

Week 9

Anomaly detection - problem motivation

Anomaly detection example:

Aircraft engine features:

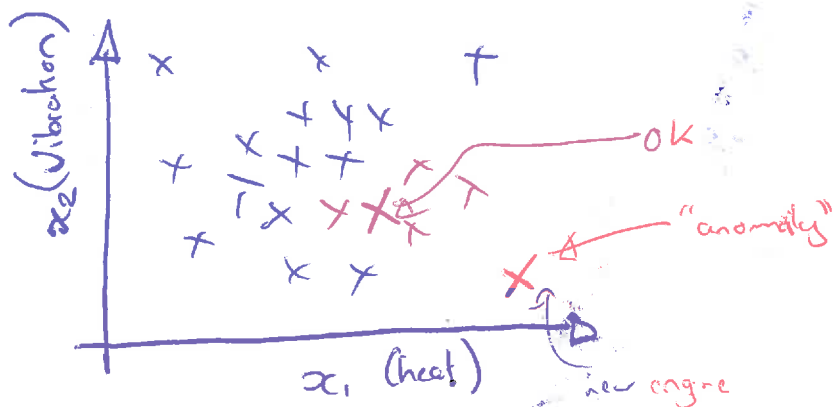
x_1 = heat generated

x_2 = vibration intensity

...

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

New engine: x_{test} .



Density estimation

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

Is x_{test} anomalous? \rightarrow build model $p(x)$



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

$p(x_{\text{test}}) \geq \epsilon \rightarrow$ OK.

Week 9

Anomaly detection - problem motivation

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 $p(x) < \epsilon$

Manufacturing

Monitoring computers in a data centre

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

x_1 = memory use, x_2 = number of disk accesses/sec,

x_3 = CPU Load, x_4 = CPU Load / network traffic.

$$p(x) < \epsilon$$

Week 9

Anomaly detection - Gaussian distribution

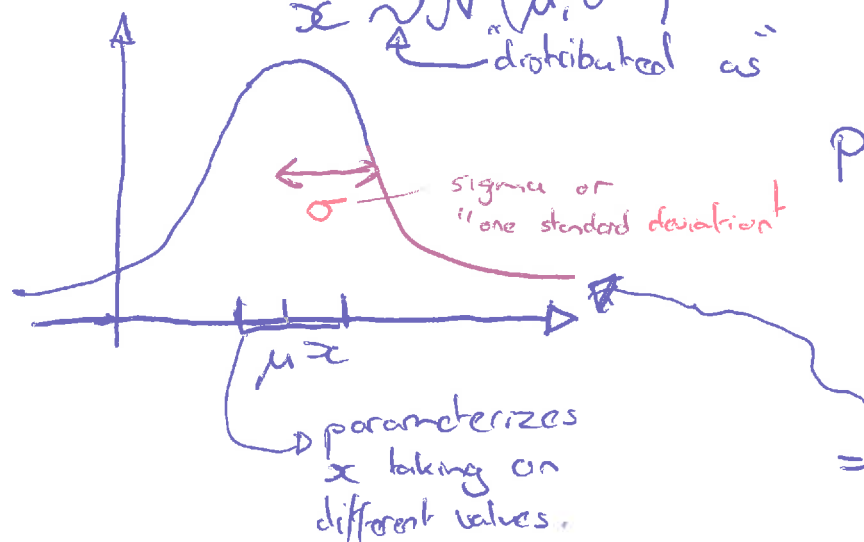
→ also called "normal distribution"

Gaussian (Normal) distribution

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

bell shape curve parameterized by:

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



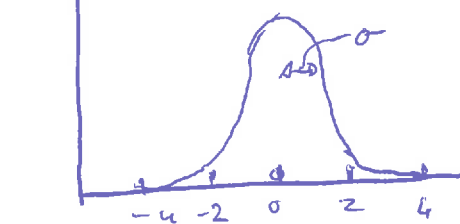
$$p(x; \mu, \sigma^2)$$

probability of x is parameterized by the two parameters μ & σ^2

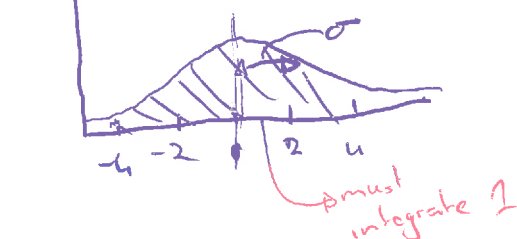
$$= \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

