

Vector
derivatives

Topics : Anomaly / Novelty detection

"Known"

✓ ① Using existing techniques

② Elliptical-envelope → MV · Gaussian

③ Isolation Forest → DT & RF

④ One-class SVMs → SVMs

⑤ Local outlier factor
↳ KNN & density

Intuition ✓

Math ✓

Code ✓

✓ Toy-datasets

Real-datasets



Anomaly \approx outlier \approx Novelty

$=$

\approx

\downarrow

not-normal datapoints

$=$

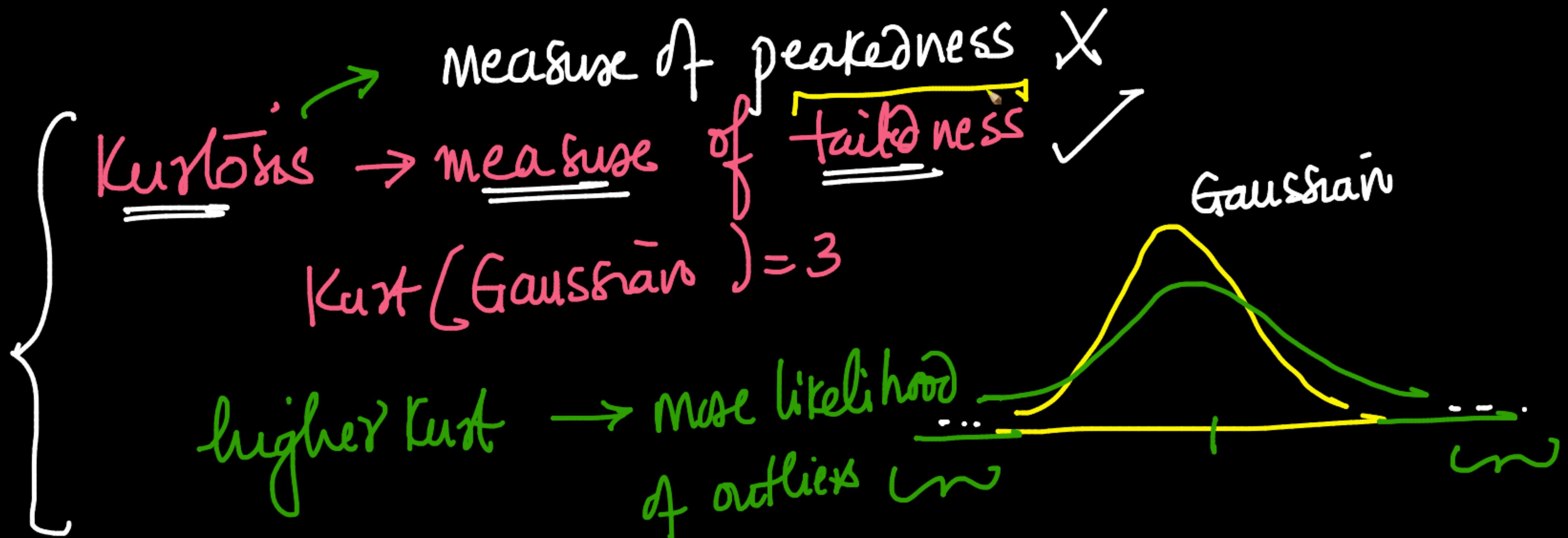
Techniques

[outlier/transliers
Novel]

- ① DBSCAN → density; outliers
- ② RANSAC → Robust → model-Specific outlier
- ③ Boxplot → IQR → Univariate → MAD, Z-score

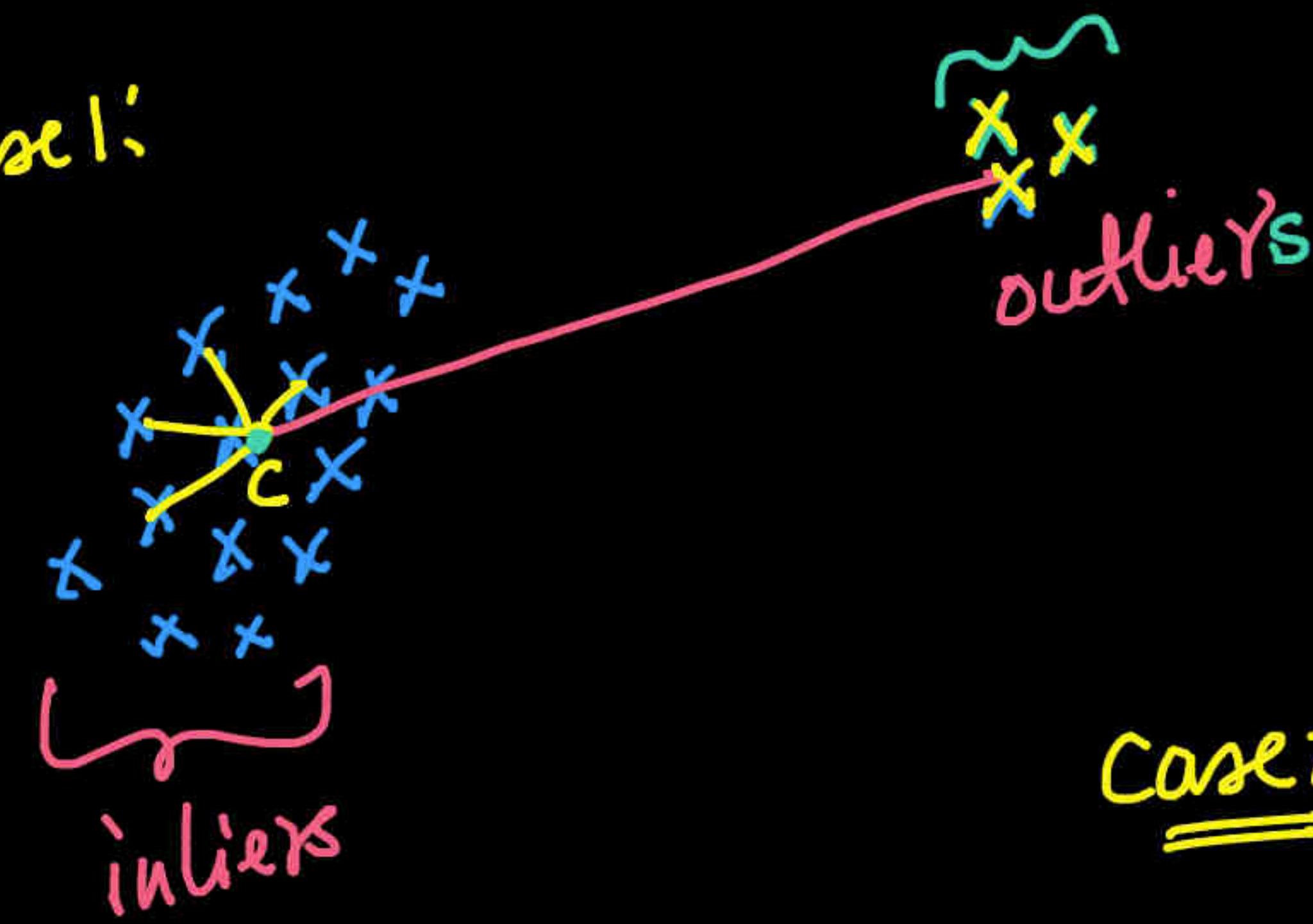
$$\frac{\text{avg-dist}}{\sigma} \quad \frac{x_i - \mu}{\sigma}$$
- ④ KNN → K-nearest neighbor
- ⑤ Viz. of e_i^2 : errors → model specific (bias-variance)

$$\frac{kurtosis}{\sigma}$$



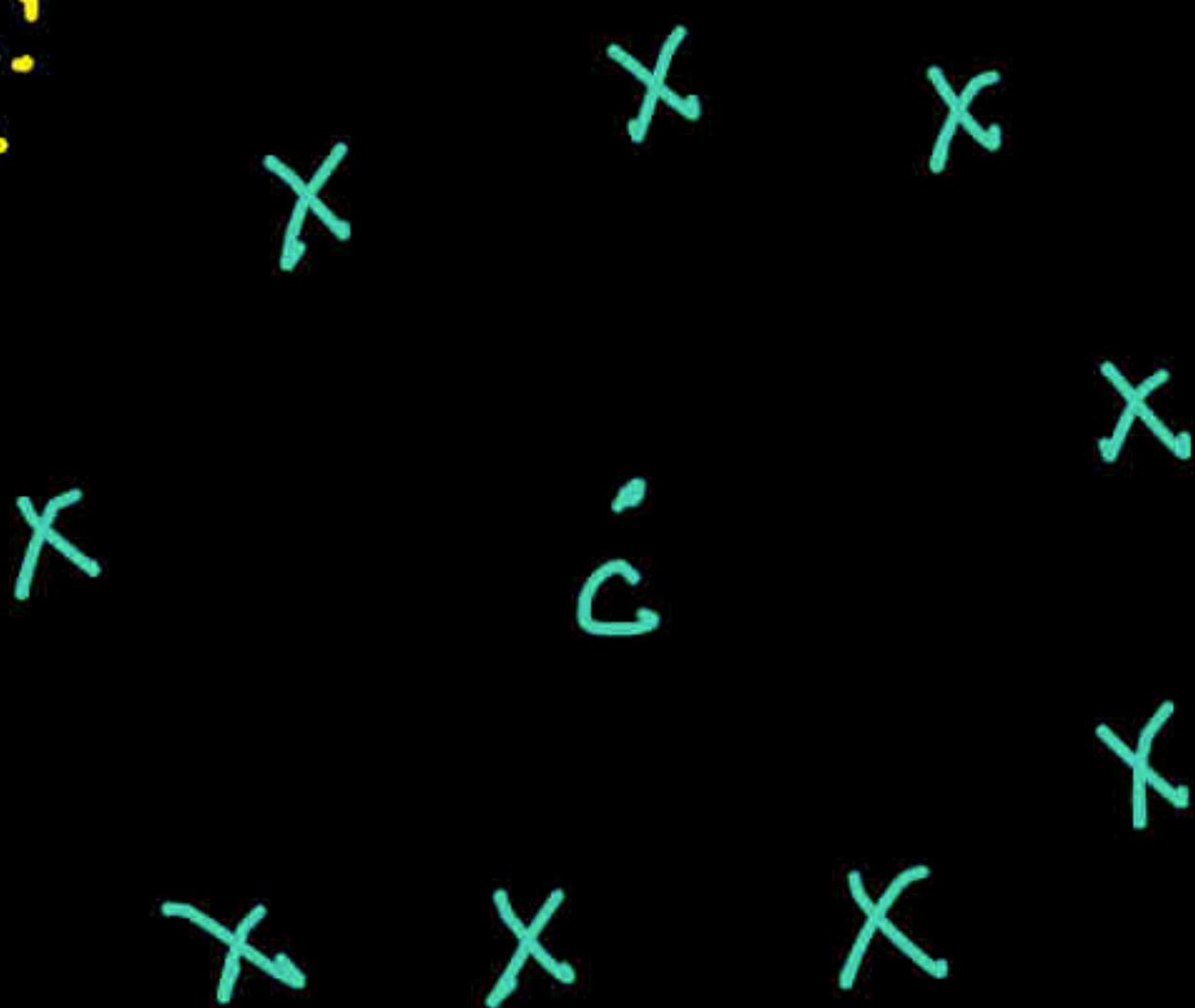


Case 1:

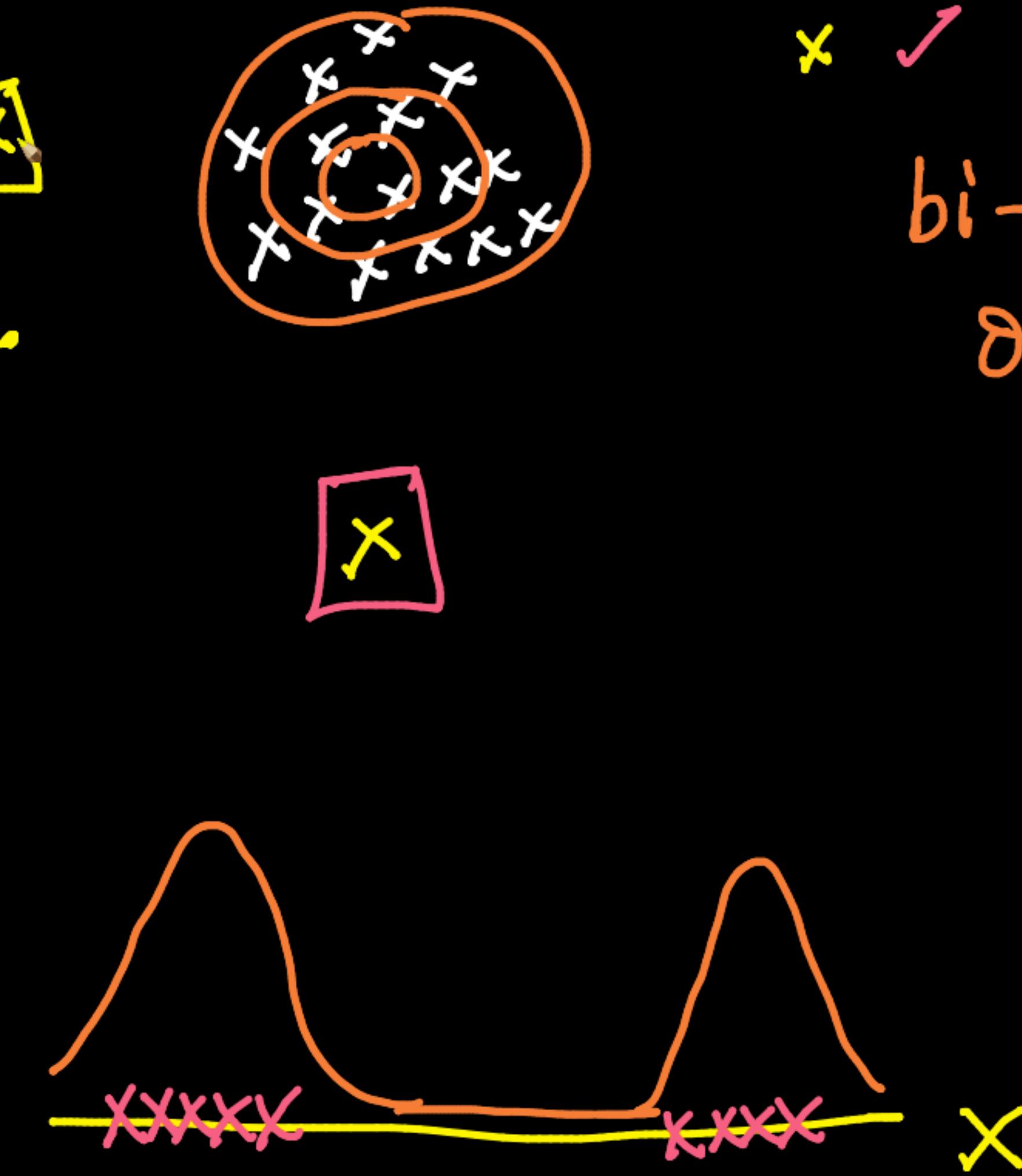
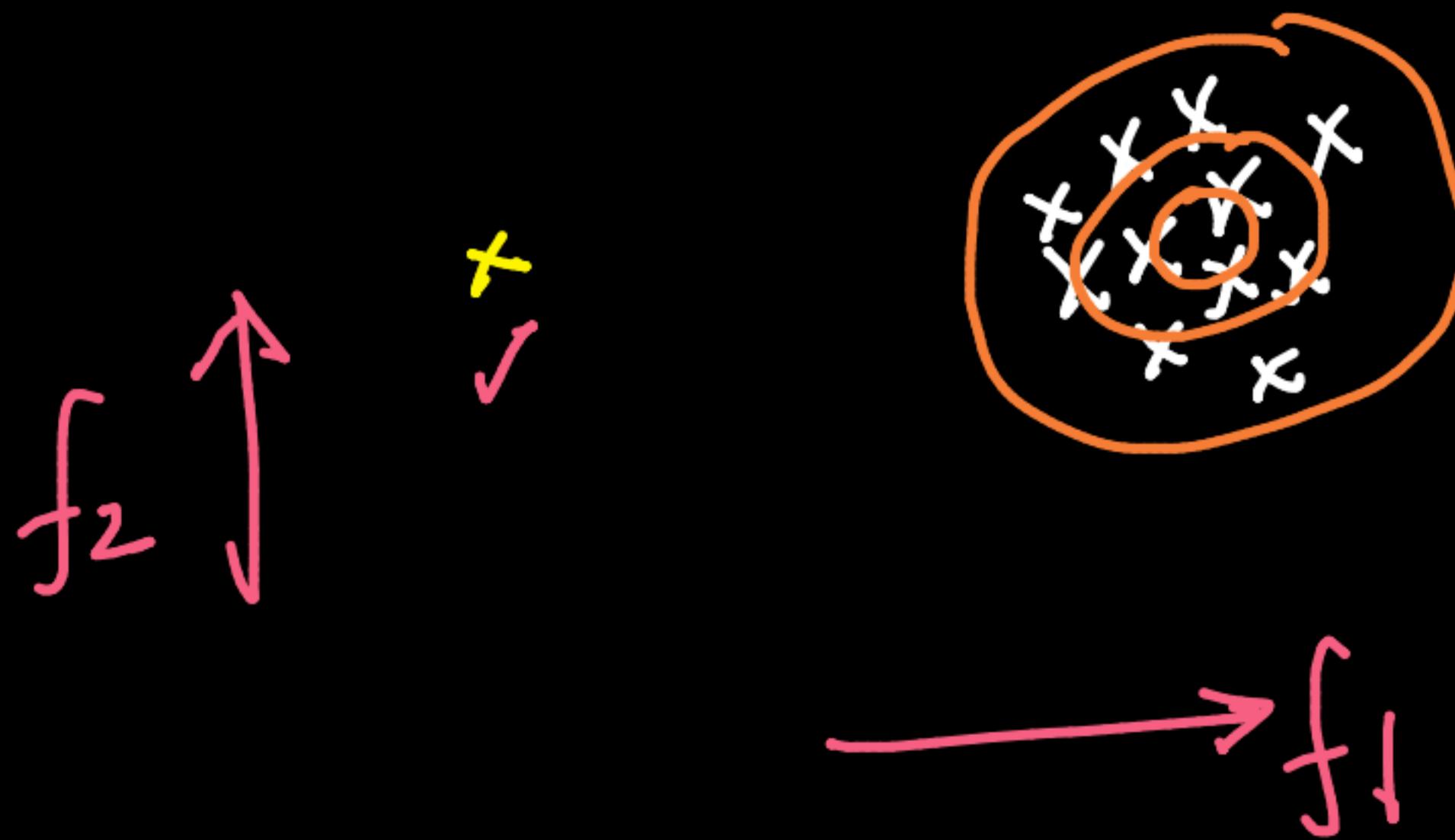


[suggested]

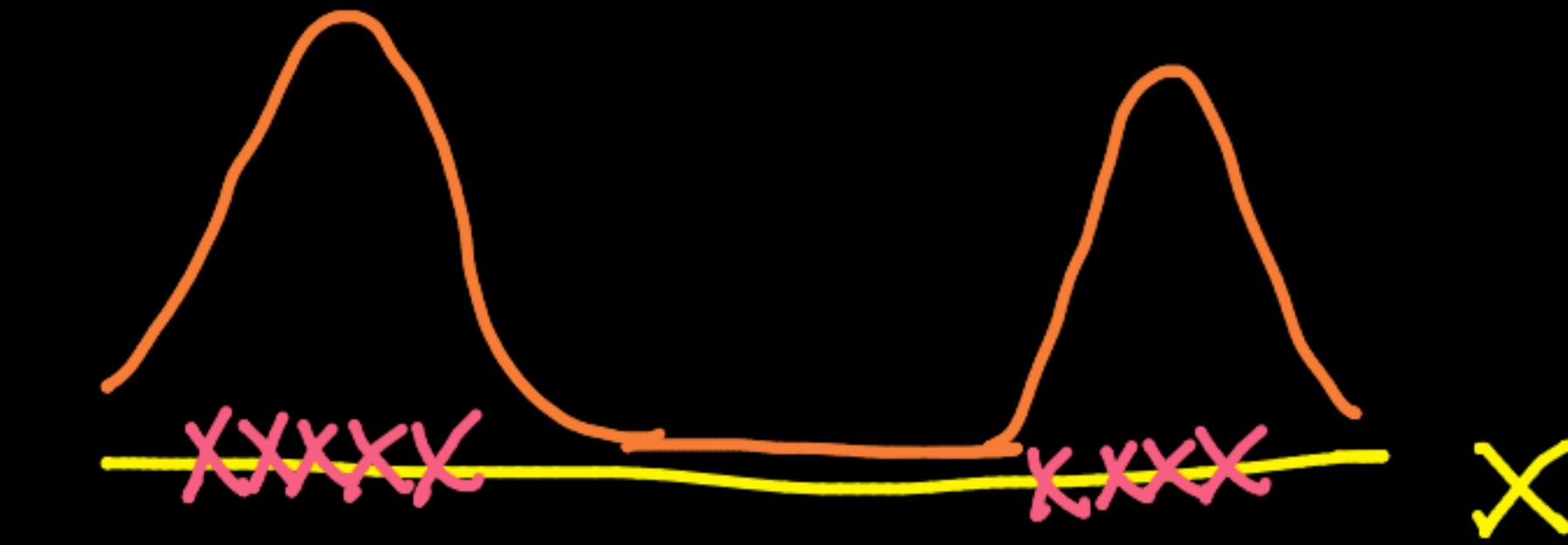
Case 2:



Case 3:



X ✓
bi-modal
data



Red scribbles across the top left.

Boxed "ID".

Known:
 $X \sim \text{disb}(\theta)$

Distribution based

f. $X : x_1 x_2 \dots x_n$

Let $X \sim N(\mu, \sigma^2)$



$\mu \approx \text{sample-mean}$

problem: param of the disb are often corrupted by outliers

Sln:

sln 1:

Drop pt

3σ - away

sln 6: Generative
models

sln 2:

5 - 95 th percentiles

sln 3:

CLT & bootstrapping → outlier-free

sln 4: IQR

sln 5:

Median ≈ Mean; std-dev ≈ MAD

HINT: $X \sim \text{Gaussian}(\mu, \sigma^2)$:  

Robustly estimate θ of my model
from z_1, z_2

$\vec{e} = \theta$ of my model
 (μ, σ^2)

RANSAG

Q

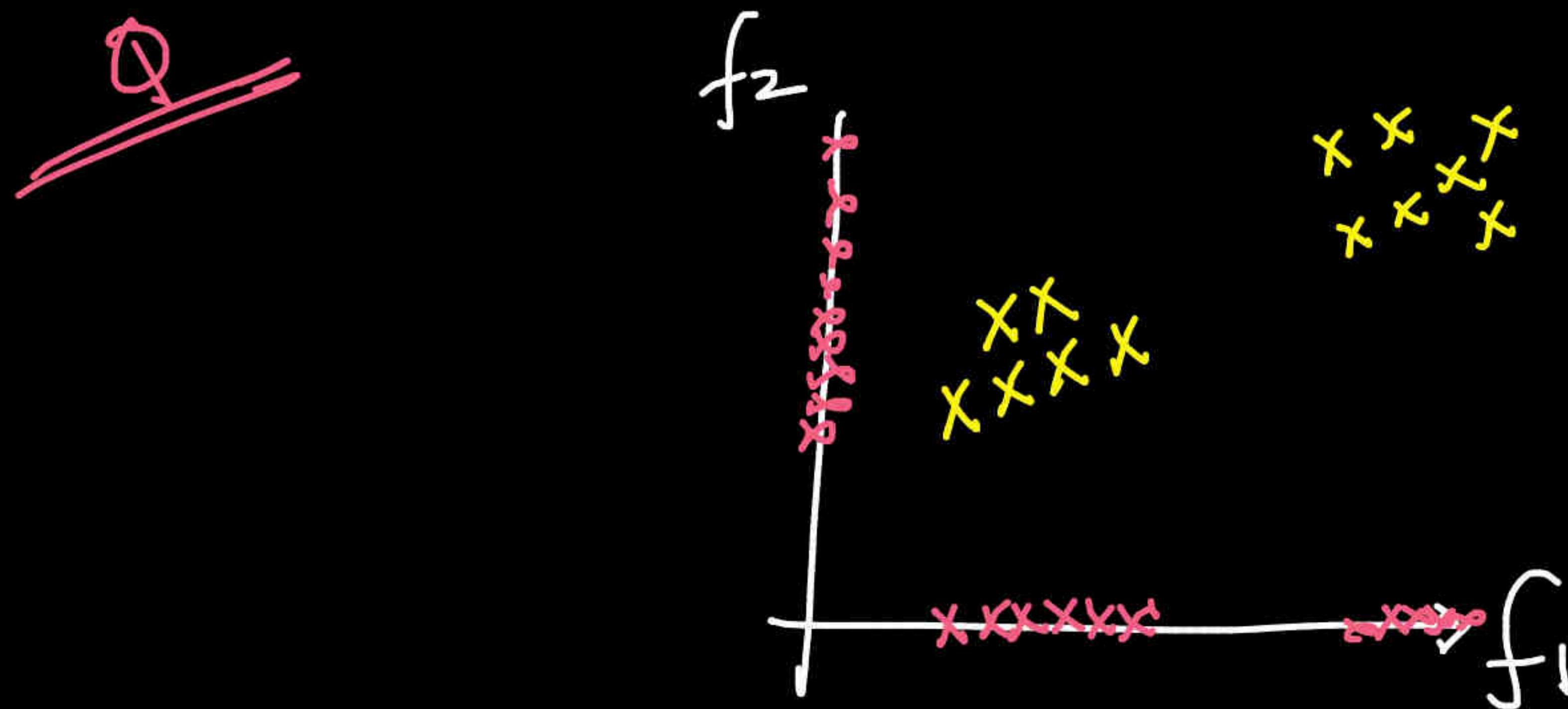
Can we use sample
std-dev (s) as a stopping criterion

s_1

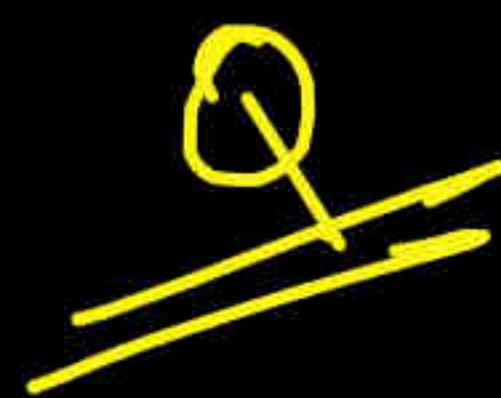
s_2

⋮

$\left[\begin{array}{c} s_t \\ s_{t+1} \end{array} \right]$.



