# The challenges of making decisions using uncertain data

Matthew D. Segall and Edmund J. Champness

Optibrium Ltd., 7221 Cambridge Research Park, Beach Drive, Cambridge, CB25 9TL, UK

## Abstract

All of the experimental compound data with which we work have significant uncertainties due to variability in experimental conditions and the imperfect correlations between experimental systems and the ultimate *in vivo* properties of compounds. When using these data to make decisions, it is essential that these uncertainties are taken into account to avoid making inappropriate decisions in the selection of compounds, which can lead to wasted effort and missed opportunities. In this paper we will consider approaches to rigorously account for uncertainties when choosing between compounds or selecting compounds against property criteria; first for an individual measurement of a single property and then for multiple measurements of a property for the same compound. We will then explore how uncertainties in multiple properties can be combined when assessing compounds against a profile of criteria, a process known as multi-parameter optimisation. This guides rigorous decision-making using complex, uncertain data to focus on compounds with the best chance of success, while avoiding missed opportunities by inappropriately rejecting compounds.

# Introduction

When working with compound data, we should be aware of the uncertainties in the values obtained from experimental measurements and consider the impact that these have on the decisions that we make based on this information. It is well established that people find it challenging to make good decisions based on uncertain information; experimental psychologists have described many so-called cognitive biases that lead to missed opportunities and inefficient use of resources [1] [2]. In this paper, we will consider ways in which uncertainties in data can affect decisions on the selection and comparison of compounds and discuss approaches to take these uncertainties into account in order to mitigate the associated risks.

There are two main sources of uncertainty in experimental measurements of compound properties:

- We know that variability is observed in the results obtained from an assay when performed multiple times, due to minor changes in the experimental conditions, instrument noise or simply the variability inherent in complex biological systems. This can be considered as noise or statistical error around the 'true' property value.
- There are also uncertainties in the *relevance* of results from experimental systems to the ultimate goal of a project; for example, in drug discovery, we must remember that all experimental systems we use are models of the human patient and these do not correlate perfectly with the *in vivo* behaviour in human.

In this paper, we will consider how we can rigorously take these sources of uncertainty into consideration when using data to make decisions about the selection or design of compounds, for example choosing between compounds or series for further investigation.

We should also remember that identifying a successful compound requires the simultaneous optimisation of multiple compound properties; for a drug discovery objective these include potency against the therapeutic target(s), selectivity over off-targets, appropriate physicochemical, absorption, distribution, metabolism and excretion (ADME) properties and safety. Therefore, we will also explore approaches to deal with the combination of uncertainties in multiple properties in this multi-parameter optimisation (MPO) challenge [3].

# Statistical Uncertainty

If an experiment is repeated multiple times under the same conditions (as far as is possible), the results will vary to some extent. Differences between experimental samples, operators, instruments and the time or location at which the experiment is conducted are among the sources of variation that can give rise to experimental variability. If we consider these variations to be sources of random errors, neglecting for now the possibility of systematic errors that consistently bias a result, we can consider the impact of this uncertainty on our confidence in choosing between two compounds.
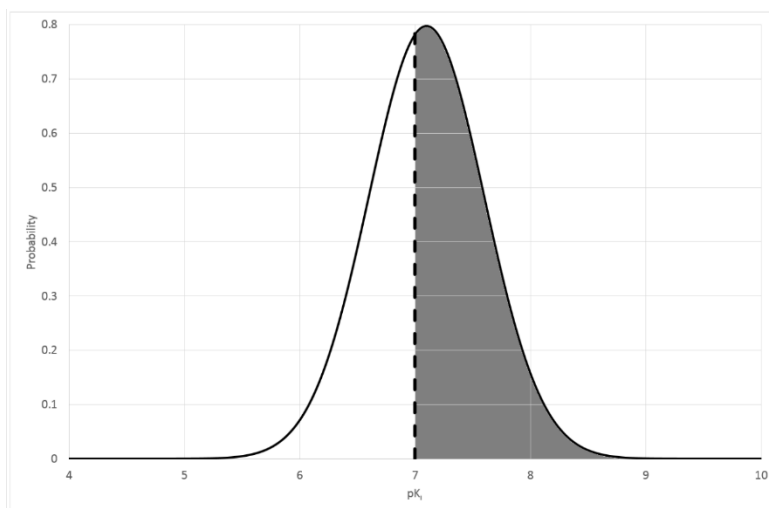
## Single point data

First, we consider the simple case where we have a single measurement, $x_A$, of a property of compound $A$, for which the 'true' (but unknown) value of the property is $X_A$. If we furthermore assume that the experimental error in this measurement is normally distributed with standard deviation $\sigma$ (a variance of $\sigma^2$), then this would be denoted:

$$X_A \sim N(x_A, \sigma^2).$$

We could, for example, estimate the standard deviation from a reference compound that has been run repeatedly through the assay (an estimate of the 'population' standard deviation).

