

Anomaly detection (異常検測)

Problem motivation

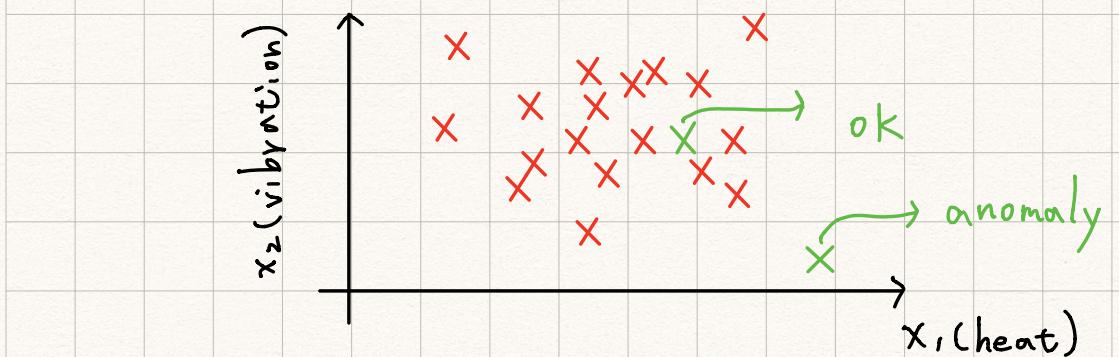
Anomaly detection example I

Aircraft engine features

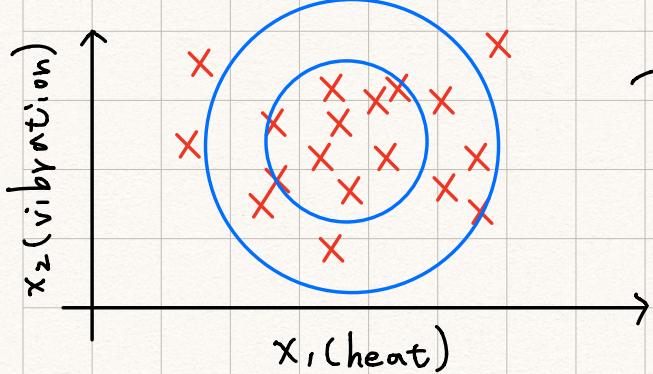
x_1 = heat generated
 x_2 = vibration intensity

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

New engine: x_{test}



Density estimation



→ 建立 model / $p(x)$

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

Anomaly detection example II

- Fraud detection:

$x^{(i)}$ = 使用者的活動特徵

Model $p(x)$ from data

Identify unusual users by checking which have $p(x) < \epsilon$

- Manufacturing

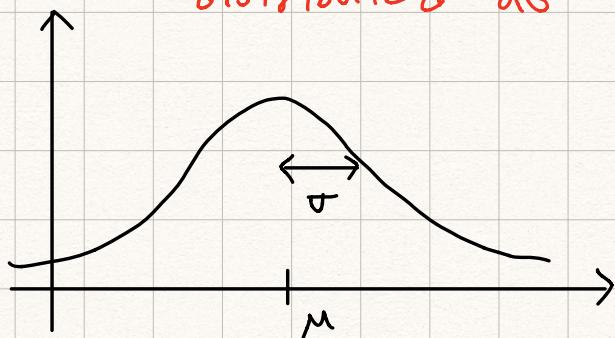
- Monitoring computers in a data center

Gaussian distribution

Say $x \in \mathbb{R}$. If x is a distributed Gaussian with mean μ , variance σ^2 (standard deviation σ)

