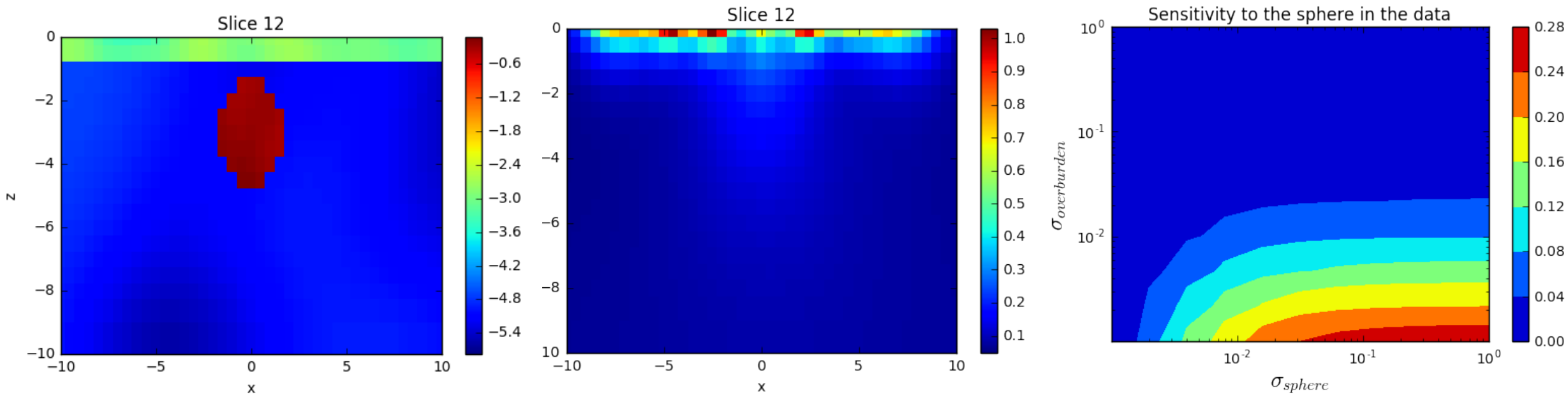


Designing objective functions: a probability approach

Thibaut Astic

SimPEG Meeting 28/02/2017

Motivation



Basics

- Consider the data d as a realization of a multivariate probability distribution

$$d_{obs} \sim P(d|A, m, \Sigma)$$

- m : the model
- A : the physics operator
- Σ : The covariance matrix of your data

Maximum Likelihood Estimation

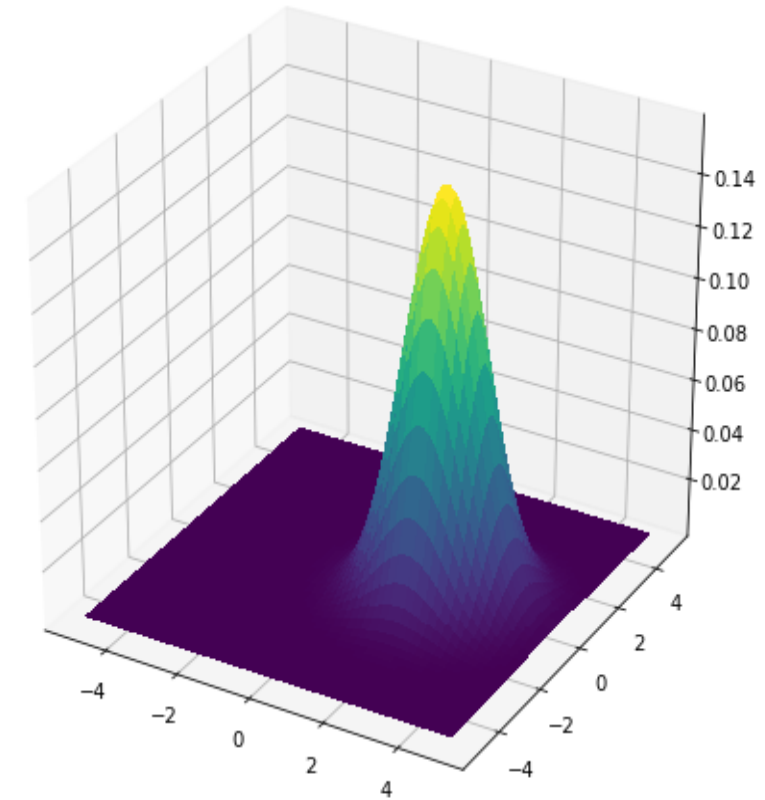
- Uncorrelated Gaussian Example

$$d_{obs} \sim P(d|A, m, \Sigma)$$

- Assume:

- P: multivariate Gaussian distribution
- Mean is $A*m$
- Covariance: uncorrelated noise:

$$\Sigma = \text{diag}(\sigma_i^2)$$



$$P \sim \mathcal{N}$$

$$P(d|A, m, \Sigma) \propto \exp\left(-\frac{1}{2}(Am - d)^T \Sigma^{-1} (Am - d)\right)$$

Maximum Likelihood Estimation

- Uncorrelated Gaussian Example

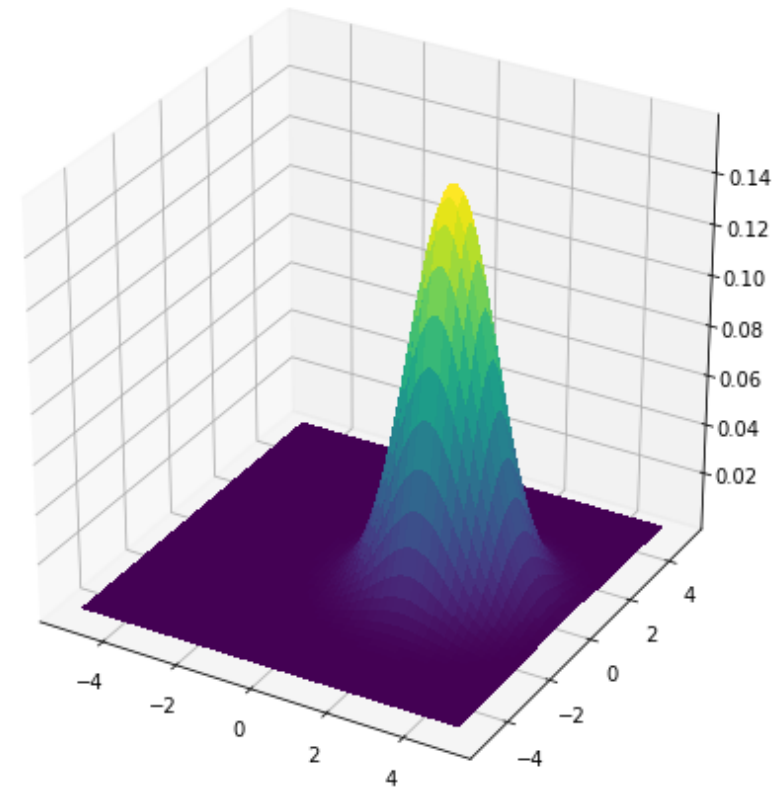
$$d_{obs} \sim P(d|A, m, \Sigma)$$

- Assume:

- P: multivariate Gaussian distribution
- Mean is $A*m$
- Covariance: uncorrelated noise:

$$\Sigma = \text{diag}(\sigma_i^2)$$

$$P(d|A, m, \Sigma) \propto \exp\left(-\frac{1}{2}(Am - d)^T \Sigma^{-1}(Am - d)\right)$$



$P \sim \mathcal{N}$
Start to look familiar



Maximum Likelihood Estimation

- Uncorrelated Gaussian Example
 - Given d , we typically want to maximize the likelihood of d by playing on the model m . This is equivalent to minimize the **negative log-likelihood**,

$$\operatorname{argmax}_m P(d|A, m, \Sigma) = \operatorname{argmin}_m -\log(P(d|A, m, \Sigma))$$

- which then give the familiar:

$$\begin{aligned} m_{opt} &= \operatorname{argmin}_m \frac{1}{2} (Am - d)^T \Sigma^{-1} (Am - d) \\ &= \operatorname{argmin}_m \frac{1}{2} \|Am - d\|_{\Sigma^{-1}}^2 \end{aligned}$$

Maximum Likelihood Estimation

- Uncorrelated Gaussian Example

$$\begin{aligned} m_{opt} &= \operatorname{argmin}_m \frac{1}{2} (Am - d)^T \Sigma^{-1} (Am - d) \\ &= \operatorname{argmin}_m \frac{1}{2} \|Am - d\|_{\Sigma^{-1}}^2 \end{aligned}$$

- But all models m are here equally good!
- The real question is then not (only) to maximize the likelihood of d , but:
- Given d , what is the most likely model? $P(m|A, d, \Sigma)$

Maximum A Posteriori Estimation

- How likely is my model?

