



optibrium

# Gaussian Processes for QSAR Modelling

Rigorously defined areas of doubt and uncertainty

Peter Hunt & Matthew Segall

[peter.hunt@optibrium.com](mailto:peter.hunt@optibrium.com), [matt.segall@optibrium.com](mailto:matt.segall@optibrium.com)



“We demand rigidly defined areas of doubt and uncertainty”

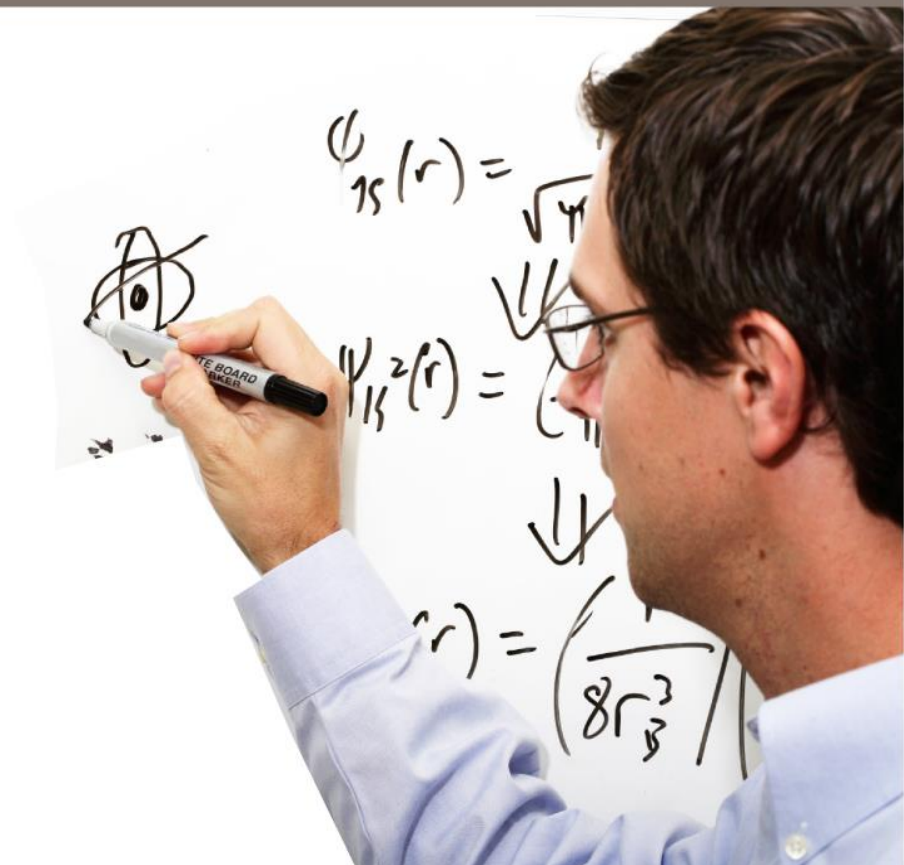
Douglas Adams, Hitchhiker's Guide to the Galaxy

# Overview

---

- Uncertainty in model predictions
  - Domain of applicability
- Introduction to Gaussian Processes (GPs)
- Simple illustration of GPs
- Automatic relevance determination
- Practical examples of GPs applied to QSAR
- Conclusions

# Uncertainty in Model Predictions



# Quantitative Structure-Activity Relationships

$$y = f(x_1, x_2, x_3, \dots) \pm \varepsilon$$

Statistical  
uncertainty



- Data
  - Quality data is essential
  - Public data needs very careful curation\* (and may not be good enough)
- Descriptors, e.g.
  - Whole molecule properties, e.g. logP, MW, PSA...
  - Structural descriptors, SMARTS, fingerprints...
- Statistical fitting or machine learning method, e.g.
  - Multiple linear regression, partial least squares
  - Artificial neural networks, support vector machines, random forest, **Gaussian processes**...

# Sources of Uncertainty in Model Predictions

---

- Experimental noise in training data
- Descriptors may not capture all sources of variation
  - Modelled property may not be ‘smooth’ in descriptor space, limiting ability to interpolate
- New compound may be ‘different’ from those used to train the model
  - ‘Domain of applicability’
  - Models often have a limited ability to extrapolate beyond the descriptor space represented by the training set

# Assessing Predictive Ability

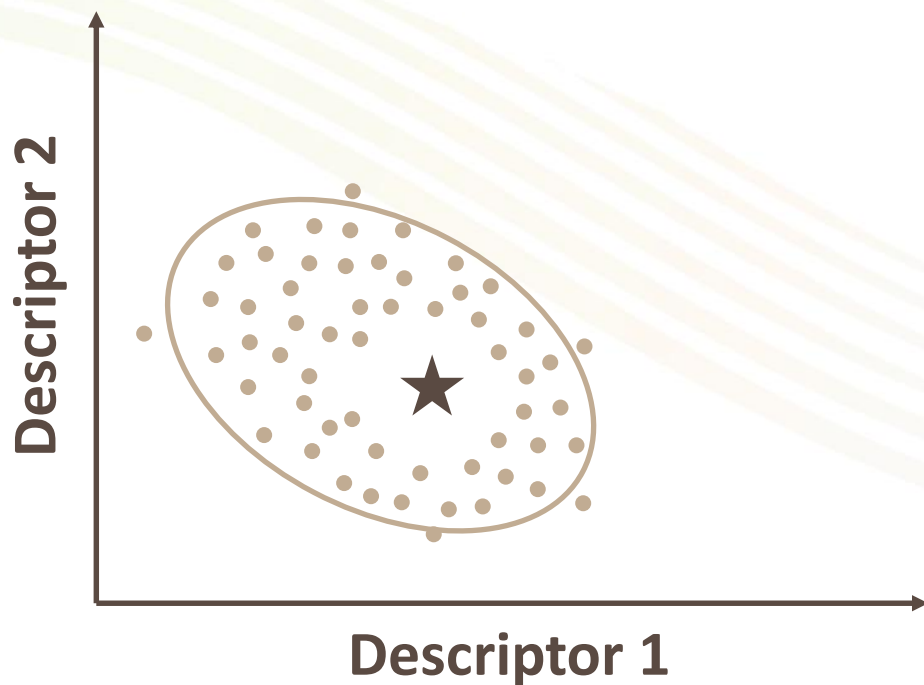
## Domain of Applicability



- The diversity of the training set defines the **domain of applicability** of the model
- The position of a new compound relative to the domain of applicability should be reflected in the reported confidence in the prediction
- Can we do better than 'in' or 'out' indication?

# Assessing Predictive Ability

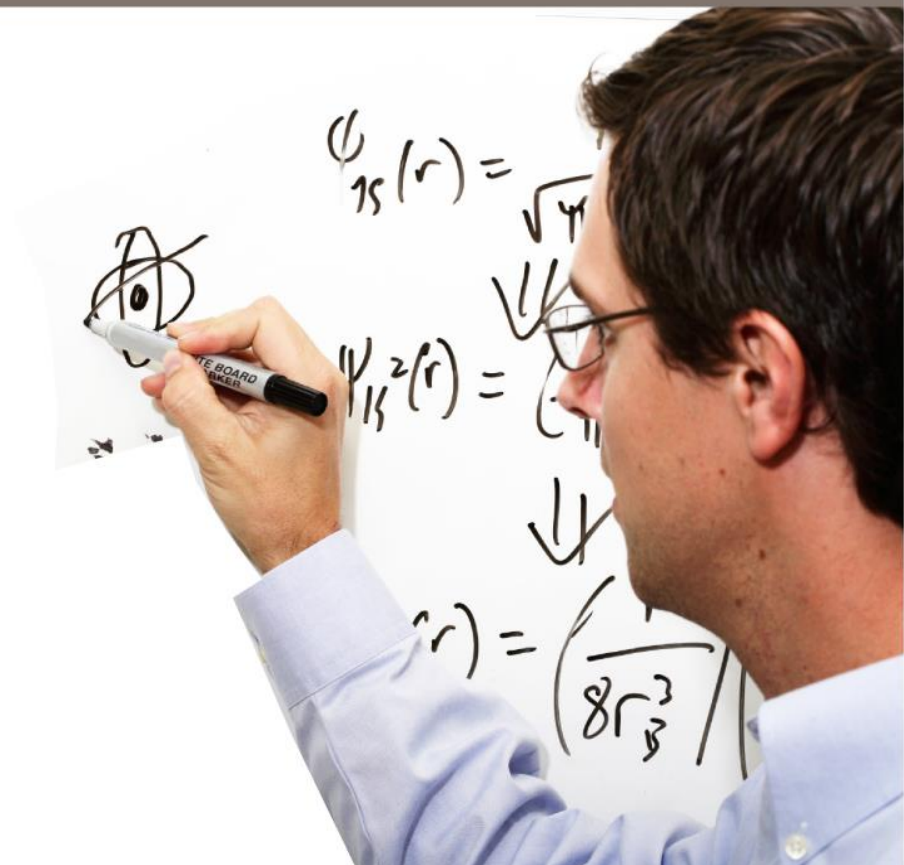
## Dealing with 'gaps' in coverage



- Distribution of the training set may not be uniform and there may be 'gaps' or sparsely sampled regions
- Is this compounds 'in' or 'out' of the domain of applicability?



