



# Anomaly detection using R

MULTIVARIATE GAUSSIAN DISTRIBUTION  
BASED ANOMALY DETECT

유 충 현

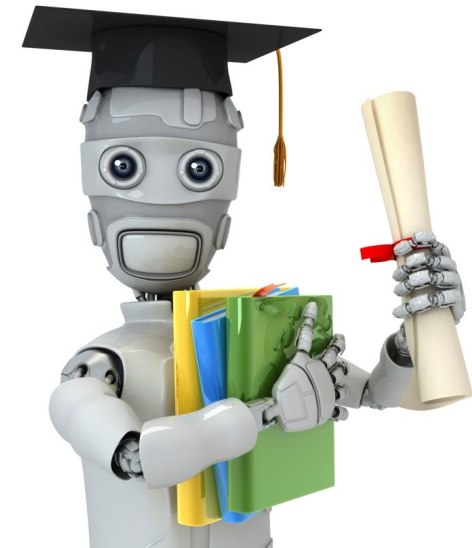
KRUG

# AGENDA

- Problem motivation
- Gaussian distribution
- Algorithm
- Anomaly detection using the multivariate Gaussian distribution
- Developing and evaluating an anomaly detection system
- `anomalyDensityEstimation` class

# 일러두기

- 슬라이드 및 이론 설명
  - Coursera 강좌 중 Machine Learning
- 구현 및 예제
  - S4를 적용하여 구현함
  - Detection function / Visualization function



Machine Learning

# Anomaly detection example

Aircraft engine features:

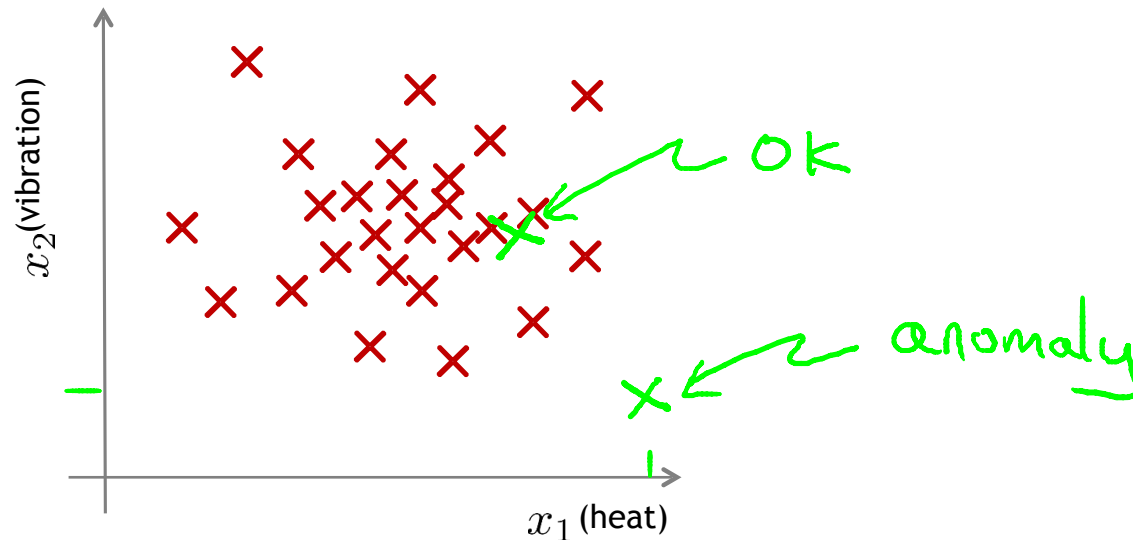
$x_1$  = heat generated

$x_2$  = vibration intensity

...

Dataset:  $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

New engine:  $x_{test}$

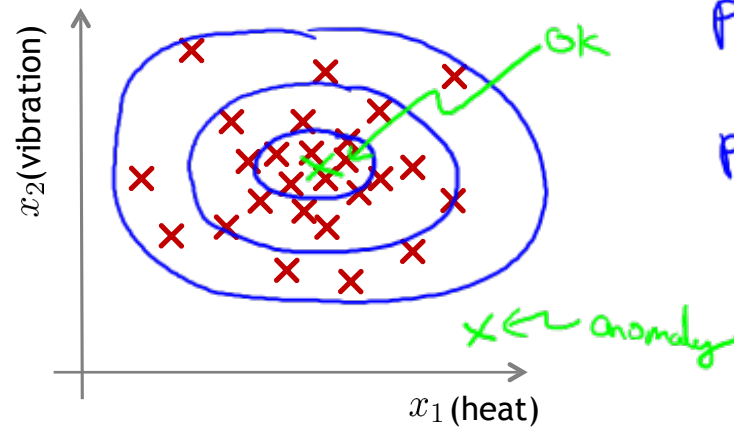


# Density estimation

→ Dataset:  $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

→ Is  $x_{test}$  anomalous?

