

The background is a dark, rustic workshop. A wooden shelf holds several cylindrical objects, possibly spools or tools. Various other tools and equipment are visible in the background, creating a sense of a traditional craft or manufacturing environment.

# **Introduction to Data Science**

***Amit Kapoor***

***Bargava Subramanian***

A person is holding a lit sparkler, with bright sparks emanating from the tip. The person's face is partially visible in the background, and they are wearing a dark, textured garment. The overall scene is dimly lit, with the primary light source being the sparkler.

**Welcome**



**Amit Kapoor**

***@amitkaps***



**Bargava**

**@*bargava***

A stack of books is shown, with the word 'Agenda' overlaid in white text. The books are of various colors and thicknesses, and the text is centered over the middle of the stack.

# Agenda









**Be Curious | Have Fun**







# How many birds remain on the tree?

- $P(\text{Hunter to hit target}) = 0.2$
- Number of birds  $n = 150$
- Shots = 3
- Birds hit = 1



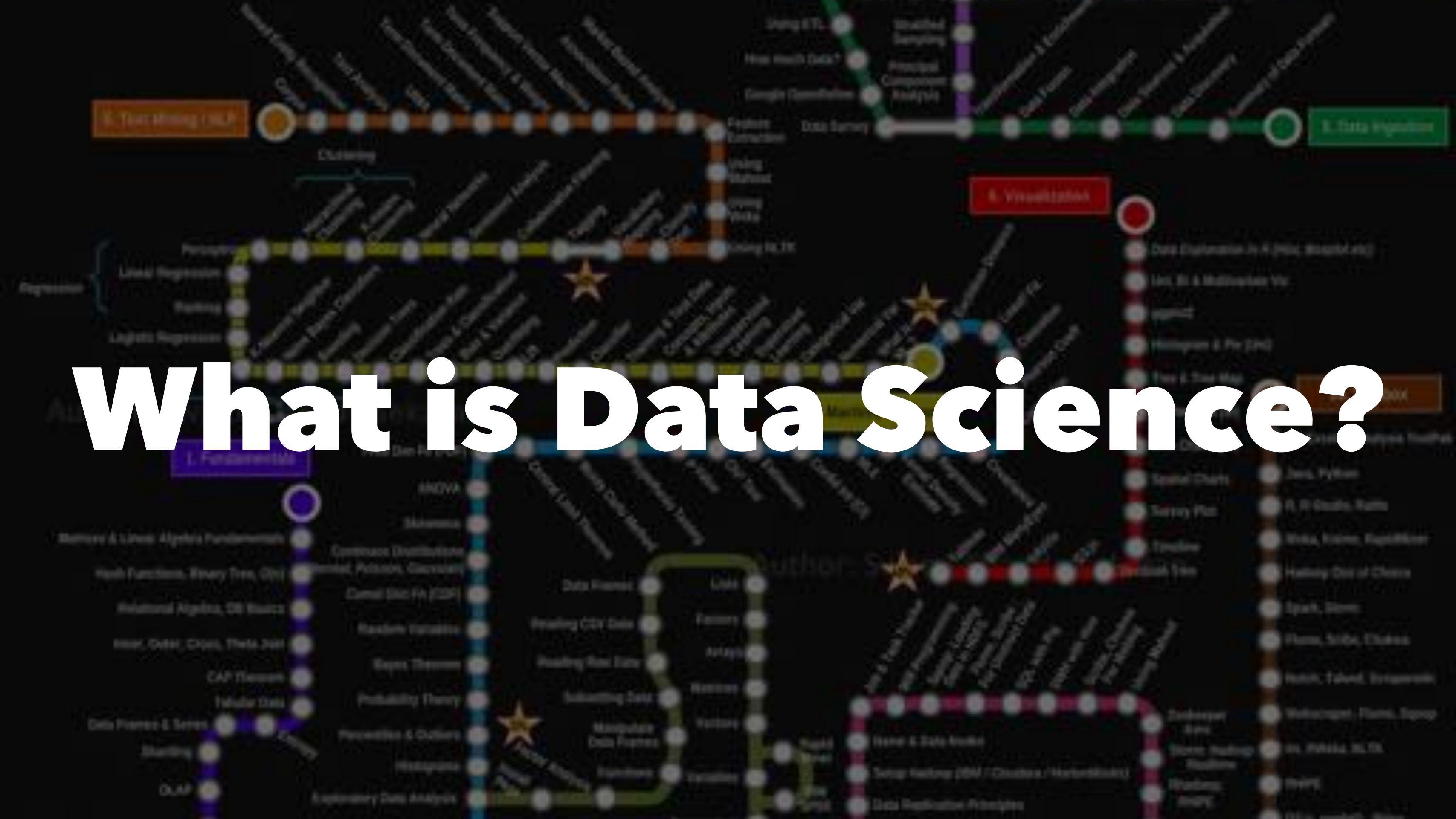


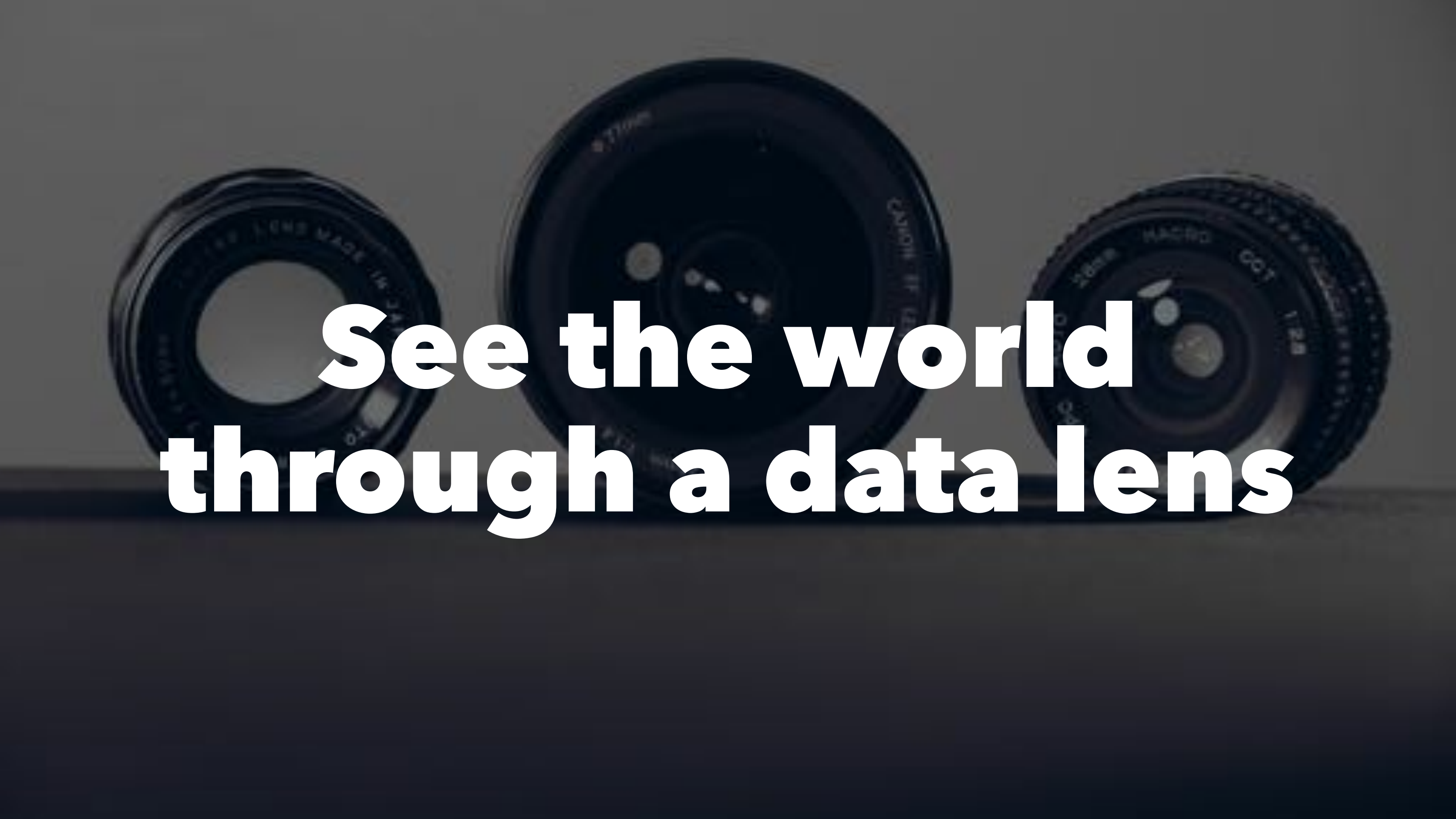
**Domain knowledge is very important**

*Don't lose the big picture*

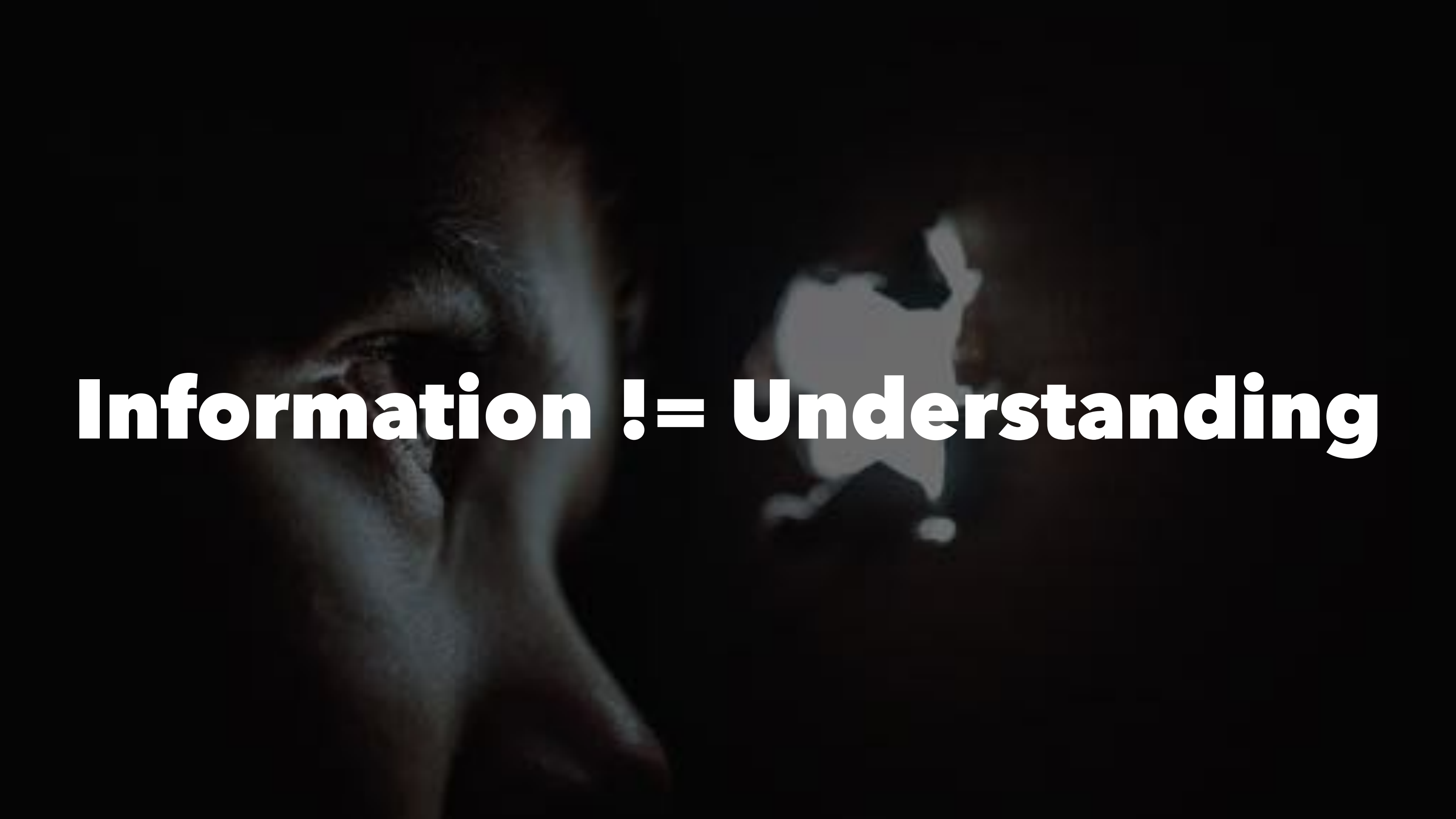


# What is Data Science?





**See the world  
through a data lens**



**Information != Understanding**



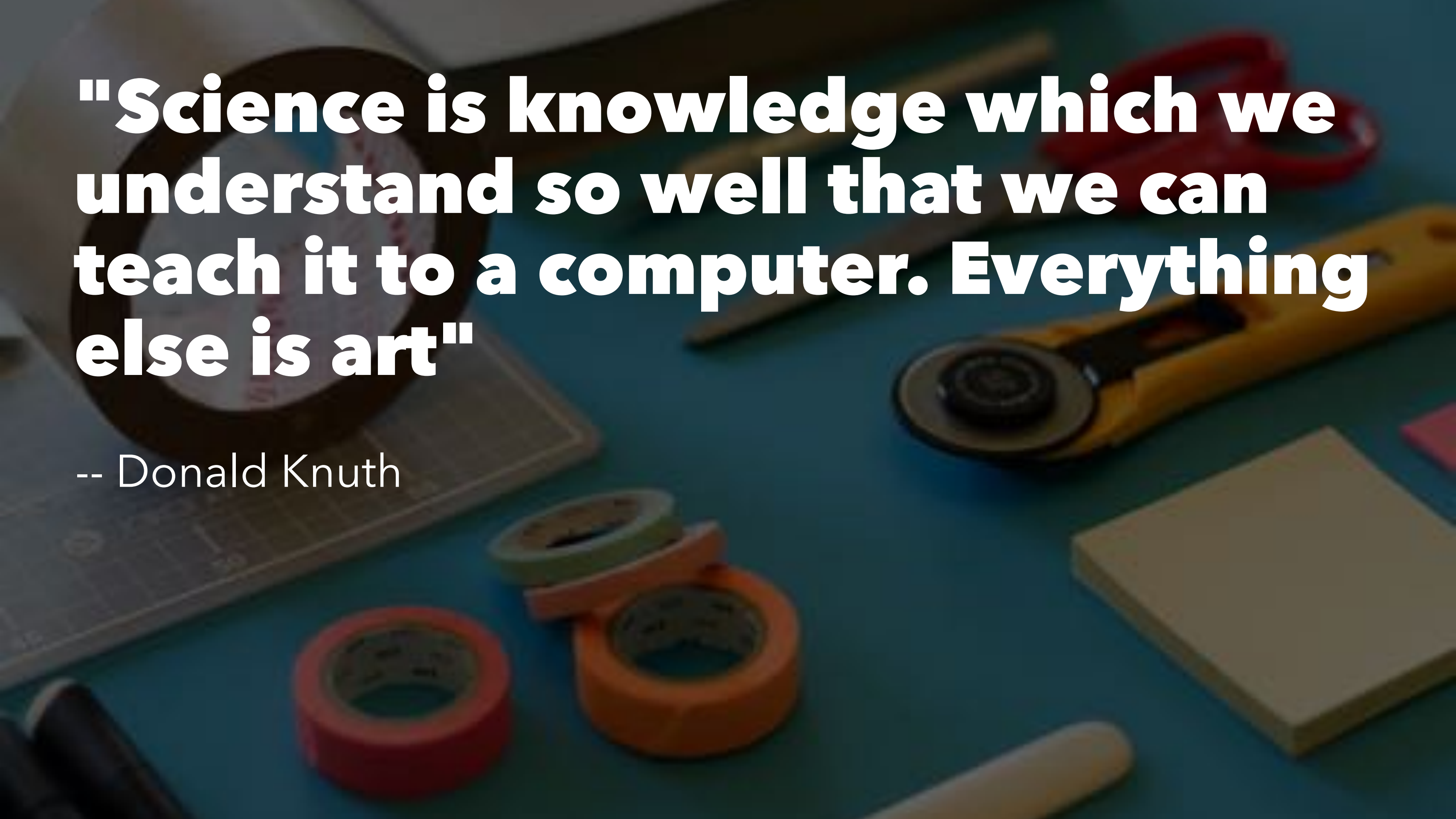


**"Data is just a clue to the end truth"**

-- *Josh Smith*



# Data Driven Decisions



**"Science is knowledge which we understand so well that we can teach it to a computer. Everything else is art"**

-- Donald Knuth





**Data Science is an Art**

A person is playing a violin, with their hand visible on the bow. The background is dark and out of focus, showing some papers or documents. The text "Hypothesis Driven Approach" is overlaid in white, bold, sans-serif font.

# **Hypothesis Driven Approach**





# Frame

**"An approximate answer to the right problem is worth a good deal"**



# Frame

- Toy Problems
- Simple Problems
- Complex Problems
- Business Problems
- Research Problems



# Types of Questions

- Descriptive
- Exploratory
- Inferential
- Predictive
- Causal
- Mechanistic





# **Descriptive**

How many people signed up for today's event?





# Exploratory

How many people turned up from each company ?





# Inferential

How many managers signed up and amongst them how many turned up? Is there a correlation with the company they come from?

