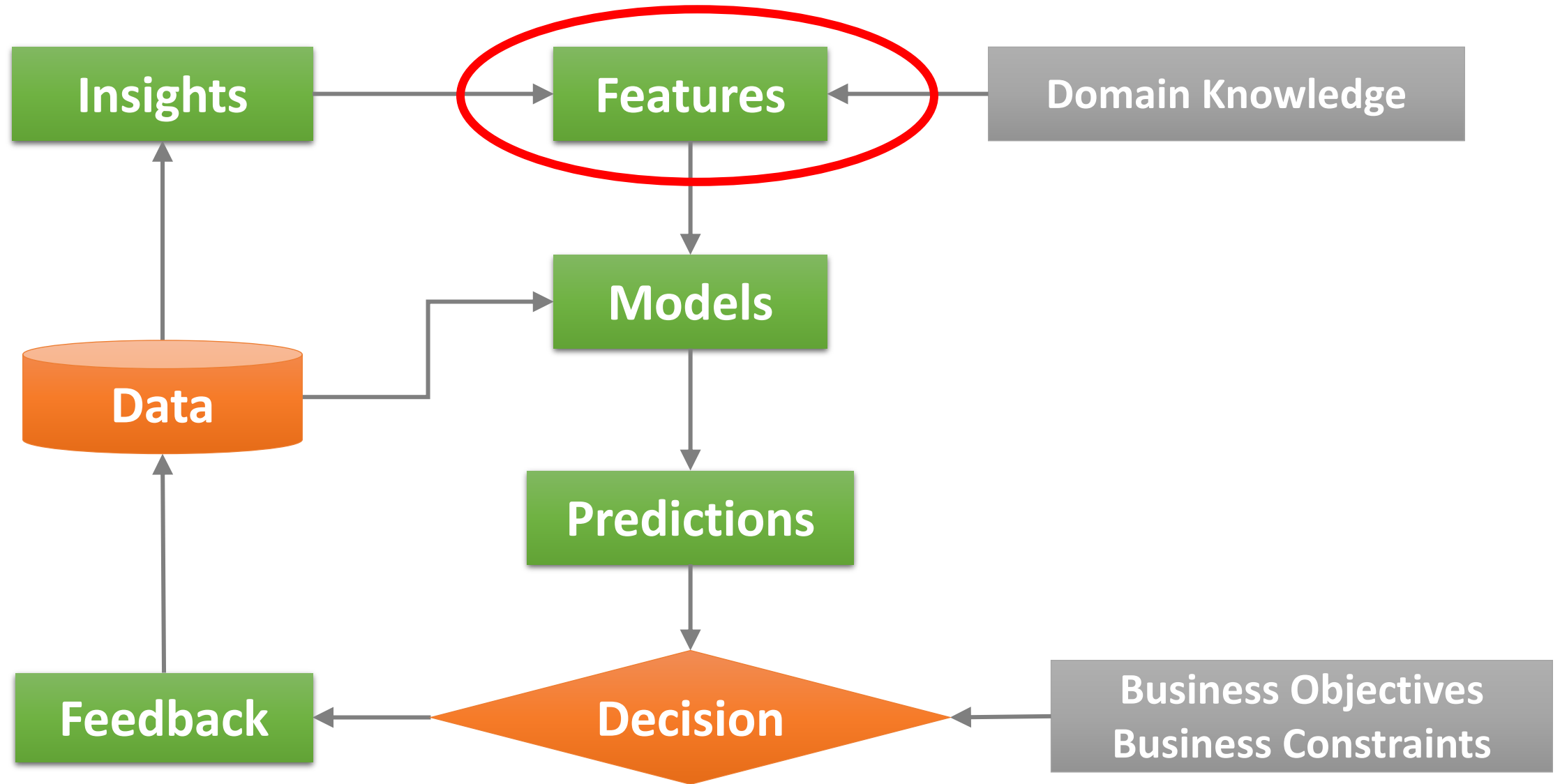


