

# Machine Learning

## Week 9

### Section 1

#### Anomaly Detection

This week will cover anomaly detection and recommender systems. Anomaly detection usually takes a large number of data points and figures out which vary significantly from the average. Recommender systems look at patterns of activities and from different sources and products and produce recommendations. Some recommender system algorithms include the collaborative filtering algorithm and low-rank matrix factorization.

#### Density Estimation

#### Problem Motivation

Anomaly Detection is a common ML problem type. It is an unsupervised problem, but there are some aspects of it that are similar to supervised learning problems as well.

Ex. Anomaly Detection in Aircraft Engine Manufacturing.

#### Anomaly detection example

Aircraft engine features:

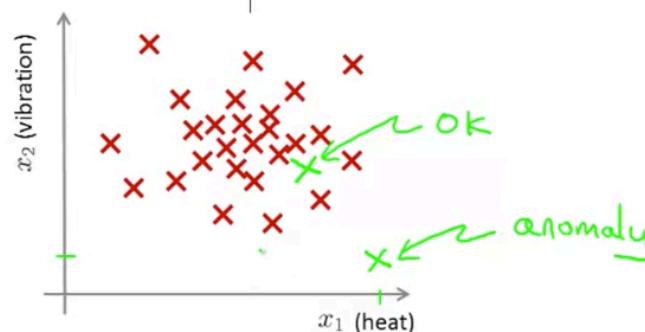
→  $x_1$  = heat generated

→  $x_2$  = vibration intensity

...

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

New engine:  $\underline{x_{test}}$



More formally, we are given a dataset and assume these examples are normal / non-anomalous.

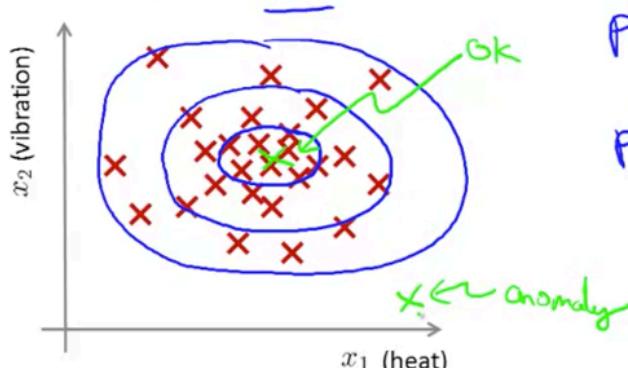
Given this unlabeled training set, we need to create a Model  $p(x)$  to predict if new cases are normal or anomalous.

#### 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

## Examples of Anomaly Detection Applications:

Here is an anomaly detection example of fraud detection.

This type of model may flag users who are acting unusually, not just fraudulently.

### 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  $\underline{p(x) < \varepsilon}$

$$\begin{array}{c} x_1 \\ x_2 \\ x_3 \\ x_4 \end{array} \quad p(x)$$

→ 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.

...

$$p(x) < \varepsilon$$

Similar to the classification problems we saw before, if there are too many false positives or true negatives being predicted by the model, you can try increasing or decreasing epsilon to tune the sensitivity.

# Gaussian Distribution

In this video we will discuss the Gaussian distribution / Normal distribution.

