

# CS 559: Machine Learning Fundamentals & Applications

## Lecture 1: Linear Algebra and Probability Theory

In Jang: [ijang@stevens.edu](mailto:ijang@stevens.edu)

Fall 2021





# Outline

- Course Information
- Python Review I – Starting Data Science with Python (see demonstration file)
- Python Review II – Visualization (see demonstration file)
- Linear Algebra Review
- Probability Theory



# Course Information

- Lectures:
  - Lecture Video will be posted on Thursday Night
  - Office Hours – Friday 10 – 11 AM or by appointment
- Tech Requirement:
  - Python (v2.7 or higher)
  - Jupyter Notebook: <https://jupyter.org/install>



# Python Review

- This review focuses on basic knowledge of python you must have to do Machine Learning and Data Science.
- Jupyter Notebook is a user-friendly computational environment for data scientists.
  - Stores work automatically
  - Markdown allows to export work in pdf, html, etc. for presentations, notes, etc.
- Why python?
  - It distinguishes itself from other languages by its flexibility, simplicity, and reliable tools.
  - It is consistent and is anchored on simplicity which makes it most appropriate for machine learning.
  - Fits the best on machine learning due to its independent platform and its popularity in the programming community.
  - We have other languages to do machine learning projects such as R and MATLAB.
    - R and MATLAB are more specialized in specific fields (statistics, mathematics, and bioinformatics)
    - Python is much convenient to do data mining and text analysis.



# Linear Algebra Review

- Scalars, Vectors, and Matrices
- Adding Matrices and Vectors
- Multiplying Matrices and Vectors
- Identity and Inverse Matrices
- Linear Dependence and Span
- Norms
- Special Matrices and Vectors
- Matrix Representation
- Linear Operators
- Trace
- Determinant
- Eigenvalues and Eigenvectors
- Subspace



# Scalars, Vectors, and Matrices

- **Scalars:** a single number. ( $x$ )
  - Let  $s \in \mathbb{R}$  be the slope of the line: real-valued scalar
  - Let  $n \in \mathbb{N}$  be the number of units: natural number scalar
- **Vectors:** an array of numbers that are arranged in order. ( $\mathbf{x}$ )
  - Elements of vectors: typically written as  $x_i$  where  $x$  is scalar and the subscript  $i$  is the order.
  - If each element is in  $\mathbb{R}$  and the vector has  $n$  elements, then the vector lies in  $\mathbb{R}^n$
  - $\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$  or  $\mathbf{x} = [x_1 \quad \cdots \quad x_n]$
- **Matrix:** a 2-D array of numbers. ( $\mathbf{A}$ )
  - $\mathbf{A} \in \mathbb{R}^{m \times n}$ : A real-valued matrix  $\mathbf{A}$  has a height of  $m$  and a width of  $n$
  - $\mathbf{A} = \begin{bmatrix} A_{1,1} & \cdots & A_{1,n} \\ \vdots & \ddots & \vdots \\ A_{m,1} & \cdots & A_{m,n} \end{bmatrix}$
  - Vectors can be considered as matrices that contain only one column.



# Scalars, Vectors, and Matrices

- **Transpose:** The mirror image of the matrix across a diagonal line, called the main diagonal, running down and to the right, starting from its upper left corner.

$$\bullet \begin{bmatrix} a_{1,1} & a_{1,2} & a_{1,3} \\ a_{2,1} & a_{2,2} & a_{2,3} \\ a_{3,1} & a_{3,2} & a_{3,3} \end{bmatrix}^T = \begin{bmatrix} a_{1,1} & a_{2,1} & a_{3,1} \\ a_{1,2} & a_{2,2} & a_{3,2} \\ a_{1,3} & a_{2,3} & a_{3,3} \end{bmatrix}$$

$$\bullet A_{i,j}^T = A_{j,i}$$

$$\bullet \text{If } \mathbf{a} = \begin{bmatrix} a_1 \\ \vdots \\ a_n \end{bmatrix}, \text{ then its transpose } \mathbf{a}^T \text{ is } \mathbf{a}^T = [a_1 \quad \cdots \quad a_n]$$

$$A \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix} \quad A^T \begin{bmatrix} 1 & 4 & 7 \\ 2 & 5 & 8 \\ 3 & 6 & 9 \end{bmatrix}$$

$$A \begin{bmatrix} 1 & 4 & 3 \\ 8 & 2 & 6 \\ 7 & 8 & 3 \\ 4 & 9 & 6 \\ 7 & 8 & 1 \end{bmatrix} \quad A^T \begin{bmatrix} 1 & 8 & 7 & 4 & 7 \\ 4 & 2 & 8 & 9 & 8 \\ 3 & 6 & 3 & 6 & 1 \end{bmatrix}$$



# Adding Matrices and Vectors

- Addition between matrices?
  - As long as the shapes are the same,
    - $C = A + B$  where  $C_{i,j} = A_{i,j} + B_{i,j}$
- A scalar to Matrix?
  - Perform the operations to each element
    - $D = a \cdot B + c$  where  $D_{i,j} = aB_{i,j} + c$
- A vector to Matrix?
  - Yields another matrix
    - $C = A + b$  where  $C_{i,j} = A_{i,j} + b_j$  ←
    - Copying of **b** to many locations is called **broadcasting**.

The vector **b** is added to each row of the matrix.



# Multiplying Matrices and Vectors

- Multiplication (**Matrix product**) between matrices?
  - **A** must have the same number of columns as **B** has rows
    - Let **A** be  $m \times n$  and **B** be  $n \times p$ , then **C** =  $\mathbf{AB}$  is a matrix with shape of  $m \times p$
    - $\mathbf{C} = \mathbf{AB}$  where  $C_{i,j} = \sum_k A_{i,k}B_{k,j}$
    - The product of individual elements is called **element-wise product**, or **Hadamard product** and is denoted as  $\mathbf{A} \odot \mathbf{B}$ .
  - The dot product between two vectors  $x$  and  $y$ 
    - $x$  and  $y$  must have the same dimension
    - Also can be considered as matrix product with the same dimensionality
    - $x \cdot y = x^T y$



# Multiplying Matrices and Vectors

$$\begin{bmatrix} -4 & 4 & 6 \\ 2 & 3 & -1 \end{bmatrix} \begin{bmatrix} 0 & -4 & -2 & 5 \\ 1 & 6 & -1 & 7 \\ 8 & 2 & 4 & 3 \end{bmatrix}$$

$$= \begin{bmatrix} -4 \times 0 + 4 \times 1 + 6 \times 8 & -4 \times -4 + 4 \times 6 + 6 \times 2 & -4 \times -2 + 4 \times -1 + 6 \times 4 & -4 \times 5 + 4 \times 7 + 6 \times 3 \\ 2 \times 0 + 3 \times 1 + -1 \times 8 & 2 \times -4 + 3 \times 6 + -1 \times 2 & 2 \times -2 + 3 \times -1 + -1 \times 4 & 2 \times 5 + 3 \times 7 + -1 \times 3 \end{bmatrix}$$

$$= \begin{bmatrix} 52 & 52 & 28 & 26 \\ -5 & 8 & -11 & 28 \end{bmatrix}$$



# Multiplying Matrices and Vectors

- **Distributive property:**
  - $A(B + C) = AB + AC$
- **Associative property:**
  - $A(BC) = (AB)C$
- **But not commutative!**
  - $AB \neq BA$
  - but in dot product between two vectors, we have  $\mathbf{x}^T \mathbf{y} = \mathbf{y}^T \mathbf{x}$ .



# Identity and Inverse Matrices

Matrix inversion: Enables to analytically solve  $\mathbf{A}\mathbf{x} = \mathbf{b}$  for many values of  $\mathbf{A}$ .

- An identity matrix: a matrix that does not change any vector when we multiply that vector by that matrix. It is denoted as  $\mathbf{I}_n$  that is  $\in \mathbb{R}^{n \times n}$  and for  $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{I}_n\mathbf{x} = \mathbf{x}$ .
  - All the entries along the main diagonal are 1 and other entries are zero.
- The matrix inverse of  $\mathbf{A}$  is denoted as  $\mathbf{A}^{-1}$  and it is defined as the matrix such that

$$\mathbf{A}^{-1}\mathbf{A} = \mathbf{I}_n$$

If we have  $\mathbf{A}\mathbf{x} = \mathbf{b}$ ,  $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$ .

- $\mathbf{A}^{-1}$  is primarily useful as a theoretical tool but *should not be used* in practice for most software applications.
- $\mathbf{A}^{-1}$  can be represented with only limited precision on a digital computer and algorithms that make use of the value of  $\mathbf{b}$  can usually obtain more accurate estimates of  $\mathbf{x}$ .



# The eight axioms

|  |  |
|--|--|
| <b>Associativity of addition</b>   | $\mathbf{u} + (\mathbf{v} + \mathbf{w}) = (\mathbf{u} + \mathbf{v}) + \mathbf{w}$  |
| <b>Commutativity of addition</b>   | $\mathbf{u} + \mathbf{v} = \mathbf{v} + \mathbf{u}$  |
| <b>Identity element of addition</b>  | There exists an element $0 \in V$ , called the zero vector, such that $\mathbf{v} + 0 = \mathbf{v}$ for all $\mathbf{v} \in V$ .   |
| <b>Inverse elements of addition</b>  | For every $\mathbf{v} \in V$ , there exists an element $-\mathbf{v} \in V$ , called the additive <i>inverse</i> of $\mathbf{v}$ , such that $\mathbf{v} + (-\mathbf{v}) = 0$ . |
| <b>Compatibility of scalar multiplication with field multiplication</b>        | $a(b\mathbf{v}) = (ab)\mathbf{v}$  |
| <b>Identity element of scalar multiplication</b>                               | $1\mathbf{v} = \mathbf{v}$ , where $1$ denotes the multiplicative identity in $F$ .  |
| <b>Distributivity of scalar multiplication with respect to vector addition</b> | $a(\mathbf{u} + \mathbf{v}) = a\mathbf{u} + a\mathbf{v}$   |
| <b>Distributivity of scalar multiplication with respect to field addition</b>  | $(a + b)\mathbf{v} = a\mathbf{v} + b\mathbf{v}$  |



# Linear Dependence and Span

Recall  $Ax = b \rightarrow x = A^{-1}b$

- One solution for every value of b but also can have no solutions or infinitely many solutions
- It is not possible to have more than one but less than infinitely many solutions for particular b.
- If both x and y are solutions, then

$$z = \alpha x + (1 - \alpha)y$$



# Linear Dependence and Span

The columns of  $\mathbf{A}$  specify directions from the origin

- Determines *how many ways* there are of reaching  $\mathbf{b}$ .
- **Linear combination**

$$\mathbf{Ax} = \sum_i x_i A_{:,i}$$

- **Span** – set of all points obtained by linear combination of the original vectors.
- Determining whether  $\mathbf{Ax} = \mathbf{b}$  has a solution is testing whether  $\mathbf{b}$  is in the *span of the columns* of  $\mathbf{A}$ .
  - The column space of  $\mathbf{A}$  be all of  $\mathbb{R}^m$  for the system to have solution for all values of  $\mathbf{b} \in \mathbb{R}^m$ !
  - $\mathbf{A}$  must have at least  $m$  columns!



# Linear Dependence and Span

Consider A have columns n and b has m elements where  $n \geq m$ ,

- If  $A$  is a  $3 \times 2$  matrix and  $x$  is  $2 \times 1$ , the target  $b$  is a  $3 \times 1$ .
- The target  $b$  is in 3-D but  $x$  is only 2-D.
- In this case, the solution is only when  $b$  lies on that plane.
- What happens if columns are identical?
  - Will fail to encompass all of  $\mathbb{R}^m$
  - **Linear independent** – if no vector in the set is a linear combination of the other vectors.
  - For the column space of the matrix to encompass all of  $\mathbb{R}^m$ , the matrix must contain at least one set of  $m$  linearly independent columns.
  - This requirement is to have exactly  $m$  linear independent columns, not at least  $m$ .
  - No set of  $m$ -D vectors can have more than  $m$  mutually independent columns, but a matrix with more than  $m$  columns may have more than one such set.



# Linear Dependence and Span

- $Ax = b \rightarrow x = A^{-1}b$

For A to have an inverse:

- Linear independency
- Must have at most one solution for each value of b  $\rightarrow$  at most m columns
  - Otherwise, there is more than one way of parameterizing each solution.
- In other words... A must be **square** matrix by being  $m \times m$  and all columns be **linearly independent**.
- A square matrix with linearly dependent columns is known as **singular**.
- For non-square or a square but singular A?
  - **We can solve using other than the method of matrix inversion.**



# Norms

Norm – In ML, we measure the size of vectors by mapping vectors to non-negative values.

$$\|x\|_p = \left( \sum_i |x_i|^p \right)^{\frac{1}{p}}$$

for  $p \in \mathbb{R}, p \geq 1$ .

Also the norm of vector  $x$  measures the distance from the origin to the point  $x$ .



# Matrix Representation

- Let  $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$  be a basis
- Every  $\mathbf{v}$  can be uniquely represented as:

$$\mathbf{v} = \alpha_1 \mathbf{v}_1 + \cdots + \alpha_k \mathbf{v}_k$$

- Denote  $\mathbf{v}$  by the column-vector:  $\mathbf{v} = \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_k \end{pmatrix}$
- Denote the basis vectors as:  $\begin{pmatrix} 1 \\ \mathbf{0} \\ \vdots \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \mathbf{0} \\ 1 \\ \vdots \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \\ \vdots \\ 1 \end{pmatrix}$

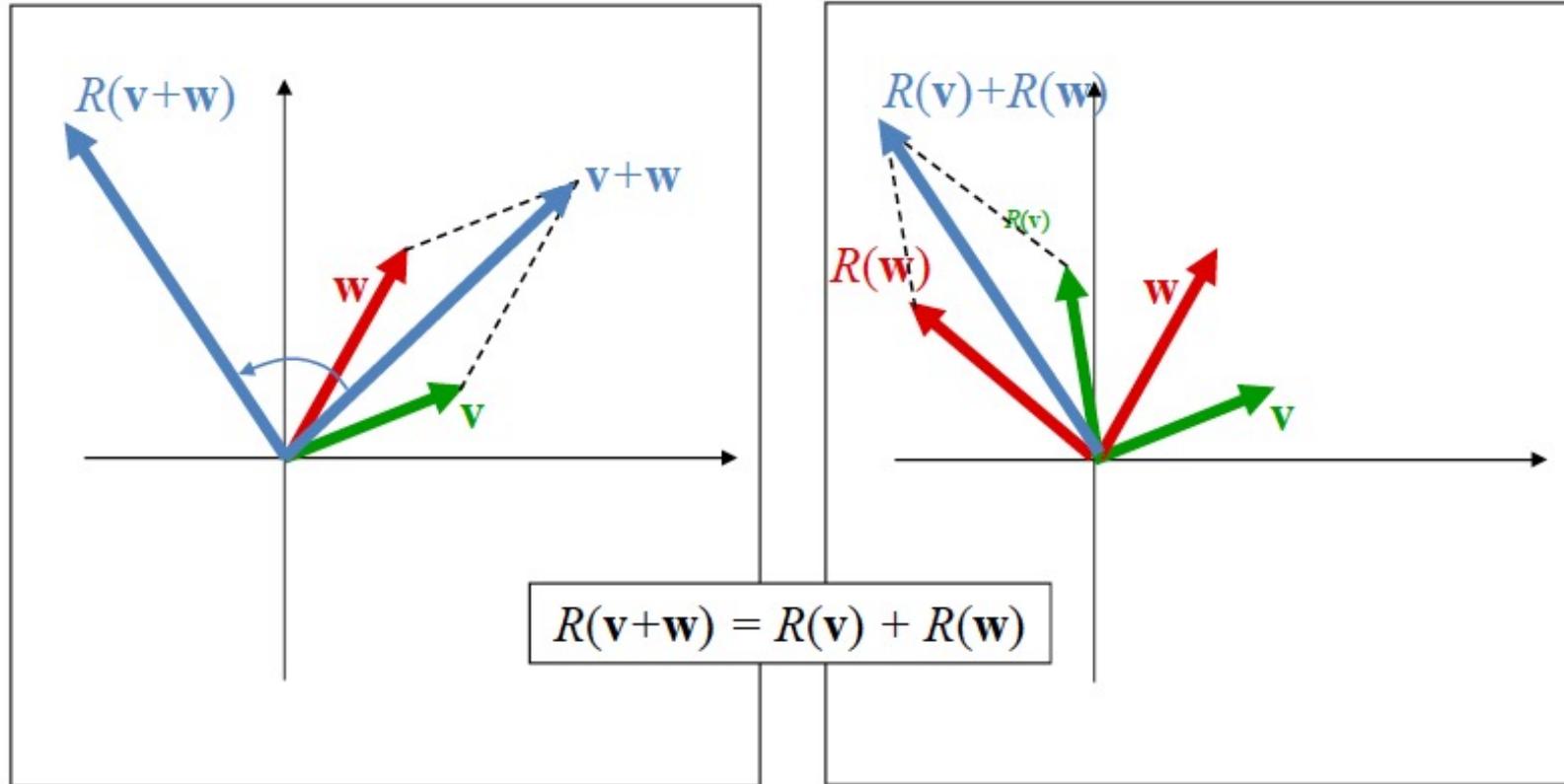


# Linear Operators

- $A : V \rightarrow W$  is called linear operator if:
  - $A(\mathbf{v} + \mathbf{w}) = A(\mathbf{v}) + A(\mathbf{w})$
  - $A(\alpha\mathbf{v}) = \alpha A(\mathbf{v})$
- In particular,  $A(\mathbf{0}) = \mathbf{0}$
- Are the following operators linear?
  - Scaling
  - Rotation
  - Translation

# Linear operators - illustration

- Rotation is a linear operator:





# Matrix properties

- Matrix  $A$  ( $n \times n$ ) is non-singular if  $\exists B, AB = BA = I$
- $B = A^{-1}$  is called the inverse of  $A$
- $A$  is non-singular  $\Leftrightarrow \det A \neq 0$
- If  $A$  is non-singular then the equation
  - $A\mathbf{x} = \mathbf{b}$  has one unique solution for each  $\mathbf{b}$
  - the rows of  $A$  are linearly independent (and so are the columns)

# Orthogonal matrices

- Matrix  $A_{(n \times n)}$  is orthogonal if  $A^{-1} = A^T$
- Follows:  $AA^T = A^TA = I$
- The rows of  $A$  are orthonormal vectors!

Proof:

$$I = A^T A = \left( \begin{array}{c} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \vdots \\ \mathbf{v}_n \end{array} \right) \left( \begin{array}{cccc} \mathbf{v}_1 & \mathbf{v}_2 & \cdots & \mathbf{v}_n \end{array} \right) = \left( \begin{array}{c} \mathbf{v}_i^T \mathbf{v}_j \end{array} \right) = \left( \begin{array}{c} \delta_{ij} \end{array} \right)$$

$$\Rightarrow \langle \mathbf{v}_i, \mathbf{v}_i \rangle = 1 \quad \Rightarrow \quad \|\mathbf{v}_i\| = 1; \quad \langle \mathbf{v}_i, \mathbf{v}_j \rangle = 0$$

Delta function  $\delta_{ij}$ :

$$\delta_{ij} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{else} \end{cases}$$



# Trace

- The trace of a square matrix denoted by  $\text{tr}(A)$  is sum of the diagonal elements

$$\text{tr}(A) = \sum_{i=1}^N A_{ii}$$

- $\text{tr} \begin{pmatrix} 1 & 2 \\ 3 & 1 \end{pmatrix} = A_{11} + A_{22} = 1 + 1 = 2$



# Determinant

- For a square matrix  $A$ , the determinant is denoted by  $|A|$  or  $\det(A)$

$$\begin{aligned} |[a_{11}]| &= a_{11} \\ \left| \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \right| &= a_{11}a_{22} - a_{12}a_{21} \\ \left| \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} \right| &= a_{11}a_{22}a_{33} + a_{12}a_{23}a_{31} + a_{13}a_{21}a_{32} \\ &\quad - a_{11}a_{23}a_{32} - a_{12}a_{21}a_{33} - a_{13}a_{22}a_{31} \end{aligned}$$



# Determinant

- $|A| = |A^T|$
- $|AB| = |A| |B|$
- $|A| = 0$ , if and only if  $A$  is singular
  - Else,  $|A^{-1}| = 1/|A|$



# Eigenvalues and Eigenvectors

- For an  $n \times n$  square matrix  $A$ ,  $\mu$  is an eigenvector with eigenvalue  $\lambda$  if

$$A\mu = \lambda\mu$$

- that is

$$(A - \lambda I)\mu = 0$$

- If  $(A - \lambda I)$  is invertible, then the only solution is  $\mu = 0$ . (trivial)
- For non-trivial solutions:

$$\det(A - \lambda I) = 0$$

- Above equation is called the “characteristic polynomial”
- Solutions are not unique:
  - If  $\mu$  is an eigenvector  $\alpha\mu$  is also an eigenvector where  $\alpha$  is any constant.



