

Lecture 12



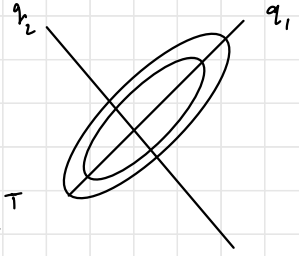
• Gaussian Distribution:

$$\begin{pmatrix} y_1 \\ \vdots \\ y_i \\ \vdots \\ y_n \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} 0 \\ \vdots \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \Sigma \right)$$

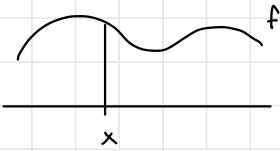
\downarrow finite index

$$\Sigma = Q \Lambda Q^T$$

$$\Sigma \geq 0$$



• Gaussian Process:



noiseless $y = f(x)$

$$y = \begin{matrix} \downarrow \\ x \\ \downarrow \end{matrix} \begin{pmatrix} y_x \end{pmatrix} = \begin{pmatrix} f(x) \end{pmatrix} = f \sim \mathcal{GP} \left(\begin{matrix} 0 \\ \vdots \\ 0 \end{matrix}, x + \sigma^2 K(x, x') \right)$$

\downarrow continuous index

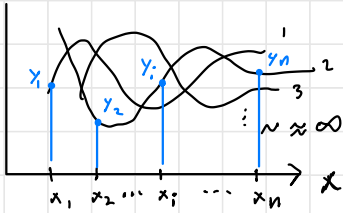
$$K \geq 0$$

$$K = Q \Lambda Q^T \quad \text{Mercer Thm.}$$

$$\begin{matrix} x' \\ \downarrow \\ \square \\ \leftarrow x \end{matrix} K(x, x') \begin{matrix} \uparrow \\ x' \end{matrix} = \begin{matrix} 1 \ 2 \ \dots \ k \ \dots \\ \square \\ \leftarrow x \end{matrix} \begin{pmatrix} \lambda_1 \\ \vdots \\ \lambda_k \end{pmatrix} \begin{matrix} \uparrow \\ x' \\ \square \\ \leftarrow x \end{matrix}$$

• Prior: $f \sim GP(0, \sigma^2 K)$

• Population:



$$y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_i = f(x_i) \\ \vdots \\ y_n \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} 1 \\ 2 \\ \vdots \\ i \\ \vdots \\ n \end{pmatrix}, \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}, \begin{matrix} 1 & 2 & \dots & j & \dots & n \\ \vdots & \vdots & & \downarrow & & \vdots \\ \vdots & \vdots & & \rightarrow \sigma^2 K_{ij} & & \vdots \\ \vdots & \vdots & & & & \vdots \\ n & & & & & \end{matrix} \right)$$

$$K_{ij} = k(x_i, x_j)$$

• Gaussian Process for Learning:

Training Data

1	x_i^T	$y_i = f(x_i)$
\vdots		
i		
\vdots		
n		
query 0	x_0^T	$y_0 = f(x_0)$

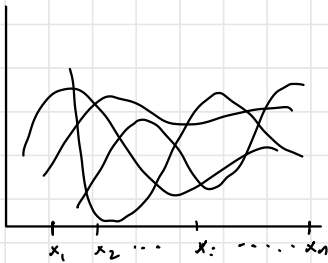
• Bayesian Inference:

$$[f \mid f(x_i) = y_i, i = 1, \dots, n]$$

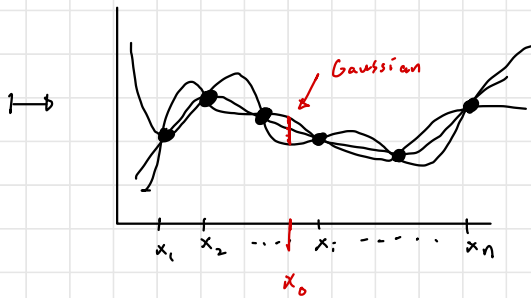
Among $N \approx \infty$ curves

$M \ll N$ curves with $f(x_i) = y_i, i = 1, \dots, n$

Prior



Posterior



(e.g. temperatures @
different locations)
spatial statistics

" Learning is memory + interpolation "