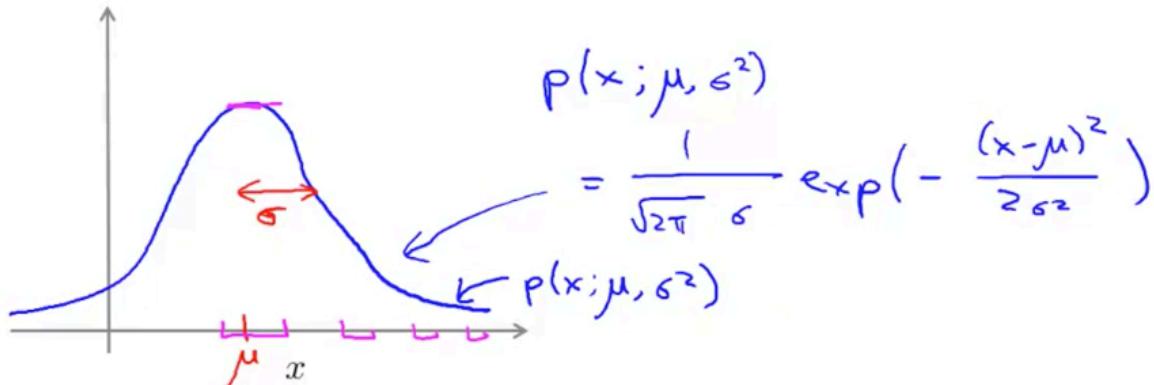
## Gaussian (Normal) distribution

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

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

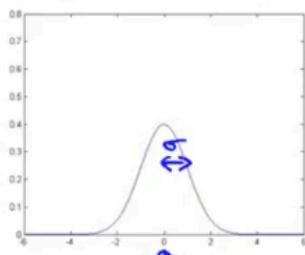
↖ "distributed as"

$\sigma$  standard deviation

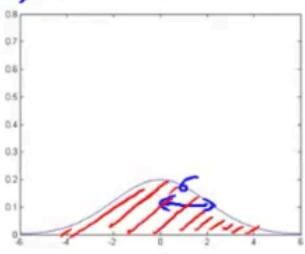


Some examples of the gaussian for different values of the mean (mu) and std (sigma)

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

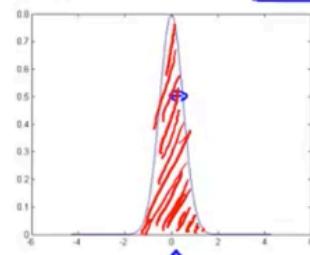


$$\rightarrow \mu = 0, \sigma = 2$$

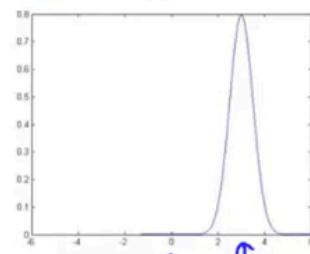


$$\rightarrow \mu = 0, \sigma = 0.5$$

$$\sigma^2 = 0.25$$



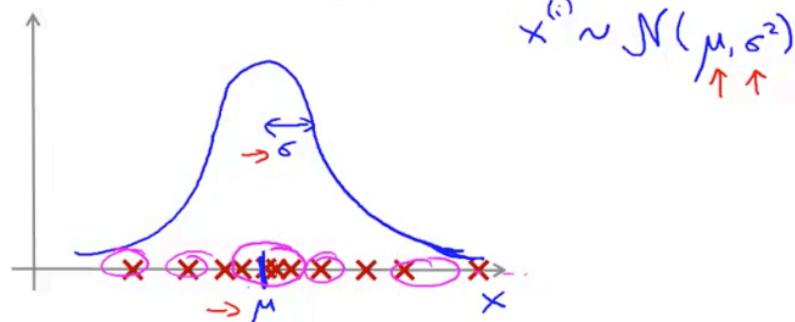
$$\rightarrow \mu = 3, \sigma = 0.5$$



Andrew N

Next lets talk about the parameter estimate problem:

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



If we suspect that our examples come from a Gaussian distribution w/ unknown mu and std. How do we estimate what those parameters could be?

The standard formulas for these estimates are as follows:

$$\rightarrow \underline{\mu} = \frac{1}{m} \sum_{i=1}^m x^{(i)} \quad \rightarrow \underline{\sigma^2} = \frac{1}{m} \sum_{i=1}^m \underline{(x^{(i)} - \underline{\mu})^2}$$

These are the 'maximum likelihood estimates' from statistics.

\*\*Note in some cases for the variances formula, people use  $1 / (m-1)$  instead of  $1 / m$ , but typically for ML it is  $1 / m$ .  
The two versions have slightly different mathematical properties but for our cases  $1 / m$  is fine.

# Anomaly Detection Algorithm

In this video we will apply the Gaussian distribution to develop an anomaly detection algorithm.

The problem of density estimation, we need to compute the gaussian for each feature on each training example.

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

$$x_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$$

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

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

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

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

$$p(x)$$

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

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

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

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

In statistics this equation corresponds to an independent assumption on the values of the features  $x_1$  through  $x_n$ .

