# Week : 9
# Part one

## Anomaly detection :

Suppose we have some training sets we want to select which of them are faulty and which are not. In that case we need to make such graph as below using their features and predict new test sets.
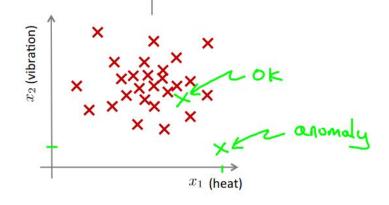
**Anomaly detection example**

Aircraft engine features:
- $x_1$ = heat generated
- $x_2$ = vibration intensity
...

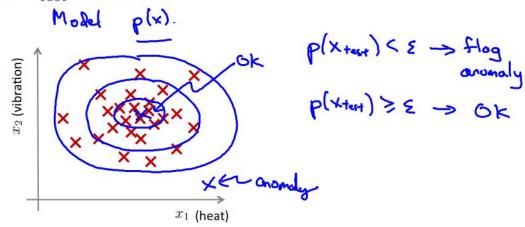Dataset: $\{x^{(1)}, x^{(2)}, \ldots, x^{(m)}\}$

New engine: $x_{test}$



Andrew Ng

## Density estimation

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

→ Is $x_{test}$ anomalous?

Model $p(x)$.

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

$$p(x_{test}) \geq \varepsilon \rightarrow ok$$

$x_2$ (vibration)

Ok

$x \leftarrow$ anomaly

$x_1$ (heat)

## 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 $p(x) < \varepsilon$

$x_1$
$x_2$      $p(x)$
$x_3$
$x_4$

→ 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.
  ...          $p(x) < \varepsilon$

# Gaussian Distribution :

Gaussian distribution is a probability distribution by which we can simulate the density of our training sets. There we used normal gaussian distribution where two perimeter are used : Mewe and sigma. There mewe is the center of the density and sigma is the width of the curve.
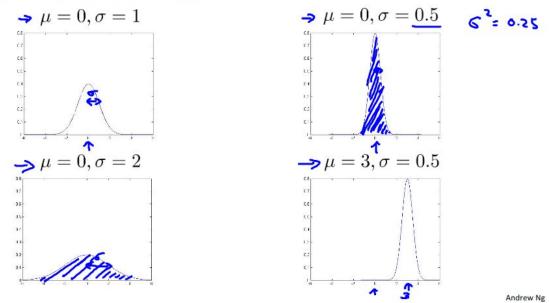
## Gaussian (Normal) distribution

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

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

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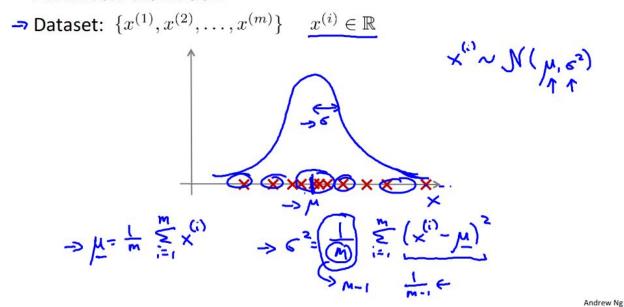
$p(x; \mu, \sigma^2)$

$\mu$  $x$

From the image below we can see that as sigma grows bigger, the width of the curve becomes thicker.

## Gaussian distribution example



$\mu = 0, \sigma = 1$

$\mu = 0, \sigma = 0.5$   $\sigma^2 = 0.25$

$\mu = 0, \sigma = 2$

$\mu = 3, \sigma = 0.5$

If we have a graph with some random training sets and want to find mewe and sigma, can apply the formula of mewe and sigma.

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

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

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

$\rightarrow m-1 \qquad \frac{1}{m-1} \leftarrow$

# Anomaly detection algorithm :

In this part we will learn an algorithm which will help us to separate test set anomaly and non-anomaly using gaussian distribution.

