

15. Anomaly Detection

>> PROBLEM MOTIVATION:

Inputs:

Anomaly detection example

Aircraft engine features:

→ x_1 = heat generated

→ x_2 = vibration intensity

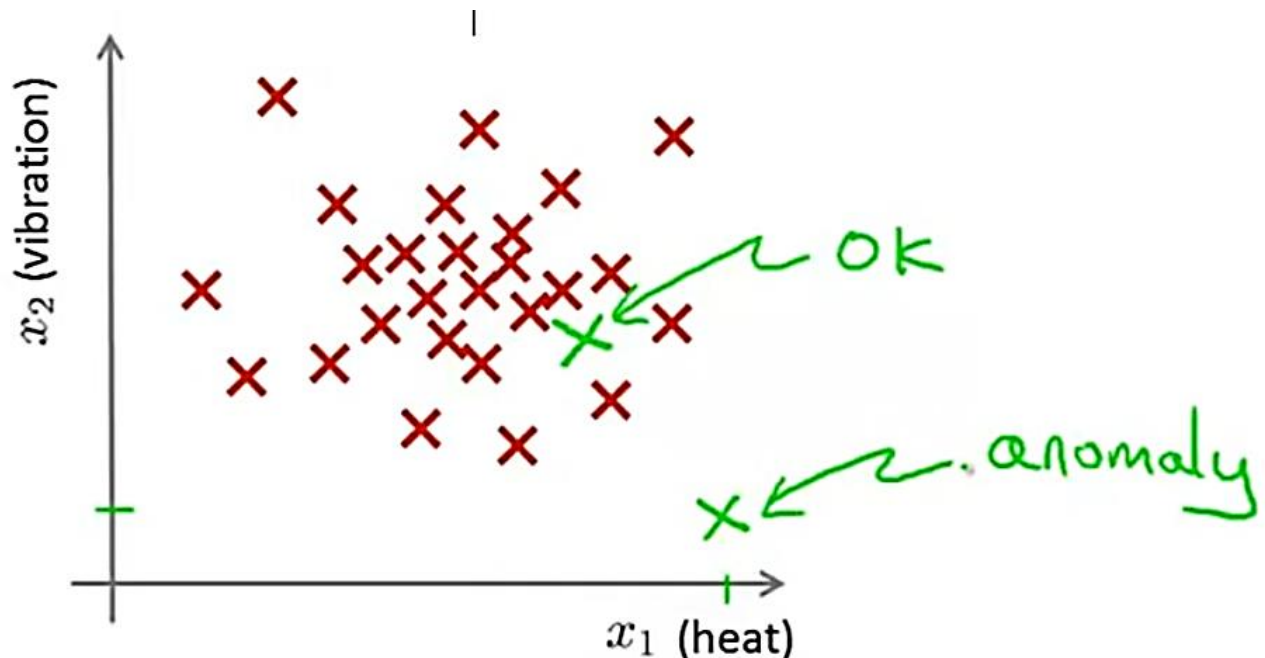
...

Test Data:

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

New engine: x_{test}

Plot:



Problem: we are given some training examples. We assume them to be NON-ANOMALOUS. So, we need to find if a new example is anomalous or not.

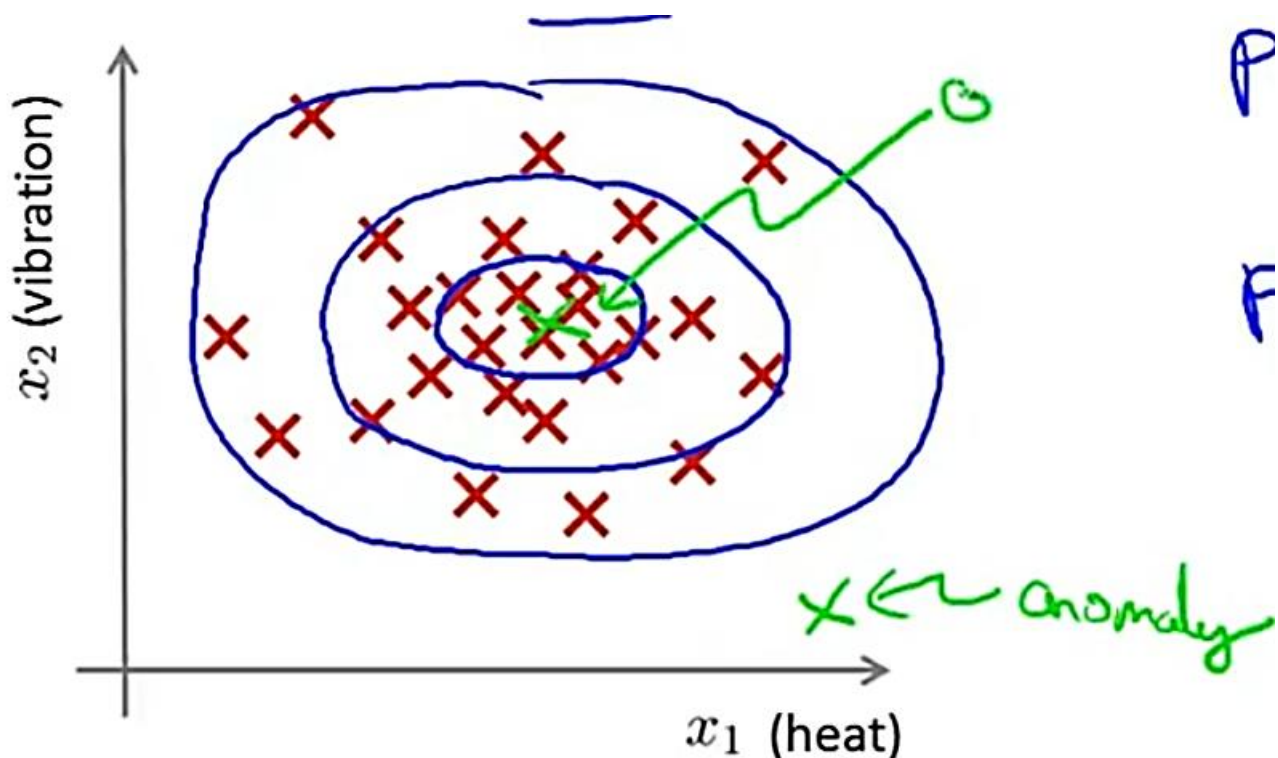
We train a PROBABILITY MODEL:

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

It divides the plot into various regions, such that each region corresponds to a level of provability.

Centre \rightarrow highest probability

Outside \rightarrow lowest P



Example:

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 $p(x) < \epsilon$

→ 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.

... $p(x) < \epsilon$

x_1
 x_2
 x_3
 x_4 $p(x)$

>> GAUSSIAN DISTRIBUTION:

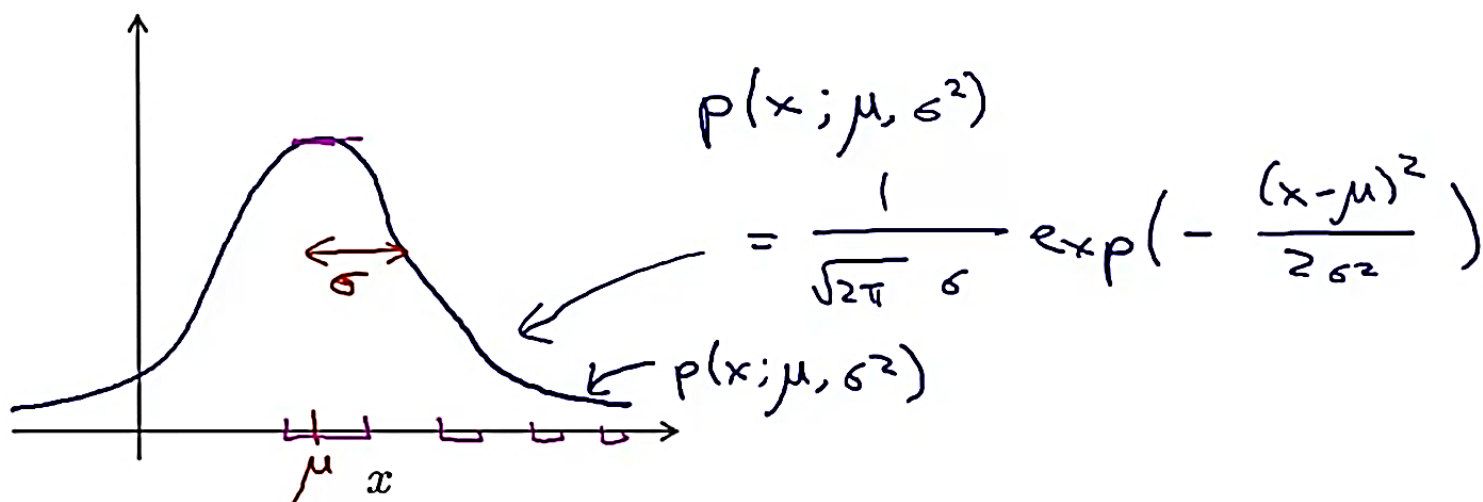
Gaussian (Normal) distribution

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

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

↗ "distributed as"

σ standard devio



Examples:

$\sigma \rightarrow$ controls the width and height of curve

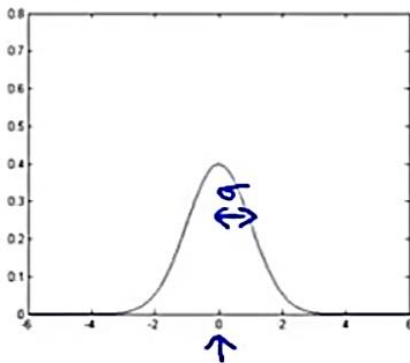
$\mu \rightarrow$ mean of all values of x : controls the center..

