



UNIVERSITY OF
CAMBRIDGE

Machine Learning and the Scientific Principle

Carl Henrik Ek - che29@cam.ac.uk

2nd of February, 2021

<http://carlhenrik.com>

Cyder Country



My Dream!

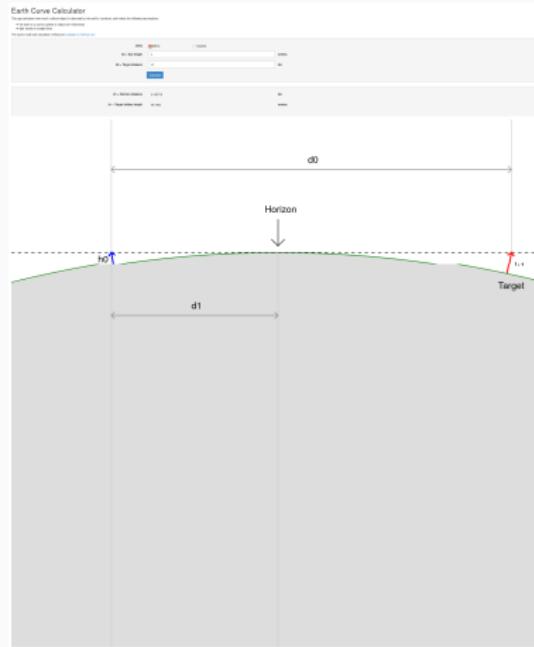








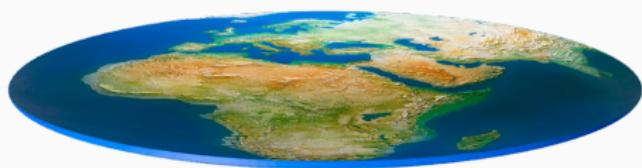




Distance to horizon 6.2km

Hidden height 125.6m



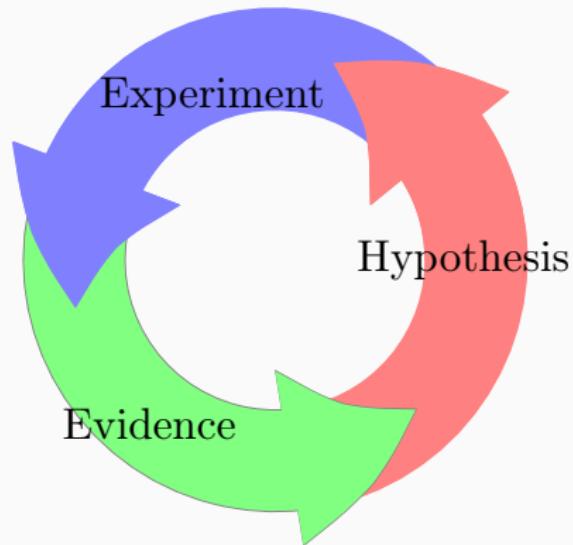




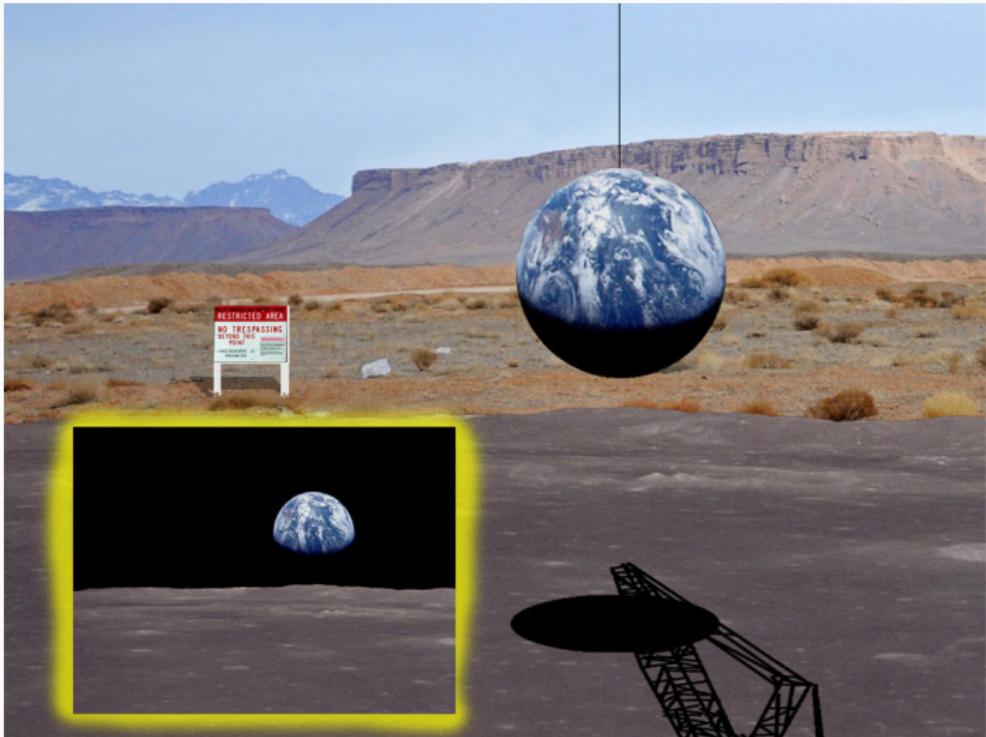
Well it kinda looks like a ball from the moon



Why does this not make sense?



Or not



Laplace Demon



All these efforts in the search for truth tend to lead the human mind back continually to the vast intelligence which we have just mentioned, but from which it will always remain infinitely removed.

– Laplace *Laplace, 1814a*



OCCAM'S RAZOR

... the simplest explanation is usually
the correct one.

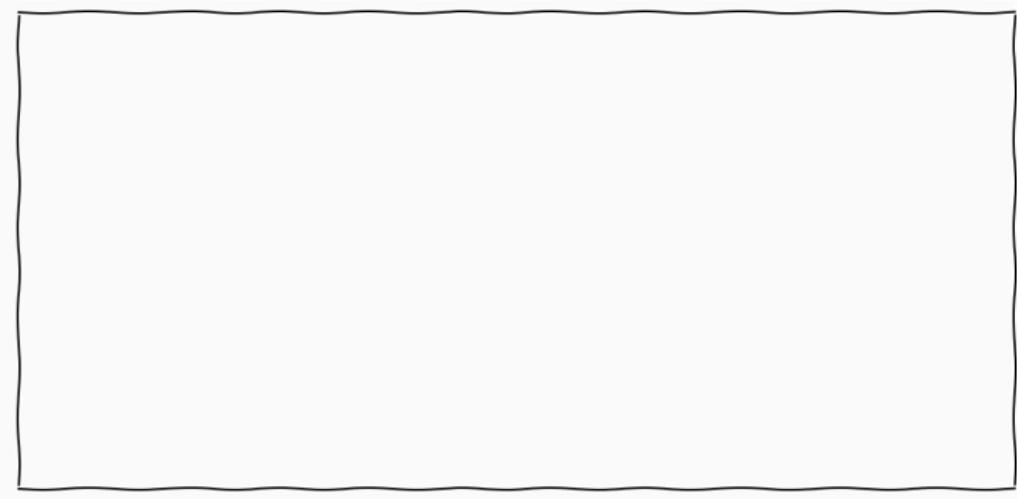
Today

Where does machine learning fit into the scientific workflow?

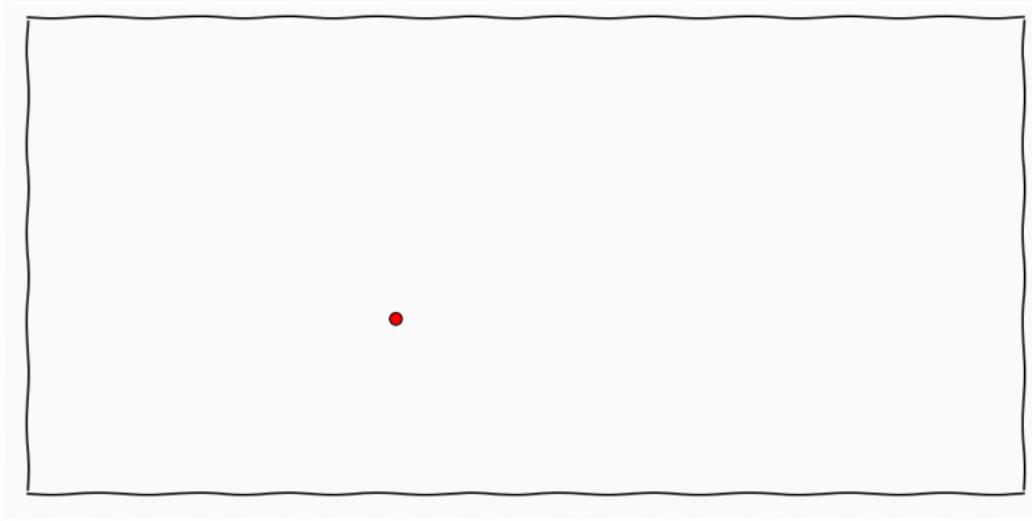
1. How do we implement the scientific principle?
2. How do we implement Occam's razor?

