



Machine Learning 101

DESY Summer Program 2025 | July 31st 2025

DESY : Stephen Jiggins

Contact Info: stephen.jiggins@desy.de

What is Machine Learning?



2018: \$432,000 painting sold at *Christie's* using a GAN method via *Obvious*

$$\min_{\mathcal{G}} \max_{\mathcal{D}} \mathbb{E}_x [\log(\mathcal{D}(x))] + \mathbb{E}_y [\log(1 - \mathcal{D}(\mathcal{G}(y)))]$$

Portrait of Edmond de Belamy

What is Machine Learning?



2018: \$432,000 painting sold at *Christie's* using a GAN method via *Obvious*

$$\min_{\mathcal{G}} \max_{\mathcal{D}} \mathbb{E}_x [\log(\mathcal{D}(x))] + \mathbb{E}_y [\log(1 - \mathcal{D}(\mathcal{G}(y)))]$$

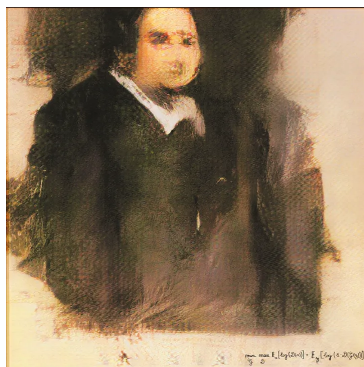
Portrait of Edmond de Belamy



2022: Colorado state fair winner using *DALL · E2*.

Sparks the question of morality in the age of AI realism

What is Machine Learning?



2018: \$432,000 painting sold at Christie's using a GAN method via *Obvious*

$$\min_{\theta} \max_{\phi} \mathbb{E}_x [\log(\mathcal{D}(x))] + \mathbb{E}_y [\log(1 - \mathcal{D}(G(y)))]$$

Portrait of Edmond de Belamy



2022: Colorado state fair winner using DALL · E2.

Sparks the question of morality in the age of AI realism

2023: Tricking the world is not hard with such technology!



I thought I was immune to being fooled online. Then I saw the pope in a coat
Joel Golby

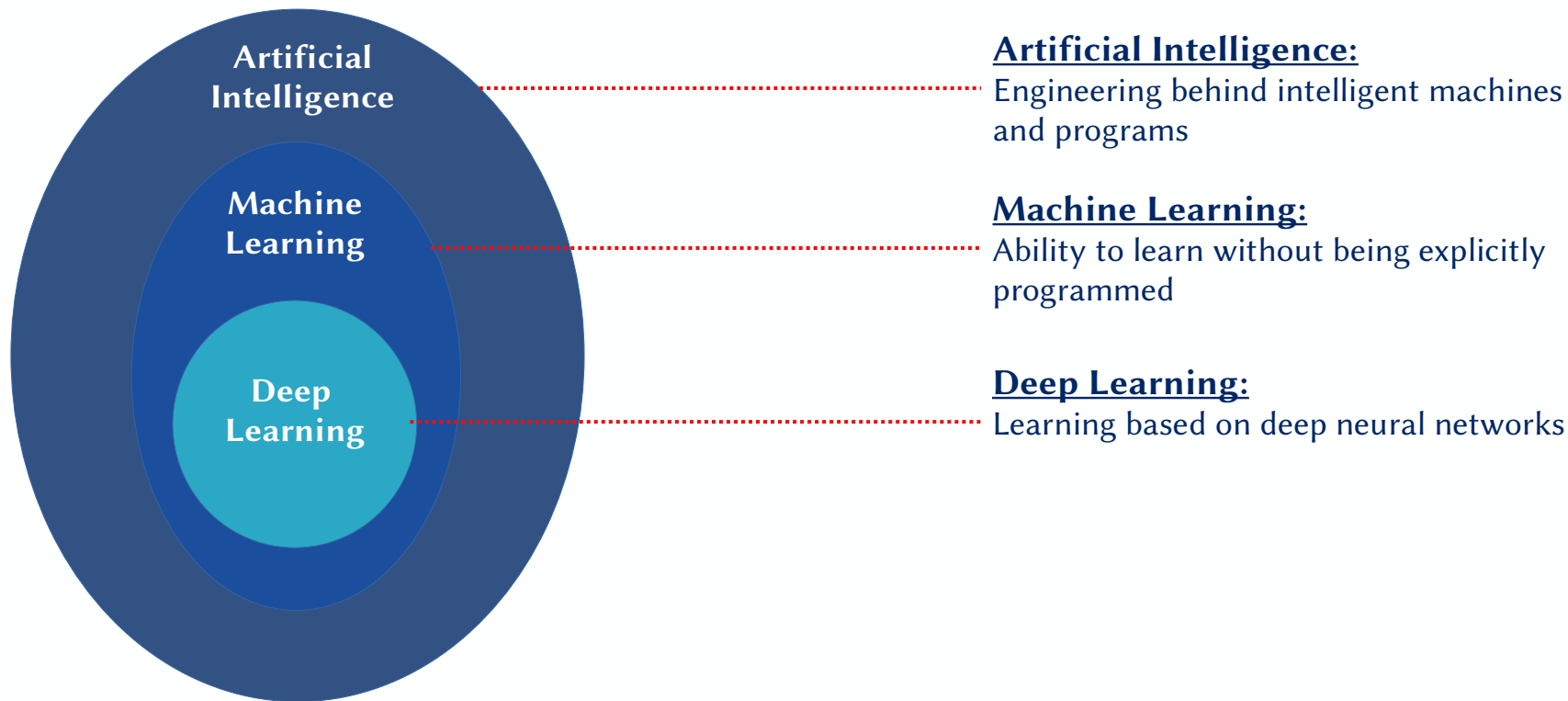
An encounter with an AI-generated image of his holiness has changed me: I now have sympathy for credulous baby boomers



📸 'I thought wearing a really big coat and looking like a Metal Gear Solid 2 boss battle might have been part of his ongoing cool guy shtick. Lord, forgive me.' Photograph: Reddit

What is Machine Learning?

5



What is Machine Learning?

*‘Giving computers the ability to learn
without explicitly programming them’*

- Arthur Samuel, 1959

What is Machine Learning?

~~'Giving computers the ability to learn without explicitly programming them'~~

7

- Arthur Samuel, 1959

Paper: [link](#)

- *'A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P , if its performance at tasks in T , as measured by P , improves with experience E .'* - Mitchell 1997

What is Machine Learning?

~~'Giving computers the ability to learn without explicitly programming them'~~

- Arthur Samuel, 1959

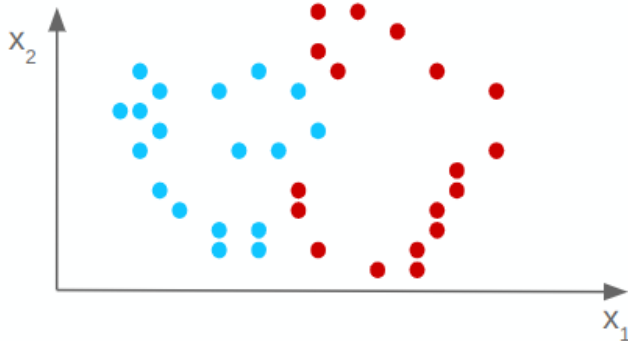
Paper: [link](#)

8

- *'A computer program is said to learn from **experience** E with respect to some class of tasks T and performance measure P , if its performance at tasks in T , as measured by P , improves with experience E .'* - Mitchell 1997

- **Experience (E):** The stimulus that drives learning, is the set of **examples** x , or a **dataset** \mathcal{D} of many examples $\mathbf{x} = \{x\}_N$:

Supervised: The dataset examples have an associated label or target $\mathbf{y} = \{y\}_N$



What is Machine Learning?

~~'Giving computers the ability to learn without explicitly programming them'~~

- Arthur Samuel, 1959

Paper: [link](#)

9

- *'A computer program is said to learn from **experience** E with respect to some class of tasks T and performance measure P , if its performance at tasks in T , as measured by P , improves with experience E .'* - Mitchell 1997

- **Experience (E):** The stimulus that drives learning, is the set of **examples** x , or a **dataset** \mathcal{D} of many examples $\mathbf{x} = \{x\}_N$:

Supervised: The dataset examples have an associated label or target $\mathbf{y} = \{y\}_N$

Unsupervised: The dataset examples have **no** labels or targets



What is Machine Learning?

~~'Giving computers the ability to learn without explicitly programming them'~~

10

- Arthur Samuel, 1959

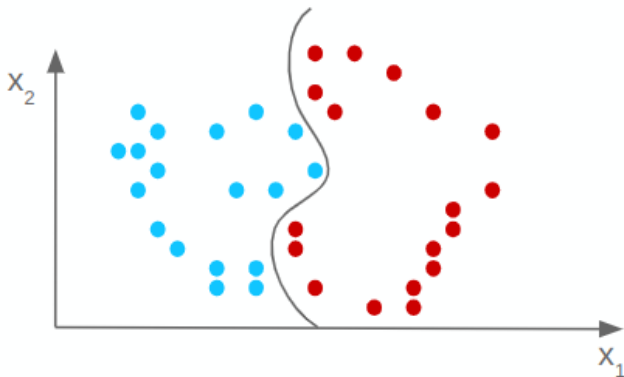
Paper: [link](#)

- 'A computer program is said to learn from experience E with respect to some class of **tasks** T and performance measure P , if its performance at tasks in T , as measured by P , improves with experience E .' - Mitchell 1997

- **Task (T):** Given an **example** $x \in \mathbb{R}^n$ the **model** f should learn a prediction $y = f(x)$. E.g.

Classification: Assign the example to a category $f: \mathbb{R}^n \rightarrow \{1, \dots, k\}$

Regression: Predict the value of a target $f: \mathbb{R}^n \rightarrow \mathbb{R}$



What is Machine Learning?

- ‘A computer program is said to learn from experience E with respect to some class of tasks T and **performance measure** P , if its performance at tasks in T , as measured by P , improves with experience E .’ - Mitchell 1997

- **Performance (P):** Evaluate the performance of the **model** f

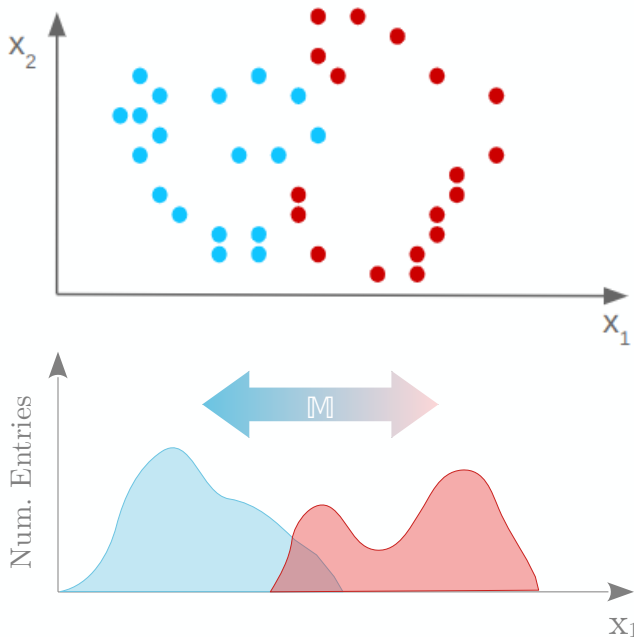
Measure Theory: Generalisation of geometric distances for a measurable space (Ω, F) , such that for a measure $\mathbb{M} : F \rightarrow [0, +\infty]$

~~‘Giving computers the ability to learn without explicitly programming them’~~

- Arthur Samuel, 1959

Paper: [link](#)

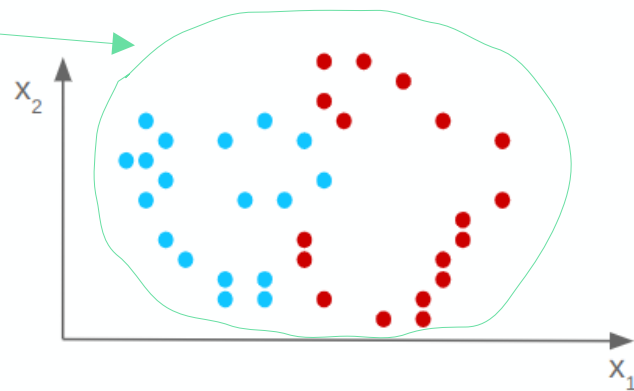
11



Basic Terminology - Day-to-day

- **Datasets:** The data from which the algorithm will need to learn from:

$$\mathcal{D} = \{ \{x_i, y_i\} \in \mathbb{R}^X \times \mathbb{R}^Y \}_N$$



Basic Terminology - Day-to-day

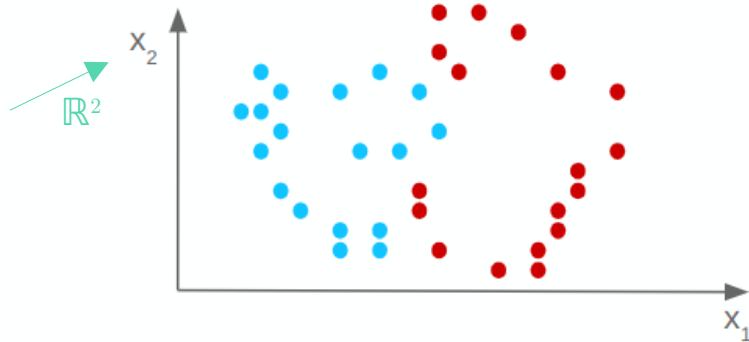
13

- **Datasets:** The data from which the algorithm will need to learn from:

$$\mathcal{D} = \{ \{x_i, y_i\} \in \mathbb{R}^X \times \mathbb{R}^Y \}_N$$

- **Features:** Each dimension of the example x , or input space \mathbb{R}^X , which represents some *feature* such as pixel color, particle energy, etc...:

$$x \in \mathbb{R}^X$$



Basic Terminology - Day-to-day

14

- **Datasets:** The data from which the algorithm will need to learn from:

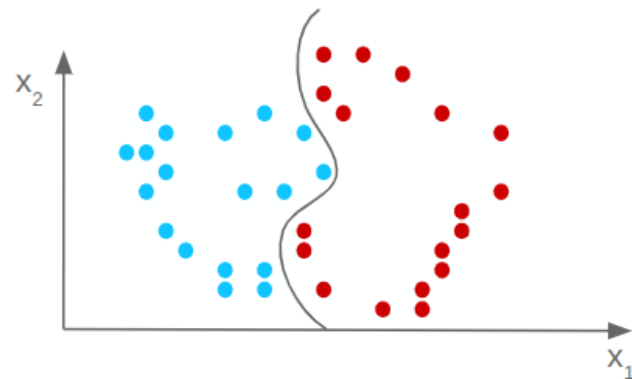
$$\mathcal{D} = \{ \{x_i, y_i\} \in \mathbb{R}^X \times \mathbb{R}^Y \}_N$$

- **Features:** Each dimension of the example x , or input space \mathbb{R}^X , which represents some *feature* such as pixel color, particle energy, etc...:

$$x \in \mathbb{R}^X$$

- **Algorithms:** The process in which the model f defined by a parameter set Φ is optimised according to some **objective function**:

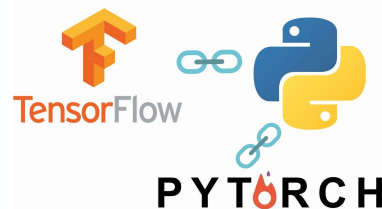
$$\Phi^* = \arg \min_{\Phi} [f]$$





Is ML useful?

- Processing Power:** The computational power of CPU/GPUs and now emerging accelerator platforms such as TPUs etc...



- ## High—Level Interpretable Languages: Julia, Python, etc...

- Bio-molecular research:

AlphaFold1.0 → 3.0 has lead to substantial advancements in predicting protein structures:

DREAMING UP PROTEINS

Researchers used deep neural networks to invent, or 'hallucinate', sequences of amino acids that could fold into proteins; in some cases they have synthesized these proteins to compare their actual structures with predictions.



'Hallucinated' protein (software prediction)



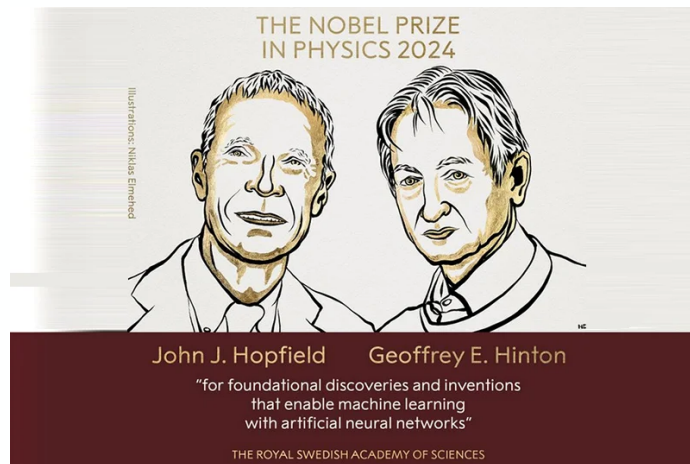
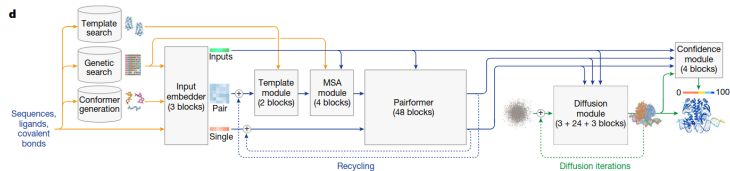
Actual structure (experimentally determined)



Superposition of hallucinated (blue) and actual (grey) structures

©nature

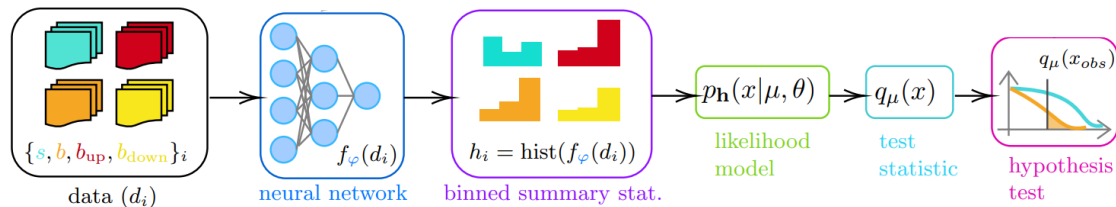
Images: Ref. 7



High Energy Physics – ML Developments

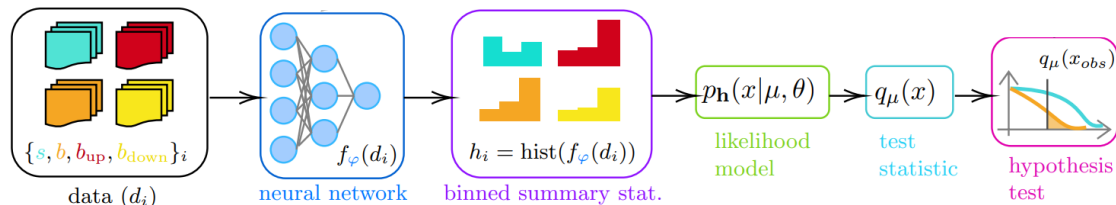
18

→ Checkout the [HEP Machine Learning Living Review](#)



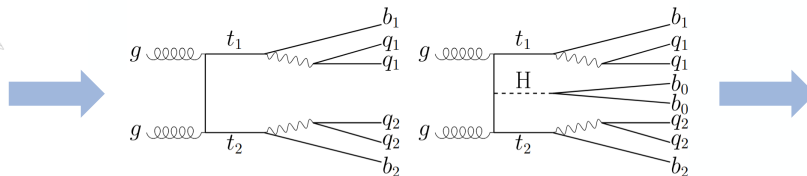
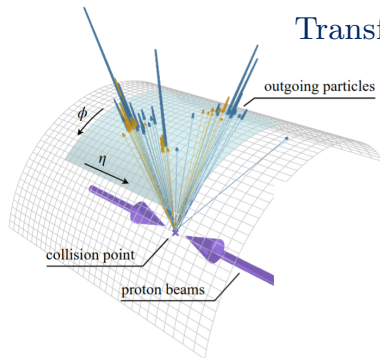
High Energy Physics – ML Developments

→ Checkout the [HEP Machine Learning Living Review](#)



SPANet Transformer - [arXiv:2106.03898](#)

Transformer model using scaled dot-product equivariant properties for top-jet reconstruction



$t\bar{t}$

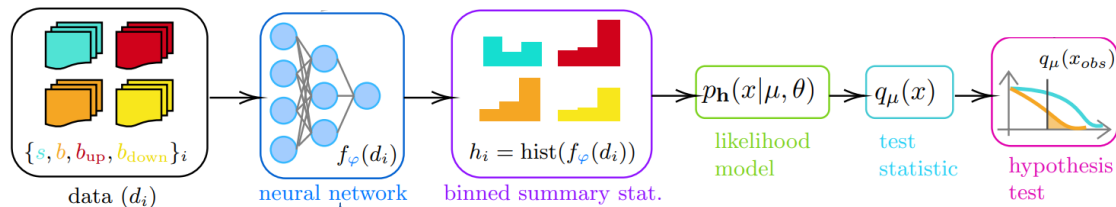
	N_{jets}	Event Fraction	SPA-NET Efficiency		χ^2 Efficiency	
			Event	Top Quark	Event	Top Quark
All Events	≤ 6	0.245	0.643	0.696	0.424	0.484
	≤ 7	0.282	0.601	0.667	0.389	0.460
	≥ 8	0.320	0.528	0.613	0.309	0.384
	Inclusive	0.848	0.586	0.653	0.392	0.457
Complete Events	≤ 6	0.074	0.803	0.837	0.593	0.643
	≤ 7	0.105	0.667	0.754	0.413	0.530
	≥ 8	0.145	0.521	0.662	0.253	0.410
	Inclusive	0.325	0.633	0.732	0.456	0.552

$t\bar{t}H$

	N_{jets}	Event Fraction	SPA-NET Efficiency			χ^2 Efficiency		
			Event	Higgs	Top	Event	Higgs	Top
All Events	≤ 8	0.261	0.370	0.497	0.540	0.044	0.151	0.053
	≤ 9	0.313	0.343	0.492	0.514	0.038	0.146	0.066
	≥ 10	0.313	0.294	0.472	0.473	0.030	0.135	0.072
	Inclusive	0.972	0.330	0.485	0.502	0.039	0.146	0.062
Complete Events	≤ 8	0.042	0.532	0.657	0.663	0.016	0.151	0.063
	≤ 9	0.070	0.422	0.601	0.596	0.013	0.146	0.076
	≥ 10	0.115	0.306	0.545	0.523	0.008	0.134	0.080
	Inclusive	0.228	0.383	0.583	0.572	0.012	0.144	0.073

High Energy Physics – ML Developments

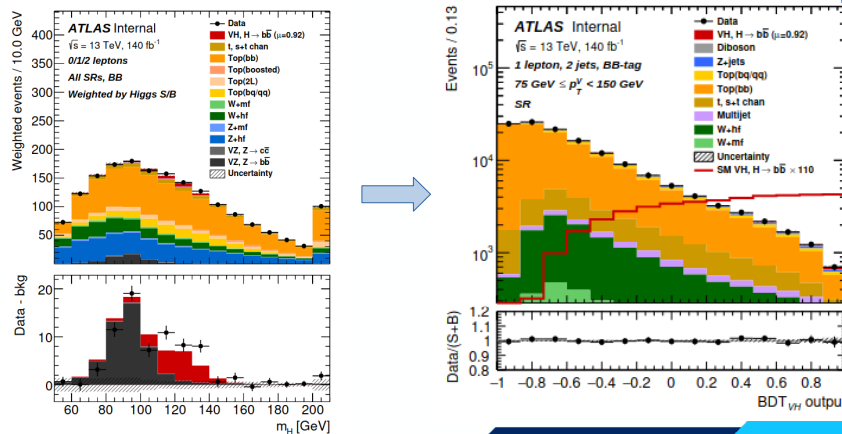
→ Checkout the [HEP Machine Learning Living Review](#)



Typical signal classification in data analyses, e.g. Higgs discovery:

Analysis	Years of data collection	Sensitivity without machine learning	Sensitivity with machine learning	Ratio of P values	Additional data required
CMS ²⁴ $H \rightarrow \gamma\gamma$	2011–2012	2.2σ , $P = 0.014$	2.7σ , $P = 0.0035$	4.0	51%
ATLAS ⁴³ $H \rightarrow \tau^+\tau^-$	2011–2012	2.5σ , $P = 0.0062$	3.4σ , $P = 0.00034$	18	85%
ATLAS ⁹⁹ $VH \rightarrow bb$	2011–2012	1.9σ , $P = 0.029$	2.5σ , $P = 0.0062$	4.7	73%
ATLAS ⁴¹ $VH \rightarrow bb$	2015–2016	2.8σ , $P = 0.0026$	3.0σ , $P = 0.00135$	1.9	15%
CMS ¹⁰⁰ $VH \rightarrow bb$	2011–2012	1.4σ , $P = 0.081$	2.1σ , $P = 0.018$	4.5	125%

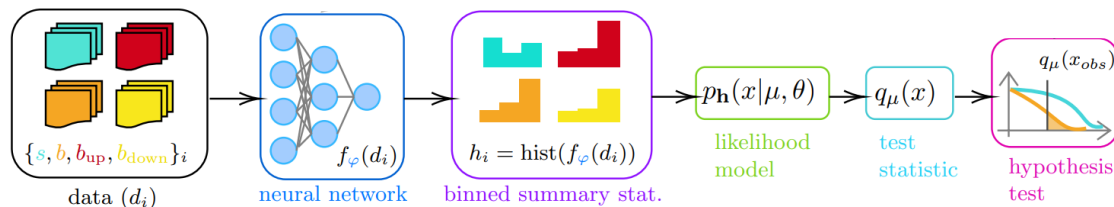
Source: A. Radovic et al., Nature 560(2018) no. 7716,41



High Energy Physics – ML Developments

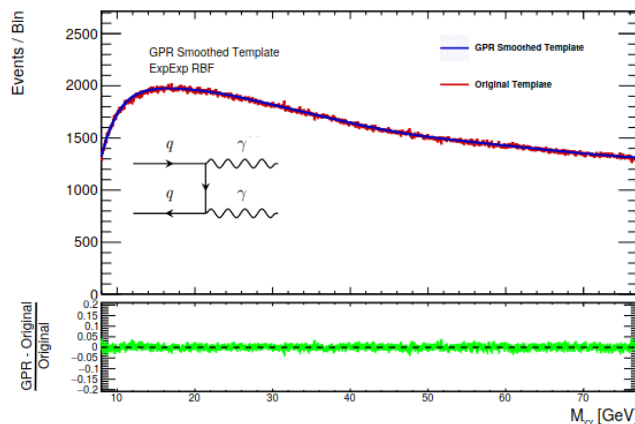
21

→ Checkout the [HEP Machine Learning Living Review](#)



Background Estimation by
smoothing data-driven
backgrounds using gaussian
processes for $H \rightarrow \gamma\gamma$

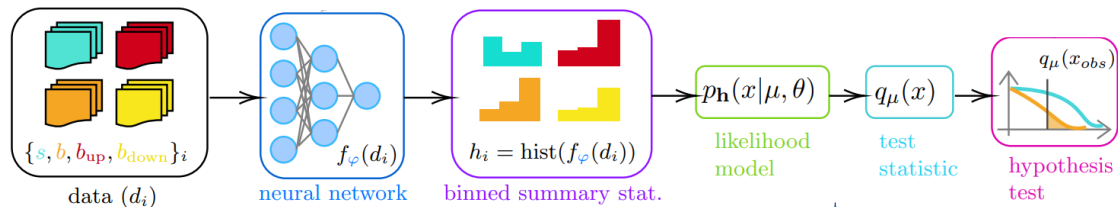
Source: [CDS-2791079](#)



High Energy Physics – ML Developments

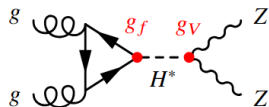
22

→ Checkout the [HEP Machine Learning Living Review](#)

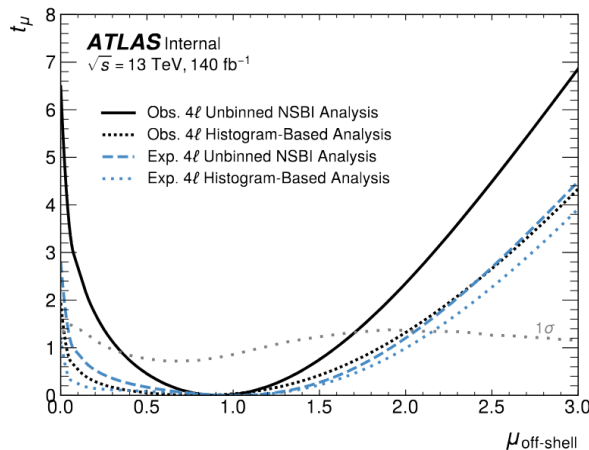


New likelihood inferencing paradigm:

H \rightarrow ZZ off-shell analysis using unbinned likelihood ratios from neural networks



Source: ATL-COM-PHYS-2024-24

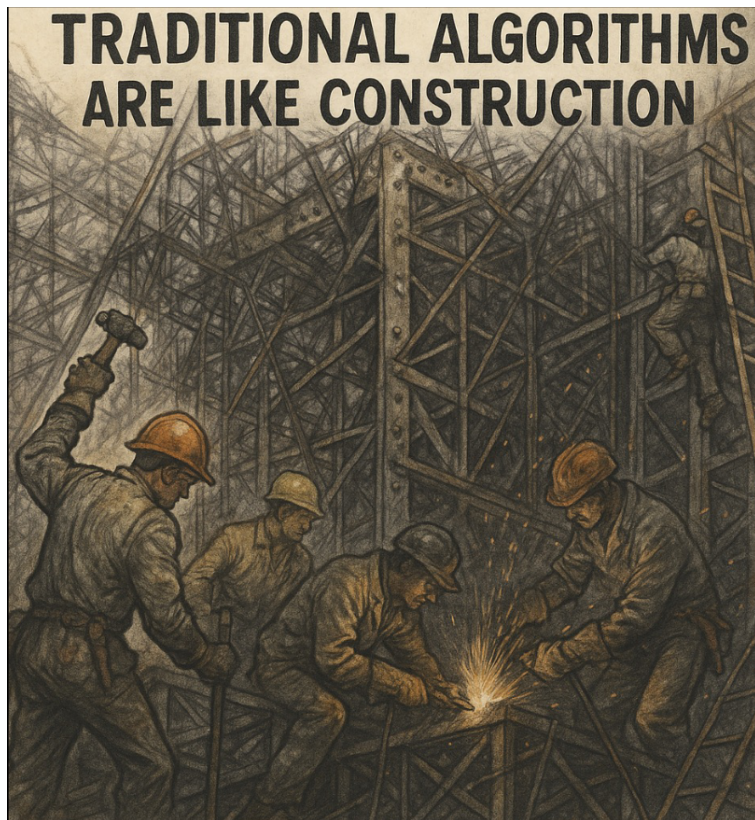




Learning Law

Machine Learning versus Algorithm Development

24



Learning Law

- **Four key parts** to the ‘Machine Learning’ process:

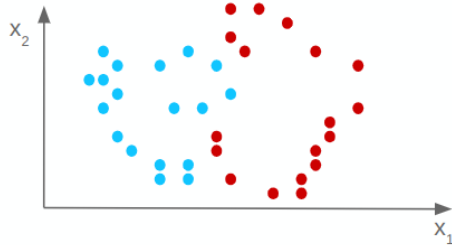
Dataset:

Supervised:

$$\mathcal{D} = \{ \{x_i, y_i\} \in \mathbb{R}^X \times \mathbb{R}^Y \}_N$$

Unsupervised:

$$\mathcal{D} = \{ \{x_i\} \in \mathbb{R}^X \}_N$$



- Four key parts to the ‘Machine Learning’ process:

Dataset:

Supervised:

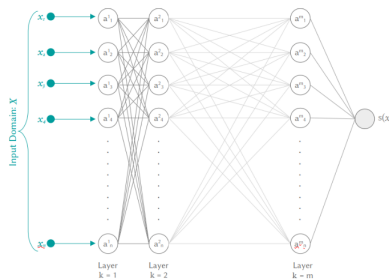
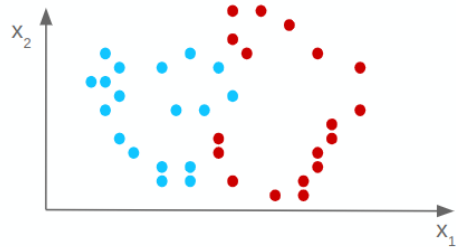
$$\mathcal{D} = \{ \{x_i, y_i\} \in \mathbb{R}^X \times \mathbb{R}^Y \}_N$$

Unsupervised:

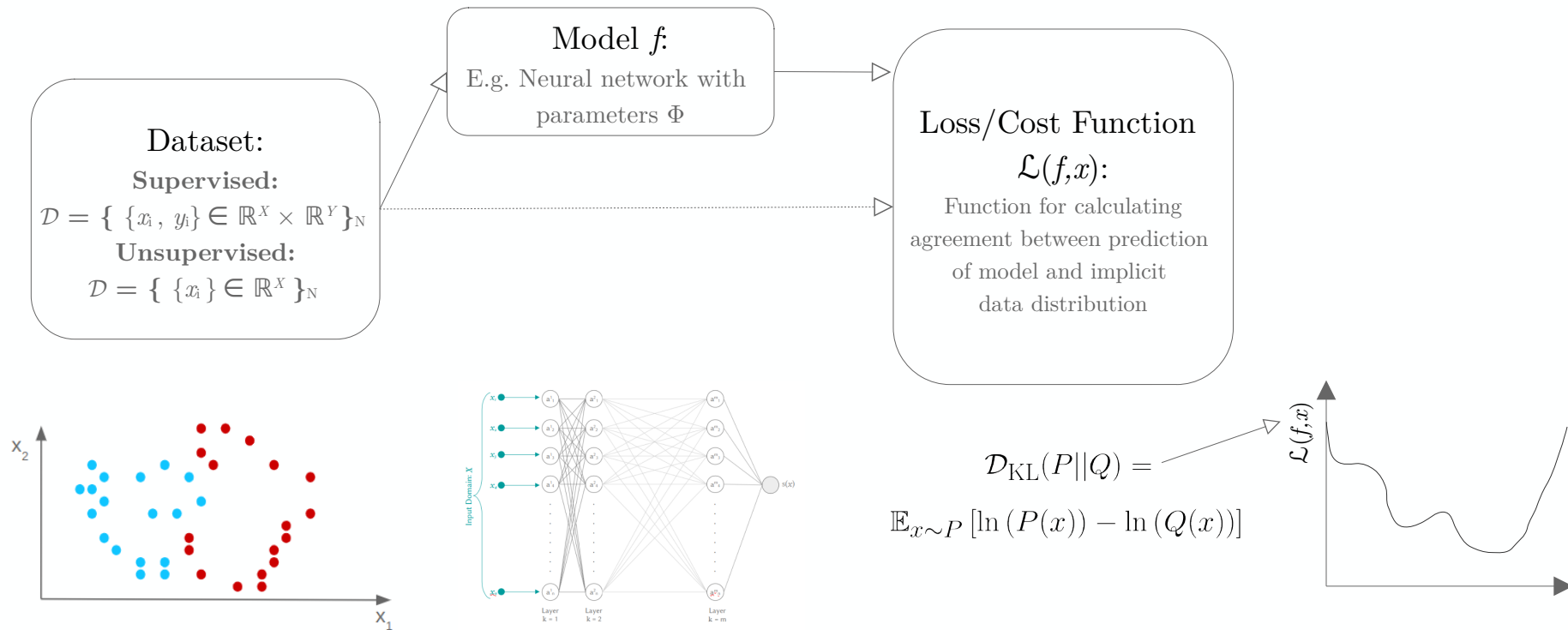
$$\mathcal{D} = \{ \{x_i\} \in \mathbb{R}^X \}_N$$