→ **Density estimation**

→ Training set: $\{x^{(1)}, \ldots, x^{(m)}\}$
Each example is $\underline{x \in \mathbb{R}^n}$

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

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

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

$\rightarrow p(x)$

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

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

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

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

## Anomaly detection algorithm

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

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

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

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

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

$\mu_1, \mu_2, \ldots, \mu_n$

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

→ 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 $p(x) < \varepsilon$

## Anomaly detection example

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

$\sigma_1^2, \sigma_2^2 = 4$

$\mu_1 = 5, \sigma_1 = 2$

$\mu_2 = 3, \sigma_2 = 1$

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

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

$p(x)$

$\varepsilon = 0.02$

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

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

# Developing and evaluating an anomaly detection system :

In order to develop an anomaly evaluation system we need to follow the same process of supervised learning algorithms. First we need some training sets with labels to train our algorithm and then we will have a graph from which we can plot a normal gaussian

distribution graph. And this gaussian distribution graph will help us to say if a test set is anomaly or not!

## 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: $x^{(1)}, x^{(2)}, \ldots, x^{(m)}$ (assume normal examples/not anomalous)

→ Cross validation set: $(x_{cv}^{(1)}, y_{cv}^{(1)}), \ldots, (x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$
→ Test set: $(x_{test}^{(1)}, y_{test}^{(1)}), \ldots, (x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$

$y = 1$

If we have 10k training sets with labels then we can divide them into 60-20-20 for training sets, cross-validation sets and test sets consecutively. Thus we can develop our algorithm.

## Aircraft engines motivating example

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

$\mu_1, \sigma_1^2, \ldots, \mu_n, \sigma_n^2.$

→ $\boxed{\text{Training set:} \boxed{6000} \text{good engines} (y = 0)}$   $p(x) = p(x_1; \mu_1, \sigma_1^2) \cdots p(x_n; \mu_n, \sigma_n^2)$
  CV: $\boxed{2000}$ good engines ($y = 0$), $\boxed{10}$ anomalous ($y = 1$)
  Test: $\boxed{2000}$ good engines ($y = 0$), $\boxed{10}$ anomalous ($y = 1$)

Alternative:
Training set: $\boxed{6000}$ good engines
→ CV: $\boxed{4000}$ good engines ($y = 0$), $\boxed{10}$ anomalous ($y = 1$)
→ Test: $\boxed{4000}$ good engines ($y = 0$), $\boxed{10}$ anomalous ($y = 1$)

For evaluating, we need to fit the unlabeled training sets into gaussian distribution and then evaluate our cross validation or test sets.

## Algorithm evaluation

→ Fit model $p(x)$ on training set $\{x^{(1)}, \ldots, x^{(m)}\}$    $(x_{test}^{(i)}, y_{test}^{(i)})$

→ On a cross validation/test example $x$ , predict

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

$y = 0$

Possible evaluation metrics:

→ - True positive, false positive, false negative, true negative

→ - Precision/Recall    CV

→ - $F_1$-score ←    Test set

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

# Anomaly vs supervised learning :

First few steps of anomaly and supervised learning are the same. So it may be confusing that when should we use anomaly detection and when supervised learning? It that case we need to evaluate some fact:
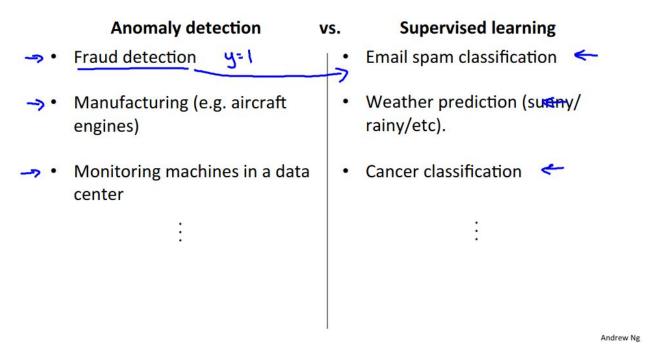
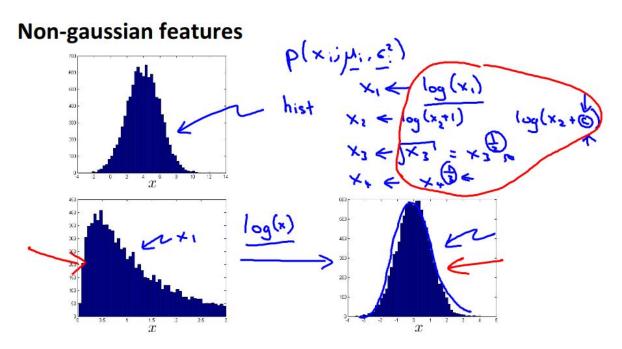| Anomaly detection | vs. | Supervised learning |
|---|---|---|
| → Very small number of positive examples ($y = 1$). (0-20 is common). | | Large number of positive and ← negative examples. |
| → Large number of negative ($y = 0$) examples. $\boxed{p(x)}$ ← | | |
| → Many different "types" of anomalies. 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 ← positive examples likely to be similar to ones in training set. |
| → future anomalies may look nothing like any of the anomalous examples we've seen so far. | | |

