



Instructor: Sael Lee

CS549 Spring – Computational Biology

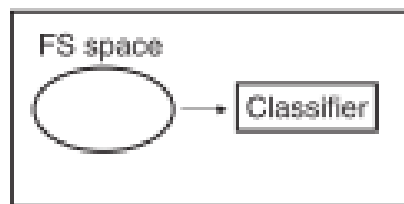
# LECTURE 12-13: FEATURE SELECTION

Ref.

1. C. M. Bishop “Pattern Recognition and Machine Learning” 2<sup>nd</sup> ed. & provided slides

# TYPES OF FEATURE SELECTION METHOD

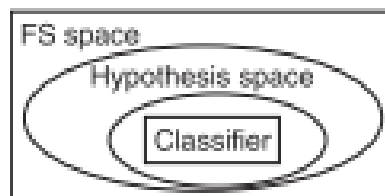
## Filtering Methods



relevance of features is evaluated by looking only at the intrinsic properties of the data

\* Often **feature relevance score** is used to evaluate each feature (gene)

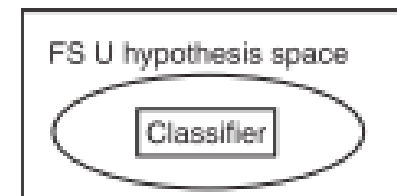
## Wrapper Methods



model hypothesis search is embed within the feature subset search

-> **various subsets of features** are generated and evaluated

## Embedded Method



optimal feature subset search is built into the classifier construction

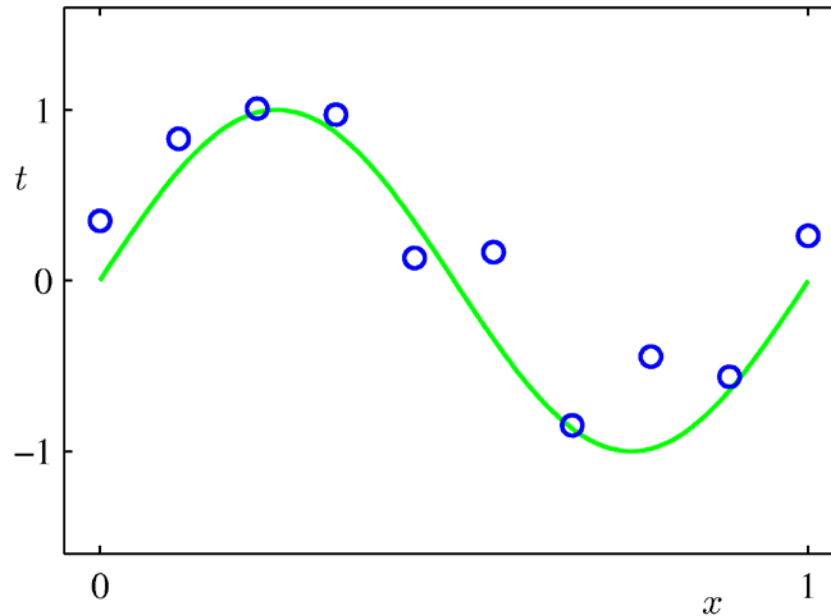
-> a search in the **combined space of feature subsets and hypotheses**

Chapter 3 of PRML

# FEATURE SELECTION WITH LASSO REGRESSION MODEL

# LINEAR BASIS FUNCTION MODELS (1)

## × Example: Polynomial Curve Fitting



$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{j=0}^M w_jx^j$$

## LINEAR BASIS FUNCTION MODELS (2)

---

× Generally

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

- × where  $\phi_j(\mathbf{x})$  are known as *basis functions*.
- × Typically,  $\phi_0(\mathbf{x}) = 1$ , so that  $w_0$  acts as a **bias**.
- × In the simplest case, we use linear basis functions :  
 $\phi_d(x) = x_d$ .

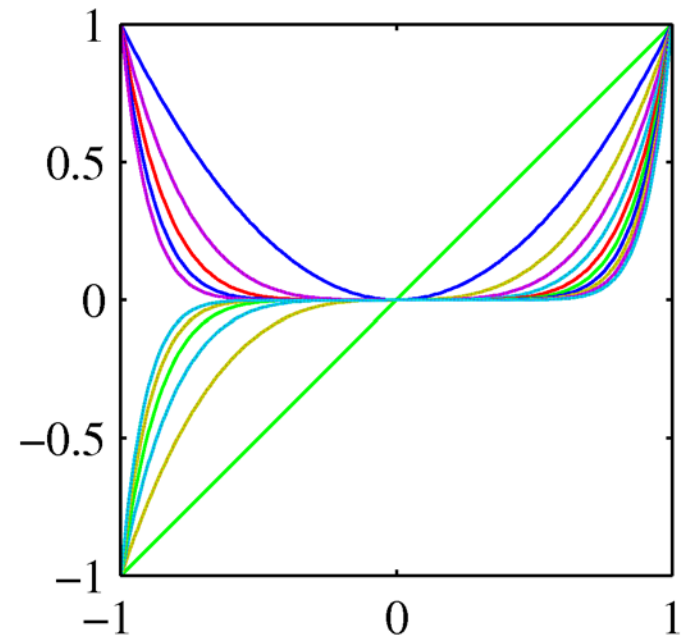
# LINEAR BASIS FUNCTION MODELS (3)

× Polynomial basis function

S:

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

× These are global; a small change in  $x$  affect all basis functions.



