



课程名称：

机器学习算法与应用

讲授人：唐晨，徐岩

电气自动化与信息工程学院



教材或参考书目

1. E. Alpaydin, **Introduction to Machine Learning**, The MIT Press, 2004
2. M. Bishop, **Pattern Recognition and Machine Learning**, Springer, 2006
3. Witten, Ian H., et al. **Data Mining: Practical machine learning tools and techniques**, Morgan Kaufmann, 2016
4. E. Alpaydin著, 范明译, **机器学习导论第三版**, 机械工业出版社, 2016



Professor Chen Tang

E-mail: tangchen@tju.edu.cn

Office: 26D-316

Associate Professor Yan Xu

E-mail: xuyan@tju.edu.cn

Office: 26D-426

MachineLearningTJU@126.com





课程设置背景



1.时代发展的需要

机器学习是新思想、新观念、新理论、新技术不断出现并迅速发展的新兴学科，是人工智能的基础。

人工智能改变未来世界。党的十九大报告指出，要推动互联网、大数据、人工智能和实体经济深度融合。天津市部署了《加快推进智能科技产业发展总体行动计划》，明确大力发展战略性新兴产业。天津大学刚刚成立了人工智能学院和智能与计算学部。

人才的质量和数量决定着人工智能发展的水平和潜力。《新一代人工智能发展规划》指出，我国人工智能尖端人才远远不能满足需求。发展人工智能，核心在于人。

开设人工智能基础课程机器学习算法和应用是十分必要，是时代发展的需要。



课程设置背景

2. 研究生培养的需要

目前几乎所有学科（工科、理科、金融和文法等等）研究生的研究方向与机器学习相关，然而现在的绝大多数研究生在本科阶段没有学过机器学习这门课程。

这门课程适合研究生公共选修课，是研究生培养的需要。





课程设置背景



近5年（2013-2017）万方以“机器学习”为关键词

万方数据 知识服务平台 V2.0 返回旧版 欢迎天津大学的朋友 | 个人登录 | 注册 | 钱包

全部 期刊 学位 会议 专利 科技报告 成果 标准 法规 地方志 视频 ▶

万方智搜 机器学习 检索

找到 5538 条结果。 重置

限定条件： 近五年 × 重置

标题 作者 关键词 学科专业 学校 导师

起始年 - 结束年 结果中检索

排序： 相关度 范围： 显示全部 显示20条

全选 清除 导出 收藏

国内外文献保障

研究趋势

外文 中文

1,500
1,250
1,000
750
500
250
0
2013 2014 2015 2016

相关热词

深度学习 机器 图像处理 特征提取 数据挖掘

1. [硕士学位论文] 基于机器学习的DNA序列分类算法研究
董云玲 计算机科学与技术；计算机应用技术 东北农业大学 2017 (学位年度)
摘要：随着人类基因组计划的完成以及测序技术的发展，产生了海量的生物数据。如何从海量的生物数据中挖掘有用的信息是摆在生命信息研究者们面前的一道难题，在这样的背景下，生物信息学应运而生。目前，研究DNA序列分类问题已经成为生物信...
关键词：生物数据库 数据挖掘 智能分类 程序语言
在线阅读 下载 导出 收藏 分享

2. [硕士学位论文] 基于机器学习的乳腺肿瘤识别
霍双红 数学 中北大学 2017 (学位年度)
摘要：乳腺肿瘤严重危害到女性的健康，目前为止还没有找到很好的预测乳腺癌的方法。目前，依照当前的医疗水平，唯一提高乳腺癌的治愈率和降低乳腺癌的死亡率的方法关键在于要提早发现，早发现可以及时治疗，早诊断，不要耽误最佳治疗时间和早治...
关键词：乳腺肿瘤诊断 神经网络 遗传算法 机器学习
在线阅读 下载 导出 收藏 分享

更多...

授予学位

硕士 (4621) 博士 (916)

在线阅读 下载 导出 收藏 分享





课程设置背景



学科分类	学位论文数
工业技术	4624
经济	267
医疗、卫生	262
交通运输	179
文化科学	160
数理科学和化学	148
生物科学	130
天文学、地球	83
农业科学	74
语言、文学	39



课程设置背景



(2013-2017) “Machine Learning” ProQuest

The screenshot shows the ProQuest search interface. The search term "Machine Learning" is highlighted in a blue box. Below it, there are two filter options: "全文文献" (Full Text) and "同行评审" (Peer-reviewed). The search results count is 50,077, with a "搜索范围" (Search Range) link. The sorting order is set to "相关性" (Relevance). The sidebar on the left includes sections for narrowing results, applying filters, and selecting document types. Under "应用的筛选器" (Applied Filters), "学位论文" (Theses) and "2013 - 2017" are selected. Under "全文文献" (Full Text), "学位论文 (50,077)" is selected. Under "出版物类型" (Publication Type), "学位论文" is also selected. Under "出版日期" (Publication Date), the range "2013 - 2017 (年)" is selected.



3.机器学习的主流算法在科学研究中广泛应用

IEEE Xplore®
Digital Library

Access provided by:
TIANJIN UNIVERSITY
» Sign Out

IEEE

Browse ▾ My Settings ▾ Get Help ▾

All SVM

Advanced Search | Other Search Options ▾

Search within result Show: All Results ▾ | Per Page: 25 ▾ | Download PDFs ▾ | Export ▾ | Set Search Alerts ▾ | Search History

Displaying results 1-25 of 11,933 for SVM ×

▼ Filters Applied: 2013-2018 ×

Conferences (10,341) Journals & Magazines (1,412) Early Access Articles (155)

Books (25)

Year

Single Year Range

From 2013 To 2018

Select All on Page Sort By: Relevance ▾

Research on Web text classification algorithm based on improved CNN and SVM
Zhiqian Wang; Zhiyi Qu
2017 IEEE 17th International Conference on Communication Technology (ICCT)
Year: 2017
Pages: 1958 – 1961
IEEE Conferences

Abstract PDF (381 Kb) ©

Featured Book

Data Mining
Concepts, Models, Methods, and Algorithms

MyXplore™ Mobile App



IEEE Xplore®
Digital Library

Access provided by:
TIANJIN UNIVERSITY
[» Sign Out](#)

IEEE

Browse ▾

My Settings ▾

Get Help ▾

All

PCA



[Advanced Search](#) | [Other Search Options ▾](#)

Search within results **Show: All Results ▾** | Per Page: 25 ▾ | Download PDFs ▾ | Export ▾ | Set Search Alerts ▾ | Search History

Displaying results 1-25 of 4,539 for **PCA**

▼ Filters Applied: 2013-2018

Conferences (3,906)

Journals & Magazines (567)

Early Access Articles (60)

Books (5)

Standards (1)

Year

Select All on Page

Sort By: **Relevance** ▾

A Comparison between ECG Beat Classifiers Using Multiclass SVM and SIMCA with Time Domain PCA Feature Reduction

Najlaa Jannah; Sillas Hadjiloucas

2017 UKSim-AMSS 19th International Conference on Computer Modelling & Simulation (UKSim)

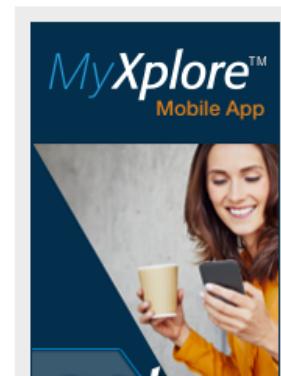
Year: 2017

Pages: 126 – 131

IEEE Conferences

► Abstract

(281 Kb)





课程设置背景



4.新工科人才培养的需要

新工科是基于国家战略发展新需求、国际竞争新形势、立德树人新要求而提出的我国工程教育改革方向。

“新工科”对应的是新兴产业，首先是指针对新兴产业的专业，如人工智能、智能制造、机器人、云计算等，也包括传统工科专业的升级改造。

统计数据显示，到2020年，新一代信息技术产业、电力装备、高档数控机床和机器人、新材料将成为人才缺口最大的几个专业，其中新一代信息技术产业人才缺口将会达到750万人。到2025年，新一代信息技术产业人才缺口将达到950万人，电力装备的人才缺口也将达到909万人。

这门课程对新工科人才的培养是十分重要的。



主要教学内容



背景

预备知识的回顾

机器学习主流算法

算法的实现

算法的应用





课程特色

(1) 注重基础理论的讲授，提高研究生基础理论的水平，为深度学习打好基础

注重

算法的意义
如何公式化
如何数值求解





(2) 注重算法的实现

标准PCA聚类：

```
%% standard PCA
disp('Performing standard PCA...');

Y1=PCA(data,d);
figure;hold on;
plot(Y1(1:2:end,1),Y1(1:2:end,2), 'b*');
plot(Y1(2:2:end,1),Y1(2:2:end,2), 'ro');
legend('class 1','class 2');
title('standard PCA');
drawnow;
```



eig

Eigenvalues and eigenvectors

```
function [Y, eigVector, eigValue]=PCA(X, d)

    %% eigenvalue analysis
    Sx=cov(X);
    [V, D]=eig(Sx);
    eigValue=diag(D);
    [eigValue, IX]=sort(eigValue, 'descend');
    eigVector=V(:, IX);

    %% normalization
    norm_eigVector=sqrt(sum(eigVector.^2));
    eigVector=eigVector./repmat(norm_eigVector, size(eigVector, 1), 1);

    %% dimensionality reduction
    eigVector=eigVector(:, 1:d);
    Y=X*eigVector;
```



课程特色

(3) 注重应用的介绍

中国人口数和大学生毕业**5**年后工资的线性回归

PCA算法在某农业生态经济系统的应用

PCA算法在男性身材分析中的应用

KPCA在人脸识别中的应用

KPCA在数据可视化中的应用

LDA算法在图像处理中的应用

KLDA在数据分类中的应用

SVM的在遥感图像识别中的应用





课程特色



(4) 注重机器学习主流算法的升级

强调基于**Kernel**的算法，并及时更新算法。

机器学习先进算法和应用





课程特色



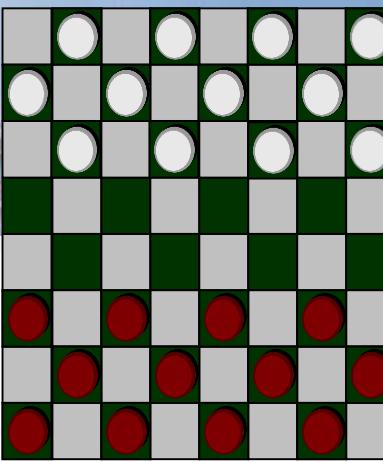
(5) 考试方式强调基础理论和综合应用与分析能力方面的考核

1.大作业（40%）

2.完成与课程相关的研究项目：60%

科技报告（文献综述、方法描述、实验测试、量化评估、结果讨论和总结）

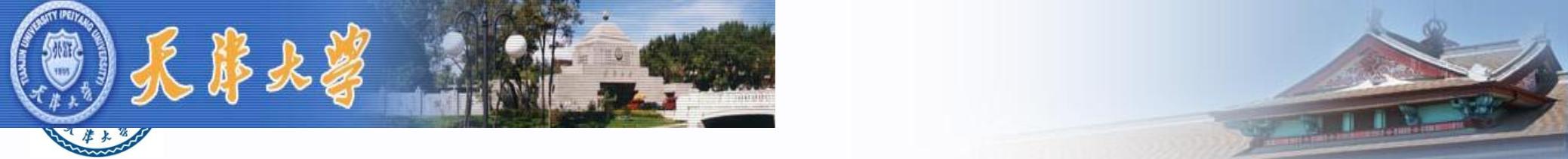
编写程序



Machine Learning definition

Arthur Samuel (1959). Machine Learning: Field of study that gives computers the ability to learn without being explicitly programmed.

机器学习是研究计算机怎样模拟或实现人类的学习行为，以获取新的知识或技能，重新组织已有的知识结构使之不断改善自身的性能。它是人工智能的核心，是使计算机具有智能的根本途径，其应用遍及人工智能的各个领域。



天津大学

Machine Learning -New capability for computers

Examples:

- Database mining

- Large datasets from growth of web.

- E.g., Web click data, medical records, biology, engineering

- Applications can't program by hand.

- E.g., Autonomous helicopter, handwriting recognition, most of Natural Language Processing (NLP), Computer Vision.

- Medical, Image classification, The prediction of an event.



天津大学



Emotion Classification on Face Images



Figure 1: Examples of face images. From left to right: Disgust (label 3), Happy (label 5) and Sadness (label 6).

SIFT (Scale-Invariant Feature Transform)



Figure 2: Examples of SIFT descriptors of our images.

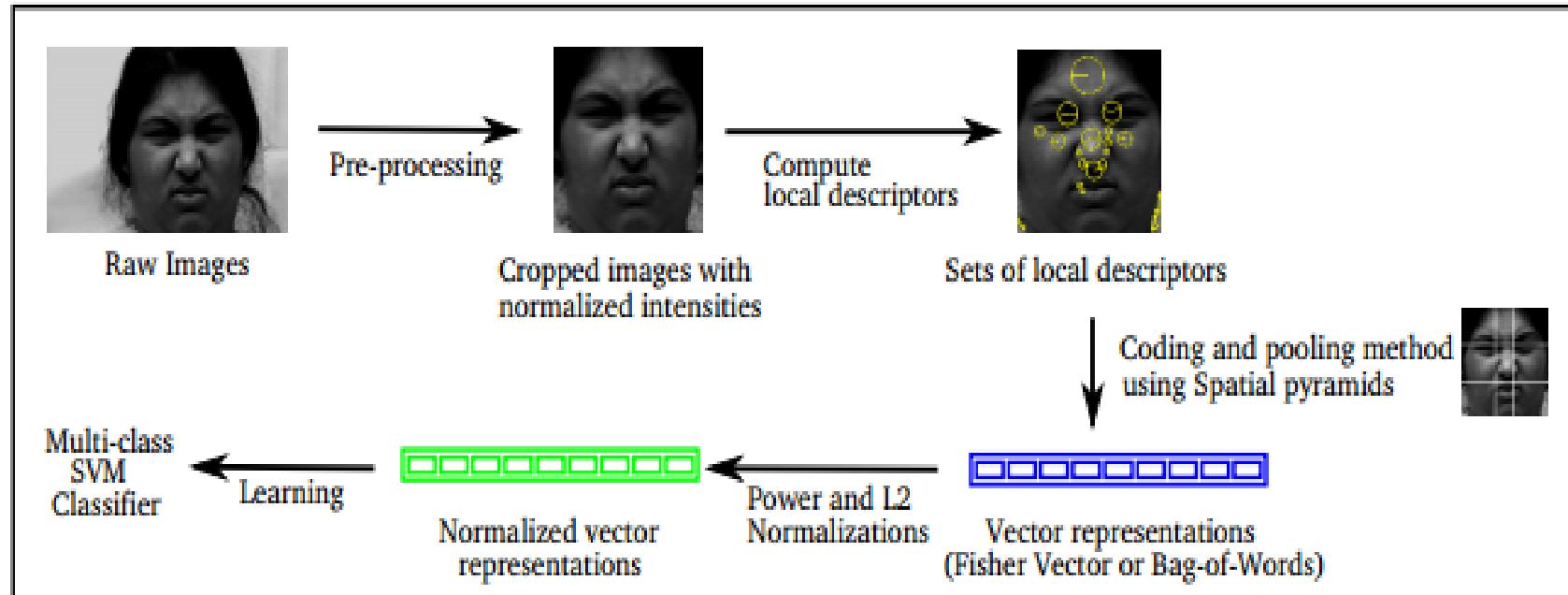


Figure 4: Pipeline of our learning algorithm.



天津大学

Calorie Estimation from Food Images

Dynamic Eye Gaze Tracking and Prediction

Predicting Short-term Market Response to Liquidity Shocks

Predicting Student Earnings After College

Predicting Stock Prices through Textual Analysis of Web News

Human Activity Recognition Using Cell Phones

Human Activity Recognition Using Wearable Devices Sensor Data

3D model classification using convolutional neural network

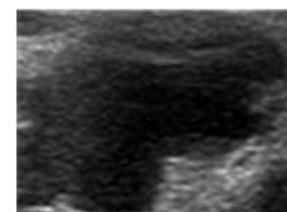


基于机器学习的甲状腺结节分级方法

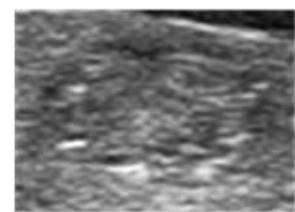
课题组采用多尺度融合方法对甲状腺结节超声图像特征进行提取，包括：实性结节内部组织，边界，回声，钙化程度。然后基于机器学习方法对甲状腺结节实现分级，目前该工作正在完善中。左图为甲状腺疾病分级图，右图为甲状腺结节超声图像。

TI-RADS ^{2.0}				
TI-RADS	评价	恶性概率	超声表现示例	建议
0	没有结节		正常甲状腺或甲状腺弥漫性增大	
1	高度提示为良性	0-7%	以囊性为主的结节，高度提示良性	如果没有临床需求超声不需随访
2	良性可能性大	8-25%	以实性为主，边界清楚，或者囊壳样，粗大钙化	如果临床需求可以长期超声随访
3	极向于良性	26-50%	以实性为主，等回声，边界尚清	3-6个月后随访
4	极向于恶性	51-90%	实�性，低回声，微钙化，边界模糊或成微分叶、纵横比>1 具有一种恶性征象	穿刺活检或手术，即使低风险的结节，也要定期随访 活检，或6个月后随访
4A				
4B			具有两种恶性征象	建议活检
4C			具有三或四种恶性征象	建议活检或手术
5	高度提示为恶性	>91%	超过四项恶性征象，尤其是有微钙化和微分叶 丝病理证实的甲状腺恶性病变	无论穿刺活检结果如何均考虑手术切除
6				

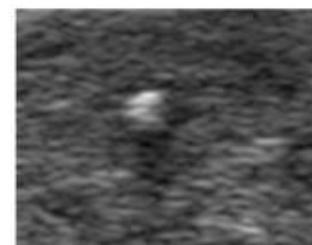
注：甲状腺可疑恶性结节同时合并颈部的转移淋巴结，TI-RADS 可定为 5 级。甲状腺乳头状癌颈部转移淋巴结具有特征性的超声表现。



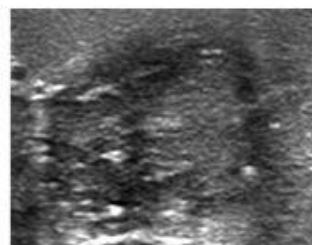
(a) 囊性结节



(b) 实性结节



(c) 粗大钙化

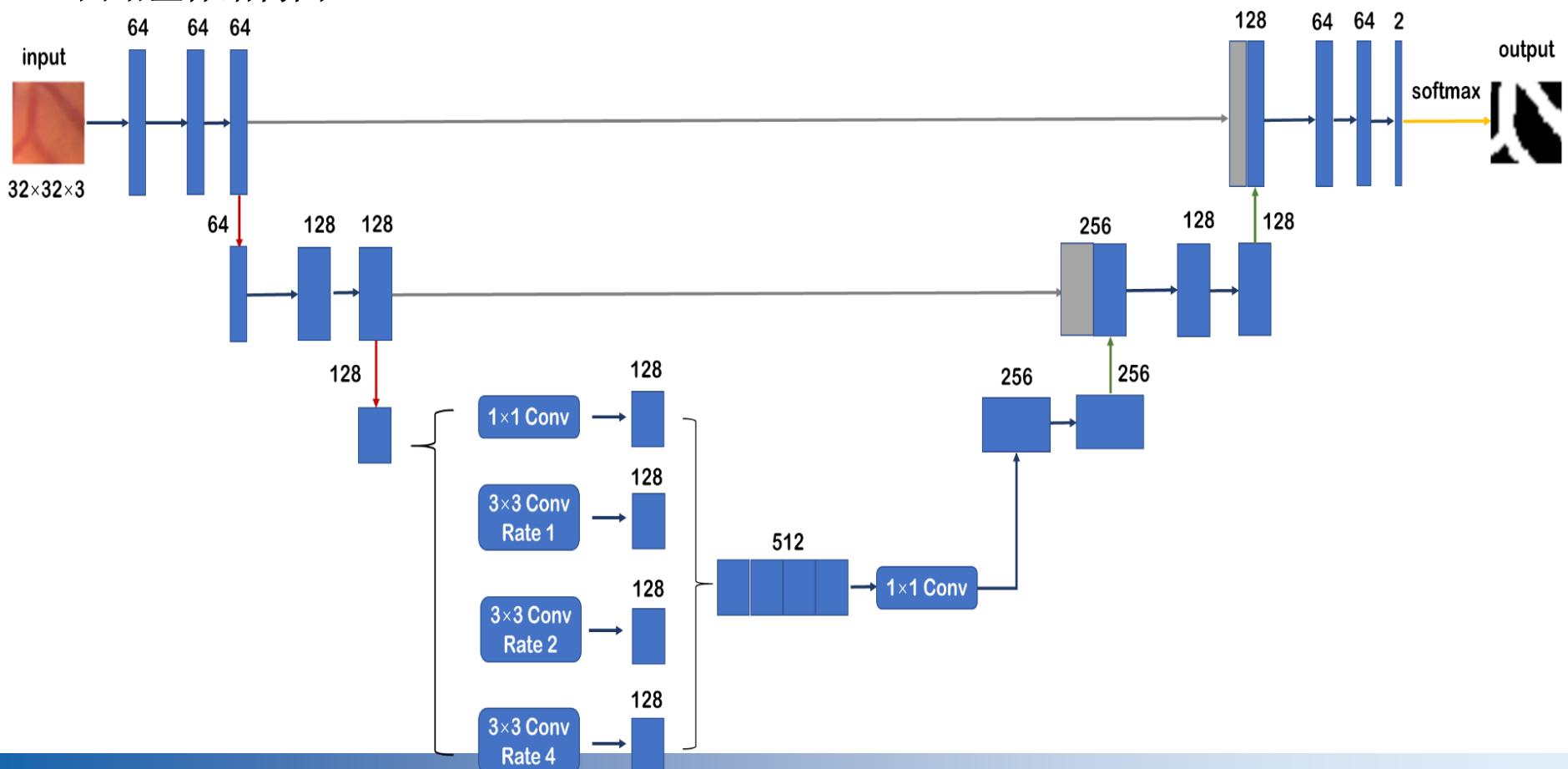


(d) 微钙化



基于结合多尺度特征的卷积神经网络的视网膜血管分割方法

➤ 网络整体结构图

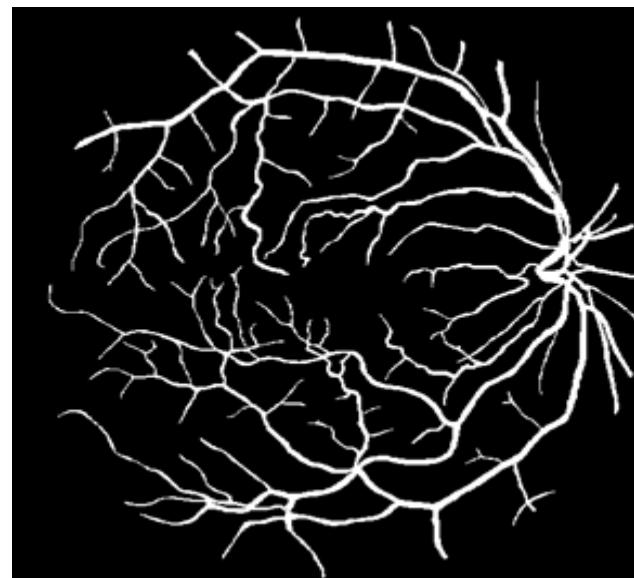




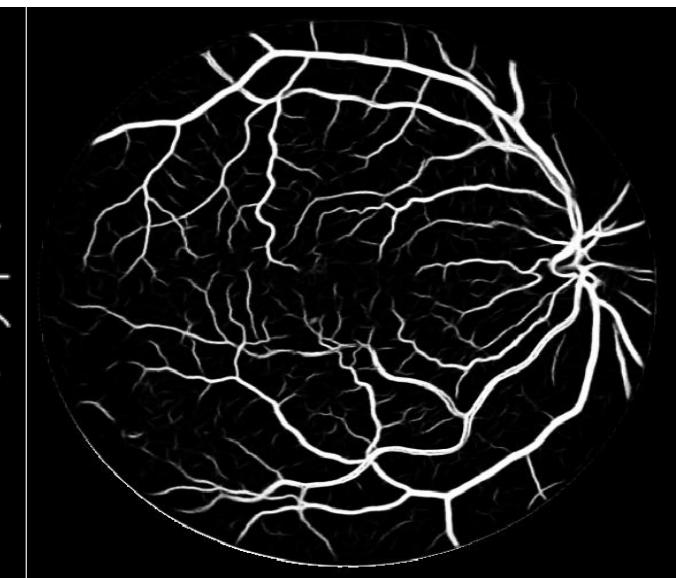
➤ 血管分割图



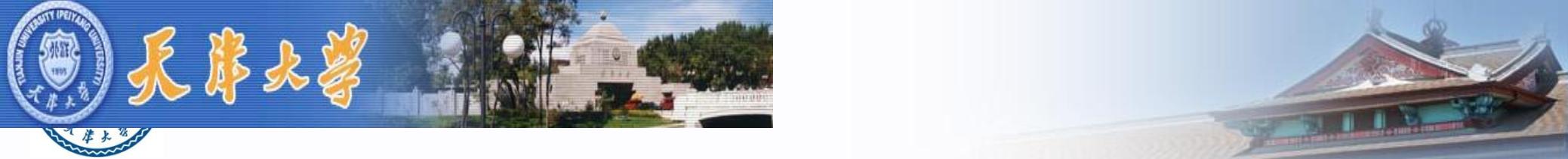
原始图像



第一专家手动分割图



视网膜血管分割图



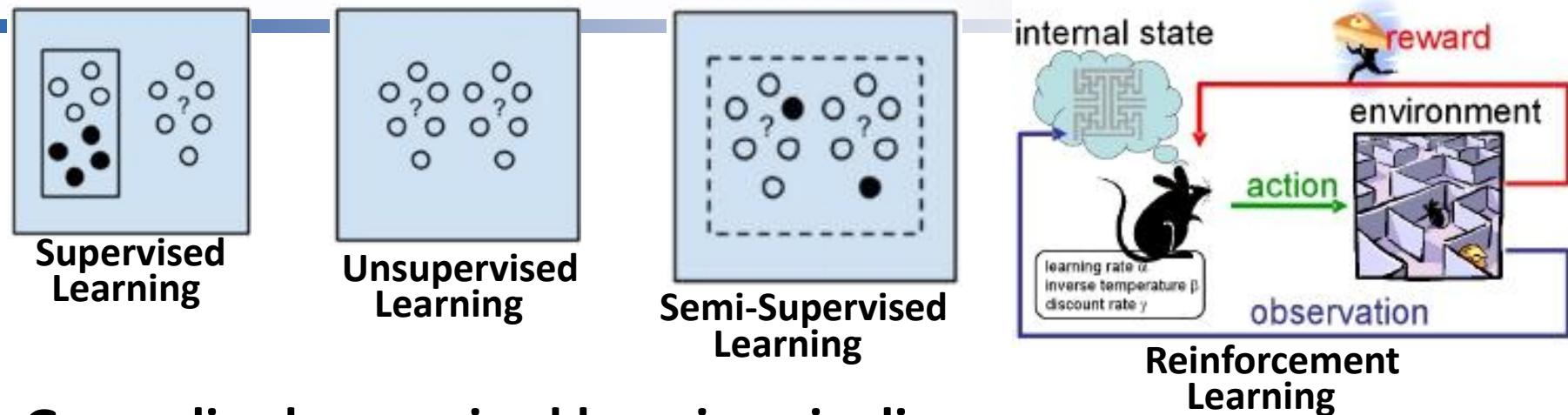
Intelligent transportation;
Intelligent home furnishing;
Intelligent home appliances;
Intelligent buildings;
Intelligent internet.



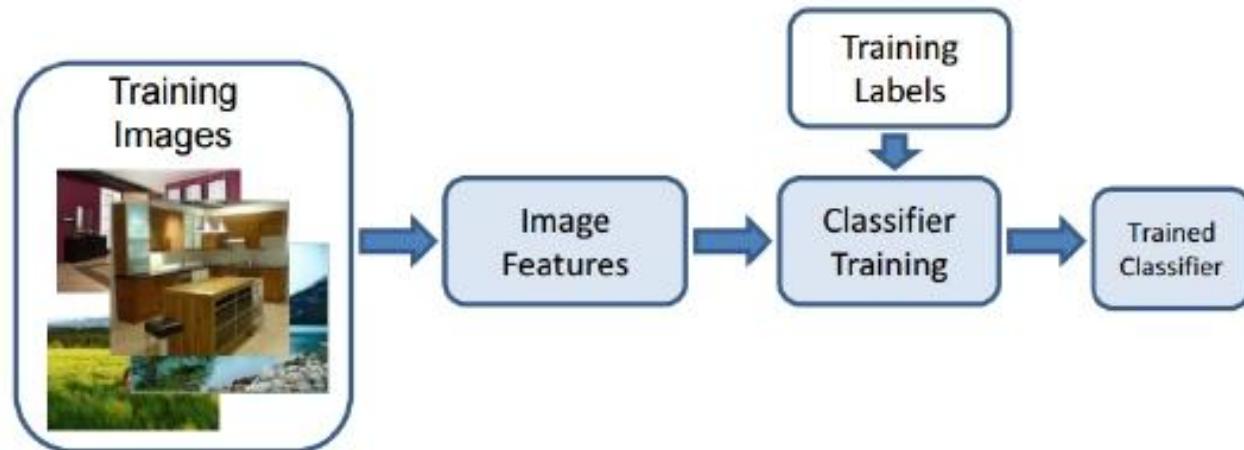
Intelligence: Simulate and extend human behavior.



Types of machine learning:



Generalized supervised learning pipeline:





Regression vs Classification

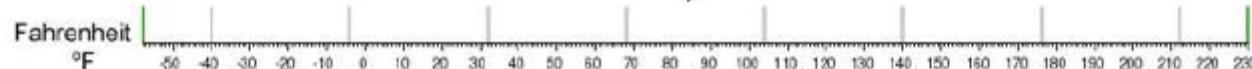


Regression

What is the temperature going to be tomorrow?

PREDICTION

84°

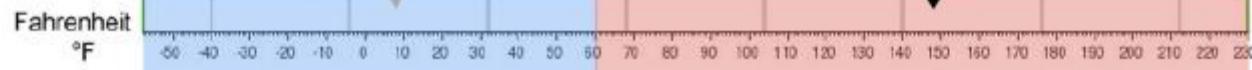


Classification

Will it be Cold or Hot tomorrow?

PREDICTION

HOT





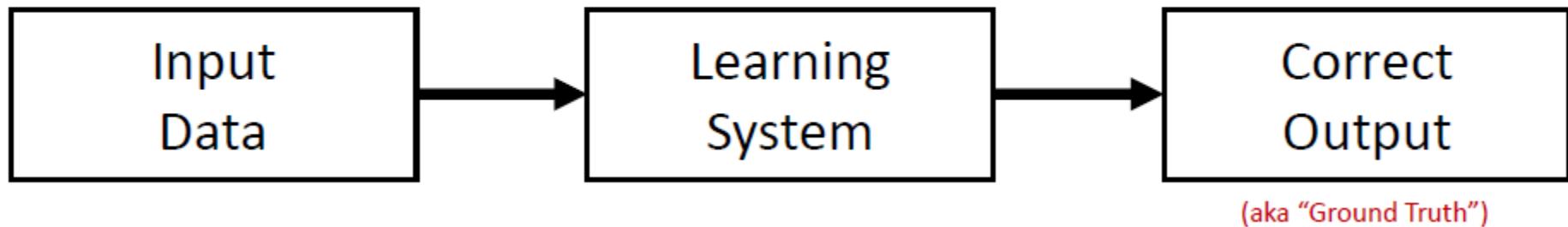
Multi-Class vs Multi-Label

	Multi-Class	Multi-Label
$C = 3$	Samples	Samples
	Labels (t)	Labels (t)
	$[0 \ 0 \ 1]$ $[1 \ 0 \ 0]$ $[0 \ 1 \ 0]$	$[1 \ 0 \ 1]$ $[0 \ 1 \ 0]$ $[1 \ 1 \ 1]$

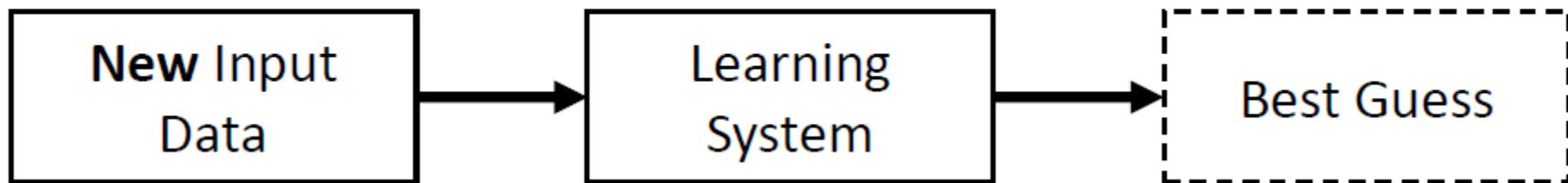


Machine Learning: Training and Testing

Training Stage:



Testing Stage:

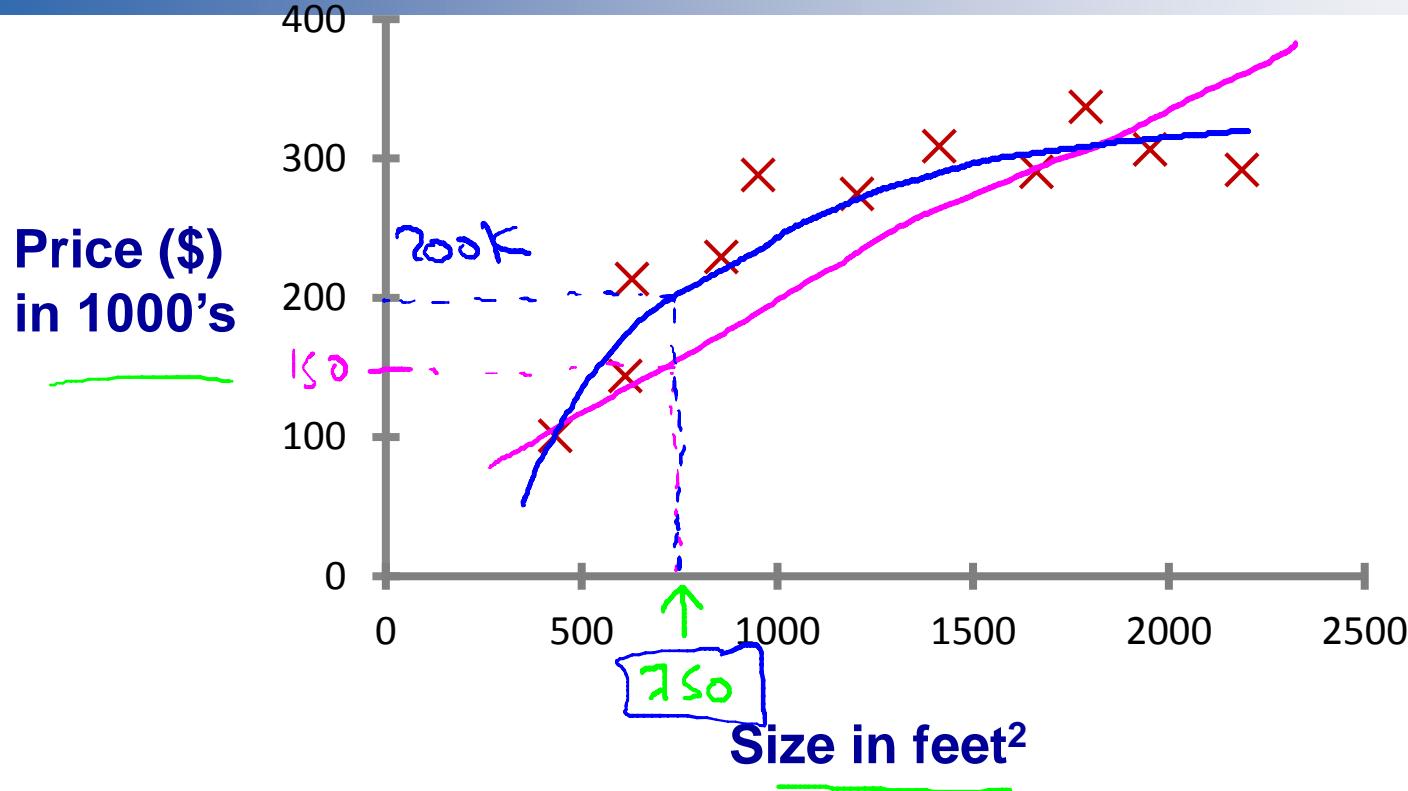


-Supervised learning

Housing price prediction.



天津大学



Supervised Learning
"right answers" given

Regression: Predict
continuous valued
output (price)

天津大学

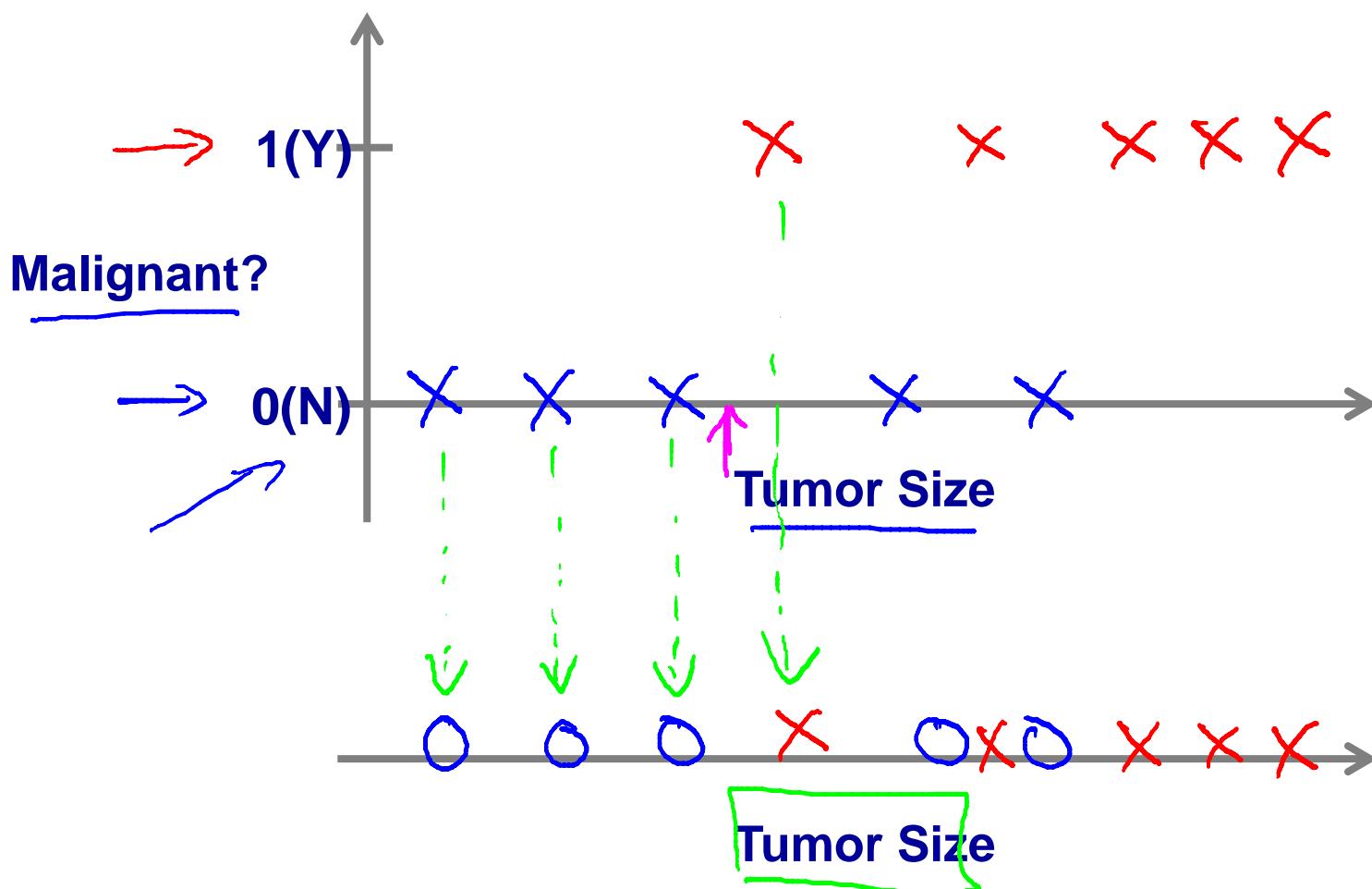
Tianjin University



天津大学

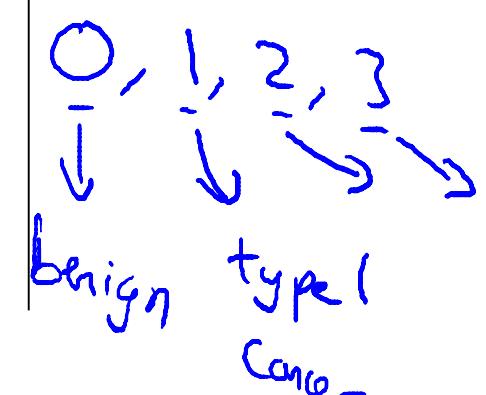


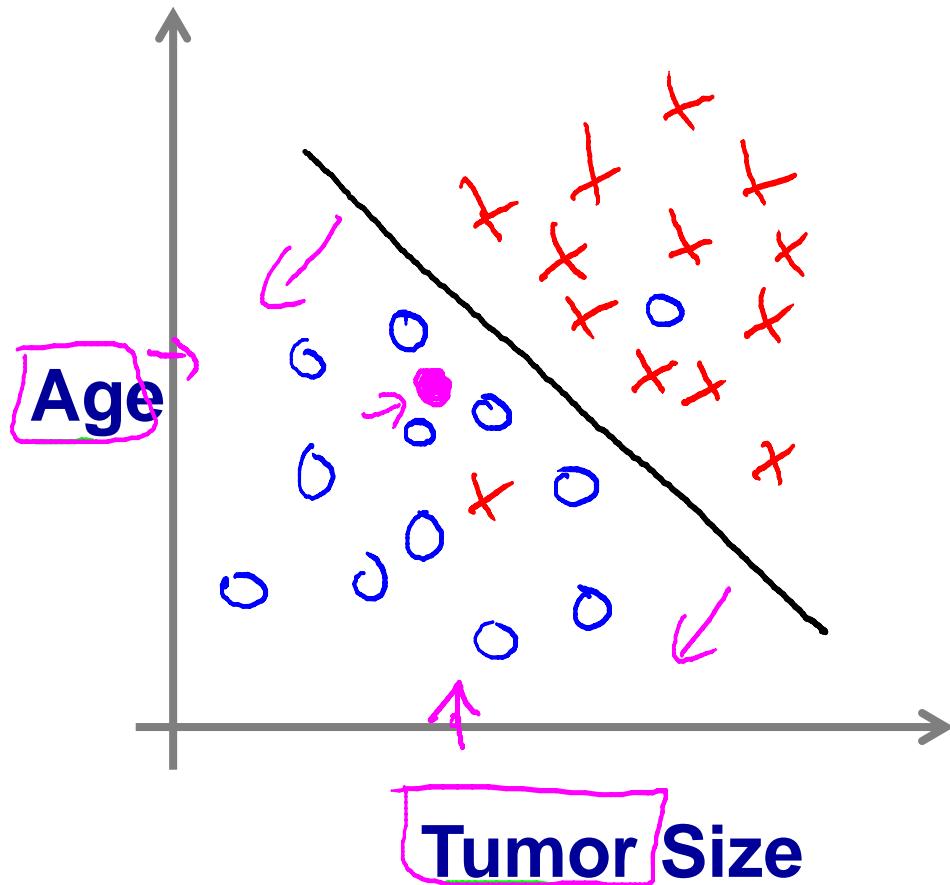
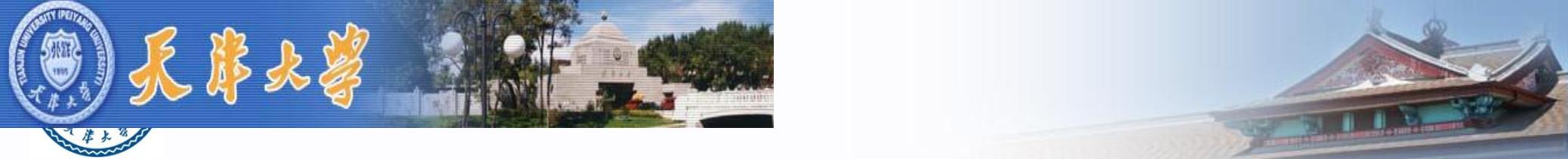
Breast cancer (malignant, benign)



Classification

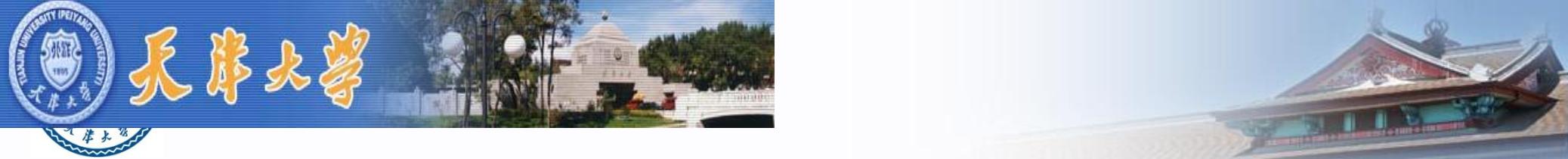
Discrete valued output
(0 or 1)





- Clump Thickness
- Uniformity of Cell Size
- Uniformity of Cell Shape

...