Model f :

E.g. Neural network with
parameters Φ



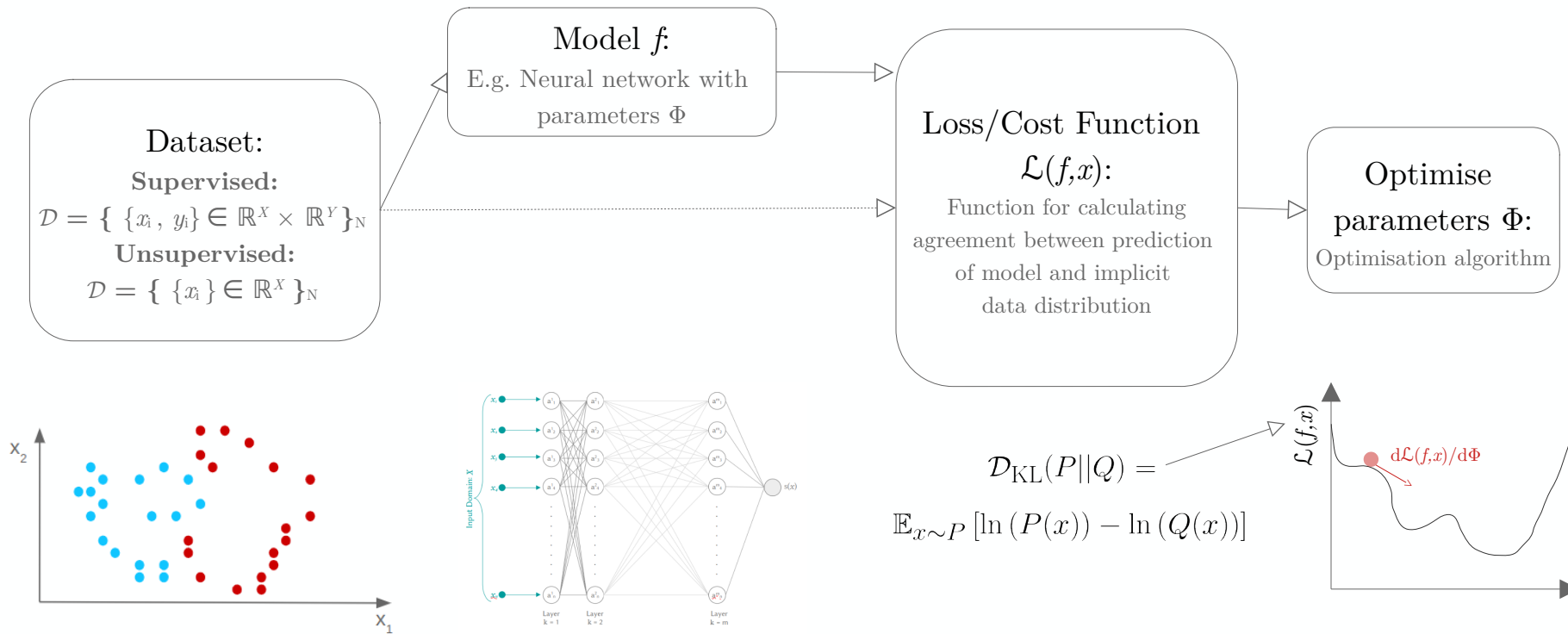
- Four key parts to the ‘Machine Learning’ process:



Learning Law - Recipe

28

- Four key parts to the 'Machine Learning' process:





Brief Statistics Review

Frequentist Statistics

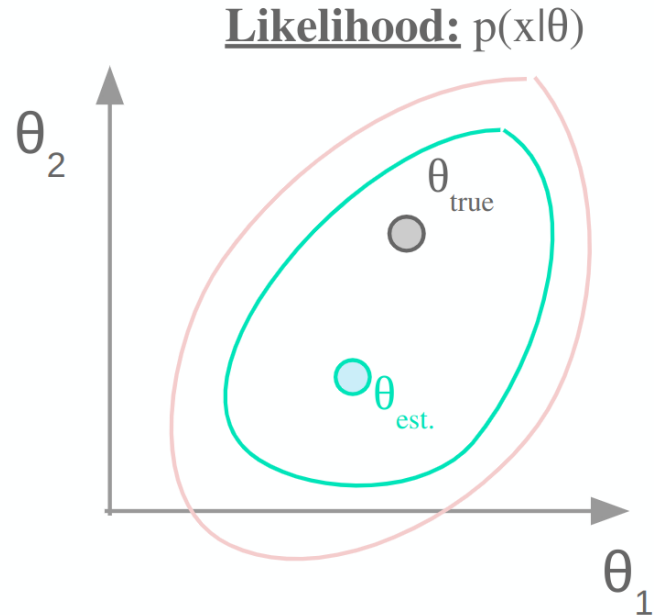
- Probability** is attributed only to the data \mathbf{x} , meaning probability of outcomes is obtained by repeatable experiments:

$$P(\mathbf{x} = x) = \lim_{n \rightarrow \infty} \frac{n_x}{N}$$

- Conditional probability**, the probability of an outcome ($\mathbf{x}=x$) conditioned on the occurrence of another random process ($\mathbf{y}=y$):

$$P(\mathbf{x} = x | \mathbf{y} = y) = \frac{P(\mathbf{x} = x, \mathbf{y} = y)}{P(\mathbf{x} = x)}$$

- Important** to note that frequentist statistics assumes that the data is drawn from $p(\mathbf{x}=x | \boldsymbol{\theta})$, with a set of parameters that characterise the underlying true distribution $\boldsymbol{\theta}$.



Frequentist Statistics

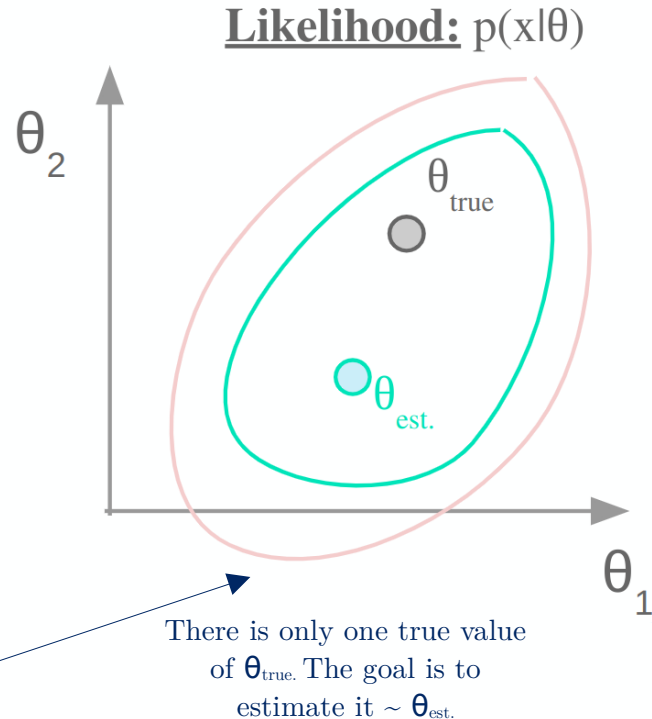
- Probability** is attributed only to the data \mathbf{x} , meaning probability of outcomes is obtained by repeatable experiments:

$$P(\mathbf{x} = x) = \lim_{n \rightarrow \infty} \frac{n_x}{N}$$

- Conditional probability**, the probability of an outcome ($\mathbf{x}=x$) conditioned on the occurrence of another random process ($\mathbf{y}=y$):

$$P(\mathbf{x} = x | \mathbf{y} = y) = \frac{P(\mathbf{x} = x, \mathbf{y} = y)}{P(\mathbf{x} = x)}$$

- Important** to note that frequentist statistics assumes that the data is drawn from $p(\mathbf{x}=x | \boldsymbol{\theta})$, with a set of parameters that characterise the underlying true distribution $\boldsymbol{\theta}$.



Bayesian Statistics

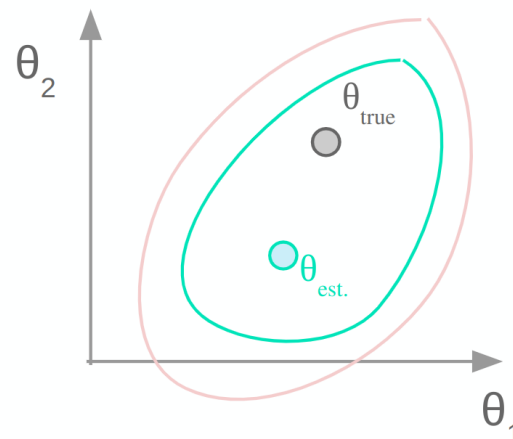
- **Probability** is a degree of belief meaning that the observed data \mathbf{x} , is not necessarily defined by repeatability:

$$P(H|\mathbf{x} = x) = \frac{P(\mathbf{x} = x|H)P(H)}{\int P(\mathbf{x} = x|H)P(H)dH}$$

Normalise over all possible
hypotheses ~ marginal
probability

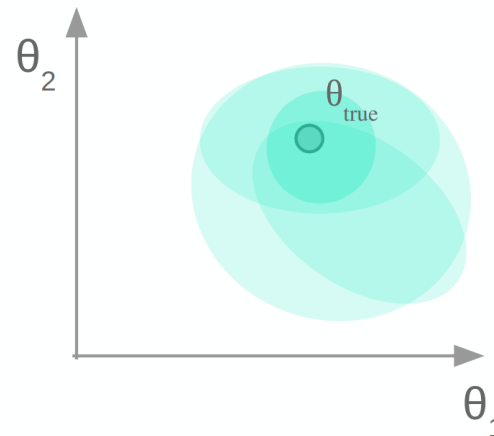
$$= \frac{P(\mathbf{x} = x|H)P(H)}{P(\mathbf{x} = x)}$$

Likelihood: $p(\mathbf{x}|\theta)$



32

Posterior: $p(\theta|\mathbf{x})$



Estimators

- **Point Estimator, or statistic**, is any function of the data that infers from the data some parameter of interest θ :

$$\hat{\theta} = g(\{x\}_m)$$

- **Point Estimator, or statistic**, is any function of the data that infers from the data some parameter of interest θ :

$$\hat{\theta} = g(\{x\}_m)$$

*Recall on slide 14 of this lecture that from a mathematical perspective a ML model is the set of optimal parameter that match **data** and **predictions**!*

Basic Terminology - Day-to-day

14

- **Datasets:** The data from which the algorithm will need to learn from:

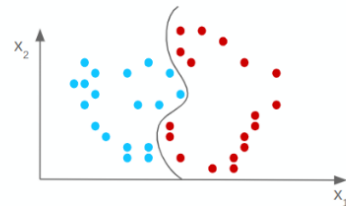
$$\mathcal{D} = \{ \{x_i, y_i\} \in \mathbb{R}^X \times \mathbb{R}^Y \}_N$$

- **Features:** Each dimension of the example x , or input space \mathbb{R}^X , which represents some *feature* such as pixel color, particle energy, etc...:

$$x \in \mathbb{R}^X$$

- **Algorithms:** The process in which the model f defined by a parameter set Φ is optimised according to some **objective function**:

$$\Phi^* = \arg \min_{\Phi} [f]$$



Estimators

- **Point Estimator, or statistic**, is any function of the data that infers from the data some parameter of interest θ :

$$\hat{\theta} = g(\{x\}_m)$$

- Is there a recipe for extracting from the data the optimal parameters or **best estimator**?

Estimators

- **Point Estimator, or statistic**, is any function of the data that infers from the data some parameter of interest θ :

$$\hat{\theta} = g(\{x\}_m)$$

- Is there a recipe for extracting from the data the optimal parameters or **best estimator**?

→ **Maximum Likelihood Method:**

General approach to estimating the point estimator

$$\hat{\theta} = \arg \max_{\theta} p(\{x\}_m | \theta)$$

$$= \arg \max_{\theta} \prod_i^m p(x_i | \theta)$$

Log-likelihood is preferred
for numerical stability:



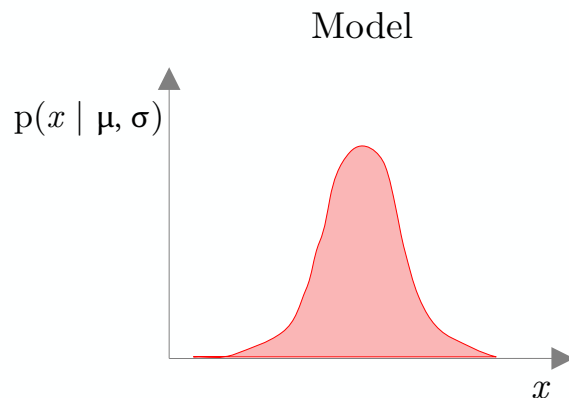
$$-\ln(\mathcal{L}(\theta | \{x\}_m)) = -\sum_i^m \ln p(x_i | \theta)$$

Guassian Estimator Example

Maximum Likelihood Method:

In a counting experiment, *model* data as Gaussian distributed:

$$p(x_i|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{1}{2} \frac{(x_i - \mu)^2}{\sigma^2} \right]$$



Guassian Estimator Example

Maximum Likelihood Method:

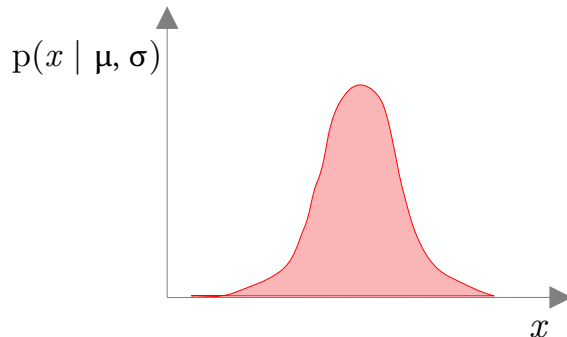
In a counting experiment, *model* data as Gaussian distributed:

$$p(x_i | \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{1}{2} \frac{(x_i - \mu)^2}{\sigma^2} \right]$$

Fill in the maximum likelihood formula:

$$-\ln(\mathcal{L}(\theta | \{x\}_m)) = -\sum_i^m \ln p(x_i | \theta) = -\sum_i^m \ln \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) - \frac{(x_i - \mu)^2}{2\sigma^2}$$

Model



Gaussian Estimator Example

Maximum Likelihood Method:

In a counting experiment, *model* data as Gaussian distributed:

$$p(x_i | \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{1}{2} \frac{(x_i - \mu)^2}{\sigma^2} \right]$$

Fill in the maximum likelihood formula:

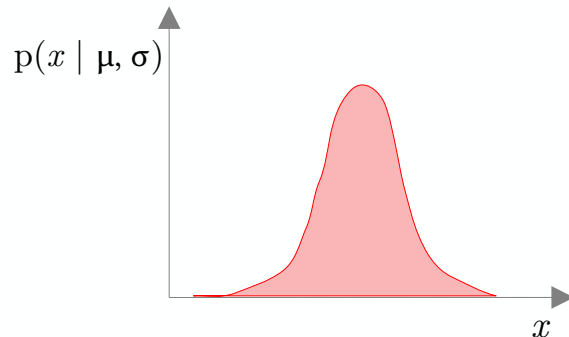
$$-\ln(\mathcal{L}(\theta | \{x\}_m)) = -\sum_i^m \ln p(x_i | \theta) = -\sum_i^m \ln \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) - \frac{(x_i - \mu)^2}{2\sigma^2}$$

We have some **data** $\{\mathbf{x}\}_m$ which forms an empirical distribution (more on this later),

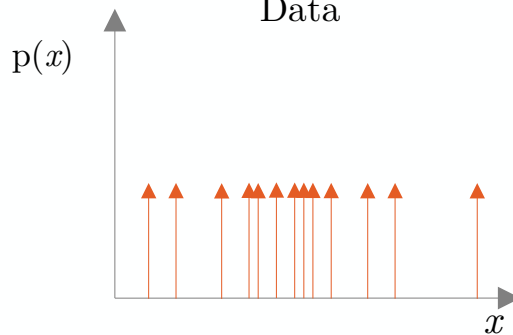
lets **maximise** the likelihood:

$$\frac{\partial}{\partial \theta} (-\ln \mathcal{L}(\theta | \{x\}_m)) \mapsto -\frac{1}{\sigma^2} \sum_i^m (x_i - \mu) = \sum_i^m x_i - m\mu = 0$$

Model



Data



Gaussian Estimator Example

Maximum Likelihood Method:

In a counting experiment, *model* data as Gaussian distributed:

$$p(x_i|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{1}{2} \frac{(x_i - \mu)^2}{\sigma^2} \right]$$

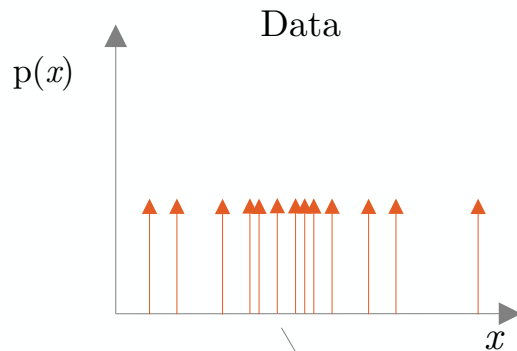
Fill in the maximum likelihood formula:

$$-\ln(\mathcal{L}(\theta|\{x\}_m)) = -\sum_i^m \ln p(x_i|\theta) = -\sum_i^m \ln \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) - \frac{(x_i - \mu)^2}{2\sigma^2}$$

We have some **data** $\{\mathbf{x}\}_m$ which forms an empirical distribution (more on this later),

lets **maximise** the likelihood:

$$\frac{\partial}{\partial \theta} (-\ln \mathcal{L}(\theta|\{x\}_m)) \mapsto -\frac{1}{\sigma^2} \sum_i^m (x_i - \mu) = \sum_i^m x_i - m\mu = 0$$



$$\hat{\mu} = \frac{1}{m} \sum_i^m x_i$$

With a bit of algebra we have derived the **sample mean**

Information Theory

- Self-Information:**

$$I(x) = -\ln(P(x))$$

- Uncertainty in an entire distribution $P(x)$ – Shannon Entropy:**

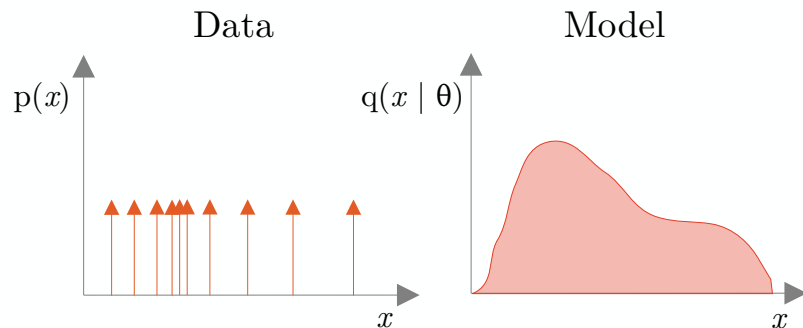
$$H(x) = \mathbb{E}_{x \sim P} [I(x)] = -\mathbb{E}_{x \sim P} [\ln(P(x))]$$

- Relative Entropy** between two distributions $P(x)$ and $Q(x)$:

$$\begin{aligned} \mathcal{D}_{\text{KL}}(P||Q) &= \mathbb{E}_{x \sim P} \left[\ln \left(\frac{P(x)}{Q(x)} \right) \right] \\ &= \mathbb{E}_{x \sim P} [\ln(P(x)) - \ln(Q(x))] \end{aligned}$$

Key Points:

- I) Data with low probability has high information content
- II) Independent data samples are additive: $I(x_1 \oplus x_2) = I(x_1) + I(x_2)$



Information Theory

- **Self-Information:**

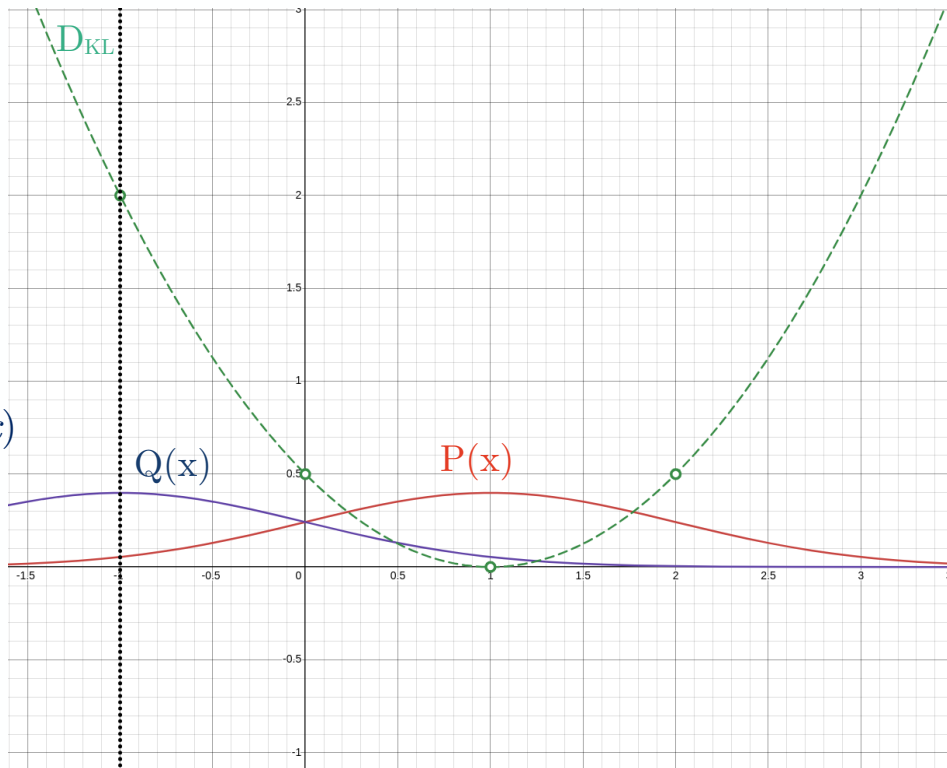
$$I(x) = -\ln(P(x))$$

- **Uncertainty in an entire distribution $P(x)$ – Shannon Entropy:**

$$H(x) = \mathbb{E}_{x \sim P} [I(x)] = -\mathbb{E}_{x \sim P} [\ln(P(x))]$$

- **Relative Entropy** between two distributions $P(x)$ and $Q(x)$:

$$\begin{aligned} \mathcal{D}_{\text{KL}}(P||Q) &= \mathbb{E}_{x \sim P} \left[\ln \left(\frac{P(x)}{Q(x)} \right) \right] \\ &= \mathbb{E}_{x \sim P} [\ln(P(x)) - \ln(Q(x))] \end{aligned}$$



Information Theory

- **Self-Information:**

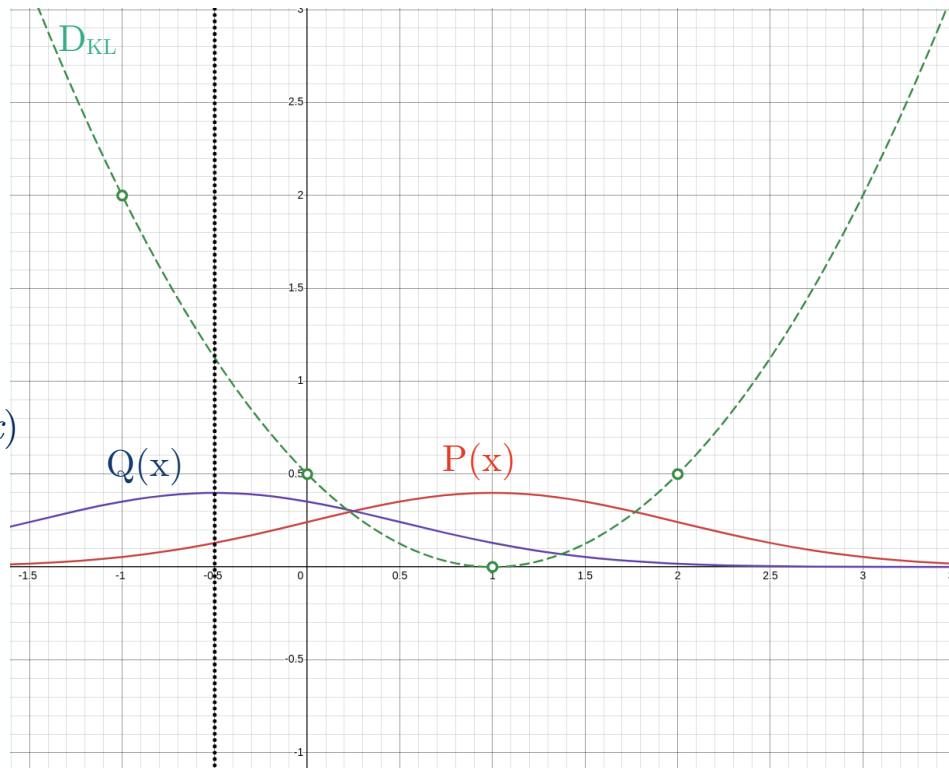
$$I(x) = -\ln(P(x))$$

- **Uncertainty in an entire distribution $P(x)$ – Shannon Entropy:**

$$H(x) = \mathbb{E}_{x \sim P} [I(x)] = -\mathbb{E}_{x \sim P} [\ln(P(x))]$$

- **Relative Entropy** between two distributions $P(x)$ and $Q(x)$:

$$\begin{aligned} \mathcal{D}_{\text{KL}}(P||Q) &= \mathbb{E}_{x \sim P} \left[\ln \left(\frac{P(x)}{Q(x)} \right) \right] \\ &= \mathbb{E}_{x \sim P} [\ln(P(x)) - \ln(Q(x))] \end{aligned}$$



Information Theory

- **Self-Information:**

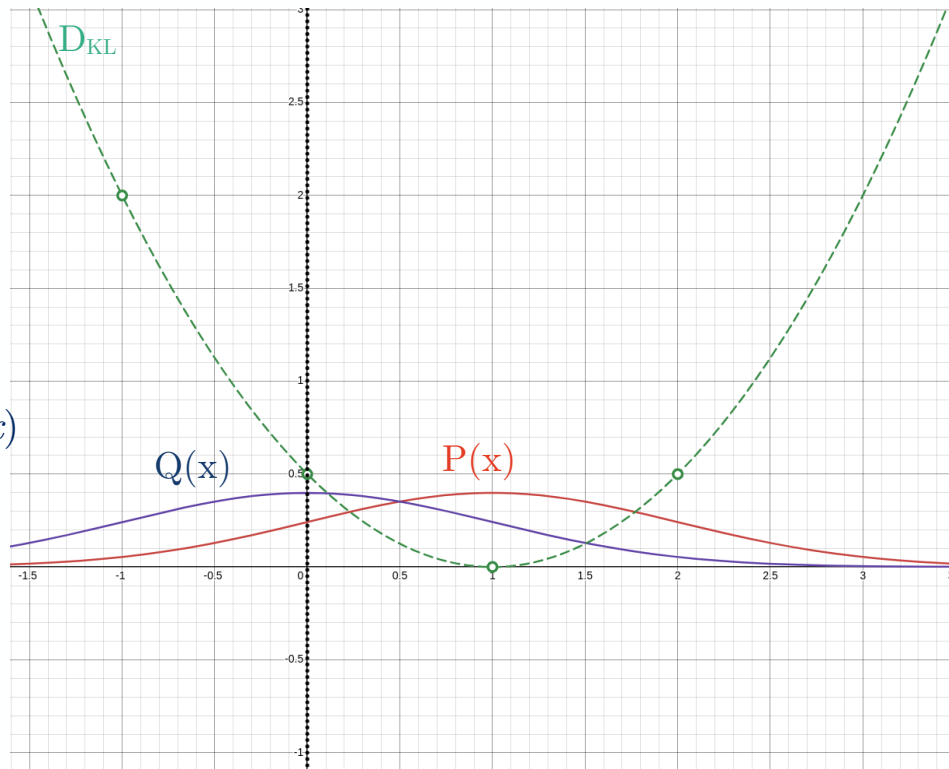
$$I(x) = -\ln(P(x))$$

- **Uncertainty in an entire distribution $P(x)$ – Shannon Entropy:**

$$H(x) = \mathbb{E}_{x \sim P} [I(x)] = -\mathbb{E}_{x \sim P} [\ln(P(x))]$$

- **Relative Entropy** between two distributions $P(x)$ and $Q(x)$:

$$\begin{aligned} \mathcal{D}_{\text{KL}}(P||Q) &= \mathbb{E}_{x \sim P} \left[\ln \left(\frac{P(x)}{Q(x)} \right) \right] \\ &= \mathbb{E}_{x \sim P} [\ln(P(x)) - \ln(Q(x))] \end{aligned}$$



- **Self-Information:**

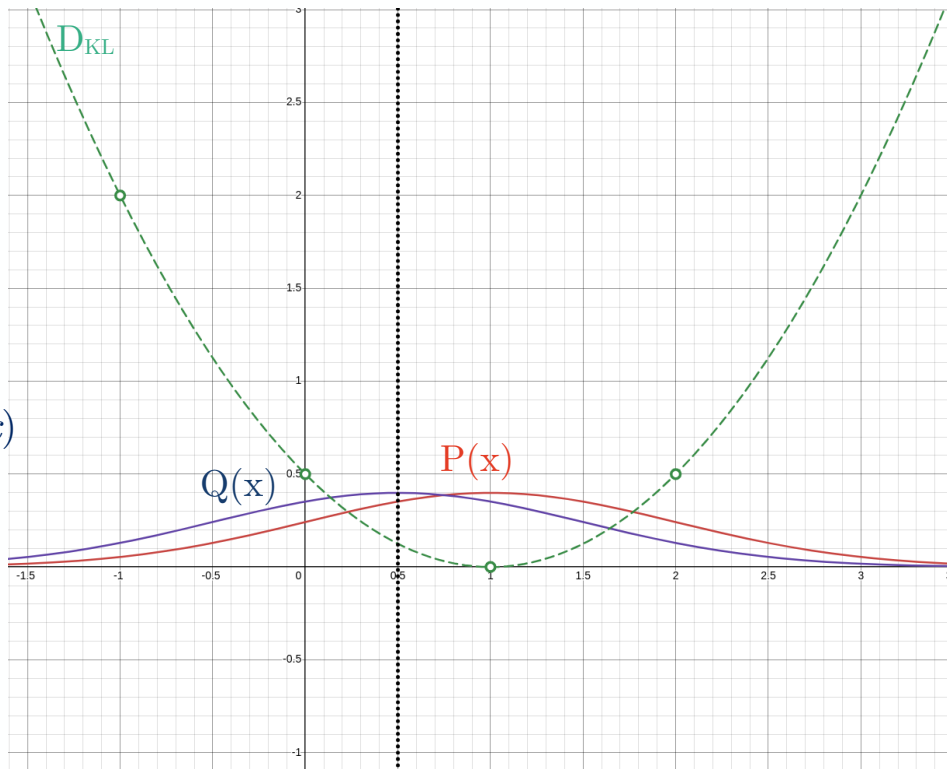
$$I(x) = -\ln(P(x))$$

- **Uncertainty in an entire distribution $P(x)$ – Shannon Entropy:**

$$H(x) = \mathbb{E}_{x \sim P} [I(x)] = -\mathbb{E}_{x \sim P} [\ln(P(x))]$$

- **Relative Entropy** between two distributions $P(x)$ and $Q(x)$:

$$\begin{aligned} \mathcal{D}_{\text{KL}}(P||Q) &= \mathbb{E}_{x \sim P} \left[\ln \left(\frac{P(x)}{Q(x)} \right) \right] \\ &= \mathbb{E}_{x \sim P} [\ln(P(x)) - \ln(Q(x))] \end{aligned}$$



- **Self-Information:**

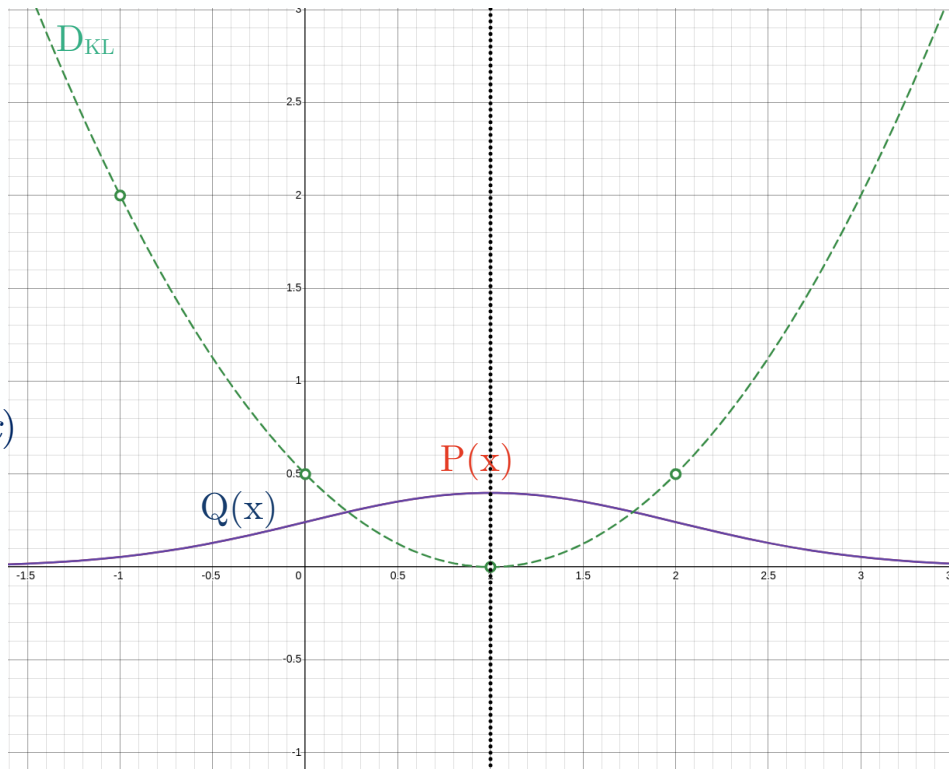
$$I(x) = -\ln(P(x))$$

- **Uncertainty in an entire distribution $P(x)$ – Shannon Entropy:**

$$H(x) = \mathbb{E}_{x \sim P} [I(x)] = -\mathbb{E}_{x \sim P} [\ln(P(x))]$$

- **Relative Entropy** between two distributions $P(x)$ and $Q(x)$:

$$\begin{aligned} \mathcal{D}_{\text{KL}}(P||Q) &= \mathbb{E}_{x \sim P} \left[\ln \left(\frac{P(x)}{Q(x)} \right) \right] \\ &= \mathbb{E}_{x \sim P} [\ln(P(x)) - \ln(Q(x))] \end{aligned}$$



