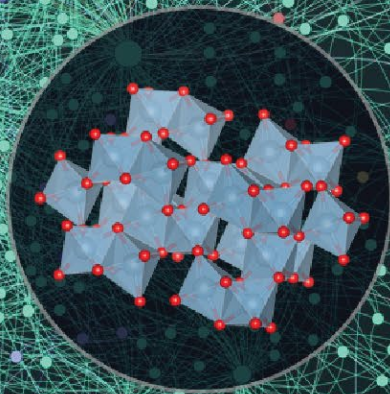
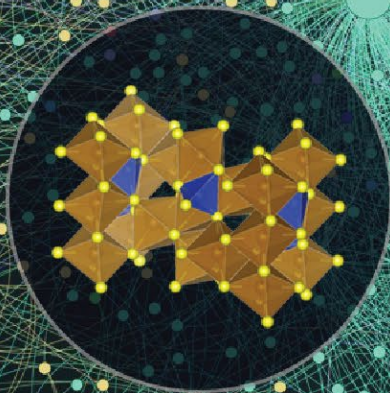


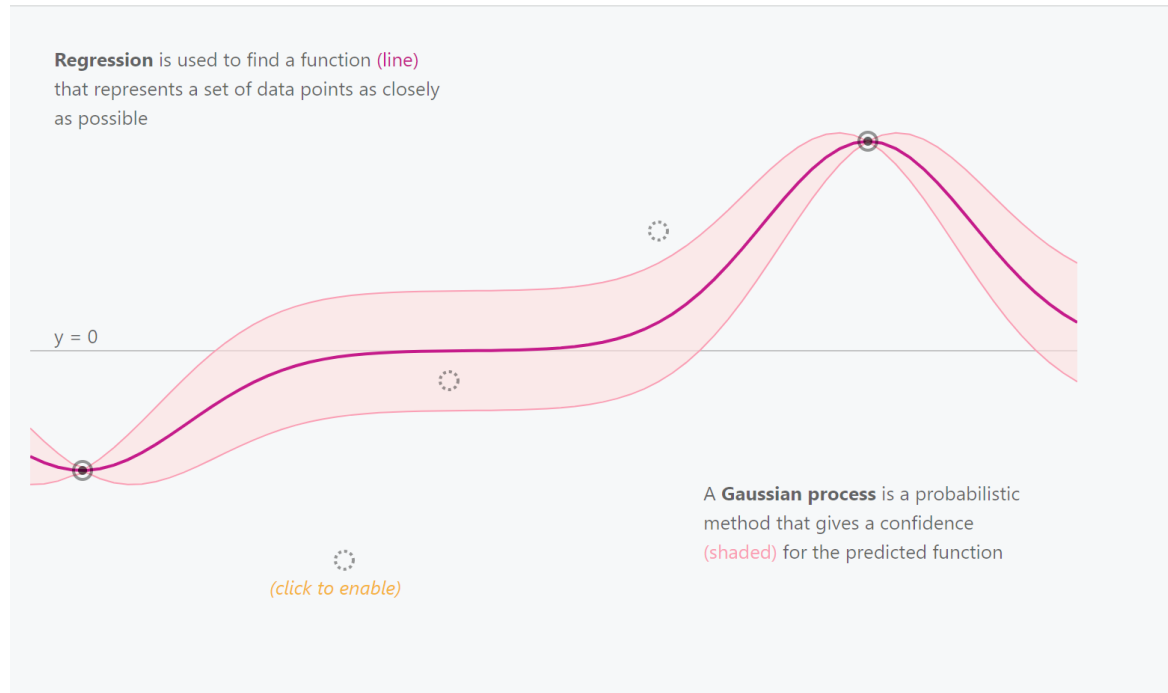
Gaussian processes



Gaussian Processes are a bit tricky to understand!

A Visual Exploration of Gaussian Processes

How to turn a collection of small building blocks into a versatile tool for solving regression problems.



AUTHORS

Jochen Görtler
Rebecca Kehlbeck
Oliver Deussen

AFFILIATIONS

University of Konstanz
University of Konstanz
University of Konstanz

PUBLISHED

April 2, 2019

DOI

10.23915/distill.00017

IntuitiveML (YouTube) has a pretty good walkthrough and summary

- Input: training data (TRAIN) and test data (TEST)
- Output: value and confidence interval at every point in TEST

How it works:

- Center the data to have mean 0
- Concatenate the TRAIN and TEST datapoints to COMBINED
- Create a covariance matrix
 - Use a kernel function to find similarity between each point in COMBINED
 - The similarity between point i and point j is the value of the matrix at (i,j)
- Calculate the joint distribution using the covariance matrix
- Calculate the conditional distribution – aka calculate the probability of TEST given TRAIN
- Marginalize the conditional distribution – aka calculate the probability of distribution of a single point in TEST given all other points
- Obtain mean and stdev for every point in TEST
- Calculate the confidence interval using stdev to the desired level (eg 95%)
- Mean is the regression value

In general, we are trying to predict Y conditional on X

$$C(Y|X) = g(X)$$

Consider linear regression

$$C(Y|X) = g(X, \beta) = \beta_0 + X_1\beta_1 + X_2\beta_2 + \dots$$

Assumptions:

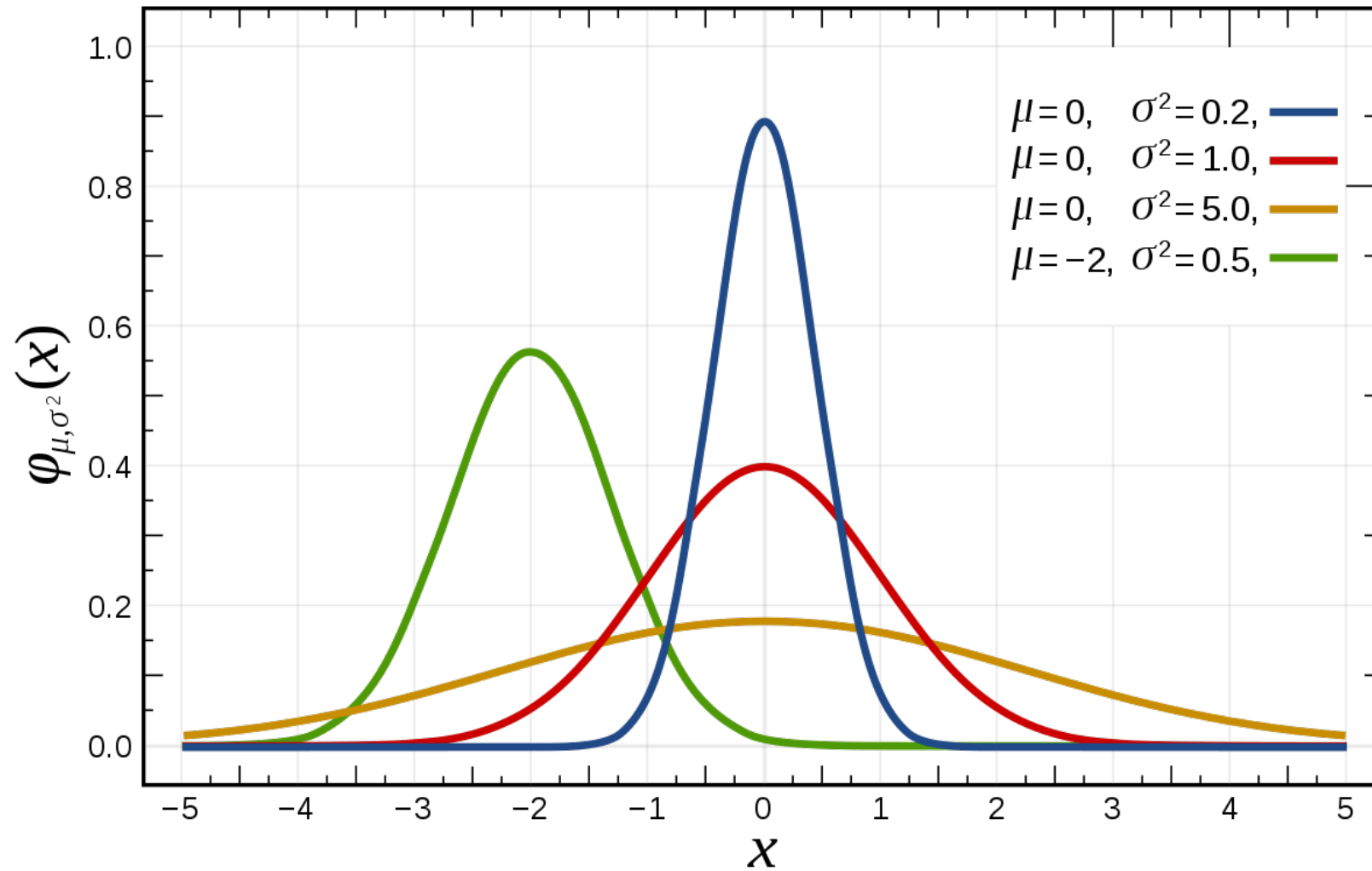
Parametric (data is normally distributed) and linear relationships between variables

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$
$$\epsilon_i \sim N(0, \sigma)$$

Logistic regression also possible

$$C(Y|X) = \frac{1}{1 + \exp(\beta X)}$$

Instead of relying on linear or non-linear approaches to fit data, let's rely on gaussians!



Instead of relying on linear or non-linear approaches to fit data, let's rely on gaussians!

Consider a bivariate Gaussian distribution

$$p(y_1, y_2 | \mu, \Sigma) = N \left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & \sigma_1 \sigma_2 \rho \\ \sigma_1 \sigma_2 \rho & \sigma_2^2 \end{bmatrix} \right)$$

Off diagonal elements in standard deviation matrix are the covariance terms which include the correlation term ($\rho = 0$, *uncorrelated*, $\rho = 1$, *correlated*)

Gaussians can be easily marginalized and conditioned

Why are Gaussians a great idea?

1. They are easy to marginalize (remove all parts of covariance matrix except what you want to keep)
 - if you want to ditch y_2 , then drop $\mu_2, \sigma_1\sigma_2\rho$, and σ_2^2
2. They are easy to condition (given the value of one value in our multivariate normal distribution, what do we expect some other unobserved value to be?)
 - $p(x|y) = N\left(\mu_x + \Sigma_{xx} \Sigma_y^{-1}(y - \mu_y)\right), \Sigma_x - \Sigma_{xy} \Sigma_y^{-1} \Sigma_{xy}^T$
Conditional mean, conditional covariance
 - Where μ_x is marginal mean (easy) and Σ_x is marginal covariance each minus some easy to calculate terms

In Gaussian Process we will generalize multivariate normal to contain infinite elements

Gaussian distribution:

$$y \sim N(\mu, \Sigma)$$

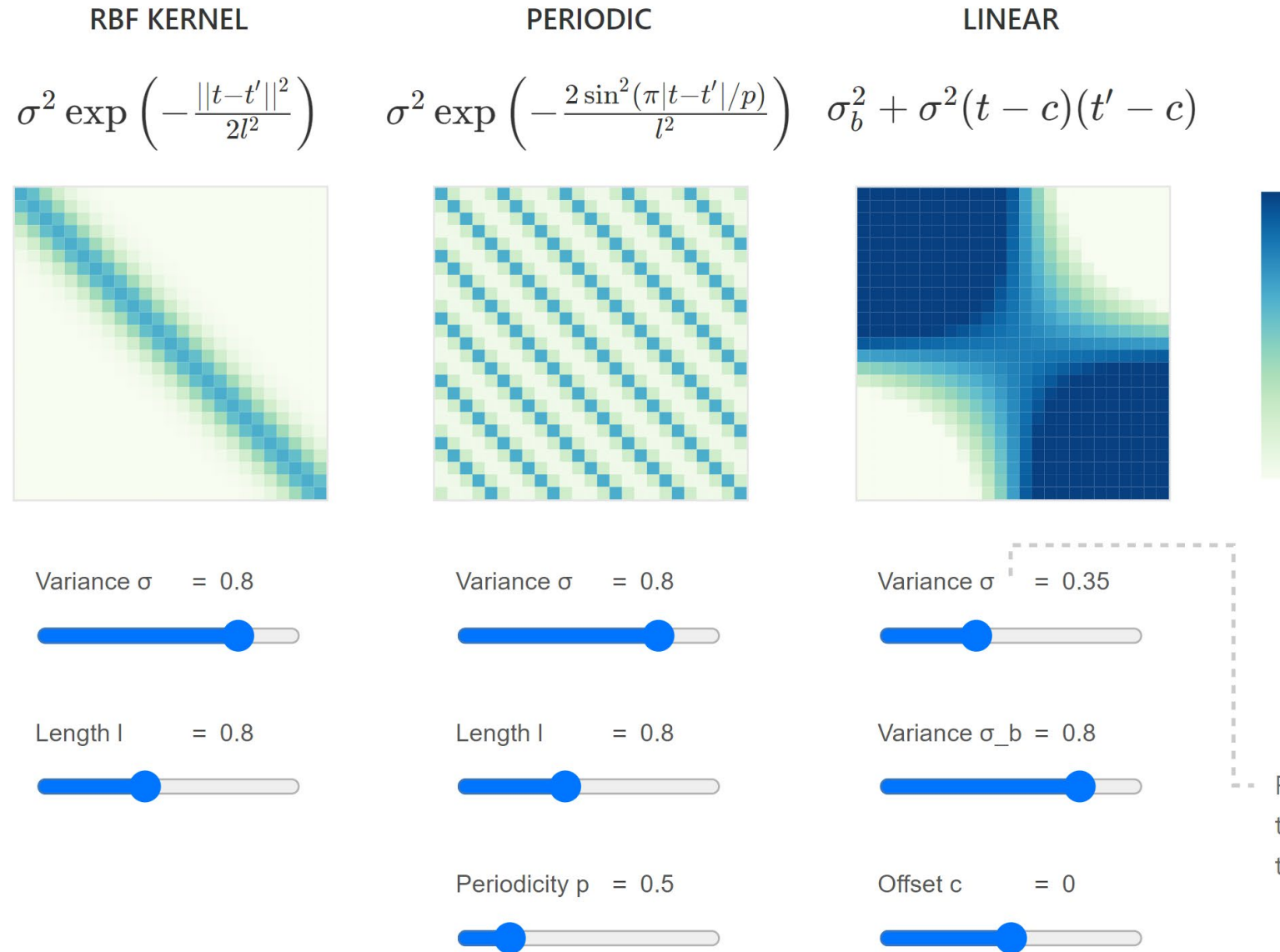
Gaussian process:

$$f \sim GP(m(x), k(x, x'))$$

Where $m(x)$ is a mean function, and $k(x, x')$ is a covariance function

Definition: Gaussian process is an infinite collection of variables, any finite subset of which has a Gaussian distribution.

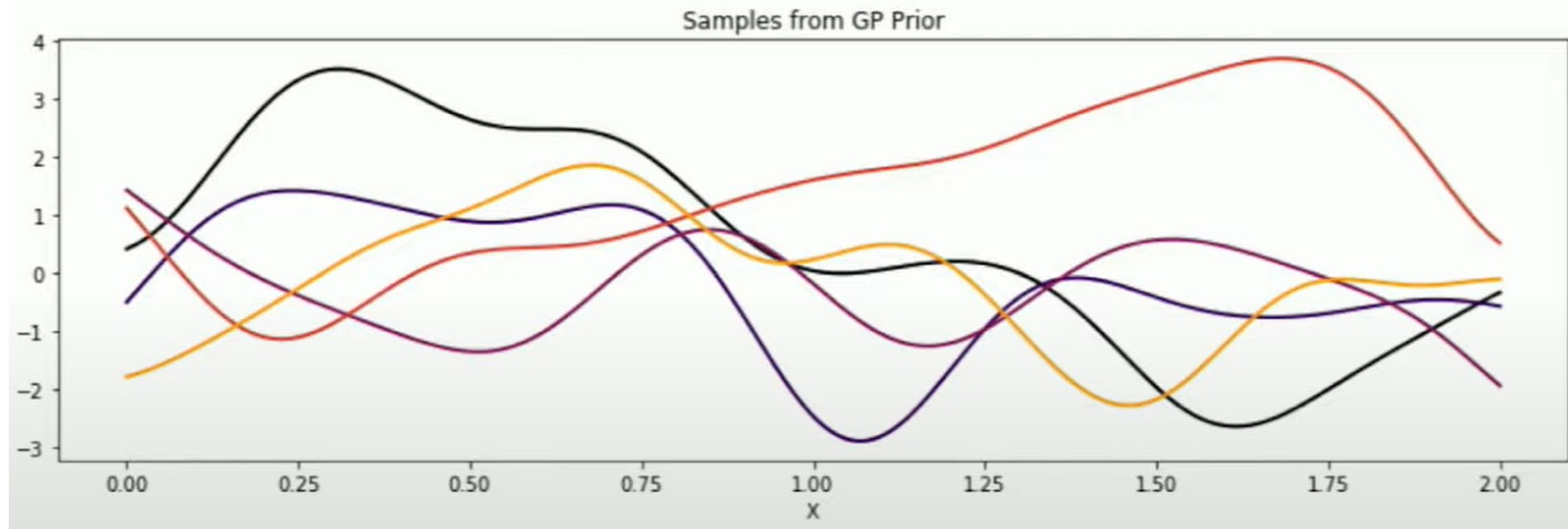
Tricky part of Gaussian process is figuring out which function to use to model covariance



For the **Linear** kernel the parameter **Variance σ** determines the average distance away from the function's mean.

Consider all the different covariance functions we can get from just one exponential quadratic

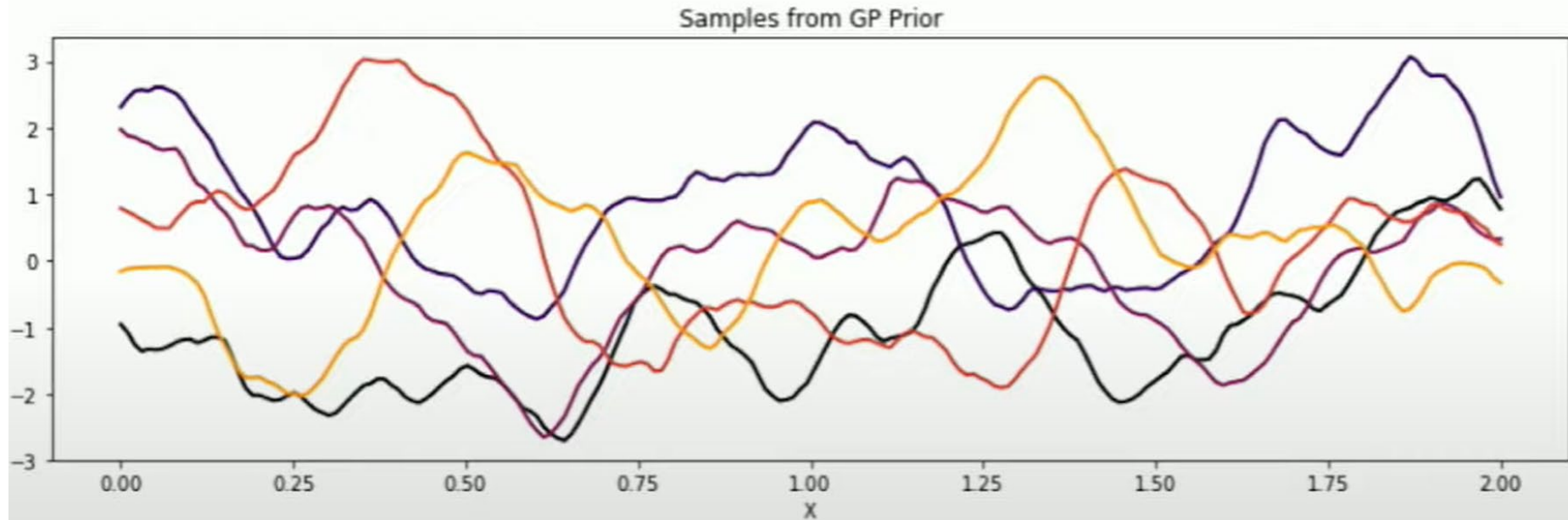
$$k(x, x') = \exp \left[-\frac{(x - x')^2}{2\ell^2} \right]$$



