# Anomaly detection

## Algorithm

Density estimation

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

Each example is $x \in \mathbb{R}$

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

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

$p(x)$

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

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

taking product of set of value

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

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

## Anomaly detection algorithm

1- Choose features $x_i$ that you think might be indicative of anomalous examples.          Then, given data set $\{x^{(1)}, \cdots, x^{(m)}\}$

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

$\mu_j = \frac{1}{m} \sum_{i=1}^{m} x_j^{(i)}$           $\longrightarrow$ Can be vectorized

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

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

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 $p(x) < \varepsilon$

# 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 labeled data, of anomalous and non-anomalous examples. ($y=0$ if normal, $y=1$ if anomalous).

Training set: $x^{(1)}, x^{(2)}, ..., x^{(m)}$ (assume normal examples / not anomalous) ⟵ unlabeled

Then, Define

Cross validation set: $(x_{cv}^{(1)}, y_{cv}^{(1)}), ..., (x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$ } to evaluate a particular anomaly detection algorithm.

Test set: $(x_{test}^{(1)}, y_{test}^{(1)}), ..., (x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$

We can include few examples into that contain examples that are known to be anomalous

### Aircraft engine motivating example

10000 good (normal) engines

20 flawed engines (anomalous)   20-50 examples  $y=1$

Training set: 6000 good engines ($y=0$) ⟹ $p(x) = p(x_1; \mu_1, \sigma_1^2) \cdots p(x_n; \mu_n, \sigma_n^2)$   (fit $p(x)$)

CV: 2000 good engines ($y=0$), 10 anomalous ($y=1$) } put anomalous just into cv and test set

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

### Alternative: (But not recommended)

Using the same data in cv and test set, not considered a good ML practice

Training set: 6000 good engines

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

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

even same set of anomalous examples in 2 sets

Put same 4000 into cv and test set, not good, as we like to think of cv and test set as completely different data sets to each other.

---

ε — a threshold that we would use to decide when to flag something as an anomaly.

## Algorithm evaluation

Fit model $p(x)$ on training set $\{x^{(1)}, ..., x^{(m)}\}$ ⟵ unlabeled examples but are those we're assuming are good and normal aircraft engines

On a cv / test example $x$, 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:

- True positive, false positive, false negative, true negative
- Precision / Recall
- $F_1$-score

keep trying many value and pick the ε that maximize $F_1$-score / or does well in cv set.

Can also use cross validation set to choose parameter $\varepsilon$

Then evaluate algorithm on cv set, and make all those decisions like what features did you use, how to set epsilon, use to evaluate algorithm on cv set, then when we've picked set of features, and value of epsilon that good, we take the final model and do final evaluation of algorithm on test sets.

Quiz: Suppose you have fit a model $p(x)$. When evaluating on cv or test set, your algorithm predict: $y = \begin{cases} 1 & \text{if } p(x) \leq \varepsilon \\ 0 & \text{if } p(x) > \varepsilon \end{cases}$

Is classification accuracy a good way to measure the algorithm's performance?

No, because of skewed classes (so an algorithm that always predict $y=0$ will have high accuracy).

# Anomaly detection

## Anomaly detection vs. supervised learning

| Anomaly detection ⟶ save for CV and test set    vs. | Supervised learning |
|---|---|
| Very small number of (positive examples)(y=1). (0-20 is common) Large number of negative (y=0) examples. (⟶ to fit $p(x)$) | Large number of positive and negative examples. |
| Many different "types" of anomalies. Hard for any algorithm to learn from small sets of positive examples what the anomalies look like; | 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. |

future anomalies may look nothing like any of the anomalous examples we've seen so far. (Maybe in your set of positive examples, you've seen 5 or 10 or 20 ways that an aircraft engine could go wrong, but maybe on another day you need to detect a totally new set/ new type of anomaly / totally new way for an aircraft engine to be broken that you've just never seen before)

So it would be better to just model negative examples(y=0) with $p(x)$ rather than try to model positive examples, as another day anomaly may be nothing like the one you've seen so far.

Application:
- Fraud detection (people out there have a lot of ways to do fraud) ⟵ or maybe you have a lot of y=1

⟶ then Fraud detection can be in supervised learning problem

Application:
- Email spam classification
- Manufacturing (e.g. aircraft engines)
- Weather prediction (sunny/rain/etc)
- Monitoring machines in a data center
- Cancer classification

⋮                    ⋮

Some anomaly detection can be shifted to supervised learning if the positive and negative examples enough much.

