

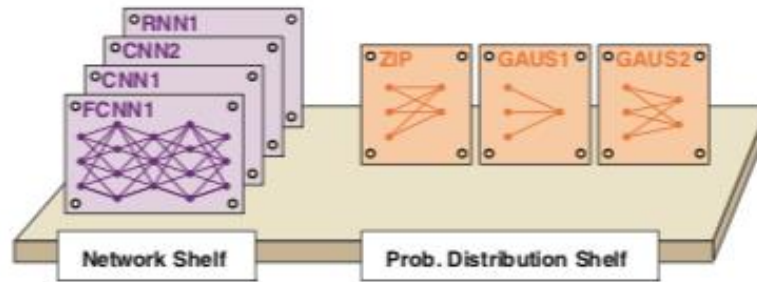
WBL Deep Learning:: Lecture 5

Beate Sick, Oliver Dürr

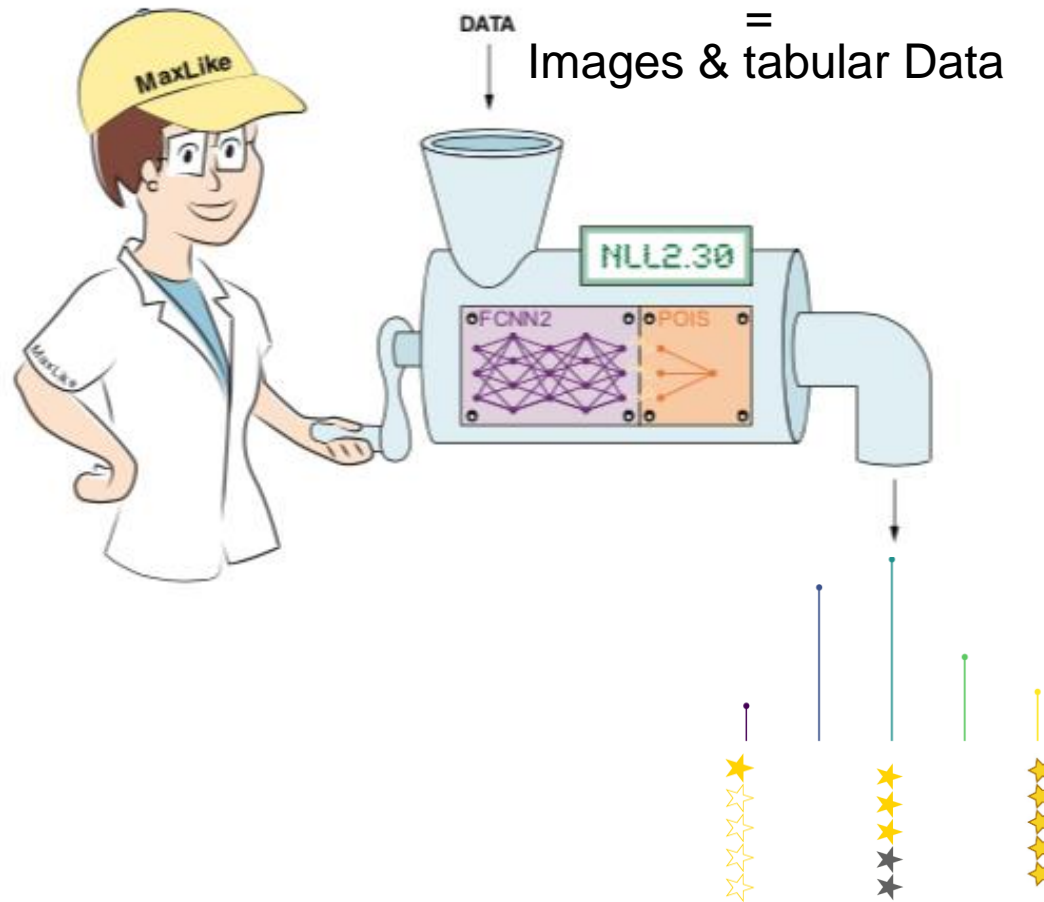
Deep Learning for interpretable semi-structured models

Zürich, 10/10/2022

Probabilistic DL with semi-structured data



semi-structured Data
=
Images & tabular Data



Aim: Interpretable models for semi-structured data

Use NN for complex data (e.g. images) combined with a statistical model which allows for interpretable coefficients

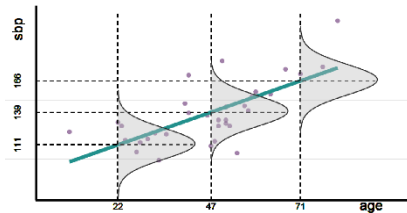
Today

- Look again at logistic regression as latent variable model
- Logistic regression with semi-structured data
- Interpretable semi-structured Ordinal Regression models (ONTRAMs)
- Interpretable Ordinal Regression Models with Neural Networks
- Formulation as transformation model
- Ensembling for improving prediction performance and quantifying parameter uncertainty

Taking the best of both worlds

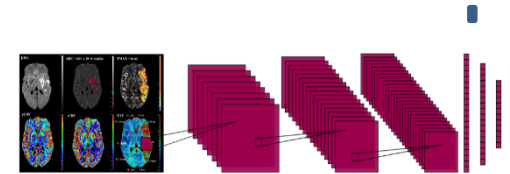
Statistics

- + transparent & interpretable
- + valid uncertainty measures
- needs tabular data



Deep Learning

- + can handle tabular (structured) & unstructured (e.g. image) data
- + high prediction performance
- black box, not interpretable