# Introduction to Gaussian Processes (GPs)



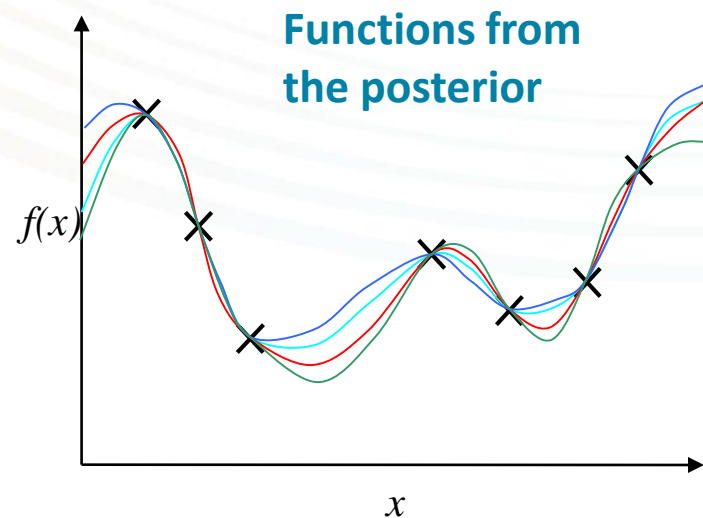
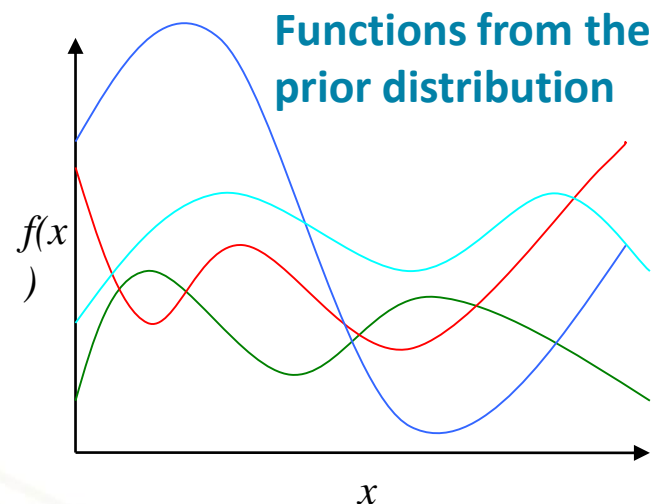
# Modelling Techniques: Gaussian Processes

---

- A machine learning method based on a Bayesian approach
- Advantages:
  - Does not require a priori determination of model parameters
  - Nonlinear relationship modelling
  - Built-in tool to prevent overtraining - no need for cross-validation
  - Inherent ability to select important descriptors
  - Provides uncertainty estimate for each prediction
- Sufficiently robust to enable automatic model generation

# Modelling Techniques: Gaussian Processes

- Define **prior distribution** over functions (controlled by hyperparameters, covariance function – ARD function)
- **Posterior distribution**: retain functions which fit experimental data
- **Prediction** is the mean of posterior distribution.
- Standard deviation of the distribution provides estimate of the **uncertainty in prediction**



# Gaussian Processes: Hyperparameters

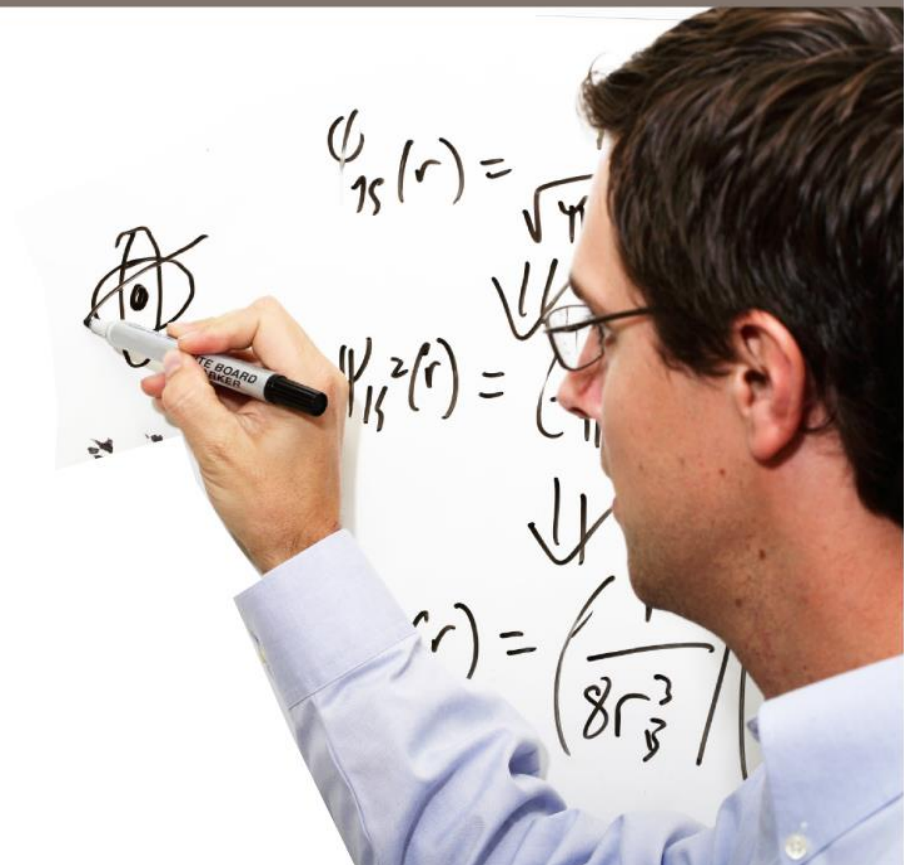
- Learning the Gaussian Process  $\sim$  finding hyperparameters
  - Optimize the marginal log-likelihood (prevents overtraining)
  - Fits parameter corresponding to estimate of noise in input data (assuming normally distributed)
- Techniques for finding hyperparameters
  - Fixed length scales (Fixed)
  - Forward variable selection (FVS)
  - Conjugate gradient optimisation (Opt)
  - Nested sampling (Nest)

