# Model Aggregation:
# Data-driven combination of black box models
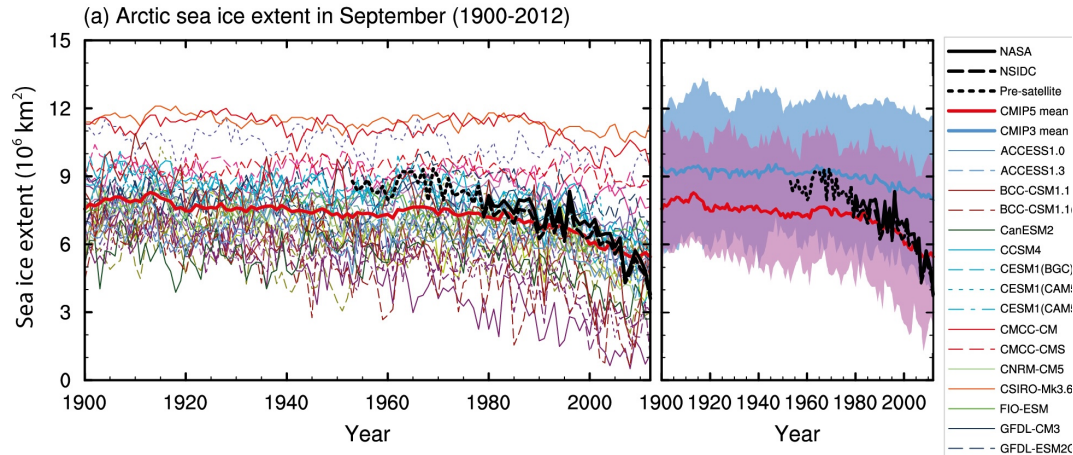
Theo Bourdais

PhD Student, Computing + Mathematical Sciences

June 15th 2025

Caltech

# Real-life example from the IPCC



(a) Arctic sea ice extent in September (1900-2012)
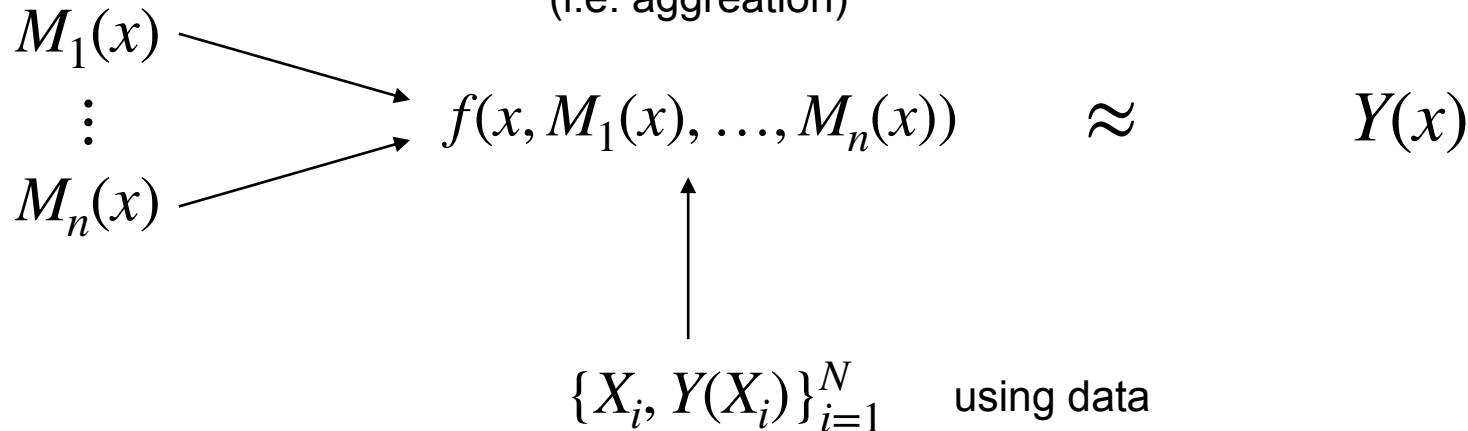
Arctic sea ice extent estimated by many models,

Coupled Model Intercomparison Project (report AR5 - figure 9.24)

Caltech

# The aggregation problem

Given the models,

create a combination
(i.e. aggreation)

that approximates the target

$$M_1(x)$$

$$\vdots$$

$$M_n(x)$$

$$f(x, M_1(x), \ldots, M_n(x)) \qquad \approx \qquad Y(x)$$

$$\{X_i, Y(X_i)\}_{i=1}^{N} \qquad \text{using data}$$

3

**Caltech**

# Best Mean Squared Error Aggregation

The best possible aggregation in Mean Squared Error is

$$M_A^*(x) := \underset{f \text{ measurable}}{\operatorname{argmin}} \mathbb{E}[|Y(x) - f(x, M_1(x), \ldots, M_n(x))|^2] = \mathbb{E}[Y(x)|M_1(x), \ldots, M_n(x)]$$

This is intractable in general

**Special Case:** $(Y(x), M_1(x), \ldots, M_n(x))$ **is Gaussian**

$$M_A^*(x) = \sum_{i=1}^n \alpha_i^*(x) M_i(x)$$

$$\alpha^*(x) = \underset{a \in \mathbb{R}^n}{\operatorname{argmin}} \mathbb{E}\left[\left|\left|Y(x) - \sum_{i=1}^n a_i M_i(x)\right|\right|^2\right] = \mathbb{E}\left[M(x)M(x)^T\right]^{-1} \mathbb{E}\left[M(x)Y(x)\right]$$

**Caltech**

# Best case aggregation: Gaussian models

To solve the Laplace equation:

$$\begin{cases} \Delta Y = f & on \ \Omega \\ Y = g & on \ \partial\Omega \end{cases}$$

We can use a Gaussian process with:

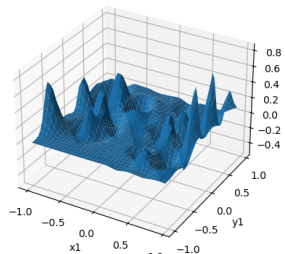- A kernel $k$

- A set of collocation points $X \subset \Omega$

To get a Gaussian approximation of the solution [Chen et al., 2021]

