

AA203

Optimal and Learning-based Control

Fundamentals of Regression



Stanford
University



Today's lecture

- Aim
 - Review regression models that will be used throughout this class
- Outline:
 - Maximum likelihood and least squares via ML
 - Bayesian inference and Bayesian linear regression
 - Gaussian process regression
 - Intro to neural networks

Regression

- Main tool in learning dynamics models, value functions, and policies
- Aims to model a continuous output as a function of the input
- Contrast to classification setting which has a finite set of outputs

Supervised learning

- We are given training dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^d$
- Want to find some f such that $\hat{y}_i = f(\mathbf{x}_i)$ is close to y_i

What does it mean to be “close”?

How do we represent f ?

How do we represent uncertainty?

Maximum likelihood

- Dominant method in statistical inference due to flexibility and computational tractability
- Assume model is parameterized by θ
- Likelihood $L(\theta) = p(\mathcal{D} | \theta)$
- Principle of maximum likelihood: choose $\theta^* = \operatorname{argmax}_{\theta} L(\theta)$
- In practice, will maximize *log likelihood*

$$\log p(\mathcal{D} | \theta) = \sum_{i=1}^d \log p(y_i | \mathbf{x}_i, \theta)$$

From max likelihood to least squares

- We have not chosen a model parameterization: how to represent f ?
- We will assume data is generated by $y = \mathbf{x}^T \boldsymbol{\theta}^* + \epsilon$
- Likelihood is

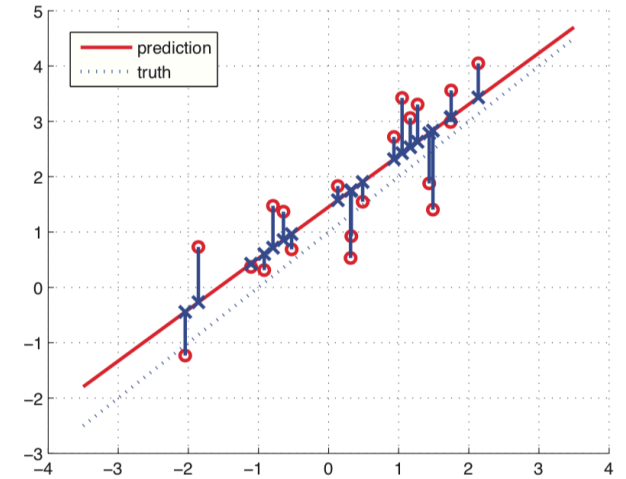
$$p(y | \mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}^T \boldsymbol{\theta}, \sigma^2) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(y - \mathbf{x}^T \boldsymbol{\theta})^2}{2\sigma^2}}$$

Least squares

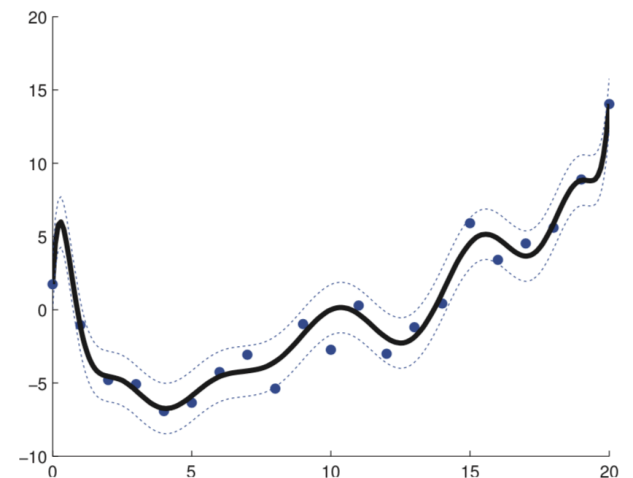
- Loss function is convex + smooth in θ , so first order necessary conditions imply global optimum

$$\nabla_{\theta} \|y - X\theta\|_2^2 = 2X^T X\theta - 2X^T y$$

- Observe \mathbf{x} without noise, so can transform nonlinearly to construct *features* $\phi(\mathbf{x})$



Murphy, 2012.



