

WBL Deep Learning:: Lecture 4

Beate Sick, Oliver Dürr

Deep Learning as probabilistic modeling

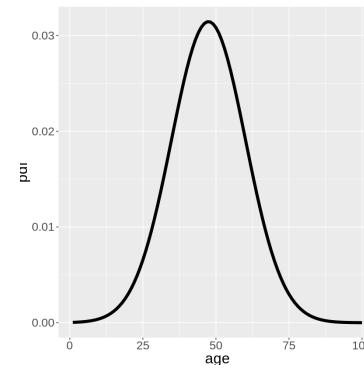
Zürich, 9/26/2022

At the end of today [probabilistic age estimation]

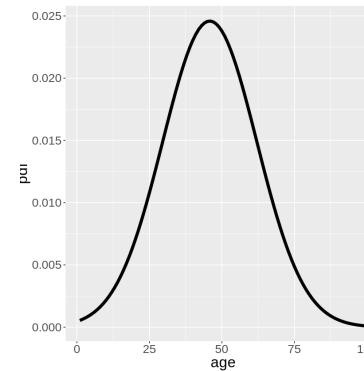
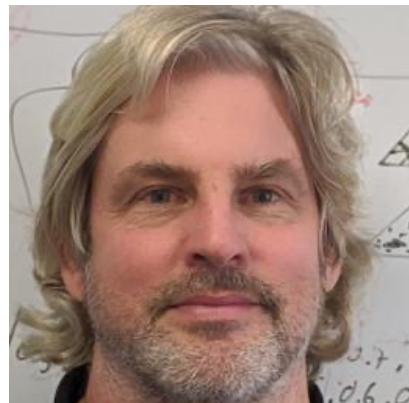
x



$p(y|x)$



NN



Outline: Probabilistic Modeling (Likelihood based)

- What is a probabilistic model?
 - Predicts whole conditional probability distribution (CPD)
- How to fit a probabilistic DL model?
 - use NLL as loss function
- How to evaluate a probabilistic DL model?
 - use NLL as loss function
- Examples using TensorFlow Probability
 - Regression via fcNN (CPD: $N(\mu, \sigma)$)
 - Count data via fcNN (CPD: Poisson, NB, ZIP)
 - Regression for Age estimation via CNN (CPD: $N(\mu, \sigma)$)

Why is it important to know about probabilities?

Philosophical reasons:

“The true logic of this world is in the calculus of probabilities”

James Clerk Maxwell

“It is scientific to say what is more likely and what is less likely...”

Richard Feynman

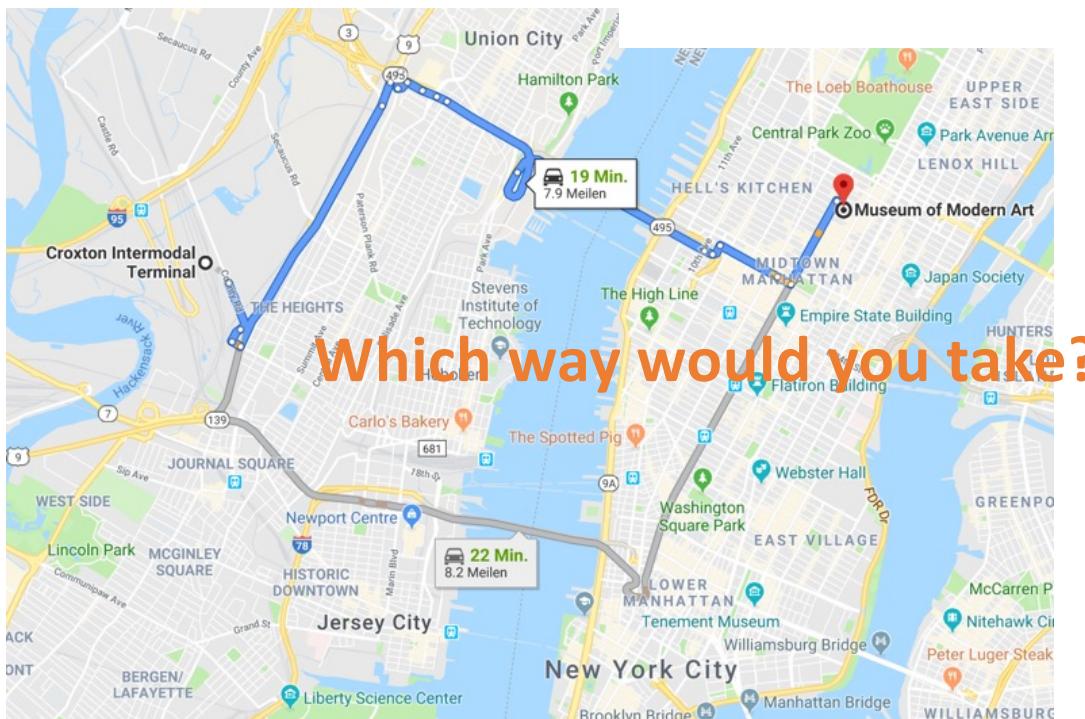
Practical reasons:

- We often want to optimize expected costs which requires CPD for computing.
- **They come for free in most DL problems**
- Choosing the right probabilistic model enhances existing DL architectures
 - Count data (using Poisson, NB, or Poisson)
 - Parallel Wavenet, PixelCNN++ (using mixtures of logistics distributions as output)

Probabilistic travel time prediction (expected cost)

You'll get 500\$ tip if I arrive at MOMA within 25 minutes!

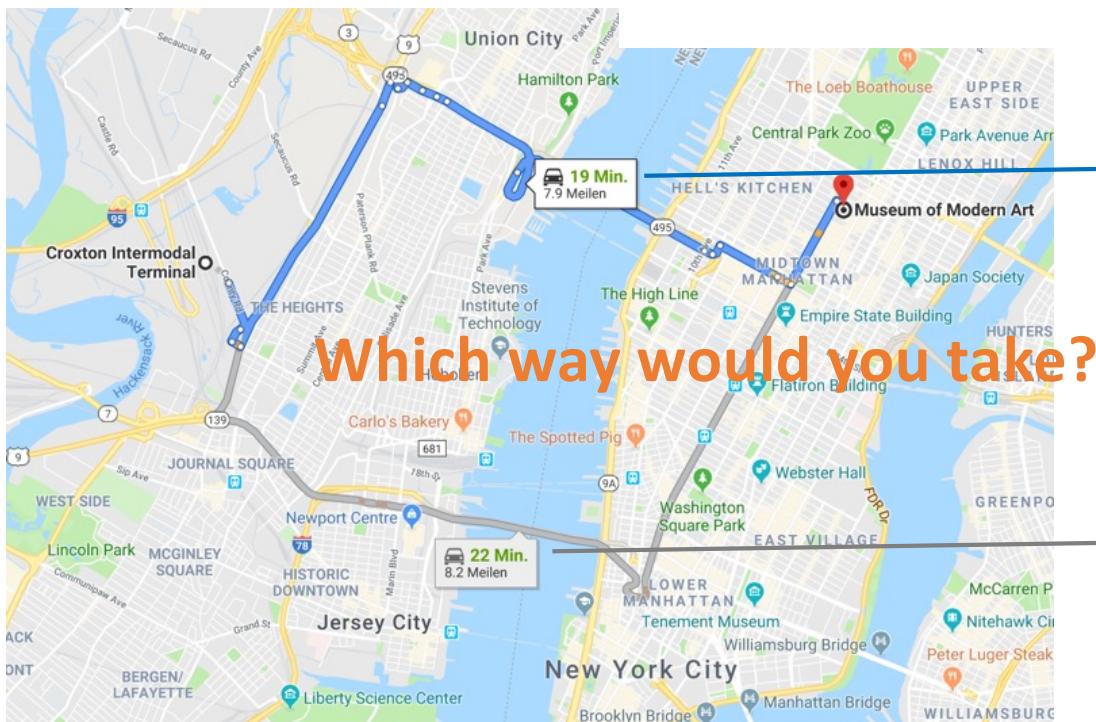
Let's use my probabilistic travel time gadget!



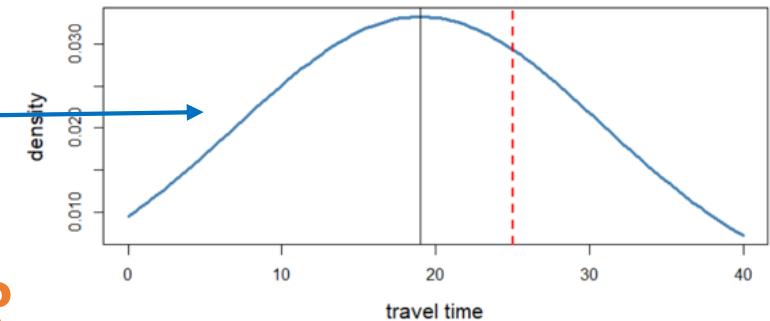
Probabilistic travel time prediction (expected cost)

You'll get 500\$ tip if I arrive at MOMA within 25 minutes!

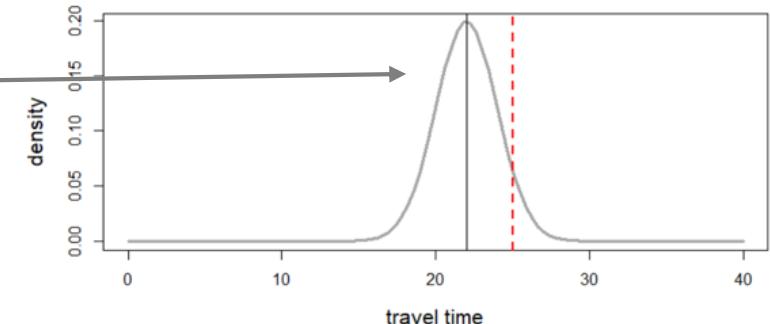
Let's use my probabilistic travel time gadget!



Chance to get tip: 69%

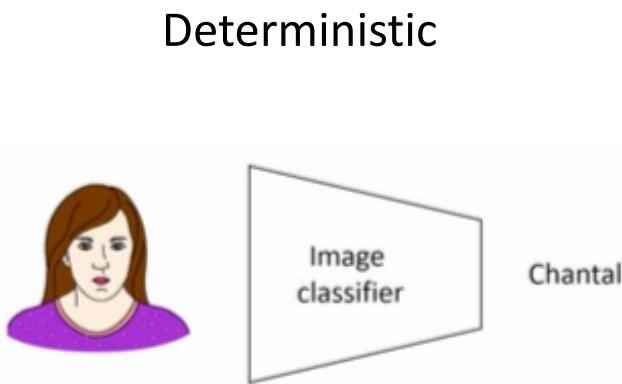


Chance to get tip: 93%

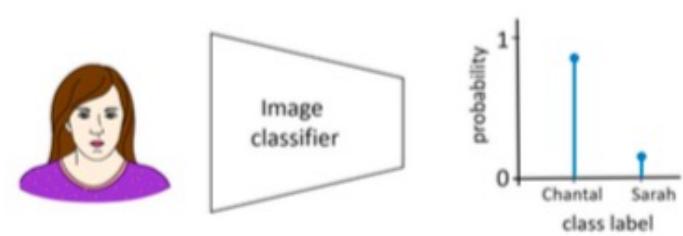


Recap Deterministic / Probabilistic

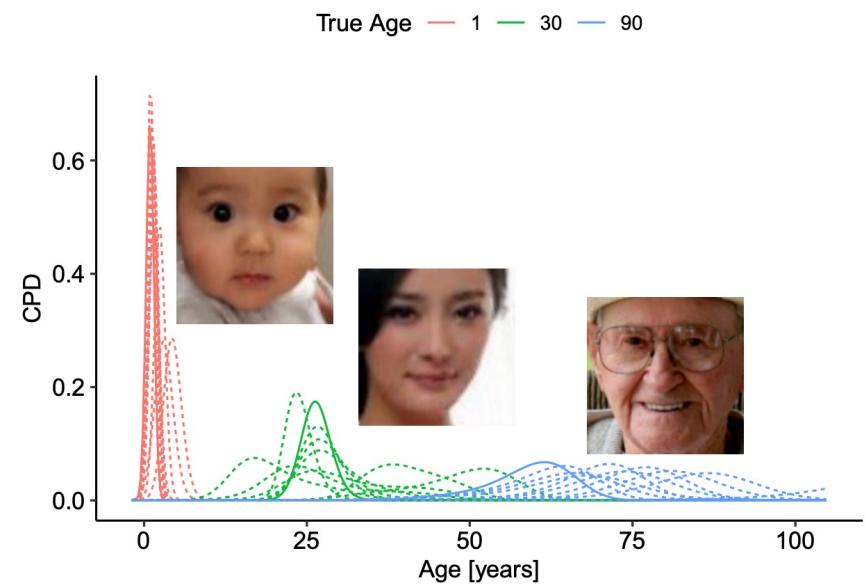
“Classification”