Putting everything together we have the following algorithm:

## Anomaly detection algorithm

→ 1. Choose features  $x_i$  that you think might be indicative of anomalous examples.  $\{x^{(1)}, \dots, x^{(m)}\}$

→ 2. Fit parameters  $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$

$$\begin{aligned} \rightarrow \mu_j &= \frac{1}{m} \sum_{i=1}^m x_j^{(i)} & p(x_j; \mu_j, \sigma_j^2) & \begin{matrix} \mu_1, \mu_2, \dots, \mu_n \\ \mu = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix} = \frac{1}{m} \sum_{i=1}^m x^{(i)} \end{matrix} \\ \rightarrow \sigma_j^2 &= \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2 & \uparrow \quad \uparrow & \uparrow \end{aligned}$$

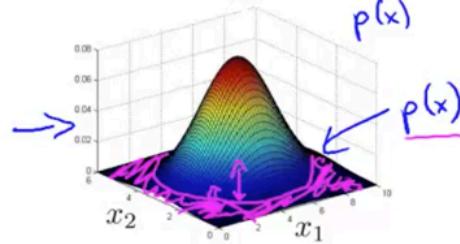
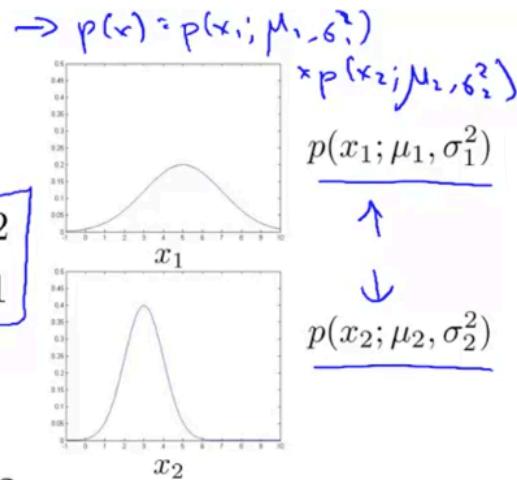
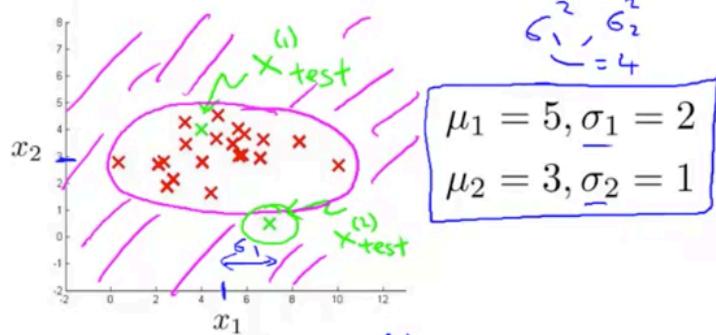
→ 3. Given new example  $x$ , compute  $p(x)$ :

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

Anomaly if  $p(x) < \varepsilon$

Example of an application of this method:

## Anomaly detection example



$$\varepsilon = 0.02$$

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

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

# Building An Anomaly Detection System

## Developing and Evaluating an Anomaly Detection System

In the last video we talked about the process of developing a specific application of anomaly detection to a particular problem. This this video we will focus on the process of evaluating an anomaly detection algorithm.

Recall from previous evaluation processes, it is easier to make decisions about these parameters and quality of your models by having a real number metric you can use.

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

\*\*It's okay if a negligible amount of anomalous data is part of the assumed normal examples in the training data set.

### Aircraft engines motivating example

- 10000 good (normal) engines
- 20 flawed engines (anomalous)  $\frac{2}{50}$   $y=1$
- Training set: 6000 good engines ( $y=0$ )  $p(x) = p(x_1; \mu_1, \sigma_1^2) \dots p(x_n; \mu_n, \sigma_n^2)$
- CV: 2000 good engines ( $y=0$ ), 10 anomalous ( $y=1$ )
- Test: 2000 good engines ( $y=0$ ), 10 anomalous ( $y=1$ )

Use the training set (normal examples, non-anomalous, but unlabeled) to fit  $p(x)$

Given the training set, cross validation set, and test set, here is how we develop an evaluation algorithm:

### Algorithm evaluation

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

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

Possible evaluation metrics:

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

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

As with the classification algorithm and skewed datasets, we use the confusion matrix and calculate the F1 score as our metric, and continually tune our parameters on the cross validation set.

