



# CS 559: Gaussian Process & Supportive Vector Machine

Lecture 6



# Outline

- Exam
- Model Selection
- Gaussian Process
- Supportive Vector Machine
- Exam Review



## Exam:

- Date: 3/23/2021 11:59 PM – 3/26/2021 11:59 PM
- Time: 1.5 hours
- Format: Online via Canvas Lockdown Browser.
- Question Types: Multiple choice, True/False, Matching, etc...
- Topics: From Lecture 1 to Lecture 6



# Empirical Risk Minimization

$$\min_w \frac{1}{n} \sum_{i=1}^N l(h_w(x_i), y_i) + \lambda r(w) \quad (1)$$

Loss                      Regularizer

- The loss function – a continuous function penalizing training error
- The regularizer – a continuous function penalizing classifier complexity



# Loss Function: Classifications

Hinge-Loss:

$$\max[1 - h_w(x_i)y_i, 0]^p$$

- Usage: SVM
- When  $p = 1$ , Standard SVM: the loss function denotes the size of the margin between linear separator and its closest points in either class.
- When  $p = 2$ , Squared Hingeless SVM: Only differentiable everywhere with  $p = 2$ .

Log-Loss:

$$\log(1 + e^{-h_w(x_i)y_i})$$

- Usage: Logistic Regression
- Its outputs are well-calibrated probabilities.
- Value range from 0 to 1.
- Goal: to minimize its value.

Zero-One Loss:

$$\delta(\text{sign}(h_w(x_i)) \neq y_i)$$

- Usage: Actual Classification Loss
- Characteristic: Not convex nor continuous.

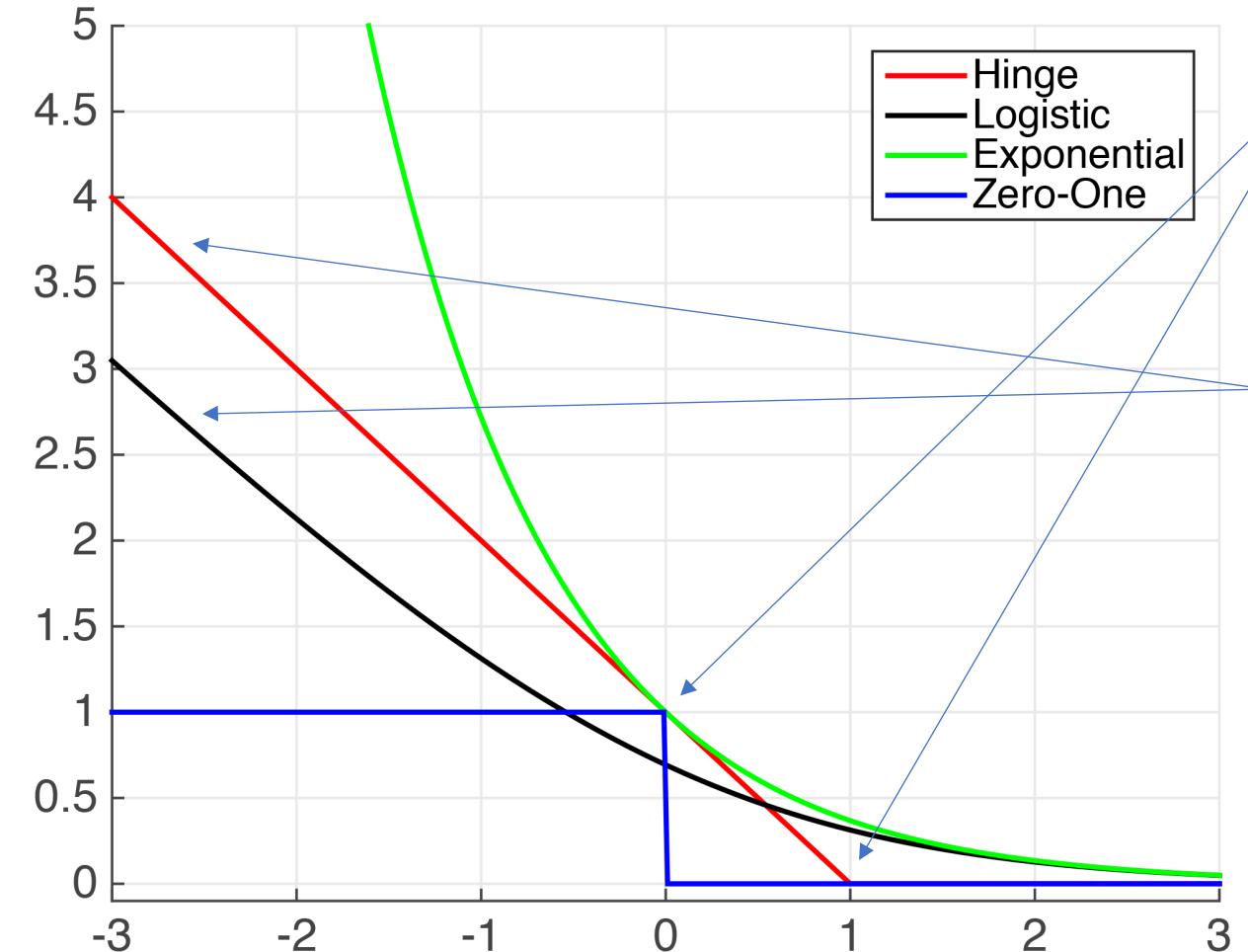
Exponential Loss:

$$\exp(-h_w(x_i)y_i)$$

- Usage: AdaBoost
- Characteristic: Very aggressive
- Misprediction increases exponentially with the value of  $-h_w(x_i)y_i$ .
- May cause problems with noisy data



# Classification Loss Functions Summary



Which functions are strict upper bounds on the 0/1-loss?

What can you say about the hinge-loss and the log-loss as  $z \rightarrow -\infty$ ?



# Loss Functions: Regressions

Squared Loss:

$$(h(x_i) - y_i)^2$$

- Most popular regression loss function.
- Estimates Mean.
- It is differentiable everywhere.
- Sensitive to outliers.

Absolute Loss:

$$|h(x_i) - y_i|$$

- Another Most popular regression loss function.
- Estimates Median.
- Less sensitive than squared loss function.
- Not differentiable at 0.

Huber Loss:

$$\begin{cases} \frac{1}{2}(h(x_i) - y_i)^2, & \text{if } |h(x_i) - y_i| < \delta \\ \delta(|h(x_i) - y_i| - \delta/2), & \text{Otherwise} \end{cases}$$

- Also known as Smooth Absolute Loss
- “Best of Both Worlds”
- Once-differentiable.
- When the loss is large, use absolute loss.
- When the loss is small, use squared loss.

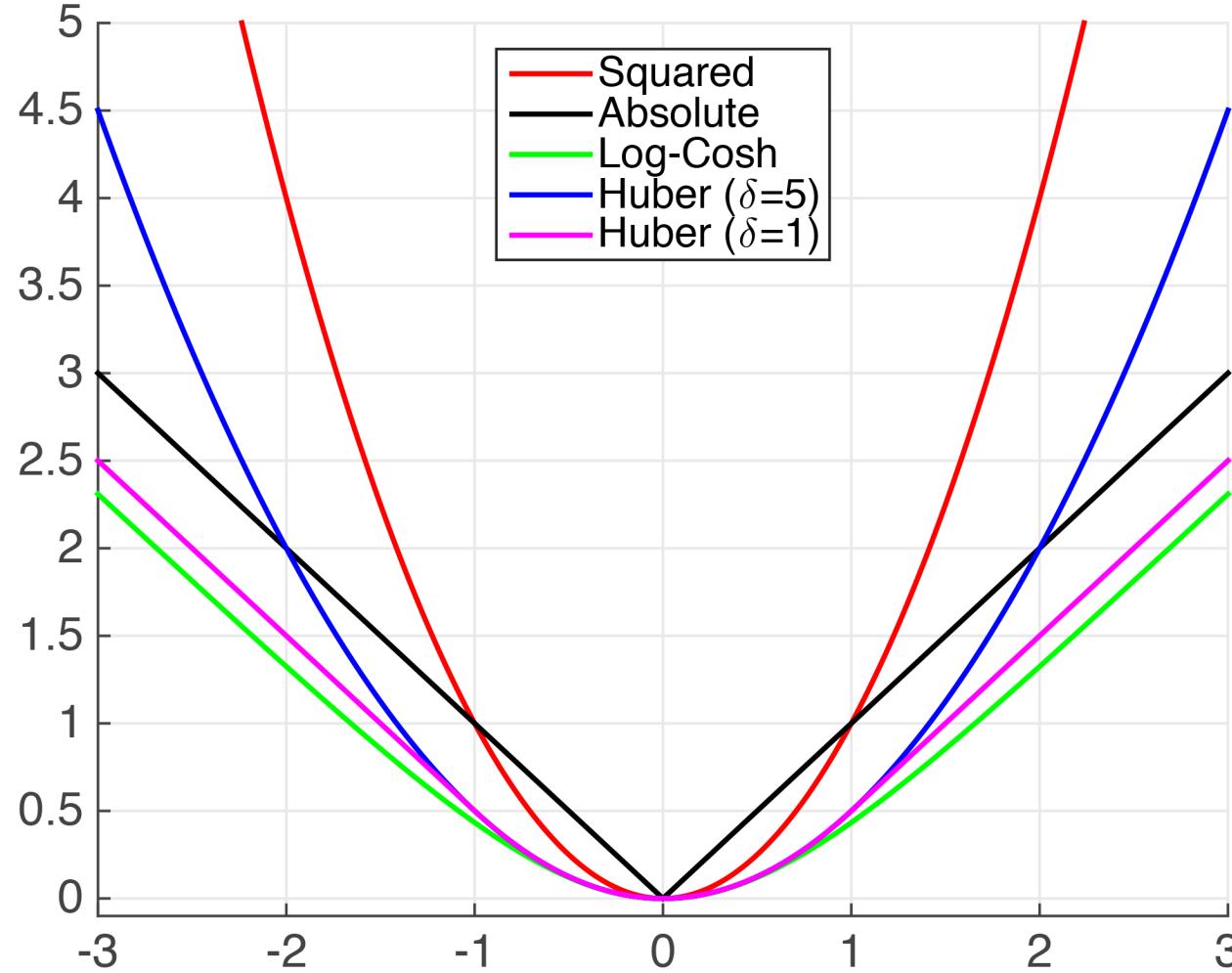
Log-Cosh Loss:

$$\log(\cosh(h(x_i) - y_i))$$
$$\cosh\left(\frac{e^x + e^{-x}}{2}\right)$$

- Similar to Huber Loss.
- Differentiable everywhere.



# Regression: Loss Functions Summary



See how sensitive functions are!



# Regularizers

We can write Equation (1) into different formation of optimization problem.

$$\begin{aligned} & \min_{w,b} \sum_{i=1}^N l(h_w(x_i), y_i) + \lambda r(w) \\ \Leftrightarrow & \min_{w,b} \sum_{i=1}^N l(h_w(x_i), y_i) \text{ s.t. } r(w) \leq B \end{aligned} \tag{2}$$

Equation (2) indicates that for each  $\lambda \geq 0$ , the regularizer  $r(w)$  is less than or equal to a constant number  $B$

# Regularizers

$l_2$  Regularization:

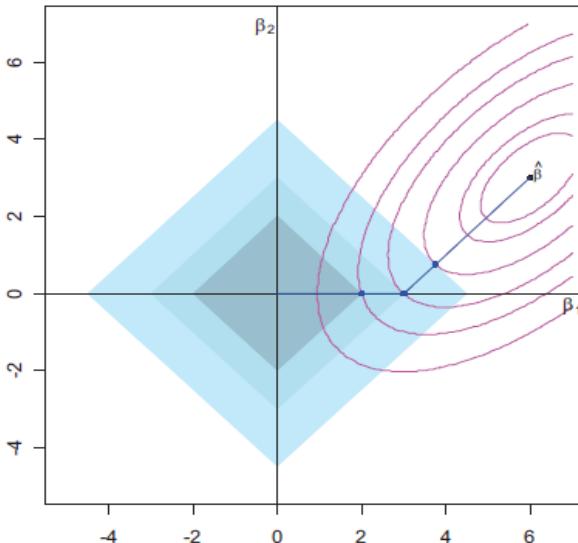
$$r(w) = w^T w = \|w\|^2$$

- Strictly convex and differentiable
- Uses weights on all features – dense solutions

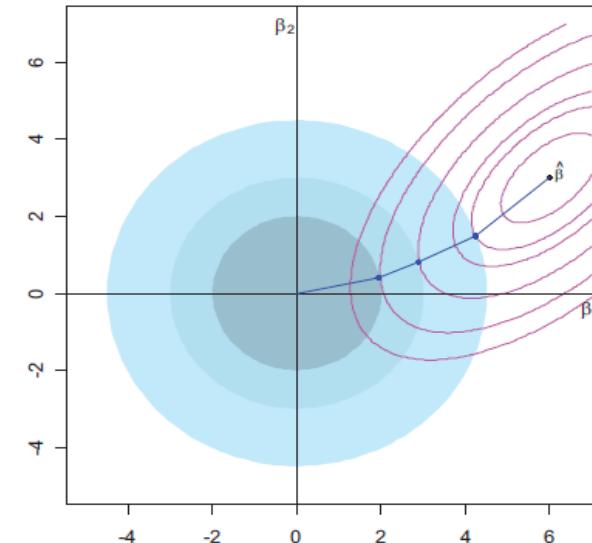
$l_1$  Regularization:

$$r(w) = \|w\|$$

- Convex but not strictly
- Not differentiable at 0
- Sparse Solutions



$l_1$



$l_2$

Lecture 6 - GP & SVM

Elastic Net Regularization  
 $\lambda\|w\| + (1 - \lambda)\|w\|^2$



## Loss and Regularizer

## Comments

Least Squares  $\frac{1}{n} \sum_{i=1}^n (w^T x_i - y_i)^2$

- Square Loss
- No Regularization
- Closed form solution:  $w = (X X^T)^{-1} X y^T$
- $X = [x_1, \dots, x_n], y = [y_1, \dots, y_n]$

Ridge Regression  $\min_w \frac{1}{n} \sum_{i=1}^n (w^T x_i - y_i)^2 + \lambda \|w\|^2$

- Squared Loss
- $l_2$  regularization
- $w = (X X^T + \lambda I)^{-1} X y^T$

Lasso Regression  $\min_w \frac{1}{n} \sum_{i=1}^n (w^T x_i - y_i)^2 + \lambda \|w\|$

- Sparsity inducing – good for feature selection
- Convex but not strict (no unique solution)
- Not differentiable at 0
- Solve with sub-gradient descent

Elastic Net

$$\min_w \frac{1}{n} \sum_{i=1}^n (w^T x_i - y_i)^2 + \lambda \|w\| + (1 - \lambda) \|w\|^2$$

- Unique solution and sparsity inducing
- Dual of squared-loss SVM
- Non-differentiable

Logistic Regression  $\min_{w,b} \frac{1}{n} \sum_{i=1}^n \log(1 - e^{-y_i(w^T x_i + b)})$