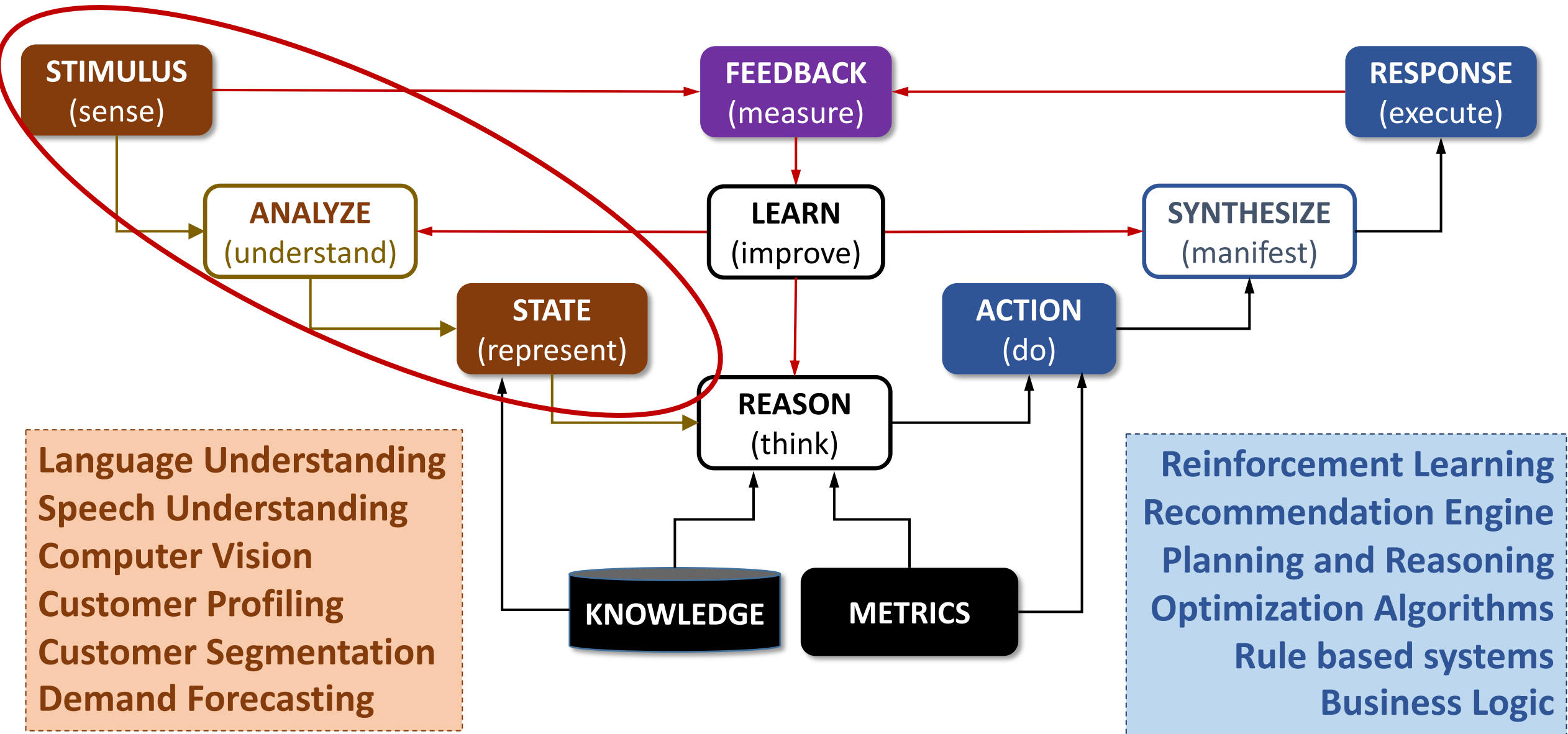
Machine Learning **FEATURE ENGINEERING**

