

Clustering Algorithm

K-means algorithm

Randomly initialize K cluster centroids

Repeat {

Assign points to cluster centroids

for $i = 1$ to m

$c^{(i)} :=$ index (from 1 to K) of cluster centroid closest to $x^{(i)}$

Move cluster centroids

for $k = 1$ to K

}

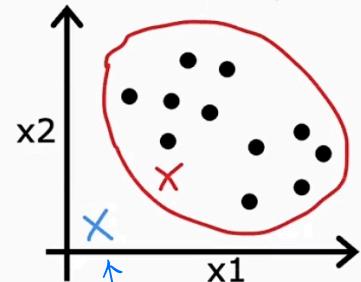
$\mu_k :=$ average (mean) of points assigned to cluster k

$$\begin{array}{l} \mu_1, \mu_2 \\ \vdots \\ \mu_1, \mu_2, \dots, \mu_K \end{array}$$

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

$n=2$ $K=2$ ~~K=2~~

$K=K-1$ ✓



① Randomly initialize centroid

DeepLearning.AI Stanford ONLINE

Andrew Ng

但在运行 k-means 时，如果没有为某簇分配点，消除该簇的做法实际上更常见

K-means optimization objective

$c^{(i)}$ = index of cluster ($1, 2, \dots, K$) to which example $x^{(i)}$ is currently assigned

μ_k = cluster centroid k

$\mu_{c^{(i)}}$ = cluster centroid of cluster to which example $x^{(i)}$ has been assigned

Cost function

$$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2$$

$\min_{c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K} J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$

Distortion

DeepLearning.AI Stanford ONLINE

Andrew Ng

Random initialization

防止陷入局部极小
因初始质心选择不当而

For $i = 1$ to 100 {

50-1000

Randomly initialize K-means.

k random examples

Run K-means. Get $c^{(1)}, \dots, c^{(m)}, \mu_1, \mu_1, \dots, \mu_k \leftarrow$

Computer cost function (distortion)

$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \mu_1, \dots, \mu_k) \leftarrow$

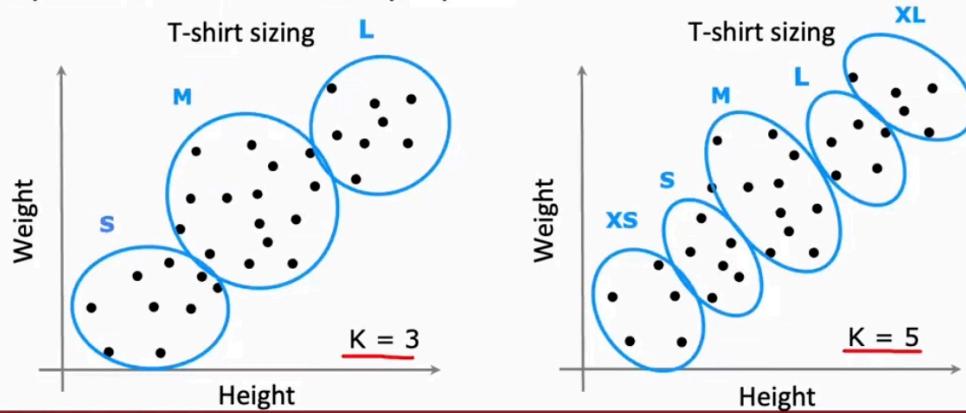
}

Pick set of clusters that gave lowest cost J

Choosing the value of K

Sometimes, you're running K-means to get clusters to use for some later/downstream purpose. Evaluate K-means based on a metric for how well it performs for that later purpose.

E.g.



异常检测

Anomaly detection example

Aircraft engine features:

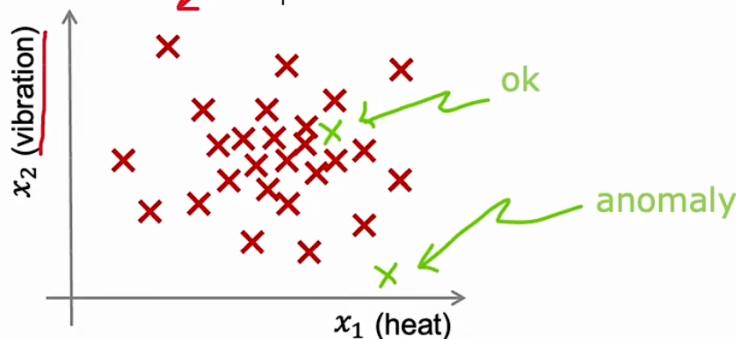
x_1 = heat generated

x_2 = vibration intensity

...

