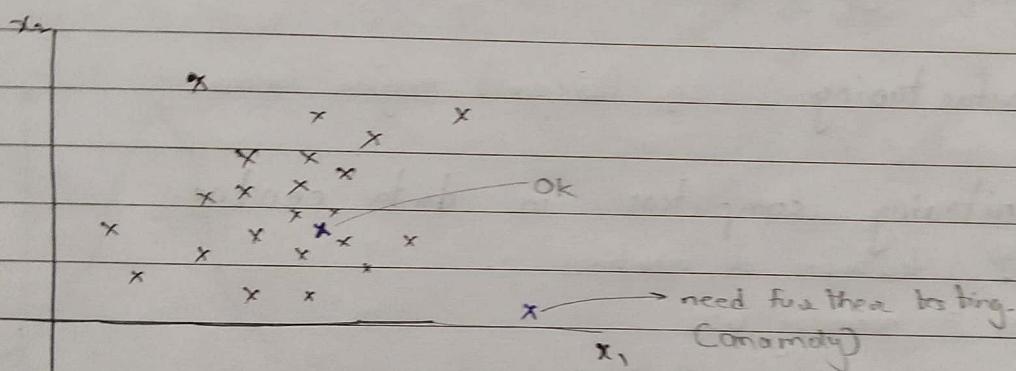


Week 9

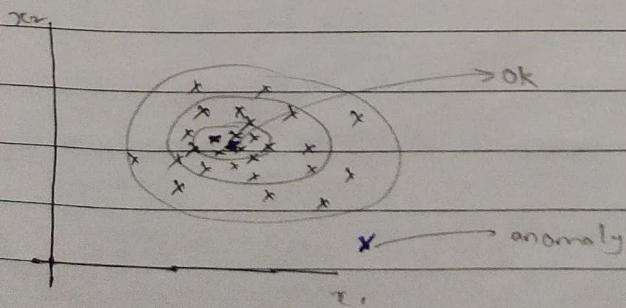
-- Problem Motivation -- -- --

Anomaly detection

Anomaly detection problem.

 $x_1 = \text{heat generated}$ $x_2 = \text{vibration intensity}$ Dataset : $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ New engine : x_{test} 

Density estimation

Dataset : $\{x^{(1)}, \dots, x^{(m)}\}$ Is x_{test} anomalous?We build a model $\rightarrow p(x)$ 

$$P(X_{test}) < \Sigma \rightarrow \text{flag anomaly.}$$

$$P(X_{test}) \geq \Sigma \rightarrow \text{ok}$$

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 $p(x) \leq \epsilon$

x_1 → frequency of login

x_2 → what pages visited

x_3 → no of posts on forum

x_4 → typing speed

- Manufacturing. e.g. Aircraft engine anomaly

- Monitoring computers in data centre.

$x_1^{(i)}$ = features of machine i

x_1 = memory use

x_2 = no of disk access/sec.

x_3 = CPU load

x_4 = CPU load/network traffic

Q1 $P(x) \leq \epsilon$, suppose system is flagging too many things [anomalous] but actually are not [false positive].

✓ Try decreasing ϵ

-- @ Gaussian Distribution --- # Gaussian Distribution / Normal distribution.

★ Gaussian (Normal) Distribution.

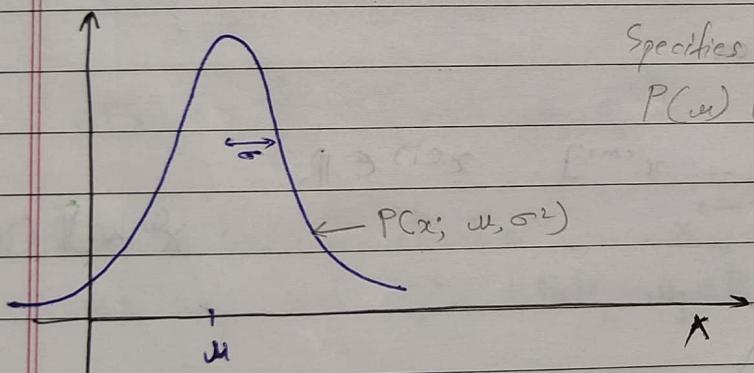
Say $x \in \mathbb{R}$. If x is a distributed Gaussian with mean μ , variance σ^2 .

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

Parameters.

★ Gaussian Distribution / Gaussian Probability density.

Curve is parameterized by μ, σ .



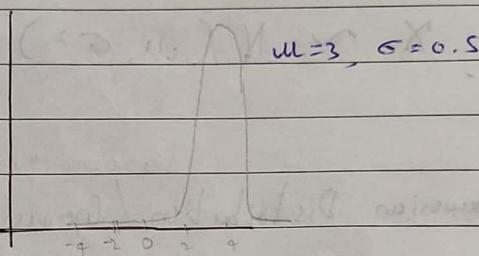
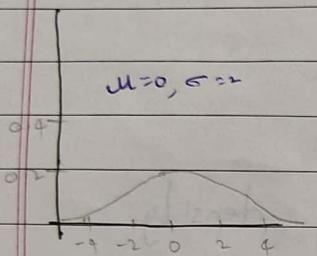
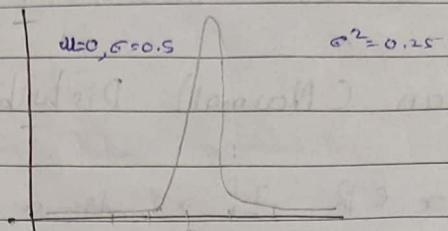
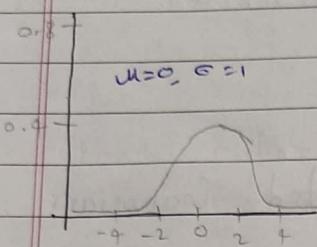
Specifies probability of x taking different values,
 $P(x)$ is high

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

σ : standard deviation.

σ^2 : variance.

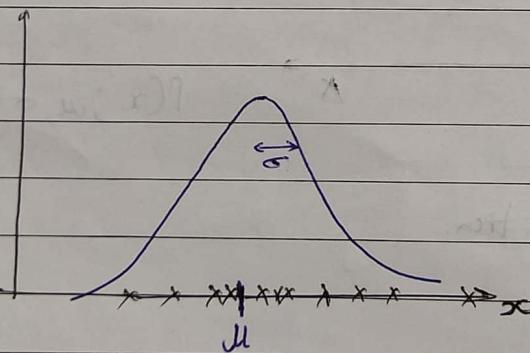
Gaussian distribution example.



