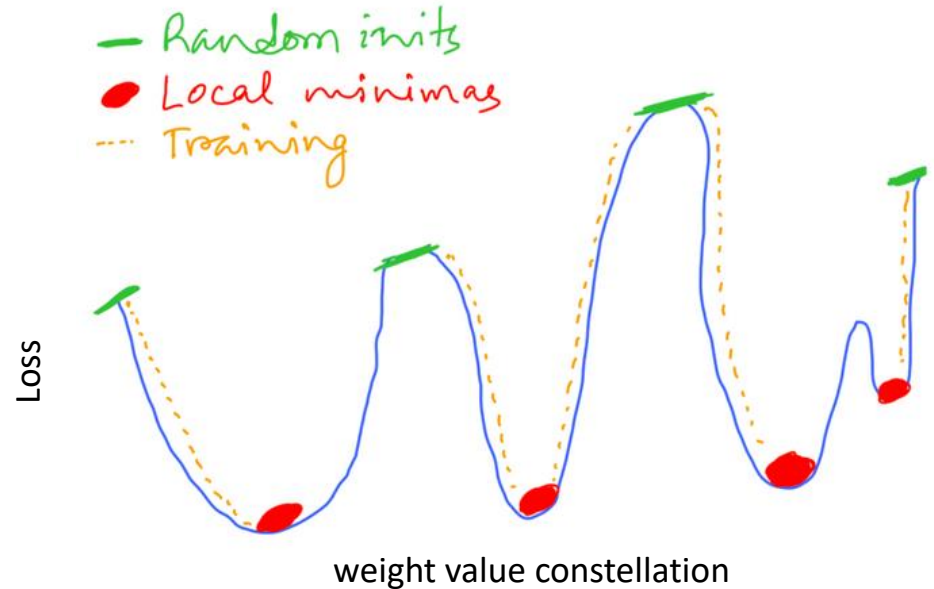
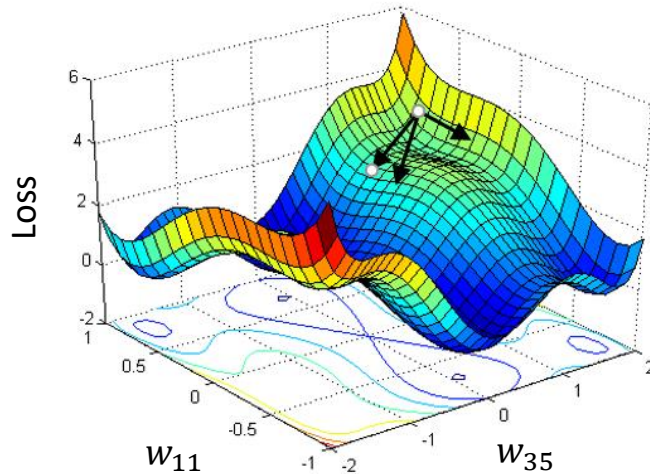
- Current DL Systems bad in out of distribution OOD situations
- Application need at least to detect OOD situations

