

# Data Science & Machine Learning

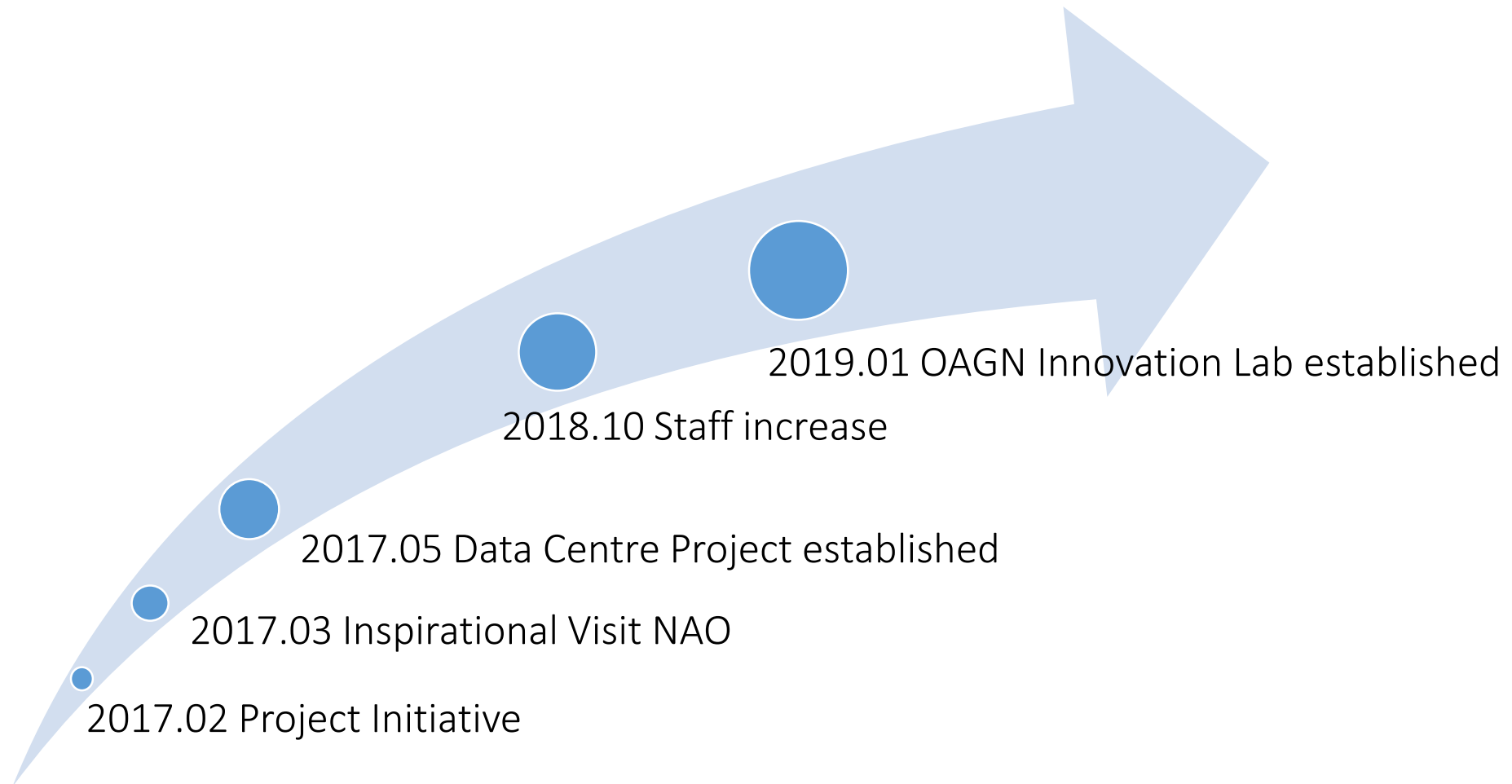
Musings and some examples from OAG Norway

Jan Roar

Chief Data Scientist

OAGN Innovation Lab

# The OAGN Data Science Project



# Purpose

1. Curating data for financial and performance audit
2. On-demand data analytics
3. Promoting data science and the use of machine learning at the OAGN
4. Experiment with new technologies and methods