$$Y \sim \mathcal{N} \left(\begin{array}{c|c} 0, & \begin{array}{c|c} \sigma^2 K_{n \times n} & \sigma^2 k_0^T \\ \hline \sigma^2 k_0 & k_{00} \end{array} \end{array} \right)$$

$$k_{0j} = K(x_0, x_j)$$

So

$$[y_0 | y_1, \dots, y_n] \sim \mathcal{N}(\sigma^2 k_0, (\sigma^2 K + \sigma^2 I_n)^{-1} y, \text{variance})$$

$$+ \sigma^2 I_n$$

$$\text{if noise } y_i = f(x_i) + \varepsilon_i$$

$$\bullet \mathbb{E}(y_0 | Y) = \sigma^2 k_0 (\sigma^2 K + \sigma^2 I_n)^{-1} Y$$

$$= k_0 (K + \underbrace{\frac{\sigma^2}{\sigma^2}}_{\lambda} I_n)^{-1} Y$$

$$\sim k_0 (K + \lambda I_n)^{-1}$$

$$= \begin{array}{c} i \\ \boxed{K(x_i, x_0)} \end{array} \begin{array}{c} c \\ \boxed{c_i} \end{array}$$

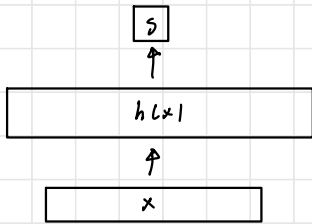
$$\hat{f}(x_0) = \sum_{i=1}^n c_i K(x_i, x_0)$$

• Another Interpretation :

$$\hat{f}(x_0) = \sum_{i=1}^n c_i \cdot k(x_i, x_0)$$

↑ affinity
↓ value
↓ key
↓ query

• Supervised Learning :



$$\sigma = f(x) = h(x)^T \beta$$

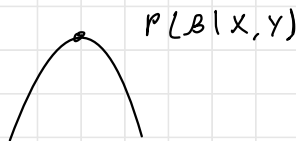
Kernel regression :

$$\min_{\beta} \sum_{i=1}^n (y_i - h(x_i)^T \beta)^2 + \lambda |\beta|^2$$

Gaussian Process :

$$\begin{aligned}
 [\beta | Y, X] &\propto P(\beta) \prod P(y_i | x_i, \beta) \\
 &\propto \exp \left(-\frac{1}{2\sigma^2} |\beta|^2 + \frac{1}{\sigma^2} (Y - h(X)^T \beta)^T \right) \\
 &\propto \frac{1}{2\sigma^2} \left(\sum_{i=1}^n (y_i - h(x_i)^T \beta)^2 + \frac{\sigma^2}{\sigma^2} |\beta|^2 \right)
 \end{aligned}$$

Gives a posterior mode $\hat{\beta}$
 $\rightarrow f(x_0) = h(x_0)^T \hat{\beta}$



But for the Gaussian Distribution, the posterior expectation is the posterior mode. $E(\beta | X, Y) = \hat{\beta}$
 $E(f(x_0) = h(x_0)^T \beta | X, Y) = \hat{f}(x_0) = h(x_0)^T \hat{\beta}$

• Advantage of GP:

(1) Variance V

$$[y_i | y] \sim \mathcal{N}(\cdot, V)$$

uncertainty quantification

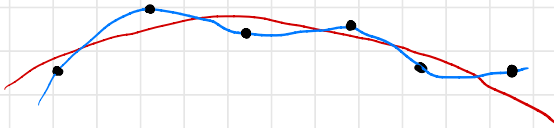
(2) hyper-parameter tuning

$$k(x, x') = \exp(-\gamma |x - x'|^2), \quad \gamma^2, \sigma^2$$

$$y \sim \mathcal{N}(0, \underbrace{\sum}_{\Sigma} k_y + \sigma^2 \mathbf{I})$$

$$\text{Likelihood}(\gamma^2, \gamma, \sigma^2) = p(y)$$

$$= \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} y^T \Sigma^{-1} y\right)$$



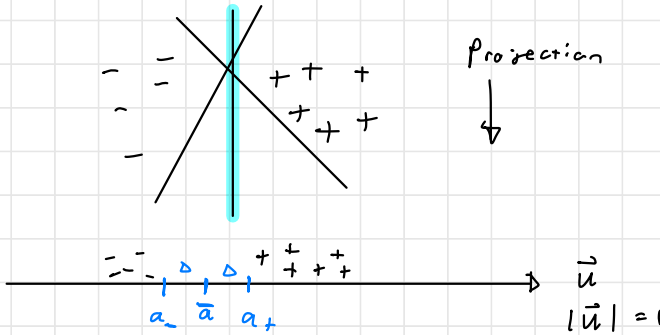
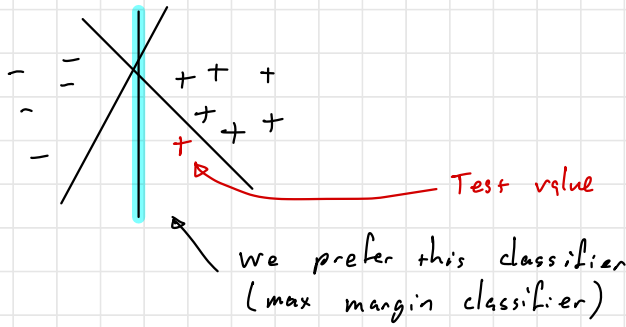
- Support Vector Machine (SVM)

- binary classification

$$s_i = \langle x_i, w \rangle + b$$

$$\text{Perceptron: } \hat{y}_i = \text{sign}(s_i) = \begin{cases} + & \text{if } s_i \geq 0 \\ - & \text{if } s_i < 0 \end{cases}$$

- Geometry (separable):



$\langle x_i, \vec{u} \rangle$ projection on \vec{u}

$$\text{Let } a_+ = \min_{i: y_i = +} \langle x_i, \vec{u} \rangle, \quad a_- = \max_{i: y_i = -} \langle x_i, \vec{u} \rangle$$

$$\bar{a} = \frac{a_+ - a_-}{2} \quad \Delta = a_+ - \bar{a} = \bar{a} - a_-$$

• We need to find \vec{u} that maximizes Δ

$$\langle x_i, \vec{u} \rangle \geq a_+ \quad \text{for } y_i = +1$$

$$\langle x_i, \vec{u} \rangle \leq a_- \quad \text{for } y_i = -1$$

" = " attainable.

• translate the optimization problem,

$$\langle x_i, \vec{u} \rangle \geq a_+ \quad \text{for } y_i = +1 \Rightarrow \langle x_i, \vec{u} \rangle - \bar{a} \geq \Delta$$

$$\langle x_i, \vec{u} \rangle \leq a_- \quad \text{for } y_i = -1 \Rightarrow \langle x_i, \vec{u} \rangle - \bar{a} \leq -\Delta$$

So,

$$\langle x_i, \frac{\vec{u}}{\Delta} \rangle - \frac{\bar{a}}{\Delta} \geq 1 \quad \text{for } y_i = +1$$

$$\langle x_i, \frac{\vec{u}}{\Delta} \rangle - \frac{\bar{a}}{\Delta} \leq -1 \quad \text{for } y_i = -1$$

$$\left. \begin{array}{l} \langle x_i, w \rangle + b \geq 1 \quad \text{for } y_i = +1 \\ \langle x_i, w \rangle + b \leq -1 \quad \text{for } y_i = -1 \end{array} \right\} \rightarrow y_i s_i = y_i (\langle x_i, w \rangle + b) \geq 1$$

