

COMP90051

Statistical Machine Learning

Workshop Week 9

Xudong Han

https://github.com/HanXudong/COMP90051_2020_S1

Bayesian Regression

- Frequentist V.S. Bayesian
- Bayesian regression with known variance
- Bayesian model selection
- Bayesian regression with unknown variance



by yourself

Frequentist V.S. Bayesian

- Frequentist (MLE)

Generally reduces to minimizing the negative log-likelihood. Returns a point-estimate.

$$\theta_{MLE} = \operatorname{argmax}_{\theta} p(X|\theta) = \operatorname{argmax}_{\theta} \prod_i^n p(x_i|\theta) = \operatorname{argmax}_{\theta} \sum_i^n \log p(x_i|\theta)$$

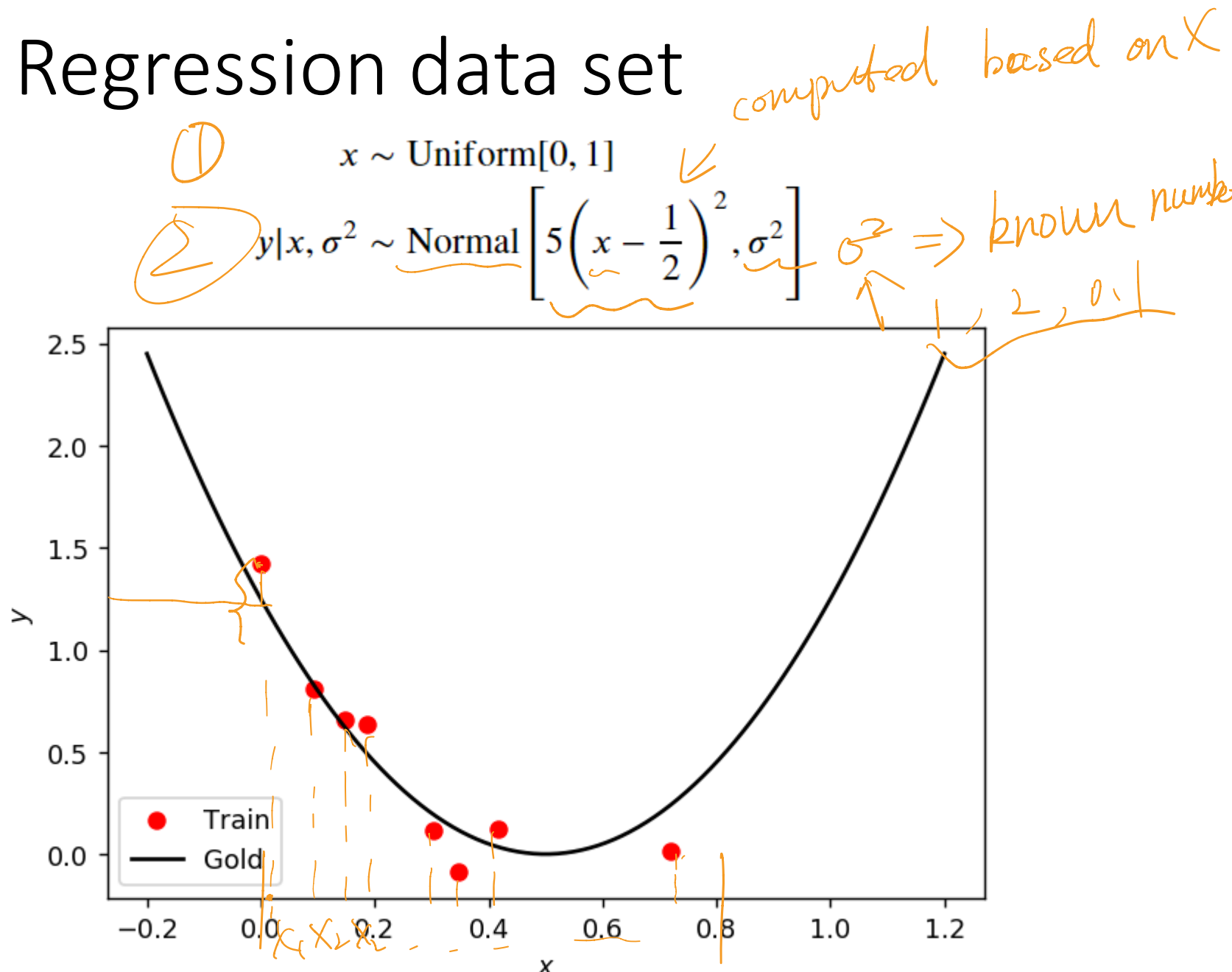
- Bayesian:

$$p(\theta|x) = \frac{\prod_i^n p(\theta|x_i)p(\theta)}{\int d\theta \prod_i^n p(\theta|x_i)p(\theta)}$$

