

ANOMALY DETECTION

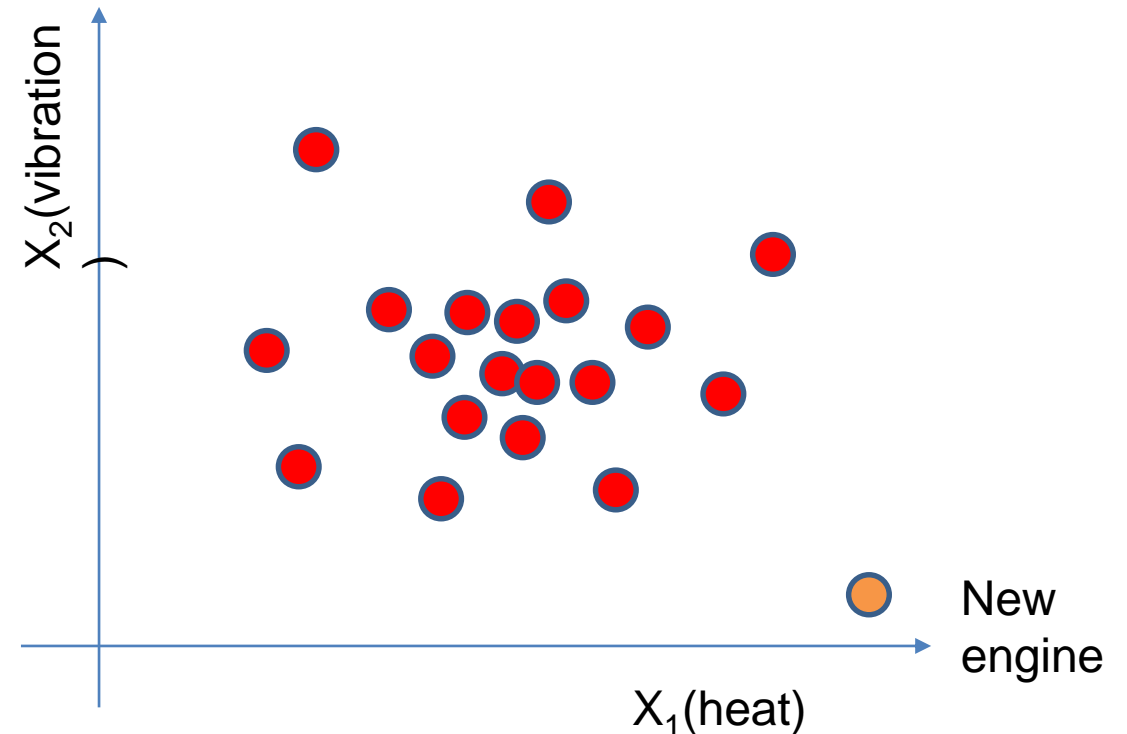
References

- Anomaly Detection – A Survey [PDF](#)
- A Comparative Evaluation of Unsupervised Anomaly Detection Algorithms for Multivariate Data [PDF](#)

What is Anomaly Detection?

- Example

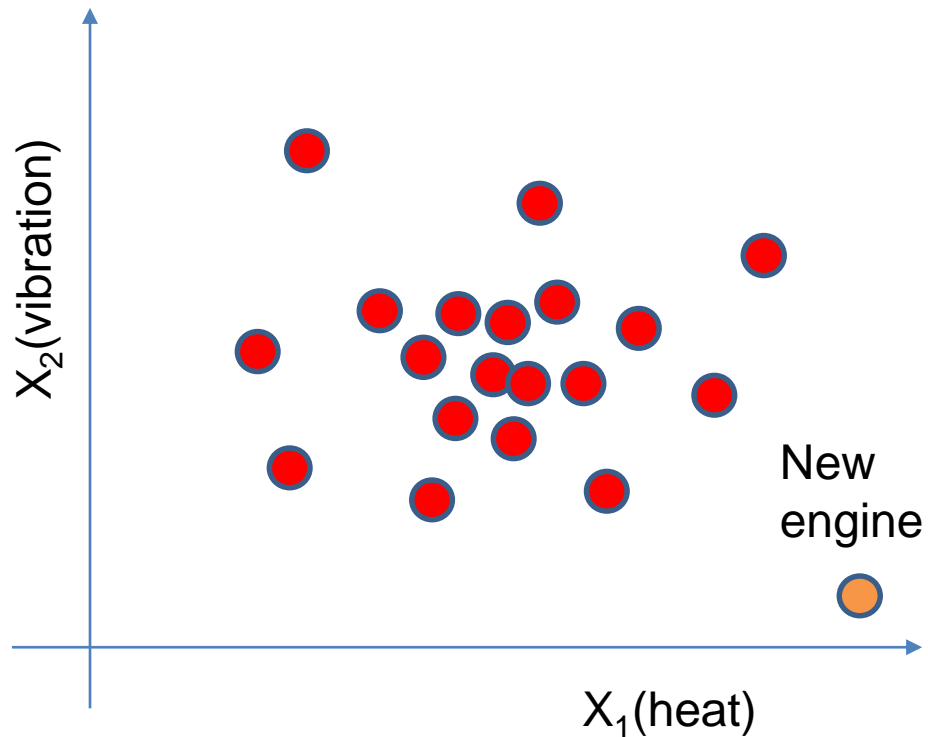
- Suppose you are a car manufacturer, and perform car engine testing.
- You are given some features
 - X_1 = heat generated
 - X_2 = vibration intensity
- Given DataSet: $x^{(1)}, x^{(2)}, \dots, x^{(n)}$
- Then for a new engine
 - We want to determine if x_{test} is an anomaly