You are running a company, and you want to develop learning algorithms to address each of two problems.

Problem1: You have a large inventory of identical items. You want to predict how many of these items will sell over the next 3 months.

Problem2: You would like software to examine individual customer accounts, and for each account decide if it has been hacked/compromised.

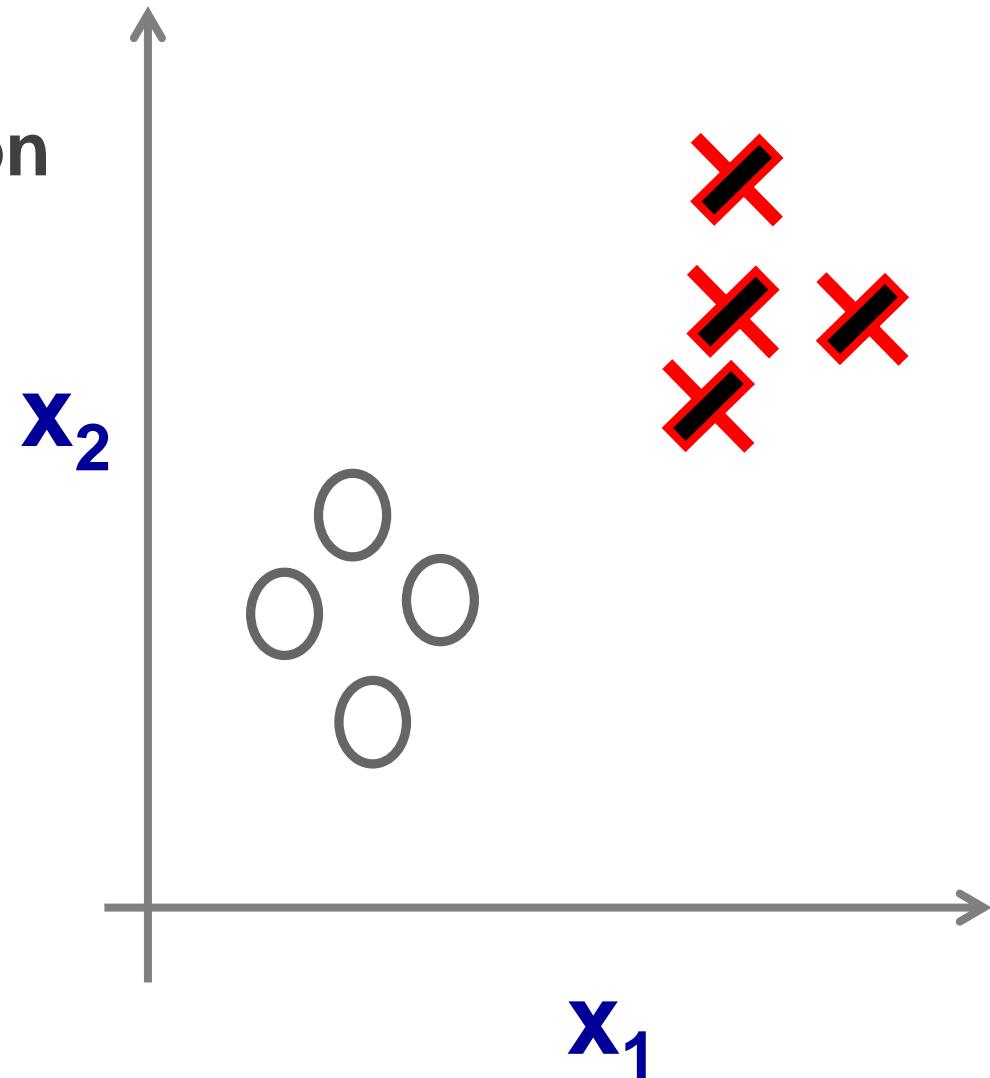
Should you treat these as classification or as regression problems?

- A. Treat both as classification problems.
- B. Treat problem 1 as a classification problem, problem 2 as a regression problem.
- C. Treat problem 1 as a regression problem , problem 2 as a classification problem.
- D. Treat both as regression problems.



Supervised Learning

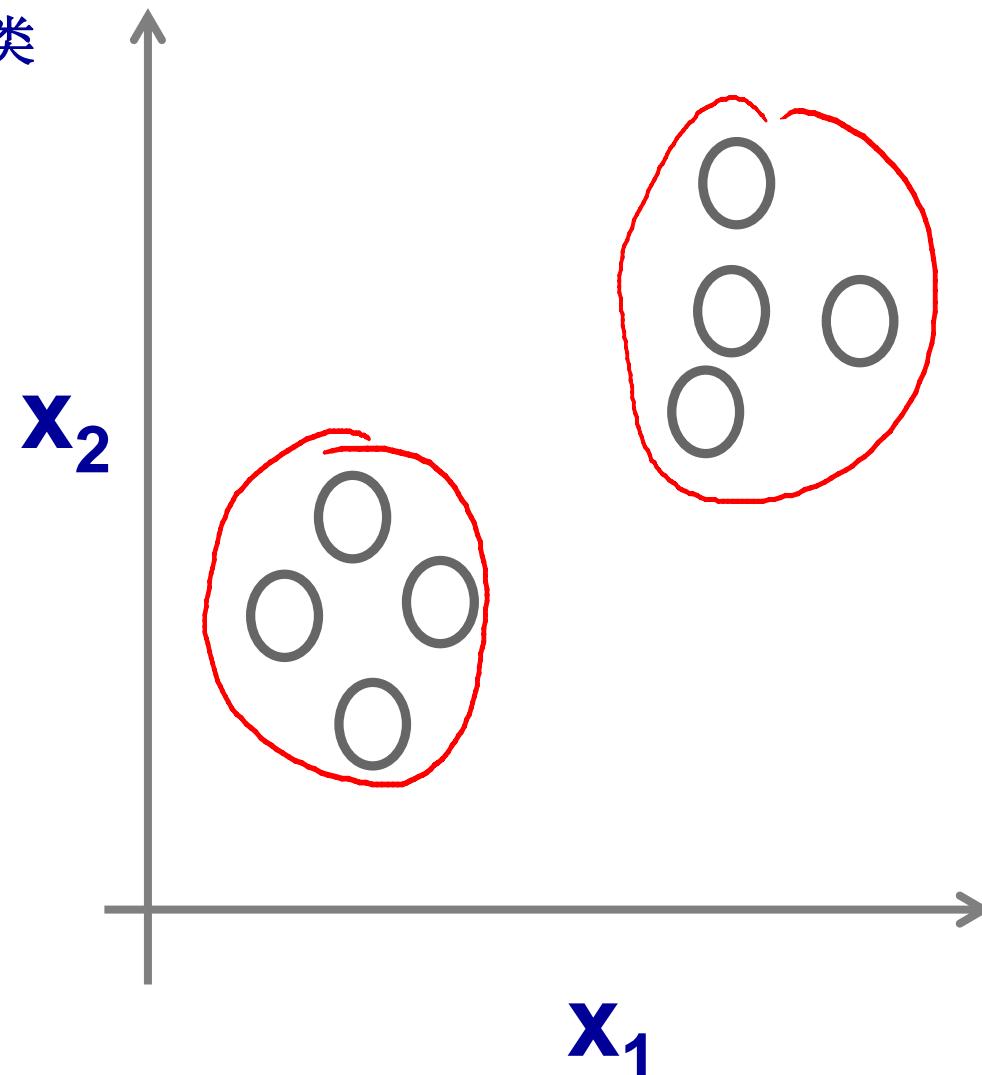
Classification
分类





Unsupervised Learning

Clustering : 聚类



news.google.com



Of the following examples, which would you address using an unsupervised learning algorithm? (Check all that apply.)

Given email labeled as spam/not spam, learn a spam filter.

Given a set of news articles found on the web, group them into set of articles about the same story.

Given a database of customer data, group customers into different market segments.

Given a dataset of patients diagnosed as either having diabetes or not, learn to classify new patients as having diabetes or not.



第一次作业：基于WIFI信号强度的室内定位

数据集: Collected in indoor space by observing signal strengths of seven WiFi signals visible on a smartphone. The decision variable is one of the four rooms.

Data Set Characteristics:	Multivariate	Number of Instances:	2000	Area:	Computer
Attribute Characteristics:	Real	Number of Attributes:	7	Date Donated	2017-12-04
Associated Tasks:	Classification	Missing Values?	N/A	Number of Web Hits:	7850

选择合适算法做实验，按作业模板完成报告

报告与实验代码打成一个压缩包，提交至MachineLearningTJU@126.com

压缩包命名：作业1-姓名-学号

提交日期：第八周周日（）



Lecture_2



Matlab Introduction

- (1) Numerical calculation(Optimization algorithm);
- (2) 2D and 3D plot;
- (3) Image processing



roots

Solving a equation

Example

$$2x^5 - 3x^3 + 71x^2 - 9x + 1 = 0$$

$$p = [2, 0, -3, 71, -9, 1];$$

$$x = \text{roots}(p)$$

x =

```
-3.4774 + 0.0000i
1.6753 + 2.7105i
1.6753 - 2.7105i
0.0634 + 0.1007i
0.0634 - 0.1007i
```



【Functional demonstration-2】 Solving a linear equation

$$\begin{cases} 2x + 3y - z = 2 \\ 8x + 2y + 3z = 4 \\ 45x + 3y + 9z = 23 \end{cases}$$

a = [2,3,-1;8,2,3;45,3,9]; %建立系数矩阵a

b = [2;4;23]; %建立列向量b

x = inv(a)*b

x =

0.5531

0.2051

-0.2784



Symbolic calculation

syms x y z %建立符号变量

[x,y,z]=solve(2*x+3*y-z-2,8*x+2*y+3*z-4,45*x+3*y+9*z-23)

x =

151/273

y =

8/39

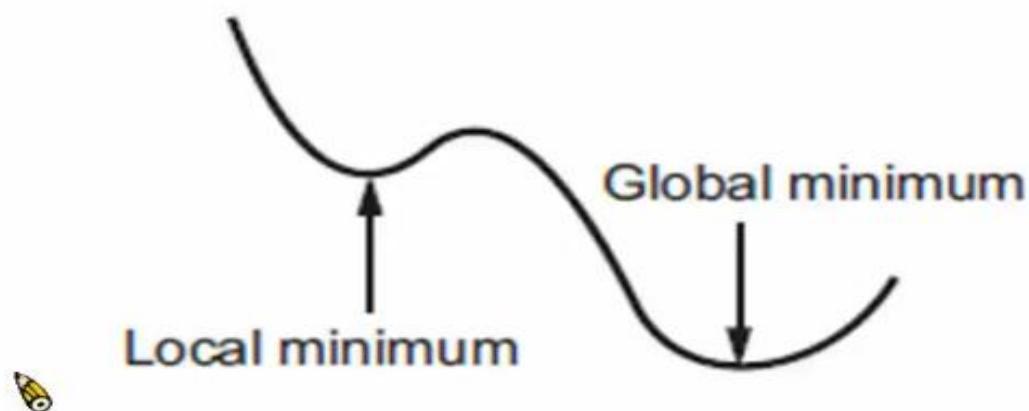
z =

-76/273



Optimization algorithm

WHAT IS GLOBAL OPTIMIZATION?





OPTIMIZATION WORKFLOW

- Local or global solution
- Write objective function (constraint functions)
- Set appropriate options
- Run the solver
- Examine the result
- If the result is unsatisfactory, change options or start points or otherwise update your optimization and rerun it.



fmincon

Find minimum of constrained nonlinear multivariable function

Finds the minimum of a problem specified by

$$\min_x f(x) \text{ such that} \begin{cases} c(x) \leq 0 \\ ceq(x) = 0 \\ A \cdot x \leq b \\ Aeq \cdot x = beq \\ lb \leq x \leq ub, \end{cases}$$

b and beq are vectors, A and Aeq are matrices, c(x) and ceq(x) are functions that return vectors, and f(x) is a function that returns a scalar. f(x), c(x), and ceq(x) can be nonlinear functions.



1. **$x = \text{fmincon}(\text{fun}, x_0, A, b)$** starts at x_0 and attempts to find a minimizer x of the function described in fun subject to the linear inequalities $A^*x \leq b$. x_0 can be a scalar, vector, or matrix.
2. **$x = \text{fmincon}(\text{fun}, x_0, A, b, A_{\text{eq}}, b_{\text{eq}})$** minimizes fun subject to the linear equalities $A_{\text{eq}}^*x = b_{\text{eq}}$ and $A^*x \leq b$. If no inequalities exist, set $A = []$ and $b = []$.
3. **$x = \text{fmincon}(\text{fun}, x_0, A, b, A_{\text{eq}}, b_{\text{eq}}, lb, ub)$** defines a set of lower and upper bounds on the design variables in x , so that the solution is always in the range $lb \leq x \leq ub$. If no equalities exist, set $A_{\text{eq}} = []$ and $b_{\text{eq}} = []$. If $x(i)$ is unbounded below, set $lb(i) = -\text{Inf}$, and if $x(i)$ is unbounded above, set $ub(i) = \text{Inf}$.



Example

Find values of x that minimize $f(x) = -x_1 x_2 x_3$, starting at the point $x = [10;10;10]$, subject to the constraints:

$$0 \leq x_1 + 2x_2 + 2x_3 \leq 72.$$

%1 Write a file that returns a scalar value f of the objective function evaluated at x:

```
function f = myfun(x)
f = -x(1) * x(2) * x(3);
```



```
1 - clear;
2
3
4
5 %Rewrite the constraints
6 - -x(1)-2*x(2)-2*x(3) <= 0;
7 - x(1) + 2*x(2) +2*x(3) <= 72;
8 % Since both constraints are linear, formulate them as the matrix inequality A · x ≤ b, where
9 - A = [-1 -2 -2; ...
10      1  2  2];
11 - b = [0;72];
12 % Supply a starting point and invoke an optimization routine:
13 - x0 = [10;10;10]; % Starting guess at the solution
14 - [x,fval] = fmincon(@myfun,x0,A,b);
```



fminunc

Find minimum of unconstrained multivariable function

Finds the minimum of a problem specified by

$$\min f(x)$$

where $f(x)$ is a function that returns a scalar. x is a vector or a matrix

Examples

minimize the function $f(x) = 3x_1^2 + 2x_1x_2 + x_2^2$



```
function f = myfun(x)
f = 3*x(1)^2 + 2*x(1)*x(2) + x(2)^2; % Cost function
```

```
%Then call fminunc to find a minimum of myfun near [1,1]:
clear
x0 = [1, 1];
[x, fval] = fminunc(@myfun, x0);

x
fval
```

minimize the function $f(x) = 3x_1^2 + 2x_1x_2 + x_2^2$



minimize the function $f(x) = \sin(x) + 3$ by using matlab software

Constructing function

```
- function f = myfun(x)
  f = sin(x)+3;
```

Calling function

```
[x,fval] = fminunc(@myfun,0)
```



Plot: 2-D line plot;

Mesh: 3-D line plot

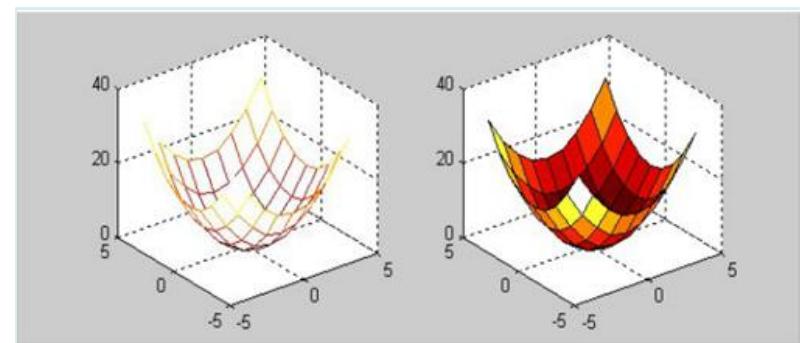


Example

Drawing the surface of function, $z=x^2+y^2$

```
x=-4:4;y=x;  
[x,y]=meshgrid(x,y);      %生成 x-y 坐标 “格点”  
                           矩阵  
z=x.^2+y.^2;              %计算格点上的函数值  
subplot(1,2,1), mesh(x,y,z); %三维网格图  
  
subplot(1,2,2), surf(x,y,z); %三维曲面图  
colormap(hot);
```

The obtained surface of function, $z=x^2+y^2$





Matlab for image processing

imread

Read image from graphics file

A = imread('ngc6543a.jpg');

imshow(A);

imwrite

Write image to graphics file



adapthisteq (adapt hist eq)

Contrast-limited adaptive histogram equalization (CLAHE)

J = adapthisteq(I)

Performs CLAHE on the intensity image I.





bwarea(BW)

total = bwarea(BW) estimates the area of the objects in binary image BW. total is a scalar whose value corresponds roughly to the total number of on pixels in the image.



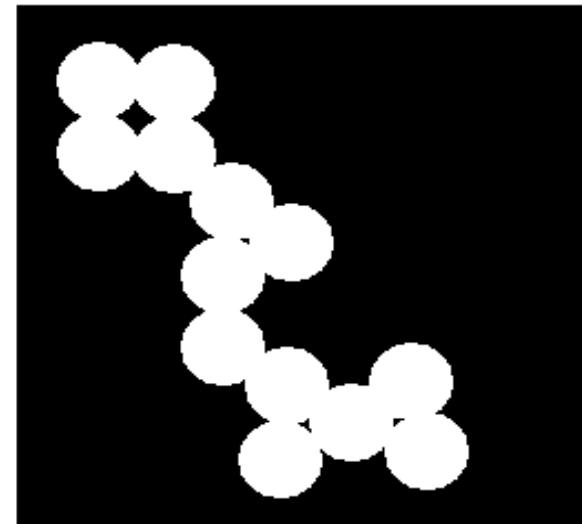
Example

```
BW = imread('circles.png');
```

```
imshow(BW);
```

```
bwarea(BW)
```

```
ans =  
1.4187e+004
```





bwareaopen

Remove small objects from binary image

Remove all objects smaller than 50 pixels

Example

```
originalBW = imread ('text.png');
```

```
imshow(originalBW)
```

```
bwAreaOpenBW =
```

```
bwareaopen(originalBW,50);
```

```
figure, imshow(bwAreaOpenBW)
```

The term watershed
refers to a ridge that ...

... divides areas
drained by different
river systems.

original BW



The term watershed
refers to a ridge that ...

... divides areas
drained by different
river systems.

original BW

The e m wa e shed
efe s o a dge ha

d v des a eas
d a ned by d ffe en
ve sys ems

bwareaopen BW



bwboundaries

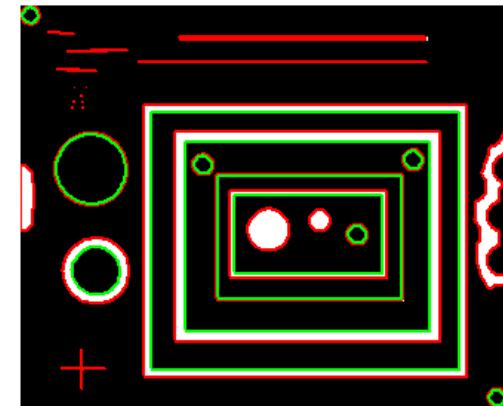
Trace region boundaries in binary image



Example

Display object boundaries in red and hole boundaries in green.

```
BW = imread('blobs.png');
[B,L,N] = bwboundaries(BW);
figure; imshow(BW); hold on;
for k=1:length(B),
    boundary = B{k};
    if(k > N)
        plot(boundary(:,2),...
              boundary(:,1),'g','LineWidth',2);
    else
        plot(boundary(:,2),...
              boundary(:,1),'r','LineWidth',2);
    end
end
```





bwmorph

Morphological operations on binary images

'clean'

Removes isolated pixels

'fill'

Fills isolated interior pixels

'remove'

Removes interior pixels.

'skel'

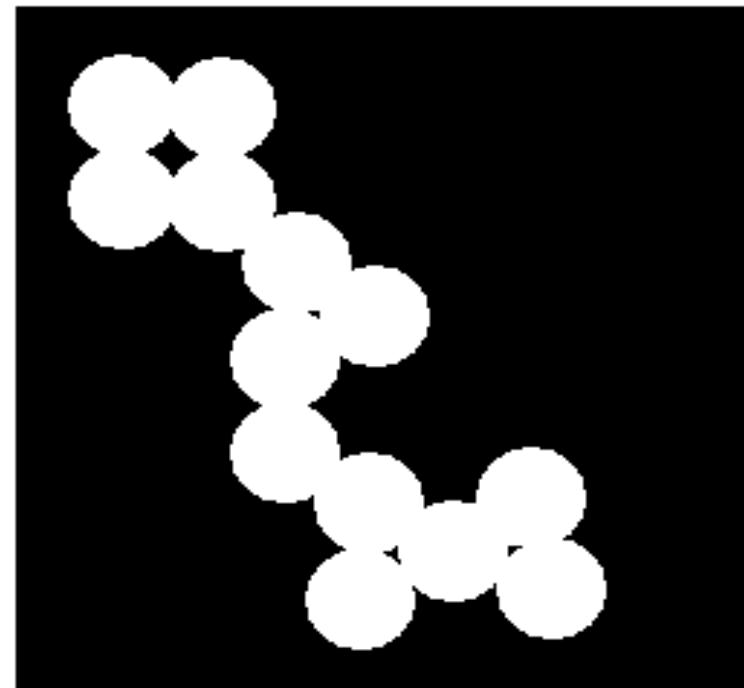
allow removes pixels on the boundaries
of objects without allowing objects to
break apart.



Example

```
BW = imread('circles.png');
```

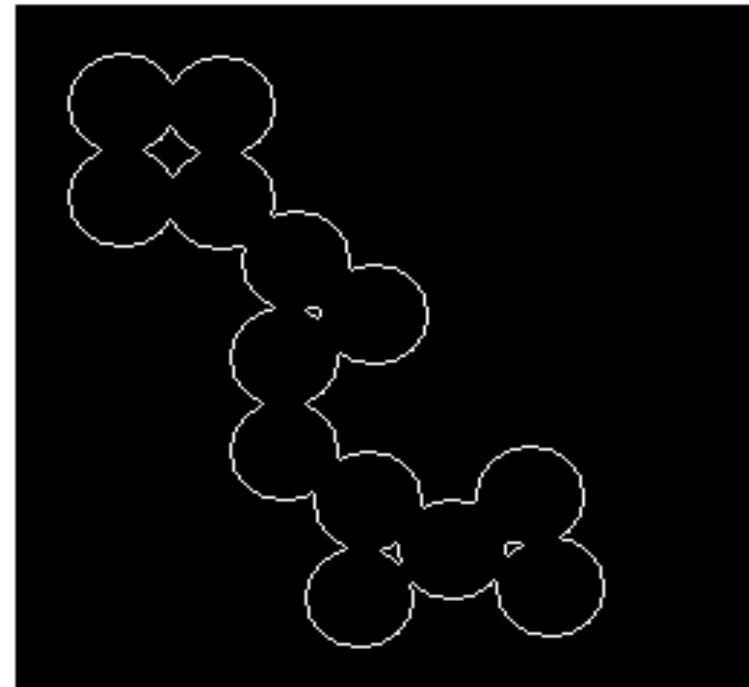
```
imshow(BW);
```





BW2 = bwmorph(BW,'remove');

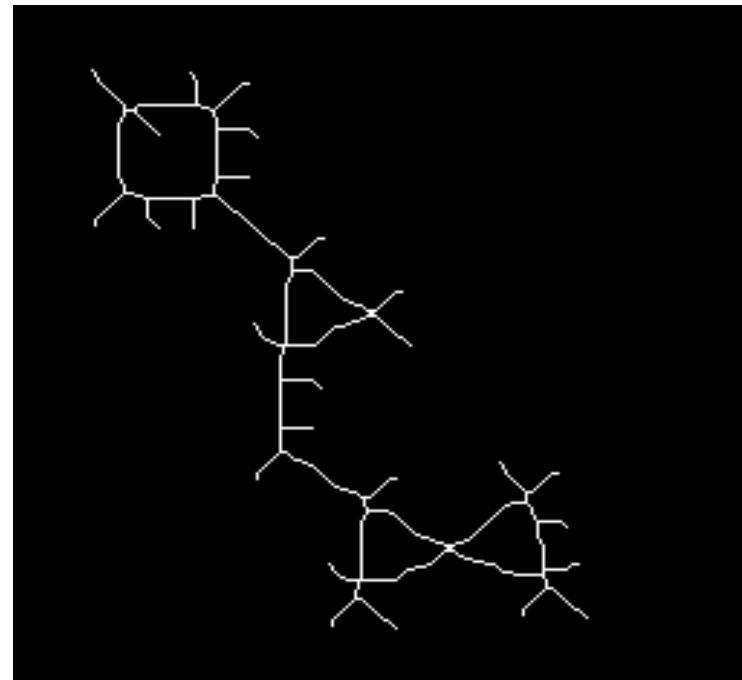
figure, imshow(BW2)





BW3 = bwmorph(BW,'skel',Inf);

figure, imshow(BW3)



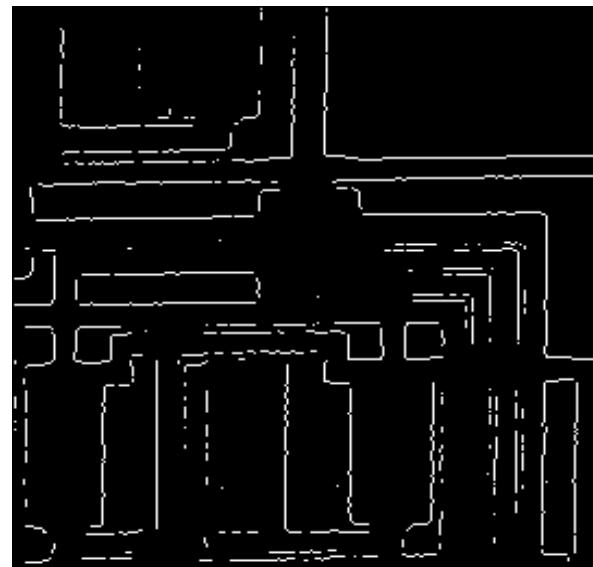


edge

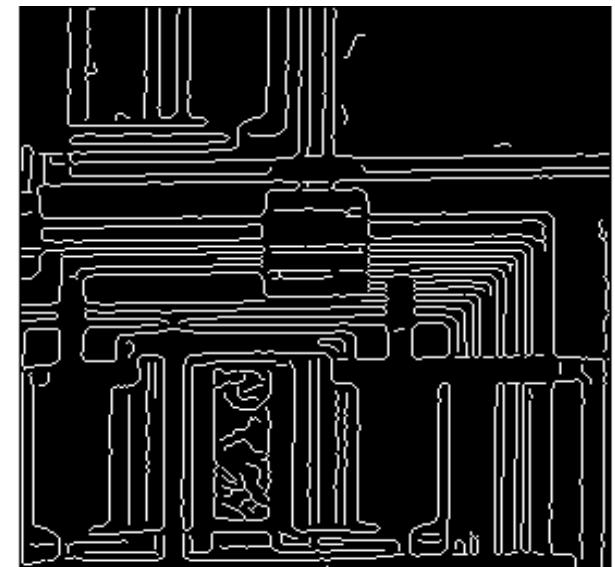
Find edges in grayscale image

Examples: Find the edges of an image using the Prewitt and Canny methods.

```
I = imread('circuit.tif');
BW1 = edge(I,'prewitt');
BW2 = edge(I,'canny');
imshow(BW1);
figure, imshow(BW2)
```



Prewitt Filter



Canny Filter



gray2ind

Convert grayscale or binary image to indexed image

gray2ind scales, then rounds, an intensity image to produce an equivalent indexed image.

Example

```
I = imread('cameraman.tif');
```

```
[X, map] = gray2ind(I, 16);
```

```
imshow(X, map);
```





imadjust

Adjust image intensity values

Examples

Adjust a low-contrast grayscale image.

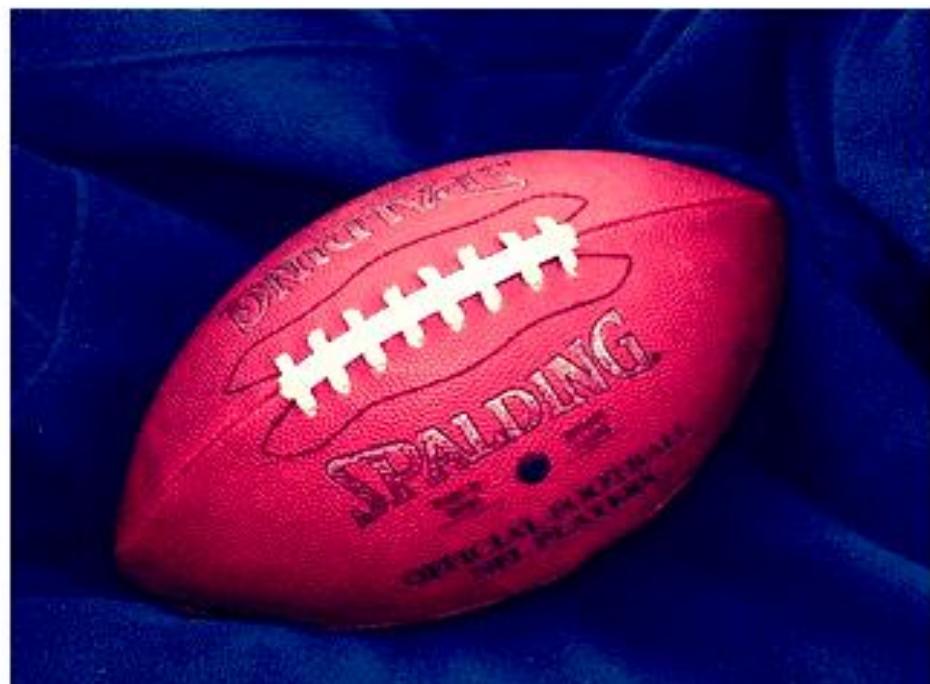
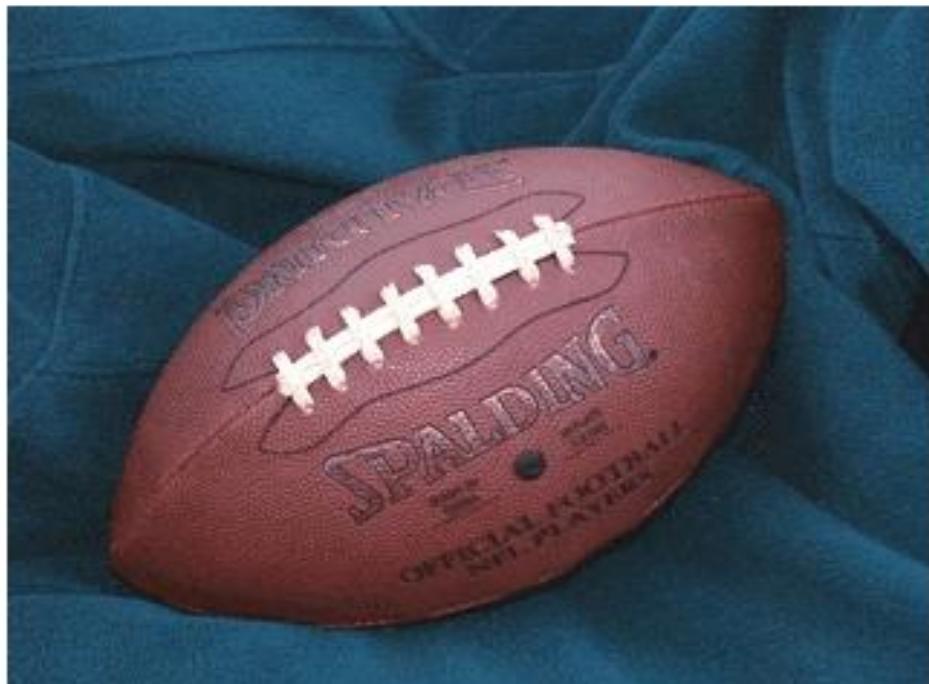
```
I = imread('pout.tif');
J = imadjust(I);
imshow(I), figure, imshow(J)
```





Adjust an RGB image.

```
RGB1 = imread('football.jpg');
RGB2 = imadjust(RGB1,[.2 .3 0; .6 .7 1],[]);
imshow(RGB1), figure, imshow(RGB2)
```





imfilter

N-D filtering of multidimensional images

imhist

Display histogram of image data

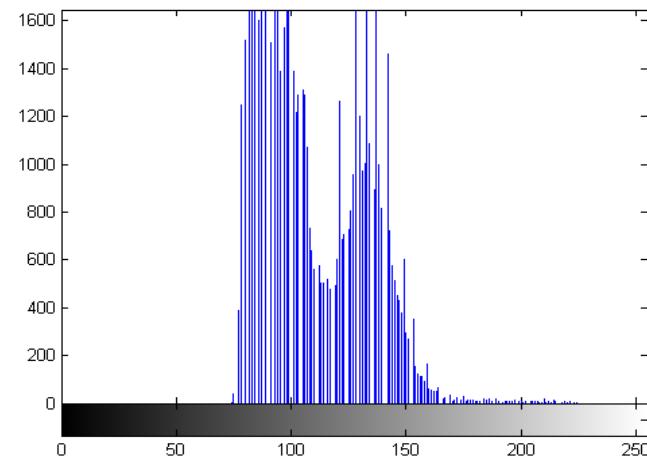
Example

```
I = imread('pout.tif');
```

```
imshow(I);
```

```
figure(2)
```

```
imhist(I)
```





imnoise

Add noise to image

Description

$J = \text{imnoise}(I, \text{type})$ adds noise of a given type to the intensity image I . type is a string that specifies any of the following types of noise. Note that certain types of noise support additional parameters. See the related syntax.

Value	Description
'gaussian'	Gaussian white noise with constant mean and variance
'localvar'	Zero-mean Gaussian white noise with an intensity-dependent variance
'poisson'	Poisson noise
'salt & pepper'	On and off pixels
'speckle'	Multiplicative noise



Examples

```
I = imread('eight.tif');
J = imnoise(I, 'salt & pepper', 0.02);
figure, imshow(I)
figure, imshow(J)
```





medfilt2

2-D median filtering

Examples

Add salt and pepper noise to an image and then restore the image using `medfilt2`.

```
I = imread('eight.tif');
J = imnoise(I,'salt & pepper',0.02);
K = medfilt2(J);
imshow(J), figure, imshow(K)
```





mat2gray

Convert matrix to grayscale image

mean2

Average or mean of matrix elements



rgb2gray

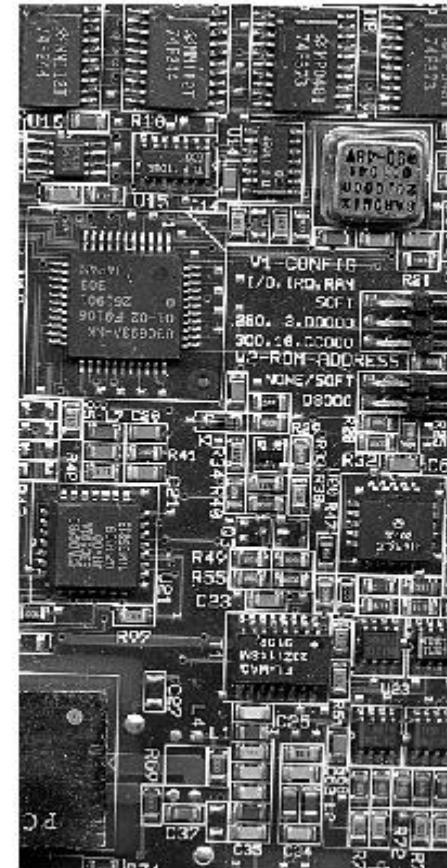
Convert RGB image or color map to grayscale

Example

I = imread('board.tif');

J = rgb2gray(I);

figure, imshow(I), figure, imshow(J);





Some basic concepts for matrices



Identity Matrix

Denoted I (or $I_{n \times n}$).

Examples of identity matrices:

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

2 x 2

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

3 x 3

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

4 x 4

For any matrix A ,

$$A \cdot \boxed{I} = \boxed{I} \cdot A = A$$

\uparrow \uparrow \uparrow \uparrow
 $m \times n$ $n \times n$ $m \times m$ $m \times n$ $m \times n$

$$I_{n \times n}$$

Note:
 $AB \neq BA$ In general
 $AI = IA$



The transpose matrix \mathbf{A}^T has elements $(\mathbf{A}^T)_{ij} = A_{ji}$.

From the definition of transpose, we have

$$(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T \quad (\text{C.1})$$

The inverse of \mathbf{A} , denoted \mathbf{A}^{-1} ,

Satisfies

$$\mathbf{AA}^{-1} = \mathbf{A}^{-1}\mathbf{A} = \mathbf{I}. \quad (\text{C.2})$$

Because $\mathbf{ABB}^{-1}\mathbf{A}^{-1} = \mathbf{I}$, we have:

$$(\mathbf{AB})^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1}. \quad (\text{C.3})$$

Because

$$\mathbf{AA}^{-1} = \mathbf{I}$$

$$(\mathbf{AA}^{-1})^T = (\mathbf{A}^{-1})^T \mathbf{A}^T = \mathbf{I}$$

So have:

$$(\mathbf{A}^T)^{-1} = (\mathbf{A}^{-1})^T \quad (\text{C.4})$$



The transpose matrix of A

A'

The inverse of A

inv(A)



Traces and Determinants

Trace and determinant apply to square matrices. The trace $\text{Tr}(A)$ of a matrix A is defined as the sum of the elements on the leading diagonal.

$$\text{Tr}(AB) = \text{Tr}(BA).$$

By applying this formula multiple times to the product of three matrices

$$\text{Tr}(ABC) = \text{Tr}(CAB) = \text{Tr}(BCA)$$



The determinant of a product of two matrices
is given by

$$|\mathbf{AB}| = |\mathbf{A}||\mathbf{B}|$$

$$|\mathbf{A}^{-1}| = 1 / |\mathbf{A}|$$

trace and det



Matrix Derivatives

Sometimes we need to consider derivatives of vectors and matrices with respect to scalars.