Parameter estimation

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

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

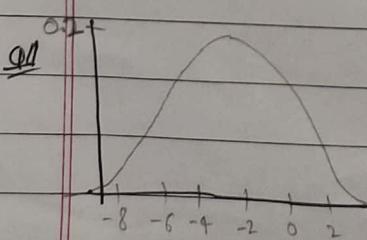


$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

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

μ, σ^2 are maximum likelihood estimates

$\rightarrow \frac{1}{m-1}$ is also used by $\frac{1}{m}$ is good practice



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

-@ Algorithm — estimate p(x)

★ Density estimation.

Training set: $\{x^{(1)}, \dots, x^{(m)}\}$

Each example is $x \in \mathbb{R}^n$.

$$x_1 \sim N(\mu_1, \sigma_1^2)$$

$$x_2 \sim N(\mu_2, \sigma_2^2)$$

$$x_3 \sim N(\mu_3, \sigma_3^2)$$

$$p(x) = p(x_1; \mu_1, \sigma_1^2) p(x_2; \mu_2, \sigma_2^2) p(x_3; \mu_3, \sigma_3^2) \dots p(x_n; \mu_n, \sigma_n^2)$$

This eqⁿ assumes independent assumption on $x_1 - x_n$,
But this works fine whether or not features are anywhere
close to independent. It always works fine.

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

@ Given a training set $\{x^{(1)}, \dots, x^{(m)}\}$, how would you estimate each μ_j and σ_j^2 (Note $\mu_j \in \mathbb{R}$, $\sigma_j^2 \in \mathbb{R}$)

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

channel

★ Anomaly detection algorithm.

1) Choose features x_i that you think might be indicative of anomalous ~~error~~ example.

Choose features which may take unusually high/low value for what an anomalous example might look like.

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

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)} \quad \text{--- } \mu_j \text{ stands for mean of } i^{\text{th}} \text{ feature.}$$

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

$$\boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}^{(i)}$$

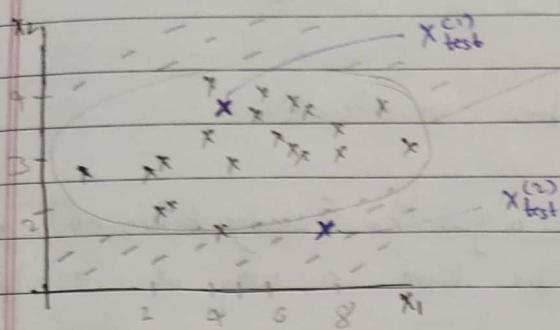
vectorized.

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 $p(x) < \Sigma$

Anomaly detection problem example.



$$\mu_1 = 5, \sigma_1 = 2$$

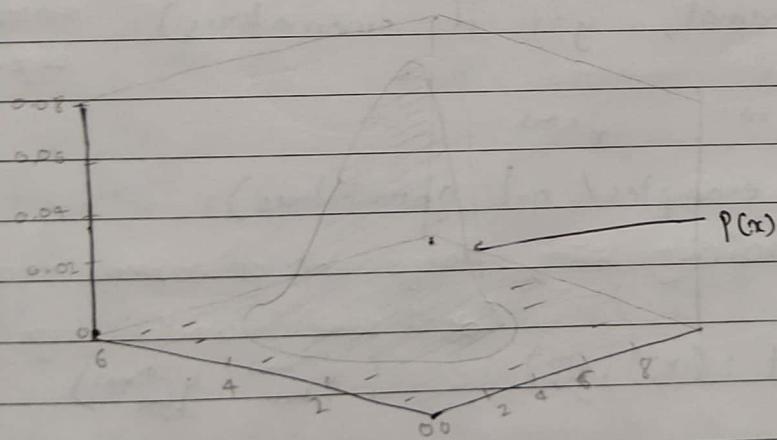
$$\mu_2 = 3, \sigma_2 = 1$$

$$\sigma_{1^2} = 4, \sigma_{2^2} = 1$$

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

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

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



All examples having high above surface corresponds to non-anomalous example or normal example.

--①

Developing and Evaluating an Anomaly detection system # how to build model.

The importance of 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.

Run a algo with & without a feature,

If we get a number back,

it is easy to say include/not the feature.

Anomaly detection is unsupervised.

Assume we have some labelled data of anomalous and non-anomalous example

($y=0$ if normal, $y=1$ if anomalous)

This brings us closer to supervised.

Training set: $x^{(0)}, \dots, x^{(m)}$

(assume normal examples/ not anomalous)

Unlabelled training set \rightsquigarrow large set of normal examples.

Cross Validation set: $(x_{cv}^{(0)}, y_{cv}^{(0)}), \dots, (x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$

(OK if some anomalies slips in training set)

Test set: $(x_{test}^{(0)}, y_{test}^{(0)}), \dots, (x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$

→ Will contain few anomalous examples, which are known

To evaluate
anomaly
detection
algo

Aircraft engines motivating examples.

10,000 good (normal) engines.

20 flawed engines (anomalous)

That's
how we
split the
data

Training set : 6000 good engines ($y=0$) $p(x) = p(x_1; \omega_1, \epsilon^2) \dots p(x_n; \omega_n, \epsilon^2)$
 CV : 2000 good engines ($y=0$), 10 anomalous ($y=1$)
 Test : 2000 good engines ($y=0$), 10 anomalous ($y=1$)

Alternative

Training set : 6000 good engines.
 CV : 4000 good engines ($y=0$), 10 anomalous ($y=1$)
 Test : 4000 good engines ($y=0$), 10 anomalous ($y=1$)

→ This is not recommended using same set of $y=0$ in both CV & Test.

Even same set of $y=1$, 20, 20 can be done for CV, test

but this is not recommended.

★ Algorithm evaluation.

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

- On a cross validation/test example x , predict y

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

 $(x_{\text{test}}^{(i)}, y_{\text{test}}^{(i)})$

New anomaly problem is predicting y , just like supervised.

