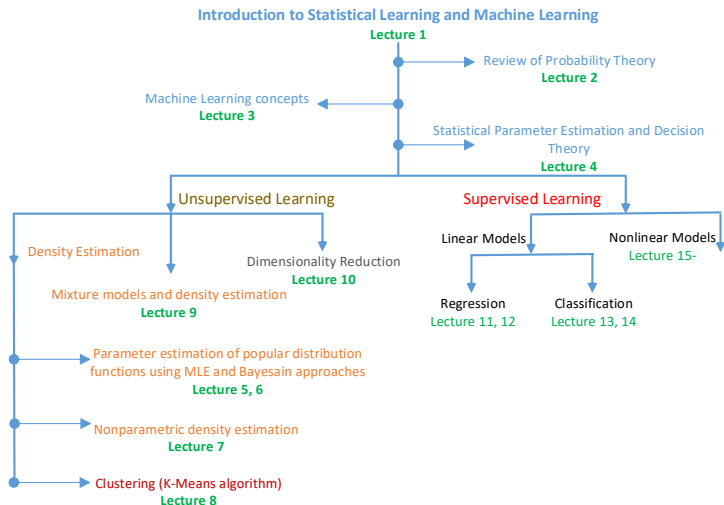


Statistical Learning and Machine Learning

Lecture 12 - Linear Models for Regression 2

October 6, 2021

Course overview and where do we stand

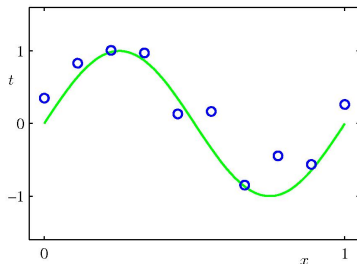


Objectives of the lecture

- Linear models for regression
 - Review of least squares approach
 - Sequential updates of weight vector in least squares
 - Multiple inputs
 - Bayesian linear regression

Goal of Regression

- The **goal of regression** is to predict the value of one or more continuous *target* variables t given the value of a D -dimensional vector \mathbf{x} of *input* variables.
- **Supervised learning**: Training data consisting of N observations $\{\mathbf{x}_n\}$ for $n = 1, \dots, N$ along with target values $\{t_n\}$ are available.
- **Output**: A function $y(\mathbf{x})$ whose values for new inputs \mathbf{x} constitute the predictions t .



Linear Basis Function Models

Two ways to define **linear models**:

- Linear w.r.t. to both input \mathbf{x} and parameters \mathbf{w}

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1x_1 + \dots, + w_Dx_D = w_0 + \sum_{j=1}^D w_jx_j = \mathbf{w}^T \mathbf{x}$$

where $\mathbf{x} = (x_0, x_1, \dots, x_D)^T$ and $x_0 = 1$ in the final form.

- **Non-linear** functions $\phi_j(\cdot), j = 1, \dots, M - 1$ (**basis functions**) w.r.t. to the input, with $\phi_0(\mathbf{x}) = 1$

$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j\phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$

where $\boldsymbol{\phi} = (\phi_0, \dots, \phi_{M-1})^T$ and $\mathbf{w} = (w_0, \dots, w_{M-1})^T$.

Examples

Linear basis function in both \mathbf{w} and \mathbf{x}

$$\phi_j(\mathbf{x}) = x_j, \text{ for } j = 1, \dots, D$$

with $\phi_0(\mathbf{x}) = 1$.

In vector form, we get

$$\phi(\mathbf{x}) = (\phi_0, \dots, \phi_{M-1})^T = \mathbf{x}$$

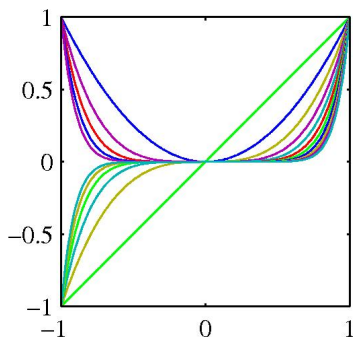
where $M - 1 = D$. In this case, the **output function** becomes

$$y(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \phi(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$$

Example basis functions

Polynomial basis function:

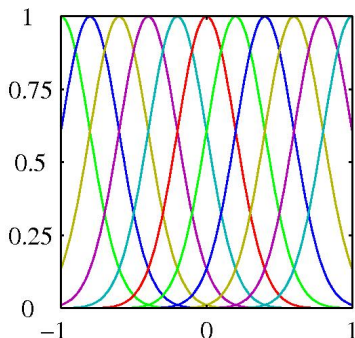
$$\phi_j(x) = x^j$$



Example basis functions

Radial basis function:

$$\phi_j(x) = \exp \left\{ -\frac{(x - \mu_j)^2}{2s^2} \right\}$$

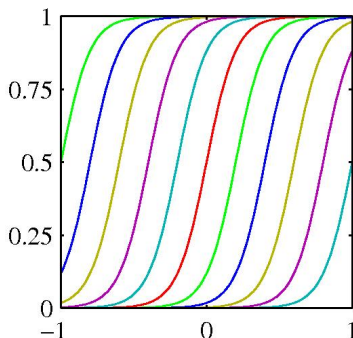


Example basis functions

Sigmoidal basis function:

$$\phi_j(x) = \sigma\left(\frac{x - \mu_j}{s}\right)$$

where $\sigma(a) = 1/(1 + \exp(-a))$.