# Predictive

Given that Ms. X has come for the previous workshop, will she turn up for today's workshop?



# Causal

Why did more people from a particular company turn up for today's workshop?

# **Mechanistic**

If more managers from a particular company turn up, how many more employees will turn up for the workshop?



# **Acquire**

**"80% perspiration, 10% great idea, 10% great output"**



# Acquire

- **Scraping** (structured, unstructured)
- **Files** (csv, xls, json, xml, pdf, ...)
- **Database** (sqlite, ...)
- APIs
- Streaming





**Refine**

**"All data is messy."**

# Refine



- Data Cleaning (inconsistent, missing, ...)
- Data Refining (derive, parse, merge, filter, convert, ...)
- Data Transformations (group by, pivot, aggregate, sample, summarise, ...)





# **Explore**

**"I don't know, what I don't  
know."**



# Explore

- Simple Vis
- Multi Dimensional Vis
- Geographic Vis
- Large Data Vis (Bin - Summarise - Smooth)
- Interactive Vis



# **Model**

**"All models are wrong, but some are useful"**



# Model - Supervised Learning

- *Continuous: Regression*
- *Discrete: Classification*

# Model - Supervised Learning

- *Continuous: Regression :*  
What will be the annual revenue for 2017 ?



# Model - Supervised Learning

- *Discrete: Classification :*  
Will company XYZ buy from us?

# Model - Supervised Learning

- *Continuous: Regression*
- *Discrete: Classification*

## Algorithms:

Linear Regression, Logistic Regression, CART, Random Forest, Gradient Boosting Machines, K-Nearest Neighbor, Support Vector Machines, Naive-Bayes, Bayesian Networks



# Model - Unsupervised Learning

- *Cluster Analysis*
- *Dimensionality Reduction*

# Model - Unsupervised Learning

- *Cluster Analysis* :  
If we segment our customers into three types, what would they look like?



# Model - Unsupervised Learning

- *Dimensionality Reduction* :  
Data is too huge to load into memory. Is there a better representation of the data?

# **Model - Unsupervised Learning**

*Cluster Analysis:* K-means, DBSCAN

*Dimensionality Reduction:* Principal Component Analysis,  
Singular Value Decomposition, MDS



# Model - Advanced / Specialized

- Deep Learning
- Network / Graph Analytics
- Optimization
- Reinforcement Learning
- Online Learning
- Applications: Time Series, Text, Image, Speech



A rustic log cabin with a chimney, nestled in a dense forest of tall trees. The scene is dimly lit, suggesting dusk or dawn, with a soft glow from the cabin's interior. The text is overlaid in white, bold font.