Area below the curve is always = 1

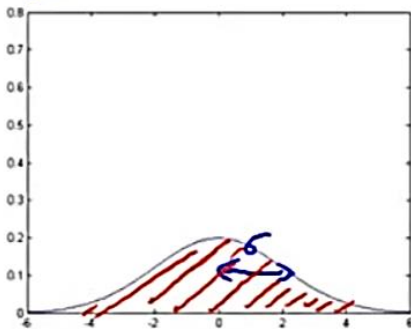
So, the curve is either taller or wider.

Gaussian distribution example

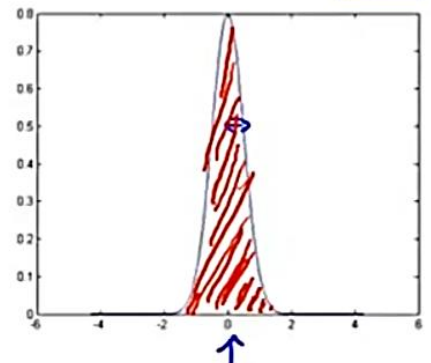
$\rightarrow \mu = 0, \sigma = 1$



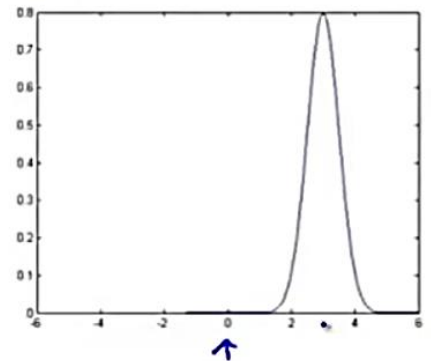
$\rightarrow \mu = 0, \sigma = 2$



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



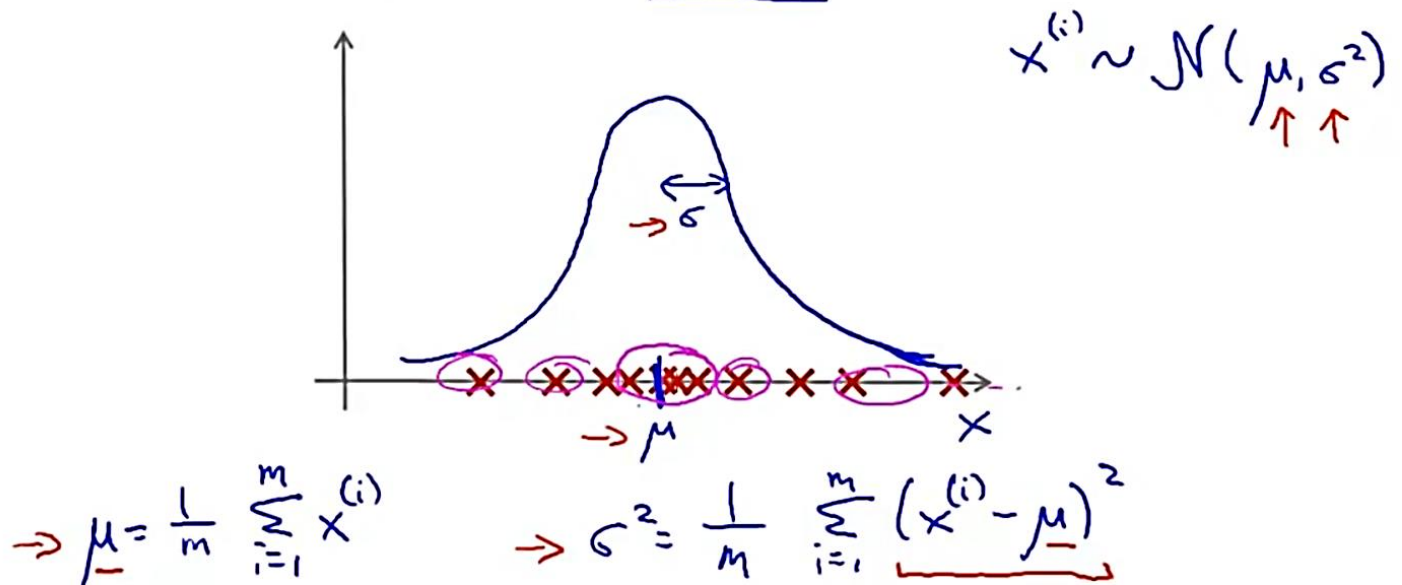
$\rightarrow \mu = 3, \sigma = 0.5$



Estimation of μ and σ :

Parameter estimation

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



>> ALGORITHM FOR ANOMALY DETECTION:

Density estimation

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

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

Each feature can be individually distributed using Normal(Gaussian) distribution:

$$x_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$$
$$x_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$$

$P(x)$ = product of probabilities of individual features

$$\begin{aligned}
 p(x) &= p(x_1; \mu_1, \sigma_1^2) p(x_2; \mu_2, \sigma_2^2) p(x_3; \mu_3, \sigma_3^2) \dots p(x_n; \mu_n, \sigma_n^2) \\
 &= \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2)
 \end{aligned}$$

Note:

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

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

Steps of algorithm:

Anomaly detection algorithm

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

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

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

$$p(x_j; \mu_j, \sigma_j^2)$$

$$\mu_1, \mu_2, \dots, \mu_n$$

$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix} = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

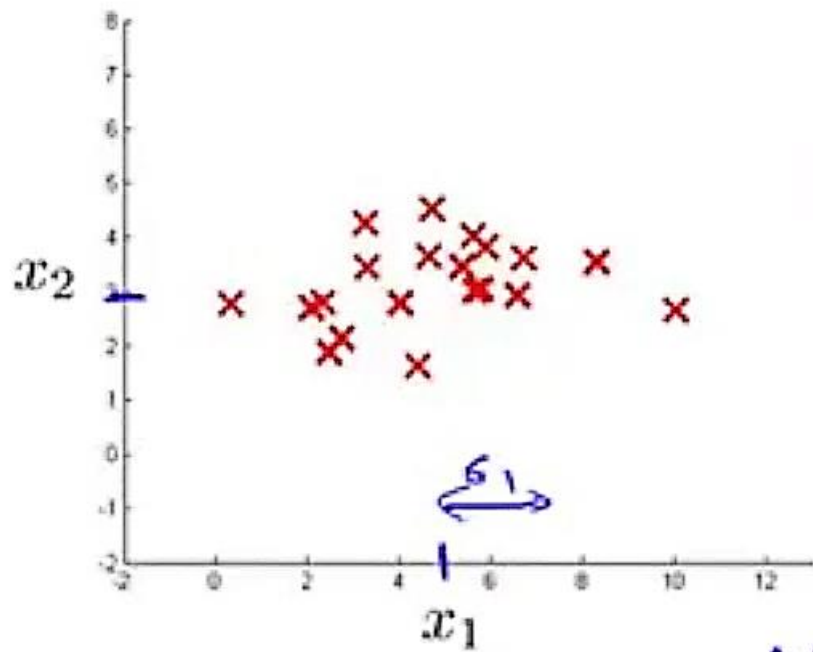
$$\rightarrow \sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2 \leftarrow$$

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 $p(x) < \varepsilon$

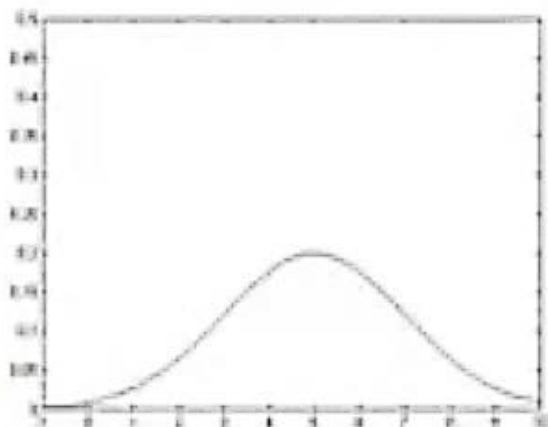
Example:



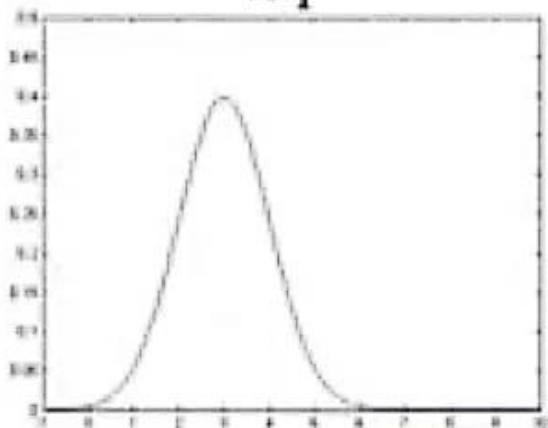
$$\mu_1 = 5, \underline{\sigma_1} = 2$$

$$\mu_2 = 3, \underline{\sigma_2} = 1$$

Gaussian Curves of both features:



x_1



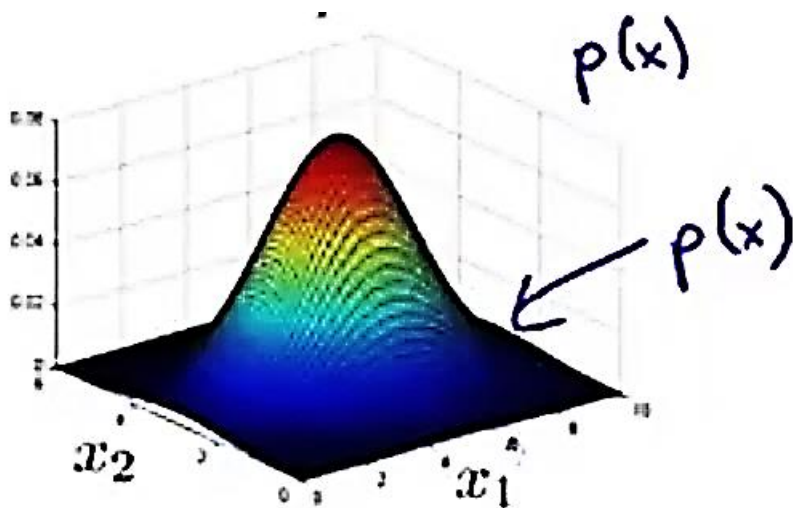
x_2

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

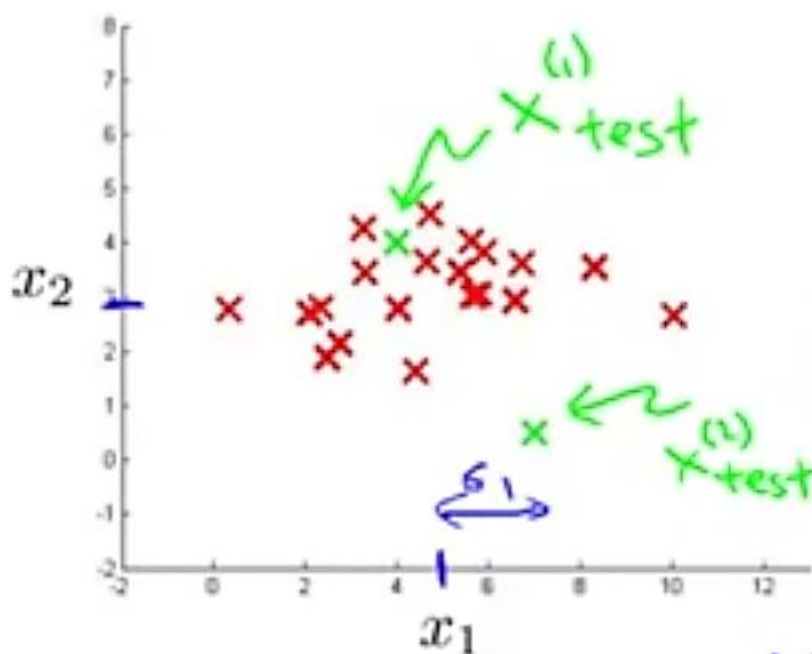


$$\underline{p(x_2; \mu_2, \sigma_2^2)}$$

$P(x) \rightarrow$ height of curve $= P(x_1; \mu_1, \sigma^2_1) \times P(x_2; \mu_2, \sigma^2_2)$



For a new example:



$$\underline{\varepsilon = 0.02}$$

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

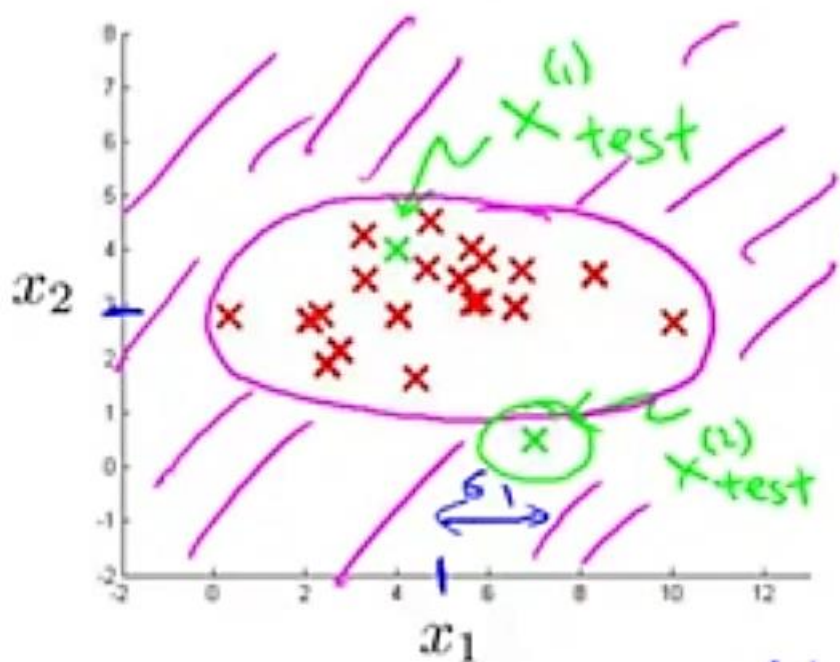
$$p(x_{test}^{(2)}) = \underline{0.0021} < \varepsilon$$

This means that anything below a particular height in the plot, given by ϵ : is an anomaly

OR

Anything outside the learned region is an anomaly:

i.e., everything outside the magenta curve:



» DEVELOPING ANOMALY DETECTION SYSTEM:

Anomaly detection system problem can be converted to a likes of Supervised Learning problem.

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.

Therefore:

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

Thus we can take:

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

Meaning all the training set examples are non-anomalous

In Cross validation set and test set, we can have some examples of **anomalous ($y=1$)** and some of **non-anomalous ($y=0$)** type:

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

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

Example:

Aircraft engines motivating example

- 10000 good (normal) engines
- 20 flawed engines (anomalous) $\underline{2-50}$ $y=1$
- 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 ($y=1$)

Algorithm evaluation

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

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

>> How to evaluate the algorithm:

Possible evaluation metrics:

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

>> How to choose ϵ :

Can also use cross validation set to choose parameter ϵ

We can choose the value of ϵ which gives the best value of F_1 score.

Anomaly detection vs Supervised Learning:

Anomaly detection	vs.	Supervised learning
→ Very small number of positive examples (<u>$y = 1$</u>). (<u>0-20</u> is common).		Large number of positive and negative examples. ←
→ Large number of negative (<u>$y = 0$</u>) examples. $p(x)$ ←		
→ <u>Many different "types" of anomalies</u> . Hard for any algorithm to learn from positive examples what the anomalies look like;		Enough positive examples for algorithm to get a sense of what positive examples are like, future ←
→ future anomalies may look nothing like any of the anomalous examples we've seen so far.		positive examples likely to be similar to ones in training set. ←

In anomaly detection, it's better to model anomalies based on negative examples, rather than positive examples... as future anomalies may be totally different.