Handwritten annotations for the Bayesian equation:

- $p(\theta|x)$: Posterior probability (circled in red)
- $\prod_i^n p(\theta|x_i)$: Likelihood (labeled "likelihood" with a red arrow)
- $p(\theta)$: Prior (labeled "prior" with an orange arrow)
- $\int d\theta$: Integration over parameters (labeled "observations" with a red arrow)

1. Regression data set



Polynomial basis functions

Since the relationship between \mathbf{y} and \mathbf{x} is non-linear,
we'll apply polynomial basis expansion to degree d .

$$\Phi = \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^d \\ 1 & x_2 & x_2^2 & \dots & x_2^d \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^d \end{bmatrix}$$

2. Bayesian regression with known variance

what we are interested with.

- Prior

$w|\gamma \sim \text{Normal}(\vec{0}, \gamma^2 I_m)$

mean $\Rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$

- Likelihood

$p(y|X, w, \sigma) = \prod_{i=1}^n p(y_i|\vec{X}_i, \vec{w}, \sigma)$

known

Since $y_i|\vec{X}_i, \vec{w}, \sigma \sim \text{Normal}(\vec{X}_i^T \vec{w}, \sigma^2)$,

$$\vec{y}|X, \vec{w}, \sigma \sim \text{Multivariate Normal}(X\vec{w}, \sigma^2 I_n)$$

Posterior distribution: ~ Multivariate Normal

$$\mathcal{N}(W_N, V_N)$$

$$P(W) \quad \swarrow \quad m = 1 + \text{features}$$

$$= (2\pi)^{-\frac{m}{2}} \det(V_N)^{-\frac{1}{2}}$$

$$\exp\left\{-\frac{1}{2} (W - W_N)^T V_N^{-1} (W - W_N)\right\}$$

$$\propto \exp\left\{-\frac{1}{2} (W - W_N)^T V_N^{-1} (W - W_N)\right\}$$

$$= \exp\left\{-\frac{1}{2} (W^T - W_N^T) V_N^{-1} (W - W_N)\right\}$$

$$= \exp\left\{-\frac{1}{2} (\underbrace{W^T V_N^{-1} - W_N^T V_N^{-1}}_{\text{merge}}) (W - W_N)\right\}$$

$$= \exp\left\{-\frac{1}{2} (W^T V_N^{-1} W - \underbrace{W^T V_N^{-1} W_N - W_N^T V_N^{-1} W}_{\text{merges}} + W_N^T V_N^{-1} W_N)\right\}$$

$$\Rightarrow W^T V_N^{-1} W_N = W_N^T V_N^{-1} W$$

Prior distribution. \sim multivariate normal

$$\Rightarrow \exp\left\{-\frac{1}{2} (W - \underline{0})^T (\gamma^2 I_m)^{-1} (W - \underline{0})\right\}$$
$$= \exp\left\{-\frac{1}{2} (\gamma^2)^{-1} W^T W\right\}$$

Likelihood \Rightarrow product of Normal distribution

$$\Rightarrow \prod_{i=1}^n \exp\left\{-\frac{1}{2} \sigma^{-2} (y_i - X_i^T W)^2\right\}$$

$$\Rightarrow \exp\left\{-\frac{1}{2} (\sigma^2)^{-1} \sum_{i=1}^n (y_i - X_i^T W)^2\right\}$$

$$\Rightarrow \exp\left\{-\frac{1}{2} (\sigma^2)^{-1} \sum_{i=1}^n [y_i^2 - 2y_i X_i^T W + (X_i^T W)^2]\right\}$$

$$\sum_{i=1}^n y_i^2 = y^T y, \quad y^T = [y_1, \dots, y_n]$$

$$\Rightarrow \exp\left\{-\frac{1}{2} (\sigma^2)^{-1} [y^T y - 2W^T X^T y + W^T X^T X W]\right\}$$