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

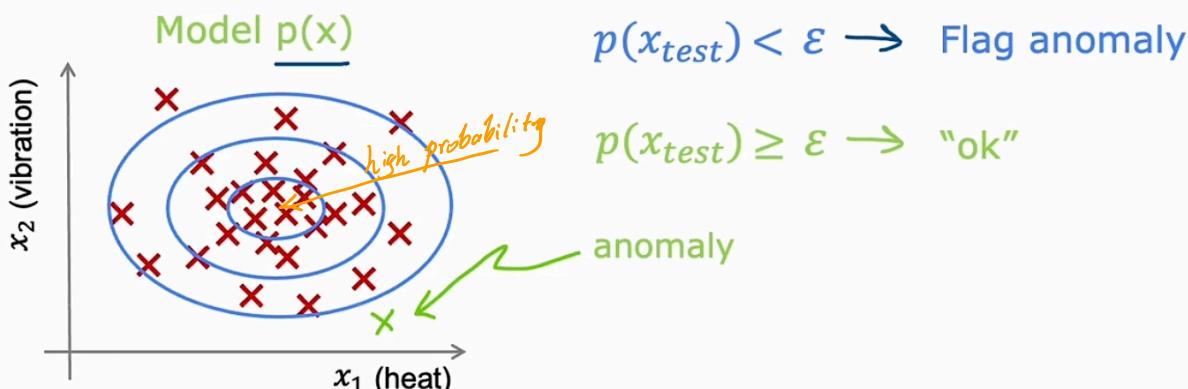
New engine: x_{test}



Density estimation

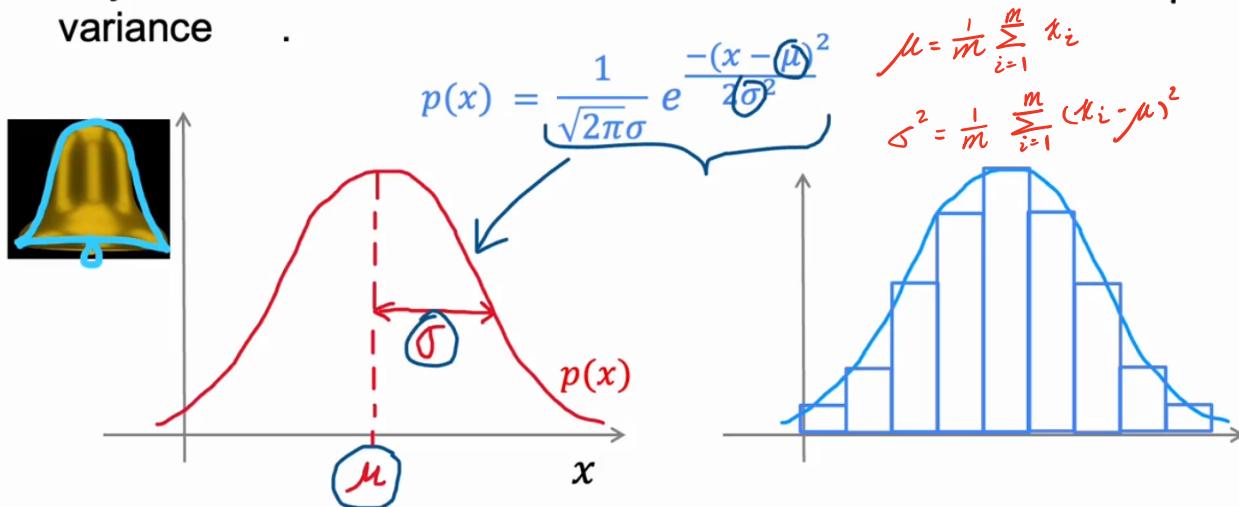
Dataset: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

Is x_{test} anomalous?



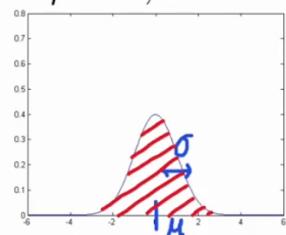
Gaussian (Normal) distribution

Say x is a number. If x is a distributed Gaussian with mean μ , σ^2 variance

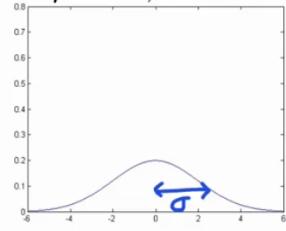


Gaussian distribution example

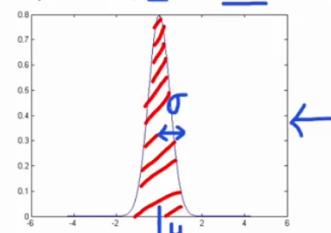
$$\mu = 0, \sigma = 1$$



$$\mu = 0, \sigma = 2 \quad \sigma^2 = 4$$

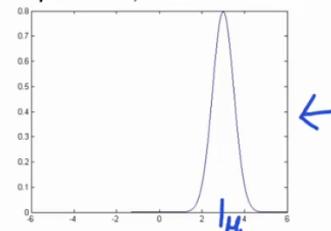


$$\mu = 0, \sigma = 0.5$$



$$\sigma^2 = 0.25$$

$$\mu = 3, \sigma = 0.5$$



Density estimation

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

Each example x_i has n features

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

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

*极大似然估计
mostly likelihood*

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

$$\sum \rightarrow + \quad \prod \rightarrow \times$$

$p(x_1 = \text{high temp}) = 1/10$
 $p(x_2 = \text{high vibra}) = 1/20$
 $p(x_1, x_2) = p(x_1) * p(x_2)$
 $= \frac{1}{10} * \frac{1}{20} = \frac{1}{200}$

Anomaly detection algorithm

1. Choose n features x_i that you think might be indicative of anomalous examples.
2. Fit parameters $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$