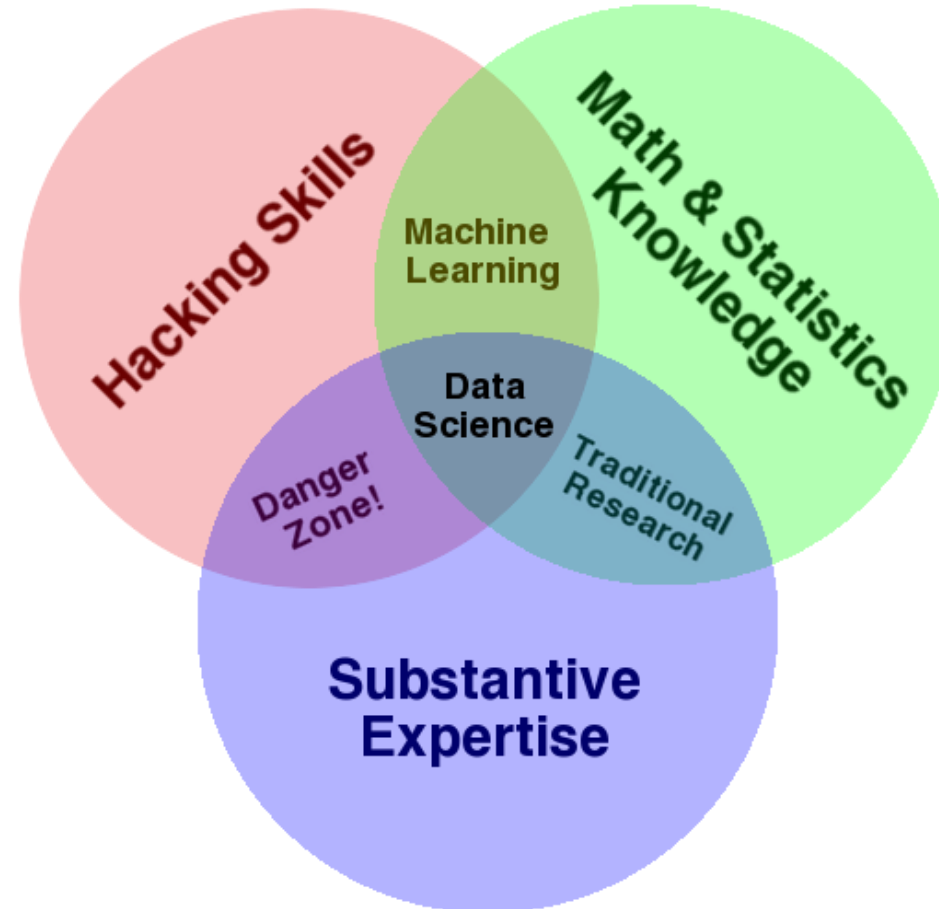
We automate the boring stuff,  
so you can audit the exiting stuff!

# Where we come from

- 2 political scientists
- 1 economist
- 1 sociologist
- 1 physicist

... we're not an IT project, we're a **data science** project

# What is Data Science



You need both:

- Coding skills
- Statistical skills
- Audit skills

<http://drewconway.com/zia/2013/3/26/the-data-science-venn-diagram>

The world is overflowing with data

This is difficult to manage for SAs

# Partly because of this

Law on The Office of the Auditor General of Norway § 12:

«The OAGN can, without restrictions of confidentiality, demand any information, any disclosure or any document, and execute any examination deemed necessary by the OAGN to fulfil audit tasks.»

