

# Physics 693: Modern ML for Physics & Astronomy

ML is revolutionizing nearly every field of science & society more broadly

Powerful new tool — enabling new analyses previously impossible



made possible by  
{}  
- Big data  
- Computing (GPU)  
- Sophisticated algorithms (NNs, -)

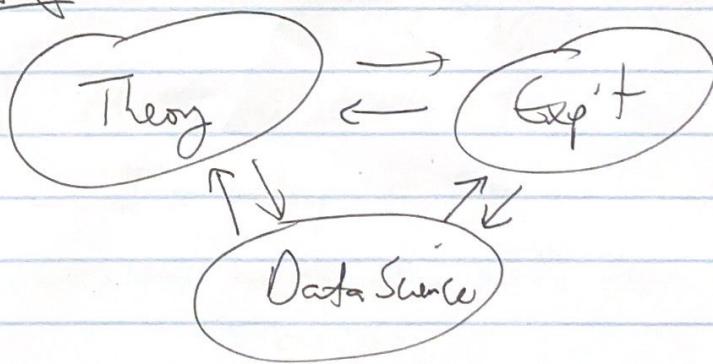
— enhancing sensitivity & precision

— accelerating simulation & inference

— unifying solutions to problems across domains

ML like a telescope!  
or microscope

~~This course:~~



~~What's new?~~

↓ formal mathematics  
novel mix of theory, statistics, numerics

in ML same paper can have all 3!

# This Course: general ML concepts

~~lectures~~

popular & state of the art architectures

(I know more → less)

applications to LHC, Astro, Gsm, Condensed Matter

Required

- No HWs but will be hands on tutorials & exercises
  - ^ & no exams
- BB + demos/shader

→ - Will ask people to pitch applications to present in class

(either you can present, or we can learn it together &  
I will present)  
need input!

- No domain knowledge of any subfield required
- Prior experience w/ python, numpy, matplotlib would be good

- lectures M-Th 10:20-11:40

- poll for makeup slot?

- in Nov some travel - either Zoom or guest lectures  
(4 lectures)

- will provide refs & reading on course website; no textbook

- will not require 568 (but will) but would help if you took it  
will try to go fast in overlapping content (SGD, backprop, NNbasics...)

please interrupt  
frequently &  
ask me lots of Q's

# What is ML?

"learning from data"      "generalized curve fitting"

- Data:  $\vec{x}_i \in \mathbb{R}^d$  d can be high!  
e.g.  $10^2, 10^3, \dots$

"low level"



vs

"high level"  
features

↳ "tabular data"

↳ a vector of features

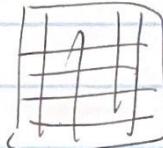
e.g.

↳ particle momenta or  $x_1, y_1, z_1, v_x, v_y, v_z, G, BR, \dots$   
@ LHC

features could include labels  
 $\vec{y}_i$ ; e.g. cat vs dog, ...

@ Gaia

or



pixel intensities in image

$N$  could be large

$10^6, 10^9, \dots$

(e.g. events @ LHC)

stars in Gaia

images (or ~~sets~~ of scalars)

- Assume in most applications: each

$$\vec{x}_i \sim p_{\text{data}}(x)$$

drawn iid

parameters

$$f(\vec{x}_i; \theta) \rightarrow \text{data}$$

- Goal: minimize some "loss" or "objective" fn

$$L(\theta) = \sum_{i=1}^N L(f(\vec{x}_i; \theta), \vec{y}_i)$$

wrt parameters  $\theta$ .

Example: maybe want  $f(x_i; \theta) = y_i$  target labels  
- discrete "classification"

could use  $\mathcal{L}(\theta) = (f(x_i; \theta) - y_i)^2$  continuous regression  
"mean squared error (MSE)"

- many more examples to follow (& how to choose loss functions in principled way)

Where do labels come from?

- if data is simulated have access to "ground truth"
- in actual data - human labeling? (common in real world apps like Imagenet, also Astro citizen science)
- no labels or labels derived from data itself

DNN, RNN, CNN, Graph NNs  
transformers, ...