$$x \sim N(\mu, \sigma^2)$$

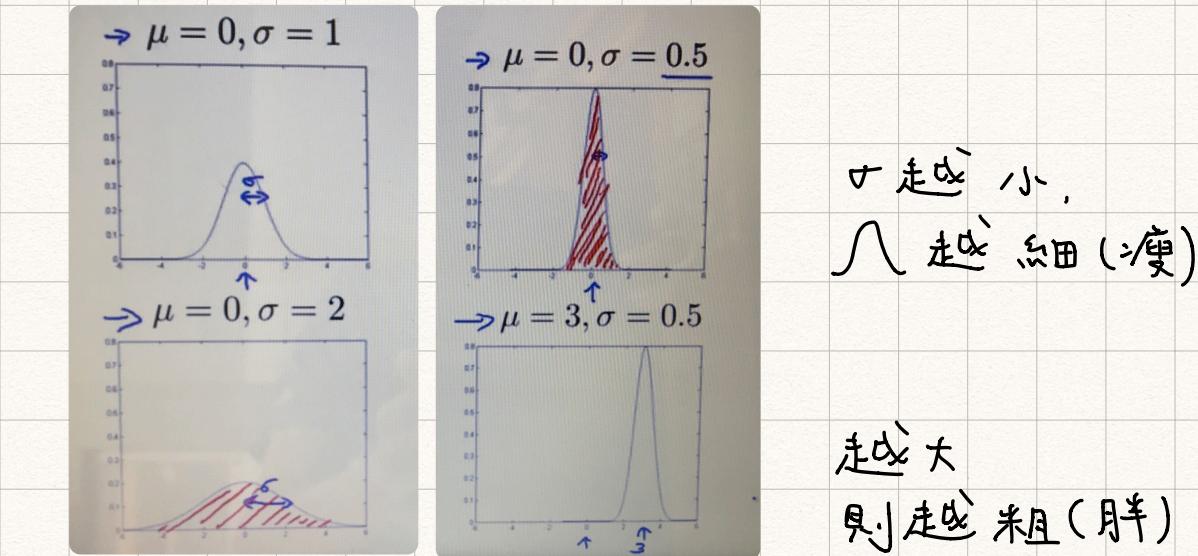
↖ "distributed as"



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

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

Gaussian distribution example

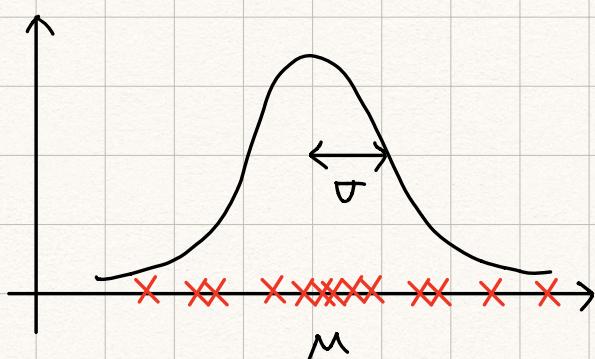


Parameter estimation

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

$x^{(i)} \in \mathbb{R}$.

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



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

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

$$m-1 \quad \frac{1}{m-1} ?$$

(ML 中，較常使用 $\frac{1}{m}$)

Algorithm

Density estimation

Training set: $\{x^{(1)}, \dots, x^{(m)}\}$, $x \in \mathbb{R}^n$

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

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

各機率 $p(x_1) \sim p(x_n)$ 獨立

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

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

Anomaly detection algorithm

1. 選擇特徵 x_i , 且這些特徵為異常例子的指標

2. Fit parameters μ_1, \dots, μ_n
 $\sigma_1^2, \dots, \sigma_n^2$

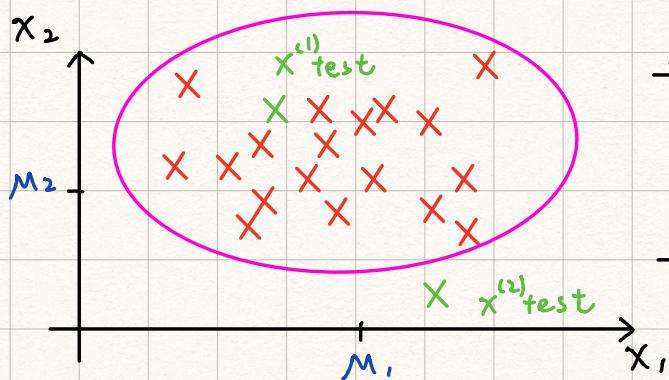
$$\hat{\mu}_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}, \quad \hat{\sigma}_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \hat{\mu}_j)^2$$

