

# Machine Learning Basics: Supervised Learning Algorithms

Sargur N. Srihari

srihari@cedar.buffalo.edu

This is part of lecture slides on [Deep Learning](http://www.cedar.buffalo.edu/~srihari/CSE676):  
<http://www.cedar.buffalo.edu/~srihari/CSE676>

# Topics

1. Learning Algorithms
2. Capacity, Overfitting and Underfitting
3. Hyperparameters and Validation Sets
4. Estimators, Bias and Variance
5. Maximum Likelihood Estimation
6. Bayesian Statistics
7. Supervised Learning Algorithms
8. Unsupervised Learning Algorithms
9. Stochastic Gradient Descent
10. Building a Machine Learning Algorithm
11. Challenges Motivating Deep Learning

# Topics in Supervised Learning

- Definition of supervised learning
  1. Probabilistic Supervised Learning
  2. Support Vector Machines
  3. Other simple supervised learning algorithms
    - $k$  - nearest neighbor
    - Decision tree

# What is Supervised Learning?

- Refers to learning algorithms that learn to associate some input with some output given a training set of inputs  $x$  and outputs  $y$
- Outputs may be collected automatically or provided by a human supervisor

# Probabilistic Supervised Learning

- Most supervised learning algorithms are based on estimating a probability distribution  $p(y|\mathbf{x})$
- We can do this by using MLE to find the best parameter vector  $\boldsymbol{\theta}$  for a parametric family of distributions  $p(y|\mathbf{x};\boldsymbol{\theta})$
- Linear regression corresponds to the family
$$p(y|\mathbf{x};\boldsymbol{\theta}) = N(y|\boldsymbol{\theta}^T \mathbf{x}, \mathbf{I})$$
- We can generalize linear regression to classification by using a different family of probability distributions

# Probabilistic Supervised Classification

- If we only have two classes we only need to specify the distribution for one of these classes
  - The probability of the other class is known

For *linear regression*, normal distribution

$$p(y|\mathbf{x};\theta) = N(y|\theta^T \mathbf{x}, I)$$

was parameterized by its mean  $\theta^T \mathbf{x}$

Any value we supply for this mean is valid

For *classification*, a distribution over a binary variable is more complex.

Its mean must be between 0 and 1.

Solved using logistic sigmoid to squash the output of a linear function into interval (0,1):

$$p(y=1|\mathbf{x};\theta) = \sigma(\theta^T \mathbf{x})$$

Known as *logistic regression*

- Linear regression has a closed-form solution
- But Logistic regression has no closed-form solution
  - Negative log-likelihood is minimized with gradient descent<sup>6</sup>

# Strategy for Supervised Learning

- Same strategy can be applied to any supervised learning problem
  - Write down a parametric family of conditional probability distributions over the right kind of input and output variables

# Support Vector Machines

- An influential approach to supervised learning
- Model is similar to logistic regression in that it is driven by a linear function  $w^T x + b$ 
  - Unlike logistic regression, SVM does not provide probabilities, but only outputs class identity
    - SVM predicts positive class when  $w^T x + b > 0$
    - SVM predicts negative class when  $w^T x + b < 0$



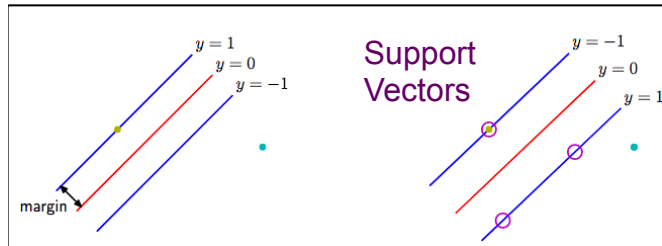
# Kernel trick

- A key innovation associated with SVM
- Many ML algorithms can be written as dot products between examples
  - Ex: linear function used by SVM can be rewritten as

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = b + \sum_{i=1}^m \alpha_i \mathbf{x}^T \mathbf{x}^{(i)}$$

- where  $\mathbf{x}^{(i)}$  is a training example,  $\boldsymbol{\alpha}$  a vector of coefficients
- Replace  $\mathbf{x}$  by a feature function  $\phi(\mathbf{x})$  and the dot product with a kernel function  $k(\mathbf{x}, \mathbf{x}^{(i)}) = \phi(\mathbf{x}) \bullet \phi(\mathbf{x}^{(i)})$ 
  - The  $\bullet$  operator represents inner product  $\phi(\mathbf{x})^T \phi(\mathbf{x}^{(i)})$
  - We may not literally use the inner product

# SVM Optimization Problem



Equation of plane:

$$g(x) = w^t \cdot x + w_0$$

Distance from  $x$  to plane:

$$r = g(x) / \|w\|$$

For each input vector, let  $z_k = \pm 1$  depending on whether input  $k$  is in class  $C_1$  or  $C_2$

Thus if  $g(y)=0$  is a separating hyperplane then  $z_k g(y_k) \geq 0, k = 1, \dots, n$   
 Since distance of a point  $y$  to hyperplane  $g(y)=0$  is  $g(y) / \|a\|$  we could require that hyperplane be such that all points are at least distant  $b$  from it, i.e.

