# Modelling phase transitions: Characterising Henry's Law

*Joshua Cheung*
*Supervised by Jo Grundy and Jeremy Frey*

https://www.psdi.ac.uk/

# Colours and Fonts
## (remove from final presentation)

▶ Headings: Open Sans Condensed

▶ Body Text: Lato

▶ Colours: if these colour boxes match up the document has applied the colour theme.

| #011e41 | #002169 | #3d7cc9 | #ff9e18 |

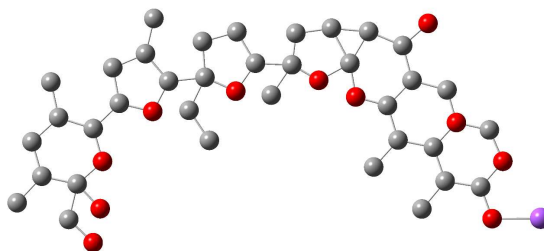| #5aa2ae | #bbbbbb | #e5e1e6 | #9d90a0 |

▶ Additional Colour:

| #993366 |

# About me



- 2nd Year Chemistry with Digital Methods and Computational Modelling
- I still don't have a photo of myself, so here's another cat photo I took
- Continuation of my project from last summer with PSDI: *Modelling Miscibility with Machine Learning*
- I enjoy reading sci fi, gaming, and playing guitar

# Presentation Outline

1. Project Description
2. Background
3. Methodology
   A. Dataset Building
   B. Data Processing
   C. Machine Learning
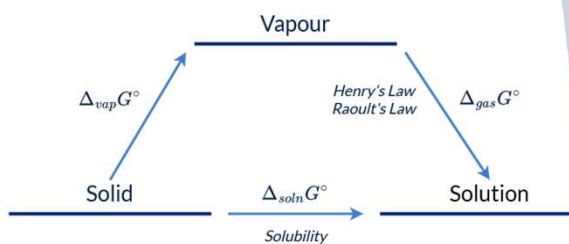   D. Results
4. Challenges
5. Conclusions and Future Work

*Monensin, an antibiotic, one of the compounds in the dataset*

# Project Description

- Exploring the links between Henry's Law constant ($k_H$) and Solubility (logS)
- Important for pharmaceuticals, synthesis in industry, electrochemistry, etc.
- Possible to calculate experimentally, but time consuming and difficult
- Extremely computationally expensive to model using quantum mechanical simulations
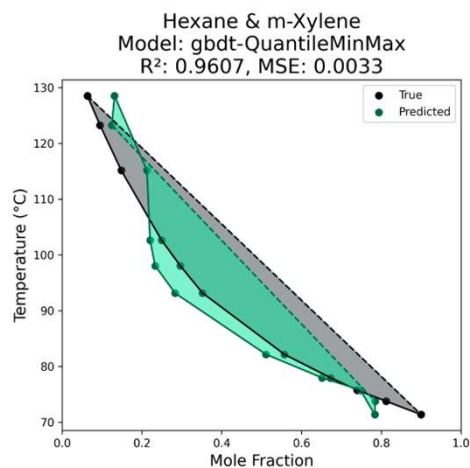
Solution: Machine Learning

Vapour

$\Delta_{vap}G^\circ$        Henry's Law        $\Delta_{gas}G^\circ$
                 Raoult's Law

Solid        $\Delta_{soln}G^\circ$        Solution
                 Solubility

Change since interim presentation:
No more CMC. Removed because there was insufficient data.

