

Week 9

Anomaly Detection

Dataset = $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ \leftarrow 本 (normal examples)

$x_{\text{test}} \leftarrow P(x_{\text{test}}) < \epsilon \rightarrow \text{anomaly}$

$P(x_{\text{test}}) \geq \epsilon \rightarrow \text{ok}$

Fraud detection

$x^{(i)}$ = features of user i 's activities $x^{(i)} = \{x_1^{(i)}, x_2^{(i)}, x_3^{(i)}, x_4^{(i)}\}$

Model $p(x)$ from data.

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

[講義中の Quiz]

異常検知システムは $p(x) \leq \epsilon$ のとき x のフラグを立てる。

システムがフラグを立てやすい時、どうする? \leftarrow false positive
(教師あり学習での場合)

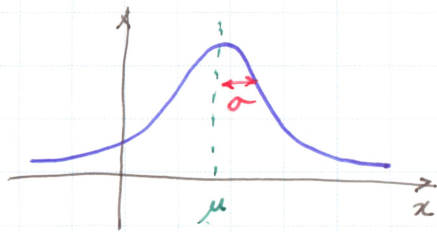
$p(x) \leq \epsilon$ を満たす x を減らしたい $\rightarrow \epsilon$ を小さく

Week 9

2

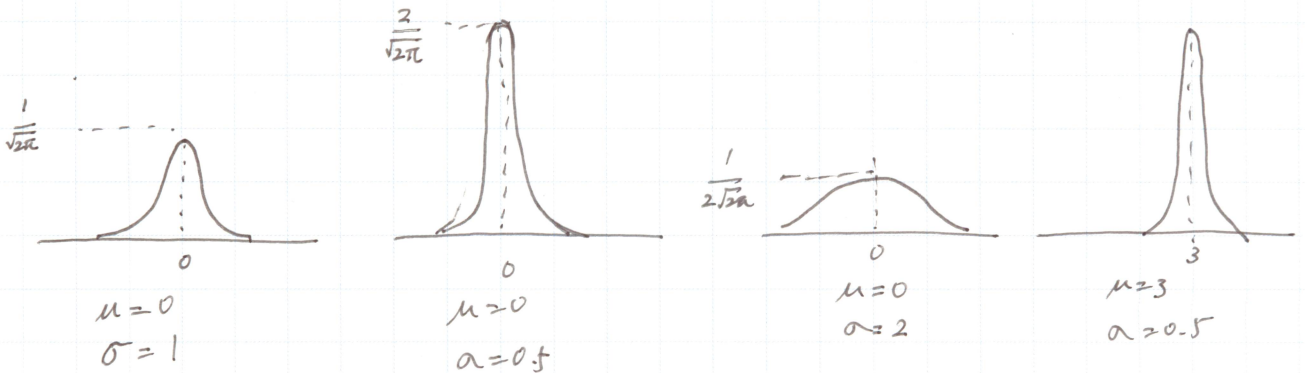
Gaussian Distribution (= Normal Distribution)
正規分布

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



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

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



Parameter estimation

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

確率変数 x が 1次元の場合

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

と表れる

$x^{(i)} \in \mathbb{R}^n$

n 次元の場合

$$x \sim \mathcal{N}_n(\mu, \Sigma)$$

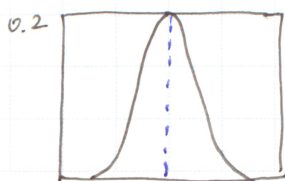
与えられた dataset に対して μ と σ^2 の値を推定したい

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

統計では $(m-1)$ を使うが、機械学習では m を使う

[講義中の Quiz]

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



$-3 = \mu$

$$\frac{1}{\sqrt{2\pi}\sigma} = 0.2 = \frac{1}{5}$$

$$\frac{6.28 \times 12500}{2512} = 12$$

$$\therefore \sqrt{2\pi}\sigma = 5$$

$$\sigma = \sqrt{\frac{25}{2\pi}} \approx \sqrt{4} = 2$$

$$\therefore p(x) = \frac{1}{\sqrt{2\pi} \cdot 2} \exp\left(-\frac{(x+3)^2}{2 \times 4}\right)$$