The derivative of a vector with respect to a scalar x is itself a vector . Whose components are given by:

$$\left(\frac{\partial \mathbf{a}}{\partial x} \right)_i = \frac{\partial \mathbf{a}_i}{\partial x}$$



with an analogous definition for the derivative of a vector. Derivatives with respect to vectors and matrices can also be defined, for instance:

$$\left(\frac{\partial x}{\partial \mathbf{a}} \right)_i = \frac{\partial x}{\partial \mathbf{a}_i}$$

$$\frac{\partial}{\partial x} (\mathbf{AB}) = \frac{\partial \mathbf{A}}{\partial x} \mathbf{B} + \mathbf{A} \frac{\partial \mathbf{B}}{\partial x}$$



The derivative of the inverse of a matrix can be expressed as

$$\frac{\partial}{\partial x} (A^{-1}) = -A^{-1} \frac{\partial A}{\partial x} A^{-1}$$

$$\begin{aligned}\frac{\partial}{\partial x} (A^{-1}) &= \frac{\partial}{\partial x} (A^{-1})(AA^{-1}) \\&= \frac{\partial A^{-1}}{\partial x} (AA^{-1}) + A^{-1} \frac{\partial AA^{-1}}{\partial x} \\&= \frac{\partial A^{-1}}{\partial x} + A^{-1} \frac{\partial A}{\partial x} A^{-1} + A^{-1} A \frac{\partial A^{-1}}{\partial x} \\&= \frac{\partial A^{-1}}{\partial x} + A^{-1} \frac{\partial A}{\partial x} A^{-1} + \frac{\partial A^{-1}}{\partial x} \\&\Rightarrow \frac{\partial}{\partial x} (A^{-1}) = -A^{-1} \frac{\partial A}{\partial x} A^{-1}\end{aligned}$$



```
>> clear
>> syms a t x; % 定义符号变量
>> f=[a, t^5; t^2*cos(x), log(x)]; % 创建符号函数数组
>> f
f =
 [ a,      t^5]
 [ t^2*cos(x), log(x)]
```



>> dfdx=diff(f) %求矩阵f对x的一阶导数

dfdx =

$$\begin{bmatrix} 0, & 0 \\ -t^2 \sin(x), & 1/x \end{bmatrix}$$





>> dfdt2=diff(f, t, 2) %求矩阵f对t的二阶导数

dfdt2 =

[0, 20*t^3]
[2*cos(x), 0]

>>





Linear Regression & Model Representation



Outline:

1、 Linear regression

2、 Gradient descent

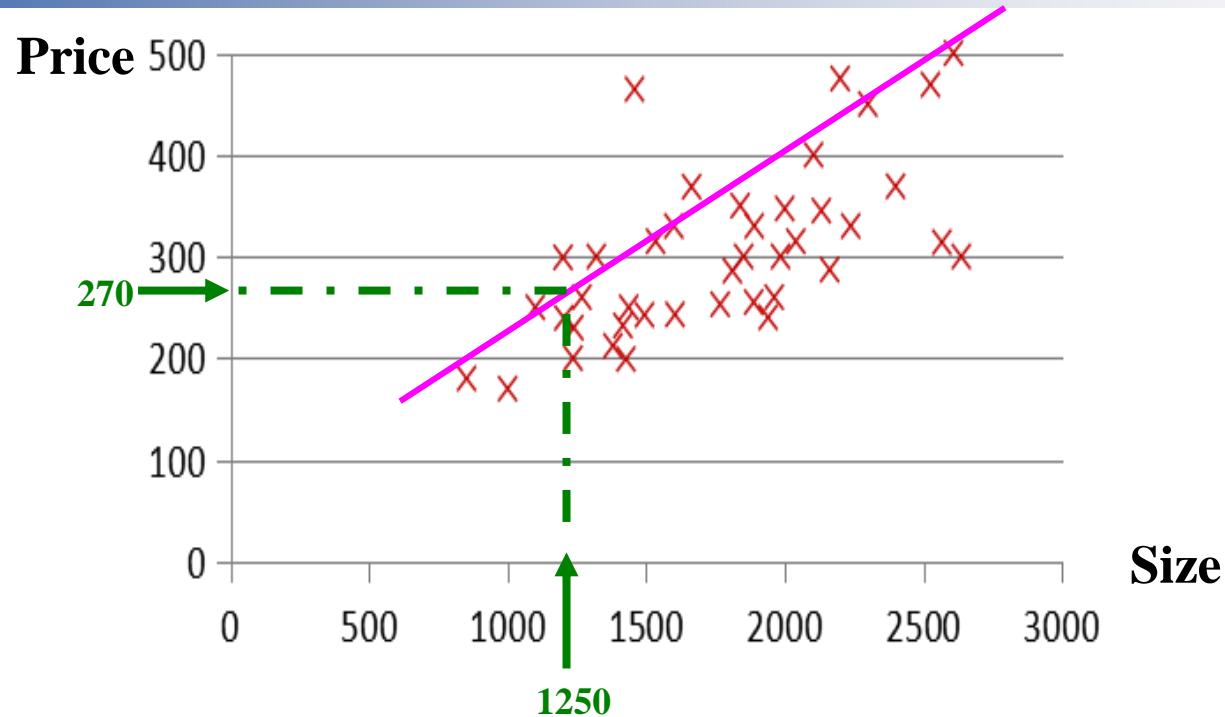
3、 Normal equations



Training set of housing prices

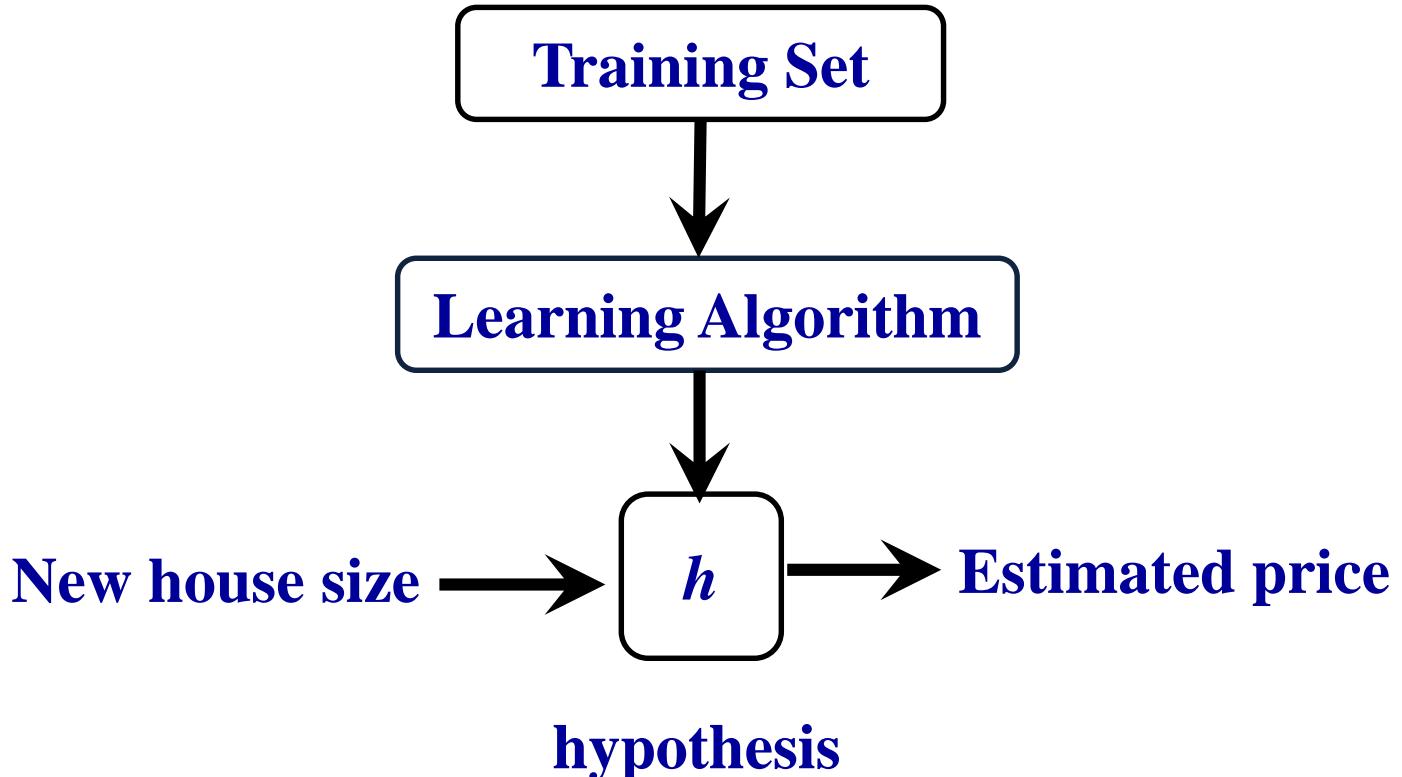
Size in feet ² (x)	Price (\$) in 1000's (y)
2104	460
1416	232
1534	315
852	178
:	:

$N=47$



Supervised Learning:
Given the “right answer” for each example in the data.

Regression Problem:
Predict real-valued output



$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

Here, $\theta_i (i=0,1)$ are the parameters



Notation

N =the number of training examples

x ='input' variables / features

y ='output' variables / 'target' variables

$(x,y) \rightarrow$ training example

i^{th} training example= (x_i, y_i)



Multiple features (variables)

Size (feet ²) $x^{(1)}$	Number of bedrooms $x^{(2)}$	Number of floors $x^{(3)}$	Age of home (years) $x^{(4)}$	Price (\$1000) y
2104	5	1	45	460
1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178
...

$N=47$

Notation:

n : number of features, here $n=4$

x_i : input (features) of i^{th} training example

$x_i^{(j)}$: the value of feature j in i^{th} training example.

$$x_2 = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 40 \end{bmatrix}$$

$$x_2^{(3)} = 2$$



Hypothesis

Previously,

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

$$h_{\theta}(x) = \theta_0 + \theta_1 x^1 + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

$$h_{\theta}(x) = \theta_0 + \theta_1 x^1 + \theta_2 x^2 + \cdots + \theta_n x^n$$

Here, n is the number of features



Hypothesis:

$$h_{\theta}(X) = \theta^T X = \theta_0 x^0 + \theta_1 x^1 + \theta_2 x^2 + \dots + \theta_n x^n$$

Parameters:

$X = (x^0, x^1, \dots, x^n)$, and $\theta = (\theta_0, \theta_1, \dots, \theta_n)$ is $n+1$ dimensional vector

Cost function:

$$J(\theta_0, \theta_1, \dots, \theta_n) = \frac{1}{2N} \sum_{i=1}^N (h_{\theta}(x_i) - y_i)^2$$

$$J(\theta) = J(\theta_0, \theta_1, \dots, \theta_n) = \frac{1}{2N} \sum_{i=1}^N (h_{\theta}(x_i) - y_i)^2$$

$$\min_{\theta} \underset{i}{\text{imize}} \frac{1}{2} \sum_{i=1}^m \left(h_{\theta}(x^{(i)}) - y^{(i)} \right)^2 = \min_{\theta} \underset{i}{\text{imize}} J(\theta)$$

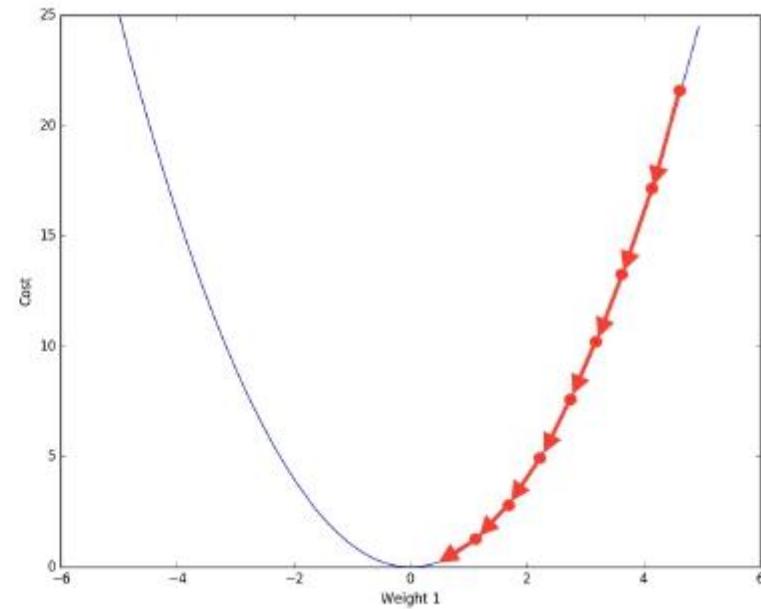
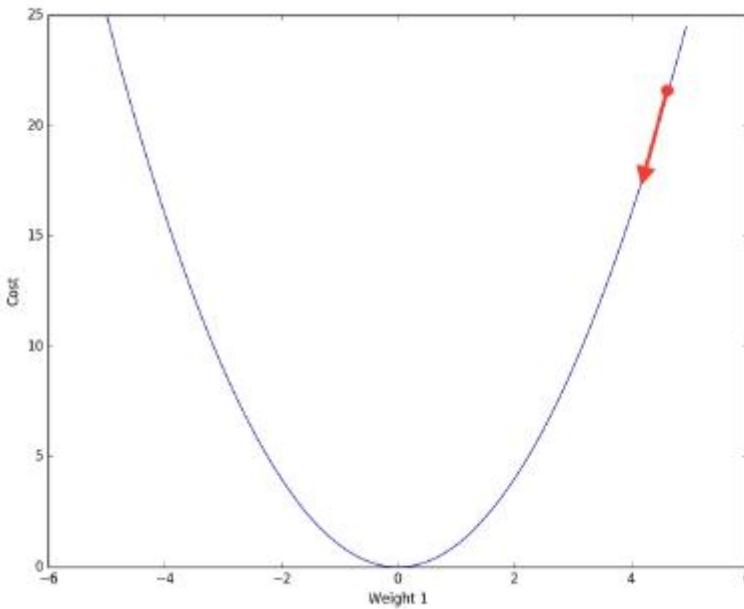


Gradient descent



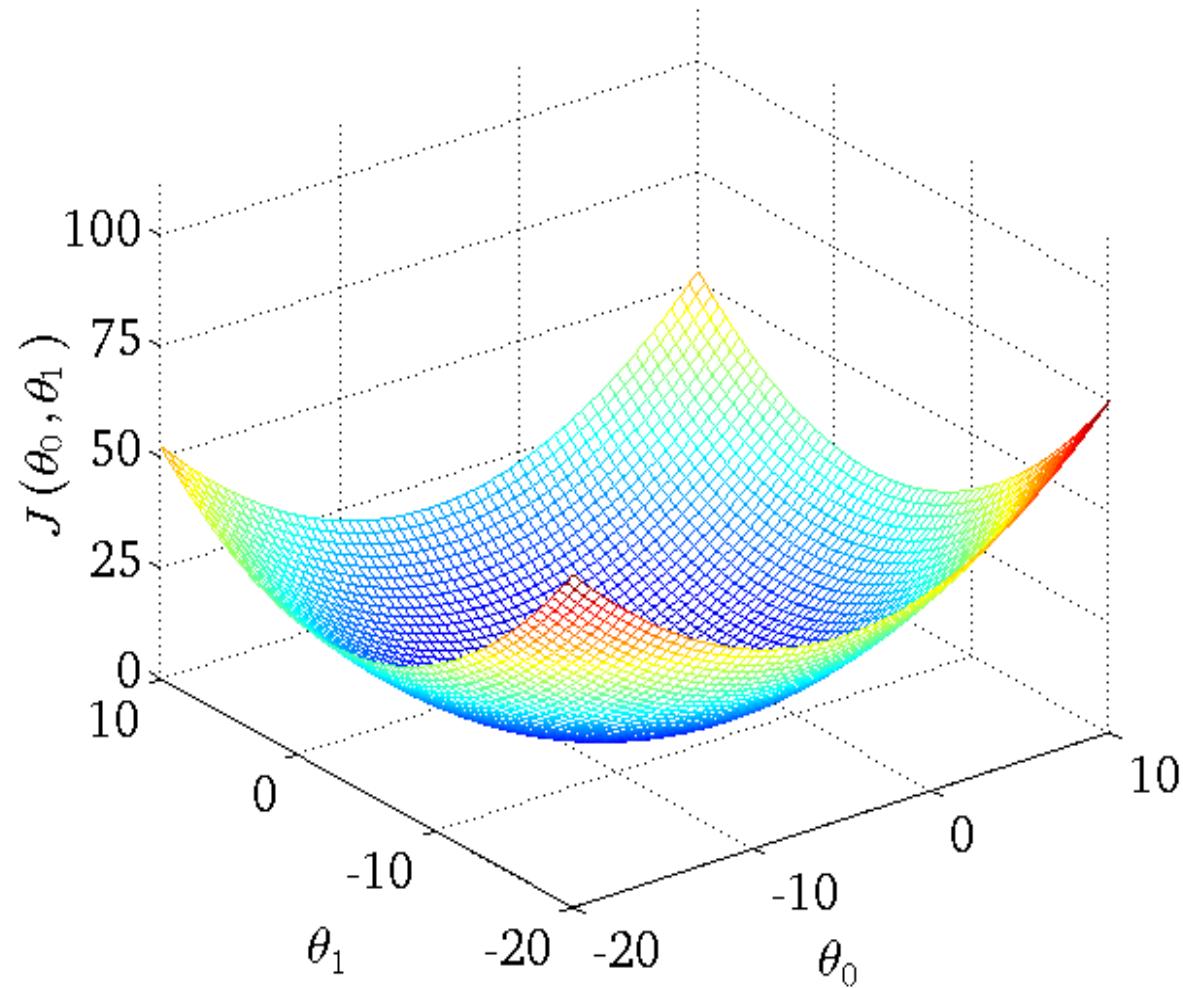
start with some θ , eg. $\theta=0$;
keep changing θ reduce $J(\theta)$;
gradient descent

$$\theta_i = \theta_i - \alpha \frac{\partial}{\partial \theta_i} J(\theta) \quad i = 0, \dots, n$$





Gradient descent

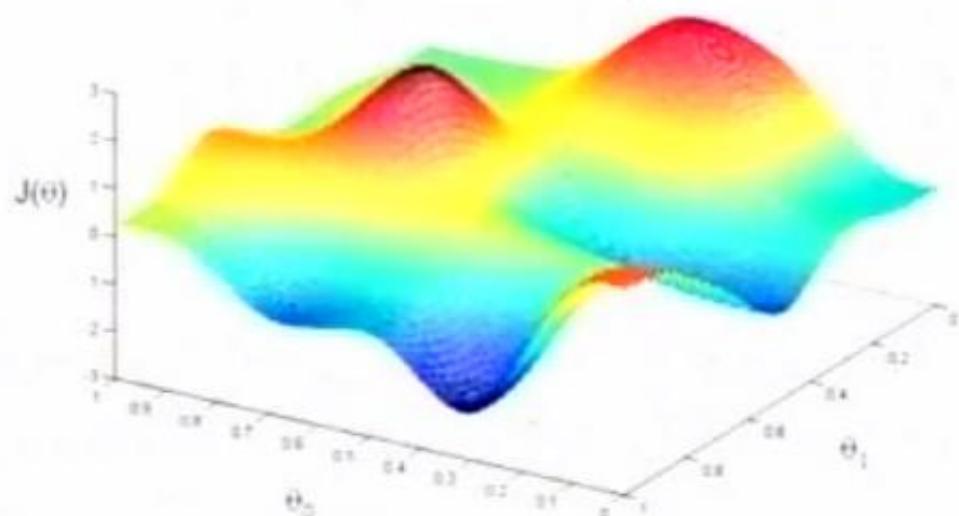


“convex function”

Bowl - shaped



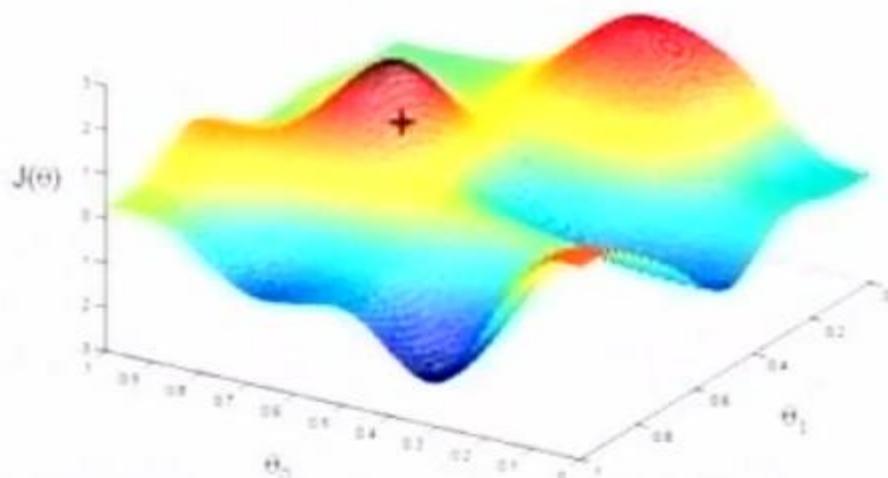
Gradient Descent



STANFOR
UNIVERSITY



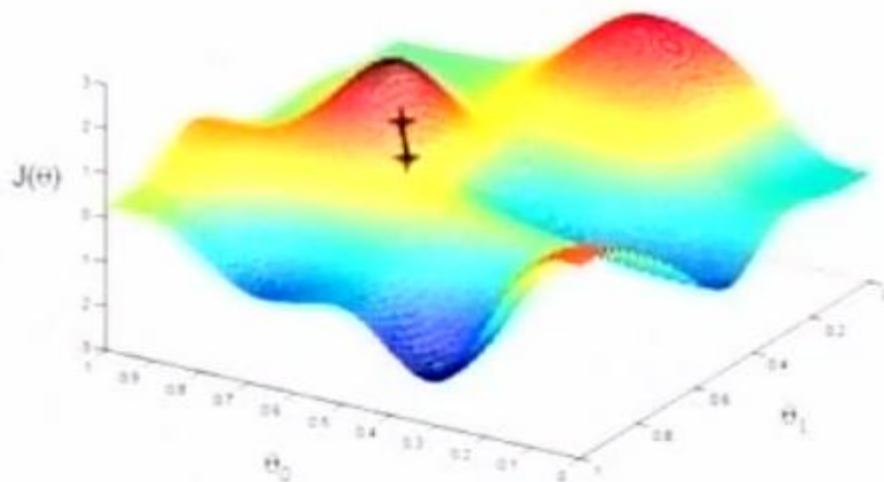
Gradient Descent



STANFORD
LEARN. EXPLORE. DISCOVER.



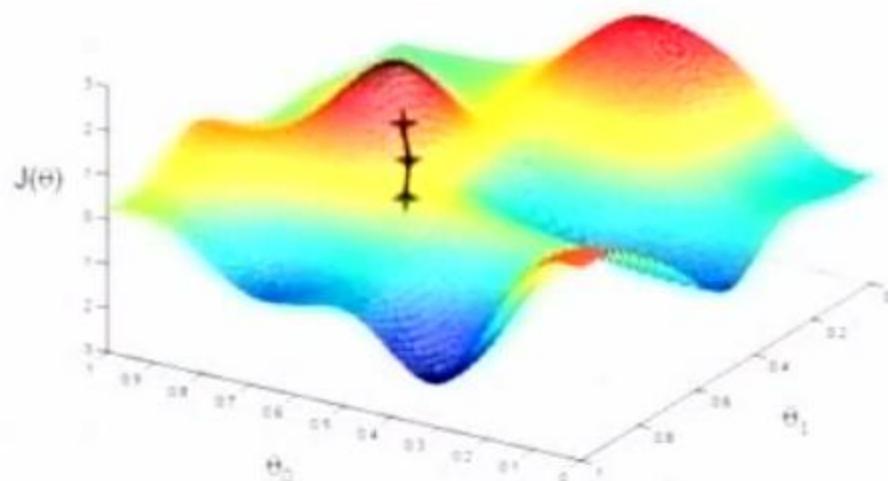
Gradient Descent



STANFORD
UNIVERSITY



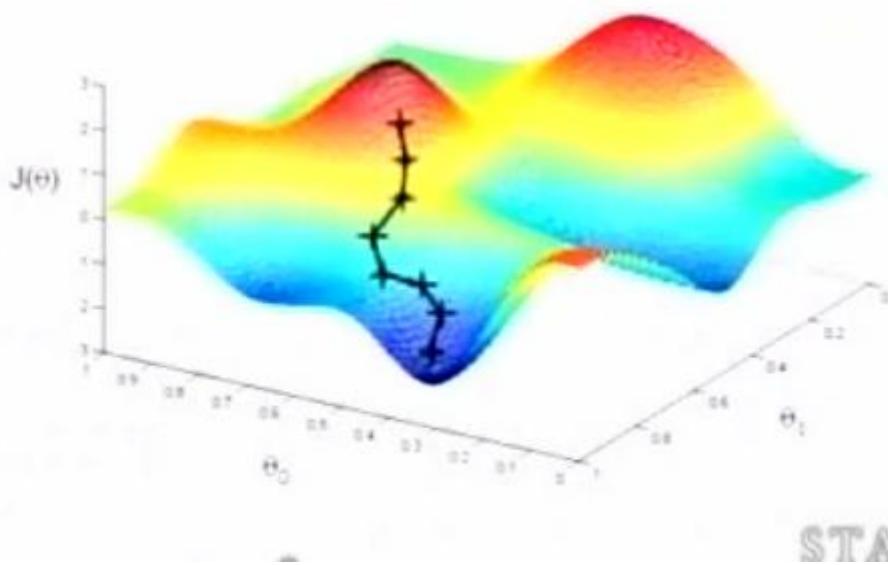
Gradient Descent



STANFOR
UNIVERSITY



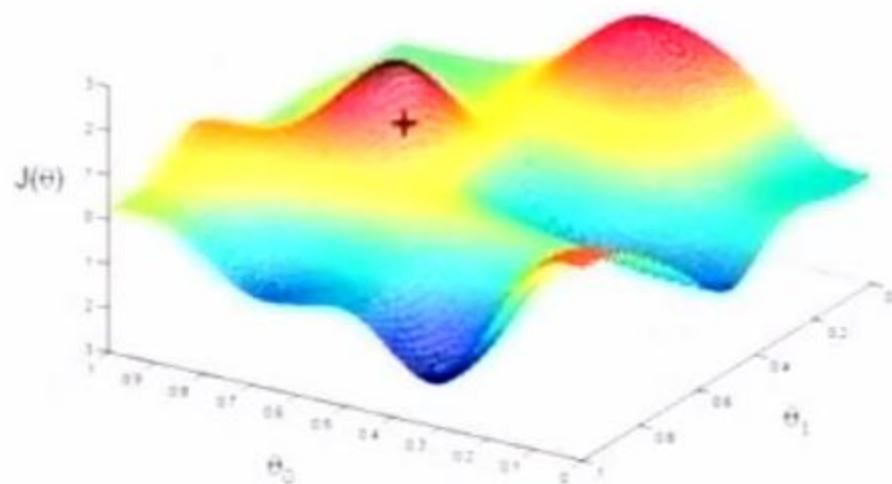
Gradient Descent



STANFOR



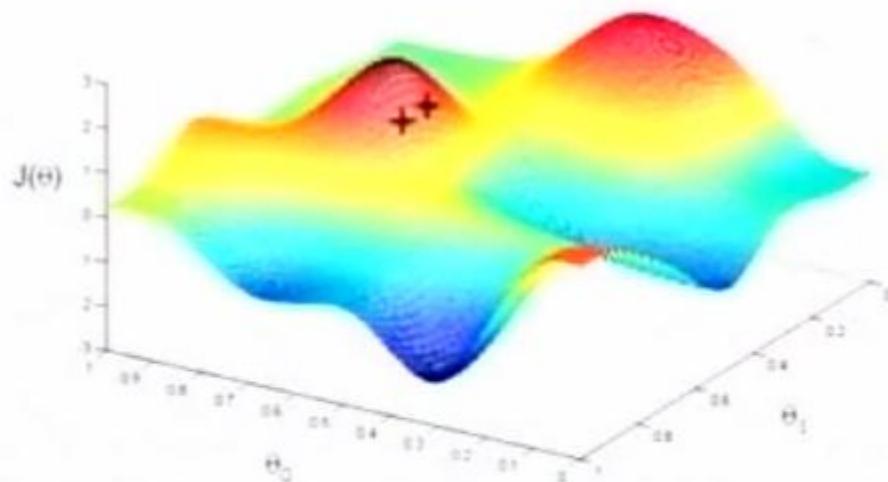
Gradient Descent



STANFORD
UNIVERSITY



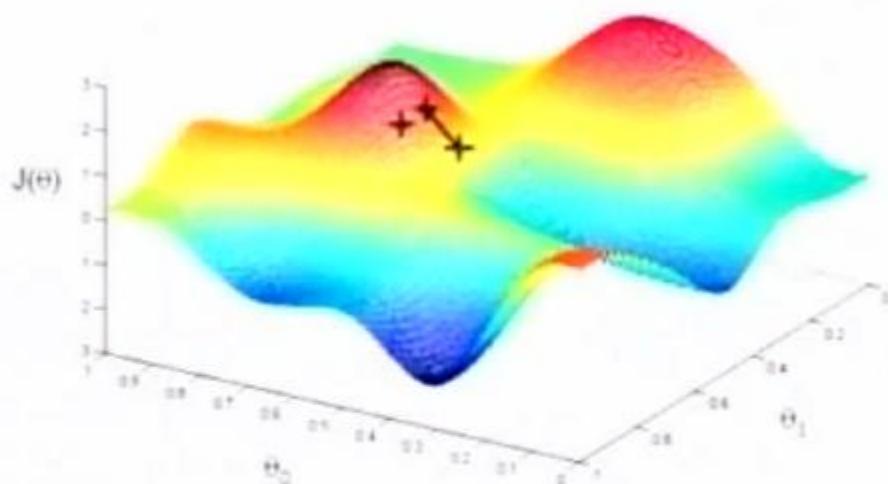
Gradient Descent



STANFORD
UNIVERSITY OF CALIFORNIA



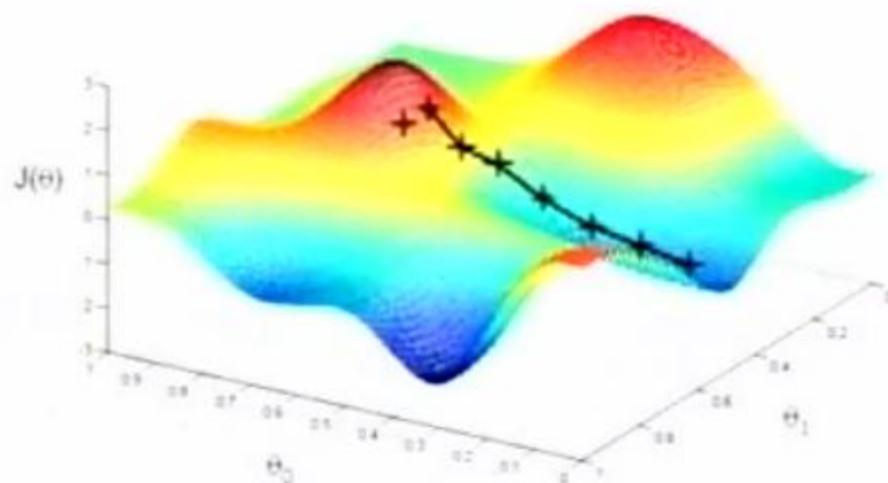
Gradient Descent



STANFORD



Gradient Descent

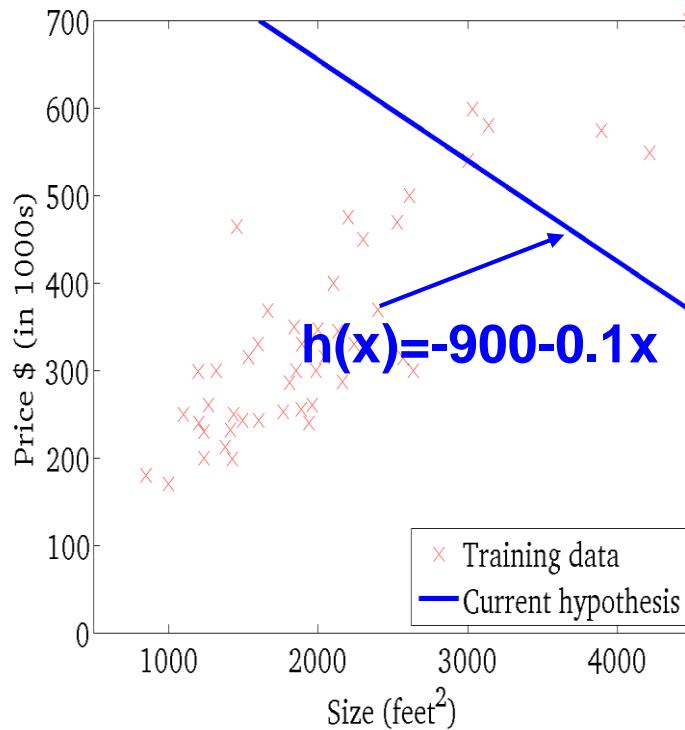


STANFOR
THE STANFORD UNIVERSITY

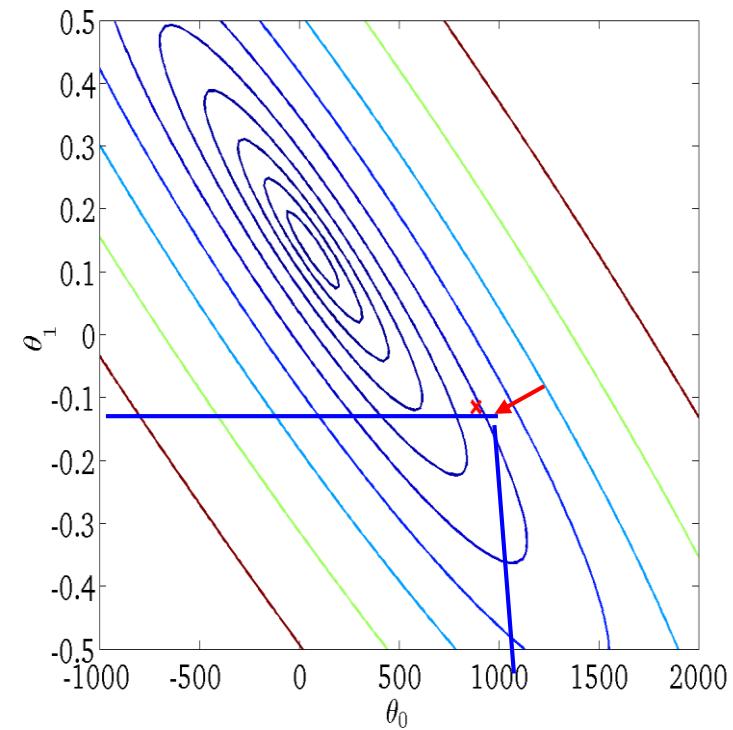


Gradient descent

$h_{\theta}(x)$
For fixed θ_0, θ_1 , this is
a function of x



$J(\theta_0, \theta_1)$
This is a function of the
parameters θ_0, θ_1



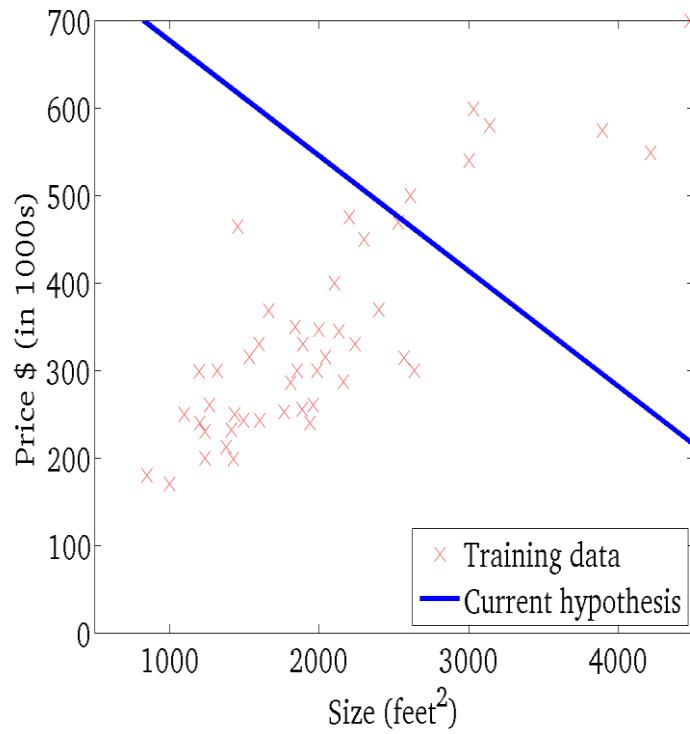


Gradient descent



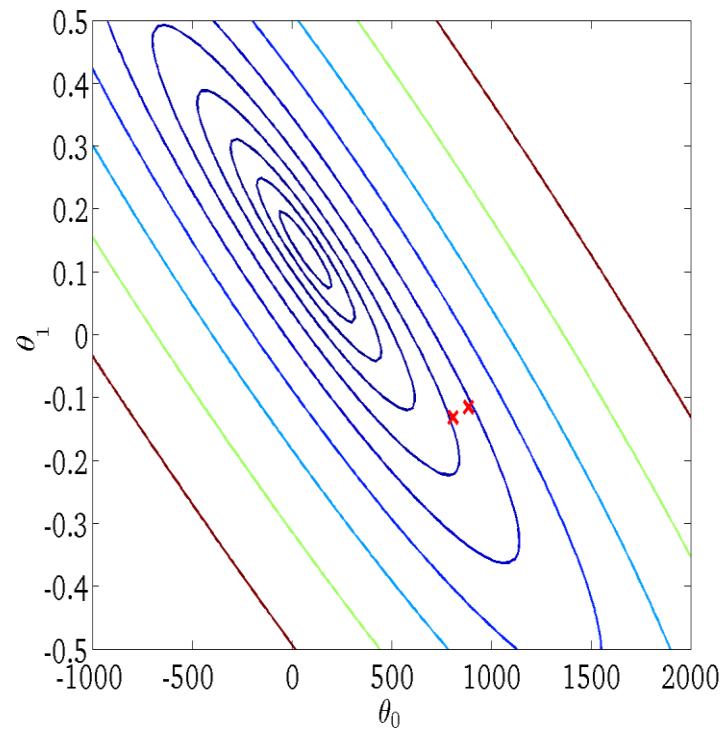
$$h_{\theta}(x)$$

For fixed θ_0, θ_1 , this is
a function of x



$$J(\theta_0, \theta_1)$$

This is a function of the
parameters θ_0, θ_1



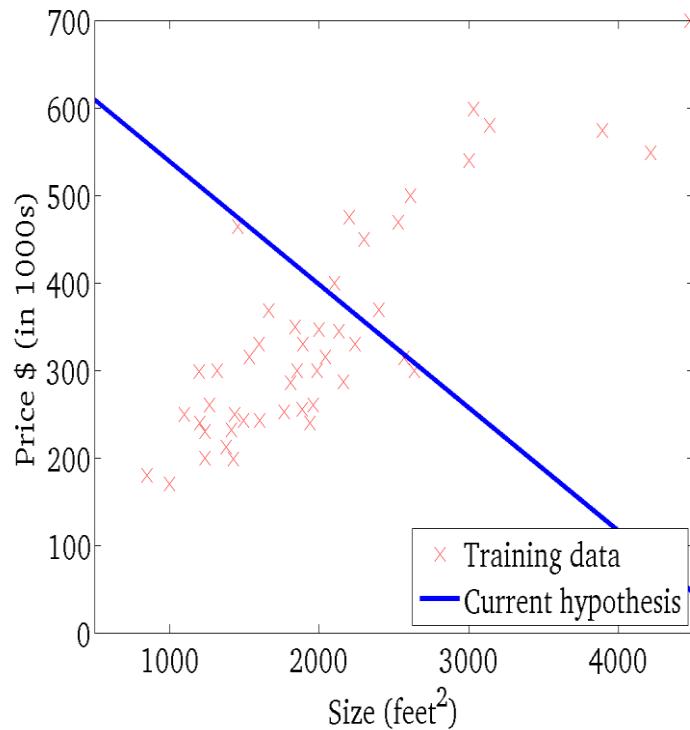


Gradient descent



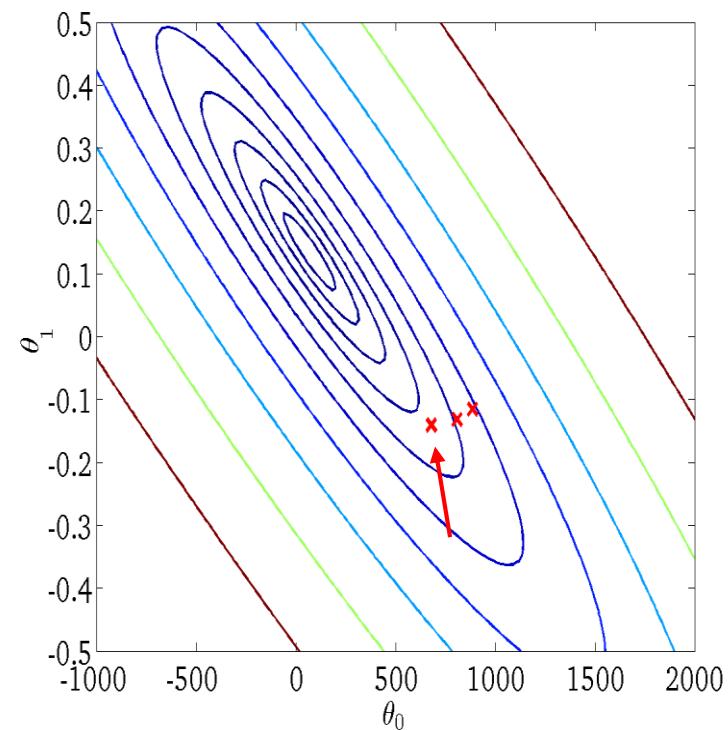
$$h_{\theta}(x)$$

For fixed θ_0, θ_1 , this is
a function of x



$$J(\theta_0, \theta_1)$$

This is a function of the
parameters θ_0, θ_1



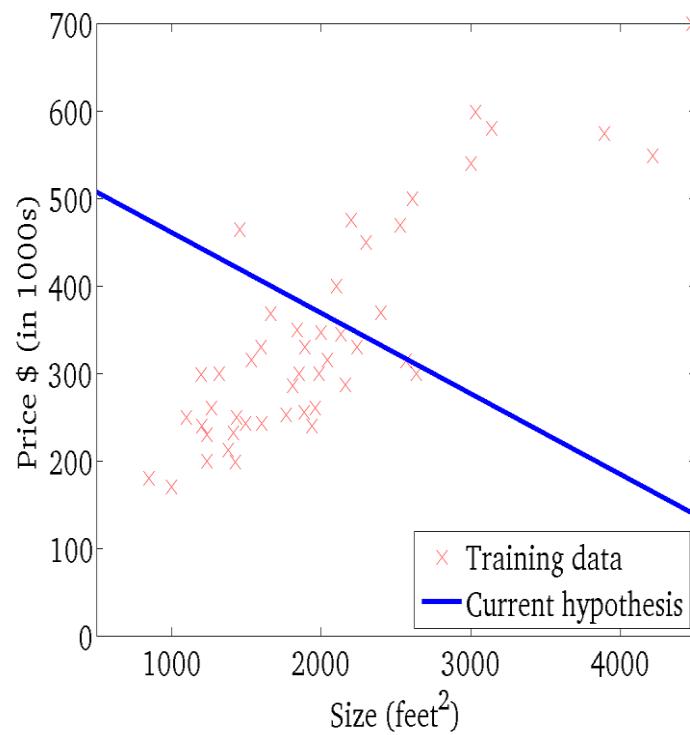


Gradient descent



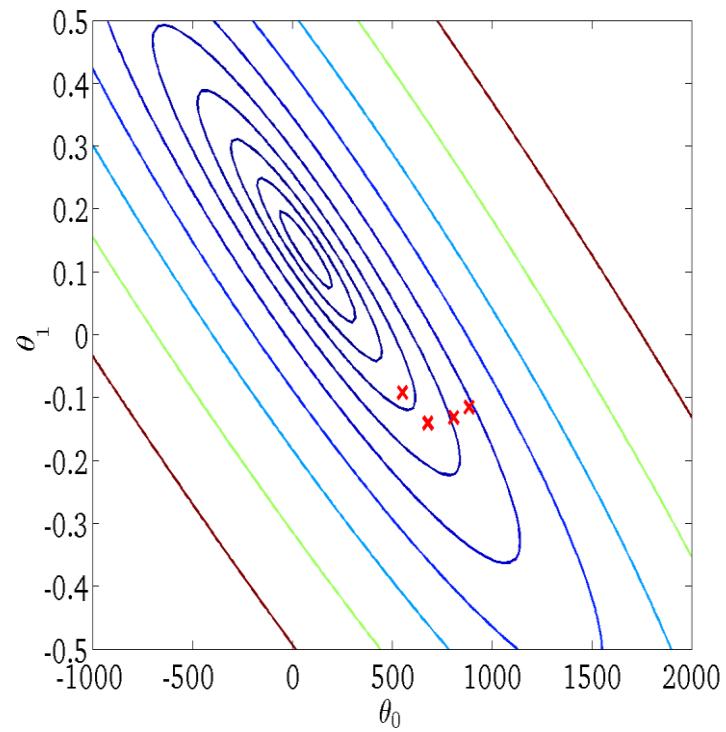
$$h_{\theta}(x)$$

For fixed θ_0, θ_1 , this is
a function of x



$$J(\theta_0, \theta_1)$$

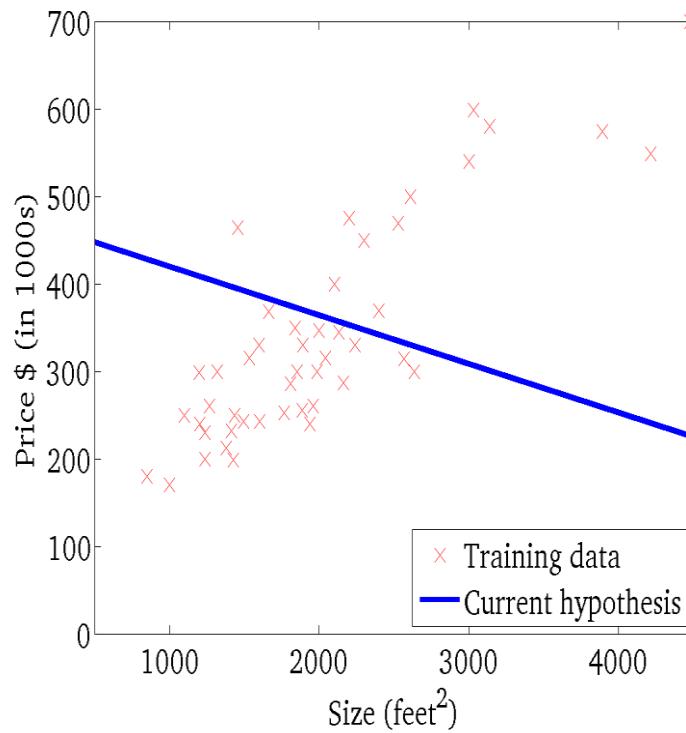
This is a function of the
parameters θ_0, θ_1