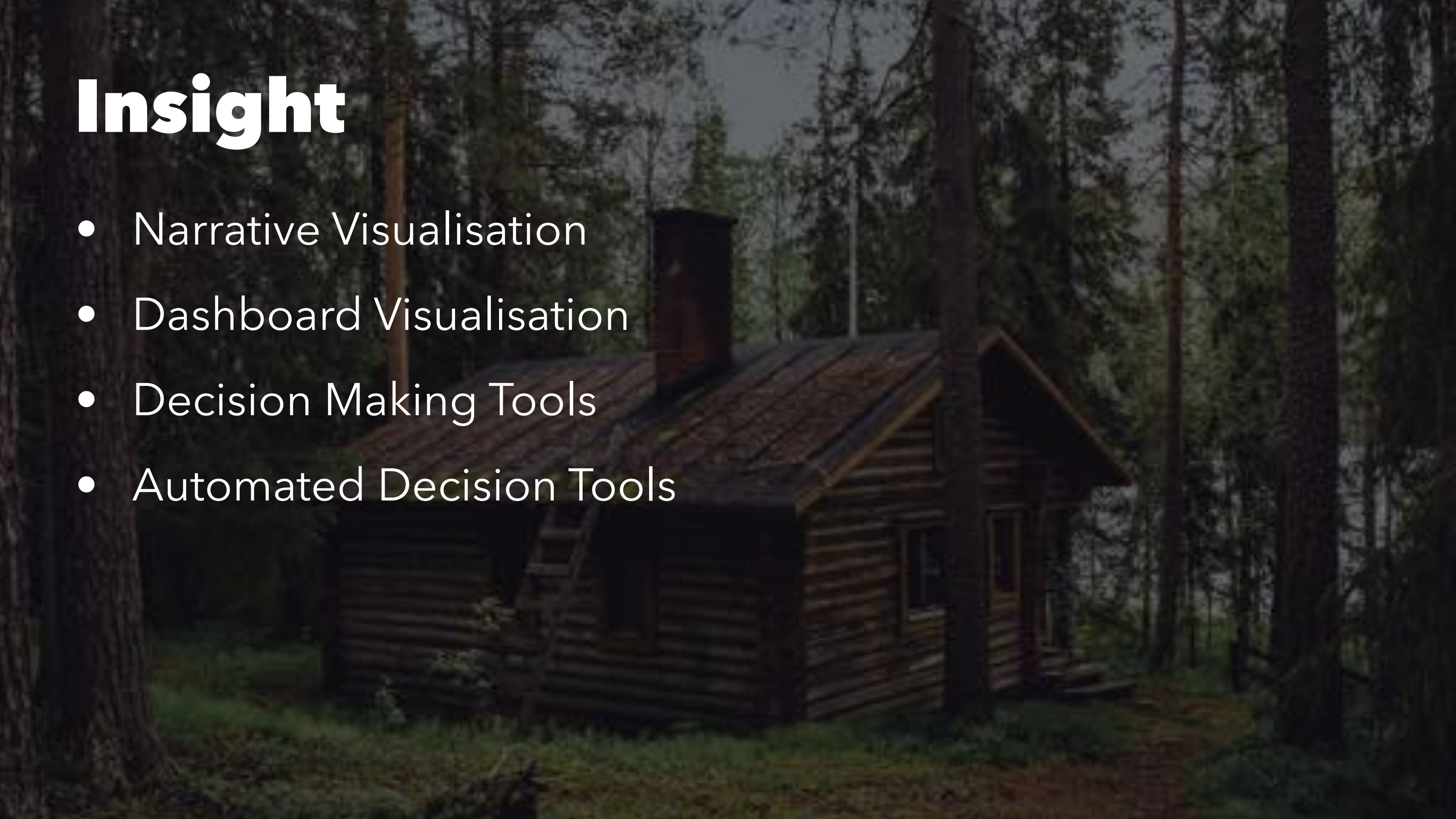
# **Insight**

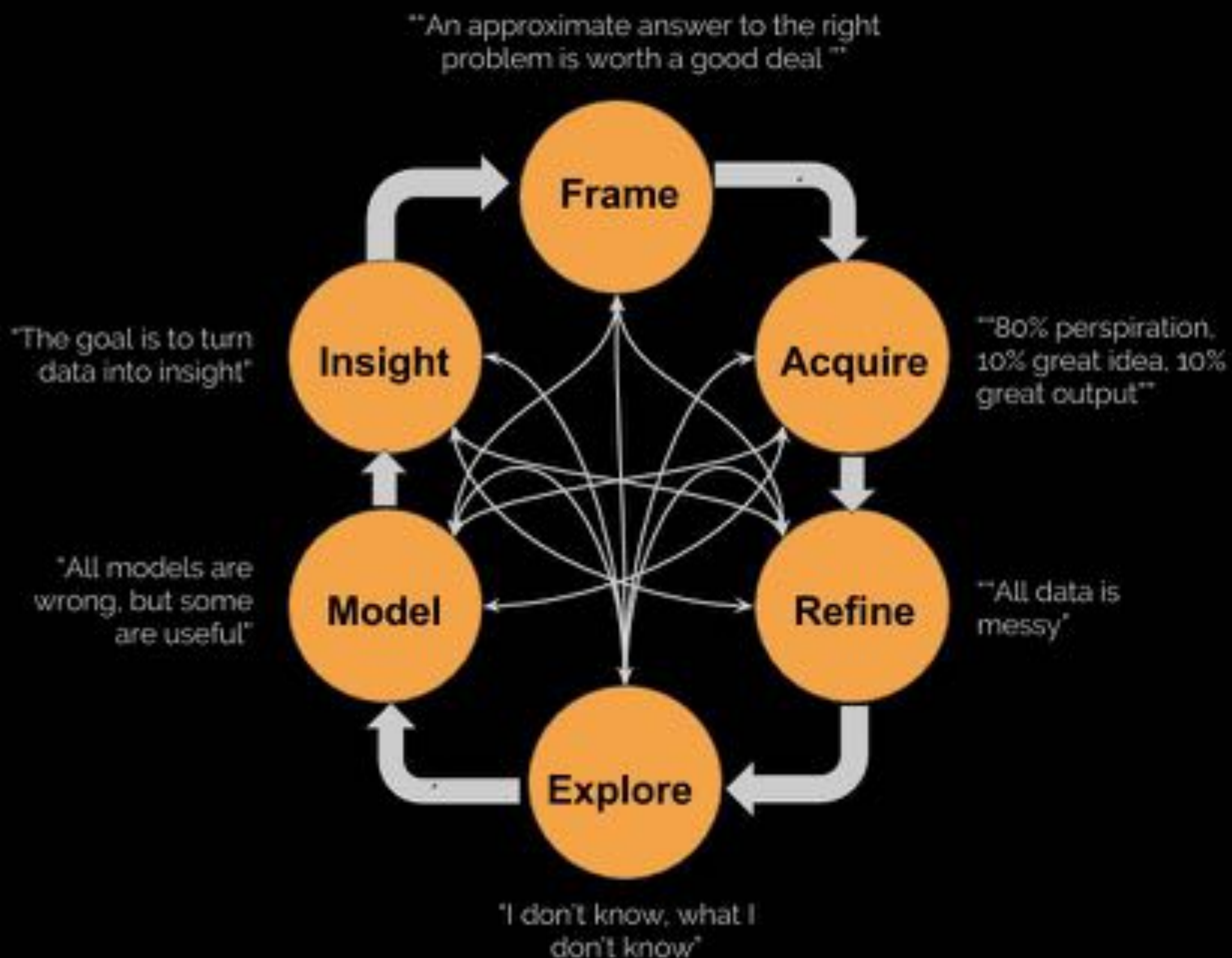
**"The goal is to turn data into insight"**



