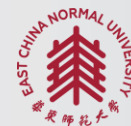


深度学习

明静思

华东师范大学统计学院



ECNU

课程内容

1. 机器学习基础
2. 全连接神经网络
3. 卷积神经网络
4. 网络优化与正则化
5. 深度生成模型
6. 循环神经网络

考核方式:

课堂表现 (10%)

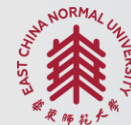
作业 (40%)

期末大作业 (50%)

学习资料

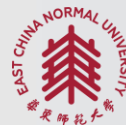
- 邱锡鹏. 神经网络与深度学习. 机械工业出版社, 2020. <https://nndl.github.io/>
- Zhang, Aston, Zachary C. Lipton, Mu Li, and Alexander J. Smola. "*Dive into deep learning*." arXiv preprint arXiv:2106.11342 (2021). <https://zh.d2l.ai> (中英文版)
- Goodfellow, Ian, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.

第一章 机器学习基础



ECNU

人工智能、机器学习与深度学习



ECNU

人工智能

- **人工智能 (Artificial Intelligence, AI)** 就是要让机器的行为看起来就像是人所表现出的智能行为一样。

“计算机控制” + “智能行为”

感知

- 计算机视觉、语音信息处理

学习

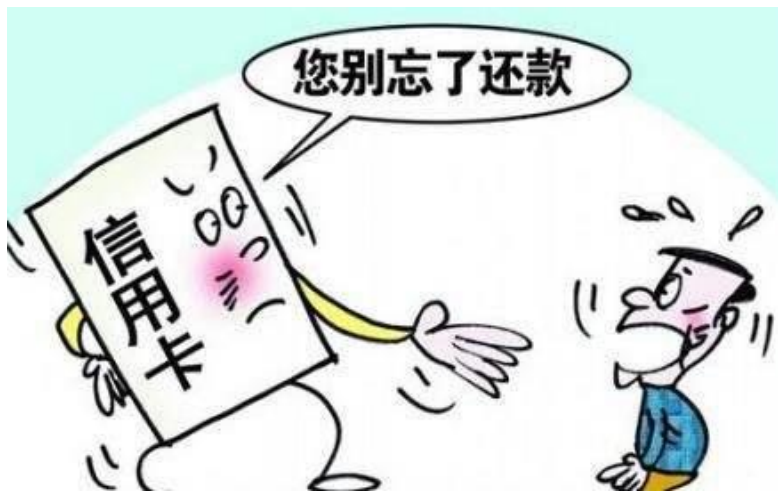
- 监督学习、无监督学习、强化学习

认知

- 知识表示、自然语言理解、推理、决策

机器学习

- 机器学习是人工智能的一个重要分支，并逐渐成为推动人工智能发展的关键因素。
- **机器学习 (Machine Learning, ML)** 是指从有限的观测数据中学习出具有一般性的规律，并利用这些规律对未知数据进行预测的方法。



任务：判断信用卡用户是否会发生逾期

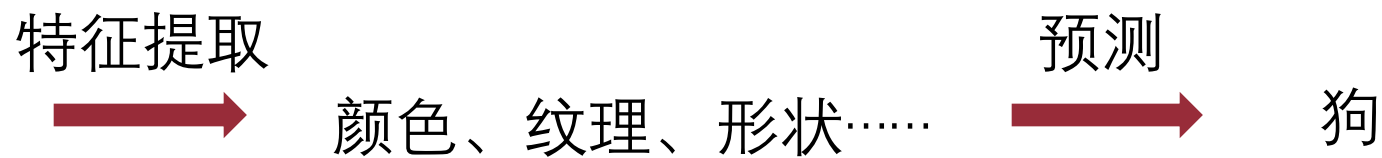
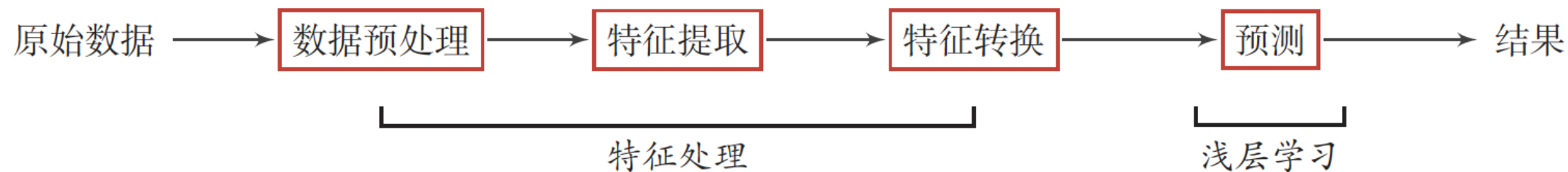
训练数据：信用卡用户的特征（年龄、性别、工作、信用卡额度、存款、是否有住房贷款等）和逾期状态

机器学习算法：学习用户特征与逾期状态之间的模型

预测：使用模型预测新用户未来的逾期情况

传统的机器学习方法

- 在实际任务中使用机器学习模型一般会包含以下几个步骤：



语义鸿沟

- 基于内容的图像检索系统中广泛存在“语义鸿沟”问题。
- 语义鸿沟问题是指输入数据的底层视觉特征和高层语义信息之间的不一致性和差异性。

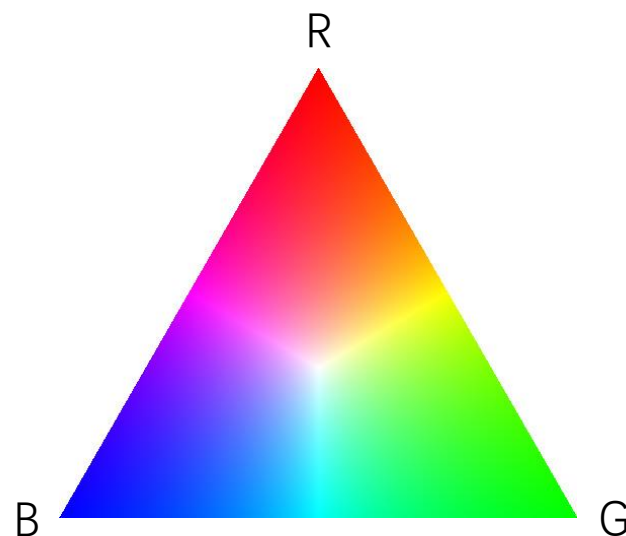


表示学习

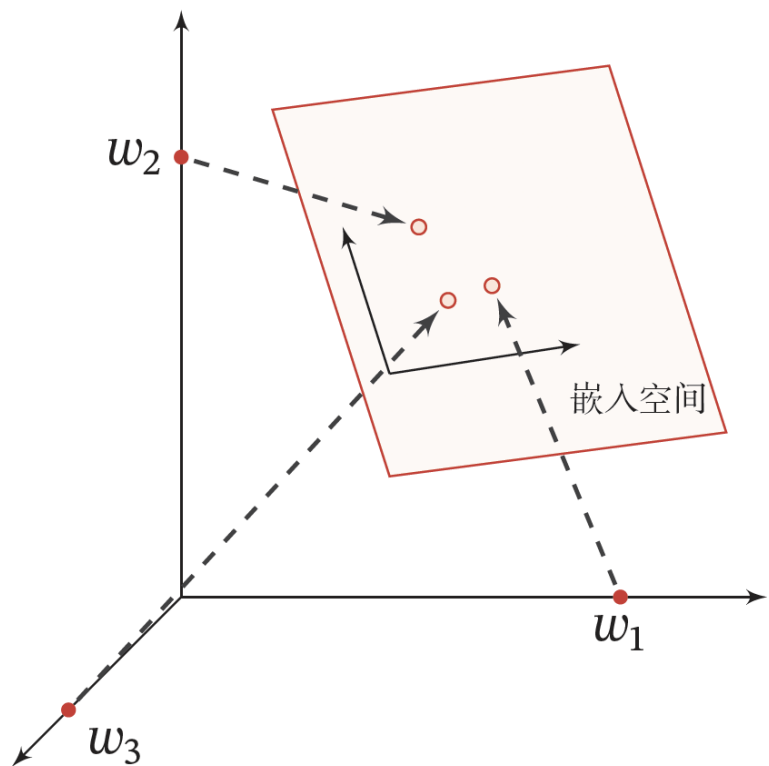
- 为了提高机器学习系统的准确率，需要将输入信息转换为有效的**特征**，或者更一般性地称为**表示**。
- **表示学习**：自动学习出有效的特征，并提高最终机器学习模型性能算法。
- 一般而言，一个好的表示具有以下几个特点：
 - 应该具有很强的表示能力，即同样大小的向量可以表示更多信息。
 - 应该使后续的学习任务变得简单，即需要包含更高层的语义信息。
 - 应该具有一般性，是任务或领域独立的。

局部表示和分布式表示

颜色	局部表示	分布式表示
琥珀色	$[1, 0, 0, 0]^T$	$[1.00, 0.75, 0.00]^T$
天蓝色	$[0, 1, 0, 0]^T$	$[0.00, 0.5, 1.00]^T$
中国红	$[0, 0, 1, 0]^T$	$[0.67, 0.22, 0.12]^T$
咖啡色	$[0, 0, 0, 1]^T$	$[0.44, 0.31, 0.22]^T$