Interpretable & probabilistic DL

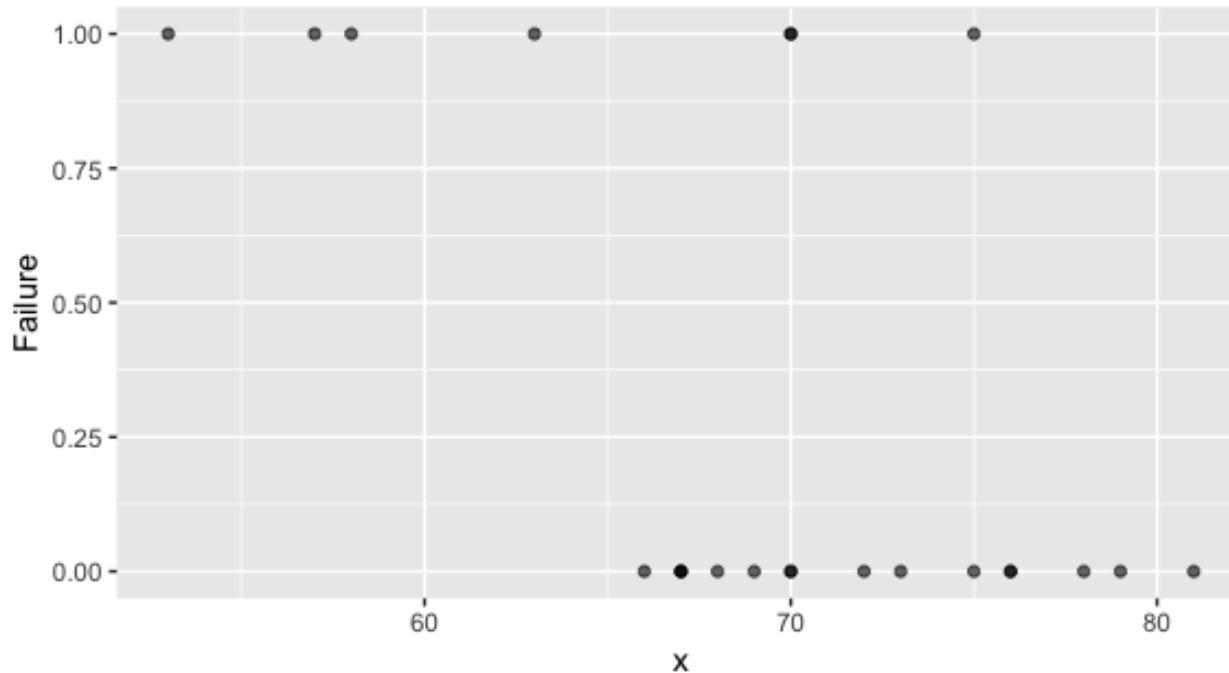
Logistic regression as latent variable model

Recap: Modelling with logistic regression

Zero / one classification

Want: $p(x) = \Pr(Y = 1|x)$

Prob. for an O-ring to be defect $Y=1$ at a given temperature X

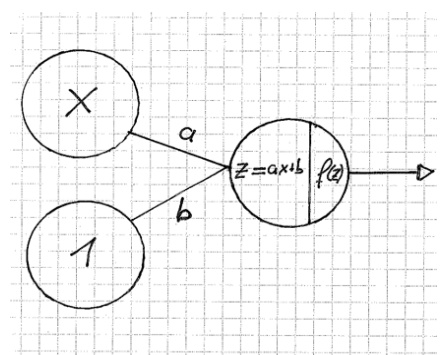


- Guess curve?

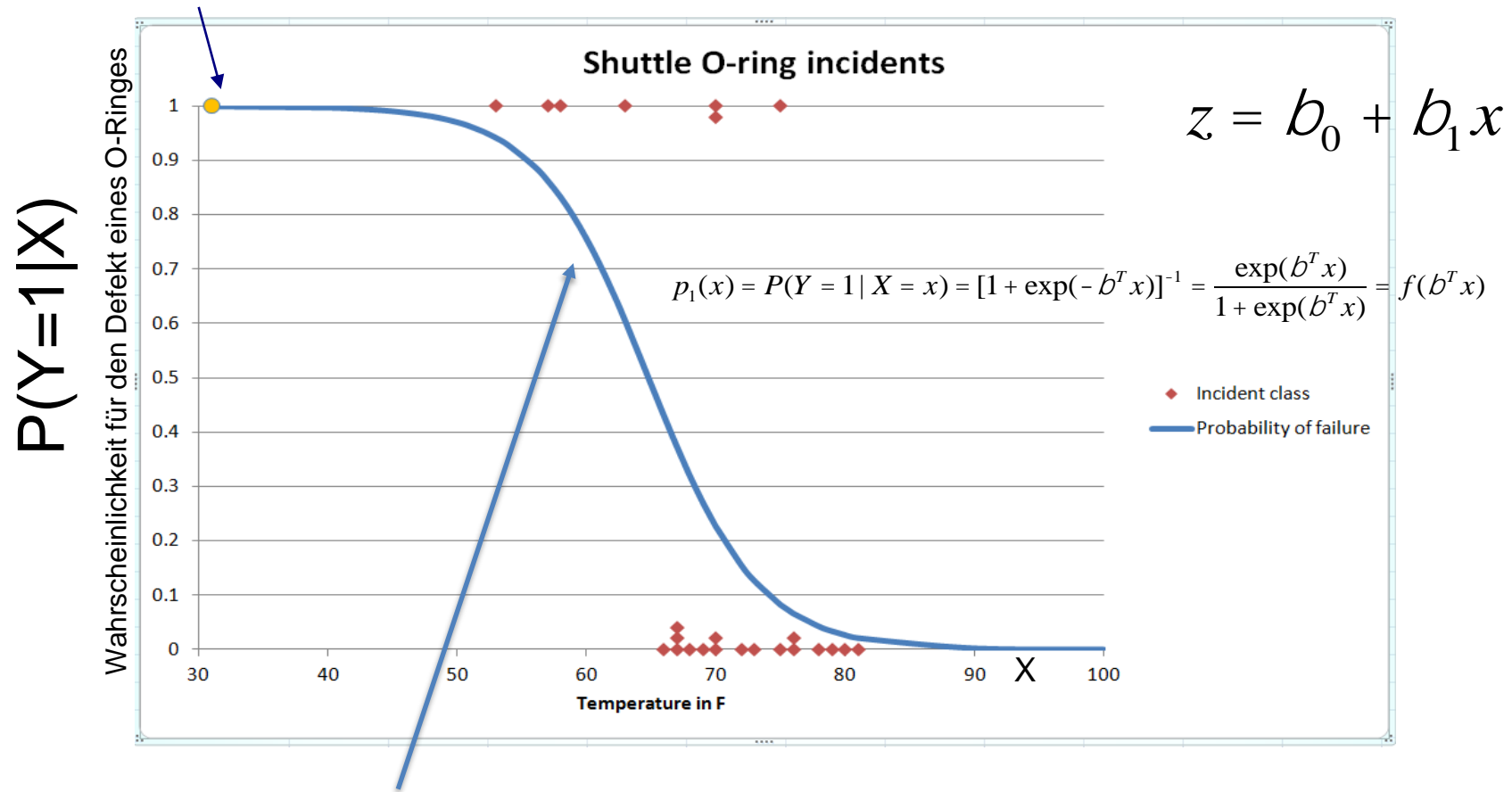


Recap: Logistic Regression

Predict if O-Ring is broken, depending on temperature



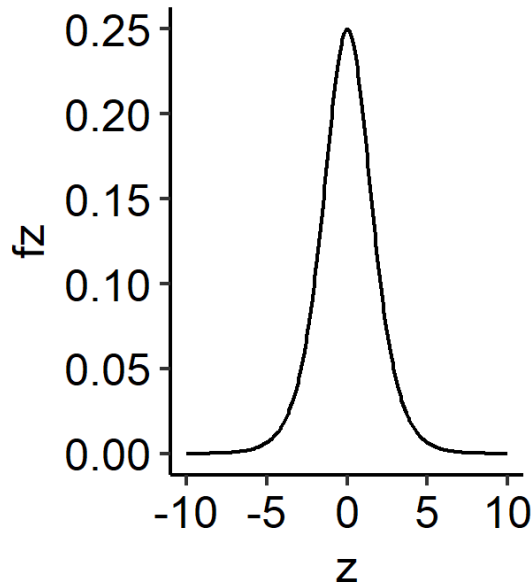
Challenger launch @31 F
Prob. of a failure=0.9997



- What is that curve?

