

Today: -handling outliers  
-Markov Chain Monte Carlo

Note: outlier material follows arxiv:1008.4686  
Section 3 (pg 11). More explanation &  
exercises there too.

## Outliers

Experimental science always encounters  
glitches. How do we deal with them?

- by hand?
- by some heuristic? (sigma-clipping)
- with statistics!

↳ So far, we have written a generative model  
for our data: eg,  $y_i = b + mx_i + e_i$

Now, we need to model how data can be bad!

We'll say: data can come from our good  
("foreground") or bad ("background")  
distribution. For the background, we'll use  
a Gaussian (whose parameters we'll fit).

- two ways of looking at this setup.  
equivalent

① For  $N$  data points, add  $N$  new parameters,

$$z_i, \quad z_i = \begin{cases} 0 & \text{if data point } i \text{ is bad} \\ 1 & \text{if data point } i \text{ is good} \end{cases}$$

Our "good" distribution is

$$P_{fg} \quad y_i \sim N(mx_i + b, \sigma_i^2)$$

our "bad" dist. is

$$P_{bg} \quad y_i \sim N(Y, V)$$

So we can write the combo as:

$$p(y_i | m, b, z_i, Y, V) = P_{fg}(y_i | m, b)^{z_i} \times P_{bg}(y_i | Y, V)^{(1-z_i)}$$

Now, we're going to marginalize out the parameter  $z_i$ :

$$\int p(y_i | m, b, z_i, Y, V) p(z_i) dz_i$$

$z_i$  only takes value 0, 1, so

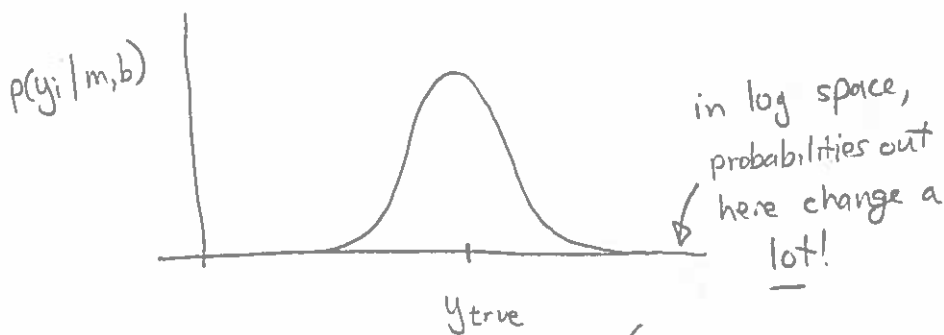
$$= \sum p(y_i | \dots z_i=0 \dots) p(z_i=0) + p(y_i | \dots z_i=1 \dots) p(z_i=1)$$

Call  $p(q_i=0)$  " $P_{\text{bad}}$ " - a priori probability that any given data point is bad.

Then  $p(q_i=1) = 1 - P_{\text{bad}}$ , and

$$p(y_i | m, b, Y, V, P_{\text{bad}}) = P_{\text{bad}} \times P_{\text{bg}}(y_i | Y, V) + (1 - P_{\text{bad}}) \times P_{\text{fg}}(y_i | m, b)$$

This is the "mixture model" or " $\text{fg-bg}$ " model version. We got rid of the  $N$  params  $q_i$ .



Ratio between 10 $\sigma$  and 9 $\sigma$  is way bigger than between 2 $\sigma$  and 1 $\sigma$   
 $\rightarrow$  Fit is dominated by the outliers.