- **Self-Information:**

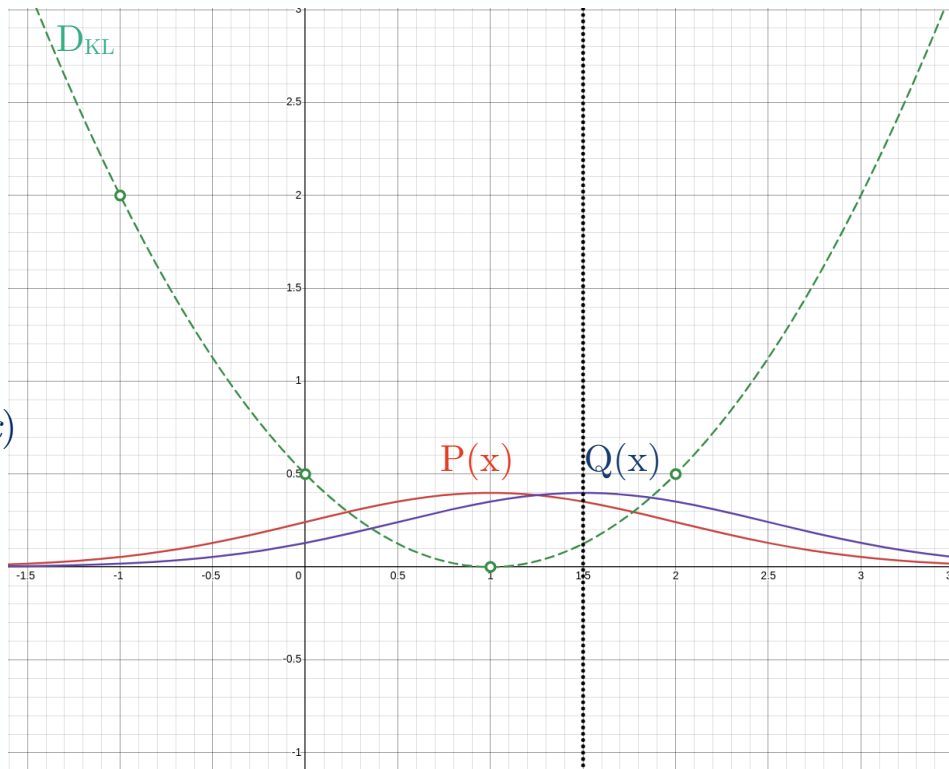
$$I(x) = -\ln(P(x))$$

- **Uncertainty in an entire distribution $P(x)$ – Shannon Entropy:**

$$H(x) = \mathbb{E}_{x \sim P} [I(x)] = -\mathbb{E}_{x \sim P} [\ln(P(x))]$$

- **Relative Entropy** between two distributions $P(x)$ and $Q(x)$:

$$\begin{aligned} \mathcal{D}_{\text{KL}}(P||Q) &= \mathbb{E}_{x \sim P} \left[\ln \left(\frac{P(x)}{Q(x)} \right) \right] \\ &= \mathbb{E}_{x \sim P} [\ln(P(x)) - \ln(Q(x))] \end{aligned}$$



Information Theory

- **Self-Information:**

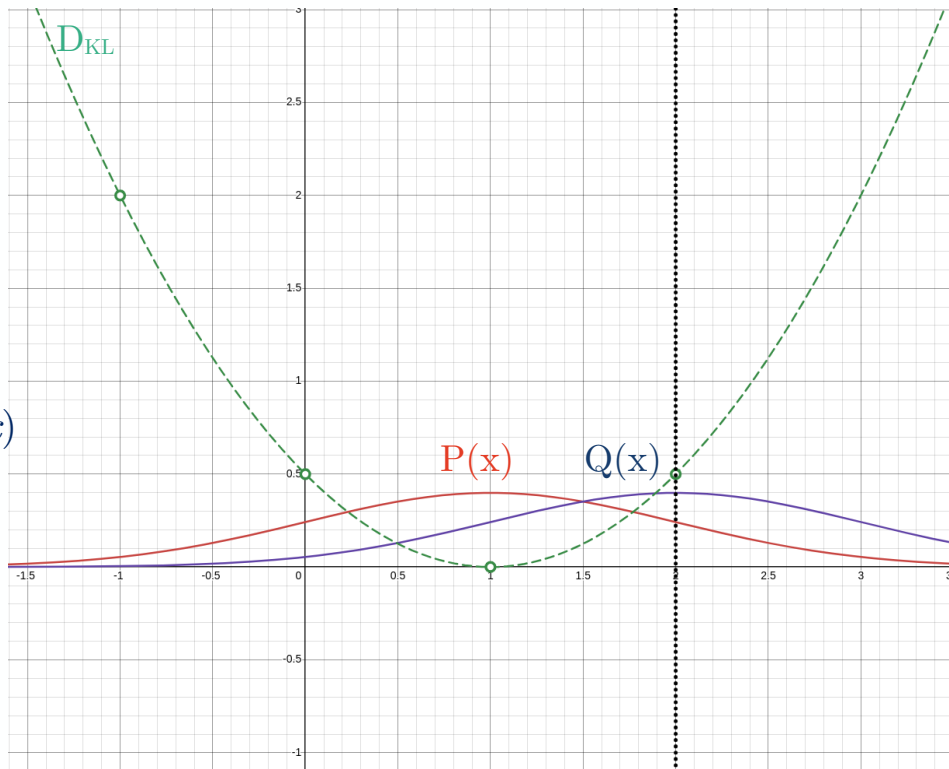
$$I(x) = -\ln(P(x))$$

- **Uncertainty in an entire distribution $P(x)$ – Shannon Entropy:**

$$H(x) = \mathbb{E}_{x \sim P} [I(x)] = -\mathbb{E}_{x \sim P} [\ln(P(x))]$$

- **Relative Entropy** between two distributions $P(x)$ and $Q(x)$:

$$\begin{aligned} \mathcal{D}_{\text{KL}}(P||Q) &= \mathbb{E}_{x \sim P} \left[\ln \left(\frac{P(x)}{Q(x)} \right) \right] \\ &= \mathbb{E}_{x \sim P} [\ln(P(x)) - \ln(Q(x))] \end{aligned}$$



Information Theory

- **Self-Information:**

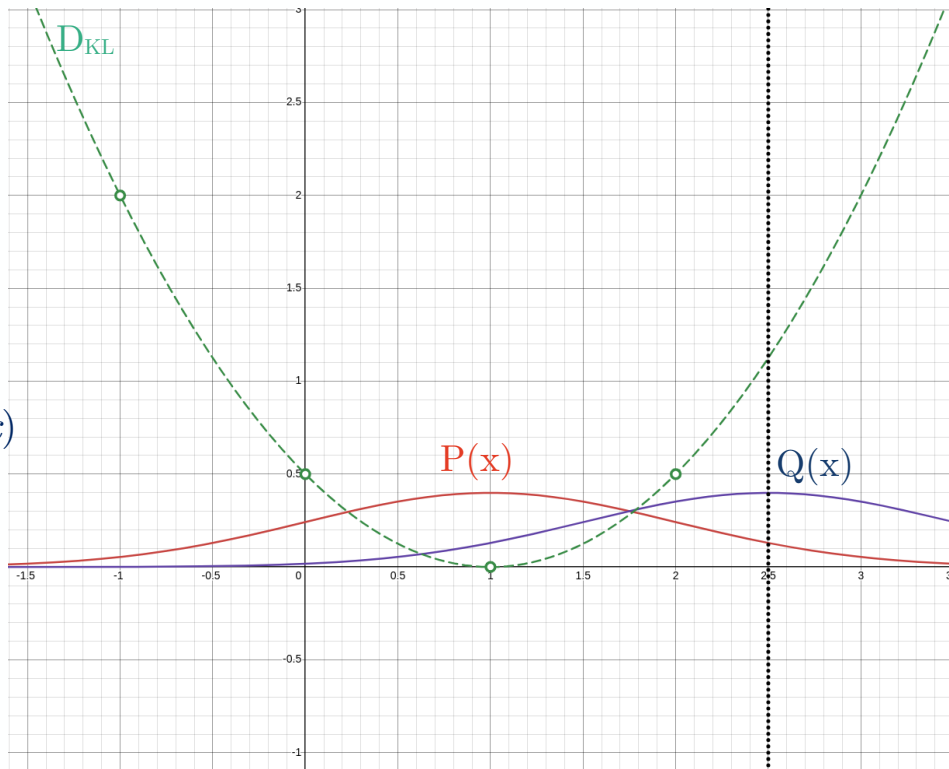
$$I(x) = -\ln(P(x))$$

- **Uncertainty in an entire distribution $P(x)$ – Shannon Entropy:**

$$H(x) = \mathbb{E}_{x \sim P} [I(x)] = -\mathbb{E}_{x \sim P} [\ln(P(x))]$$

- **Relative Entropy** between two distributions $P(x)$ and $Q(x)$:

$$\begin{aligned} \mathcal{D}_{\text{KL}}(P||Q) &= \mathbb{E}_{x \sim P} \left[\ln \left(\frac{P(x)}{Q(x)} \right) \right] \\ &= \mathbb{E}_{x \sim P} [\ln(P(x)) - \ln(Q(x))] \end{aligned}$$



Information Theory

- **Self-Information:**

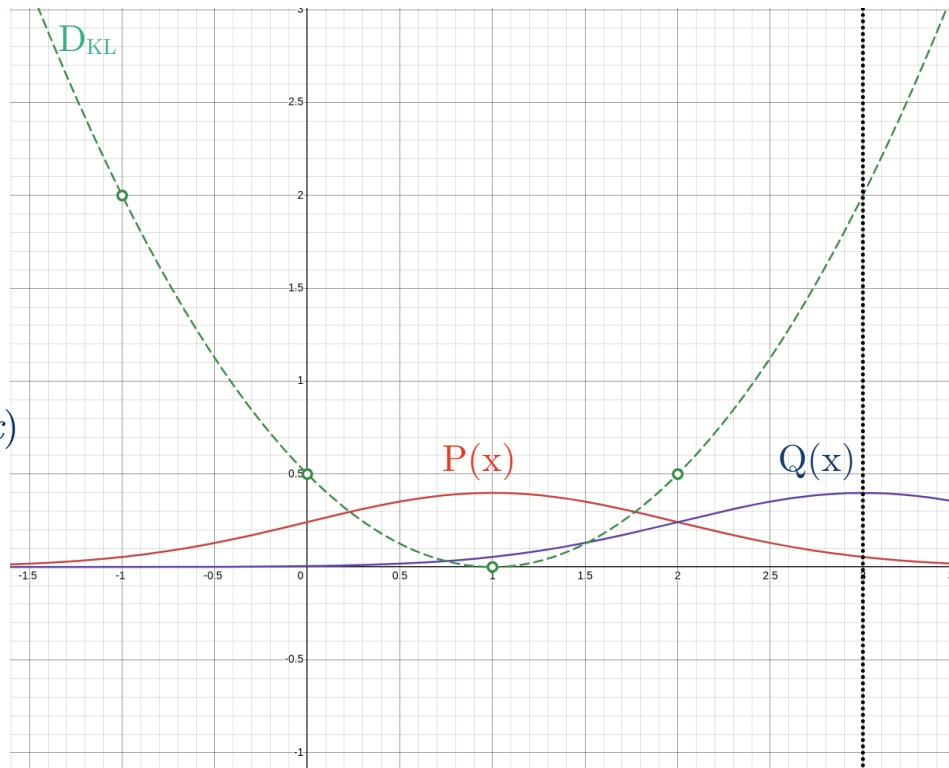
$$I(x) = -\ln(P(x))$$

- **Uncertainty in an entire distribution $P(x)$ – Shannon Entropy:**

$$H(x) = \mathbb{E}_{x \sim P} [I(x)] = -\mathbb{E}_{x \sim P} [\ln(P(x))]$$

- **Relative Entropy** between two distributions $P(x)$ and $Q(x)$:

$$\begin{aligned} \mathcal{D}_{\text{KL}}(P||Q) &= \mathbb{E}_{x \sim P} \left[\ln \left(\frac{P(x)}{Q(x)} \right) \right] \\ &= \mathbb{E}_{x \sim P} [\ln(P(x)) - \ln(Q(x))] \end{aligned}$$



Maximum Likelihood Estimator

- **Intuitive understanding** of the maximum likelihood comes from the idea of *distance* between two distribution

$$\hat{\theta} = \arg \max_{\theta} \mathbb{E}_{x \sim p_{\text{data}}} [\log(p_{\text{model}}(x|\theta))]$$

- **Kullback-Leibler Divergence** is equivalent to maximum likelihood – lets see how this works:

$$\mathcal{D}_{\text{KL}}(P||Q) = \mathbb{E}_{x \sim P} [\ln (P(x)) - \ln (Q(x))]$$

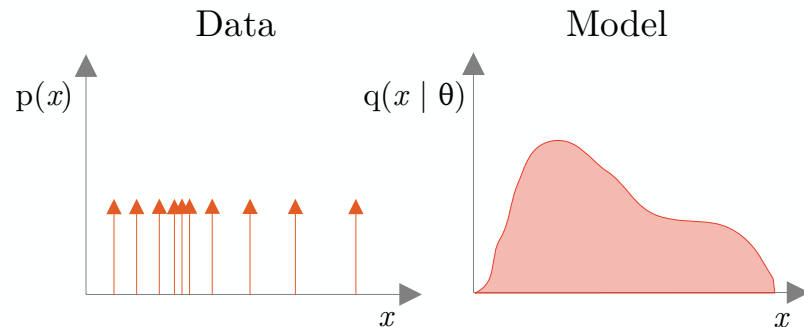
Maximum Likelihood Estimator

- **Intuitive understanding** of the maximum likelihood comes from the idea of *distance* between two distribution

$$\hat{\theta} = \arg \max_{\theta} \mathbb{E}_{x \sim p_{\text{data}}} [\log(p_{\text{model}}(x|\theta))]$$

- **Kullback-Leibler Divergence** is equivalent to maximum likelihood – lets see how this works:

$$\mathcal{D}_{\text{KL}}(P||Q) = \mathbb{E}_{x \sim P} [\ln (P(x)) - \ln (Q(x))]$$



Maximum Likelihood Estimator

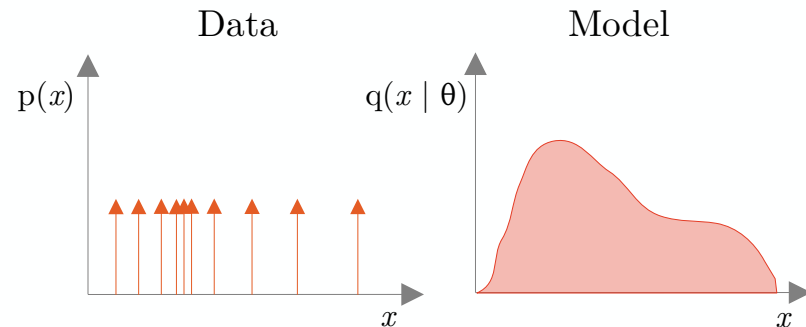
- **Intuitive understanding** of the maximum likelihood comes from the idea of *distance* between two distribution

$$\hat{\theta} = \arg \max_{\theta} \mathbb{E}_{x \sim p_{\text{data}}} [\log(p_{\text{model}}(x|\theta))]$$

- **Kullback-Leibler Divergence** is equivalent to maximum likelihood – lets see how this works:

$$\mathcal{D}_{\text{KL}}(P||Q) = \mathbb{E}_{x \sim P} [\ln (P(x)) - \ln (Q(x))]$$

$$\xrightarrow{1} \hat{\theta} = \arg \min \left[\int p(x) \log(p(x)) dx - \int p(x) \log(q(x|\theta)) dx \right]$$



Maximum Likelihood Estimator

- **Intuitive understanding** of the maximum likelihood comes from the idea of *distance* between two distribution

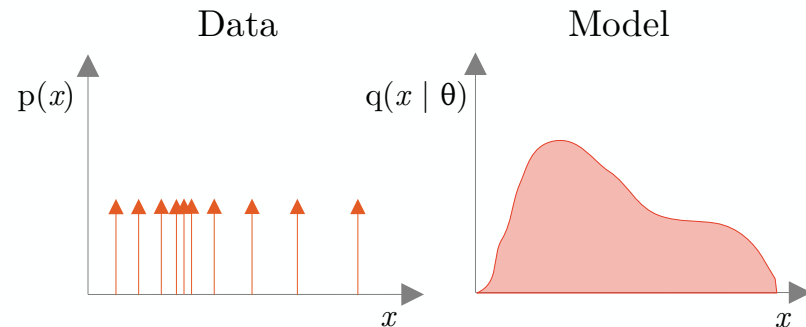
$$\hat{\theta} = \arg \max_{\theta} \mathbb{E}_{x \sim p_{\text{data}}} [\log(p_{\text{model}}(x|\theta))]$$

- **Kullback-Leibler Divergence** is equivalent to maximum likelihood – lets see how this works:

$$\mathcal{D}_{\text{KL}}(P||Q) = \mathbb{E}_{x \sim P} [\ln (P(x)) - \ln (Q(x))]$$

$$\xrightarrow{1} \hat{\theta} = \arg \min \left[\int p(x) \log(p(x)) dx - \int p(x) \log(q(x|\theta)) dx \right]$$

$$\xrightarrow{2} \hat{\theta} = \arg \min \left[- \int p(x) \log(q(x|\theta)) dx \right]$$



Maximum Likelihood Estimator

- Intuitive understanding** of the maximum likelihood comes from the idea of *distance* between two distribution

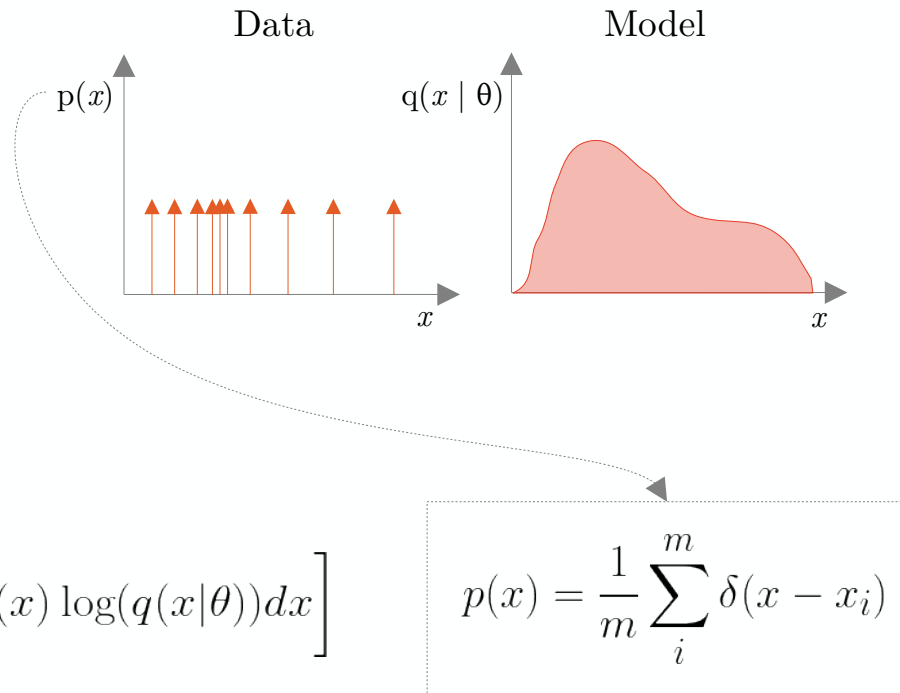
$$\hat{\theta} = \arg \max_{\theta} \mathbb{E}_{x \sim p_{\text{data}}} [\log(p_{\text{model}}(x|\theta))]$$

- Kullback-Leibler Divergence** is equivalent to maximum likelihood – lets see how this works:

$$\mathcal{D}_{\text{KL}}(P||Q) = \mathbb{E}_{x \sim P} [\ln (P(x)) - \ln (Q(x))]$$

1 →
$$\hat{\theta} = \arg \min \left[\int p(x) \log(p(x)) dx - \int p(x) \log(q(x|\theta)) dx \right]$$

2 →
$$\hat{\theta} = \arg \min \left[- \int p(x) \log(q(x|\theta)) dx \right]$$



Maximum Likelihood Estimator

- Intuitive understanding** of the maximum likelihood comes from the idea of *distance* between two distribution

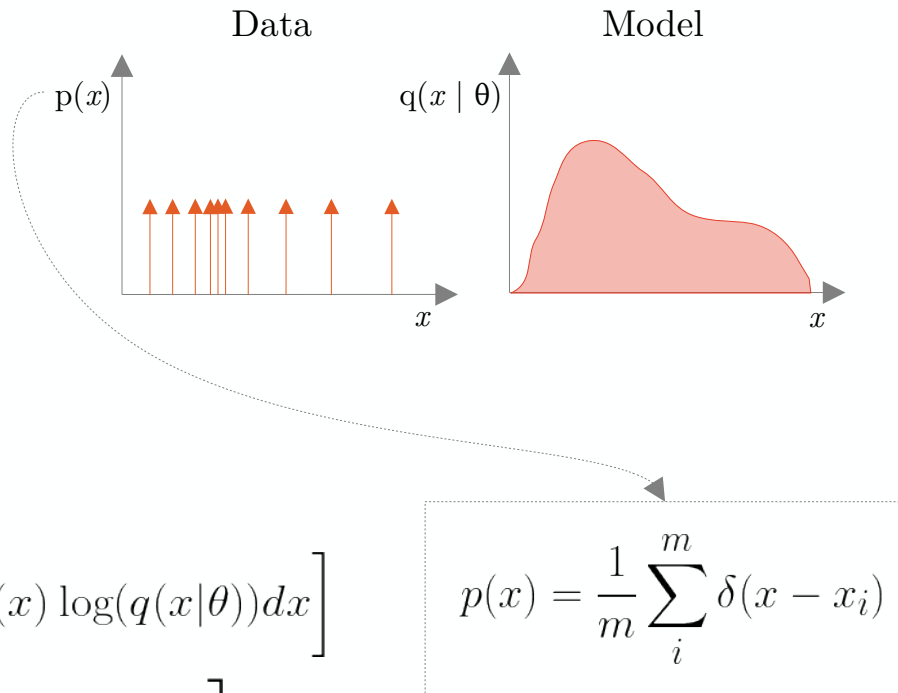
$$\hat{\theta} = \arg \max_{\theta} \mathbb{E}_{x \sim p_{\text{data}}} [\log(p_{\text{model}}(x|\theta))]$$

- Kullback-Leibler Divergence** is equivalent to maximum likelihood – lets see how this works:

$$\mathcal{D}_{\text{KL}}(P||Q) = \mathbb{E}_{x \sim P} [\ln (P(x)) - \ln (Q(x))]$$

1 →
$$\hat{\theta} = \arg \min \left[\int p(x) \log(p(x)) dx - \int p(x) \log(q(x|\theta)) dx \right]$$

3 →
$$\hat{\theta} = \arg \min \left[- \int \left[\frac{1}{m} \sum_i \delta(x - x_i) \right] \log(q(x|\theta)) dx \right]$$



Maximum Likelihood Estimator

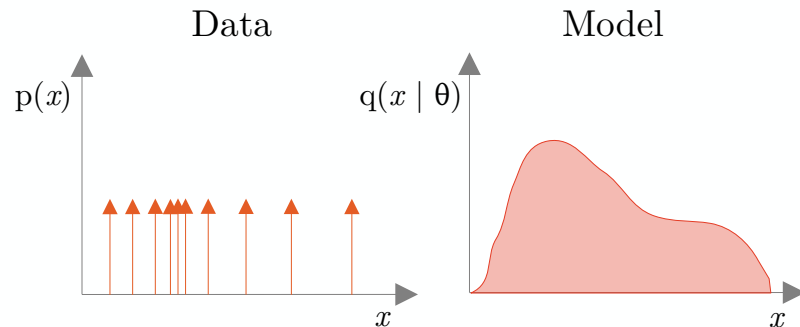
- **Intuitive understanding** of the maximum likelihood comes from the idea of *distance* between two distribution

$$\hat{\theta} = \arg \max_{\theta} \mathbb{E}_{x \sim p_{\text{data}}} [\log(p_{\text{model}}(x|\theta))]$$

- **Kullback-Leibler Divergence** is equivalent to maximum likelihood – lets see how this works:

$$\mathcal{D}_{\text{KL}}(P||Q) = \mathbb{E}_{x \sim P} [\ln (P(x)) - \ln (Q(x))]$$

$$\xrightarrow{4} \hat{\theta} = \arg \min \left[-\frac{1}{m} \sum_i^m \log(q(x_i|\theta)) \right]$$



Maximum Likelihood Estimator

- Intuitive understanding** of the maximum likelihood comes from the idea of *distance* between two distribution

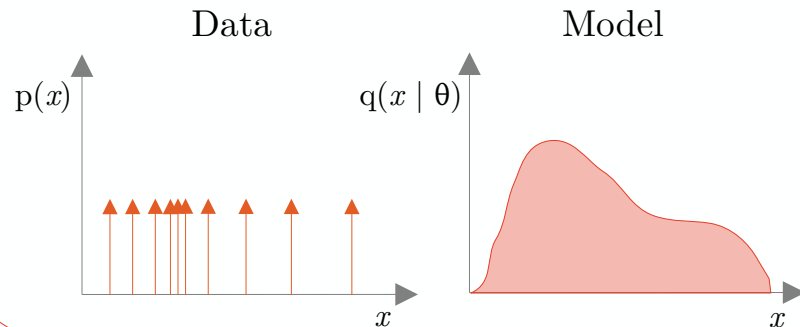
$$\hat{\theta} = \arg \max_{\theta} \mathbb{E}_{x \sim p_{\text{data}}} [\log(p_{\text{model}}(x|\theta))]$$

- Kullback-Leibler Divergence** is equivalent to maximum likelihood – lets see how this works:

$$\mathcal{D}_{\text{KL}}(P||Q) = \mathbb{E}_{x \sim P} [\ln (P(x)) - \ln (Q(x))]$$

$$\xrightarrow{4} \hat{\theta} = \arg \min \left[-\frac{1}{m} \sum_i^m \log(q(x_i|\theta)) \right]$$

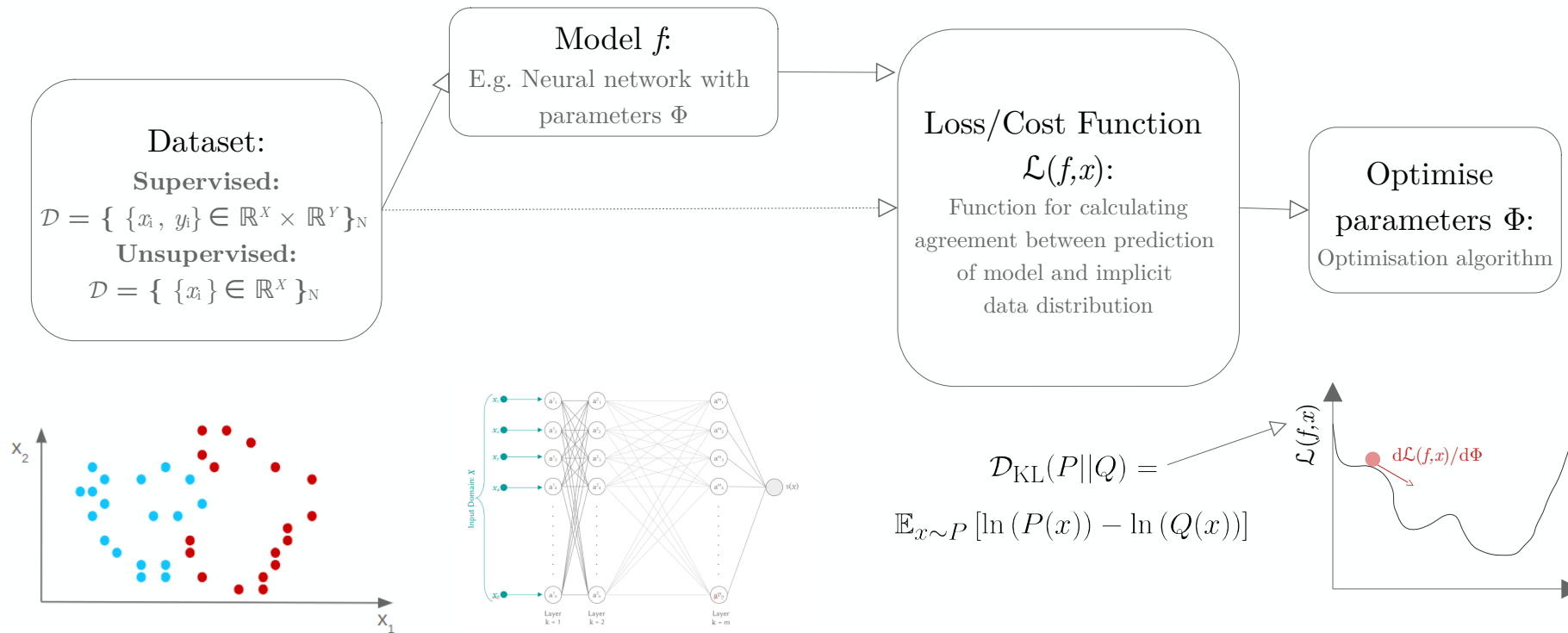
$$\xrightarrow{5} \hat{\theta} = \arg \min \mathbb{E}_{x \sim p} [-\log(q(x|\theta))]$$





Learning Process & Estimators

- Four key parts to the ‘Machine Learning’ process:



Supervised vs Unsupervised

61

Dataset:

Supervised:

$$\mathcal{D} = \{ \{x_i, y_i\} \in \mathbb{R}^X \times \mathbb{R}^Y \}_N$$

Unsupervised:

$$\mathcal{D} = \{ \{x_i\} \in \mathbb{R}^X \}_N$$

- **Supervised** learning:

Examples $\mathbf{x} = \{x\}_m$ paired with targets $y = \{y\}_m$
that instruct the algorithm on what to learn

$$p(y|x)$$

- **Unsupervised** learning:

Examples $\mathbf{x} = \{x\}_m$ with no targets

$$p(x)$$

Supervised vs Unsupervised

62

Dataset:
Supervised:
 $\mathcal{D} = \{ \{x_i, y_i\} \in \mathbb{R}^X \times \mathbb{R}^Y \}_N$
Unsupervised:
 $\mathcal{D} = \{ \{x_i\} \in \mathbb{R}^X \}_N$

- **Supervised** learning:

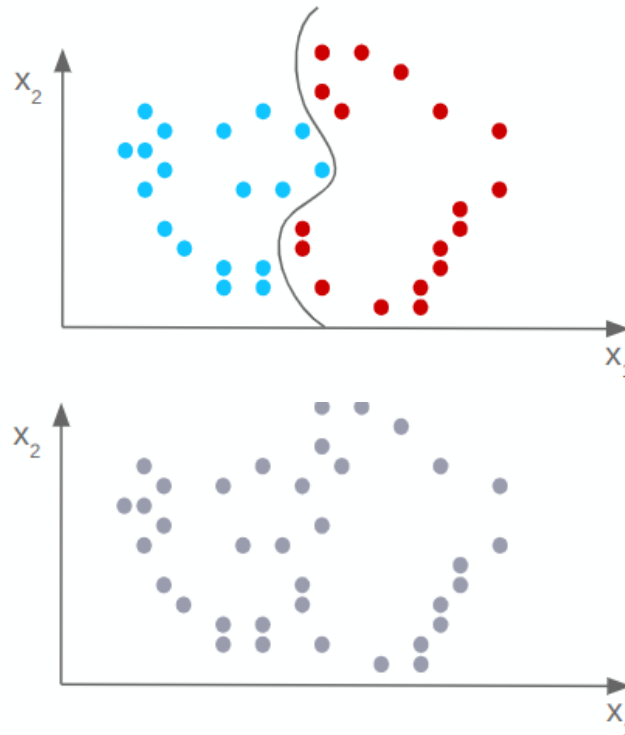
Examples $\mathbf{x} = \{x\}_m$ paired with targets $y = \{y\}_m$
that instruct the algorithm on what to learn

$$p(y|x) \quad \begin{array}{l} \text{Conditional Density} \\ \text{Estimation} \end{array}$$

- **Unsupervised** learning:

Examples $\mathbf{x} = \{x\}_m$ with no targets

$$p(x) \quad \text{Density Estimation}$$



Supervised vs Unsupervised

63

Dataset:
Supervised:
 $= \{ \{x_i, y_i\} \in \mathbb{R}^x \times \mathbb{R}^y \}_N$
Unsupervised:
 $= \{ \{x_i\} \in \mathbb{R}^x \}_N$

- **Supervised learning:**

Examples $\mathbf{x} = \{x\}_m$ paired with targets $y = \{y\}_m$
that instruct the algorithm on what to learn

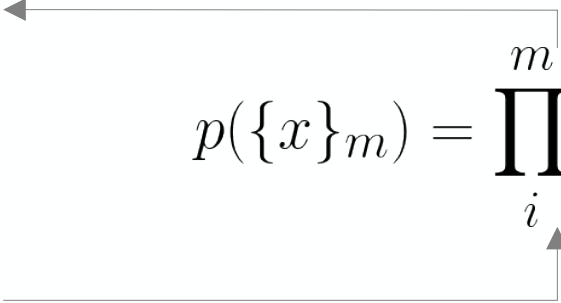
$p(y|x)$ Conditional Density
Estimation

Decomposing the problem into $m-1$ conditional density
estimation problems ~ supervised learning problems

- **Unsupervised learning:**

Examples $\mathbf{x} = \{x\}_m$ with no targets

$p(x)$ Density Estimation

$$p(\{x\}_m) = \prod_i^m p(x_i | \{x\}_{m-i})$$


Supervised vs Unsupervised

64

Dataset:
Supervised:
 $= \{ \{x_i, y_i\} \in \mathbb{R}^x \times \mathbb{R}^y \}_N$
Unsupervised:
 $= \{ \{x_i\} \in \mathbb{R}^x \}_N$

- **Supervised learning:**

Examples $\mathbf{x} = \{x\}_m$ paired with targets $y = \{y\}_m$
that instruct the algorithm on what to learn

$p(y|x)$ Conditional Density
Estimation

- **Unsupervised learning:**

Examples $\mathbf{x} = \{x\}_m$ with no targets

$p(x)$ Density Estimation

$$p(\{y\}_m | \{x\}_m) = \frac{p(\{x\}_m, \{y\}_m)}{\sum_{\tilde{y}} p(\{x\}_m, \{\tilde{y}\}_m)}$$

Learning the joint distribution $p(x,y)$ via unsupervised learning techniques, then infer.