$\begin{matrix} \nearrow & \nearrow & & \nearrow & \nearrow & \nearrow & \nearrow & \nearrow & \nearrow \\ | \times n & n \times 1 & & | \times m & m \times n & n \times 1 & & | \times m & m \times n & n \times m & m \times 1 \\ \Rightarrow | \times 1 & & & | \times 1 & & & & | \times 1 & & & | \times 1 \end{matrix}$

Posterior \propto Prior \cdot likelihood

$$= \exp\left\{-\frac{1}{2} [(\gamma^2)^{-1} W^T W + (\sigma^2)^{-1} y^T y - (\sigma^2)^{-1} 2W^T X^T y + (\sigma^2)^{-1} W^T X^T X W]\right\}$$

the kernel of posterior.

$$\exp \left\{ -\frac{1}{2} \left[(r^2)^{-1} \underbrace{W^T W}_{\textcircled{1}} + (\sigma^2)^{-1} y^T y - (\sigma^2)^{-1} \underbrace{W^T X^T X W}_{\textcircled{1}} \right] \right\}$$

$$= \exp \left\{ -\frac{1}{2} (W^T V_N^{-1} W - \underbrace{W^T V_N^{-1} W_N}_{\textcircled{1}} - W_N^T V_N^{-1} W + W_N^T V_N^{-1} W_N) \right\}$$

$$\hookrightarrow W^T V_N^{-1} W = (r^2)^{-1} W^T W + (\sigma^2)^{-1} W^T X^T X W$$

$$\Leftrightarrow V_N = \sigma^2 \left(X^T X + \frac{r^2}{\sigma^2} \right)^{-1}$$

$$\sigma^2 W^T V_N^{-1} W = \left(\frac{r^2}{\sigma^2} \right)^{-1} W^T W + W^T (X^T X) W$$

$$\Rightarrow W_N$$

$$\Rightarrow \begin{cases} \sigma^2 \sim \text{Inverse Gamma} \\ r^2 \sim \text{Inverse Gamma} \end{cases}$$

Bayesian regression with known variance

Given this formulation, the next step is to solve for the posterior over \mathbf{w}

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \sigma, \gamma) = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \sigma)p(\mathbf{w}|\gamma)}{p(\mathbf{y}|\mathbf{X}, \sigma)}$$

where $\mathbf{X} \in \mathbb{R}^{n \times m}$ is the feature matrix and $\mathbf{y} \in \mathbb{R}^n$ is the vector of target values for each instance.

In lectures, we derived the following solution:

$$\mathbf{w}|\mathbf{X}, \mathbf{y}, \sigma, \gamma \sim \text{Normal}(\mathbf{w}_N, \mathbf{V}_N)$$

where $\mathbf{V}_N = \sigma^2 \left(\mathbf{X}^T \mathbf{X} + \frac{\sigma^2}{\gamma^2} \mathbf{I}_m \right)^{-1}$ and $\mathbf{w}_N = \frac{1}{\sigma^2} \mathbf{V}_N \mathbf{X}^T \mathbf{y}$.

