# CA685 - Practicum Project

MCM - Masters in Computing (Major in Data Analytics)

## *"Interpretable Machine Learning for the Structure Odour Relationship"*

Project by:          Denis Kealy

Supervised by:          Dr. Martin Crane

# TABLE OF CONTENTS

# PROJECT INTRODUCTION

A quick word about the scope, background and goals of the project

"

*Think multidisciplinary!*

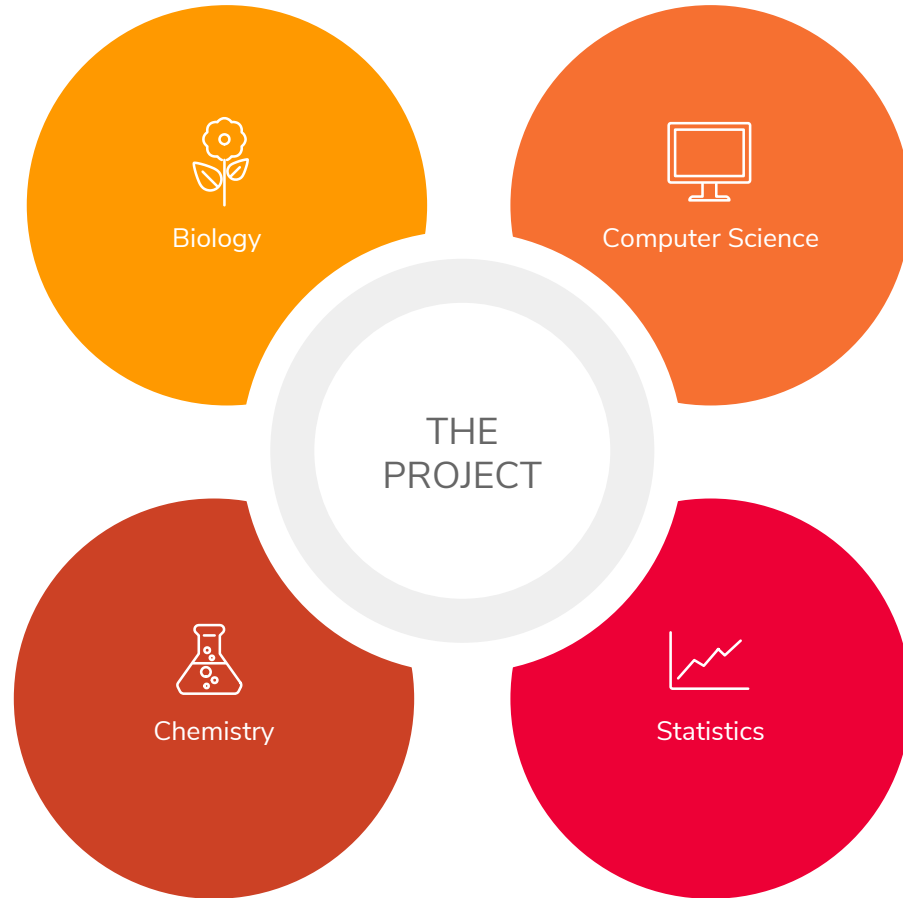Problems by definition, cross many academic disciplines.