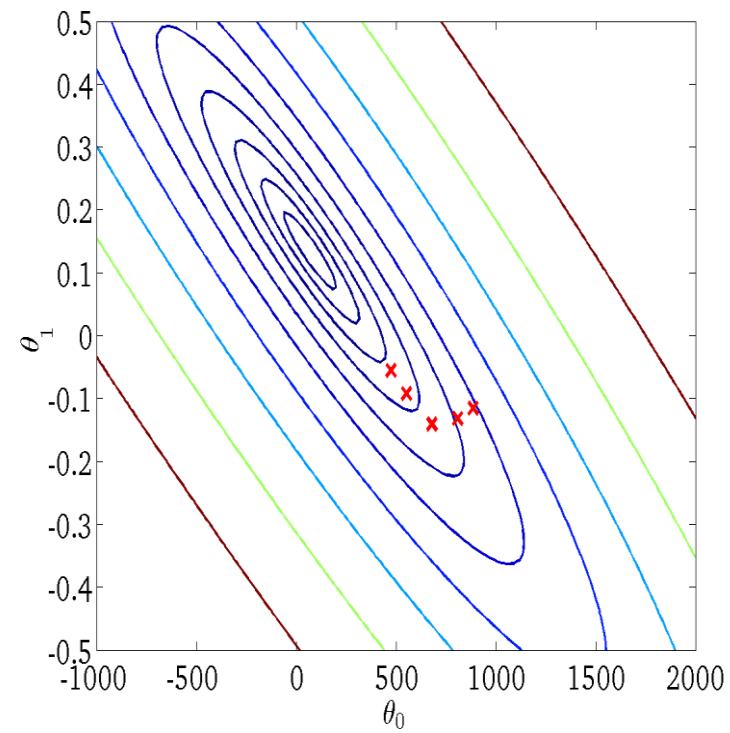


Gradient descent

$h_{\theta}(x)$
For fixed θ_0, θ_1 , this is
a function of x



$J(\theta_0, \theta_1)$
This is a function of the
parameters θ_0, θ_1



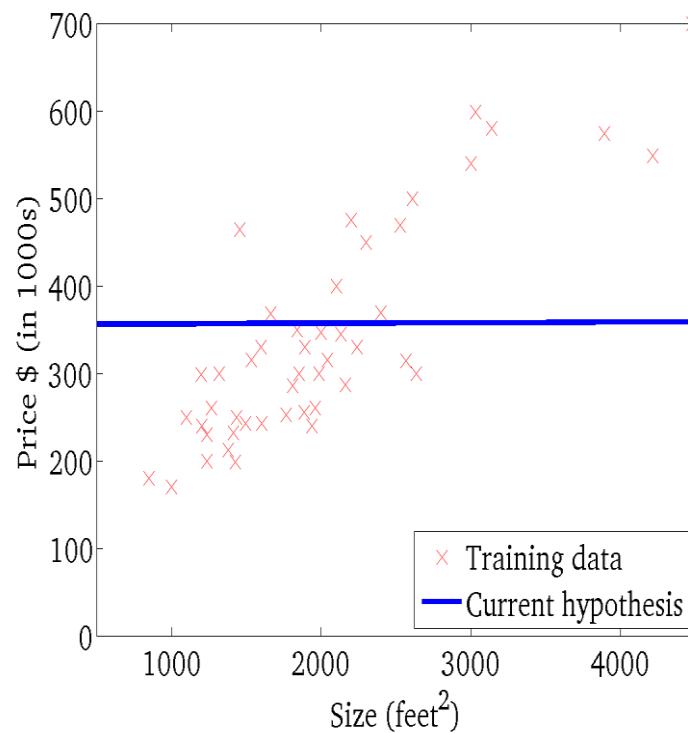


Gradient descent



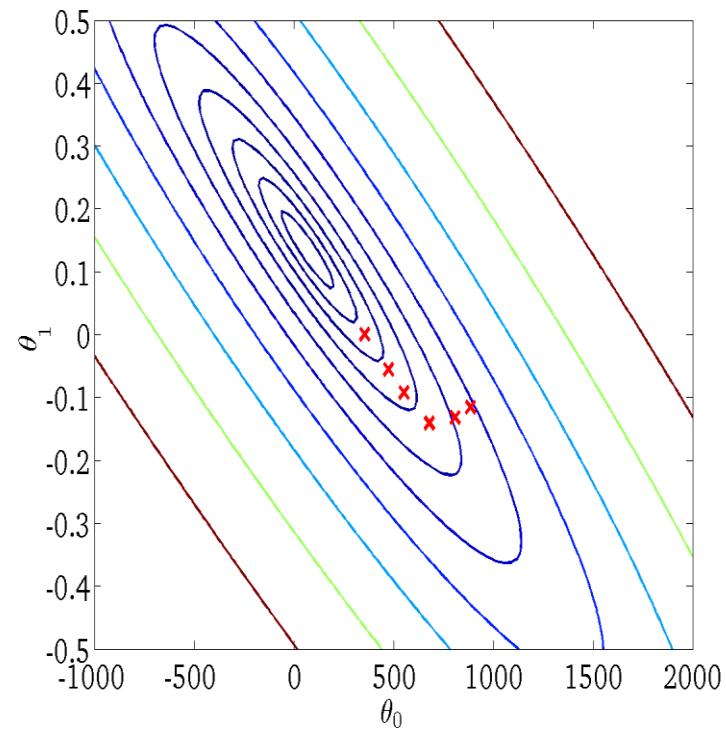
$$h_{\theta}(x)$$

For fixed θ_0, θ_1 , this is
a function of x



$$J(\theta_0, \theta_1)$$

This is a function of the
parameters θ_0, θ_1



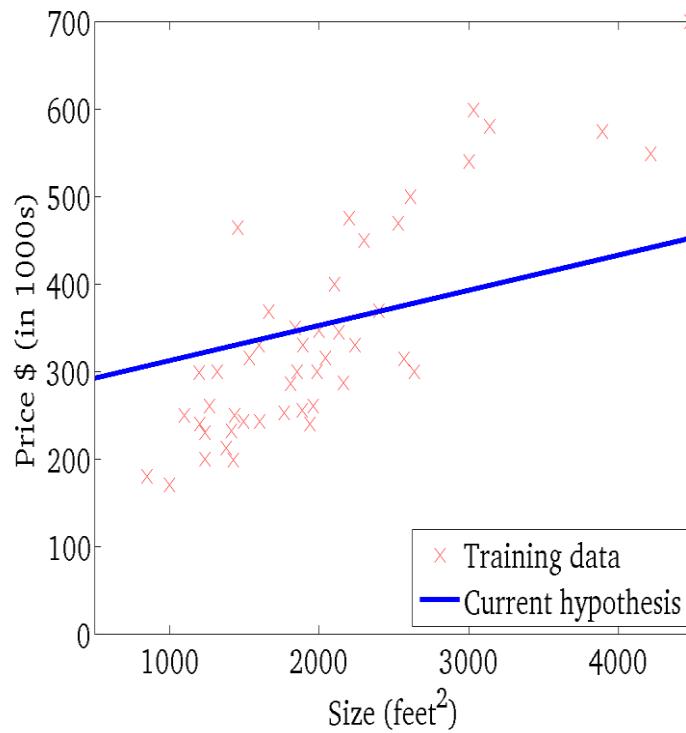


Gradient descent



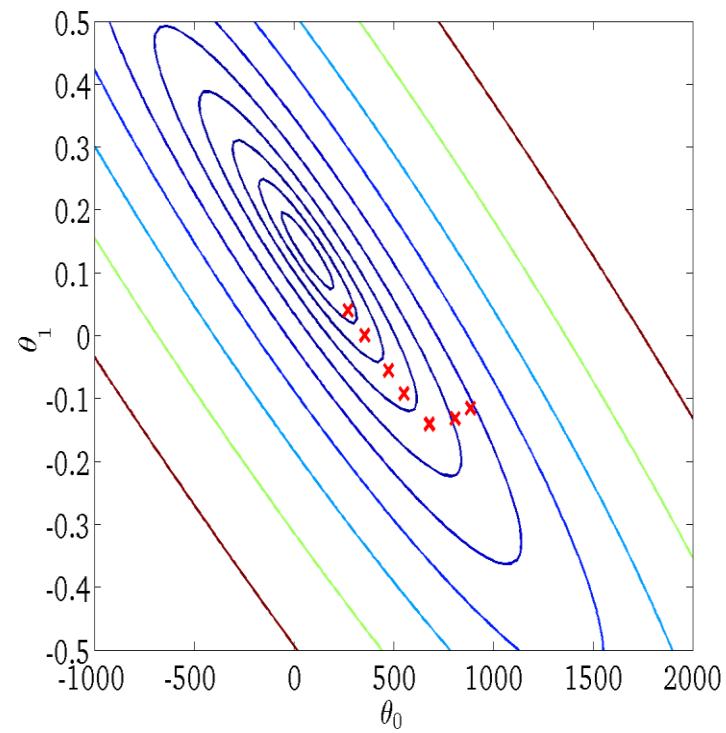
$$h_{\theta}(x)$$

For fixed θ_0, θ_1 , this is
a function of x



$$J(\theta_0, \theta_1)$$

This is a function of the
parameters θ_0, θ_1



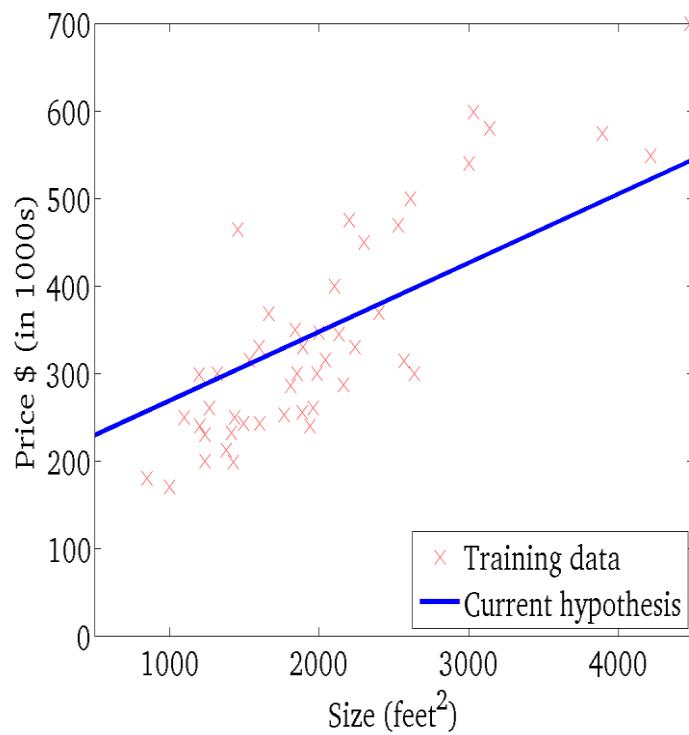


Gradient descent



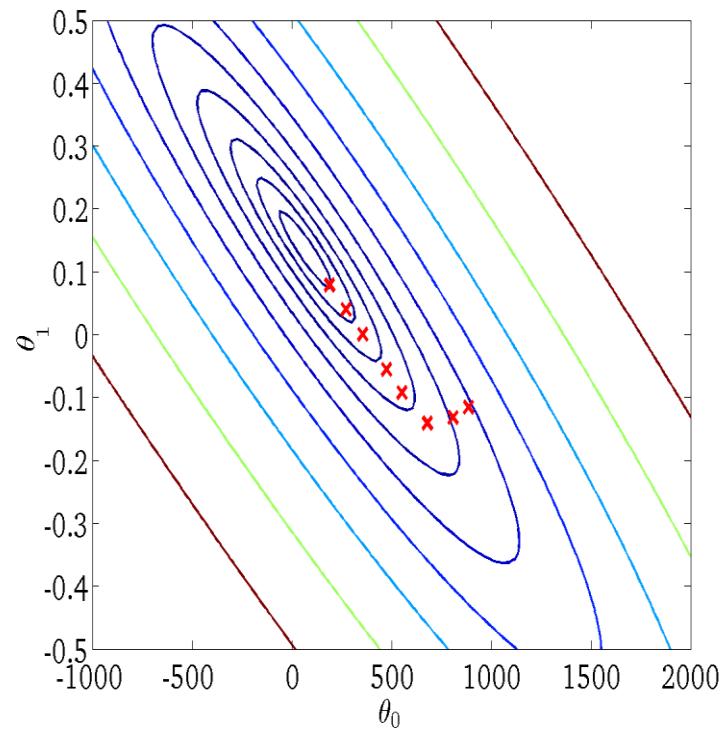
$$h_{\theta}(x)$$

For fixed θ_0, θ_1 , this is
a function of x



$$J(\theta_0, \theta_1)$$

This is a function of the
parameters θ_0, θ_1



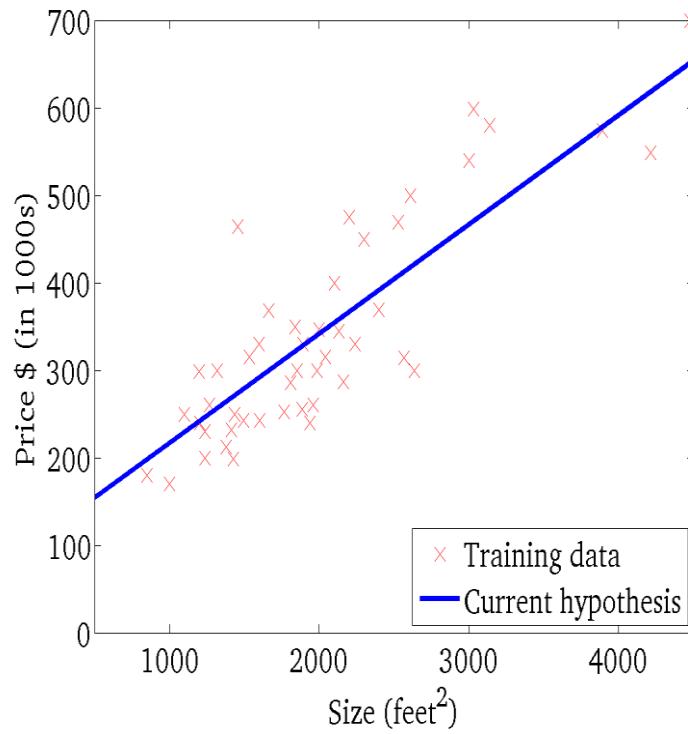


Gradient descent



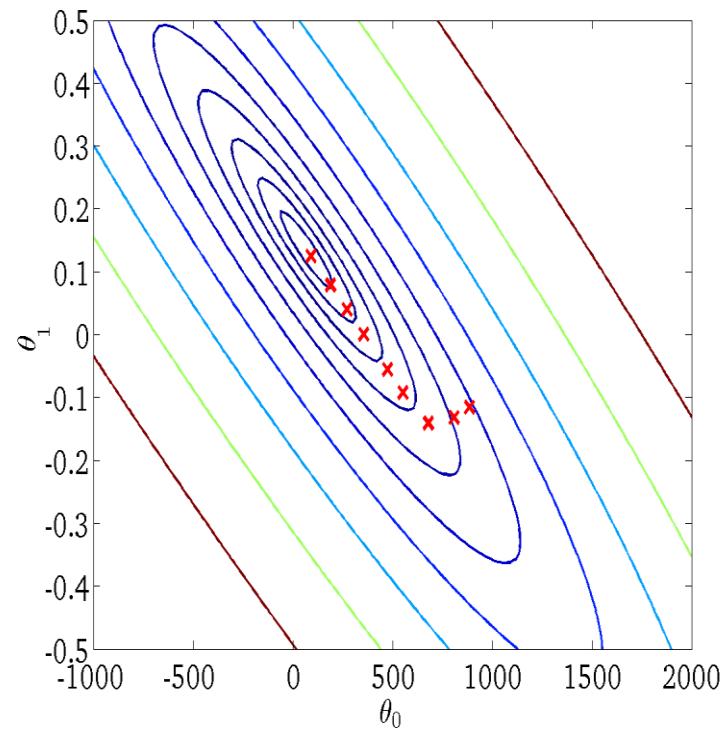
$$h_{\theta}(x)$$

For fixed θ_0, θ_1 , this is
a function of x



$$J(\theta_0, \theta_1)$$

This is a function of the
parameters θ_0, θ_1





Gradient descent



Gradient of Cost function

$$\begin{aligned}\frac{\partial}{\partial \theta_i} J(\theta) &= \frac{\partial}{\partial \theta_i} \frac{1}{2} (h_{\theta}(x) - y)^2 & N = 1 \\ &= 2 \times \frac{1}{2} (h_{\theta}(x) - y) \frac{\partial}{\partial \theta_i} (h_{\theta}(x) - y) \\ &= (h_{\theta}(x) - y) \frac{\partial}{\partial \theta_i} (\theta_0 x^0 + \dots + \theta_i x^i + \dots + \theta_n x^n - y) \\ &= (h_{\theta}(x) - y) x^i\end{aligned}$$

Updated:

$$\theta_i := \theta_i - \alpha (h_{\theta}(x) - y) x^i$$

Repeat until converge



Gradient descent

$$\theta_i := \theta_i - \alpha \sum_{j=1}^N \underbrace{\left(h_{\theta}(x_j) - y_j \right)}_{\frac{\partial}{\partial \theta_i} J(\theta)} x_j^{(i)}$$



$$\nabla_{\theta} J = \begin{bmatrix} \frac{\partial J}{\partial \theta_0} \\ \vdots \\ \frac{\partial J}{\partial \theta_n} \end{bmatrix} \in R^{n+1}$$

gradient descent:

$$\theta := \theta - \alpha \nabla_{\theta} J$$

$\searrow \qquad \qquad \qquad \searrow$

$$R^{n+1} \qquad \qquad \qquad R^{n+1}$$

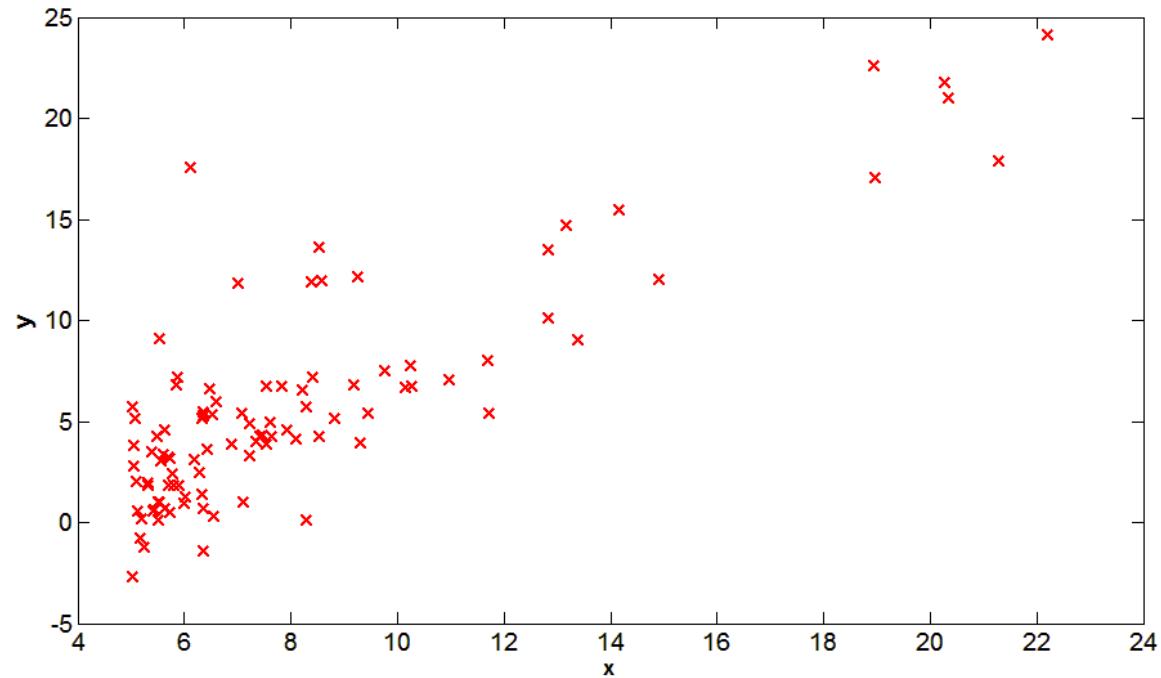


Example1_grad_desent.m

```
1 %梯度下降
2 % 加载数据和数据预处理
3 - clear;
4 - data = load ('data1.txt');
5 % data = load('data3.txt');
6 % data = data';
7 - fprintf('Running Gradient Descent ... \n');
8 %
9 - X = data(:, 1); %取X集合
10 - y = data(:, 2); %取Y集合
11 %%%%%%
12 % data = load('data3.txt');
13 % data = data';
14 % data1=load('data4.txt');
15 % X=data;
16 % X = [ones(length(data), 1), data];
17 % y=data1;
```

data1

1	6.1101, 17.592
2	5.5277, 9.1302
3	8.5186, 13.662
4	7.0032, 11.854
5	5.8598, 6.8233
6	8.3829, 11.886
7	7.4764, 4.3483
8	8.5781, 12
9	6.4862, 6.5987
10	5.0546, 3.8166
11	5.7107, 3.2522
12	14.164, 15.505
13	5.734, 3.1551
14	8.4084, 7.2258
15	5.6407, 0.71618
16	5.3794, 3.5129
17	6.3654, 5.3048
18	5.1301, 0.56077

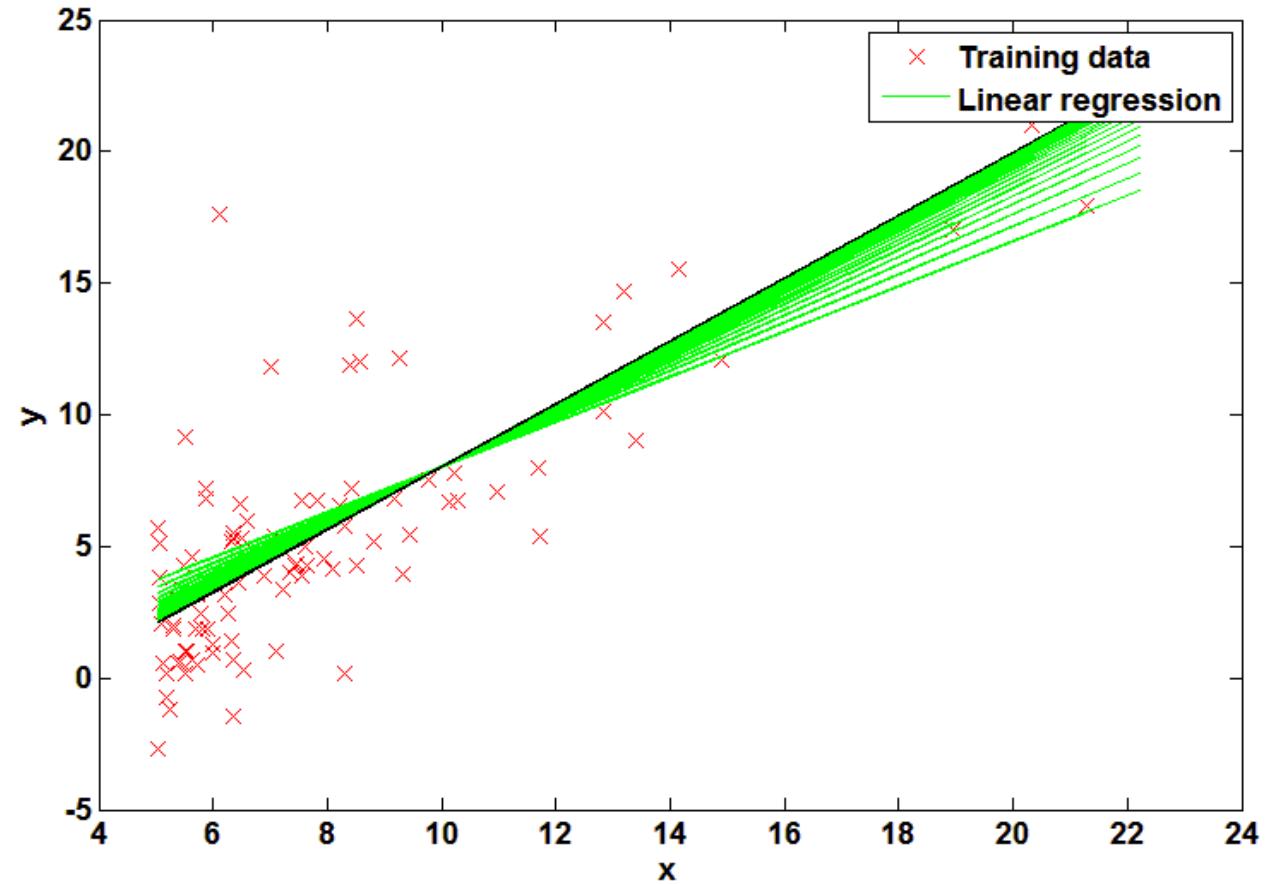




```
32 -     m = length(y); %数据组数
33 -     J_history = zeros(max_iter, 1); %初始化迭代误差变量?
34 -     % iter = 1;
35 -     for iter = 1:max_iter
36 -         %每迭代100次画一条曲线
37 -         if mod(iter, 100) == 0
38 -             plot(X(:, 2), X*theta, 'g')
39 -         end
40 -
41 -         theta = theta - alpha / m * X' * (X * theta - y); %梯度下降
42 -         J_history(iter) = sum((X * theta-y) .^ 2) / (2*m); %记录每次迭代后的全局误差
43 -         fprintf(' iter:%d ----- Error:%f\n', iter, J_history(iter));
44 -     end
45 -     %输出最终的theta值和最终的拟合曲线?disp(' Theta found by gradient descent:');
46 -     disp(theta);
47 -     plot(X(:, 2), X*theta, 'k')
48 -     legend(' Training data', 'Linear regression')
49 -     % hold off
```



```
>> grad_desent
Running Gradient Descent ...
iter:1 ----- Error:6.737190
iter:2 ----- Error:5.931594
iter:3 ----- Error:5.901155
iter:4 ----- Error:5.895229
.....
iter:4998 ----- Error:4.476971
iter:4999 ----- Error:4.476971
iter:5000 ----- Error:4.476971
-3.8953
1.1930
```





Polyfit

polynomial curve fitting

```
clear all;
close all;
data = load ('data1.txt');
x = data(:, 1);           %取x集合
y = data(:, 2);           %取y集合
% x = [ones(length(data),1),data(:,1)];
p = polyfit(x,y,1)%调用matlab自带多项式拟合, 1代表a+bx型
```



Normal equations

For the function of matrix A ,
 $f(A)$ and $A \in \mathbb{R}^{m \times n}$,

$$\nabla_A f(A) = \begin{bmatrix} \frac{\partial f}{\partial A_{11}} & \cdots & \frac{\partial f}{\partial A_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial A_{m1}} & \cdots & \frac{\partial f}{\partial A_{mn}} \end{bmatrix}$$

If $A \in \mathbb{R}^{n \times n}$, the trace of A is $\text{tr}(A) = \sum_{i=1}^n A_{ii}$



Normal equations



In fact,

$$\text{tr } (AB) = \text{tr } (BA)$$

$$\text{tr } (ABC) = \text{tr } (CAB) = \text{tr } (BCA)$$

$$\nabla_A \text{tr } (AB) = B^T$$

$$\text{tr } (A) = \text{tr } (A^T)$$

If $a \in \mathbb{R}$, $\text{tr } (a) = a$

$$\nabla_A \text{tr } (ABA^TC) = CAB + C^TAB^T$$



Normal equations

$$\begin{aligned} tr(AB) &= tr\left(\begin{bmatrix} \leftarrow & a_{1i} & \rightarrow \\ \leftarrow & a_{2i} & \rightarrow \\ \leftarrow & \vdots & \rightarrow \\ \leftarrow & a_{ni} & \rightarrow \end{bmatrix} \begin{bmatrix} \uparrow & \uparrow & \uparrow & \uparrow \\ b_{i1} & b_{i2} & \cdots & b_{in} \\ \downarrow & \downarrow & \downarrow & \downarrow \end{bmatrix}\right) \\ &= \sum_{i=1}^n a_{1i} b_{i1} + \sum_{i=1}^n a_{2i} b_{i2} + \cdots + \sum_{i=1}^n a_{ni} b_{in} \\ \Rightarrow \frac{\partial tr(AB)}{\partial a_{ij}} &= b_{ji} \\ \Rightarrow \nabla_A tr(AB) &= B^T \end{aligned}$$



$$\nabla_A \operatorname{tr}(ABA^T C) = CAB + C^T A B^T$$



$\operatorname{tr}(ABA^T C)$ 对 A 求导， A 出现过两次，相当于函数乘积求导。先固定 $BA^T C$ ，然后对第一个 A 求导；固定 AB 和 C ，对第二个 A 求导（注意，在固定 AB 和 C 时，先利用 $\operatorname{tr}(ABC) = \operatorname{tr}(BCA) = \operatorname{tr}(CAB)$ 这一结论，将 $\operatorname{tr}(ABA^T C)$ 等价变形为 $\operatorname{tr}(CABA^T)$ ；最后将两部分结果相加。具体过程如下：

$$\begin{aligned}\nabla_A \operatorname{tr}(ABA^T C) &= \nabla_A \operatorname{tr}[A(BA^T C)] + \nabla_A \operatorname{tr}[(CAB)A^T] \\&= (BA^T C)^T + \nabla_A \operatorname{tr}[(CAB)A^T]^T \quad (\text{注：利用 } \nabla_A \operatorname{tr}(AB) = B^T \text{ 和 } \operatorname{tr}(A) = \operatorname{tr}(A^T)) \\&= C^T AB^T + \nabla_A \operatorname{tr}[A(CAB)^T] \\&= C^T AB^T + ((CAB)^T)^T \\&= C^T AB^T + CAB\end{aligned}$$



Normal equations

$$X = \begin{bmatrix} 1 & (x_1)^T & - \\ 1 & (x_2)^T & - \\ \vdots & \vdots & - \\ 1 & (x_N)^T & - \end{bmatrix}$$

$$X\theta = \begin{bmatrix} 1 & (x_1)^T & - \\ 1 & (x_2)^T & - \\ \vdots & \vdots & - \\ 1 & (x_N)^T & - \end{bmatrix} \theta$$

$$X \in N \times (n+1)$$

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$

$$X\theta = \begin{bmatrix} x_1^T \theta \\ \vdots \\ x_N^T \theta \end{bmatrix} = \begin{bmatrix} h_\theta(x_1) \\ \vdots \\ h_\theta(x_N) \end{bmatrix}$$



Normal equations

$$X\theta - y = \begin{bmatrix} h_\theta(x_1) - y_1 \\ \vdots \\ h_\theta(x_N) - y_N \end{bmatrix}$$

$$\begin{aligned} & \frac{1}{2}(X\theta - y)^T(X\theta - y) \\ &= \frac{1}{2} \sum_{i=1}^N (h_\theta(x_i) - y_i)^2 \\ &= J(\theta) \end{aligned}$$



Normal equations



$$\nabla_{\theta} J(\theta) \stackrel{\text{set}}{=} 0$$

$$Z^T Z = \sum_i z_i^2 = \text{tr}(Z^T Z)$$

$$\nabla_{\theta} \frac{1}{2} (X\theta - y)^T (X\theta - y)$$

$$\text{tr}(ABC) = \text{tr}(CAB) = \text{tr}(BCA)$$

$$= \frac{1}{2} \nabla_{\theta} \text{tr}[(X\theta - y)^T (X\theta - y)]$$

$$-\nabla_{\theta} \text{tr}(\theta^T X_1^T y) = -\nabla_{\theta} \text{tr}(\theta^T X_1^T y)^T$$

$$= \frac{1}{2} \nabla_{\theta} \text{tr}[(\theta^T X^T - y^T)(X\theta - y)]$$

$$\text{tr}(A) = \text{tr}(A^T)$$

$$= \frac{1}{2} \nabla_{\theta} \text{tr}(\theta^T X^T X \theta - \theta^T X^T y - y^T X \theta + y^T y)$$

$$= \frac{1}{2} [\nabla_{\theta} \text{tr}(\theta \theta^T X^T X) - \nabla_{\theta} \text{tr}(y^T X \theta) - \nabla_{\theta} \text{tr}(y^T X \theta)]$$



Normal equations

$$\nabla_{\theta} \text{tr}(\theta \theta^T X^T X) \quad \nabla_A \text{tr}(ABA^T C) = CAB + C^T AB^T$$

$$\begin{aligned}\nabla_{\theta} \text{tr} \underbrace{\theta I \theta^T}_{\mathbf{A} \mathbf{B} \mathbf{A}^T} \underbrace{X^T X}_{\mathbf{C}} &= \underbrace{X^T X \theta I}_{\mathbf{C}} + \underbrace{X^T X \theta I}_{\mathbf{C}^T \mathbf{A} \mathbf{B}^T} \\ &= X^T X \theta + X^T X \theta\end{aligned}$$

$$trAB = trBA \longrightarrow \nabla_{\theta} \text{tr}(y^T X \theta) = \nabla_{\theta} \text{tr}(\theta y^T X)$$

$$\nabla_A \text{tr} AB = B^T \longrightarrow = (y^T X)^T = X^T y$$



Normal Equation

$$\nabla_{\theta} J(\theta) = \frac{1}{2} [X^T X \theta + X^T X \theta - X^T y - X^T y]$$

$$= X^T X \theta - X^T y \stackrel{\text{set}}{=} 0$$

$$X^T X \theta = X^T y$$

$$\theta = (X^T X)^{-1} X^T y$$



Normal Equation

```
% normal equation

clear all;
close all;
data = load ('data1.txt');
x = data(:, 1);           %取x集合
y = data(:, 2);           %取y集合
x = [ones(length(data),1),data(:,1)];
theta=zeros(2,1);
x_t=x';%矩阵转置
x1=x_t*x;
x2=inv(x1);
theta=x2*x_t*y
```



Linear Regression Extension



Example3: One-variable/multivariable linear regression

(注：该示例为血压与年龄、体重指数、吸烟习惯等的线性回归模型)

Number	Blood pressure (y)	Age (x ¹)	body weight index (x ²)	Smoking habit (x ³)	Number	Blood pressure (y)	Age (x ¹)	body weight index (x ²)	Smoking habit (x ³)
1	144	39	24.2	0	16	130	48	22.2	1
2	215	47	31.1	1	17	135	45	27.4	0
3	138	45	22.6	0	18	114	18	18.8	0
...
15	128	42	21.7	0	30	175	69	27.4	1



(1) One-variable linear regression:(注: 该一元线性回归分析为血压与体重指数的模型, 下面公式中的 X^2 指的是体重指数)

$$h_{\theta}(x) = \theta_0 + \theta_2 x^2$$

1) gradient descent:(梯度下降法实现一元线性回归及程序)

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

```
1 function J = computeCost(X, y, theta)
2
3     J = 0;
4     predictions = X * theta;
5     J = 1/(2*m)*(predictions - y)'*(predictions - y);
6 end
```



```
1 function [theta, J_history]= gradientDescent(X,y,theta,alpha, num_iters)
2 m = length(y); %number of training examples
3 J_history = zeros(num_iters,1);
4 for iter = 1:num_iters
5     temp(1) = theta(1) - (alpha / m)* sum((X * theta - y).*X(:,1));
6     temp(2) = theta(2) - (alpha / m)* sum((X * theta - y).*X(:,2));
7     theta(1)=temp(1);
8     theta(2)=temp(2);
9     J_history(iter)=computerCost(X, y, theta);
10
11 end
```

$$\theta_i := \theta_i - \alpha \sum_{j=1}^m (h_\theta(x_j) - y_j) \cdot x_j^{(i)}$$



$$y = -3.895301 + 1.192985x^2$$

```
1 -      clc
2 -      clear all
3 -      close all
4 -      data = load ('data1.txt');
5 -      X = data(:, 1);
6 -      y = data(:, 2);
7 -      X = [ones(size(X, 1), 1), X];
8 -      alpha=0.01;
9 -      num_iters=400;
10 -     theta=zeros(2, 1);
11 -     [theta, J_history]= gradientDescent(X,y,theta,alpha, num_iters);
12 -     fprintf('Theta compute from gradient descent:\n');
13 -     fprintf('%f\n',theta);
14 -     fprintf('\n');
```

