

GDR MASCOT-NUM 2019

Brief overview on
Stepwise Uncertainty Reduction for computer experiments

Emmanuel Vazquez

March 19, 2010

Contents

- 1 Stepwise Uncertainty Reduction (SUR)
- 2 SUR for Optimization
- 3 SUR for reliability
- 4 Convergence results
- 5 Concluding remarks

1 Stepwise Uncertainty Reduction

- Framework: Bayesian sequential design of **computer experiments**
- Let ϕ be a **quantity of interest** that depends on a simulator with inputs in \mathbb{X} and outputs in \mathbb{Z}
- Objective: make inference about ϕ from a finite set of experiments on the simulator (runs of the computer model)
- Usual Bayesian approach to describe such experiments: sequence of random variables $Z_1, Z_2, \dots \in \mathbb{Z}$ modeling the outcomes of the experiments at points $X_1, X_2, \dots \in \mathbb{X}$, using a **prior distribution** reflecting our belief about the simulator

- X_1 may be chosen according to an arbitrary distribution
- For all $n = 1, 2, \dots$, X_{n+1} is a decision rule, which may depend on the information

$$I_n = \{(X_1, Z_1), \dots, (X_n, Z_n)\}$$

collected after n experiments

- The Z_i s have a **prior distribution** defined through one or several random processes
- Particular case of a deterministic computer model $f : \mathbb{X} \rightarrow \mathbb{Z}$: we assume $Z_i = \xi(X_i)$, for all $i = 1, 2, \dots$, where ξ is a **random process** that models f
- The sequence of decision rules $\underline{X} = (X_1, X_2, \dots)$ is the **sampling strategy**

- Given an estimator $\hat{\phi}_n$ of ϕ that depends on I_n , how to choose \underline{X} ?
- Starting point of the construction of a **SUR strategy**:
↪ a statistic H_n based on $\hat{\phi}_n$ measuring (residual) **uncertainty about ϕ given I_n** (DeGroot, 1962, Uncertainty, information and sequential experiments)
- H_n may be viewed as a “distance” between ϕ and $\hat{\phi}_n$
- Residual uncertainty H_n is non-negative; $H_n = 0$ corresponds to the absence of uncertainty

Illustrative examples

1. **Approximation** of a deterministic computer model $f : \mathbb{X} \rightarrow \mathbb{R}$

- f modeled using a GP ξ and $Z_i = \xi(X_i)$
- Consider the kriging predictor $\hat{\xi}_n = E_n(\xi) = E(\xi \mid I_n)$
- Measures of uncertainty:
 - **IMSE:** $H_n^{(1)} = E_n \left(\|(\xi - \hat{\xi}_n)\|_{L^2(\mathbb{X}, \mu)}^2 \right) = \int_{\mathbb{X}} \sigma_n(x)^2 d\mu(x)$, where $\sigma_n(x)^2$ is the kriging variance at $x \in \mathbb{X}$
 - **MMSE:** $H_n^{(2)} = \sup_{\mathbb{X}} \sigma_n \geq 1/\mu(\mathbb{X})H_n^{(1)}$

2. **Maximization** of $f : \mathbb{X} \rightarrow \mathbb{R}$

- $M = \sup_{\mathbb{X}} \xi$, $Z_i = \xi(X_i)$ and $M_n = Z_1 \vee \dots \vee Z_n$
- Measure of uncertainty: $H_n = E_n(M - M_n)$

- A SUR strategy consists in selecting the next observation location using the rule

$$X_{n+1} = \arg \min_{x \in \mathbb{X}} J_n(x)$$

where J_n is called **sampling criterion/acquisition function**, and is defined by

$$J_n : x \mapsto \mathbf{E}_n(H_{n+1} | X_{n+1} = x)$$

(note that the expectation in J_n is with respect to Z_{n+1} when $X_{n+1} = x$)

- We can also define a notion of **information gain function**:

$$G_n : x \in \mathbb{X} \mapsto H_n - J_n(x)$$

- Bect et al. 2018 also define the notion of **quasi-SUR strategies** in the case of imperfect minimization of J_n : given a non-negative sequence $(u_n)_{n \in \mathbb{N}}$ with $u_n \searrow 0$, choose X_{n+1} such that $J(X_{n+1}) \leq \inf_{\mathbb{X}} J_n + u_n$

Side note: origins of SUR

- For **sequential design of numerical experiments**, idea proposed by myself and my co-authors for optimization and reliability problems (V. & Piera-Martinez, 2007; Villemonteix et al., 2008; V. & Bect, 2009)
- The parents: **Stepwise Entropy Reduction** for shape recognition (Geman & Jedynak, 1993), and **active learning** (MacKay 1992; Cohn et al. 1996)
- Many greedy Bayesian methods for **risk minimization/utility maximization** can actually be viewed as SUR strategies.

