(3.) Mathematics.

Given $(x, f(x))$ to find $f(x)$ that passes through all the point.

Normally we find a function to parameters which fits over data well. and then define the function with the parameters.

We have data $D$ we find $\textcircled{w} \rightarrow$ parameth then we predict $y = w^T x$. prediction directly without limiting

In $GP$ instead we try to find the ~~function directly with~~ to one function.

So we find

$$\sqrt{P(y_*|x_*, D)} = \int_{\omega} P(y_*|x_*, w) \, P(\omega|D) \, d\omega$$

We assume that our data is continuous and follows gaussian distribution.

So we $P(y_*|x_*, w) \rightarrow$ gaussian and

$$P(\omega|D) = \frac{P(D|w) \, P(w)}{2}$$

$P(w) \rightarrow$ gaussian, $P(D|w)$
                              ↓
                           gaussian

and hence $P(\omega|D)$ will be gaussian.

~~these~~ Hence $P(y_*|x_*, D)$ will be gaussian.

Now we know $P(y_*|x_*, D)$ is gaussian so we assume our data with both test and training data to be gaussian.

p.t.y.

Now $P([y_1, y_2 \cdots y_n \, y_*]|[x_1, x_2 \cdots x_n \, x_*]) \sim N(\mu, \Sigma)$

we can normalise the data to make the mean $0$

find $P([\begin{smallmatrix} y_1 \\ \vdots \\ y_n \end{smallmatrix}]|[x_1, \cdots x_n], x^*) \sim N(0, \Sigma)$

where $\Sigma$ is the covariance matrix the ~~too~~ In linear case ~~$\sharp$~~ to find covariance matrix we how to do cholesky decomposition of the variance matrix cholesky decomposition decomposes the matrix into $V^T d$ ~~$\sharp$~~ $V$ which is kind of taking square root.

$$X \longrightarrow \sqrt{V^T} \sqrt{V} \quad A \ X \longrightarrow V^T V$$

$X$ can be decomposed into $V^T V$.

Now

$X$ needs to be positive semi-definite to do cholesky decomposition because

Now $\qquad X \longrightarrow V^T V$

$$\star \quad y^T x y \longrightarrow y^T V^T V y$$

$$y^T x y \longrightarrow (Vy)^T V y$$

$$y^T x y \longrightarrow ||V y||^2$$

$$y^T x y \geq 0$$

Hence $X$ is positive semi definite.

the covariance matrix can be is a kernel because it is positive semi definite. so We can define the covariance matrix by a kernel function. We can use RBF kernel function as it is infinite dimensional kernel.

$$RBF \ kernel = K(x, z) = e^{-\frac{||x-z||^2}{2\sigma^2}} \quad \sigma \ is \ a \ free \ parameter$$

So far we have

$$P\left(\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \\ y^* \end{bmatrix} \Bigg| \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \\ x^* \end{bmatrix}\right) = N(0, \Sigma) \quad when$$

$$\Sigma = k =$$
$$k_{ij} = k(x_i, x_j)$$

We now need to

But we know $[y_1, y_2 - y_n]$ given in our data and we need to find $y^*$ at our training text. point.

$$P(y^* | y_1, y_2 - y_n, x, - x_n, x^*).$$

pf

By conditioning

$$p(y^* | y_1, y_2 - y_n \, x_1, x_2 \ldots x_n x^*) = N(\mu_{y^*|(y_1, \ldots x_n x^*)}, \Sigma_{y^*|(y_1, \ldots x_n x^*)})$$

By conditioning identity.

$$\mu_{y^* | (y_1 \ldots x_n x^*)} = K_*^T K^{-1} y$$

where $K_*$ is the row covariance matrix b/w $x$'s and $x$
Training points and Testing point.

$K$ is the covariance matrix/kernel of Training data
matrix.

$y$ is our Training data.

and

$$\Sigma_{y^*|(y_1 \ldots y_n x_1 \ldots x_n x^*)} = K_{**} - K_*^T K^{-1} K_*$$

where $K_{**}$ is the covariance matrix of of Testing point.

This completes the GP model.

---

We can define how well our model extrapolates
by the variance of the extrapolated points, if
the variance is very large then we say that
it is poorly extrapolated.
Yes we need to optimize it, we can optimize
by choosing the right parameters of the
kernel RBF kernel.