Model  $p(x)$ .



$p(x_{test}) < \varepsilon \rightarrow$  flag anomaly

$p(x_{test}) \geq \varepsilon \rightarrow$  OK

# Anomaly detection example

Fraud detection:

$x^{(i)}$  = features of user  $i$ 's activities

Model  $p(x)$  from data.

Identify unusual users by checking which have  $p(x) < \varepsilon$

Manufacturing

Monitoring computers in a data center.

$x^{(i)}$  = features of machine  $i$

$x_1$  = memory use,  $x_2$  = number of disk accesses/sec,

$x_3$  = CPU load,  $x_4$  = CPU load/network traffic.

...

# Gaussian(Normal) distribution

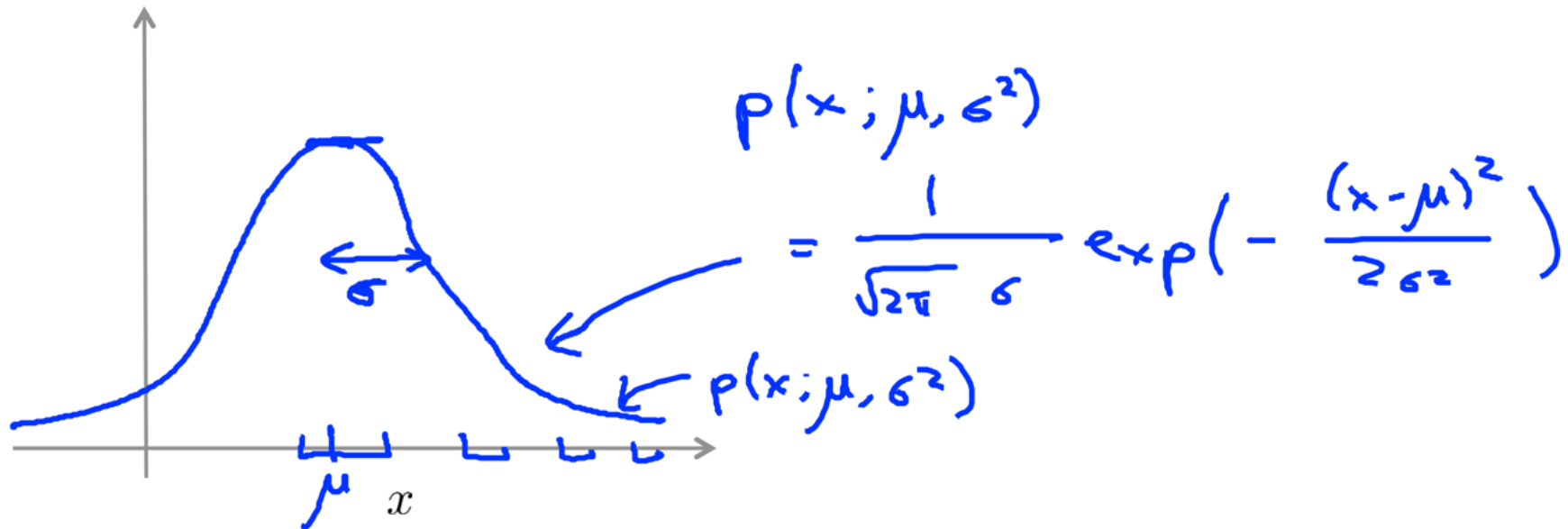
## Gaussian (Normal) distribution

Say  $x \in \mathbb{R}$ . If  $x$  is a distributed Gaussian with mean  $\mu$ , variance  $\sigma^2$ .

$$x \sim \mathcal{N}(\mu, \sigma^2)$$

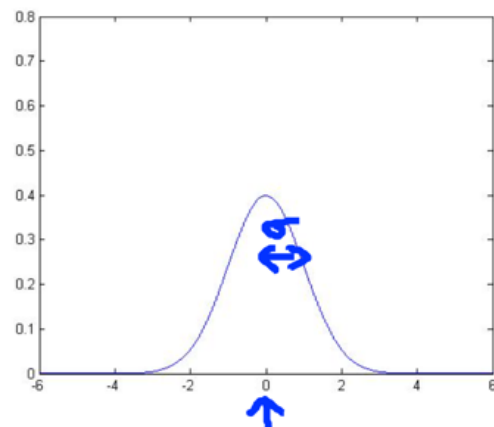
↑ "distributed as"