Deep ensembles

→ addressing algorithmic uncertainty

The loss-landscape in DL is usually not convex

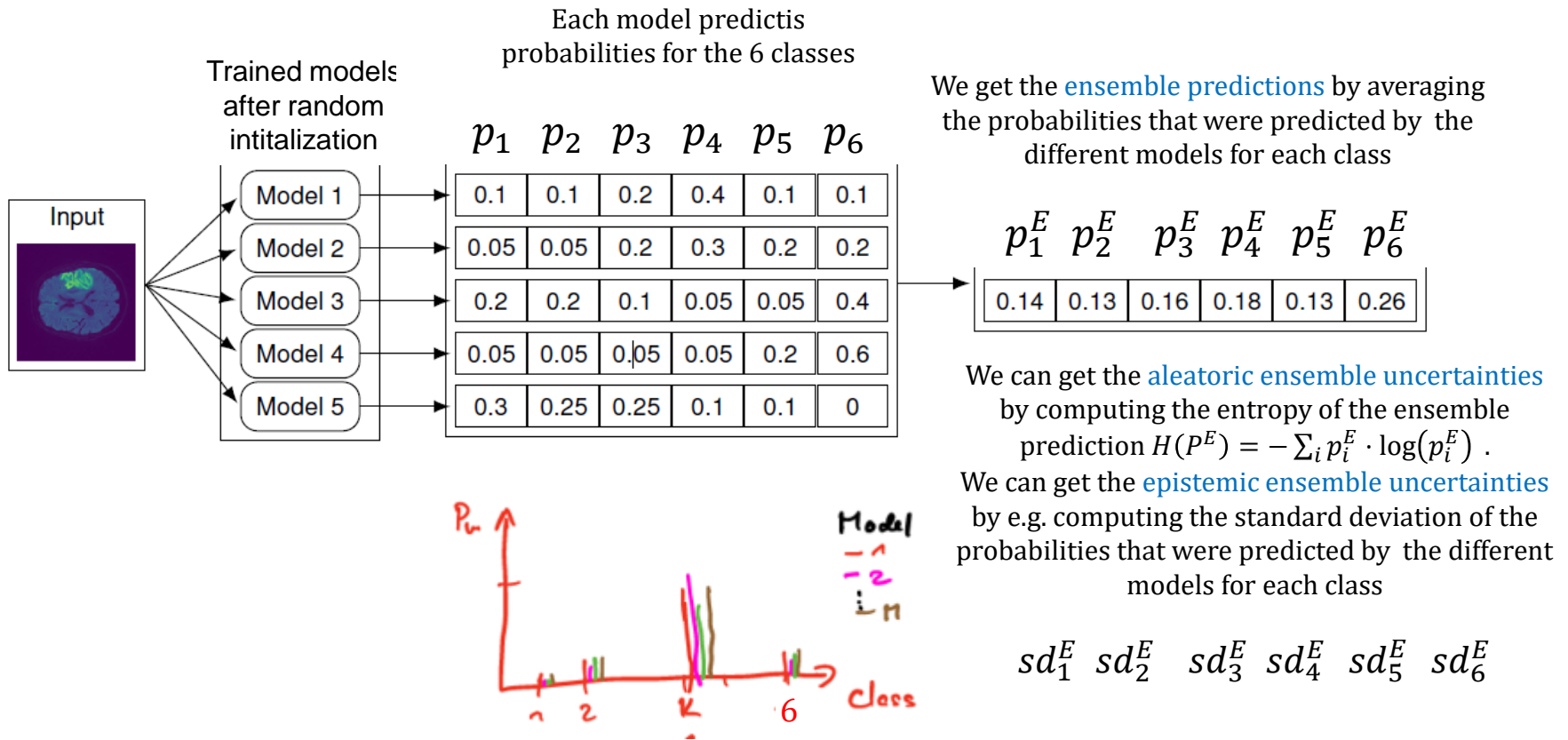
Loss-landscape



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 SGD and 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.

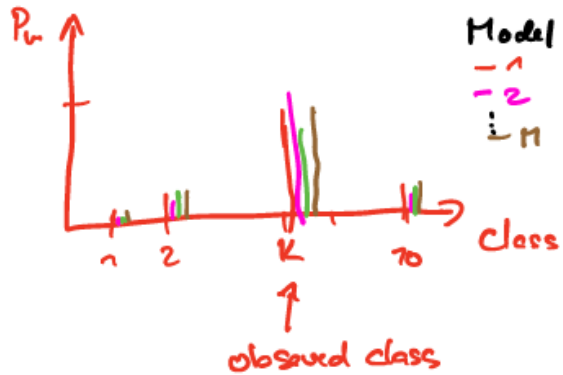
Ensembling improves the NLL performance

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

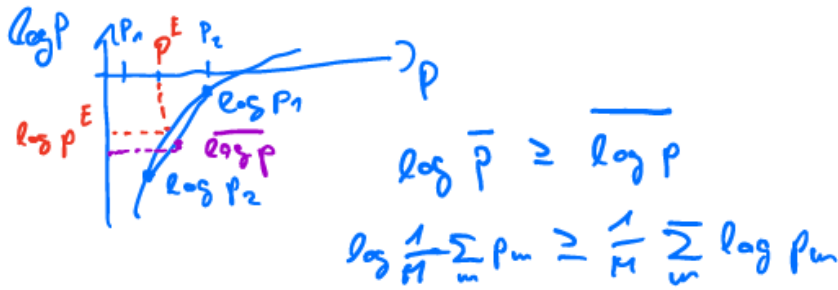
with $\bar{l} = \frac{1}{M} \sum_{m=1}^M l_m = \frac{1}{M} \sum_{m=1}^M -\log P_{km}$

