

Week 9 - Lecture 1

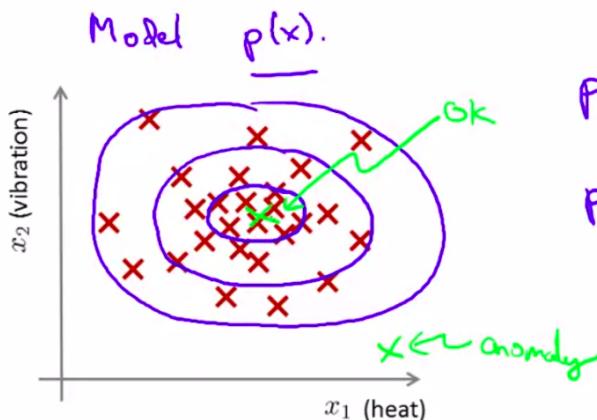
Density Estimation

Problem Motivation

Examples of anomaly detection.

Density estimation

- Dataset: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$
- Is x_{test} anomalous?



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

Andrew Ng

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.
 - ...

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

Andrew Ng

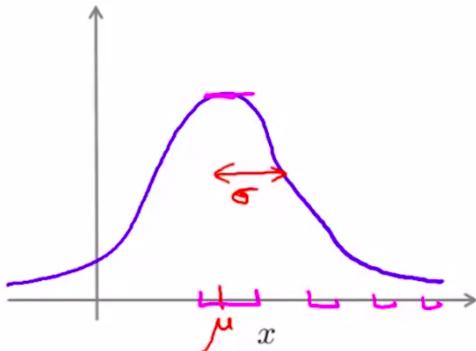
Gaussian Distribution

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"



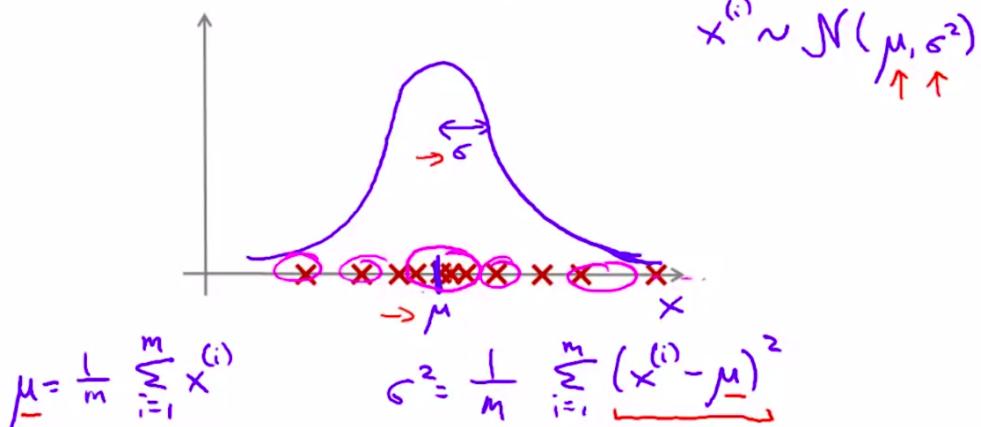
$$p(x; \mu, \sigma^2)$$

$$= \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

Andrew Ng

Parameter estimation

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



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

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

Andrew Ng

Note: in the figure above, the parameters formulas corresponds to the Maximum Likelihood Estimators for μ and σ^2 (you may remember that the variance estimator is biased. To make it unbiased, replace $1/m$ for $1/(m - 1)$). Professor Ng comments that in Machine Learning people tend to use the Max. Likelihood estimators as is pointed out in the figure and that, in practice, it shouldn't make much difference as long as m is sufficiently large.

Algorithm

Given a training set $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$, where each example $x \in \mathbb{R}^n$.

The way we are going to address anomaly detection is by modeling:

$$p(x) = p(x_1; \mu_1, \sigma_1^2) \cdot p(x_2; \mu_2, \sigma_2^2) \cdots p(x_n; \mu_n, \sigma_n^2) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2) \quad (1)$$

Note: we're assuming that $x_j \sim \mathcal{N}(\mu_j, \sigma_j)$, for all $j = 1, \dots, n$. We're also assuming that our features x 's are independent.