This scenario often occurs when we have early, single-point screening data for a compound and we could ask a question regarding whether the compound meets a selection criterion for the property. For example, if the property of interest is the activity against a target, expressed as a $pK_i$ (the negative log of the inhibition constant $K_i$ in molar concentration), we might ask if a compound with a measured $pK_i$ of 7.1 ($K_i$ = 79 nM) meets a selection criteria of $pK_i > 7$ ($K_i < 100$ nM). If we assumed the data was perfect then clearly the answer is "yes". However, a typical uncertainty in such a value might be 0.5 log units (one standard deviation), which is roughly equivalent to a factor of 3 in the $K_i$ value. Therefore, we *should* ask, "What is the *chance* that this compound meets the

**Figure 1. The probability distribution of the pK$_i$ of a compound, corresponding to a measured value of 7.1 with a standard deviation of 0.5, assuming that the error is normally distributed. The vertical dashed line corresponds to a threshold value of 7 and the shaded region corresponds to the probability that the 'true' pK$_i$ value is greater than 7.**

selection criterion?", as illustrated by the shaded region in Figure 1. Quantitatively, we can calculate the probability that the compound meets our criterion as:

$$P(X_A > 7) = P(N(x_A, \sigma^2) > 7) = P\left(Z > \frac{(7-x_A)}{\sigma}\right) = P(Z > -0.2) = 0.58,$$

where Z takes the standard Normal distribution, $N(0,1)$. Therefore, we can only say that there is a 58% chance that this compound will meet our requirements, little better than a coin toss.

If we had another compound, *B*, for which a single measurement, $x_B$, had been made for the same property, then:

$$X_B \sim N(x_B, \sigma^2).$$

We can then ask questions about the difference between the property values for compounds A and B and there are simple rules for combining the uncertainties. For example, the difference between the properties is:
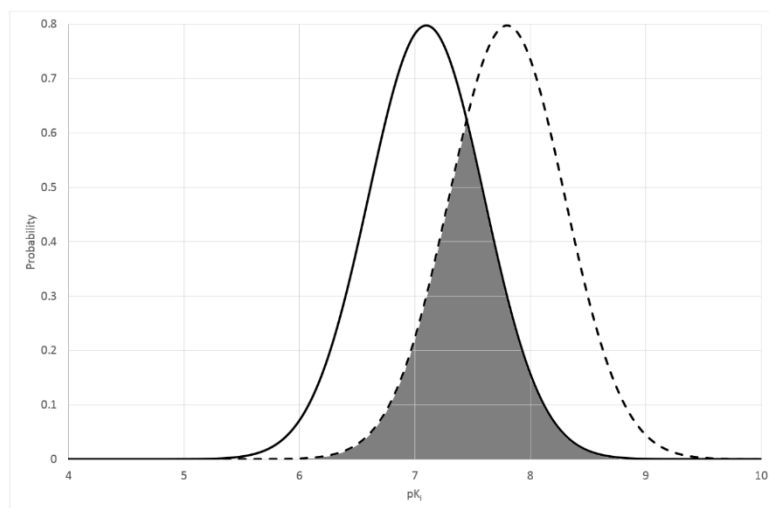
$$(X_B - X_A) \sim N(x_B - x_A, 2\sigma^2),$$

which means that the standard deviation in the difference between the property values of compounds A and B is $\sqrt{2}\sigma$.

For example, we might want to choose the more potent of the two compounds and, if the measured value of the pK$_i$ of compound B was 7.8 (K$_i$ = 16 nM), this might lead us to choose compound B over compound A. However, if this measurement had the same standard deviation, then we can't be absolutely confident that we can distinguish between these compounds, as illustrated in Figure 2. The probability that compound B is more potent than compound A is actually:

$$P(X_B - X_A > 0) = P(N(x_B - x_A, 2\sigma^2) > 0) = P\left(Z > \frac{x_A - x_B}{\sqrt{2}\sigma}\right) = P(Z > -0.99) = 0.84.$$

So, again, we can only say that there is an 84% chance that compound B is more potent than compound A. All other things being equal, we would still probably place our bets on compounds B, but would we want to take the risk of missing an opportunity in compound A?

