

Density Estimation

-W9 - Anomaly Detection-

Intro:

- Anomaly detection is widely used in fraud detection (e.g. 'has this credit card been stolen?'). Given a large number of data points we may sometimes want to figure out which ones very significantly from the average.
- For example we may want to detect defects or anomalies. We show how a dataset can be modeled using a Gaussian distribution, and how the model can be used for anomaly detection.

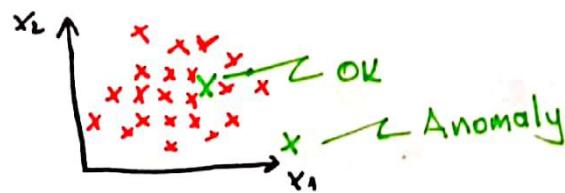
Anomaly Detection Example:

- Imagine that you are a manufacturer of aircraft engines, and you are doing some quality assurance testing and you measure some features like:

- x_1 = heat generated
- x_2 = vibration intensity

we have a dataset:

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

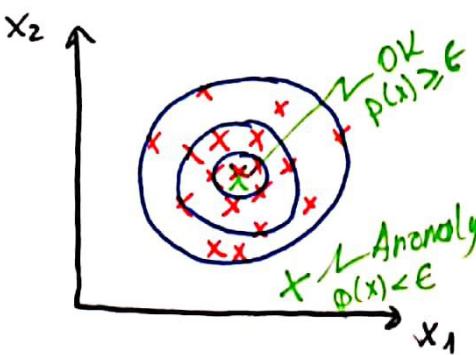


- Let's say the data-set looks like:

- Now let's say you have a new aircraft engine and it has some features: x_{test} . We want to know if this new engine is anomalous in any way. We want to know if it is an okay engine or maybe there is a problem

Density Estimation:

- For an anomaly detection problem we are given some dataset: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ and we usually assume that these m examples are normal or non-anomalous examples and we want an algorithm to tell us if some new example x_{test} is anomalous



- The approach that we are gonna take is that, given this training set (unlabeled) we are going to build a model for $p(x)$.

- Build a model for the probability of x where x are these features x_1, x_2

- Then if $p(x_{\text{test}}) < E$ flag x_{test} as an anomaly!

- If $p(x_{\text{test}}) \geq E \rightarrow \text{OK}$

Some Application Examples of Anomaly Detection:

• Fraud Detection

- $x^{(i)}$ = features of user i's activities.
- Model $p(x)$ from data.
- Identify unusual users by checking which have $p(x) < \epsilon$
- Some example features may be: x_1 = how often does user log in,
 x_2 = number of what pages visited, x_3 = number of transactions
 x_4 = number of posts, x_5 = typing speed of the user, etc ...

• Manufacturing

- Monitoring computers in a data center and detect anomaly
- $x^{(i)}$ = features of each machine
 - features might be memory use, number of disk accesses/sec, CPU load, CPU load/network traffic.
- Model $p(x)$
 - if for a test computer $p(x_{test}) \leq \epsilon$: there might be something wrong with the test computer.

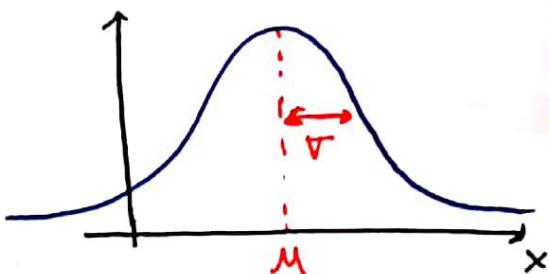
Gaussian (Normal) Distribution:

- Say x is a real valued random variable; $x \in \mathbb{R}$. If the probability distribution of x is gaussian with mean μ and variance σ^2 then it is shown as:

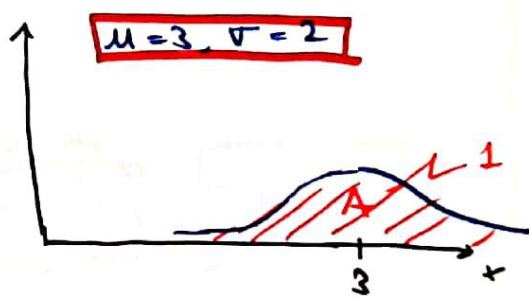
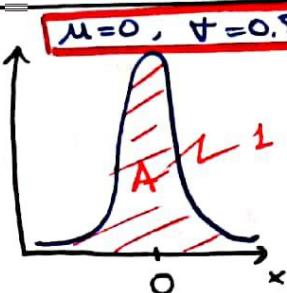
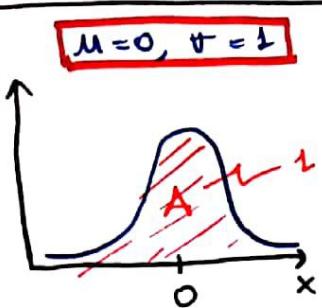
- $$x \sim N(\mu, \sigma^2)$$

distributed as

- If we plot the Gaussian Probability Density Function:

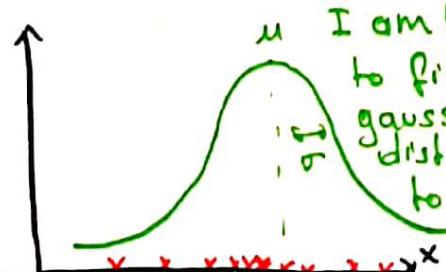


$$p(x) = p(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi} \cdot \sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

Some Gaussian Distribution Examples:Parameter Estimation Problem:

- Let's say we have a dataset of m examples and each of these examples is a real number:

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



- Let's say I suspect that these examples came from a Gaussian distribution

- I suspect each of my examples

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

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

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

- The problem of parameter estimation given my data set. I want to try to figure out what are the values of μ and σ^2 .

