UNIVERSITY OF
CAMBRIDGE

# Advanced Data Science

Lecture 5 : Introduction to Statistical Learning

Carl Henrik Ek - che29@cam.ac.uk

17th of November, 2021

http://carlhenrik.com

**Access** how to get and combine the data-sources for a potential problem

**Access** how to get and combine the data-sources for a potential problem

**Assess** things to do with data before you have a question

**Access** how to get and combine the data-sources for a potential problem

**Assess** things to do with data before you have a question

**Address** what you can do when you have data and a question to answer

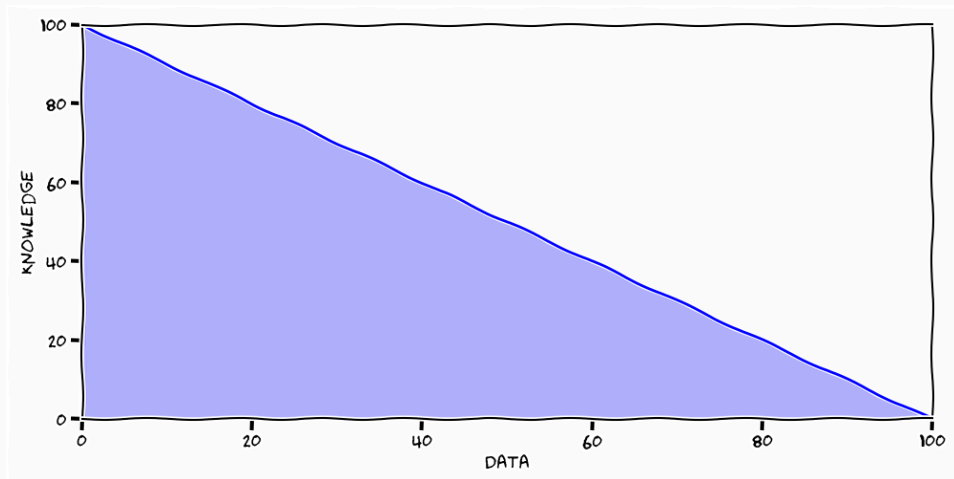| Access | Assess | Address |
|--------|--------|---------|
| • ? | • Introduction to Probability (IA) | • ML and Real-world Data (IA) |
| | • Scientific Computing (IA) | • Data Science (IB) |
| | • Cloud Computing (II) | • AI (IB) |
| | • . . . . | • ML & Bayesian Inference (II) |
| | | • Deep NN (II) |
| | | • Randomised Algorithms (II) |
| | | • . . . . |

*You need to put Machine Learning in the context of data (and humans)*
*– Neil Lawrence*

- Tasks that are too hard to program
  - speech recognition
  - image understanding

- Tasks that are too hard to program
  - speech recognition
  - image understanding
- Tasks beyond our capability
  - weather prediction
  - web search

- Tasks that are too hard to program
  - speech recognition
  - image understanding
- Tasks beyond our capability
  - weather prediction
  - web search
- Machine Learning bridges the knowledge gap by data

$$p(\theta \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \theta)p(\theta)}{p(\mathcal{D})}$$

- Inductive biases comes into the learning procedure

$$p(\theta \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \theta)p(\theta)}{p(\mathcal{D})}$$

- Inductive biases comes into the learning procedure
- *Most knowledge is introduced before we apply ML*

  **access** what data did I acquire?

  **assess** how did I prepare/treat the data?

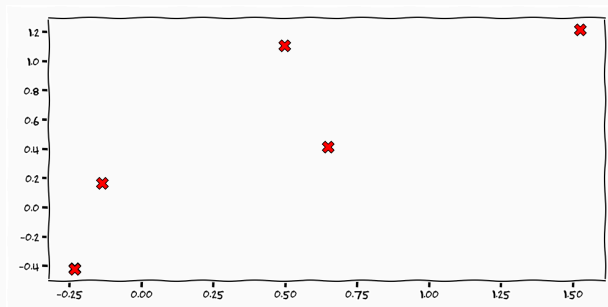$$p(\theta \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \theta)p(\theta)}{p(\mathcal{D})}$$