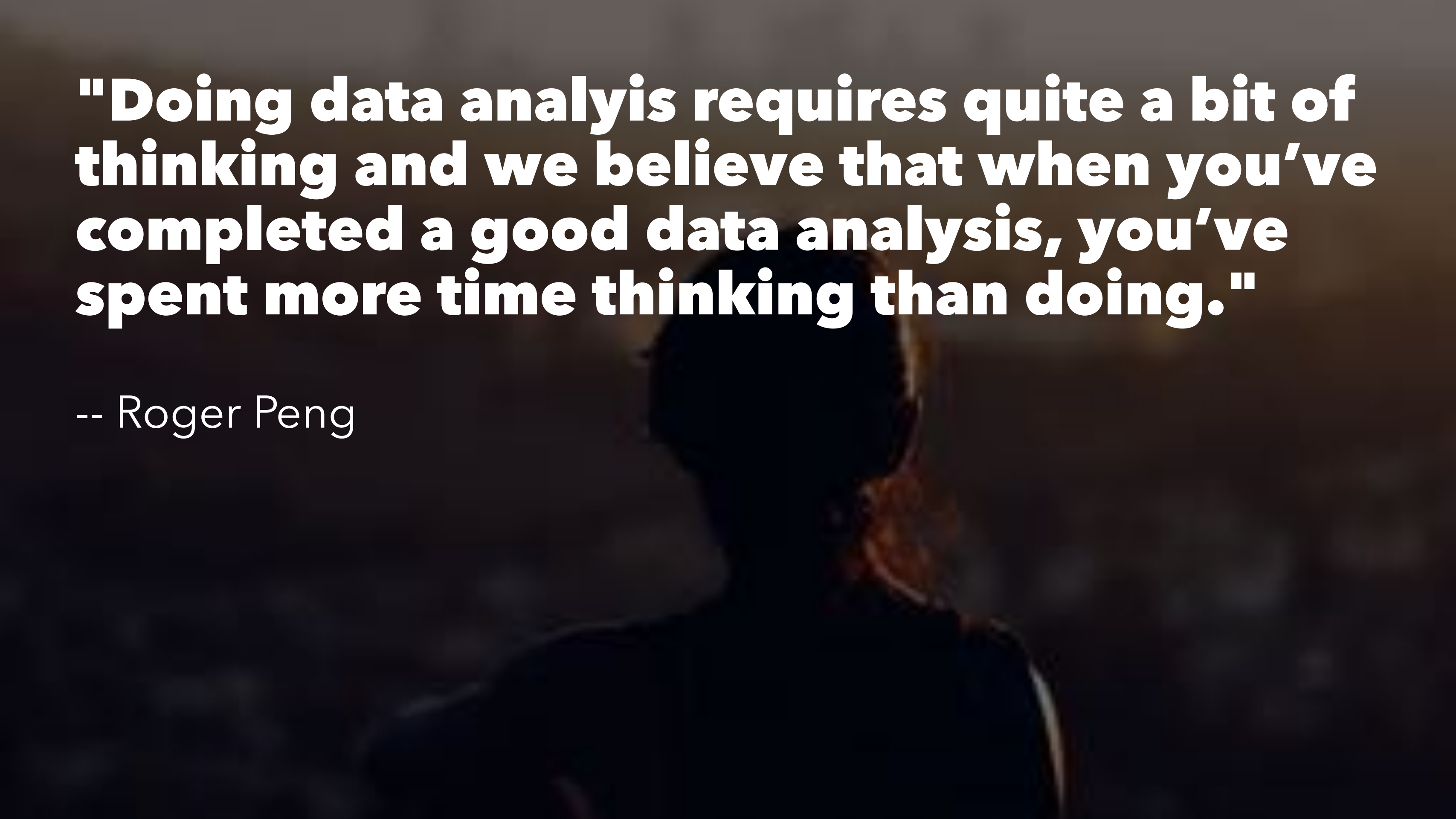
# Insight

- Narrative Visualisation
- Dashboard Visualisation
- Decision Making Tools
- Automated Decision Tools





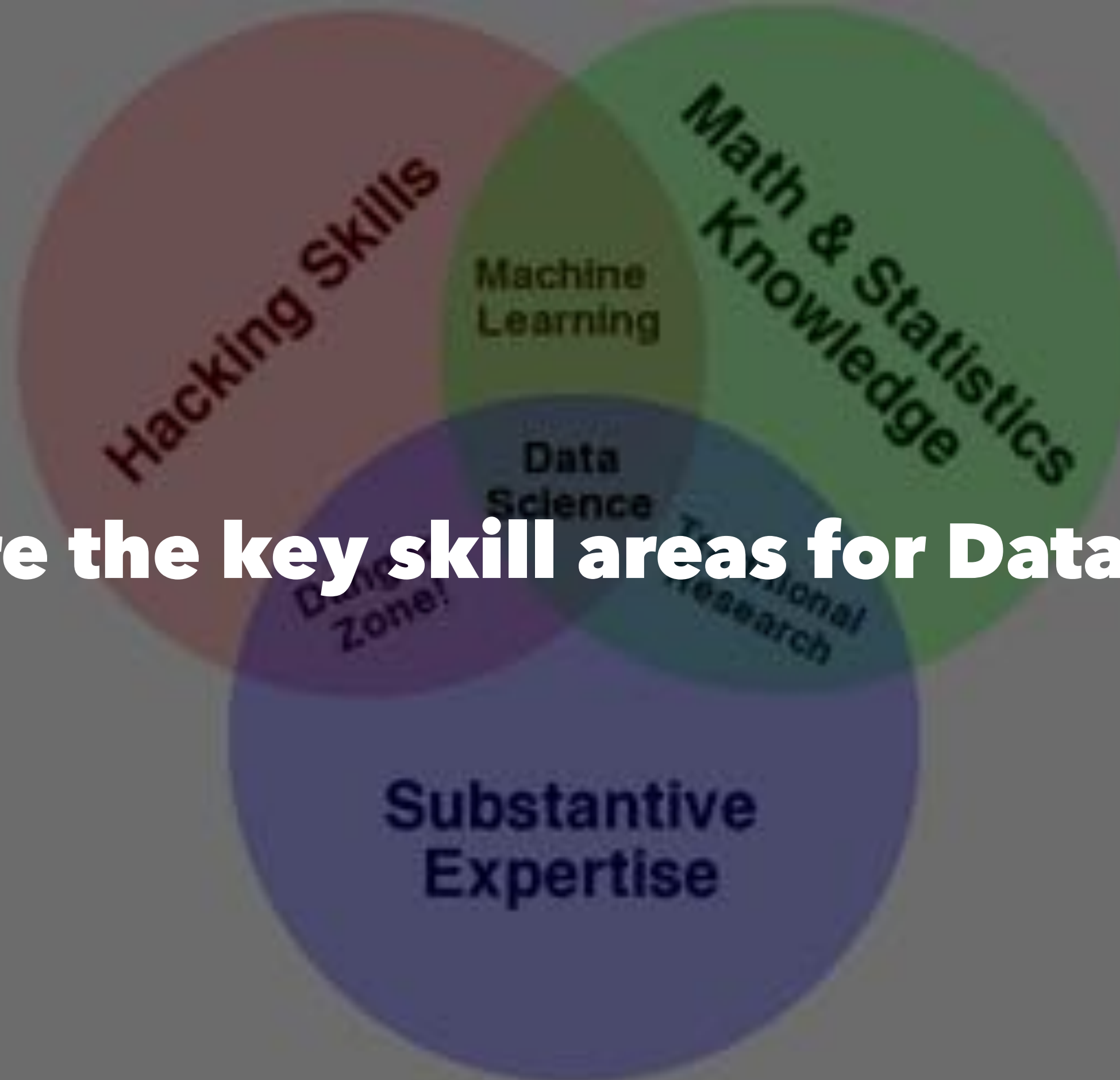


A person in a dark suit is seen from the side, looking out at a city at night. The city lights are visible in the background, creating a bokeh effect. The person's face is in shadow, and they appear to be looking towards the right side of the frame.

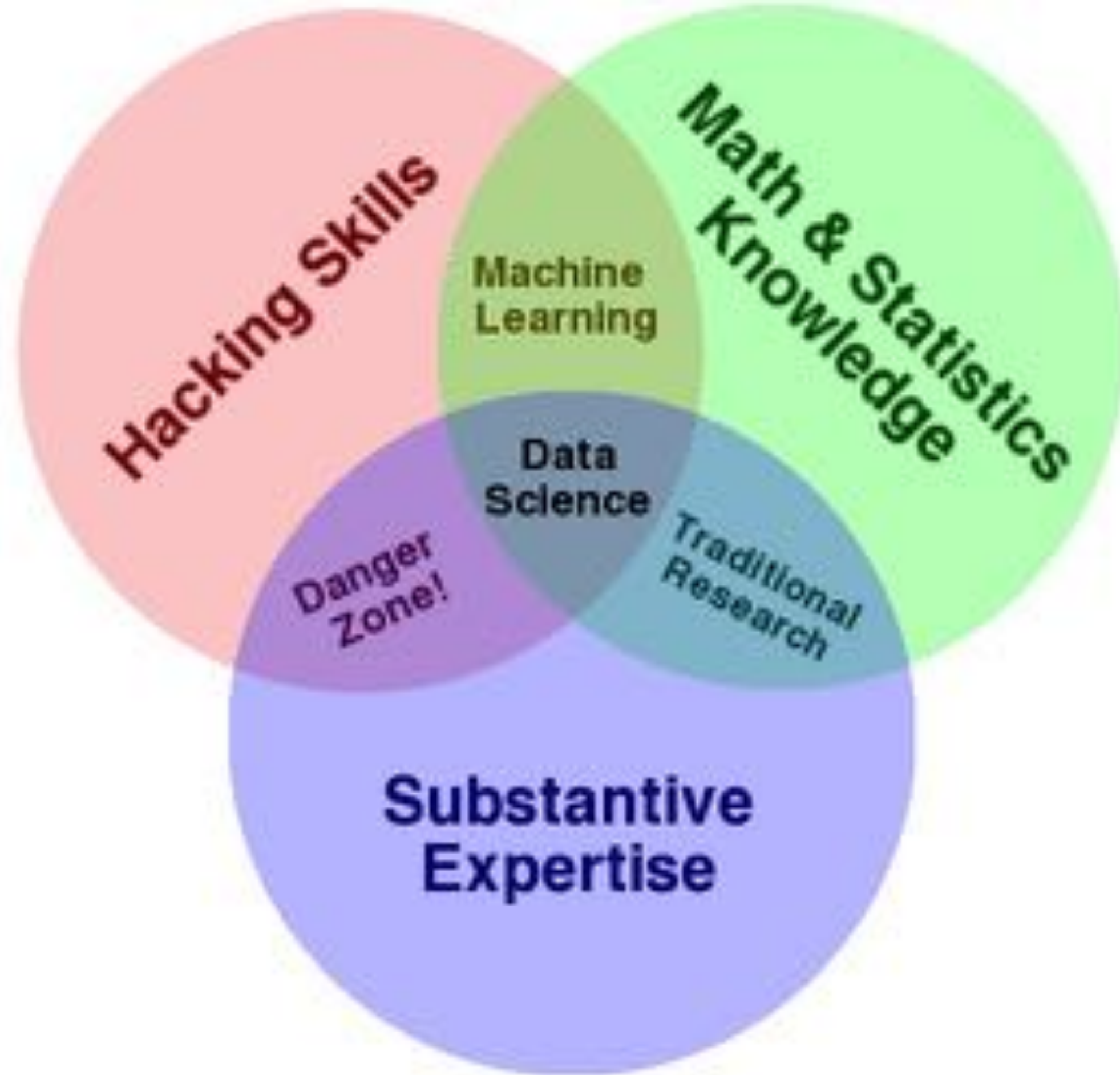
**"Doing data analysis requires quite a bit of thinking and we believe that when you've completed a good data analysis, you've spent more time thinking than doing."**

