

基础知识 什么是AI人工智能 是利用数字计算机或者数字计算机控制的机器模拟、延伸和扩展人的智能，感知环境、获取知识并使用知识获得最佳结果的理论、方法、技术及应用系统。AI的四大主流流派1.符号主义（Symbolism）例：知识图谱；2.连接主义（Connectionism）例：深度学习神经网络；3.行为主义（Behaviourism）例：机器人；4.统计主义（Statisticism）例：机器学习。AI、机器学习、深度学习三者之间的异同和关联1三者属于包含关系：AI包含机器学习，机器学习包含深度学习，深度学习学习主要指机器学习中的深度学习神经网络。2在机器学习中，需要先进行人工特征提取，再定义Model；而大部分深度学习只构建一个端到端的模型，没有人工特征提取。3在深度学习中，有时需要先进行数据清洗、格式转换、特征提取（卷积）等操作，再将数据喂给全连接神经网络。**机器学习三步曲** 1.定义Model；2.定义Loss/cost/error/objective function；3.找出最佳参数。使用反向传播不断优化参数，从而找出最佳参数。**模型** 深度学习模型Model=Neural Network Structure；2.定义这个Model好坏，使用合适的Loss Function来衡量损失；（3）找出最佳参数。使用反向传播不断优化参数，从而找出最佳参数。**机器学习项目的通用工作流程** 定义问题：软件架构设计、确定评价指标、2获取数据：自动化方式3研究数据：可视化方式，相关性研究等4准备数据：数据清理、特征选择及处理5研究模型：确定评估方法、列出可用的模型并训练、选择最有希望的3到5个模型6做模型：寻找最佳超参数，模型融合，评估泛化性能7展示解决方案：将工作进行文档化总结展示8启动、监视、维护系统；投入使用。**神经网络通用工作流程** 1定义问题、收集数据集2选择衡量成功的指标3确定评估方法4留出验证集K折验证。你应该将哪一部分数据用于验证4准备数据5开发比基准更好的模型6扩大模型规模：开发过拟合的模型7模型正则化与调节超参数。基于模型在验证数据上的性能。**深度学习和机器学习有什么不同** 机器学习：利用计算机、概率论、统计学等知识，输入数据，让计算机学会新知识。机器学习的过程，就是训练数据去优化目标函数。深度学习：是一种特殊的机器学习，具有强大的能力和灵活性。它通过学习将世界表示为嵌套的层次结构，每个表示都和更简单的特征相关，而抽象的表示则用于计算更抽象的表示。传统的机器学习需要定义一些手工特征，从而有目的的去提取目标信息，非常依赖任务的特异性以及设计特征的经验。而深度学习可以从大数据中先学习简单特征，并从此逐渐学习到更为复杂抽象的深层特征，不依赖于人工的特征工程，这也是深度学习在大数据时代受青睐的一大原因。**机器学习的定义**对于某类任务T和性能度量P，一个计算机程序被认为可以从经验E中学习是指，通过经验E改进后，它在任务T上由性能度量P衡量的性能有所提升。任务T：分类，翻译等机器学习的目标任务。性能度量P：最常用的有准确率，召回率，F值等，根据任务T不同，P也不尽相同。经验E：训练集。**什么场合下需要使用机器学习**？需要进行大量手工调整或需要拥有长串规则才能解决的问题；机器学习算法通常可以简化代码，提高性能。问题复杂，传统方法难以解决；最好的机器学习方法可以找到解决方案。环境有波动；机器学习算法可以适应新数据。洞察复杂问题和大数量数据。**模型选择的三个维度**：Scenario：数据集是否是带标签？是否是强化学习任务？要预测的值（即输出值）是否是数值型变量？是否连续？Method：选择的模型。**机器学习的分类及差异**。1**监督学习**的目标是建立一个学习过程，将预测结果与“训练数据”（即输入数据）的实际结果进行比较，不断的调整预测模型，直到模型的预测结果达到一个预期的准确率。例：手写字符识别、肿瘤分类、预测天气、支持向量机、线性判别。特点：训练样本数据和待分类的类别已知，且训练样本数据皆作为标签数据。2**非监督学习**从无标记的训练数据中推断结论，它可以在探索性分析阶段用于发现隐藏的模式或对于数据进行分组。例：聚类分析、主成分分析。特点：训练样本数据和待分类的类别已知，但训练样本数据皆作为非标签数据。3**半监督学习**的训练数据通常是少量有标记数据及大量未标记数据；介于无监督学习和监督学习之间。