Probabilistic



“Regression”



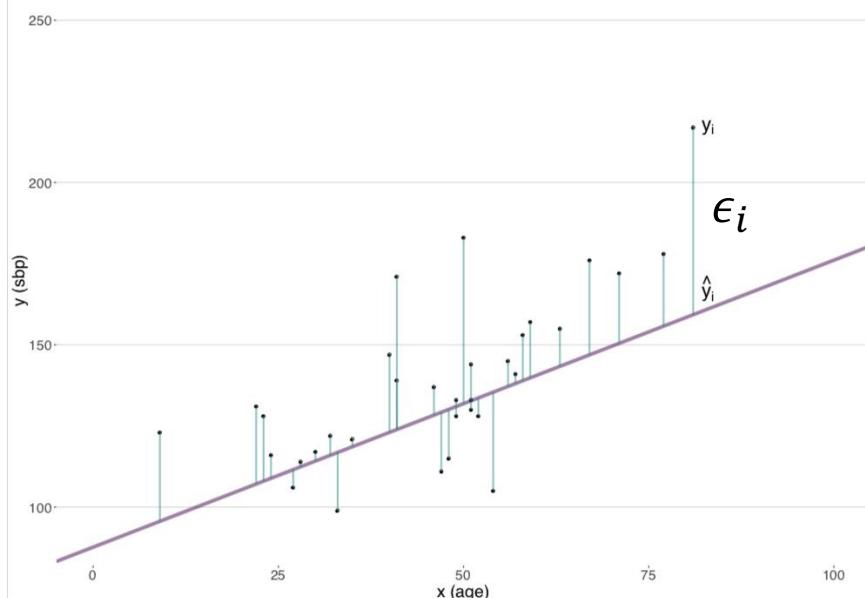
Conditional probability distribution (CPD)
 $p(y|x)$

Probabilistic models

Traditional Formulation of Regression



- Suppose we observe y_i and $x_i = (x_{i1}, \dots, x_{ip})$ for $i = 1, \dots, n$
- We can model the relationship as
 - $Y_i = f(X_i) + \epsilon_i$
 - Where ϵ_i is a random error with mean zero.



The ϵ_i are the difference to the “true values”.

This yields

$$E[Y_i] = E[Y_i|X_i] = f(X_i) \quad (\text{expression which is useful again})$$

Do not eat this poisonous fruit

Critique on the formulation



Xiao-Li Meng
@XiaoLiMeng1



The term “regression” reflects statisticians’ modesty, and perhaps also our regrets? We should not have started statistical modeling with regression, for it confuses probabilistic model fitting with deterministic line/curve fitting, building wrong intuitions for generations.

7:48 am · 1 Oct 2019 · Twitter for iPhone

Source: twitter.com/XiaoLiMeng1

Modern Formulation of regression (distributional)

- Suppose we observe y_i and $x_i = (x_{i1}, \dots, x_{ip})$ for $i = 1, \dots, n$
- We believe that there is a relationship between Y and at least one of the X's.
- We can model the relationship by an outcome distribution (dist) of Y depending on X:

$$Y_i \sim \text{dist}(z(x_i))$$

- Where z is an unknown function “controlling” the parameters of an outcome distribution dist.
- Examples
 - $\text{dist} = N(\mu = z_1(x), \sigma = z_2(x))$
 - $\text{dist} = N(\mu = z(x), \sigma)$ σ fixed does not depend on x.
 - $\text{dist} = \text{Exp}(\lambda = z(x))$

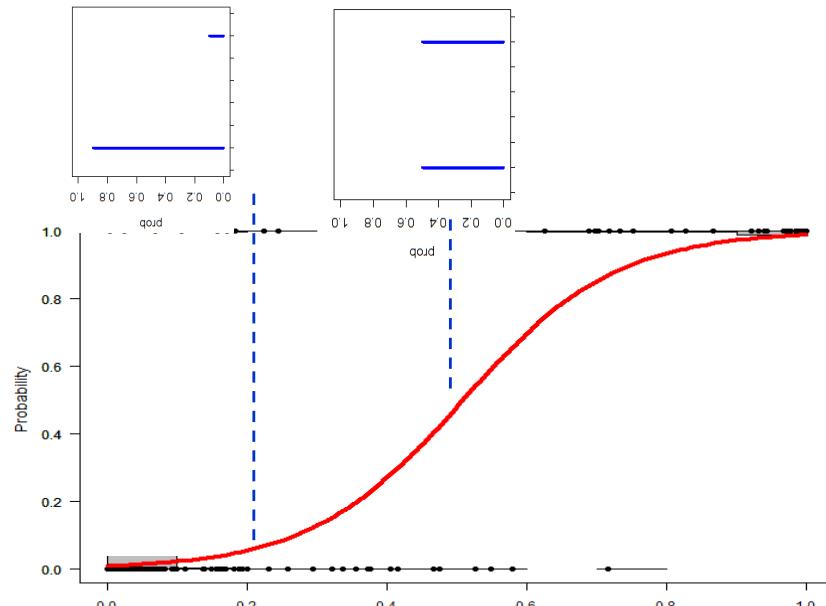
GLM examples

Logistic Regression

$$z(x) = a \cdot x + b$$

$$Y \sim \text{Bern}(p = \text{sig}(z(x)))$$

```
glm(y ~ ., binomial(logit))
```

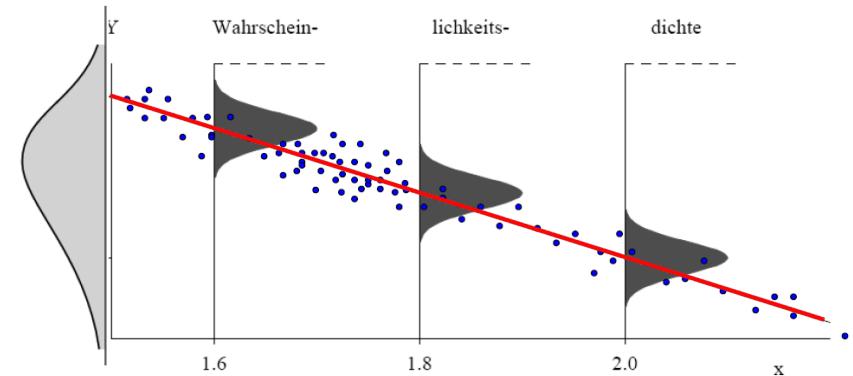


Linear Regression

$$z(x) = a \cdot x + b$$

$$Y \sim \text{Normal}(\mu = z(x), \sigma = \text{const})$$

```
glm(y ~ ., gaussian(identity))
```



Linear regression is confusingly special case
(constant variance, identity link)

GLM examples

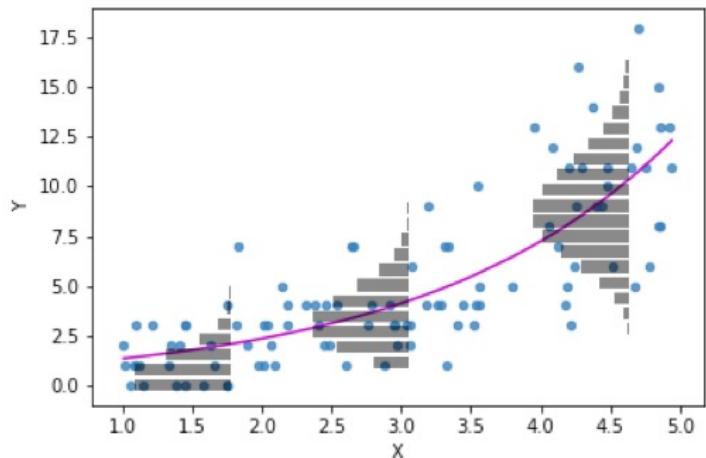


Poisson Regression

$$z(x) = a \cdot x + b$$

$$Y \sim Pois(\text{rate} = \text{Exp}(z(x)))$$

```
glm(y ~ ., poisson(log))
```



GLM

1. Linear Predictor

$$z(x) = \eta(x) = \beta_0 + \sum \beta_i x_i = \beta x^T$$

2. Inverse Link Function

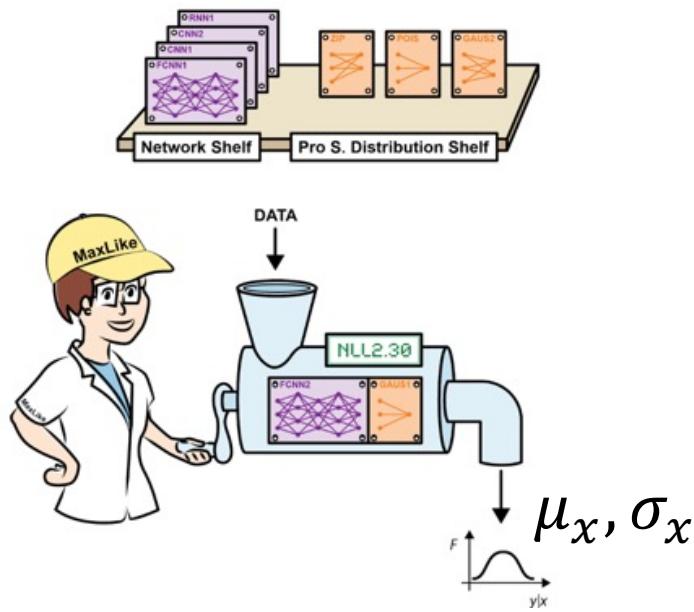
Some manipulation of $z(x)$ gives expectation of prob. distribution

Probabilistic Neural networks

- Replace linear predictor with (deep) NN
- Model not just expectation but all parameters of the CPD