- Using Bayes rule:

$$P(m|A, d, \Sigma) = \frac{\overbrace{P(d|A, m, \Sigma)}^{\text{Data Fitting}} \overbrace{P(m)}^{\text{Prior/Regularizer}}}{P(d)}$$

- Need to choose a Prior:

- For example: Multivariate Gaussian Prior

$$m \sim \mathcal{N}(m_{ref}, \frac{1}{\beta})$$

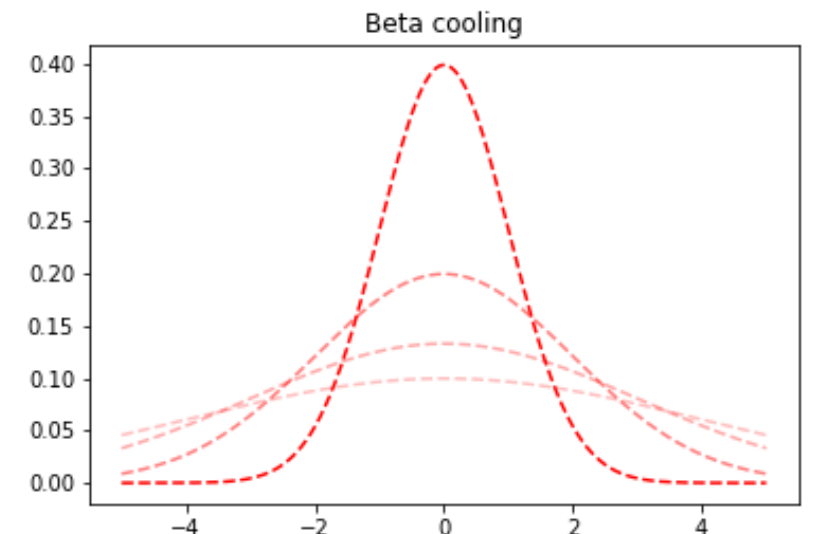
Normalization term

Maximum A Posteriori

- Uncorrelated Gaussian
 - Building on our previous example, we get the well-known Tikhonov regularization:

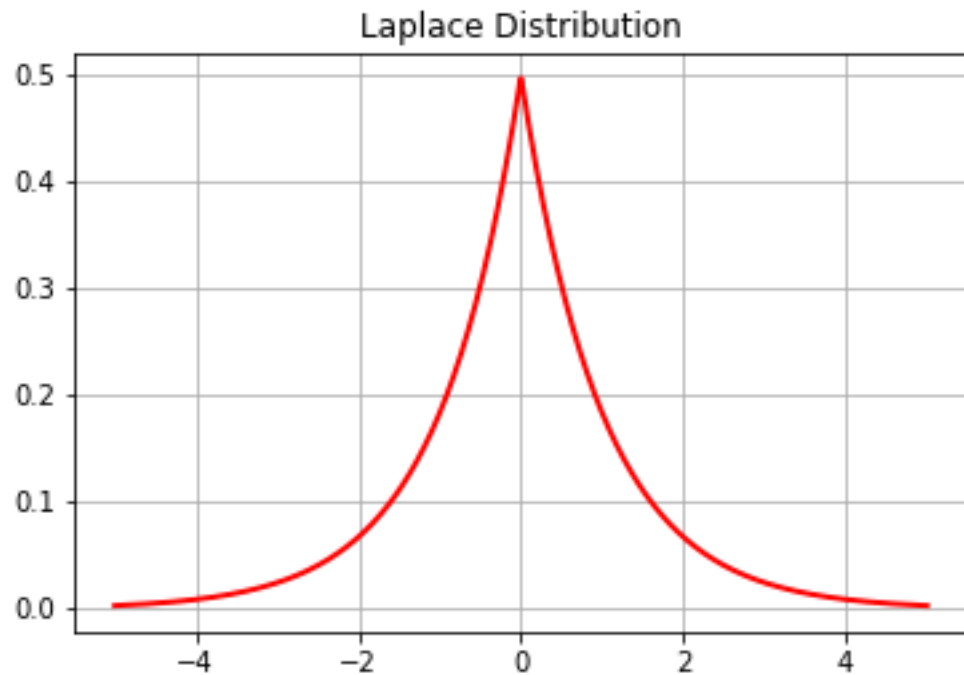
$$\operatorname{argmin}_m -\log(P(m|A, d, \Sigma, m_{ref}, \beta)) = \frac{1}{2} \|Am - d\|_{\Sigma^{-1}}^2 + \frac{\beta}{2} \|m - m_{ref}\|_2^2$$

