# Scientific Machine Learning
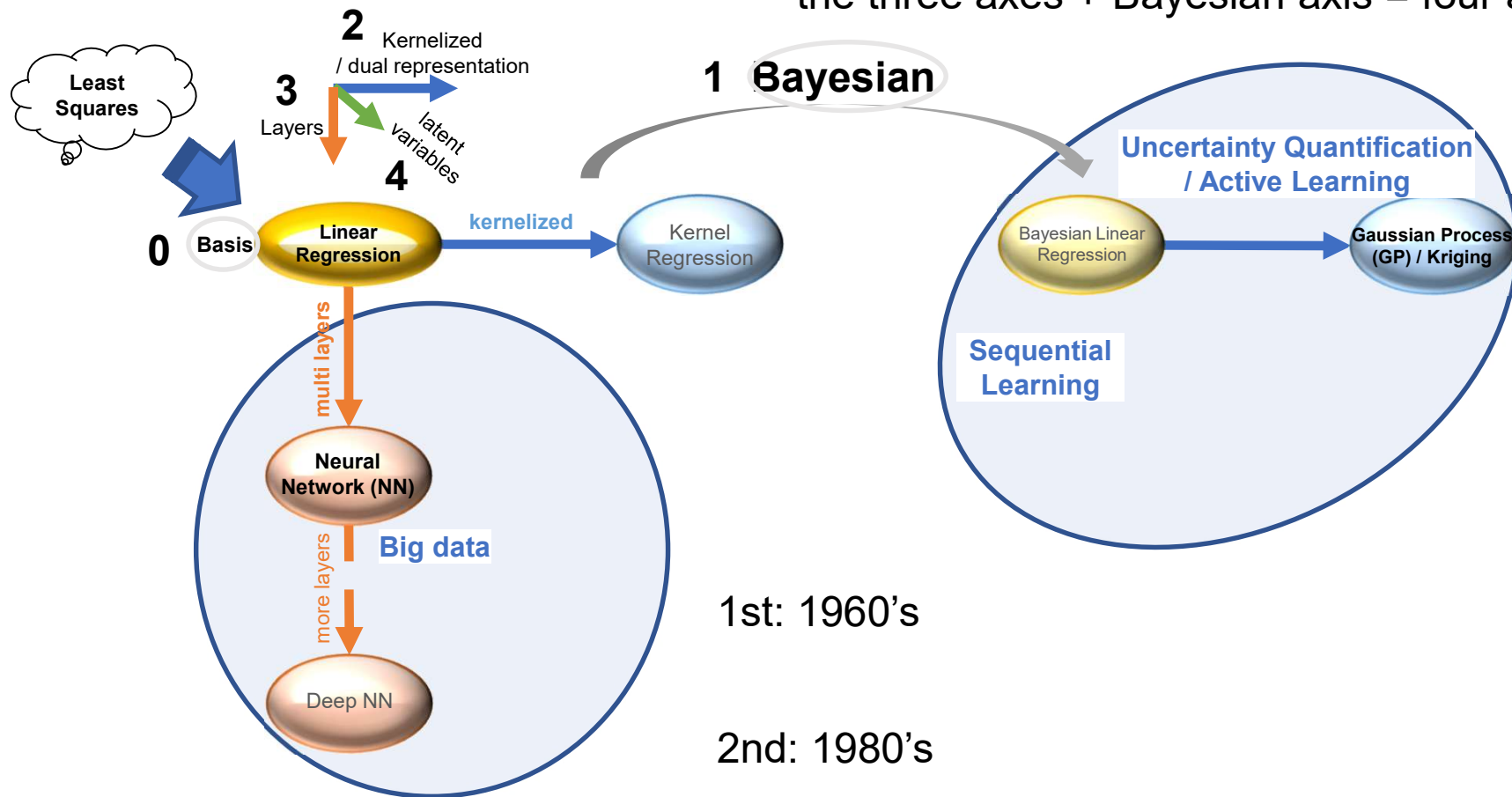*Lecture 10: Generalized Linear Model (GLM), Neural Network*

Dr. Daigo Maruyama

Prof. Dr. Ali Elham

# Key Components

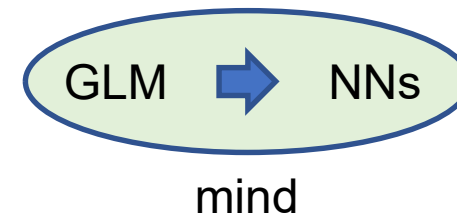the three axes + Bayesian axis = four axes



1st: 1960's

2nd: 1980's

3rd: 2010's - now

# Lecture content

- From linear regression to generalized linear model (GLM)
  - Classification (opposite: regression) in applications

- To neural networks (NNs)

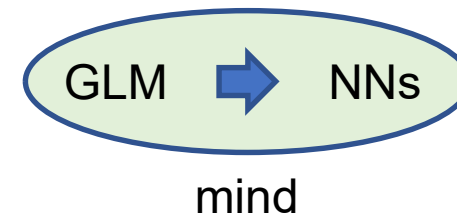- Technical Issues in neural networks

GLM ➡ NNs

mind

The lecture of this time partially follows the Chapters 4 and 5 of the book:
Christopher M. Bishop "Pattern Recognition And Machine Learning" Springer-Verlag (2006)
The name of this book is shown as "PRML" when it is referred in the slides.

The lecture slides contains a few recent topics and evidence in neural networks.

Technische
Universität
Braunschweig

IFL

# Lecture content

- From linear regression to generalized linear model (GLM)
  - Classification (opposite: regression) in applications

GLM ➡ NNs

mind

Technische
Universität
Braunschweig

IFL

# Linear Regression (REVIEW)



based on PRML, p. 29

$$\mathbf{x} = \left(x^{(1)}, x^{(2)}, \cdots, x^{(N)}\right)^{\mathrm{T}}$$
$$\mathbf{t} = \left(t^{(1)}, t^{(2)}, \cdots, t^{(N)}\right)^{\mathrm{T}}$$

$N$: sample size

Define a **regression model**

$$y(\boldsymbol{x}, \boldsymbol{w}) = \underbrace{\boldsymbol{w}^{\mathrm{T}} \boldsymbol{\phi}(\boldsymbol{x})}$$

Linear regression model

Define a **Probabilistic model**

$$p(t|\boldsymbol{x}, \mu, \sigma) = \mathcal{N}(t|\mu(\boldsymbol{x}), \sigma^2)$$

$$\mu(\boldsymbol{x}) = y(\boldsymbol{x}, \boldsymbol{w})$$

$\mu(\boldsymbol{x})$: the regression model

$$p(t|\boldsymbol{x}, \boldsymbol{w}, \sigma) = \mathcal{N}(t|y(\boldsymbol{x}, \boldsymbol{w}), \sigma^2)$$
$$= \mathcal{N}(t|\boldsymbol{w}^{\mathrm{T}}\boldsymbol{\phi}(\boldsymbol{x}), \sigma^2)$$

Technische Universität Braunschweig

IFL

# Linear Regression (REVIEW)



based on PRML, p. 29

$$\mathbf{x} = \left(\boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}, \cdots, \boldsymbol{x}^{(N)}\right)^{\mathrm{T}}$$
$$\mathbf{t} = \left(t^{(1)}, t^{(2)}, \cdots, t^{(N)}\right)^{\mathrm{T}}$$

$N$: sample size

Define a **regression model**

$$y(\boldsymbol{x}, \boldsymbol{w}) = \underbrace{\boldsymbol{w}^{\mathrm{T}}\boldsymbol{\phi}(\boldsymbol{x})}$$

Linear regression model

Define a **Probabilistic model**

$$p(t|\boldsymbol{x}, \mu, \sigma) = \mathcal{N}(t|\mu(\boldsymbol{x}), \sigma^2)$$

$$\mu(\boldsymbol{x}) = y(\boldsymbol{x}, \boldsymbol{w})$$

$\mu(\boldsymbol{x})$: the regression model

After learning $\boldsymbol{w}$ by data: **(point estimate)**

$$\hat{\mu}(\boldsymbol{x}_{new}) = \widehat{\boldsymbol{w}}^{\mathrm{T}}\boldsymbol{\phi}(\boldsymbol{x}_{new})$$

