# Unit III: Anomaly Detection

- **⊙** Generally difficult ML problem: how to detect data that is anomalous?

  ↓

  Usually "less-than-supervised" — anomalous means don't know what you're looking for
  → lack labels

- **⊙** Many scientific applications of anomaly detection
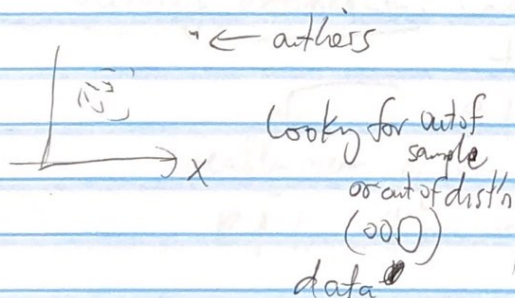  - data quality monitoring
  - triggering
  - new physics searches —— "model-agnostic" or "model-indep"
  - data-driven discovery
  
  ↳ vs "model-specific"

  → 99.99% of all current LHC searches!

- **•** Two main classes of anomaly detection:

  1. Outlier detection

  ↓

  what if we haven't found NP because we haven't searched in the right places yet?

  

  ` ← outliers`

  Looking for out of sample or out of dist'n (oo⊙) data

## Outlier detection

- A clustering problem $\longrightarrow$ many methods for unsupervised clustering
- also, a density estimation problem

(must specify $k$ in advance!)

group data $X_i$ into $k$ clusters w/ means $\vec{\mu}_a$

eg $k$ means
DBSCAN
:

$$\text{min} \sum_{a=1}^{k} \sum_{x \in C_a} \|x - \vec{\mu}_a\|^2$$

fast algorithms exist for partitioning data & minimizing w.r.t $\vec{\mu}_a$.

main drawbacks

(must specify metric in advance)

$\Rightarrow$ what if the space in which data clusters nicely is hidden? (only "see" branches in latent space?)

density estimation: want data for which $p(x) = 0$?
truly OOD

But that is self-contradictory
really mean $p_{bg}(x) = 0$?

But how will one know $p_{bg}(x)$? Not fully data-driven?
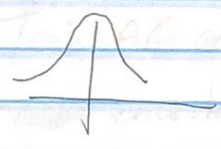Have our sample only $bg$?

In practice can never be sure $p_{bg}(x) = 0$
(finite resolution, measurement error, quantum mechanics...)

$\downarrow$

So really mean $p_{bg}(x) \ll 1$?

$\downarrow$

But this is not coordinate inv't!

→ can change $p_{bg}(x)$ to anything w/ $x \to y$

extreme example:

$p(x) = Ne^{-x^2/2}$

$x = 5$ is very anomalous?

$\downarrow$ cdf $\quad y = \int_{-\infty}^{x} dx' Ne^{-x'^2/2}$

$p(y) = 1$
(uniform) $\quad y \in [0,1)$

So no points are outliers?!

density estimate based

This issue is generally ignored in the literature...

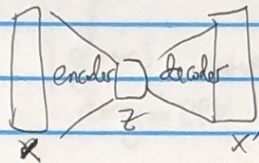Hope that there is a preferred coord frame "physically meaningful"

---

General Methods for finding latent space clustering:

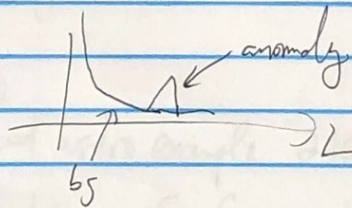— <u>Autoencoders</u> (and variational AEs)

general idea:



$$L = ||x - x'||^2 \quad \text{"reconstruction error"}$$

learn to map data back to itself through compressed
latent space
"information bottleneck"

↳ can't learn identity map $x' = x$

"one class classification" { 
- Train AE on "normal events" → learn to reconstruct well
- encounter rare, anomalous event → doesn't reconstruct well

→ can use $L$ itself as anomaly score!



Do you really need sample of "normals" to train on?
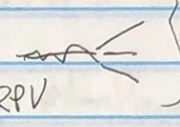then not fully unsupervised...

→ in practice AE works on bg + small amount of anomaly!
finite model capacity → still mostly learns
to reconstruct bg well

Example: Farina, Nakai & DS 1808.08992
1st appl. of AE anomaly detection to HEP!
(see also Heimel et al 1808.08979)

bg = QCD jet images

sig = {tops

{gluinos ∿∿∿<∿

decay via RPV

}

challenges:
- bg estimation?
- double AG idea
  2111.06417
- complexity bias?
  train on tops, doesn't detect
  QCD...

---

VAEs can also be used for anomaly detection
- can use recon error like AE

or can look for outliers in latent space!

"normalized AG" ↳ more like density estimation
2105.05735 (ML)
2206.14825 (HEP)

Many examples in literature
(eg VAE on MNIST missing one digit)

↓

recent astro example 2103.12102

Rubin proof-of-concept ⟶ can detect rare
transients as anomalies
using latent space clustering