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 $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$

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

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

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

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

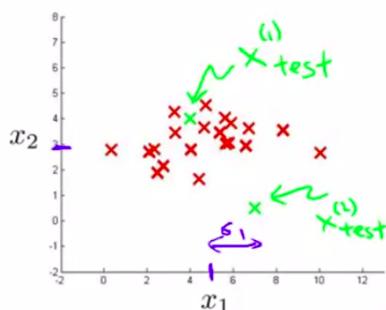
- 3. Given new example x , compute $p(x)$:

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

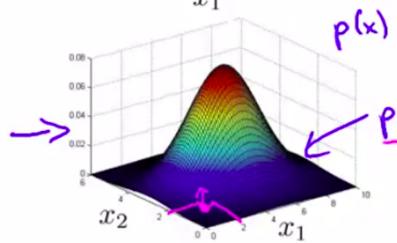
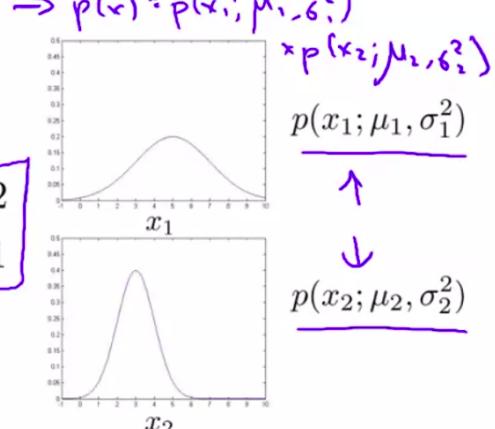
Anomaly if $p(x) < \varepsilon$

Andrew Ng

Anomaly detection example



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



$$\begin{aligned} \varepsilon &= 0.02 \\ p(x_{test}^{(1)}) &= 0.0426 > \varepsilon \\ p(x_{test}^{(2)}) &= 0.0021 < \varepsilon \end{aligned}$$

Andrew Ng

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. ($y = 0$ 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: $(x_{test}^{(1)}, y_{test}^{(1)}), \dots, (x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$

Andrew Ng

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 ($y=1$)

Alternative:

Training set: 6000 good engines

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

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

Andrew Ng

Note: commonly, this data will be very skewed (the number of anomalies is much less common than non-anomalies). So, the accuracy metric is not good for evaluating the model's performance in this case (predicting all $y = 0$ would have a good accuracy score).

Algorithm evaluation

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

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

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 ε

Andrew Ng

Anomaly Detection vs. Supervised Learning

Anomaly detection

- Very small number of positive examples ($y = 1$). (0-20 is common).
- Large number of negative ($y = 0$) examples. $p(x) \leftarrow$
- 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.

vs.

Supervised learning

Large number of positive and negative examples. \leftarrow

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.

Navigation icons

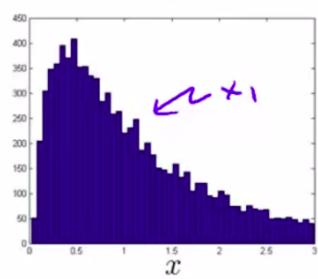
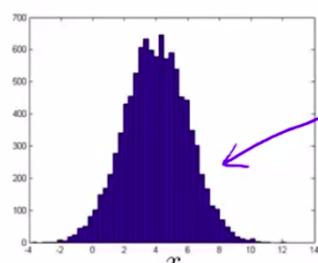
Andrew Ng