# Eigenvalues and Eigenvectors

- For a  $2 \times 2$  matrix

$$\det[\mathbf{A} - \lambda \mathbf{I}] = \begin{vmatrix} a_{11} - \lambda & a_{12} \\ a_{21} & a_{22} - \lambda \end{vmatrix} = (a_{11} - \lambda)(a_{22} - \lambda) - a_{12}a_{21} = 0$$

$$0 = a_{11}a_{22} - a_{12}a_{21} - \lambda(a_{11} + a_{22}) + \lambda^2$$

If  $\mathbf{A} = \begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix}$ , we get

$$\begin{aligned} 0 &= 1 \cdot 4 - 2 \cdot 2 - (1+4)\lambda + \lambda^2 \\ (1+4)\lambda &= \lambda^2 \end{aligned}$$

- The solutions are  $\lambda=0$  and  $\lambda=5$ .



# Eigenvalues and Eigenvectors

- The eigenvector for the first eigenvalue,  $\lambda=0$  is:

$$\mathbf{Ax} = \lambda \mathbf{x}, \quad (\mathbf{A} - \lambda \mathbf{I})\mathbf{x} = 0$$

$$\left[ \begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix} - \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \right] \cdot \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 1x + 2y \\ 2x + 4y \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

- One solution for both equations is  $x = 2, y = -1$ .
- The second eigenvalue is  $\lambda=5$

$$\left[ \begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix} - \begin{bmatrix} 5 & 0 \\ 0 & 5 \end{bmatrix} \right] \cdot \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} -4 & 2 \\ 2 & -1 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} -4x + 2y \\ 2x - 1y \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

- A solution is  $x = 1, y = 2$ .



# Properties

- The product of the eigenvalues =  $|A|$ 
  - $\lambda_1\lambda_2 = 0 \cdot 5 = 0 = (1 \cdot 4 - 2 \cdot 2) = 0$
- The sum of the eigenvalues =  $\text{trace}(A)$ 
  - $\lambda_1 + \lambda_2 = 0 + 5 = 5 = A_{11} + A_{22} = 1 + 4 = 5$
- The eigenvectors are pairwise orthogonal
  - $\begin{bmatrix} 2 \\ -1 \end{bmatrix}^T \times \begin{bmatrix} 1 \\ 2 \end{bmatrix} = (2 \cdot 1 - 1 \cdot 2) = 0$

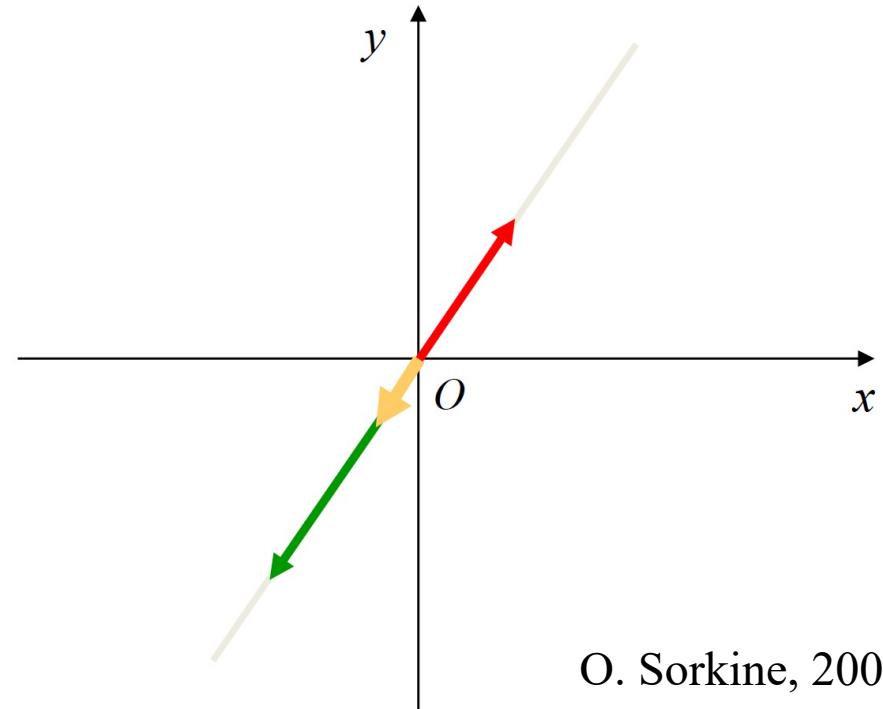


# Subspace

- Let  $F$  be a field,  $V$  be a vector space over  $F$ , and let  $W$  be a subset of  $V$ . Then  $W$  is a **subspace** if:
  - The zero vector,  $\mathbf{0}$ , is in  $W$ .
  - If  $\mathbf{u}$  and  $\mathbf{v}$  are elements of  $W$ , then the sum  $\mathbf{u} + \mathbf{v}$  is an element of  $W$ .
  - If  $\mathbf{u}$  is an element of  $W$  and  $c$  is a scalar from  $K$ , then the scalar product  $c\mathbf{u}$  is an element of  $W$ .

# Subspace Example

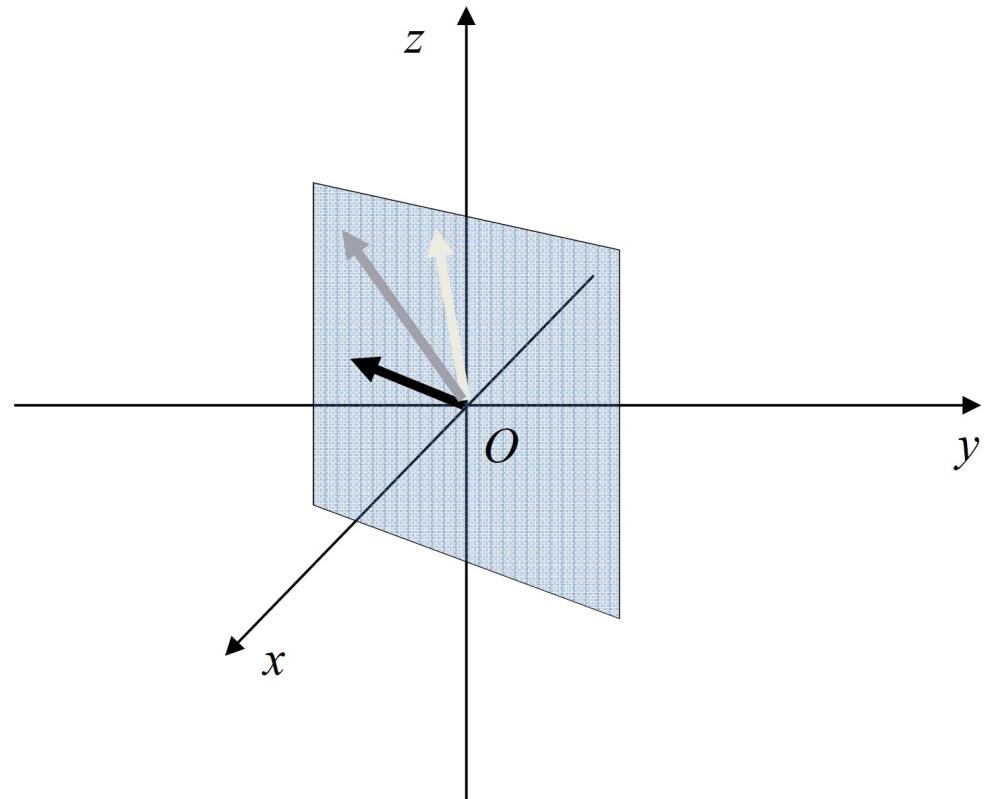
- Let  $l$  be a 2D line though the origin
- $L = \{\mathbf{p} - O \mid \mathbf{p} \in l\}$  is a linear subspace of  $R^2$



O. Sorkine, 2006

# Subspace Example

- Let  $\pi$  be a plane through the origin in 3D
- $V = \{\mathbf{p} - O \mid \mathbf{p} \in \pi\}$  is a linear subspace of  $R^3$





# Linear Independence

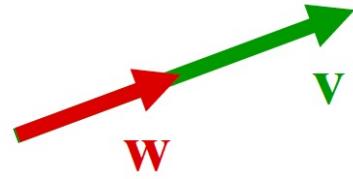
- The vectors  $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$  are a linearly independent set if:

$$\alpha_1 \mathbf{v}_1 + \cdots + \alpha_k \mathbf{v}_k = \mathbf{0} \Leftrightarrow \alpha_i = 0 \ \forall i$$

- It means that none of the vectors can be obtained as a linear combination of the others.

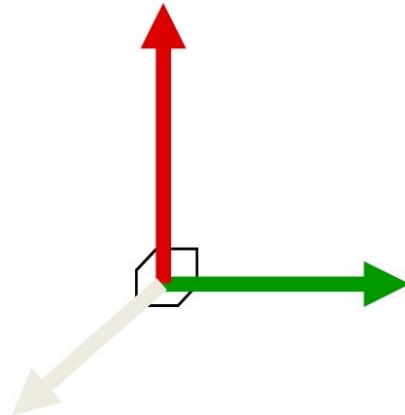
# Linear independence - example

- Parallel vectors are always dependent:



$$\mathbf{v} = 2.4 \mathbf{w} \implies \mathbf{v} + (-2.4)\mathbf{w} = 0$$

- Orthogonal vectors are always linearly independent:



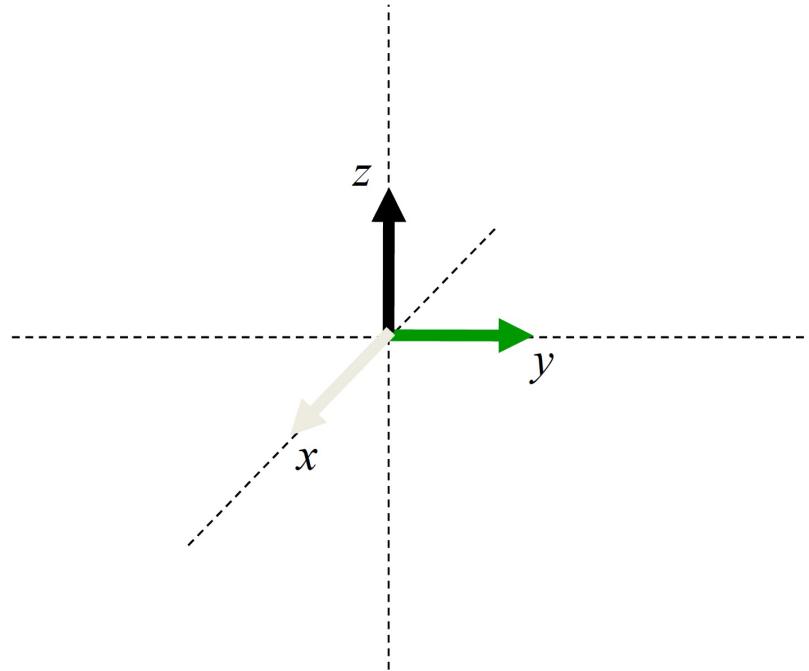


# Basis of Vector Space

- $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$  are linear independent
- $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$  span the whole vector space  $\mathbf{V}$ 
$$\mathbf{V} = \{\alpha_1 \mathbf{v}_1 + \dots + \alpha_k \mathbf{v}_k \mid \alpha_1 \in R\}$$
- Any vector in  $\mathbf{V}$  is a unique linear combination of the basis
- The number of basis vectors is called the dimension of  $\mathbf{V}$

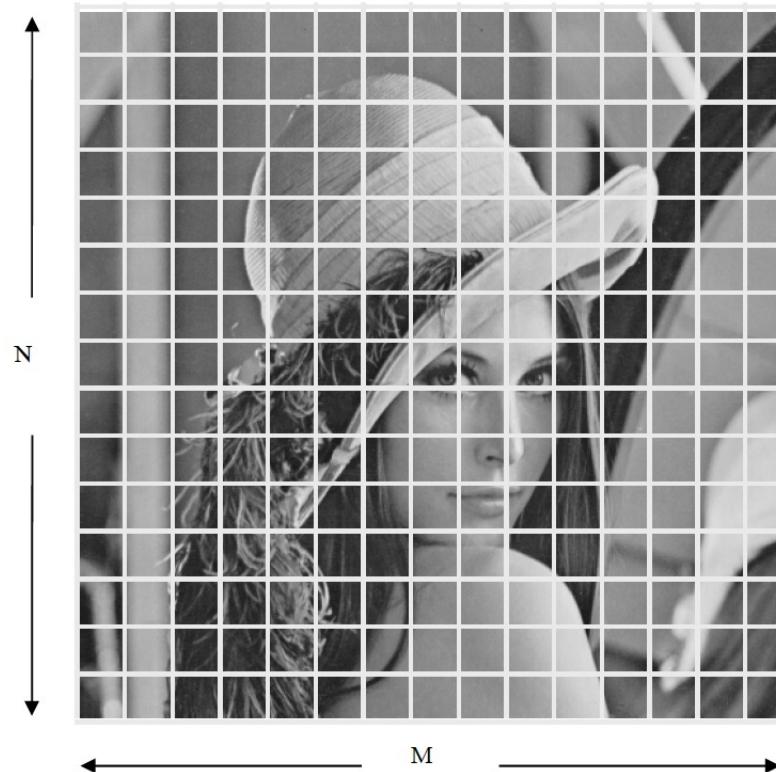
# Basis Example

- The standard basis of  $\mathbb{R}^3$



# Basis – Another Example

- Grayscale  $N \times M$  images:
  - Each pixel has value between 0 (black) and 1 (white)
  - The image can be interpreted as a vector  $\in \mathbb{R}^2$





# Probability Theory

- Axioms of Probability
- Discrete Random Variables
- Continuous Random Variables
- The Univariate Gaussian
- Maximum-likelihood Estimator



# The Axioms of Probability

- $P(A) = \# \text{ of event} / \# \text{ of possible event}$
- $0 \leq P(A) \leq 1$
- $P(\text{True}) = 1 \text{ & } P(\text{False}) = 0$
- $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$



# Discrete Random Variables

- Sample space  $\Omega$ : Possible “states”  $x$  of the random variable  $X$  (Outcomes of the experiment, output of the system, measurement)
- Either have a finite or countable number of states.
- Events: Possible combination of states ('subsets of  $\Omega$ ')



# Probability Mass Functions (PMF)

A function which tells us how likely each possible outcome is:

$$P(X = x) = P_x(x) = P(x)$$

$$\sum_{x \in \Omega} P(x) = 1$$

$$P(A) = P(x \in A) = \sum_{x \in A} P(X = x)$$



# Expectation and Variance

- Expectation (or mean):

$$E(x) = \sum_x P(X = x)x$$

- Expectation of a function:

$$E(f(x)) = \sum_x P(X = x)f(x)$$

- Moments = expectation of power of  $X$ :

$$M_k = E(X^k)$$



# Expectation and Variance

- Variance: Average (squared) fluctuation from the mean

$$\begin{aligned}Var(X) &= E\left(\left(X - E(X)\right)^2\right) = E(X^2) - E(X)^2 \\&= M_2 - M_1^2\end{aligned}$$

- Standard deviation: square root of variance  $\sqrt{Var(x)} = \sqrt{M_2 - M_1^2}$



# Bivariate Distributions

- Joint Distributions:  $P(X = x, Y = y)$ , a list of all probabilities of all possible pairs of observations
- Marginal Distribution:

$$P(X = x) = \sum_y P(X = x, Y = y)$$

- Conditional Distribution:  $P(X = x | Y = y) = P(X = x, Y = y)$ 
  - $X|Y$  has distribution  $P(X|Y)$ , where  $P(X|Y)$  specifies a “lookup-table” of all possible  $P(X = x | Y = y)$



# Expectation and Covariance of Bivariate Distributions

- Conditional distributions are just distributions which have a (conditional) mean or variance
- Covariance is the expected value of the product of fluctuations:

$$\text{Cov}(X, Y) = E \left( (X - E(X))(Y - E(Y)) \right) = E(XY) - E(X)E(Y)$$
$$\text{Var}(X) = \text{Cov}(X, X)$$

- One common way to construct bivariate random variables is to have a random variable whose parameter is another random variable.



# Expectation and Covariance of Bivariate Distributions

- Two events are independent if knowing that the first took places tells us nothing about the probability of the second:  $P(A|B) = P(A)$
- $P(A)P(B) = P(A \cap B)$
- Two random variables are independent if the joint PMF is the product of the marginals:
  - $P(X = x, Y = y) = P(X = x)P(Y = y)$
- If X and Y are independent, we write  $X \perp Y$ . Knowing the value of X does not tell us anything about Y.
- If X and Y are independent,  $Cov(X, Y) = 0$ .
- Mutual information is a measure of how “non-independent” two random variables are.



# Multivariate Distributions

- $X, x$  are vector.
- Mean:

$$E(X) = \sum_x x P(x)$$

- Covariance Matrix:

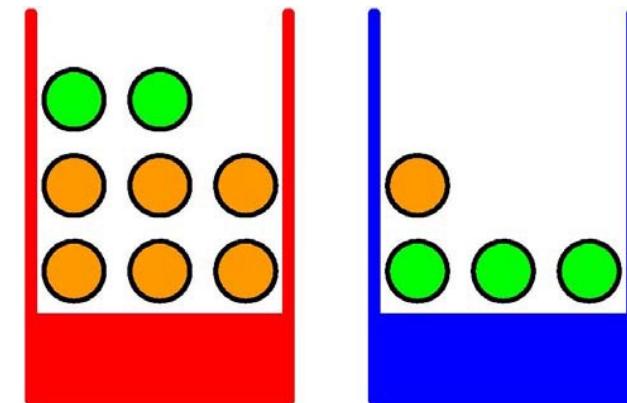
$$\begin{aligned} \text{Cov}(X_i, X_j) &= E(X_i X_j) - E(X_i) E(X_j) \\ \text{Cov}(X) &= E(X X^T) - E(X) E(X)^T \end{aligned}$$

- Conditional and marginal distributions: Can define and calculate any (multi- or single dimensional) marginals or conditional distributions we need:

$P(X_1), P(X_1, X_2), P(X_1, X_2, X_3 | X_4)$ , etc...