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

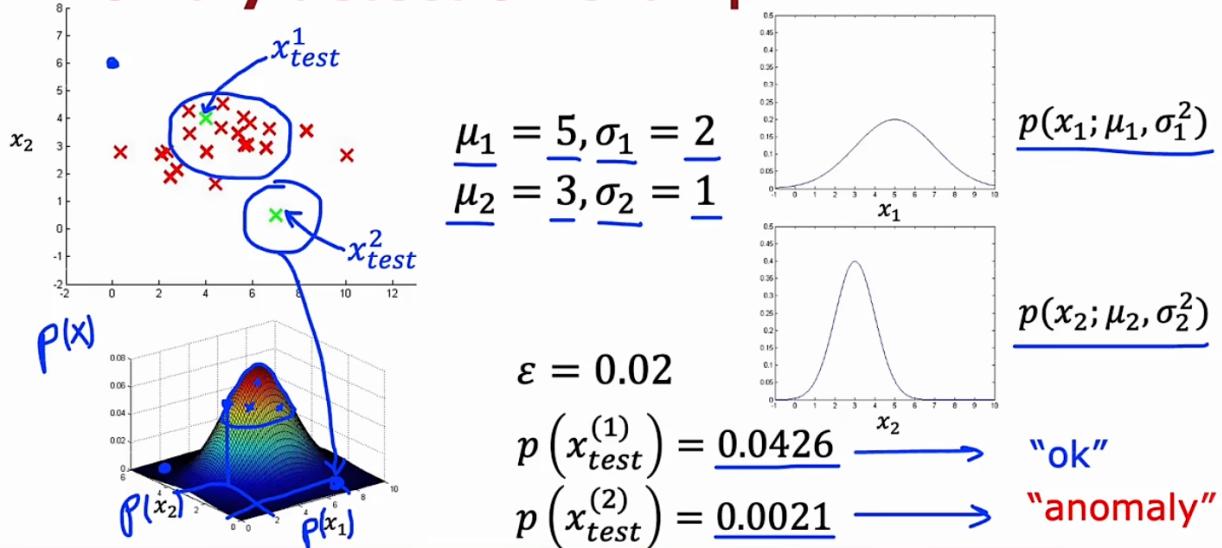
$$\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$

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 $p(x) < \varepsilon$

Anomaly detection example



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

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

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

Include a few anomalous examples
 $y=1$

Aircraft engines monitoring example

10000 good (normal) engines ↗ 2-50
→ 2000 flawed engines (anomalous) $y=1$

(Training set: 6000 good engines ↗ $y=0$)
(CV: 2000 good engines ($y=0$), 10 anomalous ($y=1$)
(Test: 2000 good engines ($y=0$), 10 anomalous ($y=1$) ←

ε x_j

Alternative:

→ Training set: 6000 good engines 2
→ CV: 4000 good engines ($y=0$), 20 anomalous ($y=1$) ←
No test set ←

Algorithm evaluation

→ Fit model $p(x)$ on training set $x^{(1)}, x^{(2)}, \dots, x^{(m)}$
On a cross validation/test example x , predict

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

10
2000

Possible evaluation metrics:

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

and that these are alternative metrics to classification accuracy that could work better when your data distribution is very skewed.

并且这些是分类准确性的替代指标，当您的数据分布非常倾斜时，它们可以更好地工作

Anomaly detection

Very small number of positive examples ($y = 1$). (0-20 is common).

Large number of negative examples ($y = 0$).

$p(x)$

$y=1$

Many different “types” of anomalies. Hard for any algorithm to learn from positive examples what the anomalies look like; future anomalies may look nothing like any of the anomalous examples we've seen so far.

Fraud

vs. Supervised learning

Large number of positive and negative examples.

20 positive examples

Enough positive examples for algorithm to get a sense of what positive examples are like, future positive examples likely to be similar to ones in training set.

Spam

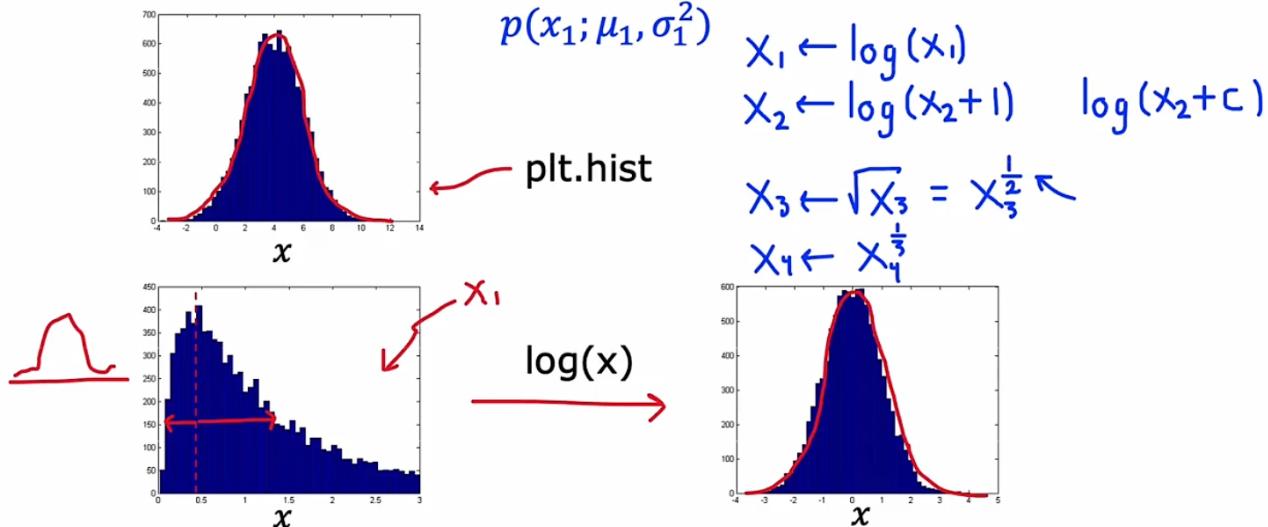
Anomaly detection

- Fraud detection
 - Manufacturing - Finding new previously unseen defects in manufacturing.(e.g. aircraft engines)
 - Monitoring machines in a data center
- ⋮