Automatic  
Relevance  
Determination

computational demand

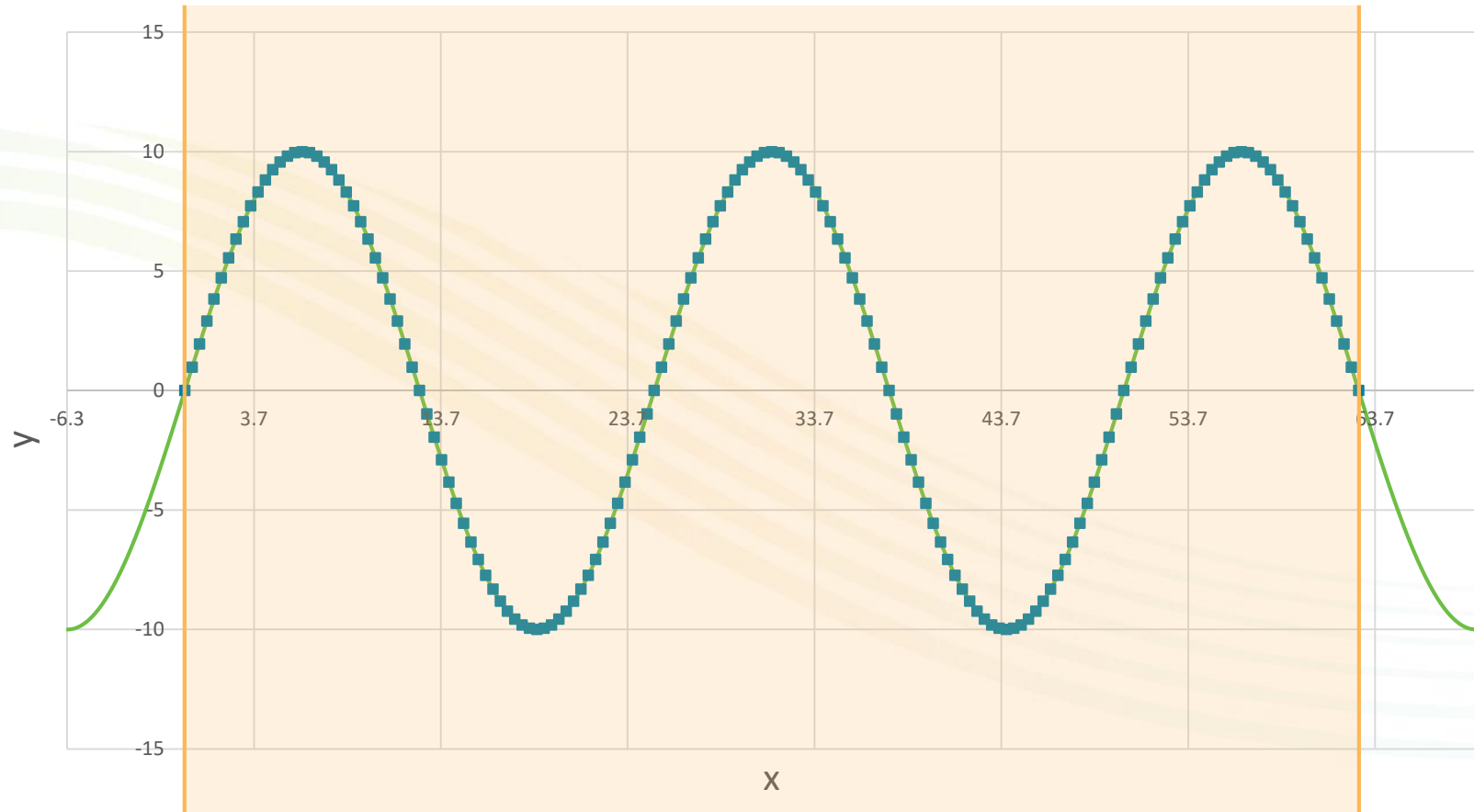


# Simple Illustration of GPs



# 'Toy' Example

## Training set from sin function in 1 dimension



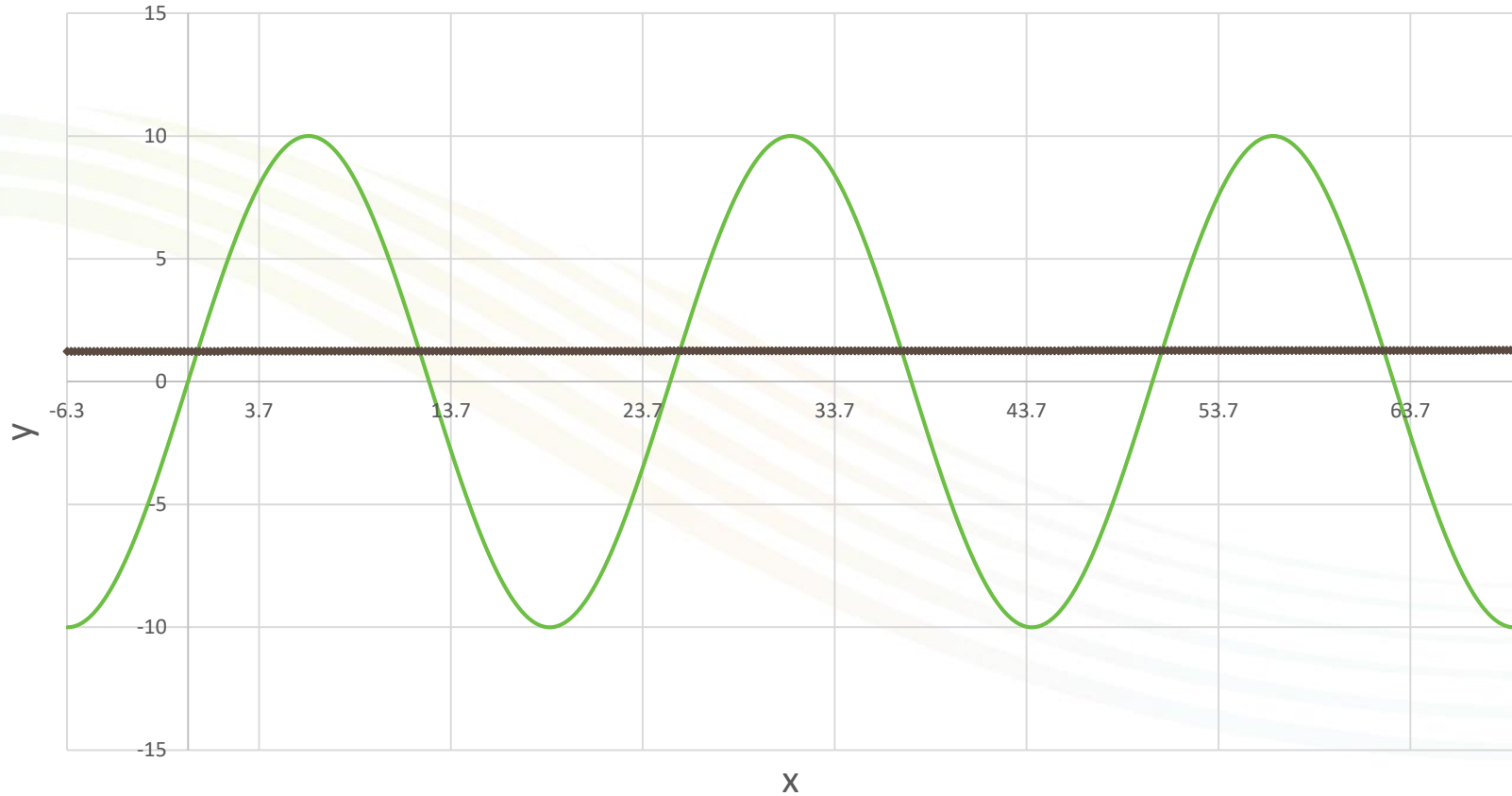
Lower boundary  
of training set

— sin function    ■ Training set

Upper boundary  
of training set

# Partial Least Squares

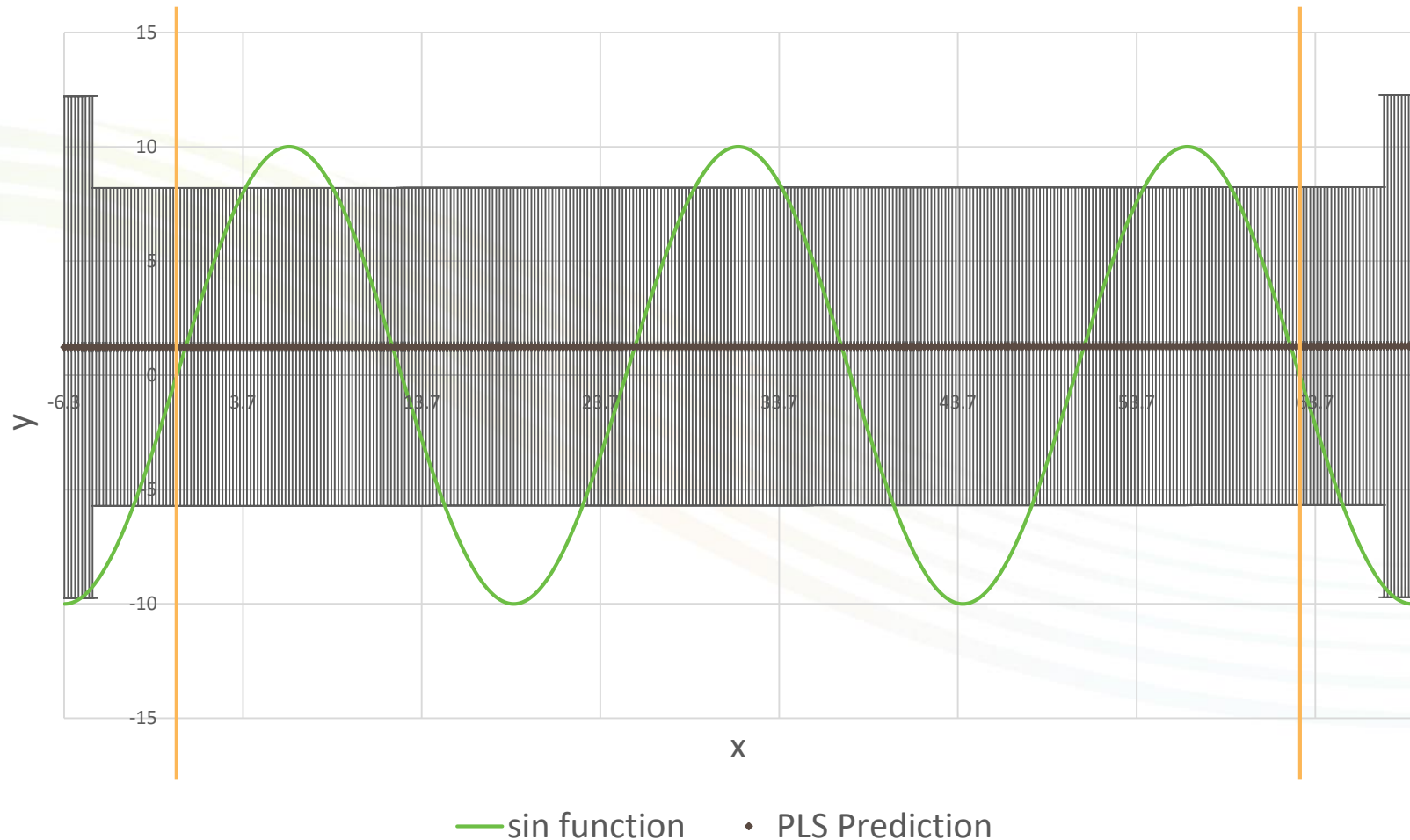
## Linear model not appropriate



— sin function    ♦ PLS Prediction

# Partial Least Squares

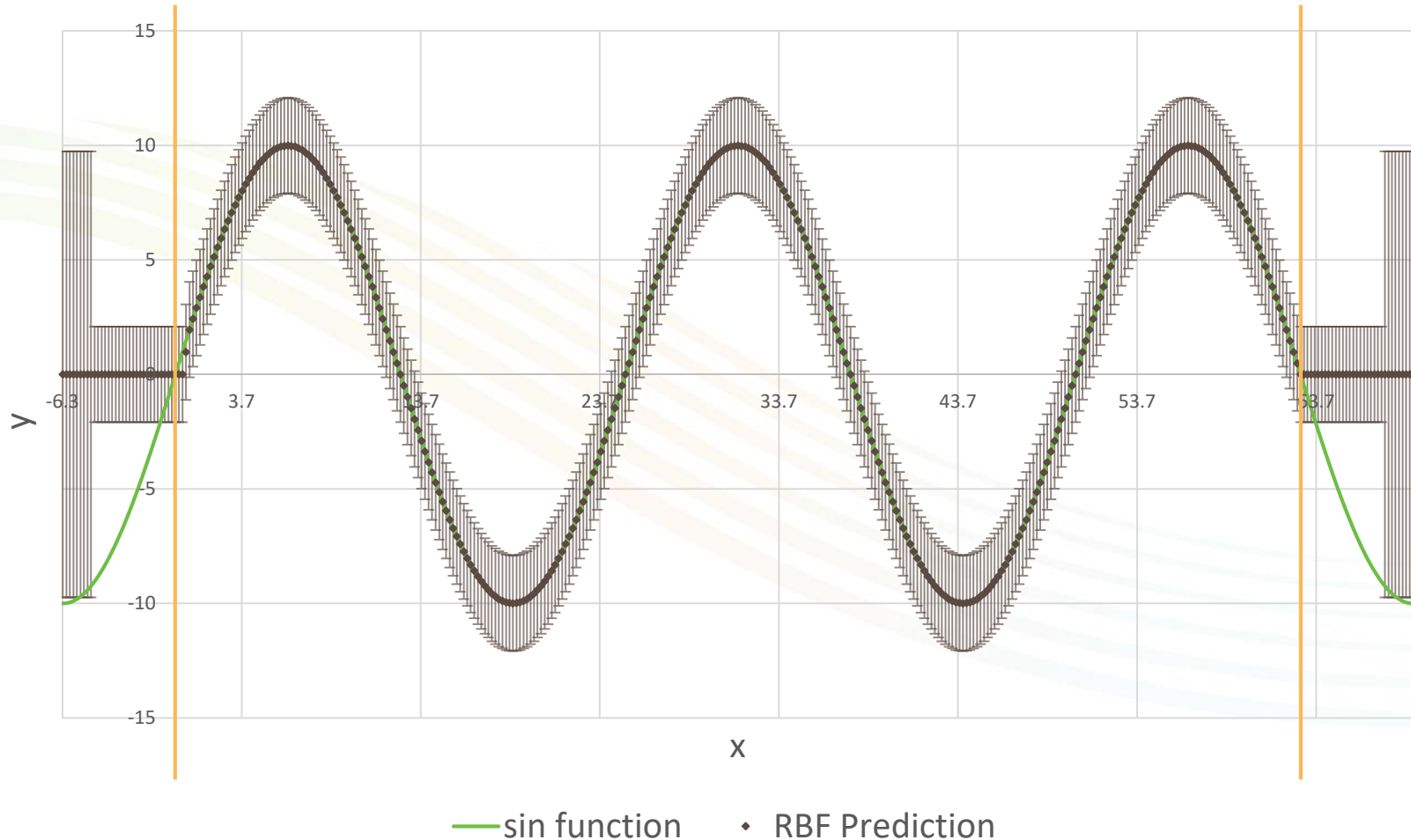
## Linear model not appropriate



Domain of applicability based on Hotelling's  $T^2$  test with 95% confidence limit  
Error bars based on RMSE error inside and outside of domain of applicability

