# Inferring Missing Data with Auto-Associators

Mark J. Embrechts (Mark.Embrechts@gmail.com)

**HMGU Munich**

HELMHOLTZ
MUNICH

**Andreas Vesalius**

| | |
|---|---|
| **Born** | 31 December 1514<br>Brussels, Habsburg Netherlands |
| **Died** | 15 October 1564 (aged 49)<br>Zakynthos, Republic of Venice |
| **Fields** | Anatomy |
| **Doctoral adviser** | Johannes Winter von Andernach<br>Gemma Frisius |
| **Doctoral students** | Matteo Realdo Colombo |
| **Known for** | *De humani corporis fabrica* or "the fabric of the human body" |
| **Influences** | Jacques Dubois<br>Jean Fernel |

Vesalius, 1543

Vesalius, Anatomy, 1543

optic radiations

Lateral geniculate nucleus (in thalamus)

Visual cortex

Optic nerve
Crossed fibres
Uncrossed fibres
Optic chiasma

Optic tract
Commissure of Gudden

Pulvinar
Lateral geniculate body
Superior colliculus
Medial geniculate body

Nucleus of oculomotor nerve

Nucleus of trochlear nerve

Nucleus of abducent nerve

Cortex of occipital lobes

Vesalius, Anatomy, 1543

- Dealing with missing data
- Auto-associators
- Some results



**HELMHOLTZ**
**MUNICH**

Missing Data

- What are missing data?
- Simple ways to deal with missing data
  - Eliminate descriptor (columns)
  - Eliminate records (rows
  - Replace by mean/median values
- More sophisticated ways to deal with missing data
  - Auto-associators, ...
- Remarks
  - There might be information in missing data
  - Categorical data need a more sophisticated approach

HELMHOLTZ
MUNICH

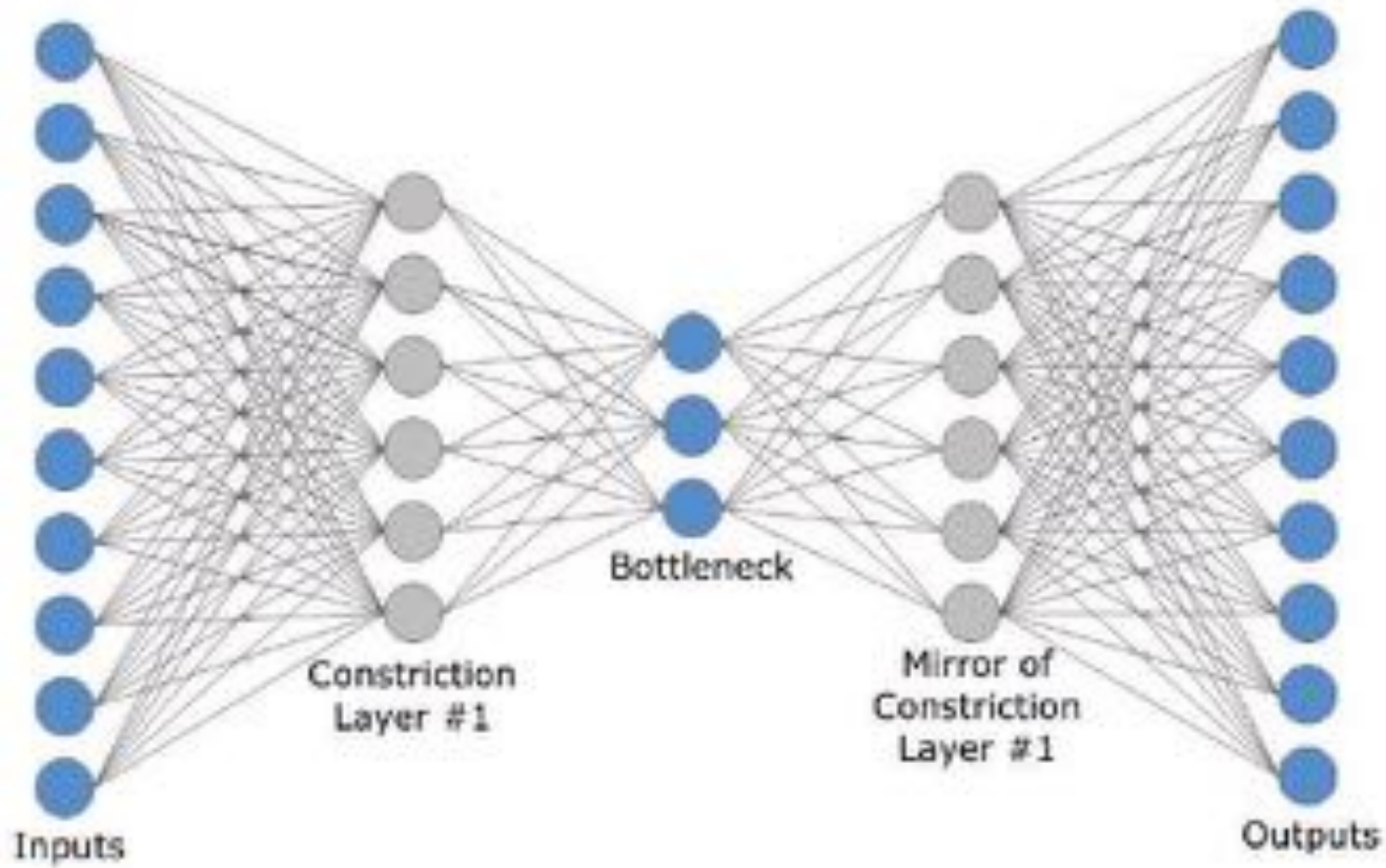| Some tricky issues related to missing data |
| --- |