What is Anomaly Detection

- Anomaly detection identifies unexpected items or events in datasets, which differ from the norm.
- Usually anomaly detection is applied to unlabeled data (so its not a standard classification task).

What is Anomaly Detection




- Given a set of n 'normal' data points, we build a model for $p(x)$.
- Then given a new data point, if
 - $p(x) < \varepsilon$ --- > anomaly
 - $p(x) \geq \varepsilon$ ----> normal

Other Applications of Anomaly Detection

- Fraud detection

- $X^{(i)}$ = features of a user's activities
- Model $p(x)$ from data
- Identify unusual users by checking which have $p(x) < \epsilon$



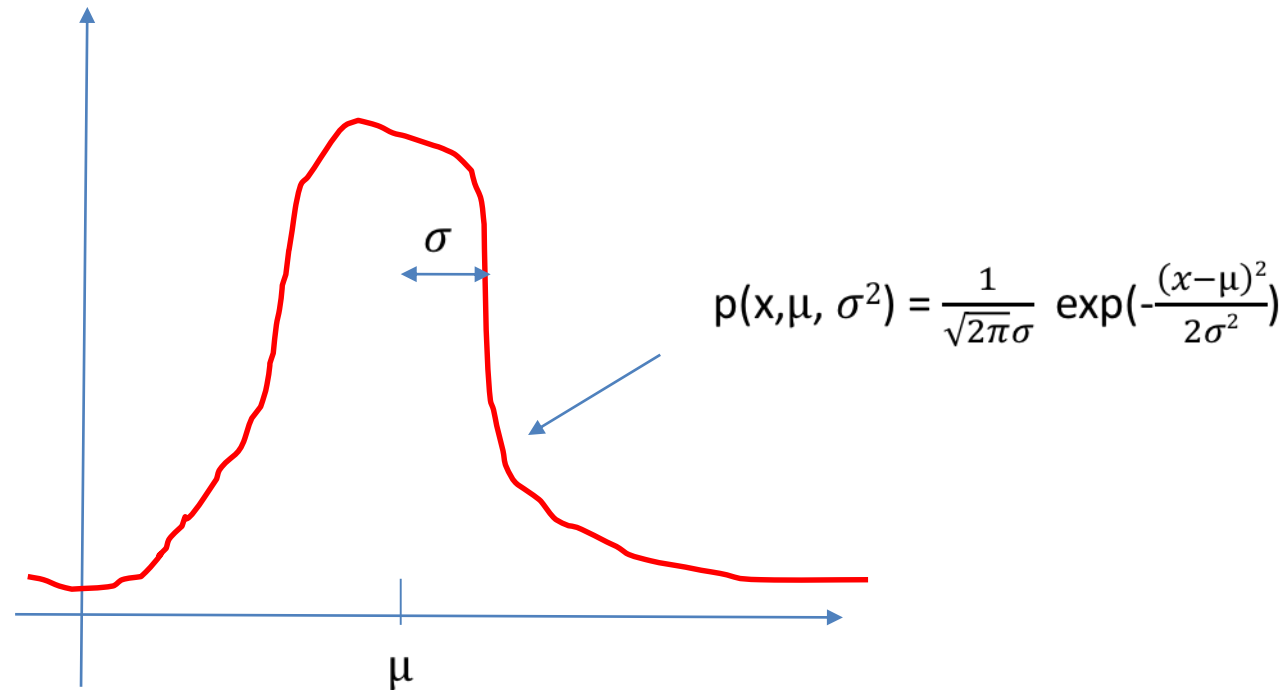
$X_1 \Rightarrow$ login frequency
 $X_2 \Rightarrow$ number of transactions
 $X_3 \Rightarrow$ typing speed