→ we use MLE
→ it is an approximator.

Density Estimation

-W9 - Anomaly Detection Algorithm-

4

- Let's say we an unlabeled training set: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$
- Each example is $x \in \mathbb{R}^n$
- We need to model the $p(x)$ where x is a vector $\in \mathbb{R}^n$

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

- we will assume each feature is distributed according to a Gaussian Distribution:

- $x_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$
- $x_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$
- $x_n \sim \mathcal{N}(\mu_n, \sigma_n^2)$

Note: The equation above corresponds to an independence assumption on the features x_1 through x_n . But in practice it turns out the algorithm which will be described works just fine whether or not these features are independent.

- $p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2)$ just a compact form of the above eq.

Put Everything Together: Anomaly Detection Algorithm

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

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

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

Find μ_1, μ_2, \dots

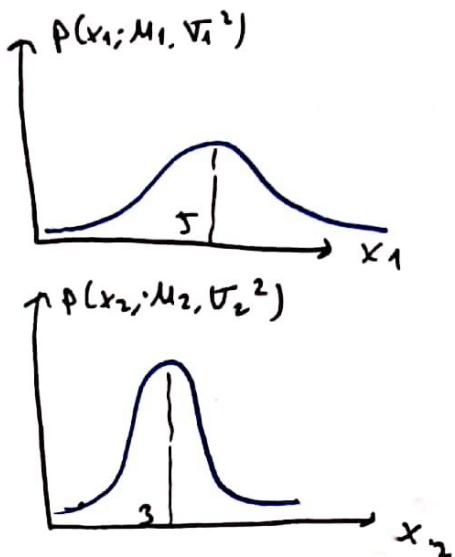
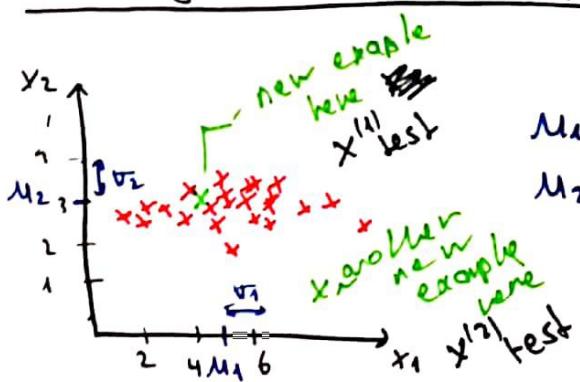
- $\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$

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

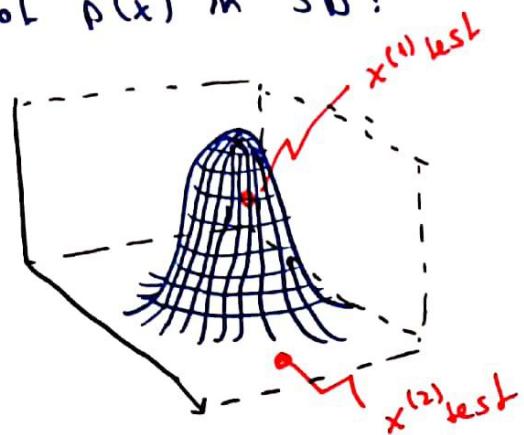
3) Given new example x , compute $p(x)$:

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

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

Anomaly Detection Example

- Plot $p(x)$ in 3D:



- Let's choose

$$\epsilon = 0.02$$

- Later we'll see how to choose ϵ as well.

- $p(x^{(1)}_{\text{test}}) = 0.0426 \geq \epsilon$ $x^{(1)}_{\text{test}}$ is not an ANOMALY
- $p(x^{(2)}_{\text{test}}) = 0.0021 < \epsilon$ $x^{(2)}_{\text{test}}$ is a ANOMALY

-W9 - Building an Anomaly Detection System -

2

Intro

- We will talk about how to develop a specific application of anomaly detection to a problem. In particular we will focus on the problem of how to evaluate an anomaly detection algorithm.