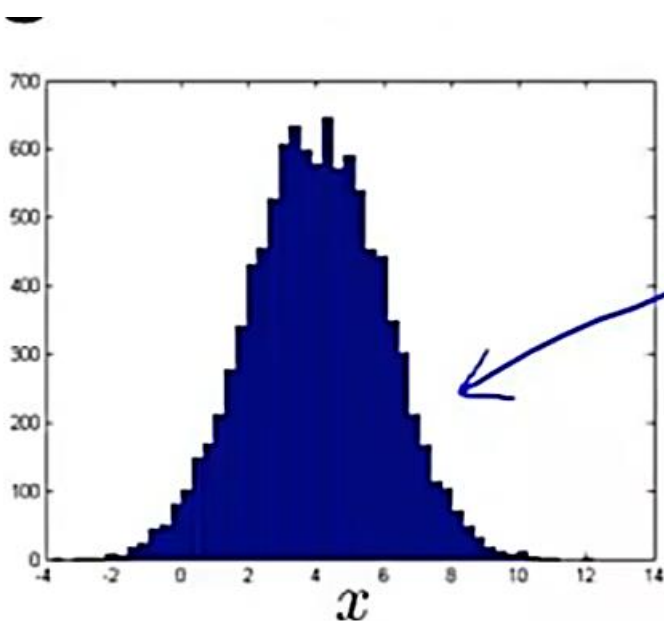
Applications of Anomaly detection vs Supervised Learning:

Anomaly detection	vs.	Supervised learning
→ • <u>Fraud detection</u> $y=1$	→	• Email spam classification ←
→ • Manufacturing (e.g. aircraft engines)		• Weather prediction (sunny/rainy/etc). ←
→ • Monitoring machines in a data center		• Cancer classification ←
⋮		⋮

Fraud detection can be a supervised learning application but only if there are a lot of people on the website who are doing fraudulent activity, i.e., most of the examples are positive, otherwise, it's an anomaly detection problem only

» Choosing what features to use:

We plot the data and the histogram looks like :

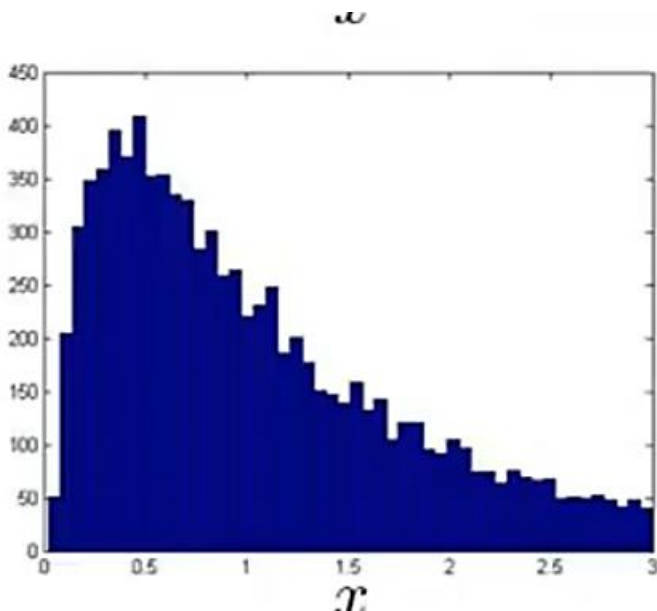


$$p(x_i; \underline{\mu_i}, \underline{\sigma_i^2})$$

hist

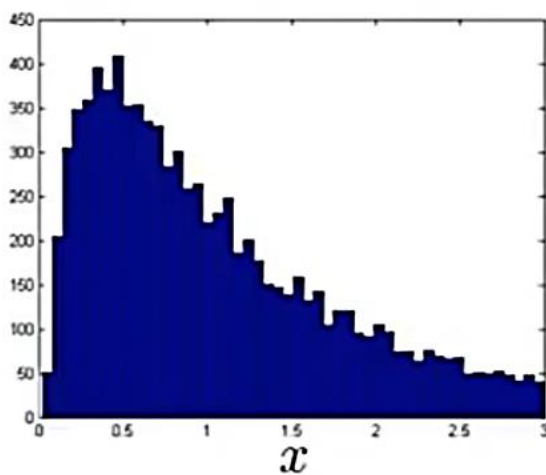
We'd be happy to see this as this means that the feature x is a gaussian feature

But: if the histogram looks like:



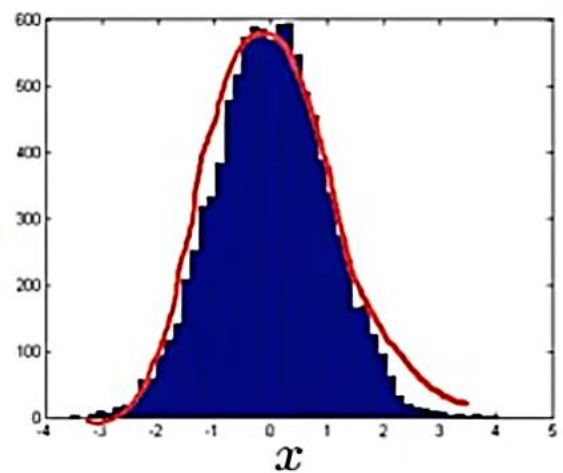
This is a non gaussian distribution

So, we use different transforms on our data to make it as close as possible to a gaussian distribution:



$\log(x)$

→



→ These are feature v/s $P(x)$

We can use different transforms for different features to make them gaussian features:

We may have to try out different transforms for the same feature to find the best (which gives the best gaussian look to the data).

$$\begin{aligned}
 x_1 &\leftarrow \frac{\log(x_1)}{\log(x_1)} \\
 x_2 &\leftarrow \log(x_2 + 1) \\
 x_3 &\leftarrow \sqrt{x_3} = x_3^{\frac{1}{2}} \\
 x_4 &\leftarrow x_4^{\frac{1}{3}}
 \end{aligned}$$

$\log(x_2 + \textcircled{c})$
 \downarrow
 \uparrow

These parameters in the red are parameters we can vary to make the data look more and more Gaussian.

>> Coming up with features:

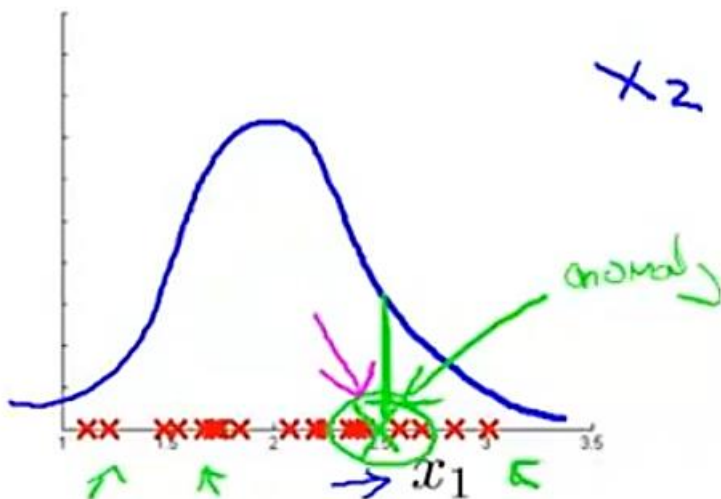
Error analysis method:

> **Error analysis for anomaly detection**

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

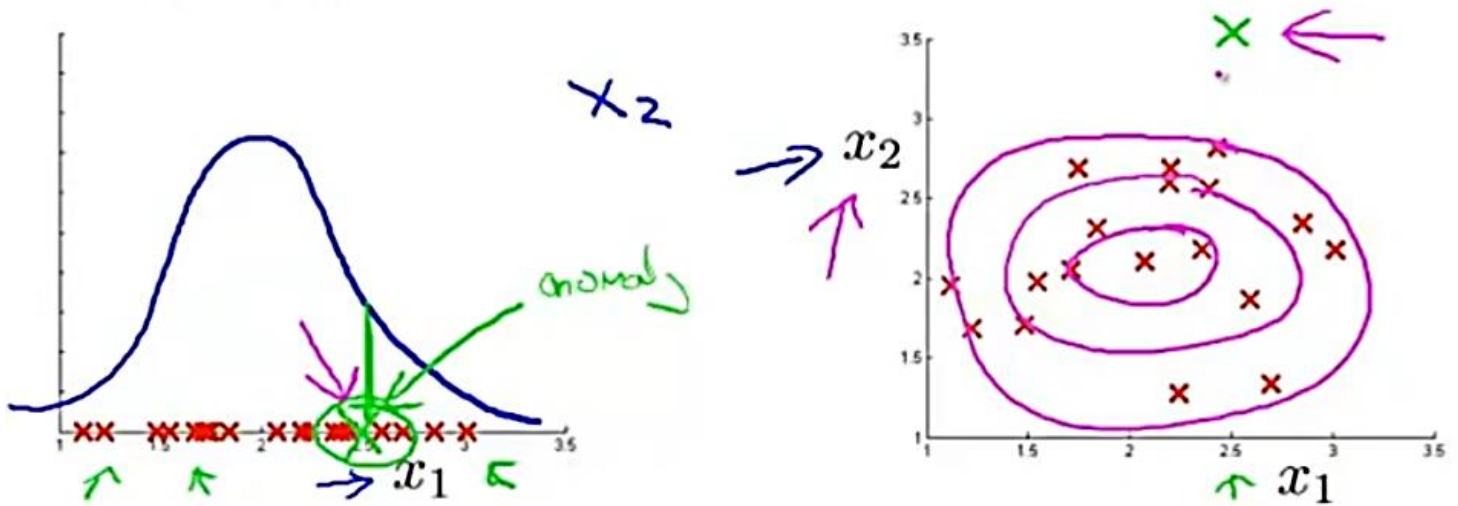
Most common problem:

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



If there is an anomalous example in middle of some non-anomalous examples, then the algo will fail.