- Solve with gradient descent
- Probability estimation

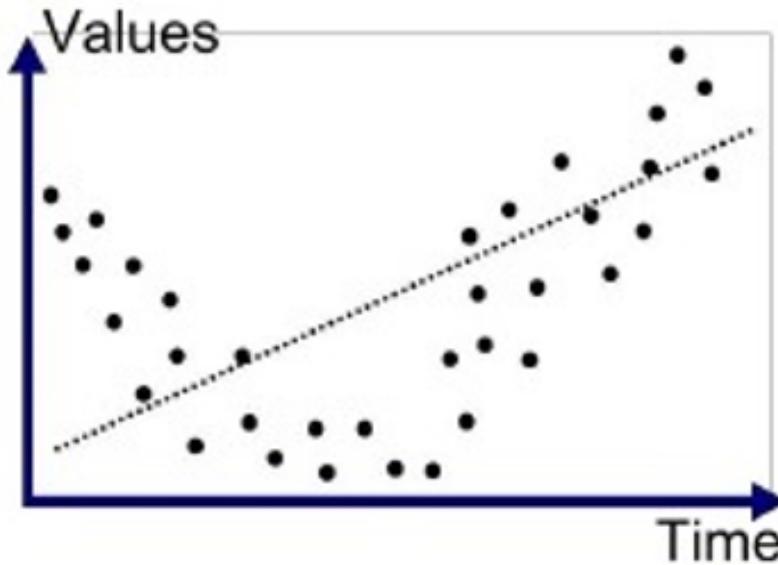
Linear SVM

$$\min_w C \sum_{i=1}^N \max[1 - y_i(w^T x_i + b), 0] + \lambda \|w\|^2$$

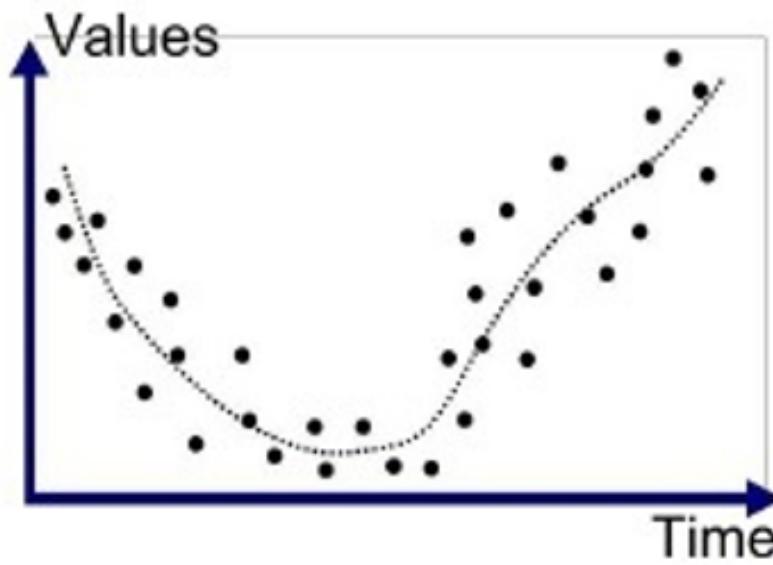
- Typically  $l_2$  regularized sometimes  $l_1$
- Quadratic program
- Kernelize – Solve very efficiently



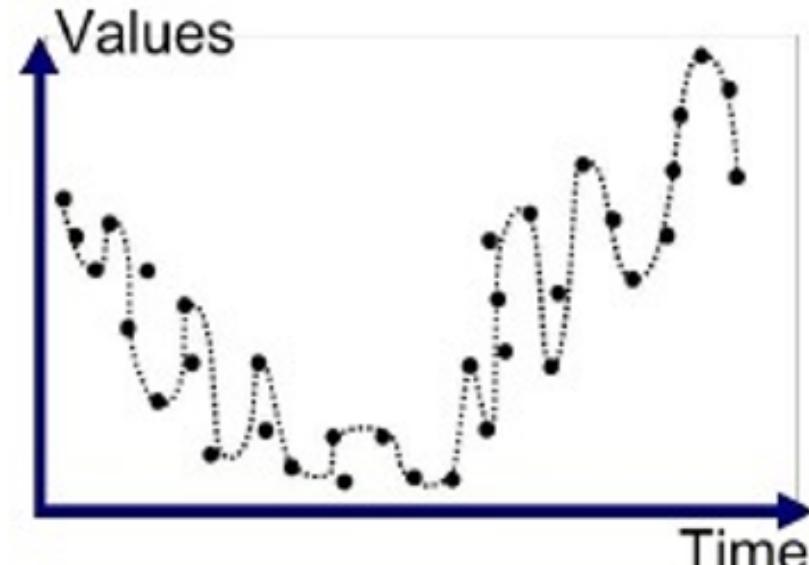
# Overfitting vs. Underfitting



Underfitted



Good Fit/Robust



Overfitted



# Overfitting vs. Underfitting

$$\min_w \frac{1}{n} \sum_{i=1}^n l(h_w(x_i), y_i) + \lambda r(w)$$

↑      ↑  
Loss      Regularizer

A question rise is on  $\lambda$  – how should we handle it?

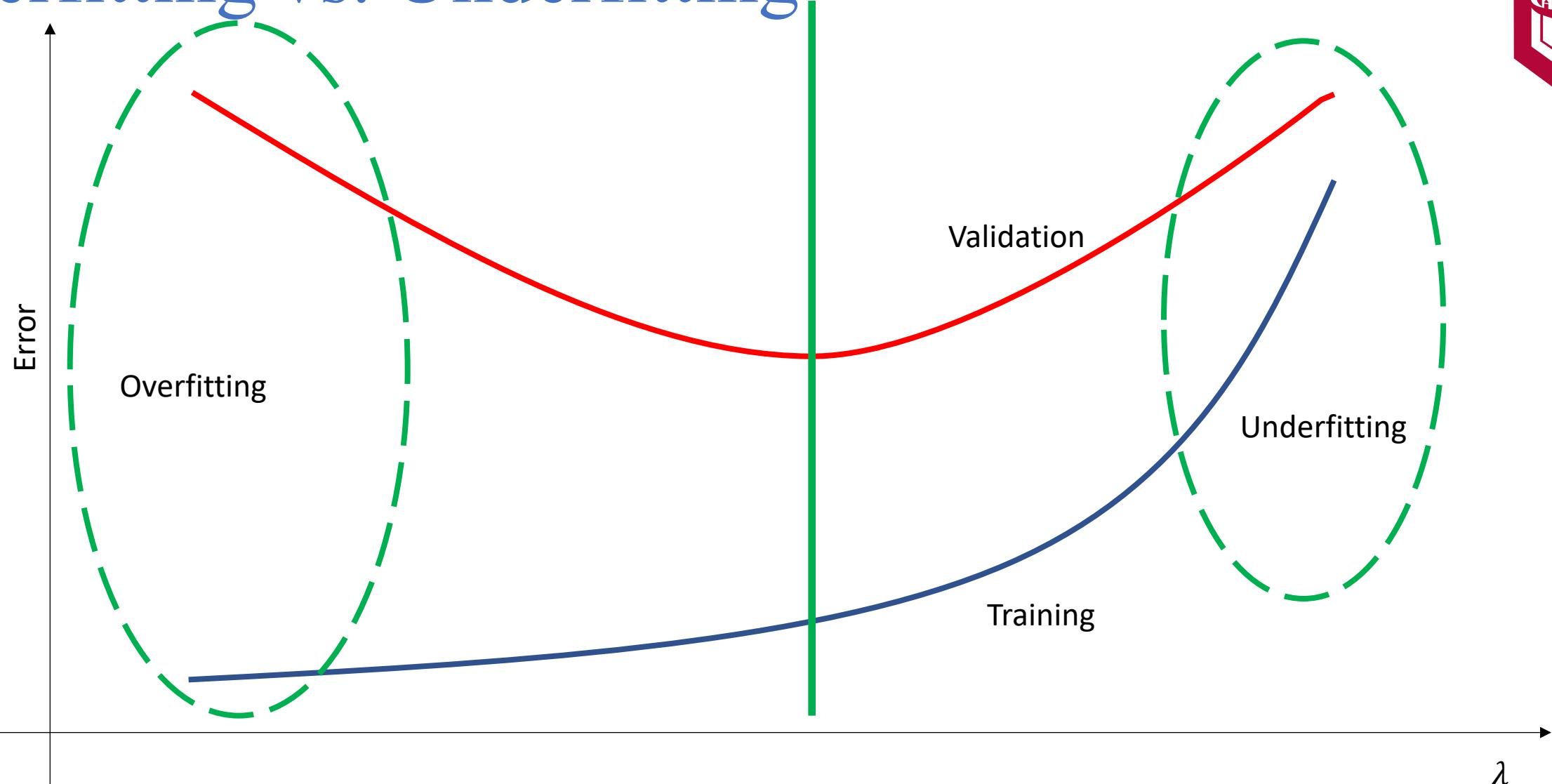
Leads to two possible scenarios: **Underfitting** or **Overfitting**

**Underfitting** – The model does not learn enough through the training process. Both errors on training and test will be high. (The model is relatively simpler than what it supposed to be.)

**Overfitting** – The model learns too much from training data set. The test error may rise because the model will reflects on the exact patterns of training data set in the test set. (The model is relatively complex than what it should be.)



# Overfitting vs. Underfitting





# Overfitting vs. Underfitting

The big question is.... **HOW do we find the exact point to stop?**

Typical range:  $10^{-5}, 10^{-4}, 10^{-3}, \dots, 10^1, 10^2, 10^3, \dots$

## **K-fold Cross Validation**

- Divide the training data  $k$  portion equally.
- Train on  $k-1$  of them and leave  $k^{\text{th}}$  one for the validation.
- Do this  $k$  times and find the best  $\lambda$  that gives the smallest error on the validation set.
  - When the data is too small, you leave only one observation for the validation.

For a specific search on  $\lambda$

- Do with  $\lambda = \{\dots, 10^{-3}, 10^{-2}, \dots, 10^1, 10^2, \dots\}$  to find the best magnitude.
- Then do for the specific value of  $\lambda$  value in the best chosen magnitude
- Example: Let's say you learned that  $\lambda = 10^{-2}$  is the best magnitude on the 1<sup>st</sup> try. Then do with  $\lambda = \{0.01, 0.015, 0.02, \dots\}$  on the 2<sup>nd</sup> try.



# Bias-Variance Tradeoff

$$\begin{aligned} E_{x,y,D}[(h_D(x) - y)^2] &= E_{x,y,D} \left[ \left[ (h_D(x) - \bar{h}(x)) + (\bar{h}(x) - y) \right]^2 \right] \\ &= E_{x,D} \left[ (h_D(x) - \bar{h}(x))^2 \right] + 2E_{x,y,D} \left[ (h_D(x) - \bar{h}(x))(\bar{h}(x) - y) \right] + E_{x,y} \left[ (\bar{h}(x) - y)^2 \right] \end{aligned} \quad (1)$$

2<sup>nd</sup> Term

$$\begin{aligned} E_{x,y,D} \left[ (h_D(x) - \bar{h}(x))(\bar{h}(x) - y) \right] &= E_{x,y} [E_D[h_D(x) - \bar{h}(x)](\bar{h}(x) - y)] \\ &= E_{x,y} \left[ (E_D[h_D(x)] - \bar{h}(x))(\bar{h}(x) - y) \right] \\ &= E_{x,y}[0] \end{aligned}$$

$E_D[h_D(x)] = \bar{h}(s)$

3<sup>rd</sup> Term

$$\begin{aligned} E \left[ (\bar{h}(x) - y)^2 \right] &= E \left[ \left( (\bar{h}(x) - \bar{y}(x)) + (\bar{y}(x) - y) \right)^2 \right] \\ &= E[(\bar{y}(x) - y)^2] + E \left[ (\bar{h}(x) - \bar{y}(x))^2 \right] + 2E \left[ (\bar{h}(x) - \bar{y}(x))(\bar{y}(x) - y) \right] \end{aligned}$$

$$E_{x,y,D}[(h_D(x) - y)^2] = E_{x,D} \left[ (h_D(x) - \bar{h}(x))^2 \right] + E_{x,y} \left[ (\bar{h}(x) - \bar{y}(x))^2 \right] + E_x[(\bar{y}(x) - y)^2] \quad (2)$$

$Var(x) = \text{Variance}$        $\text{Bias}^2$        $\epsilon = \text{Noise}$

# Bias-Variance Tradeoff

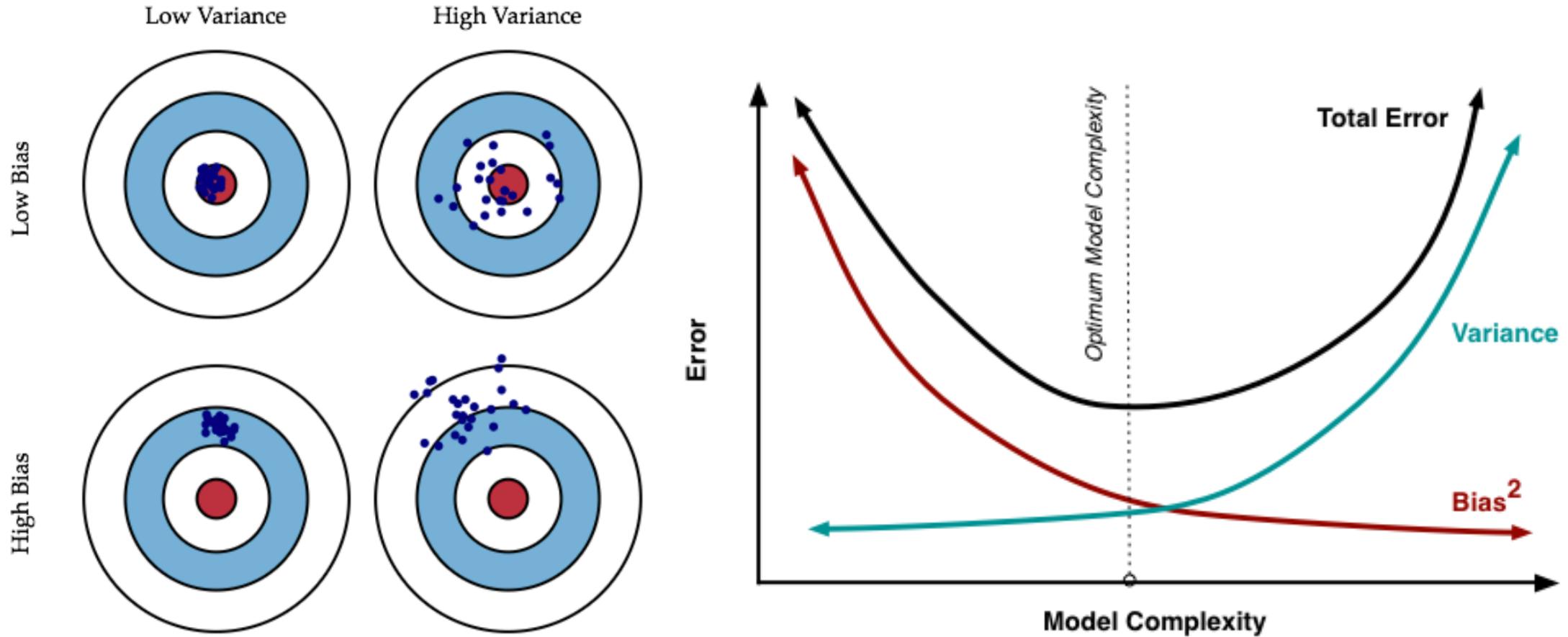


Figure: Graphical illustration of bias vs. variance



# Kernel Methods



# Kernel Function

- Given: There are a large set of fixed linear basis functions.
  - Problem: It all depends on how complex the basis functions are.
  - Solution: Use “dual trick” that depends on the amount data than the function.
- 
- We have a set of basis functions  $\phi(x)$  that map inputs  $x$  to a feature space.
  - This feature space appears in the dot product  $\phi(x)^T \phi(x')$  of input pairs,  $x, x'$ .
  - The kernel function  $k(x, x') = \phi(x)^T \phi(x')$  in feature space.
    - We are interested in kernel function not the set of basis functions.



# Dual Representations

The linear regression objective is the error.

