

# Anomaly detection using surprisals

Rob J Hyndman

26 November 2025

# Coauthors



Sevvandi  
Kandanaarachchi  
CSIRO



Kate  
Turner  
ANU



David  
Frazier  
Monash U

# Outline

1 Anomalies

2 Extreme surprisals

3 Lookout algorithm

4 Conclusions

# Outline

1 Anomalies

2 Extreme surprisals

3 Lookout algorithm

4 Conclusions

# Definitions of anomalies

*an observation (or a subset of observations) which appears to be inconsistent with the remainder of that set of data.*

(Barnett & Lewis, 1978)

# Definitions of anomalies

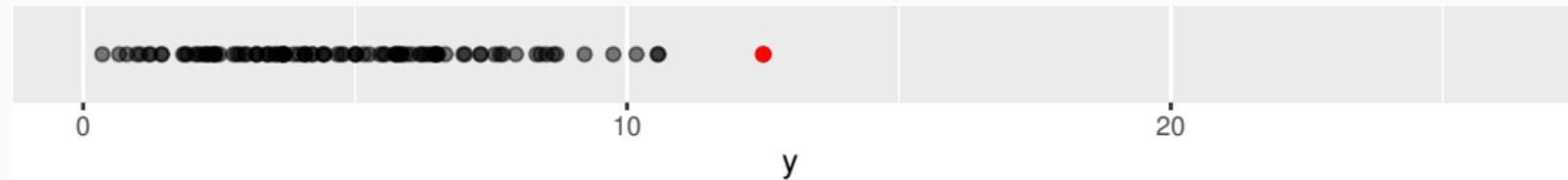
*an observation (or a subset of observations) which appears to be inconsistent with the remainder of that set of data.*

(Barnett & Lewis, 1978)

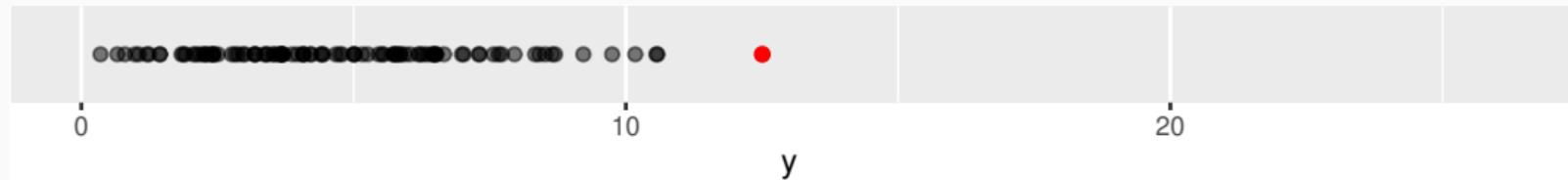
*an observation which deviates so much from other observations as to arouse suspicion it was generated by a different mechanism.*

(Hawkins, 1980)

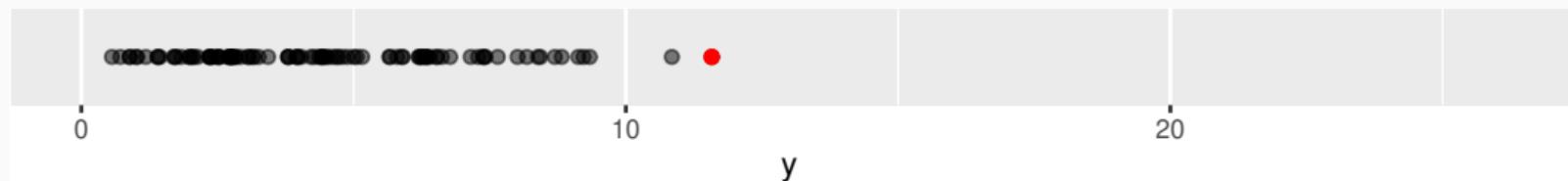
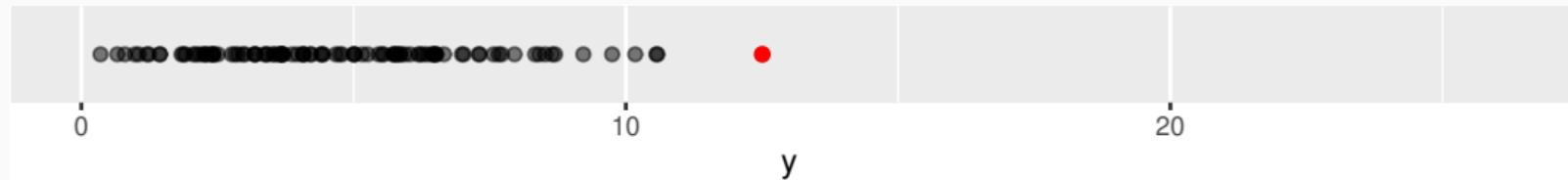
# Is this an anomaly?



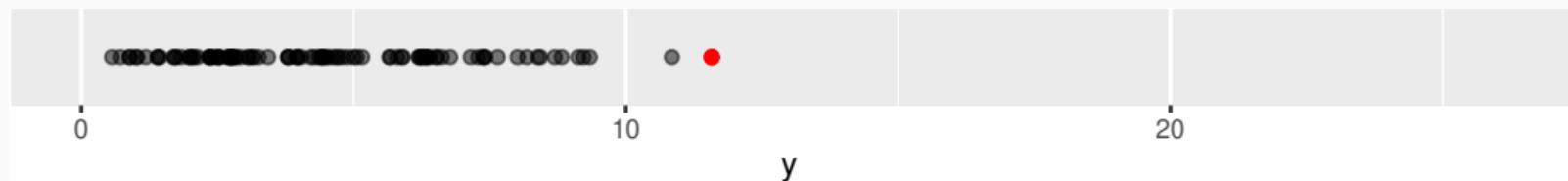
# Is this an anomaly?



# Is this an anomaly?



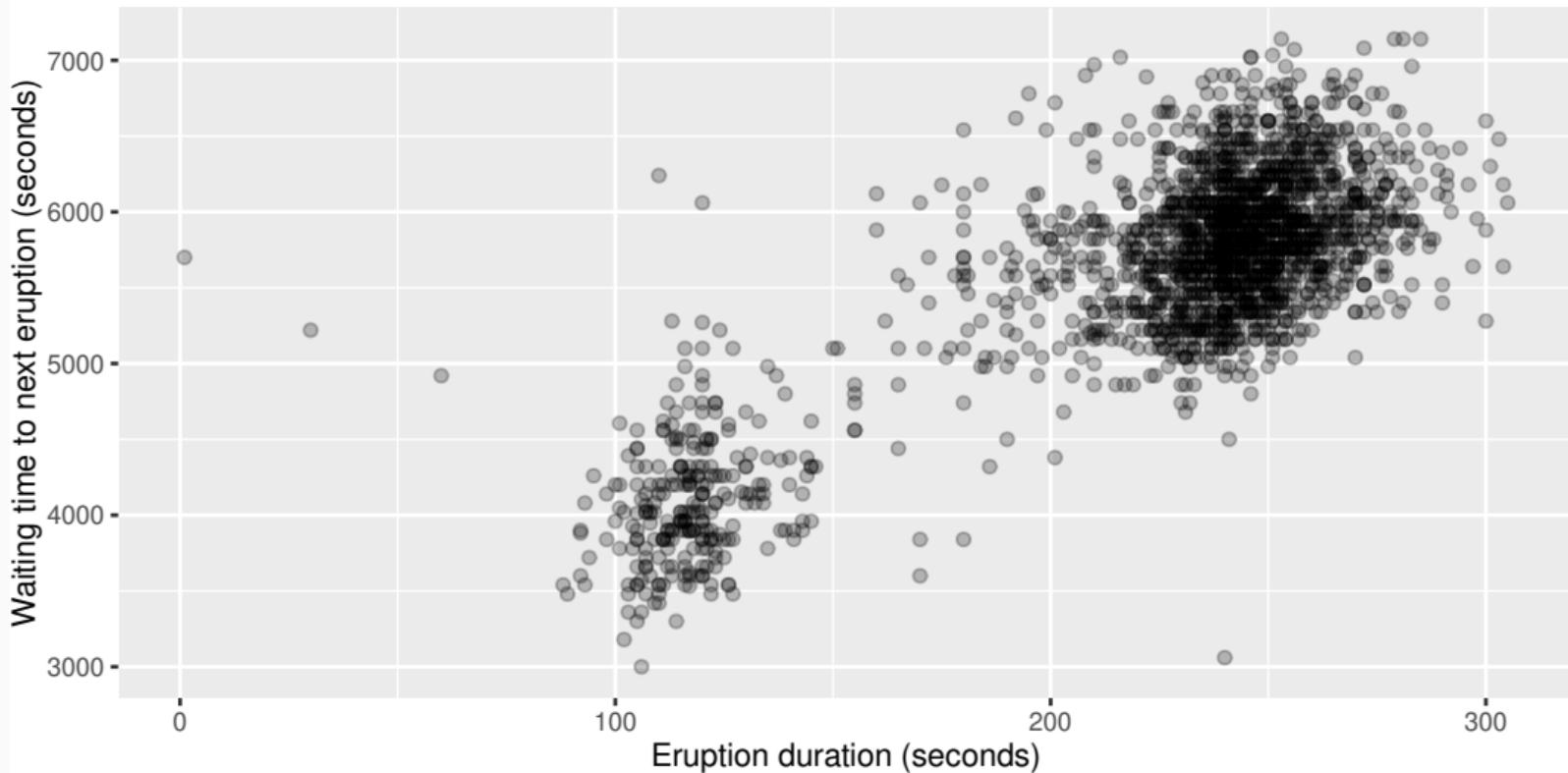
# Is this an anomaly?



All points randomly generated from a  $\chi^2_5$  distribution.

# Are there any anomalies?

Old Faithful eruptions from 14 January 2017 to 29 December 2023



# Definitions of anomalies

## Definition: Anomaly

Given a set of observations  $\{y_1, \dots, y_n\}$  and a generalized probability density  $f$ , the **anomaly score** of  $y_i$  wrt  $f$  is

$$p_i = \mathbb{P}(f(Y) \leq f(y_i))$$

where  $Y$  has density  $f$ . An observation is an **anomaly** wrt  $f$  if  $p_i < \alpha$  for some threshold  $\alpha > 0$ .