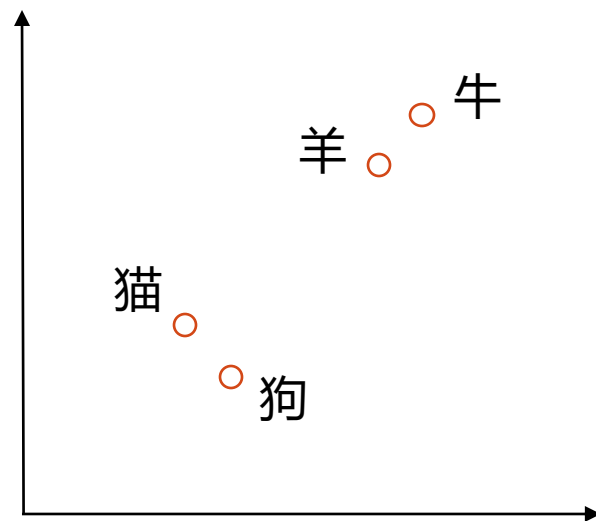


嵌入



嵌入 (Embedding)

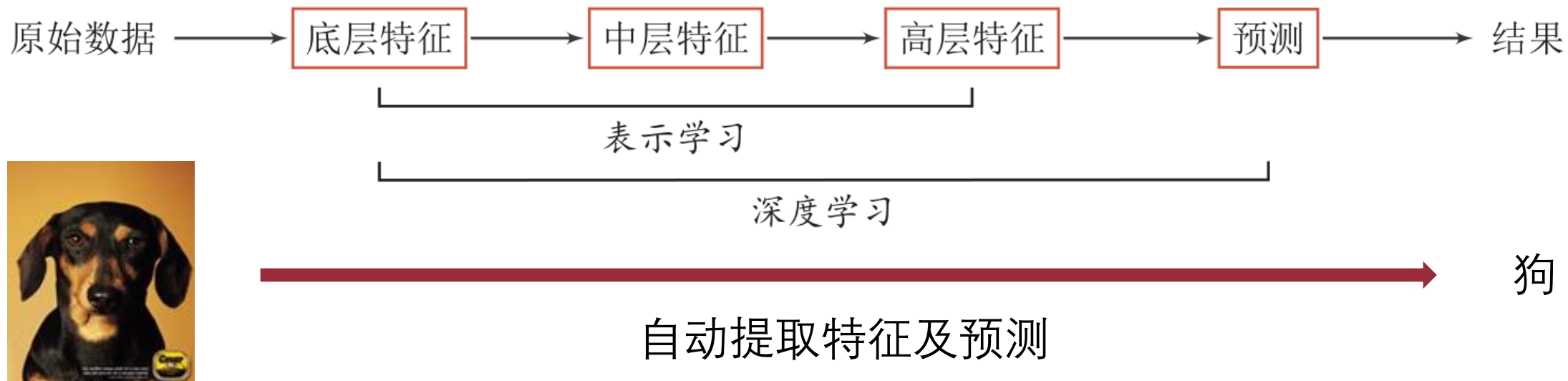
猫: [1, 0, 0, 0]
狗: [0, 1, 0, 0]
牛: [0, 0, 1, 0]
羊: [0, 0, 0, 1]



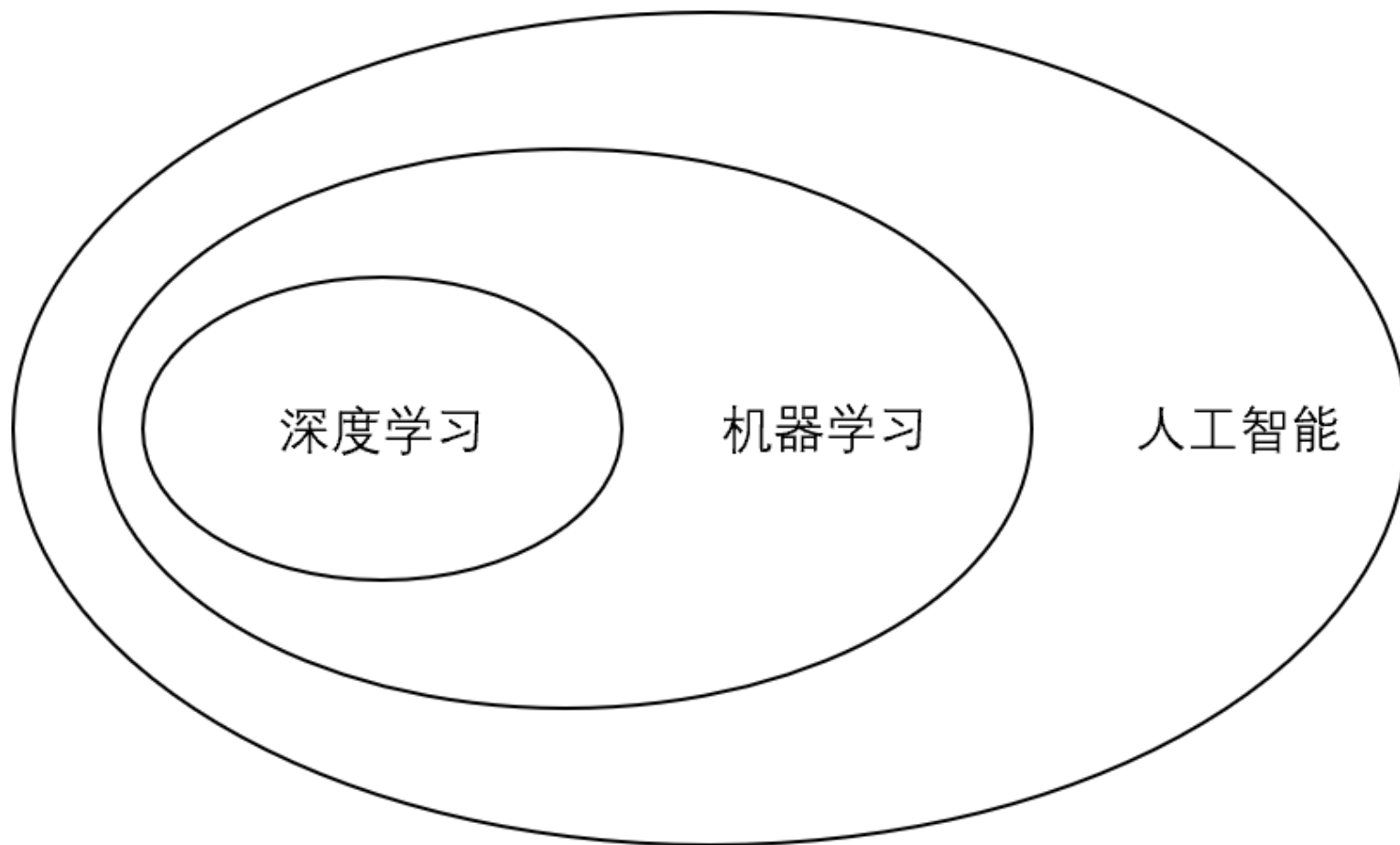
词嵌入

深度学习

- **深度学习 (Deep Learning)** 是一个复杂的机器学习算法，通过构建具有一定“深度”的模型，让模型来自动学习好的特征表示，从而最终提升预测或识别的准确性。
- 深度学习具备自动提取抽象特征的能力。



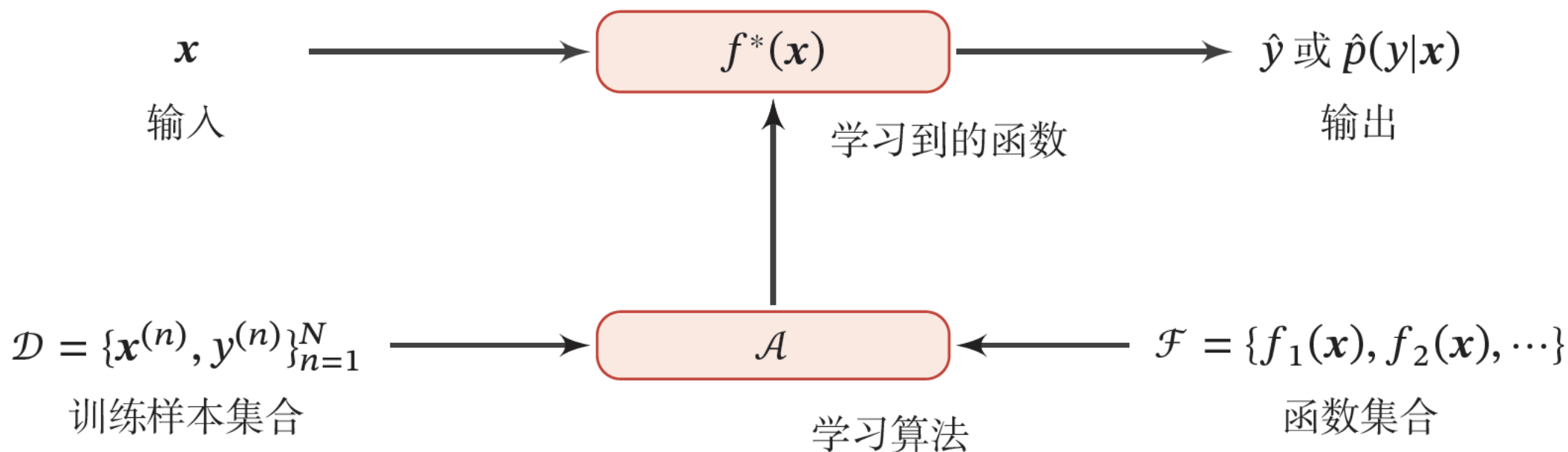
人工智能、机器学习与深度学习的关系



机器学习

模型 + 学习准则 + 优化方法

- 对于一个预测任务，输入特征向量为 \mathbf{x} ，输出标签为 y ，我们选择一个函数集合 \mathcal{F} ，通过学习算法 \mathcal{A} 和一组训练样本 \mathcal{D} ，从 \mathcal{F} 中学习得到函数 $f^*(\mathbf{x})$ ，这样对新的输入 \mathbf{x} ，就可以用函数 $f^*(\mathbf{x})$ 进行预测。