Univariate percentiles
through projections



RANSAC: model based outlier detection



Gaussian (μ, σ^2)



d-dim

$$x_i \in \mathbb{R}^n$$

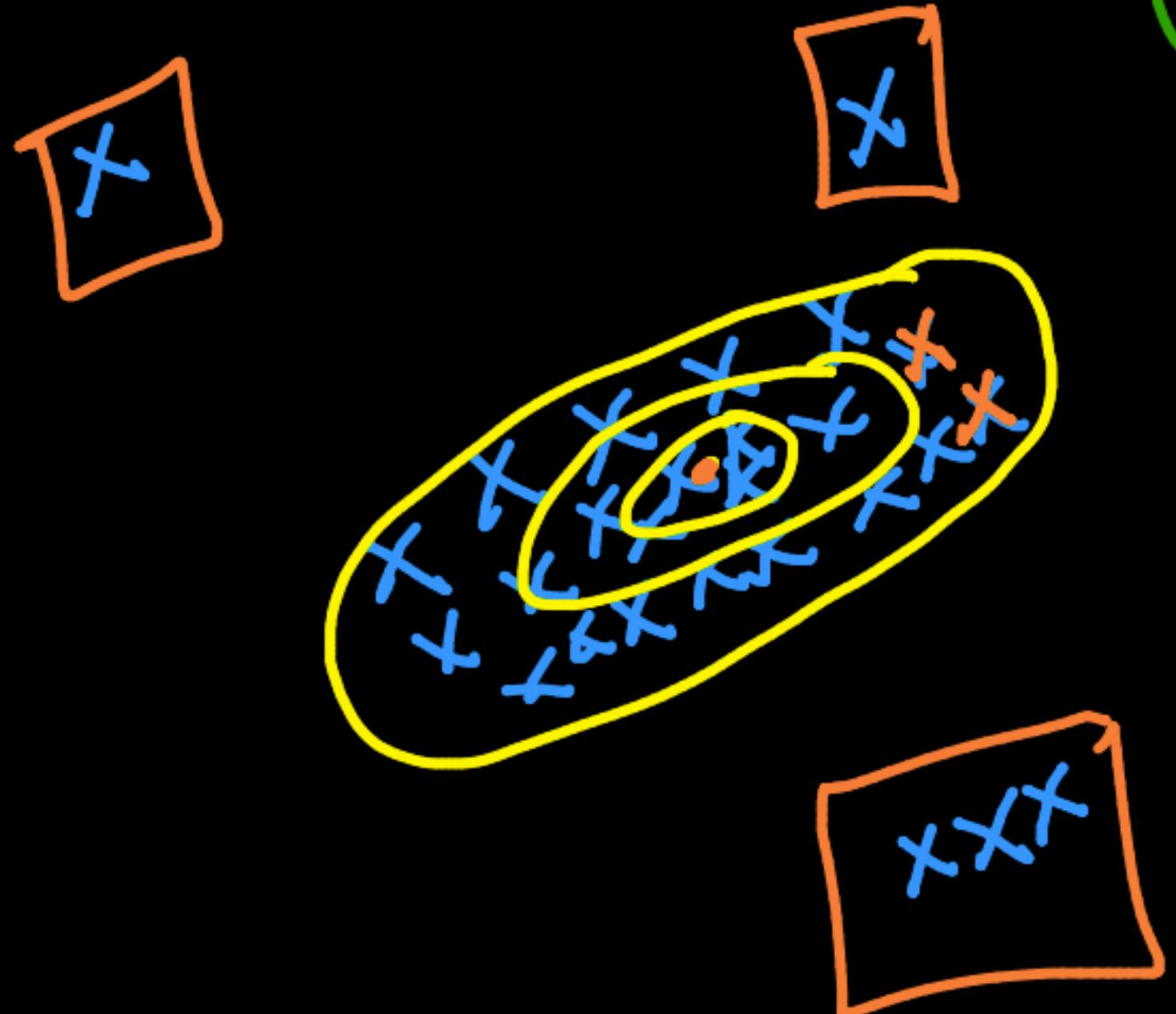
→ Known: Multivariate Gaussian

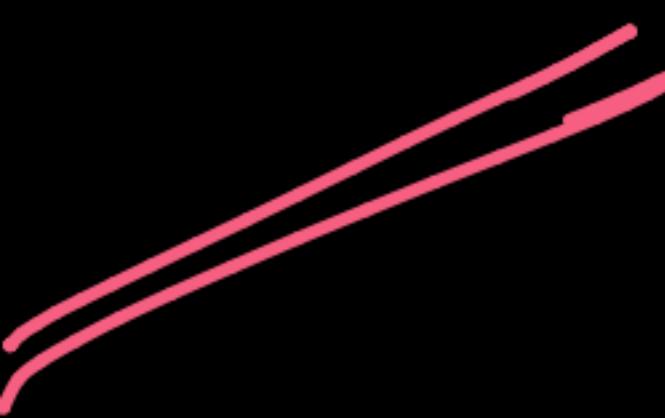
$$x \sim N(\vec{\mu}, \Sigma)$$

The diagram illustrates two primary methods for estimating camera parameters:

- RANSAC-based**: This method is labeled with a circled "1". It is associated with "L₁-Algebra" and "Gov. Matrix slacks".
- Minimum Covariance Determinant (MCD)**: This method is labeled with a circled "2". It is associated with "L₂-Algebra" and "Gov. Matrix slacks".

Both methods are shown with arrows pointing towards the center of the diagram, indicating they are the primary focus.



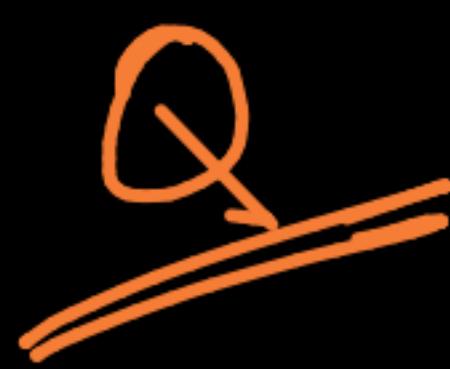


Elliptical Envelope

Given: $x_i \in \mathbb{R}^d$
 $\underline{x} \sim N(\vec{\mu}, \Sigma)$
 $\underline{x} \sim \text{unimodal}$

→ "robustly" estimate
 $\vec{\mu}, \Sigma$

→ remove points which
are outliers



Given \sim params

$$X \sim \text{Disb}(\theta) \quad \Rightarrow$$

$P(X = x_i) = \text{PMF}$ [PDF]

→ convert into Gaussian

RANSAC + estimate θ



WIKIPEDIA
The Free Encyclopedia

Main page

Contents

Current events

Random article

About Wikipedia

Contact us

Donate

Contribute

Help

Learn to edit

Community portal

Recent changes

Upload file

Tools

What links here

Related changes

Special pages

Permanent link

Page information

Multivariate normal distribution

From Wikipedia, the free encyclopedia

"MVN" redirects here. For the airport with that IATA code, see [Mount Vernon Airport](#). For the mvn build-automation software, see [Apache Maven](#).

In probability theory and statistics, the **multivariate normal distribution**, **multivariate Gaussian distribution**, or **joint normal distribution** is a generalization of the one-dimensional (univariate) normal distribution to higher dimensions. One definition is that a **random vector** is said to be k -variate normally distributed if every **linear combination** of its k components has a univariate normal distribution. Its importance derives mainly from the **multivariate central limit theorem**. The multivariate normal distribution is often used to describe, at least approximately, any set of (possibly) **correlated real-valued random variables** each of which clusters around a mean value.

Contents [hide]

1 Definitions

1.1 Notation and parameterization

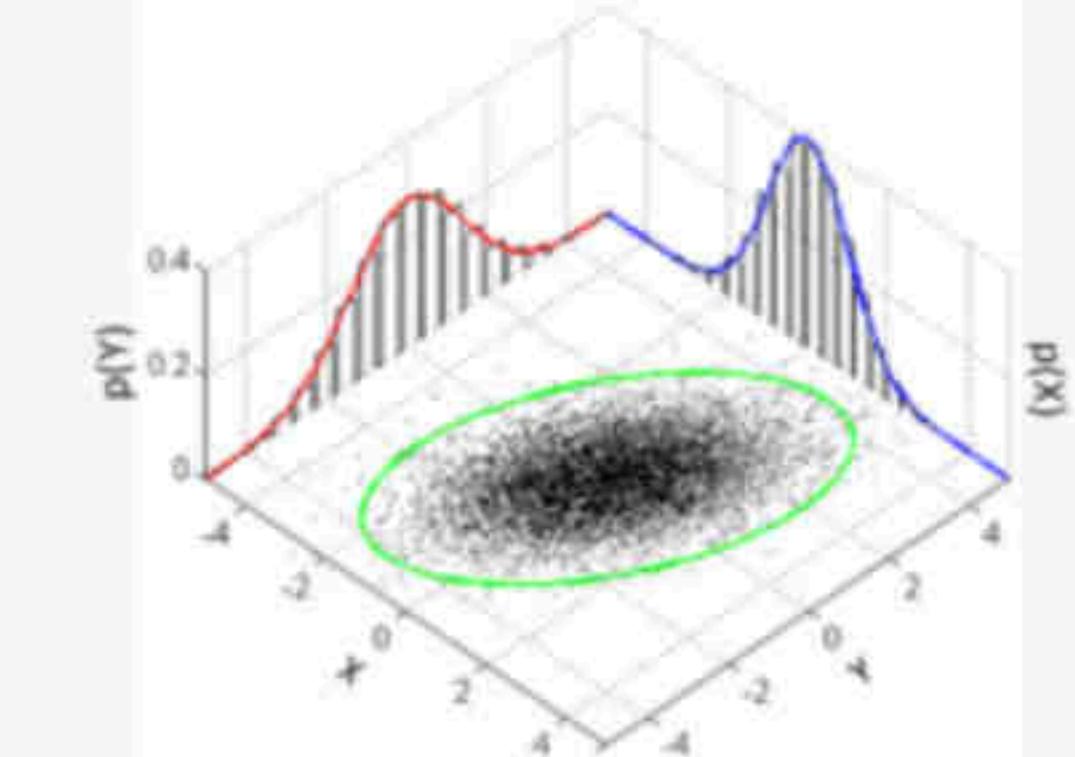
1.2 Standard normal random vector

1.3 Centered normal random vector

1.4 $N(\mu, \Sigma)$

Multivariate normal

Probability density function



Many sample points from a multivariate normal distribution with $\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$ and $\Sigma = \begin{bmatrix} 1 & 3/5 \\ 3/5 & 2 \end{bmatrix}$, shown along with the 3-sigma ellipse, the two marginal distributions, and the two 1-d histograms.

Notation $\mathcal{N}(\mu, \Sigma)$

Parameters $\mu \in \mathbf{R}^k$ — location

$\Sigma \in \mathbf{R}^{k \times k}$ — covariance (positive semi-definite matrix)

Elliptical-envelope : Gaussian + Unimodal

disadv

- non unimodal data
- MV-Gaussian

Popular techniques [edit]

Many anomaly detection techniques have been proposed in literature.^{[1][10]} Some of the popular techniques are:

- Statistical (Z-score, Tukey's range test and Grubbs's test)
- Density-based techniques (k-nearest neighbor,^{[11][12][13]} local outlier factor,^[14] isolation forests,^{[15][16]} and many more variations of this concept^[17])
- Subspace-,^[18] correlation-based^[19] and tensor-based^[20] outlier detection for high-dimensional data^[21]
- One-class support vector machines^[22]
- Replicator neural networks,^[23] autoencoders, variational autoencoders,^[24] long short-term memory neural networks^[25]
- Bayesian networks^[23]
- Hidden Markov models (HMMs)^[23]
- Minimum Covariance Determinant^{[26][27]}
- Clustering: Cluster analysis-based outlier detection^{[28][29]}
- Deviations from association rules and frequent itemsets
- Fuzzy logic-based outlier detection
- Ensemble techniques, using feature bagging,^{[30][31]} score normalization^{[32][33]} and different sources of diversity^{[34][35]}

The performance of methods depends on the data set and parameters, and methods have little systematic advantages over another when compared across many data sets and parameters.^{[36][37]}

Software [edit]

- ELKI is an open-source Java data mining toolkit that contains several anomaly detection algorithms, as well as index acceleration for them.



scikit-learn.org/stable/modules/generated/sklearn.covariance.EllipticEnvelope.html

scikit-learn Install User Guide API Examples Community More ▾ Go

Prev Up Next

scikit-learn 1.1.1

Other versions

Please [cite us](#) if you use the software.

[sklearn.covariance.EllipticEnvelope](#)

Examples using [sklearn.covariance.EllipticEnvelope](#)

sklearn.covariance.EllipticEnvelope

```
class sklearn.covariance.EllipticEnvelope(*, store_precision=True, assume_centered=False,  
support_fraction=None, contamination=0.1, random_state=None)
```

[source]

An object for detecting outliers in a Gaussian distributed dataset.

Read more in the [User Guide](#).

Parameters: **store_precision : bool, default=True**

Specify if the estimated precision is stored.

assume_centered : bool, default=False

If True, the support of robust location and covariance estimates is computed, and a covariance estimate is recomputed from it, without centering the data. Useful to work with data whose mean is significantly equal to zero but is not exactly zero. If False, the robust location and covariance are directly computed with the FastMCD algorithm without additional treatment.

support_fraction : float, default=None

19 / 19

ed in the support of the raw MCD estimate. If None, the

Toggle Menu

scikit-learn.org/stable/modules/generated/sklearn.covariance.EllipticEnvelope.html

scikit-learn Install User Guide API Examples Community More ▾ Go

Prev Up Next

scikit-learn 1.1.1

Other versions

Please [cite us](#) if you use the software.

[sklearn.covariance.EllipticEnvelope](#)

Examples using [sklearn.covariance.EllipticEnvelope](#)

sklearn.covariance.EllipticEnvelope

```
class sklearn.covariance.EllipticEnvelope(*, store_precision=True, assume_centered=False,  
support_fraction=None, contamination=0.1, random_state=None)
```

[source]

An object for detecting outliers in a Gaussian distributed dataset.

Read more in the [User Guide](#).

Parameters:

store_precision : bool, default=True

Specify if the estimated precision is stored.

$$\hat{\mu} = \bar{x}$$



assume_centered : bool, default=False

If True, the support of robust location and covariance estimates is computed, and a covariance estimate is recomputed from it, without centering the data. Useful to work with data whose mean is significantly equal to zero but is not exactly zero. If False, the robust location and covariance are directly computed with the FastMCD algorithm without additional treatment.

support_fraction : float, default=None

Used in the support of the raw MCD estimate. If None, the

scikit-learn.org/stable/modules/generated/sklearn.covariance.EllipticEnvelope.html

scikit-learn

Install User Guide API Examples Community More ▾ Go

Prev Up Next

scikit-learn 1.1.1

Other versions

Please [cite us](#) if you use the software.

[sklearn.covariance.EllipticEnvelope](#)

Examples using [sklearn.covariance.EllipticEnvelope](#)

sklearn.covariance.EllipticEnvelope

```
class sklearn.covariance.EllipticEnvelope(*, store_precision=True, assume_centered=False,  
support_fraction=None, contamination=0.1, random_state=None)
```

[source]

An object for detecting outliers in a Gaussian distributed dataset.

Read more in the [User Guide](#).

Parameters:

store_precision : bool, default=True

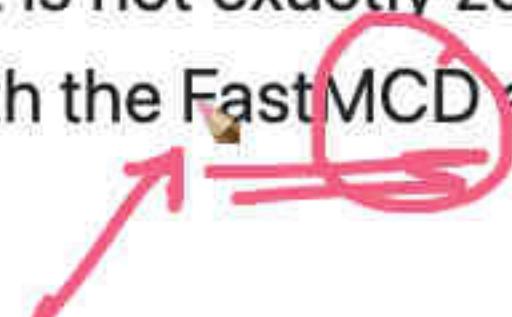
Specify if the estimated precision is stored.

assume_centered : bool, default=False

If True, the support of robust location and covariance estimates is computed, and a covariance estimate is recomputed from it, without centering the data. Useful to work with data whose mean is significantly equal to zero but is not exactly zero. If False, the robust location and covariance are directly computed with the FastMCD algorithm without additional treatment.

support_fraction : float, default=None

Used in the support of the raw MCD estimate. If None, the



scikit-learn.org/stable/modules/generated/sklearn.covariance.EllipticEnvelope.html

scikit-learn

Prev Up Next

scikit-learn 1.1.1 Other versions

Please cite us if you use the software.

sklearn.covariance.EllipticEnvelope Examples using sklearn.covariance.EllipticEnvelope

Read more in the User Guide.

Parameters:

- store_precision : bool, default=True**
Specify if the estimated precision is stored.
- assume_centered : bool, default=False**
If True, the support of robust location and covariance estimates is computed, and a covariance estimate is recomputed from it, without centering the data. Useful to work with data whose mean is significantly equal to zero but is not exactly zero. If False, the robust location and covariance are directly computed with the FastMCD algorithm without additional treatment.
- support_fraction : float, default=None**
The proportion of points to be included in the support of the raw MCD estimate. If None, the minimum value of support_fraction will be used within the algorithm: $[n_sample + n_features + 1] / 2$. Range is (0, 1).
- contamination : float, default=0.1**
The amount of contamination of the data set, i.e. the proportion of outliers in the data set. Range is (0, 0.5].
- random_state : int, RandomState instance or None, default=None**
Determines the pseudo random number generator for shuffling the data. Pass an int for junction calls. See Glossary.

500
1000

Subset $\rightarrow \bar{\mu}, \Sigma$ ✓

Toggle Menu

scikit-learn.org/stable/modules/generated/sklearn.covariance.EllipticEnvelope.html

Specify if the estimated precision is stored.

Prev Up Next

scikit-learn 1.1.1 Other versions

Please cite us if you use the software.

sklearn.covariance.EllipticEnvelope Examples using sklearn.covariance.EllipticEnvelope

assume_centered : bool, default=False

If True, the support of robust location and covariance estimates is computed, and a covariance estimate is recomputed from it, without centering the data. Useful to work with data whose mean is significantly equal to zero but is not exactly zero. If False, the robust location and covariance are directly computed with the FastMCD algorithm without additional treatment.

support_fraction : float, default=None

The proportion of points to be included in the support of the raw MCD estimate. If None, the minimum value of support_fraction will be used within the algorithm: $[n_sample + n_features + 1] / 2$. Range is (0, 1).

contamination : float, default=0.1

The amount of contamination of the data set, i.e. the proportion of outliers in the data set. Range is (0, 0.5).

random_state : int, RandomState instance or None, default=None

Determines the pseudo random number generator for shuffling the data. Pass an int for reproducible results across multiple function calls. See [Glossary](#).

Attributes:

location_ : ndarray of shape (n_features,)

Estimated robust location.



scikit-learn.org/stable/modules/generated/sklearn.covariance.EllipticEnvelope.html

scikit-learn

Prev Up Next

scikit-learn 1.1.1

Other versions

Please cite us if you use the software.

sklearn.covariance.EllipticEnvelope

Examples using sklearn.covariance.EllipticEnvelope

μ, Σ

References

1 Rousseeuw, P.J., Van Driessen, K. "A fast algorithm for the minimum covariance determinant estimator" Technometrics 41(3), 212 (1999)

Examples

```
>>> import numpy as np
>>> from sklearn.covariance import EllipticEnvelope
>>> true_cov = np.array([[.8, .3],
...                      [.3, .4]])
...
>>> X = np.random.RandomState(0).multivariate_normal(mean=[0, 0],
... cov=true_cov,
... size=500)
...
...
>>> cov = EllipticEnvelope(random_state=0).fit(X)
>>> # predict returns 1 for an inlier and -1 for an outlier
>>> cov.predict([[0, 0],
...               [3, 3]])
array([-1,  1])
>>> cov.covariance_
array([[0.7411...,  0.2535...],
       [0.2535...,  0.3053...]])
>>> cov.location_
array([0.0813...,  0.0427...])
```

Methods

Toggle Menu

scikit-learn.org/stable/auto_examples/miscellaneous/plot_anomaly_comparison.html#sphx-glr-auto-examples-

sklearn.coverage_ Comparing a... on a real data... isolation fore... scikitlearn.m... isolationFore... classification... SVM with no... Local outlier f... with Local O... Evaluation of... Multivariate n...

scikit learn Install User Guide API Examples Community More ▾

Prev Up Next

scikit-learn 1.1.1 Other versions

Please cite us if you use the software.

Comparing anomaly detection algorithms for outlier detection on toy datasets

Finally, note that parameters of the models have been here handpicked but that in practice they need to be adjusted. In the absence of labelled data, the problem is completely unsupervised so model selection can be a challenge.

Robust covariance One-Class SVM One-Class SVM (SGD) Isolation Forest Local Outlier Factor

.03s .00s .01s .12s .00s

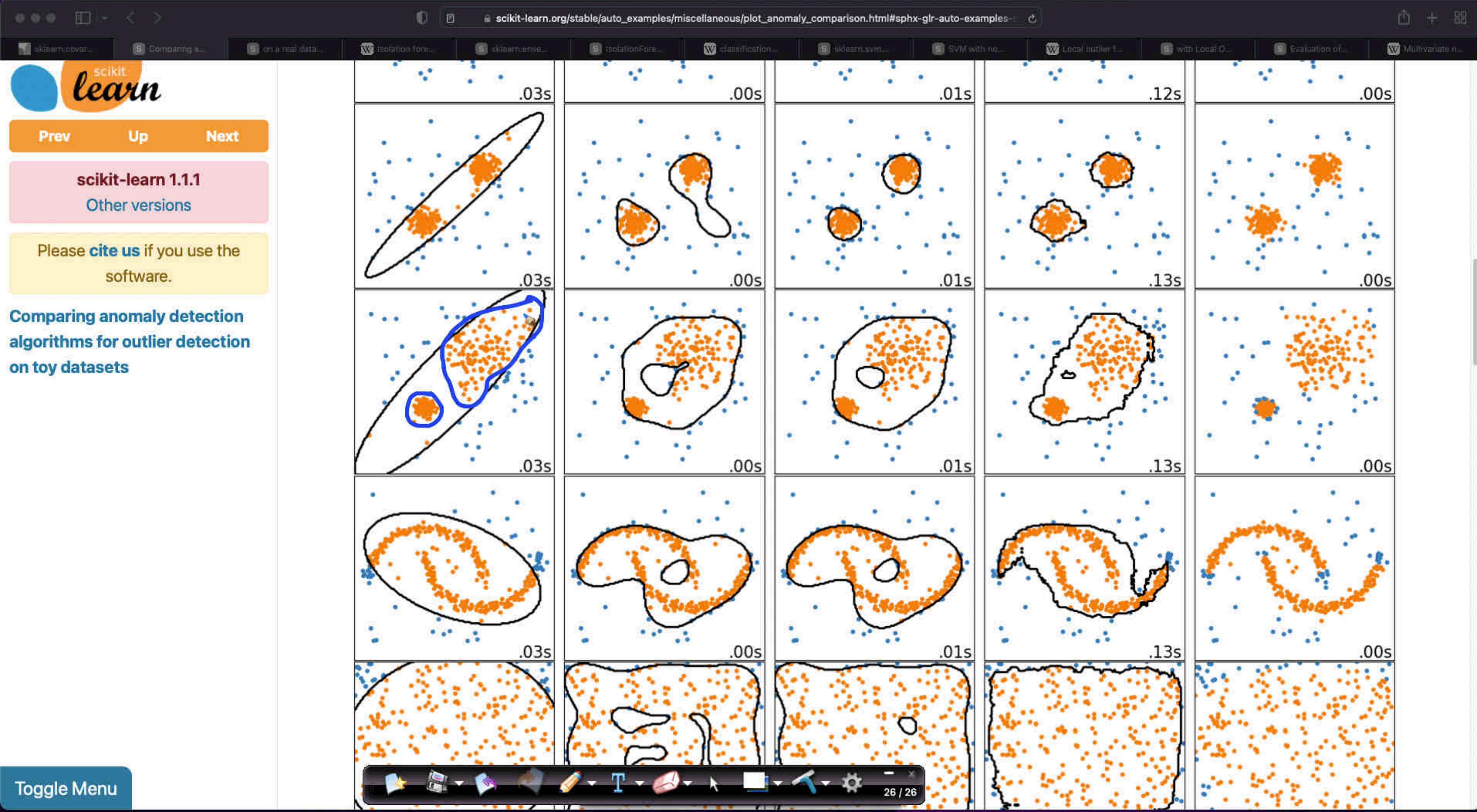
.03s .00s .01s .13s .00s

.01s .13s .00s

Toggle Menu

Handwritten annotations:

- A red checkmark is placed above the text "Finally, note that parameters of the models have been here handpicked but that in practice they need to be adjusted. In the absence of labelled data, the problem is completely unsupervised so model selection can be a challenge."
- A blue arrow points from the "Please cite us if you use the software." box to the "Robust covariance" plot.
- A blue circle highlights an outlier in the bottom-left plot (Isolation Forest).



scikit-learn.org/stable/auto_examples/miscellaneous/plot_anomaly_comparison.html#sphx-glr-auto-examples-

scikit
learn

Prev Up Next

scikit-learn 1.1.1

Other versions

Please cite us if you use the software.

Comparing anomaly detection algorithms for outlier detection on toy datasets

.03s .00s .01s .13s .00s

.03s .00s .01s .13s .00s

.05s .00s .01s .13s .00s

```
# Author: Alexandre Gramfort <alexandre.gramfort@inria.fr>
#         Albert Thomas <albert.thomas@telecom-paristech.fr>
# License: BSD 3 clause
```

Toggle Menu

scikit-learn.org/stable/auto_examples/miscellaneous/plot_anomaly_comparison.html#sphx-glr-auto-examples-

scikit learn

Prev Up Next

scikit-learn 1.1.1 Other versions

Please cite us if you use the software.