# Definitions of anomalies

## Definition: Anomaly

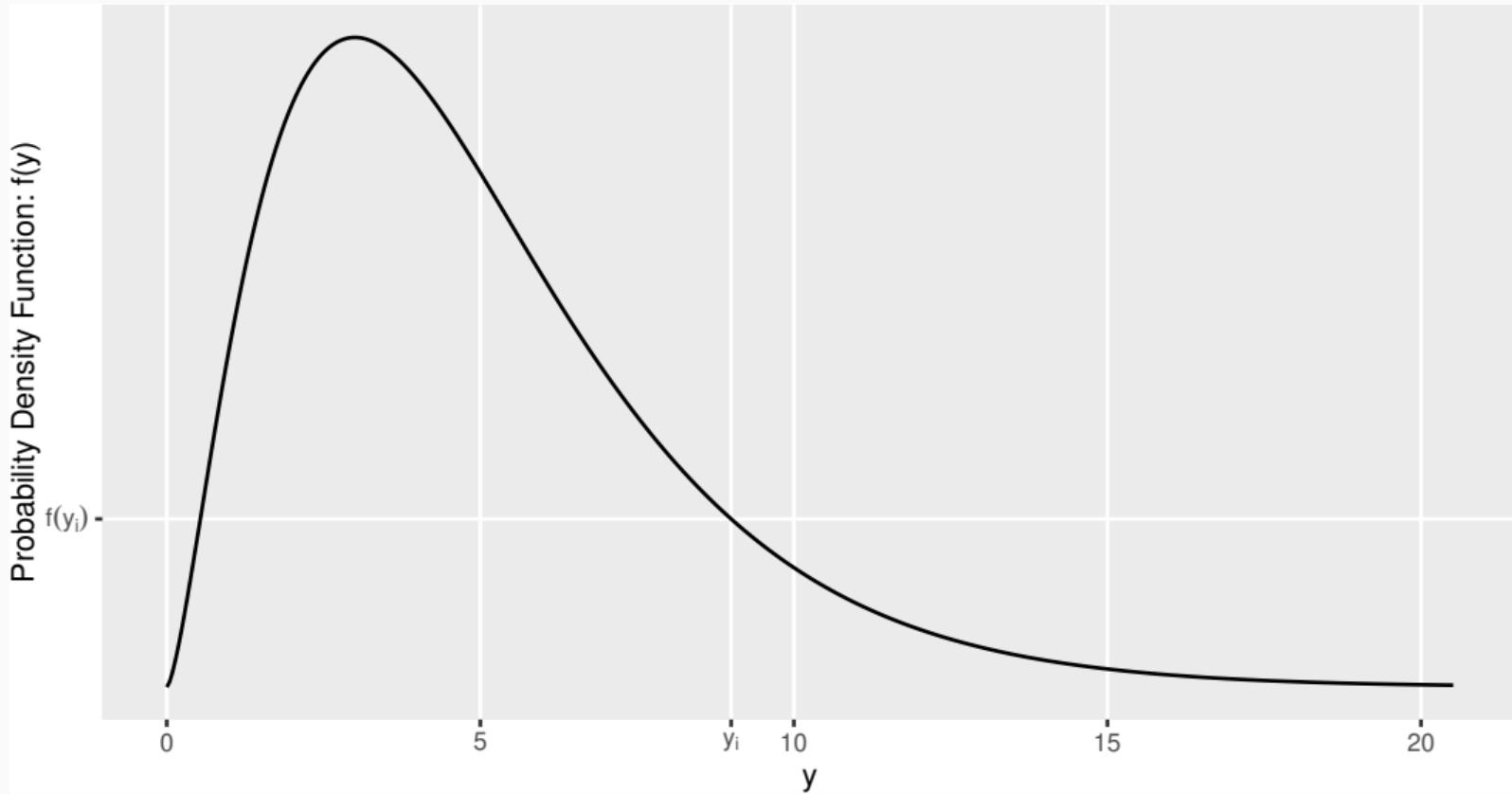
Given a set of observations  $\{y_1, \dots, y_n\}$  and a generalized probability density  $f$ , the **anomaly score** of  $y_i$  wrt  $f$  is

$$p_i = \mathbb{P}(f(Y) \leq f(y_i))$$

where  $Y$  has density  $f$ . An observation is an **anomaly** wrt  $f$  if  $p_i < \alpha$  for some threshold  $\alpha > 0$ .

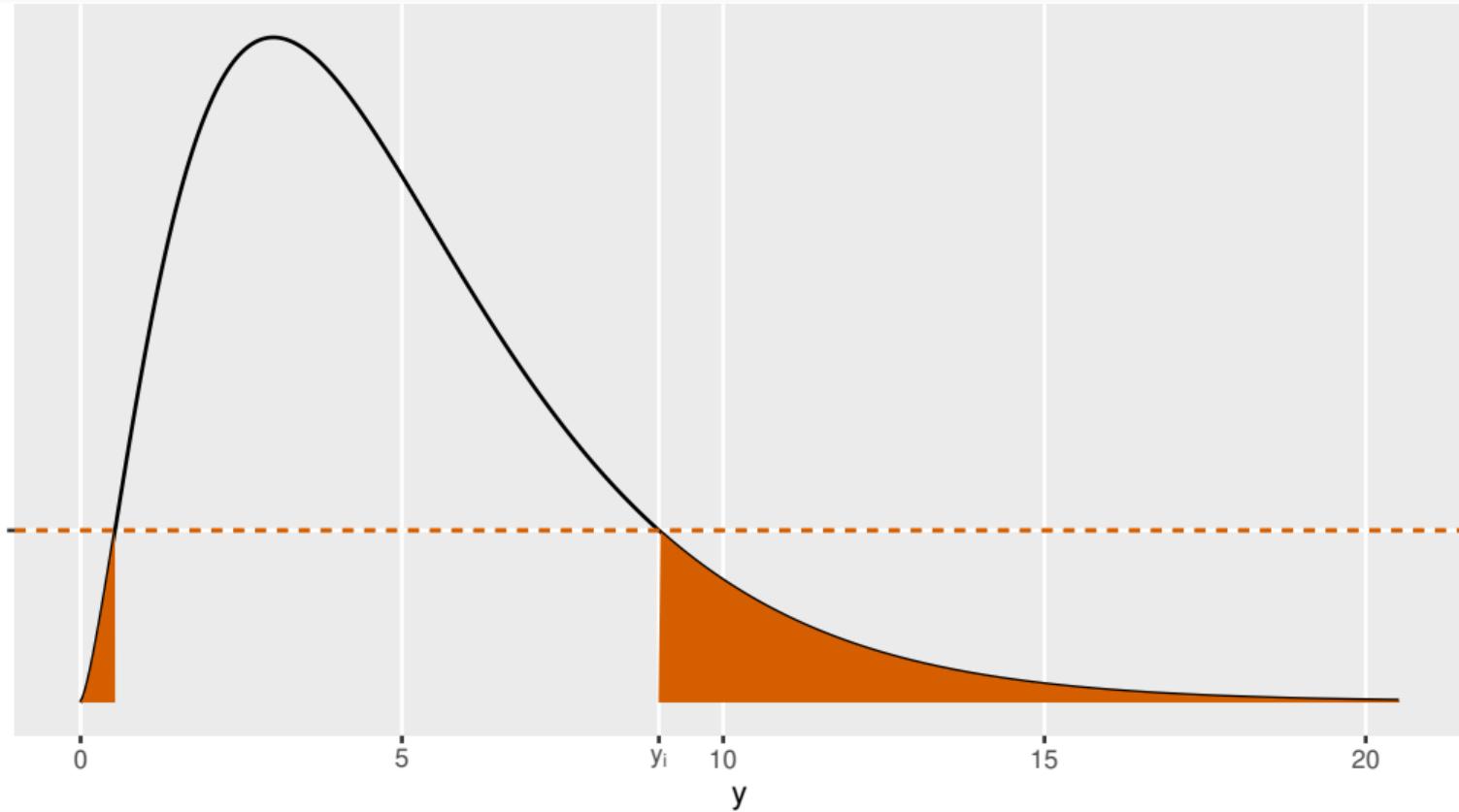
- $y_i$  can be a scalar, vector or a more complex object
- $f$  can be a conditional density, and can be known, assumed or estimated

# Definitions of anomalies



# Definitions of anomalies

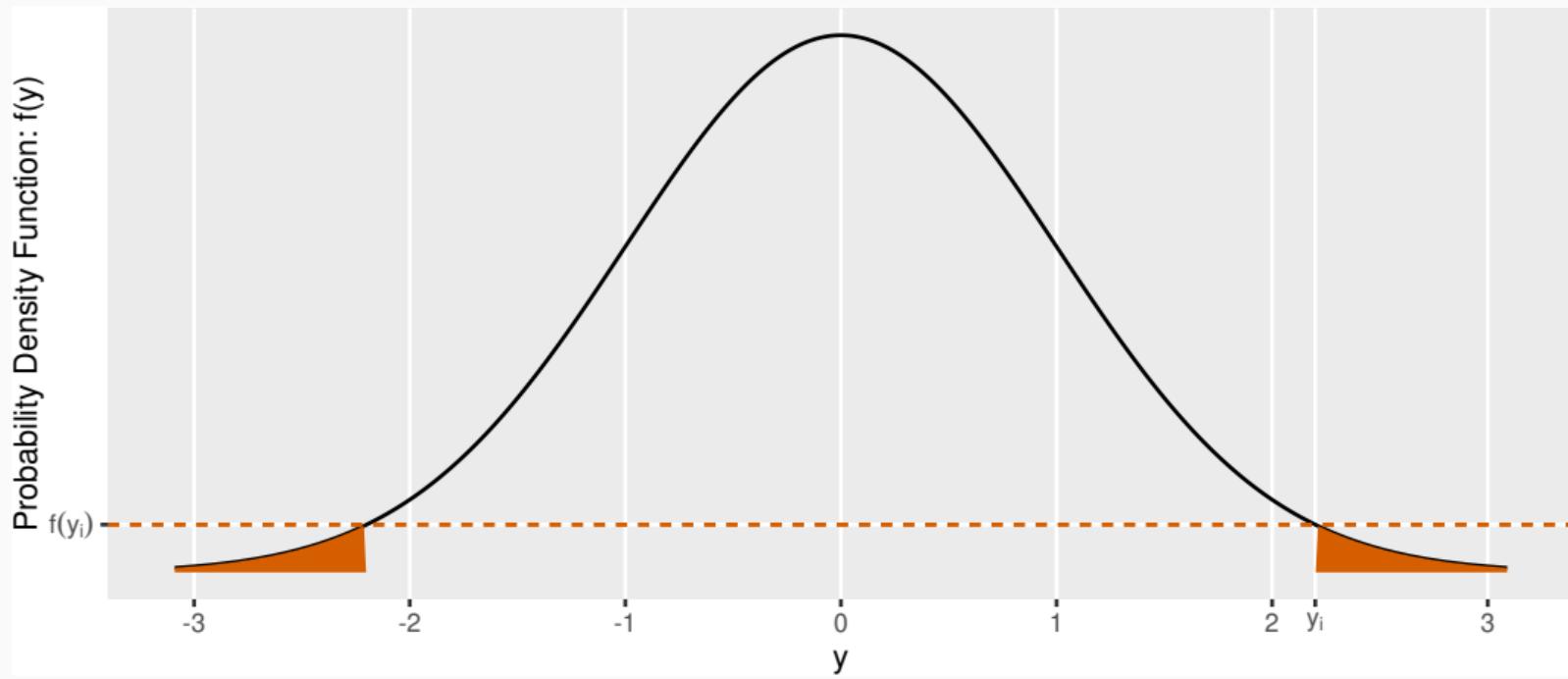
Probability Density Function:  $f(y)$



# Anomaly detection: Normal distribution

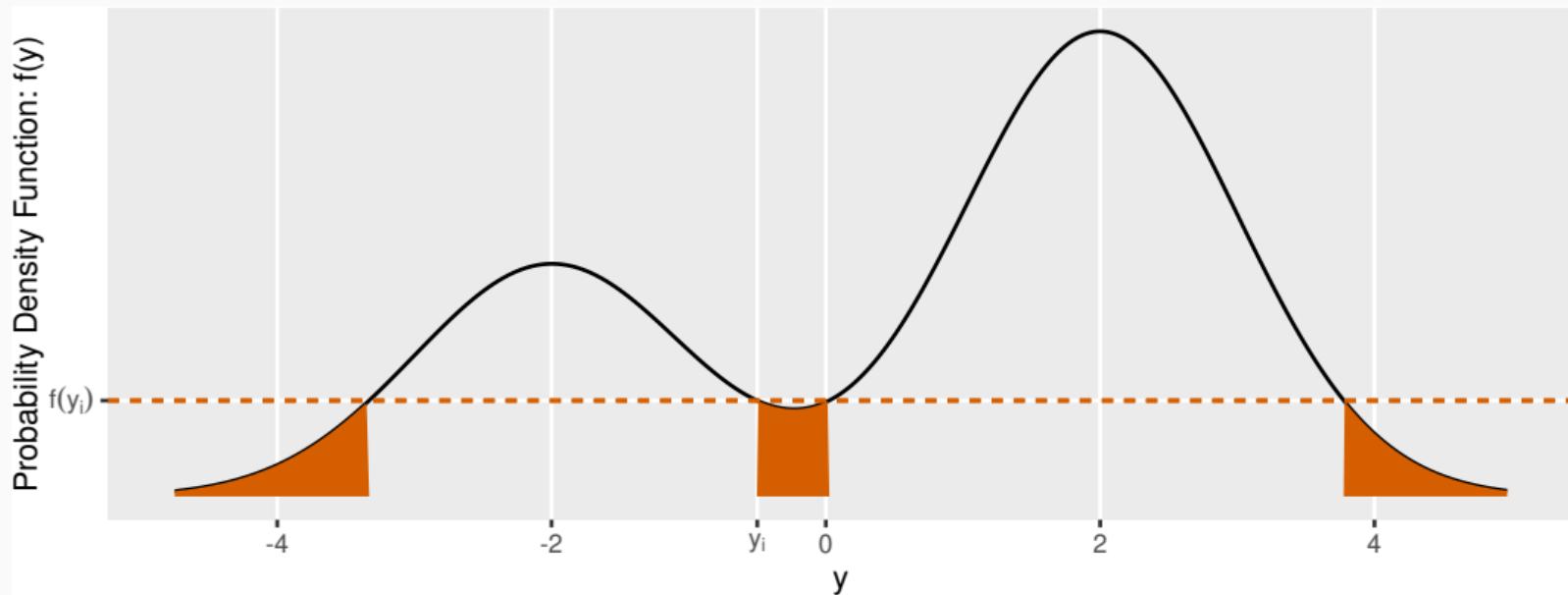
If  $f \sim N(\mu, \sigma^2)$ , then  $p_i = 2 [1 - \Phi(|y_i - \mu|/\sigma)]$