Possible evaluation metrics:

- True +ve Confusion matrix (TP, FP, TN, FN)

- Precision / recall

- F1-score \rightarrow single number

Can also use CV set to choose parameter ϵ .
 Choose ϵ that maximizes F1 score. \leftarrow Try many ϵ

Q Suppose you have fit a model $p(x)$. When evaluating on the cross validation set or test set, your algorithm predicts $y = \begin{cases} 1 & p(x) \leq \varepsilon \\ 0 & p(x) > \varepsilon \end{cases}$ Is classification accuracy a good way to measure algo's performance.

→ No, because of skewed classes

(so an algo that always predicts $y=0$ will have high accuracy)

✗ classification accuracy

So use FP, FN, TP, TN . Precision / recall is F1 score.

--① Anomaly detection vs Supervised learning.

Why didn't we use logistic reg to solve anomaly detection, if we had labelled data.
When to use logistic/anomaly DP?

Anomaly detection vs Supervised learning.

Very small number of positive examples ($y=1$) (0-20 is common) & large no of negative examples.

Large number of negative ($y=0$) examples.
→ fit $p(x)$

Many different "types" of anomalies. 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.

- Many ways to go wrong aircraft engine.

- Difficult to learn how anomaly look

like from small set of positive examples e.g. email spam classifier.

- Future anomaly may look nothing like previous anomaly.

- So model is trained on negative examples

other than hard positive e.g.

- Weather prediction

- Cancer classification.

If we have enough data for

Future anomaly may look nothing like any of the anomaly e.g. we've seen so far.

- e.g. Fraud detection

- Manufacturing (engine)

- Monitoring machines in data centres.

positive e.g.'s we can shift from anomaly to supervised learning.

~~It was skipped. easy.~~

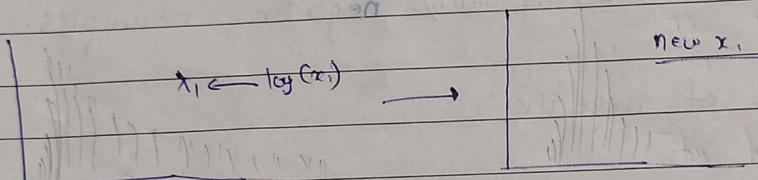
-- @ Choosing what feature to use optional
Features matters very much.

Make sure if data is gaussian / normally distributed.
If not play with it to make gaussian.

$$p(x_i; \mu_i, \sigma_i^2)$$

$$\begin{aligned}x_1 &\leftarrow \log(x_i) \\x_2 &\leftarrow \log(x_1 + t) \\&\quad \log(x_2 + t) \\x_3 &\leftarrow \sqrt{x_3} \\x_4 &\leftarrow x_4^{\frac{1}{2}}\end{aligned}$$

That's how
we play.



Error analysis for supervised.

Train model test on CV can we get extra features
to make model better than we're using in CV
Some here.

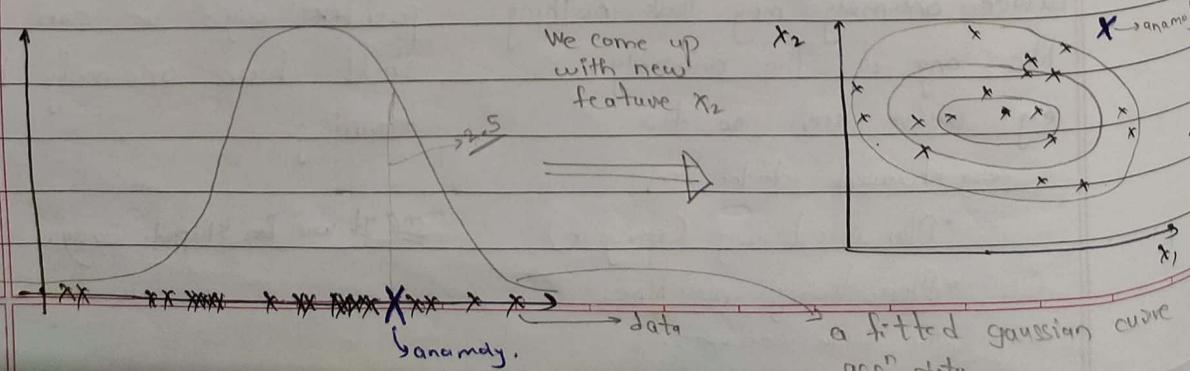
★ Error analysis for anomaly detection. Procedure.

We want [

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

Most common problem.

$\Rightarrow p(x)$ is comparable (say, both large) for normal & anomalous x .



If existing feature is failing to flag anomalies as anomaly try to find new features.

~~★~~ Choose features that will take very large / very low value, in event of an anomaly.

Monitoring Computer in data centre

x_1 = memory use of PC.

x_2 = number of disk access/sec

x_3 = CPU load

x_4 = network traffic.

New features to capture anomaly to like.

$$x_5 = \frac{\text{network load}}{\text{network traffic}} \quad \begin{matrix} (n1)^2 \\ \text{int} \end{matrix}$$

[OR]