# Example

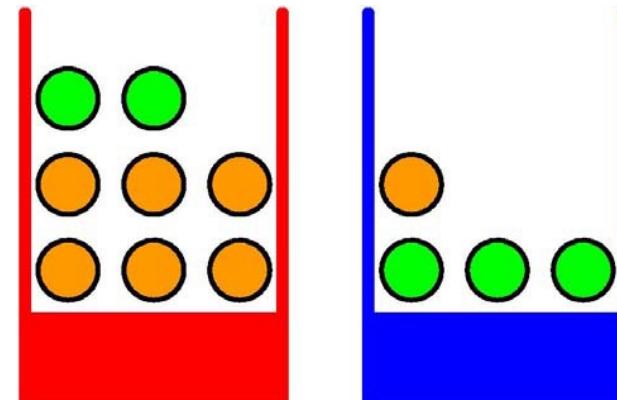
- Let us look at the following example:
  - We have two boxes, one red and one blue
  - Red box: 2 apples and 6 oranges
  - Blue box: 3 apples and 1 orange
  - Pick red box 40% of the time and blue box 60% of the time, then pick one item of fruit



C.M. Bishop, “Pattern Recognition and Machine Learning”, 2006

# Example

- Define:
  - B random variable for box picked
    - $B = \{\text{blue}(b), \text{red}(r)\}$
  - F identity of fruit
    - $F = \{\text{apple}(a), \text{orange}(o)\}$
  - $P(B=r)=0.4$  and  $P(B=b)=0.6$
  - Events are mutually exclusive and include all possible outcomes
  - Their probabilities must sum to 1

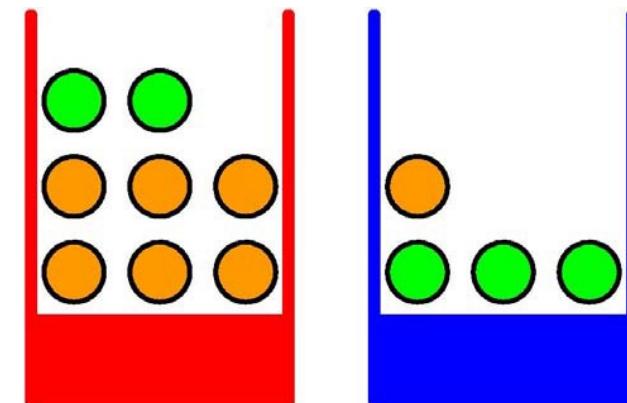


# Example

- $P(B=r) = 0.4$ ,  $P(B=b) = 0.6$
- $P(B=r) + P(B=b) = 1.0$

- Conditional Probabilities

- $P(F=a|B=r) = 2/8 = 0.25$
- $P(F=o|B=r) = 6/8 = 0.75$
- $P(F=a|B=b) = 3/4 = 0.75$
- $P(F=o|B=b) = 1/4 = 0.25$



# Example

- Note:  $P(F=a|B=r) + P(F=o|B=r) = 1$
- $P(F=a) = P(F=a|B=r)P(B=r) + P(F=a|B=b)P(B=b)$   
 $= 1/4 * 4/10 + 3/4 * 6/10 = 11/20$

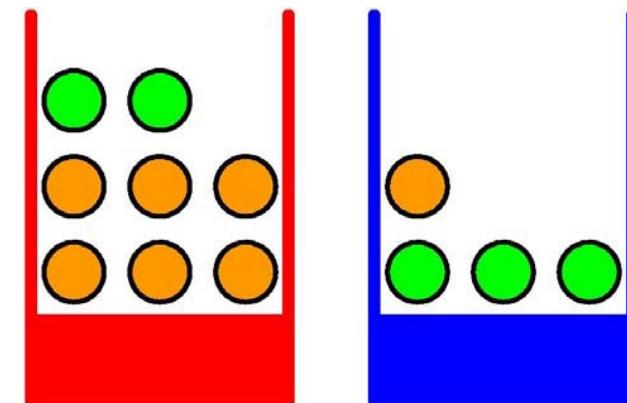
- $P(F=o) = 1 - 11/20 = 9/20 = 0.45$

$$P(F=a|B=r) = 2/8 = 0.25$$

$$P(F=o|B=r) = 6/8 = 0.75$$

$$P(F=a|B=b) = 3/4 = 0.75$$

$$P(F=o|B=b) = 1/4 = 0.25$$





# Prior vs. Posterior

- Prior Probability - If we had been asked which box had been chosen before being told the identity of the selected item of fruit, then the most complete information we have available is provided by the probability  $P(B)$ .
- Posterior Probability - Once we are told that the fruit is an orange, we can then use Bayes' theorem to compute the probability  $P(B|F)$ , which we shall call the posterior probability because it is the probability obtained after we have observed F.



# Conditional Probability

- Conditional probability: Recalculated probability of event A after someone tells you that event B happened.

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$
$$P(A \cap B) = P(A|B)P(B)$$

- Bayes Rule:

$$P(B|A) = \frac{P(A|B)P(B)}{P(A)}$$



# Bayesian Probability

- Bayesian view: probabilities provide a quantification of uncertainty. Before observing the data, the assumptions about  $w$  are captured in the form of a prior probability distribution  $P(w)$ . The effect of the observed data  $D = \{(x_1, y_1), \dots, (x_N, y_N)\}$  is expressed by  $P(D|w)$ .
- Bayes' theorem:

$$P(w|D) = \frac{P(D|w)P(w)}{P(D)}$$

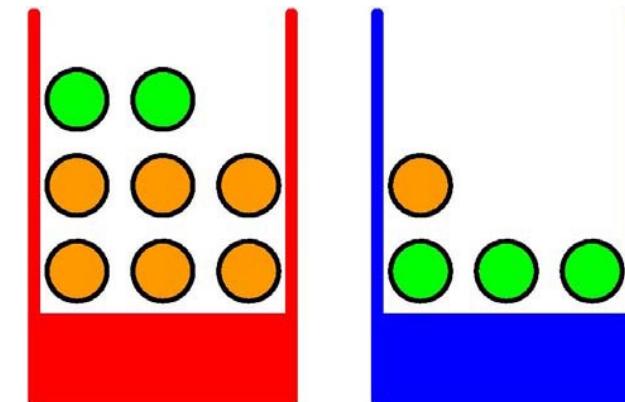
- Bayes' theorem in words: **posterior  $\propto$  likelihood  $\times$  prior**

# Bayes Rule on the Fruit Example

- Assume we have now picked an orange
- Question: which is the probability that it was from the red box?

$$P(B = r|F = o) = \frac{P(F = o|B = r)P(B = r)}{P(F = o)} = \frac{0.75 \times 0.4}{0.45} = \frac{2}{3}$$

$P(B=r) = 0.4$ ,  $P(B=b) = 0.6$   
 $P(F=o) = 0.45$ ,  
 $P(F=o|B=r) = 6/8 = 0.75$





# Continuous Random Variables

- A random variable  $X$  is continuous if its sample space  $X$  is uncountable.
- In this case,  $P(X = x) = 0$  for each  $x$ .
- If  $p_x(x)$  is a probability density function for  $X$ , then

$$P(a < X < b) = \int_a^b p(x)dx$$
$$P(a < X < a + dx) \approx p(a) \cdot dx$$



# Continuous Random Variables

- The cumulative distribution function is  $F_x(x) = P(X < x)$ . We have that  $p_x(x) = F'(x)$ , and  $F(x) = \int_{-\infty}^x p(s)ds$ .
- More generally, If  $A$  is an event, then

$$P(A) = P(X \in A) = \int_{x \in A} p(x)dx$$
$$P(\Omega) = P(X \in \Omega) = \int_{x \in \Omega} p(x)dx = 1$$



# Mean, Variance, and Conditionals

- Mean:  $E(x) = \int_x x \cdot p(x)dx$
- Variance:  $Var(X) = E(X^2) - E(X)^2$
- If  $X$  has pdf  $p(x)$ , then  $X|X \in A$  has pdf

$$p_{(x|A)}(x) = \frac{p(x)}{P(A)} = \frac{p(x)}{\int_{x \in A} p(x)dx}$$

- Only makes sense if  $P(A) > 0$  !



# Bivariate Continuous Distributions

- $p_{x,y}(x, y)$ , joints probability density function of X and Y
- $\int_x \int_y p(x, y) dx dy = 1$
- Marginal distribution:  $p_X(x) = \int_{-\infty}^{\infty} p(x, y) dy$
- Conditional distribution:  $p(x|y) = \frac{p(x,y)}{p(y)}$
- Note:  $P(Y = y) = 0!$
- Independence: X and Y are independent if  $p_{(x,y)}(x, y) = p_x(x)p_y(y)$



# The Univariate Gaussian

- Probability density function

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

- Easy to validate:

$$\int_{-\infty}^{\infty} p(x|\mu, \sigma^2) dx = 1$$

- Expectation

$$E(x) = \int_{-\infty}^{\infty} p(x|\mu, \sigma^2) x dx = \mu$$

- Variance

$$Var(x) = E(x^2) - E(x)^2 = \sigma^2$$



# Products of Gaussian pdfs

- Suppose  $p_1(x) = p\left(x, \mu_1, \frac{1}{\beta_1}\right)$  and  $p_2(x) = p\left(x, \mu_2, \frac{1}{\beta_2}\right)$ , then

$$p_1(x)p_2(x) \propto p\left(x, \mu, \frac{1}{\beta}\right)$$

$$\beta = \beta_1 + \beta_2$$

$$\mu = \frac{1}{\beta}(\beta_1\mu_1 + \beta_2\mu_2)$$

- In general,

$$p_1(x)p_2(x) \dots p_n(x) \propto p\left(x, \mu, \frac{1}{\beta}\right)$$

$$\beta = \sum \beta_n$$

$$\mu = \frac{1}{\beta} \sum_n \mu_n \beta_n$$

- This is also true for multivariate Gaussians!



# Maximum Likelihood Estimator

- Assuming data points are independent and identically distributed (i.i.d.), the probability of the data set given  $\mu$  and  $\sigma^2$  (the likelihood function):

$$P(\mathbf{x}|\mu, \sigma^2) = \prod_{n=1}^N p(x_n|\mu, \sigma^2)$$

- log-likelihood:

$$\log P(\mathbf{x}|\mu, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 - \frac{N}{2} \log \sigma^2 - \frac{N}{2} \log(2\pi)$$

- Maximizing log-likelihood with respect to  $\mu$  and  $\sigma^2$ :

$$\mu_{ML} = \frac{1}{N} \sum_{n=1}^N x_n, \sigma_{ML}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{ML})^2$$



# Maximum Likelihood Estimator

- Assuming data points are independent and identically distributed (i.i.d.), the probability of the data set given  $\mu$  and  $\sigma^2$  (the likelihood function):

$$P(\mathbf{x}|\mu, \sigma^2) = \prod_{n=1}^N p(x_n|\mu, \sigma^2)$$

- log-likelihood:

$$\log P(\mathbf{x}|\mu, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 - \frac{N}{2} \log \sigma^2 - \frac{N}{2} \log(2\pi)$$

- Maximizing log-likelihood with respect to  $\mu$  and  $\sigma^2$ :

$$\mu_{ML} = \frac{1}{N} \sum_{n=1}^N x_n, \sigma_{ML}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{ML})^2$$



# Summary

- Linear Algebra & Probability Theory
  - Linear Algebra: we reviewed definitions.
  - Linear Algebra: Orthogonality, Operations, and Eigen vectors and values will be the major concepts need to know for machine learning algorithms' understanding
  - Probability theory will be mainly used to understand the models made from machine learnings.
  - Note that not all machine learning algorithms use probability theories to produce outputs. There are also non-probabilistic algorithms that sometimes are relatively useful and robust.
- Python Review
  - From the demonstration files, we looked some of data types in python and observed how to handle them. (There will be more as we go through the semester.)
  - We observed how to computations.



# CS 559: Classification

Lecture 5

In Jang

[ijang@stevens.edu](mailto:ijang@stevens.edu)



# Lecture Outline

- Classification I
  - KNN
  - Linear Discriminant Analysis (LDA)
  - Perceptron
  - Naïve Bayes
- Classification II
  - Probabilistic Generative Models
  - Probabilistic Discriminative Models (Logistic Regression)
  - Generalization

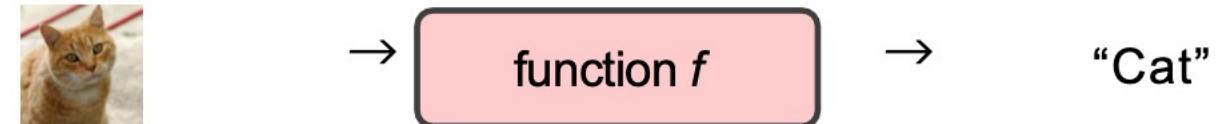


# Classification

Classification task: finding a function  $f$  that classifies examples into given set of categories  $\{C_1, C_2, \dots, C_k\}$



A classification example:





# KNN - K nearest neighbors

- Idea: average the values of the k closest observations

$$\hat{y} = \frac{1}{k} \sum_{x_i \in N_k(x)} y_i$$

where  $N_k(x)$  is the set of observations with the k smallest distances to the query point  $x$ .

- Assumption: similar inputs have similar outputs
- Rule: for a test input  $x$ , assign the most common label amongst its  $k$  most similar training inputs.
- Steps:
  - Choose the number of  $k$  and a distance metric.
  - Find the  $k$ -nearest neighbors of the sample that we want to classify.
  - Assign the class label by majority vote.



# KNN - K nearest neighbors

- What distance function to use?
  - KNN relies on a distance metric.
  - Better classification from the better metric reflecting similar similarity
  - Common choice: the Makowski distance

$$\text{dist}(\mathbf{x}, \mathbf{z}) = \left( \sum_{r=1}^d |x_r - z_r|^p \right)^{1/p}$$

- If  $p = 1$ : Manhattan distance
- If  $p = 2$ : Euclidean distance
- If  $p \rightarrow \infty$ : Max

- Advantage: the classifier adopts immediately as we collect new training data.
- Disadvantage:
  - the computational cost grows linearly with the number of samples
  - Very susceptible to overfitting due to the **curse of dimensionality**



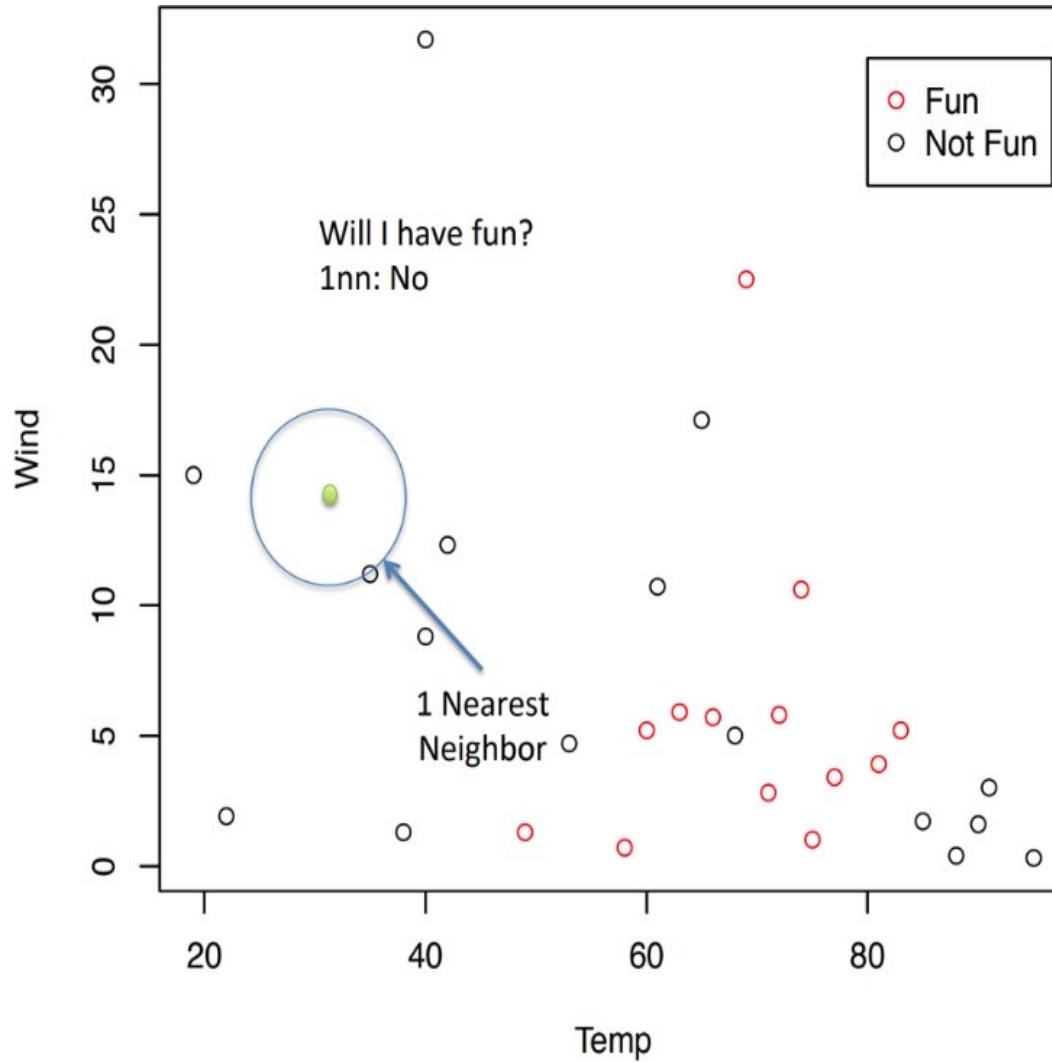
# KNN – an example

- Is it a good day to go for a run?
- A runner's data on past running. With recorded temperature, wind and whether the run was fun:
  - Temperature (degrees F)
  - Wind Speed (mph)
  - Fun (yes, no)
- It is now 65 degrees and the wind is 9 mph. Will a run be fun?