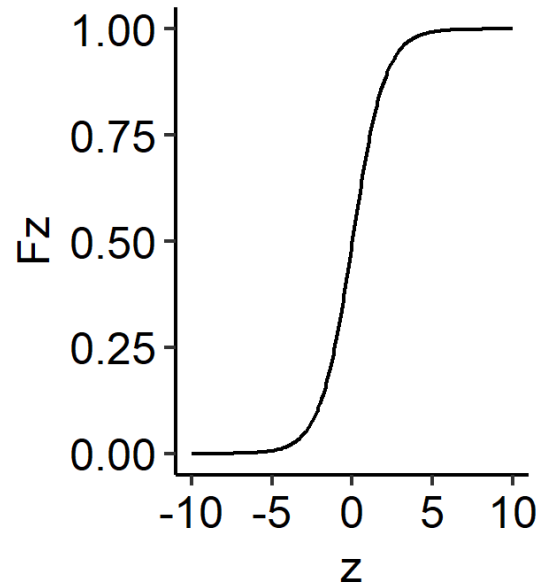
Recap: Standard Logistic Distribution

PDF (f_Z)



$$f_Z(z) = \frac{e^z}{(1 + e^z)^2}$$

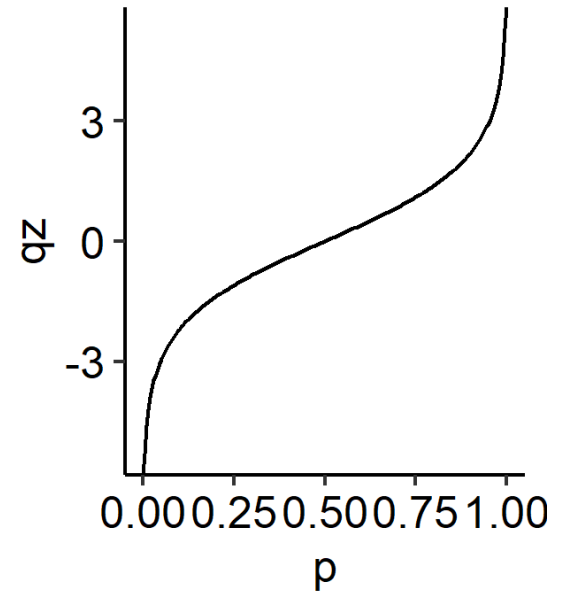
CDF (F_Z)



$$F_Z(z) = \frac{1}{1 + e^{-z}}$$

$$F_Z(z) = \text{expit}(z)$$

Quantiles function (F_Z^{-1})



$$F_Z^{-1}(p) = \log \frac{p}{1 - p}$$

$$F_Z^{-1}(p) = \log(\text{odds}) \\ = \text{logit}(p)$$

Interpretation with odds

- Probability for event

- $P(Y = 1|x) = p(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$
- $odds(x) := \frac{p(x)}{1 - p(x)} = e^{\beta_0 + \beta_1 x} \in [-\infty, +\infty]$
- $\log(odds(x)) = \beta_0 + \beta_1 x$

	1	2
bet365	1.83	1.83
1xBET	1.86	1.92
betway	1.8	2
888sport	1.8	1.95
bwin	1.8	1.95
William Hill	1.8	1.91
bet-at-home	1.78	1.94
UNIBET	1.8	1.95
betvictor	1.8	2

Bonus
Bonus
\$100
£30
£30
£10
Bonus
£20
£40
£30

<https://www.wincomparator.com/roger-federer-id2930-rafael-nadal-id2905/>

- Interpretation (odds):

- Examples

- Odds : „Prob for winning“ / „Prob for not winning“
 - Federer against Nadal 2:1 Two times more likely that Federe wins
- Odds : „Probability of broken o-ring“ vs. Non-broken
 - 8:1

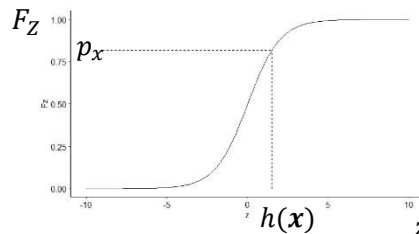
Odds ratios

- $\log(odds(x)) = \beta_0 + \beta_1 x$
- Q: How does the odds change with x
- Example Tennis (Federer against Nadal)
 - $x = 0$ (not grass) $odds(x = 0) = \exp(\beta_0) = 1.2$
 - $x = 1$ (grass) $odds(x = 1) = \exp(\beta_0 + \beta_1 \cdot 1) = 2$
 - Odds Ratio: $OR_{x=0 \rightarrow x=1} = \frac{odds(x=1)}{odds(x=0)} = \exp(\beta_1 \cdot 1) = 2/1.2$
 - $\log(OR_{x=0 \rightarrow x=1}) = \beta_1 = 2/1.2$
- Works also for continuous x important is just Δx (not absolute value)

The logistic regression model

Model for the cut-point:

$$h(\mathbf{x}) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p \rightarrow$$



$$\begin{aligned} p_{x_i} &= F_Z(h(\mathbf{x}_i)) \\ &\Updownarrow \\ h(\mathbf{x}_i) &= F_Z^{-1}(p_{x_i}) \end{aligned}$$

Using the **logistic distribution**: $F_Z(z) = F_L(z) = \frac{1}{1 + e^{-z}} \Leftrightarrow F_Z^{-1}(p) = \log \frac{p}{1-p}$

$$h(\mathbf{x}) = F_L^{-1}(p_x) = \log \left(\frac{p_x}{1-p_x} \right) = \log(\text{odds}(\mathbf{x})) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

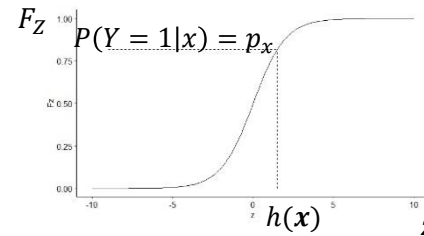
$$\text{odds}(\mathbf{x}_i) = e^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p} = e^{\beta_0} \cdot e^{\beta_1 x_1} \cdot \dots \cdot e^{\beta_k(x_k+1)} \cdot \dots \cdot e^{\beta_p x_p}$$

$$p_{x_i} = F_L(h(\mathbf{x}_i)) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p)}}$$

