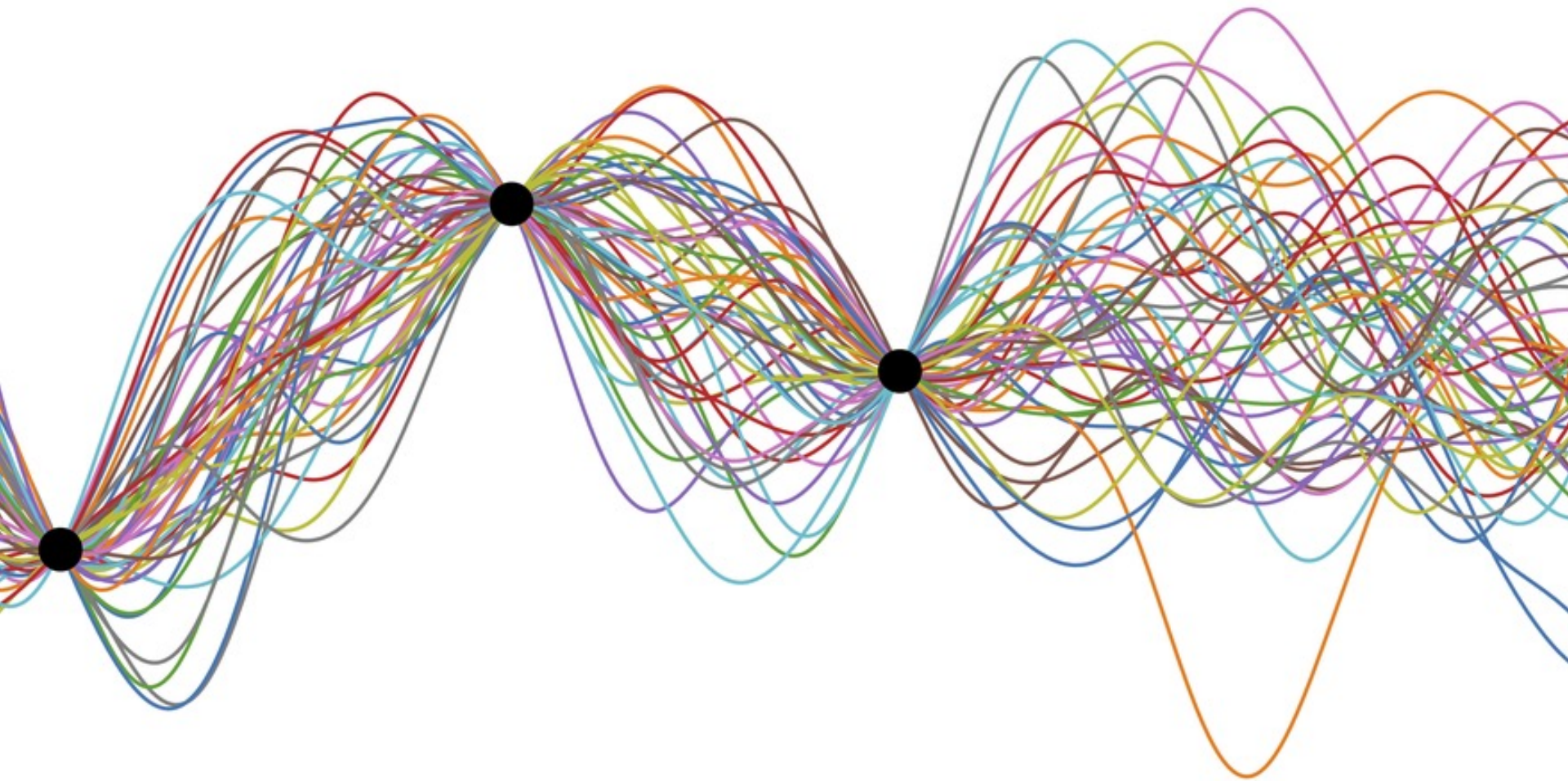


# Lecture 23: Gaussian Processes



# Goals for Today

## **Gaussian (normal) Distribution**

- 1D & multivariate
- Conditional distribution (Bayes theorem)

## **Gaussian Processes for Regression**

- Noise-free regression
- Regression with noise

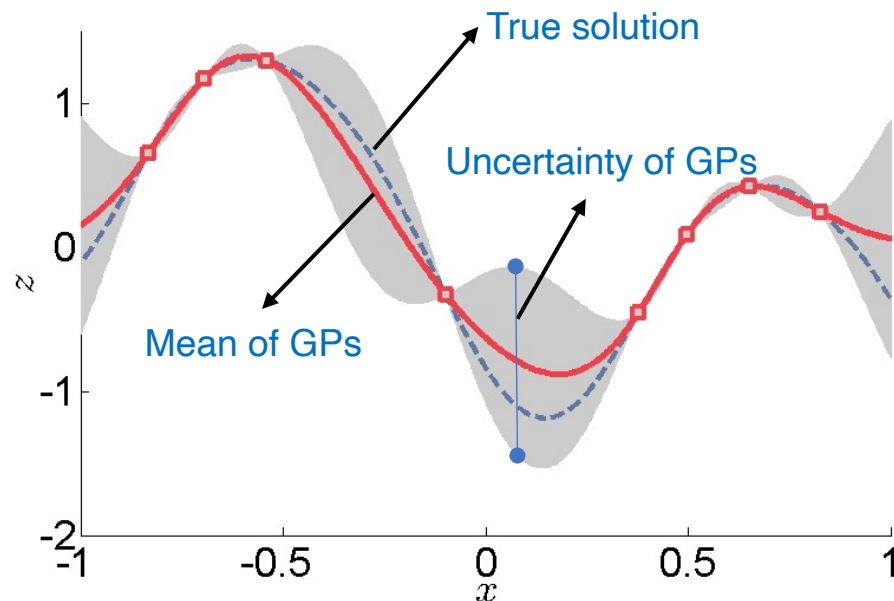
## **Gaussian Processes & Active Learning**

