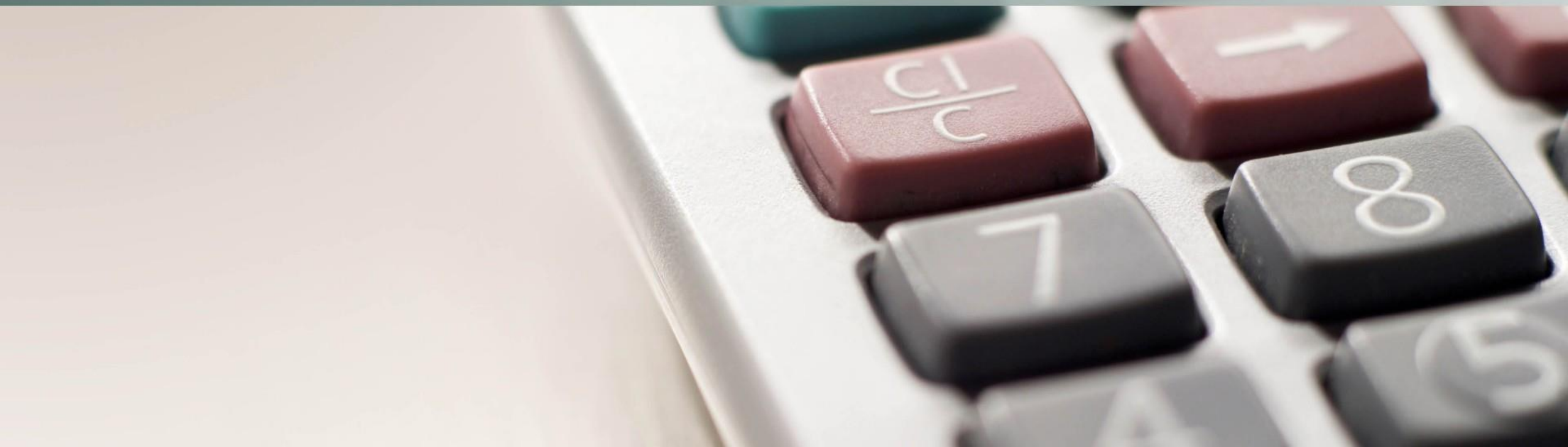


The Data Scientists: Actuary 2.0?

Doan Le FIAA & Guan Wang

AAC2018 - 18 September 2018



AT WORK

Dust Off Your Math Skills: Actuary Is Best Job of 2013

By *Lauren Weber*

Apr 22, 2013 5:00 pm ET

THE WALL STREET JOURNAL.

