



西安交通大学
XI'AN JIAOTONG UNIVERSITY

机器学习在信息安全中的应用

言湮

li.yan.88@xjtu.edu.cn

2025年11月



西安交通大学
XI'AN JIAOTONG UNIVERSITY

课程简介

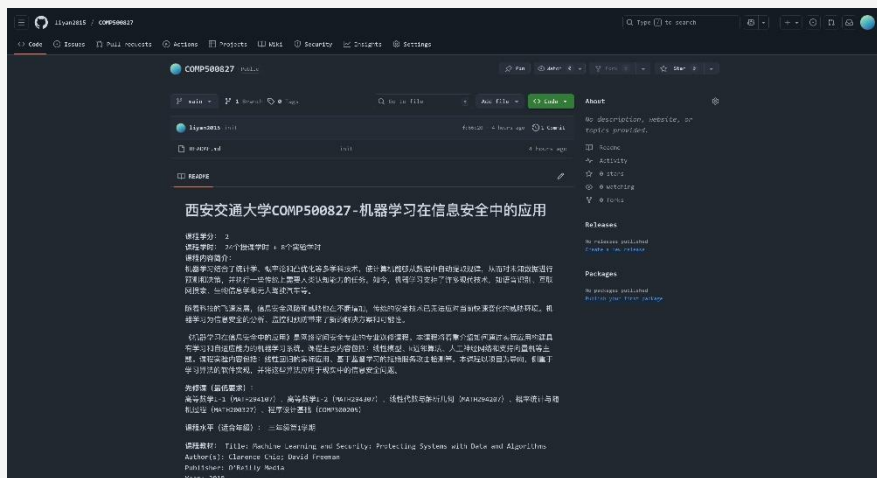
课程相关信息

课程资料地址:

<https://github.com/liyan2015/COMP500827/tree/main>

课程讨论区:

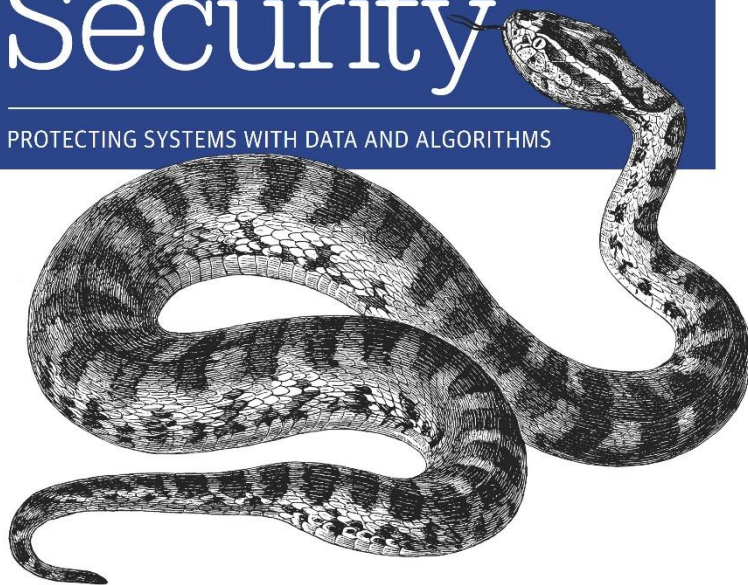
<https://class.xjtu.edu.cn/course/101578/forum#/>



O'REILLY®

Machine Learning & Security

PROTECTING SYSTEMS WITH DATA AND ALGORITHMS



Clarence Chio & David Freeman

电子书下载地址:

<http://libgen.li/edition.php?id=138042882>

书中代码示例:

<https://github.com/oreilly-mlsec/book-resources>

参考书



周志华 著. 机器学习
北京：清华大学出版社
2016年1月.

(ISBN 978-7-302-206853-6)
425页, 62.6万字

学习资料

<https://scikit-learn.org/stable/>



[Install](#) [User Guide](#) [API](#) [Examples](#) [Community](#) [More](#)



1.5.2 (stable)

scikit-learn

Machine Learning in Python

[Getting Started](#)

[Release Highlights for 1.5](#)

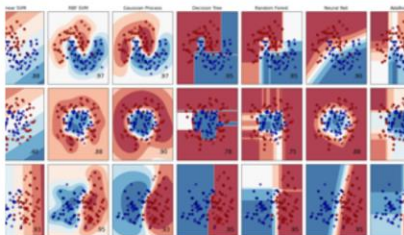
- Simple and efficient tools for predictive data analysis
- Accessible to everybody, and reusable in various contexts
- Built on NumPy, SciPy, and matplotlib
- Open source, commercially usable - BSD license

Classification

Identifying which category an object belongs to.

Applications: Spam detection, image recognition.

Algorithms: [Gradient boosting](#), [nearest neighbors](#), [random forest](#), [logistic regression](#), and [more...](#)



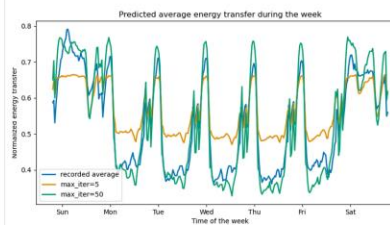
Examples

Regression

Predicting a continuous-valued attribute associated with an object.

Applications: Drug response, stock prices.

Algorithms: [Gradient boosting](#), [nearest neighbors](#), [random forest](#), [ridge](#), and [more...](#)



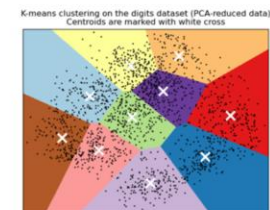
Examples

Clustering

Automatic grouping of similar objects into sets.

Applications: Customer segmentation, grouping experiment outcomes.

Algorithms: [k-Means](#), [HDBSCAN](#), [hierarchical clustering](#), and [more...](#)



Examples

Dimensionality reduction

Reducing the number of random variables to consider.

Applications: Visualization, increased efficiency.

Algorithms: [PCA](#), [feature selection](#), [non-negative matrix factorization](#), and [more...](#)

Model selection

Comparing, validating and choosing parameters and models.

Applications: Improved accuracy via parameter tuning.

Algorithms: [Grid search](#), [cross validation](#), [metrics](#), and [more...](#)

Preprocessing

Feature extraction and normalization.

Applications: Transforming input data such as text for use with machine learning algorithms.

Algorithms: [Preprocessing](#), [feature extraction](#), and [more...](#)

学习资料



coursera



考核方式

考核方式：考查

成绩比例：平时成绩20%（作业）

实验成绩占40%

期末项目考查成绩占40%

课程简介

1. 什么是学习?

“学习是系统通过经验提升性能的过程。”

--- Herbert Simon

卡内基·梅隆大学

图灵奖(1975)

人工智能, 人类认知心理学

诺贝尔经济学奖(1978)

经济组织内的决策过程

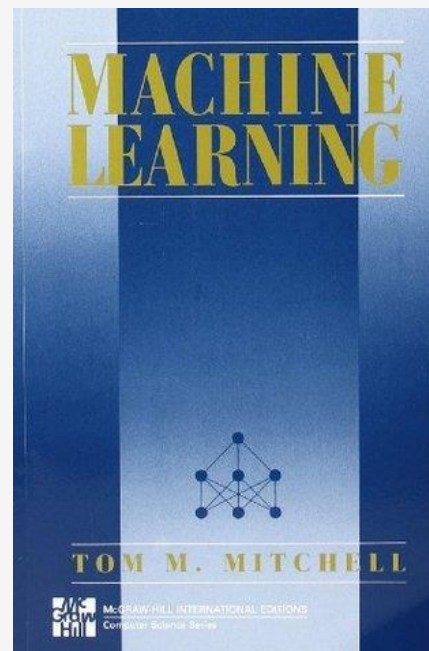


课程简介

由Tom Mitchell 给出的更加数学化的定义

- Ability for machines to learn without being *explicitly* programmed

"A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P if its performance at tasks in T , as measured by P , improves with experience E ." --- Mitchell, T. (1997). Machine Learning. McGraw Hill. p. 2.