模型

- 对于样本空间中的样本 $(x, y) \in \mathcal{X} \times \mathcal{Y}$, 机器学习的目标是找到一个模型来近似 x 和 y 之间的真实映射函数 $g(x)$ 或真实条件概率分布 $p_r(y|x)$ 。
- 假设一个函数集合 \mathcal{F} , 称为假设空间, 然后通过观测其在训练集 \mathcal{D} 上的特性, 从中选择一个理想的假设 $f^* \in \mathcal{F}$ 。
- 假设空间 \mathcal{F} 通常为一个参数化的函数族

$$\mathcal{F} = \{f(x; \theta) | \theta \in \mathbb{R}^D\},$$

其中 $f(x; \theta)$ 是参数为 θ 的函数, 也称为模型, D 为参数的数量。

机器学习——模型

- 线性模型

$$f(\boldsymbol{x}; \theta) = \boldsymbol{w}^\top \boldsymbol{x} + b,$$

其中参数 θ 包含了权重向量 \boldsymbol{w} 和偏置 b 。

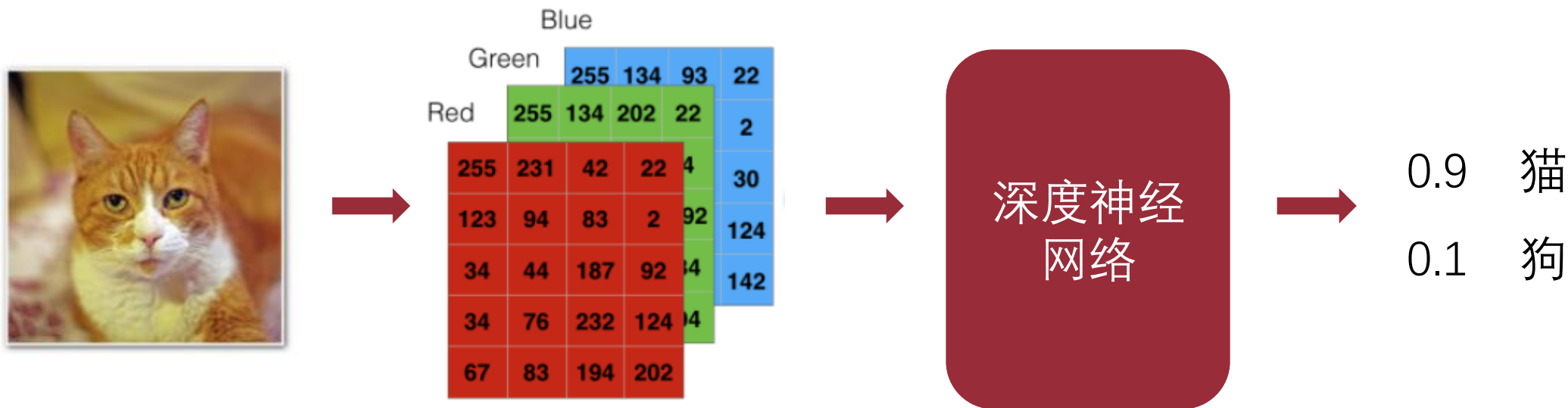
- 非线性模型

$$f(\boldsymbol{x}; \theta) = \boldsymbol{w}^\top \phi(\boldsymbol{x}) + b,$$

广义的非线性模型可以写为多个非线性基函数 $\phi(\boldsymbol{x})$ 的线性组合，参数 θ 包含了权重向量 \boldsymbol{w} 和偏置 b 。

深度学习——模型

- 大多数的深度学习可以看成是一个高度复杂的非线性回归模型。给定一个输入 X ，就可以给出一个预测值 $\hat{y} = f(X; \theta)$ ，其中 θ 是待估参数。在实际应用中，构建的神经网络极其复杂，因此 θ 的维度也非常高。



深度学习——模型

- 深度学习采用的模型主要是神经网络模型。
- 广义的非线性模型可以写为多个非线性基函数 $\phi(x)$ 的线性组合

$$f(x; \theta) = \mathbf{w}^\top \phi(x) + b$$

- 如果 $\phi(x)$ 本身为可学习的基函数, 比如

$$\phi_k(x) = h(\mathbf{w}_k^\top \phi'(x) + b_k), \forall 1 \leq k \leq K$$

其中 $h(\cdot)$ 为非线性函数, $\phi'(x)$ 为另一组基函数, \mathbf{w}_k 和 b_k 为可学习的参数, 则

$f(x; \theta)$ 就等价于神经网络模型。

学习准则

- 一个好的模型 $f(x, \theta^*)$ 应该在所有 (x, y) 的可能取值上都与真实映射函数 $y = g(x)$ 或与真实条件概率分布 $p_r(y|x)$ 一致。
- 模型 $f(x; \theta)$ 的好坏可以通过期望风险 $\mathcal{R}(\theta)$ 来衡量

$$\mathcal{R}(\theta) = \mathbb{E}_{(x,y) \sim p_r(y|x)}[\mathcal{L}(y, f(x; \theta))]$$

其中 $p_r(x, y)$ 为真实的数据分布, $\mathcal{L}(y, f(x; \theta))$ 为损失函数, 用来量化两个变量之间的差异。

损失函数

- 0-1损失函数

$$\mathcal{L}(y, f(\mathbf{x}; \theta)) = \begin{cases} 0 & \text{if } y = f(\mathbf{x}; \theta) \\ 1 & \text{if } y \neq f(\mathbf{x}; \theta) \end{cases} = I(y \neq f(\mathbf{x}; \theta))$$

- 不连续且导数为0，难以优化

- 平方损失函数

$$\mathcal{L}(y, f(\mathbf{x}; \theta)) = \frac{1}{2} (y - f(\mathbf{x}; \theta))^2$$

- 用在预测标签 y 为实数值的任务中，一般不适用于分类问题。

损失函数

- **交叉熵损失函数**：一般用于分类问题

- 假设样本的标签 $y \in \{1, \dots, C\}$ 为离散的类别，模型 $f(\mathbf{x}; \theta) \in [0, 1]^C$ 的输出为类别标签的条件概率分布 $p(y = c | \mathbf{x}; \theta) = f_c(\mathbf{x}; \theta)$ 。
- 用一个 C 维的one-hot向量 \mathbf{y} 来表示样本标签，看作样本标签的真实条件概率分布 $p_r(\mathbf{y} | \mathbf{x})$ 。
- 用交叉熵衡量两个概率分布的差异

$$\mathcal{L}(\mathbf{y}, f(\mathbf{x}; \theta)) = -\mathbf{y}^T \log f(\mathbf{x}; \theta) = -\sum_{c=1}^C y_c \log f_c(\mathbf{x}; \theta) = -\log f_{\mathbf{y}}(\mathbf{x}; \theta)$$

风险最小化准则

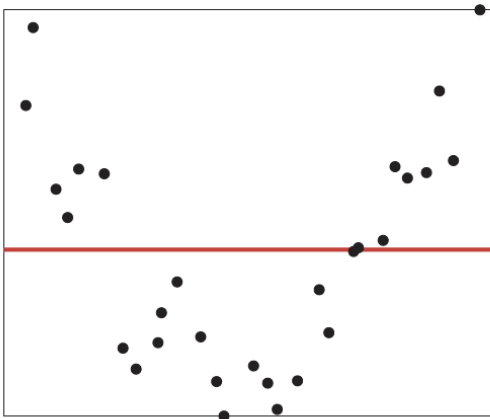
- 经验风险：在训练集 \mathcal{D} 上的平均损失

$$\mathcal{R}_D^{emp}(\theta) = \frac{1}{N} \sum_{n=1}^N \mathcal{L}(y^{(n)}, f(x^{(n)}; \theta))$$

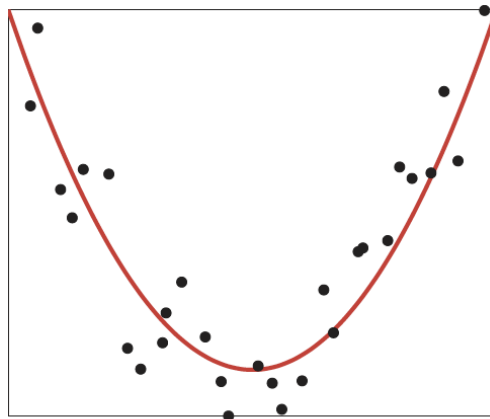
- 经验风险最小化

$$\theta^* = \underset{\theta}{\operatorname{argmin}} \mathcal{R}_D^{emp}(\theta)$$