→ building blocks

Categories of ML

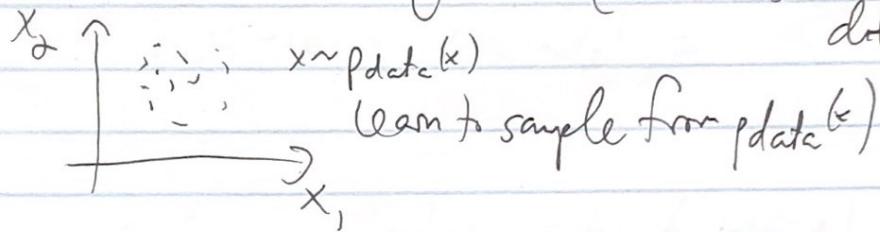
- pretty much solved problems in ML
  - fully supervised - all data labeled & used in training
    - classification, regression
  - unsupervised - no labels!
    - generative models, density estimation, some kinds of anomaly detection
  - weakly supervised - noisy labels
  - semi supervised - mix of labeled & unlabeled data
  - self supervised - labels generated from data e.g. contrastive loss

this is where  
the active ML  
development  
is happening

*(less than expensive)*  
Also really important for science — strive to be as  
data-driven as possible  
(simulation-free)

In more detail:

- Generative modeling (GANs, VAEs, flow-based models, ...)

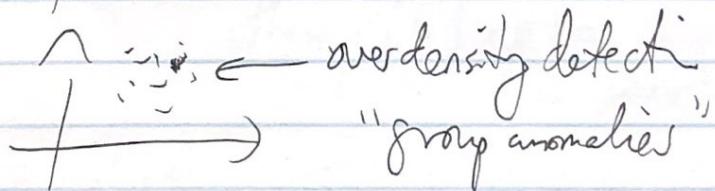
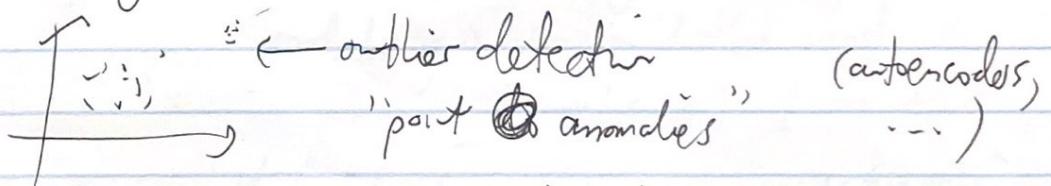


- Density estimation — learn  $p_{\text{data}}(x)$  (flow-based models, ...)

(Gen. model  $\leftrightarrow$  Density est.)

one doesn't guarantee the other

- Anomaly detection



General principle for loss fns: Maximum Likelihood Estimation (MLE)

want to maximize  $P(\text{data} | \text{model})$

if data iid

$$\prod_{i=1}^N P(x_i | \theta)$$

$$L = -\log P(\text{data} | \text{model}) = -\sum_{i=1}^N \log P(x_i | \theta)$$

- MLE has properties in lot of large  $N$ :

- consistency (as  $N \rightarrow \infty$ , estimated  $\theta \rightarrow$  true  $\theta$ )

- efficiency (as  $N \rightarrow \infty$ , minimum variance estimator of  $\theta$ )

Example:

Suppose want to predict  $y$  given  $x$  (regression)

model:  $y$  is gaussian dist'ed around

$$f(x; \theta) \sim \text{some dist'}$$

$$\text{then } P(y | x, \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-f(x; \theta))^2}{2\sigma^2}}$$

$$L = -\log P = \sum \frac{(y_i - f(x_i; \theta))^2}{2\sigma^2} + \log \sqrt{2\pi\sigma^2}$$

If  $\sigma$  known  $\rightarrow$  MSE! (otherwise abs. fit or  $\sigma$ )