DATA

Data Scientist: The Sexiest Job of the 21st Century

by Thomas H. Davenport and D.J. Patil

FROM THE OCTOBER 2012 ISSUE

Harvard
Business
Review

Actuaries



What my friends think I do



What my mom thinks I do



What society thinks I do



What other actuaries think I do



What I think I do



What I actually do

Data Scientist



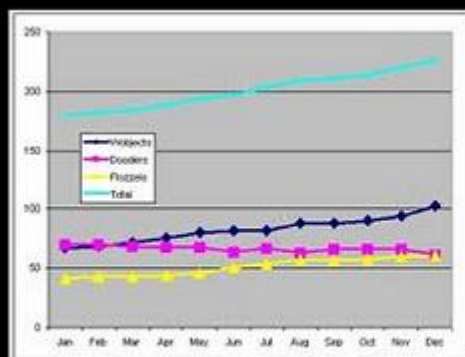
What my friends think I do



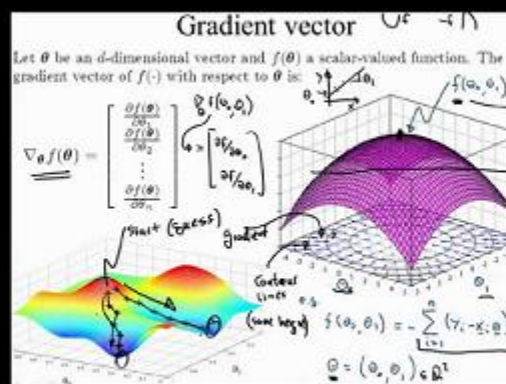
What my mom thinks I do



What society thinks I do



What my boss thinks I do

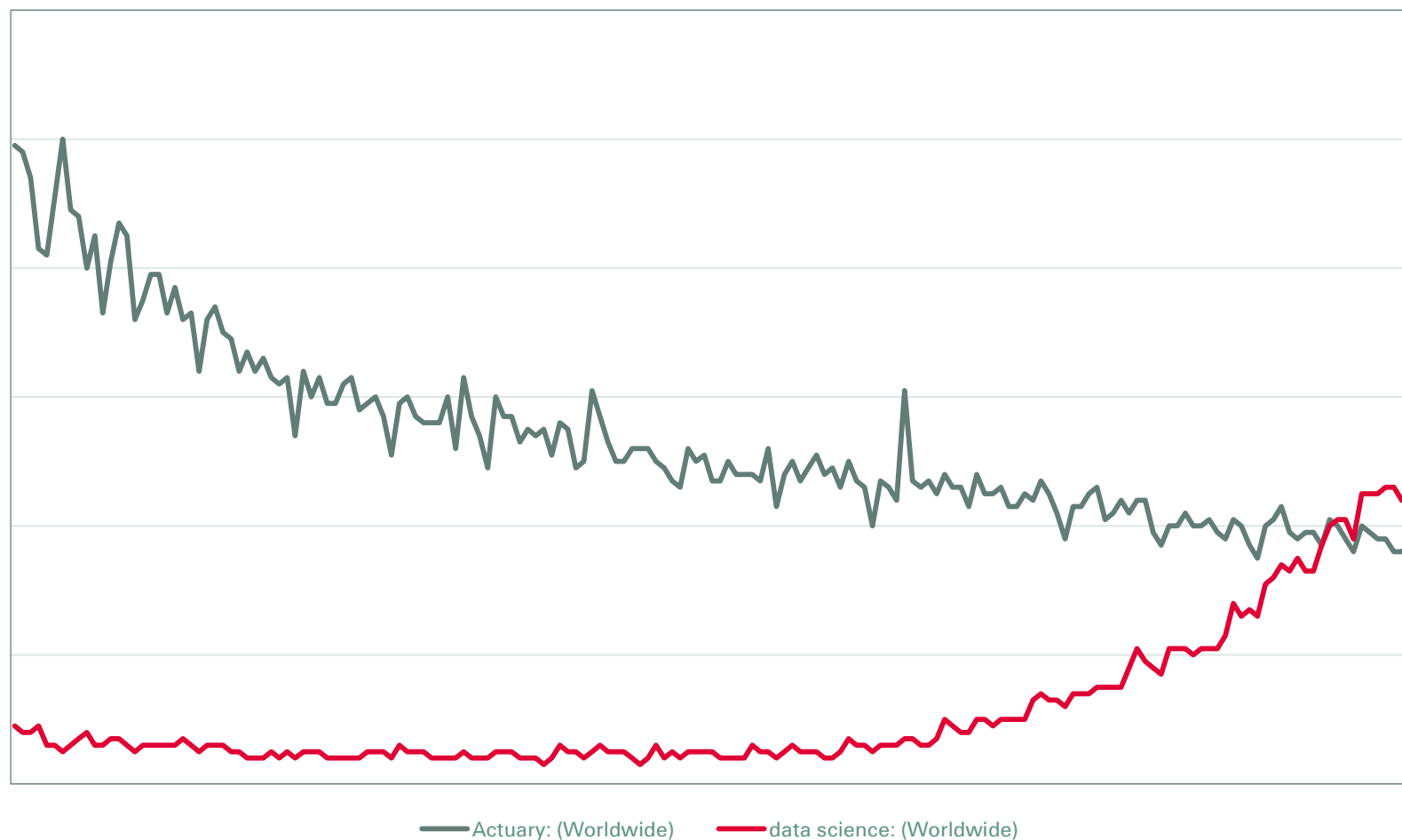


What I think I do

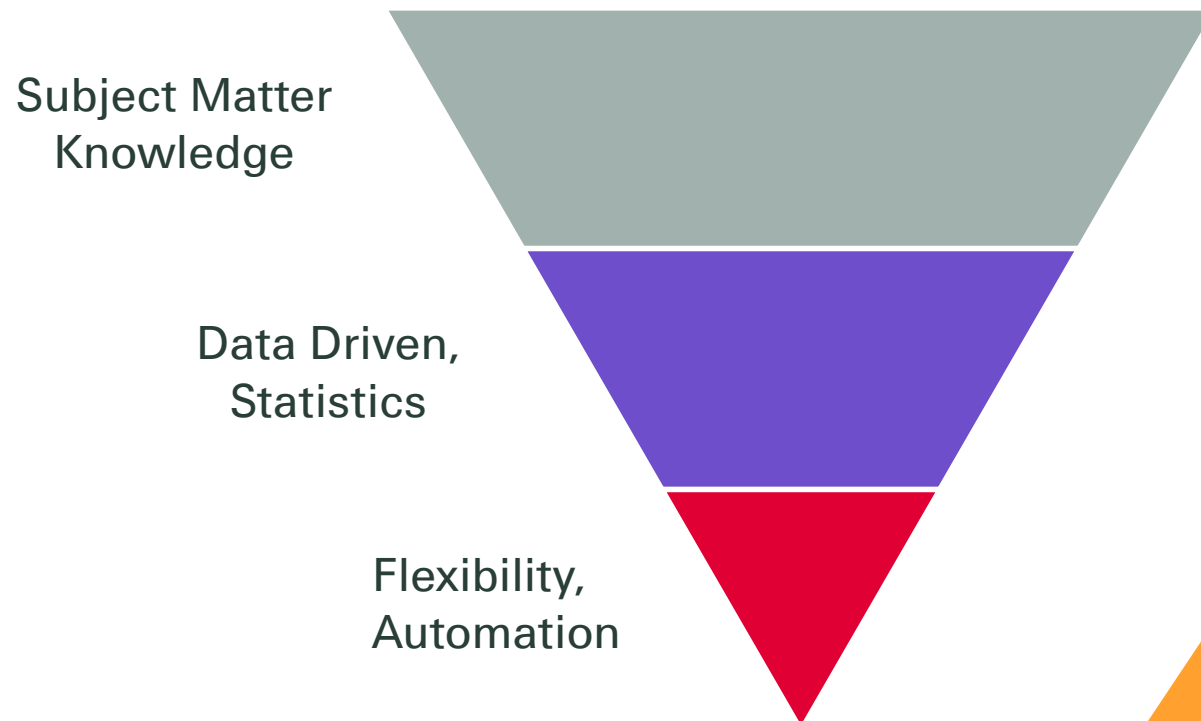


What I actually do

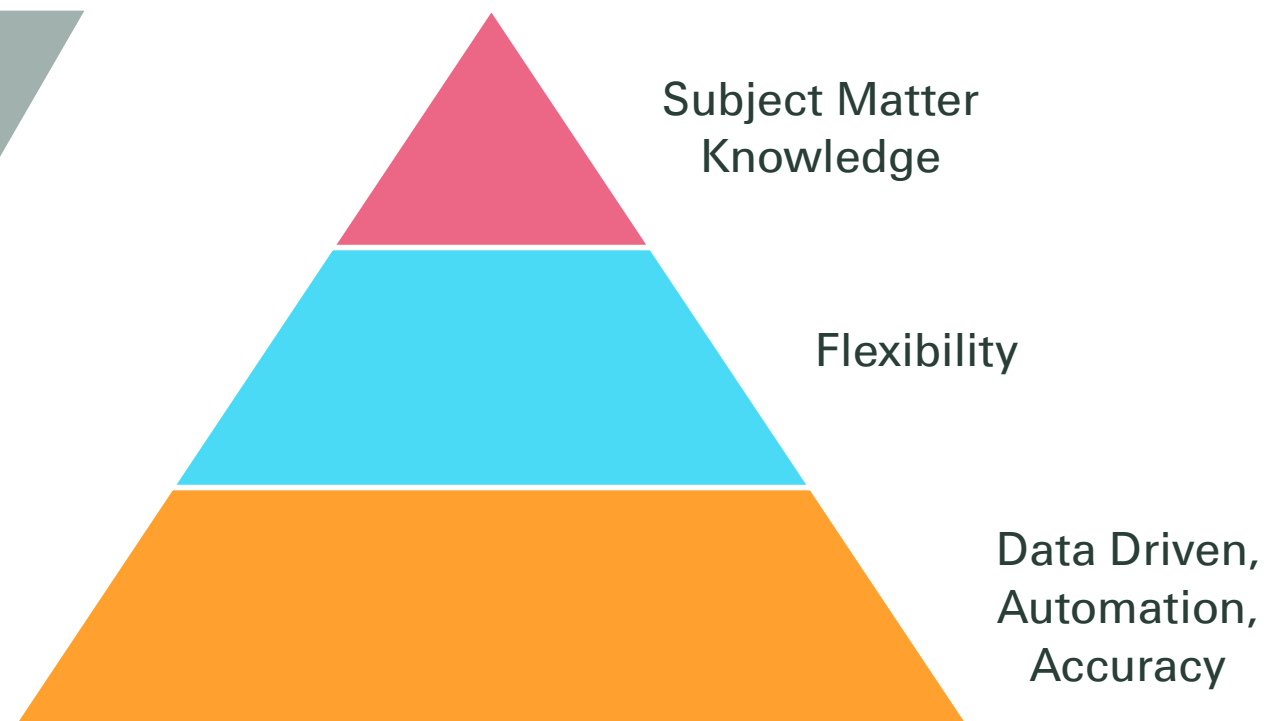
“Actuary” vs “Data Science” Google Searches (2004 – 2018)