Week 9 3

Density estimation (確率密度関数の推定)

Training set : $\{x^{(1)}, \dots, x^{(m)}\}$
 $x \in \mathbb{R}^n$

$$\begin{cases} x_1 \sim N(\mu_1, \sigma_1^2) \\ \vdots \\ x_n \sim N(\mu_n, \sigma_n^2) \end{cases}$$

$$\begin{aligned} P(x) &= P(x_1; \mu_1, \sigma_1^2) P(x_2; \mu_2, \sigma_2^2) \dots P(x_n; \mu_n, \sigma_n^2) \\ &= \prod_{j=1}^n P(x_j; \mu_j, \sigma_j^2) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi} \sigma_j} \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right) \end{aligned}$$

[講義中の Quiz]

Anomaly if $P(x) \leq \epsilon$

$$\{x^{(1)}, \dots, x^{(m)}\}, \quad \mu_j \in \mathbb{R}, \quad \sigma_j^2 \in \mathbb{R}$$

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

Week 9 4

Building an Anomaly Detection System

real-number evaluation は重要

labeled data (教師ありデータ)

$\begin{cases} y=0 & \text{if normal} \\ y=1 & \text{if anomaly} \end{cases}$

Training set $\{x^{(1)}, \dots, x^{(m)}\}$ ← 全て normal データ

Cross validation set $(x_{cv}^{(1)}, y_{cv}^{(1)}), \dots, (x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$ ← anomaly データを含む

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

航空機のエンジンの例

$\begin{cases} 10000 & \leftarrow \text{good engine (normal)} \\ 20 & \leftarrow \text{flawed engine (anomaly)} \end{cases}$

Training set: 6000 good

CV set: 2000 good ($y=0$), 10 anomaly ($y=1$)

Test set: 2000 good ($y=0$), 10 anomaly ($y=1$)

Algorithm evaluation

モデル $p(x)$ を training set $\{x^{(1)}, \dots, x^{(m)}\}$ に fit させる。

cross validation set に対して

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

評価基準

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

各パラメータを決定するに CV set を使う。

		実際	
		1	0
予測	1	true positive	false positive
	0	false negative	true negative

注意
 - accuracyは特定の
 計算方法を指している

$$\frac{(1) + (4)}{(1) + (2) + (3) + (4)}$$

	actual	
	1	0
予 1	(1)	(2)
予 0	(3)	(4)

[講義中の Quiz]

$$y = \begin{cases} 1 & \text{if } p(x) \leq \epsilon \\ 0 & \text{if } p(x) > \epsilon \end{cases}$$

モデルの性能を評価するのに classification accuracyはどの方法か?

- ① Yes. cv & test set 両方 labelが同じ X
(答, 教師)
- ② No. cv & test set 両方 labelが同じ X
- ③ No. skewed class があるため $y=0$ が多い結果となる。
- ④ No. for the cv, yes for test. X

Anomaly Detection vs. Supervised Learning

Anomaly detection

positive 事例がほとんど少ない (0~20個)
 $y=1$

negative 事例 ($y=0$) はたくさんある。

anomaly はいろいろな type があるので
positive 事例が少くは anomaly 検出が難しい。
学習するのは難しい。

将来出会う anomaly は、今まで出会った
anomalies とは全く異なっている。

Fraud detection

Manufacturing

Monitoring machines in data center

vs. Supervised learning

positive 事例も negative 事例もたくさんある。

positive と学習するための十分な量の
positive 事例が与えられる。
未来に出会う positive 事例も、ほとんど
過去の positive と似ている。

Email spam classification

Weather prediction (sunny/rainy/etc)