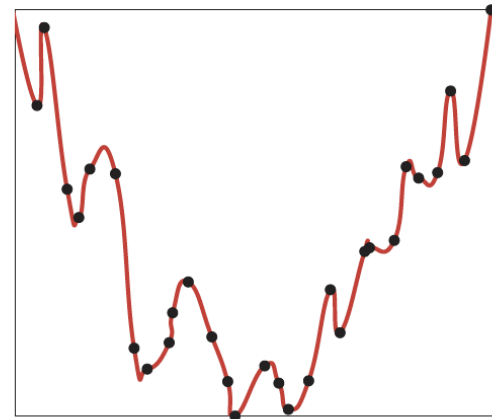
欠拟合



正常

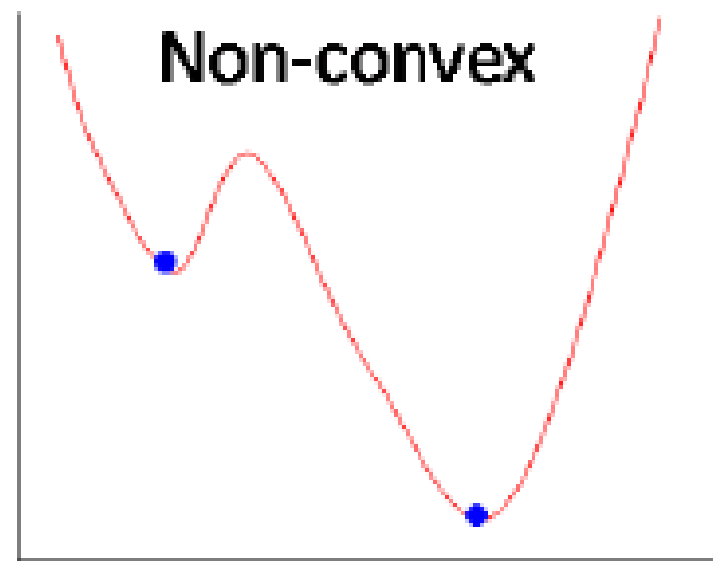
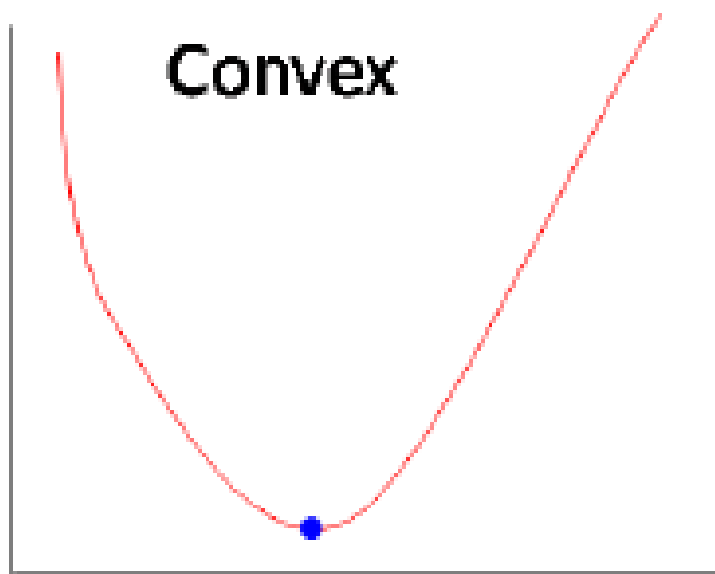


过拟合



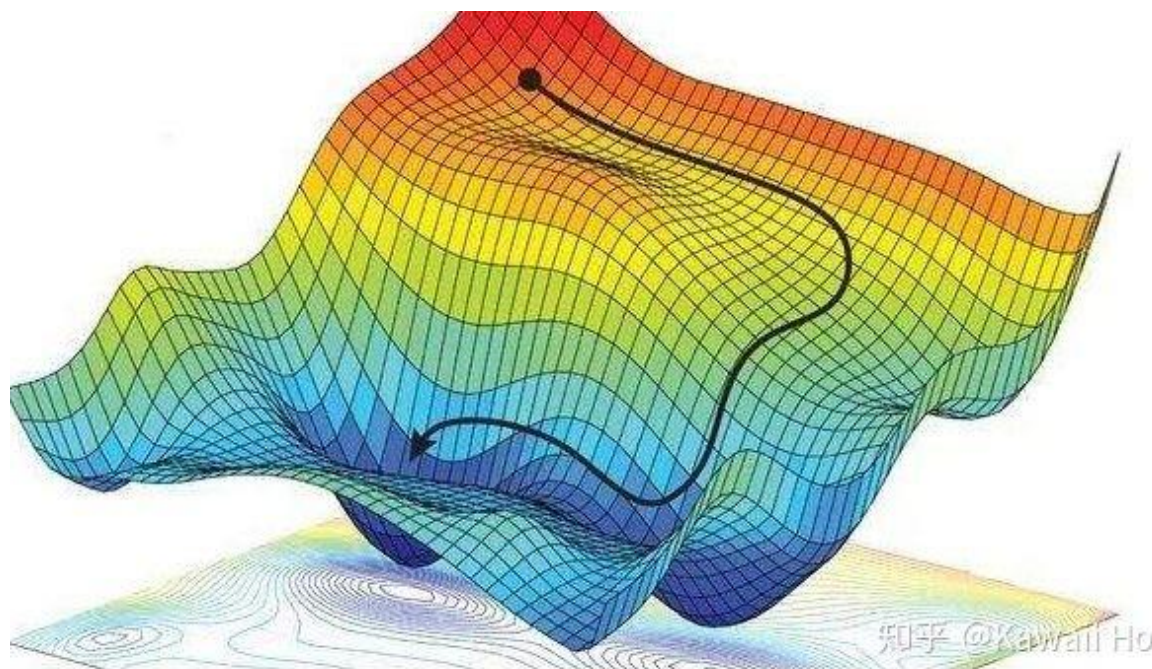
优化问题

- 机器学习问题转化为一个最优化问题



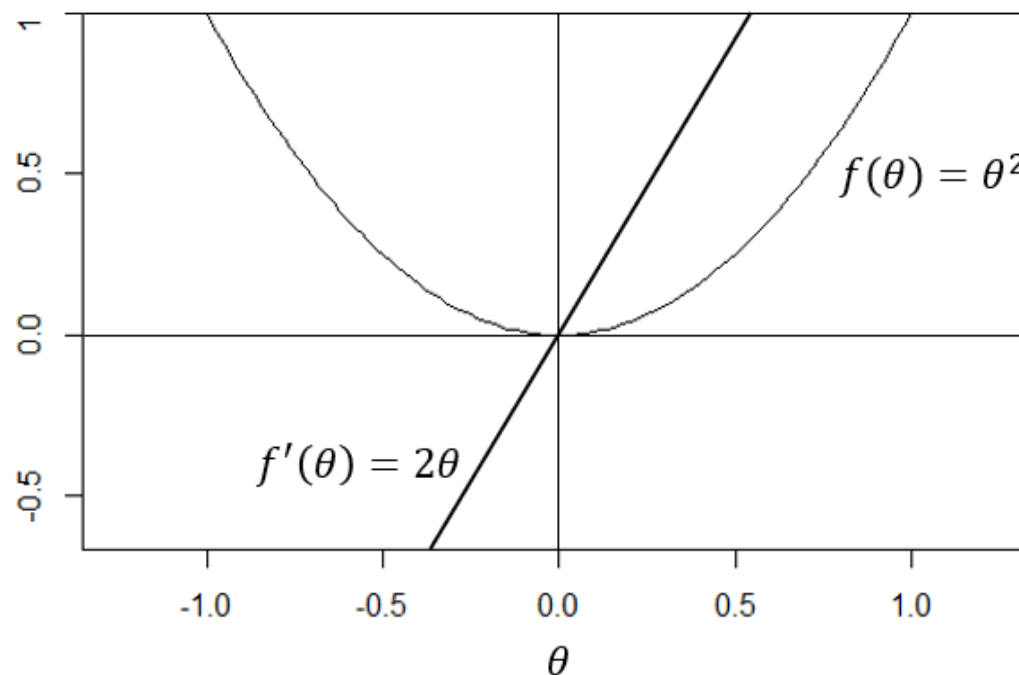
深度学习——优化算法

- 深度学习中，大多流行的优化算法通常基于梯度下降（gradient descent）。
- 梯度下降法：沿函数值下降最快的方向，改变 θ 而获得更小的 $f(\theta)$ 的技术。
- 梯度下降算法的直观理解：梯度下降就是寻找最陡峭的方向，也就是负梯度方向。



梯度下降法的理解

$$f(\theta) = \theta^2$$

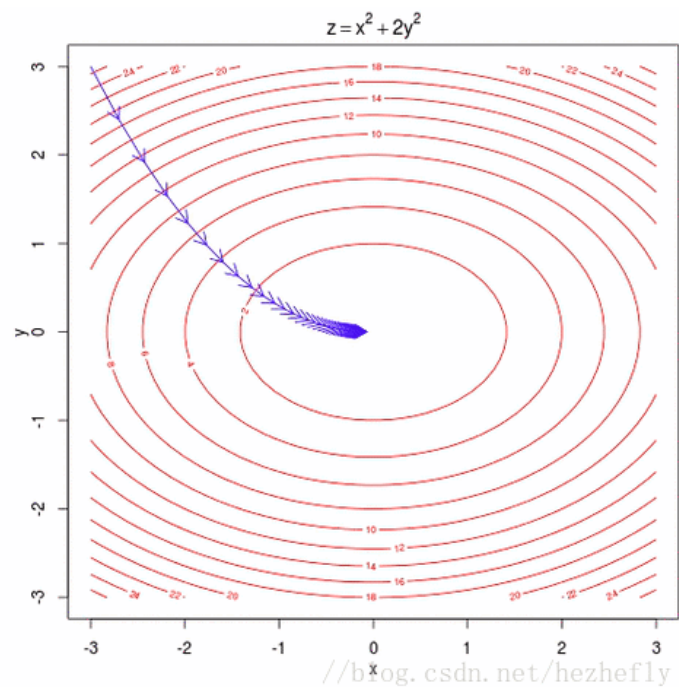
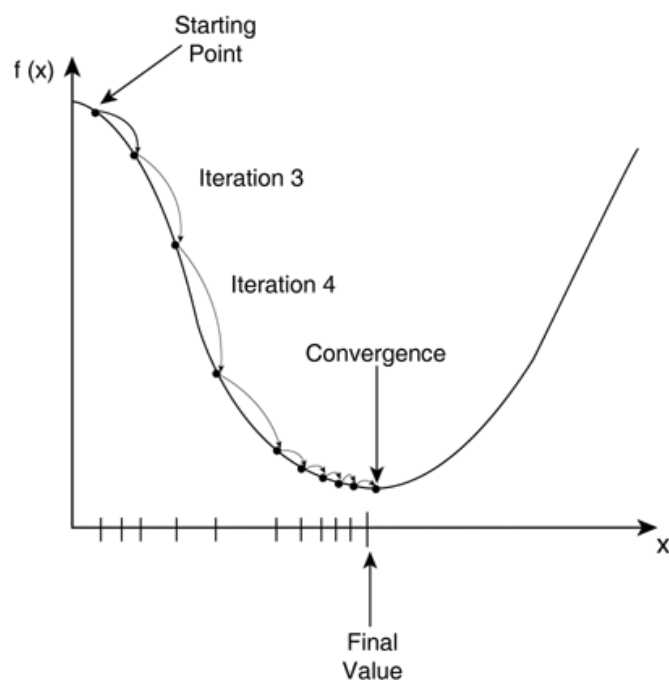