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N [w^T \phi(x_n) - y_n]^2 + \frac{1}{2} \lambda \mathbf{w}^T \mathbf{w} \quad (30)$$

We set the gradient to 0.

$$E(\mathbf{w}) = \sum_n (\mathbf{w}^T \phi(x_n) - y_n) \phi(x_n) + \lambda \mathbf{w} = 0 \quad (31)$$

$$\mathbf{w} = -\frac{1}{\lambda} \sum_n (\mathbf{w}^T \phi(x_n) - y_n) \phi(x_n) \quad (32)$$

$\mathbf{w}$  is linear combination of inputs in feature space

$$\{\phi(x_n) | 1 \leq n \leq N\}$$



# Dual Representations

Substitute  $\mathbf{w} = \Phi\mathbf{a}$

Where  $\Phi = [\phi(x_1), \phi(x_2), \dots, \phi(x_N)]$

$$\mathbf{a} = \begin{bmatrix} a_1 \\ \vdots \\ a_N \end{bmatrix} \text{ and } a_n = -\frac{1}{\lambda}(\mathbf{w}^T \phi(x_n) - y_n)$$

Dual objective: minimize E with respect to a

$$E(a) = \frac{1}{2} \mathbf{a}^T \Phi^T \Phi \Phi^T \Phi \mathbf{a} - \mathbf{a}^T \Phi^T \Phi \mathbf{y} + \frac{\mathbf{y}^T \mathbf{y}}{2} + \frac{\lambda}{2} \mathbf{a}^T \Phi^T \Phi \mathbf{a} \quad (33)$$



# Gram Matrix

The Gram Matrix:  $K = \Phi^T \Phi$

Substitution K into (31)

$$E(a) = \frac{1}{2} a^T K K a - a^T K y + \frac{y^T y}{2} + \frac{\lambda}{2} a^T K a \quad (34)$$

$$\nabla E(a) = K K a - K y + \lambda K a$$

$$K y = K(K + \lambda I)a$$

$$a = (K - \lambda I)^{-1}y$$

Prediction  $y_* = \phi(x_*)^T w = \phi(x_*)^T \Phi a = k(x_*, X)(K + \lambda I)^{-1}y$  (35)

where  $(X, y)$  is the training set and  $(x_*, y_*)$  is a test instance

- Primal Solution: depends on # of basis functions
- Dual solution: depends on amount of data
  - Advantage: can use very large # of basis functions

Just need to know  $k$



# Constructing Kernels

Possibilities:

- Find mapping  $\Phi$  to feature space and let  $K = \phi^T \phi$
- Directly specify  $K$

A valid kernel must be positive semi-definite:

- $k$  must factor into the product of a transpose matrix by itself
- or all eigenvalues  $\geq 0$
- problem? It is not always easy to verify.



# Rules to construct Kernels

Let  $k_1(x, x')$  and  $k_2(x, x')$  be valid kernels.

The following kernels are also valid:

$$= ck_1(x, x') \text{ c is the constant } > 0$$

$$= f(x)k_1(x, x')f(x') \text{ f is the function}$$

$$= q(k_1(x, x')) \text{ q is the polynomial}$$

$$= \exp(k_1(x, x'))$$

$$k(x, x') = k_1(x, x') + k_2(x, x')$$

$$= k_3(\phi(x), \phi(x'))$$

$$x^T Ax' \text{ A is symmetric positive semi-definite}$$

$$= k_a(x_a, x'_a) + k_b(x_b, x'_b) \text{ } k_a \text{ & } k_b \text{ are valid kernel}$$

$$= k_a(x_a, x'_a)k_b(x_b, x'_b)$$

sets, strings, graphs, etc...

Example -

Lodhi, Saunders, Shawe-Taylor, Christianini, Watkins,

**Text Classification Using String Kernels**, JMLR, p. 419-444, 2002

Other Application?



# Kernel Construction Example

Consider a 2-D input space  $\mathbf{x} = (x_1, x_2)$  and suppose we are going to map into a feature space in the degree of 2 using

$$k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z})^2$$

$$\begin{aligned} k(\mathbf{x}, \mathbf{z}) &= (\mathbf{x}^T \mathbf{z})^2 = (x_1 z_1 + x_2 z_2)^2 \\ &= x_1^2 z_1^2 + 2x_1 z_1 x_2 z_2 + x_2^2 z_2^2 \\ &= (x_1^2, \sqrt{2}x_1 x_2, x_2^2)(z_1^2, \sqrt{2}z_1 z_2, z_2^2) \\ &= \phi(\mathbf{x})^T \phi(\mathbf{z}) \end{aligned}$$

The feature mapping comprises all possible second order terms.

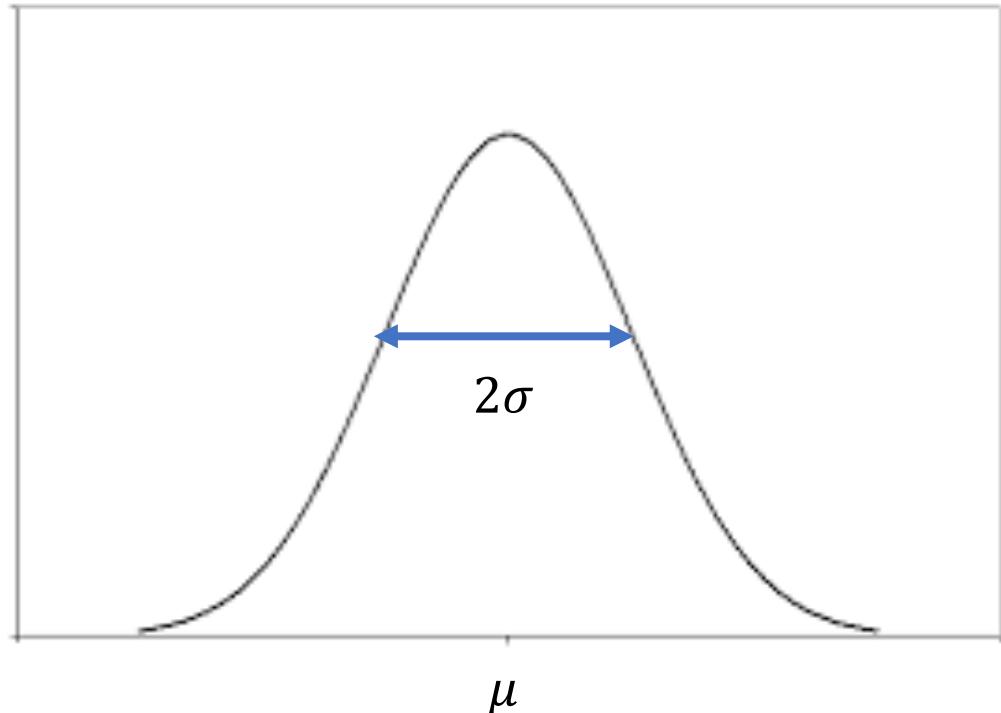
Which kernel functions to use? We can always use the functions whose validities are already confirmed.



# Gaussian Processes

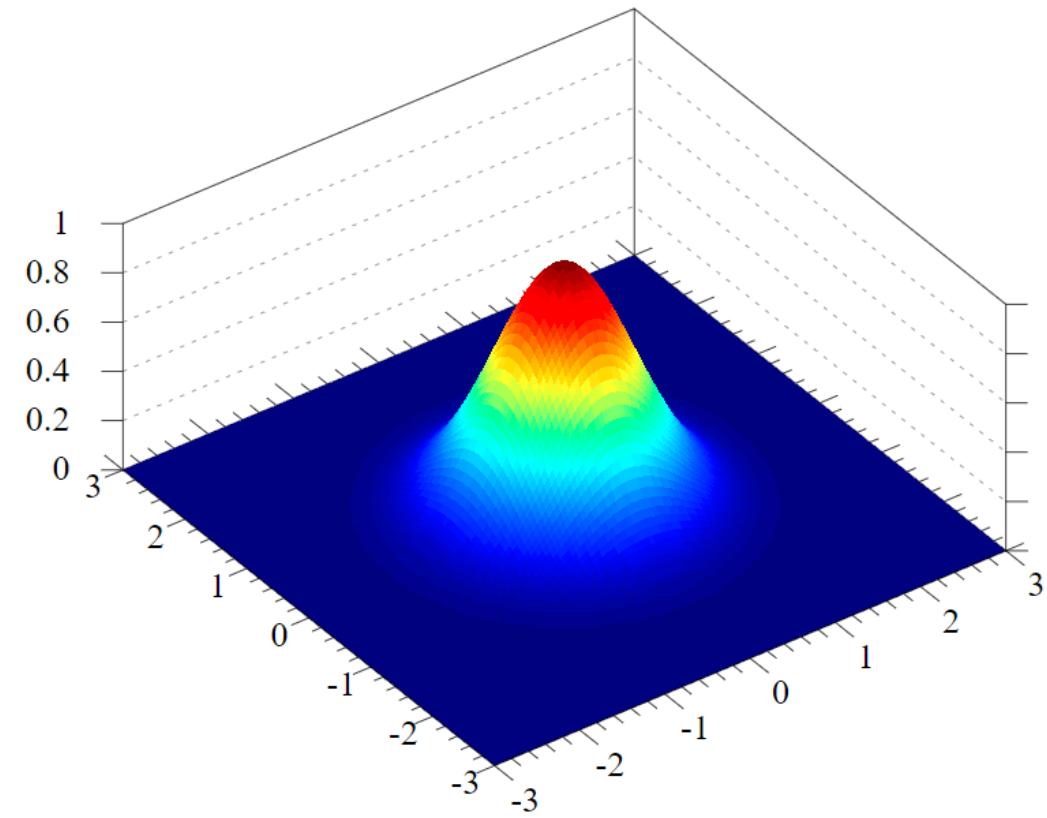
# Gaussian Distributions

**1-D Gaussian**



$$p(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

**Multivariate Gaussian**



$$p(x|\mu, \Sigma) = \frac{1}{\sqrt{|2\pi\Sigma|}} \exp \left\{ -\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

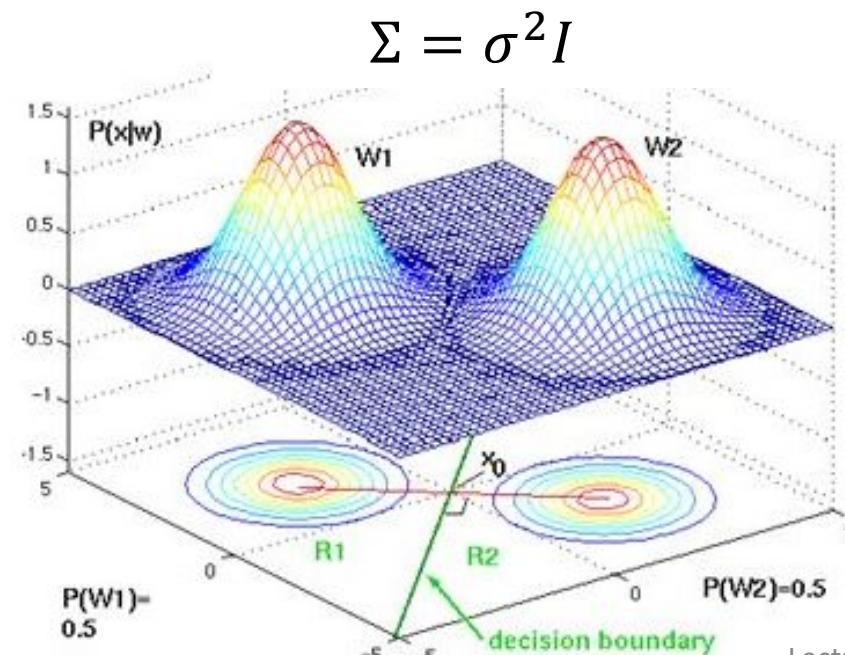
# Gaussian Distributions

## Multivariate Gaussian Example: A 2-D Gaussian

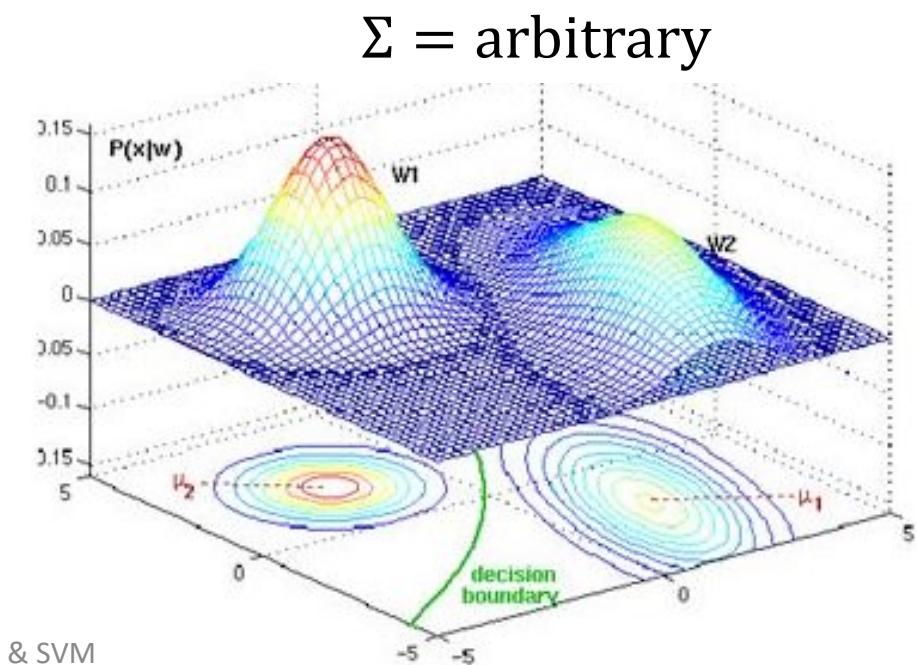
- A mean vector:  $\mu = [\mu_1, \mu_2]$

- A covariance matrix:  $\Sigma = \begin{bmatrix} \sigma_{11}^2 & \sigma_{21}^2 \\ \sigma_{12}^2 & \sigma_{22}^2 \end{bmatrix}$

where  $\sigma_{ij}^2 = E[(x_i - \mu_i)(x_j - \mu_j)]$  is covariance and  $\Sigma$  is symmetric “positive semi-definite”.



Lecture 6 - GP & SVM

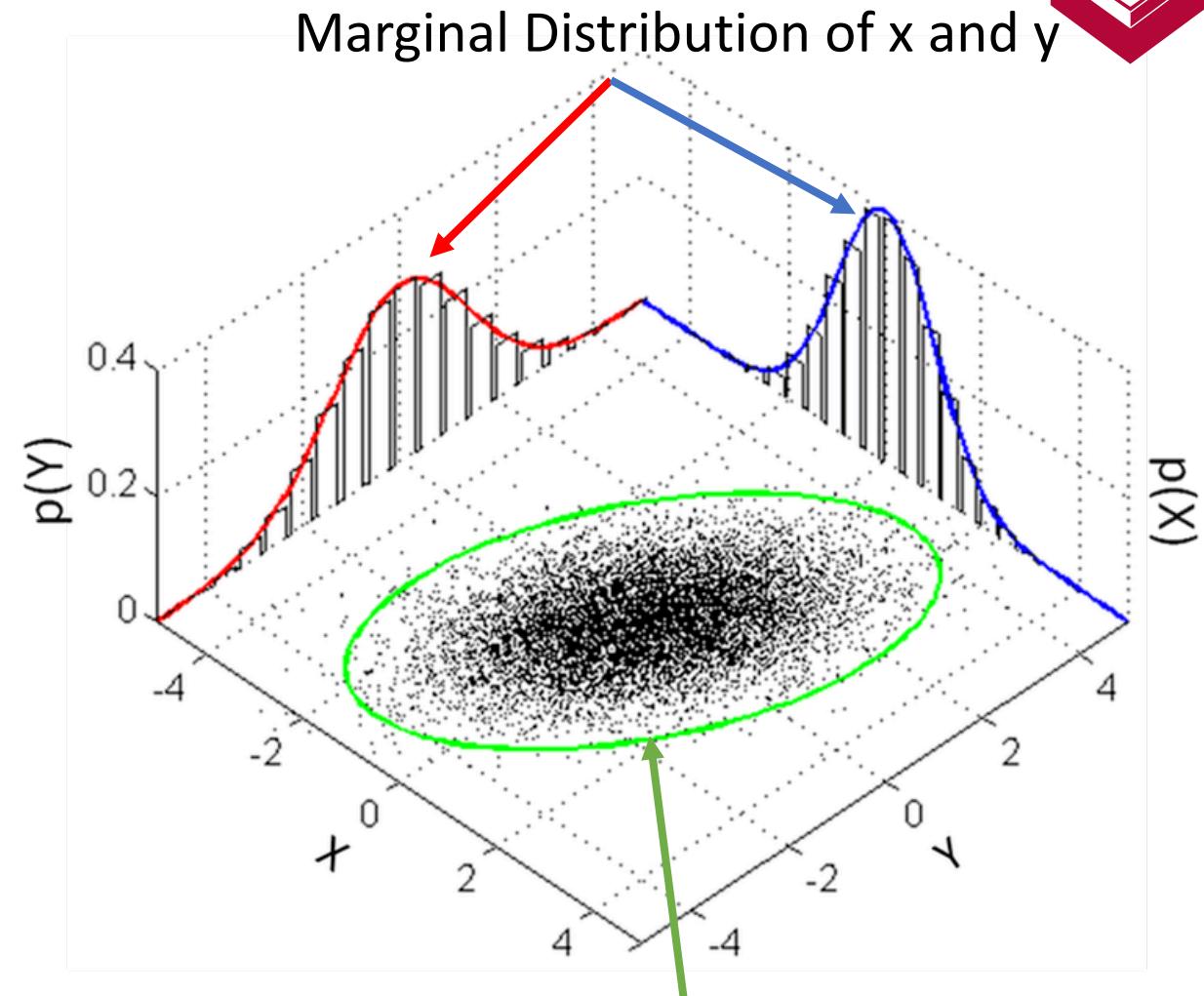


# Gaussian Distributions

## Marginal Distributions of Gaussians

- Given:  $x = (x_a, x_b)$ ,  $\mu = (\mu_a, \mu_b)$ ,
- $\Sigma = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}$
- Marginal distributions of Gaussians are Gaussian
- Marginal distribution - probability of variable w/o reference to other variables.