Cancer classification

[講義中のQuiz]

anomaly detection を用いる (supervised learning ではない) 1つ問題を選び

① 電源設備を動作させている、動作がおかしいのモニターする ○

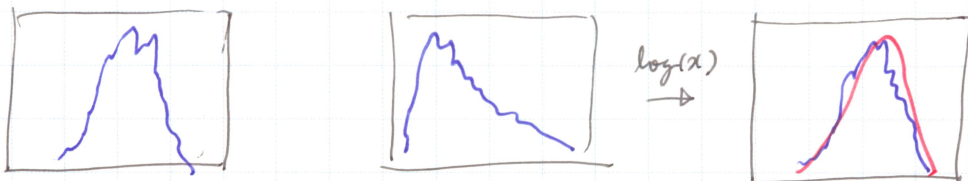
② = 明日の必要電力を予測する ×

④ computer vision / security アプリケーション, 店に入ってきた人の男か女の
1つ-2つを予測する ×

③ = 店の駐車場の普通では無い使い方をとっている
video から予測する ○

Choosing what features to use

Non - Gaussian features



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

ガウス分布(正規分布)に似せる。

Week 9

7

Error analysis for anomaly detection

Want $p(x)$ $\left\{ \begin{array}{l} \text{large for normal example } x \\ \text{small for anomalous } \end{array} \right. =$

Most Common Problem

$p(x)$ is comparable $\left\{ \begin{array}{l} \text{normal example} \\ \text{anomalous } \end{array} \right. =$

[講義中の Quiz]

anomaly detection \rightarrow 性能悪い

normal example に対し $p(x)$ 大

anomalies in CV に対し $p(x)$ 大

どう変更するよいか?

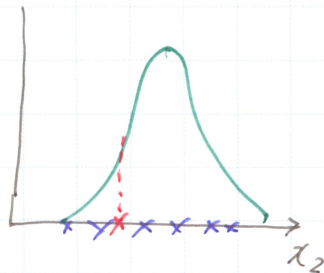
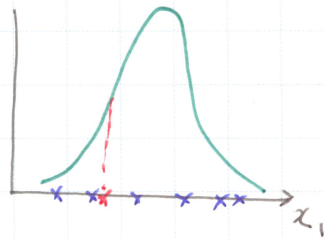
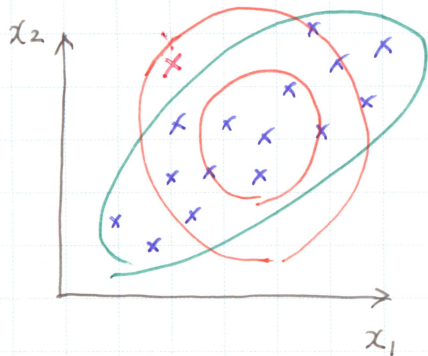
① 少い feature \rightarrow X

② normal example と abnormal example とを区別する為に \rightarrow ϵ 多 (a feature ϵ 増える) X
○ 不正解

③ 大量の (normal example の) training set ϵ 追加する X

④ ϵ をかえる ~~X~~ ϵ を大きくする

Multi-variate Gaussian Distribution (Optional)



x_1, x_2, \dots それぞれで正規分布を判断しては上の \times は anomaly detection できない

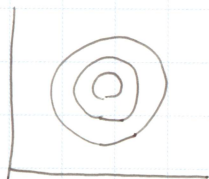
$x \in \mathbb{R}^n$. $P(x_1), P(x_2), \dots$ を個別に model 化してはいけない.

全ての x について一度にモデル化 $P(x)$ が正しい.

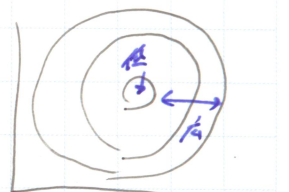
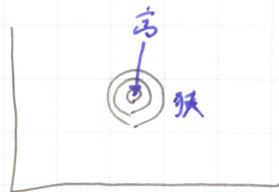
パラメータ: $\mu \in \mathbb{R}^n$, $\Sigma \in \mathbb{R}^{n \times n}$ (covariance matrix)
共分散行列

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

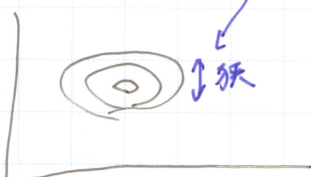
$$\mu = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \Sigma = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$



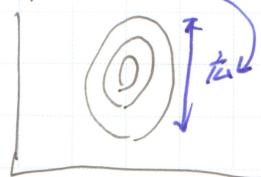
$$\mu = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \Sigma = \begin{pmatrix} 0.6 & 0 \\ 0 & 0.6 \end{pmatrix}, \mu = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \Sigma = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}$$