CPD in ML / deep learning setting

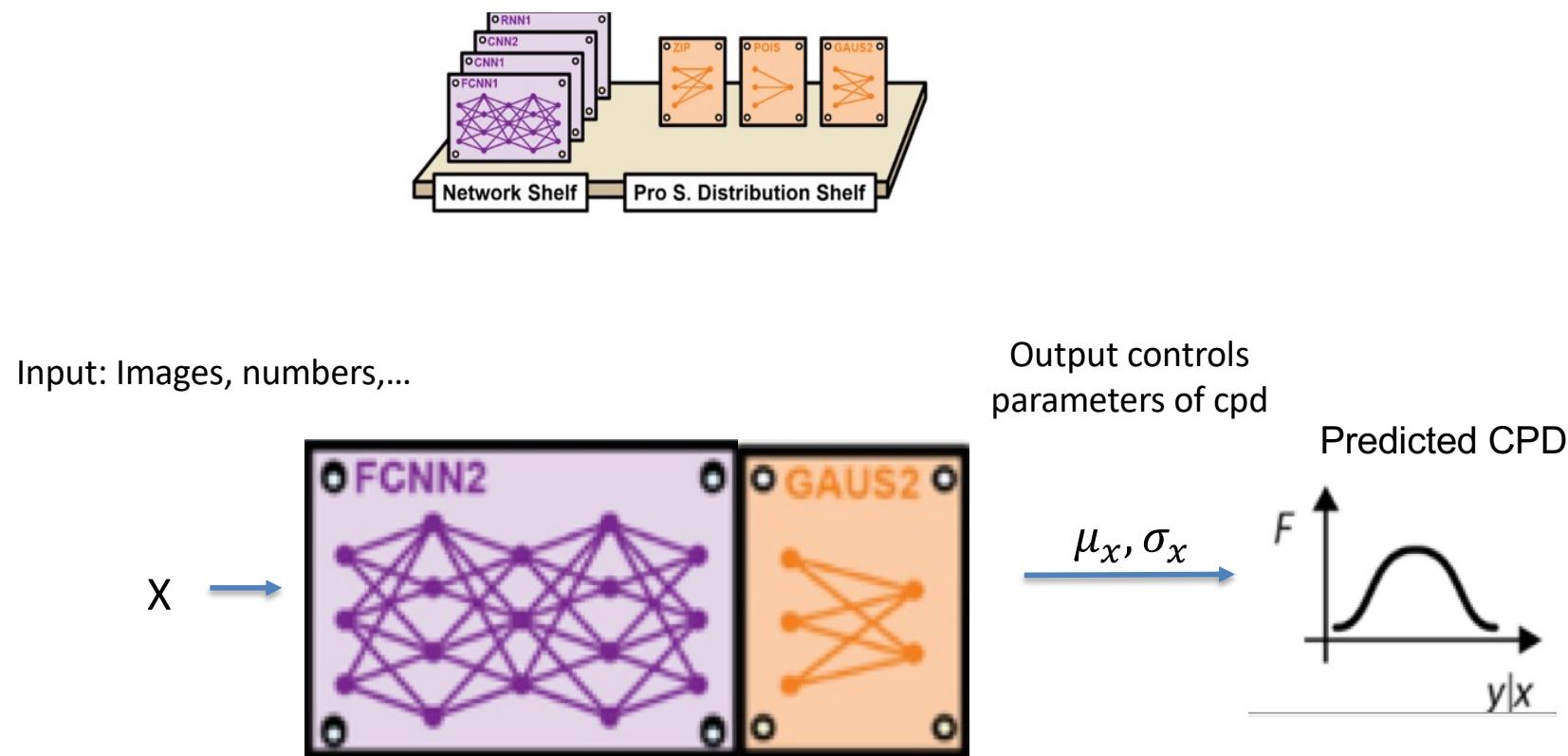
- In machine learning, we learn from N examples $D = \{x^{(i)}, y^{(i)}\}$ with $i = 1, \dots, N$.
- After training all information is in the optimally chosen weights \hat{w}
 $p(y|x, \hat{w})$
- Deep Learning recipe for probabilistic models



The network determines the parameters of a CPD given x .

Regression output as CPD

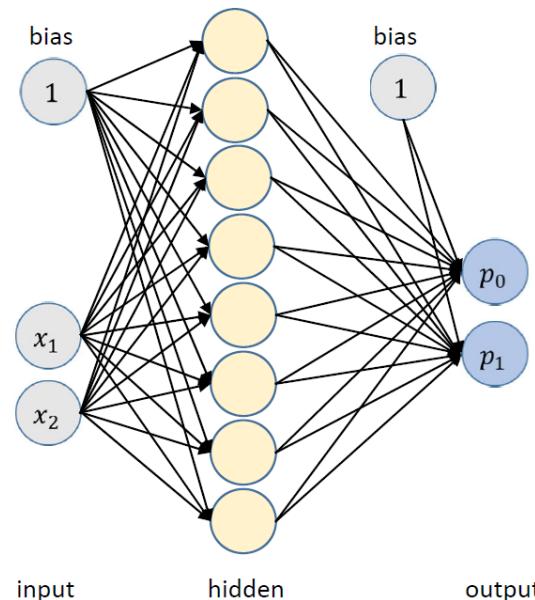
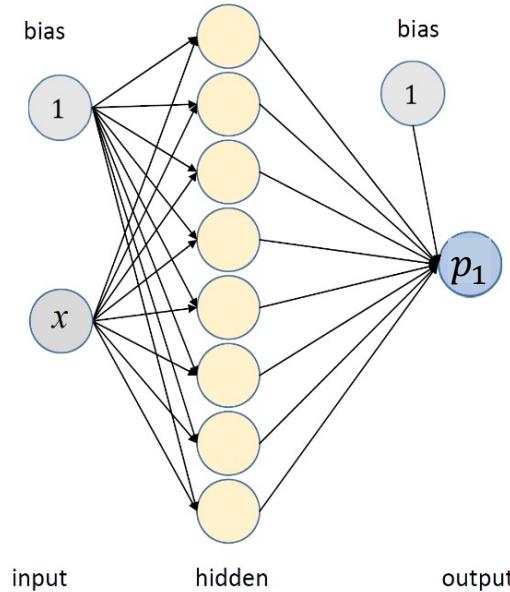
- A trained network outputs for each x the parameters of a distribution.
- For example, network controls μ_x, σ_x of $N(\mu_x, \sigma_x)$ for a given x .



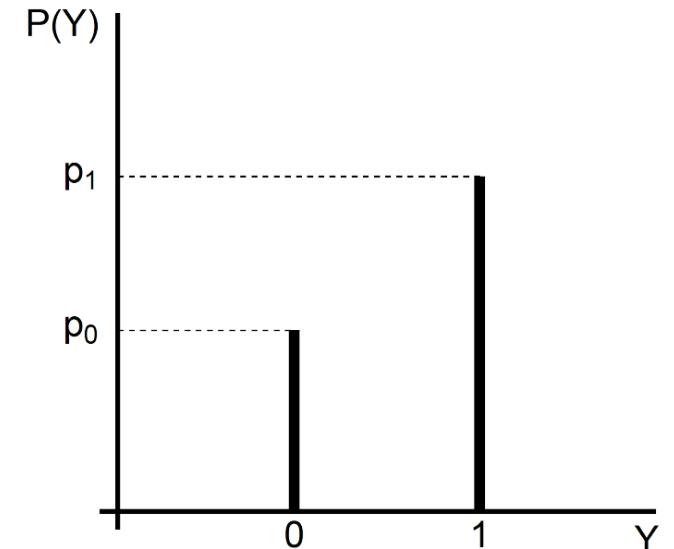
Binary Classification Output as CPD

- Binary Classification single output is sufficient (prob. p_1 for class 1)
- This can be seen as parameter p_1 for Bernoulli distribution
 - $y \sim \text{Bern}(p_1)$

Two possible NN architectures



Predicted CPD

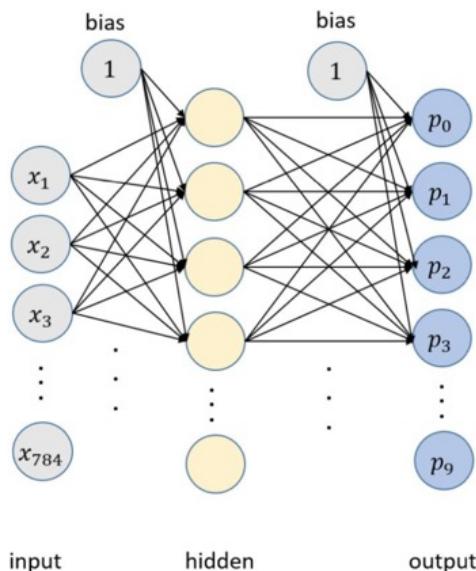


Which activation function is used in last layer?

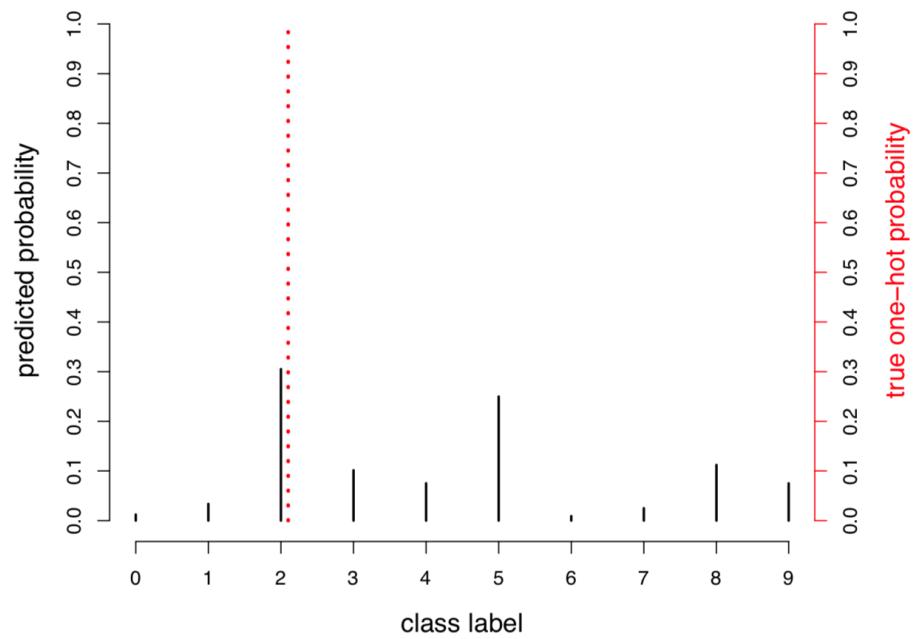
$$P(Y = k) = \begin{cases} p_0 & \text{for } k = 0 \\ p_1 & \text{for } k = 1 \end{cases} \quad \text{with } \sum p_i = 1$$

Multi-class Classification Output as CPD

- In multi-class classification settings the output at node k can be interpreted as probability for that class **or alternatively as parameters p_1, p_2, \dots, p_k for a categorical CDF with k categories.**



Predicted CPD



$$p(Y = k | x, w) = \begin{cases} p_0(x, w) & \text{for } k=0 \\ p_1(x, w) & \text{for } k=1 \\ \vdots & \vdots \\ p_9(x, w) & \text{for } k=9 \end{cases} \quad \text{with } \sum_{i=0}^9 p_i(x, w) = 1$$

Conditional **P**robability **D**istribution:
Unifying framework for
classification and regression.

Network controls the parameters of
a CPD for given value x .

How to fit a probabilistic model?

Recap: Maximum Likelihood (one of the most beautiful ideas in statistics)



Likelihood / “probability”
(often known)
 $M(\theta) \longrightarrow \text{Data}$

Ronald Fisher in 1913
Also used before by
Gauss, Laplace

Tune the parameter(s) θ of the model M
so that (observed) data is most likely

Maximum Likelihood rules Machine Learning



Tune the parameters weights of the network, so that observed data (training data) is most likely. We assume iid training data:

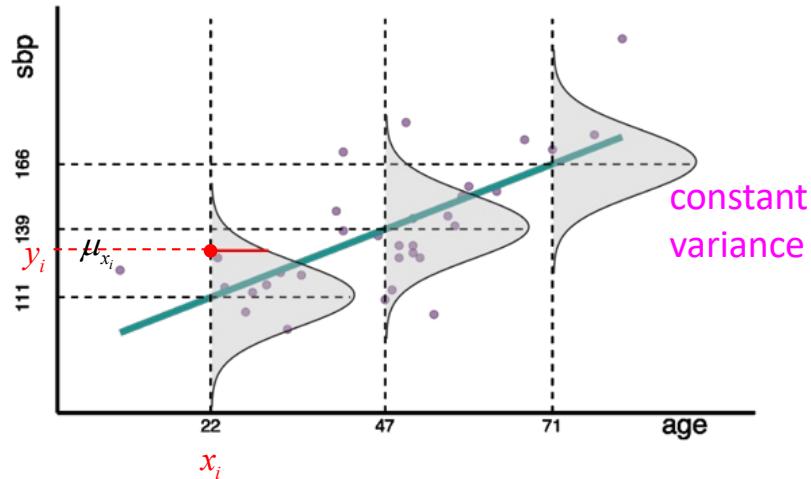
$$\hat{w} = \operatorname{argmax} \prod_{i=1}^N p(y_i | x_i, w)$$