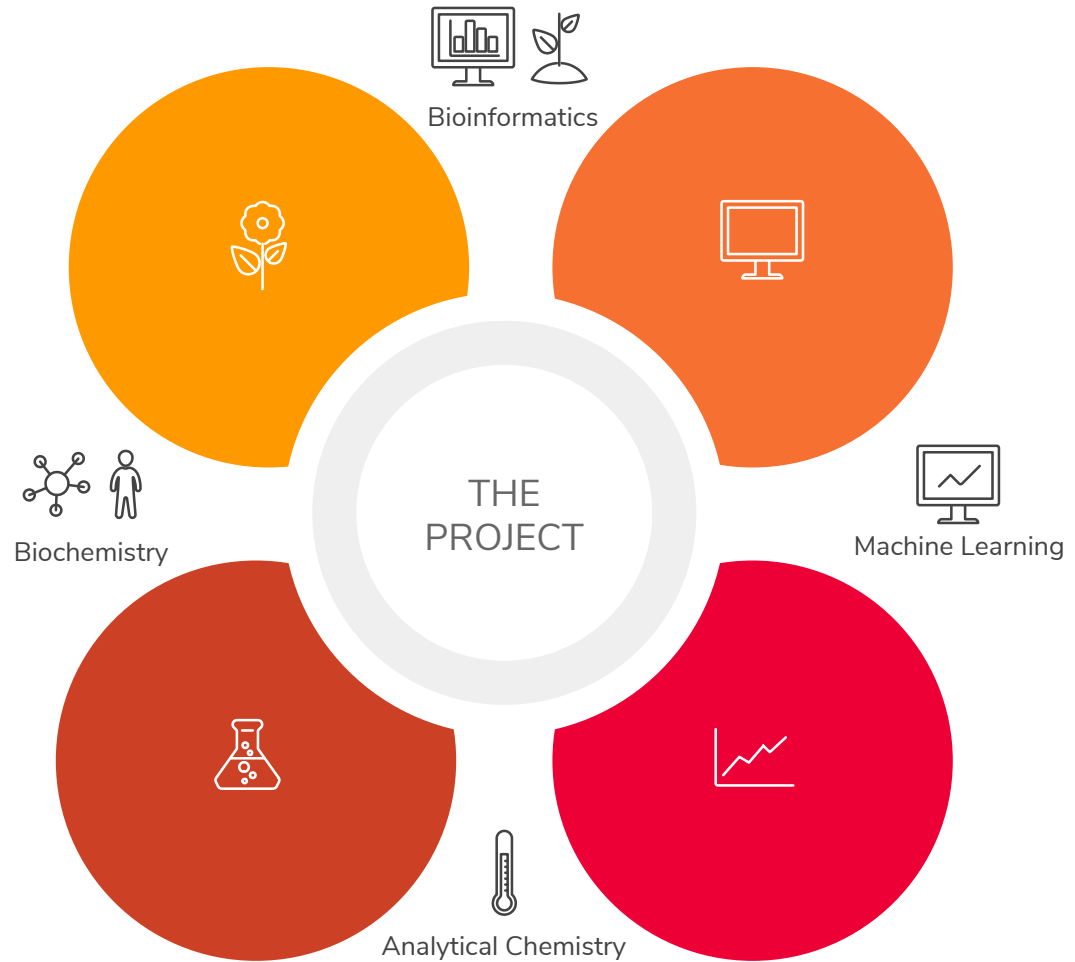- LUCAS REMERSWAAL

# DOMAIN

- Biology
- Chemistry

# APPROACH

- Statistics
- Computer Science



Biology

Computer Science

THE PROJECT

Chemistry

Statistics

# STRUCTURE ODOUR PROBLEM

... At the intersection of many disciplines

Bioinformatics

Biochemistry

Machine Learning

Analytical Chemistry

THE PROJECT

# So we are on the same page!

- Definitions

- Aims of the Project

- Previous Work

# DEFINING A FEW TERMS

## Machine Learning

- Interpretable
- Explainable / XAI
- Deep Learning
- Hyperparameter
- XNN
- Completeness & Accuracy
- AI, IA & AIA

## Our Sense of Smell

- Olfaction
- Neuron
- Receptor
- GPCR
- Limbic system
- Neo-cortex
- Swipe card model

## Chemical Structure

- SOR vs SAR
- Odorant
- Ligand
- Ligand-Receptor
- Binding
- Chemical Descriptors
- Chirality

# PROJECT AIMS

The goals of the project & the contribution to knowledge.