One parameter, length scale, that determines how far away you have to go away from two points to make them independent

Matern(3/2) is another popular covariance function

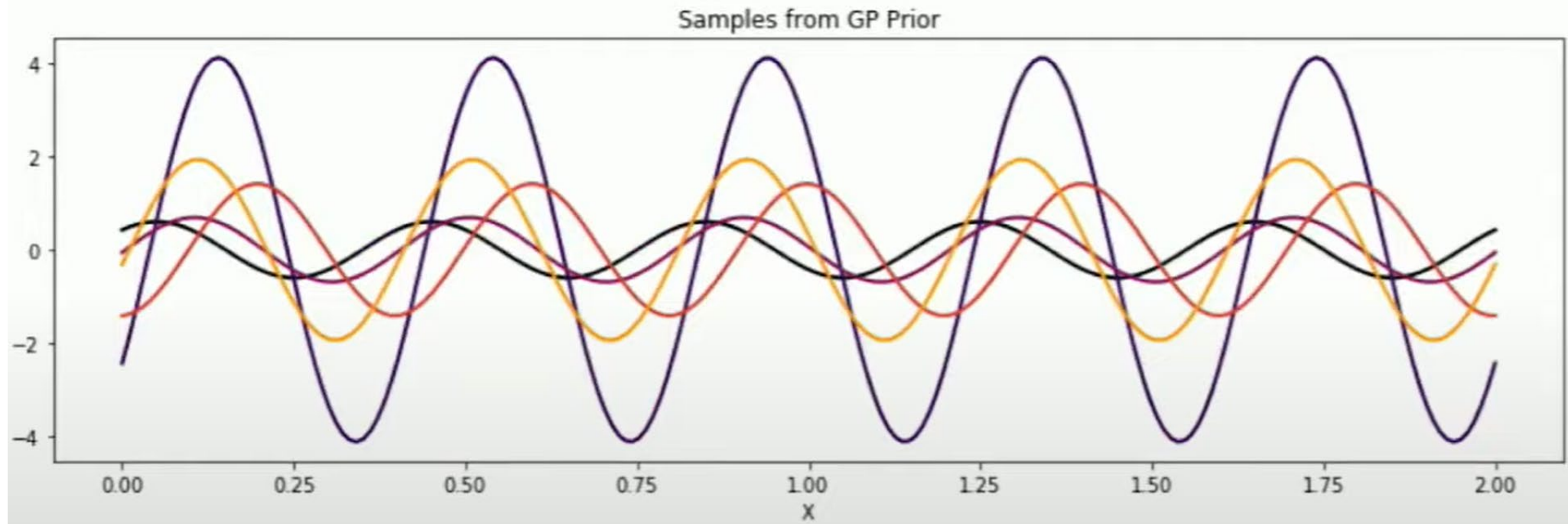
$$k(x, x') = \left(1 + \frac{\sqrt{3(x - x')^2}}{\ell} \right) \exp \left[-\frac{\sqrt{3(x - x')^2}}{\ell} \right]$$



Length scale parameter, plus a scaling term to scale roughness

Cosine is another popular covariance function

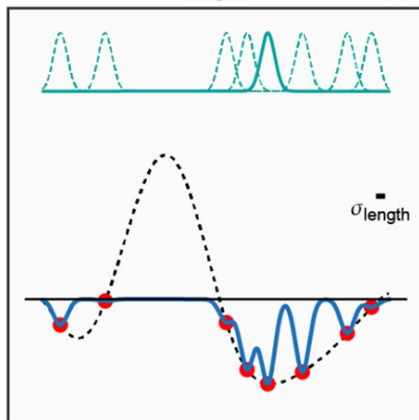
$$k(x, x') = \cos \left(\frac{||x - x'||}{\ell^2} \right)$$



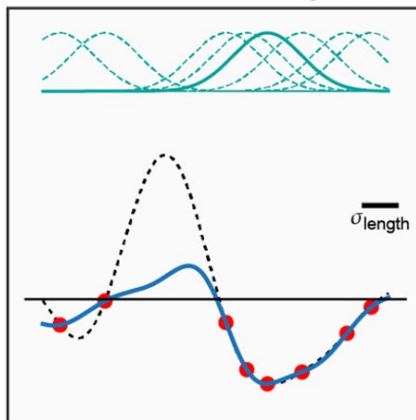
Getting the kernel hyperparameters right is important

Learning from function values

Too small σ_{length} (overfitting)

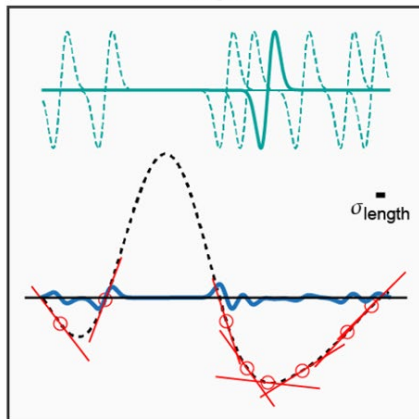


Appropriate σ_{length}

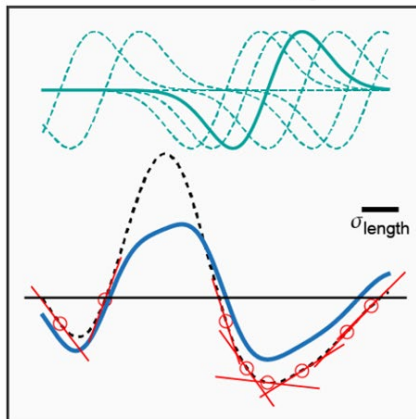


Learning from derivative values

Too small σ_{length} (overfitting)