Spam ←

|  | Anomaly detection | vs. | Supervised learning |
|---|---|---|---|

→ • Fraud detection   $y=1$

• Email spam classification ←

→ • Manufacturing (e.g. aircraft engines)

• Weather prediction (sunny/ rainy/etc).

→ • Monitoring machines in a data center

• Cancer classification ←

⋮

⋮

Choosing what feature to use:

Sometimes it may be necessary that we need to choose the right features or debug some feature. In that case we need to follow some technique:

1. If our datasets don't give a bowl shape graph then we need to customize our features.

## Non-gaussian features



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

hist

$x_1 \leftarrow \log(x_1)$

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

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

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

$\log(x)$

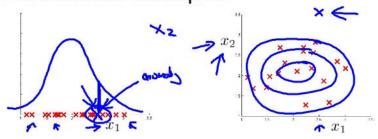2. Sometimes we need to choose a new feature to evaluate new test case:

## → 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



3. If you think that one or more features become faulty then you need to customize them with new features. In the picture below, we get x3,x4 linearly high. It that case, to detect which feature is faulty we customize them with x5 or x6 feature.

## → 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 ←

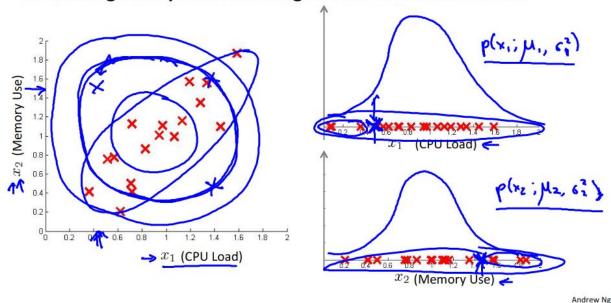$$X_5 = \frac{CPU\ load}{network\ traffic}$$

$$X_6 = \frac{(CPU\ load)^2}{network\ traffic}$$

Multivariate Gaussian Distribution :
Sometimes normal gaussian distribution may mark some test cases as non anomaly cases though they are not. It happens because of doing multiplication. In that case, we

can use multivariate gaussian distribution to get more perfect result.

## Motivating example: Monitoring machines in a data center

## Multivariate Gaussian (Normal) distribution

$x \in \mathbb{R}^n$. Don't model $p(x_1), p(x_2), \ldots$, etc. separately.
Model $p(x)$ all in one go.
Parameters: $\mu \in \mathbb{R}^n, \Sigma \in \mathbb{R}^{n \times n}$ (covariance matrix)

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

$|\Sigma| = $ determinant of $\Sigma$   |   det (Sigma)

Here are some graphs of multivariate gaussian distribution. As we notice, the graphs change with the values of mewe and sigma.