Least Squares

- **Goal:** Given a set of i.i.d. data points $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ and the corresponding t_n , $n = 1, \dots, N$, we want to estimate the parameters \mathbf{w} of the regression model.
- Which model?

$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$

Least Squares

- **Cost function to be minimized:** Sum of the squares of the individual errors:

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left(y(\mathbf{x}_n, \mathbf{w}) - t_n \right)^2 = \frac{1}{2} \sum_{n=1}^N \left(\mathbf{w}^T \phi(\mathbf{x}_n) - t_n \right)^2$$

- **Minimization:** By setting $\frac{\partial E_D(\mathbf{w})}{\partial \mathbf{w}} = 0$:

$$\mathbf{w} = \left(\Phi^T \Phi \right)^{-1} \Phi^T \mathbf{t}$$

where $\Phi \in \mathbb{R}^{N \times M}$ is formed by using the vectors $\phi(\mathbf{x}_n)$ as rows and $\mathbf{t} = [t_1, \dots, t_N]^T$.

Linear Regression: Sequential updates

We obtain a **sequential (or on-line)** learning algorithm for updating \mathbf{w} by applying **stochastic gradient descent (SGD)**:

- If the error function has the form $E(\mathbf{w}) = \sum_n E_n(\mathbf{w})$ then:

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n$$

where τ denotes the iteration number, η is a learning rate parameter and ∇ is the gradient operator.

- For the least-squares error case:

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} + \eta \left(t_n - \mathbf{w}^{(\tau)T} \phi(\mathbf{x}_n) \right) \phi(\mathbf{x}_n).$$

The value of η needs to be chosen appropriately to ensure convergence of the algorithm.

Linear Regression: Sequential updates

A more effective sequential learning algorithm for updating \mathbf{w} is based on the [Newton-Raphson](#) iterative optimization scheme, which uses a local quadratic approximation of the log-likelihood function:

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \mathbf{H}^{-1} \nabla E_n(\mathbf{w})$$

where \mathbf{H} is the [Hessian matrix](#) having elements $H_{ij} = \frac{\partial^2 E_n(\mathbf{w})}{\partial w_i \partial w_j}$.

$$\nabla E_n(\mathbf{w}) = \sum_{n=1}^N (\mathbf{w}^T \phi(\mathbf{x}_n) - t_n) \phi(\mathbf{x}_n) = \Phi^T \Phi \mathbf{w} - \Phi^T \mathbf{t}$$

$$\mathbf{H} = \nabla \nabla E_n(\mathbf{w}) = \sum_{n=1}^N \phi_n \phi_n^T = \Phi^T \Phi$$

Linear Regression: Sequential updates

Because $E_n(\mathbf{w})$ is a quadratic function, the Newton-Raphson method gives the exact solution in one step if applied on the entire data set.

The update takes the form:

$$\begin{aligned}\mathbf{w}^{(\tau+1)} &= \mathbf{w}^{(\tau)} - (\Phi^T \Phi)^{-1} \left(\Phi^T \Phi \mathbf{w}^{(\tau)} - \Phi^T \mathbf{t} \right) \\ &= (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}\end{aligned}$$

Linear Regression: Multiple outputs

When we want to regress to multiple target values $\mathbf{t}_n \in \mathbb{R}^K$:

$$\mathbf{y}(\mathbf{x}_n, \mathbf{W}) = \mathbf{W}^T \phi(\mathbf{x}_n)$$

where $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_K] \in \mathbb{R}^{M \times K}$.

When the targets are multi-dimensional random variables \mathbf{t} , we need to estimate $p(\mathbf{t}|\mathbf{x}, \mathbf{W}, \beta)$. We will consider the case:

$$p(\mathbf{t}|\mathbf{x}, \mathbf{W}, \beta) = \mathcal{N}(\mathbf{t}|\mathbf{W}^T \phi(\mathbf{x}), \beta^{-1} \mathbf{I}).$$

Linear Regression: Multiple outputs

Given a set of i.i.d. data points $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$ and the corresponding target vectors $\mathbf{T} \in \mathbb{R}^{N \times K}$, having as its n -th row the target vector \mathbf{t}_n .