The Importance of real-number evaluation

- We've already talked about 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.
- Thus, in order to be able to develop an anomaly detection system quickly, it would be really helpful to have a way of evaluating an anomaly detection system.
- In order to evaluate the A.D. System we need a labeled data of anomalous and non-anomalous examples. ($y=0$ if normal, $y=1$ if anomalous).
- We think our training set: $x^{(1)}, x^{(2)}, \dots, x^{(m)}$ as a ~~one~~ large collection of normal / non-anomalous examples. But it's okay if there are few anomalies slip into your unlabeled training set.
- Next we are going to define our Cross-Validation and Test sets:
 - Cross Validation Set: $(x_{cv}^{(1)}, y_{cv}^{(1)}), \dots, (x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$
 - Test Set: $(x_{test}^{(1)}, y_{test}^{(1)}), \dots, (x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$
 - For both CV and Test sets we include some examples that are known to be anomalous

Aircraft Engines Motivating Example

- Let's say we have a dataset like:
 - 10 000-as far as we know- good (normal) engines
 - 20 flawed engines (anomalous) ~ Typically b/w 2-50
- Given this dataset a typical way to split it into training - CV and test sets:
 - 6000 ^(unlabeled) good engines \leftarrow Training Set \rightarrow Used for fitting $p(x)$
 - 2000 good engines ($y=0$), 10 anomalous ($y=1$) \leftarrow CV set
 - 2000 good engines ($y=0$), 10 anomalous ($y=1$) \leftarrow Test Set
- So we are using Training Set to fit the $p(x)$. Then we can use CV to select features or select the f , etc.. Finally after we fit our complete model we can test the performance of the final model on the test set.

~ it is okay if there are couple of bad engines inside 10 000 since it is not going to effect the $p(x)$

Algorithm Evaluation

- Fit model $p(x)$ on training set $\{x^{(1)}, \dots, x^{(m)}\}$
- On cross-validation on test sets, 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:

~~✗~~ Since the data is skewed, because $y=0$ is much more common compared to $y=1$. Classification accuracy would not be a good evaluation metric.

- Instead we can use evaluation metrics like:

- True positive, false positive, false negative, true negative
- Precision/Recall
- F₁-score

→ We can use CV set to choose ε . We can try many different values of ε and pick the value that maximizes F₁ score on the CV set.

→ We make decisions like "what features to use", "what value of ε should I use" by evaluating the algorithm on the CV set. When we find the final model we can evaluate it on the Test Set.

WG - Anomaly Detection vs Supervised Learning

3

- In the last part we talked about the process of evaluating an anomaly detection algorithm. There, we used some labeled data with examples that we knew were either anomalous or not anomalous ($y=0$ or $y=1$)

- So the question then arises: If we have a labeled data why don't we just use a supervised learning algorithm? why don't we just use Logistic Regression or a Neural Network to try to learn directly from the labeled data?

Some Guidelines for when we should probably use an A.D. algorithm

Anomaly Detection

vs

Supervised Learning

- If we have very small number positive examples ($y=1$). (0-20 is common).

- We will have relatively large number of negative ($y=0$) examples
↳ use it to fit $p(x)$

- There might be 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.

- For the spam classifier example, there were many different "types" of anomalies (spam mails) but we had enough positive examples (spam mails) so we've used a supervised learning.

- Typically we have 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 the training set.

Note: Assume we've used a supervised learning alg. instead of A.D. Since $y=1$ (blue) examples are only three there is not enough info. what if a new anomaly occurs like the green one. Then supervised model is not going to work properly.

Application Areas:

Anomaly Detection

Supervised Learning

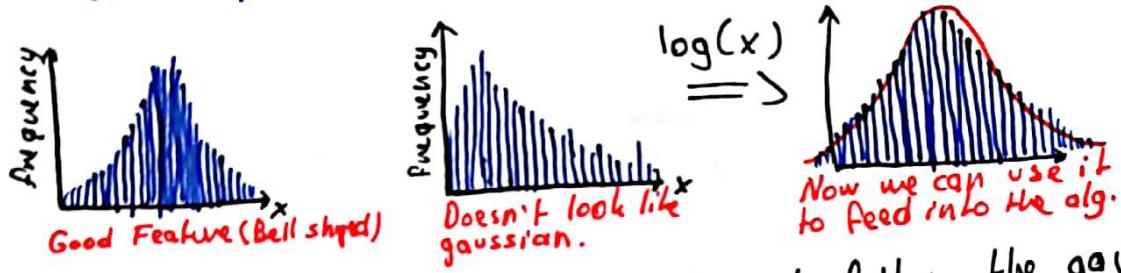
- Fraud detection
- Manufacturing
- Monitoring machines in a data center
- Email spam classification
- Weather prediction
- Cancer classification

W9 - Choosing What Features to Use -

~~(*)~~ When we are applying anomaly detection, one of the things that has a huge effect on how well it does, is what features we choose.

Guidelines for Feature Selection:

→ One thing we did was to model the features using gaussian distribution, $p(x; \mu, \sigma^2)$. By plotting the histogram of every feature we can be sure that the distribution of the data follows the gaussian distribution. ~~(*)~~ It usually works okay even if it's not gaussian but this is a better practice. → Use hist command on matlab's



~~(*)~~ So if I have some features looks like doesn't follow the gaussian dist. I can manipulate these features and obtain new ones then I can use the new features for my Anomaly Detection Algorithm.
 play with c to make it look more gaussian ↪ $x_1 \leftarrow \log(x_1 + c)$ ↪ $x_2 \leftarrow \log(x_2 + c)$ ↪ $x_3 \leftarrow x_3$