$\sigma$  standard deviation



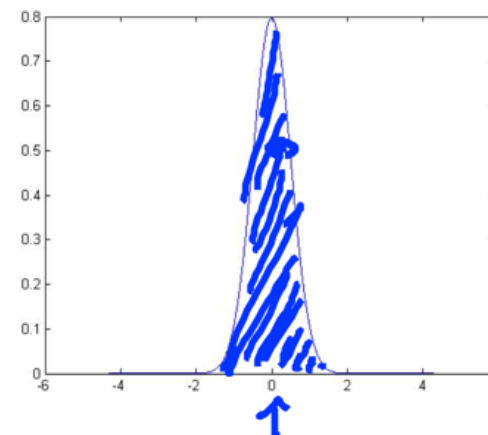
# Gaussian distribution example

→  $\mu = 0, \sigma = 1$

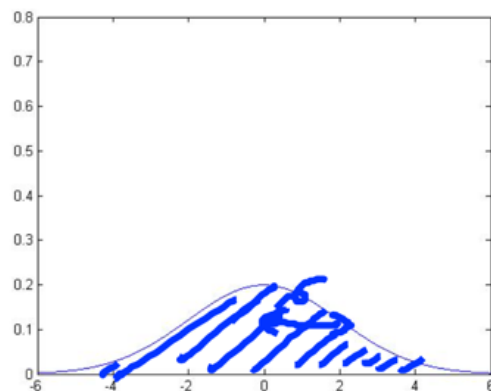


I

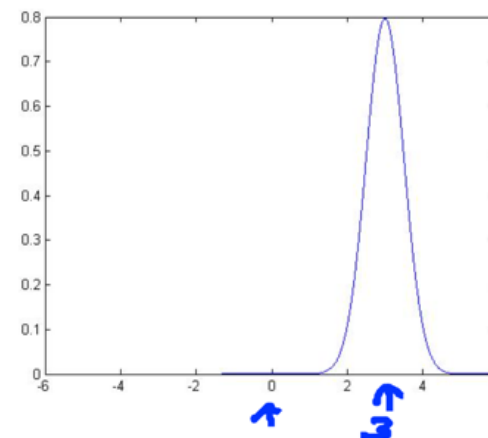
→  $\mu = 0, \sigma = \underline{0.5}$



→  $\mu = 0, \sigma = 2$



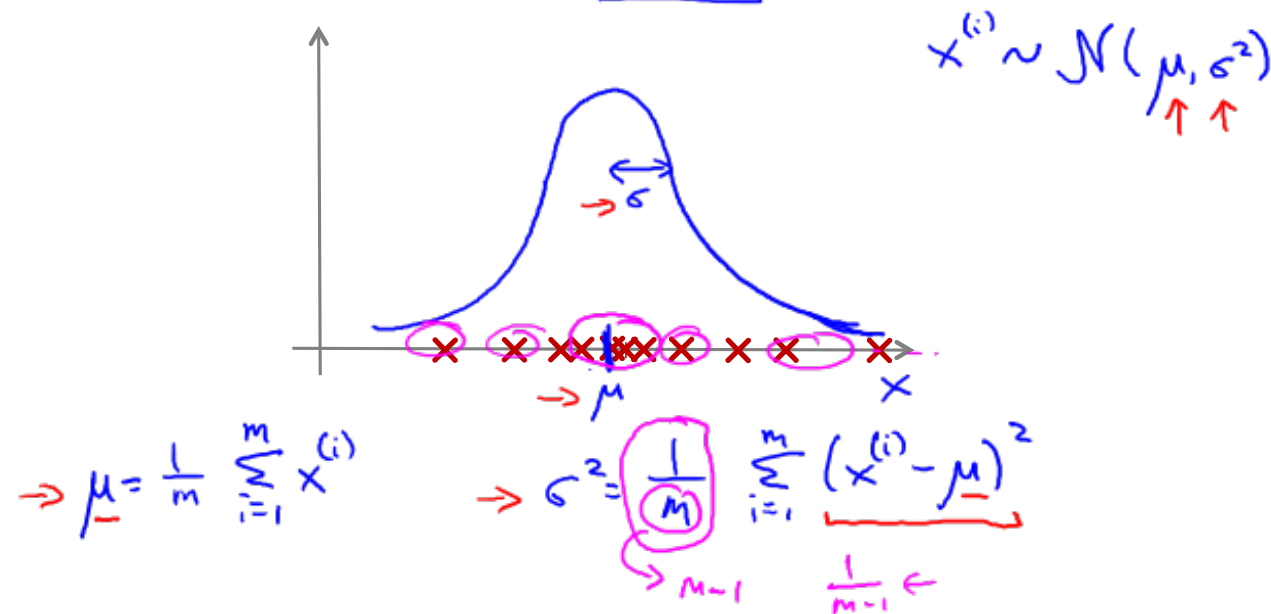
→  $\mu = 3, \sigma = 0.5$





# Parameter estimation

→ Dataset:  $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$   $x^{(i)} \in \mathbb{R}$



