

Machine Learning

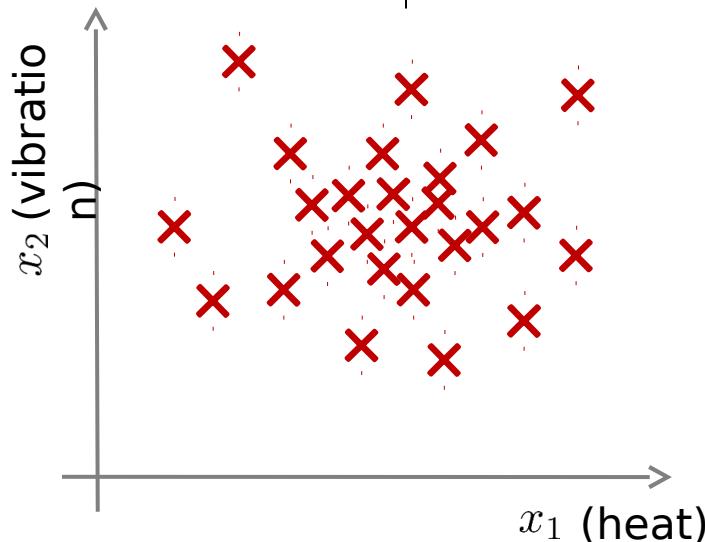
# Anomaly detection Problem motivation

# Anomaly detection example

Aircraft engine  
features:  
 $x_1$  = heat generated  
= vibration  
intensity

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

New engine  $x_{test}$

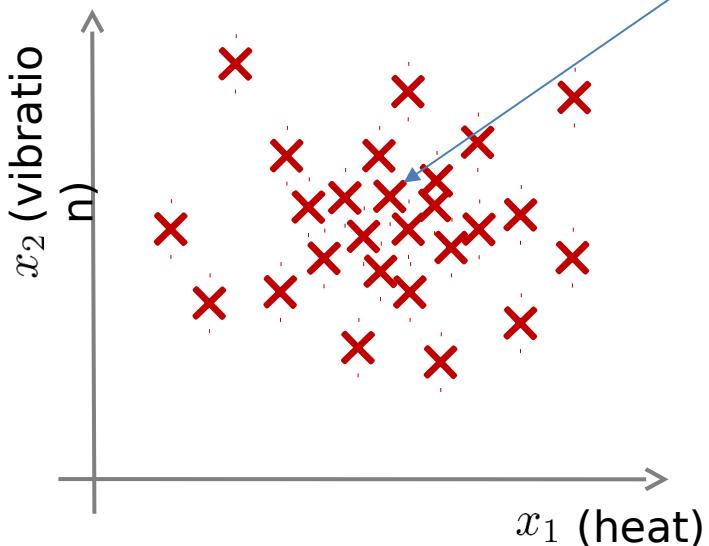


# Density estimation

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

Is  $x_{test}$  anomalous?

(Probability is high at the centre)



# 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  
Manufacturing

Monitoring computers in a data center.

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

$x_1$  = memory use,  $x_2$  = number of disk  
accesses/sec,  $x_3$  = CPU load,  $x_4$  = CPU load/network traffic.

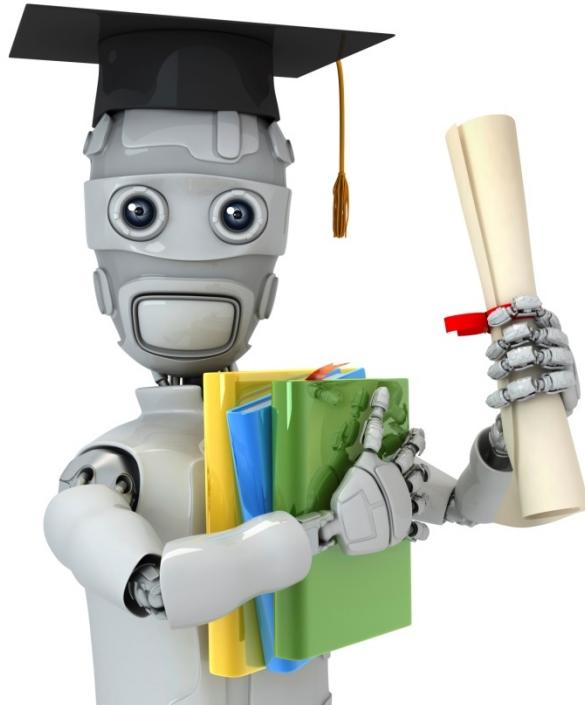
...

Your anomaly detection system flags  $x$  as anomalous whenever  $p(x) \leq \epsilon$ . Suppose your system is flagging too many things as anomalous that are not actually so (similar to supervised learning, these mistakes are called false positives). What should you do?

---

- Try increasing  $\epsilon$ .
- Try decreasing  $\epsilon$ .