- So, we can look at that particular example and try to come up with a new feature that can tell what went wrong with that example



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

Example:

Monitoring computers in a data center

x_1 = memory use of computer

x_2 = number of disk accesses/sec

x_3 = CPU load ←

x_4 = network traffic ←

It may occur that if one of the computers is stuck in an infinite loop, the CPU load grows but the network traffic doesn't:

Then we can come up with a new feature:

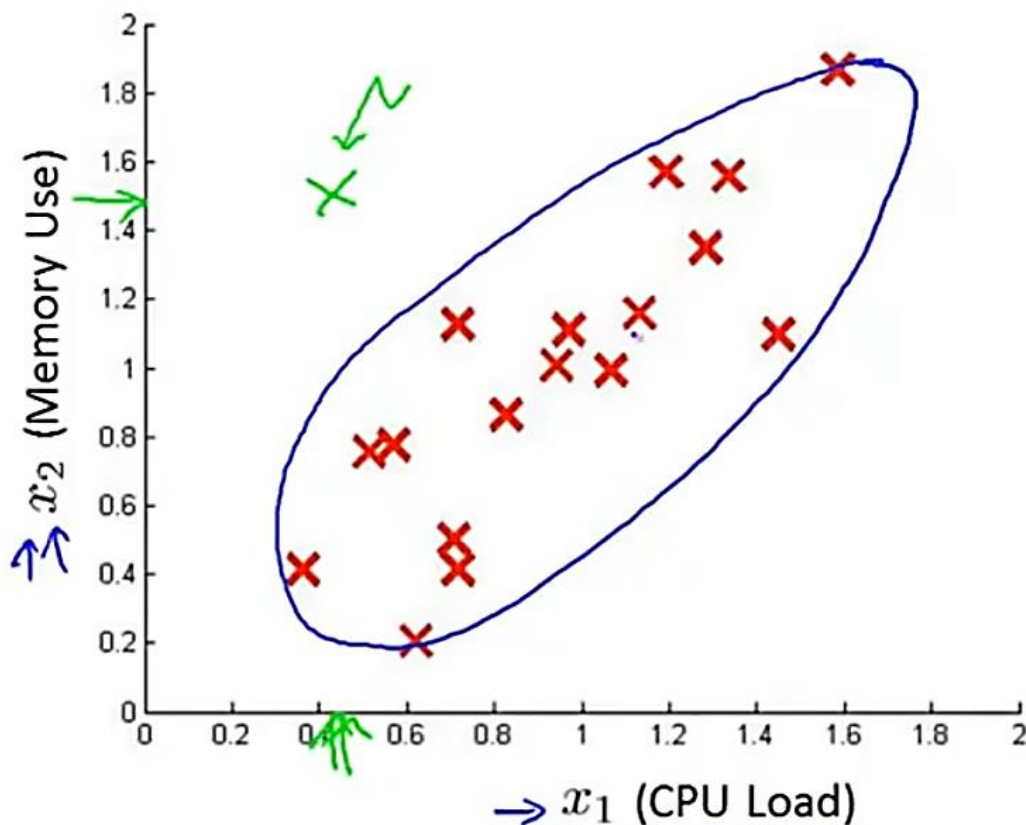
$$x_5 = \frac{\text{CPU load}}{\text{network traffic}}$$

OR

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

» MULTIVARIATE GAUSSIAN DISTRIBUTION:

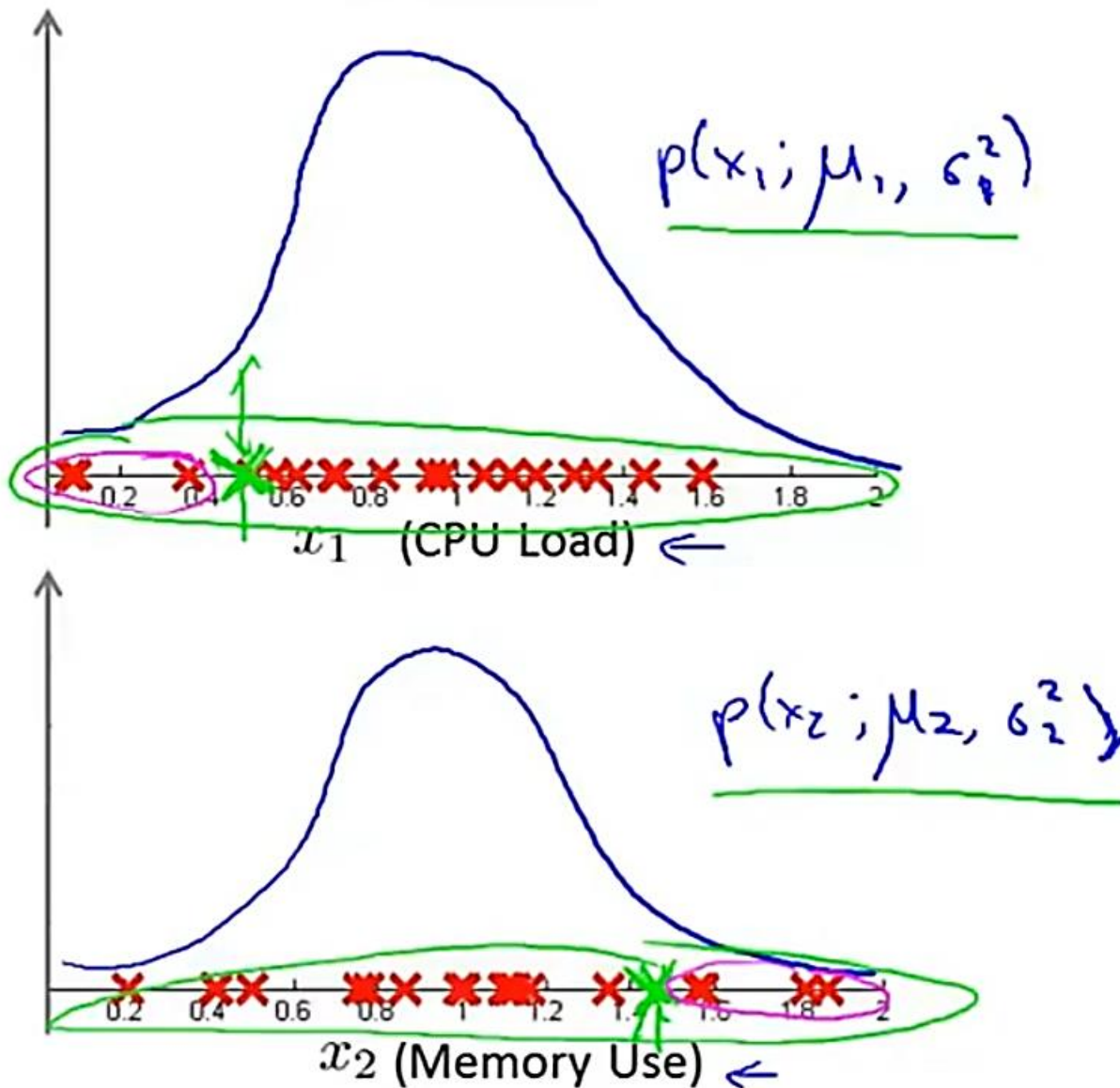
Motivating example: Monitoring machines in a data center