- 对于 $\theta > 0$, 存在 $f'(\theta) > 0$, 因此可以向左移动来减小 $f(\theta)$ 。
- 对于 $\theta < 0$, 存在 $f'(\theta) < 0$, 因此可以向右移动来减小 $f(\theta)$ 。
- 对于 $\theta = 0$, 存在 $f'(\theta) = 0$, 出现全局最小点, 梯度下降到这里停止。

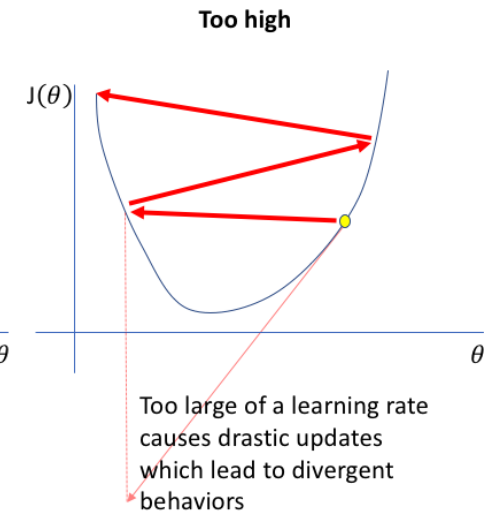
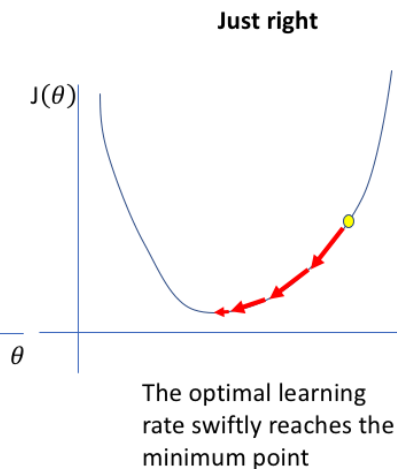
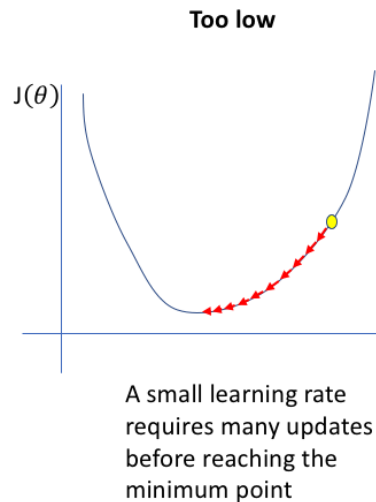
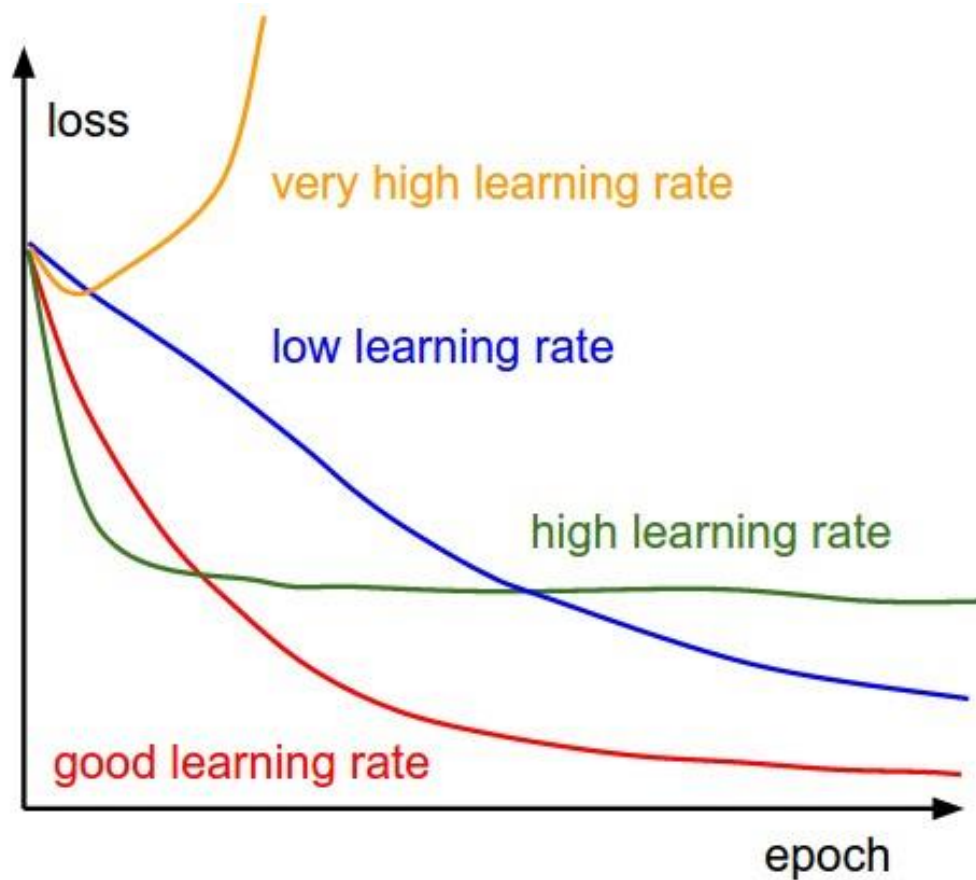
梯度下降法

- 利用学习率 α (Learning Rate) 来定义每次参数更新的幅度。
- 更新参数公式:

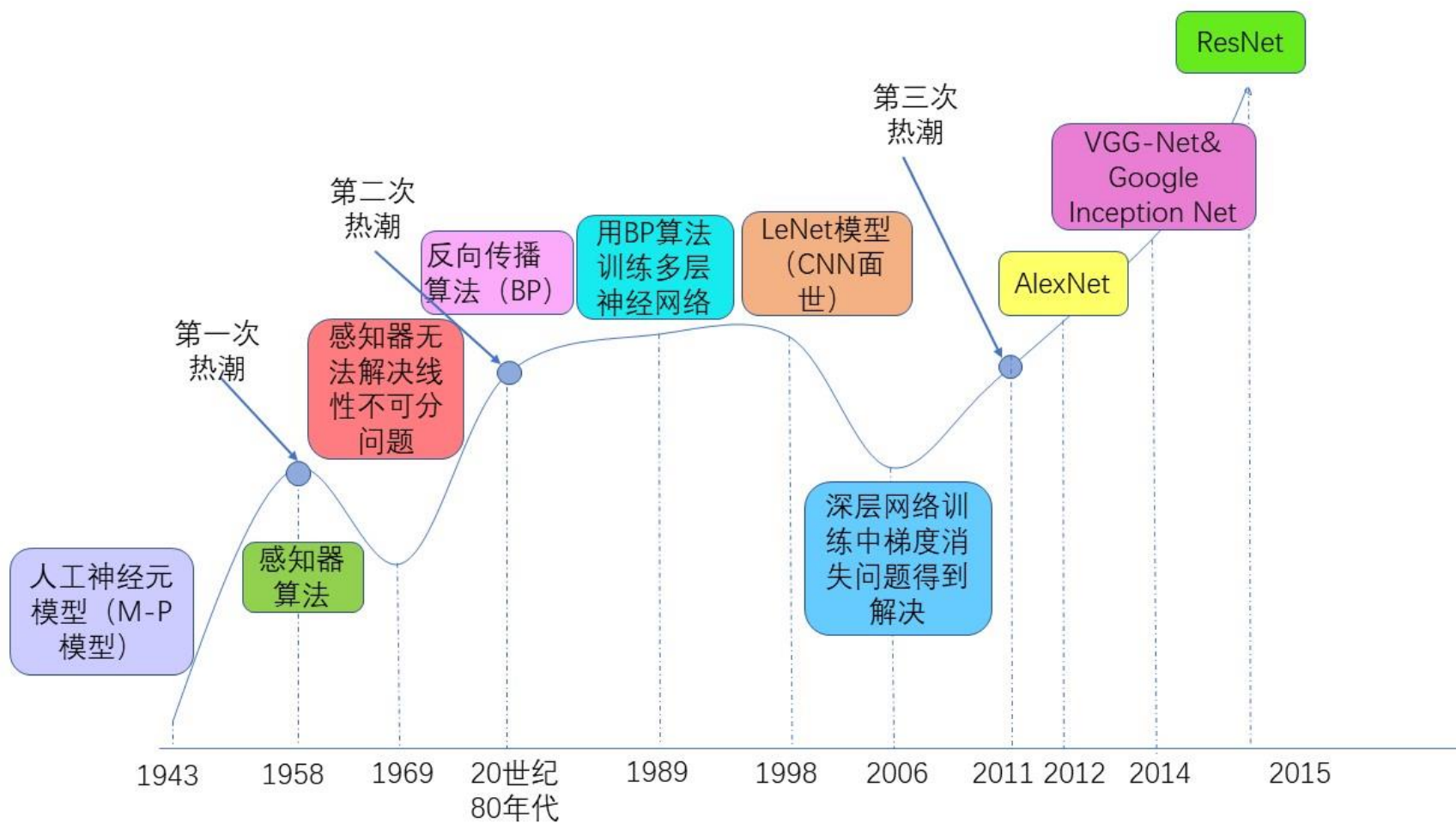
$$\theta_{n+1} = \theta_n - \alpha \frac{\partial f(\theta_n)}{\partial \theta}$$



