

Applications

- (Bayesian) Regression where we have joint distribution over multiple predictions
- + Optimizing functions for which gradients are not available?
 - E.g. Hyper-parameter optimization of deep models

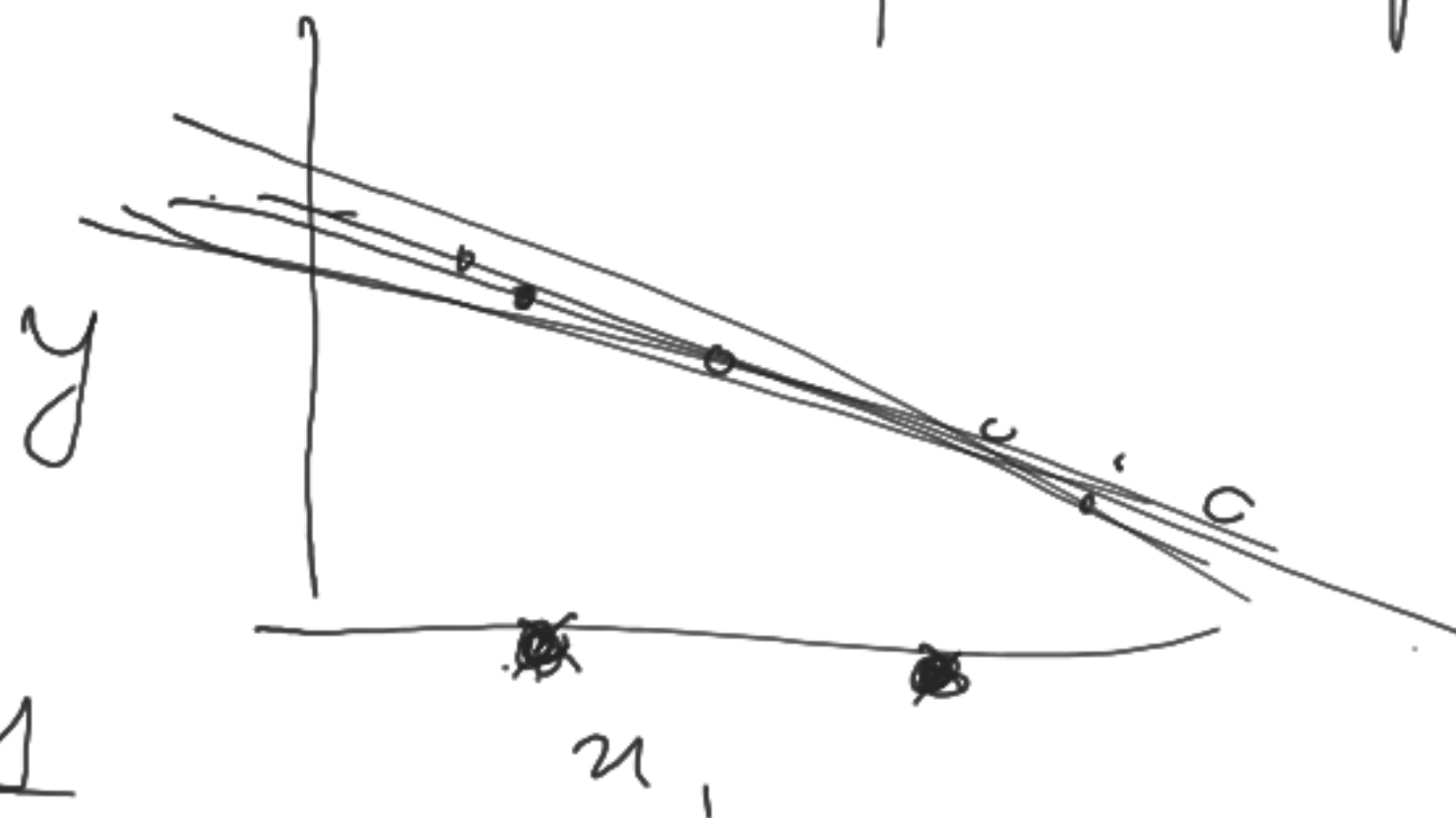
Normal Regression

$$f(x) = w_1 x_1 + w_2 x_2 + \dots + w_d x_d + b$$

$x \in \mathbb{R}^d$

w_1, \dots, w_d, b are parameters

$$y \sim \mathcal{N}(f(x), \sigma^2)$$



Gaussian Processes:

$$x \in \mathbb{R}^d$$

$$f(\cdot)$$

$$x^1, x^2, x^3, \dots, x^N$$

$$P \left(\begin{bmatrix} f(x^1) \\ f(x^2) \\ \vdots \\ f(x^N) \end{bmatrix} \right)$$

$$\sim \mathcal{N} \left(\begin{bmatrix} \underline{f(x^1)} & \underline{f(x^2)} & \dots & \underline{f(x^N)} \end{bmatrix} \sigma^2 \right)$$

$$\begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$k(x^1, x^1)$$

$$k(x^i, x^j)$$

$$k(x^N, x^N)$$

$$P \left(\begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} \right) \sim \mathcal{N}(\mu; \Sigma) =$$

$$\Sigma_{ij} = k(x^i, x^j)$$

$$\mu_{N \times 1}$$

$$\Sigma_{N \times N}$$

$$N\left(\underbrace{\begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}}_Y; \mu, \Sigma_{N \times N}\right) = \frac{1}{(2\pi)^{k/2}} e^{-\frac{(Y-\mu)^T \Sigma^{-1} (Y-\mu)}{2}}$$

Kernel function:

RBF kernel:

$$e^{-\frac{\|x^i - x^j\|^2}{2l^2}}$$

$l \leftarrow$ length parameter

$x^i, x^j \in \mathbb{R}^d$

Properties of multivariate Gaussians.

$$\textcircled{1} Y_A \sim N(\mu_A; \Sigma_{AA}) ; Y_B \sim N(\mu_B; \Sigma_{BB}) \quad Y_A \perp Y_B$$

$$Y_A + Y_B \sim N(\mu_A + \mu_B; \Sigma_{AA} + \Sigma_{BB})$$

$$Y \equiv \begin{bmatrix} Y_A \\ Y_B \end{bmatrix} \sim N(\mu; \Sigma) \equiv N\left(\begin{bmatrix} \mu_A \\ \mu_B \end{bmatrix}; \begin{bmatrix} \Sigma_{AA} & \Sigma_{AB} \\ \Sigma_{BA} & \Sigma_{BB} \end{bmatrix}\right)$$

$$Y \equiv \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \begin{matrix} Y_A \\ \\ \\ Y_B \end{matrix}$$

$$P(Y_A | Y_B = \underline{0}_B) \equiv N(\mu_{A|B}; \Sigma_{A|B})$$

$$\begin{aligned} \mu_{A|B} &\equiv \mu_A + \Sigma_{AB} \Sigma_{BB}^{-1} (\underline{0}_B - \mu_B) \\ \Sigma_{A|B} &\equiv \Sigma_{AA} - \Sigma_{AB} \Sigma_{BB}^{-1} \Sigma_{BA} \end{aligned}$$

Posterior distribution of the function:

$$Y \equiv \begin{bmatrix} y_1 \\ \vdots \\ y_N \\ y^* \end{bmatrix} \quad \begin{bmatrix} x \\ \vdots \\ x^N \\ x^* \end{bmatrix} \times P(y^* | \left[\begin{array}{l} f(x^1) = y_1 \\ f(x^2) = y_2 \\ \vdots \\ f(x^N) = y_N \end{array} \right] \text{ training data}, x^*)$$

$N(\mu_*, \sigma_*^2)$ where

$$\mu_* = 0 + k(x^*, x) [K(x, x)]^{-1} \vec{y}$$

$$\sigma_*^2 = \underline{k(x^*, x^*)} - k(x^*, x) [\underline{K(x, x)}]^{-1} k(x, x^*)$$

Le's say we want to get distribution over

x^{N+1}

\vdots

\vdots

x^{N+M}

x^*

$P(Y^*$

$m \times 1$

$f(x^1) = y_1$

$f(x^N) = y_N$

$\sim N(u^*, \Sigma^*)$

$$u^* = 0 + K(x^*, x) [K(x, x)]^{-1} Y$$

$$\Sigma^* \approx K(x^*, x^*) - K(x^*, x) [K(x, x)]^{-1} K(x, x^*)$$

Hyper-parameter optimization of an expansion deep model M

$X \equiv$ space of hyper-parameters -

eg. # layers, learning rate, # dimensions of hidden unit.
vocab. size.

$k(x^i, x^j)$ for .

$f(x)$ \equiv validation loss ^{or error} of M with hyper-parameters x .

Goal: find the hyper-parameters for which $f(x)$ is minimum.

$f(x)$ is not differentiable.

For more on use of GPs for hyper-parameter optimization see: <https://www.cs.cornell.edu/courses/cs4787/2019sp/notes/lecture16.pdf>