Supervised vs Unsupervised

65

Dataset:

Supervised:

$$\mathcal{D} = \{ \{x_i, y_i\} \in \mathbb{R}^X \times \mathbb{R}^Y \}_N$$

Unsupervised:

$$\mathcal{D} = \{ \{x_i\} \in \mathbb{R}^X \}_N$$

- **Supervised** learning:

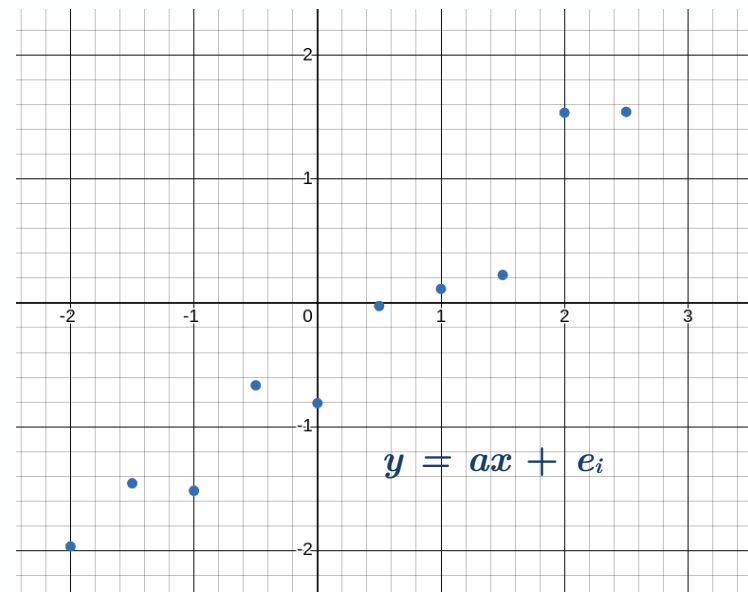
Examples $\mathbf{x} = \{x\}_m$ paired with targets $y = \{y\}_m$
that instruct the algorithm on what to learn

$$p(y|x)$$

- **Unsupervised** learning:

Examples $\mathbf{x} = \{x\}_m$ with no targets

$$p(x)$$



Linear Regression Example

Model

Model f :
E.g. Neural network with
parameters Φ

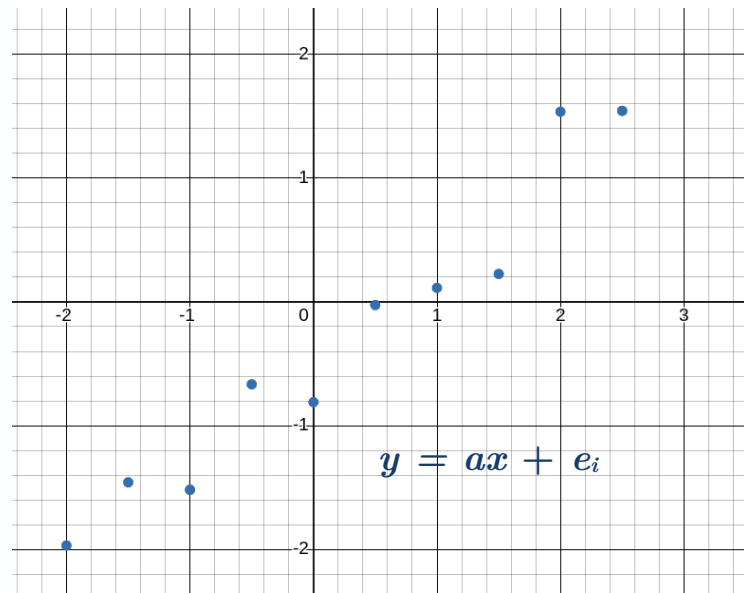
66

- Parametric linear model:

$$\hat{y} = \mathbf{w}^T \mathbf{x}$$

- Configurable parameters of model $\Phi = \mathbf{w}$

$$y = \begin{pmatrix} w_1, \dots, w_n \end{pmatrix} \begin{pmatrix} x_1, \\ \cdot \\ \cdot \\ x_n \end{pmatrix}$$



Linear Regression Example

Performance Measure: Loss Function

Loss/Cost Function

67

$\mathcal{L}(f, x)$:

Function for calculating
agreement between prediction
of model and implicit
data distribution

- **Parametric linear model:**

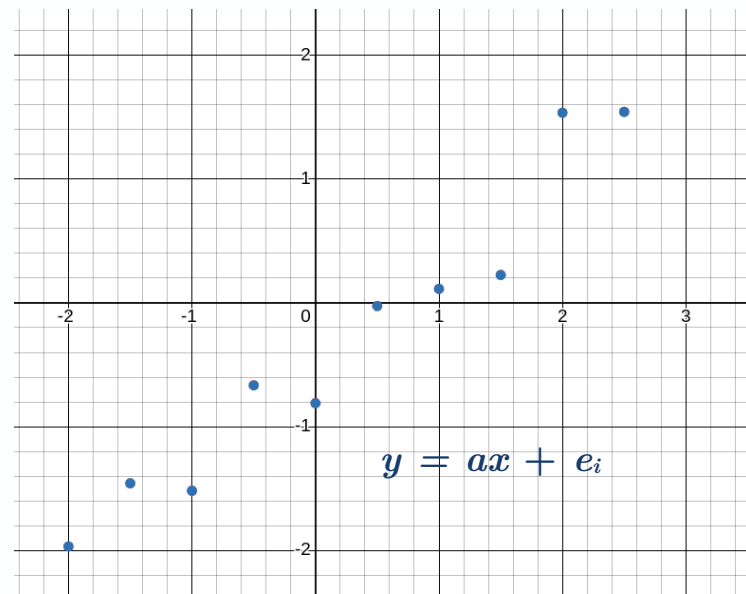
$$\hat{y} = \mathbf{w}^T \mathbf{x}$$

- **Configurable parameters of model** $\Phi = \mathbf{w}$

$$y = \begin{pmatrix} w_1, \dots, w_n \end{pmatrix} \begin{pmatrix} x_1, \\ \vdots \\ x_n \end{pmatrix}$$

- **Loss:** Mean squared error

$$\mathcal{L} = \frac{1}{m} \sum_i (\hat{y} - y)^2$$



Linear Regression Example

Performance Measure: Loss Function

Loss/Cost Function

68

$\mathcal{L}(f, x)$:

Function for calculating
agreement between prediction
of model and implicit
data distribution

- **Parametric linear model:**

$$\hat{y} = \mathbf{w}^T \mathbf{x}$$

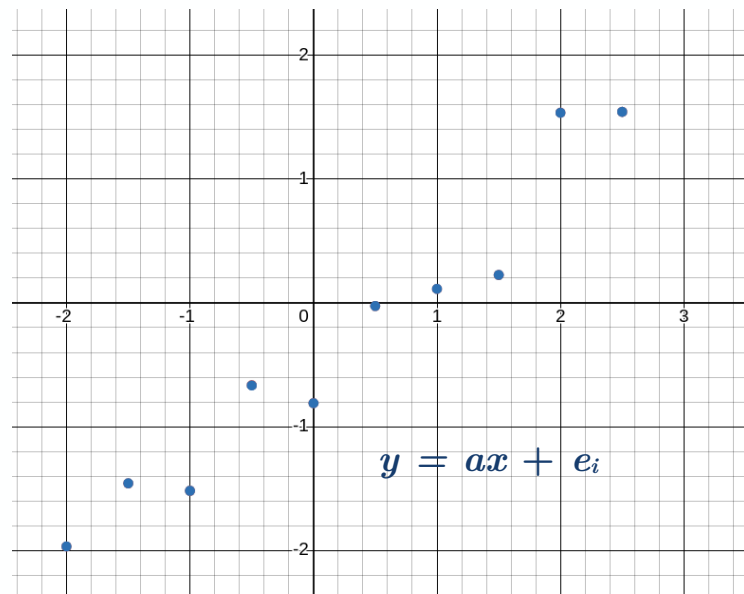
- **Configurable parameters of model $\Phi = \mathbf{w}$**

$$y = \begin{pmatrix} w_1, \dots, w_n \end{pmatrix} \begin{pmatrix} x_1, \\ \vdots \\ x_n \end{pmatrix}$$

- **Loss: Mean squared error**

$$\mathcal{L} = \frac{1}{m} \sum_i (\hat{y} - y)^2$$

Why the mean
squared error?



Linear Regression Example

Maximum Likelihood Approach

69

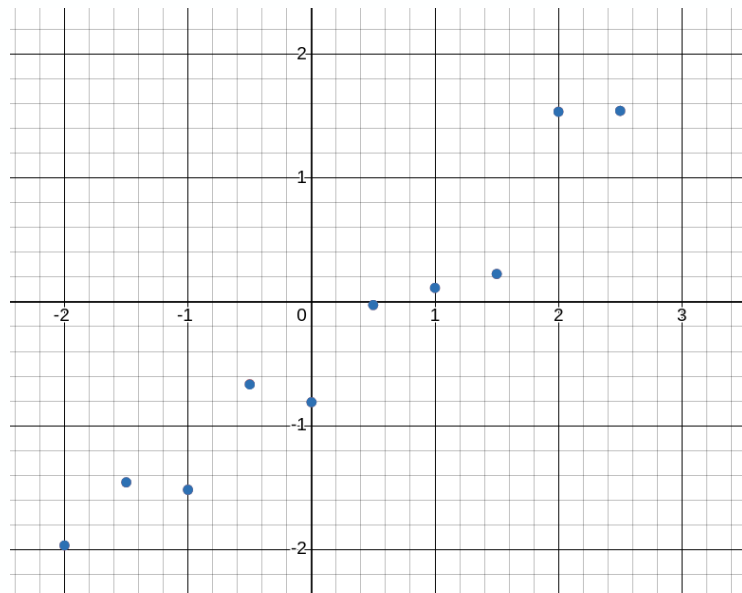
Optimise
parameters Φ :
Optimisation algorithm

- I said there was a general likelihood approach to *estimating/learning* parameters:

$$\hat{\theta} = \arg \max_{\theta} \sum_i^m \ln p(x_i | \theta)$$

- General MLE** estimate also applies to conditional probability:

$$\hat{\theta} = \arg \max_{\theta} \sum_i^m \ln (p(y_i | x_i, \theta))$$



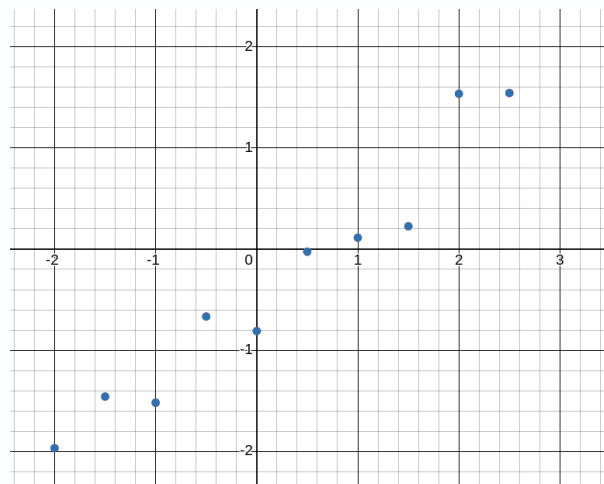
Linear Regression Example

Maximum Likelihood Approach

70

Optimise
parameters Φ :
Optimisation algorithm

Linear Regression Example



$$y_i = \alpha x_i + e_i \quad \longrightarrow \quad y_i \sim \mathcal{N}(\alpha x_i, \sigma)$$

$e_i \sim \mathcal{N}(0, \sigma)$

Maximum Likelihood Approach

71

Optimise
parameters Φ :
Optimisation algorithm

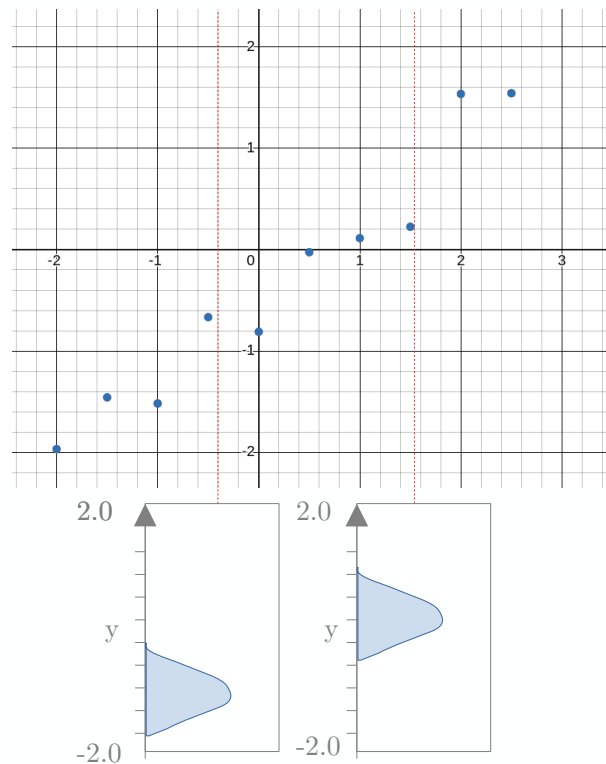
Linear Regression Example

- Probabilistic view of regression:**

Data generated with an error term that is gaussian distributed, meaning repeated experiments have some noise:

$$y_i = \alpha x_i + e_i \quad \longrightarrow \quad y_i \sim \mathcal{N}(\alpha x_i, \sigma)$$

$e_i \sim \mathcal{N}(0, \sigma)$



Maximum Likelihood Approach

72

Optimise
parameters Φ :
Optimisation algorithm

Probabilistic view of regression:

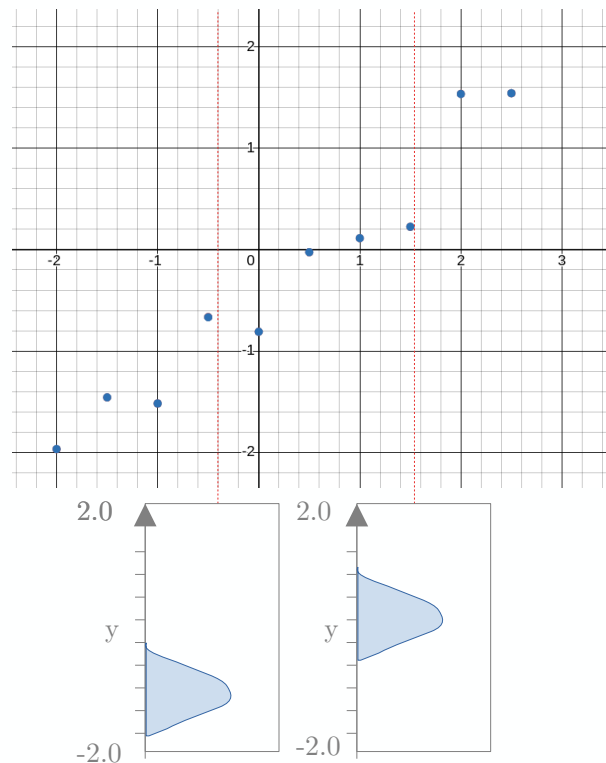
Data generated with an error term that is gaussian distributed, meaning repeated experiments have some noise:

$$y_i = \alpha x_i + e_i \quad \longrightarrow \quad y_i \sim \mathcal{N}(\alpha x_i, \sigma)$$

$$e_i \sim \mathcal{N}(0, \sigma)$$

$$p(y_i | x_i, \alpha) = \mathcal{N}(y_i | \hat{y}_i, \sigma) = \mathcal{N}(y_i | \alpha x_i, \sigma)$$

Linear Regression Example



Maximum Likelihood Approach

73

Optimise
parameters Φ :
Optimisation algorithm

- Generalise MLE estimate the conditional probability:

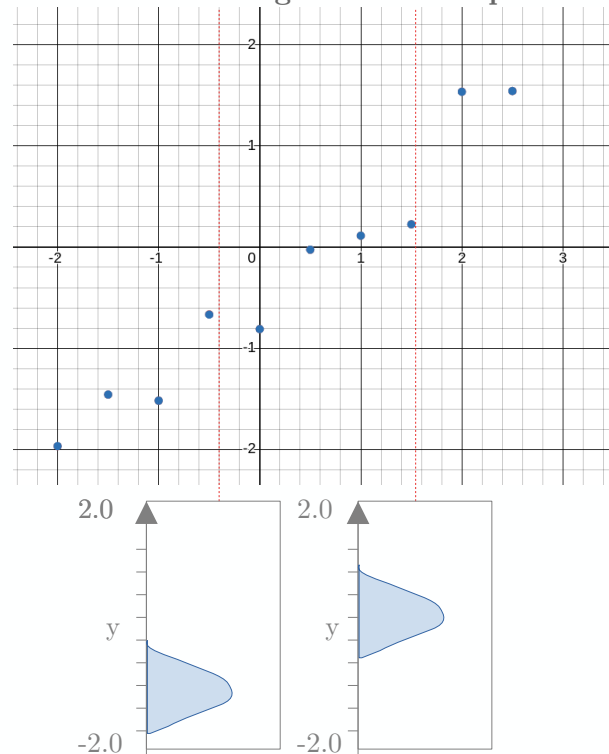
$$\hat{\theta} = \arg \max_{\theta} \sum_i^m \ln (p(y_i|x_i, \theta))$$

- Probabilistic view of regression:

$$\hat{\alpha} = \arg \max_{\theta} \sum_i^m \ln \left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\alpha x_i - \mu)^2}{2\sigma^2}} \right)$$

$$p(y_i|x_i, \alpha) = \mathcal{N}(y_i|\alpha x_i, \sigma)$$

Linear Regression Example



Maximum Likelihood Approach

74

Optimise
parameters Φ :
Optimisation algorithm

- Generalise MLE estimate the conditional probability:

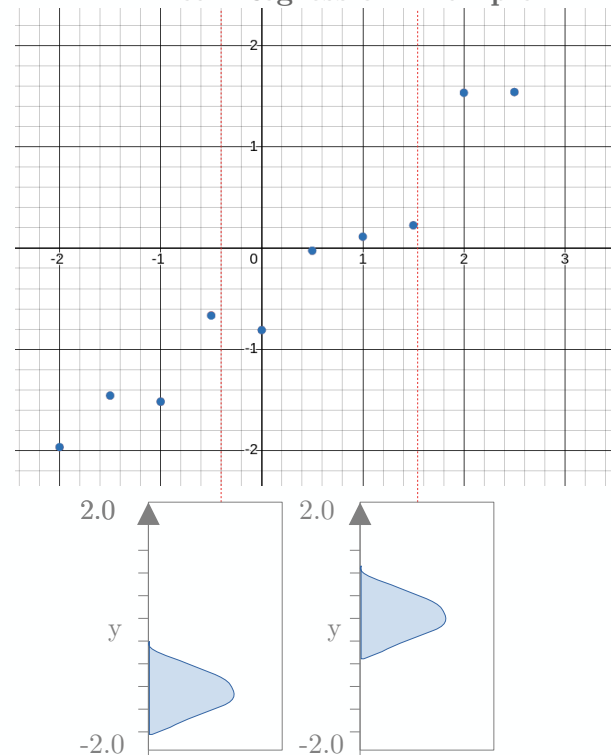
$$\hat{\theta} = \arg \max_{\theta} \sum_i^m \ln (p(y_i|x_i, \theta))$$

- Probabilistic view of regression:

$$\begin{aligned} \hat{\alpha} &= \arg \max_{\theta} \sum_i^m \ln \left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\alpha x_i - \mu)^2}{2\sigma^2}} \right) \\ &= -m \ln(\sigma) - \frac{m}{2} \ln(2\pi) - \sum \frac{(\alpha x_i - \mu)^2}{2\sigma^2} \end{aligned}$$

$$\begin{aligned} p(y_i|x_i, \alpha) \\ = \mathcal{N}(y_i|\alpha x_i, \sigma) \end{aligned}$$

Linear Regression Example



Maximum Likelihood Approach

75

Optimise
parameters Φ :
Optimisation algorithm

- Generalise MLE estimate the conditional probability:

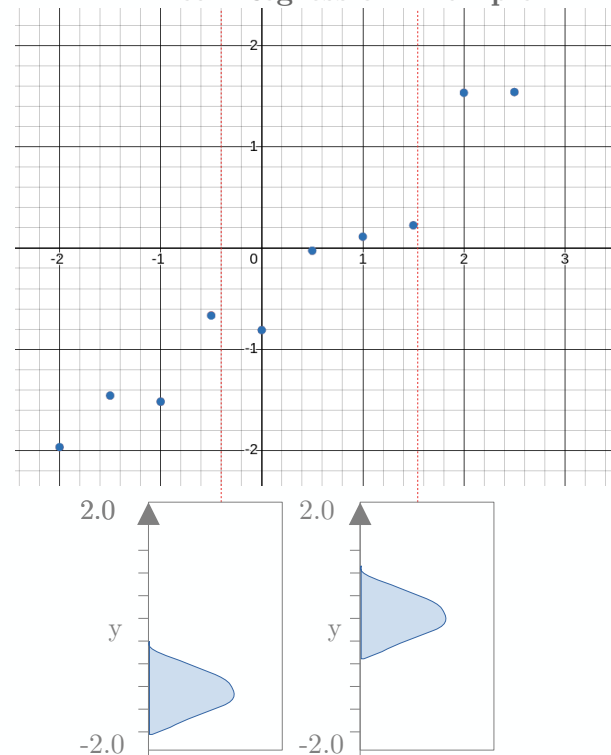
$$\hat{\theta} = \arg \max_{\theta} \sum_i^m \ln (p(y_i|x_i, \theta))$$

- Probabilistic view of regression:

$$\begin{aligned} \hat{\alpha} &= \arg \max_{\theta} \sum_i^m \ln \left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\alpha x_i - \mu)^2}{2\sigma^2}} \right) \\ &= -m \ln(\sigma) - \frac{m}{2} \ln(2\pi) - \sum_i^m \frac{(\alpha x_i - \mu)^2}{2\sigma^2} \end{aligned}$$

- Maximising the log-likelihood is equivalent to Mean Squared Error minimisation

Linear Regression Example



Gradient Descent

76

Optimise
parameters Φ :
Optimisation algorithm

- First Order Gradient Descent:

$$\nabla_w \text{MSE}_{\text{train}} = 0$$

$$\Rightarrow \nabla_w \frac{1}{m} \|\hat{\mathbf{y}}^{(\text{train})} - \mathbf{y}^{(\text{train})}\|_2^2 = 0$$

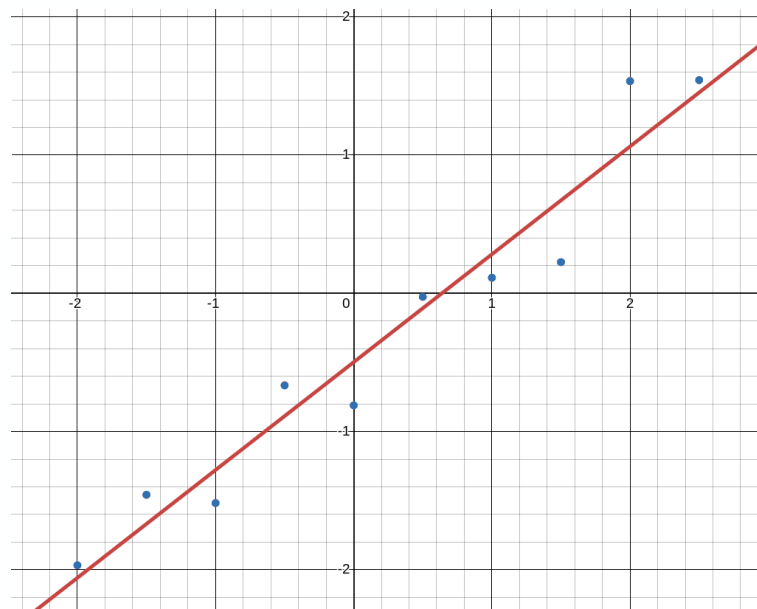
$$\Rightarrow \frac{1}{m} \nabla_w \|\mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})}\|_2^2 = 0$$

$$\Rightarrow \nabla_w \left(\mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})} \right)^\top \left(\mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})} \right) = 0$$

$$\Rightarrow \nabla_w \left(\mathbf{w}^\top \mathbf{X}^{(\text{train})\top} \mathbf{X}^{(\text{train})} \mathbf{w} - 2 \mathbf{w}^\top \mathbf{X}^{(\text{train})\top} \mathbf{y}^{(\text{train})} + \mathbf{y}^{(\text{train})\top} \mathbf{y}^{(\text{train})} \right) = 0$$

$$\Rightarrow 2 \mathbf{X}^{(\text{train})\top} \mathbf{X}^{(\text{train})} \mathbf{w} - 2 \mathbf{X}^{(\text{train})\top} \mathbf{y}^{(\text{train})} = 0$$

$$\Rightarrow \mathbf{w} = \left(\mathbf{X}^{(\text{train})\top} \mathbf{X}^{(\text{train})} \right)^{-1} \mathbf{X}^{(\text{train})\top} \mathbf{y}^{(\text{train})}$$



Linear Regression Example

Capacity & Bias-Variance Trade-Off

77

- **Capacity** of a model \sim the number of degrees of freedom that can be tuned

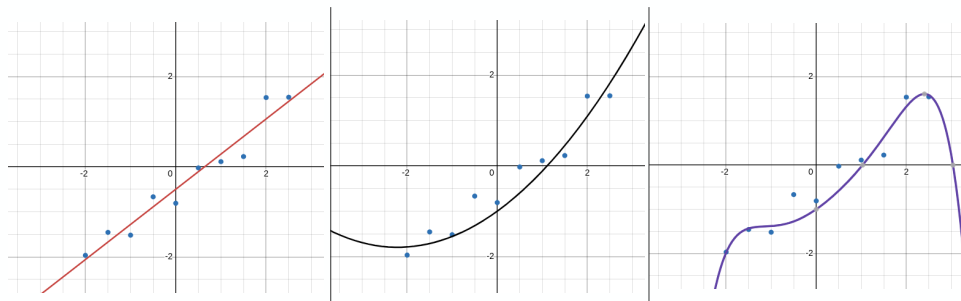
1st Order Polynomial: $\hat{y} = wx + b$

2nd Order Polynomial: $\hat{y} = w_1x + w_2x^2 + b$

•

•

9th Order Polynomial: $\hat{y} = \sum_i^9 w_i x^i + b$



Capacity & Bias-Variance Trade-Off

78

- **Capacity** of a model \sim the number of degrees of freedom that can be tuned

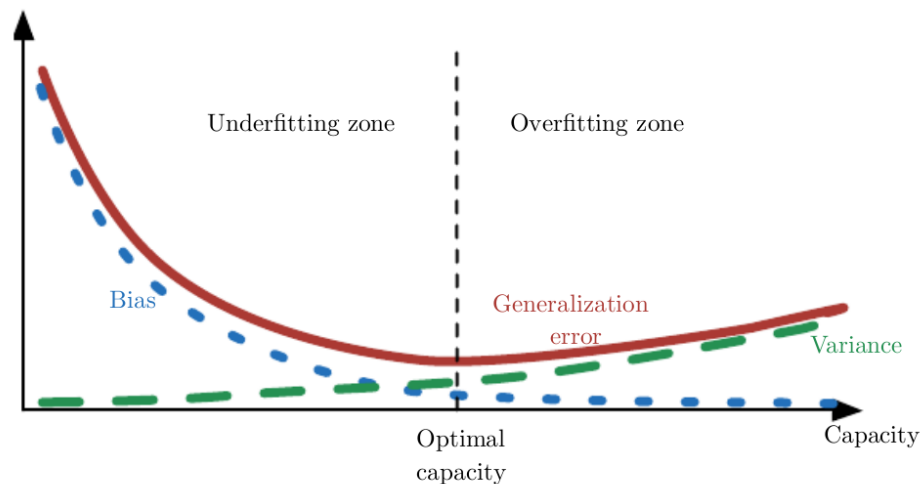
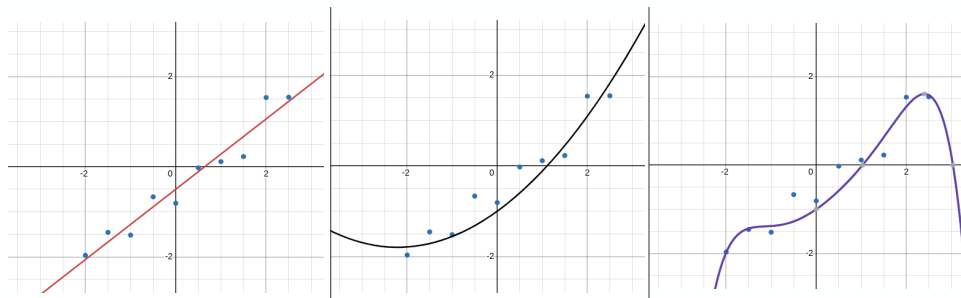
1st Order Polynomial: $\hat{y} = wx + b$

2nd Order Polynomial: $\hat{y} = w_1x + w_2x^2 + b$

.

9th Order Polynomial: $\hat{y} = \sum_i^9 w_i x^i + b$

- **Under/Over-fitting** avoided by matching capacity of model to complexity of the problem



Estimators

- **Point Estimator, or statistic**, is any function of the data that infers from the data some parameter of interest θ :

$$\hat{\theta} = g(\{x\}_m)$$

- **Bias** of an estimator is given by:

$$\text{bias}(\hat{\theta}) = \mathbb{E}(\hat{\theta}) - \theta$$

- **Variance** of an estimator:

$$\text{Var}(\hat{\theta})$$

Estimators

- **Point Estimator, or statistic**, is any function of the data that infers from the data some parameter of interest θ :

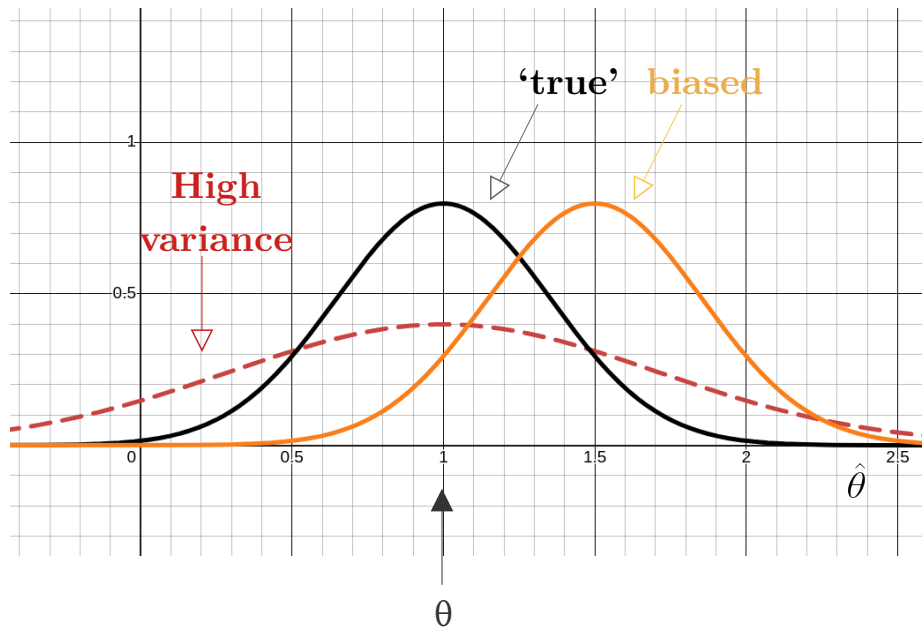
$$\hat{\theta} = g(\{x\}_m)$$

- **Bias** of an estimator is given by:

$$\text{bias}(\hat{\theta}) = \mathbb{E}(\hat{\theta}) - \theta$$

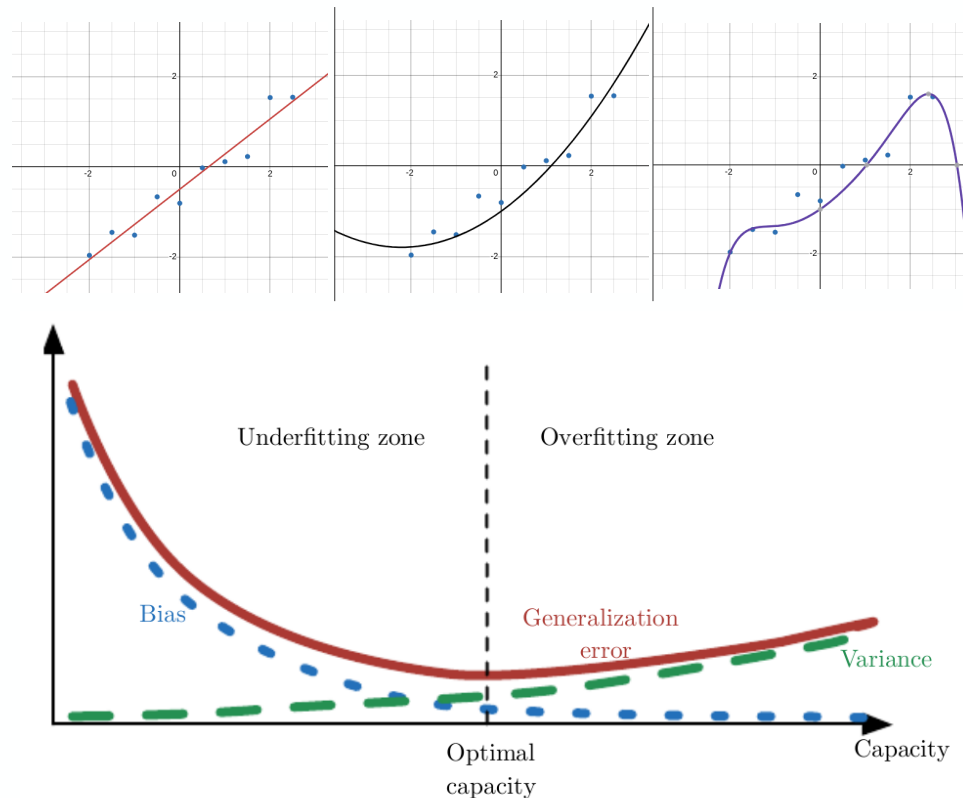
- **Variance** of an estimator:

$$\text{Var}(\hat{\theta})$$



- Bias-Variance trade-off:**

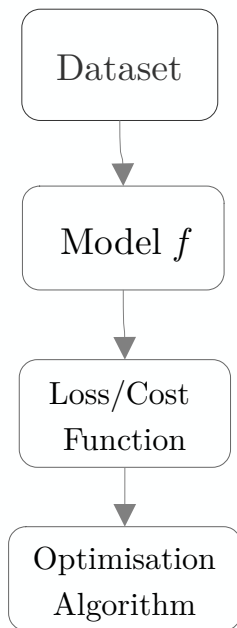
$$\begin{aligned}\text{MSE}(\hat{\theta}, \theta) &= \mathbb{E} \left[(\hat{\theta} - \theta)^2 \right] \\ &= \mathbb{E} \left[\hat{\theta}^2 - 2\hat{\theta}\theta + \theta^2 \right] \\ &= \mathbb{E} \left[\hat{\theta}^2 \right] - 2\mathbb{E} \left[\hat{\theta} \right] \theta + \theta^2 \\ &= \text{Var}(\hat{\theta}) + \mathbb{E} \left[\hat{\theta} \right]^2 - 2\mathbb{E} \left[\hat{\theta} \right] \theta + \theta^2 \\ &= \text{Var}(\hat{\theta}) + \text{bias}(\hat{\theta})^2\end{aligned}$$



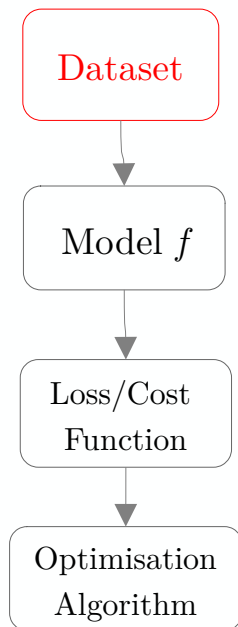


Learning Recipe

General recipe for *most* ML algorithms or learning processes:



General recipe for *most* ML algorithms or learning processes:



Experience

- **Supervised learning:**

Examples $\mathbf{x} = \{x\}_m$ paired with targets $y = \{y\}_m$ that instruct the algorithm on what to learn

$$p(y|x) \quad \begin{array}{l} \text{Conditional Density} \\ \text{Estimation} \end{array}$$

- **Unsupervised learning:**

Examples $\mathbf{x} = \{x\}_m$ with no targets

$$p(x) \quad \text{Density Estimation}$$

Dataset:

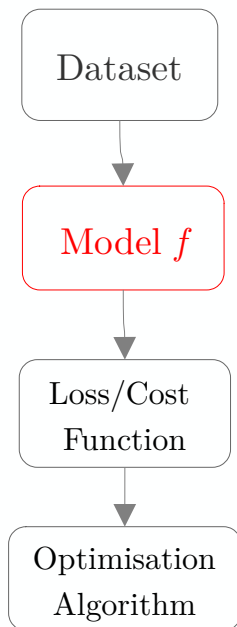
Supervised:

$$\mathcal{D} = \{ \{x_i, y_i\} \in \mathbb{R}^X \times \mathbb{R}^Y \}_N$$

Unsupervised:

$$\mathcal{D} = \{ \{x_i\} \in \mathbb{R}^X \}_N$$

General recipe for *most* ML algorithms or learning processes:



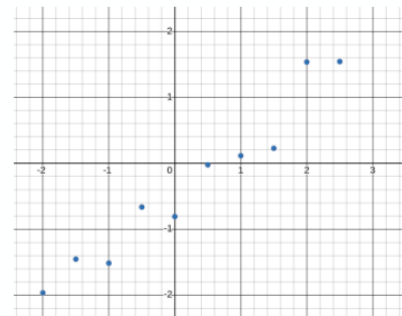
Task

- Parametric linear model:

$$\hat{y} = \mathbf{w}^T \mathbf{x}$$

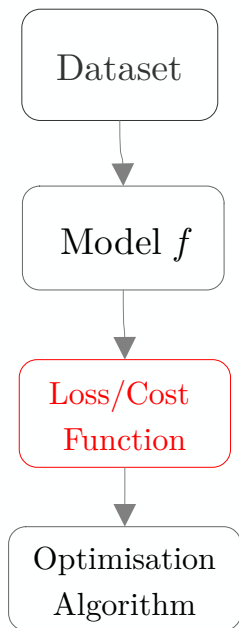
- Configurable parameters of model $\Phi = \mathbf{w}$

$$\mathbf{y} = \begin{pmatrix} w_1, \dots, w_n \end{pmatrix} \quad \begin{pmatrix} x_1, \\ \vdots \\ x_n \end{pmatrix}$$



Linear Regression Example

General recipe for *most* ML algorithms or learning processes:



Performance

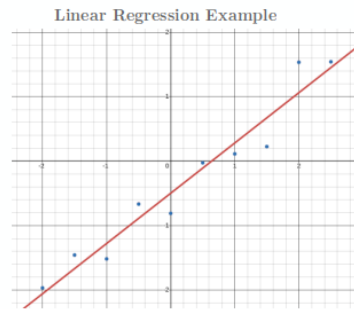
- Generalise MLE estimate the conditional probability:

$$\hat{\theta} = \arg \max_{\theta} \sum_i^m \ln(p(y_i|x_i, \theta))$$

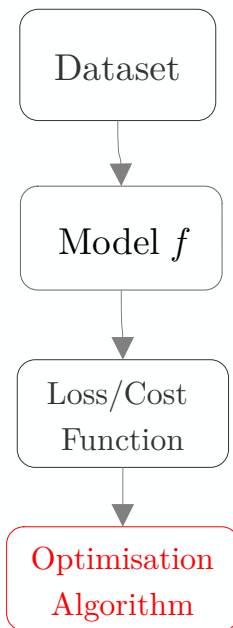
- Probabilistic view of regression:

$$\begin{aligned} \hat{m} &= \arg \max_{\theta} \sum_i^m \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(\alpha x_i - y_i)^2}{2\sigma^2}} \\ &= -m \log(\sigma) - \frac{m}{2} \log(2\pi) - \sum_i^m \frac{(\alpha x_i - \mu)^2}{2\sigma^2} \end{aligned}$$

- Maximising the log-likelihood is equivalent to MSE minimisation



General recipe for *most* ML algorithms or learning processes:



- First Order Gradient Descent:

$$\nabla_w \text{MSE}_{\text{train}} = 0 \quad (5.6)$$

$$\Rightarrow \nabla_w \frac{1}{m} \|\mathbf{y}^{(\text{train})} - \mathbf{y}^{(\text{train})}\|_2^2 = 0 \quad (5.7)$$

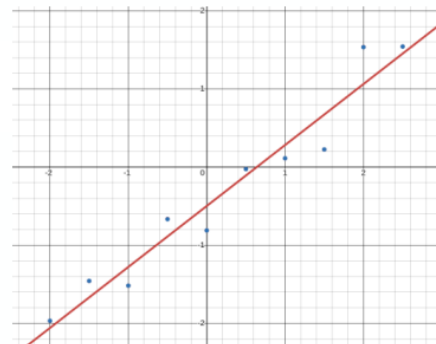
$$\Rightarrow \frac{1}{m} \nabla_w \|\mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})}\|_2^2 = 0 \quad (5.8)$$

$$\Rightarrow \nabla_w \left(\mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})} \right)^\top \left(\mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})} \right) = 0 \quad (5.9)$$

$$\Rightarrow \nabla_w \left(\mathbf{w}^\top \mathbf{X}^{(\text{train})\top} \mathbf{X}^{(\text{train})} \mathbf{w} - 2 \mathbf{w}^\top \mathbf{X}^{(\text{train})\top} \mathbf{y}^{(\text{train})} + \mathbf{y}^{(\text{train})\top} \mathbf{y}^{(\text{train})} \right) = 0 \quad (5.10)$$

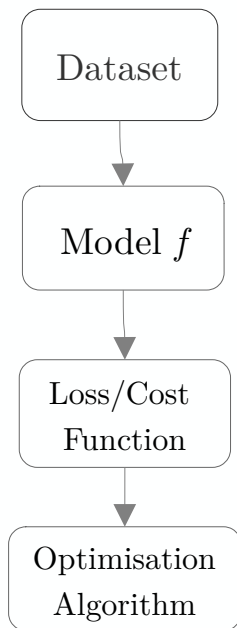
$$\Rightarrow 2 \mathbf{X}^{(\text{train})\top} \mathbf{X}^{(\text{train})} \mathbf{w} - 2 \mathbf{X}^{(\text{train})\top} \mathbf{y}^{(\text{train})} = 0 \quad (5.11)$$

$$\Rightarrow \mathbf{w} = \left(\mathbf{X}^{(\text{train})\top} \mathbf{X}^{(\text{train})} \right)^{-1} \mathbf{X}^{(\text{train})\top} \mathbf{y}^{(\text{train})} \quad (5.12)$$



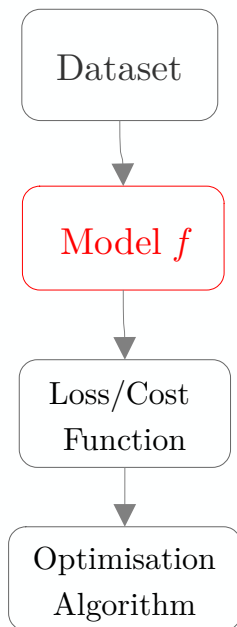
Linear Regression Example

General recipe for *most* ML algorithms or learning processes:



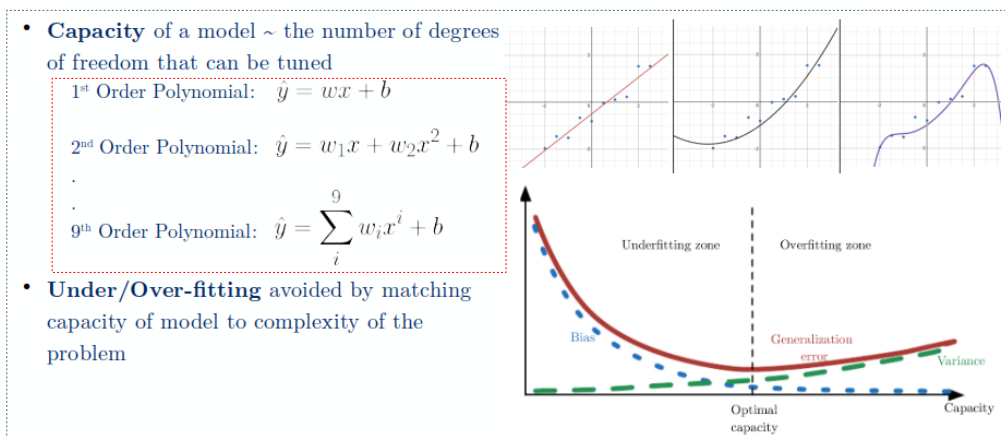
Are we done?

General recipe for *most* ML algorithms or learning processes:



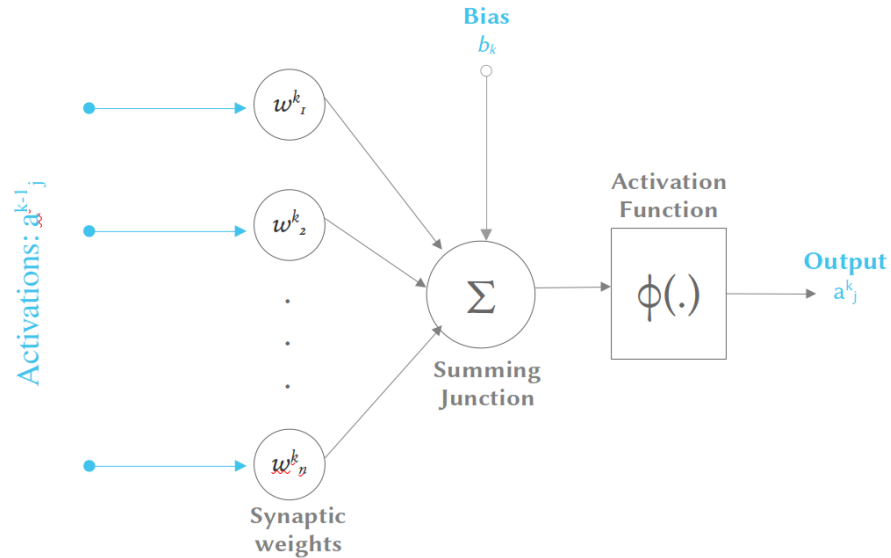
Are we done?

- How do we change the **family of functions** that a model learns? Scale capacity?





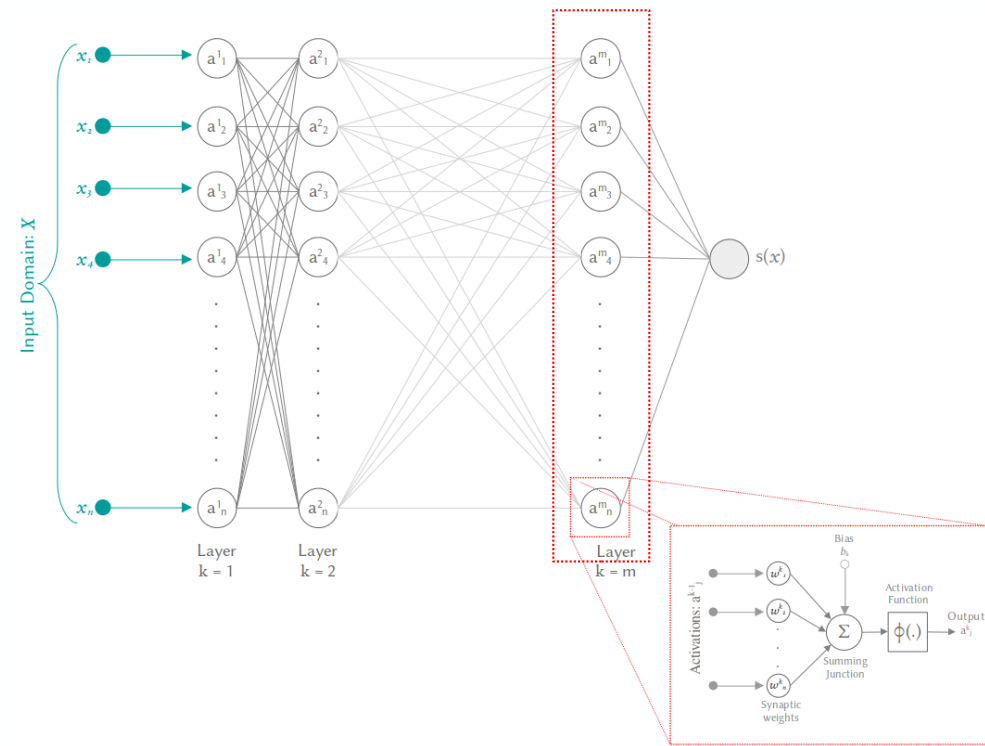
Neural Networks



- **Building block** of neural networks:

$$a_i^{(k+1)} = \phi \left(\sum_j^n w_{i,j}^{(k)} a_j^{(k)} + b^{(k)} \right)$$

Neuron \rightarrow Layer



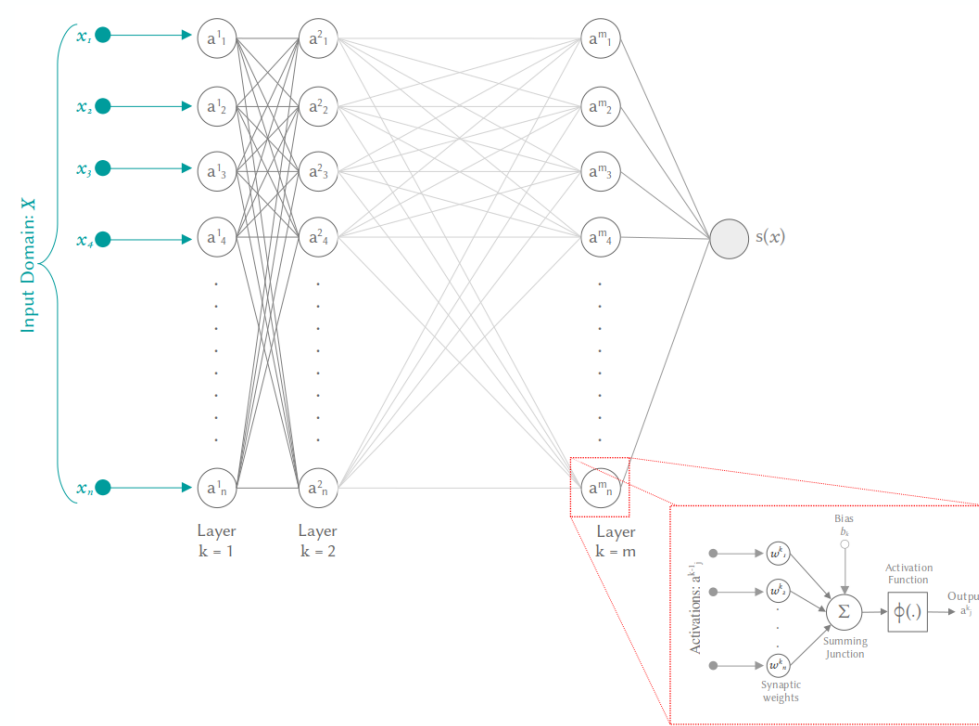
- Building block** of neural networks:

$$a_i^{(k+1)} = \phi \left(\sum_j^n w_{i,j}^{(k)} a_j^{(k)} + b^{(k)} \right)$$

Where, $a^{(k)}_i$ is the activation value, and the $\Phi = \{w^{(k)}_i, b^{(k)}\}$ are the configurable weights & biases

- Neural Network** layer is the combination of many neurons densely connected (mostly...):

$$\begin{bmatrix} a_0^{(k+1)} \\ \vdots \\ a_n^{(k+1)} \end{bmatrix} = \phi \left(\begin{bmatrix} w_0^{(0)} & \dots & w_n^{(0)} \\ \vdots & \ddots & \vdots \\ w_0^{(k)} & \dots & w_n^{(k)} \end{bmatrix} \begin{bmatrix} a_0^{(k)} \\ \vdots \\ a_n^{(k)} \end{bmatrix} + \begin{bmatrix} b_0^{(k)} \\ \vdots \\ b_n^{(k)} \end{bmatrix} \right)$$



- A **total network** is therefore:

$$\mathbf{a}^{(1)} = \phi \left(\mathbf{W}^{(0)} \mathbf{x} + \beta^{(0)} \right)$$

$$\vdots$$

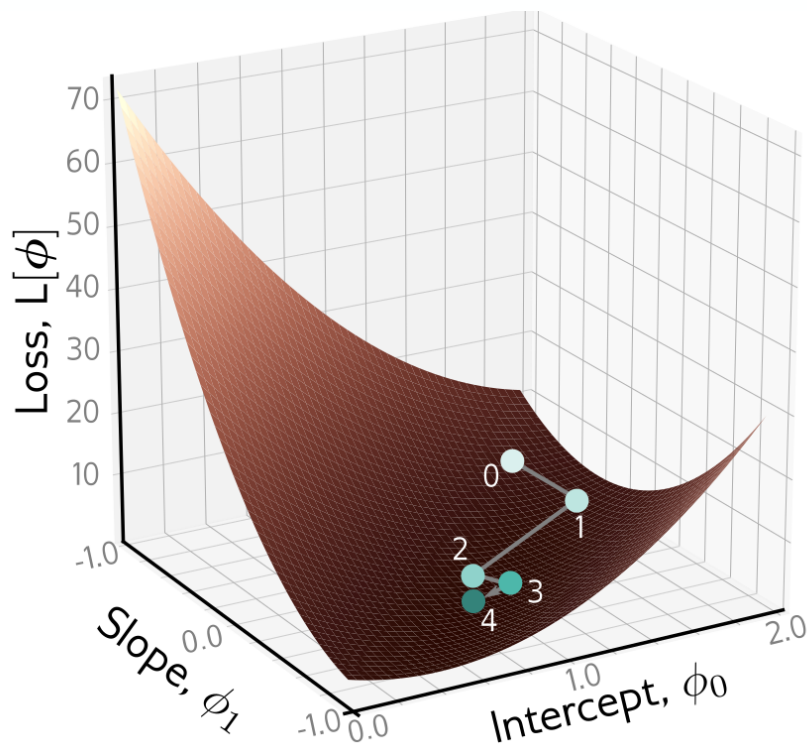
$$\mathbf{a}^{(k)} = \phi \left(\mathbf{W}^{(k-1)} \mathbf{a}^{(k-1)} + \beta^{(k-1)} \right)$$

$$\vdots$$

$$s(\mathbf{x}) = \mathbf{W}^{(k)} \mathbf{a}^{(k)} + \beta^{(k)}$$

- ϕ is the activation function, which can take many forms
- Optimisation is done via typically gradient based descent algorithms

Stochastic Gradient Descent



- **Stochastic Gradient Descent** is an adaptation of gradient descent:

1) Calculate gradients:

$$\frac{\partial \mathcal{L}}{\partial \Phi} = \begin{bmatrix} \frac{d\mathcal{L}}{d\phi_0} \\ \dots \\ \frac{d\mathcal{L}}{d\phi_h} \end{bmatrix}$$

2) Update the parameters in direction that minimises loss:

$$\phi_i^{(t+1)} \leftarrow \phi_i^{(t)} - \alpha \cdot \sum_{j \in \mathcal{B}_t} \frac{\partial \mathcal{L}_j}{\partial \phi_{i,j}^{(t)}}$$

- **Batch based** training with iterative updates to the gradients based on an **epoch** ‘*t*’, but the key goal is to learn:

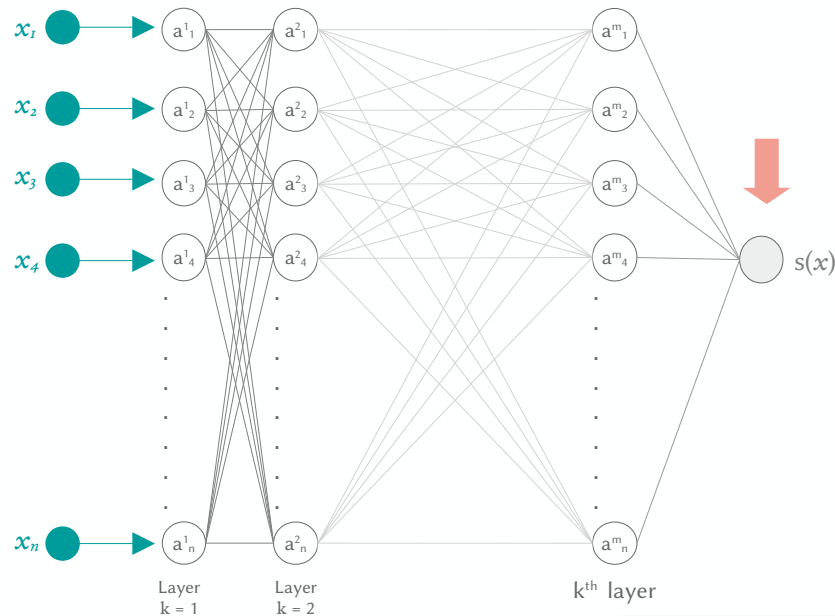
$$\frac{\partial \mathcal{L}_j}{d\beta_j^{(k)}} \quad \frac{\partial \mathcal{L}_j}{d\mathbf{W}_j^{(k)}}$$

Backwards Propagation – Optimisation Algorithm

95

$$s(\mathbf{x}) = \mathbf{W}^{(k)} \mathbf{a}^{(k)} + \beta^{(k)}$$

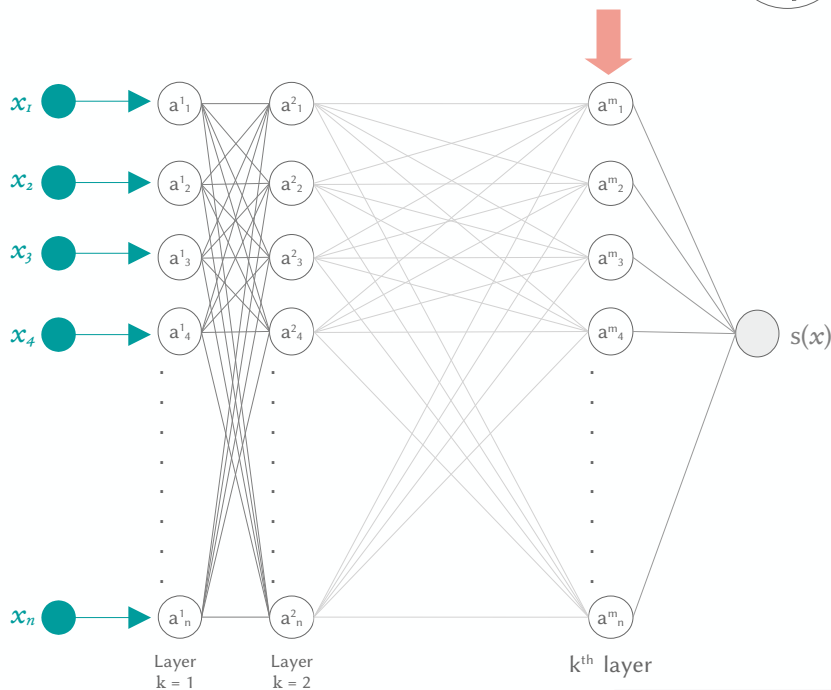
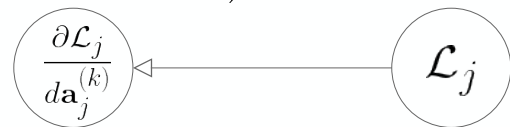
\mathcal{L}_j



Backwards Propagation – Optimisation Algorithm

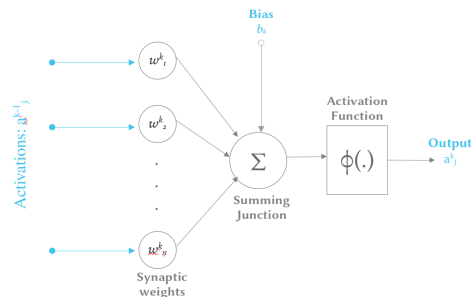
96

$$\mathbf{a}^{(k)} = \phi \left(\mathbf{W}^{(k-1)} \mathbf{a}^{(k-1)} + \beta^{(k-1)} \right) \quad \mathbf{s}(\mathbf{x}) = \mathbf{W}^{(k)} \mathbf{a}^{(k)} + \beta^{(k)}$$

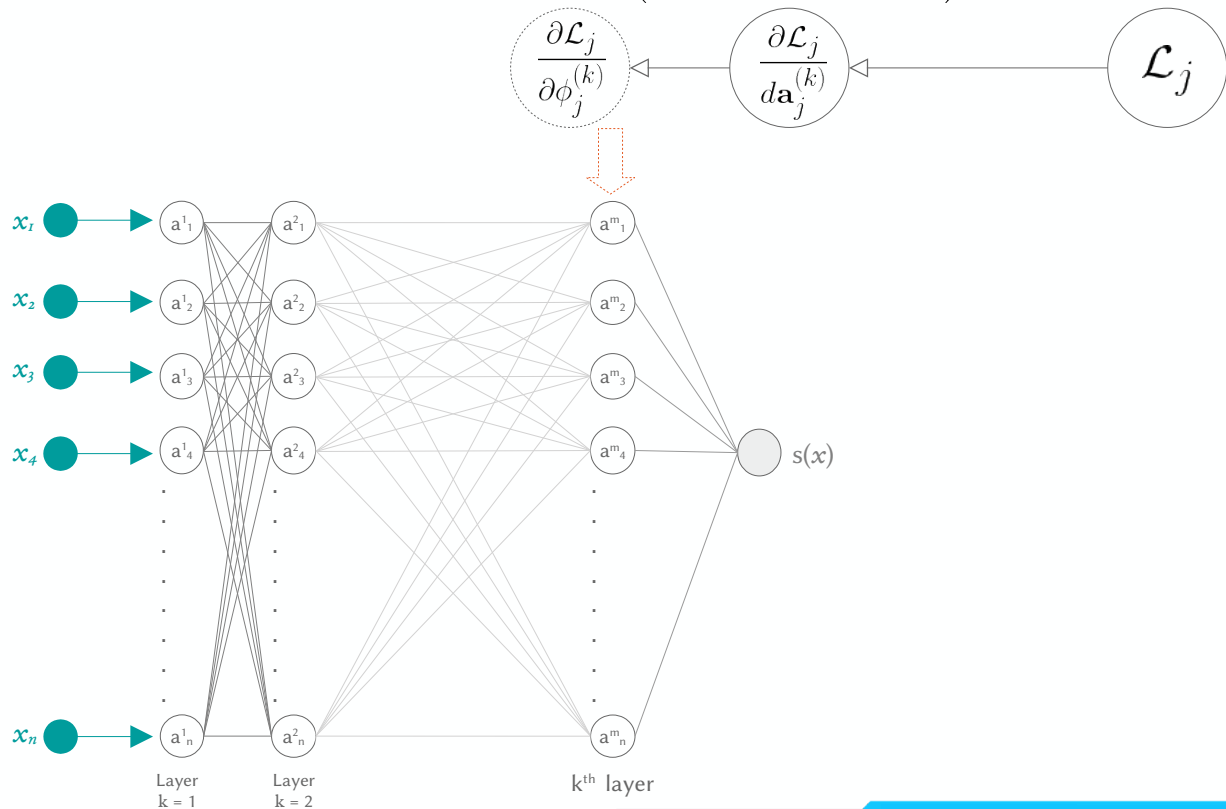


Backwards Propagation – Optimisation Algorithm

97

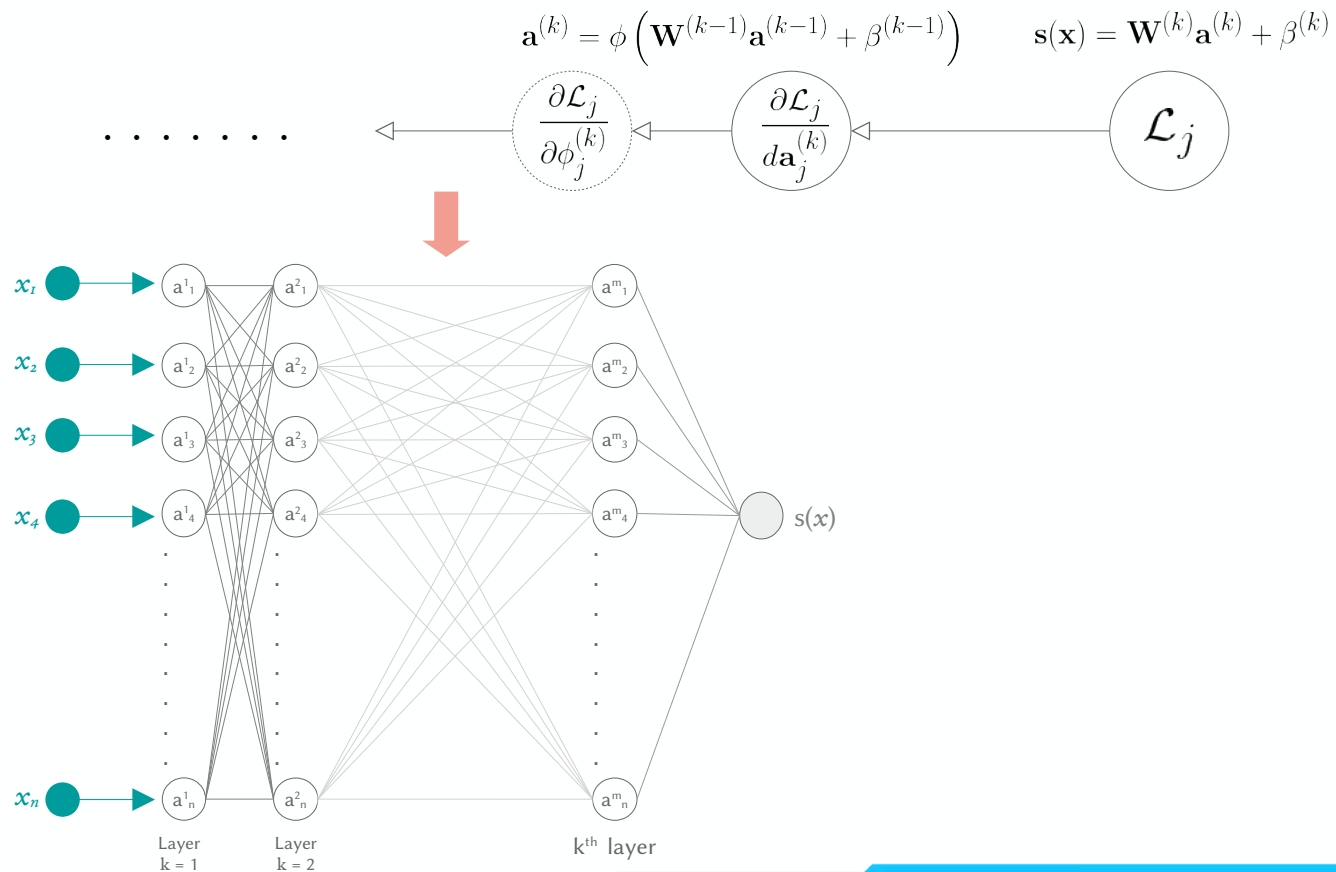


$$\mathbf{a}^{(k)} = \phi \left(\mathbf{W}^{(k-1)} \mathbf{a}^{(k-1)} + \beta^{(k-1)} \right) \quad \mathbf{s}(\mathbf{x}) = \mathbf{W}^{(k)} \mathbf{a}^{(k)} + \beta^{(k)}$$



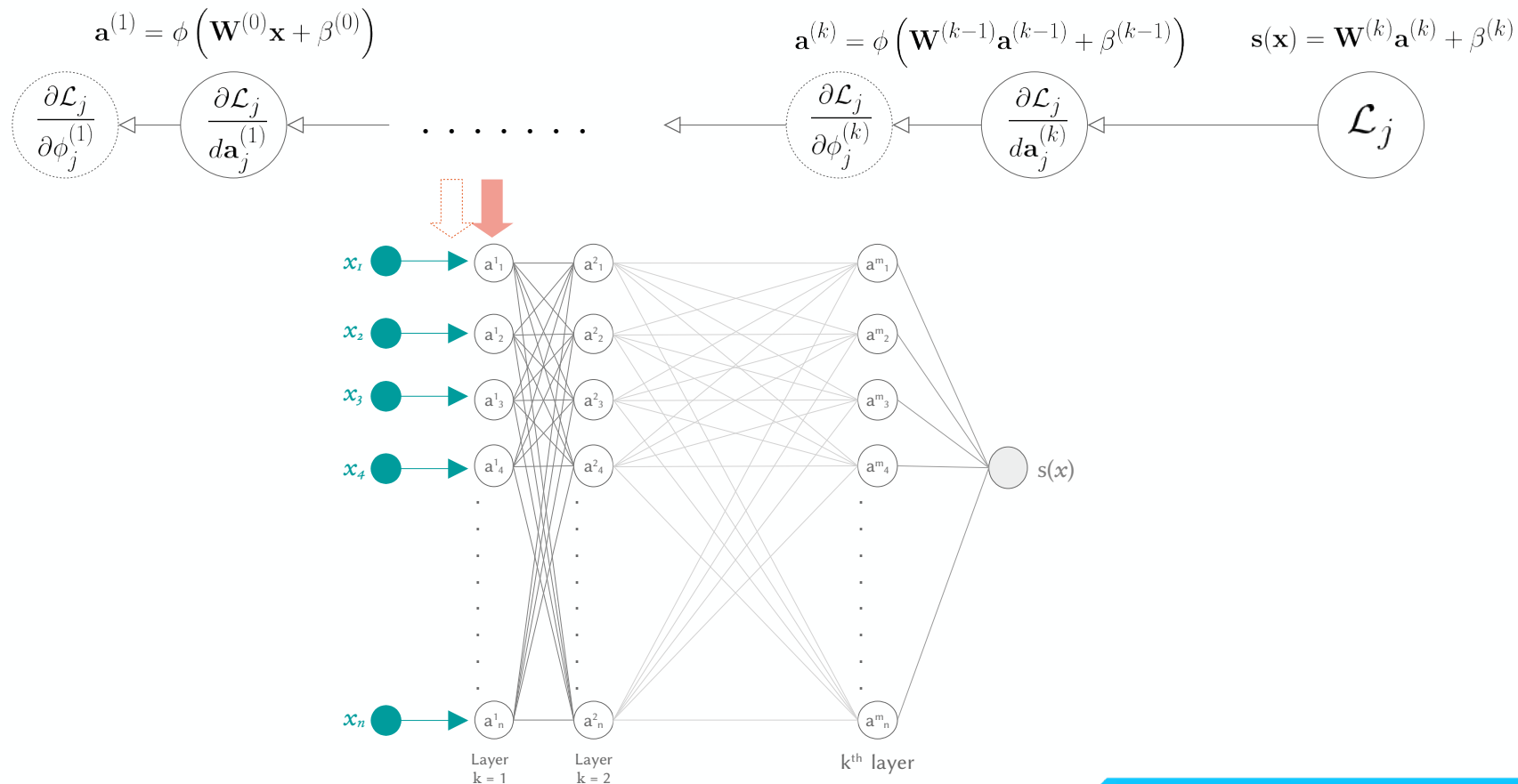
Backwards Propagation – Optimisation Algorithm

98



Backwards Propagation – Optimisation Algorithm

99

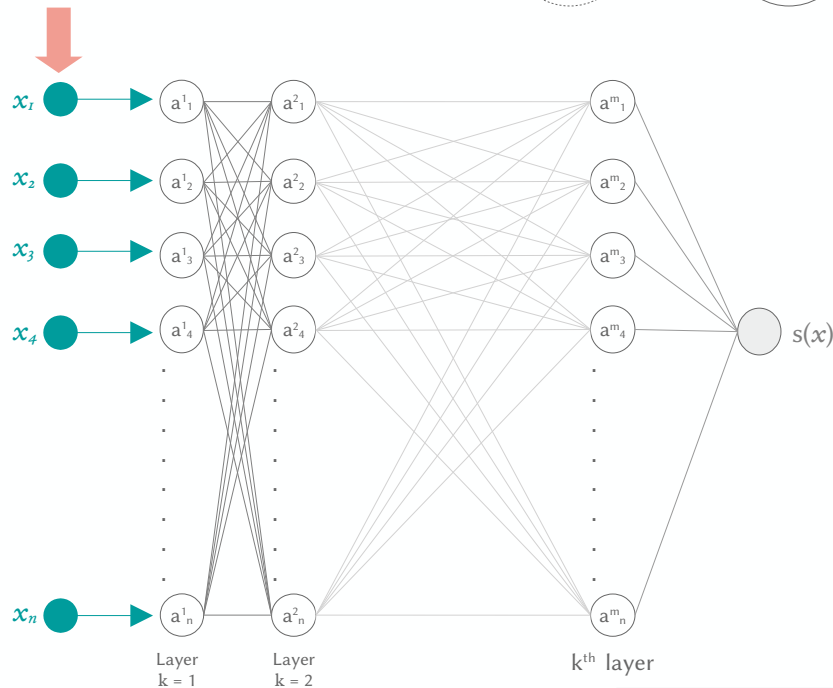


Backwards Propagation – Optimisation Algorithm

100

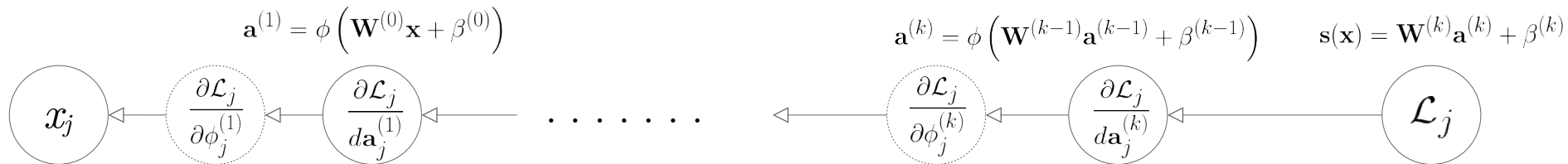
$$\mathbf{a}^{(1)} = \phi \left(\mathbf{W}^{(0)} \mathbf{x} + \beta^{(0)} \right)$$

$$\mathbf{a}^{(k)} = \phi \left(\mathbf{W}^{(k-1)} \mathbf{a}^{(k-1)} + \beta^{(k-1)} \right) \quad \mathbf{s}(\mathbf{x}) = \mathbf{W}^{(k)} \mathbf{a}^{(k)} + \beta^{(k)}$$



