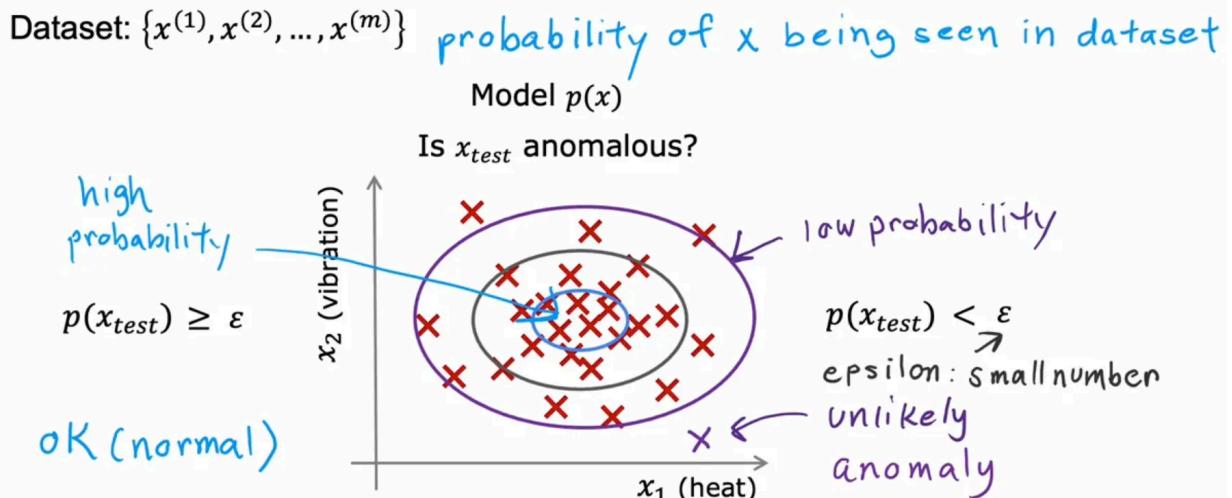


# Unsupervised Learning, Recommenders, Reinforcement Learning

## Anomaly Detection

Anomaly detection is a type of unsupervised learning. It identifies unusual data points in a dataset that consists mostly of normal examples. The algorithm learns the pattern of normal data and flags any example that may be an anomaly.

### Density estimation



## 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  $p(x) < \epsilon$

how often log in?

how many web pages visited?  
transactions?

posts? typing speed?

perform additional checks to identify real fraud vs. false alarms

Manufacturing:

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

airplane engine  
circuit board  
smartphone

ratios

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.

## Gaussian Distribution

For anomaly detection, we use the Gaussian distribution, which is defined by two parameters: the mean ( $\mu$ ) and the variance ( $\sigma^2$ ). The mean determines the center of the distribution, while the standard deviation ( $\sigma$ ) controls the width of the curve.

If a data point has a high probability density under this curve, it is considered normal; otherwise, it is considered an anomaly.

## Gaussian (Normal) distribution

Say  $x$  is a number.

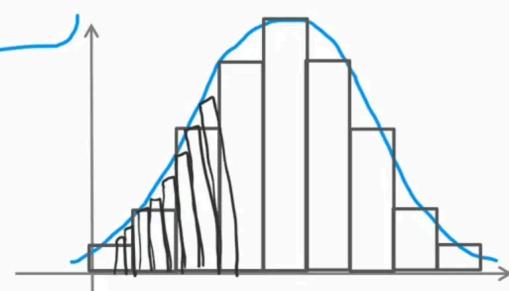
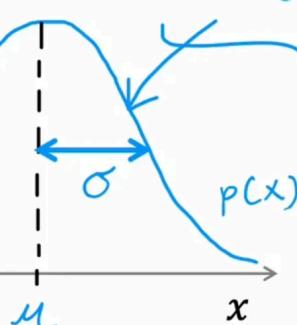
Probability of  $x$  is determined by a Gaussian with mean  $\mu$ , variance  $\sigma^2$ .

