



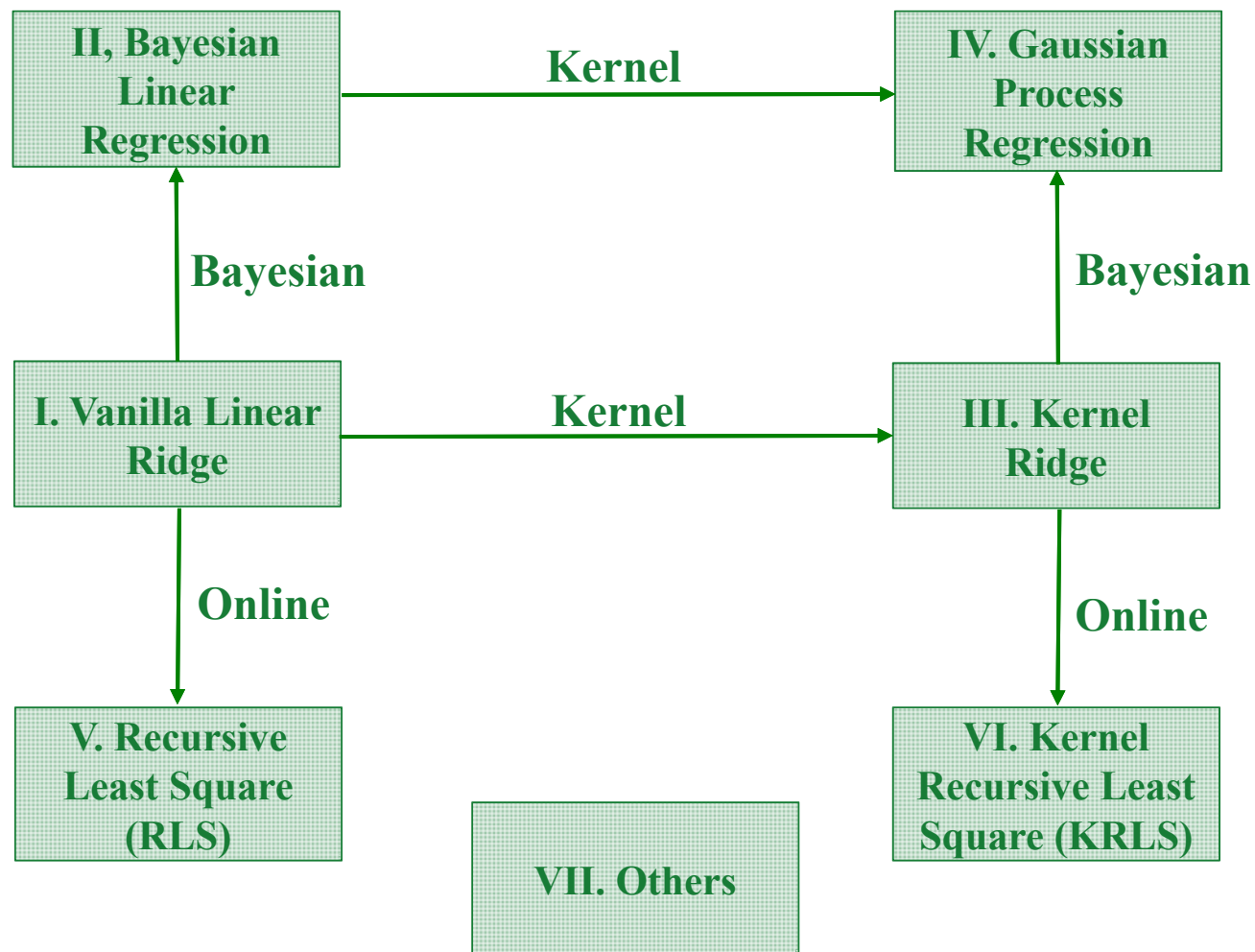
# **Advanced Linear Regression: Bayesian, Kernel and Online**