-- Roger Peng

# What are the key skill areas for Data Science







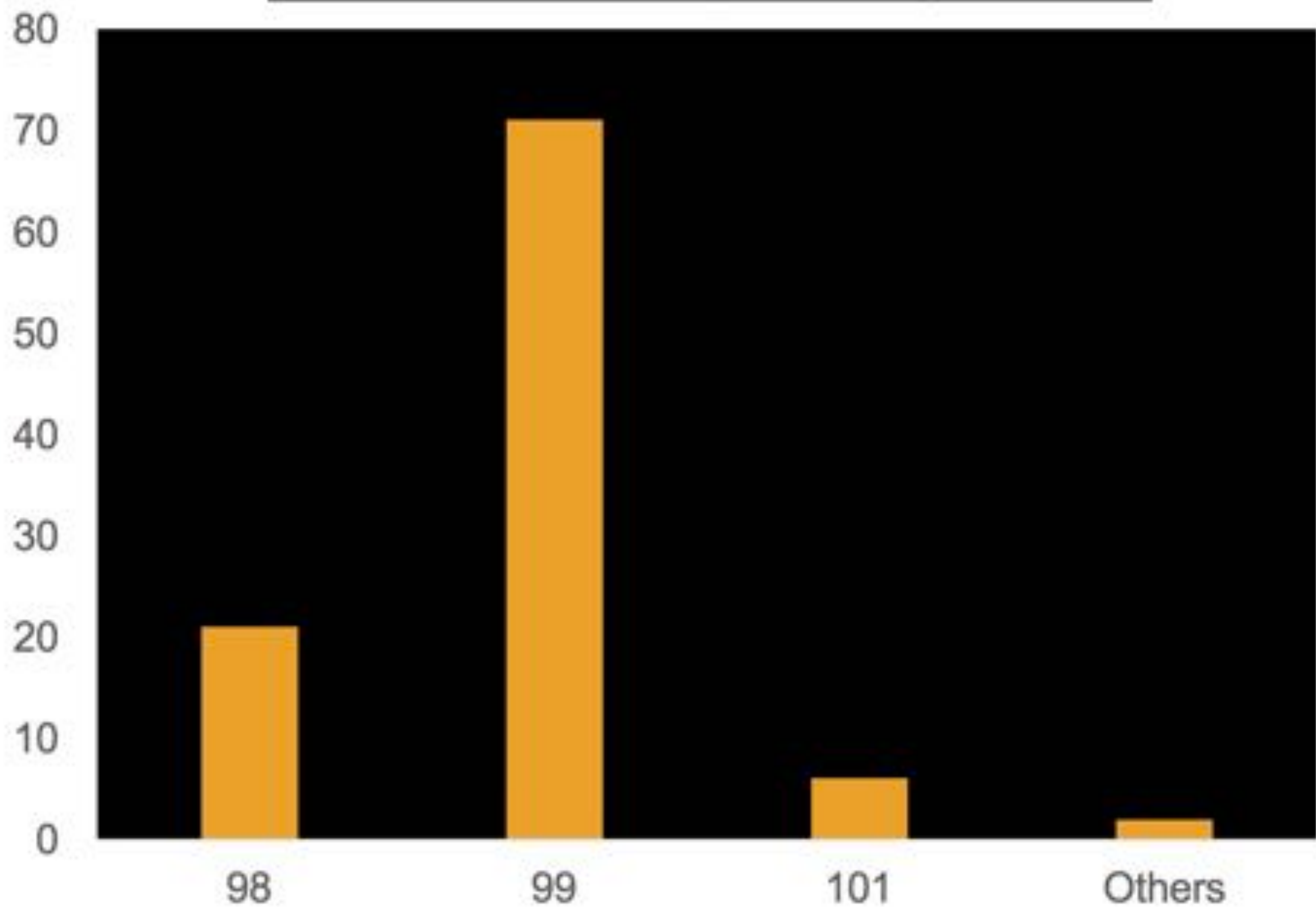
# SOLVE THE PUZZLE



3	2	=	7
5	4	=	23
7	6	=	47
9	8	=	79
10	9	=	?



**Distribution of Responses by Answer**



# SOLVE THE PUZZLE

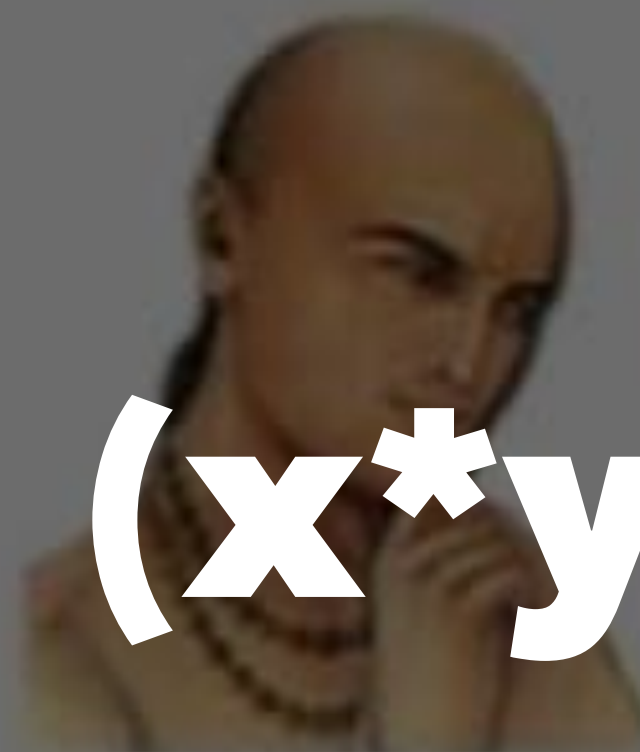
3	2	=	7
5	4	=	23
7	6	=	47
9	8	=	79
10	9	=	?

**99**  
**( $x*y$  + odd counter)**



# SOLVE THE PUZZLE

3	2	=	7
8	4	=	23
7	6	=	47
9	8	=	79
10	9	=	?



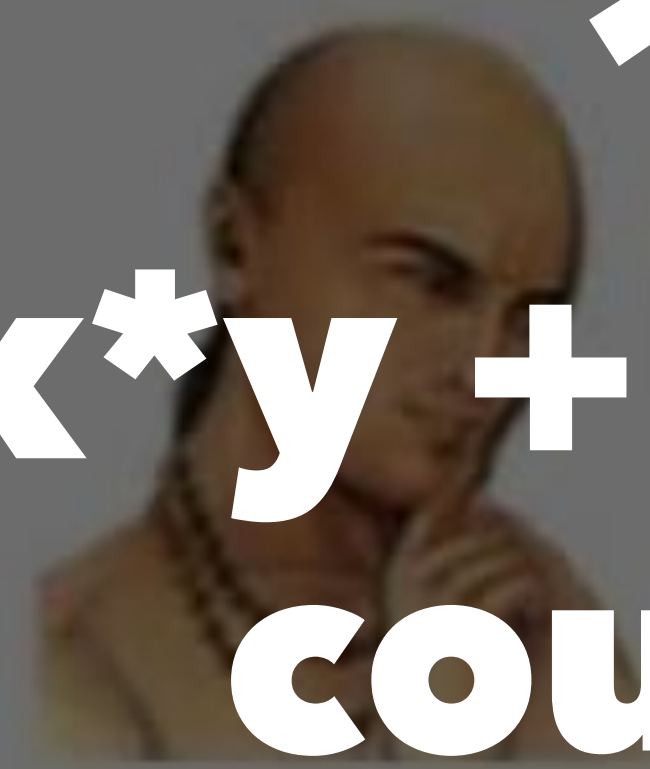
98

$(x * y + y = 1)$

# SOLVE THE PUZZLE

101

( $x * y + \{1, \text{prime counter}\}$ )



3	2	=	7
5	4	=	23
7	6	=	47
9	8	=	79
10	9	=	?



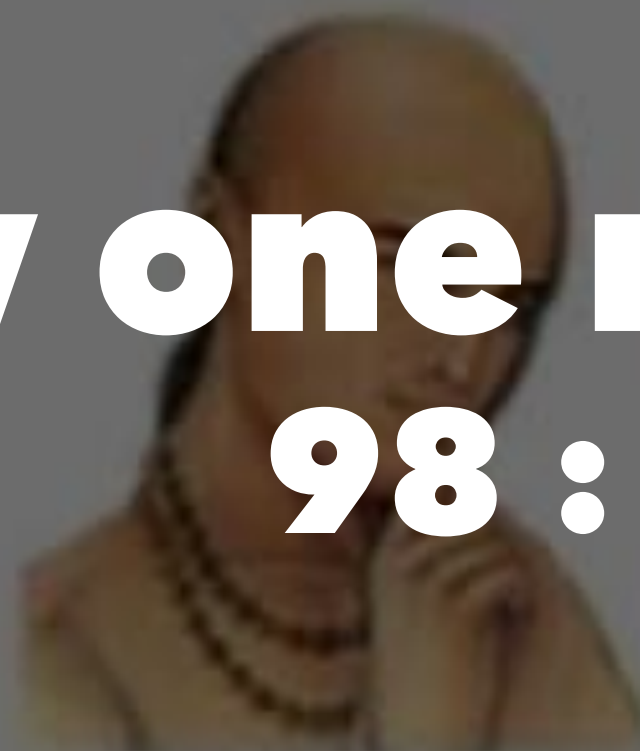


***Stakeholder/Client wants an answer***  
**Wisdom of crowds?**

# SOLVE THE PUZZLE

3	2	=	7
5	4	=	23
7	6	=	47
9	8	=	79
10	9	=	?

**Throw one more to the fire**  
 **$98 : (x^2 - 2)$**





A pair of dark, textured razors are crossed at their handles, forming a large triangle. The blades are pointed towards the top center of the frame. The background is a solid dark gray.