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

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

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

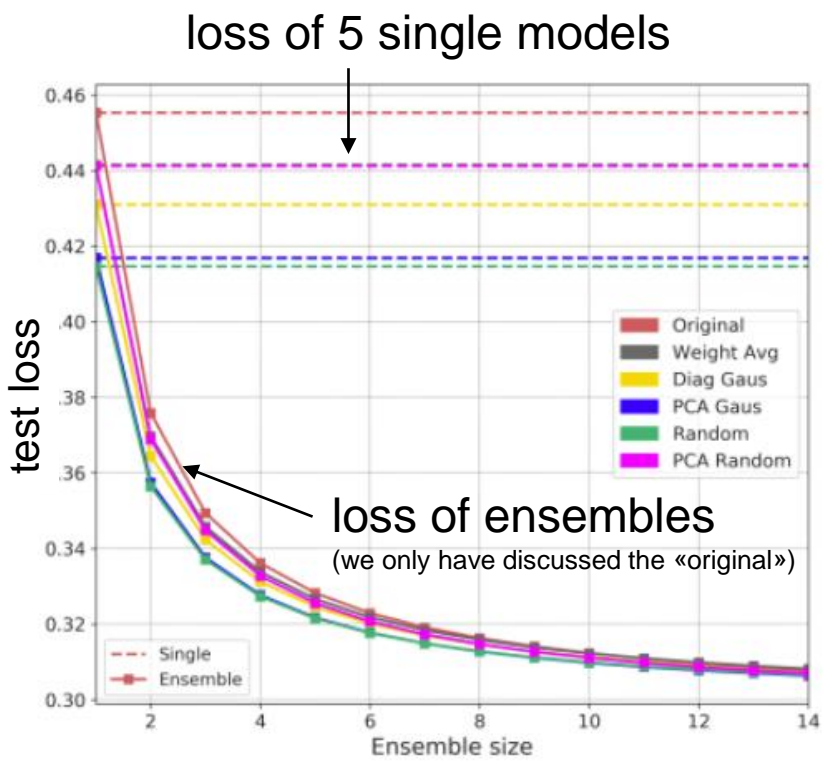
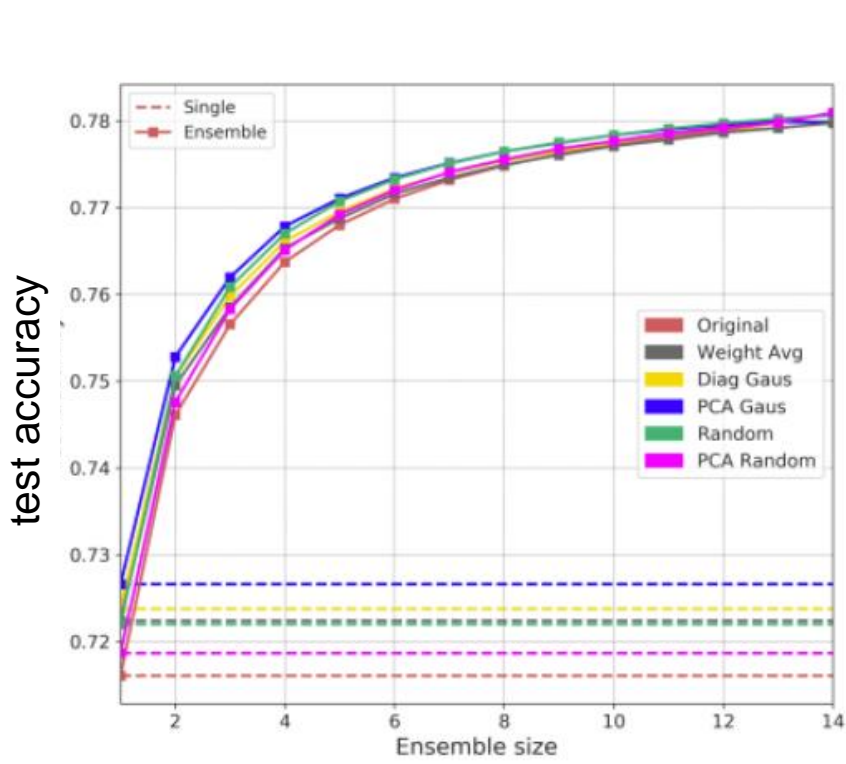
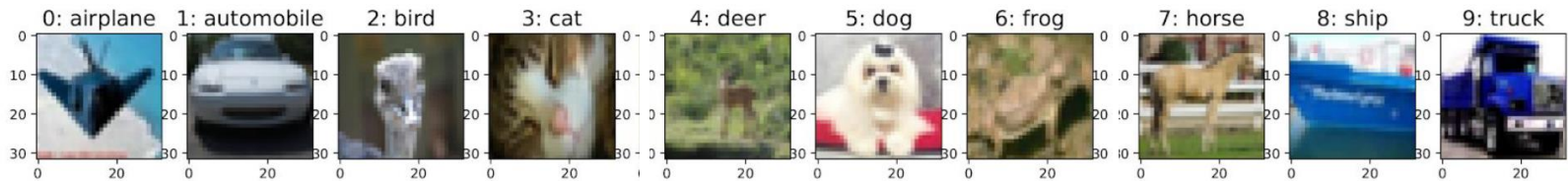
use Jensen inequality



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

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

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.