**Correct Response**



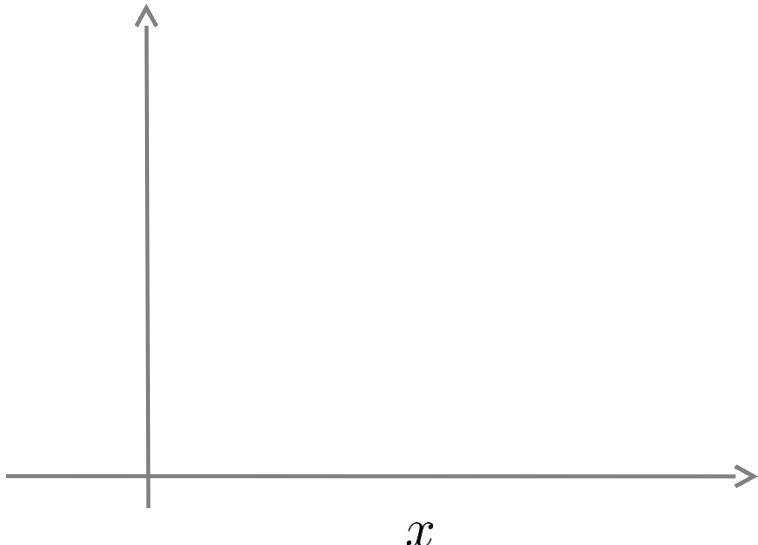
Machine Learning

# Anomaly detection

## Gaussian distribution

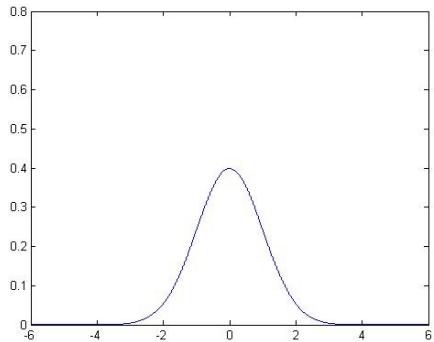
# Gaussian (Normal) distribution

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

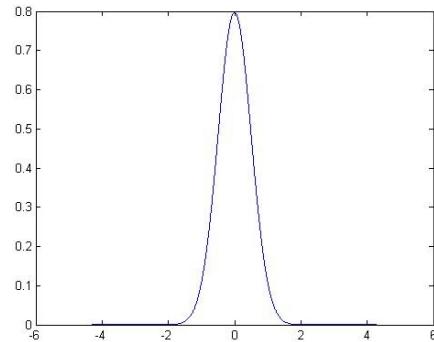


# Gaussian distribution example

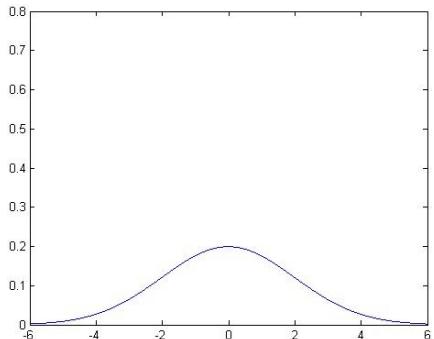
$$\mu = 0, \sigma = 1$$



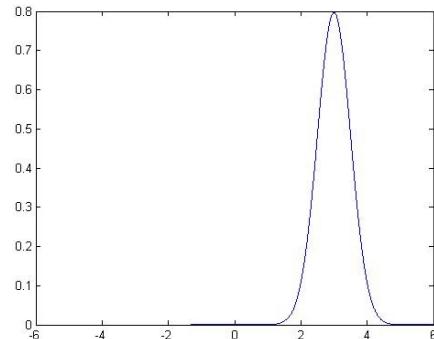
$$\mu = 0, \sigma = 0.5$$



$$\mu = 0, \sigma = 2$$



$$\mu = 3, \sigma = 0.5$$

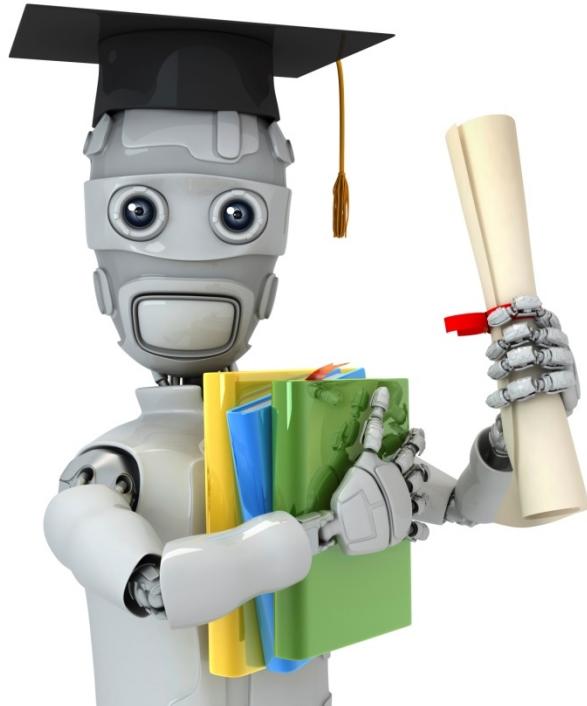


# Parameter estimation

(Given the dataset we want to estimate  
the value of mu and sigma)

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





Machine Learning

# Anomaly detection Algorithm

# Density estimation

Training set  $\{x^{(1)}, \dots, x^{(m)}\}$

Each example is  $\mathbb{R}^n$

Given a training set  $\{x^{(1)}, \dots, x^{(m)}\}$ , how would you estimate each  $\mu_j$  and  $\sigma_j^2$  (Note  $\mu_j \in \mathbb{R}, \sigma_j^2 \in \mathbb{R}$ .)

---

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

**Correct Response**

# Anomaly detection algorithm

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

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

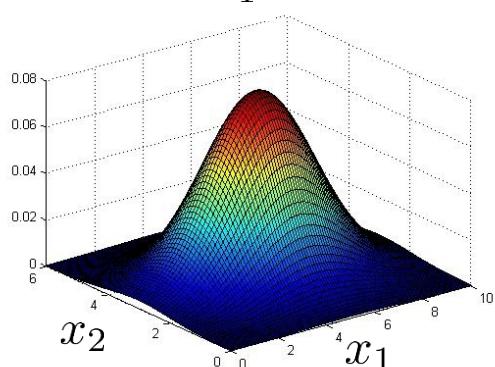
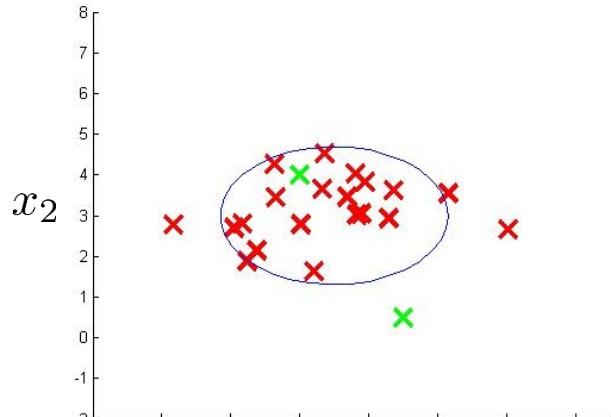
$$\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$

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

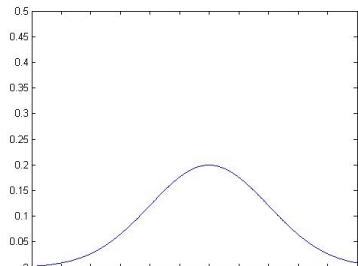
$$p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi}\sigma_j} \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right)$$

Anomaly if  $p(x) < \varepsilon$

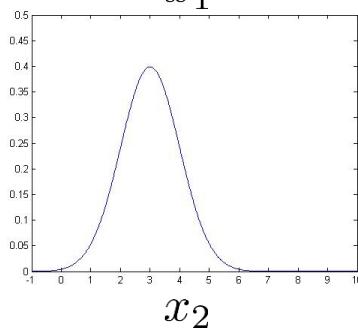
# Anomaly detection exam



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



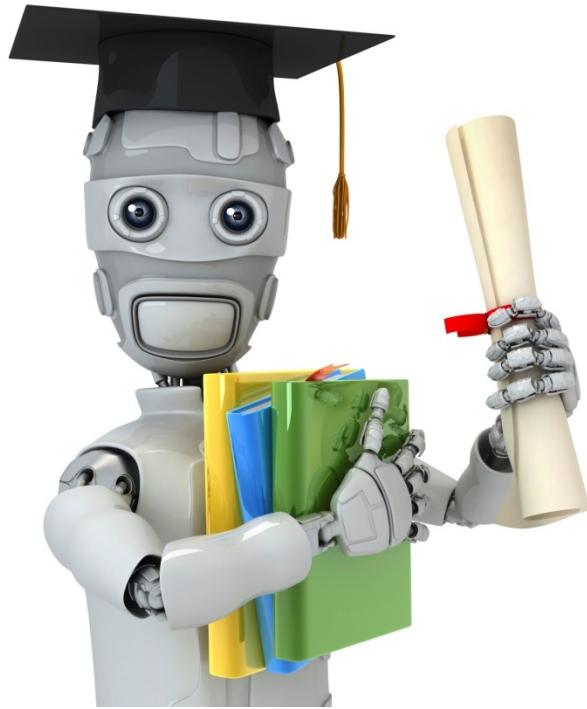
$$p(x_1; \mu_1, \sigma_1^2)$$



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

$$\varepsilon = 0.02$$

$$\begin{aligned}p(x_{test}^{(1)}) &= 0.0426 \\ p(x_{test}^{(2)}) &= 0.0021\end{aligned}$$



Machine Learning

# Anomaly detection

---

Developing and  
evaluating an  
anomaly detection  
system

# 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. ( $y = 1$  if normal,

if anomalous).