Appropriate σ_{length}



A mean function is also necessary, but not as hard to pick

Mean function generates mean vectors

$$\mu \sim m(x|\phi)$$

Simplest is just zero!

$$m(x) = 0$$

Or constant

$$m(x) = C$$

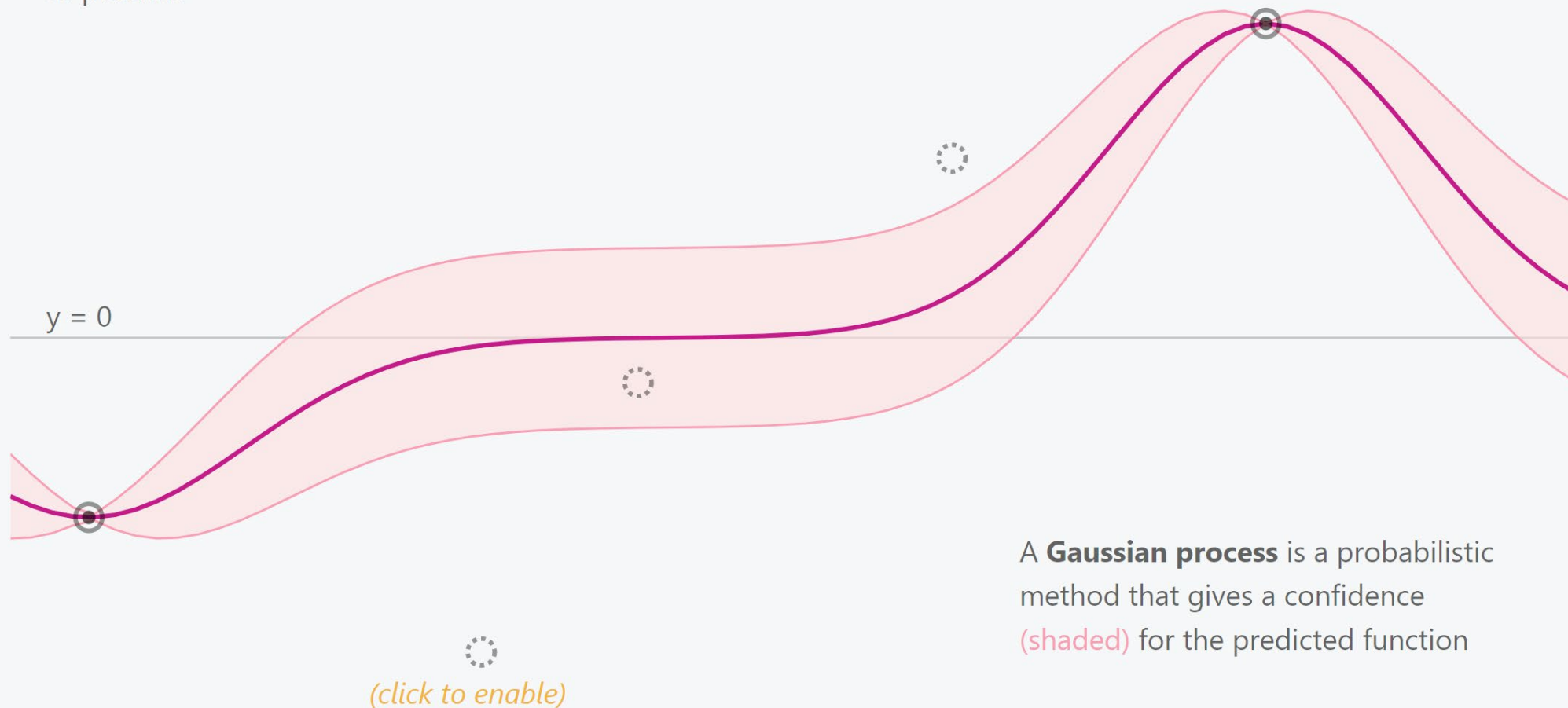
Or linear

$$m(x) = Ax + b$$

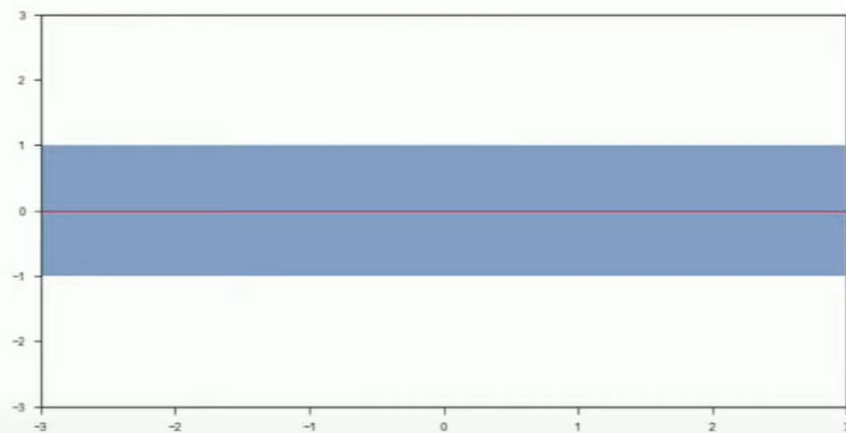
Posterior does not use the mean, it is just used as a guess for the prior because where you have no data, it trends towards the mean

In the absence of data, the regression trends towards the mean

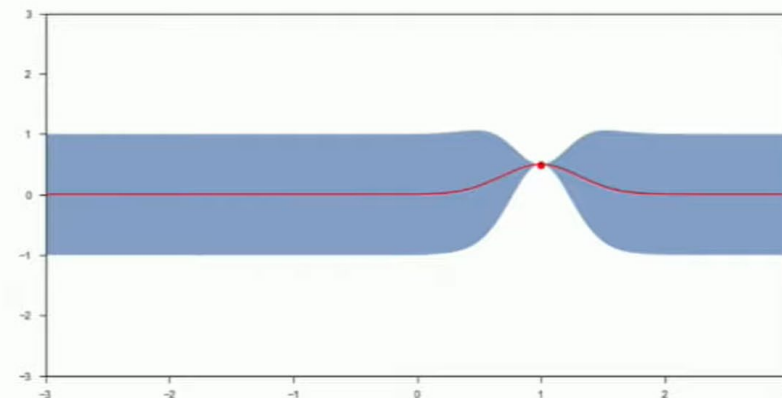
Regression is used to find a function (line) that represents a set of data points as closely as possible