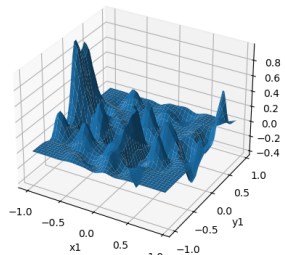$$\xi \sim \mathcal{N}(0,k) \quad \hat{Y} = \mathbb{E}[\xi \,|\, \Delta\xi(X) = f(X)]$$

Caltech

Predictions

$M_1(x)$

$M_2(x)$

$\vdots$

$M_n(x)$

Laplace equation $\begin{cases} \Delta Y = f & on \ \Omega \\ Y = g & on \ \partial\Omega \end{cases}$

Average

$\not\approx$

True solution

$Y(x)$

Caltech

Predictions

$M_1(x)$

$M_2(x)$

$\vdots$

$M_n(x)$

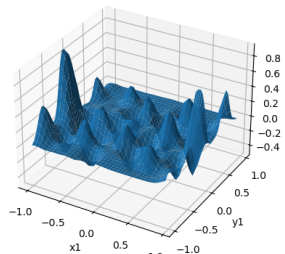Laplace equation $\begin{cases} \Delta Y = f & on \ \Omega \\ Y = g & on \ \partial\Omega \end{cases}$
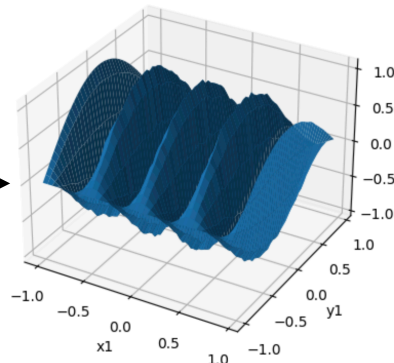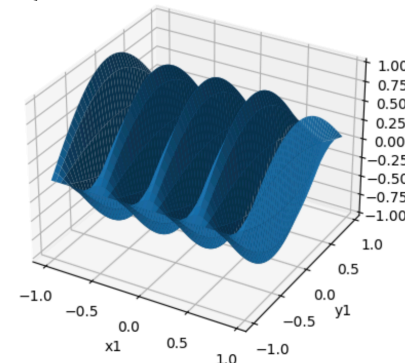
$$\sum_{i=1}^{n} \alpha^*(x) M_i(x)$$

$\approx$

True solution

$Y(x)$

Caltech

# Minimal Error Aggregation

This does not work!

$\alpha^*$ is defined as:

$$\alpha^*(x) = \operatorname*{argmin}_{a \in \mathbb{R}^n} \mathbb{E}\left[\left\|Y(x) - \sum_{i=1}^{n} a_i M_i(x)\right\|^2\right]$$

And we only have access to data $\{X_i, Y(X_i)\}_{i=1}^{N}$. So we could pick a Machine Learning Method, learn over the training set and extrapolate for all $x$

$$\hat{\alpha}_E = \operatorname*{argmin}_{a} \sum_{k=1}^{N}\left[\left\|Y(X_k) - \sum_{i=1}^{n} a_i(X_k) M_i(X_k)\right\|^2\right]$$

(This is Mixture-of-Experts with frozen experts)
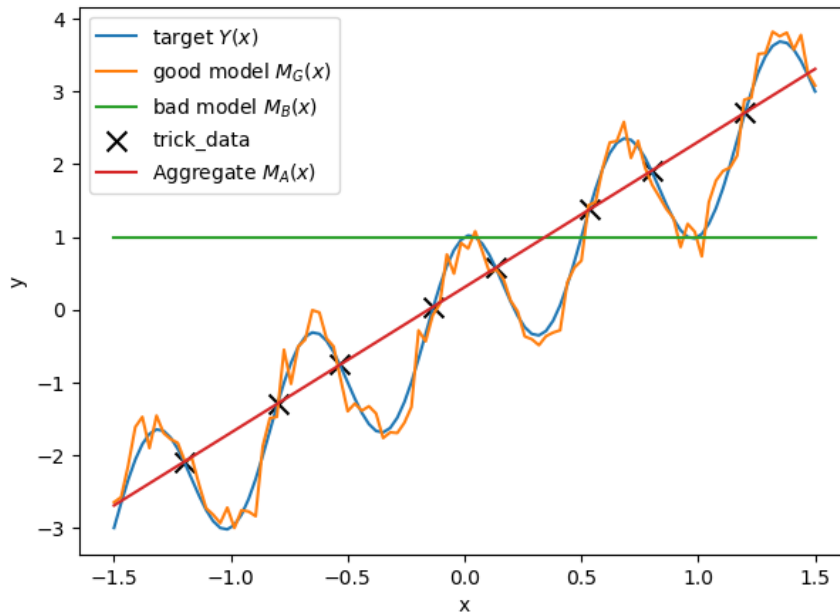
8

**Caltech**

# A pathological example

Given the target, models and data:

- Take $\alpha$ linear

  $$\alpha(x) = (a_G x + b_G, a_B x + b_B)$$

- Train using empirical MSE

Notice that:

- For each data point, the good model performs better than the bad model



- The aggregate ignores the good model and interpolates the data
- Aggregation uses models as features, not approximations of Y

9

**Caltech**

# Minimal Variance Aggregation

Problem: we don't have enough constraints / we didn't define what a good model is.
Let

$$\begin{cases} M_1(x) = Y(x) + \epsilon_1(x) \\ \qquad \vdots \\ M_n(x) = Y(x) + \epsilon_n(x) \end{cases} \text{ where } \begin{cases} \text{(For simplicity)} \quad \epsilon_i \text{ are independent} \\ \qquad \text{(Write)} \qquad Var[\epsilon_i(x)] = V_i(x) \\ \text{(Assumption)} \qquad \mathbb{E}[\epsilon_i(x)] = 0 \end{cases}$$