例：聚类预测、流形推理。特点：训练样本数据和待分类的类别已知，然而训练样本既有标签数据，也有非标签数据。4**强化学习**的输入数据作为对模型的反馈，强调如何在未知的环境下行动，以取得最大化的长期利益。与监督式学习之间的区别在于，它并不需要知道正确的输入/输出，也不需要精确校正优化行为。强化学习更加专注于在线规划，需要在探索（在未知的领域）和遵从（取得现有知识）之间找到平衡。例：学习下围棋、打星际争霸、DotA2、双人德州扑克。（全部都是1v1的场景，目前强化学习领域对于多人博弈研究的很少）特点：决策流程，激励系统，学习一系列的行动。**5迁移学习**是通过从已学习的相关任务中迁移其知识来对需要学习的新任务进行提高。例：牛津的VGG模型、谷歌的Inception模型和word2vec模型、微软的ResNet模型。特点：需求的训练数据集较小，训练时间较小，可以方便的进行迁移以满足个性化。6**弱监督学习**，特点：弱监督学习可以看成是有多个标记的数据集合，交集可以是空集，单个元素，或包含多种情况（没有标记，有一个标记，和有两个标记）的多个元素。数据集的标签是不可靠的，这里的不可靠可以是标记不正确，多种标记，标记不充分，局部标记等。已知数据和其一一对应的弱标签，训练一个智能算法，将输入数据映射到一组更弱的标签的过程。标签的强到弱指的是标签蕴含的信息量的多少，比如相对于分标的标签来说，分类的标签就是弱标签。算法举例：举例，给出一张包含气球的图片，需要得出气球在图片中的位置及气球和背景的分割线，这就是已知弱标签学习强标签的问题。**机器学习三类数据**：训练集：利用训练集上的误差来训练模型，学习参数。验证集：利用验证数据上的误差来挑出各选项模型，确定各个备选模型的最佳性能，并从中挑选出最佳模型。测试集：用于估计泛化误差，评估模型的表现。**机器学习两个计算过程**：1学习训练：实例化；目标是找到最佳函数来代入输入输出之间的关系，理想的最佳函数，应在优化和泛化的平衡点上，应刚好在欠拟合和过拟合的交界线上。2推理：预测泛化：将找到的最佳函数应用到测试数据，真实数据上，做出推理预测。**机器学习两类函数**：最佳函数：（模型训练的终极目标是找到这个函数）成本/代价/损失函数：用来定义训练过程中找到的函数的好坏程度。**机器学习两类参数**：参数：模型训练过程中自动找到，超参数：程序员来选择，指定半自动：网格搜索、随机搜索。手动：程序员设定。**训练过程中，模型的状态**：欠拟合过拟合正好。机器学习两个评价指标：模型评估标准：在训练数据集和验证数据集上；优化指标，即损失函数。通常用于模型的训练中。损失函数的取值 越小，则越好。模型评估标准：在测试集上；验收指标：用于测试模型的泛化性能。比如：样本平衡的问题中常用的验收指标 ROC AUC但是，ROC AUC不能直接被优化，故不能作为损失函数，而是将 ROC AUC的替代指标（交叉熵作为优化指标，即损失函数。**机器学习评估方法**：留出验证集，交叉验证（CrossValidation）。**Epoch**：对整个数据集进行一次forward和backward过程，被称为一个Epoch。**Batch size**：每一批的大小，在一次前向传播、反向传播中，批量的大小。**Batch**：使用训练集中一小部分样本对模型权重进行一次反向传播的参数更新，这一小部分样本被称为一个Batch。**Iteration**：梯度更新的次数，也就是总共已经使用了多少个batch。**Batch num**：每次使用Batch size个数据，期望预测与真实标记的误差（偏差）越大偏离理论值越大，在一个训练集D上模型f对测试样本x的预测输出为f(x;D)，那么学。使用样本数相同的不同训练集产生的方差为：var(X)=E[Df(X;D)](这里用偏差的平方来表示偏差的计算公式Bias 2(x)=E(f(x)-y) 2。2.方差variance 预测模型的离散程度，方差越大离散程度越大。那么学。使用样本数相同的不同训练集产生的方差为：var(X)=E[Df(X;D)](这里用偏差的平方来表示偏差)期望误差和经验误差之间的差异。用来表示泛化性能的优劣。训练集验证于无穷大时，泛化误差趋向下0。**什么是交叉验证** 将数据集D划分成k个大小相似的互斥子集。每次用k-1个子集作为训练集，余下的1个子集做测试集，最终返回k个训练结果的平均值。交叉验证法评估结果的稳定性和真实性很大程度上取决于k的取值。适用于数据集不是特别大时。**几种正则化** L1-Norm：参数绝对值之和。L1(θ)=L(θ)+λ||θ||1，其中||θ||1=Σ_i=1^n|θ_i|作用：特征筛选，把得分最低的参数为稀疏性向量L2-Norm：参数平方的和。L2(θ)=L(θ)+λ||θ||2，其中||θ||2=Σ_i=1^nθ_i2作用：通过限制模型复杂度，从而避免过拟合，提高泛化能力。Elastic Net:(L1+L2) L(θ)=(L1(θ)+L2(θ))/2。

