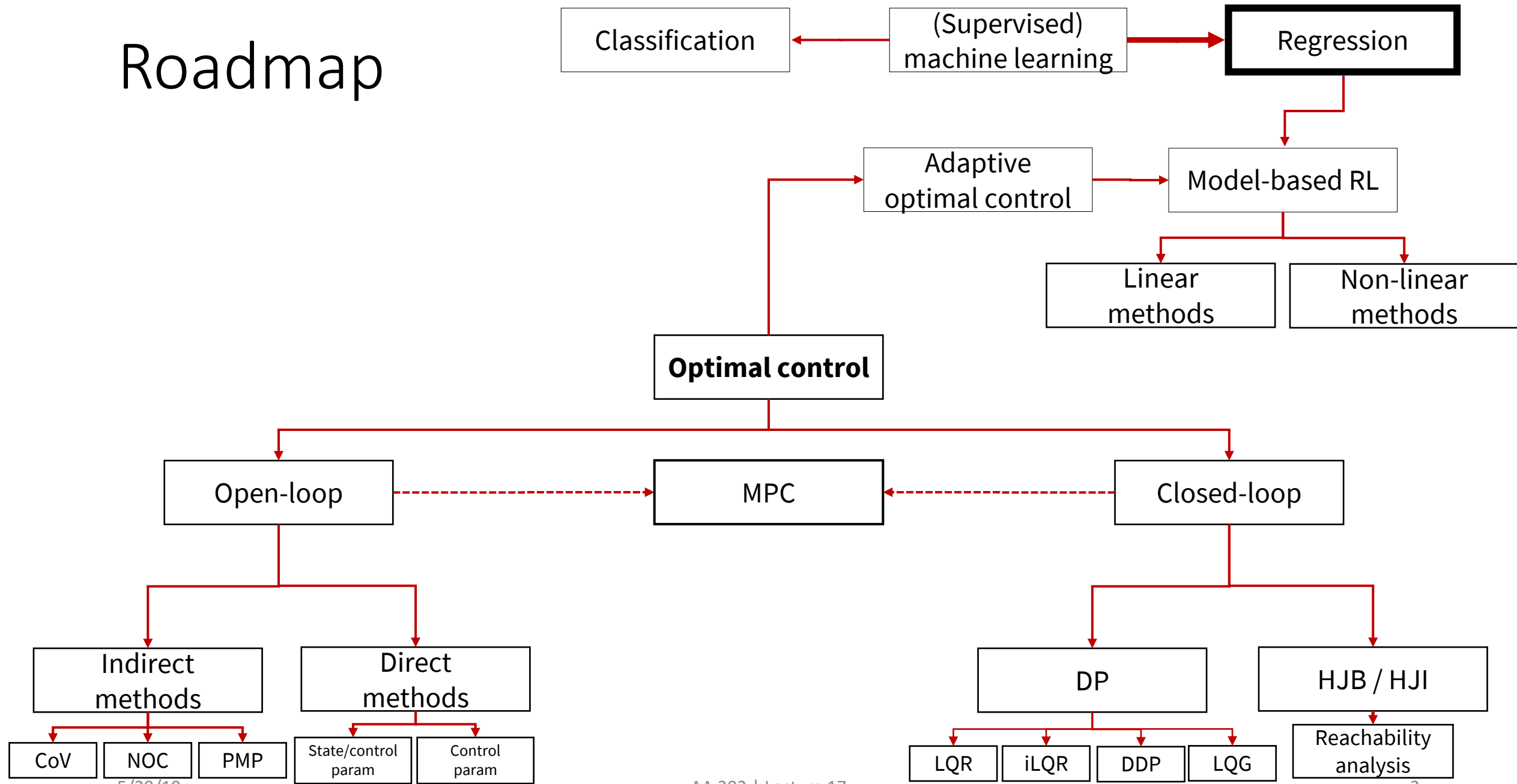


AA203

Optimal and Learning-based Control

Nonlinear Regression

Roadmap



Linear Regression

- We will begin by revisiting linear regression (in the scalar output case, for simplicity).
- We will assume a model of the form $y_i = \boldsymbol{\theta}^T \mathbf{x}_i + \epsilon_i$, where ϵ_i is zero-mean Gaussian with known covariance σ_ϵ^2 .
- Then, assuming data of the form $y_1, \dots, y_n, \mathbf{x}_1, \dots, \mathbf{x}_n$, we can write the least squares problem as

$$\min_{\boldsymbol{\theta}} ||\mathbf{y} - X\boldsymbol{\theta}||^2$$

where $\mathbf{y} = [y_1, \dots, y_n]^T$, $X = [\mathbf{x}_1, \dots, \mathbf{x}_n]^T$.