- Monitoring computers in a data center

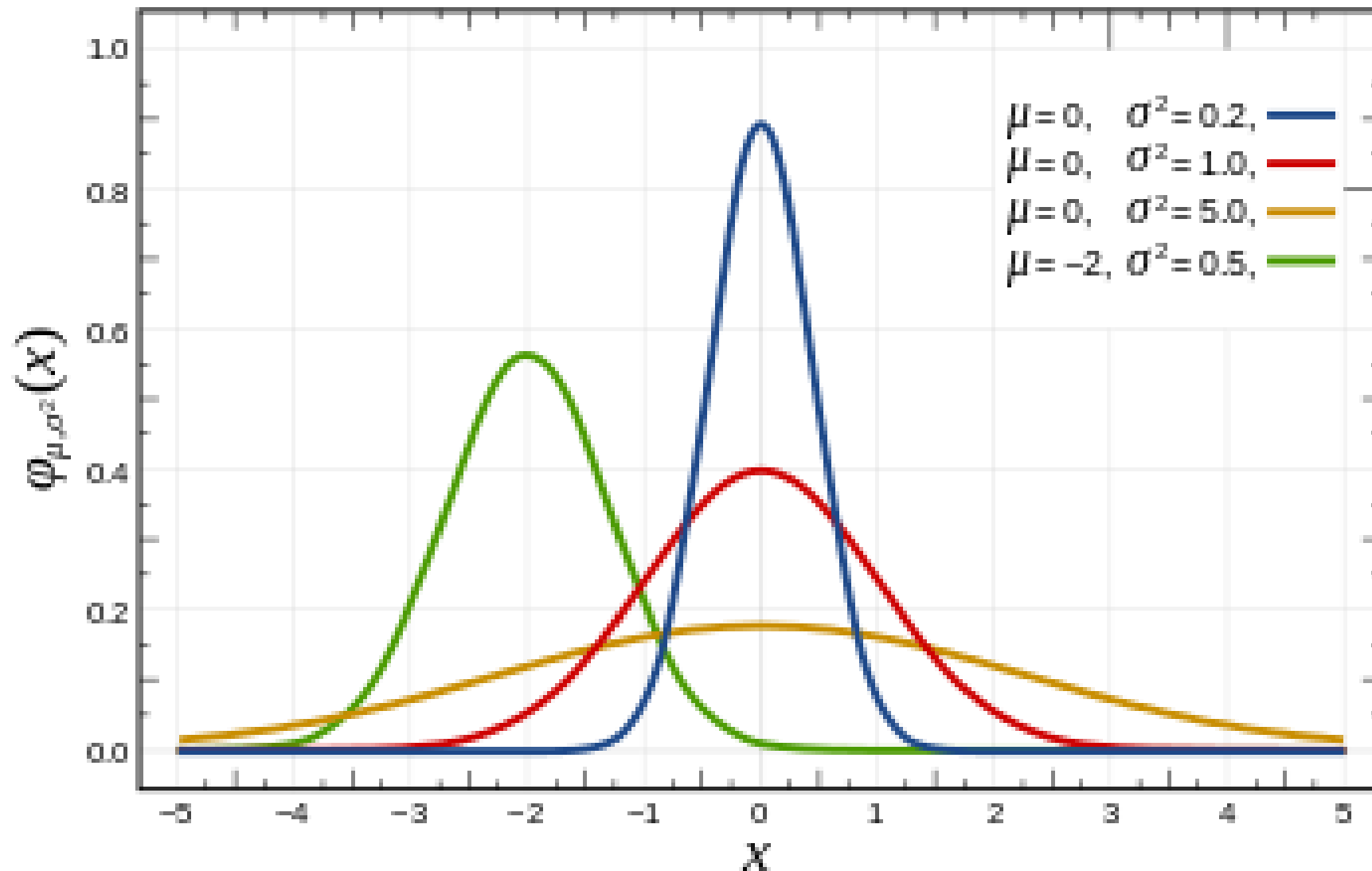
- $X^{(i)}$ = features of machine i
- Features :
 - x_1 = memory use,
 - x_2 = number of disk accesses,
 - x_3 = CPU load,
 - x_4 = CPU load / network traffic
- Model $p(x)$ and then identify abnormal machines that may indicate bad nodes in the network