Technische Universität Braunschweig

IFL

# Linear Regression (extension to multiple output)

$$y(x, w) = w^{\mathrm{T}}\phi(x) = w_1\phi_1(x) + w_2\phi_2(x) + \cdots + w_M\phi_M(x)$$

input    output

$\phi_1(x)$ $\quad w_1$

$\phi_2(x)$ $\quad w_2$

$\phi_3(x)$ $\quad w_3$

$\phi_4(x)$ $\quad w_4$

$\phi_5(x)$ $\quad w_5$

$\boldsymbol{w}$

$\phi(x)$ $\qquad y$

$M = 5$

$$\widehat{w} = (\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\Phi})^{-1}\boldsymbol{\Phi}^{\mathrm{T}}\mathbf{t}$$

$$= \begin{pmatrix} & \cdots & \\ \vdots & \ddots & \vdots \\ & \cdots & \end{pmatrix}\begin{pmatrix} t^{(1)} \\ \vdots \\ t^{(N)} \end{pmatrix}$$

The input $x$ could be multidimensional.
How about the output $t$?
(we have been considering always one dimensional output)

$y(x_0, \mathbf{w})$

$2\sigma$

● : Data

$$\mathbf{x} = \left(x^{(1)}, x^{(2)}, \cdots, x^{(N)}\right)^{\mathrm{T}}$$

$$\mathbf{t} = \left(t^{(1)}, t^{(2)}, \cdots, t^{(N)}\right)^{\mathrm{T}}$$

$x_0$ $\qquad x$

# Linear Regression (multiple output)

$$y(x, w_1) = w_1^T \phi(x) = w_{11}\phi_1(x) + w_{21}\phi_2(x) + \cdots + w_{M1}\phi_M(x)$$

input     output

$\phi_1(x)$    $w_{11}$      $y_1$

$\phi_2(x)$    $w_{21}$

     $w_{31}$

$\phi_3(x)$    $w_{41}$

$\phi_4(x)$    $w_{51}$

$\phi_5(x)$

       $w_1$

$\phi(x)$        $y$

$$\hat{w}_1 = \left(\Phi^T\Phi\right)^{-1}\Phi^T t_1$$

$$= \begin{pmatrix} & \cdots & \\ \vdots & \ddots & \vdots \\ & \cdots & \end{pmatrix} \begin{pmatrix} t_1^{(1)} \\ \vdots \\ t_1^{(N)} \end{pmatrix} = \begin{pmatrix} \hat{w}_{11} \\ \vdots \\ \hat{w}_{51} \end{pmatrix}$$

Consider a multiple output case:

$$\mathbf{T} = \left(t^{(1)}, t^{(2)}, \cdots, t^{(N)}\right)^T$$
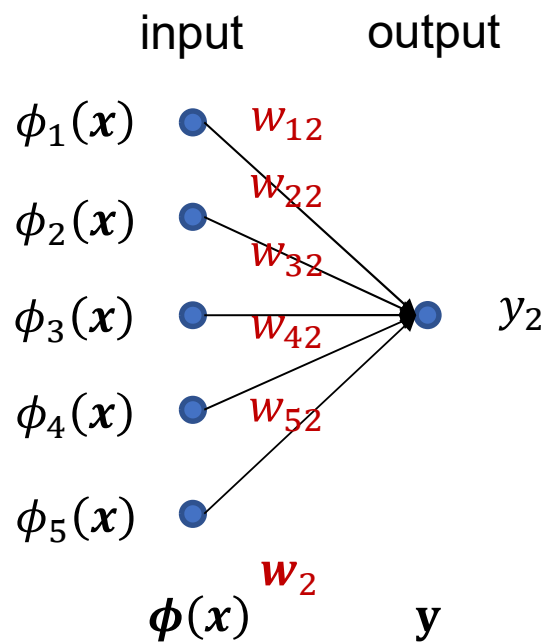
by focusing on the $i$ th component

$$\mathbf{t}_i = \left(t_i^{(1)}, t_i^{(2)}, \cdots, t_i^{(N)}\right)^T$$

At first, focusing on the first component $y_1$ in $\mathbf{y}$

Technische Universität Braunschweig

IFL

# Linear Regression (multiple output)

$$y(x, \boldsymbol{w}_2) = \boldsymbol{w}_2{}^{\mathrm{T}}\boldsymbol{\phi}(x) = w_{12}\phi_1(x) + w_{22}\phi_2(x) + \cdots + w_{M2}\phi_M(x)$$

input    output

$\phi_1(x)$   $w_{12}$

    $w_{22}$

$\phi_2(x)$

    $w_{32}$

$\phi_3(x)$   $w_{42}$    $y_2$

$\phi_4(x)$   $w_{52}$

$\phi_5(x)$

    $\boldsymbol{w}_2$

$\boldsymbol{\phi}(x)$     $\mathbf{y}$

$$\widehat{\boldsymbol{w}}_2 = \left(\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\Phi}\right)^{-1}\boldsymbol{\Phi}^{\mathrm{T}}\mathbf{t}_2$$

$$= \begin{pmatrix} & \cdots & \\ \vdots & \ddots & \vdots \\ & \cdots & \end{pmatrix}\begin{pmatrix} t_2{}^{(1)} \\ \vdots \\ t_2{}^{(N)} \end{pmatrix} = \begin{pmatrix} \widehat{w}_{12} \\ \vdots \\ \widehat{w}_{52} \end{pmatrix}$$

Consider a multiple output case:

$$\mathbf{T} = \left(\boldsymbol{t}^{(1)}, \boldsymbol{t}^{(2)}, \cdots, \boldsymbol{t}^{(N)}\right)^{\mathrm{T}}$$

by focusing on the $i$ th component

$$\mathbf{t}_i = \left(t_i{}^{(1)}, t_i{}^{(2)}, \cdots, t_i{}^{(N)}\right)^{\mathrm{T}}$$

Then, focusing on the second component $y_2$ in $\mathbf{y}$

# Linear Regression (multiple output)

$$y(x, w_3) = w_3^{\mathrm{T}} \phi(x) = w_{13}\phi_1(x) + w_{23}\phi_2(x) + \cdots + w_{M3}\phi_M(x)$$

input    output

$\phi_1(x)$   $w_{13}$

$\phi_2(x)$   $w_{23}$

    $w_{33}$

$\phi_3(x)$   $w_{43}$

$\phi_4(x)$   $w_{53}$

$\phi_5(x)$     $y_3$

$\boldsymbol{w_3}$

$\boldsymbol{\phi(x)}$     $\mathbf{y}$

$$\widehat{w}_3 = \left(\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\Phi}\right)^{-1}\boldsymbol{\Phi}^{\mathrm{T}}\mathbf{t}_3$$

$$= \begin{pmatrix} & \cdots & \\ \vdots & \ddots & \vdots \\ & \cdots & \end{pmatrix} \begin{pmatrix} t_3^{(1)} \\ \vdots \\ t_3^{(N)} \end{pmatrix} = \begin{pmatrix} \widehat{w}_{13} \\ \vdots \\ \widehat{w}_{53} \end{pmatrix}$$

Consider a multiple output case:

$$\mathbf{T} = \left(t^{(1)}, t^{(2)}, \cdots, t^{(N)}\right)^{\mathrm{T}}$$

by focusing on the $i$ th component

$$\mathbf{t}_i = \left(t_i^{(1)}, t_i^{(2)}, \cdots, t_i^{(N)}\right)^{\mathrm{T}}$$

<u>The output $\mathbf{y}$ is three dimensional in this example.</u>

Technische Universität Braunschweig

IFL

# Linear Regression (multiple output)

$$\mathbf{y}(x, \mathbf{W}) = \mathbf{W}^{\mathrm{T}} \boldsymbol{\phi}(x)$$

The tuned parameter $\widehat{w}_i$ ($i$=1,2,3) can be obtained at the same time as $\widehat{\mathbf{W}}$.

input        output

$\phi_1(x)$    $y_1$

$\phi_2(x)$

$\phi_3(x)$   $\mathbf{W}$   $y_2$

$\phi_4(x)$

$\phi_5(x)$    $y_3$

$\boldsymbol{\phi}(x)$    $\mathbf{y}$

Looks like a neural network?

$$\widehat{\mathbf{W}} = \left(\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\Phi}\right)^{-1}\boldsymbol{\Phi}^{\mathrm{T}}\mathbf{T}$$

$$= \begin{pmatrix} & \cdots & \\ \vdots & \ddots & \vdots \\ & \cdots & \end{pmatrix} \begin{pmatrix} t_1^{(1)} & \cdots & t_3^{(1)} \\ \vdots & \ddots & \vdots \\ t_1^{(N)} & \cdots & t_3^{(N)} \end{pmatrix} = \begin{pmatrix} \widehat{w}_{11} & \cdots & \widehat{w}_{13} \\ \vdots & \ddots & \vdots \\ \widehat{w}_{51} & \cdots & \widehat{w}_{53} \end{pmatrix}$$

$$\widehat{\mathbf{W}} = (\widehat{w}_1, \widehat{w}_2, \widehat{w}_3)$$

Consider a multiple output case:

$$\mathbf{T} = \left(t^{(1)}, t^{(2)}, \cdots, t^{(N)}\right)^{\mathrm{T}}$$

This is still the **linear regression model**.