- Missing data can by systemic:
  e.g., in a survey a person might not want to mention alcohol use.
- It might be helpful to add a column indicating whether a descriptor is missing
- Data curation: suspect data might be flagged as missing (-999)
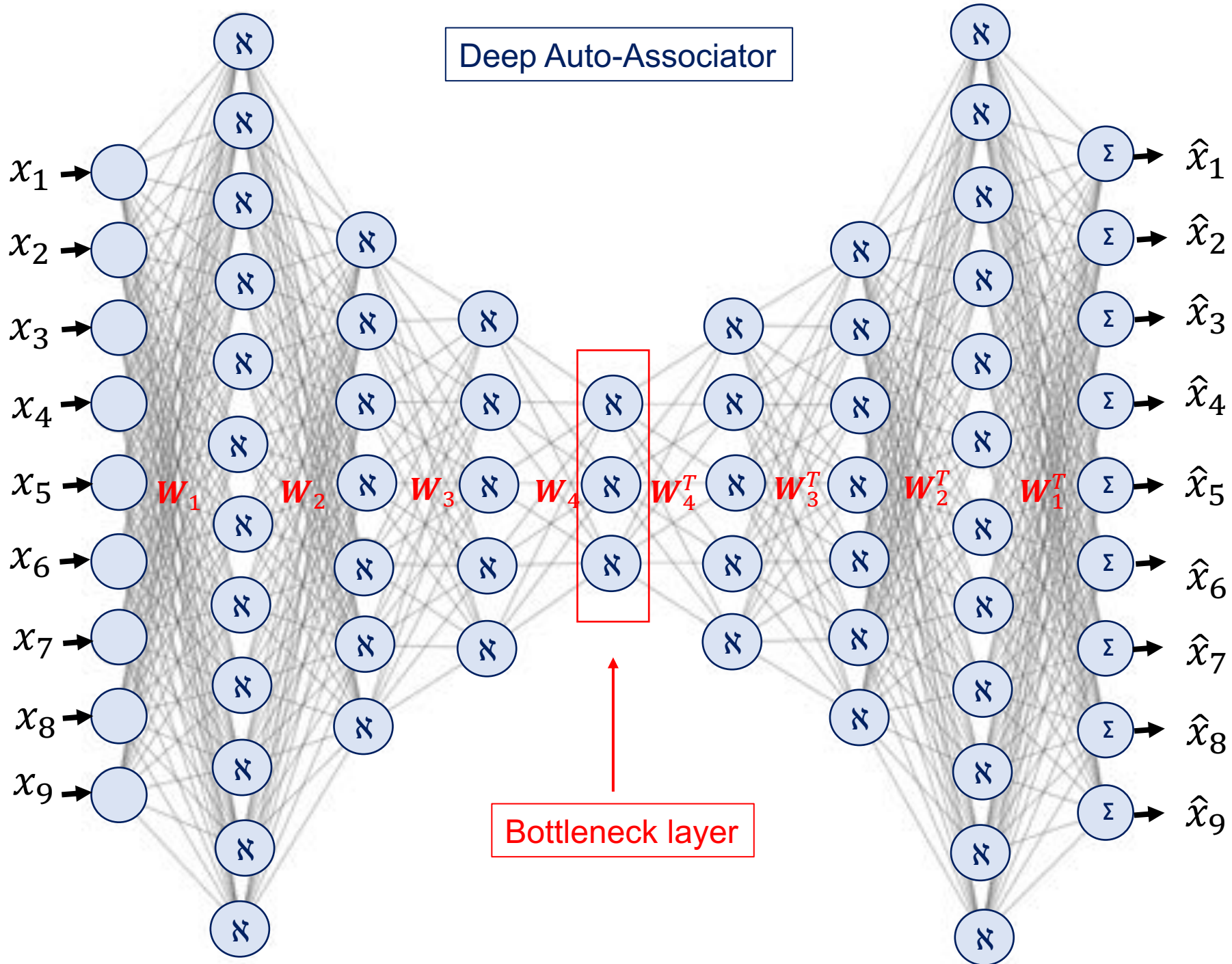- Example of systemic issues with missing data: 9 class Italian Olive Oil Data