# Anomaly Detection vs. Supervised Learning

For anomaly detection algorithms, we saw in the previous video that we used some labelled data for evaluating performance, even though this is an unsupervised learning algorithm.

Why don't we just use a supervised learning algorithm to perform anomaly detection?

Let's look at some guidelines for deciding when to use which algorithm:

<b>Anomaly detection</b>	<b>vs.</b>	<b>Supervised learning</b>
<ul style="list-style-type: none"><li>→ Very small number of positive examples (<math>y = 1</math>). (0-20 is common).</li><li>→ Large number of negative (<math>y = 0</math>) examples. <math>p(x)</math></li><li>→ Many different "types" of anomalies. Hard for any algorithm to learn from positive examples what the anomalies look like;</li><li>→ future anomalies may look nothing like any of the anomalous examples we've seen so far.</li></ul>		<p>Large number of positive and negative examples.</p> <p>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.</p>

\*\* Negative examples -> Non-anomalous

\*\* Note that spam is a supervised learning problem because we have such a large collection of spam and non-spam emails even though there are so many different types of spam.

There are some cases where anomaly detection could be a supervised problem but it depends on the size and ratio of the dataset.

Example of different applications:

<b>Anomaly detection</b>	<b>vs.</b>	<b>Supervised learning</b>
<ul style="list-style-type: none"><li>→ • Fraud detection <math>y=1</math></li><li>→ • Manufacturing (e.g. aircraft engines)</li><li>→ • Monitoring machines in a data center</li><li>⋮</li></ul>		<ul style="list-style-type: none"><li>• Email spam classification</li><li>• Weather prediction (sunny/rainy/etc).</li><li>• Cancer classification</li><li>⋮</li></ul>

# Choosing What Features to Use

One of the biggest factors on performance for anomaly detection algorithms are the features we use/choose to give to the anomaly detection algorithm.

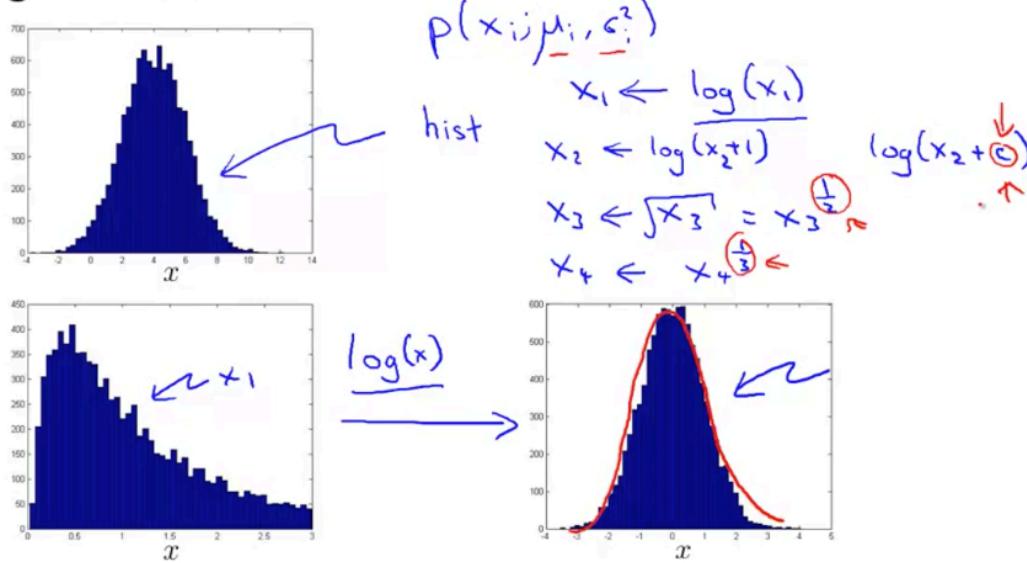
What are the best practices when choosing features to give to an anomaly detection algorithm?

Ideally we want to use gaussian features.

If we have non-gaussian data, it's not the end of the world it can still work.

- We can try playing around with different transformations on the data to make it look more gaussian.
  - There are many different transforms we can try out and play with.

