Two Use-cases of Uncertainty: Active Learning and Bayesian Optimization

> CS772A: Probabilistic Machine Learning Piyush Rai

Active Learning



Passive Learning

- Standard supervised learning is passive
 - Learner has no control over what labelled training examples it gets to learn from



Supervised Passive Learner

to determine labels



Active Learning

It is therefore also a sequential learning strategy (training data is not given all at once)

- In Active Learning, the learner can request specific labelled examples as it trains
 - In particular, examples that the learner thinks will be most useful to learn the underlying function





Active Learning

The figure below is another illustration of AL





Measuring Usefulness in AL

Given the acquisition function, the most useful example can be selected from the pool as



- Various ways to measure the usefulness of an unlabeled example \boldsymbol{x}_* Note: We will use shorthand $A(\boldsymbol{x}_*)$
 - Defined by an "acquisition function" $A(x_*)$ (high value for most useful unlabeled examples)
- Approach 1: For x_* , look at uncertainty in output y_* predicted by the current model
 - Can use variance in the posterior predictive distribution: $A(x_*) = var(y_*)$
 - More generally, can use entropy of the PPD: $A(x_*) = \mathbb{H}[p(y_*|x_*, D)]$



Low entropy



new example in our training set

Note that this is the "marginal entropy" of the output distribution since the posterior predictive of the output is obtained by marginalizing over the posterior

CS772A: PML

0 input. x

• Approach 2: Look at how much our model will <u>improve</u> if we add this unlabeled example with its true label, to our training set, and <u>retrain</u> the model $A(x_*) = \mathbb{H}[p(\theta|\mathcal{D})] - \mathbb{E}_{p(y_*|x_*} \mathcal{D}) \mathbb{H}[p(\theta|\mathcal{D} \cup (x_*, y_*))] = \mathbb{I}[\theta; y_*|\mathcal{D}, x_*]$

here since y_* is not known

$$x_*) = \mathbb{H}[p(\theta|\mathcal{D})] - \mathbb{E}_{p(y_*|x_*,\mathcal{D})} \mathbb{H}[p(\theta|\mathcal{D} \cup (x_*, y_*))] = \mathbb{I}[\theta; y_*|\mathcal{D}, x_*]$$
Entropy of the current posterior Meed to use expectation Entropy of the new posterior after including the

Batch Active Learning

- Approaches we saw work by querying and adding one example at a time
- Expensive in practice since we have to retrain every time after including a new example
 - Especially true for deep learning models which are computationally expensive to train
- In practice, we want to use AL to jointly query the labels of B > 1 examples $(\hat{x}_1, \hat{x}_2, ..., \hat{x}_B) = \operatorname{argmax}_{(x_1, x_2, ..., x_B) \in X_{pool}} A(x_1, x_2, ..., x_B | p(\theta | D))$
- Difficult to construct such joint acquisition function and maximize them
- A greedy scheme is to simply select the B highest scoring points

$$A(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_B | p(\theta | \mathcal{D})) = \sum_{b=1}^B \mathbb{I}[\theta; y_b | \mathcal{D}, \mathbf{x}_b]$$

The above however is myopic and ignores correlations among the selected points

CS772A: PML

Some recent works have addresses this issue¹



Bayesian Optimization: The Basic Formulation

• Consider finding the optima x_* (say minima) of a function f(x)



- Caveat: We don't know the form of the function; can't get its gradient, Hessian, etc
- Can only query the function's values at certain points (i.e., only "black-box" access)
 - The values may or may not be noisy (i.e., we may be given f(x) or $f(x) + \epsilon$)

CS772A: PML

Bayesian Optimization: Some Applications

- Drug Design: Want to find the optimal chemical composition for a drug
 - Optimal composition will be the one that has the best efficacy
 - But we don't know the efficacy function
 - Can only know the efficacy via doing clinical trials
 - Each trial is expensive; can't do too many trials

Hyperparameter Optimization: Want to find the optimal hyperparameters for a model

- Optimal hyperparam values will be those that give the lowest test error
- Don't know the true "test error" function
- Need to train the model each time with different h.p. values and compute test error
- Training every time will be expensive (e.g., for deep nets)
- Note: Hyperparams here can even refer to the structure of a deep net (depth, width, etc)
- Many other applications: Website design via A/B testing, material design, optimizing physics based models (e.g., aircraft design), etc

CS772A: PML

- Can use BO to find maxima or minima
- Would like to locate the optima by querying the function's values (say, from an oracle)



- We would like to do so using as few queries as possible
- Reason: The function's evaluation may be time-consuming or costly



- Suppose we are allowed to make the queries sequentially
- This information will be available to us in form of query-function value pairs

By solving a

0.8

Queries so far can help us estimate the function

1.0

0.5

-0.5

-1.0

-1.5

(× 0.0

Dotted curve: True function Green curve: Current estimate ("surrogate") of the function Shaded region: Uncertainty in the function's estimate

- BO uses past queries + function's estimate+uncertainty to decide where to query next
- Similar to Active Learning but the goal is to learn f as well as finds its optima

Note: Function values can be noisy too, e.g.,

CS772A: PML

 $f(x_n) + \epsilon_n$

 $\{(x_n, f(x_n))\}_{n=1}^N$

- BO requires two ingredients
 - A regression model to learn a surrogate of f(x) given previous queries $\{(x_n, f(x_n))\}_{n=1}^N$
 - An acquisition function A(x) to tell us where to query next