Comparing anomaly detection algorithms for outlier detection on toy datasets

Author: Alexandre Gramfort <alexandre.gramfort@inria.fr>
Albert Thomas <albert.thomas@telecom-paristech.fr>
License: BSD 3 clause

Toggle Menu

scikit-learn.org/stable/auto_examples/miscellaneous/plot_anomaly_comparison.html#sphx-glr-auto-examples-

scikit
learn

Prev Up Next

scikit-learn 1.1.1

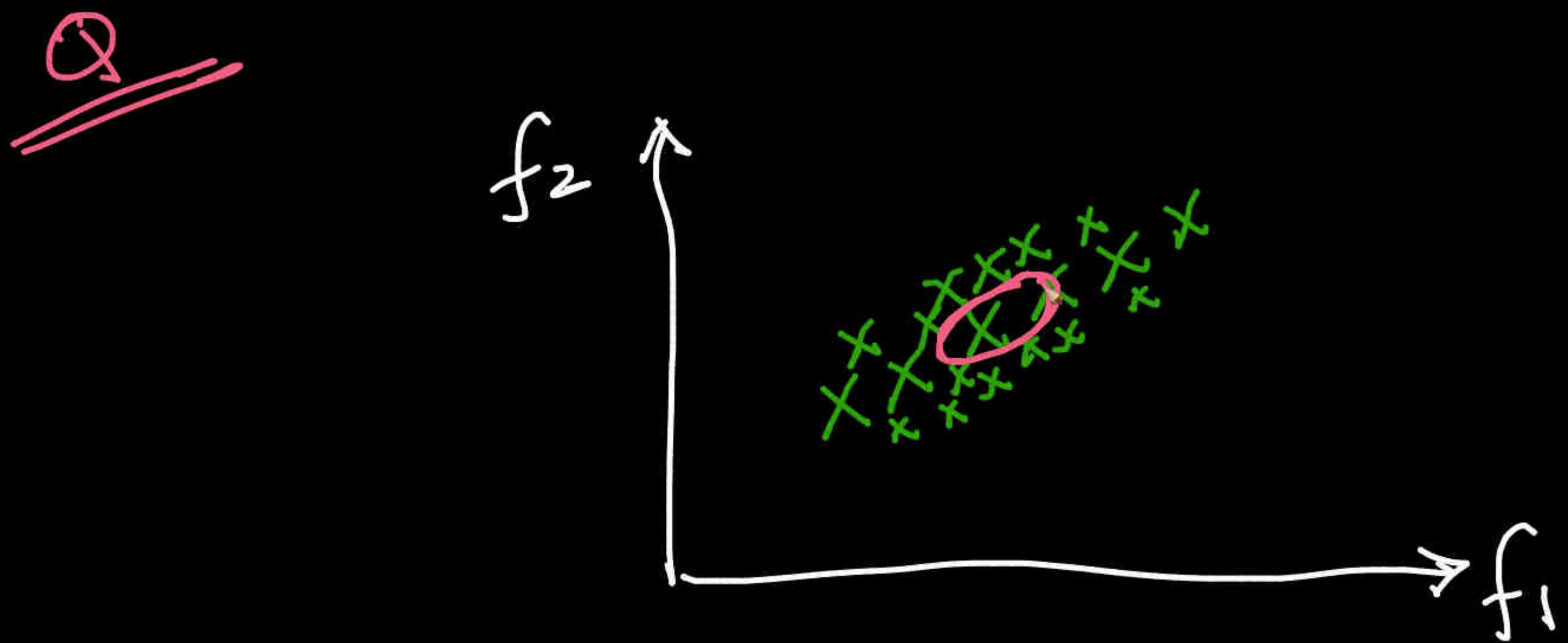
Other versions

Please cite us if you use the software.

Comparing anomaly detection algorithms for outlier detection on toy datasets

Author: Alexandre Gramfort <alexandre.gramfort@inria.fr>
Albert Thomas <albert.thomas@telecom-paristech.fr>
License: BSD 3 clause

Toggle Menu

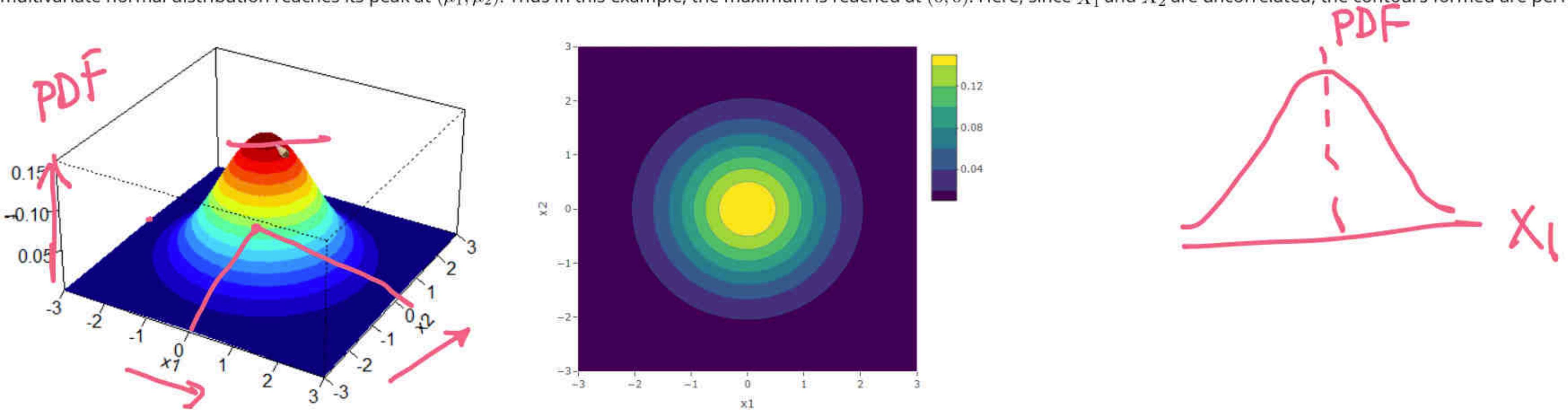




Case 1: $X_1, X_2 \sim \mathcal{N}(0, 1)$ with correlation 0.

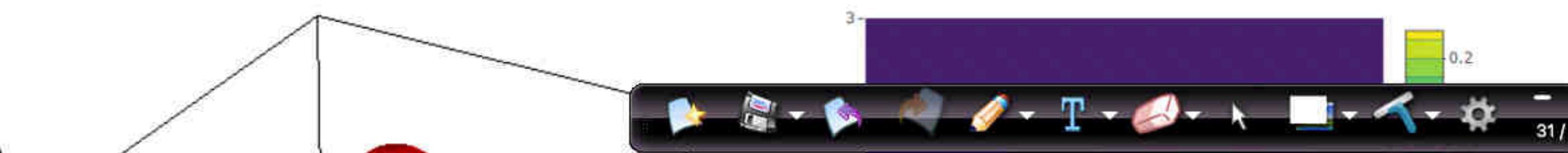
The equation for the correlation ρ is given by $\rho = \frac{\sigma_{12}}{\sqrt{\sigma_{11}\sigma_{22}}}$. Hence $\sigma_{12} = \rho\sqrt{\sigma_{11}\sigma_{22}}$. When $\rho = 0$, $\sigma_{11} = 1$ and $\sigma_{22} = 1$, $\sigma_{12} = \rho = 0$.

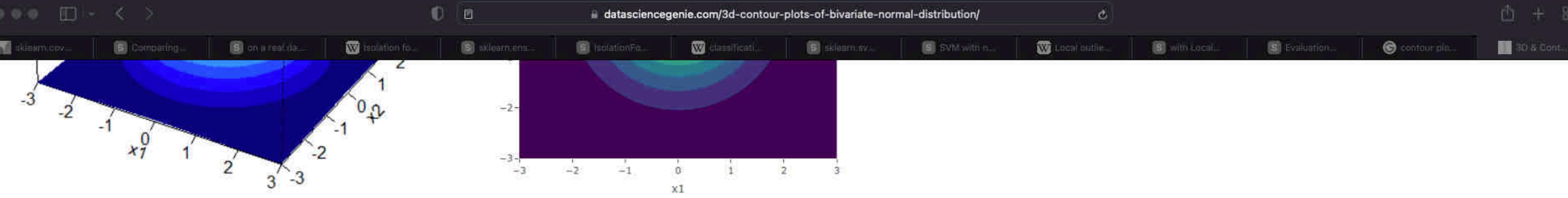
We substitute $\mu_1, \mu_2 = 0$, $\sigma_{11}, \sigma_{22} = 1$ and $\sigma_{12} = 0$ in Equation (1) and obtain the following 3D plot and contour plot. The shape of the bivariate normal distribution is again similar to that of a bell. The multivariate normal distribution reaches its peak at (μ_1, μ_2) . Thus in this example, the maximum is reached at $(0, 0)$. Here, since X_1 and X_2 are uncorrelated, the contours formed are perfect circles.



Case 2: $X_1, X_2 \sim \mathcal{N}(0, 1)$ with correlation 0.7.

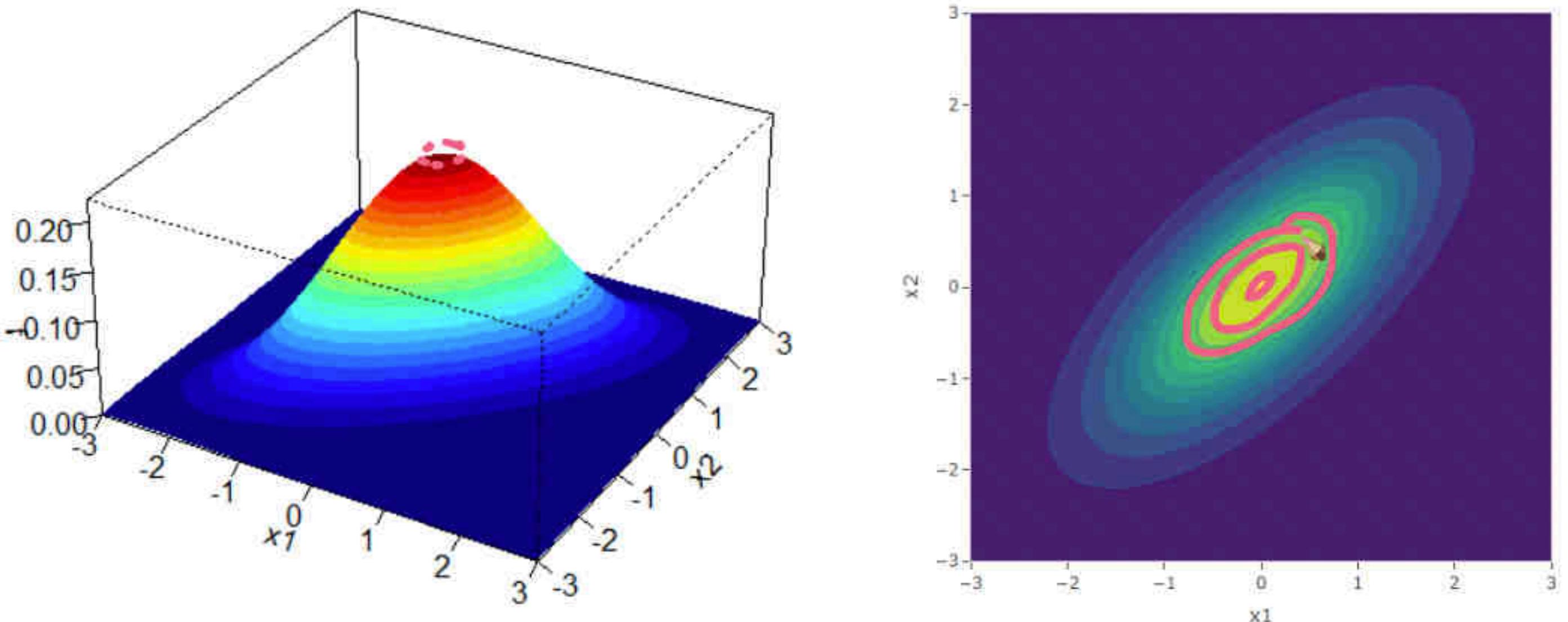
Since $\rho = 0.7$ and $\sigma_{12} = \rho\sqrt{\sigma_{11}\sigma_{22}}$, then $\sigma_{12} = 0.7$. By inputting the values of the means, variances and covariances in Equation (1), we obtain the following plots. Again we obtain a bell-shaped bivariate distribution. Since the correlation between X_1 and X_2 is positive, we obtain elliptical contours. It is like taking the circular contours of the uncorrelated case and elongate them along the diagonal $x_2 = x_1$.





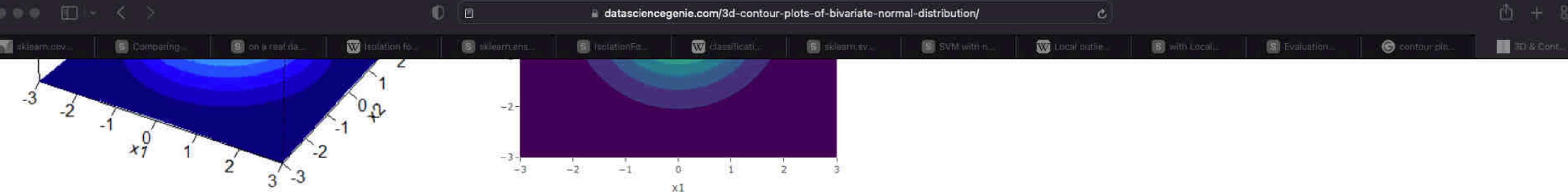
Case 2: $X_1, X_2 \sim \mathcal{N}(0, 1)$ with correlation 0.7.

Since $\rho = 0.7$ and $\sigma_{12} = \rho\sqrt{\sigma_{11}\sigma_{22}}$, then $\sigma_{12} = 0.7$. By inputting the values of the means, variances and covariances in Equation (1), we obtain the following plots. Again we obtain a bell-shaped bivariate distribution. Since the correlation between X_1 and X_2 is positive, we obtain elliptical contours. It is like taking the circular contours of the uncorrelated case and elongate them along the diagonal $x_2 = x_1$.



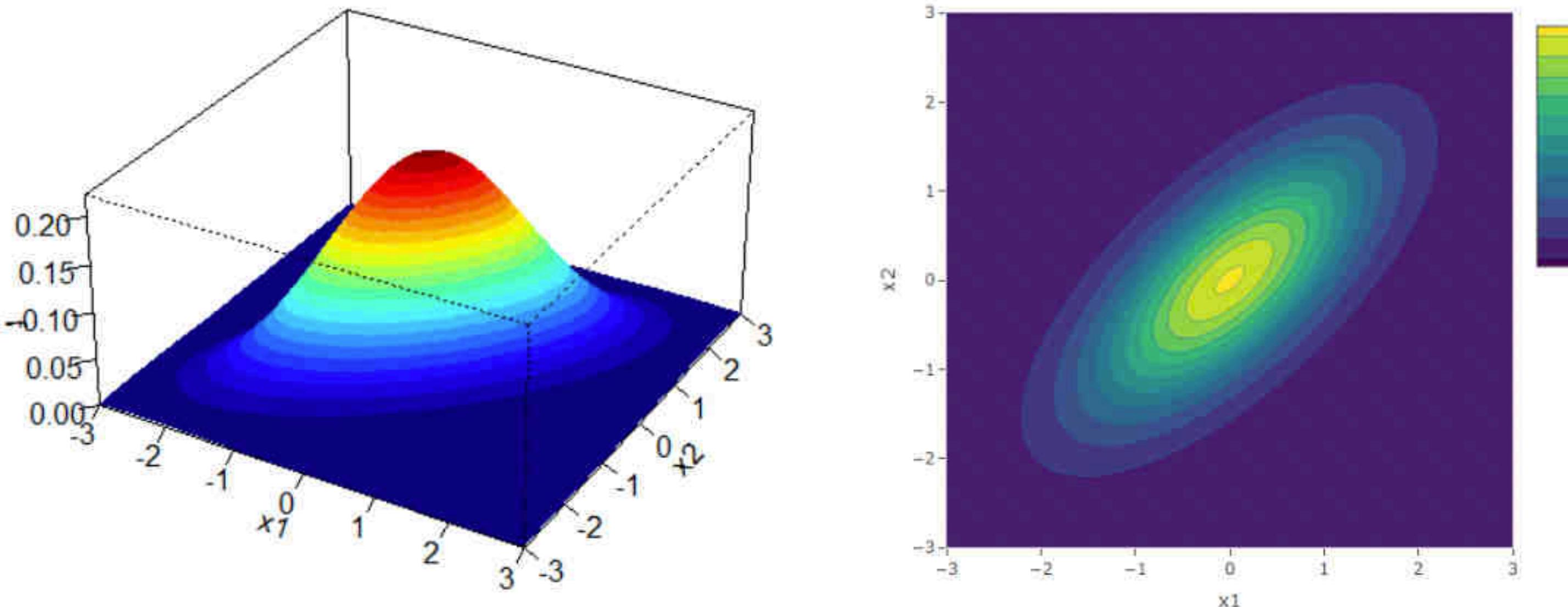
Case 3: $X_1, X_2 \sim \mathcal{N}(0, 1)$ with correlation -0.7.

The following are the plots for the case when the correlation is negative. Similar to the second case, we have elliptical contours. However in this case, the circular contours are elongated along the second diagonal, that is, the line $x_2 = -x_1$.



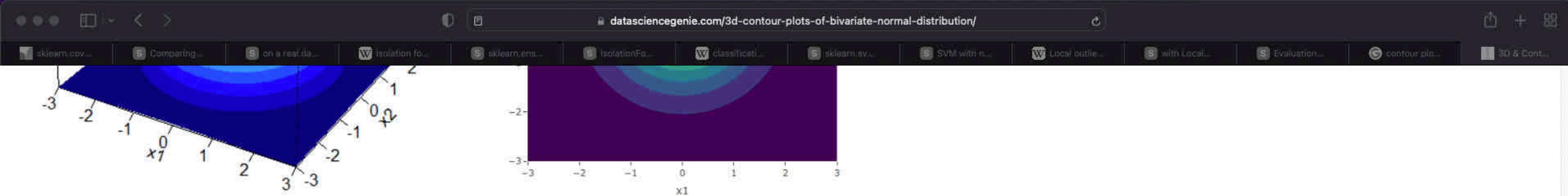
Case 2: $X_1, X_2 \sim \mathcal{N}(0, 1)$ with correlation 0.7.

Since $\rho = 0.7$ and $\sigma_{12} = \rho\sqrt{\sigma_{11}\sigma_{22}}$, then $\sigma_{12} = 0.7$. By inputting the values of the means, variances and covariances in Equation (1), we obtain the following plots. Again we obtain a bell-shaped bivariate distribution. Since the correlation between X_1 and X_2 is positive, we obtain elliptical contours. It is like taking the circular contours of the uncorrelated case and elongate them along the diagonal $x_2 = x_1$.



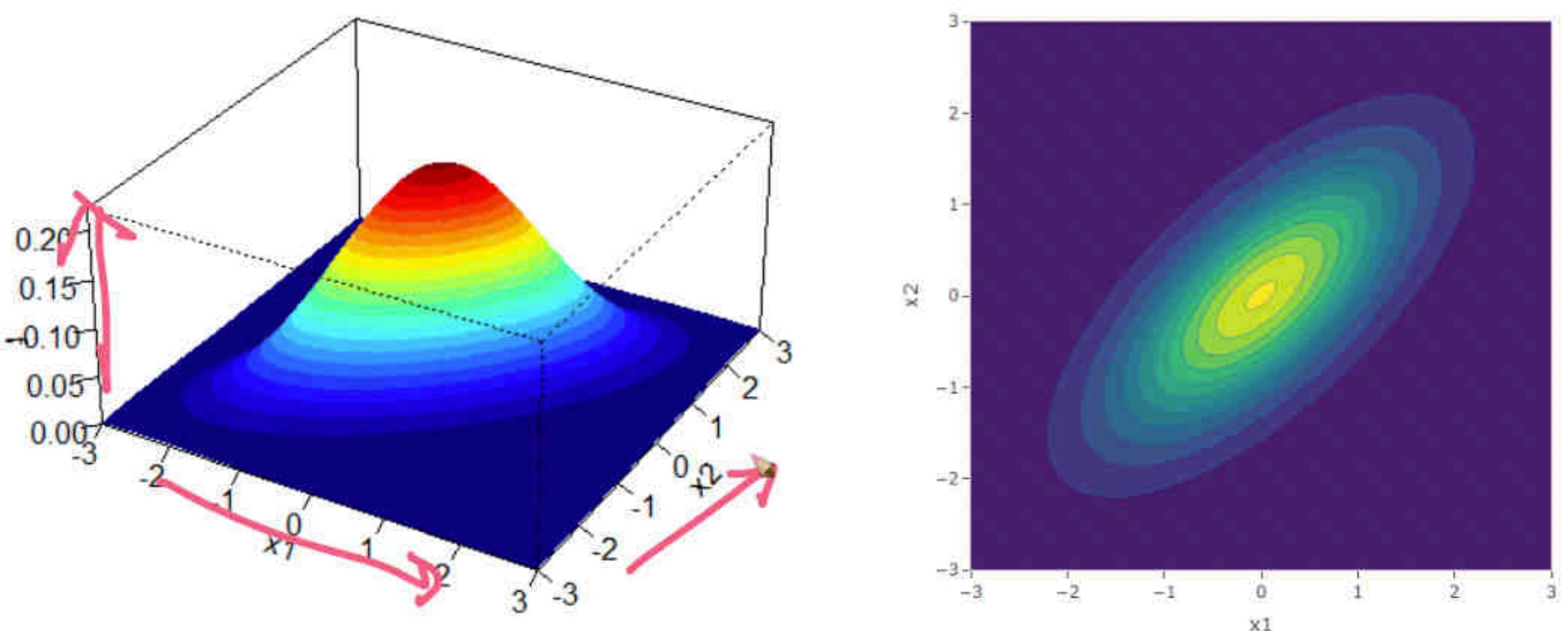
Case 3: $X_1, X_2 \sim \mathcal{N}(0, 1)$ with correlation -0.7.

The following are the plots for the case when the correlation is negative. Similar to the second case, we have elliptical contours. However in this case, the circular contours are elongated along the second diagonal, that is, the line $x_2 = -x_1$.



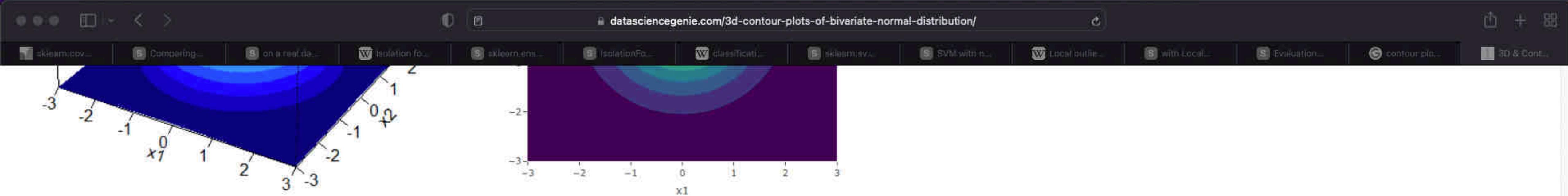
Case 2: $X_1, X_2 \sim \mathcal{N}(0, 1)$ with correlation 0.7.

Since $\rho = 0.7$ and $\sigma_{12} = \rho\sqrt{\sigma_{11}\sigma_{22}}$, then $\sigma_{12} = 0.7$. By inputting the values of the means, variances and covariances in Equation (1), we obtain the following plots. Again we obtain a bell-shaped bivariate distribution. Since the correlation between X_1 and X_2 is positive, we obtain elliptical contours. It is like taking the circular contours of the uncorrelated case and elongate them along the diagonal $x_2 = x_1$.



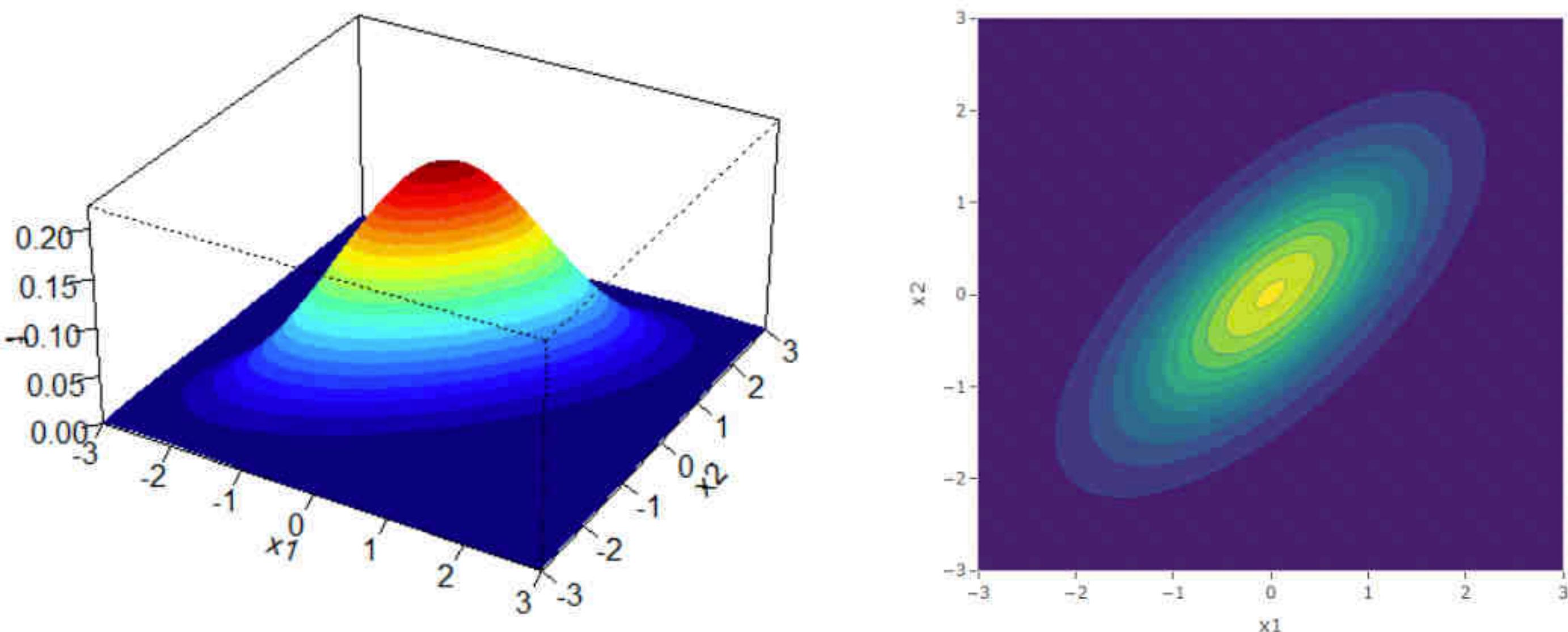
Case 3: $X_1, X_2 \sim \mathcal{N}(0, 1)$ with correlation -0.7.

The following are the plots for the case when the correlation is negative. Similar to the second case, we have elliptical contours. However in this case, the circular contours are elongated along the second diagonal, that is, the line $x_2 = -x_1$.



Case 2: $X_1, X_2 \sim \mathcal{N}(0, 1)$ with correlation 0.7.

Since $\rho = 0.7$ and $\sigma_{12} = \rho\sqrt{\sigma_{11}\sigma_{22}}$, then $\sigma_{12} = 0.7$. By inputting the values of the means, variances and covariances in Equation (1), we obtain the following plots. Again we obtain a bell-shaped bivariate distribution. Since the correlation between X_1 and X_2 is positive, we obtain elliptical contours. It is like taking the circular contours of the uncorrelated case and elongate them along the diagonal $x_2 = x_1$.



Case 3: $X_1, X_2 \sim \mathcal{N}(0, 1)$ with correlation -0.7.

The following are the plots for the case when the correlation is negative. Similar to the second case, we have elliptical contours. However in this case, the circular contours are elongated along the second diagonal, that is, the line $x_2 = -x_1$.