Ex: binary classification

~~what about bin~~  
labels  $y_i = 0$  or  $1$  from data  $x_i$

model:  $f(x_i; \theta) = \text{prob of } x_i \text{ being label 1.}$

~~P(data)~~  
data

$\begin{matrix} x_i & \xrightarrow{i=1, \dots, N_1} & x_1 \\ \downarrow & \text{label 1} & | \\ x_i & \xrightarrow{i=N+1, \dots, N+N_0} & x_2 \\ \downarrow & \text{label 0} & | \\ x_i & & x_N \end{matrix}$

$$P(\text{data} | \text{model}) = \prod_{i=1}^{N_1} f(x_i; \theta) \prod_{i=N+1}^{N+N_0} (1 - f(x_i; \theta))$$

$$L = -\log P = -\sum_{i \in 1} \log f(x_i; \theta) - \sum_{i \in 0} \log (1 - f(x_i; \theta))$$

$$= -\sum_i (y_i \log f(x_i; \theta) + (1-y_i) \log (1-f(x_i; \theta)))$$

"binary cross entropy loss"  $\rightarrow$  best loss for classification

$\rightarrow$  could use MSE

but it would be suboptimal

## Lecture II

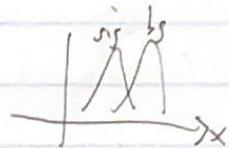
- email list
- slack

- pitches for applications

### Neyman-Pearson Lemma

$$L = - \left( \sum_{i \in S} \log f(x_i) + \sum_{i \in B} \log (1-f(x_i)) \right)$$

(what  $f$  minimizes  $L$  in ideal case?)



↓ continuum limit ( $\infty$  data)

$$L = - \int dx \left[ p_S(x) \log f(x) + p_B(x) \log (1-f(x)) \right]$$

$$\int dx p(x) \approx \sum_{x \sim p(x)} \quad f \rightarrow f + \delta f, \quad \frac{\partial L}{\partial \delta f} \Big|_{f=f_0} = 0 \quad \begin{matrix} \downarrow \\ \text{assume } N_S = N_B \\ \text{WLOG} \end{matrix}$$

$$0 = \frac{p_S(x)}{f_*} - \frac{p_B(x)}{1-f_*} \Rightarrow f_* = \frac{p_S(x)}{p_S(x) + p_B(x)}$$

so optimal classifier  $\hat{f}_* = \frac{R(x)}{R(x)+1}$

$$R(x) = \frac{p_S(x)}{p_B(x)} \quad \frac{f_*}{R(x)} = 1$$

$f_*$  is monotonic w/  $R(x)$

$$\Rightarrow f_* > f_C \iff R(x) > R_C \quad \begin{matrix} \text{cut on } f \text{ equals} \\ \text{cut on } R \end{matrix}$$

so optimal classifier is likelihood ratio

"Neyman-Pearson Lemma"

- powerful result! fundamental result in statistics  
(LR uniformly most powerful test for simple hypothesis)

I later will see:

- Cart and random forest → learn <sup>ML</sup> classifier from samples  
→ approach LR.  
"likelihood ratio trick"

- Another perspective: Bayes Thm

$$f_x = \frac{p_S(x)}{p_S(x) + p_B(x)} = \frac{p(x|S)}{p(x|S) + p(x|B)} = \frac{\underbrace{p(x|S)p(S)}_{\sim}^{\frac{1}{2}}}{p(x)} \\ = \frac{p(x|S) + p(x|B)}{2}$$

so  $f_x$  is the prob of signal given  $x$  which is where we started!

i.e. BCE is minimized when our model for this prob  
= true prob.