Theta compute from gradient descent:
-1.963957
0.998961 **num_iters=400**

Theta compute from gradient descent:
-3.895301
1.192985 **num_iters=5000**



2) normal equation:(注: 梯度下降法实现一元线性回归及程序)

$$\theta = (X^T X)^{-1} X^T y$$

```
1 -> function [theta] = normalEqn(X, y)
2 -> theta = zeros(size(X, 2), 1);
3 ->
4 -> theta = pinv( (X' * X) ) * X' * y;
5 -> end
```



```
1 -      clc
2 -      clear all
3 -      close all
4 -      data = csvread('data1.txt');
5 -      X = data(:, 1);
6 -      y = data(:, 2);
7 -      m=length(y);
8 -      X=[ones(m, 1) X];
9 -      [theta]=normalEqn(X, y);
10 -     fprintf(' Theta compute from the normal equations:\n');
11 -     fprintf('%f\n', theta);
12 -     fprintf('\n');
```

Theta compute from the normal equations:

-3.895781

1.193034

$$y = -3.895781 + 1.193034x^2$$



regress

Multiple linear regression



(2) **Multiple Linear Regression:** (注：使用梯度法和正则方程法实现多元线性回归分析时，程序代码与一元线性回归相同，只需调主程序第五、六行为`X = data(:,1:3); y = data(:,4)`即可。)

$$y = \theta_0 + \theta_1 x^1 + \theta_2 x^2 + \theta_3 x^3$$

```
1 -      clc
2 -      clear all
3 -      close all
4 -      x1=[39 47 45 47 65 46 67 42 67 56 64 56 59 34 42 48 45 18 20 19 36 50 39 21 44 53 63 29 25 69];
5 -      x2=[24.2 31.1 22.6 24.0 25.9 25.1 29.5 19.7 27.2 19.3 28.0 25.8 27.3 20.1 21.7 22.2 27.4 18.8 22.6...
6 -      21.5 25.0 26.2 23.5 20.3 27.1 28.6 28.3 22.0 25.3 27.4];
7 -      x3=[0 1 0 1 1 0 1 0 1 0 0 0 0 1 0 0 0 0 1 0 0 1 1 0 1 0 1];
8 -      y=[144 215 138 145 162 142 170 124 158 154 162 150 140 110 128 130 135 114 116 124 136 142 120 120 160 15
9 -      X=[ones(30,1),x1',x2',x3'];
10 -     [b,bint,r,rint,stats]=regress(y',X);
11 -     b, bint, stats
```



b =

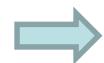
45.3636
0.3604
3.0906
11.8246

bint =

3.5537	87.1736
-0.0758	0.7965
1.0530	5.1281
-0.1482	23.7973



regression coefficient	Estimates of regression coefficient	confidence interval of regression coefficient
θ_0	45.3636	[3.5537 87.1736]
θ_1	0.3604	[-0.0758 0.7965]
θ_2	3.0906	[1.0530 5.1281]
θ_3	11.8246	[-0.1482 23.7973]



$$y = 45.3636 + 0.3604x^1 + 3.0906x^2 + 11.8246x^3$$



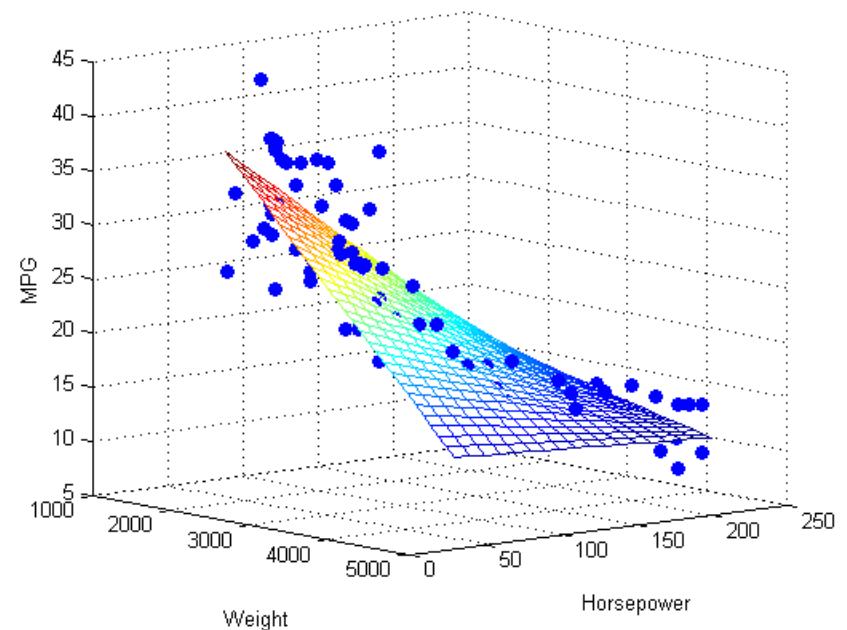
Example4: Load the sample data. Identify weight and horsepower as predictors, and mileage as the response.

```
load carsmall  
x1 = Weight;  
x2 = Horsepower;      % Contains NaN data  
y = MPG;
```

```
X = [ones(size(x1)) x1 x2 x1.*x2];  
b = regress(y,X)      % Removes NaN data
```



```
scatter3(x1,x2,y,'filled')
hold on
x1fit = min(x1):100:max(x1);
x2fit = min(x2):10:max(x2);
[X1FIT,X2FIT] = meshgrid(x1fit,x2fit);
YFIT = b(1) + b(2)*X1FIT + b(3)*X2FIT + b(4)*X1FIT.*X2FIT;
mesh(X1FIT,X2FIT,YFIT)
xlabel('Weight')
ylabel('Horsepower')
zlabel('MPG')
view(50,10)
```





Improved gradient descend:

$$\begin{cases} \theta^{(0)} \\ \theta^{(m+1)} = \theta^{(m)} - \eta^{(m)} \left(\frac{\partial J^{(m)}}{\partial \theta_i^m} \right) + \alpha^{(m)} (\theta^{(m)} - \theta^{(m-1)}) \\ \eta^{(m+1)} = \begin{cases} \eta^{(m)} \phi & J^{(m+1)} - J^{(m)} \leq 0 \\ \eta^{(m)} \beta & J^{(m+1)} - J^{(m)} > 0 \end{cases} \\ \alpha^{(m+1)} = \begin{cases} \alpha_0 & J^{(m+1)} - J^{(m)} \leq 0 \\ 0 & J^{(m+1)} - J^{(m)} > 0 \end{cases} \\ i, m = 0, 1, \dots \end{cases}$$

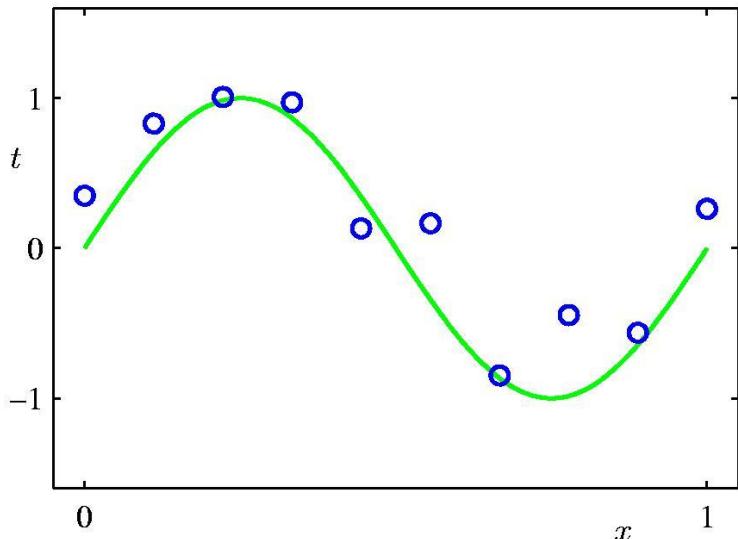


In this equation, the parameter η and α are the step length and the inertia coefficient of the iterative process respectively, which are changing with the changing number of iteration.

Besides, the parameter Φ and β are the weight coefficient for η , used in the iterative process. And $\Phi > 1$, $0 < \beta < 1$.



Overfitting underfitting



- Training set: x, t
- Goal: exploit this training set in order to make predictions of the value \hat{t} of the target variable for some new value \hat{x}
- Method: find a function $t = f(x)$

fit the data using a polynomial function

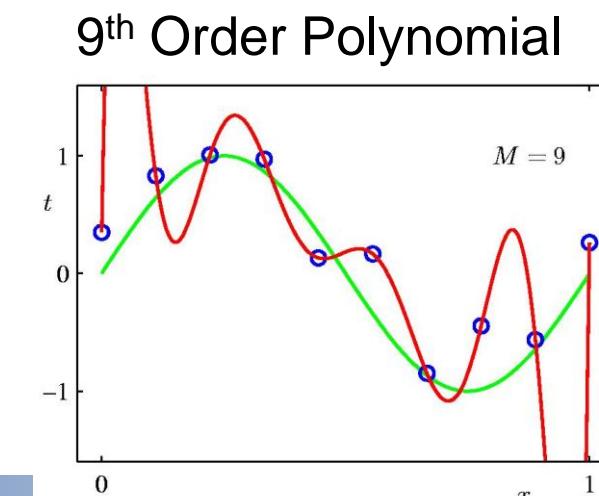
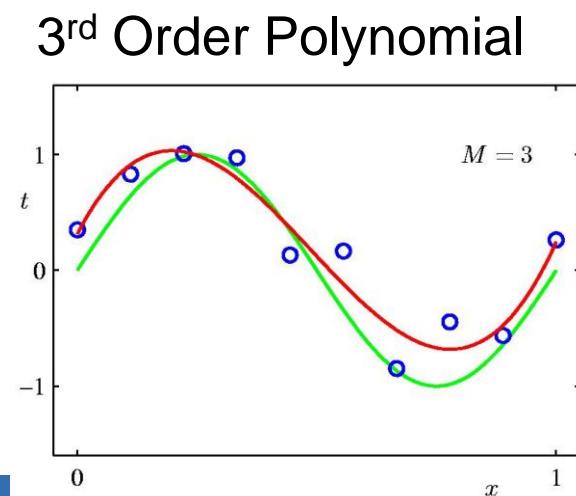
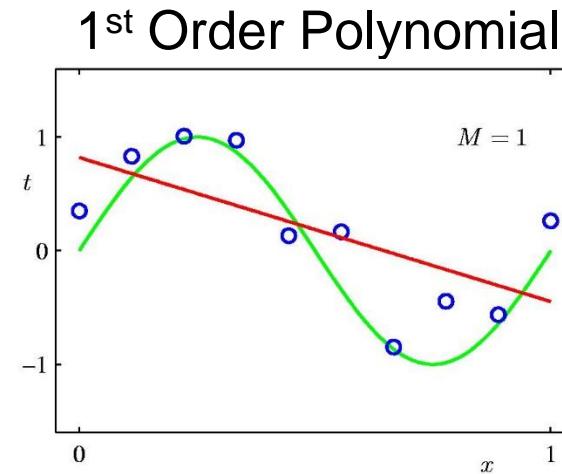
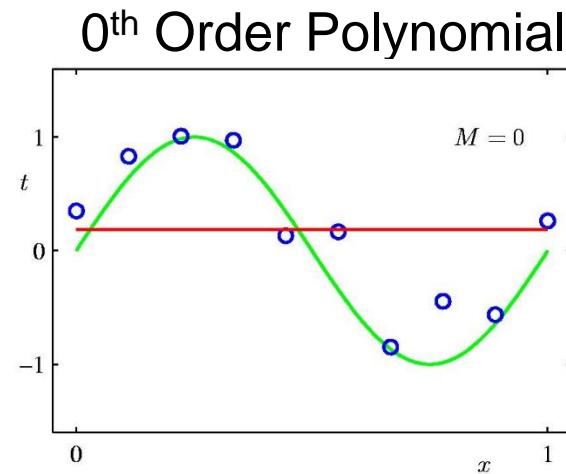
$$y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \dots + w_M x^M = \sum_{j=0}^M w_j x^j$$

Linear model: although the polynomial function $y(x, \mathbf{w})$ is a nonlinear function of x , it is a linear function of the coefficients w .



Overfitting underfitting

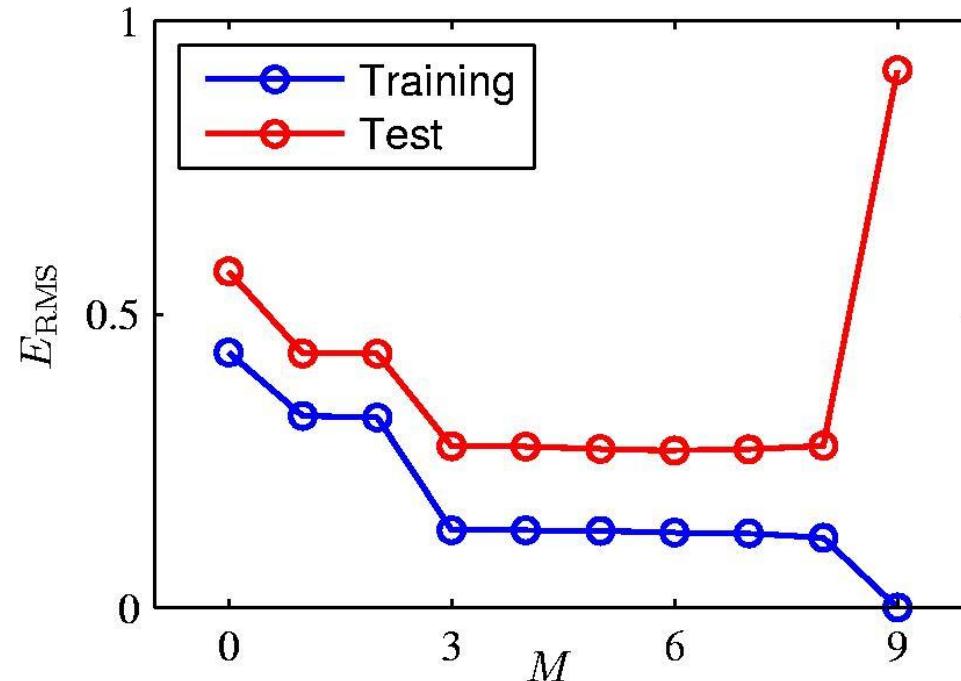
- Problem: choosing the order M of the polynomial





Over-fitting

Goal: is to achieve good generalization by making accurate predictions for new data.



Root-Mean-Square (RMS) Error: $E_{\text{RMS}} = \sqrt{2E(\mathbf{w}^*)/N}$



Polynomial Coefficients

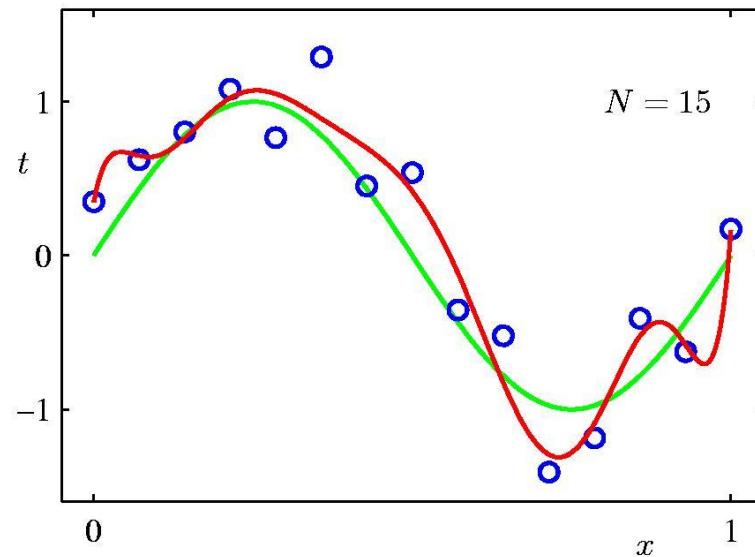
	$M = 0$	$M = 1$	$M = 3$	$M = 9$
w_0^*	0.19	0.82	0.31	0.35
w_1^*		-1.27	7.99	232.37
w_2^*			-25.43	-5321.83
w_3^*			17.37	48568.31
w_4^*				-231639.30
w_5^*				640042.26
w_6^*				-1061800.52
w_7^*				1042400.18
w_8^*				-557682.99
w_9^*				125201.43



Data Set Size:

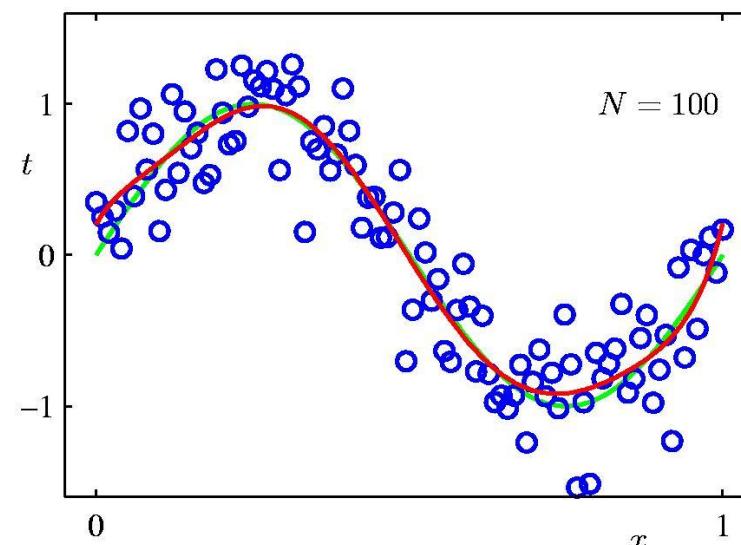
$$N = 15$$

9th Order Polynomial



$$N = 100$$

9th Order Polynomial



the larger the data set, the more complex (in other words more flexible) the model that we can afford to fit to the data.



3 solutions to over-fitting problem

- Rough heuristic: data-set-size \geq some multiple (say 5 or 10) of the number of adaptive parameters in the model.
- Regularization: adding a penalty term to the error function in order to discourage the coefficients from reaching large values.
- Local weighted regression(LWR)



Regularization

Penalize large coefficient values

The simplest such penalty term: sum of squares of all of the coefficients

$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

where $\|\mathbf{w}\|^2 = \mathbf{w}^T \mathbf{w} = w_1^2 + \dots + w_M^2$ (w_0 is omitted, or with its own regularization coefficient)

coefficient λ : governs the relative importance of the regularization term compared with the sum-of-squares error term.



Gradient Descent with Regularization

- Repeat {

$$\theta_0 \coloneqq \theta_0 - \alpha \sum_{i=1}^N (h_\theta(x_i) - y_i) x_i^0$$

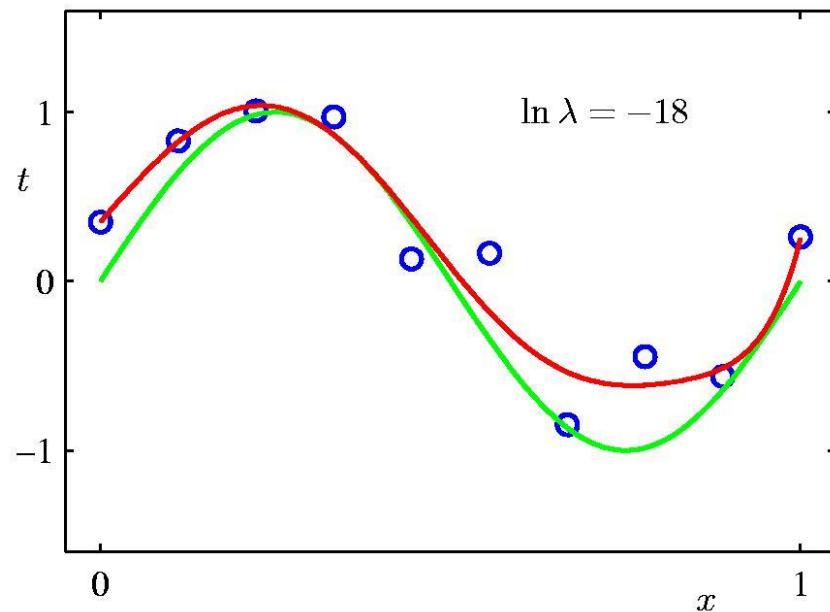
$$\theta_j \coloneqq \theta_j - \alpha \left[\sum_{i=1}^N (h_\theta(x_i) - y_i) x_i^j + \lambda \theta_j \right]$$

- }

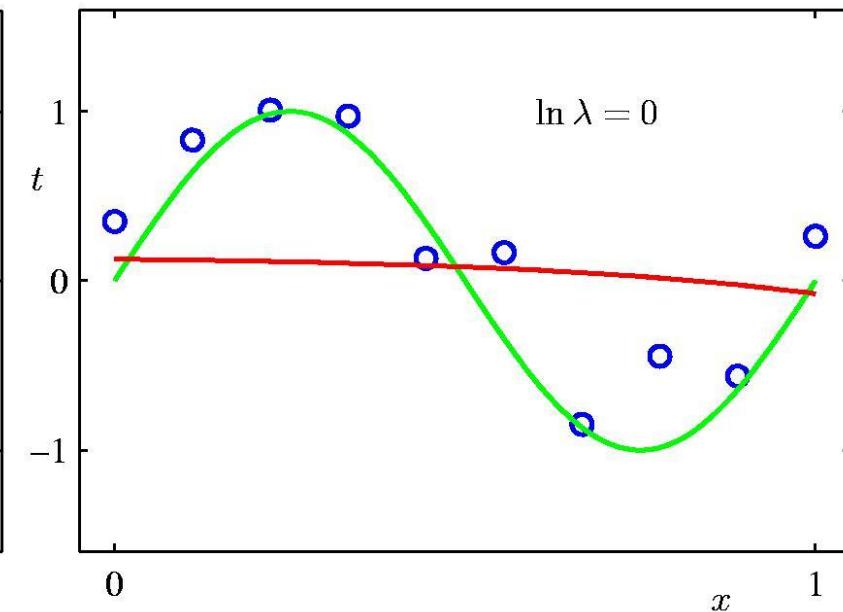


Regularization:

$$\ln \lambda = -18$$



$$\ln \lambda = 0$$





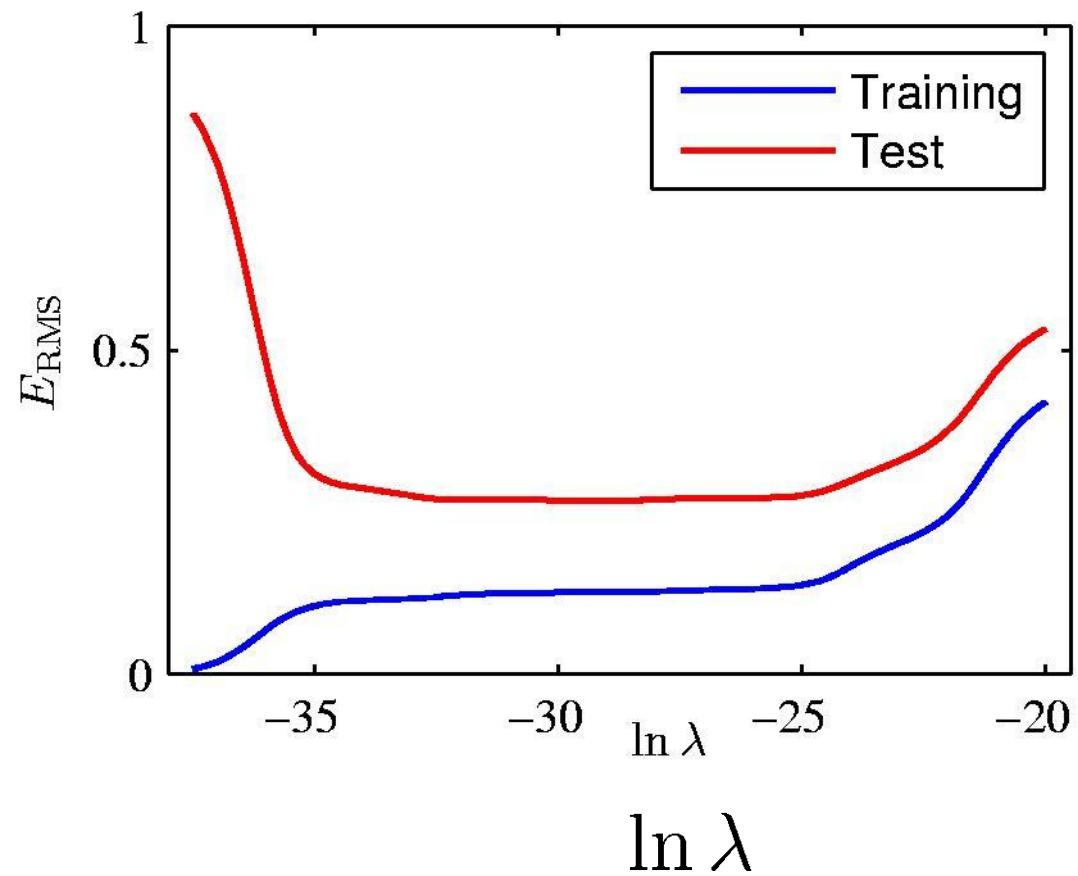
Polynomial Coefficients

	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
w_0^*	0.35	0.35	0.13
w_1^*	232.37	4.74	-0.05
w_2^*	-5321.83	-0.77	-0.06
w_3^*	48568.31	-31.97	-0.05
w_4^*	-231639.30	-3.89	-0.03
w_5^*	640042.26	55.28	-0.02
w_6^*	-1061800.52	41.32	-0.01
w_7^*	1042400.18	-45.95	-0.00
w_8^*	-557682.99	-91.53	0.00
w_9^*	125201.43	72.68	0.01



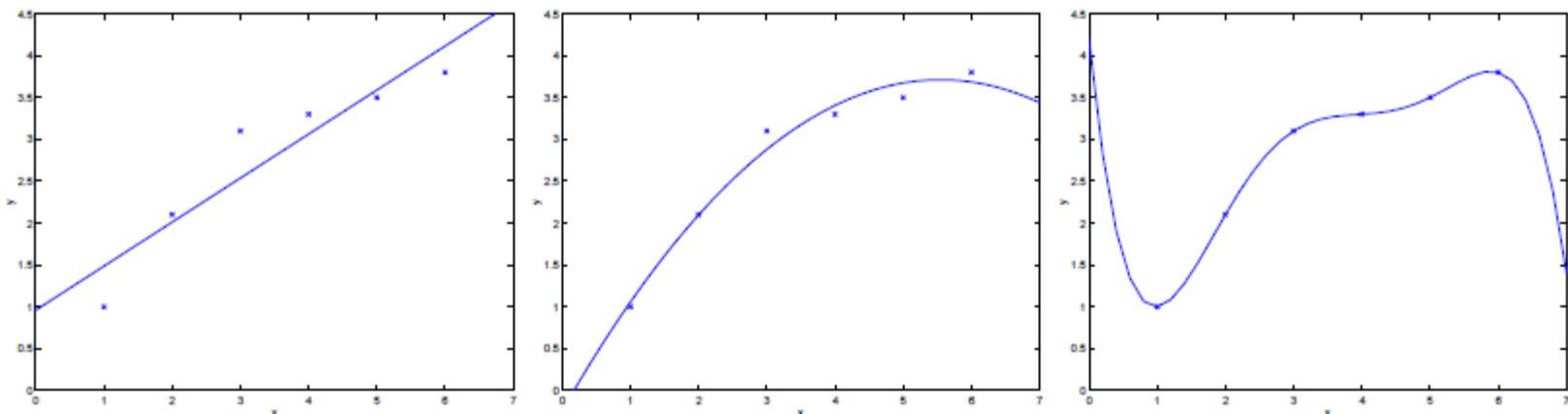
Regularization

E_{RMS}





Locally weighted linear regression



- Left: linear, underfitting
- Middle: quadra $y = \theta_0 + \theta_1 x + \theta_2 x^2$,
- Right: fitting a 5-th order polynomial $y = \sum_{j=0}^5 \theta_j x^j$.
 - overfitting
- Naively, it might seem that the more features we add, the better. However, there is also a danger in adding too many features
- LWR: makes the choice of features less critical.



In the original linear regression algorithm, to make a prediction at a query point x (i.e., to evaluate $h(x)$), we would:

1. Fit θ to minimize $\sum_i (y^{(i)} - \theta^T x^{(i)})^2$.
2. Output $\theta^T x$.

In contrast, the locally weighted linear regression algorithm does the following:

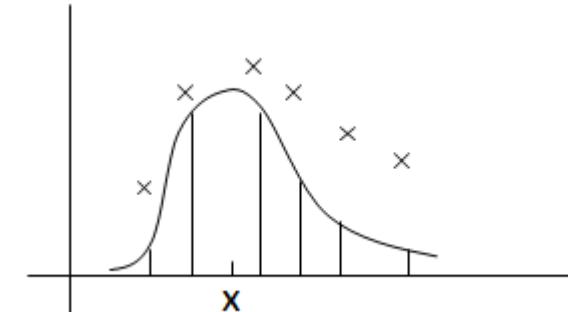
1. Fit θ to minimize $\sum_i w^{(i)} (y^{(i)} - \underline{\theta^T x^{(i)}})^2$.
2. Output $\theta^T x$.



$w^{(i)}$'s Non-negative valued weights

- If small, the error term will be pretty much ignored in the fit
- Standard choice

$$w^{(i)} = \exp\left(-\frac{(x^{(i)} - x)^2}{2\tau^2}\right)$$



- Depends on the particular point x at which we are trying to evaluate.
- If $|x^{(i)} - x|$ is small, $w^{(i)}$ is close to 1
- If $|x^{(i)} - x|$ is large, $w^{(i)}$ is small
- θ is chosen giving a much higher “weight” to the (errors on) training examples close to the query point x .
- τ : bandwidth parameter, controls how quickly the weight of a training example falls off with distance of its $x^{(i)}$ from the query point x



- Locally weighted linear regression: non-parameteric algorithm
 - Keep the entire training set around
- (unweighted) linear regression : parametric learning algorithm
 - Has a fixed, finite number of parameters, which fit to the data
 - Once fit the θ , stored them away, no longer need to keep the training data around to make future predictions



Logistic Regression



Classification

- The goal in classification is to take an input vector \mathbf{x} and to assign it to one of K discrete classes C_k where $k = 1, \dots, K$.
- The input space is divided into K decision regions whose boundaries are called decision boundaries or decision surfaces.
- Each region represents one class and each input is classified into one of regions (classes).



Linear Classification

- Linear classification means that the decision surfaces are linear functions of the input vector x and hence are defined by $(D - 1)$ - dimensional hyperplanes within the D -dimensional input space.
- Data sets whose classes can be separated exactly by linear decision surfaces are said to be linearly separable.



Linear Regression vs. Linear Classification

- Linear regression

$$y(x, w) = \sum_{j=0}^{M-1} w_j \phi_j(x)$$

- Linear classification

- Regression $y(x, w) = \sum_{j=0}^{M-1} w_j x_j$, $\phi_j(x) = x_j$, where x_j is the jth element of x

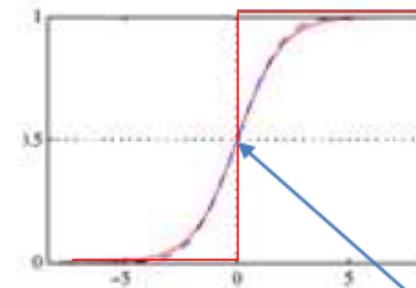
- Classification $c = f(y(x))$, where f is a function (activation function) that maps x into different classes C .



Activation Functions

- The simplest activity function is just y itself, i.e., $f(y)=y$.
- The most commonly used activity function is the sigmoid (S-shape) function

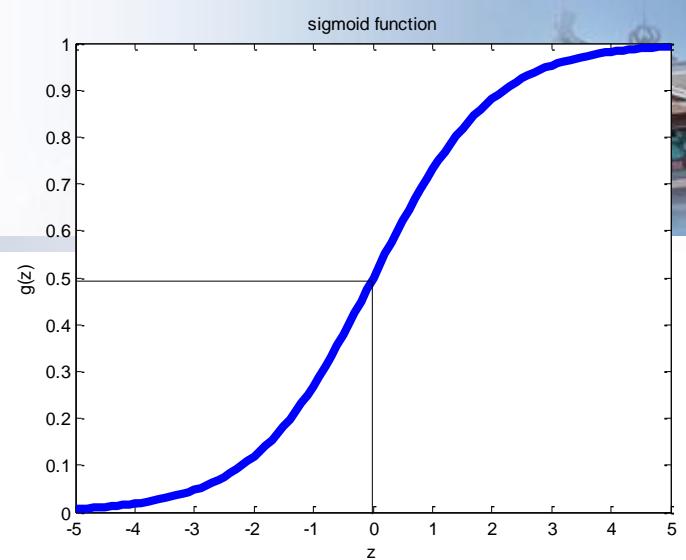
$$f(y) = \frac{1}{1 + e^{-y}}$$



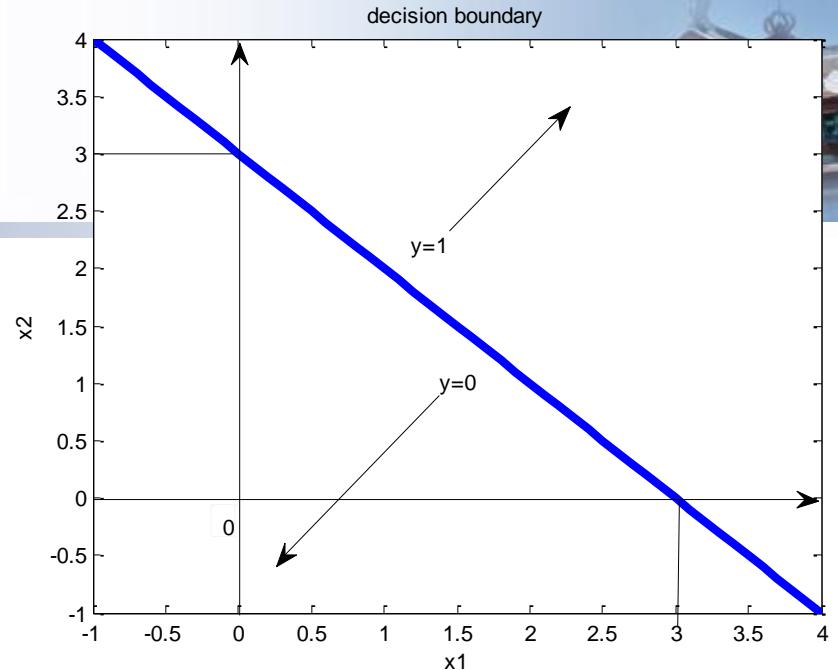
which maps (squashes) the real value y into a finite intervals (classes).

Step
threshold
function

- Classification with sigmoid activation function is often called logistic regression.



- Sigmoid function:
 - $Y > 0, f(y) > 0.5; y \rightarrow \infty, f(y) \rightarrow 1$
 - $y < 0, f(y) < 0.5; y \rightarrow -\infty, f(y) \rightarrow 0$
- $y = w^T x, f(y) = f(w^T x) = P(y=1|x; w)$
 - $w^T x > 0, P(y=1|x; w) > 0.5$, predict “ $y=1$ ”
 - $w^T x < 0, P(y=1|x; w) < 0.5$, predict “ $y=0$ ”
 - $w^T x = 0$, decision boundary



decision boundary:

In logistic regression: input = 0; $w^T x = 0$

Linear decision boundary: $w^T x = \sum w_i x_i$

non-linear decision boundary

Example1: linear decision boundary

$$X = [1 \ x_1 \ x_2], w = [-3 \ 1 \ 1]^T ; -3 + x_1 + x_2 = 0$$

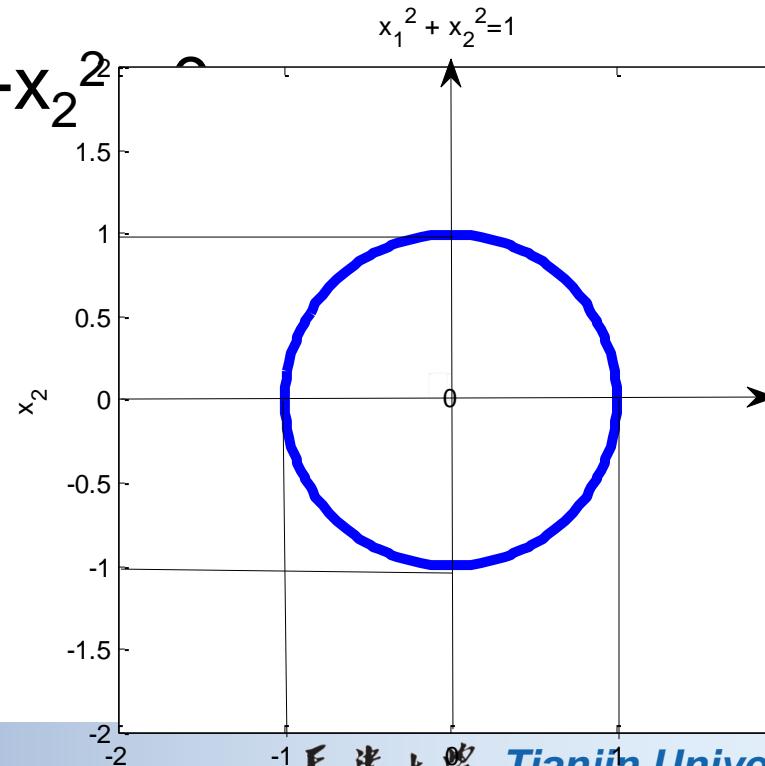
Predict “y=1” if $-3 + x_1 + x_2 > 0$



- Example 2: nonlinear decision boundary

- $X = [1 \ x_1 \ x_2 \ x_1^2 \ x_2^2]$, $w = [-1 \ 0 \ 0 \ 1 \ 1]$, $-1 + x_1^2 + x_2^2 = 0$

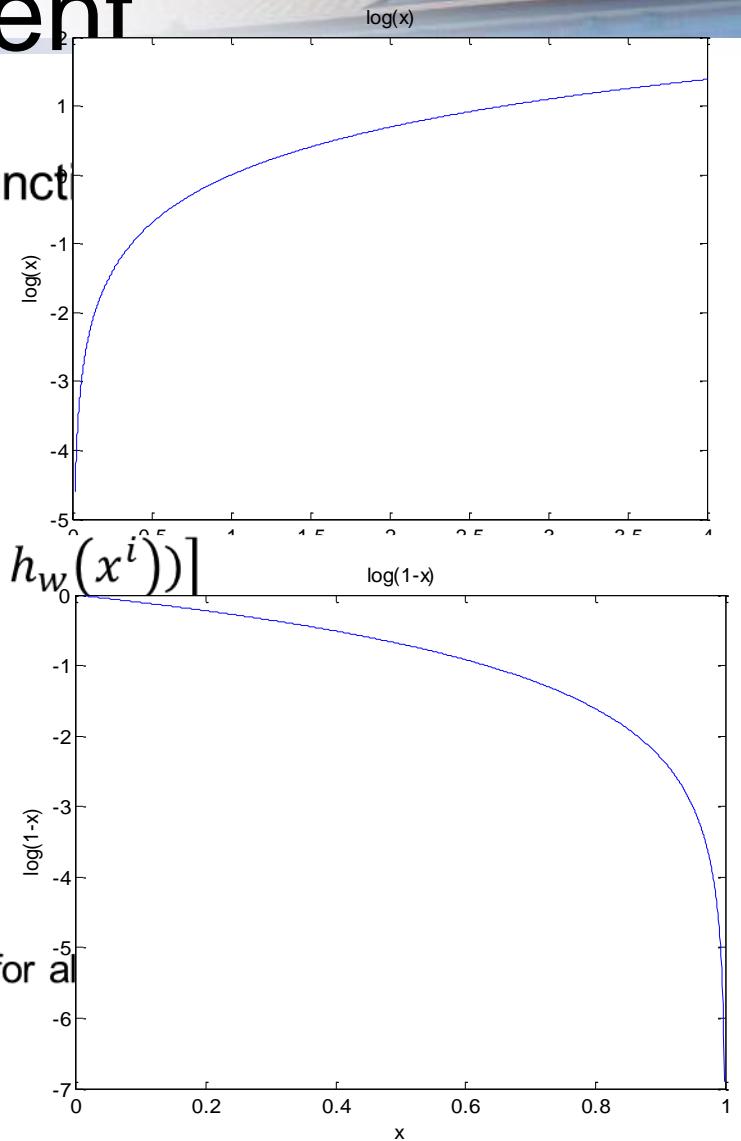
- Predict “y=1” if $-1 + x_1^2 + x_2^2 \geq 0$





Update rule 1: cost function & gradient descent

- Cost function $J(w) = \frac{1}{N} \sum_{i=1}^N Cost(h_w(x^i), y^i)$
- $Cost(h_w(x^i), y^i)$: can't use LMS(not a convex function with minima)
- $Cost(h_w(x^i), y^i) = \begin{cases} -\log(h_w(x^i)) & y^i = 1 \\ -\log(1 - h_w(x^i)) & y^i = 0 \end{cases}$
- Simplified cost function:
- $J(w) = \frac{1}{N} \sum_{i=1}^N [-y^i \log(h_w(x^i)) - (1 - y^i) \log(1 - h_w(x^i))]$
- Gradient descent
- $w_j := w_j - \alpha \frac{\partial}{\partial w_j} J(w)$
- $\frac{\partial}{\partial w_j} J(w) = \frac{1}{N} \sum_{i=1}^N (h_w(x^i) - y^i) x_j$
- Repeat until convergence{
 - $w_j := w_j - \alpha \frac{1}{N} \sum_{i=1}^N (h_w(x^i) - y^i) x_j$ (simultaneously update for all w_j)
 - }