Least Squares via MLE

- This minimization problem is intuitive – we want to reduce the error of our model predictions, so we minimize the error. But formally, where does it come from?
- One answer: Maximum Likelihood Estimation (MLE)
 - Define the **likelihood function**,

$$\mathcal{L}(\boldsymbol{\theta} \mid X, \mathbf{y}) = p(\mathbf{y} \mid X, \boldsymbol{\theta}) = \prod_{i=1}^n p(y_i \mid \mathbf{x}_i, \boldsymbol{\theta})$$

where because we assumed zero mean gaussian errors,

$$p(y_i \mid \mathbf{x}_i, \boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}^T \mathbf{x}_i, \sigma_\epsilon^2).$$

- This is the probability of our output variable as a function of the parameter

Least Squares via MLE

- MLE maximizes the likelihood function (i.e. the likelihood of our data).
 - For convenience, we will minimize **negative log likelihood** – the optimization problem yields the same solution

- So, $-\log \mathcal{L}(\boldsymbol{\theta} \mid X, \mathbf{y}) = -\sum_{i=1}^n \log \frac{1}{\sqrt{2\pi}\sigma_{\epsilon}} e^{-\frac{(y_i - \boldsymbol{\theta}^T \mathbf{x}_i)^2}{2\sigma_{\epsilon}^2}}$
- Can drop $\log \frac{1}{\sqrt{2\pi}\sigma_{\epsilon}}$ from optimization problem, which yields

$$-\log \mathcal{L}(\boldsymbol{\theta} \mid X, \mathbf{y}) = \sum_{i=1}^n \frac{(y_i - \boldsymbol{\theta}^T \mathbf{x}_i)^2}{2\sigma_{\epsilon}^2}$$

Which, ignoring the covariance scaling factor, yields the least squares problem.

Least Squares: Bayesian perspective

- MLE gives a point estimate of $\boldsymbol{\theta}$. What if we instead take a Bayesian approach, in which we specify a **prior** over $\boldsymbol{\theta}$, and compute the posterior distribution after observing data?
- This is based on **Bayes' rule**:

$$p(\boldsymbol{\theta} \mid X, \mathbf{y}) = \frac{p(\mathbf{y} \mid X, \boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathbf{y} \mid X)}$$

- Ignoring $p(\mathbf{y} \mid X)$ for now, we will fix a Gaussian prior, $\boldsymbol{\theta} \sim \mathcal{N}(0, \Sigma_0)$, thus

$$\begin{aligned} p(\mathbf{y} \mid X, \boldsymbol{\theta})p(\boldsymbol{\theta}) &= e^{-\frac{1}{2\sigma_\epsilon^2}(\mathbf{y}-X^T\boldsymbol{\theta})^T(\mathbf{y}-X^T\boldsymbol{\theta})} e^{-\frac{1}{2}\boldsymbol{\theta}^T\Sigma_0^{-1}\boldsymbol{\theta}} \\ &= e^{-\frac{1}{2}(\boldsymbol{\theta}-\bar{\boldsymbol{\theta}})^T(\sigma_\epsilon^{-2}X X^T + \Sigma_0^{-1})(\boldsymbol{\theta}-\bar{\boldsymbol{\theta}})} \end{aligned}$$

where $\bar{\boldsymbol{\theta}} = \sigma_\epsilon^{-2}(\sigma_\epsilon^{-2}X X^T + \Sigma_0^{-1})^{-1}X\mathbf{y}$