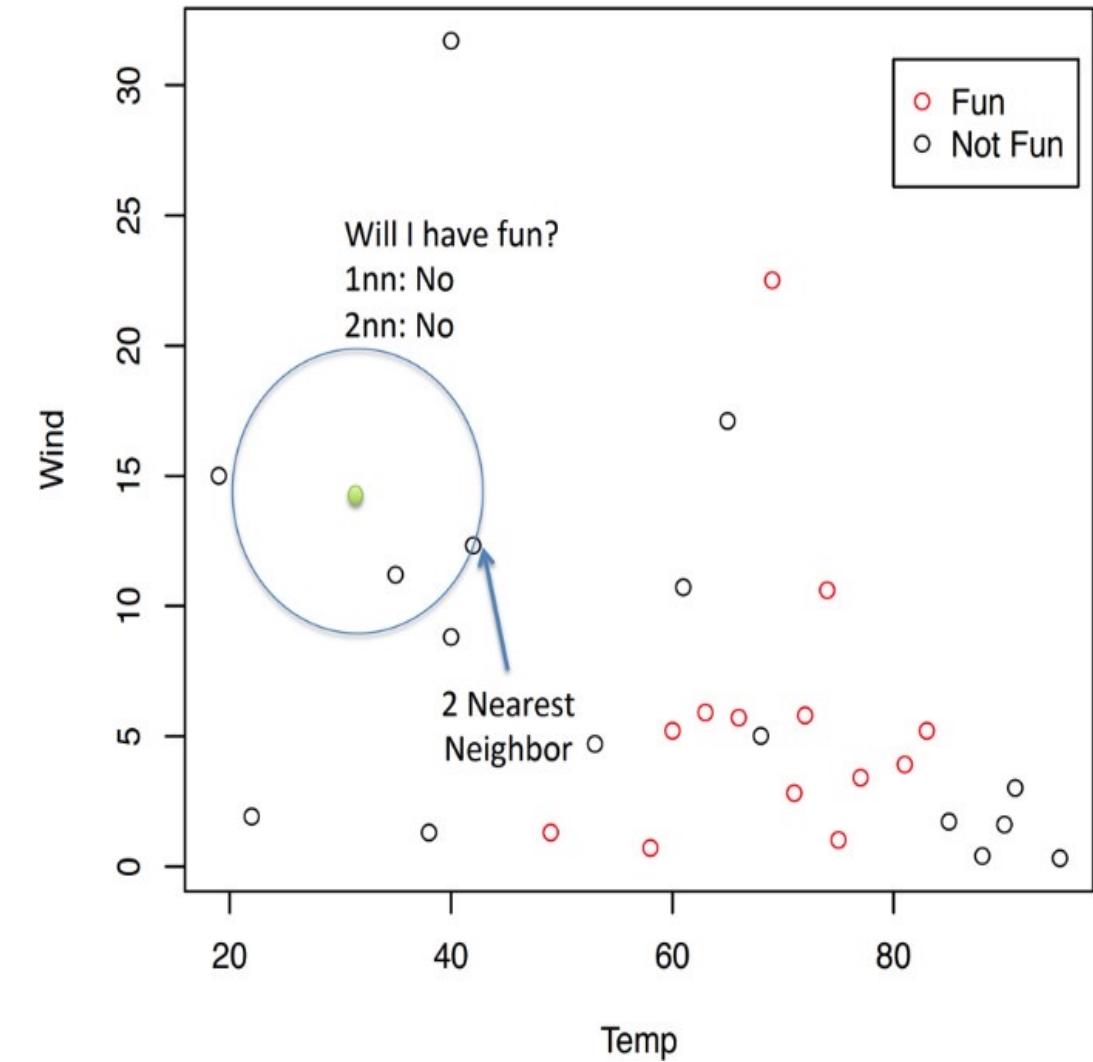
# KNN classification



KNN classification: k=1

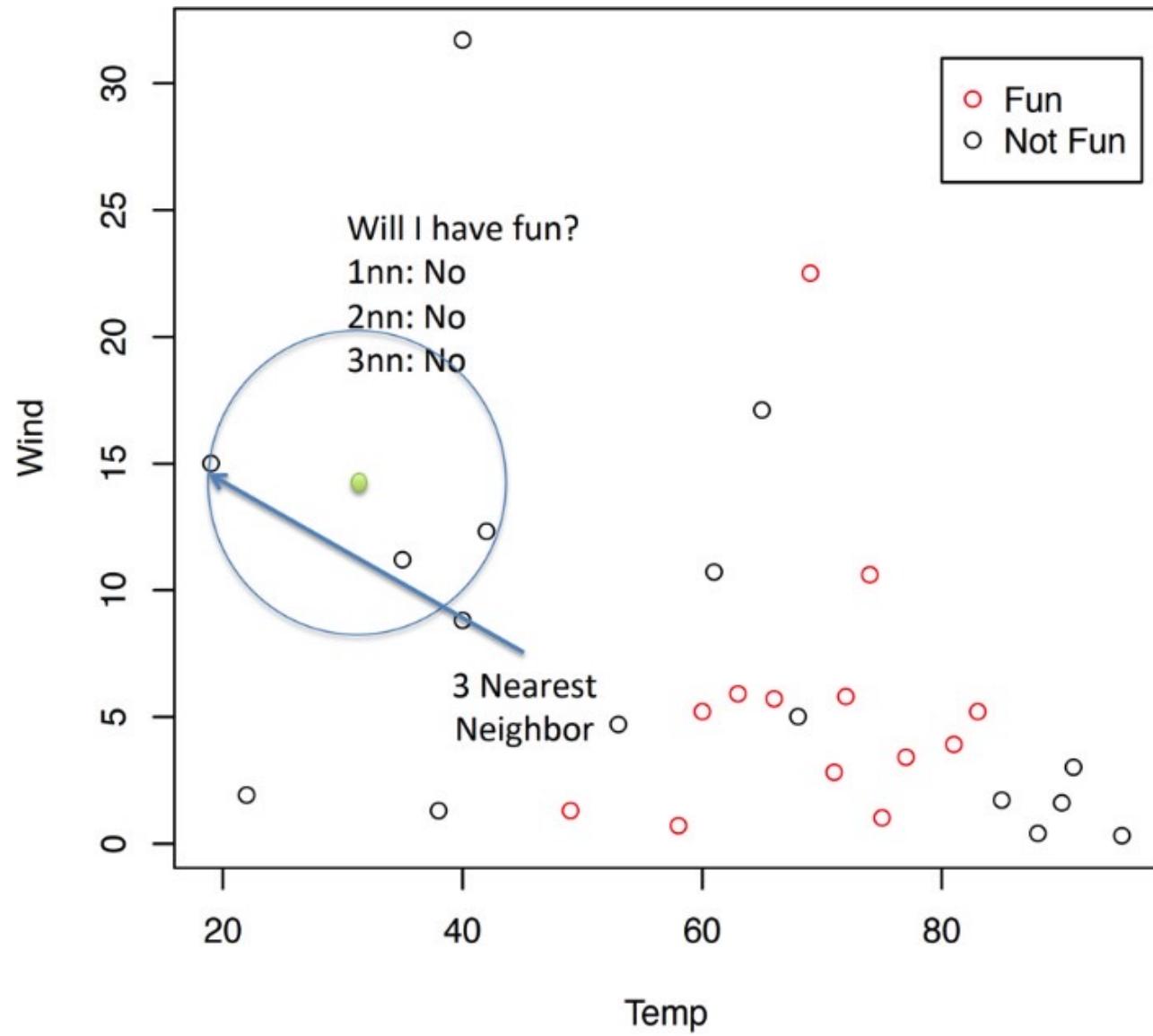


KNN classification: k=2



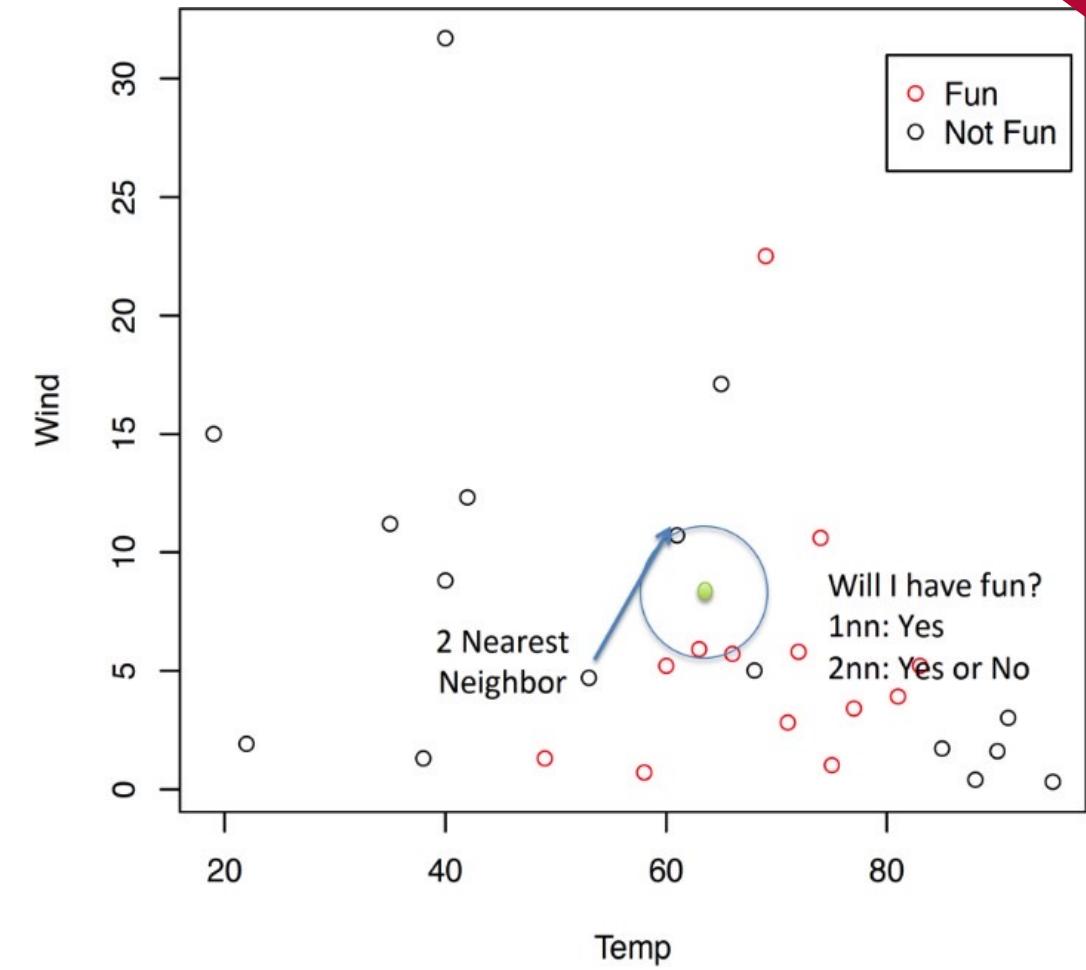
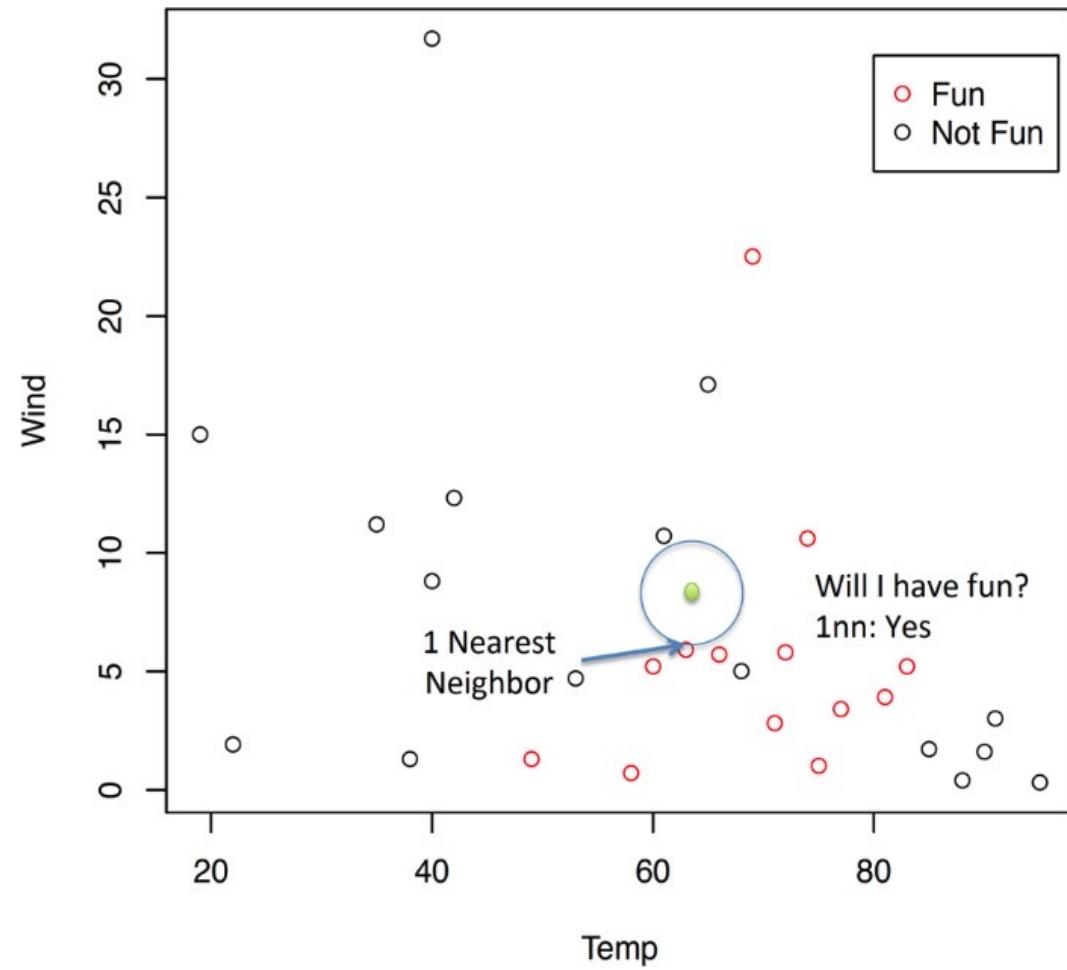


# KNN classification: k=3



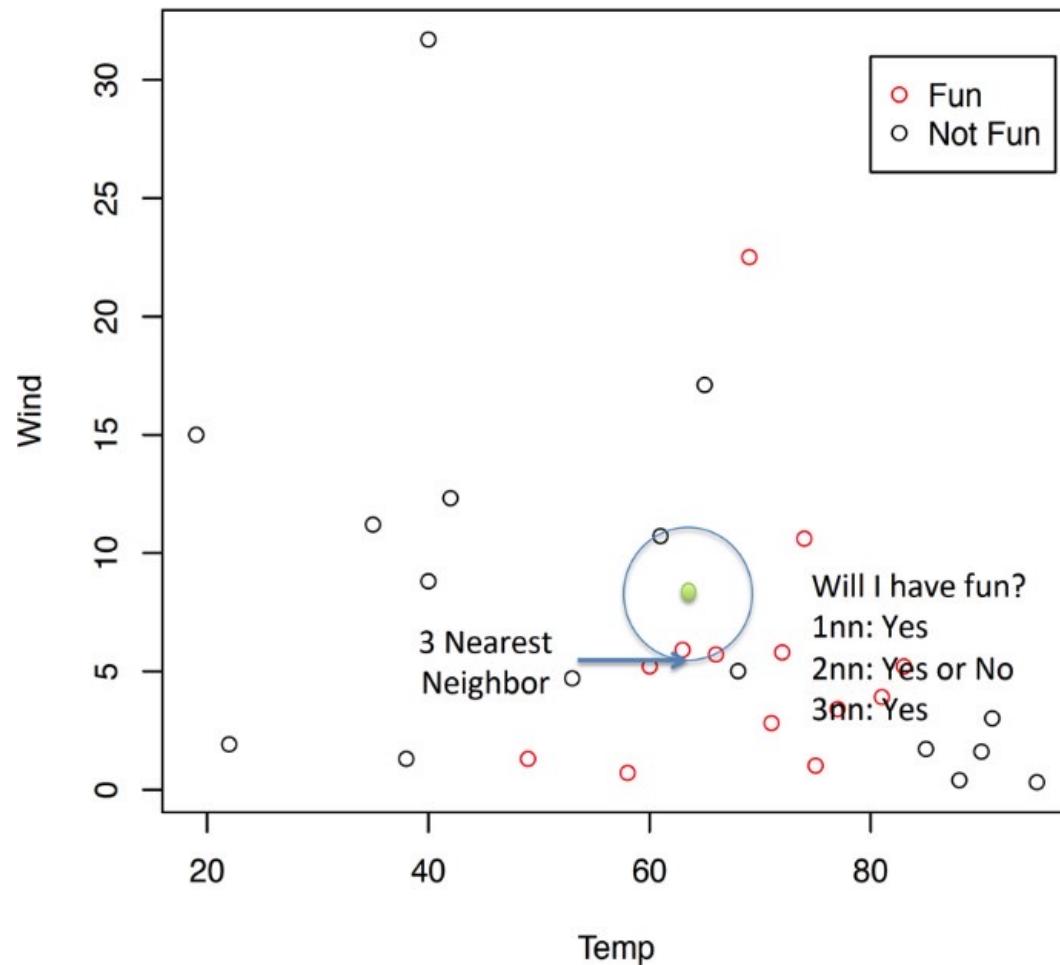


# KNN classification: k=1





# KNN classification: k=3





# Decision Theory for Classification

Decision theory, when combined with probability theory, allows us to make optimal decisions in situations involving uncertainty.

- Training data: input values  $X$  and target values  $y$
- Inference stage: use the training data to learn a model for  $p(C_k|x)$
- Decision stage: use the given posterior probabilities to make optimal class assignments.



# Generative Methods

- Solve the inference problem of estimating the **class-conditional densities**  $p(x|C_k)$  for each class  $C_k$
- Infer the **prior class probabilities**  $p(C_k)$
- Use Bayes' theorem to find the **class posterior probabilities**:

$$p(C_k|x) = \frac{p(x|C_k)p(C_k)}{p(x)}$$

where

$$p(x) = \sum_k p(x|C_k)p(C_k)$$

- Use decision theory to determine class membership for each new input  $x$ .

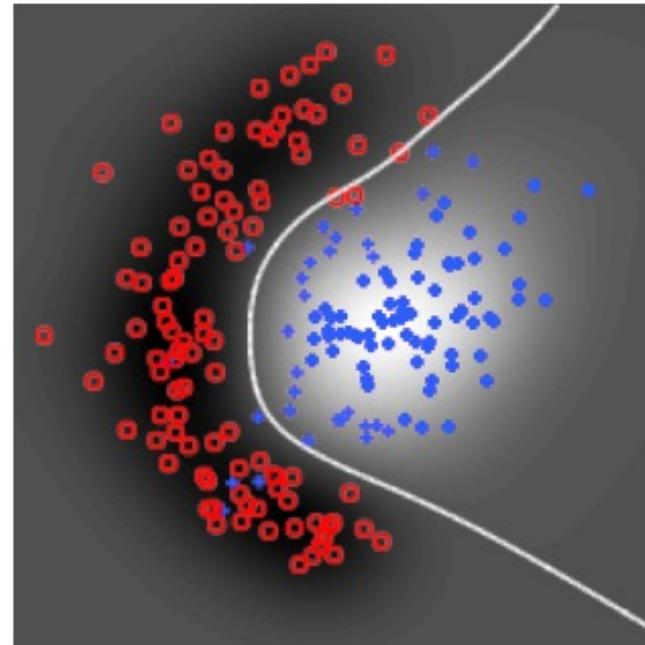
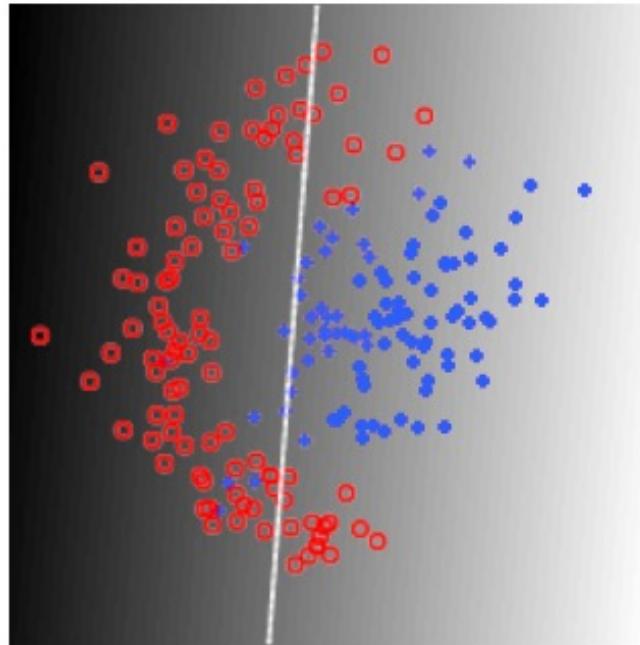


# Discriminative Methods

- Solve directly the inference problem of estimating the **class posterior probabilities**  $p(C_k|x)$ .
- Discriminative Functions: Find a function  $f(x)$  which maps each input directly onto a class label. Probabilities play no role here.
- Use decision theory to determine class membership for each new input  $x$ .



# Linear Discriminant Functions



Of course, linear algorithms can be used together with **nonlinear feature spaces** or **nonlinear basis functions** in order to solve nonlinear classification problems!



# Linear Discriminant Functions

- Linear discriminants separate the space by a hyperplane, and the parameters define its normal vector.
- Decision function:  $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$
- Classification:
  - if  $f(x) > 0$  say  $x$  belongs to class 1 if  $f(x) < 0$  say  $x$  belongs to class -1
- The decision-surface has equation  $f(x) = 0$ , and is a hyperplane of dimensionality  $D - 1$ .
- $\mathbf{w}$  is the normal vector to the hyperplane, and points into the positive class or negative class.
- $w_0$  determines the location of the decision-surface
- $|f(\mathbf{x})|$  is proportional to the perpendicular distance to the decision-surface (with factor 1 if  $\|\mathbf{w}\| = 1$ ).



# Linear Discriminant Functions-Geometrical Properties

- Decision boundary:

$$f(x) = w^T x + \omega_0 = 0$$

- Let  $x_1, x_2$  be two points which lie on the decision boundary

$$\begin{aligned} f(x_1) &= w^T x_1 + \omega_0 = 0, f(x_2) = w^T x_2 + \omega_0 = 0 \\ &\Rightarrow w^T(x_1 - x_2) = 0 \end{aligned}$$

- $w$  represents the orthogonal direction to the decision boundary.