---

## Anomaly detection algorithm
- You run a power utility (supplying electricity to customers) and want to monitor your electric plants to see if any one of them might be behaving strangely.
- A computer vision/security application, where you examine video images to see if anyone in your company's parking lot is acting in an unusual way.

## Anomaly detection - Choosing what features to use
If the graph looks non-gaussian, take the features to log it, make gaussian
$$x_1 \to \log(x_1) \quad \text{or} \quad x_2 \to \log(x_2+1) \;/\; \log(x_2+c) \;/\; x_4 \to x_4^{\frac{1}{2}} \;/\; x_3 \to x_3^{\frac{1}{3}}$$
(⟶ constant)
So to make it looks as Gaussian as possible.

Plot data in Octave ⟶ hist(X) (histogram) ← check whether Gaussian or not
if not then, hist(x.^0.1, 50)
(⟶ to make the graph show in 50 bins)
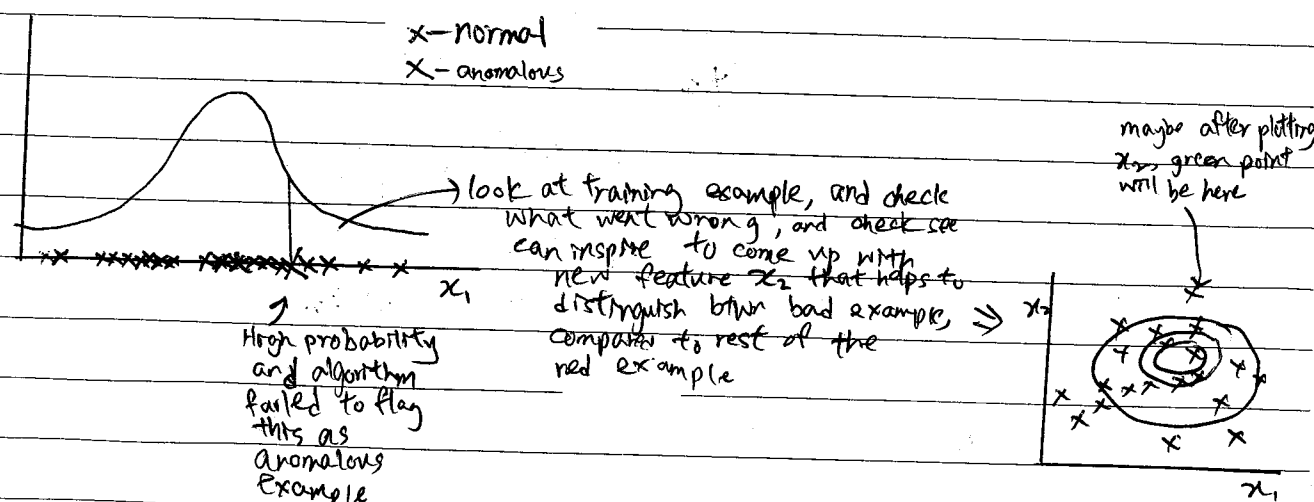If it is Gaussian enough,
then xNew = x.^0.1;

## 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.

x - normal
x - anomalous



⟶ look at training example, and check what went wrong, and check see can inspire to come up with new feature $x_2$ that helps to distinguish btwn bad example, ≫ compared to rest of the red example

High probability and algorithm failed to flag this as anomalous example

maybe after plotting $x_2$, green point will be here

## Monitoring Computers in a data center

Choose features that might take on unusually large or unusually 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 ←  ⎫ this 2 work linearly, but maybe one
$x_4$ = network traffic ← ⎭ failure case which code stuck in an infinite load, CPU load grow and network traffic doesn't

new feature $x_5 = \dfrac{CPU\ load}{network\ traffic}$ ⟹ help in anomaly detection

OR

$x_6 = \dfrac{(CPU\ load)^2}{network\ traffic}$   (can come out with new features to capture different sort of anomaly examples.)

Quiz: Suppose your anomaly detection algorithm is performing poorly and outputs a large value of $p(x)$ for many normal examples and for many anomalous examples in your cross validation dataset. Which of the following changes to your algorithm is most likely to help?