Isolation Forest

$$x_i \in \mathbb{R}^d$$

Isolation Forest

==>

2007 - ...

$$\mathcal{D} = \{x_i\}_{i=1}^n$$

- Idea:
- build many trees (like RF)
 - each tree:
 - randomly pick a feature
 - randomly threshold it
 - build each tree till the leafnode consists of just one datapoint

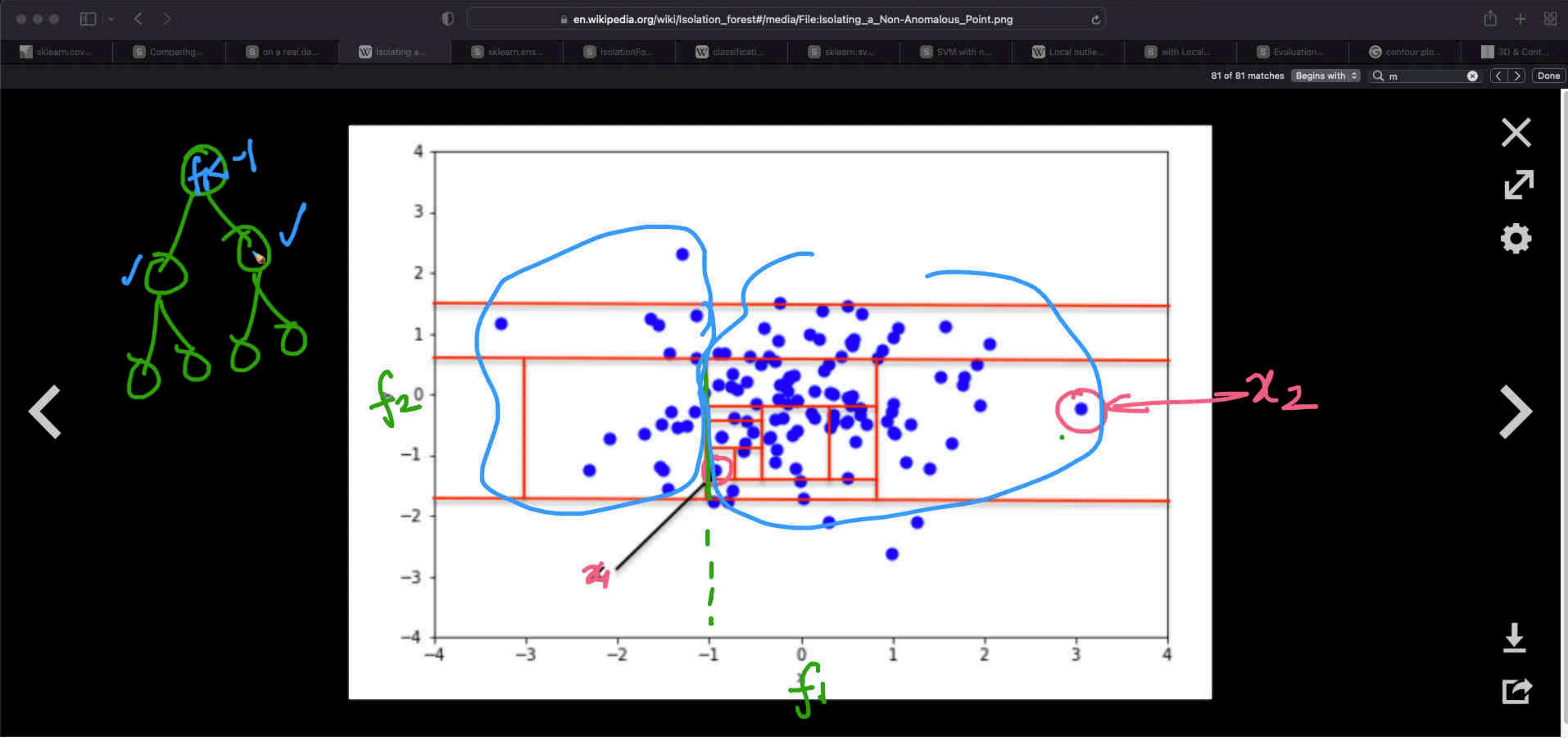


Fig. 2 - an example of isolating a non-anomalous point in a 2D Gaussian distribution.

More details

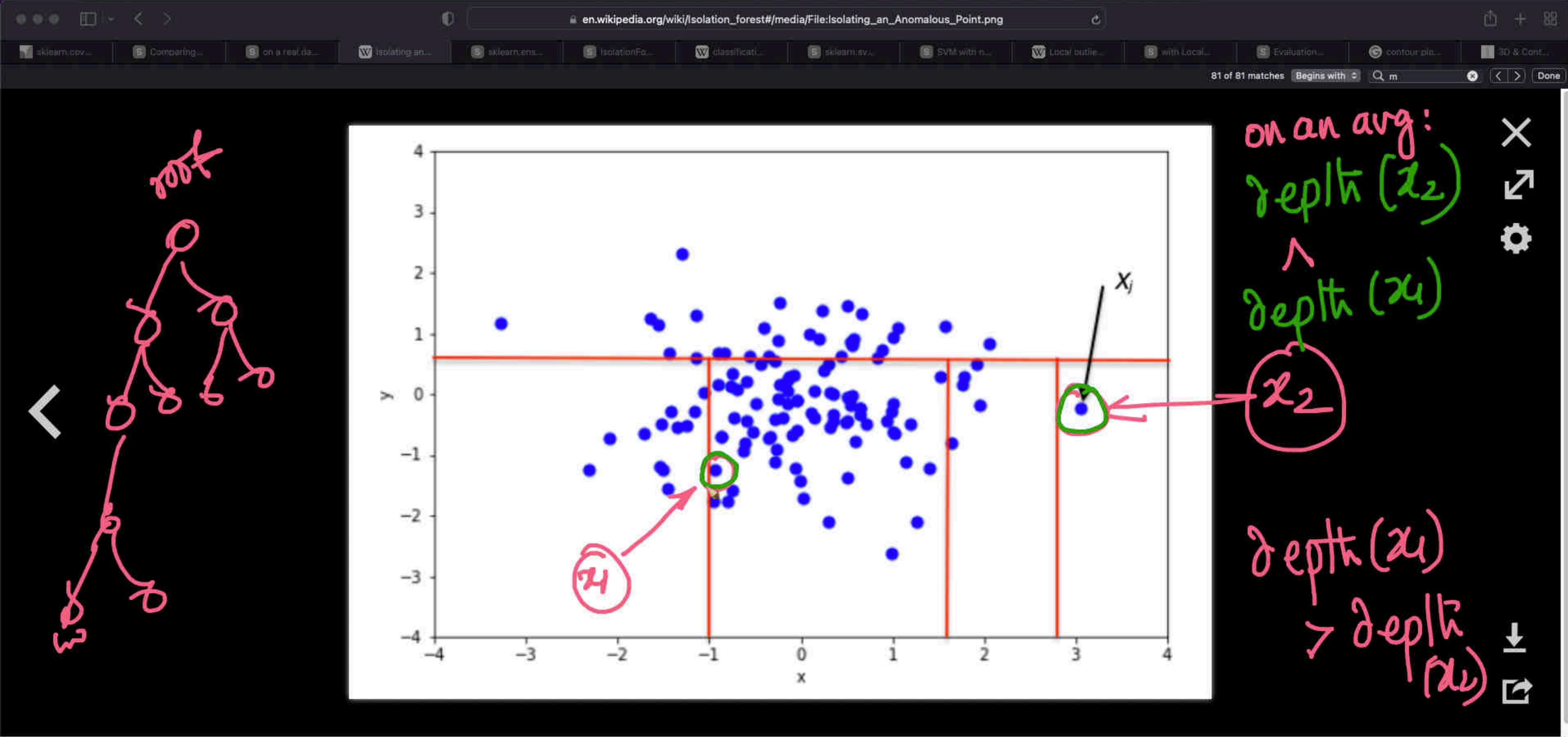


Fig. 3 - an example of isolating an anomalous point in a 2D Gaussian distribution.

More details

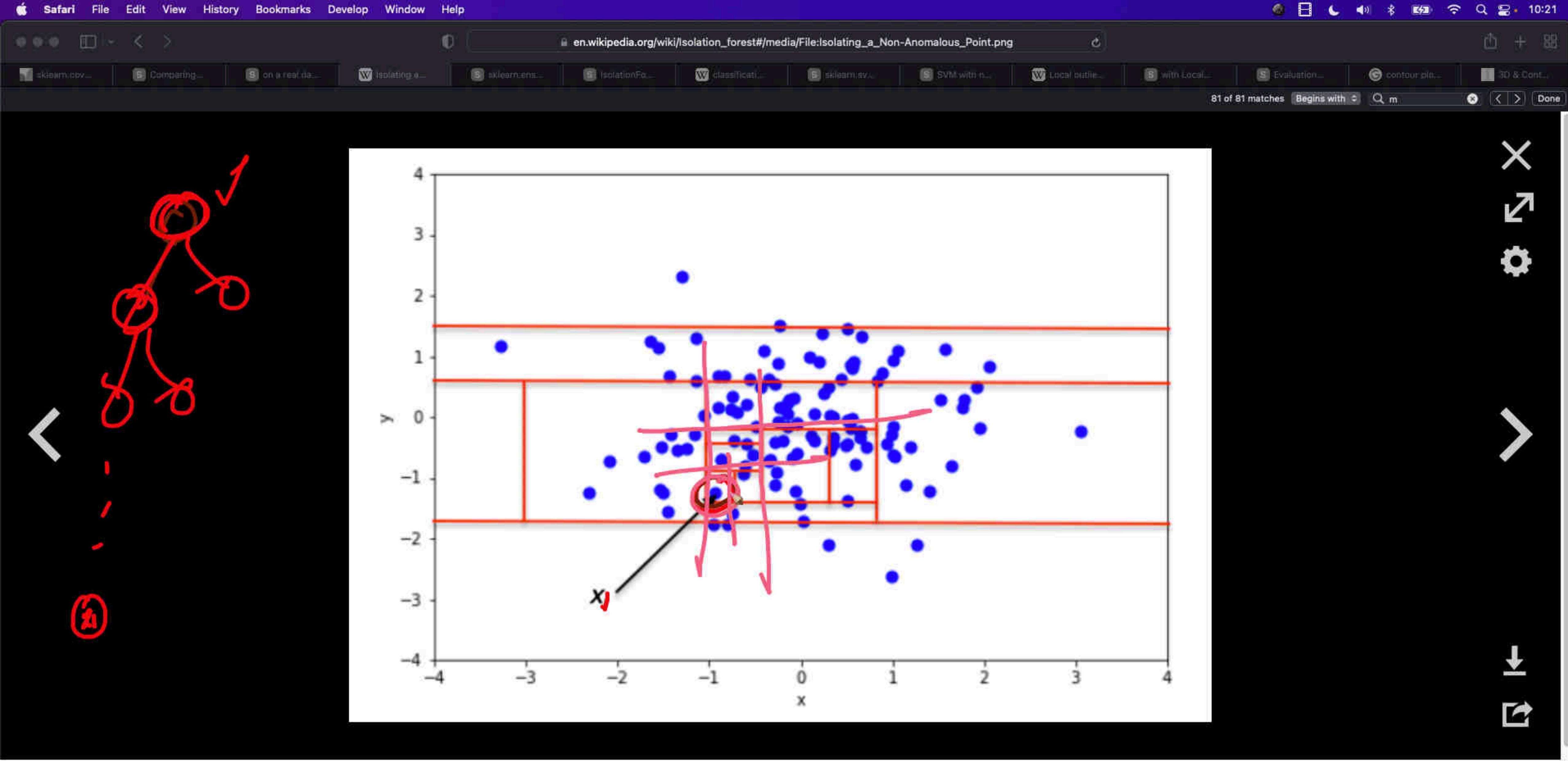


Fig. 2 - an example of isolating a non-anomalous point in a 2D Gaussian distribution.

More details

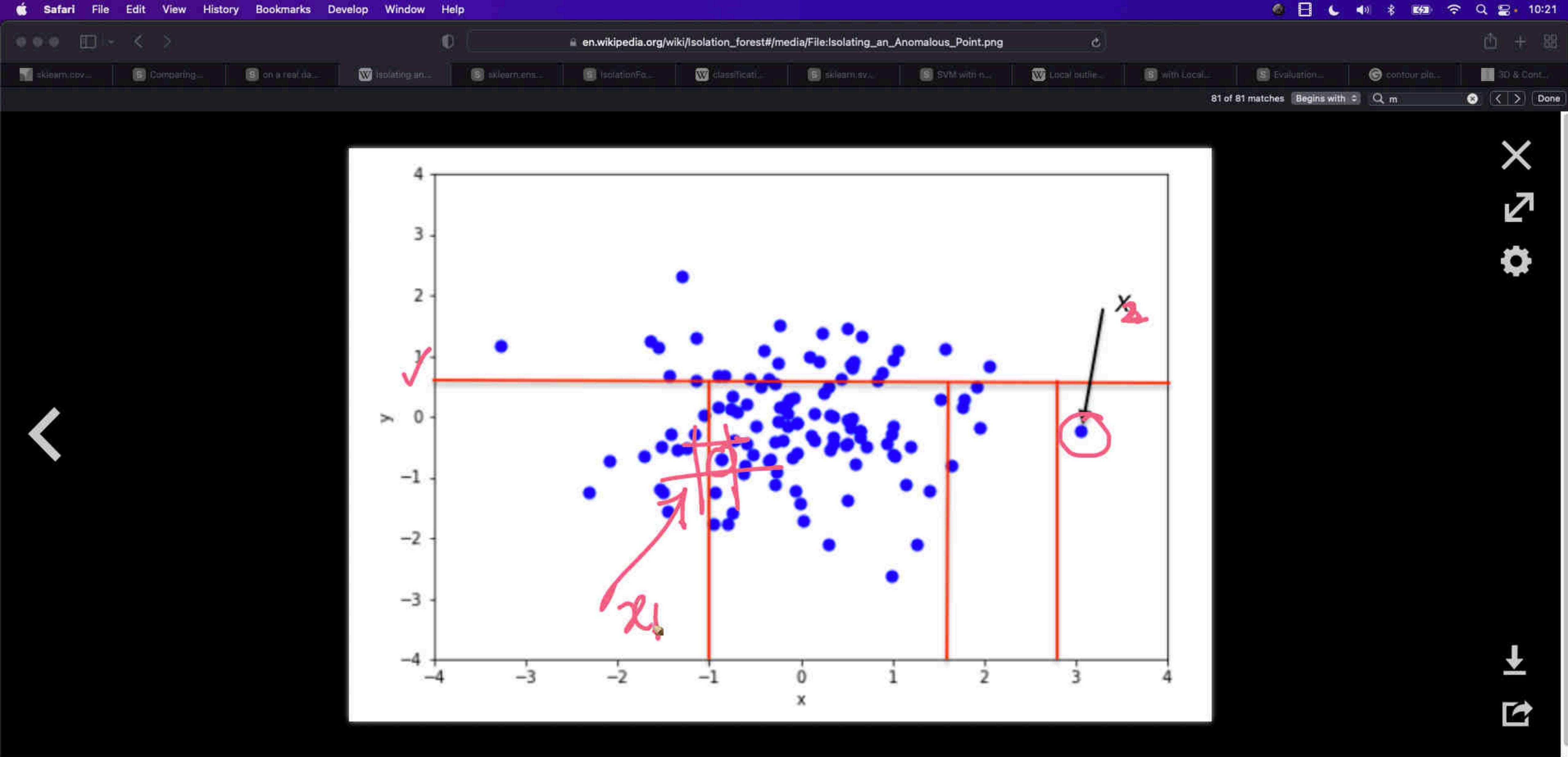


Fig. 3 - an example of isolating an anomalous point in a 2D Gaussian distribution.

More details

Intuition

→ outliers → lower depth on an arg
(fewer cuts) in the random trees

→ inliers → more depth on an arg
(more cuts)

Isolate each pt
Isolation forest ✓

100 → random-trees

$y_i ; x_i \rightarrow$ avg depth →
across all
100 trees

one setup:
 $\mathcal{D} = \{x_i\}_{i=1}^n$

↓
metric (lots of metrics)
✓

2007 - 2014

Different - Setup:-

✓ $\mathcal{D}_{Tr} = \{x_i\}_{i=1}^n$ no y_i 's \rightarrow iForest ✓

✓ $\{ \mathcal{D}_{Te} = \{x_i\}_{i=1}^m \}$ outliers
 x_q → avg. depth of x_q across all trees
iForest



$\uparrow n = \# \text{ data points}$

$\underbrace{}$

\downarrow

Sample of

DIV

$n < n_{\text{max}}$

iForest

: density
↑
depth when
isolation

DTree

As a consequence, the estimation of average $h(x)$ for external node terminations is the same as that of the unsuccessful searches in BST, that is^[11]

$$c(m) = \begin{cases} 2H(m-1) - \frac{2(m-1)}{n} & \text{for } m > 2 \\ 1 & \text{for } m = 2 \\ 0 & \text{otherwise} \end{cases}$$

$$\frac{-E[h(x)]}{2 c(m)}$$

where n is the testing data size, m is the size of the sample set and H is the harmonic number, which can be estimated by $H(i) = \ln(i) + \gamma$, where $\gamma = 0.5772156649$ is the Euler-Mascheroni constant.

The value of $c(m)$ above represents the average of $h(x)$ given m , so we can use it to normalise $h(x)$ and get an estimation of the anomaly score for a given instance x :

$$s(x, m) = 2^{\frac{-E(h(x))}{c(m)}}$$

$$E[h(x)] = \text{mean-height} + z$$

where $E(h(x))$ is the average value of $h(x)$ from a collection of iTrees. It is interesting to note that for any given instance x :

- if s is close to 1 then x is very likely to be an anomaly
- if s is smaller than 0.5 then x is likely to be a normal value
- if for a given sample all instances are assigned an anomaly score of around 0.5, then it is safe to assume that the sample doesn't have any anomaly

0.9

Open source implementations [edit]

Original implementation:

- [Isolation Forest](#) an algorithm that detects data anomalies using binary trees written in R. Released by the paper's first author Liu,



avg. depth of x_i



Outliers

asely

labelled outliers \rightarrow
 \equiv
 $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$
y_i : inliers

highly imbalanced
binary
classfn.



scikit-learn.org/stable/modules/generated/sklearn.ensemble.IsolationForest.html

Install User Guide API Examples Community More ▾ Go

Prev Up Next

scikit-learn 1.1.1

Other versions

Please [cite us](#) if you use the software.

sklearn.ensemble.IsolationForest

```
class sklearn.ensemble.IsolationForest(*, n_estimators=100, max_samples='auto', contamination='auto',  
max_features=1.0, bootstrap=False, n_jobs=None, random_state=None, verbose=0,  
warm_start=False)
```

[source]

Isolation Forest Algorithm.

Return the anomaly score of each sample using the IsolationForest algorithm

The IsolationForest ‘isolates’ observations by randomly selecting a feature and then randomly selecting a split value between the maximum and minimum values of the selected feature.

Since recursive partitioning can be represented by a tree structure, the number of splittings required to isolate a sample is equivalent to the path length from the root node to the terminating node.

This path length, averaged over a forest of such random trees, is a measure of normality and our decision function.

Random partitioning produces noticeably shorter paths for anomalies. Hence, when a forest of random trees

samples, they are highly likely to be anomalies.

Toggle Menu

scikit-learn.org/stable/modules/generated/sklearn.ensemble.IsolationForest.html

scikit-learn

Prev Up Next

scikit-learn 1.1.1

Other versions

Please cite us if you use the software.

sklearn.ensemble.IsolationForest

Examples using sklearn.ensemble.IsolationForest

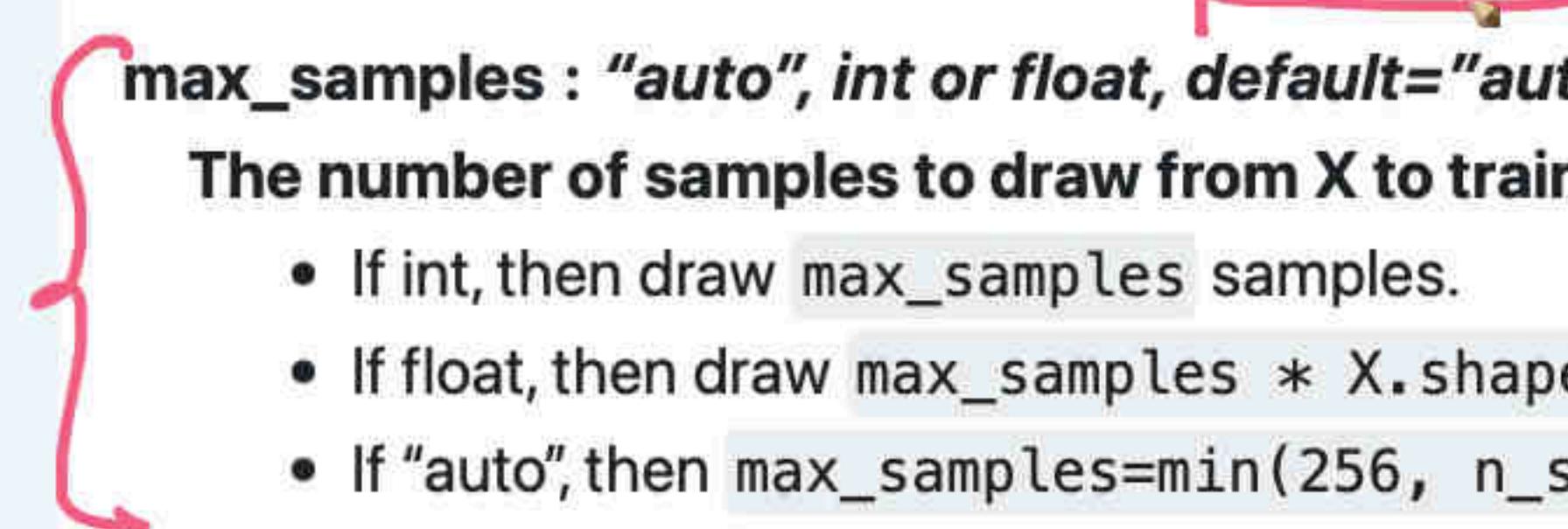
collectively produce shorter path lengths for particular samples, they are highly likely to be anomalies.

Read more in the [User Guide](#).

New in version 0.18.

Parameters:

n_estimators : int, default=100
The number of base estimators in the ensemble.

max_samples : "auto", int or float, default="auto"
The number of samples to draw from X to train each base estimator.
• If int, then draw max_samples samples.
• If float, then draw max_samples * X.shape[0] samples.
• If "auto", then max_samples=min(256, n_samples).


If max_samples is larger than the number of samples provided, all samples will be used for all trees (no sampling).

contamination : 'auto' or float, default='auto'
The amount of contamination of the data set, i.e. the proportion of outliers in the data set.
Used when fitting to define the threshold on the scores of the samples.
• If 'auto', the threshold is determined as in the original paper.
• If float, the contamination should be in the range (0, 0.5].

Toggle Menu

scikit-learn.org/stable/modules/generated/sklearn.ensemble.IsolationForest.html

View in version 0.16.

scikit learn

Prev Up Next

scikit-learn 1.1.1

Other versions

Please cite us if you use the software.

sklearn.ensemble.IsolationForest

Examples using sklearn.ensemble.IsolationForest

Parameters:

n_estimators : int, default=100

The number of base estimators in the ensemble.

max_samples : "auto", int or float, default="auto"

The number of samples to draw from X to train each base estimator.

- If int, then draw max_samples samples.
- If float, then draw max_samples * X.shape[0] samples.
- If "auto", then max_samples=min(256, n_samples).

If max_samples is larger than the number of samples provided, all samples will be used for all trees (no sampling).

contamination : 'auto' or float, default='auto'

The amount of contamination of the data set, i.e. the proportion of outliers in the data set. Used when fitting to define the threshold on the scores of the samples.

- If 'auto', the threshold is determined as in the original paper.
- If float, the contamination should be in the range (0, 0.5].

Changed in version 0.22: The default value of contamination changed from 0.1 to 'auto'.

Toggle Menu

1.0

49 / 49

scikit-learn.org/stable/modules/generated/sklearn.ensemble.IsolationForest.html

scikit-learn

Prev Up Next

scikit-learn 1.1.1

Other versions

Please cite us if you use the software.

sklearn.ensemble.IsolationForest Examples using sklearn.ensemble.IsolationForest

1 Liu, Fei Tony, Ting, Kai Ming and Zhou, Zhi-Hua. "Isolation forest." Data Mining, 2008. ICDM'08. Eighth IEEE International Conference on.

2 Liu, Fei Tony, Ting, Kai Ming and Zhou, Zhi-Hua. "Isolation-based anomaly detection." ACM Transactions on Knowledge Discovery from Data (TKDD) 6.1 (2012): 3.

Examples

```
>>> from sklearn.ensemble import IsolationForest  
>>> X = [[-1.1], [0.3], [0.5], [100]]  
>>> clf = IsolationForest(random_state=0).fit(X)  
>>> clf.predict([[0.1], [0], [90]])  
array([ 1,  1, -1])
```

Methods

<code>decision_function(X)</code>	Average anomaly score of X of the base classifiers.
<code>fit(X[, y, sample_weight])</code>	Fit estimator.
<code>fit_predict(X[, y])</code>	Perform fit on X and returns labels for X.
<code>get_params([deep])</code>	Get parameters for this estimator.
<code>predict(X)</code>	Predict if a particular sample is an outlier or not.
<code>score_samples(X)</code>	Opposite of the anomaly score defined in the original paper.
<code>set_params(**params)</code>	Set the parameters of this estimator.

[source]

scikit-learn.org/stable/auto_examples/miscellaneous/plot_anomaly_comparison.html#sphx-glr-auto-examples-

scikit learn

Prev Up Next

scikit-learn 1.1.1

Other versions

Please cite us if you use the software.