**Figure 2.** Probability distributions of the $pK_i$ values of two compounds: Compound A (solid line) has a measured value of 7.1 and compound B (dashed line) has a measured value of 7.8. Both measurements have a standard deviations of 0.5 and we have assumed that the errors are normally distributed. The shaded region highlights the region in which there is a significant probability that the 'true' $pK_i$ of compound A is higher than that of compound B.

## Multiple measurements

Given the inherent variability in biological systems, it is common to compare the average (or mean) property values of a compound. Therefore, as a project progresses, experiments will be often performed in replicate, to generate several data points for the same compound. If we had a very large number of replicates we could obtain a precise estimate of the 'true' mean, $\bar{X}$. However, in practice, only a handful of measurements may be made, which limits the accuracy of our estimate of this mean, $\bar{x}$, made from the sample. However, the accuracy of the estimate, the standard error in the mean, $SE_{\bar{x}}$, made from the limited sample, can be estimated from the standard deviation of the sample, $s$, as follows:

$$\bar{x} = \frac{1}{N}\sum_{i=1}^{N} x_i,$$

$$s = \sqrt{\frac{1}{N-1}\sum_{i=1}^{N}(x_i - \bar{x})^2},$$

$$SE_{\bar{x}} = \frac{s}{\sqrt{N}},$$

where, $N$ is the number of measurements and $x_i$ is the $i$'th measurement of the property.

In these scenarios, the mean, $\bar{x}$, takes the Student's t distribution [4] with $N-1$ degrees of freedom and this enables us to rigorously estimate probabilities in a similar manner to the single-measurement case above. For example, we may have two compounds C and D, for which the measurements in Table 1 have been made.

**Table 1 Example samples of data for two compound and calculated sample statistics.**

| Compound | Measurements | | | | | Sample Mean ($\bar{x}$) | Sample Standard Deviation (s) | Standard Error in Mean ($SE_{\bar{x}}$) |
|---|---|---|---|---|---|---|---|---|
| C | 1.5 | 3.2 | 2.4 | 3.5 | 4 | 2.92 | 0.98 | 0.44 |
| D | 3.5 | 4.3 | 5.5 | 4.9 | 3.2 | 4.28 | 0.95 | 0.43 |

We could then ask the question, what is the chance that the mean property value for compound C, $\bar{X}_C$, is greater than 3? Which may be calculated as follows:

$$P(\bar{X}_C > 3) = P\left(t_4 > \frac{\bar{x}_C - 3}{SE_{\bar{x}_C}}\right) = P(t_4 > -0.18) = 0.43,$$

where $t_4$ takes the Student's t distribution with 4 degrees of freedom. In this case, we might have been tempted to reject compound C, based on a measured average of 2.92 and a selection criterion of >3, when there is actually a 43% chance that the true average value meets this criterion.

When comparing two compounds, based on their average measured values, the formulae become even more complex [4], but, for example, we can ask if compound D has a higher mean property than compound C, as follows:

$$P(\bar{X}_D - \bar{X}_C > 0) = P\left(t > \frac{\bar{x}_C - \bar{x}_D}{\sqrt{\frac{s_C^2 + s_D^2}{N}}}\right) = P(t_8 > -2.22) = 0.97.$$

Therefore, in this case, we can be 97% confident that compound D has a higher average property value than compound C.

### Combining measurements

When combining data for different properties, to calculate a derived value such as target selectivity, the uncertainties in the individual measurements also combine. For example, ligand efficiency indices are currently popular metrics for comparing the 'quality' of compounds [5] and the ligand lipophilicity efficiency is defined as:

$$LLE = pK_i - \log P,$$

where logP is the logarithm of the octanol:water partition coefficient. However, as we've seen there will be experimental error in the $pK_i$ value and, similarly, there will be uncertainty in the logP value, particularly if a predicted value is used. Therefore, assuming both errors are normally distributed, the standard deviation in the LLE will be given by:

$$\sigma_{LLE} = \sqrt{\sigma_{pK_i}^2 + \sigma_{\log P}^2},$$
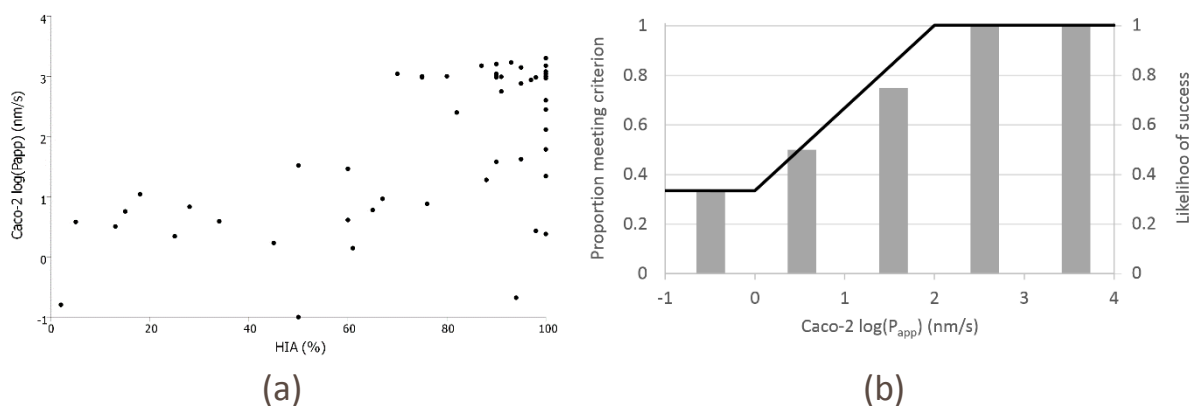
where $\sigma_{pK_i}$ is the standard deviation in the $pK_i$ value and $\sigma_{\log P}$ is the standard deviation in the logP.

It is important to keep this in mind because uncertainties can accumulate quickly and make it difficult to distinguish between compounds based on these derived properties.

## Relevance

Even if we knew the results of an experimental measurement precisely, this would not necessarily mean that we could make a confident decision. This is because the properties that are commonly measured early in a project are often models of the ultimate objective; in particular, in drug discovery, all experimental systems, whether *in vitro* or *in vivo,* are models of the human patient. These models do not correlate exactly with the *in vivo* outcome in human and therefore it may not be appropriate to apply hard criteria when selecting compounds, because this may lead to rejecting good compounds inappropriately.

Take, for example, the data in Figure 3(a) comparing permeability across the human epithelial colorectal adenocarcinoma (Caco-2) cell line [6], a commonly used model of permeation across the human intestine, with clinically measured human intestinal absorption (HIA), as published by Irvine *et* al. [7]. Here we can see that a high measured Caco-2 permeability would give us confidence that the compound would be well absorbed, but a low permeation does not strongly indicate poor absorption, although we could say that the *chance* of achieving good oral absorption in humans would be lower. Therefore, it would not be appropriate to reject a compound outright based on a low Caco-2 permeability, particularly if other properties of the compound were good.

**Figure 3. (a) scatter plot of experimentally measured Caco-2 P_app against clinical human intestinal absorption for 52 compounds published by Irvine et al. [7]. The histogram in (b) shows the proportion of compounds achieving a human intestinal absorption greater than 50% for Caco-2 P_app values binned in one log-unit ranges. The solid line corresponds to a desirability function approximately representing the likelihood of success of compounds for this objective against Caco-2 P_app.**

One approach to avoiding hard cut-offs is to use a 'desirability function' [8] that relates the value of a measurement to its desirability, on a scale between 1 (ideal) to 0 (reject absolutely). These can reflect the impact of a property value on the chance of success of a compound to give a measurement appropriate weight in a decision. This is important because we know that some property criteria are critical, while it may be appropriate to compromise or trade-off other properties to achieve better results for critical factors.

An example of this is shown in Figure 3(b) for the objective of achieving a HIA greater than 50%. The histogram bars indicate the chance of a compound achieving this objective for measured Caco-2 permeability in one log-unit ranges, according to the data in Figure 3(a). Here we can see that, even for compounds with the lowest measured Caco-2 permeability, one third have a clinical HIA greater than 50%. The corresponding desirability function is shown, indicating that the ideal outcome would be a Caco-2 permeability above 100 nm/s (log(P_app) > 2) while the worst outcome would be a Caco-2 permeability below 1 nm/S (log(P_app) < 0), where the chance of success is still approximately one third. Between these values, the desirability increases approximately linearly.