# Multivariate Gaussian (Normal) examples

# Multivariate Gaussian (Normal) examples

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

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

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

Andrew Ng

# Multivariate Gaussian (Normal) examples

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

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

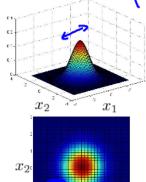$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$

Andrew Ng

# Multivariate Gaussian (Normal) examples

$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$  $\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix}$  $\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & -0.8 \\ -0.8 & 1 \end{bmatrix}$

# Multivariate Gaussian (Normal) examples

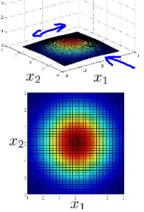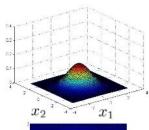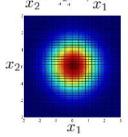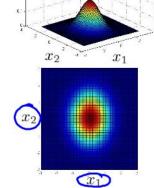$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$  $\mu = \begin{bmatrix} 0 \\ 0.5 \end{bmatrix} \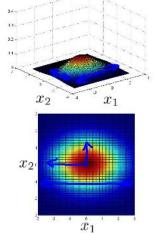Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$  $\mu = \begin{bmatrix} 1.5 \\ -0.5 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$

# Anomaly direction using multivariate gaussian distribution :

## Multivariate Gaussian (Normal) distribution

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

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



Parameter fitting:

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

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

## Anomaly detection with the multivariate Gaussian

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

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

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



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 $p(x) < \varepsilon$

## Relationship to original model

Original model: $p(x) = p(x_1; \mu_1, \sigma_1^2) \times p(x_2; \mu_2, \sigma_2^2) \times \cdots \times p(x_n; \mu_n, \sigma_n^2)$

$p(x)$



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

Corresponds to multivariate Gaussian

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

where $\Sigma = $

---

### → Original model    vs.    → Multivariate Gaussian

| Original model | Multivariate Gaussian |
|---|---|
| $p(x_1; \mu_1, \sigma_1^2) \times \cdots \times p(x_n; \mu_n, \sigma_n^2)$ | $p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left( -\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu) \right)$ |

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

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

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

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

→ Automatically captures correlations between features

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

Computationally more expensive

→ $\Sigma \sim \frac{n^2}{2}$

$\to X_1 = X_2$
$X_3 = X_4 + X_5$

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

# Part two

## Recommender system :

Online tvs (e.g : netflix, amazon, HVO) use this system to suggest movies or series to their users. In this system they let their users rate movies or series from zero to five stars and then predict their want and choice.



**Example: Predicting movie ratings**

→ User rates movies using ~~one~~ to five stars
zero

| Movie | Alice (1) | Bob (2) | Carol (3) | Dave (4) |
|---|---|---|---|---|
| Love at last | 5 | 5 | 0 | 0 |
| Romance forever | 5 | ? 4.5 | ? 0 | 0 |
| Cute puppies of love | ? 5 | 4 | 0 | ? 0 |
| Nonstop car chases | 0 | 0 | 5 | 4 |
| Swords vs. karate | 0 | 0 | 5 | ? 4 |

$n_u = 4$    $n_M = 5$

→ $n_u$ = no. users
→ $n_m$ = no. movies
→ $r(i,j)$ = 1 if user $j$ has rated movie $i$
→ $y^{(i,j)}$ = rating given by user $j$ to movie $i$ (defined only if $r(i,j) = 1$)

$0,..., 5$

Andrew Ng

# Content based recommendation :

Previously we saw that users rate movies from zero to five stars. But they don't rate all of the movies. In this section we want to make a prediction by which we can say how much a user would rate an unrated movie based on their other movie rating.

**Content-based recommender systems**

$n_u = 4, \quad n_m = 5$

$x_0 = 1$

$x^{(1)} = \begin{bmatrix} 1 \\ 0.9 \\ 0 \end{bmatrix}$

| Movie | Alice (1) $\theta^{(1)}$ | Bob (2) $\theta^{(2)}$ | Carol (3) $\theta^{(3)}$ | Dave (4) $\theta^{(4)}$ | | |
|---|---|---|---|---|---|---|
| Love at last  1 | 5 | 5 | 0 | 0 | → | → |
| Romance forever  2 | 5 | ? | ? | 0 | → | → |
| Cute puppies of love  3 | ? 4.95 | 4 | 0 | ? | → | → |
| Nonstop car chases  4 | 0 | 0 | 5 | 4 | → | → |
| Swords vs. karate  5 | 0 | 0 | 5 | ? | → | → |