如何区分回归问题与分类问题 输出值可以有限约束性出来就是分类问题，否则就是回归问题。**为什么不可以用线性回归的模型解决分类问题** 当样本分布比较复杂时，线性回归法做出准确的分类。有些问题不是线性可分的（异或问题）。可以使用线性回归解决二分类问题，但是无法解决多分类的问题。原因是：需要人为的划定区间，引入误差。在损失函数上，MSE无法很好地评估模型。举例：假如有5元分类问题，class0 class4，如果正确分类是class4，但是预测到了class1，和预测到了class3。根据MSE，损失函数是不同的。这就无法很好地区分了。**在逻辑回归中，为何使用Cross Entropy作为损失函数而不使用MSE**，如果是交叉熵，距离target越远，微分值就越大，就可以做到距离target越远，更新参数越快。而平方误差在距离target很远的時候，微分值非常小，会造成收敛的速度非常慢，这就是很差的后果了。离目标目标值相差的绝对值， ΔL 很小 MSE 更新很慢。逻辑回归公式 $\sigma(z) = 1 / (1 + e^{-z})$ ， $z(x) = wx + b$ ，分别对w和b求导，有 $\partial \sigma(z) / \partial w = \sigma(z)(1 - \sigma(z)) * x$ 和 $\partial \sigma(z) / \partial b = \sigma(z)(1 - \sigma(z))$ 。如果使用MSE作为损失函数的话，那写出来是这样的 $(L(w, b) = 1/2m \sum_{i=1}^m (y^{(i)} - \sigma(w^{(i)}))^2)$ 。使用梯度下降来求导时，需要计算 $L(w, b)$ 对w和b的偏导数： $\partial L(w, b) / \partial w = 1/m \sum_{i=1}^m (\sigma(w^{(i)}) - y^{(i)}) \sigma(w^{(i)}) (1 - \sigma(w^{(i)}))$ 和 $\partial L(w, b) / \partial b = 1/m \sum_{i=1}^m (\sigma(w^{(i)}) - y^{(i)})$ 。从之前的w换成b，最后的x不要，所以在 $\sigma(x)$ 接近于1或者0的时候，也就是预测的结果和真实结果很相近或者很不相近的时候， $\sigma(x)$ 和 $1 - \sigma(x)$ 中总有一个会特别小，这样会导致梯度很小，从而使使得优化速度大大减缓。而当使用Cross Entropy作为损失函数时，损失函数为： $L(w, b) = 1/m \sum_{i=1}^m (-y^{(i)} \log(\sigma(w^{(i)})) - (1 - y^{(i)}) \log(1 - \sigma(w^{(i)})))$ 。(w,b)分别对w和b求偏导，结果如下 $\partial L(w, b) / \partial w = 1/m \sum_{i=1}^m (\sigma(w^{(i)}) - y^{(i)})$ ， $\partial L(w, b) / \partial b = 1/m \sum_{i=1}^m (\sigma(w^{(i)}) - y^{(i)})$ 。这样梯度始终和预测值与真实值之差挂钩，预测值与真实值相差很大时，梯度也会很大，偏离很小时，梯度会很小。所以我们更倾向于使用Cross Entropy而不使用MSE。4. 级联逻辑回归(Cascading logistic regression model)模型是啥？为什么要引入这个概念？级联逻辑回归模型是将很多的逻辑回归接到一起，以进行特征转换再用一个逻辑回归来进行分类。