vs. Supervised learning

- Email spam classification
 - Manufacturing - Finding known, previously seen defects $y=1$ scratches
 - Weather prediction (sunny/rainy/etc.)
 - Diseases classification
- ⋮

Non-gaussian features



DeepLearning.AI Stanford ONLINE

Andrew Ng

我有时会观察我的特征，如果画完直方图发现某个特征非常的不高斯，我可能会选择用这样或其他的转换让它变得更高斯

I'll sometimes take a look at my features and if I see any of the highly non Gaussian by plotting a histogram, I might choose transformations like these or others in order to try to make it more Gaussian.

提醒一下，无论你对训练集使用什么转换

And just as a reminder, whatever transformations you apply to the training set

请记住对交叉验证和测试集数据也进行相同的转换

Please remember to apply the same transformation to your cross validation and test set data as well.

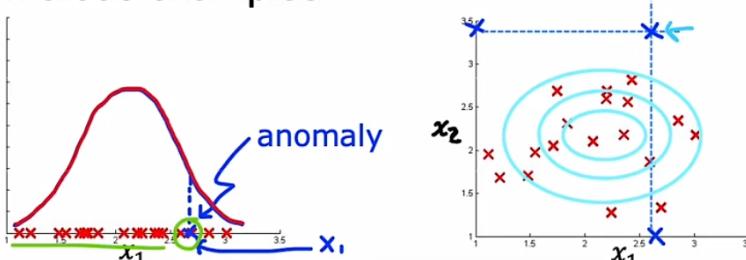
Error analysis for anomaly detection

Want $p(x) \geq \epsilon$ large for normal examples x .

$p(x) < \epsilon$ small for anomalous examples x .

Most common problem:

$p(x)$ is comparable (say, both large) for normal and anomalous examples



先训练模型，然后查看算法在交叉验证集中没能检测出的异常，然后查看这些样本，考虑是否有必要新建一个特征，从而使算法能够发现这些异常

Andrew Ng

What if we have features of the movies?

Movie	Alice(1)	Bob(2)	Carol(3)	Dave(4)	x_1 (romance)	x_2 (action)	$n_u = 4$	$n_m = 5$	$n = 2$
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			
Nonstop car chases	0	0	5	4	0.1	1.0			
Swords vs. karate	0	0	5	?	0	0.9			

For user 1: Predict rating for movie i as: $w^{(1)} \cdot x^{(i)} + b^{(1)}$ ← just linear regression

$$w^{(1)} = \begin{bmatrix} 5 \\ 0 \end{bmatrix} \quad b^{(1)} = 0 \quad x^{(3)} = \begin{bmatrix} 0.9 \\ 0 \end{bmatrix} \quad w^{(1)} \cdot x^{(3)} + b^{(1)} = 4.95$$

→ For user j : Predict user j 's rating for movie i as $w^{(j)} \cdot x^{(i)} + b^{(j)}$

所以这里的参数 w_j 和 b_j 是用于预测用户 j 给电影 i 的评分的参数，它是 x_i 的函数，是电影 i 的一个特征

Cost function

Notation:

- $r(i,j) = 1$ if user j has rated movie i (0 otherwise)
- $y^{(i,j)}$ = rating given by user j on movie i (if defined)
- $w^{(j)}, b^{(j)}$ = parameters for user j
- $x^{(i)}$ = feature vector for movie i

For user j and movie i , predict rating: $w^{(j)} \cdot x^{(i)} + b^{(j)}$

- $m^{(j)}$ = no. of movies rated by user j
- To learn $w^{(j)}, b^{(j)}$

$$\min_{w^{(j)}, b^{(j)}} J(w^{(j)}, b^{(j)}) = \frac{1}{2m^{(j)}} \sum_{i:r(i,j)=1} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2m^{(j)}} \sum_{k=1}^n (w_k^{(j)})^2$$

number of features
Regularization term

Cost function

① 电影特征 \rightarrow 用户评价

To learn parameters $w^{(j)}, b^{(j)}$ for user j :

$$J(w^{(j)}, b^{(j)}) = \frac{1}{2} \sum_{i:r(i,j)=1} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{k=1}^n (w_k^{(j)})^2$$

To learn parameters $w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}, \dots, w^{(n_u)}, b^{(n_u)}$ for all users :

$$J(w^{(1)}, b^{(1)}, \dots, w^{(n_u)}, b^{(n_u)}) = \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i:r(i,j)=1} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (w_k^{(j)})^2$$