The objective is to obtain $\widehat{\mathbf{W}}$ or $p(\mathbf{W})$ (by using data).

Technische Universität Braunschweig

IFL

# Generalized Linear Model (GLM)

Linear regression

$$y(x, w) = w^{\mathrm{T}} \phi(x)$$

the single output form

$$\hat{w} = \left(\Phi^{\mathrm{T}}\Phi\right)^{-1} \Phi^{\mathrm{T}} t$$

**Analytical solution**:
The important property of linear regression

Generalized linear model (GLM)

$$y(x, w) = f\left(w^{\mathrm{T}} \phi(x)\right)$$

$f$: **nonlinear function**

To obtain $\hat{w}$, we need
**numerical optimization** tools.

Even if we lose this nice property (analytical solution), we will obtain further benefit!

## Generalized Linear Model (GLM)

$$y(x, w) = f\left(w^{\mathrm{T}}\phi(x)\right)$$

General ideas of using this form:

1. Transformation to probability output (0-1)

   **Classification**

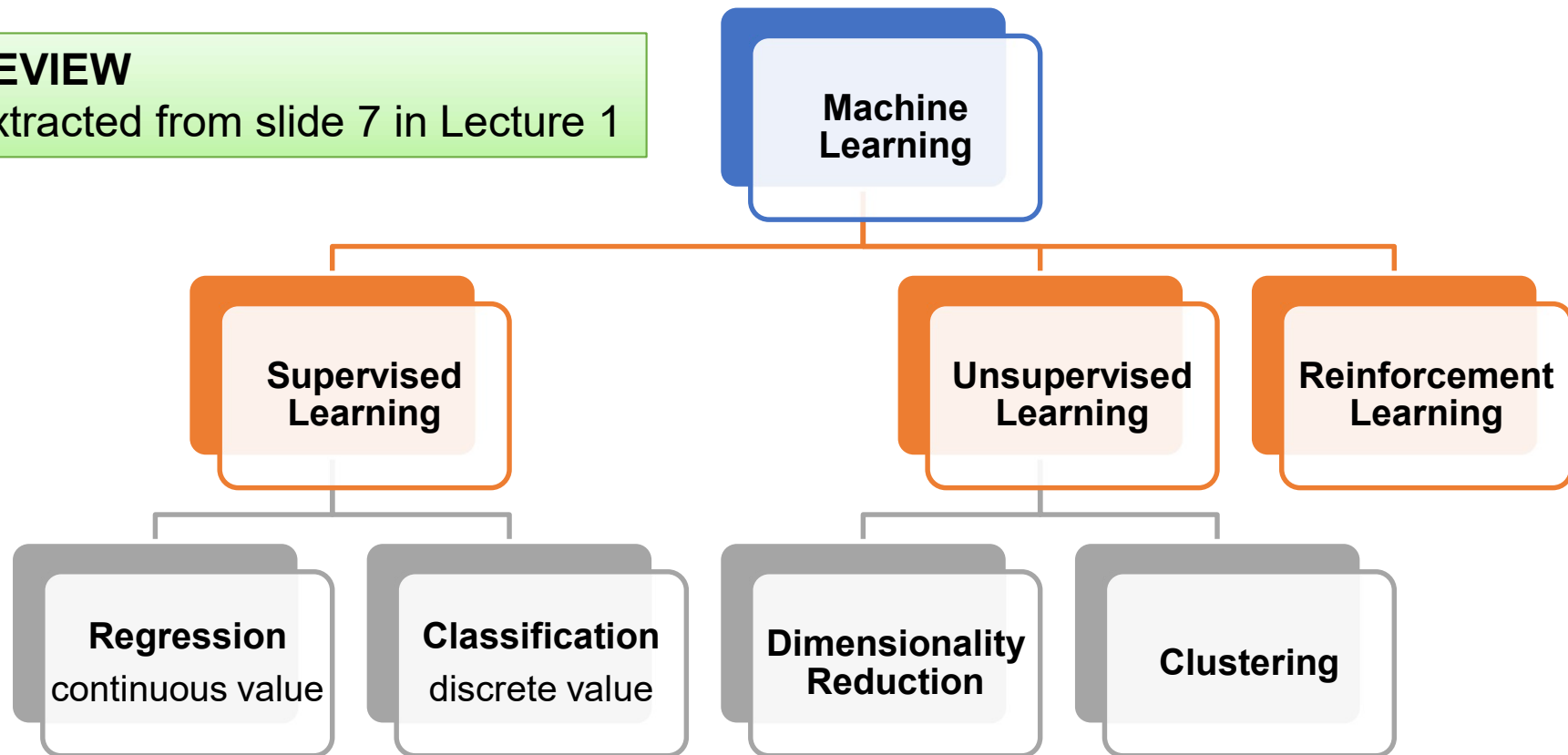2. Efficient construction of a complex function          1+2 is possible.

   **Neural Networks**

Technische
Universität
Braunschweig

IFL

# Machine Learning Classification by Use/Application



**REVIEW**
Extracted from slide 7 in Lecture 1

**Machine Learning**

- **Supervised Learning**
  - **Regression** continuous value
  - **Classification** discrete value
- **Unsupervised Learning**
  - **Dimensionality Reduction**
  - **Clustering**
- **Reinforcement Learning**

In this course, machine learning classification is done by **methods and their concepts**.

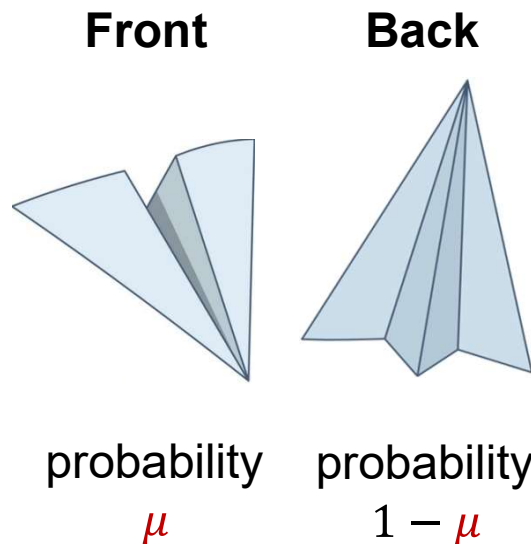➡ Then the use/application is naturally derived/understood.

Technische Universität Braunschweig

IFL
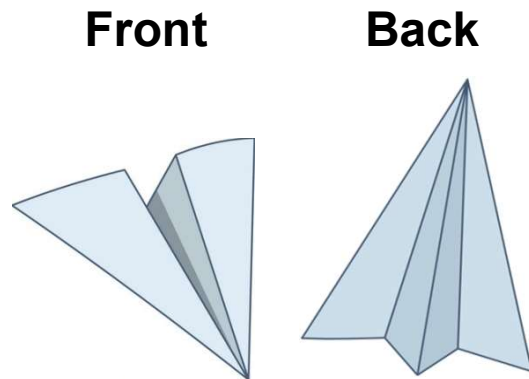
# Classification (Example: Binomial case)

Another example (the problem slightly change):

You threw three different paper planes.

**Observed Data**: all the three planes were front.

Shape1　Shape2　Shape3
$x^{(1)}$　$x^{(2)}$　$x^{(3)}$

**Front**　　**Back**

What is the probability $\mu(x)$?