## **Extended Topics on Gaussian Processes**

- Kernel parameters estimation
- Relationship to other regression methods

# Gaussian Processes - Introduction

- In supervised learning, given some inputs  $\mathbf{x}_i$  and outputs  $\mathbf{y}_i$ , we assume  $\mathbf{y}_i = \mathbf{f}(\mathbf{x}_i)$ . The optimal approach is to infer a distribution over the functions given the data, which is called Gaussian processes (GPs) developed by Rasmussen (2004).
- A GP defines a prior over functions, which can be converted into a posterior over functions once we have seen some data. (Murphy 2012).
- A GP assumes that  $\mathbf{p}(\mathbf{f}|\mathbf{X}_i, \mathbf{y}_i)$  is a joint Gaussian with some mean  $\boldsymbol{\mu}$  and variance  $\boldsymbol{\Sigma}$ .



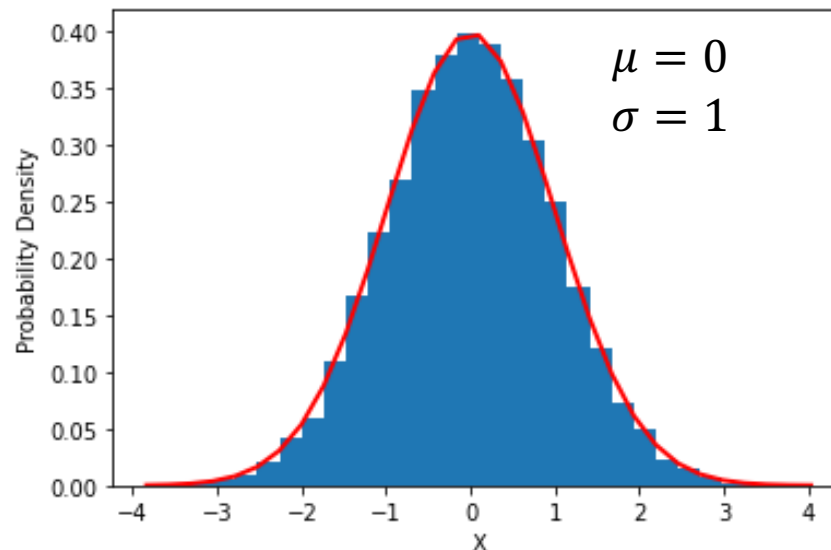
# Gaussian Processes - Basics

- In probability theory, a Gaussian (or normal) distribution is a type of continuous probability distribution for a real-valued random variable. Conventionally, we write the Gaussian distribution of mean  $\mu$  and covariance  $\Sigma$  as

$$f \sim \mathcal{N}(\mu, \Sigma)$$

- In 1D Gaussian distribution, we have

$$f \sim \mathcal{N}(\mu, \Sigma = \sigma^2)$$



$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

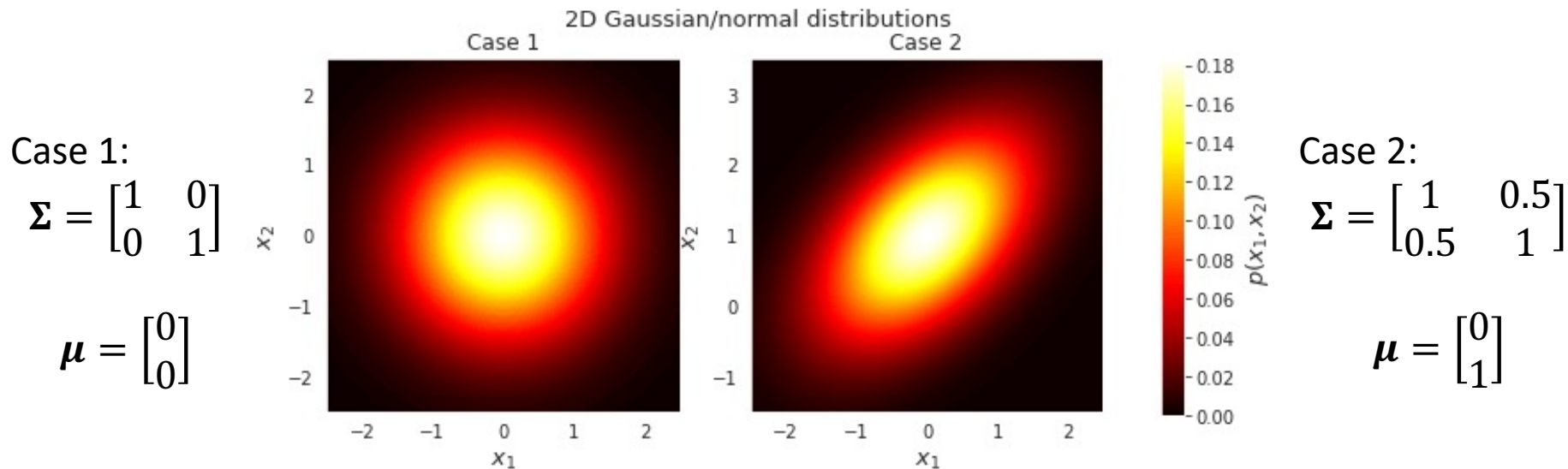
↓

$$\int_{-\infty}^{\infty} f(x) dx = 1$$

# Gaussian Processes - Basics

- For multivariate Gaussian distribution, the joint probability of dimension  $d$  is given by

$$p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}|}} \exp \left( -\frac{1}{2} \underbrace{(\mathbf{x} - \boldsymbol{\mu})^T}_{d \times 1} \cdot \underbrace{\boldsymbol{\Sigma}^{-1}}_{d \times d} \cdot \underbrace{(\mathbf{x} - \boldsymbol{\mu})}_{d \times 1} \right)$$