Bayesian inference

- Previously, we have computed a point estimate of model parameters, have not incorporated prior information
- Bayesian inference: specify prior belief over parameter, use Bayes' rule to compute the posterior distribution

$$p(\boldsymbol{\theta} \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathcal{D})}$$

- As opposed to point estimate, gives full distribution over parameter
- Allows incorporation of prior information

Bayesian inference: difficulties

$$p(\boldsymbol{\theta} | \mathcal{D}) = \frac{p(\mathcal{D} | \boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathcal{D})}$$

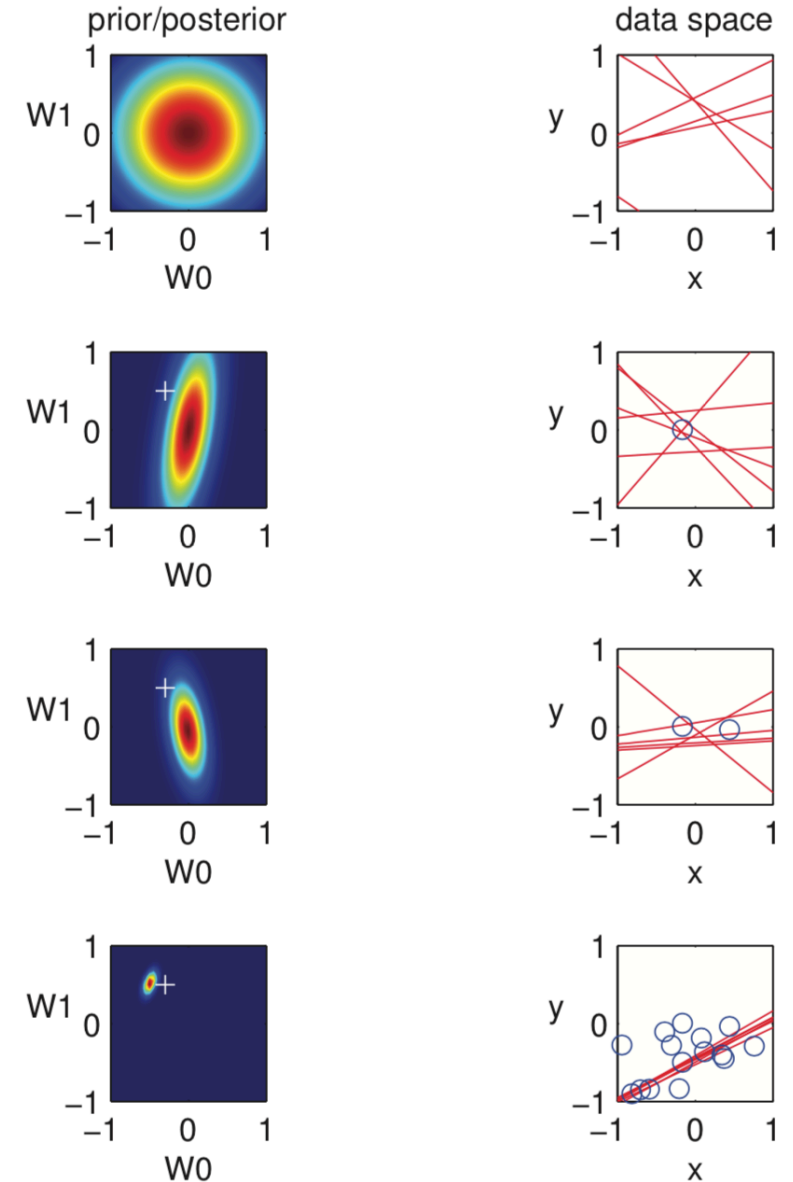
- Requires computing marginal likelihood $p(\mathcal{D})$
- Typically results in intractable integral

$$\int p(\mathcal{D} | \boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta}$$

- In general, we will not be able to exactly compute posterior
 - Instead, turn to sampling-based approximations
- We will examine special cases where posterior computation is exact

Bayesian least squares

- Assume linear model with iid Gaussian noise
- Choose prior $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}_0, \Sigma_0)$



Murphy, 2012.

Maximum a Posteriori Estimation

- Want to factor in prior information but computing full posterior is irrelevant/too difficult
- Maximum a Posteriori (MAP) estimator

$$\hat{\boldsymbol{\theta}}_{MAP} = \operatorname{argmax}_{\boldsymbol{\theta}} p(\boldsymbol{\theta} \mid \mathcal{D}) = \operatorname{argmax}_{\boldsymbol{\theta}} p(\mathcal{D} \mid \boldsymbol{\theta})p(\boldsymbol{\theta})$$

- As prior becomes flatter, approaches maximum likelihood estimator

Gaussian process regression

- Gaussian process: can think of as infinite dimensional Gaussian distribution, or Gaussian distribution over functions
- Fully specified by mean and variance