Cost function

② 用户评价 \rightarrow 电影特征

Given $w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}, \dots, w^{(n_u)}, b^{(n_u)}$

to learn $x^{(i)}$:

$$J(x^{(i)}) = \frac{1}{2} \sum_{j:r(i,j)=1} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{k=1}^n (x_k^{(i)})^2$$

→ To learn $x^{(1)}, x^{(2)}, \dots, x^{(n_m)}$:

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

Collaborative filtering

Cost function to learn $w^{(1)}, b^{(1)}, \dots, w^{(n_u)}, b^{(n_u)}$:

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

Cost function to learn $x^{(1)}, \dots, x^{(n_m)}$:

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

Put them together:

$$\min_{\substack{w^{(1)}, \dots, w^{(n_u)} \\ b^{(1)}, \dots, b^{(n_u)} \\ x^{(1)}, \dots, x^{(n_m)}}} J(w, b, x) = \frac{1}{2} \sum_{(i,j):r(i,j)=1} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (w_k^{(j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2$$

Gradient Descent

collaborative filtering

Linear regression (course 1)

repeat {

$$\underline{w_i = w_i - \alpha \frac{\partial}{\partial w_i} J(w, b)}$$

$$\underline{b = b - \alpha \frac{\partial}{\partial b} J(w, b)}$$

$$w_i^{(j)} = w_i^{(j)} - \alpha \frac{\partial}{\partial w_i^{(j)}} J(w, b, x)$$

$$b^{(j)} = b^{(j)} - \alpha \frac{\partial}{\partial b^{(j)}} J(w, b, x)$$

$$x_k^{(i)} = x_k^{(i)} - \alpha \frac{\partial}{\partial x_k^{(i)}} J(w, b, x)$$

}

parameters w, b, x

\times is also a parameter

Cost function for binary application

Previous cost function:

$$\frac{1}{2} \sum_{(i,j):r(i,j)=1} \underbrace{(w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2}_{f(x)} + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (w_k^{(j)})^2$$

Loss for binary labels $y^{(i,j)} : f_{(w,b,x)}(x) = g(w^{(j)} \cdot x^{(i)} + b^{(j)})$

$$L(f_{(w,b,x)}(x), y^{(i,j)}) = -y^{(i,j)} \log(f_{(w,b,x)}(x)) - (1 - y^{(i,j)}) \log(1 - f_{(w,b,x)}(x)) \quad \leftarrow \text{Loss for single example}$$

$$J(w, b, x) = \sum_{(i,j):r(i,j)=1} L(f_{(w,b,x)}(x), y^{(i,j)})$$

$$g(w^{(j)} \cdot x^{(i)} + b^{(j)})$$

然后，这就给你了可用于二元标签协同过滤的代价函数

Mean Normalization

→	5	5	0	0	?	2.5
→	5	?	?	0	?	2.5
?	4	0	?	?	?	2
0	0	5	4	?	?	2.25
0	0	5	0	?	?	1.25

$$\mu = \begin{bmatrix} 2.5 \\ 2.5 \\ 2 \\ 2.25 \\ 1.25 \end{bmatrix}$$

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	?

For user j , on movie i predict:

$$w^{(j)} \cdot x^{(i)} + b^{(j)} + \mu_i$$

$$y^{(i,j)} = w^{(j)} \cdot b^{(j)} \cdot x^{(i)}$$

User 5 (Eve):

$$w^{(5)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad b^{(5)} = 0$$

$$w^{(5)} \cdot x^{(1)} + b^{(5)} + \mu_1 = 2.5$$

将所有用户标为均值

Collaborative filtering vs Content-based filtering

→ Collaborative filtering:

Recommend items to you based on rating of users who gave similar ratings as you

→ Content-based filtering:

Recommend items to you based on features of user and item to find good match

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

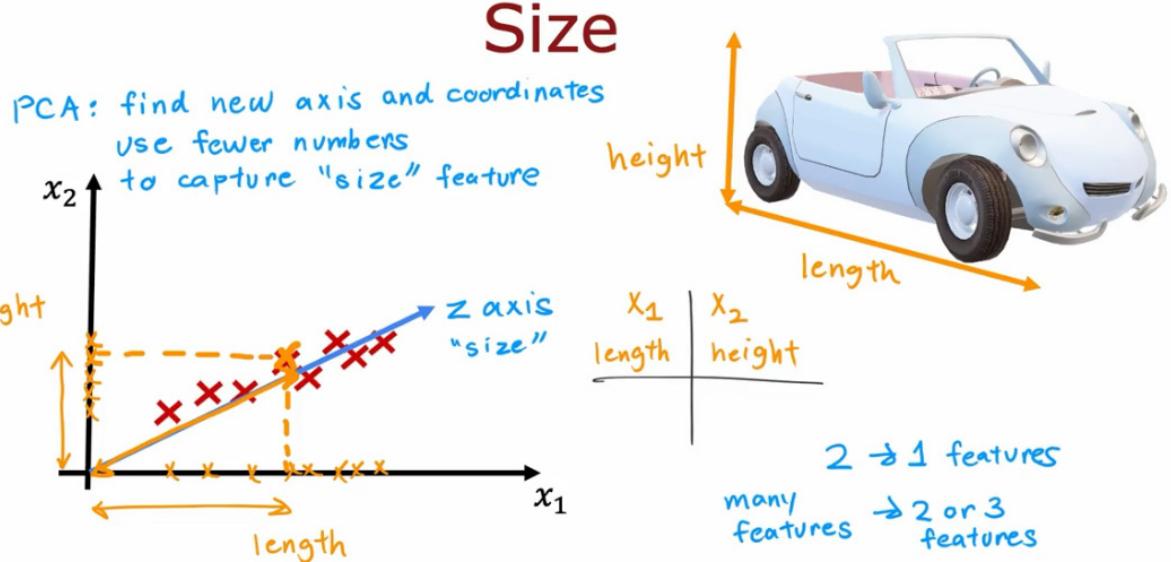
$y^{(i,j)}$ rating given by user j on item i (if defined)