# Occam's Razor

*Problem solving principle: Amongst competing hypotheses, the one with the fewest assumptions should be selected*



# **Different Profiles in Data Science**



# Different Profiles in Data Science

- Data Analyst
- Data Engineer
- Data Visualization
- Data SME
- Data Scientist



# Data Analyst

A budding/junior data scientist. Supports EDA and data wrangling.

*Typical skills:*

- Good with Excel.
- Basic knowledge of R/Python/SQL



# Data Engineer

Builds and supports the data pipeline. Data Architect.

*Typical skills:*

- SQL
- Spark
- Hadoop/Cassandra
- Data Orchestration(Eg: Luigi)

# Data Visualization



Builds visualization

*Typical skills:*

- Tableau/Qlik/D3.js
- Basics of HTML/JavaScript



# Data SME

The data guru. Understands business impact of each of the attribute stored in the system.

*Typical skills:*

- Domain knowledge + system architecture

**Data Scientist**







# SO WHAT DOES A DATA SCIENTIST DO?

## Data Scientist



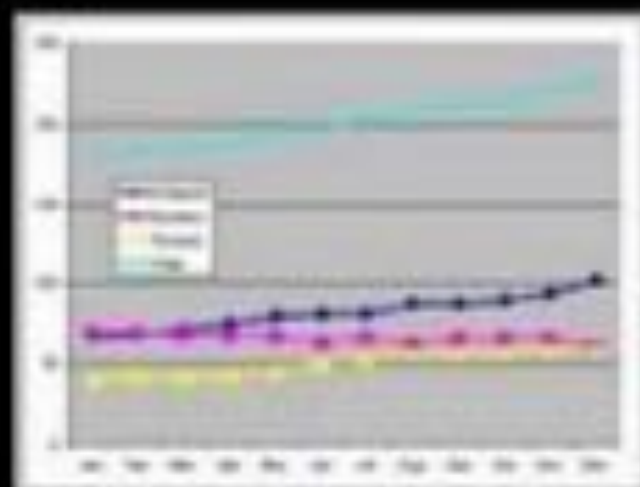
What my friends think I do



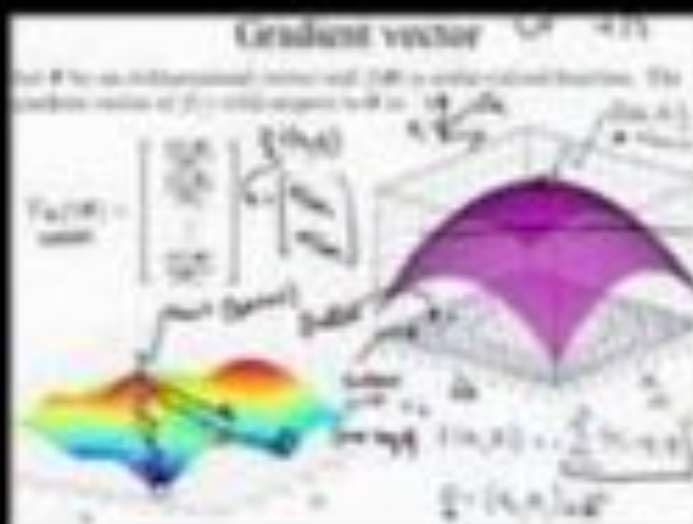
What my mom thinks I do



What society thinks I do



What my boss thinks I do



What I think I do



What I actually do



# Data Scientist

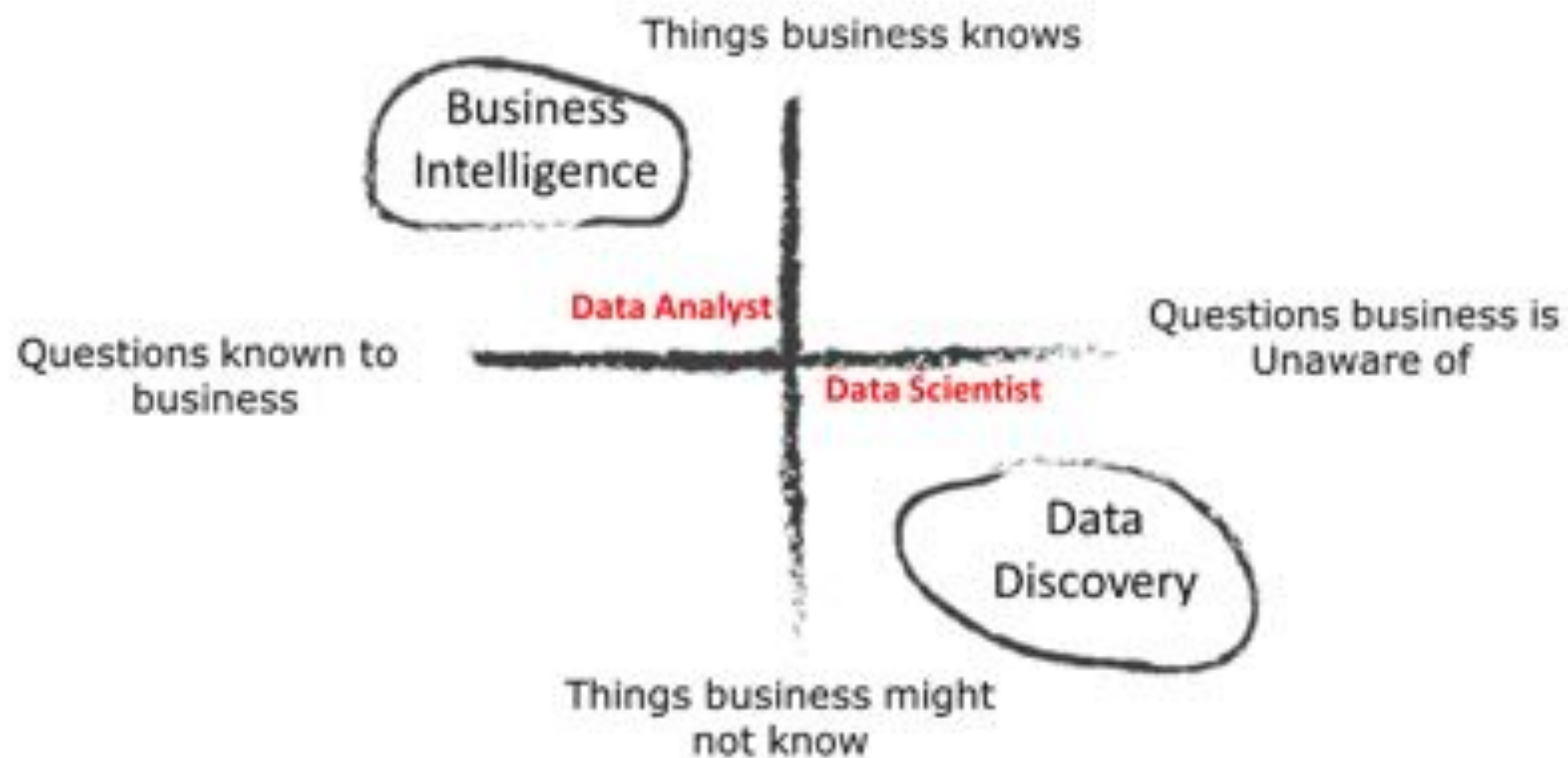
Builds models and find insights.

*Typical skills:*

- R
- Python
- Spark
- Machine Learning
- Math/Statistics

# **Data Analyst Vs Data Scientist?**







Chicago Daily Tribune  
**DEWEY**

**DEFEATS**

**TRUMAN**



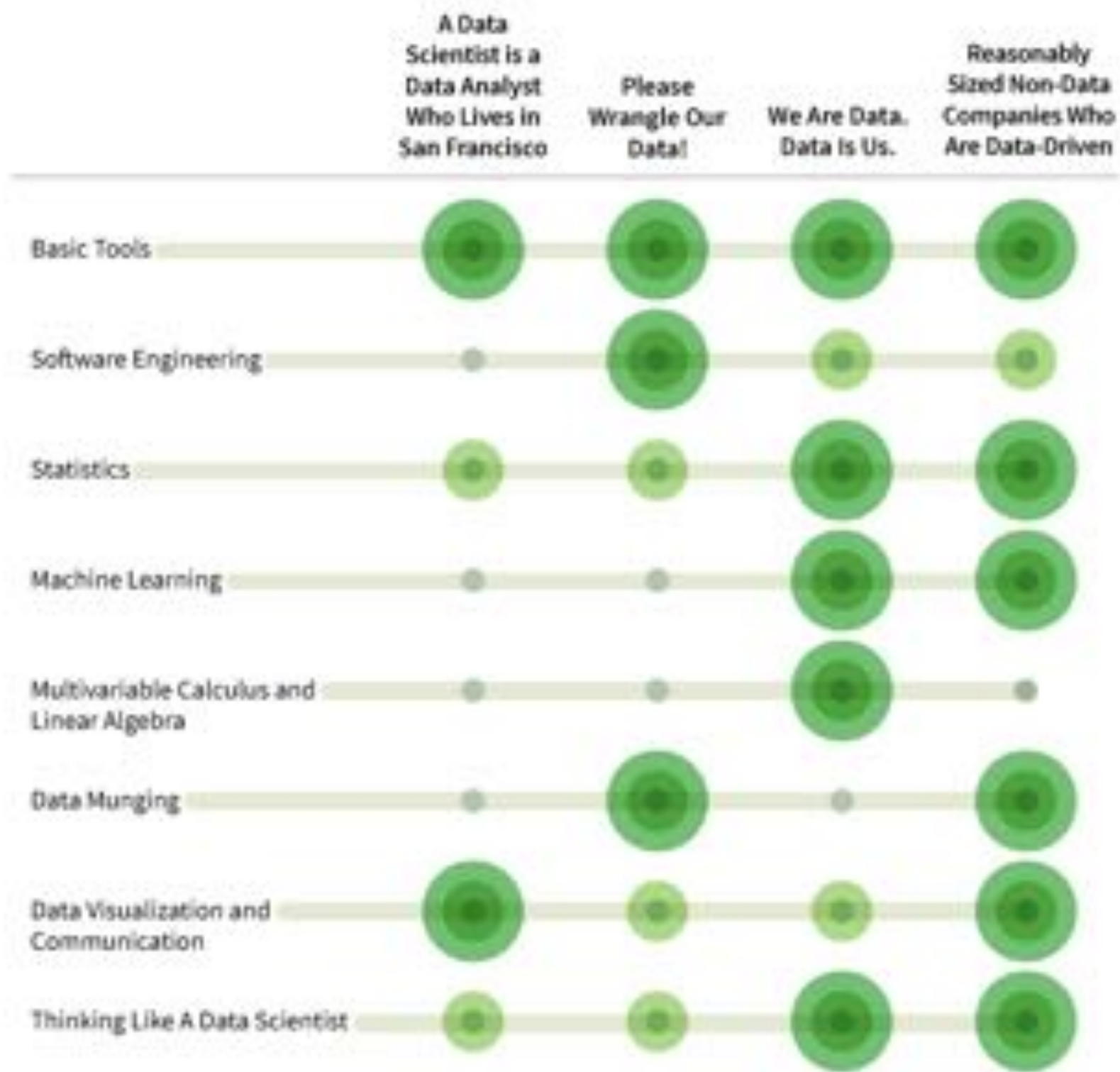
A black and white photograph of Dwight D. Eisenhower, smiling and holding a newspaper. The newspaper's headline reads "DEWEY DEFEATS TRUMAN". The image is used to illustrate the concept of selection bias.

