# L3 : Anomaly Detection

Jan 13, 2025

Data Mining : Jeff Phillips

# Review

Distribution   $D(\Theta)$
    $unif([m])$

- Events   # trials  until  collision

    Median:   $\approx \sqrt{2m}$

    Birthday Paradox

- Event:  # trials  until  see  all

    Expected  # $= (0.577 + \ln(n)) \cdot n$

    Coupon Collectors

    $O(n \cdot \log n)$

# Anomalies

Data  $X \underset{iid}{\sim} D(\theta)$     $X = \{x_1, x_2 \dots x_n\}$

$x_i \sim D(\theta)$

What is likely?

If we have X, how to generate similar X?

## Anomalies

1. what is distribution of data? (model)

2. what would an anomaly look like.
   - score (likelihood, LLR)
   - shape.

3. How interesting?   Quantits

# Likelihood : unnormalized probability

$$L(X; \theta)$$

parameters $\theta$    e.g. $\theta = p \in \mathbb{R}$

potential anomaly $S \subset X$

Data $\overset{S}{[\circ \; \infty \; \circ]} \cdots \cdots X$

$\underline{X \setminus S}$

likelihood of anomaly

$$L(S, X \setminus S; \theta')$$

$\theta' = p', \; \theta' \in \mathbb{R}$

score   Log-likelihood ratio   $LLR(S, X) = \log\left(\dfrac{L(S, X \setminus S; \theta')}{L(X; \theta)}\right)$
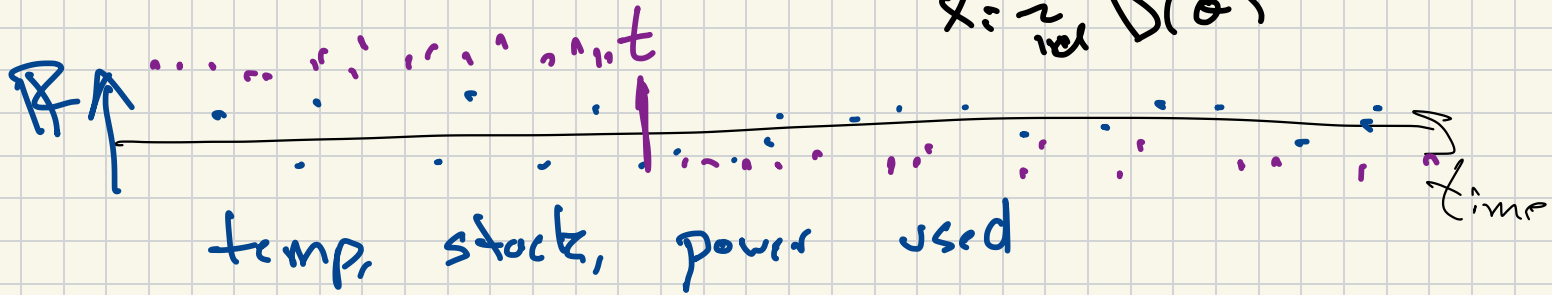
# Change Point Detection

Data $X$ : sequence $n$ real values

$$\langle x_1, x_2, \ldots x_n \rangle \qquad x_i \in \mathbb{R}$$
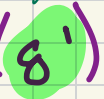
$$x_i \underset{iid}{\sim} D(\theta)$$



temp, stock, power used

$D(\theta)$  no change  $x_i \underset{iid}{\sim} D(\theta)$

mean $p$  Normal noise

$D(\theta')$ : change at time $t_a$

$x_1 \ldots x_t \underset{iid}{\sim} D(p')$ ,  $x_{t+1} \ldots x_n \underset{iid}{\sim} D(\theta')$

mean
$x_{t+1} \ldots x_n$
$\downarrow$

Litzelihood $L(X; \theta = p)$



$$P_r(x_i; \theta = p) = \frac{1}{\underline{\sqrt{2\sigma\pi}}} \cdot \exp\left(-\frac{(x_i - p)^2}{\sigma^2}\right)$$