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



$$\mu = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \Sigma = \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}, \mu = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \Sigma = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 1 \end{pmatrix}$$

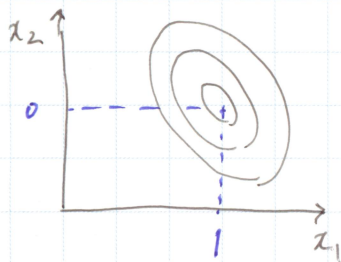


Week 9 9

$$\mu = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \Sigma = \begin{pmatrix} 1 & -0.5 \\ -0.5 & 1 \end{pmatrix}, \quad \mu = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \Sigma = \begin{pmatrix} 1 & -0.8 \\ -0.8 & 0 \end{pmatrix}$$



[講義中の Quiz]



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

Anomaly Detection using the Multivariate Gaussian Distribution
多変量正規分布

パラメータ μ, Σ

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

パラメータ fitting

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

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

これら得られた μ, Σ を anomaly detection に利用する。

Anomaly detection with the multivariate Gaussian

1. model $p(x)$ を fit する

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

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

2. 新しい example x が与えられるので $p(x)$ を計算する

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

if $p(x) < \epsilon$, フラグを立てる.

Relationship to original model

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

Corresponds to multivariate Gaussian

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

where

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

↑ feature 間の相関を自動的に扱うことができる

○ $m > n$ が必須 ($m > 10n$ が理想)
 x が与えられたときの逆行列が存在しない。

○ feature 間の相関を扱うための $x_3 = \frac{x_2}{x_1}$ のように自分で考えた相関をうまく表す feature を作って追加する必要はある

○ 計算量は小さい ($n=10000$ 以上でも OK)

○ m (training set の #) が小さくても OK

特異行列 Σ が singular なのは $\left\{ \begin{array}{l} m < n \text{ の時 } \text{ or } \\ \text{feature が冗長の場合} \end{array} \right.$ 楕円に属している
 非可逆

[講義中の Quiz]

$\{x^{(1)}, \dots, x^{(n)}\}$ where $x^{(i)} \in \mathbb{R}^n$, anomaly detection を適用する。

① original model $\prod_{i=1}^n p(x_i; \mu_i, \sigma_i^2)$ は $p(x; \mu, \Sigma)$ の輪郭が一軸 aligned

なるは multivariate Gaussian に対応している

~~*~~ feature の間に相関は一切なければ

~~*~~ Σ 軸に沿っているということは
相関がないということ

Σ が非可逆に子子の下

② multivariate Gaussian は $m < n$ の時 有利? ある \times

③ multivariate Gaussian は feature 間の相関を自動的に扱ってくれる \circ

④ original model は 計算量的に少なくて済む。 n が大きい時に有利 \circ

Predicting Movie Rating

Problem Formulation

Recommendation System ... 機械学習のよい応用例
 featureを自動的に選択する例になる。

例) movie rating

movie \ user	user1	2	3	4
movie 1	5	5	0	0
2	5	?	?	0
3	2	4	0	?
4	0	0	5	4
5	0	0	5	?

← 星の数

n_u : ユーザー数

n_m : 映画の数

$r(i, j) = 1$ if ユーザーjが映画iをratingしている

$r^{(i, j)}$ = rating (ユーザーj, 映画i)

[講義中の Quiz]

	user		
movie	1	2	3
1	0	1	?
2	?	5	5

movie user

↓ ↓
 $r(2, 1) = 0$

$r^{(2, 1)} = ?$

Content Based Recommendations

例

movie \ user	1	2	3	4	romance x_1	action x_2
movie 1	5	5	0	0	0.9	0
2	5	?	?	0	1.0	0.01
3	?	4	0	?	0.99	0
4	0	0	5	4	0.1	1.0
5	0	0	5	?	0	0.9

$x^{(1)} = \begin{pmatrix} 1 \\ 0.9 \\ 0 \end{pmatrix}$ Bias?