级联逻辑回归模型是神经网络的雏形。5. Loss Function - 损失函数 回归任务假设函数： $h_{w,b}(x) = wx + b$ 二分类任务假设函数： $h_{w,b}(x) = \sigma(wx + b)$ 。在回归任务中，多使用均方误差作为损失函数；上面有。在分类任务中，多使用交叉熵作为损失函数，上面有。6. Sigmoid和Softmax Sigmoid Function 不具体表示哪一个函数，而是表示S型函数，常用的有逻辑函数 $\sigma(z)$ 上面有。而Softmax，或称归一化指数函数，是逻辑函数的一种推广，二分类情况下，Softmax退化到逻辑函数。该函数的形式通常按下面的式子给出： $\sigma_j(z) = e^{z_j} / \sum_{k=1}^K e^{z_k}$ ， $j = 1, \dots, K$ 。**Logistic 回归与Softmax 回归联系与区别** Logistic 回归与 Softmax 回归是两个基础的分类模型，虽然softmax字像是回归模型，实际上并非如此如。Logistic 回归，Softmax 回归以及线性回归都是基于线性模型。其实Softmax 就是 Logistic 的推广，Logistic 一般用于二分类，而softmax 是多分类。softmax是lr的升级版，可以用来求多类问题划分，lr是二分项问题。softmax的权重w是矩阵，lr是向量。softmax输出是属于每个类的概率值，概率最大的就是其对应类。lr就两个值，是和不是（0-1 || -1-1 一般是这样）。只是应用场景不同，其他的没啥区别。

分类 分类器的性能评估指标 True Positive（真正 TP）被模型预测为正的样本数。True Negative（真负，TN）被模型预测为负的样本数。False Positive（假正，FP）被模型预测为正的样本数。False Negative（假负，FN）被模型预测为负的正样本。混淆矩阵竖轴为预测类别，横轴为实际类别。**评价指标** 1 准确率或者正确率（accuracy）正确率是我们最常见的评价指标，accuracy = (TP+TN)/(P+N)，正确率是被分对的样本数在所有样本数中的占比，通常来说，正确率越高，分类器越好。2 错误率（error rate）错误率则与正确率相反，描述被分类器错分的比例，error rate = (FP+FN)/(P+N)，对某一个实例来说，对分与错是互斥事件，所以accuracy = 1 - error rate。3 灵敏度（sensitivity）sensitivity = TP/P，表示的是所有正例中被分对的比例，衡量了分类器对正例的识别能力。4 特异性（specificity）specificity = TN/N，表示的是所有负例中被分对的比例，衡量了分类器对负例的识别能力。5 精确率或者查准率或者精度（precision）precision=TP/(TP+FP)，精度是精确性的度量，表示被分为正例的实例中实际为正例的比例。6 查全率或者召回率（recall）召回率是覆盖面的度量，度量有多个正例被分为正例，recall=TP/(TP+FN)=TP/P=sensitivity，可以看到召回率与灵敏度是一样的。7 其他评价指标，计算速度：分类器训练和预测需要的时间；鲁棒性：处理缺失值和异常值的能力；可扩展性：处理大数据集的能力；可解释性：分类器的预测标准的可理解性。像决策树产生的规则就很容易被理解，而神经网络的一堆参数就不好理解，我们只好把它看成一个黑盒子。8) 精度和召回率反映了分类器分类性能的两个方面，如果综合考虑查准率与查全率，可以得到新的评价指标F1-score，也称为综合F分数。F1=2*precision*recall/(precision+recall)，为综合多个类别的分类情况，评测系统整体性能。经常采用的还有微平均F1（micro-averaging）和宏平均F1（macro-averaging）两个指标：（1）宏平均F1与微平均F1是以两种不同的平均方式求的全局F1指标。（2）宏平均F1的计算方法先对每个类别单独计算F1值，再取这些F1值的算术平均值作为全局指标。（3）微平均F1的计算方法是先累加计算各个类别的a、b、c、d的值，再由这些值求出F1值。