Note: $|w| = \left| \frac{\vec{u}}{\Delta} \right| = \frac{1}{\Delta}$

- So $\Delta = \frac{1}{|w|}$

- To maximize Δ we must:

$$\begin{aligned} \min \quad & \frac{1}{2} |w|^2 \\ \text{s.t.} \quad & y_i (\langle x_i, w \rangle + b) \geq 1 \\ & i = 1, \dots, n \end{aligned}$$

Optimization Problem

• Lagrangian:

$$L(w, b, \alpha) = \underbrace{\frac{1}{2} |w|^2}_{\text{primal}} + \sum_{i=1}^n \underbrace{\alpha_i (1 - y_i (\langle x_i, w \rangle + b))}_{\text{dual}} \quad \text{Primal Problem}$$

$$\alpha = \begin{pmatrix} \alpha_1 \geq 0 \\ \vdots \\ \alpha_n \end{pmatrix}$$

$$\min_{(w, b)} \max_{\alpha > 0} L(w, b, \alpha)$$

↑
new constraint, much simpler

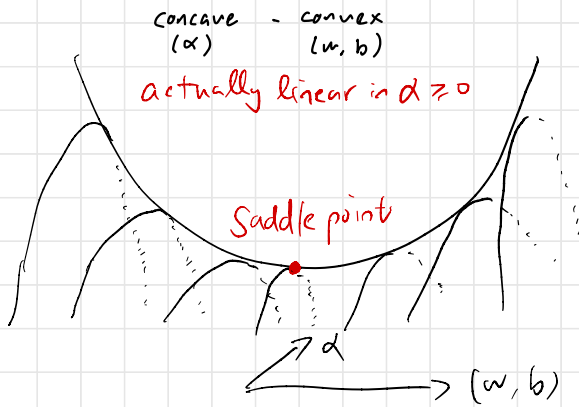
Claim: all old constraints $y_i (\langle x_i, w \rangle + b) \geq 1 \quad i=1 \dots n$ will be automatically satisfied

Proof by contradiction:

Otherwise $\exists i$ s.t. $y_i (\langle x_i, w \rangle + b) < 1$

$$\max_{\alpha_i > 0} \alpha_i (1 - y_i (\langle x_i, w \rangle + b)) \rightarrow \infty$$

↓
 ∞ can not be min (w, b) !



Saddle point

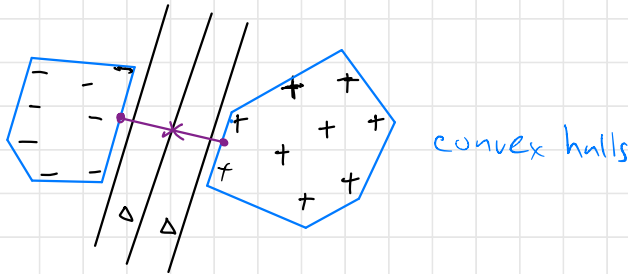
$$\min_{(w, b)} \max_{d \geq 0} = \max_{d \geq 0} \min_{(w, b)}$$

Von Neumann Zero sum game

$$\max_{\alpha \geq 0} \min_{(w, b)} \mathcal{L}(w, b, \alpha)$$

\downarrow closed form, representer
 $\max_{\alpha \geq 0} Q(\alpha)$ Dual Problem

o Geometry :



primal: max margin

dual: min distance

o non-separable: slackness relax to $\sum_{i=1}^n \xi_i$

$$\min \frac{1}{2} \|w\|^2 + c \sum_{i=1}^n \xi_i$$

\rightarrow hinge loss

s.t. $\langle x_i, w \rangle + b \geq 1 - \xi_i, \quad i=1, \dots, n$