Shailesh Kumar

Feature Engineering



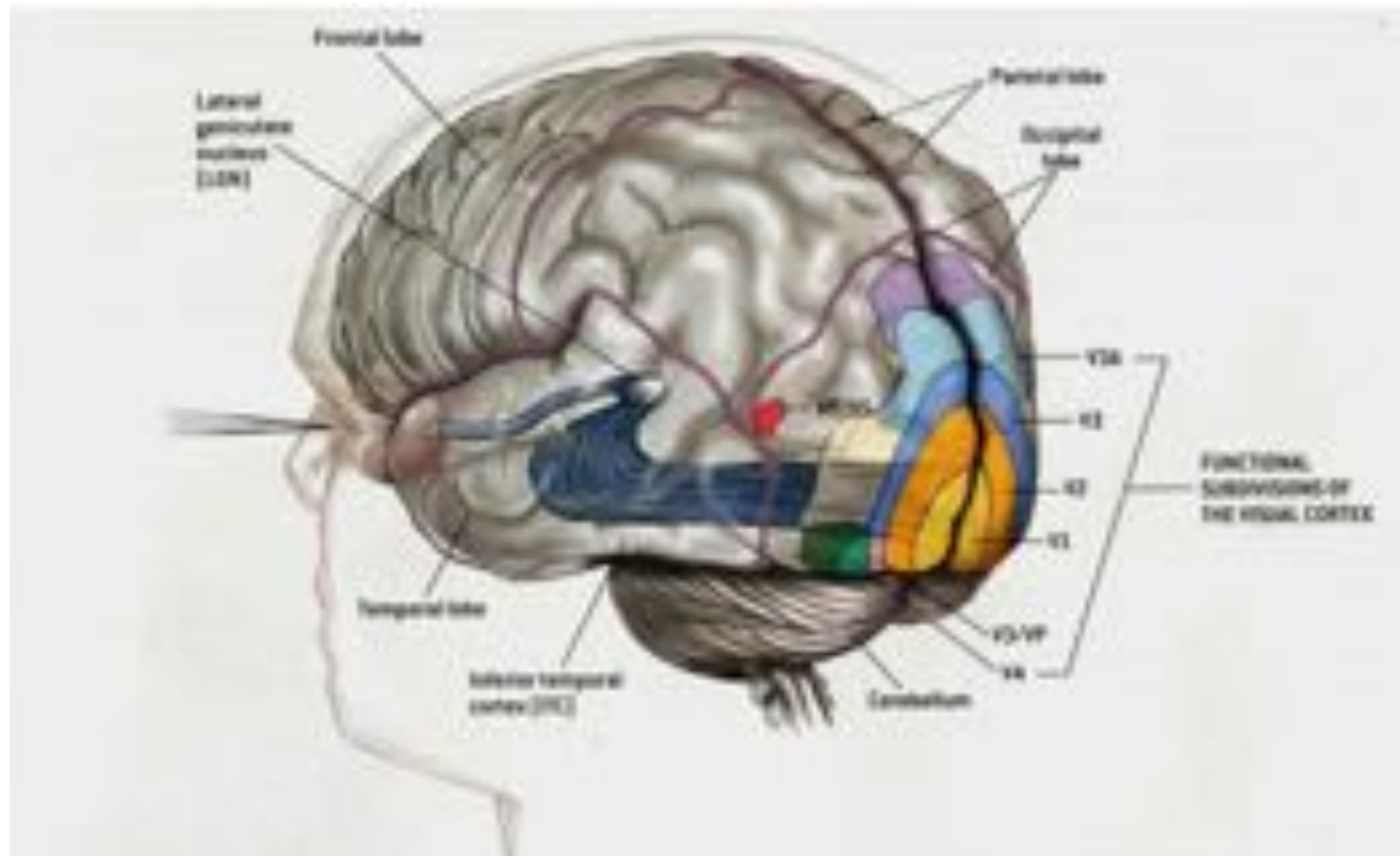
Feature Engineering