## Non-gaussian features



We can come up with a choice of features by using an error analysis procedure (similar to the error analysis procedure we used for supervised learning):

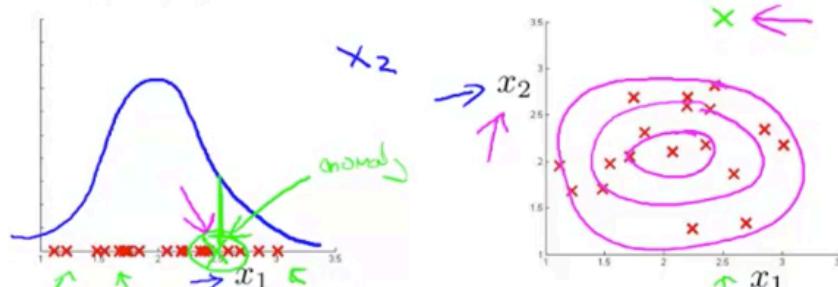
What if the green example is incorrectly classified? We should explore that feature manually and see if there is a problem, ex.  $x_2$  is abnormally high. Maybe we can create a new feature that would work better.

## → Error analysis for anomaly detection

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

Most common problem:

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



What are some other good guidelines for picking features?

Ex.

CPU load and Network traffic should grow linearly, so if one is high and one is low that may indicate an anomaly.

→ **Monitoring computers in a data center**

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

→  $x_1$  = memory use of computer

→  $x_2$  = number of disk accesses/sec

→  $x_3$  = CPU load ←

→  $x_4$  = network traffic ←

$$x_5 = \frac{\text{CPU load}}{\text{network traffic}}$$

$$x_6 = \frac{(\text{CPU load})^2}{\text{network traffic}}$$

# Multivariate Gaussian Distribution

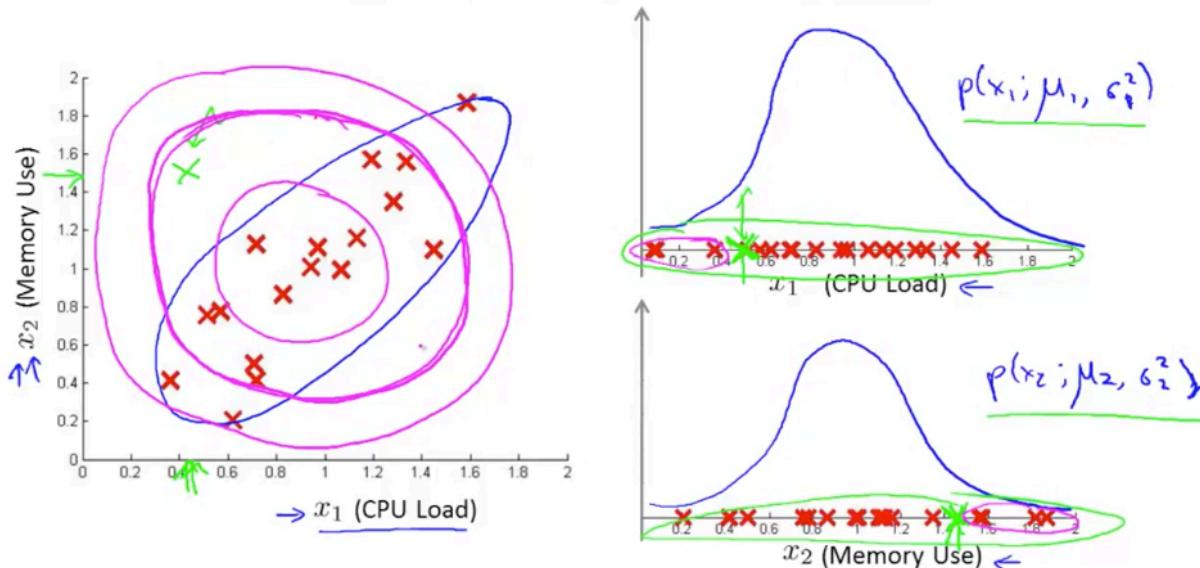
## Multivariate Gaussian Distribution

In this video we will discuss an extension to the Anomaly Detection Algorithm we've developed so far.

This extension uses the Multivariate Gaussian Distribution and has some advantages and disadvantages for Anomaly Detection over the previous algorithm.

Ex. Here we can see if the CPU Load and Memory Use grow linearly, but if we look at the individual values of  $x_1$  and  $x_2$  for the green example, both still lie within the Gaussian distribution (not too far out).