Practically: Minimize Negative Log-Likelihood NLL of the CPD (take log; minimize after multiplication with -1)

$$\hat{w} = \operatorname{argmin} \sum_{i=1}^N -\log(p(y_i | x_i, w))$$

Example linear regression

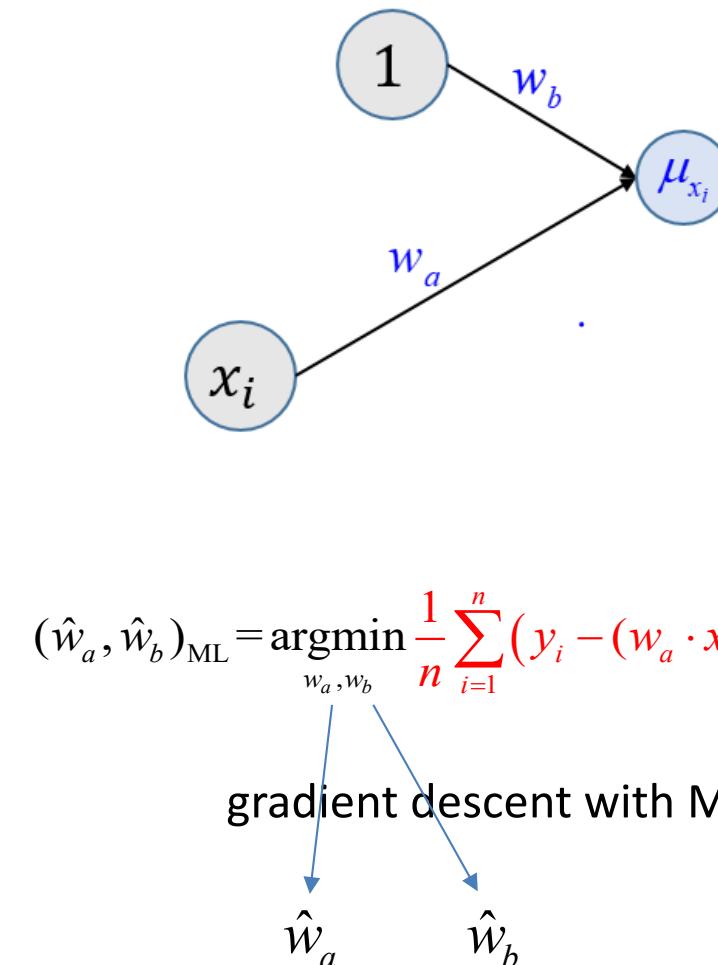
$$Y_{X_i} \sim N(\mu_{x_i}, \sigma^2)$$



Maximum likelihood:

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

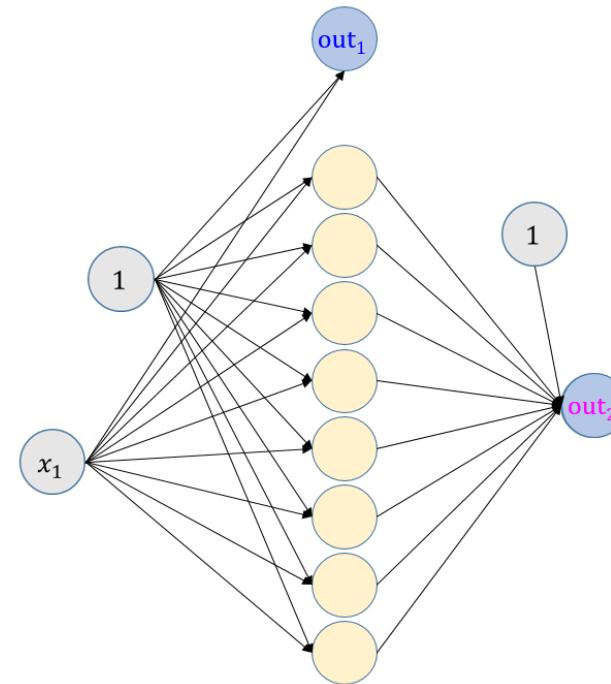
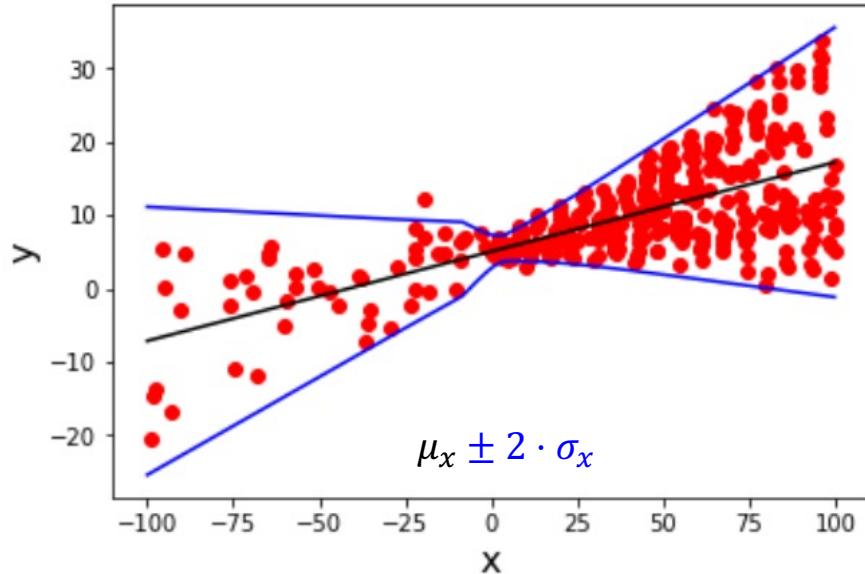
Negative Log-Likelihood (NLL)



Minimizing NLL loss $\stackrel{\sigma \text{ const}}{\longrightarrow}$ Minimizing MSE loss

Fit a probabilistic regression with flexible non-constant variance

$$Y_{X_i} = (\mathbf{Y} | \mathbf{X}_i) \sim N(\mu_{x_i}, \sigma_x^2)$$



$$\mu_{x_i} = out_{1_i}$$

$$\sigma_{x_i} = e^{out_{2_i}}$$

Minimize the mean negative log-likelihood (NLL) on train data:

$$NLL(w) = \frac{1}{n} \sum_{train-data} -\log(p(y_i | x_i, w))$$

gradient descent with NLL loss

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

Note: we do not need to know the “ground truth for σ ” – the likelihood does the job!

Modelling the standard deviation (positive values)

- The variance or standard deviation are both positive
- Neural networks output is not constrained.
- Two common approaches to fix this (exp or softplus)

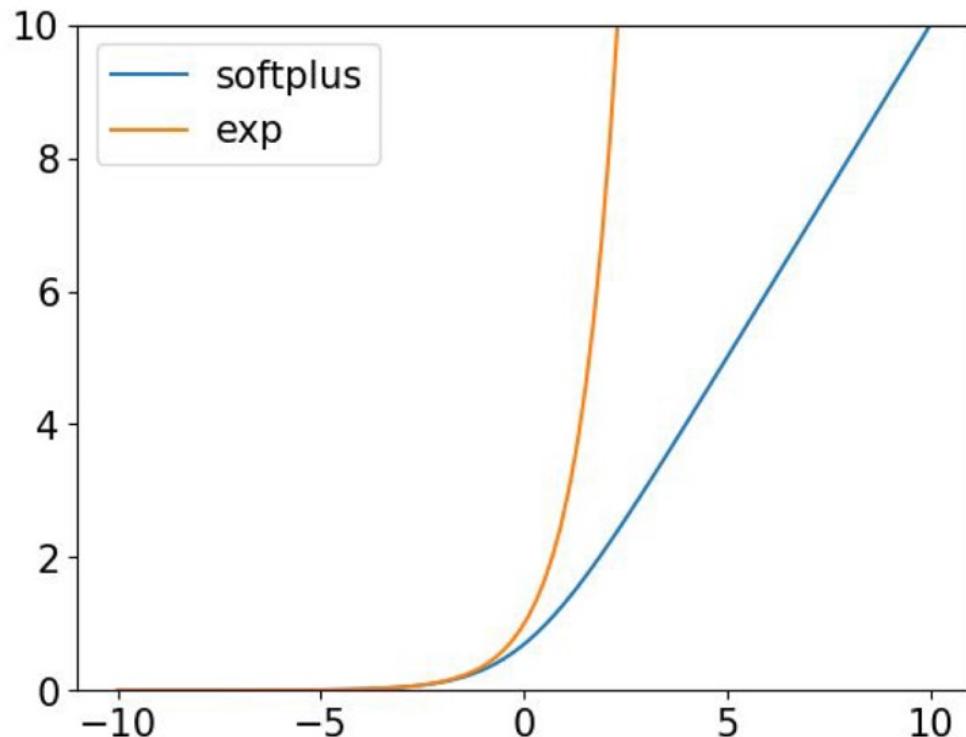
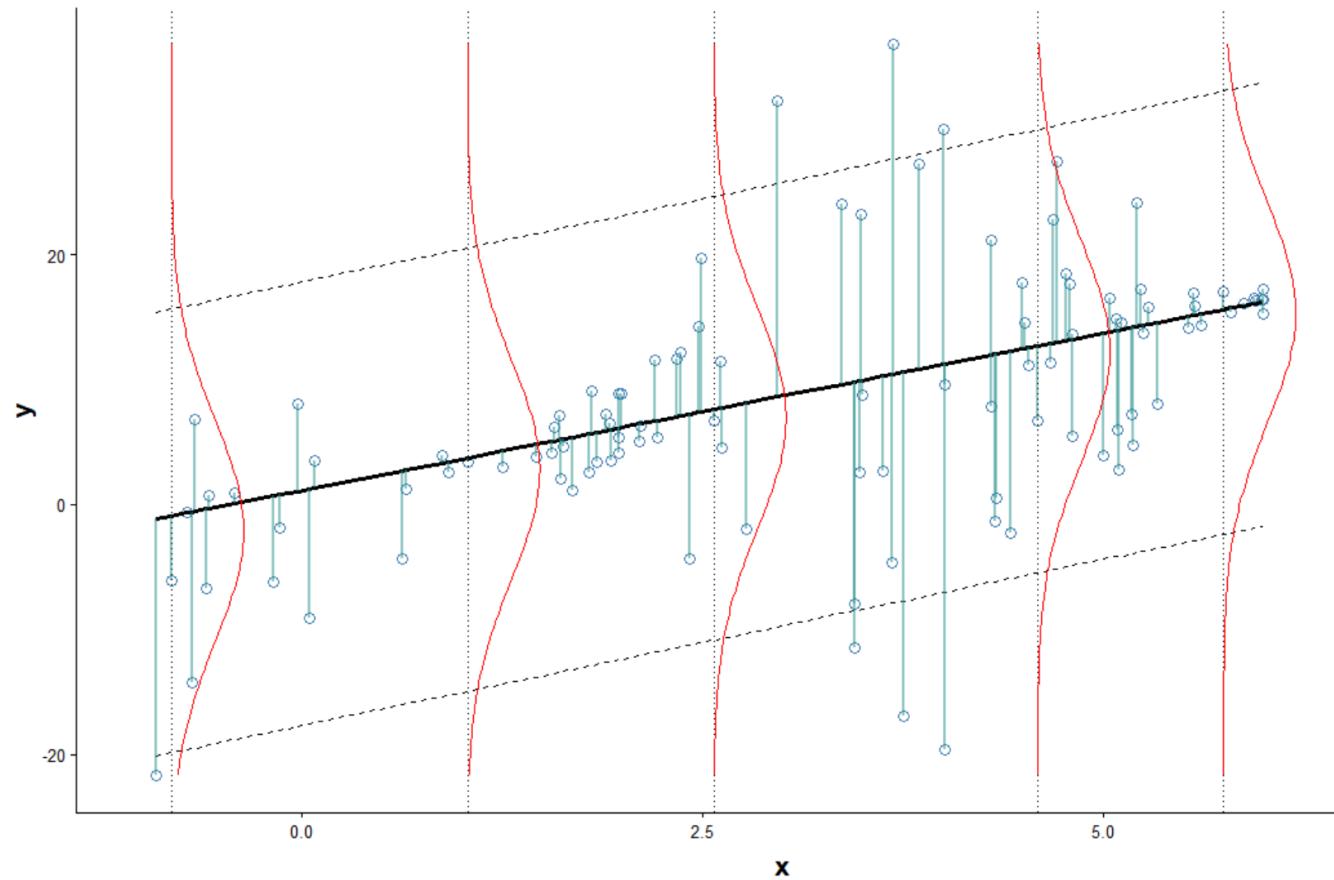


Figure 5.sp: The softplus function compared with the exponential function. Both functions map arbitrary values to positive values.

How to evaluate a probabilistic prediction model?