## Combining Multiple Properties, Relevance and Uncertainty

So far, we have explored approaches to considering uncertainties in data for a single property. However, when optimising a compound against a profile of property criteria, we should also consider how the uncertainties in the data combine to affect our ability to distinguish between compounds.

If we consider a naïve approach of applying a series of hard cut-offs, or filters, we can see the issues that can arise. For example, if we apply filters for 5 different properties that are each 80% accurate in distinguishing 'good' from 'bad' outcomes, the probability of an ideal compound passing all 5 filters is only 33%, i.e. we are twice as likely to reject a perfect compound than to take it forward. Given that perfect compounds are usually rare, the opportunity cost of these errors can be high.

Using desirability functions softens the impact of hard cut-offs and the desirabilities of multiple properties can be combined to calculate a 'desirability index', representing the overall quality of a compound against a required property profile. The most common approaches for combining the individual property desirabilities use additive or multiplicative approaches:

Additive: $D(x_1, x_2, \ldots, x_N) = \sum_{i=1}^{N} d_i(x_i)$

Multiplicative: $D(x_1, x_2, \ldots, x_M) = \prod_{i=1}^{N} d_i(x_i)$

where $x_i$ are the values of N compound properties and $d_i$ are the desirability functions for the properties. These are sometimes normalised by the number of properties by taking the arithmetic or geometric mean, for the additive or multiplicative approaches respectively, and the individual desirabilities can be weighted to reflect different degrees of importance of each property. The relative strengths and weaknesses of these approaches and some other alternatives are discussed in more detail in [9].
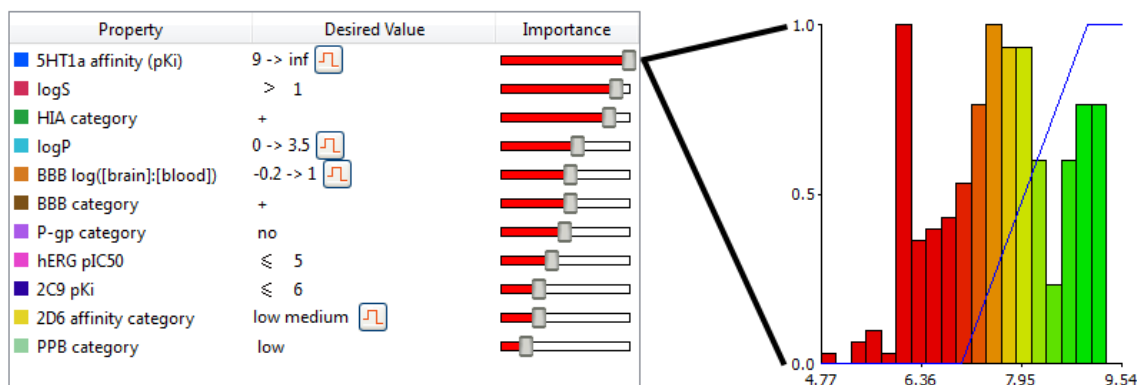
**Figure 4. An example of a profile of property criteria suitable for identifying a compound that is a potent inhibitor of the serotonin 5-Hydroxytryptamine (5-HT$_{1A}$) receptor and has suitable physicochemical and ADME properties for oral dosing and a target in the central nervous system. Underlying each of the criteria are desirability functions, as illustrated for the pK$_i$ against 5-HT$_{1A}$. The histogram behind the desirability function shows the distribution of pK$_i$ values for the compounds an example data set.**

