# BAYESIAN OPTIMISATION

## IS PROBABILISTIC NUMERICS

Michael A Osborne, @maosbot

# Global optimisation is proper optimisation.

# Exploitation

# Exploration

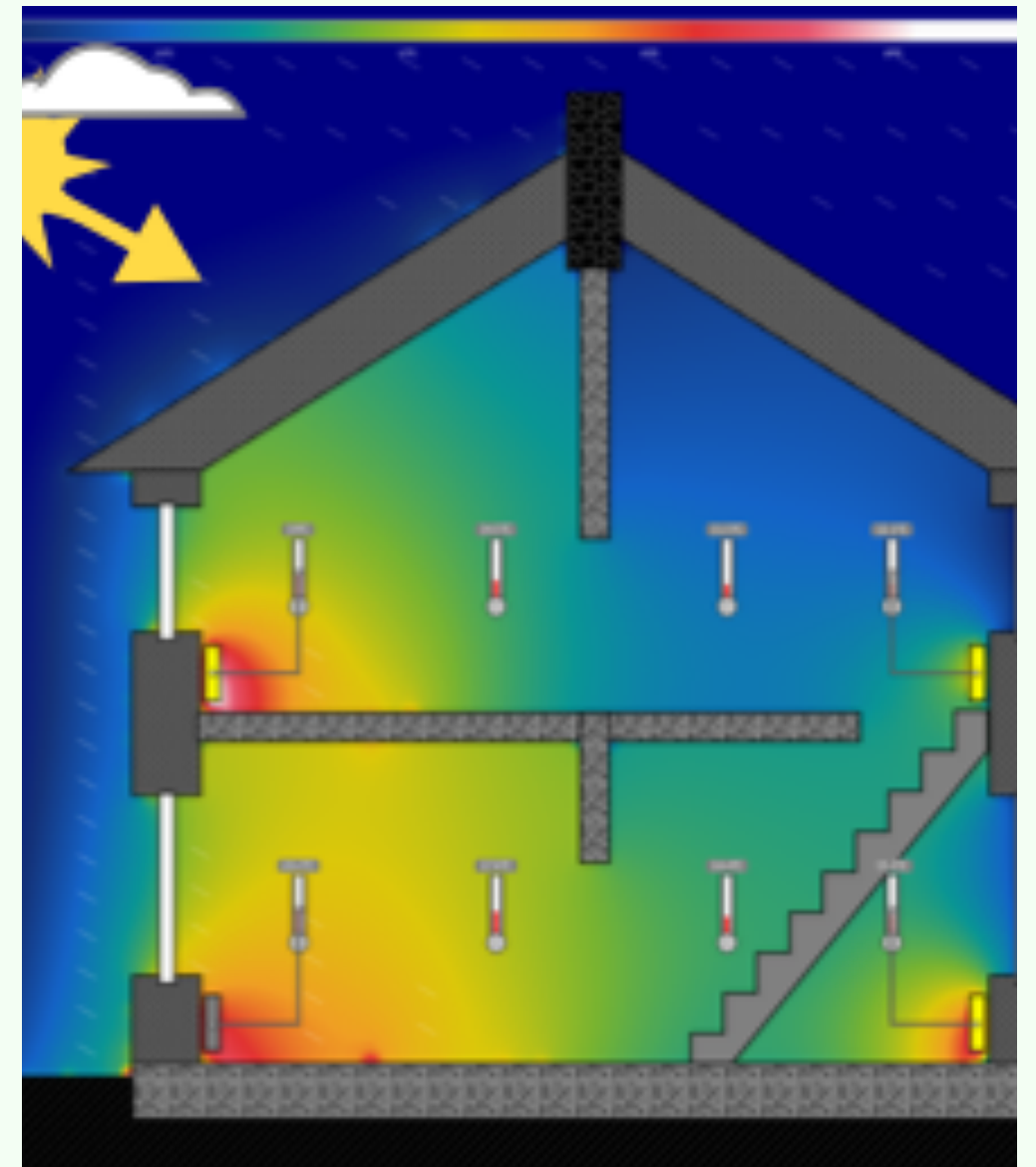# Relative to local optimisation, global optimisation:

1. is less amenable to **theory;**

2. requires higher **overhead;** and

3. overhead costs scale more poorly in **dimension.**

# Global optimisation is widely used.

**Machine learning** treats **algorithms** as agents.

**Probabilistic numerics** treats **numeric algorithms** as agents.

An agent receives **data**, **predicts**, **& then makes decisions**.

# In global optimisation:

data = ?;

predictand = ?; &

decisions = ?.

# In global optimisation:
**data = evaluations;**
**predictand = minimiser; &**
**decisions = locations.**

# **Bayesian optimisation** is probabilistic numerics for global optimisation.

# An agent is defined by its <span style="color:red">prior</span> and loss function.

**The surrogate is the prior for the objective:** options include

Gaussian processes,

random forests,

tree-structured Parzen (density) estimators and

Bayesian neural networks.

# TO IMPROVE OPTIMISATION, IMPROVE YOUR SURROGATE.

**Improving calibration** is as important as **improving** accuracy.

# What should we pick as the loss function for optimisation?

# The loss for optimisation could be:

1. the lowest evaluation **(value);** or

2. the uncertainty in the minimiser **(location-information);** or

3. the uncertainty in the minimum **(value-information).**

**1. Value:** $\lambda_{\mathrm{VL}} := y_N$.

**2. Location-information:**
$\lambda_{\mathrm{LIL}} := \mathbb{H}(\boldsymbol{x}_* \mid \boldsymbol{x}_N, y_N, \mathcal{D}_N)$.

**2. Value-information:**
$\lambda_{\mathrm{VIL}} := \mathbb{H}(y_* \mid \boldsymbol{x}_N, y_N, \mathcal{D}_N)$.

The minimiser is $\boldsymbol{x}_*$ and the minimum $y_*$.

# An acquisition function is an expected loss function.

# Most Bayesian optimisation is myopic, in ignoring all but the next evaluation.

# Myopia can lead to insufficient exploration.



# On the other hand, any flaws of a surrogate are magnified by non-myopia.

Myopic expected loss | 2 steps ahead | 3 steps ahead
5 steps ahead | 10 steps ahead | 20 steps ahead

With a myopic strategy, the acquisition function is

$$\alpha(\boldsymbol{x}_n \mid \mathcal{D}_n) = \mathbb{E}\big(\lambda(\boldsymbol{x}_n, y_n, \mathcal{D}_n)\big)$$

$$= \int \lambda(\boldsymbol{x}_n, y_n, \mathcal{D}_n)\, p(y_n \mid \mathcal{D}_n)\, \mathrm{d}y_n.$$

The next evaluation location will be

$$\boldsymbol{x}_n = \arg\min_{\boldsymbol{x}} \alpha(\boldsymbol{x} \mid \mathcal{D}_n).$$

$$x_n = \arg\min_{x} \alpha(x \mid \mathcal{D}_n).$$

# We have succeeded in turning optimisation into optimisation.

# The acquisition function:

is **less expensive** than the objective;

gives us **gradients and Hessians;** and

need not be **optimised exactly.**

# Expected improvement

is a myopic approximation to the value loss:

$$
\begin{aligned}
& \lambda_{\mathrm{VL}}(\boldsymbol{x}_N, f(\boldsymbol{x}_N), \mathcal{D}_N) \\
\simeq\ & \lambda_{\mathrm{EI}}(\mathcal{D}_{n+1}) \\
:=\ & \min_{i \in \{0, \ldots, n\}} f(\boldsymbol{x}_i).
\end{aligned}
$$

Defining the lowest function value available at the $n$th step as

$$
\eta := \min_{i \in \{0, \ldots, n-1\}} f(\boldsymbol{x}_i),
$$

we can simply rewrite the loss as
$\lambda_{\mathrm{EI}}(\mathcal{D}_{n+1}) = \min\{\eta, f(\boldsymbol{x}_n)\}$.