Gaussian (Normal) Distribution

- Say $x \in \mathbb{R}$.
- If x is a distributed Gaussian with mean μ , *variance* σ^2 .
- $X \sim \mathcal{N}(\mu, \sigma^2)$

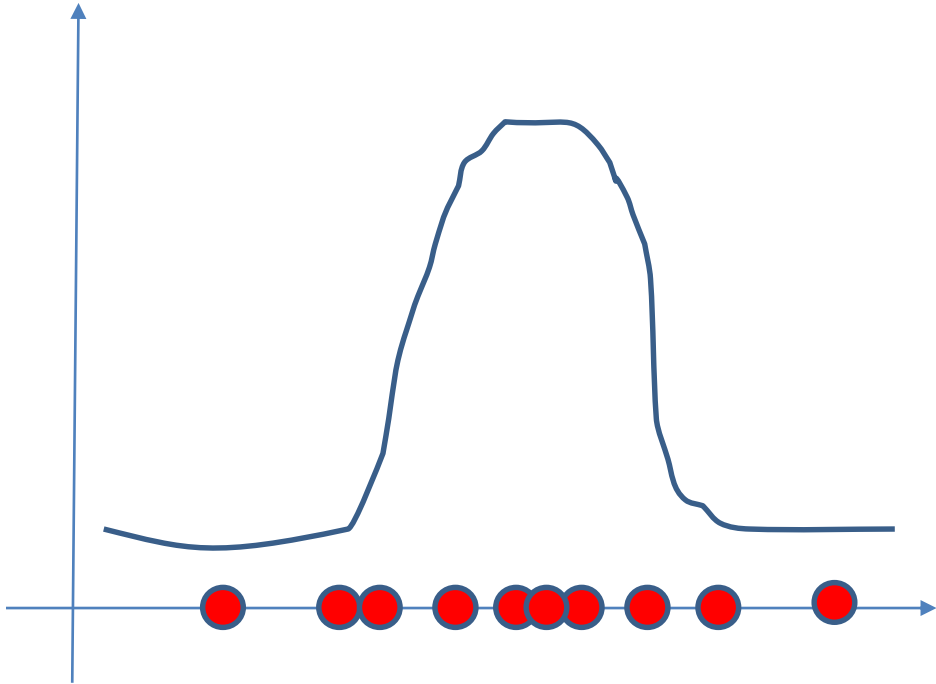


Gaussian Distribution Examples



Parameter Estimation

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



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

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

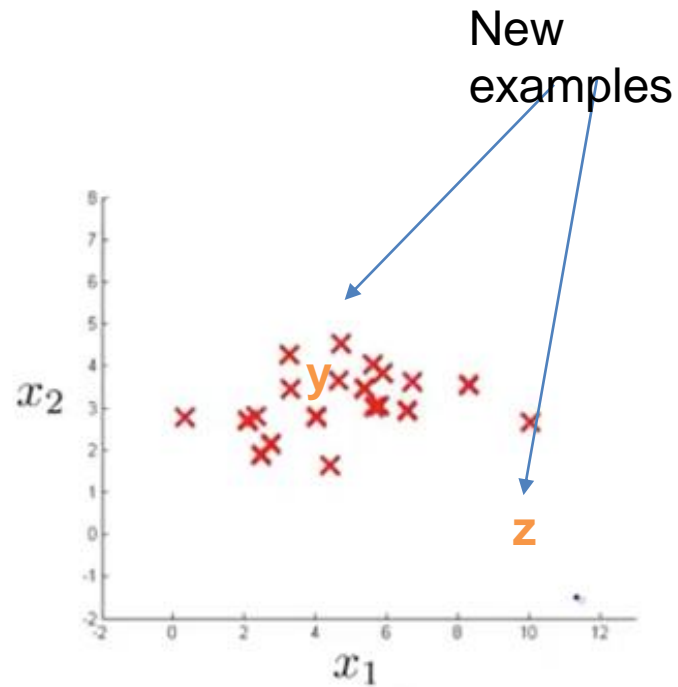
Density Estimation

- Given a Training Set: $\{x^{(1)}, \dots, x^{(m)}\}$
- Each example is $x \in \mathbb{R}^n$
- Model the probability of x
 - $p(x) = p(x_1; \mu_1, \sigma^2_1) p(x_2; \mu_2, \sigma^2_2) p(x_3; \mu_3, \sigma^2_3) \dots p(x_n; \mu_n, \sigma^2_n) = \prod_{j=1}^n p(x_j; \mu_j, \sigma^2_j)$
- How to model these terms?
 - Assume that each feature follows a Gaussian distribution with some mean and variation.
 - So, for example assume $x_1 \sim \mathcal{N}(\mu_1, \sigma^2_1)$