The log-likelihood function is given by:

$$\begin{aligned}\ln p(\mathbf{T}|\mathbf{X}, \mathbf{W}, \beta) &= \sum_{n=1}^N \ln \mathcal{N}(\mathbf{t}_n | \mathbf{W}^T \phi(\mathbf{x}), \beta^{-1} \mathbf{I}) \\ &= \frac{NK}{2} \ln \left(\frac{\beta}{2\pi} \right) - \frac{\beta}{2} \sum_{n=1}^N \|\mathbf{t}_n - \mathbf{W}^T \phi(\mathbf{x})\|^2.\end{aligned}$$

Setting $\frac{\partial \ln p(\mathbf{T}|\mathbf{X}, \mathbf{W}, \beta)}{\partial \mathbf{W}} = 0$:

$$\mathbf{W}_{ML} = \Phi^\dagger \mathbf{T}.$$

Bayesian Linear Regression

We consider the model parameters \mathbf{w} as parameters drawn by a distribution $p(\mathbf{w})$ (we assume β is known).

The **likelihood function** $p(\mathbf{t}|\mathbf{w})$ is an exponential of quadratic function of \mathbf{w} :

$$p(\mathbf{t}|\mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | \mathbf{w}^T \phi(\mathbf{x}), \beta^{-1}).$$

The corresponding **conjugate prior** has the form:

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0).$$

Bayesian Linear Regression

Reminder: posterior \propto likelihood \times prior

The posterior distribution is:

$$p(\mathbf{w}|\mathbf{t}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N)$$

where

$$\begin{aligned}\mathbf{m}_N &= \mathbf{S}_N(\mathbf{S}_0^{-1}\mathbf{m}_0 + \beta\Phi^T\mathbf{t}) \\ \mathbf{S}_N^{-1} &= \mathbf{S}_0^{-1} + \beta\Phi^T\Phi.\end{aligned}$$

Cases:

- If $\mathbf{S}_0 = \alpha^{-1}\mathbf{I}$ with $\alpha \rightarrow 0$, then $\mathbf{m}_N \rightarrow \mathbf{w}_{ML}$
- If $N = 0$ then $p(\mathbf{w}|\mathbf{t})$ becomes $p(\mathbf{w})$

Bayesian Linear Regression

When $p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|0, \alpha^{-1}\mathbf{I})$:

$$\begin{aligned}\mathbf{m}_N &= \beta \mathbf{S}_N \Phi^T \mathbf{t} \\ \mathbf{S}_N^{-1} &= \alpha \mathbf{I} + \beta \Phi^T \Phi\end{aligned}$$

and

$$\ln p(\mathbf{w}|\mathbf{t}) = -\frac{\beta}{2} \sum_{n=1}^N \left(t_n - \mathbf{w}^T \phi(\mathbf{x}_n) \right)^2 - \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} + \text{const.}$$

Maximization of the posterior distribution w.r.t. \mathbf{w} is equivalent to the minimization of the sum-of-squares error with an additional quadratic regularization term with regularization term $\lambda = \alpha/\beta$.

Bayesian Linear Regression: Predictive distribution

In practice, we want to make predictions of t for new values of \mathbf{x} :

$$p(t|\mathbf{t}, \alpha, \beta) = \int p(t|\mathbf{w}, \beta) p(\mathbf{w}|\mathbf{t}, \alpha, \beta) d\mathbf{w}$$

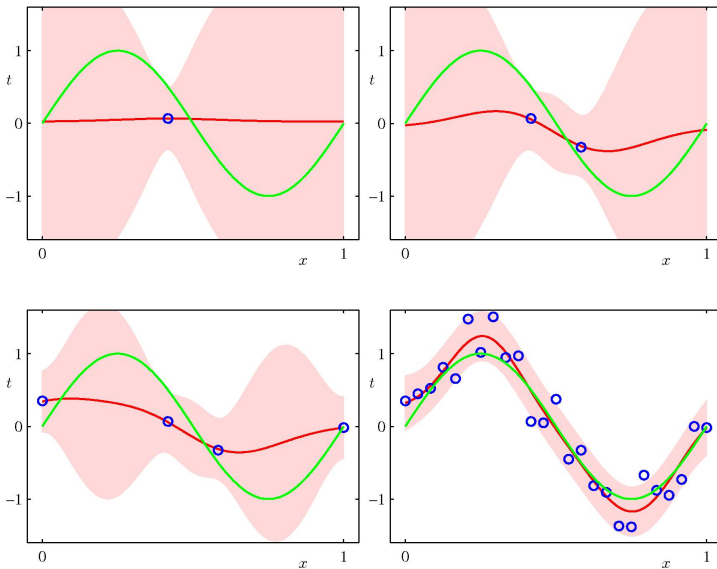
- **Conditional distribution:** $p(t|\mathbf{w}, \beta) = \mathcal{N}(t|y(\mathbf{x}, \mathbf{w}), \beta^{-1})$
- **Posterior:** $p(\mathbf{w}|\mathbf{t}, \alpha, \beta) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N)$

The result is a Gaussian:

$$p(t|\mathbf{x}, \mathbf{t}, \alpha, \beta) = \mathcal{N}\left(t|\mathbf{m}_N^T \phi(\mathbf{x}), \sigma_N^2(\mathbf{x})\right)$$

where: $\sigma_N^2(\mathbf{x}) = \beta^{-1} + \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x})$.

Bayesian Linear Regression: Predictive distribution



Bayesian Linear Regression: Predictive distribution