Equivalent to a two-sided p-value from a z-score test.



# Anomaly detection: Highest density regions

HDR with probability  $1 - \alpha$  is  $R_\alpha = \{y : f(y) \geq c_\alpha\}$  where  $c_\alpha$  is largest constant s.t.  $\mathbb{P}(Y \in R_\alpha) \geq 1 - \alpha$ .  
An observation is an anomaly if  $y_i \notin R_\alpha$ .



# Outline

1 Anomalies

2 Extreme surprisals

3 Lookout algorithm

4 Conclusions

# Surprises

## Definition: Surprisal

The **surprisal** of an observation  $y_i$  drawn from a probability distribution with density  $f$  is defined as

$$s_i = -\log f(y_i)$$

- Better known as “log scores” in statistics.
- “Surprisal” coined by Tribus (1961).
- Expected surprisal = entropy of random variable
- Sum of surprisals = negative log likelihood

# Anomaly detection using surprisals

Let  $G(s) = \mathbb{P}(S \leq s)$  be the **surprisal distribution** where  $S = -\log f(Y)$  and  $Y$  has density  $f$ .

$$G(s) = \mathbb{P}(-\log f(Y) \leq s) = \mathbb{P}(f(Y) \geq e^{-s})$$

Then  $p_i = 1 - G(s_i)$ .

# Anomaly detection using surprisals

Let  $G(s) = \mathbb{P}(S \leq s)$  be the **surprisal distribution** where  $S = -\log f(Y)$  and  $Y$  has density  $f$ .

$$G(s) = \mathbb{P}(-\log f(Y) \leq s) = \mathbb{P}(f(Y) \geq e^{-s})$$

Then  $p_i = 1 - G(s_i)$ .

- An anomaly is an extreme value of the surprisal distribution.
- It is not necessarily extreme in the sample space of  $f$ .

# Three-type theorem for surprises

**A1: Sub-Gaussian:**  $S = -\log f(Y)$  satisfies, for all  $\lambda \in \mathbb{R}$ , and some  $\nu > 0$ ,  $\mathbb{E} \exp\{\lambda(S - \mathbb{E}[S])\} \leq \exp\{\lambda^2 \nu^2 / 2\}$ .

**A2: Sub-exponential:**  $S$  is sub-exponential with parameters  $\nu$  and  $b$ , i.e.,  $\mathbb{E} \exp\{\lambda(S - \mathbb{E}[S])\} \leq \exp\{\lambda^2 \nu^2 / 2\}$  for all  $|\lambda| < 1/b$ .

**A3: Polynomial:**  $|S|$  has polynomial moments of order  $p \geq 1$ ; i.e.,  $\mathbb{E}[|S|^p] \leq C^p$  for some  $C > 0$  such that  $C^p - 1 > 0$ .

- A1 satisfied when  $f$  has bounded support
- A2 satisfied when  $\log f$  unbounded below, and light tails (e.g., Gaussian)
- A3 satisfied when  $f$  has heavy tails (e.g., t with df  $\geq 3$ )

# Three-type theorem for surprises

Let  $y_1, \dots, y_n$  be an iid sequence from density  $f$ ,  $s_i = -\log f(y_i)$ ,  $M_n = \max\{s_1, \dots, s_n\}$ , and  $S = -\log f(Y)$  where  $Y \sim f$ .

1 Under A1:

$$\sup_{s:s>0} \left| \mathbb{P} \left\{ |M_n - \mathbb{E}[S]| \geq \sqrt{2\nu^2 s} + \sqrt{2\nu^2 \log(2n)} \right\} - e^{-s} \right| = o(1).$$

2 Under A2:

$$\sup_{s:s>1/b} \left| \mathbb{P} \left\{ |M_n - \mathbb{E}[S]| \geq (2b)s + (2b)\log(2n) \right\} - e^{-e^{-s}} \right| = o(1).$$

3 Under A3:

$$\sup_{s:s>c} \left| \mathbb{P} \left\{ |M_n - \mathbb{E}[S]| \geq (Csn^{1/p}) \right\} - e^{-s^{-p}} \right| = o(1).$$

# Three-type theorem for surprisals

- If surprisal has Gaussian like tails, then maximum surprisal is a reversed Weibull;
- If surprisal only has an exponential tail, then maximum surprisal is Gumbel;
- If surprisal only has a polynomial moment, then maximum surprisal is Fréchet.

## Corollary (due to Pickands theorem)

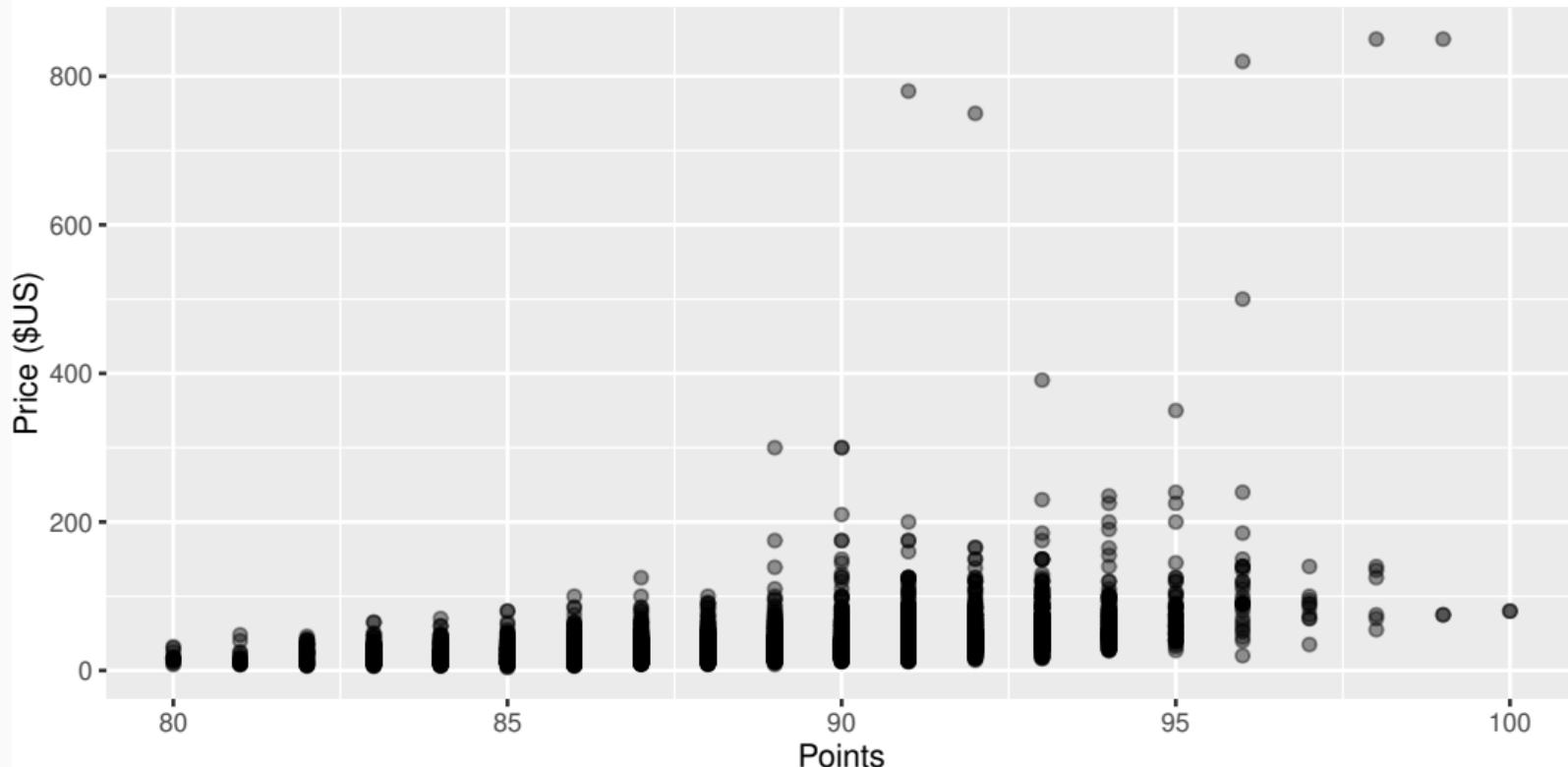
upper tail of the surprisal distribution can be approximated by a Generalized Pareto Distribution, even if the assumed density  $f$  is incorrect, provided one of A1–A3 is satisfied.

# Surprises and EVT

- Suppose we have  $n$  iid observations  $Y_1, \dots, Y_n$  and a density  $f$ .
- Let  $S_i = -\log f(Y_i)$  be the surprisal of  $Y_i$  wrt  $f$
- Then  $S_1, \dots, S_n$  are iid from the surprisal distribution  $G(s) = \mathbb{P}(S \leq s)$ .
- For almost all  $f$ , we can approximate the upper tail of  $G$  by a Generalized Pareto Distribution fitted to the top  $1 - \beta$  of the surprisal values.
- In practice, we typically use  $\beta = 0.9$ .

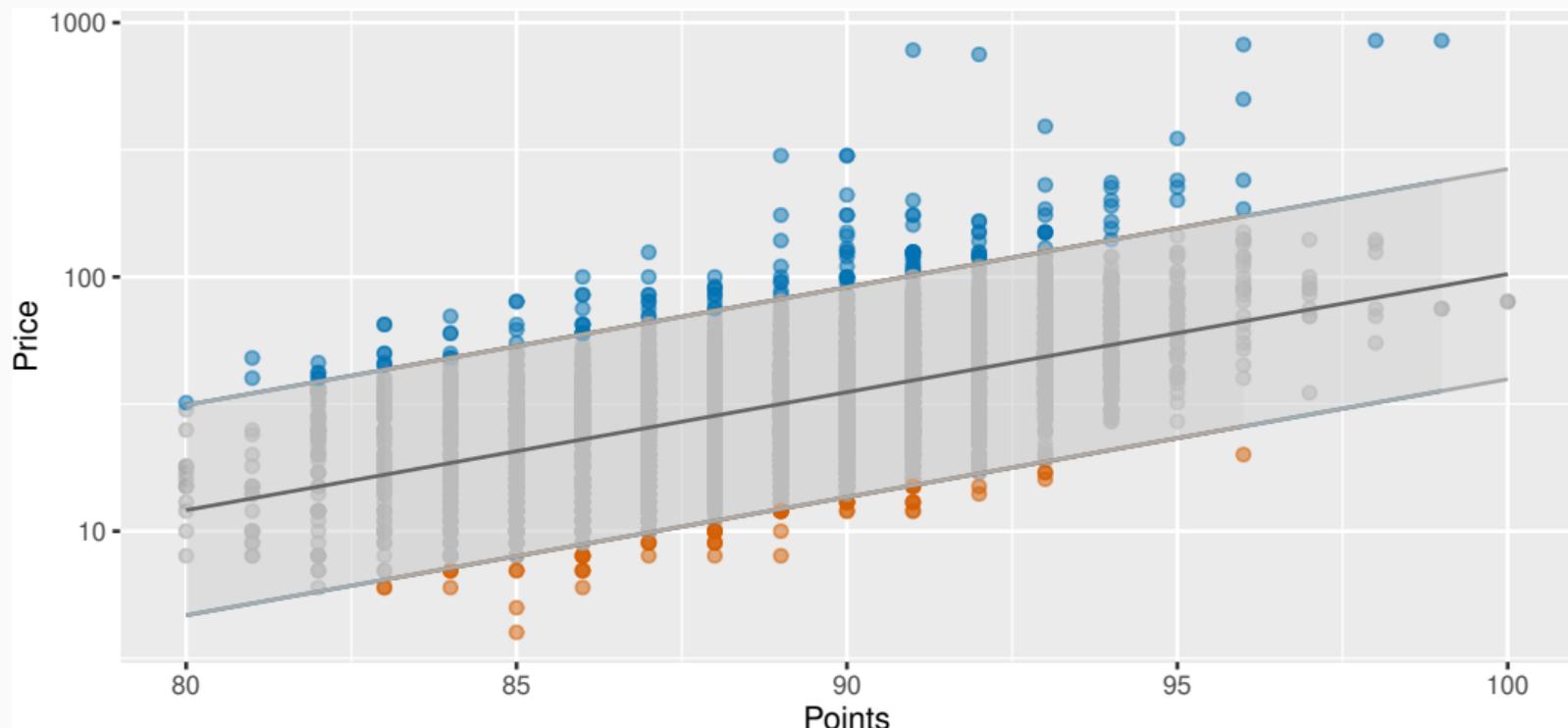
# Application to wine quality and prices

Reviews of 4496 Shiraz/Syrah wines from 'Wine Enthusiast', 15 June 2017



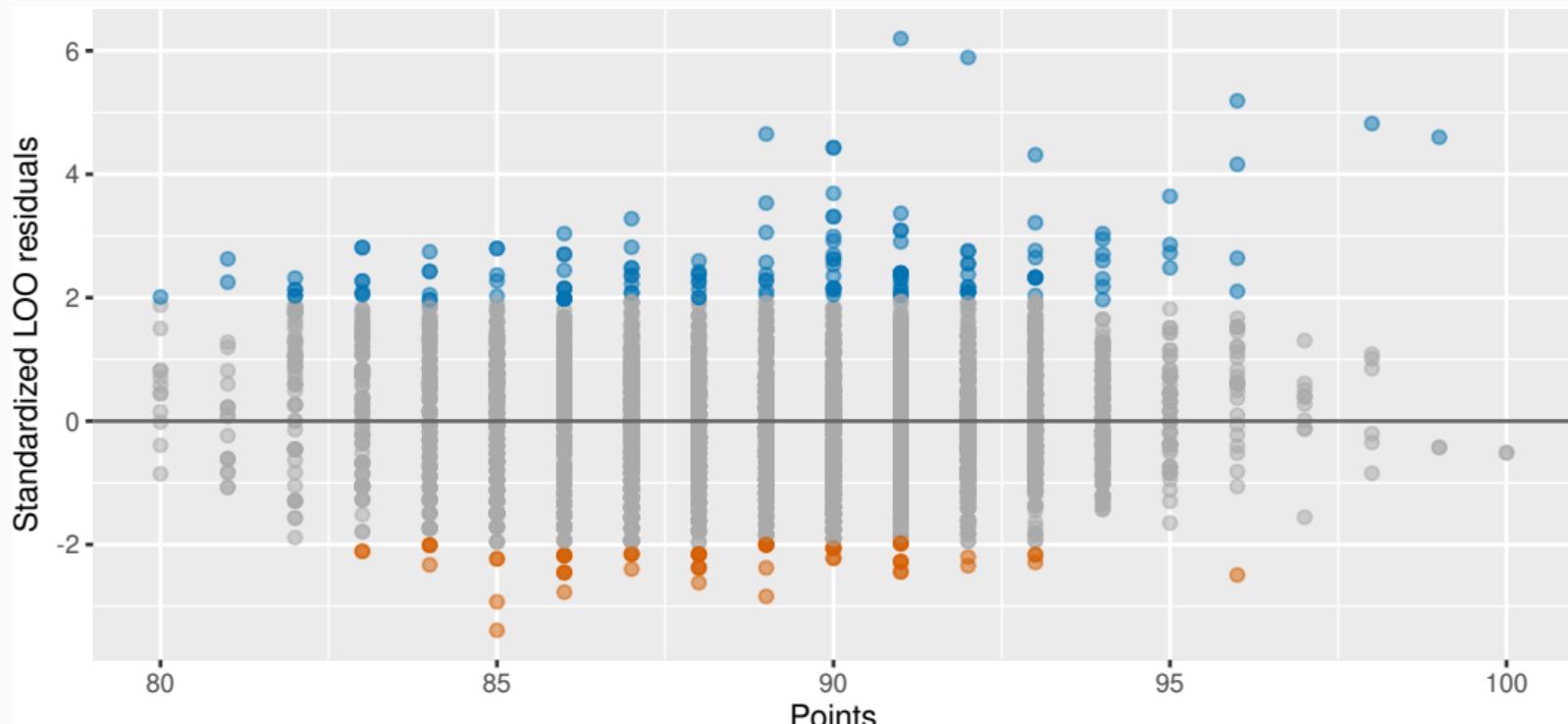
# Application to wine quality and prices

Proposed model:  $\log \text{Price} | \text{Points} \sim N(a + b\text{Points}, \sigma^2)$ .



# Application to wine quality and prices

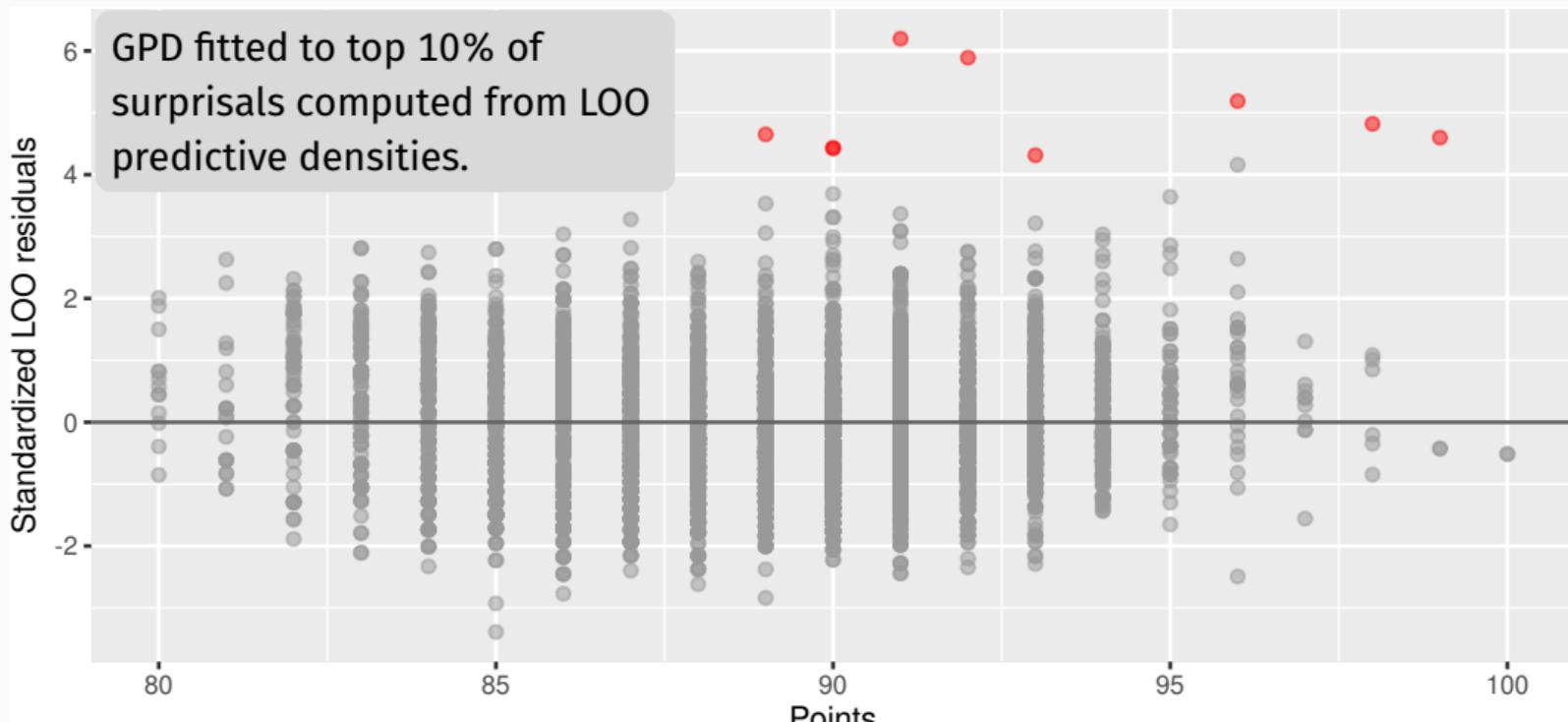
Proposed model:  $\log \text{Price} | \text{Points} \sim N(a + b\text{Points}, \sigma^2)$ .



# Application to wine quality and prices

Proposed model:  $\log \text{Price} | \text{Points} \sim N(a + b\text{Points}, \sigma^2)$ .

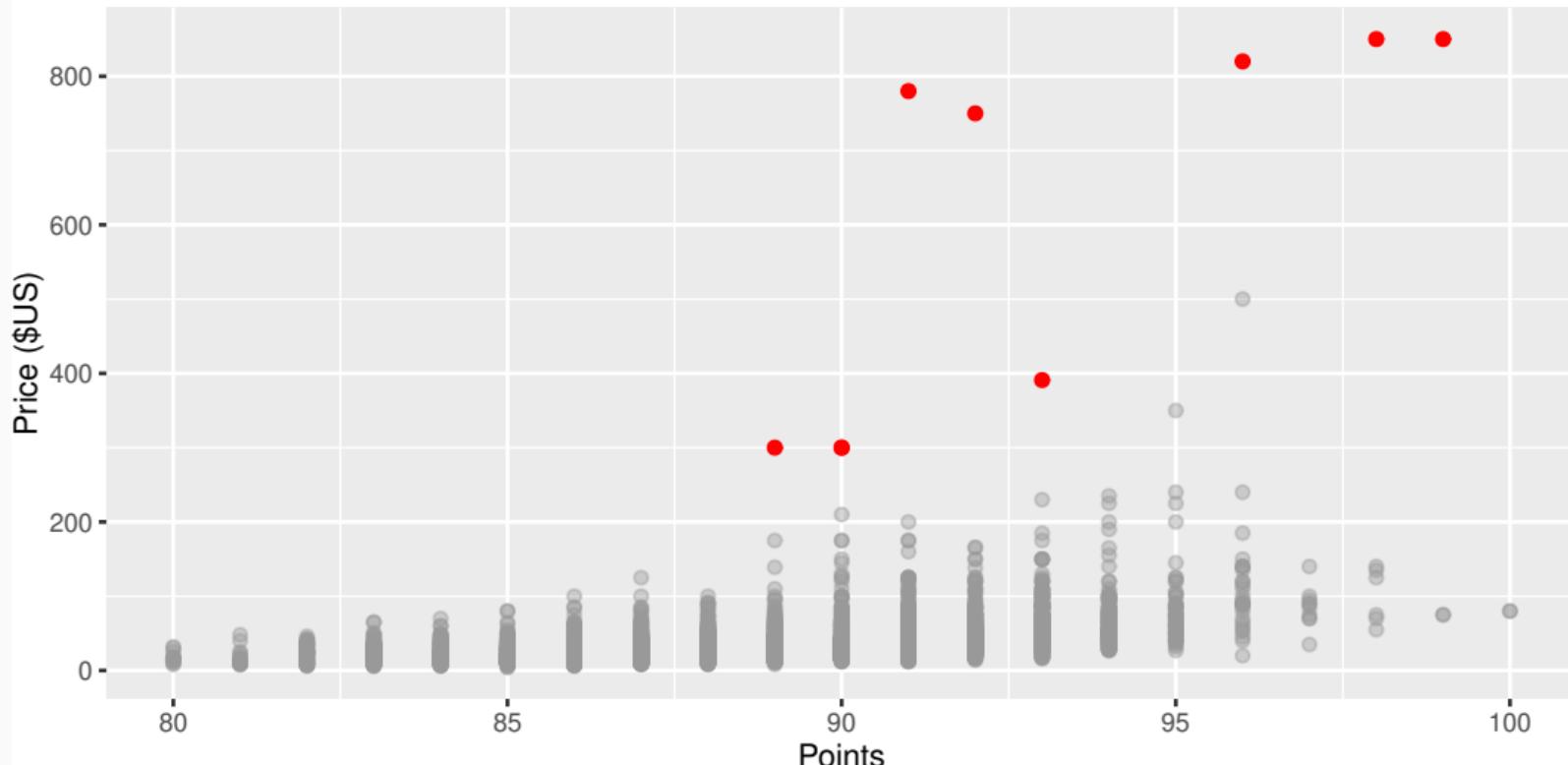
$\alpha = 0.001$



# Application to wine quality and prices

Reviews of 4496 Shiraz/Syrah wines from 'Wine Enthusiast', 15 June 2017

$\alpha = 0.001$



# Application to wine quality and prices

Anomalies detected ( $\alpha = 0.001$ ):

Area	Winery	Year	Points	Price
South Australia	Henschke	2009	91	780
California	Law	2013	92	750
South Australia	Henschke	2010	96	820
South Australia	Penfolds	2008	98	850
Tuscany	Tua Rita	2011	89	300
South Australia	Penfolds	2010	99	850
Tuscany	Tua Rita	2013	90	300
Tuscany	Tua Rita	2012	90	300
Rhône Valley	Domaine Jean-Michel Gerin	2013	93	391

# Univariate experiment

What happens when the distribution used to compute surprisals is mis-specified?

## Data $N(0,1)$

- 1000 observations from a  $N(0,1)$  distribution
- Surprisals computed using a  $t(4)$  distribution.
- Estimate surprisal probabilities using  $N(0,1)$ ,  $t(4)$ , GPD

# Univariate experiment

What happens when the distribution used to compute surprisals is mis-specified?

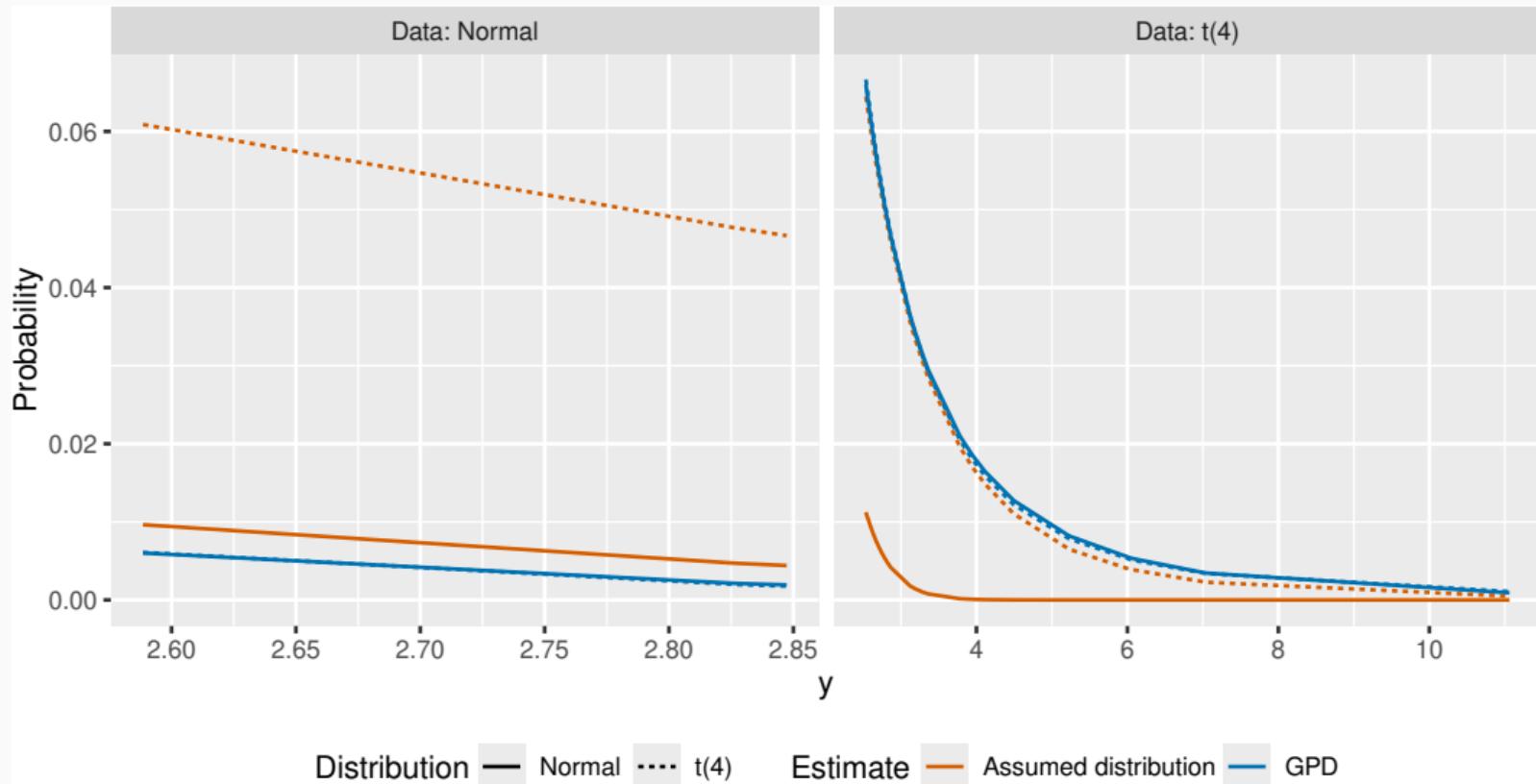
## Data $N(0,1)$

- 1000 observations from a  $N(0,1)$  distribution
- Surprisals computed using a  $t(4)$  distribution.
- Estimate surprisal probabilities using  $N(0,1)$ ,  $t(4)$ , GPD

## Data $t(4)$

- 1000 observations from a  $t(4)$  distribution
- Surprisals computed using a  $N(0,1)$  distribution.
- Estimate surprisal probabilities using  $N(0,1)$ ,  $t(4)$ , GPD

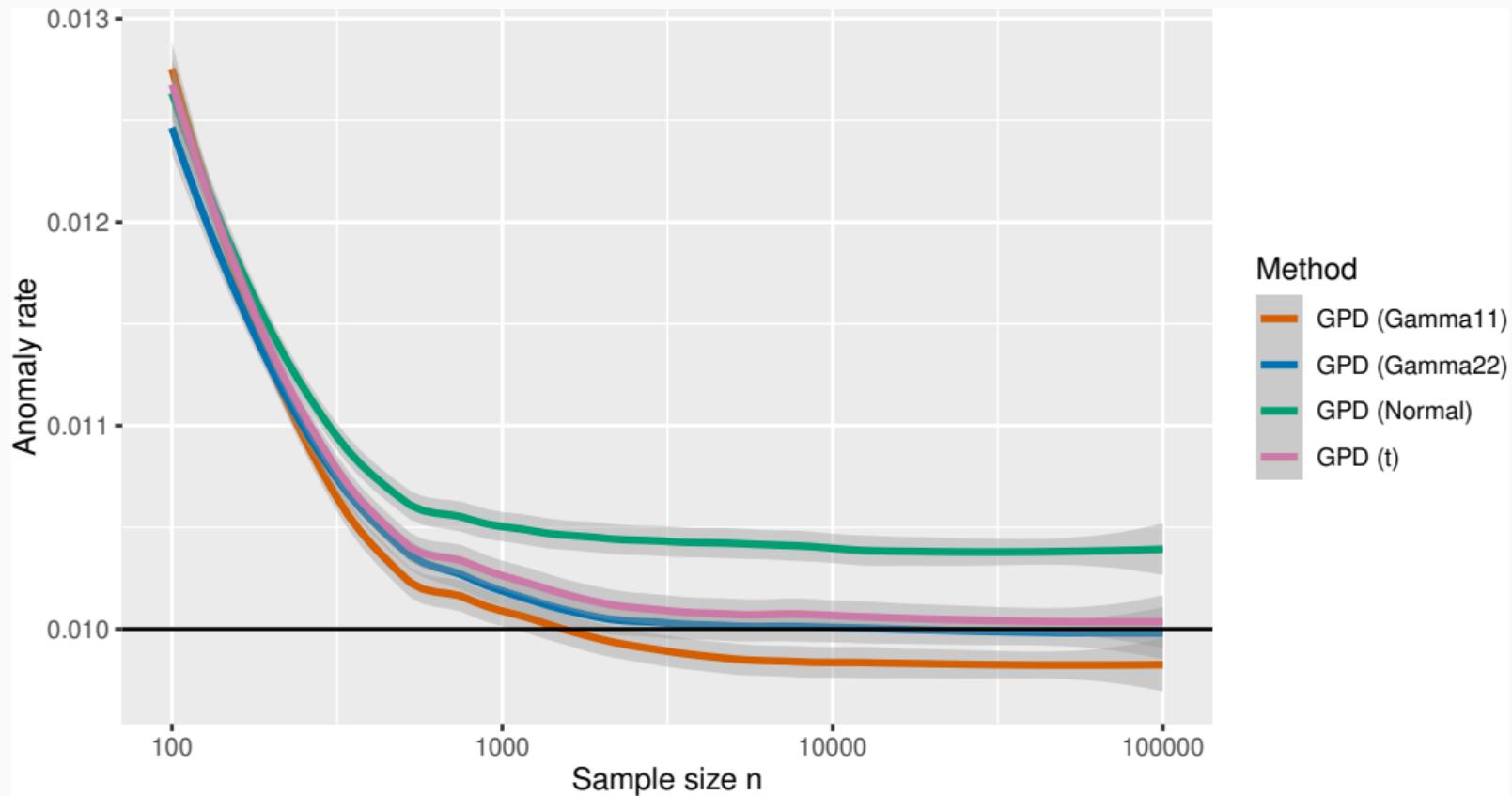
# Univariate experiment



# Bivariate experiment

- Data: 2 independent Gamma(2,2) variables
- Sample size:  $n = 100, \dots, 100000$
- Surprisals computed using:
  - ▶ Gamma(2,2) x 2
  - ▶ Gamma(1,1) x 2
  - ▶ Bivariate normal with correct mean and variance
  - ▶ Bivariate non-central t(4) with correct mean

# Bivariate experiment



# Outline

1 Anomalies

2 Extreme surprisals

3 Lookout algorithm

4 Conclusions

# Kernel density estimation

Observations:  $\mathbf{y}_i \in \mathbb{R}^m$  for  $i \in \{1, \dots, n\}$ .

## KDE

$$\hat{f}(\mathbf{y}) = \frac{1}{n} \sum_{i=1}^n |\mathbf{H}|^{-1/2} K(\mathbf{H}^{-1/2}(\mathbf{y} - \mathbf{y}_i)),$$

- $K$  is a square-integrable spherically-symmetric function, bounded below by 0, with a finite second-order moment and unit integral.
- $\mathbf{H}$  is a symmetric  $m \times m$  positive-definite matrix.

# Kernel density estimation

Observations:  $\mathbf{y}_i \in \mathbb{R}^m$  for  $i \in \{1, \dots, n\}$ .

## KDE

$$\hat{f}(\mathbf{y}) = \frac{1}{n} \sum_{i=1}^n |\mathbf{H}|^{-1/2} K(\mathbf{H}^{-1/2}(\mathbf{y} - \mathbf{y}_i)),$$

- $K$  is a square-integrable spherically-symmetric function, bounded below by 0, with a finite second-order moment and unit integral.
- $\mathbf{H}$  is a symmetric  $m \times m$  positive-definite matrix.

## LOO KDE values

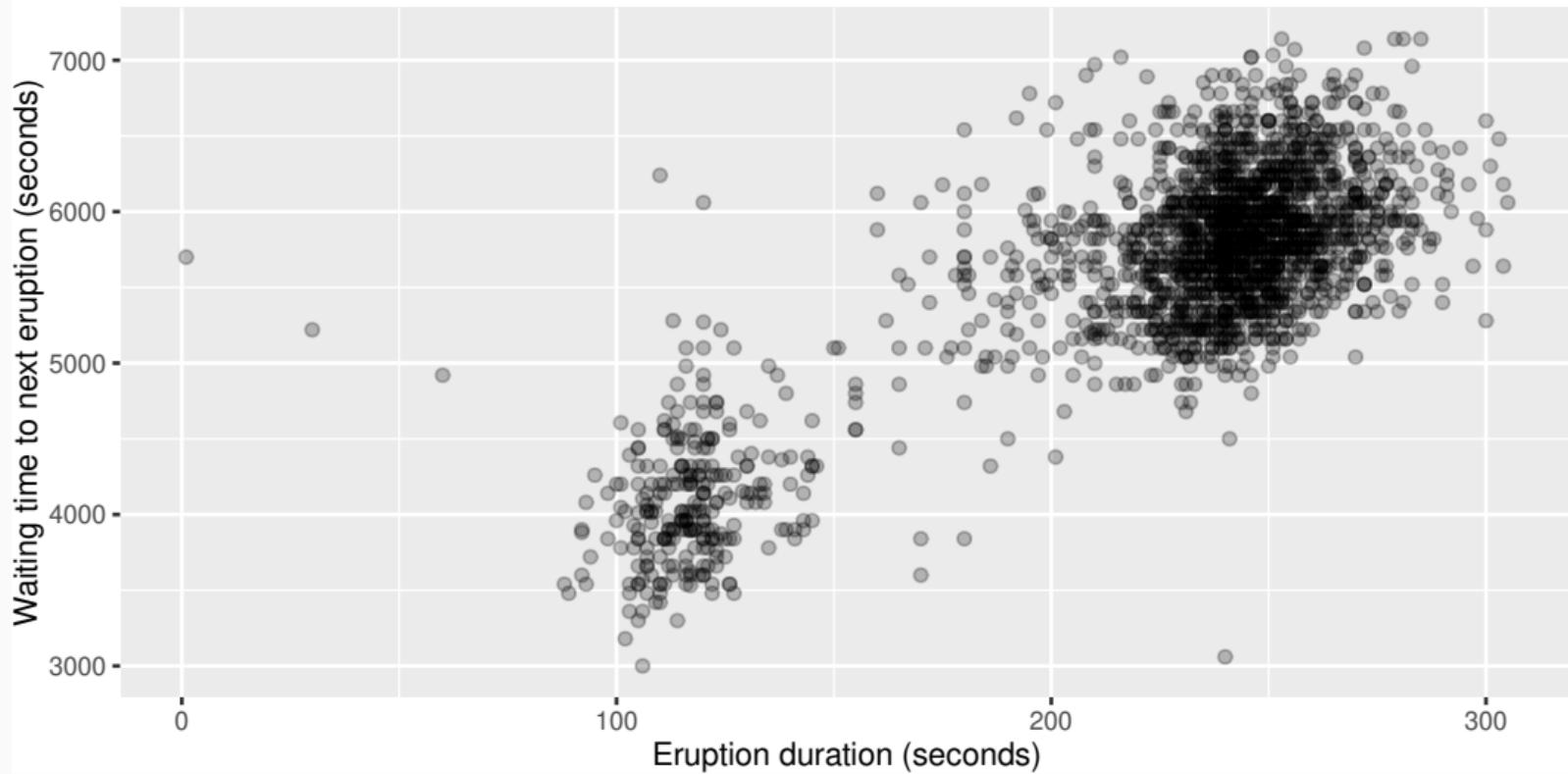
$$f_{-i} = \frac{1}{n-1} (n\hat{f}(\mathbf{y}_i) - \mathbf{H}^{-1/2}K(\mathbf{0}))$$

# Lookout algorithm

- 1  $\mathbf{Y}$  = data matrix with rows  $\mathbf{y}_1, \dots, \mathbf{y}_n$ .
- 2  $\hat{\Sigma}$  = orthogonalized Gnanadesikan-Kettenring estimate of  $\text{Cov}(\mathbf{Y})$ , with eigendecomposition  $\hat{\Sigma} = \mathbf{U}\mathbf{D}\mathbf{U}^\top$ .
- 3 Rotate and scale the data:  $\mathbf{Z} = \mathbf{U}\mathbf{Y}$ .
- 4 Compute persistence homology barcode of  $\mathbf{Z}$  for dim zero using Vietoris-Rips diameter; obtain ordered death diameters  $\{d_i\}_{i=1}^n$ .
- 5  $d_\gamma^* = \gamma$  sample quantile computed from  $\{d_i\}_{i=1}^n$ .
- 6 Compute kde:  $f_i = \hat{f}(\mathbf{z}_i)$ ,  $i = 1, \dots, n$ , where  $\mathbf{H} = (d_\gamma^*)^{2/m} \mathbf{I}_m$ .
- 7 Compute LOO kde values  $f_{-i} = \frac{1}{n-1} (nf_i - \mathbf{H}^{-1/2} K(\mathbf{0}))$ ,  $i = 1, \dots, n$ .
- 8 Fit GPD to largest  $1 - \beta$  of surprisals  $\{-\log f_i\}_{i=1}^n$ , constraining shape parameter to be non-positive.
- 9  $p_i = (1 - \beta)P(-\log f_{-i} | \hat{\mu}, \hat{\sigma}, \hat{\xi})$ ,  $P$  = GPD cdf.