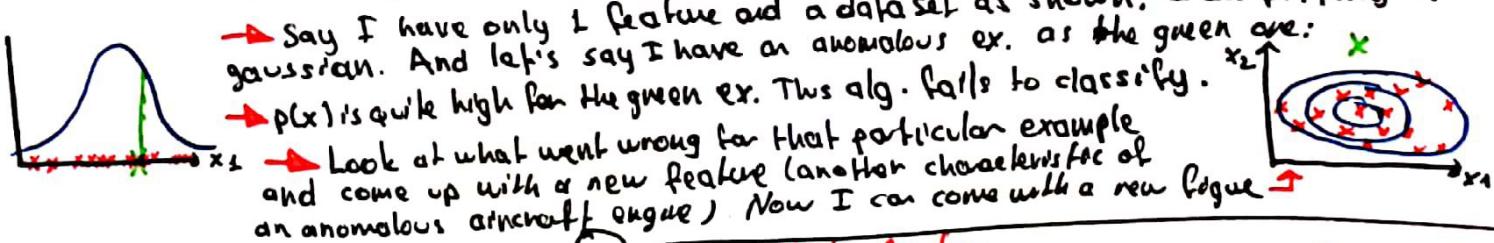
How Do we Come Up With Features for an Anomaly Detection Algorithm

→ What we want is large $p(x)$ values for normal examples x . small $p(x)$ values for anomalous examples x .

~~(*)~~ The most common problem is: $p(x)$ is comparable (say, both large) for normal and anomalous examples.

~~(*)~~ Here we can benefit from the Error Analysis procedure. It's very similar to the error analysis procedure we have seen for the supervised learning.

→ Train a complete algorithm. Run it on the CV set and look at the example it gets wrong and see if we can come up with extra features to help the algorithm do better on the examples that it got wrong.

Choosing Features for Monitoring Computers in a Data Center

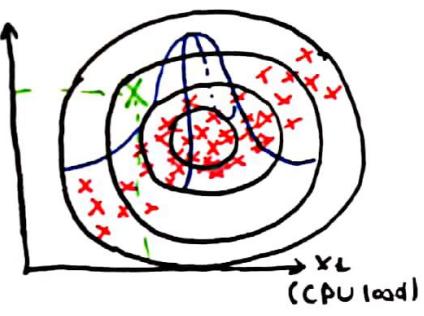
~~(*)~~ Choose features that might take on unusually large or small values in the event of anomaly,
 • memory use of computer (x_1) • x_2 = number of disk accessed • x_3 = CPU load • x_4 = Network Traffic
~~(*)~~ Maybe x_3 and x_4 are linearly growth with each other then add another feature to distinguish anomaly
 • x_5 = CPU load / Network Traffic or • x_6 = $(\text{CPU load})^2 / \text{Network Traffic}$

1. K. Lantuejoul
-wg - Multivariate Gaussian Distribution -
 - OPTIONAL -

→ We'll talk about one possible extension of the Anomaly Detection Algorithm we have developed so far. This extension uses multi-variate gaussian distribution and it has some adv.s and disadv.s. It can sometimes catch some anomalies that the earlier algorithm didn't.

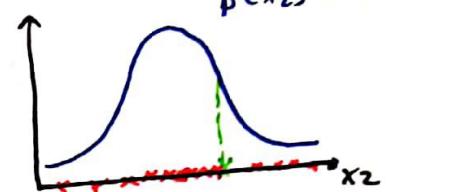
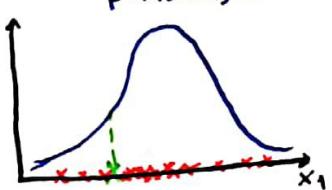
Motivating Example: Monitoring Machines in a Data Center

→ Let's say I have my unlabeled data as plotted below:



→ If I try to fit a Gaussian $p(x)$ to this dataset:

$$p(x_1; \mu_1, \Sigma_1^2)$$



$$p(x) \text{ will be } p(x_1) \cdot p(x_2).$$

Thus by defining $p(x) = p(x_1) \cdot p(x_2)$, the Gaussian will always be ~~axis-aligned~~ but not ellipses. ~~Thus~~ we can not model the data as we want.

For ex. the green example is actually an anomaly but since the gaussian is ~~axis-aligned~~ our model fails to classify it as an anomaly.

To solve this problem and shape the gaussians more flexible way we will use Multivariate Gaussian (Normal) Distribution.

Multivariate Gaussian (Normal) Distribution

For features $x \in \mathbb{R}^n$.

Instead of modeling $p(x) = p(x_1) \cdot p(x_2) \dots$ separately. We will model $p(x)$ all in one go.

For this new modeling the parameters will be:

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

mean vector covariance matrix

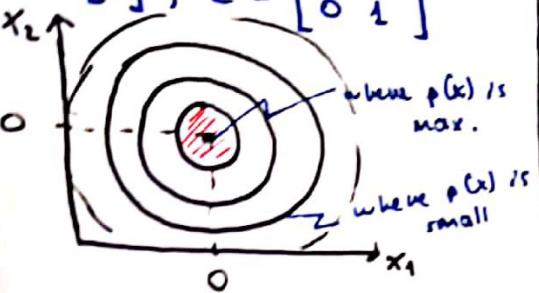
$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{n/2} \cdot |\Sigma|^{1/2}} \cdot \exp \left(-\frac{1}{2} (x - \mu)^T \cdot \Sigma^{-1} \cdot (x - \mu) \right)$$

$|\Sigma|$: Determinant of Σ

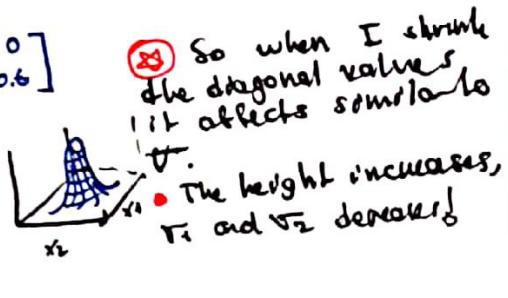
Now we know how to define a Multivariate Gaussian Distribution $P(x; \Sigma, \mu)$. Let's look at how do they look for different values of μ and Σ .