**(Jan 11<sup>th</sup>, 2017)**

**YANG Jiancheng**

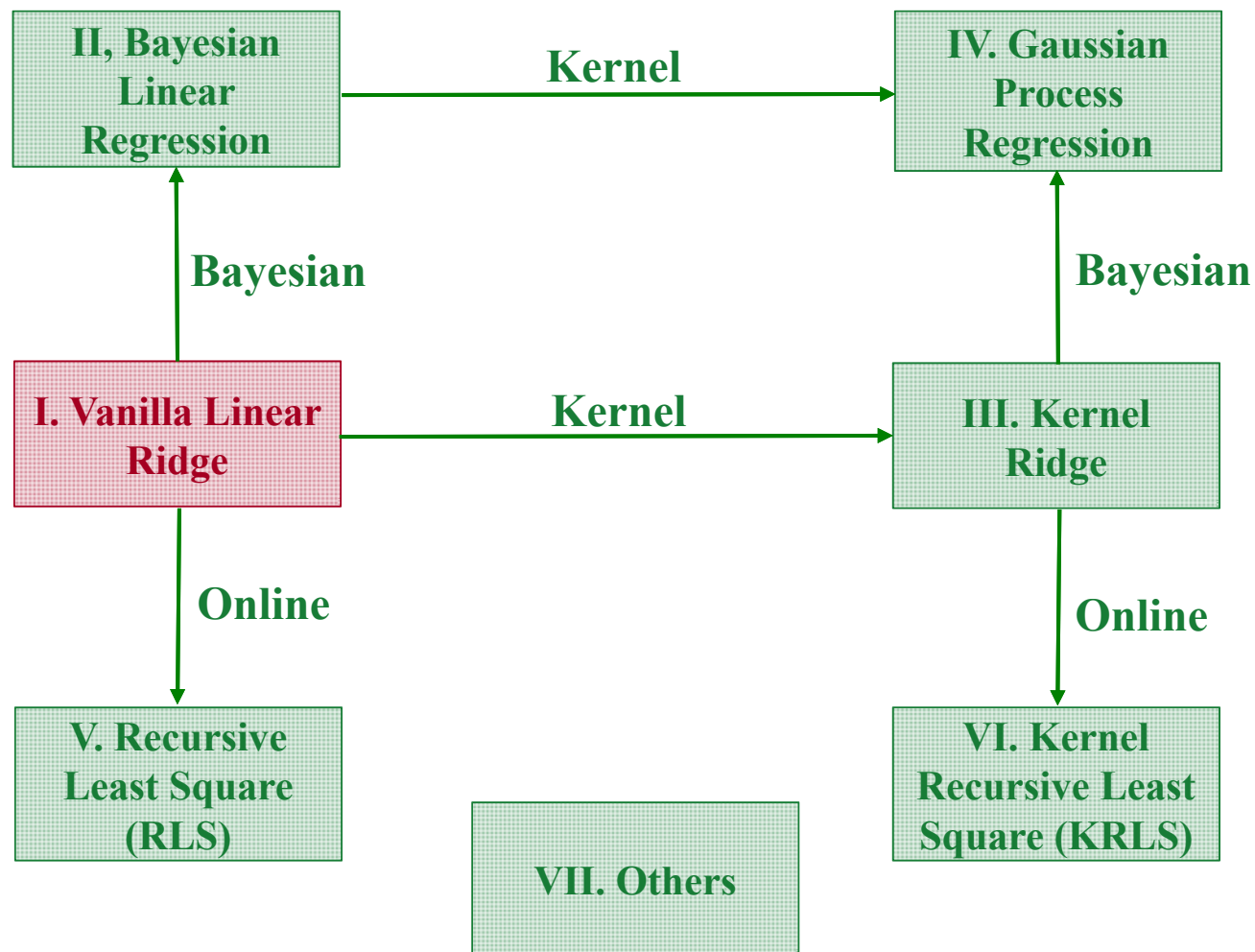


# Outline





# Outline





## • I. Vanilla Linear Ridge

- Matrix Form

$$\underset{n \times m}{Y} = \underset{n \times d}{X} \cdot \underset{d \times m}{W} + e$$

$d \rightarrow$  features

$m \rightarrow$  multiple outputs

$$\begin{aligned} \text{loss} &= \|e\|^2 = \|Y - XW\|^2 = (Y - XW)^T (Y - XW) \\ &= Y^T Y - Y^T XW - W^T X^T Y + W^T X^T XW \end{aligned}$$

$$\frac{\partial}{\partial W} \text{loss} = -2X^T Y + 2X^T XW = 0$$

$$W = (X^T X)^{-1} X^T Y = X^+ Y$$

$\hookrightarrow$  pseudo-inverse



- **I. Vanilla Linear Ridge**

- Ridge

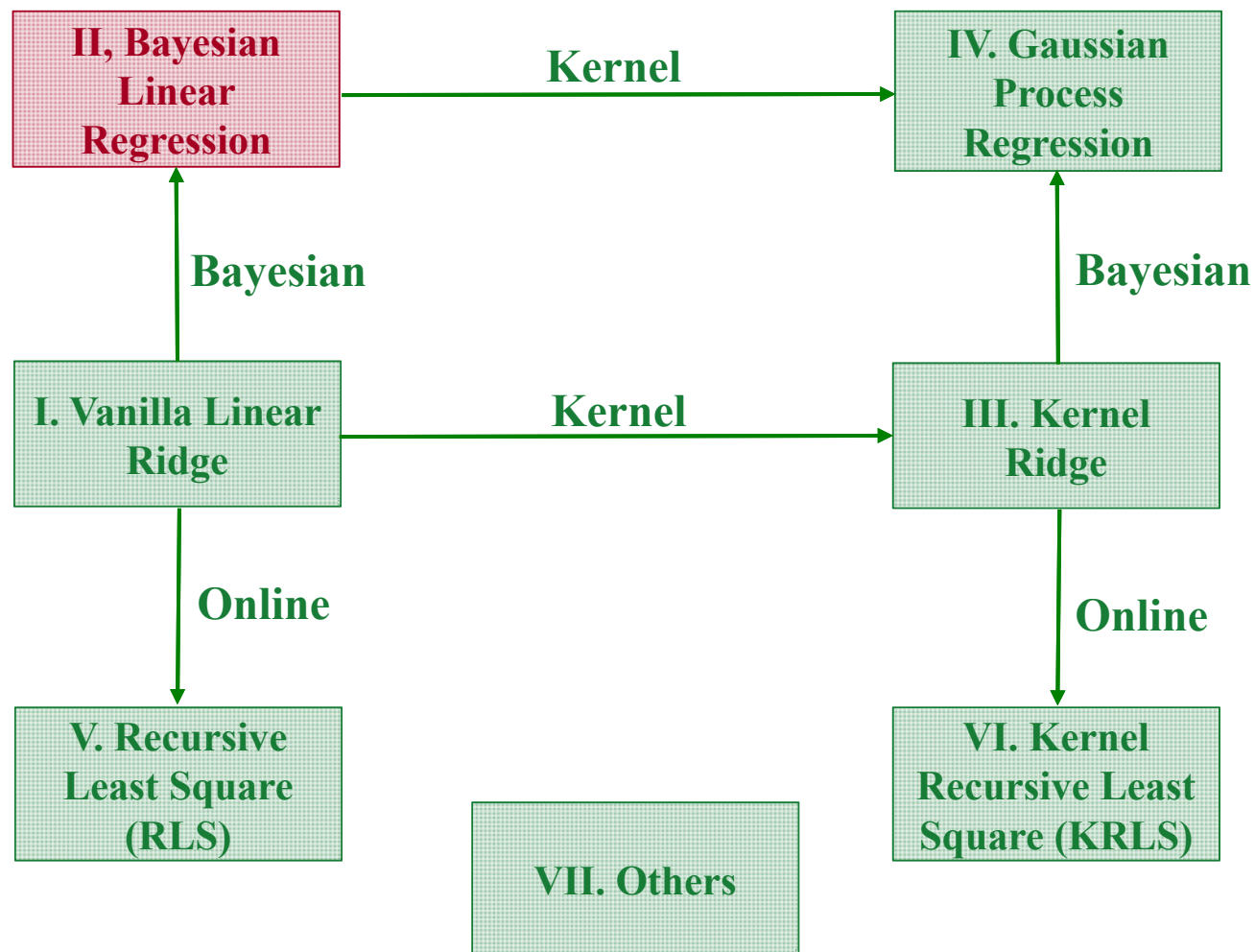
$$\text{loss} = \|y - Xw\|^2 + \lambda \|w\|^2$$

$$\frac{\partial}{\partial w} \text{loss} = -2X^T y + 2X^T X w + 2\lambda w = 0$$