Let's draw from a prior Gaussian process and we'll do it point by point

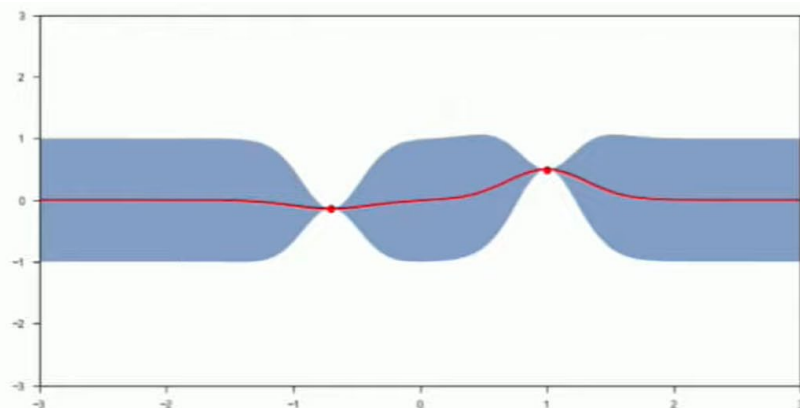


$$y \sim GP(m = 0, k = \text{Quad}(1))$$



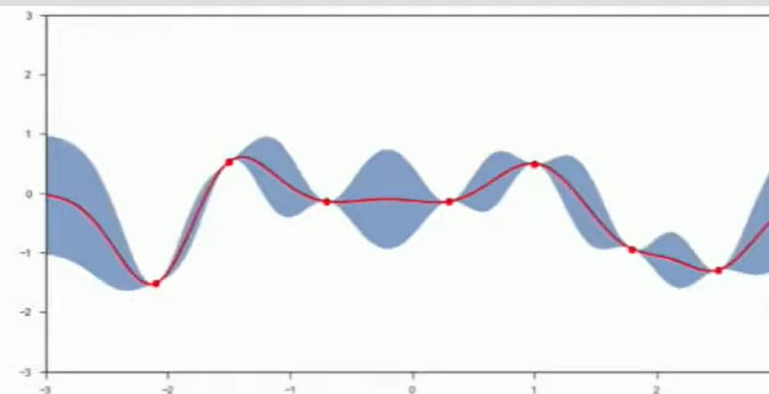
```
x = [1.]
y = [np.random.normal(scale=σ_θ)]
y
```

```
[0.4967141530112327]
```



```
x = -0.7
m, s = conditional([-0.7], x, y, θ)
y.append(np.random.normal(m, s))
y
```

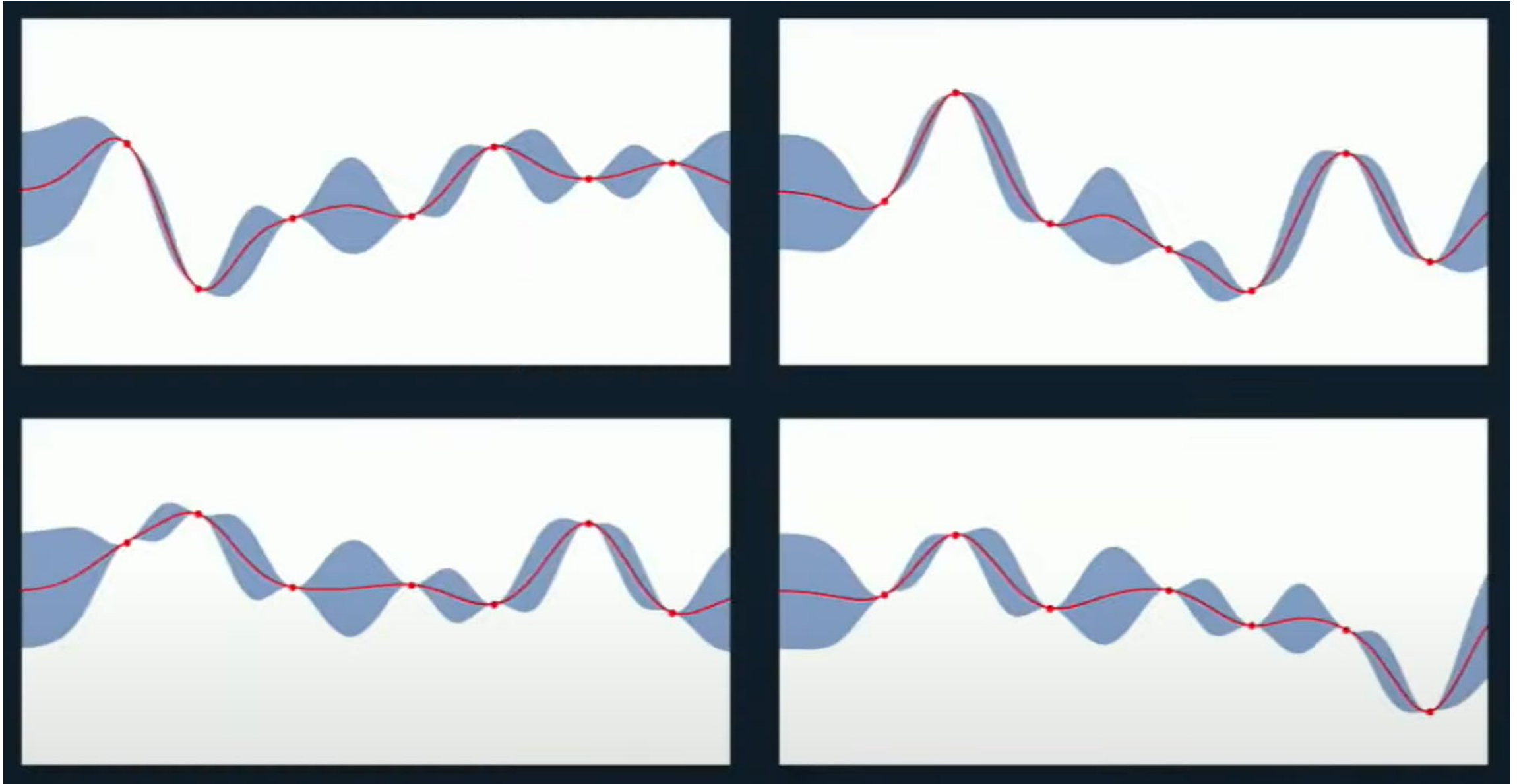
```
[0.4967141530112327, -0.1382640378102619]
```



```
x = [-2.1, -1.5, 0.3, 1.8, 2.5]
mu, s = conditional(x_more, x, y, θ)
y += np.random.multivariate_normal(mu, s).tolist()
y
```

```
[0.4967141530112327, -0.1382640378102619, -1.5128756, 0.52371713,
-0.13952425, -0.93665367, -1.29343995]
```