Multivariate Gaussian Examples

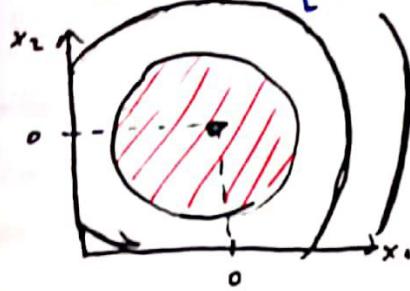
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



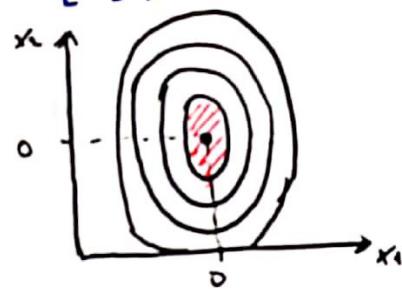
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 0.4 & 0 \\ 0 & 0.6 \end{bmatrix}$$



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



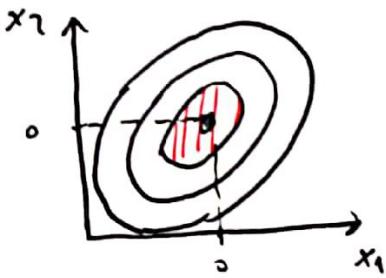
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 0.6 & 0 \\ 0 & 1 \end{bmatrix}$$



What happens if we change the off-diagonal values: We can model the correlated data.

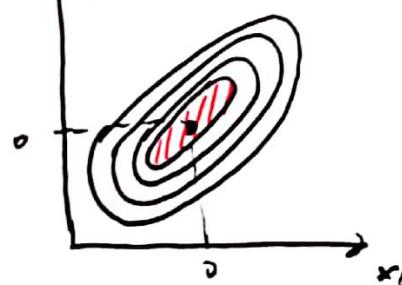
If we increase the off-diagonal values it will affect the shape like we are pulling the gaussian from the ends in $x=y$ line.

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



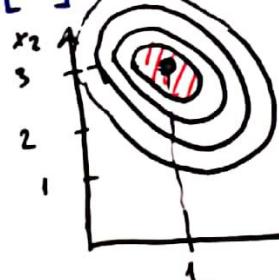
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & 0.7 \\ 0.9 & 1 \end{bmatrix}$$

Off-diagonals must be equal to the same value



$$\Sigma = \begin{bmatrix} \sigma_{x_1}^2 & \sigma_{x_1} \cdot \sigma_{x_2} \\ \sigma_{x_2} \cdot \sigma_{x_1} & \sigma_{x_2}^2 \end{bmatrix}$$

$$\mu = \begin{bmatrix} 1 \\ 3 \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & -0.7 \\ -0.7 & 1 \end{bmatrix}$$



If off-diags are negative then we are pulling in the $x=y$ line

Recap Multivariate Gaussian (Normal) Distribution

Parameters: μ, Σ where $\mu \in \mathbb{R}^n$: mean vector, $\Sigma \in \mathbb{R}^{n \times n}$: covariance matrix

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

Parameter Estimation Problem

- Given training set: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ where $x^{(i)} \in \mathbb{R}^n$

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

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

→ Just like the sigma (Σ) in the PCA part.

Put All This Together: Anomaly Detection with the Multivariate Gaussian:

- 1) Fit model $p(x)$ by setting μ and Σ using the training set.
- 2) When given a new example x_{new} , compute $p(x_{\text{new}})$. Flag an anomaly if $p(x) < \epsilon$

Relationship Between Multivariate Gaussian Dist. Model and the Original Model

→ Original model was: $p(x) = p(x_1; \mu_1, \sigma_1^2) \cdot p(x_2; \mu_2, \sigma_2^2) \cdot \dots \cdot p(x_n; \mu_n, \sigma_n^2)$

→ It turned out that the original model is a special case of Multivariate Gaussian when off-diagonals are always zero! So for the original model the contours of the gaussian is always axis-aligned

Axis Aligned:



Not Axis-Aligned:



→ Original model is the same as a Multivariate Gaus. Model but with a constraint and the constraint is that: $\Sigma = \begin{bmatrix} \sigma_1^2 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_n^2 \end{bmatrix}$

When Do we use Original Model and When Gaussian Mixture Model:

Original Model (^{used more often}) vs

Multivariate Gaussian Model

→ If anomaly occurs with the unusual combination of features like the example (2 pages before)

→ We can use original model but we need to create (manually) features to capture anomalies. Like in the example of CPU load (3 pages before)

→ Computationally cheaper (alternatively scales better to large n)

→ Works okay even if m is small.