# LINEAR BASIS FUNCTION MODELS (4)

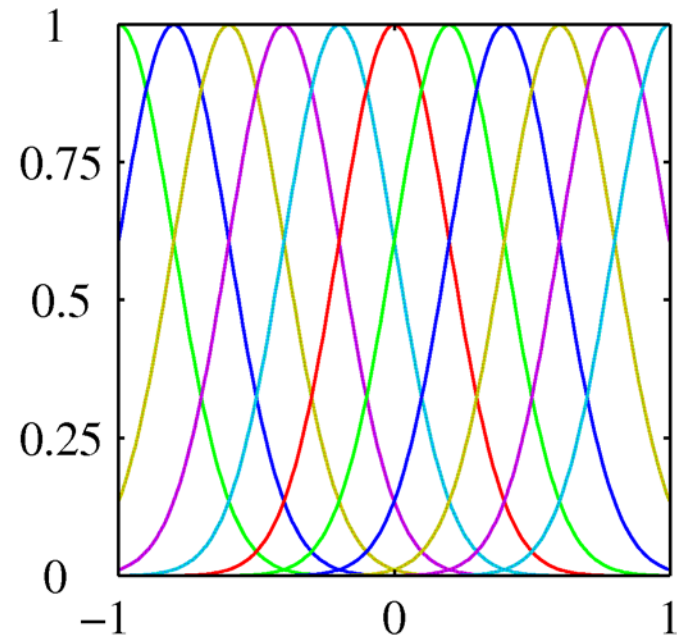
✖ Gaussian basis functions:

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

✖ These are local;

+ a small change in  $x$  only affect nearby basis functions.

+  $\mu_j$  and  $s$  control location and scale (width).



# LINEAR BASIS FUNCTION MODELS (5)

× Sigmoidal basis functions:

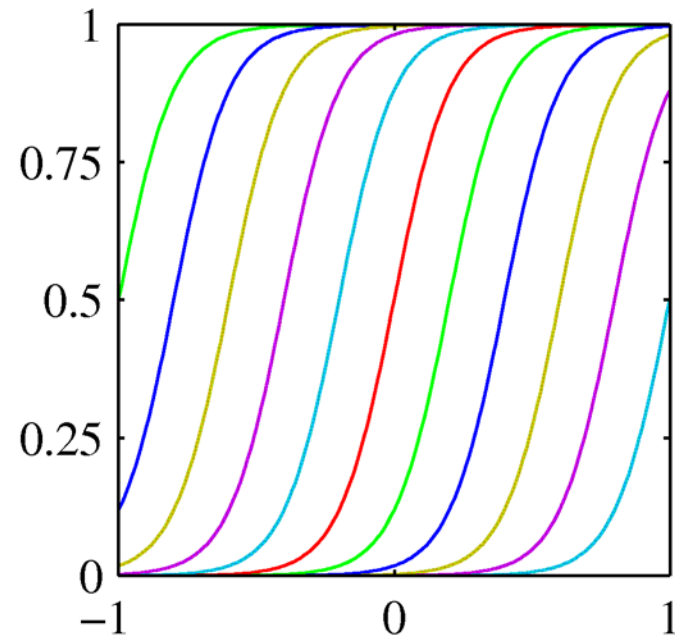
where 
$$\phi_j(x) = \sigma\left(\frac{x - \mu_j}{s}\right)$$

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

× Also these are local;

+ a small change in  $x$  only affect nearby basis functions.

+  $\mu_j$  and  $s$  control location and scale (width).





# MAXIMUM LIKELIHOOD AND LEAST SQUARES (1)

- × Assume observations from a deterministic function with added Gaussian noise:

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon \quad \text{where} \quad p(\epsilon|\beta) = \mathcal{N}(\epsilon|0, \beta^{-1})$$

- × which is the same as saying,

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

- × Given observed inputs,  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  and targets,  $\mathbf{t} = [t_1, \dots, t_N]^T$ , we obtain the likelihood function

likelihood function

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

# MAXIMUM LIKELIHOOD AND LEAST SQUARES (2)

× Log likelihood:

$$\begin{aligned}\ln p(\mathbf{t}|\mathbf{w}, \beta) &= \sum_{n=1}^N \ln \mathcal{N}(t_n | \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1}) \\ &= \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta E_D(\mathbf{w})\end{aligned}$$

where

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

is the **sum-of-squares error**.

$$\begin{aligned}N(t_n | \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1}) &= \\ & \left(\frac{\beta}{2\pi}\right)^{\frac{1}{2}} \exp\left(-\frac{\beta}{2} (t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n))^2\right)\end{aligned}$$

Relationship of log likelihood and sum-of-squares error in univariate Gaussian noise model.