Here, red points are training data

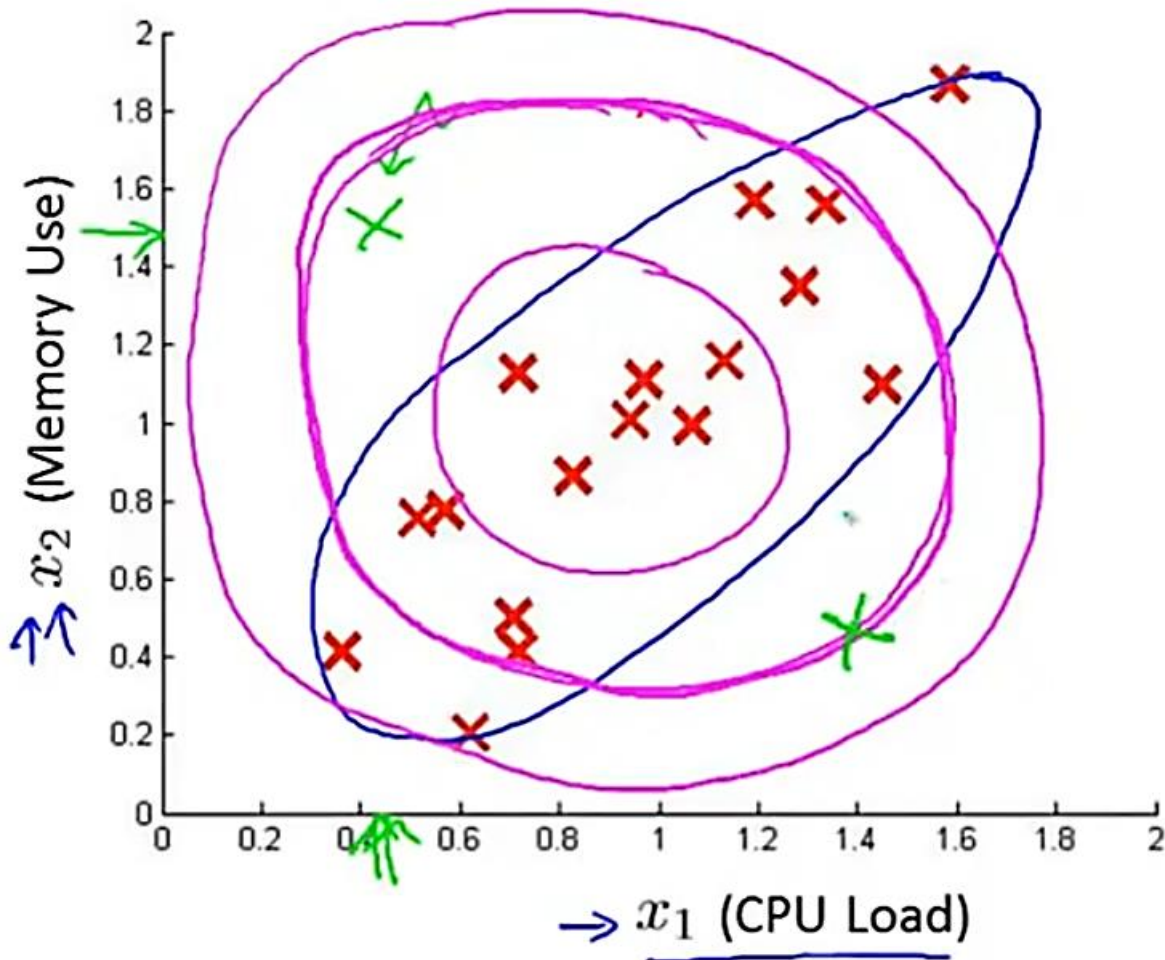
Green point is our test data

➤ Lets look at both features individually:



The algo will not predict the right o/p ... since the i/p data is distributed on the whole axis, so all the points have some probability of being correct.

Probability Contours:



To solve this:

Multivariate Gaussian (Normal) distribution

- $x \in \mathbb{R}^n$. 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)

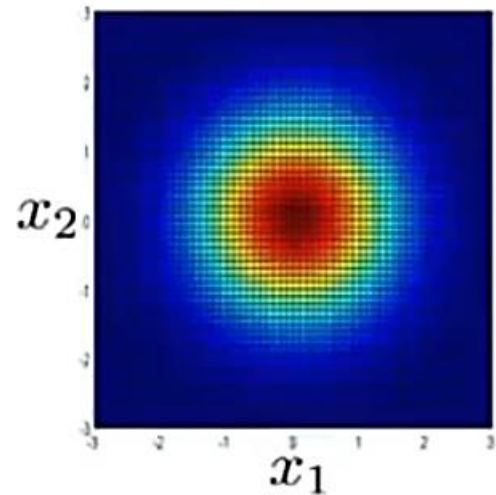
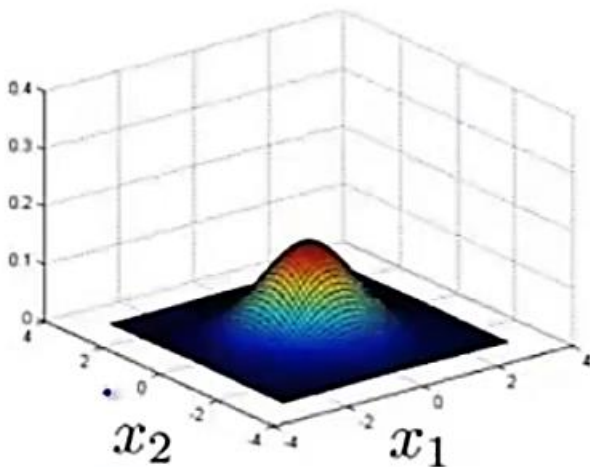
$$p(x; \mu, \Sigma) =$$

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

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

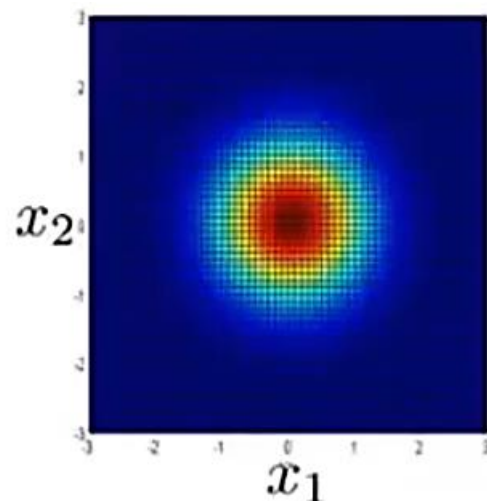
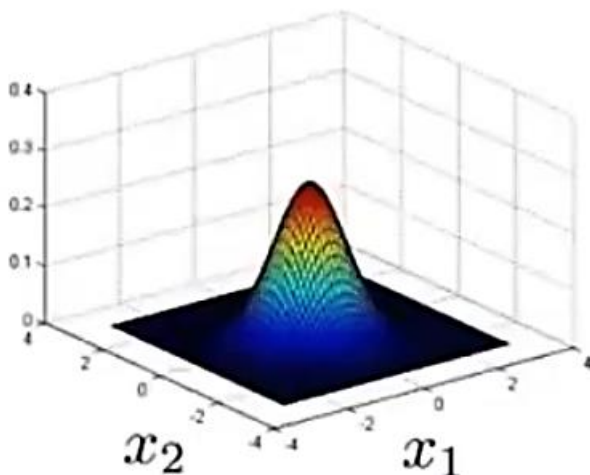
Examples:

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



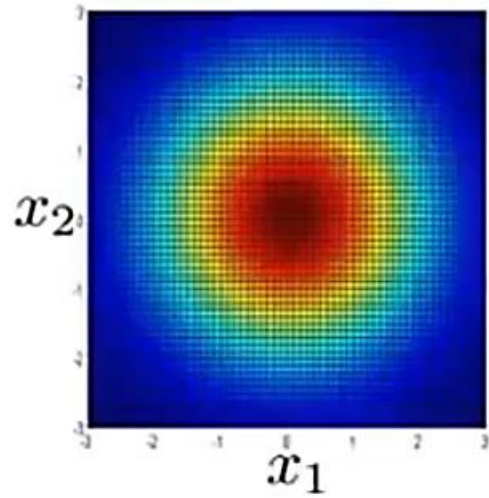
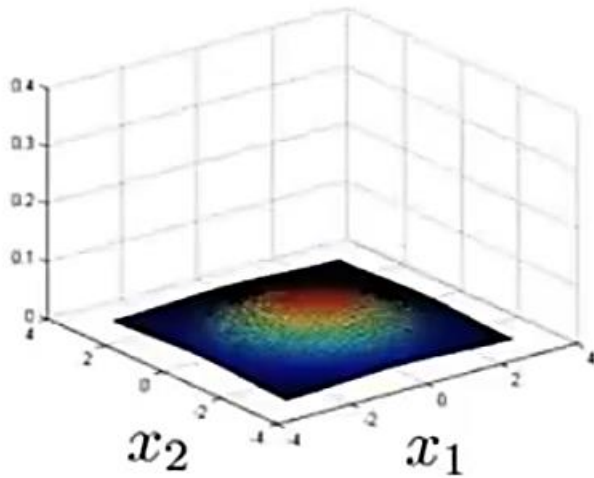
If we decr Σ :

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



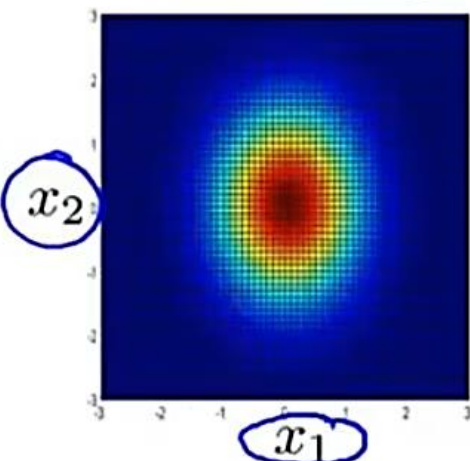
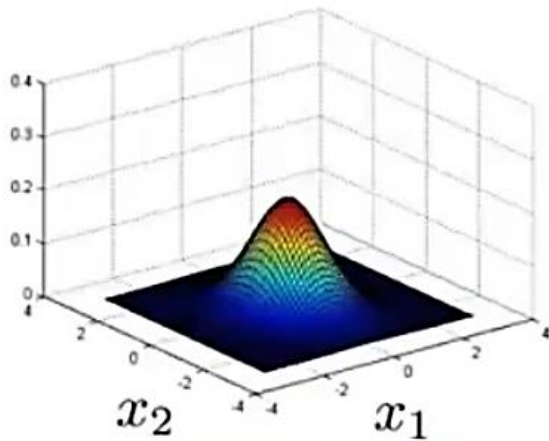
If we incr Σ :