We want to predict $p(\mu(x))$.

probability　probability
$\mu(x)$　　$1 - \mu(x)$　　$0 \leq \mu(x) \leq 1$　　in the regression (curve fitting):

$$\mu(x) = y(x, w) = w^{\mathrm{T}}\phi(x)$$

Technische
Universität
Braunschweig

IFL

# Classification (Example: Binomial case)

$f$: sigmoid function

$$\mu(x) = y(x, w) = f\left(w^{\mathrm{T}}\phi(x)\right)$$

$$w^{\mathrm{T}}\phi(x) \in [-\infty, \infty]$$

$$f(a) = \frac{1}{1 + \exp(-a)}$$

nonlinear transformation

$$f\left(w^{\mathrm{T}}\phi(x)\right) \in [0,1]$$

to a probability like
$$p\left(w^{\mathrm{T}}\phi(x)\right) \in [0,1]$$



https://en.wikipedia.org/wiki/Sigmoid_function

$\phi(x)$: **Input** - shape of the paper plane

$w$: unknown (we want to determine by data)

$f$: sigmoid function

$y$: **Output** - front or back (1 or 0)

$$\hat{\mu}(x) = f\left(\hat{w}^{\mathrm{T}}\phi(x)\right)$$

The output $\hat{\mu}(x_{new})$ is interpreted as a probability when front.

Another example:

$\boldsymbol{\phi(x)}$: **Input** - weight, lung capacity, etc.

$w$: unknown (we want to determine by data)

$f$: sigmoid function

$y$: **Output** - disease or not (1 or 0)

## Classification (Example: Multiple case)

A new picture was given.

You want to classify it from the following three possibilities:

$$p(\text{"}dog\text{"}) = A?$$

$$p(\text{"}cat\text{"}) = B?$$

$$p(\text{"}fox\text{"}) = C?$$

$A + B + C$ has to be 1.

Technische
Universität
Braunschweig

IFL

# Classification (Example: Multiple case)

$f$: softmax function

$$\mu(\boldsymbol{x}) = y(\boldsymbol{x}, \boldsymbol{w}) = f\left(\boldsymbol{w}^{\mathrm{T}} \boldsymbol{\phi}(\boldsymbol{x})\right)$$

$$f(\boldsymbol{a})_k = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

$$\boldsymbol{w}^{\mathrm{T}} \boldsymbol{\phi}(\boldsymbol{x}) \in [-\infty, \infty]$$

nonlinear transformation

$$f\left(\boldsymbol{w}^{\mathrm{T}} \boldsymbol{\phi}(\boldsymbol{x})\right) \in [0,1]$$

to a probability like
$$p\left(\boldsymbol{w}^{\mathrm{T}} \boldsymbol{\phi}(\boldsymbol{x})\right) \in [0,1]$$

$\boldsymbol{\phi}(\boldsymbol{x})$: **Input** – RGB numbers at each pixel

$\boldsymbol{w}$: unknown (we want to determine by data)

$f$: softmax function

$y$: **Output** – "dog" or "cat" or "fox" (e.g.: 1 or 2 or 3)

## Classification

Let's summarize the processes by following **the process as usual**.

**1. Define a regression model**

**2. Define probabilistic model**

⬇

**A likelihood function**
(or posterior) ⬅ Data (and prior)

⬇

*The point estimate* (by MLE) of $w$

**optimization**

We do not consider the Bayesian approach as
*the probability distribution* of $w$
(Because analytical solutions are not expected).

Technische
Universität
Braunschweig

IFL

# Classification (2 classes: binary)

**The regression model**

$$\mu(\boldsymbol{x}) = y(\boldsymbol{x}, \boldsymbol{w}) = sigmoid\left(\boldsymbol{w}^{\mathrm{T}}\boldsymbol{\phi}(\boldsymbol{x})\right)$$

called logistic regression
(even though this is classification…)

**The probabilistic model**

$$p(t|\mu) = \mathrm{Bern}(t|\mu) = \mu^t(1-\mu)^{1-t}$$

$$= p(t|\boldsymbol{w})$$

> **Bernoulli distribution**
> (see Lecture 3, slide 19)

because <u>the output is discrete</u>
(see Lecture 3, slide 18).

The likelihood function

$$p(\mathbf{t}|\boldsymbol{w}) = \prod_{n=1}^{N} y(\boldsymbol{x}, \boldsymbol{w})^{t^{(n)}}\left(1 - y(\boldsymbol{x}, \boldsymbol{w})\right)^{1-t^{(n)}}$$

Data
$$\mathbf{x} = \left(\boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}, \cdots, \boldsymbol{x}^{(N)}\right)^{\mathrm{T}}$$
$$\mathbf{t} = \left(t^{(1)}, t^{(2)}, \cdots, t^{(N)}\right)^{\mathrm{T}}$$

$$\Rightarrow \quad E(\boldsymbol{w}) = -\ln p(\mathbf{t}|\boldsymbol{w})$$

Data: generated "<u>i.i.d.</u>"
(see Lecture 2, slide 27)

Take neg. log as usual $\quad \widehat{\boldsymbol{w}} = \underset{\boldsymbol{w}}{\mathrm{argmin}}\, E(\boldsymbol{w})$

# Classification (multiple classes)

**The regression model**

$$\mu(x) = y(x, w) = softmax\left(w^{\mathrm{T}}\phi(x)\right)$$

When "dog", "cat", "fox", $K = 3$.

**The probabilistic model**

> **Categorical distribution**
> (see Lecture 3, slide 21)

$$p(t|\mu) = \mathrm{Cat}(t|\mu) = \prod_{k=1}^{K} \mu_k{}^{x_k}$$

because <u>the output is discrete and multiple</u>
(see Lecture 3, slide 18).

The likelihood function

$$p(\mathbf{t}|w) = \prod_{n=1}^{N} y(x, w)^{t^{(n)}}\left(1 - y(x, w)\right)^{1 - t^{(n)}}$$

Data
$$\mathbf{x} = \left(x^{(1)}, x^{(2)}, \cdots, x^{(N)}\right)^{\mathrm{T}}$$
$$\mathbf{t} = \left(t^{(1)}, t^{(2)}, \cdots, t^{(N)}\right)^{\mathrm{T}}$$

$$E(w) = -\ln p(\mathbf{t}|w)$$

Data: generated "<u>i.i.d.</u>"
(see Lecture 2, slide 27)

Take neg. log as usual
$$\widehat{w} = \underset{w}{\mathrm{argmin}}\, E(w)$$

Technische
Universität
Braunschweig

IFL

# Neural Networks (for regression – e.g. curve fitting)

**The regression model**

$$\mu(\boldsymbol{x}) = y(\boldsymbol{x}, \boldsymbol{w}) = f\left(\boldsymbol{w}^{\mathrm{T}}\boldsymbol{\phi}(\boldsymbol{x})\right)$$

**The probabilistic model**

> **(Isotropic) Gaussian distribution**
> (used often since previous lectures)

$$p(t|\boldsymbol{x}, \mu, \sigma) = \mathcal{N}(t|\mu(\boldsymbol{x}), \sigma^2)$$

because <u>the output is continuous</u>
(see Lecture 3, slide 18).

The likelihood function

$$p(\mathbf{t}|\mathbf{X}, \boldsymbol{w}, \sigma) = \prod_{n=1}^{N} \mathcal{N}\left(t^{(n)}|y(\boldsymbol{x}^{(n)}, \boldsymbol{w}), \sigma^2\right)$$

Data
$$\mathbf{X} = \left(\boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}, \cdots, \boldsymbol{x}^{(N)}\right)^{\mathrm{T}}$$
$$\mathbf{t} = \left(t^{(1)}, t^{(2)}, \cdots, t^{(N)}\right)^{\mathrm{T}}$$

$$\Rightarrow \quad E(\boldsymbol{w}) = -\ln p(\mathbf{t}|\boldsymbol{w}) \quad \textbf{Least squares}$$

Data: generated "<u>i.i.d.</u>"
(see Lecture 2, slide 27)

Take neg. log as usual
$$\widehat{\boldsymbol{w}} = \underset{\boldsymbol{w}}{\mathrm{argmin}}\, E(\boldsymbol{w})$$

# Neural Networks (for regression: multiple output)

## The regression model

$$\boldsymbol{\mu}(\boldsymbol{x}) = \mathbf{y}(\boldsymbol{x}, \mathbf{W}) = f.\left(\mathbf{W}^{\mathrm{T}}\boldsymbol{\phi}(\boldsymbol{x})\right)$$

$f. : f$ is applied to each component

## The probabilistic model

**(Isotropic) Gaussian distribution**
(used often since previous lectures)

$$p(\boldsymbol{t}|\boldsymbol{x}, \mu, \sigma) = \mathcal{N}(\boldsymbol{t}|\boldsymbol{\mu}(\boldsymbol{x}), \sigma^2 \mathbf{I})$$

because the output is continuous
(see Lecture 3, slide 18).