Then the aggregation is unbiased if

$$\sum_{i=1}^{n} \alpha_i = 1$$

**Caltech**

# Minimal Variance Aggregation

$$\alpha_V(x) = \operatorname{argmin} \mathbb{E}\left[\left|\left|Y(x) - \sum_{i=1}^{n} a_i M_i(x)\right|\right|^2\right]$$

$\sum_{i=1}^{n} \alpha_i = 1$

$\Longrightarrow$

$$\alpha_V(x)^T M(x) = \frac{\sum_{i=1}^{n} \frac{1}{V_i(x)} M_i(x)}{\sum_{i=1}^{n} \frac{1}{V_i(x)}}$$

Note that:

- If $\epsilon_i(x) \sim \mathcal{N}(0, V_i(x))$, this is MLE
- If $\mathbb{E}[\epsilon_i(x)] = 0$, $\mathbb{E}[\epsilon_i(x)^2] = V_i(x)$, this is BLUE
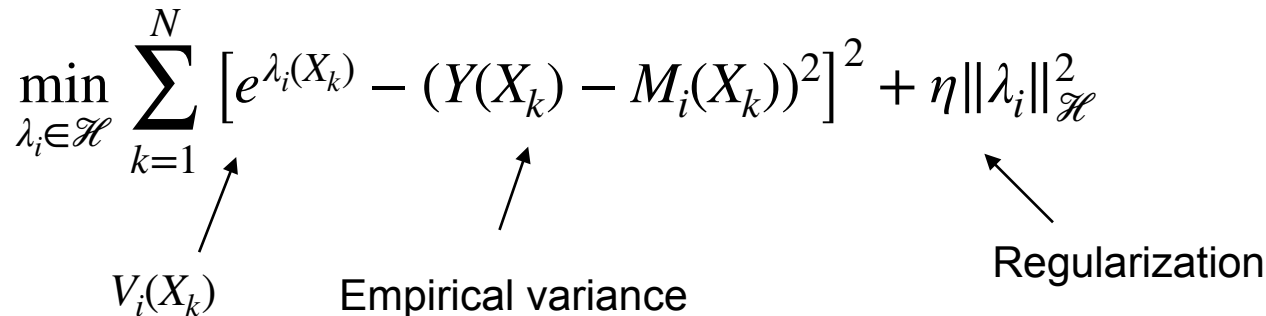- If no assumption, best convex combination

We just need to learn $V_i(x)$, the expected error of each model

11

**Caltech**

# Learning the variance/error

To predict the variance, we:

- Write $V_i(x) = e^{\lambda_i(x)}$ where $\lambda_i$ is a Machine Learning method (Gaussian process, neural network…) to ensure positivity
- Then the aggregation is a softmax
- Use the loss

$$\min_{\lambda_i \in \mathcal{H}} \sum_{k=1}^{N} \left[ e^{\lambda_i(X_k)} - (Y(X_k) - M_i(X_k))^2 \right]^2 + \eta \|\lambda_i\|_{\mathcal{H}}^2$$

$V_i(X_k)$          Empirical variance                    Regularization

This is different from minimizing the error with a softmax        **Caltech**

# Theorem on linear regression:

Assume samples $(M_j, Y_j)_{j=1}^N$, which one has the best loss $\mathscr{L}(\alpha) = \mathbb{E}[|Y - \alpha^T M|^2]$?

$$\hat{\alpha}_E(x) = \underset{a \in \mathbb{R}^n}{\text{argmin}} \sum_{j=1}^N \left[ \left| Y_j - a^T M_j \right|^2 \right]$$

$$\hat{\alpha}_V(x) = \underset{a \in \mathbb{R}^n}{\text{argmin}} \begin{cases} \sum_{j=1}^N \left[ \left| Y_j - a^T M_j \right|^2 \right] \\ \text{such that } \sum_{i=1}^n a_i = 1 \end{cases}$$

Minimal (Empirical) Error Aggregation

Minimal (Empirical) Variance Aggregation

There exists $\lambda \in [0,1]$ s.t.:

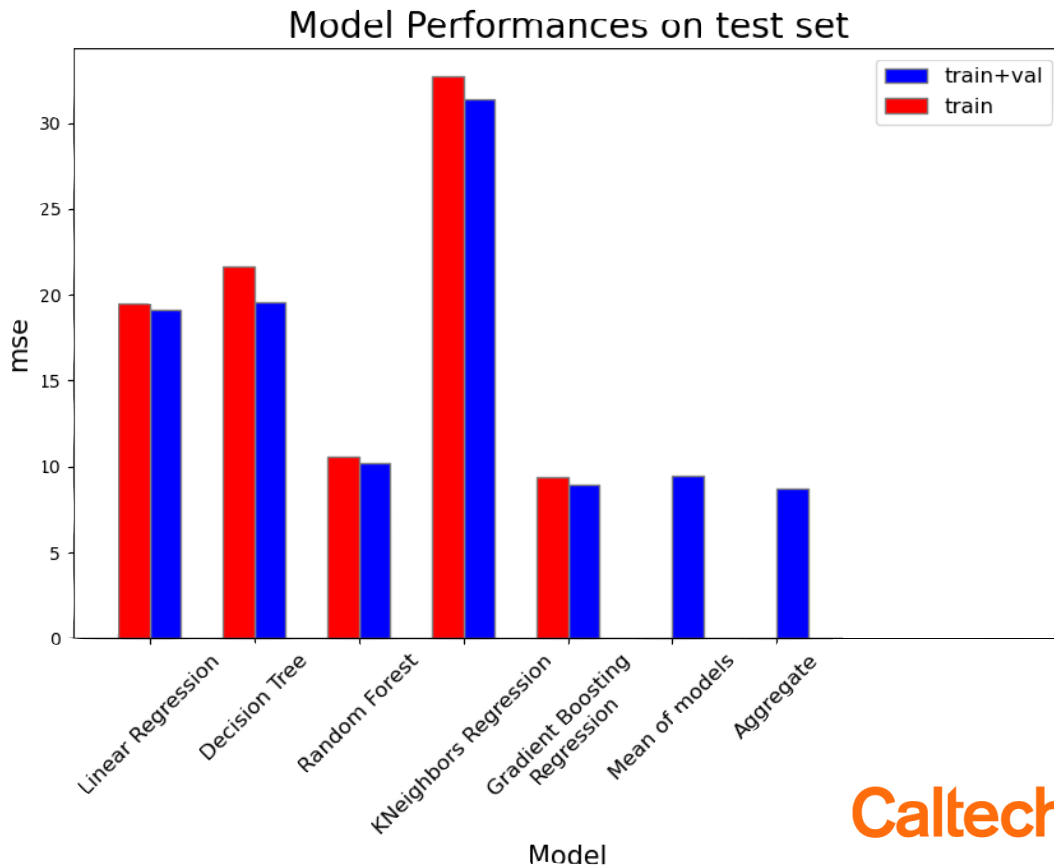$$\mathscr{L}(\hat{\alpha}_E) = \mathscr{L}(\alpha^*) + \mathscr{O}\left( \frac{1}{\sqrt{N}} \right)$$