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

$$p(x) = \prod_{j=1}^n p(x_j; \hat{\mu}_j, \hat{\sigma}_j^2) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi} \cdot \hat{\sigma}_j} \exp\left(-\frac{(x_j - \hat{\mu}_j)^2}{2\hat{\sigma}_j^2}\right)$$

Anomaly : if $p(x) < \epsilon$

Anomaly detection example



找出 M_1, M_2 的
高斯分布
→ 算出 $p(x)$ 並畫出
三維圖形 $(x_1, x_2, p(x))$

$$\varepsilon = 0.02, \quad p(x^{(1)\text{test}}) = 0.0426 \geq \varepsilon$$
$$p(x^{(2)\text{test}}) = 0.021 < \varepsilon$$

Building an Anomaly Detection System

Developing and Evaluating an Anomaly Detection System

The importance of real-number evaluation

when developing a learning algorithm, making decisions is much easier if we have a way of evaluating our learning algorithm.

assume we have some labeled data, of anomalies and non-anomalous examples.

Aircraft engines motivating example

10000 good (normal) engines → 儘管裡頭有些
20 flawed engines (anomalous) 異常樣本也 OK

6 Train: 6000 good engines (unlabeled)

:

2 CV : 2000 good engines, 10 anomalous
 $y=0$ $y=1$

:

2 Test : 2000 good engines, 10 anomalous
 $y=0$ $y=1$

Algorithm evaluation

- Fit model $p(x)$ on training set
- On a cv / 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:

Precision / Recall / F₁-Score

挑選不同的 ε , 並在驗證集上確認

那個 ε 有最好的 ε

Anomaly Detection v.s. Supervised learning

Anomaly detection v.s. Supervised learning

- 十分少的正樣本數量
($y=1$) ($0-20$ is common)

- 大數量的負樣本數量
($y=0$)

• 存在許多不同種類異常狀況，未來遇到的異常可能是以前從沒見過的

- 大數量的正負樣本

- 許多的正樣本讓算法了解正樣本長得如何

Email spam classification

Weather prediction

Fraud detection

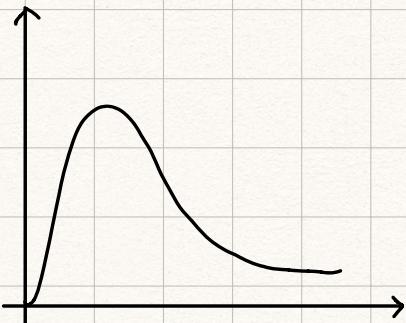
Cancer classification

Manufacturing

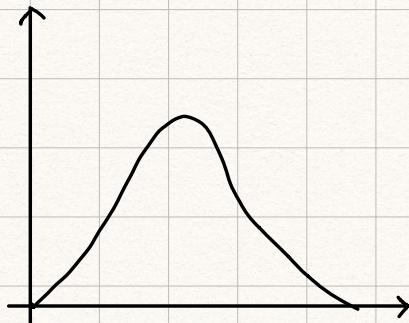
Monitoring machines in a data center

Choosing what features to use

Non-gaussian features



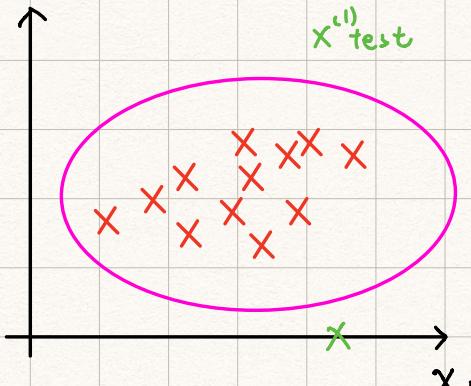
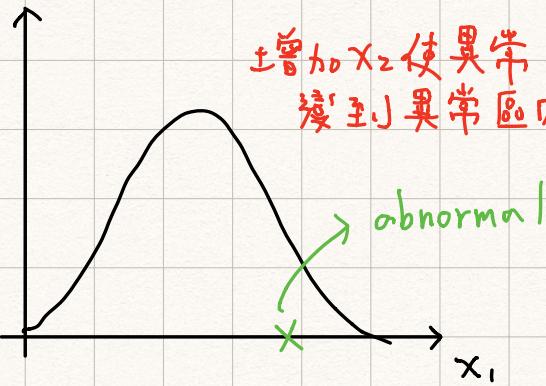
$\log(x)$