Comparing anomaly detection algorithms for outlier detection on toy datasets

Robust covariance One-Class SVM One-Class SVM (SGD) Isolation Forest Local Outlier Factor

.03s .00s .01s .12s .00s

.03s .00s .01s .13s .00s

.03s .00s .01s .13s .00s

Toggle Menu

scikit-learn.org/stable/auto_examples/miscellaneous/plot_anomaly_comparison.html#sphx-glr-auto-examples-

scikit learn

Prev Up Next

scikit-learn 1.1.1

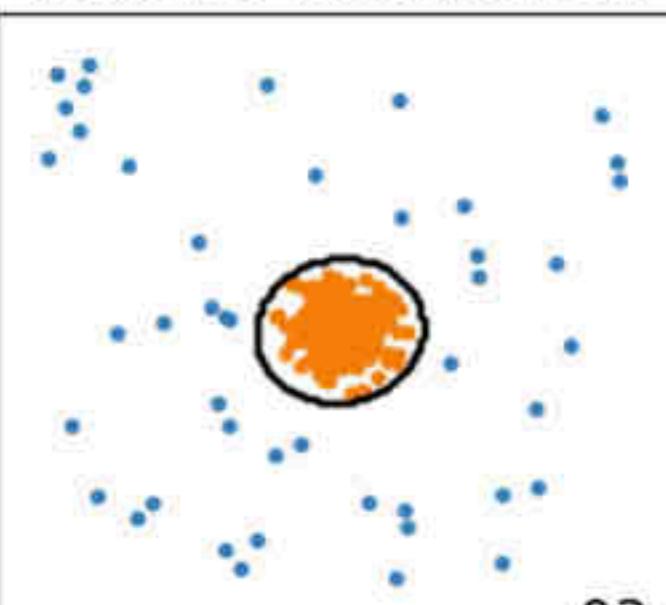
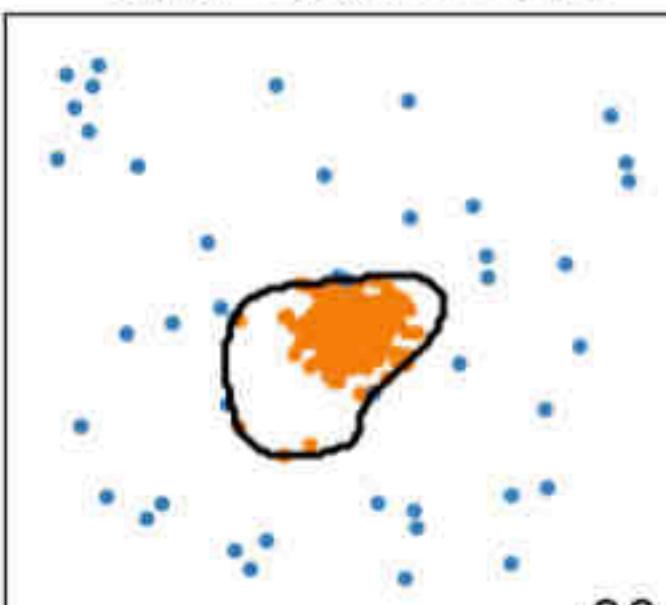
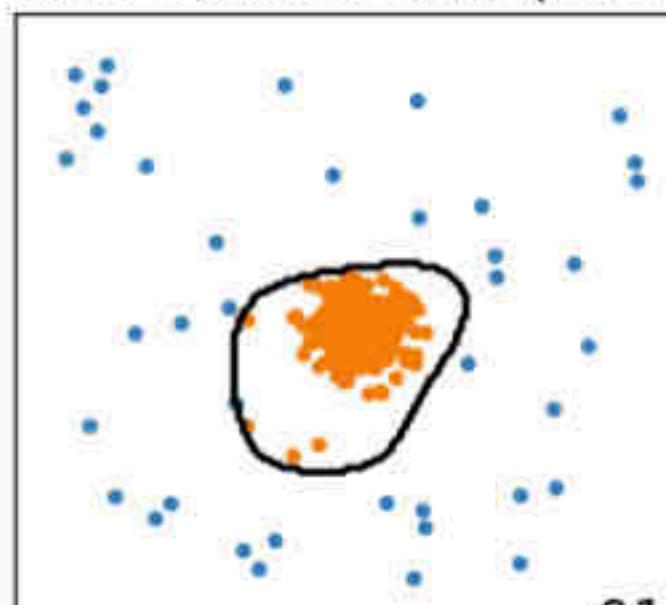
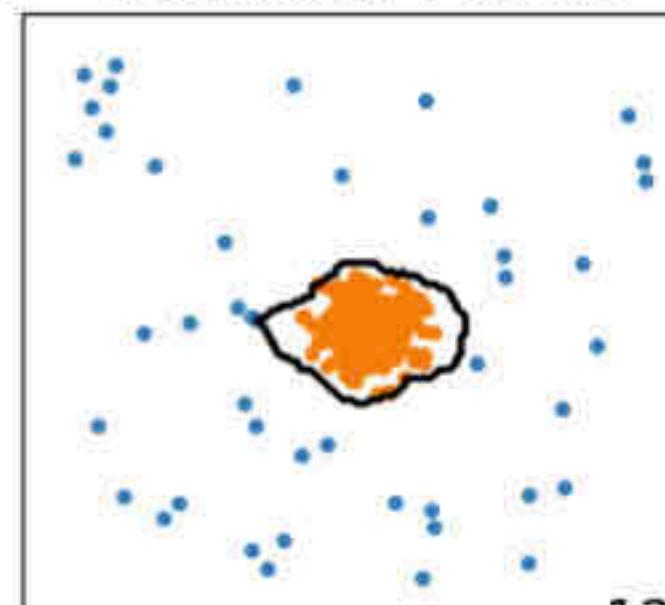
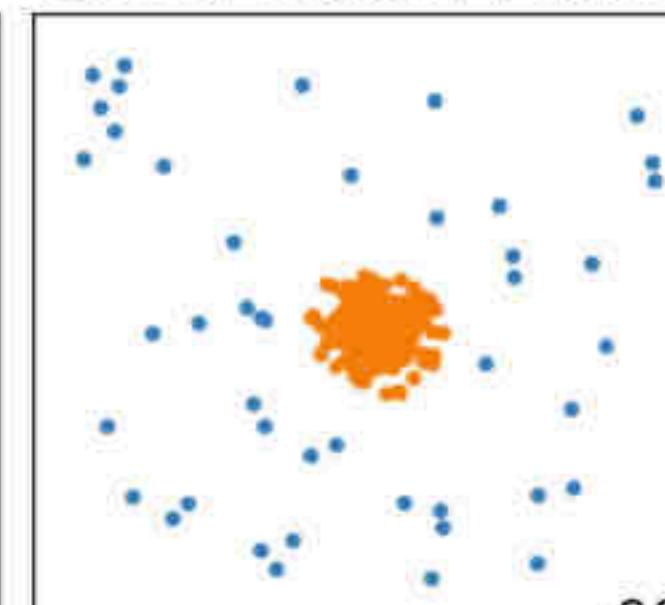
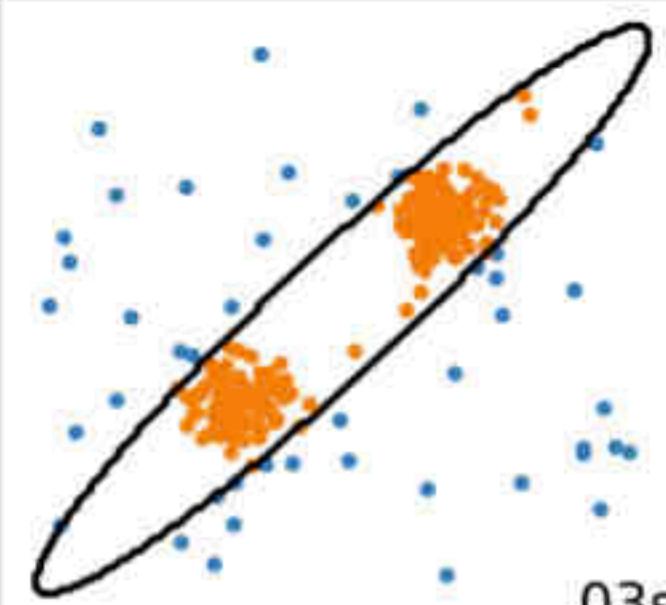
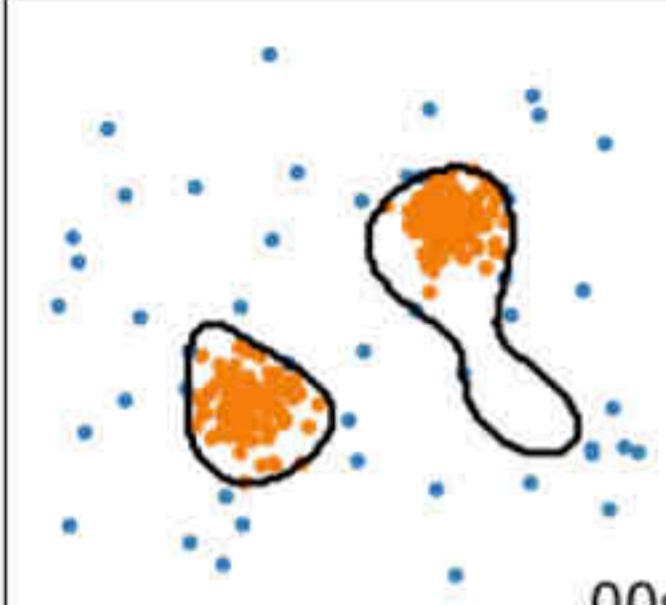
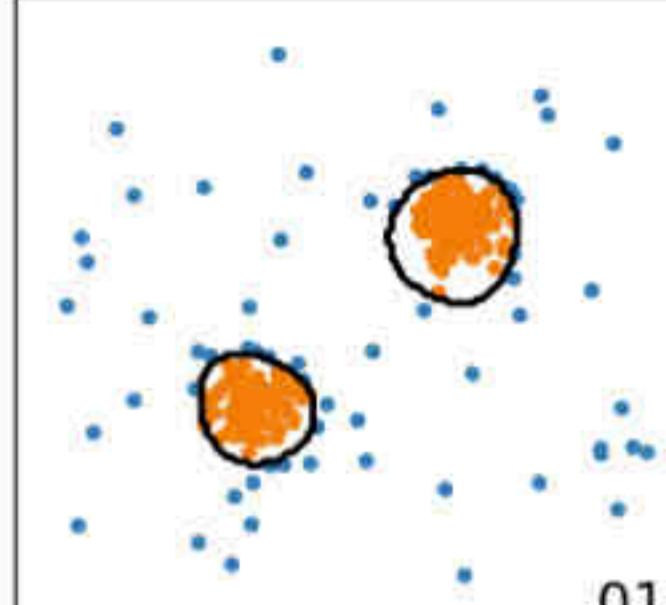
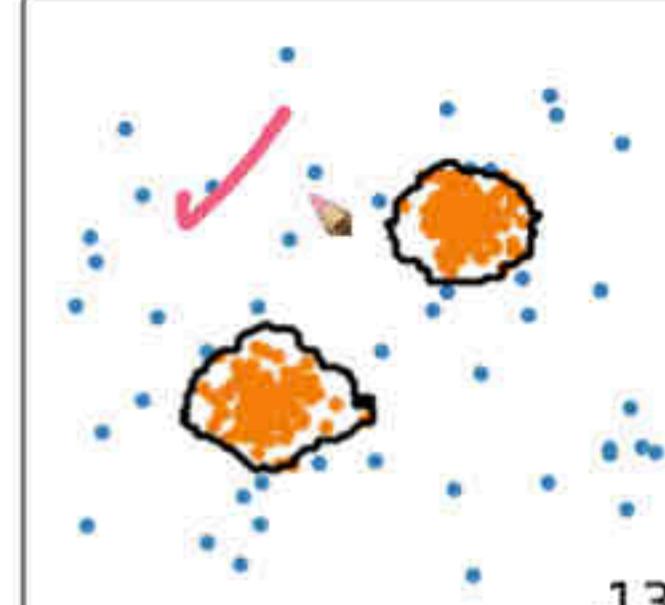
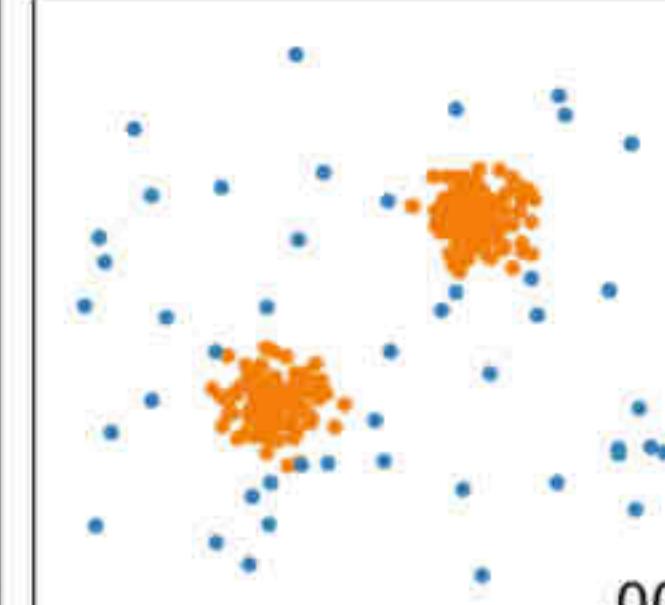
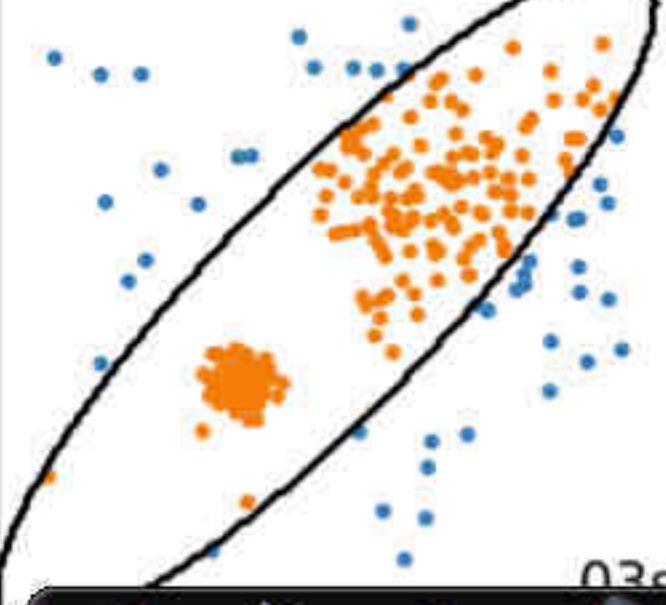
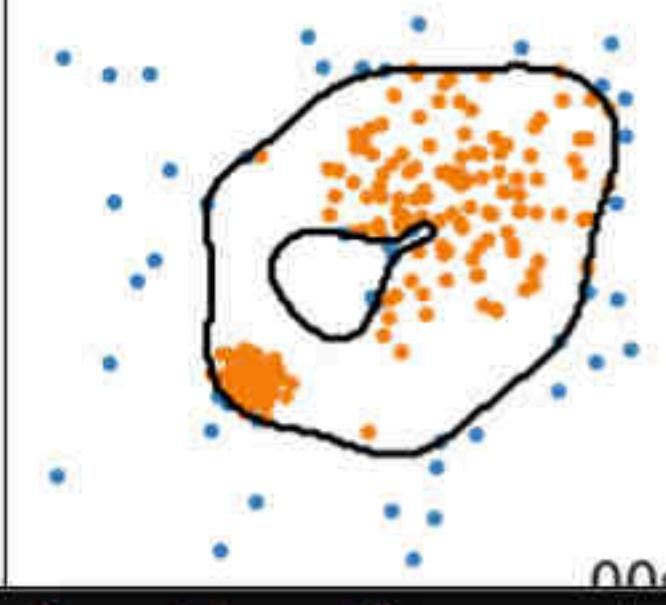
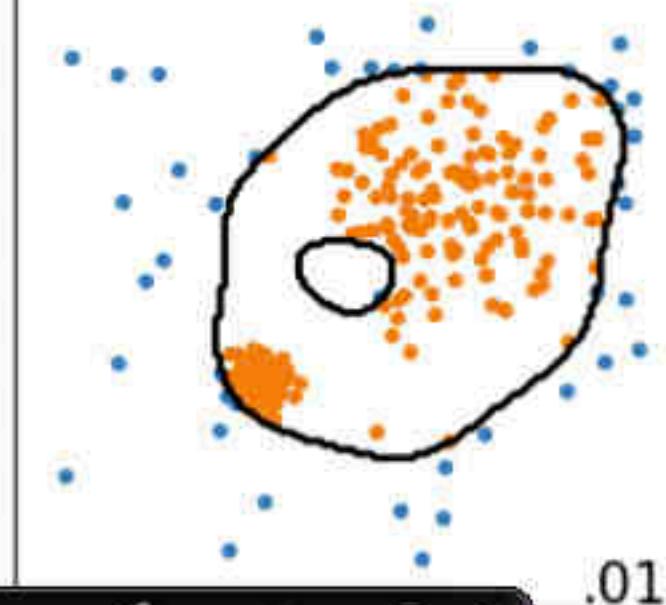
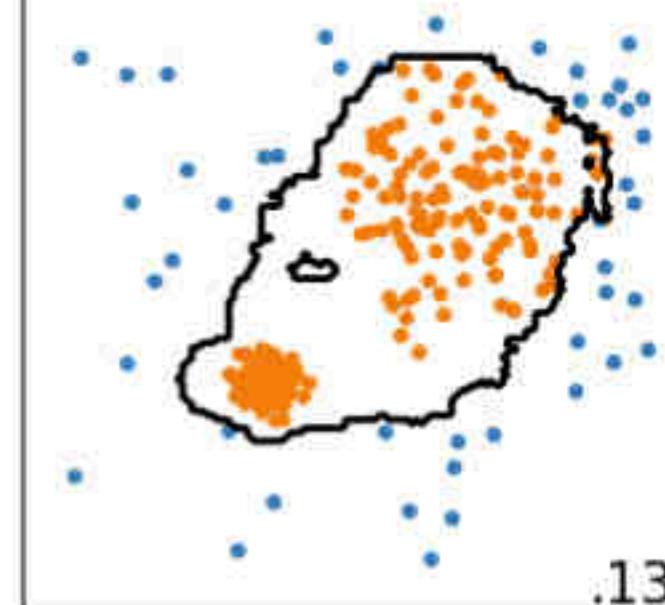
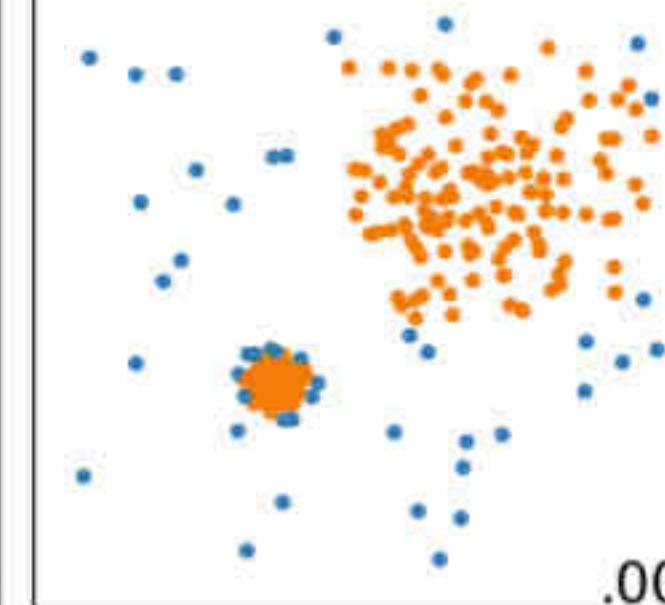
Other versions

Please cite us if you use the software.

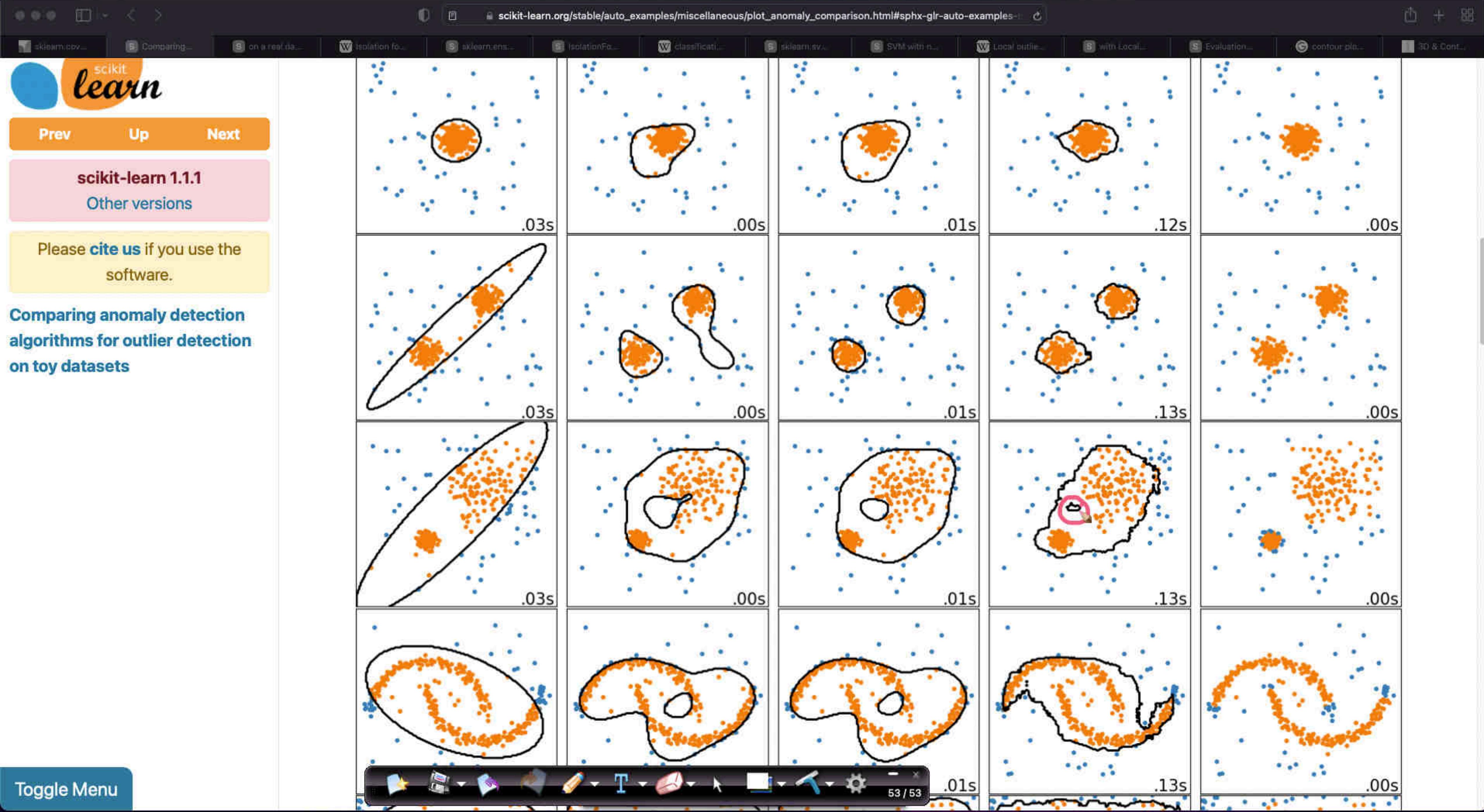
Comparing anomaly detection algorithms for outlier detection on toy datasets

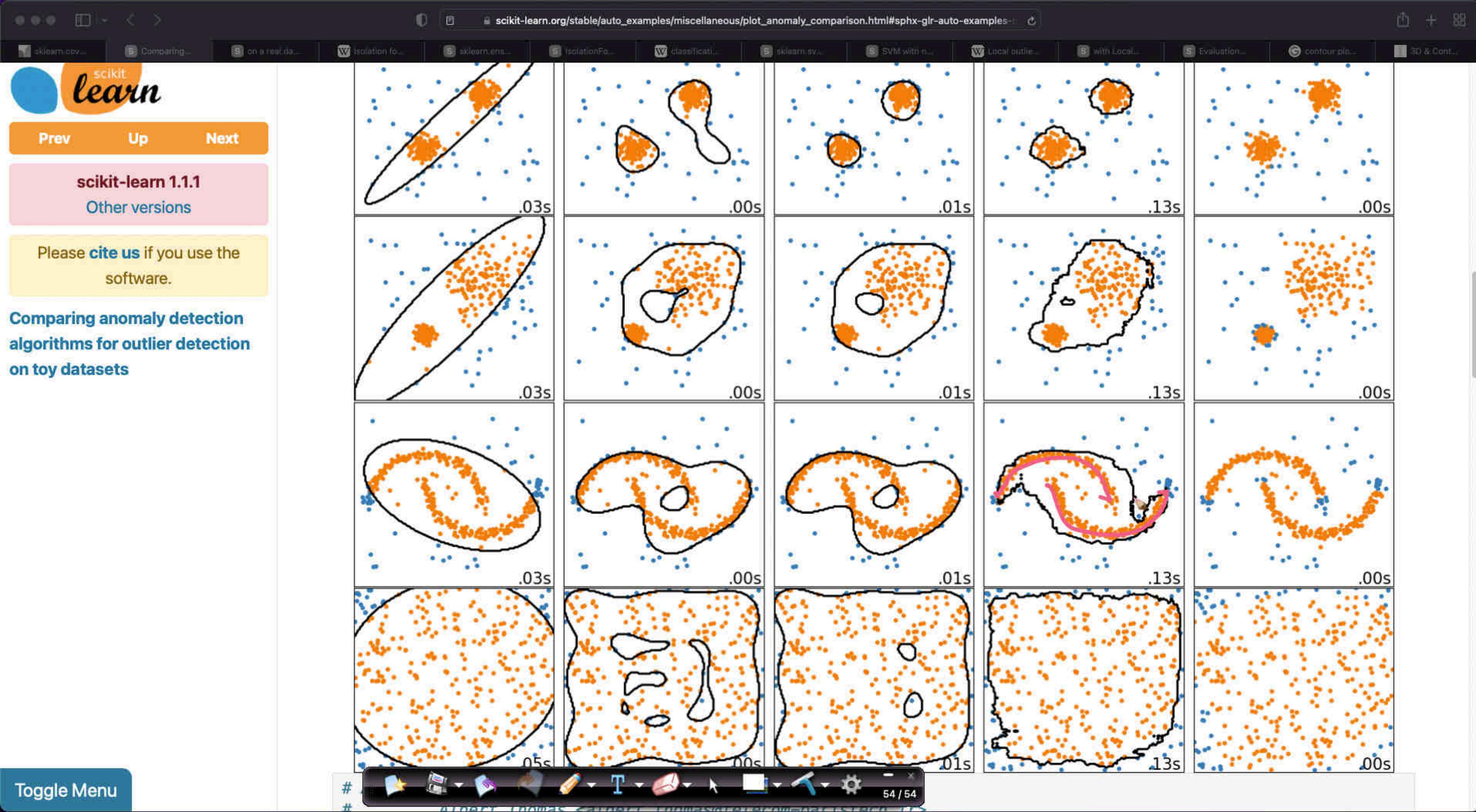
While these examples give some intuition about the algorithms, this intuition might not apply to very high dimensional data.