# **Selection Bias !**

# skills







Very important



Somewhat important



Not that important

# R Stack

- **Acquire:** `rvest`, `XML`, `jsonlite`, `httr`, `RSQLite`, `RPostgreSQL`, `readxl`, `haven`, `readr`, `data.table`
- **Refine:** `dplyr`, `tidyr`, `lubridate`, `stringr`
- **Explore:** `graphics`, `ggplot2`, `ggvis`, `ggmap`, `map`, `vcd`, `rgl`, `htmlwidgets`, `leaflet`, `choroplethr`, `plotly`
- **Model:** `stats`, `caret`, `ranger`, `glmnet`, `xgboost`, `party`, `mxnet`, `forecast`
- **Insight:** `OpenCPU`, `Rserve`, `shiny`, `RMarkdown`, `knitr`



# PyData Stack

- **Acquire / Refine:** Pandas, BeautifulSoup, Selenium, Requests, SQL Alchemy, Numpy, Blaze
- **Explore:** Matplotlib, Seaborn, Bokeh, Plotly, Vega, Folium
- **Model:** Scikit-Learn, StatsModels, SciPy, Gensim, Keras, Tensor Flow, PySpark
- **Insight:** Django, Flask

# **A day in the life of a Data Scientist**



$$\sum_{k=1}^n p_k \quad y = y(x) \quad (x \in \mathbb{R}) \quad S(\alpha, t) = \frac{1}{t} \int_0^t \frac{1}{e} d\alpha \quad P(\eta_0 < x) = F(x)$$

$$W_k = \binom{n}{k} p^k (1-p)^{n-k} \quad P(\eta < y | \xi = x) = \sup_{y' < y, y' \neq 0} P(\eta < y' | \xi = x)$$

$$S_n = A_n U \Gamma A_n \quad \int_0^1 f(x) \log_2 \frac{1}{f(x)} dx < \varepsilon \quad g^{-1} \cdot g = e \quad \gamma = \sqrt{\frac{2n}{\gamma_n}} \left( \frac{\gamma_{2n}}{\gamma_n} + \frac{\gamma_n - \gamma_{2n}}{\gamma_{2n}} \right) \quad f(t|y) = \frac{2e^{\frac{y^2}{2}}}{12\pi} \int_{\frac{y}{\sqrt{2}}}^{\frac{y}{\sqrt{2}}} \frac{e^{-\frac{u^2}{2}} du}{\left(1 - \frac{y^2}{2u^2}\right)^{\frac{3}{2}}}$$

$$\Delta V = \sum_{n=1}^N \frac{E_n}{u} \quad \sum_{k=1}^n e^{-\frac{k^2 \pi^2}{n^2}} = H(n) \quad \prod_{k \leq b} \bigcup_{i=1}^{n-1} H_i; \bigcap_{n=0}^{\infty} X_n \quad f_n(t) = \frac{2^{n-1} (n-1)! e^{-2t}}{(n-1)!} \quad H_r(x) = \frac{G_r(x)}{1 + G_r(x)}$$

$$U_n^{+0} = \binom{2n}{n} - \binom{2n}{n-1} \quad \int_0^1 dG_k(x) \geq \frac{1}{2} \quad \lim_{t \rightarrow \infty} \frac{f(t)}{t} = P_e \quad R = \int_{-\infty}^{\infty} p(t) dt \quad \frac{\sin t_k}{t_k} [\varphi(t) e^{-itx} + \varphi(-t) e^{itx}]$$

$$\frac{1}{n} \varphi\left(\frac{n}{m}t\right) = \varphi\left(c\left(\frac{n}{m}\right)t\right) \quad \log \varphi(t) = i\gamma t - c|t|^k \left[1 + i\beta \frac{k}{|t|} \omega(t)\right] \quad \beta(u) = \sum_{k=1}^r \Psi^*(b_k u) \quad \lim_{n \rightarrow \infty} P\left(\frac{\sum_{i=1}^n a_{ij} b_{ij}}{\sqrt{\frac{1-g}{g}}}\right) C_n(x) \geq \frac{n!}{\prod_{k=1}^n n_k(k)!}$$

$$\int_0^{\infty} e^{-\frac{u^2}{2}} du = F(x) \left(\frac{1}{\sqrt{2\pi}}\right)^{-1} \quad |\Psi_S(H)| = \left| \int_0^{\infty} e^{itx} dF(x) \right| \leq \int_0^{\infty} e^{-ux} dF(x) = \varphi_S(iu) \quad g^{-1}Ng = \{g^{-1}ng | n \in N\} \quad Q = F^{-1}(q) \quad q_A(x) = \sum_{j=1}^r p_j^x \quad P(\Pi_2 =$$

$$\Pi_m = \Pi_r \Pi_{m-r} \quad |X \cup Y| = |X| + |Y| - |X \cap Y| \quad \lim_{n \rightarrow \infty} \frac{1}{n} k_n\left(\frac{x}{\sqrt{n}}\right) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} \quad P_n(k) = \frac{c_n!}{k!} \quad P\left(\limsup_{n \rightarrow \infty} \frac{|h_n|}{\sqrt{2n \log \log n}} \leq 1\right) = 1 \quad (A \cap B) = 1 - \sqrt{1 - e^{2B}}$$

$$f: X \rightarrow X \cap W \quad \mathcal{H}(A) = \int_A \mathcal{H}(\omega) d\mathbb{P} \quad l'(x) = -\log_2 \left( \frac{\sum_{k=1}^r P_k^x \log_2 \frac{1}{P_k}}{\sum_{k=1}^r P_k^x} - \left( \frac{\sum_{k=1}^r P_k^x \log_2 \frac{1}{P_k}}{\sum_{k=1}^r P_k^x} \right)^2 \right) \quad f g(u_i) = f\left(\sum_{j=1}^{\dim V_k} a_{ji} v_j\right) = \sum_{j=1}^{\dim V_k} a_{ji} \left(\sum_{k=1}^{\dim V_k} b_{kj} u_k\right) \frac{\binom{2k}{k}}{2^{2k}} \approx \frac{1}{\sqrt{2k}}$$

$$P_{j,k}^{(m)} = \sum_{c=0}^{\infty} P_{j,c}^{(r)} P_{c,k}^{(m-r)} \quad \frac{1}{2\pi} \int_{-\infty}^{\infty} \operatorname{Re} \left\{ \varphi(t) \frac{e^{itx} - e^{-itx}}{it} \right\} dt \quad \frac{1}{2^{2k}} \approx \frac{1}{\sqrt{2k}} \quad \frac{1}{2\pi} \int_{-\infty}^{\infty} \operatorname{Re} \left\{ \varphi(t) \frac{e^{itx} - e^{-itx}}{it} \right\} dt$$

$$\liminf_{N \rightarrow \infty} \int_0^N f_N(x) dx \geq \int_0^{\infty} f(x) dx \quad M(1, \delta_j - 1)^2 = \int_0^{\infty} (x-1)^2 e^{-x} dx \quad \lim_{N \rightarrow \infty} \int_{-A}^A f_N(x) \log_2 \frac{1}{f_N(x)} dx = \int_{-A}^A f(x) \log_2 \frac{1}{f(x)} dx$$