（4）由两种平均F1的计算方式不难看出，宏平均F1平等对待每一个类别，所以它的值主要受到稀有类别的影响，而微平均F1平等考虑文档中的每一个文档，所以它的值受到常见类别的影响比较大。**ROC 曲线 AUC 面积和 PRC ROC (Receiver Operating Characteristic) 受试者工作特征曲线**，在不同的分类阈值下，分别计算其对应的“真正型率”TPR 和“假正型率”FPR，然后以TPR为Y轴，以FPR为X轴，画出不同分类阈值下的（FPR，TPR）值对。**TPR = TP/(TP + FN)**，**FPR = FP/(TN + FP)**。ROC 曲线越靠近左上角，说明其对应模型越可靠。若两个学习器的ROC曲线发生交叉，则比较ROC曲线下的面积（Area Under Curve，AUC）来评价模型，AUC越大，模型越可靠。PRC曲线是Precision Recall Curve的简称，描述的是precision和recall之间的关系，以recall为横坐标，precision为纵坐标画出所有阈值情况下的（R，P）值对。该曲线下面面积实际上是目标检测中常用的评价指标平均精度（Average Precision，AP）。AP越高，说明模型性能越好。实践中，引入“平衡点”（BEP）来度量，BEP表示“查准率=查全率”时的数据点，值越大表明分类性能越好。**多分类任务的常用评价指标有哪些？** mAP (Mean Average Precision))：平均AP所有类别的AP的算术平均值。Kappa系数：kappa系数是在统计学中评估一致性的一种方法，取值范围是[-1,1]实际应用中，一般是[-0.1,1]，在ROC曲线中一般不会出现下凸形曲线的原理类似。这个系数的值越高，则代表模型实现的分分类准确度越高。铰链损失（Hinge loss）：一般用来使“边缘最大化”（maximal margin）。损失取值在0-1之间，当取值为0，表示多分类模型分类完全正确，取值为1表明完全不起作用。**多标签分类的评价指标** 1. Hamming loss (汉明损失)，表示所有label中错误样本的比例，所以该值越小则模型的多分类能力 越强。计算公式：HammingLoss(X(i),Y(i)) = (1/|D|)Σ_i=1^n|D|∩X(i)∩Y(i)|/|L|其中：|D|表示样本总数，|L|表示标签总数，xi和yi分别表示预测结果和ground truth，xor表示异或运算。2. Jaccard index(卡卡指数)，概念挺陌生的，公式是1/(|D|Σ_i=1^n|D|∩X(i)∩Y(i)|)/|L|其中：|D|表示样本总数，|L|表示标签总数，xi和yi分别表示预测结果和ground truth，xor表示异或运算。3. 精度、召回率和F1值。其中精度计算公式为T/|D|/|P|，召回率计算公式为T/|D|/|P|，F1值的 计算为精度和召回率的平均和平均数。4. 准确匹配。这个是最严格的标准了，是预测结果和ground truth完全一致时的样本数与总的样本 数值的比例。正确率能很好的评估分类算法吗不同算法有不同特点，在不同数据情况下有不同的表现效果，根据特定的任务选择不同的算法。如何评价分类算法的好坏，要做具体任务具体分析。对于决策树，主要用正确率来评估，但是其他算法，只用正确率能很好的评估吗？答案是否定的。正确率确实是一个很直观很好的评价指标，但是有时候正确率高并不能完全代表一个算法好坏。比如对某个地区进行地震预测，地震发生属性分0：不发生重大地震、1发生重大地震。我们都知道，不发生的概率是极大的，对于分类而言，如果分类器不加思考，对每一个测试样例的类别都划分为0，达到99%的正确率，但是，问题来了，如果真的发生地震时，这个分类器毫无察觉，那带来的后果将是巨大的。很显然，99%正确率的分类器并不是我们想要的。出现这种现象的原因主要是数据分布不均匀，类别为1的数据太少，错分了类别1但达到了很高的正确率就忽视了研究者本身最为关注的情况。**生成模型和判别模型的区别**生成模型：由数据学习联合概率密度分布P(X,Y)，然后求出条件概率分布P(Y|X)作为预测的模型，即判别模型：典型的判别模型包括k近邻，感知机，决策树，支持向量机等。这些模型的特点都是输入属性X可以直接得到后验概率P(Y|X)，输出条件概率最大的作为最终的类别（对于二分类任务来说，实际得到是一个score，当score大于threshold时则为正类，否则为负类）。