$\sigma$  standard deviation  
 $\sigma^2$  variance



$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$\pi = 3.14$$

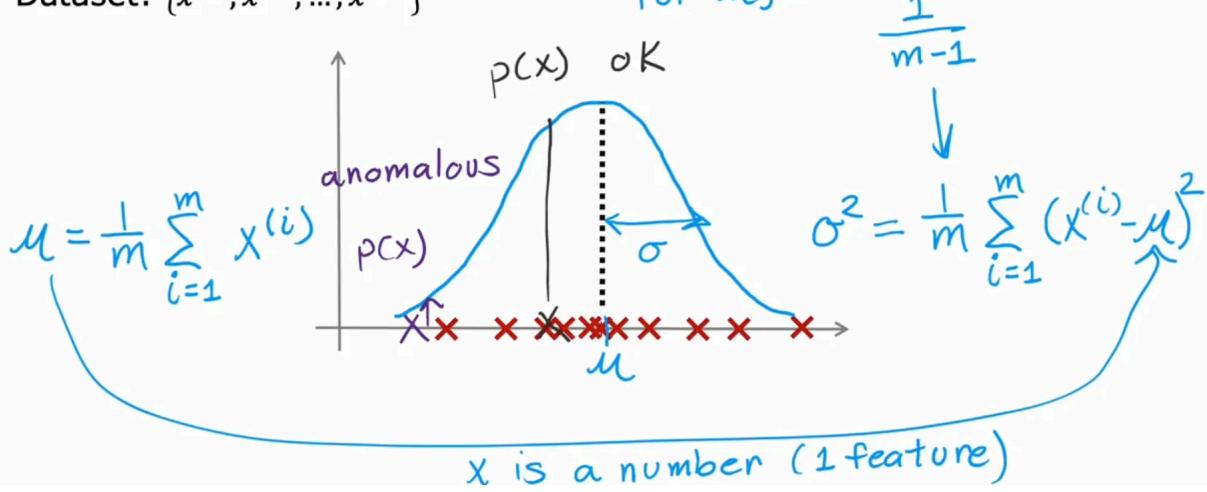


## Parameter estimation

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

maximum likelihood

for  $\mu, \sigma$



Now, how do we apply the Gaussian distribution to datasets with multiple features?

For each feature, we calculate its mean ( $\mu$ ) and variance ( $\sigma^2$ ). The total probability of an example  $x$ , denoted as  $p(x)$ , is computed as the product of the individual probabilities of each feature. If this total probability falls below a threshold  $\epsilon$ , the example is flagged as an anomaly.

Since the total probability is a product of individual probabilities, even a single feature with a low probability can significantly reduce  $p(x)$ , allowing the model to detect anomalies effectively.

## Density estimation

Training set:  $\{\vec{x}^{(1)}, \vec{x}^{(2)}, \dots, \vec{x}^{(m)}\}$

Each example  $\vec{x}^{(i)}$  has  $n$  features

$$\vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$p(\vec{x}) = p(x_1; \mu_1, \sigma_1^2) * p(x_2; \mu_2, \sigma_2^2) * p(x_3; \mu_3, \sigma_3^2) * \dots * p(x_n; \mu_n, \sigma_n^2)$$

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

"add" "multiply"

$$p(x_1 = \text{high temp}) = 1/10$$

$$p(x_2 = \text{high vibra}) = 1/20$$

$$p(x_1, x_2) = p(x_1) * p(x_2)$$

$$= \frac{1}{10} * \frac{1}{20} = \frac{1}{200}$$

## Anomaly detection algorithm

1. Choose  $n$  features  $x_i$  that you think might be indicative of anomalous examples.

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

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

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

Vectorized formula

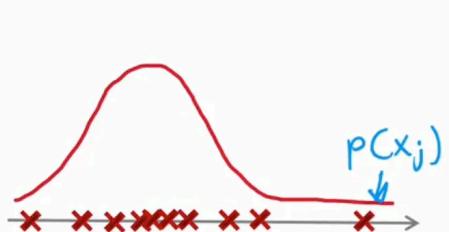
$$\vec{\mu} = \frac{1}{m} \sum_{i=1}^m \vec{x}^{(i)}$$

$$\vec{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix}$$

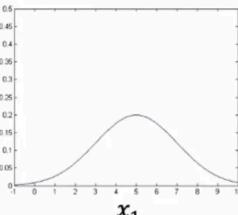
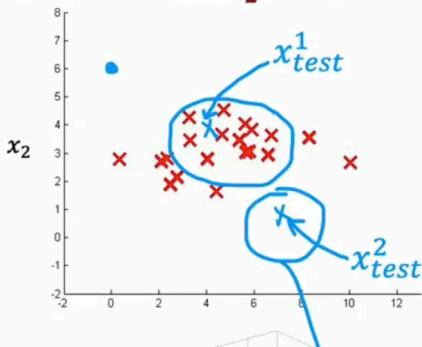
3. Given new example  $x$ , compute  $p(x)$ :