$$w = (X^T X + \lambda I_d)^{-1} X^T y$$



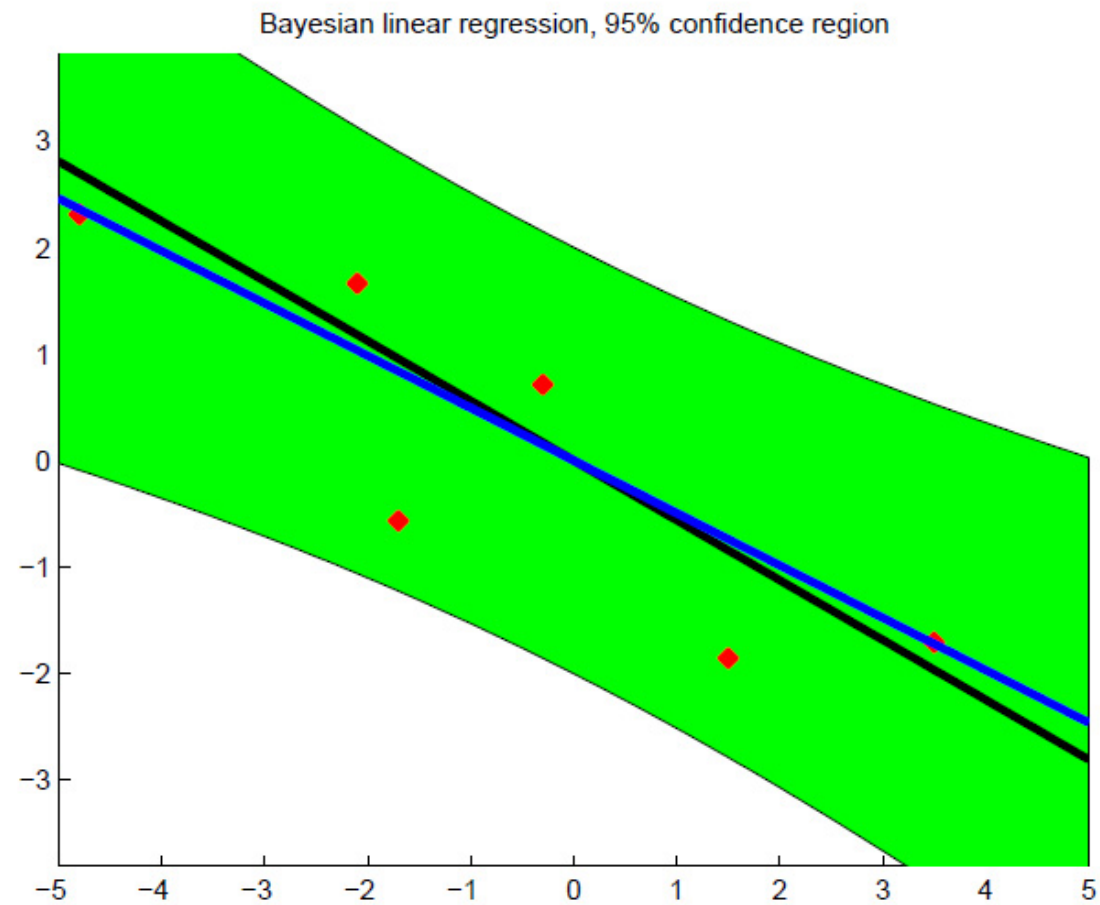
# Outline





## • II, Bayesian Linear Regression

- Effect of Bayesian





## • II, Bayesian Linear Regression

- Bayesian Interpretation

For simplicity,  $m=1$  (single target)

$$\underset{n \times 1}{Y} = \underset{n \times d}{X} \underset{d \times 1}{W} + \underset{n \times 1}{e}, \quad \boxed{e_i \sim N(0, \sigma^2)}$$

$$\Rightarrow Y \sim N(XW, \sigma^2 I_n)$$

Plus, we assume =

$$\boxed{W \sim N(0, \tau^2 I_d)}$$





## • II, Bayesian Linear Regression

- Posterior

$$\phi(w|x, y) \propto P(x, y|w) P(w)$$

$$\propto \exp\left(-\frac{(y - Xw)^T (y - Xw)}{2\sigma^2} - \frac{w^T w}{2\tau^2}\right)$$

↓ + Const

$$(w - \mu)^T \Lambda (w - \mu)$$

$$\Lambda_{d \times d} = \left(\frac{1}{\sigma^2} X^T X + \frac{1}{\tau^2} I_d\right)$$

$$\mu = \frac{1}{\sigma^2} \Lambda^{-1} X^T y = \left(X^T X + \frac{\sigma^2}{\tau^2} I\right) X^T y$$

$$P(w|D) \sim \mathcal{N}(\mu, \Lambda)$$

$$\text{MAP} : w = \mu$$

$\Leftrightarrow$  ridge



## • II, Bayesian Linear Regression

- Predictive (1)

Given a new sample  $(x_0, y_0)$   $\rightarrow ()_{dx_1} \rightarrow \text{Scalar}$

We want

$$\underline{P(y_0 | x_0, D)}, D = \{(x_1, y_1), \dots, (x_n, y_n)\}$$

Use Bayesian Average:

$$P(y_0 | x_0, D) = \int P(y_0, w | x_0, D) dw$$

$$= \int \underbrace{P(y_0 | x_0, D, w)}_{\substack{\downarrow \\ y_0 | D | w}} \underbrace{P(w | x_0, D)}_{\substack{\downarrow \\ P(w | D)}} dw$$

$$= \int N(w^T x_0, \sigma^2) N(w, \Lambda^{-1}) dw$$



## • II, Bayesian Linear Regression

### • Predictive (2)

$$p(y_0 | x_0, D) \propto \int \exp\left( \underbrace{\frac{(y_0 - w^T x_0)^2}{2\sigma^2} - \frac{1}{2} (w - \mu)^T L (w - \mu)}_{(w - m)^T L (w - m) + \dots} \right) dw$$