Actuary



Data Scientist



	 Actuary	Data Scientist
 Underlying principle	<ul style="list-style-type: none"> • Pooling is a key concept • Interested in cohort experience of similar risk profiles. • Balance between data insight vs. business knowledge 	<ul style="list-style-type: none"> • Data science generally looks at individuals. • A major goal is to achieve a good balance of high accuracy and good generalization
 Application	<ul style="list-style-type: none"> • Actuarial modelling combines data analysis and subject matter expertise • Problems solved are almost always financial • Risk-oriented 	<ul style="list-style-type: none"> • Data Science modelling's raw outcome is usually a probabilistic score • Involved in wide range of problems not just financial e.g. optimizing customer experience • Opportunity-focused
 Data	<ul style="list-style-type: none"> • Actuaries predominantly focus on attributes directly applicable to problem statement • Reliant on data but has to fill in gaps with expert knowledge 	<ul style="list-style-type: none"> • Data Science modelling is exploratory and usually looks at as many attributes as available • Structured, unstructured, labelled, unlabelled • Completely dependent on data
 Modelling technique	<ul style="list-style-type: none"> • Actuarial modelling tends to rely on established methods and are more explanatory • A global explanation is important e.g. male is x% more likely to claim than female 	<ul style="list-style-type: none"> • A wider selection of modelling techniques are used. Some are extremely complicated • Selection of model depends on data • Strong command of programming and automation skills

A COMMON ACTUARIAL PROBLEM

How much premium should we charge to cover
policyholder death?

A COMMON ACTUARIAL PROBLEM

How much premium should we charge to cover policyholder death?



A COMMON ACTUARIAL PROBLEM

How much premium should we charge to cover policyholder death?





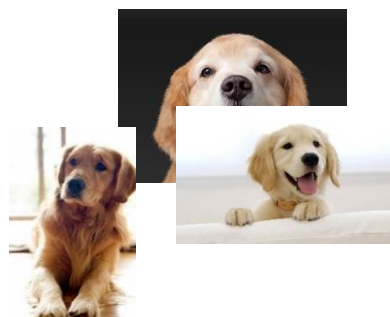
Data Sourcing: Why do we need label data?



Supervised

Learn a function by observing examples containing the input and the expected output.

- Classification
- Regression



This looks like the dogs I saw before, it should be a dog too

Unsupervised

Find underlying relations in data by observing the raw data only (without the expected output)

- Clustering
- Dimensionality reduction



I can see two types of animals. But I don't know what they are. Can someone tell me what they are?

Data Sourcing: How do we get label data?



- Most Machine Learning Applications are **Supervised** Learning

- If lucky: automatically generated

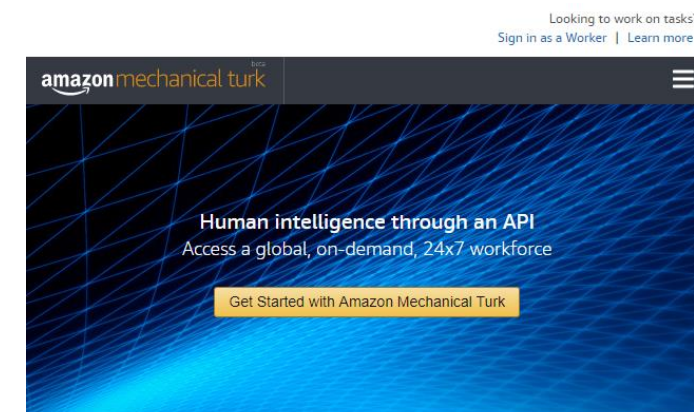
- Click Google ad on a webpage
- Add stuff in shopping cart on Amazon
- Listen to music on Spotify

- Not so lucky, but there are smart tricks

- Crawl from Internet
- Distant Supervision

- Not lucky, no tricks (usually)

- Label by hands



Amazon Mechanical Turk (MTurk) operates a marketplace for work that requires human intelligence. The MTurk web service enables companies to programmatically access this marketplace and a diverse, on-demand workforce. Developers can leverage this service to build an intelligence directly into their applications.