→ Automatically captures correlations between features.

→ Computationally more expensive

→ Let's say $n = 1000000$ original model will work OK. But for this one $\Sigma \in \mathbb{R}^{n \times n}$ computing Σ^{-1} will be comp. exp.

→ Must have $m > n$, or else Σ is non-invertible.

→ In practice use Mult. G.M. when $m \geq 10 \cdot n$

→ If Σ is non-invertible either $m > n$ or we have redundant features: $x_3 = x_4 + x_5$ get rid of x_4 and x_5

-WG - Part 2 - Recommender Systems

→ Recommender systems look at patterns of activities between different users and different products to produce recommendations. In this 2nd module of WG, we introduce recommender algorithms such as the collaborative filtering algorithm and low-rank matrix factorization.

→ Applications on websites like Amazon, Netflix, Ebay etc. looks at what movies/products you've rated/purchased in the past and recommend new movies/products.

✖✖✖ During this module we will also talk about something quite important. We've already seen that features have a big effect on the performance of the learning algorithm. ⇒ For some problems there are algorithms that can try to automatically learn a good set of features for you. Recommender systems are just one example of these problems.

Example: Predicting Movie Ratings

→ Let's say you are a company that sells movies and you let your users to rate your movies from 0 stars to 5 stars

→ Assume I have a data set as below:

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	?
Car Chases	0	0	5	4
Swords vs. Karate	0	/	5	?

↳ Looks like Alice and Bob like romance

↳ Looks like Carol and Dave likes action movies

→ Let's introduce a new notation:

• n_u : # of users (4 here) • n_m : # of movies (5 here)

• $r(i,j)$: If user j has rated movie i then it's equal to 1. o/w 0.

• $y(i,j)$: Rating given by user j to movie i (defined only if $r(i,j)=1$)

→ Our purpose is to guess the values of ? data.

→ Here I have a small number of users and movies and only a small portion of the movies are unrated (?) for a user. In reality this fraction will be much bigger.

→ Our first approach to building a recommender system will be the content based recommendations.

Content Based Recommender Systems

Movie	Alice(1)	Bob(2)	Carol(3)	Dave(4)	x_1 (romance)	x_2 (action)	
Love	5	5	0	0	0.9	0	• $n_u = 4$
Romance	5	?	?	0	1.0	0.02	• $n_m = 5$
Cute Love	?	4	0	?	0.99	0	• $n = 2$
Car Chase	0	0	5	?	0.1	1.0	
Karate	0	0	5	?	0	0.9	

→ Let's say I have 2 features for each movie. Now for each user we consider a linear regression problem. For each user we have some data set and we need to find a set of parameters θ that ~~line fitting~~ $\theta_0 \cdot x_0 + \theta_1 \cdot x_1 + \theta_2 \cdot x_2$ fits our data set.

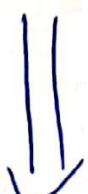
→ Each movie can be represented with a feature vector. Just like houses.
• The first movie: Love can be represented as $x^{(1)} = \begin{bmatrix} x_0^{(1)} \\ x_1^{(1)} \\ x_2^{(1)} \end{bmatrix} = \begin{bmatrix} 1 \\ 0.9 \\ 0 \end{bmatrix}$

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

→ Dijelim ki Alice'in θ parametresini yani $\theta^{(1)}$ 'i elde ettik. (Bunu nasıl elde edeceğimizi sonra göreceğiz.) Elde edilen parametre yardımıyla da 3. film olan Cute Love'in Alice'den kaç puan alacağını bulmak istiyoruz:

$$\bullet x^{(3)} = \begin{bmatrix} 1 \\ 0.99 \\ 0 \end{bmatrix} \quad \bullet \theta^{(1)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix} \quad \Rightarrow (\theta^{(1)})^T \cdot x^{(3)} = 5 \cdot 0.99 = 4.95 //$$

→ Özette olarak movie rating'i tahmin etme problemimizi linear regression problemine çevirdik. Her filmi 2 feature ile represente ediyoruz. Her user'in aynı bir linear regression problemi olduğunu düşün. Line fitting yapacağımız ama her user'in dataset farklı, dolayısıyla farklı bir line oluşturmalıyız.



Devam ediyor...

WS - Part 2 - Recommender Systems III -

Problem Formulation:

- $r(i, j) = 1$ if the user has rated movie i (0 otherwise)
- $y(i, j)$ = rating by user j on movie i (if defined)
- $\theta^{(j)}$ = parameter vector of user j
- $x^{(i)}$ = feature vector of movie i $\rightarrow (\theta^{(j)})^T \cdot x^{(i)}$: Rate for movie i , user j .
- $m^{(j)}$ = number of movies rated by user j .

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

(4): TO LEARN $\theta^{(j)}$: sum over every rated movie

minimum cost function just like in the part of linear regression

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

reg. term

→ To simplify the eq. I will get rid of the $m^{(j)}$ terms. The result will be same.

Optimization Objective:

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

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

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

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

$J(\theta^{(1)}, \dots, \theta^{(n)})$

We can minimize the Cost Function J by Gradient Descent:

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

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

$\underbrace{\frac{\partial}{\partial \theta_k^{(j)}} J(\theta^{(1)}, \dots, \theta^{(n)})}_{\theta_k^{(j)}}$

→ We can use some advanced opt. alg.s instead of Grad. Desc.