If we have a Gaussian posterior for the next evaluation,

$$p\big(f(\boldsymbol{x}_n) \mid \mathcal{D}_n\big) := \mathcal{N}\big(f(\boldsymbol{x}_n); m(\boldsymbol{x}_n), V(\boldsymbol{x}_n)\big),$$

the **expected improvement acquisition function** is

$$
\begin{aligned}
\alpha_{\mathrm{EI}}(\boldsymbol{x}_n) :=\quad & \mathbb{E}\big(\lambda_{\mathrm{EI}}\big)(\boldsymbol{x}_n) - \eta \\
=\quad & \int_{-\infty}^{\eta} \big(f(\boldsymbol{x}_n) - \eta\big) p\big(f(\boldsymbol{x}_n) \mid \mathcal{D}_n\big)\,\mathrm{d}f(\boldsymbol{x}_n) \\
=\quad & -V(\boldsymbol{x}_n)\mathcal{N}\big(\eta; m(\boldsymbol{x}_n), V(\boldsymbol{x}_n)\big) \\
& + \big(m(\boldsymbol{x}_n) - \eta\big)\,\Phi\big(\eta; m(\boldsymbol{x}_n), V(\boldsymbol{x}_n)\big).
\end{aligned}
$$

$$\alpha_{\mathrm{EI}}(\boldsymbol{x}_n) = \int_{-\infty}^{\eta} \bigl(f(\boldsymbol{x}_n) - \eta\bigr)p\bigl(f(\boldsymbol{x}_n) \mid \mathcal{D}_n\bigr)\,\mathrm{d}f(\boldsymbol{x}_n)$$

# Function Evaluation 1



Legend:
- **·······** Objective function
- **+** Observation
- **———** (red) Mean
- **▬** (pink) ± 1SD
- **———** (blue) Expected loss
- **◆** Chosen position of next observation

Function Evaluation 2

Legend:
- ⋯⋯⋯ Objective function
- + Observation
- —— Mean
- �merk ± 1SD
- —— Expected loss
- ◆ Chosen position of next observation

# Function Evaluation 3



Legend:
- ⋯⋯ Objective function
- + Observation
- — Mean
- ± 1SD
- — Expected loss
- ◆ Chosen position of next observation

Function Evaluation 4

Legend:
- Objective function
- $+$ Observation
- Mean
- $\pm$ 1SD
- Expected loss
- Chosen position of next observation

# Function Evaluation 5



Legend:
- ⋯⋯ Objective function
- + Observation
- — Mean
- ▮ ± 1SD
- — Expected loss
- ◆ Chosen position of next observation

# Function Evaluation 6



Legend:
- **Objective function** (dotted line)
- **Observation** (+)
- **Mean** (red line)
- **± 1SD** (pink shaded region)
- **Expected loss** (blue line)
- **Chosen position of next observation** (blue diamond)

Function Evaluation 7

Legend:
- Objective function (dotted)
- Observation (+)
- Mean (red line)
- ± 1SD (pink shaded)
- Expected loss (blue line)
- Chosen position of next observation (blue diamond)

Function Evaluation 8

Legend:
- Objective function
- + Observation
- Mean
- ± 1SD
- Expected loss
- ◆ Chosen position of next observation

Function Evaluation 9

Legend:
- ⋯⋯ Objective function
- + Observation
- — Mean (red)
- ± 1SD (pink shaded)
- — Expected loss (blue)
- ◆ Chosen position of next observation

# If our evaluations are noisy, the best evaluation ($\eta$) is also probably the most noise-corrupted.

# Probability of improvement

defines (for $\mathbb{I}$ the indicator function) the myopic loss

$$\lambda_{n,\mathrm{PI}}(\mathcal{D}_{n+1}) := \mathbb{I}\big(f(\boldsymbol{x}_n) \geq \eta\big).$$

The probability of improvement acquisition function is hence

$$\alpha_{n,\mathrm{PI}}(\boldsymbol{x}_n) := \mathbb{E}\big(\lambda_{n,\mathrm{PI}}(\mathcal{D}_{n+1})\big) = P\big(f(\boldsymbol{x}_n) \geq \eta \mid \mathcal{D}_n\big).$$
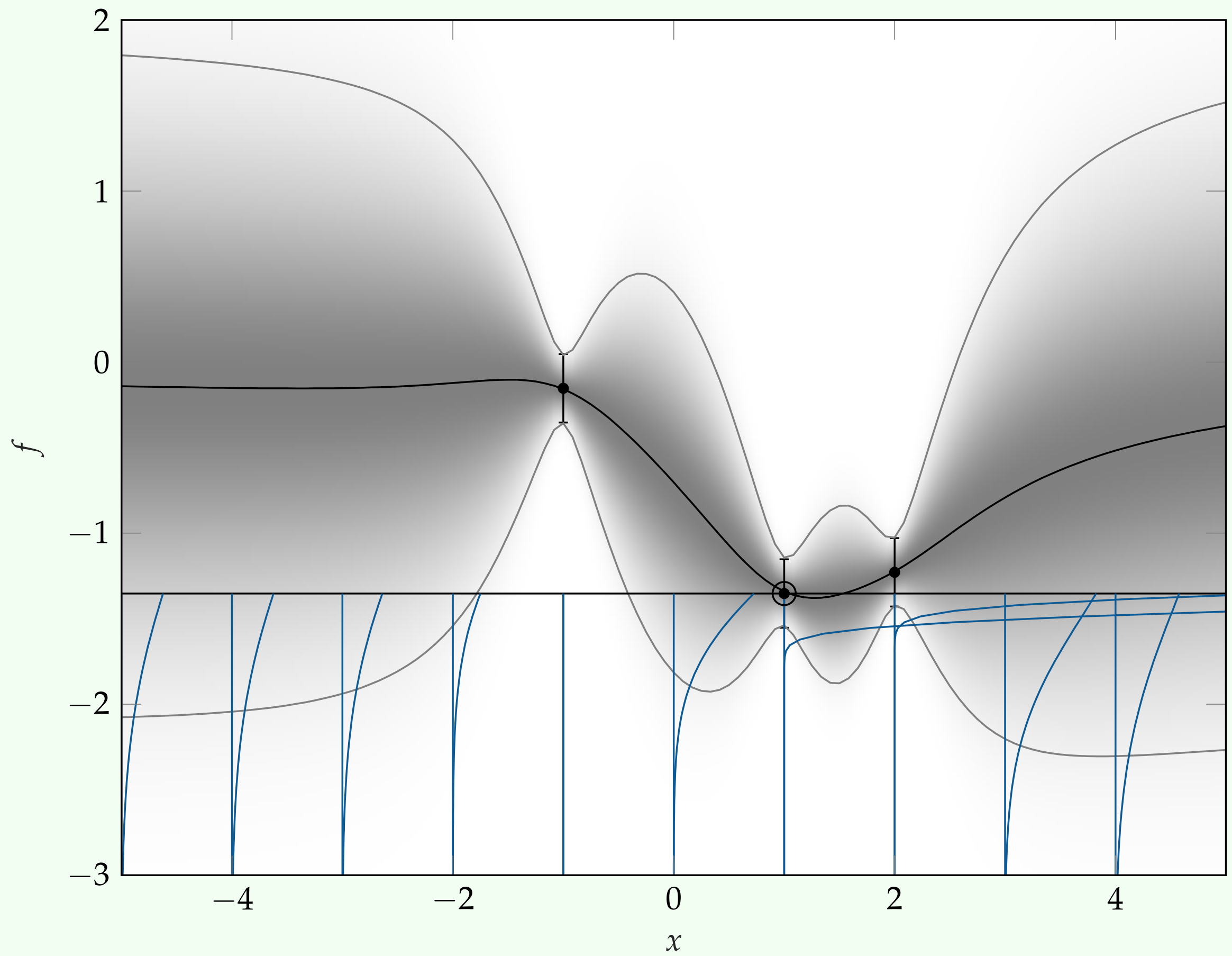
# Probability of improvement