| | | | | | | | | class | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| 32.9 | 14.3 | 45.7 | 67.3 | 35.9 | 58.3 | 65 | 80.4 | 4 | 318 |
| 31.7 | 15.1 | 54.3 | 66 | 37.7 | 59.7 | 59.2 | 80.4 | 4 | 319 |
| 46 | 21.9 | 90.6 | 58.2 | 27.4 | 58.3 | 71.8 | 64.3 | 4 | 320 |
| 51.4 | 26.4 | 54.3 | 54.3 | 35.8 | 58.3 | 68.9 | 71.4 | 4 | 321 |
| 57.5 | 30.2 | 42.6 | 52.4 | 35.6 | 56.9 | 66 | 58.9 | 4 | 322 |
| 58 | 27.5 | 66.8 | 55.5 | 25 | 56.9 | 67 | 51.8 | 4 | 323 |
| 46.1 | 39.6 | 31.4 | 46.1 | 65 | 56.9 | 93.2 | 0 | 5 | 324 |
| 38.4 | 45.3 | 26 | 51 | 65.4 | 45.8 | 85.4 | 0 | 5 | 325 |
| 43.8 | 30.6 | 26 | 51.2 | 62.3 | 41.7 | 89.3 | 0 | 5 | 326 |
| 45.2 | 30.9 | 30.9 | 46.4 | 69.1 | 45.8 | 89.3 | 0 | 5 | 327 |

- Descriptors for certain classes are missing, but set to zero in the original data set
- We will replace the zero settings by -999, indicating they are really missing

# Auto-Encoders or Auto-Associative Networks



Inputs — Constriction Layer #1 — Bottleneck — Mirror of Constriction Layer #1 — Outputs

HELMHOLTZ MUNICH

Deep Auto-Associator

$x_1$ $x_2$ $x_3$ $x_4$ $x_5$ $x_6$ $x_7$ $x_8$ $x_9$

$W_1$ $W_2$ $W_3$ $W_4$ $W_4^T$ $W_3^T$ $W_2^T$ $W_1^T$

$\hat{x}_1$ $\hat{x}_2$ $\hat{x}_3$ $\hat{x}_4$ $\hat{x}_5$ $\hat{x}_6$ $\hat{x}_7$ $\hat{x}_8$ $\hat{x}_9$

