

WEEK - 9

ANAMOLY DETECTION

Example:

Aircraft engine features:

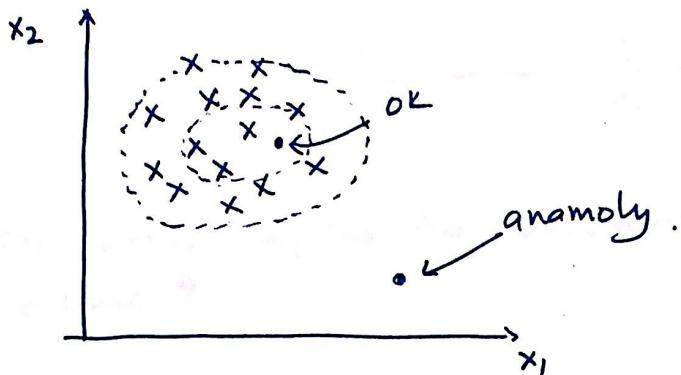
x_1 = heat generated

x_2 = vibration intensity

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

New engine: x_{test}

is x_{test} is anomalous?



$\left. \begin{array}{l} \text{Build a} \\ \text{Model } p(x). \end{array} \right. \quad \left. \begin{array}{l} p(x_{\text{test}}) < \epsilon \rightarrow \text{flag anomaly} \\ p(x_{\text{test}}) \geq \epsilon \rightarrow \text{OK} \end{array} \right\}$

Fraud detection example:

→ x^i = features of user i 's activities

→ Model $p(x)$

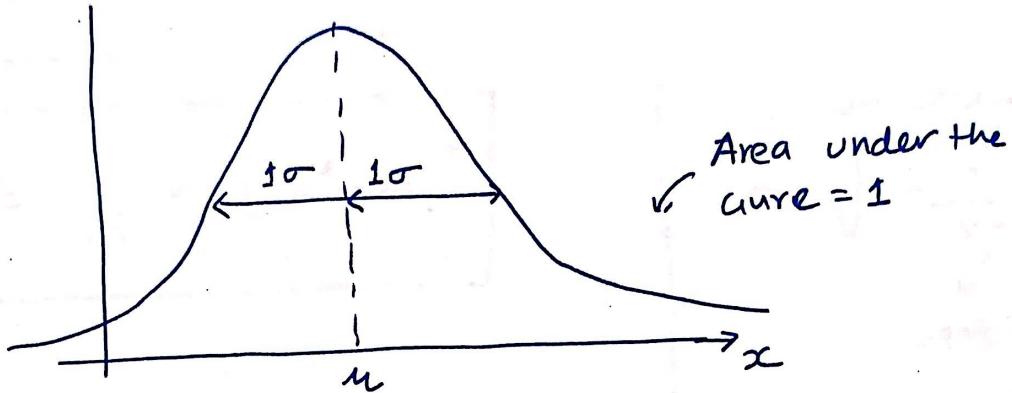
→ Identify unusual users by checking which have $p(x) < \epsilon$

- Gaussian distribution
(Normal)

$$x \sim N(\mu, \sigma^2) : x \in \mathbb{R}$$

σ^2 = variance

σ = S.D



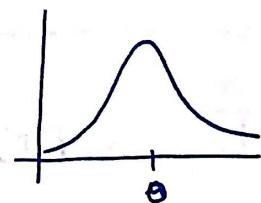
- The curve defines the probability of x taking different values

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

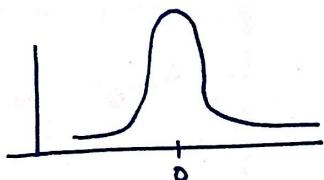
$$\Downarrow$$

$$\begin{aligned} & \approx \text{simil}(x, \ell) \\ & = \exp\left(-\frac{\|x - \ell\|^2}{2\sigma^2}\right) \end{aligned}$$

① $\mu=0, \sigma=1$



② $\mu=0, \sigma=0.5$



• Parameter Estimation :-

(111)

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

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

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

i. in statistics.