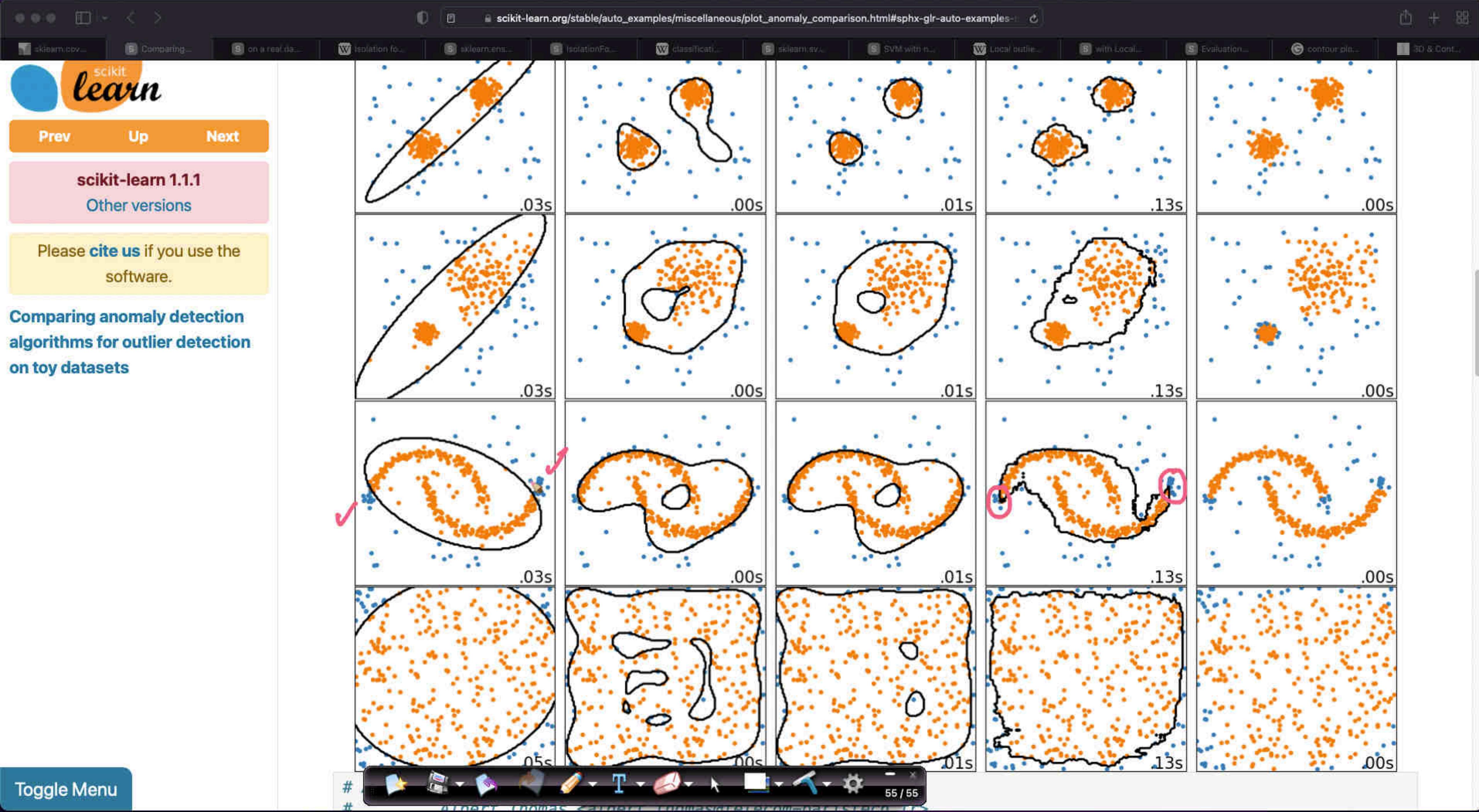
Finally, note that parameters of the models have been here handpicked but that in practice they need to be adjusted. In the absence of labelled data, the problem is completely unsupervised so model selection can be a challenge.

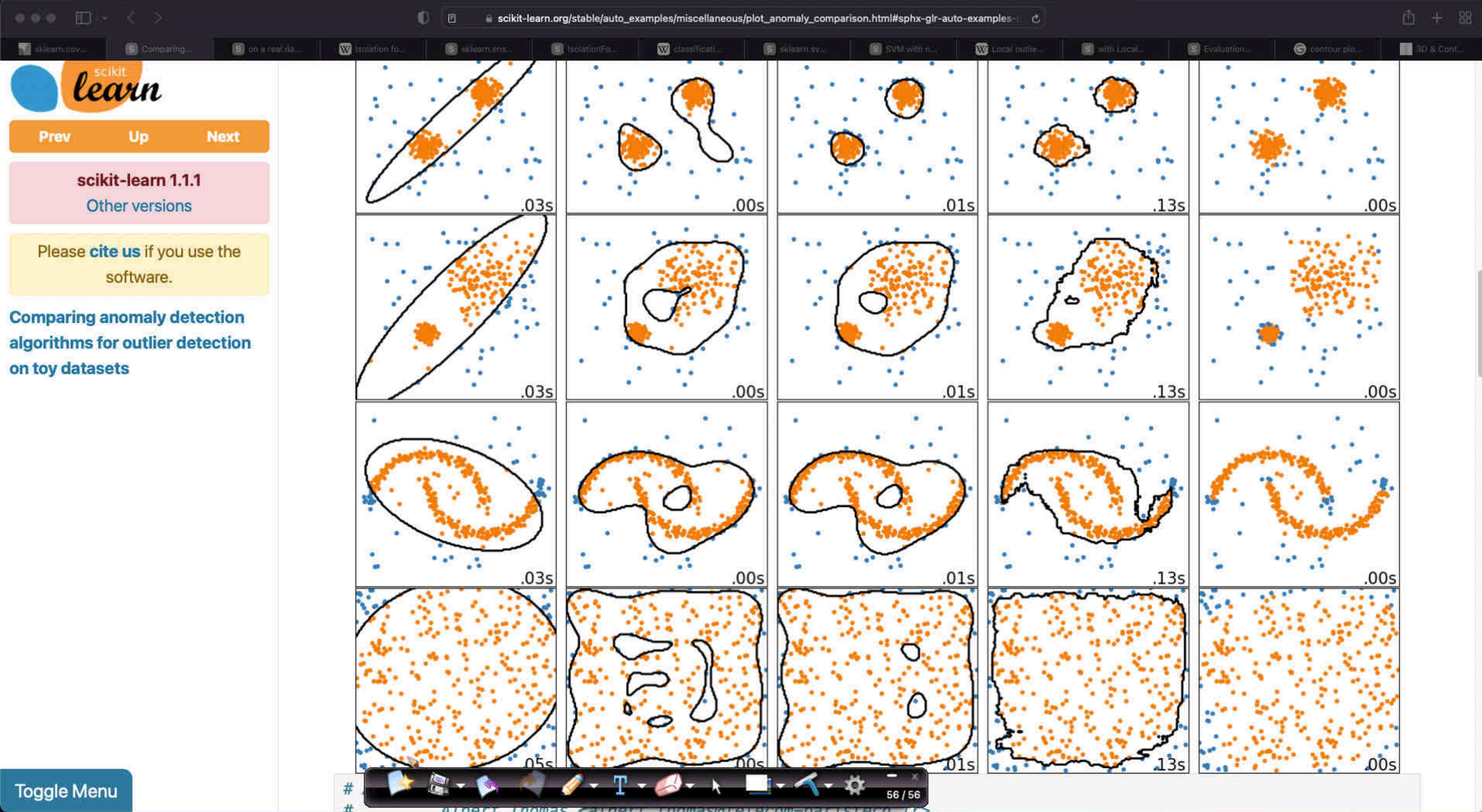
	Robust covariance	One-Class SVM	One-Class SVM (SGD)	Isolation Forest	Local Outlier Factor
.03s					
.00s					
.03s					

Toggle Menu









disadv

↳ contamination-param → heuristics

↳ axis parallel splits



Prev Up Next

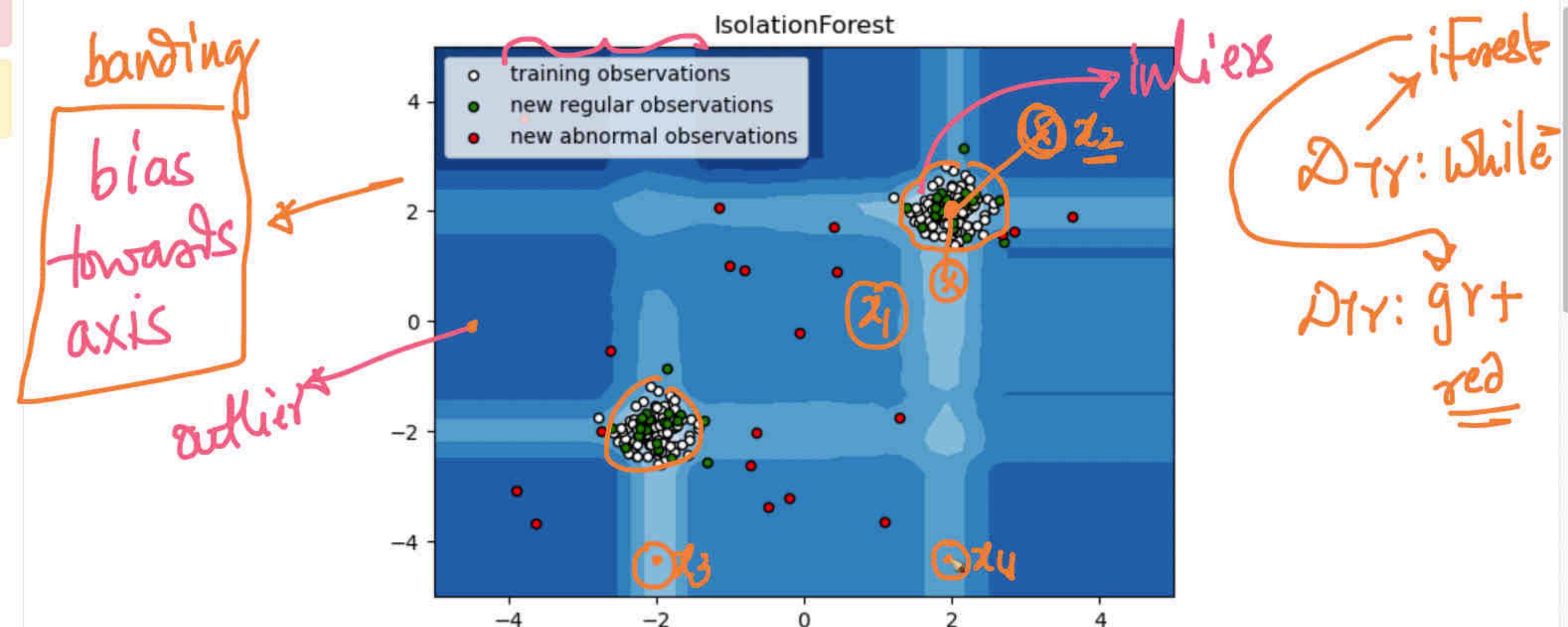
scikit-learn 1.1.1

Other versions

Please cite us if you use the software.

IsolationForest example

Random partitioning produces noticeable shorter paths for anomalies. Hence, when a forest of random trees collectively produce shorter path lengths for particular samples, they are highly likely to be anomalies.



```
import numpy as np  
import matplotlib.pyplot as plt  
from
```

Toggle Menu

hyper-param

slack

$$\min_{C, \gamma, \xi_i} \frac{1}{2} \gamma^2 + \lambda \sum_{i=1}^n \xi_i$$

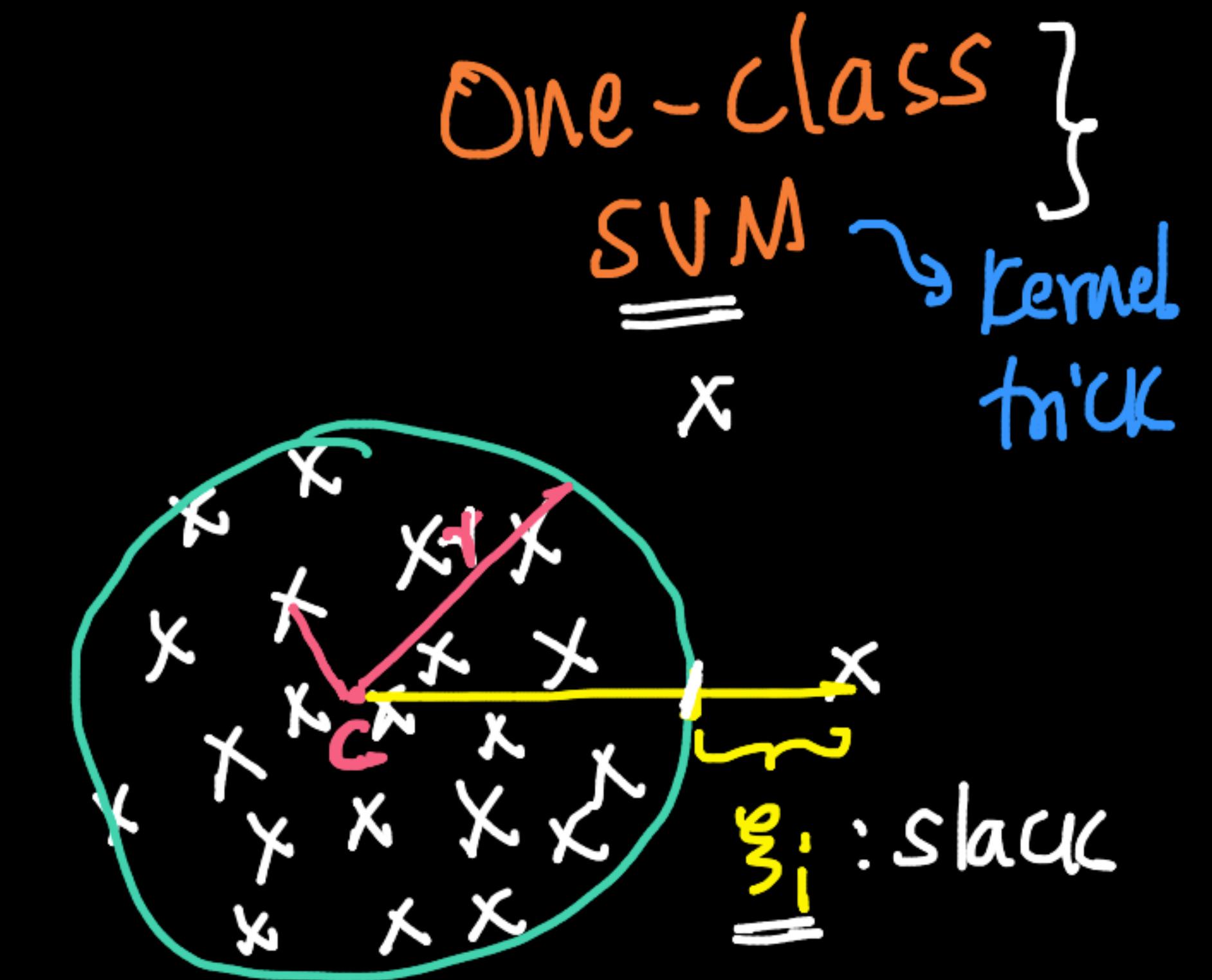
s.t.

$$\text{dist}(x_i, C) \leq \gamma + \xi_i \quad \forall i$$

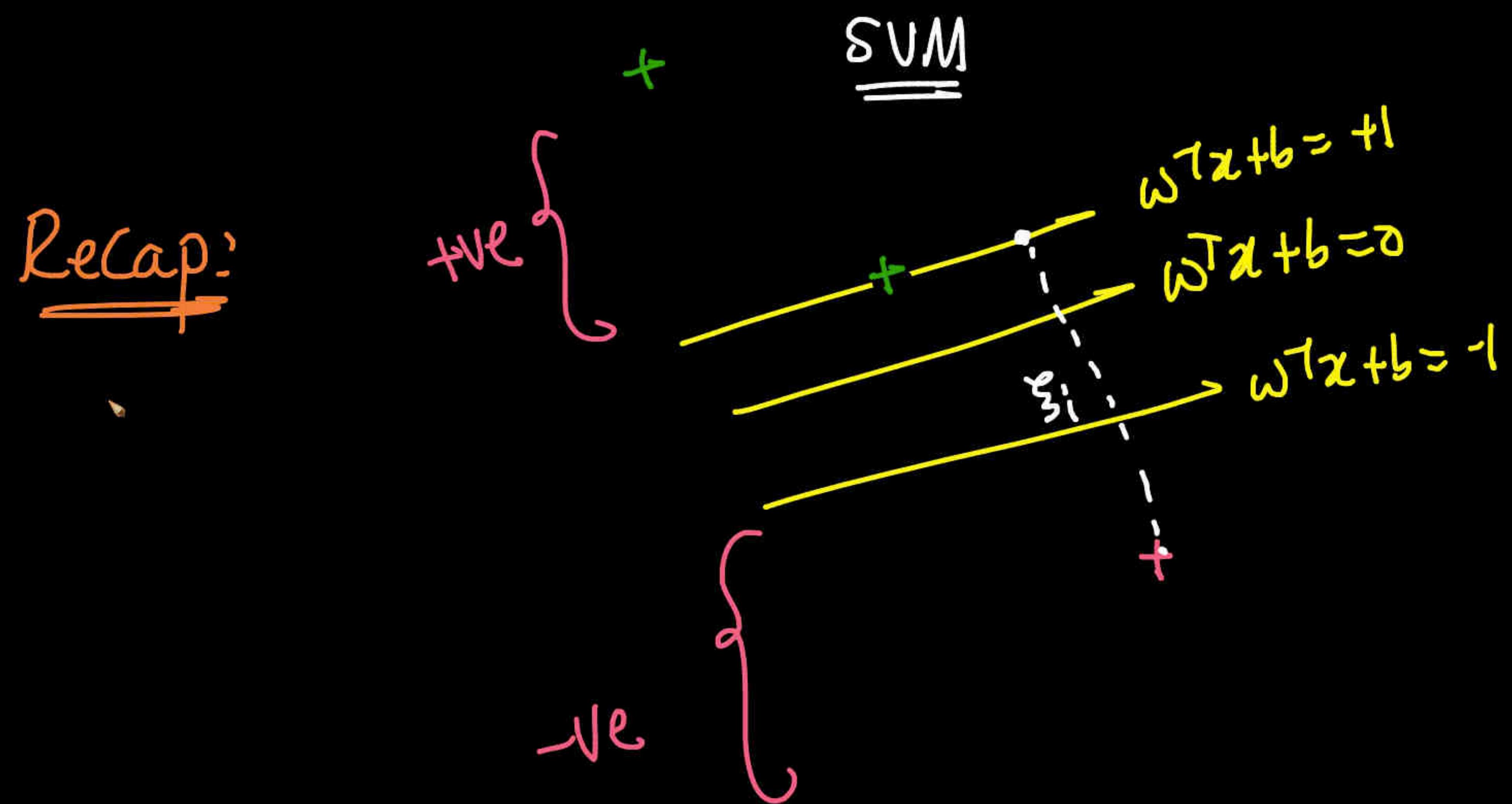
$$\xi_i \geq 0$$

SVM-Margin

RBF Kernelize This



$$\xi_i = \begin{cases} 0 & \text{if } d(x_i, C) \leq \gamma \\ d(x_i, C) - \gamma & \text{otherwise} \end{cases}$$



Original-Daten \rightarrow 2D

RBF-Kernel one-class SVM \rightarrow any shape

Kernel Space: hyper spheres

disadv

→ Kernel Selection (RBF: default)

→ as $n \uparrow$: time complex increases → SGD

→ SVMs disadvantages

scikit-learn.org/stable/modules/generated/sklearn.svm.OneClassSVM.html

scikit-learn Install User Guide API Examples Community More ▾ Go

Prev Up Next

scikit-learn 1.1.1 Other versions

Please cite us if you use the software.

sklearn.svm.OneClassSVM Examples using sklearn.svm.OneClassSVM

sklearn.svm.OneClassSVM

```
class sklearn.svm.OneClassSVM(*, kernel='rbf', degree=3, gamma='scale', coef0=0.0, tol=0.001, nu=0.5, shrinking=True, cache_size=200, verbose=False, max_iter=-1)
```

[source]

Unsupervised Outlier Detection.

Estimate the support of a high-dimensional distribution.

The implementation is based on libsvm.

Read more in the [User Guide](#).

Parameters:

kernel : {‘linear’, ‘poly’, ‘rbf’, ‘sigmoid’, ‘precomputed’} or callable, default=‘rbf’
Specifies the kernel type to be used in the algorithm. If none is given, ‘rbf’ will be used. If a callable is given it is used to precompute the kernel matrix.

degree : int, default=3
Degree of the polynomial kernel function (‘poly’). Ignored by all other kernels.

gamma : {‘scale’, ‘auto’} or float, default=‘scale’
Kernel coefficient for ‘rbf’, ‘poly’ and ‘sigmoid’.

- if gamma=‘scale’ (default) is passed then it uses $1 / (\text{n_features} * \text{X.var()})$ as value of gamma,
- if ‘auto’ uses $1 / \text{n_features}$

Toggle Menu

scikit-learn.org/stable/modules/generated/sklearn.svm.OneClassSVM.html

scikit-learn

Prev Up Next

scikit-learn 1.1.1

Other versions

Please cite us if you use the software.

sklearn.svm.OneClassSVM

Examples using sklearn.svm.OneClassSVM

gamma : float, default='scale'

- if `gamma='scale'` (default) is passed then it uses $1 / (\text{n_features} * \text{X.var()})$ as value of gamma,
- if 'auto', uses $1 / \text{n_features}$.

Changed in version 0.22: The default value of `gamma` changed from 'auto' to 'scale'.

coef0 : float, default=0.0

Independent term in kernel function. It is only significant in 'poly' and 'sigmoid'.

tol : float, default=1e-3

Tolerance for stopping criterion.

nu : float, default=0.5

An upper bound on the fraction of training errors and a lower bound of the fraction of support vectors. Should be in the interval $(0, 1]$. By default 0.5 will be taken.

shinking : bool, default=True

Whether to use the shrinking heuristic. See the [User Guide](#).

cache_size : float, default=200

Specify the size of the kernel cache (in MB).

verbose : bool, default=False

Enable verbose output. Note that this setting takes advantage of a per-process runtime setting in libsvm that, if enabled, may not work properly in a multithreaded context.

max_iter : int, default=-1

Hard limit on iterations within solver, or -1 for no limit.

ξ_i > 0

Toggle Menu

[Prev](#) [Up](#) [Next](#)**scikit-learn 1.1.1**[Other versions](#)

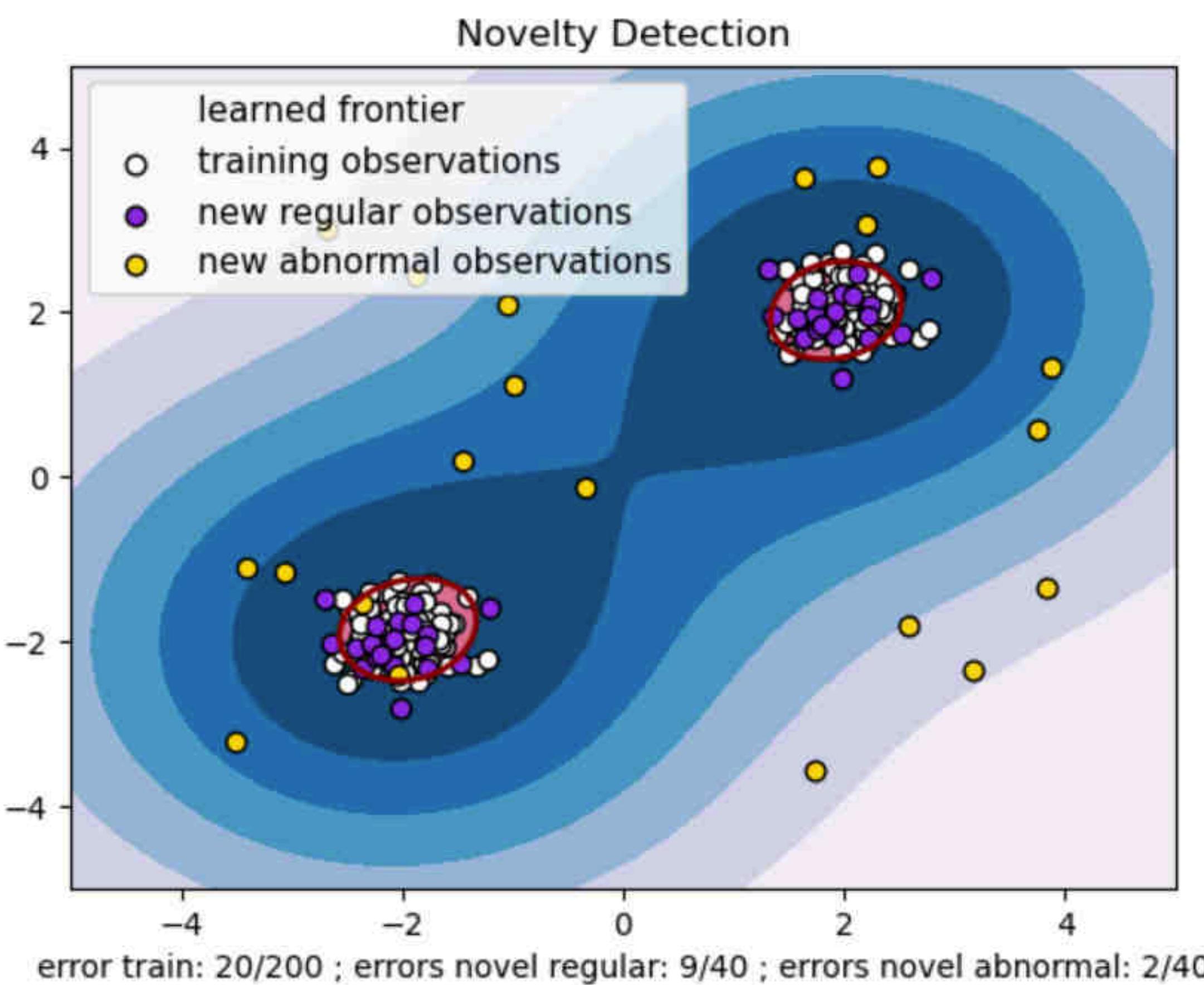
Please [cite us](#) if you use the software.

[One-class SVM with non-linear kernel \(RBF\)](#)

RBF

An example using a one-class SVM for novelty detection.

One-class SVM is an unsupervised algorithm that learns a decision function for novelty detection: classifying new data as similar or different to the training set.



D_r → 1 class SVM
D_{te}: purple + yellow

scikit-learn.org/stable/auto_examples/svm/plot_oneclass.html#sphx-glr-auto-examples-svm-plot-oneclass-py

sklearn cov... Comparing... on a real da... Isolation fo... scikitlearn... IsolationFo... classificat... SVM with n... Local outlier... with local... Evaluation... contour plot... 3D & Cont...

scikit learn Install User Guide API Examples Community More ▾ Go

Prev Up Next

scikit-learn 1.1.1 Other versions

Please cite us if you use the software.