学习率是十分重要的超参数



深度学习的发展历程

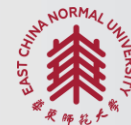


深度学习的发展历程

- 大数据+大规模算力

年代	数据规模	内存	每秒浮点运算
1970	100 （鸢尾花卉）	1 KB	100 KF (Intel 8080)
1980	1 K （波士顿房价）	100 KB	1 MF (Intel 80186)
1990	10 K （光学字符识别）	10 MB	10 MF (Intel 80486)
2000	10 M （网页）	100 MB	1 GF (Intel Core)
2010	10 G （广告）	1 GB	1 TF (Nvidia C2050)
2020	1 T （社交网络）	100 GB	1 PF (Nvidia DGX-2)

深度学习擅长的领域



ECNU

深度学习擅长的领域

- **图像处理**

- 图像分类、目标检测、图像分割

- **语音识别**

- **自然语言处理**

- 机器翻译、文本分类、自动问答

- **视频处理**

- **棋牌竞技**

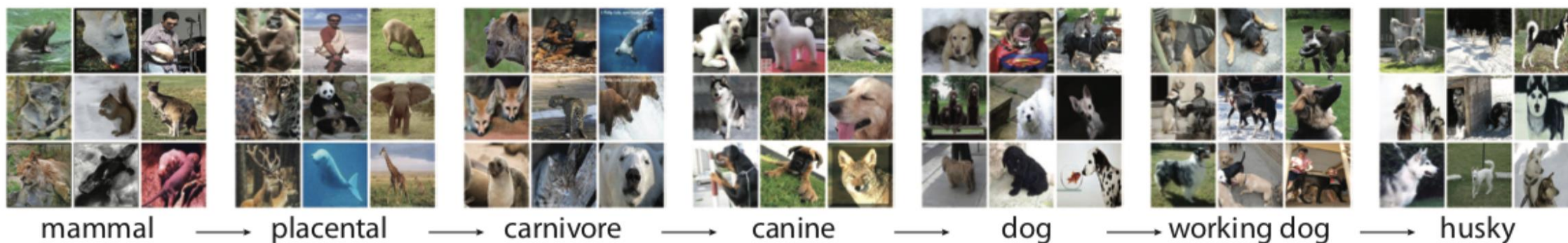
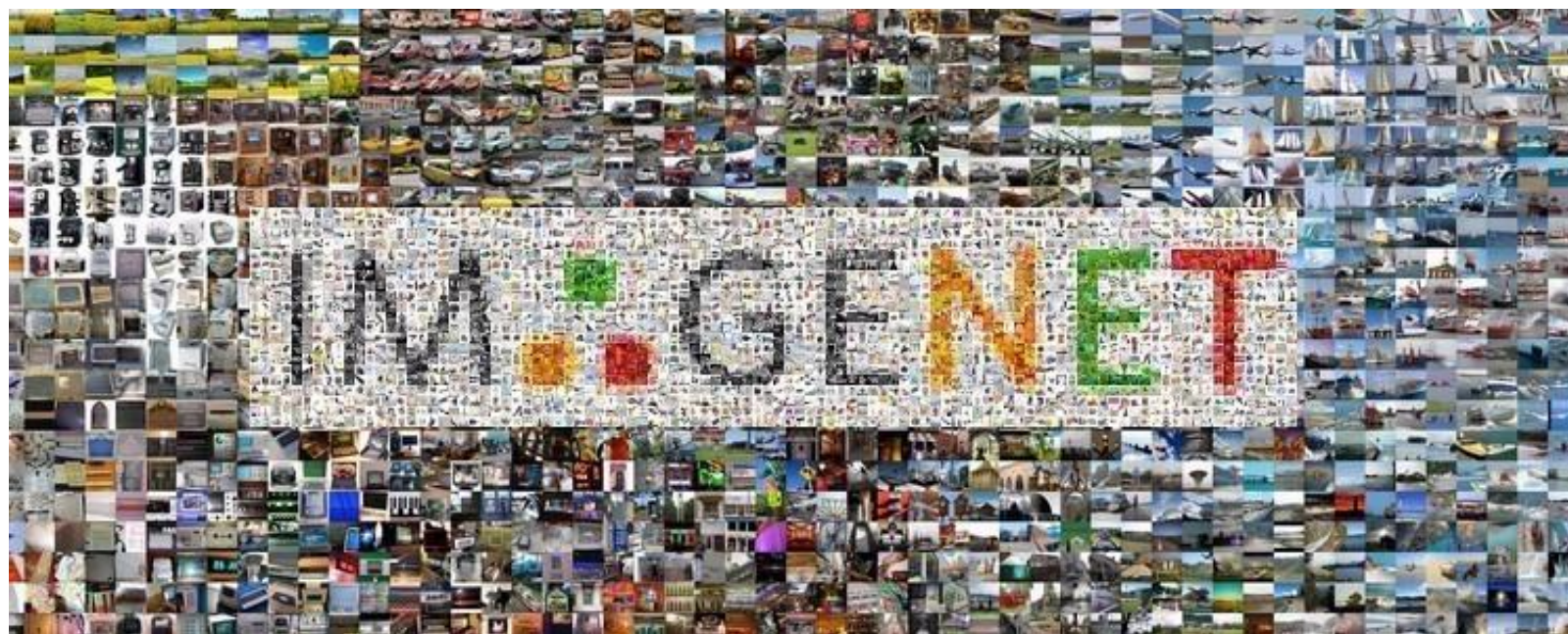
图像分类

Imagenet: A large-scale hierarchical image database

[J Deng](#), [W Dong](#), [R Socher](#), [LJ Li](#), [K Li](#)... - 2009 IEEE conference ..., 2009 - [ieeexplore.ieee.org](#)

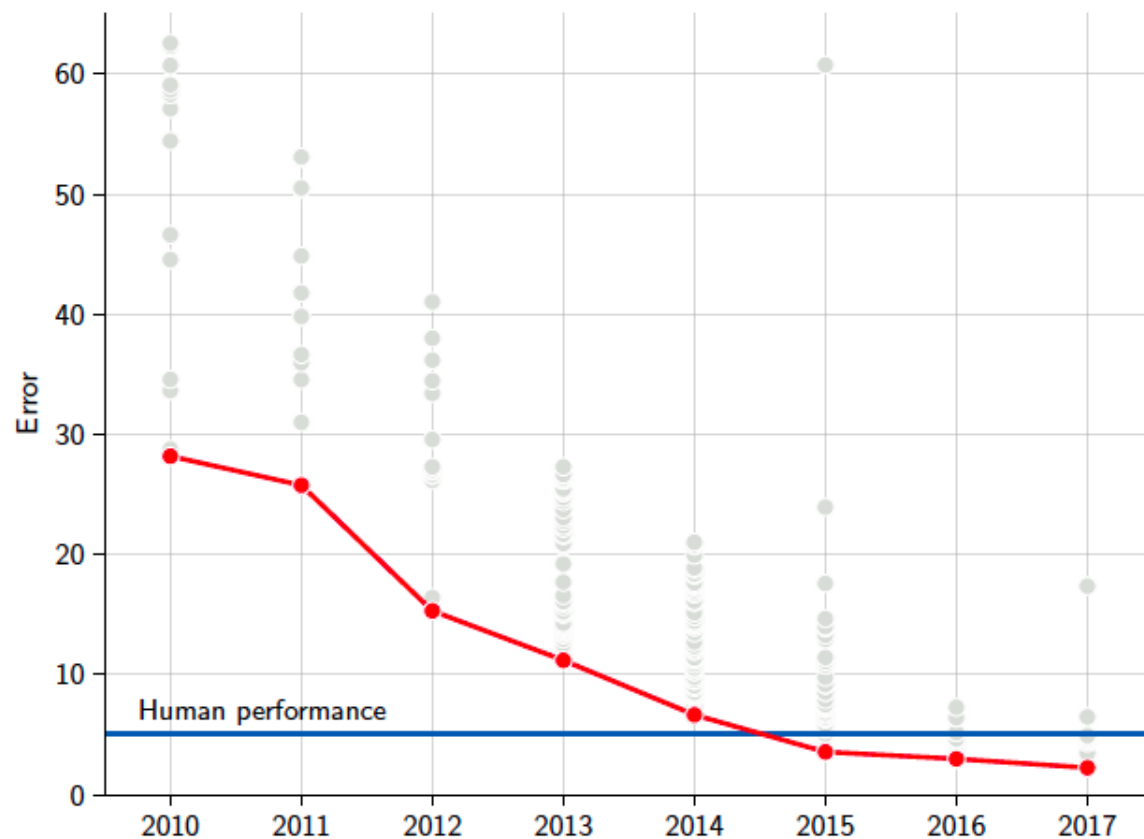
The explosion of image data on the Internet has the potential to foster more sophisticated and robust models and algorithms to index, retrieve, organize and interact with images and multimedia data. But exactly how such data can be harnessed and organized remains a ...

