

预测与 MoE Forecasting and MoE

李丰

北京大学光华管理学院

https://feng.li/forecasting-with-ai

Mixture of Experts



为什么需要 Mixture of Experts (MoE) ?

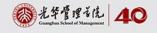
• 大模型越来越大, 但不是所有任务都需要一个"巨无霸"

- 企业真实问题: 需求预测、用户画像、营销投放、风险控制——场景差异巨大
- 一个模型难以擅长所有任务("No Free Lunch")

- MoE 的核心思想:
 - "把复杂问题分给更擅长的小专家,高效又精准。"

时间序列预测的 No Free Lunch 理论

- No Free Lunch 理论 (NFL)
 - Wolpert, D. H., & Macready, W. G. (1997). <u>No free lunch theorems for optimization</u>. IEEE Transactions on Evolutionary Computation, 1(1), 67–82.
 - 没有任何预测模型能在所有场景下表现最好。
 - 每一种模型都只对某些类型的数据/业务特别擅长。
- 为什么预测中会出现 No Free Lunch?
 - 时间序列预测的世界很复杂:
 - 有些产品 → 稳定有趋势(如:咖啡豆)
 - 有些产品 → 明显季节性(如:空调)
 - 有些产品 → 大促跳点(如: 电商快消)
 - 有些 → 完全随机 (如: 短信验证码)
 - 有些 → 长尾断货 (intermittent demand)
 - 不同的模式 → 不同的最优模型
 - ARIMA 擅长趋势
 - ETS 擅长季节
 - Prophet 擅长节假日
 - XGBoost/LSTM 擅长复杂非线性
 - 大模型 (LLM/TimeGPT) 擅长跨序列迁移
 - 没有哪个模型可以全都擅长。



No Free Lunch 在预测中的形式化

- 如果你把所有可能的数据集平均起来:
 - 所有预测模型的平均表现都是一样的。
 - 不存在一个永远最优的"万能预测模型"。
- 在真实企业预测里:不能指望一个模型吃遍所有 SKU、区域、品类、促销场景。
- 很多 人 认为: "用 TimeGPT / DeepSeek / GPT-5 不就能解决了吗?"
 - 答案是: ★。大模型 (LLM/FMs) 也无法逃避 No Free Lunch
- 原因:
 - 虽然大模型非常强,但仍然不可能在所有类别的时间序列上同时最佳。
 - 趋势型序列、节季型序列、噪声型序列、促销型序列 → 需要不同的专长。
- MoE(混合专家)就是大模型用来应对 NFL 的策略之一。



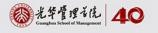
MoE(Mixture of Experts)解决方案

- 不同时间序列 → 激活不同专家
 - 趋势专家
 - 季节专家
 - 促销专家
 - 外生变量专家
 - 异常专家
- Gating 让模型自动判断: "用谁做预测?"
- 这等于建立一个自动化的预测团队,而不是靠一个"巨无霸模型",由多个小模型(仍然可以是LLM, 但是规模和参数很小)组成。



MoE 从经典模型到现代深度学习

- Jacobs, R. A., Jordan, M. I., Nowlan, S. J., & Hinton, G. E. (1991). <u>Adaptive mixtures of local experts</u>. Neural Computation, 3(1), 79–87.
 - MoE 概念首次提出
 - 使用 Gating network 选择专家
 - 基础思想:不同专家处理不同输入区域
- Shazeer, N., Mirhoseini, Azalia, Maziarz, Krzysztof, Davis, A., Le, Q., Hinton, G., & Dean, J. (2017).
 Outrageously Large Neural Networks: The Sparsely-Gated Mixture-of-Experts Layer. International Conference on Learning Representations.
 - Google 深度 MoE 的奠基之作
 - 引入 Top-K gating (稀疏激活)
 - 每次只激活部分专家 → 巨大规模 + 低成本
 - 后来 GPT-4、DeepSeek、LLaMA 都受它启发
- Fedus, W., Zoph, B., & Shazeer, N. (2022). <u>Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity</u>. The Journal of Machine Learning Research, 23(1), 120:5232-120:5270.
 - Google 最著名的 MoE 模型
 - 每层只选 1 个专家(Switch gating)
 - 训练速度极快,参数上千亿,成本可控



通用大模型(LLM)中的 MoE

• GPT-4 (推测)

- 公认为采用 Sparse MoE
- 多层专家结构
- 参数量大但激活参数少 → 性能高、成本 低

DeepSeek MoE (DeepSeek-V2/V3/R1)

- Hybrid MoE (Dense + MoE)
- 有共享专家 + 专家组
- Token-level routing
- 更稳定、更具工业化特点
- 特别适合预测任务(heterogeneous time series)
- 引发了 DeepSeek 做空英伟达

LLaMA-3.1 70B MoE (Meta)

- 最新开源 MoE LLM
- Top-k MoE,每次激活两个专家
- 非常高效

Mixtral (Mistral Al, 2023)