Bottleneck layer

HELMHOLTZ MUNICH

# More Recent History of Neural Networks

# 2015 Adam: Adaptive moment estimation

**PARAMETER SETTINGS:**
η    -- stepsize (i.e., learning parameter η 0.0001)
$β_1$ -- exponential decay rate gradient (0.9)
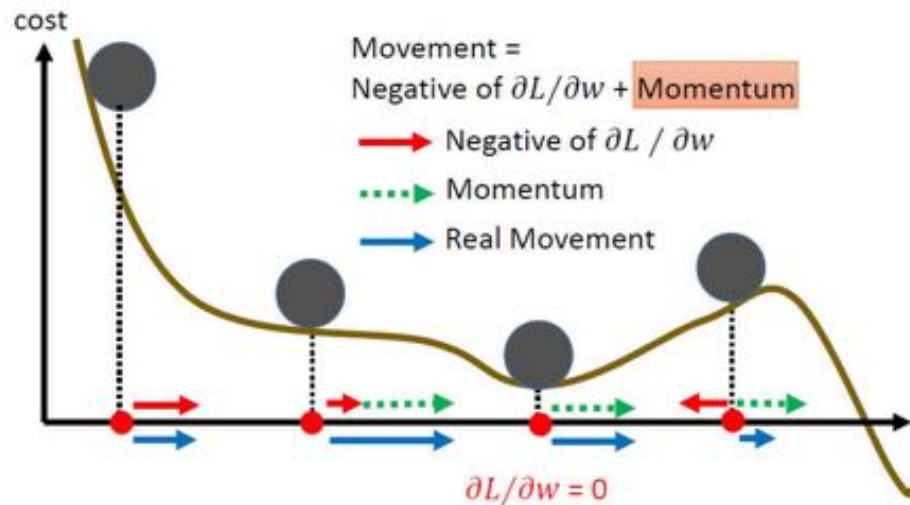$β_2$ -- exponential decay rate 2nd moment (0.999)

**INITIALIZATION:**
$m_0 \leftarrow 0$ *(gradient tensor)*
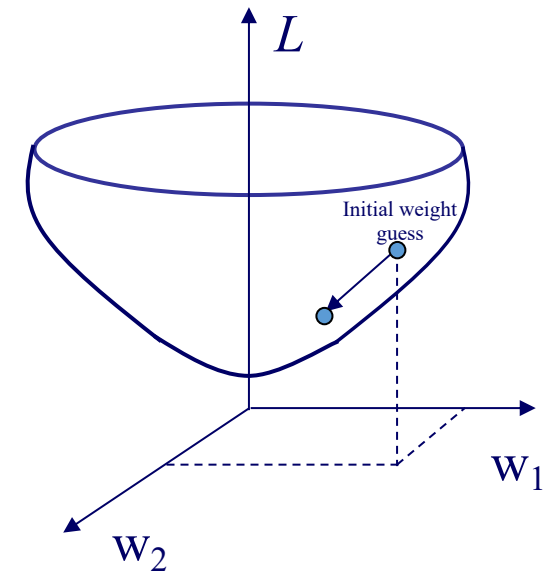$v_0 \leftarrow 0$ *(2nd moment tensor)*
$w_0 \leftarrow 0$ *(weight tensor)*

## Momentum

cost

Movement =
Negative of $\partial L/\partial w$ + Momentum

→ Negative of $\partial L / \partial w$
····▶ Momentum
→ Real Movement

$\partial L/\partial w = 0$

**UPDATING RULE:**

$$g_t \leftarrow \nabla L_t(w_t)$$
$$m_t \leftarrow β_1 m_{t-1} + (1 - β_1) g_t$$
$$v_t \leftarrow β_2 v_{t-1} + (1 - β_2) g_t^2$$
$$\hat{m}_t \leftarrow {m_t}/{(1-β_1^t)}$$
$$\hat{v}_t \leftarrow {v_t}/{(1-β_2^t)}$$
$$w_t \leftarrow w_{t-1} - η \frac{\hat{m}_t}{\sqrt{\hat{v}_t}+ε}$$

$$\Delta_t = \alpha \cdot \hat{m}_t / \sqrt{\hat{v}_t} + \epsilon$$
$$|\Delta_t| \lessgtr \alpha.$$
**Trust-region**
$$\hat{m}_t / \sqrt{\hat{v}_t} \approx E[g_t] / \sqrt{E[g_t^2]} \leq 1$$

$L$

Initial weight guess

$W_1$

$W_2$

Diederick P. Kingman, and Jimmy Lei Ba [2015] Adam: A method for stochastic optimization. Proceedings of the 3rd International Conference for Learning Representations ICLR15, San Diego, 2015.