☆ 被引用 25514 相关文章 所有 31 版本



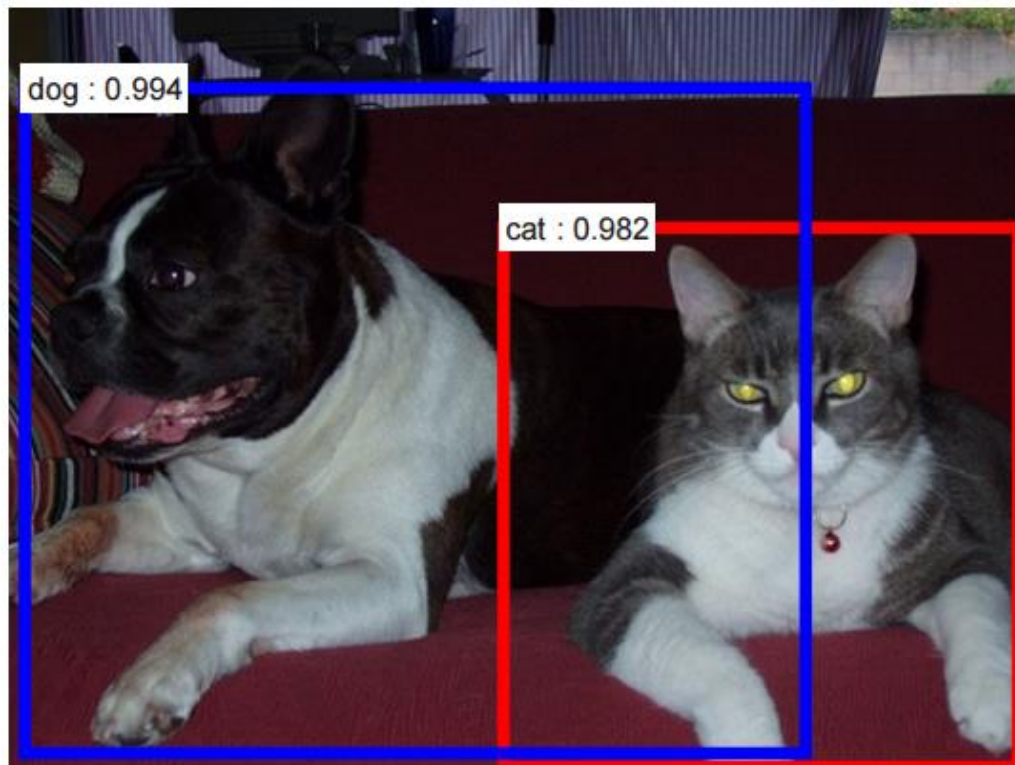
图像分类

- ImageNet 大规模视觉识别挑战 (ILSVRC)



D. Gershgorin. The data that transformed AI research|and possibly the world, July 2017.

目标检测和分割

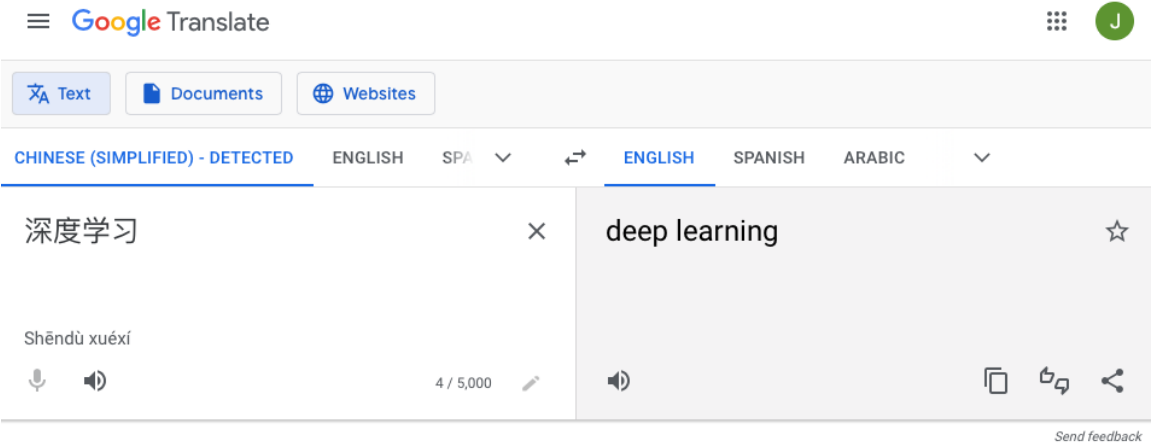


语音识别



自然语言处理

机器翻译



文本分类

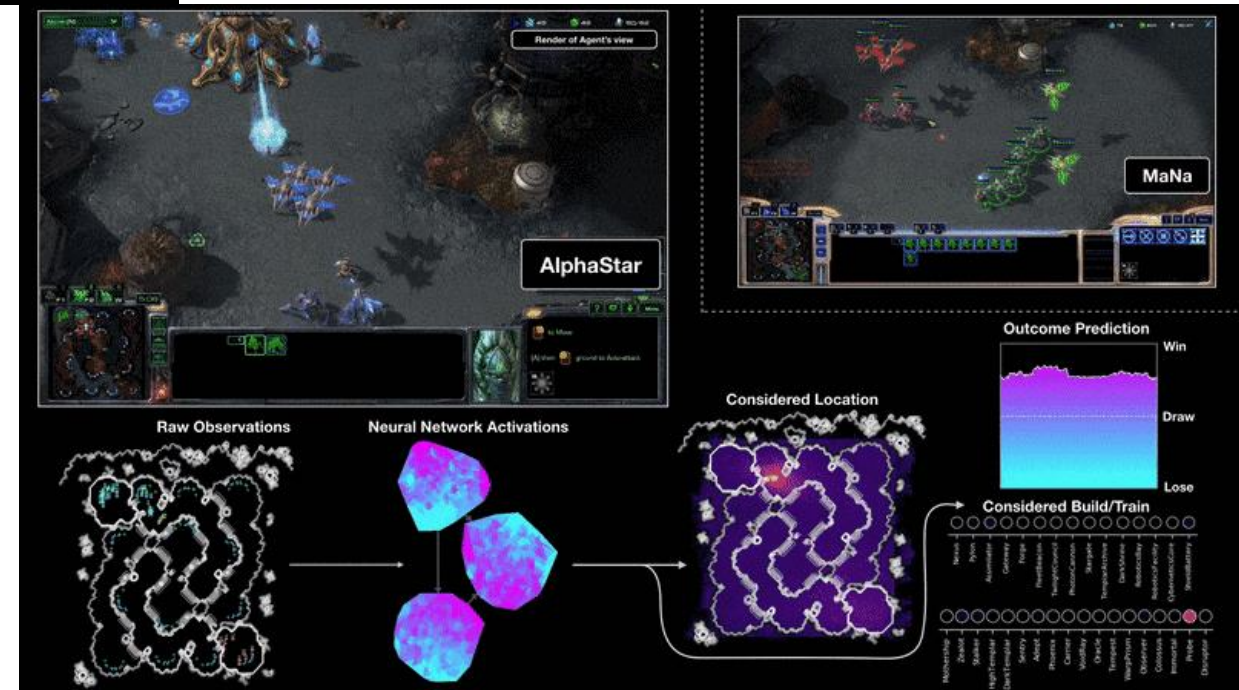
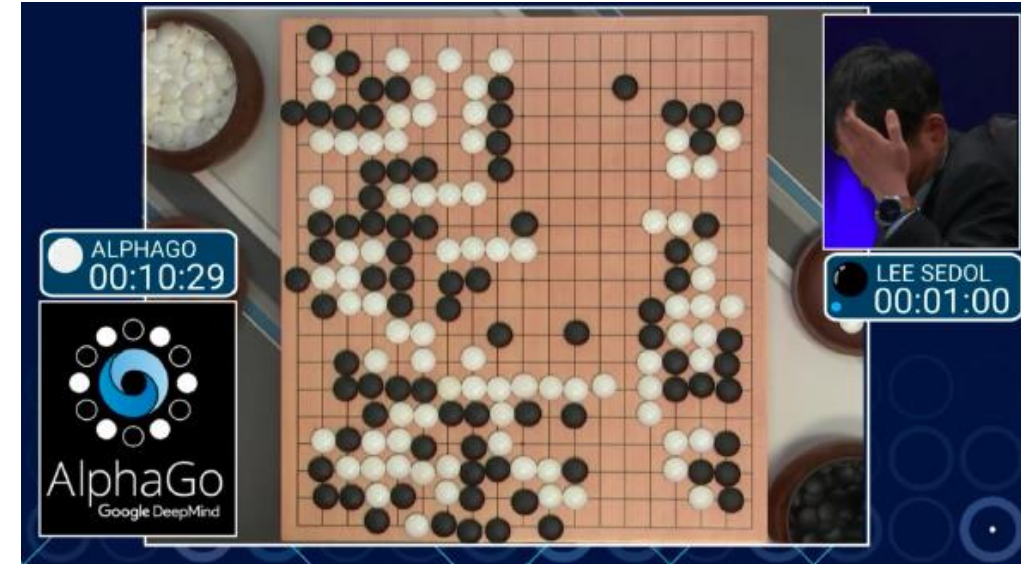