to find better matches than potentially a pure collaborative filtering

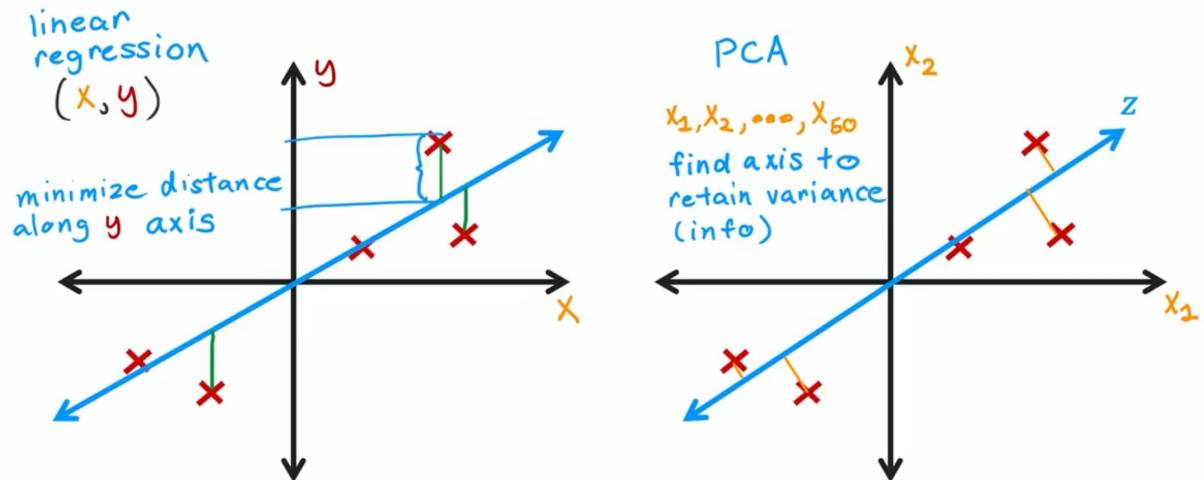
Stanford ONLINE approach might be able to.

