# Machine Learning Part V

## Anomaly Detection/ Recommender System

# 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,
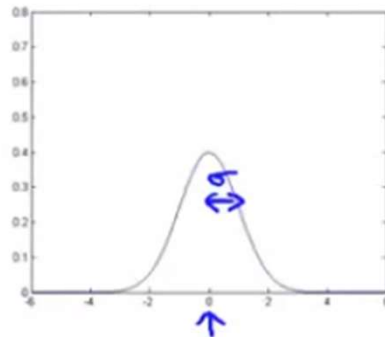
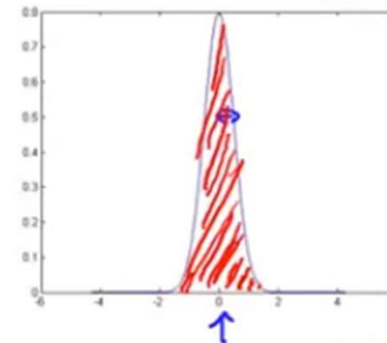  $x_3$ = CPU load, $x_4$ = CPU load/network traffic.

  ...

$x_1$
$x_2$
$x_3$
$x_4$

$p(x)$

$p(x) < \varepsilon$

# Gaussian (Normal) Distribution

- Probability distribution (add up to 1)

# Parameter estimation

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

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



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

# Anomaly detection algorithm

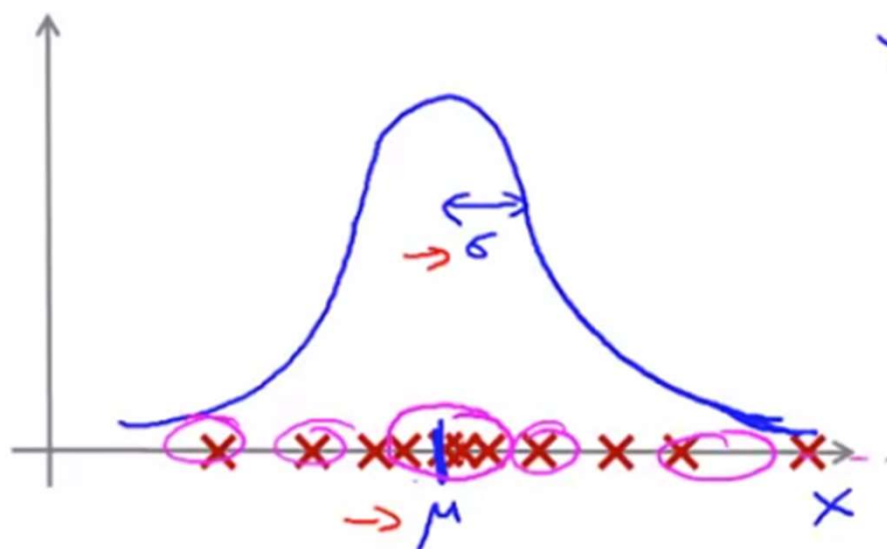1. Choose features $x_i$ that you think might be indicative of anomalous examples. $\{x^{(1)}, \ldots, x^{(m)}\}$

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

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

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

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

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

$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix} = \frac{1}{m} \sum_{i=1}^{m} x^{(i)}$$
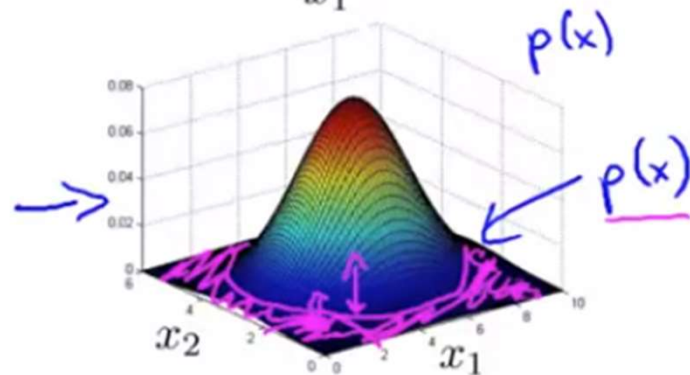
3. Given new example $x$, compute $p(x)$:

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

Anomaly if $p(x) < \varepsilon$

# Visualization



**Anomaly detection example**

$$\mu_1 = 5, \sigma_1 = 2$$
$$\mu_2 = 3, \sigma_2 = 1$$

$$\sigma_1^2, \sigma_2^2$$
$$\underbrace{\quad} = 4$$

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

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

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

$$\varepsilon = 0.02$$

$$p(x_{test}^{(1)}) = 0.0426 \quad > \varepsilon$$
$$p(x_{test}^{(2)}) = 0.0021 \quad < \varepsilon$$

Full Sc

# Evaluation (supervised learning)

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: $(x_{test}^{(1)}, y_{test}^{(1)}), \ldots, (x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$

**Aircraft engines motivating example**

→ 10000 good (normal) engines

→ 20     flawed engines (anomalous)   2-50             $y=1$

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

→ Training set: 6000 good engines $(y=0)$   $p(x) = p(x_1; \mu_1 \sigma_1^2) \cdots p(x_n; \mu_n, \sigma_n^2)$

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 $\{x^{(1)}, \ldots, x^{(m)}\}$
> On a cross validation/test example $x$, predict

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

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 $\varepsilon$

Also, use validation set to decide what features to include (square co?)

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

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

→ 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;

→ future anomalies may look nothing like any of the anomalous examples we've seen so far.

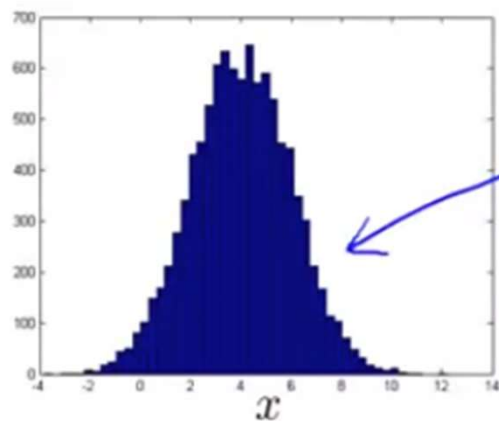Large number of positive and ← negative examples.

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.

Spam ←

# Non-gaussian features

→ **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



Adding features that can distinguish normal and anomalous samples

**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)

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

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

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

# Multivariate Gaussian (Normal) distribution

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

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



$x_1$    $x_2$        $x_1$    $x_2$        $x_1$    $x_2$

Parameter fitting:

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

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

# 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 $p(x) < \varepsilon$

## Original model → vs. → Multivariate Gaussian

$$p(x_1; \mu_1, \sigma_1^2) \times \cdots \times p(x_n; \mu_n, \sigma_n^2)$$

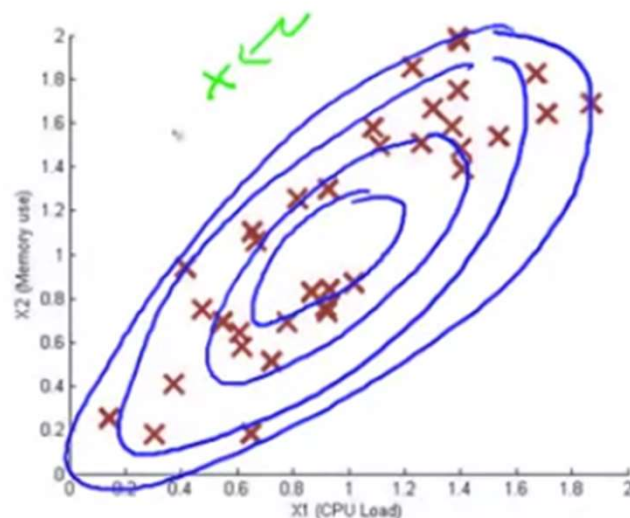$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

**Cool !**

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

→ Automatically captures correlations between features

$$x_3 = \frac{x_1}{x_2} = \frac{CPU\ load}{memory}$$

> Computationally cheaper (alternatively, scales better to large $n$) 