如果某特徵分佈得不像高斯函數，可對其取 $\log(x)$, $\log(x+c)$, x^c
調整 c 值 \rightarrow hist \rightarrow 見察
使特徵為高斯分佈故可讓算法工作良好

Error analysis for anomaly detection

Want $p(x)$ large for normal examples x
 $p(x)$ small for anomalous examples x



Monitoring computers in a data center

Choosing 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

電腦陷入死循環， x_3 大 but x_4 小，故設計 x_5, x_6 feature

x_5 = CPU / network traffic

x_6 = (CPU load)² / network traffic

Recommender Systems

reason: 許多硅谷工程師最想改進的
便是推薦系統

(Netflix, Amazon, 電商都有在用)

在學術界，推薦系統的討論卻滿少的

Predicting Movie Ratings

Problem Formulation

Example: Predicting movie rating

User rates movies using zero to five stars

| movie | people 1 | 己 | 3 | 4 | x_1 | x_2 |
|---------|----------|---|---|---|-------|-------|
| movie 1 | 5 | 5 | 0 | 0 | 0.9 | 0 |
| 己 | 5 | ? | ? | 0 | 1 | 0.01 |
| 3 | ? | 4 | 0 | ? | 0.99 | 0 |
| 4 | 0 | 0 | 5 | 4 | 0.1 | 1.0 |
| 5 | 0 | 0 | 5 | ? | 0 | 0.9 |

根據上述 data, 預測「？」處是多少

Content Based Recommendations

設定特徵 x_1 (romance) . x_2 (action)

| | | | |
|---|---------------|------|------|
| $x_0 = 1$ | \rightarrow | 0.9 | 0 |
| $x^{(1)} = \begin{bmatrix} 1 \\ 0.9 \\ 0 \end{bmatrix}$ | | 1.0 | 0.01 |
| | | 0.99 | 0 |
| | | 0.1 | 1.0 |
| | | 0 | 0.9 |

For each user j , learn a parameter $\theta^{(j)} \in \mathbb{R}^3$.
 Predict user j 's rating for movie i , with
 $(\theta^{(j)})^T x^{(i)}$ stars.

$$x^{(3)} = \begin{bmatrix} 1 \\ 0.99 \\ 0 \end{bmatrix} \leftrightarrow \theta^{(1)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix}, (\theta^{(1)})^T x^{(1)} = 4.95$$

對每個用戶應用不同的線性迴歸模型
 → 每個人都有一個 θ

Problem formulation

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

$\theta^{(j)}$ = parameter vector for user j
 $x^{(i)}$ = feature vector for movie i

$\theta^{(j)}$ = parameter vector for user j
 $x^{(i)}$ = feature vector for movie i
 For user j , movie i , predicted rating: $(\theta^{(j)})^T x^{(i)}$

Optimization objective

To learn $\theta^{(j)}$. $\theta^{(j)} \in \mathbb{R}^{n+1}$:

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

To learn $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(n)}$: $J(\theta^{(1)} + \dots + \theta^{(n)})$

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

Gradient descent update

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

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

如果對 x_1, x_2 特徵又沒有掌握，怎麼辦？

Collaborative Filtering (協同過濾)

能自動替人學習到優良的特徵

Problem motivation

x_1 (romance) x_2 (action) $x_0 = 1$

? ?
?
?
?
?
?
?

找尋方 user，問他們喜
歡 romance or action 的程度，藉
此得到 θ 值

user 1~4

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

由 θ 和 user 對產品的評分

推測該產品的特徵 (to movie 屬 x_1 v x_2 的程度)

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

5. 5. 0. 0. 為評分

Optimization algorithm

Given $\theta^{(1)}, \dots, \theta^{(n)}$, 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 $\theta^{(1)}, \dots, \theta^{(n)}$, to learn $x^{(1)}, \dots, x^{(n)}$:

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

Collaborative filtering

Given $x^{(1)}, \dots, x^{(n)}$ (and movie ratings)
can estimate $\theta^{(1)}, \dots, \theta^{(n)}$