$$L(x_i; \theta = p) = P_r[x_i; \theta = p]$$

proportional to $\propto \exp\left(\frac{-(x_i - p)^2}{\sigma^2}\right)$

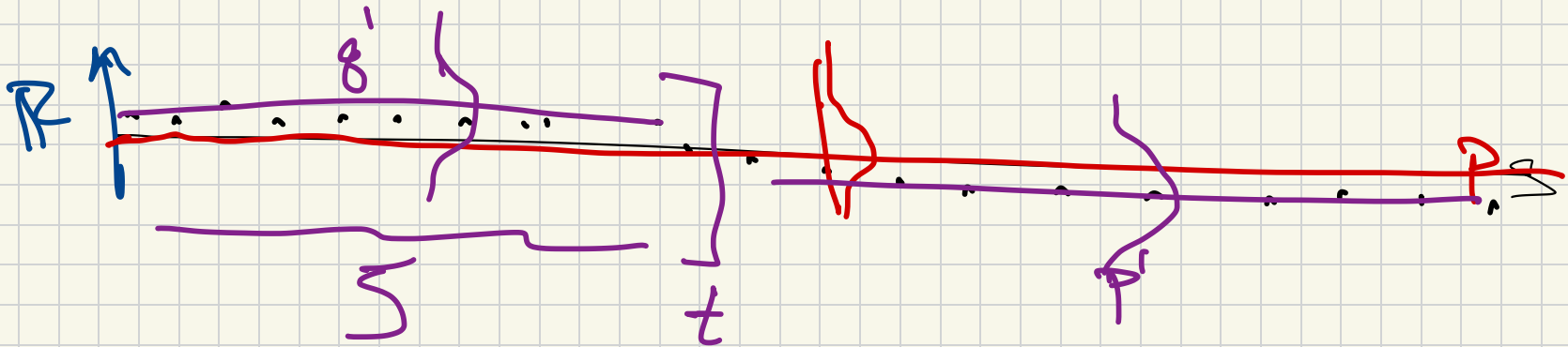$$L(X; \theta = p) \propto \prod_{i=1}^n \exp\left(\frac{-(x_i - p)^2}{\sigma^2}\right)$$

$$\ln\left(L(X; \theta = p)\right) = \ln\left(\prod_{i=1}^n \exp\left(\frac{-(x_i - p)^2}{\sigma^2}\right)\right)$$

$$= \sum_{i=1}^n \ln\exp\left(-\frac{(x_i - p)^2}{\sigma^2}\right) = -\frac{1}{\sigma^2}\sum_{i=1}^n (x_i - p)^2$$

$$\frac{\partial}{\partial p} \quad p = \frac{1}{n}\sum_{i=1}^n x_i$$

$$\ln\left(L(X; P)\right) = -\frac{1}{\sigma^2} \sum_{i=1}^{n} \left(x_i - \textcolor{red}{\boxed{P}}\right)^2$$

$$\ln\left(L(S, X\backslash S; P', g')\right) = -\frac{1}{\sigma^2}\left(\sum_{i=1}^{t}\left(x_i - g'\right)^2 + \sum_{i=t+1}^{n}\left(x_i - P'\right)^2\right)$$

$$\text{score} \quad \gamma(s) = \max_{\hat{p}, \hat{g}} \ln\left(L(s, x \cap s; \hat{p}, \hat{g})\right)$$

$$- \max_{\hat{p}} \ln\left(L(x; \hat{p})\right)$$

score $\gamma(s)$

$LLR(s, x)$

Find $S^* = \underset{S_t}{\text{argmax}} \underbrace{LLR(S_t, X)}$

$LLR(0) \to LLR(t=1)$
$O(1)$ time $\quad \dots \quad O(n) + n \cdot O(1) = O(n)$
$t=0$
$t=1$



$S_t$

$t$

LLR

$t^*$

) does not
work

$\mathbb{R}$

$[s \overset{S}{\underset{\mid}{}} t]$

now $n^2$ anomalies $S_{s,t}$

Two distribution $X \sim D(\theta)$, $X' \sim D(\theta')$

baseline          anomalous

way to measure          $(X, \mathcal{R})$          range space

$[-\infty, t]$   change point

$(s, t]$   two-sided

$\mathbb{R}^2$

balls

$$d_{\mathcal{R}}\Big(D(\theta), D(\theta')\Big) = \sup_{R \in \mathcal{R}} \left| \underset{x \sim D(\theta)}{\mathbb{E}}[R(x)] - \underset{x' \sim D(\theta')}{\mathbb{E}}[R(x')] \right|$$

max

integral probability metric

So...  I found   $S^* = \underset{S}{argmax} \; LLR(S;X)$

is  it  interesting?

- if  I totally trust  $D(\theta), D(\theta')$
  then  LLR  $\approx$  t-score  related
  $$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad to \; p\text{-values.}$$
  $\approx$  if LLR $\approx$ 2 $\Rightarrow$ p $\approx$ 0.05

What if LLR is just a "score"?

$2_a$ Draw more data from $D(\theta)$

$$X_1, X_2 \ldots X_n \sim D(\theta)$$

compute score $\quad \gamma(X_i) = \gamma_1$

$$\gamma_1, \gamma_2, \ldots \gamma_n$$

compare $\quad \gamma(X)$ to $\gamma_1 \ldots \gamma_n$

↑
input

what fraction of $\gamma_1 \ldots \gamma_n > \gamma(X)$

if $\rho = \dfrac{\# >}{n} < 0.05$ interesting

3. What if I do not know $\phi(\theta)$?

Permute existing data

$x_1, x_2, \dots x_n \in X$     randomly permute order

new $\langle x_7, x_{17}, x_3, \dots x_8 \rangle = \begin{matrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{matrix}$

$(n \times n)$     $\vdots$

Permutation test.

Calc n values $\gamma_i = \gamma(X_i)$

compare $\gamma(X)$ to $\gamma_1, \gamma_2 \dots \gamma_n$