# Density estimation

→

→ Training set:  $\{x^{(1)}, \dots, x^{(m)}\}$   
 Each example is  $x \in \mathbb{R}^n$

$$x_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$$

$$x_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$$

$$x_3 \sim \mathcal{N}(\mu_3, \sigma_3^2)$$

→  $p(x)$

$$= p(x_1; \mu_1, \sigma_1^2) p(x_2; \mu_2, \sigma_2^2) p(x_3; \mu_3, \sigma_3^2) \dots p(x_n; \mu_n, \sigma_n^2) \leftarrow$$

$$= \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2)$$

$$\sum_{i=1}^n i = 1+2+3+\dots+n$$

$$\prod_{i=1}^n i = 1 \times 2 \times 3 \times \dots \times n$$

# Anomaly detection algorithm

1. Choose features  $x_i$  that you think might be indicative of anomalous examples.
2. Fit parameters  $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$$
$$\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$

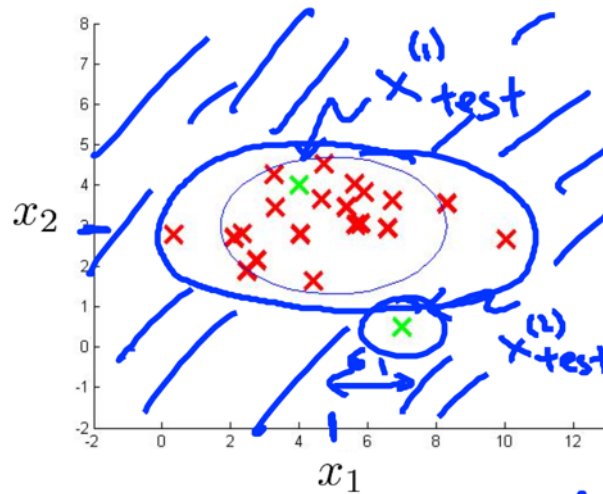
3. Given new example  $x$ , compute  $p(x)$  :

$$p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi}\sigma_j} \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right)$$

Anomaly if  $p(x) < \varepsilon$

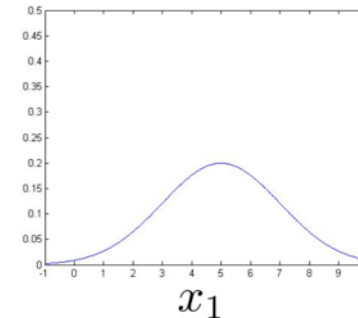
# Anomaly detection example

## Anomaly detection example

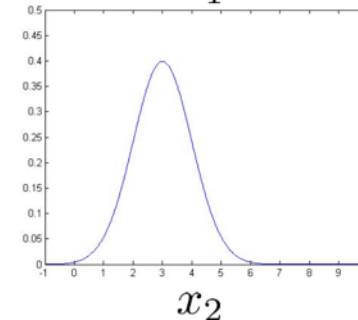


$$\begin{aligned} \mu_1 &= 5, \sigma_1 = 2 \\ \mu_2 &= 3, \sigma_2 = 1 \end{aligned}$$

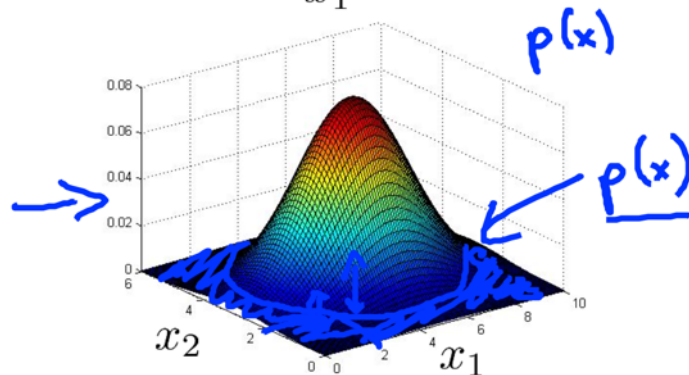
$$\rightarrow p(x) = p(x_1; \mu_1, \sigma_1^2) \times p(x_2; \mu_2, \sigma_2^2)$$



$$p(x_1; \mu_1, \sigma_1^2)$$



$$p(x_2; \mu_2, \sigma_2^2)$$



$$\varepsilon = 0.02$$

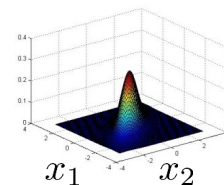
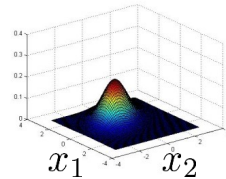
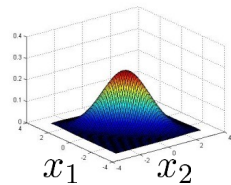
$$p(x_{test}^{(1)}) = 0.0426 \geq \varepsilon$$

$$p(x_{test}^{(2)}) = 0.0021 < \varepsilon$$

# Multivariate Gaussian distribution

Parameters  $\mu, \Sigma$

$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left( -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right)$$