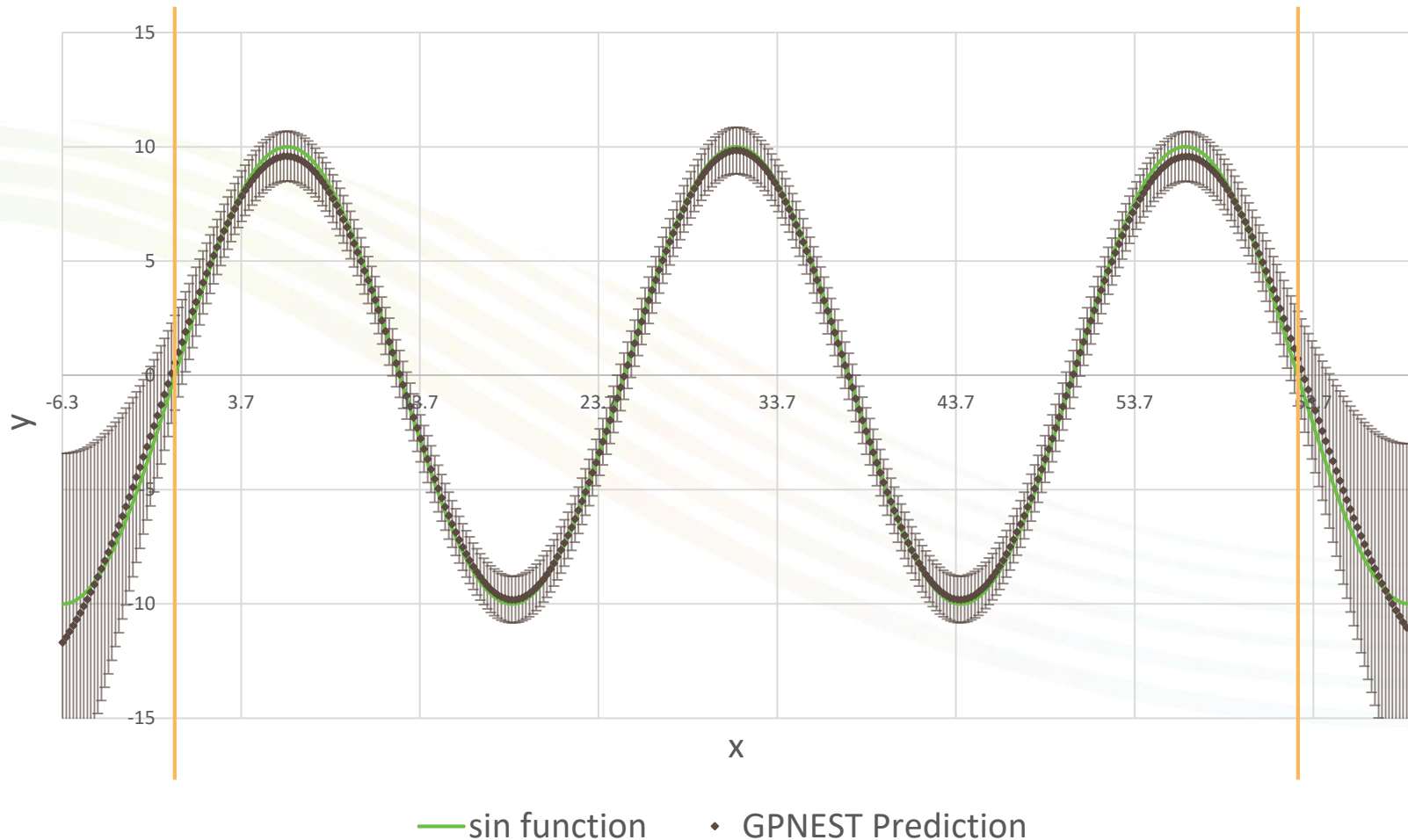


# Radial Basis Function Model



Domain of applicability based on Hotelling's  $T^2$  test with 95% confidence limit  
Error bars based on RMSE error inside and outside of domain of applicability

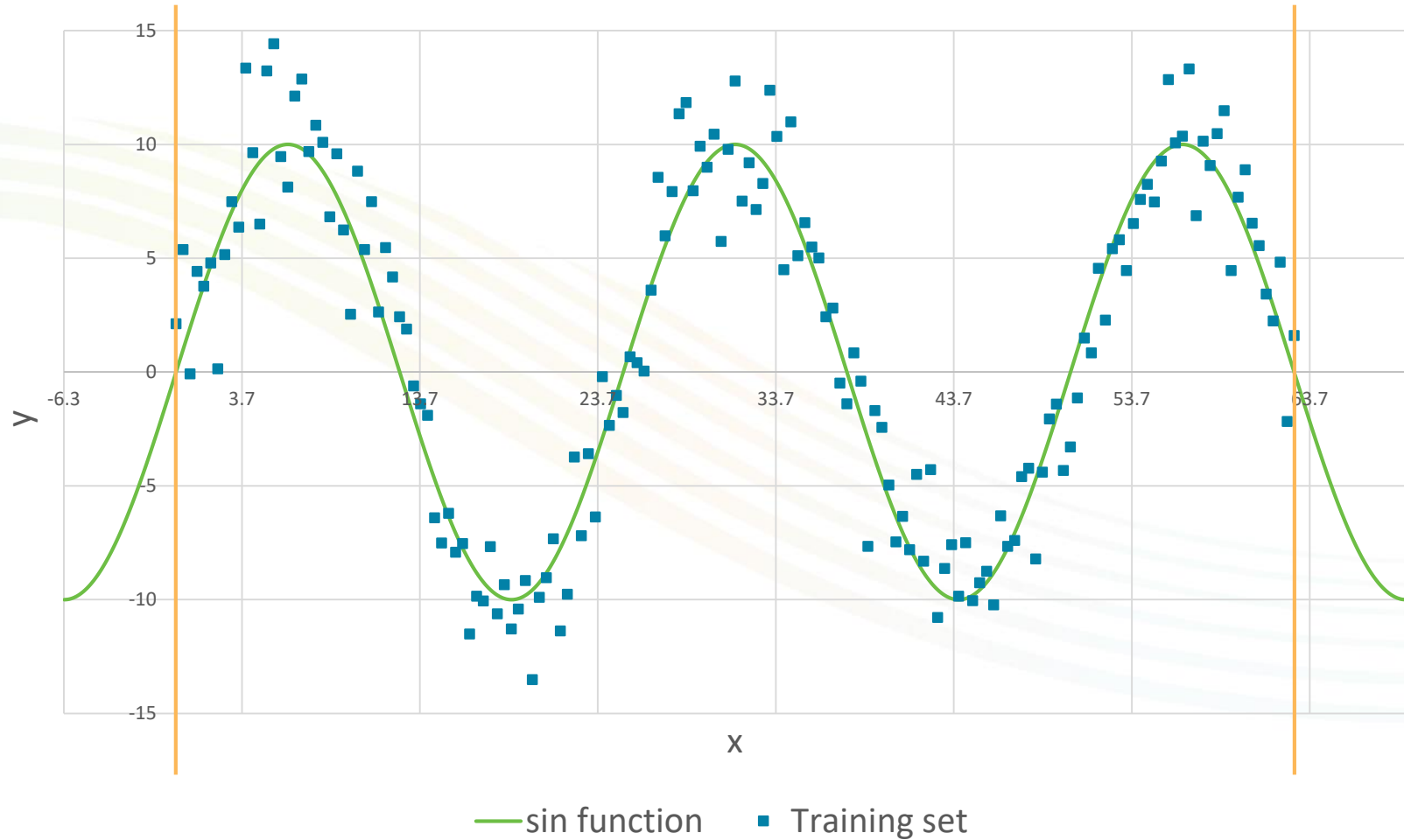
# Gaussian Processes (Nested Sampling)



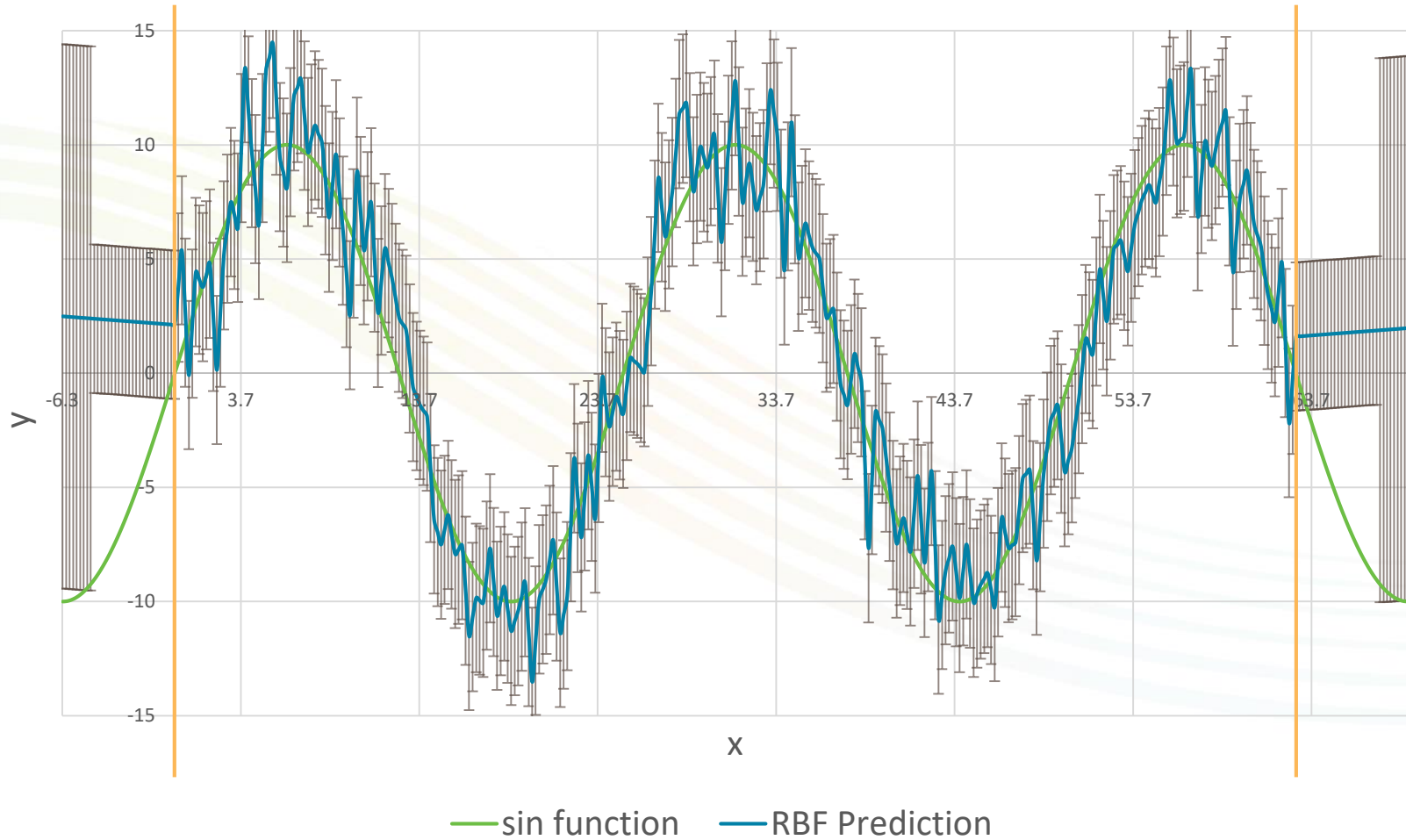
Error bars estimated by standard deviation of Gaussian process for each predicted value