VCOV of posterior

mean vector of posterior

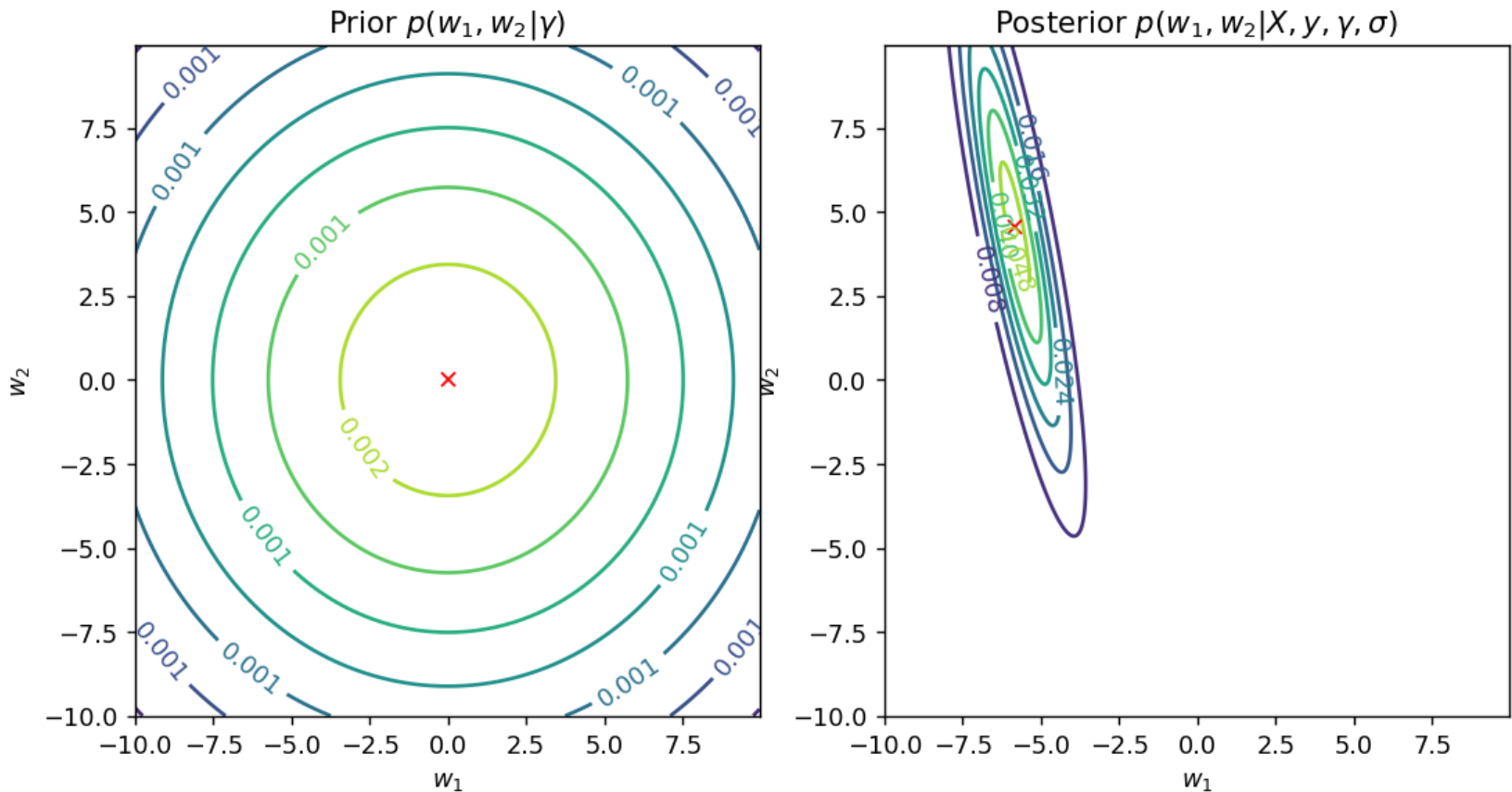
`numpy.linalg.inv()` Compute the (multiplicative) inverse of a matrix.

<https://docs.scipy.org/doc/numpy/reference/generated/numpy.linalg.inv.html#numpy.linalg.inv>

`np.Identity / np.eye` Return a 2-D array with ones on the diagonal and zeros elsewhere.