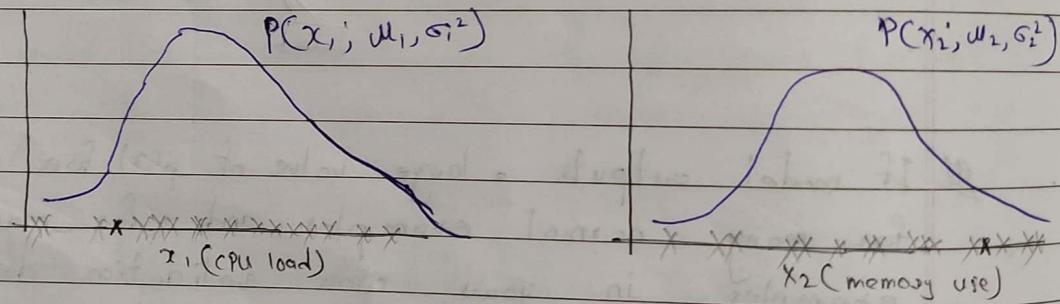
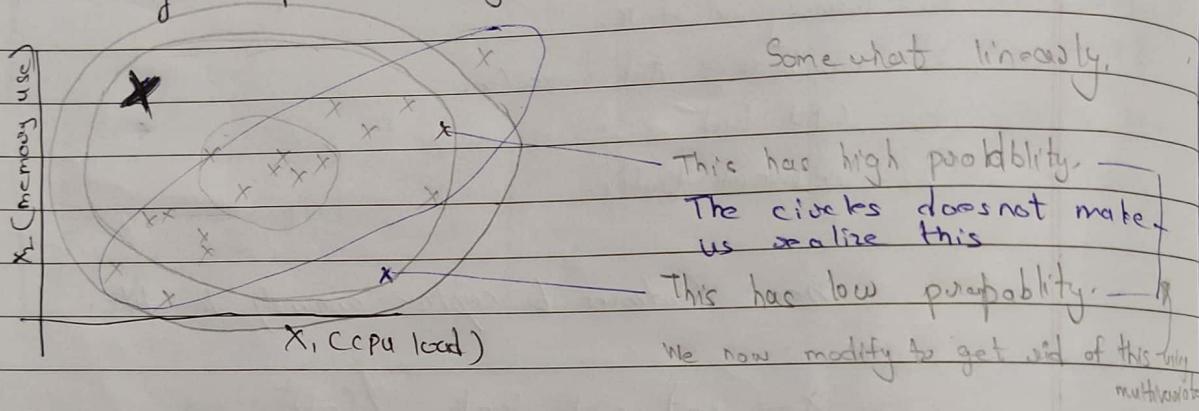
Q1 If model outputs a large value of $p(x)$ for many values of $p(x)$ for many normal examples and for many anomalous examples in your cross validation dataset.

Which of the following changes to your algo. what is most likely to help

- ✓ Try coming up with more features b/w normal & anomalous e.g.
- Try using fewer features.

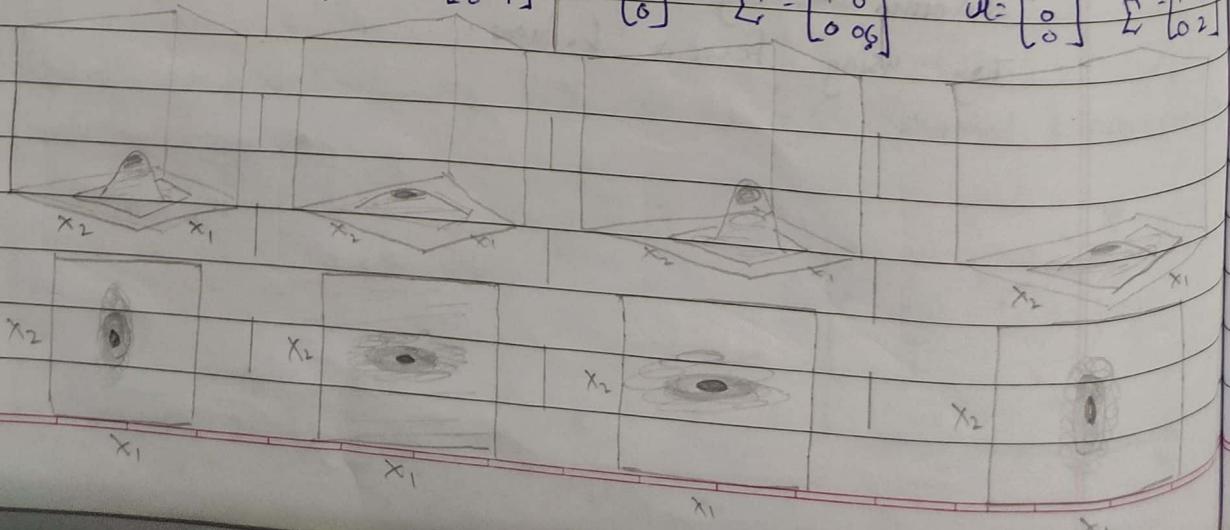
-- ① Multivariate Gaussian Distribution
 # extension of anomaly detection
 ↗ uses multivariate GD
 Advantages
 Disadvantages

motivating example: Monitoring machine in data centre.



x_1 vs x_2 example looks anomalous.

But in x_1, x_2 separately, e.g. not look like anomalous.
 $\mu = [0]$, $\Sigma = [0 \ 0]$ $\mu = [0]$ $\Sigma = [2 \ 0]$ $\mu = [0]$ $\Sigma = [1 \ 0]$ $\mu = [0]$ $\Sigma = [0 \ 0.5]$



★ Multivariate Gaussian (Normal) distribution.

$x \in \mathbb{R}^n$. 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)
 ↪ same form PCA.

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

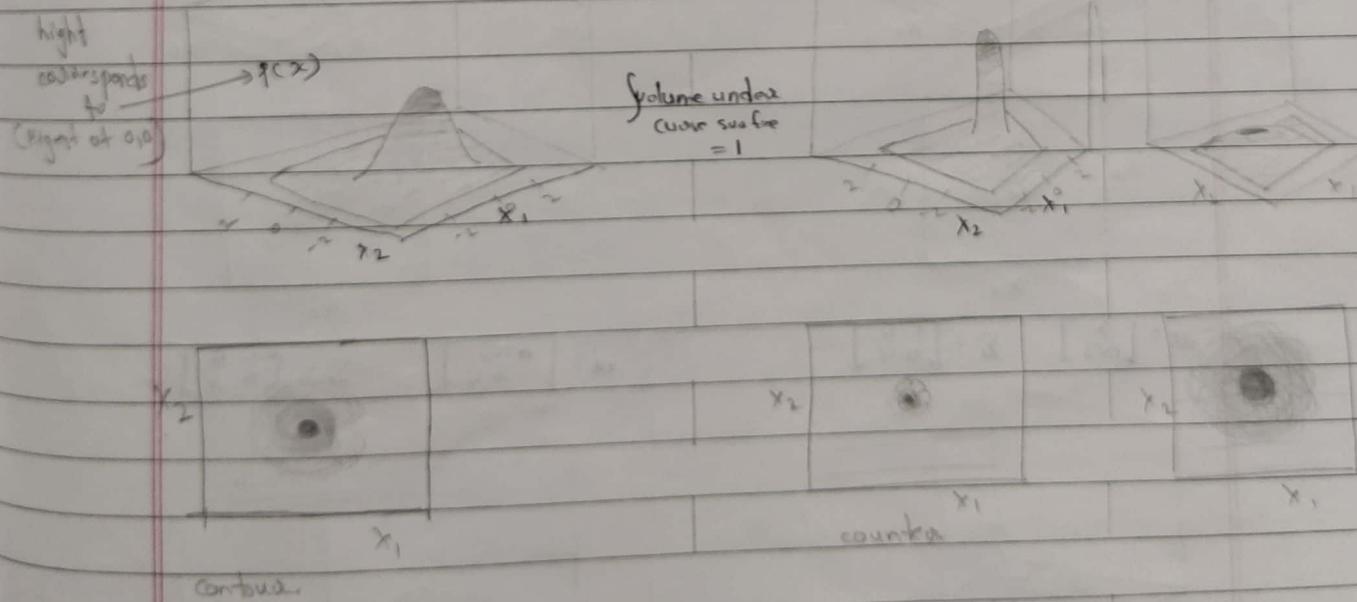
$\rightarrow |\Sigma| = \text{determinant of } \Sigma \quad | \det(\text{Sigma})|$