# Training Set with Noise

Normally distributed error with standard deviation of 2

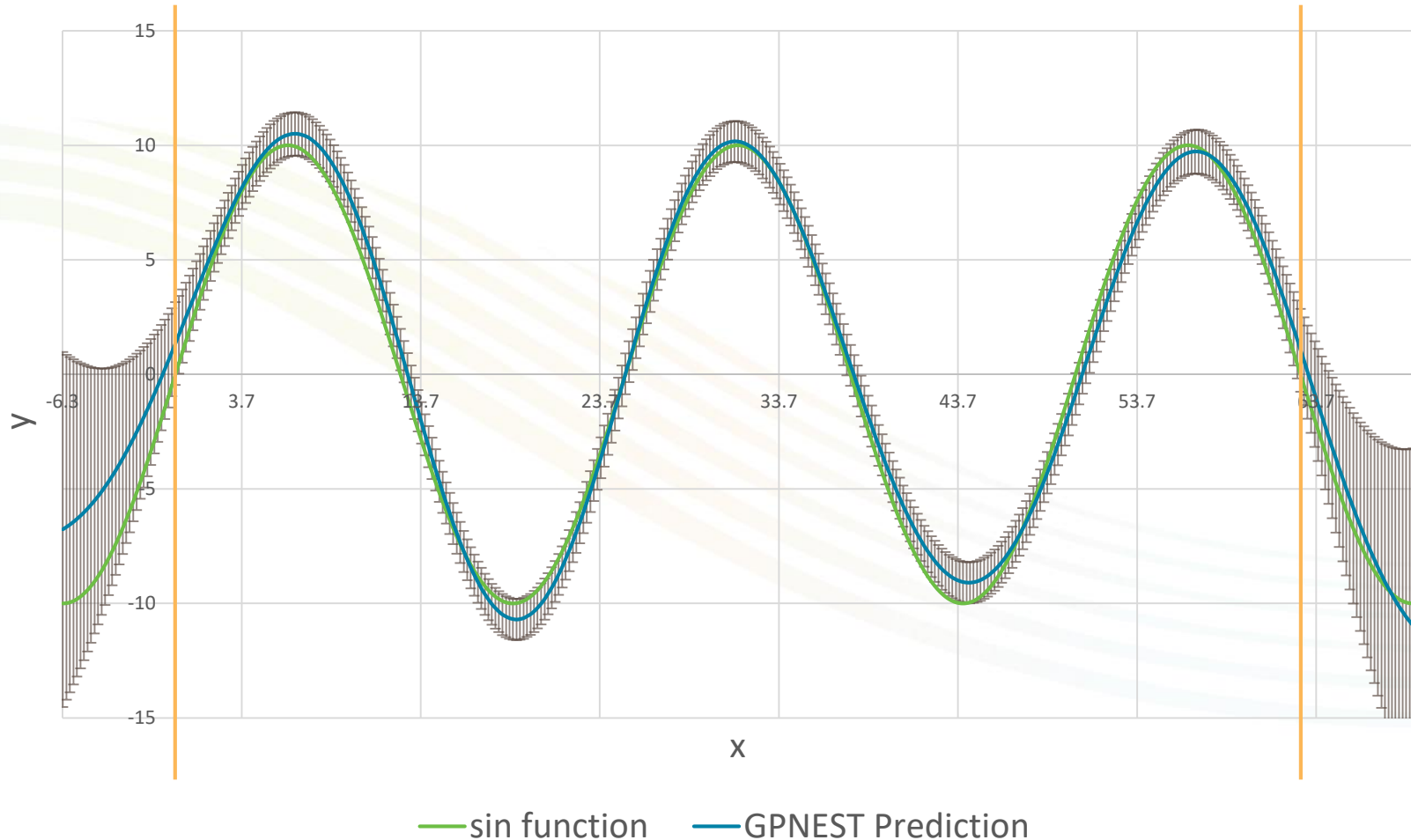


# RBF Model of Noisy Data



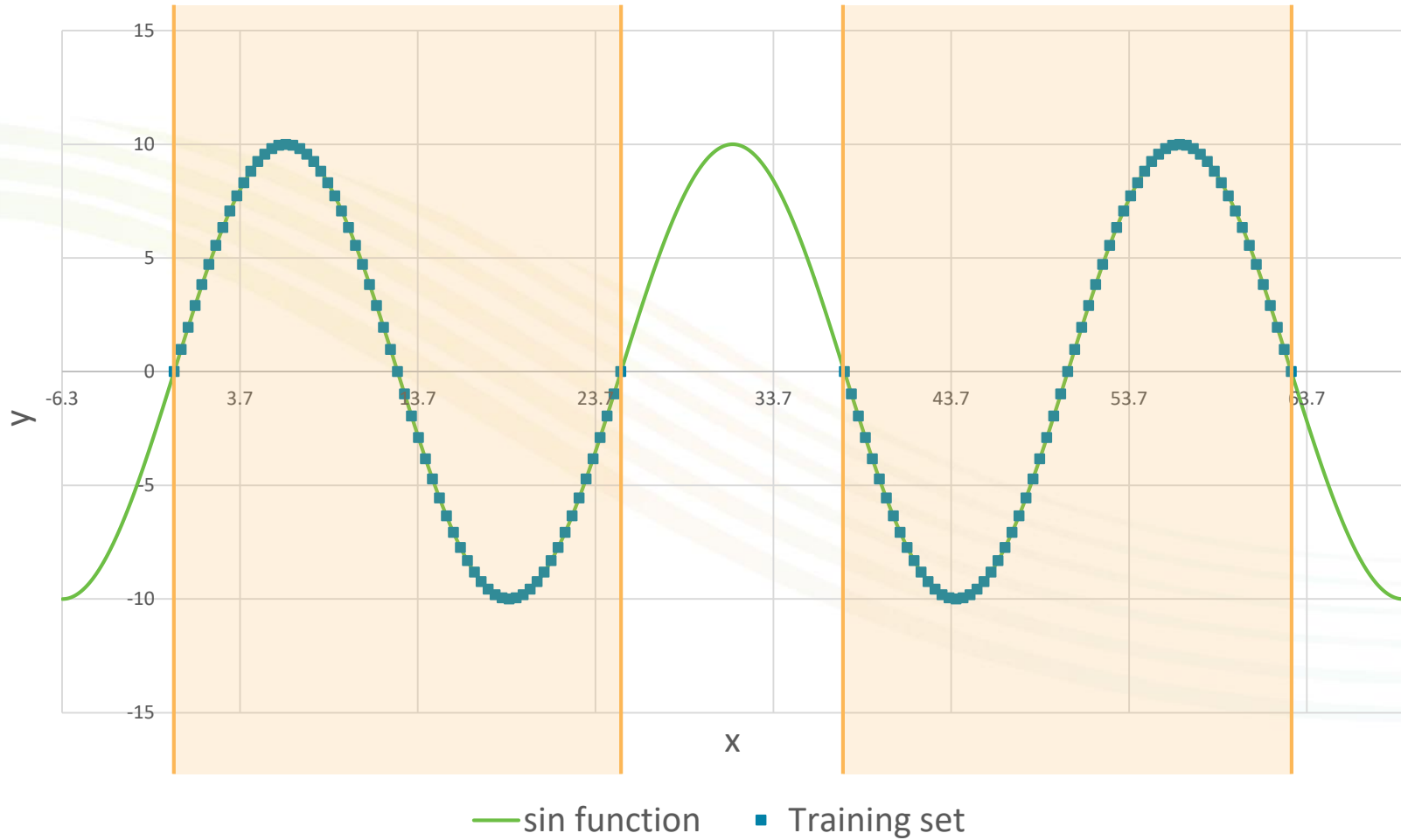
Greater error in prediction in domain of applicability, so error bars increase accordingly