Machine Learning

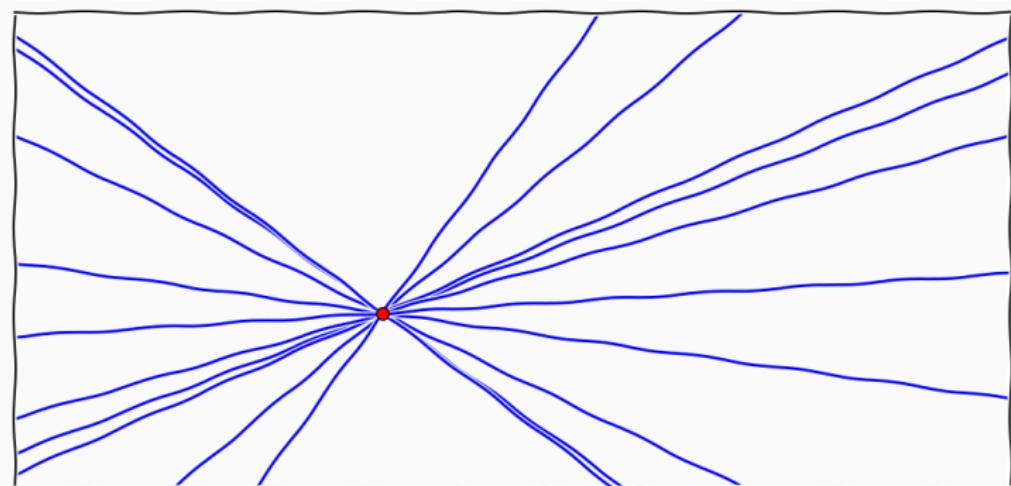
Linear Regression



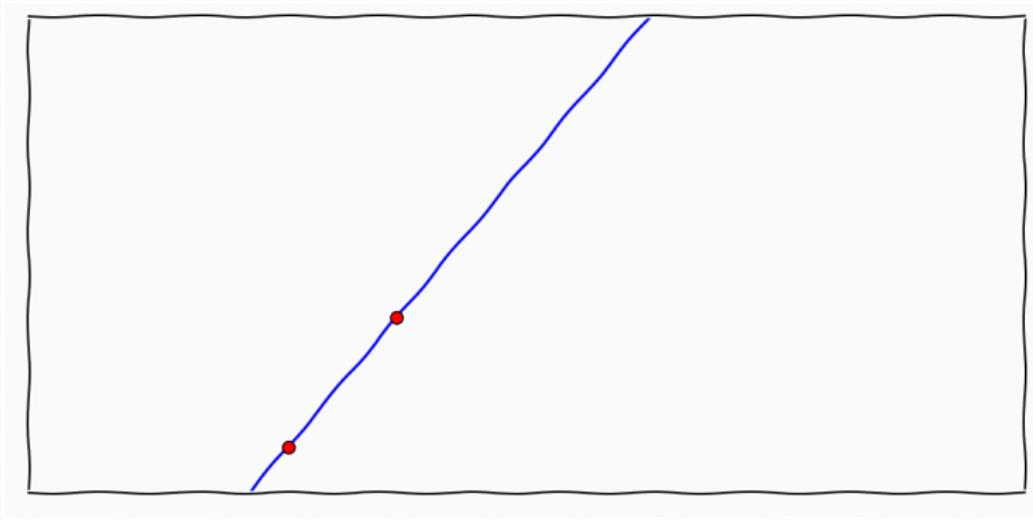
Linear Regression



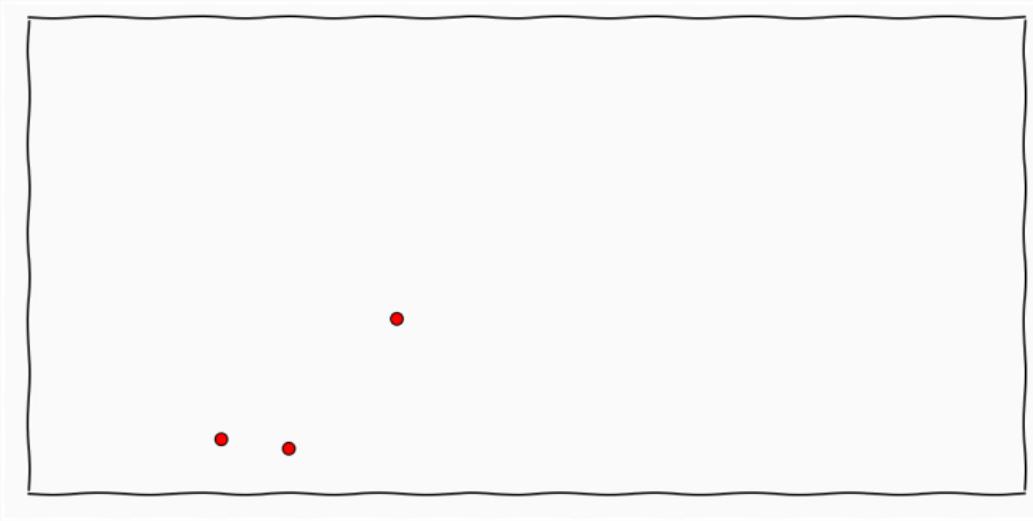
Linear Regression



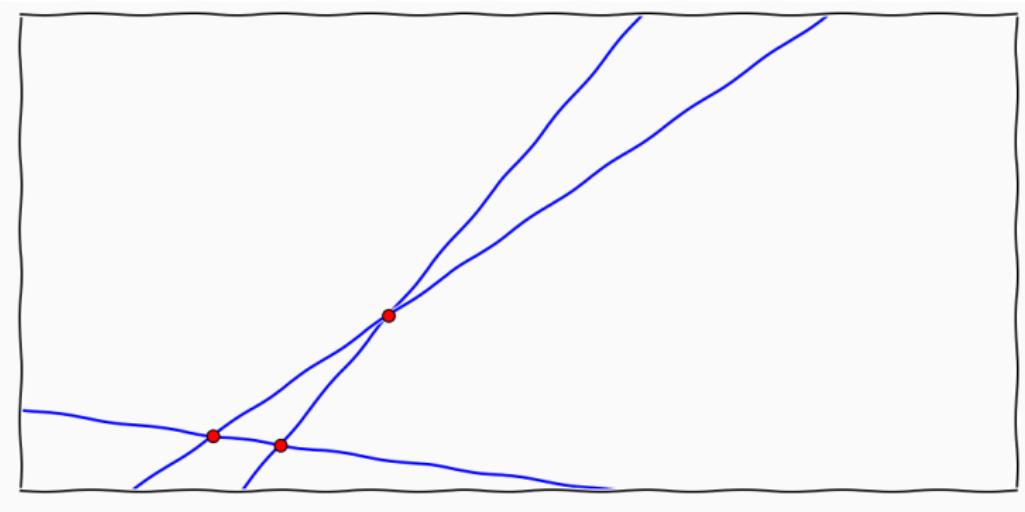
Linear Regression



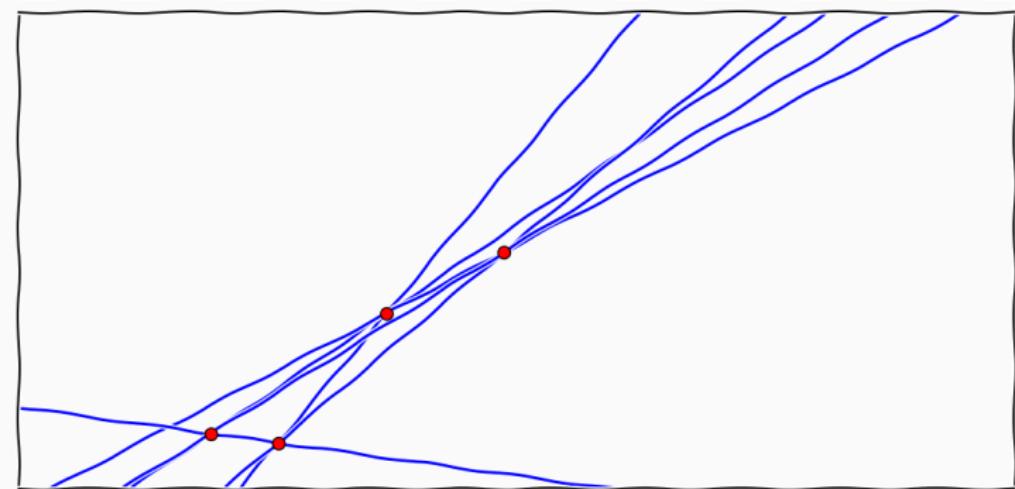
Linear Regression



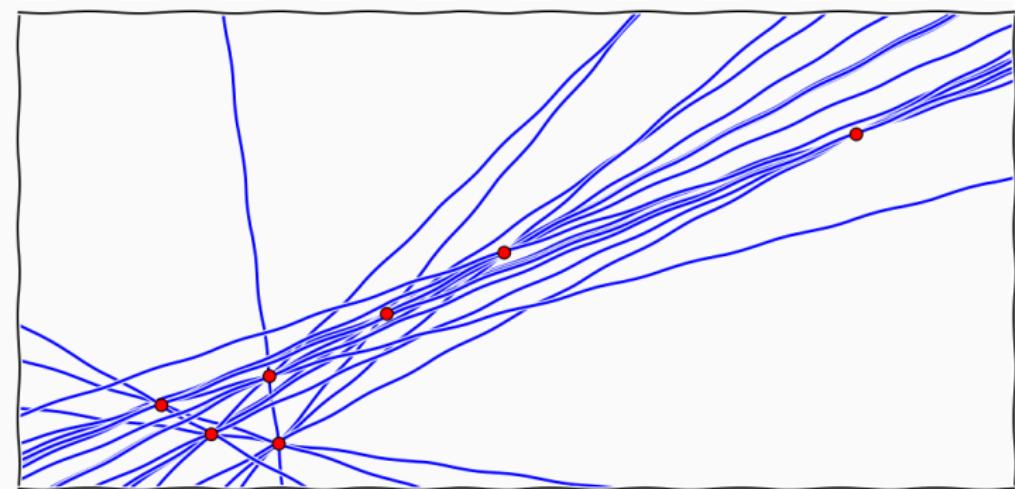
Linear Regression



Linear Regression



Linear Regression



Linear Regression: High School

$$\underbrace{\mathbf{A}}_{m \times n} \underbrace{\mathbf{x}}_{n \times 1} = \underbrace{\mathbf{b}}_{m \times 1}$$

- Over-determined $m > n$

$$\hat{\mathbf{x}} = \underset{\mathbf{x}}{\operatorname{argmin}} \sum_i^m (b_i - \mathbf{A}_{i:} \mathbf{x})^2$$

Linear Regression: High School

$$\underbrace{\mathbf{A}}_{m \times n} \underbrace{\mathbf{x}}_{n \times 1} = \underbrace{\mathbf{b}}_{m \times 1}$$

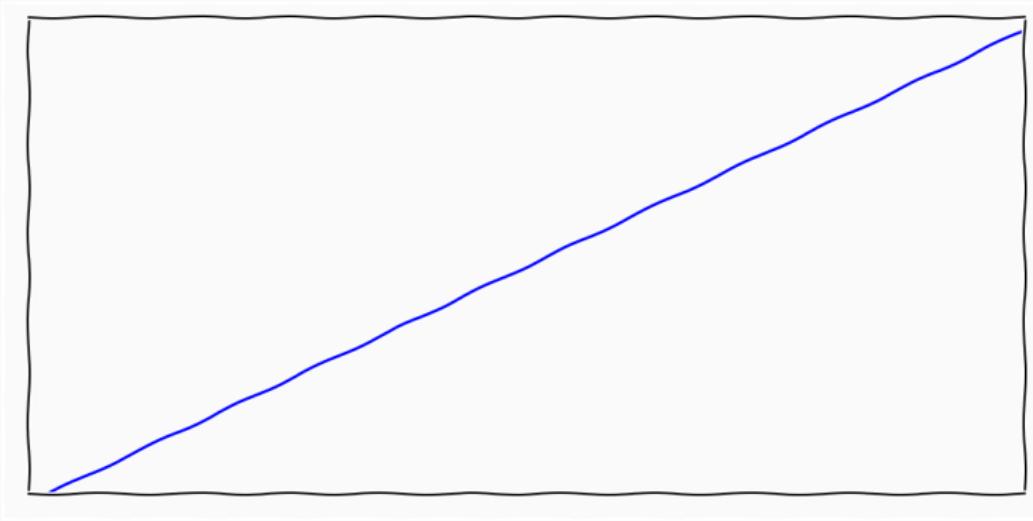
- Over-determined $m > n$

$$\hat{\mathbf{x}} = \underset{\mathbf{x}}{\operatorname{argmin}} \sum_i^m (b_i - \mathbf{A}_{i:} \mathbf{x})^2$$

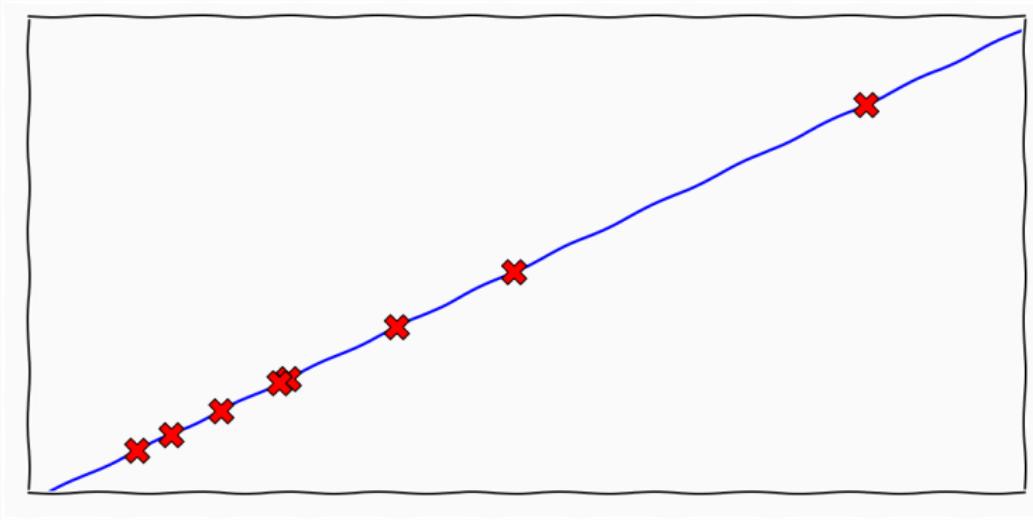
- Under-determined $m < n$

$$\hat{\mathbf{x}} = \underset{\mathbf{x}}{\operatorname{argmin}} \sum_i^m (b_i - \mathbf{A}_{i:} \mathbf{x})^2 + \lambda \mathbf{x}^T \mathbf{x}$$