Multivariate Gaussian (Normal) examples.

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

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

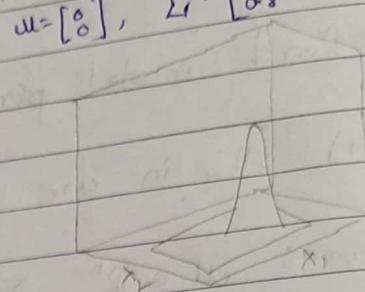
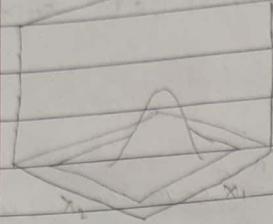
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$



Go To Previous
Page.
and Read
Skip This →

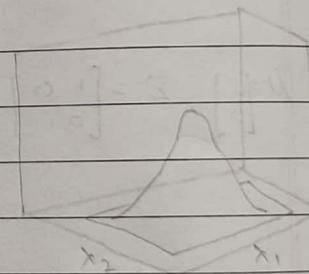
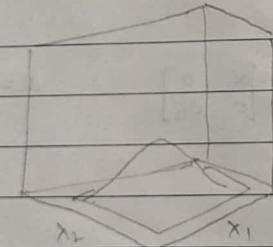
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$$

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$



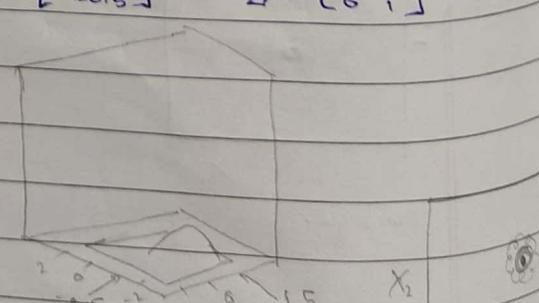
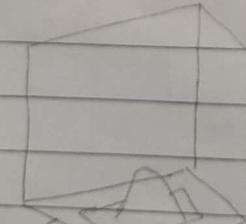
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix}$$

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & -0.8 \\ -0.8 & 1 \end{bmatrix}$$



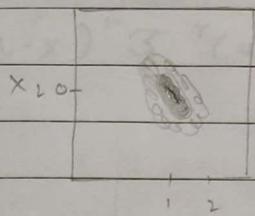
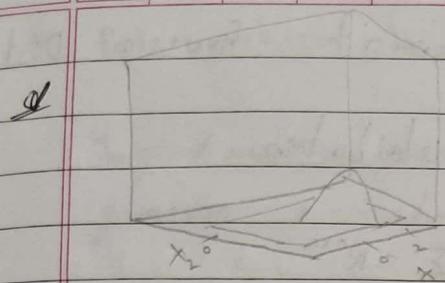
$$\mu = \begin{bmatrix} 0 \\ 0.5 \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

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



x_1

x_2



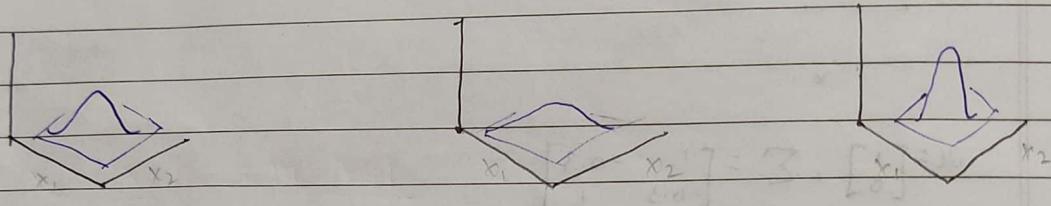
$$\mu = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & -0.3 \\ -0.3 & 1 \end{bmatrix}$$

-- @ Anomaly Detection using Multivariate Gaussian Distribution.

Multivariate Gaussian (Normal) Distribution.

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

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



Parameter fitting:

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

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

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

* Anomaly detection using multivariate.

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.

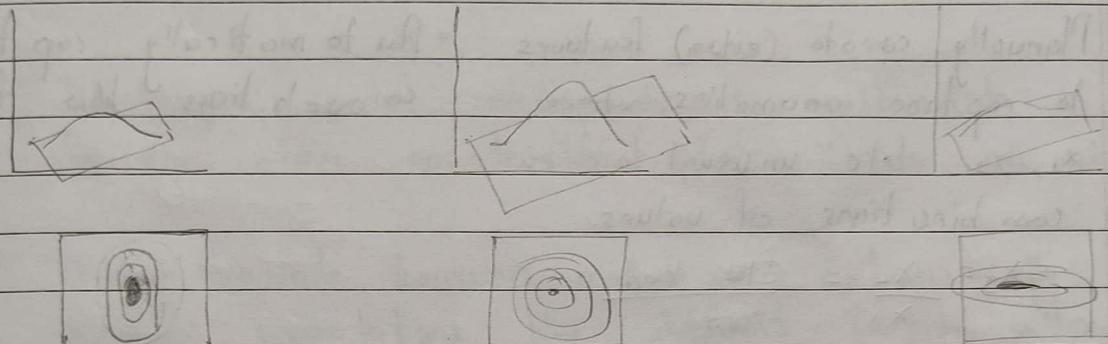
$$p(x) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

Flag anomaly if $p(x) < \epsilon$

* Relationship to original model.

