

# Noisy-input classification of Fermi-LAT unidentified point-like sources

### Bryan Zaldivar

IFT/UAM Madrid

#### work in progress with:



Javier Coronado-Blázquez, Viviana Gammaldi Miguel A. Sánchez-Conde Machine Learning Group



Carlos Villacampa-Calvo, Eduardo Garrido Merchán Daniel Hernández-Lobato

# **Motivation and Data origin**

#### **4FGL Fermi-catalog**

circa 5000 point-like sources, out of which ~ 1500 are unidentified (unID)

#### Can we classify those unIDs alla supervised learning?



### What if some of the unIDs are better classified as dark matter?

**Include the dark matter into the**  $\beta$ **-plane!** 

## **Data visualization**



# **Data visualization**



### - unID's seem to be distributed similarly to the ID's

#### - error bars on $\beta$ partially correlated also with $\sigma_{ m curv}$

Machine Learning procedure Step # 1

# Standard classification without input uncertainties

warm up: want to know what the simplest thing to do can give you

Considered classifiers:

- Naive Bayes
- Logistic regression
- Random Forest

work in progress, but conceptually trivial...

### Next steps:

- search for an out-of-the-box classifier dealing with noisy inputs
- search for a paper addressing the classification with noisy inputs
- call your ML-expert colleagues, ask them for references!
- do it ourselves!! 🗸

Machine Learning procedure Step # 2: incorporating input uncertainties

## **Bayesian classification with parametric models**

$$\mathcal{D}_{\text{train}} = \{\mathbf{x}_i, y_i\}_{i=1}^N \qquad y_i = 1, 2, ..., K$$

 $t_{ik}$  : one-hot-encoding of  $y_i$ 

Parametric models assume a specific form for the Likelihood of data

$$-\log p(\mathcal{D}_{\text{train}}|\mathbf{w}) = \left[-\sum_{i=1}^{N} \sum_{k=1}^{K} t_{ik} \log p(C_k|\mathbf{x}_i, \mathbf{w})\right] \text{ Cross-entropy}$$

$$p(C_k|\mathbf{x}_i, \mathbf{w}) = \sigma_k(f(\mathbf{x}_i, \mathbf{w})) \text{ (softmax function)}$$

and assume a specific form for the function  $f(\mathbf{x}_i, \mathbf{w})$  (e.g. a neural network)

In Bayesian approach, we build the predictive distribution for a new point  $\mathbf{x}_*$ 

$$p(y_*|\mathbf{x}_*, \mathcal{D}_{\text{train}}) = \int p(y_*|\mathbf{x}_*, \mathbf{w}) p(\mathbf{w}|\mathcal{D}_{\text{train}}) d\mathbf{w}$$
$$p(\mathbf{w}|\mathcal{D}_{\text{train}}) = \frac{p(\mathcal{D}_{\text{train}}|\mathbf{w})p(\mathbf{w})}{p(\mathcal{D}_{\text{train}})}$$



## **Gaussian Process**

- Rasmussen & Williams, 2006

GP approach is non-parametric: no predefined form for  $f(\mathbf{x}_i, \mathbf{w})$ 

Instead you have a Gaussian *distribution over functions* (in case of regression)



## **Classification with noisy input using Gaussian Processes**

$$\mathcal{D}_{\text{train}} = \{\mathbf{x}_i + \epsilon_i, y_i\}_{i=1}^N \qquad y_i = 1, 2, ..., K$$

- As usual: introduce one output latent variable  $f_i^k$  per point *i* per class *k*,
- NEW: introduce one input latent variable  $\tilde{\mathbf{x}}_i$  per point *i*

$$\mathbf{x}_i = \tilde{\mathbf{x}}_i + \epsilon_i$$

The predictive distribution for a class  $y_*$  at a test point  $\mathbf{x}_*$ 

$$p(y_*|\mathbf{x}_*, \mathcal{D}_{\text{train}}) = \int p(y_*|\mathbf{f}_*) p(\mathbf{f}_*|\mathbf{\tilde{x}}_*, \mathcal{D}_{\text{train}}) p(\mathbf{\tilde{x}}_*|\mathbf{x}_*) \ d\mathbf{f}_* d\mathbf{\tilde{x}}_*$$

usual term

Gaussian posterior (new term)

$$p(\mathbf{f}_*|\tilde{\mathbf{x}}_*, \mathcal{D}_{\text{train}}) = \int d\mathbf{F} d\tilde{\mathbf{X}} \ p(\mathbf{f}_*|\tilde{\mathbf{x}}_*, \mathbf{F}) p(\mathbf{F}|\tilde{\mathbf{X}}, \mathbf{y}) p(\tilde{\mathbf{X}}|\mathbf{X}, \mathbf{y})$$