# Linear Discriminant Functions-Geometrical Properties Cont.

$$x_1 = x_2 + r \frac{w}{\|w\|}$$

$$f(x_2) = f\left(x_1 - r \frac{w}{\|w\|}\right) = (w^T x_1 + w_0) - r\|w\|$$

$$\rightarrow r = \frac{f(x_1)}{\|w\|}$$

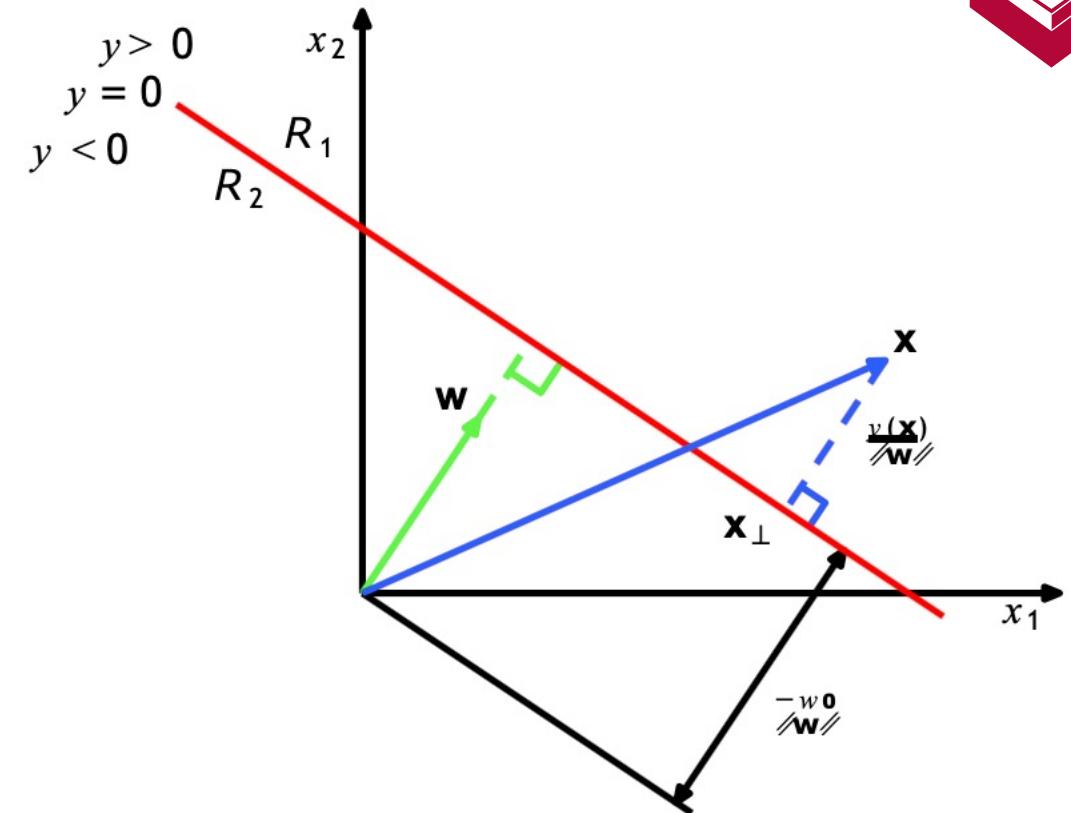
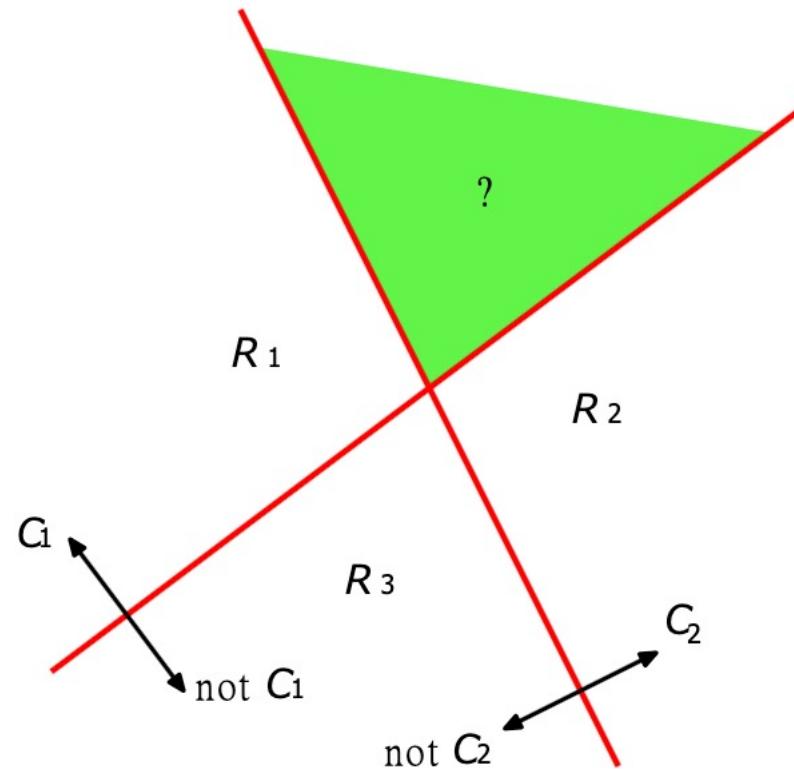


Figure: Signed orthogonal distance of the origin from the decision



# Linear Discriminant Functions: Multiple classes

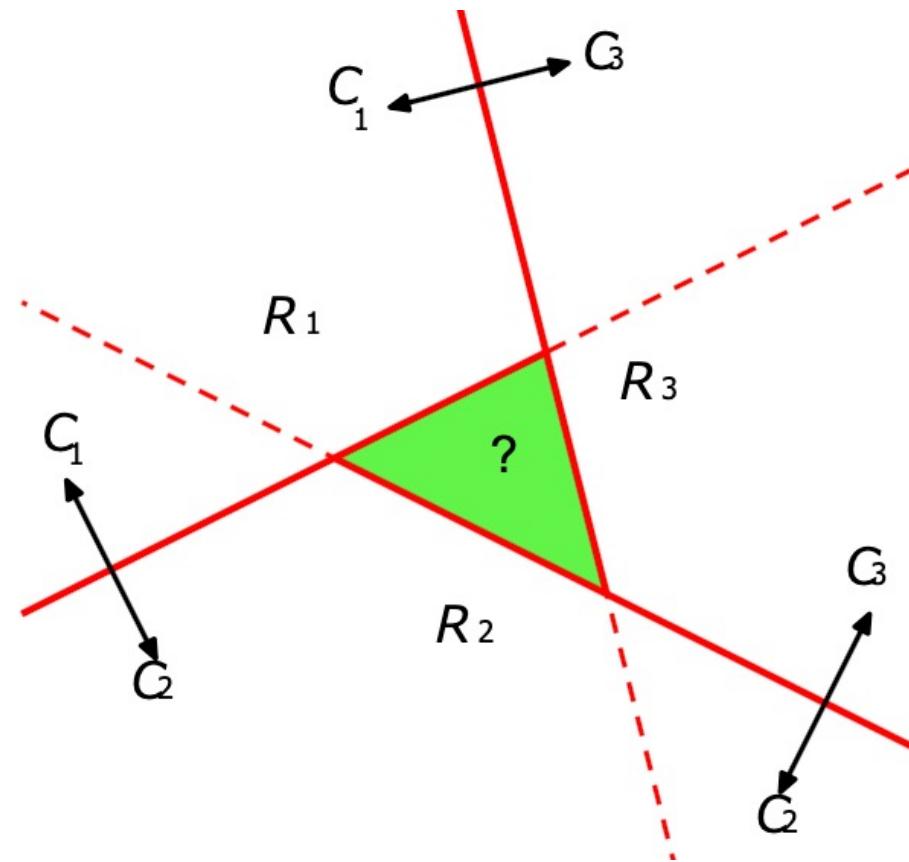
one-versus-the-rest:  $K - 1$  classifiers each of which solves a two-class problem of separating points of  $C_k$  from points not in that class.





# Linear Discriminant Functions: Multiple classes

one-versus-one:  $K(K-1)/2$  binary discriminant functions, one for every possible pair of class. Each point is then classified according to a majority vote amongst the discriminant functions.





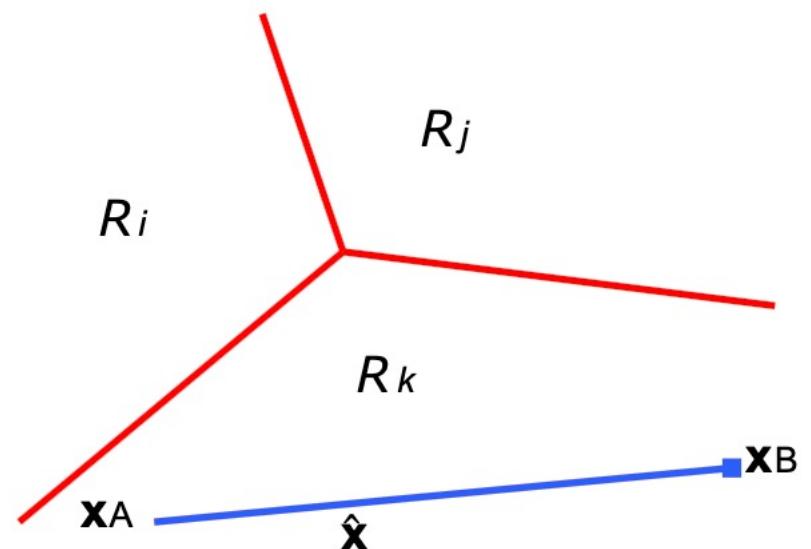
# Linear Discriminant Functions: Multiple classes

- Solution: consider a single K-class discriminant comprising K linear functions of the form

$$f_k(x) = w^T x + w_{k0}$$

- Assign a point  $x$  to class  $C_k$  if  $f_k(x) > f_j(x) \forall j, x = k$
- The decision boundary between class  $C_k$  and class  $C_j$  is given by:

$$f_k(x) = f_j(x) \Rightarrow (w_k - w_j)^T x + (w_{k0} - w_{j0}) = 0$$

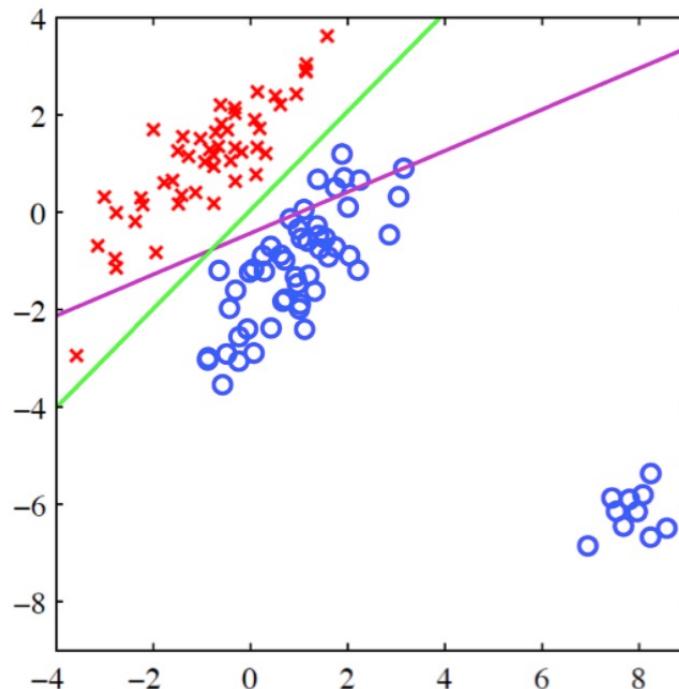
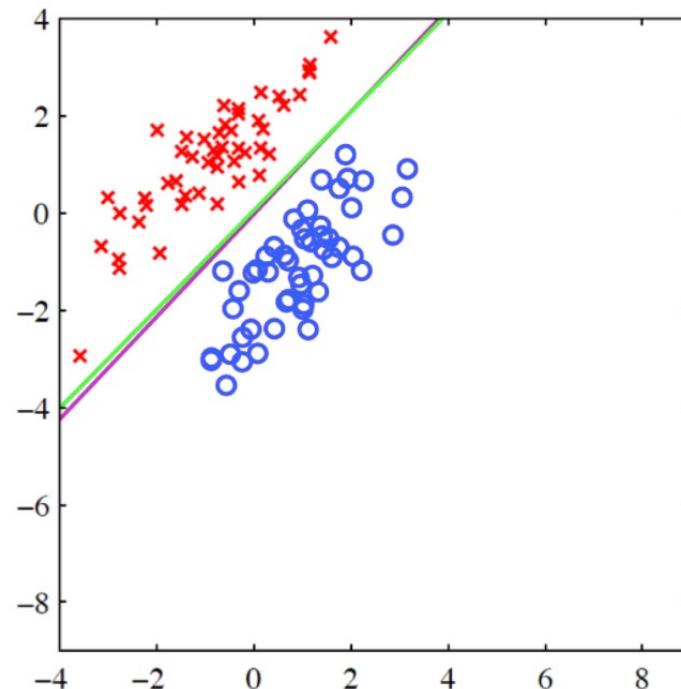


# Least square classification

- We have to fit the function  $f(x) = w^T x + \omega_0$  to data.
- Simply do a linear regression from  $x$  to  $y$  by minimizing the sum-of-squared errors  $\sum_n(f(x_n) - y_n)^2$ .

$$w_{reg} = \left( \sum_n x_n x_n^T \right)^{-1} \sum_n x_n y_n$$

- Q: In what situations might this be a bad idea?



# Least square classification

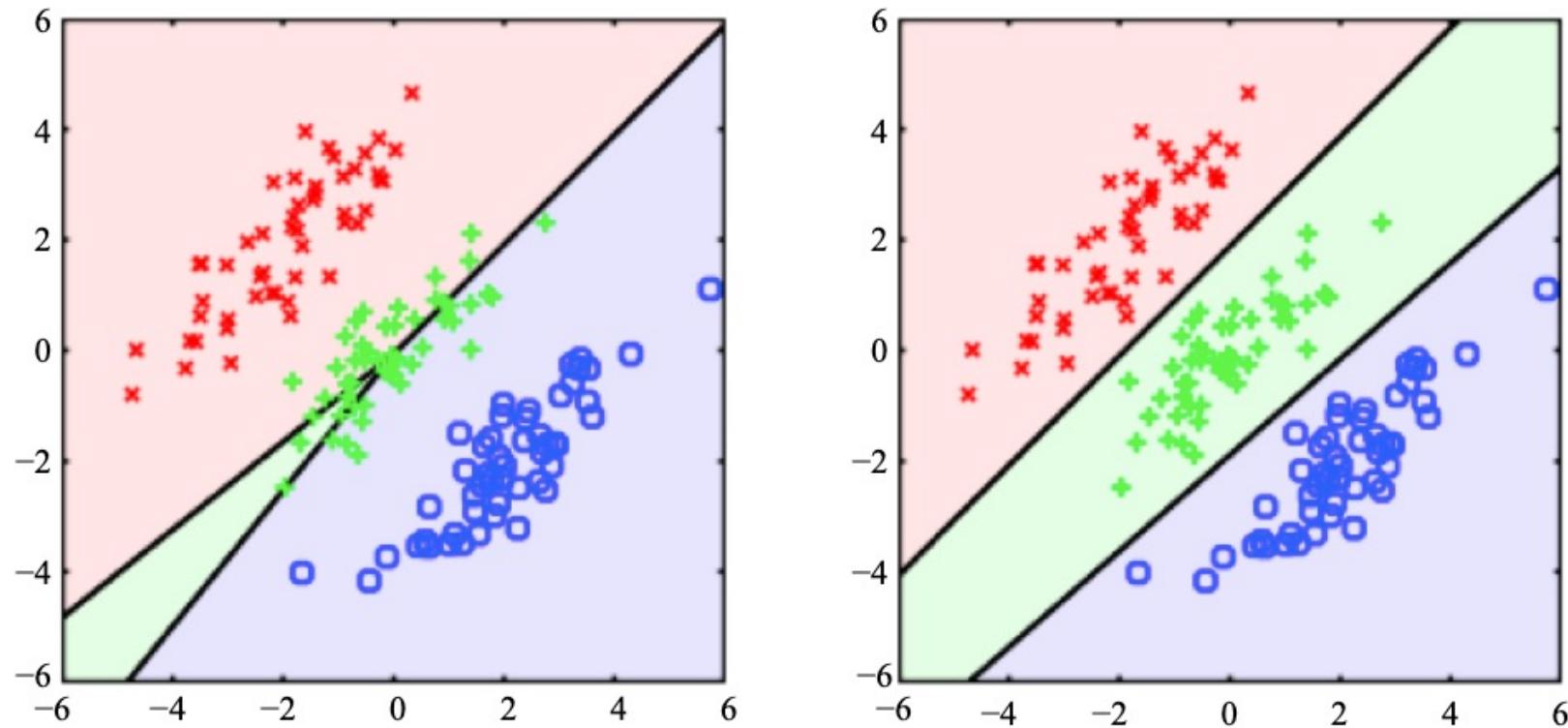
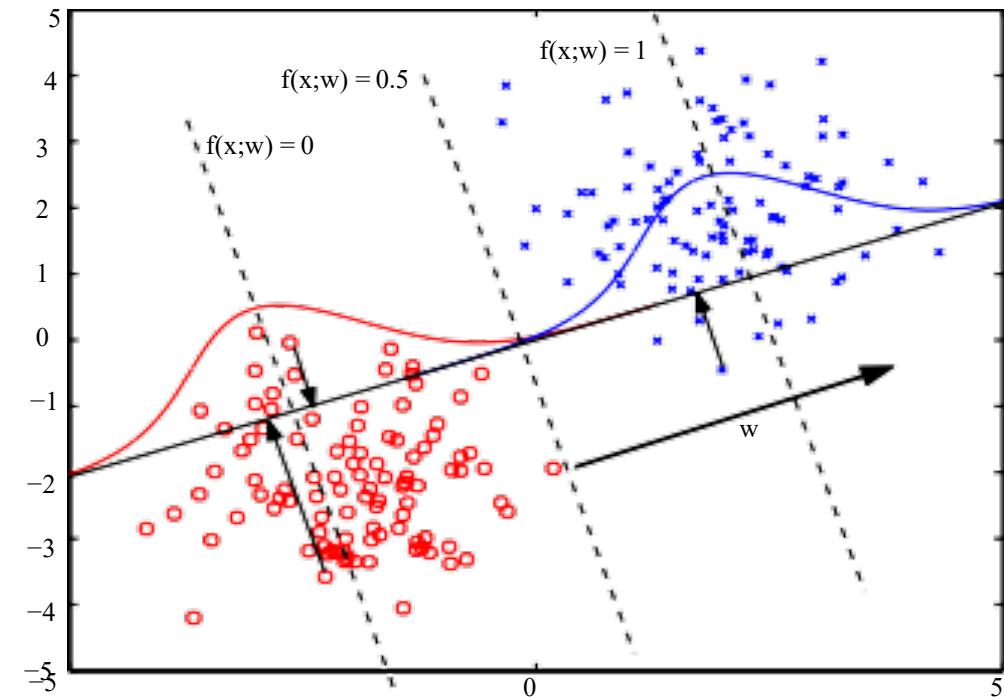


Figure: Left: using a least-squares discriminant; Right: using logistic regression

# Classification via projection

- A linear function:  $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$  assuming in 2D, projects each point  $\mathbf{x} = [x_1, x_2]^T$  to a line parallel to  $\mathbf{w}$ :

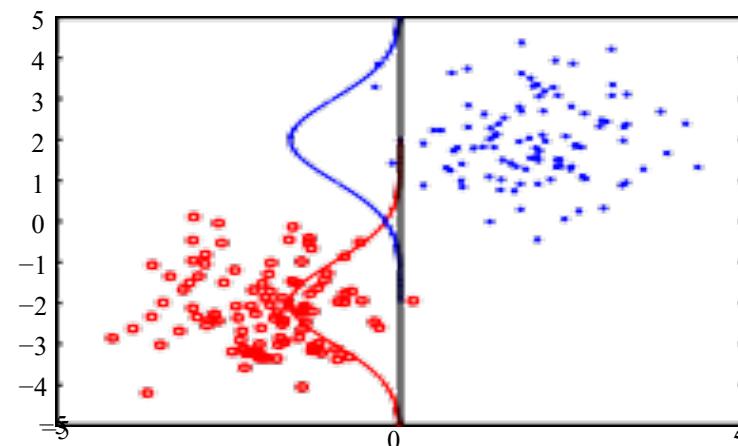
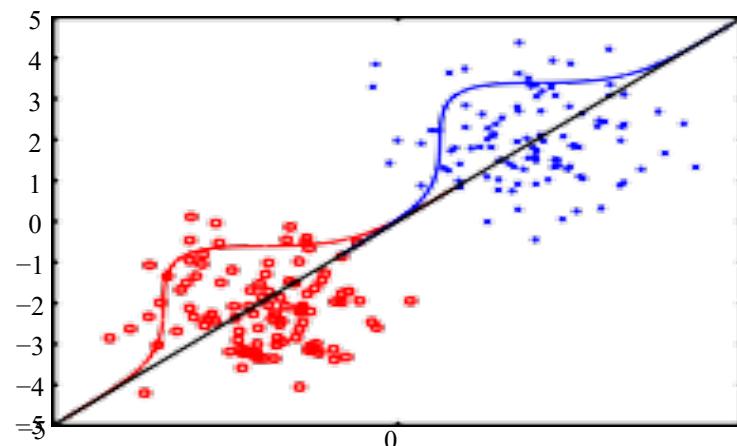
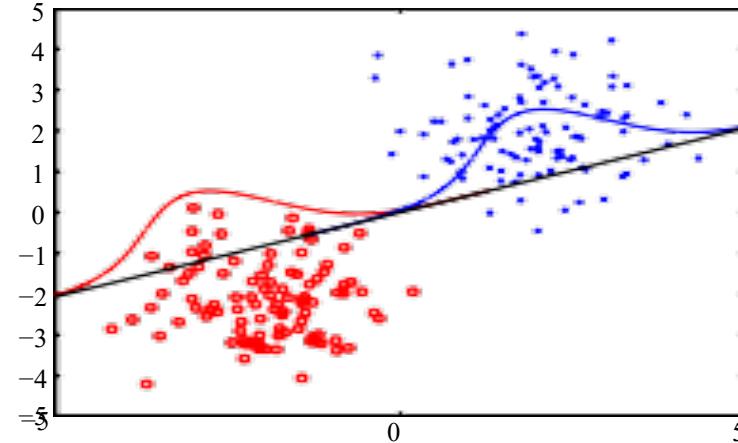
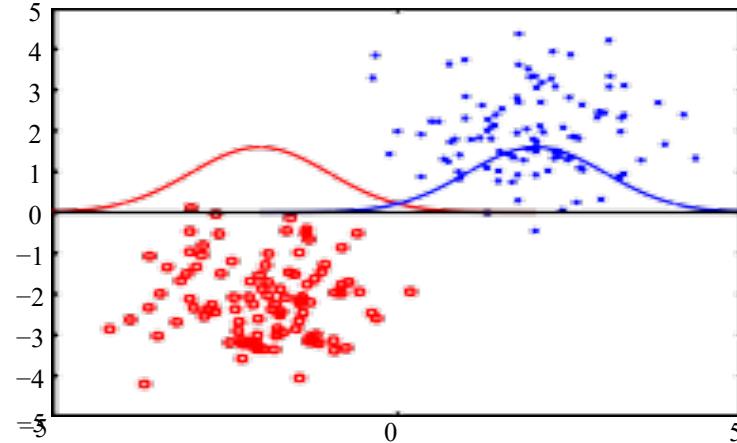
| Point in $R^d$ | Projected point in $R$            |
|----------------|-----------------------------------|
| $x_1$          | $z_1 = \mathbf{w}^T \mathbf{x}_1$ |
| $x_2$          | $z_2 = \mathbf{w}^T \mathbf{x}_2$ |
| $\vdots$       | $\vdots$                          |
| $x_n$          | $z_n = \mathbf{w}^T \mathbf{x}_n$ |



- We can study how well the projected points  $z_1, \dots, z_n$  viewed as functions of  $w$  are separated across the classes.