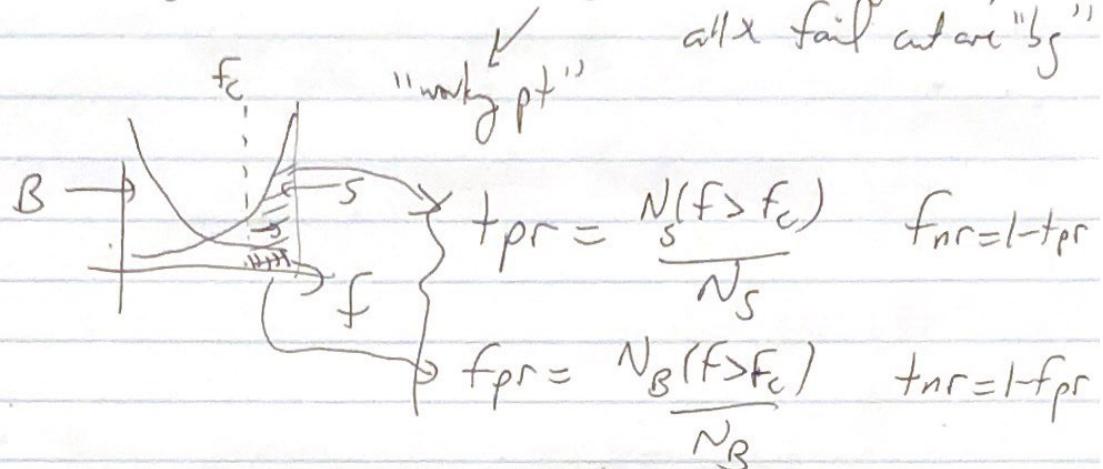
Back to NP Lemma: what does "most powerful" test mean?

→ metrics for binary classification

How do we use binary classifiers?

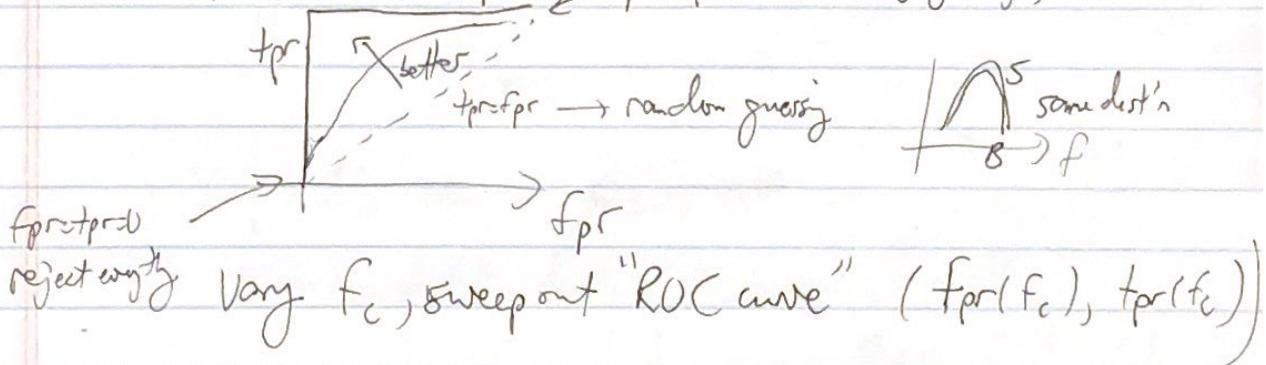
classified as

Usually: cut  $f(x) > f_c$  all  $x$  satisfy are "S" all  $x$  fail cut are "B"



- accuracy =  $\max_{f_c} \frac{tpr + (1-fpr)}{2}$

- ROC curve  $tpr=1 \leftarrow fpr=tpr=1$  (let everythg thru)



→ AUC "Area under the curve"

$$\begin{cases} AUC=1 \text{ perfect} \\ AUC=0.5 \text{ random} \end{cases} \quad \begin{matrix} \text{(same for acc} \\ \text{but } AUC \neq ACC \end{matrix}$$

- AUC, ACC mostly sensitive to  $tpr, fpr \sim 0(1)$  part of ROC curve

when  $N_{\text{bg}} \gg N_{\text{sig}}$  (as in many cases, e.g. LHC)

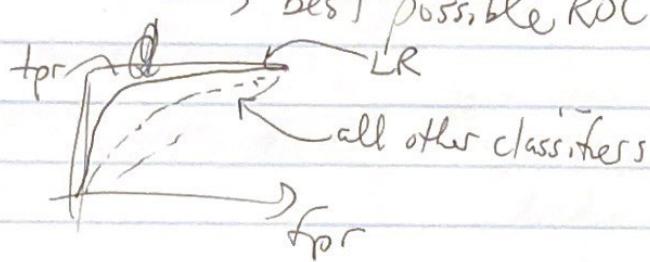
(are more about  $f_{\text{pr}} \ll 1$ ) ~~work~~ part of ROC curve