The likelihood function

$$p(\mathbf{T}|\mathbf{X}, \mathbf{W}, \sigma) = \prod_{n=1}^{N} \mathcal{N}\left(\boldsymbol{t}^{(n)}|\mathbf{y}(\boldsymbol{x}^{(n)}, \mathbf{W}), \sigma^2 \mathbf{I}\right)$$

Data
$$\mathbf{X} = \left(\boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}, \cdots, \boldsymbol{x}^{(N)}\right)^{\mathrm{T}}$$
$$\mathbf{T} = \left(\boldsymbol{t}^{(1)}, \boldsymbol{t}^{(2)}, \cdots, \boldsymbol{t}^{(N)}\right)^{\mathrm{T}}$$

$$E(\mathbf{W}) = -\ln p(\mathbf{T}|\mathbf{W})$$ **Least squares**

Take neg. log as usual    $\widehat{\mathbf{W}}$ is obtained.

Data: generated "i.i.d."
(see Lecture 2, slide 27)

$\mathbf{W}$ can be expressed by a vector $\boldsymbol{w}$ not by a matrix.

IFL

# Lecture content

- To neural networks (NNs)
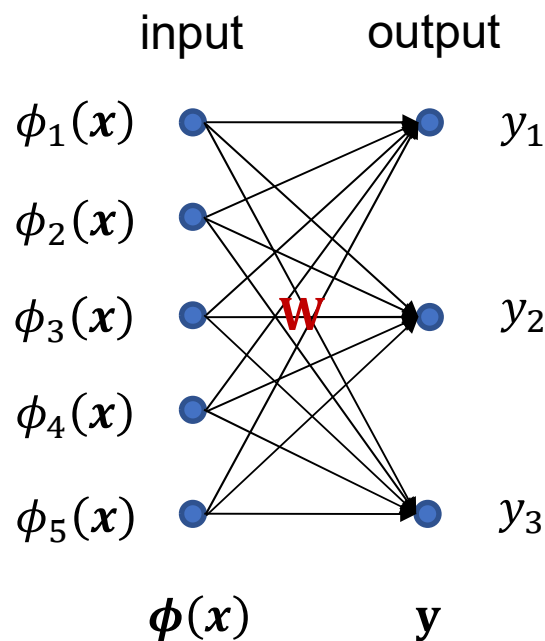
GLM → NNs

mind

Technische
Universität
Braunschweig

IFL

# Linear Regression (multiple output) - Repost

$$\mathbf{y}(x, \mathbf{W}) = \mathbf{W}^{\mathrm{T}}\boldsymbol{\phi}(x)$$

The tuned parameter $\hat{w}_i$ ($i$=1,2,3) can be obtained at the same time as $\hat{\mathbf{W}}$.

input      output

$\phi_1(x)$     $y_1$

$\phi_2(x)$

$\phi_3(x)$    **W**    $y_2$

$\phi_4(x)$

$\phi_5(x)$     $y_3$

$\boldsymbol{\phi}(x)$      **y**

$$\hat{\mathbf{W}} = \left(\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\Phi}\right)^{-1}\boldsymbol{\Phi}^{\mathrm{T}}\mathbf{T}$$

$$= \begin{pmatrix} & \cdots & \\ \vdots & \ddots & \vdots \\ & \cdots & \end{pmatrix} \begin{pmatrix} t_1^{(1)} & \cdots & t_3^{(1)} \\ \vdots & \ddots & \vdots \\ t_1^{(N)} & \cdots & t_3^{(N)} \end{pmatrix} = \begin{pmatrix} \hat{w}_{11} & \cdots & \hat{w}_{13} \\ \vdots & \ddots & \vdots \\ \hat{w}_{51} & \cdots & \hat{w}_{53} \end{pmatrix}$$

$$\hat{\mathbf{W}} = (\hat{w}_1, \hat{w}_2, \hat{w}_3)$$

Consider a multiple output case:

$$\mathbf{T} = \left(t^{(1)}, t^{(2)}, \cdots, t^{(N)}\right)^{\mathrm{T}}$$

This is still the **linear regression model**.

Looks like a neural network?

The objective is to obtain $\hat{\mathbf{W}}$ or $p(\mathbf{W})$ (by using data).

Technische Universität Braunschweig

IFL

# From Linear Regression to Neural Networks

A Linear regression (in general)

$$\mathbf{y}(x, \mathbf{W}) = \mathbf{W}^{\mathrm{T}}\boldsymbol{\phi}(x)$$

Simple extension of the linear regression

Then, nonlinear transformation

$$\mathbf{y}(x, \mathbf{W}) = f.\left(\mathbf{W}^{\mathrm{T}}\boldsymbol{\phi}(x)\right)$$

$f.$ :
a common nonlinear function $f$
is applied to each component

$f$ is a function of a **scalar** input.

e.g. $y_n(x, w_n) = f\left(w_n^{\mathrm{T}}\boldsymbol{\phi}(x)\right)$

$n = 1,2,3$

$$\mathbf{W} = (w_1, w_2, w_3)$$

input     output

$\phi_1(x)$   $f$   $y_1$

$\phi_2(x)$

$\phi_3(x)$   $\mathbf{W}$ $f$   $y_2$

$\phi_4(x)$

$\phi_5(x)$   $f$   $y_3$

$\boldsymbol{\phi}(x)$     $\mathbf{y}$

# The nonlinear function $f$
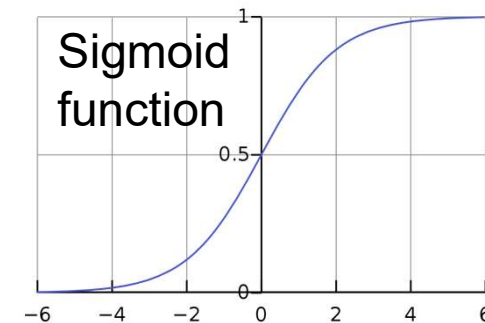
If the application is **regression** (the outputs are <u>continuous</u> values)

- Identity (no mapping)
- ReLU
- Sigmoid function
- …

ReLU

$R(z) = max(0, \ z)$

If the application is **classification** (the outputs are <u>discrete</u> values)

- Sigmoid function (for 2 classes)
- Softmax function (for multi classes)

Sigmoid function

# Important Properties in GLM in General

$$y(x, w) = f\left(w^{\mathrm{T}}\phi(x)\right)$$

**disadvantage**

- $\widehat{w}$ cannot be analytically obtained anymore.          $\widehat{w} = \underset{w}{\operatorname{argmin}}\, E(w)$

   Optimizer is required.          by optimizer

- Neither the posterior distributions $p(w)$ nor predictive distributions can be analytically obtained $p(t)$.