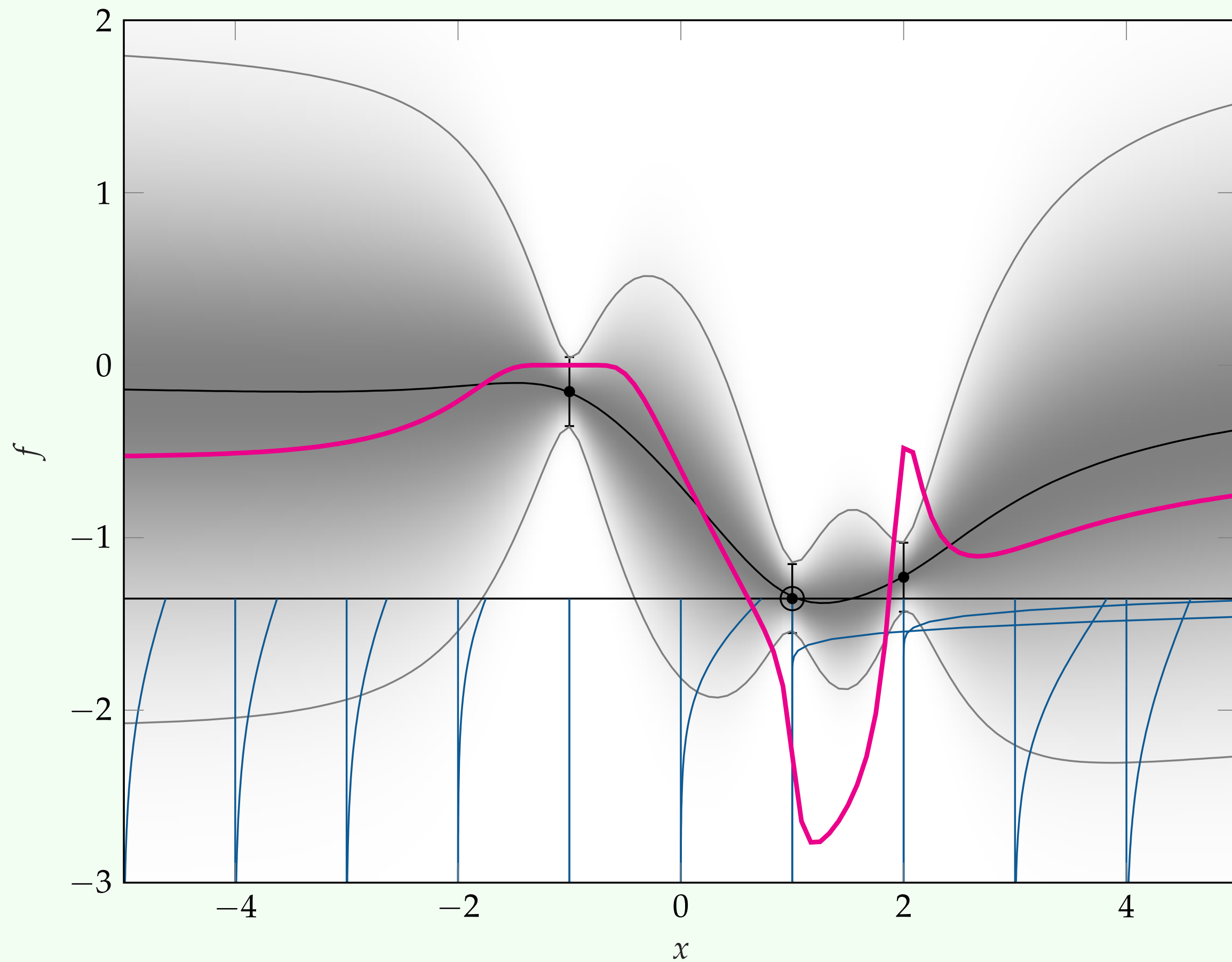defines a myopic loss (for $\mathbb{I}$ the indicator function)

$$\lambda_{n,\mathrm{PI}}(\mathcal{D}_{n+1}) := \mathbb{I}\big(f(\boldsymbol{x}_n) \geq \eta\big).$$

The probability of improvement acquisition function is hence

$$\alpha_{n,\mathrm{PI}}(\boldsymbol{x}_n) := \mathbb{E}\big(\lambda_{n,\mathrm{PI}}(\mathcal{D}_{n+1})\big) = P\big(f(\boldsymbol{x}_n) \geq \eta \mid \mathcal{D}_n\big).$$

**PI values incremental improvement every step.**

# Upper confidence bound

is the myopic acquisition function

$$\alpha_{\mathrm{UCB}}(\boldsymbol{x}_n) := m(\boldsymbol{x}_n) - \beta_n V(\boldsymbol{x}_n)^{\frac{1}{2}}.$$

given a surrogate with mean $m(\boldsymbol{x}_n)$ and variance $V(\boldsymbol{x}_n)$.