$$p(x_a) = \mathcal{N}(x_a | \mu_a, \Sigma_{aa})$$





# Gaussian Distributions

## Block Matrix Inversion

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix}^{-1} = \begin{bmatrix} (A - BD^{-1}C)^{-1} & -A^{-1}B(D - CA^{-1}B)^{-1} \\ -D^{-1}C(A - BD^{-1}C)^{-1} & (D - CA^{-1}B)^{-1} \end{bmatrix} = \begin{bmatrix} S_D^{-1} & -A^{-1}BS_A^{-1} \\ -D^{-1}CS_D^{-1} & S_A^{-1} \end{bmatrix}$$

Schur complements:

$$S_A = D - CA^{-1}B$$

$$S_D = A - BD^{-1}C$$

## Conditional Distribution

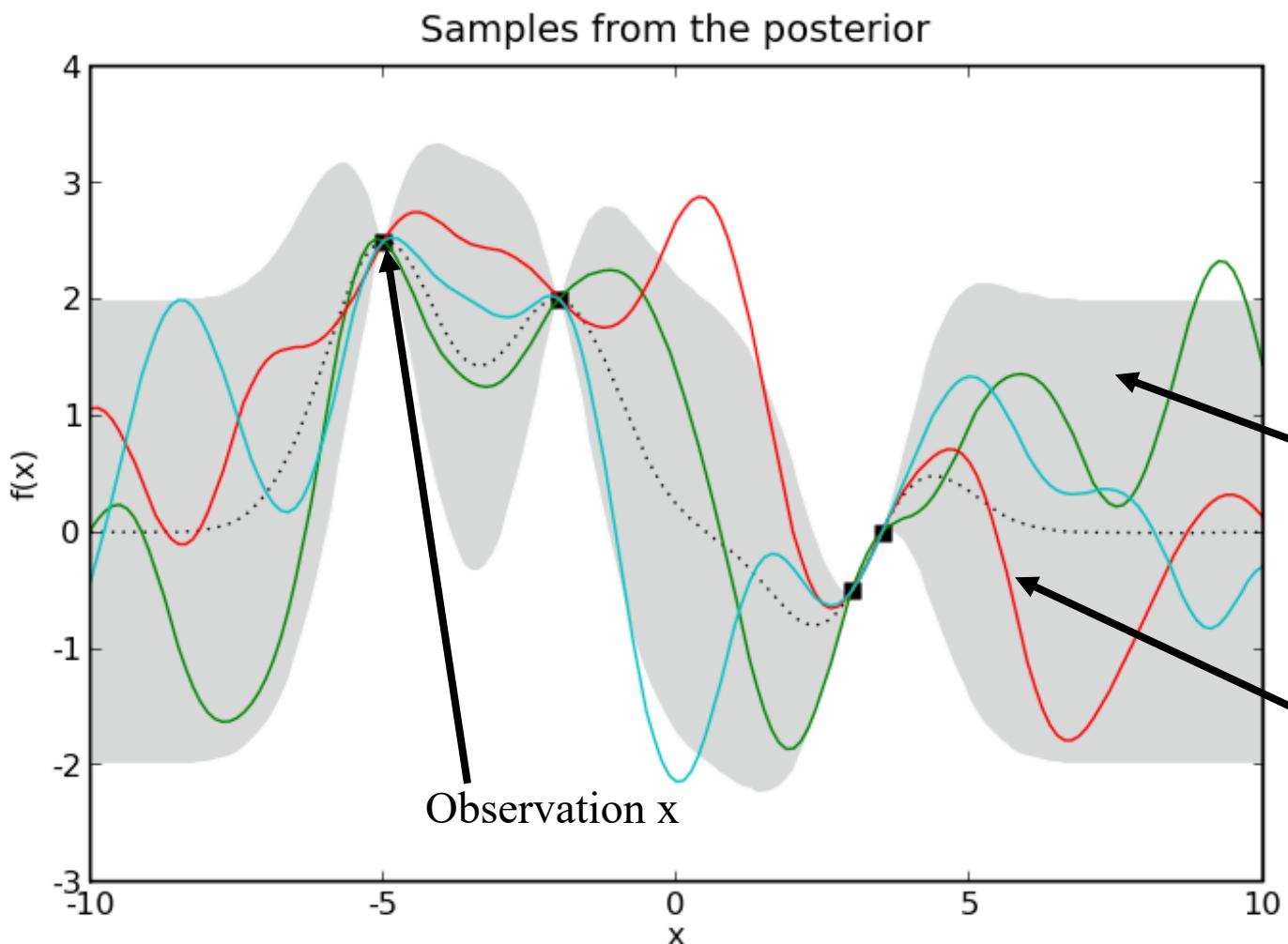
Notation:  $\Sigma = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}$   $\Lambda = \Sigma^{-1} = \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix}$

Distribution:  $p(x_a | x_b) = \mathcal{N}(x_a | \mu_{(a|b)}, \Lambda_{aa}^{-1})$

where  $\mu_{(a|b)} = \mu_a - \Lambda_{aa}^{-1}\Lambda_{ab}(x_b - \mu_a) = \mu_a - \Sigma_{ab}\Sigma_{bb}^{-1}(x_b - \mu_a)$

and Schur complement  $\Lambda_{aa}^{-1} = \Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba}$

# Regression



## Linear Regression:

$$y = w^T x + \epsilon$$

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

Linear regression gives a predictive model for one particular parameter  $w$ .

## How can we do regression?

- We would like a method that could provide confidence during the estimation.
- How about kernelize linear regression?



# Gaussian Process

Definition:

- Probability distribution indexed by an arbitrary set.
- Each element gets a Gaussian distribution over the reals with mean  $\mu$
- Distributions are dependent/correlated as defined by  $k(x,z)$
- Any finite subset of indices defines a multivariate Gaussian distribution
- $f \sim GP(\mu, k)$

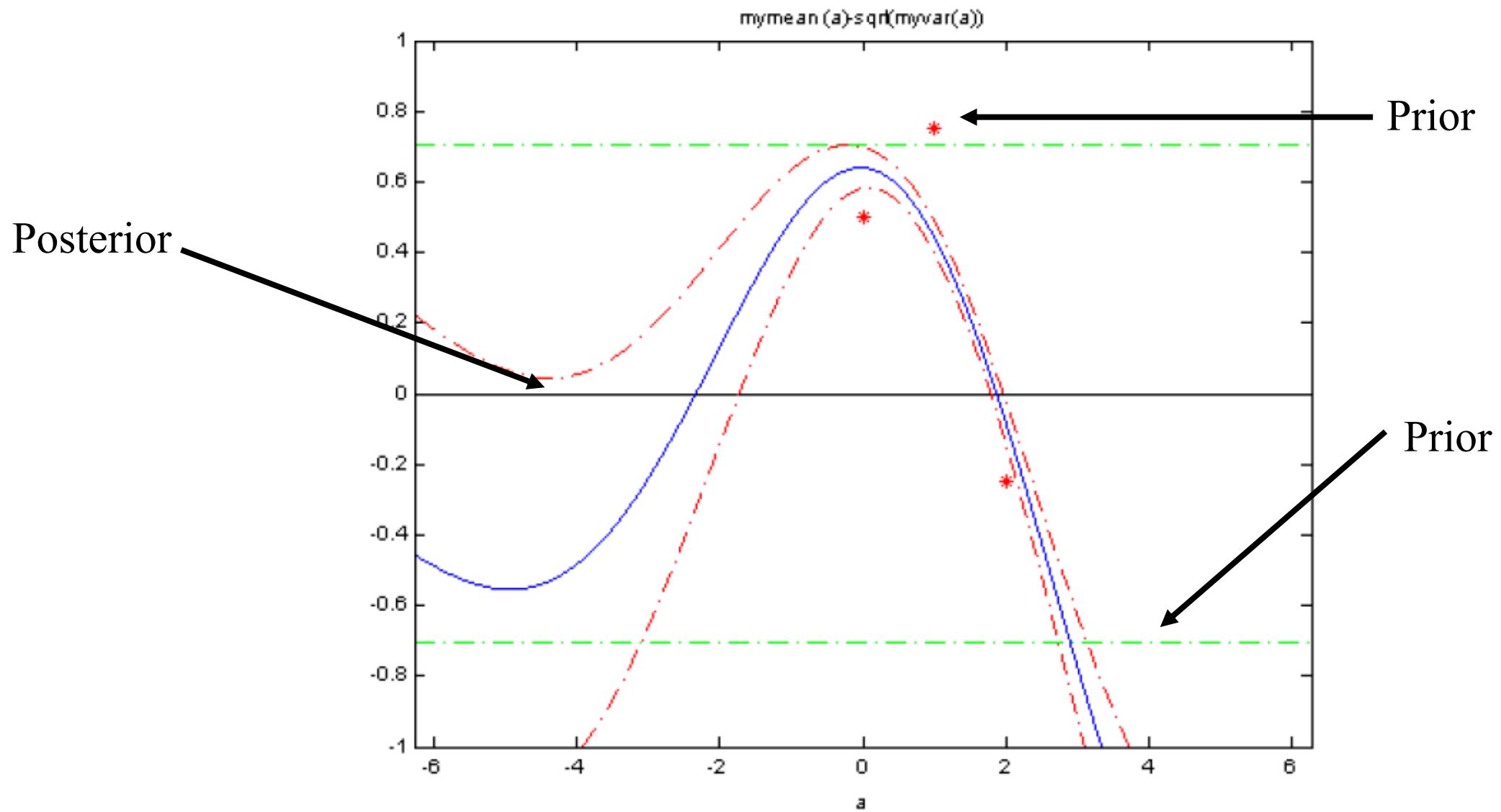
What we get...

- Distribution over functions and functions for any domain index sets.

General Procedure:

- Start with GP prior
- Get some data
- Compute a posterior

# Gaussian Process





# GP – Weight Space View

**GP = Bayesian ridge regression in feature space + kernel trick to carry out computations**

**Training Set:**  $D = \{(x_i, y_i) | i = 1, \dots, n\}$

- $x_i$  is the input vector of D-dimension
- $y_i$  is a scalar output or target

**Training Data: X and Y**

Design matrix  $X = [x_1 | \dots | x_n] \in \mathbb{R}^{D \times n}$

Target vector:  $Y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \in \mathbb{R}^n$

$f(x) = x^T w \in \mathbb{R}, x, w \in \mathbb{R}^D$

$y = f(x) + \epsilon = x^T w + \epsilon \in \mathbb{R}$  where  $\epsilon \sim \mathcal{N}(0, \sigma^2) \in \mathbb{R}$



# GP – Weight Space View LR w/ Gaussian Noise

Calculate the Likelihood and put w prior over parameters w,  $w \sim \mathcal{N}(0, \Sigma)$ :

$$p(Y|X, w) = \prod_{i=1}^n p(y_i|x_i^T w) \quad (1)$$

$$= \prod_{i=1}^n \mathcal{N}(x_i^T w, \sigma^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(y_i - x_i^T w)^2}{2\sigma^2}\right]$$

$$= \frac{1}{(2\pi\sigma^2)^{\frac{n}{2}}} \exp\left[-\frac{1}{2\sigma^2} \|Y - X^T w\|^2\right]$$

$$= \mathcal{N}(X^T w, \sigma^2 I) \quad (2)$$



# GP – Weight Space View LR w/ Gaussian Noise

The prior:  $w \sim \mathcal{N}(0, \Sigma)$

The posterior:

$$p(w|X, Y) = \frac{p(Y|X, w)p(w)}{p(Y|X)} = \frac{p(Y|X, w)p(w)}{\int p(Y|X, w)dw} \quad (3)$$

$$\begin{aligned} &= \frac{\mathcal{N}(X^T w, \sigma^2) \mathcal{N}(0, \Sigma)}{\int \mathcal{N}(X^T w, \sigma^2) \mathcal{N}(0, \Sigma) dw} \\ &\sim \mathcal{N}(X^T w, \sigma^2) \mathcal{N}(0, \Sigma) \end{aligned} \quad (4)$$



# GP – Weight Space View LR w/ Gaussian Noise

$$\begin{aligned} p(w|X, Y) &\sim \mathcal{N}(X^T w, \sigma^2) \mathcal{N}(0, \Sigma) \\ &\sim \exp\left[-\frac{1}{2\sigma^2}(Y - X^T w)^T(Y - X^T w)\right] \exp\left[-\frac{1}{2}w^T \Sigma^{-1} w\right] \quad (5) \\ &\sim \exp\left[-\frac{1}{2}(w - \bar{w})^T \left(\frac{1}{\sigma^2}XX^T + \Sigma^{-1}\right)(w - \bar{w})\right] \\ &\sim \mathcal{N}(\bar{w}, A^{-1}) \end{aligned}$$

Ridge Regression

After “completing  
the square”      (6)

where  $\bar{w} = \sigma^{-2}(\sigma^{-2}XX^T + \Sigma^{-1})^{-1}XY \in \mathbb{R}^D$  and  $A = (\sigma^{-2}XX^T + \Sigma^{-1}) \in \mathbb{R}^{D \times D}$

$A^{-1} \in \mathbb{R}^{D \times D}$

# GP – Weight Space View LR w/ Gaussian Noise



Use the posterior to predict  $f$  in a test point  $x_*$ :

$$f(x) = x^T w \in \mathbb{R}, x, w \in \mathbb{R}^D$$

$$y = f(x) + \epsilon = x^T w + \epsilon \in \mathbb{R} \text{ where } \epsilon \sim \mathcal{N}(0, \sigma^2) \in \mathbb{R}$$

$$p(f_* | x_*, X, Y) = \int p(f_* | x_*, w) p(w | Y, X) dw \quad (7)$$

$\underbrace{\delta(f_*, x_*^T, w)}_{\mathcal{N}(\bar{w}, A^{-1})}$

$$= N(x_*^T \bar{w}, x_*^T A^{-1} x_*) \quad (8)$$



# GP – Explicit Feature Space

$$\begin{aligned}\phi(x) &= [x_1, x_1 x_2^2, \dots]^T \in \mathbb{R}^N \\ \Phi(X) &= [\phi(x_1) | \phi(x_2) | \cdots | \phi(x_n)] \in \mathbb{R}^{N \times n}\end{aligned}$$

$\Phi(X)$  is the aggregation of columns  $\phi(x)$  for all cases in the training set.

The model is  $f(x) = \phi(x)^T w$  and  $\phi(x), w \in \mathbb{R}^N$

$$y = f(x) + \epsilon = \phi(x)^T w + \epsilon \in \mathbb{R}$$



# GP – Explicit Feature Space

The predictive distribution after feature map:

$$p(f_*|x_*, X, Y) = \mathcal{N}(\phi(x_*)^T \bar{w}, \phi(x_*)^T A^{-1} \phi(x_*)) \quad (9)$$

where  $\bar{w} = \sigma^{-2} (\sigma^{-2} \Phi(X) \Phi(X)^T + \Sigma^{-1}) \Phi(X) Y \in \mathbb{R}^N$  and

$A^{-1} \in \mathbb{R}^{N \times N}$        $\in \mathbb{R}^{N \times n}$        $\in \mathbb{R}^{n \times 1}$

$$A = (\sigma^{-2} \Phi(X) \Phi(X)^T + \Sigma^{-1}) \in \mathbb{R}^{N \times N}$$

A problem is we need  $N \times N$  matrix inversion...