   No spots of analytical solutions in the Bayesian approach. The Bayesian approach is very challenging.
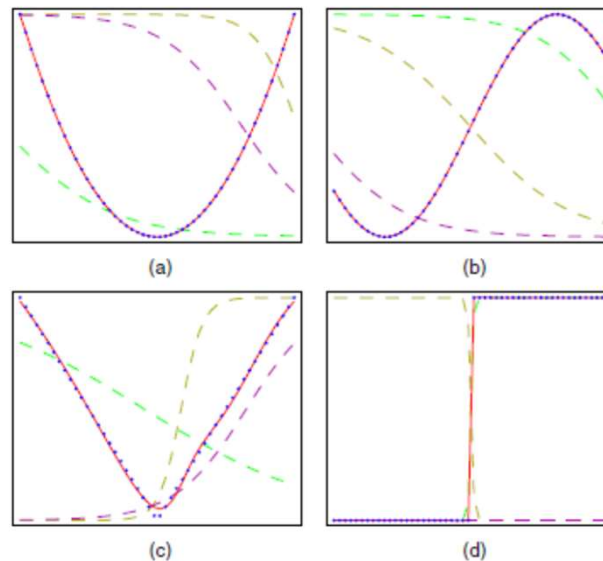
# Important Properties in GLM in General

$$\mu(x) = y(x, w) = f\left(w^{\mathrm{T}}\phi(x)\right)$$

**advantage**

- A variety of efficient expression of the function (the regression model)
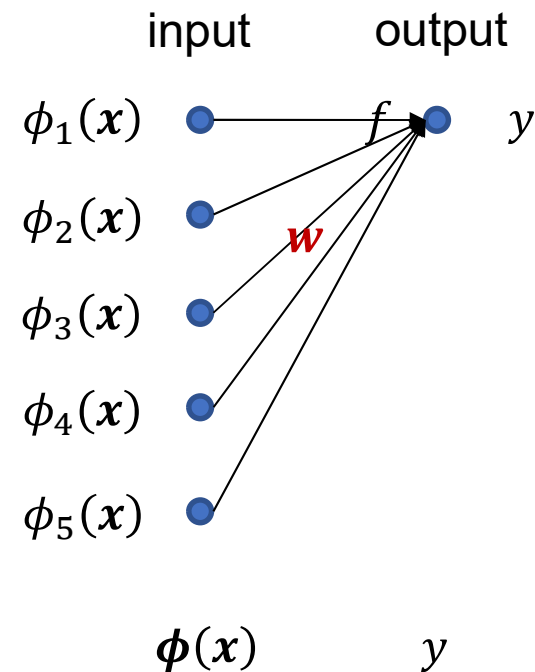
One of the strongest points in neural networks



(a)　(b)

(c)　(d)

A variety of functions

PRML, Fig. 5.3

Technische
Universität
Braunschweig

IFL

# Important Properties in GLM in General

Once we lost the nice properties (analytical solutions), we can focus on the trade-off between the complexity of the function (the regression model) and efficiency (computational time in the learning process)

$$y(\boldsymbol{x}, \boldsymbol{w}) = f.\left(\boldsymbol{w}^{\mathrm{T}}\boldsymbol{\phi}(\boldsymbol{x})\right)$$

input      output

$\phi_1(\boldsymbol{x})$      $f$   $y$

$\phi_2(\boldsymbol{x})$

$\boldsymbol{w}$

$\phi_3(\boldsymbol{x})$

$\phi_4(\boldsymbol{x})$

$\phi_5(\boldsymbol{x})$

$\boldsymbol{\phi}(\boldsymbol{x})$      $y$

Technische Universität Braunschweig

IFL

# Important Properties in GLM in General

Once we lost the nice properties (analytical solutions), we can focus on the trade-off between the complexity of the function (the regression model) and efficiency (computational time in the learning process)
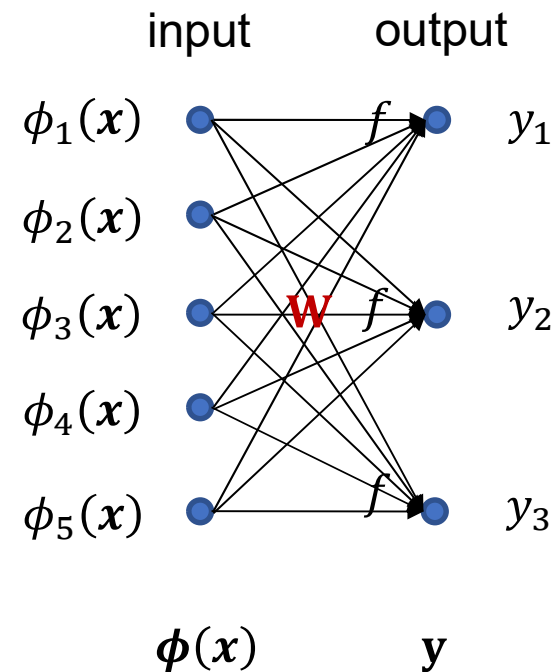
$$\mathbf{y}(\mathbf{x}, \mathbf{W}) = f \cdot \big(\mathbf{W}\boldsymbol{\phi}(\mathbf{x})\big)$$

Technische
Universität
Braunschweig

IFL

# Important Properties in GLM in General

Once we lost the nice properties (analytical solutions), we can focus on the trade-off between the complexity of the function (the regression model) and efficiency (computational time in the learning process)
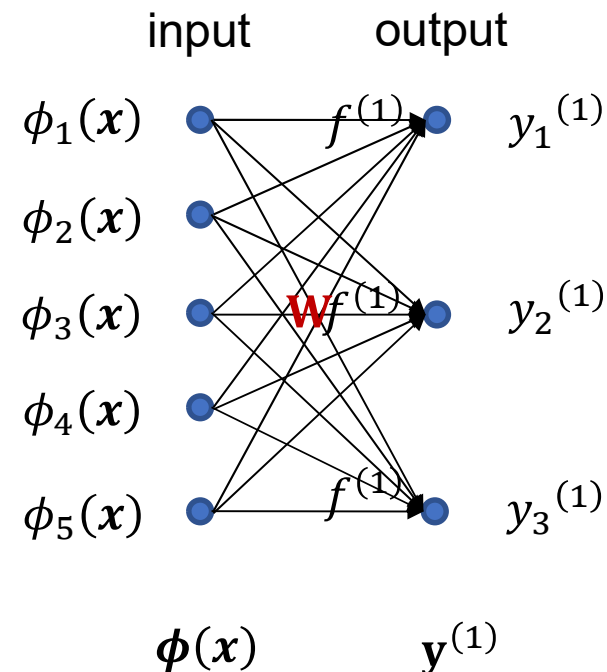
$$\mathbf{y}^{(1)}\left(x, \mathbf{W}^{(1)}\right) = f^{(1)} \cdot \left(\mathbf{W}^{(1)} \boldsymbol{\phi}(x)\right)$$

input     output

$\phi_1(x)$    $f^{(1)}$    $y_1{}^{(1)}$

$\phi_2(x)$

$\phi_3(x)$    $\mathbf{W} f^{(1)}$    $y_2{}^{(1)}$

$\phi_4(x)$

$\phi_5(x)$    $f^{(1)}$    $y_3{}^{(1)}$

$\boldsymbol{\phi}(x)$      $\mathbf{y}^{(1)}$

# Important Properties in GLM in General

Once we lost the nice properties (analytical solutions), we can focus on the trade-off between the complexity of the function (the regression model) and efficiency (computational time in the learning process)
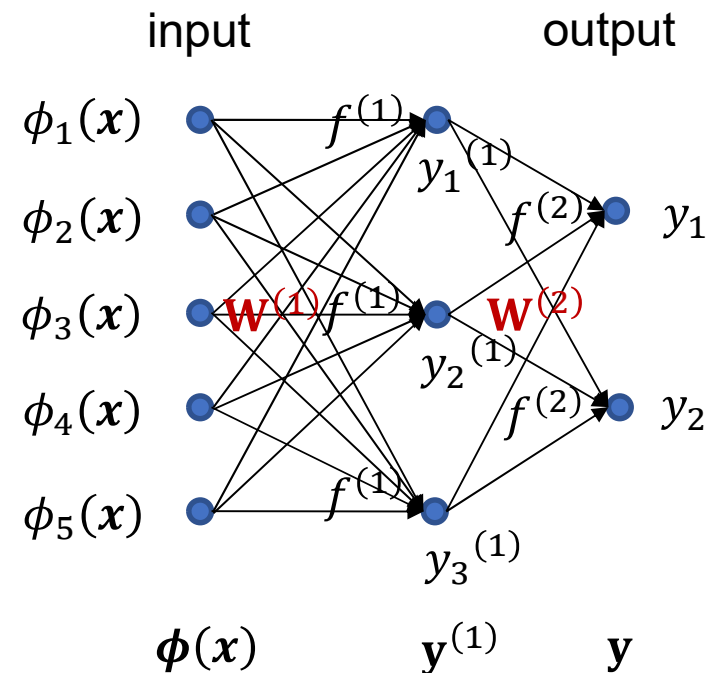
$$\mathbf{y}^{(1)}\big(x, \mathbf{W}^{(1)}\big) = f^{(1)} \cdot \big(\mathbf{W}^{(1)}\boldsymbol{\phi}(x)\big)$$

$$\mathbf{y}(x, \mathbf{W}) = \mathbf{W}^{(2)} f^{(1)} \cdot \big(\mathbf{W}^{(1)}\boldsymbol{\phi}(x)\big)$$

$$\underbrace{\phantom{\mathbf{W}^{(2)} f^{(1)} \cdot \big(\mathbf{W}^{(1)}\boldsymbol{\phi}(x)\big)}}_{\mathbf{y}^{(1)}(x, \mathbf{W})}$$