- Take home:
 - Everything on the right side of 'P(m |...)' are our assumptions
 - L2-norm assumes a Gaussian Distribution of the parameters values
 - Beta-cooling: the smaller is Beta, the more likely become values far from the mean



What can we do with it?

- Use a different probability distribution



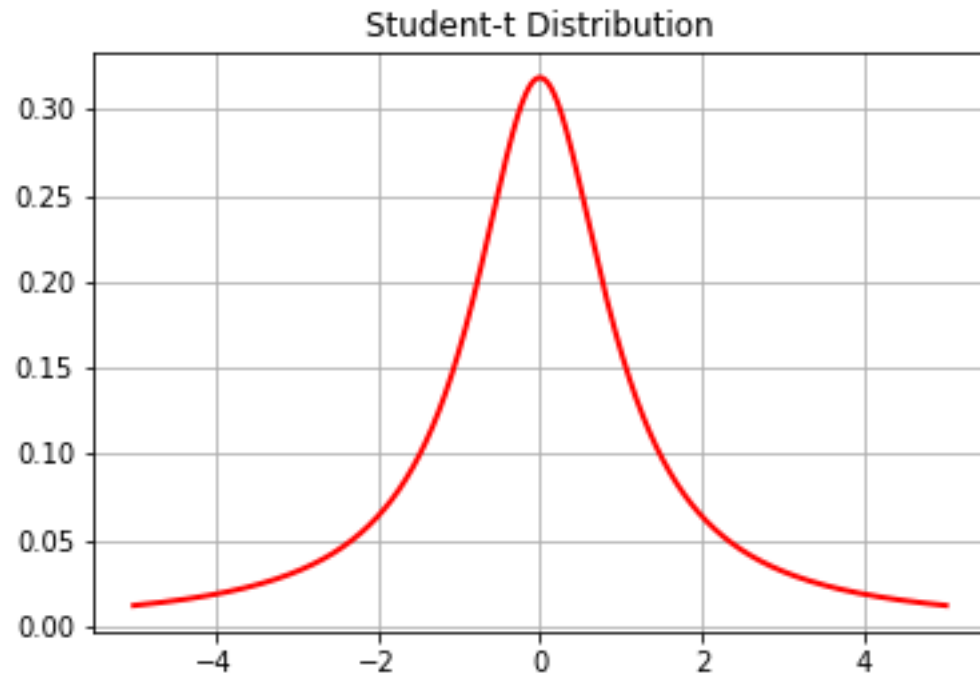
L1-Norm

Example: Regularizer: $||m||_1$

Very useful if your model is sparse
in some basis (Space, Fourier...)

What can we do with it?