# GP – Explicit Feature Space

We can shorten this equation using

$$\begin{aligned}\phi_* &= \phi(x) \in \mathbb{R}^N, \Phi = \Phi(X) \in \mathbb{R}^{N \times n}, \\ \bar{w} &= \sigma^{-2}(\sigma^{-2}\Phi\Phi^T + \Sigma^{-1})\Phi Y \in \mathbb{R}^N, \\ A &= (\sigma^{-2}\Phi\Phi^T + \Sigma^{-1}) \in \mathbb{R}^{N \times N}\end{aligned}$$

We only need to work with  **$n \times n$  matrices**, not  $N \times N$ !

$$\begin{aligned}p(f_*|x_*, X, Y) &= \mathcal{N}(\phi_*^T \bar{w}, \phi_*^T A^{-1} \phi) \\ &= \mathcal{N}(\sigma^{-2} \phi_*^T [\sigma^{-2} \Phi \Phi^T + \Sigma^{-1}]^{-1} \Phi Y, \phi_*^T [\sigma^{-2} \Phi \Phi^T + \Sigma^{-1}]^{-1} \phi_*)\end{aligned}\tag{10}$$

Let  $K = \Phi(X)^T \Sigma \Phi(X)$ ,

$$p(f_*|x_*, X, Y)$$

$$= \mathcal{N}(\phi_*^T \Sigma \Phi(K + \sigma^2 I)^{-1} Y, \phi_*^T \Sigma \phi_* - \phi_*^T \Sigma \Phi(K + \sigma^2 I)^{-1} \Phi^T \Sigma \phi_*)\tag{11}$$

$$\mathbb{R}^{1 \times n}$$

$$\mathbb{R}^{n \times n}$$

$$\mathbb{R}^{n \times 1}$$



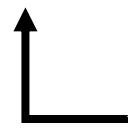
# GP – Explicit Feature Space

$$\begin{aligned} p(f_* | x_*, X, Y) & \\ = \mathcal{N}\left(\phi_*^T \Sigma \Phi(K + \sigma^2 I)^{-1} Y, \phi_*^T \Sigma \phi_* - \phi_*^T \Sigma \Phi(K + \sigma^2 I)^{-1} \Phi^T \Sigma \phi_*\right) & (11) \end{aligned}$$

The feature space **always** enters in the form of

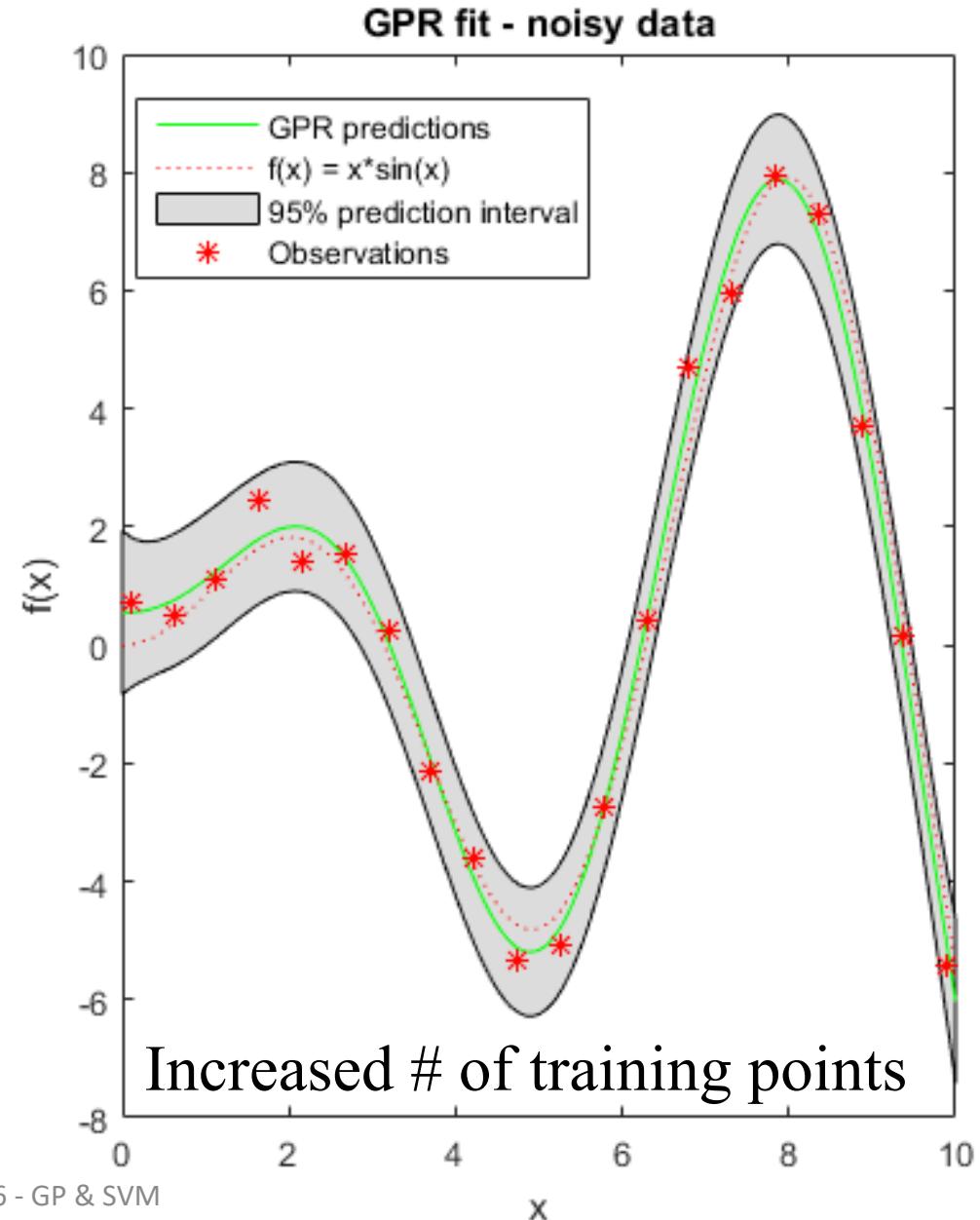
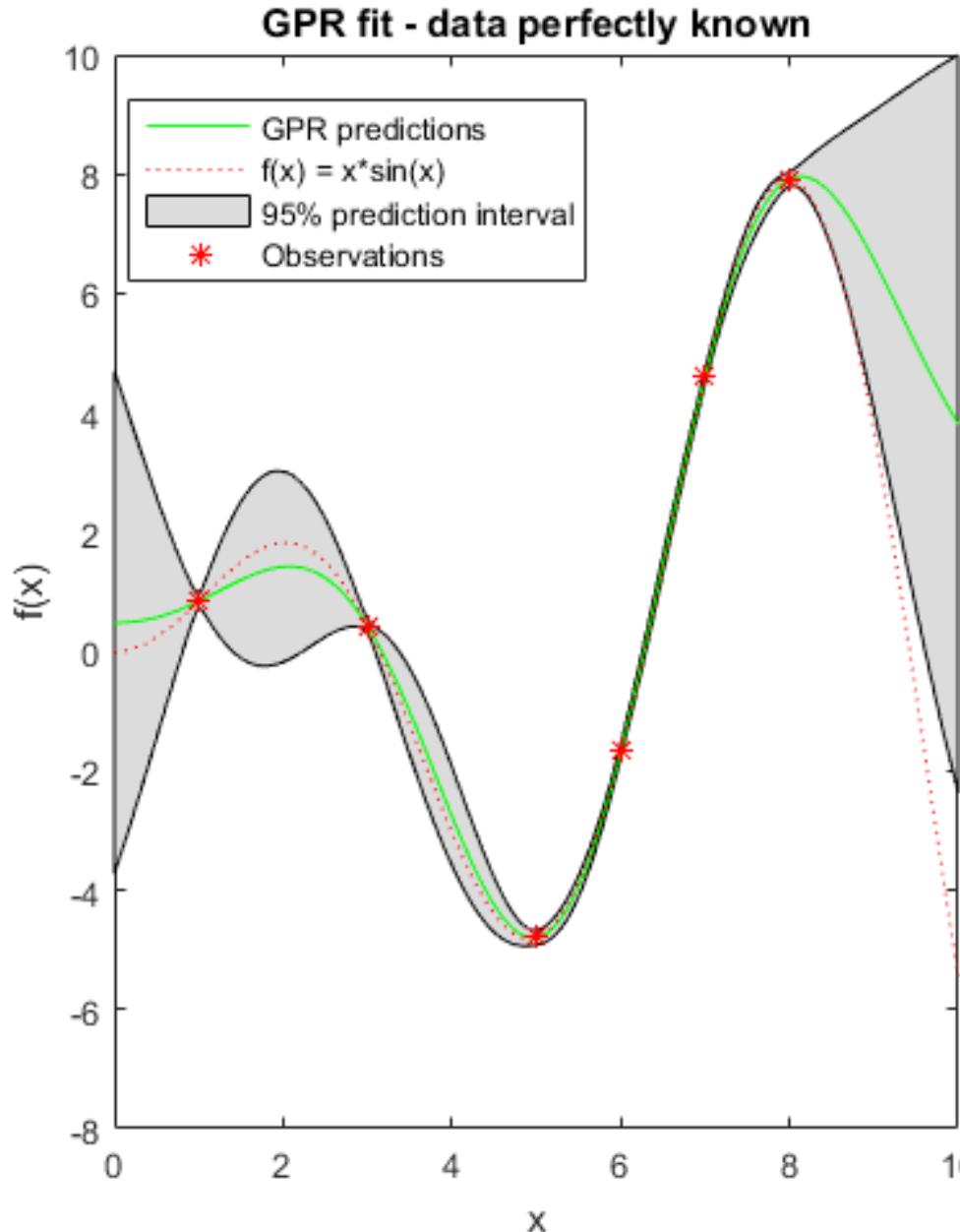
$$\phi_*^T \Sigma \phi_*, \phi_*^T \Sigma \phi, \phi^T \Sigma \phi \in \mathbb{R}^{n \times n}$$

Let form of  $k = \phi(x)^T \Sigma \phi(x')$  - a covariance function or kernel.



**Inner product** in the feature space:  $\psi(x) = \Sigma^{1/2} \phi(x)$

# GP – Results





# GP – Function Space View

Definition: GP is a collection of random variables, any finite number of which have a joint Gaussian distribution.

$$m(x) = E[f(x)] \quad (\text{mean function})$$

$$k(x, x') = E[(f(x) - m(x))(f(x') - m(x'))^T] \in \mathbb{R} \quad (\text{covariance func.})$$

GP is completely specified by its mean function  $m(x)$  and covariance function  $k(x, x')$ .

$$f(x) \sim GP(m(x), k(x, x')) \in \mathbb{R}, x \in \mathbb{R}^D \quad (12)$$



# GP – Function Space View

For each  $x \in \mathbb{R}^D$  we associate a Gaussian variable  $f(x)$  such that  $f(x) \in \mathcal{N}(m(x), k(x, x))$  and its correlation with other  $f(x')$  variables is  $k(x, x')$ .

$$\mathbb{R} \ni f(x) \sim \mathcal{N}(m(x), k(x, x))$$

$$\begin{bmatrix} f(x) \\ f(x') \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} m(x) \\ m(x') \end{bmatrix}, \begin{bmatrix} k(x, x) & k(x', x) \\ k(x, x') & k(x', x') \end{bmatrix}\right) \quad (13)$$



# GP – noise free observations

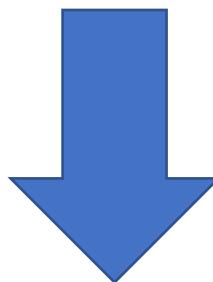
Training set:  $D = \{(x_i, f_i) | i = 1, \dots, n\}$

$$X = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix} \in \mathbb{R}^{n \times D}, n \text{ training inputs.}$$

$$f = \begin{bmatrix} f_1 \\ \vdots \\ f_n \end{bmatrix} \in \mathbb{R}^{n \times D}, n \text{ training inputs}$$

$$X_* = \begin{bmatrix} x_{*1}^T \\ \vdots \\ x_{*m}^T \end{bmatrix} \in \mathbb{R}^{m \times D}, m \text{ test inputs}$$

$$f_* = \begin{bmatrix} f_{*1} \\ \vdots \\ f_{*m} \end{bmatrix} \in \mathbb{R}^{m \times D}, m \text{ test inputs}$$



$$\begin{bmatrix} f \\ f_* \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} 0_n \\ 0_m \end{bmatrix}, \begin{bmatrix} k(X, X) & k(X, X_*) \\ k(X_*, X) & k(X_*, X_*) \end{bmatrix} \right)$$

(14)



# GP – noise free observations

$$\begin{bmatrix} f \\ f_* \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} 0_n \\ 0_m \end{bmatrix}, \begin{bmatrix} k(X, X) & k(X, X_*) \\ k(X_*, X) & k(X_*, X_*) \end{bmatrix} \right)$$

$$\begin{aligned}\Sigma &= \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix} \quad \Lambda = \Sigma^{-1} = \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix} \\ p(x_a | x_b) &= \mathcal{N} \left( x_a \middle| \mu_a|_b, \Lambda_{aa}^{-1} \right) \\ \mu_a|_b &= \mu_a - \Lambda_{aa}^{-1} \Lambda_{ab} (x_b - \mu_a) = \mu_a - \Sigma_{aa} \Sigma_{bb}^{-1} (x_b - \mu_a) \\ \Lambda_{aa}^{-1} &= \Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba}\end{aligned}$$

$$\begin{aligned}p(f_* | X_*, X, f) \\ = \mathcal{N} \left( k(X_*, X) k(X, X)^{-1} f, k(X_*, X_*) - k(X_*, X) k(X, X)^{-1} k(X, X_*) \right)\end{aligned}\tag{15}$$

If  $X_* = X \Rightarrow f_* = f$  and the covariance is 0.

$$\begin{aligned}p(f_* | X_*, X, f) \\ = \mathcal{N} \left( \phi_*^T \Sigma \Phi (K + \sigma^2 I)^{-1} f, \phi_*^T \Sigma \phi_* - \phi_*^T \Sigma \Phi (K + \sigma^2 I)^{-1} \Phi^T \Sigma \phi_* \right)\end{aligned}\tag{16}$$

# GP – noise observations

$$\begin{aligned} y &= f(x) + \epsilon \in \mathbb{R} \\ \epsilon &\sim \mathcal{N}(0, \sigma^2) \in \mathbb{R} \end{aligned}$$

$$\begin{bmatrix} f \\ f_* \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} 0_n \\ 0_m \end{bmatrix}, \begin{bmatrix} k(X, X) & k(X, X_*) \\ k(X_*, X) & k(X_*, X_*) \end{bmatrix} \right)$$

$$\Rightarrow cov(y_p, y_q) = k(x_p, x_q) + \sigma^2 \delta_{p,q} \quad (\text{Dirac Delta Function})$$

$$\Rightarrow cov([y_1, \dots, y_n]) = k(X, X) + \sigma^2 I \in \mathbb{R}^{n \times n}$$

$$\Rightarrow \begin{bmatrix} y \\ f_* \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} 0_n \\ 0_m \end{bmatrix}, \begin{bmatrix} k(X, X) + \sigma^2 I & k(X, X_*) \\ k(X_*, X) & k(X_*, X_*) \end{bmatrix} \right)$$

(New joint distribution)  
(17)

The posterior for the noisy observations:

$$p(f_* | X, Y, X_*) = \mathcal{N} \left( \bar{f}_*, cov(f_*) \right) \quad (18)$$