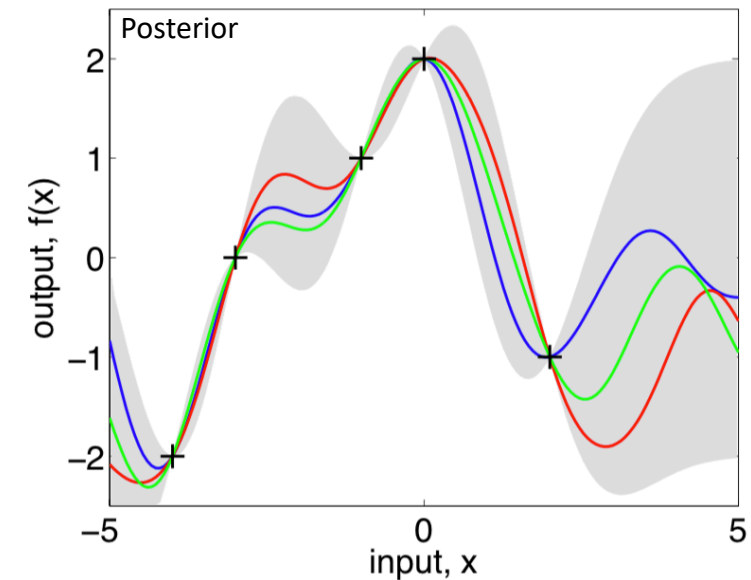
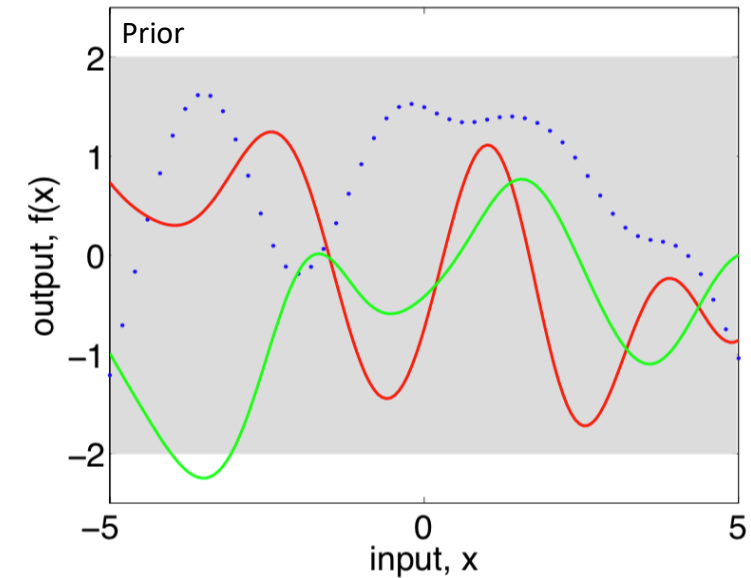
$$\mu(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})]$$

$$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - \mu(\mathbf{x}))(f(\mathbf{x}') - \mu(\mathbf{x}'))]$$

- The function $k(\mathbf{x}, \mathbf{x}')$ is referred to as the *kernel*; must be

- Symmetric: $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}', \mathbf{x})$

- Positive definite:
$$\begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \dots & k(\mathbf{x}_1, \mathbf{x}_k) \\ \vdots & & \vdots \\ k(\mathbf{x}_k, \mathbf{x}_1) & \dots & k(\mathbf{x}_k, \mathbf{x}_k) \end{bmatrix} \succ 0$$



Rasmussen and Williams, 2006.

Kernel functions

- Consider Bayesian linear regression with prior $\boldsymbol{\theta} \sim \mathcal{N}(0, \Sigma_0)$