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

Anomaly if  $p(x) < \varepsilon$

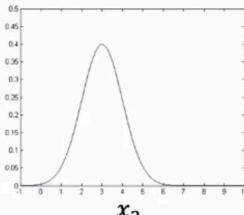


## Anomaly detection example



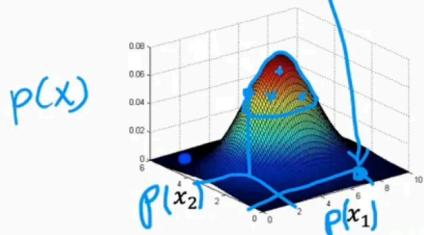
$$\mu_1 = 5, \sigma_1 = 2$$

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



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

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



$$\varepsilon = 0.02$$

$$p(x_1^{(1)}; \mu_1, \sigma_1^2) = 0.0426 \rightarrow "ok"$$

$$p(x_2^{(2)}; \mu_2, \sigma_2^2) = 0.0021 \rightarrow \text{anomaly}$$

Not only for anomaly detection algorithms but for any model development, it's important and helpful to evaluate the algorithm during development. This allows you to notice the impact of adjustments to features and parameters.

## Anomaly detection

Very small number of positive examples ( $y = 1$ ). (0-20 is common).

Large number of negative ( $y = 0$ ) examples.

p(x)

y = 1

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.

Fraud

## vs. Supervised learning

Large number of positive and negative examples.

20 positive 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.

Spam

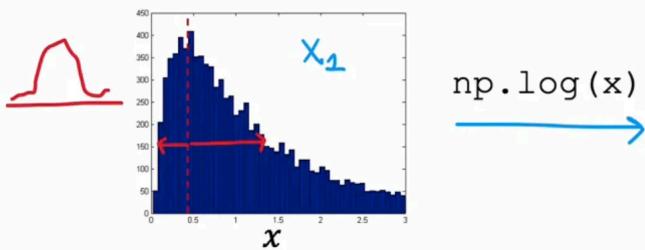
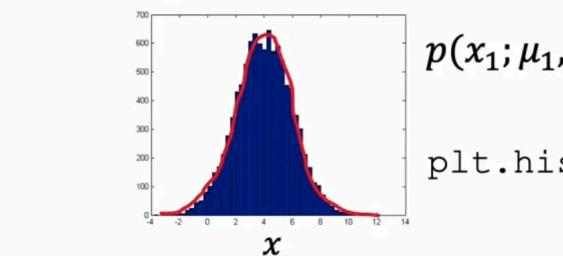
Feature selection is very important when building an anomaly detection model. This is because anomaly detection relies on unlabeled data, and unhelpful features make it harder to detect anomalies.

Anomaly detection algorithms often assume that features are normally (Gaussian) distributed, so it's important to transform non-Gaussian features into Gaussian ones. This can be done using transformations like logarithms, square roots, or exponents.

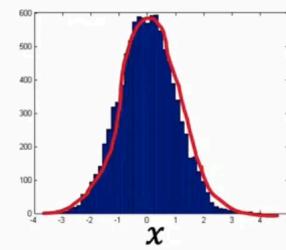
Another helpful approach is to create new features, such as combinations or ratios of existing ones, to better capture anomalous patterns.

Overall, feature engineering and error analysis are essential steps for preparing high-quality data that improves the performance of anomaly detection systems.

## Non-gaussian features



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

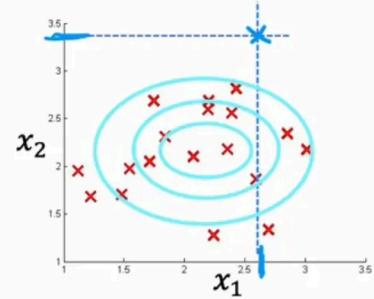
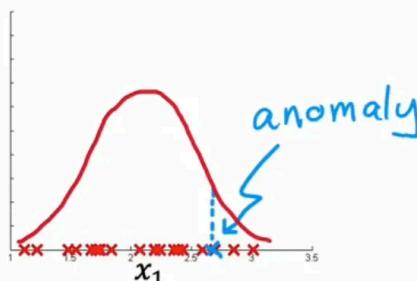


## Error analysis for anomaly detection

Want  $p(x) \geq \epsilon$  large for normal examples  $x$ .  
 $p(x) < \epsilon$  small for anomalous examples  $x$ .

Most common problem:

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



# Recommender Systems — Collaborative Filtering