# PROJECT AIMS

## Investigate Smell
Prior interest & promising new results

## Generalised Methods
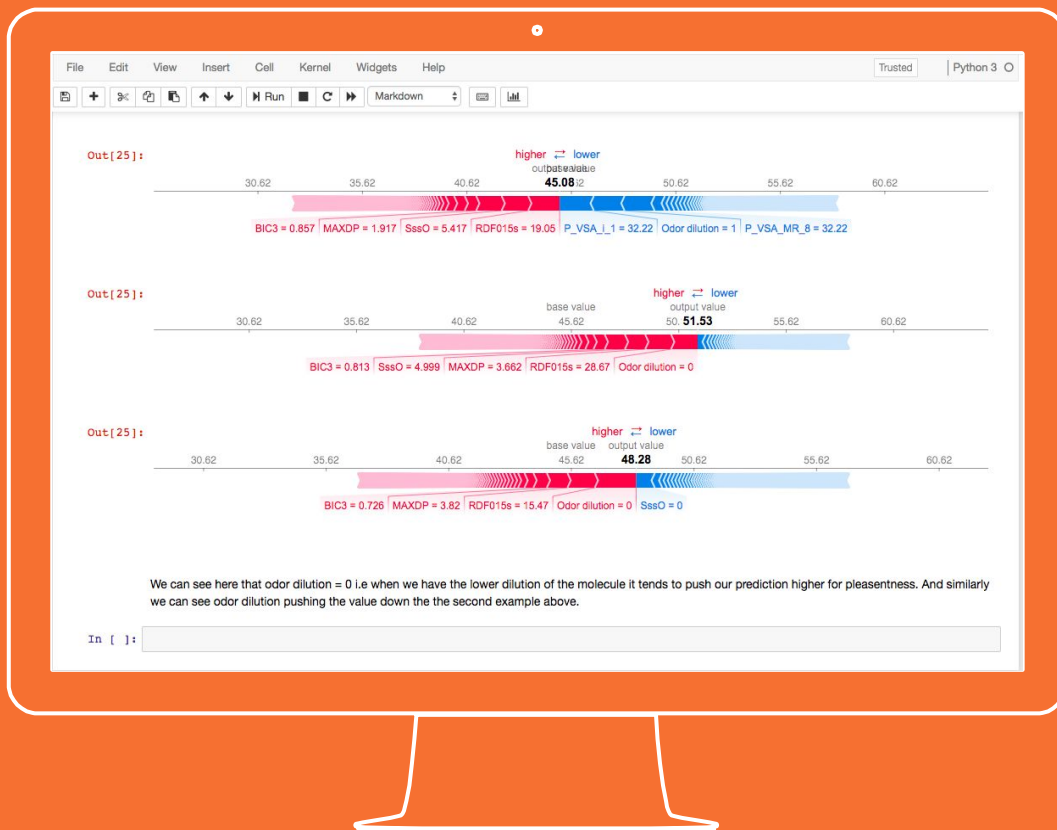A solution that is not specific to this problem

## Learn
To investigate new techniques and gain
Experience with existing ones

# FINAL GOAL

Demonstrate a generalised methodology for investigating this problem and similar problems.
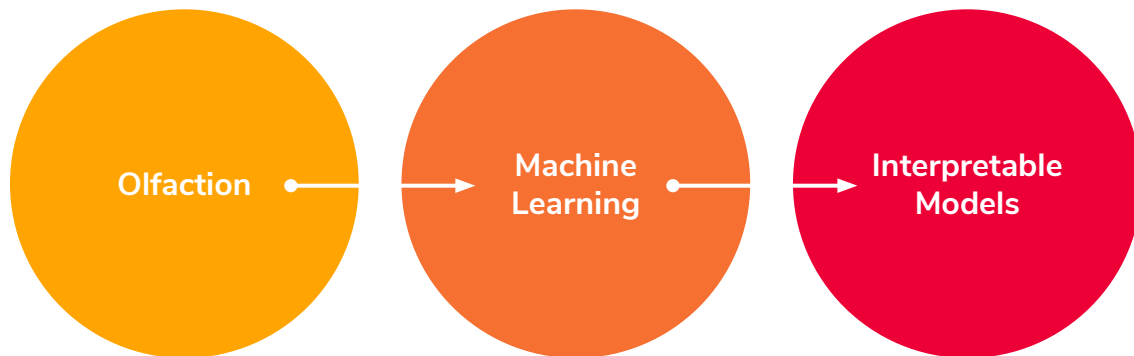
# PREVIOUS WORK

If I have seen further than others, it is by standing upon the **shoulders of giants**.
- Isaac Newton

# PREVIOUS WORK

**Domain Knowledge**

- Olfactory System
- SOR Studies
- Geonomics
- Proteomics
- New imaging technologies...

**Previous Methods**

- Odour Networks
- DREAM Machine Learning Tasks
- Neural Networks
  - E-nose
  - Fly Brain

**Novel Interpretations**

- LIME
- SHAP
- AIA

**Olfaction** → **Machine Learning** → **Interpretable Models**

Future

Previous

# BACKGROUND

Some background about the problem, the data, and the available techniques
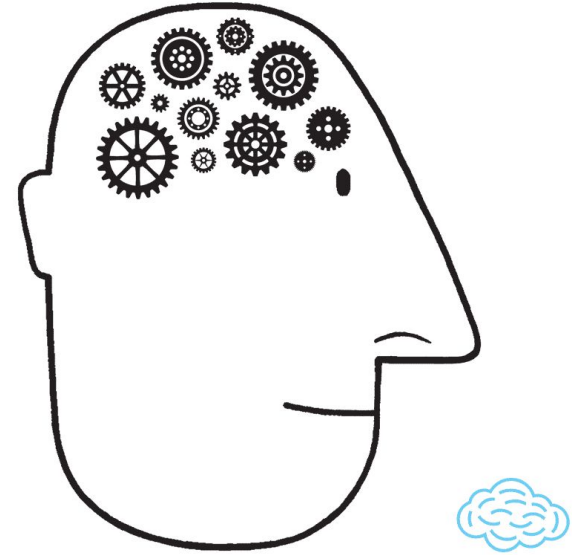
# OUR SENSE OF SMELL

INTRINSICALLY **COMPLEX**

FULL OF **COMPLEX** INTRICACIES

**INFLUENCES** OUR CONSCIOUSNESS

**INFLUENCES** OTHER SENSES

PREDICTIVE OF **COGNITIVE DECLINE**

INSPIRING **ARTIFICIAL OLFACTION**

POTENTIAL **BENEFITS** FOR:

- HEALTH
- FRAGRANCE **&** FLAVOUR
- DIGITAL SCENT
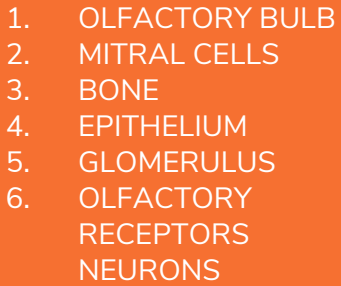- VIRTUAL REALITY

# SMELL FACTS

## What affects our sense of smell?

- Retro vs Ortho-nasal

- Mixing Compounds

- Culture & Language

- Familiarity

- Our Genes

- Time

# HUMAN OLFACTORY SYSTEM

1. OLFACTORY BULB
2. MITRAL CELLS
3. BONE
4. EPITHELIUM
5. GLOMERULUS
6. OLFACTORY RECEPTORS NEURONS

# STRUCTURE - ODOUR RELATIONSHIP

Can we predict:

- what a single molecule will smell of?
- what a SET of molecules will smell of?

**Quantum Smell**

- Structure doesn't tell the whole story
- Quanum properties seem to affect our sense of smell
- E.g. deuterated odorants & Drosophila Melanogaster

**Swipe Card Model**

- Structure plays a role, certainly, but other factors are at play
- Both chemical & physical properties of molecule are considered
- Combinatorial Encoding of Receptor Responses
- Cannot account for Chirality

# STRUCTURE - ODOUR RELATIONSHIP (SOR)

# VS

# STRUCTURE - ACTIVITY RELATIONSHIP (SAR)

## SOR

- Psychological Data

- Odour Panels
- Odour Descriptors
- Similarity Ratings
- Measure of conscious perception

- Many public databases
- Recent comprehensive dataset
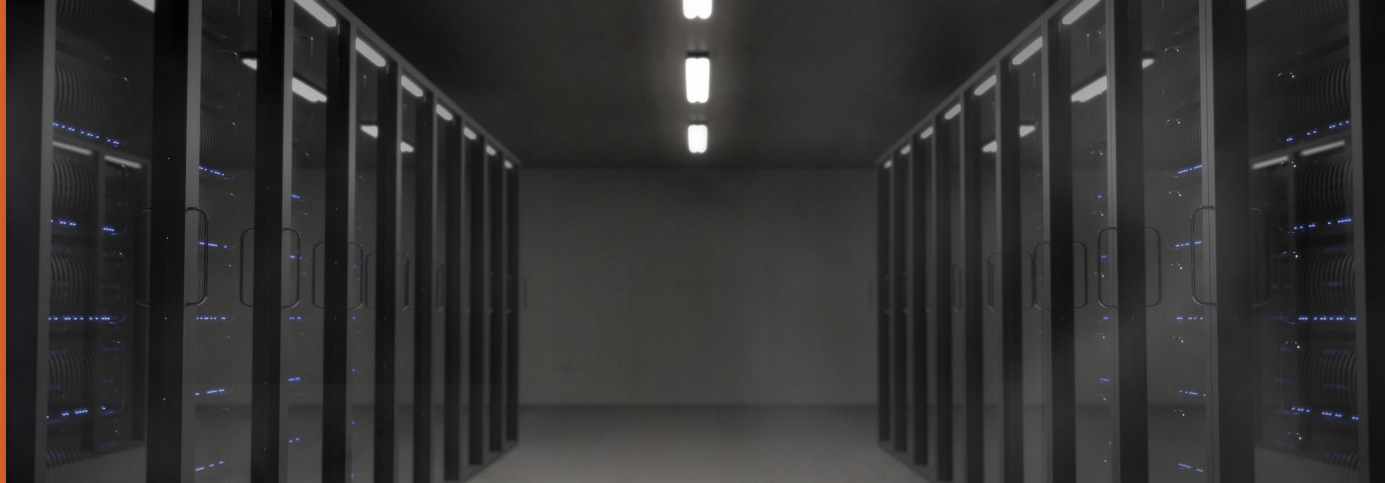- Highly predictive models recently published w/ best performance on this task to date

## SAR

- Physiological Data

- Receptor Level (depolarisation of cell)
- Human Brain Imaging - e.g. fMRI, PET, CT scans
- Mice, Flies and Worms - in vivo, high granularity imaging
- Response levels & haplotypes for Human ORs
- Measure of biological activity

- Not enough data to study population level responses
- Some available data from outdated studies

# Olfaction DREAM Challenge

- **The Rockefeller University:** Data collected in 2014
- **Team GuanLab** - winner of individual sub challenge
- **Team IKW** - winner of population sub challenge
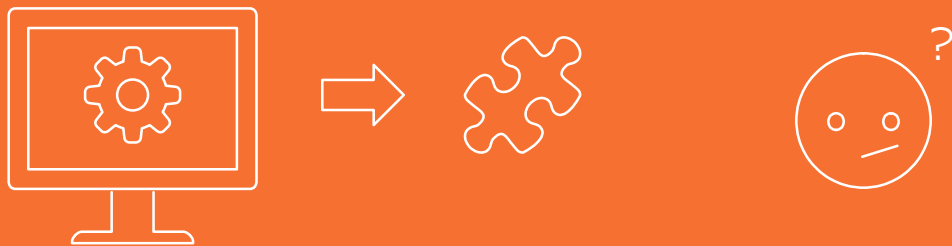- **Team BioLab** - runner up for population sub challenge



IBM Research, NY

Team GuanLab, MI

Team BioLab, Ljubljana

Semmelweis University, Budapest

National Taiwan University

Team IKW, AZ

The Rockefeller University, NY

# DATA

**PSYCHOLOGICAL DATA**

- 55 Subjects
- 480 Molecules
- 2 Dilutions
- Responses (0-100)
- 3 Mandatory Fields
    - Pleasantness
    - Intensity
    - Familiarity
- 19 Optional Fields e.g.
    - Fish
    - Fruit

**CHEMICAL/PHYSICAL DATA**

- 4884 features
- Describes Molecular Shape
- Describes Molecular Vibrational Frequencies
- 476 molecules (4 missing)
- Compiled using DRAGON software
- Top 20 descriptive features for each target published

# INTERPRETABLE MODELS

- Interpretability vs Completeness
- Natural Interpretation
- Post-hoc
- Local vs Global
- Model Specific/Agnostic
- Improve Trust & Discovery

OUTPUTTING A PREDICTION

OUTPUTTING A JUSTIFICATION OR EXPLANATION

# METHODS

Pre-processing, training & evaluation methods

# METHODS OVERVIEW

**Research**

- Olfactory System

- Machine Learning Methods

- Interpretability Methods

**Experimentation**

- Data Collection

- Network Analysis

- Training, training, training...

- Evaluation

**Documentation**

- Code - Jupyter Notebooks

- Blog

- Presentation & Report

Research → Experiment → Write up

# METHODS

Preprocessing

**PREPROCESSING**

- Cleaning Data

- Dilution/Concentration

- Calculating Mean Responses

- Combining Data

- Hidden Test Set Split

- Persists Datasets

# METHODS

Training Overview

**TRAINING MODELS**

- Baseline Modelling - 22 Targets

    - Default Hyperparameters

    - 10 Algorithms Tested

    - 5-fold Cross Validation

- Hyperparameter Tuning

    - Random Grid Search

    - 5 Algorithms Tested

    - 3-Fold Cross Validation

- Final Approach - 8 Targets, 3 set of models

    - Regularized Linear Models (L1)

    - Random Forest with & without reduced features

# METHODS

Baseline Predictions

Prediction Accuracy for Baseline - INTENSITY

# METHODS

Baseline Predictions

**BASELINE MODELS - PLEASANTNESS**

# METHODS

Baseline Predictions

**BASELINE MODELS - SEMANTIC DESCRIPTORS**



Prediction Accuracy for Baseline - Average of Semantic Descriptors

# METHODS

**TRAINING MODELS**

- Hyperparameter Tuning - Random Grid Search

- Testing Feature Spaces

  - Reduced Features

  - Principal Components

  - Mean vs Raw Responses

  - Imputation & Masking

- Reducing the considered Targets and Algorithms

- Evaluation & Comparison

  - Pearson R for each target/model pair

  - More details in Results

# METHODS

Top 20 published features

Delta Error Method

## INTENSITY

- Predictions worse across the board
- Intensity predictions are difficult with reduced descriptors - DREAM

## PLEASANTNESS

- Predictions worse for all except Random Forest
- Predictions improve for Random Forest Baseline



Prediction Accuracy for Baseline - INTENSITY



Prediction Accuracy for Baseline - PLEASANTNESS

# RESULTS

Visualising and interpreting our predictions

# AI

## PREDICTION
## ACCURACY

# RESULTS

Predictions

**PREDICTION CORRELATION - POPULATION**

# RESULTS

Residual Analysis

- shows that our error is largest around the mean
- This conforms to the observed data and previous analysis