(my own highly informal translation)



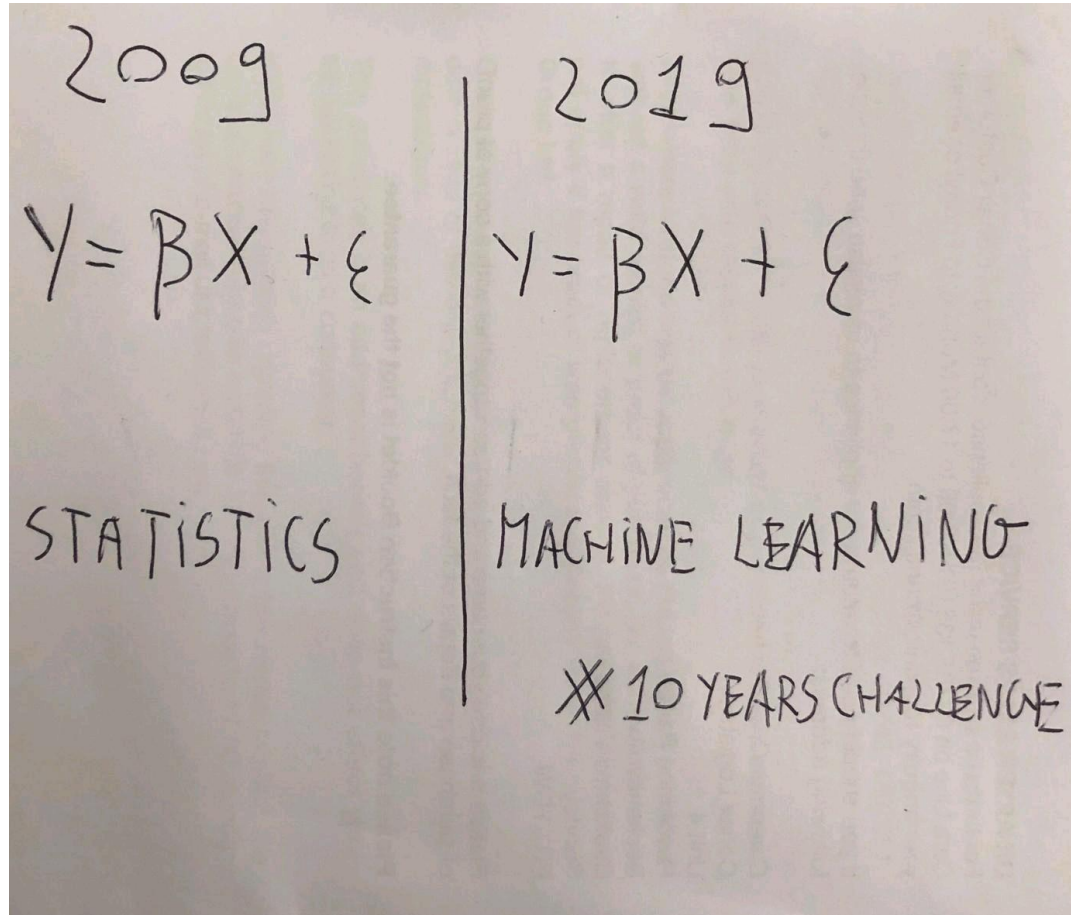
Data → Information → Knowledge → Wisdom

# What's all this about machine learning?

«Everybody» keeps fussing about  
machine learning and artificial intelligence these days  
... and you can get a bit worried...  
... when seeing stuff like this...



# However...



Found on LinkedIn

# What's the point – an example

- You have 10 000 pages of procurement contract documents (text)
- And you have regulations for procurement
- Then you want to compare the 10 000 pages and the regulations
- To uncover illegal purchases
- Can machines do this for us?

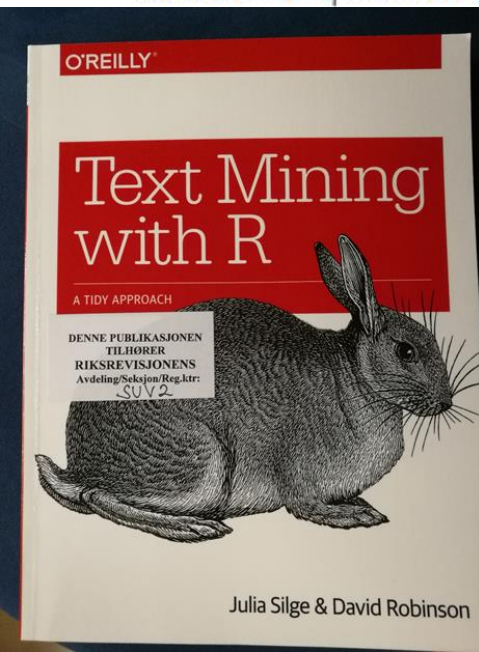
# Natural Language Processing

What it is and why it matters

Natural language processing (NLP) is a branch of **artificial intelligence** that helps computers understand, interpret and manipulate human language. NLP draws from many disciplines, including computer science and computational linguistics, in its pursuit to fill the gap between human communication and computer understanding.

