

Bayesian Optimization With Censored Response Data

Frank Hutter, Holger Hoos, and Kevin Leyton-Brown. University of British Columbia, CS department, Vancouver, Canada.

Cost-varying function minimization

Want to minimize a blackbox function $f: \Theta \rightarrow R$ Observation: *cost* of evaluating $f(\theta)$ often depends on θ

- E.g.: clinical studies (θ specifies a drug)
- E.g.: learning deep networks (θ includes #layers, #neurons/layer, etc)

Def. A *cost-varying* function minimization problem is a tuple (f,c) of blackbox functions $f: \Theta \rightarrow R$, $c: \Theta \rightarrow R$. The *budget* for minimizing f is a limit on the cumulative cost of function evaluations.

Cost-monotonicity

Often, we can cut off function evaluations early and get a lower bound on the function value

- E.g.: clinical studies (e.g. side effects too bad)

Def. A cost-varying function minimization problem (f,c) is *cost-monotonic* iff

 $\forall \boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \boldsymbol{\Theta}. \left(f(\boldsymbol{\theta}_1) < f(\boldsymbol{\theta}_2) \Leftrightarrow c(\boldsymbol{\theta}_1) < c(\boldsymbol{\theta}_2) \right).$

Bayesian Optimization under censoring

- Adaptively censor costly function evaluations
- Integrate censored data points in the model



Base Regression Model:

Due to our high-dimensional, predominantly discrete inputs, we use random forests [Breiman 2001], ensembles of regression trees like this:



Regression Models Under Censoring

Our data: $\{(\theta_i, y_i, c_i)\}_{i=1,...n}$, where c_i is a *censoring indicator*: $y_i = f(\theta_i)$ if $c_i = 0$ and $y_i \le f(\theta_i)$ if $c_i = 1$ Truncated distribution $N(\mu_{\theta}, \sigma_{\theta}^2)_{>\kappa}$ is defined by pdf:

$$p(x) = \begin{cases} 0\\ \frac{1}{\sigma_{\theta}}\varphi(\frac{x-\mu_{\theta}}{\sigma_{\theta}})/(1-\Phi) \end{cases}$$

Direct adaptation of previous EM algorithm [Schmee & Hahn, 1979]: Fit initial random forest, then iterate:

E. For each tree T and each i s.t. $c_i = 1$: $\hat{y}_i^{(T)} \leftarrow \text{mean of } \mathcal{N}(\mu_{\boldsymbol{\theta}_i}, \sigma_{\boldsymbol{\theta}_i}^2) \geq y_i.$ M. Re-fit the random forest u $(\boldsymbol{\theta}, \hat{\boldsymbol{u}})^{(T)} c_{\tau})^{n}$, as the ba

$$(\boldsymbol{v}_i, g_i, c_i)_{i=1}$$
 as the basis

To preserve our uncertainty, change E step to:

E. For each tree
$$T$$
 and each i s.t. $c_i = 1$:
 $\hat{y}_i^{(T)} \leftarrow \text{sample from } \mathcal{N}(\mu_{\theta_i}, \sigma_{\theta_i}^2) \ge y_i.$





 $x < \kappa$ $\Phi(\frac{\mu_{\theta}-\kappa}{\sigma_{\theta}})) \quad x \ge \kappa$

asis for tree
$$T$$

Algorithm configuration (AC)

Def. Given a parameterized algorithm A, a distribution D of problem instances $\pi \in \mathcal{I}$, and a performance metric $m(\theta, \pi)$, let $f(\boldsymbol{\theta}) = \mathbb{E}_{\pi \sim D}[m(\boldsymbol{\theta}, \pi)]$. The algorithm configuration (AC) problem is then to find a parameter setting $\boldsymbol{\theta}$ of A that solves $\arg \min_{\boldsymbol{\theta}} f(\boldsymbol{\theta})$.

A challenge for Bayesian optimization

- High dimensions (e.g. 76 for optimizing CPLEX)
- Mixed discrete/continuous parameters
- Tens of thousands of data points
- Very large non-Gaussian noise
- Time budget (learning & EI opt. counts as part of it!)
- Marginal optimization over heterogeneous instances
- Massive parallelism possible
- Cost-varying problem (sometimes drastic variance)

Sequential Model-based Algorithm Configuration (SMAC) [HHL-B, LION'11]

- State-of-the-art AC procedure
- Handles issues above (using random forests & heuristics)
- Room for improvement in future work:
- + Uncertainty estimates can be overconfident
- + Heuristics: number of runs per setting, which

instances (& which censoring time)

Here: exploit cost monotonicity in SMAC

- Censor runs just above best observed cost: SF * $f(\theta_{inc})$

Improvements of state of the art in AC solving

1.4 1.2 1.2 1.2 1.2 1.2 1.2 1.2 1.2	2.5 T 2.5 T 1.5 L 1.5 S SS CPLEX12-		0 T T SE 1.3 Nº 0 CPLEX 12-M	1.4 1.2 1 0.8 1 1 0.8 1 1 0.8 1 1 0.8	A A A A A A A A A A A A A A A A A A A	SF13 No. SPEAR-SW	4.16 4.15 4.14 4.14 4.13 4.12 5F 1.3 No cens W SPEAR-IBM
Scenario	Unit	SF 1	Medi SF 1.1	an of mea SF 1.3	n runtim SF 1.5	es on tes SF 2	t set No censoring
CPLEX12-CLS	$[\cdot 10^{\circ} s]$	5.27	6.21	6.47	8.3	6.66	21.4
CPLEX12-MASS	$\left[\cdot 10^2 s\right]$	6.39	1.94	2.02	1.94	1.97	2.33
CPLEX12-CORLAT	$[\cdot 10^0 s]$	17.6	9.52	20.5	15.4	16.9	826
CPLEX12-MIK	$[\cdot 10^{-1}s]$	8.88	9.3	9.54	9.45	9.86	23.9
CPLEX12-Regions200	$[\cdot 10^0 s]$	6.93	6.65	6.85	7.21	8.07	12
SPEAR-SWV	$[\cdot 10^{0}s]$	67.2	521	8.15	7.78	290	1030
SPEAR-IBM	$[\cdot 10^4 s]$	1.36	1.36	1.36	1.36	1.36	1.36