# Why Predict Henry's Law Constant?

Hexane & m-Xylene
Model: gbdt-QuantileMinMax
R²: 0.9607, MSE: 0.0033



*Temperature dependent phase-equilibrium diagram predicted using machine learning (2023)*

- Highly important properties with a wide range of applications in research and synthesis
- Existing data is limited and derived indirectly via experiment using an equation
- Existing models and papers use semi-empirical methods -> lack of experimentation with machine learning
- **Can be calculated if there is phase-equilibria data for the compound pair at a given temperature**

# Background

Henry's Law: *The abundance of a volatile solute dissolved in a liquid is proportional to its abundance in the gas phase.*

Solubility: *The analytical composition of a saturated solution, expressed in terms of the proportion of a designated solute in a designated solvent, is the solubility of that solute.*

**Or more simply: How much of the substance (liquid or gas) can be dissolved in the solvent?**
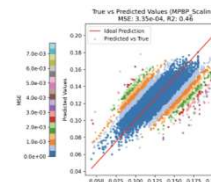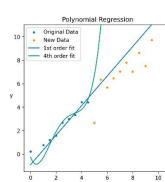


*https://iupac.org/recommendation/henrys-law-constants/*
*https://goldbook.iupac.org/terms/view/S05740*
*https://www.youtube.com/watch?v=zMaTrgUKC1w*

Overall process

# Dataset Building

# Dataset Building

| Feature | Entry Count | Unique Compound Count | Most abundant |
|---|---|---|---|
| logS (solubility) | 11 703 | 11 089 | Cyclohexanol (24) |
| Henry's law | 12 167 | 9 532 | Toluene (24) |
| Total Entries | 21 291 | 18 006 | **Cyclohexanol (26), Pentan-1-ol** |

Scaling affects distribution of data, makes it easier for model to effectively differentiate between values
Feature selection tested variety of different ways, eventually just removed anything with <1% distinct values

# Machine Learning and Results

- 2 types of machine learning models: **Regression** and Classification

- **Prediction of continuous values** -> not a classification problem

- No way of predicting which algorithm will work best, so multiple tested

  - How to score models and compare them to figure out what works best?

| $R^2$ | Mean Squared Error (MSE) |
|---|---|
| A measure of how well the predictions and true values correlate with each other. A perfect correlation would have a score of 1. $$R^2 = 1 - \frac{\Sigma_i(y_i - \hat{y}_i)^2}{\Sigma_i(y_i - \bar{y})^2}$$ | A metric of the Euclidean distance between the predicted value and the true value. The higher the error, the worse the prediction. $$\text{MSE} = \frac{\Sigma(y_i - \hat{y}_i)^2}{n}$$ |

# Choosing the Algorithms

## logS

| Algorithm | MSE | $R^2$ |
|---|---|---|
| **LightGBM** | **1.01** | **0.80** |
| KRR | 1.06 | 0.79 |
| KNN | 1.23 | 0.76 |
| RandomForest | 1.28 | 0.75 |
| AdaBoost | 1.53 | 0.70 |
| SVR | 1.65 | 0.68 |

## log($k_H$)

| Algorithm | MSE | R2 |
|---|---|---|
| **LightGBM** | **5.04** | **0.89** |
| KRR | 5.92 | 0.87 |
| RandomForest | 8.74 | 0.81 |
| KNN | 9.31 | 0.81 |
| AdaBoost | 12.76 | 0.73 |
| SVR | 68.92 | -0.45 |

Due to evolving datasets, these results are not for the version of the dataset that was used for the final results.

# Results and Analysis

- Trained model to predict logS using data that has no $k_H$ value
- Used model to predict logS and fill in dataset gaps
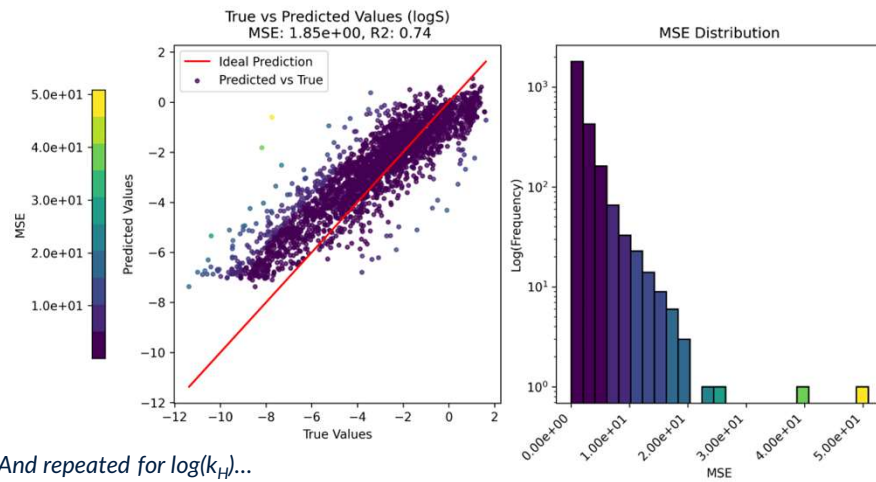- Withheld test set of data that contains logS and Henry's constant value

# Test Set: logS



0.2.2 logS LGBM MRobust

*And repeated for log(k_H)...*

# Test Set Metrics: logS



MSE Distribution for LogS Predictions
(LGBM, 0.2.2, MRobust)
Max = 5.09e+01

| R$^2$ | 0.73 |
|---|---|
| MSE | 1.85 |
| % MSE > 1 | 45.52 % |



Lowest error: $6.00 \times 10^{-7}$    Highest error: 50.80

- 6[th] highest error in the test set is ammonia, $NH_3$, with an MSE of 18.51.
- Only bottom 5 have MSE > 20

# Test Set: log(k$_H$)

### 0.2.2 logS LGBM MRobust

The Antilog: Test Set Metrics

MSE Distribution for Henry's Constant Predictions
(LGBM, 0.2.2, MRobust)
Max = 3.54e+24

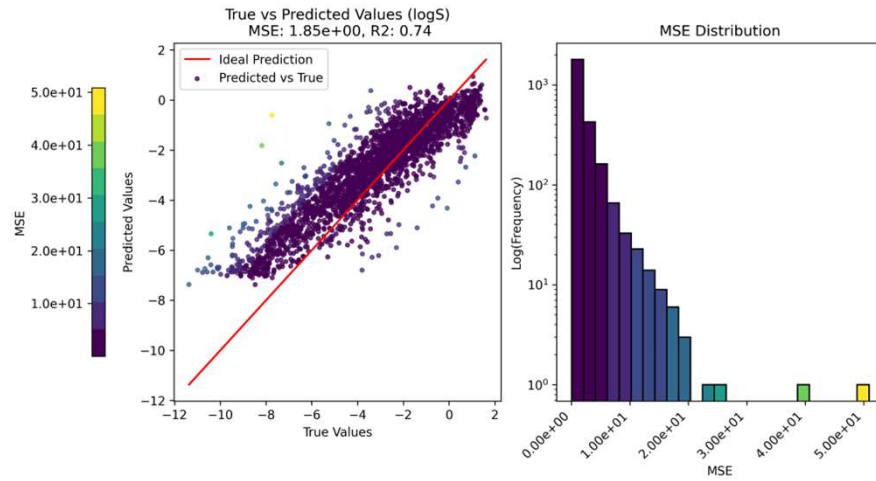MSE Distribution for Log Henry's Constant Predictions
(LGBM, 0.2.2, MRobust)
Max = 1.64e+03

No overlap in 10 worst predictions, but overlap in best 10.

|  | Log Scaled | Unscaled |
|---|---|---|
| $R^2$ | 0.66 | -68.67 |
| MSE | 19.11 | $2.13 \times 10^{21}$ |
| % MSE > 1 | 76.53 % | 52.93 % |

Highly likely that these results could be improved by rerunning hyperparameter optimisation, etc, as dataset size doubled with bug fix