<https://www.mturk.com/>

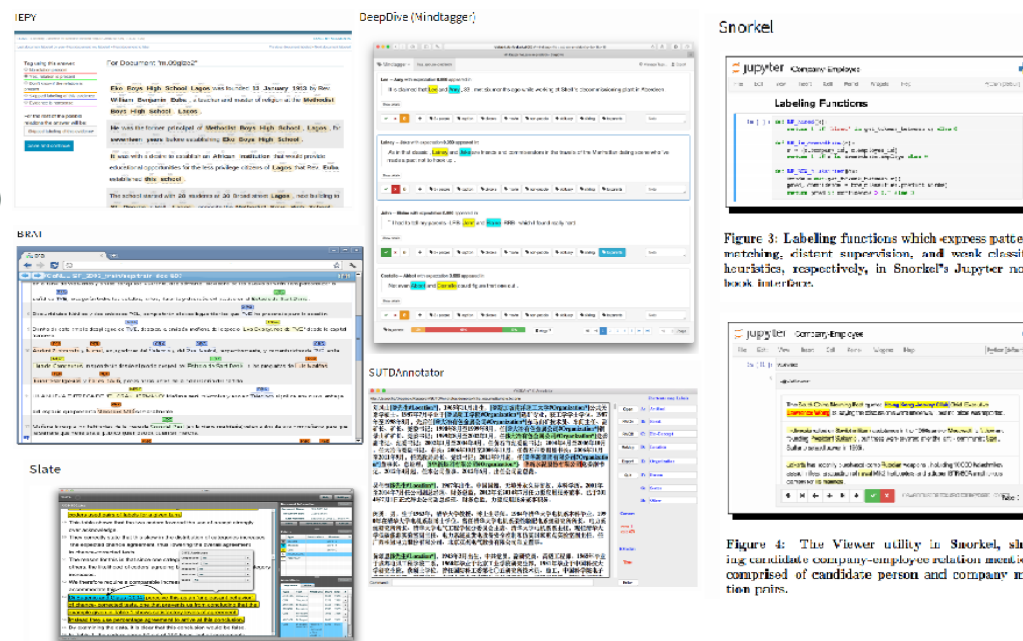
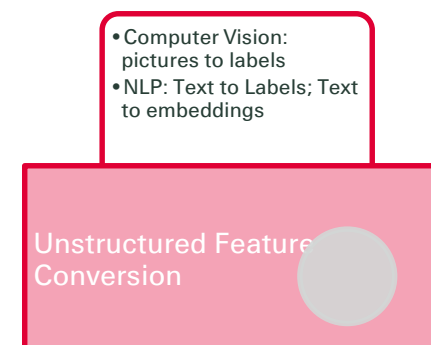
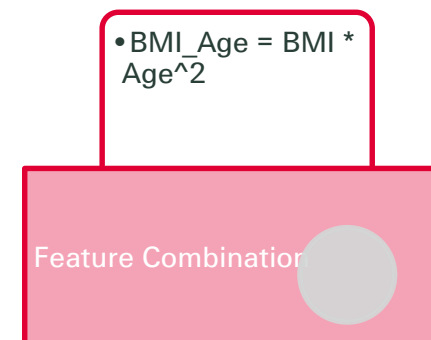
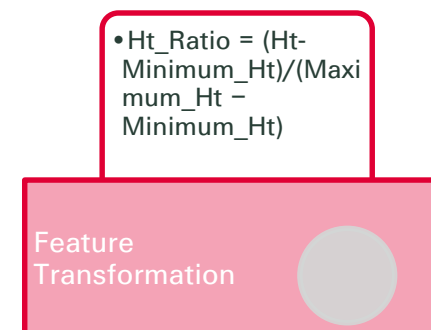
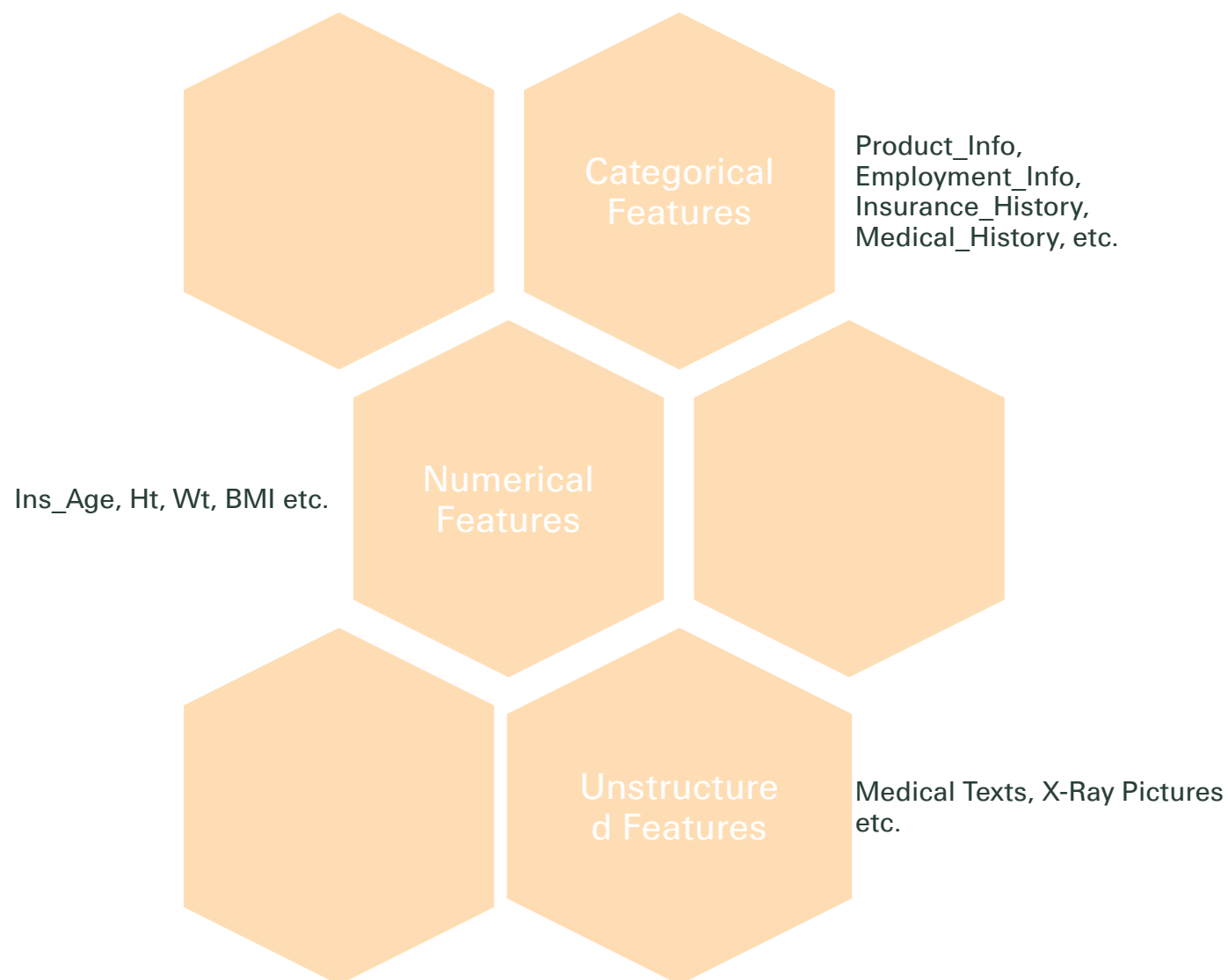


Figure 3: Labelling functions which express pattern-matching, distant supervision, and weak classifier heuristics, respectively, in Snorkel's Jupyter notebook interface.

Figure 4: The Viewer utility in Snorkel, showing candidate company-employee relation mentions, comprised of candidate person and company mention pairs.



Cleansing and Feature Engineering



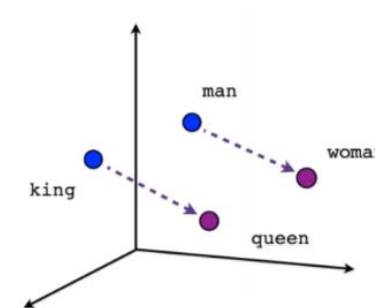
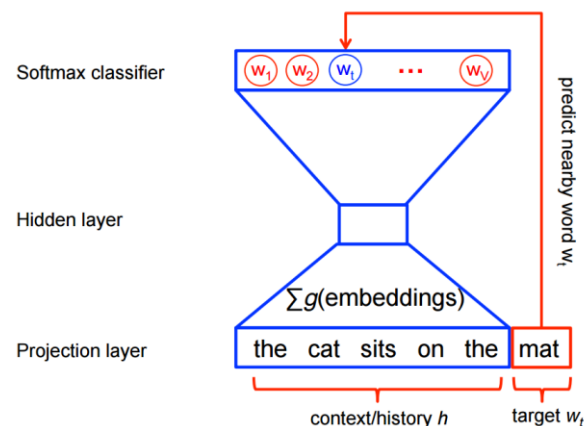
Feature Engineering: Representation Learning