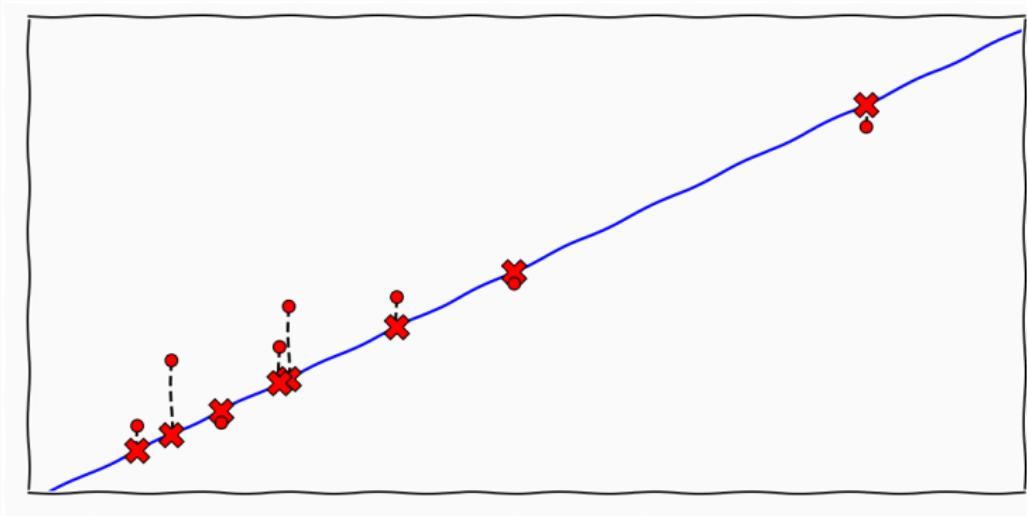
Linear Regression: Modelling



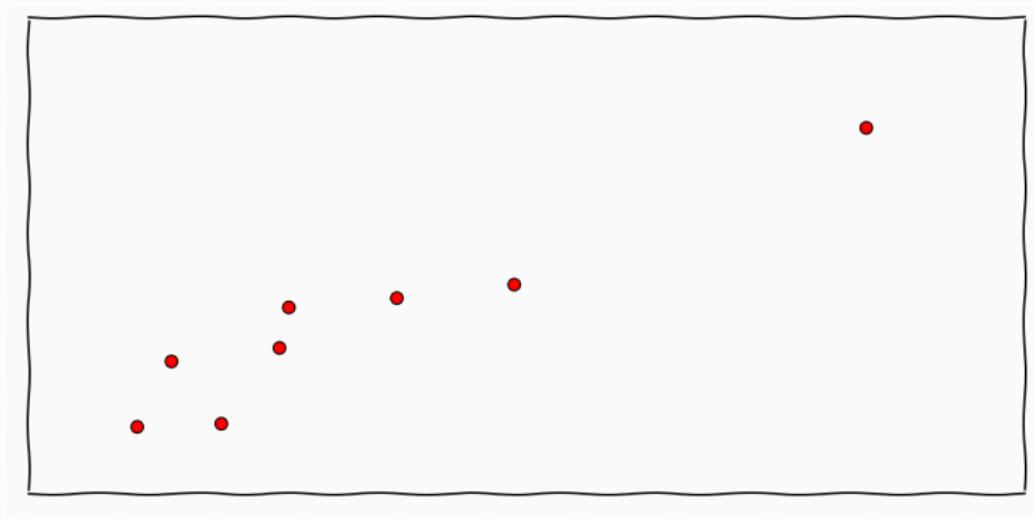
Linear Regression: Modelling



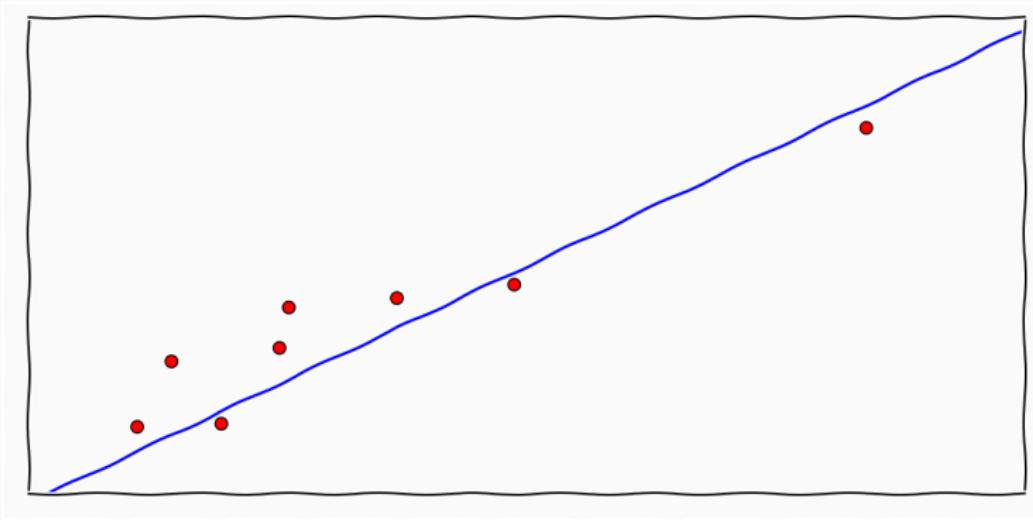
Linear Regression: Modelling



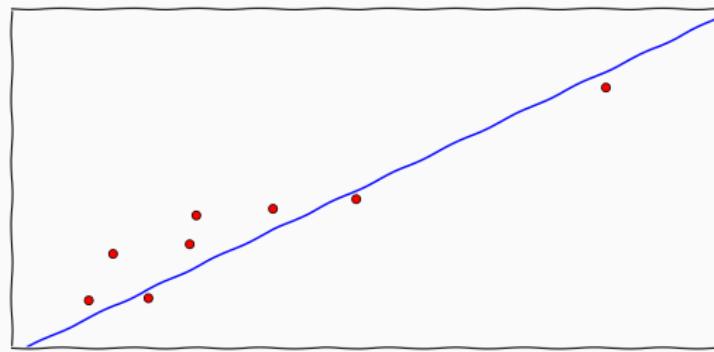
Linear Regression: Modelling



Linear Regression: Modelling



Linear Regression: Modelling



$$y_i = \mathbf{w}^T \mathbf{x}_i + \epsilon_i$$

$$\epsilon_i \sim \mathcal{N}(0, \beta^{-1})$$

Belief/Hypothesis

$$y = \mathbf{w}^T \mathbf{x} + \epsilon$$

Belief/Hypothesis

$$y = \mathbf{w}^T \mathbf{x} + \epsilon$$

$$y - \mathbf{w}^T \mathbf{x} = \epsilon$$

Belief/Hypothesis

$$y = \mathbf{w}^T \mathbf{x} + \epsilon$$

$$y - \mathbf{w}^T \mathbf{x} = \epsilon$$

$$y - \mathbf{w}^T \mathbf{x} \sim \mathcal{N}(\epsilon | 0, \beta^{-1} I) = \left(\frac{\beta}{2\pi} \right)^{\frac{1}{2}} e^{-\frac{1}{2}(\epsilon-0)\beta(\epsilon-0)}$$

Belief/Hypothesis

$$y = \mathbf{w}^T \mathbf{x} + \epsilon$$

$$y - \mathbf{w}^T \mathbf{x} = \epsilon$$

$$y - \mathbf{w}^T \mathbf{x} \sim \mathcal{N}(\epsilon | 0, \beta^{-1} I) = \left(\frac{\beta}{2\pi} \right)^{\frac{1}{2}} e^{-\frac{1}{2}(\epsilon-0)\beta(\epsilon-0)}$$

$$\Rightarrow \mathcal{N}(y - \mathbf{w}^T \mathbf{x} | 0, \beta^{-1} I) = \left(\frac{\beta}{2\pi} \right)^{\frac{1}{2}} e^{-\frac{1}{2}(y-\mathbf{w}^T \mathbf{x})\beta(y-\mathbf{w}^T \mathbf{x})}$$

Belief/Hypothesis

$$y = \mathbf{w}^T \mathbf{x} + \epsilon$$

$$y - \mathbf{w}^T \mathbf{x} = \epsilon$$

$$y - \mathbf{w}^T \mathbf{x} \sim \mathcal{N}(\epsilon | 0, \beta^{-1} I) = \left(\frac{\beta}{2\pi} \right)^{\frac{1}{2}} e^{-\frac{1}{2}(\epsilon-0)\beta(\epsilon-0)}$$

$$\Rightarrow \mathcal{N}(y - \mathbf{w}^T \mathbf{x} | 0, \beta^{-1} I) = \left(\frac{\beta}{2\pi} \right)^{\frac{1}{2}} e^{-\frac{1}{2}(y-\mathbf{w}^T \mathbf{x})\beta(y-\mathbf{w}^T \mathbf{x})}$$

$$\Rightarrow \mathcal{N}(y - \mathbf{w}^T \mathbf{x} | 0, \beta^{-1} I) = \mathcal{N}(y | \mathbf{w}^T \mathbf{x}, \beta^{-1} I)$$

Belief/Hypothesis

$$y = \mathbf{w}^T \mathbf{x} + \epsilon$$

$$y - \mathbf{w}^T \mathbf{x} = \epsilon$$

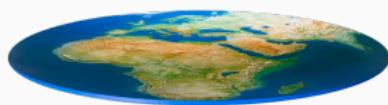
$$y - \mathbf{w}^T \mathbf{x} \sim \mathcal{N}(\epsilon | 0, \beta^{-1} I) = \left(\frac{\beta}{2\pi} \right)^{\frac{1}{2}} e^{-\frac{1}{2}(\epsilon-0)\beta(\epsilon-0)}$$

$$\Rightarrow \mathcal{N}(y - \mathbf{w}^T \mathbf{x} | 0, \beta^{-1} I) = \left(\frac{\beta}{2\pi} \right)^{\frac{1}{2}} e^{-\frac{1}{2}(y-\mathbf{w}^T \mathbf{x})\beta(y-\mathbf{w}^T \mathbf{x})}$$

$$\Rightarrow \mathcal{N}(y - \mathbf{w}^T \mathbf{x} | 0, \beta^{-1} I) = \mathcal{N}(y | \mathbf{w}^T \mathbf{x}, \beta^{-1} I)$$

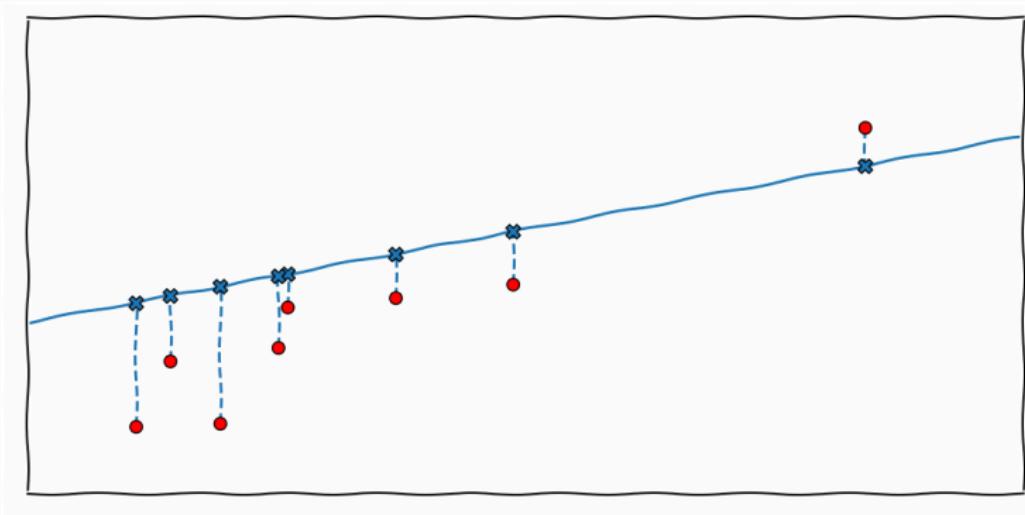
$$\Rightarrow p(y | \mathbf{w}, \mathbf{x}) = \mathcal{N}(y | \mathbf{w}^T \mathbf{x}, \beta^{-1} I)$$

Likelihood

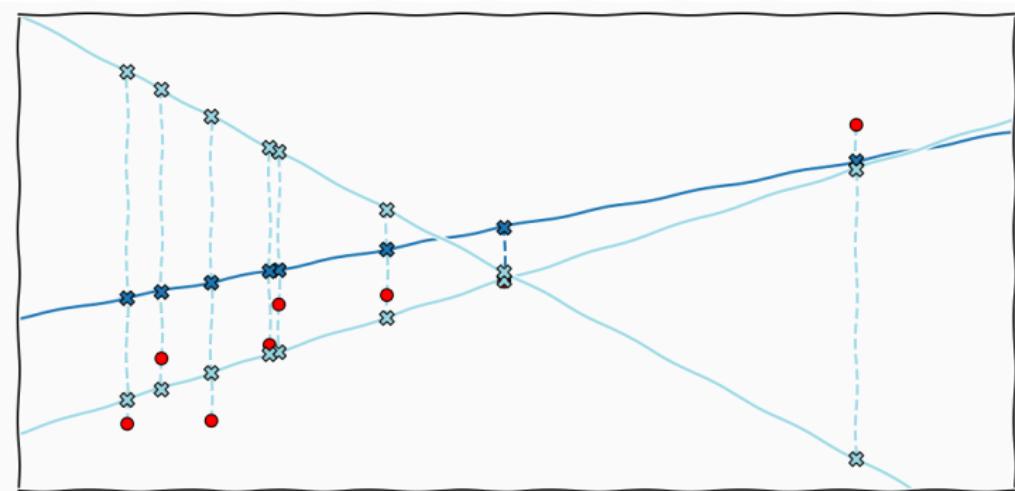


"Given the premise that the earth is flat, how supportive do I believe observations y are of this?"