One-class SVM with non-linear kernel (RBF)

Novelty Detection

learned frontier

- training observations
- new regular observations
- new abnormal observations

error train: 20/200 ; errors novel regular: 9/40 ; errors novel abnormal: 2/40

```
import numpy as np
import matplotlib.pyplot as plt
import matplotlib.font_manager
from sklearn import svm
```

XX, yy = np.meshgrid(np.linspace(-5, 5, 500), np.linspace(-5, 5, 500))

#

Toggle Menu

scikit-learn.org/stable/auto_examples/svm/plot_onesclass.html#sphx-glr-auto-examples-svm-plot-onesclass-py

sklearn cov... Comparing... on a real da... Isolation fo... scikitlearn... IsolationFo... classificat... SVM with n... Local outlier... with local... Evaluation... contour plot... 3D & Cont...

scikit learn Install User Guide API Examples Community More ▾ Go

Prev Up Next

scikit-learn 1.1.1 Other versions

Please cite us if you use the software.

One-class SVM with non-linear kernel (RBF)

Novelty Detection

learned frontier

- training observations
- new regular observations
- new abnormal observations

error train: 20/200 ; errors novel regular: 9/40 ; errors novel abnormal: 2/40

```
import numpy as np
import matplotlib.pyplot as plt
import matplotlib.font_manager
from sklearn import svm
```

XX, yy = np.meshgrid(np.linspace(-5, 5, 500), np.linspace(-5, 5, 500))

#

Toggle Menu

scikit-learn.org/stable/auto_examples/miscellaneous/plot_anomaly_comparison.html#sphx-glr-auto-examples-

While these examples give some intuition about the algorithms, this intuition might not apply to very high dimensional data.

Finally, note that parameters of the models have been here handpicked but that in practice they need to be adjusted. In the absence of labelled data, the problem is completely unsupervised so model selection can be a challenge.

Prev Up Next

scikit-learn 1.1.1

Other versions

Please cite us if you use the software.

Comparing anomaly detection algorithms for outlier detection on toy datasets

Robust covariance One-Class SVM One-Class SVM (SGD) Isolation Forest Local Outlier Factor

.03s .00s .01s .12s .00s

.03s .00s .01s .13s .00s

.03s .00s .01s .13s .00s

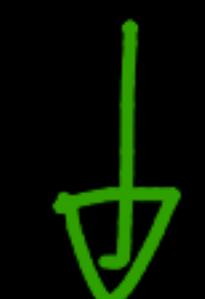
Toggle Menu

RBF

68 / 68

Local outlier

factor



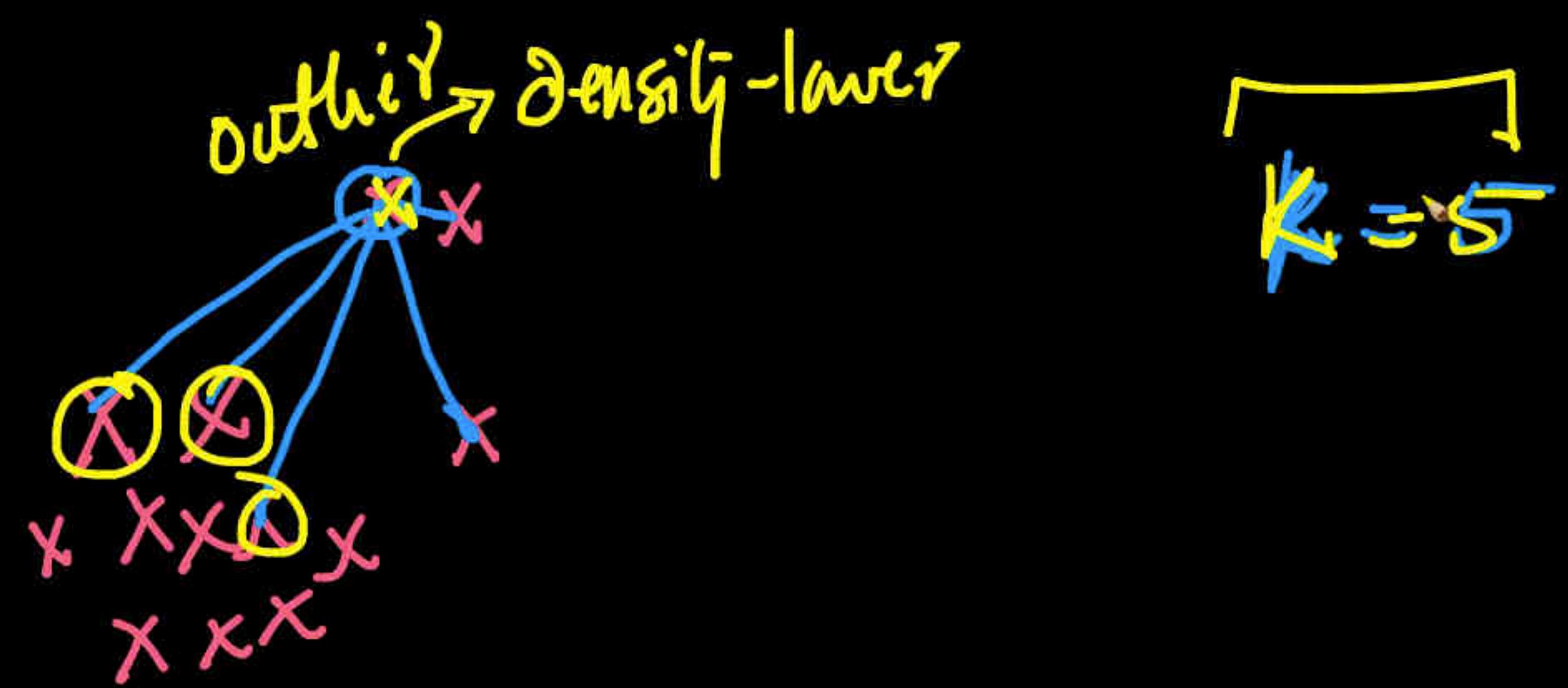
[
KNN
density]

Intuition

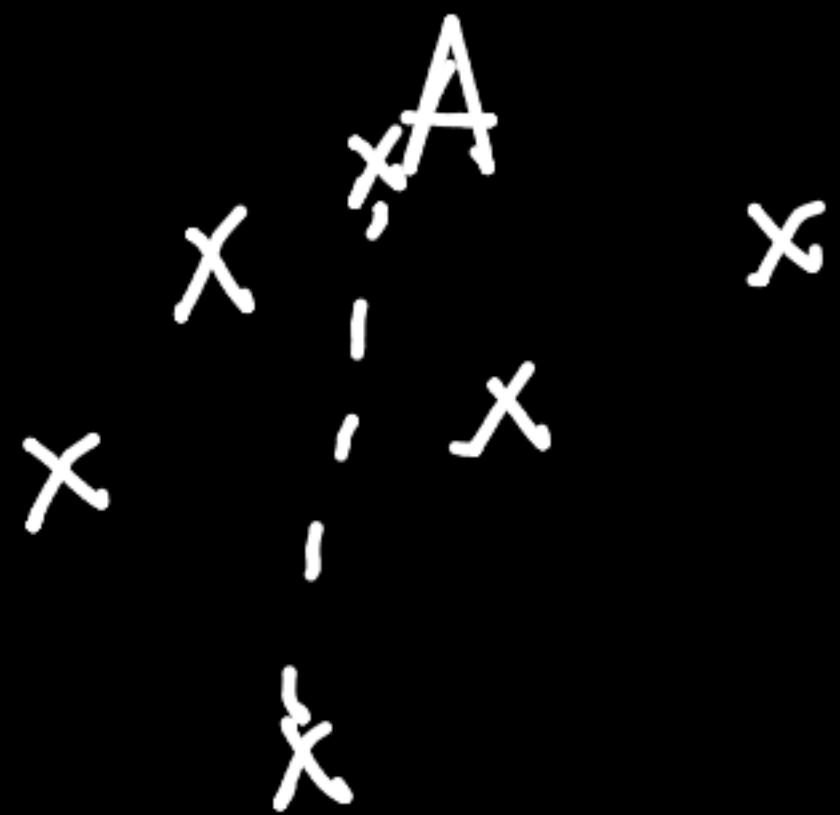
Compare a ~~pts~~ point with its K neighbors

low

high



① $k\text{-dist}(A)$ = distance from A to its k^{th} -NN
✓ ↳ high



② $\underline{N_k(A)}$ = set of k -nearest neighbors of A

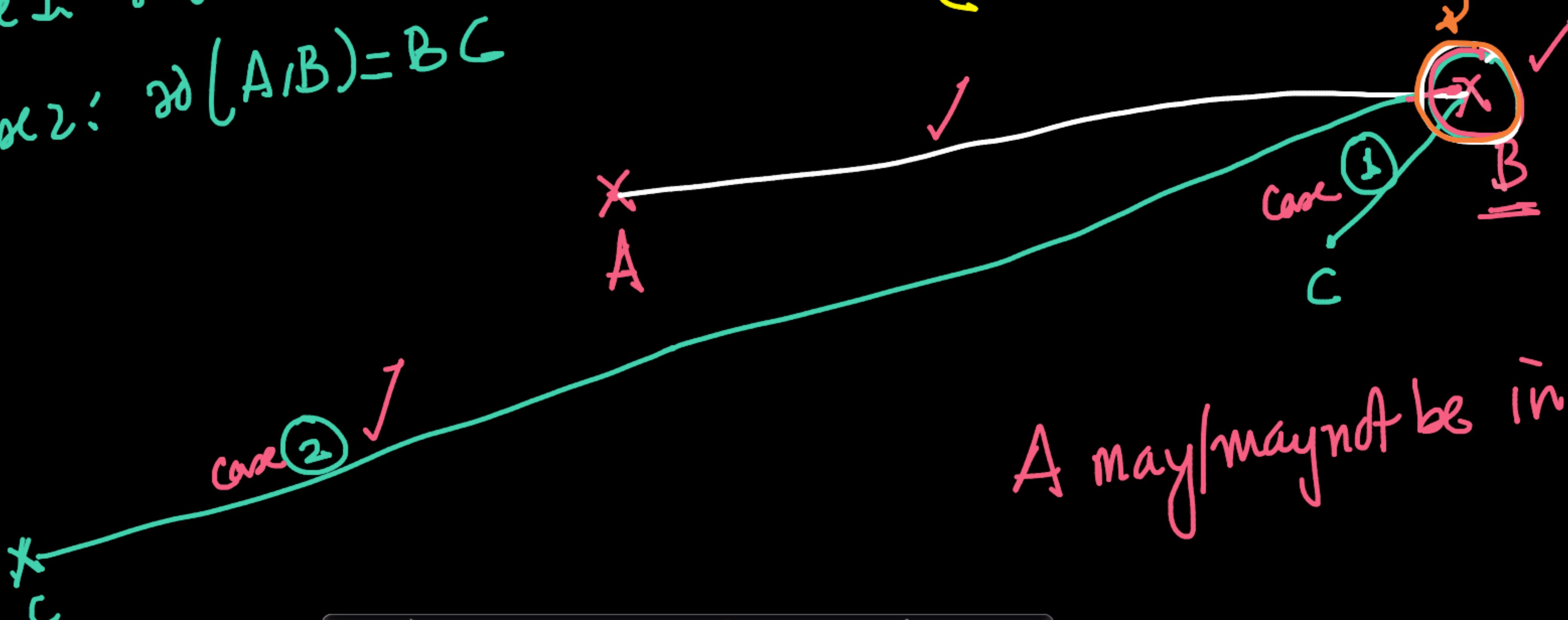
dist
②

$$\text{reachability-dist} \stackrel{\text{def}}{=} \text{dist}(A, B) = \text{rd}(A, B)$$

$$= \max \left\{ \begin{array}{l} \text{dist}(A, B), \\ \text{dist}(B) \end{array} \right\}$$

Case 1: $\text{rd}(A, B) = \text{dist}(A, B)$

Case 2: $\text{rd}(A, B) = BC$



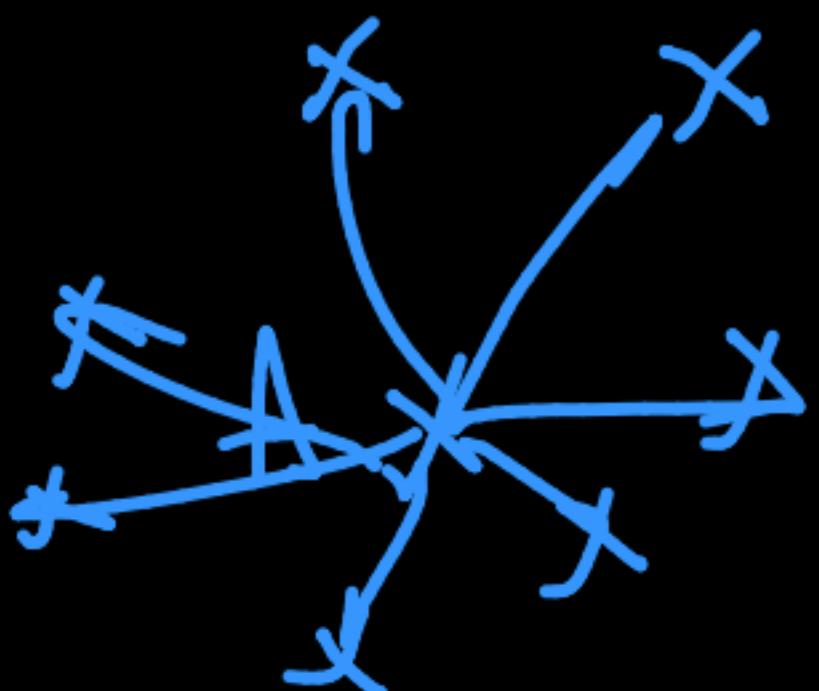
A may/may not be in $N_k(B)$



* density
local-reachability = density $\approx \text{lrdr}_k(A)$

$$\text{lrdr}_k(A) = \frac{\sum_{B \in \text{EN}_k(A)} \gamma d_k(A, B)}{|N_k(A)|}$$

: avg γd between A & k-neighbors



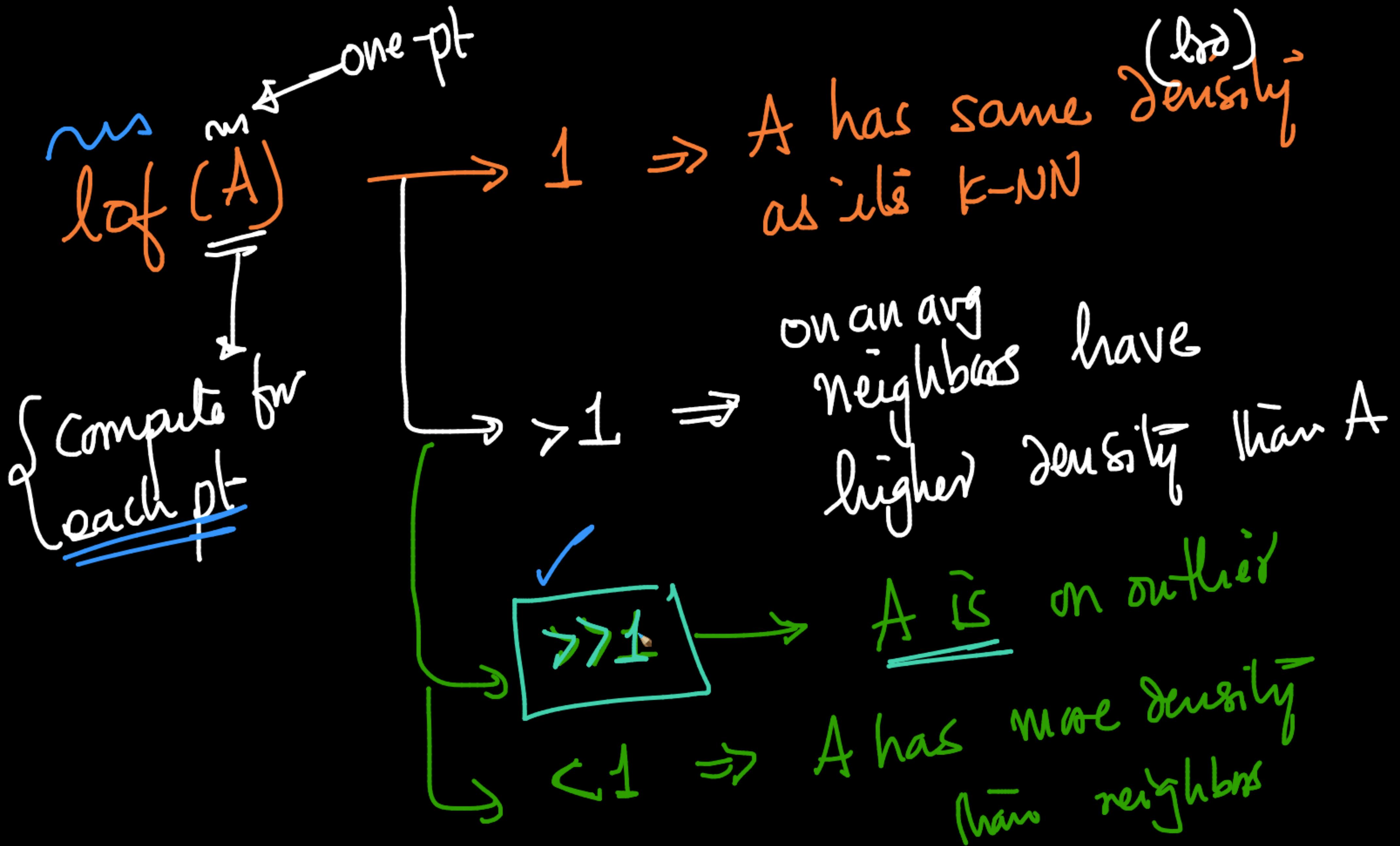
$$\text{LOF}_k(A) = \frac{\sum_{B \in N_k(A)} \text{lof}_k(B)}{|N_k(A)|} \cdot \text{lof}_k(A)$$

avg. neighbour
hood
density

avg-neighbourhood-density of A

density of A

(< 1)
(> 1 ?)
(1)

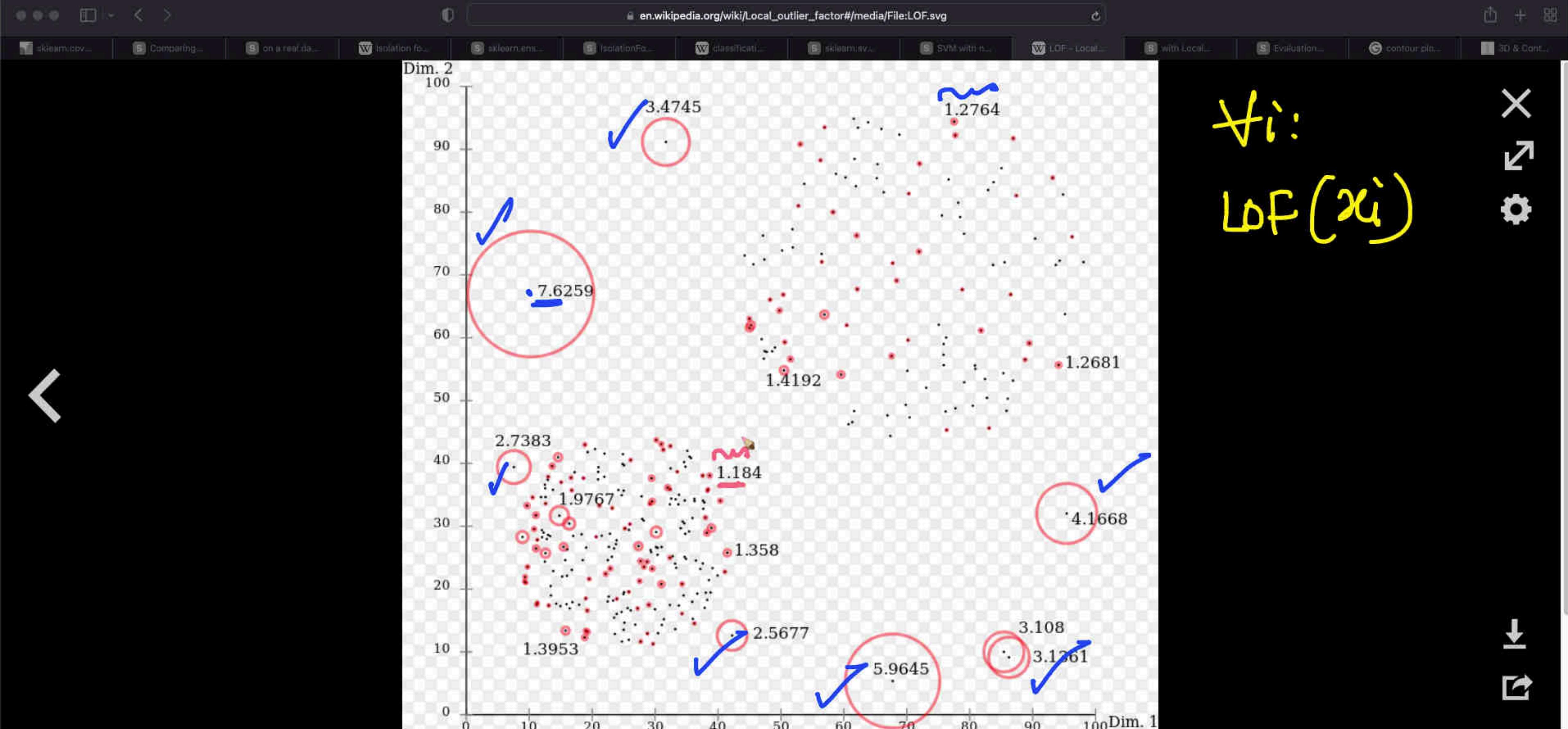


disadv

$$\rightarrow \underset{=}{{\cancel{K}}} \quad ; \quad \text{LOF}(A) >> 1 \quad \underset{=}{{\cancel{(\text{th})}}} \quad \checkmark$$

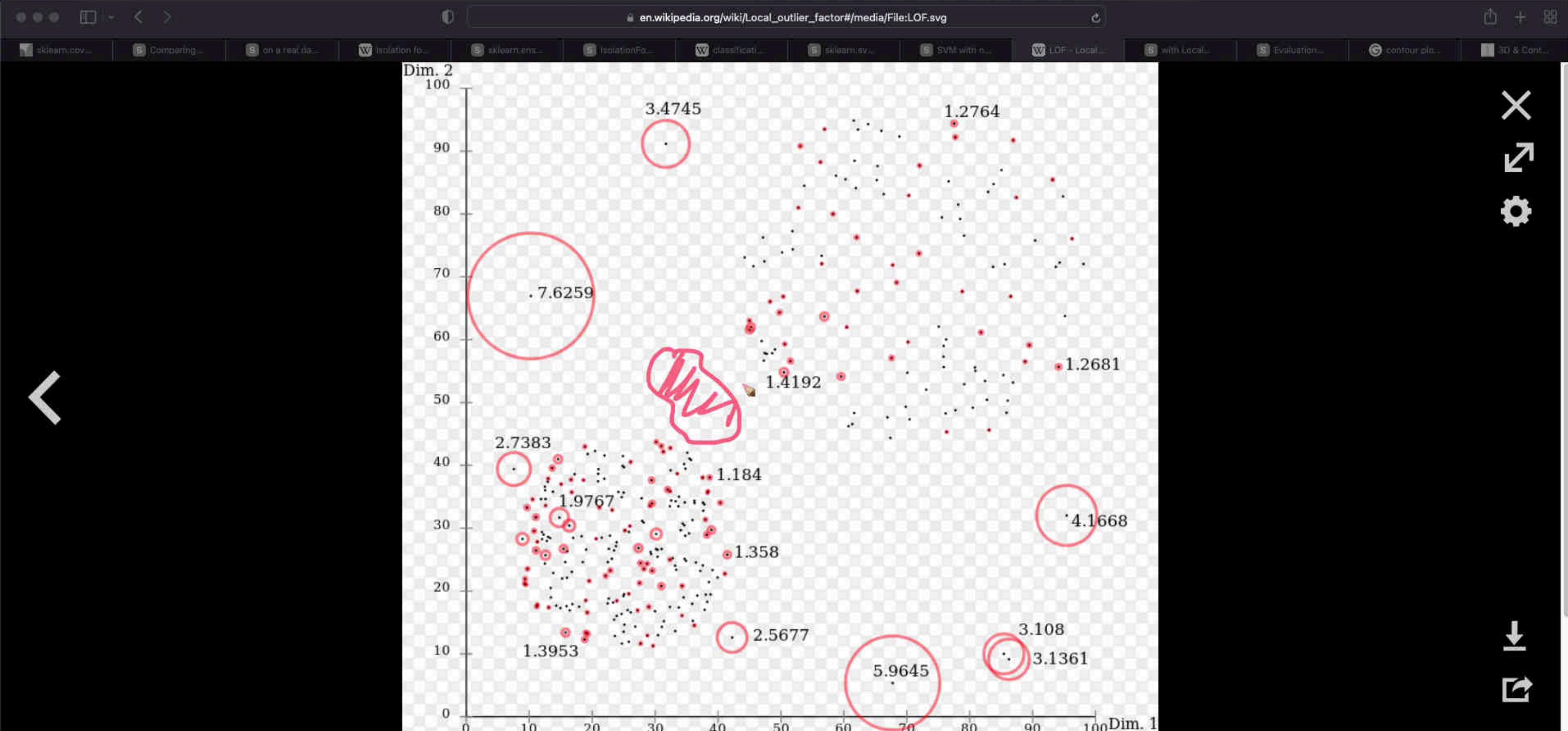
\rightarrow high - dim \rightarrow (non-euclidean)

\rightarrow time-complex (KNN fw each x_i)



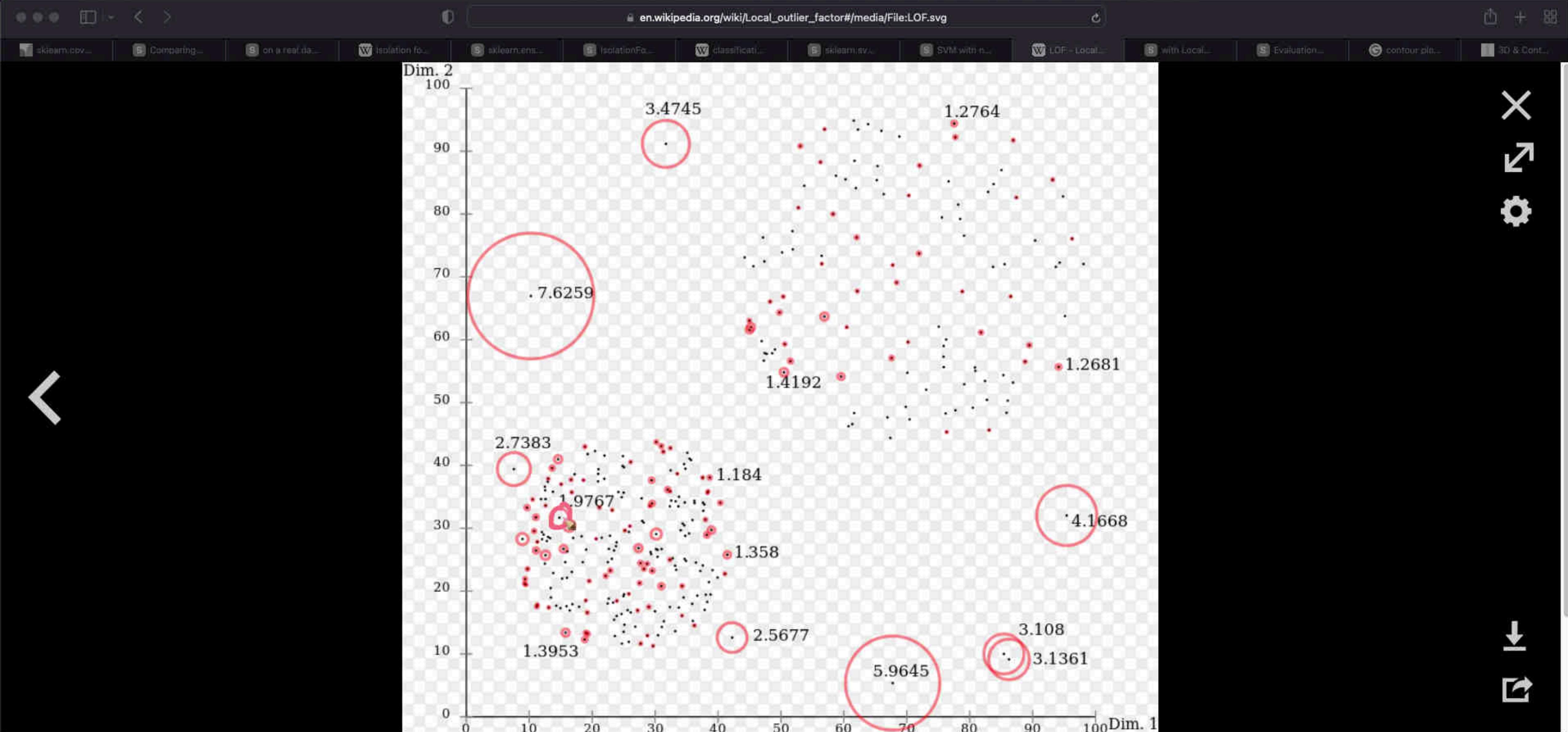
LOF scores as visualized by ELKI. While the upper right cluster has a comparable density to the outliers close to the bottom left cluster, they are detected correctly.