$$\mathscr{L}(\hat{\alpha}_V) = \frac{1}{\lambda} \mathscr{L}(\alpha^*) + \mathscr{O}\left( \frac{1}{N} \right)$$

In model aggregation, $N$ is small and $\lambda \to 1$

**Caltech**
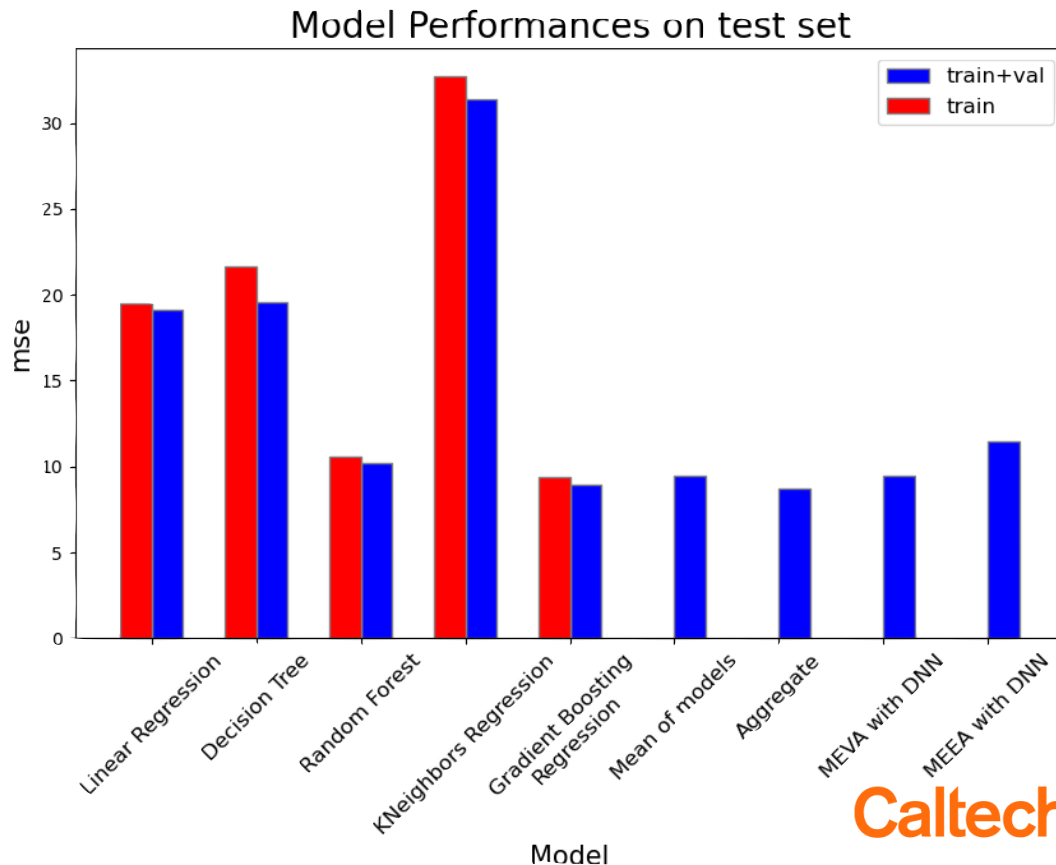
# Applications

Caltech

# The Boston housing dataset

- Data: 506 samples $\{X_i, Y_i\}$

- Data is split into train-test-val

- Aggregation of red models using val data

  - Red models only see train data

  - Blue models for comparison see train+val

- Aggregation is:

  - Better than models aggregated

  - Better than the mean

  - Better than all models



Model Performances on test set

Caltech

# The Boston housing dataset

A comparison with minimal error aggregation:

- Take two identical Neural networks
- Train:
  - To minimize error (bad loss)
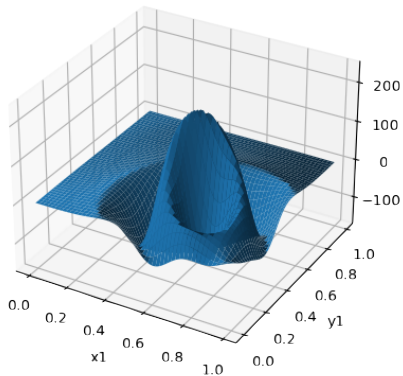  - To estimate variance (our loss)



Model Performances on test set

16

# PDE examples

Given a PDE, we may have multiple solvers/approximations giving a solution.
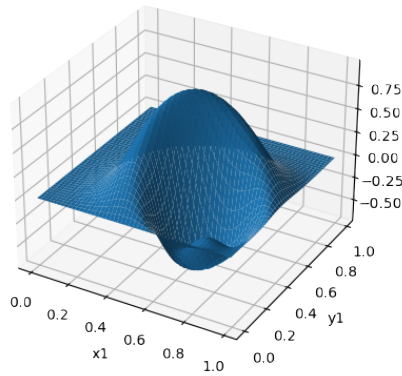
For example:

Laplace equation: $\begin{cases} \Delta u = f & on \ \Omega \\ u = 0 & on \ \partial\Omega \end{cases}$