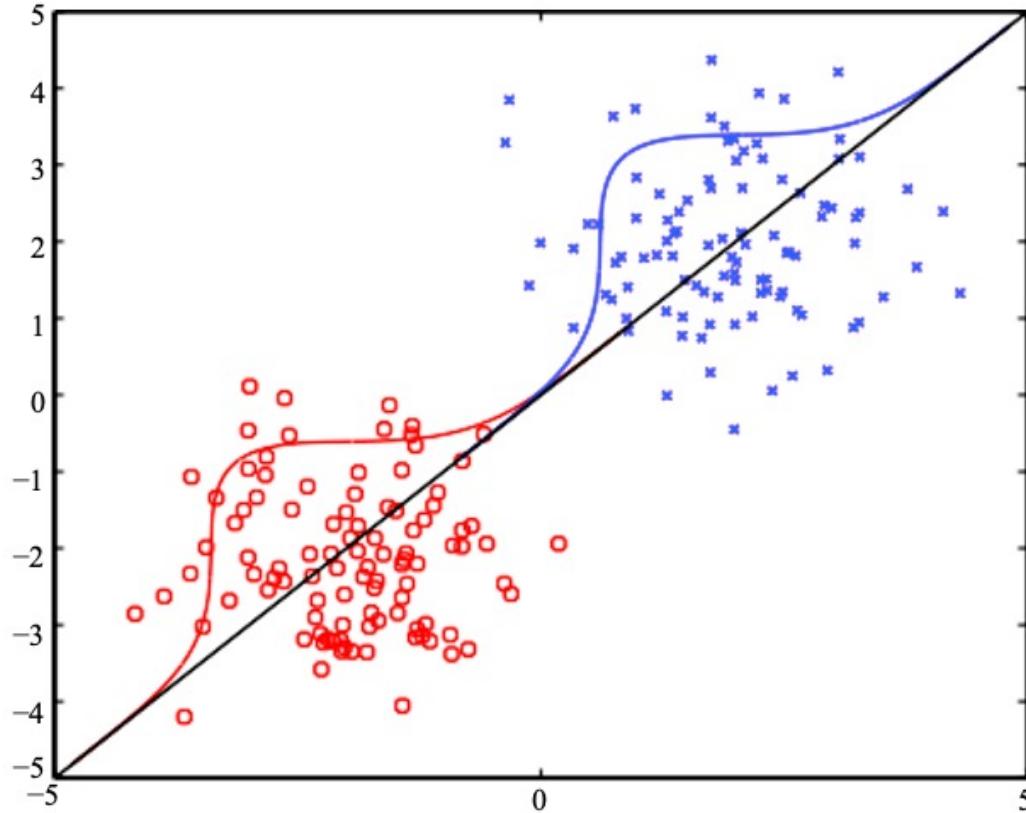
# Classification via projection

By varying  $\mathbf{w}$  we get different levels of separation between the projected points



# Optimizing the projection

- We would like to find  $\mathbf{w}$  that somehow maximizes the separation of the projected points across classes.



- We can quantify the separation (overlap) in terms of means and variances of the resulting 1-dimensional class distributions



# Fisher's linear discriminant

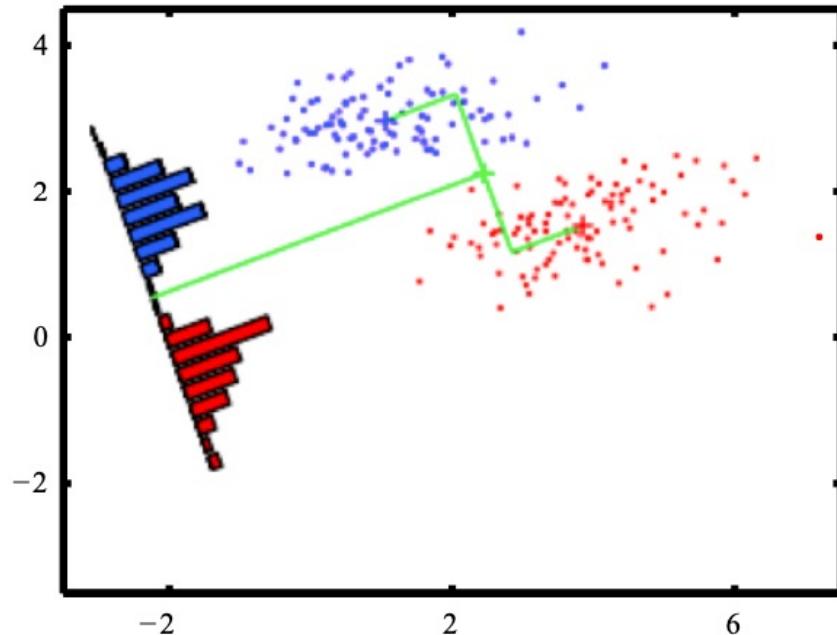
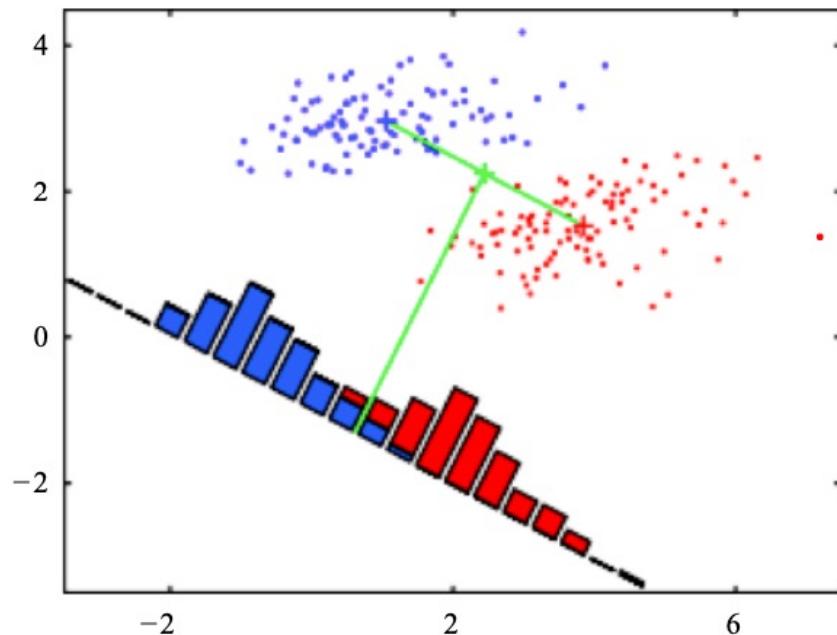
- One way to view a linear classification model is in terms of dimensionality reduction.
- Two class case: suppose we project  $x$  onto one dimension:

$$f = \mathbf{w}^T \mathbf{x}$$

- Set a threshold  $t$ :

|               |                              |
|---------------|------------------------------|
| if $f \leq t$ | assign $C_1$ to $\mathbf{x}$ |
| otherwise     | assign $C_2$ to $\mathbf{x}$ |

# Fisher's linear discriminant



- Find an orientation along which the projected samples are well separated;
- This is exactly the goal of linear discriminant analysis (LDA);
- In other words: we are after the linear projection that best separates the data, i.e. best discriminates data of different classes.



# Fisher's linear discriminant

- Two classes:  $\{C_1, C_2\}$
- $N_1$  samples of class  $C_1$
- $N_2$  samples of class  $C_2$
- Consider  $\mathbf{w} \in R^d$  with  $\|\mathbf{w}\| = 1$
- Then:  $\mathbf{w}^T \mathbf{x}$  is the projection of  $\mathbf{x}$  along the direction of  $\mathbf{w}$ .
- We want the projections  $\mathbf{w}^T \mathbf{x}$  where  $\mathbf{x} \in C_1$  separated from the projections  $\mathbf{w}^T \mathbf{x}$  where  $\mathbf{x} \in C_2$



# Fisher's linear discriminant

- A measure of the separation between the projected points is the difference of the sample means:
- Sample mean of class  $C_1$ :

$$\mathbf{m}_1 = \frac{1}{N_1} \sum_{n \in C_1} x_n$$

- Sample mean for the projected points:

$$m_1 = \frac{1}{N_1} \sum_{n \in C_1} \mathbf{w}^T x_n = \mathbf{w}^T \mathbf{m}_1$$
$$\rightarrow |m_1 - m_2| = \mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2)$$

- We wish to make the above difference as large as we can. In addition, ...



# Fisher's linear discriminant

To obtain good separation of the projected data, we really want the difference between the means to be large relative to some measure of the standard deviation of each class:

- Scatter of the projected samples of class  $C_1$ :

$$s_1^2 = \sum_{n \in C_1} (w^T x_n - m_1)^2$$

- Total within-class scatter of the projected samples:

$$s_1^2 + s_2^2$$

- Fisher linear discriminant analysis:

$$\arg \max_w \frac{|m_1 - m_2|^2}{s_1^2 + s_2^2}$$



# Fisher's linear discriminant

$$\text{Fisher's criterion } J(w): J(w) = \frac{|m_1 - m_2|^2}{s_1^2 + s_2^2}$$

To obtain  $J(w)$  as an explicit function of  $w$ , we define the following matrices:

$$S_1 = \sum_{n \in C_1} (x_n - \mathbf{m}_1)(x_n - \mathbf{m}_1)^T$$

Within-class scatter matrix:

$$S_w = S_1 + S_2$$

Then:

$$\begin{aligned} s_1^2 &= \sum_{x \in C_1} (\mathbf{w}^T x - m_1)^2 = \sum_{x \in C_1} (\mathbf{w}^T x - \mathbf{w}^T \mathbf{m}_1)^2 \\ &= \sum_{x \in C_1} \mathbf{w}^T (x - \mathbf{m}_1)(x - \mathbf{m}_1)^T \mathbf{w} = \mathbf{w}^T S_1 \mathbf{w} \end{aligned}$$



# Fisher's linear discriminant

So,  $s_1^2 = \mathbf{w}^T \mathbf{S}_1 \mathbf{w}$  and  $s_2^2 = \mathbf{w}^T \mathbf{S}_2 \mathbf{w}$

Thus,

$$\begin{aligned}s_w &= s_1^2 + s_2^2 = \mathbf{w}^T \mathbf{S}_1 \mathbf{w} + \mathbf{w}^T \mathbf{S}_2 \mathbf{w} \\&= \mathbf{w}^T (\mathbf{S}_1 + \mathbf{S}_2) \mathbf{w} \\&= \mathbf{w}^T \mathbf{S}_w \mathbf{w}\end{aligned}$$

where  $\mathbf{S}_w = \sum_{n \in C_1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^T + \sum_{n \in C_2} (\mathbf{x}_n - \mathbf{m}_2)(\mathbf{x}_n - \mathbf{m}_2)^T$ .

Similarly:

$$\begin{aligned}(m_1 - m_2)^2 &= (\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_2)^2 \\&= \mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w} \\&= \mathbf{w}^T \mathbf{S}_B \mathbf{w}.\end{aligned}$$

where  $\mathbf{S}_B = (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^T$  is the between-class scatter matrix.



# Fisher's linear discriminant

We have obtained:

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}}$$

$J(\mathbf{w})$  is maximized when

$$(\mathbf{w}^T \mathbf{S}_B \mathbf{w}) \mathbf{S}_w \mathbf{w} = (\mathbf{w}^T \mathbf{S}_w \mathbf{w}) \mathbf{S}_B \mathbf{w}$$

We observe that

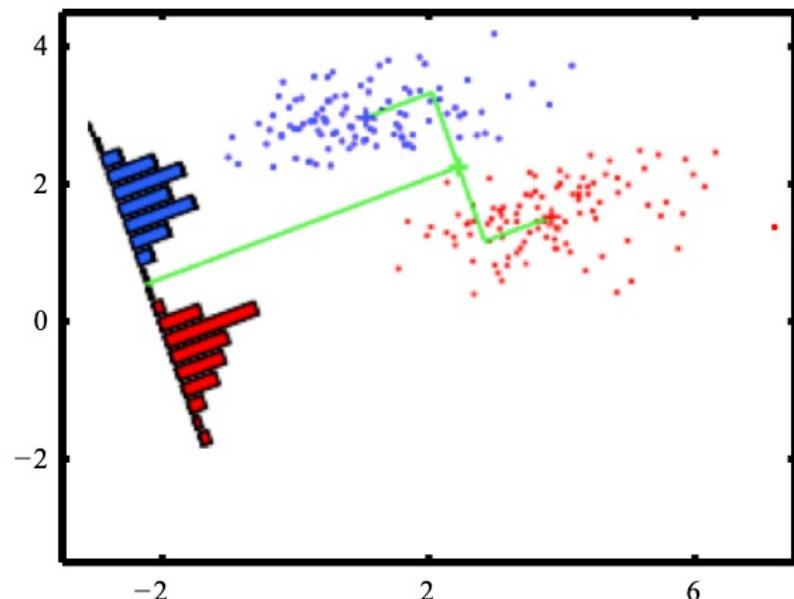
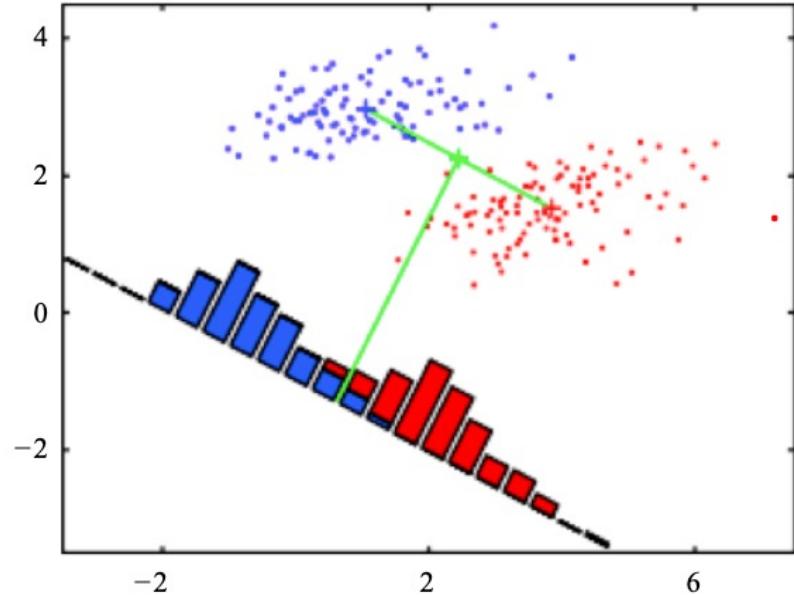
$$\mathbf{S}_B \mathbf{w} = (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w}$$

where  $(\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w}$  is a scalar and always in the direction of  $(\mathbf{m}_1 - \mathbf{m}_2)$ .

Solution:

$$\mathbf{w} \propto \mathbf{S}_w^{-1}(\mathbf{m}_1 - \mathbf{m}_2)$$

# Fisher's linear discriminant Summary



- $\mathbf{m}_1 = \frac{1}{N} \sum_{n \in C_1} \mathbf{x}_n$  &  $\mathbf{m}_2 = \frac{1}{N_1} \sum_{n \in C_2} \mathbf{x}_n$
- Maximize Projection-distance of class means  $\mathbf{w}_{simple} \propto \mathbf{m}_1 - \mathbf{m}_2$
- Maximizing distance between means ignores that the projected variances might also be big.
- Fix: Maximize the ratio of between-class variance to within-class variance ('signal to noise'). Fisher criterion:

$$J_w = \frac{(\mathbf{m}_1 - \mathbf{m}_2)^2}{S_1^2 + S_2^2}$$

$$\mathbf{w}_{Fisher} = S_w^{-1}(\mathbf{m}_1 - \mathbf{m}_2)$$



# Fisher's linear discriminant Summary: Multi-Class

- The analysis can be extended to multiple classes.
- $S_w = \sum_{k=1}^K \sum_{x_i \in C_k} (x_i - m_k)(x_i - m_k)^T$
- $S_B = \sum_{k=1}^K m_k (m_k - m)(m_k - m)^T$  where  $m$  is the global mean and  $m_k$  is the number of samples in class  $k$ .
- Solve:  $S_B \mathbf{v} = \lambda S_W \mathbf{v}$  the generalized eigenvalue problem
- At most  $K-1$  distinct solution eigenvalues
- The optimal projection matrix  $\mathbf{V}$  to a subspace of dimension  $k$  is given by the eigenvectors corresponding to the largest  $k$  eigenvalues



# Fisher's linear discriminant Summary: Multi-Class

- LDA is a linear technique for **dimensionality** reduction: it projects the data along directions that can be expressed as linear combination of the input features.
- The “appropriate” transformation depends on the data and on the task we want to perform on the data. Note that LDA uses class labels.
- **Non-linear** extensions of LDA exist (e.g., generalized LDA).



# The Perceptron Algorithm (Frank Rosenblatt, 1957)

- First learning algorithm for neural networks.
- Originally introduced for character classification, where each character is represented as an image.
- Total input to output node:

$$\sum_j w_j x_j$$

- Output unit performs the function (activation function):

$$H(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ -1 & \text{if } x < 0 \end{cases}$$



# Perceptron: Learning Algorithm

- **Goal:** compute a mapping from inputs to the outputs.
- Example: two class character recognition problem.
  - Training set: set of images representing either the character ‘a’ or the character ‘b’ (supervised learning);
  - Learning task: learn the weights so that when a new unlabelled image comes in, the network can predict its label.
  - Setting:  $d$  input units (intensity level of a pixel), 1 output unit.
- The algorithm proceeds as follows:
  - Initial random setting of weights;
  - The input is a random sequence  $\{x_k\}$
  - For each element of class  $C_1$ , if output = 1 (correct), **do nothing**; otherwise, **update weights**;
  - For each element of class  $C_2$ , if output = 0 (correct), **do nothing**; otherwise, **update weights**;



# Perceptron: Learning Algorithm

- More formally:  $x = (x_1, x_2, \dots, x_d)^T, w = (w_1, w_2, \dots, w_d)^T$
- $\theta$ : Threshold of the output unit
- Unit output:  $w^T x = w_1 x_1 + w_2 x_2 + \dots + w_d x_d$
- Output class 1 if  $w^T x - \theta \geq 0$
- To eliminate the explicit dependence on  $\theta$ : Output class 1 if:  $w^T x \geq 0$



# Perceptron: Learning Algorithm

- We want to learn values of the weights so that the perceptron correctly discriminate elements of  $C_1$  from elements of  $C_2$ .
- Given  $x$  in input, if  $x$  is classified correctly, weights are unchanged, otherwise:

$$w = \begin{cases} w + x & \text{if an element of class } C_1 \text{ was classified as in } C_2 \\ w - x & \text{if an element of class } C_2 \text{ was classified as in } C_1 \end{cases}$$



# Perceptron: Learning Algorithm

- **1<sup>st</sup> case:**  $\mathbf{x} \in C_1$  and was classified in  $C_2$ . The correct answer is 1, which corresponds to:  $\mathbf{w}^T \mathbf{x} \geq 0$ , we have  $\mathbf{w}^T \mathbf{x} < 0$ . We want to get closer to the correct answer:  $\mathbf{w}^T \mathbf{x} < \mathbf{w}'^T \mathbf{x}$ .

$$\mathbf{w}^T \mathbf{x} < \mathbf{w}'^T \mathbf{x}, \text{ iff } \mathbf{w}^T \mathbf{x} < (\mathbf{w} + \mathbf{x})^T \mathbf{x}$$

$$(\mathbf{w} + \mathbf{x})^T \mathbf{x} = \mathbf{w}^T \mathbf{x} + \mathbf{x}^T \mathbf{x} = \mathbf{w}^T \mathbf{x} + \|\mathbf{x}\|^2$$

because  $\|\mathbf{x}\|^2 > 0$ , the condition is verified.