General Public Release

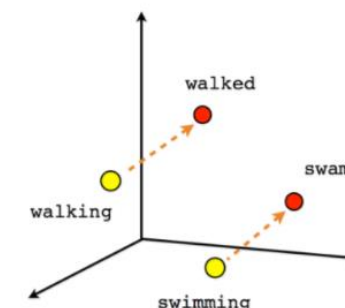


Word Embedding: From text to vectors

King - Man + Woman = Queen

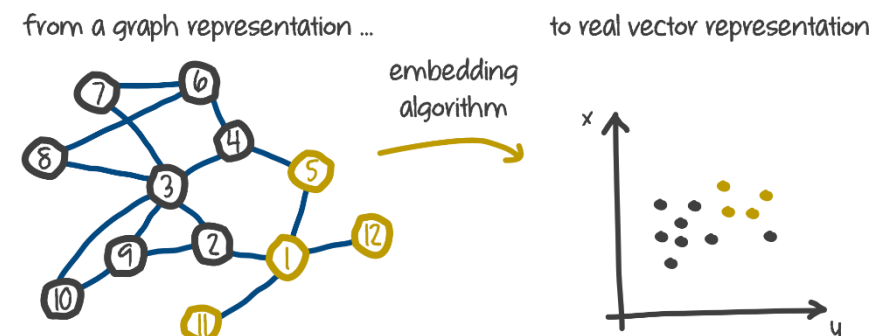


Male-Female



Verb tense

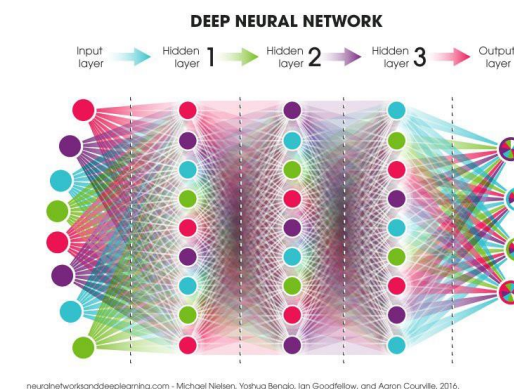
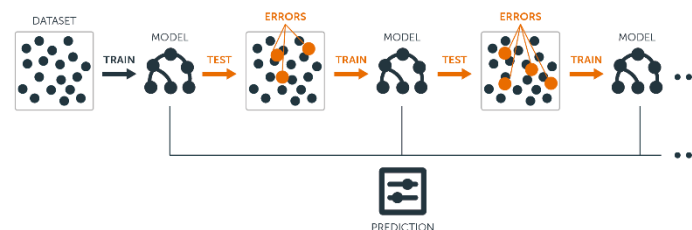
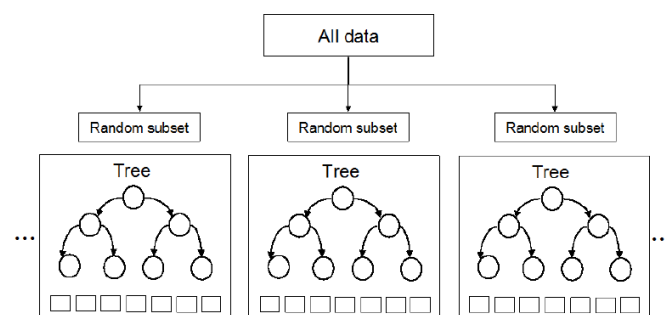
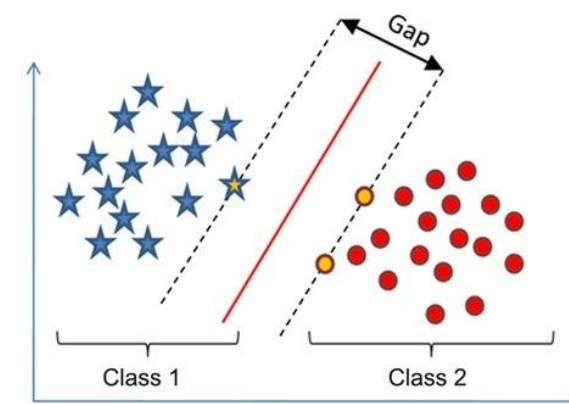
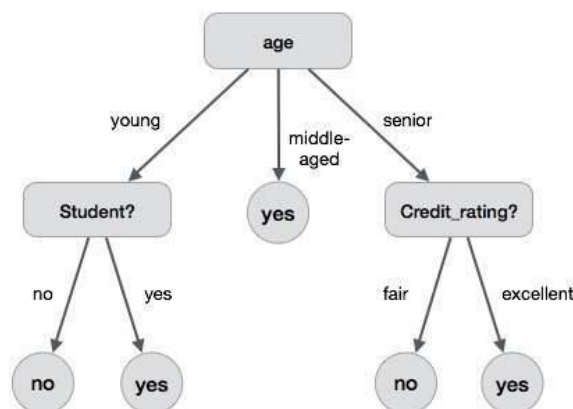
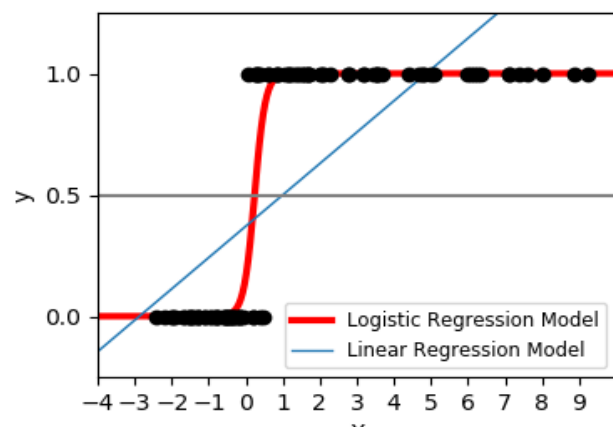
Graph Embedding: From graphs to vectors