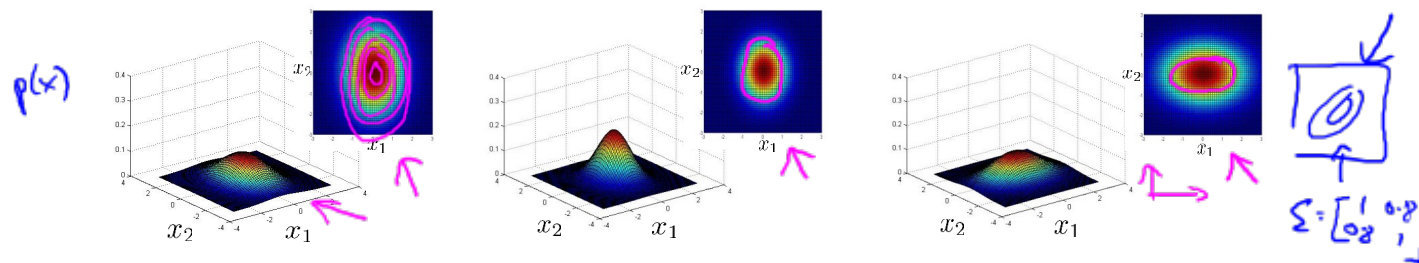
Parameter fitting:

Given training set  $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)} \quad \Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)(x^{(i)} - \mu)^T$$

# Relationship to original model

Original model  $p(x) = p(x_1; \mu_1, \sigma_1^2) \times p(x_2; \mu_2, \sigma_2^2) \times \dots \times p(x_n; \mu_n, \sigma_n^2)$



Corresponds to multivariate Gaussian

$$\rightarrow p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left( -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right)$$

where

A handwritten representation of the covariance matrix  $\Sigma$  as a 2x2 matrix with elements  $\sigma_{11}$ ,  $\sigma_{12}$ ,  $\sigma_{21}$ , and  $\sigma_{22}$  on the diagonal and off-diagonal positions respectively. The matrix is enclosed in a box with arrows pointing to it.

# The importance of real-number evaluation

When developing a learning algorithm (choosing features, etc.), making decisions is much easier if we have a way of evaluating our learning algorithm.

Assume we have some labeled data, of anomalous and non-anomalous examples. ( $y = 0$  if normal,  $y = 1$  if anomalous).

Training set:  $x^{(1)}, x^{(2)}, \dots, x^{(m)}$  (assume normal examples/ not anomalous)

Cross validation set:  $(x_{cv}^{(1)}, y_{cv}^{(1)}), \dots, (x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$

Test set:  $(x_{test}^{(1)}, y_{test}^{(1)}), \dots, (x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$

# Aircraft engines example

## Aircraft engines motivating example

10000 good (normal) engines

20 flawed engines (anomalous)

Training set: 6000 good engines

CV: 2000 good engines (  $y = 0$  ), 10 anomalous (  $y = 1$  )

Test: 2000 good engines (  $y = 0$  ), 10 anomalous (  $y = 1$  )

Alternative:

Training set: 6000 good engines

CV: 4000 good engines (  $y = 0$  ), 10 anomalous (  $y = 1$  )

Test: 4000 good engines (  $y = 0$  ), 10 anomalous (  $y = 1$  )



# Algorithm evaluation

- Fit model  $p(x)$  on training set  $\{x^{(1)}, \dots, x^{(m)}\}$
- On a cross validation/test example  $x$ , predict

$(x_{\text{test}}^{(i)}, y_{\text{test}}^{(i)})$   
↑

$$y = \begin{cases} 1 & \text{if } p(x) < \epsilon \text{ (anomaly)} \\ 0 & \text{if } p(x) \geq \epsilon \text{ (normal)} \end{cases}$$

$y = 0$

Possible evaluation metrics:

- - True positive, false positive, false negative, true negative
- - Precision/Recall
- -  $F_1$ -score ←