— Try coming up with more features to distinguish between the normal and the anomalous examples.

## Anomaly detection — Multivariate Gaussian distribution



this is what we think the probability in

But computer will think inside circle has probability

But actually $x_1$ has much more higher probability to be normal than $x_2$

We are going to fix this with Multivariate Gaussian

## Multivariate Gaussian (Normal) distribution

$x \in \mathbb{R}^n$. Don't model $p(x_1), p(x_2), \ldots,$ etc. separately.

Model $p(x)$ all in one go.

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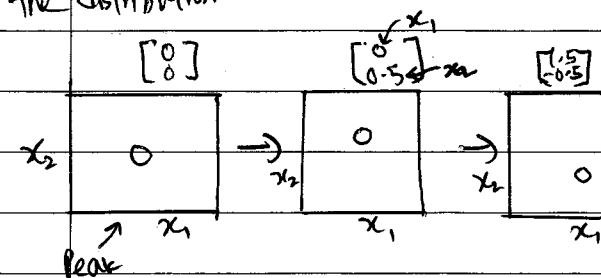
$$p(x; \mu, \Sigma)$$

Probability

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

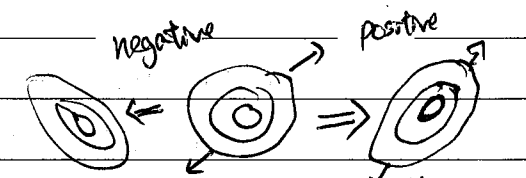$|\Sigma|$ = determinant of $\Sigma$  | Octave: det(sigma)



negative    positive

increasing this can change the circle

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

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

control variance of $x_1$

control variance of $x_2$

The distribution become narrower and taller.

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

The distribution become wider and shorter (positive)

Change the peak of the distribution

normal distribution

$x_2$

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

$\begin{bmatrix} 0.5 \\ 0 \end{bmatrix}$  $x_1$

$\begin{bmatrix} 1.5 \\ 0.5 \end{bmatrix}$

peak

# Anomaly detection

## Anomaly detection using the multivariate Gaussian distribution

Multivariate Gaussian (Normal) distribution

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

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

$\Rightarrow$ Can vary $\mu$ and $\Sigma$, can get a range of different distributions

Parameter fitting:

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

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

Anomaly detection with the multivariate Gaussian

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

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

2- Given a new example $x$, compute

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

Flag an anomaly if $p(x) < \varepsilon$

Relationship to original model

Original model: $p(x) = p(x_1; \mu_1, \sigma_1^2) \times p(x_2; \mu_2, \sigma_2^2) \times \ldots \times p(x_n; \mu_n, \sigma_n^2)$
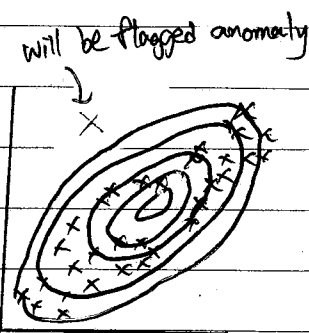
- The contours of $p(x)$ are axis aligned

Corresponds to multivariate Gaussian

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

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

all the numbers in the off diagonal are "0" zero

will be flagged anomaly



| | Original model | vs. | Multivariate Gaussian |
|---|---|---|---|
| | $p(x_1; \mu_1, \sigma_1^2) \times \ldots \times p(x_n; \mu_n, \sigma_n^2)$ | | $p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{n}{2}}|\Sigma|^{\frac{1}{2}}}\exp\left(-\frac{1}{2}(x-\mu)^T\Sigma^{-1}(x-\mu)\right)$ |
| | Manually create features to capture anomalies where $x_1, x_2$ take unusual combinations of values. | | Automatically captures correlations between features |
| | Create new feature $x_3 = \frac{x_1}{x_2} = \frac{CPU\ load}{memory}$ | | If $n=10,000$   $\Sigma \in \mathbb{R}^{n \times n}$ $\Longleftarrow$ will be |
| | Computationally cheaper (alternatively, scales better to large $n$) | | Computationally more expensive |
| | e.g. $n=10,000$, $n=100,000$ still work fine | | $m < n$ / Redundant features like $x_1 = x_2$, $x_3 = x_4 + x_5$   Can make $\Sigma$ non-invertible |
| | OK even if $m$ (training set size) is small | | Must have $m > n$, or else $\Sigma$ is non-invertible. *Will only use when $m \gg n$ / $m \geq 10n$ |