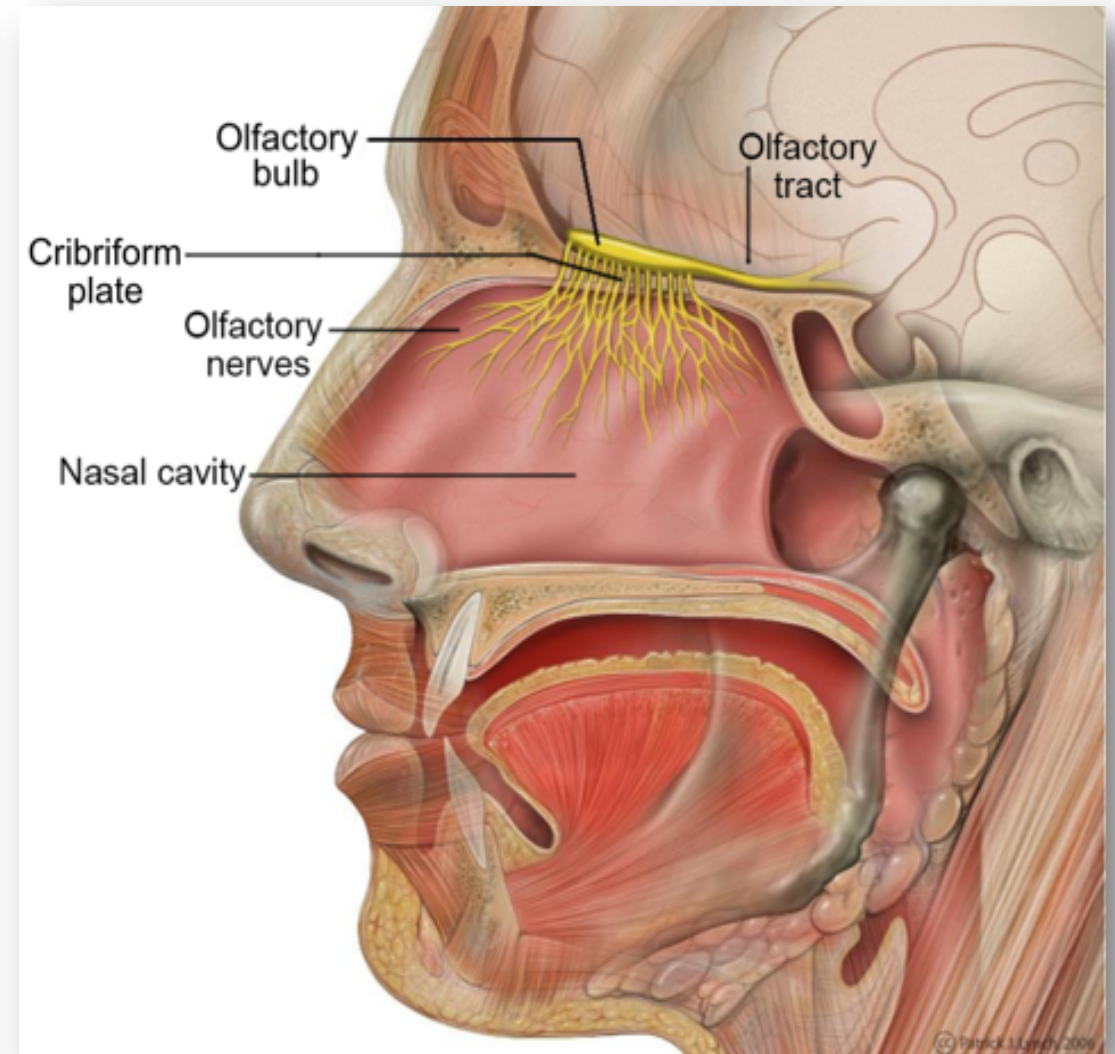
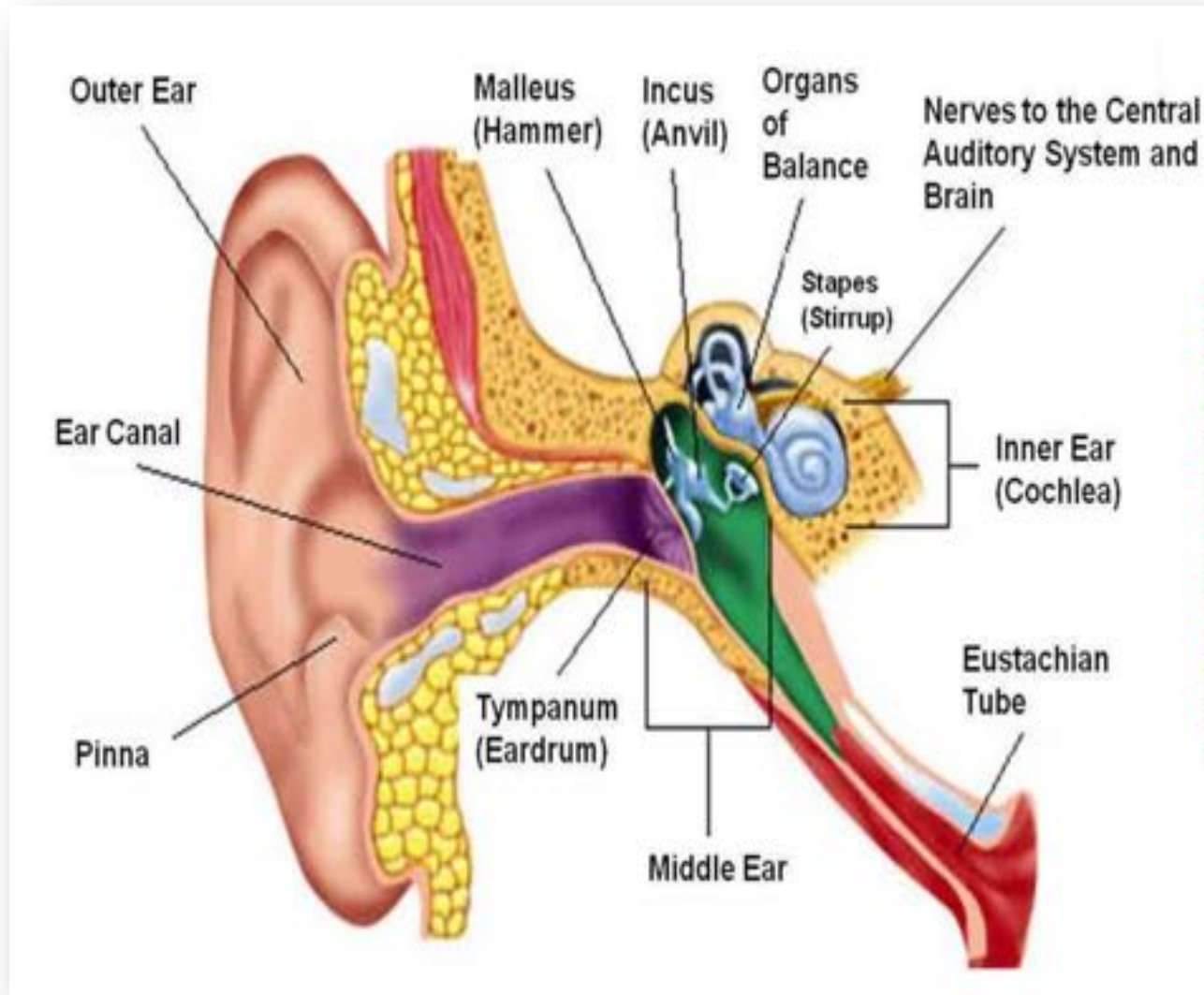
Visual Cortex | The Most Complex Feature Engineering