**HELMHOLTZ MUNICH**

Tricks to make learning faster or more effective: Adam (Adaptive moment estimation)

$$w_t = w_{t-1} - \eta \, {}^{\widehat{m}_t}\!/_{\sqrt{\widehat{v}_t} + \varepsilon}$$

$$\widehat{m}_t = {}^{m_t}\!/_{(1-\beta_1^t)}$$

$$w_t \leftarrow w_{t-1} - \frac{\nabla E}{H}$$

$$w_t \leftarrow w_{t-1} - \eta \, {}^{\widehat{m}_t}\!/_{\sqrt{\widehat{v}_t} + \varepsilon}$$

$$\widehat{v}_t = {}^{v_t}\!/_{(1-\beta_2^t)}$$

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

$\beta_1$ is the exponential decay rate gradient (typically, 0.9) and $\beta_2$ is the exponential decay rate for the 2nd momentum (typically 0.999)

HELMHOLTZ
MUNICH

- I typically use a Mx800x400x200x100x50x10x50x100x200x800xM structure (12 hidden layers)
- Train by policy:
   - Use Adam
   - mini-batches of 30 data and 30 passes through the data
- For missing data, outputs with missing data are not backpropagated
- More details
   - tanh activation function

# Case Study #1:Toxicity challenge data as an example of real-world QSAR data

- 963 training data and 120 test data
- 2223 descriptors
- Consider 6 sigma outliers as missing data

```
REM GET DATA
mje tox --TOX
REM EXTRACT MOE DESCRIPTORS
REM execute tox (255 1)
tox
```

Typical phenotype of a zebrafish embryo incubated from 24-hpf to 96 hpf in (a) embryo medium as a Negative control, in (b) 10 mm diethyl-aminobenzaldehyde (DEAB), and in (c) 100 mm DEAB. Note the Deformed embryos in DEAB: short size, scoliosis, yolk, and heart edema (black arrows).

Younes, Fatima Mraiche, Sahar I. Da'as, and Husevin C. Yalkin [2018] Using zebrafish for investigating the Molecular Mechanisms of Drug-induced cardiotoxicity. Biomedical Research International (BMRI), Vol. 2018, Article ID 1642684.

## 2008 Toxicity challenge data as an example of real-world QSAR data

- ICANN 2008 QSAR Toxicity Data with 2223 descriptors:
    255 MOE descriptors (255 1)
    1664 DRAGON descriptors (1664 256)
    221 SIMULATIOM PLUS descriptors (221 1920)
    60 ELECTRONIC STATE descriptors (60 2141)
    23 QUANTUM CHEMISTRY descriptors (23 2201)
- There are 644+339+110 = 1093 training and 120 test molecules

[1] Hao Zhu, Alexander Tropsha, Denis Fourches, Alexandre Varnek, Ester Papa, Paola Gramatica, Tomas Oberg, Phuong Dao, Artem Cherkasov, and Igor V. Tetko [2008] Combinatorial QSAR Modeling of Chemical Toxicants Tested against Tetrahymena pyriformis. Journal of Chemical Information and Modeling, Vol. 48 pp. 766-784.