Recommender Systems — Problem formulation

Example: Predicting movie ratings

User rates movies using zero to five stars

? – Didn't watch and didn't rate, doesn't have rating

| Movie | Alice(1) | Bob(2) | Carol(3) | Dave(4) |
|---|---|---|---|---|
| Love at last | 5 | 5 | 0 | 0 |
| Romance forever | 5 (predict maybe) | ? → 4/5 | ? 0 | 0 |
| Cute puppies of love | ? 5 | 4 | 0 | ? 0 |
| Nonstop car chases | 0 | 0 | 5 | 4 |
| Swords vs. karate | 0 | 0 | 5 | ? 4 ← predict |

Romantic comedy movies {Love at last, Romance forever, Cute puppies of love}

Action movies {Nonstop car chases, Swords vs. karate}

$n_u$ = no. users     $n_u = 4$   $n_m = 5$

$n_m$ = 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 $r(i,j) = 1$)   $0, \ldots, 5$

$\Rightarrow$ Recommender system problem is given this dataset, that give these $r(i,j)$ and $y^{(i,j)}$ to look through data and look at all movie rating that are missing, try to predict what value of "?" should be.

$\Rightarrow$ Automatically fill in missing value, so can recommend new movie that users haven't watch to them, predict what can be interesting to users.

# Recommender Systems

## Content-based recommendations

### Content-based recommender systems $n_u=4$, $n_m=5$, $x_0=1$

| Movie | | Alice (1) $\theta^{(1)}$ | Bob(2) $\theta^{(2)}$ | Carol(3) $\theta^{(3)}$ | Dave(4) $\theta^{(4)}$ | $x_1$ (romance) | $x_2$ (action) |
|---|---|---|---|---|---|---|---|
| $x^{(1)}$ | Love at last 1 | 5 | 5 | 0 | 0 | 0.9 | 0 |
| $x^{(2)}$ | Romance forever 2 | 5 | (?) | (?) | 0 | 1.0 | 0.01 |
| $x^{(3)}$ | Cute puppies of love 3 | (?) 4.95 | 4 | 0 | (?) | 0.99 | 0 |
| $x^{(4)}$ | Nonstop car chases 4 | 0 | 0 | 5 | 4 | 0.1 | 1.0 |
| $x^{(5)}$ | Swords vs. karate 5 | 0 | 0 | 5 | (?) | 0 | 0.9 |

$x_1$ — measures degree to which a movie is a romantic movie

$x_2$ — measures degree to which a movie is an action movie

Have features, each movie can be represented with a feature vector.

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

$n$ — number of features not counting $x_0$

In this case, $n=2$

$\mathbb{R}^{n+1}$ ($+x_0$)

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.

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

For this example, let's say we found out parameter vector for Alice $\theta^{(1)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix}$

Our prediction for "?" in Alice $= (\theta^{(1)})^T x^{(3)}$
$= 5 \times 0.99$
$= 4.95$

## Problem formulation

$r(i,j) = 1$ if user $j$ has rated movie $i$ (0 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 $i$, predicted rating: $(\theta^{(j)})^T (x^{(i)})$

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

number of features we have per movie

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

To learn $\theta^{(j)}$:  To simplify the subsequent math  prediction of user $j$ rating on movie $i$

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

Like ⟹ linear regression

Just a constant, doesn't change $J(\theta)$

sum over all movies user $j$ has rated (summation over all values of $i$, so $r(i,j)$ is equal to 1)

actual user rating

Don't regularize bias term

## Optimization objective

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

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

To learn $\theta^{(1)}, \theta^{(2)}, \ldots, \theta^{(n_u)}$:

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

## Optimization algorithm:

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

$J(\theta^{(1)}, \ldots, \theta^{(n_u)})$

## Gradient descent update:

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

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

$\frac{\partial}{\partial \theta_k^{(j)}} J(\theta^{(1)}, \ldots, \theta^{(n_u)})$