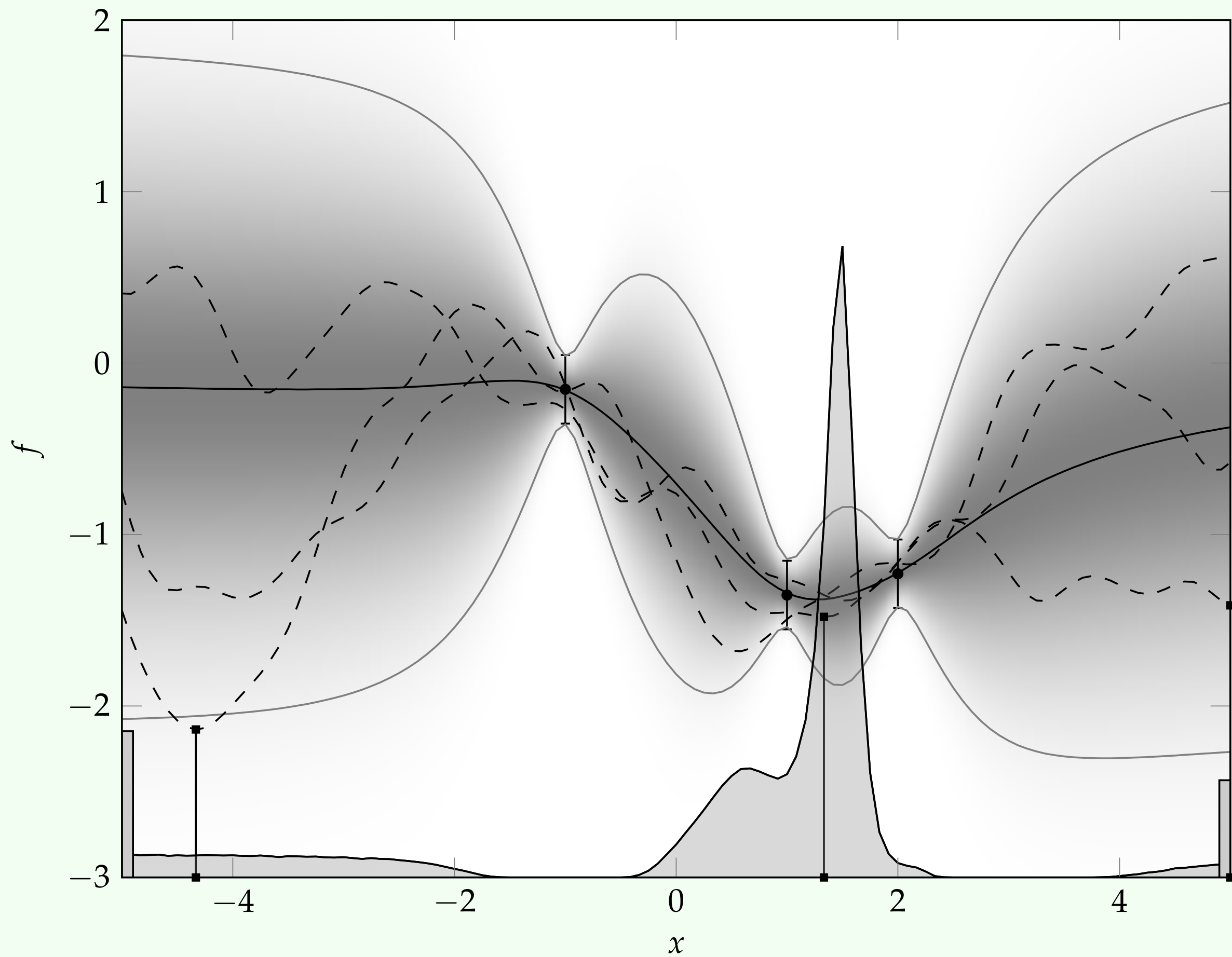**It is difficult to reconcile UCB with a defensible loss function.**
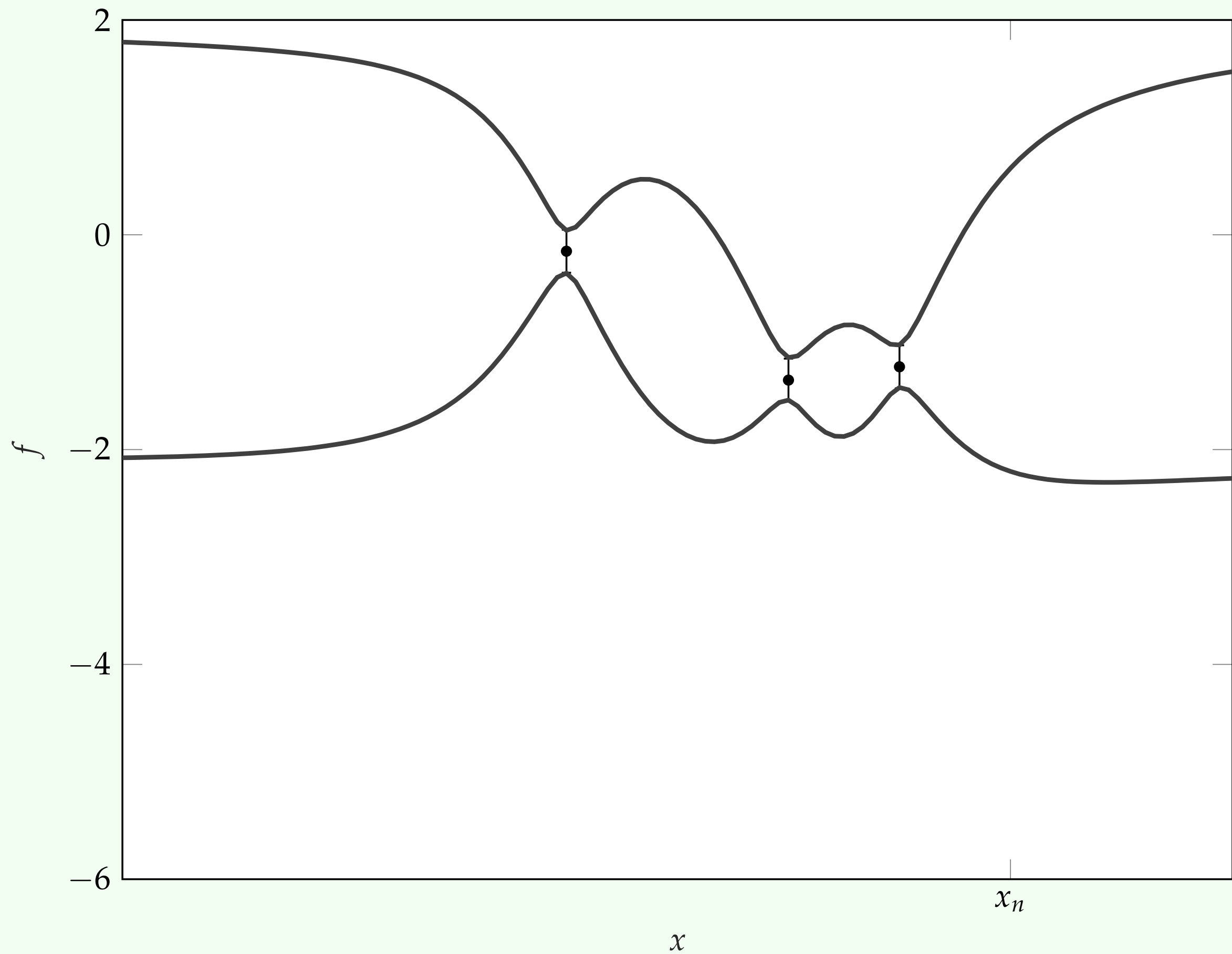
# Information-theoretic methods