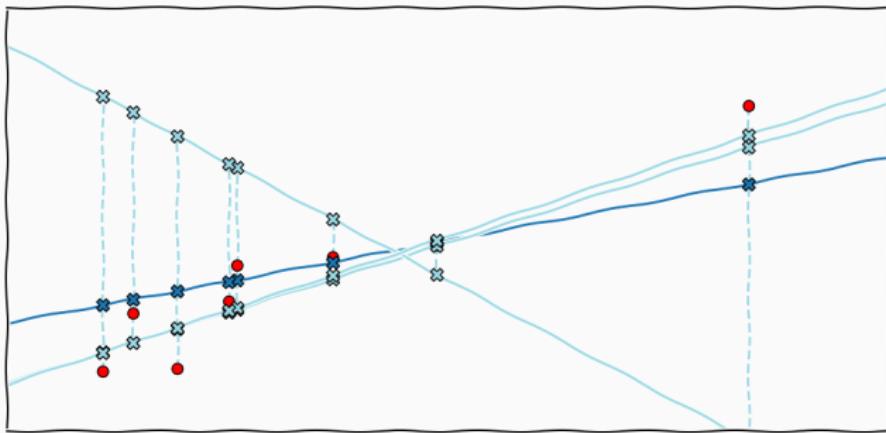
Linear Regression: Modelling



Linear Regression: Modelling

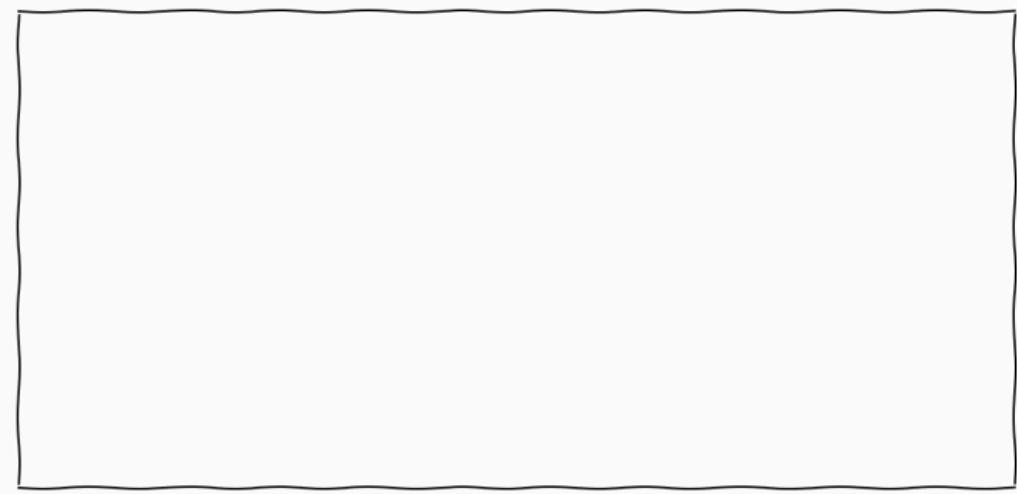


Linear Regression: Likelihood

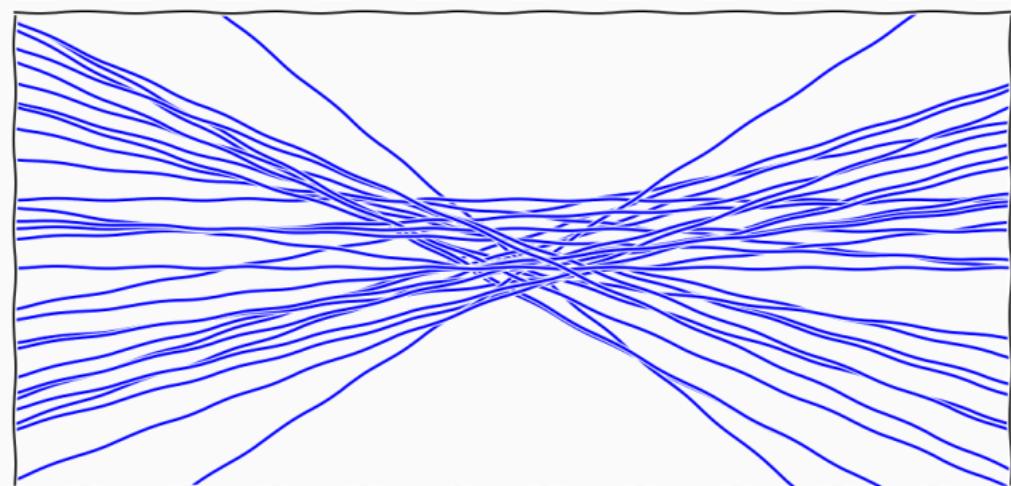


$$p(y|\mathbf{w}, \mathbf{x}) = \mathcal{N}(y|\mathbf{w}^T \mathbf{x}, \beta^{-1} I)$$

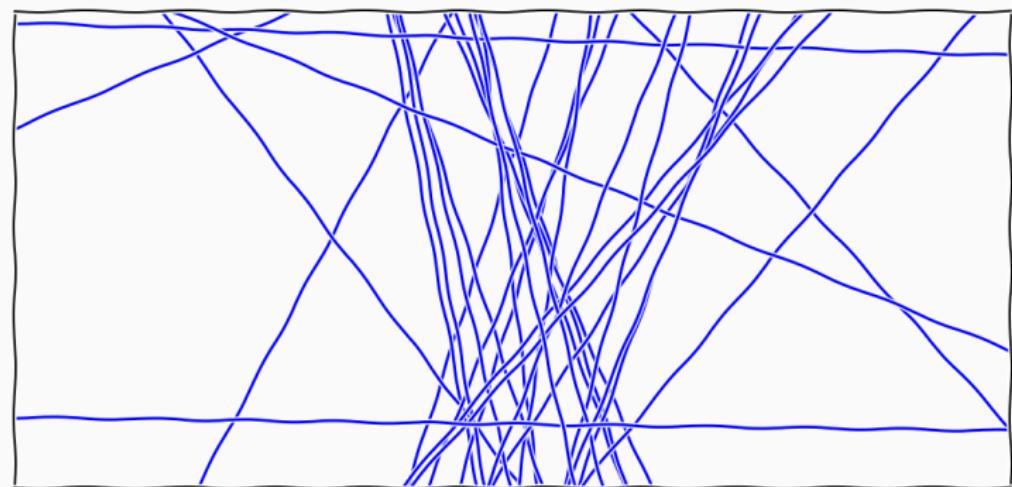
But I don't believe in a flat earth



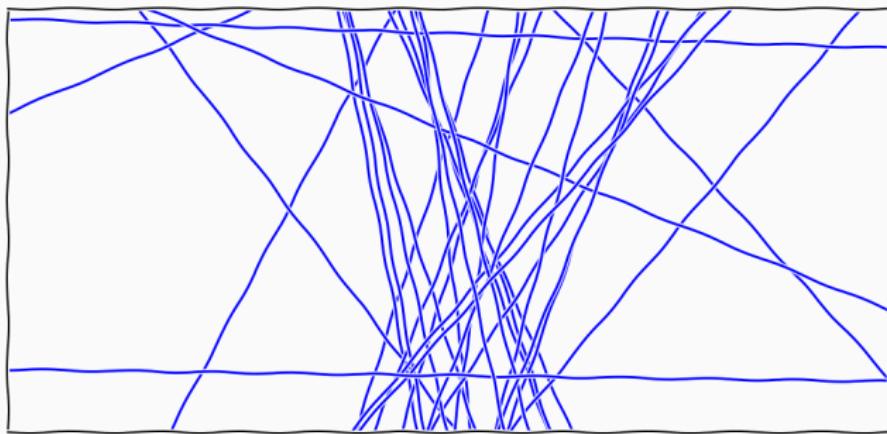
Prior



Prior



Prior

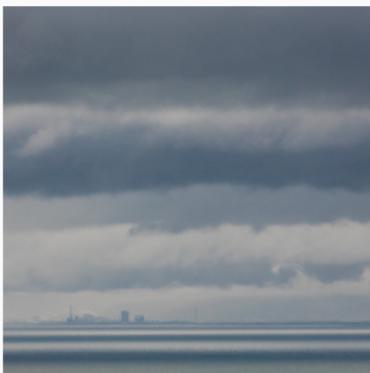


$$w \sim \mathcal{N}(0, 2)$$



"Well this is how much credibility **belief** I give to the hypothesis
that the earth is flat"

Two opposing theories



That French Dude Again



It is remarkable that a science which began with the consideration of games of chance should have become the most important object of human knowledge.
– Laplace *Laplace, 1814b*

Bayes' "Rule"

$$p(y, w) = p(y|w)p(w)$$

Bayes' "Rule"

$$p(y, w) = p(y|w)p(w)$$

$$p(y, w) = p(w|y)p(y)$$

Bayes' "Rule"

$$p(y, w) = p(y|w)p(w)$$

$$p(y, w) = p(w|y)p(y)$$

$$p(w|y)p(y) = p(y|w)p(w)$$

Bayes' "Rule"

$$p(y, w) = p(y|w)p(w)$$

$$p(y, w) = p(w|y)p(y)$$

$$p(w|y)p(y) = p(y|w)p(w)$$

$$p(w|y) = \frac{p(y|w)p(w)}{p(y)}$$

Bayes' "Rule"

$$p(y, w) = p(y|w)p(w)$$

$$p(y, w) = p(w|y)p(y)$$

$$p(w|y)p(y) = p(y|w)p(w)$$

$$p(w|y) = \frac{p(y|w)p(w)}{p(y)}$$

$$= \frac{p(y|w)p(w)}{\int p(y|w)p(w)dw}$$

$$p(w | y) = \frac{p(y | w)p(w)}{\int p(y | w)p(w)dw}$$

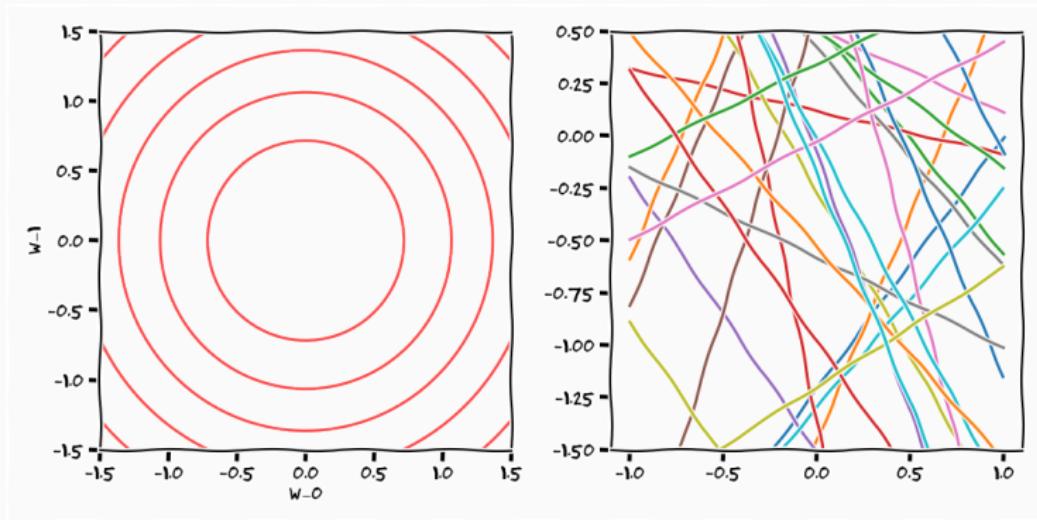
Likelihood How much **evidence** is there in the data for a specific hypothesis