- Feature scaling also applies to gradient descent in logistic regression.
- Differences with implementing linear regression:
- $$h_w(x^i) = \frac{1}{1+e^{-w^T x^i}}$$



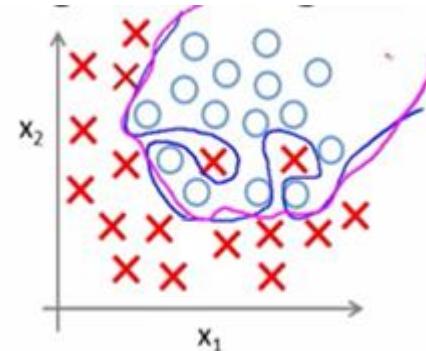
Update rule 2: log-likelihood gradient ascent

- Assume that: $P(y=1|x; w) = h_w(\mathbf{x}^i)$
 $P(y=0|x; w) = 1 - h_w(\mathbf{x}^i)$
- can be written more compactly as
- $P(y|x; w) = h_w(\mathbf{x}^i)^{y^i} (1 - h_w(\mathbf{x}^i))^{1-y^i}$
- likelihood function of the parameters
- $L_w = \prod_{i=1}^N p(y^i|\mathbf{x}^i, w) = \prod_{i=1}^N h_w(\mathbf{x}^i)^{y^i} (1 - h_w(\mathbf{x}^i))^{1-y^i}$
- maximize the log likelihood:
- $l_w = \log L_w = \sum_{i=1}^N \left[y^i \log(h_w(\mathbf{x}^i)) + (1 - y^i) \log(1 - h_w(\mathbf{x}^i)) \right]$
- Update rule: $w_j := w_j + \alpha \frac{\partial}{\partial w_j} l_w$ $\frac{\partial}{\partial w_j} l_w = \sum_{i=1}^N (y^i - h_w(\mathbf{x}^i)) x_j$



Regularized logistic regression

- $h_w(\mathbf{x}^i) = \frac{1}{1+e^{-w^T\phi}}$
- Cost function:



- $J(w)$
 $= \frac{1}{N} \sum_{i=1}^N [-y^i \log(h_w(\mathbf{x}^i)) - (1 - y^i) \log(1 - h_w(\mathbf{x}^i))] + \frac{\lambda}{2M} \sum_{j=1}^M w_j^2$
- Gradient descent:

Repeat until convergence{

$$w_0 := w_0 - \alpha \frac{1}{N} \sum_{i=1}^N (h_w(\mathbf{x}^i) - y^i) x_0$$

$$w_j := w_j - \alpha \frac{1}{N} \sum_{i=1}^N (h_w(\mathbf{x}^i) - y^i) x_j + \frac{\lambda}{M} w_j$$

}



Multiple Class Discriminant

- A single K -class discriminant consists of K linear functions of the form

$$y_k(x) = w_k^T x + w_{k0}$$

for $k=1, 2, \dots, K$

- assigning \mathbf{x} to class C_k if $y_k(\mathbf{x}) > y_j(\mathbf{x})$ for all $j \neq k$.

$$k^* = \arg \max_k y_k(x)$$



Discriminant Parameter Learning

- Learn the weights w , given input a set of input x and output t .
- Learning methods
 - Least-squares method
 - Fisher discriminant learning



Least-squares method (ML method)

- Estimate w such that the model predictions as close as possible to a set of target values
- Each class C_k is described by its own linear model so that

$$y_k(x) = w_k^T x + w_{k0} = [w_{k0} \ w_k^T] \begin{bmatrix} 1 \\ x \end{bmatrix}$$

$$k = 1, 2, \dots, K$$



Least-squares method

- In vector format, we have

$$Y^{(Kx)}(x) = \begin{bmatrix} y_1(x) \\ y_2(x) \\ \vdots \\ y_K(x) \end{bmatrix}$$
$$W^{(D+1) \times K} = \begin{bmatrix} w_{10} & w_{20} & \dots & w_{K0} \\ w_1 & w_2 & \dots & w_K \end{bmatrix}$$
$$Y(x) = W^T \begin{bmatrix} 1 \\ x \end{bmatrix} = W^T x'$$



Least-squares method (cont'd)

- Given N samples x_1, x_2, \dots, x_N and their corresponding class labels t_1, t_2, \dots, t_N , where t_n is a $K \times 1$ vector that follows 1-of-K coding, i.e., the elements t_n are all zeros, except for the class that x_n belongs to. Least-squares solution is to find W that minimizes the sum-of-squares errors E , i.e.,

$$\begin{aligned} E &= \sum_{n=1}^N [t_n - Y(x_n)]^T [t_n - Y(x_n)] \\ &= \sum_{n=1}^N [t_n - W^T x_n^{' }]^T [t_n - W^T x_n^{'}] \end{aligned}$$



1-of-K Coding

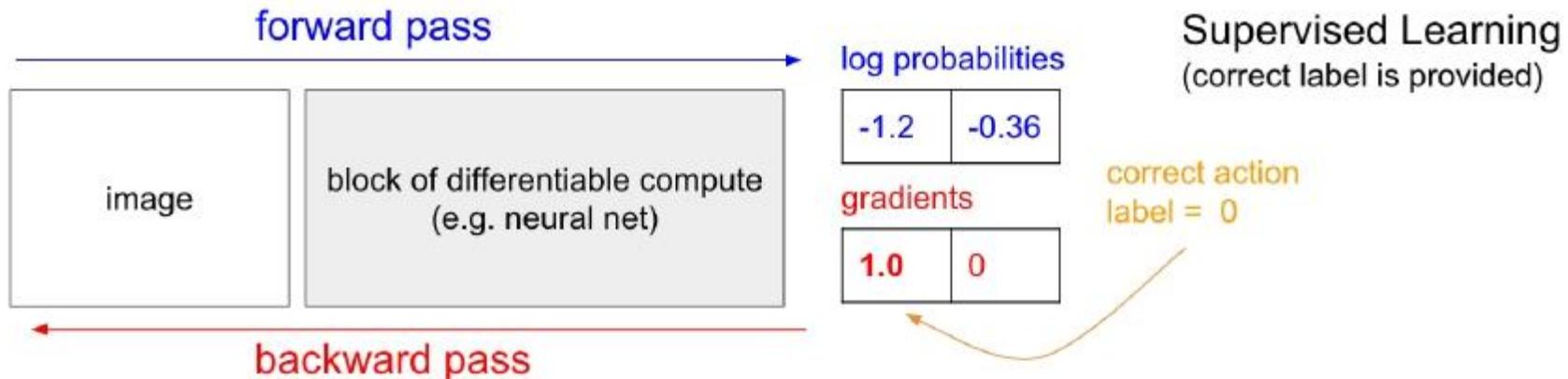
- discrete variables that can take on one of K possible mutually exclusive states.
- convenient representation: 1-of-K scheme
 - variable is represented by a K-dimensional vector y in which one of the elements y_k equals 1, and all remaining elements equal 0.
- E.g.

g	y_1	y_2	y_3	y_4
3	0	0	1	0
1	1	0	0	0
2	0	1	0	0
4	0	0	0	1
1	1	0	0	0

$$y^1 = [0, 0, 1, 0]^T$$



Task: Classify an Image of a Number

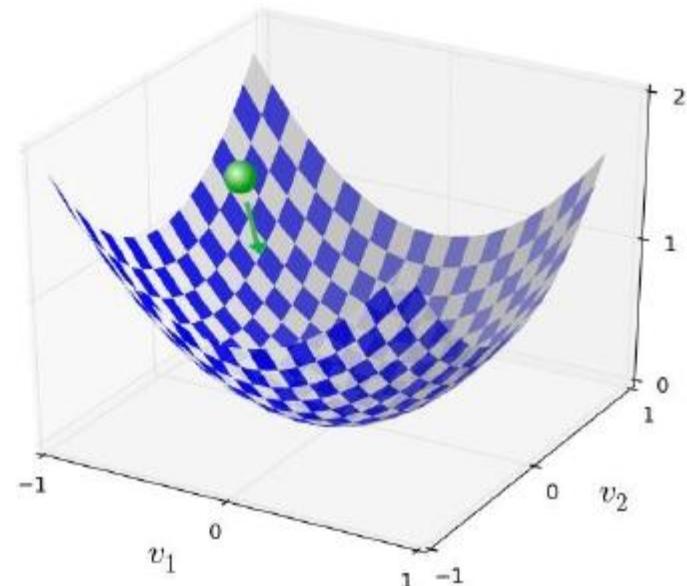


Ground truth for “6”:

$$y(x) = (0, 0, 0, 0, 0, 0, 1, 0, 0, 0)^T$$

“Loss” function:

$$C(w, b) \equiv \frac{1}{2n} \sum_x \|y(x) - a\|^2$$





Least-squares method (cont'd)

Take the derivatives of E w.r.t to W and setting it to zero, we can solve W as follows

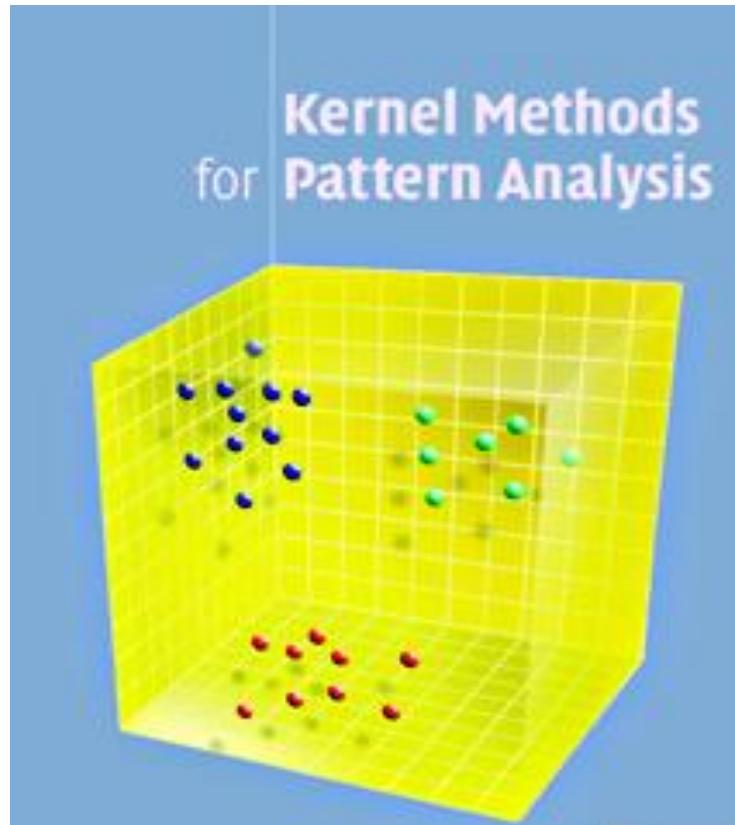
$$W = (X^T X)^{-1} X^T T$$

where

$$X^{Nx(D+1)} = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{bmatrix} \quad T^{Nxk} = \begin{bmatrix} t_1^T \\ t_2^T \\ \vdots \\ t_N^T \end{bmatrix}$$



Kernel Method



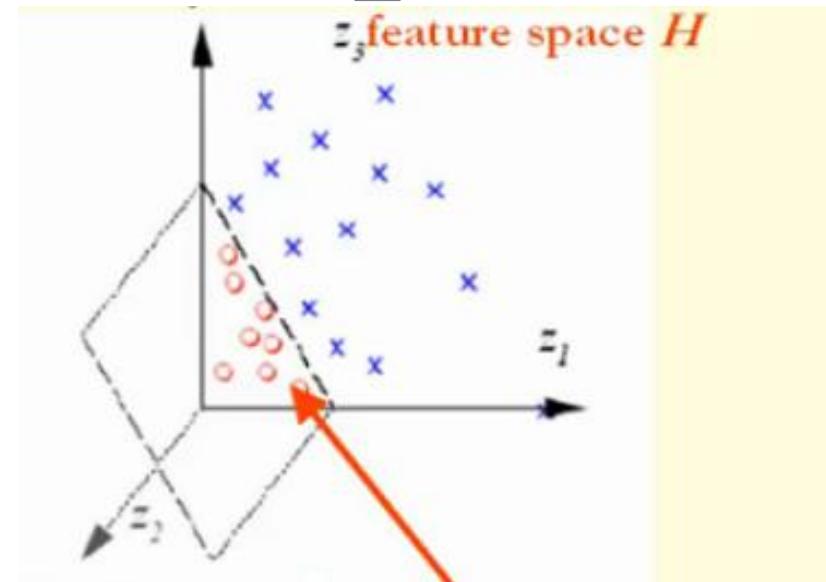
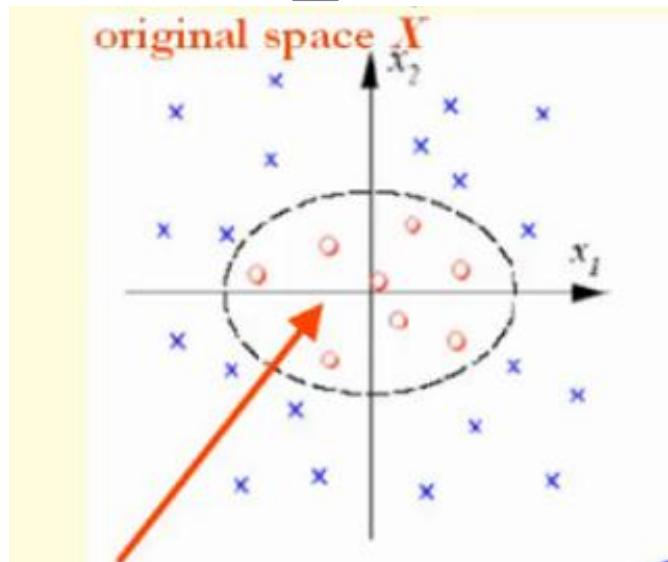


The Basic Idea

feature mapping

$$\phi: R^2 \rightarrow R^3$$

$$(x_1, x_2) \mapsto (z_1, z_2, z_3) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$$



$$\frac{x_1^2}{a^2} + \frac{x_2^2}{b^2} = 1$$

$$\frac{x_1^2}{a^2} + \frac{x_2^2}{b^2} = 1 \Rightarrow \frac{1}{a^2}z_1 + 0 \cdot z_2 + \frac{1}{a^2}z_3 = 1$$

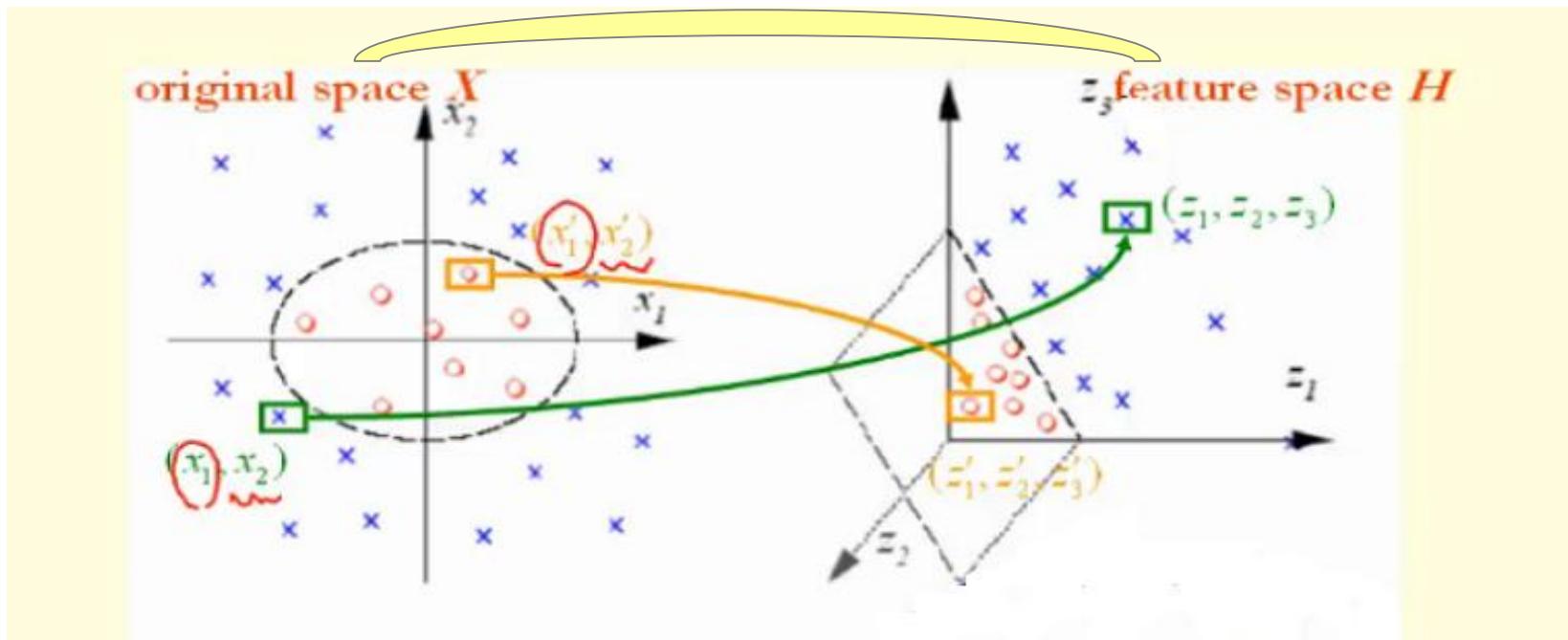


The Basic Idea

feature mapping

$$\phi: R^2 \rightarrow R^3$$

$$(x_1, x_2) \mapsto (z_1, z_2, z_3) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$$



$$\begin{aligned} & \langle \phi(x_1, x_2), \phi(x'_1, x'_2) \rangle = \langle (z_1, z_2, z_3), (z'_1, z'_2, z'_3) \rangle = \langle (x_1^2, \sqrt{2}x_1x_2, x_2^2), (x'_1^2, \sqrt{2}x'_1x'_2, x'_2^2) \rangle \\ &= x_1^2 x'_1^2 + 2x_1 x_2 x'_1 x'_2 + x_2^2 x'_2^2 = (x_1 x'_1 + x_2 x'_2)^2 = (\langle x, x' \rangle)^2 = K(x, x') \end{aligned}$$



The Basic Idea



Kernel Methods work by

- embedding data in a vector space (with more high dimensionality)
- looking for (linear) relations in such space

If map is chosen suitably, complex relations can be simplified, and easily detected.



The Basic Idea

Distance in the Feature Space

$$\begin{aligned}\|\phi(x) - \phi(x')\|^2 &= (\phi(x) - \phi(x'))^T (\phi(x) - \phi(x')) \\&= \phi(x)^T \phi(x) - 2\phi(x)^T \phi(x') + \phi(x')^T \phi(x') \\&= \langle \phi(x), \phi(x) \rangle - 2 \langle \phi(x), \phi(x') \rangle + \langle \phi(x'), \phi(x') \rangle \\&= K(x, x) - 2K(x, x') + K(x', x')\end{aligned}$$



The Basic Idea

Angle in the Feature Space

$$\langle \phi(x), \phi(x') \rangle = \|\phi(x)\| \|\phi(x')\| \cos \theta$$

$$\Rightarrow \cos \theta = \frac{\langle \phi(x), \phi(x') \rangle}{\|\phi(x)\| \|\phi(x')\|} = \frac{\langle \phi(x), \phi(x') \rangle}{\sqrt{\langle \phi(x), \phi(x) \rangle} \sqrt{\langle \phi(x'), \phi(x') \rangle}}$$

$$= \frac{K(x, x')}{\sqrt{K(x, x)} \sqrt{K(x', x')}}$$



The Basic Idea

Recall that

$$\begin{aligned} & \langle \phi(x_1, x_2), \phi(x'_1, x'_2) \rangle = \langle (z_1, z_2, z_3), (z'_1, z'_2, z'_3) \rangle = \langle (x_1^2, \sqrt{2}x_1x_2, x_2^2), (x'_1{}^2, \sqrt{2}x'_1x'_2, x'_2{}^2) \rangle \\ & = x_1^2 x'_1{}^2 + 2x_1 x_2 x'_1 x'_2 + x_2^2 x'_2{}^2 = (x_1 x'_1 + x_2 x'_2)^2 = (\langle x, x' \rangle)^2 := K(x, x') \end{aligned}$$

Distance in the Feature Space

$$\begin{aligned} & \|\phi(x) - \phi(x')\|^2 = (\phi(x) - \phi(x'))^T (\phi(x) - \phi(x')) \\ & = \phi(x)^T \phi(x) - 2\phi(x)^T \phi(x') + \phi(x')^T \phi(x') \\ & = \langle \phi(x), \phi(x) \rangle - 2 \langle \phi(x), \phi(x') \rangle + \langle \phi(x'), \phi(x') \rangle \\ & = K(x, x) - 2K(x, x') + K(x', x') \end{aligned}$$

Angle in the Feature Space

$$\begin{aligned} & \langle \phi(x), \phi(x') \rangle = \|\phi(x)\| \|\phi(x')\| \cos \theta \\ \Rightarrow & \cos \theta = \frac{\langle \phi(x), \phi(x') \rangle}{\|\phi(x)\| \|\phi(x')\|} = \frac{\langle \phi(x), \phi(x') \rangle}{\sqrt{\langle \phi(x), \phi(x) \rangle} \sqrt{\langle \phi(x'), \phi(x') \rangle}} = \frac{K(x, x')}{\sqrt{K(x, x)} \sqrt{K(x', x')}} \end{aligned}$$



The Basic Idea

Inner Product Matrix/Gram Matrix/Kernel Matrix

$$K = \begin{bmatrix} <\phi(x_1), \phi(x_1)> & \cdots & <\phi(x_1), \phi(x_N)> \\ \vdots & \ddots & \vdots \\ <\phi(x_N), \phi(x_1)> & \cdots & <\phi(x_1), \phi(x_N)> \end{bmatrix}$$
$$= \begin{bmatrix} K(x_1, x_1) & \cdots & K(x_1, \dots, x_N) \\ \vdots & \ddots & \vdots \\ K(x_N, x_1) & \cdots & K(x_N, x_N) \end{bmatrix}$$



Kernel Function and Feature Mapping

One definition:

Finitely positive semi-definite functions

One theorem:

Characterization of kernels



Finitely Positive Semi-definite Functions

A function

$$K : X \times X \rightarrow R$$

satisfies the finitely positive semi-definite property if it is a symmetric function for which the matrices formed by restriction to any finite subset of the space are positive semi-definite.

Given any

$$\{x_1, \dots, x_n\} \longrightarrow K = \begin{bmatrix} k(x_1, x_1) & \cdots & k(x_1, x_n) \\ \vdots & & \vdots \\ k(x_n, x_1) & \cdots & k(x_n, x_n) \end{bmatrix} \quad \text{any } y \quad y^T K y \geq 0$$



Finitely Positive Semi-definite Functions

Example: Show that

$$K(x, z) = \langle x, z \rangle$$

is a finitely positive semi-definite function.

Given any $\{x_1, \dots, x_n\}$ $\Rightarrow K = \begin{bmatrix} k(x_1, x_1) & \cdots & k(x_1, x_n) \\ \vdots & & \vdots \\ k(x_n, x_1) & \cdots & k(x_n, x_n) \end{bmatrix} = \begin{bmatrix} \langle x_1, x_1 \rangle & \cdots & \langle x_1, x_n \rangle \\ \vdots & & \vdots \\ \langle x_n, x_1 \rangle & \cdots & \langle x_n, x_n \rangle \end{bmatrix}$

$$= \begin{bmatrix} x_1^T x_1 & \cdots & x_1^T x_n \\ \vdots & & \vdots \\ x_n^T x_1 & \cdots & x_n^T x_n \end{bmatrix} = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix} \begin{bmatrix} x_1 & \cdots & x_n \end{bmatrix}$$

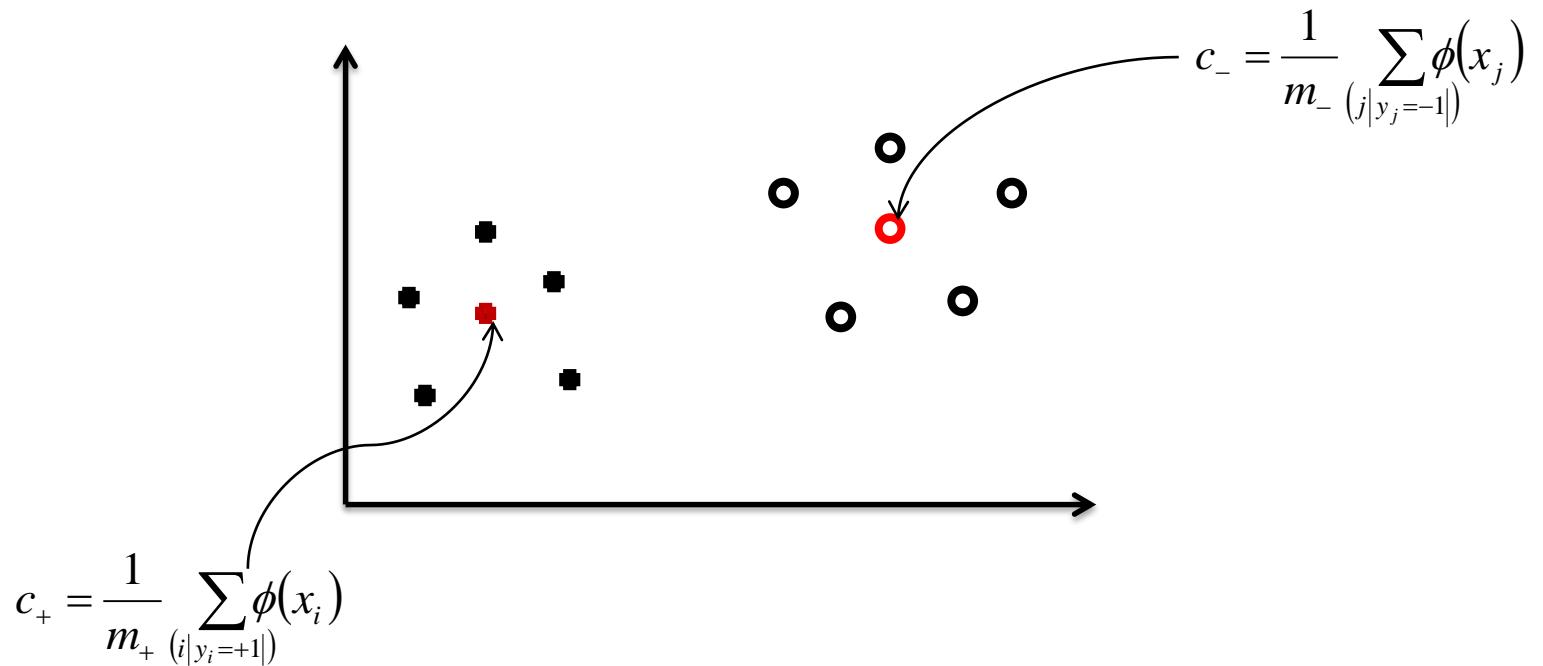
$$\begin{aligned} (AB)^T &= B^T A^T & \Rightarrow y^T K y &= y^T \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix} \begin{bmatrix} x_1 & \cdots & x_n \end{bmatrix} y = (\begin{bmatrix} x_1 & \cdots & x_n \end{bmatrix} y)^T (\begin{bmatrix} x_1 & \cdots & x_n \end{bmatrix} y) \\ &&&= \| \begin{bmatrix} x_1 & \cdots & x_n \end{bmatrix} y \|^2 \geq 0 \end{aligned}$$



A Simple Recognition Algorithm

Training Samples: feature mapping

$$\{(x_1, y_1), \dots, (x_n, y_n)\} \subset R^d \times \{1, -1\} \xrightarrow{\phi} \{(\phi(x_1), y_1), \dots, (\phi(x_n), y_n)\} \subset H \times \{1, -1\}$$



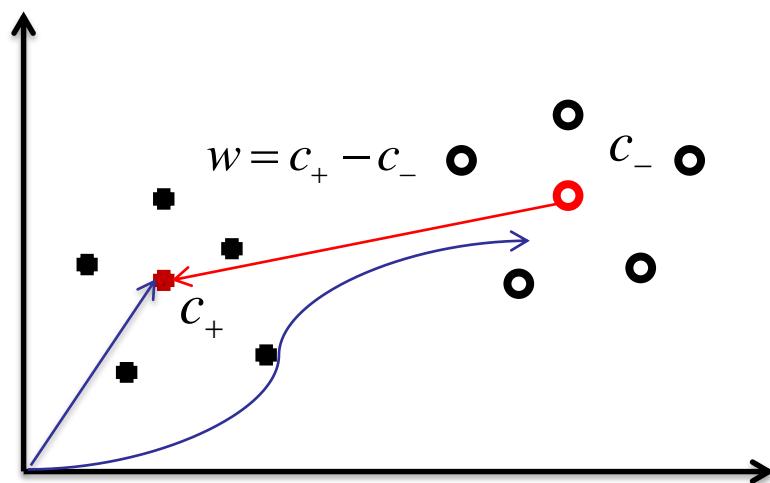


A Simple Recognition Algorithm

Training Samples:

feature mapping
 ϕ

$$\{(x_1, y_1), \dots, (x_n, y_n)\} \subset R^d \times \{1, -1\} \rightarrow \{(\phi(x_1), y_1), \dots, (\phi(x_n), y_n)\} \subset H \times \{1, -1\}$$

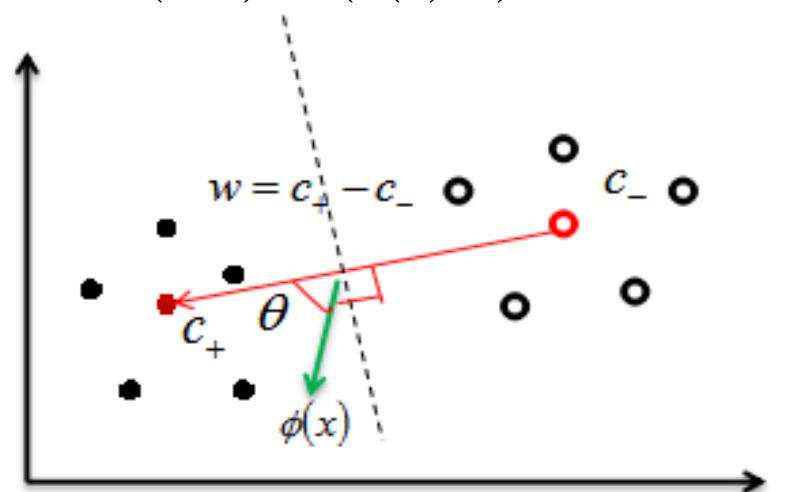




A Simple Recognition Algorithm

Testing Sample:

$$(x, y) \xrightarrow{\phi} (\phi(x), y) \rightarrow y = ?$$



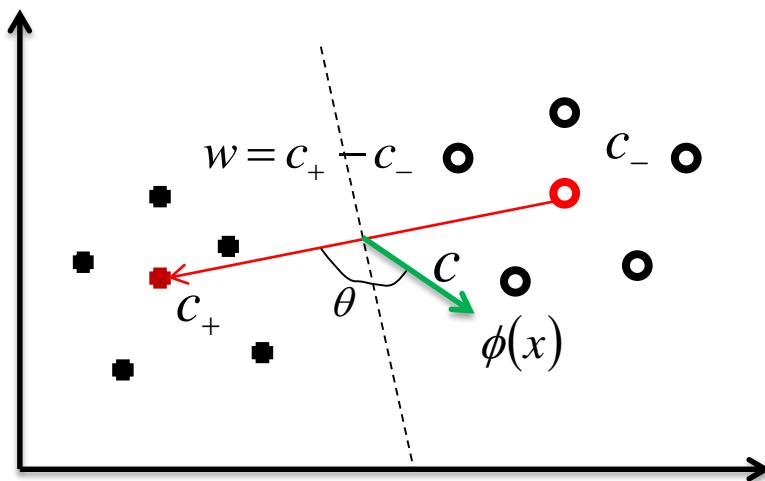
$$y = 1 \Leftrightarrow 0 \leq \theta < \frac{\pi}{2} \Leftrightarrow 0 < \cos \theta \leq 1 \Leftrightarrow 0 < \frac{<\phi(x) - c, w>}{\sqrt{\|\phi(x) - c\|} \sqrt{\|w\|}} \leq 1 \Leftrightarrow <\phi(x) - c, w> \geq 0$$



A Simple Recognition Algorithm

Testing Sample:

$$(x, y) \xrightarrow{\phi} (\phi(x), y) \rightarrow y = ?$$



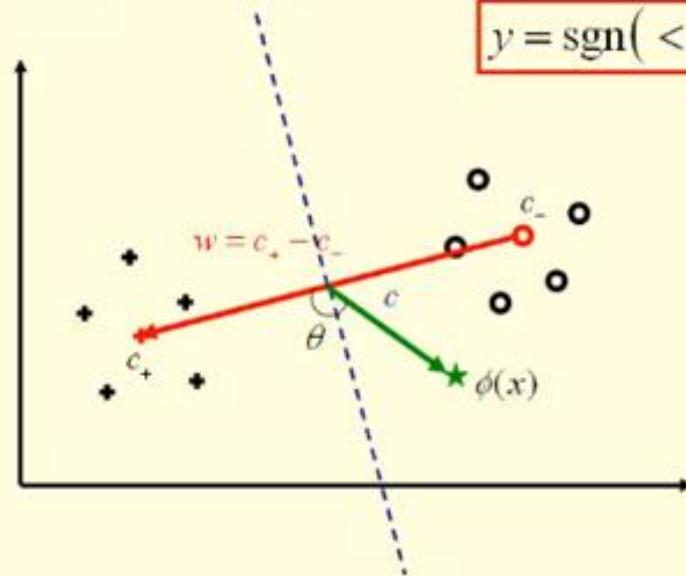
$$y = -1 \Leftrightarrow \frac{\pi}{2} < \theta \leq \pi \Leftrightarrow -1 \leq \cos \theta < 0 \Leftrightarrow -1 \leq \frac{\langle \phi(x) - c, w \rangle}{\sqrt{\|\phi(x) - c\|} \sqrt{\|w\|}} < 0 \Leftrightarrow \langle \phi(x) - c, w \rangle < 0$$



A Simple Recognition Algorithm

Testing Sample: $(x, y) \xrightarrow{\phi} (\phi(x), y) \rightarrow y = ?$

$$y = \text{sgn}(\langle \phi(x) - c, w \rangle)$$



$$y = 1 \Leftrightarrow 0 \leq \theta < \frac{\pi}{2} \Leftrightarrow 0 < \cos \theta \leq 1 \Leftrightarrow 0 < \frac{\langle \phi(x) - c, w \rangle}{\sqrt{\|\phi(x) - c\|} \sqrt{\|w\|}} \leq 1 \Leftrightarrow \langle \phi(x) - c, w \rangle > 0$$

$$y = -1 \Leftrightarrow \frac{\pi}{2} < \theta \leq \pi \Leftrightarrow -1 \leq \cos \theta < 0 \Leftrightarrow -1 \leq \frac{\langle \phi(x) - c, w \rangle}{\sqrt{\|\phi(x) - c\|} \sqrt{\|w\|}} < 0 \Leftrightarrow \langle \phi(x) - c, w \rangle < 0$$



A Simple Recognition Algorithm

- How to compute $y = \text{sgn}(\langle \phi(x) - c, w \rangle)$?
- One way:



$$\{(x_1, y_1), \dots, (x_n, y_n)\} \subset R^d \times \{1, -1\} \xrightarrow{\phi} \{(\phi(x_1), y_1), \dots, (\phi(x_n), y_n)\} \subset H \times \{1, -1\}$$



$$(x, y) \xrightarrow{\phi} (\phi(x), y) \rightarrow y = ?$$

$$y = \text{sgn}(\langle \phi(x) - c, w \rangle)$$



A Simple Recognition Algorithm

How to compute $y = \text{sgn}(\langle \phi(x) - c, w \rangle)?$

The other way

$$\begin{aligned}\langle \phi(x) - c, w \rangle &= w^T (\phi(x) - c) = w^T \phi(x) - w^T c = (c_+ - c_-)^T \phi(x) - \frac{1}{2} (c_+ - c_-)^T (c_+ + c_-) \\ &= \left(\frac{1}{m_+} \sum_{(i|y_i=1)} \phi(x_i) - \frac{1}{m_-} \sum_{(i|y_i=-1)} \phi(x_i) \right)^T \phi(x) - \left(\frac{1}{2} c_+^T c_+ - c_+^T c_- + \frac{1}{2} c_-^T c_- \right) \\ &= \left(\frac{1}{m_+} \sum_{(i|y_i=1)} \phi(x_i)^T \phi(x) - \frac{1}{m_-} \sum_{(i|y_i=-1)} \phi(x_i)^T \phi(x) \right) - b \\ &= \left(\frac{1}{m_+} \sum_{(i|y_i=1)} K(x_i, x) - \frac{1}{m_-} \sum_{(i|y_i=-1)} K(x_i, x) \right) - b\end{aligned}$$

What is b? Exercise!

Note that if we only have the Kernel Function K, we still can find y.