→ This approach is called content based approach because we assume that we have features available to us. Features capture what is the content of these movies.

W9 - Collaborative Filtering -

- The algorithm that we are going to talk about has very interesting property that it does: feature learning.
- Dataset'i hatırlarsan her filmi x_1 ve x_2 feature ile temsil etmişlik. Ayrıca bu feature'ların her film için değerlerini biliyoruz kabul ettik. Ancak gerçekle bu feature'ları bilmemiz zor. Birilerinin filmleri tek tek izleyip feature'ları öğrenmesi genellikle bu da zor ve zaman alıcı. O zaman data set şöyle olsun:

Movie	Alice(1) $\theta^{(1)}$	Bob(2) $\theta^{(2)}$	Carol(3) $\theta^{(3)}$	Dave(4) $\theta^{(4)}$	x_1 (romance)	x_2 (action)
Love	5	5	0	0	?	?
Romance	5	?	?	0	?	?
Cute Love	?	4	0	?	?	?
Car Chase	0	0	5	4	?	?
Karate	0	0	5	?	?	?

- Let's assume each user has told us how much they like romantic movies and how much they like action ~~not~~ movies:

• $\theta^{(1)} = [0 \ 5 \ 0]^T$ • $\theta^{(2)} = [0 \ 5 \ 0]^T$ • $\theta^{(3)} = [0 \ 0 \ 5]^T$ • $\theta^{(4)} = [0 \ 0 \ 5]^T$

• Carol for ex. likes action movies thus $\theta_3^{(2)} = 5$ and doesn't like romance.

→ Bu durumda problem verilen θ 'ları kullanarak $x^{(1)}, x^{(2)}, \dots, x^{(nm)}$, bulmak oluyor. Mesela 1. filme Alice ve Bob 5 vermiş o zaman $x_1^{(1)}$ 1 ~~not~~ çünkü Alice ve Bob Romance çok seviyor. Matematiksel olarak yazarsak:
 \hookrightarrow Ana olasılıkta bilir yani romance'i 5'lik seviyor ama bu filme 3 vermiş olabilir.

→ What feature vector should $x^{(1)}$ be so that $(\theta^{(1)})^T \cdot x^{(1)} \approx 5$ and $(\theta^{(2)})^T \cdot x^{(1)} \approx 5$ and $(\theta^{(3)})^T \cdot x^{(1)} \approx 0$ and $(\theta^{(4)})^T \cdot x^{(1)} \approx 0$. $x = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$

Optimization Algorithm: Let's formalize above problem:

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

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

\hookrightarrow choose $x^{(1)}$ to minimize error.

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

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

W9 - Collaborative Filtering - II

→ Putting everything together: the alg. we talked about 2 pages ago and the algorithm we talked 2 page ago:

Collaborative Filtering

→ Previously we have seen: Given $x^{(1)}, \dots, x^{(nm)}$ (and movie ratings) we can estimate $\Theta^{(1)}, \dots, \Theta^{(nu)}$

→ What we've also seen: Given $\Theta^{(1)}, \dots, \Theta^{(nu)}$ (and movie ratings) we can estimate $x^{(1)}, \dots, x^{(nm)}$

★ Here is how collaborative filtering works:

- ① Randomly guess Θ
- ② Based on guessed Θ estimate x values for different movies
- ③ Then using these x features estimate Θ (a better one)
- ④ Then using these x features estimate Θ (a better one)
- ⑤ So on... Keep iterating $\Theta \rightarrow x \rightarrow \Theta \rightarrow x \rightarrow \Theta \rightarrow x \dots$

→ For this problem this process is possible when each user makes multiple movies and each movie is rated by multiple users. So we can do this back and forth process.

★ This is just an initial collaborative filtering algorithm. Daha sonra data efficient bir versiyonunu görelim. Collaborative demetinin sebebi bir kullanıcıyı bir kişiin评级in kullanarak herkesin电影i öneLEYEBİLİR. Every user is helping the algorithm to learn better features and these features will be used to make better predictions for everyone else.

→ Özetle: Feature'ları bilişik sadece bir satırda tahmin yapabiliyoruz. Θ değerlerini linear reg. ile elde edebiliyoruz. Birer satırda Θ 'ları bilişikte sadece bir satırda tahmin yapabiliyoruz. x değerleri elde edilebilir. Yani bir kullanıcının oyladığı filmlere bakarak yeri bir filme kaç puan vererğini veya bir film için hangi kullanıcılardan ne oy verdigine bakarak yeri bir kullanıcının ne oy vereceğini tahmin ediyoruz. Random Θ ile başlayıp $x \rightarrow \Theta \rightarrow x \rightarrow \Theta$ pattern'i izleyerek Θ re x optimise edebiliyoruz.

W9 - Collaborative Filtering Alg. III -

→ By putting previous ideas together we'll come up with a collaborative filtering algorithm:

Collaborative Filtering Optimization Algorithm

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

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

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

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

→ Collaborative Filtering Alg. ıçın swaklık θ ve x arasında git-gel yapabileceğimizi görmüştük böylece random bir θ ile başlayıp θ ve x 'i optimize edebiliyoruz.