Bayesian Least Squares

- So, the posterior over $\boldsymbol{\theta}$ is a Gaussian of the form

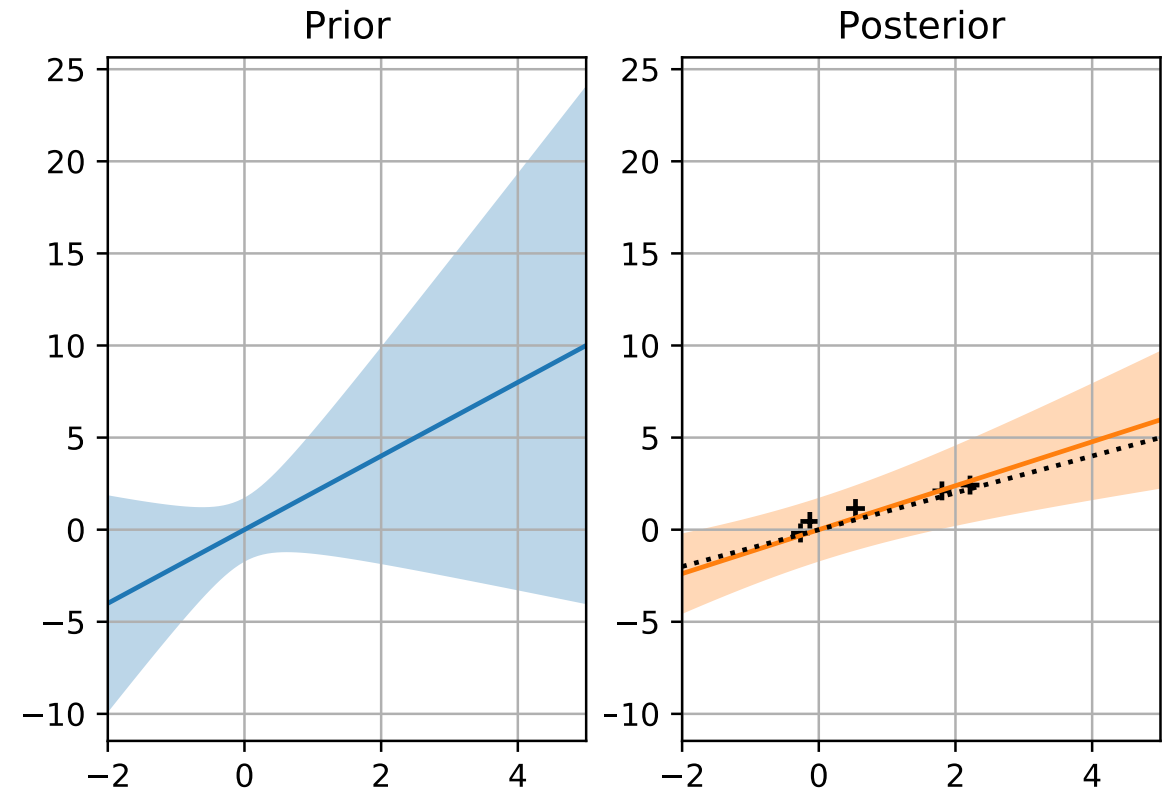
$$p(\boldsymbol{\theta} \mid X, \mathbf{y}) = \mathcal{N}(\sigma_{\epsilon}^{-1} \Sigma_n X \mathbf{y}, \Sigma_n)$$

$$\text{where } \Sigma_n = (\sigma_{\epsilon}^{-1} X X^T + \Sigma_0^{-1})^{-1}$$

- What just happened? If the prior is Gaussian and the data likelihood is Gaussian, the posterior over the weights is Gaussian
- What does this mean about the predictive distribution, $p(\mathbf{y}_* \mid \mathbf{x}_*, X, \mathbf{y})$?
- Note that for constant \mathbf{a} , $E[\mathbf{a}^T \mathbf{x}] = \mathbf{a}^T E[\mathbf{x}]$ and $\text{Var}(\mathbf{a}^T \mathbf{x}) = \mathbf{a}^T \text{Var}(\mathbf{x}) \mathbf{a}$
- Thus, $p(\mathbf{y}_* \mid \mathbf{x}_*, X, \mathbf{y}) = \mathcal{N}(\sigma_{\epsilon}^{-1} \mathbf{x}_*^T \Sigma_n X \mathbf{y}, \mathbf{x}_*^T \Sigma_n \mathbf{x}_* + \sigma_{\epsilon}^2)$

Bayesian Least Squares

- What features does this model have?
 - Our posterior over weights is Gaussian
 - Our posterior predictive distribution is Gaussian
 - We have a closed form representation of the posterior density $p(\mathbf{y}_* | \mathbf{x}_*, X, \mathbf{y})$; **useful for capturing model confidence!**
 - Note that in this derivation, we could replace \mathbf{x} with $\phi(\mathbf{x})$, derivation proceeds the same
 - Note: Tikhonov regularization (e.g. ridge regression) is a consequence of this model!