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

Andrew Ng

Choosing What Features to Use

Non-gaussian features



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

hist

$$x_1 \leftarrow \log(x_i)$$

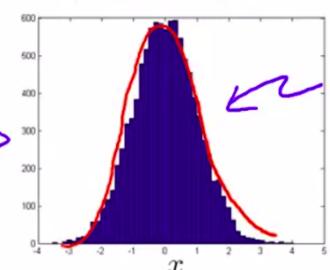
$$x_2 \leftarrow \log(x_1 + 1)$$

$$x_3 \leftarrow \sqrt{x_2} = x_2^{\frac{1}{2}}$$

$$x_4 \leftarrow \frac{1}{x_3^3}$$

$$\log(x_2 + \epsilon)$$

$$\log(x)$$



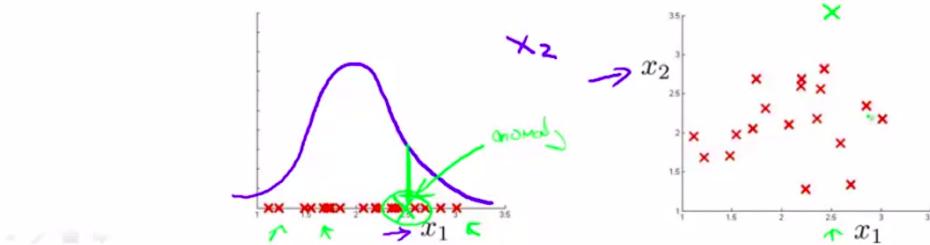
Note: in some cases where the features are non gaussian, the algorithm still works just fine (professor Ng words). However, if a feature doesn't look gaussian at all, we can try to find a transformation to make it look like gaussian.

→ 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



What we can do in this case is to look at this anomaly example and try to figure out new features which may help us to distinguish this "bad" example from the others (like the x_2 variable shown in the figure above. Hopefully, this "bad" example will take on a really unusual value for this new feature).

Example: to try to identify an anomaly in a Data Center, we could try to identify situations in which the CPU load is pretty high, but the network traffic is low (signalling that the CPU might have gotten stuck in a loop). To identify this kind of anomaly, we can come up with new features (x_5 and x_6) like shown in the next figure.

→ 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 ↘

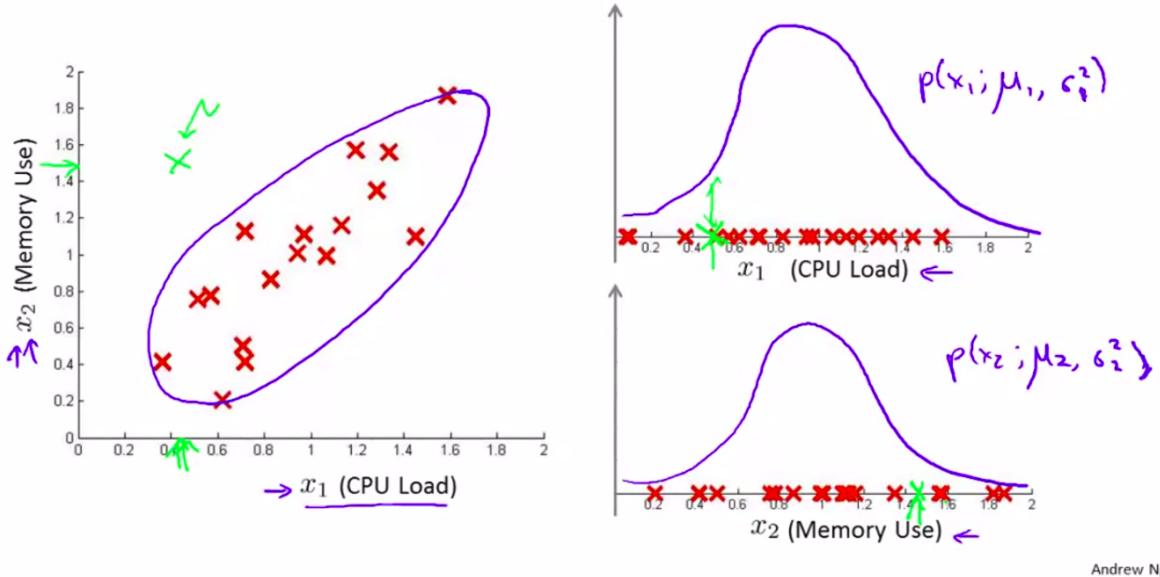
$$x_5 = \frac{\text{CPU load}}{\text{network traffic}}$$

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

Multivariate Gaussian Distribution

Now we'll start to make an extension to the anomaly detection algorithm we discussed so far.