$n_u = 4, n_m = 5$

$n = 2$

各ユーザー j に対して, パラメータ $\theta^{(j)} \in \mathbb{R}^3$ を学習させる.
 ユーザー j について, 映画 i の rating を $(\theta^{(j)})^T x^{(i)}$ と予測する.

[使用例] $x^{(3)} = \begin{pmatrix} 1 \\ 0.99 \\ 0 \end{pmatrix}$, $\therefore \theta^{(1)} = \begin{pmatrix} 0 \\ 5 \\ 0 \end{pmatrix}$ と学習したのであれば

$y^{(1,3)} = \theta^{(1)T} x^{(3)} = (0 \ 5 \ 0) \begin{pmatrix} 1 \\ 0.99 \\ 0 \end{pmatrix} = 4.95$

[講義中の Quiz]

Movie \ user	1	2	3	4	romance x_1	action x_2
Movie 1	5	5	0	0	0.9	0
2	5	?	?	0	1.0	0.01
3	?	4	0	?	0.99	0
4	0	0	5	4	0.1	1.0
5	0	0	5	?	0	0.9

$i = 1, 2$
 $x^{(i)} = \begin{pmatrix} 1 \\ \square \\ \text{小} \end{pmatrix}$

$i = 3, 4$
 $x^{(i)} = \begin{pmatrix} 1 \\ \text{小} \\ \square \end{pmatrix}$

この rating を predict する $\theta^{(j)}$

$x_0 = 1$. $\theta^{(3)}$ として適切なのはどれ

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

\therefore ユーザー 3 は ユーザー 4 と同じクラスで、これを高く評価するのは x_2 .

$\theta^{(i)T} x^{(j)}$ で x_2 の係数は $\begin{pmatrix} 0 \\ \text{小} \\ \square \end{pmatrix}$

$x_2 = 1$ のとき, 値として 5 を出力したいので ④ が正解.

Problem Formulation

$$r(i, j) = \begin{cases} 1 & \text{if user } j \text{ has rated movie } i \\ 0 & \text{otherwise} \end{cases}$$

$y^{(i,j)}$ = ユーザー j が映画 i に与えた rating の値

$\theta^{(j)}$ = ユーザー j についての パラメータ ベクトル

$x^{(i)}$ = 映画 i に対する feature ベクトル

ユーザー j が映画 i について与える rating の予測値 = $(\theta^{(j)})^T x^{(i)}$

$m^{(j)}$ = ユーザー j が rating を与えた映画の数

$\theta^{(j)}$ を学習する. $\theta^{(j)} \in \mathbb{R}^{n+1}$ Regression (回帰) で解ける

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

これは無してよい. 定数倍だから

$\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(n_u)}$ を学習する

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

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

Gradient descent update

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

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

Collaborative Filtering

[講義中の Quiz]

	user 1	2	3	name
Movie 1	0	1.5	2.5	x_1 ? ← 0.5

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

$$x = \begin{pmatrix} 1 \\ x_1 \end{pmatrix}$$

$$y^{(1,3)} = (0 \ 5) \begin{pmatrix} 1 \\ x_1 \end{pmatrix} = 5x_1 = 2.5$$

$$y^{(2,3)} = (0 \ 3) \begin{pmatrix} 1 \\ x_1 \end{pmatrix} = 3x_1 = 1.5$$

$$\therefore x_1 = 0.5 \quad (\text{答})$$

Optimization algorithm

$\theta^{(1)}, \dots, \theta^{(n_u)}$ が与えられた時, $x^{(i)}$ を学習する

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

$\theta^{(1)}, \dots, \theta^{(n_u)}$ が与えられた時, $x^{(1)}, x^{(2)}, \dots, x^{(n_m)}$ を学習する

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

Collaborative filtering

[講義中の Quiz]

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

$i \neq 0$ に対する gradient descent の正しい更新法?