Model parameters are interpretable as log-odds ratio

Model for cut-point:

$$\log(\text{odds}(x)) = \beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p$$



The **coefficient** β_k as the **log-odds-ratio** for $Y = 1$ when comparing a situation where x_k is increases by 1 unit (while fixing all other variables) with the situation before increasing x_k

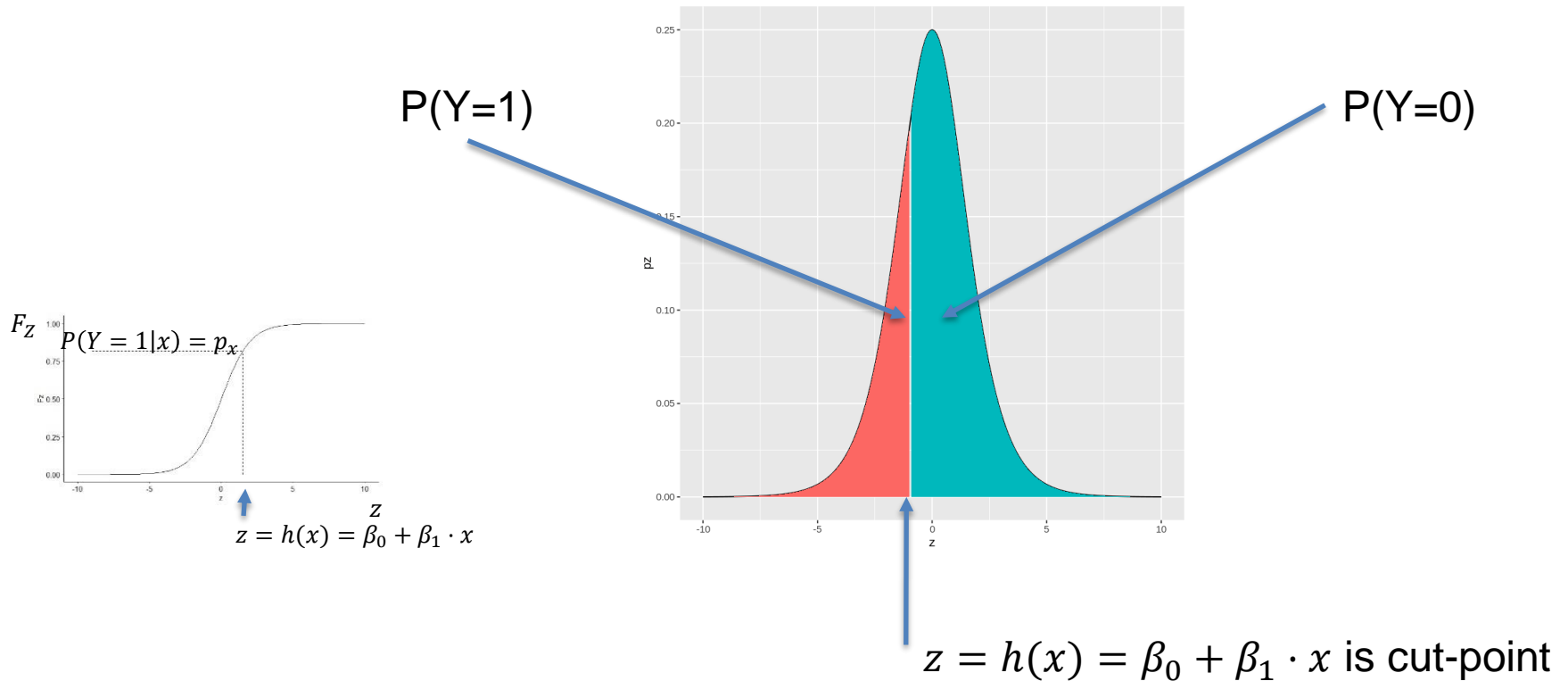
$$\log(\text{OR}_k) = \log\left(\frac{\text{odds}(x_1, \dots, x_k+1, \dots, x_p)}{\text{odds}(x_1, \dots, x_k, \dots, x_p)}\right) = \log\left(\frac{e^{\beta_0} \cdot e^{\beta_1 x_1} \cdots e^{\beta_k(x_k+1)} \cdots e^{\beta_p x_p}}{e^{\beta_0} \cdot e^{\beta_1 x_1} \cdots e^{\beta_k x_k} \cdots e^{\beta_p x_p}}\right) = \log(e^{\beta_k}) = \beta_k$$

$$\Rightarrow e^{\beta_k} = \text{OR}_{x_k \rightarrow x_k+1}$$

Logistic Regression as latent variable model

Idea:

A continuous latent (unobserved) variable z determines the probability to observe $Y = 1$



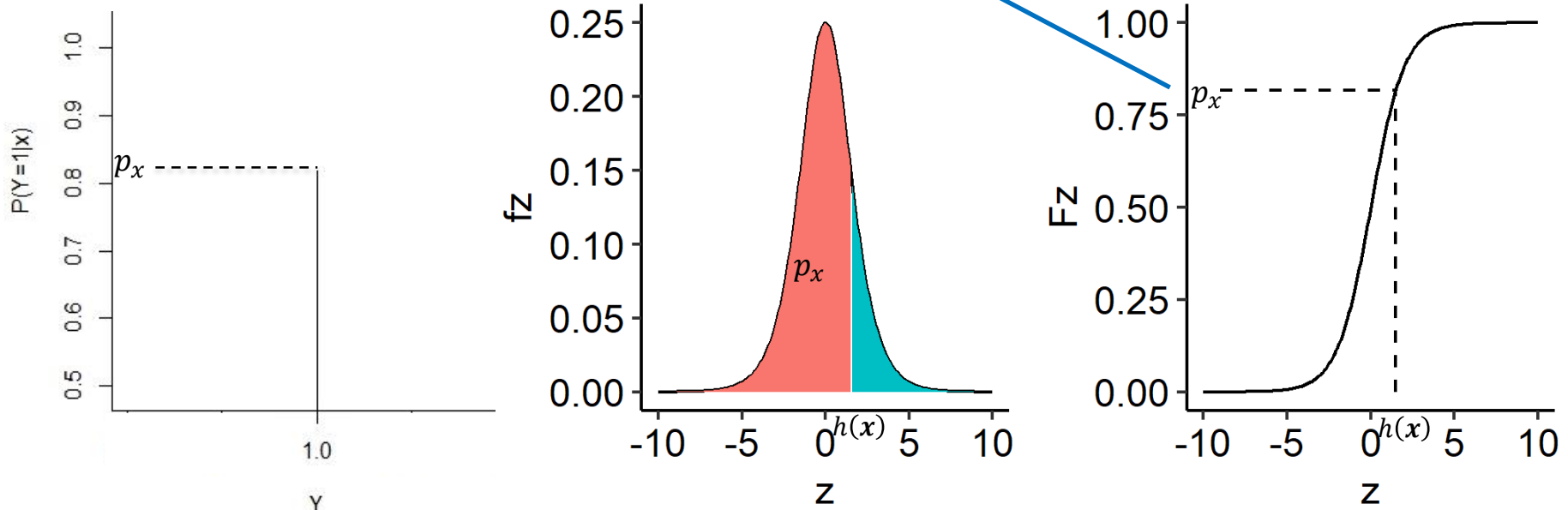
Logistic Regression as latent variable model

