

# 主题：时序大模型读书会第一期笔记

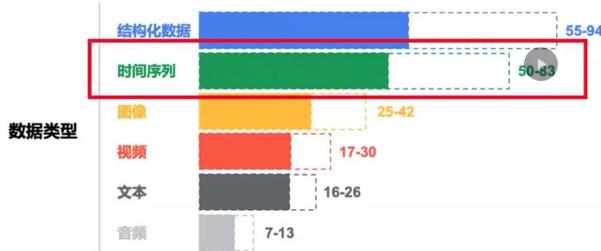
主要内容：对读书会专题一时序大模型的整体讲解，包括脉络与框架。分享从三个部分展开，分别是研究背景、基于自然语言大模型的时序分析和时序数据大模型/基座模型。

## 一. 时间序列数据概况

**定义**

- 1. 时间序列指将同一统计指标的数值按其发生的**时间先后顺序**排列而成的「数列」
- 2. 时序数据的**顺序**和**时间依赖性**特征使得它在预测和趋势分析中具有独特的价值。

不同类型的潜在总价值中的百分比 %

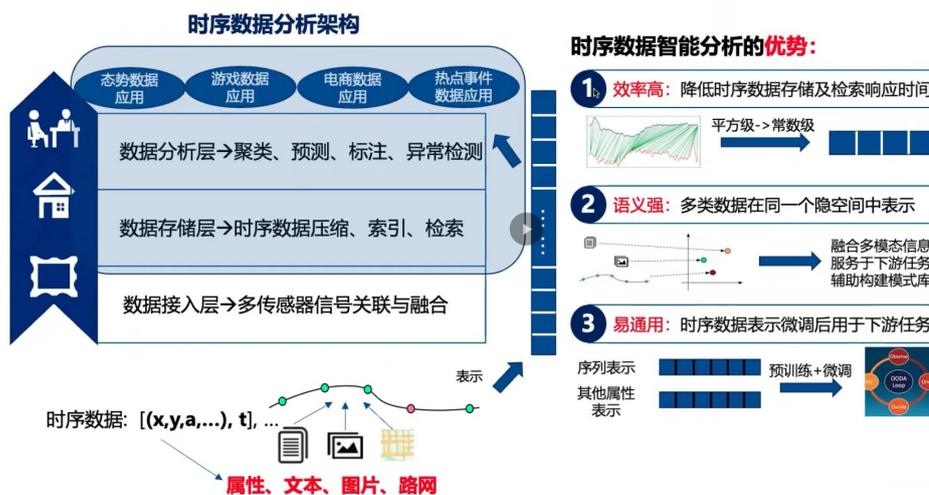


Source: McKinsey Global Institute



首先，姚迪老师给出了时间序列的定义与两大特性，即顺序性和时间依赖性，并指出时间序列数据的研究价值。接着介绍了目前主流的时序分析任务可以分为时序预测、异常检测、分类和插值，在不同应用场景有不同的定义和解决方案。

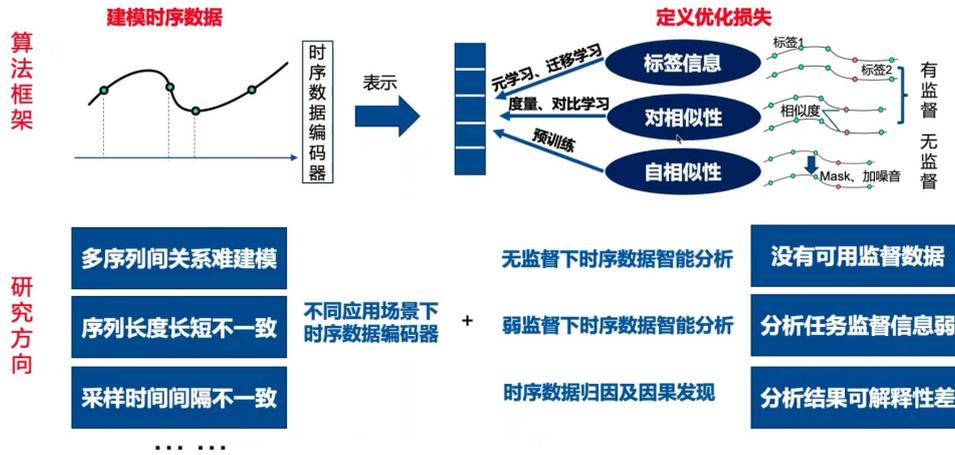
## 时序数据智能分析



时序数据分析框架：姚老师从数据接入层、存储层和分析层三个方面对时序数据分析的任务和观测重点进行总结，研究的总体策略是基于深度学习或智能算法来解决时序分析中的

应用，并阐明该框架分析的优势，主要包括效率高、语义强和易通用等特点。

## 时序数据智能分析



姚老师分享了自己研究过程的整体范式，可以分为两阶段：首先对时序数据编码，将时序数据编码成固定长度的向量，然后在基于不同的下游任务再定义向量的优化方式，如可以利用标签信息做有监督学习、利用对比学习的方式、用自监督的方式做预训练等，同时也分享了自己在这两个阶段中的工作。