Prior What are my beliefs about different hypothesis

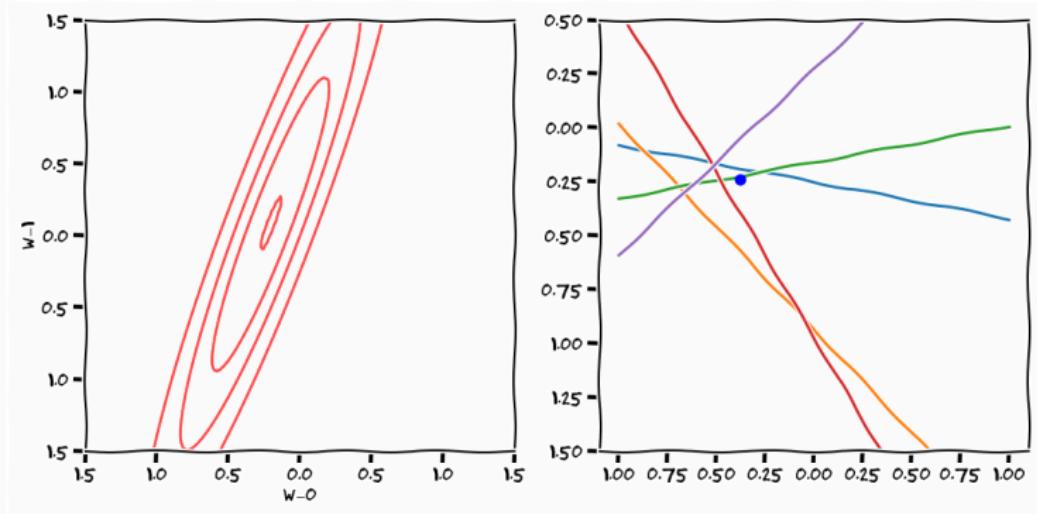
Posterior What is my **updated** belief after having seen data