Root mean square error (RMSE) or mean absolute error (MAE)

Validation data



$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{\mu}_{x_i})^2}$$

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{\mu}_{x_i}|$$

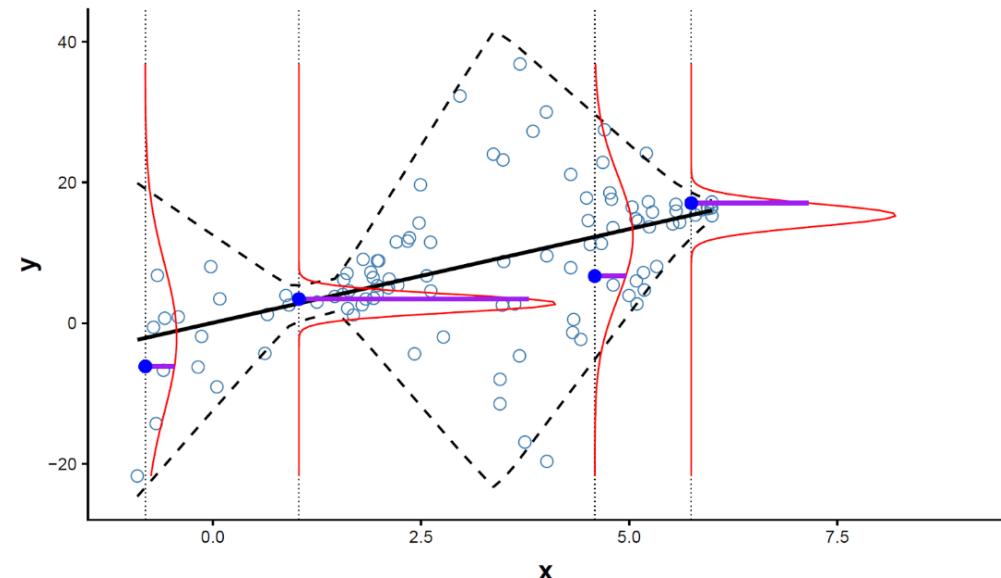
RMSE and MAE alone do not capture performance for probabilistic models!

Both only depend on the mean (μ) of the CPD, but not on its shape or spread (σ) and are not appropriate to evaluate the quality of the predicted distribution of a probabilistic model.

Use the NLL on test data to assess the prediction performance

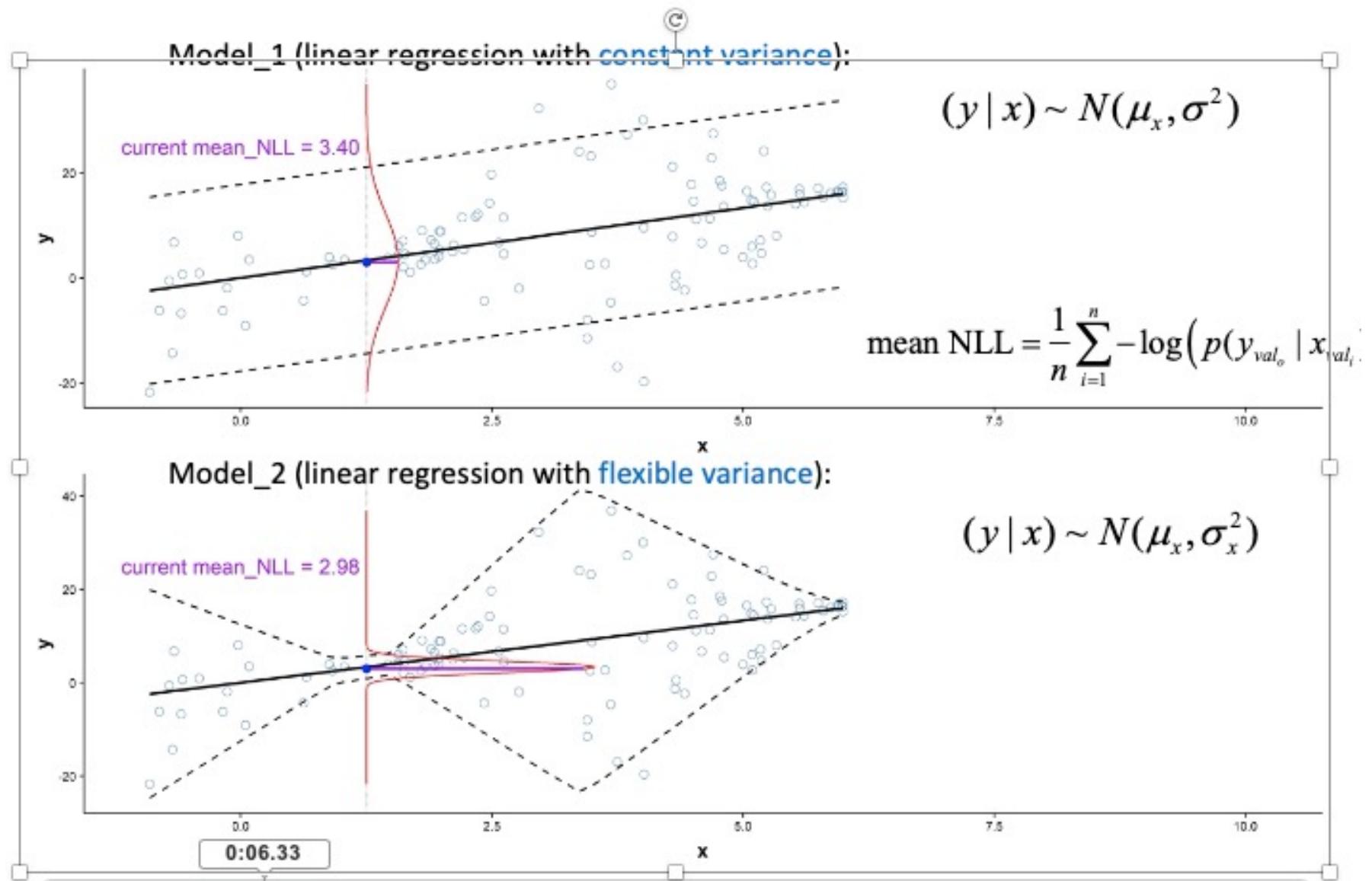
- Use the model on **test data** to **predict** for each input x_i a CPD p_{pred}
- Evaluate the predicted CPD at the position of the **observed outcome** y_i :

$$NLL = \frac{1}{n} \sum_{test\text{-data}} -\log(p(y_i|x_i))$$



[Tensorchief's youtu.be channel](#)

Live Demo



https://www.youtube.com/watch?v=QeWd0g_UzZY

Notes

- Log-Score a.k.a. NLL was first mentioned in connection to weather forecasting by Good (1952), who went so far as to suggest that the funding of the Met Office should depend on it.
- “If you see an analysis using something else, either it is a special case of the log scoring rule or it is possibly much worse”.
 - Richard McElreath in [Statistical Rethinking](#)
- Widely used in DL-community as NLL on validation set to compare different models (sometimes different names) cross-entropy, “bits per pixel”, perplexity $\exp(\text{NLL})$

Some hint to detailed treatment

- Statistical View: NLL is a so-called *proper scoring rule*
 - It's minimal (in expectation) if predicted distribution is data-generating distribution
 - Under certain assumptions (locality, smoothness, continuous outcome) NLL is the only one (strictly proper)^{*}
 - For NLL: expected score = $E_{(x',y') \sim p_{true}}[-\log(p_{pred}(y'|x'))]$
 - Gneiting, Raftery “Strictly Proper Scoring Rules, Prediction, and Estimation”
- Alternative View: Information Theory
 - The expectation of NLL is call *cross entropy*
 - For $p = p_{true} = p_{pred}$ it reaches its minimal value
 - $E_{(x',y') \sim p}[-\log(p(y'|x'))] = E_p[-\log(p)] =: H(p)$ The Entropy
 - Information theory
 - Entropy coding length under true distribution
 - Cross Entropy coding length under wrong distribution p_{pred}
 - Difference is the KL-Divergence
 - See e.g. McElreath “Statistical Rethinking”

^{*}(Bernardo, J. M., 1979: Expected information as expected utility. *Ann. Stat.*, 7, 686–690)

NLL as general cure-all in probabilistic modeling

Training of probabilistic models

- Maximize likelihood \leftrightarrow minimize negative log-likelihood (NLL)

Evaluation / model comparison:

- The log-score (NLL) is strictly proper score for regression.
 - The log-score (NLL) is also strictly proper for classification models.
- => To evaluate or compare probabilistic models: use the validation NLL!

The NLL can be used for scoring and training

Probabilistic NN (technical details)

TensorFlow Probability

TensorFlow Probability (TFP) is an add-on to TF/Keras

TFP can model Distributions with `tfd_...`

- Different distributions all have the same API.

```
d = tfd_normal(loc = 1.0, scale = 0.1) # Create a normal distribution with mu=1 and sigma=0.1
(s = tfd_sample(d,2)) # Draw two random points, note we have tensor
as.numeric(s) # Now, we have numeric values
s$numpy()      # Alternative use when as.numeric() makes troubles
tfd_prob(d, 3) #Compute density/mass.
tfd_cdf(d, 3) #Compute CDF
tfd_mean(d) #Compute mean (expectation)
tfd_stddev(d) #Compute std
````
```

```
tf.Tensor([1.0305088 0.98513895], shape=(2,), dtype=float32)
[1] 1.030509 0.985139
[1] 1.030509 0.985139
tf.Tensor(0.0, shape=(), dtype=float32)
tf.Tensor(1.0, shape=(), dtype=float32)
tf.Tensor(1.0, shape=(), dtype=float32)
tf.Tensor(0.1, shape=(), dtype=float32)
```

## Exercise 1



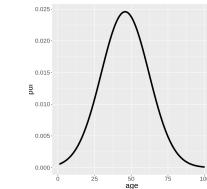
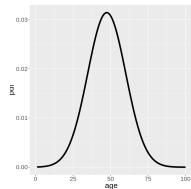
Notebook: [06-into-tfp.ipynb](#)

- Getting Started with TFP

## Excercise UK\_Face

- Build an CNN with outputs a Gaussian describing the age

$x$



- The CNN controls  $\mu$  and  $\sigma$

## Code

```
model <- keras_model_sequential() %>%
 ## First Block
 layer_conv_2d(filters = 16, kernel_size = c(3, 3)...) %>%
 layer_conv_2d(filters = 16, kernel_size = c(3, 3)... %>%
 layer_max_pooling_2d(pool_size = c(2, 2)) %>%
...
 layer_flatten() %>%
 layer_dense(units = 256, activation = "relu") %>%
 layer_dropout(dropout_rate) %>%
 layer_dense(units = 64, activation = "relu") %>%
 layer_dropout(dropout_rate) %>%
 layer_dense(units = 2) %>%
 layer_distribution_lambda(make_distribution_fn = CPD)
```

The CPD (conditional probability distribution) is a to be defined functions, which takes the last layer of the network (here two nodes) and returns a TFP function.

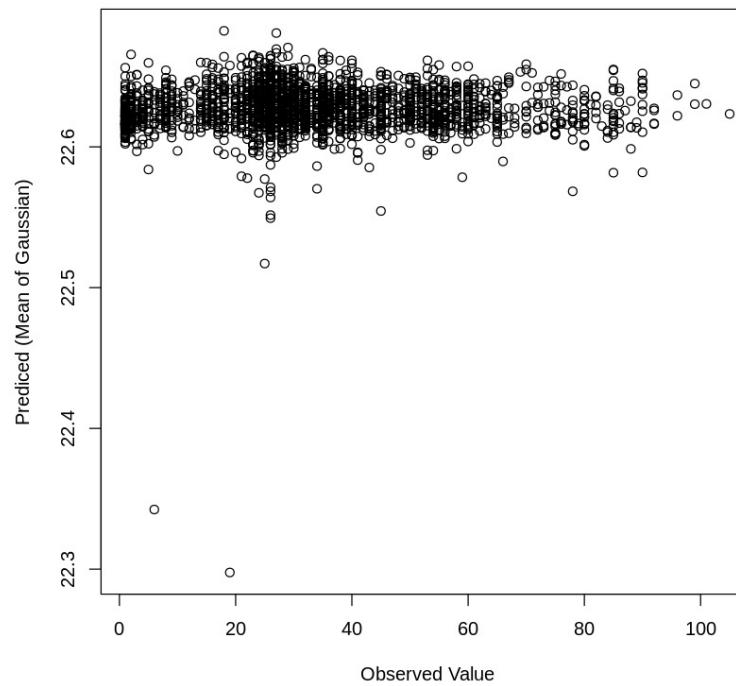
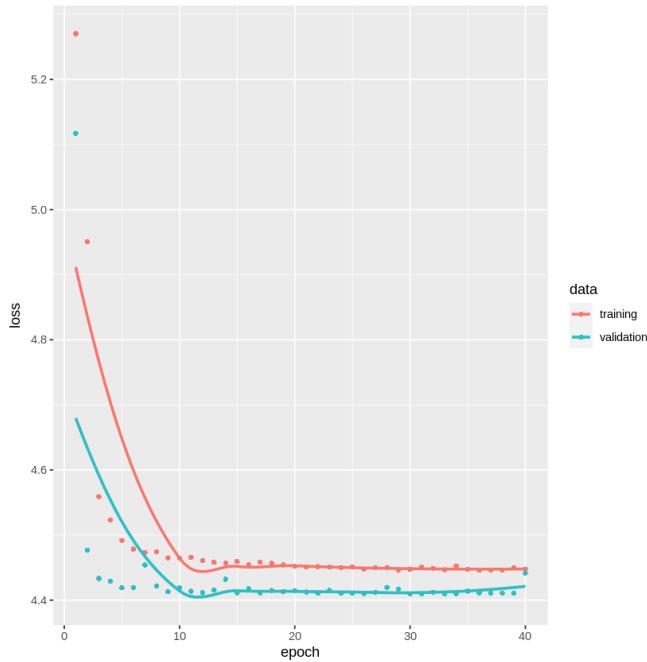
# Attention (How I spend my last Friday)