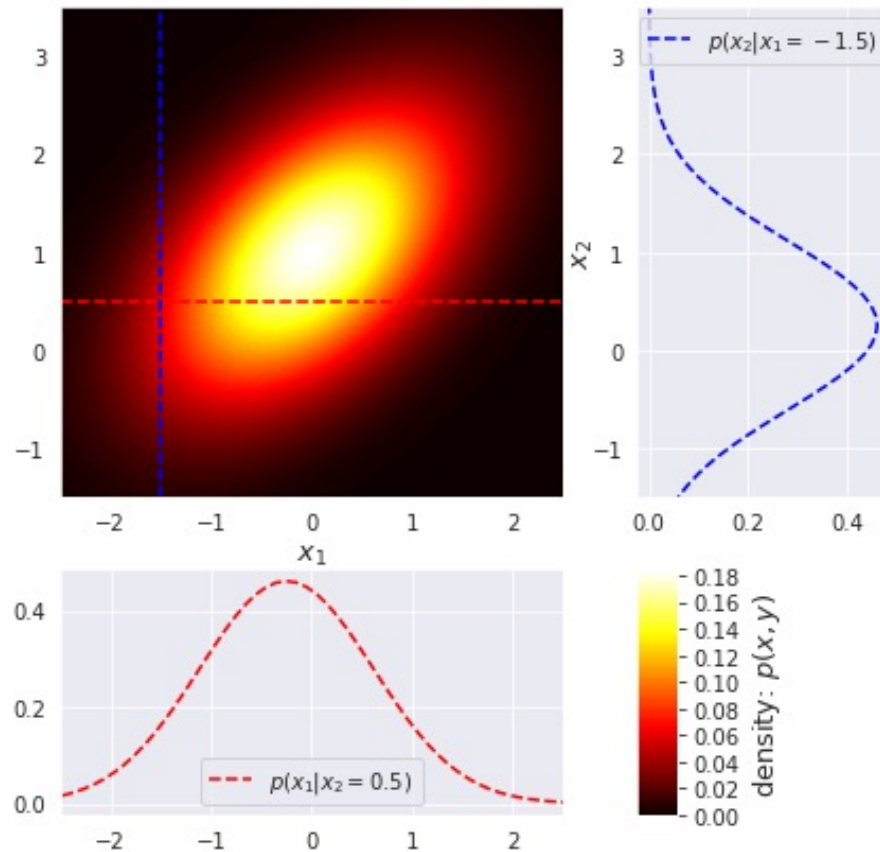


Covariance  $\boldsymbol{\Sigma}$  is a symmetric, positive definite matrix!

# Gaussian Processes - Basics

- Conditional distribution of  $x_1$  given  $x_2$  is the probability distribution of  $x_1$  when  $x_2$  is known to be a particular value

Posterior distribution  $\leftarrow p(x_1 | x_2) = \frac{p(x_2 | x_1) * p(x_1)}{p(x_2)}$   $\rightarrow$  Prior distribution



$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \right)$$

$$p(x_1 | x_2) = p_{1|2} = \mathcal{N}(\mu_{1|2}, \Sigma_{1|2})$$

$$\mu_{1|2} = \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (x_2 - \mu_2)$$

$$\Sigma_{1|2} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$$

Similarly, one can derive multidimensional conditional probability distribution (e.g. 2d)

$$\begin{aligned} p(x_1, x_2 | x_3, x_4, \dots) \\ = \mathcal{N}(\mu_{1,2|3,4,\dots}, \Sigma_{1,2|3,4,\dots}) \end{aligned}$$

# Gaussian Processes for Regression

- Now we discuss GPs for regression. Let the prior on the regression function be a Gaussian process such that

$$f(\mathbf{x}) \sim GP(m(\mathbf{x}), \kappa(\mathbf{x}, \mathbf{x})) = N(f | \boldsymbol{\mu}, K)$$

Mean function:

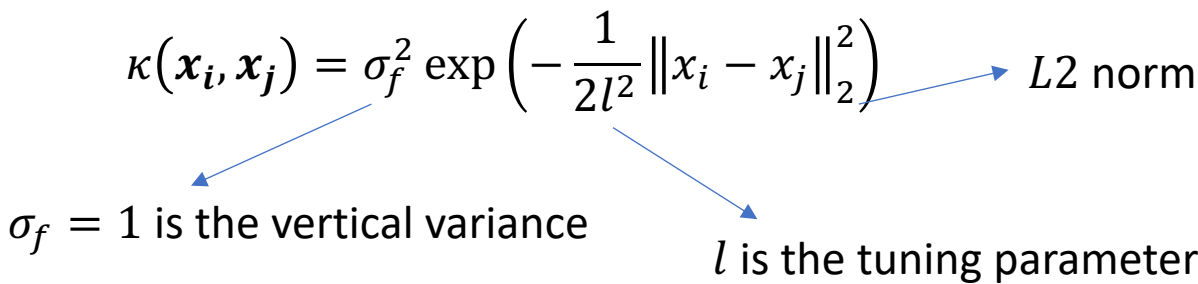
$$m(\mathbf{x}) = E[f(\mathbf{x})], \boldsymbol{\mu} = (m(\mathbf{x}_1), \dots, m(\mathbf{x}_N))$$

Covariance function kernel:

$$\kappa(\mathbf{x}, \mathbf{x}) = E[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}) - m(\mathbf{x}))]$$

$$K_{ij} = \kappa(\mathbf{x}_i, \mathbf{x}_j)$$

- It is common to use a mean **function** of  $m(\mathbf{x}) = 0$  for the prior.
- The kernel selection highly depend on the considered problem. Typical choice is a **squared exponential kernel** (Gaussian kernel or RBF kernel)

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \sigma_f^2 \exp\left(-\frac{1}{2l^2} \|\mathbf{x}_i - \mathbf{x}_j\|_2^2\right)$$