Modelling: Model Selection

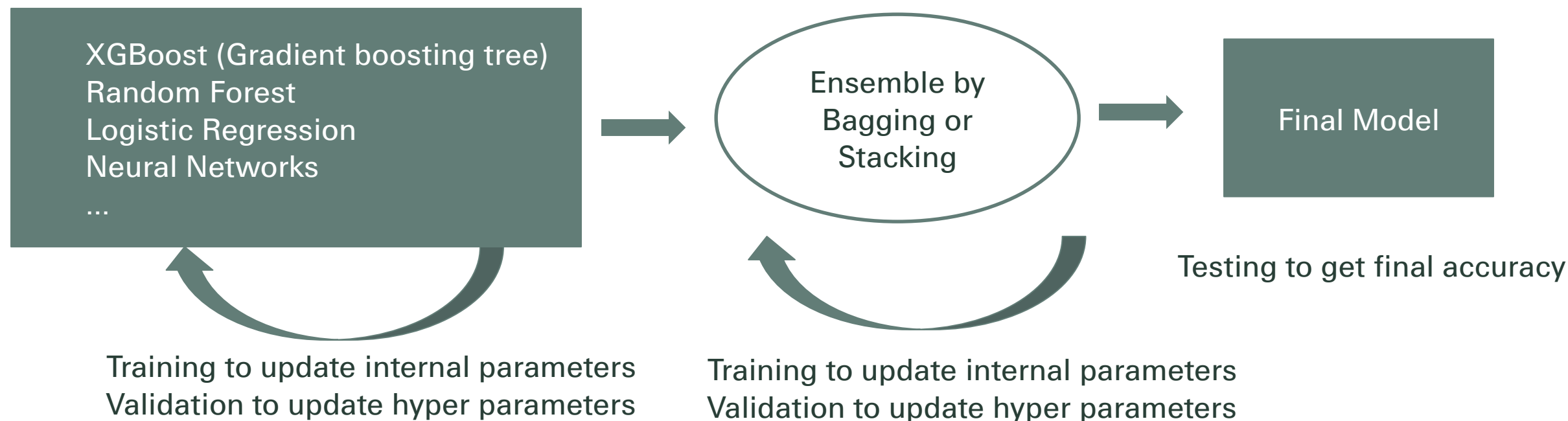
• Single Model Algorithms





Modelling: Model Ensembling

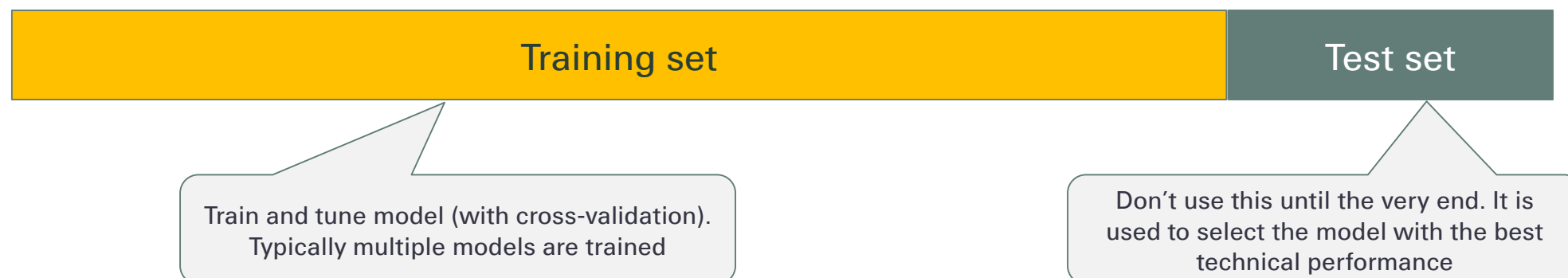
- Ensemble single models to get better results (in the cost of more complex models)





Modelling: Data Splitting

- Data Splitting



In model build, our goal is to apply machine learning techniques to learn the pattern of labelled data, and this pattern is able to generalize well to unseen data.

The data is split into two sets, training and test set.

Training set is used to train and tune model (with cross-validation). The trained model contains the pattern of labelled data. The test set acts as unseen data, which is used to select the one with the best technical performance among the multiple models trained.