Motivating example: Monitoring machines in a data center



In the example above, it is not easy to distinguish the anomaly (in green) from the rest of normal points. To address this, we are going to develop a Multivariate Gaussian Anomaly Detection Algorithm.

Note: in the previous algorithm, we treated each x_j as independent and identically distributed features and took the joint distribution of them by multiplying all $p(x_j; \mu_j, \sigma_j^2)$. In the multivariate Gaussian case, we're going to model all $p(x)$ vector entries in one go.

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)

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

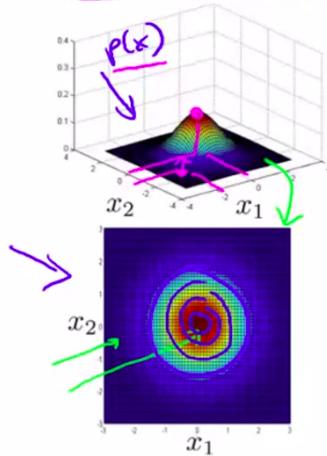
$$\frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp(-\frac{1}{2} (x - \mu)^\top \Sigma^{-1} (x - \mu))$$

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

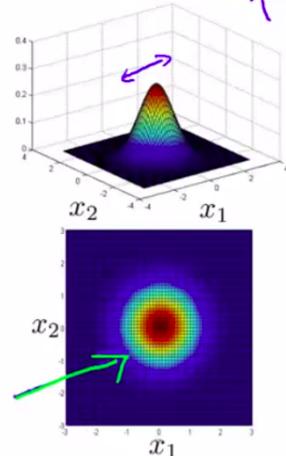
Andrew Ng

Multivariate Gaussian (Normal) examples

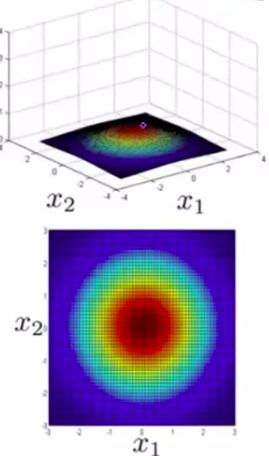
$$\rightarrow \mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



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



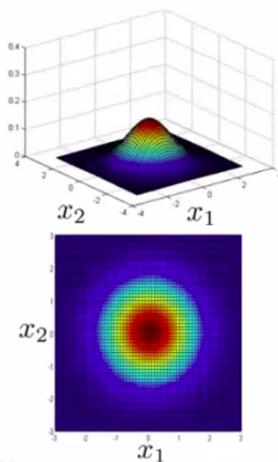
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$



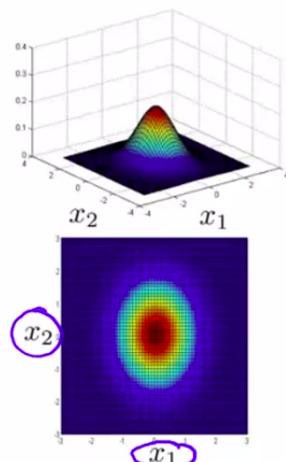
Andrew Ng

Multivariate Gaussian (Normal) examples

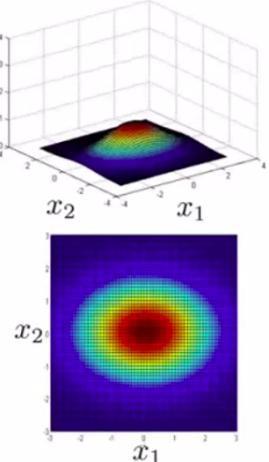
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