Gaussian distribution example

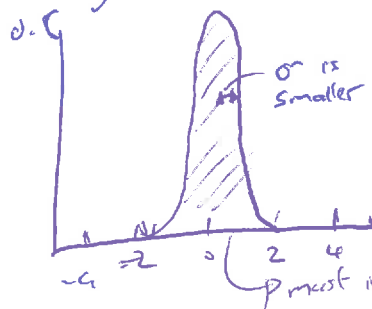
$\mu=0, \sigma=1$



$\mu=0, \sigma=2$

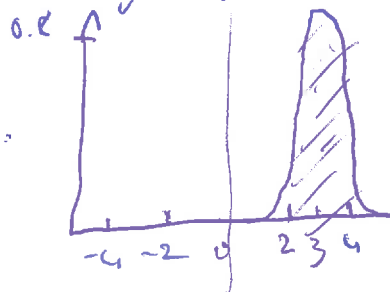


$\mu=0, \sigma=0.5$



$\sigma^2 = 0.25$

$\mu=3, \sigma=0.5$



Week 9

Parameter estimation

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

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



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

ML

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

both are the "maximum likely estimations"

Week 9

Anomaly detection - Algorithm

Density estimation

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

each example $x \in \mathbb{R}^n$

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

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

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

$p(x)$

$$= p(x_1 | \mu_1, \sigma_1^2) p(x_2 | \mu_2, \sigma_2^2) p(x_3) \dots p(x_n)$$

$$= \prod_{j=1}^n p(x_j | \mu_j, \sigma_j^2) \quad \text{compact form}$$

product of a set of values

similar to \sum except product...

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

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

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

Anomaly detection algorithm

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

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

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)} \quad \text{parameterized by } p(x_j | \mu_j, \sigma_j^2)$$

estimate all the values for μ simultaneously

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

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

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

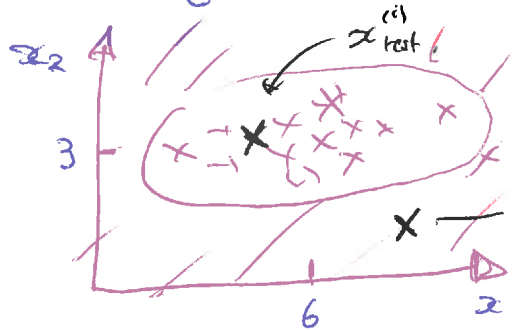
$$p(x) = \prod_{j=1}^n p(x_j | \mu_j, \sigma_j^2) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi}\sigma_j} \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right)$$

Anomaly if $p(x) < \epsilon$

Week 9

Anomaly detection - Algorithm

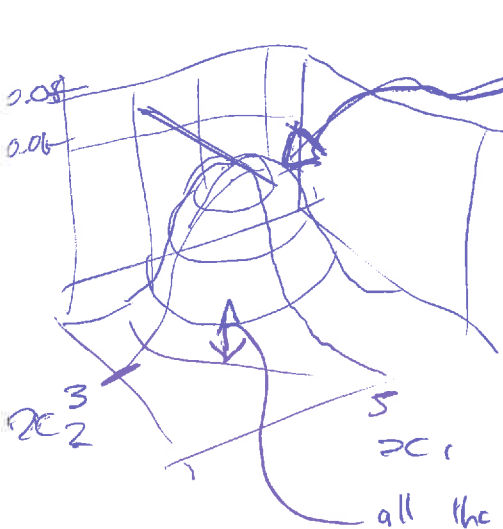
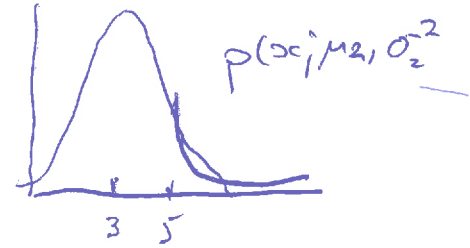
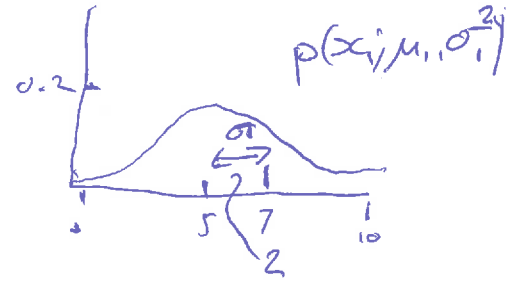
Anomaly detection example:



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

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

$$\epsilon = 0.02$$



$p(x)$ which can be given by:

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

$$\epsilon = 0.02$$

$$p(x_{test}^{(1)}) = 0.0426 \geq \epsilon \therefore \text{not an anomaly}$$

$$p(x_{test}^{(2)}) = 0.0021 < \epsilon \checkmark \Rightarrow \text{anomaly}$$

all the points with a lot of height correspond to "non-anomalous" areas.

Week 9

Anomaly detection - developing and evaluating an anomaly detection system

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.

⇒ Assume we have some labelled data, of anomalous and non-anomalous examples ($y=0$ if normal, $y=1$ if anomalous)

Training set: $x^{(1)}, x^{(2)}, \dots, x^{(n)}$ (assume normal examples / not anomalous)

cross validation set: $(x_{cv}^{(1)}, y_{cv}^{(1)}, \dots)$

test set: $(x_{test}^{(1)}, y_{test}^{(1)}, \dots)$

have a few examples where $y=1$ (anomalous)

Aircraft engines motivating example

10000 good (normal) engines
20 flawed engines (anomalous)

$y=1$ (20... 50)
↑ typical
anomalous example
number

Training set 6000 good engines ($y=0$)

ex: 2000 good engines ($y=0$), 10 anomalous ($y=1$)

Test: 2000 good engines ($y=0$), 10 " ($y=1$)

Use these to fit $p(x) = \prod p(x_i; \mu_i, \sigma_i^2) \times p(x_2, \dots)$

Alternative:

Training set: 6000 good engines

CV: 4000 good engines ($y=0$), 10 anomalous ($y=1$)

Test: 4000 " engine ($y=0$), 10 anomalous ($y=1$)

→ some data → not recommended.

Week 9

Anomaly detection - developing...

Algorithm evaluation

fit model $p(x)$ on training set $\{x^{(1)}, \dots, x^{(n)}\}$
on a cross validation / test example x , predict.

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

on the cross validation set we are going to use our model to predict y .

Possible evaluation metrics:

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

- True positive, false positive, false negative, true negative
- precision / recall
- F₁-score

continually evaluate on the cross-validation sets

Can also use cross-validation set to choose parameter ϵ

final evaluation on the test sets

Week 9

algorithm.

Anomaly detection: anomaly detection vs. supervised learning

Anomaly detection

vs.

Supervised learning

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 like any of the anomalies we have seen so far.

Large number of positive and negative examples. \leftarrow

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.

spam

Ex: $y=1$

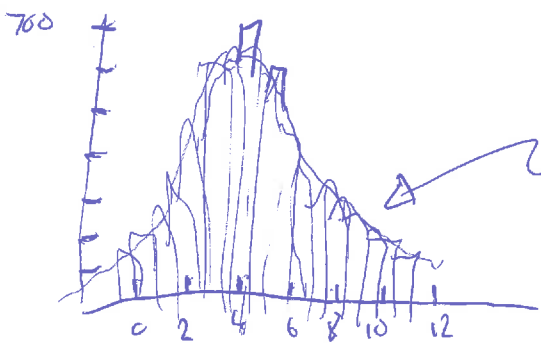
- fraud detection
- manufacturing (aircraft engines)
- monitoring machines in a data centre

- email spam classification
- weather prediction
- cancer classification

Week 9

Anomaly detection - choosing what features to use

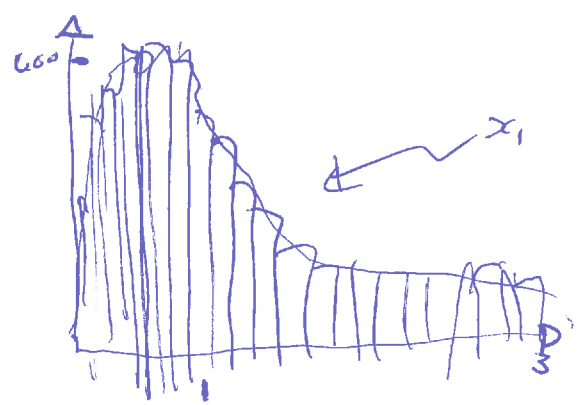
Non-gaussian features



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

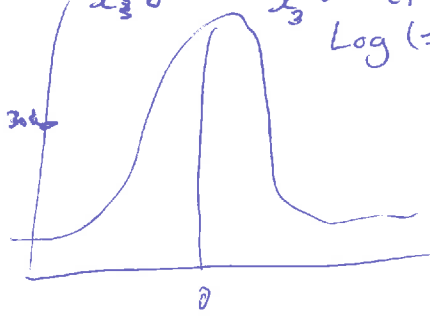
hist command to plot a histogram in octave

hist(x, 50) 50 bins.