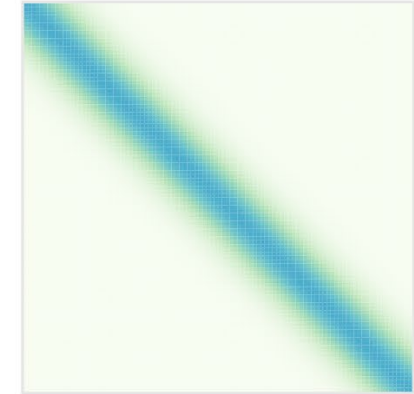
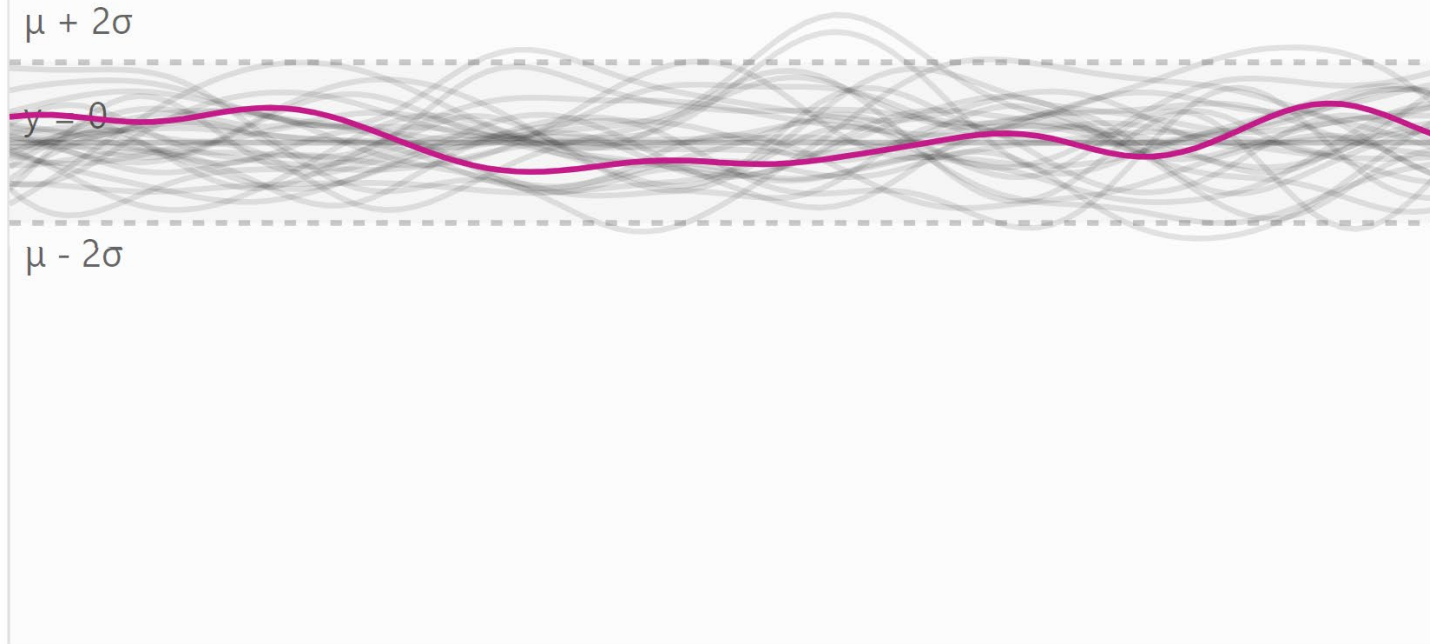

We could create as many different functions as we like by sampling different parameters



Nice animation of the sampling of distributions

☒ RBF ☐ Periodic ☐ Linear

(click to pause)