Can plug into more advanced optimization algorithm

# Collaborative filtering

## Problem motivation

| Movie | Alice(1) $\theta^{(1)}$ | Bob(2) $\theta^{(2)}$ | Carol(3) $\theta^{(3)}$ | Dave(4) $\theta^{(4)}$ | $x_1$ (romance) | $x_2$ (action) |
|---|---|---|---|---|---|---|
| Love at last | 5 | 5 | 0 | 0 | 0.9 ? | 0 ? |
| Romance forever | 5 | ? | ? | 0 | 1.0 ? | 0.01 ? |
| Cute puppies of love | ? | 4 | 0 | ? | 0.99 ? | 0 ? |
| Non stop car chases | 0 | 0 | 5 | 4 | 0.1 ? | 1.0 ? |
| Swords vs karate | 0 | 0 | 5 | ? | 0 ? | 0.9 ? |

$x^{(1)}$ (left label)

↑ Suppose we don't have the features data set

let's say we've gone to each of
our users, each of the users told us
how much they like the romantic movies and
how much they like action packed movies. } If can get these, can infer value of $x_1$ and $x_2$

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

Both of them love romantic movie but hate action movie

Both of them love action movie but hates romantic movie

What $x^{(1)}$ should be so that,
$(\theta^{(1)})^T x^{(1)} \approx 5$
$(\theta^{(2)})^T x^{(1)} \approx 5$
$(\theta^{(3)})^T x^{(1)} \approx 0$
$(\theta^{(4)})^T x^{(1)} \approx 0$

↳ can conclude $x^{(1)}$ as $[\,X_1 = 1.0, \; X_2 = 0.0\,]$

$$x^{(1)} = \begin{bmatrix} 1.0 \\ 0.0 \end{bmatrix}$$

and can keep on to figure out $x_1, x_2$ in movie $x^{(2)} x^{(3)}, x^{(4)}, x^{(5)}$

**Quiz**

| | User 1 | User 2 | User 3 | (romance) |
|---|---|---|---|---|
| Movie 1 | 0 | 1.5 | 2.5 | ? |

Note that there is only one feature $x_1$, suppose that.
$$\theta^{(1)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \theta^{(2)} = \begin{bmatrix} 0 \\ 3 \end{bmatrix} \quad \theta^{(3)} = \begin{bmatrix} 0 \\ 5 \end{bmatrix}$$

What would be a reasonable value for $x_1^{(1)}$?
$\Rightarrow 0.5$

## Optimization algorithm

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

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

**Quiz** Suppose you use gradient descent to minimize:
$$\min_{x^{(1)},...,x^{(n_m)}} \frac{1}{2} \sum_{i=1}^{n_m} \sum_{j:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^{n} (x_k^{(i)})^2$$

Correct gradient descent update rule for $i \neq 0$
$$\Rightarrow x_k^{(i)} := x_k^{(i)} - \alpha \left( \sum_{j:r(i,j)=1} ((\theta^{(j)})^T (x^{(i)}) - y^{(i,j)}) \theta_k^{(j)} + \lambda x_k^{(i)} \right)$$

## Collaborative filtering

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

**Parameters (By user)** Given $\theta^{(1)}, ..., \theta^{(n_u)}$
    can estimate $x^{(1)}, ..., x^{(n_m)}$

Can start by making initial guess on $\theta \xrightarrow{fit} x \xrightarrow{\substack{Make \\ better \\ changes}} \theta \xrightarrow{\substack{Better \\ fit}} x \rightarrow \theta \rightarrow x \rightarrow ...$

keep iterating and optimizing
$\Rightarrow$ will cause the algorithm to converge

Recommender Systems

## Collaborative Filtering algorithm

Collaborative filtering optimization objective

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

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

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

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

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

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

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

instead of going
$\theta \rightarrow x \rightarrow \theta \rightarrow x \rightarrow \theta$ ----- we minimize with respect to both sets of parameters simultaneously

~~$x_0$~~ $x \in \mathbb{R}^n$ ~~$\theta \in \mathbb{R}^{n+1}$~~
$\theta \in \mathbb{R}^n$

no need $x_0 = 1$ as it the algorithm wants the feature
to be 1, it can learn itself and can set $x_1 = 1$

## Collaborative filtering algorithm   ~~$x_0$~~ ~~$x_1$~~ $x \in \mathbb{R}^n$ $\theta \in \mathbb{R}^n$

1. Initialize $x^{(1)}, ..., x^{(n_m)}, \theta^{(1)}, ..., \theta^{(n_u)}$ to small random values.

2. Minimize $J(x^{(1)}, ..., x^{(n_m)}, \theta^{(1)}, ..., \theta^{(n_u)})$ using gradient descent (or an advanced optimization algorithm).