### **Motivating example: Monitoring machines in a data center**



The algorithm doesn't understand that outside of the blue ellipse, the probability of an anomaly is much higher.

As the circles increase, there is higher probability of an anomaly.

As such, we need to create a modified version of the anomaly detection algorithm using the multivariate Gaussian (normal) distribution.

### **Multivariate Gaussian (Normal) distribution**

$\rightarrow x \in \mathbb{R}^n$ . Don't model  $p(x_1), p(x_2), \dots$ , etc. separately.

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

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

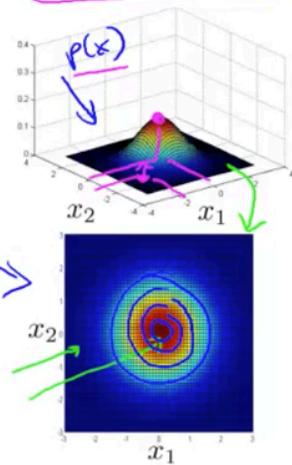
$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x - \mu)^\top \Sigma^{-1} (x - \mu)\right)$$

$|\Sigma| = \text{determinant of } \Sigma \quad |\det(\Sigma)|$

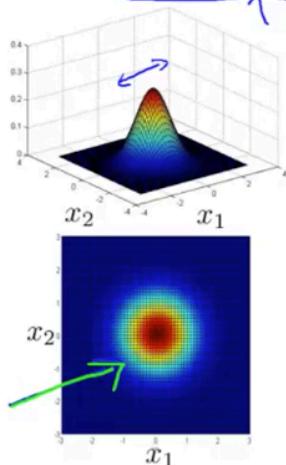
What does this 'new/modified'  $p(x)$  look like? Let us look at a 2D example.

## Multivariate Gaussian (Normal) examples

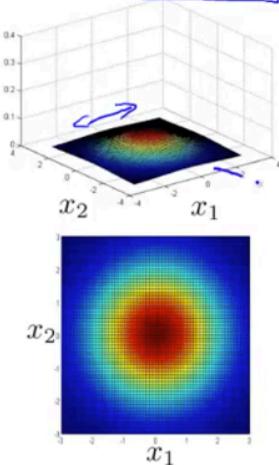
$$\rightarrow \mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



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

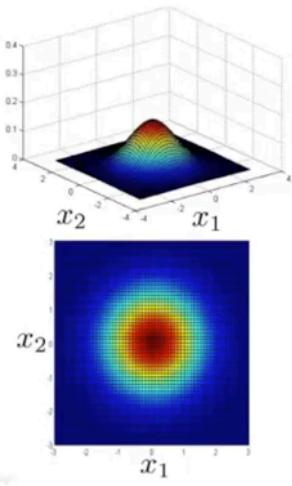


$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$

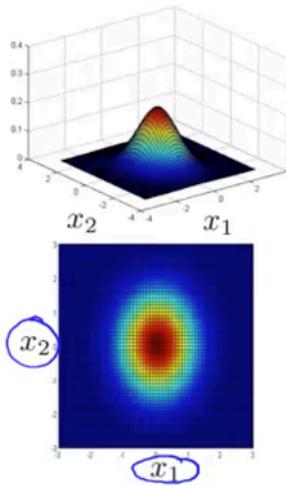


## Multivariate Gaussian (Normal) examples

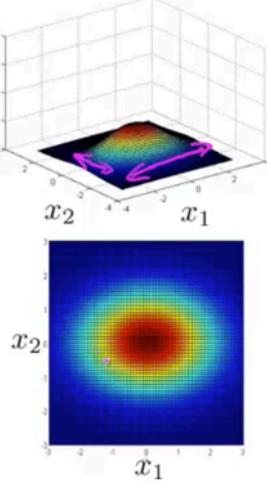
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



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



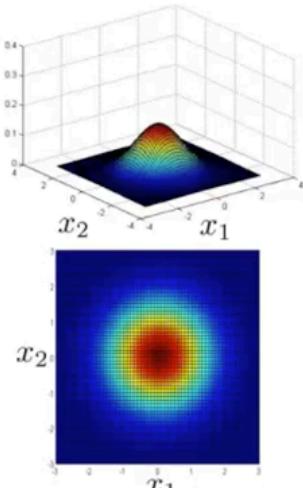
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 0 \end{bmatrix}$$