Source:

[https://www.sas.com/en\\_us/insights/analytics/what-is-natural-language-processing-nlp.html](https://www.sas.com/en_us/insights/analytics/what-is-natural-language-processing-nlp.html)



# NLP in practice – from speech to text

## Cloud Speech-to-Text

Speech-to-text conversion powered by machine learning and available for short-form or long-form audio.

[GO TO CONSOLE](#)

View [documentation](#) for this product.

## Powerful speech recognition

Google Cloud Speech-to-Text enables developers to convert audio to text by applying powerful neural network models in an easy-to-use API. The API recognizes 120 languages and variants to support your global user base. You can enable voice command-and-control, transcribe audio from call centers, and more. It can process real-time streaming or prerecorded audio, using Google's machine learning technology.



Convert your speech to text right now

Select a language and click "Start Now" to begin recording

Input type

☒ Microphone   ☐ File upload

Language

English (United States)

Speaker diarization BETA

Off

Speakers

1 speaker

Punctuation



Google Home



# Algorithms – Economic, but Effective?

- government have started using machine learning, for example in awarding grants
  - even for decisions involving discretionary considerations
- we will need performance auditors that can look into “the black box of algorithms”
- to evaluate and conclude on decisions made by machines



If you understand this,  
you have understood machine learning

$$Y = f(x) + e$$

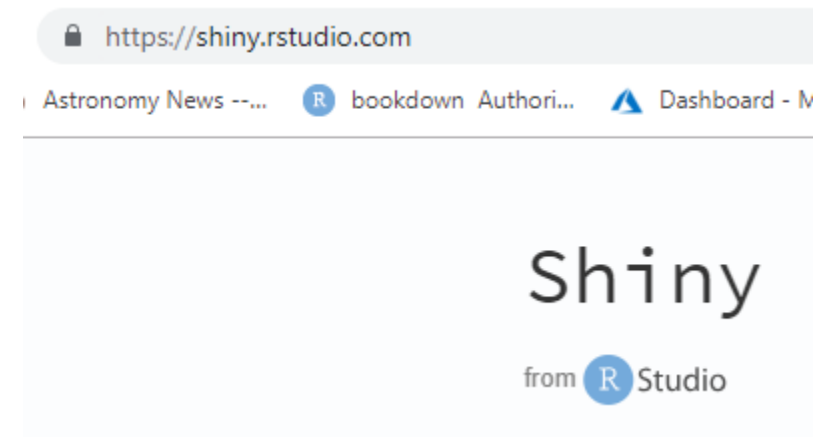
Some concrete examples

# Example: Webapp for DEA-analysis

- DEA = Data Envelopment Analysis
- Benchmark analysis based on Output vs. Input
- Example:
  - District courts – how many judges/office clerks vs. how many cases handled
- Productivity Analysis

# DEA – efficiency in the form of productivity

- DEA an important method for performance audit, but
- Most software for DEA is old and/or user unfriendly
- Several packages in R for doing DEA # but not easy to use
- So: we wrote a Shiny app for DEA
- We named it pioneerR



# We're open!

pioneerR is open source (free to use and free to edit the source code)!

You can find us on GitHub\*:

[github.com/Riksrevisjonen/pioneerR](https://github.com/Riksrevisjonen/pioneerR)

\*The largest collection of open source software on the Internet

# Takeaway:

- Program code is universal
- As such very easy to share

# Search: 19 000 PDFs – what to do

- We have 19 000 PDF documents from Norwegian hospital boards
- 25 hospital boards, 12-14 board meetings a year, last 5 years)
- All documents published on the web (25 different websites)



elastic

Products

Cloud

Services

Customers

Learn

downloads

contact



EN

Elasticsearch

Download

As A Service



## The Heart of the Elastic Stack

Elasticsearch is a distributed, RESTful search and analytics engine capable of solving a growing number of use cases. As the heart of the Elastic Stack, it centrally stores your data so you can discover the expected and uncover the unexpected.