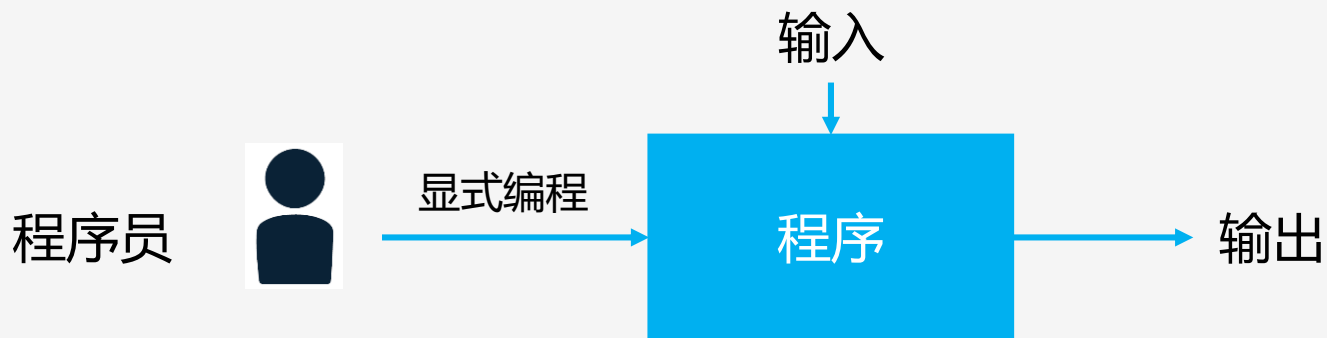
课程简介

由Tom Mitchell 给出的更加数学化的定义

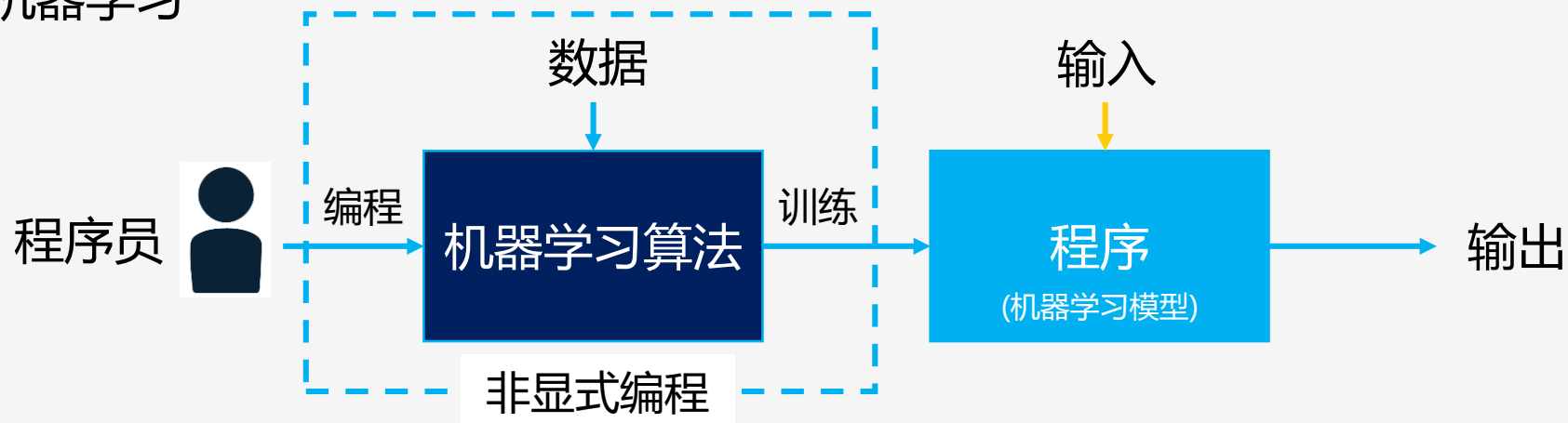
- 机器学习是一门研究学习算法的学科，这些算法能够：
 - 在某些任务 T 上
 - 通过经验 E
 - 提升性能 P
 - 非显式编程
- 一个学习任务可以由三元组 $\langle T, P, E \rangle$ 明确定义
 - 为何不能用知识、经验或者专业技能来训练机器呢？
 - 人类总能够解释他们的专业技能吗？

课程简介

▣ 传统编程



▣ 机器学习



课程简介

机器学习在什么情况下具有优势？

应用情形：

- 模型基于大量数据
 - 例子：Google 网络搜索，垃圾邮件识别
- 输出必须是个性化的
 - 例子：新闻/物品/广告推荐
- 人类不能解释专业知识
 - 例子：语音/人脸识别，异常行为检测
- 人类的专业知识不存在
 - 例子：在火星上导航

课程简介

两种机器学习类型

□ 预测

- 根据数据预测所需的输出（监督学习）
- 生成数据实例（无监督学习）

□ 决策

- 在动态环境中采取行动（强化学习）
 - 转变到新的状态
 - 获得即时奖励
 - 随着时间的推移最大化累积奖励

课程简介

2. 机器学习的历史

□ 1950年代:

- Samuel 的跳棋程序
- 创建机器学习术语

□ 1960年代:

- 神经网络, 感知机
- 模式识别
- Minsky和Papert证明感知机的局限性



Arthur Samuel在 1959年创造了"机器学习"这个词

□ 1970年代:

- 符号概念归纳
- Winston的结构学习系统
- 专家系统和知识获取瓶颈
- Quinlan的ID3算法
- 使用 Automatic Mathematician 的数学发现

□ 1980年代:

- 高级决策树和规则学习
- 基于解释的学习 (EBL)
- 学习、规划、解决问题
- 三间小屋问题
- 类比
- 认知架构
- 神经网络复苏 (反向传播)
- Valiant的PAC学习理论
- 注重实验方法

课程简介

2. 机器学习的历史

□ 1990年代

- 数据挖掘
- 自适应软件代理和网络应用程序
- 文本学习
- 强化学习 (RL)
- 归纳逻辑编程 (ILP)
- 组合: 装袋、提升和堆叠
- 贝叶斯网络学习
- 支持向量机
- 核方法

□ 2000年代

- 图模型
- 变分推理
- 统计关系学习
- 迁移学习
- 序列标记
- 集体分类和结构化输出
- 计算机系统应用
- 电子邮件管理
- 学习的个性化助手
- 机器人和视觉的学习

课程简介

2. 机器学习的历史

□ 2010年代

- 深度学习
- 从大数据中学习
- 结合知识图谱的机器学习
- 使用 GPU 或 HPC 学习
- 多任务 + 终身学习
- 深度强化学习, AlphaGo
- 视觉、语音、文本、网络、行为等的庞大应用

□ 2020年代

- 结合逻辑推理的深度学习
- 超大规模预训练模型 (GPT-3等)
- 无人驾驶
- AI生物制药技术
- 生成式模型的爆发 (ChatGPT、DALL-E、Diffusion)

课程简介

3. 机器学习基本思想

机器学习类型:

□ 监督学习

- 给定数据和标签，预测所需的输出

□ 无监督学习

- 分析和利用隐式数据模式/结构

□ 强化学习

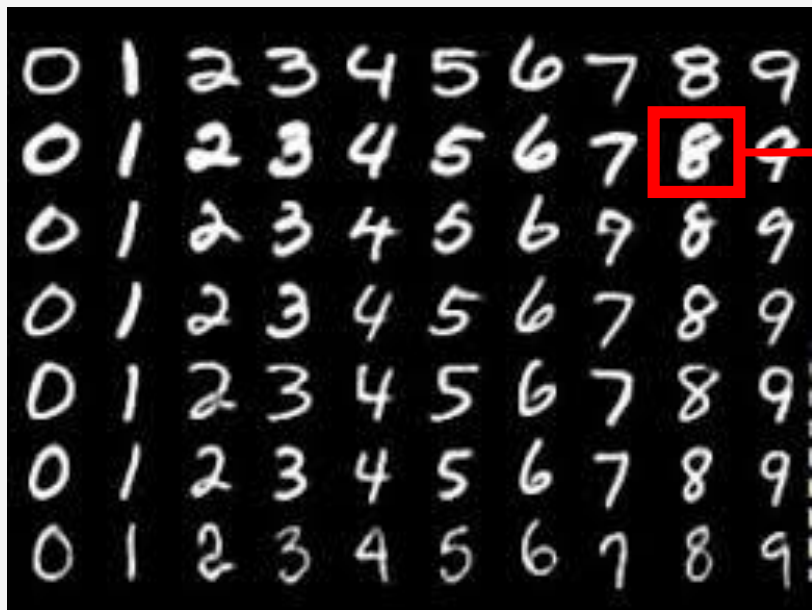
- 学习在动态环境中动作执行的决策，并获得尽可能多的奖励值

课程简介

任务 (T) :

- 给定一个手写数字图片集合 $x \in [0,255]^{28 \times 28}$ 和其对应的标签 $y \in [0,9]$ 找到一个映射函数

$$f: x \rightarrow y$$



经验 (E) :

- 一个用于训练的标注好0~9标签的图片集

性能 (P) :

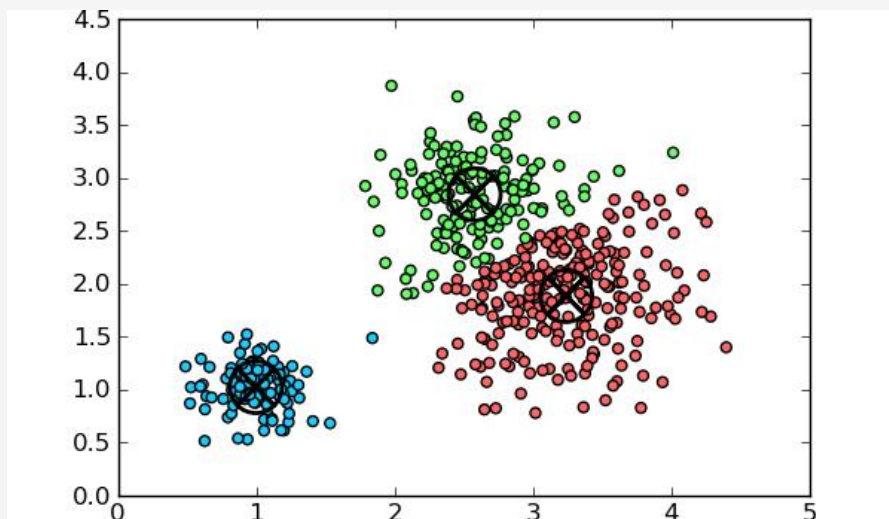
- 这个映射函数在未标注标签的测试集上的识别准确率

“**监督学习 (Supervised Learning)**”

课程简介

任务 (T) :

- 如何将一组文档 “聚类” (Cluster) 成 k 个组, 使得属性 “相似” 的文档出现在同一组中?