Our goal:

$$\int g(y_0) \cdot N(w | \dots) dw = g(y_0)$$

$$p(y_0 | x_0, D) \propto \int \exp\left( \frac{(w - m)^T L (w - m)}{2} \right) \underbrace{\exp\left( \frac{m^T L m}{2} - \frac{y_0^2}{2\sigma^2} \right)}_{\downarrow g(y_0)} dw$$

$$g(y_0) = \exp\left( -\frac{\lambda}{2} (y_0 - \hat{y}_0)^2 \right)$$

$$\hat{y}_0 = \mu^T x_0$$

$$\frac{1}{\lambda} = \sigma^2 + x_0^T \Lambda^{-1} x_0$$

$$L = \frac{x_0 x_0^T}{\sigma^2} + \Lambda$$

$$m = L^{-1} \left( \frac{y_0 x_0}{\sigma^2} + \Lambda \mu \right)$$



## • II, Bayesian Linear Regression

- Predictive (3)

$$N(y_0 | x_0, D) \sim N(\overset{\text{Same!}}{\mu^T x_0}, \sigma^2 + x_0^T \Lambda^{-1} x_0)$$

$$\varepsilon_i \sim N(0, \sigma^2)$$

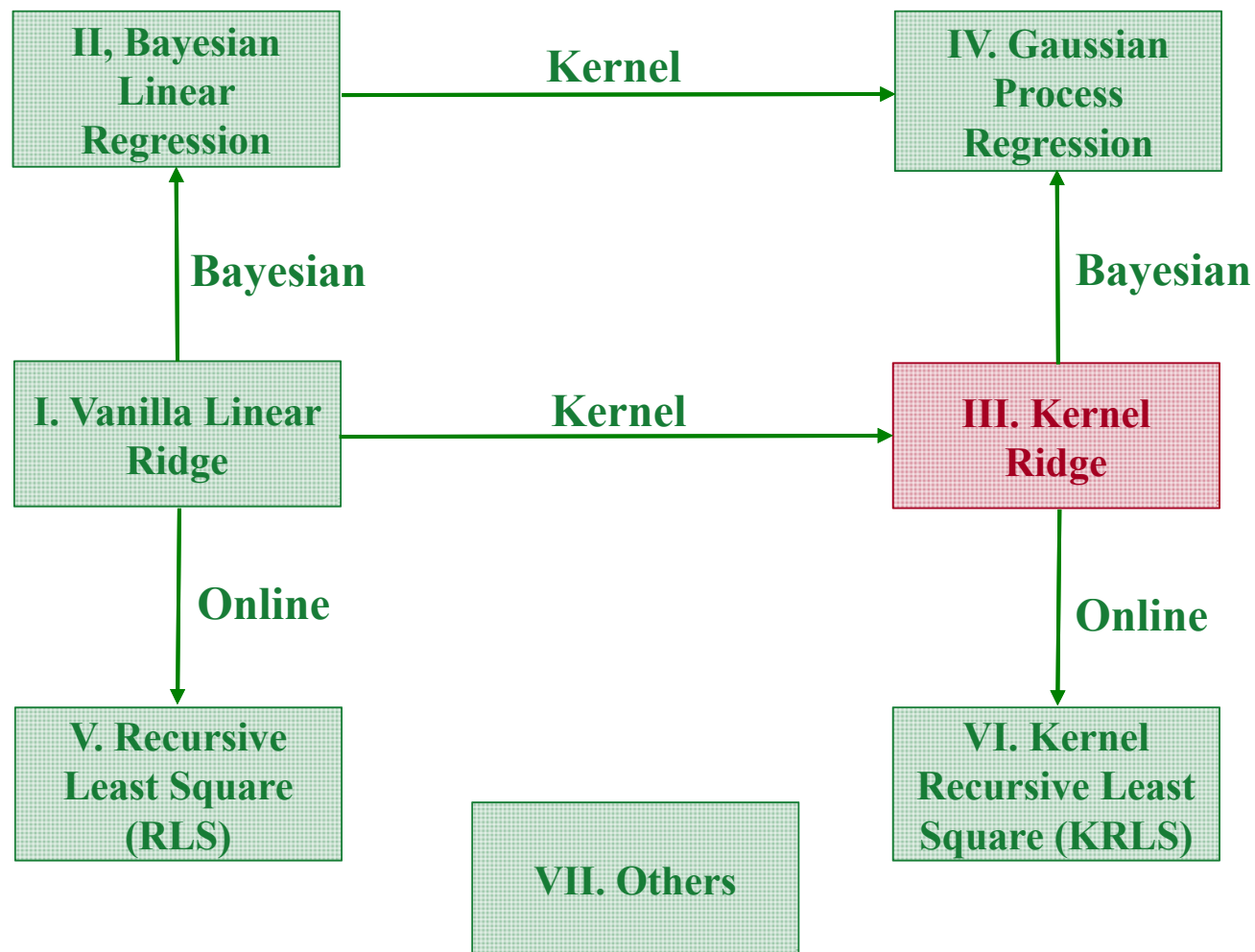
$$w \sim N(0, \tau^2 I_d)$$

$$\mu = \left( X^T X + \frac{\sigma^2}{\tau^2} I \right) X^T y$$

$$\Lambda = \frac{1}{\sigma^2} X^T X + \frac{1}{\tau^2} I_d$$



# Outline





### • III. Kernel Ridge

- Dual Form

$$(X^T X + \lambda I) w = X^T y$$

$$\lambda w = X^T (y - X w)$$

$$w = X^T \cdot \boxed{\frac{y - X w}{\lambda}} \quad \underline{\underline{\text{def}}} \quad w = X^T \alpha$$

the optimal  $w$  is

the linear combination of  $X$



### • III. Kernel Ridge

- Kernelized Ridge

$$w = x^T \alpha$$

def  $K$   
 $\nearrow$  or  $K = \phi^T(x) \phi(x)$

$$y = x w + \varepsilon = \underbrace{x \cdot x^T}_{K} \alpha + \varepsilon$$

$$L_2 = ||\varepsilon||_2 + \lambda ||w||_2 = (y - K\alpha)^T (y - K\alpha) + \lambda \underbrace{\alpha^T x x^T \alpha}_{K}$$

$$L_2 = y^T y - \alpha^T K^T y - y^T K \alpha + \alpha^T K^T K \alpha + \lambda \alpha^T K \alpha$$

$$\frac{\partial}{\partial \alpha} L_2 = -2K^T y + 2K^T K \alpha + 2\lambda K^T = 0$$