→ often report  $f_{\text{pr}}$  @ fixed  $t_{\text{pr}}$  e.g.  $t_{\text{pr}}=50\%$   
or 30%

or rejection factor  $\left\{ \begin{array}{l} R_{50} = \frac{1}{f_{\text{pr}} @ t_{\text{pr}} = 50\%} \\ R_{30} \text{ etc} \end{array} \right.$

- Neyman-Pearson LR classifier is optimal

→ best possible ROC curve



- In practice cannot achieve NP optimality

— exact likelihoods unknown

— finite training data → Bias of big detcs

— limited model capacity (expressivity)

↳ NNs very ~~poorly~~ expressive

↳ can maybe get very close to  
optimal!

Then possibly biased on val set  
→ final metrics on another held out dataset  
"test set"

Example 1: Gaussian toy model (10d)

$$\begin{cases} b_j(0) & p(x) = N(0, I)^{10} \\ \bar{s}_j(1) & p_{\bar{s}_j}(x) = N(0.5, I)^{10} \end{cases}$$

- can compare w/ NP Lemma
- train simple DNN

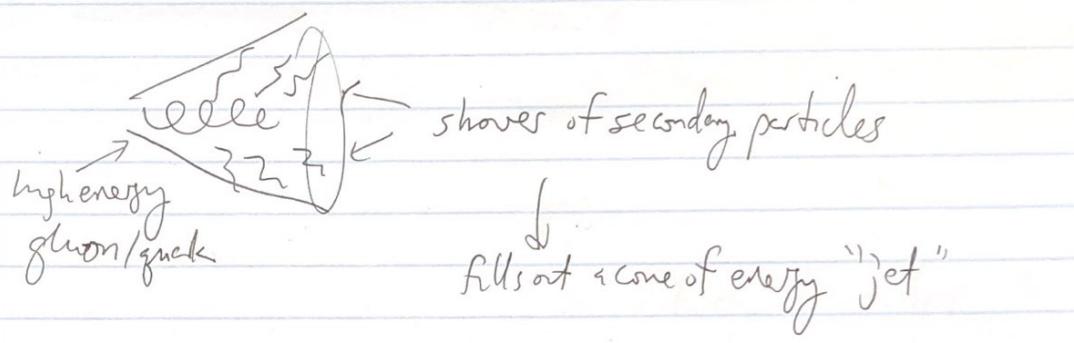
$$P(s_j|x) = \frac{p_{s_j}(x)}{p_j(x) + p_{\bar{s}_j}(x)}$$

- on samples
- - plot losses
- - ROC curve
- - AUC
- - classifier output

- Jupyter nb demo

## Example 2: Jet classification @ LHC

- Jets are essential objects @ LHC  
Product of strong interaction ( $\text{QCD}$ )  $\rightarrow$  showers of hadrons

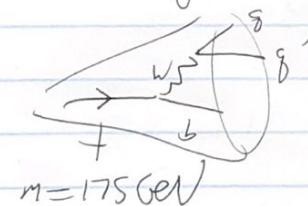


- Many different flavors of jets - important to tell them apart!

$g, g, b, W, Z, \text{top}, \dots$

↓      ↓      ↓  
5M measurements      BSM searches  
Typically QCD jets ( $g, g$ ) are background & everything else is signal

- Especially rich example: boosted top jets



• Decays to  $b\bar{b}$

$W$  can decay to neutrinos

• If top is very energetic  $\rightarrow$  boosted decay products collimated "fat jet"

Note: QCD

xsec thousands  
times more than  
top (or more)  
dependence  
energy  
 $(\sim 1000)$   
 $@ m_j \sim 100\text{ GeV}$ )

- Top jet has substructure  $\rightarrow$  expect 3 "prongs"  
of energy  $\approx b, g, g'$   
Expect  $\beta_{\text{jet}} \approx m_{\text{top}}$
- Expect  $gg' \rightarrow$  mass  $\approx m_W$ .
- expect  $b_j$  jet (<sup>on its own</sup> (wasn't use here))  
w/ special properties
- Meanwhile QCD jets are relatively structureless.  
 $\rightarrow$  can use these properties to classify tops vs QCD!