$x^{(1)}, x^{(2)}, x^{(3)}, x^{(4)}, x^{(5)}$

$n = 2$

→ For each user $j$, learn a parameter $\theta^{(j)} \in \mathbb{R}^3$. Predict user $j$ as rating movie $(\theta^{(j)})^T x^{(i)}$ stars. $\theta^{(j)} \in \mathbb{R}^{n+1}$

$x^{(3)} = \begin{bmatrix} 1 \\ 0.99 \\ 0 \end{bmatrix} \longleftrightarrow \theta^{(1)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix}$

$(\theta^{(1)})^T x^{(3)} = 5 \times 0.99$
$= 4.95$

The prediction is likely to linear regression.

## Problem formulation

→ $r(i,j) = 1$ if user $j$ has rated movie $i$ (0 otherwise)

→ $y^{(i,j)} =$ rating by user $j$ on movie $i$ (if defined)

→ $\theta^{(j)}$ = parameter vector for user $j$

→ $x^{(i)}$ = feature vector for movie $i$

→ For user $j$, movie $i$, predicted rating: $(\theta^{(j)})^T (x^{(i)})$    $\theta^{(j)} \in \mathbb{R}^{n+1}$

→ $m^{(j)}$ = no. of movies rated by user $j$

To learn $\theta^{(j)}$:

$$\min_{\theta^{(j)}} \frac{1}{2m^{(j)}} \sum_{i:r(i,j)=1} \left( (\theta^{(j)})^T (x^{(i)}) - y^{(i,j)} \right)^2 + \frac{\lambda}{2m^{(j)}} \sum_{k=1}^{n} (\theta_k^{(j)})^2$$

We want to predict for all users as we calculate for theta 1,2,3,...

**Optimization objective:**

To learn $\theta^{(j)}$ (parameter for user $j$):

$$\min_{\theta^{(j)}} \frac{1}{2} \sum_{i:r(i,j)=1} \left((\theta^{(j)})^T x^{(i)} - y^{(i,j)}\right)^2 + \frac{\lambda}{2} \sum_{k=1}^{n} (\theta_k^{(j)})^2$$

To learn $\theta^{(1)}, \theta^{(2)}, \ldots, \theta^{(n_u)}$:

$$\min_{\theta^{(1)},\ldots,\theta^{(n_u)}} \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i:r(i,j)=1} \left((\theta^{(j)})^T x^{(i)} - y^{(i,j)}\right)^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^{n} (\theta_k^{(j)})^2$$

$\theta^{(1)}, \ldots, \theta^{(n_u)}$

**Optimization algorithm:**

$$\min_{\theta^{(1)},\ldots,\theta^{(n_u)}} \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i:r(i,j)=1} \left((\theta^{(j)})^T x^{(i)} - y^{(i,j)}\right)^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^{n} (\theta_k^{(j)})^2$$

$J(\theta^{(1)}, \ldots, \theta^{(n_u)})$

Gradient descent update:

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} \quad (\text{for } k = 0)$$

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \left( \sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} + \lambda \theta_k^{(j)} \right) \quad (\text{for } k \neq 0)$$

$$\frac{\partial}{\partial \theta_k^{(j)}} J(\theta^{(1)}, \ldots, \theta^{(n_u)})$$

# Collaborative filtering :

Sometimes it may be difficult to set specific feature values for all movies. Instead of doing that we can let our users rate their taste on different features. Using those inputs we can calculate the specific values for different movies.

# Problem motivation

| Movie | Alice (1) | Bob (2) | Carol (3) | Dave (4) | $x_1$ (romance) | $x_2$ (action) |
|---|---|---|---|---|---|---|
| Love at last | 5 | 5 | 0 | 0 | 0.9 | 0 |
| Romance forever | 5 | ? | ? | 0 | 1.0 | 0.01 |
| Cute puppies of love | ? | 4 | 0 | ? | 0.99 | 0 |
| Nonstop car chases | 0 | 0 | 5 | 4 | 0.1 | 1.0 |
| Swords vs. karate | 0 | 0 | 5 | ? | 0 | 0.9 |