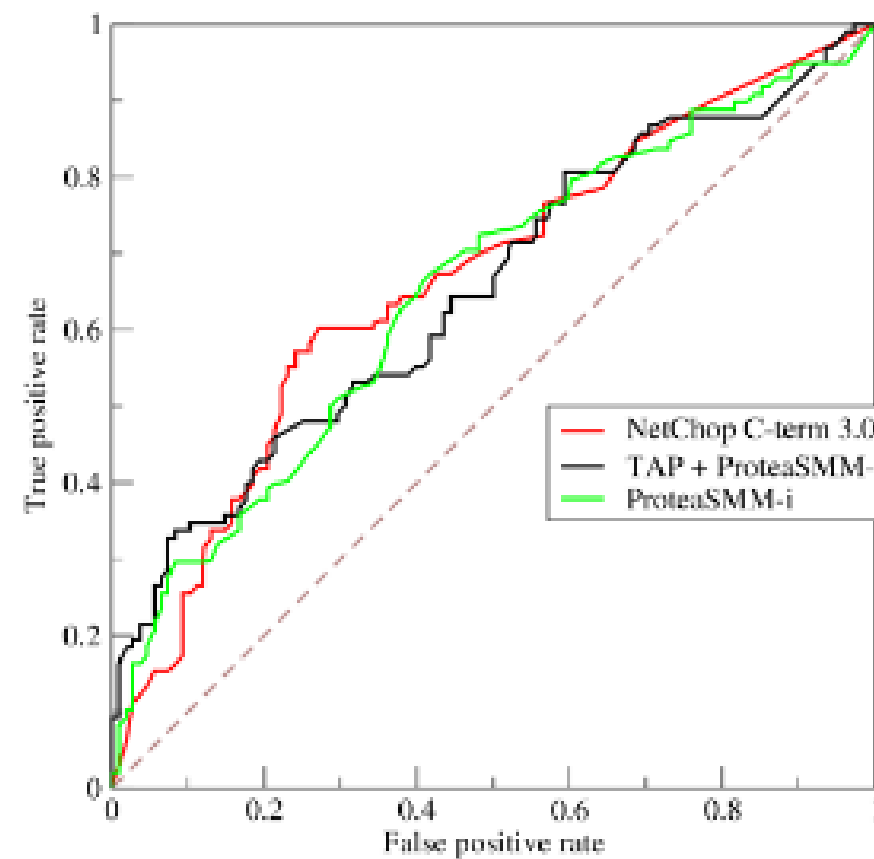


Modelling: Testing

– Confusion Matrix

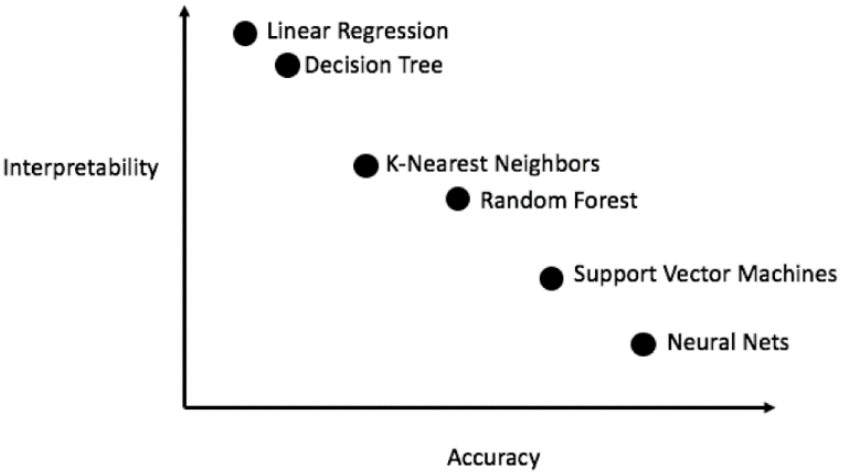
		True condition				
		<u>Total population</u>	Condition positive	Condition negative	<u>Prevalence</u> = $\frac{\Sigma \text{ Condition positive}}{\Sigma \text{ Total population}}$	<u>Accuracy</u> (ACC) = $\frac{\Sigma \text{ True positive} + \Sigma \text{ True negative}}{\Sigma \text{ Total population}}$
Predicted condition	Predicted condition positive	<u>True positive, Power</u>	<u>False positive, Type I error</u>	<u>Positive predictive value</u> (PPV), <u>Precision</u> = $\frac{\Sigma \text{ True positive}}{\Sigma \text{ Predicted condition positive}}$	<u>False discovery rate</u> (FDR) = $\frac{\Sigma \text{ False positive}}{\Sigma \text{ Predicted condition positive}}$	
	Predicted condition negative	<u>False negative, Type II error</u>	<u>True negative</u>	<u>False omission rate</u> (FOR) = $\frac{\Sigma \text{ False negative}}{\Sigma \text{ Predicted condition negative}}$	<u>Negative predictive value</u> (NPV) = $\frac{\Sigma \text{ True negative}}{\Sigma \text{ Predicted condition negative}}$	
		<u>True positive rate</u> (TPR), <u>Recall</u> , <u>Sensitivity</u> , probability of detection = $\frac{\Sigma \text{ True positive}}{\Sigma \text{ Condition positive}}$	<u>False positive rate</u> (FPR), <u>Fall-out</u> , probability of false alarm = $\frac{\Sigma \text{ False positive}}{\Sigma \text{ Condition negative}}$	<u>Positive likelihood ratio</u> (LR+) = $\frac{\text{TPR}}{\text{FPR}}$	<u>Diagnostic odds ratio</u> (DOR) = $\frac{\text{LR+}}{\text{LR-}}$	
	<u>False negative rate</u> (FNR), Miss rate = $\frac{\Sigma \text{ False negative}}{\Sigma \text{ Condition positive}}$	<u>True negative rate</u> (TNR), <u>Specificity</u> (SPC) = $\frac{\Sigma \text{ True negative}}{\Sigma \text{ Condition negative}}$	<u>Negative likelihood ratio</u> (LR-) = $\frac{\text{FNR}}{\text{TNR}}$	<u>F₁ score</u> = $\frac{2}{\frac{1}{\text{Recall}} + \frac{1}{\text{Precision}}}$		

– Receiver Operating Characteristic (ROC) and Area Under Curve (AUC)



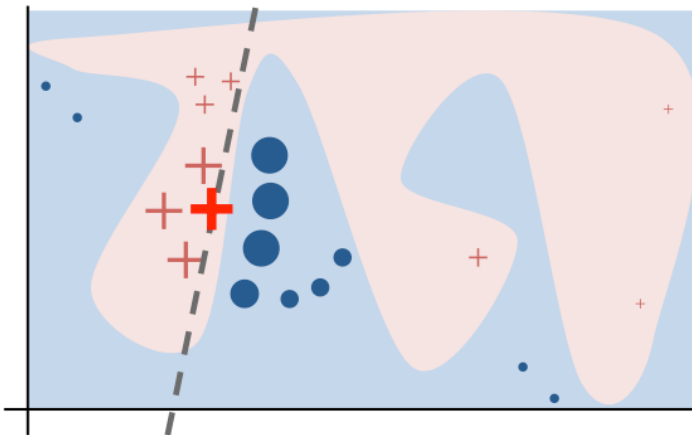
Modelling: Interpretation

General Public Release



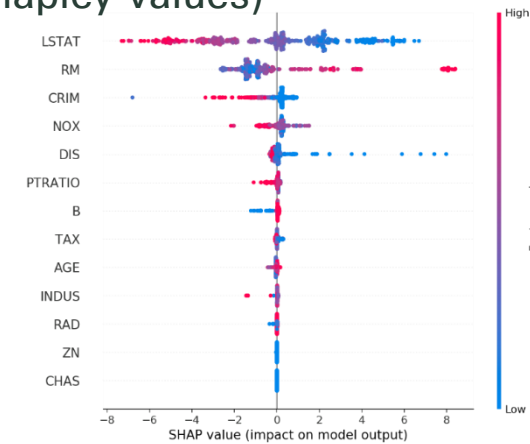
<https://medium.com/ansaro-blog/interpreting-machine-learning-models-1234d735d6c9>

- Many ways to interpret the “Black Box”
 - Fit local prediction into a linear model (LIME)



<https://github.com/marcotcr/lime>

- Switch on and off feature combinations to get feature importance (Shapley Values)