We fit a logistic regression model $(Y|x) \sim \text{Ber}(p_x)$ by minimizing the NLL

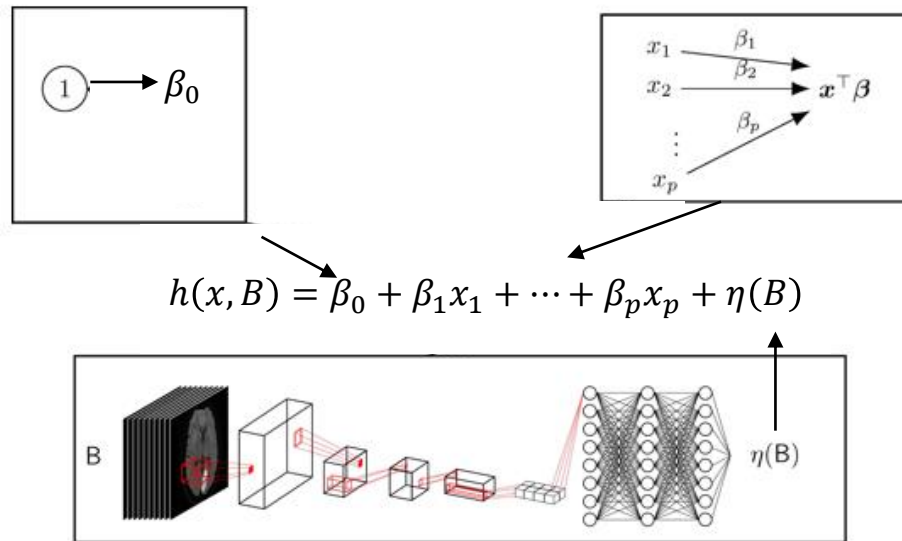
$$\text{NLL} = -\frac{1}{n} \sum_i \log(L_i), \quad \log(L_i) = \log(P(Y = y_i | x_i)) = y_i \cdot \log(p_{x_i}) + (1 - y_i) \cdot \log(1 - p_{x_i})$$

In DL we minimize the NLL by finding the optimal β via SGD (stochastic gradient decent)

$$h(x) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p \rightarrow p_{x_i} = P(Y = 1|x) = F_Z(h(x_i))$$



Semi-structured logistic regression via DL approach



Jointly train all NNs by minimizing the NLL loss:

$$\text{NLL} = -\frac{1}{n} \sum_i \log(L_i), \quad L_i = F_L(h(x, B)) = \frac{1}{1 + e^{-h(x, B)}}$$

Parameters are still interpretable as log-odds ratios: $\hat{\beta}_i = \log(\text{OR}_{x_i \rightarrow (x_i+1)})$

Ordinal regression

(more than two levels)

Ordinal Data (Examples)

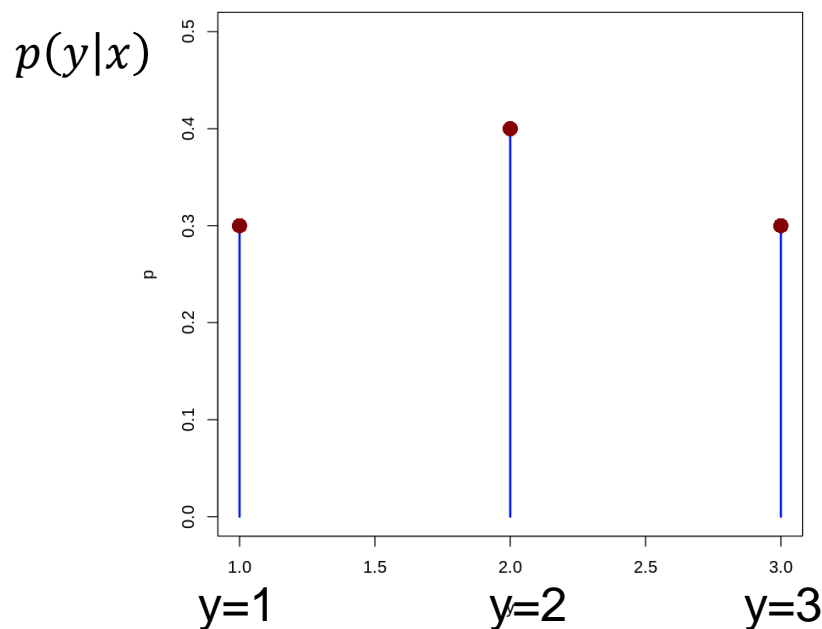
- Amazon Ratings ★★☆☆☆
- Medical ratings
 - Reopatiy Score
 - Neuroscore
- Wine quality (score between 0 and 10)
 - UCI-Data* set (N=4898) from (Cortez et. al. 2009) with 11 covariates like:
 - fixed acidity, residual sugar, sulphates, ..., alcohol
- Age group (UKTFace)
 - Images as high dimensional features



«Less than numerical data, more than just classes / categories»

Goal of ordinal regression

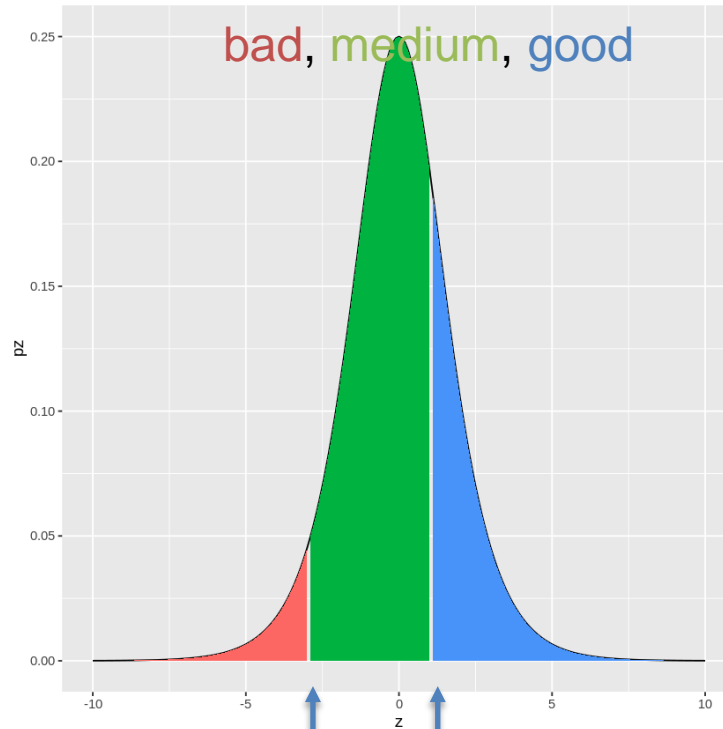
- Get the conditional probability distribution CPD for a given x .
 - $p(y|x)$
- Simple example
 - x alcohol content in wine
 - $y=1,2,3$ corresponds to bad, medium, good