BLR with nonlinear basis functions

- Let $\boldsymbol{\phi}_* = \boldsymbol{\phi}(\mathbf{x}_*)$, and $f_* = \boldsymbol{\phi}_*^T \boldsymbol{\theta}$. Then, $f_* \mid \mathbf{x}_*, X, \mathbf{y}$ is Gaussian with

$$\begin{aligned} E[f_* \mid \mathbf{x}_*, X, \mathbf{y}] &= \boldsymbol{\phi}_*^T \Sigma_0 \Phi (K + \sigma_\epsilon^2 I)^{-1} \mathbf{y} \\ \text{Var}(f_* \mid \mathbf{x}_*, X, \mathbf{y}) &= \boldsymbol{\phi}_*^T \Sigma_0 \boldsymbol{\phi}_* - \boldsymbol{\phi}_*^T \Sigma_0 \Phi (K + \sigma_\epsilon^2 I)^{-1} \Phi^T \Sigma_0 \boldsymbol{\phi}_* \\ \text{where } K &= \Phi^T \Sigma_0 \Phi \end{aligned}$$

- This is a mess, but one important fact: the basis function $\boldsymbol{\phi}$ always enters through $\boldsymbol{\phi}_*^T \Sigma_0 \boldsymbol{\phi}_*$, $\boldsymbol{\phi}_*^T \Sigma_0 \Phi$, or $\Phi^T \Sigma_0 \Phi$, so entries are always of the form $\boldsymbol{\phi}(\mathbf{x})^T \Sigma_0 \boldsymbol{\phi}(\mathbf{x}')$
- We will define **kernel function** $k(\mathbf{x}, \mathbf{x}') = \boldsymbol{\phi}(\mathbf{x})^T \Sigma_0 \boldsymbol{\phi}(\mathbf{x}')$
- Kernel trick: we can skip computing $\boldsymbol{\phi}(\mathbf{x})$, only ever have to compute $k(\mathbf{x}, \mathbf{x}')$!

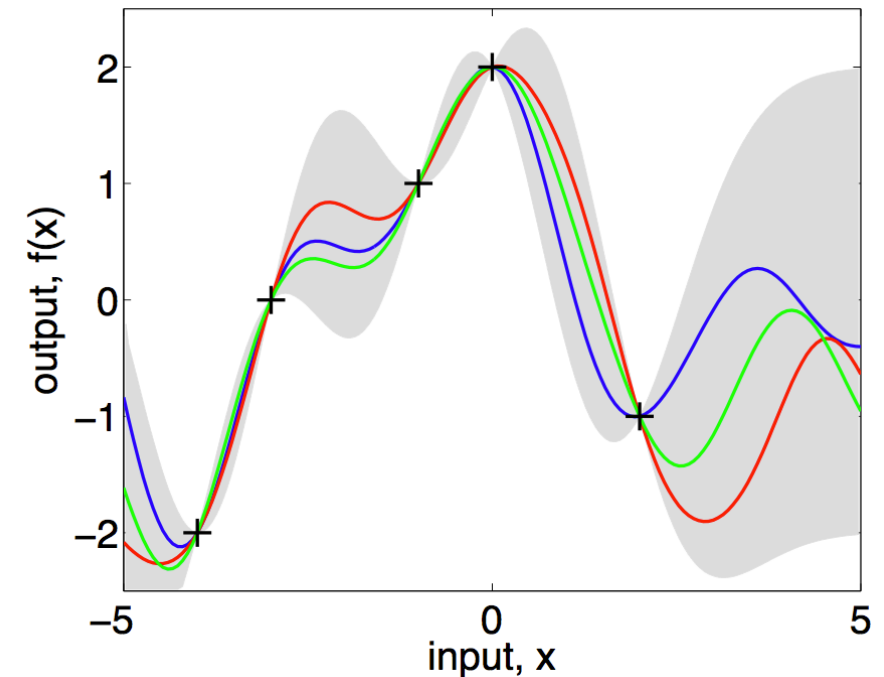
Gaussian processes

- What does the kernel trick buy us? We can use kernels that correspond to very high dimensional (even infinite-dimensional) basis functions, since we never have to compute them!