- Inductive biases comes into the learning procedure
- *Most knowledge is introduced before we apply ML*
    - **access** what data did I acquire?
    - **assess** how did I prepare/treat the data?
- The idea of the $80/20$

- what can machine learning actually do?
- what role does the "data scientist" play in the machine learning loop?
- put machine learning into context

# Statistical Learning

*Learning is the process of converting experience into expertise or knowledge.*
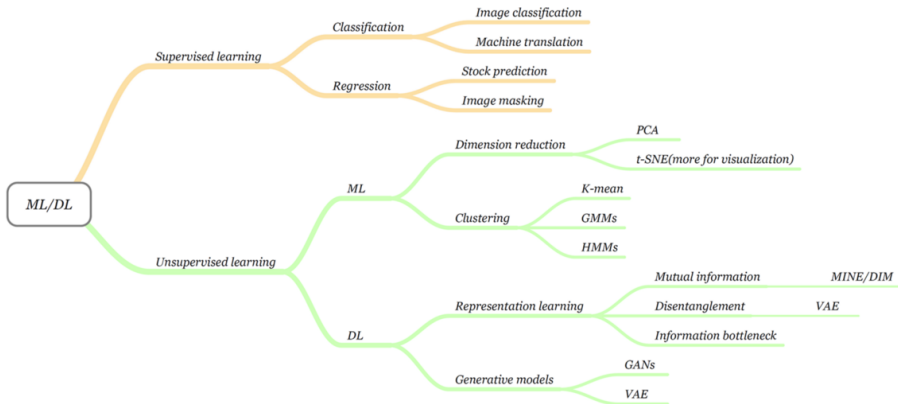*– Sha Ben-David*

**Supervised Learning** $p(y \mid x)$

**"Unsupervised" Learning** $p(y)$

**Reinforcement Learning** $p(\pi, f \mid \mathcal{L})$

**Domain Set** $\mathcal{X}$ the set of measurements/objects that we want to label (input)

**Domain Set** $\mathcal{X}$ the set of measurements/objects that we want to label (input)

**Label Set** $\mathcal{Y}$ the set of outputs

**Domain Set** $\mathcal{X}$ the set of measurements/objects that we want to label (input)

**Label Set** $\mathcal{Y}$ the set of outputs

**Training Data** $\mathcal{S}$ a finite sequence of pairs in $\mathcal{X} \times \mathcal{Y}$

**Data Distribution** $\mathcal{D}$ probability distribution governing the measurements
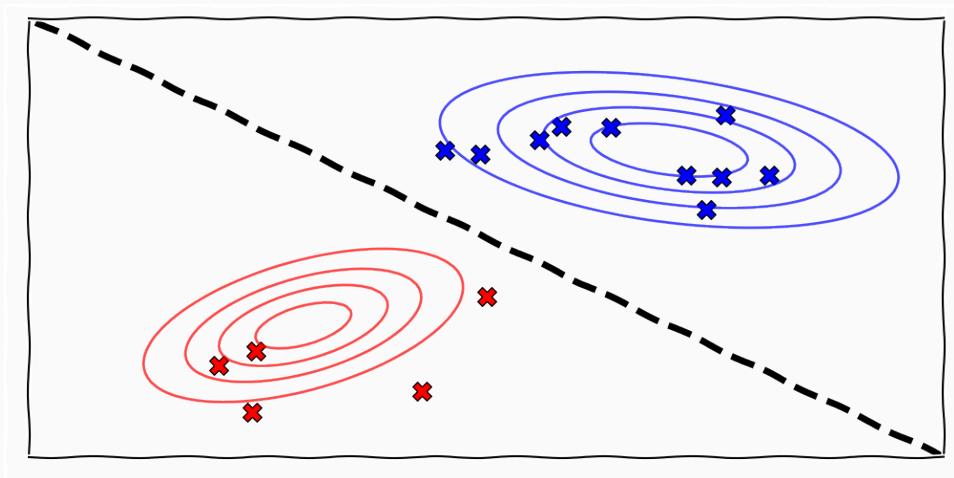
**Data Distribution** $\mathcal{D}$ probability distribution governing the measurements