$w$ is composed of $\mathbf{W}^{(1)}, \mathbf{W}^{(2)}$.

$$\mathbf{y}(x, w) = f^{(2)} \cdot \Big(\mathbf{W}^{(2)} f^{(1)} \cdot \big(\mathbf{W}^{(1)}\boldsymbol{\phi}(x)\big)\Big)$$
$$= f^{(2)} \cdot \big(\mathbf{W}^{(2)}\mathbf{y}^{(1)}\big)$$

$$\big(\mathbf{y}^{(1)} = \mathbf{y}^{(1)}(x, \mathbf{W}^{(1)})\big)$$

input        output

$\phi_1(x)$   $f^{(1)}$
$\phi_2(x)$   $y_1^{(1)}$   $f^{(2)}$   $y_1$
$\phi_3(x)$   $\mathbf{W}^{(1)} f^{(1)}$   $\mathbf{W}^{(2)}$
$\phi_4(x)$   $y_2^{(1)}$   $f^{(2)}$   $y_2$
$\phi_5(x)$   $f^{(1)}$

$y_3^{(1)}$

$\boldsymbol{\phi}(x)$     $\mathbf{y}^{(1)}$     $\mathbf{y}$

Technische Universität Braunschweig

IFL

# Important Properties in GLM in General



e.g. one layer with three units

hidden layers

input

output

$\phi_1(x)$
$\phi_2(x)$
$\phi_3(x)$
$\phi_4(x)$
$\phi_5(x)$

$f^{(1)}$
$y_1^{(1)}$
$\mathbf{W}^{(1)} f^{(1)}$
$y_2^{(1)}$
$f^{(1)}$
$y_3^{(1)}$

...

$y_1^{(n)}$
$f^{(n)}$
$\mathbf{W}^{(n)}$
$y_2^{(n)}$
$f^{(n)}$
$y_3^{(n)}$

$f^{(n)}$
$y_1$
$f^{(n)}$
$y_2$

$\phi(x)$
$\mathbf{y}^{(1)}$
...
$\mathbf{y}^{(n)}$
$\mathbf{y}$

If your objective is **classification**, you can set $f^{(n)}$ to the sigmoid/softmax function.

**Regression model**
$\mathbf{y}(x, \mathbf{w})$

If many layers, the model is called **deep** neural network.

**The dimensionality** of $\mathbf{w}$, which is composed of $\mathbf{W}^{(1)}$, $\cdots$, $\mathbf{W}^{(n)}$, tends to be **large**.

Technische Universität Braunschweig

IFL

# Example: Gaussian Processes (GPs) for Classification

**The regression model**

- $\mu(\boldsymbol{x}) = y(\boldsymbol{x}, \boldsymbol{w}) = sigmoid\left(\boldsymbol{w}^{\mathrm{T}}\boldsymbol{\phi}(\boldsymbol{x})\right)$
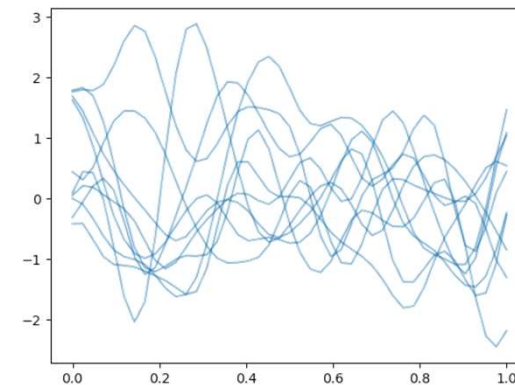
- $\mu(\boldsymbol{x}) = sigmoid(\mathbf{GP})$

**The probabilistic model**

$$p(t|\mu) = \mathrm{Bern}(t|\mu) = \mu^t(1-\mu)^{1-t}$$

The same process then… (showing the process skipped)

**GP**  $\qquad\qquad$ $sigmoid(\mathbf{GP})$

PRML, Fig. 6.11

Technische
Universität
Braunschweig

IFL

# Strength of Neural Networks

- Why not using complex linear regression models?
- Why not using other nonlinear regression models?
- Why not using Gaussian processes?

Prerequieste:

$\phi$ in linear regression models can be nonlinear functions but need to be determined in advance.

$\phi: x \rightarrow s$

$f(s_1, \cdots, s_M)$

$$f(w_1 s_1 + \cdots + w_M s_M) = f(\mathbf{w}^T \mathbf{s}) \qquad \mathbf{w}: M \text{ dim.}$$

$$w_1 s_1^2 + w_2 s_1 s_2 + w_3 s_2^2 + \cdots \qquad \mathbf{w}: M^2/2 \text{ dim.}$$

> The increase of the number of the parameters from exponential to linear by using the ideas of layers.

Technische Universität Braunschweig

IFL

# Issues in Neural Networks

The difference between **NNs** and **Linear regression models**:

**Optimization** or **not** (in the **learning process**)

Note:
We saw that optimization was used in Gaussian processes (GPs) to determine the hyperprameters $\theta$ by MLE. The problems in the optimization are different between GPs (see Lecture 9) and NNs.

The error function to be minimize
(when the probabilistic model is the Gaussian)

$$E(\boldsymbol{w}) = \sum_{n=1}^{N} \{\boldsymbol{t}^{(n)} - \mathbf{y}(\boldsymbol{x}^{(n)}, \boldsymbol{w})\}^2$$

$N$: a huge number    millions
(**big data**)

$$\hat{w} = \underset{\boldsymbol{w}}{\operatorname{argmin}} E(\boldsymbol{w})$$

The regularization can be naturally used

$$E_{reg}(\boldsymbol{w}) = E(\boldsymbol{w}) + \lambda\|\boldsymbol{w}\|^q$$

See Lecture 4, slide 35

Especially when **LASSO as sparse model** is very useful to decrease unnecessary parameters from the plenty of parameters in $\boldsymbol{w}$.

Technische
Universität
Braunschweig

IFL

**How efficiently the learning process (the optimization process) is done**

- Optimization techniques
- The properties of the NN model

Human's learning/experimental process from the data

Then, the goal:
As far as the prediction is good, the model and the chosen techniques are ok.

Details of all the theories behind are not clarified yet.

But now you know how to do it (**Bayesian neural network**)!

**The Bayesian approach** is not used (too expensive),
usually we split data into training data and validation data.

Simply the posterior and the predictive distribution cannot be obtained analytically.

(Lecture 12:
numerical approaches)

(See Lecture 2, slide 7)

Technische
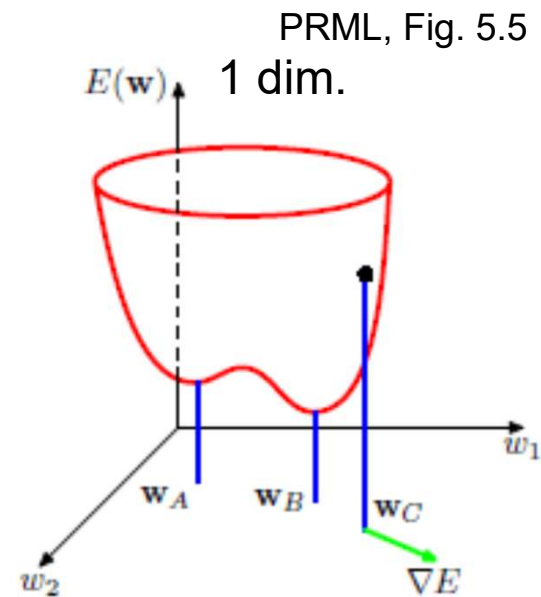Universität
Braunschweig

IFL

# Lecture content

- Technical Issues in neural networks

IFL

# Technical Issues in Optimization

We need to search for the global minimum of the function $E(\mathbf{w})$ in a very high-dimensional space $\mathbf{w}$.