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

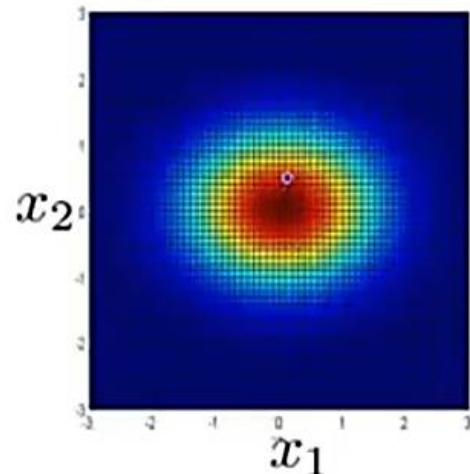
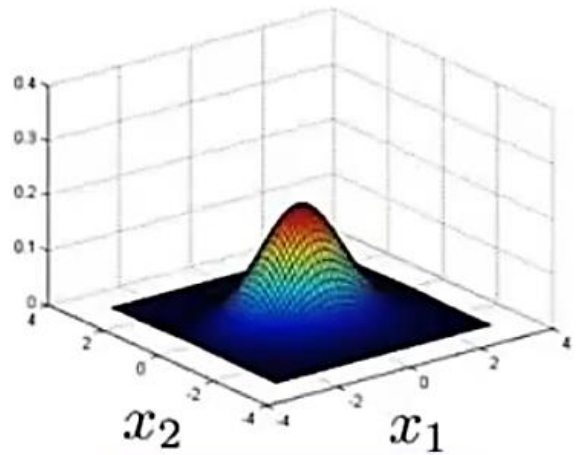


If we decr the variance of only one of the features:

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



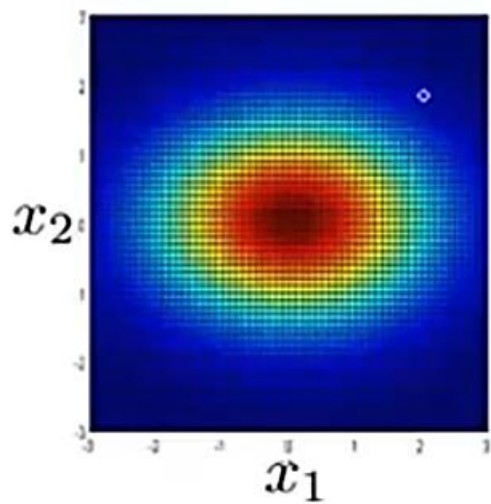
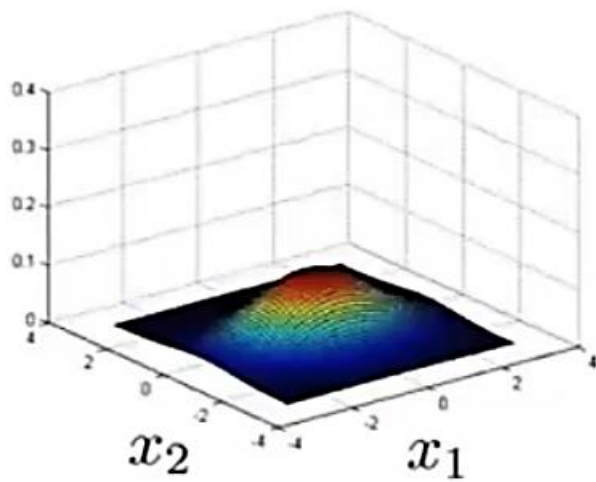
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 0.6 \end{bmatrix}$$



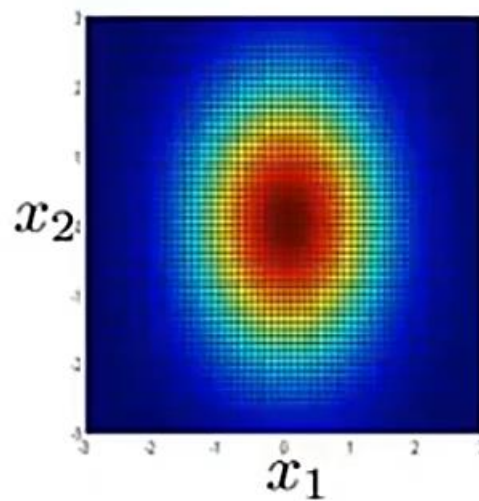
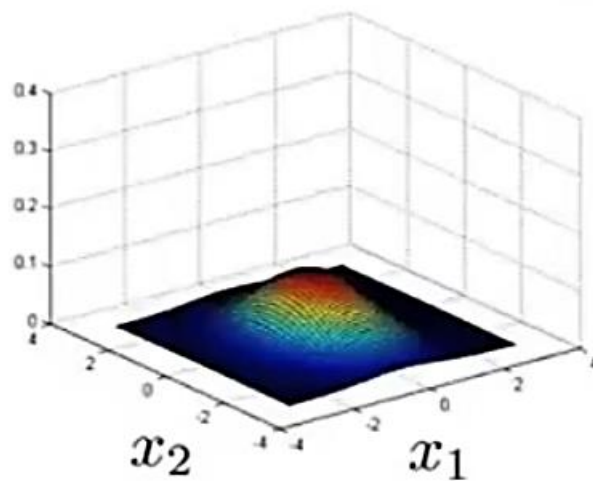
OR

If we incr the variance of only one of the features:

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



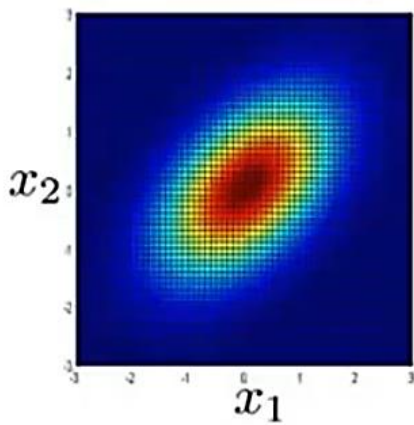
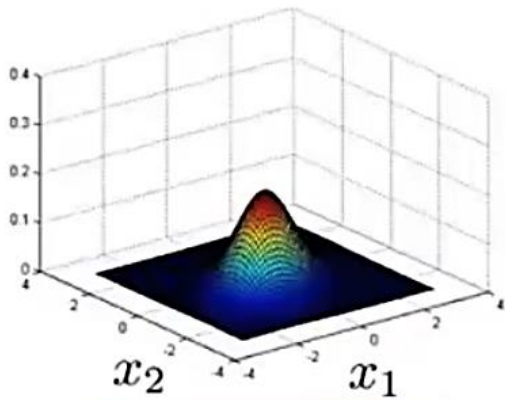
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$$



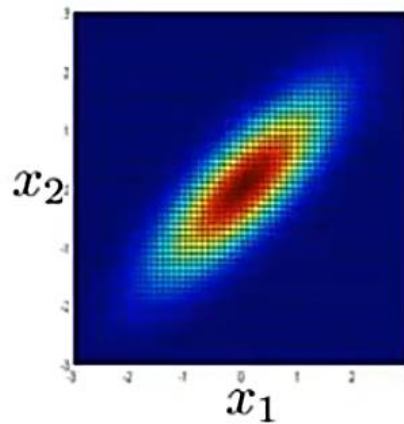
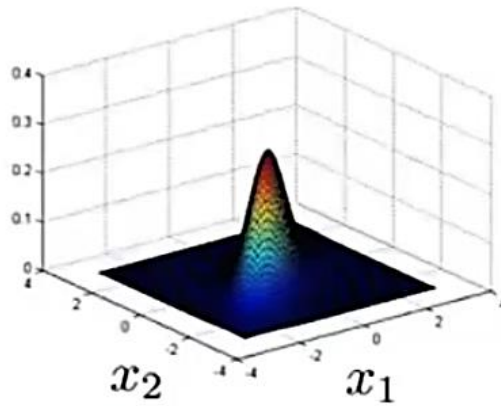
OR

Other variations in Σ :

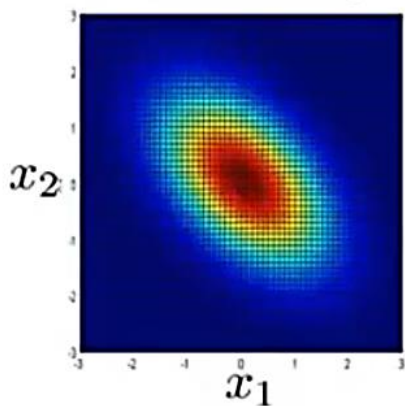
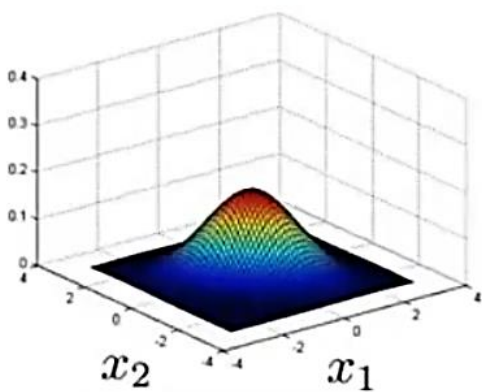
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$$