集成学习 什么是模型融合集成学习 将一组弱 基预测器的预测结果进行融合集成，以实现一个强预测器，从而获得比单个预测器更好的泛化能力。**模型融合的策略集合？** 1平均法：一般用于回归预测模型中。平均法包括一般的平均和加权平均融合。投票回归器（Voting Regressor Boosting 系列融合模型 2投票法：一般用于对绝对多数投票（得票超过一半），相对多数投票（得票最多），加权投票。例如bagging 模型投票分类器（Voting Classifier）3学习法：通过另一个预测器称为混合器或元学习器来实现融合。常见的有Stacking和Blending两种。stacking一般使用交叉验证的方式，Blending是建立一个Holdout集。**模型融合的分类**：按照弱预测器之间的依赖关系，模型融合方法可分为两类。1弱预测器间不存在强依赖关系、可同时学习的并行化方法，如果弱预测器只有一种，即所有的弱预测器都是基于同一类模型，代表是Bagging和“随机森林”。如果弱预测器的种类多，比如投票分类器，投票回归器，Stacking。2弱预测器间存在强依赖关系、必须串行学习的序列化方法，代表是Boosting方法有下面几类模型AdaBoost, Gradient Boosting: GBRt, XGBoost, LightGBM, 注意区分：强依赖≈强相关性，弱预测器之间的存在强依赖，是指单个弱预测器的训练存在必需的先后关系；弱预测器之间若存在强相关性，会影响融合后的性能。**模型融合集成学习一定有效吗？**可以通过数学证明，随着弱预测器数目T的增加，集成后得到的强预测器的错误率将呈指数级下降最终趋近于零。弱预测器数目T过大，可能导致强预测器出现过拟合。集成学习有效，必须：弱预测器之间的性能表现差不多，不能差距太大，同时弱预测器之间的相关性要尽可能的少。**误差、偏差和方差有什么区别和联系**对于Error：误差（error）：一般地，我们把学习器的实际预测输出与样本的真是输出之间的差异称为“误差”。Error = Bias + Variance + Noise。Error反映的是整个模型的准确度。Bias：描述了在当前任务上任何学习算法所能达到的期望泛化误差的下界，即刻画了学习问题本身的难度。对于Bias：Bias 衡量模型拟合训练数据的能力（训练数据不一定是真正的 training dataset，而是只用于训练它的某一部分数据，例如：mini-batch），Bias反映的是模型在样本上的输出与真实值之间的误差，即模型本身精确度。Bias 越小，拟合能力越高（可能产生overfitting）；反之，拟合能力越低（可能产生underfitting）。偏差越大，越偏离真实数据，如图第一行所示。对于Variance：方差公式：S_y² = 1/N Σ (x_i - x̄)²。Variance描述的是预测值的变化范围，离散程度，也就是离其期望值的距离。方差越大，数据的分布越分散，模型的拟合程度越低。Variance反映的是模型每一次输出结果与模型输出期望之间的误差，即模型的稳定性。Variance越小，模型的泛化的能力越高；反之，模型的泛化的能力越低。如果模型在训练集上拟合效果比较优秀，但是在测试集上拟合效果比较差劣，则方差较大，说明模型的稳定性较差，出现这种现象可能是由于模型对训练集过拟合造成的。**如何判断欠拟合和过拟合？** 1利用训练误差和验证误差，训练误差和验证误差都很大，并且很接近，此时模型存在欠拟合问题；当训练误差很小，但验证误差还是很大时，此时模型存在过拟合问题。2利用训练集和验证集上的学习曲线。欠拟合：高 bias 误差，低 variance 误差过拟合：低 bias 误差，高 variance 误差。好：低 bias 误差，低 variance 误差。学习曲线就是不同训练集大小时训练集和交叉验证的准确率。当训练集和测试集的误差收敛但依旧很高时，为高偏差。当训练集和测试集的误差收敛但依旧很高时，为高偏差。当训练集的准确率比其其他独立数据集上的测试结果的准确率要高时，一般都是过拟合。**如何解决欠拟合** 1添加其他特征项。组合、泛化、相关性、上下文特征、平台特征等特征是特征添加的重要手段，有时候特征项不够会考虑模型欠拟合。2添加多项式特征。