give alternative myopic implementations of value-information and location-information losses:
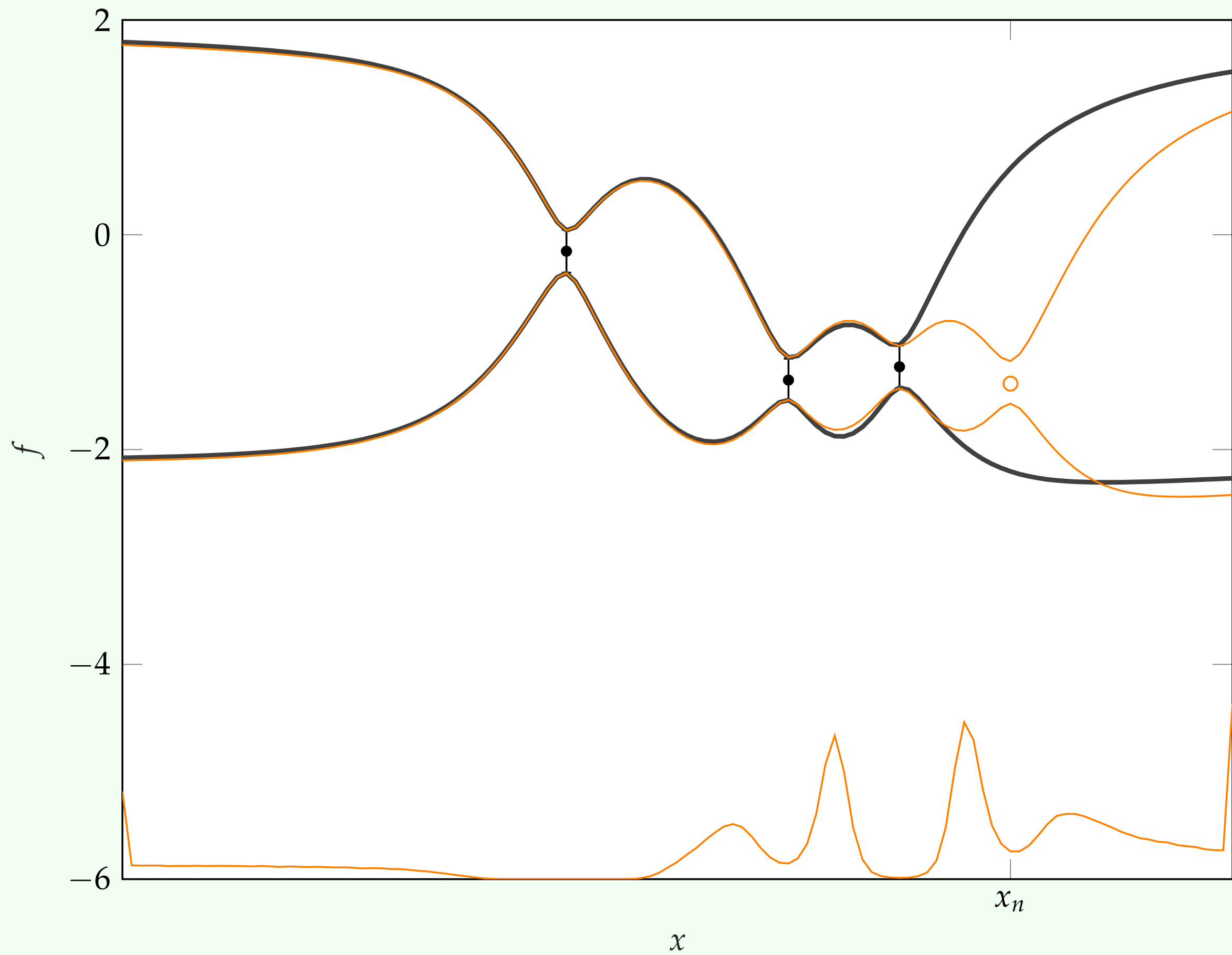
$$\alpha_{\mathrm{LIL}} := \quad \mathbb{E}_{y_n} \; \mathbb{H}(\boldsymbol{x}_* \mid y_n, \boldsymbol{x}_n, \mathcal{D}_n) \quad \text{and}$$