40% of brain activity is focused on visual processing alone!

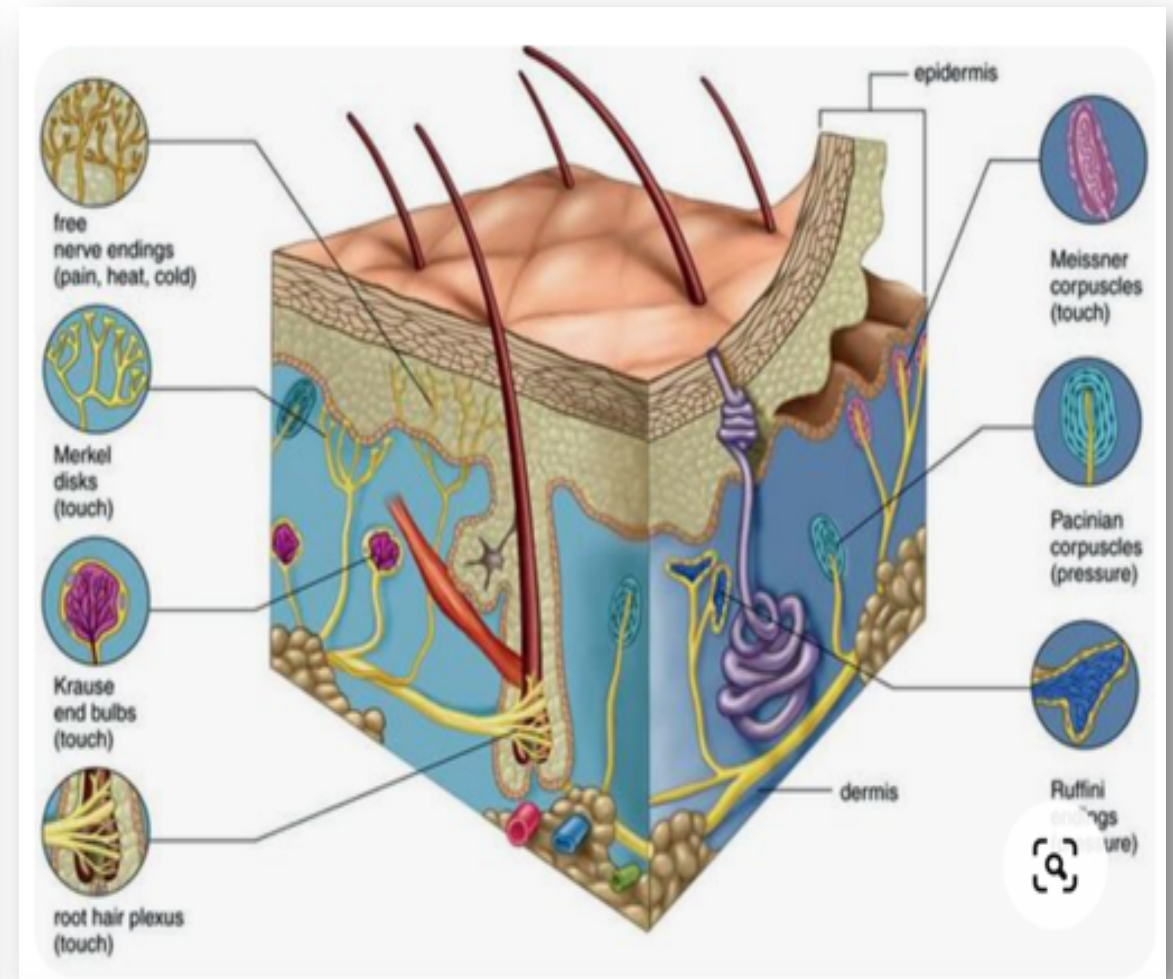