**Data Generation** $f : \mathcal{X} \rightarrow \mathcal{Y}$ the underlying generating process that we wish to recover

**Data Distribution** $\mathcal{D}$ probability distribution governing the measurements

**Data Generation** $f : \mathcal{X} \to \mathcal{Y}$ the underlying generating process that we wish to recover

**Prediction Rule** $h : \mathcal{X} \to \mathcal{Y}$ what we wish to recover, the object that encodes the recovered knowledge

$$L_{\mathcal{D},f}(h) := \mathcal{D}(\{x : h(x) \neq f(x)\})$$

- measure of success as probability of misclassified points (true risk)

$$L_{\mathcal{D},f}(h) := \mathcal{D}(\{x : h(x) \neq f(x)\})$$
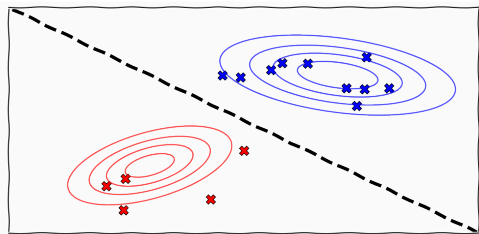
- measure of success as probability of misclassified points (true risk)
- we do not have access to $\mathcal{D}$

$$L_{\mathcal{D},f}(h) := \mathcal{D}(\{x : h(x) \neq f(x)\})$$

- measure of success as probability of misclassified points (true risk)
- we do not have access to $\mathcal{D}$
- we do not have access to $f$

$$L_{\mathcal{S}}(h) := \frac{|\{i \in [m] : h(x_i) \neq y_i\}|}{m}$$

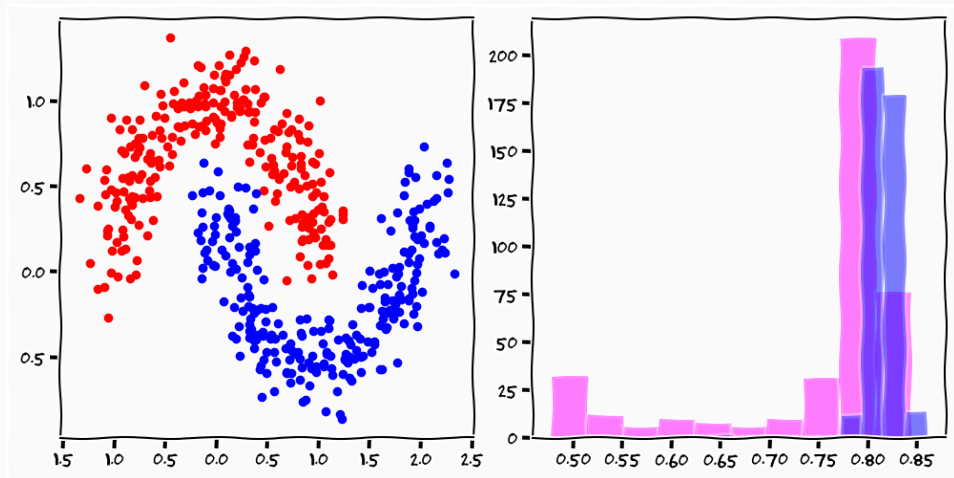- We assume that $\mathcal{S} \sim \mathcal{D}$
- Empirical measure of risk

$$h_{\mathcal{S}}(x) = \begin{cases} y_i & \text{if } \exists i \in [m] \text{s.t. } x_i = x \\ 0 & \text{otherwise} \end{cases}$$

- $L_{\mathcal{S}}(h_{\mathcal{S}}) = 0$ for all training data-sets
- if label $0$ corresponds to red $L_{\mathcal{D}}(h_{\mathcal{S}}) = \frac{1}{3}$
- if label $0$ corresponds to blue $L_{\mathcal{D}}(h_{\mathcal{S}}) = \frac{2}{3}$

$$\mathcal{D} = \frac{1}{3}\mathcal{N}(\cdot, \cdot) + \frac{2}{3}\mathcal{N}(\cdot, \cdot)$$