$\sigma_f = 1$  is the vertical variance

$l$  is the tuning parameter



# Gaussian Processes – Noise-free Regression

- Given a training data set  $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{f}_i), i = 1, \dots, N\}$ , where  $\mathbf{f}_i$  is the noise-free observation (output function), we want to predict  $\mathbf{f}_*$  for a test set  $\mathbf{x}_*$
- By the definition of GPs, the joint distribution has following distribution

$$\begin{pmatrix} \mathbf{f} \\ \mathbf{f}_* \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} \boldsymbol{\mu} \\ \boldsymbol{\mu}_* \end{pmatrix}, \begin{pmatrix} \mathbf{K} & \mathbf{K}_* \\ \mathbf{K}_*^T & \mathbf{K}_{**} \end{pmatrix} \right)$$

Prior distribution

$$p(\mathbf{f}_* | \mathbf{x}_*, \mathbf{x}, \mathbf{f}) = \mathcal{N}(\mathbf{f}_* | \boldsymbol{\mu}_*, \boldsymbol{\Sigma}_*)$$

Posterior distribution

$$\boldsymbol{\mu}_* = \boldsymbol{\mu}(\mathbf{x}_*) + \mathbf{K}_*^T \mathbf{K}^{-1} (\mathbf{f} - \boldsymbol{\mu}(\mathbf{x}))$$

$$\boldsymbol{\Sigma}_* = \mathbf{K}_{**} - \mathbf{K}_*^T \mathbf{K}^{-1} \mathbf{K}_*$$

$$\mathbf{K}_{**} = \kappa(\mathbf{x}_*, \mathbf{x}_*) \quad \mathbf{K}_* = \kappa(\mathbf{x}, \mathbf{x}_*) \quad \mathbf{K} = \kappa(\mathbf{x}, \mathbf{x})$$

Typically, we let  $\boldsymbol{\mu}=0$ ; then  $\boldsymbol{\mu}_* = \mathbf{K}_*^T \mathbf{K}^{-1} \mathbf{f}$



# Gaussian Processes – Implementation

- GPs is a rather computationally expensive method; the major computation time is consumed for the inversion of  $\mathbf{K}$ , which is about  $\mathcal{O}(N^3)$ .
- In addition, a simple inversion operation  $\mathbf{K}^{-1}$  is unstable. There are two solutions:
  1. First compute the Cholesky decomposition,  $\mathbf{K} = \mathbf{L}\mathbf{L}^T$ ; next,  $\mathbf{K}^{-1}\mathbf{y} = \mathbf{L}^{-T}\mathbf{L}^{-1}\mathbf{y}$ . ( $\mathbf{L}$  is lower triangular matrix)

---

**Algorithm 15.1:** GP regression

---

```
1  $\mathbf{L} = \text{cholesky}(\mathbf{K} + \sigma_y^2 \mathbf{I});$   
2  $\boldsymbol{\alpha} = \mathbf{L}^T \setminus (\mathbf{L} \setminus \mathbf{y});$   
3  $\mathbb{E}[f_*] = \mathbf{k}_*^T \boldsymbol{\alpha};$   
4  $\mathbf{v} = \mathbf{L} \setminus \mathbf{k}_*;$   
5  $\text{var}[f_*] = \kappa(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{v}^T \mathbf{v};$   
6  $\log p(\mathbf{y}|\mathbf{X}) = -\frac{1}{2}\mathbf{y}^T \boldsymbol{\alpha} - \sum_i \log L_{ii} - \frac{N}{2} \log(2\pi)$ 
```

---

2. Another approach is to solve the linear equation, namely  $\mathbf{K}\boldsymbol{\alpha} = \mathbf{y}$  (e.g. conjugate gradient)