$n = 10,000, \quad n = 100,000$

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

$\Sigma^{-1}$

Computationally more expensive

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

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

$$\rightarrow x_1 = x_2$$
$$x_3 = x_4$$
$$+ x_5$$

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

# Recommender Systems

# Content-based recommender systems

$n_u = 4$, $n_m = 5$

$x_0 = 1$

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

| Movie | Alice (1) $\theta^{(1)}$ | Bob (2) $\theta^{(2)}$ | Carol (3) $\theta^{(3)}$ | Dave (4) $\theta^{(4)}$ | $x_1$ (romance) | $x_2$ (action) |
|---|---|---|---|---|---|---|
| Love at last 1 | 5 | 5 | 0 | 0 | 0.9 | 0 |
| Romance forever 2 | 5 | ? | ? | 0 | 1.0 | 0.01 |
| Cute puppies of love 3 | ? 4.95 | 4 | 0 | ? | 0.99 | 0 |
| Nonstop car chases 4 | 0 | 0 | 5 | 4 | 0.1 | 1.0 |
| Swords vs. karate 5 | 0 | 0 | 5 | ? | 0 | 0.9 |

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

$n = 2$

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

$\theta^{(j)} \in \mathbb{R}^{n+1}$

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

$$= 4.95$$

# Content Based Recommendations

- content means we have the features to describe the product
- It is essentially a linear regression problem, only that we train a set of parameters for each user.

**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:**

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

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

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

13:24 / 14:31

# Collaborative Filtering

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

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

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

Each movie is a vector of genre features
Each audience is also a vector of preference features

**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:

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

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

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

5:54 / 8:26

Andrew Ng

Learn x and θ in simultaneously!

**Collaborative filtering algorithm**

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 $j = 1, \ldots, n_u, i = 1, \ldots, n_m$ :

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

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

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

*Handwritten annotations:*

$x_0 = 1$ (crossed out)

$x \in \mathbb{R}^n, \; \theta \in \mathbb{R}^n$

$\theta_0$ (crossed out)

$\theta_1$

$\vdots$

$\theta_n$

$\dfrac{\partial}{\partial x_k^{(i)}} J(\cdots)$

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

# Vectorization

**Collaborative filtering**

$$X \, \Theta^T \leftarrow$$

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

**Predicted ratings:**

$$(i,j) \nearrow$$

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

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

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

$$\rightarrow \text{Low rank matrix factorization}$$

**Finding related movies**

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

$\rightarrow x_1 = \text{romance}, \; x_2 = \text{action}, \; x_3 = \text{comedy}, \; x_4 = \ldots$
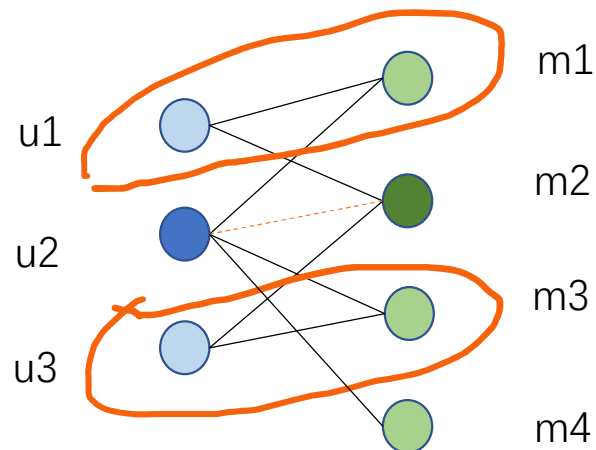
How to find movies $j$ related to movie $i$?

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

5 most similar movies to movie $i$:

$\rightarrow$ Find the 5 movies $j$ with the smallest $\|x^{(i)} - x^{(j)}\|$.

# Movie rating to network setting?

- Bipartite network (users, movies)
- Objective: **link weight prediction** (predict user rating for unseen movies).
- common **3 hop neighbors** between a user and a movie

# Ex8-2.2.2 of Collaborative Filtering

```matlab
for i = 1:num_movies
    idx = find(R(i, :) == 1);
    Theta_tmp = Theta(idx, :);
    Y_tmp = Y(i, idx);
    X_grad(i, :) = (X(i, :) * Theta_tmp' - Y_tmp) * Theta_tmp;
end


for j = 1:num_users
    idx = find(R(:, j) == 1);
    X_tmp = X(idx, :);
    Y_tmp = Y(idx, j);
    Theta_grad(j, :) = (Theta(j, :) * X_tmp' - Y_tmp') * X_tmp;
end
```