- Use a different probability distribution



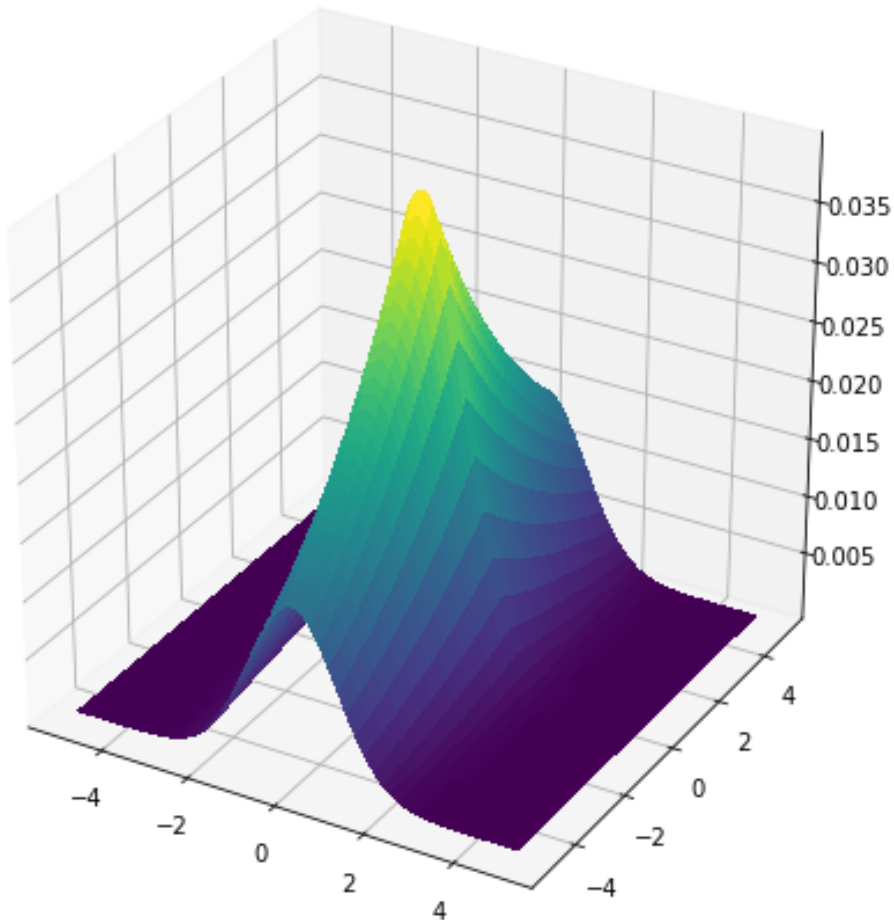
Very robust to outliers.

Example: Data fitting

$$\operatorname{argmin}_m \frac{v+1}{2} \sum_{i=1}^n \log\left(1 + \frac{(m^T a_i - d_i)^2}{v}\right)$$

What can we do with it?