# Gaussian Processes – Implementation

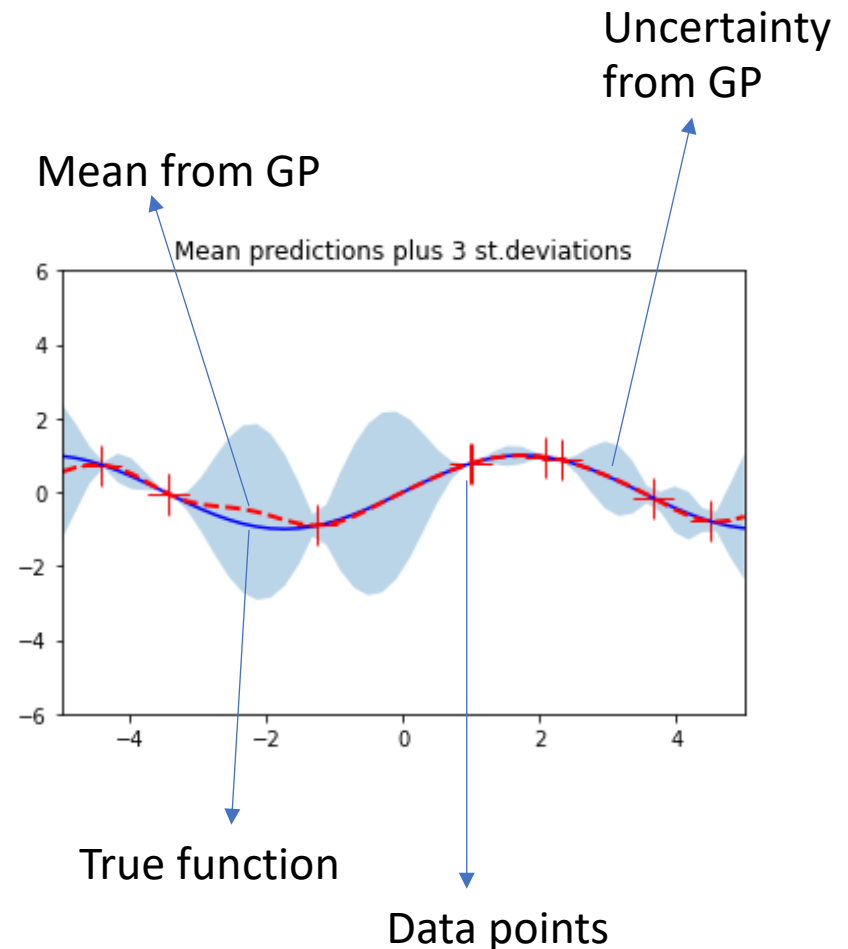
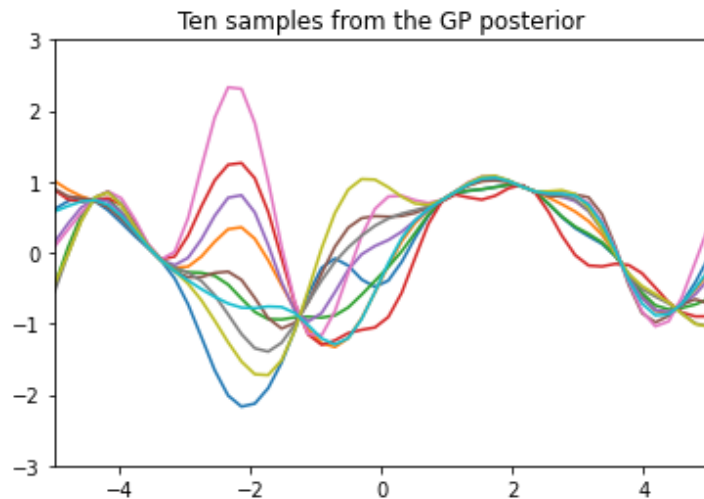
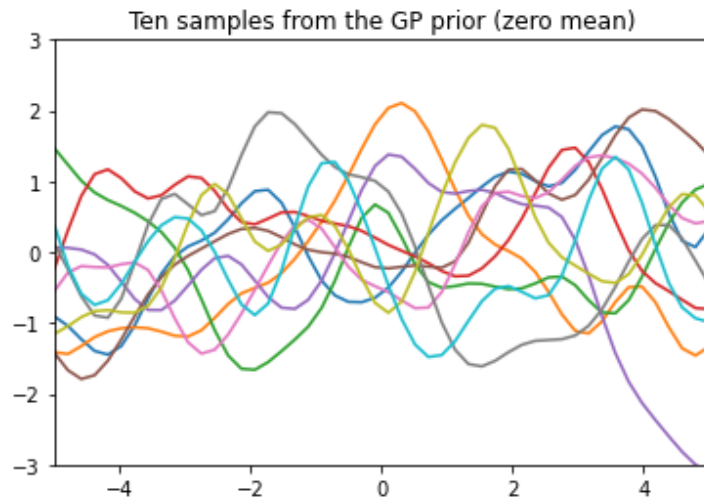
```
K = kernel(X, X) # compute the kernel function with training data X
L = np.linalg.cholesky(K) # Cholesky decomposition
Lk = np.linalg.solve(L, kernel(X, Xtest)) # Inv(L)*kernel(X, Xtest)

# compute the mean at our test points.
mu = np.dot(Lk.T, np.linalg.solve(L, y)) # y is training outcome
# compute the variance at our test points.
K_ = kernel(Xtest, Xtest)
s2 = np.diag(K_) - np.sum(Lk**2, axis=0) # variance
s = np.sqrt(s2) # standard deviation
```

- (1) Using Cholesky decomposition to compute  $L$  for  $K = LL^T$ ; it is much easier (more stable) for invert  $L$  than  $K$
- (2) 'mu' is the mean function for test data point, which is the curve we try to predict;
- (3) 's' is the standard deviation for the uncertainty of the predicted curve. The power of GPs comes from the ability to predict the reliable estimate of uncertainty!

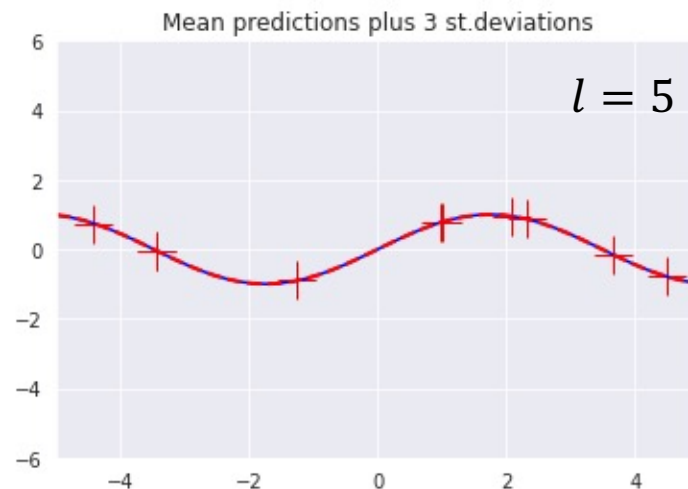
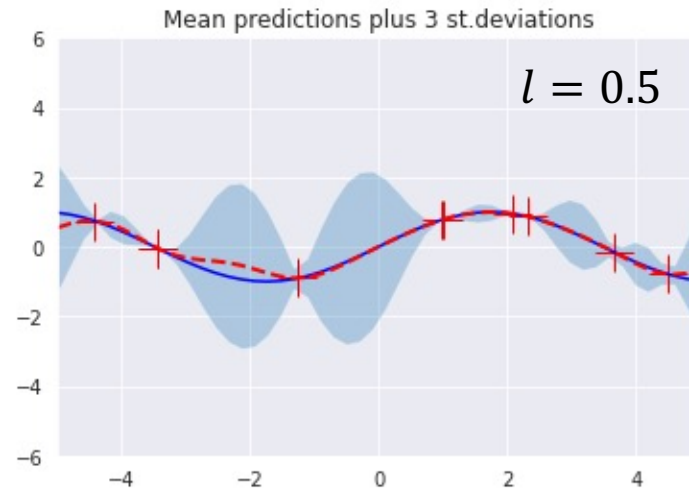
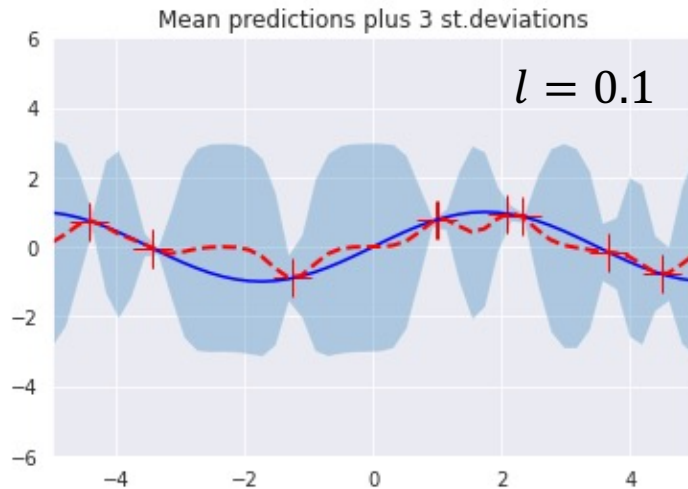
# Gaussian Processes – Example (noise-free)