Variance σ = 0.8



Length l = 0.8



We now use these priors to calculate posteriors given some data observations

Recall Bayes' formula

$$p(\theta|y) \propto \prod_{i=1}^N p(y_i|\theta) p(\theta)$$

For Gaussian process, this becomes

$$\textit{Closed form GP} = \textit{Gaussian data} * \textit{GP prior}$$

When we fit the model, we are learning the kernel and mean hyperparameters

Predictions are done using Bayes' formula as well

Posterior predictive distribution

$$p(y^{new}|y) \propto \int p(y^{new}|\theta) p(\theta|y) d\theta$$

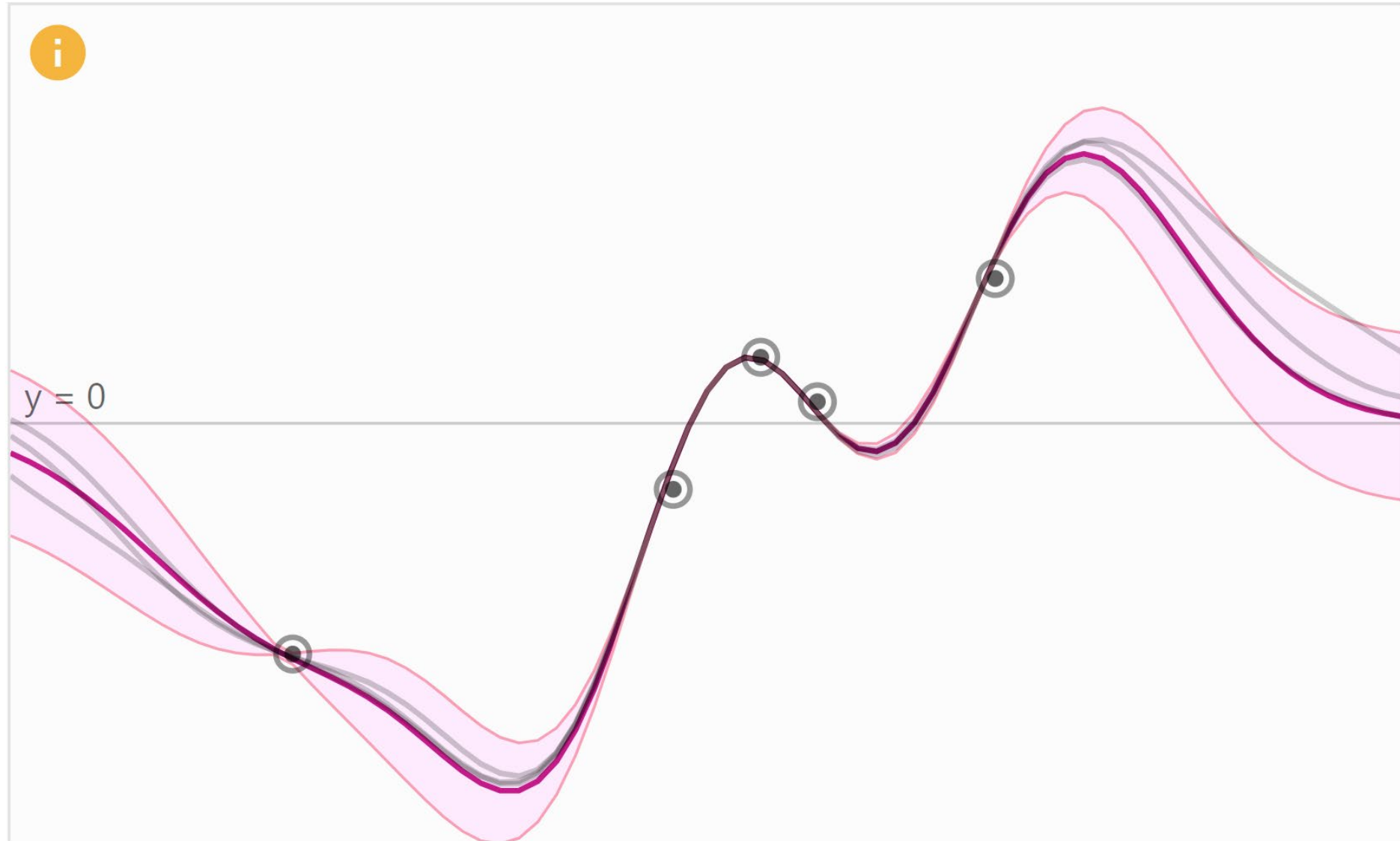
Likelihood is aleatoric uncertainty
Or noise in data itself

Posterior is epistemic uncertainty
Or our unknowns

Predictions are simply:

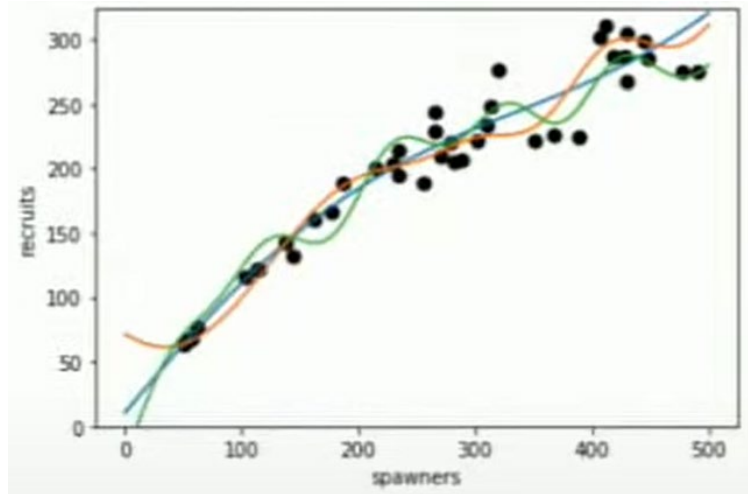
$$p(y^{new}|x^{new}, y, x) = N(\mu^{new}, \Sigma^{new})$$

Confidence intervals and predictions come from the distribution statistics μ, Σ

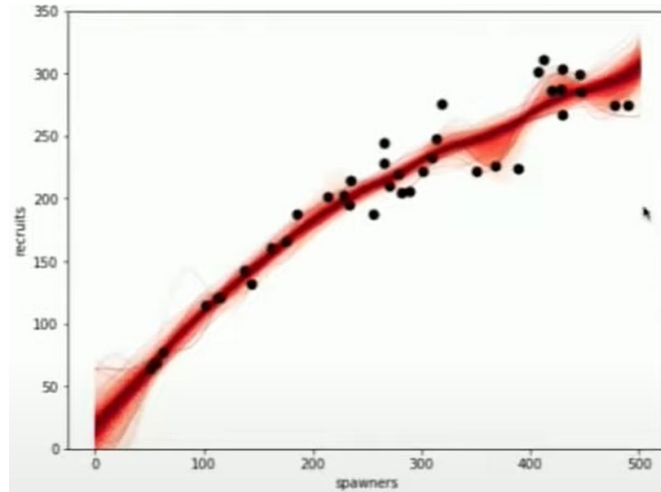


Confidence intervals and predictions come from the distribution statistics μ, Σ

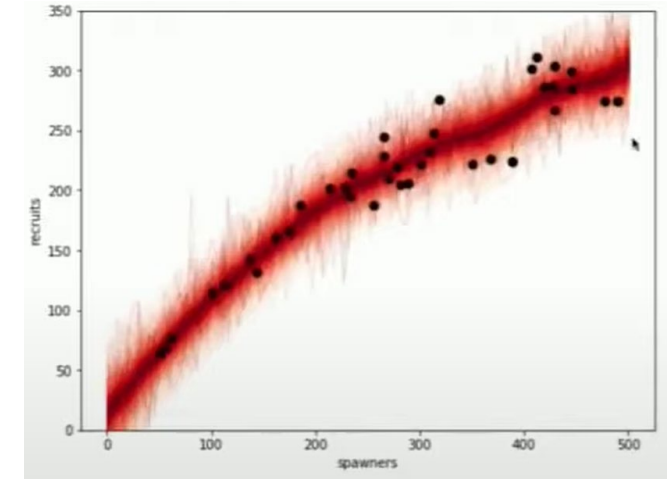
3 attempts at posteriors



100 attempts at posteriors



With prediction noise



Limitations of Gaussian Processes

Extrapolation does really poorly with Gaussian processes (depends on length scale of kernel)

They are not sparse (they use the whole samples/features information to perform the prediction)

They lose efficiency in high-dimensional spaces (features > a few dozen)

Picking and tuning kernel is tricky