Distribution for fixed x

Proportional odds-model

Latent variable model can be extended for more than 2 levels



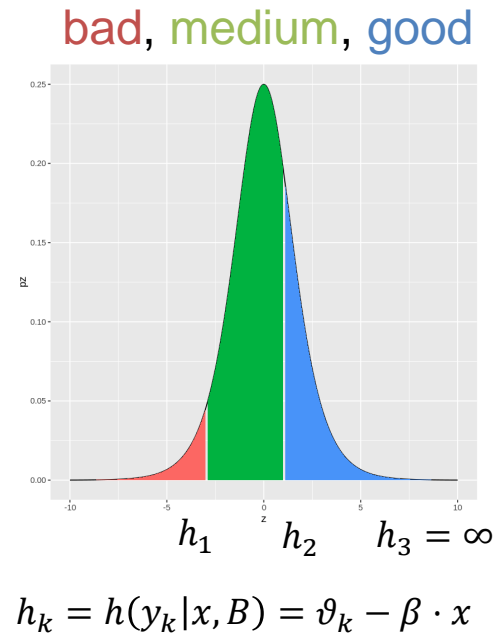
If x changes, all cut points are moving together.

$$h_1 = \vartheta_1 - x \cdot \beta \quad h_2 = \vartheta_2 - x \cdot \beta$$

Cut-points modeled via: $h_i = \vartheta_i - x \cdot \beta$

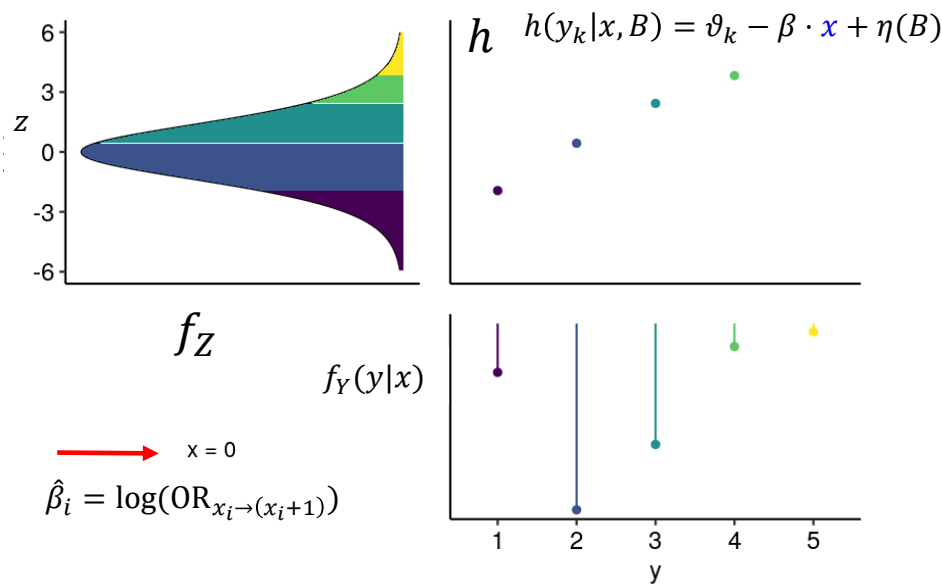
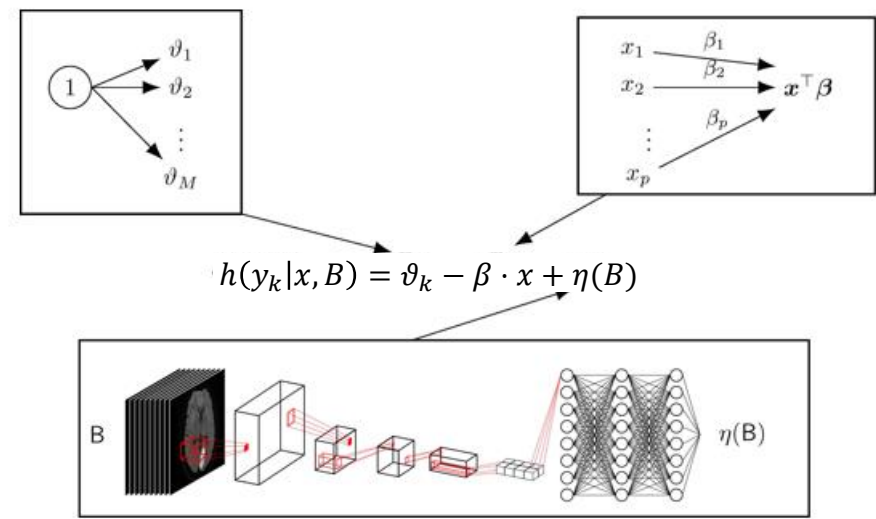
Proportional odds model (interpretation)

- Since we have more levels, odds are now defined as
 - $odds(Y > y_k | x) = \frac{P(Y > y_k | x)}{P(Y \leq y_k | x)}$
 - Example $k = 2$ odds prob for good / prob less than good
- Changes in odds with x again with odds ratio
 - $OR_{x \rightarrow x+1} = \frac{odds(Y > y_k | x+1)}{odds(Y > y_k | x)} = e^\beta$ (not depending on k)
- Example $x = alcohol$



- **Questions:** Does alcohol make the wine taste better?
- What is β in $h_i = \vartheta_i - x \cdot \beta$

Semi-structured ordinal regression (ONTRAMs) via DL



Jointly train NNs with NLL loss:

$$\text{NLL} = -\frac{1}{n} \sum_i \log(L_i)$$

$$\mathcal{L}_i(h; y_{ki}, \mathbf{x}_i) = \mathbb{P}(Y = y_{ki} | \mathbf{x}_i) = F_Y(y_{ki} | \mathbf{x}_i) - F_Y(y_{(k-1)i} | \mathbf{x}_i)$$

$$= F_Z(h(y_{ki} | \mathbf{x}_i)) - F_Z(h(y_{(k-1)i} | \mathbf{x}_i)).$$

UKT-Face

Ordinal neural network regression to model ordinal outcome is the **age-category**.