$x_0 = 1$

| Movie | Alice (1) $\theta^{(1)}$ | Bob (2) $\theta^{(2)}$ | Carol (3) $\theta^{(3)}$ | Dave (4) $\theta^{(4)}$ | $x_1$ (romance) | $x_2$ (action) |
|---|---|---|---|---|---|---|
| Love at last $x^{(1)}$ | 5 | 5 | 0 | 0 | 1.0 | 0.0 |
| Romance forever | 5 | ? | ? | 0 | ? | ? |
| Cute puppies of love | ? | 4 | 0 | ? | ? | ? |
| Nonstop car chases | 0 | 0 | 5 | 4 | ? | ? |
| Swords vs. karate | 0 | 0 | 5 | ? | ? | ? |

$$x^{(1)} = \begin{bmatrix} 1 \\ 1.0 \\ 0.0 \end{bmatrix}$$

$$\theta^{(1)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix}, \theta^{(2)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix}, \theta^{(3)} = \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix}, \theta^{(4)} = \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix}$$

$\theta^{(j)}$

$x^{(1)}$

$(\theta^{(1)})^T x^{(1)} \approx 5$
$(\theta^{(2)})^T x^{(1)} \approx 5$
$(\theta^{(3)})^T x^{(1)} \approx 0$
$(\theta^{(4)})^T x^{(1)} \approx 0$

## Optimization algorithm

Given $\theta^{(1)}, \ldots, \theta^{(n_u)}$, to learn $x^{(i)}$:

$$\min_{x^{(i)}} \frac{1}{2} \sum_{j:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{k=1}^{n} (x_k^{(i)})^2$$

Given $\theta^{(1)}, \ldots, \theta^{(n_u)}$, to learn $x^{(1)}, \ldots, x^{(n_m)}$:

$$\min_{x^{(1)}, \ldots, x^{(n_m)}} \frac{1}{2} \sum_{i=1}^{n_m} \sum_{j:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^{n} (x_k^{(i)})^2$$

## Collaborative filtering

Given $x^{(1)}, \ldots, x^{(n_m)}$ (and movie ratings),
can estimate $\theta^{(1)}, \ldots, \theta^{(n_u)}$

$r^{(i,j)}$
$y^{(i,j)}$

Given $\theta^{(1)}, \ldots, \theta^{(n_u)}$,
can estimate $x^{(1)}, \ldots, x^{(n_m)}$

Guess $\theta \rightarrow x \rightarrow \theta \rightarrow x \rightarrow \theta \rightarrow x \rightarrow \cdots$

# Collaborative filtering algorithm :

In the previous section we learned that by collaborative filtering we can train theta and X one after another. That means we need to use X to get min theta and to get min X need to use theta. But in collaborative filtering algorithm we can minimize our perimeter theta and X simultaneously.

**Collaborative filtering optimization objective**

$(i,j) : r(i,j) = 1$

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

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

$x_1 = 1$

→ Given $x^{(1)}, \ldots, x^{(n_m)}$, estimate $\theta^{(1)}, \ldots, \theta^{(n_u)}$:

$$\min_{\theta^{(1)}, \ldots, \theta^{(n_u)}} \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^{n} (\theta_k^{(j)})^2$$

→ Given $\theta^{(1)}, \ldots, \theta^{(n_u)}$, estimate $x^{(1)}, \ldots, x^{(n_m)}$:

$$\min_{x^{(1)}, \ldots, x^{(n_m)}} \frac{1}{2} \sum_{i=1}^{n_m} \sum_{j:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^{n} (x_k^{(i)})^2$$

Minimizing $x^{(1)}, \ldots, x^{(n_m)}$ and $\theta^{(1)}, \ldots, \theta^{(n_u)}$ simultaneously:

$$J(x^{(1)}, \ldots, x^{(n_m)}, \theta^{(1)}, \ldots, \theta^{(n_u)}) = \frac{1}{2} \sum_{(i,j):r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^{n} (x_k^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^{n} (\theta_k^{(j)})^2$$

$$\min_{\substack{x^{(1)}, \ldots, x^{(n_m)} \\ \theta^{(1)}, \ldots, \theta^{(n_u)}}} J(x^{(1)}, \ldots, x^{(n_m)}, \theta^{(1)}, \ldots, \theta^{(n_u)})$$

$\theta \to x \to \theta \to x \to \ldots$

**Collaborative filtering algorithm**

$x_0 \neq 1$     $x \in \mathbb{R}^n, \theta \in \mathbb{R}^n$

→ 1. Initialize $x^{(1)}, \ldots, x^{(n_m)}, \theta^{(1)}, \ldots, \theta^{(n_u)}$ to small random values.