The Probabilistic Scoring method [10] builds on desirability functions to explicitly account for the uncertainties in the underlying data. Using this approach, a profile of property criteria can be defined, as illustrated in Figure 4, to reflect the requirements of a specific project. Underlying each of these criteria is a desirability function that reflects the importance and acceptable trade-offs for each property. By assessing the data and associated uncertainties, a score is calculated that represents the chance of success against the desired profile, i.e. the probability of achieving the ideal property criteria. Furthermore, the uncertainty in the overall score can be calculated, which indicates when compounds can be confidently distinguished or, conversely, when the data do not support this decision, as illustrated in Figure 5. This helps to avoid missed opportunities caused by giving too much weight to uncertain data. Furthermore, the impact of missing data, where a property of a compound has not yet been measured, can be accounted for rigorously. The impact of the missing data on the priority given to the compound can be assessed to identify when it would be valuable to 'fill in' the missing data point.
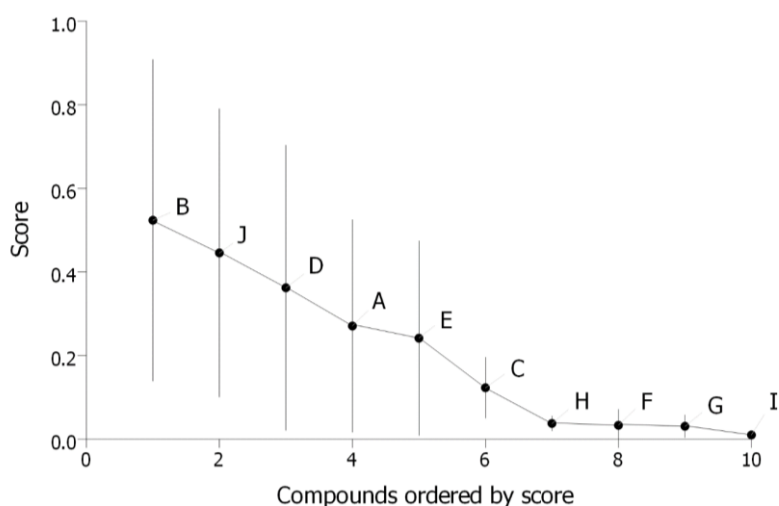


**Figure 5. The results of probabilistic scoring for the 10 compounds and associated data shown in Figure 6. The compounds are ordered from left to right along the x-axis in order of their score and the overall score for each compound is plotted on the y-axis. The uncertainty in each score (one standard deviation), due to the uncertainty in the underlying data, is shown by error bars around the corresponding point. From this it can be seen that compounds B, J, D, A and E cannot be confidently distinguished based on the available data, while H, F, G and I can be confidently rejected. The probability that compound C is equivalent to the highest scoring compound is small although the difference is not statistically significant.**

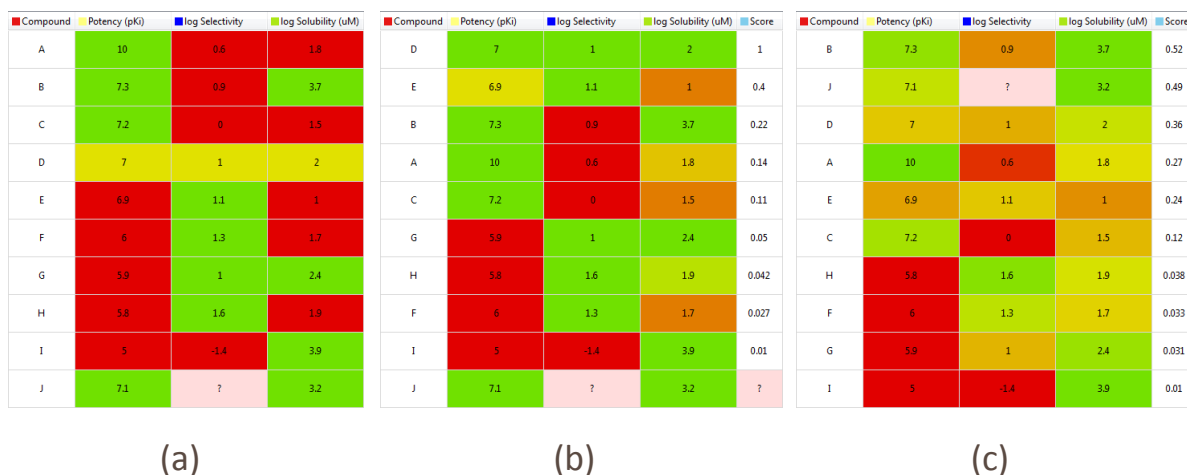|     |     |     |
| --- | --- | --- |
| (a) | (b) | (c) |

**Figure 6. A simple, hypothetical example of prioritisation of 10 compounds (labelled A through J) with data for potency (pKi), log selectivity and log solubility (µM) using three methods.**

a) The results of applying filters corresponding to pKi > 7, log selectivity > 1 and log solubility > 2. A green cell indicates that the property passes the criterion and red that it fails. Compound D is coloured yellow because it lies exactly on the thresholds for all three properties.

b) The results of calculating a score corresponding to a multiplicative desirability index using the desirability functions for the three properties shown in Figure 7. The compounds are sorted by score and the cells are coloured by the desirability of each property value from red (0) to green (1).

c) The results of applying Probabilistic Scoring using desirability functions shown in Figure 7 and the following uncertainties (1 standard deviation): pKi ± 0.3; log selectivity ± 0.4; log solubility ± 0.6 log units. The compounds are ordered by score and the cells are coloured by the likelihood of achieving the ideal outcome for the corresponding property from red (0) to green (1).