Training set:  $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$  (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 motivating example

10000 good (normal) engines

20 flawed engines (anomalous)

Training set: 6000 good engines

CV: 2000 good engines( ), 10 anomalous ( )

Test: 2000 good engines( ), 10 anomalous ( )

Alternative:

Training set: 6000 good engines

CV: 4000 good engines( ), 10 anomalous ( )

Test: 4000 good engines( ), 10 anomalous ( )

# Algorithm evaluation

Fit model  $p(x)$  on training set  $\{x^{(1)}, \dots, x^{(m)}\}$

On a cross validation/test example  $x$ , predict

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

Possible evaluation metrics:

- True positive, false positive, false negative, true negative

- Precision/Recall

- F<sub>1</sub>-score

Can also use cross validation set to choose  $\varepsilon$  parameter

Suppose you have fit a model  $p(x)$ . When evaluating on the cross validation set or test set, your algorithm predicts:

$$y = \begin{cases} 1 & \text{if } p(x) \leq \epsilon \\ 0 & \text{if } p(x) > \epsilon \end{cases}$$

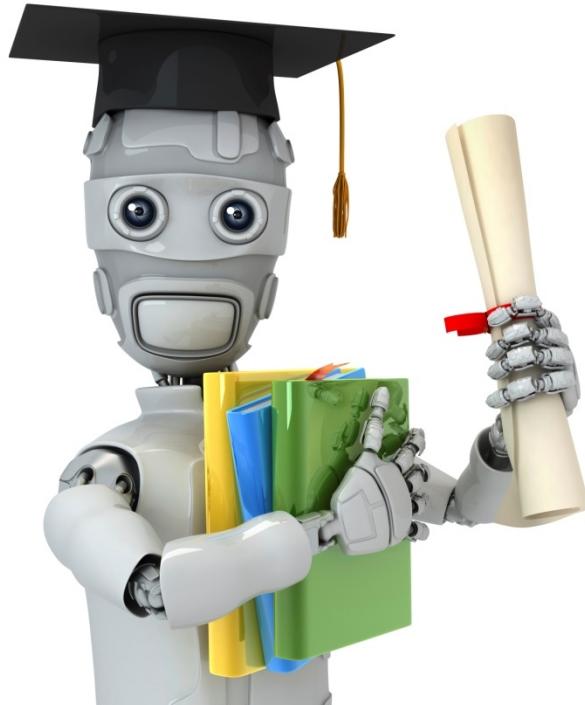
Is classification accuracy a good way to measure the algorithm's performance?

---

- Yes, because we have labels in the cross validation / test sets.
- No, because we do not have labels in the cross validation / test sets.
- No, because of skewed classes (so an algorithm that always predicts  $y = 0$  will have high accuracy).

**Correct Response**

- No for the cross validation set; yes for the test set.



Machine Learning

# Anomaly detection detection vs. supervised learning

## **Anomaly detection vs. Supervised learning**

Very small number of positive examples ( ). (0-20 is common).

Large number of negative ( ) examples.

Many different “types” of anomalies. Hard for any algorithm to learn from positive examples what the anomalies look like; future anomalies may look nothing like any of the anomalous examples we’ve seen so far.

Large number of positive and negative examples.

Enough positive examples for algorithm to get a sense of what positive examples are like, future positive examples likely to be similar to ones in training set.

## Anomaly detection

## vs. Supervised learning

- Fraud detection
- Manufacturing (e.g. aircraft engines)
- Monitoring machines in a data center

⋮

- Email spam classification
- Weather prediction (sunny/rainy/etc).
- Cancer classification

⋮

Which of the following problems would you approach with an anomaly detection algorithm (rather than a supervised learning algorithm)? Check all that apply.

- You run a power utility (supplying electricity to customers) and want to monitor your electric plants to see if any one of them might be behaving strangely.

**Correct Response**

- You run a power utility and want to predict tomorrow's expected demand for electricity (so that you can plan to ramp up an appropriate amount of generation capacity).

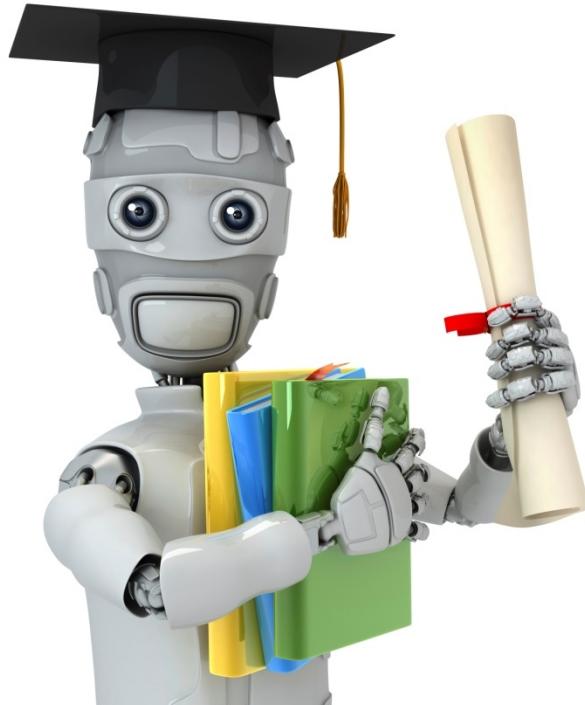
**Correct Response**

- A computer vision / security application, where you examine video images to see if anyone in your company's parking lot is acting in an unusual way.

**Correct Response**

- A computer vision application, where you examine an image of a person entering your retail store to determine if the person is male or female.

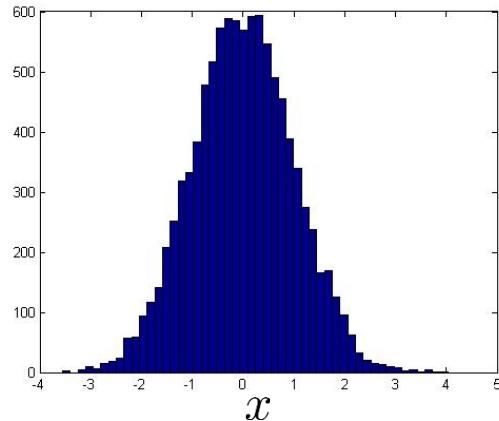
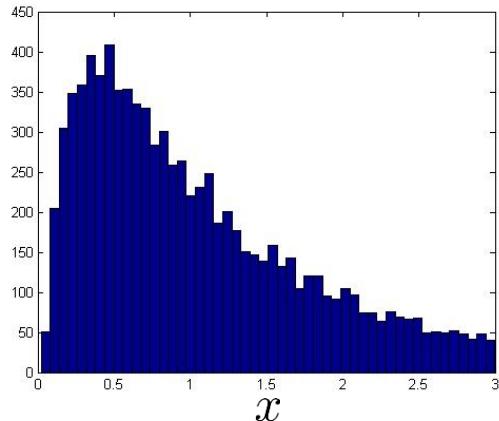
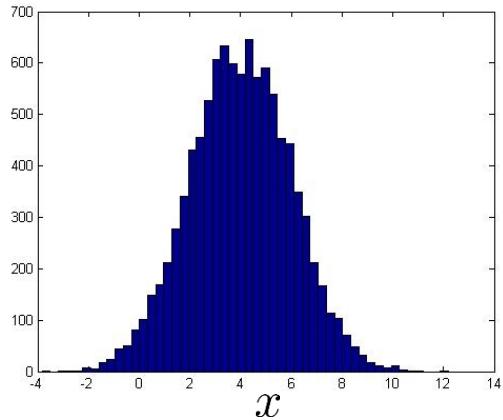
**Correct Response**