经验 (E) :

- 一个用于训练的无标签的文档集合

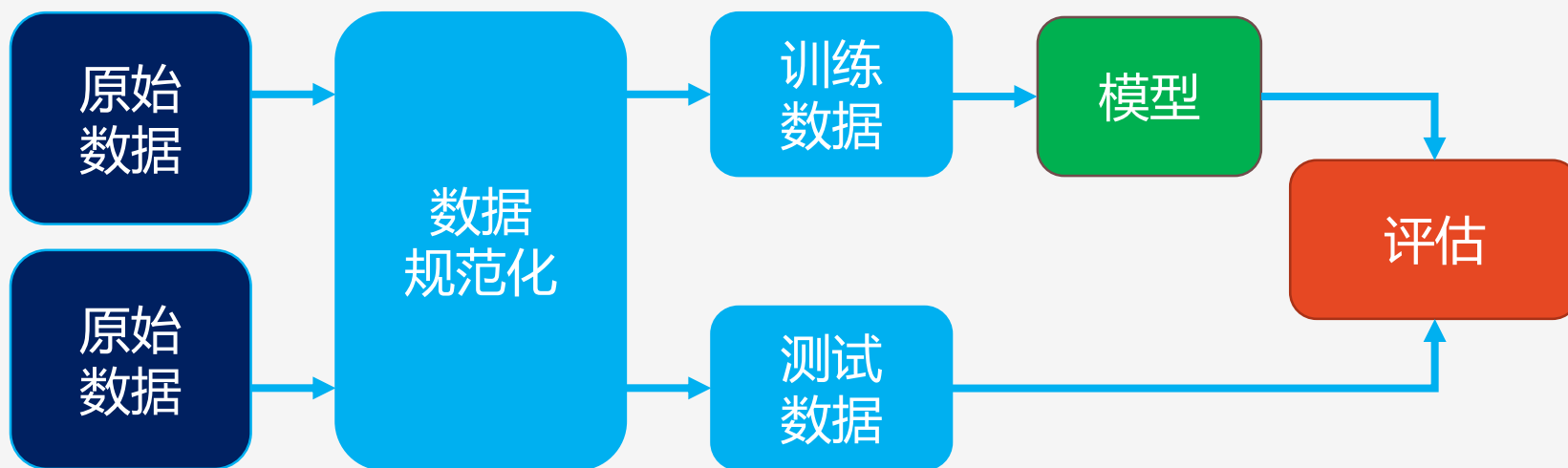
性能 (P) :

- 所有文档的属性坐标到聚类中心的平均距离

“无监督学习 (Unsupervised Learning) ”

课程简介

机器学习过程:



- 基本假设：在训练和测试数据中存在相同的模式（**pattern**）

课程简介

监督学习:

定义

- 给定带标签的训练数据集: $D = \{(x_i, y_i)\}_{i=1,2,\dots,N}$, 其中 x_i 为特征数据, y_i 为其对应的标签, 让机器学习一个从特征数据映射到标签的函数映射

$$y_i \simeq f_{\theta}(x_i)$$

- 函数集 $\{f_{\theta}(\cdot)\}$ 被称为假设空间
- 学习的过程即为参数 θ 的更新

如何学习?

- 更新参数以使预测结果接近真实的标签
 - 学习目标是什么?
 - 如何更新参数?

课程简介

监督学习:

学习目标

- ▣ 使预测结果接近真实的标签

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N \mathcal{L}(y_i, f_{\theta}(x_i))$$

- ▣ 损失函数 $\mathcal{L}(y_i, f_{\theta}(x_i))$ 用来衡量标签和预测结果之间的误差
- ▣ 损失函数的定义取决于数据和任务
- ▣ 最常见的损失函数：平方误差 (squared loss)

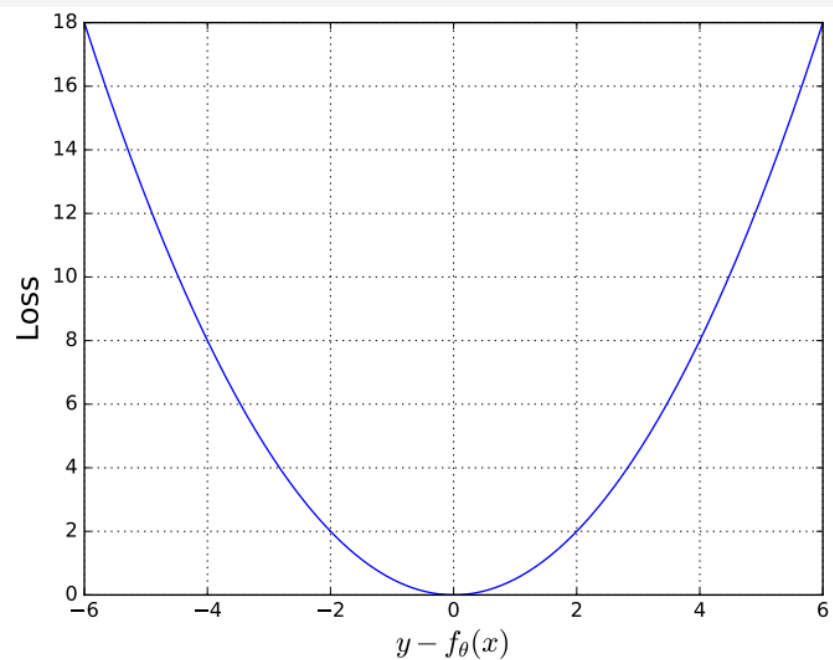
$$\mathcal{L}(y_i, f_{\theta}(x_i)) = \frac{1}{2} (y_i - f_{\theta}(x_i))^2$$

课程简介

监督学习:

平方误差

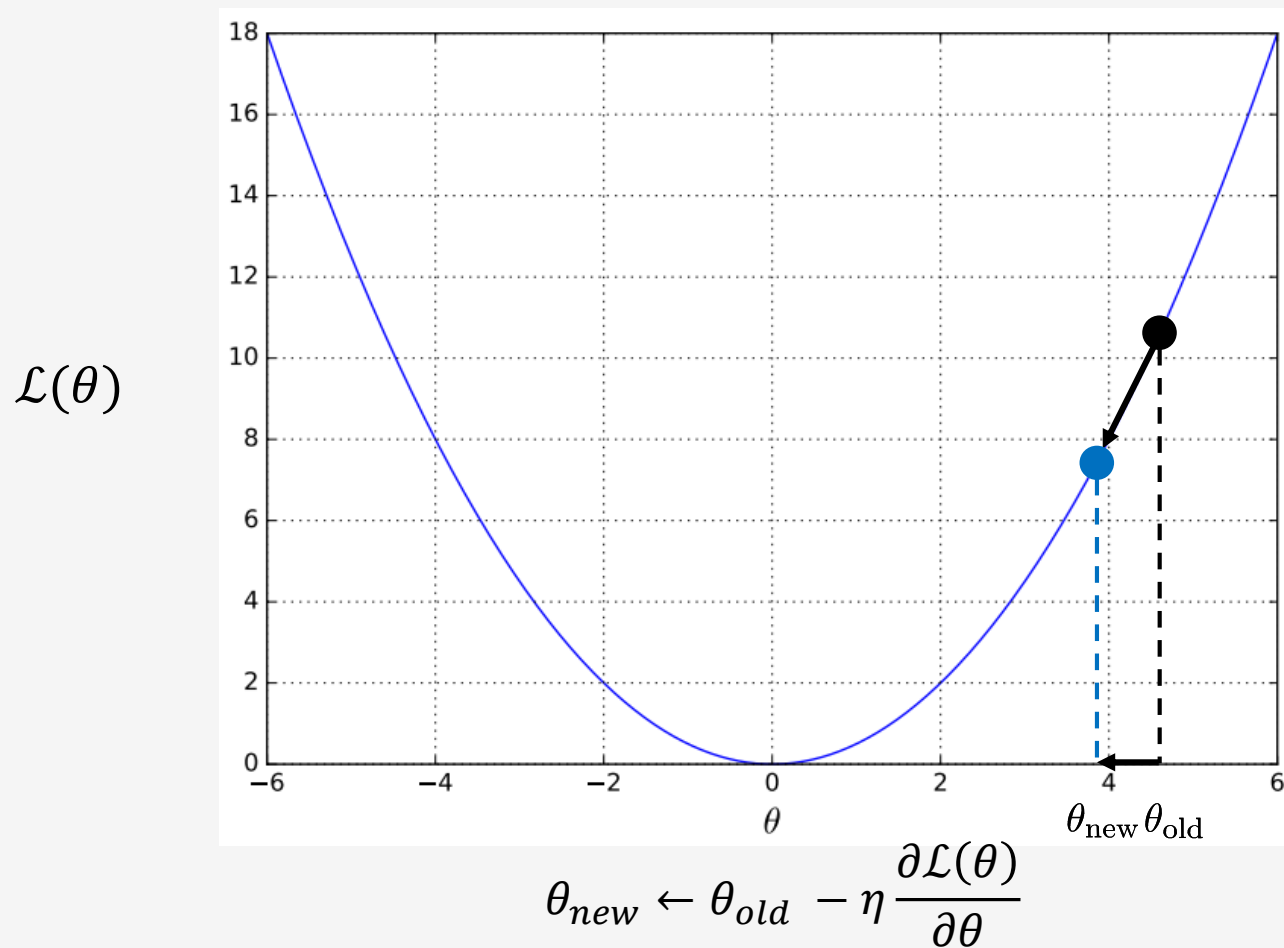
$$\mathcal{L}(y_i, f_{\theta}(x_i)) = \frac{1}{2} (y_i - f_{\theta}(x_i))^2$$