```
layer_dense(units = 2) %>%
 layer_distribution_lambda(make_distribution_fn = CPD)
```

CPD = `function(out){`

**What is wrong?**

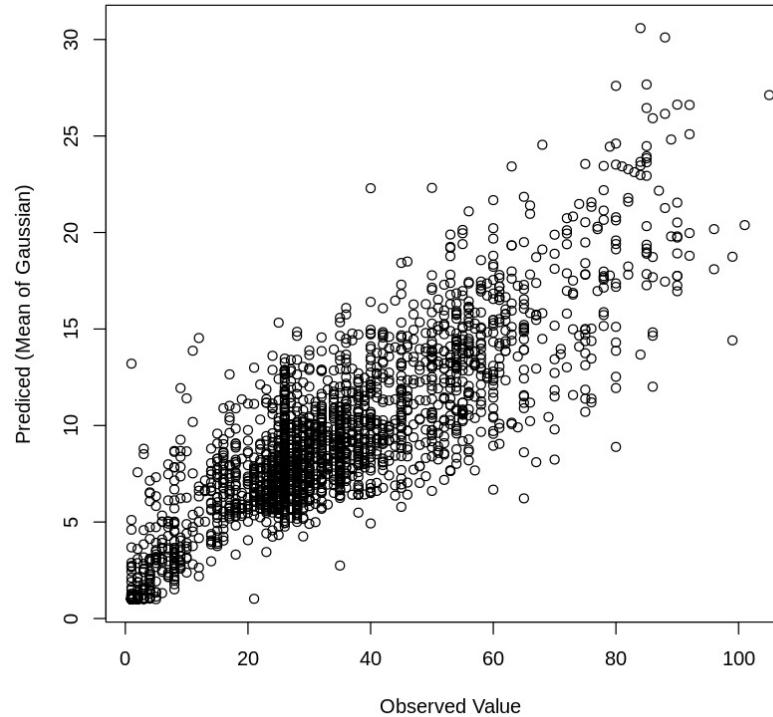
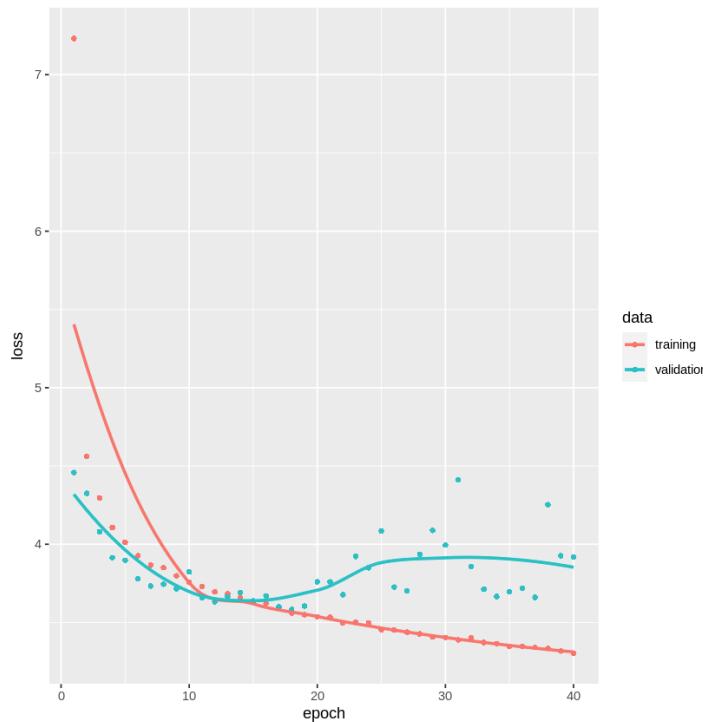
```
 mu = k_softplus(out[, 1]) #Need to be positive as well
 s = 5 + 0.1 * k_softplus(out[, 2]) # 0.1 for numerical stability
 return (tfd_normal(loc = mu, scale = s))
}
```



# Attention (How I spend my last Friday)

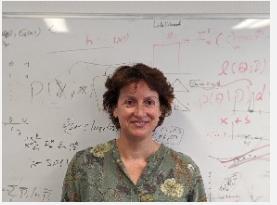
```
layer_dense(units = 2) %>%
 layer_distribution_lambda(make_distribution_fn = CPD)
```