## 二. 时间序列大模型概况

姚迪老师介绍自己之前的工作主要是小模型在专用任务上的处理，但当前随着大模型的不兴起，在各个领域都展现出了突破性进展，因此姚迪老师对于开源大模型和闭源大模型的背景和技术原理进行了介绍，并指出开源大模型对未来研究的重要意义，如可以通过对模型内部结构和模型中的参数做微调或定制化的操作使其能更好地适用于时序数据分析任务。

## 将大模型和时间序列数据结合已经成为当下的学术热点



### 大模型和时间序列数据结合面临的难点与挑战：

- **跨模态:** 时序数据模态与自然语言不同，难以统一处理
- **Tokenizing/ Packing:** 数值型数据难以统一分割、离散化
- **数据规模:** 公开可获取时序数据规模有限

随着时间序列与大模型的进展，将两者结合进行研究已经成为当下的学术热点。对时间序列分析的研究归纳为三个阶段。20 世纪 60 年代是传统时间序列分析阶段，主要通过统计相关的方法，更倾向处理平稳序列或已知非平稳模式的序列。在 2010 年前后，深度学习技术被提出，使得面向时间序列分析的任务被越来越多深度学习的方法所主导，如 RNN, LSTM, Transformer 等序列建模的框架来建模时间序列数据，但是此时主要还是聚焦于单一任务，即针对一个特定任务来设计一个特定模型解决，同时训练数据也是特定的，因此泛化能力一般。从 2023 年开始，以 Time-GPT-1 和 Lag-Llama 等面向时序数据的基座模型和大模型为主的研究工作出现，利用了大模型处理不同类型时间序列数据的能力，因此能够适应数据的变化和异常情况，从而提高预测的准确性和稳定性，但是将大模型和时间序列数据结合也面临一些难点与挑战：如时序数据与自然语言模态不同，难以统一处理；目前大语言模型基本都是基于单词的 embedding 作为模型的基石，对于自然语言来说是闭集，但由于时序是连续性数值数据，如何放在语言模型框架下进行分析处理也是一个挑战；时间序列数据规模相比自然语言数据来说有限，因此在训练过程中可能不足以支撑参数量特别大的模型的训练。

### 三. 当前时序大模型的主要工作

该部分概括介绍了当前时序大模型的主要工作，可以分为两类，并对两类工作进行概括，交代未来两期时序数据大模型的整体结构。

	基于自然语言大模型的时序分析 (第二期: LLM赋能时序数据分析)	时序数据大模型/基座模型 (第三期: 时序数据基座模型)
主要思路	<ol style="list-style-type: none"> <li>1. 将时间序列数据转换为适合自然语言处理模型处理的格式</li> <li>2. 这种方法利用了NLP模型的强大处理能力和灵活性, 以及它们在处理序列数据方面的潜在优势</li> </ol>	<ol style="list-style-type: none"> <li>1. 定制化模型训练方法, 提高模型对时序数据的理解预测能力</li> <li>2. 通过这种专用化的训练, 模型能更好地理解时间序列数据的独特性质和模式</li> </ol>

### 1. 基于自然语言大模型的时序分析。

这类工作主要利用已有的自然语言大模型的能力赋能时序数据分析。主要思路是将时间序列数据转换为适合自然语言处理模型处理的格式。这种方法充分利用了 NLP 模型的强大处理能力和灵活性, 以及它们在处理序列数据方面的潜在优势。这一类的工作主要将自然语言大模型当成工具使用, 大部分参数固定。

在未来的读书会中, 将会回答以下问题: 在数据、模型和使用层面, 如何利用 LLM 已有的能力去赋能时间序列分析? 如何处理自然语言数据与时间序列数据的模态对齐挑战? 同时也会请到该方向下相关工作的讲者来分享若干典型工作的细节, 包括 Time-LLM、PromptCast、GPT4TS、TEST、LLM4TS 等。

### 2. 时序数据大模型/基座模型。

#### 传统深度学习方法在处理时间序列数据时面临一系列挑战



这类工作解决的问题为如何从头开始去训练一个时间序列数据专用的大模型, 即 train

ing from scratch。它的主要思路为从设计数据、模型结构、定义优化损失等方面训练面向时间序列的基座模型，该模型需要具有对于不同类型的时序数据的泛化能力，使其能够专业地理解时间序列数据的特性和模式，并能处理多种时序分析任务，但这些模型就不一定需要具备对自然语言数据的分析能力。这样做的优势有，模型可以学习到时间序列相关的通用模式和特征，因此在时序数据的理解上会更有优势，典型的工作有 Time-GPT-1、Lag-Llama、Moment 等，也将会在后续的读书会进行详细介绍。