- Example: squared exponential kernel

$$k(\mathbf{x}, \mathbf{x}') = e^{-\frac{1}{2l^2} \|\mathbf{x} - \mathbf{x}'\|^2}$$

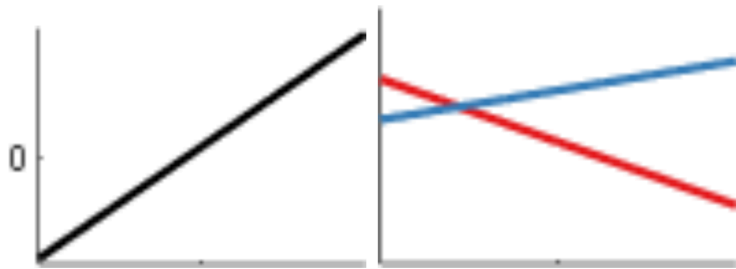
- Often a default for GP regression
- Captures smooth local correlations
- Corresponds to infinite-dimensional basis function



Choice of kernel function

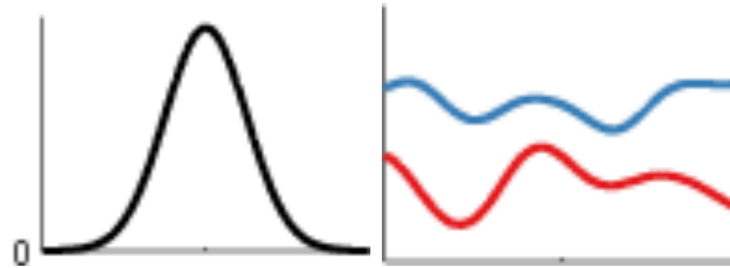
- The choice of kernel function allows us to encode prior knowledge into regression problem
- Can we choose any function as a kernel function?
 - Short answer: no. Kernel must correspond to inner product between some basis functions
- Let k_1, k_2 be valid kernels. Then, the following are also valid kernels:
 - $a k_1(\mathbf{x}, \mathbf{x}') + b k_2(\mathbf{x}, \mathbf{x}')$ for $a, b \geq 0$
 - $k_1(\mathbf{x}, \mathbf{x}') k_2(\mathbf{x}, \mathbf{x}')$
 - $k_1(f(\mathbf{x}), f(\mathbf{x}'))$

Kernel function examples



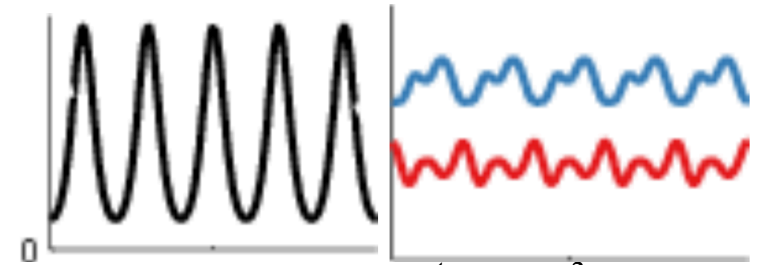
$$k_{\text{Lin}}(x, x') = \sigma_b^2 + \sigma_v^2(x - c)(x' - c)$$

Linear



$$k_{\text{SE}}(x, x') = \sigma^2 \exp\left(-\frac{(x-x')^2}{2\ell^2}\right)$$

Squared Exponential



$$k_{\text{Per}}(x, x') = \sigma^2 \exp\left(-\frac{2 \sin^2(\pi|x-x'|/p)}{\ell^2}\right)$$