- We assume the prior to be zero mean ( $\mu = 0$ ) and kernel function to be RBF ( $l = 0.5$ ).



# Gaussian Processes – Kernel Parameter

- We examine how kernel parameter affect the mean prediction as well as uncertainty.



# Gaussian Processes – Regression with Noise

- Now let us consider the case where what we observe is a noisy function,  $\mathbf{y} = \mathbf{f}(\mathbf{x}) + \epsilon$ , where  $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma_y^2)$ .
- The covariance of the noisy data is  $\text{cov}[\mathbf{y}|\mathbf{X}] = \mathbf{K} + \sigma_y^2 \mathbf{I} \equiv \mathbf{K}_y$ .
- The joint density of the observed data and test points ( $\mathbf{X}_*$ ); again, we typically assume the mean of the prior to be zero ( $\mu = \mathbf{0}$ )

$$\begin{pmatrix} \mathbf{y} \\ \mathbf{f}_* \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} \mu \\ \mu_* \end{pmatrix}, \begin{pmatrix} \mathbf{K}_y & \mathbf{K}_* \\ \mathbf{K}_*^T & \mathbf{K}_{**} \end{pmatrix} \right)$$

$$p(\mathbf{f}_* | \mathbf{X}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(\mathbf{f}_* | \mu_*, \Sigma_*)$$

$$\mu_* = \mu(\mathbf{X}_*) + \mathbf{K}_*^T \mathbf{K}_y^{-1} (\mathbf{f} - \mu(\mathbf{X}))$$

$$\Sigma_* = \mathbf{K}_{**} - \mathbf{K}_*^T \mathbf{K}_y^{-1} \mathbf{K}_*$$

$\mu = \mathbf{0}$

$$\begin{pmatrix} \mathbf{y} \\ \mathbf{f}_* \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} \mathbf{0} \\ \mu_* \end{pmatrix}, \begin{pmatrix} \mathbf{K}_y & \mathbf{K}_* \\ \mathbf{K}_*^T & \mathbf{K}_{**} \end{pmatrix} \right)$$

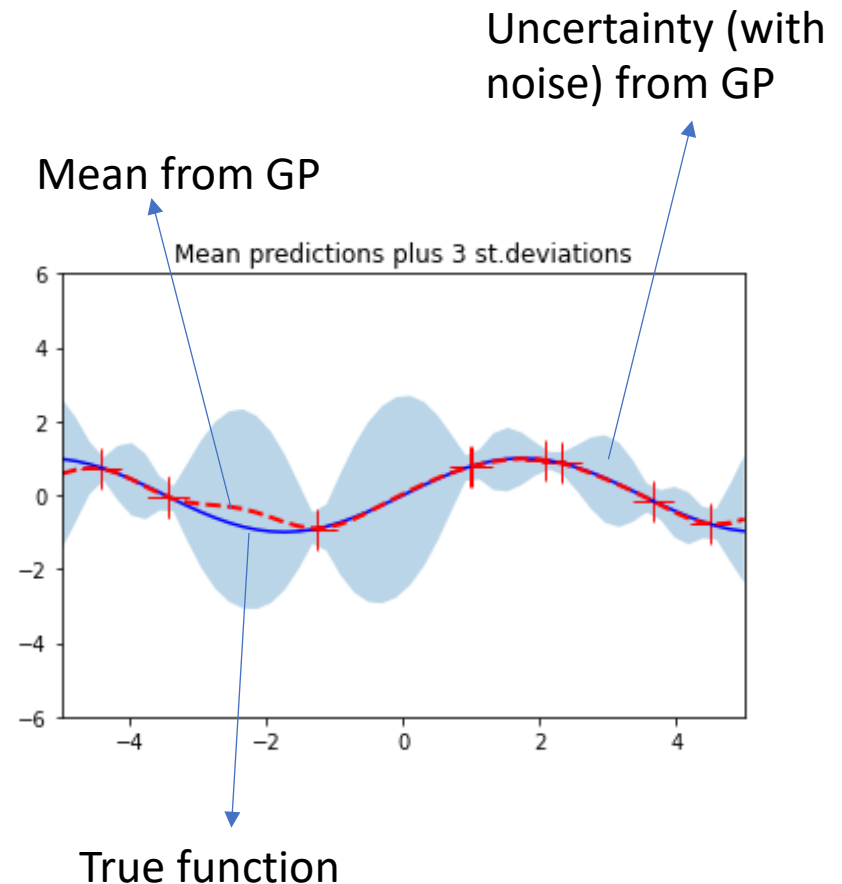
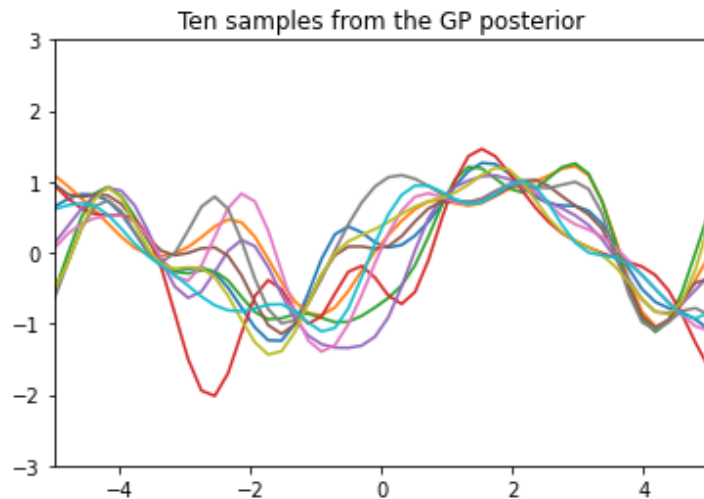
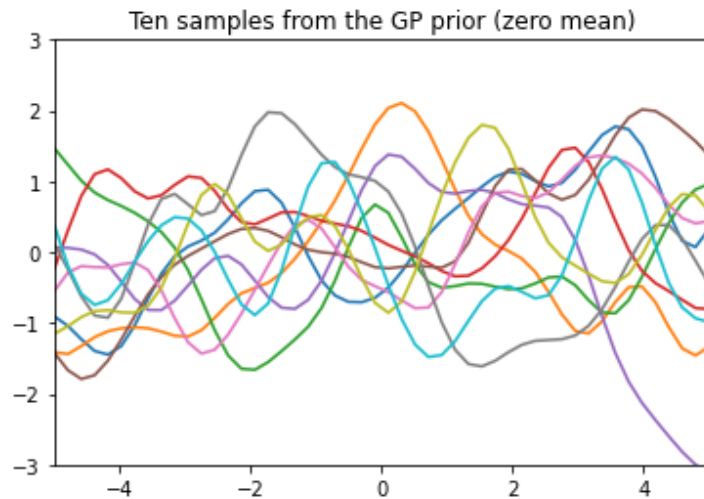
$$p(\mathbf{f}_* | \mathbf{X}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(\mathbf{f}_* | \mu_*, \Sigma_*)$$