情感分析



棋牌竞技

- AlphaGo
- AlphaGo Zero
- AlphaZero
- MuZero
- AlphaStar



机器学习≈构建一个映射函数

- 图像识别

$$f(\text{9}) = \text{"9"}$$

- 语音识别

$$f(\text{语音波形}) = \text{"你好"}$$

- 机器翻译

$$f(\text{"你好!"}) = \text{"Hello!"}$$

- 围棋

$$f(\text{围棋棋盘}) = \text{"6-5" (落子位置)}$$

生成式人工智能 (AIGC)

Stable Diffusion

stability.ai



失眠的达芬奇71015

zzZ 我的第9个梦境

夜晚森林中的巨型南瓜屋，灯火通明，动态照明，概念艺术，幻想，震撼的视觉效果，8K



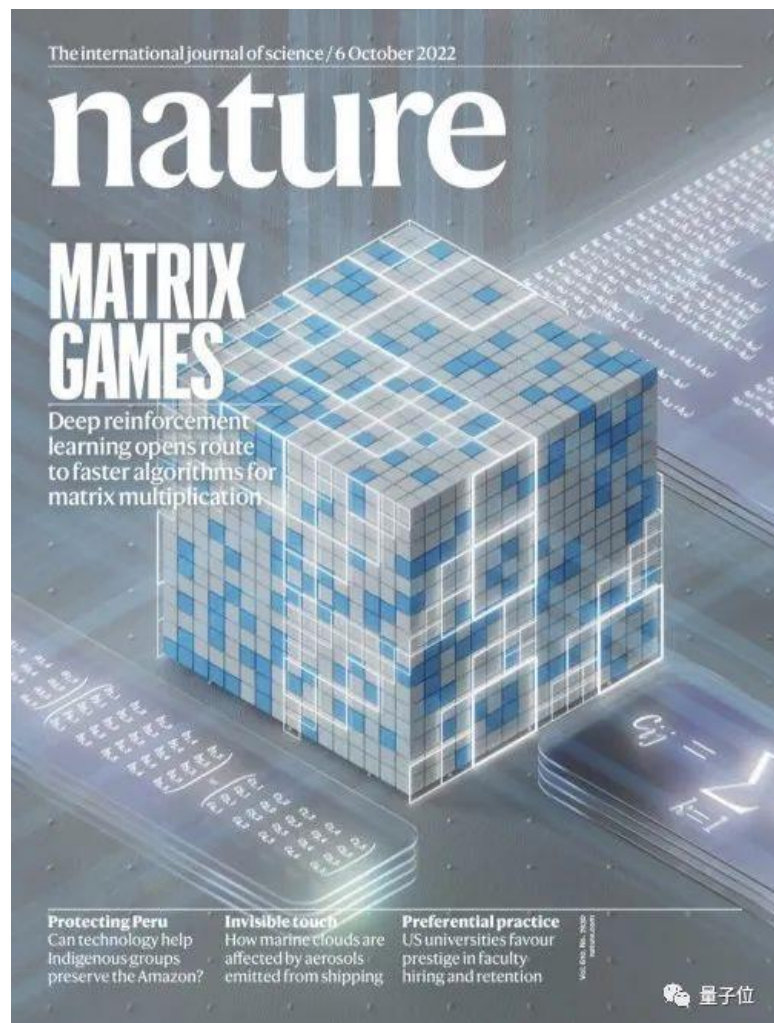
长按图片识别二维码
和我一起体验 AI智能作画

盗梦师

Sora

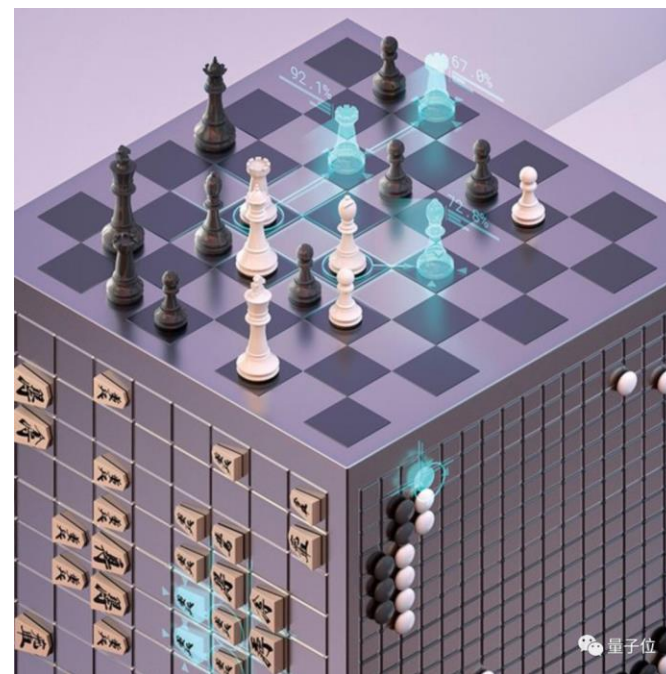


AlphaTensor



$$\begin{matrix} & \begin{matrix} 1 & 2 & \dots & n \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ \vdots \\ m \end{matrix} & \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ a_{31} & a_{32} & \dots & a_{3n} \\ \vdots & \vdots & \vdots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix} \end{matrix}$$

量子位



AlphaCode



2022-2-2

Competition-Level Code Generation with AlphaCode

Yujia Li*, David Choi*, Junyoung Chung*, Nate Kushman*, Julian Schrittwieser*, Rémi Leblond*, Tom Eccles*, James Keeling*, Felix Gimeno*, Agustin Dal Lago*, Thomas Hubert*, Peter Choy*, Cyprien de Masson d'Autume*, Igor Babuschkin, Xinyun Chen, Po-Sen Huang, Johannes Welbl, Sven Gowal, Alexey Cherepanov, James Molloy, Daniel J. Mankowitz, Esme Sutherland Robson, Pushmeet Kohli, Nando de Freitas, Koray Kavukcuoglu and Oriol Vinyals

*Joint first authors

Programming is a powerful and ubiquitous problem-solving tool. Developing systems that can assist programmers or even generate programs independently could make programming more productive and accessible, yet so far incorporating innovations in AI has proven challenging. Recent large-scale language models have demonstrated an impressive ability to generate code, and are now able to complete simple programming tasks. However, these models still perform poorly when evaluated on more complex, unseen problems that require problem-solving skills beyond simply translating instructions into code. For example, competitive programming problems which require an understanding of algorithms and complex natural language remain extremely challenging. To address this gap, we introduce AlphaCode, a system for code generation that can create novel solutions to these problems that require deeper reasoning. Evaluated on recent programming competitions on the Codeforces platform, AlphaCode achieved on average a ranking of top 54.3% in programming competitions with more than 5,000 participants. We found that three key components were critical to achieve good and reliable performance: (1) an extensive and clean competitive programming dataset for training and evaluation, (2) large and efficient-to-sample transformer-based architectures, and (3) large-scale model sampling to explore the search space, followed by filtering based on program behavior to a small set of submissions.

The screenshot shows a Codeforces problem page for '1553_D. Backspace'. The problem description is on the left, and the solution is on the right. The solution is written in Python and uses a stack to simulate the backspace operation.

Problem Description

1553_D. Backspace python correct Layer 18 play

Head 1 Head 2 Head 3 Head 4 Head 5 Head 6 Head 7 Head 8 Head 9 Head 10 Head 11 all none

Problem Description

// RATING: 1200
// TAGS: *special, implementation-
// LANGUAGE IS: cpp-
// CORRECT SOLUTION-
// You are given two strings s and t, both consisting of lowercase English-
// letters. You are going to type the string s character by character, from the-
// first character to the last one.-
// -
// When typing a character, instead of pressing the button corresponding to it,-
// you can press the "Backspace" button. It deletes the last character you have-
// typed among those that aren't deleted yet (or does nothing if there are no-
// characters in the current string). For example, if s is "abcbd" and you press-
// Backspace instead of typing the first and the fourth characters, you will get-
// the string "bd" (the first press of Backspace deletes no character, and the-
// second press deletes the character "c"). Another example, if s is "abcac" and-
// you press Backspace instead of the last two letters, then the resulting text-
// is "a" -
// -
// Your task is to determine whether you can obtain the string t, if you type-
// the string s and press "Backspace" instead of typing several (maybe zero)-
// characters of s.-
// Input-
// -
// The first line contains a single integer q (1 ≤ q ≤ 10^5) - the number of-
// test cases.-
// -
// The first line of each test case contains the string s (1 ≤ |s| ≤ 10^5). Each-
// character of s is a lowercase English letter.-
// -
// The second line of each test case contains the string t (1 ≤ |t| ≤ 10^5).-
// Each character of t is a lowercase English letter.-
// -
// It is guaranteed that the total number of characters in the strings over all-
// test cases does not exceed 2 · 10^5.-
// -
// Output-
// -
// For each test case, print "YES" if you can obtain the string t by typing the-
// string s and replacing some characters with presses of "Backspace" button, or-
// "NO" if you cannot.-
// -
// You may print each letter in any case (YES, yes, Yes will all be recognized-
// as positive answer, NO, no and nO will all be recognized as negative answer).-

Solution

```
t=int(input())
for i in range(t):
    s=input()
    t=input()
    a=[]
    b=[]
    for j in s:
        a.append(j)
    for j in t:
        b.append(j)
    a.reverse()
    b.reverse()
    c=[]
    while len(b)!=0 and len(a)!=0:
        if a[0]==b[0]:
            c.append(b.pop())
            a.pop()
        elif a[0]!=b[0] and len(a)!=1:
            a.pop()
            a.pop()
        elif a[0]!=b[0] and len(a)==1:
            a.pop()
    if len(b)!=0:
        print("YES")
    else:
        print("NO")
```