Anomaly Detection Algorithm

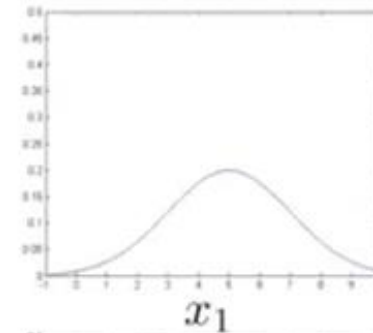
- Given a dataset of examples, $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$
- Choose features x_i that you think might be indicative of anomalous examples.
- Fit parameters $\mu_1, \dots, \mu_n; \sigma^2_1, \dots, \sigma^2_n$
 - $\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$
 - $\sigma^2_j = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$
- Given new example x , compute $p(x)$:
 - $p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma^2_j) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi}\sigma_j} \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right)$
 - Anomaly if $p(x) < \epsilon$

Example

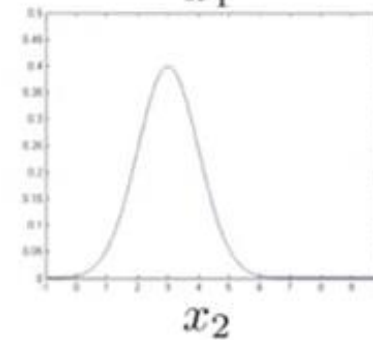


$$\mu_1 = 5, \sigma_1 = 2$$

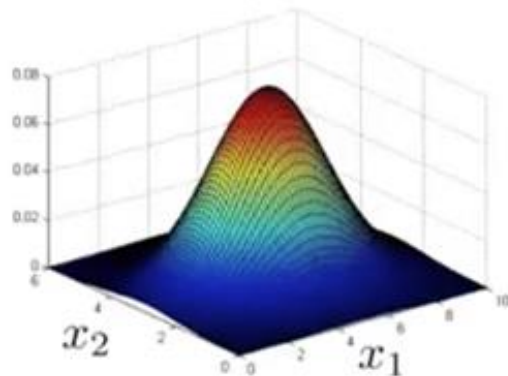
$$\mu_2 = 3, \sigma_2 = 1$$



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



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



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

$$\varepsilon = 0.02$$

$$p(y) = 0.0426 > \varepsilon$$

$$p(z) = 0.0021 < \varepsilon$$

How to evaluate the algorithm?

- We have been treating anomaly detection as an unsupervised learning problem, using non-labeled dataset.
- However, if we have some labeled data, of anomalous and non-anomalous examples ($y=0$ if normal, $y=1$ if anomalous), then we can use this information to make decisions about our learning algorithm.
 - Such as which features to choose
- Training set: $x^{(1)}, x^{(2)}, \dots, x^{(n)}$ assume normal examples / not anomalous
- Cross Validation set : normal and anomalous examples
- Test Set : normal and anomalous examples

Example (1/2)

- Dataset

- 1000 good (normal) examples
- 20 flawed examples (anomalous)

- Splitting Dataset

$$y=0 \quad p(x) = p(x_1, \mu_1, \sigma_1^2) * p(x_2, \mu_2, \sigma_2^2) \dots$$


- **Training set (60%)** : 6000 good examples (unlabeled training set)
- **Cross validation set(20%)** : 2000 good examples ($y=0$) & 10 anomalous ($y=1$)
- **Test set(20%)** : 2000 good examples ($y=0$) & 10 anomalous ($y=1$)

Example (2/2)

- Fit model $p(x)$ on training set $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$
- On a cross validation / test example x , predict

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

- Possible evaluation metrics:
 - True positive, false positive, false negative, true negative
 - Precision / recall
 - F_1 -recall
- Can you also use cross validation set to choose parameter ε

Anomaly Detection VS Supervised Learning

- Very small number of positive examples ($y=1$) maybe 0 – 20 anomalous examples.
- Large number of negative ($y=0$) examples
- Many different “types” of anomalies. Hard for any algorithm to learn from positive examples what future anomalies may look like
- 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.