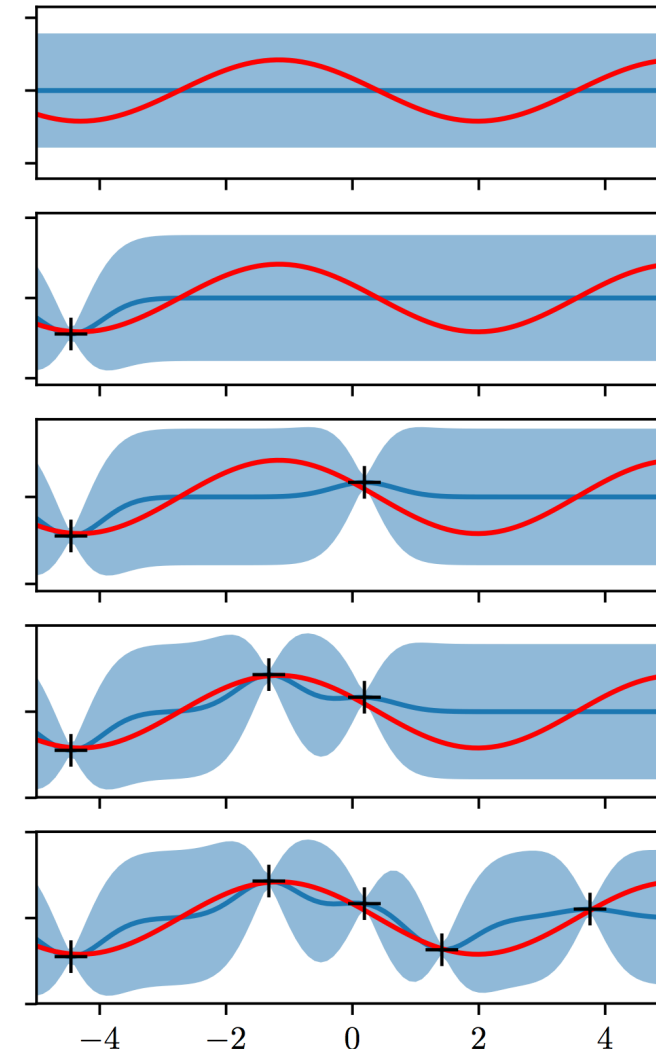
Periodic

Plus many more (e.g. polynomial, Matern, rational quadratic, etc.).

See: <https://www.cs.toronto.edu/~duvenaud/cookbook/>

Function space view of GPs

- A GP is a collection of random variables, any finite number of which follow a multivariate Gaussian distribution
- GPs generalize the notion of Gaussian distributions to the function space (i.e. the infinite-dimensional case); thus, they are distributions over functions (where this distribution is implied by choice of covariance function/kernel)
- Finite-dimensional distributional view: distribution of data points and test point(s) is jointly Gaussian with correlation governed by the chosen kernel



Limitations of Gaussian process regression

- The largest limitation of GP regression is computational complexity
 - Computation of posterior for a point requires computation of $(K + \sigma_\epsilon^2 I)^{-1}$, where $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$
 - Thus, K has dimension equal to the number of data points conditioned on
 - Inversion of dense matrix has complexity $O(n^3)$
 - Thus GP regression scales extremely poorly with the amount of data
 - One approach: use small collection of “inducing points”, ignore the rest of the data
- Other limitations:
 - encode prior structure directly via kernels (possibly limiting the capacity of features)
 - Relies on known noise covariance (relaxing this assumption leads to t processes)

Further reading on GP methods

- *Gaussian Processes for Machine Learning*, Rasmussen and Williams
 - Great, comprehensive coverage, free online
- <https://distill.pub/2019/visual-exploration-gaussian-processes/>
 - Interactive GP visualizations