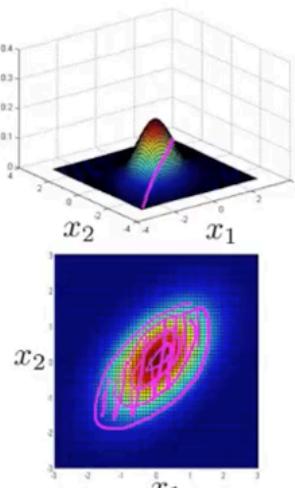
We can also use this to model correlations:

## Multivariate Gaussian (Normal) examples

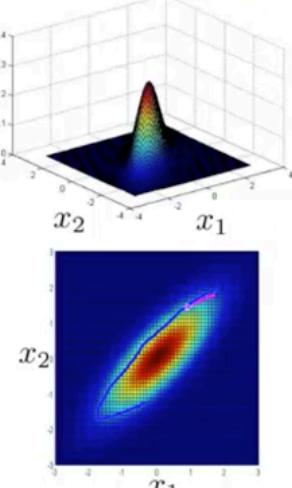
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



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



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

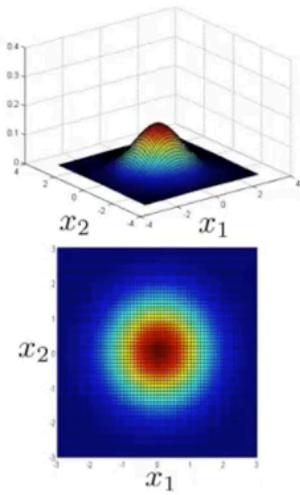


Ar

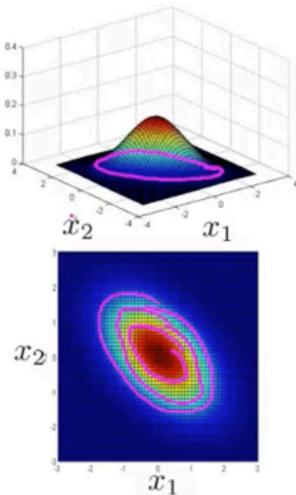
We can also use negative values (negative correlations):

## Multivariate Gaussian (Normal) examples

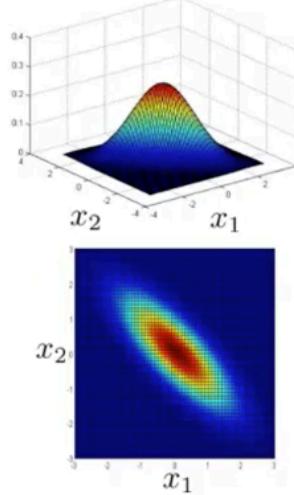
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



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



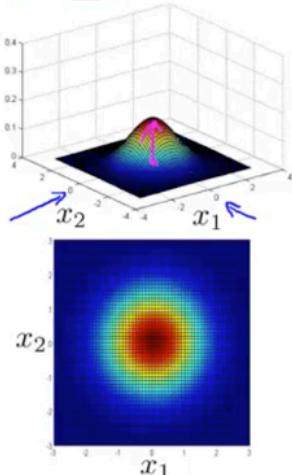
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & -0.8 \\ -0.8 & 1 \end{bmatrix}$$



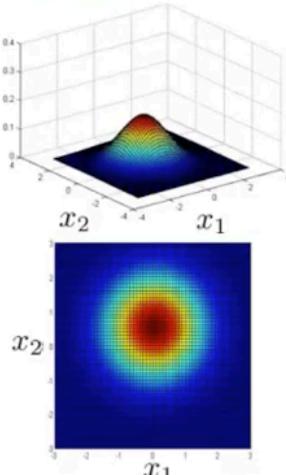
Above we are only changing the covariance matrix, but we can also change the values for mu:

## Multivariate Gaussian (Normal) examples

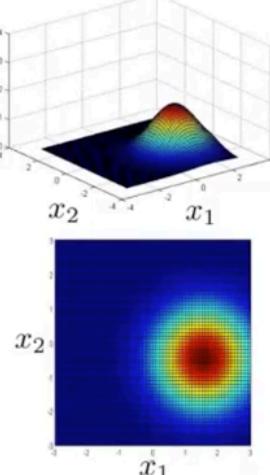
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



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



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



As such with this approach we can capture when we expect two different features to be positively or negatively correlated.