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



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

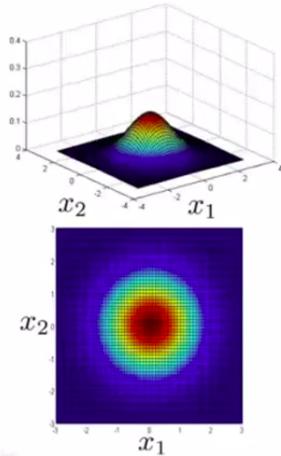


Andrew Ng

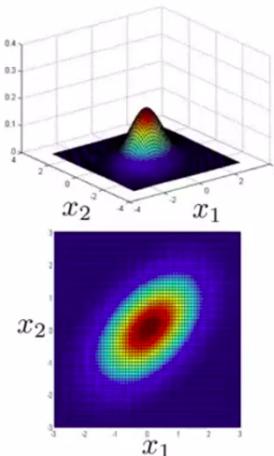
The Multivariate Gaussian Distribution is also good for modelling correlation between the variables:

Multivariate Gaussian (Normal) examples

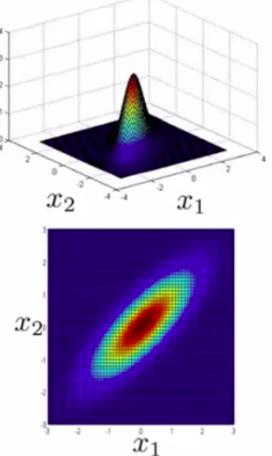
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



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



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



Andrew Ng

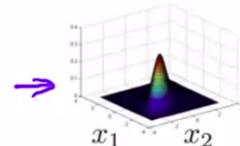
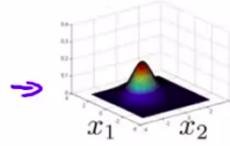
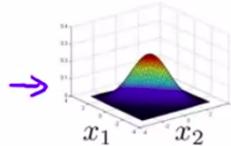
Anomaly Detection Using Multivariate Gaussian Distribution

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

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

Andrew Ng

Anomaly detection with the multivariate Gaussian

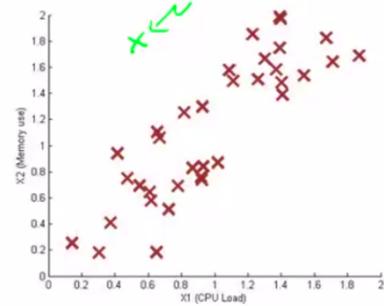
1. Fit model $p(x)$ by setting

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

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

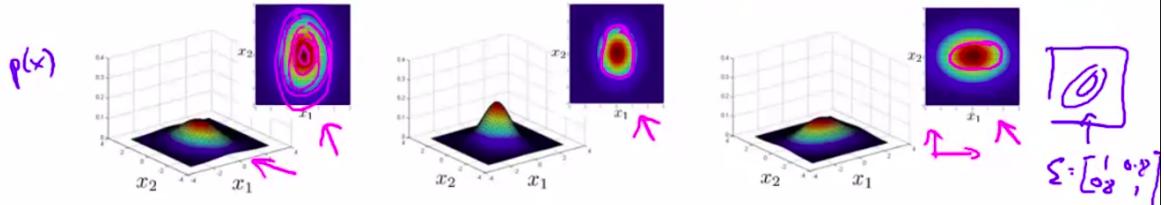


Andrew Ng

In fact, the Multivariate Gaussian model is a generalization of the previous model, where the latter corresponds to the former if the contours of the distribution are axis-aligned (i.e., the covariance matrix has all elements outside the main diagonal equal to 0).

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 $\Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \dots & \sigma_{1n} \\ \sigma_{21} & \sigma_2^2 & \dots & \sigma_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{n1} & \sigma_{n2} & \dots & \sigma_n^2 \end{bmatrix}$

Andrew Ng

But where do we use one model instead of the other?

→ 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) $n=10,000, n=100,000$

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

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

$$\Sigma^{-1}$$

Computationally more expensive

$$\Sigma \approx n \frac{n^2}{2}$$

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

$$m \geq 10n$$

Andrew Ng

References

[1] Machine Learning - Stanford University (<https://www.coursera.org/learn/machine-learning>)