Given $\theta^{(1)}, \dots, \theta^{(n)}$
can estimate $x^{(1)}, \dots, x^{(n)}$

Guess $\theta \rightarrow x \rightarrow \theta \rightarrow x \rightarrow \theta \rightarrow x \rightarrow \dots$

Collaborative Filtering Algorithm

Given $x^{(1)} \cup x^{(nm)}$, estimate $\theta^{(1)} \cup \theta^{(nu)}$

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

Given $\theta^{(1)} \cup \theta^{(nu)}$, estimate $x^{(1)} \cup x^{(nm)}$

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

* 合併以上兩種算法，便不需 $x \rightarrow \theta \rightarrow x \rightarrow \theta \dots$

minimizing $x^{(1)}, \dots, x^{(nm)}$ and $\theta^{(1)}, \dots, \theta^{(nu)}$ simultaneously

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

$$\min_{\substack{x^{(1)} \dots x^{(nm)} \\ \theta^{(1)} \dots \theta^{(nu)}}} J(x^{(1)}, \dots, x^{(nm)}, \theta^{(1)}, \dots, \theta^{(nu)})$$

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

$$x_0 = 1$$

Collaborative filtering algorithm

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

2. Minimize $J(x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)})$

using gradient descent (or an advanced optimization algorithm). E.g. for every $j = 1, \dots, n_u$, $i = 1, \dots, n_m$:

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

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

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

Low Rank Matrix Factorization

Vectorization

$$n_m = 5, n_u = 4 \rightarrow$$

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

$$X = \begin{bmatrix} \dots (x^{(1)})^T \dots \\ \dots (x^{(2)})^T \dots \\ \vdots \\ \dots (x^{(n)})^T \dots \end{bmatrix} \quad \theta = \begin{bmatrix} \dots (\theta^{(1)})^T \dots \\ \dots (\theta^{(2)})^T \dots \\ \vdots \\ \dots (\theta^{(n)})^T \dots \end{bmatrix}$$

Predicted ratings: $X\theta^T$



• 特徵 X, θ 向量化

Finding related movies

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

$X_1 = \text{romance}, X_2 = \text{action}, X_3 = \text{comedy}, X_4 = \dots$

How to find movie j related to movie i ?

Small $\|x^{(i)} - x^{(j)}\|$ → movie j and i are similar

根據用戶看過電影 i 的特徵 x .

去找尋是否有其他雷同的 x

5 most similar movies to movie i :

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

Implementational Detail: Mean Normalization

能讓演算法
運作的更好

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

見此行，根據 $J()$

公式由於不存在任何
 $r(i, j) = 1$, 故只剩下
黑色框起來的這項,

$$J = \frac{\lambda}{2} \left((\theta_1^{(s)})^2 + (\theta_2^{(s)})^2 \right)$$

$$\theta^{(s)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$J(x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_w)}) = \frac{1}{2} \sum_{(i,j): r(i,j)=1} ((\theta^{(j)})^\top x^{(i)} - y^{(i),j})^2$$

$$+ \frac{\lambda}{2} \sum_{j=1}^{n_m} \sum_{k=1}^n (\theta_k^{(j)})^2$$

$$\boxed{\frac{\lambda}{2} \sum_{j=1}^{n_m} \sum_{k=1}^n (\theta_k^{(j)})^2}$$

Mean Normalization

$$M_Y = \begin{pmatrix} 2.5 \\ 2.5 \\ 2 \\ 2.5 \\ 1.25 \end{pmatrix} \rightarrow Y - M = \begin{pmatrix} 2.5 & 2.5 & -2.5 & -2.5 & ? \\ 2.5 & ? & ? & -2.5 & ? \\ ? & 2 & -2 & ? & ? \\ -2.5 & -2.25 & 2.75 & 1.75 & ? \\ -2.5 & -1.25 & 3.75 & -1.25 & ? \end{pmatrix}$$

For user j , on movie i predict:

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

User 5:

$$\theta^{(s)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\underbrace{(\theta^{(s)})^T (x^{(i)})}_{= 0} + \boxed{\mu;}$$

$$\therefore = \mu;$$