$$\sigma = \sqrt{\frac{\sum (x - \mu)^2}{n-1}}$$

sample

As in machine learning $n \gg 25$

$\therefore n-1$ is not used

• Density Estimation :

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

$\underbrace{\qquad\qquad\qquad}_{\neq n \text{ features.}}$

\therefore

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

Anomaly detection algorithm

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

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

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

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

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

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

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

- ④ Anomaly if $p(x) < \epsilon$

• Developing & Evaluation of Anomaly detection System

(113)

- 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 anomalous & non-anomalous examples ($y \in \{0, 1, \dots\}$)
- Training set: $\{x^1, x^2, \dots, x^m\} : x_{train}$
- $\{x_{test}, x_{cv}\}_{\text{and } \cancel{x_{train}}}$
 $\{y_{test}, y_{cv}\}$

example: 10000 good engines

20 anomalous.

good : Training set ($y=0$)	CV		test	
	2000 good engines	10 anomalous ($y=1$)	2000 good ($y=0$)	10 anomalous ($y=1$)

- Algorithm Evaluation :

→ Fit model $p(x)$ on training set $\{x^1, x^2, \dots, x^m\}$

→ on a cv/test x , predict

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

→ Possible evaluation metrics

		1	0
1	True positive	False positive	
	False negative	True negative	

② precision / Recall

③ F1-score

$$\underline{\text{precision}} = \frac{\text{True positive}}{\text{True pos} + \text{False posi}}$$

$$\underline{\text{Recall}} = \frac{\text{True pos}}{\text{False neg} + \text{True pos.}}$$

$$\underline{\text{Accuracy}} = \frac{\text{True pos} + \text{True neg}}{\text{All examples}}$$

Anomaly Detection vs. Supervised Learning

(115)

- ① Very small number of positive examples ($y=1$)
- ② Large number of negative examples ($y=0$)

Large number of positive & negative examples

- * Many different "type" 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.

* 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

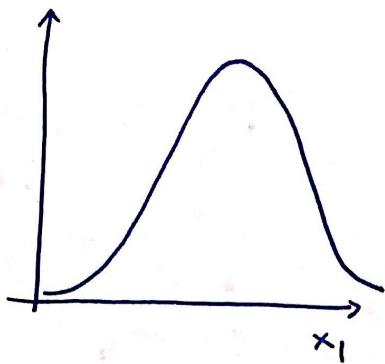
Examples:

- fraud detection
- Manufacturing (e.g. aircraft)
- Monitoring machines in a data center

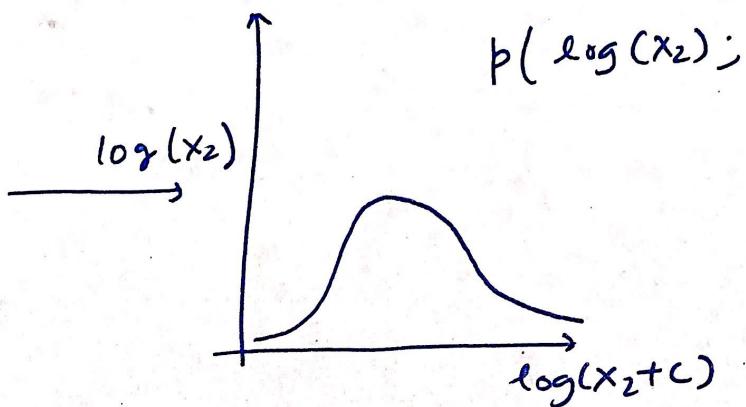
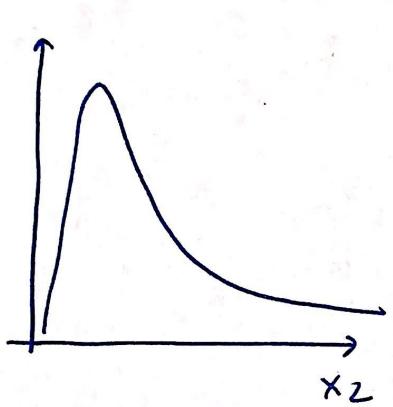
- Email classification
- weather prediction
- Cancer classification