More details



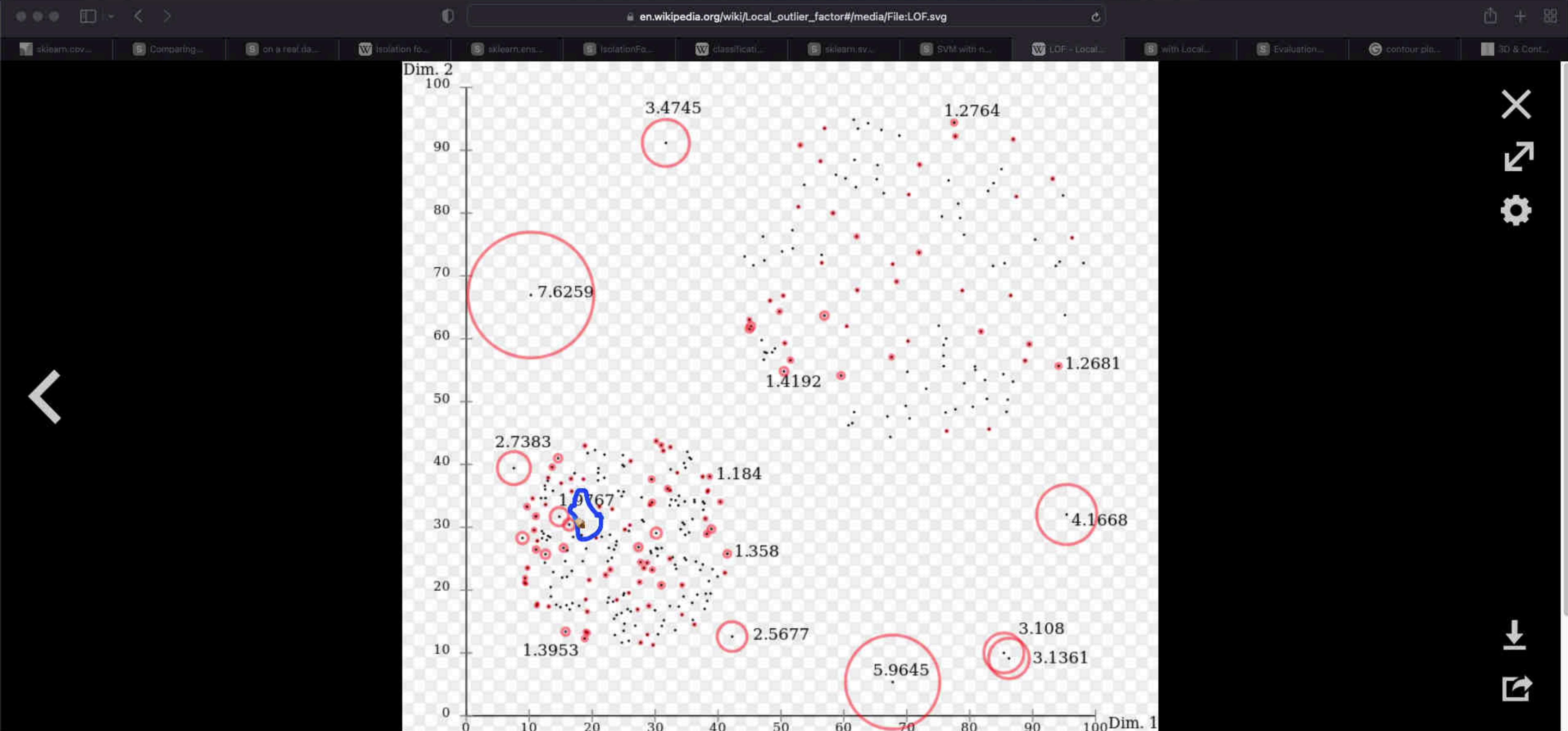
LOF scores as visualized by ELKI. While the upper right cluster has a comparable density to the outliers close to the bottom left cluster, they are detected correctly.

More details



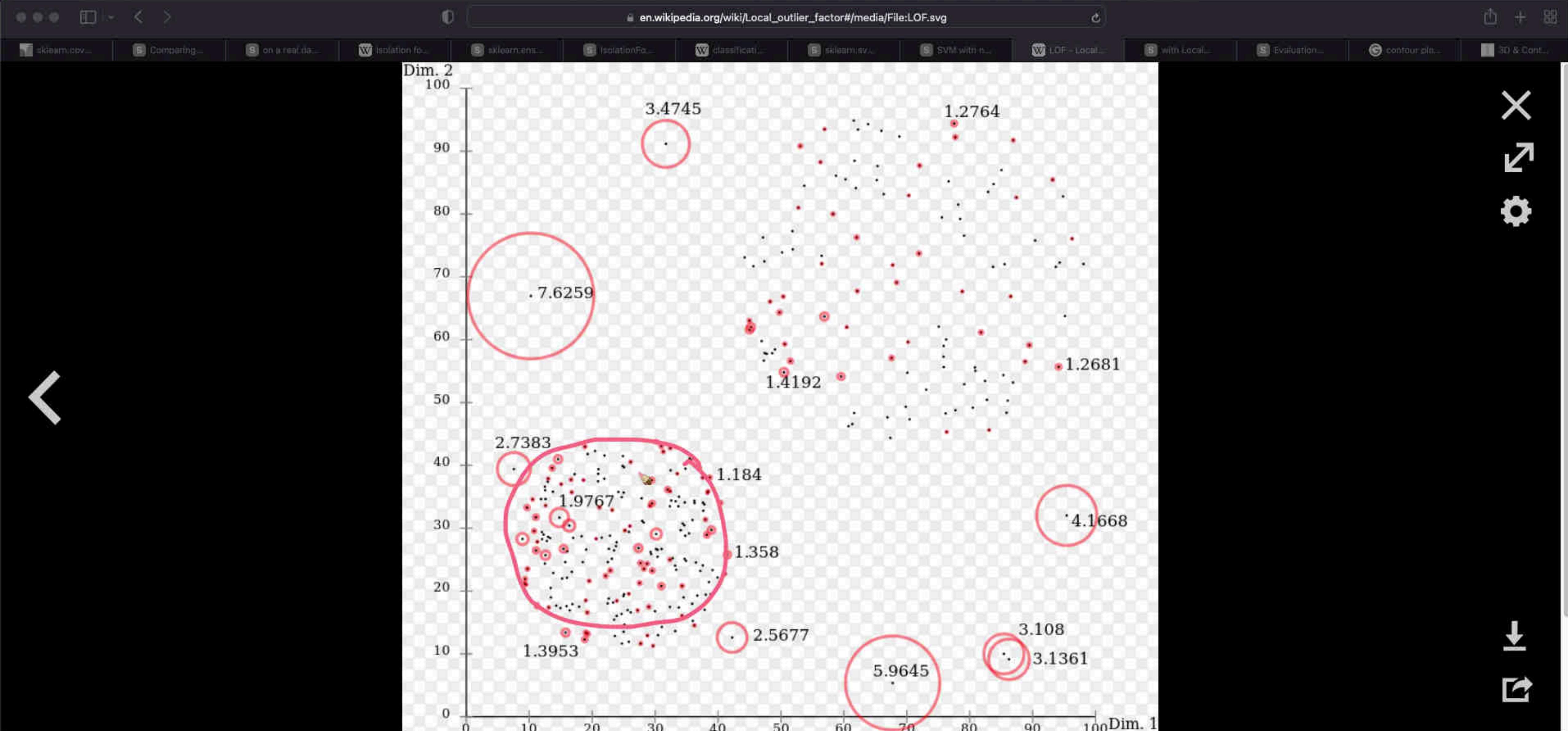
LOF scores as visualized by ELKI. While the upper right cluster has a comparable density to the outliers close to the bottom left cluster, they are detected correctly.

More details



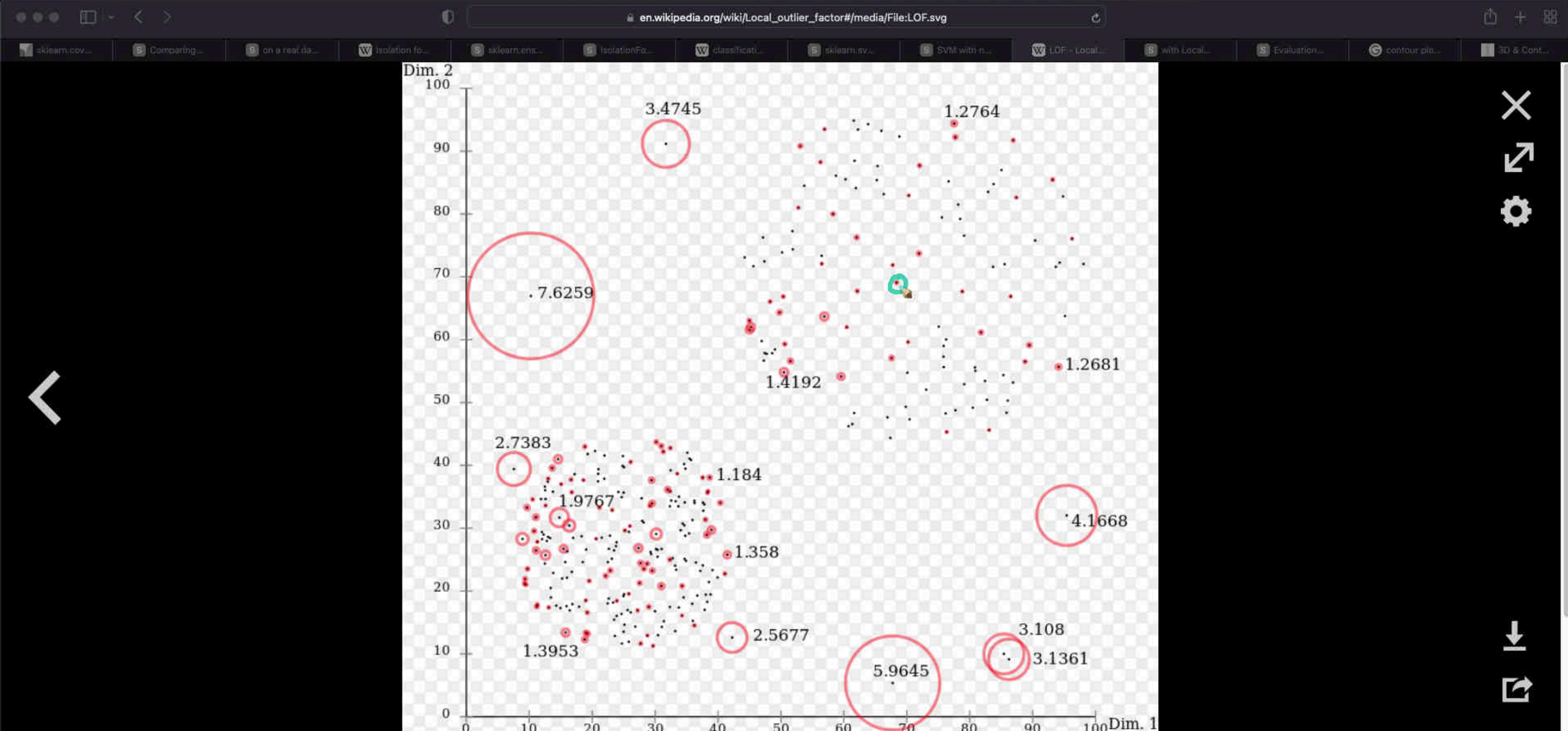
LOF scores as visualized by ELKI. While the upper right cluster has a comparable density to the outliers close to the bottom left cluster, they are detected correctly.

More details



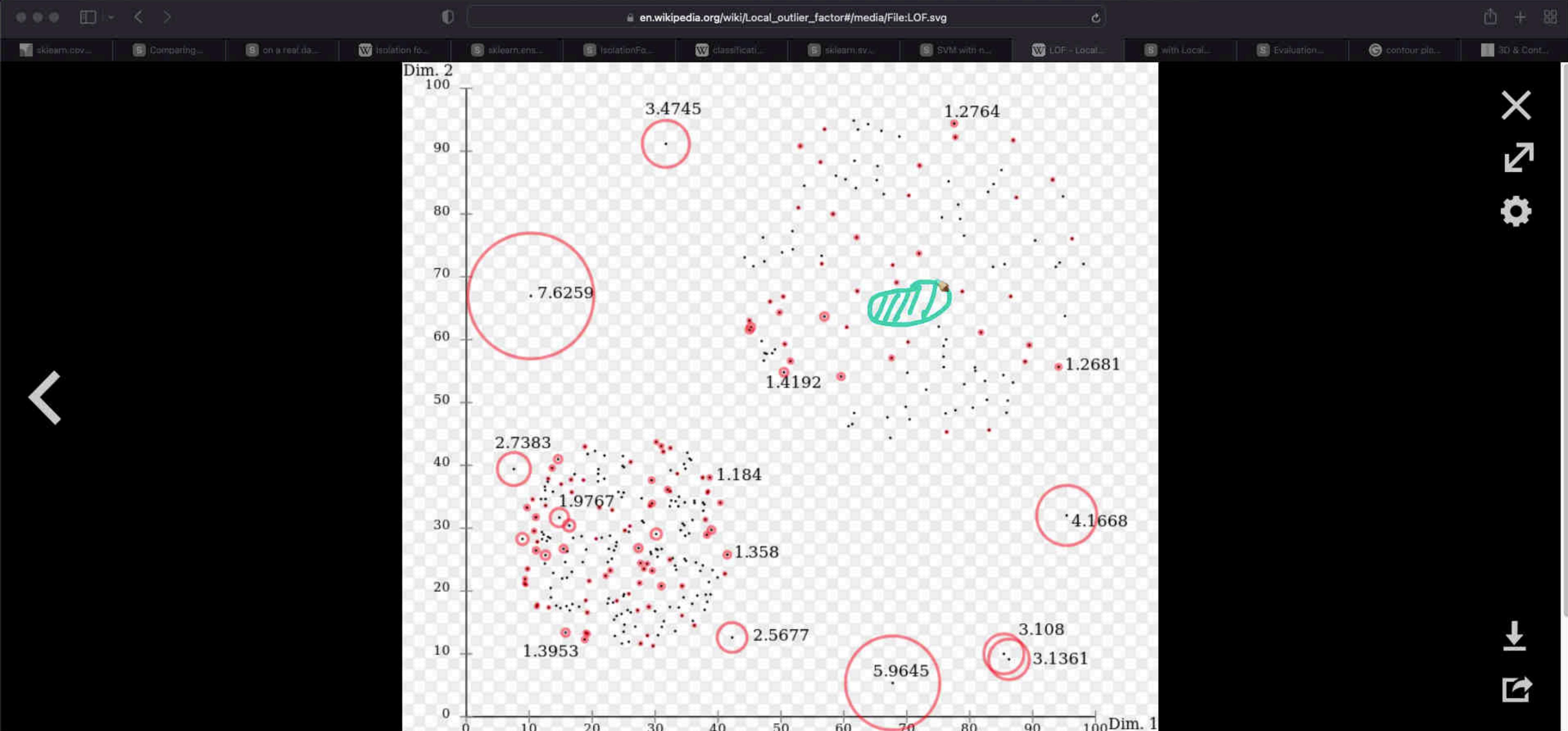
LOF scores as visualized by ELKI. While the upper right cluster has a comparable density to the outliers close to the bottom left cluster, they are detected correctly.

More details



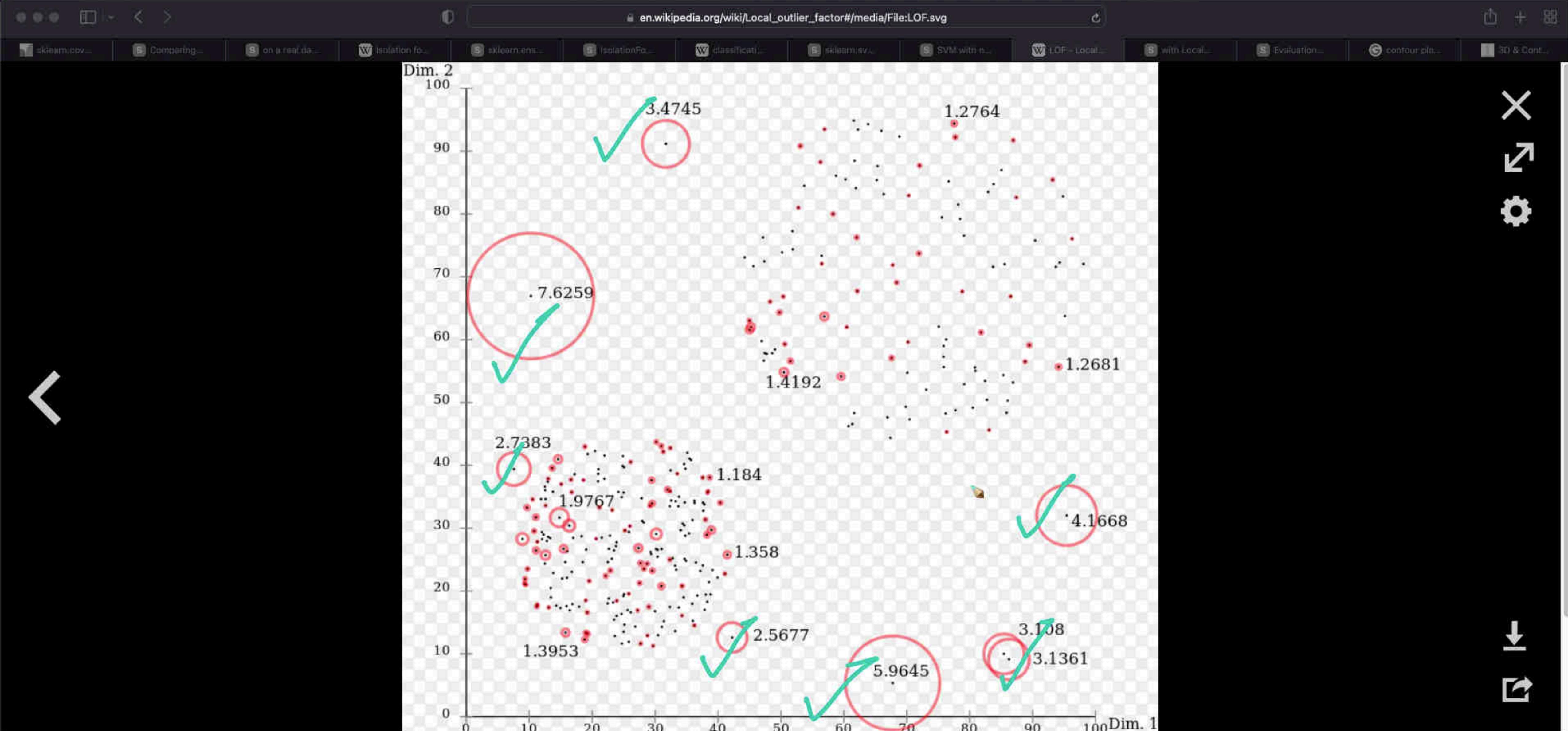
LOF scores as visualized by ELKI. While the upper right cluster has a comparable density to the outliers close to the bottom left cluster, they are detected correctly.

More details



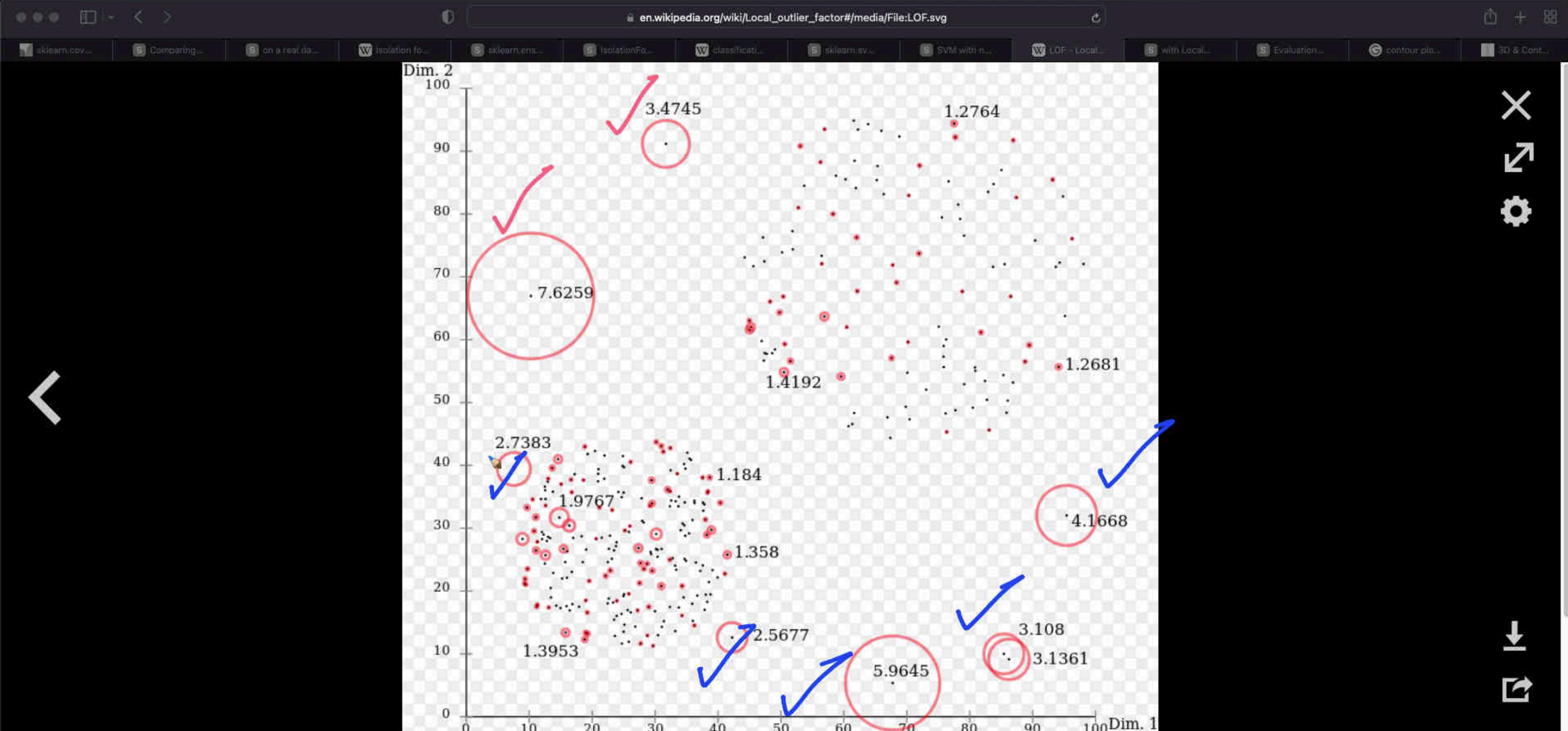
LOF scores as visualized by ELKI. While the upper right cluster has a comparable density to the outliers close to the bottom left cluster, they are detected correctly.

More details



LOF scores as visualized by ELKI. While the upper right cluster has a comparable density to the outliers close to the bottom left cluster, they are detected correctly.

 More details



LOF scores as visualized by ELKI. While the upper right cluster has a comparable density to the outliers close to the bottom left cluster, they are detected correctly.

More details

https://scikit-learn.org/stable/modules/outlier_detection.html

2.7. Novelty... Comparing... on a real da... Isolation fo... sklearns... IsolationFo... classificati... sklearns... SVM with n... LOP - Local... with local... Evaluation... contour plot... 3D & Cont...

scikit-learn Install User Guide API Examples Community More ▾ Go

Prev Up Next

scikit-learn 1.1.1

Other versions

Please [cite us](#) if you use the software.

2.7. Novelty and Outlier Detection

2.7.1. Overview of outlier detection methods

2.7.2. Novelty Detection

2.7.3. Outlier Detection

2.7.4. Novelty detection with Local Outlier Factor

2.7. Novelty and Outlier Detection

Many applications require being able to decide whether a new observation belongs to the same distribution as existing observations (it is an *inlier*), or should be considered as different (it is an *outlier*). Often, this ability is used to clean real data sets. Two important distinctions must be made:

outlier detection:	The training data contains outliers which are defined as observations that are far from the others. Outlier detection estimators thus try to fit the regions where the training data is the most concentrated, ignoring the deviant observations.
novelty detection:	The training data is not polluted by outliers and we are interested in detecting whether a new observation is an outlier. In this context an outlier is also called a novelty.

Outlier detection and novelty detection are both used for anomaly detection, where one is interested in detecting abnormal or unusual observations. Outlier detection is then also known as unsupervised anomaly detection and novelty detection as semi-supervised anomaly detection. In the context of outlier detection, the outliers/anomalies cannot form a dense cluster as available estimators assume that the outliers/anomalies are located in low density regions. On the contrary, in the context of novelty detection, novelties/anomalies can form a dense cluster as long as they are in a low density region of the training data, considered as normal in this context.

Toggle Menu