[2] Igor V. Tetko, Iurii Sushko, Anil Kumar Pandey, Hao Zhu, Alexander Tropsha, Ester Papa, Tomas Oberg, Robrto Todeschini, Denis Fourches, and Alexandre Varnek [2008] Critical Assessment of QSAR Models of Environmental Toxicity against Tetrahymena pyriformis: Focusing on Applicability Domain of Overfitting by Variable Selection. Journal of Chemical Information and Modeling, Vol. 48, pp. 1733 – 1746.
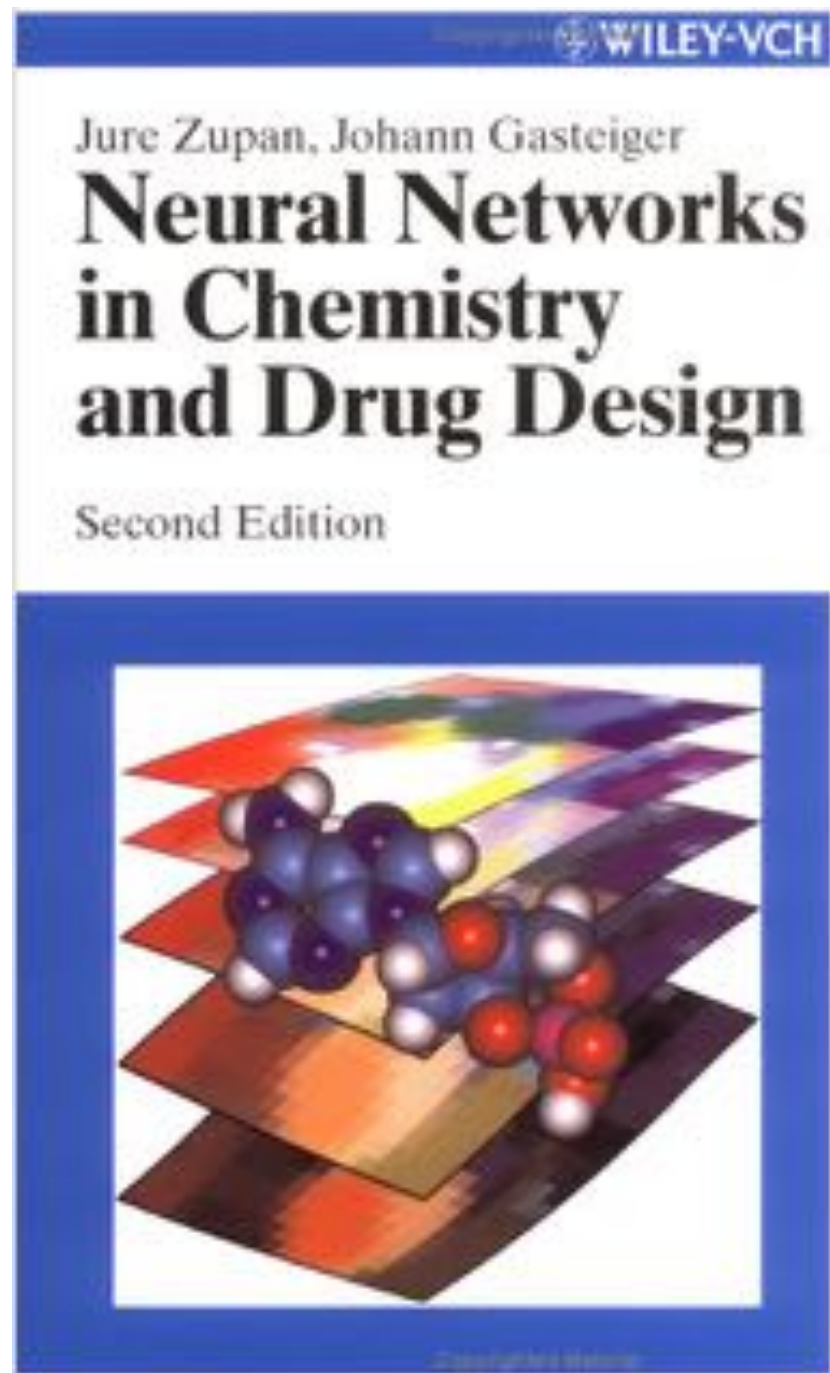
# 2008 Toxicity challenge data as an example of real-world QSAR data



| q2 | Q2 | MSE | MAE | #inputs | net | Note |
|---|---|---|---|---|---|---|
| 0.4885 | 0.5489 | 0.835 | 0.598 | 2223 | 800-400-200-110-50 | |
| 0.4583 | 0.5103 | 0.805 | 0.579 | 629 | same | 95% cousin removed |
| 0.4653 | 0.5300 | 0.820 | 0.610 | 629 | same | cousin + six sigma missing |
| 0.5776 | 0.6613 | 0.916 | 0.675 | 2223 | same | six sigma missing |
| 0.4951 | 0.5452 | 0.832 | 0.593 | 111 | same | + cousins removed |
| 0.4421 | 0.4663 | 0.769 | 0.573 | 2223 | same | six sigma corrected |
| 0.5076 | 0.5681 | 0.849 | 0.607 | 1429 | | + 97% cousin removed |