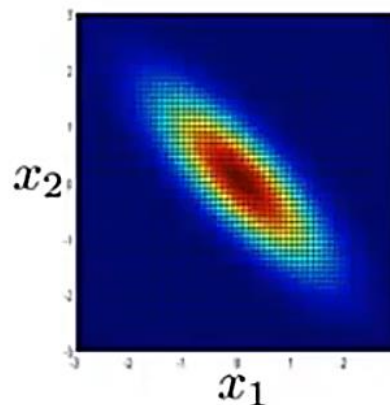
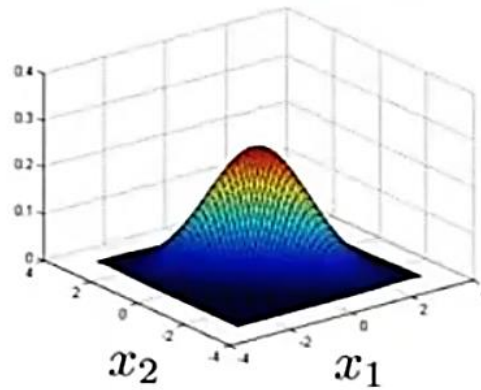
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$



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

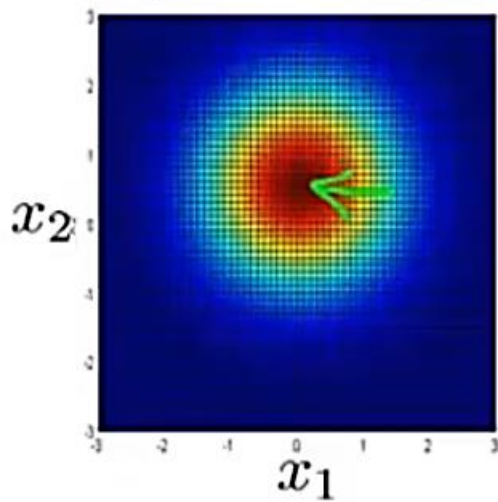
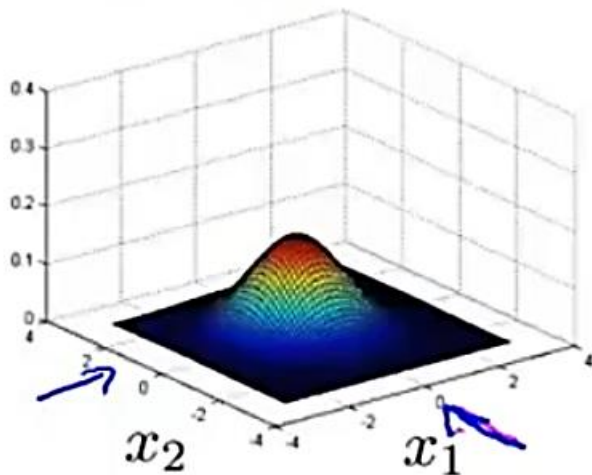


$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & -0.8 \\ -0.8 & 1 \end{bmatrix}$$

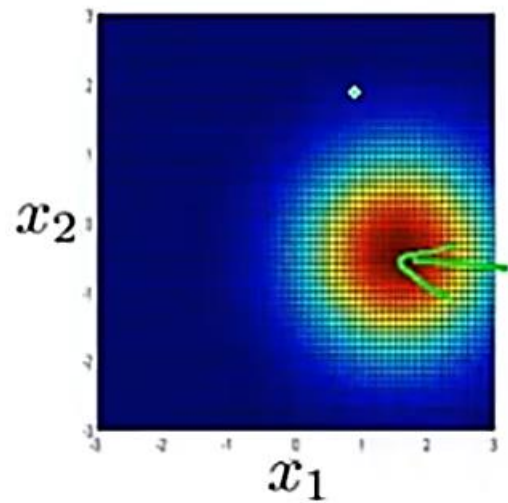
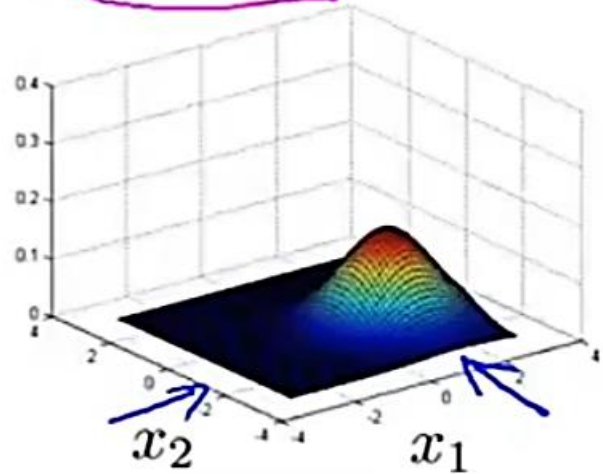


Varying the mean (μ): it moves the centre of contours, where the probability ($P(x)$) is highest.

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



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



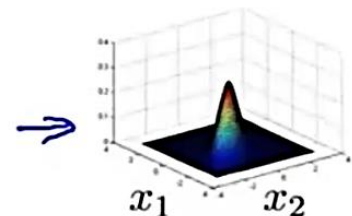
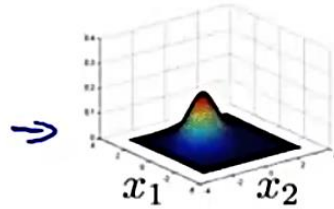
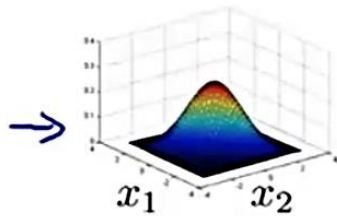
>> Multivariate Gaussian Distribution Algorithm:

Multivariate Gaussian (Normal) distribution

Parameters μ, Σ

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

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



Parameter fitting:

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

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

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

Algorithm:

Anomaly detection with the multivariate Gaussian

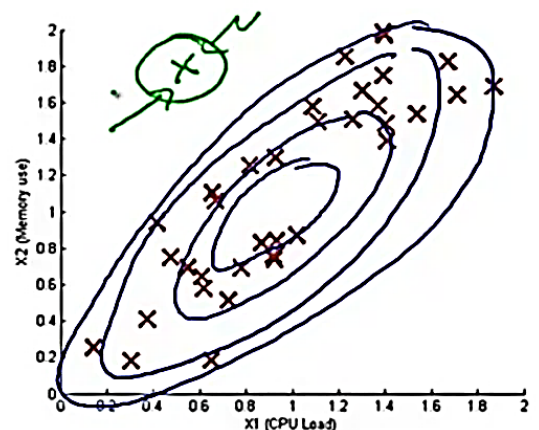
1. Fit model $p(x)$ by setting

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

2. Given a new example x , compute

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

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

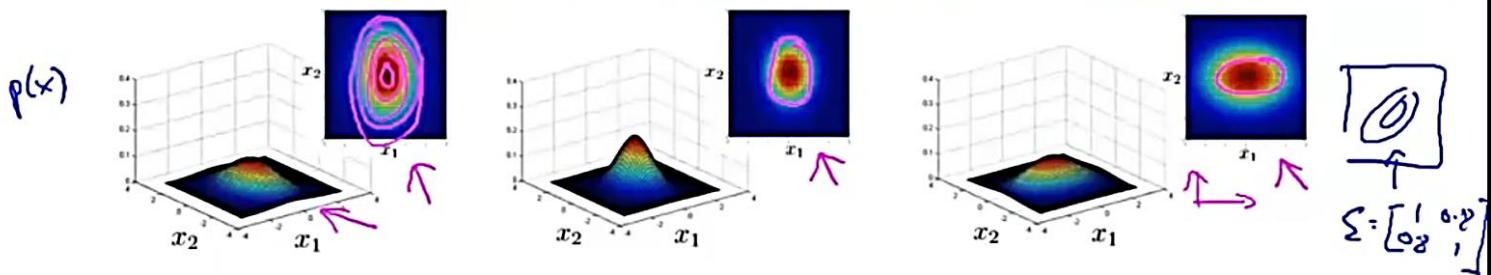


>> Relationship of multivariate Gaussian Model with Original Gaussian Model:

Original gaussian model is actually a special case of multivariate model, in which, the contours have their axes aligned with the features axes, i.e., the contours are not at nay angles:

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

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

Original model is mathematically the multivariate model with a constraint, that is:

$$\Sigma = \begin{bmatrix} \sigma_1^2 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_n^2 \end{bmatrix}$$

>> WHEN TO USE ORIGINAL MODEL vs MULTIVARIATE MODEL:

→ 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 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) $n=10,000, \quad n=100,000$

OK even if 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

$$\Sigma \in \mathbb{R}^{n \times n} \quad \Sigma^{-1}$$

Computationally more expensive

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

Must have $m > n$ or else Σ is non-invertible. $m \geq 10n$

- In some cases, in original model, we may require to manually create extra features, so that the model can work fine.
- In case of multivariate model, its important to get rid of redundant features, o/w the algo is very expensive, and Σ may even be non-invertible