$$\Rightarrow \boxed{\alpha = (K + \lambda I_n)^{-1} y}$$



### • III. Kernel Ridge

- Big-O Analysis

For a new entry  $x$

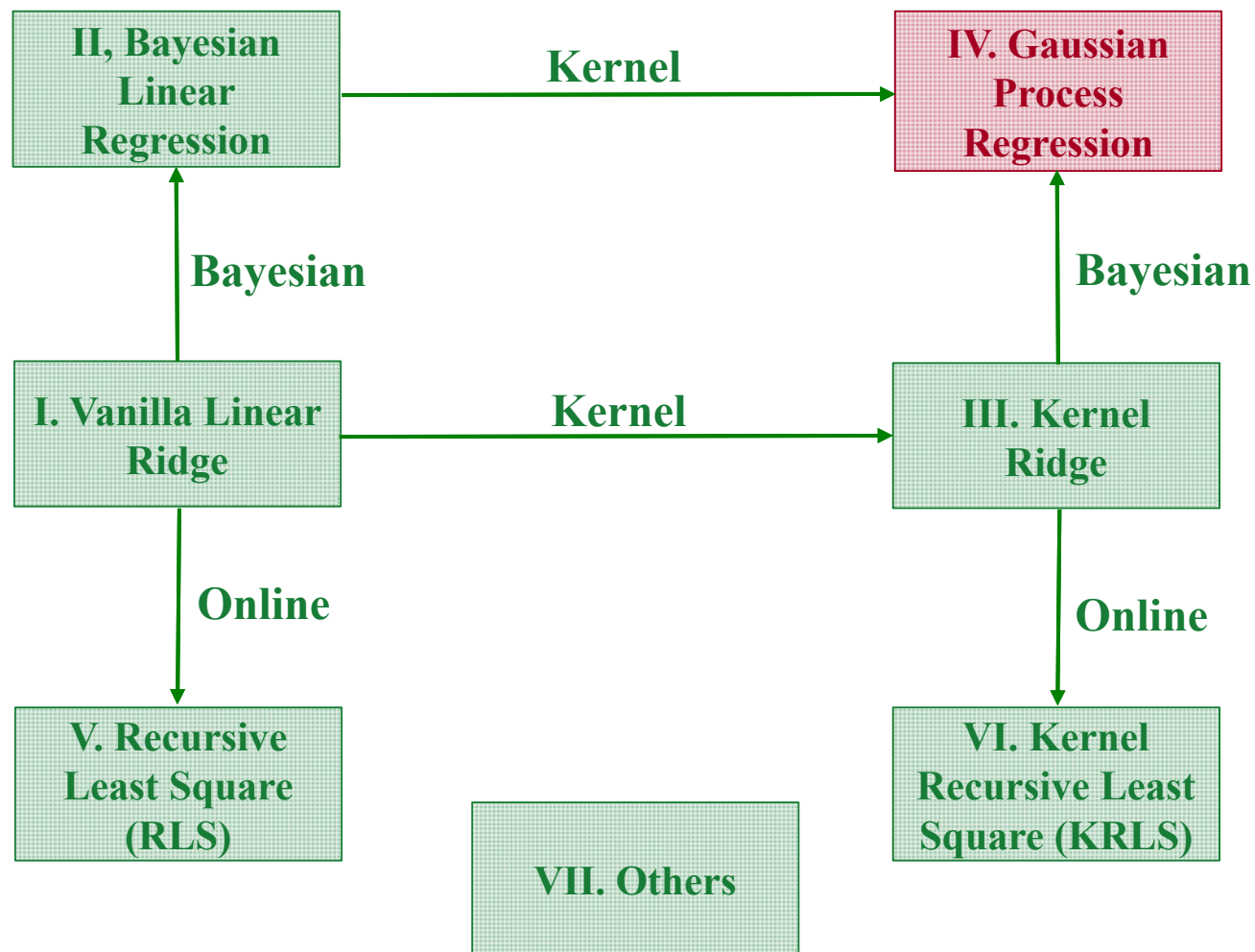
$$\underset{1 \times m}{\underline{y}} = \underset{1 \times n}{\underbrace{K(x, X)}} \cdot \underset{n \times m}{\underbrace{\underline{\alpha}}} + \varepsilon$$

	Linear Ridge	Kernel Ridge
Training Time	$O(d^3 + Nd^2)$	$O(N^3)$
Predicting Time	$O(d)$	$O(d)$
Model Size	$O(d)$	$O(d) + O(Nd) = O(Nd)$





# Outline





## • IV. Gaussian Process Regression

- Idea

### a) Kernel as Covariance function

$$k(x, x') = \sigma_f^2 \exp \left[ \frac{-(x - x')^2}{2l^2} \right] + \sigma_n^2 \delta(x, x')$$

$$K = \begin{bmatrix} k(x_1, x_1) & k(x_1, x_2) & \cdots & k(x_1, x_n) \\ k(x_2, x_1) & k(x_2, x_2) & \cdots & k(x_2, x_n) \\ \vdots & \vdots & \ddots & \vdots \\ k(x_n, x_1) & k(x_n, x_2) & \cdots & k(x_n, x_n) \end{bmatrix}$$

$$K_* = [k(x_*, x_1) \quad k(x_*, x_2) \quad \cdots \quad k(x_*, x_n)] \quad K_{**} = k(x_*, x_*)$$

### b) Sample from a multivariate Gaussian

$$\begin{bmatrix} y \\ y_* \end{bmatrix} \sim \mathcal{N} \left( 0, \begin{bmatrix} K & K_*^T \\ K_* & K_{**} \end{bmatrix} \right)$$

$$y_* | y \sim \mathcal{N}(K_* K^{-1} y, K_{**} - K_* K^{-1} K_*^T)$$

### c) Fit the **distribution of functions**



## • IV. Gaussian Process Regression

- Hyperparameter as Optimization

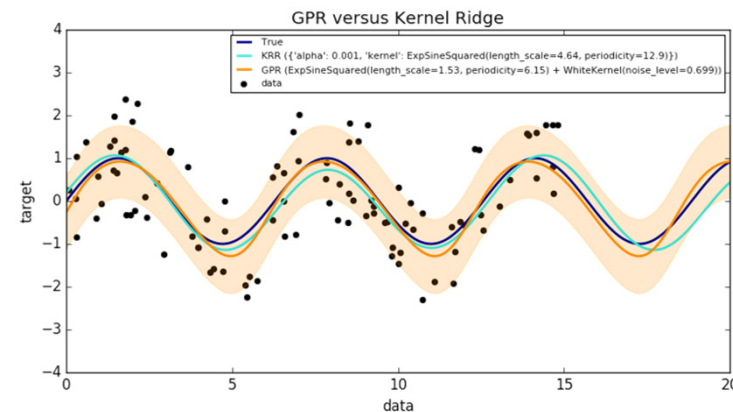
- a) Use Conjugate Gradients or other optimizer to optimize the likelihood => **Bayesian Optimization**

$$\theta = \{l, \sigma_f, \sigma_n\}$$

$$\log p(\mathbf{y}|\mathbf{x}, \theta) = -\frac{1}{2}\mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} - \frac{1}{2} \log |\mathbf{K}| - \frac{n}{2} \log 2\pi$$