As an example, consider the simple, hypothetical data set shown in Figure 6, showing values for potency, selectivity and solubility for 10 compounds labelled A through J. In Figure 6(a) the results are shown for filtering the compounds based on the following cut-offs:

- Potency (pK$_i$) > 7 (better than 100 nM)
- Log selectivity > 1 (better than a factor of 10)
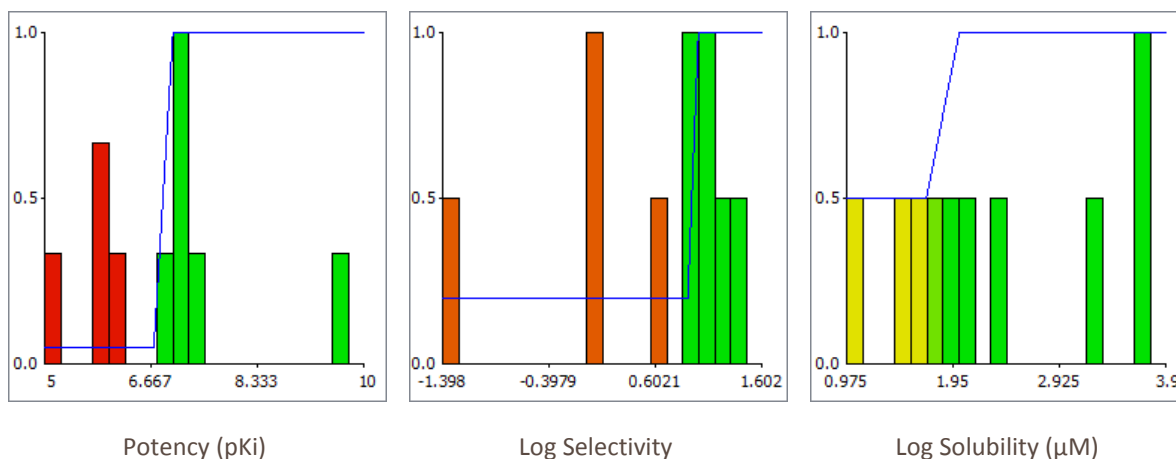- Log solubility (µM) > 2 (better than 100 µM)

Note that all of the compounds fail on one or more criteria, except for one which is on the threshold for all three properties. Also note that the selectivity value for compound J is missing, so it is not known if this compound would pass all of the criteria. Therefore, if we were to choose any compound on this basis, it would be compound D, but none of the compounds clearly meet all of the criteria.

If we apply the desirability functions shown in Figure 7 to the property values of these compounds, using a multiplicative scheme to calculate the overall scores, the results are shown in Figure 6(b). Here we see that compound D would still be ranked highest, but the remaining compounds can also be prioritised according to the importance of each property and the desirability of property values close to the ideal cut-offs. It is still not clear how to prioritise compound J due to the missing data for selectivity.

However, there are uncertainties in the property values, as follows:

- Potency (pK$_i$): ±0.3 log units (a factor of 2 in the K$_i$)
- Log Selectivity: ±0.4 log units (a factor of 2.6, derived from the ratio of two potencies each with a factor of 2 uncertainty)
- Log Solubility: ±0.6 log units

Potency (pKi)          Log Selectivity          Log Solubility (µM)

**Figure 7. Desirability functions for potency (pKi), log selectivity and log solubility (µM), as applied to the example compounds in Figure 6. The desirability function for potency corresponds to an ideal $pK_i$ greater than 7 and a linearly increasing likelihood of success from a minimum of 0.05 for $pK_i$ values less than 6.7. The desirability function for log selectivity corresponds to an ideal value greater than 1 and a linearly increasing likelihood of success from a minimum of 0.3 for values less than 1.7. The desirability function for log solubility corresponds to an ideal value greater than 2 and a linearly increasing likelihood of success from a minimum of 0.3 for values less than 0.9. The histograms in each case shows the distribution of the corresponding property for the data set in Figure 6.**