Andrew Ng

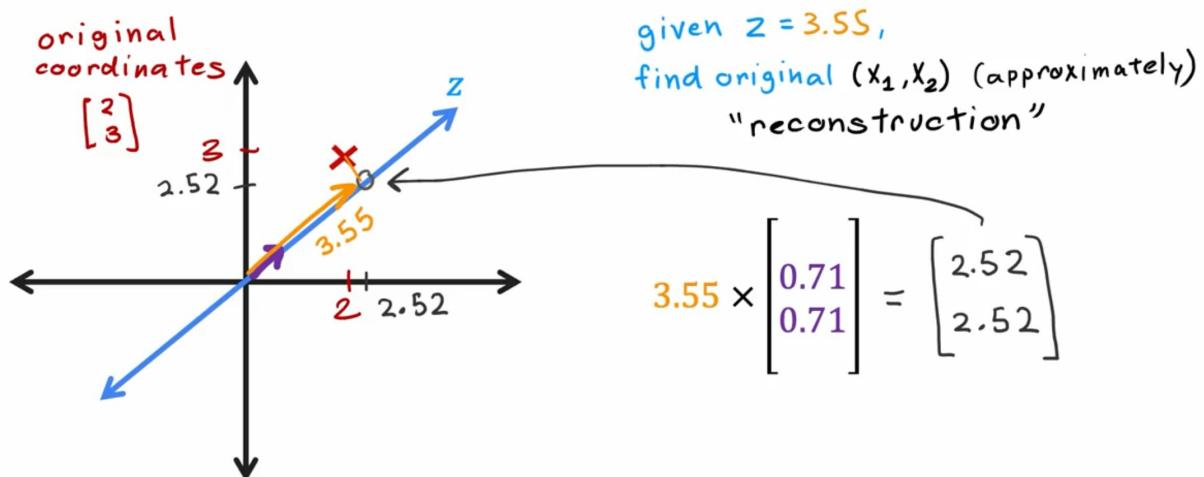
PCA



PCA is not linear regression



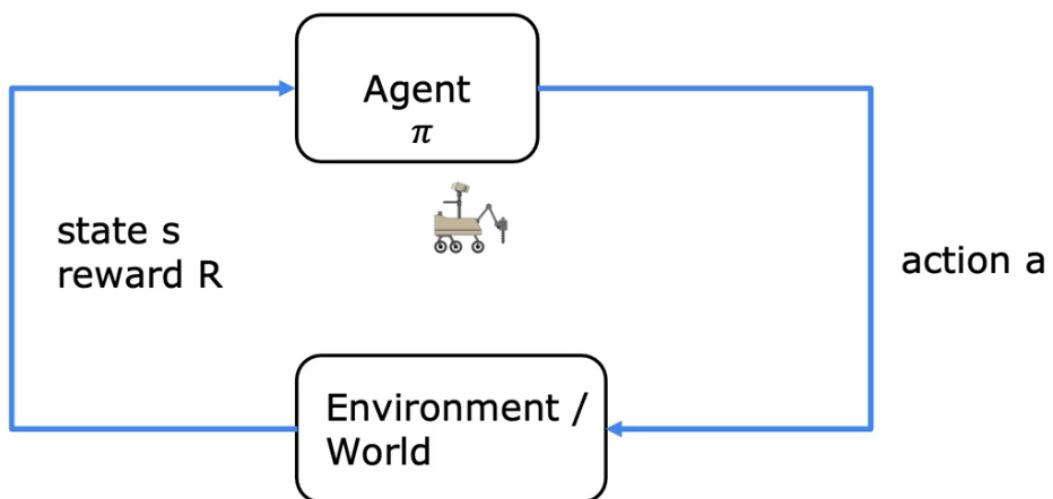
Approximation to the original data



Reinforcement Learning

	Mars rover	Helicopter	Chess
states	6 states	position of helicopter	pieces on board
actions	↔ ↔	how to move control stick	possible move
rewards	100, 0, 40	+1, -1000	+1, 0, -1
discount factor γ	0.5	0.99	0.995
return	$R_1 + \gamma R_2 + \gamma^2 R_3 + \dots$	$R_1 + \gamma R_2 + \gamma^2 R_3 + \dots$	$R_1 + \gamma R_2 + \gamma^2 R_3 + \dots$
policy π	100 ← ← ← → 40	Find $\pi(s) = a$	Find $\pi(s) = a$

Markov Decision Process (MDP)



State action value function

$Q(s, a) =$ Return if you
 • start in state s . Given Policy
 • take action a (once).
 • then behave optimally after that.

$Q(s, a)$

100	50	25	12.5	20	40
100	0	0	0	0	40
100	0	0	0	0	40

$Q(2, \leftarrow) Q(2, \rightarrow)$

100	50	25	12.5	20	40
100	0	0	0	0	40
100	0	0	0	0	40

1 2 3 4 5 6

← return
 ← action
 ← reward

$$Q(2, \rightarrow) = 12.5 \\ 0 + (0.5)0 + (0.5)^2 0 + (0.5)^3 100$$

$$Q(2, \leftarrow) = 50 \\ 0 + (0.5)100$$

$$Q(4, \leftarrow) = 12.5 \\ 0 + (0.5)0 + (0.5)^2 0 + (0.5)^3 100$$

Bellman Equation

$Q(s, a) =$ Return if you
 • start in state s .
 • take action a (once).
 • then behave optimally after that.



s : current state
 a : current action

$R(s)$ = reward of current state

s' : state you get to after taking action a
 a' : action that you take in state s'

$$Q(s, a) = R(s) + \gamma \max_{a'} Q(s', a')$$

Expected Return

Goal of Reinforcement Learning:

Choose a policy $\pi(s) = a$ that will tell us what action a to take in state s so as to maximize the expected return.

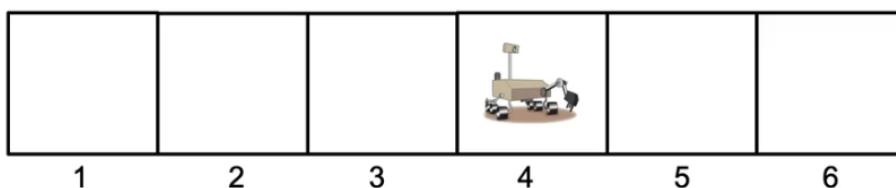
Bellman
Equation:

$$Q(s, a) = R(s) + \gamma E[\max_{a'} Q(s', a')]$$

↑ ↑
3 ←
 ↑
 2 or 4

Discrete vs Continuous State

Discrete State:

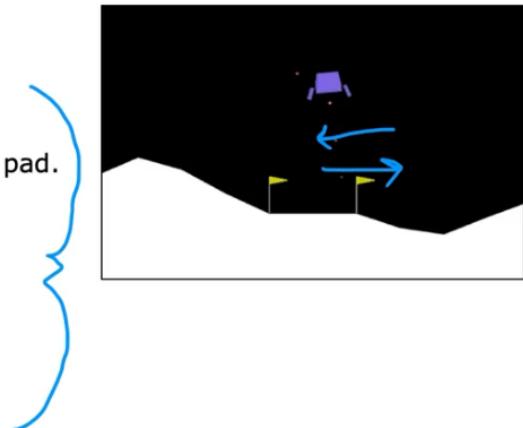


Continuous State:



Reward Function

- Getting to landing pad: 100 – 140
- Additional reward for moving toward/away from pad.
- Crash: -100
- Soft landing: +100
- Leg grounded: +10
- Fire main engine: -0.3
- Fire side thruster: -0.03