$$\mu(\mathbf{x}) = \mathbb{E}[\boldsymbol{\theta}^T \phi(\mathbf{x})] = 0$$

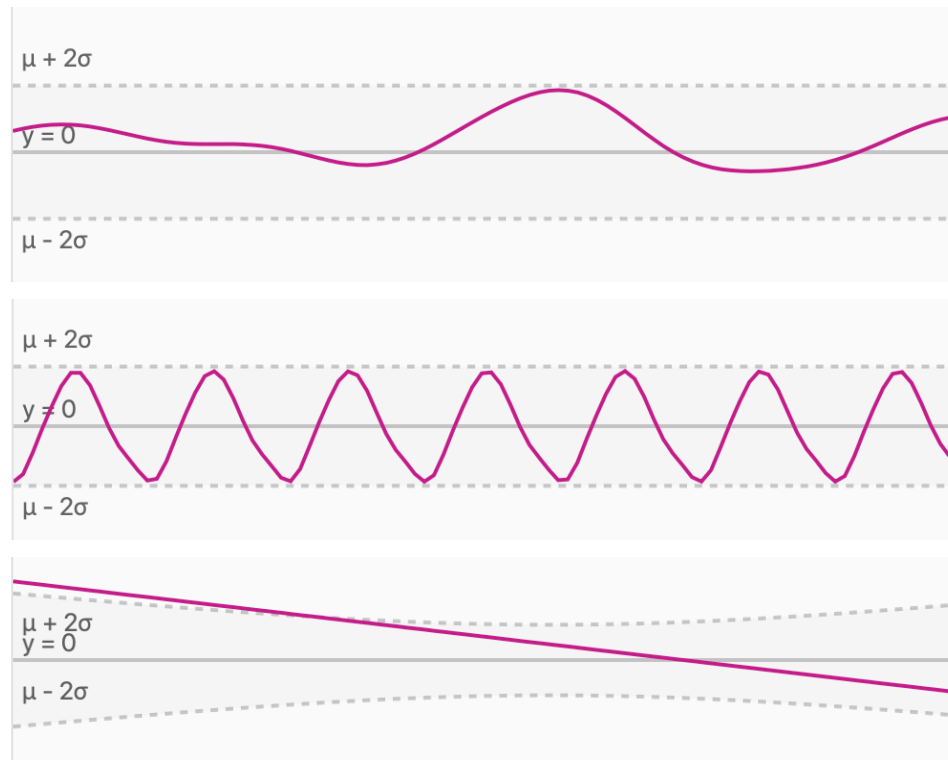
$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \mathbb{E}[\boldsymbol{\theta}\boldsymbol{\theta}^T] \phi(\mathbf{x}') = \phi(\mathbf{x})^T \Sigma_0 \phi(\mathbf{x}')$$

- So kernel is a simple inner product of feature vectors; can rewrite BLR in this form
- In general: every positive definite kernel corresponds to a (possibly infinite dimensional) set of features; e.g

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2\ell} \|\mathbf{x} - \mathbf{x}'\|_2^2\right)$$

Kernel function

- Choice of kernel determines features/shape of the posterior belief over function

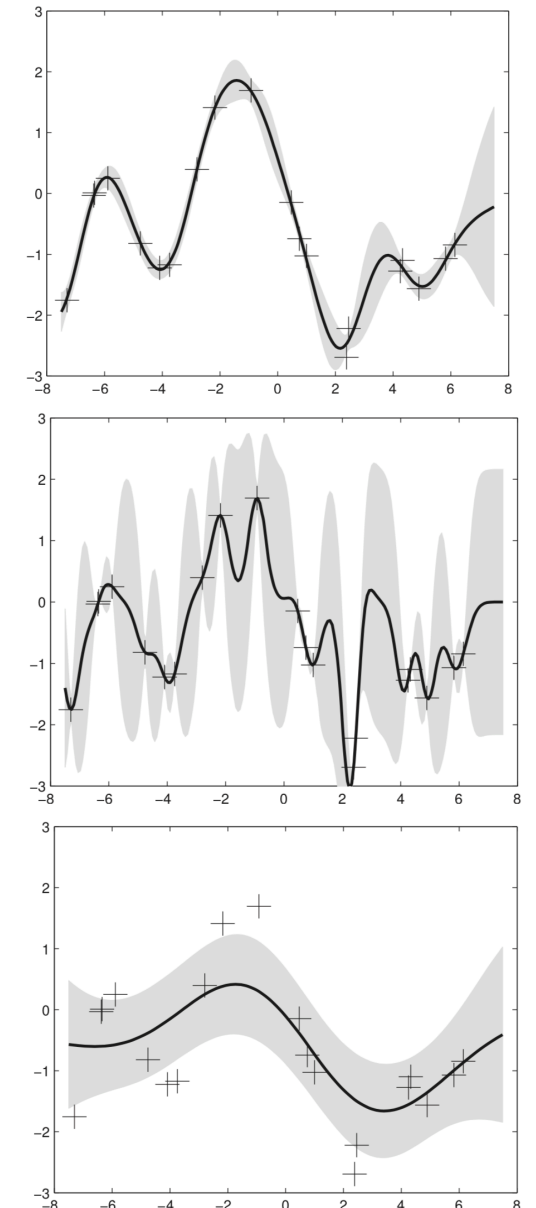


Gortler et al.,
Distill, 2019

Squared
exponential
kernel

Periodic
kernel

Linear
kernel



Murphy, 2012.