PRML, Fig. 5.5



$$\mathbf{w} = (w_1, w_2, \cdots, w_{10000} \;)^{\mathrm{T}}$$

millions of dim.

**Topics of optimization** (using some nice properties of the problem setting in NNs):

In such a high-dimensional space,
- **Gradient-based optimization**
    - Techniques to efficiently compute the gradients (**chain rule**) – called <u>back propagation</u>
    - Improved gradient opt. algorithms
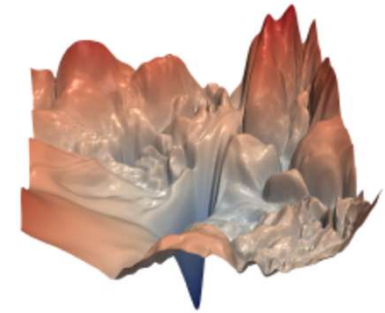        - AdaGrad
        - Adam
        - …many others

Most of open-source libraries in Python nowadays have these functionalities as standard.

Technische Universität Braunschweig

IFL

# Technical Issues in Optimization

To avoid large computational costs by using all data, a subset of the data is used to compute the gradient at each iteration of the optimization process.

- **Stochastic** gradient descent (SGD) algorithm
  - Drop-out
  - Over-parametrization
  - Parallelization (power of GPU)
  - …many others



Li, H., Xu, Z., Taylor, G., Studer C. amd Goldstein, T., "visualizing the loss landscape of neural nets", 2018.

As a result, the algorithm also contributes to avoid local minimum!

- Convolutional neural network (CNN) – parameters $w$ can be drastically reduced.
- …many others

A variety of methods has been developed for these several years!

Technische Universität Braunschweig
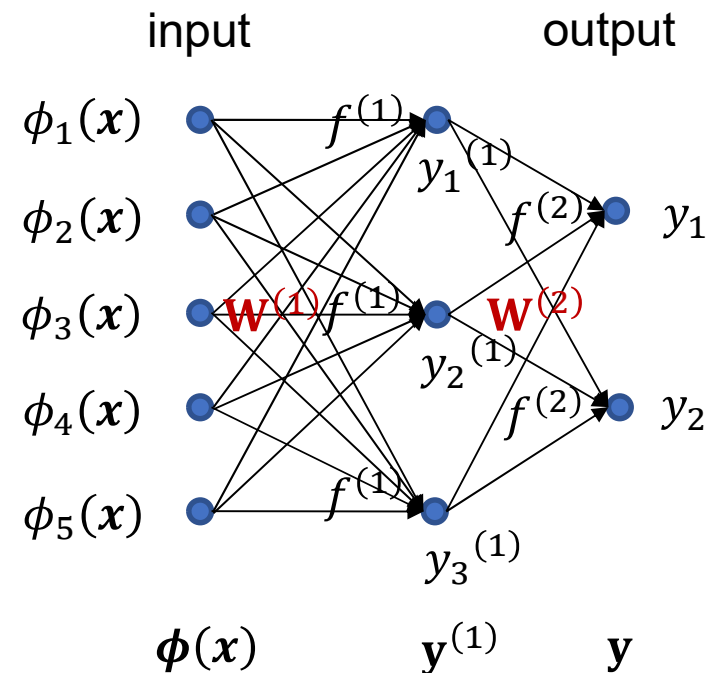
IFL

# Recent Advanced Methods

- AlexNet (2012)
  - 8 layers
  - 60 million parameters


- VGG19 Net (2014)
  - 46 layers
  - more than 100 million parameters


- Residual Network (ResNet) (2015)
  - 152 layers
  - 20 million parameters (fewer parameters per layer)

# Hyperparameter Tuning

- How many layers do we set?
- How many parameters (the number of units per layer) do we set?
- ...

input                      output

Hyperparameter tuning

By changing these parameters, we try to find
a combination of the parameters with the best
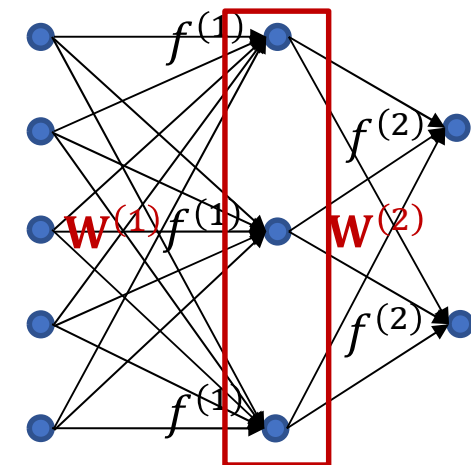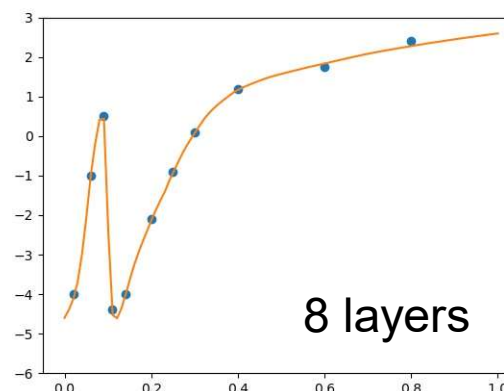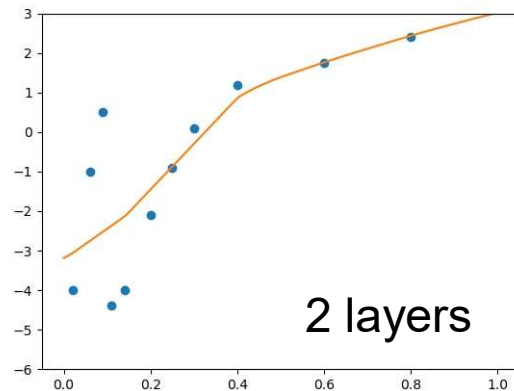score (good prediction of a validation data).

$\phi_1(x)$    $f^{(1)}$

$y_1^{(1)}$   $f^{(2)}$   $y_1$

$\phi_2(x)$

$\phi_3(x)$   $W^{(1)}$ $f^{(1)}$   $W^{(2)}$

$y_2^{(1)}$

$\phi_4(x)$    $f^{(2)}$   $y_2$

$\phi_5(x)$    $f^{(1)}$

$y_3^{(1)}$

$\boldsymbol{\phi}(x)$     $\mathbf{y}^{(1)}$    $\mathbf{y}$

Technische
Universität
Braunschweig

IFL

- Any <u>smooth functions</u> can be represented by 2 layers NNs
  (kernel methods such as GPs as well).
  - No other models can improve AER.

- When the target function is a <u>step</u> / <u>non-uniform</u> function, NNs with 4 layers
  or more achieve the best AER.
  - No other models can exceed the rate.

Example: A non-uniform function



2 layers



8 layers



$f^{(1)}$

$f^{(2)}$

$\mathbf{W}^{(1)} f^{(1)}$  $\mathbf{W}^{(2)}$

$f^{(2)}$

$f^{(1)}$

Note:
When # of the units → ∞,
the model → a GP.

Technische
Universität
Braunschweig

IFL

# Perspective of Neural Networks

**Mystery of Hundred millions of parameters**

training data size

When the number of the parameters $w > N$ ⟹ overfitting

See Lecture 2, slide 6
(Fig. 1.8 in PRML)

More parameters can reduce the error.

Avoiding the overfitting



Neyshabur, B., Tomioka, R., Salakhutdinov
, R. and Srebro, N., "Geometry of Optimization and
Implicit Regularization in Deep Learning", 2017.

Belkin, M., Hsu, D., Ma, S. and Mandal, S., "reconciling modern
machine learning and the bias-variance trade-off", 2019.

# Brief Summary

|  | *Point estimate* | *Bayesian* |
|---|---|---|
| Linear regression | Analytical solutions | Some useful properties available (Gaussian processes) |
| Nonlinear regression GLM (Classification, NNs) | **A variety of functions** | hard |

focusing on this strength

Technische Universität Braunschweig

IFL

# Next Step

Neural networks (not only them but supervised learning techniques in general) are often used with unsupervised learning techniques.

Without these, the state of the art techniques are not discussed.

➡️ Unsupervised learning methods in the next lecture