Therefore, applying Probabilistic Scoring to these compounds, using the same desirability functions shown in Figure 7 and taking into account these uncertainties, gives rise to the scores shown in Figure 6(c). Now we can see that compound B has the best overall chance of success because it lies confidently above the ideal criteria for potency and solubility and close to the cut-off for selectivity. Also note that compound J can now be ranked alongside the other compounds because the missing data can be rigorously considered as a very uncertain value and the impact of this uncertainty assessed. Due to the good values compound J achieved for potency and solubility it ranks higher than compounds which fail these property criteria with confidence and those that confidently fail the criterion for selectivity; it is better to have an uncertain result than a value that is known to be poor.

The overall impact of the uncertainties is shown in Figure 5, where the error bars indicate the uncertainties in the overall score. Here we can see that, in fact, the top 5 or 6 compounds cannot be confidently distinguished from the highest ranked compound based on the available data. Only compounds F, G, H and I can be rejected with confidence, while the probability that C is equivalent to the highest scoring compound is small, although the difference is not statistically significant.

## Conclusion

We have explored approaches to account for the uncertainties in compound data and the impact these have on decisions regarding the selection and prioritisation of compounds. Neglecting uncertainties can lead to poor decisions, resulting in wasted time and effort and missed opportunities. The last of these is possibly the most insidious because, once rejected, it is rare to return to a compound or series, so the lost value is unlikely to be discovered.

Further analysis of the impact of uncertainty on decisions can yield answers to strategic questions regarding the value of different sources of data to decision-making, in light of the confidence they provide [11]. This analysis uses Bayesian probability theory [12] which requires a knowledge of the prior probability distribution, or underlying distribution of the property, in question. However, the scope of this analysis is currently limited because priors for the most prevalent risk factors for compound optimisation are not generally known.

Finally, as we have seen, the mathematics involved in assessing the impact of uncertainties can be quite daunting, which leads to the temptation to ignore uncertainty and hope for the best! Therefore, it essential that chemistry software can automatically propagate uncertainties through data analyses and present the results in an intuitive way to guide effective decisions on compound optimisation.

# References

1   Tversky, A, and Kahneman, D. (1974) Judgment under Uncertainty: Heuristics and Biases. Science 185:1124-1131.

2   Chadwick, A.T. and Segall, M.D. (2010) Overcoming psychological barriers to good discovery decisions. Drug Discov. Today 15:561-569.

3   Segall, M.D. (2012) Multi-Parameter Optimization: Identifying high quality compounds with a balance of properties. Curr. Pharm. Des. 18:1292-1310.

4   Rice, J (2007) Mathematical Statistics and Data Analysis (Third ed.). Duxbury Press, Belmont, CA.

5   Abad-Zapatero, c. and Metz, J.M. (2005) Ligand Efficiency Indices as Guideposts for Drug Discovery. Drug. Discov. Today 10:464-469.

6   Hidalgo, I.J., Raub, , T.J., and Borchardt, R.T. (1989) Characterization of the human colon carcinoma cell line (Caco-2) as a model system for intestinal epithelial permeability. Gastroenterology 96:736-749.

7   Irvine, J.D., Takahashi, L., Lockhart, K., Cheong, J., Tolan, J.W., Selick, H.E., and Grove, J.R. (1999) MDCK (Madin–Darby Canine Kidney) Cells: A Tool for Membrane Permeability Screening. J. Pharm. Sci. 88:28-33.

8   Harrington, E.C. (1965) The desirability function. Ind. Qual. Control 21:494-498.

9   Segall, M.D. (2014) Advances in multiparameter optimization methods for de novo drug design. Expert Opin. Drug Discov. 9:803-817.

10  Segall, MD, Beresford, AP, Gola, JMR, Hawksley, D, and Tarbit, MH. (2006) Focus on Success: Using in silico optimisation to achieve an optimal balance of properties. Expert Opin. Drug Metab. Toxicol. 2:325-337.

11  Segall, M.D. and Chadwick, A. (2010) Making Priors a Priority. J. Comp.-Aided Mol. Des. 24:957-960.

12  Jaynes, E.T. (2003) Probability Theory: The Logic of Science: Principles and Elementary Applications Vol. 1. Cambridge University Press, Cambridge.