Backwards Propagation – Optimisation Algorithm

101

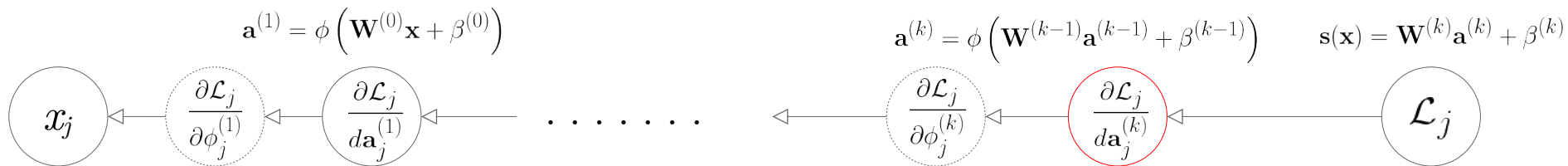


- **Chain rule** used to decompose the differential behaviour of the loss \mathcal{L} per-layer (intermediate) variables in reverse order:

$$\frac{\partial y}{\partial x^{(1)}} = \frac{\partial x^{(2)}}{\partial x^{(1)}} \frac{\partial x^{(3)}}{\partial x^{(2)}} \frac{\partial y}{\partial x^{(3)}}$$

Backwards Propagation – Optimisation Algorithm

102



- **Chain rule** used to decompose the differential behaviour of the loss \mathcal{L} per-layer (intermediate) variables in reverse order:

$$\frac{\partial y}{\partial x^{(1)}} = \frac{\partial x^{(2)}}{\partial x^{(1)}} \frac{\partial x^{(3)}}{\partial x^{(2)}} \frac{\partial y}{\partial x^{(3)}}$$

$$\frac{\partial \mathcal{L}_j}{\partial \mathbf{a}_j^{(k)}} \quad \leftarrow$$

$$\mathbf{a}^{(1)} = \phi \left(\mathbf{W}^{(0)} \mathbf{x} + \beta^{(0)} \right)$$

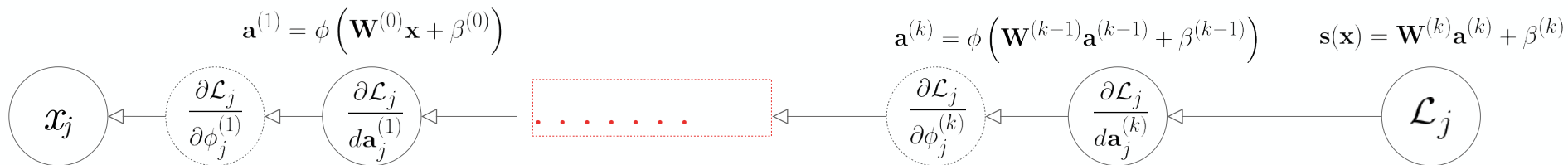
$$\mathbf{a}^{(k)} = \phi \left(\mathbf{W}^{(k-1)} \mathbf{a}^{(k-1)} + \beta^{(k-1)} \right) \quad \mathbf{s}(\mathbf{x}) = \mathbf{W}^{(k)} \mathbf{a}^{(k)} + \beta^{(k)}$$



- **Chain rule** used to decompose the differential behaviour of the loss \mathcal{L} per-layer (intermediate) variables in reverse order:

$$\frac{\partial \mathcal{L}_j}{\partial \phi_j^{(k)}} = \frac{\partial \mathbf{a}_j^{(k)}}{\partial \phi^{(k)}} \frac{\partial \mathcal{L}_j}{\partial \mathbf{a}_j^{(k)}} \quad \leftarrow$$

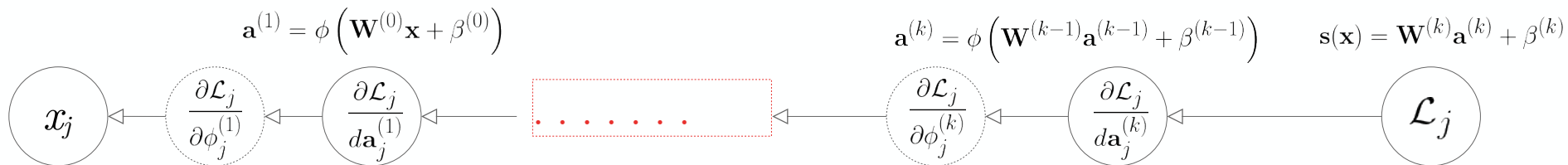
$$\frac{\partial y}{\partial x^{(1)}} = \frac{\partial x^{(2)}}{\partial x^{(1)}} \frac{\partial x^{(3)}}{\partial x^{(2)}} \frac{\partial y}{\partial x^{(3)}}$$



- **Chain rule** used to decompose the differential behaviour of the loss \mathcal{L} per-layer (intermediate) variables in reverse order:

$$\frac{\partial y}{\partial x^{(1)}} = \frac{\partial x^{(2)}}{\partial x^{(1)}} \frac{\partial x^{(3)}}{\partial x^{(2)}} \frac{\partial y}{\partial x^{(3)}}$$

$$\frac{\partial \mathcal{L}_j}{\partial \mathbf{a}_j^{(k-1)}} = \frac{\partial \phi_j^{(k)}}{\partial \mathbf{a}_j^{(k-1)}} \left(\frac{\partial \mathbf{a}_j^{(k)}}{\partial \phi_j^{(k)}} \frac{\partial \mathcal{L}_j}{\partial \mathbf{a}_j^{(k)}} \right) \quad \leftarrow$$



- **Chain rule** used to decompose the differential behaviour of the loss \mathcal{L} per-layer (intermediate) variables in reverse order:

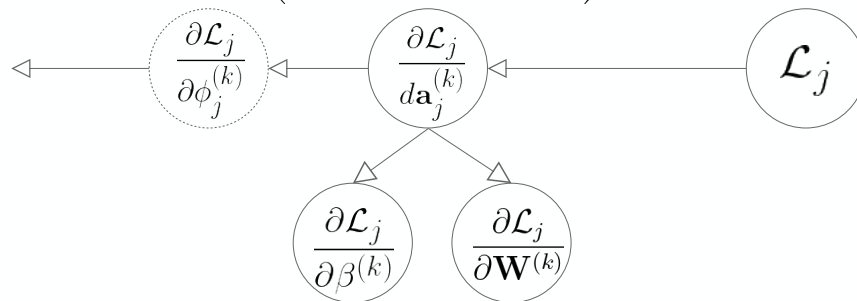
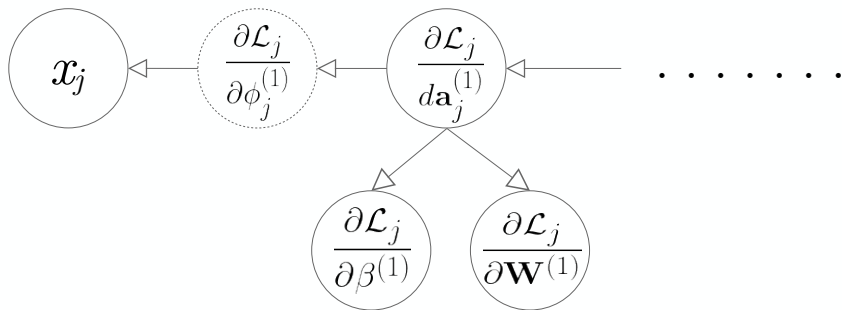
$$\frac{\partial y}{\partial x^{(1)}} = \frac{\partial x^{(2)}}{\partial x^{(1)}} \frac{\partial x^{(3)}}{\partial x^{(2)}} \frac{\partial y}{\partial x^{(3)}}$$

$$\frac{\partial \mathcal{L}_j}{\partial \mathbf{a}_j^{(k-1)}} = \frac{\partial \mathbf{a}_j^{(k-1)}}{\partial \phi_j^{(k-1)}} \left(\frac{\partial \phi_j^{(k)}}{\partial \mathbf{a}_j^{(k-1)}} \frac{\partial \mathbf{a}_j^{(k)}}{\partial \phi_j^{(k)}} \frac{\partial \mathcal{L}_j}{\partial \mathbf{a}_j^{(k)}} \right)$$

$$\frac{\partial \mathcal{L}_j}{\partial \phi_j^{(k-1)}} = \frac{\partial \mathbf{a}_j^{(k-1)}}{\partial \phi_j^{(k-1)}} \left(\frac{\partial \phi_j^{(k)}}{\partial \mathbf{a}_j^{(k-1)}} \frac{\partial \mathbf{a}_j^{(k)}}{\partial \phi_j^{(k)}} \frac{\partial \mathcal{L}_j}{\partial \mathbf{a}_j^{(k)}} \right) \quad \leftarrow$$

$$\mathbf{a}^{(1)} = \phi \left(\mathbf{W}^{(0)} \mathbf{x} + \beta^{(0)} \right)$$

$$\mathbf{a}^{(k)} = \phi \left(\mathbf{W}^{(k-1)} \mathbf{a}^{(k-1)} + \beta^{(k-1)} \right) \quad \mathbf{s}(\mathbf{x}) = \mathbf{W}^{(k)} \mathbf{a}^{(k)} + \beta^{(k)}$$



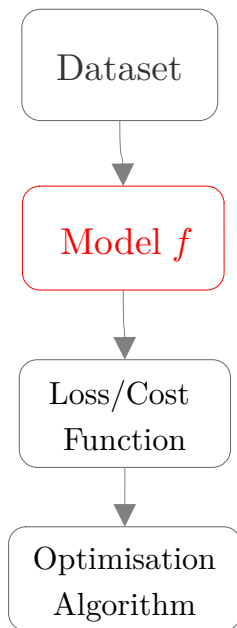
- **Chain rule** again used to decompose the differential behaviour of the loss \mathcal{L} per-parameter $\Phi = \{\mathbf{W}, \beta\}$:

$$\frac{\partial \mathcal{L}_j}{\partial \mathbf{W}_j^{(k)}} = \frac{\partial \mathbf{a}_j^{(k)}}{\partial \mathbf{W}^{(k)}} \frac{\partial \mathcal{L}_j}{\partial \mathbf{a}_j^{(k)}} \quad \frac{\partial \mathcal{L}_j}{\partial \mathbf{a}_j^{(k-1)}} = \frac{\partial \phi_j^{(k)}}{\partial \mathbf{a}_j^{(k-1)}} \left(\frac{\partial \mathbf{a}_j^{(k)}}{\partial \phi_j^{(k)}} \frac{\partial \mathcal{L}_j}{\partial \mathbf{a}_j^{(k)}} \right)$$

$$\frac{\partial \mathcal{L}_j}{\partial \beta_j^{(k)}} = \frac{\partial \mathbf{a}_j^{(k)}}{\partial \beta^{(k)}} \frac{\partial \mathcal{L}_j}{\partial \mathbf{a}_j^{(k)}} \quad \frac{\partial \mathcal{L}_j}{\partial \beta_j^{(k-1)}} = \frac{\partial \mathbf{a}_j^{(k)}}{\partial \beta^{(k-1)}} \frac{\partial \mathcal{L}_j}{\partial \mathbf{a}_j^{(k-1)}} \left(\frac{\partial \mathbf{a}_j^{(k)}}{\partial \phi_j^{(k)}} \frac{\partial \mathcal{L}_j}{\partial \mathbf{a}_j^{(k)}} \right)$$

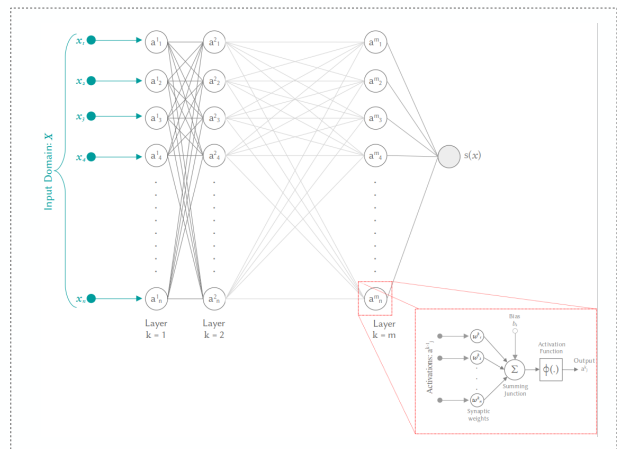
Red arrows indicate the flow of gradients from the right-hand side of the equations to the corresponding terms in the left-hand side equations.

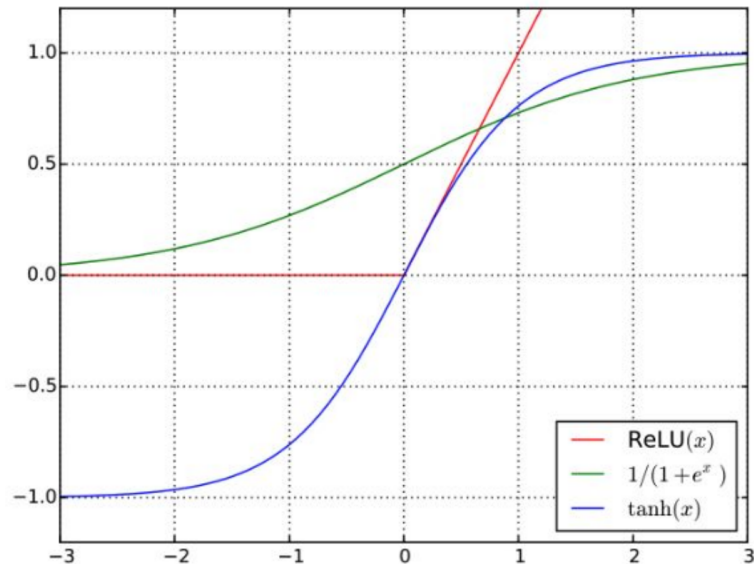
General recipe for *most* ML algorithms or learning processes:



Are we done?

- How do we change the **family of functions** that a model learns? Scale capacity? ☒
- What if we want to predict a discrete value and not a continuous value (classification)? ☐





- **Choice** of ϕ activation functions is pathological
- Decisions based on an array of factors
- The most common is the concept of **vanishing gradients**:

$$\frac{\partial \mathcal{L}_j}{\partial \mathbf{W}_j^{(k-2)}} = \frac{\partial \mathbf{a}_j^{(k-2)}}{\partial \mathbf{W}_j^{(k-2)}} \left(\underbrace{\frac{\partial \phi_j^{(k-1)}}{\partial \mathbf{a}_j^{(k-2)}} \frac{\partial \mathbf{a}_j^{(k-1)}}{\partial \phi^{(k-1)}}}_{\text{Sigmoid} \leq 1/4} \underbrace{\frac{\partial \phi_j^{(k)}}{\partial \mathbf{a}_j^{(k-1)}} \frac{\partial \mathbf{a}_j^{(k)}}{\partial \phi^{(k)}}}_{\text{ReLU} \leq 1} \frac{\partial \mathcal{L}_j}{\partial \mathbf{a}_j^{(k)}} \right)$$

Sigmoid	ReLU	Sigmoid	ReLU
$\leq 1/4$	≤ 1	$\leq 1/4$	≤ 1

Classification – Discrete Value Predictions

- **Task:** Classify using a dataset drawn from a joint distribution

$p(x, y)$:

Features $x \in \mathbb{R}^n$

Labels $y \in \mathbb{R}$

- **Goal** is to predict y given an instance of x : $p(y = 1 | \mathbf{x}) = \frac{p(\mathbf{x} | y)p(y = 1)}{p(\mathbf{x})}$

Classification – Discrete Value Predictions

- **Task:** Classify using a dataset drawn from a joint distribution

$p(x, y)$:

Features $\mathbf{x} \in \mathbb{R}^n$

Labels $y \in \mathbb{R}$

- **Goal** is to predict y given an instance of x :

$$\begin{aligned}
 p(y = 1 | \mathbf{x}) &= \frac{p(\mathbf{x} | y)p(y = 1)}{p(\mathbf{x})} \\
 &= \frac{1}{\frac{p(\mathbf{x} | y=0)p(y=0)}{p(\mathbf{x} | y)p(y=1)} + 1} \\
 &= \frac{p(\mathbf{x} | y)p(y = 1)}{p(\mathbf{x} | y = 0)p(y = 0) + p(\mathbf{x} | y = 1)p(y = 1)} \\
 &= \frac{1}{1 + \exp \left(\log \left(\frac{p(\mathbf{x} | y=0)p(y=0)}{p(\mathbf{x} | y=1)p(y=1)} \right) \right)}
 \end{aligned}$$

Classification – Discrete Value Predictions

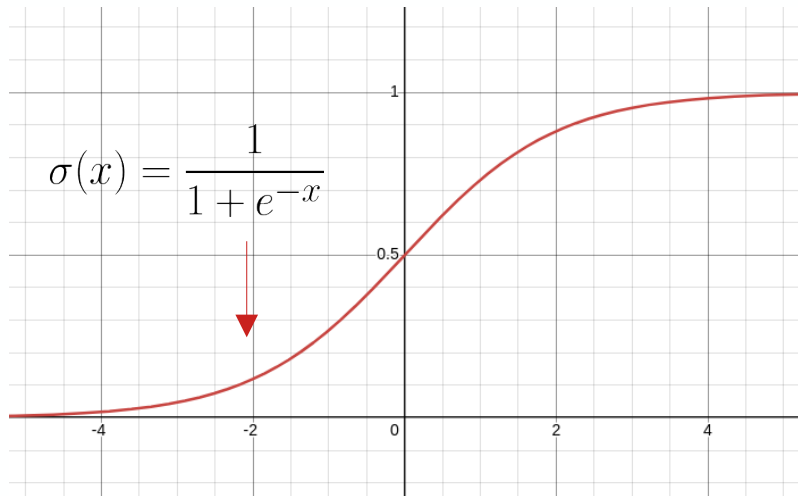
- Task:** Classify using a dataset drawn from a joint distribution

$p(x, y)$:

Features $x \in \mathbb{R}^n$

Labels $y \in \mathbb{R}$

- Goal** is to predict y given an instance of x : $p(y = 1 | \mathbf{x}) = \frac{p(\mathbf{x}|y)p(y = 1)}{p(\mathbf{x})}$



$$\begin{aligned}
 p(y = 1 | \mathbf{x}) &= \frac{p(\mathbf{x}|y)p(y = 1)}{p(\mathbf{x})} \\
 &= \frac{1}{\frac{p(\mathbf{x}|y=0)p(y=0)}{p(\mathbf{x}|y)p(y=1)} + 1} \\
 &= \frac{p(\mathbf{x}|y)p(y = 1)}{p(\mathbf{x}|y = 0)p(y = 0) + p(\mathbf{x}|y = 1)p(y = 1)} \\
 &= \frac{1}{1 + \exp \left(\log \left(\frac{p(\mathbf{x}|y=0)p(y=0)}{p(\mathbf{x}|y=1)p(y=1)} \right) \right)}
 \end{aligned}$$

Classification – Discrete Value Predictions

- Task:** Classify using a dataset drawn from a joint distribution

$p(x, y)$:

Features $x \in \mathbb{R}^n$

Labels $y \in \mathbb{R}$

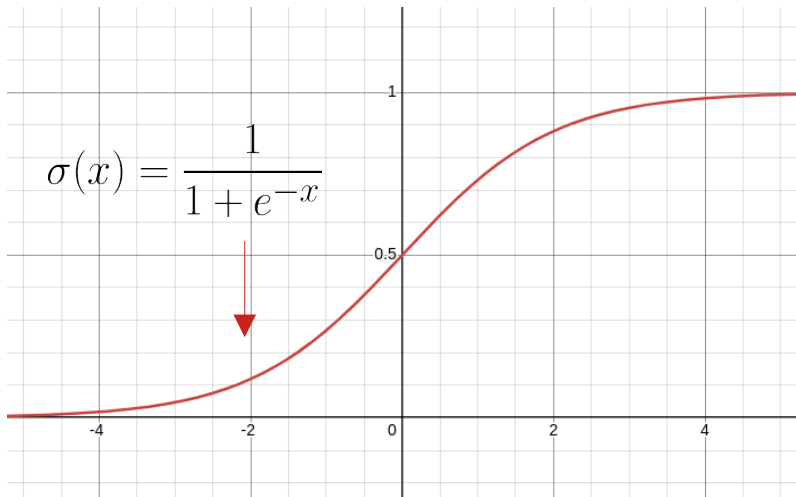
- Goal** is to predict y given an instance of x :

Log-Likelihood Ratio

$$p(y = 1 | \mathbf{x}) = \sigma \left(\log \left(\frac{p(\mathbf{x} | y = 1)}{p(\mathbf{x} | y = 0)} \right) \right) + \log \left(\frac{p(y = 1)}{p(y = 0)} \right)$$

Class Marginal

(does not depend on x)



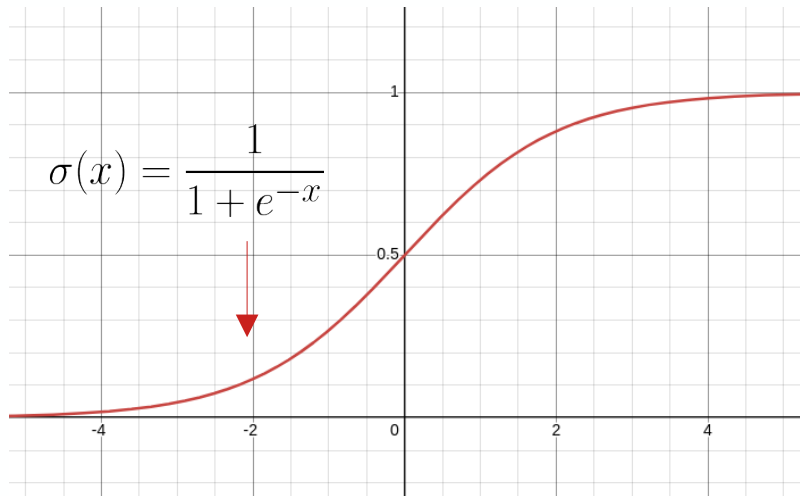
- **Task:** Classify using a dataset drawn from a joint distribution

$p(x, y)$:

Features $x \in \mathbb{R}^n$

Labels $y \in \mathbb{R}$

- **Goal** is to predict y given an instance of x :

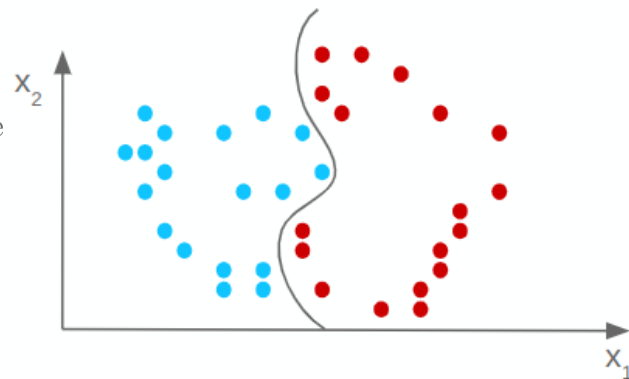


The sigmoid converts the distance from a normal regression problem to a probabilistic decision boundary

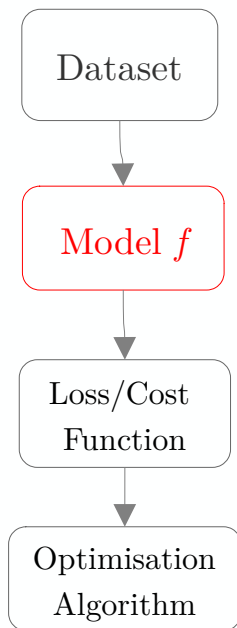
Log-Likelihood Ratio

$$p(y = 1 | \mathbf{x}) = \sigma \left(\log \left(\frac{p(\mathbf{x} | y = 1)}{p(\mathbf{x} | y = 0)} \right) \right) + \log \left(\frac{p(y = 1)}{p(y = 0)} \right)$$

Class Marginal
(does not depend on x)

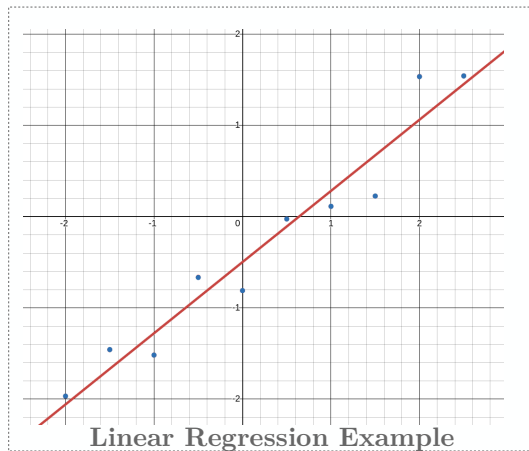


General recipe for *most* ML algorithms or learning processes:



Are we done?

- How do we change the **family of functions** that a model learns? Scale capacity? ✓
- What if we want to predict a discrete value and not a continuous value (classification)? ✓



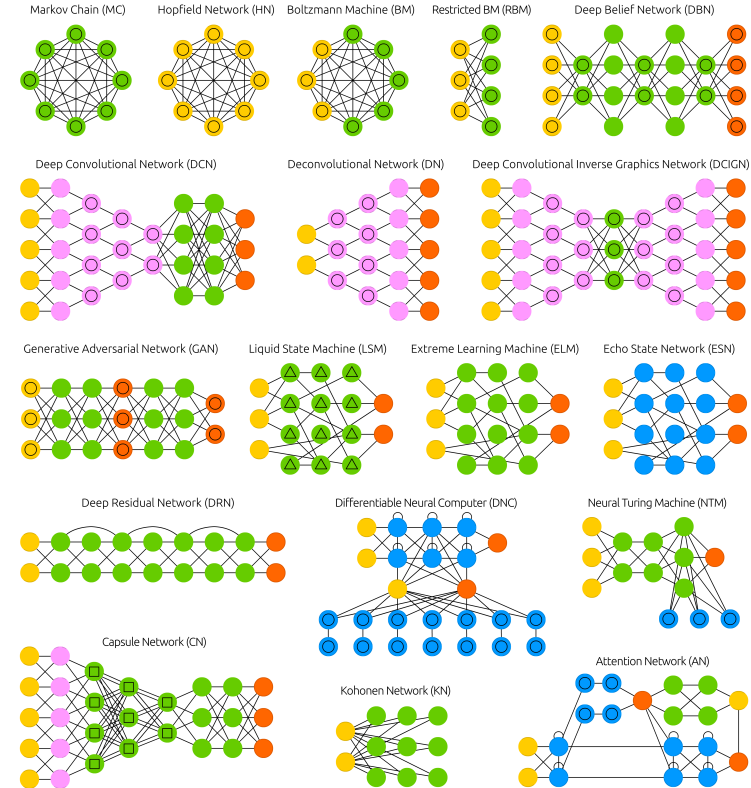
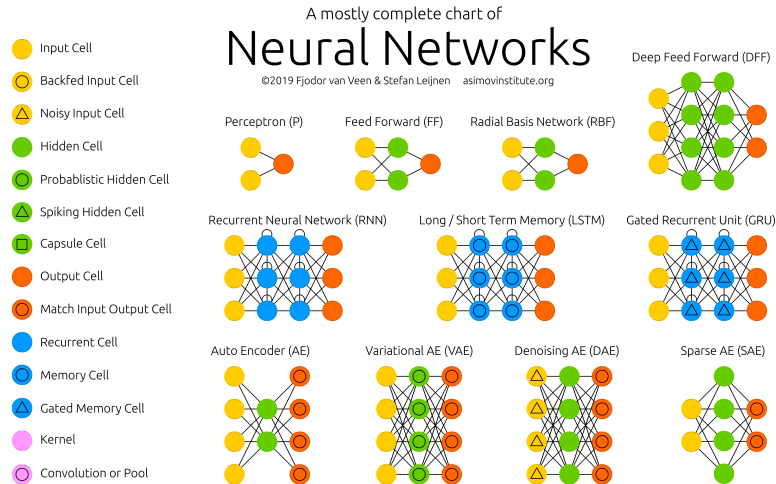


Neural Zoo

Diversity of Neural Networks

116

- **Ever growing** number of fundamental building blocks to neural based systems
- **No universally best block** for all uses cases, the various building blocks have different strengths and weaknesses



Source: [link](#)



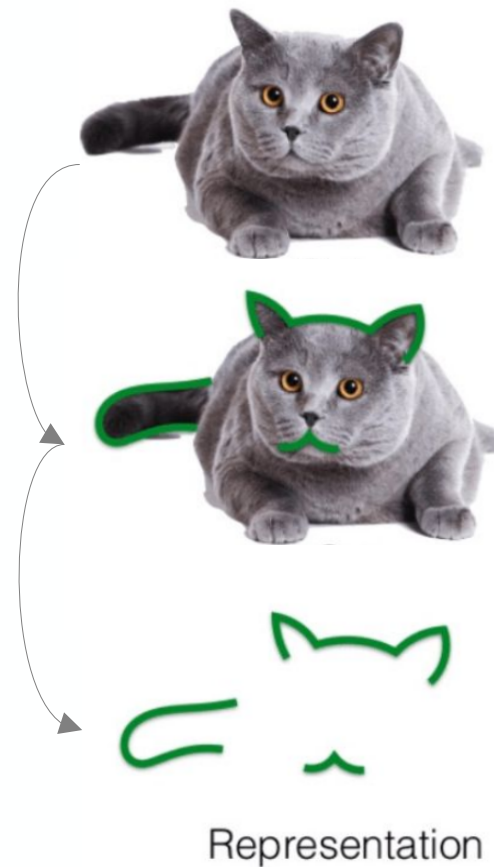
Neural Zoo

Convolutional Neural Networks

Convolutional Neural Network

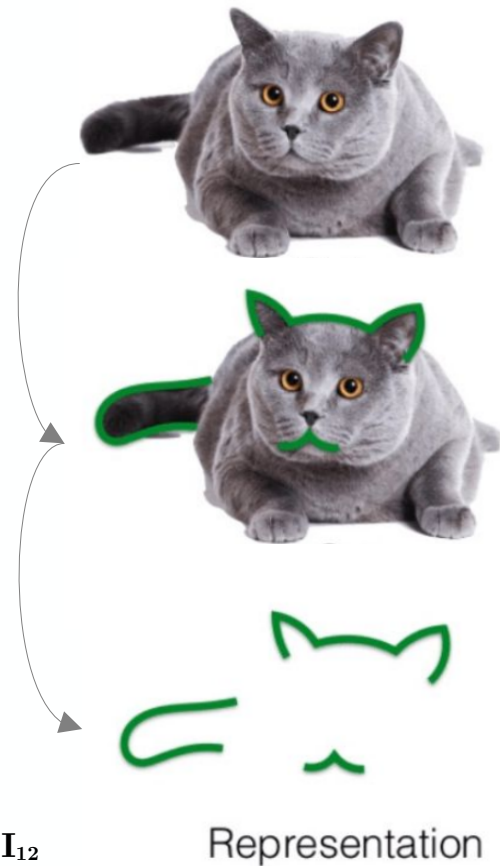
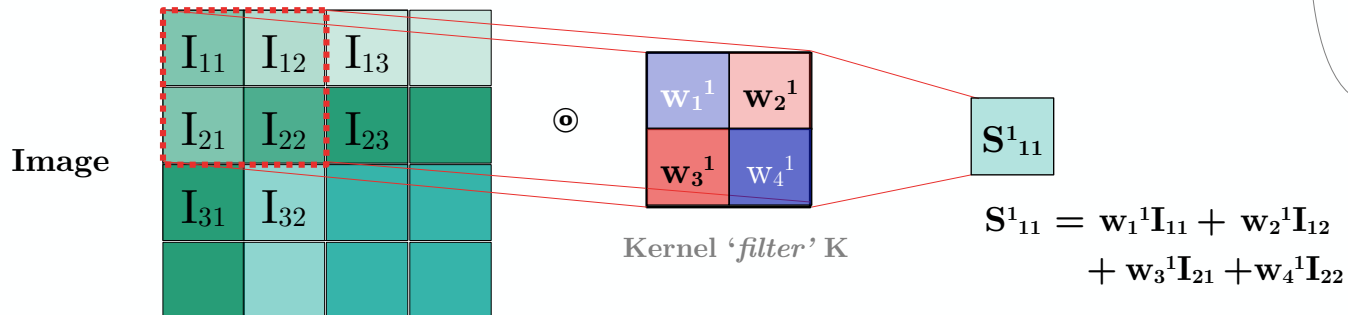
118

- **Extract** relevant **features** from a high dimensional input domain space given by the pixel space of an image of height and width $\mathbf{H} \times \mathbf{W} \rightarrow x \in \mathbb{R}^{\mathbf{H} \times \mathbf{W}}$
- **Salient features** ~ e.g. the lines of contrast between the foreground and the background



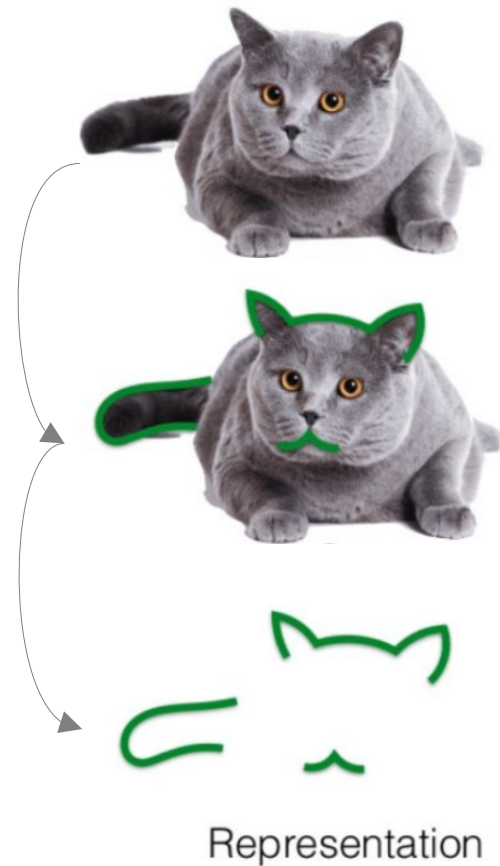
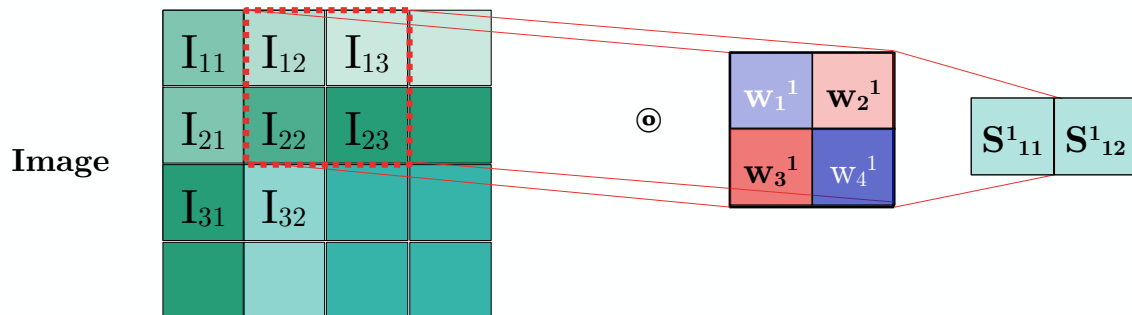
Convolutional Neural Network