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

i=0 が必要

Collaborative filtering (協調的フィルタリング)

$x^{(1)}, x^{(2)}, \dots, x^{(n_m)}$ \rightarrow $\theta^{(1)}, \dots, \theta^{(n_u)}$ を評価する

$\theta^{(1)}, \dots, \theta^{(n_u)}$ \rightarrow $x^{(1)}, \dots, x^{(n_m)}$ を評価する

Collaborative filtering algorithm

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

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

θ と x を同時に解くことが出来る

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

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

と $J(\dots)$ を定義して

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

もしも常に必要なfeatureがあれば、それを学習して feature n 選んでくれるはず

このとき feature を学習することによってバイアス項 λ はいらない。

$\therefore x \in \mathbb{R}^n$

θ は x と同じ次元であるので $\theta \in \mathbb{R}^n$

$x \in \mathbb{R}^{n+1}, \theta \in \mathbb{R}^{n+1}$ ではない
ことに注意せよ

Collaborative Filtering Algorithm

1. $x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)}$ をランダムに小さい値で初期化する

2. gradient descent (or 他の方法でも可) を用いて,
 $J(x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)})$ を最小化する

for every $j=1, \dots, n_u, i=1, \dots, n_m$

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

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

バイアス項 ($x_0=1$) が存在しないので, 全ての $x_k^{(i)}$ について regularization (これは θ に関しても同様)

3.

[講義中の Quiz]

$x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)}$ を小さい random な値で初期化した理由

- ① この手順は オフショタル (なくてもよい)。0 で初期化しても動く X
- ② この問題についても gradient descent を使う時には random initialization が必要 \Rightarrow
放物線のように local minimum が存在しない問題では初期値は all 0 でもよい X
- ③ 各 i, j について $x^{(i)} \neq \theta^{(j)}$ を保証する X
- ④ symmetry breaking のために必要。 $x^{(1)}, \dots, x^{(n_m)}$ ~~がそれぞれ異なり~~ \Rightarrow それぞれ異なり \Rightarrow 2つを学習する \Rightarrow これを学習する X

Low Rank Matrix Factorization

$$Y = \begin{pmatrix} 5 & 5 & 0 & 0 \\ 5 & ? & ? & 0 \\ ? & 4 & 0 & ? \\ 0 & 0 & 5 & 4 \\ 0 & 0 & 5 & 0 \end{pmatrix} \quad 5 \times 4 \text{ 行列}$$

Predict ratings

$$i \begin{pmatrix} (\theta^{(1)})^T x^{(1)} & (\theta^{(2)})^T x^{(1)} & \dots & (\theta^{(n_u)})^T x^{(1)} \\ \vdots & \vdots & \vdots & \vdots \\ (\theta^{(1)})^T x^{(n_m)} & (\theta^{(2)})^T x^{(n_m)} & \dots & (\theta^{(n_u)})^T x^{(n_m)} \end{pmatrix}$$

↑ $(\theta^{(j)})^T x^{(i)}$

$$X = \begin{pmatrix} -(x^{(1)})^T \\ \vdots \\ -(x^{(n_m)})^T \end{pmatrix}, \quad H = \begin{pmatrix} -(\theta^{(1)})^T \\ \vdots \\ -(\theta^{(n_u)})^T \end{pmatrix}$$

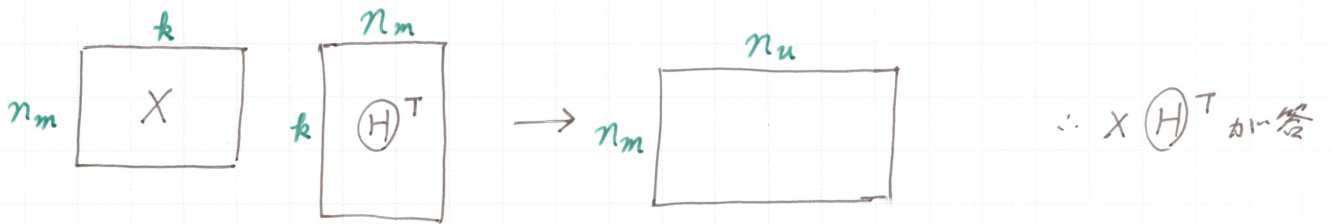
$$X \circ H^T$$

Low Rank Matrix Factorization
低ランク行列分解

[講義中の Quiz]