Machine Learning

# Anomaly detection what features to use

# Non-gaussian features



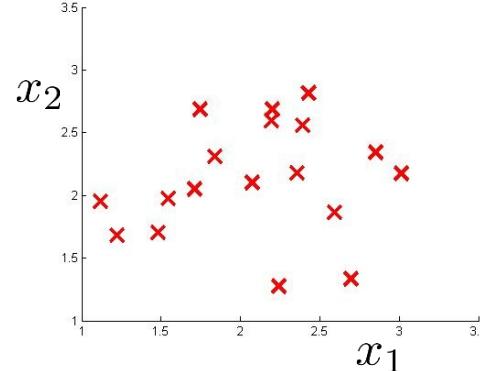
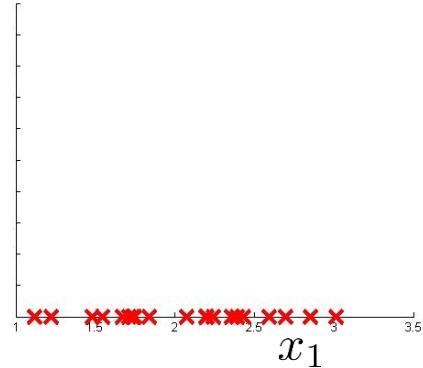
# Error analysis for anomaly detection

Want  $p(x)$  large for normal examples  $x$ .

small for anomalous examples  $x$ .

Most common problem:

$p(x)$  comparable (say, both large) for normal and anomalous examples



# **Monitoring computers in a data center**

Choose features that might take on unusually large or small values in the event of an anomaly.

- =<sub>2</sub> memory use of computer
- =<sub>3</sub> number of disk accesses/sec
- =<sub>4</sub> CPU load
- = network traffic

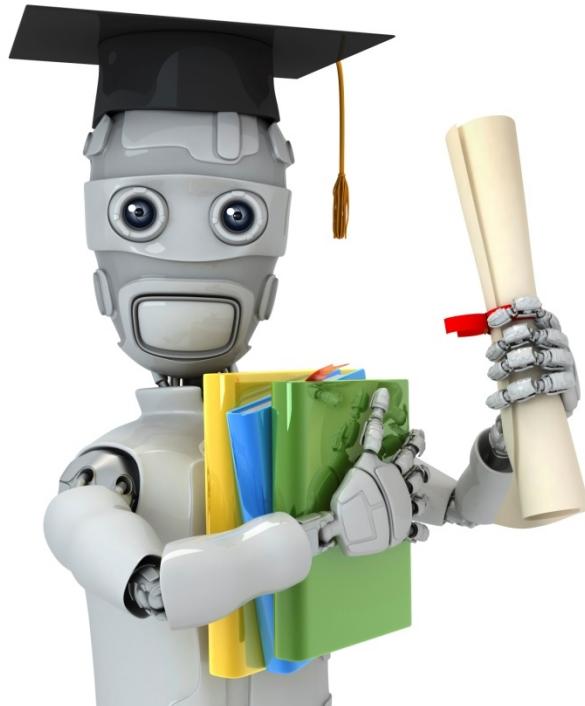
Suppose your anomaly detection algorithm is performing poorly and outputs a large value of  $p(x)$  for many normal examples and for many anomalous examples in your cross validation dataset. Which of the following changes to your algorithm is most likely to help?

---

- Try using fewer features.
- Try coming up with more features to distinguish between the normal and the anomalous examples.

**Correct Response**

- Get a larger training set (of normal examples) with which to fit  $p(x)$ .
- Try changing  $\epsilon$ .



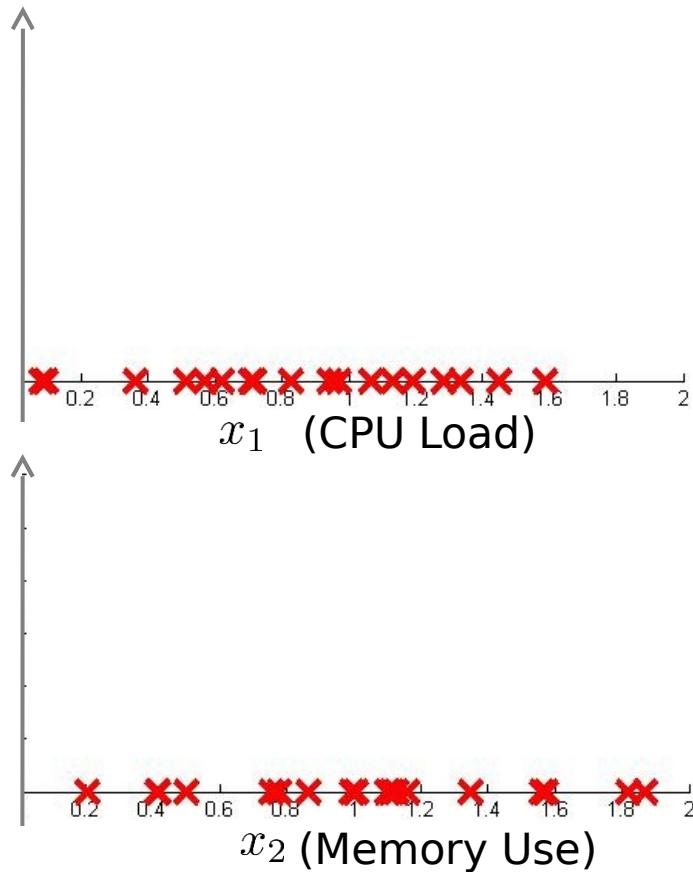
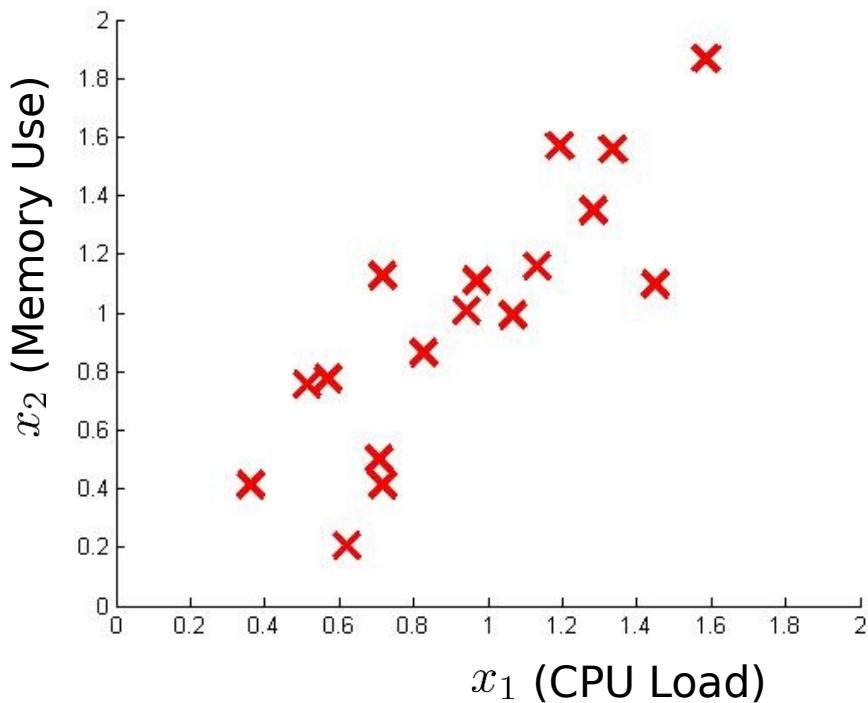
Machine Learning

# Anomaly detection

---

## Multivariate Gaussian distribution

# Motivating example: Monitoring machines in a data center



# Multivariate Gaussian (Normal)

~~distribution~~ Don't model  $p(x_1)$ ,  $p(x_2)$ , .. etc. separately.

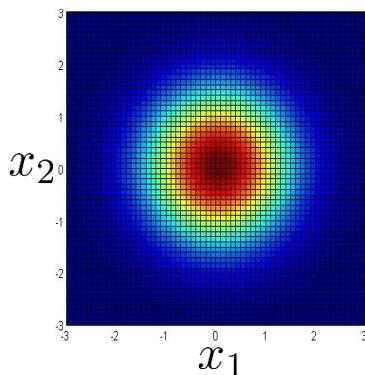
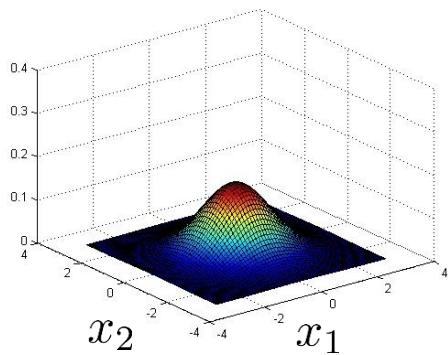
Model  $p(x)$  all in one go.

Parameters  $\mu \in \mathbb{R}^n$ ,  $\Sigma \in \mathbb{R}^{n \times n}$  Covariance matrix)

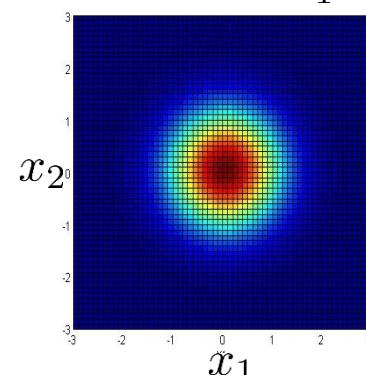
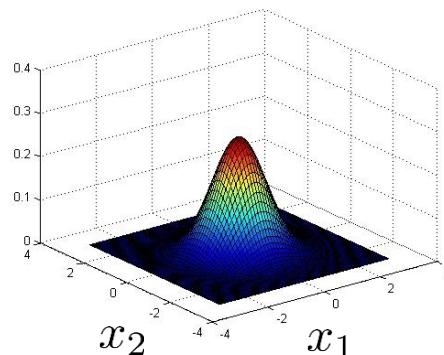
# Multivariate Gaussian (Normal)

examples

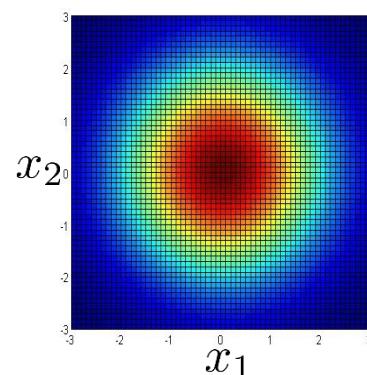
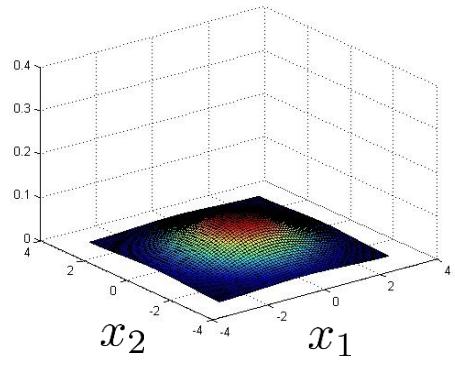
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}$$



$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 0.6 & 0 \\ 0 & 0.6 \end{bmatrix}$$



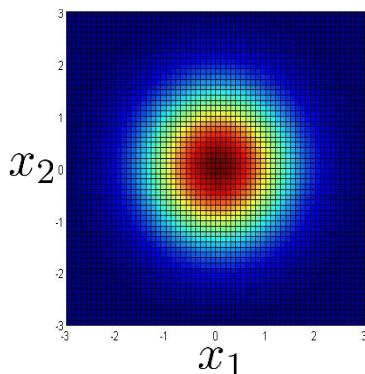
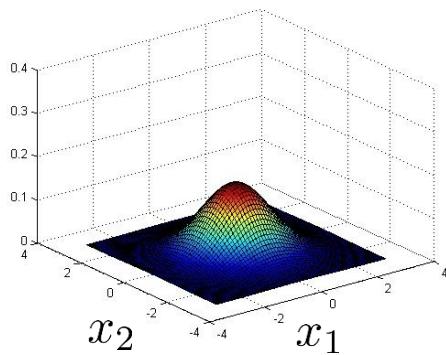
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$



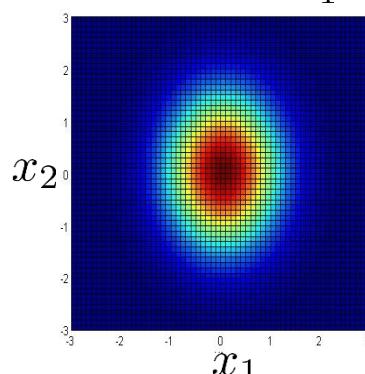
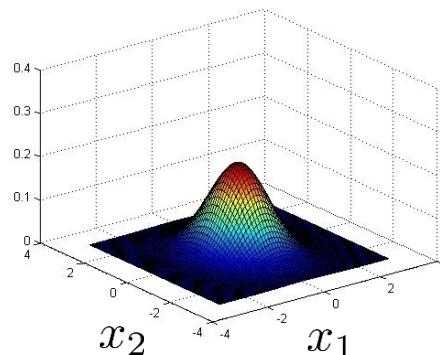
# Multivariate Gaussian (Normal)

examples

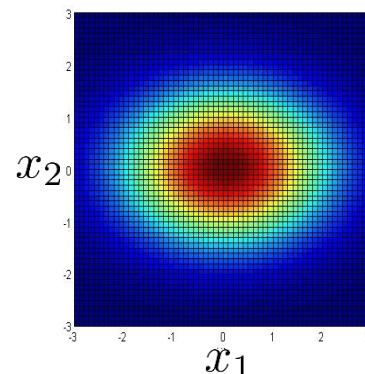
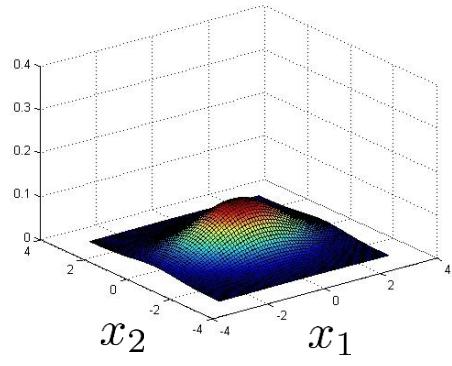
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}$$



$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 0.6 & 0 \\ 0 & 1 \end{bmatrix}$$



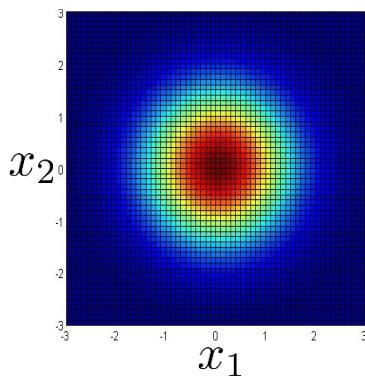
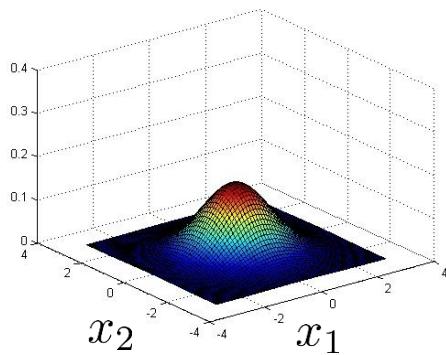
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}$$



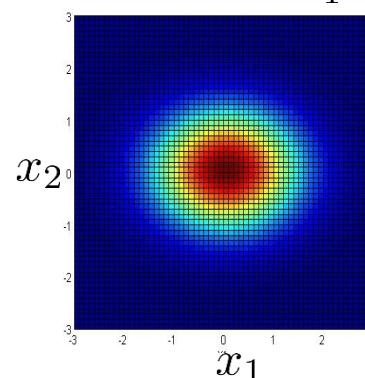
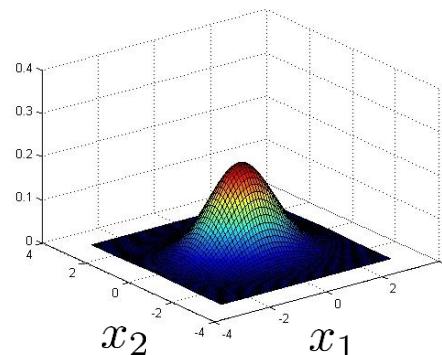
# Multivariate Gaussian (Normal)

## examples

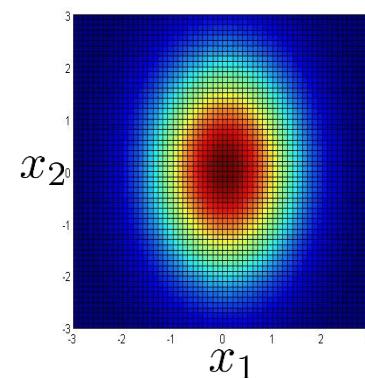
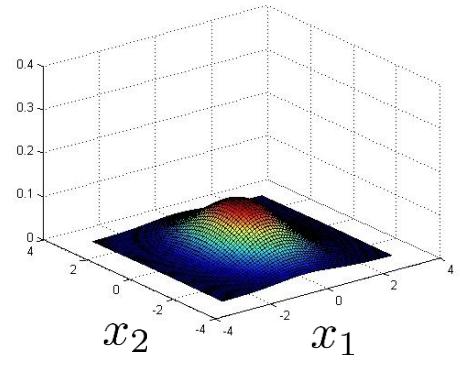
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}$$



$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 0.6 \end{bmatrix}$$



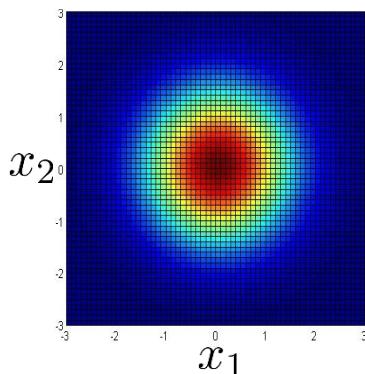
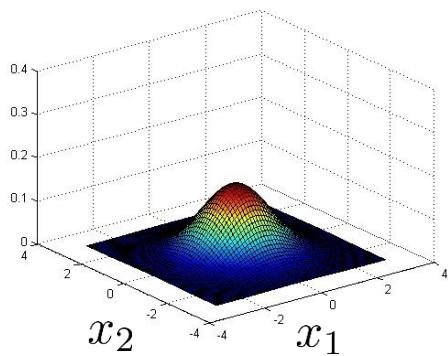
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$$



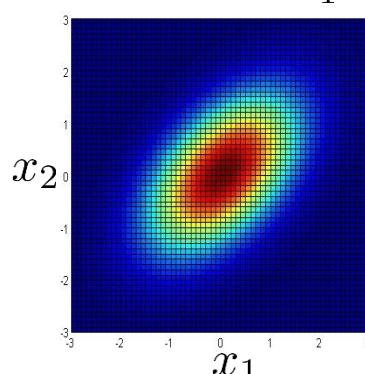
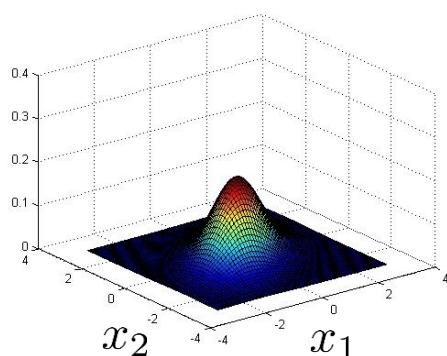
# Multivariate Gaussian (Normal)

examples

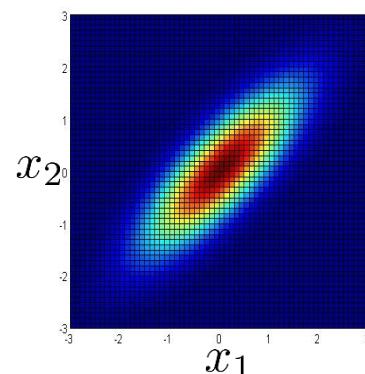
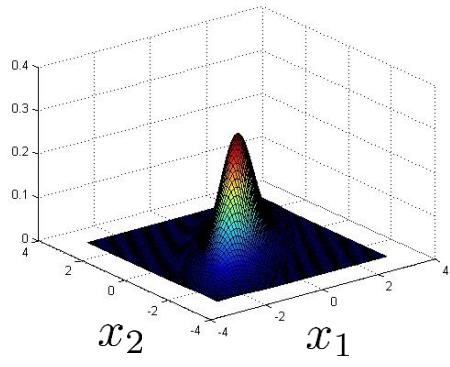
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}$$