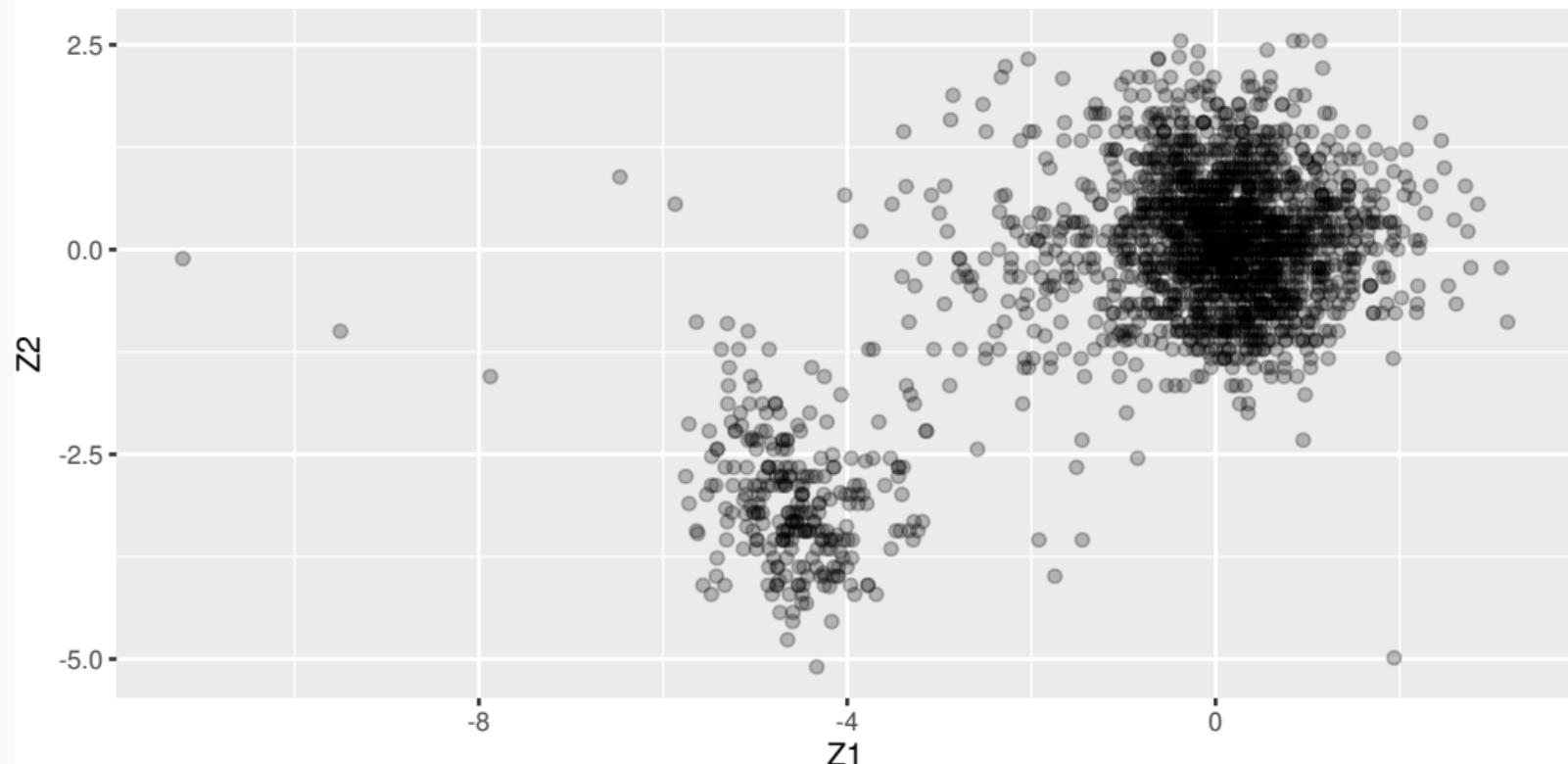
# Old Faithful eruptions

Old Faithful eruptions from 14 January 2017 to 29 December 2023



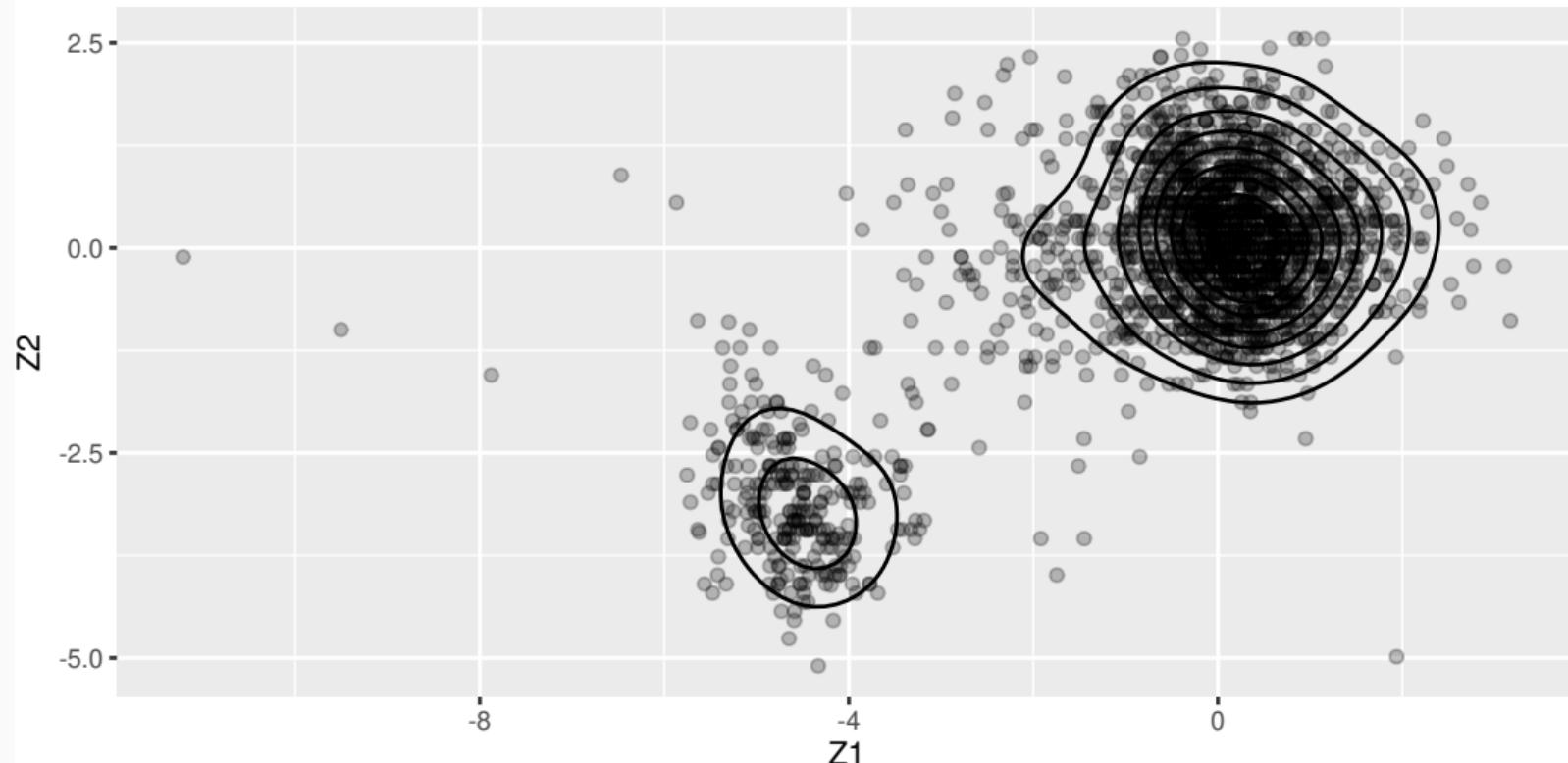
# Old Faithful eruptions

Standardized Old Faithful eruptions from 14 January 2017 to 29 December 2023



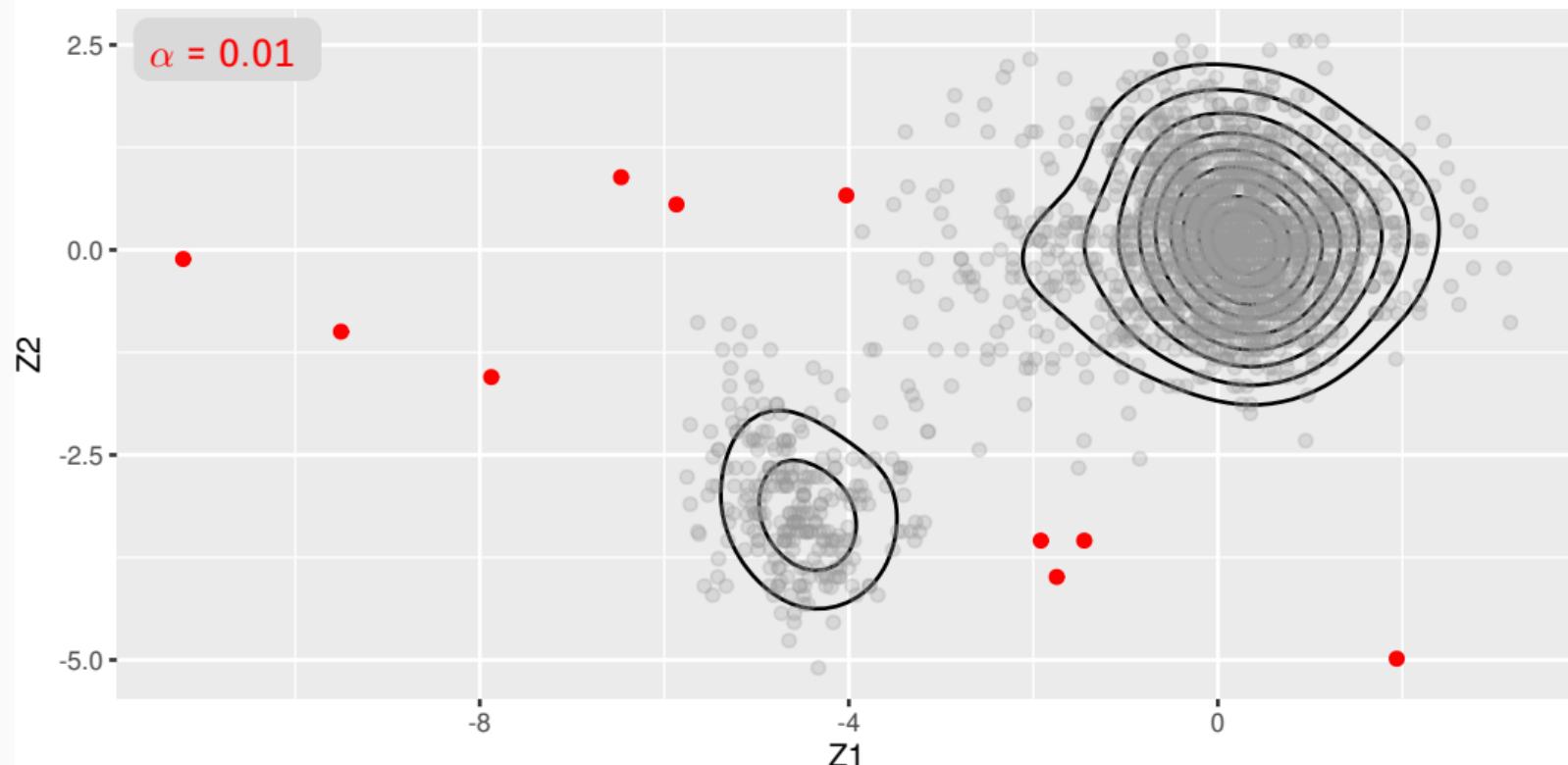
# Old Faithful eruptions

Standardized Old Faithful eruptions from 14 January 2017 to 29 December 2023



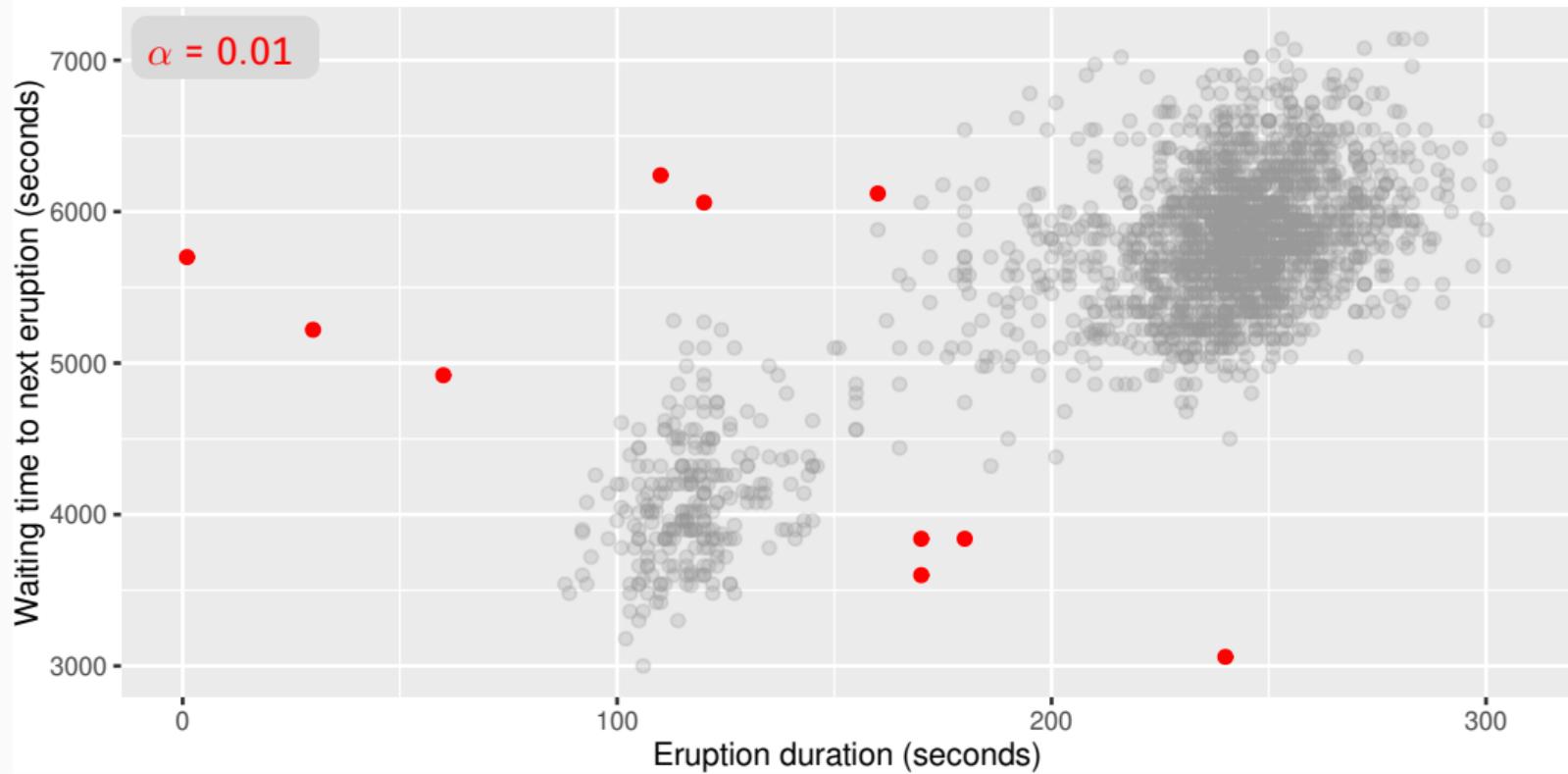
# Old Faithful eruptions

Standardized Old Faithful eruptions from 14 January 2017 to 29 December 2023



# Old Faithful eruptions

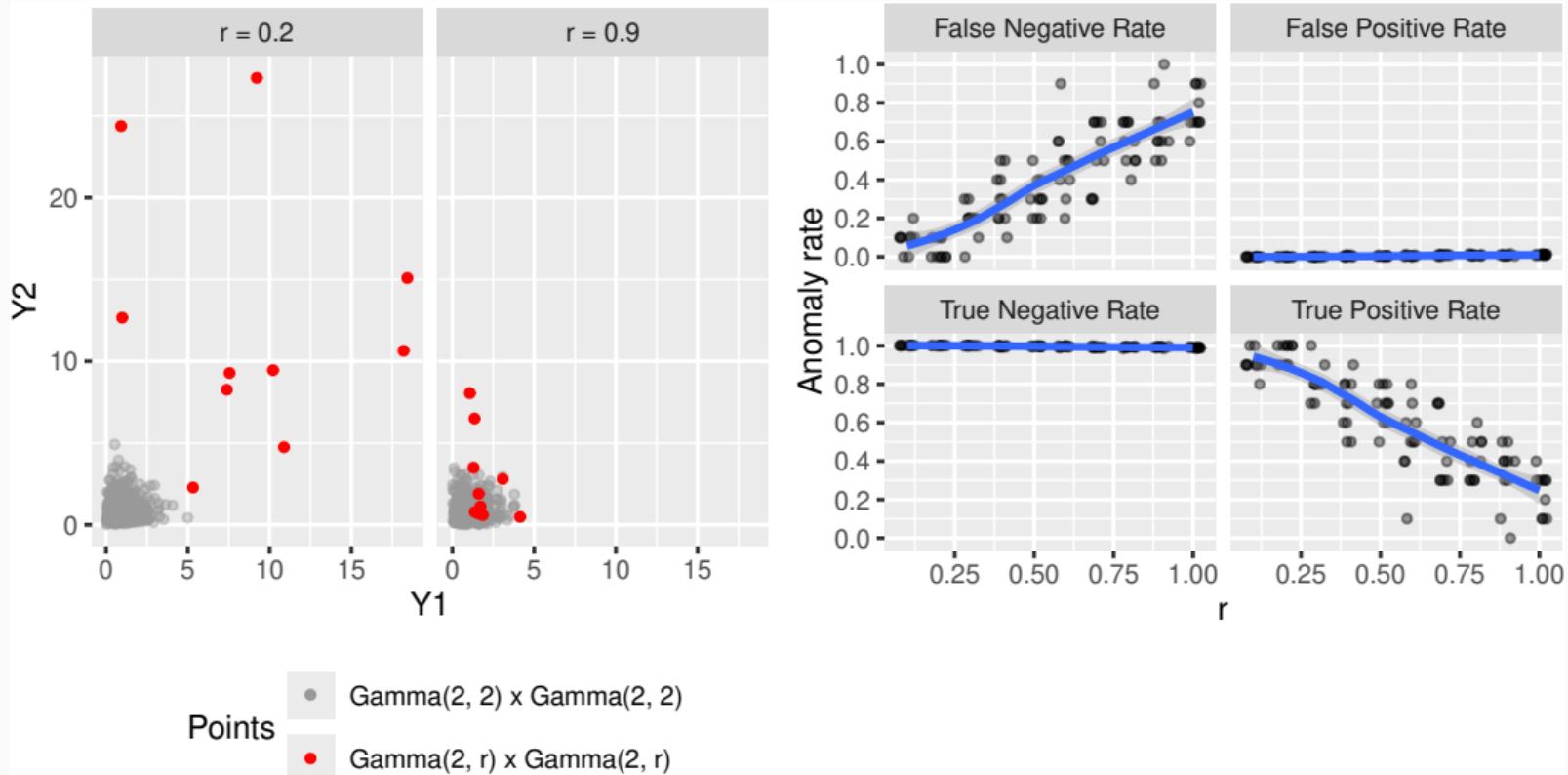
Old Faithful eruptions from 14 January 2017 to 29 December 2023



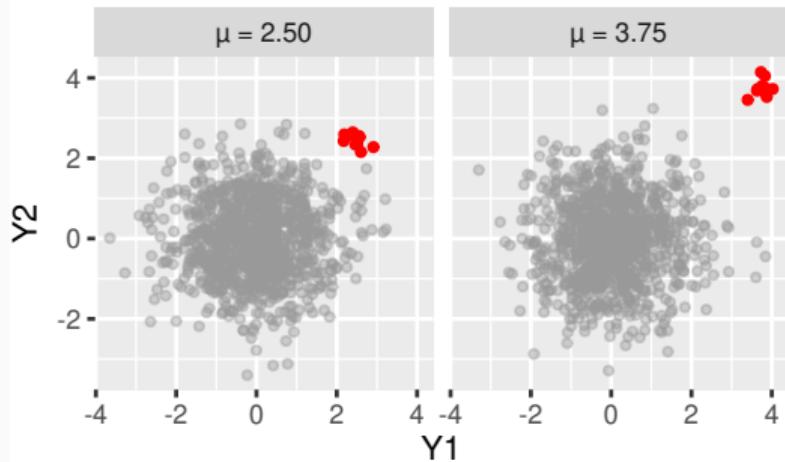
# Old Faithful eruptions

time	recorded_duration	duration	waiting	prob
2018-04-25 19:08:00	1s	1	5700	0.0000000
2022-12-07 17:19:00	~4 30s	30	5220	0.0000000
2023-07-04 12:03:00	~1 minute 55ish seconds	60	4920	0.0000612
2020-09-04 01:38:00	>1m 50s	110	6240	0.0000012
2020-06-01 21:04:00	2 minutes	120	6060	0.0001390
2020-09-16 14:44:00	>2m40s	160	6120	0.0078057
2020-08-31 09:56:00	~2m50s	170	3840	0.0076247
2021-01-22 18:35:00	2m50s	170	3600	0.0029265
2022-11-29 14:51:00	~3m	180	3840	0.0033112