- **2<sup>nd</sup> case:**  $\mathbf{x} \in C_2$  and was classified in  $C_1$ . The correct answer is 0, which corresponds to:  $\mathbf{w}^T \mathbf{x} < 0$ , we have  $\mathbf{w}^T \mathbf{x} \geq 0$ . We want to get closer to the correct answer:  $\mathbf{w}^T \mathbf{x} > \mathbf{w}'^T \mathbf{x}$ .

$$\mathbf{w}^T \mathbf{x} > \mathbf{w}'^T \mathbf{x}, \text{ iff } \mathbf{w}^T \mathbf{x} > (\mathbf{w} - \mathbf{x})^T \mathbf{x}$$

$$(\mathbf{w} - \mathbf{x})^T \mathbf{x} = \mathbf{w}^T \mathbf{x} - \mathbf{x}^T \mathbf{x} = \mathbf{w}^T \mathbf{x} - \|\mathbf{x}\|^2$$

because  $\|\mathbf{x}\|^2 > 0$ , the condition is verified.



# Perceptron: Learning Algorithm

In summary:

- A random sequence  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k$  is generated such that  $x_i \in C_1 \cup C_2$
- If  $\mathbf{x}_k$  is correctly classified, then  $\mathbf{w}_{k+1} = \mathbf{w}_k$  otherwise:

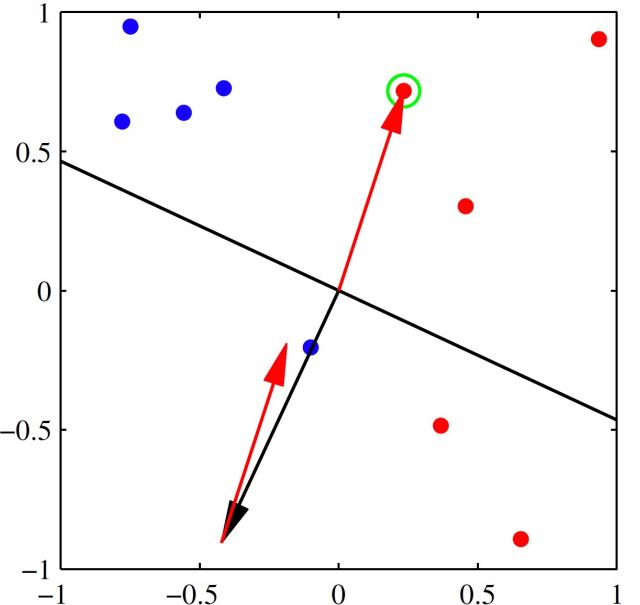
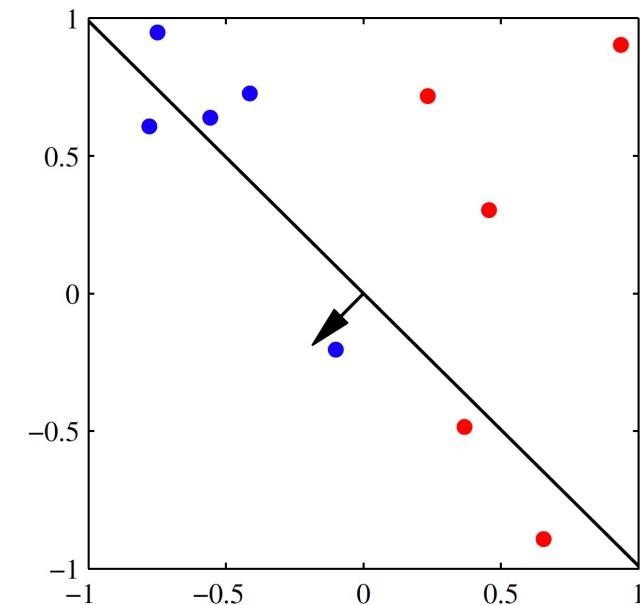
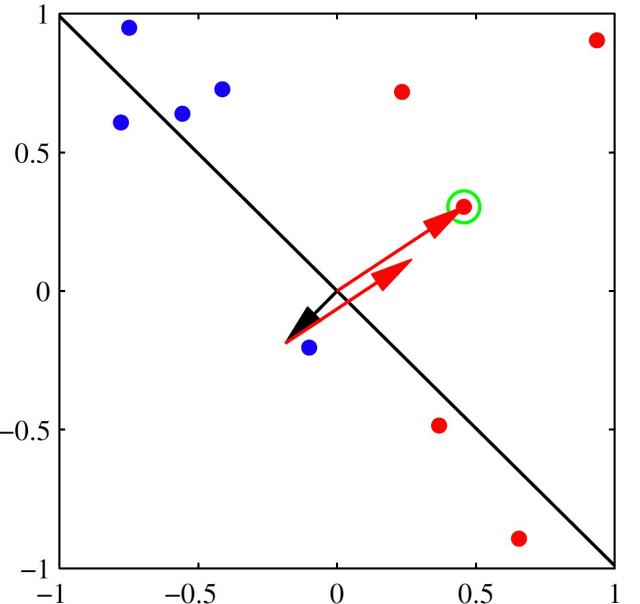
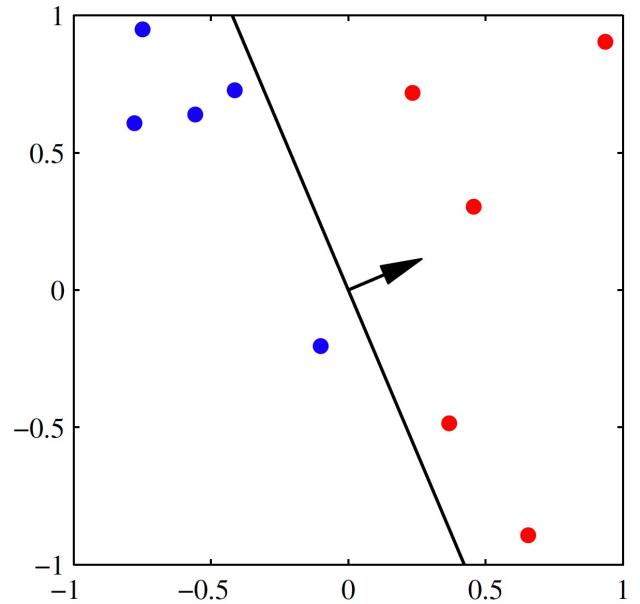
$$\mathbf{w}_{k+1} = \begin{cases} \mathbf{w}_k + \mathbf{x}_k & \text{if } \mathbf{x}_k \in C_1 \\ \mathbf{w}_k - \mathbf{x}_k & \text{if } \mathbf{x}_k \in C_2 \end{cases}$$

- Convergence theorem: regardless of the initial choice of weights, if the two classes are linearly separable, there exists  $\mathbf{w}$  such that:

$$= \begin{cases} \mathbf{w}^T \mathbf{x} \geq 0 \\ \mathbf{w}^T \mathbf{x} < 0 \end{cases}$$

then the learning rule will find such solution after a finite number of steps.

# Perceptron: Example





# Naive Bayes: not (necessarily) a Bayesian method

- $A$  and  $B$  are independent iff  $p(A, B) = p(A)p(B)$

- $A$  and  $B$  are conditionally independent given  $C$  iff

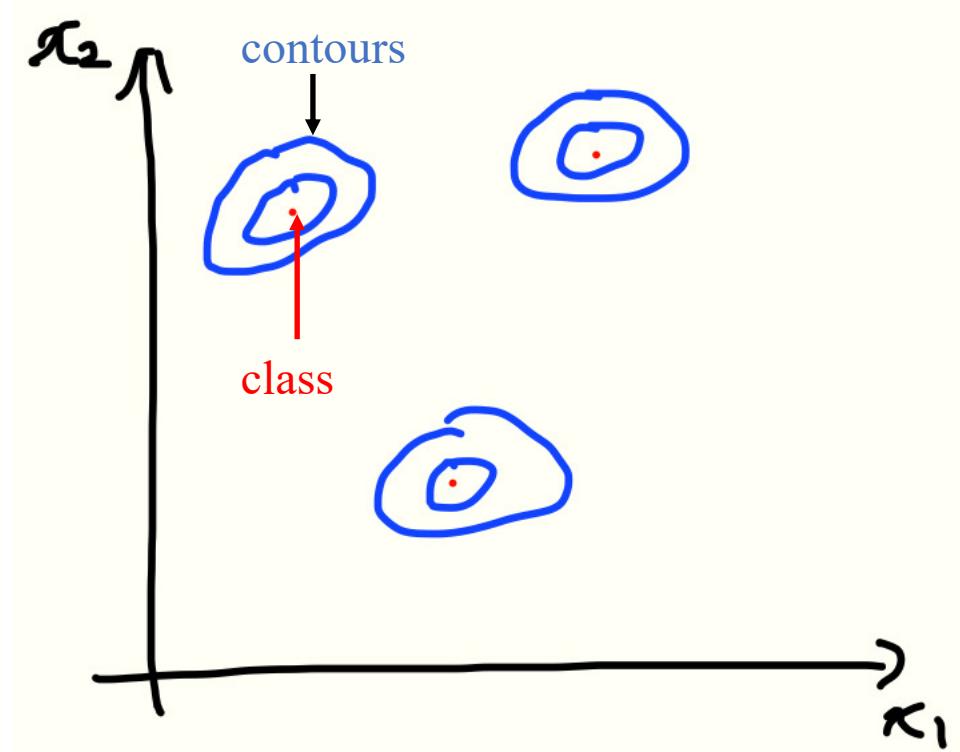
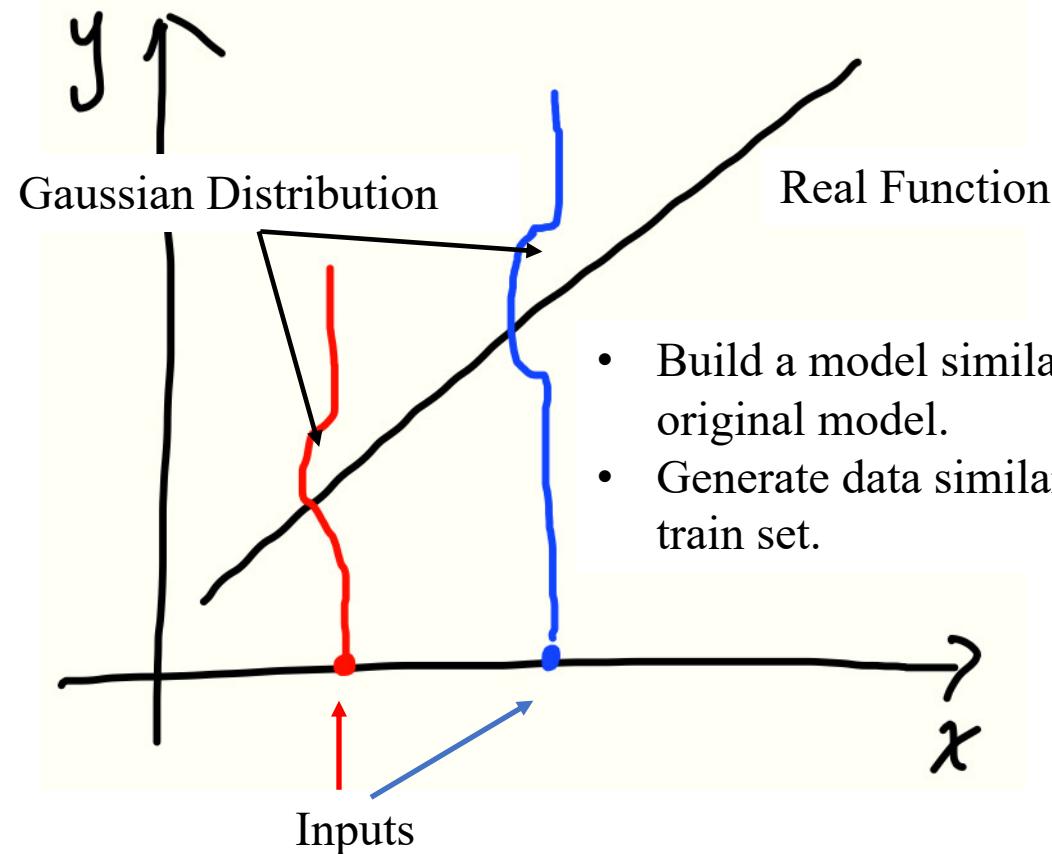
$$p(A, B|C) = p(A|C)p(B|C)$$



# Linear Classification II

# Probabilistic Generative Models

- We now turn to a probabilistic approach to classification.
- How models with linear decision boundaries arise from simple assumptions about the distribution of the data.





# Generative Approach

- Infer the **prior class probabilities**  $p(C_k)$ .
- Solve the inference problem of estimating the **class-conditional densities**  $p(\mathbf{x}|C_k)$  for each class  $C_k$ .
- Use Bayes' theorem to find the **class posterior probabilities**:

$$p(C_k|\mathbf{x}) = \frac{p(\mathbf{x}|C_k)p(C_k)}{p(\mathbf{x})} \quad (1)$$

where

$$p(\mathbf{x}) = \sum_k p(\mathbf{x}|C_k)p(C_k) \quad (2)$$

- Use decision theory to determine class membership for each new input  $\mathbf{x}$ .



- In classification, the number of classes is finite, so natural prior  $p(C)$  is the multinomial  
(e.g., coin or dice)

$$p(C = c_k) = \pi_k$$

- When  $x \in \mathbb{R}^D$ , then it is okay to assume that  $p(x|C)$  is Gaussian.
- The **class-conditional densities** are Gaussian with the same covariance matrix:

$$p(x|C_k) = \frac{1}{(2\pi)^{D/2}} \left( \frac{1}{|\Sigma|^{1/2}} \right) \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\} \quad (3)$$



# General Posterior Distribution

$$p(C_k | \mathbf{x}) = \frac{p(\mathbf{x} | C_k) p(C_k)}{\sum_k p(\mathbf{x} | C_k) p(C_k)} \quad (1)$$

$\frac{1}{(2\pi)^{D/2}} \left( \frac{1}{|\Sigma|^{1/2}} \right)$  vanishes

$$= \frac{\pi_k \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\}}{\sum_k \pi_k \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\}} \quad (4-1)$$

$\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x}$  is independent from  $C_k$  !

$$= \frac{\pi_k \exp \left\{ -\frac{1}{2} (\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - 2\boldsymbol{\mu}_k^T \boldsymbol{\Sigma}^{-1} \mathbf{x} + \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_k) \right\}}{\sum_k \pi_k \exp \left\{ -\frac{1}{2} (\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - 2\boldsymbol{\mu}_k^T \boldsymbol{\Sigma}^{-1} \mathbf{x} + \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_k) \right\}} \quad (4-2)$$

$$= \frac{\pi_k \exp \left\{ \left( \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_k \right) \right\}}{\sum_k \pi_k \exp \left\{ \left( \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_k \right) \right\}} \quad (4-3)$$



## Two Classes - $C_k$ and $C_j$

$$p(C_k | \mathbf{x}) = \frac{\pi_k \exp \left\{ \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_k \right\}}{\pi_j \exp \left\{ \boldsymbol{\mu}_j^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_j^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_j \right\} + \pi_k \exp \left\{ \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_k \right\}} \quad (5)$$

Using  $\frac{a}{a+b} = \frac{1}{1+b/a}$  and  $x = e^{\ln x}$

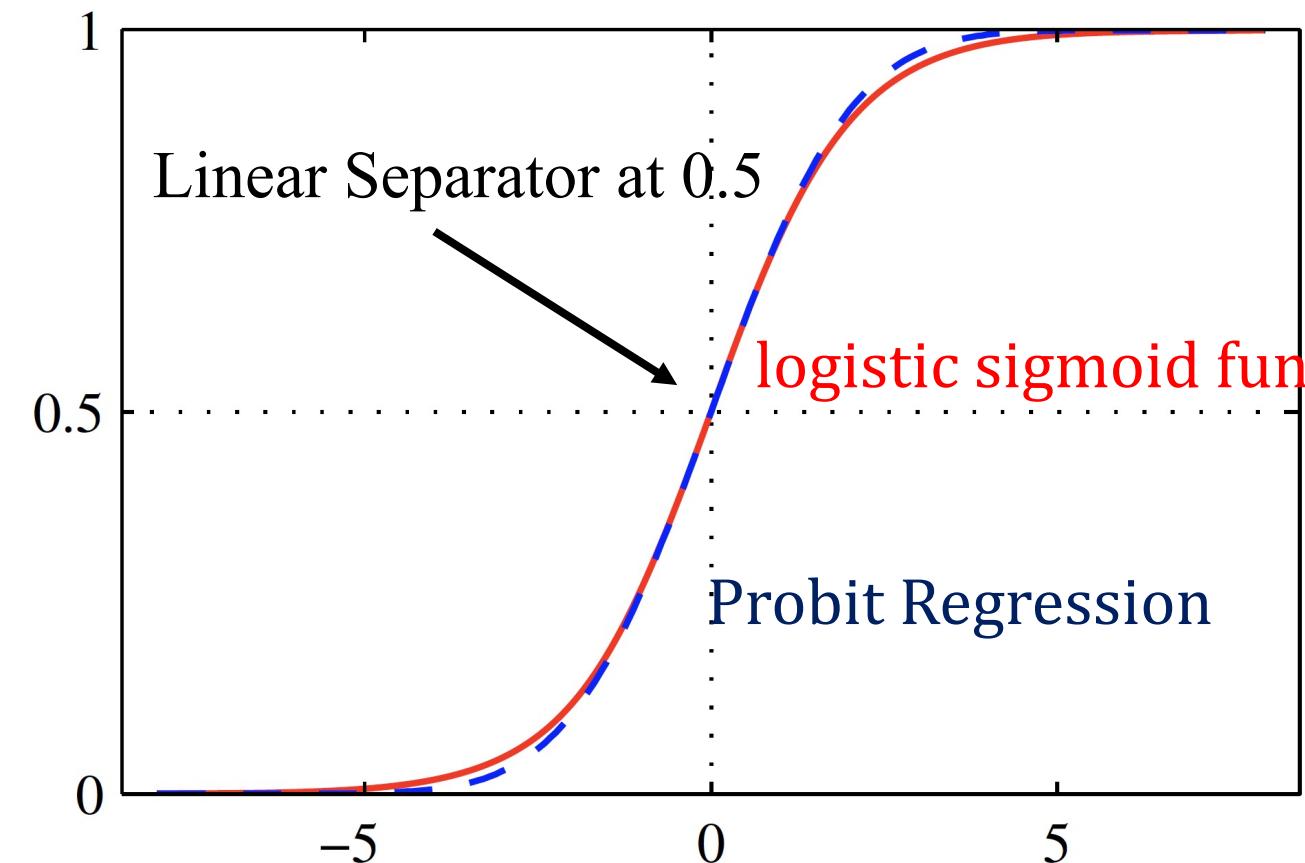
$$\begin{aligned} &= \frac{1}{1 + \pi_j \exp \left\{ \boldsymbol{\mu}_j^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_j^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_j \right\} / \pi_k \exp \left\{ \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_k \right\}} \\ &= \frac{1}{1 + \exp \left\{ - \underbrace{(\boldsymbol{\mu}_k^T - \boldsymbol{\mu}_j^T) \boldsymbol{\Sigma}^{-1} \mathbf{x}}_{= \mathbf{w}^T} - \underbrace{\frac{1}{2} \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_k + \frac{1}{2} \boldsymbol{\mu}_j^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_j + \ln \left( \frac{\pi_k}{\pi_j} \right)}_{= w_0} \right\}} \end{aligned} \quad (6)$$

Equation (6) simply becomes

$$= \frac{1}{1 + \exp \{-(\mathbf{w}^T \mathbf{x} + w_0)\}} \quad \text{the logistic sigmoid function} \quad (7)$$