$$\alpha_{\mathrm{VIL}} := \quad \mathbb{E}_{y_n} \; \mathbb{H}(y_* \mid y_n, \boldsymbol{x}_n, \mathcal{D}_n).$$
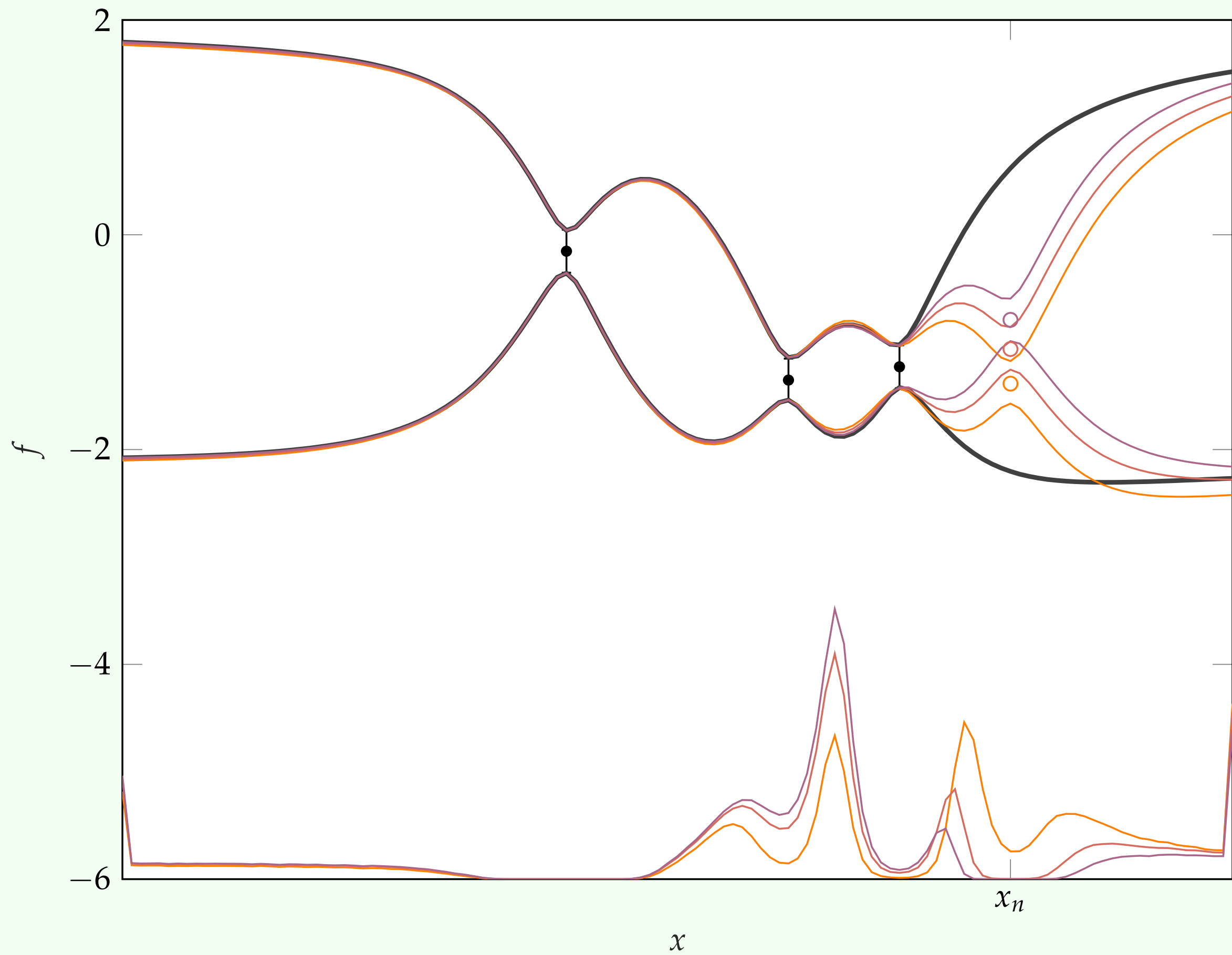
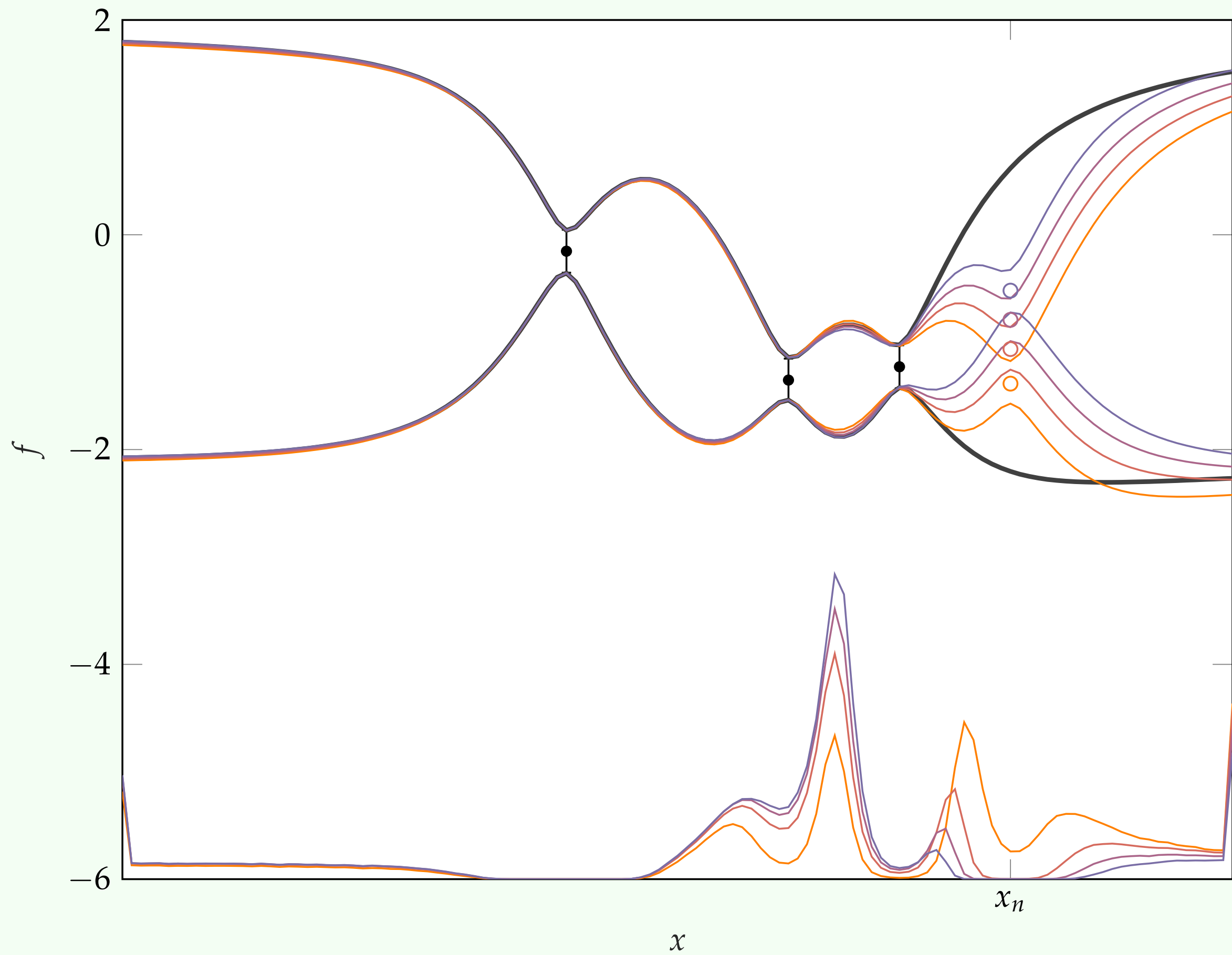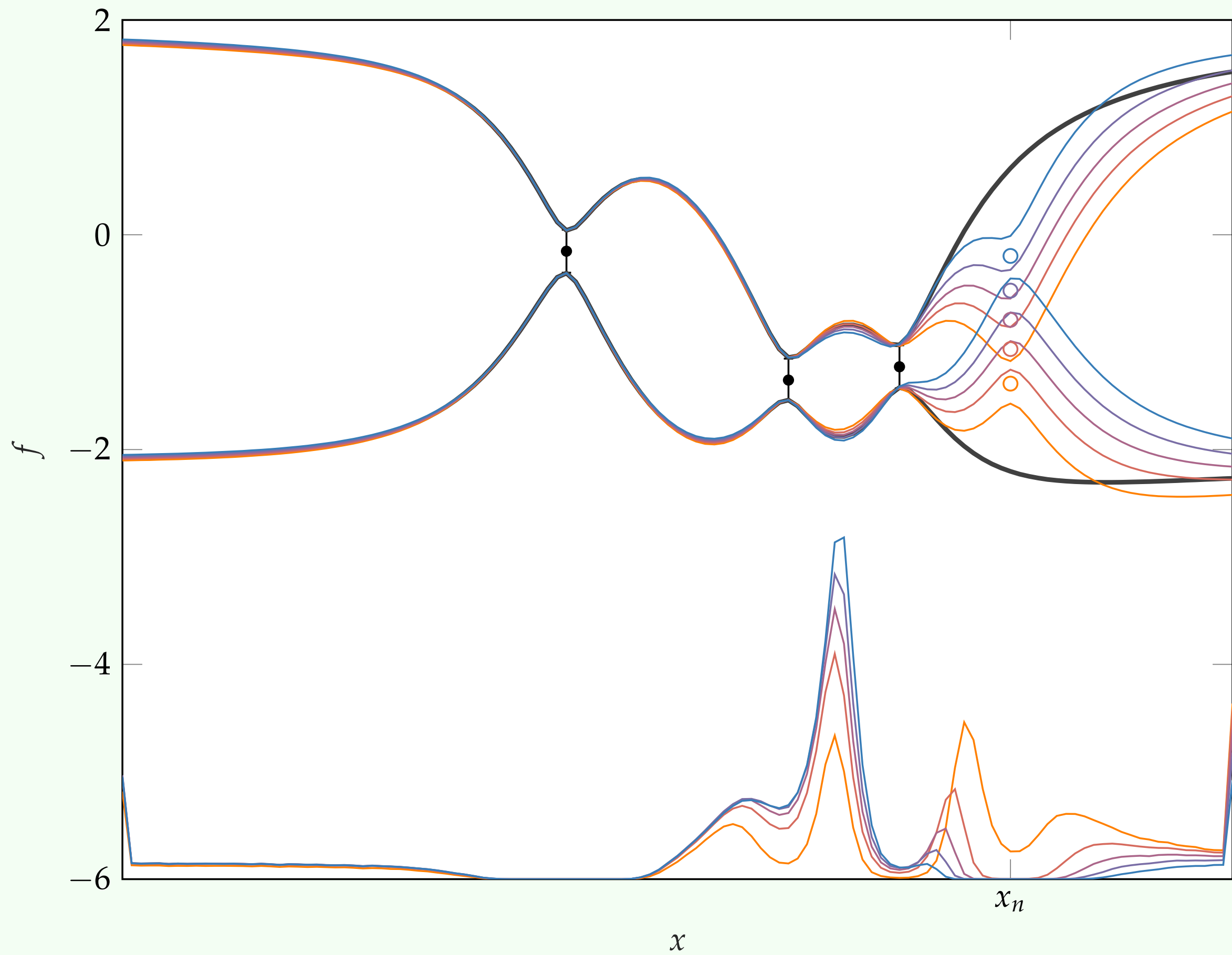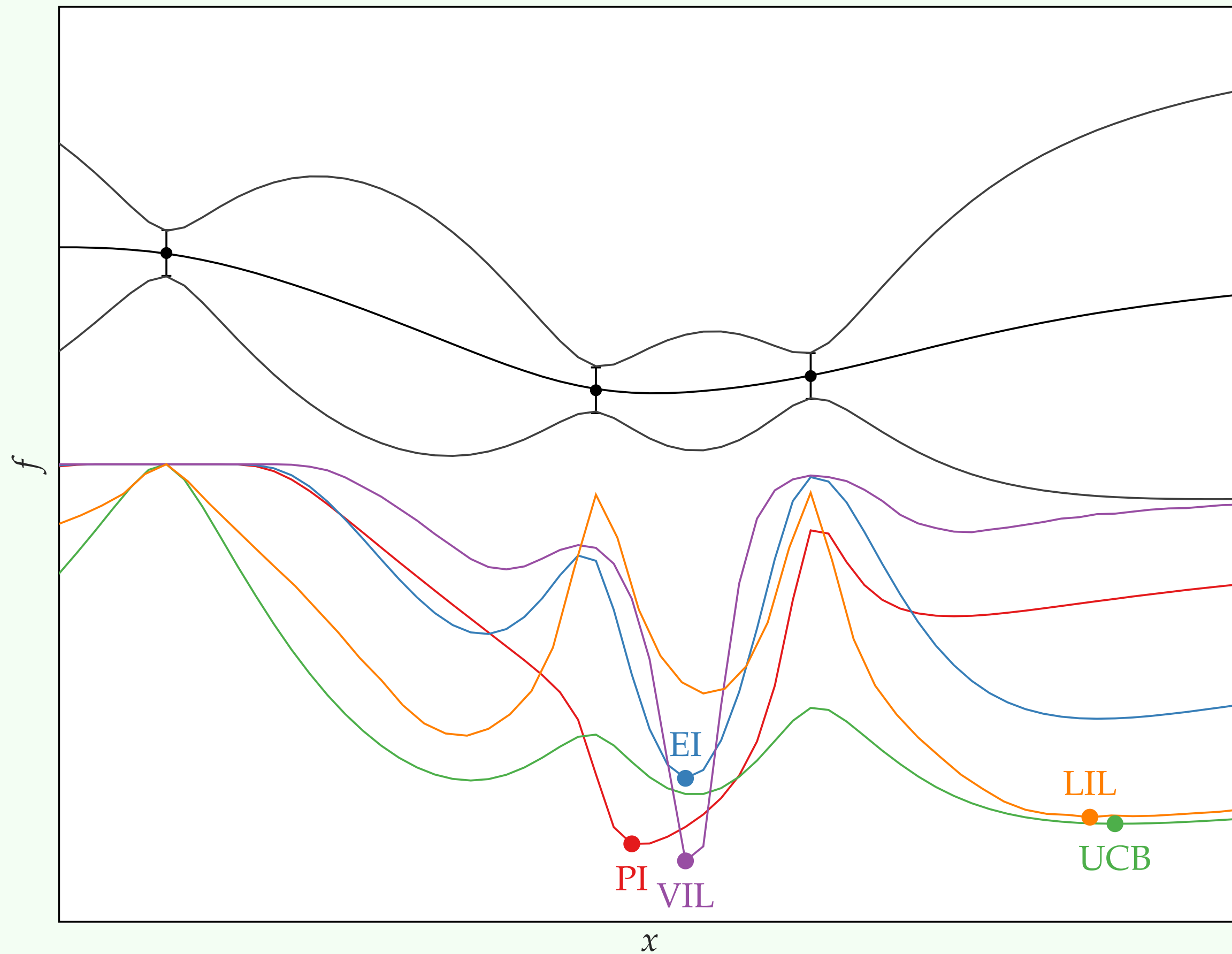These methods tend to be more exploratory, helping performance.