例如将线性模型添加二次项或三次项增加模型的泛化能力更强。例如，FM (Factorization Machine) 模型，FFM (Field-aware Factorization Machine) 模型，其实都是线性模型，增加了二项或多项式，保证了模型一定的拟合程度。3可以添加模型的复杂程度。4减小正则化系数。正则化的目的是用来防止过拟合的，但是现在模型出现了欠拟合，则需要减少正则化参数。例如解决过拟合：1重新清洗数据，数据不会导致过拟合，此类情况需要重新清洗数据。2增加训练样本数量。3降低模型复杂程度。4增大 dropout 方法，dropout 方法，通俗的讲就是在训练的时候让神经元以一定的概率不工作。6. early stopping。7. 减少迭代次数。8. 增大学习率。9. 添加噪声数据。10. 树结构中，可以对树进行剪枝。11. 减少特征项。**如何加快模型的训练** 1特征缩放/标准化 (1) Feature Scaling - 特征缩放/归一化。输入值减去样本本均值，然后除以样本标准差。式子如下：x=(x-min(x))/(max(x)-min(x)) (2) Mean Normalization - 均值归一化。输入值减去样本本均值，然后除以样本标准差。式子如下：x=(x-avg(x))/(max(x)-min(x)) (3) Standardization 标准化：让输入的值减去样本平均数μ，再除以样本标准差σ。经过这样的处理，数据符合标准正态分布，即均值为0，标准差为1。x=(x-μ)/σ。2**梯度下降** (1) Gradient Descent - 梯度下降 如果需要找到一个函数的局部极小值，必须朝着当前点上当前点所对应梯度（或者是近似梯度）的反方向，向前规定步长的距离进行迭代搜索。θ^{t+1}=θ^t-η∇L(θ^t)。为什么要以梯度的反方向为更新方向？因为梯度方向是函数方向导数最大的方向，所以沿着梯度方向反方向更新的话，函数下降变化率最大。优点：能保证证明梯度下降的方向去优化，易于并行实现。缺点：样本数很多时，会很慢。(2) Stochastic Gradient Descent - 随机梯度下降每次更新只用到了的一个样本，如果这个训练集有m个样本，那么梯度下降更新一次参数，则离局部最优值已经更新了m次参数了。优点：快。BGD更新一个次的时候，SGD会更新的次数=Batch size。缺点：引入了随机噪声，可能会出现震荡的现象。当前目标函数非凸时，反而可以使其逃离局部最优值。无法充分利用计算机的并行能力。(3) Mini-batch Gradient Descent Mini-batch 梯度下降是梯度下降和随机梯度下降的中和版本，Mini-batch 梯度下降每次更新所考虑的样本是可以被指定的，如果总共有m个样本，那就可以在1/m中任意指定。如果每次更新时所参的样本数合适，那么既兼顾了随机梯度下降更新速度的特性，又兼顾了梯度下降更新后的稳定性。优点：收敛快，计算开销小。3. **调整学习率** 当我们在训练过程中，发现 loss 下降的很慢时，可以适当增大学习率；发现 loss 不降反增的时候，要降低学习率。(1) Adagrad 算法的学习率会根据迭代次数来调整学习率，从而让学习率越来越小的方向。通过公式可以看到，学习率会一直除以前面所有梯度的平方和并开根号的话，这一定是一个大于0的数，所以学习率会越来越小，但是防止一开始的时候梯度是0，从而导致学习率为0会导致错误的，所以后面还要跟一个很小的正数ε，最终的式子是这样的：w^{t+1}=w^t-η/(Σ_{i=1}^tg_i²+ε)₀ Adagrad 算法也有很多不足：a) 如果初始的学习率设置过大的话，这个学习率要除以一个较大梯度，那么此算法会对梯度的调节太大；b) 在训练的中后期，分母上梯度平方的累加将会越来越大，使 gradient>0，使得训练提前结束。(2) RMSProp Adagrad 算法的改进版RMSProp 算法：w^{t+1}=w^t-ηg^t/σ^t σ^t=√(α(σ^{t-1})²+(1-α)(g^t)²)Adagrad 和 RMSprop 算法这两个算法很相近，不同之处在于 RMSprop 算法增加了一个衰减系数α来控制历史信息的获取多少。(3) SGD with Momentum (SGD-M) SGD 在遇到沟