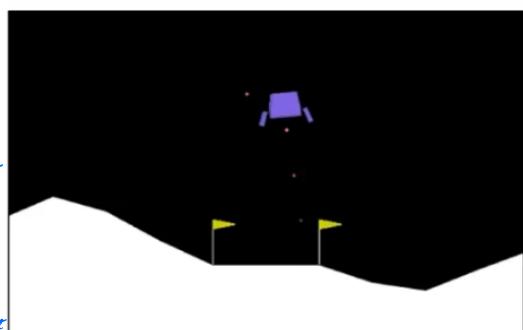
Lunar Lander Problem

Learn a policy π that, given

$$s = \begin{bmatrix} x \\ y \\ \dot{x} \\ \dot{y} \\ \theta \\ \dot{\theta} \\ l \\ r \end{bmatrix}$$

with annotations:

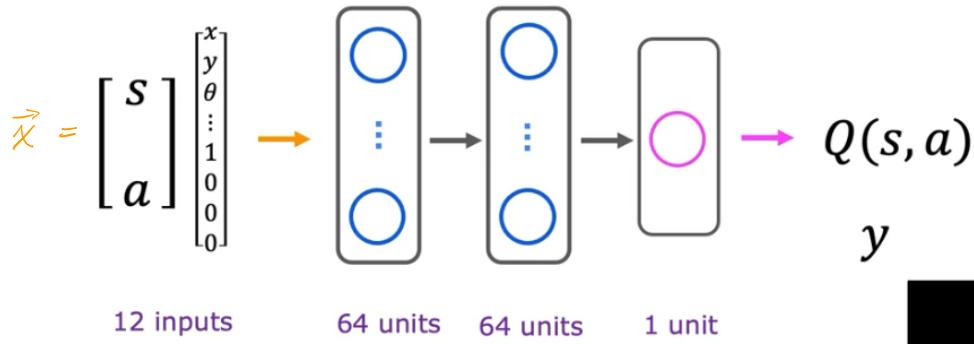
- x, y : Position
- \dot{x}, \dot{y} : Speed in direction
- θ : angle
- $\dot{\theta}$: angle changes speed
- l : whether left/right leg has touched ground



picks action $a = \pi(s)$ so as to maximize the return.

$$Y = 0.985$$

Deep Reinforcement Learning



In a state s , use neural network to compute

$Q(s, \text{nothing}), Q(s, \text{left}), Q(s, \text{main}), Q(s, \text{right})$

Pick the action a that maximizes $Q(s, a)$



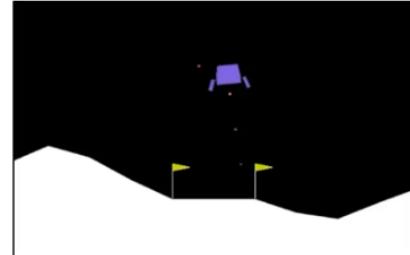
Learning Algorithm

Initialize neural network randomly as guess of $Q(s, a)$.

Repeat {

Take actions in the lunar lander. Get $(s, a, R(s), s')$.

Store 10,000 most recent $(s, a, R(s), s')$ tuples.



Replay Buffer

Train neural network:

Create training set of 10,000 examples using

$$x = (s, a) \text{ and } y = R(s) + \gamma \max_{a'} Q(s', a')$$

Train Q_{new} such that $Q_{\text{new}}(s, a) \approx y$.

Set $Q = Q_{\text{new}}$.

$$f_{w, b}(x) \approx y$$

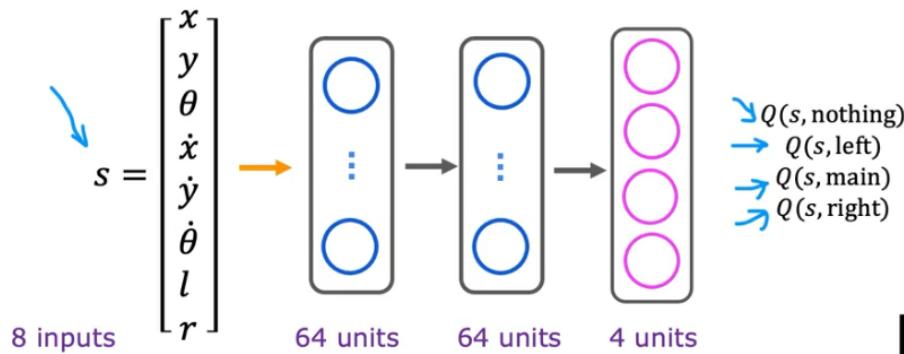
x, y

$x^{(1)}, y^{(1)}$

:

x^{10000}, y^{10000}

Deep Reinforcement Learning



In a state s , input s to neural network.

Pick the action a that maximizes $\underline{Q(s, a)}$. $R(s) + \gamma \max_{a'} Q(s', a')$



How to choose actions while still learning?

In some state s

Option 1:

Pick the action a that maximizes $\underline{Q(s, a)}$.

Option 2:

- With probability 0.95, pick the action a that maximizes $\underline{Q(s, a)}$. Greedy, "Exploitation"
- With probability 0.05, pick an action a randomly. "Exploration"

ε -greedy policy ($\varepsilon = 0.05$)
0.95

why exploration:

if an action always takes low Q , the model may never choose this action, but the action can be effective sometimes.

$Q(s, \text{main})$ is low



Start ε high
1.0 → 0.01
Gradually decrease

Limitations of Reinforcement Learning

- Much easier to get to work in a simulation than a real robot!
- Far fewer applications than supervised and unsupervised learning.
- But ... exciting research direction with potential for future applications.