Evidence What is my belief about **any** observations

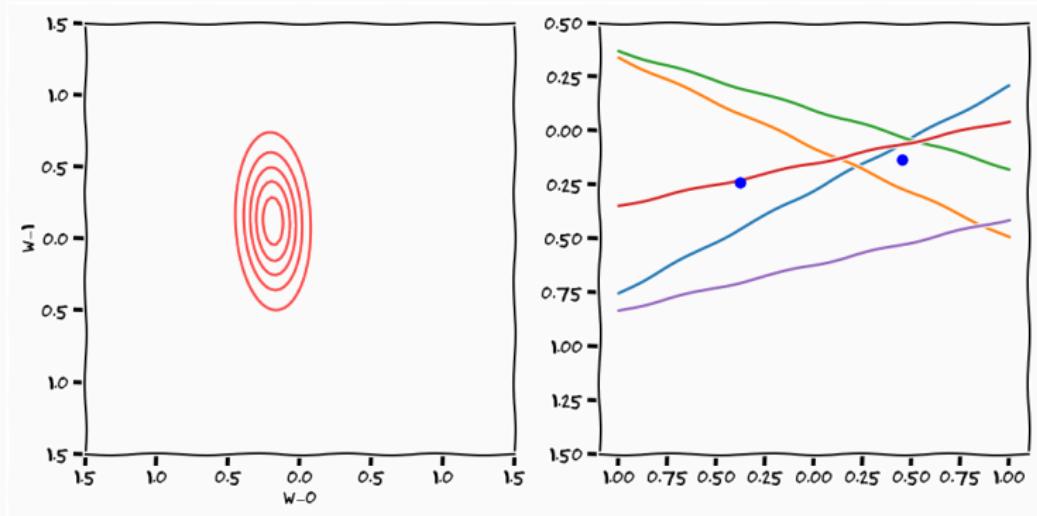
Linear Regression Example



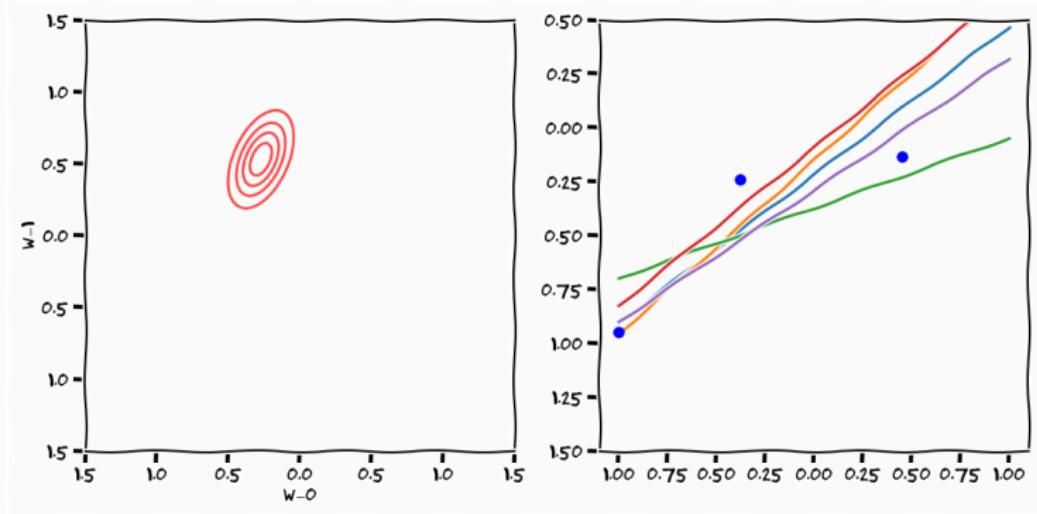
Linear Regression Example



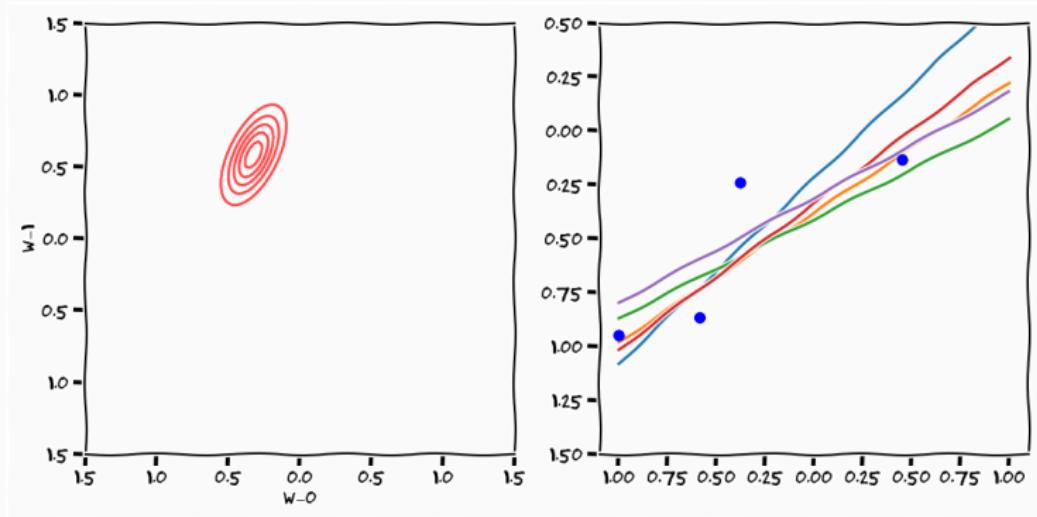
Linear Regression Example



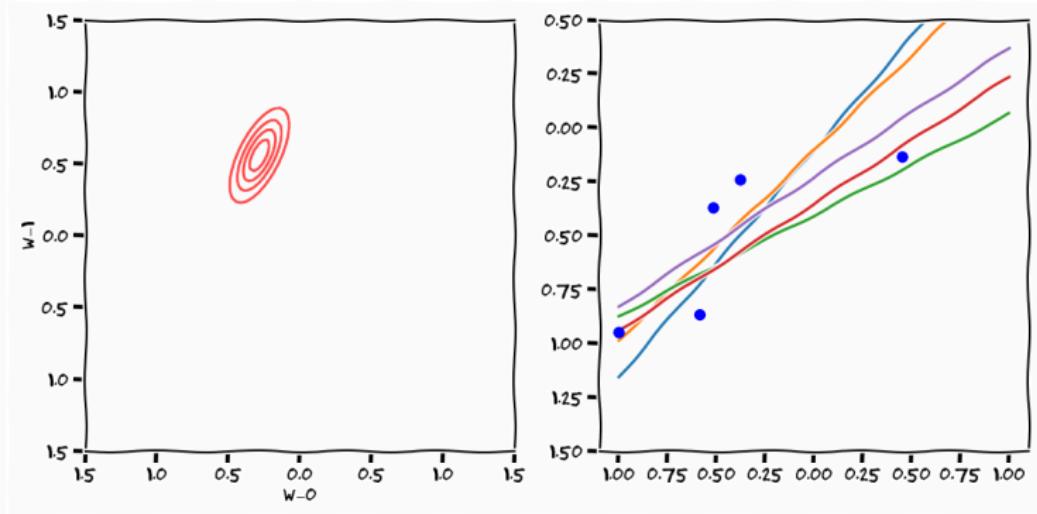
Linear Regression Example



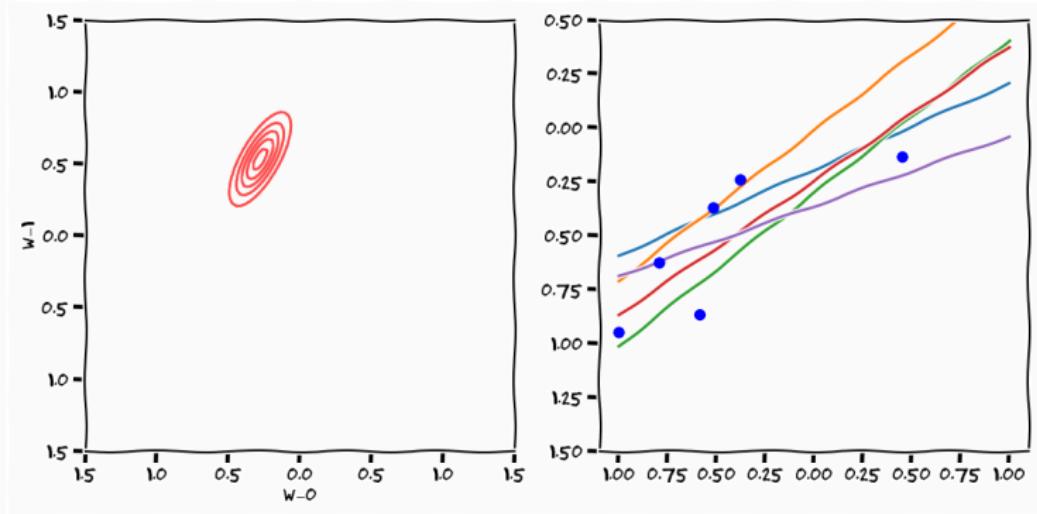
Linear Regression Example



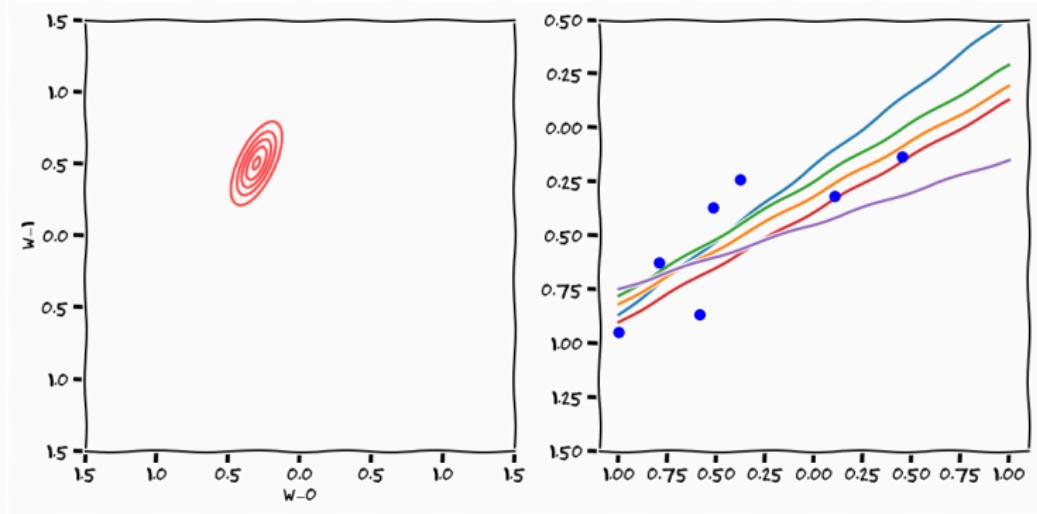
Linear Regression Example



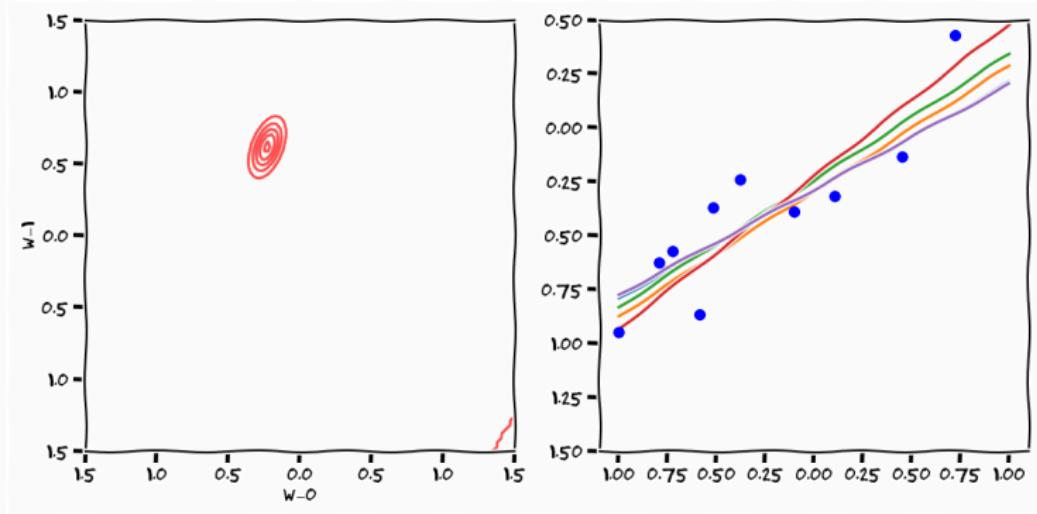
Linear Regression Example



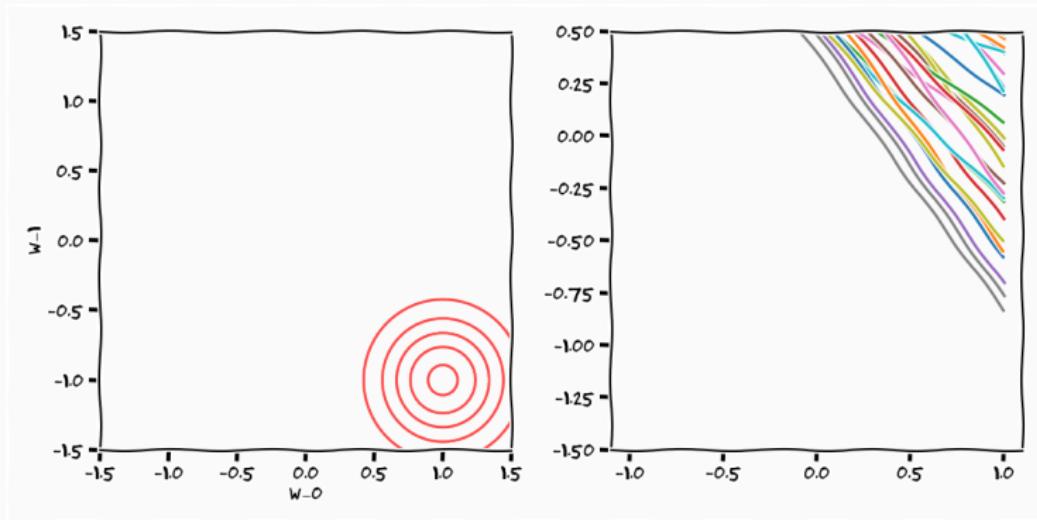
Linear Regression Example



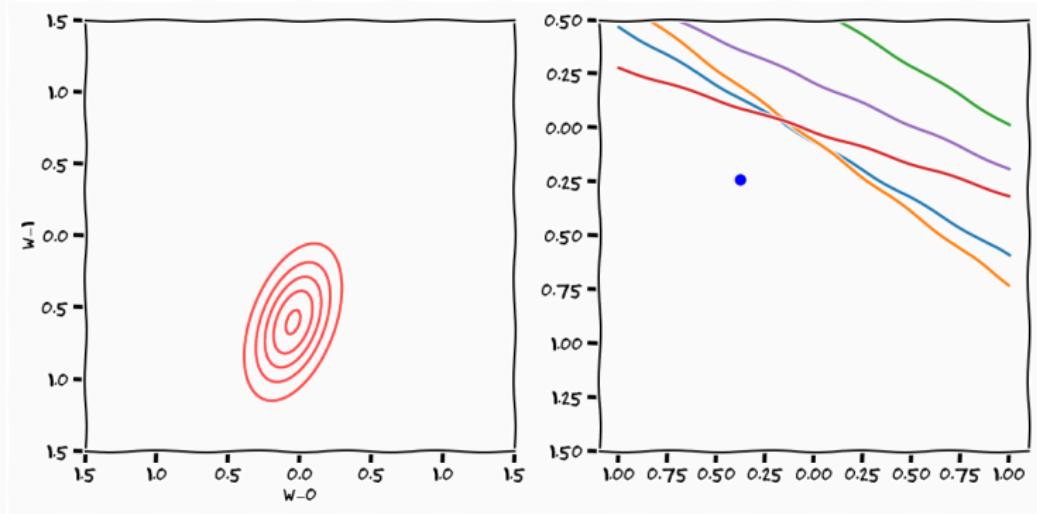
Linear Regression Example



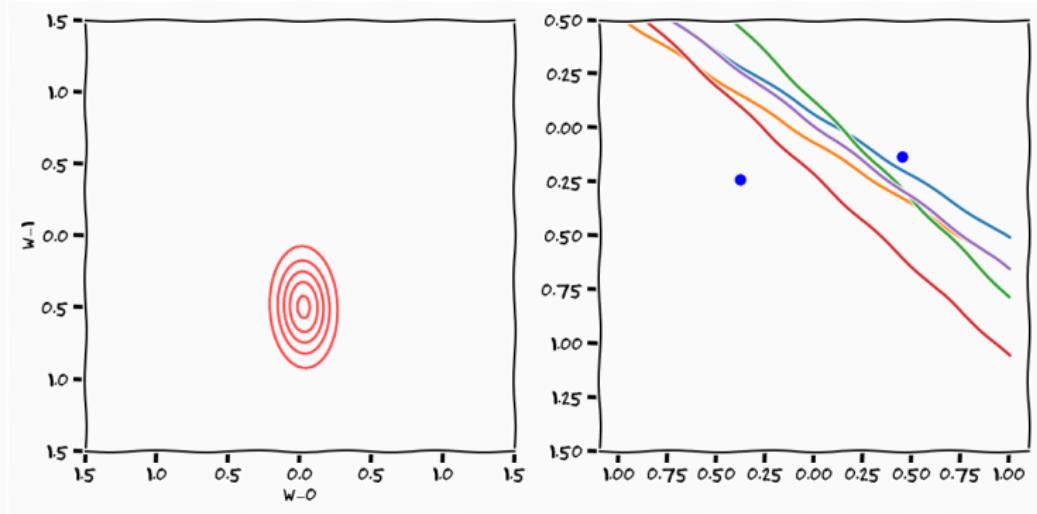
Linear Regression Example



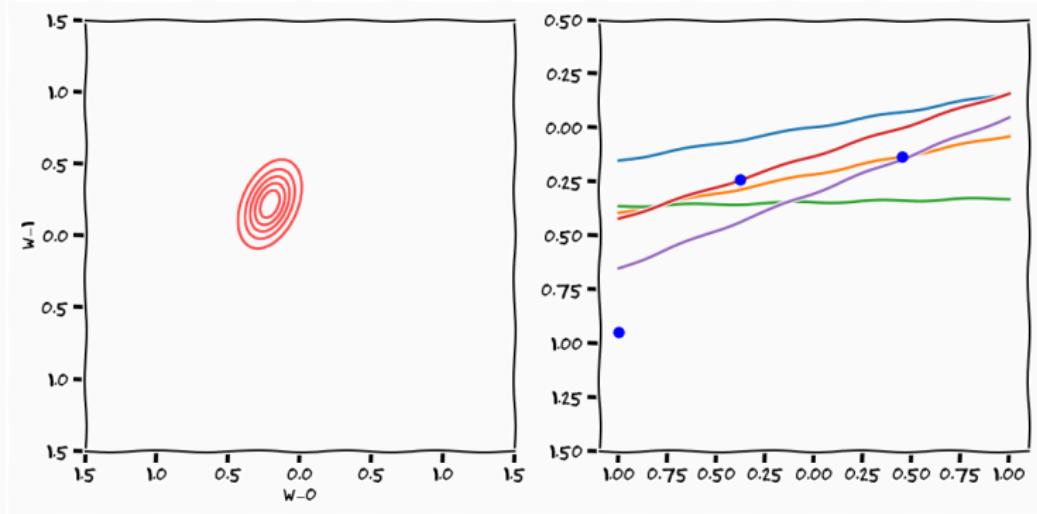
Linear Regression Example



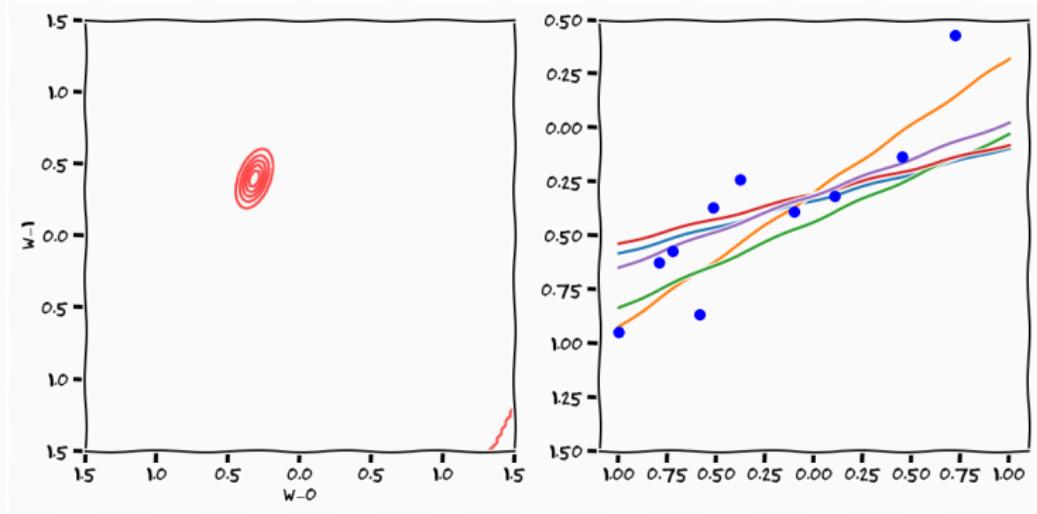
Linear Regression Example



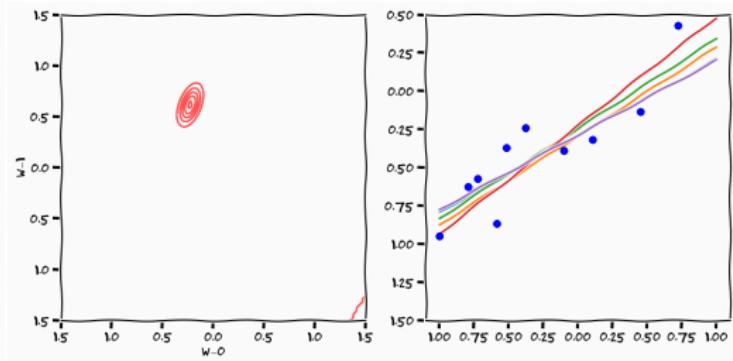