$$X = \begin{pmatrix} \dots & (x^{(1)})^T & \dots \\ \vdots & \vdots & \vdots \\ \dots & (x^{(n_m)})^T & \dots \end{pmatrix}, \quad \Theta = \begin{pmatrix} \dots & (\theta^{(1)})^T & \dots \\ \vdots & \vdots & \vdots \\ \dots & (\theta^{(n_u)})^T & \dots \end{pmatrix}$$

$$n_m \begin{pmatrix} (x^{(1)})^T \theta^{(1)} & \dots & (x^{(1)})^T \theta^{(n_u)} \\ \vdots & \ddots & \vdots \\ (x^{(n_m)})^T \theta^{(1)} & \dots & (x^{(n_m)})^T \theta^{(n_u)} \end{pmatrix} \quad \text{の別の表記法は?}$$



Finding related movies

movie i について, feature への $x^{(i)} \in \mathbb{R}^k$ を学習した時
 movie i に関連した movie j を見つけるには

$$\text{small } \|x^{(i)} - x^{(j)}\| \rightarrow j \text{ と } i \text{ は 似ている}$$

Implementation detail : Mean Normalization
平均標準化

$n_u = 5$

user \ movie	user 1	2	3	4	5
movie 1	5	5	0	0	
2	5			0	
$n_m = 5$ 3		4	0		
4	0	0	5	4	
5	0	0	5		

$n = 2$ とする. $\theta^{(5)} \in \mathbb{R}^2$

$$\min_{\theta^{(1)}, \dots, \theta^{(5)}} \frac{1}{2} \sum_{(i,j): r(i,j)=1}^{n_u} ((\theta^{(i)})^T x^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \underbrace{\sum_{i=1}^{n_m} \sum_{k=1}^2 (x_k^{(i)})^2}_{(2)} + \underbrace{\sum_{j=1}^5 \sum_{k=1}^2 (\theta_k^{(j)})^2}_{(3)}$$

user 5 は何も rating していないので ① は 0 となる. ② も 0 となる.

③ の $\frac{\lambda}{2} ((\theta_1^{(5)})^2 + (\theta_2^{(5)})^2)$ を min (最小化) しようとする.

その結果 $\theta^{(5)} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$ となる.

$(\theta^{(5)})^T x^{(i)}$ を計算すると、この i についても 0 となる. これは意味が変い.

$$Y = \begin{pmatrix} 5 & 5 & 0 & 0 & ? \\ 5 & ? & ? & 0 & ? \\ ? & 4 & 0 & ? & ? \\ 0 & 0 & 5 & 4 & ? \\ 0 & 0 & 5 & 0 & ? \end{pmatrix}, \quad \begin{matrix} \text{平均を} \\ \text{計算する} \\ \text{(?は除く)} \end{matrix} \mu = \begin{pmatrix} 2.5 \\ 2.5 \\ 2 \\ 2.25 \\ 1.25 \end{pmatrix} \quad \begin{matrix} Y \text{ から} \\ \text{引く} \end{matrix} Y = \begin{pmatrix} 2.5 & 2.5 & -2.5 & -2.5 & ? \\ 2.5 & ? & ? & -2.5 & ? \\ ? & 2 & -2 & ? & ? \\ -2.25 & -2.25 & 2.75 & 1.25 & ? \\ -1.25 & -1.25 & 3.75 & -1.25 & ? \end{pmatrix}$$

$i = j$, movie j に関しては

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

↓
 $(H)^{(i)}, x^{(i)}$ を学習

$i = 5$ について $\theta^{(5)} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$ だったので $\underbrace{(\theta^{(5)})^T x^{(i)}}_{=0} + \mu_i = \mu_i$

Week 9 21

mean normalisation (平均標準化) について話した。
[講義中の Quiz]

feature scaling の他のアプローチとは異なり, movie の rating では
(max - min) で割るというスケールングを行わなかった。

- ① predict されるのは 実際の値なので、この種の scaling は役に立たない。 X
- ② 全ての rating は 比較可能である (0~5)。したがって 既に同じスケールである。 O
- ③ 平均を引くことは 範囲で割ると 数学的に同じだ。 X
- ④ これにより 全部のアルゴリズムは 非常に効果的になる。 X