# MAXIMUM LIKELIHOOD AND LEAST SQUARES (3)

- × Computing the gradient and setting it to zero yields

$$\nabla_{\mathbf{w}} \ln p(\mathbf{t}|\mathbf{w}, \beta) = \beta \sum_{n=1}^N \{t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n)\} \boldsymbol{\phi}(\mathbf{x}_n)^T = \mathbf{0}.$$

- × Solving for  $w$ , we get

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

The Moore-Penrose pseudo-inverse,  $\boldsymbol{\Phi}^\dagger$ .

- × where

Design matrix

$$\boldsymbol{\Phi} = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \cdots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \cdots & \phi_{M-1}(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \cdots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}.$$

# MAXIMUM LIKELIHOOD AND LEAST SQUARES (4)

- × Maximizing with respect to the bias,  $w_0$ , alone, we see that

$$\begin{aligned}w_0 &= \bar{t} - \sum_{j=1}^{M-1} w_j \bar{\phi}_j \\ &= \frac{1}{N} \sum_{n=1}^N t_n - \sum_{j=1}^{M-1} w_j \frac{1}{N} \sum_{n=1}^N \phi_j(\mathbf{x}_n).\end{aligned}$$

- × We can also maximize with respect to  $\beta$ , giving

$$\frac{1}{\beta_{\text{ML}}} = \frac{1}{N} \sum_{n=1}^N \underbrace{\{t_n - \mathbf{w}_{\text{ML}}^T \boldsymbol{\phi}(\mathbf{x}_n)\}^2}$$

Residual variance of the target value  
around the regression function

# REGULARIZED LEAST SQUARES (1)

- × Consider the error function:

$$E_D(\mathbf{w}) + \lambda E_W(\mathbf{w})$$

Data term + Regularization term

- × With the **sum-of-squares error (SSE)** function and a quadratic regularizer, we get

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

- × which is minimized by

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

$\lambda$  is called the regularization coefficient.

## REGULARIZED LEAST SQUARES (2)

- × With a more general regularizer, we have

$$\frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \sum_{j=1}^M |w_j|^q$$

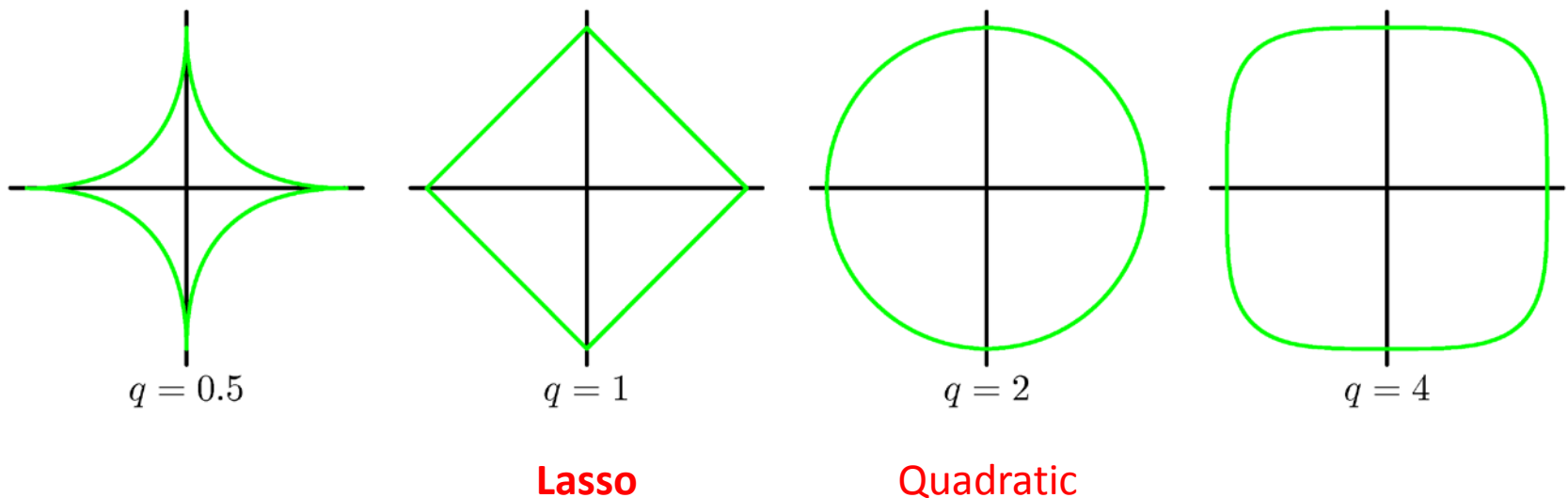


Fig: Contours of the regularization terms

# USING LASSO FOR FEATURE SELECTION

## Lasso tends to generate sparser solutions

- + If  $\lambda$  is sufficiently large, some of the coefficients  $w_j$  are driven to zero, leading to a sparse model in which the corresponding basis function plays no role.