- choosing what features to use :-

116



$$p(x_1; u, \sigma^2)$$



$$p(\log(x_2); u, \sigma^2)$$

Octave

`hist(x.^0.5, 50)`

`hist(log(x), 50)`

• Error analysis for anomaly detection:

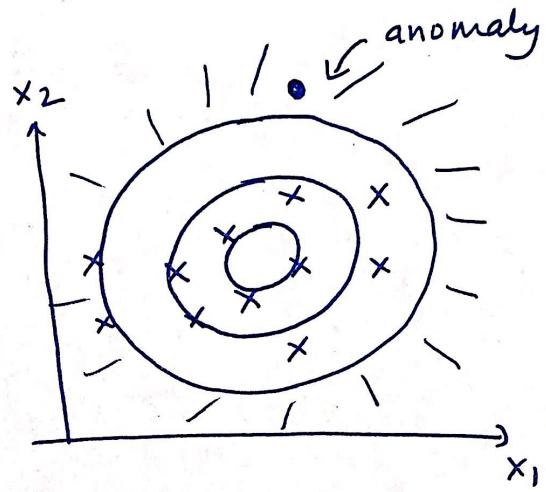
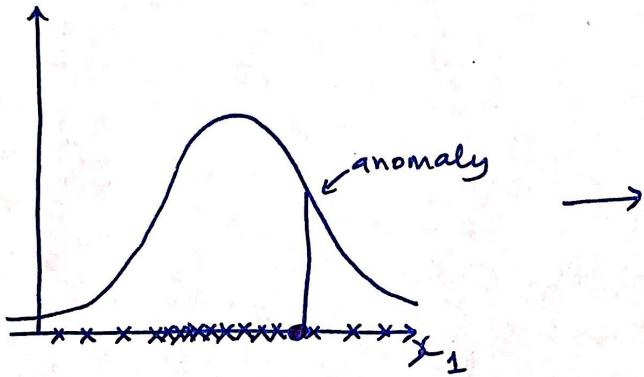
(117)

Want $p(x) \geq \epsilon$: large for normal examples, x

$p(x) \leq \epsilon$: small for anomalous examples, x

Most common problem:

$p(x)$ is comparable for normal & anomalous examples.



plot against
a new feature,
 x_2 (or try/create
a new feature)

Example: Monitoring computers in a data center (118)

x_1 = memory use of a computer

x_2 = number of disk accesses/sec

x_3 = CPU load

x_4 = network traffic.

- choose features that might take on unusually large or small values in the event of an anomaly.

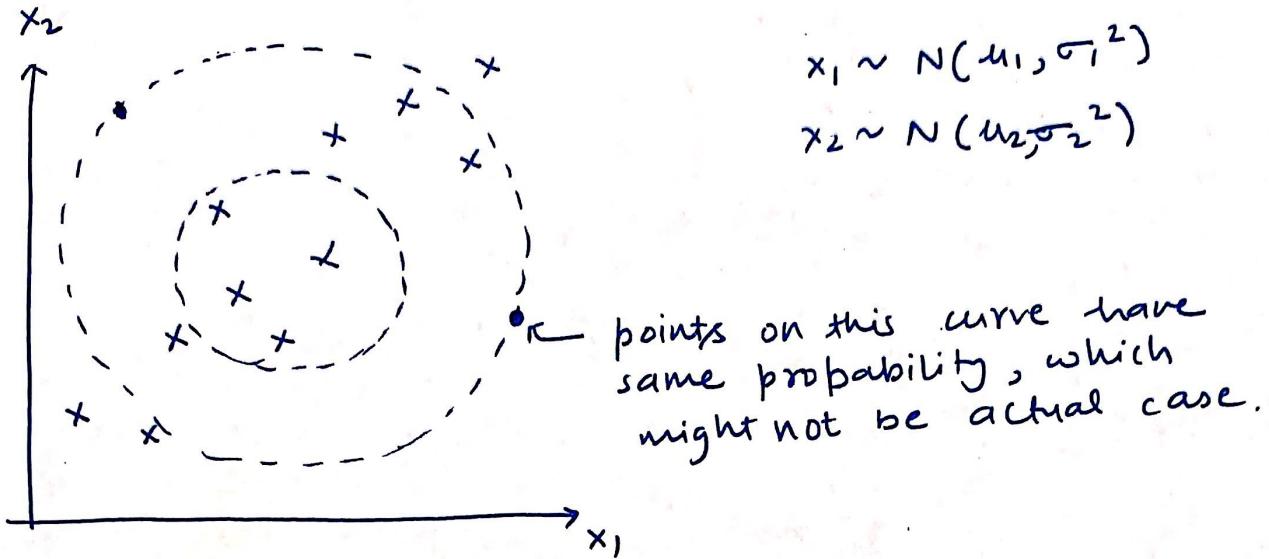
$$x_5 = \frac{\text{CPU load}}{\text{N/W traffic}}$$