```
CPD = function(out){
 mu = k_softplus(out[, 1, drop=FALSE]) #Need to be positive as well
 s = 5 + 0.1 * k_softplus(out[, 2, drop=FALSE]) # 0.1 for numerical stability
 return (tfd_normal(loc = mu, scale = s))
}
```

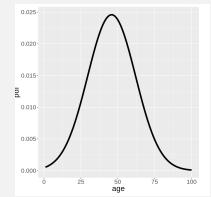
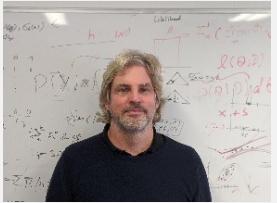
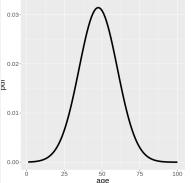


## Exercise 2 (07\_uk\_faces)

$x$



$$p(y|x)$$



Make sure to run on GPU!  
Runtime -> Change runtime type



If you want your own images, they need to be aligned

# Count Data

# Modeling the CPD for Count Data

Principle idea: Let network model the parameters of a distribution (same as above)

- Counts are non-negative integer numbers from 0 to  $\infty$
- Examples
  - Number of products bought
  - Number of customers coming
  - Number of fish caught
- Similar data (count data with upper bound)
  - rgb values of pixels bounded count data [0,255] e.g. PixelCNN
  - Amplitude of digitized sound e.g. [0,65535] e.g. Wavenet
  - Number of distinct types of crimes an individual has committed

# 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 of life bait was used

X4=camper: indicates if camper is brought



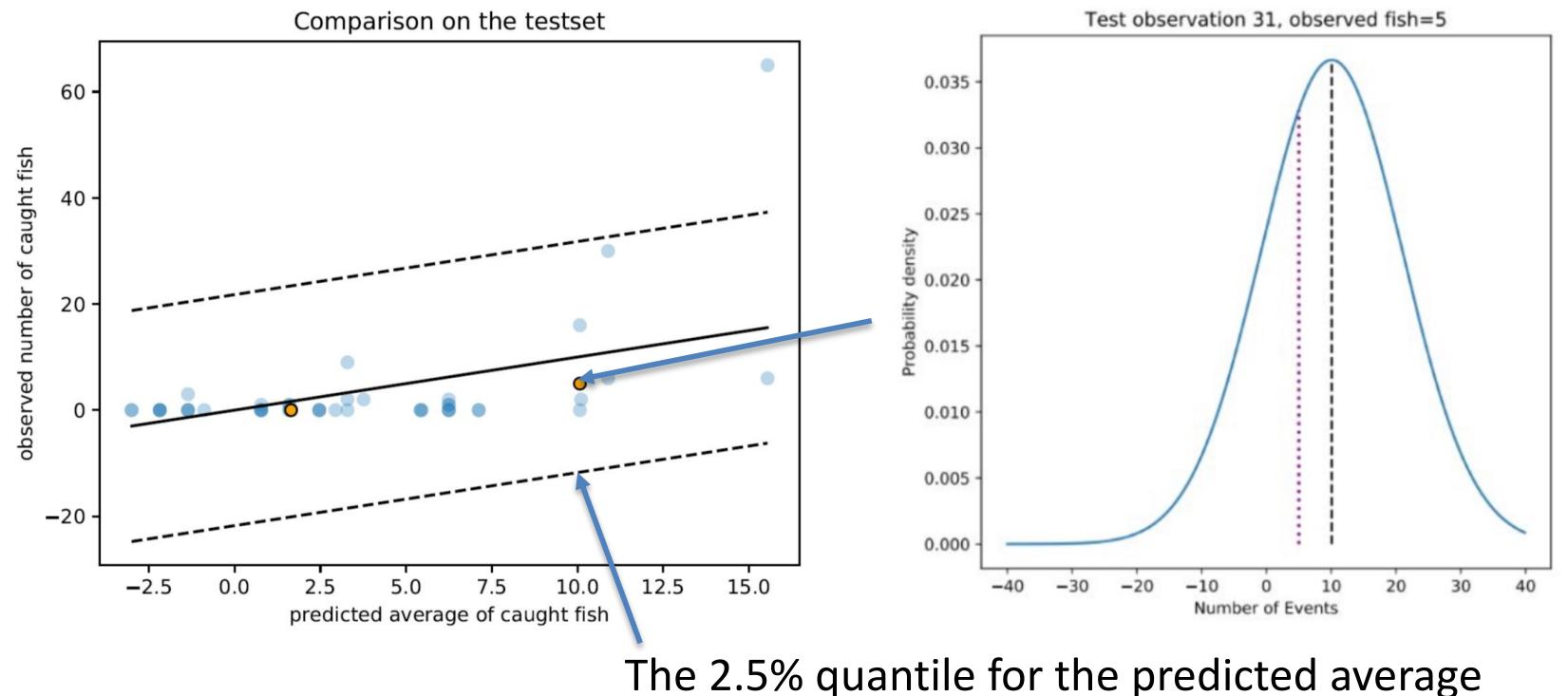
Data: <https://stats.idre.ucla.edu/r/dae/zip>

# Modeling count data:

## M1: Linear regression

# Linear Regression

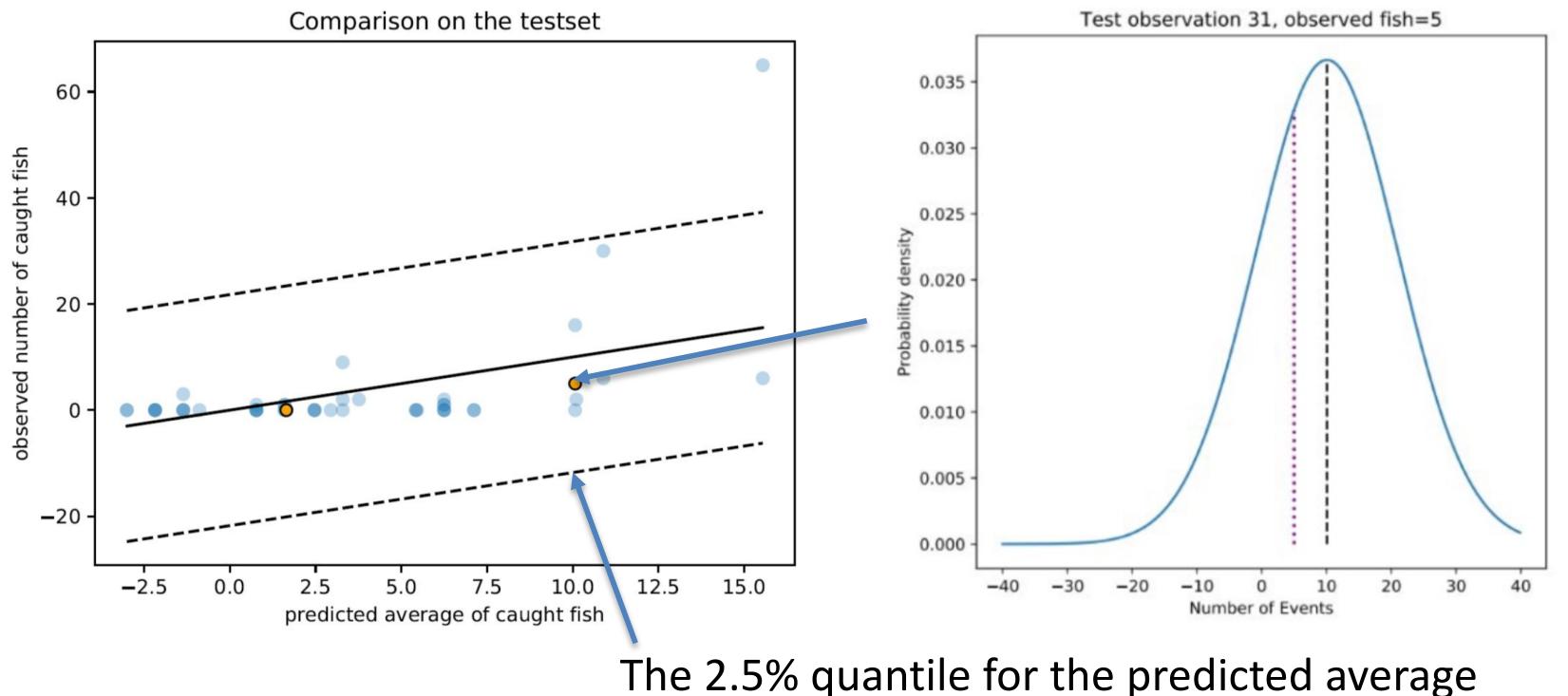
- First idea: use linear regression (constant variance)
- Result for two observations



Why is this not a good model for count data?

# Linear Regression

- First idea: use linear regression (constant variance)
- Result for two observations



Why is this not a good model for count data?

- Predicts negative number of fish
- Predicts non-integer number of fish

# Is a Gaussian CPD appropriate to model count data?

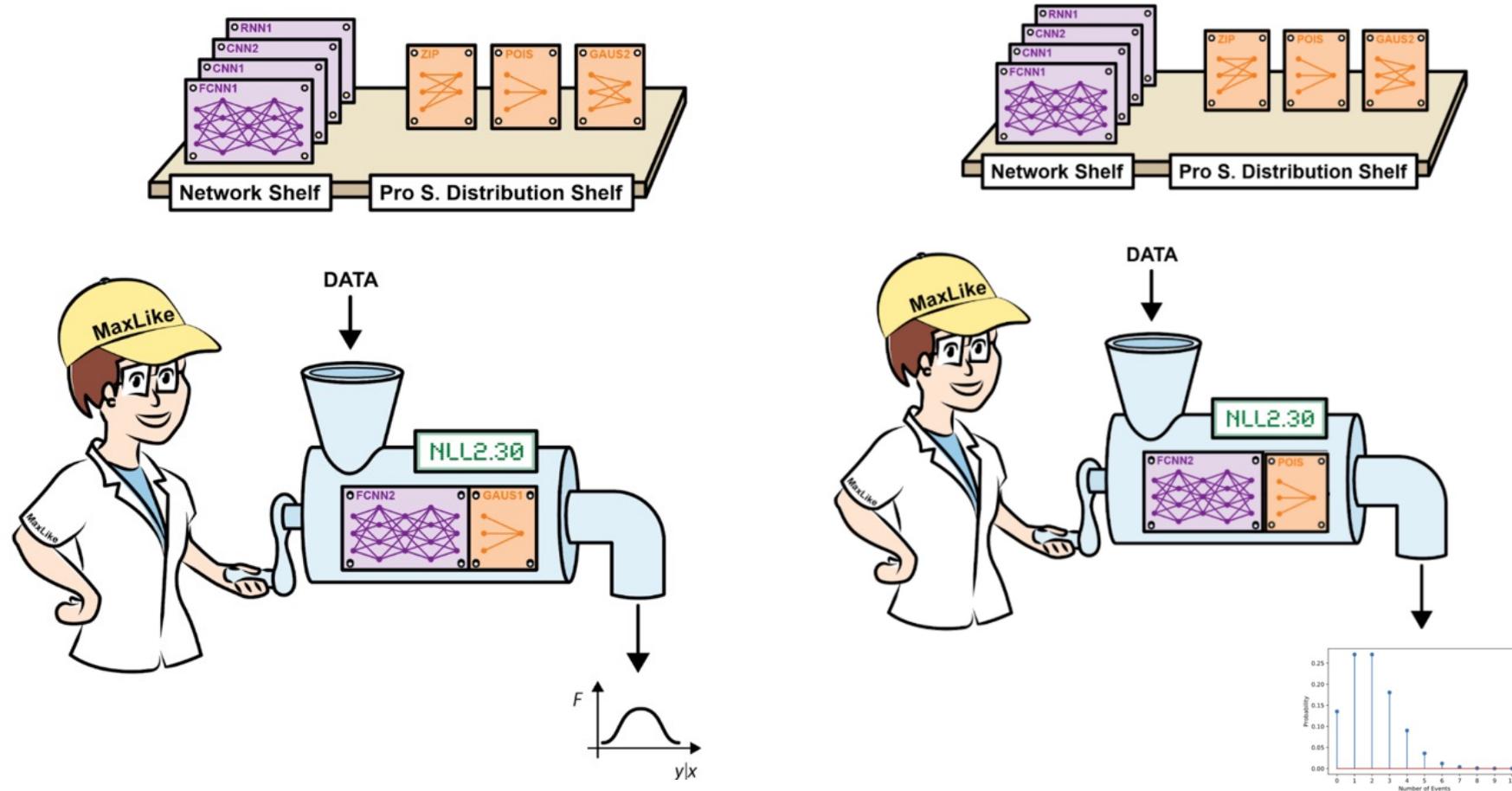
**A Gaussian CPD is not ideal, because:**

- Counts can only take integer values, but a Gaussian assigns likelihoods also to values between two integers
- If counts result from a Poisson process then the mean of the counts is equal to the variance, meaning the variance is not constant, which is not taken into account by linear regression with a Gaussian CPD  $N(\mu_x, \sigma)$

**A Gaussian CPD works often quite ok, because:**

- For large counts the discrete nature of the counts do not matter so much (looks like normal rounding of continuous values)
- When log-transforming count outcome  $Y$ , then the variance of  $\log(Y)$  is approximately constant and can therefore be modeled by a Gaussian CPD  $N(\mu_x, \sigma)$  (Tukey's first Aid Transformations)

# Probabilistic regression: from continuous outcome to counts → just change the tfp head of your DL model



Just change the distribution head (orange), pink architectural part (CNN,...) can stay the same!

Modeling count data:

M2: Poisson regression

(GLM)

# Poisson regression for count data

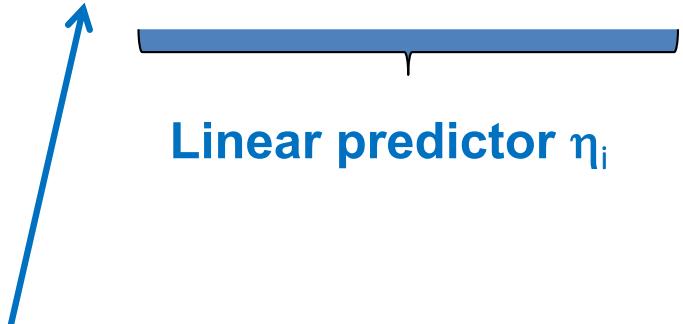
Goal: Predict a Poissonian CPD for  $(Y|X=x)$  which depends on predictor values

CPD:  $Y_{X_i} = (Y|X_i) \sim \text{Pois}(\lambda_{x_i})$

We only need to model  $\lambda_x$  to fix the Poissonian CPD!

Model:

$$\lambda_i = \exp(\beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip})$$



Inverse link-function: ensures positive counts

## Model 2: Poisson regression via NNs in keras

We use a NN without hidden layer to control the rate  $\lambda$ .

```
model_p = keras_model_sequential() %>%
 layer_dense(units=1L, activation=tfb_exp()) %>%
 layer_distribution_lambda(tfd_poisson)

model_p(as.matrix(X_train[1:3,])) #Returns distribution
model_p(as.matrix(X_train[1:3,]))$mean() #Returns mean value of distribution
```

Using the exp-activation ensures that  
the rate  $\lambda$  is a positive number

Glueing the NN and the output layer together.  
tfd\_poisson is a distribution