- Note: The regression model must also have estimate of function's uncertainty
 - Bayesian nonlinear regression, such as GP, Bayesian Neural network, etc would be ideal



Note: Function values can be noisy too, e.g.,

 $f(x_n) + \epsilon_n$

Bayesian Optimization: An Illustration

• Suppose our goal is to find the maxima of f(x) using BO





Pic source: http://krasserm.github.io/2018/03/21/bayesian-optimization/

Some Basic Acquisition Functions for BO (assuming we are finding the minima)



Acquisition Functions: Probability of Improvement¹⁶

- Assume past queries $\mathcal{D}_N = (X, f) = (x_n, f(x_n))_{n=1}^N$ and suppose $f_{min} = \min f$
- Suppose f_{new} denotes the function's value at the next query point x_{new}
- We have an improvement if $f_{new} < f_{min}$ (recall we are doing minimization)
- Assuming the function is real-valued, suppose the posterior predictive for x_{new} is $p(f_{new}|x_{new}, \mathcal{D}_N) = \mathcal{N}(f_{new}|\mu(x_{new}), \sigma^2(x_{new}))$
- We can define a probability of improvement based acquisition function

$$A_{PI}(x_{new}) = p(f_{new} \le f_{min}) = \int_{-\infty}^{f_{min}} \mathcal{N}(f_{new} | \mu(x_{new}), \sigma^2(x_{new})) df_{new} = \Phi\left(\frac{f_{min} - \mu(x_{new})}{\sigma(x_{new})}\right)$$

Exercise: Verify

CS772A: PML

of $\mathcal{N}(0,1)$

• The optimal query point will be one that maximizes $A_{PI}(x_{new})$

$$x_* = \operatorname{argmax}_{x_{new}} A_{PI}(x_{new})$$

Acquisition Functions: Expected Improvement

- PI doesn't take into account the amount of improvement
- Expected Improvement (EI) takes this into account and is defined as

$$A_{EI}(x_{new}) = \mathbb{E}[f_{min} - f_{new}] = \int_{-\infty}^{f_{min}} (f_{min} - f_{new}) \mathcal{N}(f_{new} | \mu(x_{new}), \sigma^2(x_{new})) df_{new}$$

Exercise: Prove
this result = $(f_{min} - \mu(x_{new})) \Phi\left(\frac{f_{min} - \mu(x_{new})}{\sigma(x_{new})}\right) + \sigma(x_{new}) \mathcal{N}\left(\frac{f_{min} - \mu(x_{new})}{\sigma(x_{new})}; 0, 1\right)$

• The optimal query point will be one that maximizes $A_{EI}(x_{new})$

$$x_* = \operatorname{argmax}_{x_{new}} A_{EI}(x_{new})$$

Focus on points where the function has high uncertainty (so that including them improves our estimate of the function)

CS772A: PML

Focus on points where

the function has small

looking for its minima)

values (since we are

17

- Note that the above acquisition function trades off exploitation vs exploration
 - Will prefer points with small predictive mean $\mu(x_{new})$: Exploitation
 - Will prefer points with large predictive variance $\sigma(x_{new})$: Exploration

Acquisition Functions: Lower Confidence Bound

- Lower Confidence Bound (LCB) also takes into account exploitation vs exploration
- Used when the regression model is a Gaussian Process (GP)
- Assume the posterior predictive for a new point to be

 $p(f_{new}|x_{new}, \mathcal{D}_N) = \mathcal{N}(f_{new}|\mu(x_{new}), \sigma^2(x_{new}))$

The LCB based acquisition function is defined as

 $A_{LCB}(x_{new}) = \mu(x_{new}) - \kappa \,\sigma(x_{new}) \checkmark$

Point with the smallest LCB is selected as the next query point

$$x_* = \operatorname{argmin}_{x_{new}} A_{LCB}(x_{new})$$

- κ is a parameter to trade-off exploitation (low mean) and exploration (high variance)
- Under certain conditions, the iterative application of this acquisition function will converge to the true global optima of f (Srinivas et al. 2010)





CS772A: PML

18

Thus prefer points at which the function has low mean but high variance

When using BO for maximization, we use

Bayesian Optimization: The Overall Algo

- Initialize $\mathcal{D} = \{\}$
- For n = 1, 2, ..., N (or until the budget doesn't exhaust)
 - Select the next query point x_n by optimizing the acquisition function

$$x_n = \operatorname{argopt}_x A(x)$$

• Get function's value from the black-box oracle: $f_n = f(x_n)$

• $\mathcal{D} = \{\mathcal{D} \cup (x_n, f_n)\}$ Can get the function's minima from this set of function's values

- Update the regression model for f using data ${\mathcal D}$



BO: Some Challenges/Open Problems

- Learning the regression model for the function
 - GPs are flexible but can be expensive as N grows
 - Bayesian neural networks can be a more efficient alternative to GPs (Snoek et al, 2015)
 - Hyperparams of the regression model itself (e.g., GP cov. function, Bayesian NN hyperparam)
- High-dimensional Bayesian Optimization (optimizing functions of many variables)
 - Most existing methods work well only for a moderate-dimensional x
 - Number of function evaluations required would be quite large in high dimensions
 - Lot of recent work on this (e.g., based on dimensionality reduction)
- Multitask Bayesian Optimization (joint BO for several related functions)
 - Basic idea: If two functions are similar their optima would also be nearby