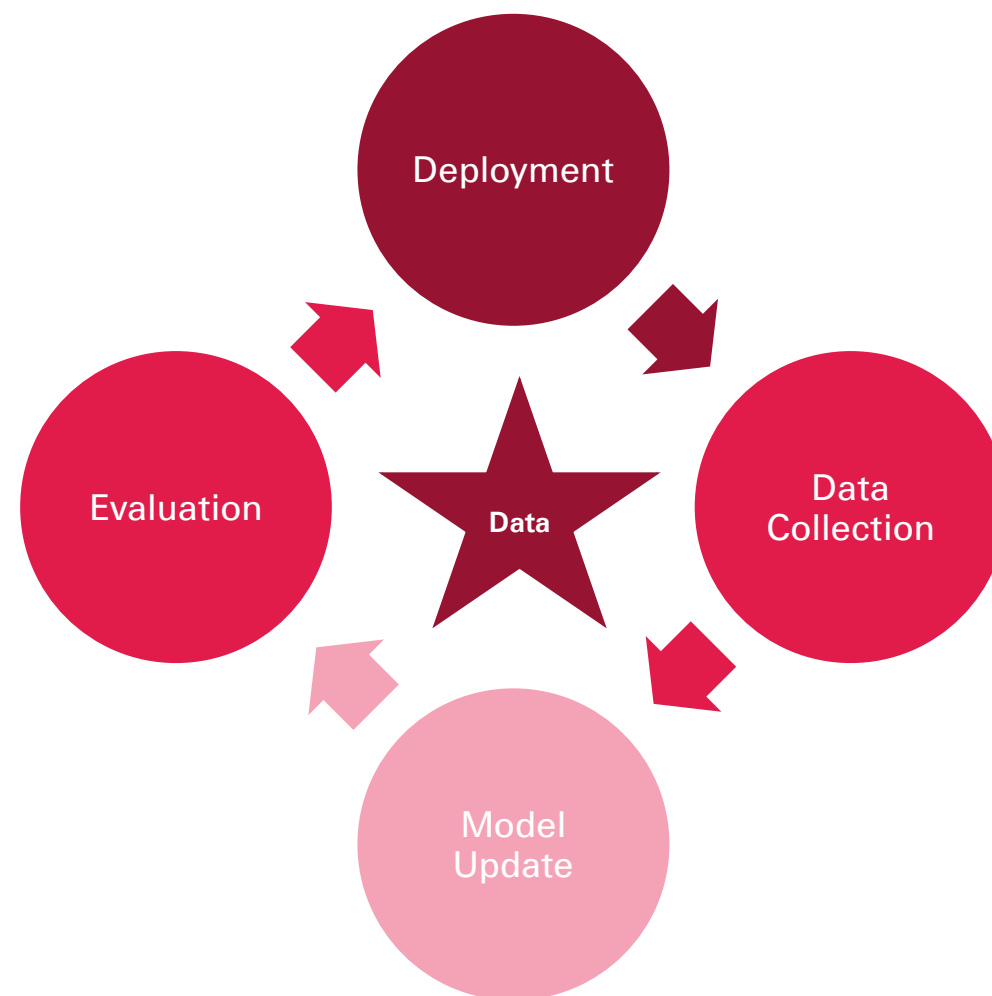
<https://github.com/slundberg/shap>





Model Deployment

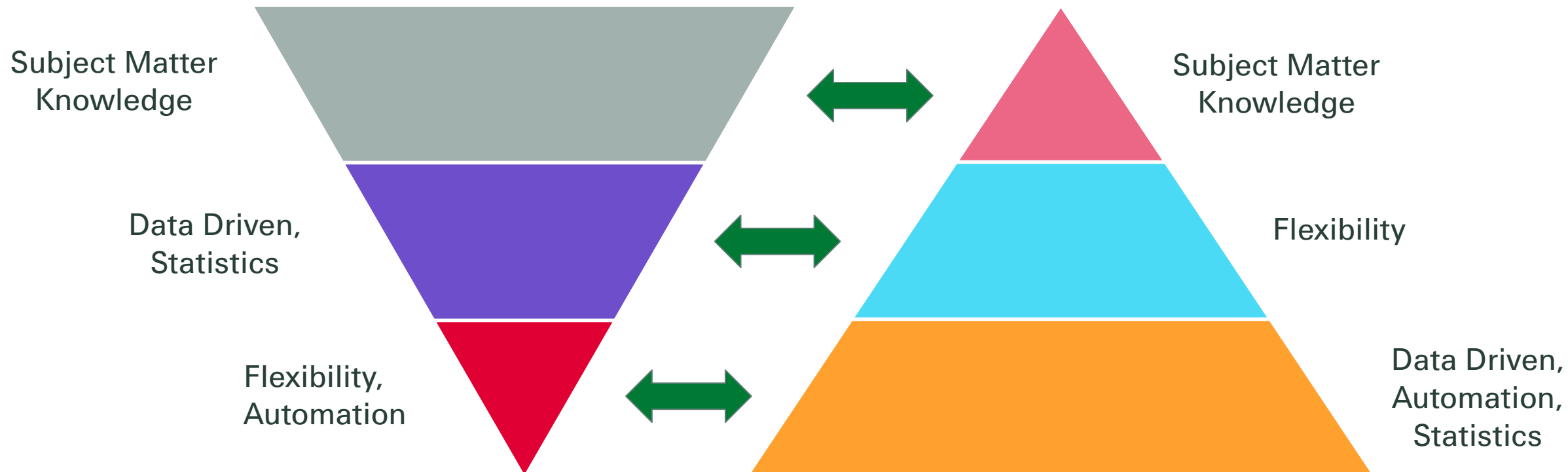
- Model access through API
- Easy integration into business process
- Constant Model Update
- A/B Testing
- Automation



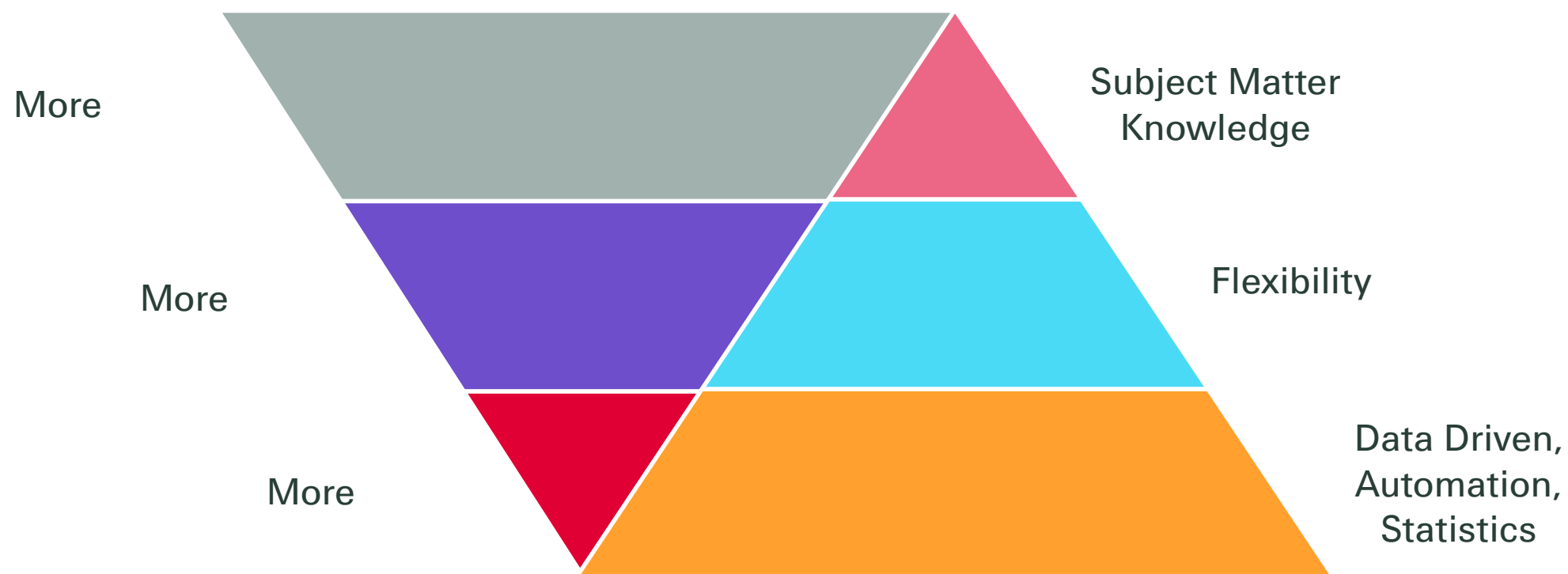
Summary

Actuary

Data Scientist



Actuaries together with Data Scientists



How would you build a system to provide personalized insurance product recommendations to potential customers, just like advertisement on Google, books on Amazon, movies on Netflix and music on Spotify?

Questions?



Legal notice

©2018 Swiss Re. All rights reserved. You are not permitted to create any modifications or derivative works of this presentation or to use it for commercial or other public purposes without the prior written permission of Swiss Re.

The information and opinions contained in the presentation are provided as at the date of the presentation and are subject to change without notice. Although the information used was taken from reliable sources, Swiss Re does not accept any responsibility for the accuracy or comprehensiveness of the details given. All liability for the accuracy and completeness thereof or for any damage or loss resulting from the use of the information contained in this presentation is expressly excluded. Under no circumstances shall Swiss Re or its Group companies be liable for any financial or consequential loss relating to this presentation.