- Would have ranked 6th in competition
- All training was done by policy

HELMHOLTZ
MUNICH

Jure Zupan, Johann Gasteiger

# Neural Networks in Chemistry and Drug Design

## Second Edition

## Case study #2: Italian Olive Oil Data

| Class | Region | # samples |
|---|---|---|
| 1 | North Apulia | 25 |
| 2 | Calabria | 56 |
| 3 | South Apulia | 206 |
| 4 | Sicily | 36 |
| 5 | Inner Sardinia | 65 |
| 6 | Coastal Sardinia | 33 |
| 7 | East Liguria | 50 |
| 8 | West Liguria | 50 |
| 9 | Umbria | 51 |
| | | 572 |

- 9-class Italian olive oil data for 572 Italian olive oils
  - 8 fatty acid indicators by 9 regions
  - 9-class data are not balanced by region
  - Gasteiger used 250 training data (we do same)

[1] M. Forina and C. Armanino [1981] Eigenvector projection and simplified nonlinear mapping of fatty acid content of Italian olive oils. Ann. Chem, Vol. 72, pp. 125-127.
[2] Jure Zapan and Johann Gasteiger [1999] Neural networks in chemistry and drug design (2nd edition). Wiley - VCH

# Issue related to Italian Olive Oil Data

- Missing data are systemic: some classes have descriptors missing over entire class ➜ gives away answer
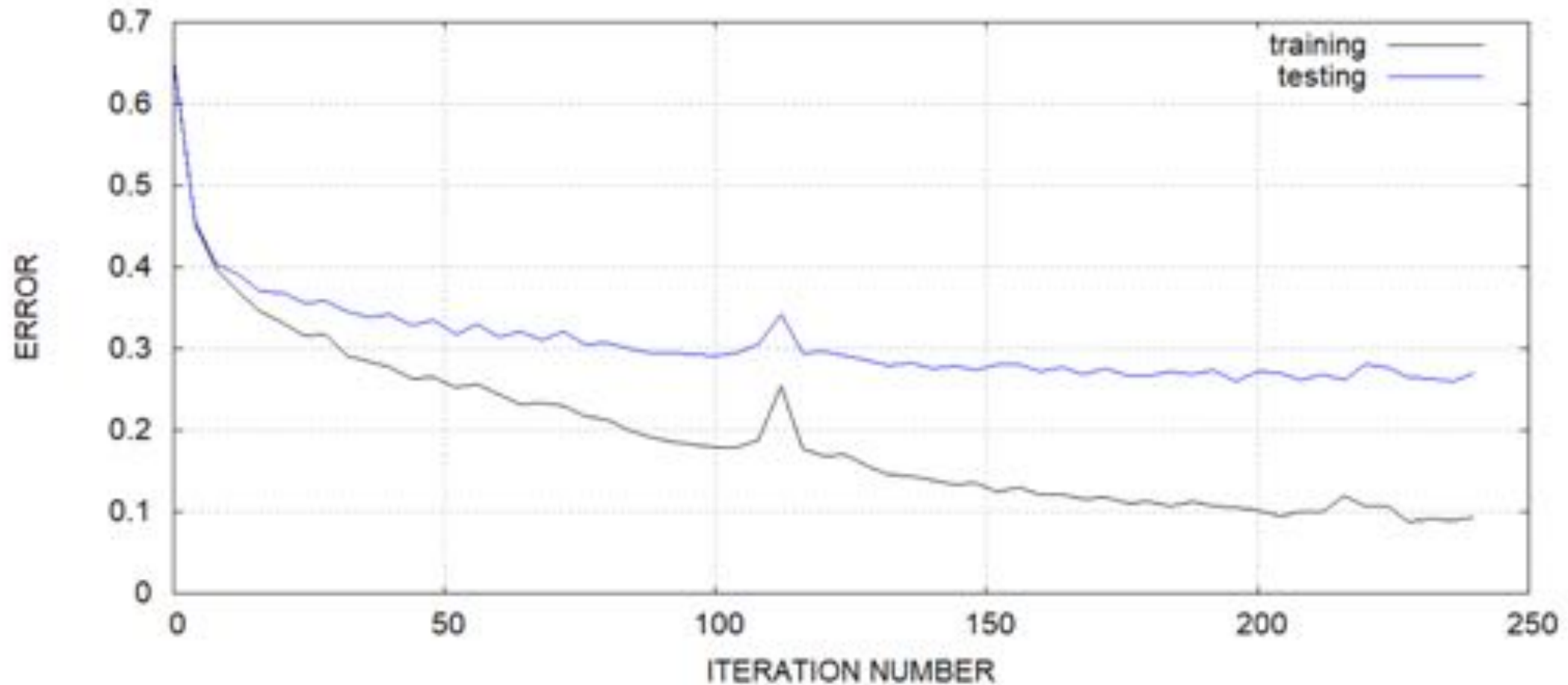- Data curation: flag missing data by -999