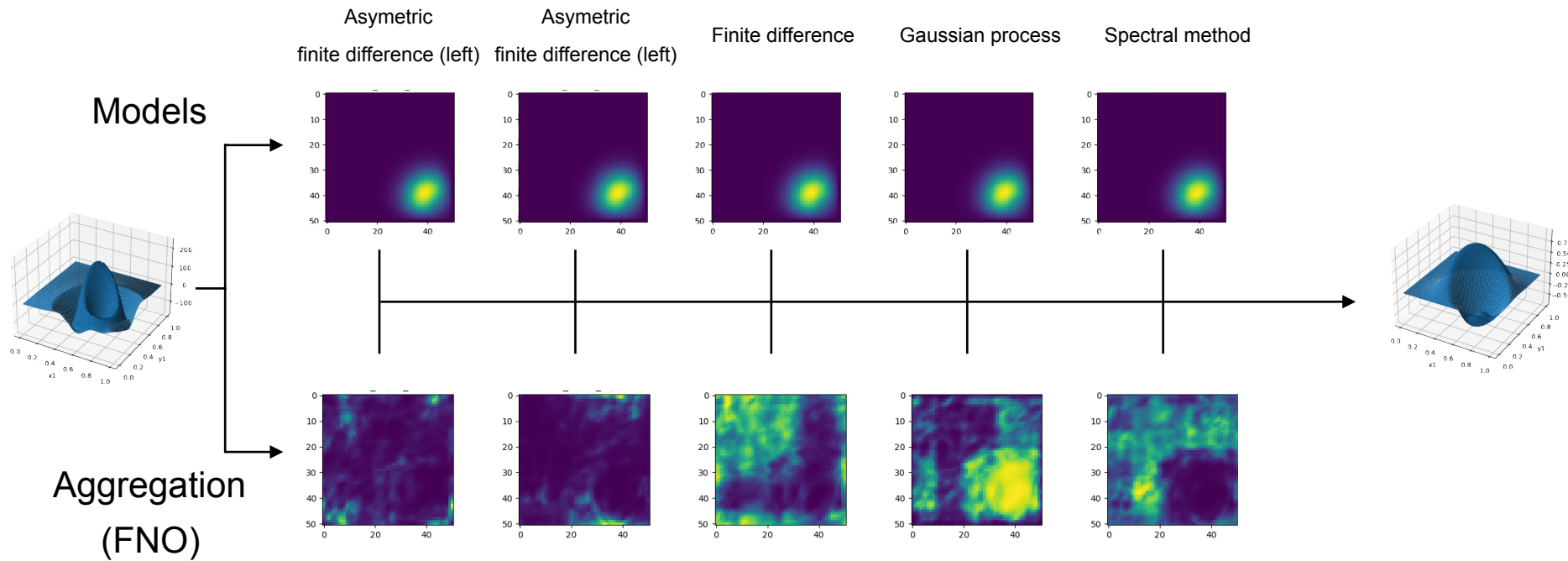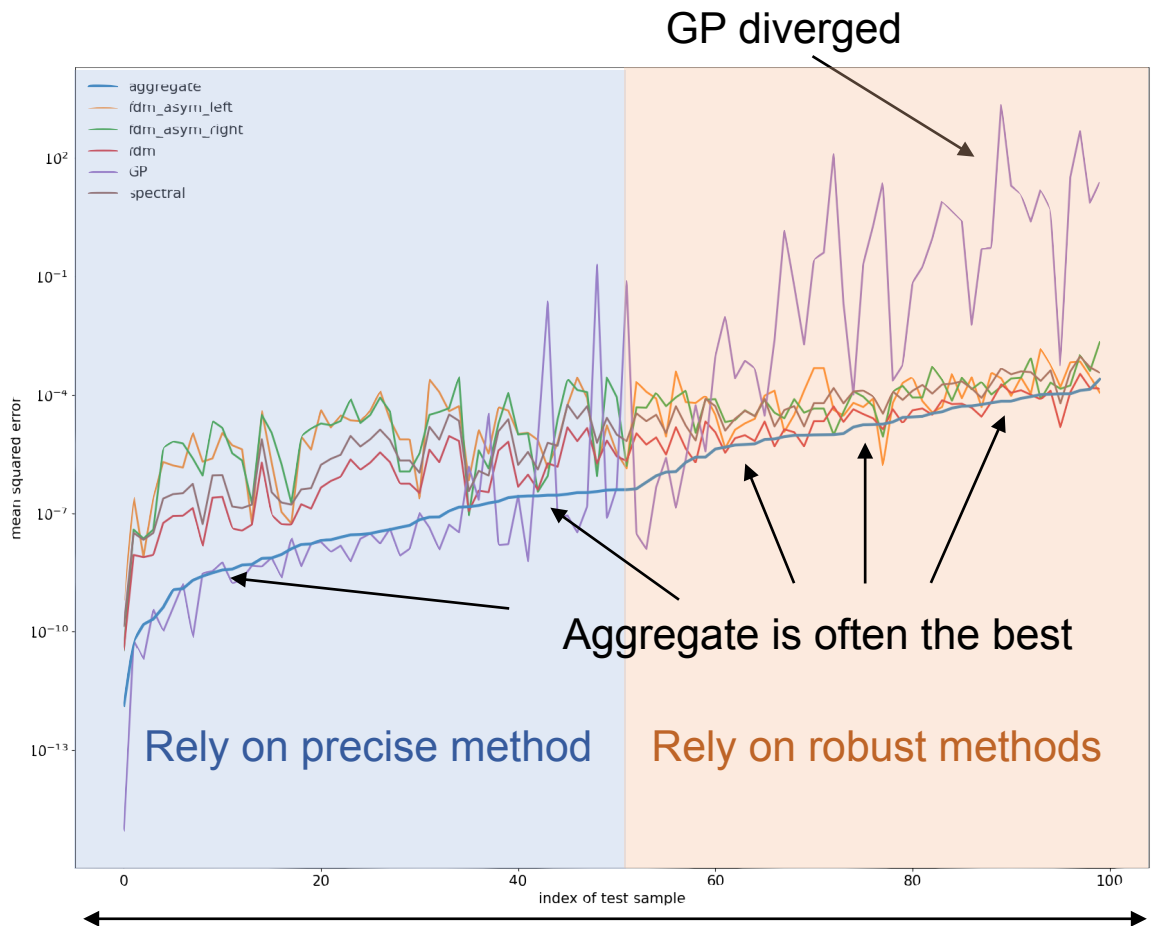Given models $M_i(f) \approx u$, we want to learn the aggregation operator $\alpha(f)$



