# CS485/685 Machine Learning
# Lecture 4: Jan 12, 2012

Linear Regression

[B] Section 3.1

CS485/685 (c) 2012 P. Poupart 1

# Linear model for regression

- Simplest form of regression
- Picture:

CS485/685 (c) 2012 P. Poupart 2

# Problem

- Data: $\{(\boldsymbol{x_1}, t_1), (\boldsymbol{x_2}, t_2), \dots, (\boldsymbol{x_N}, t_N)\}$
  - $\boldsymbol{x} = <x_1, x_2, \dots, x_D>$: input vector
  - $t$: target (continuous value)
- Problem: find hypothesis $h$ that maps $\boldsymbol{x}$ to $t$
  - Assume that $h$ is linear:

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

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

# Optimization

- Find best $w$ that minimizes Euclidean loss

$$\boldsymbol{w}^* = argmin_{\boldsymbol{w}} \frac{1}{2}\sum_{n=1}^{N}\left(t_n - \boldsymbol{w}^T\begin{pmatrix} 1 \\ \boldsymbol{x_n} \end{pmatrix}\right)^2$$

- Convex optimization problem
  $\implies$ unique optimum (global)

# Solution

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

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

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

- This is a linear system in $w$, therefore we rewrite it as $Aw = b$

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

# Solution

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

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

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

# Picture

# Regularization

- Least square solution may not be stable
  - i.e., slight perturbation of the input may cause a dramatic change in the output
  - Form of **overfitting**

# Example 1

- Training data: $\bar{x}_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$    $\bar{x}_2 = \begin{pmatrix} 1 \\ \epsilon \end{pmatrix}$

$$t_1 = 1 \qquad t_2 = 1$$

- $A =$

- $A^{-1} =$            $b =$

- $w =$

# Example 2

- Training data: $\bar{x}_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$    $\bar{x}_2 = \begin{pmatrix} 1 \\ \epsilon \end{pmatrix}$

$$t_1 = 1 + \epsilon \quad t_2 = 1$$

- $A =$

- $A^{-1} =$            $b =$

- $w =$

# Picture

# Regularization

- Idea: favor smaller values
- Tikhonov regularization: add $\left|\left|w\right|\right|_2^2$ as a penalty term
- Ridge regression:

$$w^* = argmin_w \frac{1}{2}\sum_{n=1}^{N} \left(t_n - w^T\bar{x}_n\right)^2 + \frac{\lambda}{2}\left|\left|w\right|\right|_2^2$$

where $\lambda$ is a weight to adjust the importance of the penalty

# Regularization

- Solution: $(\lambda I + A)w = b$

- Notes
  - Without regularization: eigenvalues of linear system may be arbitrarily close to 0 and the inverse may have arbitrarily large eigenvalues.
  - With Tikhonov regularization, eigenvalues of linear system are $\geq \lambda$ and therefore bounded away from 0. Similarly, eigenvalues of inverse are bounded above by $1/\lambda$.

CS485/685 (c) 2012 P. Poupart                    13

# Regularized Examples

Example 1                         Example 2

CS485/685 (c) 2012 P. Poupart                    14

# Generalized Linear Regression

- How can we do non-linear regression while using the same machinery?

- 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
$$H = \{x \to w_0\phi_0(x) + w_1\phi_1(x) + w_2\phi_2(x)|w_i \in \Re\}$$

- 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.

# Next class

- Linear regression by
  - Maximum likelihood estimation (ML)
  - Maximum a posteriori estimation (MAP)
  - Bayesian learning

CS485/685 (c) 2012 P. Poupart                    19