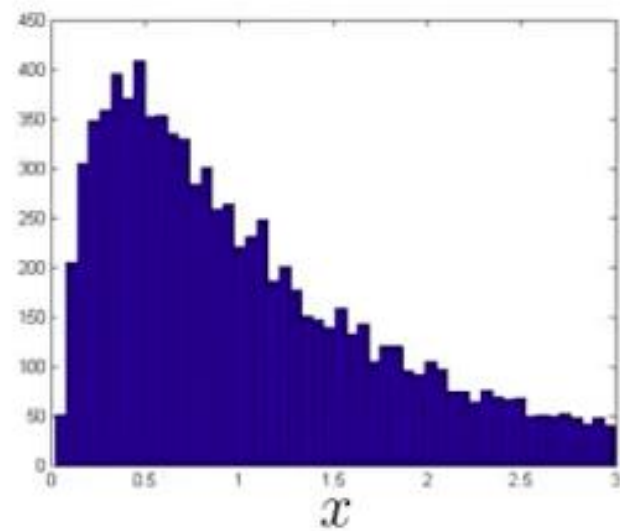
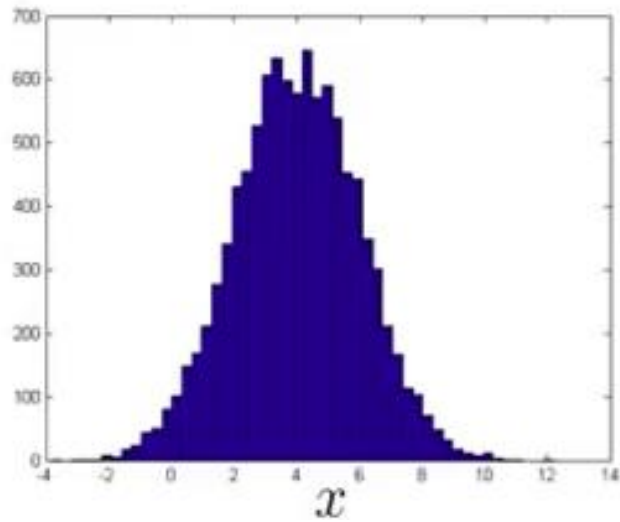
Anomaly Detection