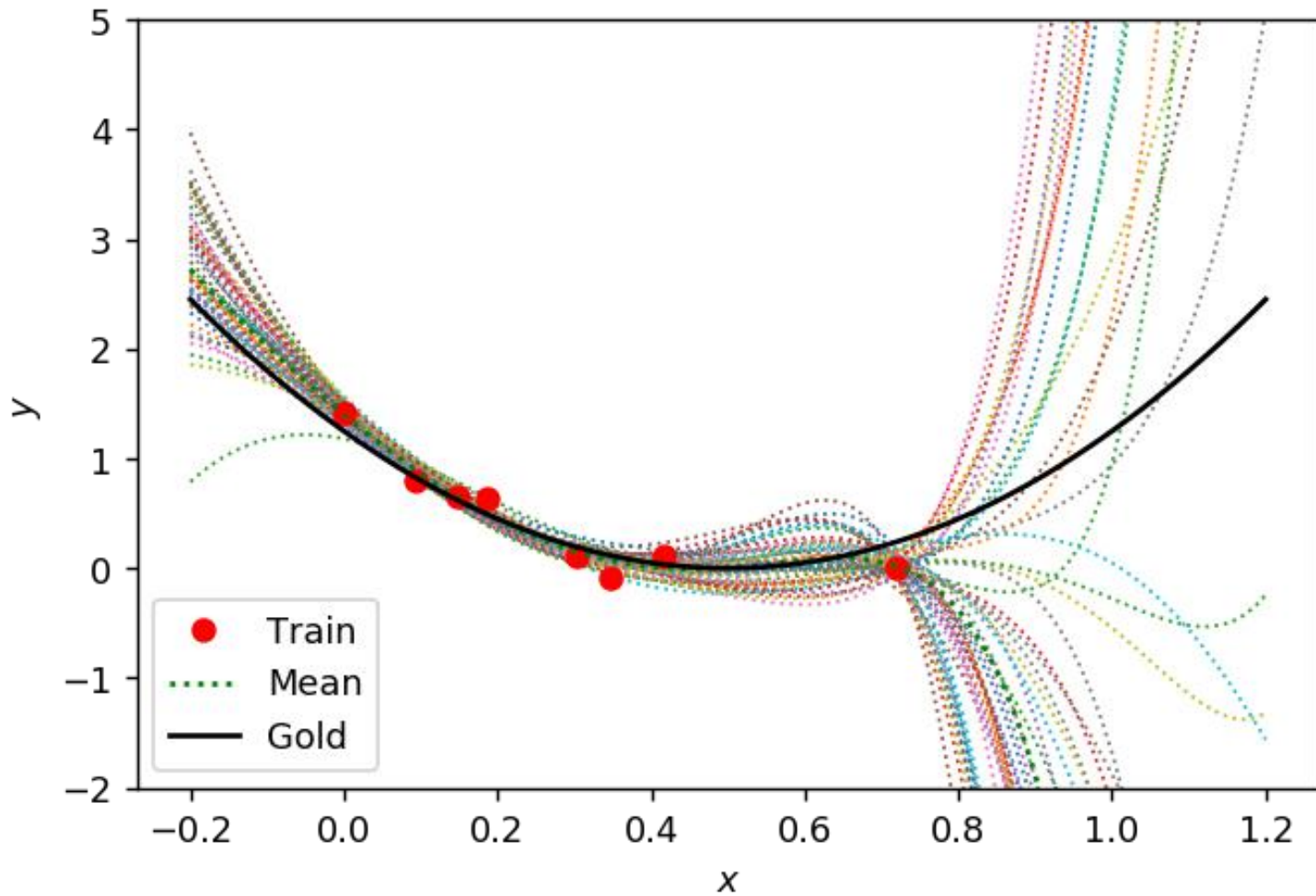
<https://github.com/numpy/numpy/blob/v1.9.1/numpy/core/numeric.py#L2125>

plot the prior and posterior over w_1, w_2



Discussion question: Can you explain why the prior and the posterior are so different? How is this related to the dataset? Why are the ellipses in the posterior not aligned to the axes? *You might want to change the parameter indices from 0,1 to other pairs to get a better idea of the full posterior.*

Bayesian inference



The Bayesian Predictive Distribution

Thanks to conjugacy, the predictive distribution can be found in closed form in our toy problem.

$$y_* | \mathbf{x}_*, \mathbf{w}_N, \mathbf{V}_N, \sigma = \text{Normal}[\langle \mathbf{x}^*, \mathbf{w}_N \rangle, \sigma_N^2(\mathbf{x}^*)]$$