| | | | | | | | | class | |
|---|---|---|---|---|---|---|---|---|---|
| 32.9 | 14.3 | 45.7 | 67.3 | 35.9 | 58.3 | 65 | 80.4 | 4 | 318 |
| 31.7 | 15.1 | 54.3 | 66 | 37.7 | 59.7 | 59.2 | 80.4 | 4 | 319 |
| 46 | 21.9 | 90.6 | 58.2 | 27.4 | 58.3 | 71.8 | 64.3 | 4 | 320 |
| 51.4 | 26.4 | 54.3 | 54.3 | 35.8 | 58.3 | 68.9 | 71.4 | 4 | 321 |
| 57.5 | 30.2 | 42.6 | 52.4 | 35.6 | 56.9 | 66 | 58.9 | 4 | 322 |
| 58 | 27.5 | 66.8 | 55.5 | 25 | 56.9 | 67 | 51.8 | 4 | 323 |
| 46.1 | 39.6 | 31.4 | 46.1 | 65 | 56.9 | 93.2 | **0** | 5 | 324 |
| 38.4 | 45.3 | 26 | 51 | 65.4 | 45.8 | 85.4 | **0** | 5 | 325 |
| 43.8 | 30.6 | 26 | 51.2 | 62.3 | 41.7 | 89.3 | **0** | 5 | 326 |
| 45.2 | 30.9 | 30.9 | 46.4 | 69.1 | 45.8 | 89.3 | **0** | 5 | 327 |

- Descriptors for certain classes are missing, but set to zero in the original data set

# Case study #2: Italian Olive Oil Data

| q2 | Q2 | %COR | BER(%) | F1 | Comment |
|---|---|---|---|---|---|
| 0.0269 | 0.0274 | 94.099 | 92.59 | 0.92 | Original Data (8x400x200x110x50x23x9 net) – **systemic issue** |
| 0.1192 | 0.1230 | 88.820 | 85.73 | 0.85 | Missing data replaced by average (8x400x200x110x50x23x9 net) |
| 0.1372 | 0.1440 | 89.441 | 84.82 | 0.84 | Missing data inferred with auto-associator |
| 0.1018 | 0.1050 | 90.062 | 87.15 | 0.86 | same but bottleneck is now 7 rather than 8 neurons |
| 0.0975 | 0.1015 | 90.373 | 86.52 | 0.86 | same but 6 neurons in bottleneck |

## Conclusions

- Missing data are often systemic
- Missing data are usually represented by -999 (as a tag)
- Faulty data (e.g., six-sigma outliers in MOE descriptors can be flagged as missing and then inferred
- Auto-associators are an effective trick to infer missing data
- This approach can also be used on the NETFLIX challenge

HELMHOLTZ
MUNICH