# Cloud Platform: The big 3

Microsoft Azure

Kontakt salgssavdelingen: 800-62-116

Hvorfor Azure

Løsninger

Produkter

Dokumentasjon

Priser

Opplæring

Markedsplass

Partnere

Støtte

Blogg

Mer

Opprett den gratis Azure-kontoen din i dag

Kom igang med utviklingen av din neste ide med Azure


Start gratis

A screenshot of the AWS website. The top navigation bar is dark blue with white text for links: 'Contact Sales', 'Products' (with a dropdown arrow), 'Solutions', 'Pricing', 'Getting Started', 'Documentation', 'AWS Marketplace', 'Support', and a partially visible 'C'. The AWS logo is on the left. The main content area has a dark blue background. It features the heading 'Introducing Amazon EC2 C5d Instances' in large white text. Below the heading is a paragraph: 'Amazon EC2 C5 instances backed by high performance block-level storage for compute-intensive workloads'. At the bottom of this section is a link 'Learn more »' in orange. In the bottom right corner, there are five small white dots, with the third dot from the left being filled, indicating the current slide in a sequence.

The image shows the Google Cloud Platform logo, which is a stylized cloud composed of four interlocking rings in red, yellow, green, and blue. To the right of the logo is a diagram illustrating the Google Cloud architecture. The diagram shows a central 'Google Network' (represented by a hexagon) connected to various services. Above the network, there are labels for 'Google Cloud Storage' and 'Google Cloud Compute Engine'. Below the network, there are labels for 'Google Cloud SQL' and 'Google Cloud Pub/Sub'. The diagram also includes a 'Premium Tier' label and a 'Pay' label with a dollar sign icon.

The screenshot displays the Azure portal's monitoring interface for a specific resource group. The top navigation bar includes links for Dashboard, Home, Notifications, Diagnostics, and Settings. The main content area is organized into a grid of monitoring widgets:



- Web Front End:**
  - HTTP 200: HTTP 200 (msec):** A line chart showing response time over the last 10 minutes, with a peak of 5.2K.
  - CPU percentage and Memory Percentage Past week:** A line chart showing CPU and memory usage over the last 7 days.
  - Average Response Time and CPU Time Today:** A line chart showing average response time and CPU usage over the last 24 hours.
- Database:**
  - CPU percentage and Database size (percentage, GB):** A line chart showing CPU usage and database size over the last 24 hours.
- Processes:**
  - All runs:** A table listing processes with columns for Name, Status, Time, and Duration.
- Cognitive Services:**
  - Total Calls past week:** A line chart showing the total number of calls over the last 7 days.
- Summary Section:**
  - Summary:** A table showing the total CPU usage (32%) and total database size (55000 GB).
  - SQL Azure (by instance) (percentage):** A line chart showing SQL Azure usage over the last 24 hours.
  - SQL Azure:** A table showing the total CPU usage (32%) and total database size (55000 GB).



**Lightsail**  
Everything you need to get started on  
AWS—for a low, predictable price



**Amazon EC2 I3 Bare Metal Instances**  
 Deliver direct access to next-generation Nitro-based AWS hardware infrastructure