# GP Model of Noisy Data

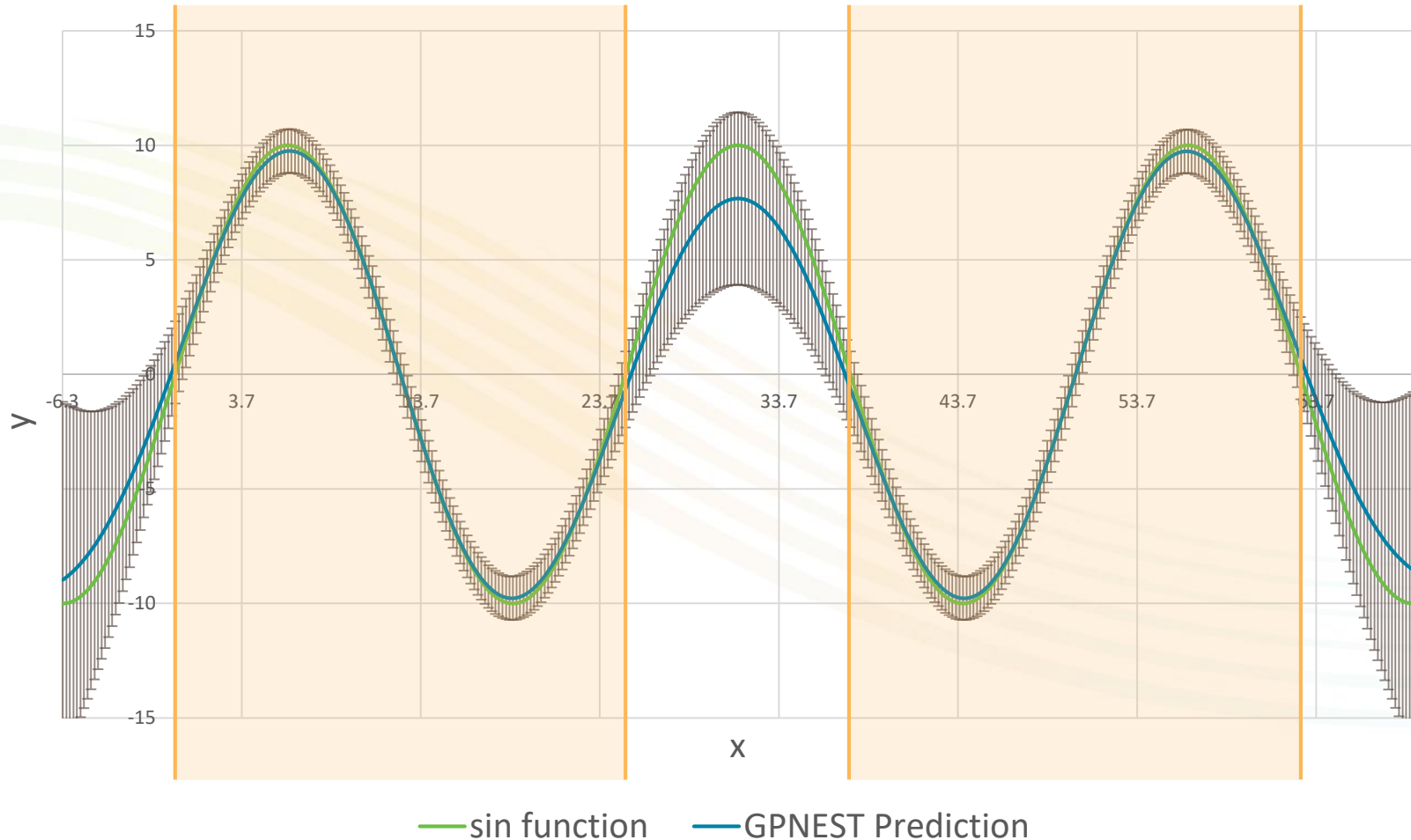


GP infers underlying functional form, fits model of noise and corrects error bars accordingly  
More difficult to extrapolate, but this is accounted for outside domain of applicability

# Training Set with Missing Data



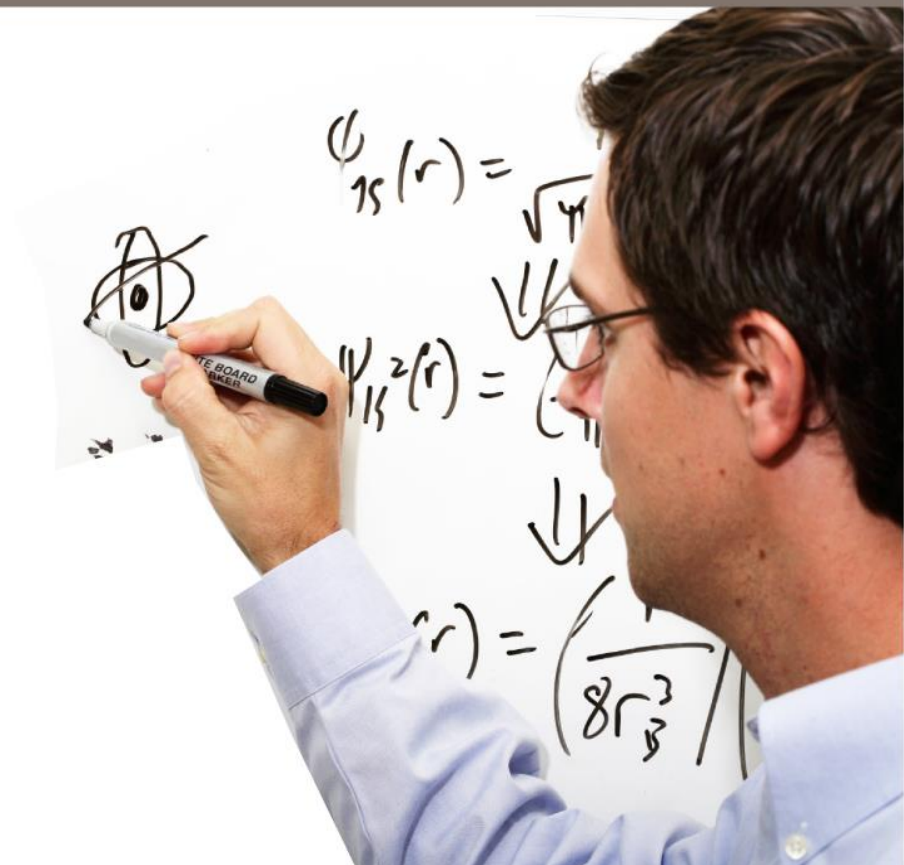
# GP Model Built with Missing Data



Where data is sparse or missing uncertainty in prediction is higher, as reflected by larger error bars

# Automatic Relevance Determination

## Identifying most important descriptors

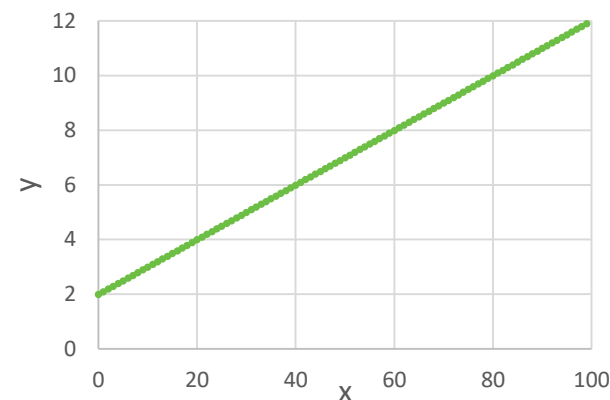




# Experiment

## Detecting relevant descriptors

- 100 training data points
- One descriptor (x) with perfect linear correlation with property (y)
- Hide this descriptor in a data set containing  $N_{\text{random}}$  randomly generated descriptors

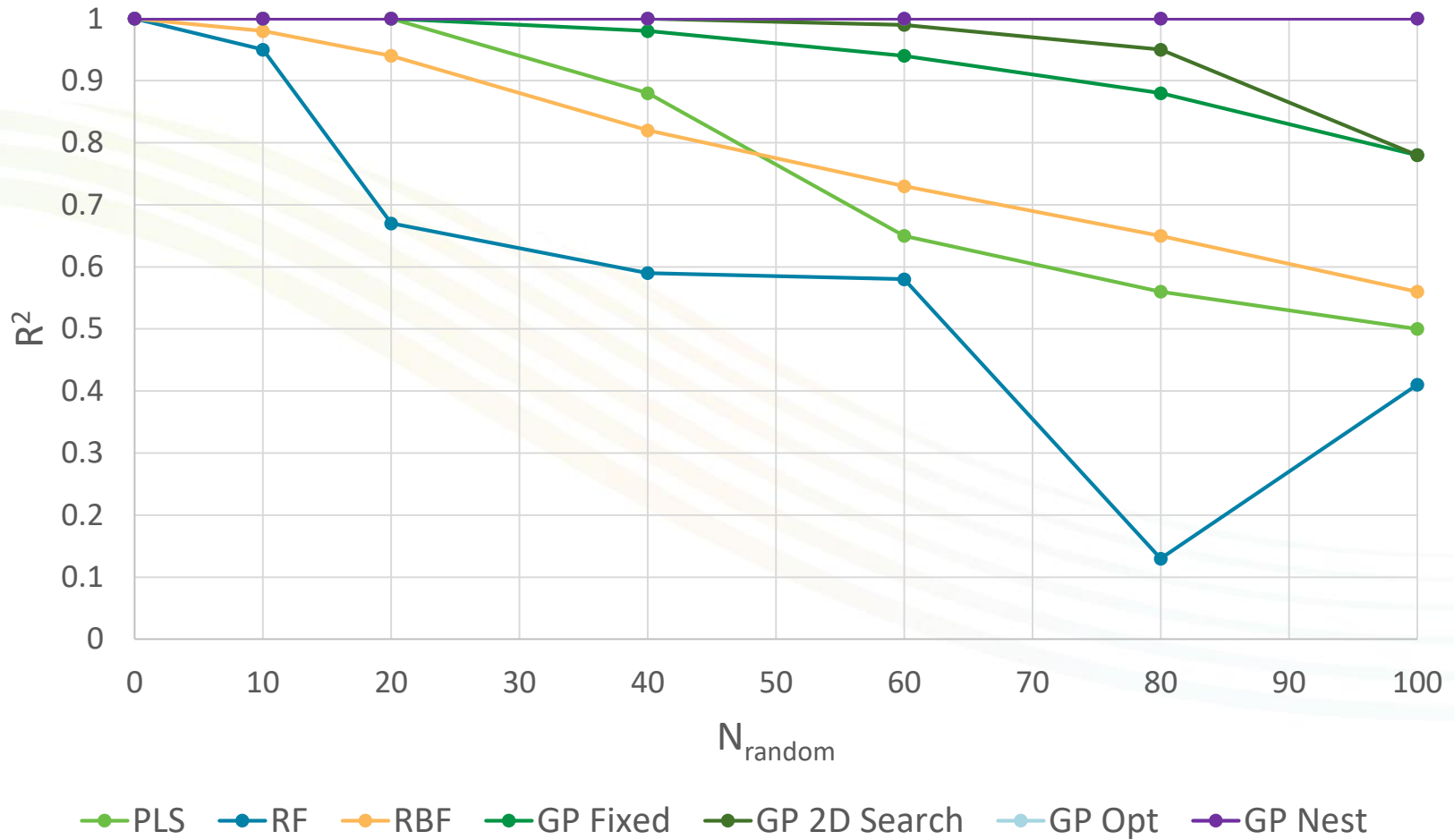


Identifier	y	x	Rnd 1	Rnd 2	Rnd 3	Rnd 4	...	Rnd $N_{\text{random}}$
Compound 1	0	2	1.252494	4.741985	2.14597	9.457343	3.958759	4.780421
Compound 2	1	2.1	8.64592	1.747653	6.76429	3.99527	7.626024	8.251326
Compound 3	2	2.2	7.064635	5.553097	0.355306	9.588649	2.791829	9.871042
Compound 4	3	2.3	4.783329	5.126768	3.525646	2.39005	0.392087	7.550868
...	4	2.4	7.077723	1.938555	2.028159	7.487378	9.672227	9.300353
Compound 100	99	11.9	0.588886	7.136536	9.538188	1.295742	3.522841	9.480185

- Can a method find the relevant descriptor?