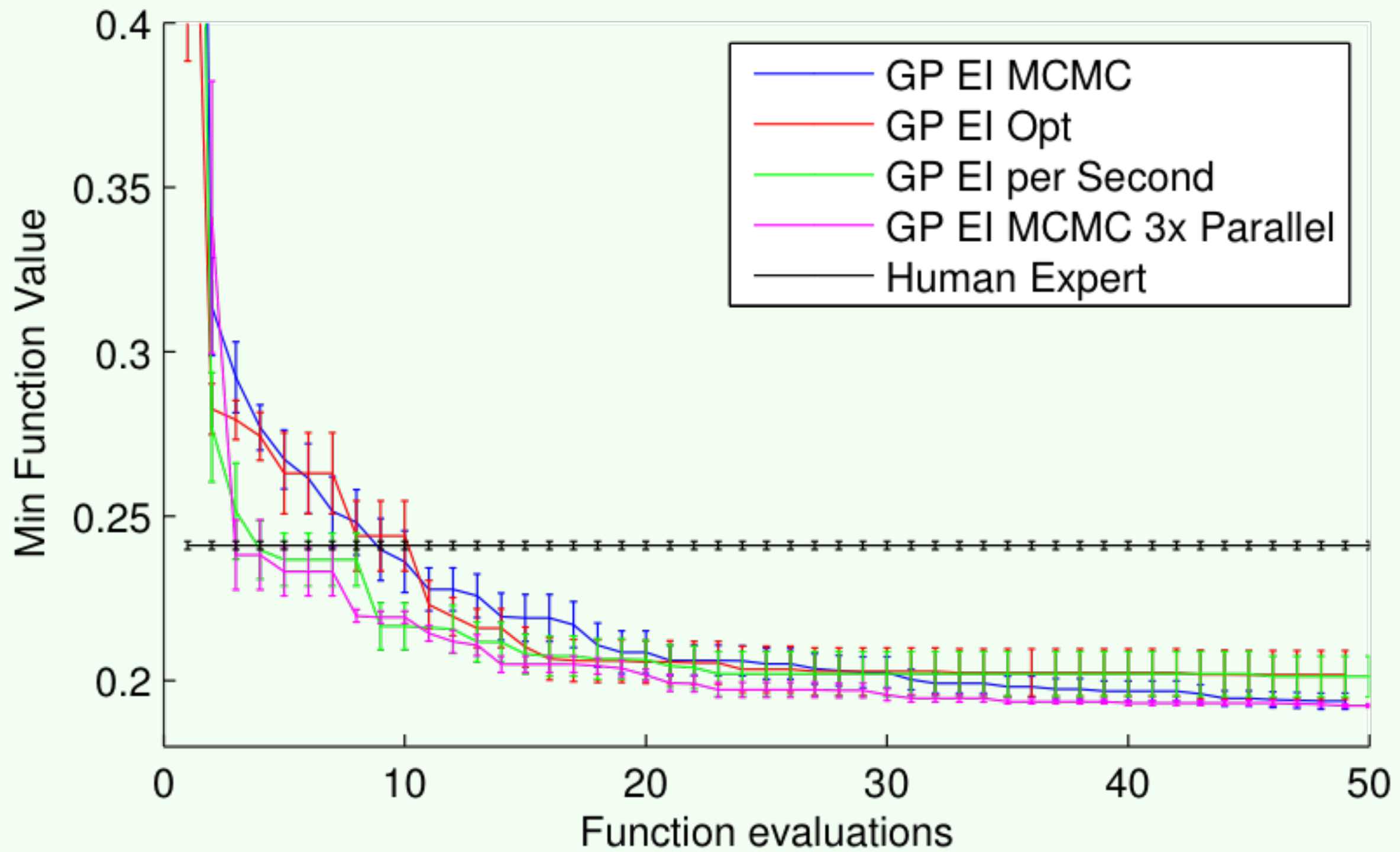# Bayesian optimisation of hyperparameters is used in AutoML.