$$\sigma_N^2(\mathbf{x}^*) = \sigma^2 + (\mathbf{x}^*)^T \mathbf{V}_N \mathbf{x}^*$$

```
def target_std(X, V_N, sigma):
    """
    Compute the predictive standard deviation for the target variable, given X, V_N and sigma

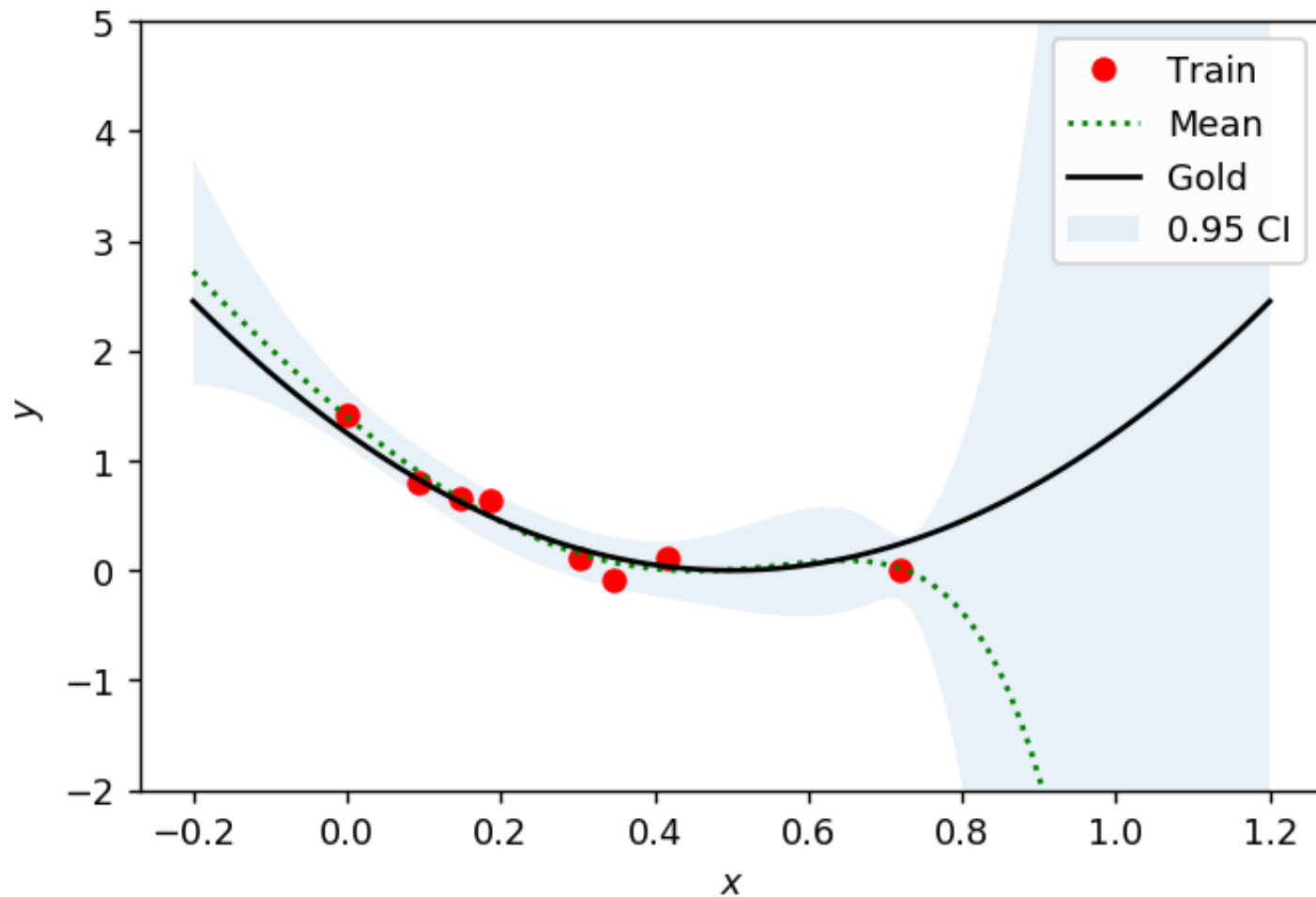
    Arguments
    =====
    X : numpy array, shape: (n_instances, n_features)
        feature matrix
    V_N : numpy array, shape: (n_features, n_features)
        covariance parameter

    Returns
    =====
    std : numpy array, shape: (n_instances,)
        predictive standard deviation for each instance in X
    """
    # your code here #

    variance = ...
    std = np.sqrt(variance)

    return std
```

Bayesian inference



Bayesian model selection