→ It turns out that there is a more efficient algorithm that doesn't need to go back and forth between x 's and the θ 's but that can solve for x and θ simultaneously.

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

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

$$\min_{x^{(1)}, \dots, x^{(n)}, \theta^{(1)}, \dots, \theta^{(n)}} J(x^{(1)}, \dots, x^{(n)}, \theta^{(1)}, \dots, \theta^{(n)})$$

→ Instead of sequentially going back and forth we are just minimizing the cost function w.r.t. both sets of parameters simultaneously

→ Now we are not using $x_0 = 1$ convention thus $x \in \mathbb{R}^n$, $\theta \in \mathbb{R}^n$. Since we are learning our features algorithm can learn a feature = 1 always if it's needed.

✖ Aslında burada olayı söyle deşemebilim: Line fitting yaparken cost için (puan tahmini - gerçek puan)² kullanıyoruz. Puan tahmini de etkileyen x ve θ normalde sadece θ 'yı minimize ederken burada x 'yi de minimize ediyorsun. x de bir parametre gibi düşünüyor olsaydı

ikisi de aynı her naked movie user i,j patrı'nın toplama yapılıyordu. Bu i sahnelerin diğerini sunsun eklentiler.

- W9 - Collaborative Filtering Alg. IV -

Collaborative Filtering Algorithm: Guideline

- 1) Initialize $x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)}$ to small/random values.
↑ like NN algorithm.
- 2) Minimize $J(x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)})$ using Gradient Descent
(or an Adv. Opt. Alg.) E.g. for every $j=1, \dots, n_u$, $i=1, \dots, n_m$:
 - ① $x_k^{(i)} := x_k^{(i)} - \alpha \left(\sum_{j: r(i,j)=1} ((\theta^{(j)})^T \cdot x^{(i)} - y^{(i,j)}) \theta_k^{(j)} + \lambda x_k^{(i)} \right)$ $\frac{\partial}{\partial x_k^{(i)}} J \rightarrow$
 - ② $\theta_k^{(j)} := \theta_k^{(j)} - \alpha \left(\sum_{i: r(i,j)=1} ((\theta^{(j)})^T \cdot x^{(i)} - y^{(i,j)}) \cdot x_k^{(i)} + \lambda \theta_k^{(j)} \right)$ $\frac{\partial}{\partial \theta_k^{(j)}} J \rightarrow$
- 3) For a user with parameter θ and a movie with learned features x , predict a star rating of $\theta^T \cdot x$

~~*)~~ $x_0 = 1$ arılık yolu! $x \in \mathbb{R}^n$, $\theta \in \mathbb{R}^n$ böyle olunca arılık form
 x ve θ 'ları regularize ediyoruz θ_0 için ayrı bir cate yok!

W9 - Low Rank Matrix Factorization

→ Bu kısımda algoritmanın vectorsalınan implementation'ından bahsedileceğiz.
Bunun yanında verilen bin filme ve ya ona benzeyen diğer film/urunleri
nasıl bulabileceğimden bahsedeceğiz.

Alternative Way of Writing out the Predictions of the Collaborative Filtering Algorithm:

→ Film rating tablosunu alıp matrix selinde yazabiliriz. i : filmler, j : kullanıcılar.

$$y = \begin{bmatrix} 5 & 5 & 0 & 0 \\ 5 & ? & ? & 0 \\ ? & 4 & 0 & ? \\ 0 & 0 & 5 & 4 \\ 0 & 0 & 5 & 0 \end{bmatrix} \quad \text{Predicted Ratings: } \begin{bmatrix} (\Theta^{(1)})^T \cdot (x^{(1)}) \\ (\Theta^{(2)})^T \cdot (x^{(2)}) \\ \vdots \\ (\Theta^{(n)})^T \cdot (x^{(n)}) \end{bmatrix} \dots \dots \begin{bmatrix} (\Theta^{(m)})^T \cdot (x^{(m)}) \end{bmatrix}$$

• $n_m = 5, n_u = 4$

→ There is a vectorized way:
each row is a movie

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

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

Given vectorized X and Θ , the matrix of predictions will be:

$$X \cdot \Theta^T$$

Using Learned Features to Find Related Movies:

→ For each movie/product i , we learn a feature vector $x^{(i)} \in \mathbb{R}^n$.
→ Bu algoritmayla feature öğrenince bu feature'ların ne olduğunu bilmemiz
sadece filmleri karakterize eden bazı özellikler olduğunu söylemektedir.
mesela $x_1 = \text{romance}, x_2 = \text{action}, x_3 = \text{comedy}, x_4 = \dots$ gibi birsey olabilir.
Ancak bunları böyle anlayamayačı insanın anlayabilecegi bir feature değil.

→ How to find movies j related to movie i ?

- ☒ Bunu feature values ile yapabiliriz. Eger iki film n dimensional uzayda birbirine yakinsa demek ki binalar asegi yukarı aynı film/urundur.
- ☒ if $\|x^{(i)} - x^{(j)}\|$ is small movie j and i are "similar"
- ☒ If I need to recommend 5 similar movies to i then find the 5 movies j with the smallest $\|x^{(i)} - x^{(j)}\|$.