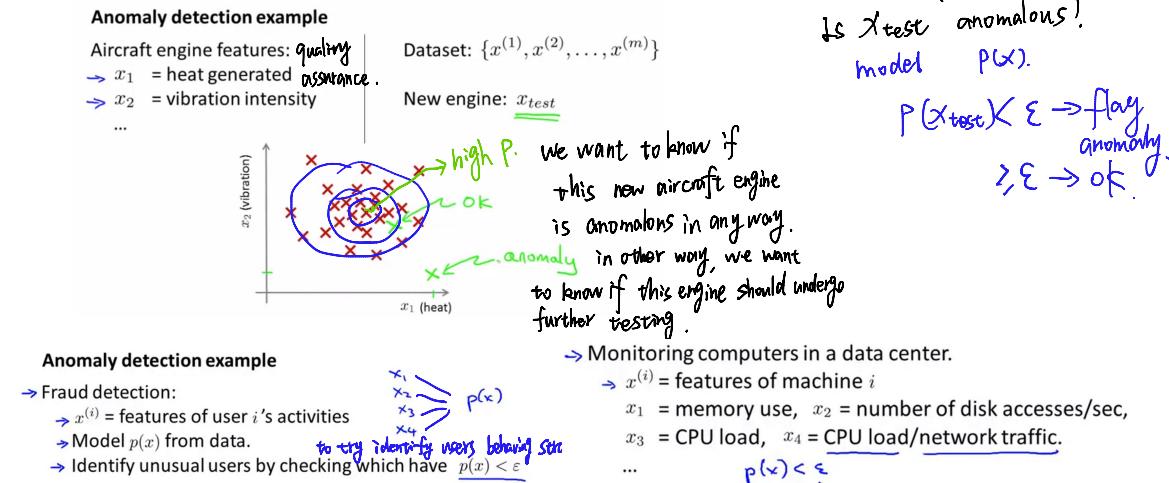


Week 9

• Density Estimation:

a. Problem Motivation:

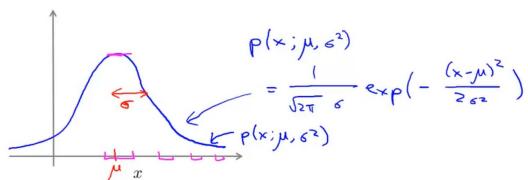


b. Gaussian Distribution:

Gaussian (Normal) distribution

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

$$x \sim \mathcal{N}(\mu, \sigma^2) \quad \text{"distributed as"}$$



$$\rightarrow \mu = 0, \sigma = 1$$

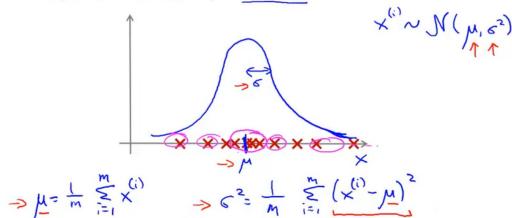
$$\rightarrow \mu = 0, \sigma = 2$$

$$\rightarrow \mu = 0, \sigma = 0.5$$

$$\rightarrow \mu = 3, \sigma = 0.5$$

Parameter estimation

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



c. Algorithm:

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

$$\begin{aligned} p(x) &= p(x_1; \mu_1, \sigma_1^2) p(x_2; \mu_2, \sigma_2^2) p(x_3; \mu_3, \sigma_3^2) \dots p(x_n; \mu_n, \sigma_n^2) \\ &= \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2) \end{aligned}$$

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

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

① choose features x_i that you think might be indicative of anomalous examples.