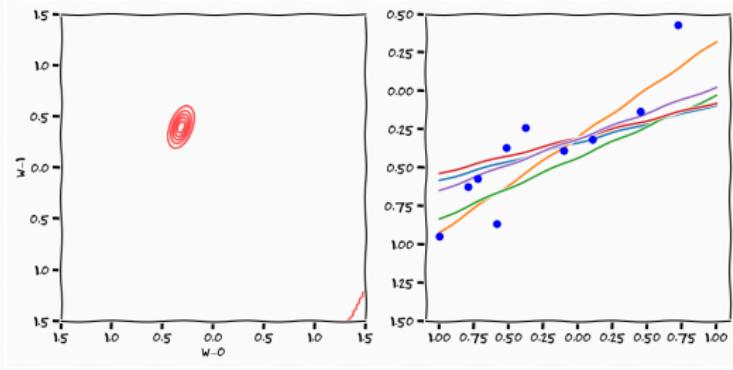
Linear Regression Example



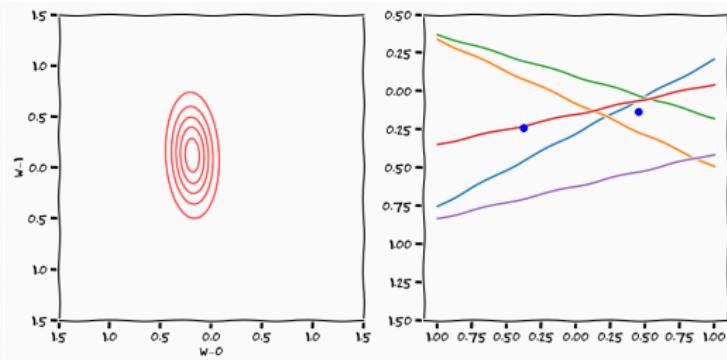
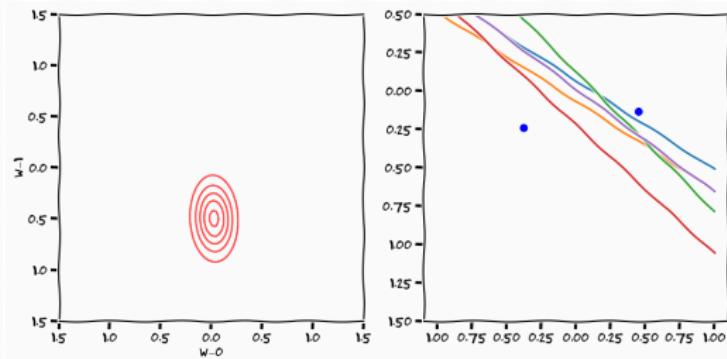
Linear Regression Example



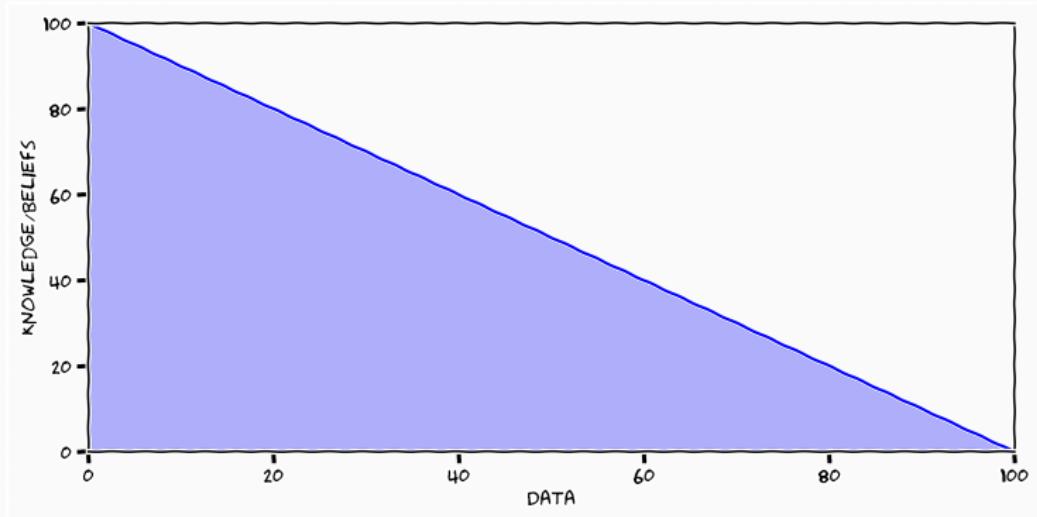
There is overwhelming evidence



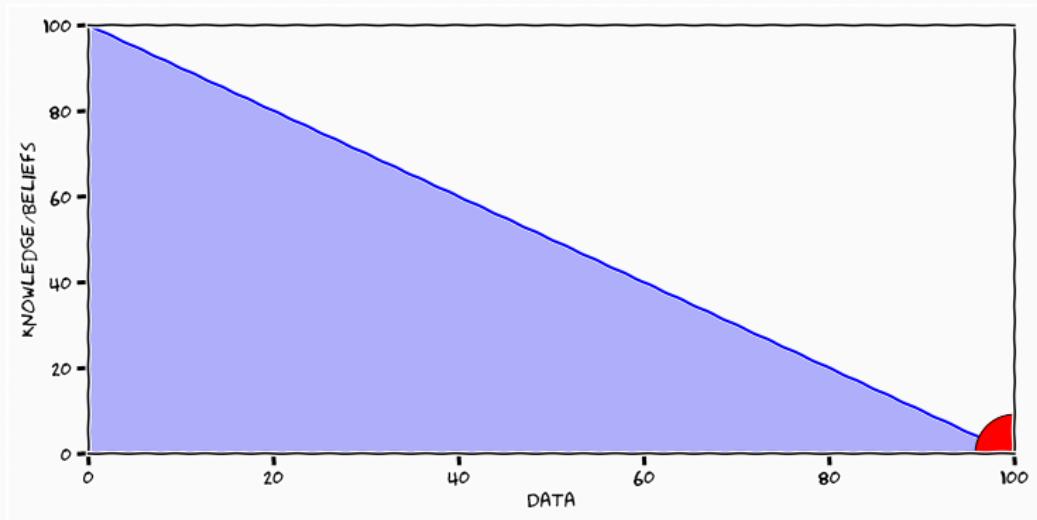
Knowledge is Relative



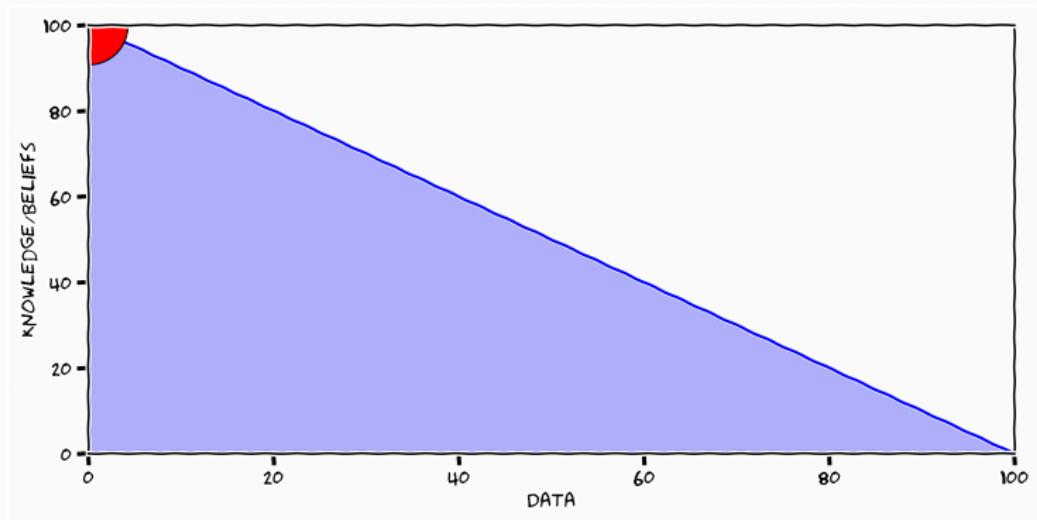
Data and Knowledge



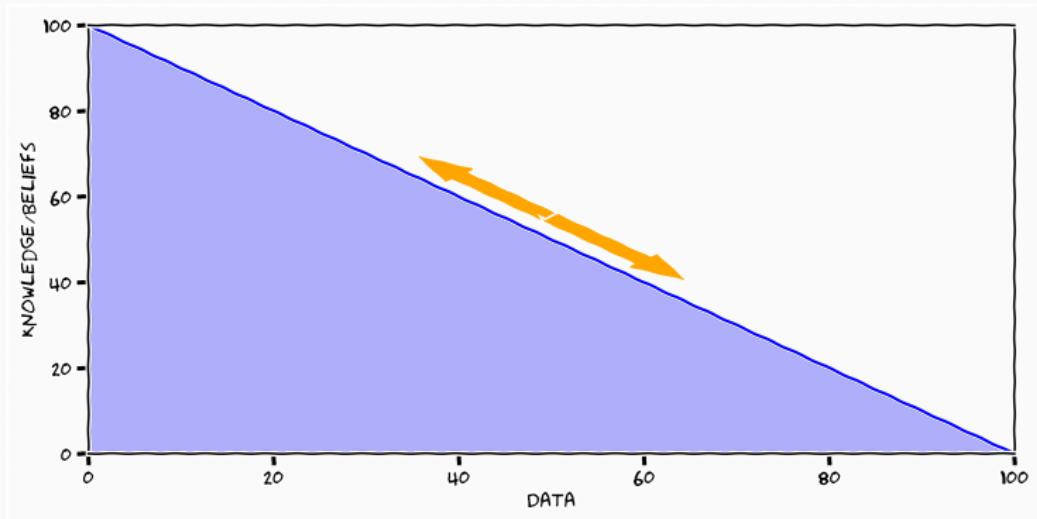
Data and Knowledge



Data and Knowledge

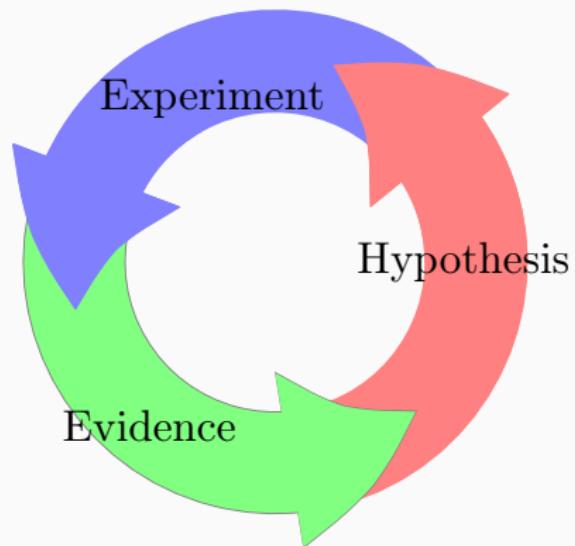


Data and Knowledge



Scientific Principle

Scientific Principle



Where do we start

- Inductive Reasoning

Observation → Pattern → Hypothesis → Theory

Where do we start

- Inductive Reasoning

Observation → Pattern → Hypothesis → Theory

- Deductive Reasoning

Theory → Hypothesis → Observation → Confirmation

Deductive Science

"Science should attempt to disprove a theory, rather than attempt to continually support theoretical hypotheses."

