

# CS475/CS675

## Lecture 22: July 14, 2016

Application: regression in machine  
learning

# Machine Learning

- Arthur Samuel (1959): Machine learning is the field of study that gives computers the ability to learn without being explicitly programmed.
- Tom Mitchell (1998): A computer program is said to *learn* from **experience E** with respect to some class of **tasks T** and performance **measure P**, if its performance at tasks in T, as measured by P, improves with experience E.

# Supervised Learning

- Inductive learning:
  - Given a **training set of examples** of the form  $(x, f(x))$ 
    - $x$  is the input,  $f(x)$  is the output
  - Return a function  $h$  that **approximates**  $f$ 
    - $h$  is called the hypothesis

# Linear model for regression

- Simplest form of regression
- Picture:

# Problem

- Data:  $\{(\mathbf{x}_1, t_1), (\mathbf{x}_2, t_2), \dots, (\mathbf{x}_N, t_N)\}$ 
  - $\mathbf{x} = \langle x_1, x_2, \dots, x_D \rangle$ : input vector
  - $t$ : target (continuous value)
- Problem: find hypothesis  $h$  that maps  $\mathbf{x}$  to  $t$ 
  - Assume that  $h$  is linear:

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x_1 + \dots + w_D x_D = \mathbf{w}^T \begin{pmatrix} 1 \\ \mathbf{x} \end{pmatrix}$$

- Objective: minimize some loss function
  - Euclidean loss:  $L_2(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (y(\mathbf{x}_n, \mathbf{w}) - t_n)^2$

# Optimization

- Find best  $w$  that minimizes Euclidean loss

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \frac{1}{2} \sum_{n=1}^N \left( t_n - \mathbf{w}^T \begin{pmatrix} 1 \\ \mathbf{x}_n \end{pmatrix} \right)^2$$

- Convex optimization problem  
⇒ unique optimum (global)

# Solution

- Let  $\bar{\mathbf{x}} = \begin{pmatrix} 1 \\ \mathbf{x} \end{pmatrix}$  then  $\min_{\mathbf{w}} \frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}^T \bar{\mathbf{x}}_n)^2$
- Find  $\mathbf{w}^*$  by setting the derivative to 0

$$\frac{\partial L_2}{\partial w_j} = \sum_{n=1}^N (t_n - \mathbf{w}^T \bar{\mathbf{x}}_n) \bar{x}_{nj} = 0 \quad \forall j$$

$$\Rightarrow \sum_{n=1}^N (t_n - \mathbf{w}^T \bar{\mathbf{x}}_n) \bar{\mathbf{x}}_n = 0$$

- This is a linear system in  $\mathbf{w}$ , therefore we rewrite it as  $\mathbf{A}\mathbf{w} = \mathbf{b}$

$$\text{where } \mathbf{A} = \sum_{n=1}^N \bar{\mathbf{x}}_n \bar{\mathbf{x}}_n^T \text{ and } \mathbf{b} = \sum_{n=1}^N t_n \bar{\mathbf{x}}_n$$

# Solution

- If training instances span  $\mathfrak{R}^{D+1}$  then  $\mathbf{A}$  is invertible:

$$\mathbf{w} = \mathbf{A}^{-1}\mathbf{b}$$

- In practice it is faster to solve the linear system  $\mathbf{Aw} = \mathbf{b}$  directly instead of inverting  $\mathbf{A}$ 
  - Gaussian elimination
  - Conjugate gradient
  - Iterative methods



# Generalized Linear Models

- How can we do non-linear regression with linear algebra?
- Idea: map inputs to a different space and do linear regression in that space

# Example

- Suppose the underlying function is quadratic

# Basis functions

- Use non-linear basis functions:

- Let  $\phi_i$  denote a basis function

$$\phi_0(x) = 1$$

$$\phi_1(x) = x$$

$$\phi_2(x) = x^2$$

- Let the hypothesis space  $H$  be

$$H = \{x \rightarrow w_0\phi_0(x) + w_1\phi_1(x) + w_2\phi_2(x) \mid w_i \in \mathbb{R}\}$$

- If the basis functions are non-linear in  $x$ , then a non-linear hypothesis can still be found by linear regression

# Common basis functions

- Polynomial:  $\phi_j(x) = x^j$
- Gaussian:  $\phi_j(x) = e^{-\frac{(x-\mu_j)^2}{2s^2}}$
- Sigmoid:  $\phi_j(x) = \sigma\left(\frac{x-\mu_j}{s}\right)$   
where  $\sigma(a) = \frac{1}{1+e^{-a}}$
- Also Fourier basis functions, wavelets, etc.