Some Widely Used Kernel Functions

Linear Kernel

$$k(x, z) = \langle x, z \rangle$$

Polynomial Kernel

$$k(x, z) = (\langle x, z \rangle + 1)^r, r \in \mathbb{Z}^+$$

RBF(Gaussian) Kernel

$$k(x, z) = \exp\left(\frac{-\|x - z\|^2}{2\sigma^2}\right), \sigma \in R - \{0\}$$



Kernel Function and Feature Mapping

- Is ϕ necessary? **Not necessary!**
- Can we only employ κ ? **Yes!**
- What kind of κ can be used? **Finitely positive semi-definite!**
- Given a feature mapping ϕ , can we find a corresponding kernel function κ computing the dot product in feature space? **Yes!**
- Given a kernel function κ , can we construct a feature space H (i.e., a feature mapping ϕ) such that κ computes the dot product in H ? **Yes!**



The review of the Linear regression

Given m observation x_1, \dots, x_N and the corresponding responses y_1, \dots, y_N ,

$$X_1 = \begin{bmatrix} 1 & (x_1)^T \\ \vdots & \vdots \\ 1 & (x_N)^T \end{bmatrix}$$





The review of the Linear regression

Normal Equation

$$J(\theta) = \frac{1}{2} (X_1 \theta - y)^T (X_1 \theta - y)$$

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \frac{1}{2} [X_1^T X_1 \theta + X_1^T X_1 \theta - X_1^T y - X_1^T y] \\ &= X_1^T X_1 \theta - X_1^T y \stackrel{\text{set}}{=} 0\end{aligned}$$

$$X_1^T X_1 \theta = X_1^T y \quad \theta = (X_1^T X_1)^{-1} X_1^T y$$



Kernel-based Linear Regression

In addition, a feature mapping ϕ and the corresponding kernel function κ are also given

$$X_1 = \begin{bmatrix} 1 & \phi(x_1)^T \\ \vdots & \vdots \\ 1 & \phi(x_N)^T \end{bmatrix} = \begin{bmatrix} \mathbf{1}_{N \times 1}, X \end{bmatrix}$$
$$X = \begin{bmatrix} \phi(x_1)^T \\ \vdots \\ \phi(x_N)^T \end{bmatrix}$$

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$



Kernel-based Linear Regression

Again,

$$X_1^T X_1 \theta = X_1^T y \Rightarrow \theta = (X_1^T X_1)^{-1} X_1^T y$$

Note that

$$\theta = (X_1^T X_1) (X_1^T X_1)^{-1} (X_1^T X_1)^{-1} X_1^T y = (X_1^T X_1) (X_1^T X_1)^{-2} X_1^T y = X_1^T \alpha,$$

where $\alpha = X_1 (X_1^T X_1)^{-2} X_1^T y$

$$X_1 = \begin{bmatrix} 1 & \phi(x_1)^T \\ \vdots & \vdots \\ 1 & \phi(x_N)^T \end{bmatrix} = \begin{bmatrix} \mathbf{1}_{N \times 1}, \mathbf{X} \end{bmatrix}$$

$$X = \begin{bmatrix} \phi(x_1)^T \\ \vdots \\ \phi(x_N)^T \end{bmatrix}$$



Kernel-based Linear Regression

Now we want to find α

$$X_1^T X_1 \theta = X_1^T y$$

Substitute θ by $X_1^T \alpha$

$$X_1^T X_1 X_1^T \alpha = X_1^T y$$

and multiply both sides by X

$$X_1 X_1^T X_1 X_1^T \alpha = X_1 X_1^T y$$

$$\Rightarrow K_1^2 \alpha = K_1 y, \quad \text{where } K_1 = X_1 X_1^T$$

$$\Rightarrow K_1 \alpha = y, \quad \text{if } \det(K) \neq 0$$

$$\Rightarrow \alpha = K_1^{-1} y$$



Kernel-based Linear Regression

$$K_1 = X_1 X_1^T = [\mathbf{1}_{N \times 1}, X] \begin{bmatrix} \mathbf{1}_{1 \times N} \\ X^T \end{bmatrix} = \mathbf{1}_{N \times N} + XX^T$$
$$= \mathbf{1}_{N \times N} + K$$

$$X = \begin{bmatrix} \phi(x_1)^T \\ \vdots \\ \phi(x_N)^T \end{bmatrix}$$

$$XX^T = \begin{bmatrix} \phi(x_1^{(1)}) & \cdots & \phi(x_1^{(n)}) \\ \phi(x_2^{(1)}) & \cdots & \phi(x_2^{(n)}) \\ \vdots & \vdots & \vdots \\ \phi(x_N^{(1)}) & \cdots & \phi(x_N^{(n)}) \end{bmatrix} \begin{bmatrix} \phi(x_1^{(1)}) & \cdots & \phi(x_N^{(1)}) \\ \vdots & \cdots & \vdots \\ \vdots & \vdots & \vdots \\ \phi(x_1^{(n)}) & \cdots & \phi(x_N^{(n)}) \end{bmatrix}$$



Kernel-based Linear Regression

$$\begin{aligned} &= \begin{bmatrix} \langle \phi(x_1), \phi(x_1) \rangle & \cdots & \langle \phi(x_1), \phi(x_N) \rangle \\ \vdots & \cdots & \vdots \\ \vdots & \vdots & \vdots \\ \langle \phi(x_N), \phi(x_1) \rangle & \cdots & \langle \phi(x_N), \phi(x_N) \rangle \end{bmatrix} \\ &= \begin{bmatrix} K(x_1, x_1) & \cdots & K(x_1, x_N) \\ \vdots & \cdots & \vdots \\ \vdots & \vdots & \vdots \\ K(x_N, x_1) & \cdots & K(x_N, x_N) \end{bmatrix} = K \end{aligned}$$



Kernel-based Linear Regression



Hence,

$$\alpha = (\mathbf{1}_{N \times N} + K)^{-1} y$$

$$\theta = X^T \alpha = X^T (\mathbf{1}_{N \times N} + K)^{-1} y$$

$$X = \begin{bmatrix} \phi(x_1)^T \\ \vdots \\ \phi(x_N)^T \end{bmatrix}$$

$$X^T \text{ is } (n+1) \times N$$
$$K = XX^T \quad K_1 = \mathbf{1}_{N \times N} + K \quad \mathbf{1}_{N \times N} + K \text{ is } N \times N$$
$$y \text{ is } N \times 1$$
$$\theta \text{ is } (n+1) \times 1$$
$$\alpha \text{ is } N \times 1$$

$$X_1 = [\mathbf{1}_{N \times 1}, X]$$



Kernel-based Linear Regression

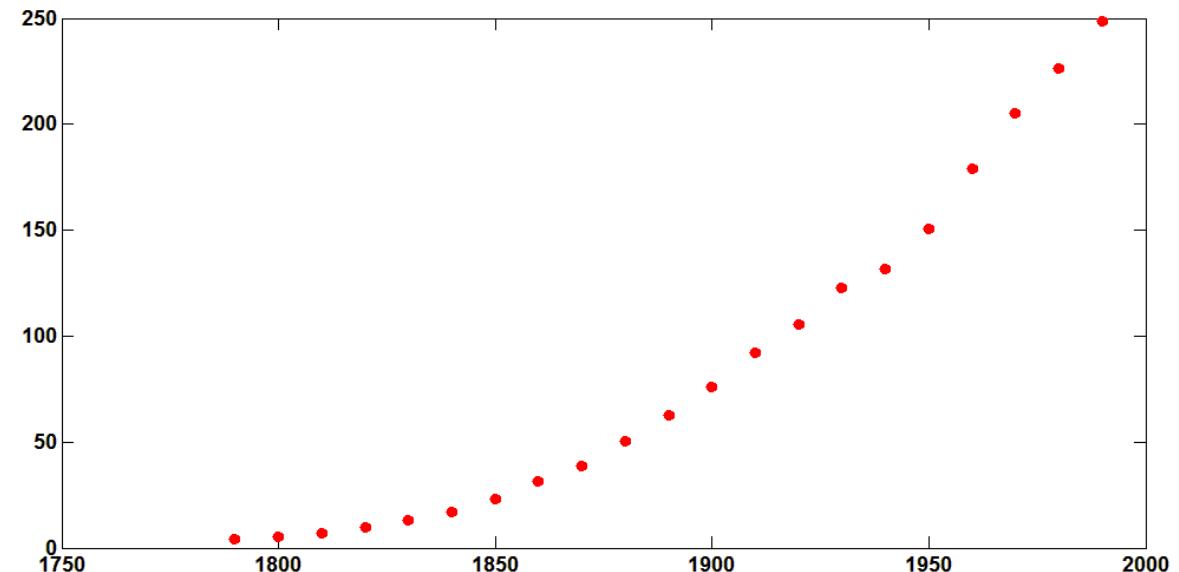
For a input x , we have

$$\begin{aligned}f(x) &= \left[1, \phi(x)^T \right] \theta = \left[1, \phi(x)^T \right] X_1^T \alpha \\&= \left[1, \phi(x)^T \right] \begin{bmatrix} \mathbf{1}_{1 \times N} \\ X^T \end{bmatrix} \alpha \\&= \left(\mathbf{1}_{1 \times N} + (X\phi(x))^T \right) \alpha \\&= \left[1 + \kappa(x_1, x), \dots, 1 + \kappa(x_N, x) \right] \alpha \\&= (\alpha_1 + \dots + \alpha_N) + \alpha_1 \kappa(x_1, x) + \dots + \alpha_N \kappa(x_N, x)\end{aligned}$$



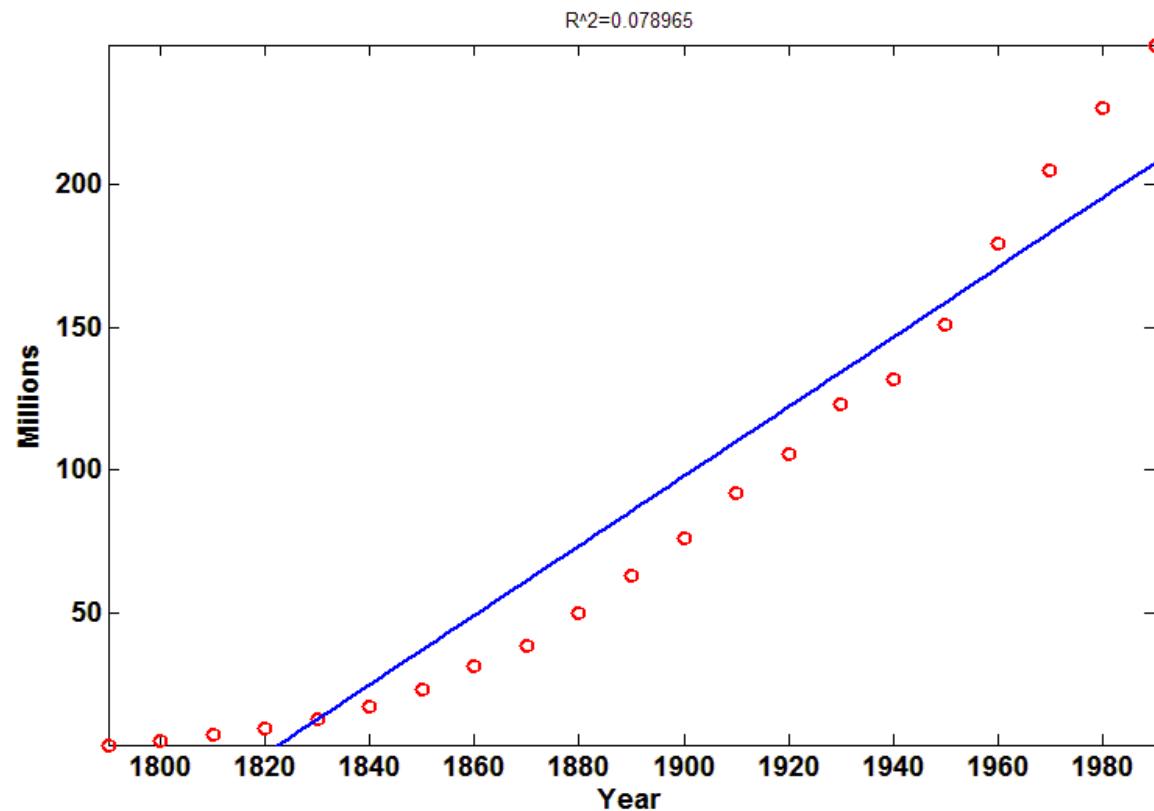
Example: Predicting the US Population

```
1 -      clear
2 -      clc
3 -      close all
4 -
5 -      load census;
6 -
7 -      x = cdate;
8 -      y = pop;
9 -      plot(x,y, 'ro');
10 -     hold on
```





```
12 %% Linear regression
13
14 - [N, d]=size(x);
15 - X=[ones(N, 1), x];
16 - hw=pinv(X'*X)*X'*y;
17 - yhat=X*hw;
18 - R2=norm(y-yhat)^2/norm(y-mean(y))^2;
19
20 - xfit=linspace(min(x), max(x), 100)';
21 - yfit=[ones(size(xfit, 1), 1), xfit]*hw;
22 - plot(xfit, yfit, 'b-');
23 - xlabel('Year');
24 - ylabel('Millions');
25 - axis([min(x), max(x), min(y), max(y)]);
26 - title(['R^2=' num2str(R2)]);
27 - pause
```





```
% %%Kernel-based linear regression  
sigma=[1:1:10, 20:10:100, 200:100:1000];  
figure  
for i=1:length(sigma)  
    K=rbf_kernel(x,x,sigma(i));  
    ha=pinv(ones(size(K))+K)*y;  
    yhat=(ones(size(K))+K)*ha;  
    R2=norm(y-yhat)^2/norm(y-mean(y))^2;  
end
```

$$\alpha = (\mathbf{1}_{N \times N} + K)^{-1} y$$

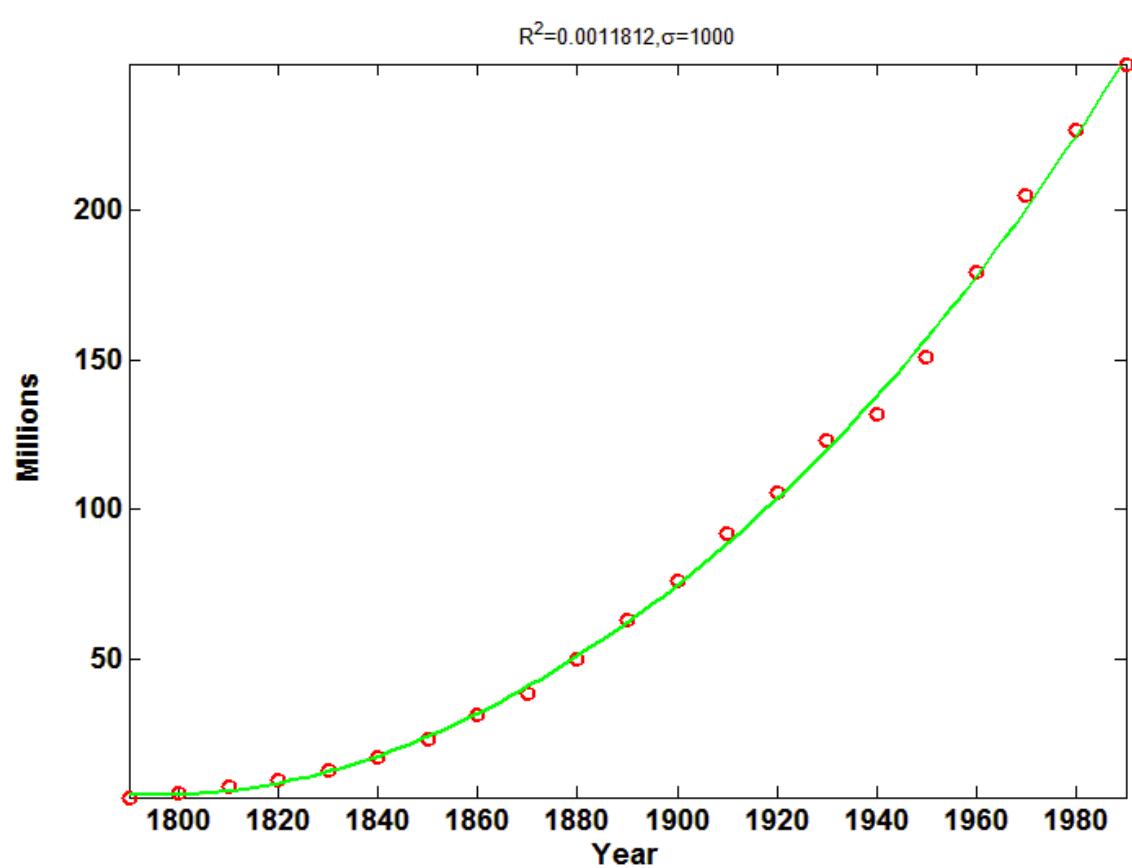
$$K_1 \alpha = y$$

$$\begin{aligned} K_1 &= X_1 X_1^T = [\mathbf{1}_{N \times 1}, X] \begin{bmatrix} \mathbf{1}_{1 \times N} \\ X^T \end{bmatrix} = \mathbf{1}_{N \times N} + XX^T \\ &= \mathbf{1}_{N \times N} + K \end{aligned}$$



```
38 - Ktest=rbf_kernel(xfit,x,sigma(i));
39 - yfit=(ones(size(Ktest))+Ktest)*ha;
40 - plot(x,y,'ro');
41 - hold on
42 - plot(xfit,yfit,'g-');
43 - xlabel('Year');
44 - ylabel('Millions');
45 - axis([min(x),max(x),min(y),max(y)]);
46 - title(['R^2=' num2str(R2) ', \sigma=' , num2str(sigma(i))]);
47 - hold off
48 - pause(0.5)
49 - end
```

$$\begin{aligned} f(x) &= \left[1, \phi(x)^T \right] \theta = \left[1, \phi(x)^T \right] X_1^T \alpha \\ &= \left[1 + \kappa(x_1, x), \dots, 1 + \kappa(x_N, x) \right] \alpha \end{aligned}$$





Principal Component Analysis (PCA)

Kernel –Based PCA (KPCA)

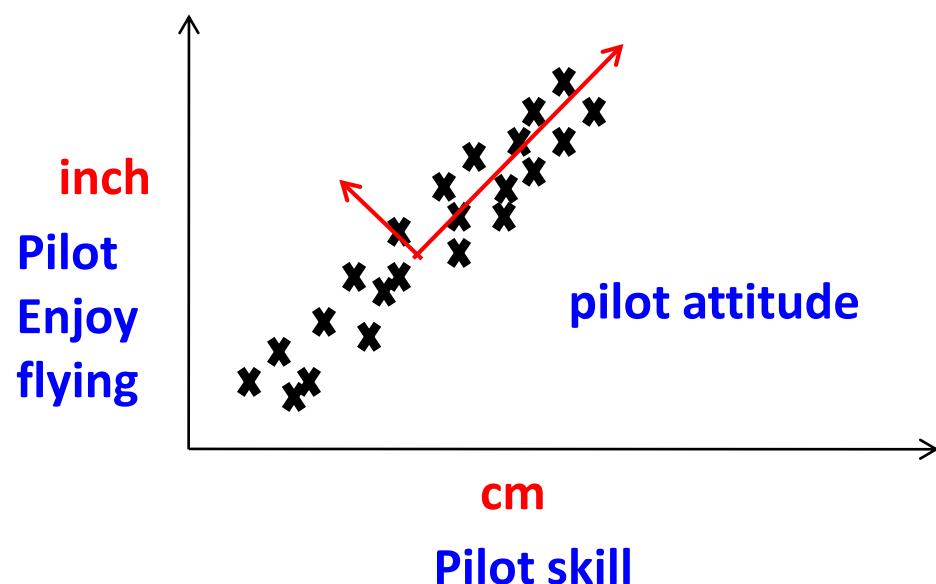


Principal Component Analysis

Given $\{x^{(1)}, \dots, x^{(m)}\}$ $\{x_1, \dots, x_N\}$

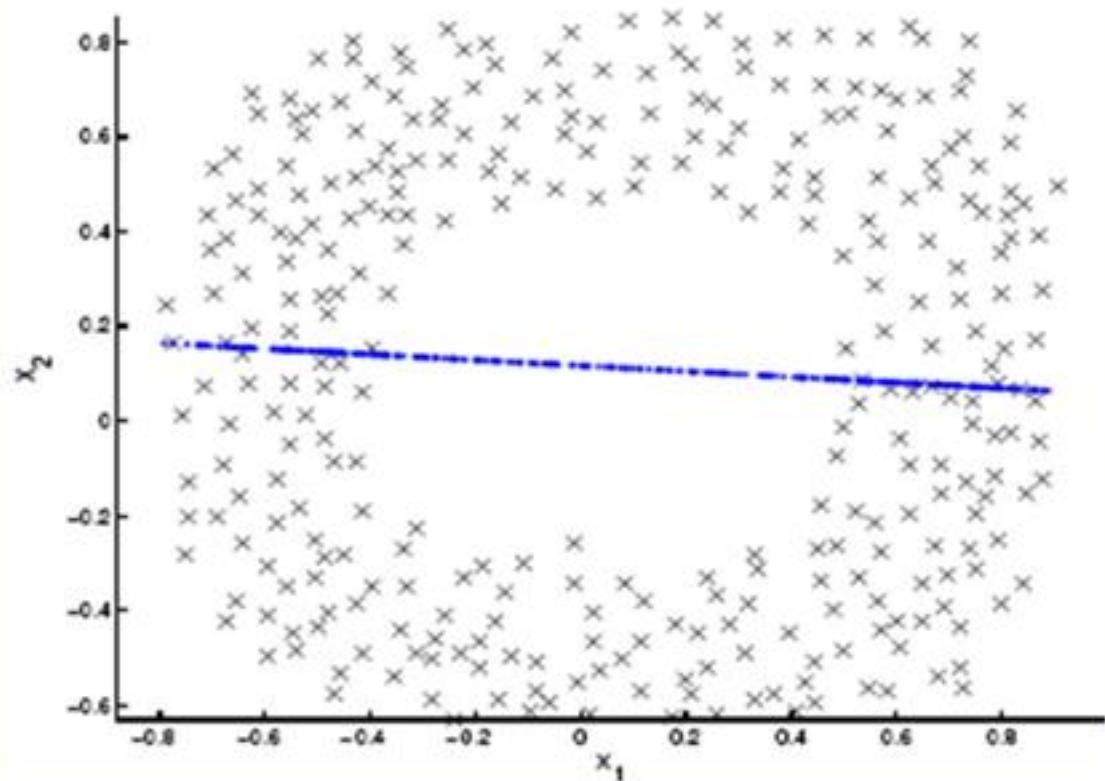
$$x^{(i)} \in R^n \quad x_i \in R^n$$

Reduce it to p -dimensional, maximum variance $(p < n)$



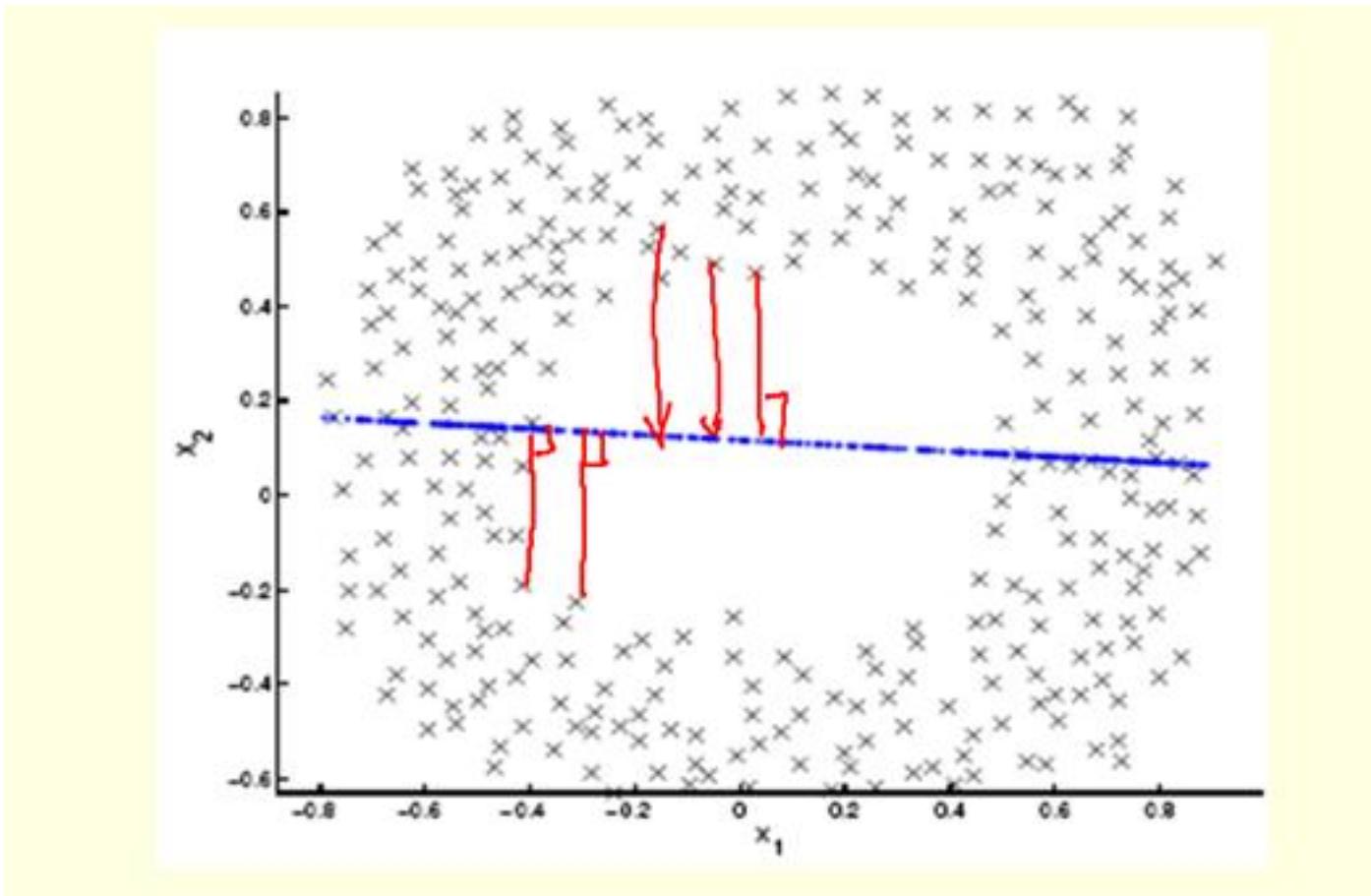


Principal Component Analysis



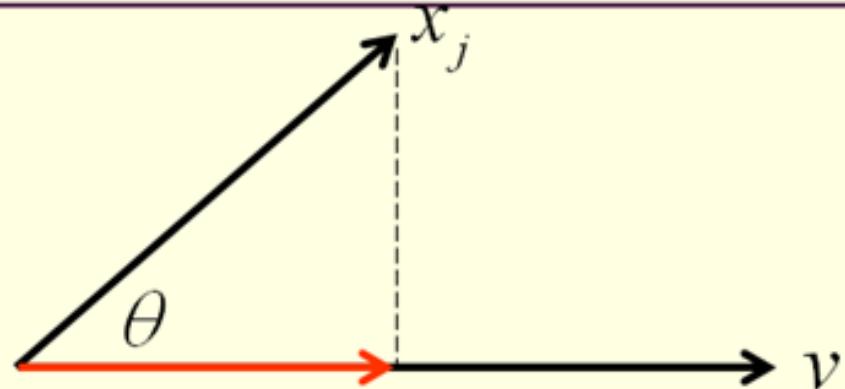


Principal Component Analysis





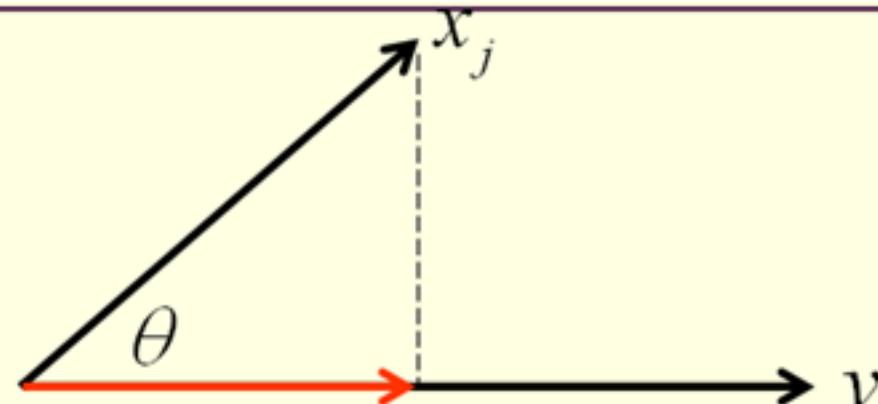
Principal Component Analysis



$$\langle x_j, v \rangle = \|x_j\| \cdot \|v\| \cos \theta \Rightarrow \begin{cases} \langle x_j, v \rangle \geq 0, & \text{if } 0 \leq \theta \leq \frac{\pi}{2} \\ \langle x_j, v \rangle < 0, & \text{if } \frac{\pi}{2} < \theta \leq \pi \end{cases}$$



Principal Component Analysis



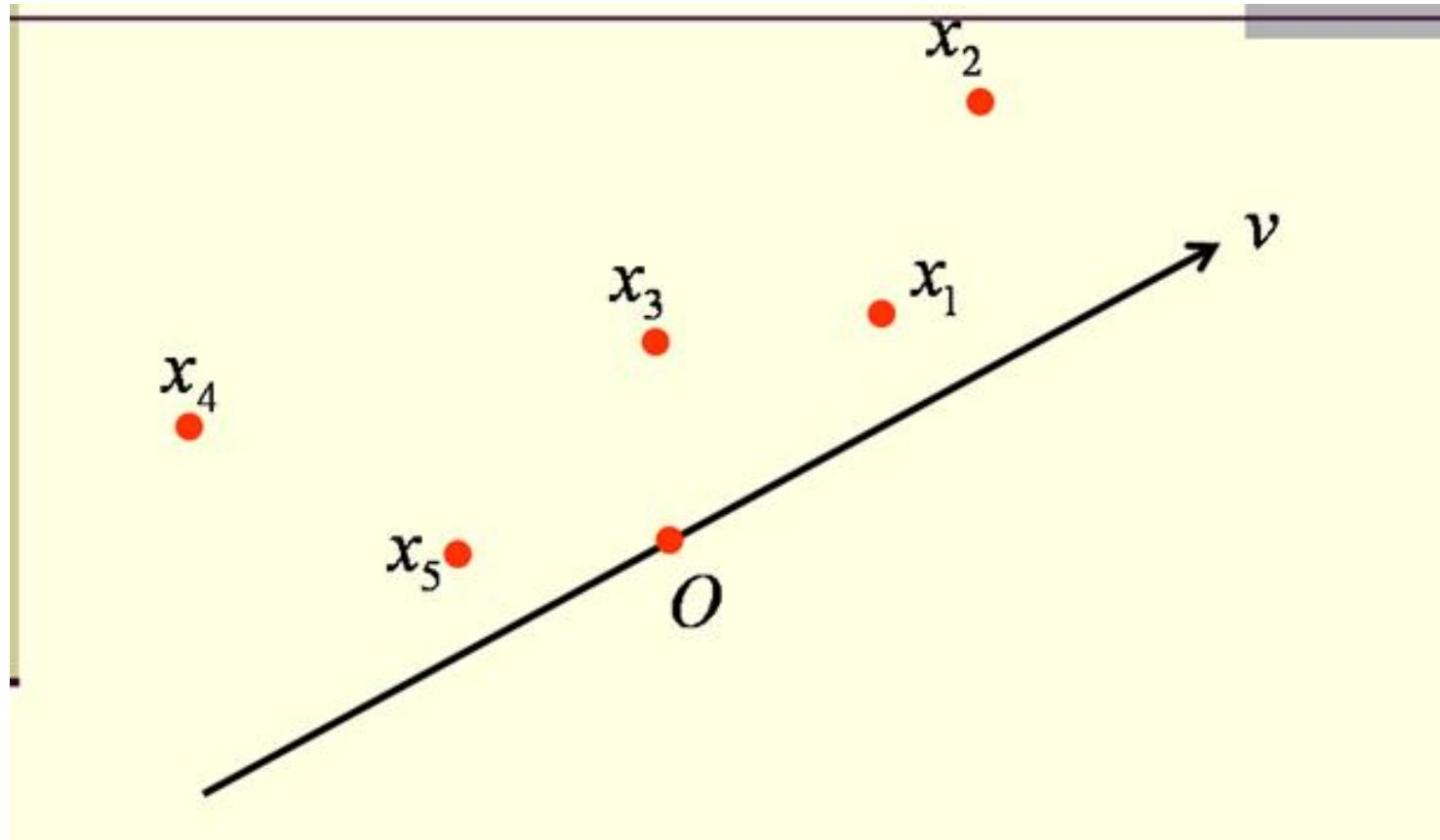
$$(\|x_j\| \cos \theta) \frac{v}{\|v\|} = \|x_j\| \frac{\langle x_j, v \rangle}{\|x_j\| \cdot \|v\|} \frac{v}{\|v\|} = \frac{\langle x_j, v \rangle}{\|v\|^2} v$$

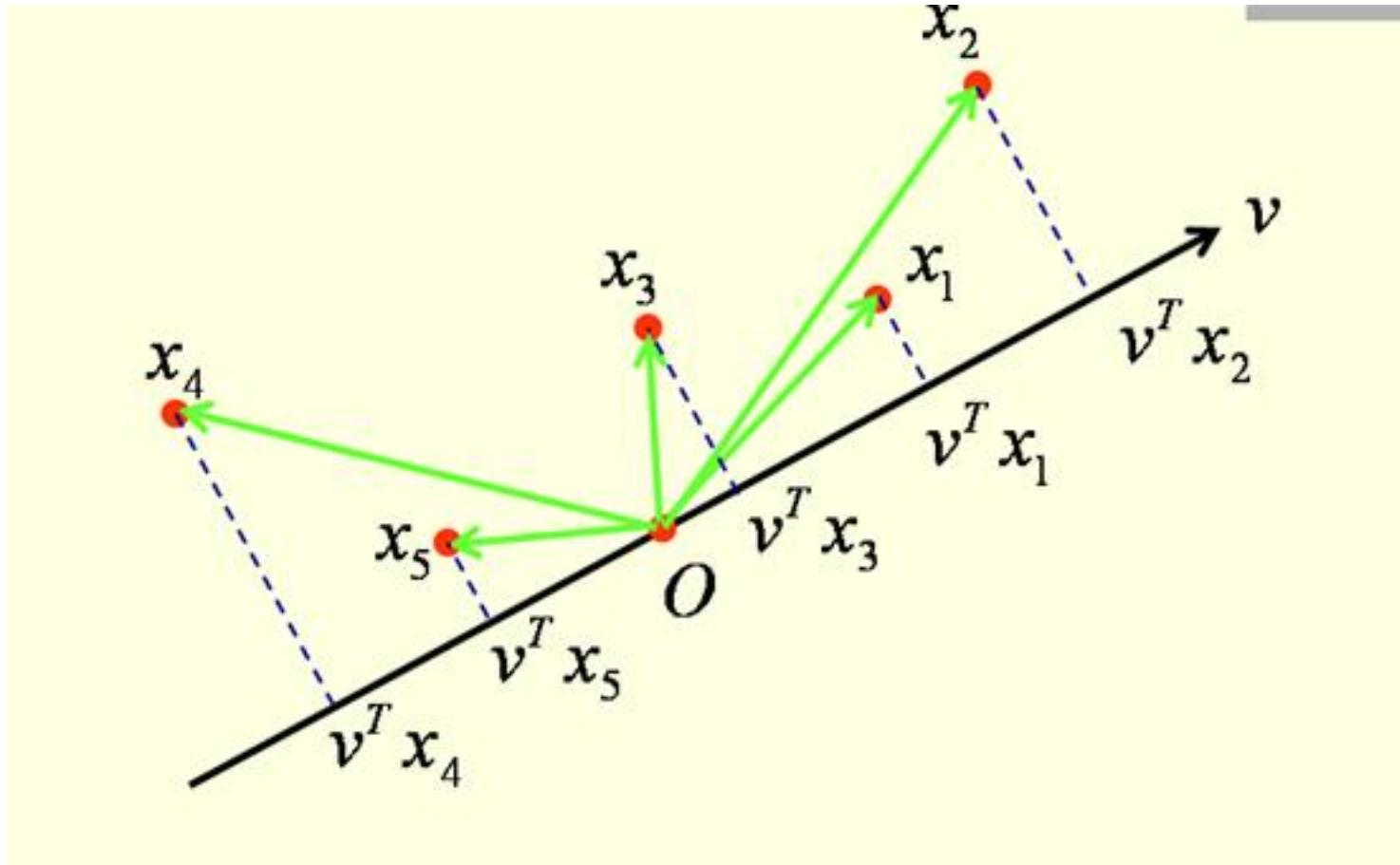
If $\|v\|=1$, then $(\|x_j\| \cos \theta) \frac{v}{\|v\|} = \langle x_j, v \rangle v$,

i.e., the coefficient is $\langle x_j, v \rangle = x_j^T v = v^T x_j$.



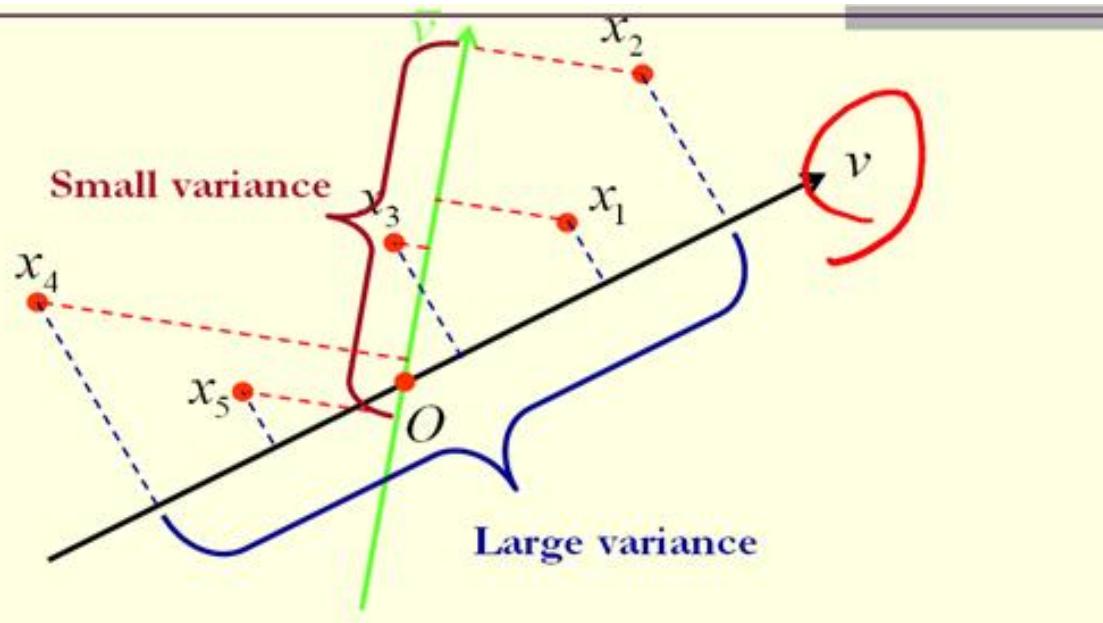
Principal Component Analysis







Principal Component Analysis



Pre-processing:

1. set $\mu = \frac{1}{N} \sum_{i=1}^N x_i$

Zero out
mean

2. Replace x_i with $x_i - \mu$



Principal Component Analysis



- The projections of the all pixels x_j onto this normalized direction v are

$$v^T x_1, \dots, v^T x_N.$$

- The variance of the projections is

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (v^T x_i - 0)^2 = \frac{1}{N} \sum_{i=1}^N (v^T x_i)(v^T x_i) = \frac{1}{N} \sum_{i=1}^N (v^T x_i)(v^T x_i)^T$$

$$= \frac{1}{N} \sum_{i=1}^N v^T x_i x_i^T v = v^T \left(\frac{1}{N} \sum_{i=1}^N x_i x_i^T \right) v = v^T C v$$

where

$$C = \frac{1}{N} \sum_{i=1}^N x_i x_i^T.$$



Principal Component Analysis

- Note that

$$C = \frac{1}{N} \sum_{i=1}^N x_i x_i^T = \frac{1}{N} [x_1, \dots, x_N] \begin{bmatrix} x_1^T \\ \vdots \\ x_N^T \end{bmatrix}.$$

- If $X^T = [x_1, \dots, x_N]$,
then

$$C = \frac{1}{N} X^T X.$$



Principal Component Analysis

- The first principal vector can be found by the following equation:

$$v = \arg \max_{v \in R^d, \|v\|=1} v^T Cv.$$



Principal Component Analysis

- The first principal vector can be found by the following equation:

$$v = \arg \max_{v \in R^d, \|v\|=1} v^T C v.$$

Lagrangian: $f(v, \lambda) = v^T C v - \lambda(v^T v - 1)$

$$\frac{\partial f}{\partial v} = 2Cv - 2\lambda v = 0 \Rightarrow Cv = \lambda v$$

$$\frac{\partial f}{\partial \lambda} = v^T v - 1 = 0 \Rightarrow v^T v = 1$$

$$\begin{aligned} v^T C v &= v^T \lambda v \\ &= \lambda v^T v = \lambda = \sigma^2 \end{aligned}$$



Principal Component Analysis

- The first principal vector can be found by the following equation:

$$v = \arg \max_{v \in R^d, \|v\|=1} v^T Cv.$$

- This is equivalent to find the largest eigenvalue of the following eigenvalue problem:

$$\begin{cases} Cv = \lambda v \\ \|v\|=1 \end{cases}$$



Principal Component Analysis

Steps

Step.1 Pre-processing:

1. set $\mu = \frac{1}{N} \sum_{i=1}^N x_i$
 2. Replace x_i with $x_i - \mu$
- } Zero out mean



Step.2 Calculate matrix (correlation coefficient matrix)

$$C = \frac{1}{N} \sum_{i=1}^N x_i x_i^T$$



Principal Component Analysis



Step.3

Calculate eigenvalues and eigenvectors :

$$Cv = \lambda v$$

$$|\lambda I - C| = 0$$

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$$



Principal Component Analysis



Step.4

计算主成分贡献率及累计贡献率

▲ 贡献率: $\frac{\lambda_i}{\sum_{i=1}^n \lambda_i} \quad (i = 1, 2, \dots, n)$

▲ 累计贡献率:

$$\frac{\sum_{j=1}^i \lambda_j}{\sum_{j=1}^n \lambda_j} \quad (i = 1, 2, \dots, k)$$

一般取累计贡献率达85—95%的特征值
所对应的第一、第二、...、第k ($k \leq n$) 个主成分。



Principal Component Analysis

Now x_i can be represented as

$$x_i = \sum_{j=1}^n (x_i^T v_j) v_j$$

x_i can be also approximated by

$$x_i \approx \sum_{j=1}^k (x_i^T v_j) v_j$$





Principal Component Analysis



-Visualization

-Compression

-Anomaly detection

-Matching/Distance calculations



Principal Component Analysis



Training set

STANFORD



Principal Component Analysis



Principal components learned (eigenfaces)

STANF



Principal Component Analysis



三、主成分分析方法应用实例（一）

下面，我们根据表1给出的数据，对某农业生态经济系统做主成分分析：

表1 某农业生态经济系统各区域单元的有关数据

样本序号	x_1 : 人口密度 (人/km ²)	x_2 : 人均耕地面积 (ha)	x_3 : 森林覆盖率 (%)	x_4 : 农民人均纯收入(元/人)	x_5 : 人均粮食产量 (kg/人)	x_6 : 经济作物占农作物播种面 比例(%)	x_7 : 耕地占土地面 积比率(%)	x_8 : 果园与林 地面积之比(%)	x_9 : 灌溉田占耕地 面积之比(%)
1	363.912	0.352	16.101	192.11	295.34	26.724	18.492	2.231	26.262
2	141.503	1.684	24.301	1752.35	452.26	32.314	14.464	1.455	27.066
3	100.695	1.067	65.601	1181.54	270.12	18.266	0.162	7.474	12.489
4	143.739	1.336	33.205	1436.12	354.26	17.486	11.805	1.892	17.534
5	131.412	1.623	16.607	1405.09	586.59	40.683	14.401	0.303	22.932



Principal Component Analysis

6	68.337	2.032	76.204	1540.29	216.39	8.128	4.065	0.011	4.861
7	95.416	0.801	71.106	926.35	291.52	8.135	4.063	0.012	4.862
8	62.901	1.652	73.307	1501.24	225.25	18.352	2.645	0.034	3.201
9	86.624	0.841	68.904	897.36	196.37	16.861	5.176	0.055	6.167
10	91.394	0.812	66.502	911.24	226.51	18.279	5.643	0.076	4.477
11	76.912	0.858	50.302	103.52	217.09	19.793	4.881	0.001	6.165
12	51.274	1.041	64.609	968.33	181.38	4.005	4.066	0.015	5.402
13	68.831	0.836	62.804	957.14	194.04	9.11	4.484	0.002	5.79
14	77.301	0.623	60.102	824.37	188.09	19.409	5.721	5.055	8.413
15	76.948	1.022	68.001	1255.42	211.55	11.102	3.133	0.01	3.425
16	99.265	0.654	60.702	1251.03	220.91	4.383	4.615	0.011	5.593
17	118.505	0.661	63.304	1246.47	242.16	10.706	6.053	0.154	8.701
18	141.473	0.737	54.206	814.21	193.46	11.419	6.442	0.012	12.945
19	137.761	0.598	55.901	1124.05	228.44	9.521	7.881	0.069	12.654
20	117.612	1.245	54.503	805.67	175.23	18.106	5.789	0.048	8.461
21	122.781	0.731	49.102	1313.11	236.29	26.724	7.162	0.092	10.078