→ 2. Minimize $J(x^{(1)}, \ldots, x^{(n_m)}, \theta^{(1)}, \ldots, \theta^{(n_u)})$ using gradient descent (or an advanced optimization algorithm). E.g. for every $j = 1, \ldots, n_u, i = 1, \ldots, n_m$:

$\theta_0$

$\theta_n$

$$x_k^{(i)} := x_k^{(i)} - \alpha \left( \sum_{j:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) \theta_k^{(j)} + \lambda x_k^{(i)} \right)$$
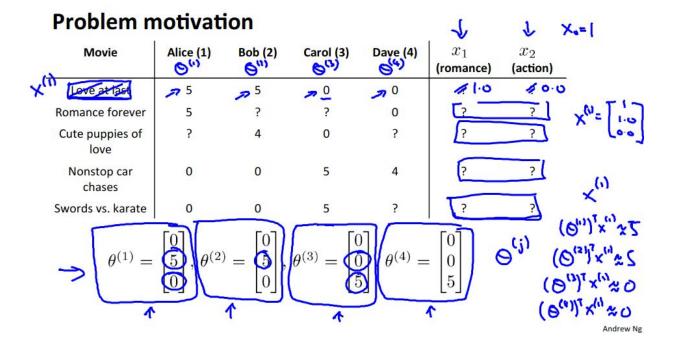
$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \left( \sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} + \lambda \theta_k^{(j)} \right)$$

$\frac{\partial}{\partial x_k^{(i)}} J(\ldots)$

3. For a user with parameters $\theta$ and a movie with (learned) features $x$, predict a star rating of $\theta^T x$.

$$(\theta^{(j)})^T (x^{(i)})$$

Vectorize low rank matrix factorization :
Instead of using loops to calculate collaborative filtering we can use vectorize implementation. It will make computation easier.

## Collaborative filtering

$n_m = 5$
$n_u = 4$

| Movie | Alice (1) | Bob (2) | Carol (3) | Dave (4) |
|---|---|---|---|---|
| Love at last | 5 | 5 | 0 | 0 |
| Romance forever | 5 | ? | ? | 0 |
| Cute puppies of love | ? | 4 | 0 | ? |
| Nonstop car chases | 0 | 0 | 5 | 4 |
| Swords vs. karate | 0 | 0 | 5 | ? |

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 \\ 5 & ? & ? & 0 \\ ? & 4 & 0 & ? \\ 0 & 0 & 5 & 4 \\ 0 & 0 & 5 & 0 \end{bmatrix}$$

$y^{(i,j)}$

## Collaborative filtering

$X(\Theta)^T$

$(\Theta^{(j)})^T(x^{(i)})$

$(i,j)$

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 \\ 5 & ? & ? & 0 \\ ? & 4 & 0 & ? \\ 0 & 0 & 5 & 4 \\ 0 & 0 & 5 & 0 \end{bmatrix}$$

Predicted ratings:

$$\begin{bmatrix} (\theta^{(1)})^T(x^{(1)}) & (\theta^{(2)})^T(x^{(1)}) & \cdots & (\theta^{(n_u)})^T(x^{(1)}) \\ (\theta^{(1)})^T(x^{(2)}) & (\theta^{(2)})^T(x^{(2)}) & \cdots & (\theta^{(n_u)})^T(x^{(2)}) \\ \vdots & \vdots & \vdots & \vdots \\ (\theta^{(1)})^T(x^{(n_m)}) & (\theta^{(2)})^T(x^{(n_m)}) & \cdots & (\theta^{(n_u)})^T(x^{(n_m)}) \end{bmatrix}$$

$$X = \begin{bmatrix} - (x^{(1)})^T - \\ - (x^{(2)})^T - \\ \vdots \\ - (x^{(n_m)})^T - \end{bmatrix} \qquad \Theta = \begin{bmatrix} - (\Theta^{(1)})^T - \\ - (\Theta^{(2)})^T - \\ \vdots \\ - (\Theta^{(n_u)})^T - \end{bmatrix}$$