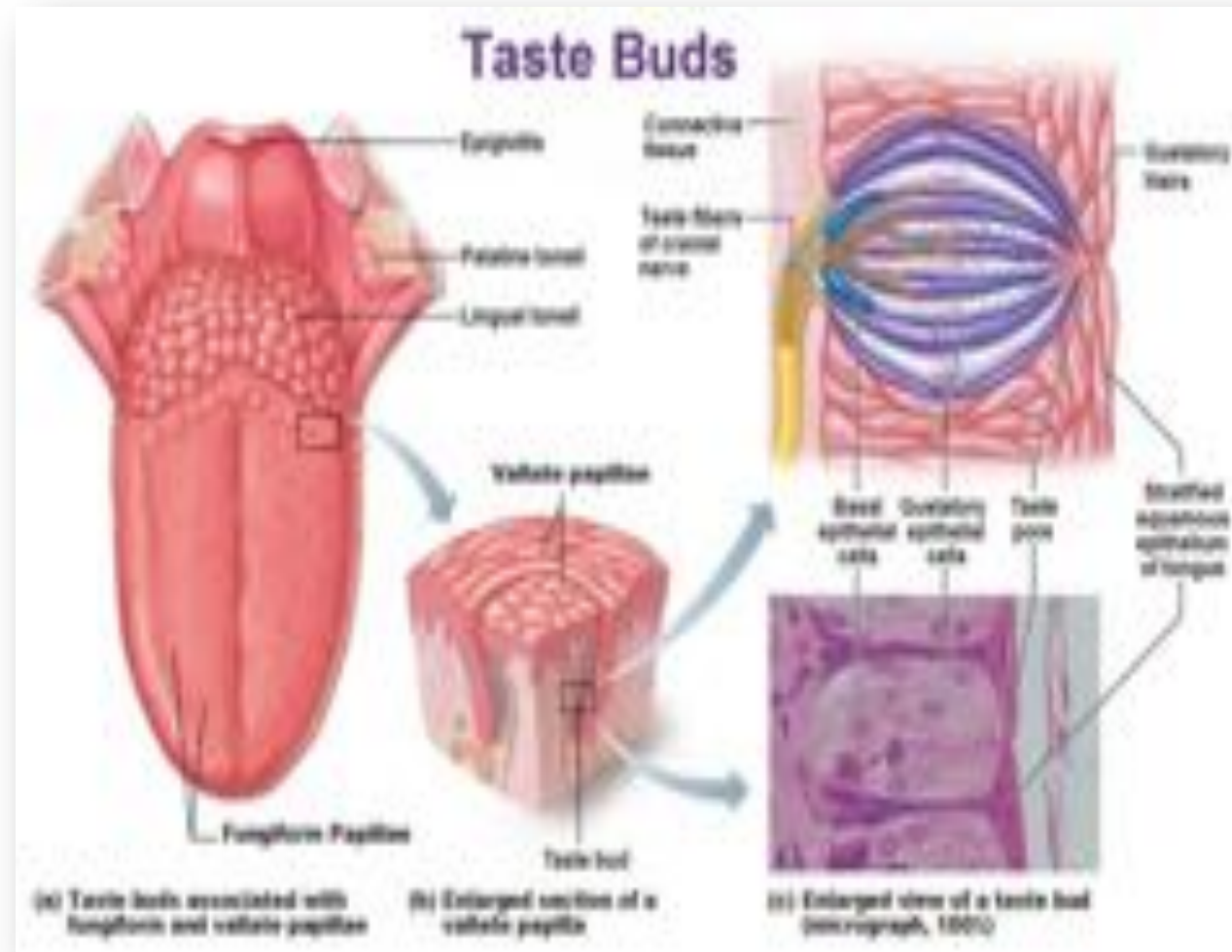
Modern Deep Learning for Vision is inspired in part by the Visual Cortex