E.g. for every $j = 1, ..., n_u, i = 1, ..., n_m$:
$$x_k^{(i)} := x_k^{(i)} - \alpha \left( \sum_{j:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) \theta_k^{(j)} + \lambda x_k^{(i)} \right) \quad \text{Partial derivative value}$$
$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \left( \sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} + \lambda \theta_k^{(j)} \right)$$

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

Predict User $j$ is going to rate movie $i$ with $(\theta^{(j)})^T (x^{(i)})$

**Quiz:** In algorithm, we initialized $x^{(1)}, ..., x^{(n_m)}$ and $\theta^{(1)}, ..., \theta^{(n_m)}$ to small random values. Why?
This serves as symmetry breaking (similar to the random initialization of a neural network's parameters) and ensures the algorithm learns features $x^{(1)}, ..., x^{(n_m)}$ that are different from each other.

## Vectorization: Low rank matrix factorization

| 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 | ? |
| Nonstop car chases | 0 | 0 | 5 | 4 |
| Swords vs. karate | 0 | 0 | 5 | ? |

$n_m = 5$
$n_u = 4$

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

$5 \times 4$

$\uparrow y^{(i,j)}$

## Collaborative filtering

Predicted rating: $(\theta^{(j)})^T (x^{(i)})_{(i,j)} \downarrow$

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

$$\begin{bmatrix} (\theta^{(1)})^T(x^{(1)}) & (\theta^{(2)})^T(x^{(1)}) & \cdots & (\theta^{(n_u)})^T(x^{(1)}) \\ (\theta^{(1)})^T(x^{(2)}) & (\theta^{(2)})^T(x^{(2)}) & \cdots & (\theta^{(n_u)})^T(x^{(2)}) \\ \vdots & \vdots & & \vdots \\ (\theta^{(1)})^T(x^{(n_m)}) & (\theta^{(2)})^T(x^{(n_m)}) & \cdots & (\theta^{(n_u)})^T(x^{(n_m)}) \end{bmatrix}$$

$X \Theta^T \nearrow$

$$X = \begin{bmatrix} -(x^{(1)})^T- \\ -(x^{(2)})^T- \\ \vdots \\ -(x^{(n_m)})^T- \end{bmatrix} \qquad \textcircled{H} = \begin{bmatrix} -(\theta^{(1)})^T- \\ -(\theta^{(2)})^T- \\ \vdots \\ -(\theta^{(n_u)})^T- \end{bmatrix}$$

$\uparrow$ Low rank matrix factorization $\nearrow$

### Finding related movies

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

$\Rightarrow$ Let's say we have $x_1 = $ romance, $x_2 = $ action, $x_3 = $ comedy, ....

How to find movies $j$ related to movie $i$?

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

5 most similar movies to movie $i$:

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

$\nearrow$ smallest distance between features between these different movies

---

## Implementational detail: Mean normalization

Users who have not rated any movies

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

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

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

$\underbrace{\qquad}$ Don't have value, so going to be 0

$n=2 \qquad \theta^{(5)} \in \mathbb{R}^2$

$\underbrace{\qquad} \frac{\lambda}{2}\left[(\theta_1^{(5)})^2 + (\theta_2^{(5)})^2\right]$

$\theta^{(5)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \longleftarrow$ will regularize and make it near to 0 as possible

$(\theta^{(5)})^T x^{(1)} = 0$

$\xrightarrow{\qquad}$ Not good.

### Mean Normalization:

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

$\downarrow$ Use as actual rating
Learn $\theta^{(j)}, x^{(i)}$

For user $j$, on movie $i$ predict:

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

$\therefore$ Movie rating are comparable (0-5) are on similar scales so no need to do feature scaling.

User 5 (Eve):

$\theta^{(5)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \qquad (\theta^{(5)})^T(x^{(i)}) + \mu_i$
$$= 0 + \mu_i$$

$\Rightarrow$ If user hasn't rated any movies, we just predict for each of the movies, the average rating that those movies got