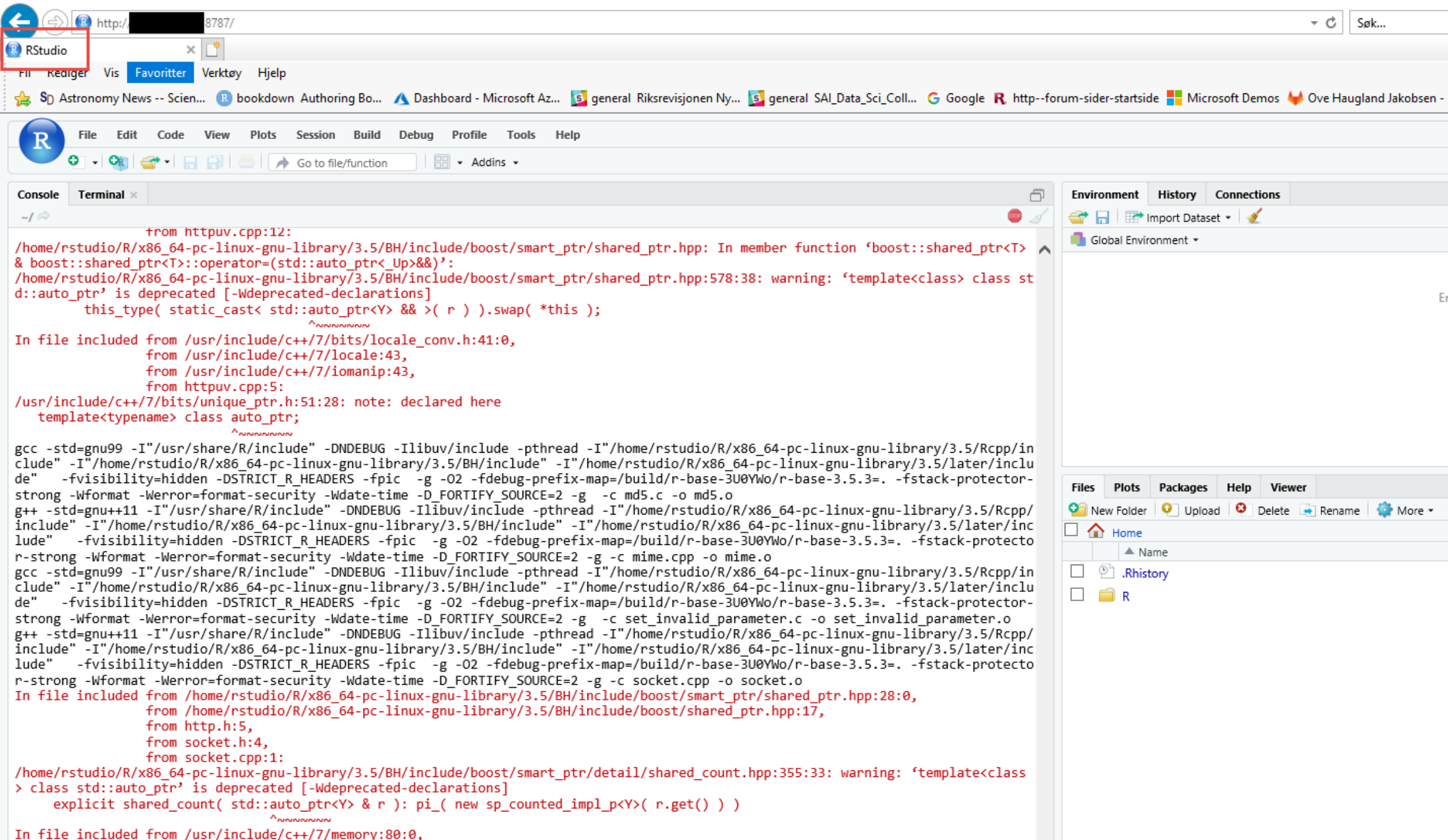
The logo for Amazon Aurora Backtrack features a stylized database icon on the left, consisting of three stacked cylinders with horizontal lines and a small cloud icon above them. To the right of the database icon is a circular clock face with a large 'L' in the center, indicating a time-related function. The entire logo is rendered in a light blue color.

**AWS Database Migration Service**  
Join the 70,000+ databases already migrated and converted

# Cloud computing

- You need heavy computing machinery to run machine learning on large datasets
- Buy an on-prem machine for \$ 10 000?
- Or rent one in the cloud for \$ 50 a week?





# Summary: Some success factors

- We were given full freedom to experiment
- Recruited people from audit, not IT
- There were no detailed planning — only the “what”, not the “how”
- Full support from top management
- Short development cycles # agile development
- Prioritise solutions to long standing issues

[www.saireports.org](http://www.saireports.org)

A final snack before we adjourn