- 距离越远，得到的惩罚更多
- 容忍小距离（误差）
 - 观察噪声等
 - 泛化性

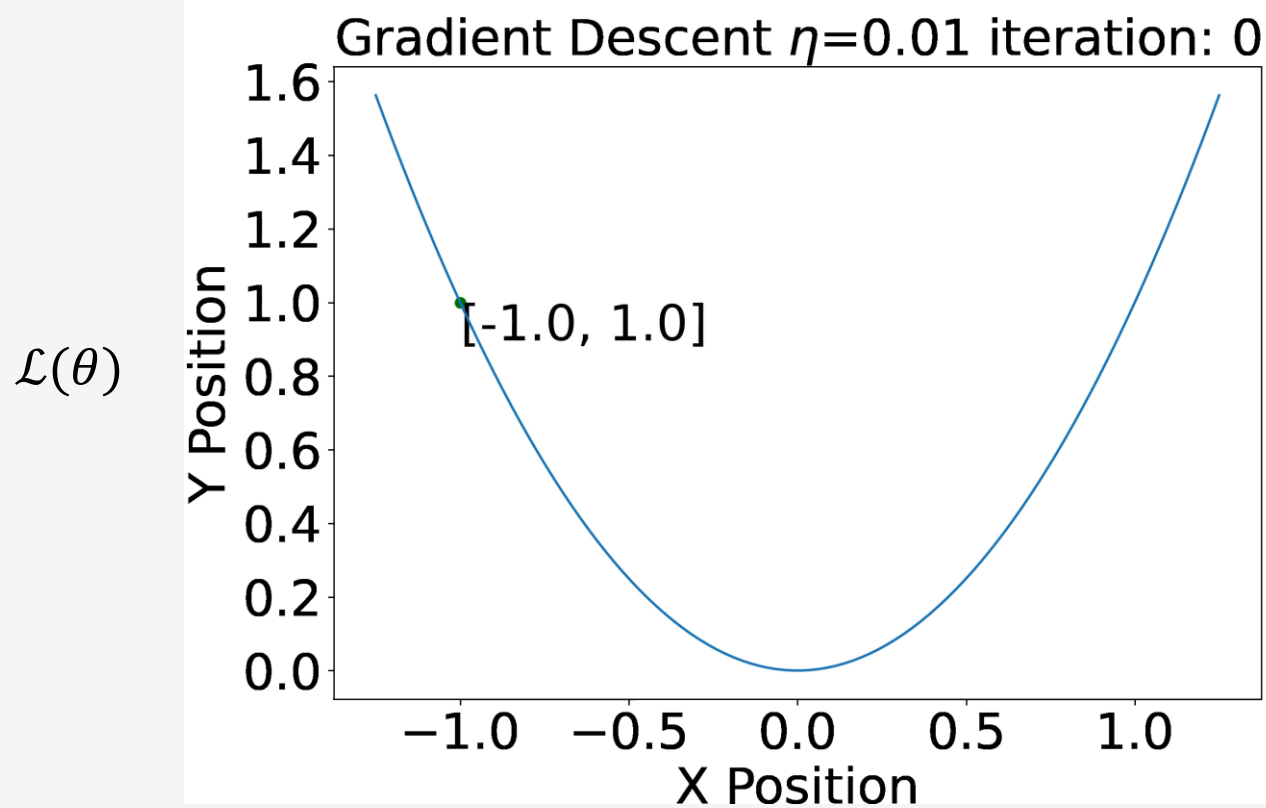
课程简介

梯度学习方法:



课程简介

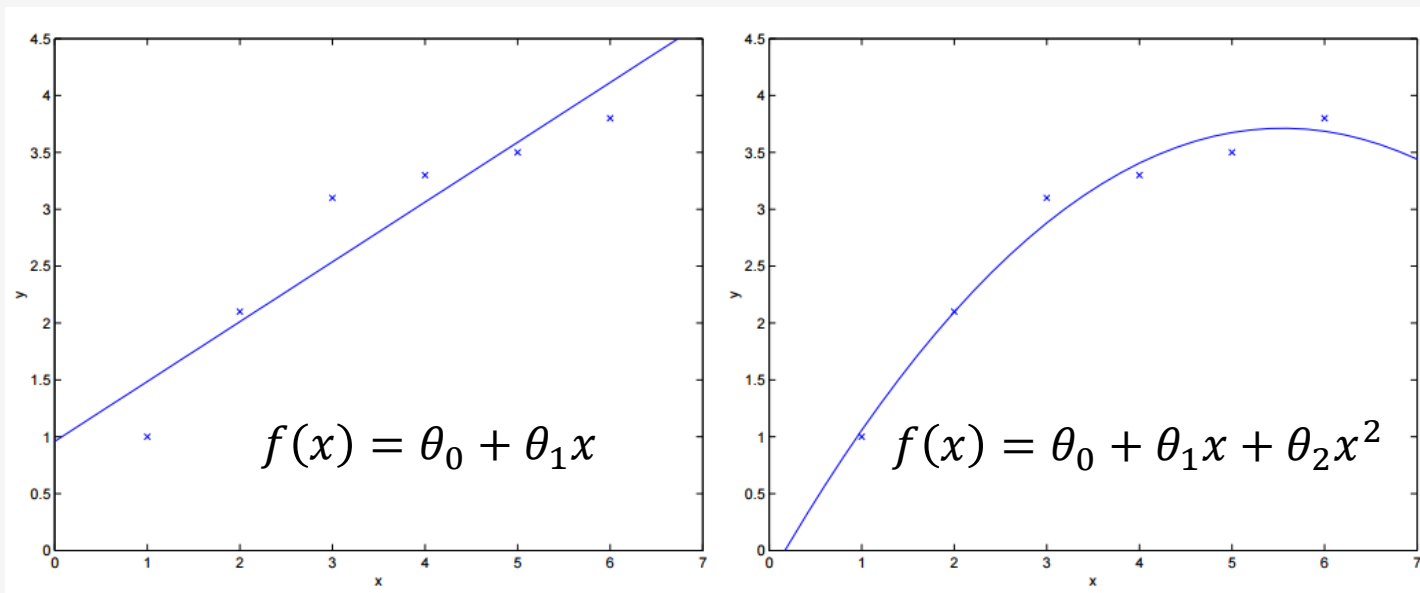
梯度学习方法:



$$\theta_{new} \leftarrow \theta_{old} - \eta \frac{\partial \mathcal{L}(\theta)}{\partial \theta}$$

课程简介

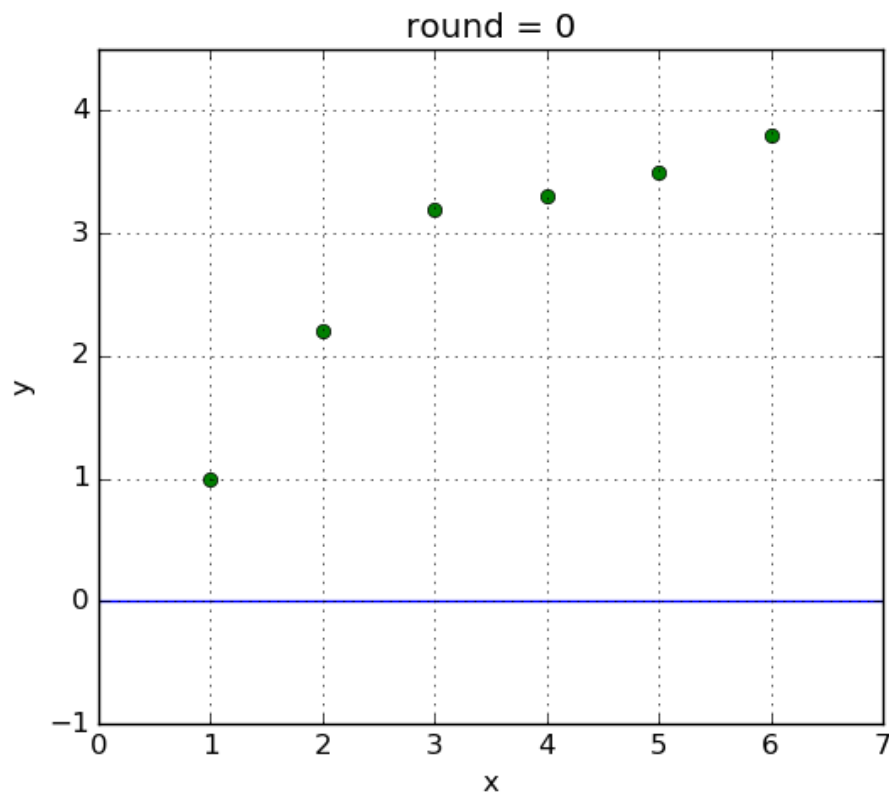
一个简单的例子:



- 观察数据 $\{(x_i, y_i)\}_{i=1,2,\dots,N}$, 我们可以使用不同的模型 (假设空间) 来学习
 - 模型选择 (线性或二次)
 - 参数学习

课程简介

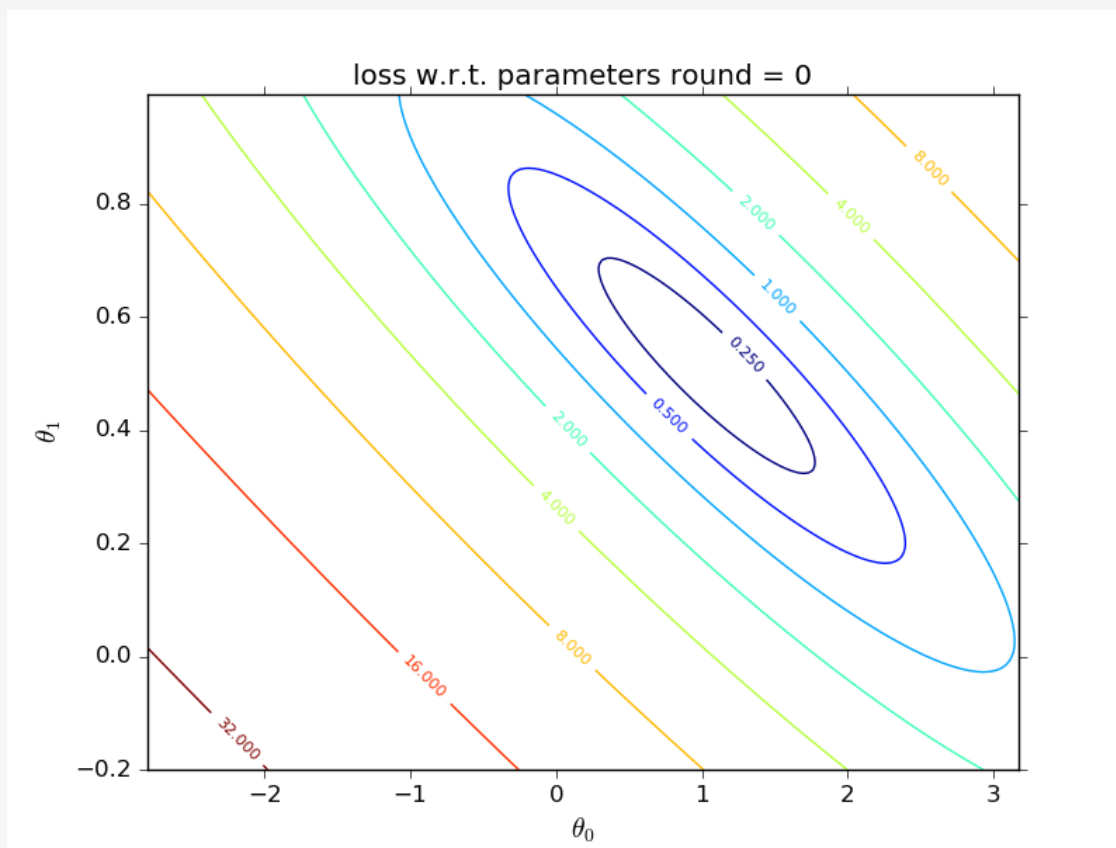
利用线性模型学习-曲线拟合



$$f(x) = \theta_0 + \theta_1 x$$

课程简介

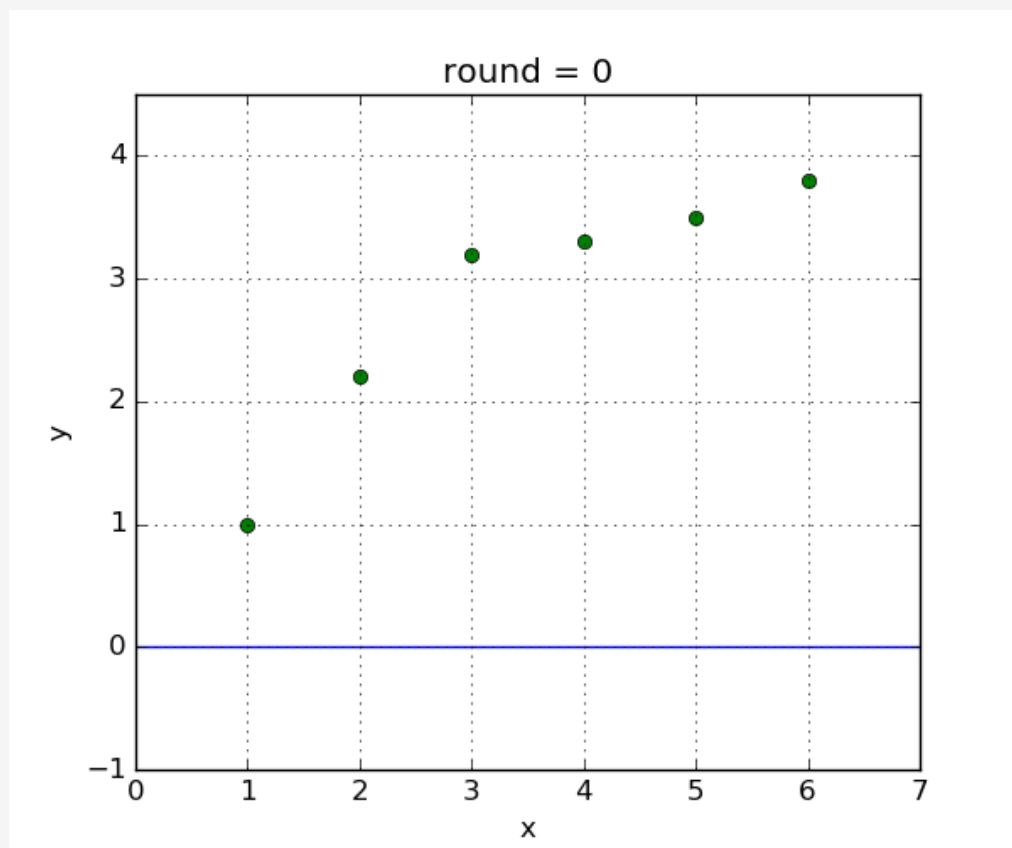
利用线性模型学习-权重分布



$$f(x) = \theta_0 + \theta_1 x$$

课程简介

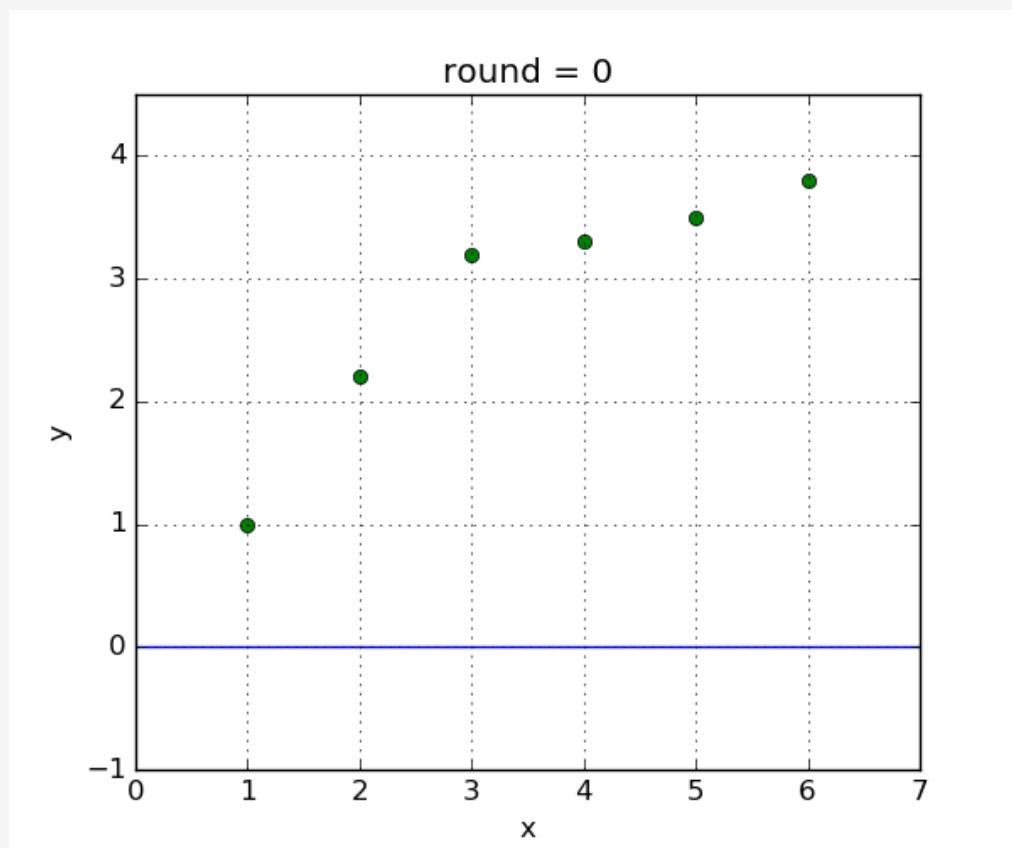
利用二次模型学习



$$f(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$

课程简介

利用三次模型学习



$$f(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$

课程简介

机器学习应用举例：垃圾邮件分类

任务：

- 邮件 $x \in$ 所有邮件, $y \in \{\text{垃圾邮件, 非垃圾邮件}\}$



垃圾邮件?

$$f: x \rightarrow y$$

经验：

- 所有已被标识的“垃圾邮件”和“非垃圾邮件”，即“有标签的训练数据集”

性能：

- “垃圾邮件”的识别准确率