# GP – noise observations

$$\Rightarrow \begin{bmatrix} y \\ f_* \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} 0_n \\ 0_m \end{bmatrix}, \begin{bmatrix} k(X, X) + \sigma^2 I & k(X, X_*) \\ k(X_*, X) & k(X_*, X_*) \end{bmatrix} \right)$$

$$\Sigma = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix} \quad \Lambda = \Sigma^{-1} = \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix}$$

$$p(x_a | x_b) = \mathcal{N} \left( x_a \mid \mu_a|_b, \Lambda_{aa}^{-1} \right)$$

$$\mu_a|_b = \mu_a - \Lambda_{aa}^{-1} \Lambda_{ab} (x_b - \mu_a) = \mu_a - \Sigma_{aa} \Sigma_{bb}^{-1} (x_b - \mu_a)$$

$$\Lambda_{aa}^{-1} = \Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba}$$

$$p(f_* | X, Y, X_*) = \mathcal{N} \left( \bar{f}_*, cov(f_*) \right)$$

where

$$\bar{f}_* = E[f_* | X, Y, X_*] = k(X_*, X)[k(X, X) + \sigma^2 I]^{-1} Y \in \mathbb{R}^m \quad (19)$$

$$cov(f_*) = k(X_*, X_*) - k(X_*, X)[k(X, X) + \sigma^2 I]^{-1} K(X, X_*) \in \mathbb{R}^{m \times m} \quad (20)$$

$$\mathcal{N}(\phi_*^T \Sigma \Phi (K + \sigma^2 I)^{-1} Y, \phi_*^T \Sigma \phi_* - \phi_*^T \Sigma \Phi (K + \sigma^2 I)^{-1} \Phi \Sigma \phi_*) \xleftarrow{\text{Weight Space View}}$$

$$\text{If } k = \phi(x)^T \Sigma \phi(x')$$

# GP – noise observations

**Short Notation** for a single test point  $x_*$ :

$$K = k(X, X) \in \mathbb{R}^{n \times n} \quad k(x_*) = k_* = k(X, x_*) = \begin{bmatrix} k(x_1, x_*) \\ \vdots \\ k(x_n, x_*) \end{bmatrix} \in \mathbb{R}^n$$

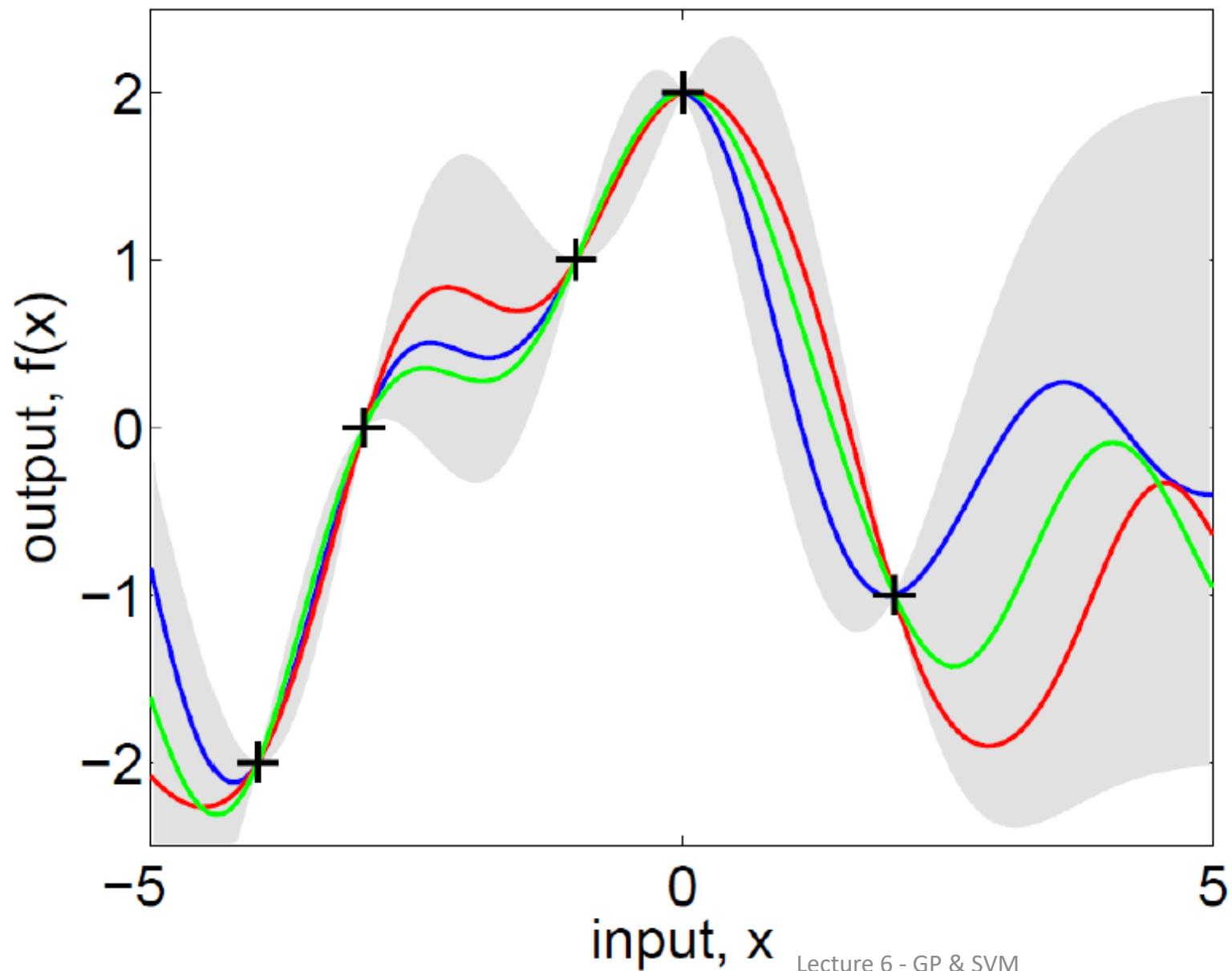
$$K_* = k(X, X_*) \in \mathbb{R}^{n \times m}$$

$$\bar{f}_* = k_*^T [K + \sigma^2 I]^{-1} Y \in \mathbb{R} \quad (21)$$

$\mathbb{R}^{1 \times n}$        $\mathbb{R}^{n \times n}$        $\mathbb{R}^n$

$$cov(f_*) = k(x_*, x_*) - k_*^T [K + \sigma^2 I]^{-1} k_* \in \mathbb{R} \quad (22)$$

$\mathbb{R}$        $\mathbb{R}^{1 \times n}$        $\mathbb{R}^{n \times n}$        $\mathbb{R}^n$





# GP – noise observations

Linear combination of observation Y - linear predictor

$$\text{Let } \beta^T = k_*^T [K + \sigma^2 I]^{-1} \in \mathbb{R}^{1 \times n}$$

$$\bar{f}_* = \sum_{i=1}^N \beta_i^T y_i$$

$$\bar{f}_* = \underbrace{k_*^T}_{\mathbb{R}^{1 \times n}} \underbrace{[K + \sigma^2 I]^{-1}}_{\mathbb{R}^n} \underbrace{Y}_{\mathbb{R}^n} \in \mathbb{R}$$

A linear combination of n kernel function - representer theorem

$$\text{Let } \alpha = [K + \sigma^2 I]^{-1} Y$$

$$\bar{f}_* = \sum_{i=1}^n \alpha_i^T k(x_i, x_*)$$

$$\bar{f}_* = \underbrace{\{[K + \sigma^2 I]^{-1} Y\}^T}_{\mathbb{R}^{n \times n}} \underbrace{k_*}_{\mathbb{R}^{n \times 1}} \in \mathbb{R}$$



# Kernel Construction Example

Consider a 2-D input space  $\mathbf{x} = (x_1, x_2)$  and suppose we are going to map into a feature space in the degree of 2 using

$$k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z})^2$$

$$\begin{aligned} k(\mathbf{x}, \mathbf{z}) &= (\mathbf{x}^T \mathbf{z})^2 = (x_1 z_1 + x_2 z_2)^2 \\ &= x_1^2 z_1^2 + 2x_1 z_1 x_2 z_2 + x_2^2 z_2^2 \\ &= (x_1^2, \sqrt{2}x_1 x_2, x_2^2)(z_1^2, \sqrt{2}z_1 z_2, z_2^2) \\ &= \phi(\mathbf{x})^T \phi(\mathbf{z}) \end{aligned}$$

The feature mapping comprises all possible second order terms.

Which kernel functions to use? We can always use the functions whose validities are already confirmed.



# Supportive Vector Machine

- Maximum Margin Classifiers
- Overlapping class distributions
- Classification
- Regression
- Kernel SVM - Regression



# SVM

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

- $\phi(x)$  a fixed feature-space transformation
- $b$  the bias parameter
- training set - linearly separable in feature space
- Data:  $\{(x_i, t_i) | n = 1, \dots, n\}$
- $t_i \in \{-1, 1\}$

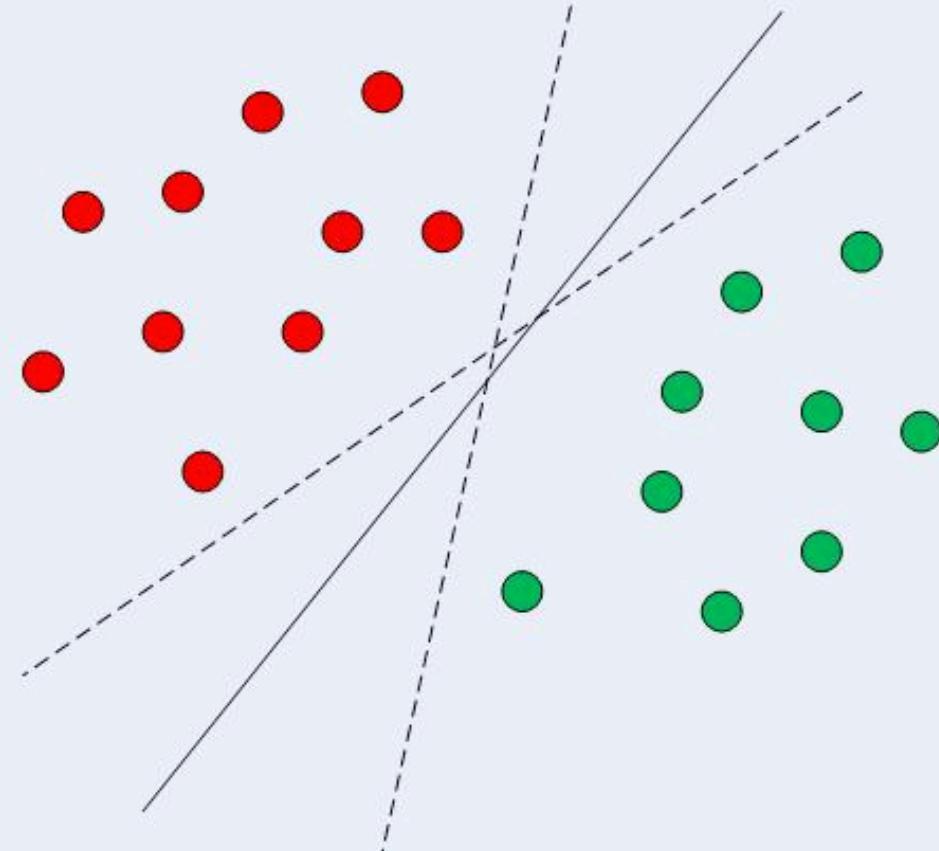
# SVM vs. Perceptron Algorithm

## Perceptron Algorithm

- guarantees to find a solution in a finite number of steps
- initial value for  $w$  and  $b$  starts from arbitrary chosen value..

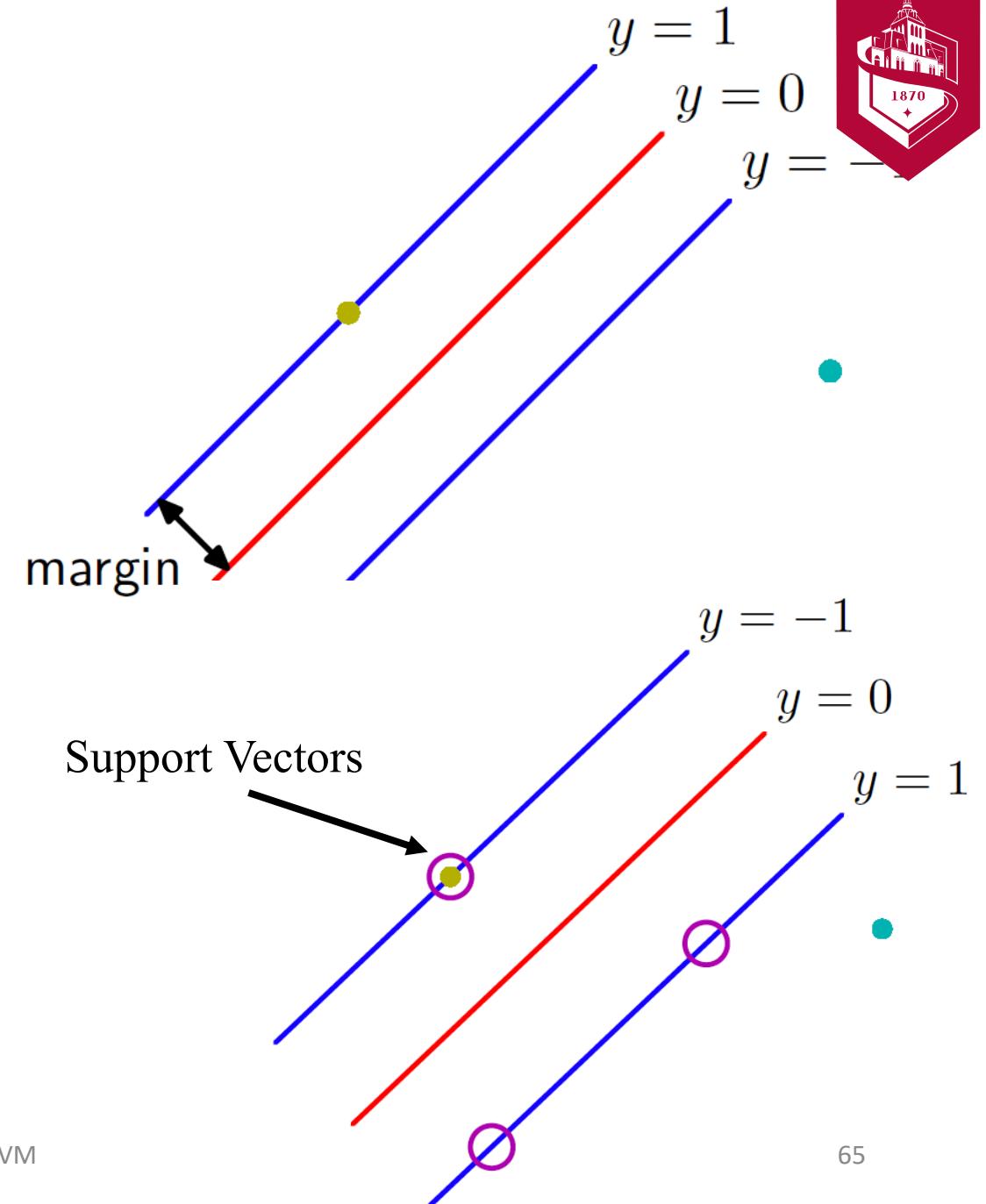
If there are multiple solutions of all of which classify the training data set exactly, then we should try to find the one that give **the smallest generalization error!**

Which is the better classification?



# SVM

- An extension of the Perceptron developed by Rosenblatt in 1958.
- Concept of the margin - the smallest distance between the decision boundary and any of the samples.
- Choose the max margin
- Using the optimal boundary, determine the best ***hyperplane*** by minimizing the probability of error relative to the learned density mode.
- hyperplane becomes independent of data points that are not ***support vectors (SV)***.