$$x_6 = \frac{(\text{CPU load})^2}{\text{N/W traffic}}$$

#

Multi-variate Gaussian Distribution Anomaly Detection

(119)



$\therefore x \in \mathbb{R}^n$. Don't model $p(x_1), p(x_2) \dots$ etc. separately
 Model $p(x)$ all in one go

Parameters: $\mu \in \mathbb{R}^n$

$\Sigma \in \mathbb{R}^{n \times n}$: Covariance matrix

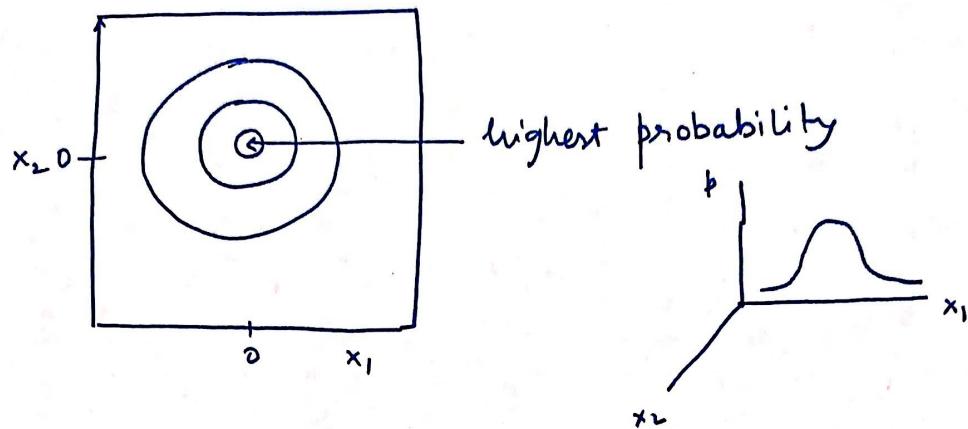
$$\therefore p(x) \Rightarrow p(x; \mu; \Sigma)$$

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

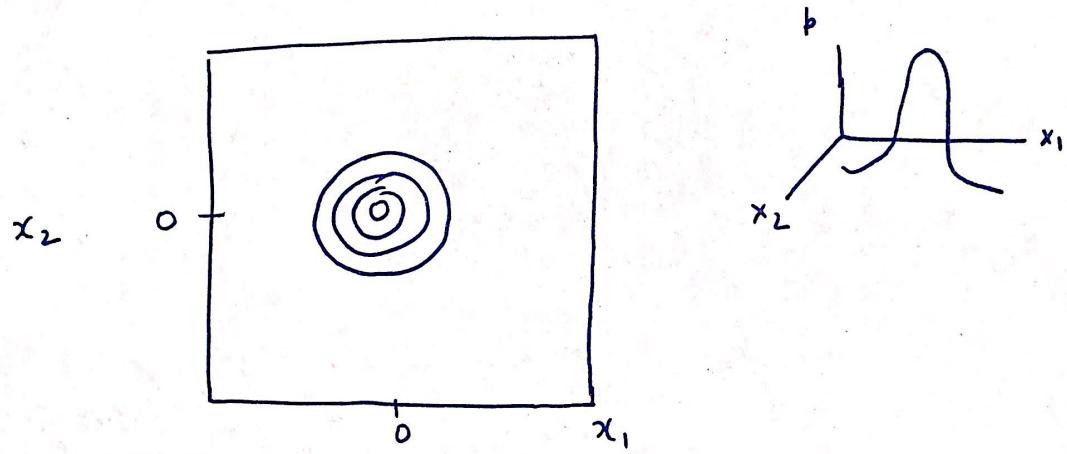
$$|\Sigma| = \text{determinant}$$

Example:

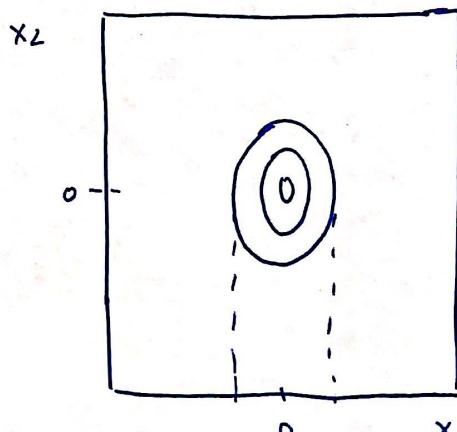
$$\textcircled{1} \quad \mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



$$\textcircled{2} \quad \mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} 0.6 & 0 \\ 0 & 0.6 \end{bmatrix}$$



$$\textcircled{3} \quad \mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} 0.6 & 0 \\ 0 & 1 \end{bmatrix}$$

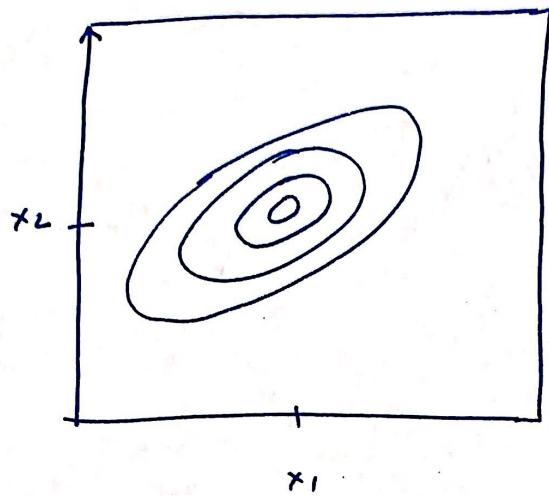


(121)

$$④ \quad u = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

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

high/medium/ low (+)vely correlation

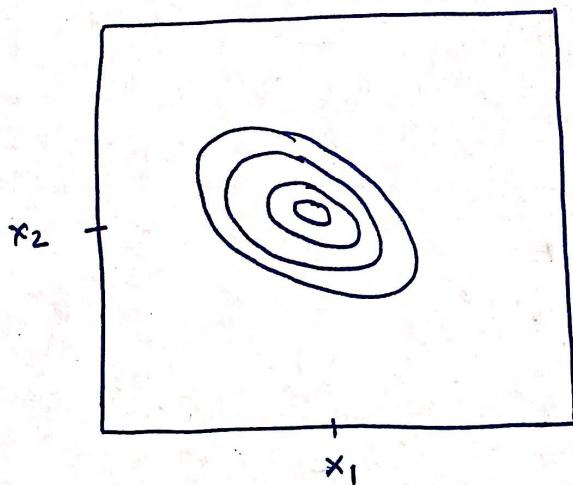


(5)

$$u = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

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

-vely correlated



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

parameters : μ, Σ

Parameter fitting:

①

Given training set : $\{x^1, x^2, \dots, x^m\}$

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

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

Fix model $p(x)$

② Given a new example x , compute

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

flag an anomaly if $p(x) < \epsilon$

- Relationship to Original Model

(123)

Original model:

$$p(x) = p(x_1; u_1, \sigma_1^2) \times p(x_2; u_2, \sigma_2^2) \times \dots \times p(x_n; u_n, \sigma_n^2)$$

corresponds to Multi-variate Gaussian model:

$$p(x, u, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp(-\frac{1}{2}(x-u)^\top \Sigma^{-1} (x-u))$$

s.t.

$$\Sigma = \begin{bmatrix} \sigma_1^2 & & & \\ & \sigma_2^2 & & \\ & & \ddots & \\ & & & \sigma_n^2 \end{bmatrix}$$

This implies that ~~we~~ we aren't allowed to model correlation between features.

