

这是一份根据你提供的PDF课件整理的详细机器学习复习笔记。我将内容按逻辑模块进行了分类，并提炼了核心要点、公式和复习重点，方便你记忆和理解。

第一部分：机器学习的基础概念 (P1 - P9)

1. 机器学习的定义 (核心考点)

- **Samuel (1959)**: 在不进行显式编程的情况下，赋予计算机学习的能力。
- **Mitchell (1998) - 形式化定义**:
 - 一个程序被认为从**经验 E (Experience)** 中学习，解决**任务 T (Task)**，并用**性能度量 P (Performance)** 来衡量。
 - 如果在任务 T 上的性能 P 随着经验 E 的增加而提高，则称该程序具有学习能力。
 - **例子 (邮件过滤)**：
 - **T**: 分类邮件是否为垃圾邮件。
 - **E**: 用户标记邮件的历史数据。
 - **P**: 正确分类邮件的比例 (准确率)。

2. 大规模机器学习 (Large-scale ML)

- **Banko & Brill (2001) 的研究**: 在复杂的语言处理任务中，随着数据量 (Millions of Words) 的增加，各种算法 (朴素贝叶斯、感知机等) 的表现趋于一致且都在提升。
- **核心金句**: "It's not who has the best algorithm that wins. It's who has the most data." (决定胜负的往往不是谁的算法最好，而是谁的数据最多)。
- **深度学习 vs 传统机器学习**:
 - 传统机器学习在数据量达到一定程度后性能趋于平缓。
 - 深度学习 (Deep Learning) 能更好地利用海量数据，性能随数据量增加持续提升。

第二部分：机器学习算法分类 (P10 - P19)

1. 监督学习 (Supervised Learning)

- **特征**: 训练数据包含“正确答案” (即标签 Label)。
- **两大任务**:
 - **回归 (Regression)**: 预测**连续值** (如：房价预测、产品销量预测)。
 - **分类 (Classification)**: 预测**离散类别** (如：肿瘤良性/恶性、账号是否被盗、垃圾邮

件检测)。

2. 无监督学习 (Unsupervised Learning)

- **特征:** 训练数据不含标签，由算法自动发现数据中的结构。
- **核心任务:**
 - **聚类 (Clustering):** 自动将相似的数据归为一类 (如：市场细分、新闻自动分组)。
 - **概率密度估计:** 估计数据的分布。
 - **降维/可视化 (Dimensionality Reduction):** 压缩数据并保留关键特征。

3. 其他分类 (仅提及)

- 强化学习 (Reinforcement Learning)
- 推荐系统 (Recommender System)

第三部分：机器学习核心问题——曲线拟合与正则化 (P20 - P38)

1. 曲线拟合 (Curve Fitting)

- **目标:** 找到一个多项式函数 $y(x, \mathbf{w}) = \sum_{j=0}^M w_j x^j$ 来拟合给定的数据。
- **误差函数 (Error Function):** 常用的平方和误差 (Sum of Squares Error):
$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N [y(x_n, \mathbf{w}) - t_n]^2$$
 - 目标是最小化 $E(\mathbf{w})$ ，即让预测值尽可能接近真实目标值 t_n 。

2. 模型阶数 M 与过拟合 (Overfitting)

- $M = 0, 1$ (**低阶**): 欠拟合 (Underfitting)，模型太简单，无法捕捉数据规律。
- $M = 3$: 拟合良好，能很好地泛化到新数据。
- $M = 9$ (**高阶**): **过拟合 (Overfitting)**。
 - **表现:** 训练集误差几乎为0 (曲线穿过了每个点)，但测试集误差剧增 (曲线剧烈震荡)。
 - **本质:** 模型参数 w 的数值会变得非常巨大 (见P33表格)。
- **缓解过拟合的方法:**
 1. **增加数据量 (N):** 即使是高阶模型，如果有海量数据 (如 $N = 100$)，也能减轻震荡。

2. 正则化 (Regularization).

3. 正则化 (Regularization) - 核心技术

- **目的:** 在误差函数中加入对参数 w 大小的惩罚，防止参数过大导致震荡。

- **新的目标函数:**

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

- λ : 正则化系数，控制惩罚力度。

- **正则化效果:**

- $\ln \lambda = -18$ (适中): 完美平衡，消除过拟合。
- $\ln \lambda = 0$ (太大): 惩罚过度，模型趋向于平滑导致欠拟合。

第四部分：总结与复习重点 (P39 - P41)

1. 三大核心概念:

- **机器学习定义:** T (任务), E (经验), P (性能)。
- **监督学习:** 必须有标签，分为回归和分类。
- **无监督学习:** 无标签，典型代表是聚类。

2. 模型选择与优化:

- **过拟合:** 模型阶数过高、数据量太少导致。
- **正则化:** 通过增加惩罚项 λ 来限制模型复杂度，是解决过拟合的关键手段。

3. 重要参考文献:

- Samuel (1959) - checkers (西洋跳棋)研究。
- Mitchell (1998) - 经典教材。
- Banko & Brill (2001) - 数据量的重要性研究。

复习建议:

- 熟记 Mitchell 关于 ML 的定义，考试常考 T/E/P 的辨析。
- 理解并能手写平方和误差公式及带正则化项的误差公式。
- 记住过拟合的直观表现 (训练误差小，测试误差大，参数值极大) 及其两种解决方法。