# Non-linear regression

- Same as linear regression, but replace  $\bar{\mathbf{x}}$  by  $\phi(\mathbf{x})$
- Problem:

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2$$

- Solution:  $\mathbf{A}\mathbf{w} = \mathbf{b}$

where  $\mathbf{A} = \sum_{n=1}^N \phi(\mathbf{x}_n)\phi(\mathbf{x}_n)^T$  and  $\mathbf{b} = \sum_{n=1}^N t_n \phi(\mathbf{x}_n)$

# Case Study: AIBO Gait Optimization



# Gait Optimization

- Problem: find best parameter setting of the gait controller to maximize walking speed
  - Why?: Fast robots have a better chance of winning in robotic soccer
- Solutions:
  - Stochastic hill climbing
  - **Gaussian Process Regression**
    - Lizotte, Wang, Bowling, Schuurmans (2007) Automatic Gait Optimization with Gaussian Processes, *International Joint Conferences on Artificial Intelligence (IJCAI)*.

# Search Problem

- Let  $\mathbf{x} \in \mathfrak{R}^{15}$ , be a vector of 15 parameters that defines a controller for gait
- Let  $f: \mathbf{x} \rightarrow \mathfrak{R}$  be a mapping from controller parameters to gait speed
- Problem: find parameters  $\mathbf{x}^*$  that yield highest speed.

$$\mathbf{x}^* \leftarrow \operatorname{argmax}_{\mathbf{x}} f(\mathbf{x})$$

But  $f$  is unknown...



# Approach

- Picture

# Approach

- Initialize  $f(\cdot) \sim GP(m(\cdot), k(\cdot, \cdot))$
- Repeat:

– Select new  $\mathbf{x}$ :

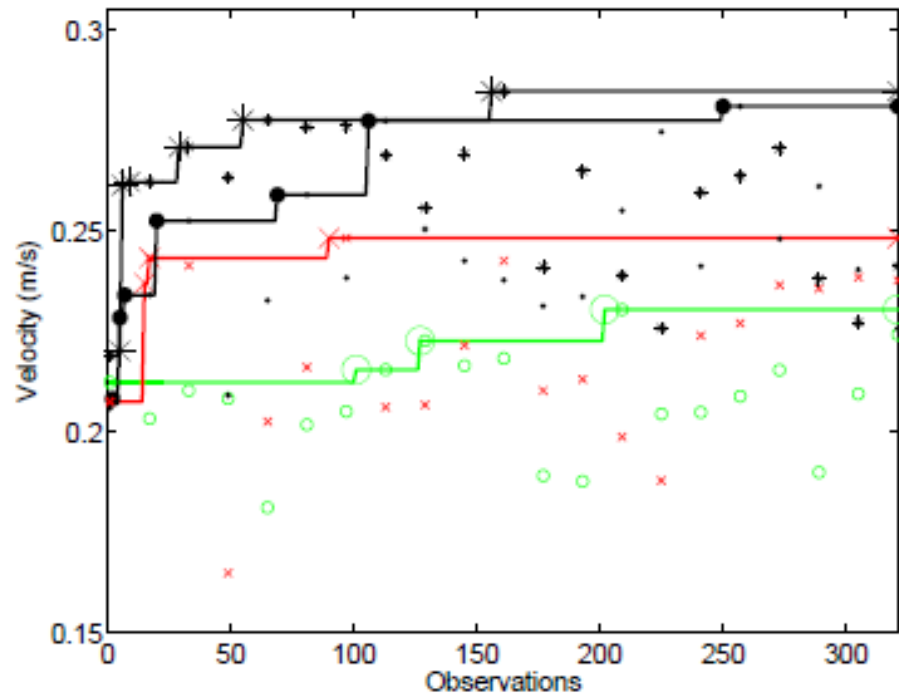
$$\mathbf{x}_{new} \leftarrow \operatorname{argmax}_{\mathbf{x}} \frac{k(\mathbf{x}, \mathbf{x})}{\max_{\mathbf{x}' \in X} f(\mathbf{x}') - m(\mathbf{x})}$$

– Evaluate  $f(\mathbf{x}_{new})$  by observing speed of robot with parameters set to  $\mathbf{x}_{new}$

– Update Gaussian process:

- $\mathbf{X} \leftarrow \mathbf{X} \cup \{\mathbf{x}_{new}\}$  and  $\mathbf{y} \leftarrow \mathbf{y} \cup f(\mathbf{x}_{new})$
- $m(\cdot) \leftarrow k(\cdot, \mathbf{X})(\mathbf{K} + \sigma^{-2}\mathbf{I})^{-1}\mathbf{y}$
- $k(\cdot, \cdot) \leftarrow k(\cdot, \cdot) - k(\cdot, \mathbf{X})(\mathbf{K} + \sigma^2\mathbf{I})^{-1}k(\mathbf{X}, \cdot)$

# Results



(●) GP w/MPI	(*) GP w/MPI	(×) H.C1mb	(○) U.Rand
0.281 m/s	0.285 m/s	0.248 m/s	0.230 m/s
$\sigma_f^2 = 0.06$	$\sigma_f^2 = 0.6$		

Kernel (covariance function):

$$k(\mathbf{x}, \mathbf{x}') = \sigma_f^2 e^{-\frac{1}{2}(\mathbf{x}-\mathbf{x}')^T S (\mathbf{x}-\mathbf{x}')}$$