# Anomaly Detection using the Multivariate Gaussian distribution

In this video we will apply the ideas we learn from the previous video to developing a different Anomaly Detection Algorithm.

Recall we have two parameters, mu and sigma (the covariance matrix)

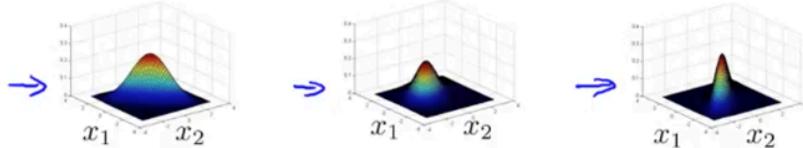
Lets look at the parameter fitting / estimation problem.

Parameters  $\mu, \Sigma$

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

$$\Sigma \in \mathbb{R}^{n \times n}$$

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



Parameter fitting:

Given training set  $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\} \leftarrow x \in \mathbb{R}^n$

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

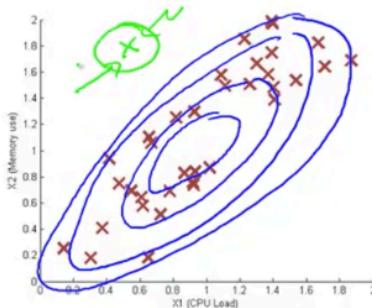
## Anomaly detection with the multivariate Gaussian

1. Fit model  $p(x)$  by setting

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

2. Given a new example  $x$ , compute

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



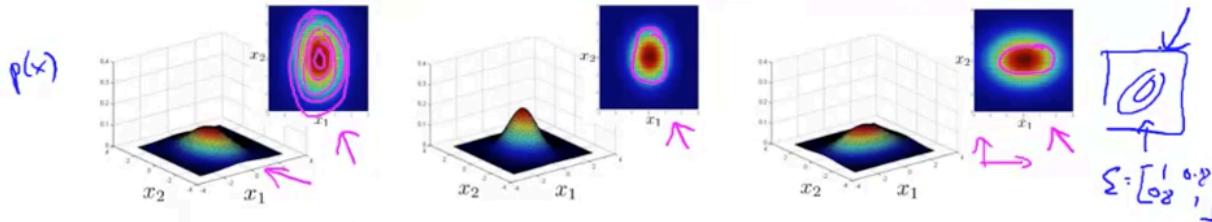
Flag an anomaly if  $p(x) < \varepsilon$

What is the relationship between the multivariate Gaussian distribution model and the original model?

The original model turns out to be a special case of the multivariate Gaussian distribution model where the distribution of  $p(x)$  is constrained so that the contours of the probability distribution functions are 'axis aligned'.

## 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

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

where  $\Sigma = \begin{bmatrix} \sigma_1^2 & \dots & \dots \\ \dots & \dots & \dots \\ \dots & \dots & \sigma_n^2 \end{bmatrix}$

i.e. horizontal and vertical, not diagonal.

Mathematically this means the covariance matrix sigma is constrained as shown.

So when would you use either of these two models?

→ Original model

vs. → Multivariate Gaussian

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

Manually create features to capture anomalies where  $x_1, x_2$  take unusual combinations of values.

$$\rightarrow X_3 = \frac{x_1}{x_2} = \frac{\text{CPU load}}{\text{memory}}$$

→ Computationally cheaper (alternatively, scales better to large  $n$ )  $n=10,000, m=100,000$

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

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

→ Automatically captures correlations between features

$$\Sigma \in \mathbb{R}^{n \times n}$$

$$\Sigma^{-1}$$

Computationally more expensive

$$\rightarrow \Sigma \sim \frac{n^2}{2}$$

$$\begin{cases} \rightarrow x_1 = x_2 \\ \cancel{x_2} = x_4 + x_5 \end{cases}$$

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

$$m \geq 10n$$

Andrew Ng

The original model is probably used more often, whereas the multivariate model is used less often but is able to automatically capture the correlations between features.

\*\* Note redundant features == linearly dependent.