Original model is special case of multivariate with contours of gaussian & constrained to be always axis aligned.

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



$$\Sigma = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$

This will give diagonal covariance.

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 $\Sigma = \begin{bmatrix} \sigma_1^2 & 0 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & 0 & \dots & 0 \\ 0 & 0 & \ddots & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 & \sigma_n^2 \end{bmatrix}$

Original model

$$p(x_1; \mu_1, \Sigma_1^2) \times \dots \times p(x_n; \mu_n, \Sigma_n^2)$$

used more often

Multivariate Gaussian.

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

used somewhat, but it can capture correlations b/w features

- Manually create (extra) features to capture anomalies where x_1, x_2 take unusual combinations of values.
- Automatically captures correlations b/w features

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

If we got time to manually write extra feature, original works.

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

- Computationally cheap (computationally scales better to large n). $n=10,000/100,000$
- Computationally more expensive.
- Ok even if m (training set size) is small.
- Must have $m > n$ or else Σ is non-invertible

Recommended

$$m \geq 10n$$

$$\Sigma \sim \frac{n^2}{2} \text{ no. of parameters}$$

In multivariate if Σ is singular/non-invertible

→ Fails $m > n$ condⁿ

→ We got redundant features

$$x_1 = x_2 \quad | \quad x_3 = x_4 + x_5$$

All consider anomaly detection using training set $\{x^{(i)}\}_{i=1}^m \in \mathbb{R}^n$
where $x^{(i)} \in \mathbb{R}^n$.

- ✓ The original model corresponds to multivariate Gaussian where contours of $p(x; \mu, \Sigma)$ are axis aligned.
- ✓ The multivariate Gaussian model can automatically capture correlations b/w different features in x .
- ✓ The original model can be more computationally efficient than the multivariate Gaussian model, and thus might scale better to very large value of n (no. of features)

(Recommender system)

--@ Problem formulation

Why ↗ Important -

↗ Big idea → Algorithm, just to learn what features to use.

algo, just to learn what feature to use, ~~and~~ by itself
recommender system is one of them.

e.g. Predicting movie ratings.

(simple)

User rating movies using + to 5 stars

movie	Alice(1)	Bob(2)	Carol(3)	Dave(4)
Lol	5	5	6	0
RF	5	2	?	0
CPL	2	4	0	?
NCC	6	0	5	4
Sysk	0	0	5	?

n_u = no of users.

n_m = no of movies.

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

$y(i,j) = \text{rating given by user } j \text{ to movie } i$
(defined only if $a(i,j) = 1$)

$n_u = 4$, $n_m = 5$

Try to predict 2.

Q1

$$n_m = 2, n_u = 3$$

	User1	User2	User3
movie 1	0	1	2
movie 2	3	5	5

→ ✓ $a(2,1) = 0, g(2,3) = \text{undefined} \leftarrow \text{Recommendation system goal is to find this.}$

Content based recommendations - - - - -

Content based recommendation systems.

Movie	$\alpha^{(1)}$	$\alpha^{(2)}$	$\alpha^{(3)}$	$\alpha^{(4)}$
Lal	1	5	5	0
RF	2	5	2	0
CPOL	3	4.95	4	2
NBCC	4	0	0	5
S vs k	5	0	0	2

$$x_0 = 1 \quad x^{(1)} = \begin{bmatrix} 1 \\ 0.9 \\ 0 \end{bmatrix}$$

$n=2$
excluding $x_0=1$

$$x^{(1)} = \begin{bmatrix} x_1^{(1)} = 1 \\ x_2^{(1)} = 0.9 \\ x_3^{(1)} = 0 \end{bmatrix}$$

$$n_u = 4, n_m = 5$$

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

$$x^{(3)} = \begin{bmatrix} 1 \\ 0.99 \\ 0 \end{bmatrix} \longleftrightarrow \alpha^{(1)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix}$$

$$(\alpha^{(1)})^T x^{(3)} = 5 \times 0.99 = 4.95 \quad \therefore \text{Alice, CPOL} = 4.95$$

We are applying a different copy of this linear page for each user, and we're saying what it does is A has parameter vector θ^A that he uses that we use to predict his ratings as a function of how romantic & action packed a movie is.

B, C, D each has diff LR fun of x_1, x_2 that's how we are going to predict their ratings.

All consider

	A	B	C	D	romance	action
	x_1	x_2			x_1	x_2
-	5	5	0	6	0.9	0
-	5	2	2	0	1.0	0.9
-	3	4	0	3	0.99	0
-	0	0	5	4	0.1	1.0
-	0	0	5	7	0	0.9

$$\text{minimizing } Q(\theta) = \frac{1}{2}$$

$$\begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix} - \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 4 \end{bmatrix} \quad \checkmark \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix} \rightarrow \text{because we want } \theta \text{ to be small.}$$

* Problem formulation.

$a_{ij} = 1$ if user j has rated movie i (0 otherwise)

$y^{(i,j)}$ = rating by user j on movie i (if defined)

$\alpha^{(j)}$ = parameter vector for user j .

$x^{(i)}$ = feature vector for movie i .

For user j , movie i , predicted rating: $(\alpha^{(j)})^T (x^{(i)})$

$m^{(j)}$ = no of movies rated by user j .

To learn $\alpha^{(j)}$:

- Basically a linear regression problem.
- choose $\alpha^{(j)}$ as predicted value are close as possible to values that we observed in training sets & the values we observed in our data.

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

$$\min_{\alpha^{(j)}} \frac{1}{2m^{(j)}} \sum_{i: a_{ij}=1} ((\alpha^{(j)})^T (x^{(i)}) - y^{(i,j)})^2 + \frac{\lambda}{2m^{(j)}} \sum_{k=1}^n (\alpha^{(j)})_k^2$$

just to make both simple.

* optimization objective.

To learn $\alpha^{(j)}$ (parameter for user j): (only for specific user)

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

To learn $\alpha^{(1)}, \alpha^{(2)}, \dots, \alpha^{(n_u)}$: (For all users)

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

 Optimization algorithm:

$$\min_{\alpha^{(1)}, \dots, \alpha^{(n_u)}} \frac{1}{2} \sum_{j=1}^{n_u}$$

Gradient descent update.