The goal is to find the weight vector  $a$  that satisfies

$$\frac{z_k g(y_k)}{\|a\|} \geq b$$

$$\frac{z_k g(y_k)}{\|a\|} \geq b, \quad k = 1, \dots, n$$

while maximizing  $b$ .

To ensure uniqueness we impose the constraint  $b \|a\| = 1$  or  $b = 1 / \|a\|$  which implies that  $\|a\|^2$  is to be minimized

Support vectors are training vectors which represent equality. A quadratic optimization problem: minimize a quadratic subject to linear inequality constraints.

optimize

$$\arg \min_{a,b} \frac{1}{2} \|a\|^2$$

subject to constraints

$$z_k a^t y_k \geq 1, \quad k = 1, \dots, n$$

Can be cast as unconstrained problem by introducing Lagrange multipliers with one multiplier  $\alpha_k$  for each constraint. The Lagrange function is

$$L(a, \alpha) = \frac{1}{2} \|a\|^2 - \sum_{k=1}^n \alpha_k [z_k a^t y_k - 1]$$

Dual problem:

$$L(\alpha) = \sum_{k=1}^n \alpha_k - \frac{1}{2} \sum_{k=1}^n \sum_{j=1}^n \alpha_k \alpha_j z_k z_j k(y_j, y_k)$$

subject to constraint:

$$\sum_{k=1}^n \alpha_k z_k = 0 \quad \alpha_k \geq 0, \quad k = 1, \dots, n$$

where kernel is defined as:

$$k(y_j, y_k) = y_j^t \cdot y_k = \phi(x_j)^t \cdot \phi(x_k)$$

# Prediction using SVM

- After replacing dot products with kernel evaluations, we can make predictions using

$$f(\mathbf{x}) = b + \sum_i \alpha_i k(\mathbf{x}, \mathbf{x}^{(i)})$$

- Depending on whether  $f(\mathbf{x}) > 0$ 
  - The function is nonlinear in  $\mathbf{x}$  but the relationship between  $f(\mathbf{x})$  and  $\phi(\mathbf{x})$  is linear
    - Also the relationship between  $\alpha$  and  $f(\mathbf{x})$  is linear
- The kernel-based function is exactly equivalent to preprocessing the data by applying  $\phi(\mathbf{x})$  to all inputs, then learning a linear model in the new transformed space

# Efficacy & Efficiency of kernel

1. Kernel trick allows us to learn models that are nonlinear as a function of  $\mathbf{x}$  using convex optimization guaranteed to converge efficiently
  - Possible because  $\phi$  is fixed and optimize only  $\alpha$
2. Kernel  $k$  implemented more efficiently than constructing  $\phi$  vectors and taking dot product
  - $k(\mathbf{x}, \mathbf{x}')$  tractable even when  $\phi(\mathbf{x})$  is intractable

# Gaussian Kernel

- Most commonly used kernel is  $k(\mathbf{u}, \mathbf{v}) = N(\mathbf{u} - \mathbf{v}; 0, \sigma^2 \mathbf{I})$ 
  - Called radial basis: decreases along lines radiating from  $\mathbf{u}$
- To see that it is a valid kernel
  - Consider  $k(\mathbf{u}, \mathbf{v}) = \exp(-\|\mathbf{u} - \mathbf{v}\|^2 / 2\sigma^2)$ 
    - By expanding the square  $\|\mathbf{u} - \mathbf{v}\|^2 = \mathbf{u}^T \mathbf{u} + \mathbf{v}^T \mathbf{v} - 2\mathbf{u}^T \mathbf{v}$
    - we get  $k(\mathbf{u}, \mathbf{v}) = \exp(-\mathbf{u}^T \mathbf{u} / 2\sigma^2) \exp(-\mathbf{u}^T \mathbf{v} / \sigma^2) \exp(-\mathbf{v}^T \mathbf{v} / 2\sigma^2)$
  - Validity follows from kernel construction rules
    - If  $k_1(\mathbf{x}, \mathbf{x}')$  is a valid kernel, so are
    - $k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x}) k_1(\mathbf{x}, \mathbf{x}') f(\mathbf{x}')$  and  $k(\mathbf{x}, \mathbf{x}') = \exp(k_1(\mathbf{x}, \mathbf{x}'))$
    - together with validity of linear kernel  $k(\mathbf{u}, \mathbf{v}) = \mathbf{u}^T \mathbf{v}$

# Intuition of Gaussian kernel

- It performs a kind of template matching
- When a test point  $x'$  is near a template  $x$  its response is high, putting a large weight on the associated training label  $y$
- Overall, the prediction combines many such training labels weighted by the similarity of the corresponding training examples
- SVM is not the only algorithm enhanced by the kernel trick
  - Methods employing the kernel trick are kernel methods