Random f



Random u

Caltech

# PDE example 1 - Laplace equation

Models

Asymetric
finite difference (left)

Asymetric
finite difference (left)

Finite difference

Gaussian process

Spectral method

Aggregation
(FNO)

Note: FNO uses both $f$ and model outputs to predict aggregation

**Caltech**

GP diverged

Aggregate is often the best

Rely on precise method

Rely on robust methods

| Method | Geometric mean of MSE (log scale) |
|---|---|
| **Aggregate** | **-6.282** |
| FDM | -5.523 |
| Spectral | -4.988 |
| Gaussian process | -4.739 |
| FDM asymetric (right) | -4.685 |
| FDM asymetric (left) | -4.699 |

Easy input                                Hard input

19

**Caltech**

# PDE example 2 - Burger's equation

Consider Burger's equation on $\Omega = [0,1]^2$:

$$\begin{cases} \partial_t u + u\partial_x u = \nu\partial_{xx}u & \text{for} \quad (x,t) \in \Omega \\ \quad u(0,x) = f(x) & \text{for} \quad x \in [0,1] \\ \quad u(t,0) = u(t,1) & \text{for} \quad t \in [0,1] \end{cases}$$
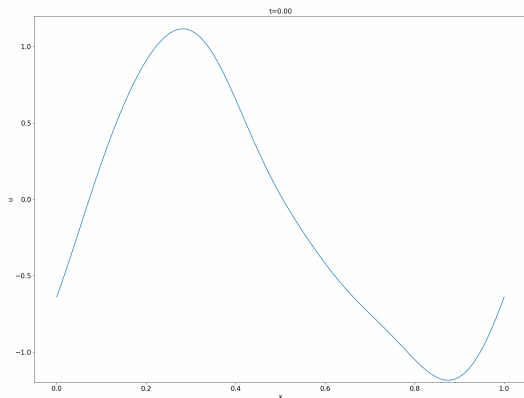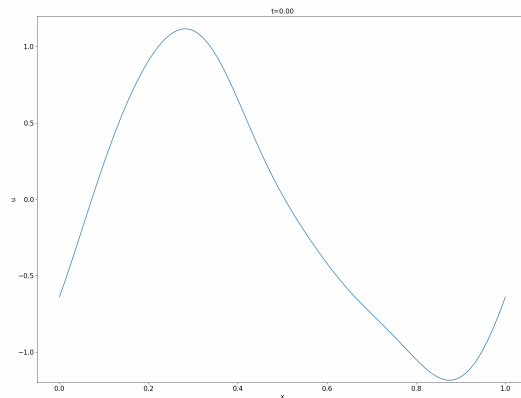
Choose:

- $\nu$ to be small

- $f \sim \mathcal{N}(0,K)$ where $K(x,y) = \exp\left(-\dfrac{2}{l^2}\sin^2\left(\pi|x_i - x_j|^2\right)\right)$
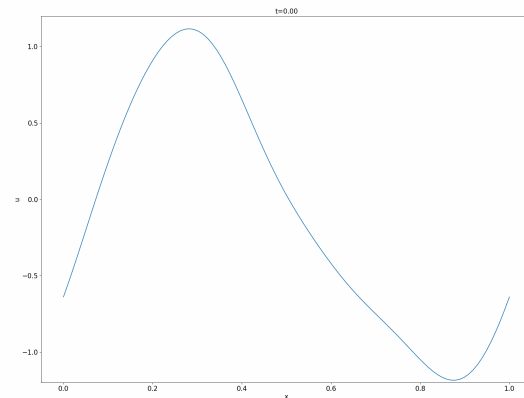
  - i.e. $f$ is periodic and infinitely differentiable