Principal Component Analysis

步骤如下：（1）将表1中的数据作标准差标准化处理，然后将它们代入公式计算相关系数矩阵（见表2）。

2 相关系数矩阵

	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9
x_1	1	-0.327	-0.714	-0.336	0.309	0.408	0.79	0.156	0.744
x_2	-0.33	1	-0.035	0.644	0.42	0.255	0.009	-0.078	0.094
x_3	-0.71	-0.035	1	0.07	-0.74	-0.755	-0.93	-0.109	-0.924
x_4	-0.34	0.644	0.07	1	0.383	0.069	-0.05	-0.031	0.073
x_5	0.309	0.42	-0.74	0.383	1	0.734	0.672	0.098	0.747
x_6	0.408	0.255	-0.755	0.069	0.734	1	0.658	0.222	0.707
x_7	0.79	0.009	-0.93	-0.046	0.672	0.658	1	-0.03	0.89
x_8	0.156	-0.078	-0.109	-0.031	0.098	0.222	-0.03	1	0.29
x_9	0.744	0.094	-0.924	0.073	0.747	0.707	0.89	0.29	1



Principal Component Analysis



(2) 由相关系数矩阵计算特征值，以及各个主成分的贡献率与累计贡献率（见表3）。



Principal Component Analysis

表3. 特征值及主成分贡献率

主成分	特征值	贡献率(%)	累积贡献率(%)
z_1	4.661	51.791	51.791
z_2	2.089	23.216	75.007
z_3	1.043	11.589	86.596
z_4	0.507	5.638	92.234
z_5	0.315	3.502	95.736
z_6	0.193	2.14	97.876
z_7	0.114	1.271	99.147
z_8	0.0453	0.504	99.65
z_9	0.0315	0.35	100

由表3可知，第一，第二，第三主成分的累计贡献率已高达**86.596%**（大于85%），故只需要求出第一、第二、第三主成分 z_1 , z_2 , z_3 即可。



Principal Component Analysis



(3) 对于特征值=4. 6610, =2. 0890, =1. 0430分别求出其特征向量 z_1 , z_2 , z_3 , 计算各变量 x_1 , x_2 , ..., x_9 在主成分 z_1 , z_2 , z_3 上的投影。



Principal Component Analysis

	1	2	3	4
1 '变量'	'Z3'	'Z2'	'Z1'	
2 'x1:人口密度'	-0.0599	0.3679	0.3421	
3 'x2:人均耕地面积'	-0.0276	-0.6135	0.0572	
4 'x3:森林覆盖率'	0.0929	-0.0661	-0.4464	
5 'x4:人均纯收入'	0.0362	-0.6006	0.0193	
6 'x5:人均粮食产量'	-0.0107	-0.3068	0.3765	
7 'x6:经济作物所占比例'	0.1222	-0.1241	0.3793	
8 'x7:耕地占土地面积比率'	-0.2461	0.0920	0.4322	
9 'x8:果园与林地面积之比'	0.9497	0.0695	0.0914	
10 'x9:灌溉田所占比率'	0.0898	0.0173	0.4464	
11				



Principal Component Analysis



分析：

- ① 第一主成分 z_1 与 x_1, x_5, x_6, x_7, x_9 呈显出较强的正相关，与 x_3 呈显出较强的负相关，而这几个变量则综合反映了生态经济结构状况，因此可以认为第一主成分 z_1 是生态经济结构的代表。
- ② 第二主成分 z_2 与 x_2, x_4, x_5 呈显出较强的负相关，与 x_1 呈显出较强的正相关，其中，除了 x_1 为人口总数外， x_2, x_4, x_5 都反映了人均占有资源量的情况，因此可以认为第二主成分 z_2 代表了人均资源量。



Principal Component Analysis

③第三主成分 z_3 ，与 x_8 呈显出的正相关程度最高，其次是 x_6 ，而与 x_7 呈负相关，因此可以认为第三主成分在一定程度上代表了经济作物和果园对农业经济的影响。

显然，用三个主成分 z_1 、 z_2 、 z_3 代替原来9个变量 (x_1, x_2, \dots, x_9) ，描述农业生态经济系统，可以使问题更进一步简化、明了。



Principal Component Analysis

data - 记事本

文件(F) 编辑(E) 格式(O) 查看(V) 帮助(H)

```
363.912,0.352,16.101,192.11,295.34,26.724,18.492,2.231,26.262
141.503,1.684,24.301,1752.35,452.26,32.314,14.464,1.455,27.066
100.695,1.067,65.601,1181.54,270.12,18.266,0.162,7.474,12.489
143.739,1.336,33.205,1436.12,354.26,17.486,11.805,1.892,17.534
131.412,1.623,16.607,1405.09,586.59,40.683,14.401,0.303,22.932
68.337,2.032,76.204,1540.29,216.39,8.128,4.065,0.011,4.861
95.416,0.801,71.106,926.35,291.52,8.135,4.063,0.012,4.862
62.901,1.652,73.307,1501.24,225.25,18.352,2.645,0.034,3.201
86.624,0.841,68.904,897.36,196.37,16.861,5.176,0.055,6.167
91.394,0.812,66.502,911.24,226.51,18.279,5.643,0.076,4.477
76.912,0.858,50.302,103.52,217.09,19.793,4.881,0.001,6.165
51.274,1.041,64.609,968.33,181.38,4.005,4.066,0.015,5.402
68.831,0.836,62.804,957.14,194.04,9.11,4.484,0.002,5.79
77.301,0.623,60.102,824.37,188.09,19.409,5.721,5.055,8.413
76.948,1.022,68.001,1255.42,211.55,11.102,3.133,0.01,3.425
99.265,0.654,60.702,1251.03,220.91,4.383,4.615,0.011,5.593
118.505,0.661,63.304,1246.47,242.16,10.706,6.053,0.154,8.701
141.473,0.737,54.206,814.21,193.46,11.419,6.442,0.012,12.945
137.761,0.598,55.901,1124.05,228.44,9.521,7.881,0.069,12.654
117.612,1.245,54.503,805.67,175.23,18.106,5.789,0.048,8.461
122.781,0.731,49.102,1313.11,236.29,26.724,7.162,0.092,10.078
```



Principal Component Analysis

```
clc
clear all
close all
data=load('data.txt');%读入数据，即表3.4.5中的数据
[Z,MU,SIGMA] = zscore(data);%将读入的数据做标准化处理
covMat = cov(Z);%计算相关系数矩阵
[eigVect,eigVal] = eig(covMat);%计算相关系数矩阵的特征值和特征向量
d=diag(eigVal);%提取特征值矩阵的对角线元素，即特征值
D=sort(d,'descend');%将各个特征值进行降序排列
rat1=D./sum(D);%计算贡献率
rat2=cumsum(D)./sum(D);%计算累积贡献率
```



Principal Component Analysis



%%导出特征值，贡献率及累积贡献率%%

```
result1(1, :)={‘特征值’, ‘贡献率’, ‘累积贡献率’};
```

```
result1(2:10, 1)=num2cell(D);
```

```
result1(2:10, 2)=num2cell(rat1);
```

```
result1(2:10, 3)=num2cell(rat2);
```

%%主成分载荷%%

```
result2(:, 1)={‘变量’; ‘x1:人口密度’; ‘x2:人均耕地面积’; ‘x3:森林覆盖率’;
```

```
‘x4:人均纯收入’; ‘x5:人均粮食产量’; ‘x6:经济作物所占比例’;
```

```
‘x7:耕地占土地面积比率’; ‘x8:果园与林地面积之比’; ‘x9:灌溉田所占比率’};
```

```
result2(1, 2:4)={‘z3’, ‘z2’, ‘z1’}; %由累积贡献率可知，只需分析三个主成分，z1, z2, z3
```

```
result2(2:10, 2:4)=num2cell(eigVect (:, 7:9)); %z1, z2, z3分别为特征值4.6611, 2.0895, 1.0430对应的特征向量
```



Principal Component Analysis



三、主成分分析方法应用实例（二）

128名男子身材的六项指标的样本相关系数矩阵

变量	身 高 (x_1)	座 高 (x_2)	胸 围 (x_3)	手 臂 长 (x_4)	肋 围 (x_5)	腰 围 (x_6)
身高 (x_1)	1	0.79	0.36	0.76	0.25	0.51
座高 (x_2)	0.79	1	0.31	0.55	0.17	0.35
胸围 (x_3)	0.36	0.31	1	0.35	0.64	0.58
手臂长 (x_4)	0.76	0.55	0.35	1	0.64	0.38
肋围 (x_5)	0.25	0.17	0.64	0.16	1	0.63
腰围 (x_6)	0.51	0.35	0.58	0.38	0.63	1



Principal Component Analysis



相关系数矩阵：

pho =

1.0000	0.7900	0.3600	0.7600	0.2500	0.5100
0.7900	1.0000	0.3100	0.5500	0.1700	0.3500
0.3600	0.3100	1.0000	0.3500	0.6400	0.5800
0.7600	0.5500	0.3500	1.0000	0.1600	0.3800
0.2500	0.1700	0.6400	0.1600	1.0000	0.6300
0.5100	0.3500	0.5800	0.3800	0.6300	1.0000



Principal Component Analysis

调用pcacov函数作主成分分析：

```
>> [COEFF, latent, explained]=pcacov(pho)
```

COEFF =					
-0.4689	-0.3648	0.0922	-0.1224	-0.0797	0.7856
-0.4037	-0.3966	0.6130	0.3264	0.0270	-0.4434
-0.3936	0.3968	-0.2789	0.6557	0.4052	0.1253
-0.4076	-0.3648	-0.7048	-0.1078	-0.2346	-0.3706
-0.3375	0.5692	0.1643	-0.0193	-0.7305	-0.0335
-0.4268	0.3084	0.1193	-0.6607	0.4899	-0.1788

特征值对应的特征向量

latent =	explained =	>> cumsum(explained)
3.2872	54.7867	ans =
1.4062	23.4373	54.7867
0.4591	7.6516	78.2240
0.4263	7.1057	85.8756
0.2948	4.9133	92.9813
0.1263	2.1054	97.8946
		100.0000

特征值

贡献率

累积贡献率



Principal Component Analysis

从上述结果来看，前3个主成分的累积贡献率达到了**85.8756%**，因此可以用前3个主成分进行后续的分析；这样做虽然会有一定的信息损失，但损失不大，不影响大局。

	z1	z2	z3
x1	-0.4689	-0.3648	0.0922
x2	-0.4037	-0.3966	0.6130
x3	-0.3936	0.3968	-0.2789
x4	-0.4076	-0.3648	-0.7048
x5	-0.3375	0.5692	0.1643
x6	-0.4268	0.3084	0.1193



Principal Component Analysis



结果分析：

1. 从第一主成分 z_1 来看，它在每个变量上有近似的负载荷，说明每个变量对 z_1 的重要性都差不多。当一个人的身材“五大三粗”，也就是说又高又胖时， x_1, x_2, \dots, x_6 都比较大，此时 z_1 的值就比较小；反之，当一个人又矮又瘦时 x_1, x_2, \dots, x_6 都比较小，此时 z_1 的值就比较大，所以可以认为第一主成分是身材的综合成分（或魁梧成分）。
2. 从第二主成分 z_2 来看，它在变量 x_1, x_2 和 x_4 上有近似的负载荷，在 x_3, x_5 和 x_6 上有近似的正载荷，说明当 x_1, x_2 和 x_4 增大时， z_2 的值减小，当 x_3, x_5 和 x_6 增大时， z_2 的值增大。当一个人的身材瘦高时， z_2 的值比较小；当一个人的身材矮胖时， z_2 的值比较大，所以可以认为第二主成分是身材的高矮和胖瘦的协调成分。
3. 从第三主成分 z_3 的表达式来看，它在变量 x_2 上有比较大的正载荷，在 x_4 上有比较大的负载荷，在其它变量上的载荷比较小，说明 x_2 （坐高）和 x_4 （手臂长）对 z_3 的影响比较大，也就是说 z_3 反映了坐高（即上半身）与手臂长之间的协调关系，这对做长袖上衣时制定衣服和袖子的长短提供了参考；所以可认为第三主成分 z_3 是臂长成分。



Principal Component Analysis



featureNormalize

%FeatureNormalize.m Normalizes the features in X

```
mu = mean(X);  
X_norm = bsxfun(@minus, X, mu);
```

```
sigma = std(X_norm);  
X_norm = bsxfun(@rdivide, X_norm, sigma);
```

Standard deviation

corr(X,Y)



Principal Component Analysis

A = magic(5)

A =

```
17  24   1   8   15
23   5   7  14  16
 4   6  13  20  22
10  12  19  21   3
11  18  25   2   9
```

mean(A)

ans =

```
13  13  13  13  13
```

X_norm = bsxfun(@rdivide, X_norm, sigma)'

X_norm = bsxfun(@minus, A, mean(A))

X_norm =

```
 4   11  -12  -5   2
 10  -8   -6   1   3
 -9  -7    0   7   9
 -3  -1    6   8  -10
 -2   5   12  -11  -4
```

sigma = std(X_norm)'

sigma =

```
7.245688373094719e+000
8.062257748298549e+000
9.486832980505138e+000
8.062257748298549e+000
7.245688373094719e+000
```



Principal Component Analysis (PCA) algorithm

Reduce data from n -dimensions to k -dimensions

Compute “covariance matrix”:

$$\sum_{ij} = \frac{1}{M-1} X_i^T X_j \quad \sum = \frac{1}{M-1} X^T X$$

$\text{cov}(X)$



Compute “eigenvectors” of matrix Σ :

`[U, S, V] = svd(Sigma);`

Covariance matrix always
Satisfy “symmetric semi-
positive Definite.

$$U = \begin{bmatrix} u^{(1)} & u^{(2)} & \dots & u^{(n)} \\ | & | & & | \end{bmatrix} \in \mathbb{R}^{n \times n}$$

The first k columns $u^{(1)}$ to $u^{(k)}$ are the k dimensional vectors
that we want

$$U_{reduced} = \begin{bmatrix} \vdots & \vdots & & \vdots \\ u^{(1)} & u^{(2)} & \dots & u^{(k)} \\ \vdots & \vdots & & \vdots \\ \vdots & \vdots & & \vdots \end{bmatrix}$$



$Z \in R^k$

$$Z^{(i)} = X^{(i)} U_{reduced}$$

Choosing k (number of principal components)

Average squared projection error:

Total variation in the data:

Typically, choose k to be smallest value so that

$$\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{approx}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq 0.01 \quad (1\%)$$

“99% of variance is retained”



90%~99%



Choosing k (number of principal components)

[$\mathbf{U}, \mathbf{S}, \mathbf{V}$] = $\text{svd}(\Sigma)$

Pick smallest value of k for which

$$\frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \geq 0.99$$

(99% of variance retained)

QUESTION

0
0

$$\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{approx}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2}$$



$$1 - \frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}}$$





Kernel-based PCA



Note that

$$C = \frac{1}{N} \sum_{i=1}^N \phi(x_i) \phi(x_i)^T = \frac{1}{N} [\phi(x_1), \dots, \phi(x_N)] \begin{bmatrix} \phi(x_1)^T \\ \vdots \\ \phi(x_N)^T \end{bmatrix}$$

If $X^T = [\phi(x_1), \dots, \phi(x_N)]$,

then

$$C = \frac{1}{N} X^T X$$



Kernel-based PCA

Note that the kernel matrix

$$\begin{aligned} K &= \begin{bmatrix} \kappa(x_1, x_1) & \cdots & \kappa(x_1, x_N) \\ \vdots & \ddots & \vdots \\ \kappa(x_N, x_1) & \cdots & \kappa(x_N, x_N) \end{bmatrix} = \\ &\quad \begin{bmatrix} \phi(x_1)^T \phi(x_1) & \cdots & \phi(x_1)^T \phi(x_N) \\ \vdots & \ddots & \vdots \\ \phi(x_N)^T \phi(x_1) & \cdots & \phi(x_N)^T \phi(x_N) \end{bmatrix} \\ &= \begin{bmatrix} \phi(x_1)^T \\ \vdots \\ \phi(x_N)^T \end{bmatrix} [\phi(x_1), \dots, \phi(x_N)] = XX^T \end{aligned}$$

Note that $X^T X \neq XX^T = K$



Kernel-based PCA



The eigenvalue problem of $K = XX^T$ is

$$(XX^T)u = \beta u.$$

Multiply both sides by X^T :

$$X^T(XX^T)u = \beta X^T u \Rightarrow C(X^T u) = \beta(X^T u).$$

This means that

$$X^T u$$

is a eigenvectors of C

However its norm is not sure 1.



Kernel-based PCA

$$\begin{aligned}v &= \frac{1}{\|X^T u\|} X^T u = \frac{1}{\sqrt{u^T X X^T u}} X^T u \\&= \frac{1}{\sqrt{u^T (\beta u)}} X^T u = \frac{1}{\sqrt{\beta}} X^T u,\end{aligned}$$

Where β is the corresponding eigenvalue of K .



Kernel-based PCA



The projection of the testing sample $\phi(x')$ on the i -th eigenvector can be computed by

$$\begin{aligned} v_i^T \phi(x') &= \left(\frac{1}{\sqrt{\beta_i}} X^T u_i \right)^T \phi(x') = \frac{1}{\sqrt{\beta_i}} u_i^T X \phi(x') \\ &= \frac{1}{\sqrt{\beta_i}} u_i^T \begin{bmatrix} \phi(x_1)^T \\ \vdots \\ \phi(x_N)^T \end{bmatrix} \phi(x') = \frac{1}{\sqrt{\beta_i}} u_i^T \begin{bmatrix} \kappa(x_1, x') \\ \vdots \\ \kappa(x_N, x') \end{bmatrix}. \end{aligned}$$



Summarization of Kernel-based PCA

□ Solve the following eigenvalue problem:

$$Ku_i = \beta_i u_i, \quad \beta_1 \geq \beta_2 \geq \dots \geq \beta_k$$

□ The projection of the testing sample $\phi(x')$ on the i-th eigenvector can be computed by

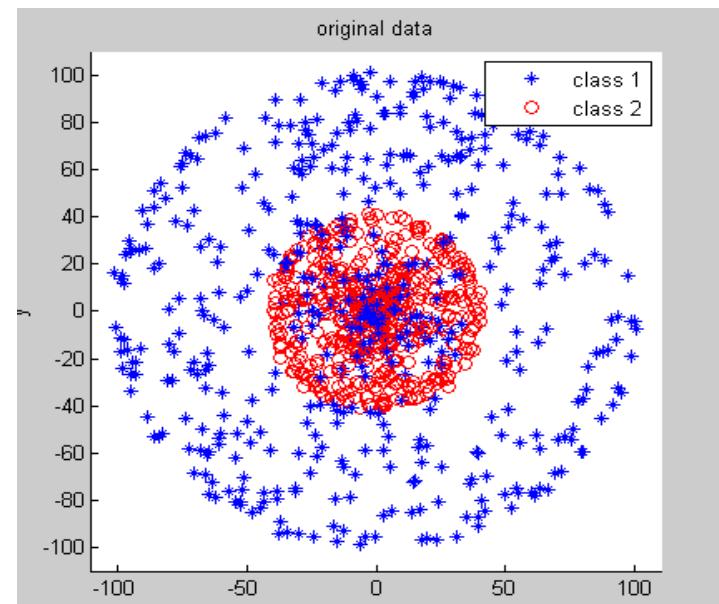
$$v_i^T \phi(x') = \frac{1}{\sqrt{\beta_i}} u_i^T \begin{bmatrix} k(x_1, x') \\ \vdots \\ k(x_k, x') \end{bmatrix}.$$



Pattern Classification for Synthetic Data

3.1.1. DATA DESCRIPTION

We assume that we have equal numbers of data points distributed on two concentric sphere surfaces. If N is the total number of all data points, then we have $N/2$ class 1 points on a sphere of radius r_1 , and $N/2$ class 2 points on a sphere of radius r_2 . In the spherical coordinate system, the inclination (polar angle) θ is uniformly distributed in $[0, \pi]$, and the azimuth (azimuthal angle) ϕ is uniformly distributed in $[0, 2\pi]$ for both classes. Our observations of the data points are the (x, y, z) coordinates in the Cartesian coordinate system, and all the three coordinates are perturbed by a Gaussian noise of standard deviation σ_{noise} . We set $N = 1000$, $r_1 = 40$, $r_2 = 100$, $\sigma_{\text{noise}} = 1$, and give a 3D plot of the data in Figure 1.





原始数据为1000*3的矩阵，表示两类不同的点。

为节约篇幅，从中截取部分表示如下：

	1	2	3
1	-88.2783	-44.9607	7.0588
2	18.7660	12.1832	33.9881
3	85.6706	51.3761	-3.8543
4	2.8330	-1.9369	38.9479
5	24.5787	-5.8881	-95.7886
6	16.4186	34.3755	13.1516
7	2.1680	-2.9920	-99.4955
8	-16.8696	-11.3769	-34.3947
9	-84.0856	20.0495	-52.8139
10	-0.2452	0.8573	-39.3020
11	-2.7936	46.4088	-87.0521
12	-1.3023	38.2229	-1.4955
13	-4.2918	-0.4123	-98.7998
14	-16.8444	8.2598	-36.8473
15	85.4805	-17.6086	-48.4263
16	20.1288	-27.8644	21.7843
17	1.0819	-0.5867	-100.5959
18	15.7314	20.3911	29.5500
19	-8.1825	83.1212	54.0508
20	2.8475	-30.0695	24.6174
21	4.1690	-12.7177	97.3090
22	4.0096	-11.5281	-39.9845
23	39.0154	91.4269	-4.3496
24	27.8044	27.4301	-9.7474
25	-64.7173	-59.3353	-49.3511
26	-2.2166	-39.8347	1.6756
27	65.6895	-7.2869	-75.4613
28	-38.5712	-9.8804	3.5479
29	54.4305	-8.6069	83.8852
30	-3.9201	-10.1438	-38.8079
31	28.0910	63.9337	72.2697
32	16.3860	-36.3944	5.0192

	1	2	3
32	16.3860	-36.3944	5.0192
33	-100.6940	-6.4452	-1.0960
34	-28.5941	-26.7088	7.9973
35	-29.2244	89.7650	-32.8902
36	27.9354	27.9749	-2.8863
37	-11.5172	65.3030	75.4915
38	-31.6205	-25.1088	0.8250
39	-66.5808	-1.4838	-75.4468
40	5.8847	31.1994	23.7942
41	40.9660	-80.3497	43.7461
42	25.4686	-29.5762	9.2897
43	-5.3068	-23.7835	96.8680
44	-10.0345	38.9043	3.7991
45	-3.7578	10.0168	98.6781
46	1.9958	-5.7128	40.2225
47	16.5785	-57.6323	79.1247
48	28.1717	28.1596	-1.3051
49	12.3425	19.8227	-97.7523
50	38.8445	5.2272	-3.6825
51	-97.5163	-11.7493	-20.4343
52	-34.2409	18.8793	11.8706
53	89.7897	-16.2756	40.6405
54	-32.3305	10.0440	19.9357
55	42.4553	-84.7449	-34.1944
56	0.3350	-18.3883	35.6729
57	-5.8878	-7.3094	99.8485
58	-2.8967	-13.0541	-38.9086
59	-51.2399	68.8863	-50.2910
60	-16.8810	16.1630	-32.0438
61	-25.0192	-5.8718	-95.5164
62	7.9160	-3.4730	39.9748
63	-82.7399	-51.8764	-24.1630

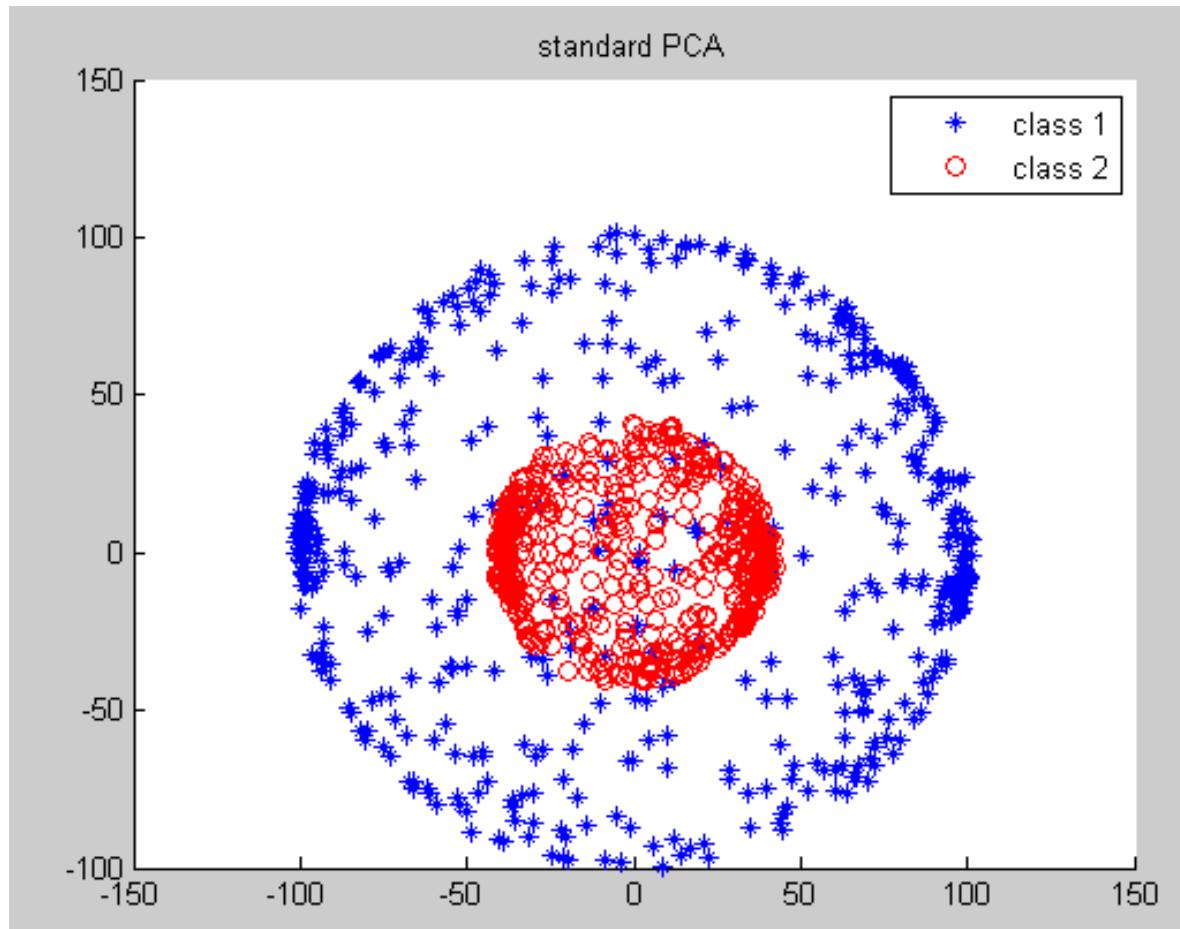
	1	2	3
64	33.7954	4.3451	-24.3365
65	-9.6023	-40.9501	-93.1028
66	-0.5545	33.6727	19.4823
67	-88.1515	37.2810	-27.5918
68	32.1920	6.7318	22.6564
69	50.1551	-77.5785	-36.1045
70	-13.1424	-14.0452	-34.9507
71	-20.3627	68.3021	72.3880
72	-11.3068	-34.1125	15.2281
73	27.6767	-12.1305	94.0013
74	8.1707	-20.6474	33.3189
75	-61.8277	-78.5864	-7.2388
76	11.6272	15.5509	36.9328
77	78.0188	-38.4314	-47.0786
78	-33.1136	22.4965	1.1917
79	56.7623	39.1851	-72.9892
80	1.8105	5.9147	-41.2094
81	-14.8765	-47.3224	-87.1255
82	-14.2931	14.6900	-33.5905
83	-22.3948	-1.5863	-96.3321
84	-0.1045	-0.6660	-40.3545
85	94.6064	-34.4596	-6.2075
86	-20.0563	32.1719	16.4020
87	-35.3545	20.5063	91.4720
88	4.4381	7.8252	39.9853
89	-39.1146	90.0463	-13.3213
90	32.8661	-2.4591	-24.8136
91	-4.0100	-58.9698	-78.8484
92	16.8515	9.3305	33.6024
93	-10.9411	35.8924	95.7406
94	-18.0028	-25.2695	-24.6476
95	-40.0947	-33.2454	83.8713



To visualize our results, we project the original 3-imensional data onto a 2-dimensional feature space by using both standard PCA and kernel PCA.

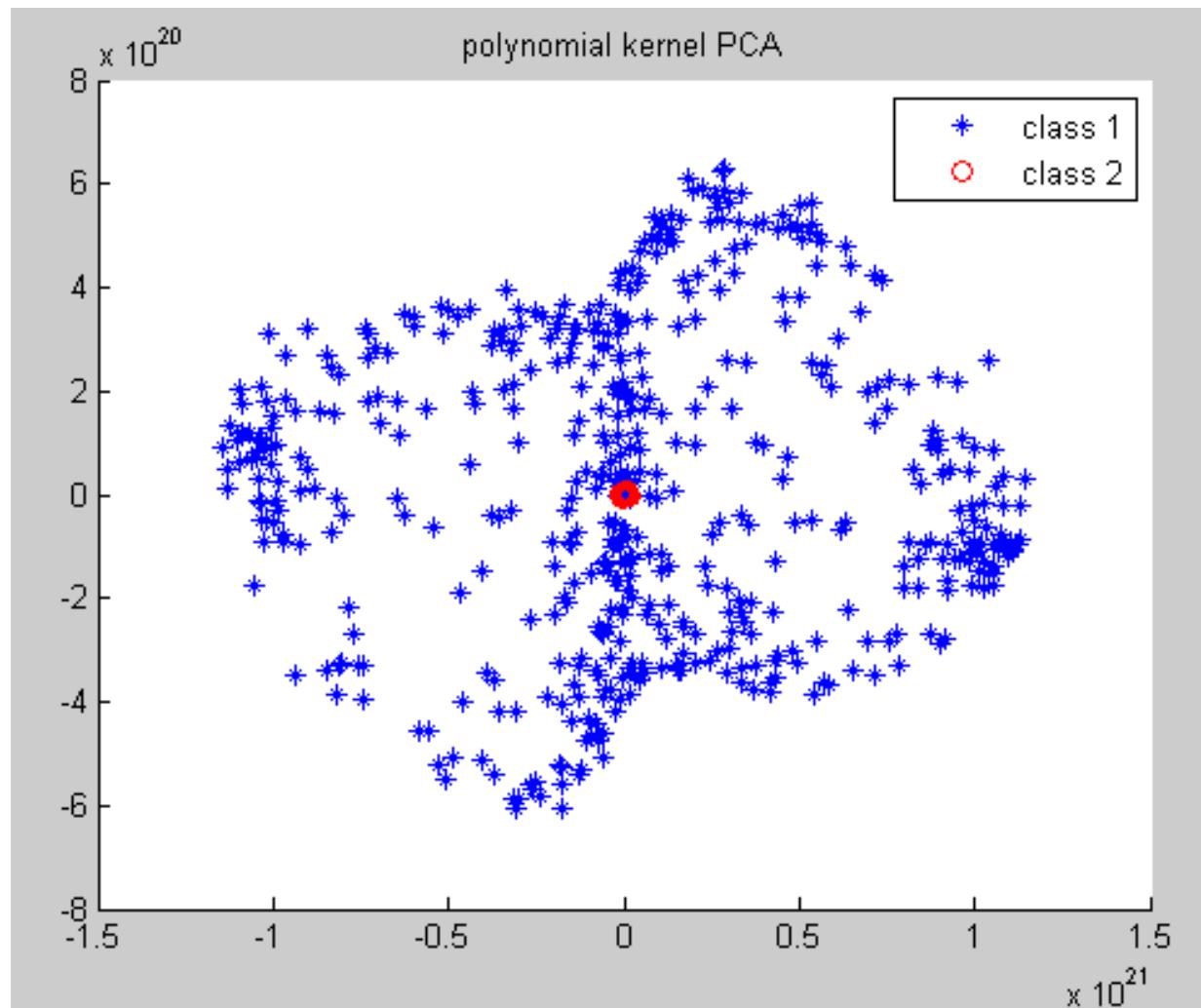


使用标准PCA将上述数据聚类的结果：



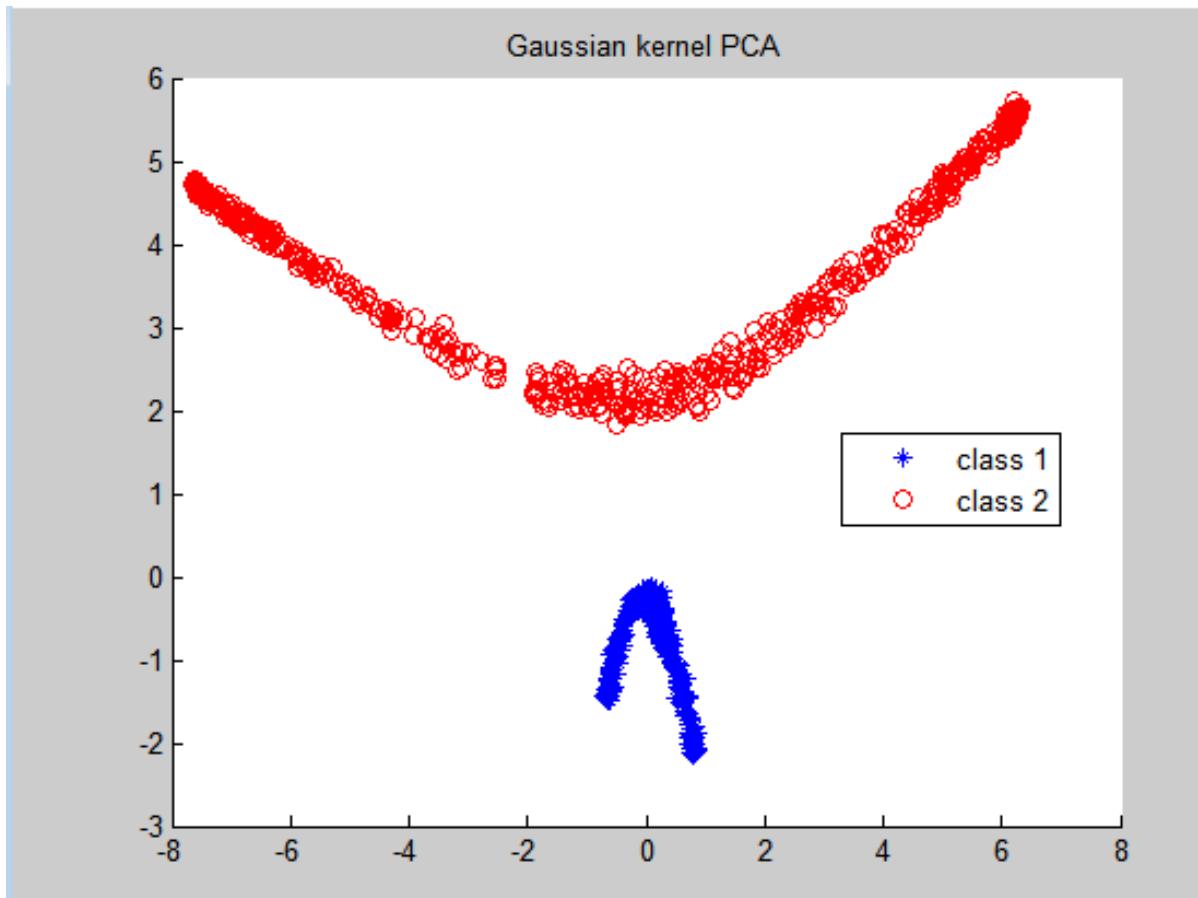


使用**polynomial kernel PCA** 将上述数据聚类的结果：





使用**Gaussian kernel PCA** 将上述数据聚类的结果：





In the resulting figures we can see that, standard PCA does not reveal any structural information of the original data.

For polynomial kernel PCA, in the new feature space, class 2 data points are clustered while class 1 data points are scattered. But they are still not linearly separable.

For Gaussian kernel PCA, the two classes are completely linearly separable.

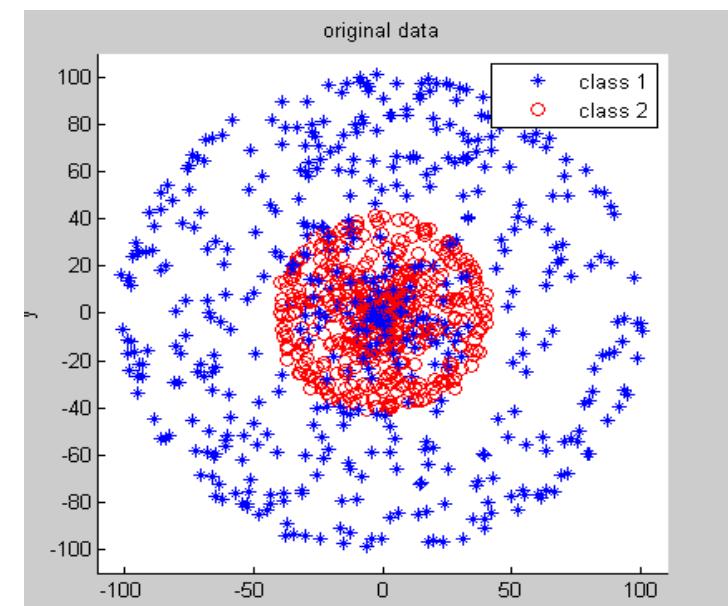


实现上述可视化效果的matlab代码

原始数据可视化：

```
%% original data
figure; hold on;
plot3(data(1:2:end, 1), data(1:2:end, 2), data(1:2:end, 3), 'b*');
plot3(data(2:2:end, 1), data(2:2:end, 2), data(2:2:end, 3), 'ro');
legend('class 1', 'class 2');
axis equal;
xlabel('x');
ylabel('y');
zlabel('z');
axis([-110 110 -110 110 -110 110]);
title('original data');
drawnow;
```

Data里1,3, ...行是蓝色数据；
2,4, ...行是红色数据





标准PCA聚类：

eig Eigenvalues and eigenvectors

```
%% standard PCA
```

```
disp('Performing standard PCA...');
```

```
Y1=PCA(data,d);
```

```
figure;hold on;
```

```
plot(Y1(1:2:end,1),Y1(1:2:end,2),'b*');
```

```
plot(Y1(2:2:end,1),Y1(2:2:end,2),'ro');
```

```
legend('class 1','class 2');
```

```
title('standard PCA');
```

```
drawnow;
```

```
function [Y, eigVector, eigValue]=PCA(X, d)
```

```
%% eigenvalue analysis
Sx=cov(X);
[V, D]=eig(Sx);
eigValue=diag(D);
[eigValue, IX]=sort(eigValue, 'descend');
eigVector=V(:, IX);
```

```
%% normalization
norm_eigVector=sqrt(sum(eigVector.^2));
eigVector=eigVector./repmat(norm_eigVector, size(eigVector, 1), 1);
```

```
%% dimensionality reduction
eigVector=eigVector(:, 1:d);
Y=X*eigVector;
```



polynomial kernel PCA聚类:

```
%% polynomial kernel PCA
para=5;
disp('Performing polynomial kernel PCA... ');
[Y2]=kPCA(data, d, 'poly', para);
figure;hold on;
plot(Y2(1:2:end, 1), Y2(1:2:end, 2), 'b*');
plot(Y2(2:2:end, 1), Y2(2:2:end, 2), 'ro');
legend('class 1', 'class 2');
title('polynomial kernel PCA');
drawnow;
```



Gaussian kernel PCA聚类:

```
%% Gaussian kernel PCA
DIST=distanceMatrix(data);
DIST(DIST==0)=inf;
DIST=min(DIST);
para=5*mean(DIST);
disp('Performing Gaussian kernel PCA... ');
[Y3, eigVector]=kPCA(data, d, 'gaussian', para);
figure;hold on;
plot(Y3(1:2:end, 1), Y3(1:2:end, 2), 'b*');
plot(Y3(2:2:end, 1), Y3(2:2:end, 2), 'ro');
legend('class 1', 'class 2');
title('Gaussian kernel PCA');
drawnow;
```



```
function [Y, eigVector, eigValue]=kPCA(X, d, type, para)
```

```
%% check input
```

```
if ( strcmp(type,'simple') || strcmp(type,'poly') || ...  
    strcmp(type,'gaussian') ) == 0
```

```
Y=[];
```

```
eigVector=[];
```

```
fprintf(['\nError: Kernel type \' type \' is not supported. \n']);
```

```
return;
```

```
end
```

```
N=size(X, 1);
```

```
%% kernel PCA
```

```
K0=kernel(X,type,para);
```

```
oneN=ones(N,N)/N;
```

```
K=K0-oneN*K0-oneN+oneN*K0*oneN;
```

```
%% eigenvalue analysis
```

```
[V,D]=eig(K/N);
```

```
eigValue=diag(D);
```

```
[~, IX]=sort(eigValue, 'descend');
```

```
eigVector=V(:,IX);
```

```
eigValue=eigValue(IX);
```

```
%% normalization
```

```
norm_eigVector=sqrt(sum(eigVector.^2));
```

```
eigVector=eigVector./repmat(norm_eigVector, size(eigVector, 1), 1);
```

```
%% dimensionality reduction
```

```
eigVector=eigVector(:,1:d);
```

```
Y=K0*eigVector;
```

X: data matrix, each row is one observation, each column is one feature
d: reduced dimension
type: type of kernel, can be 'simple', 'poly', or 'gaussian'
para: parameter for computing the 'poly' and 'gaussian' kernel,
 for 'simple' it will be ignored
Y: dimensionality-reduced data
eigVector: eigen-vector, will later be used for pre-image
 reconstruction

If the projected dataset $\{\phi(x_i)\}$ does not have zero mean, we can use the Gram matrix \tilde{K} to substitute the kernel matrix K. The Gram matrix is given by

$$\tilde{K} = K - \mathbf{1}_N K - K \mathbf{1}_N + \mathbf{1}_N K \mathbf{1}_N,$$

where $\mathbf{1}_N$ is the $N \times N$ matrix with all elements equal to $1/N$