```
NLL = function(y_true, dist) -tfd_log_prob(dist, y_true)
model_p = compile(model_p, loss=NLL, optimizer=
 optimizer_adam(lr = 0.01))
```

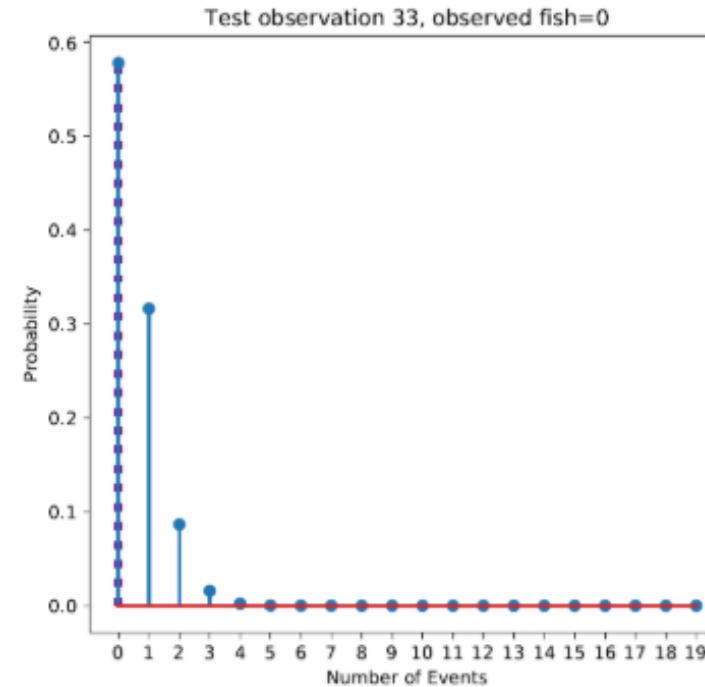
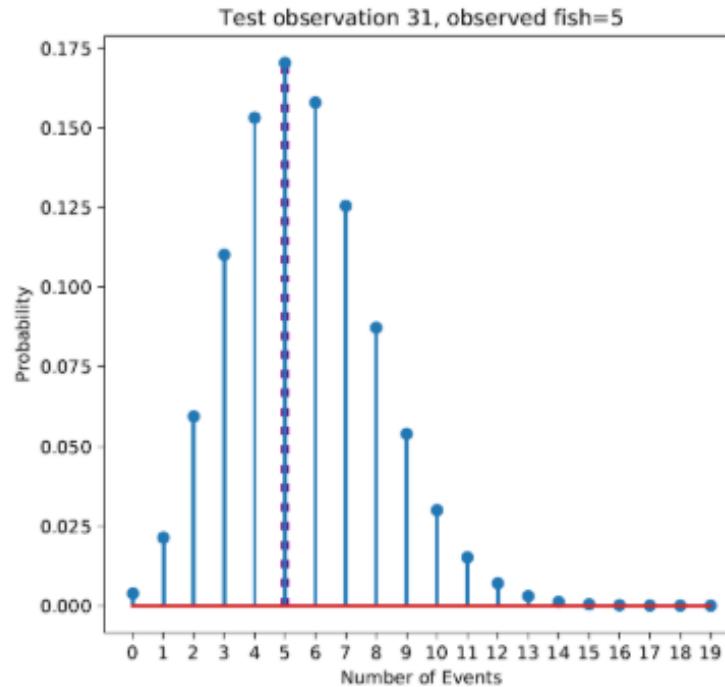
The second argument is the output of the model and thus a TFP distribution. It's as simple as calling log\_prob to calculate the log probability of the observation that's needed to calculate the NLL

## Model 2: Poisson regression, get test NLL from Gaussian CPD

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.

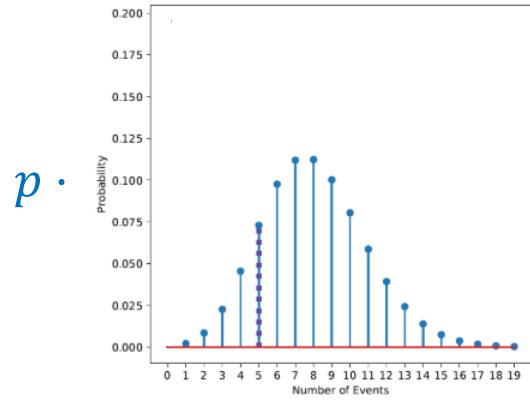


What is the likelihood of the observed outcome in test obs 31 and 33?

# Modeling count data: M3: ZIP regression

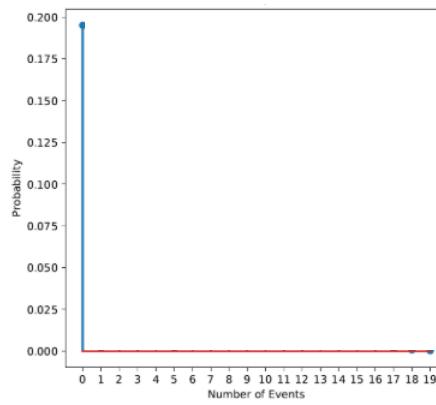
# Zero-Inflated Poisson (ZIP) can be seen as Mixture Distribution

Poisson-Process



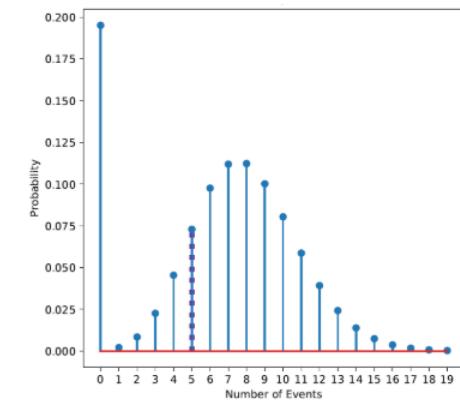
$p \cdot$

Zero-Process



=

Zero-inflated Poisson



# Custom distribution for a ZIP distribution Just for Reference

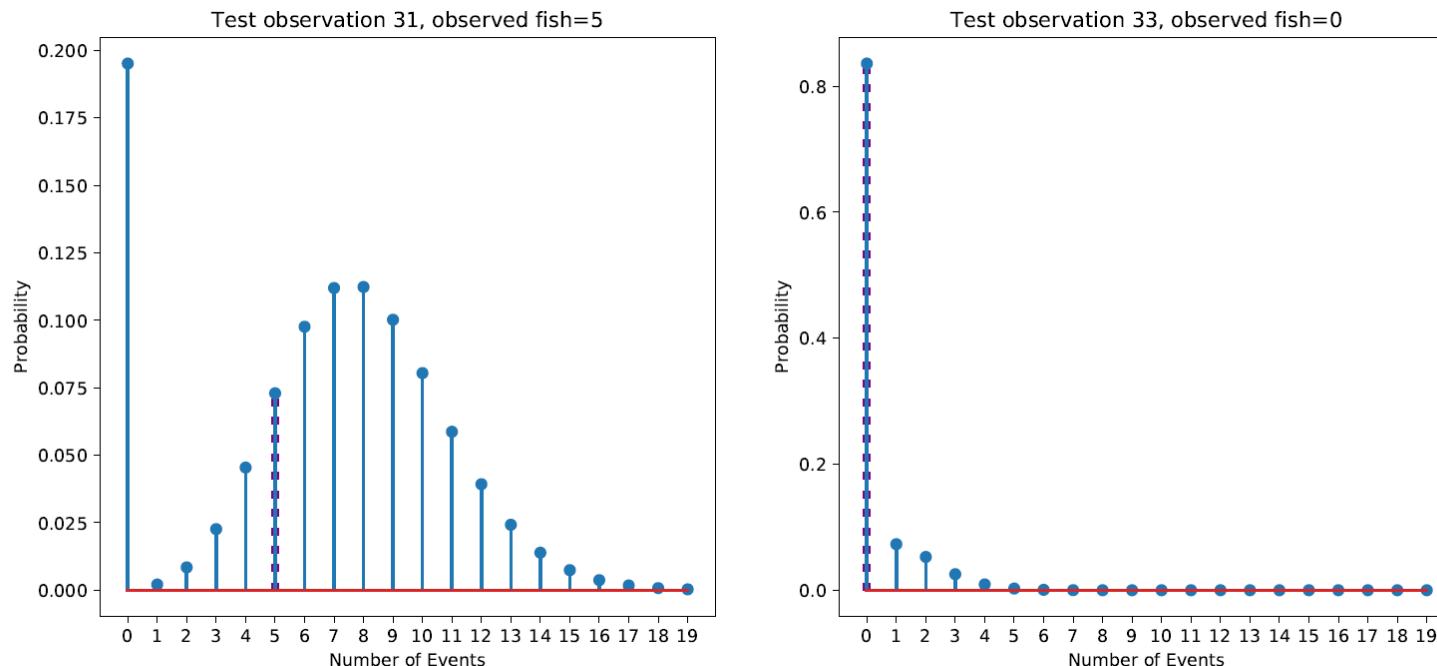
```
1 zero_inf = function(out){
2 rate = k_exp(out[, 1, drop=FALSE]) First NN output controls lambda. Exp guarantee value >0
3 s = k_sigmoid(out[, 2, drop=FALSE]) Second NN output controls p; sigmoid guarantees value in [0,1]
4 probs = k_concatenate(list(k_expand_dims(1-s),
5 k_expand_dims(s)), axis=-1) #C
6 The two probabilities for 0's or Poissonian distribution
7 return(
8 tfd_mixture(
9 #creating a mixture of two components:
10 cat = tfd_categorical(probs=probs) tfd.Categorical allows creating a mixture of two components
11 components = list(
12 #Zero as a deterministic va Zero as a deterministic value
13 tfd_deterministic(loc=k_zeros_like(rate)),
14 #Value drawn from a Poissonian distribution
15 tfd_poisson(rate=rate) Value drawn from a Poissonian distribution
16)
17)
18)
19}
20
```

## Model 3: ZIP regression two examples

Predict CPD for outcome in test data:

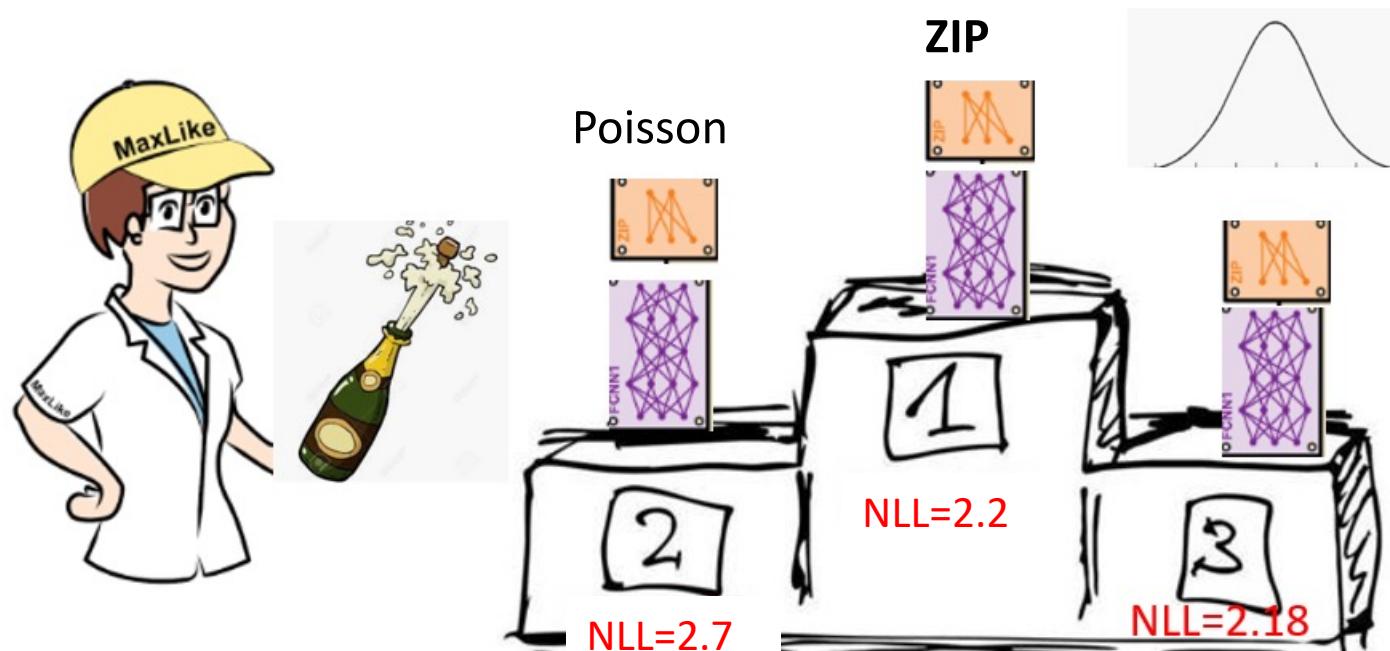
Group 31 used livebait, had a camper and were 4 persons with one child.  $Y=5$  fish.

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?

# Validation NLL allows to rank different probabilistic models



Note that NLL only can be directly compared among models which are both discrete or continuous (we could code continuous outcomes as interval-scaled values to get on same scale).

## Exercise 2

- Model Count Data 08



## Take home messages

- A probabilistic model predicts for each input a whole outcome CPD
- Probabilistic models are the models of choice when dealing with uncertainty
- Probabilistic models fit naturally in DL framework
- Good Probabilistic modeling improves SOTA in deep learning models
- Use the NLL for **training, evaluating and comparing** probabilistic models

# Technical Notes / Exercise

# Data Sets for Projekt

- Kaggle
  - Suited for DL (Computer Vision + Classification)
    - <https://www.kaggle.com/datasets?tags=13207-Computer+Vision%2C13302-Classification>
    - Good with code examples
  - Suited for DL (Regression let us know!)
  - Competitions
    - <https://www.kaggle.com/competitions>
- Own Data
  - See notebook 7\_uk\_face at the end of this presentation

# Upload data to kaggle

In Kaggle

08\_UTKFace Draft saved  
File Edit View Run Add-ons Help



Markdown ▾

Draft Session (11m)



Share

Save Version 0



## Probabilistic Age Estimation

In this NB we develop a CNN which can take images of faces and outputs a CPD, which is here a conditional Gaussian distribution, that assigns to each possible age a likelihood.

## Installing TensorFlow and Keras

The deep learning packages, TensorFlow and Keras are not installed by default, but can be installed as follows:

+ Code + Markdown

[ 6 ]:

```
COLAB = FALSE
KAGGLE = TRUE
```

Data

+ Add Data

Input

Output (280.1MB / 19.5GB)

/kaggle/working

Settings

ACCELERATOR

None

LANGUAGE

R

## Press Add Data

```
system("tar xvzf ../input/ukface-split/ukface_split.tgz") #Unpacking the images
```

Create a New Dataset

ukface\_split

www.kaggle.com/oduerr/ukface-split

Uploaded files

ukface\_split.tgz (106.52 MB)

+ Upload Files

✓

Uploaded 1 of 1 file.

Private Create

## Output (178.6MB / 19.5GB)

- ▼  /kaggle/working
  - ▶  valid
  - ▶  train
  - ▶  test



## Excercise



- 06 Getting started with TFP
- 07 Continuous Data (UK\_Face)
- 08 Count Data