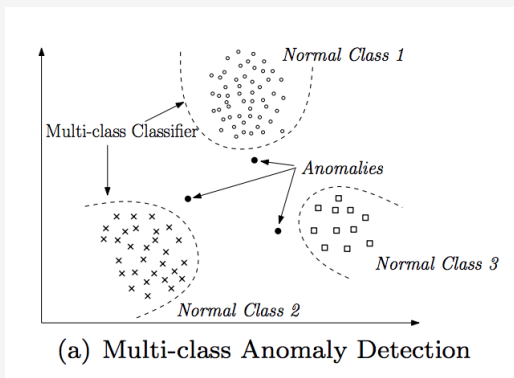
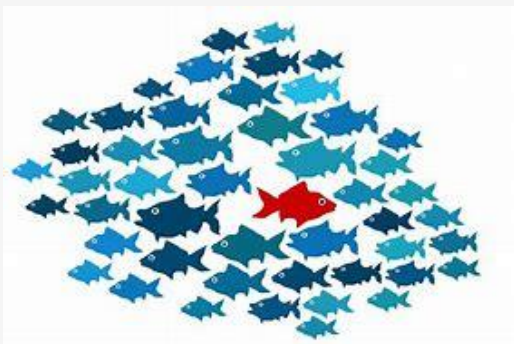
“监督学习 (分类问题)”

课程简介

机器学习应用举例：异常事件检测

任务：

- 筛选出和其他事件最不相似的实例



经验：

- 所有没有明确标识的事件样本，即“无标签的训练数据集”

性能：

- 异常事件的识别准确率

“无监督学习（聚类问题）”

第一章：模型选择



1.1 欠拟合与过拟合

1.2 正则化

1.3 奥卡姆剃刀原则

1.4 交叉验证

1.5 模型泛化性

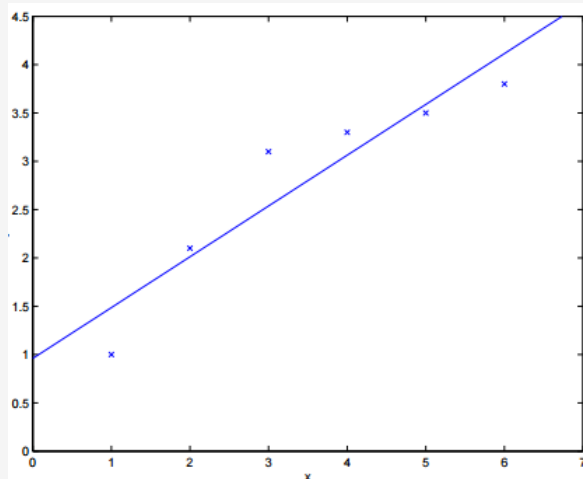


西安交通大学
XI'AN JIAOTONG UNIVERSITY

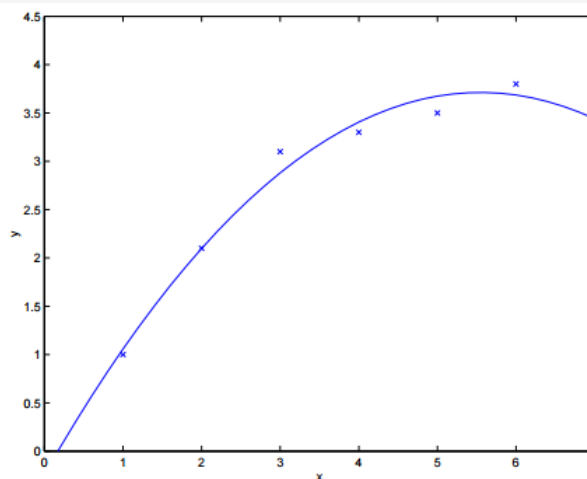
1.1 欠拟合与过拟合

1.1 欠拟合与过拟合

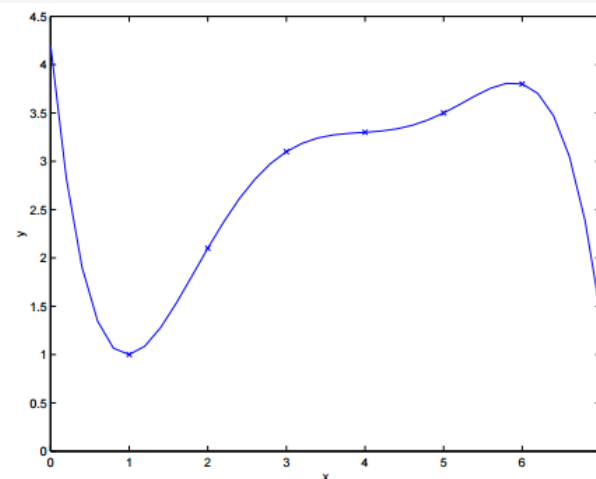
□ 下面哪个模型是最好的？



线性模型：欠拟合



二次模型：合适

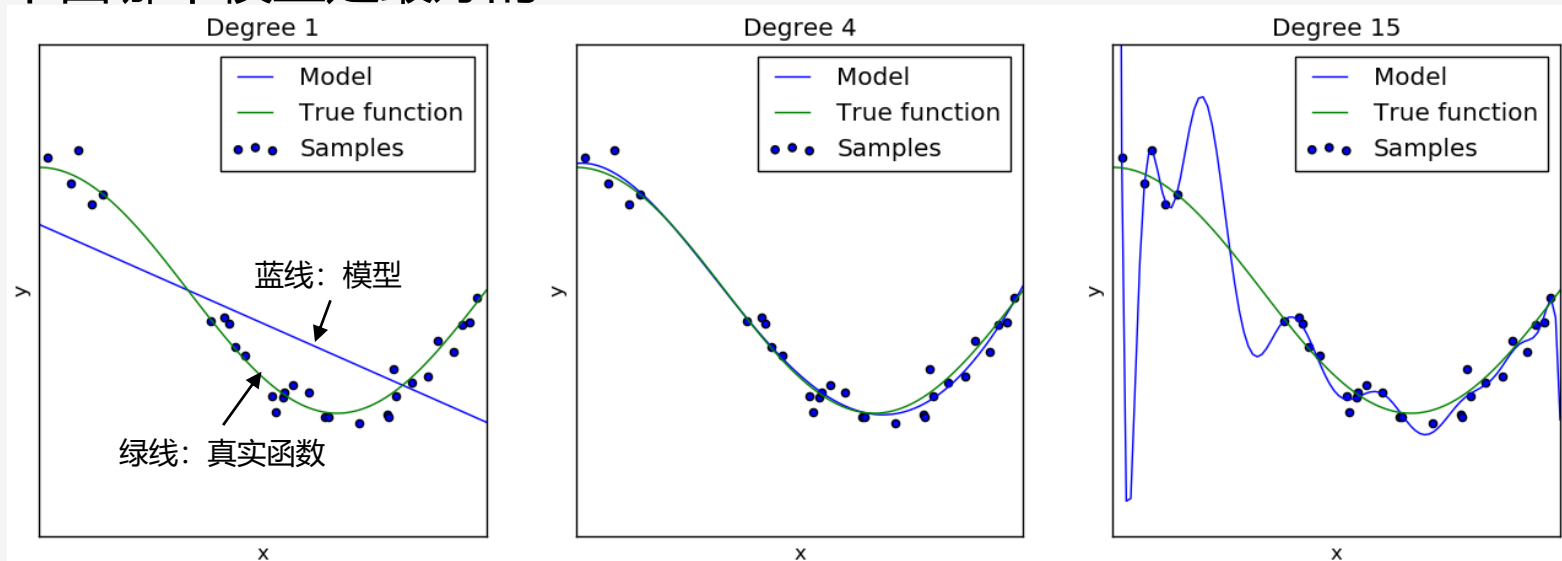


五阶模型：过拟合

- 当统计模型或机器学习算法无法捕捉数据的基本变化趋势时,就会出现**欠拟合**。
- 当统计模型把随机误差和噪声也考虑进去而不仅仅是考虑数据的基础关联时,就会出现**过拟合**。