18

$$L_{\mathcal{S}}(A(\mathcal{S})) := \frac{|\{i \in [m] : h(x_i) \neq y_i\}|}{m}$$

- We use an algorithm $A : \mathcal{S} \to h$ to find a hypothesis

$$h_{\mathcal{S}} \in \operatorname*{argmin}_{h \in \mathcal{H}} L_{\mathcal{S}}(h)$$

- We cannot parametrise all possible hypothesis

$m$ How much data do I need?

$m$  How much data do I need?

$A$  How much does my solution depend on what I find?

$m$ How much data do I need?

$A$ How much does my solution depend on what I find?

$H$ How does my solution depend on the hypothesis class I choose?

- We will assume the probability to get a non-representative sample to be $\delta$

- We will assume the probability to get a non-representative sample to be $\delta$
  - $(1 - \delta)$ is the confidence in our prediction

- We will assume the probability to get a non-representative sample to be $\delta$
  - $(1 - \delta)$ is the confidence in our prediction
- We have a threshold $\epsilon$ where we accept the hypothesis

$$L_{\mathcal{D},f}(h_{\mathcal{S}}) \leq \epsilon$$

- We will assume the probability to get a non-representative sample to be $\delta$
  - $(1 - \delta)$ is the confidence in our prediction
- We have a threshold $\epsilon$ where we accept the hypothesis

$$L_{\mathcal{D},f}(h_{\mathcal{S}}) \leq \epsilon$$

- We are looking for a *hypothesis* that will be probably (with confidence $(1 - \delta)$) approximately (up to an error $\epsilon$) correct.
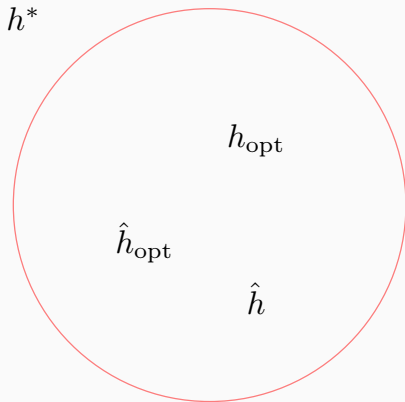
$$\mathcal{S} \sim \mathcal{D}^m$$

- each sample is *independent* from the other
- the *order* of samples does not effect the data distribution
- the sampling process does not effect the "world"

$$m \geq \frac{\log\left(\frac{|\mathcal{H}|}{\delta}\right)}{\epsilon}$$

- PAC learning allows us to provide bounds on the learning procedure
- How much data do we need?
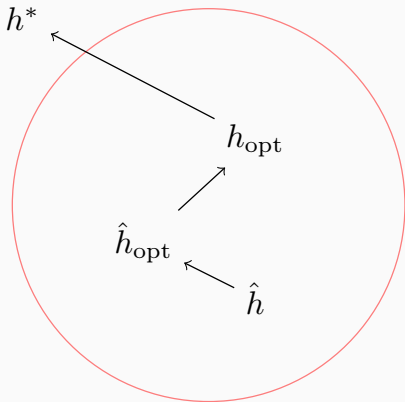- How large hypothesis class can we allow?

$h^*$ the optimal predictor

$h_{\mathbf{opt}}$ the optimal hypothesis

$\hat{h}_{\mathbf{opt}}$ the optimal hypothesis on training data

$\hat{h}$ the hypothesis found by learning algorithm

$$\epsilon(\hat{h}) - \epsilon(h^*)$$

$$= \underbrace{\epsilon(h_{\mathsf{opt}}) - \epsilon(h^*)}_{\text{Approximation}}$$

$$+ \underbrace{\epsilon(\hat{h}_{\mathsf{opt}}) - \epsilon(h_{\mathsf{opt}})}_{\text{Estimation}}$$

$$+ \underbrace{\epsilon(\hat{h}) - \epsilon(\hat{h}_{\mathsf{opt}})}_{\text{Optimisation}}$$

**High Complexity** low bias ($\epsilon_{\mathsf{app}}$ small), but high risk of overfitting ($\epsilon_{\mathsf{est}}$ large)

**Low Complexity** high bias ($\epsilon_{\mathsf{app}}$ large), low risk of overfitting ($\epsilon_{\mathsf{est}}$ small)
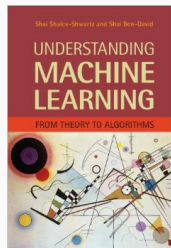
**Theorem (The No-Free-Lunch Theorem)**
*Let $A$ be any learning algorithm fo the task of binary classification with respect to $0 - 1$ loss over the domain $\mathcal{X}$. Let $m$ be any number smaller than $\frac{|\mathcal{X}|}{2}$. Then there exists a distribution $\mathcal{D}(\{\mathcal{X} \times \{0,1\}\})$ such that,*

- *There exists a function $f : \mathcal{X} \to \{0,1\}$ with $L_{\mathcal{D}}(f) = 0$*
- *With probability at least $\frac{1}{7}$ over the choice of $\mathcal{S} \sim \mathcal{D}^m$ we have $L_{\mathcal{D}}(A(S)) \geq \frac{1}{8}$*
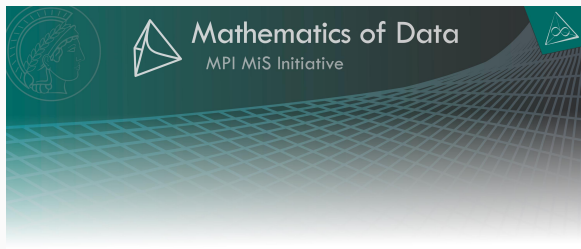
- There exists no universal learner
- For every learner there exist a task on which it fails
- Every algorithm that learns something useful does so by assumptions

- There exists no universal learner
- For every learner there exist a task on which it fails
- Every algorithm that learns something useful does so by assumptions
- There is no free lunch algorithm

- We can never have sufficient data
- We can never find a method that will guarantee to find the right solution
- We can never be certain about the true risk of our outcome

- Shai Shalev-Shwartz et al. (2014). *Understanding Machine Learning: From Theory to Algorithms*. New York, NY, USA: Cambridge University Press https://www.cs.huji.ac.il/~shais/UnderstandingMachineLearning/
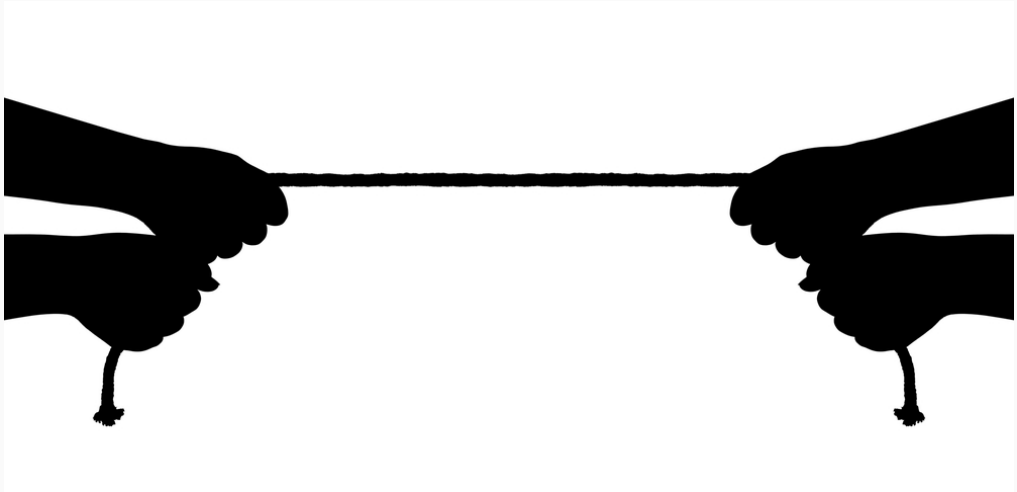
- O. Bousquet et al. (2004). "Introduction to Statistical Learning Theory". In: vol. Lecture Notes in Artificial Intelligence 3176. Heidelberg, Germany: Springer, pp. 169–207, http://www.econ.upf.edu/~lugosi/mlss_slt.pdf

HOW TO MAKE AN **OMELETTE**

3 EGGS · WHISK · HEAT OLIVE OIL

POUR EGGS INTO PAN · USE A SPATULA · TO PULL SOLIDIFIED EDGES TOWARDS THE MIDDLE

KEEP PULLING UNTIL THERE'S JUST A BIT OF FLUID LEFT · FOLD IN THREE AND SERVE

**Access** enormous inductive bias in what data to acquire

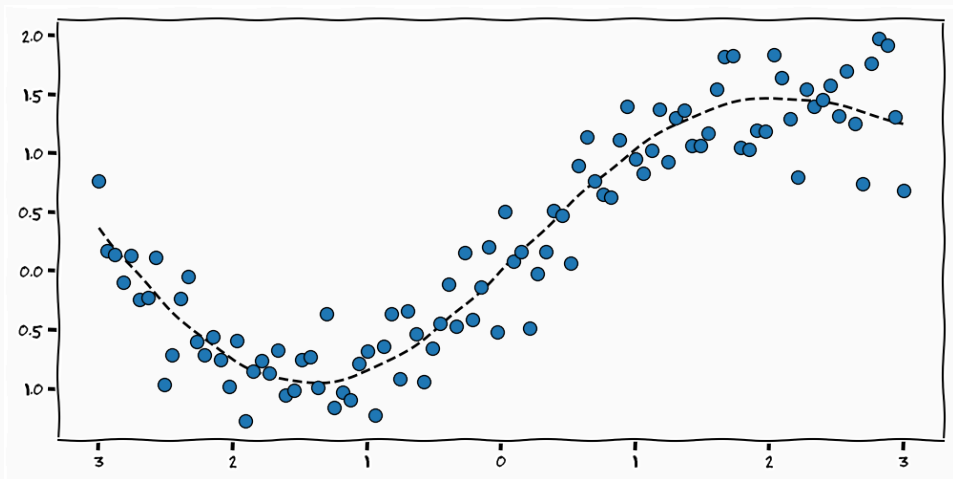**Access** enormous inductive bias in what data to acquire

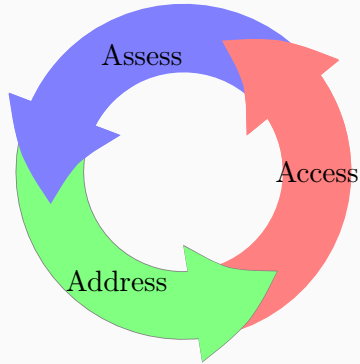**Assess** human bias in what questions will probably be asked

**Access** enormous inductive bias in what data to acquire

**Assess** human bias in what questions will probably be asked

**Address** "it is just curve fitting"

**Today/Mon** Generalised Linear Models

**Today/Mon** Generalised Linear Models

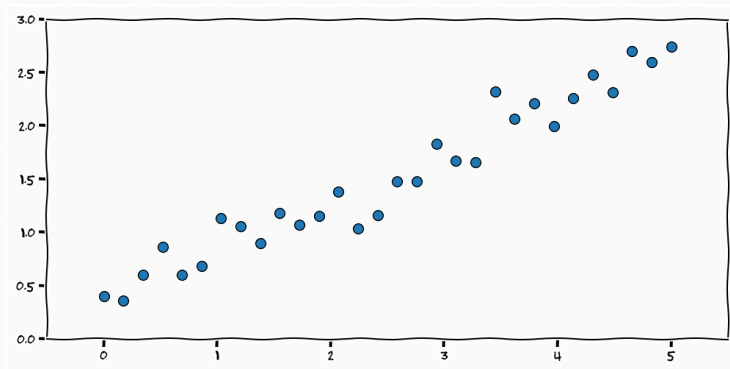**Wen** Unsupervised Learning

**Today/Mon** Generalised Linear Models

**Wen** Unsupervised Learning

**Fri** Approximate Inference
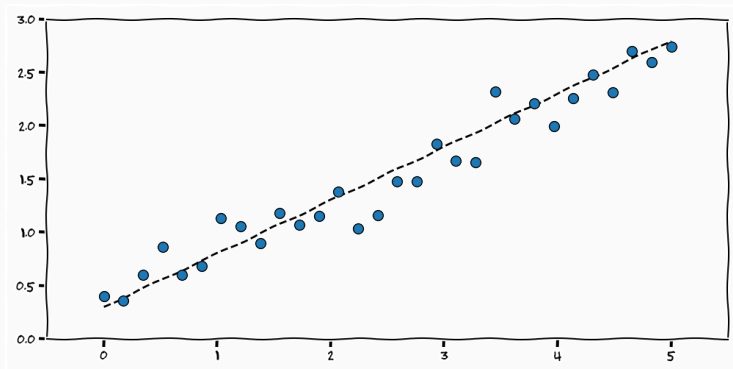
# Generalised Linear Models

$$h \in \mathcal{H}$$

$\mathbf{x} \in \mathcal{X}$ explanatory variable

$y \in \mathcal{Y}$ response variable

**Task** *explain the response by the explanatory variables*

$$y_i = \sum_{j=1}^{d} \beta_j x_{ij} + \epsilon, \qquad \epsilon \sim \mathcal{N}(0, \sigma^2).$$

44

$$\mathbb{E}[y_i \mid \mathbf{x}_i] = \mathbb{E}\left[\sum_{j=1}^{d} \beta_j x_{ij} + \epsilon\right]$$

$$\mathbb{E}[y_i \mid \mathbf{x}_i] = \mathbb{E}\left[\sum_{j=1}^{d} \beta_j x_{ij} + \epsilon\right]$$

$$= \mathbb{E}\left[\sum_{j=1}^{d} \beta_j x_{ij}\right] + \mathbb{E}[\epsilon]$$

$$\mathbb{E}[y_i \mid \mathbf{x}_i] = \mathbb{E}\left[\sum_{j=1}^{d} \beta_j x_{ij} + \epsilon\right]$$

$$= \mathbb{E}\left[\sum_{j=1}^{d} \beta_j x_{ij}\right] + \mathbb{E}[\epsilon]$$

$$= \sum_{j=1}^{d} \beta_j x_{ij} + 0.$$

$$y_i = \sum_{j=1}^{d} \beta_j x_{ij} + \epsilon,$$

$$y_i = \sum_{j=1}^{d} \beta_j x_{ij} + \epsilon,$$

$$y_i + \epsilon = \sum_{j=1}^{d} \beta_j x_{ij},$$

# Linear Regression

$$y_i = \sum_{j=1}^{d} \beta_j x_{ij} + \epsilon,$$

$$y_i + \epsilon = \sum_{j=1}^{d} \beta_j x_{ij},$$

$$\hat{y}_i = \sum_{j=1}^{d} \beta_j x_{ij},$$

## Linear Regression

$$y_i = \sum_{j=1}^{d} \beta_j x_{ij} + \epsilon,$$

$$y_i + \epsilon = \sum_{j=1}^{d} \beta_j x_{ij},$$

$$\hat{y}_i = \sum_{j=1}^{d} \beta_j x_{ij},$$

$$\hat{y}_i \sim \mathcal{N}(y_i, \sigma^2) = \mathcal{N}\left(\sum_{j=1}^{d} \beta_j x_{ij}, \sigma^2\right),$$

$$g(\mathbb{E}[y_i \mid \mathbf{x}_i]) = \sum_{j=1}^{d} \beta_j x_{ij},$$

$g(\cdot)$ link function

$y \sim \mathcal{D}$ Exponential Dispersion Family

$\sum_{j=1}^{d} \beta_j x_{ij}$ Linear predictor

$$\mathbb{E}[y_i \mid \mathbf{x}_i] = g^{-1}(\sum_{j=1}^{d} \beta_j x_{ij}),$$

- The inverse of the *link* maps the linear predictor to the first moment of the response

- Linear regression the link is identity

[1]https://towardsdatascience.com/
glms-part-iii-deep-neural-networks-as-recursive-generalized-linear-URL

$$\mathbb{E}[y_i \mid \mathbf{x}_i] = g^{-1}(\sum_{j=1}^{d} \beta_j x_{ij}),$$

- The inverse of the *link* maps the linear predictor to the first moment of the response

- Linear regression the link is identity

- Looks an awful lot like a neural network[1]

---

[1] https://towardsdatascience.com/
glms-part-iii-deep-neural-networks-as-recursive-generalized-linear-URL

$$f(y; \theta, \phi) = e^{\frac{\theta y - b(\theta)}{a(\phi)} + c(y, \phi)},$$

$\theta$ location parameter

$\phi$ scale parameter

[2] https://en.wikipedia.org/wiki/Exponential_dispersion_model

$$f(y; \mu, \sigma^2) = e^{\frac{\mu y - \frac{1}{2}\mu^2}{\sigma^2} - \frac{y^2}{2\sigma^2} - \frac{1}{2}\ln(2\pi\sigma^2)}$$

- $\theta = \mu$
- $\phi = \sigma^2$
- $b(\theta) = \frac{1}{2}\mu^2$

- $a(\phi) = \sigma^2$
- $c(y, \phi) = \frac{y^2}{2\sigma^2} - \frac{1}{2}\ln(2\pi\sigma^2)$

$$\mathbb{E}[y \mid \mathbf{x}] = \frac{\partial}{\partial \theta} b(\theta)$$
$$\mathbb{V}[y \mid \mathbf{x}] = a(\phi) \frac{\partial^2}{\partial \theta^2} b(\theta).$$

- Through a consistent parametrisation we can generalise the moment calculations

| Model | Response Variable | Link | Explanatory Variable |
|---|---|---|---|
| Linear Regression | Normal | Identity | Continuous |
| Logistic Regression | Binomial | Logit | Mixed |
| Poisson Regression | Poisson | Log | Mixed |
| ANOVA | Normal | Identity | Categorical |
| ANCOVA | Normal | Identity | Mixed |
| Loglinear | Poisson | Log | Categorical |
| Multinomial response | Multinomial | Generalized Logit | Mixed |

# Summary

- Brief introduction to statistical learning theory
- Take home
  - ML models and algorithms is only a small part of the story
  - we are doing a lot better than we should be
  - we are not sure what we are doing but it somehow works

- Brief introduction to statistical learning theory
- Take home
  - ML models and algorithms is only a small part of the story
  - we are doing a lot better than we should be
  - we are not sure what we are doing but it somehow works
  - it is not explicit knowledge that pushes data-science forward, it is tacit and implicit

- Generalised Linear Models
  - "main" tool for statisticians
  - very well studied, excellent literature
  - linearity provides means of interpretability
  - excellent software packages `statsmodels`

## The rest

**Monday (22/11)** Lecture: Generalised Linear Models

**Tuesday (23/11)** Lab: Generalised Linear Models

**Wendesday (24/11)** Lecture: Unsupervised Learning/Visualisation

**Thursday (25/11)** Tick: Generalised Linear Models

**Friday (26/11)** Lecture: Statistical Inference

eof

# References

📄 Bishop, Christopher M. (2006). *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Secaucus, NJ, USA: Springer-Verlag New York, Inc.

📄 Bousquet, O., S. Boucheron, and G. Lugosi (2004). "Introduction to Statistical Learning Theory". In: vol. Lecture Notes in Artificial Intelligence 3176. Heidelberg, Germany: Springer, pp. 169–207.

📄 McCullagh, P. and J. A. Nelder (1989). *Generalized Linear Models*. London, UK: Chapman Hall / CRC: Chapman Hall / CRC.

📄 Shalev-Shwartz, Shai and Shai Ben-David (2014). *Understanding Machine Learning: From Theory to Algorithms*. New York, NY, USA: Cambridge University Press.