→ Low rank matrix factorization

To improve user experience online tvs need to suggest related movies or videos. In that case they use existing user data to do that. Suppose, movie x1 is watched by a user. Now if you want to test if x2 movie is related to x1 then you need to following steps.

## Finding related movies

For each product $i$, we learn a feature vector $x^{(i)} \in \mathbb{R}^n$.

$\rightarrow$ $x_1$ = romance, $x_2$ = action, $x_3$ = comedy, $x_4$ = ....

How to find movies $j$ related to movie $i$?

small $\|x^{(i)} - x^{(j)}\|$ $\rightarrow$ movie $j$ and $i$ are "similar"

5 most similar movies to movie $i$:
Find the 5 movies $j$ with the smallest $\|x^{(i)} - x^{(j)}\|$.

Implementation details : mean normalization :
Suppose a new user just signed up. He hasn't rated any movie yet. How do you suggest movies or other stuff? In that case, mean normalization will help us.

## Users who have not rated any movies

| Movie | Alice (1) | Bob (2) | Carol (3) | Dave (4) | Eve (5) |
|---|---|---|---|---|---|
| Love at last | 5 | 5 | 0 | 0 | ? 0 |
| Romance forever | 5 | ? | ? | 0 | ? 0 |
| Cute puppies of love | ? | 4 | 0 | ? | ? 0 |
| Nonstop car chases | 0 | 0 | 5 | 4 | ? 0 |
| Swords vs. karate | 0 | 0 | 5 | ? | ? 0 |

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 & ? \\ 5 & ? & ? & 0 & ? \\ ? & 4 & 0 & ? & ? \\ 0 & 0 & 5 & 4 & ? \\ 0 & 0 & 5 & 0 & ? \end{bmatrix}$$

$$\min_{\substack{x^{(1)},\ldots,x^{(n_m)} \\ \theta^{(1)},\ldots,\theta^{(n_u)}}} \frac{1}{2} \sum_{(i,j):r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^{n} (x_k^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^{n} (\theta_k^{(j)})^2$$

$n=2$  $\theta^{(5)} \in \mathbb{R}^2$  $\theta^{(5)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$  $\frac{\lambda}{2}\left[(\theta_1^{(5)})^2 + (\theta_2^{(5)})^2\right] \leftarrow$

$(\theta^{(5)})^T x^{(i)} = 0$

To apply this algorithm you need to compute the average rating for every single movie. Then subtract the avg value from original value. Then use this formula for user j, on movie i predict:

## Mean Normalization:

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 & ?=2.5 \\ 5 & ? & ? & 0 & ?=2.5 \\ ? & 4 & 0 & ? & ?=2 \\ 0 & 0 & 5 & 4 & ? \\ 0 & 0 & 5 & 0 & ? \end{bmatrix}$$

$$\mu = \begin{bmatrix} 2.5 \\ 2.5 \\ 2 \\ 2.25 \\ 1.25 \end{bmatrix} \rightarrow Y = \begin{bmatrix} 2.5 & 2.5 & -2.5 & -2.5 & ? \\ 2.5 & ? & ? & -2.5 & ? \\ ? & 2 & -2 & ? & ? \\ -2.25 & -2.25 & 2.75 & 1.75 & ? \\ -1.25 & -1.25 & 3.75 & -1.25 & ? \end{bmatrix}$$

learn $\Theta^{(j)}, x^{(i)}$

For user $j$, on movie $i$ predict:

$$\rightarrow (\Theta^{(j)})^T (x^{(i)}) + \mu_i$$

User 5 (Eve):

$$\Theta^{(5)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\underbrace{(\Theta^{(5)})^T (x^{(i)})}_{\to 0} + \boxed{\mu_i}$$