(unlike Gaussian toy model, here we do not know  
true likelihoods, so do not know NP optimal classifier)

- Jupyter nb example dems
  - $(m_j, m_W, \tau_{32})$   $\rightarrow$  cuts  $\rightarrow$  DNN
  - jet constituents  $\rightarrow$  DNN

Metric: R30  
for  $\text{C}_t \text{pr} = 30\%$   
more informative than  
AUC/ACC

A jet = collection of p<sub>T</sub>l 4-vectors "LLFs"

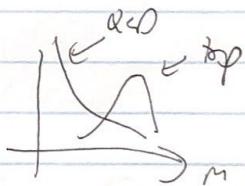
$$\vec{J}_i^{\text{ll}} = \begin{pmatrix} p_T^i, p_T^i, p_T^i, \epsilon^i \\ \vdots \\ p_T^{N_i}, p_T^{N_i}, p_T^{N_i}, \epsilon^{N_i} \end{pmatrix} \quad N_i \text{ diff from jet to jet}$$

$$(\epsilon^a)^2 = (\vec{p}^a)^2$$

HLFs like jet mass  $m_i^2 = \left( \sum_{a=1}^{N_i} \epsilon^a \right)^2 - \left( \sum_{a=1}^{N_i} \vec{p}^a \right)^2$

↓ and N-subjettiness  $\tau_{32}$

physics motivated fits of LLFs



We have seen:

cut based classifier on HLFs  $\approx$  DNN on HLFs

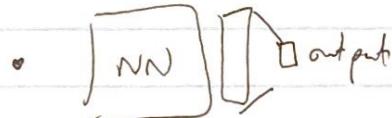
$\approx$  DNN on LLFs (30 highest  $p_T$  constituents)

"Automated feature engineering"

- DNN can construct HLFs from LLFs that perform as well as physics HLFs!

- Q: is it learning the physics HLFs or something else?  
How would we know?

A: Could add physics HLFs to inputs - perf. improves?



last layer → then or like NN engineered HLFs!  
 ↴  
 automated  
 the HLF engineer

not necessarily very last layer



↳ could interpret these as HLFs

- Is this all or can we surpass physics HLF performance?

→ yes but need more powerful NN architectures!

general lesson: → leverage symmetry / structure of data!

on DNN p<sub>i</sub> ordered constituents  
 selected 30 hardest

↑  
 like getting more data for free  
 "inductive bias" → helps machine learn more effectively

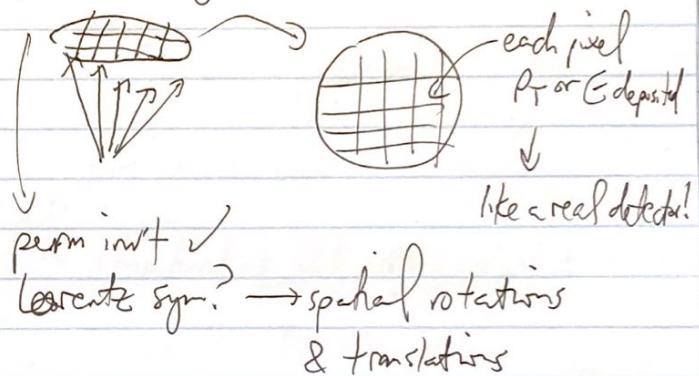
But jets have permutation invariance!

also: Lorentz symmetry → spatial rotations  
& boosts

"point cloud"

A rich playground for different ML architectures

- CNNs — jets as images



- RNNs, LSTMs, ... — jets as sequences

— not perm inv! what ordering?

- Graph NNs — perm inv!

(can add Lorentz)

- Deepsets — a perm inv! DNN

- Transformers — a perm inv! sequence NN

) plan for next few lectures + (detour to Astro w/ CNNs)