- Fraud Detection
- Manufacturing
- Monitoring machines in a data center

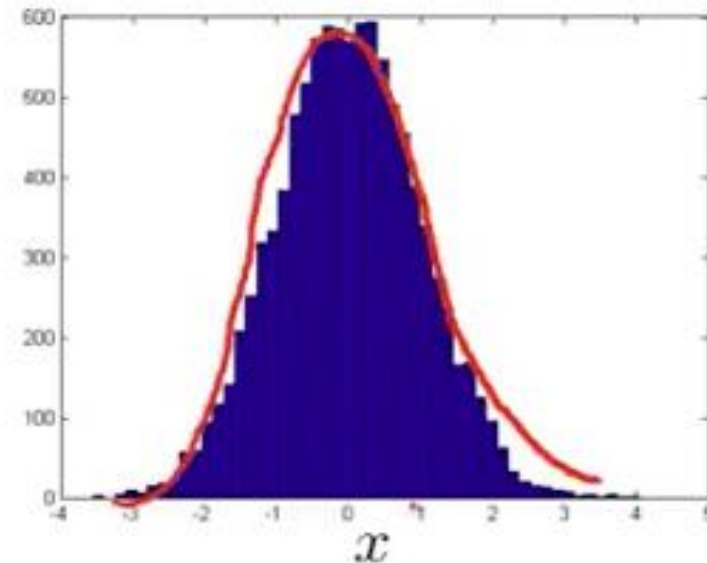
Supervised Learning

- Email Spam Classification
- Weather prediction
- Cancer classification

Non-gaussian Features



$\log(x)$



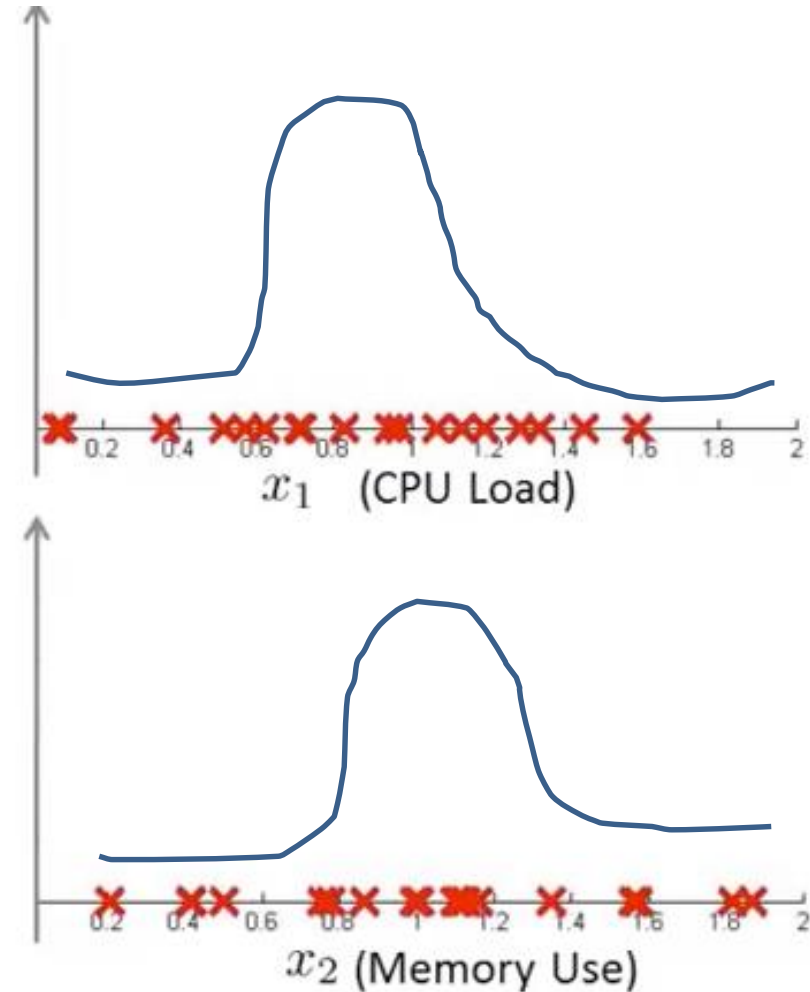
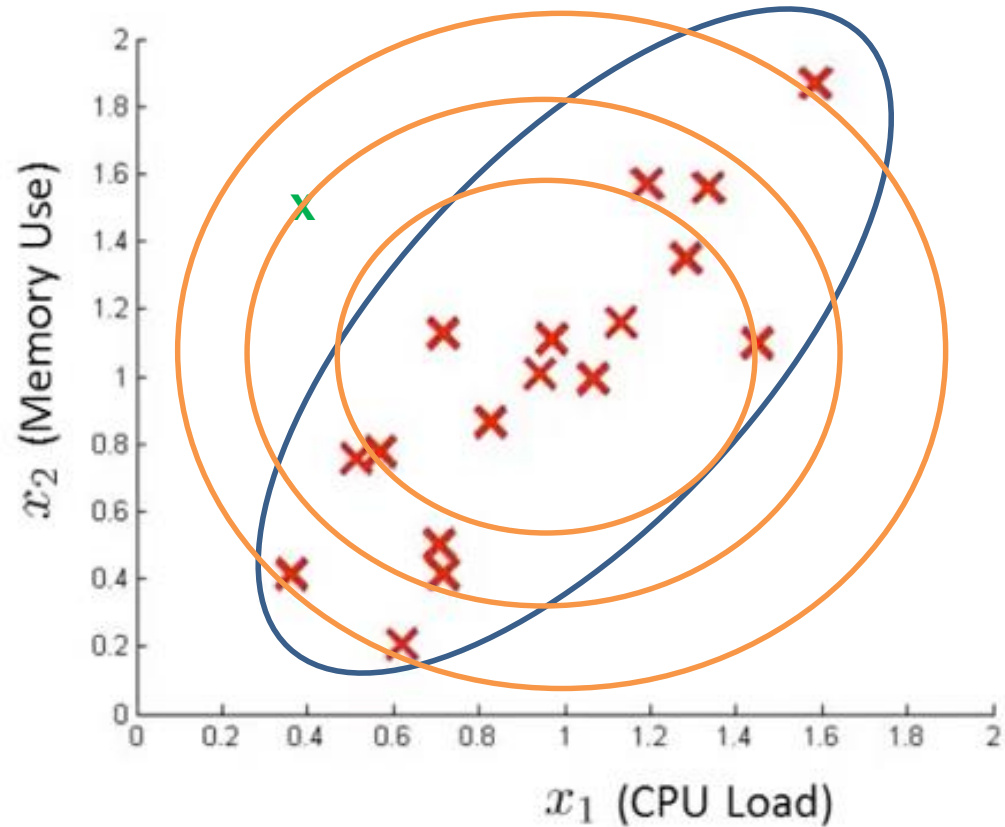
Error analysis for anomaly detection