log(x)

$x_1 \leftarrow \log(x_1)$
 $x_2 \leftarrow \log(x_2 + 1)$ *only with*
 $x_3 \leftarrow x_3$ or $\log(x_2 + c)$

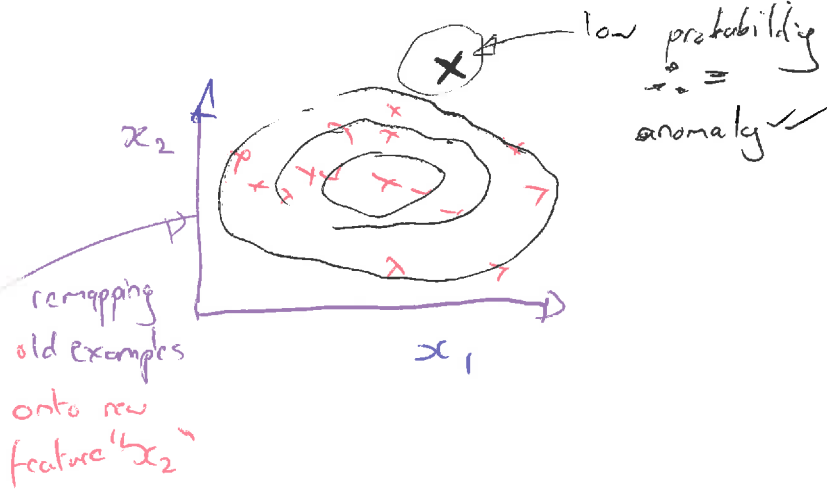
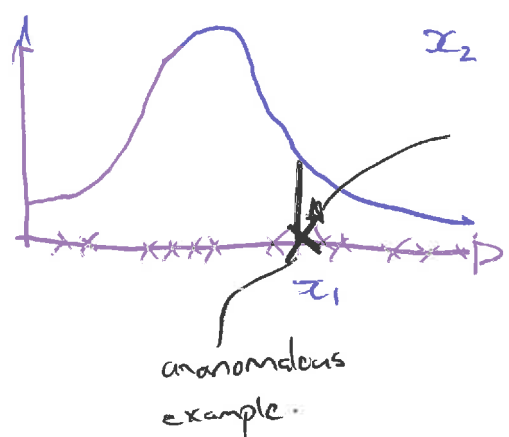


error Analysis for anomaly detection

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.



Week 9

Anomaly detection - choosing what features to use

Monitoring computers in a data centre

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 %

$$x_5 = \frac{\text{CPU LOAD}}{\text{Network traffic}}$$

(in the case of an infinite loop CPU LOAD $\rightarrow \infty$ whilst network traffic $\rightarrow 0$)

\therefore seems like a likely failure case

$$x_6 = \frac{(\text{CPU load})^2}{\text{network traffic}}$$

Week 9

Anomaly detection - Anomaly detection using the multivariate gaussian distribution

Multivariate Gaussian (normal) distribution

Parameter fitting:

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

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

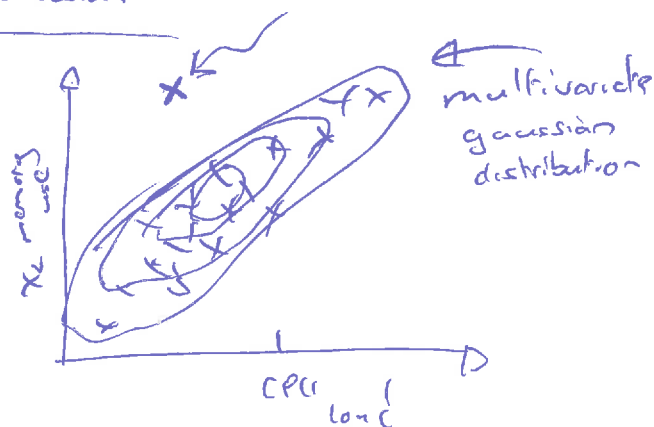
Anomaly detection with the multivariate Gaussian