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

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

Content based recommendation.

Someone else giving us all of these features for all of the movie in our data set.

-- @

Collaborative filtering - - -

property

Feature learning → algo that can learn for itself what features to use

Previously we had assume for each movie someone had come and told us how romantic that movie was and how much action it was.

We had data about how the movie is.

But we want to learn more features than just two.

Problem motivation.

Movie	Alicia	Bob	Candice	Dave	John	Karen	Liam	Mia	Natalie	Olivia	Paul	Rachel	Sarah	Taylor	Uma	Vivian	Wendy	Xavier	Yvonne	Zoe
	$\alpha^{(1)}$	$\alpha^{(2)}$	$\alpha^{(3)}$	$\alpha^{(4)}$																
Lal	5	5	0	0			5	2												
RF	5	3	3	0			2	2												
CPOL	3	4	0	2			3	2												
WCC	0	0	0	0			0	0												
Susk	0	0	5	3			2	2												

Now let make a diff assumption

assumption → We're gone to each of users, & each of our users has told us how much they like romantic movie and how much they like action.

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

↳ likes romantic ↳ like romantic ↳ like action ↳ no like action

no like action no like action no like action no like romantic

Each user tells us value of $\alpha^{(j)}$

specifies how much they like difl types of movies.

$$\Omega^T x = y$$

$$\Omega^T \quad x \quad = \quad y$$

previously

$$\begin{bmatrix} 0 \\ 0 \\ 2 \\ 0 \end{bmatrix}$$

To find

given

given

content
based

Partial Ω was

Features

given,

were given

Dataset

we found

x_1, x_2

had ratings.

whole by $\Omega^T x = y$

LR

Train Given

To find

Given

Now

Collaborative
filtering.

$$\text{For } X^{(1)} : (\Omega^{(1)})^T x^{(1)} \approx s$$

$$(\Omega^{(2)})^T x^{(1)} \approx s$$

$$(\Omega^{(3)})^T x^{(1)} \approx s$$

$$(\Omega^{(4)})^T x^{(1)} \approx 0$$

so $x_1^{(1)}$ might be $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$

Q1 Consider following movie ratings:

	user1	user2	user3	user4 (comedy)
movie 1	0	1.5	2.5	?

only one feature x_1 , suppose $\Omega^{(1)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$, $\Omega^{(2)} = \begin{bmatrix} 0 \\ 3 \end{bmatrix}$, $\Omega^{(3)} = \begin{bmatrix} 0 \\ 5 \end{bmatrix}$, $x_1 = ?$

→ 0.5, $\Omega^T x$ should be $\approx y$ (dataset recordings)

* Optimization algorithm.

- Given $\alpha^{(1)}, \dots, \alpha^{(n_u)}$, to learn $x^{(i)}$:

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

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

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

Q Suppose you use gradient descent to minimize,
this update for $i \neq 0$

$$\rightarrow x_k^{(i)} := x_k^{(i)} - \alpha \left(\sum_{j: a_{ij}=1} ((\alpha^{(j)})^T x^{(i)} - y^{(i,j)}) \alpha_k^{(j)} + \lambda x_k^{(i)} \right)$$

* Collaborative filtering

1) content based recommendation: Given $x^{(1)}, \dots, x^{(n_m)}$ (and movie rating)
can estimate $\alpha^{(1)}, \dots, \alpha^{(n_u)}$

2) collaborative filtering

Given $\alpha^{(1)}, \dots, \alpha^{(n_u)}$

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

first guess randomly.

Every new entry make algo stronger, \rightarrow "collaborative"

Collaborative filter algorithm.

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

$$\min_{\theta^{(1)}, \dots, \theta^{(n_u)}} \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i: x^{(i)} \neq 0} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \lambda \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2$$

Sum over all users Sum over all movies rated by that user Summing over all users that have rated movie i that corresponds to movie i started by user j

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

$$\min_{x^{(1)}, \dots, x^{(n_m)}} \frac{1}{2} \sum_{i=1}^{n_m} \sum_{j: \theta^{(j)} \neq 0} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \lambda \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2$$

for a movie i ; sum over all users j that have rated movie i .

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

$$J(x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)}) = \frac{1}{2} \sum_{(i,j): x^{(i)} \neq 0} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2$$

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

Same
as $\Rightarrow \theta \rightarrow x \rightarrow \theta \rightarrow x \rightarrow \theta \rightarrow \dots$
but
cheaper

$x_0 = 1$ $x \in \mathbb{R}^n$ If algo

θ_0 $\theta \in \mathbb{R}^n$ want a feature always

$x \in \mathbb{R}^{n+1}$ It will choose to learn itself
e.g. $x^{(1)} = 1$

This has got some properties like.

If we hold x 's constant & just minimize with respect to thetas $\rightarrow 1$

because either of regularization term is constant

If we hold theta constant & just minimize $J()$ with respect to x 's $\rightarrow 2$

♦ Collaborative filtering algo.

- 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 algo).

e.g. for every $j=1, \dots, n_u$, $i=1, \dots, n_m$:

$$x_k^{(i)} := x_k^{(i)} - \alpha \left(\sum_{j: x_{k,j} \neq 1} ((\Theta^{(j)})^T x^{(i)} - y^{(i,j)}) \Theta_k^{(j)} + \lambda x_k^{(i)} \right)$$

$$\Theta_k^{(j)} := \Theta_k^{(j)} - \alpha \left(\sum_{i: x_{k,i} \neq 1} ((\Theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} + \lambda \Theta_k^{(j)} \right)$$

- 3) For a user with parameters Θ and a movie with (known) features x , predict a star rating of $\Theta^T x$.

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

Q1 Why we initialized $x^{(1)}, \dots, x^{(n_m)}$ & $\Theta^{(1)}, \dots, \Theta^{(n_u)}$ to small random values.

→ This avoids a symmetry breaking (similar to random initialization of a neural network's parameters) and ensures the algo learns features $x^{(1)}, \dots, x^{(n_m)}$ that are different from each other.

Collaborative filtering

Vectorization: Low Rank matrix Factorization.
Vectorization of collaborative filtering.