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

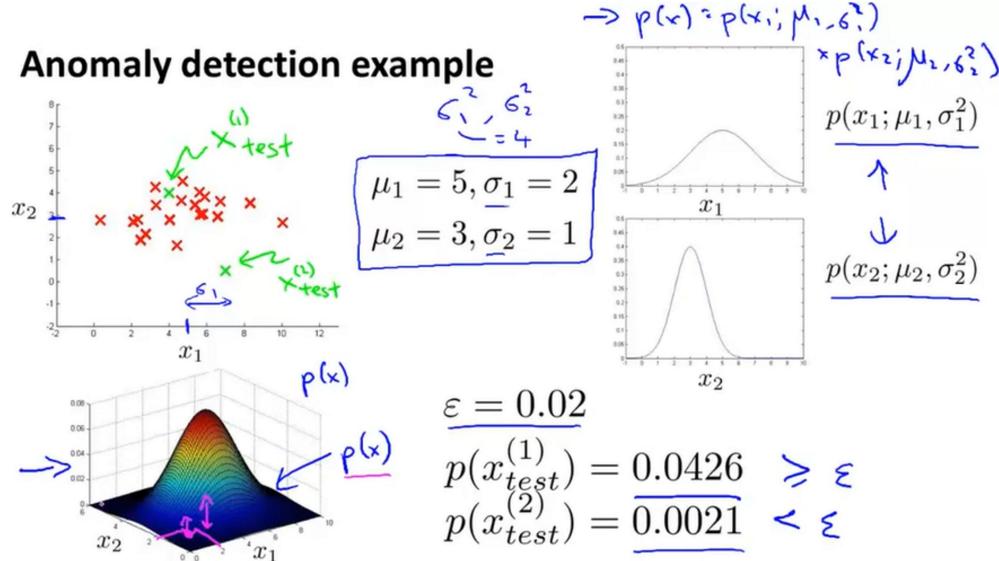
$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$$

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

③ Give new example x , compute $p(x)$:

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

Anomaly $p(x) < \epsilon$.



• Building Anomaly Detection System:

a. Developing and Evaluating an Anomaly Detection System:

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: $x^{(1)}, x^{(2)}, \dots, x^{(m)}$ (assume normal examples/not anomalous) *still unsupervised learning*.
- Cross validation set: $(x_{cv}^{(1)}, y_{cv}^{(1)}), \dots, (x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$
- Test set: $(x_{test}^{(1)}, y_{test}^{(1)}), \dots, (x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$

Do I need to add one feature?
I want to run the algorithm including this feature or not. It will return two real numbers to evaluate.

Aircraft engines motivating example

- 10000 good (normal) engines
- 20 flawed engines (anomalous) *fit the parameters* $\mu_1, \sigma_1^2, \dots, \mu_n, \sigma_n^2$ $y=1$
- Training set: 6000 good engines ($y=0$) $p(x) = p(x_1; \mu_1, \sigma_1^2) \dots p(x_n; \mu_n, \sigma_n^2)$
CV: 2000 good engines ($y=0$), 10 anomalous ($y=1$)
Test: 2000 good engines ($y=0$), 10 anomalous ($y=1$)

- Fit model $p(x)$ on training set $\{x^{(1)}, \dots, x^{(m)}\}$
 - On a cross validation/test example x , predict $(x_{test}^{(1)}, y_{test}^{(1)})$
- $$y = \begin{cases} 1 & \text{if } p(x) < \epsilon \text{ (anomaly)} \\ 0 & \text{if } p(x) \geq \epsilon \text{ (normal)} \end{cases} \quad y=0$$

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 ϵ *choose ϵ that max the F-score*.

b. Anomaly Detection vs. Supervised Learning:

when use an anomaly detection algorithm, and whether it might be more fruitful instead of using a supervisor in the algorithm.

Anomaly detection

- Very small number of positive examples ($y = 1$). (0-20 is common).
 - Large number of negative ($y = 0$) examples. $p(x)$
 - 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.
try to hard to model the positive examples because tomorrow's anomaly may be nothing like the ones you've seen so far.
- application:** Fraud detection
Manufacturing (e.g. aircraft engines)
Monitoring machines in a data center

vs.

Supervised learning

Large number of positive and negative examples.