Caltech

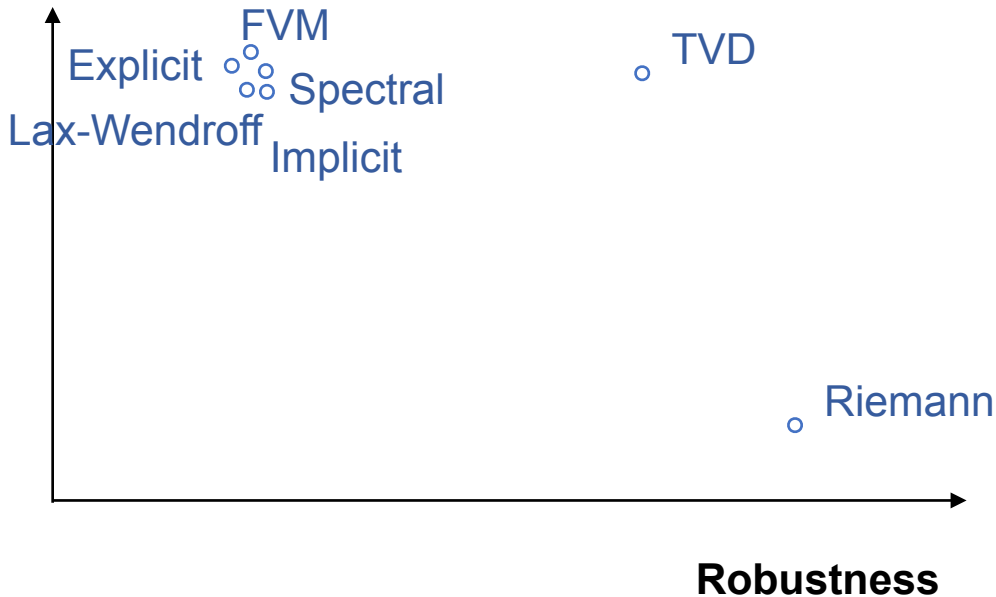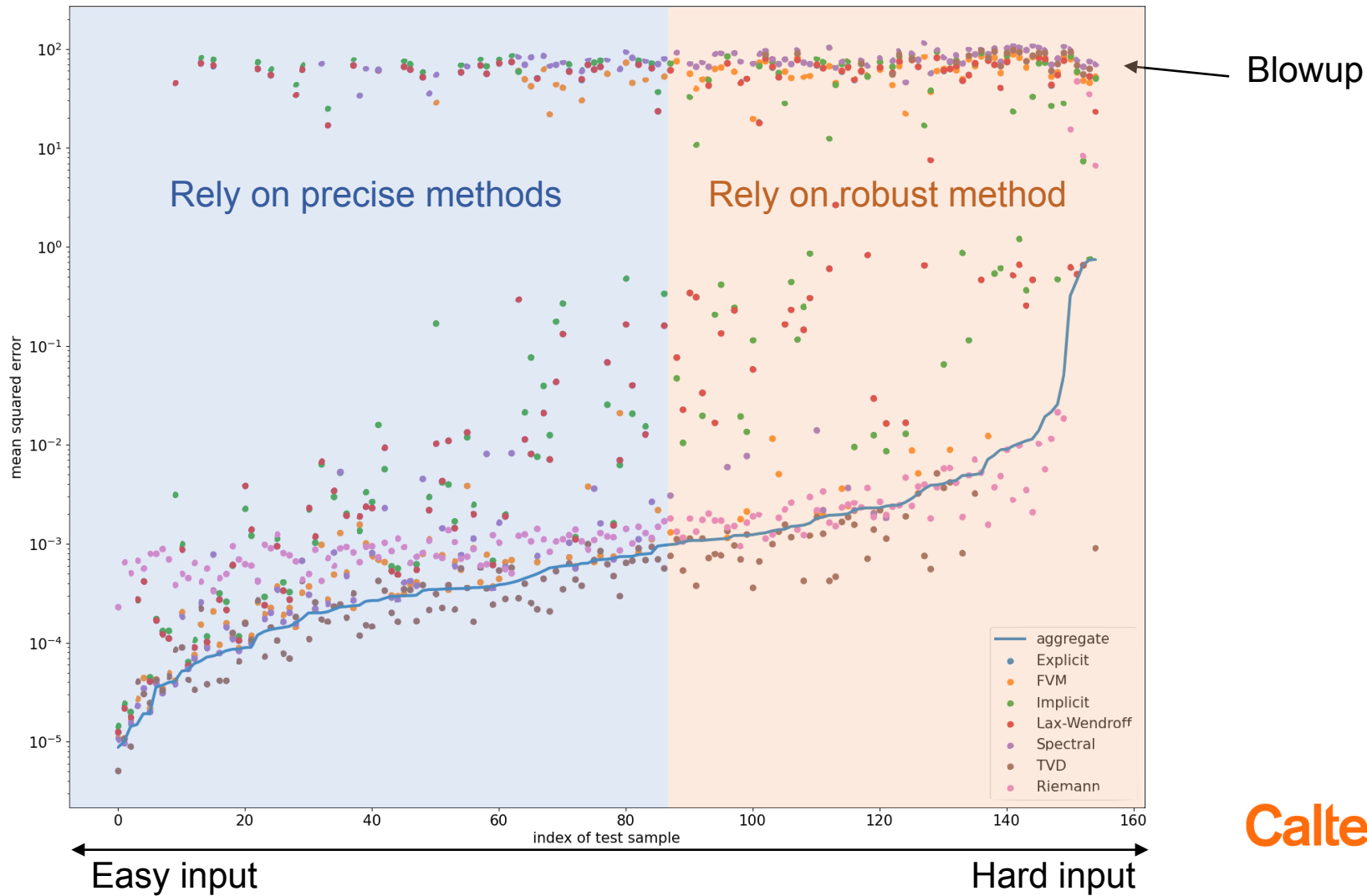# PDE example 2 - Burger's equation



Correct



Blowup



Oscillations

Caltech

# PDE example 2 - Burger's equation

**Accuracy**

FVM

Explicit ○ ○

Spectral ○ ○ ○

Lax-Wendroff

Implicit

○ TVD

○ Riemann

**Robustness**

| Method | Geometric mean of MSE (log scale) |
|--------|-----------------------------------|
| **Aggregate** | **-3.106** |
| Riemann | -2.734 |
| TVD | -2.568 |
| FDM | -1.228 |
| Spectral | -0.625 |
| Implicit | -0.488 |
| Explicit | -0.455 |
| Lax-Wendroff | -0.455 |

Caltech

Rely on precise methods

Rely on robust method

Blowup

Easy input

Hard input

23

Caltech

# Conclusion



We introduce a simple framework to aggregate existing models

- Only requires model output (no assumption, non intrusive)

- Most useful in scientific computing settings with legacy models

- Aggregate any type of methods (ML, solvers…)

**Bourdais, T.,** & Owhadi, H. (2025).
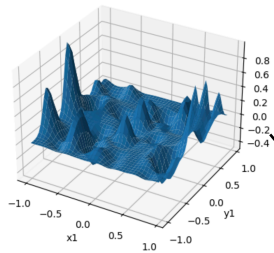*Minimal Variance Model Aggregation: A principled, non-intrusive, and versatile integration of black box models*
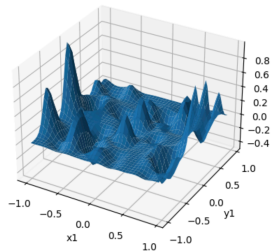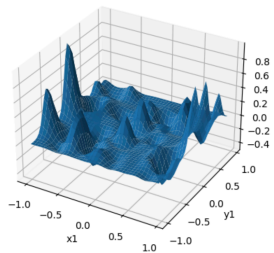ICLR 2025
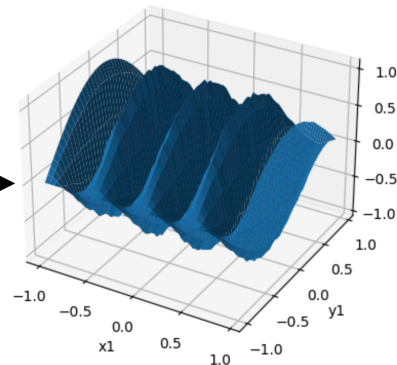
**Caltech**

Predictions

$M_1(x)$

$M_2(x)$

$M_n(x)$

We want to create this !