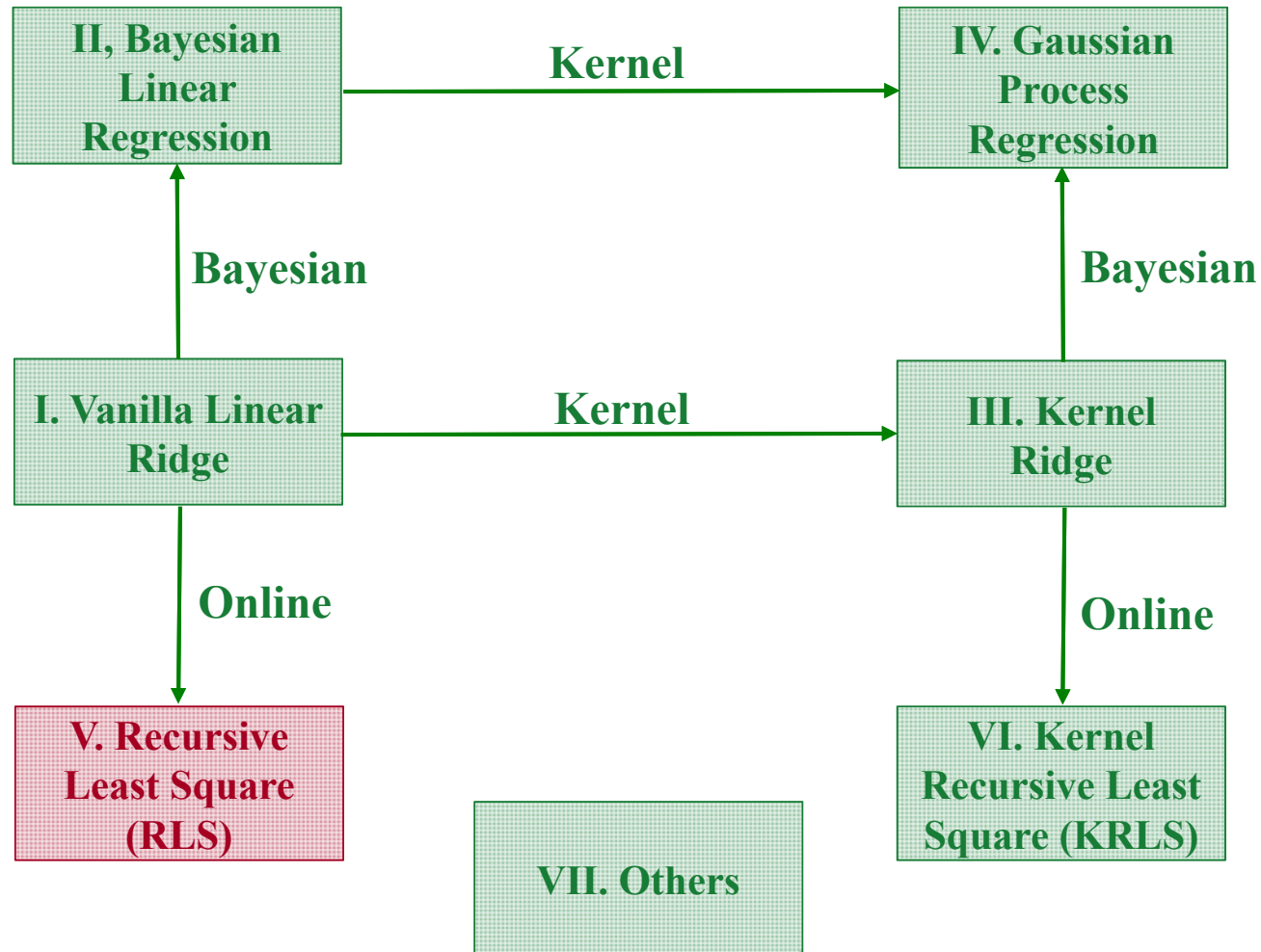
b) GPR and Kernel Ridge

- Time for KRR fitting: 10.243
- Time for GPR fitting: 0.223
- Time for KRR prediction: 0.083
- Time for GPR prediction: 0.091
- Time for GPR prediction with standard-deviation: 0.386





# Outline





## • V. Recursive Least Square (RLS)

### • Formalization (1)

$$y_N = x_N^T \cdot w_N + \varepsilon$$

$$w_N = \underbrace{(x_N^T x_N + \lambda I_d)^{-1}}_{S_N} x_N^T y_N$$

$$S_N = x_N^T x_N + \lambda I_d$$

$$S_{N+1} = x_{N+1}^T x_{N+1} + \lambda I_d = \begin{bmatrix} x_N^T & x_{N+1}^T \end{bmatrix} \begin{bmatrix} x_N \\ x_{N+1} \end{bmatrix} + \lambda I_d$$

$$\Rightarrow \boxed{S_{N+1} = S_N + x_{N+1}^T x_{N+1}}$$

$$S_N w_N = x_N^T y_N$$

$$S_{N+1} w_{N+1} = x_{N+1}^T y_{N+1} = x_N^T y_N + x_{N+1}^T y_{N+1}$$

$$\Rightarrow w_{N+1} = w_N + S_{N+1}^{-1} x_{N+1}^T (y_{N+1} - x_{N+1} \cdot w_N)$$



- V. Recursive Least Square (RLS)

- Formalization (2)

$$e_n = y_{n+1} - x_{n+1}^T w_n$$

$$S_{n+1} = S_n + x_{n+1} x_{n+1}^T$$

$$w_{n+1} = w_n + \underbrace{S_{n+1}^{-1} x_{n+1}^T}_{\text{gain}} e_n$$

Can it be better?