Our Sensory organs have in-built feature engineering



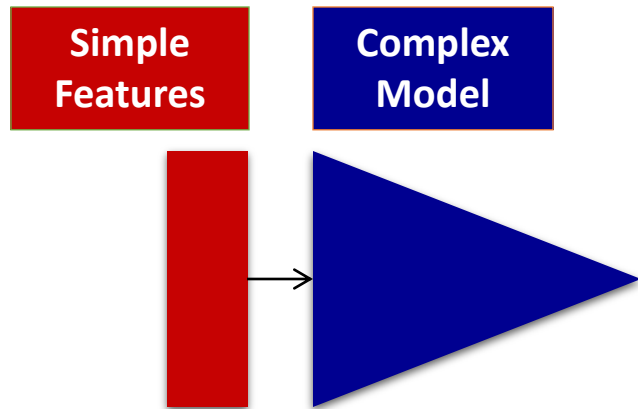
Our Sensory organs have in-built feature engineering



Two Mindsets to Modelling

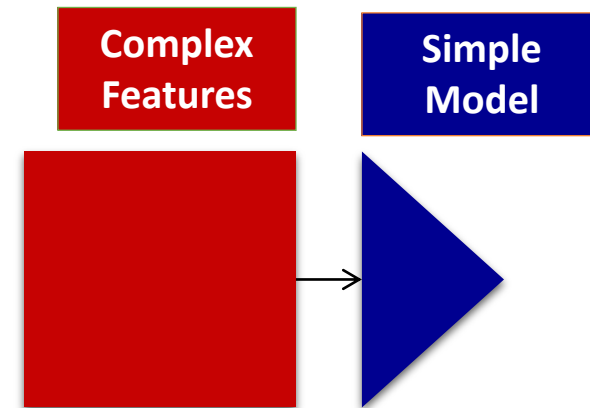
Model-Centric

- Throw all features in!
- Have enough data
- Build Complex models



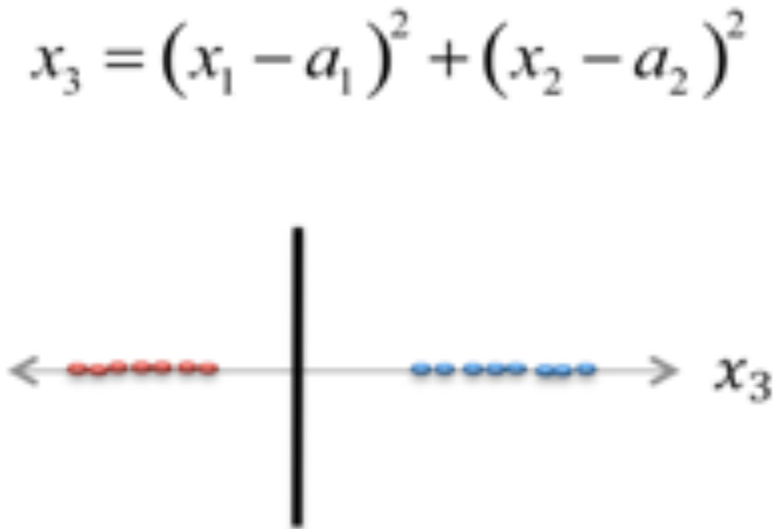
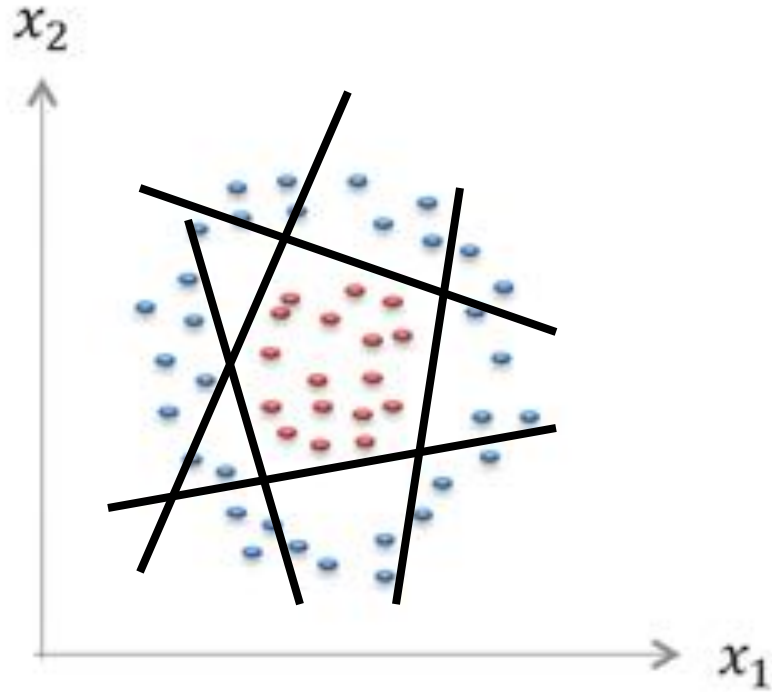
Feature-centric