2 Examples of SUR criteria from the literature of optimization

Expected improvement (Mockus, Zilinksas \sim 1970)

- Objective: given a compact set $\mathbb{X} \subset \mathbb{R}^d$ and a smooth function $f : \mathbb{X} \rightarrow \mathbb{R}$, estimate $M = \sup_{\mathbb{X}} f$ using the estimator $M_n = f(X_1) \vee \dots \vee f(X_n)$
- The efficiency of the optimization strategy \underline{X} at iteration n can be measured using the **loss function**

$$\varepsilon_n(\underline{X}, f) = M - M_n$$

- Measure of **residual uncertainty** about M :

$$H_n = \mathbf{E}_n(\varepsilon_n(\underline{X}, \xi)) = \mathbf{E}_n(M - M_n)$$

- SUR strategy:

$$X_{n+1} = \arg \min_{x \in \mathbb{X}} \mathbf{E}_n(H_{n+1} \mid X_{n+1} = x)$$

- In terms of expected gain/information gain:

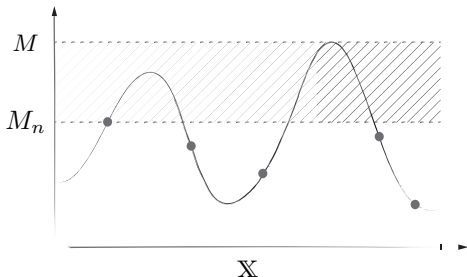
$$\begin{aligned} X_{n+1} &= \arg \max_{x \in \mathbb{X}} H_n - \mathbf{E}_n(H_{n+1} \mid X_{n+1} = x) \\ &= \arg \max_{x \in \mathbb{X}} \mathbf{E}_n(M - M_n) - \mathbf{E}_n(M - M_{n+1} \mid X_{n+1} = x) \\ &= \arg \max_{x \in \mathbb{X}} \mathbf{E}_n(M_{n+1} - M_n \mid X_{n+1} = x) \end{aligned}$$

- $\rho_n(x) = \mathbf{E}_n(M_{n+1} - M_n \mid X_{n+1} = x)$ is called the **expected improvement (EI)** sampling criterion (Mockus, Zilinskas, 1970-1980)

Expected integrated expected improvement

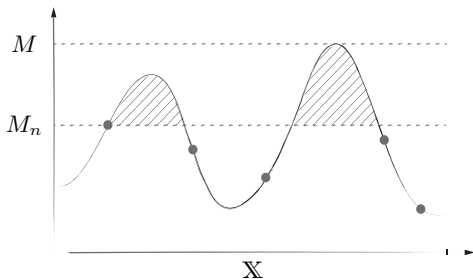
- In a global optimization problem, it is generally of interest to obtain a good approximation of **both** $M = \max_{\mathbb{X}} f$ and $x^* = \arg \max_{\mathbb{X}} f$
- The loss function $\varepsilon_n(\underline{X}, f) = M - M_n$ does not measure directly the distance of x_n^* to x^* (with x_n^* such that $\xi(x_n^*) = M_n$)

- However, note that $\varepsilon_n(\underline{X}, f) = M - M_n \propto \lambda(\mathbb{X})(M - M_n)$



→ coarse measure of the uncertainty about the pair (M, x^*)

- Consider instead the **integral loss** $\varepsilon_n(\underline{X}, f) = \int_{\mathbf{X}} (f(x) - M_n)_+ \lambda(dx)$



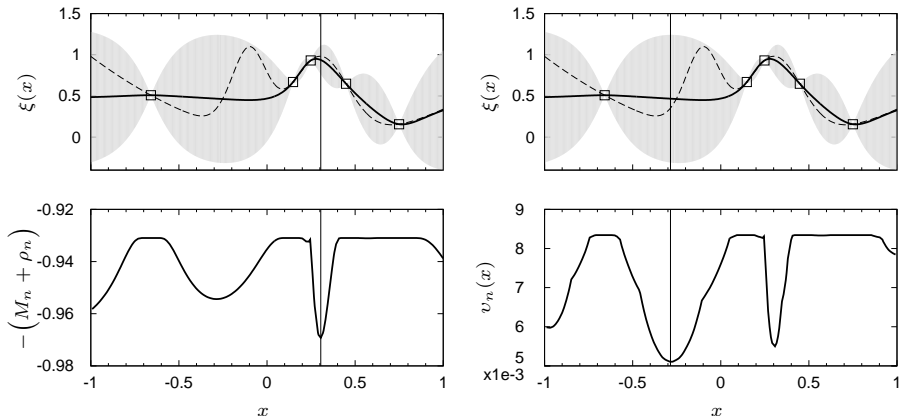
→ ε_n gets smaller when uncertainty about x^* decreases

- Then, a SUR decision rule for choosing a new evaluation point X_{n+1} can be written as

$$\begin{aligned}
 X_{n+1} &= \arg \min_{x \in \mathbb{X}} \mathbf{E}_n \left(\int_{\mathbb{X}} (\xi(y) - M_{n+1})_+ dy \mid X_{n+1} = x \right) \\
 &= \arg \min_{x \in \mathbb{X}} \mathbf{E}_n \left(\int_{\mathbb{X}} \mathbf{E}_{n+1}((\xi(y) - M_{n+1})_+) dy \mid X_{n+1} = x \right) \\
 &= \arg \min_{x \in \mathbb{X}} v_n(x) := \mathbf{E}_n \left(\int_{\mathbb{X}} \rho_{n+1}(y) dy \mid X_{n+1} = x \right)
 \end{aligned}$$

- We call v_n the **Expected Integrated Expected Improvement (EI²)** (V. & Bect 2014)

Large expected improvement in a small region, smaller expected improvement over a large region of the search domain \rightarrow here, v_n favors better exploration than ρ_n



Knowledge Gradient (Frazier et al. 2008)

- Consider a **stochastic computer model** with scalar output \rightsquigarrow observation model:

$$Z_i \mid \xi \stackrel{\text{iid}}{\sim} \mathcal{N}(\xi(X_i), \sigma^2), \quad i = 1, 2, \dots$$

with $\xi \mid m, k \sim \text{GP}(m, k)$

- Objective: estimate $M = \sup_{\mathbb{X}} \xi$
- Which loss function for this problem?

$$\varepsilon_n(\underline{X}, \xi) = M - M_n \leftarrow \xi(X_n^*)$$

where X_n^* is an estimator of $X^* = \arg \max_{\mathbb{X}} \xi(x)$, for instance:

$$X_n^* = \arg \max_{\mathbb{X}} \hat{\xi}_n(x)$$

- Measure of uncertainty: $H_n = E_n(M - \xi(X_n^*))$
- Then, a SUR strategy for this loss function may be written as

$$\begin{aligned}
 X_{n+1} &= \arg \max_{\mathbf{X}} H_n - E_n(H_{n+1} \mid X_{n+1} = x) \\
 &= \arg \max_{\mathbf{X}} E_n(\xi(X_{n+1}^*) - \xi(X_n^*) \mid X_{n+1} = x) \\
 &= \arg \max_{\mathbf{X}} E_n(\xi(X_{n+1}^*) \mid X_{n+1} = x) \\
 &= \arg \max_{\mathbf{X}} E_n(E_{n+1}\{\xi(X_{n+1}^*)\} \mid X_{n+1} = x) \\
 &= \arg \max_{\mathbf{X}} E_n(\widehat{\xi}_{n+1}(X_{n+1}^*) \mid X_{n+1} = x)
 \end{aligned}$$

- The sampling criterion $\rho_n(x) = E_n(\widehat{\xi}_{n+1}(X_{n+1}^*) \mid X_{n+1} = x) - \widehat{\xi}_n(X_n^*)$ is called the **Knowledge Gradient** (Frazier et al. 2008)

Expected hyper-volume improvement (Emmerich 2005)

- **Multi-objective optimization**: consider a set of functions $f_j : \mathbb{X} \rightarrow \mathbb{R}$, $j = 1, \dots, p$, to be minimized
- Objective: build an approximation of the **Pareto front** and of the set of corresponding solutions

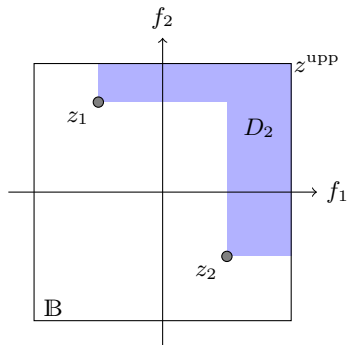
$$\Gamma = \{x \in \mathbb{X} : \nexists x' \in \mathbb{X} \text{ such that } f(x') \prec f(x)\},$$

where \prec stands for the Pareto domination rule defined by

$$y = (y_1, \dots, y_p) \prec z = (z_1, \dots, z_p) \iff \begin{cases} \forall i \leq p, & y_i \leq z_i, \\ \exists j \leq p, & y_j < z_j. \end{cases}$$

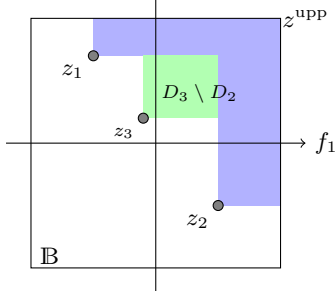
- Given evaluation results $z_1 = f(X_1), \dots, z_n = f(X_n) \in \mathbb{R}^p$, define a set of dominated solutions:

$$D_n = \{y \in \mathbb{B}; \exists i \leq n, f(X_i) \prec y\},$$



- Use the increase of the volume of the dominated region as information gain

$$G_n(X_{n+1}) = |D_{n+1} \setminus D_n| = |D_{n+1}| - |D_n|,$$



\rightsquigarrow expected hyper-volume improvement (Emmerich 2005)

- Extended in (Feliot et al. 2015) to deal with **constrained multi-objective problems** and in (Feliot et al. 2018) to deal with **user preferences for the exploration of a Pareto front**

3 Examples from the literature of reliability

- Let $f : \mathbb{X} \rightarrow \mathbb{R}$ be a deterministic computer model, $u \in \mathbb{R}$, and consider the **excursion set**

$$\Gamma = \{x \in \mathbb{X} : f(x) > u\}$$

- An estimator $\widehat{\Gamma}_n$ of Γ can be obtained by building an **approximation** $\eta_n : \mathbb{X} \rightarrow \{0, 1\}$ of the excess indicator $\mathbb{1}_{f>u}$ from I_n
- For instance, we could use a support vector machine to build η_n (SMART method, Deheeger and Lemaire, 2007)

- Or we can use a GP ξ as a model of f and set

$$\eta_n(x) = \mathbb{1}_{p_n(x) > 1/2} = \mathbb{1}_{\hat{\xi}_n(x) > u}$$

where p_n is the **posterior excursion probability**:

$$p_n : x \in \mathbb{X} \mapsto P_n\{\xi(x) > u\}$$

- Or we can use the notion of **conservative estimates** (Azzimonti et al. 2018): for a high probability β , define

$$\hat{\Gamma}_n^\beta = \arg \min_{Q \in \mathcal{Q}_{n,\beta}} \mu(Q)$$

where $\mathcal{Q}_{n,\beta}$ is the family of Vorob'ev quantiles $Q_{n,\rho} = \{p_n \geq \rho\}$, $\rho \in [0, 1]$, such that $P_n(Q_{n,\rho} \supset \Gamma) \geq \beta$

- Related objectives:

- given a measure/probability μ over \mathbb{X} , estimate the **volume of the excursion set/probability of failure** $\alpha = \mu(\Gamma)$, in which case a natural estimator for α is the posterior mean

$$\hat{\alpha}_n = E_n(\alpha) = E_n\left(\int_{\mathbb{X}} \mathbb{1}_{\xi > u} d\mu(x)\right) = \int_{\mathbb{X}} p_n d\mu$$

- **quantile estimation**: given a random input $X \in \mathbb{X}$ and $\alpha \in [0, 1]$ estimate

$$q_\alpha = \inf \{z \in \mathbb{R}; P(f(X) \leq z) \geq \alpha\},$$

- Given a **stochastic simulator** $x \mapsto P_x = \mathcal{N}(\xi(x), \sigma^2)$, estimate $\alpha = \int_{\mathbb{X}} P_x(\cdot]u, \infty]) d\mu(x)$ (Stroh et al. 2017)

- Uncertainty measures?
- Assuming f is modeled using a random process ξ , V. and Piera-Martinez 2007, Bect et al 2012, Chevalier et al. 2014:
 - $H_n^{(1)} = \text{var}_n(\alpha) = \mathbf{E}_n((\alpha - \hat{\alpha}_n)^2)$
 - $H_n^{(2)} = \int_{\mathbf{X}} \mathbf{E}_n [(\mathbb{1}_{\xi > u} - \mathbb{1}_{\hat{\xi}_n > u})^2]^{1/2} d\mu$
 - $H_n^{(3)} = \mathbf{E}_n(\mu(\Gamma \Delta \hat{\Gamma}_n)) = \mathbf{E}_n\left(\int_{\mathbf{X}} (\mathbb{1}_{\xi > u} - \mathbb{1}_{\hat{\xi}_n > u})^2 d\mu\right) = \int_{\mathbf{X}} p_n(1 - p_n) d\mu$
 - ...
- Note that $(H_n^{(1)})^{1/2} \leq H_n^{(2)} \leq (H_n^{(3)})^{1/2}$ by Minkowski and Cauchy-Schwarz inequalities
- If flagging “unsafe” regions as “safe” is bad $\rightsquigarrow H_n^{(4)} = \mathbf{E}(\mu(\Gamma \setminus \hat{\Gamma}_n))$ (Azzimonti et al. 2018)

4 Convergence results

- V. & Bect 2011 prove the consistency of the EI algorithm using two properties:

1. $\liminf_{n \rightarrow \infty} \sup_{\mathbb{X}} G_n = 0$

2. $\inf H_n > 0 \implies \liminf_{n \rightarrow \infty} \sup_{\mathbb{X}} G_n > 0$

- Bect et al. 2018 improve the result and show the consistency of several other SUR strategies. In particular, if $\forall x \in \mathbb{X}$

$$J_n(x) = \mathbf{E}_n(H_{n+1} \mid X_{n+1} = x) \leq H_n$$

then $\sup_{\mathbb{X}} G_n \rightarrow 0$ a.s.

Rates?

- Almost unknown!
- For the optimization problem, Bull (2011) constructs an **upper-bound of the convergence rate of the expected improvement strategy**:

$$\sup_{\|f\|_{\mathcal{R}} \leq 1} M - M_n = O(n^{-(\nu \wedge 1)/d} (\log n)^\beta)$$

- For the optimization of analytic functions using expected improvement, Yarotsky (2012) prove exponential convergence rates
- Open question: could we recover the optimal convergence rate of approximation for SUR strategies?

- In more details, let $\xi \sim \text{GP}(0, k)$ and assume there exists $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}$ with Fourier transform $\tilde{\Phi}$ satisfying

$$c_1(1 + \|u\|_2^2)^{-\nu-d/2} \leq \tilde{\Phi}(u) \leq c_2(1 + \|u\|_2^2)^{-\nu-d/2}, \quad u \in \mathbb{R}^d,$$

such that $k(x, y) = \Phi(x - y)$ ($\nu > 0, 0 < c_1 < c_2$)

- Then, Ritter (2010) shows that $e_n = \sup_{\mathbb{X}} \sigma_n(x) \geq Cn^{-2\nu/d}$
- Moreover, V. & Bect (2011) show that the classical sequential non-adaptive strategy

$$X_{n+1} = \arg \max_{\mathbb{X}} \sigma_n(x)$$

achieves $\sup_{\mathbb{X}} \sigma_n(x) = O(n^{2\nu/d})$

5 Concluding remarks

Computational cost

- SUR sampling criteria mostly often based on the computation of many posterior distributions \rightsquigarrow non negligible computational additional cost
- Alleviate computational costs:
 - Bettinger et al. 2009 \rightsquigarrow Tsallis entropy to maximize response diversity
 - Chevalier et al. 2014, Stroh et al. 2017 \rightsquigarrow fast formulaes for uncertainty reduction about excursion sets
 - Hernandez-Lobato et al. 2014 \rightsquigarrow predictive entropy search for optimization
 - Couckuyt et al. 2014, Hupkens et al. 2015, Zhao et al. 2019: fast computation of the EHVI

- Labopin-Richard & Picheny, 2017 \rightsquigarrow quantile estimation
- ...
- Batch evaluations: Ginsbourger et al. 2010, Chevalier et al. 2014, Dutrieux et al. 2015...

About priors

- Non stationary models?
 - Gaussian trees (Gramacy 2007)
 - Warped GP (Snelson et al 2004)...
 - Deep GP (Snoek et al. 2015, Hebbal et al. 2019...)
 - ...

Fully Bayesian vs empirical Bayes vs alternatives

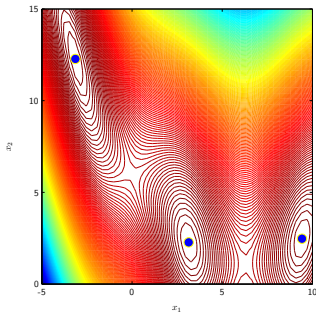
↪ use a prior

$$\xi : \begin{cases} \xi \mid \beta_1, \dots, \beta_q, \theta \sim \text{GP}\left(\sum_{i=1}^q \beta_i p_i, k_\theta\right), \\ \beta_1, \dots, \beta_q \sim \text{N}(0, \infty), \\ \theta \sim \pi_0, \end{cases}$$

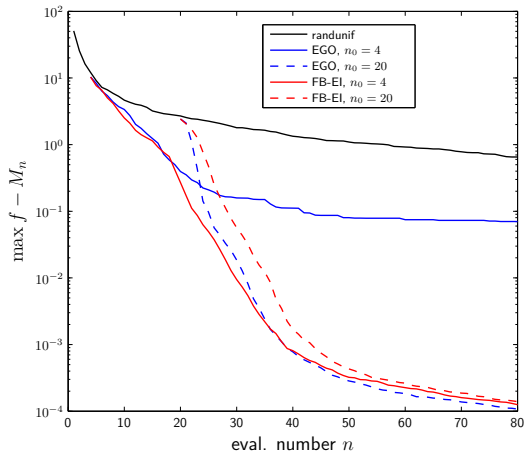
or **estimate the prior** from data inside a family of prior processes $\{\xi_\theta; \theta \in \Theta\}$ of the type

$$\xi_\theta : \begin{cases} \xi_\theta \mid \beta_1, \dots, \beta_q \sim \text{GP}\left(\sum_{i=1}^q \beta_i p_i, k_\theta\right), \\ \beta_1, \dots, \beta_q \sim \text{N}(0, \infty) \end{cases} \quad ?$$

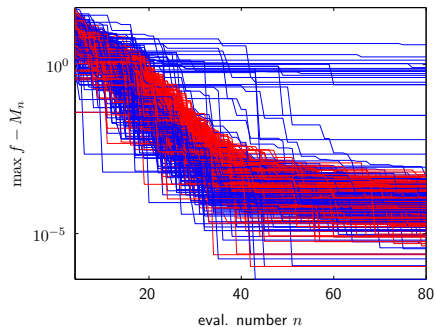
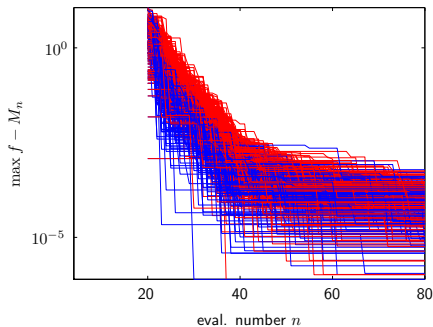
Possible shortcoming of the empirical Bayes approach (Benassi et al. 2011)



Branin function (three local maxima)



Average error (200 optim.) with random init. designs

**(a)** $n_0 = 4$ **(b)** $n_0 = 20$

Errors on each optimization run
(EGO in blue, FB-EI in red)

- Conclusion:
 - use a plugin method with a sufficiently large initial design
 - or a fully Bayesian approach
 - or maybe “modify” the kriging variance (Harville & Jeske, 1992, Abt 1999, Zimmerman 2006, Muller et al. 2010, Pronzato & Muller 2012. . .)

Future work?

- Research on SUR-like methods is under active development in the DACE/UQ community and also in machine learning
- Development of uncertainty measures for specific tasks (robust optimization, multi-fidelity. . .)
- Development of “good” models
- Better understanding of the properties of SUR strategies

GDR MASCOT-NUM 2019

Brief overview on
Stepwise Uncertainty Reduction for computer experiments

Emmanuel Vazquez

March 19, 2010