$$\mu_* = \mathbf{K}_*^T \mathbf{K}_y^{-1} \mathbf{y}$$

$$\Sigma_* = \mathbf{K}_{**} - \mathbf{K}_*^T \mathbf{K}_y^{-1} \mathbf{K}_*$$

# Gaussian Processes – Example (with noise)

- We assume the prior to be zero mean ( $\mu = 0$ ) and kernel function to be RBF ( $l = 0.5, \sigma_y = 0.01$ ).



# Gaussian Processes – Kernel Parameters Estimation

- To estimate the kernel parameters, we could exhaust the search of a discrete values with validation as an objective. However, this can be very slow.

- The marginal likelihood is  $p(\mathbf{f}|\mathbf{X}) = \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K})$

$$\log p(\mathbf{y}|\mathbf{X}) = \log \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K}_y) = -\frac{1}{2}\mathbf{y}^T \mathbf{K}_y^{-1} \mathbf{y} - \frac{1}{2} \log |\mathbf{K}_y| - \frac{N}{2} \log(2\pi)$$

To maximize the likelihood, let the kernel be the hyper-parameters denoted as  $\theta_i = (l, \sigma_y)$  and then apply gradient based approach

$$\frac{\partial}{\partial \theta_i} \log p(\mathbf{y}|\mathbf{X}) = \frac{1}{2} \mathbf{y}^T \frac{\partial \mathbf{K}_y^{-1}}{\partial \theta_i} \mathbf{y} - \frac{1}{2} \text{tr} \left( \mathbf{K}_y^{-1} \frac{\partial \mathbf{K}_y}{\partial \theta_i} \right)$$

- In case, the above method is not convex, other optional methods include:
  - (1) Optimize the posterior probability
  - (2) Multiple kernel learning (optimize weight instead of kernel parameters)



# Gaussian Processes – Connection to Other Methods

- Connecting back to Ridge Regression:

$$\boldsymbol{\beta} = (\mathbf{X}^T \mathbf{X} + \sigma^2 \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\hat{\mathbf{y}} = \mathbf{X} \boldsymbol{\beta} = \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$$

Let  $\mathbf{K} = \mathbf{X}^T \mathbf{X}$ , we have  $\hat{\mathbf{y}} = \mathbf{K}(\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$ , which is exactly the posterior form of GPs.

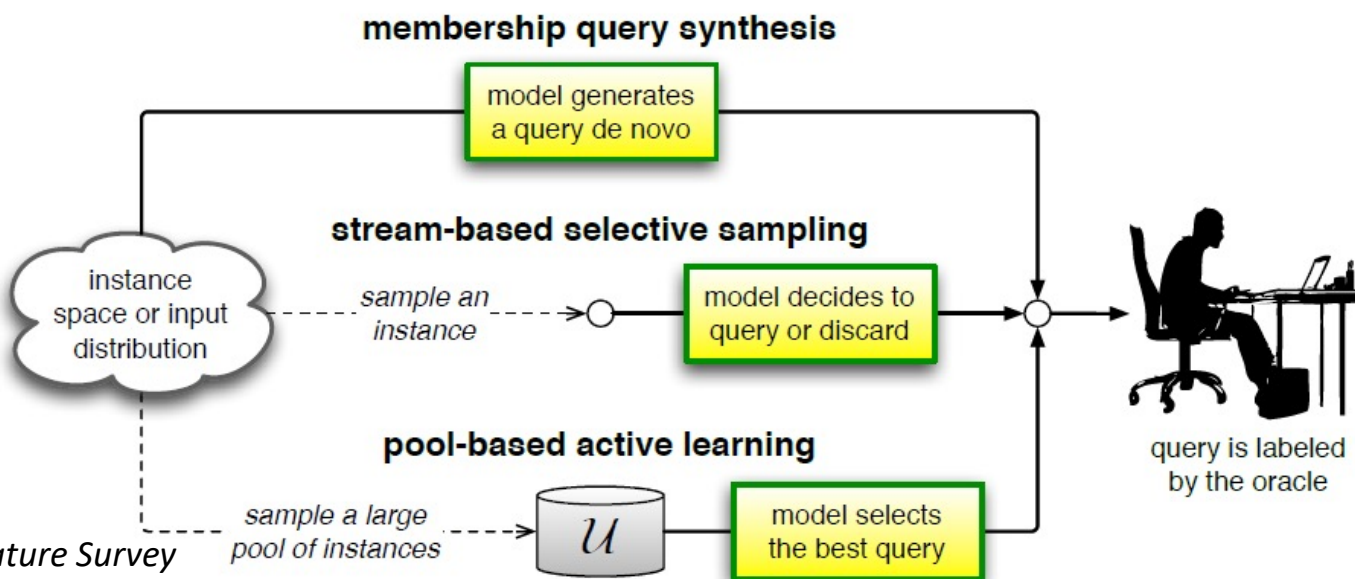
- Infinitely large Neural Nets (one hidden layer)  $\rightarrow$  GPs