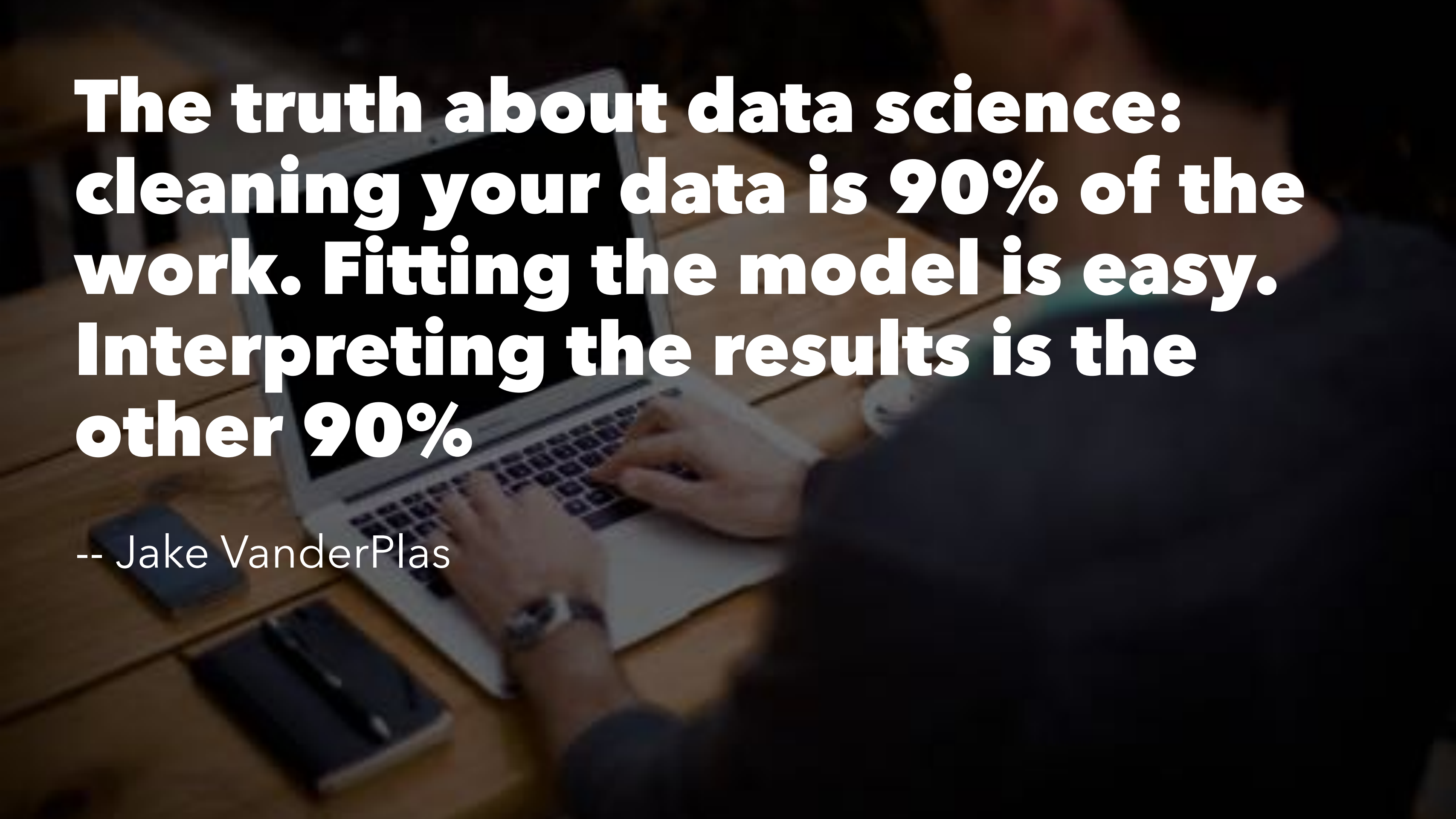
$$M_{n-k_k} = \binom{2n}{n+k_k} = \binom{2n}{n-k_k} \quad \lim_{N \rightarrow \infty} \int_{-A}^A f_N(x) \log_2 \frac{1}{f_N(x)} dx = \int_{-A}^A f(x) \log_2 \frac{1}{f(x)} dx \quad \lim_{N \rightarrow \infty} \int_{-A}^A f_N(x) \log_2 \frac{1}{f_N(x)} dx = \int_{-A}^A f(x) \log_2 \frac{1}{f(x)} dx$$





$$\begin{aligned}
 &= p(\sqrt{15}x^2 + 10x + 6) \quad \sum_{k=1}^{\infty} \int_{\mathbb{H}^n} \left( \int_0^1 \tilde{p}_k(\tau) d\tau \right) dt = x \int_0^1 \tilde{\psi}_k(\tau) d\tau = \frac{x^2}{2} B(x) + \int_0^1 (x-u) \sum_{k=1}^{\infty} \tilde{p}_k(u) du \\
 &(\lambda_1 \sigma_k^2 = \lambda; c_{ik} \quad \eta_1 = \sum_{k=1}^n a_{ik} \tilde{\eta}_k \quad \log \varphi(u) = -\frac{\sigma^2 u^2}{2} \quad i^2 = -1; j^2 = -1; k \in \mathbb{N}) \\
 &y = \phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{t^2}{2}} dt \quad S(\alpha, T) = \frac{2}{\pi} \int_0^{\pi} \frac{\sin \alpha t}{t} dt \quad P(\eta_m < x) = F(x) \\
 &\binom{n}{k} p^k (1-p)^{n-k} \quad P(\eta < y | \xi = x) = \sup_{r < y, y < r} P(\eta < y | \xi = x) \\
 &\log \left| \frac{1}{f(x)} dx \right| < \varepsilon \quad g^{-1} \cdot g = e \quad \gamma = \left[ \frac{2u}{\sqrt{n}} \left( \frac{\gamma_{2n}}{\sqrt{2n}} + \frac{\gamma_n - \gamma_{2n}}{\sqrt{2n}} \right) \right] \quad f(t|y) = \frac{2e^{\frac{y^2}{2}}}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \frac{e^{-\frac{t^2}{2}}}{\sqrt{1-t^2}} dt \\
 &e^{-\frac{h^2 \eta^2}{2}} = H(h) \quad \prod_{k \leq b}; \bigcup_{i=1}^{n-1} M_i; \bigcap_{n=0}^{\infty} X_n \quad f_n(t) = \frac{2^{-(n-1)} e^{-2t}}{(n-1)!} \quad H_r(x) = \frac{G_r(x)}{1+G_r(x)} \\
 &(e-u) du = \frac{2^{n+1} e^n e^{-2b}}{n!} \quad \lim_{t \rightarrow 0} (e^t) = 0 \quad \lim_{n \rightarrow \infty} \frac{3(n)}{n} = P_e \quad R = \int_{-\infty}^{\infty} p(t) dt \\
 &|t|^k \left[ 1 + i \beta \frac{k}{|t|} \omega(t) \right] B(\omega) = \sum_{k=1}^r \Psi^*(b_k v) \quad C_{iv} = \sum_{j=1}^n a_{ij} b_{jv} \quad \lim_{n \rightarrow \infty} P \left( \frac{3n - k(n) - \log \frac{1}{q}}{\sqrt{1-\frac{q}{2}}} \right) C_n(x) \geq \frac{n!}{\prod_{k=1}^n n_k(x)!} \\
 &\left( \frac{1}{\sqrt{2\pi}} \right)^n |\Psi_S(H)| = \left| \int_{-\infty}^{\infty} e^{itx} dF(x) \right| \leq \int_{-\infty}^{\infty} e^{-ux} dF(x) = \varphi_S(iv) \quad g^{-1}Ng = \{g^{-1}ng | n \in N\} \quad Q = F^{-1}(q) \\
 &|-|x \cap \Psi| \quad \lim_{n \rightarrow \infty} \frac{1}{n} k_n \left( \frac{x}{\sqrt{n}} \right) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} \quad P_n(k) = \frac{C(n)}{2^n} P \left( \lim_{n \rightarrow \infty} \sup \frac{|h_n|}{\sqrt{2n \log \log n}} \leq 1 \right) \\
 &x \cap w \quad \log_2 \left( \frac{\sum_{k=1}^r P_k^* \log_2 \frac{1}{P_k}}{\sum_{k=1}^r P_k^*} - \left( \frac{\sum_{k=1}^r P_k^* \log_2 \frac{1}{P_k}}{\sum_{k=1}^r P_k^*} \right)^2 \right) \quad f_g(u_i) = f \left( \sum_{j=1}^{dim V_k} a_{ji} v_j \right) = \sum_{j=1}^{dim V_k} a_{ji} \left( \sum_{k=1}^{dim V_k} b_{kj} \right) \\
 &\sqrt{\frac{q(1-q)}{n}} + o\left(\frac{1}{\sqrt{n}}\right) \quad \prod_{k=1}^r \left[ g_k \left( \frac{t}{\sqrt{12}} \right) \right]^{N_k \alpha_k} = e^{-\frac{t^2}{2}} \quad P_{jk}^{(m)} = \sum_{r=0}^{\infty} P_{jr}^{(r)} P_{ek}^{(m-r)} \quad \frac{1}{2\pi} \int_{-\infty}^{\infty} \operatorname{Re} \left\{ \varphi(t) \frac{e^{-itx}}{t} \right\} dt \\
 &\int_{\mathbb{R}} f(x) dx \geq \int_{\mathbb{R}} f(x)^n dx \quad \lim_{N \rightarrow \infty} \int_{-1}^1 f_N(x) \log_2 \frac{1}{f_N(x)} dx = \int_{-1}^1 f(x) \log_2 \frac{1}{f(x)} dx \\
 &\frac{1}{k} \sum_{k=1}^n R(k) \quad M(\|\delta_j - 1\|^2) = \int_0^1 (x-1)^2 e^{-x} dx \quad \lim_{N \rightarrow \infty} \int_{-1}^1 f_N(x) \log_2 \frac{1}{f_N(x)} dx = \int_{-1}^1 f(x) \log_2 \frac{1}{f(x)} dx \\
 &\det(M') = \det(M) + \det(M^*) = \det(M) \quad h(xy) = \frac{1}{2\pi} \left[ \sqrt{2} e^{-\frac{x^2}{2}} - e^{-x^2} \right] \operatorname{Im}
 \end{aligned}$$



A person is seen from behind, sitting at a wooden desk and typing on a silver laptop. The person is wearing a dark long-sleeved shirt and a watch on their left wrist. On the desk, there is also a smartphone and a tablet. The background is slightly blurred, showing other people in a room. Overlaid on the image is a large, bold, white text quote.

**The truth about data science:  
cleaning your data is 90% of the  
work. Fitting the model is easy.  
Interpreting the results is the  
other 90%**

-- Jake VanderPlas

# Key Challenges for a Data Scientist

- Data Cleaning/wrangling
- Feature Engineering
- Hyperparameter Optimization
- Insights



A dark, atmospheric background featuring a microphone on a stand and a spotlight beam. The microphone is positioned on the left side, and a bright spotlight beam cuts through the darkness from the upper right. The overall mood is mysterious and focused.

**Questions?**



**Thank you**