$$\textcircled{1} (A+B)^{-1} = A^{-1} - \frac{1}{1+g} A^{-1} B A^{-1}$$

$$g = \text{trace}(B A^{-1})$$

② trace trick

$$\text{tr}(uvw) = \text{tr}(wuv) = \text{tr}(vwu)$$



## • V. Recursive Least Square (RLS)

- Computation Optimization

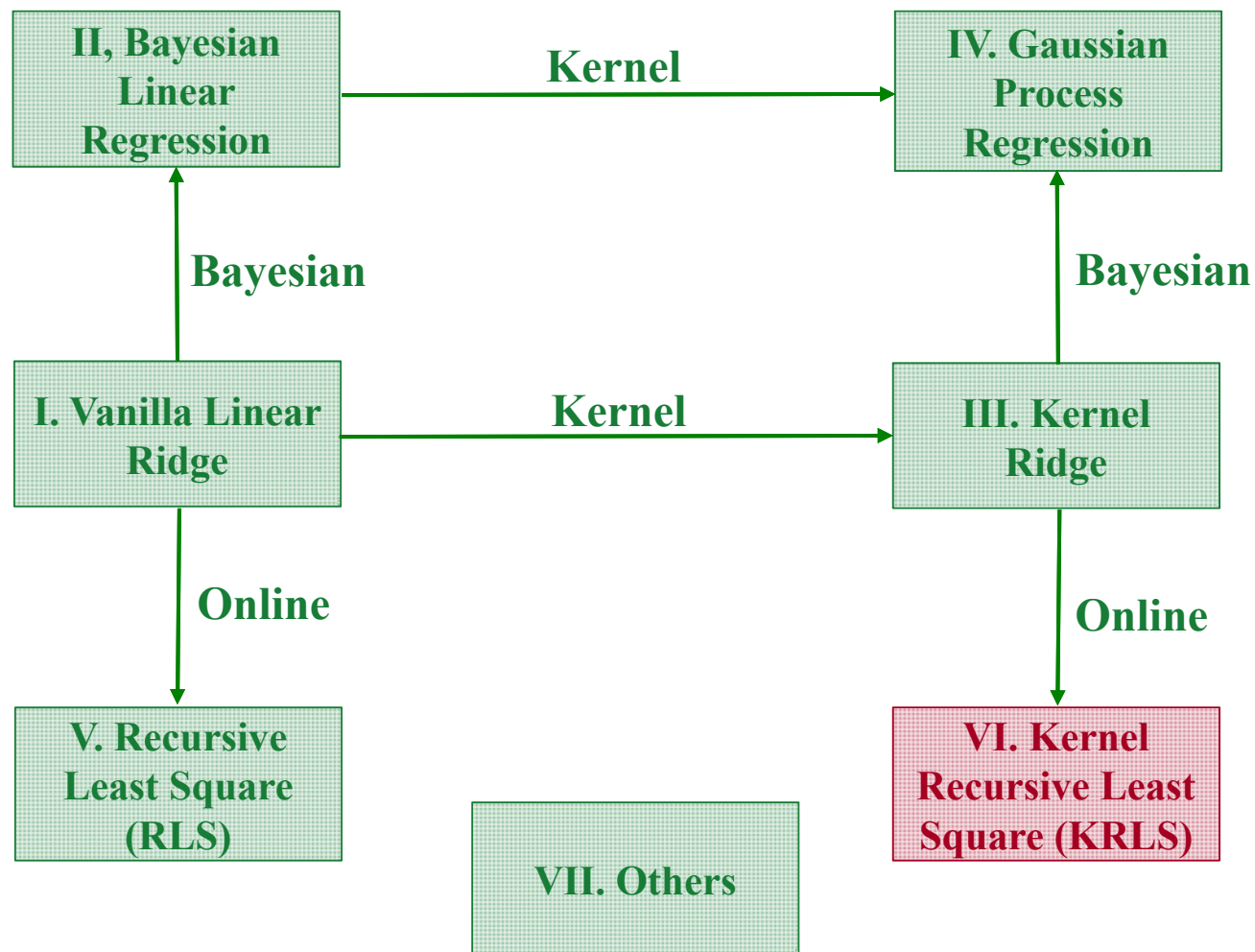
$$\begin{aligned} S_{N+1}^{-1} &= (S_N + x_{N+1}^T x_{N+1})^{-1} \\ &= S_N^{-1} - \frac{1}{1 + \text{tr}(x_{N+1}^T x_{N+1} S_N^{-1})} S_N^{-1} x_{N+1}^T \cdot x_{N+1} S_N^{-1} \end{aligned}$$

$O(d^3)$  to  $O(d^2)$

$$\begin{aligned} \text{tr}(x_{N+1}^T x_{N+1} S_N^{-1}) &= \text{tr}(x_{N+1} S_N^{-1} x_{N+1}^T) \\ &= x_{N+1} S_N^{-1} x_{N+1}^T \quad \text{scalar} \end{aligned}$$



# Outline







## • VI. Kernel Recursive Least Square (KRLS)

- Infinite Dictionary

- The core step is to get the inverse of new kernel  $O(N^2)$

$$\begin{aligned}\dot{\mathbf{K}} &= \mathbf{K} + c\mathbf{I}, & \dot{\mathbf{K}}_n &= \begin{bmatrix} \dot{\mathbf{K}}_{n-1} & \mathbf{k}_n \\ \mathbf{k}_n^\top & k_{nn} + c \end{bmatrix} \\ \mathbf{a}_n &= \dot{\mathbf{K}}_{n-1}^{-1} \mathbf{k}_n & \gamma_n &= k_{nn} + c - \mathbf{k}_n^\top \mathbf{a}_n \\ \dot{\mathbf{K}}_n^{-1} &= \frac{1}{\gamma_n} \begin{bmatrix} \gamma_n \dot{\mathbf{K}}_{n-1}^{-1} + \mathbf{a}_n \mathbf{a}_n^\top & -\mathbf{a}_n \\ -\mathbf{a}_n & 1 \end{bmatrix}\end{aligned}$$

- The new instance parameters will be easy to get  $O(N^2)$

$$\boldsymbol{\alpha}_{n-1} = \dot{\mathbf{K}}_{n-1}^{-1} \mathbf{y}_{n-1} \quad \hat{y}_n = \mathbf{k}_n^\top \boldsymbol{\alpha}_{n-1} \quad e_n = y_n - \hat{y}_n$$

$$\boldsymbol{\alpha}_n = \begin{bmatrix} \boldsymbol{\alpha}_{n-1} - \mathbf{a}_n e_n / \gamma_n \\ e_n / \gamma_n \end{bmatrix}$$