Gaussian process: machinery

• Let $\begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}\right)$ then $p(\mathbf{y}_1 | \mathbf{y}_2) = \mathcal{N}(\boldsymbol{\mu}_1 + \Sigma_{21}\Sigma_{11}^{-1}(\mathbf{y}_2 - \boldsymbol{\mu}_2), \Sigma_{11} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12})$

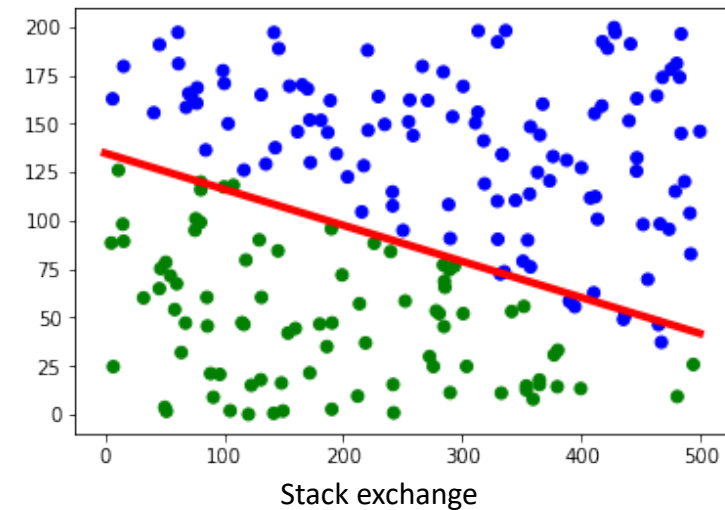
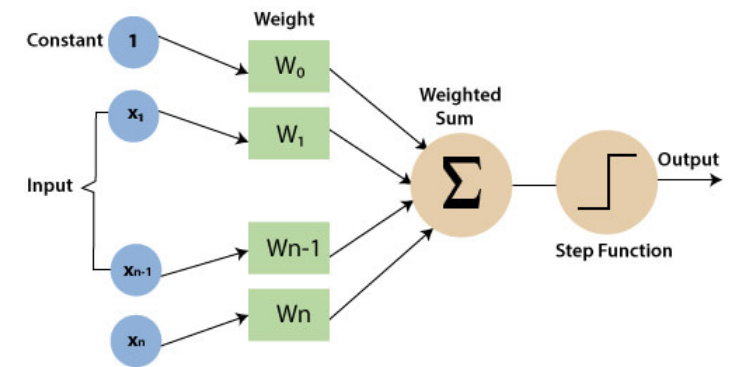
• So have $\begin{bmatrix} f(\mathbf{x}^*) \\ \mathbf{y} \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} k(\mathbf{x}^*, \mathbf{x}^*) & K(X, \mathbf{x}^*) \\ K(\mathbf{x}^*, X) & K(X, X) + \sigma^2 I \end{bmatrix}\right)$

• With $K(\mathbf{x}^*, X) = \begin{bmatrix} k(\mathbf{x}^*, \mathbf{x}_1) \\ \vdots \\ k(\mathbf{x}^*, \mathbf{x}_d) \end{bmatrix}$ $K(X, X) = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \dots & k(\mathbf{x}_1, \mathbf{x}_d) \\ \vdots & & \vdots \\ k(\mathbf{x}_d, \mathbf{x}_1) & \dots & k(\mathbf{x}_d, \mathbf{x}_d) \end{bmatrix}$

• Gives $\mathbb{E}[f(\mathbf{x}^*)] = K(\mathbf{x}^*, X)(K(X, X) + \sigma^2 I)^{-1}\mathbf{y}$
 $\text{var}(f(\mathbf{x}^*)) = k(\mathbf{x}^*, \mathbf{x}^*) - K(\mathbf{x}^*, X)(K(X, X) + \sigma^2 I)^{-1}K(X, \mathbf{x}^*)$

Neural networks: perceptron

- Have so far pre-specified features
- Since 2012, significant progress in learning features via neural networks
- Perceptron: $h(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$
- The nonlinearity $\sigma(\cdot)$ is the *threshold* function
- Example: sigmoid
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$
- Results in linear binary classifier



Neural networks

- Feed-forward neural network: stack hidden units

- Single hidden layer:

$$h = \sigma(W_1x)$$

$$y = W_2h$$

- Can choose loss function
 - Squared error common for regression

- Problem: optimization no longer convex
 - Bottleneck until empirical performance became state of the art