<u>Original Model</u>	<u>vs.</u>	<u>Multi-variate</u>
→ Manually create features to capture anomalies where x_1, x_2 take unusual combination of values.		→ Automatically captures co-relations b/w features
→ computationally cheaper (alternatively, scales better to large, n)		→ computationally more expensive
→ OK even if m is small		→ Must-have <u>$m > n$</u> or else Σ is non-invertible ($m \geq 10n$)

Recommender Systems

/ Example: Predicting movie ratings

Movie	Alice(1)	Bob(2)	Carol(3)	Dave(4)
Movie 1	5	5	0	0
Movie 2	5	?	?	?
Movie 3	?	4	0	
Movie 4	0	0	5	4
Movie 5	0	0	5	?

n_u = no. of users.

n_m = no. of 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 $r(i,j) = 1$)

Problem: develop a recommender system to come up with a learning algorithm that can automatically fill in the missing values. And then try to predict what else might be interesting to a user.

Approach 1: Content-based recommender systems.

(126)

Movie	x_1	x_2	Feature vector
1	0.9	0	
2	1	0.02	
3	0.99	0	
4	0.1	2.0	
5	0	0.9	

{ Degree to which the movie is romantic } { Degree to which the movie is Action }

\times^i : feature vector for movie $i \Rightarrow x_i^i = \begin{bmatrix} 1 \\ 0.9 \\ 0 \end{bmatrix}$ $\Rightarrow x_0 = 1$

n = number of features (not counting x_0)
 $\Rightarrow n=2$

- for each user j , learn a parameter $\theta^j \in \mathbb{R}^{n+1}$
 - Predict user j 's rating for movie i with $(\theta^j)^T \cdot x^i$
- Starts.

x^3 = feature vector for movie 3.

$$= \begin{bmatrix} 1 \\ 0.99 \\ 0 \end{bmatrix}$$

$$\theta^1 = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix}$$

(somehow got parameter θ^j
for user j ; explored later)

$$\therefore (\theta^1)^T x^3 = 4.95$$

Problem formulation
(sort of linear regression)

$r(i, j) = 1$: if user j has rated movie i

$y(i, j)$ = rating by user j for movie i (if defined)

θ^j = parameter vector for user j

x^i = ~~parameter~~ feature vector for movie i

m^j = no. of movies rated by user j

For user j predict rating for movie i : $(\theta^j)^T x^i$

- Optimization Objective

④ To learn parameter θ^j (parameter for user j)

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

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

⑤ To learn $\theta^1, \theta^2, \theta^3 \dots \theta^{n_u}$

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

- **Approach 2**: Collaborative filtering

Movie	x_1	x_2
1	?	?
2	?	?
3	?	?
4	?	?
5	?	?

→ We do not have the feature vector for each movie

and, we have gotten our users to provide us the parameter vector θ_j [based on their interest]

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

so, we are trying to find out x^i for user j

$$(\theta^j)^T (x^i) = r(i, j)$$

ex. $(\theta^1)^T x^1 \approx 5$

[Because Alice rated movie 1 to be 5]

Optimization Algorithm:

→ Given $\theta^1, \dots, \theta^{n_u}$, to learn x^i :

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

prediction actual
 min the error

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

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

~~Given~~ Guess: $\theta \rightarrow x \rightarrow \theta \rightarrow x \rightarrow \dots$

New optimization Objective :

Minimize $\{x^1, \dots, x^{nm}\}$ and $\{\theta^1, \dots, \theta^{nu}\}$

simultaneously :

$$\begin{aligned} x &\in \mathbb{R}^n \\ \theta &\in \mathbb{R}^n \end{aligned}$$

$$\Rightarrow \left[J(\theta^1, \dots, \theta^{nu}, x^1, \dots, x^{nm}) \right. \\ = \frac{1}{2} \sum_{i,j: r(i,j)=1} ((\theta^j)^T x^i - y^{ij})^2 + \frac{\lambda}{2} \sum_{j=1}^{nm} \sum_{k=1}^{n_u} (\theta_k^j)^2 \\ \left. + \frac{\lambda}{2} \sum_{i=1}^{nm} \sum_{k=1}^n (x_{ik}^i)^2 \right]$$

$$\Rightarrow \left[\min_{\substack{x^1, \dots, x^{nm} \\ \theta^1, \dots, \theta^{nu}}} J(x^1, \dots, x^{nm}, \theta^1, \dots, \theta^{nu}) \right]$$