- **Extract** relevant **features** from a high dimensional input domain space given by the pixel space of an image of height and width $\mathbf{H} \times \mathbf{W} \rightarrow x \in \mathbb{R}^{\mathbf{H} \times \mathbf{W}}$
- **Salient features** ~ e.g. the lines of contrast between the foreground and the background
- The key building block is the **kernel filter**:



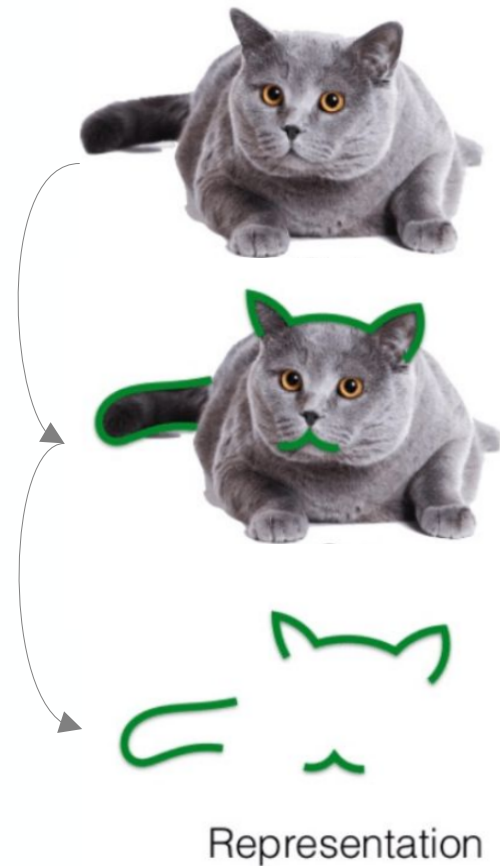
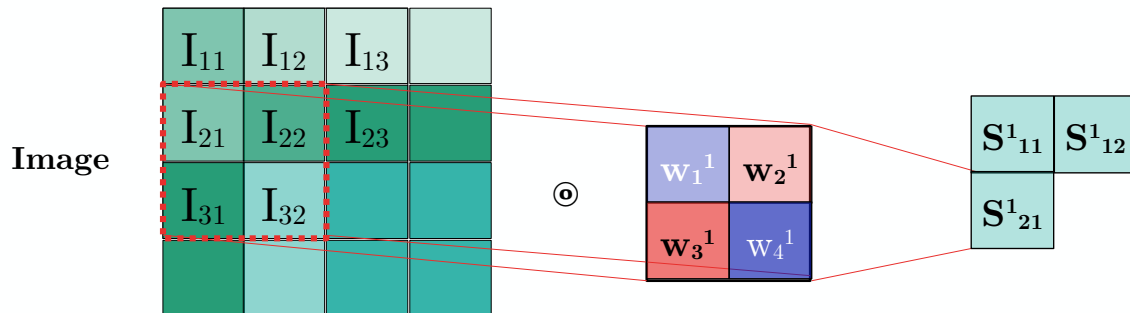
Convolutional Neural Network

- **Extract** relevant **features** from a high dimensional input domain space given by the pixel space of an image of height and width $\mathbf{H} \times \mathbf{W} \rightarrow x \in \mathbb{R}^{\mathbf{H} \times \mathbf{W}}$
- **Salient features** ~ e.g. the lines of contrast between the foreground and the background
- The key building block is the **kernel filter**:



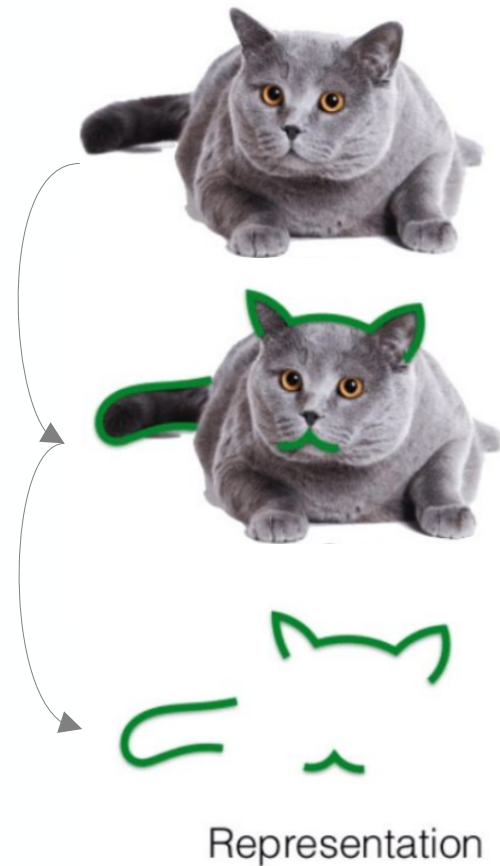
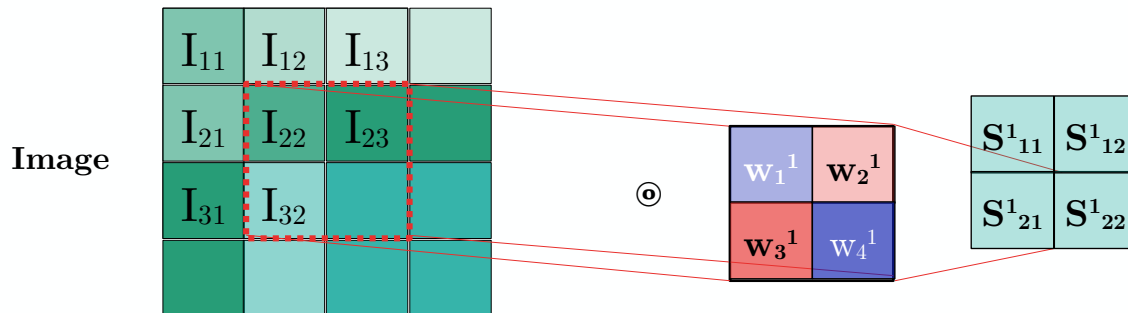
Convolutional Neural Network

- **Extract** relevant **features** from a high dimensional input domain space given by the pixel space of an image of height and width $\mathbf{H} \times \mathbf{W} \rightarrow x \in \mathbb{R}^{\mathbf{H} \times \mathbf{W}}$
- **Salient features** ~ e.g. the lines of contrast between the foreground and the background
- The key building block is the **kernel filter**:



Convolutional Neural Network

- **Extract** relevant **features** from a high dimensional input domain space given by the pixel space of an image of height and width $\mathbf{H} \times \mathbf{W} \rightarrow x \in \mathbb{R}^{\mathbf{H} \times \mathbf{W}}$
- **Salient features** ~ e.g. the lines of contrast between the foreground and the background
- The key building block is the **kernel filter**:



Convolutional Neural Network

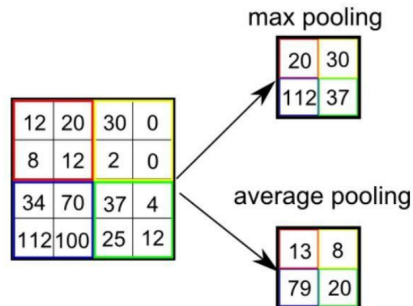
- **Formally** the kernel convolution operation:

$$S(i, j) = (I * K)(i, j) = \sum_m \sum_n I(m, n) K(i - m, j - n)$$

- **Salient feature map 'S' is by default smaller:**

$$S \in \mathbb{R}^{(H-K+2P/S+1) \times (W-K+2P/S+1)}$$

- **Compressing** image size further to reduce input dimension by **max/average pooling**



Convolutional Neural Network

124

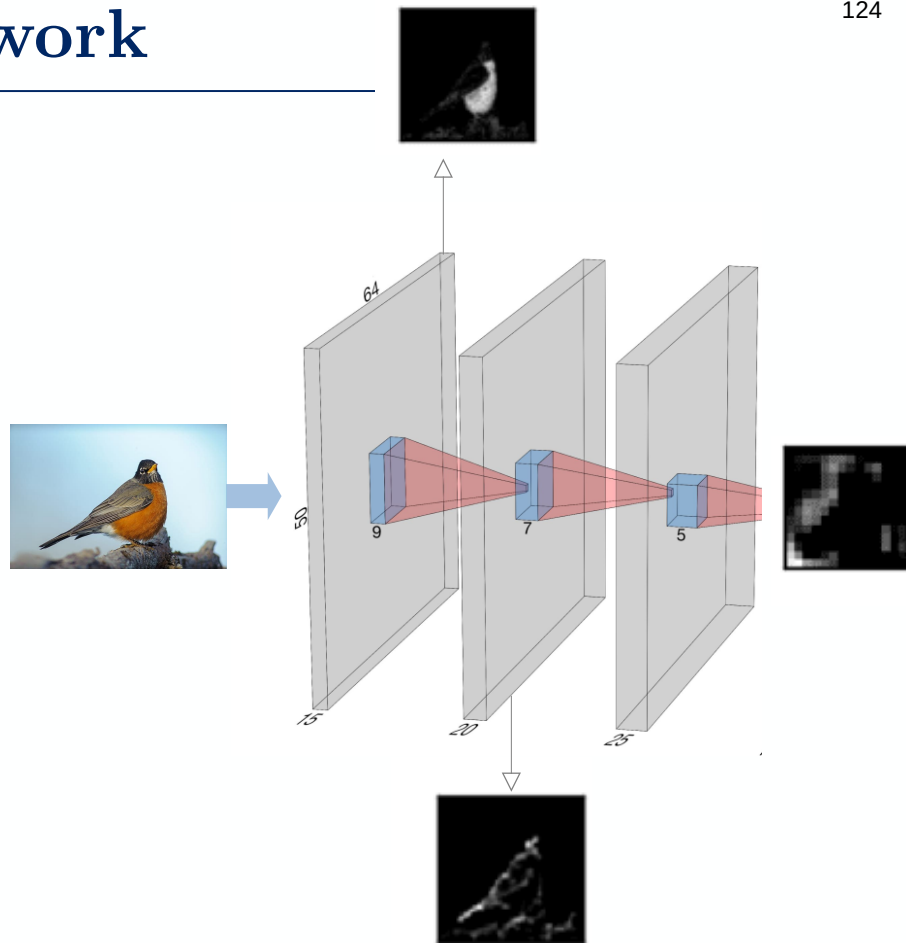
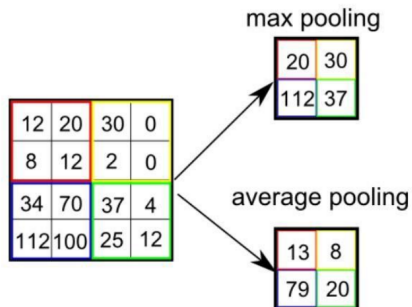
- Formally the kernel convolution operation:

$$S(i, j) = (I * K)(i, j) = \sum_m \sum_n I(m, n) K(i - m, j - n)$$

- Salient feature map 'S' is by default smaller:

$$S \in \mathbb{R}^{(H-K+2P/S+1) \times (W-K+2P/S+1)}$$

- Compressing image size further to reduce input dimension by **max/average pooling**



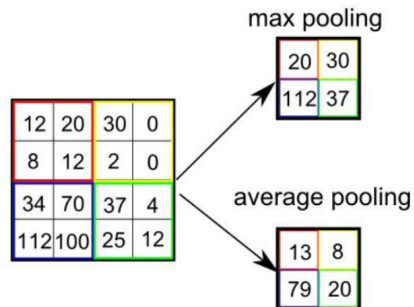
- Formally the kernel convolution operation:

$$S(i, j) = (I * K)(i, j) = \sum_m \sum_n I(m, n) K(i - m, j - n)$$

- Salient feature map 'S' is by default smaller:

$$S \in \mathbb{R}^{(H-K+2P/S+1) \times (W-K+2P/S+1)}$$

- Compressing image size further to reduce input dimension by **max/average pooling**



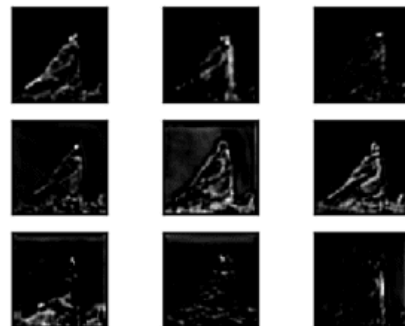
CNN Layer 1

x6 Kernels



CNN Layer 2

x6 Kernels



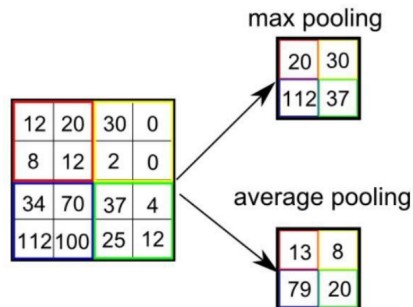
- Formally the kernel convolution operation:

$$S(i, j) = (I * K)(i, j) = \sum_m \sum_n I(m, n) K(i - m, j - n)$$

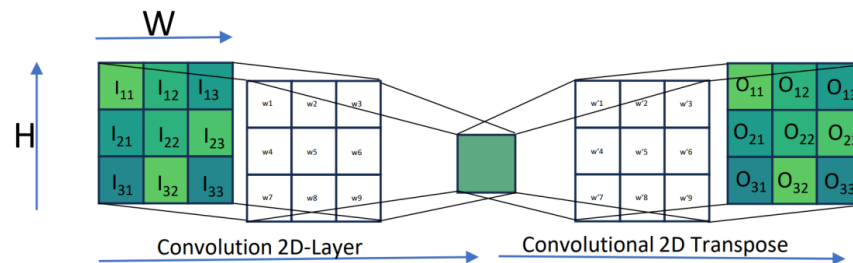
- Salient feature map 'S' is by default smaller:

$$S \in \mathbb{R}^{(H-K+2P/S+1) \times (W-K+2P/S+1)}$$

- Compressing image size further to reduce input dimension by **max/average pooling**



- The **inverse transpose operation** expands the image:



Neural Zoo Transformers



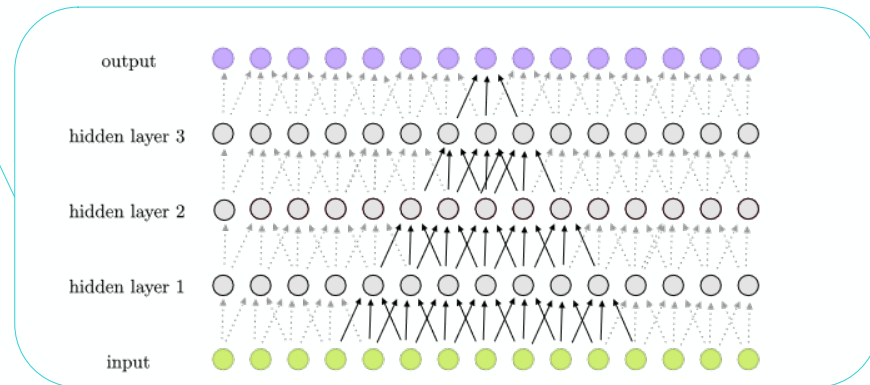
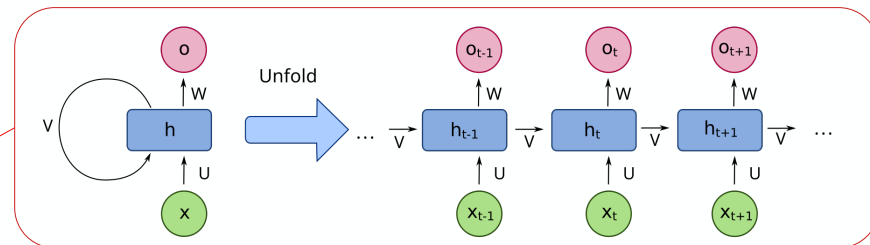
- **Before Transformers:**

- **Recurrent Neural Networks:**

Sequential models utilise recurrent connections, in which the output of the network is passed from time step t to $t+1$ via a recurrent unit – great sequences

- **Convolution Neural Networks:**

Model that utilises kernel filters to learn ‘salient’ features by convoluting nearby activity into increasingly abstract features – great for local information extraction in sequences



Transformers

- **Before Transformers:**

- **Recurrent Neural Networks:**

Sequential models utilise recurrent connections, in which the output of the network is passed from time step t to $t+1$ via a recurrent unit – great sequences

- **Convolution Neural Networks:**

Model that utilises kernel filters to learn ‘salient’ features by convoluting nearby activity into increasingly abstract features – great for local information extraction in sequences

- **Transformers:**

- Introduced **attention mechanisms**
 - Provides information about all positions **simultaneously**
 - Great for **sequences**, and geometrical data such as **images**

Provided proper attribution is provided, Google hereby grants permission to reproduce the tables and figures in this paper solely for use in journalistic or scholarly works.

Attention Is All You Need

Ashish Vaswani [*] Google Brain avaswani@google.com	Noam Shazeer [*] Google Brain noam@google.com	Niki Parmar [*] Google Research nikip@google.com	Jakob Uszkoreit [*] Google Research usz@google.com
Llion Jones [*] Google Research llion@google.com	Aidan N. Gomez [†] University of Toronto aidan@cs.toronto.edu	Lukas Kaiser [*] Google Brain lukaszkaizer@google.com	
Illia Polosukhin [†] illia.polosukhin@gmail.com			

Source: [arXiv:1706.03762](https://arxiv.org/abs/1706.03762)

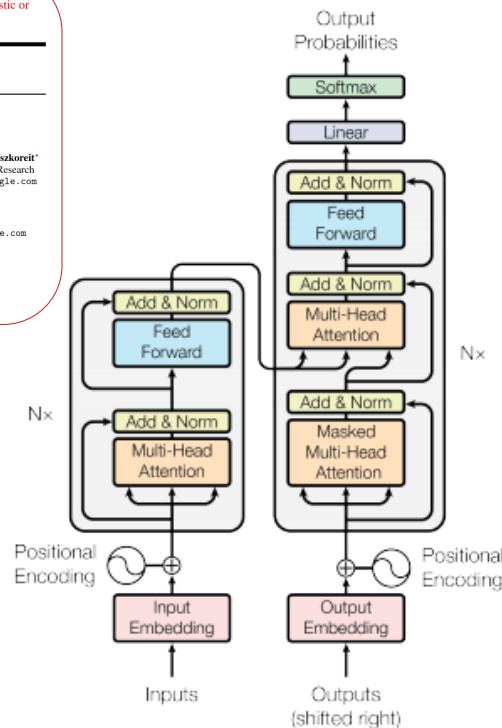


Figure 1: The Transformer - model architecture.

Context Matters – Global Attention



The distribution of pixel values
can be generated to reflect a
conditional probability:

$$p(x_i|\mathcal{C} = \text{'Small cute animal with wide eyes'})$$



Context Matters – Global Attention



The distribution of pixel values
can be generated to reflect a
conditional probability:

$$p(x_i|\mathcal{C} = \text{'Small cute animal with wide eyes'})$$



But what we want is each pixel
to be dependent on the context
and the rest of the pixel values:

$$p(x_i|\mathcal{C} = \text{'Small cute...'}, \sum_{j \neq i}^{H \times W} x_j)$$



Transformers – Scaled dot-product attention

- **Advantages of Transformers (physics focus):**

- Handle variable length data
- Capture **long-range dependencies** in data
- Permutation/order invariant
- Highly parallelisable

- **Attention Mechanism:**

- The success of transformers resides at *first order* in the introduction of the attention mechanism
- Understanding attention needs you to understand the idea of:
 - **vector embeddings**
 - **attention**

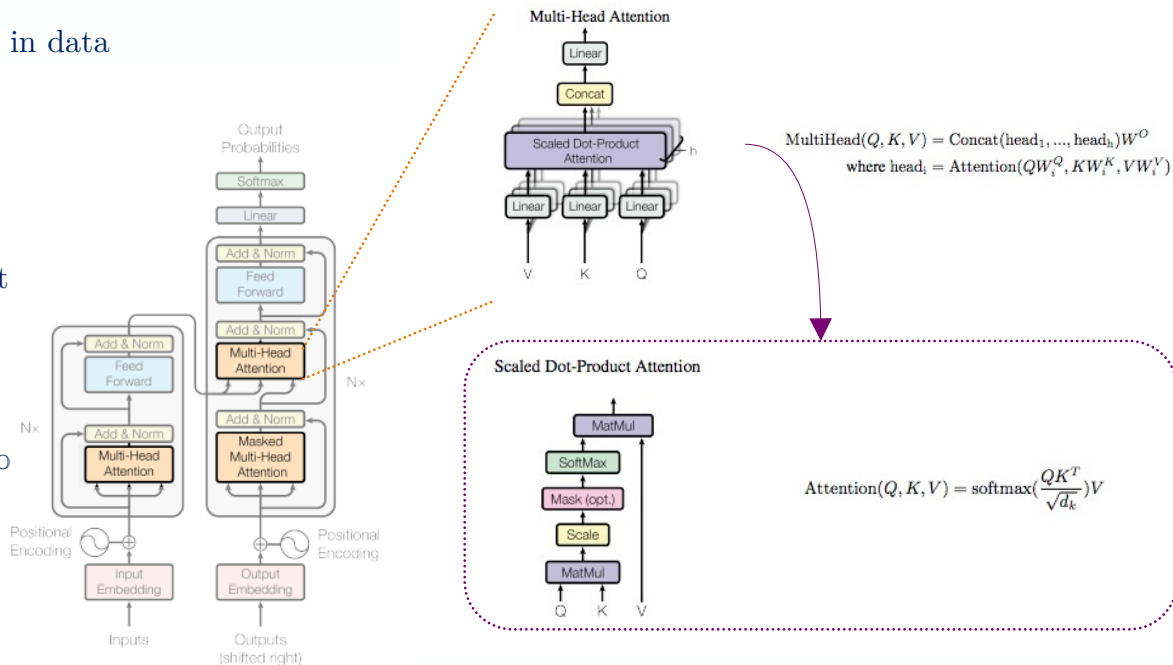


Figure 1: The Transformer - model architecture.

Transformers – Vector Embeddings

- Transformers act on vectors:
 - Tokenisation+embedding** turns each element of the sequence into a vector that resides in a vector space:

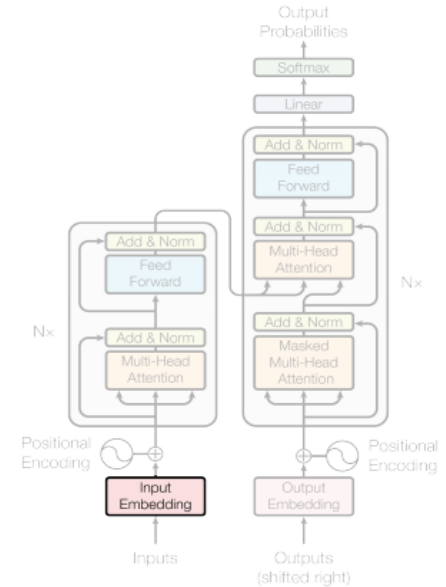
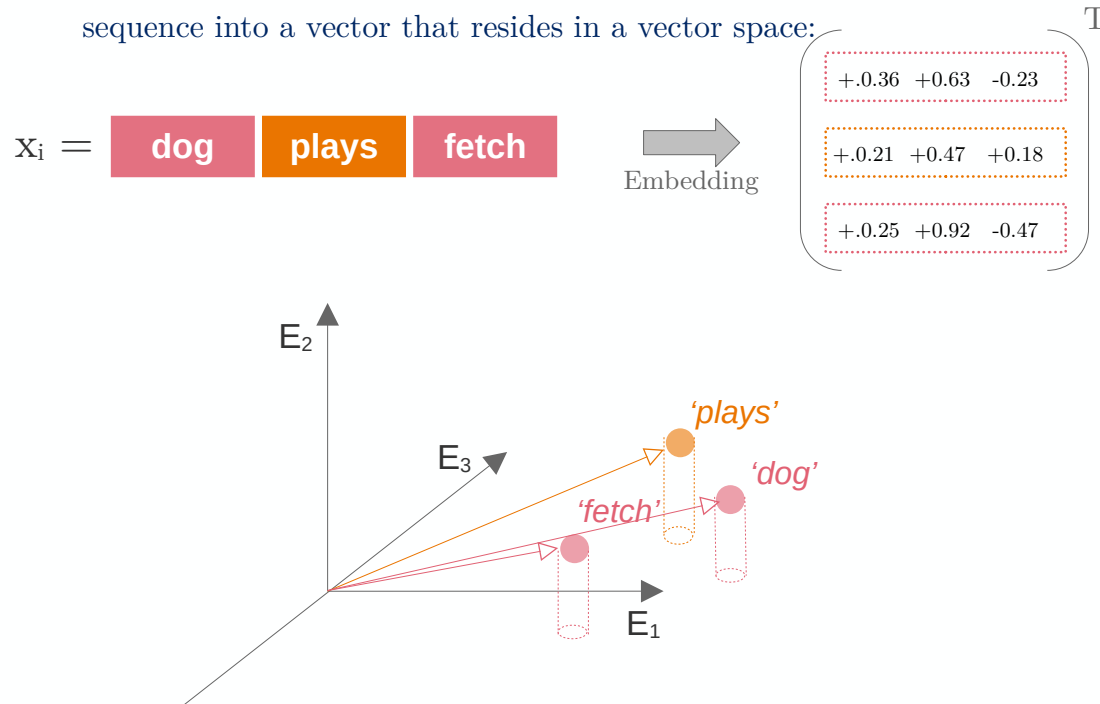


Figure 1: The Transformer - model architecture.

Transformers – Vector Embeddings

- Transformers act on vectors:
 - Tokenisation+embedding** turns each element of the sequence into a vector that resides in a vector space:

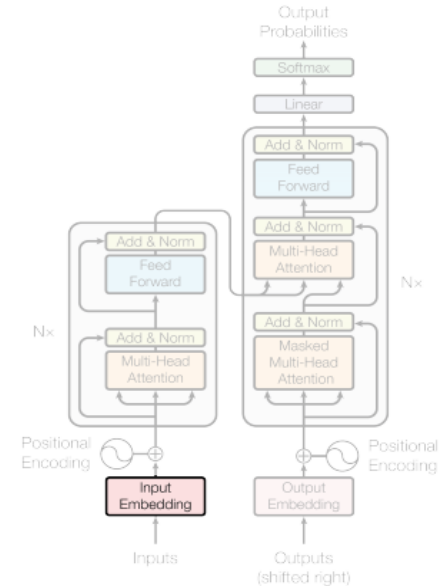
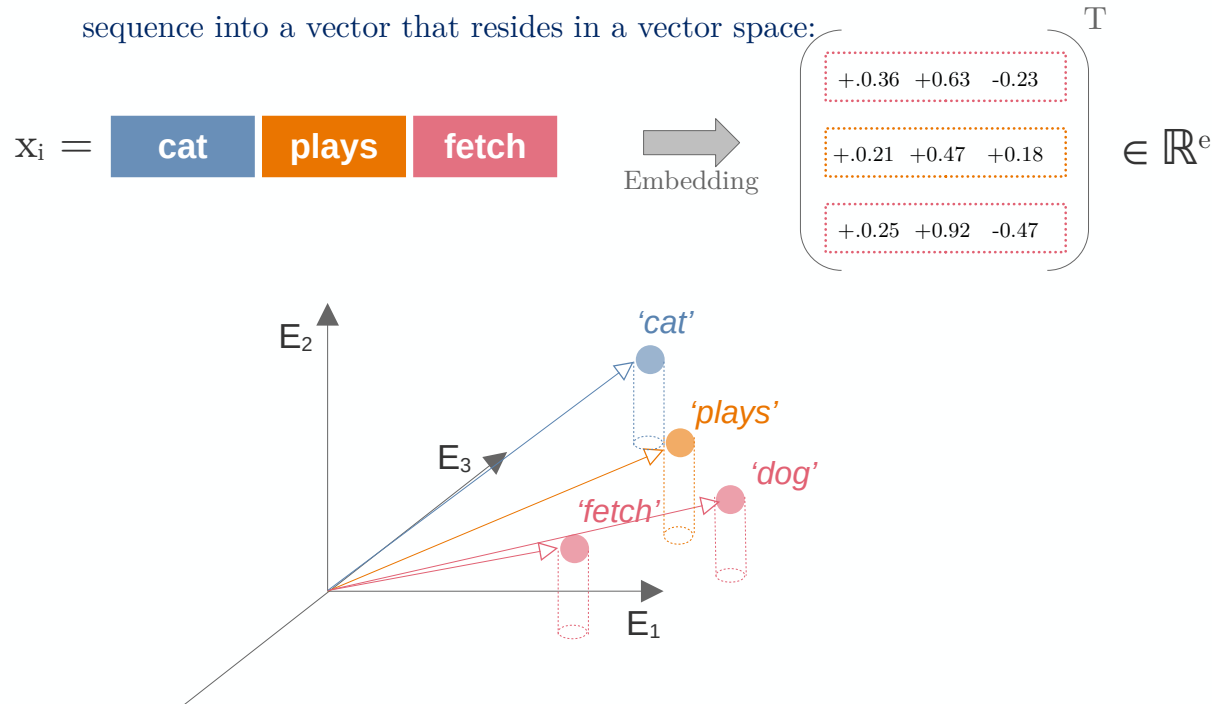
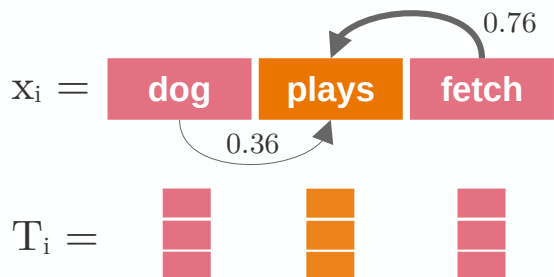


Figure 1: The Transformer - model architecture.

- **Attention = Dynamic flow of information**

- Each **embedded token** stores which token matters and by how much



- **Scaled-dot product** attention is one type of attention that achieves this goal:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

- Three key ingredients:
 - **Query (Q)** : What am I looking for?
 - **Key (K)** : What do I have to offer?
 - **Value (V)** : What do I share if picked?

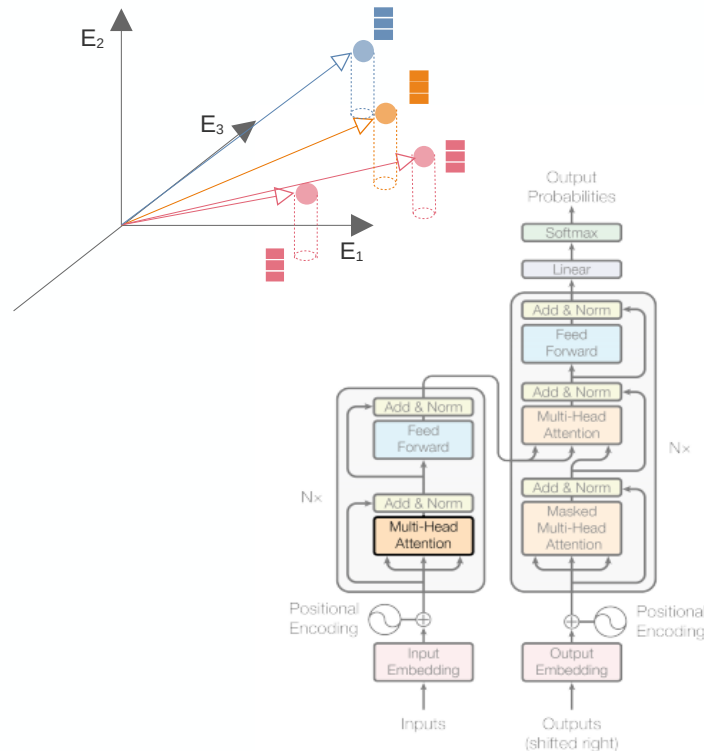
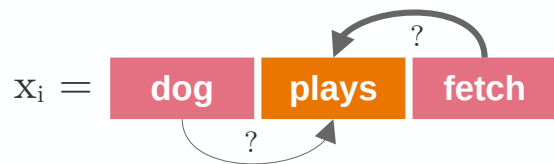


Figure 1: The Transformer - model architecture.

- Lets see how this works for a simple case:



- 1) Define scaled-dot product attention weights for optimisation:

- W_Q : Matrix of weights for queries, shape = $[n, e]$



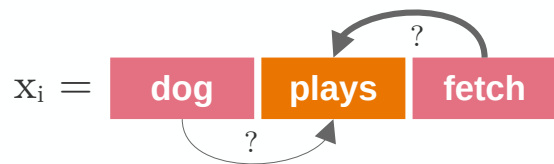
- W_K : Matrix of weights for keys, shape = $[n, e]$



- W_V : Matrix of weights for values, shape = $[n, e]$



- Lets see how this works for a simple case:

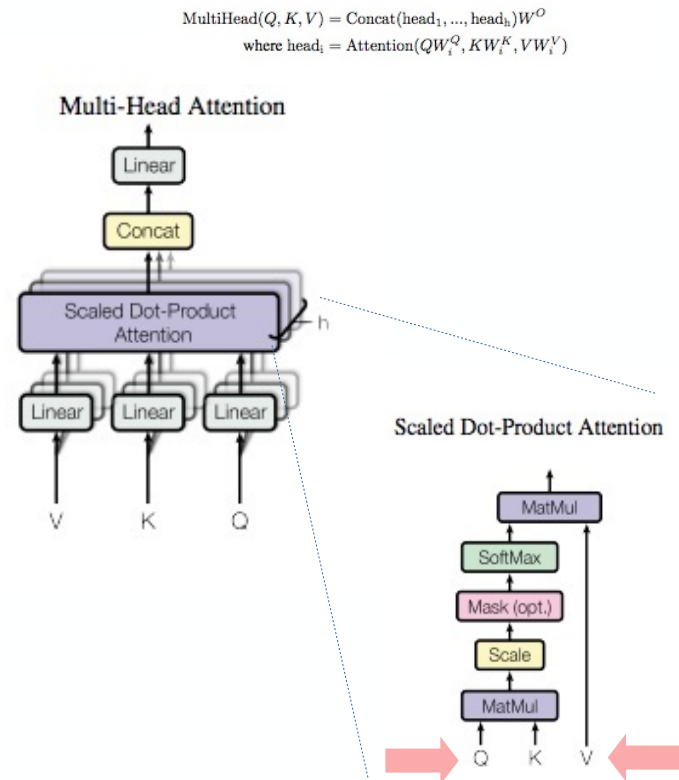


- 2) Compute the query, key and value triplet (**Q,K,V**) for ‘plays’

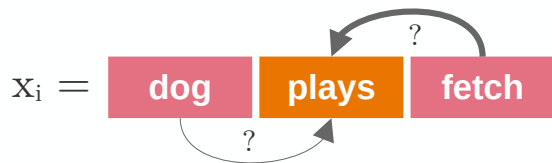
- $\mathbf{Q}_{\text{plays}} : W_Q \cdot T_p =$  \cdot  $=$  $[1, e]$

- $\mathbf{K}_{d(f)} : W_K \cdot T_{d(f)} =$  \cdot  $=$  $[1, e]$

- $\mathbf{V}_{d(f)} : W_V \cdot T_{d(f)} =$  \cdot  $=$  $[1, e]$



- Lets see how this works for a simple case:



- 2) Compute the **similarity** between ‘dog(fetch)’ to ‘play’

$$S_{\text{plays-dog}} : Q_{\text{plays}} \cdot K_d^T = \begin{bmatrix} \text{orange} \\ \text{orange} \\ \text{orange} \end{bmatrix} \cdot \begin{bmatrix} \text{pink} \\ \text{pink} \\ \text{pink} \end{bmatrix}^T = \begin{bmatrix} \text{green} \end{bmatrix} [1,1]$$

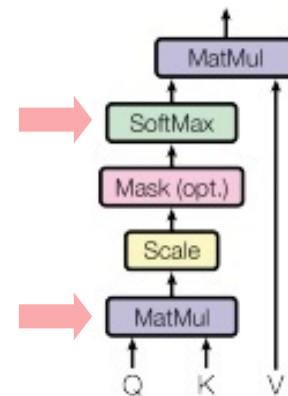
$$S_{\text{plays-fetch}} : Q_{\text{plays}} \cdot K_f^T = \begin{bmatrix} \text{orange} \\ \text{orange} \\ \text{orange} \end{bmatrix} \cdot \begin{bmatrix} \text{pink} \\ \text{pink} \\ \text{pink} \end{bmatrix}^T = \begin{bmatrix} \text{green} \end{bmatrix} [1,1]$$

- 3) Convert to **Similarity** to ‘scaled weights’:

$$a^p_{d(f)} : \sigma(S_{\text{plays-dog(fetch)}}) = \sigma(\begin{bmatrix} \text{green} \end{bmatrix}) = \begin{bmatrix} \text{dashed green} \end{bmatrix}$$

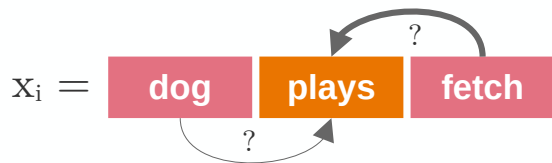
Softmax forces values to be in range $[0,1]$

Scaled Dot-Product Attention



Transformers – Visualising Attention

- Lets see how this works for a simple case:



- 4) Form the context of 'plays' relative to 'dog' & 'fetch':

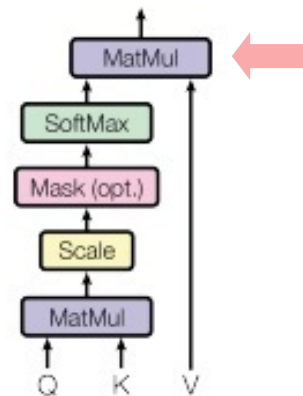
$$C_{\text{plays-dog}} : a_d^p \cdot V_d = \begin{bmatrix} \text{yellow} \\ \text{dashed} \end{bmatrix} \cdot \begin{bmatrix} \text{red} \\ \text{red} \\ \text{red} \end{bmatrix} = \begin{bmatrix} \text{grey} \\ \text{grey} \\ \text{grey} \end{bmatrix} [1, e]$$

$$C_{\text{plays-fetch}} : a_f^p \cdot V_f = \begin{bmatrix} \text{yellow} \\ \text{dashed} \end{bmatrix} \cdot \begin{bmatrix} \text{red} \\ \text{red} \\ \text{red} \end{bmatrix} = \begin{bmatrix} \text{grey} \\ \text{grey} \\ \text{grey} \end{bmatrix} [1, e]$$

- 5) What is the total context of 'plays'?

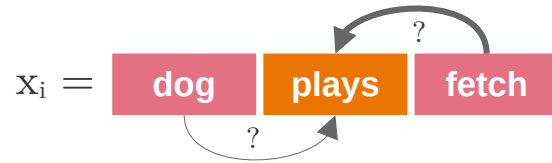
$$C_{\text{plays}} : C_{\text{plays-dog}} + C_{\text{plays-fetch}} = \begin{bmatrix} \text{grey} \\ \text{grey} \\ \text{grey} \end{bmatrix} + \begin{bmatrix} \text{grey} \\ \text{grey} \\ \text{grey} \end{bmatrix} = \begin{bmatrix} \text{green} \\ \text{green} \\ \text{green} \end{bmatrix} [1, e]$$

Scaled Dot-Product Attention



Transformers – Visualising Attention

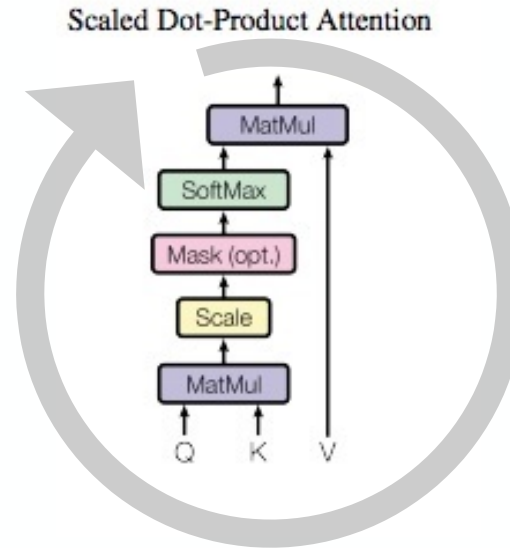
- Lets see how this works for a simple case:



- 6) Repeat now for 'dog' and 'fetch':

$$C = \begin{bmatrix} C_d & C_p & C_f \\ \text{[green box]} & \text{[green box]} & \text{[green box]} \\ \text{[green box]} & \text{[green box]} & \text{[green box]} \\ \text{[green box]} & \text{[green box]} & \text{[green box]} \end{bmatrix} [n,e]$$

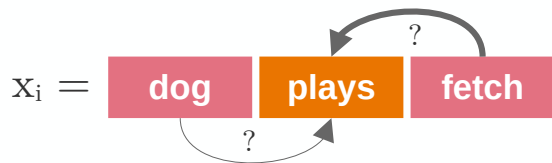
Diagram illustrating the attention matrix C with dimensions $[n,e]$. The matrix is composed of three columns labeled C_d , C_p , and C_f , each containing three green boxes representing attention weights.



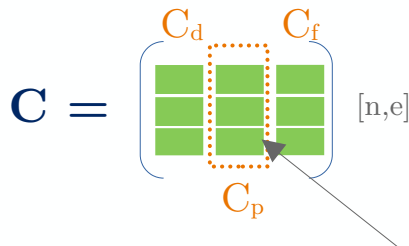
Transformers – Visualising Attention

141

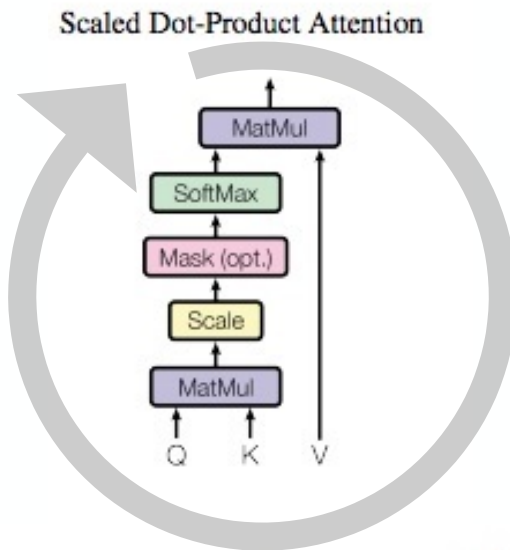
- Lets see how this works for a simple case:



- 6) Repeat now for 'dog' and 'fetch':



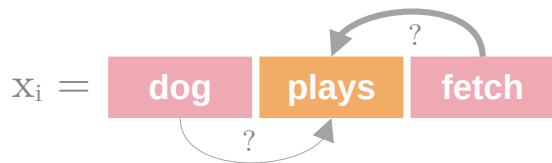
What is the context of all embedding vectors to the word 'play'?



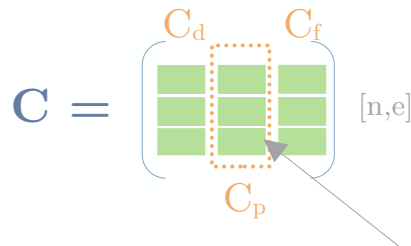
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Transformers – Visualising Attention

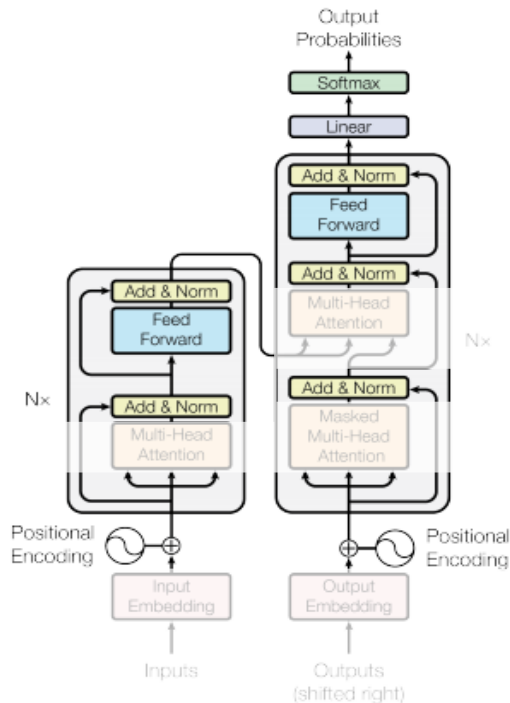
- Lets see how this works for a simple case



- 6) Repeat now for 'dog' and 'fetch':



What is the context of all embedding vectors to the word 'play'?



The rest of the transformer is actually quite simple – just MLPs and adding/normalisation

Figure 1: The Transformer - model architecture.



The End
Thank you!

- **Books:**

- [1] I. Goodfellow, Y. Bengio, and A. Courville, '*Deep Learning*', MIT Press, 2016, <http://www.deeplearningbook.org>
- [2] Simon J.D. Prince, '*Understanding Deep Learning*', MIT Press, 2023, <http://udlbook.com>
- [3] Kevin P. Murphy, "Probabilistic Machine Learning: An introduction", MIT Press, 2022, probml.ai

- **Lectures:**

- [1] ATLAS-D 2023, F. Meloni, <https://indico.cern.ch/event/1263122/>
- [2] G. Louppe, Info8010 Deep Learning, 2024, <https://github.com/glouppe/info8010-deep-learning>

- **Misc. :**

- [1] Kaare Petersen, Michael Pedersen, '*Matrix Cookbook*', 2012, <https://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf>



Backup

Point Estimators: Technicality of continuous variables

- **Covariance:** Linear correlation between two variables:

$$\text{Cov}(f(x), g(y)) = \mathbb{E}[(f(x) - \mathbb{E}[f(x)])(g(y) - \mathbb{E}[g(y)])]$$

- **Deterministically Correlated Random Variables:**

→ Two random variables X & Y , such that $y = g(x)$:

$$p_x(x) = p_y(g(x)) \left| \frac{\partial g(x)}{\partial x} \right|$$

→ In higher dimensions \mathbf{x} & \mathbf{y} , for $\mathbf{x} = g(\mathbf{y})$:

$$p_x(\mathbf{x}) = p_y(g(\mathbf{x})) \left| \det \left(\frac{\partial g(\mathbf{x})}{\partial \mathbf{x}} \right) \right|$$

Point Estimators: Technicality of continuous variables

- **Covariance:** Linear correlation between two variables:

$$\text{Cov}(f(x), g(y)) = \mathbb{E}[(f(x) - \mathbb{E}[f(x)])(g(y) - \mathbb{E}[g(y)])]$$

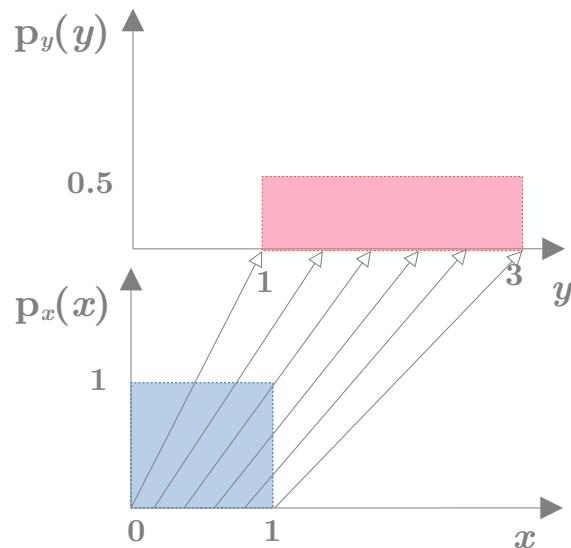
- **Deterministically Correlated Random Variables:**

→ Two random variables X & Y , such that $y = g(x)$:

$$p_x(x) = p_y(g(x)) \left| \frac{\partial g(x)}{\partial x} \right|$$

→ In higher dimensions \mathbf{x} & \mathbf{y} , for $\mathbf{x} = g(\mathbf{y})$:

$$p_x(\mathbf{x}) = p_y(g(\mathbf{x})) \left| \det \left(\frac{\partial g(\mathbf{x})}{\partial \mathbf{x}} \right) \right|$$



Estimator Properties

- **Point Estimator, or statistic**, is any function of the data that infers from the data some parameter of interest θ :

$$\hat{\theta} = g(\{x\}_m)$$

- **Bias** of an estimator is given by:

$$\text{bias}(\hat{\theta}) = \mathbb{E}(\hat{\theta}) - \theta$$

- **Variance** of an estimator:

$$\text{Var}(\hat{\theta})$$

- **Example – Gaussian:**

Sample $\{x\}_m$ generated by a Gaussian *pdf*:

$$\hat{\sigma}^2 = \frac{1}{m} \sum \left[(x_i - \hat{\mu})^2 \right] \quad \leftarrow \quad p(x_i | \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{1}{2} \frac{(x_i - \mu)^2}{\sigma^2} \right]$$

→ Bias calculation:

$$\text{bias}(\hat{\sigma}^2) = \underbrace{\mathbb{E}(\hat{\sigma}^2)} - \sigma^2 = -\frac{\sigma^2}{m}$$

$$\begin{aligned} \mathbb{E}(\hat{\sigma}^2) &= \mathbb{E} \left[\frac{1}{m} \sum (x_i - \hat{\mu})^2 \right] \\ &= \frac{m-1}{m} \sigma^2 \end{aligned}$$

Estimator Properties

- **Point Estimator, or statistic**, is any function of the data that infers from the data some parameter of interest θ :

$$\hat{\theta} = g(\{x\}_m)$$

- **Bias** of an estimator is given by:

$$\text{bias}(\hat{\theta}) = \mathbb{E}(\hat{\theta}) - \theta$$

- **Variance** of an estimator:

$$\text{Var}(\hat{\theta})$$

- **Example – Gaussian:**

Sample $\{x\}_m$ generated by a Gaussian *pdf*:

$$\hat{\sigma}^2 = \frac{1}{m-1} \sum [(x_i - \hat{\mu})^2] \quad \leftarrow \quad p(x_i | \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left[-\frac{1}{2} \frac{(x_i - \mu)^2}{\sigma^2} \right]$$

→ Bias calculation:

$$\text{bias}(\hat{\sigma}^2) = \mathbb{E}(\hat{\sigma}^2) - \sigma^2 = 0$$

Estimator Properties

- **Point Estimator, or statistic**, is any function of the data that infers from the data some parameter of interest θ :

$$\hat{\theta} = g(\{x\}_m)$$

- **Bias** of an estimator is given by:

$$\text{bias}(\hat{\theta}) = \mathbb{E}(\hat{\theta}) - \theta$$

- **Variance** of an estimator:

$$\text{Var}(\hat{\theta})$$

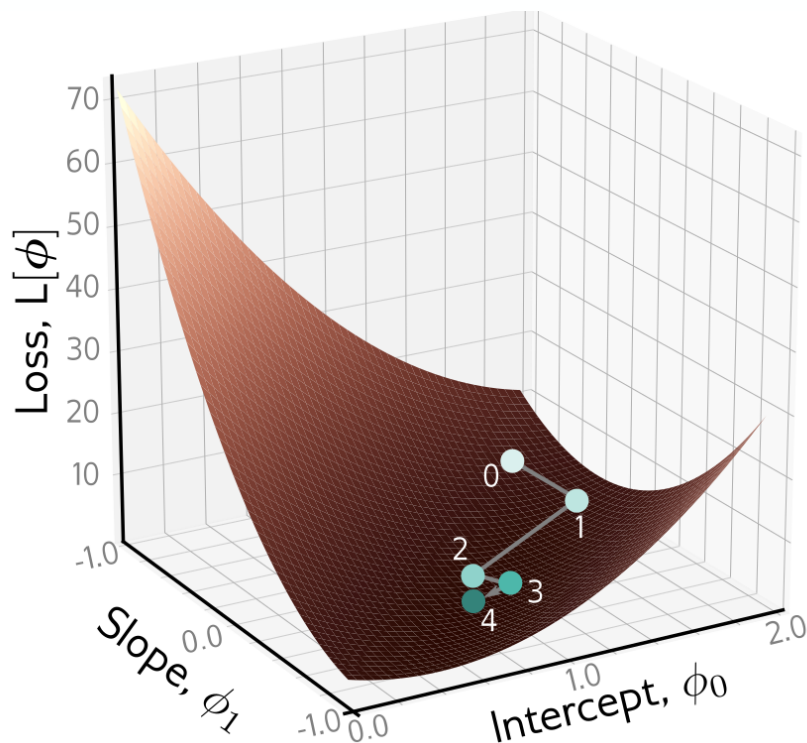
- **Example – Gaussian:**

Sample $\{x\}_m$ generated by a Gaussian *pdf*:

$$\hat{\mu} = \frac{1}{m} \sum_i^m x_i \quad \leftarrow \quad p(x_i|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{1}{2} \frac{(x_i - \mu)^2}{\sigma^2} \right]$$

→ Variance calculation:

$$\begin{aligned} \text{Var}(\hat{\mu}) &= \text{Var}\left(\frac{1}{m} \sum_i^m x_i\right) = \mathbb{E} \left[\left(\frac{1}{m} \left(\sum_i^m x_i - \mathbb{E}\left(\sum_i^m x_i\right) \right) \right)^2 \right] \\ &= \frac{1}{m^2} \sum_i^m \text{Var}(x_i) = \frac{\sigma^2}{m^2} \end{aligned}$$



- Optimising the model f is achieved by minimising the **loss/cost function**:

$$\hat{\Phi} = \arg \min [\mathcal{L}(f(\mathbf{x}|\Phi), \mathbf{x})]$$

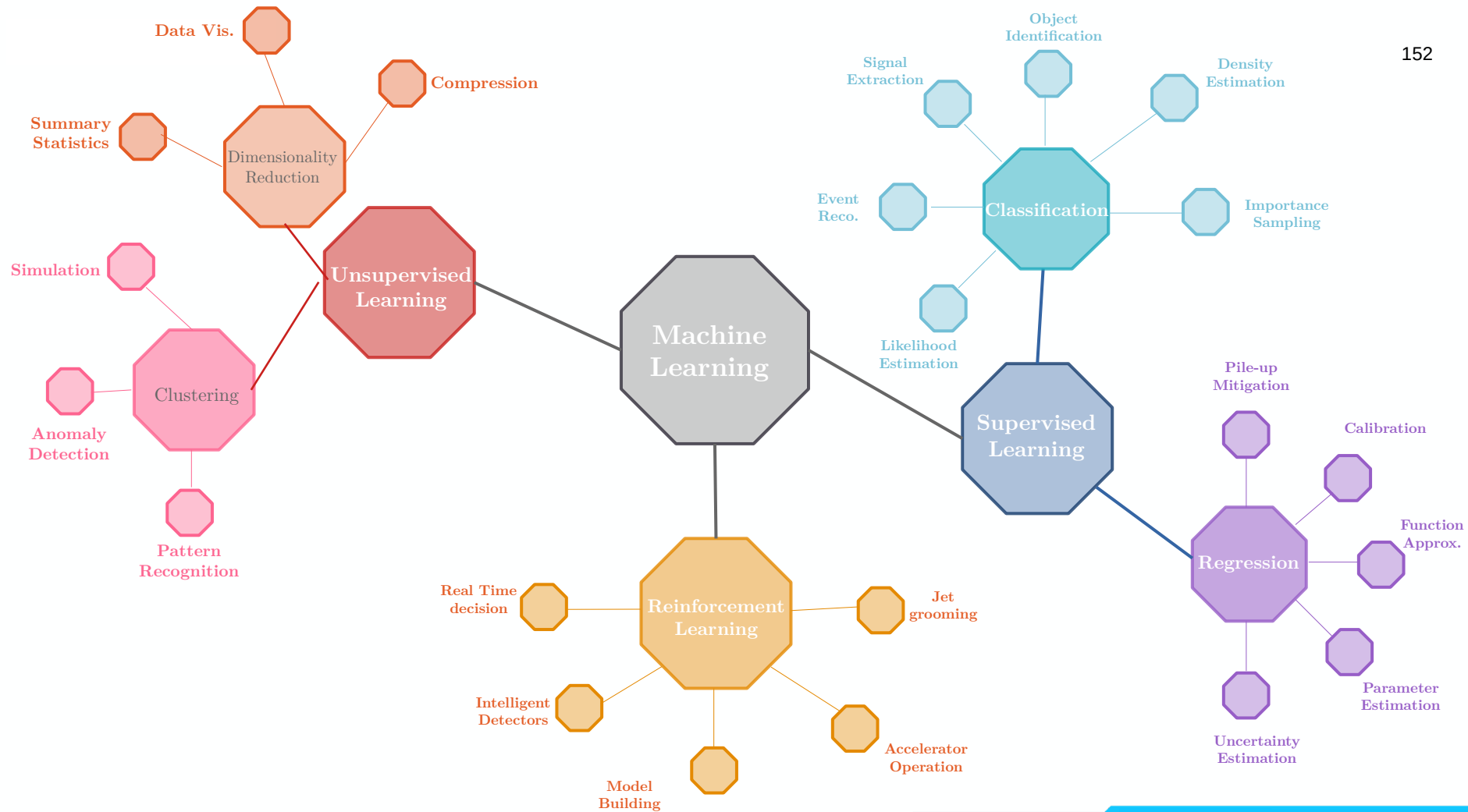
- Gradient descent** has primarily two steps:

1) Calculate gradients:

$$\frac{\partial \mathcal{L}}{\partial \Phi} = \begin{bmatrix} \frac{d\mathcal{L}}{d\phi_0} \\ \vdots \\ \frac{d\mathcal{L}}{d\phi_h} \end{bmatrix}$$

2) Update the parameters in direction that minimises loss:

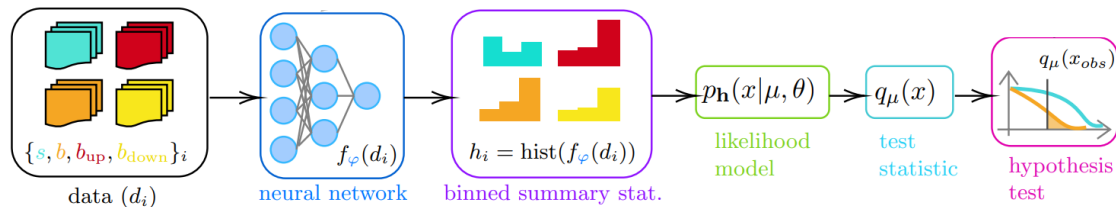
$$\phi_i \leftarrow \phi_i - \alpha \frac{\partial \mathcal{L}}{\partial \phi_i}$$



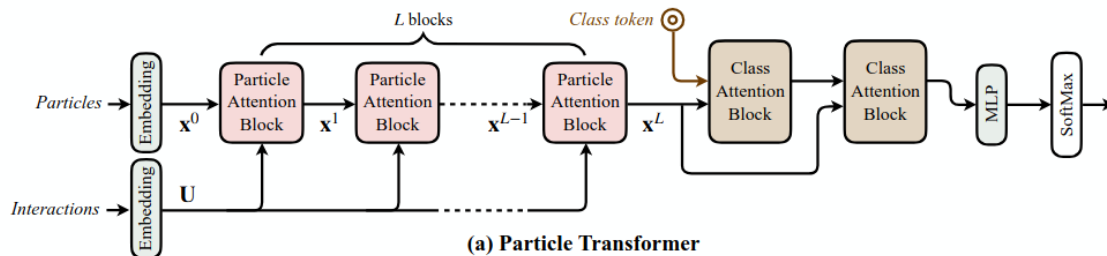
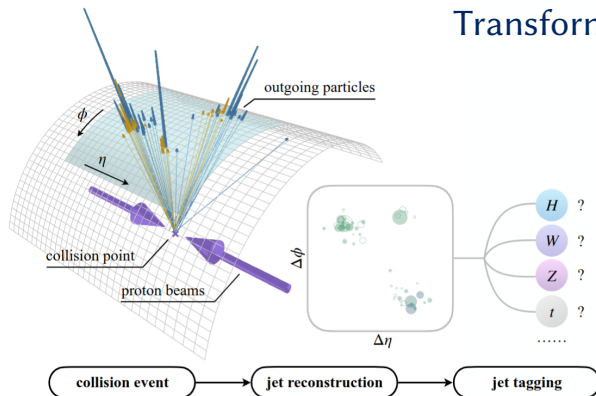
High Energy Physics – ML Developments

153

→ Checkout the [HEP Machine Learning Living Review](#)



Particle Transformer - [arXiv:2202.03772](#) Transformer model for tagging jets at the LHC





Probability Triplet

Can likely skip

Continuous Probability

- **Absolute continuous probability** is the infinitesimal sum^[1] over the **probability density function** $p : \mathbb{R} \rightarrow [0, \infty]$:

$$P(X \in A) = \int_A p(x) dx$$

- **Marginal probability** density:

$$p(x) = \int p(x, y) dy$$

- **Conditional probability** density:

$$p(y|x) = \frac{p(x, y)}{p(x)}$$

[1] *kind of...see uncountable sets*

Continuous Probability

- **Absolute continuous probability** is the infinitesimal sum^[1] over the **probability density function** $p : \mathbb{R} \rightarrow [0, \infty]$:

$$P(X \in A) = \int_A p(x) dx$$

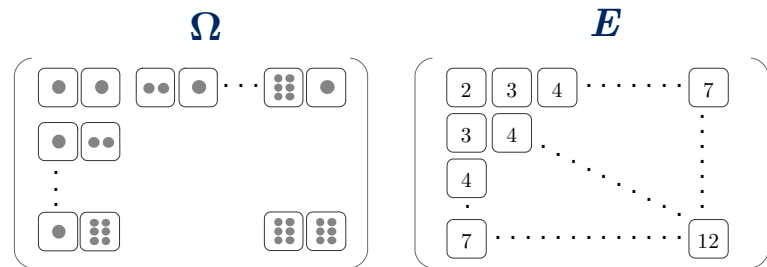
- **Marginal probability** density:

$$p(x) = \int p(x, y) dy$$

- **Conditional probability** density:

$$p(y|x) = \frac{p(x, y)}{p(x)}$$

- **Random Variable:** Measurable function $X : \Omega \rightarrow E$ from sample to measurable space:



- **Realisation:** An instance sampled from the random variable distribution:

$$x \sim X$$

[1] kind of...see uncountable sets

Continuous Probability

- Absolute continuous probability** is the infinitesimal sum^[1] over the **probability density function** $p : \mathbb{R} \rightarrow [0, \infty]$:

$$P(X \in A) = \int_A p(x) dx$$

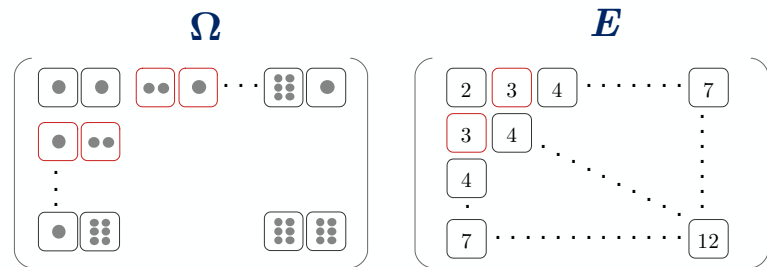
- Marginal probability** density:

$$p(x) = \int p(x, y) dy$$

- Conditional probability** density:

$$p(y|x) = \frac{p(x, y)}{p(x)}$$

- Random Variable:** Measurable function $X : \Omega \rightarrow E$ from sample to measurable space:



- Realisation:** An instance sampled from the random variable distribution:

$$x \sim X$$

[1] kind of...see uncountable sets

Continuous Probability

- Absolute continuous probability** is the infinitesimal sum^[1] over the **probability density function** $p : \mathbb{R} \rightarrow [0, \infty]$:

$$P(X \in A) = \int_A p(x) dx$$

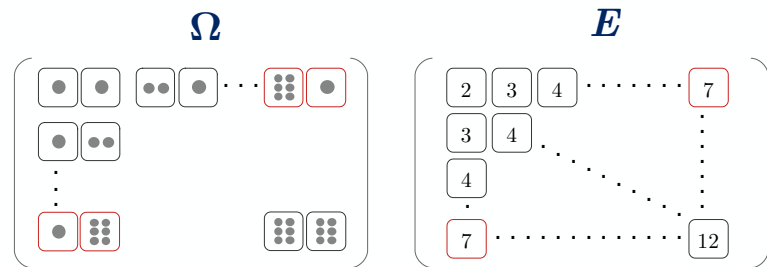
- Marginal probability** density:

$$p(x) = \int p(x, y) dy$$

- Conditional probability** density:

$$p(y|x) = \frac{p(x, y)}{p(x)}$$

- Random Variable:** Measurable function $X : \Omega \rightarrow E$ from sample to measurable space:



- Realisation:** An instance sampled from the random variable distribution:

$$x \sim X$$

[1] kind of...see uncountable sets

Continuous Probability

- Absolute continuous probability** is the infinitesimal sum^[1] over the **probability density function** $p : \mathbb{R} \rightarrow [0, \infty]$:

$$P(X \in A) = \int_A p(x) dx$$

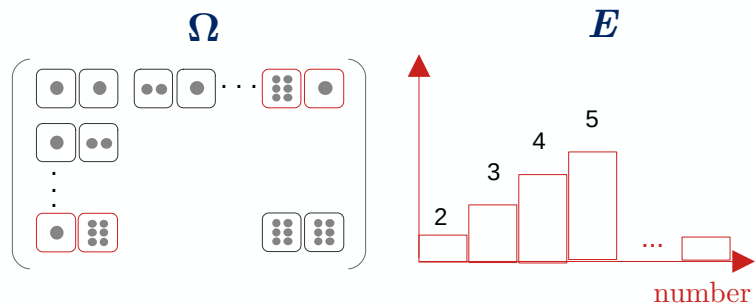
- Marginal probability** density:

$$p(x) = \int p(x, y) dy$$

- Conditional probability** density:

$$p(y|x) = \frac{p(x, y)}{p(x)}$$

- Random Variable:** Measurable function $X : \Omega \rightarrow E$ from sample to measurable space:



- Realisation:** An instance sampled from the random variable distribution:

$$x \sim X$$

[1] kind of...see uncountable sets

Continuous Probability

- Absolute continuous probability** is the infinitesimal sum^[1] over the **probability density function** $p : \mathbb{R} \rightarrow [0, \infty]$:

$$P(X \in A) = \int_A p(x) dx$$

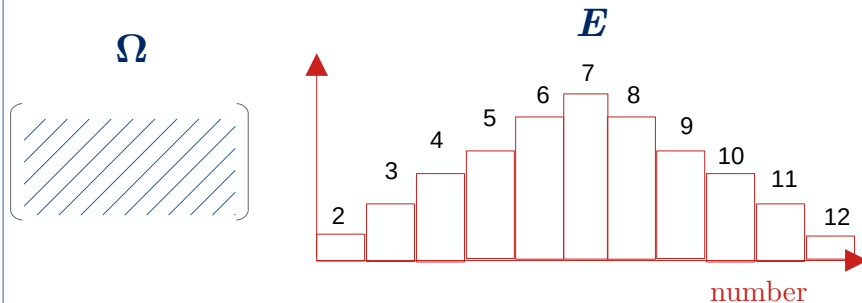
- Marginal probability** density:

$$p(x) = \int p(x, y) dy$$

- Conditional probability** density:

$$p(y|x) = \frac{p(x, y)}{p(x)}$$

- Random Variable:** Measurable function $X : \Omega \rightarrow E$ from sample to measurable space:



- Realisation:** An instance sampled from the random variable distribution:

$$x \sim X$$

[1] kind of...see uncountable sets

Continuous Probability

- Absolute continuous probability** is the infinitesimal sum^[1] over the **probability density function** $p : \mathbb{R} \rightarrow [0, \infty]$:

$$P(X \in A) = \int_A p(x) dx$$

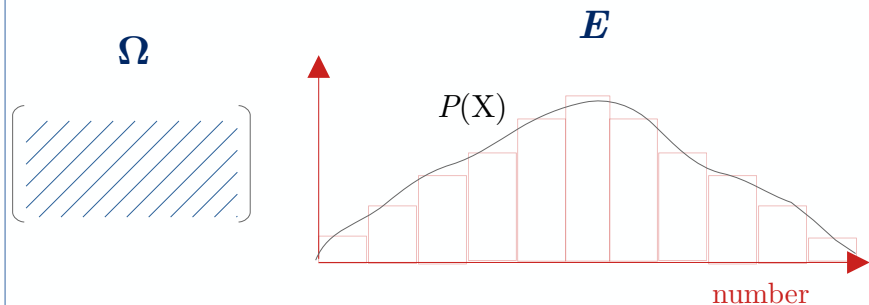
- Marginal probability** density:

$$p(x) = \int p(x, y) dy$$

- Conditional probability** density:

$$p(y|x) = \frac{p(x, y)}{p(x)}$$

- Random Variable:** Measurable function $X : \Omega \rightarrow E$ from sample to measurable space:



- Realisation:** An instance sampled from the random variable distribution:

$$x \sim X$$

[1] kind of...see uncountable sets

Continuous Probability

- **Absolute continuous probability** is the infinitesimal sum^[1] over the **probability density function** $p : \mathbb{R} \rightarrow [0, \infty]$:

$$P(X \in A) = \int_A p(x) dx$$

- **Marginal probability** density:

$$p(x) = \int p(x, y) dy$$

- **Conditional probability** density:

$$p(y|x) = \frac{p(x, y)}{p(x)}$$

- **(Non-)Conditional Independence:**

Two random variables are independent if:

$$p(\mathbf{x} = x, \mathbf{y} = y) = p(\mathbf{x} = x)p(\mathbf{y} = y)$$

$$\forall x \in \mathbf{x}, \forall y \in \mathbf{y}$$

[1] kind of...see uncountable sets

Point Estimators: Expectation, bias, variance

- **Expectation:** Expected value of a function of random variables

$$\mathbb{E}_{x \sim p} [f(x)] = \int p(x) f(x) dx$$

- **Variance:** Variation of samples from a random variable:

$$\text{Var}(f(x)) = \mathbb{E} \left[(f(x) - \mathbb{E}[f(x)])^2 \right]$$

- **Covariance:** Linear correlation between two variables:

$$\text{Cov}(f(x), g(y)) = \mathbb{E} [(f(x) - \mathbb{E}[f(x)])(g(y) - \mathbb{E}[g(y)])]$$