## Minimizing

general regularizer

$$\frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \sum_{j=1}^M |w_j|^q$$

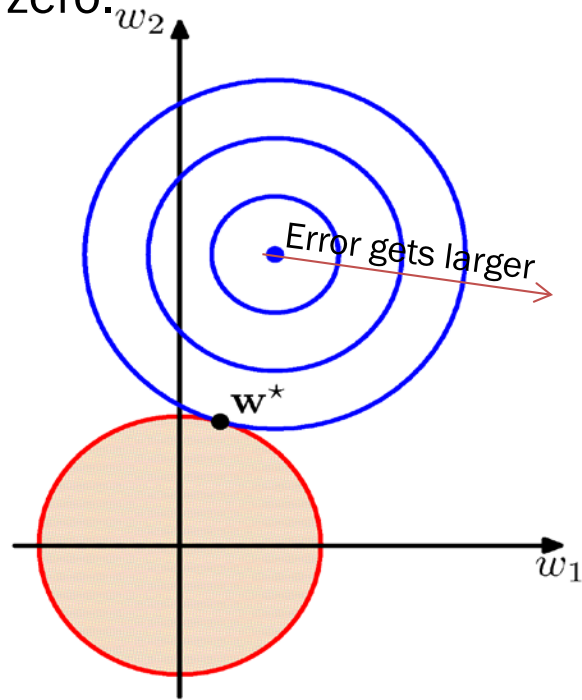
is equivalent to minimizing the unregularized SSE subjected to constraint

Lagrangian Multiplier

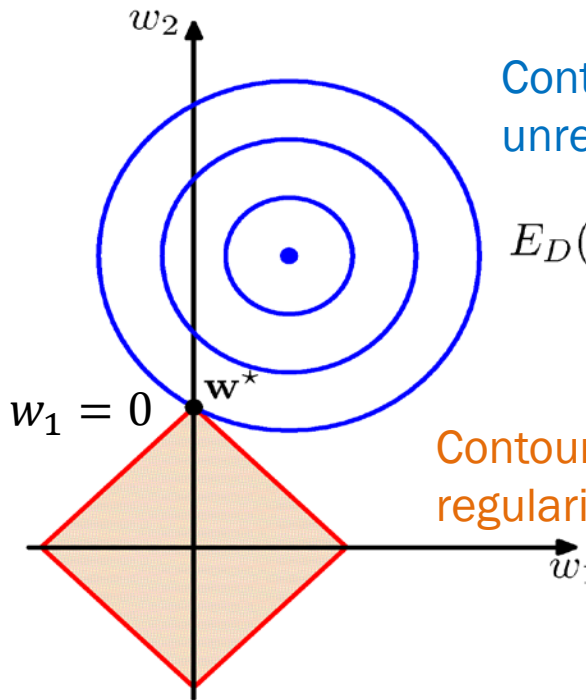
$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 \quad \text{Subjected to} \quad \sum_{j=1}^M |w_j|^q \leq \eta$$

# REGULARIZED LEAST SQUARES (3)

Figure shows the minimum of the error function, subjected to constraint. As  $\lambda$  is increased, so an increasing number of parameters are driven to zero.



Quadratic



Lasso

Contours of unregularized SSE

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

Contours of the regularization terms

$$\sum_{j=1}^M |w_j|^q \leq \eta$$

Lasso give sparse solution in which  $w^* = 0$ .

Q: So, how do we find the right  $\lambda$ ?



# LIMITATIONS OF LASSO BASED FEATURE SELECTION

- × Linear feature space : inadequate to capture non-linear dependencies from features to output

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

# SUPPORT VECTOR MACHINES

---

# REVIEW

---

linear regression

$$h_{\mathbf{w}}(\mathbf{x}) = (\mathbf{w}^T \mathbf{x})$$

Logistic regression

Estimated prob. that  $y=1$   
on input  $\mathbf{x}$

$$h_{\mathbf{w}}(\mathbf{x}) = \text{Logistic}(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x})}}$$

# LEARNING THE WEIGHTS

$$w^* = \operatorname{argmin}_w \operatorname{Loss}(h_w)$$

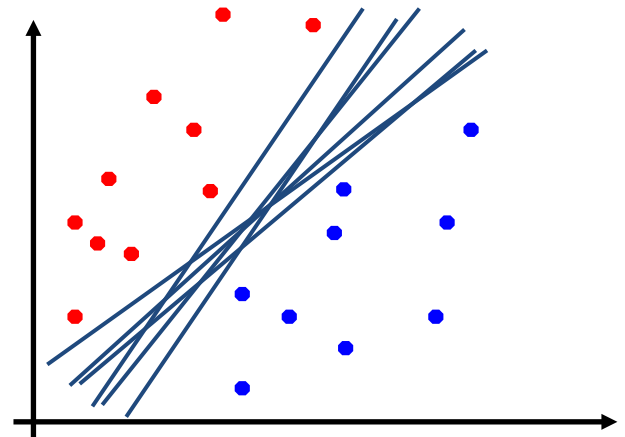
linear regression:

$$\operatorname{Loss}(h_w) = \sum_{j=1}^N (y_j - (w^T x_j))^2$$

logistic regression:

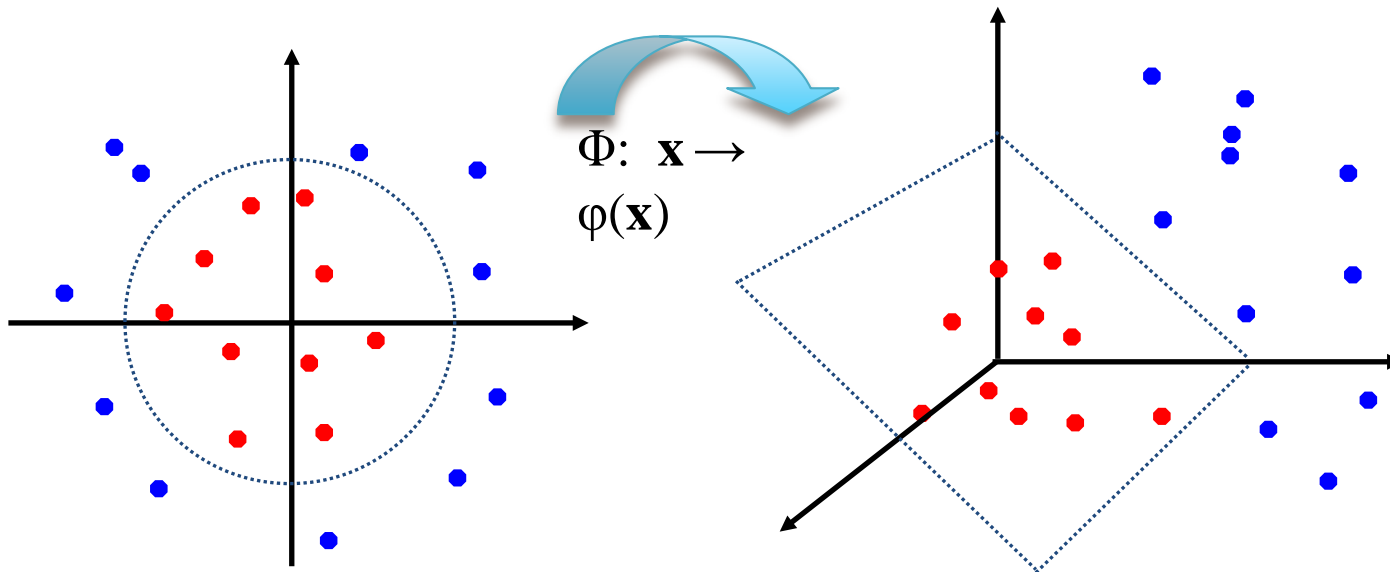
$$\operatorname{Loss}(h_w) = \sum_{j=1}^N -y_j(\log(h_w(x_j))) - (y_j - 1)(\log(1 - h_w(x_j)))$$

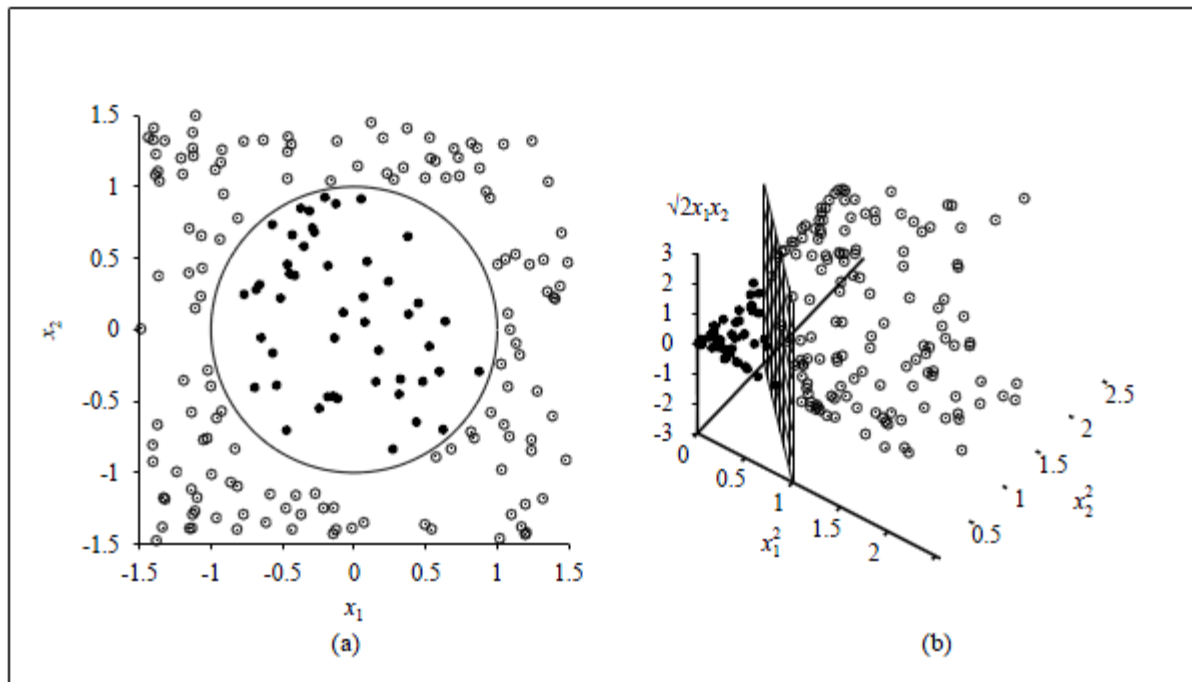
- $y$  is classification label in logistics regression (0 or 1)
- $y$  is scalar values in linear regression