→ Steps for Collaborative filtering Algorithm

①: Initialize $x, \theta \in \mathbb{R}^n$ to small random values.

② Min. $J(\dots)$ using gradient descent (or any advanced optimization algo). $\forall i, j$

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

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

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

$\frac{\partial J(\dots)}{\partial \theta_k^j}$

③ For a user with parameter θ & movie with feature x , predict a star rating of $\theta^T x$

- Vectorization : low rank matrix factorization
(collaborative filtering)

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

Predicted ratings: $(\Theta^j)^T x^i$

$$\begin{bmatrix} (\Theta^1)^T x^1 & (\Theta^2)^T x^1 & \dots & (\Theta^{n_u})^T x^1 \\ (\Theta^1)^T x^2 \\ \vdots \\ (\Theta^1)^T x^{n_m} \end{bmatrix}_{n_m \times n_u} = x^* \Theta^T$$

$$x = \begin{bmatrix} -(x^1)^T \\ -(x^2)^T \\ \vdots \\ -(x^{n_m})^T \end{bmatrix}_{(n_m \times 1)} ; \quad \Theta = \begin{bmatrix} -(\Theta^1)^T \\ -(\Theta^2)^T \\ \vdots \\ -(\Theta^{n_u})^T \end{bmatrix}_{(n_u \times 1)}$$

• Application: Recommender System

(134)

for each product i / movie, we learn a feature vector $x^i \in \mathbb{R}^{nm}$

→ How to find movies j related to movie i?

small $\|x^i - x^j\|$: movie j and i are similar

→ 5 most similar movies to movie i

↳ Find the 5 movies j with the smallest $\|x^j - x^i\|$

⇒ Users who have not rated any movies!

$$J(\theta, \dots, x) = \min_{\theta, x} \frac{\lambda}{2} \left(\sum_k x_k^i - \sum_k \theta_k^j \right)^2$$

$$\theta^5 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\leftarrow \frac{\lambda}{2} \left[\theta_1^5 + \theta_2^5 \right]^2$$

$$\therefore (\theta^5)^T x^i = 0.$$

Q How to solve this problem?

↳ Mean Normalization

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 & 2 \\ 5 & ? & ? & 0 & ? \\ 5 & ? & 0 & ? & ? \\ ? & 4 & 0 & ? & ? \\ 0 & 0 & 5 & 4 & ? \\ 0 & 0 & 5 & 0 & ? \end{bmatrix}; \quad \bar{Y} = \begin{bmatrix} 2.5 \\ 2.5 \\ 2 \\ 2.25 \\ 1.25 \end{bmatrix} \quad \begin{array}{l} 10/4 \\ 5/2 \\ 4/2 \\ 9/4 \\ 5/4 \end{array}$$

$$\therefore Y - \bar{Y} = \begin{bmatrix} 2.5 & 2.5 & -2.5 & -2.5 & ? \\ 2.5 & ? & ? & -2.5 & ? \\ ? & 2 & -2 & ? & ? \\ -2.5 & -2.25 & 2.75 & 1.75 & ? \\ -1.25 & -1.25 & 3.75 & -1.25 & ? \end{bmatrix} \quad \text{Essentially making mean of all movies = 0}$$

↓
learn θ^j, x^i from this data

for user j , on movie i predict

$$(\theta^j)^T x^i + \mu_i$$

for user 5 (Eve):

$$\theta^5 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\therefore \underbrace{(\theta^5)^T x^i}_{0} + \mu^i$$

* In above case: we have normalized row: user not having rated movie.

* Alternatively, if we have a movie with no ratings
we can normalize column: