

Gaussian processes and bayesian optimization

Stanisław Jastrzębski

🏠 kudkudak.github.io 🐙 kudkudak

Plan

Goal: talk about modern hyperparameter optimization algorithms

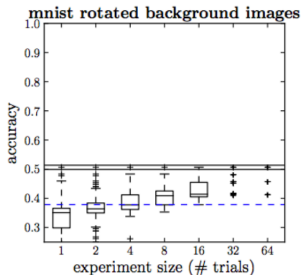
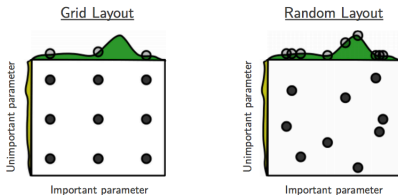
- Bayes reminder: equivalent linear regression formulations
- Gaussian Processes
- Hyperparameter optimization via Gaussian Processes (+code and simple experiment)

What do we want to achieve

- Select best hyperparameters for our model
- We will model directly $f(x)$, where x are our hyperparameters
- We will ask question like “how probably it is that accuracy is $> 60\%$ for $C=100$ ”
- We will need accuracy and variance estimate for each point

Why care? We have grid search

[Bergstra & Bengio, 2012] show that random search is **usually better**!



Linear Regression :)

Linear regression [Bishop, ch.3]

- Assume model $y(x) = \theta^T x$
- Optimal θ under MSE error $J(\theta) = \sum_i (y_n - \theta^T x_n)^2$ is

Definition (Least Squares estimate)

$$\hat{\theta} = (X^T X)^{-1} X^T y$$

Ridge regression [Bishop, ch.3]

Definition (Ridge regression cost)

$$J(\theta) = \sum_i (y_n - \theta^T x_n)^2 + \lambda \|w\|^2$$

Definition (Ridge regression estimate)

$$\hat{\theta} = (\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Kernels [Bishop, ch.3]

Just replace \mathbf{x} with $\phi_i(\mathbf{x})$, each ϕ_i is a new feature. Examples:

- $\phi_i(\mathbf{x}) = x_i$
- $\phi_i(\mathbf{x}) = \exp\left\{\frac{-(x-\mu_j)^2}{2s^2}\right\}$
- $\phi_i(\mathbf{x}) = x_i^2$

Kernels [Bishop, ch.3]

Just replace \mathbf{x} with $\phi_i(\mathbf{x})$, each ϕ_i is a new feature. Examples:

- $\phi_i(\mathbf{x}) = x_i$
- $\phi_i(\mathbf{x}) = \exp\left\{\frac{-(x-\mu_j)^2}{2s^2}\right\}$
- $\phi_i(\mathbf{x}) = x_i^2$

$$y(x) = \theta_1 \phi_1(x) + \theta_2 \phi_2(x) + \theta_3 \phi_3(x)$$

Kernels [Bishop, ch.3]

Just replace \mathbf{x} with $\phi_i(\mathbf{x})$, each ϕ_i is a new feature. Examples:

- $\phi_i(\mathbf{x}) = x_i$
- $\phi_i(\mathbf{x}) = \exp\left\{\frac{-(x-\mu_j)^2}{2s^2}\right\}$
- $\phi_i(\mathbf{x}) = x_i^2$

$$y(x) = \theta_1 \phi_1(x) + \theta_2 \phi_2(x) + \theta_3 \phi_3(x)$$

$$y(x) = 1.2 \exp\{-(x-1)^2\} + 2.3 \exp\{-(x-2)^2\} + 0.3 \exp\{-(x-3)^2\}$$

Kernels [Bishop, ch.3]

$$\Phi_{ij} = \phi_i(x_j)$$

Definition (Kernelized ridge regression estimate)

$$\hat{\theta} = (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$$

Probabilistic formulation [Bishop, ch.3]

- Standard probabilistic model $y = \theta^T x + \nu$, where $\nu \sim \mathcal{N}(0, \sigma^2)$
- Wow! It follows $y - \theta^T x \sim \mathcal{N}(0, \sigma^2)$, or equivalently

$$P(y|x, \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y - \theta^T x)^2}{2\sigma^2}\right)$$

Probabilistic formulation [Bishop, ch.3]

Everything will be a Gauss from now on :)

Definition (Least-squares likelihood function)

$$p(y|\mathbf{x}, \boldsymbol{\theta}, \sigma) = \mathcal{N}(\boldsymbol{\theta}^T \mathbf{x}, \sigma^2)$$

[Derive on whiteboard MLE estimate]

Probabilistic formulation [Bishop, ch.3]

Everything will be a Gauss from now on :)

Definition (Least-squares likelihood function)

$$p(y|\mathbf{x}, \boldsymbol{\theta}, \sigma) = \mathcal{N}(\boldsymbol{\theta}^T \mathbf{x}, \sigma^2)$$

[Derive on whiteboard MLE estimate]

Definition (MLE estimate)

$$\hat{\boldsymbol{\theta}} = \operatorname{argmax}_{\boldsymbol{\theta}} p(\mathbf{t}|\mathbf{x}, \boldsymbol{\theta}, \sigma) = (\boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^T \mathbf{y}$$

Bayesian approach

Bayesian approach is all about priors. Every analysis makes prior assumption, but here we state them in an explicit manner.

Bayesian approach

Bayesian approach is all about priors. Every analysis makes prior assumption, but here we state them in an explicit manner. In general we have 3 things in a bayesian model (everything can be formulated as a bayes model, even SVM):

- likelihood $p(\mathbf{y}|x, \theta)$
- prior $p(\theta)$
- postertior $p(\theta|\mathbf{X}, \mathbf{y})$

Bayesian linear regression

Goal: formulate linear regression with explicit prior on weights

- likelihood $p(y|x, \boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}^T \phi(x), \beta^{-1})$

Bayesian linear regression

Goal: formulate linear regression with explicit prior on weights

- likelihood $p(y|x, \boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}^T \phi(x), \beta^{-1})$
- prior $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}|\boldsymbol{\theta}_0, \boldsymbol{\Sigma}_0)$

Bayesian linear regression

Goal: formulate linear regression with explicit prior on weights

- likelihood $p(y|x, \theta) = \mathcal{N}(\theta^T \phi(x), \beta^{-1})$
- prior $p(\theta) = \mathcal{N}(\theta|\theta_0, \Sigma_0)$
- posterior $p(\theta|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{X}, \theta)p(\theta)}{p(\mathbf{X})} = ?$

Gaussian ju-jitsu

Gaussian are “closed” with respect to many operations

Theorem (Marginal and Conditional Gaussian)

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$$
$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\mathbf{x}, \mathbf{L}^{-1})$$

Then

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\mu}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^T)$$
$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\Sigma}\{\mathbf{A}^T\mathbf{L}\mathbf{y} + \boldsymbol{\Lambda}\boldsymbol{\mu}\}, \boldsymbol{\Sigma})$$

, where $\boldsymbol{\Sigma} = (\boldsymbol{\Lambda} + \mathbf{A}^T\mathbf{A})^{-1}$

Bayesian linear regression

Goal: formulate linear regression with explicit prior on weights

- likelihood $p(y|x, \theta) = \mathcal{N}(\theta^T \phi(x), \beta^{-1})$
- prior $p(\theta) = \mathcal{N}(\theta|\theta_0, \Sigma_0)$
- posterior $p(\theta|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{X}, \theta)p(\theta)}{p(\mathbf{X})} = \mathcal{N}(\theta|\theta_N, \Sigma_N)$

, where

$$\begin{aligned}\theta_N &= \Sigma_N(\Sigma_0^{-1}\theta_0 + \beta\Phi^T\mathbf{y}) \\ \Sigma_N^{-1} &= \Sigma_0^{-1} + \beta\Phi^T\Phi\end{aligned}$$

Let's assume $\Sigma_0 = \alpha^{-1}I$. **What happens when $\alpha \rightarrow 0$?** . Infinitely broad prior gives 0 regularization.

Bayesian linear regression

If $\alpha \neq 0$:

$$\ln p(\boldsymbol{\theta}|\mathbf{y}) = -\frac{\beta}{2} \sum_{n=1}^N \{y_n - \boldsymbol{\theta}^T \phi(x_n)\} - \frac{\alpha}{2} \boldsymbol{\theta}^T \boldsymbol{\theta} + \text{const} \quad (1)$$

Bayesian linear regression

If $\alpha \neq 0$:

$$\ln p(\boldsymbol{\theta}|\mathbf{y}) = -\frac{\beta}{2} \sum_{n=1}^N \{y_n - \boldsymbol{\theta}^T \phi(x_n)\} - \frac{\alpha}{2} \boldsymbol{\theta}^T \boldsymbol{\theta} + \text{const} \quad (1)$$

$\ln p(\boldsymbol{\theta}|\mathbf{y})$ (MAP - maximum a posteriori estimate) *is same as ridge regression* objective. Magic?

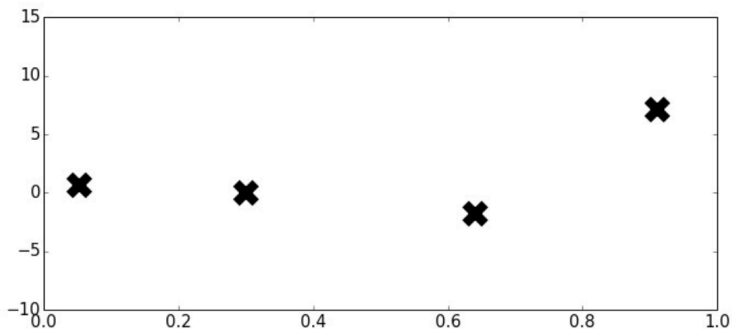
Bayes a'int magic

Putting up your Bayes Hat is just about making your assumptions clear. In the end it is just different road to the same goal (objective).



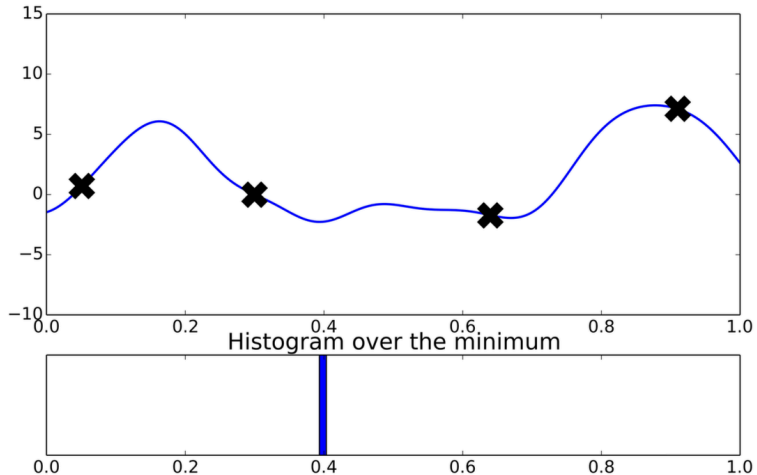
Gaussian Processes

Intuition - we want to build this somehow



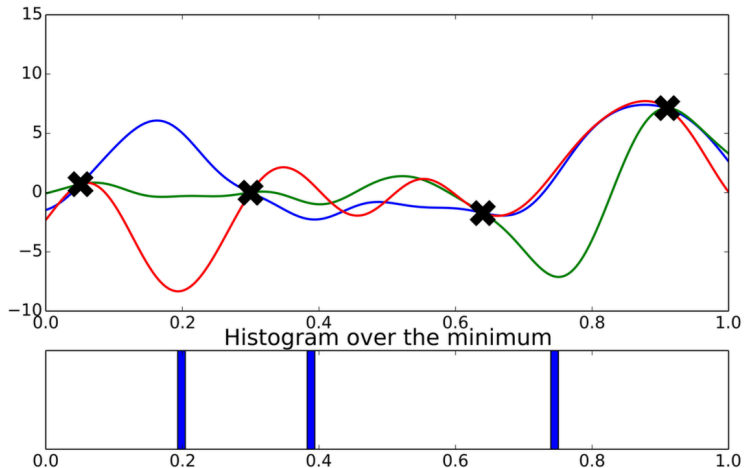
Where is the minimum of f ?
Where should the take the next evaluation?

Intuition - we want to build this somehow



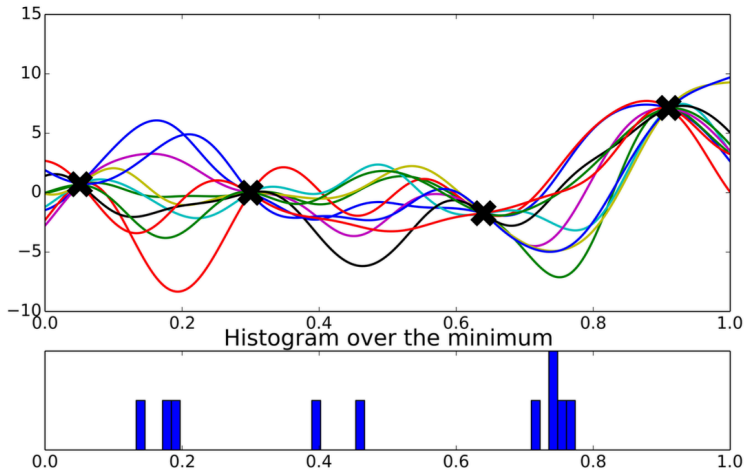
We want to build this

Intuition - we want to build this somehow



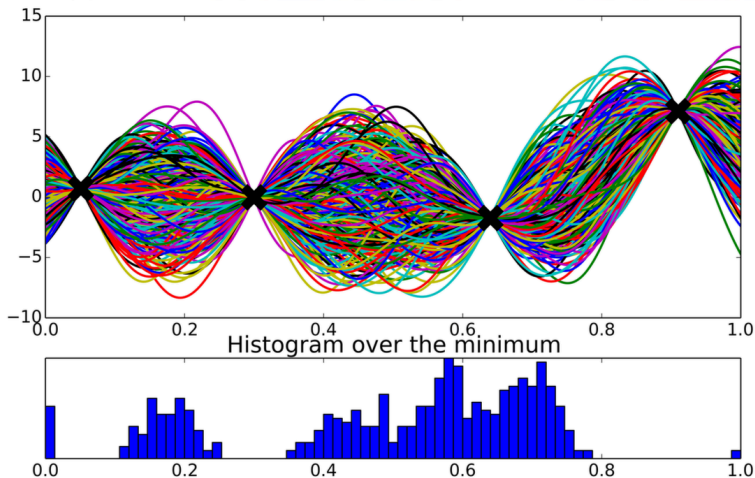
We want to build this

Intuition - we want to build this somehow



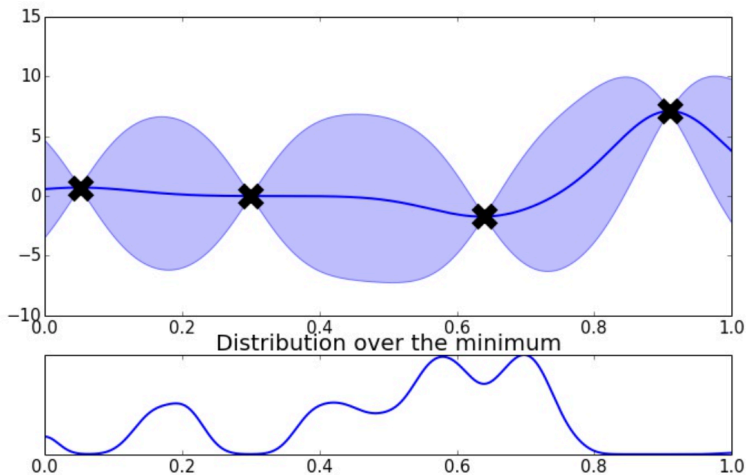
We want to build this

Intuition - we want to build this somehow



We want to build this

Intuition - we want to build this somehow



We want to build this

Gauss family provides inference \Leftrightarrow linear algebra bridge

$$p(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{n/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

$$\mathbf{x} = \begin{bmatrix} x_A \\ x_B \end{bmatrix} \quad \boldsymbol{\mu} = \begin{bmatrix} \mu_A \\ \mu_B \end{bmatrix} \quad \boldsymbol{\Sigma} = \begin{bmatrix} \Sigma_{AA} & \Sigma_{AB} \\ \Sigma_{BA} & \Sigma_{BB} \end{bmatrix}.$$

1. **Normalization.** The density function normalizes, i.e.,

$$\int_{\mathbf{x}} p(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x} = 1.$$

2. **Marginalization.** The marginal densities,

$$p(x_A) = \int_{x_B} p(x_A, x_B; \boldsymbol{\mu}, \boldsymbol{\Sigma}) dx_B$$

$$p(x_B) = \int_{x_A} p(x_A, x_B; \boldsymbol{\mu}, \boldsymbol{\Sigma}) dx_A$$

$$x_A \sim \mathcal{N}(\mu_A, \Sigma_{AA})$$

$$x_B \sim \mathcal{N}(\mu_B, \Sigma_{BB}).$$

Gauss family provides inference \Leftrightarrow linear algebra bridge

3. **Conditioning.** The conditional densities

$$p(x_A | x_B) = \frac{p(x_A, x_B; \mu, \Sigma)}{\int_{x_A} p(x_A, x_B; \mu, \Sigma) dx_A}$$
$$p(x_B | x_A) = \frac{p(x_A, x_B; \mu, \Sigma)}{\int_{x_B} p(x_A, x_B; \mu, \Sigma) dx_B}$$

are also Gaussian:

$$x_A | x_B \sim \mathcal{N}(\mu_A + \Sigma_{AB}\Sigma_{BB}^{-1}(x_B - \mu_B), \Sigma_{AA} - \Sigma_{AB}\Sigma_{BB}^{-1}\Sigma_{BA})$$
$$x_B | x_A \sim \mathcal{N}(\mu_B + \Sigma_{BA}\Sigma_{AA}^{-1}(x_A - \mu_A), \Sigma_{BB} - \Sigma_{BA}\Sigma_{AA}^{-1}\Sigma_{AB}).$$

Gaussian process - formally

$f(x_0) = 2.5, f(x_1) = 3 \dots f(x_m) = 4$ Let \mathcal{H} be a set of all functions over m points (it is finite). Let's exploit correspondence between function $f(\cdot) \in \mathcal{H}$

$$p(\vec{f}) = \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(f(x_i) - \mu_i)^2\right) \quad (\text{The only important eq.})$$

Gaussian process - formally

$f(x_0) = 2.5, f(x_1) = 3 \dots f(x_m) = 4$ Let \mathcal{H} be a set of all functions over m points (it is finite). Let's exploit correspondence between function $f(\cdot) \in \mathcal{H}$

$$p(\vec{f}) = \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(f(x_i) - \mu_i)^2\right) \quad (\text{The only important eq.})$$

Gaussian Processes - more formally

A **Gaussian process** is a stochastic process such that any finite subcollection of R.V has a multivariate Gaussian distribution

$$\begin{bmatrix} f(x_1) \\ \vdots \\ f(x_m) \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} f(m(x_1)) \\ \vdots \\ f(m(x_m)) \end{bmatrix}, \begin{bmatrix} k(x_1, x_1) & \dots & k(x_1, x_m) \\ \vdots & \ddots & \dots \\ k(x_m, x_1) & \dots & k(x_m, x_m) \end{bmatrix} \right) \quad (2)$$

Theorem (Mean and covariance of point)

$$\begin{aligned} m(x) &= E[x] \\ k(x, x') &= E[(x - m(x))(x' - m(x'))] \end{aligned}$$

Kernel matrix

Kernels are very powerful, a kernel function could do complicated parsing during sentence comparison, or perform graph similarity when comparing compounds, and all of this is compatible with GP.

Techniques for Constructing New Kernels.

Given valid kernels $k_1(\mathbf{x}, \mathbf{x}')$ and $k_2(\mathbf{x}, \mathbf{x}')$, the following new kernels will also be valid:

$$k(\mathbf{x}, \mathbf{x}') = ck_1(\mathbf{x}, \mathbf{x}') \quad (6.13)$$

$$k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}') \quad (6.14)$$

$$k(\mathbf{x}, \mathbf{x}') = g(k_1(\mathbf{x}, \mathbf{x}')) \quad (6.15)$$

$$k(\mathbf{x}, \mathbf{x}') = \exp(k_1(\mathbf{x}, \mathbf{x}')) \quad (6.16)$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}') \quad (6.17)$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}') \quad (6.18)$$

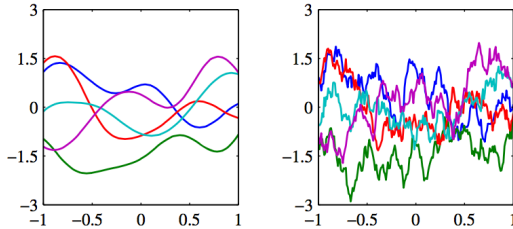
$$k(\mathbf{x}, \mathbf{x}') = k_3(\phi(\mathbf{x}), \phi(\mathbf{x}')) \quad (6.19)$$

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{A} \mathbf{x}' \quad (6.20)$$

$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a) + k_b(\mathbf{x}_b, \mathbf{x}'_b) \quad (6.21)$$

$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a)k_b(\mathbf{x}_b, \mathbf{x}'_b) \quad (6.22)$$

Kernel matrix - intuition



[Code presentation]

Gaussian Process regression

Definition (GP regression)

$$y^{(i)} = f(x^{(i)}) + \epsilon^{(i)}, \text{ where } f(\cdot) \sim \mathcal{GP}(0, k(\cdot, \cdot))$$

Same as before, just not $\theta^T x$ but $f(x)$

Inference in GP - just use Gaussian ju-jitsu

We have trained on \mathbf{X}, \mathbf{y} and want to make predictions on $\mathbf{X}_*, \mathbf{y}_*$ (to calculate new hyperparameters, for instance)

$$\begin{bmatrix} f \\ \vec{f} \\ f^* \end{bmatrix} | \mathbf{X}, \mathbf{X}_* \sim \mathcal{N}\left(\vec{0}, \begin{bmatrix} k(\mathbf{X}, \mathbf{X}) & k(\mathbf{X}, \mathbf{X}_*) \\ k(\mathbf{X}_*, \mathbf{X}) & k(\mathbf{X}_*, \mathbf{X}_*) \end{bmatrix}\right) \quad (3)$$

Add noise and we find distribution of y^*

Inference in GP - just use Gaussian ju-jitsu

We have trained on \mathbf{X}, \mathbf{y} and want to make predictions on $\mathbf{X}_*, \mathbf{y}_*$ (to calculate new hyperparameters, for instance)

$$\begin{bmatrix} f \\ \vec{f} \\ f^* \end{bmatrix} | \mathbf{X}, \mathbf{X}_* \sim \mathcal{N}\left(\vec{0}, \begin{bmatrix} k(\mathbf{X}, \mathbf{X}) & k(\mathbf{X}, \mathbf{X}_*) \\ k(\mathbf{X}_*, \mathbf{X}) & k(\mathbf{X}_*, \mathbf{X}_*) \end{bmatrix}\right) \quad (3)$$

Add noise and we find distribution of y^*

$$\begin{bmatrix} y \\ y^* \end{bmatrix} | \mathbf{X}, \mathbf{X}_* \sim \begin{bmatrix} \vec{f} \\ \vec{f}^* \end{bmatrix} + \begin{bmatrix} \epsilon \\ \epsilon^* \end{bmatrix} \sim \mathcal{N}\left(\vec{0}, \begin{bmatrix} k(\mathbf{X}, \mathbf{X}) + \sigma^2 I & k(\mathbf{X}, \mathbf{X}_*) \\ k(\mathbf{X}_*, \mathbf{X}) & k(\mathbf{X}_*, \mathbf{X}_*) + \sigma^2 I \end{bmatrix}\right) \quad (4)$$

Inference in GP - just use Gaussian ju-jitsu #2

Now we want:

$$\mathbf{y}_* | \mathbf{y} = ? \quad (5)$$

Inference in GP - just use Gaussian ju-jitsu #2

Now we want:

$$\mathbf{y}_* | \mathbf{y} = ? \quad (5)$$

If we substitute into previous equations we will get:

$$\mathbf{y}_* | \mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}_*, \boldsymbol{\Sigma}_*) \quad (6)$$

, where

$$\boldsymbol{\mu}_* = K(\mathbf{X}_*, \mathbf{X})(K(\mathbf{X}, \mathbf{X}) + \sigma^2 I)^{-1} \mathbf{y}$$

$$\boldsymbol{\Sigma}_* = K(\mathbf{X}_*, \mathbf{X}_*) + \sigma^2 I - K(\mathbf{X}_*, \mathbf{X})(K(\mathbf{X}, \mathbf{X}) + \sigma^2 I)^{-1} K(\mathbf{X}, \mathbf{X}_*)$$

Inference in GP - just use Gaussian ju-jitsu #2

Now we want:

$$\mathbf{y}_* | \mathbf{y} = ? \quad (5)$$

If we substitute into previous equations we will get:

$$\mathbf{y}_* | \mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}_*, \boldsymbol{\Sigma}_*) \quad (6)$$

, where

$$\boldsymbol{\mu}_* = K(\mathbf{X}_*, \mathbf{X})(K(\mathbf{X}, \mathbf{X}) + \sigma^2 I)^{-1} \mathbf{y}$$

$$\boldsymbol{\Sigma}_* = K(\mathbf{X}_*, \mathbf{X}_*) + \sigma^2 I - K(\mathbf{X}_*, \mathbf{X})(K(\mathbf{X}, \mathbf{X}) + \sigma^2 I)^{-1} K(\mathbf{X}, \mathbf{X}_*)$$

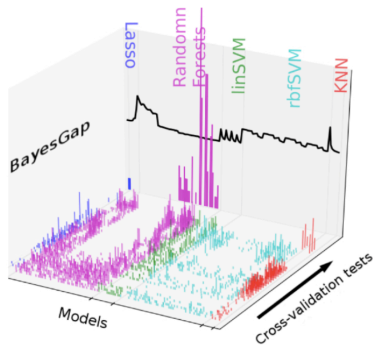
Note: in practice we fit $\boldsymbol{\theta}_{MLE} = \operatorname{argmax}_{\boldsymbol{\theta}} p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta})$ (for instance kernel width)

[Show me some code]

Bayesian Optimization

Taking human out of the loop. Applications

- Picking parameters in different solvers (CPLEX has 72 parameters)
- Optimizing hyperparameters in ML algoirhtms
- A/B testing of website
- ...



Bayesian optimization algorithm

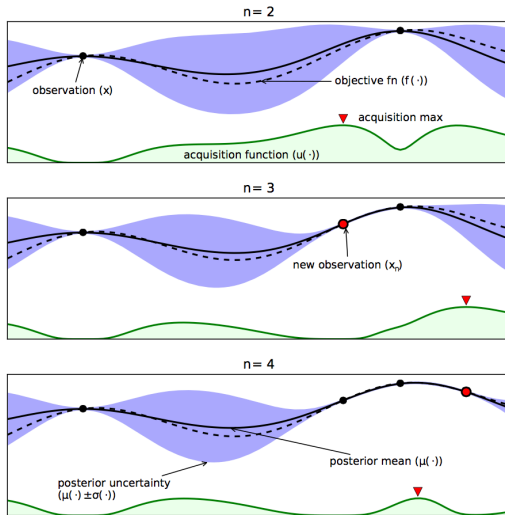
Algorithm 1 Bayesian optimization

- 1: **for** $n = 1, 2, \dots$ **do**
- 2: select new \mathbf{x}_{n+1} by optimizing acquisition function α

$$\mathbf{x}_{n+1} = \arg \max_{\mathbf{x}} \alpha(\mathbf{x}; \mathcal{D}_n)$$

- 3: query objective function to obtain y_{n+1}
 - 4: augment data $\mathcal{D}_{n+1} = \{\mathcal{D}_n, (\mathbf{x}_{n+1}, y_{n+1})\}$
 - 5: update statistical model
 - 6: **end for**
-

Bayesian optimization algorithm



Acquisition functions

Having full probabilistic model you are only limited by your creativity as to how to pick next point.

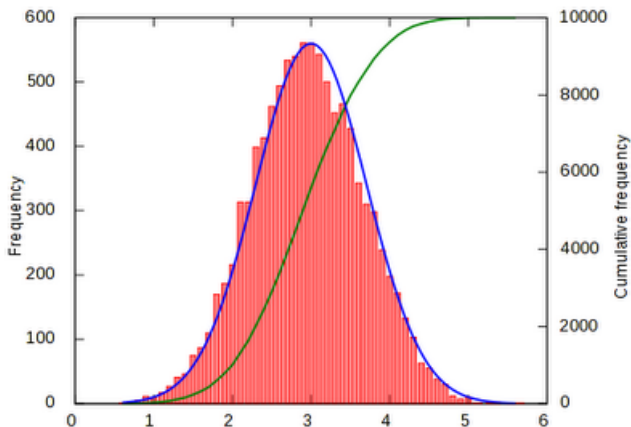
$$\alpha(x; \mathcal{D}_n) = \mathbb{E}_{v|x, \theta}[U(x, v, \theta)] \quad (7)$$

Let's assume we have observed some \mathbf{X} and \mathbf{y} .

$$\begin{aligned}\mu_n(\mathbf{x}) &= \mu_0(\mathbf{x}) + k(\mathbf{x})^T (K + \sigma^2 I)^{-1} (\mathbf{y} - m) \\ \sigma_n^2(\mathbf{x}) &= k(\mathbf{x})^T (K + \sigma^2 I)^{-1} k(\mathbf{x})\end{aligned}$$

Improvement-based

$$\alpha_{PI}(x; \mathcal{D}_n) = \Phi\left(\frac{\mu_n(x) - \tau}{\sigma_n(x)}\right) \quad (8)$$



Optimistic :)

$$\alpha_{UCB}(x; \mathcal{D}_n) = \mu_n(\mathbf{x}) + \beta \sigma_n^2(\mathbf{x}) \quad (9)$$

Information theoretical

$$\alpha_{ES}(x; \mathcal{D}_n) = H(\mathbf{x}^* | \mathcal{D}_n) - H_{y | \mathcal{D}_n, x}(\mathbf{x}^* | \mathcal{D}_n \cup \{(\mathbf{x}, \mathbf{y})\}) \quad (10)$$

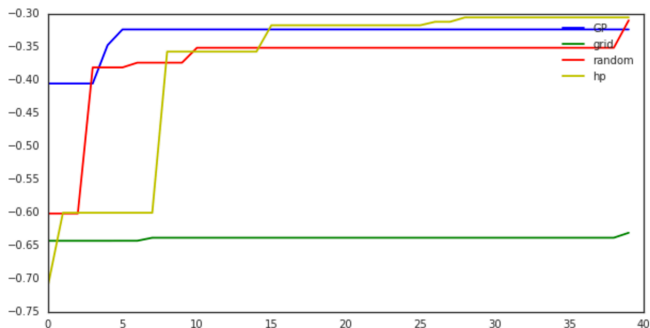
Note: we need here posterior distribution $p_*(\mathbf{x} | \mathcal{D}_n)$ which is very hard to compute.

Caveats and practical considerations of GPs in BO

- Hard to handle conditional hyperparameters (hyperopt uses tree based prob. estimation)
- $O(n^3)$ inference complexity (reducible using many tricks)
- Most BO algorithms are sequential
- No overall *best* acquisition function

Simple experiment on virtual screening

Fitting SVM RBF for 10 values of γ and 12 values of C to 5-HT2a with KRFP fingerprint (using *dr.GMUM*)



Thank you