# SVM

Hyperplane is  $y(x) = 0$

Right classification  $t_n y(x_n) > 0$

$$\frac{t_n y(x_n)}{\|w\|} = \frac{t_n [w^T \phi(x_n) + b]}{\|w\|}$$

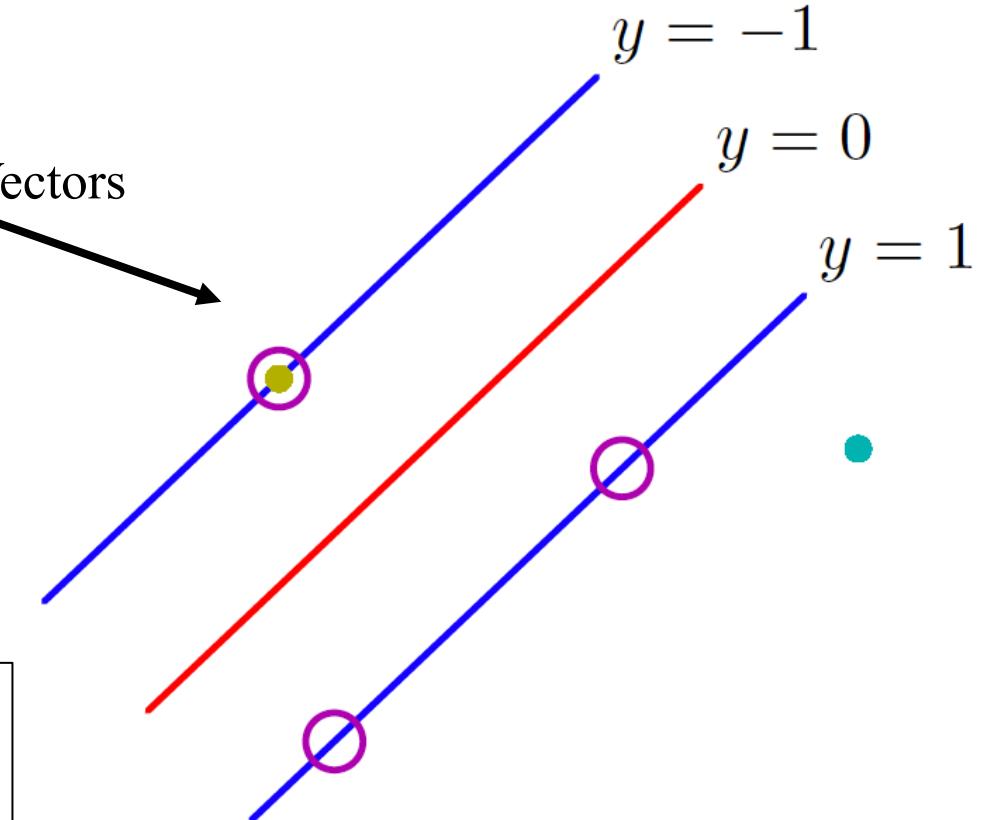
Support Vectors

Optimize the parameters w and b to maximize the distance.

Maximum margin solution is

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

w does not depend on n. canonical representation  
of the decision hyperplane.





# SVM

Direction solution is quite complex, we convert it into an equivalent problem (rescale)

$$t_n(w^T \phi(x) + b) = 1$$

for the point that is closest to the surface. All data points will satisfy the constraints

$$t_n(w^T \phi(x) + b) \geq 1$$

once the margin has been maximized there will be at least two active constraints.

So we maximize  $\|w\|^{-1}$  and this is same as minimizing  $\|w\|^2$

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



# Lagrange multipliers

We introduce Lagrange Multipliers  $a_n \geq 0$

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N a_n \{t_n (\mathbf{w}^T \phi(x_n) + b) - 1\} \quad (1)$$

- a strategy for finding the local maxima and minima of a function subject to equality constraints

Derivative:

$$\begin{aligned} \mathbf{w} &= \sum_{n=1}^N a_n t_n \phi(x_n) \\ 0 &= \sum_{n=1}^N a_n t_n \end{aligned}$$



# Dual Representation

Dual representation of the maximum margin problem

$$L(a) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{m=1}^N \sum_{n=1}^N a_n a_m t_n t_m k(x_n, x_m) \quad (2)$$

Use the kernel function.

This takes the form of a quadratic programming problem

$$\begin{aligned} a_n &\geq 0 \\ 0 &= \sum_{n=1}^N a_n t_n \end{aligned}$$



# Prediction

$$y(\mathbf{x}) = \sum_{n=1}^N a_n t_n k(\mathbf{x}, \mathbf{x}_n) + b \quad (3)$$

Satisfies the following conditions:

$$a_n \geq 0$$

$$t_n y(x_n) - 1 \geq 0$$

$$a_n \{t_n y(x_n) - 1\} = 0$$

for point with  $a_n = 0$  will not appear in the sum and plays no role making predictions for the new point.

The remaining data points are called support vectors because they satisfy  $t_n y(x_n) = 1$  and correspond to points that lie on the maximum margin hyperplanes in feature space.



# Prediction

Using Eq (3)

$$t_n \left( \sum_{m \in S} a_m t_m k(x_n, x_m) + b \right) = 1 \quad (4)$$

$S$  denotes the set of indices of the support vectors (SV).

Making a use of  $t_n^2 = 1$  averaging over all SV and solve for  $b$

$$b = \frac{1}{N_s} \left( \sum_{n \in S} \left( t_n - \sum_{m \in S} a_m t_m k(x_n, x_m) \right) \right) \quad (5)$$

$N_s$  is the total number of SV.

The maximum margin classifier in terms of the minimization of an error function

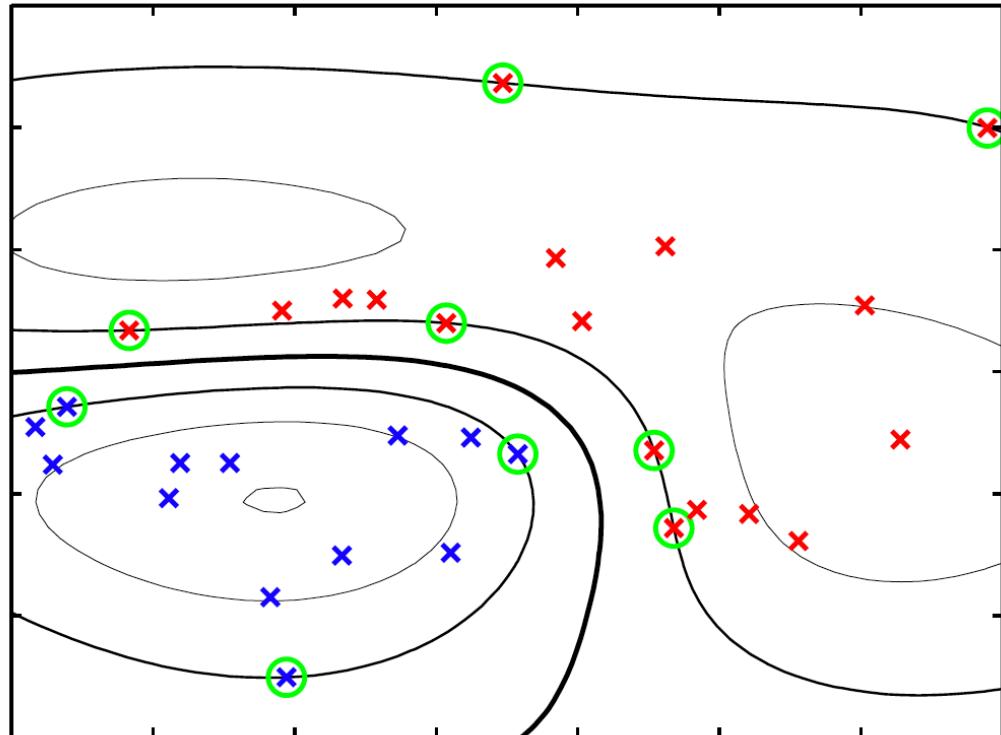
$$\sum_{n=1}^N E_\infty(y(x_n)t_n - 1) + \lambda \left\| w \right\|^2 \quad (6)$$

# Prediction

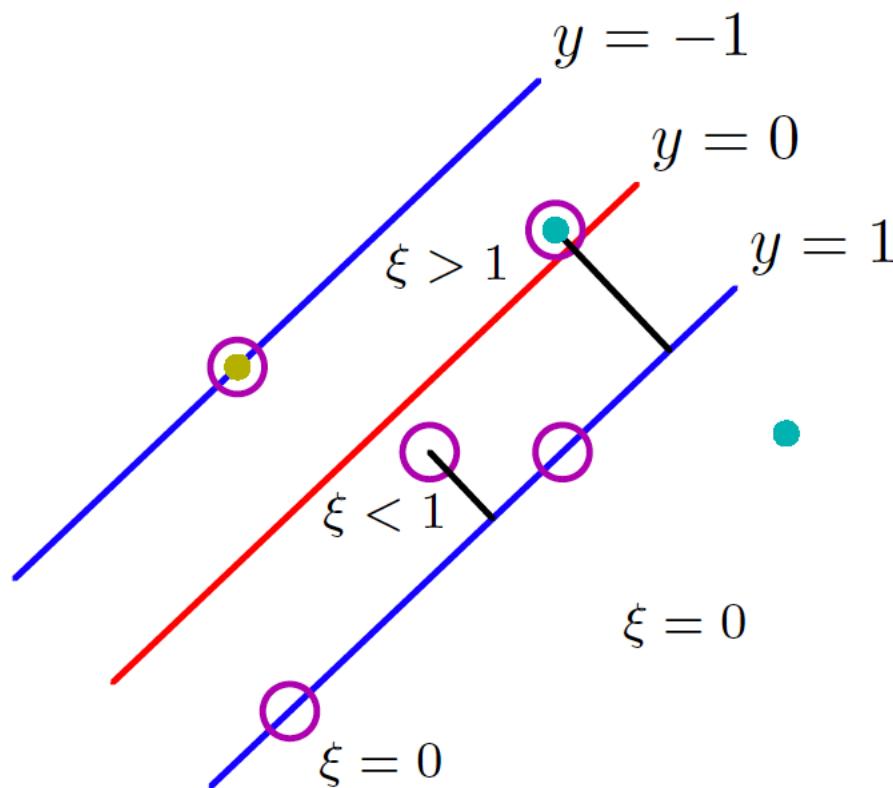
The maximum margin classifier in terms of the minimization of an error function

$$\sum_{n=1}^N E_\infty(y(x_n)t_n - 1) + \lambda \|w\|^2 \quad (6)$$

where  $E$  is a function that is zero if  $z \geq 0$  and  $\infty$  otherwise to ensure that the constraints satisfies.



# Overlapping Class Distributions



The class-conditional distributions may overlap, in which case exact separation of the training data can lead to poor generalization.

## Need to modify SVM

Modify so that data points are allowed to be on the wrong side of the margin boundary, but with penalty that increases with the distance from that boundary.

# Overlapping Class Distributions – soft SVM



Make the penalty a linear function of the distance.

slack variable  $\xi_n \geq 0$

$\xi = 0$  for data points that are on or inside the correct margin boundary and  $\xi = |t_n - y(x_n)|$  for other points.

On the decision boundary,  $y(x_n) = 0$  and have  $\xi = 1$  and with  $\xi > 1$  will be misclassified

$$t_n y(x_n) \geq 1 - \xi_n$$

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

# Overlapping Class Distributions – soft SVM



$$C \sum_{n=1}^N \xi_n + \frac{1}{2} \|w\|^2 \quad (7)$$

where parameter  $C > 0$  controls the trade-off between the slack variable penalty and the margin.

$$\begin{aligned} L(w, b, \xi, a, \mu) \\ = \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N \xi_n + \sum_{n=1}^N a_n \{t_n y_n(x_n) - 1 + \xi_n\} - \sum_{n=1}^N \mu_n \xi_n \end{aligned} \quad (8)$$

Lagrange Multipliers -  $\{a_n \geq 0\}$  and  $\{\mu_n \geq 0\}$ :

$$\begin{aligned} a_n = 0, t_n y_n(x_n) - 1 + \xi_n \geq 0, a_n \{t_n y_n(x_n) - 1 + \xi_n\} = 0 \\ \mu_n \geq 0, \xi_n \geq 0, \mu_n \xi_n = 0 \end{aligned}$$

# Overlapping Class Distributions – soft SVM



Optimize out  $w$ ,  $b$ , and  $\xi$

$$\frac{\partial L}{\partial w} = 0 \Rightarrow w = \sum_{n=1}^N a_n t_n \phi(x_n)$$

$$\frac{\partial L}{\partial b} = 0 \Rightarrow \sum_{n=1}^N a_n t_n = 0$$

$$\frac{\partial L}{\partial \xi_N} = 0 \Rightarrow a_n = C - \mu_n$$

$0 \leq a_n \leq C$  Box Constraints

$$L(a) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{m=1}^N \sum_{n=1}^N a_n a_m t_n t_m k(x_n, x_m) \quad (9)$$

# Overlapping Class Distributions – soft SVM



Solution

$$t_n \left( \sum_{m \in S} a_m t_m k(x_n, x_m) + b \right) = 1 \quad (10-1)$$

$$b = \frac{1}{N_M} \left( \sum_{n \in M} \left( t_n - \sum_{m \in S} a_m t_m k(x_n, x_m) \right) \right) \quad (10-2)$$

where  $M$  is the set of indices of data points having  $0 < a_n < C$ .

# Overlapping Class Distributions – soft SVM



IF we wish to use the SVM as a module in a larger probabilistic system

Platt (2000) has proposed fitting a logistic sigmoid to the outputs of a trained SVM.

$$p(t = 1|x) = \sigma(Ay(x) + B)$$

A & B are found by minimizing the cross-entropy error function defined by a training set consisting of pairs of values y and t.

- The data used to fit the sigmoid needs to be independent to avoid sever over-fitting.
- Give a poor approximation to the posterior prob.

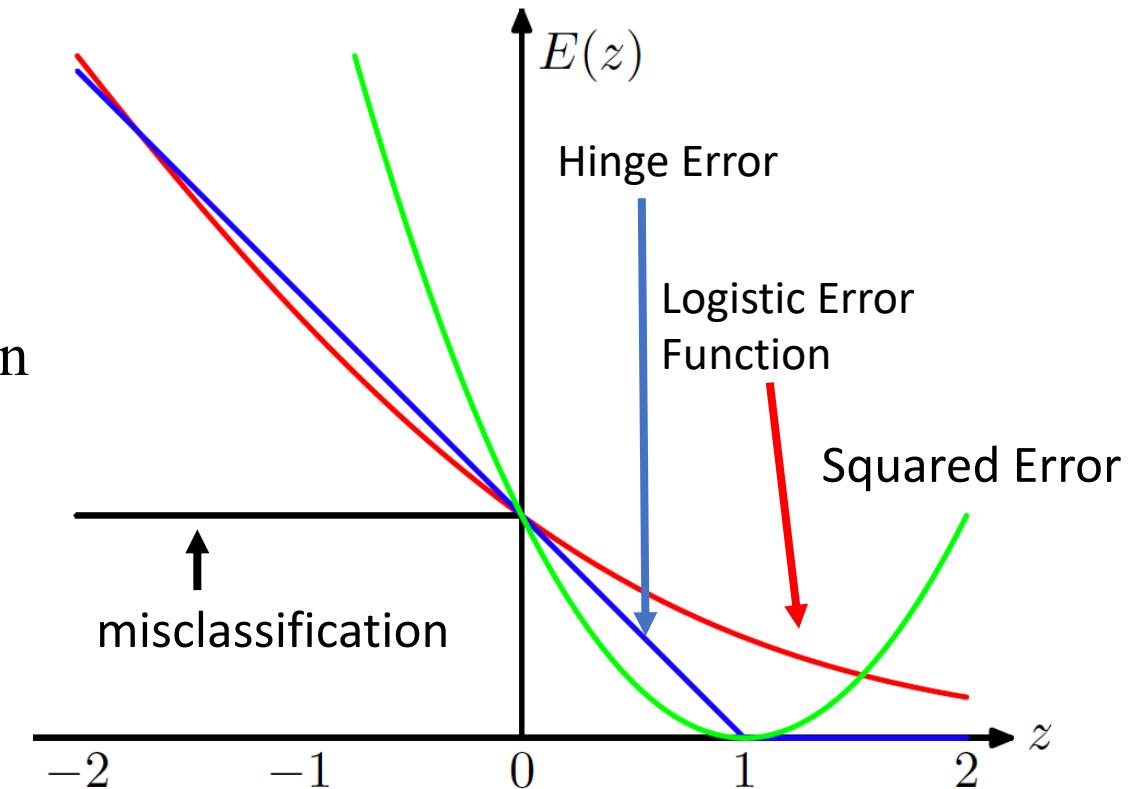
# Relation to logistic regression – Hinge Error

Update Eq (7)

$$\sum_{n=1}^N E_{SV}(y_n t_n) + \lambda \|w\|^2 \quad (11)$$

where  $\lambda = (2C)^{-1}$  and  $E_{SV}$  is the **hinge error** function

$$E_{SV}(y_n t_n) = [1 - y_n t_n]_+ \quad (12)$$





# Relation to logistic regression

	SVMs	Logistic Regression
<b>Loss Function</b>	Hinge Loss	Log-loss
<b>High dimensional features</b>	Yes	No
<b>Solution Sparse</b>	Yes	Almost no
<b>Output</b>	Margin	Real Probabilities!