1) Fit model $p(x)$ by setting

$$\mu = \dots$$

$$\Sigma = \dots$$

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



Relationship to original model

corresponds to multivariate Gaussian;

where $\Sigma = \begin{bmatrix} \sigma_1^2 & 0 & 0 \\ 0 & \sigma_2^2 & 0 \\ 0 & 0 & \sigma_3^2 \\ \vdots & \vdots & \vdots & \ddots \end{bmatrix}$

& exactly the same as the Gaussian model, however.

Original Model

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

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

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

Computationally scales better for large n (alternatively, cheaper)

OK even if m (training set size) is small

Multivariate Gaussian

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

$\Sigma \rightarrow \mathbb{R}^{n \times n}$
 Σ^{-1} computationally very expensive

Automatically captures correlations between features

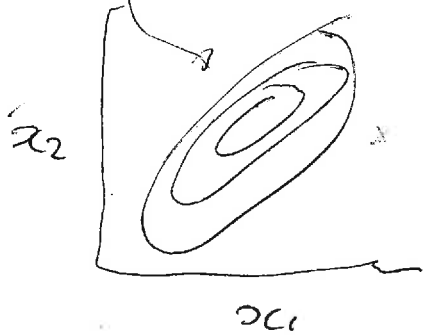
Computationally more expensive, look for redundant features

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

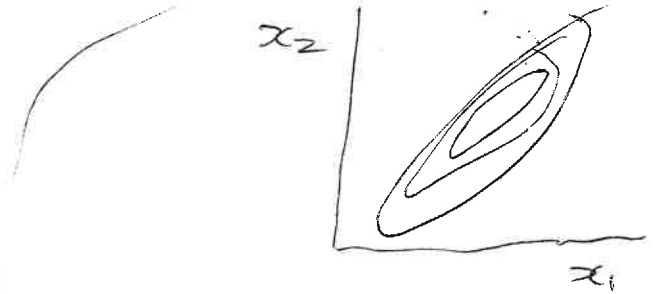
$\rightarrow m \geq 10n$ rule of thumb.

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 1 & 0.8 \\ 0.5 & 1 \end{bmatrix}$$

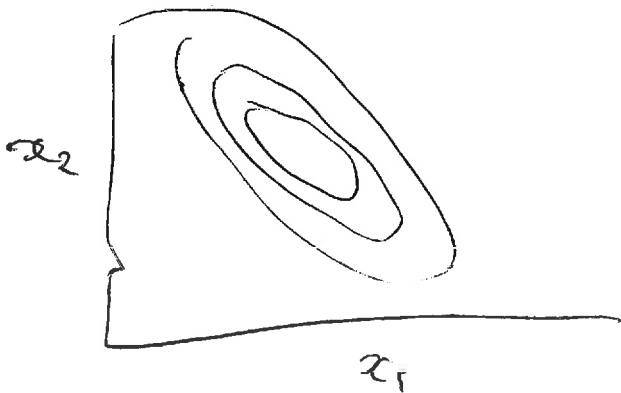
if the covariance matrix
skews the relationships between the
variables



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



$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 1 & -0.8 \\ -0.8 & 1 \end{bmatrix} \quad \text{---ve correlation}$$