Applications :-

- Given one product can you find other product that are related to this

	Marie	Alice(1)	Bob(2)	Catal(3)	Dave(4)
LaL	S	S	0	0	
RF	S	2	2	0	
CPOL	2	4	0	2	
NCC	0	0	5	4	
Susk	0	0	5	2	

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 \\ 5 & 2 & 2 & 0 \\ 3 & 4 & 0 & 2 \\ 0 & 0 & 5 & 4 \\ 0 & 0 & 5 & 0 \end{bmatrix}$$

y (ij)

~~→ Predicted eatings:~~

$$\begin{array}{ccccccc}
 & & \text{year, month} & & & & \\
 & (\mathbb{Q}^{(1)})^T(x^{(1)}) & & (\mathbb{Q}^{(2)})^T(x^{(1)}) & \cdots & (\mathbb{Q}^{(n_u)})^T(x^{(1)}) & \\
 \\
 (\mathbb{Q}^{(1)})^T(x^{(2)}) & (\mathbb{Q}^{(2)})^T(x^{(2)}) & \cdots & (\mathbb{Q}^{(n_u)})^T(x^{(2)}) & & & \\
 \\
 & \vdots & & \vdots & & \vdots & \\
 \\
 (\mathbb{Q}^{(1)})^T(x^{(n_m)}) & (\mathbb{Q}^{(2)})^T(x^{(n_m)}) & \cdots & (\mathbb{Q}^{(n_u)})^T(x^{(n_m)}) & & &
 \end{array}$$

movie last

(matrix way)

$$c_{ij} \rightarrow (Q^{(j)})^T(x^{(i)})$$

$$X @^T$$

(vectorized way)

Vectorization

$$X = \begin{bmatrix} -(x^{(1)})^\top \\ -(x^{(2)})^\top \\ \vdots \\ -(x^{(n)})^\top \end{bmatrix} \quad W = \begin{bmatrix} -(\alpha^{(1)})^\top \\ -(\alpha^{(2)})^\top \\ \vdots \\ -(\alpha^{(n)})^\top \end{bmatrix}$$

XW^\top has property name low rank matrix

So algo is called " " factorization.

Collaborative filtering algo here, you can use the learned features in order to find related movies.

* Finding related movies.

- For each product i , we learn a feature vector $x^{(i)} \in \mathbb{R}^n$, $x_1 = \text{romance}$, $x_2 = \text{action}$, $x_3 = \text{comedy}$, $x_4 = \dots$

It can hard to figure out what these features are, But it will learn features that are very meaningful for capturing whatever are most important properties of movie

- How to find movies j related to movie i ?

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

5 most similar movies to movie i :

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

Implementation Details: Mean normalization. -

* Users who have not rated any movie.

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

Lol 5 5 0 6 3.0

RF 5 2 3.0 0 2.0

CPOL 4 0 2.0 2.0 2.0 $\gamma = \begin{bmatrix} 5 & 5 & 0 & 0 & 2 \\ 5 & 2 & 2 & 0 & 2 \\ 2 & 4 & 0 & 2 & 2 \\ 0 & 0 & 5 & 4 & 2 \\ 0 & 0 & 5 & 0 & 2 \end{bmatrix}$

NCC 0 5 4 2.0 2.0

Sink 0 0 5 2.0 2.0

$$\min_{x^{(1)}, \dots, x^{(n_u)}, \alpha^{(c_s)}, \dots, \alpha^{(c_n)}} \frac{1}{2} \sum_{i,j} ((\alpha^{(c_j)})^T x^{(i)} - y^{(c_i,j)})^2 + \lambda \sum_{i=1}^m \sum_{k=1}^{n_u} (x_{ik}^{(c_i)})^2 + \lambda \sum_{j=1}^n \sum_{k=1}^{n_u} (\alpha_{jk}^{(c_k)})^2$$

This term does not matter
 $\alpha^{(c_i)} \cdot \alpha^{(c_j)} = 1$

only this term matters in estimation $\alpha^{(c_s)}$

$n=2$, (\because We want to learn two features)

$$\lambda [(\alpha_1^{(c_s)})^2 + (\alpha_2^{(c_s)})^2]$$

minimizing this gives

$$\alpha^{(c_s)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Rating given by Eve $\Rightarrow (\alpha^{(c_s)})^T X = 0$ --- (so Eve will rate all movies with zero stars)

This doesn't feel useful. We cannot recommend her any movie.

Mean normalization will fix this problem.



Mean normalization. (Mean normalization as a sort of pre-processing step for collaborative filtering)

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

$$\bar{u} = \frac{2.5 + 2.5 + 2 + 2.25 + 1.25}{5} = 2$$

$$Y = \begin{bmatrix} 2.5 & 2.5 & -2.5 & -2.5 & ? \\ 2.5 & ? & ? & -2.5 & ? \\ ? & 2 & -2 & ? & ? \\ -2.25 & -2.25 & 2.75 & -1.75 & ? \\ -1.25 & -1.25 & 3.75 & -1.25 & ? \end{bmatrix}$$

compute average rating that each movie obtained.

$$\bar{u} = \frac{\text{sum of observation}}{\text{no of observation}}$$

(who has rated the movie)

Subtract off mean rating.

Normalizing each movie to have average rating of zero.

(? stays ?.)

Now each movie has average rating of 0.

Use this as dataset.

keep $\alpha^{(i)}, x^{(i)}$

For user j on movie i predict:

$$\rightarrow (\alpha^{(ij)})^\top (x^{(i)}) + u_i$$

User S (Eve):

$$\alpha^{(s)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \quad \underbrace{(\alpha^{(s)})^\top (x^{(i)})}_{=0} + u_i$$

So, we are predicting rating as u_i for user Eve.
movie1 = 2.5, 2 → 2.5, 3 → 2, 4 → 2.5
(This makes sense)

If we have movie with no rating. --

movie with no

user who hasn't rated anything.

analogous to

Normalize different columns (instead of normalizing rows)

Maybe shouldn't recommend that movie to anyone, anyway.

If you have a movie with no rating.

- ④ We talked about mean normalization. However, unlike some other applications of feature scaling, we did not scale the movie ratings by dividing by the range ($\max - \min$) value. This is because:-
- ✓ All movie ratings are already comparable (0-5 stars), so they are already on similar scale.