# Results

## Detecting relevant descriptors

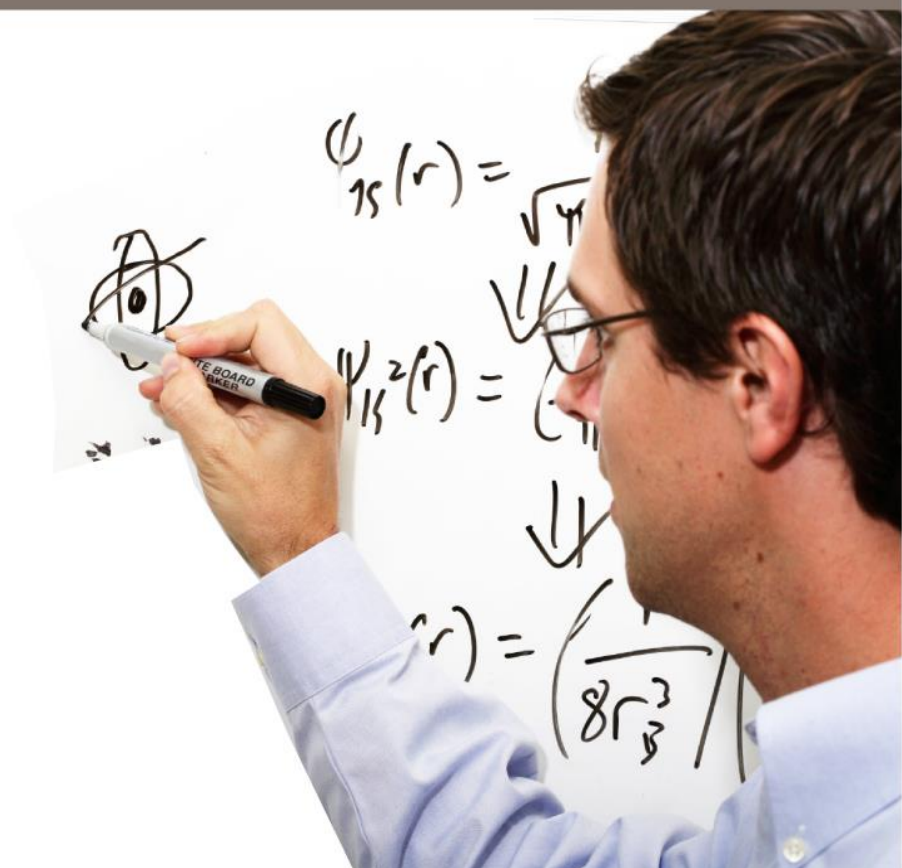


# Number of Descriptors Use in Model

---

- Often quoted rule of thumb... at least 5 compounds in training set per descriptor
- This is relevant for simple models where the only complexity control is the number of descriptors in the model
- But, GP Nest model includes all descriptors
  - Influence of random descriptor on model is negligible
  - Posterior probability of complex models is low
- Including additional non-influential descriptors in the model can be valuable
  - Better definition of the domain of applicability
  - Detect when new compound differs significantly from training set
  - Uncertainty in prediction will increase

# Practical Application to QSAR Modelling



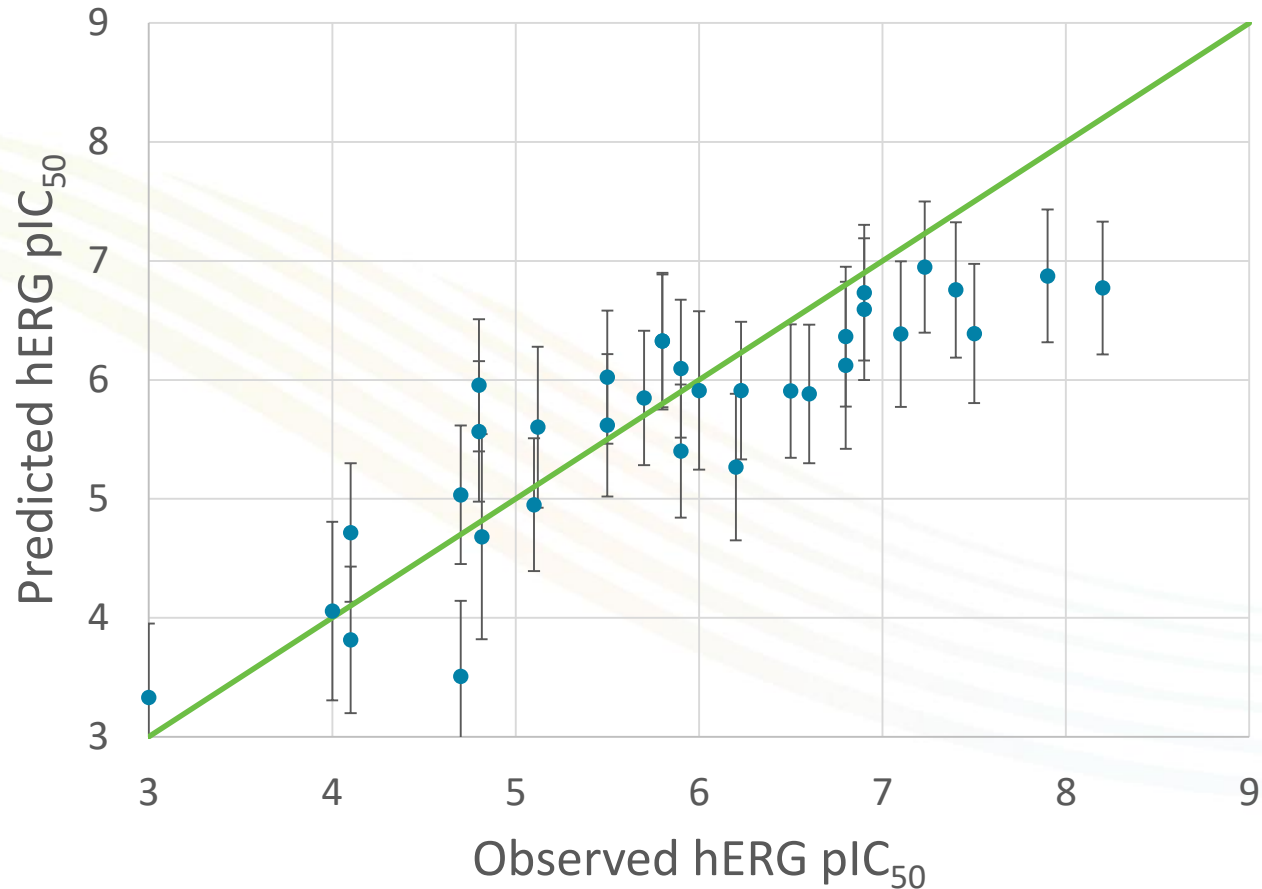
# hERG pIC<sub>50</sub>

---

- Diverse data set of 168 compounds
  - All manual patch clamp measurements in mammalian cells
- Divided into training (135) and external test (33) sets
- Descriptors including
  - Whole molecule properties: logP,  $V_x$ , TPSA, MW, flexibility...
  - 156 structural descriptors expressed as SMARTS