# KERNELS

- The original feature space can always be mapped to some higher-dimensional feature space (even infinite) where the training set is separable





**Figure 18.31** FILES: . (a) A two-dimensional training set with positive examples as black circles and negative examples as white circles. The true decision boundary,  $x_1^2 + x_2^2 \leq 1$ , is also shown. (b) The same data after mapping into a three-dimensional input space  $(x_1^2, x_2^2, \sqrt{2}x_1x_2)$ . The circular decision boundary in (a) becomes a linear decision boundary in three dimensions. Figure 18.29(b) gives a closeup of the separator in (b).

# KERNELS

- The linear classifier relies on an inner product between vectors  $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$
- If every data point is mapped into high-dimensional space via some transformation  $\Phi: \mathbf{x} \rightarrow \varphi(\mathbf{x})$ , the inner product becomes:

$$K(\mathbf{x}_i, \mathbf{x}_j) = \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j)$$

- A *kernel function* is some function that corresponds to an inner product in some expanded feature space.
  
- Kernel function should measure some **similarity** between data
- kernel must be **positive semi-definite**
  
- You should **scale the features** to have same scale!!
  
- Most widely used is **linear kernels** and **Gaussian kernels**

# GAUSSIAN KERNELS

$$k(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right) = \exp\left(-\frac{\sum_{k=1}^n (x_{ik} - x_{jk})^2}{2\sigma^2}\right)$$

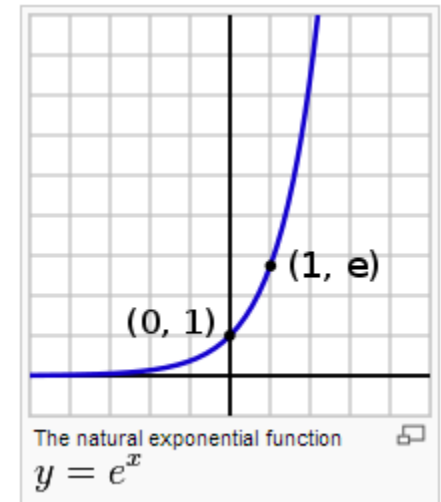
If  $x_i$  and  $x_j$  is similar:

$$k(x_i, x_j) \approx \exp\left(-\frac{0^2}{2\sigma^2}\right) \approx 1$$

If  $x_i$  and  $x_j$  is different:

$$k(x_i, x_j) \approx \exp\left(-\frac{(\text{large number})^2}{2\sigma^2}\right) \approx 0$$