*– Karl Popper *The Logic of Scientific Discovery**

1. Facilitate viewing implications of Hypothesis in observation space

$$p(w) \rightarrow p(w \mid y)$$

2. Facilitate selection procedure of Hypothesis preference

$$w_1 \succ w_2 \quad p(y \mid w_1) = p(y \mid w_2)$$

What is a good hypothesis?

"In so far as a scientific statement speaks about reality, it must be falsifiable: and in so far as it is not falsifiable, it does not speak about reality."

*– Karl Popper *The Logic of Scientific Discovery**

What is a good hypothesis?

"A theory that explains everything, explains nothing"
– Karl Popper *The Logic of Scientific Discovery*

Logic vs. Probability

$$P \rightarrow \neg Q$$

$$\neg P$$

$$Q$$

$$p(y) = \int p(y|\theta)p(\theta)d\theta$$

Falsifiability and Occams' Razor

- A hypothesis should be judged based on how easily it can be falsified

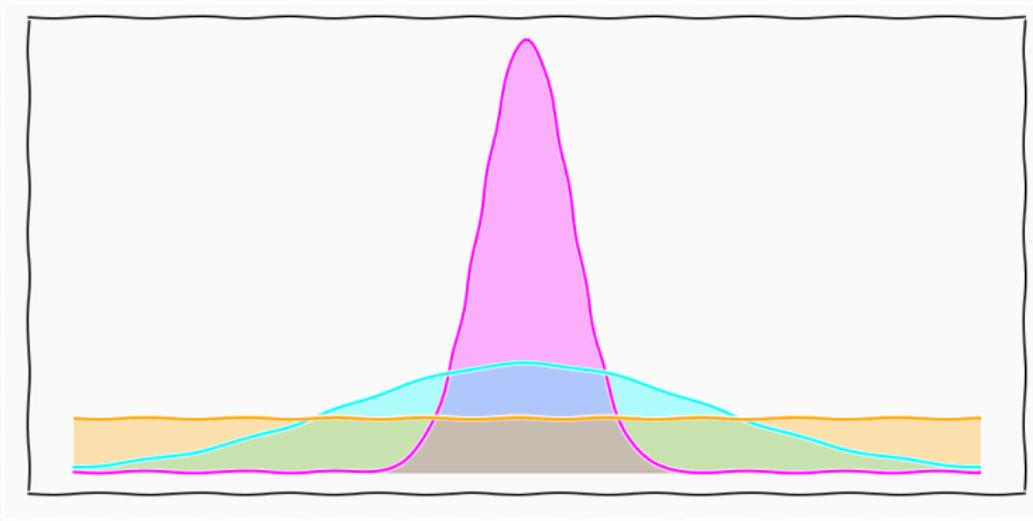
Falsifiability and Occams' Razor

- A hypothesis should be judged based on how easily it can be falsified
- The more general a theory is the more cases/possibilities it allows for falsification

Falsifiability and Occams' Razor

- A hypothesis should be judged based on how easily it can be falsified
- The more general a theory is the more cases/possibilities it allows for falsification
- "The more strongly our framework can **differentiate** different hypothesis the better it is for falsification"

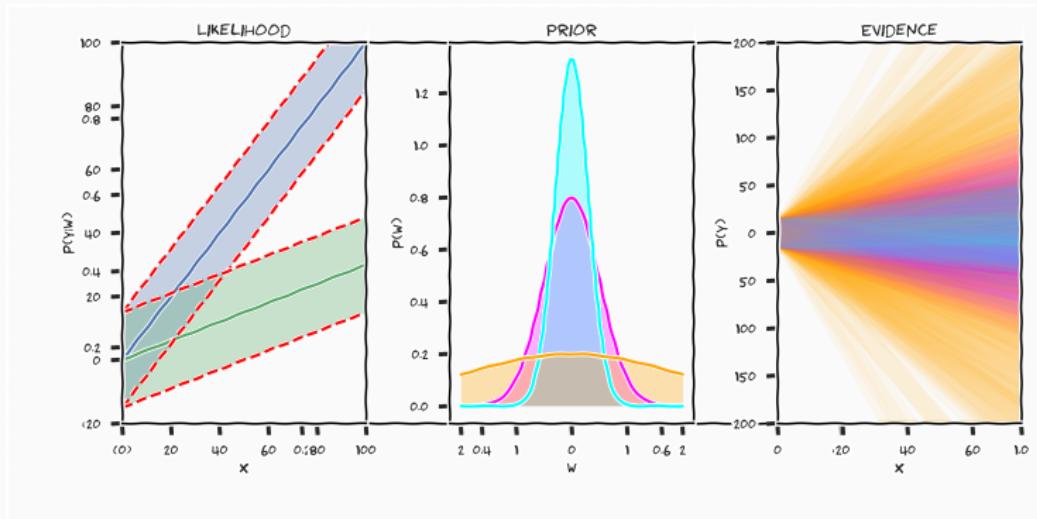
Distributions



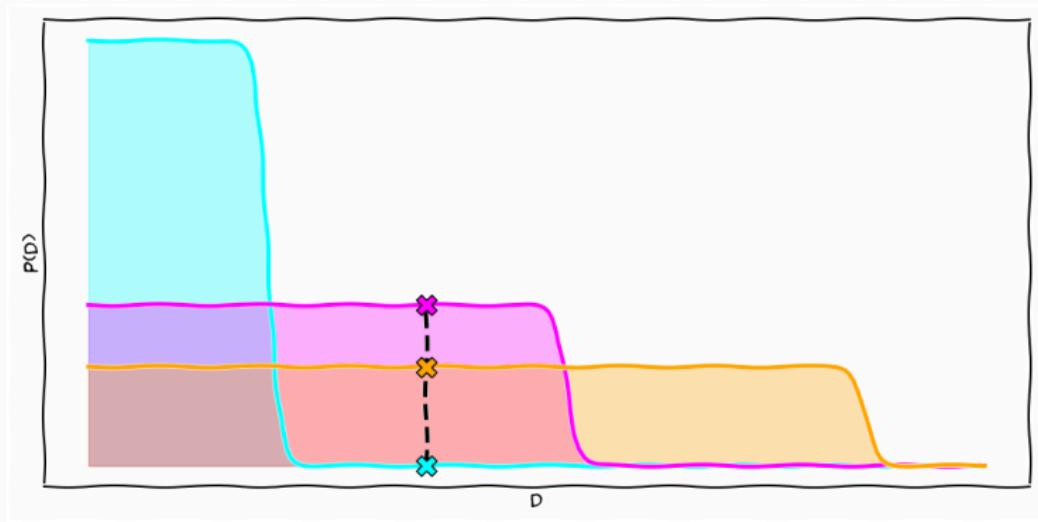
Model Evidence and Occams' Razor

$$p(y) = \int p(y | w)p(w)dw$$

What is can be falsified?



The MacKay Plot Mackay, 1991



Is Machine Learning a Science?

- How to build mathematical models of hypothesis

$$\text{hypothesis} \approx p(w)$$

Is Machine Learning a Science?

- How to build mathematical models of hypothesis

$$\text{hypothesis} \approx p(w)$$

- How to mathematically update our knowledge with data

$$p(w | y) \approx \frac{p(y | w)p(w)}{\int p(y | w)p(w)dw}$$

Summary

Summary

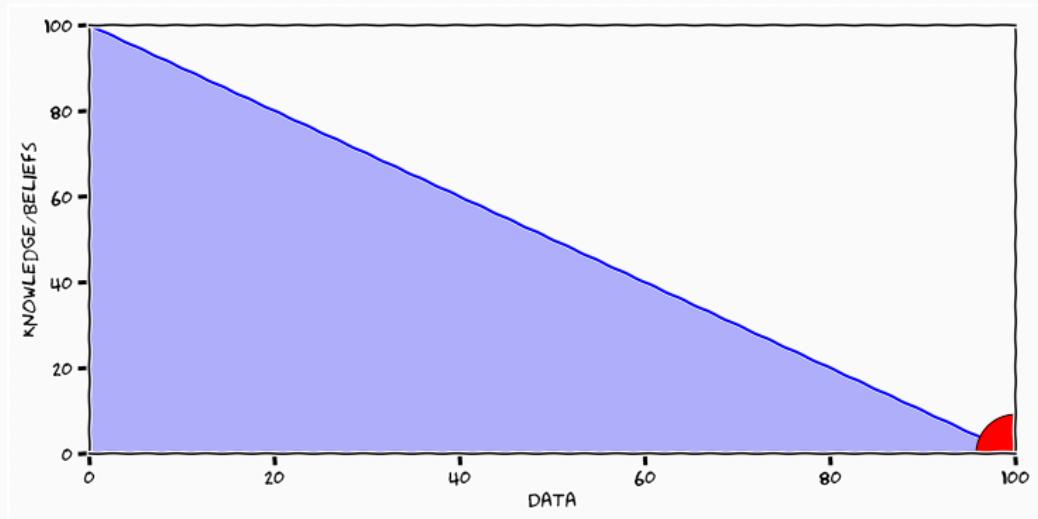
Machine Learning is a framework for combining knowledge with data to recover an interpretation of the data in light of the knowledge.

Where does Knowledge Come From?





Don't believe the hype

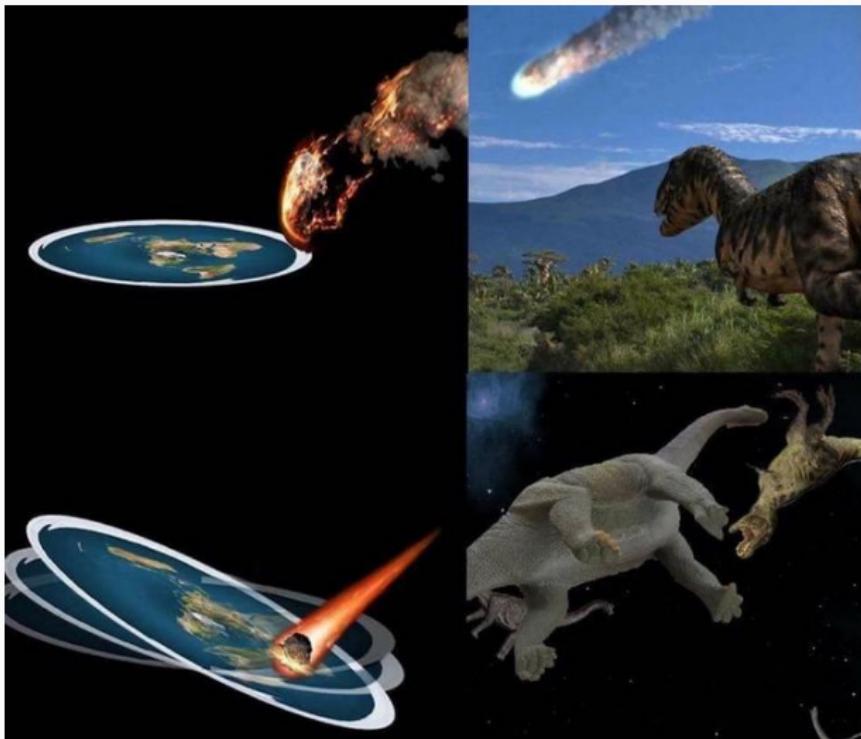


Summary

"There is no logical paths leading to "these laws" they can only be reached by intuition based on something like and intellectual love of the objects of experience"

– Albert Einstein

Thats why they disappeared



eof

References

-  Laplace, Pierre Simon (1814a). *A philosophical essay on probabilities*.
-  — (1814b). *Théorie analytique des probabilités*. Mme. Ve. Sourcier.
-  Mackay, David J C (Dec. 1991). “Bayesian methods for adaptive models”. PhD thesis. California Institute of Technology: California Institute of Technology.
-  Popper, K.R. (1959). *The Logic of Scientific Discovery*. ISSR library. Routledge.