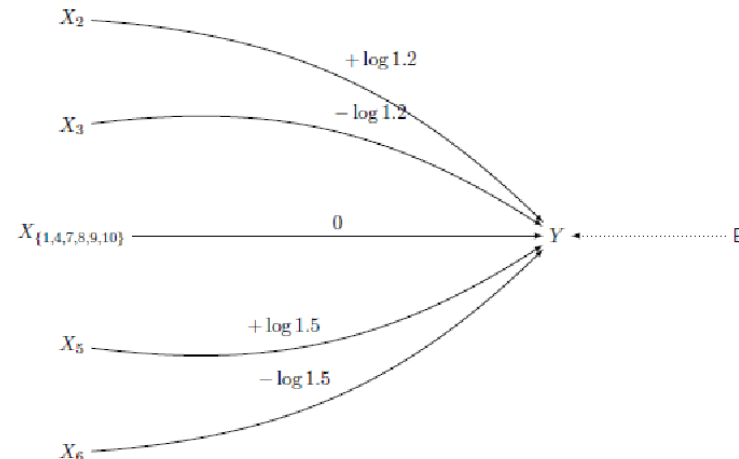
Modeled transformation function: $h(y_k|x, B) = \vartheta_k - x^\top \beta - \eta(B)$

Image input Data (one random example for each of the 7 age-categories):

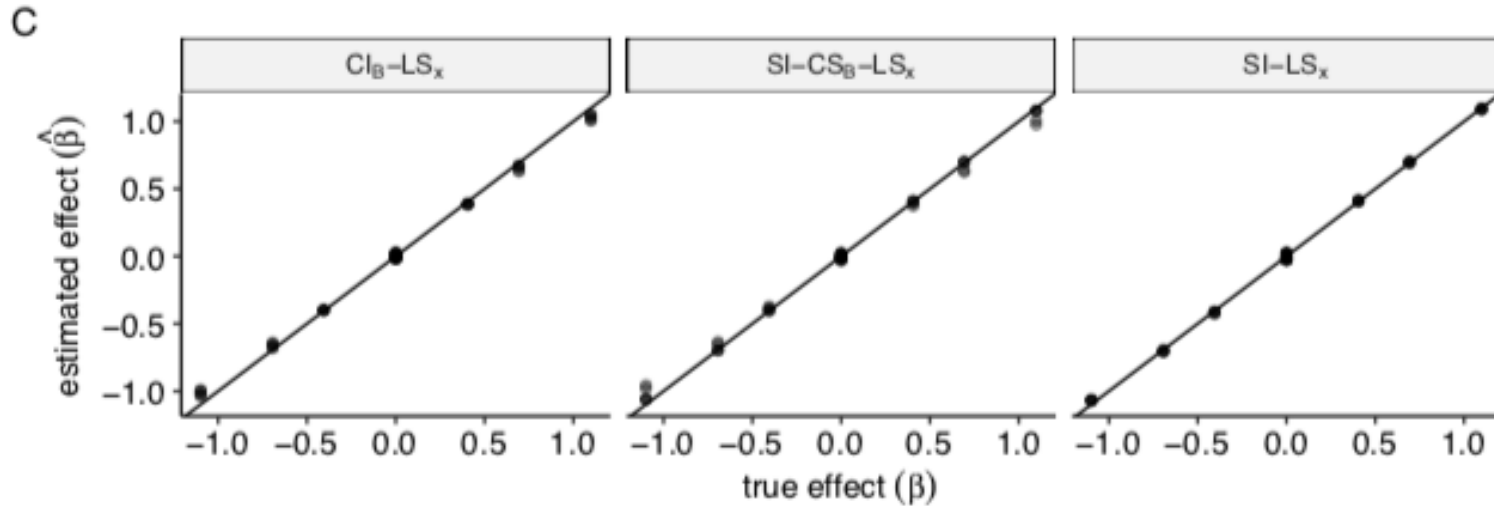
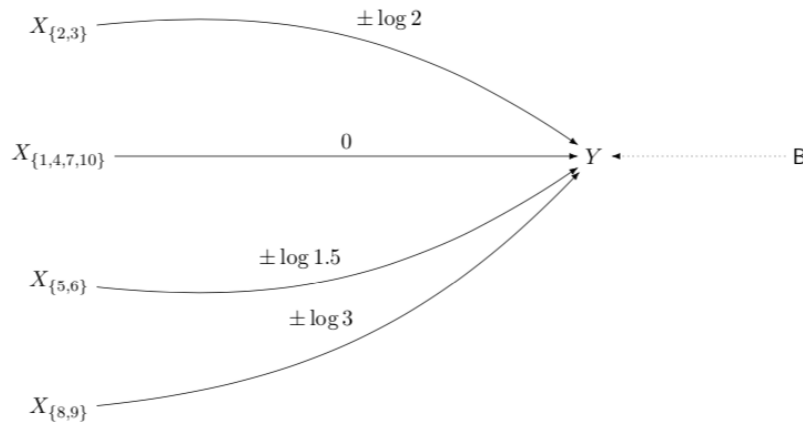


Tabular Co-Variates:

- Gender
- Ethnicity
- 10 simulated covariates with known effect sizes



UKT-Face: Can, we recover simulated tabular data?

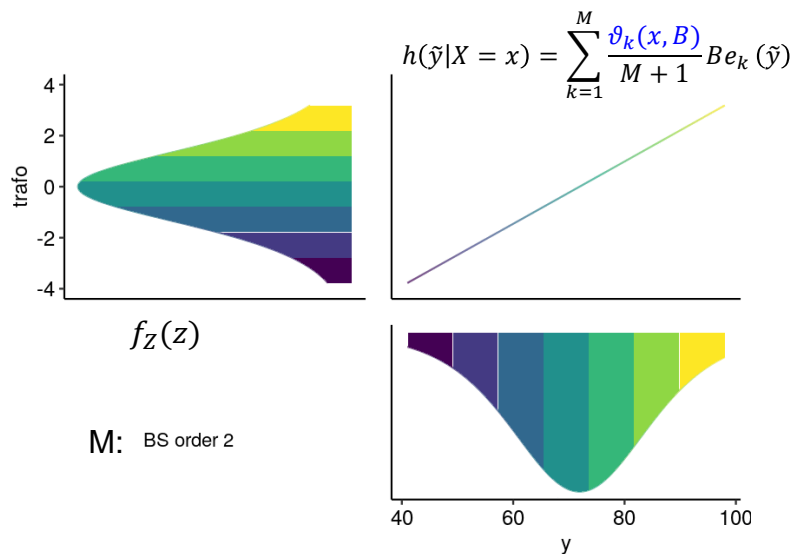


Recovering the simulated effects.
Practical example. Ongoing...