# SVM for regression

$$\frac{1}{2} \sum_{n=1}^N \{y_n - t_n\}^2 + \frac{\lambda}{2} \|w\|^2 \quad (13)$$

$\epsilon$ -insensitive error function (Vapnik, 1995)

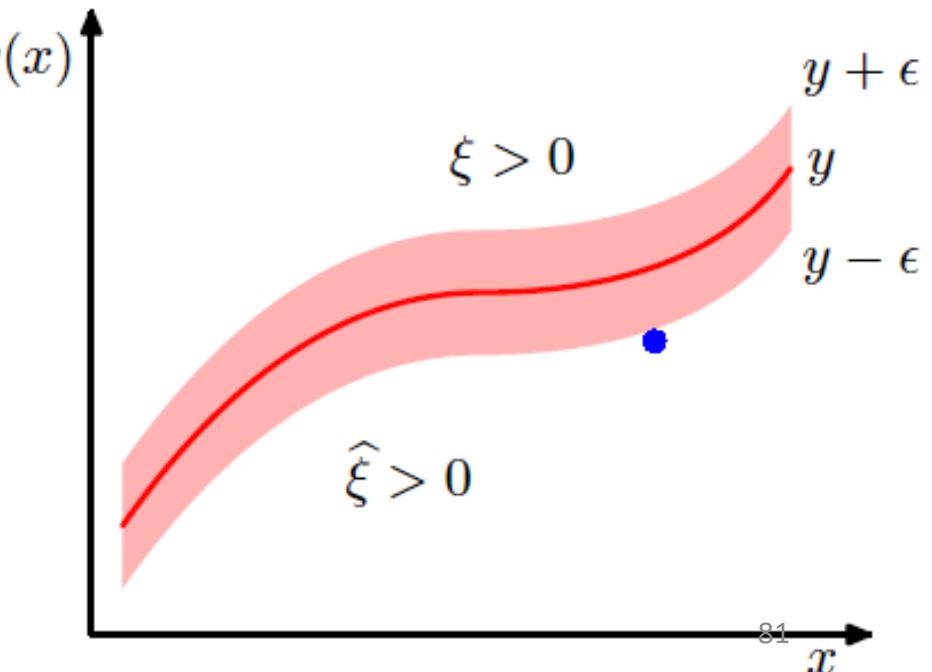
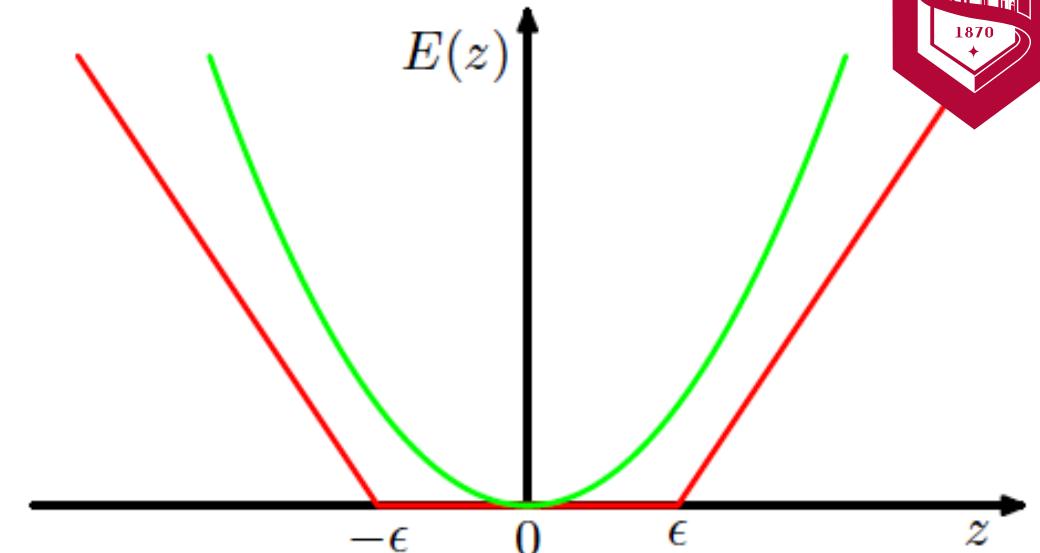
$$E_\epsilon(y(x) - t) = \begin{cases} 0 & \text{if } |y(x) - t| < \epsilon \\ |y(x) - t| & \text{Otherwise} \end{cases}$$

$$C \sum_{n=1}^N E_\epsilon(y_n(x_n) - t_n) + \frac{\lambda}{2} \|w\|^2 \quad (14)$$

$\xi_n > 0$  corresponds to a point  $t_n > y(x_n) + \epsilon$  and

$\hat{\xi}_n \geq 0$  for  $t_n < y(x_n) - \epsilon$

$$C \sum_{n=1}^N (\xi_n + \hat{\xi}_n) + \frac{1}{2} \|w\|^2 \quad (15)$$



# SVM for regression - Optimize by slack variable



$$\begin{aligned} L \\ = C \sum_{n=1}^N (\xi_n + \hat{\xi}_n) + \frac{1}{2} \|w\|^2 - \sum_{n=1}^N (\mu_n \xi_n + \hat{\mu}_n \hat{\xi}_n) - \sum_{n=1}^N a_n (\epsilon + \xi_n + y_n - t_n) \\ - \sum_{n=1}^N \hat{a}_n (\epsilon + \hat{\xi}_n - y_n + t_n) \end{aligned} \tag{16}$$

The derivative w.r.t.  $w$ ,  $b$ ,  $\xi_n$ , and  $\hat{\xi}_n$

$$\frac{\partial L}{\partial w} = 0 \Rightarrow w = \sum_{n=1}^N (a_n - \hat{a}_n) \phi(x_n), \frac{\partial L}{\partial b} = 0 \Rightarrow \sum_{n=1}^N (a_n - \hat{a}_n) = 0$$

$$\frac{\partial L}{\partial \xi_n} = 0 \Rightarrow a_n = C - \mu_n, \frac{\partial L}{\partial \hat{\xi}_n} = 0 \Rightarrow \hat{a}_n = C - \hat{\mu}_n$$

# SVM for regression - Optimize by slack variable



$$\begin{aligned}
 L(a, \hat{a}) &= -\frac{1}{2} \sum_{m=1}^N \sum_{n=1}^N (a_n - \hat{a}_n)(a_m - \hat{a}_m) k(x_n, x_m) - \epsilon \sum_{n=1}^N (a_n + \hat{a}_n) + \sum_{n=1}^N (a_n - \hat{a}_n) t_n
 \end{aligned} \tag{17}$$

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

$$\begin{aligned}
 a_n(\epsilon + \xi_n + y_n - t_n) &= 0 \\
 \hat{a}_n(\epsilon + \hat{\xi}_n - y_n + t_n) &= 0 \\
 (C - a_n)\xi_n &= 0 \\
 (C - \hat{a}_n)\hat{\xi}_n &= 0
 \end{aligned}$$

$$b = t_n - \epsilon - w^T \phi(x_n) = t_n - \epsilon - \sum_{m=1}^N (a_m - \hat{a}_m) k(x_n, x_m) \tag{18}$$

$$L(a, \hat{a}) = -\frac{1}{2} \sum_{m=1}^N \sum_{n=1}^N (a_n - \hat{a}_n)(a_m - \hat{a}_m) k(x_n, x_m) + \sum_{n=1}^N (a_n - \hat{a}_n) t_n \tag{19}$$



# Kernel SVM

- Outputs represents decisions (can be an advantage but also can be a disadvantage)
- Originally formulated for two classes
- Predictions are linear combination of kernel functions **but centered on training data points and that are required to be positive-definite.**
- How can we overcome the limitation of centered kernel?
- The relevance vector machine or RVM (Tipping 2001)
  - Faster performance on test data while the generalization error is maintained.



# RVM - Regression

The modified prior:

$$p(t|x, w, \beta) = \mathcal{N}(t|y(x), \beta^{-1}) \quad (20)$$

where  $\beta = \sigma^{-2}$  is the noise precision.

The mean of the linear model

$$y(x) = \sum_{i=1}^M w_i \phi_i(x) = w^T \phi(x)$$

In RVM, we use kernels with one kernel associated with each of the data points in training set.

$$y(x) = \sum_{i=1}^N w_n k(x, x_n) + b \quad (21)$$

where  $b$  is the bias parameter and  $M = N + 1$ .

Here, **there is no restriction of positive-definite and the number nor the location to the training data sets.**



# RVM - Regression

Suppose we have N observations of data matrix X. The likelihood function is

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N p(t_n|x_n, \mathbf{w}, \beta)$$

and a weight prior distribution is

$$p(\mathbf{w}|\boldsymbol{\alpha}) = \prod_{n=1}^M \mathcal{N}(w_i|0, \alpha_i^{-1})$$

$\alpha_i$  is a separate hyperparameter  
for each of weight parameter.

If the posterior distribution for the weights is Gaussian, then it is

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \mathcal{N}(\mathbf{w}|\mathbf{m}, \Sigma)$$

Mean:  $\beta \Sigma \Phi^T \mathbf{t} = \beta \Sigma \mathbf{K}^T \mathbf{t}$

Covariance:  $\Sigma = (\mathbf{A} + \beta \Phi^T \Phi)^{-1}$   
 $= (\mathbf{A} + \beta \mathbf{K}^T \mathbf{K})^{-1}$

The maximization of likelihood function by integration of the weight parameters

$$p(\mathbf{t}|\mathbf{X}, \boldsymbol{\alpha}, \beta) = \int p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta)p(\mathbf{w}|\boldsymbol{\alpha})d\mathbf{w} \quad (22)$$



# RVM - Regression

If the posterior distribution for the weights is Gaussian, then it is

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \mathcal{N}(\mathbf{w}|\mathbf{m}, \Sigma)$$

The maximization of likelihood function by integration of the weight parameters

$$p(\mathbf{t}|\mathbf{X}, \boldsymbol{\alpha}, \beta) = \int p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta)p(\mathbf{w}|\boldsymbol{\alpha})d\mathbf{w} \quad (22)$$

The log likelihood also can be formatted as

$$\ln p(\mathbf{t}|\mathbf{X}, \boldsymbol{\alpha}, \beta) = \ln \mathcal{N}(\mathbf{t}|\mathbf{0}, \mathbf{C}) - \frac{1}{2}\{N \ln(2\pi) + \ln|\mathbf{C}| + \mathbf{t}^T \mathbf{C}^{-1} \mathbf{t}\} \quad (23)$$

where  $\mathbf{C} = \beta^{-1} \mathbf{I} + \Phi \mathbf{A}^{-1} \Phi^T$ .

Set the required derivatives of the likelihood to zero and find the hyperparameters  $\boldsymbol{\alpha}$  and  $\beta$ :

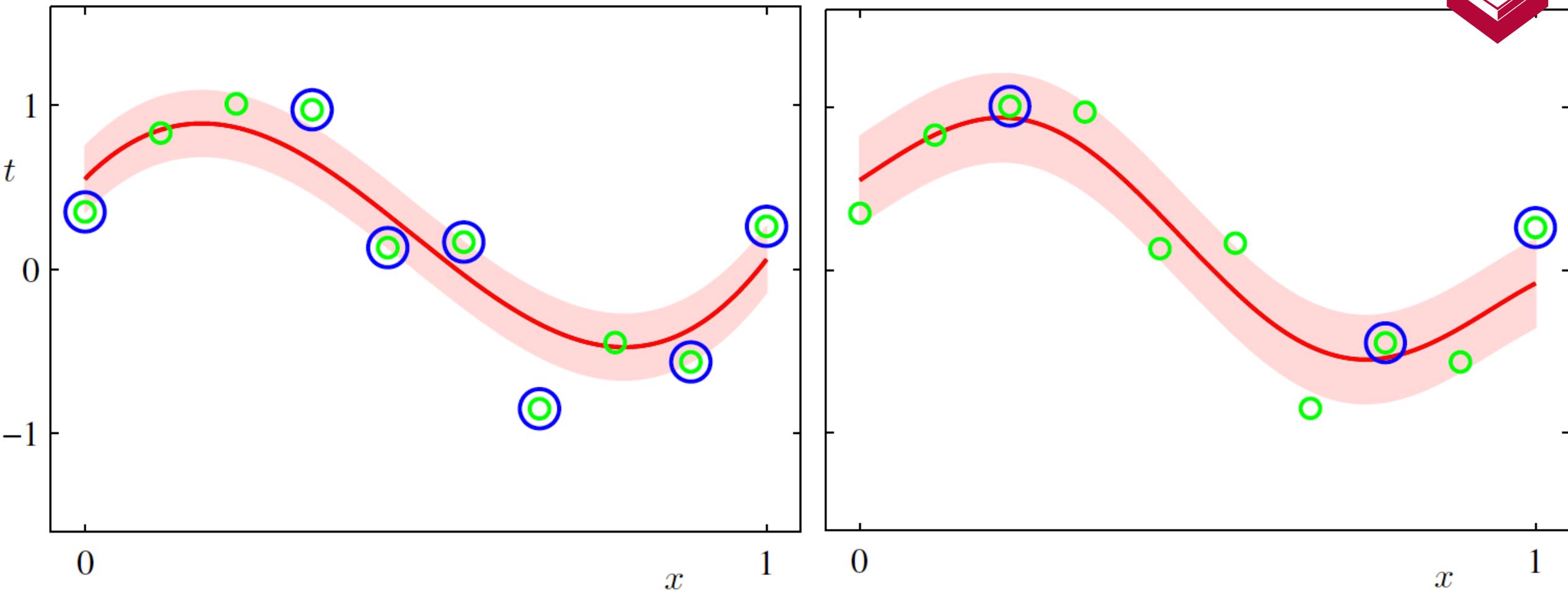
$$\alpha_i^{new} = \frac{\gamma_i}{m_i^2}$$

$$(\beta^{new})^{-1} = \frac{\|\mathbf{t} - \Phi \mathbf{m}\|^2}{N - \sum_i \gamma_i}$$

where  $m_i$  is the  $i^{th}$  component of the poster mean  $\mathbf{m}$  and the quantity  $\gamma_i = 1 - \alpha_i \Sigma_{ii}$  is the measurement of how well  $w_i$  is measured.



# RVM - Regression





# RVM - Regression

Linear Regression – the predictive variance becomes small in regions of input space.

SVM Regression -  $\phi$  is centered on data points and the model will become increasingly certain of its predictions when extrapolating outside the domain of the data.

Gaussian Process – the computation cost is high.

RVM – A significant improvement in the speed of processing on test data. This greater sparsity is achieved with little or no reduction in generalization error compared with SVM.

But... Training involves optimizing a nonconvex function and takes longer.



# RVM - Classification

Let's begin from the two-class problem with a binary target  $t \in \{0,1\}$ .

If the model takes the linear combination format of transformation function  $\phi(\cdot)$  by a logistic sigmoid function

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

The posterior distribution over  $\mathbf{w}$  (assume it is a Gaussian prior) is

$$p(\mathbf{w}|t) \propto p(\mathbf{w})p(t|\mathbf{w}).$$

We can solve for  $\alpha$  by following the way we saw in RVM for regression using Hessian matrix.



# RVM - Classification

$$\begin{aligned}\ln p(\mathbf{w}|\mathbf{t}, \boldsymbol{\alpha}) &= \ln\{p(\mathbf{t}|\mathbf{w})p(\mathbf{w}|\boldsymbol{\alpha})\} - \ln p(\mathbf{t}|\boldsymbol{\alpha}) \\ &= \sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\} - \frac{1}{2} \mathbf{w}^T \mathbf{A} \mathbf{w} + \text{const.}\end{aligned}$$

From logistic regression  
lecture,  
Refer to equation 22.

Regularization comes from the  
prior in the previous slide,  
 $p(\mathbf{w}|\mathbf{t}) \propto p(\mathbf{w})p(\mathbf{t}|\mathbf{w})$

Note:  $\ln p(\mathbf{w}|\mathbf{t}) = -\frac{1}{2}(\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1}(\mathbf{w} - \mathbf{m}_0) + \sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\} + \text{Const.}$  when the Gaussian prior of  $\mathbf{w}$  is  $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0)$

From here, we can take the gradient of  $\ln p(\mathbf{w}|\mathbf{t}, \boldsymbol{\alpha})$  be zero to obtain the mean and  $\alpha_i^{new}$  is

$$\alpha_i^{new} = \frac{\gamma_i}{m_i^2}$$