Neural networks

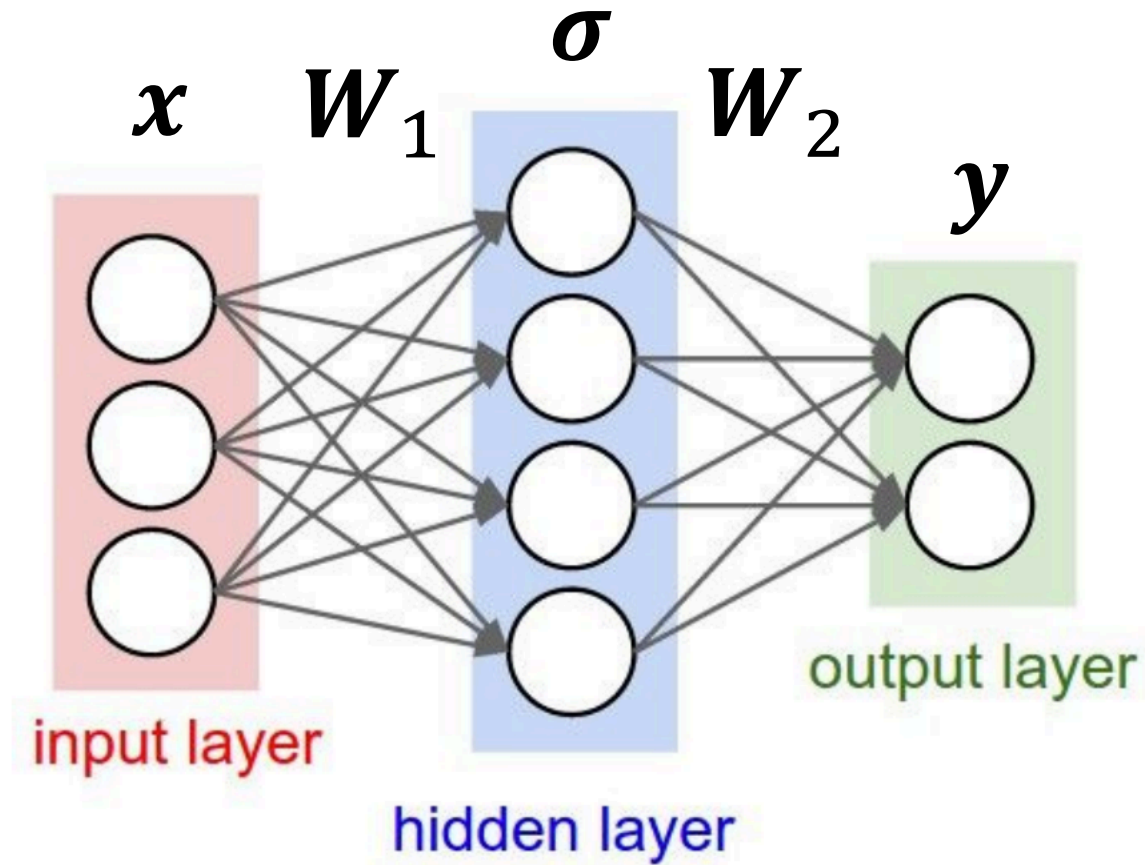
- Linear models are restrictive; while nonlinear basis functions are more expressive, GPs are limited in terms of choice of kernels
- Neural networks directly optimize nonlinear features
- Example: 2-layer neural network

$$\mathbf{y} = f(\mathbf{x}) = W_2 \boldsymbol{\sigma}(W_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2$$

where $\boldsymbol{\sigma}$ is some nonlinear function (applied element-wise), typically called the **activation function**. Examples:

- Relu: $\sigma(x_i) = \max(0, x_i)$
- Tanh: $\sigma(x_i) = \tanh(x_i)$
- Sigmoid/logistic: $\sigma(x_i) = \frac{1}{1+e^{-x_i}}$
- and many many more
- Can also have nonlinearities that are not applied element-wise, but they are recently less popular

Neural networks



Training neural networks

- Given this model, we wish to solve the optimization problem

$$\min_{W_1, W_2} \sum_{i=1}^n L(\mathbf{x}_i, \mathbf{y}_i)$$

where L is some loss function.

- In least squares, we had an analytical solution to minimize our empirical loss
- Because of nonlinear dependence on parameters, neural network optimization problem is non-convex, and so we must turn to gradient descent
- Gradient descent can be done efficiently via **error backpropagation**: basically using chain rule backwards on the network
 - Network activation functions (and losses) are usually chosen to be cheap/easy to differentiate

Computation graphs and deep learning

- Calculating gradients for networks can be done automatically/efficiently for networks by building computation graphs and performing automatic differentiation
 - This is what deep learning frameworks such as e.g. tensorflow do
- What is deep learning? Stacking more layers/more exciting “neural network architectures”/collection of tricks to improve training
 - Examples: convolutional networks (popular in vision), recurrent networks (popular in language modelling), etc.
- These architectures have, for the most part, not been used very much in model learning for RL and control

Limitations of neural networks

- Data hungry: may require a large amount of data to train effectively
- Hard to incorporate prior knowledge into (compared to GPs), e.g. smoothness
- Possibly unstable to train; definitely time/compute intensive
- No confidence measure/expression of uncertainty
 - there are some architectures that address this but fully Bayesian treatment of neural networks are still an active research area

Further reading on neural networks

- *Deep Learning*, Goodfellow, Bengio and Courville
 - The standard on deep learning methods
- CS231n lecture notes and slides
 - Introduction to neural networks + vision models
- CS229 notes on deep learning
- <http://playground.tensorflow.org>

Next time

- Nonlinear model-based reinforcement learning