- 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
- Solution
 - Find new features that expose the difference between anomalous and normal examples

Deriving features

- Consider the following example
 - Lets assume we want to detect anomalies in a data center
 - We are given a set of features
 - X_1 = memory use of computer
 - X_2 = number of disk accesses / sec
 - X_3 = CPU load
 - X_4 – network traffic
 - The goal is to choose features that might take on unusually large or small values in the event of an anomaly.
 - May derive features such as
 - x_5 = cpu load / network traffic

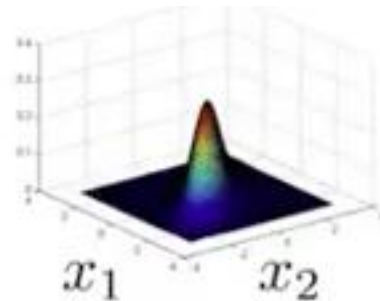
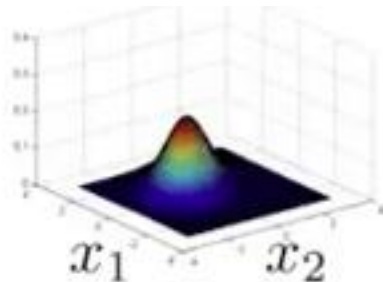
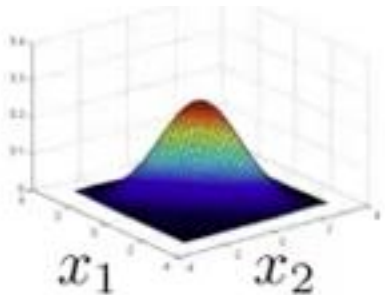
Multivariate Gaussian : Motivating Example



Multivariate Gaussian (Normal) Distribution

- $X \in \mathbb{R}^n$, don't model $p(x_1)$, $p(x_2)$, ..., $p(x_n)$. separately.
- Model $p(x)$ collectively.
- Parameters : $\mu \in \mathbb{R}^n$, $\Sigma \in \mathbb{R}^{n \times n}$ (*covariance matrix*)

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



We get different distributions as we vary these parameters

Parameter fitting

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

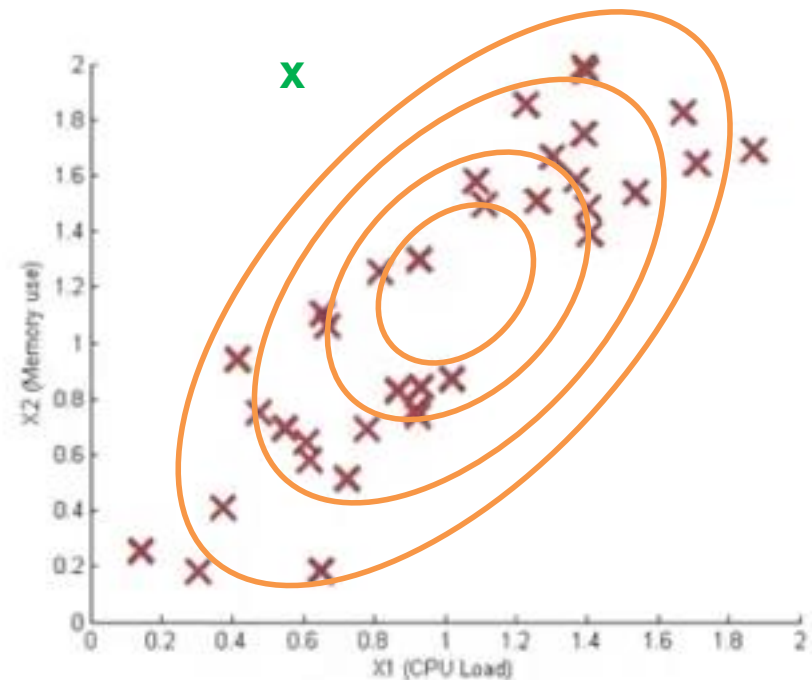
- Fit model $p(x)$ by setting

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

- 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 $p(x) < \varepsilon$



Original Model

- $p(x_1; \mu_1; \sigma^2_1) \times p(x_2; \mu_2; \sigma^2_2) \times \dots \times p(x_n; \mu_n; \sigma^2_n)$
- Manually create features to capture anomalies where x_1, x_2 take unusual combinations of values
- Computationally cheaper (alternatively, scales better to large n)
- Okay even if m (training set size) is small

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 Σ is non invertible