- Compose different distributions to promote structures



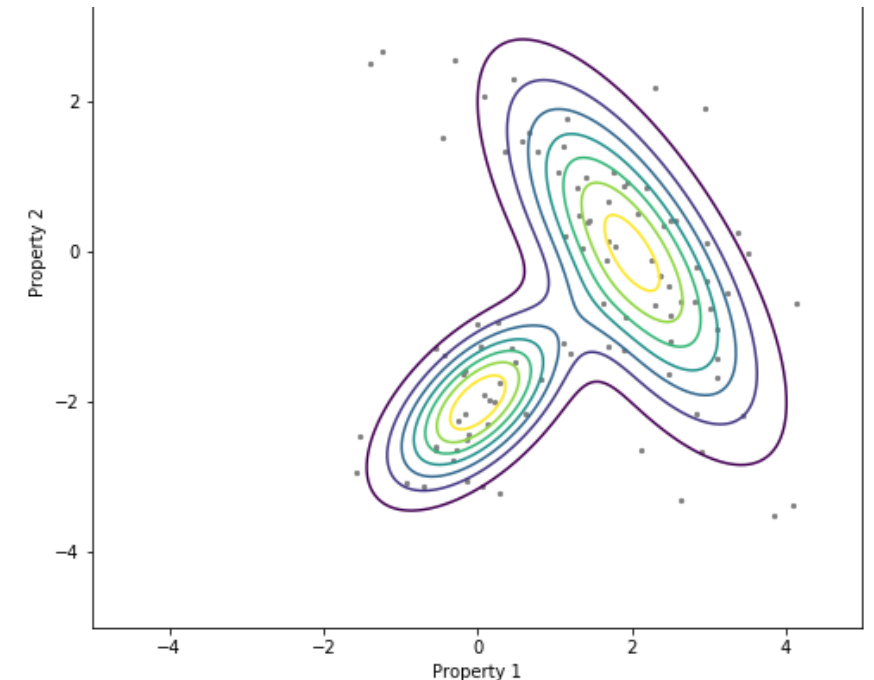
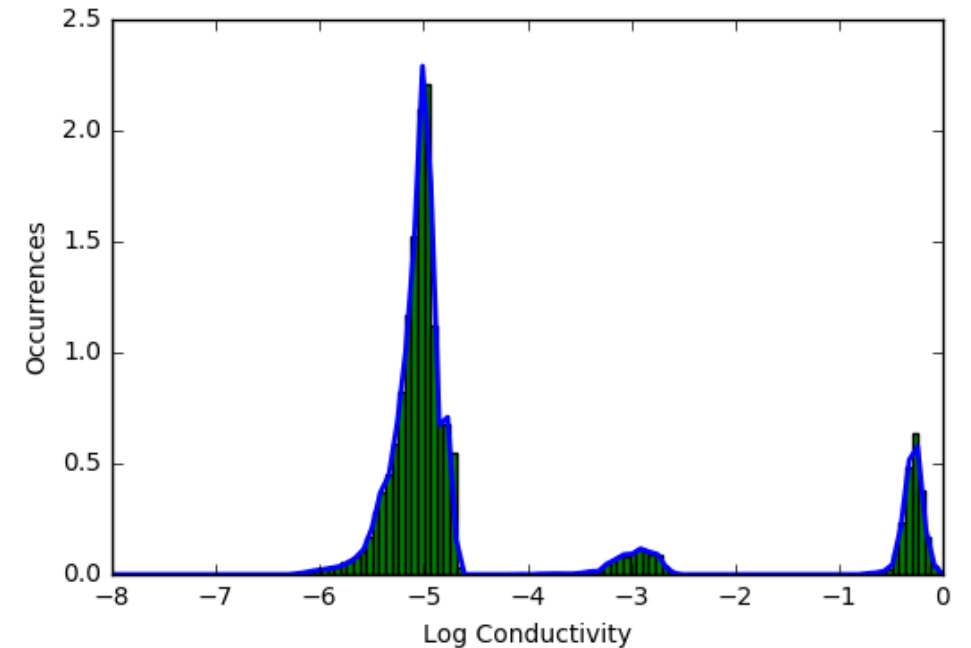
- Some parts of the model will be sparse, others will be smooth
- Can also define dependencies between variables
 - Ex: If $m_1=0$ then $m_2=0$, else both are smoothed

What can we do with it?

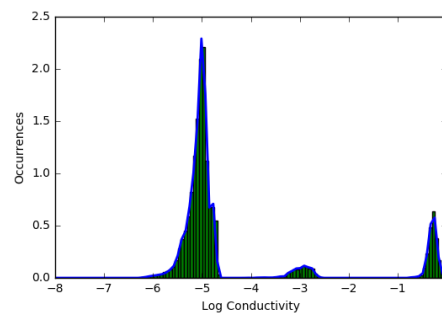
- My interest these days:
- Mixture of Gaussians prior

$$P(m) = \sum_{i=1}^k \pi_i p(m | z_i = c)$$

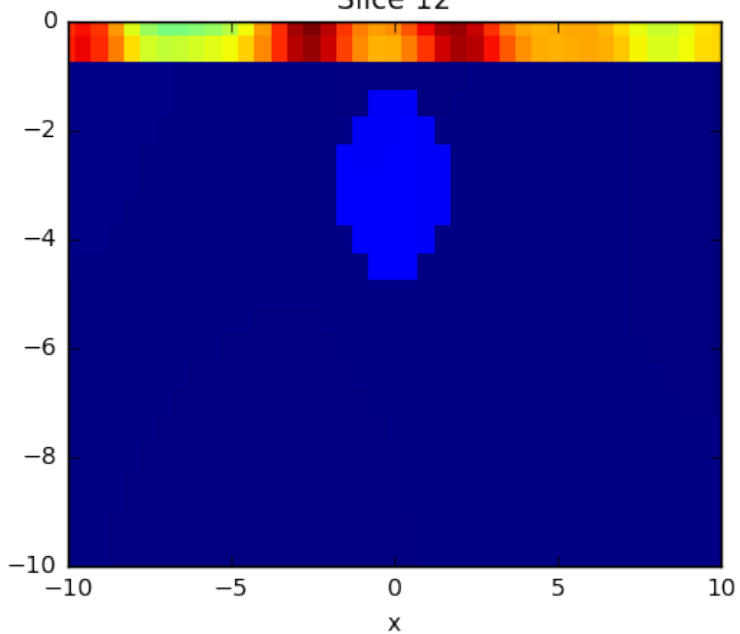
- Can fit any parameter distributions
- Can describe complex dependencies
- Easily scalable (as long as data keep up with dimensionality)
- Englobe k-means, Parzen Windows as special cases