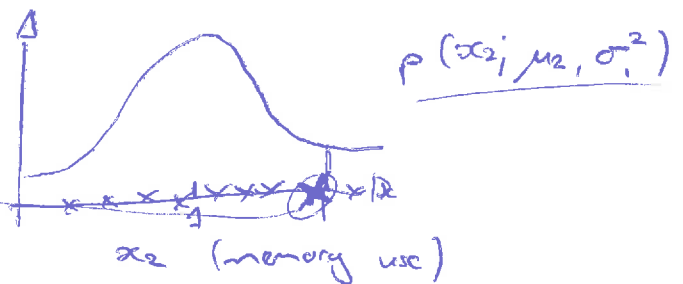
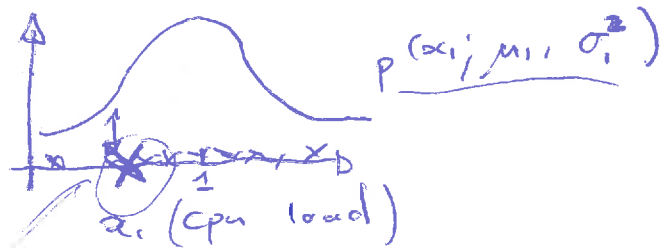
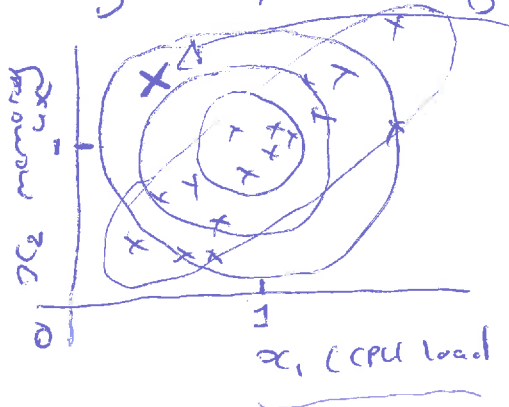


→ can also vary μ such that we move the
peak or centre (shift it).

Week 9

Anomaly Detection - Multivariate Gaussian distribution

motivating example: monitoring machines in a data center



doesn't
look too bad

∴ anomaly detection algorithm fails to classify data point.

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.

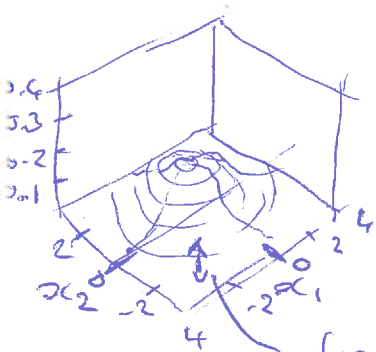
Parameters: $\mu \in \mathbb{R}^n, \Sigma \in \mathbb{R}^{n \times n}$ (covariance matrix)

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

↑
determinant of Σ

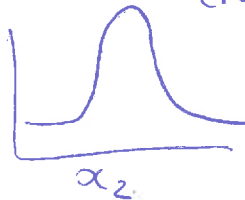
Multivariate Gaussian (Normal) examples

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



for a specific
value of μ & Σ
the height is $p(x)$.

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



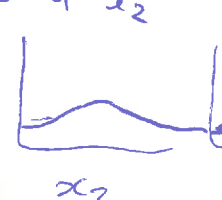
changes the variance of x_1

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

changes the variance of x_2



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



Week 9

Recommender Systems: problem formulation

example: predicting movie ratings
user ratings zero to five stars.

one star ~~★~~ → ★★★★★ five stars.

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)
Love at last	5	5	0	0
Romance forever	5	4.5	0	0
Cute puppies of bee	5	4	0	0
Nonstop car chases	0	0	5	4
Swords vs Karate	0	0	5	4

$n_u = 4$ $n_m = 5$

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

Week 9

Recommender systems:

Content based recommendations

Content based recommendation system

Movie	$\theta^{(1)}$ Alice (1)	$\theta^{(2)}$ Bob (2)	$\theta^{(3)}$ Carol (3)	$\theta^{(4)}$ Dave (4)	x_1 Romance	x_2 action
Love at first sight	5	5	0	0	0.9	0.1
Romance forever	5	?	?	0	1.0	0.0
Cute puppies of lag	?	4	0	?	0.99	0
Nonstop car chase	0	0	5	4	0.1	1.0
Swords vs karate	0	0	5	?	0	0.9

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

feature vector.

For each user j , learn a parameter $\theta^{(j)} \in \mathbb{R}^3$. Predict user j as rating movie i with $(\theta^{(j)})^T x^{(i)}$ stars $\hookrightarrow \theta^{(j)} \in \mathbb{R}^{n+1}$

$$n=2$$

number of features

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

(degree of romance, degree of action)

Problem formulation

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

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

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

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

for user j , movie i , predicting rating: $(\theta^{(j)})^T x^{(i)}$

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

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

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

summing over all the values that user j has rated.

don't regularize over the bias term.

Week 9

Recommender System: Content based recommendations

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

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

To learn $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(n_u)}$ for all users

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

↑
extra summation.

↳ in a nutshell gradient descent for recommendation. ..