CV

Test set

Can also use cross validation set to choose parameter  $\epsilon$  ←

# Algorithm evaluation

$$prec = \frac{tp}{tp + fp} \quad (4)$$

$$rec = \frac{tp}{tp + fn}, \quad (5)$$

where

- $tp$  is the number of true positives: the ground truth label says it's an anomaly and our algorithm correctly classified it as an anomaly.
- $fp$  is the number of false positives: the ground truth label says it's not an anomaly, but our algorithm incorrectly classified it as an anomaly.
- $fn$  is the number of false negatives: the ground truth label says it's an anomaly, but our algorithm incorrectly classified it as not being anomalous.

$$F_1 = \frac{2 \cdot prec \cdot rec}{prec + rec},$$

# class

- **anomalyDensityEastimation**
  - S4 class
  - detection & visualization
- **slots**
  - X : data. matrix object
  - param : density parameter. list object (mu, sigma2)
  - p : probability
  - y : training set's anomaly flag. vector (0=normality, 1=anomaly)
  - threshold : epsilon. list object
  - anomaly : detected observations

# methods

- estimateGaussian
  - 확률밀도함수의 모수추정
- multivariateGaussian
  - multivariate Gaussian 확률 계산
- select Threshold
  - best epsilon, best F1 계산
- findAnomaly
  - anomaly detection
- plot
  - 2D multivariate Gaussian plotting

# parameter estimation

```
> anomaly.trian <- new("anomalyDensityEstimation", X = X)
> anomaly.trian@param <- estimateGaussian(anomaly.trian)
> anomaly.trian@param
$mu
throughput    latency
    14.11223    14.99771

$sigma2
      [,1] [,2]
[1,] 1.832631 0.000000
[2,] 0.000000 1.709745
```

## select threshold

```
> anomaly.cv@threshold <- selectThreshold(anomaly.cv)
> anomaly.cv@threshold
$eps
[1] 8.990853e-05

$F1
[1] 0.875
```

# anomaly detect

```
> anomaly.trian@p <- multivariateGaussian(anomaly.trian)
> anomaly.trian@anomaly <- findAnomaly(anomaly.trian, anomaly.cv@threshold)
> anomaly.trian@anomaly
```

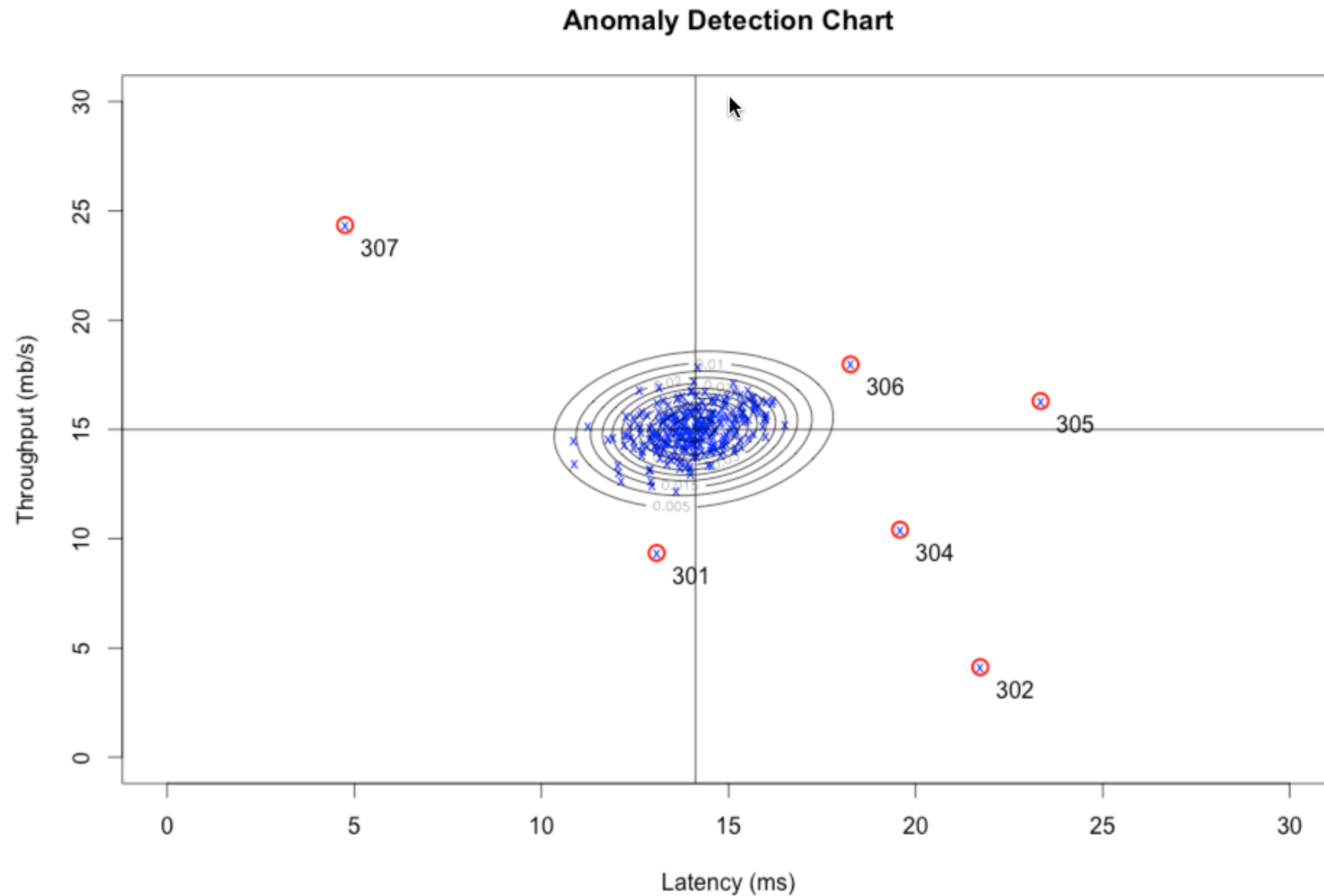
\$idx

```
[1] 301 302 304 305 306 307
```