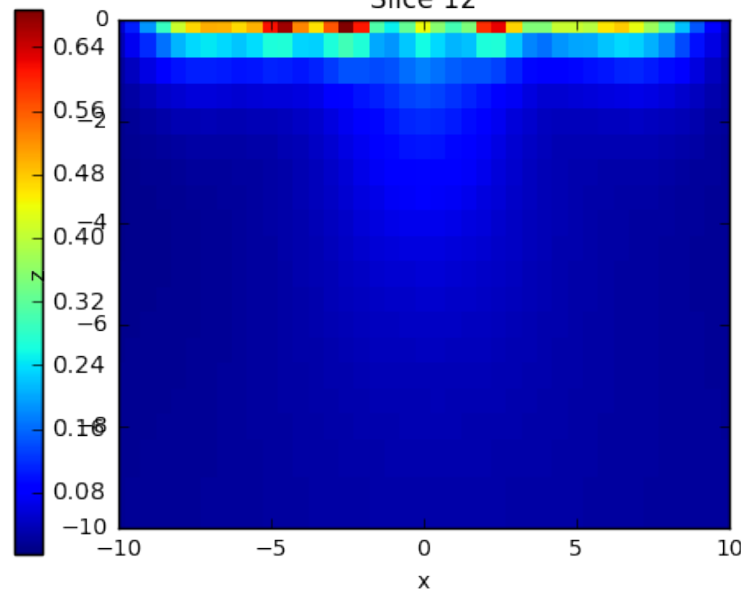
Motivation



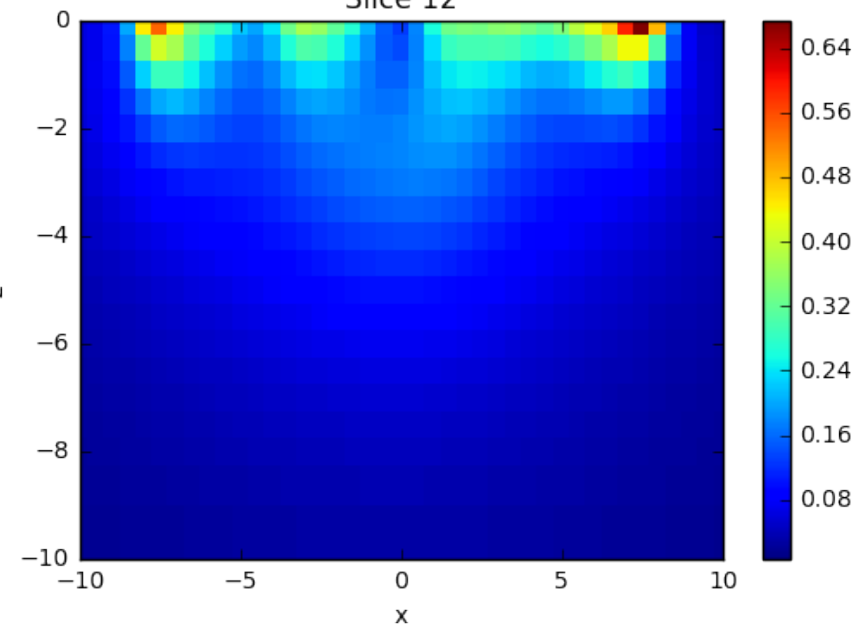
Slice 12



Slice 12



Slice 12



Mixture of Gaussians as prior

- MAP estimator

$$\operatorname{argmin}_m \frac{1}{2} \|Am - d\|_{\Sigma^{-1}}^2 - \log\left(\sum_{i=1}^k \pi_i p(m|z_i = c)\right)$$

- Where things get complicated numerically:
 - $\log(\sum(\exp(\dots)))$: well-known to underflow or overflow easily (there are some tricks fortunately)
 - Gradient and Hessian define everywhere ($\exp > 0$) but can get arbitrary large
 - No guarantee to have a positive negative log-likelihood

Mixture of Gaussians as prior