# Disadvantages of Kernel Methods

- Cost of decision function evaluation: linear in  $m$ 
  - Because the  $i^{\text{th}}$  example contributes a term  $\alpha_i k(\mathbf{x}, \mathbf{x}^{(i)})$  to the decision function
  - Can mitigate this by learning an  $\alpha$  with mostly zeros
    - Classification requires evaluating the kernel function only for training examples that have a nonzero  $\alpha_i$
    - These are known as *support vectors*
- Cost of training: high with large data sets
- Generic kernels struggle to generalize well
  - Neural network outperformed RBF kernel SVM on MNIST benchmark

# Other simple supervised learning

- K-nearest neighbor
- Decision trees



# K-Nearest Neighbors

- A family of techniques that can be used for regression or classification
- As a nonparametric learning algorithm:
  - it is not restricted to a fixed no of parameters
- A simple function of the training data
  - When we want to produce output  $y$  for input  $x$ , we find  $k$ -nearest neighbors to  $x$  in the training data  $\mathbf{X}$ . We then return the average of the corresponding  $y$  values in the training set
    - Works for any kind of supervised learning where we can define average over  $y$  values

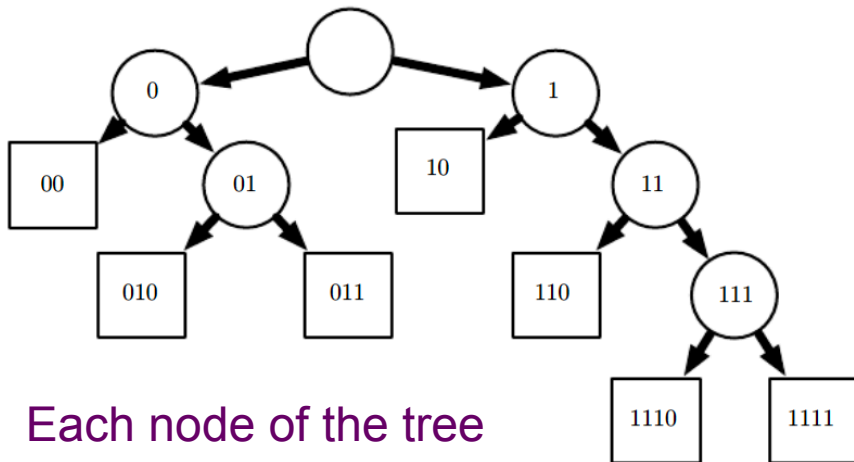
# K-Nearest Neighbor Classification

- We can average over one-hot-vectors  $\mathbf{c}$  with  $c_y=1$  and  $c_i=0$  for all other values of  $I$ 
  - Can interpret the average values as giving a probability distribution over classes
- K-nn classification has high capacity
  - High accuracy given a large training set
    - 1-nearest neighbor error rate approaches twice Bayes error rate, with training set size, assuming 0-1 loss
      - Error in excess of Bayes rate because of choosing a neighbor by breaking ties between equally distant neighbors randomly
      - If we chooses all neighbors to vote, we get Bayes rate
- A weakness: cannot learn one feature more discriminative than another

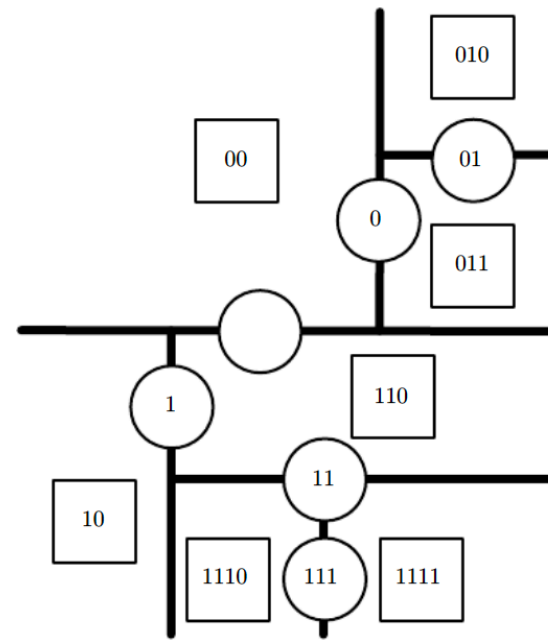
# Decision Tree

- A learning algorithm that breaks the input space into regions and has separate parameters for each region
- Each node is associated with with a region of input space
- Internal nodes break that region into one or more subregion for each child

# How a decision tree works



Each node of the tree chooses to send the input example to the child node on the left (0) or the child node on the right (1). Internal nodes are circles. Leaf nodes are squares.



The tree divides space into regions. This 2-D plane shows how a decision tree might divide  $R^2$ .

The nodes are plotted in this plane with each internal node drawn along with the dividing line it uses to categorize examples and leaf nodes drawn in the center of the region of examples they receive. The result is a piecewise-constant function, with one piece per leaf.