$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$$



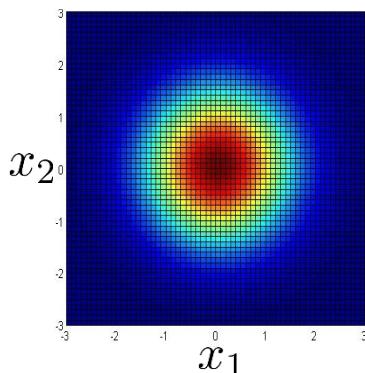
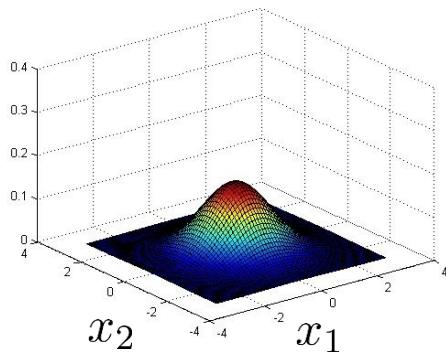
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$



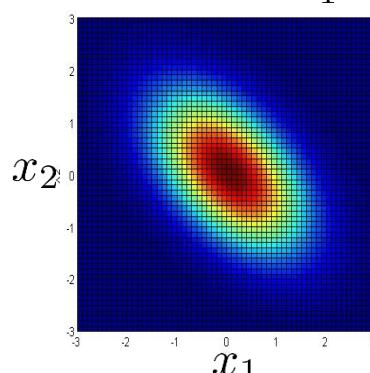
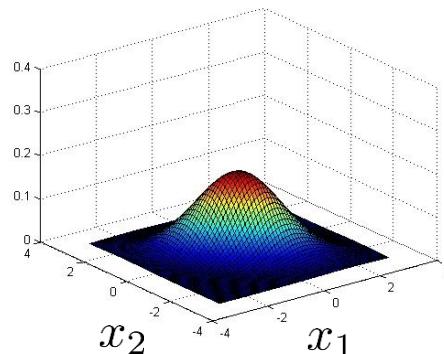
# Multivariate Gaussian (Normal)

## examples

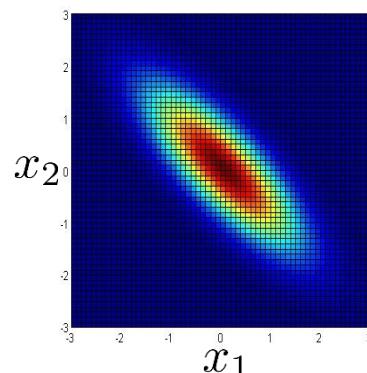
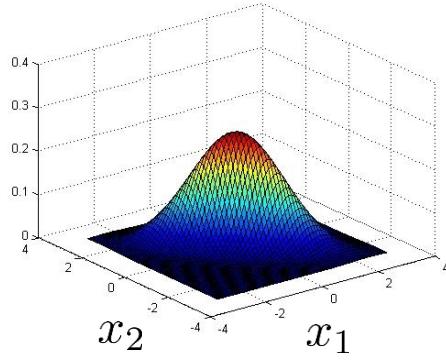
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}$$



$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix}$$



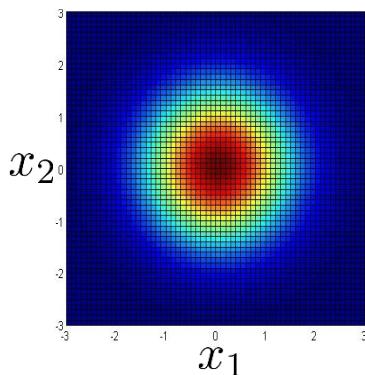
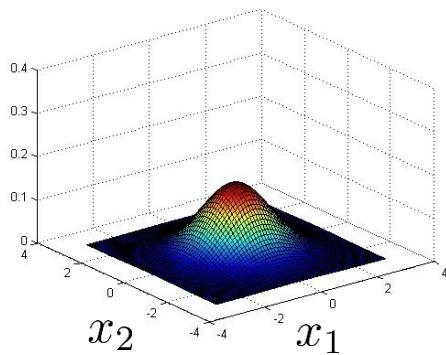
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 1 & -0.8 \\ -0.8 & 1 \end{bmatrix}$$



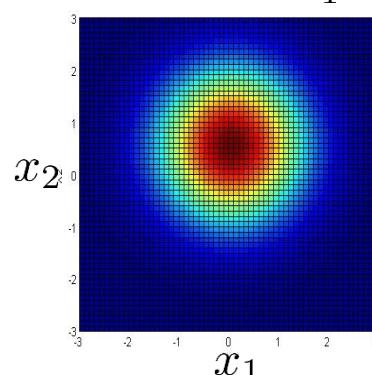
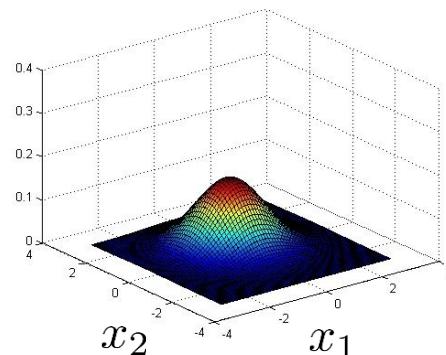
# Multivariate Gaussian (Normal)

## examples

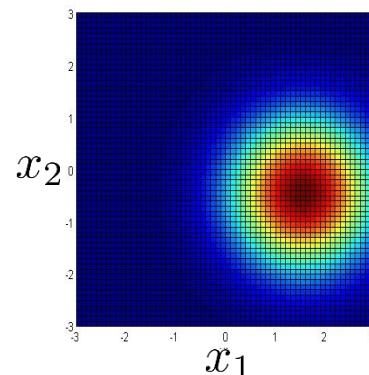
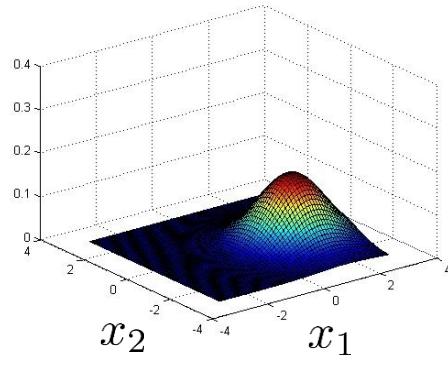
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}$$



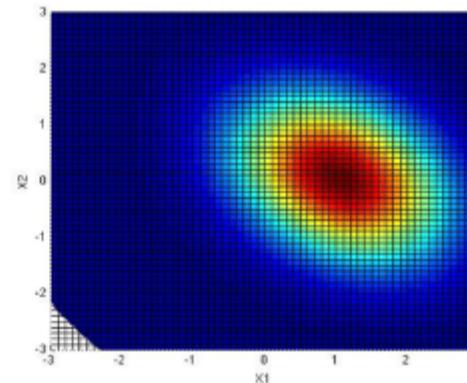
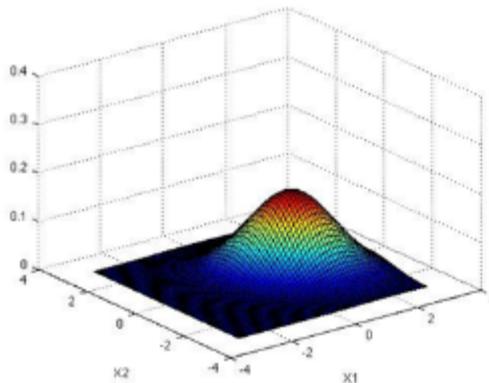
$$\mu = \begin{bmatrix} 0 \\ 0.5 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



$$\mu = \begin{bmatrix} 1.5 \\ -0.5 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