- Output function:  $z^1(\mathbf{x}) = \beta_0^1 + \sum_{j=1}^N x_j^1(\mathbf{x}) \beta_j^1$
- $x_j^1(\mathbf{x}) = \phi(\beta_j^0 + \sum_{k=1}^N x_k \beta_{jk}^0)$  is the activation function
- Assume that  $\beta_0^1$  and  $\beta_j^1$  are independent and normal distribution variable
- Central limit theorem:

$$\lim_{N \rightarrow \infty} \sum_{j=1}^N x_j^1(\mathbf{x}) \beta_j^1 = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{j=1}^N x_j^1(\mathbf{x}) \tilde{\beta}_j^1 \sim \mathcal{N}(\mathbf{0}, \mathbf{K})$$

# Active Learning

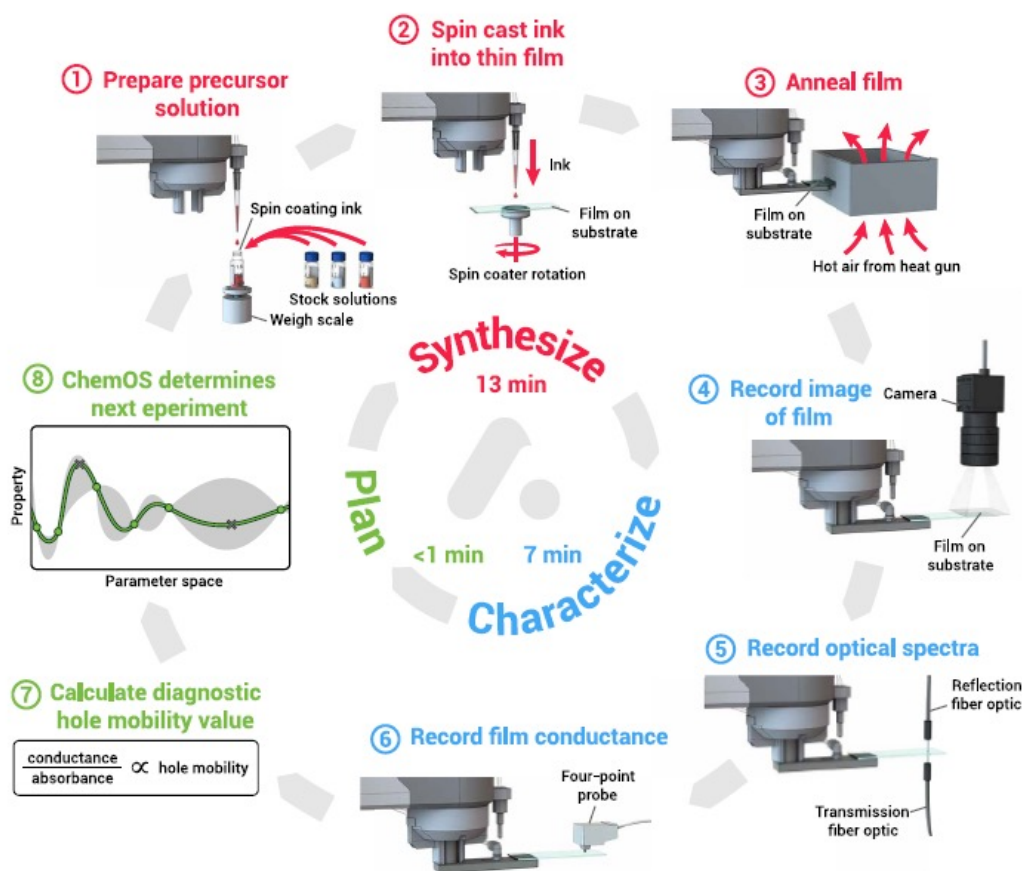
- Active learning is a machine learning algorithm can achieve greater accuracy with fewer labeled training instances if it is allowed to choose the data from which it learns.
- An active learner may ask queries in the form of unlabeled instances to be labeled by an oracle (e.g., a human annotator or a machine)
- Active learning is particularly suitable when unlabeled data may be abundant but labels are difficult, time-consuming, or expensive to obtain



Burr Settles (2009)  
Active Learning Literature Survey  
Technical Report #1648

# Active Learning - Gaussian Processes

- Because GPs offer the uncertainty of the prediction, the algorithm can provide an effective query mechanism for active learning.
- A recent interesting example: Self-driving laboratory for accelerated discovery of thin-film materials



- Discovering and optimizing commercially viable materials (e.g. for clean energy applications) typically takes more than a decade

- Material synthesis and characterization offer the data point for planning

- Active learning combining with GPs completes the query, thereof entire process is automatic.