\$value

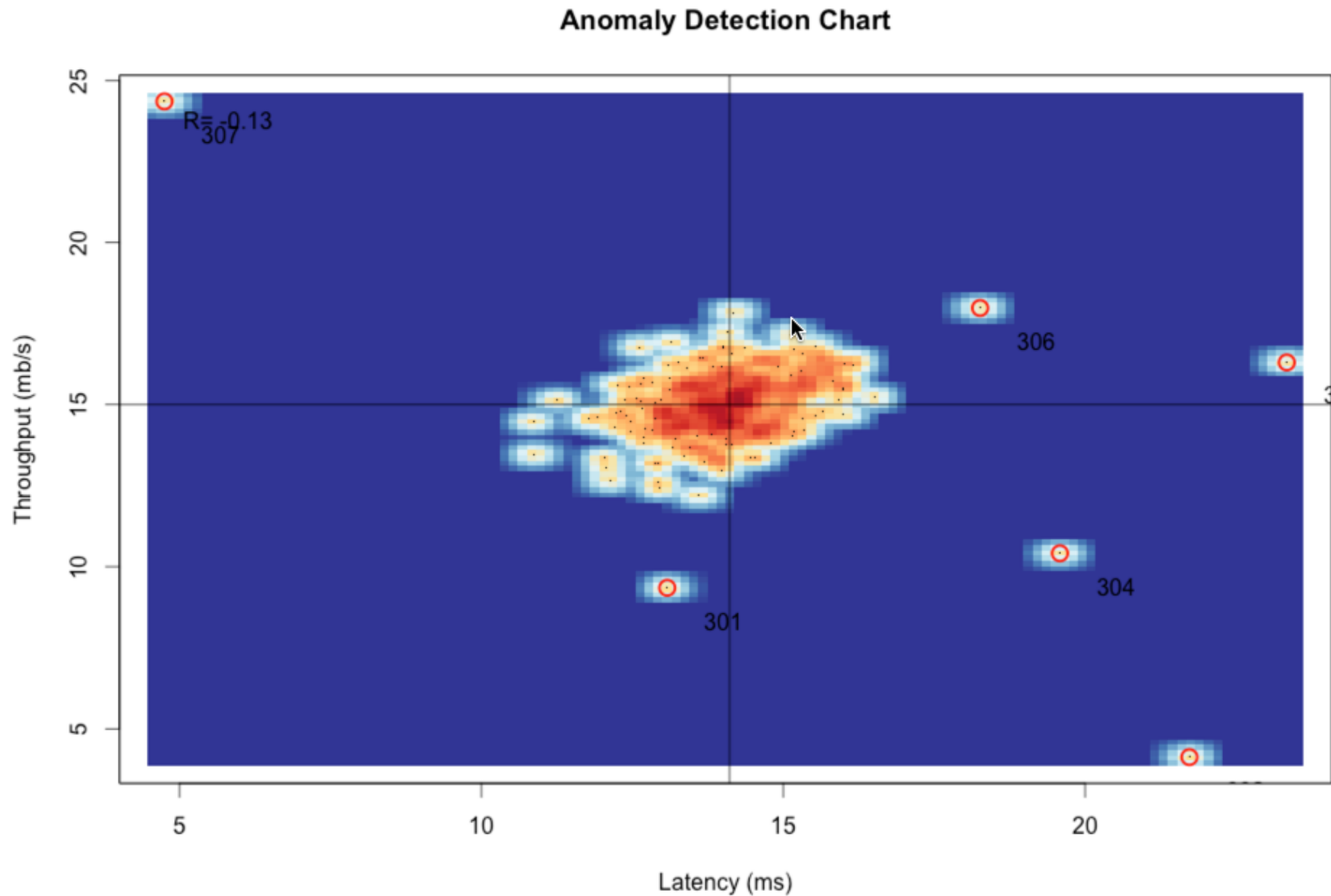
	throughput	latency
[1,]	13.079310	9.347878
[2,]	21.727134	4.126232
[3,]	19.582573	10.411619
[4,]	23.339868	16.298874
[5,]	18.261188	17.978309
[6,]	4.752613	24.350407