1.1 欠拟合与过拟合

□ 下面哪个模型是最好的？



线性模型：欠拟合

四阶模型：合适

十五阶模型：过拟合

- 当统计模型或机器学习算法无法捕捉数据的基础变化趋势时,就会出现**欠拟合**。
- 当统计模型把随机误差和噪声也考虑进去而不仅仅是考虑数据的基础关联时, 就会出现**过拟合**。



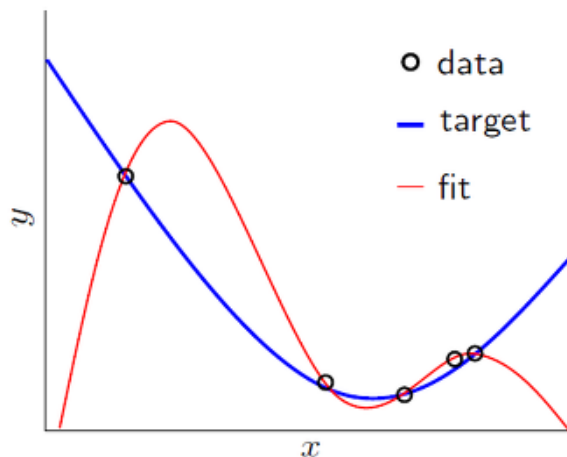
西安交通大学
XI'AN JIAOTONG UNIVERSITY

1.2 正则化

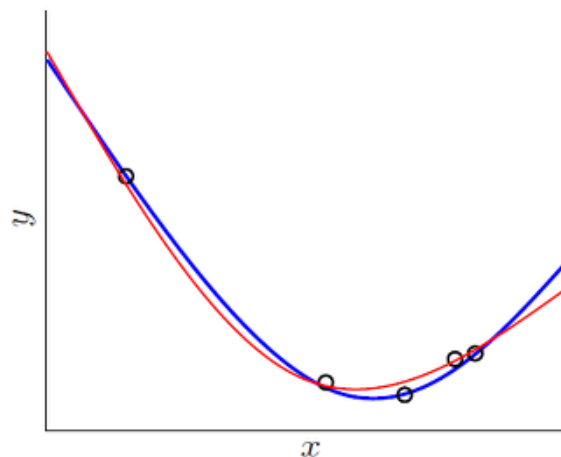
1.2 正则化

- 添加参数的惩罚项，防止模型对数据过拟合

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N \mathcal{L}(y_i, f_{\theta}(x_i)) + \lambda \Omega(\theta)$$



(a) without regularization



(b) with regularization

1.2 正则化

经典正则化方法

□ L2正则化 (岭回归Ridge)

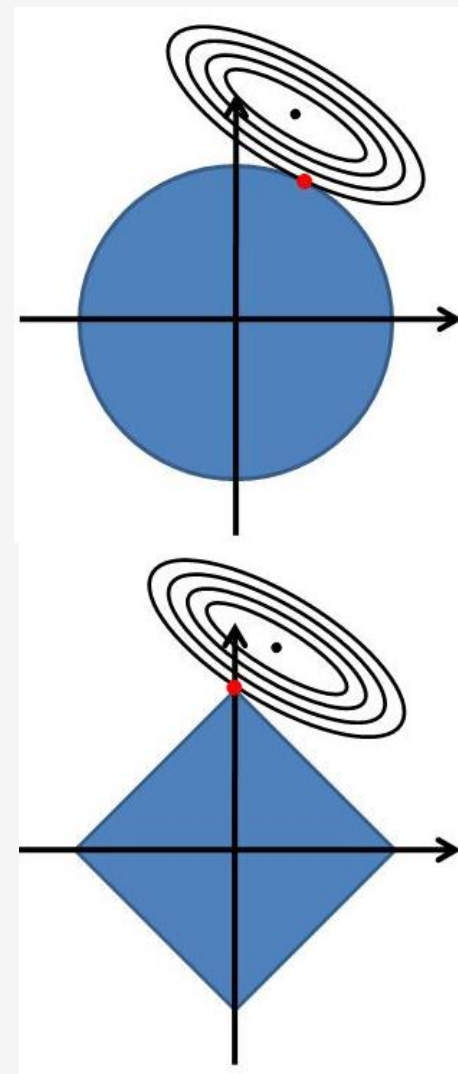
$$\Omega(\theta) = \|\theta\|_2^2 = \sum_{m=1}^M \theta_m^2$$

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N \mathcal{L}(y_i, f_{\theta}(x_i)) + \lambda \|\theta\|_2^2$$

□ L1正则化 (拉索回归LASSO)

$$\Omega(\theta) = \|\theta\|_1 = \sum_{m=1}^M |\theta_m|$$

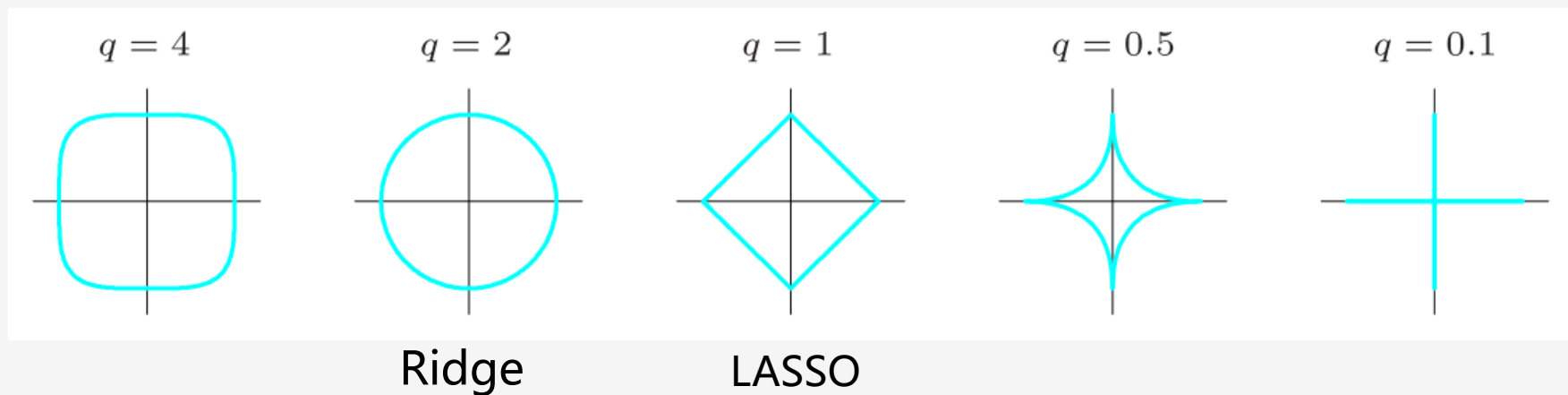
$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N \mathcal{L}(y_i, f_{\theta}(x_i)) + \lambda \|\theta\|_1$$



1.2 正则化

经典正则化方法

□ 常值 $\sum_j |\theta_j|^q$ 的数值分布图



- 当 $q \leq 1$ 的时候，模型进行稀疏性学习
- 很少会用 $q > 2$ 来进行正则化
- 99%的情形下都取 $q = 1$ 或 2



西安交通大学
XI'AN JIAOTONG UNIVERSITY

1.3 奥卡姆剃刀原则

1.3 奥卡姆剃刀原则

□ 有多个假设模型时，我们应该选择假设条件最少的建模方法。

□ 函数集 $\{f_{\theta}(\cdot)\}$ 被称作假设空间

$$\min_{\theta} \left[\frac{1}{N} \sum_{i=1}^N \mathcal{L}(y_i, f_{\theta}(x_i)) \right] + \lambda \Omega(\theta)$$

原始损失

基于假设的罚值

1.3 奥卡姆剃刀原则

模型选择

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N \mathcal{L}(y_i, f_{\theta}(x_i)) + \lambda \|\theta\|_2^2$$