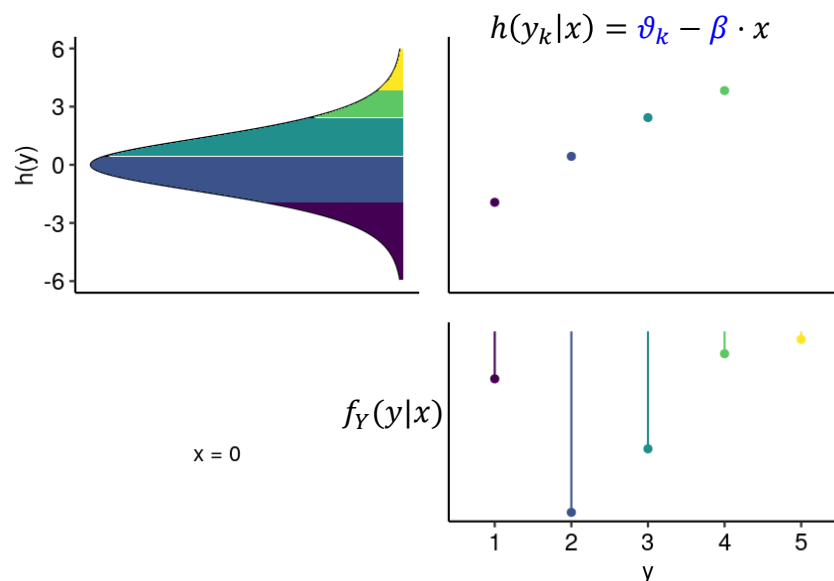
Transformation models for flexible distributional regression

- In traditional regression, the family of the conditional outcome distribution $F_Y(y|x, B \dots)$ is predefined, and the parameters of this distribution are fitted.
- In TMs the conditional outcome distribution $F_Y(y|x, B \dots)$ is achieved by transforming a parameter-free distribution $F_Z(z = h(y|x, B \dots))$ requiring to **estimate the parameters of the conditional transformation function $h(y|x, B \dots)$** .

Continuous outcome \rightarrow continuous h



Discrete outcome \rightarrow discrete h

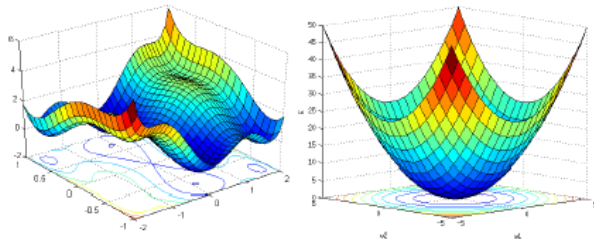


Ensembling

Classical “Deep Ensembles” as used in deep learning

Refitting a deep NN with same data but new random initialization yields slightly different parameter estimates.

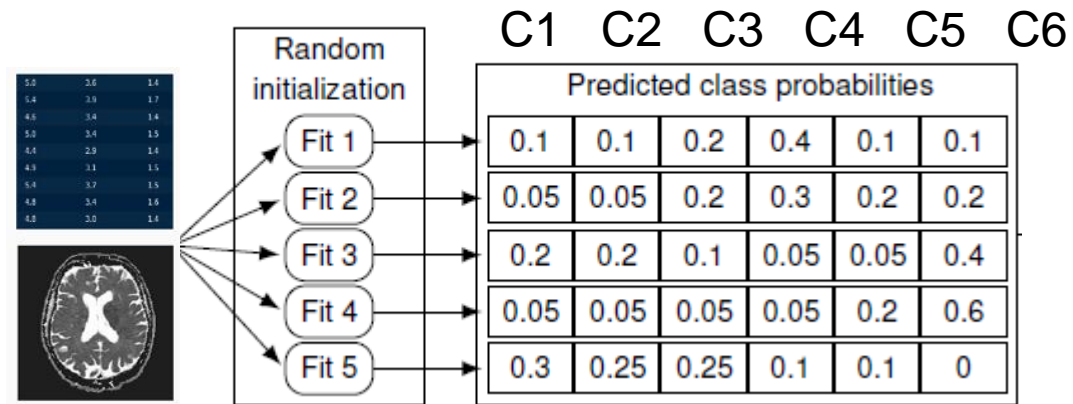
Optimization is non-convex:



Reasons:

- Over-parametrization
- Training on mini-batches
- Random weight initialization

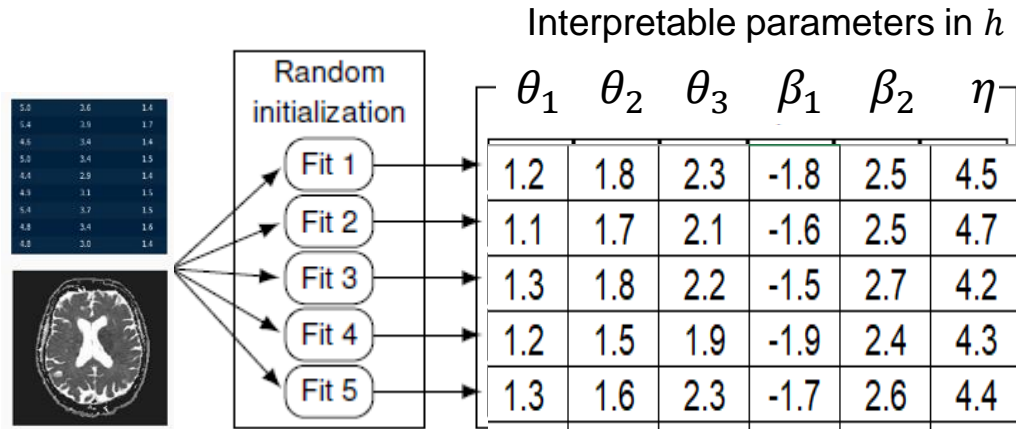
Constructing deep ensembles:



- Column mean:
point estimate for class probability
- Column standard deviation:
uncertainty of point estimate

Deep Ensembles have a higher prediction performances than an average single NN.

Interpretable transformation Ensembles



- Column mean:
point estimate for parameters
- Column standard deviation:
uncertainty for parameters

Transformation Ensembles yield interpretable parameters along with uncertainty estimates and have higher performance than single deep transformation models.

$$\bar{F}_M^t = F_Z \left(\frac{1}{M} \sum_{m=1 \dots M} h_m(y|x, B) \right) = F_Z(\bar{\theta}_k - x^\top \bar{\beta} - \bar{\eta}(B))$$

Summary on ordinal neural transformation networks (ONTRAMs)

- ONTRAM allows to work with image data and tabular predictors
- ONTRAM allows for the same interpretability than statistical ordinal regression
 - $F_Z = \text{logistic}$ and $h(y_k|x) = \vartheta_k + \sum_{l=1}^p \beta_l \cdot x_l$ with log-odds interpretation of β
- ONTRAM has the high prediction performance of DL models
- Via transformation ensembles, the prediction performance can be further improved and the uncertainty of the interpretable parameters are quantified.