# plotting visualization

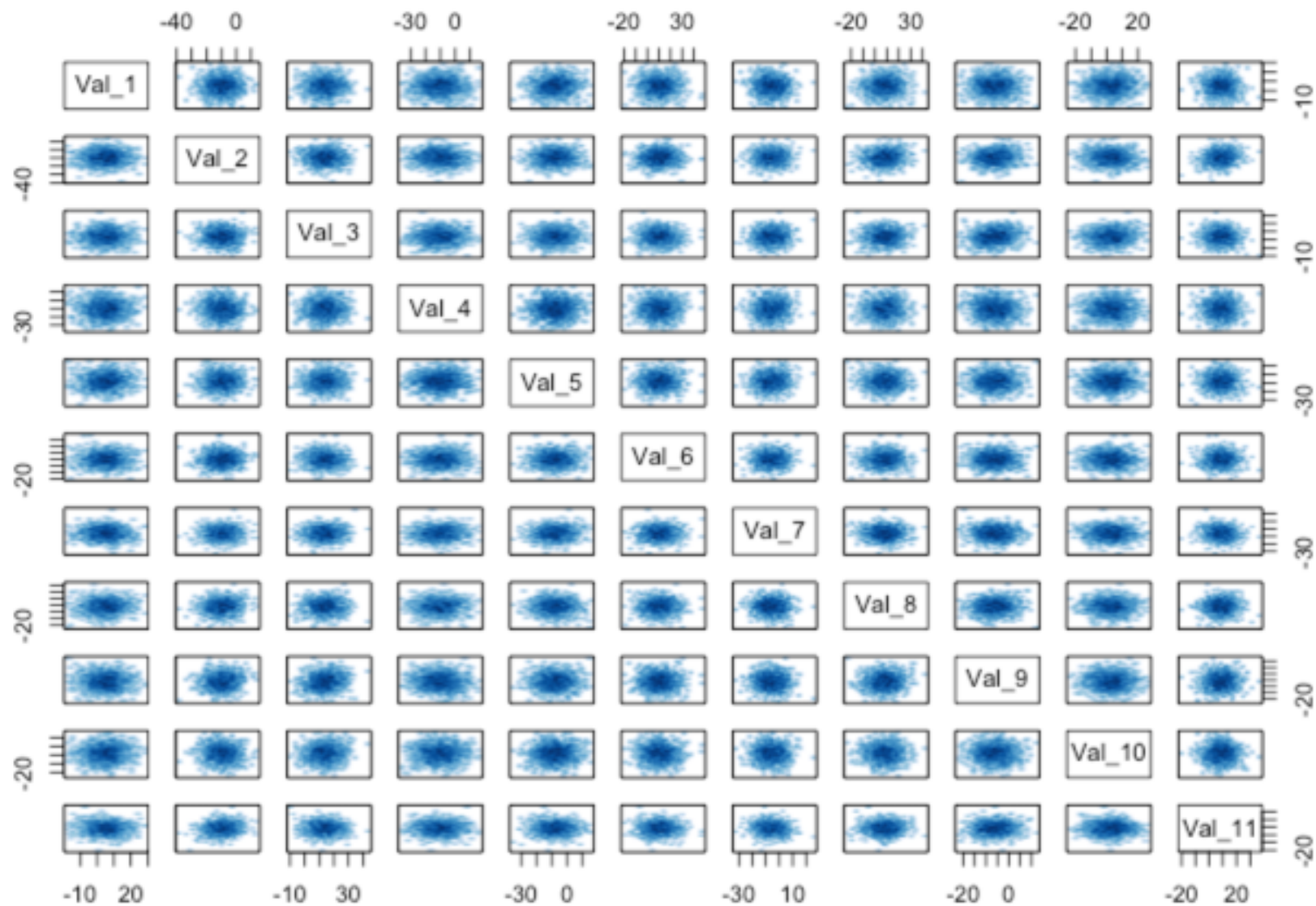




# plotting visualization



# plotting visualization



**Demos**

감사합니다