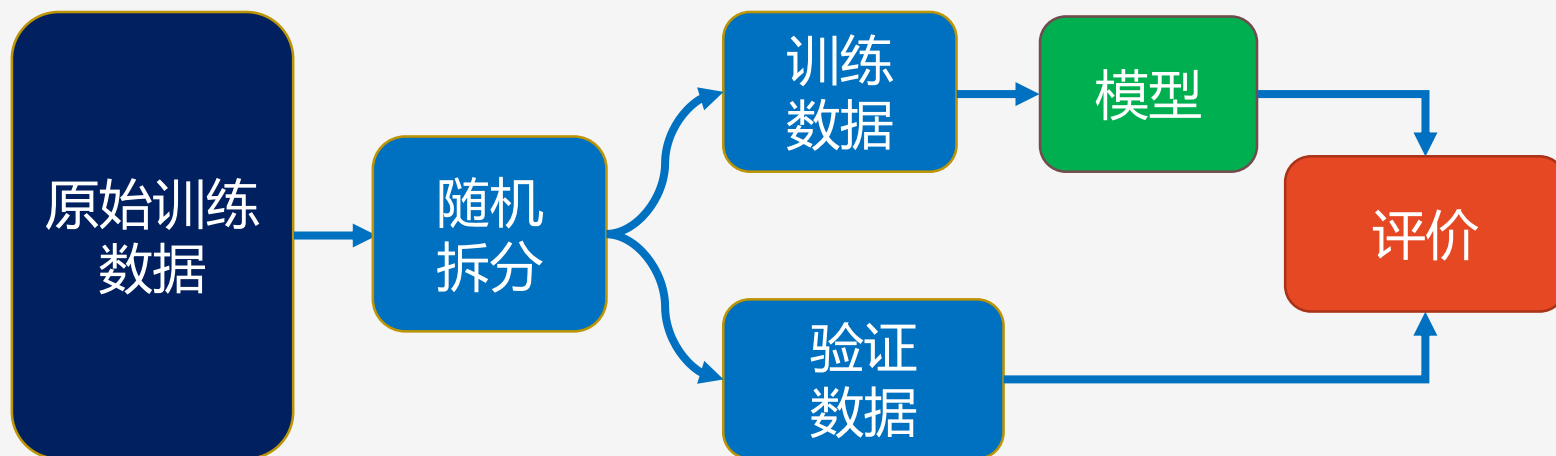
- 一个机器学习的解决方案的模型包含参数 θ 和超参数 λ
- 超参数
 - 定义模型的更高层次的概念，如复杂性或学习能力。
 - 在标准模型训练过程中**无法直接从数据中学习**，需要预先定义。
 - 可以通过不同的参数设置、训练不同的模型，以及选择最好的测试结果来进行超参数选择
- 模型选择（或超参数优化）关注如何选择最佳超参数



西安交通大学
XI'AN JIAOTONG UNIVERSITY

1.4 交叉验证

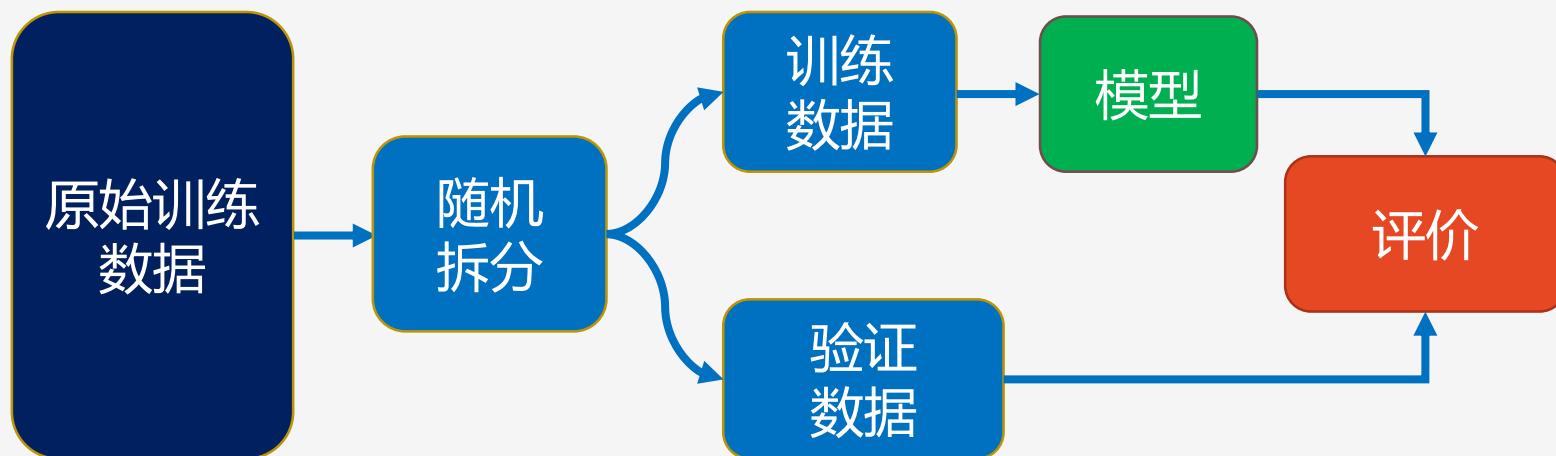
1.4 交叉验证



K-折交叉验证

1. 设置超参数
2. 将原始训练数据随机拆分为K份
3. 重复K次:
 - 若当前为第 i 次重复 ($i=1,\dots,K$) , 选择第 i 份数据作为验证数据集, 其余 $K-1$ 份作为训练数据集
 - 对训练数据进行建模,并在验证数据上对其进行评估,从而获得评估分数
4. 对K个评估分数取平均作为模型性能

1.4 交叉验证



- 选择了“好的”超参数后，对整个训练数据进行模型训练，然后用测试数据对模型进行测试。



西安交通大学
XI'AN JIAOTONG UNIVERSITY

1.5 模型泛化性

1.5 模型泛化性

泛化能力

□ 泛化能力指的是模型对未观测数据的预测能力

- 可以通过泛化误差来评估，定义如下：

$$R(f) = \mathbb{E}[\mathcal{L}(Y, f(X))] = \int_{X \times Y} \mathcal{L}(y, f(x)) p(x, y) dx dy$$

- $p(x, y)$ 是潜在的（可能是未知的）联合数据分布

□ 在训练数据集上对泛化能力的经验估计为：

$$\hat{R}(f) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(y_i, f(x_i))$$

1.5 模型泛化性

泛化误差

□ 有限假设集 $\mathcal{F} = \{f_1, f_2, \dots, f_d\}$

□ 泛化误差约束定理:

对任意函数 $f \in \mathcal{F}$, 以不小于 $1 - \delta$ 的概率满足下式:

$$R(f) \leq \hat{R}(f) + \epsilon(d, N, \delta)$$

其中,

$$\epsilon(d, N, \delta) = \sqrt{\frac{1}{2N} \left(\log d + \log \frac{1}{\delta} \right)}$$

- N : 训练实例个数
- d : 假设集的函数个数

小测验

- 下列有关机器学习基本术语说法错误的是（ D ）
- A. 从训练数据中学得模型的过程称为“学习”或“训练”
 - B. 训练过程中使用的数据称为“训练数据”，每一个样本称为一个“训练样本”，训练样本组成的集合称为“训练集”
 - C. 学习得到的模型对应了关于数据的某种潜在规律，称为“假设”
 - D. 学习过程就是为了找出数据的某种潜在规律，这个规律自身，一般称为“数据特征”

小测验

□ 下列说法错误的是（ D ）

- A. 模型是通过学习算法得到的
- B. 机器学习通常解决高度不确定性和复杂性的问题
- C. 分类和回归是监督学习的代表
- D. 机器学习一定需要类别标记

小测验

□ 下列说法错误的是 (C)

- A. 学习获得的模型适用于新样本的能力, 称为“泛化”能力
- B. 机器学习问题一般有“独立同分布”假设
- C. 机器学习在只要见过的数据上做好了就行, 未见过样本的性能不重要
- D. 机器学习问题通常假设拿到的所有数据都来自一个潜在的分布

小测验

□ 下列说法错误的是（ C ）

- A. “色泽” 取值为 “青绿” ， 这里的 “青绿” 是属性值
- B. 输出是离散值的学习任务为分类任务
- C. 模型找出的规律一定是正确的
- D. 一般假设正类和反类是可交换的

小测验

- 以下哪个选项不是指“奥卡姆剃刀”原则？（A, B, D）
- A. 若有多个假设与观察一致,则随机选一个
 - B. 若有多个假设与观察一致,则选即不简单又不复杂的那个
 - C. 若有多个假设与观察一致,则选最简单的那个
 - D. 若有多个假设与观察一致,则选最复杂的那个

小测验

□ 以下关于机器学习预测任务的说法正确的是 (A, B)

- A. 一般地, 预测任务是希望对训练集进行学习, 建立一个从输入空间 x 到输出空间 y 的映射 $f: x \rightarrow y$
- B. 对于二分类任务, 一般令 $y = \{-1, +1\}$ 或 $\{0, 1\}$
- C. 在任何情况下, 总有一个最优的学习算法
- D. 预测任务不需要训练样本的标记信息

小测验

□ 把未见过的汽车分为若干组,这是一个__C__任务

A. 分类

B. 回归

C. 聚类

小测验

- 两种算法在某种情况下取得评估结果后不能直接比较以评判优劣的原因中，正确的是（ D ）
- A. 测试性能不等于泛化性能
 - B. 测试性能随着测试集的变化而变化
 - C. 很多机器学习算法本身有一定随机性
 - D. 以上均正确



西安交通大学
XI'AN JIAOTONG UNIVERSITY

谢谢大家！