Method	R <sup>2</sup>	RMSE
GP Nest	0.72	0.64
RF	0.68	0.68
RBF	0.70	0.66

# hERG pIC<sub>50</sub> GP Nested Sampling Results

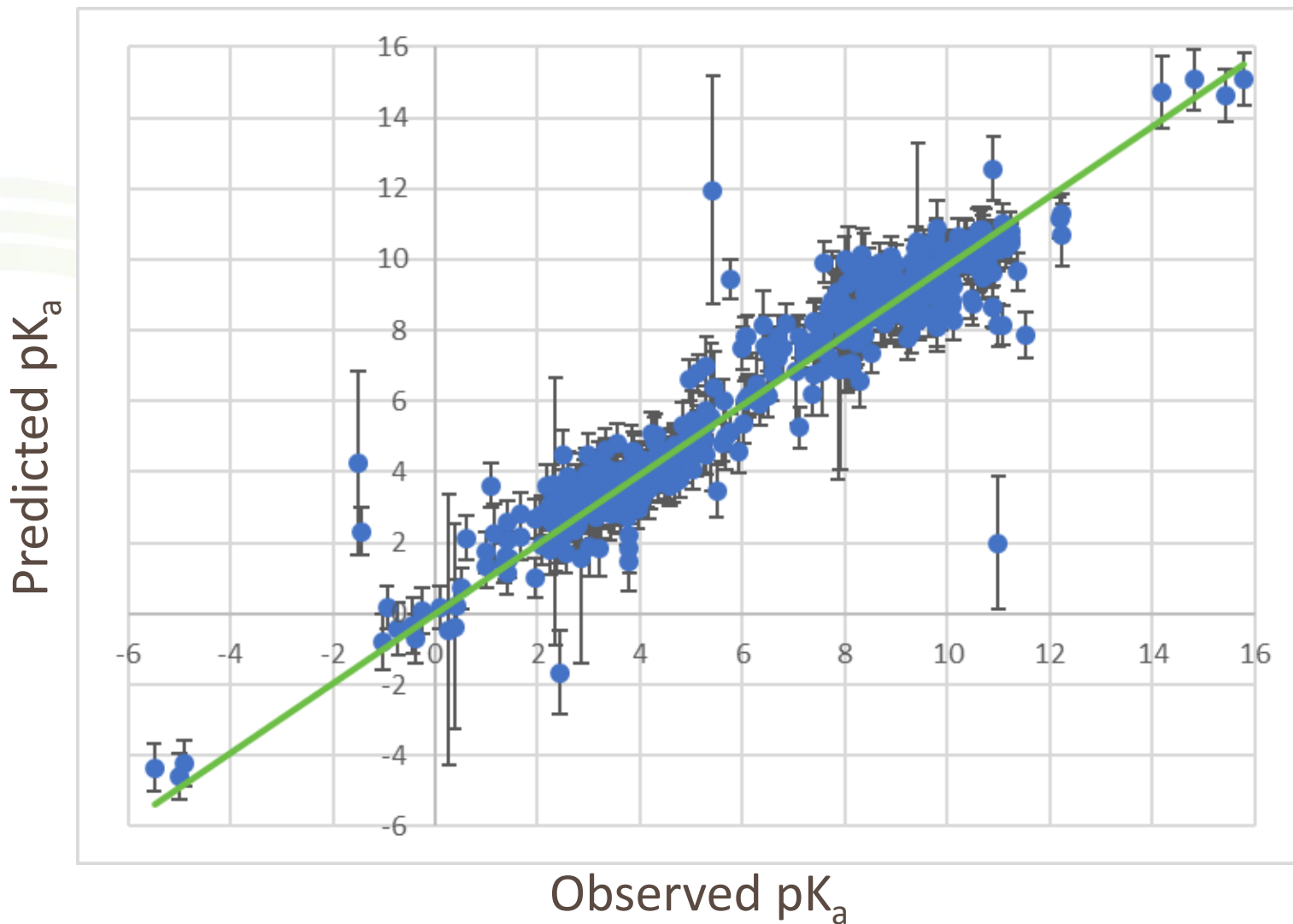


64% within 1 SD, 94% within 2 estimated SD (assuming 0.5 log units uncertainty in expt. data)

- Diverse data set of 1,849 compounds
  - pK<sub>a</sub> range of -5 to +16
  - Range of ionisable groups, preferentially monoprotic
  - Split into 70:30%, training vs test sets
- 35 descriptors derived from quantum mechanical (AM1) calculations

Method	R <sup>2</sup>	RMSE
GP Opt (24 descs)	0.91	0.98
RBF	0.94	0.81
RF	0.92	0.91
GP FIXED	0.88	1.13
PLS	0.76	1.62

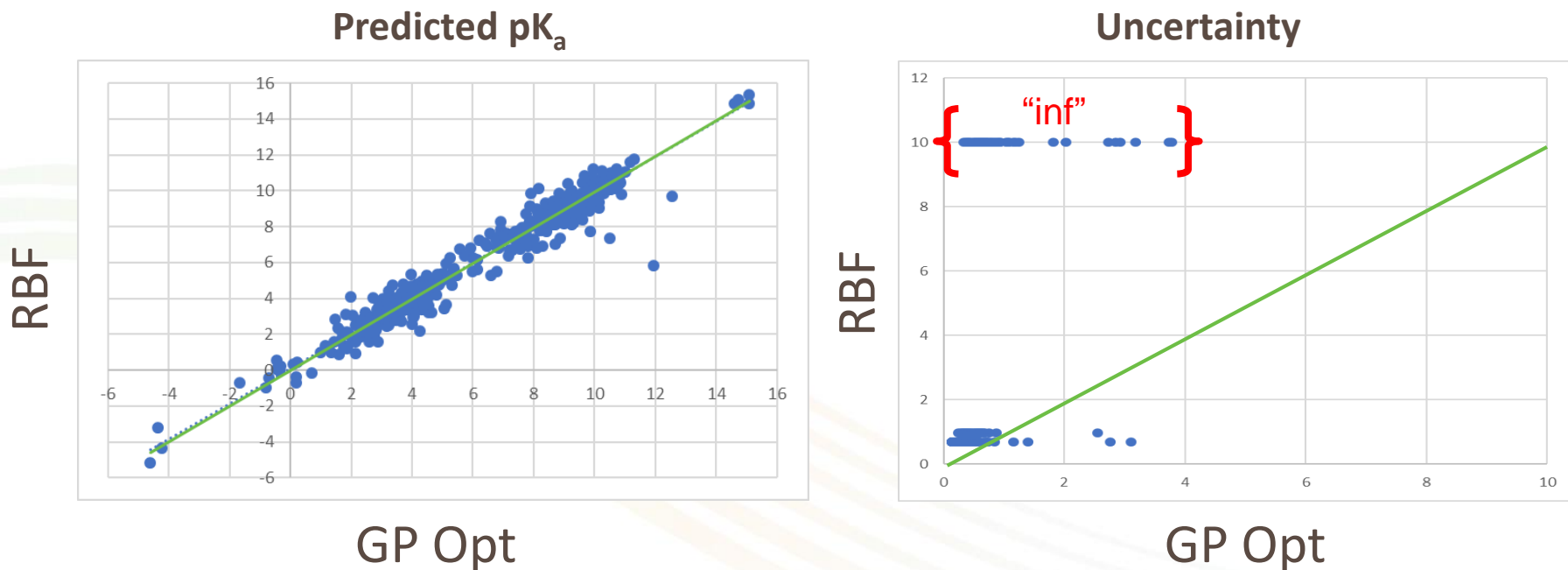
# pK<sub>a</sub> GP Conjugate Gradient Opt. Results



63% within 1 SD, 89% within 2 estimated SD (assuming 0.5 log units uncertainty in expt. data)



# pK<sub>a</sub> GP Opt vs RBF results



- Both models have utility but the GP Opt model has a better estimate of uncertainty
  - Also able to obtain descriptor relevance from the GP Opt model
  - Re-running with the 13 relevant descriptors produces RBF, RF, and GP Opt models with equivalent predictive performance.

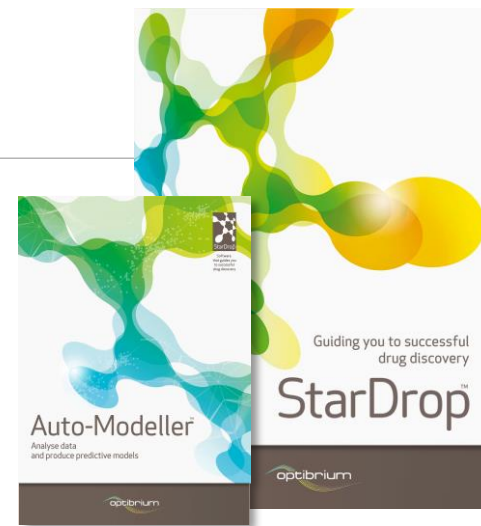
# Conclusions

- Gaussian Processes

- Bayesian non-linear modelling technique
- Generates a probability distribution over possible models
- Explicitly calculates uncertainties for each prediction
- Similar performance to methods such as random forests, radial basis functions...
- Can also be used for classification

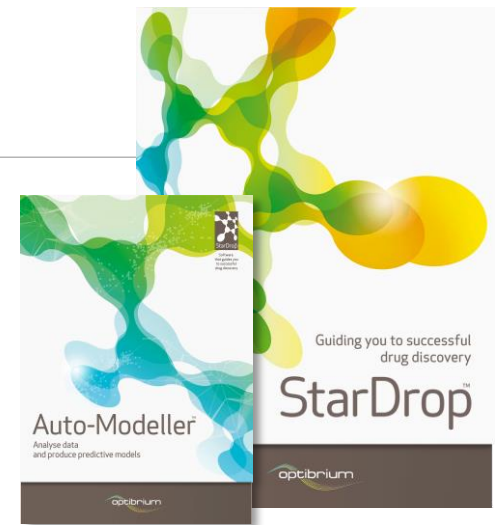
- Limitations

- Most expensive optimisations methods are computationally expensive (e.g. nested sampling  $O(N^4)$ )
- Can't deal with potential sources of variability (e.g. structural features) not captured by descriptors



# Uncertainty

- Rowan Atkinson as Sir Marcus Browning MP



*Life is uncertain.*

*My life certainly has a certain uncertainty about it and I'm certain yours does too.*

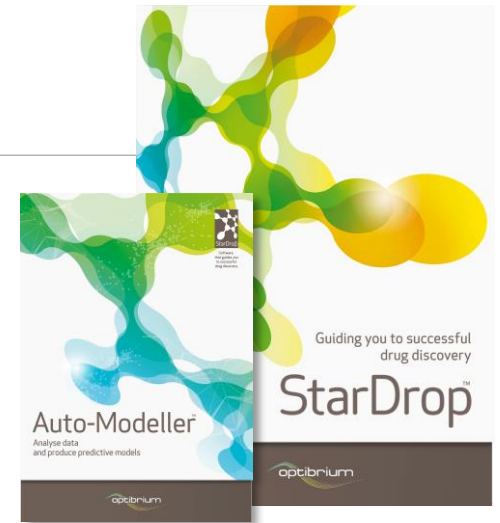
*So with my uncertainty, and your uncertainty, there's certainly a certain degree of uncertainty about.*

*Of that we can be quite cer..., sure.*

# Acknowledgements

---

- Olga Obrezoanova
- Iskander Yusof
- Ed Champness
- All our other colleagues at Optibrium



# References

- Gaussian processes
  - Obrezanova *et al.* JCIM (2007) **47**(5) pp. 1847-57
  - Obrezanova *et al.* JCAMD (2008) **22**(6-7) pp. 431-440
  - Obrezanova *et al.* JCIM (2010) **50**(6), pp. 1053-1061
- Download from [www.optibrium.com/community](http://www.optibrium.com/community)