Model Aggregation

True solution

$\approx$

$Y(x)$

Caltech

# **Summary**

$$M_1(x)$$

Given $\qquad\vdots\qquad f(x, M_1(x), \ldots, M_n(x)) \;\approx\; Y(x)$

$$M_n(x)$$

**Caltech**

# **Summary**

Given

$M_1(x)$

$\vdots$

$M_n(x)$

$$\sum_{i=1}^{n} \alpha_i(x) M_i(x) \approx Y(x)$$

Where

$$\alpha^*(x) = \operatorname*{argmin}_{a \in \mathbb{R}^n} \mathbb{E}\left[\left\|Y(x) - \sum_{i=1}^{n} a_i M_i(x)\right\|^2\right]$$

1. Simplification + Gaussian ideal case

**Caltech**

# **Summary**

Given

$$M_1(x)$$
$$\vdots$$
$$M_n(x)$$

No assumption

$$\sum_{i=1}^{n} \alpha_i(x) M_i(x) \approx Y(x)$$

Where

$$\alpha^*(x) = \arg\min_{a \in \mathbb{R}^n} \sum_{k=1}^{N} \left[ \left\| Y(X_k) - \sum_{i=1}^{n} \alpha_i(X_k) M_i(X_k) \right\|^2 \right]$$
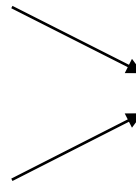
$$\{X_i, Y(X_i)\}_{i=1}^{N}$$

1. Simplification + Gaussian ideal case

2. Directly minimize error
   Does not work

**Caltech**

# Summary

Given

$$M_1(x)$$
$$\vdots$$
$$M_n(x)$$

$$\sum_{i=1}^{n} \alpha_i(x) M_i(x) \quad \approx \quad Y(x)$$

$$\mathbb{E}[M_i(x)] = 0$$

Where $\quad \alpha_i(x) = \dfrac{\frac{1}{Var[M_i(x)]}}{\sum_{k=1}^{n} \frac{1}{Var[M_k(x)]}}$

1. Simplification + Gaussian ideal case

2. Directly minimize error
   Does not work

3. Assume unbiased models

29

**Caltech**

# Summary

Given
$$M_1(x)$$
$$\vdots$$
$$M_n(x)$$

$$\sum_{i=1}^{n} \alpha_i(x) M_i(x) \quad \approx \quad Y(x)$$

$$\mathbb{E}[M_i(x)] = 0$$

Where $\quad \alpha_i(x) = \dfrac{\frac{1}{Var[M_i(x)]}}{\sum_{k=1}^{n} \frac{e^{-\lambda_k(x)}}{Var[M_k(x)]}}$

1. Simplification + Gaussian ideal case
2. Directly minimize error
   Does not work
3. Assume unbiased models
4. Learn
   $$e^{\lambda_i(x)} \approx Var[M_i(x)]$$

$$\lambda_i = \underset{l \in \mathcal{H}}{\operatorname{argmin}} \sum_{k=1}^{N} \left[ e^{l(X_k)} - (Y(X_k) - M_i(X_k))^2 \right]^2 + \eta \|l\|_{\mathcal{H}}^2$$

$\uparrow$ ML regression

$\uparrow$ $\{X_i, Y(X_i)\}_{i=1}^{N}$

(Neural network, Gaussian process...)

30

**Caltech**

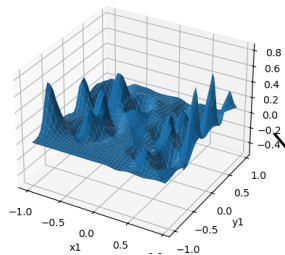# Minimal Variance Aggregation

Let:

$$
\begin{cases}
M_1(x) = Y(x) + \epsilon_1(x) \\
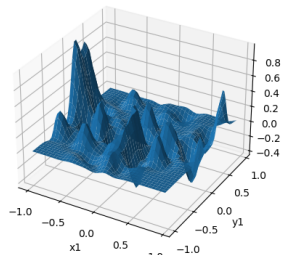\qquad\qquad \vdots \\
M_n(x) = Y(x) + \epsilon_n(x)
\end{cases}
$$

Where

- $\epsilon_i$ are independent (ease of presentation)
- We write $\mathbb{E}\left[|Y(x) - M_i(x)|^2\right] = V_i(x)$
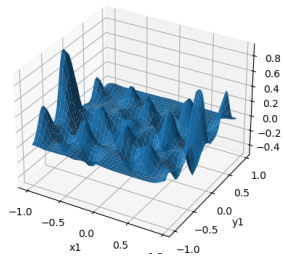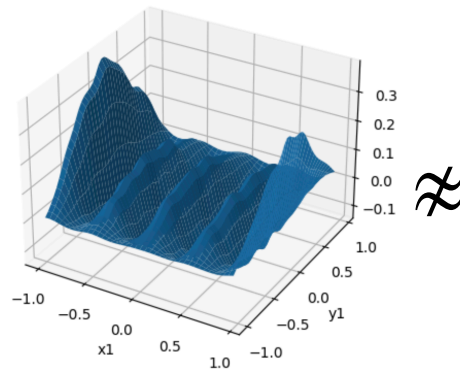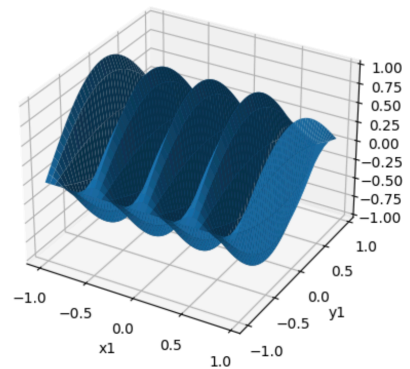
Caltech

Predictions

$M_1(x)$

$M_2(x)$

$M_n(x)$

Average

True solution

$\approx$

$Y(x)$

Caltech