# Machine Learning

Lecture 4,5,6,7

Indian Institute of Information Technology Dharwad

# Decision Theory

- Using probability theory to make optimal decisions
- Input vector $\mathbf{x}$, target vector $\mathbf{t}$
    - *Regression:* $\mathbf{t}$ is continuous
    - *Classification:* $\mathbf{t}$ will consist of class labels
- Summary of uncertainty associated is given by $p(\mathbf{x},t)$
- *Inference problem* is to obtain $p(\mathbf{x},t)$ from data
- *Decision:* make specific prediction for value of $\mathbf{t}$ and take specific actions based on $\mathbf{t}$

# Medical Diagnosis Problem

- X-ray image of patient
- Whether patient has cancer or not
- Input vector $\mathbf{x}$ is set of pixel intensities
- Output variable $t$ represents whether cancer or not $C_1$ is cancer and $C_2$ is absence of cancer
- General inference problem is to determine $p(x, C_k)$ which gives most complete description of situation
- In the end we need to decide whether to give treatment or not. Decision theory helps do this

# Bayes Decision

- How do probabilities play a role in making a decision?

- Given input $\mathbf{x}$ and classes $C_k$ using Bayes theorem

$$p(C_k \mid \mathbf{x}) = \frac{p(\mathbf{x} \mid C_k) p(C_k)}{p(\mathbf{x})}$$

- Quantities in Bayes theorem can be obtained from $p(x, C_k)$ either by marginalizing or conditioning wrt appropriate variable
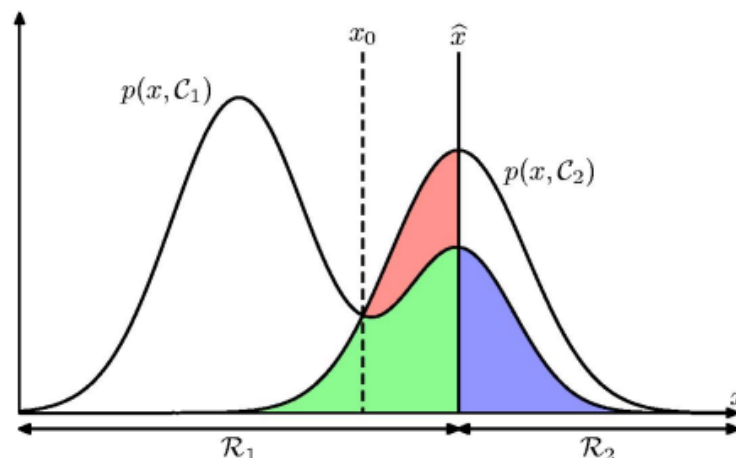
# Minimizing Expected Error

- Probability of mistake (2-class)

$$P(error) = p(x \, \varepsilon \, R_1, C_2) + p(x \, \varepsilon \, R_2, C_1)$$

$$= \int_{R_1} p(\mathrm{x}, C_2) d\mathrm{x} + \int_{R_2} p(\mathrm{x}, C_1) d\mathrm{x}$$



- Minimum error decision rule
  - For a given $\mathrm{x}$ choose class for which integrand is smaller
  - Since $p(\mathrm{x}, C_k) = p(C_k | \mathrm{x}) p(\mathrm{x})$, choose class for which *a posteriori* probability is highest
  - Called Bayes Classifier

Single input variable $x$

If priors are equal, decision is based on class-conditional densities $p(x|C_k)$

# Minimizing Expected Loss

- Unequal importance of mistakes
- Medical Diagnosis
- Loss or Cost Function given by Loss Matrix
- Utility is negative of Loss
- Minimize Average Loss

$$E[L] = \sum_k \sum_j \int_{R_j} L_{kj} p(\mathbf{x}, C_k) d\mathbf{x}$$
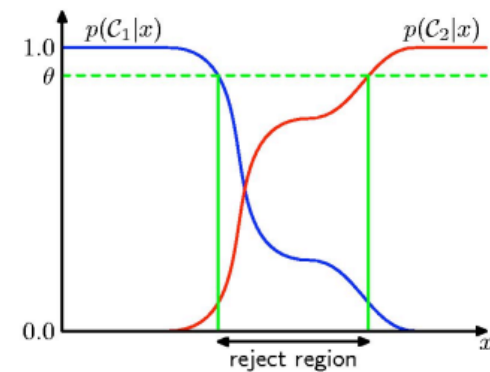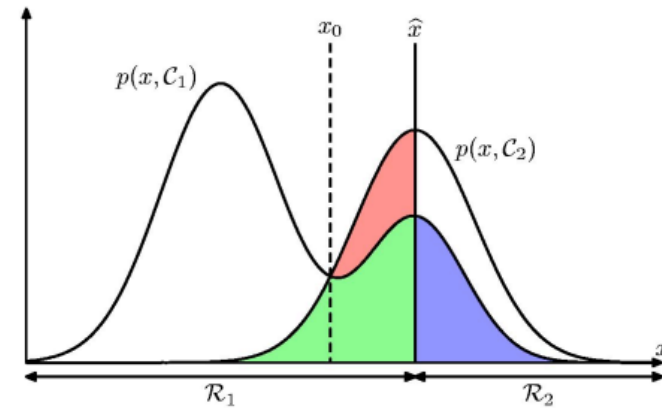
- Minimum Loss Decision Rule

    – Choose class for which $\sum_k L_{kj} p(C_k \mid \mathbf{x})$ is minimum

    – Trivial once we know a posteriori probabilities

## Loss Function for Cancer Decision

Decision Made

|  |  | cancer | normal |
|---|---|---|---|
| True Class | cancer | 0 | 1000 |
|  | normal | 1 | 0 |

# Reject Option

- Decisions can be made when *a posteriori* probabilities are significantly less than unity or joint probabilities have comparable values

- Avoid making decisions on difficult cases

# Inference and Decision

- Classification problem broken into two separate stages
  - Inference, where training data is used to learn a model for $p(C_k, x)$
  - Decision, use posterior probabilities to make optimal class assignments
- Alternatively can learn a function that maps inputs directly into labels
- Three distinct approaches to Decision Problems
  1. Generative
  2. Discriminative
  3. Discriminant Function

# 1. Generative Models

- First solve inference problem of determining class-conditional densities $p(\mathbf{x}|C_k)$ for each class separately

- Then use Bayes theorem to determine posterior probabilities

$$p(C_k \mid \mathbf{x}) = \frac{p(\mathbf{x} \mid C_k)\, p(C_k)}{p(\mathbf{x})}$$

- Then use decision theory to determine class membership

# 2. Discriminative Models

- First solve inference problem to determine posterior class probabilities $p(C_k|x)$

- Use decision theory to determine class membership

# 3. Discriminant Functions

- Find a function $f(\mathbf{x})$ that maps each input $\mathbf{x}$ directly to class label
  - In two-class problem, $f(.)$ is binary valued
    - $f=0$ represents class $C_1$ and $f=1$ represents class $C_2$

- Probabilities play no role
  - No access to posterior probabilities $p(C_k|\mathbf{x})$

# Need for Posterior Probabilities

- Minimizing risk
  - Loss matrix may be revised periodically as in a financial application
- Reject option
  - Minimize misclassification rate, or expected loss for a given fraction of rejected points
- Compensating for class priors
  - When far more samples from one class compared to another, we use a balanced data set (otherwise we may have 99.9% accuracy always classifying into one class)
  - Take posterior probabilities from balanced data set, divide by class fractions in the data set and multiply by class fractions in population to which the model is applied
  - Cannot be done if posterior probabilities are unavailable
- Combining models
  - X-ray images ($x_I$) and Blood tests ($x_B$)
  - When posterior probabilities are available they can be combined using rules of probability
  - Assume feature independence $p(x_I, x_B | C_k) = p(x_I | C_k) \, p(x_B, | C_k)$  [Naïve Bayes Assumption]
  - Then

$$p(C_k | x_I, x_B) \; \alpha \;\; p(x_I, x_B | C_k) p(C_k)$$

$$\alpha \;\; p(x_I | C_k) \, p(x_B, | C_k) \, p(C_k)$$

$$\alpha \;\; p(C_k | x_I) \, p(C_k | x_B) / p(C_k)$$

  - Need $p(C_k)$ which can be determined from fraction of data points in each class. Then need to normalize resulting probabilities to sum to one
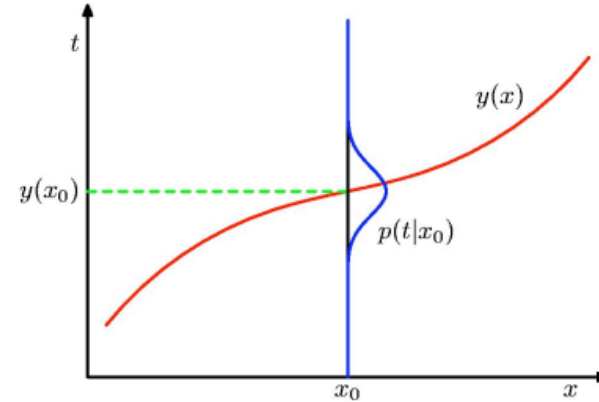
# Loss Functions for Regression

- Curve fitting can also use a loss function
- Regression decision is to choose a specific estimate $y(\mathbf{x})$ of $t$ for a given $\mathbf{x}$
- Incur loss $L(t,y(\mathbf{x}))$
- Squared loss function
  $L(t,y(x))=\{y(x)-t\}^2$
- Minimize expected loss
  $E[L]=\displaystyle\iint L(t,y(\mathbf{x}))p(\mathbf{x},t)dxdt$

  Taking derivative and setting equal to zero yields a solution
  $y(x)=E_t[t|x]$



Regression function $y(x)$, which minimizes the expected squared loss, is given by the mean of the conditional distribution $p(t|x)$

# Loss function for Regression

$$EPE(f) = E\left[(Y - f(X))^2\right]$$

$$= \int \int [y - f(x)]^2 \Pr(y, x) \, dy \, dx$$

$$= \int_x \int_y [y - f(x)]^2 \Pr(y|x) \Pr(x) \, dy \, dx \qquad \Pr(X, Y) = \Pr(Y|X) \Pr(X)$$

$$EPE(f) = E_X E_{Y|X}([Y - f(X)]^2 | X).$$

# Loss function for Regression

Notice that by conditioning on $X$, we have freed the dependency of the function $f$ on $X$ and since

the quantity $[Y - f]^2$ is convex, there is a unique solution. We can now minimize to solve for $f$

$$f(x) = \arg\min_f E_{Y|X}([Y - f]^2 | X = x)$$

$$\Rightarrow \frac{\partial}{\partial f} \int [Y - f]^2 \Pr(y|x)\, dy = 0$$

$$= \int \frac{\partial}{\partial f}[y - f]^2 \Pr(y|X)\, dy = 0$$

$$= 2\int yPr(y|x)dy = 2f \int Pr(y|x)\, dy = 0$$

$$\Rightarrow 2E[Y|X] = 2f$$

$$\Rightarrow f = E[Y|X = x].$$
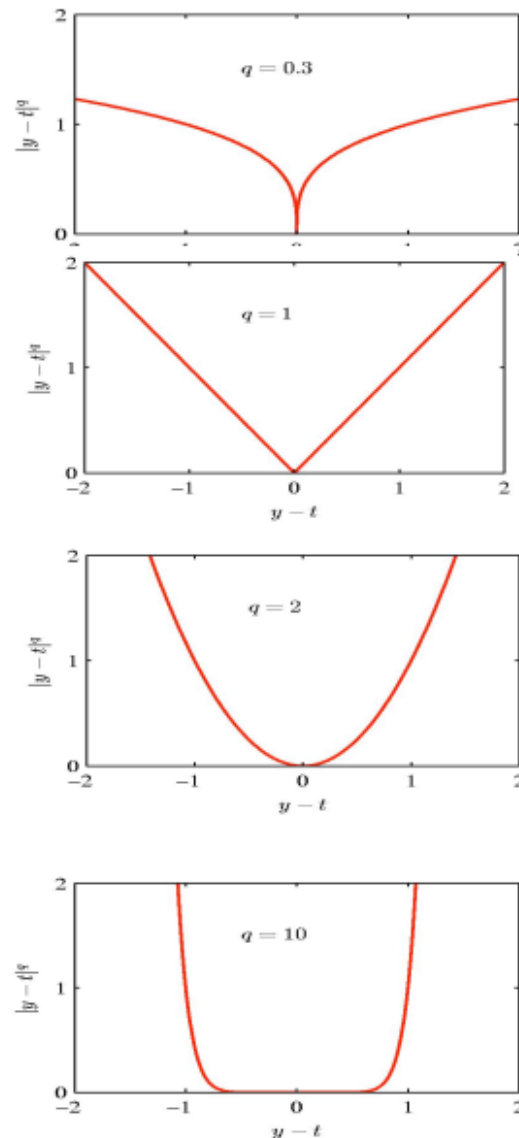
# Inference and Decision for Regression

- Three distinct approaches (decreasing complexity)
- Analogous to those for classifiction
  1. Determine joint density $p(\mathbf{x},t)$

     Then normalize to find conditional density $p(t|\mathbf{x})$

     Finally marginalize to find conditional mean $E_t[t|\mathbf{x}]$
  2. Solve inference problem of determining conditional density $p(t|\mathbf{x})$

     Marginalize to find conditional mean
  3. Find regression function $y(\mathbf{x})$ directly from training data

# Minkowski Loss Function



- Squared Loss is not only possible choice for regression
- Important example concerns multimodal $p(t|\mathbf{x})$
- Minkowski Loss

$$\mathbb{E}[L_q] = \iint |y(\mathbf{x}) - t|^q p(\mathbf{x}, t) \, d\mathbf{x} \, dt$$

- Minimum of $\mathbb{E}[L_q]$ is given by
  - conditional mean for $q=2$,
  - conditional median for $q=1$ and
  - conditional mode for $q \rightarrow 0$

# Linear Regression with Basis Function

# The regression task

- It is a supervised learning task
- Goal of regression:
  - predict value of one or more <u>target</u> variables $t$
  - given <u>$d$-dimensional</u> vector $\mathbf{x}$ of input variables
  - With dataset of known inputs and outputs
    - $(\mathbf{x}_1,t_1),\ ..(\mathbf{x}_N,t_N)$
    - Where $\mathbf{x}_i$ is an input (possibly a vector) known as the predictor
    - $t_i$ is the target output (or response) for case $i$ which is real-valued
  - Goal is to predict $t$ from $\mathbf{x}$ for some future test case
    - We are not trying to model the distribution of $\mathbf{x}$
  - We dont expect predictor to be a linear function of $\mathbf{x}$
    - So ordinary linear regression of inputs will not work
    - We need to allow for a nonlinear function of $\mathbf{x}$
    - We don't have a theory of what form this function to take $_3$

# ML Terminology

- ## Regression
  - Predict a numerical value $t$ given some input
    - Learning algorithm has to output function $f : \mathrm{R}^n \rightarrow \mathrm{R}$
      - where $n =$no of input variables

- ## Classification
  - If $t$ value is a label (categories): $f : \mathrm{R}^n \rightarrow \{1,..,k\}$

- ## Ordinal Regression
  - Discrete values, ordered categories

# Polynomial Curve Fitting with a Scalar
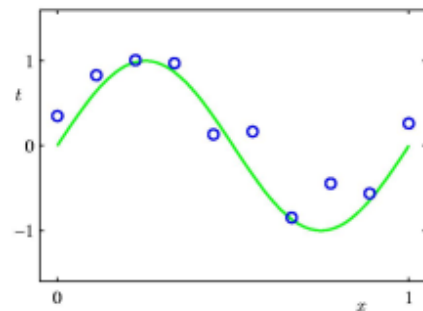
– With a <u>single</u> input variable $x$

$$y(x,\mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \ldots + w_M x^M = \sum_{j=0}^{M} w_j x^j$$

$M$ is the order of the polynomial,

$x^j$ denotes $x$ raised to the power $j$,

Coefficients $w_0, \ldots, w_M$ are collectively denoted by vector $\mathbf{w}$

Training data set
$N=10$, Input $x$, target $t$

– Task: Learn $\mathbf{w}$ from training data $D = \{(x_i, t_i)\}, \ i = 1, .., N$

- Can be done by minimizing an error function that minimizes the misfit between $y(x,\mathbf{w})$ for any given $\mathbf{w}$ and training data

- One simple choice of error function is sum of squares of error between predictions $y(x_n,\mathbf{w})$ for each data point $x_n$ and corresponding target values $t_n$ so that we minimize

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} \left\{ y(x_n, \mathbf{w}) - t_n \right\}^2$$

- It is zero when function $y(x,\mathbf{w})$ passes exactly through each training data point

# Regression with multiple inputs

- Generalization
  - Predict value of continuous target variable $t$ given value of $D$ input variables $\mathrm{x}=[x_1,..x_D]$
  - $t$ can also be a set of variables (multiple regression)
  - Linear functions of adjustable parameters
    - Specifically linear combinations of <u>nonlinear</u> functions of input variable
- Polynomial curve fitting is good only for:
  - Single input variable scalar $x$
  - It cannot be easily generalized to several variables, as we will see

# Simplest Linear Model with $D$ inputs

- Regression with $D$ input variables

$$y(\mathrm{x,w}) = w_0 + w_1 x_1 + .. + w_d\, x_D = \mathrm{w}^T\mathrm{x}$$
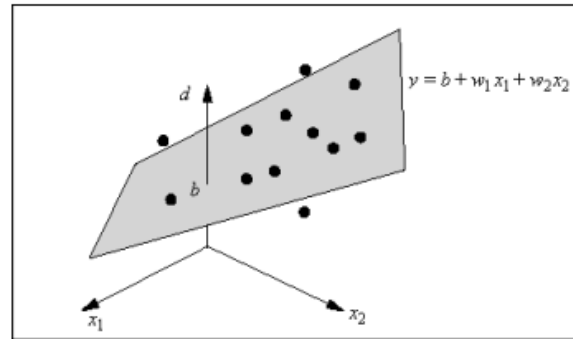
This differs from
Linear Regression with <u>one</u> variable
and Polynomial Reg with <u>one</u> variable

*where* $\mathrm{x} = (x_1,..,x_D)^\mathrm{T}$ are the input variables

- Called Linear Regression since it is a linear function of
  - parameters $w_0,..,w_D$
  - input variables $x_1,..,x_D$

- Significant limitation since it is a linear function of input variables
  - In the one-dimensional case this amounts a straight-line fit (degree-one polynomial)
  - $y(x,\mathrm{w}) = w_0 + w_1 x$

# Fitting a Regression Plane

- Assume $t$ is a function of inputs $x_1, x_2, \ldots x_D$
  Goal: find best linear regressor of $t$ on all inputs
  - Fitting a hyperplane through $N$ input samples
  - For $D = 2$:



$y = b + w_1 x_1 + w_2 x_2$

| $x_1$ | $x_2$ | $t$ |
|---|---|---|
| 1 | 2 | 2 |
| 2 | 5 | 1 |
| 2 | 3 | 2 |
| 2 | 2 | 2 |
| 3 | 4 | 1 |
| 3 | 5 | 3 |
| 4 | 6 | 2 |
| 5 | 5 | 3 |
| 5 | 6 | 4 |
| 5 | 7 | 3 |
| 6 | 8 | 4 |
| 7 | 6 | 2 |
| 8 | 4 | 4 |
| 8 | 9 | 3 |
| 9 | 8 | 4 |

- Being a linear function of input variables imposes limitations on the model
  - Can extend class of models by considering fixed nonlinear functions of input variables

# Basis Functions

- In many applications, we apply some form of fixed-preprocessing, or feature extraction, to the original data variables

- If the original variables comprise the vector $\mathbf{x}$, then the features can be expressed in terms of basis functions $\{\phi_j(\mathbf{x})\}$

  - By using nonlinear basis functions we allow the function $y(\mathbf{x},\mathbf{w})$ to be a nonlinear function of the input vector $\mathbf{x}$

    - They are linear functions of parameters (gives them simple analytical properties), yet are nonlinear wrt input variables

# Linear Regression with $M$ Basis Functions

- Extended by considering nonlinear functions of input variables

$$y(\mathbf{x},\mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x})$$

  - where $\phi_j(\mathbf{x})$ are called Basis functions
  - We now need $M$ weights for basis functions instead of $D$ weights for features
  - With a dummy basis function $\phi_0(\mathbf{x}){=}1$ corresponding to the bias parameter $w_0$, we can write

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

  - where $\mathbf{w}{=}(w_0, w_1, .., w_{M-1})$ and $\Phi{=}(\phi_0, \phi_1, .., \phi_{M-1})^T$

- Basis functions allow non-linearity with $D$ input variables

# Choice of Basis Functions

- Many possible choices for basis function:

    1. Polynomial regression

        - Good only if there is only one input variable

    2. Gaussian basis functions

    3. Sigmoidal basis functions

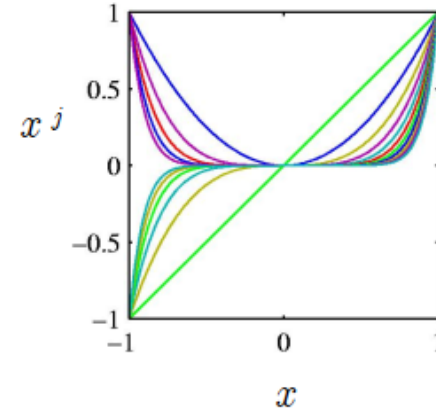    4. Fourier basis functions

    5. Wavelets

# 1. Polynomial Basis for one variable

- **Linear Basis Function Model**

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



- **Polynomial Basis (for single variable $x$)**
  $\phi_j(x) = x^j$ with degree $M$-1 polynomial

- **Disadvantage**
  - Global:
    - changes in one region of input space affects others
  - Difficult to formulate
    - Number of polynomials increases exponentially with $M$
  - Can divide input space into regions
    - use different polynomials in each region:
    - equivalent to spline functions

# Can we use Polynomial with $D$ variables?
## (Not practical!)

- **Consider (for a vector $\mathbf{x}$) the basis:** $\phi_j(\mathbf{x}) = \| \mathbf{x} \|^j = \left[ \sqrt{x_1^2 + x_2^2 + .. + x_d^2} \right]^j$
  - $\mathbf{x}=(2,1)$ and $\mathbf{x}=(1,2)$ have the same squared sum, so it is unsatisfactory
  - Vector is being converted into a scalar value thereby losing information
- **Better polynomial approach:**
  - Polynomial of degree $M$-1 has terms with variables taken none, one, two… $M$-1 at a time.
  - Use multi-index $j=(j_1,j_2,..j_D)$ such that $j_1+j_2+..j_D \leq M$-1
  - For a quadratic $(M=3)$ with three variables $(D=3)$

$$y(\mathbf{x}, \mathbf{w}) = \sum_{(j_1,j_2,j_3)} w_j \phi_j(\mathbf{x}) = w_0 + w_{1,0,0}x_1 + w_{0,1,0}x_2 + w_{0,0,1}x_3 + w_{1,1,0}x_1 x_2 + w_{1,0,1}x_1 x_3 +$$

$$w_{0,1,1}x_2 x_3 + w_{2,0,0}x_1^2 + w_{0,2,0}x_2^2 + w_{0,0,2}x_3^2$$

  - Number of quadratic terms is $1+D+D(D\text{-}1)/2+D$
  - For $D=46$, it is $1128$
- **Better to use Gaussian kernel, discussed next**

# Disadvantage of Polynomial

- Polynomials are *global* basis functions
  - Each affecting the prediction over the whole input space
- Often local basis functions are more appropriate

# Review and Derivations

# Linear Regression

- The simplest linear model for regression is one that involves a linear combination of the input variables

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x_1 + \ldots + w_D x_D \quad = \sum_{j=0}^{D} w_j x^j$$

where $\mathbf{x} = (x_1, \ldots, x_D)^{\mathrm{T}}$

- Extension of linear regression models

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

# Linear Regression

- Error Function:

$$E(w) = \frac{1}{2} \sum_{n=1}^{N} \left( t_n - \sum_{i=0}^{D} w_i x_i \right)^2$$

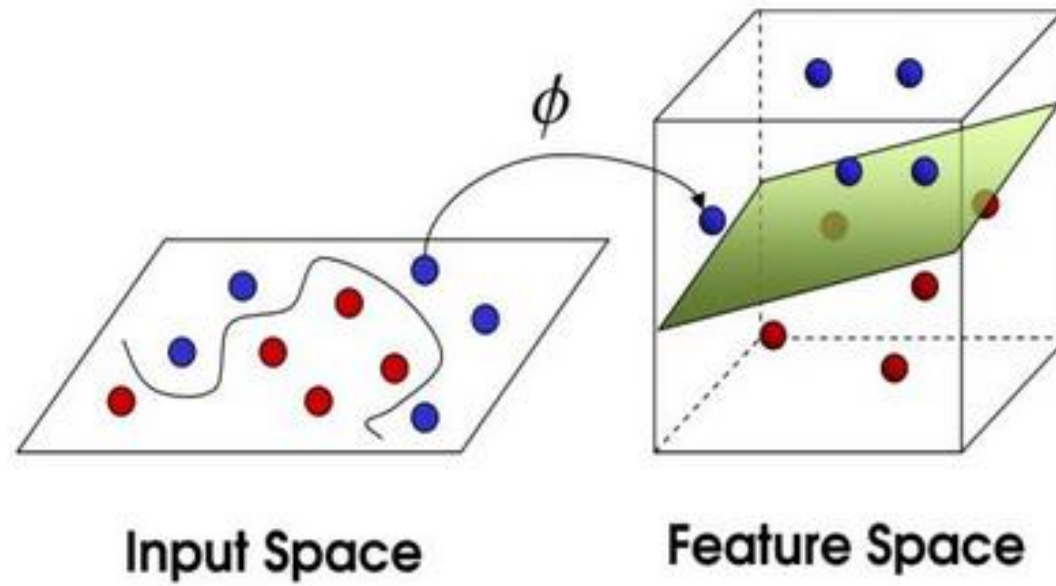$$E(w) = \frac{1}{2} \sum_{n=1}^{N} \left( t_n - w^T x \right)^2$$

$$E(w) = (t - Xw)^T (t - Xw)$$

- To find $w$ which maximizes this function, we differentiate it w.r.t. to $w$; we get

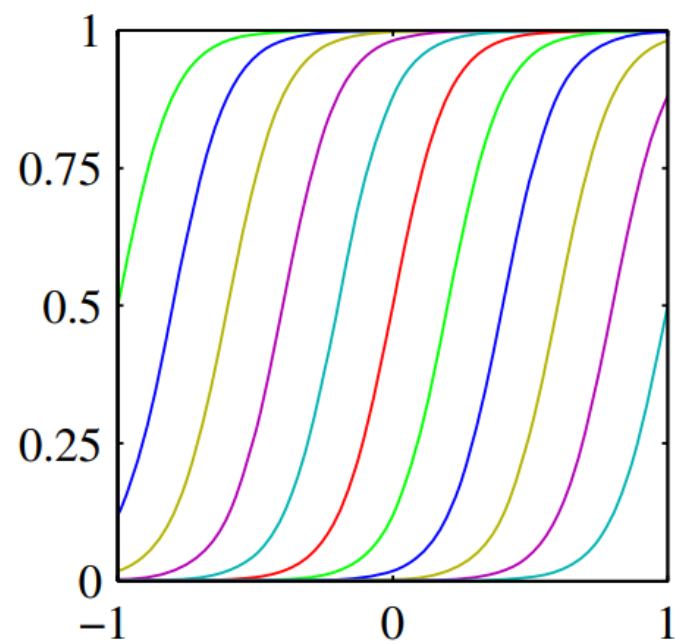$$X^T (t - Xw) = 0$$

$$w = (X^T X)^{-1} X^T t$$

# Basis Function Motivation



Input Space       Feature Space

# Basis Functions

# Basis Functions

- Gaussian Basis Function

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

- Sigmoid Function

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

$$\sigma(a) = \frac{1}{1 + \exp(-a)}$$

# Example 1: Linear Regression

- We specify the set of functions $\phi_1$, $\phi_2$, ...., $\phi_M$ from X to $\mathbb{R}$ and look for function in the form of linear combination

$$y(x) = \sum_{i=0}^{M} w_i \emptyset_i(x)$$

- Performing the regression then reduces to finding the real parameters $w_0$, $w_1$, ..., $w_M$
- When $x$ is only one dimension the **simplest basis function** would be:
$$\phi_0(x) = 1 \text{ and } \phi_1(x) = x$$

- This gives:

$$y(x) = \sum_{i=0}^{1} w_i \emptyset_i(x) = w_0 + w_1 x$$

- In the multidimensional case, where $x = < x_1, ..., x_D >$ we would take:
$$\phi_0(x) = 1, \phi_1(x) = x_1, \phi_2(x) = x_2, ......, \phi_D(x) = x_D$$

$$y(x) = \sum_{i=0}^{D} w_i \emptyset_i(x) = w_0 + w_1 x_1 + w_2 x_2 + ....... + w_D x_D$$

# Example 2: Polynomial Regression

- **Polynomial Regression (when $x$ is 1-dimensional):**
- Another possible choice to set $\quad \phi_i(\text{x}) = x^i$ for $\quad i = 1,2,\ldots,\text{M} \quad$ where M is the degree of the polynomial

$$y(x) = \sum_{i=0}^{M} w_i \emptyset_i(x) = \text{w}_0 + \text{w}_1 \text{ x}^1 + \text{w}_2 \text{ x}^2 + \ldots\ldots + \text{w}_M \text{ x}^M$$

- **Polynomial Regression (when $x$ is D-dimensional):** For example D = 3 and M = 2
- X = $<\text{x}_1,\text{x}_2,\text{x}_3>$ and $\text{x}_1 + \text{x}_2 + \text{x}_3 <= 2$

- $y(x) = \sum_{(x1,x2,x3)} w_i \emptyset_i(x) = \text{w}_0 + \text{w}_{(1,0,0)} \text{ x}_1 + \text{w}_{(0,1,0)} \text{ x}_2 + \text{w}_{(0,0,1)} \text{ x}_3 + \text{w}_{(1,1,0)} \text{ x}_1\text{x}_2$

$$+ \text{w}_{(1,1,0)} \text{ x}_1\text{x}_2 + \text{w}_{(0,1,1)} \text{ x}_2\text{x}_2 + \text{w}_{(1,0,1)} \text{ x}_1\text{x}_3 + \text{w}_{(2,0,0)} x_1^2 + w_{(0,2,0)}x_2^2 + w_{(0,0,2)}x_3^2$$

# Maximum likelihood and least squares

Deterministic function with Gaussian noise

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon$$

Alternatively, we can write

$$p(t|\mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t|y(\mathbf{x}, \mathbf{w}), \beta^{-1})$$

If we assume squares loss function, then optimal prediction for **x** ….???

$$\mathbb{E}[t|\mathbf{x}] = \int tp(t|\mathbf{x})\,\mathrm{d}t = y(\mathbf{x}, \mathbf{w})$$

# Likelihood Function for regression

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

$$\ln p(\mathbf{t}|\mathbf{w}, \beta) = \sum_{n=1}^{N} \ln \mathcal{N}(t_n|\mathbf{w}^{\mathrm{T}}\boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1})$$

$$= \frac{N}{2}\ln\beta - \frac{N}{2}\ln(2\pi) - \beta E_D(\mathbf{w})$$

$$E_D(\mathbf{w}) = \frac{1}{2}\sum_{n=1}^{N}\{t_n - \mathbf{w}^{\mathrm{T}}\boldsymbol{\phi}(\mathbf{x}_n)\}^2$$

# Linear Regression using basis function

- Error Function:

$$E(w) = \frac{1}{2} \sum_{n=1}^{N} (t_n - \sum_{i=0}^{M} w_i \emptyset(x_i))^2$$

$$E(w) = \frac{1}{2} \sum_{n=1}^{N} (t_n - w^T \emptyset(x))^2$$

- Where,     $w = (w_0, w_1, \ldots, w_M)^T$
- To represent the error function in vector form we get:

$$E(w) = \frac{1}{2} (t - Qw)^T (t - Qw)$$

# Linear Regression using basis function

- To represent the error function in vector form we get:

$$E(w) = \frac{1}{2}(t - Qw)^T(t - Qw)]$$

Where $t = (t_1, t_2, ...., t_N)^T$, $w = (w_0, w_1, ..., w_M)^T$ and $Q$ is as follows:

$$Q = \begin{pmatrix} \phi_0(x_1) & \phi_1(x_1) & \cdots & \phi_M(x_1) \\ \phi_0(x_2) & \phi_1(x_2) & & \phi_M(x_2) \\ \vdots & & \ddots & \vdots \\ \phi_0(x_N) & \phi_1(x_N) & \cdots & \phi_M(x_N) \end{pmatrix}$$

- To find $w^*$ (which minimize the error function), set the derivatives of the error function with respect to each $w_i$ equal to zero

$$\frac{d}{d(w_i)}\left[\frac{1}{2}(t - Qw)^T(t - Qw)]\right] = 0$$

# Linear Regression using basis function

- In gradient vector form it can be written as:

$$\mathbf{0} = \nabla_w[\frac{1}{2}(\mathbf{t} - \mathbf{Qw})^T(\mathbf{t} - \mathbf{Qw})]$$

$$\mathbf{0} = \nabla_w[\frac{1}{2}(\mathbf{t}^T\mathbf{t} - 2\mathbf{t}^T\mathbf{Qw} + \mathbf{w}^T\mathbf{Q}^T\mathbf{Qw})]$$

$$\mathbf{0} = [\frac{1}{2}(-2\mathbf{Q}^T\mathbf{t} + 2\mathbf{Q}^T\mathbf{Qw})]$$

- After simplification

$$\mathbf{Q}^T\mathbf{Qw} = \mathbf{Q}^T\mathbf{t}$$

$$\mathbf{w} = (\mathbf{Q}^T\mathbf{Q})^{-1}\mathbf{Q}^T\mathbf{t}$$

- No matter whether we do linear, polynomial or Gaussian, the only thing that changes is the definition of the matrix **Q**

# Regularized least squares

- $E(w) + \lambda E_w (w)$     where    $E_w (w) = \dfrac{1}{2} \boldsymbol{w}^T \boldsymbol{w}$

$$E(w) = \frac{1}{2}(\boldsymbol{t} - \boldsymbol{Qw})^T(\boldsymbol{t} - \boldsymbol{Qw}) + \frac{\lambda}{2}\boldsymbol{w}^T\boldsymbol{w}$$

$$\boldsymbol{0} = \nabla_{\boldsymbol{w}}[\frac{1}{2}(\boldsymbol{t} - \boldsymbol{Qw})^T(\boldsymbol{t} - \boldsymbol{Qw}) + \frac{\lambda}{2}\boldsymbol{w}^T\boldsymbol{w}]$$

$$\boldsymbol{0} = \nabla_{\boldsymbol{w}}[\frac{1}{2}(\boldsymbol{t}^T\boldsymbol{t} - 2\boldsymbol{t}^T\boldsymbol{Qw} + \boldsymbol{w}^T\boldsymbol{Q}^T\boldsymbol{Qw}) + \frac{\lambda}{2}\boldsymbol{w}^T\boldsymbol{w}]$$

$$\boldsymbol{0} = [-\boldsymbol{Q}^T\boldsymbol{t} + \boldsymbol{Q}^T\boldsymbol{Qw} + \lambda\boldsymbol{w}]$$

$$\boldsymbol{Q}^T\boldsymbol{Qw} + \lambda\boldsymbol{w} = \boldsymbol{Q}^T\boldsymbol{t}$$

$$(\boldsymbol{Q}^T\boldsymbol{Q} + \lambda\boldsymbol{I})\boldsymbol{w} = \boldsymbol{Q}^T\boldsymbol{t}$$

$$\textcolor{cyan}{\boldsymbol{w} = (\lambda\boldsymbol{I} + \boldsymbol{Q}^T\boldsymbol{Q})^{-1}\boldsymbol{Q}^T\boldsymbol{t}}$$

| likelihood | prior/posterior | data space |

# Bayesian Linear Regression

- $D = \{(x_1,y_1), (x_2,y_2),\ldots,(x_N,y_N)\}$     where   $x_i \in \mathbb{R}^D$ and $y \in \mathbb{R}$
- Model : $Y_1, Y_2, \ldots, Y_N$ independent given $\mathbf{w}$,   $Y \sim \mathcal{N}(w^T x_i, \beta)$ [$\beta$ is precision; $\beta = 1/\sigma^2$]
- $\mathbf{w} \sim \mathcal{N}(0, \alpha^{-1}I)$ where $\mathbf{w} = (w_0, w_1, \ldots, w_D)^T$

- Assume $\beta$ and $\alpha$ are known
  - ✓ therefore only unknown parameter is $\mathbf{w}$
- **Likelihood:**

$$p(\boldsymbol{D}|\boldsymbol{w})\, \alpha\; exp(-\frac{\beta}{2}(\boldsymbol{y} - \boldsymbol{Q}\boldsymbol{w})^T(\boldsymbol{y} - \boldsymbol{Q}\boldsymbol{w}))$$

- **Posterior:**

$$p(\boldsymbol{w}|\boldsymbol{D})\, \alpha\; p(\boldsymbol{D}|\boldsymbol{w})\, \boldsymbol{p}(\boldsymbol{w})$$

# Posterior of $w$

$$p(\boldsymbol{w}|\boldsymbol{D}) \; \alpha \; p(\boldsymbol{D}|\boldsymbol{w}) \; \boldsymbol{p}(\boldsymbol{w})$$

$$p(\boldsymbol{w}|\boldsymbol{D}) \; \alpha \; exp\left(-\frac{\beta}{2}(\boldsymbol{t}-\boldsymbol{Q}\boldsymbol{w})^T(\boldsymbol{t}-\boldsymbol{Q}\boldsymbol{w})\right) \boldsymbol{exp}(-\frac{\alpha}{2}\boldsymbol{w}^T\boldsymbol{w})$$

$$p(\boldsymbol{w}|\boldsymbol{D}) \; \alpha \; exp\left(-\frac{\beta}{2}(\boldsymbol{t}-\boldsymbol{Q}\boldsymbol{w})^T(\boldsymbol{t}-\boldsymbol{Q}\boldsymbol{w}) - \frac{\alpha}{2}\boldsymbol{w}^T\boldsymbol{w}\right)$$

$$p(\boldsymbol{w}|\boldsymbol{D}) \; \alpha \; exp\left(-\frac{1}{2}\left(\beta(\boldsymbol{t}-\boldsymbol{Q}\boldsymbol{w})^T(\boldsymbol{t}-\boldsymbol{Q}\boldsymbol{w}) + \alpha\boldsymbol{w}^T\boldsymbol{w}\right)\right)$$

$$p(\boldsymbol{w}|\boldsymbol{D}) \; \alpha \; exp\left(-\frac{1}{2}\left(\beta(\boldsymbol{t}^T\boldsymbol{t} - 2\boldsymbol{w}^T\boldsymbol{Q}^T\boldsymbol{t} + \boldsymbol{w}^T\boldsymbol{Q}^T\boldsymbol{Q}\boldsymbol{w})\right) + \alpha\boldsymbol{w}^T\boldsymbol{w}\right) \; [-2\boldsymbol{w}^T\boldsymbol{Q}^T\boldsymbol{t} = -\boldsymbol{t}^T\boldsymbol{Q}\boldsymbol{w} - (\boldsymbol{Q}\boldsymbol{w})^T\boldsymbol{t}\,]$$

$$p(\boldsymbol{w}|\boldsymbol{D}) \; \alpha \; exp\left(-\frac{1}{2}\left((\beta\boldsymbol{t}^T\boldsymbol{t} - 2\beta\boldsymbol{w}^T\boldsymbol{Q}^T\boldsymbol{t} + \beta\boldsymbol{w}^T\boldsymbol{Q}^T\boldsymbol{Q}\boldsymbol{w})\right) + \alpha\boldsymbol{w}^T\boldsymbol{w}\right)$$

$$p(\boldsymbol{w}|\boldsymbol{D}) \; \alpha \; exp\left(-\frac{1}{2}\left((\beta\boldsymbol{t}^T\boldsymbol{t} - 2\beta\boldsymbol{w}^T\boldsymbol{Q}^T\boldsymbol{t} + \boldsymbol{w}^T(\beta\boldsymbol{Q}^T\boldsymbol{Q} + \alpha\boldsymbol{I})\boldsymbol{w}\right)\right)$$

## Posterior of $w$

$$p(\boldsymbol{w}|\boldsymbol{D}) \; \alpha \; exp\left(-\frac{1}{2}\left(\textcolor{green}{\beta t^T t} - \textcolor{blue}{2}\beta \textcolor{gold}{w^T Q^T t} + \textcolor{blue}{w^T}(\textcolor{red}{\beta Q^T Q} + \alpha I)w\right)\right)$$

Completing the square:

$$\mathscr{N}(\boldsymbol{\mu}, \Lambda^{-1}) \; \alpha \; \boldsymbol{exp}(-\frac{1}{2}(\boldsymbol{w} - \boldsymbol{\mu})^T \Lambda (\boldsymbol{w} - \boldsymbol{\mu}))$$

$$\mathscr{N}(\boldsymbol{\mu}, \Lambda^{-1}) \; \alpha \; \boldsymbol{exp}(-\frac{1}{2}(\boldsymbol{w^T} \Lambda \boldsymbol{w} - \boldsymbol{w^T} \Lambda \boldsymbol{\mu} - \boldsymbol{\mu^T} \Lambda \boldsymbol{w} + \boldsymbol{\mu^T} \Lambda \boldsymbol{\mu}))$$

$$\mathscr{N}(\boldsymbol{\mu}, \Lambda^{-1}) \; \alpha \; \boldsymbol{exp}(-\frac{1}{2}(\textcolor{blue}{w^T \Lambda w} - \textcolor{blue}{2}\textcolor{red}{w^T} \Lambda \textcolor{gold}{\mu} + \textcolor{green}{\mu^T \Lambda \mu}))$$

Hence:

$$\textcolor{red}{\Lambda = \beta Q^T Q + \alpha I}$$

$$\textcolor{gold}{\mu = \beta \Lambda^{-1} Q^T t} \qquad (\textcolor{gold}{\cancel{w^T} \Lambda \mu} = \beta \textcolor{gold}{\cancel{w^T} Q^T t})$$

# Posterior of $w$

$$\Lambda = \beta Q^T Q + \alpha I$$

$$\mu = \beta \Lambda^{-1} Q^T t$$

$$\mu = \beta(\beta Q^T Q + \alpha I)^{-1} Q^T t = (Q^T Q + \frac{\alpha}{\beta} I)^{-1} Q^T t$$

- $p(w|D) \sim \mathcal{N}(\mu, \Lambda^{-1})$
- MAP estimate of $w$

$$w_{MAP} = (Q^T Q + \frac{\alpha}{\beta} I)^{-1} Q^T t \qquad \text{[From regularized least square slide}$$

$$w = (\lambda I + Q^T Q)^{-1} Q^T t]$$

# Prediction

$$p(y|\boldsymbol{x}, \boldsymbol{D}) = \int p(y|\boldsymbol{x}, \boldsymbol{D}, \boldsymbol{w})p(\boldsymbol{w}|\boldsymbol{x}, \boldsymbol{D})d\boldsymbol{w}$$

$$p(y|\boldsymbol{x}, \boldsymbol{D}) = \int p(y|\boldsymbol{x}, \boldsymbol{w})p(\boldsymbol{w}|\boldsymbol{D})d\boldsymbol{w} \quad \text{[Note: } p(\boldsymbol{w}|\boldsymbol{D}) \, \alpha \, exp\left(-\frac{1}{2}\left((\beta\boldsymbol{y^T}\boldsymbol{y} - 2\beta\boldsymbol{w^T}\boldsymbol{Q^T}\boldsymbol{y} + \boldsymbol{w^T}(\beta\boldsymbol{Q^T}\boldsymbol{Q} + \alpha\boldsymbol{I})\boldsymbol{w}\right)\text{]}$$

$$p(y|\boldsymbol{x}, \boldsymbol{D}) \, \alpha \int \exp\left(-\frac{\beta}{2}(y - \boldsymbol{w^T}\boldsymbol{x})^2\right)\exp\left(-\frac{1}{2}(\boldsymbol{w} - \boldsymbol{\mu})^T\Lambda(\boldsymbol{w} - \boldsymbol{\mu})\right)dw$$

$$p(y|\boldsymbol{x}, \boldsymbol{D}) \, \alpha \, \exp(-\frac{\gamma}{2}(y - \upsilon)^2)$$

where $\quad \upsilon = \boldsymbol{\mu}^{\text{T}}\boldsymbol{x}$

$$\frac{1}{\gamma} = \frac{1}{\beta} + \boldsymbol{x}^T \Lambda^{-1}\boldsymbol{x}$$

$$\boldsymbol{\mu} = \beta\Lambda^{-1}\boldsymbol{Q^T}\boldsymbol{y}$$
$$\Lambda = \beta\boldsymbol{Q^T}\boldsymbol{Q} + \alpha\boldsymbol{I}$$

# References

- Chapter 3, Pattern Recognition and Machine Learning, C. Bishop