# Choosing the Algorithms

**log($k_H$)**

| Algorithm | MSE | R2 |
|---|---|---|
| **LightGBM** | **5.04** | **0.89** |
| KRR | 5.92 | 0.87 |
| RandomForest | 8.74 | 0.81 |
| KNN | 9.31 | 0.81 |
| AdaBoost | 12.76 | 0.73 |
| SVR | 68.92 | -0.45 |

- Henry's law dataset doubled in size after fixing unresolved identifiers.
- Best algorithm, scaling methods, and hyperparameters likely to have changed.
- Results from previous slide could almost definitely be improved (time limitations in project)
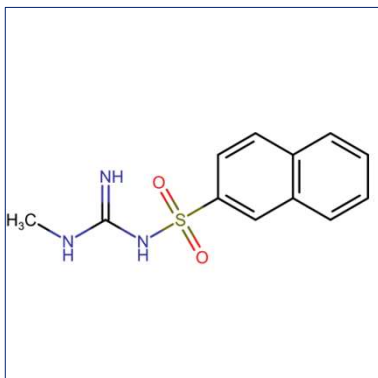
# Challenges

- Questionable quality of data, e.g. range of values, number of significant figures, data sources
  - Lack of temperature data for Henry's Law constant (all assumed to be at 25c instead)
- Identifier conversion (CAS to SMILES, SMILES to InChI) leading to incorrect data in Henry's Law dataset.
  - ~33 compounds had their SMILES replaced with O, leading to an abundance of chemicals mistakenly being identified as water
  - Over 2000 had unresolved InChI keys which had to be resolved using alternative methods
- My own programming skills. Dataset generation with melting points takes over 24 hours to run, and is probably very poorly optimized.
  - I could run these programs on Iridis, but errors are likely with file paths and saving.
  - Runs a lot faster when I stop adding melting point data (~250 000 data points) for no good reason.
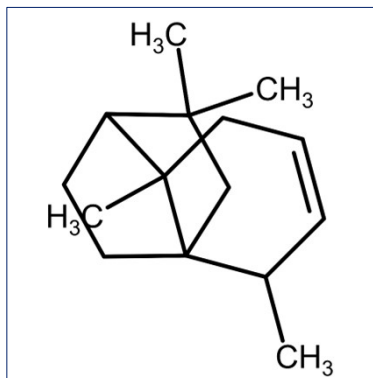
# Questionable Identifiers



| N-(N-Methylcarbamimidoyl)-2-naphthalenesulfonamide |
| --- |
| CNC(=N)NS(=O)(=O)C1=CC2=C(C=CC=C2)C=C1 |
| InChI=1S/Mo |



| γ-neoclovene |
| --- |
| O |
| InChI=1S/H2O/h1H2 |

Resolved by using PubChem API instead of NIST, and scraping UoY Master Chemical Mechanism site

# What would you do differently?

- Better data sanitization to identify compounds which have erroneous values or identifiers
  - Actually looking at the output of my programs properly instead of just assuming it's correct because the first 100 worked!
- Identify a more solid aim and purpose at the beginning of the project so less time is wasted
  - Scope was too wide at the beginning of the project, including melting/boiling point, and CMC data.

# Conclusions and Future Work

It is possible to predict solubility and Henry's Law constant for an aqueous system

Future Work:

- Try using recursive feature elimination to see why no scaling or normalisation had good results
- Explore links with melting and boiling points
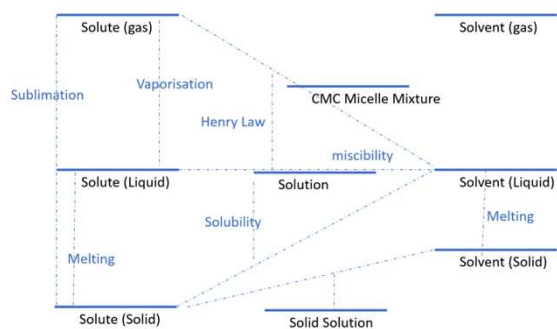- Look into why the models predicts better for certain types of compound



*Diagram by Jeremy Frey*