- Carefully craft features
- Use Domain Knowledge
- Build Simpler Models



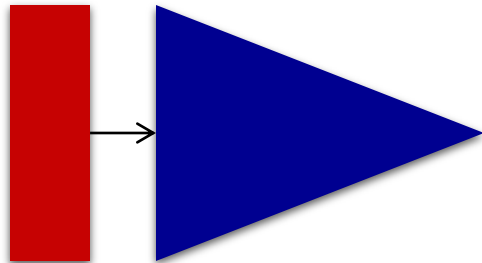
The Law of Conservation of Complexity

Two Mindsets of Modelling



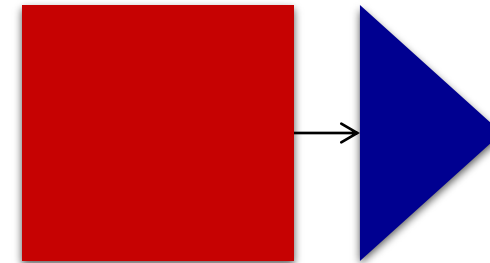
Simple
Features

Complex
Model



Complex
Features

Simple
Model



Feature Engineering using **Domain Knowledge**

Raw Input

- **Time** of **current** trans.
- **Place** of **current** trans.
- **Time** of **prev.** trans.
- **Place** of **prev.** trans.



Derived Feature - 1

- ▶ **Distance**(Prev→Current)
- ▶ **TimeLag**(Prev→Current)

Derived Feature - 2

- ▶ **Velocity**(Prev→Current)

$$\begin{aligned} &\text{Velocity}(\text{Prev} \rightarrow \text{Current}) \\ &= \frac{\text{Distance}(\text{Prev} \rightarrow \text{Current})}{\text{TimeLag}(\text{Prev} \rightarrow \text{Current})} \end{aligned}$$

Model the **Variance** that **Matters**

- $Match(\mathbf{field} \mid \mathbf{query})$ = how relevant is a field to a query

$$Match(\mathbf{field} \mid \mathbf{query}) = \sum_{token \in \mathbf{query}} Weight(token) \times Match(\mathbf{field} \mid token)$$

- Does **query length** matter to overall relevance?

$$Match(\mathbf{field} \mid \mathbf{query}) = \frac{\sum_{token \in \mathbf{query}} Weight(token) \times Match(\mathbf{field} \mid token)}{\sum_{token \in \mathbf{query}} Weight(token)}$$

- What about **field length**?

Model the **Variance** that **Matters**

Raw Feature

Normalized Feature

TotalCardBalance

$$\frac{TotalCardBalance}{TotalCreditLimit}$$

TotalCardPayment

$$\frac{TotalCardPayment}{TotalCardBalance}$$

TotalDebt

$$\frac{TotalDebt}{AnnualIncome}$$

$$\frac{\log(TotalDebt)}{\log(AnnualIncome)}$$

Model Deviations from Expected

$$\text{TotalSales}(\text{Context}) \rightarrow \log \left(\frac{\text{TotalSales}(\text{Context})}{\text{ExpectedSales}(\text{Context})} \right)$$

$$\text{CTR}(\text{query}, \text{url}, \text{position}) \rightarrow \log \left(\frac{\text{CTR}(\text{query}, \text{url}, \text{position})}{\text{ExpectedCTR}(\text{position})} \right)$$

“Bugs” in Feature Engineering

- Observation: Model is **Unexpectedly Complex**
- Hypothesis: It is **Compensating for some “bug”**

field = *the quick brown fox jumped over a lazy dog*

*TermFrequency(query|**field**)*

*TermFrequency(quick|**field**) = 1*

*TermFrequency(brown|**field**) = 1*

*TermFrequency(dog|**field**) = 1*

*TermFrequency(cat|**field**) = 0*

*FirstOccurence(query|**field**)*

*FirstOccurence(quick|**field**) = 1*

*FirstOccurence(brown|**field**) = 2*

*FirstOccurence(dog|**field**) = 8*

*FirstOccurence(cat|**field**) = 0*

Is there anything wrong here?

Careful with those “Defaults”

field = *the quick brown fox jumped over a lazy dog*

$FirstOccurrence(quick|\mathbf{field}) = 1$

$FirstOccurrence(brown|\mathbf{field}) = 2$

$FirstOccurrence(dog|\mathbf{field}) = 8$

$FirstOccurrence(cat|\mathbf{field}) = 0$

- If query term **present** in field → LOWER is BETTER
- If query term **absent** in field → **No Match = Best Match**
- What is the correct DEFAULT? How about **-1**?
- If query term **absent** in field → **field_length + K**