- There are many criteria for growing and pruning dictionary



## • VI. Kernel Recursive Least Square (KRLS)

- Dictionary Update

criterion	type	complexity
all	growing	—
coherence	growing	$\mathcal{O}(m)$
ALD	growing	$\mathcal{O}(m^2)$
oldest	pruning	—
least weight	pruning	$\mathcal{O}(m)$
least a posteriori SE	pruning	$\mathcal{O}(m^2)$

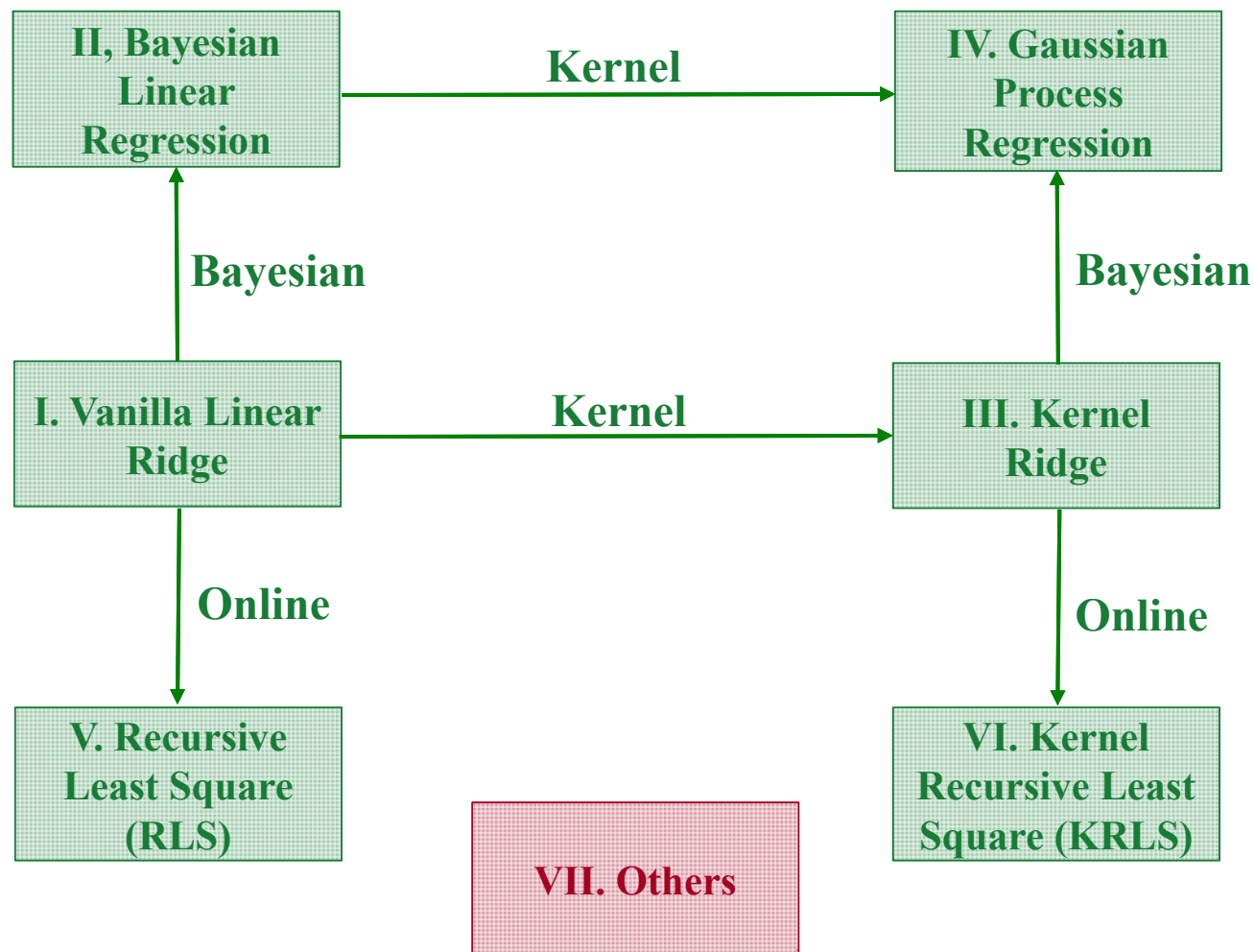


## • VI. Kernel Recursive Least Square (KRLS)

- Sliding-Window
- Very simple update to the dictionary: remove the oldest
- Suitable for Time Series
- Inverse Kernel
$$\dot{\mathbf{K}}_{n-1} = \begin{bmatrix} a & \mathbf{b}^T \\ \mathbf{b} & \mathbf{D} \end{bmatrix} \quad \dot{\mathbf{K}}_{n-1}^{-1} = \begin{bmatrix} e & \mathbf{f}^T \\ \mathbf{f} & \mathbf{G} \end{bmatrix}$$
$$\mathbf{D}^{-1} = \mathbf{G} - \mathbf{f}\mathbf{f}^T/e.$$
- Remove the oldest instance, and its instance parameter



# Outline





## • VII. Others

- Kalman Filter for Online Linear Regression

$\Theta_t$  is the **state**,  $y_t$  is the **observation** -

$$\begin{cases} \Theta_t = G_t \Theta_{t-1} + W_t \\ y_t = F_t^T \Theta_t + V_t \end{cases}$$

---

Online linear regression?

$$\Theta_t \leftarrow \beta_t, \quad y_t = \beta_t^T x_t + \varepsilon$$

$$\begin{cases} \beta_t = \overset{G_t}{\boxed{I}} \cdot \beta_{t-1} + W_t \\ y_t = \underset{F_t^T}{\boxed{x_t}} \cdot \beta_t + V_t \end{cases}$$



# Bibliography

- [Online regression with kernel](#)
- [Gaussian Processes: A Quick Introduction](#)
- [CS229 Lecture note of Gaussian Process](#)
- [Bayesian Linear Regression YouTube](#)
- [Recursive Least Square \(RLS\) Otexts](#)
- [Dynamic Hedge Ratio Between ETF Pairs Using the Kalman Filter](#)



**Thanks for listening!**