Inference

Costly:<br/> $\mathcal{O}(N^3)$ Sparse GP<br/> $\mathcal{O}(M^2N)$  $p(\mathbf{F}|\tilde{\mathbf{X}},\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{F},\tilde{\mathbf{X}})p(\mathbf{F}|\tilde{\mathbf{X}})}{p(\mathbf{y})}$ intractableVariationalNon-Gaussian Likelihood

## **Sparse Gaussian Process**

- here inspired in Titsias (2009)

$$p(\mathbf{f}_*|\tilde{\mathbf{x}}_*, \mathcal{D}_{\text{train}}) = \int d\mathbf{F} d\tilde{\mathbf{X}} \ p(\mathbf{f}_*|\tilde{\mathbf{x}}_*, \mathbf{F}) p(\mathbf{F}|\tilde{\mathbf{X}}, \mathbf{y}) p(\tilde{\mathbf{X}}|\mathbf{X}, \mathbf{y})$$
  
involves inverting an *N*x*N* matrix, **cost**  $\mathcal{O}(N^3)$ 

**Idea** is to make inference on a smaller M < N set of function points, which represent approximately the entire posterior over the *N* function points.

$$p(\mathbf{f}_*|\tilde{\mathbf{x}}_*, \mathcal{D}_{\text{train}}) = \int d\mathbf{F} d\tilde{\mathbf{X}} d\mathbf{U} \ p(\mathbf{f}_*|\tilde{\mathbf{x}}_*, \mathbf{F}, \mathbf{U}) p(\mathbf{F}, \mathbf{U}|\tilde{\mathbf{X}}, \mathbf{y}) p(\tilde{\mathbf{X}}|\mathbf{X}, \mathbf{y}) \qquad \begin{array}{l} \mathbf{F} : \ \dim = N \times K \\ \mathbf{U} : \ \dim = M \times K \end{array}$$
$$= \int d\mathbf{F} d\tilde{\mathbf{X}} d\mathbf{U} \ p(\mathbf{f}_*|\tilde{\mathbf{x}}_*, \mathbf{F}, \mathbf{U}) p(\mathbf{F}|\mathbf{U}, \tilde{\mathbf{X}}, \mathbf{y}) p(\mathbf{U}|\tilde{\mathbf{X}}, \mathbf{y}) \mathbf{p}(\tilde{\mathbf{X}}|\mathbf{X}, \mathbf{y})$$

If U were sufficient statistics for  $f_*$  , we were left with

$$p(\mathbf{f}_*|\tilde{\mathbf{x}}_*, \mathcal{D}_{\text{train}}) \approx \int d\tilde{\mathbf{X}} d\mathbf{U} \ p(\mathbf{f}_*|\tilde{\mathbf{x}}_*, \mathbf{U}) p(\mathbf{U}|\tilde{\mathbf{X}}, \mathbf{y}) p(\tilde{\mathbf{X}}|\mathbf{X}, \mathbf{y})$$
  
Cost:  $\mathcal{O}(NM^2)$ 

# **Variational Inference**

- Jordan, Ghahramani, Jaakkola & Saul, 1999

**Idea** is approximate the exact posterior distribution by an easier one (e.g. Gaussians) according to the *variational principle* 

$$p(\mathbf{F}, \tilde{\mathbf{X}}, \mathbf{U} | \mathbf{X}, \mathbf{y}) = p(\mathbf{U} | \mathbf{F}, \tilde{\mathbf{X}}) p(\mathbf{F}, \tilde{\mathbf{X}} | \mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{F} | \mathbf{U}, \tilde{\mathbf{X}}) p(\mathbf{U}) p(\mathbf{y} | \mathbf{F}) p(\tilde{\mathbf{X}} | \mathbf{X}, \mathbf{y})}{p(\mathbf{y})}$$

$$q(\mathbf{F}, \mathbf{\tilde{X}}, \mathbf{U}) = p(\mathbf{F} | \mathbf{U}, \mathbf{\tilde{X}}) q(\mathbf{U}) q(\mathbf{\tilde{X}})$$

# Minimize KL $\left[q(\mathbf{F}, \tilde{\mathbf{X}}, \mathbf{U}) || p(\mathbf{F}, \tilde{\mathbf{X}}, \mathbf{U} | \mathbf{X}, \mathbf{y})\right]$ w.r.t. $q(\mathbf{F}, \tilde{\mathbf{X}}, \mathbf{U})$

Kullback-Leibler divergence:

$$\mathrm{KL}[q(z)||p(z|x)] = \int dz \ q(z) \ln\left(\frac{q(z)}{p(z|x)}\right) \ge 0$$

# Likelihood of the model

 $N \quad K$ 

 $p(\mathcal{D}_{\text{train}}|\mathbf{w}) = \prod \prod \sigma_k (f(\mathbf{\tilde{x}}_i, \mathbf{w}))^{t_{ik}}$ 

i = 1 k = 1

Common form in parametric models:

"generalized Bernoulli" (the -log of which is the cross entropy)

$$\sigma_k(f(\mathbf{\tilde{x}}_i, \mathbf{w})) = p(y_i | f(\mathbf{\tilde{x}}_i, \mathbf{w})) \text{ e.g. if 3 classes: } \mathbf{0.05} \mathbf{0.80} \\ y_i = 1 \mathbf{0} \mathbf{0.15} \\ y_i = 2 \mathbf{0} \mathbf{0.15} \\ y_i = 3 \mathbf{0} \mathbf{0.15} \\ y_i = 3 \mathbf{0} \mathbf{0.15} \\ \mathbf{0.80} \\ \mathbf{0.15} \\ \mathbf{0.80} \\ \mathbf{0.15} \\ \mathbf{0.80} \\ \mathbf{0.15} \\ \mathbf{0.15$$

Instead here Misclassification noise included in the prior for  $f_k(\tilde{\mathbf{x}}_i)$ Labelling noise (with probability e) also included: Labelling rule:  $y_i = \underset{k}{\operatorname{argmax}} f_k(\tilde{\mathbf{x}}_i)$ Likelihood for label at point i:  $p(y_i | \mathbf{f}_i) = \prod_{k \neq y_i} \Theta(f_{y_i}(\tilde{\mathbf{x}}_i) - f_k(\tilde{\mathbf{x}}_i))$ (noiseless)  $p(y_i | \mathbf{f}_i) = (1 - e) \prod_{k \neq y_i} \Theta(f_{y_i}(\tilde{\mathbf{x}}_i) - f_k(\tilde{\mathbf{x}}_i)) + \frac{e}{K - 1} \left[ 1 - \prod_{k \neq y_i} \Theta(f_{y_i}(\tilde{\mathbf{x}}_i) - f_k(\tilde{\mathbf{x}}_i)) \right]$  12

D. Hernandez-Lobato, J.M. Hernandez-Lobato & P. Dupont, 2011

### Results

(python + TensorFlow)

# **Results on toy data**

ex. of dataset



- we compare with a standard GP classif. without noise
- we modify an existing GP noise model for regression McHutchon & Rasmussen, 2011

Generate a set (~100) of synthetic datasets to evaluate average performance



	Input noise level		Input noise level		Input noise level	
	$-\mathrm{ln}\mathcal{L}_{\mathrm{test}}$	Err. rate	$-ln\mathcal{L}_{test}$	Err. rate	$-ln\mathcal{L}_{test}$	Err. rate
Noiseless model		0.110				
NOISCIESS IIIOUEI	0.76	0.113	1.14	0.164	1.54	0.218
Rasmussen-like	0.321	0.109	0.53	0.158	0.77	0.209
This work	0.259	0.108	0.37	0.158	0.50	0.210

# **Conclusions/work in progress**

- Unidentified point-like sources can be classified among predefined known classes (including the potential dark matter class)
- Interestingly, including the dark matter class into the well-known beta-plane for point-like sources results in a reasonably good separability
- Only non-straightforward issue with this problem: inputs come with their own error bars surprisingly not yet explicitly addressed in the context of ML classification!
- A warm-up classification exercise w/o error bars is being conducted
- Error bars are incorporated in a Gaussian Process model for multiclass classification, by treating the input as a noisy realization of extra latent variables to be learned.
- Very satisfactory preliminary results with synthetic data
- Time to apply it to real Fermi-LAT data!

Thank you!

bckp

# Classification with error bars in the input (parametric approach)

Suppose you have data 
$$\{\mathbf{x}_i \pm \Delta \mathbf{x}_i; \mathbf{t}_i\}_{i=1}^N$$
  
 $\dim(\mathbf{x}_i) = D$   $\mathbf{t}_i = \{t_{i1}, \dots, t_{iK}\}$  e.g. If  $\mathbf{X}_i$  in class 2  $K = 3$   
"one-hot-encoding"  $\mathbf{t}_i = \{0, 1, 0\}$ 

Are noisy samples from unknown means  $\tilde{\mathbf{x}}_i$  assume  $\mathcal{N}(\mathbf{x}_i|\tilde{\mathbf{x}}_i,\Delta\mathbf{x}_i)$ 

Then the (- log) joint Likelihood of data can be written as

$$-\log p(\mathbf{X}, \mathbf{T} | \tilde{\mathbf{X}}, \mathbf{w}) = -\sum_{i=1}^{N} \sum_{k=1}^{K} t_{ik} \log p(C_k | \tilde{\mathbf{x}}_i, \mathbf{w}) - \sum_{i=1}^{N} \log \mathcal{N}(\mathbf{x}_i | \tilde{\mathbf{x}}_i, \Delta \mathbf{x}_i)$$

$$p(C_k|\tilde{\mathbf{x}}_i, \mathbf{w}) = \sigma(f(\tilde{\mathbf{x}}_i, \mathbf{w}))$$

e.g. a linear model, or a NN model