If you use Gaussian kernel,  
You will need to pick  $\sigma$





# SUPPORT VECTOR MACHINES

---

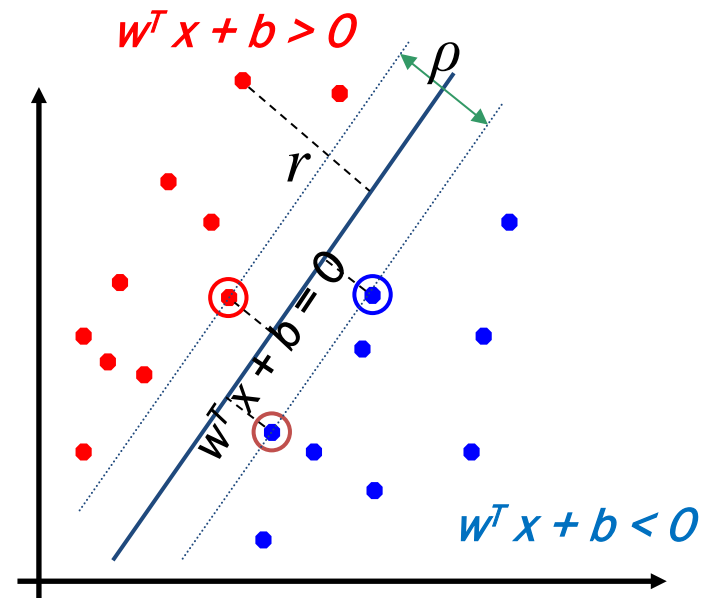
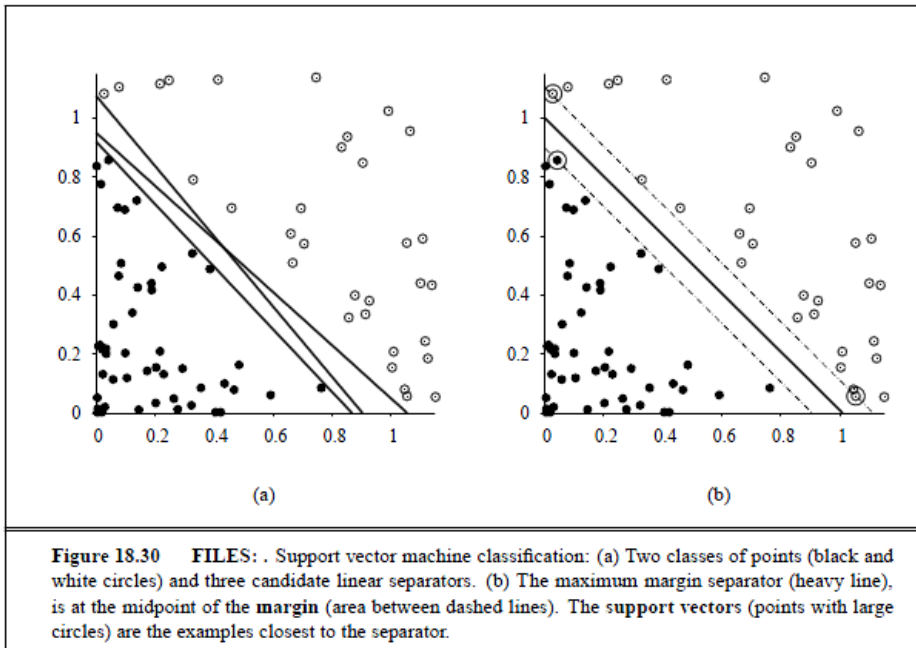
- × SVMs constructs a **maximum margin separator**
- × SVMs create a linear separating hyperplane
  - + But have ability to embed that in to higher-dimensional space (via **Kernel trick**)
- × SVM are a nonparametric method
  - + Retain training examples and potentially need to store all or part of the data
  - + Some example are more important than others (**support vectors**)

# SVM TERMS

- Distance from example  $x_i$  to the separator is

$$r = \frac{(w^T x + b)}{\|w\|}$$

- Examples closest to the hyperplane are *support vectors*.
- Margin*  $\rho$  of the separator is the distance between support vectors



# MARGINS

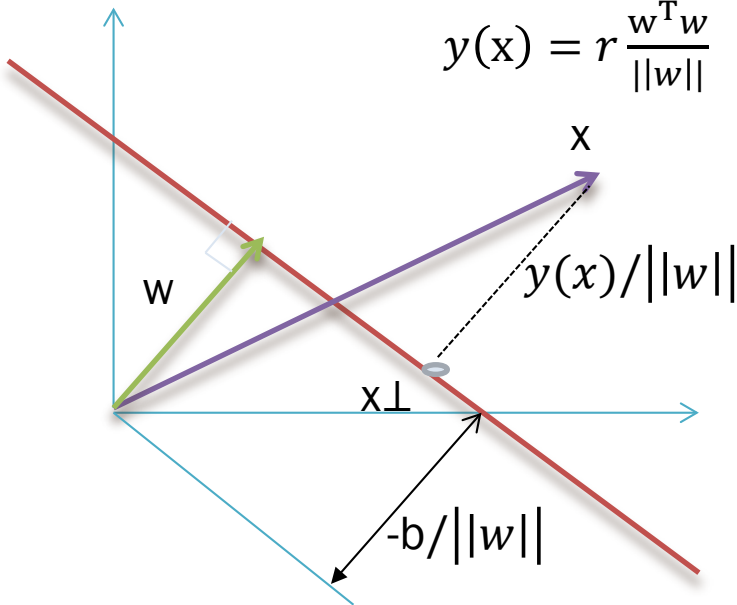
Instead of minimizing expected **empirical loss** in the training data, SVM attempts to minimize expected **generalization loss**.

$$y(x) = w^T x + b \text{ where } w \text{ is weight vector and } b \text{ is bias}$$
$$x = x_{\perp} + r \frac{w}{\|w\|} \quad (\text{multiply } w^T \text{ and add } b)$$

$$w^T x + b = w^T \left( x_{\perp} + r \frac{w}{\|w\|} \right) + b \quad (y(x) = w^T x + b)$$