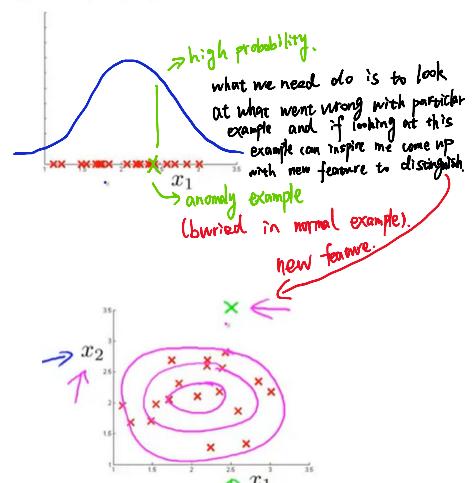
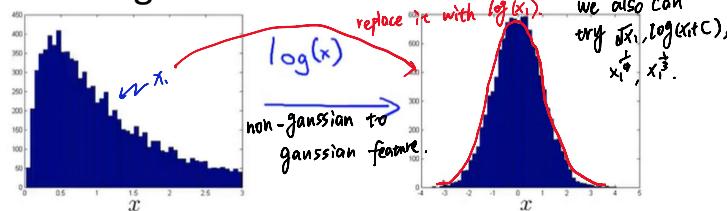
Enough positive examples for algorithm to get a sense of what positive examples are like, future positive examples likely to be similar to ones in training set.

Email spam classification.

Weather prediction

Cancer classification.

c. Choosing What Features to Use:



→ 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

- 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

so x_5 will be large.

$$x_5 = \frac{\text{CPU load}}{\text{network traffic}}$$

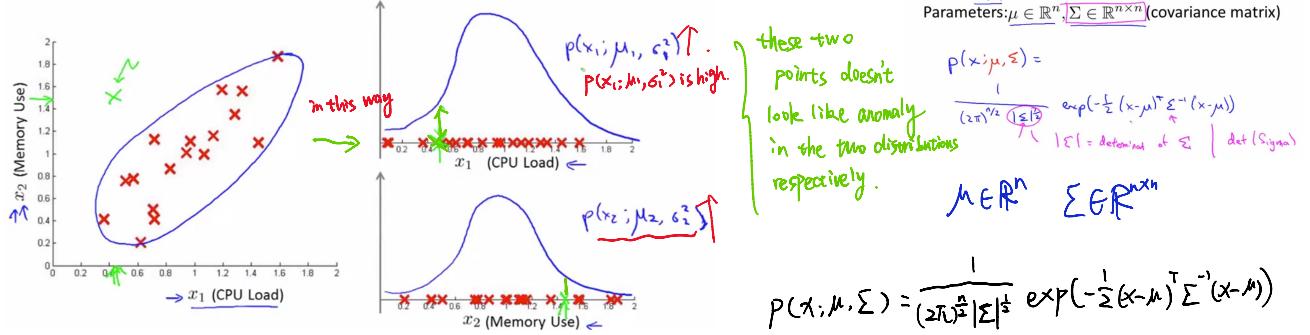
more feature

monitor some special case: computer gets stuck in infinite loop

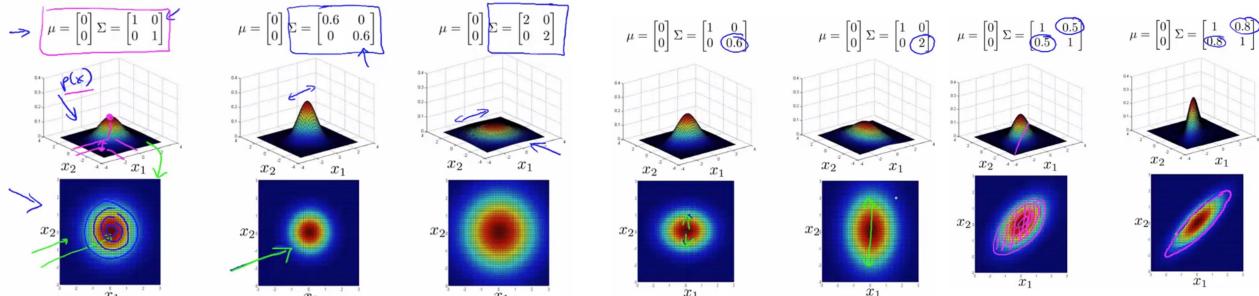
high CPU load but low network traffic

• Multivariate Gaussian Distribution:

a. Multivariate Gaussian Distribution:



Multivariate Gaussian (Normal) examples



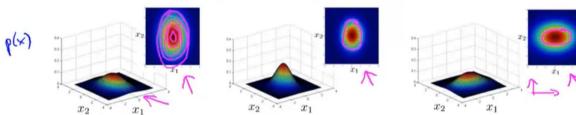
b. Anomaly Detection using the Multivariate Gaussian Distribution:

Parameter fitting:

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

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

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



For the original model, its contours of the Probability density function are axis aligned.

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

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \cdots & \sigma_{n-1}^2 \\ \vdots & \ddots & \vdots \\ \sigma_{n-1}^2 & \cdots & \sigma_n^2 \end{bmatrix}$$

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

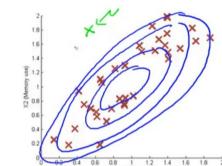
$$p(x_1; \mu_1, \sigma_1^2) \times \dots \times p(x_n; \mu_n, \sigma_n^2)$$

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

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

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

OK even if m (training set size) is small



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

→ Automatically captures correlations between features

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

Computationally more expensive

→ $\Sigma \approx \frac{n}{2}$ Σ^{-1} is redundant features.

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

• Predicting Movie Ratings:

a. Problem Formulation:

Pick a set of features automatically.

Example: Predicting movie ratings

→ User rates movies using one to five stars

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)
Love at last	5	5	6	6
Romance forever	5	5	6	6
Cute puppies of love	?	4	?	5
Nonstop car chases	0	0	0	5
Swords vs. karate	0	0	0	5

$$n_u = 4$$

$$n_m = 5$$

- $n_u = \text{no. users}$
- $n_m = \text{no. movies}$
- $r(i, j) = 1$ if user j has rated movie i
- $y^{(i,j)} = \text{rating given by user } j \text{ to movie } i$ (defined only if $r(i, j) = 1$)
- $0, \dots, 5$

What recommender system do is to predict "?" area.

b. Content Based Recommendations:

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)	$x_0=1$
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	?	

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.

$$x^{(3)} = \begin{bmatrix} 1 \\ 0.9 \\ 0 \end{bmatrix} \text{ suppose } \theta^{(1)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix} \rightarrow (\theta^{(1)})^T x^{(3)} = 4.95.$$

$\rightarrow r(i,j) = 1$ if user j has rated movie i (0 otherwise)
 $\rightarrow y^{(i,j)}$ = rating by user j on movie i (if defined)

$\rightarrow \theta^{(j)}$ = parameter vector for user j
 $\rightarrow x^{(i)}$ = feature vector for movie i
 \rightarrow For user j , movie i , predicted rating: $\underline{(\theta^{(j)})^T x^{(i)}} \quad \theta^{(j)} \in \mathbb{R}^3$

$\rightarrow m^{(j)}$ = no. of movies rated by user j
To learn $\theta^{(j)}$:

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

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

Gradient descent update:

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

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

• Collaborative Filtering:

a. Collaborative Filtering:

It can start to learn for itself what features to use.

Problem motivation

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)	$x_0=1$
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	?	

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

how much they like romantic movies
and how much they like action movies.

know these features is time consuming.

we can conclude from Alice and Bob like romantic movie and they give a high rate to

"Love at last".

$$(\theta^{(1)})^T x^{(1)} \approx 10, (\theta^{(2)})^T x^{(1)} \approx 20 \quad \left\{ x^{(1)} = \begin{bmatrix} 1 \\ 1.0 \\ 0.0 \end{bmatrix} \right.$$

$$(\theta^{(1)})^T x^{(1)} \approx 10, (\theta^{(2)})^T x^{(1)} \approx 20 \quad \left\{ x^{(1)} = \begin{bmatrix} 1 \\ 1.0 \\ 0.0 \end{bmatrix} \right.$$

Given $\theta^{(1)}, \dots, \theta^{(n_u)}$, to learn $x^{(i)}$ \rightarrow i th movie.

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

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

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

Collaborative filtering

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

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

Every users are helping algorithm to learn a better features, and these features can be used by the system to make better movies predictions for everyone else.

Guess $\Theta \rightarrow x \rightarrow \Theta \rightarrow x \rightarrow \Theta \rightarrow x \rightarrow \dots$

b. Collaborative Filtering Algorithm:

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

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

sum over all users j and then sum over all movies rated by that user.

drop!

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

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

for every movie i , sum over all the users j that have rated that movie.

when we learn features in this way, convention $x_0 = 1$, but in this way $x \in \mathbb{R}^n$.

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

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

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

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

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

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

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

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

~~On~~

• Low Rank Matrix Factorization:

a. Vectorization: Low Rank Matrix Factorization:

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

Predicted ratings: $(\theta^{(1)})^T (x^{(1)})$... $(\theta^{(n_u)})^T (x^{(1)})$

$(\theta^{(1)})^T (x^{(2)})$... $(\theta^{(n_u)})^T (x^{(2)})$

\vdots

$(\theta^{(1)})^T (x^{(n_m)})$... $(\theta^{(n_u)})^T (x^{(n_m)})$

$X = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(n_m)} \end{bmatrix}$

$\Theta = \begin{bmatrix} \theta^{(1)} \\ \theta^{(2)} \\ \vdots \\ \theta^{(n_u)} \end{bmatrix}$

$\Rightarrow X \Theta^T$

Low rank matrix factorization.

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 movies j related to movie i ?

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

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

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

\Rightarrow 5 most similar to movie i :

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

Finding related movies

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

$\rightarrow x_1 = \text{romance}, x_2 = \text{action}, x_3 = \text{comedy}, x_4 = \dots$

How to find movies j related to movie i ?

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

$L - (\theta^{(i)})^\top -$

$\Rightarrow 5$ most similar to movie i :

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

b. Implementational Detail: Mean Normalization:

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)	Eve (5)
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	?	?

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

$\Theta^{(s)} \in \mathbb{R}^n$ equals 0 for Eve.

$$\frac{\lambda}{2} [(\theta_1^{(s)})^2 + (\theta_2^{(s)})^2] \Rightarrow \theta^{(s)} = [0, 0]$$

Mean Normalization:

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 & ? \\ 5 & ? & ? & 0 & ? \\ ? & 4 & 0 & ? & ? \\ 0 & 0 & 5 & 4 & ? \\ 0 & 0 & 5 & 0 & ? \end{bmatrix} \rightarrow \mu = \begin{bmatrix} 2.5 \\ 2.5 \\ 2.25 \\ 1.25 \end{bmatrix} \rightarrow Y = \begin{bmatrix} 2.5 & 2.5 & -2.5 & -2.5 & ? \\ 2.5 & ? & ? & -2.5 & ? \\ ? & 2 & -2 & ? & ? \\ -2.25 & -2.25 & 2.75 & 1.75 & ? \\ -1.25 & -1.25 & 3.75 & -1.25 & ? \end{bmatrix}$$

For user j , on movie i predict:

$$\rightarrow (\Theta^{(s)})^\top (x^{(i)}) + \mu_i$$

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

User 5 (Eve):

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

$$\underbrace{(\Theta^{(s)})^\top (x^{(i)})}_{\approx 0} + \mu_i$$