2022-12-03 16:20:00	~4m	240	3060	0.0000000

# Experiment 1

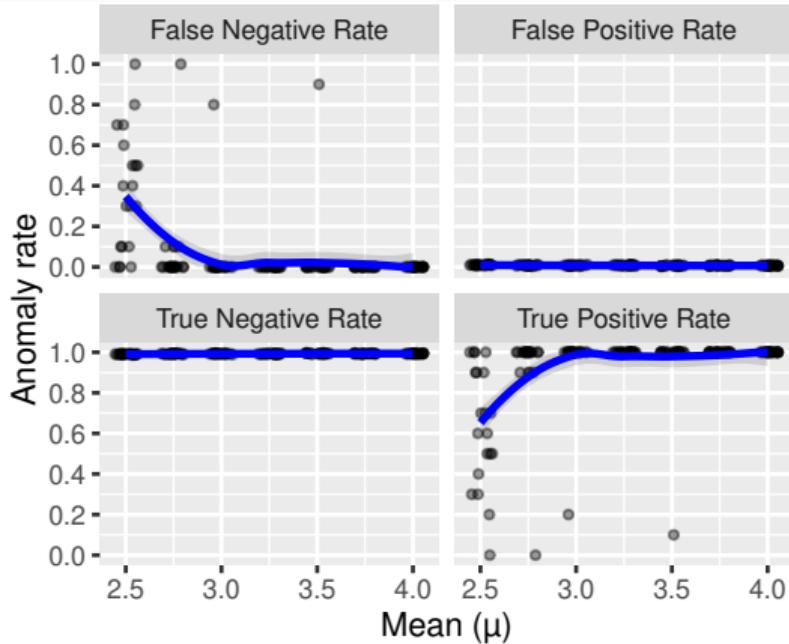


# Experiment 2



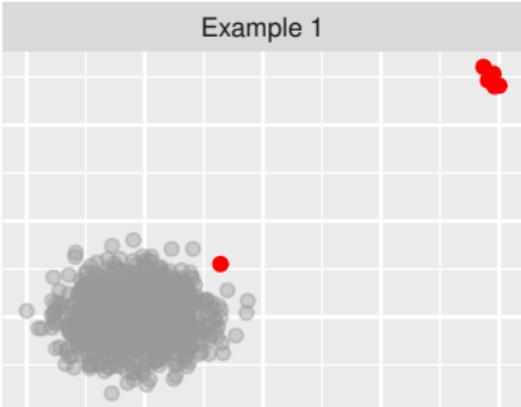
Points

- $N(0, 1) \times N(0, 1)$
- $N(\mu, 0.04) \times N(\mu, 0.04)$

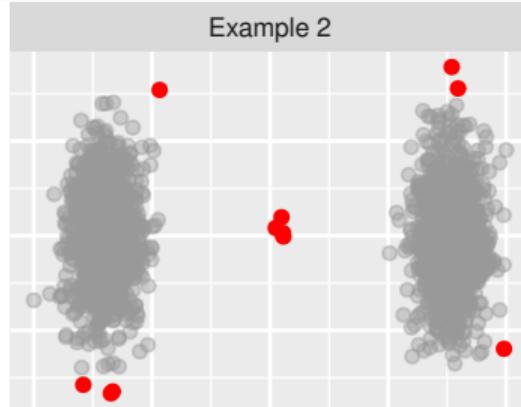


# Experiment 3

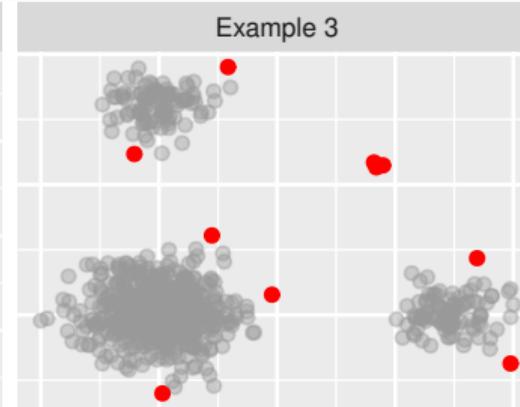
Example 1



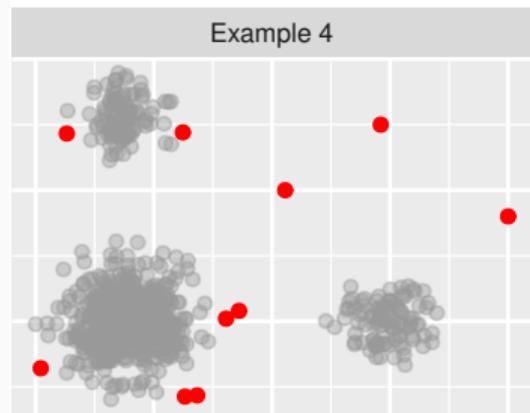
Example 2



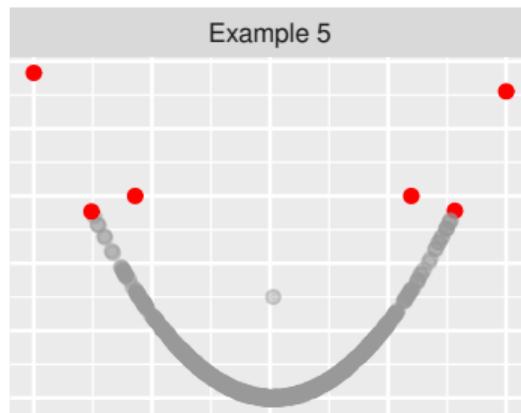
Example 3



Example 4



Example 5



# Outline

1 Anomalies

2 Extreme surprisals

3 Lookout algorithm

4 Conclusions

# Conclusions

- Surprisal-based anomaly detection is a flexible, probabilistic approach that can be used in any context where a probability distribution can be defined on the space of observations
- EVT theory on surprisals due to Hyndman & Frazier (in preparation)
- Original lookout algorithm due to Kandanaarachichi & Hyndman (*JCGS*, 2022). <https://robjhyndman.com/publications/lookout>
- Modified lookout algorithm due to Hyndman, Kandanaarachichi & Turner (in preparation)
- R packages `lookout` and `weird` available on CRAN
- Book in preparation at <https://OTexts.com/weird>
- Slides and links: <https://robjhyndman.com/fsi2025>