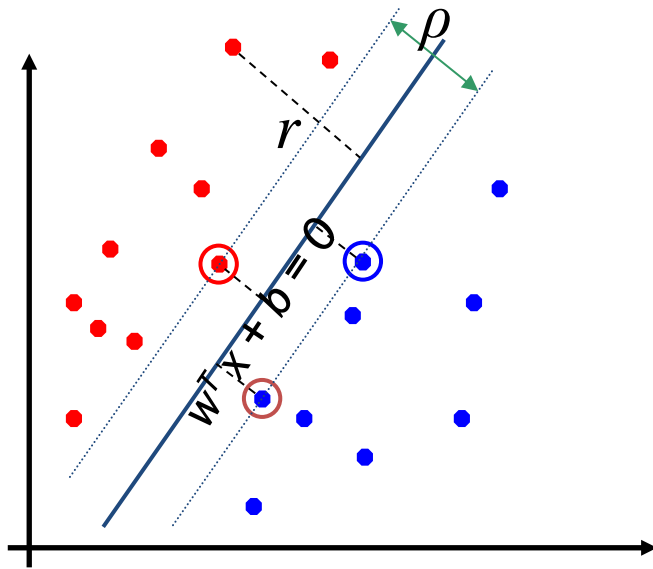
$$y(x) = w^T x_{\perp} + r \frac{w^T w}{\|w\|} + b \quad (y(x_{\perp}) = w^T x_{\perp} + b = 0)$$

$$y(x) = r \frac{w^T w}{\|w\|} = |$$



$$\text{or } r = \frac{(w^T x + b)}{\|w\|}$$

# MAXIMUM MARGINS



Solving this is non-trivial and will not be discussed in class

$\phi(x_n)$  in the feature space

$$r = \frac{(w^T x + b)}{\|w\|}$$

$$\operatorname{argmax}_{w,b} \left\{ \frac{1}{\|w\|} \min_n [t_n (w^T x_n + b)] \right\}$$



$$\operatorname{argmin}_{w,b} \frac{1}{2} \|w\|^2$$



$$w = \sum_{n=1}^N a_n t_n \phi(x_n)$$

$$\sum_{n=1}^N a_n t_n = 0$$

# SOFT MARGINS

**Idea:** Allow data point to be in the wrong side of the margin boundary, but with a penalty that increases with the distance from that boundary.

**Penalty** for each data point : **slack variable  $\xi$**

$\xi_n = 0$  if point is on the right side

$\xi_n = |t_n - y(x_n)|$  if point is on the wrong side

Such that

$t_n y(x_n) \geq 1 - \xi_n$  for  $n = 1, \dots, N$  and  $\xi_n \geq 0$

- $0 < \xi_n \leq 1$  for points inside the margin
- $\xi_n = 1$  for points on the margin
- $\xi_n > 1$  for points that are on the wrong side

**Goal** now is to maximize the margin while softly penalizing points that lie on the wrong side of the margin boundary

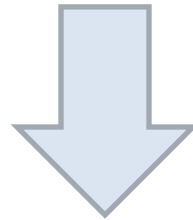
$$\operatorname{argmin}_{w,b} C \sum_n^N \xi_n + \frac{1}{2} \|w\|^2$$

# OPTIMIZATION ON SOFT MARGINS

$$\operatorname{argmin}_{w,b} C \sum_n^N \xi_n + \frac{1}{2} \|w\|^2$$

subjected to  $t_n y(x_n) \geq 1 - \xi_n$  for  $n = 1, \dots, N$  and  $\xi_n \geq 0$

$\xi_n$ : slack variable for training data  $x_n$



Complex calculations  
Lagrangian  
Etc.

$$w = \sum_{n=1}^N a_n t_n \phi(x_n)$$

$$\sum_{n=1}^N a_n t_n = 0$$

$$a_n = C - \mu_n$$

$a_n$  is Lagrangian multiplier related to  $w_n$

$\mu_n$  is Lagrangian multiplier related to  $\xi_n$



$$b = \frac{1}{N_M} \sum_{n \in M} (t_n - \sum_{n \in S} (a_n t_n k(x_n x_m)))$$

# PREDICTION USING KERNELS

$$y(x) = w^T \phi(x_n) + b$$

$$w = \sum_{n=1}^N a_n t_n \phi(x_n)$$

$a_n$  is a Lagrangian multiplier

$$y(x) = \sum_{n=1}^N a_n t_n k(x, x_n) + b$$

Any data point  $a_n = 0$  will not appear in the sum

New  
Data

Training data

Training data target (-1,1)

# METHODS OF USING SVM OF FEATURE SELECTION

---

## 1. Linear SVM:

Evaluation of learned **weights** in Linear SVM

$$y(x) = w^T \phi(x_n) + b \qquad w = \sum_{n=1}^N a_n t_n \phi(x_n)$$

## 2. Recursive Feature Elimination:

Iterative evaluation of significance of features in SVM classification:

## 3. Feature Vector Machines

Variation of Lasso like SVM that generate kernels on features not samples.