- 模型: Mixtral 8×7B
- 8 个专家, 每次激活 2 个
- 开源性能极强(可媲美 GPT-3.5)
- MoE 实用化的重要标志

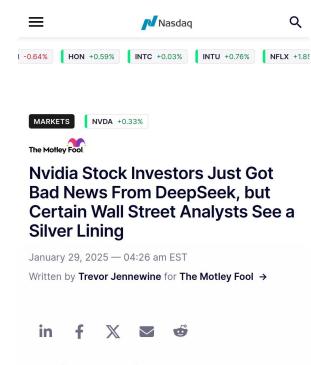
小故事: DeepSeek 事件做空英伟达

事件概况

- 2025年1月27日,DeepSeek发布其新模型,引起市场强烈反应,英伟达股价单日下跌约 17%,市值蒸发数千亿美元。
- 短线做空英伟达的市场参与者据称从此次下跌中获利约 66 亿美元。
- 媒体将其称作"英伟达被 DeepSeek 恐慌性冲击"的一次典型事件。
- 市场情绪被触发。在 AI 热潮中,英伟达估值已经非常高。任何对其增长路径或基础设施模式的怀疑,都可能成为"做空"或风险重估的触发器。

做空英伟达的大致机制

- 虽然具体做空交易人的全部细节不可得,但从公开报道可以推断如下流程:
- 做空者观察到 DeepSeek 发布其模型,同时市场对 AI 基础设施依赖的假设受到挑战。
- 建立英伟达空仓或买入看跌期权,或通过衍生工具押注英伟达股价下跌。
- 随着英伟达股价因为 DeepSeek 的冲击而快速下跌, 空头获利。
- 这一过程也伴随市场做空成本增加、监管介入(比如 "uptick rule" 等)的问题。



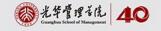
Nvidia (NASDAQ: NVDA) made stock market history on Monday, Jan. 27, but not the good kind. The chipmaker saw its share price decline 17%, due to concerns about an artificial intelligence (AI) model from Chinese start-up DeepSeek. That nosedive erased \$589 billion of its market value, the largest single-day loss for any company on record.

Mixture of Experts (MoE) 自动化预测团队

现实团队	MoE 中的对应
每个分析师(专家)对某类数据很擅长	Expert (趋势、季节、促销、异常·····)
团队负责人分配任务	Gating (路由器)
每个任务可能由不同人完成	Token-level / sample-level routing
最终由负责人综合意见	Expert 加权融合

MoE 的本质就是:

- 一个能自动分工协作的预测团队结构。
- 而不是一个"什么都试图做、但什么都不精"的巨无霸单模型。



预测(Forecasting)领域中的 MoE

- TimeGPT (Nixtla)
 - 使用 Mixture-of-Experts + Transformer 架构
 - 负责分配不同时间序列结构的专家
 - 特别擅长跨序列泛化
- Chronos-Bolt (Amazon, 2024)
 - Amazon 的大规模时间序列模型
 - MoE 结构 + Token routing
 - 为百万级 SKU 提供预测

- Google Moirai(2024, 世界最大 TS 模型)
 - 100B+ 参数 + Sparse MoE
 - 面向时序预测
 - Gating 自动为不同 regions / patterns 分派专家
 - 强调 scaling laws 与推理效率



MoE 的数学表示

MoE in Math



输入表示(时间序列 Embedding)

对时间序列 (y_t) 进行编码:

$$x_t = f_{\text{embed}}(y_{t-L:t}, \mathbf{s}_t, \mathbf{z}_t)$$

其中:

• y_{t-L:t}: 过去(L)个历史点

• \mathbf{s}_t : 时间特征 (weekday, holiday, season)

• \mathbf{z}_t : 外生变量(价格、天气、宏观等)

• x_t : Embedding



路由器(Gating Network)

- MoE 的核心是**给每个时间点选择不同专家**。
- 对每一个时间点(t), 路由器给出对(K)个专家的权重分布:

$$\mathbf{g}_{t} = \text{softmax}(W_{g}\mathbf{x}_{t} + b_{g}), g_{t}^{(k)} = \frac{\exp(\alpha_{t}^{(k)})}{\sum_{j=1}^{K} \exp(\alpha_{t}^{(j)})}$$

其中:

- $g_t^{(k)}$: 第 (k) 个专家在时间点 (t) 的权重
- W_g , b_g : Gating 参数
- *K* < ∞: 专家数量

Sparse MoE(DeepSeek、Moirai、Mixtral 等)采用 Top-(m) 激活:

$$g_t^{(k)} = 0$$
 if $k \notin \text{Top-}m(\mathbf{g}_t)$

专家 (Expert) 模块

每个专家 (E_k) 是一个预测子模型,例如:

- 小 Transformer block
- /J\ MLP (MLP expert)
- 小 LSTM / Dilated CNN expert
- Domain expert (趋势 / 季节 / 促销 / 异常)

其形式为:

$$\mathbf{h}_t^{(k)} = E_k(\mathbf{x}_t)$$

其中: $\mathbf{h}_t^{(k)}$ 是第 (k) 个专家的隐藏输出。可以看成:

$$\mathbf{h}_t^{(k)} = E_k(\mathbf{x}_t) = \sigma(W_k \mathbf{x}_t + b_k)$$

或 Transformer block:

$$\mathbf{h}^{(k)} = FFN_k(MHSA_k(\mathbf{x}))$$

其中MHSA_k为多头自注意力 (Multi-Head Self-Attention),FFN_k为前馈网络(Feed-Forward Network)

MoE 输出融合(Mixture Aggregation)

MoE 的融合是加权专家输出:

$$\mathbf{h}_t = \sum_{k=1}^K g_t^{(k)}, \mathbf{h}_t^{(k)}$$

如果是 Sparse MoE:

$$\mathbf{h}_t = \sum_{k \in \text{Top-}m(t)}^K g_t^{(k)}, \mathbf{h}_t^{(k)}$$

最终预测头(Forecasting Head)

输出未来(H)期预测:

$$\hat{y}_{t+1:t+H} = f_{\text{forecast}}(\mathbf{h}_t)$$

或逐步预测:

$$\hat{y}_{t+h} = W_o^{(h)} \mathbf{h}_t + b_o^{(h)}, \quad h = 1, ..., H$$

把所有步骤合并:

•
$$\left| \hat{y}_{t+h} = f_{\text{forecast}} \left(\sum_{k=1}^{K} g_t^{(k)} E_k \left(f_{\text{embed}}(y_{t-L:t}, \mathbf{s}_t, \mathbf{z}_t) \right) \right) \right|$$

DeepSeek、TimeGPT、Moirai 的 MoE 特点

DeepSeekMoE (Hybrid MoE)

在 MoE 之前加入 dense FFN (稳定性):

$$\mathbf{u}_{t} = \text{FFN}_{dense}(\mathbf{x}_{t}),$$

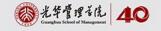
$$\mathbf{h}_{t} = \mathbf{u}_{t} + \sum_{k \in \text{Top-}m(t)}^{K} g_{t}^{(k)} E_{k}(\mathbf{x}_{t})$$

TimeGPT / Moirai (Forecasting FM)

• 使用跨序列共享专家: $E_k(\mathbf{x}_t; \theta_k)$, 其中 θ_k 在不同 SKU / 城市 / 国家共享。

什么是大模型的"蒸馏"(Distillation)?

- 蒸馏就是让一个大模型(老师)教出一个更小、更快、更便宜的模型(学生), 而且保留主要能力。
- 就像把一锅汤熬浓,把精华提取出来。
- 大模型(teacher)包含大量知识。
 - 在蒸馏过程中把不重要的部分"挥发掉", 把关键信息"提纯"
- 得到一个更轻、更高效的小模型(student)
- 它不是简单裁剪,而是"知识迁移+模型压缩"。



蒸馏的三个核心步骤

老师模型生成大量高质量数据或答案

学生模型学习老师的"思考方式"、判断边界、输出概率分布。

- 学生模型模仿老师的行为
- 优化目标从模仿正确答案变成模仿老师的输出分布(更细腻)。
- 学生模型学习这种"软判断"。得到一个更小、更便宜能跑在边缘设备的模型
- 如: Teacher = 70B, Student = 7B
- 但保留了 Teacher 70%~90% 的能力。

蒸馏的主要类型

- Logit Distillation(经典蒸馏)
 - 学习 Teacher 的输出概率。
- Response Distillation (答案蒸馏)
 - 学生学习 Teacher 的回答、推理链。
- Self-Distillation(大模型自我蒸馏)
 - 如 DeepSeek-R1、OpenAl O1
 - 老师是自己生成的"更优解"或"延长推理"
 - 学生学习更短、更高效的推理方式。
- Preference Distillation (偏好蒸馏)
 - 把 RLHF 的偏好传给学生,提升对齐程度。



蒸馏与 MoE、预测有什么关系?

- 让 MoE 的专家更轻量
 - 大模型 MoE 专家往往 1-7B 规模,蒸馏让它们轻但保持能力。
 - 让推理(预测)更快
- 时间序列预测需要速度:
 - 存货盘点
 - 电商分钟级更新
 - 金融高频风险识别
 - 蒸馏后的模型能实现「10 倍到 100 倍加速」。
- 大模型界最重要的蒸馏应用
 - GPT-4 → GPT-4-mini / GPT-3.5
 - Mixtral MoE → Distilled Mixtral
 - BERT → DistilBERT (最早的成功案例)
 - LLM → LLM for Forecasting (TimeGPT 蒸馏小模型)