# Batch Bayesian optimisation is run in parallel.



OEI
QEI
LP-EI
QEI-CL

Rontsis, Osborne, and Goulart. "Distributionally Robust Optimization Techniques in Batch Bayesian Optimization" (2017).

```
ea = params[]_

wa = params[2:3]

secw, sesw = np.sqrt(ea)*np.cos(wa), np.sqrt(ea)*np.sin(wa)
```

# Hyperparameter optimisation is often treated as a black-box optimisation problem.

```
pv_base = np.array([0.1, 0, 0, 0, 0, 0, 0, 0])

                                    1.21, 0.80])

pv_2 = np.array([5.5, -8.23, secw[1], sesw[1], 2.3, 0.79, 1.38])

bx, by = af.baseline_m(df.time, pv_base, df.px, df.py, nthr)

px, py = af.am_model_em(df.time, np.r_[pv_base, pv_1, pv_2], 2, 1

mx, my = 1e6*(bx+px), 1e6*(by+py)
```

It is difficult to imagine a more **white-box problem** than one where you have full access to the problem's **source code.**

```
ea = params[]
wa = params[2:3]

secw, se    = np.sqrt(ea)*np.cos(wa)  np.sqrt(ea)*np.sin(wa)

pv_1 = np.array([6.2, -7.26, secw[0], sesw[0], 1.1, 1.21, 0.80])
pv_2 = np.array([5.5, 8.23, secw[1], sesw[1], 2.3, 0.79, 1.38])

bx, by = af.baseline_m(df.time, pv_base, df.px, df.py, nthr)

px, py = af.am_model_em(df.time, np.r_[pv_base, pv_1, pv_2], 2, 1

mx, my = 1e6*(bx+px), 1e6*(by+py)
```
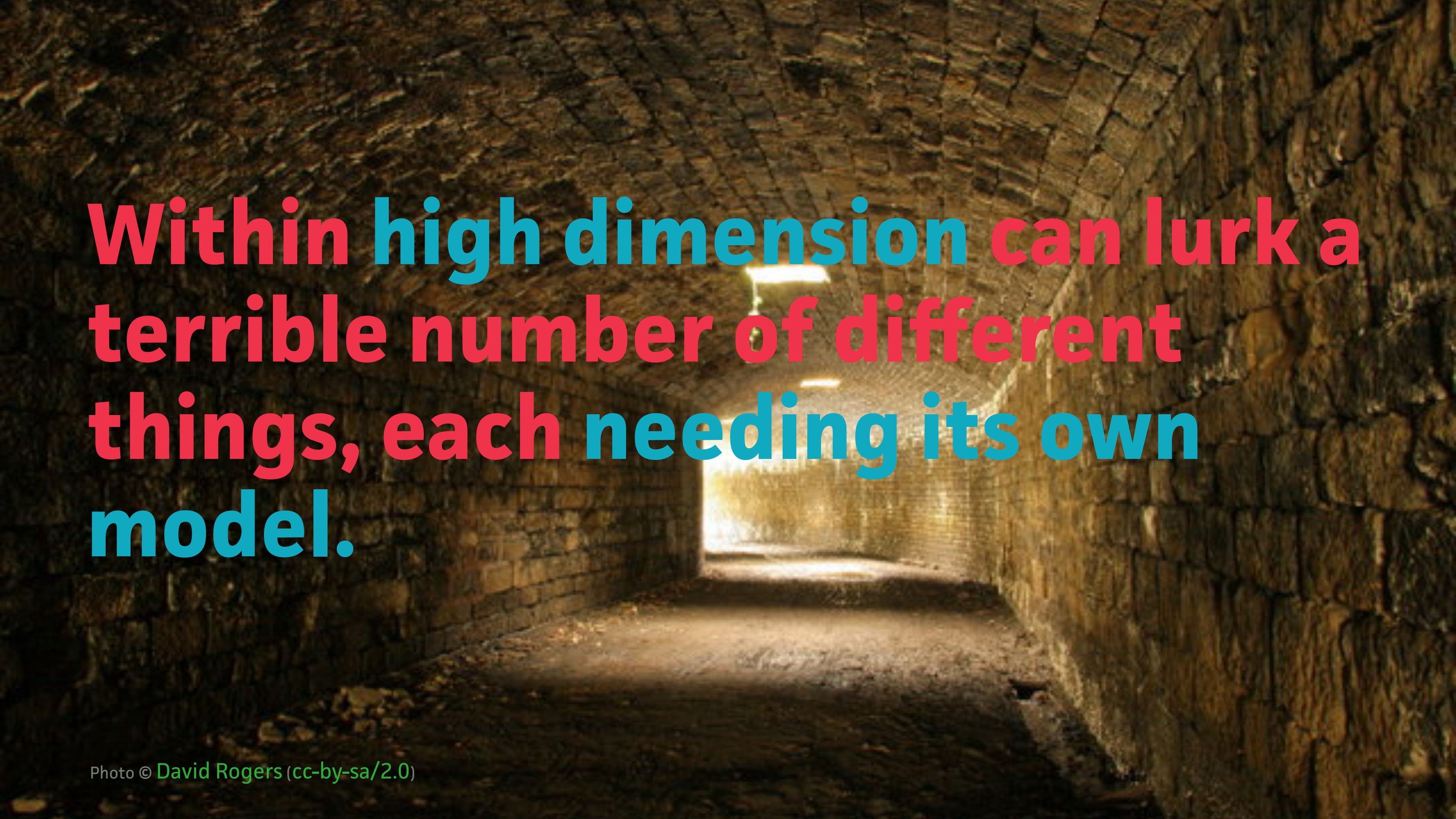
**Within high dimension can lurk a terrible number of different things, each needing its own model.**

# Hyperparameters should usually be **marginalised**, not **optimised**.

# Huge thanks to Roman Garnett & Philipp Hennig.