## OBSERVED MEAN VS STANDARD DEVIATION

INTERPRET
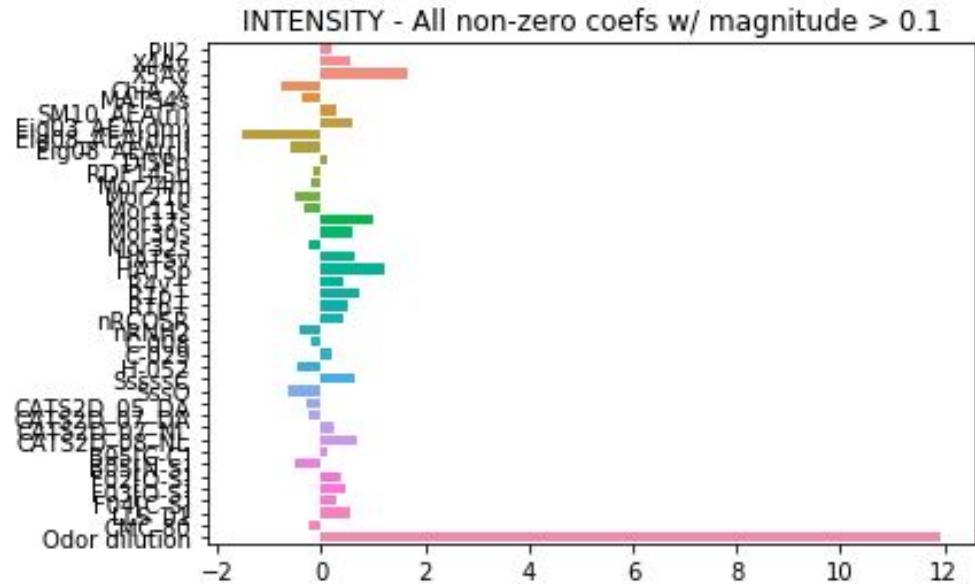MODELS

# RESULTS

Naturally Interpretable
Global Explanation

PLEASANTNESS - All non-zero coefs w/ magnitude > 0.1

# RESULTS

Naturally Interpretable
Global Explanation

INTENSITY - All non-zero coefs w/ magnitude > 0.1

# RESULTS

Naturally Interpretable
Global Explanation

## L1 REGULARISED LINEAR MODEL - FRUIT



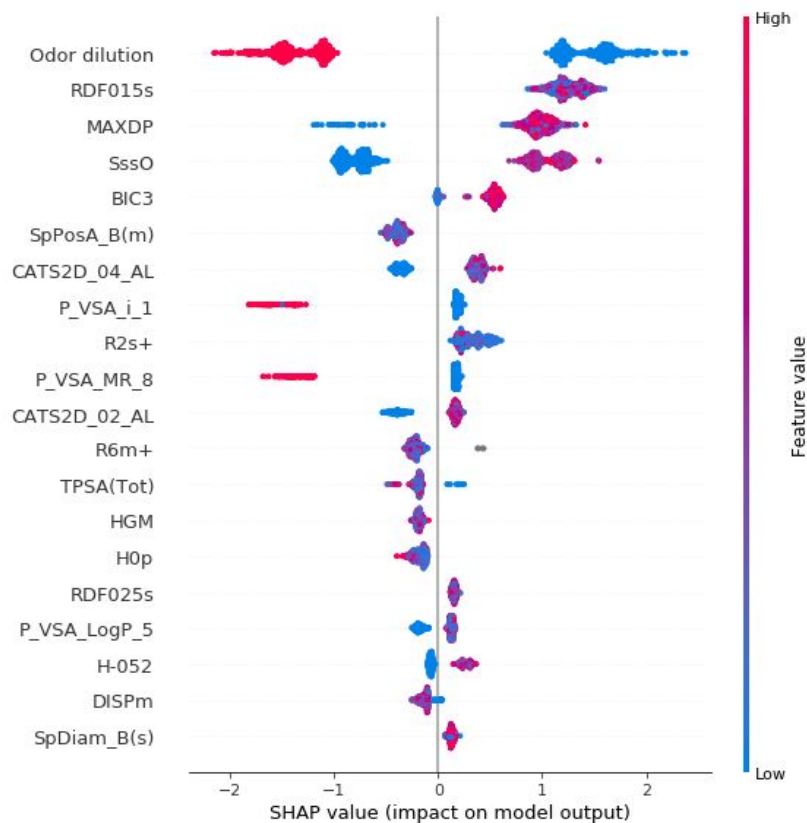FRUIT - All non-zero coefs w/ magnitude > 0.1

# RESULTS

Local, Post-hoc
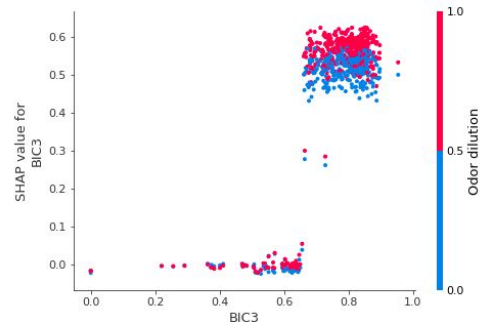Explanation

**RANDOM FOREST & SHAP**

# RESULTS
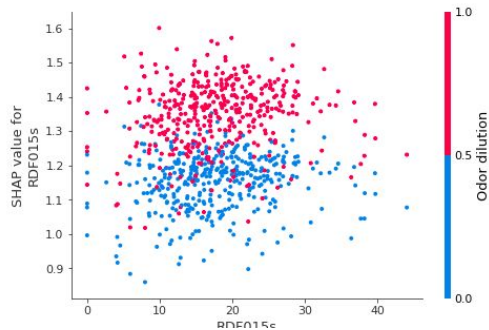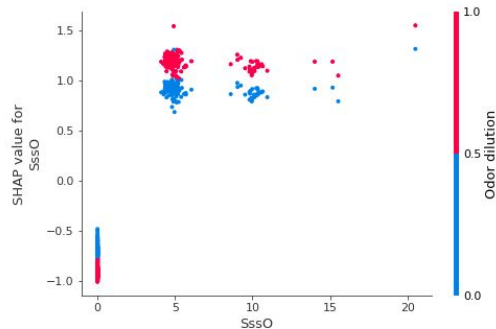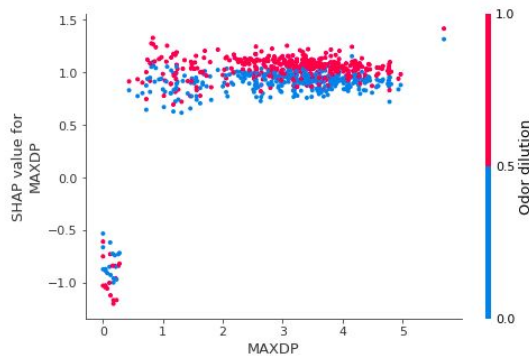
Global, Post-hoc
Explanation



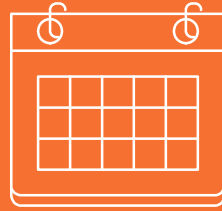RANDOM FOREST & SHAP - PLEASANTNESS

# RESULTS

Post-hoc Explanations
- Feature Dependence

**RANDOM FOREST & SHAP - PLEASANTNESS**

# DISCUSSION

Questions & Clarifications

# REFERENCES

Similar work and resources for presentation