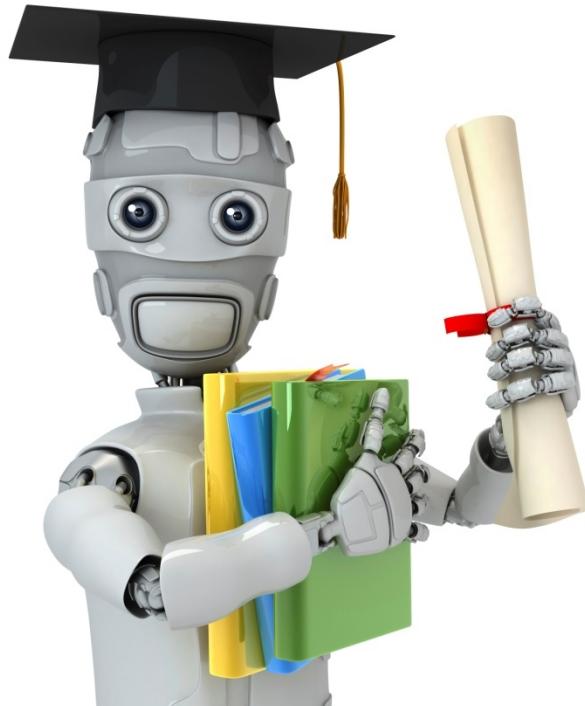
Consider the following multivariate Gaussian:



Which of the following are the  $\mu$  and  $\Sigma$  for this distribution?

- $\mu = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & 0.3 \\ 0.3 & 1 \end{bmatrix}$
- $\mu = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & -0.3 \\ -0.3 & 1 \end{bmatrix}$
- $\mu = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & 0.3 \\ 0.3 & 1 \end{bmatrix}$
- $\mu = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & -0.3 \\ -0.3 & 1 \end{bmatrix}$

Correct Response



Machine Learning

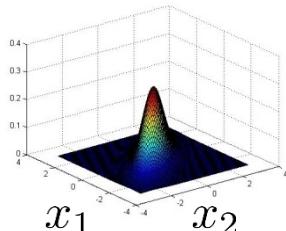
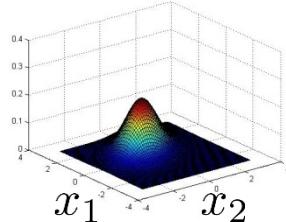
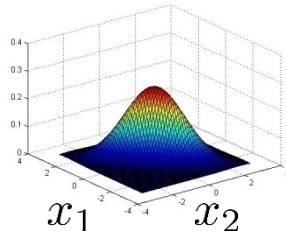
# Anomaly detection

Anomaly detection  
using the  
multivariate  
Gaussian  
distribution

# Multivariate Gaussian (Normal) distribution

Parameters,  $\Sigma$

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



Parameter fitting:

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

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

# Anomaly detection with the multivariate Gaussian

Given model  $p(x)$  by setting

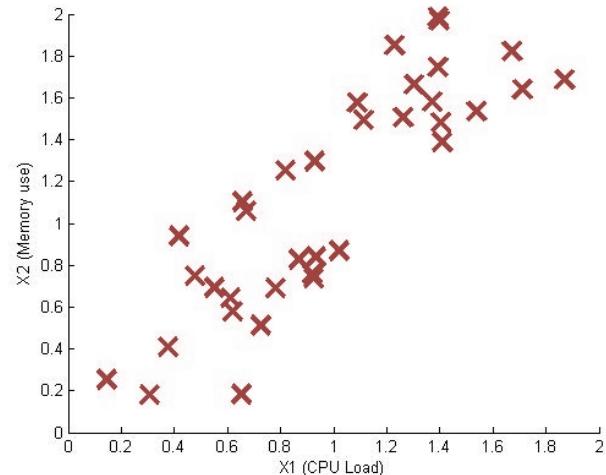
$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

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

2. Given a new example  $x$ , compute

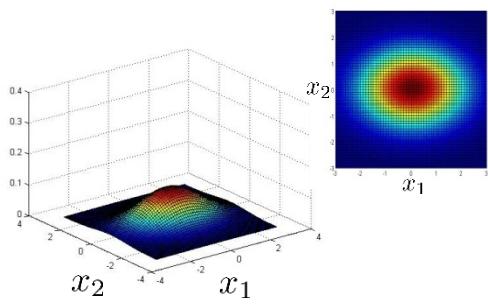
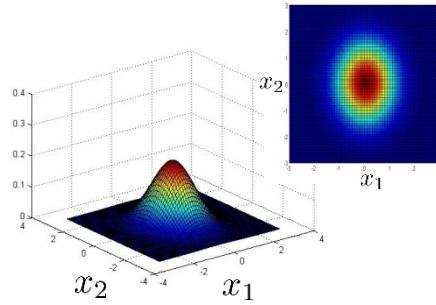
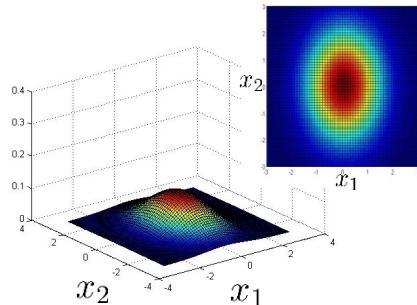
$$p(x) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left( -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right)$$

Flag an anomaly if  $p(x) < \varepsilon$



# Relationship to original model

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



Corresponds to multivariate Gaussian

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

where

## Original model

$$p(x_1; \mu_1, \sigma_1^2) \times \cdots \times p(x_n; \mu_n, \sigma_n^2)$$

Manually create features to capture anomalies where take unusual combinations of values.

Computationally cheaper (alternatively, scales better to large

OK even if  $n > m$  (training set size) is small

## vs. Multivariate Gaussian

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

Automatically captures correlations between features

Computationally more expensive

Must have  $m > n$ , or else is non-invertible.

Consider applying anomaly detection using a training set  $\{x^{(1)}, \dots, x^{(m)}\}$  where  $x^{(i)} \in \mathbb{R}^n$ . Which of the following statements are true? Check all that apply.

- The original model  $p(x_1; \mu_1, \sigma_1^2) \times \dots \times p(x_n; \mu_n, \sigma_n^2)$  corresponds to a multivariate Gaussian where the contours of  $p(x; \mu, \Sigma)$  are axis-aligned.

**Correct Response**

- Using the multivariate Gaussian model is advantageous when m (the training set size) is very small ( $m < n$ ).

**Correct Response**

- The multivariate Gaussian model can automatically capture correlations between different features in  $x$ .

**Correct Response**

- The original model can be more computationally efficient than the multivariate Gaussian model, and thus might scale better to very large values of  $n$  (number of features).

**Correct Response**