# Logistic sigmoid function



The logistic sigmoid function  $\sigma(a)$  is

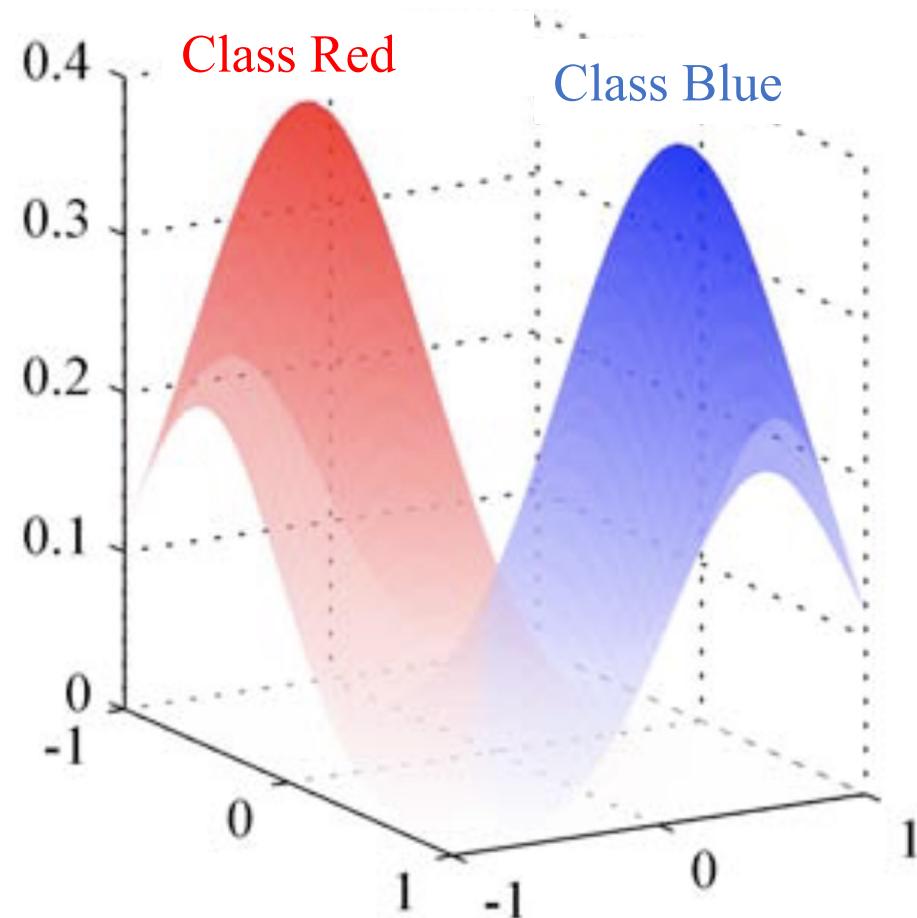
$$\sigma(a) = \frac{1}{1 + e^{-a}}$$

Then  $p(C_k | \mathbf{x}) = \sigma(-\mathbf{w}^T \mathbf{x} + w_0)$

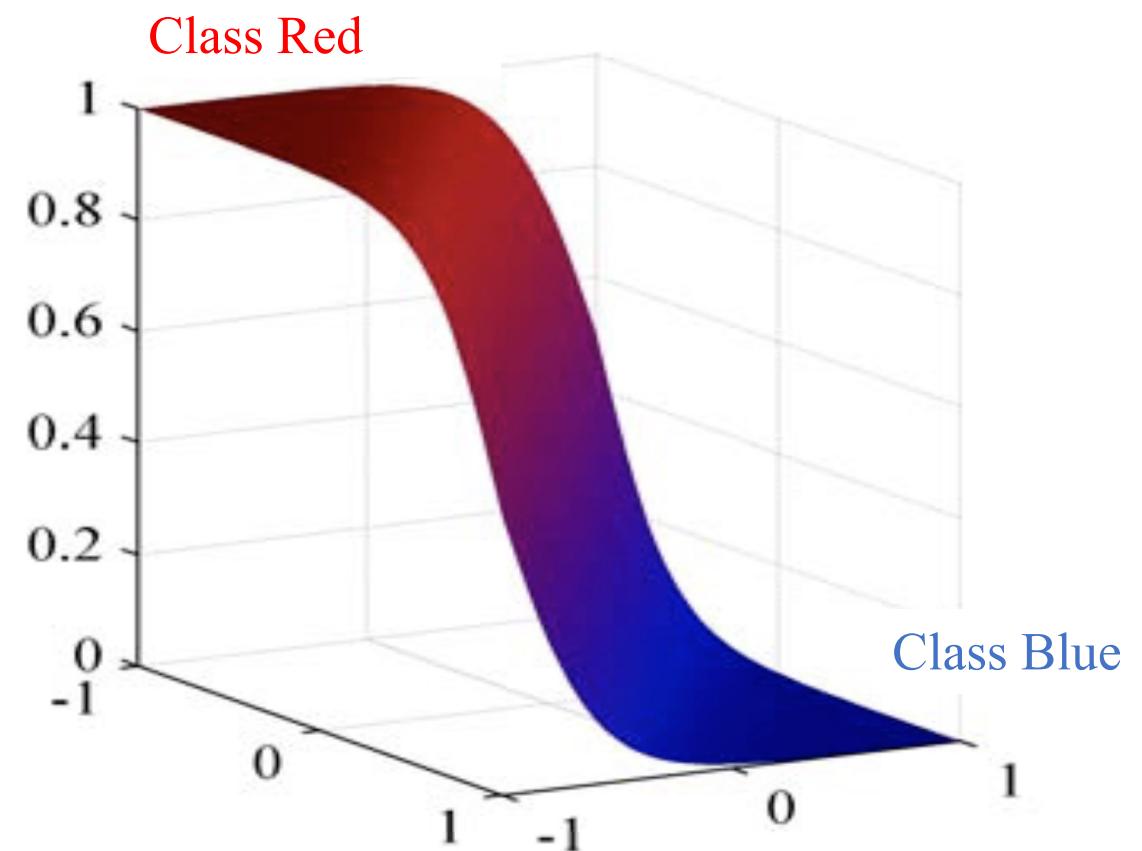
# Logistic sigmoid function

Objective: find the probability of inputs likely being classified as Red.

Class Conditional

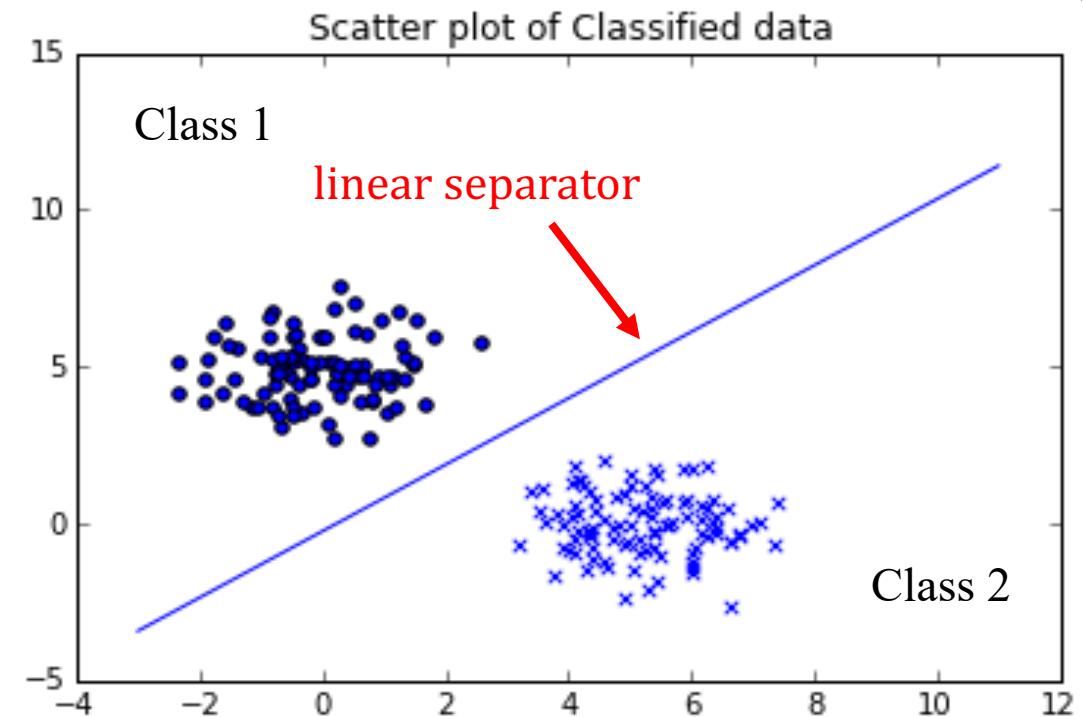


Posterior



# Two Class Prediction

- best class =  $\operatorname{argmax}_k P(C_k | \mathbf{x})$
- $= \begin{cases} c_1 & \sigma(a) \geq 0.5 \\ c_2 & \text{otherwise} \end{cases}$
- Class Boundary:  $\sigma(\mathbf{w}_k^T \mathbf{x} + w_0) = 0.5$   
 $\Rightarrow e^{-(\mathbf{w}^T \mathbf{x} + w_0)} = 1$   
 $\Rightarrow \mathbf{w}^T \mathbf{x} + w_0 = 0$   
 $\therefore$  linear separator





# Posterior Distribution for Multi-Class: $k > 2$

- The normalized exponential:

$$p(C_k | \mathbf{x}) = \frac{\pi_k \exp \left\{ \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_k \right\}}{\sum_k \pi_k \exp \left\{ \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_k \right\}} \quad (4-3)$$

$$= \frac{\exp \left\{ \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_k + \ln \pi_k \right\}}{\sum_j \exp \left\{ \boldsymbol{\mu}_j^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_j^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_j + \ln \pi_j \right\}}$$

$$\text{softmax} \Rightarrow = \frac{e^{\mathbf{w}_k^T \bar{\mathbf{x}}}}{\sum_j e^{\mathbf{w}_j^T \bar{\mathbf{x}}}} \quad (8)$$

where  $\mathbf{w}_k = \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}^{-1}$  and  $w_{0k} = -\frac{1}{2} \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_k + \ln \pi_k$



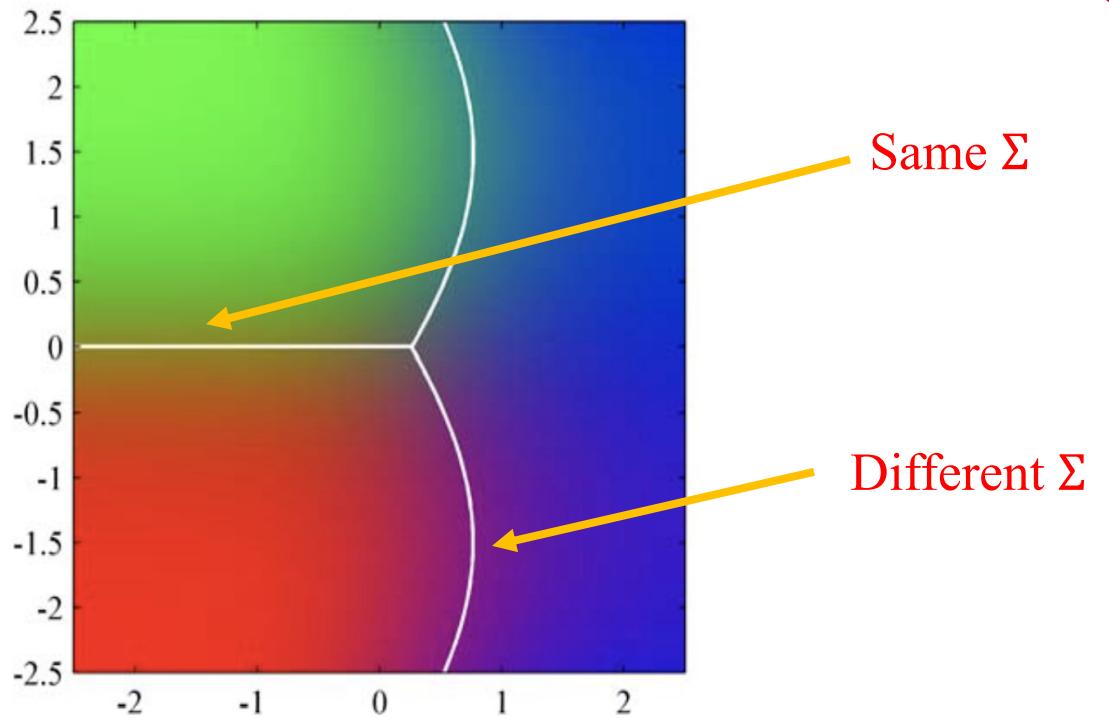
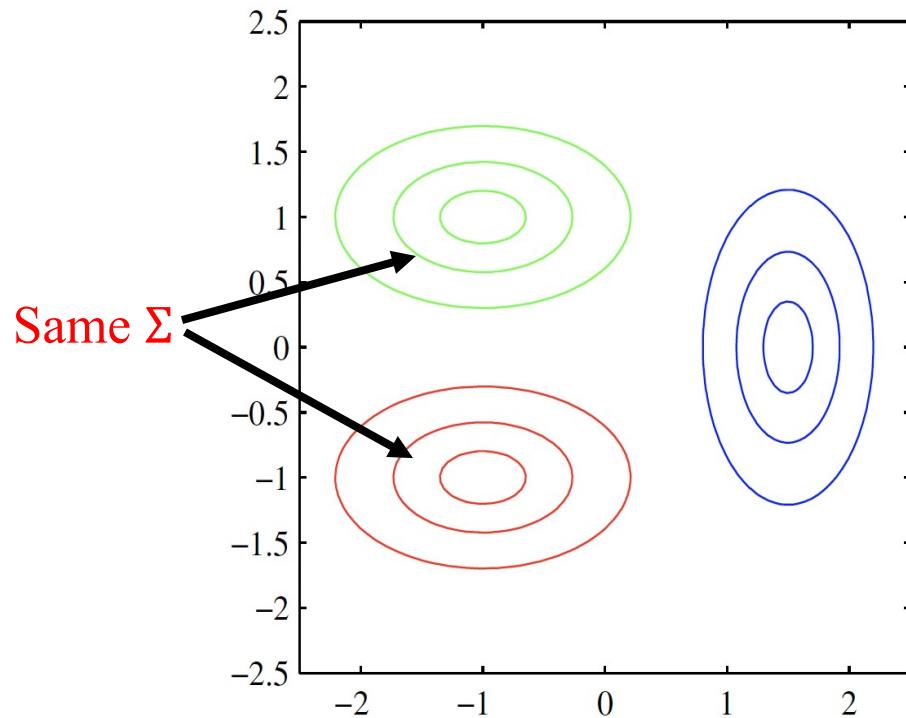
- The posterior is a softmax (generalization of the sigmoid)
- softmax distribution:

$$p(C_k | x) = \frac{\exp(f_k(x))}{\sum_j \exp(f_j(x))}$$

- argmax distribution

$$\begin{aligned} p(C_k | x) &= \begin{cases} 1 & \text{if } k = \operatorname{argmax}_j f_j(x) \\ 0 & \text{otherwise} \end{cases} \\ &= \lim_{base \rightarrow \infty} \frac{base^{f_k(x)}}{\sum_j base^{f_j(x)}} \\ &\approx \frac{\exp(f_k(x))}{\sum_j \exp(f_j(x))} \end{aligned} \tag{9}$$

# Softmax



$$p(c_k | \mathbf{x}) = \frac{\pi_k \exp \left\{ -\frac{1}{2} (\mathbf{x}^T \Sigma_k^{-1} \mathbf{x} - 2\mu_k^T \Sigma_k^{-1} \mathbf{x} + \mu_k^T \Sigma_k^{-1} \mu_k) \right\}}{\sum_j \pi_j \exp \left\{ -\frac{1}{2} (\mathbf{x}^T \Sigma_j^{-1} \mathbf{x} - 2\mu_j^T \Sigma_j^{-1} \mathbf{x} + \mu_j^T \Sigma_j^{-1} \mu_j) \right\}} \quad (4-3)$$

Depends on class number  $j$



# Parameter Estimation

- We have a parametric function form for the class-conditional densities:

$$p(x|C_k) = \frac{1}{(2\pi)^{D/2}} \left( \frac{1}{|\Sigma|^{1/2}} \right) \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\} \quad (3)$$

- We can estimate the parameters and the prior class probabilities using maximum likelihood.

- Two class case with shared covariance matrix.

- Training data:

- $\{x_n, y_n\}, n = 1, \dots, N$

- $y_n = 1$  denotes class  $C_1$ ;  $y_n = 0$  denotes class  $C_2$ ;

- Priors:  $p(C_1) = \pi, p(C_2) = 1 - \pi$

- For a data point  $x_n$  from class  $C_1$ , we have  $y_n = 1$  and therefore:

$$p(\mathbf{x}_n, C_1) = p(\mathbf{x}|C_1)p(C_1) = \pi \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}) \quad (9)$$

- For a data point  $x_n$  from class  $C_2$ , we have  $y_n = 0$  and therefore:

$$p(\mathbf{x}_n, C_2) = p(\mathbf{x}|C_2)p(C_2) = (1 - \pi) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) \quad (10)$$



# Maximum Likelihood Solution

- Assuming observations are drawn independently, the likelihood function is as below:

$$L(\mathbf{x}, \mathbf{y}) = p(\mathbf{x}, \mathbf{y} | \pi, \mu_1, \mu_2, \Sigma) = \prod_{n=1}^N [p(\mathbf{x}_n, C_1)]^{y_n} [p(\mathbf{x}_n, C_2)]^{1-y_n} \quad (11-1)$$

Very common step

$$= \prod_{n=1}^N [\pi \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \Sigma)]^{y_n} [(1 - \pi) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \Sigma)]^{1-y_n} \quad (11-2)$$

↓

$$\ln L(\mathbf{x}, \mathbf{y}) = \sum_{n=1}^N [y_n \ln \pi + y_n \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \Sigma) + (1 - y_n) \ln(1 - \pi) - (1 - y_n) \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \Sigma)] \quad (12)$$



# Maximum Likelihood Solution - parameter $\pi$

- We first maximize the log-likelihood with respect to  $\pi$  (set derivative to 0)

$$\frac{\partial}{\partial \pi} \left( \sum_{n=1}^N [y_n \ln \pi + (1 - y_n) \ln(1 - \pi)] \right) = 0 \quad (13)$$

- The maximum likelihood estimate of  $\pi$  is the fraction of points in class  $C_1$ .
- For multi-class: maximum likelihood estimate for  $p(C_k)$  is given by the fraction of points in the training set in  $C_k$ .

$$\Rightarrow \pi = \frac{1}{N} \sum_{n=1}^N y_n = \frac{N_1}{N} = \frac{N_1}{N_1 + N_2} \quad (14)$$

# Maximum Likelihood Solution - parameter $\mu$

- We then maximize the log-likelihood with respect to  $\mu_1$  (set derivative to 0)

$$\frac{\partial}{\partial \mu_1} \left( \sum_{n=1}^N y_n \ln \mathcal{N}(\mathbf{x}_n | \mu_1, \Sigma) \right) = 0$$

$$\frac{\partial}{\partial \mu_1} \left( \frac{1}{2} \sum_{n=1}^N y_n (\mathbf{x}_n - \mu_1)^T \Sigma^{-1} (\mathbf{x}_n - \mu_1) + \dots \right) = 0 \quad (15)$$

$$\mu_1 = \frac{1}{N_1} \sum_{n=1}^N y_n \mathbf{x}_n \quad (16)$$

- The maximum likelihood estimate of  $\mu_1$  is the sample mean of all input  $x_n$  in class C1.
- The maximum likelihood estimate of  $\mu_2$  is:

$$\mu_2 = \frac{1}{N_2} \sum_{n=1}^N (1 - y_n) \mathbf{x}_n \quad (17)$$



# Maximum Likelihood Solution - parameter $\Sigma$

Maximize the log-likelihood w.r.t.  $\Sigma$ , we obtain

$$\Sigma = \frac{N_1}{N} S_1 + \frac{N_2}{N} S_2 \quad (18)$$

where

$$S_1 = \frac{1}{N_1} \sum_{n \in C_1} (x_n - \mu_1)(x_n - \mu_1)^T \quad (19)$$

$$S_2 = \frac{1}{N_2} \sum_{n \in C_2} (x_n - \mu_2)(x_n - \mu_2)^T \quad (20)$$



# Exponential Family

- Exponential family members - e.g., Gaussian, Bernoulli, Poisson, Beta, Dirichlet, Gamma, etc...

- For all, the posterior is in the format of:

$$P(x|\theta_k) = \exp\left(\theta_k^T T(x) - A(\theta_k) + B(x)\right) \quad (21)$$

- The posterior is a sigmoid logistic linear function x

$$p(c_k|x) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$



- Probabilistic Generative Model - Learn prior and posterior by maximum likelihood and find posterior using Bayesian Theorem.
- The posterior is known - either logistic sigmoid or softmax
- Probabilistic Discriminative Model - Learn posterior directly by maximum likelihood.



# Logistic Regression

$$L(\mathbf{x}, \mathbf{y}) = p(\mathbf{x}, \mathbf{y} | \pi, \mu_1, \mu_2, \Sigma) = \prod_{n=1}^N [\pi \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \Sigma)]^{y_n} [(1 - \pi) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \Sigma)]^{1-y_n} \quad (11-2)$$

- The log-likelihood function is

$$\ln L(\mathbf{x}, \mathbf{y}) = \sum_{n=1}^N [y_n \ln \pi + y_n \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \Sigma) + (1 - y_n) \ln(1 - \pi) - (1 - y_n) \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \Sigma)] \quad (22)$$

- Use the negative log-likelihood function (cross-entropy error function) to find the parameters

$$E(\mathbf{w}) = -\ln L(\mathbf{w}) = -\sum_{n=1}^N [y_n \ln \sigma(\mathbf{w}^T \bar{\mathbf{x}}) + (1 - y_n) \ln (1 - \sigma(\mathbf{w}^T \bar{\mathbf{x}}))] \quad (23)$$

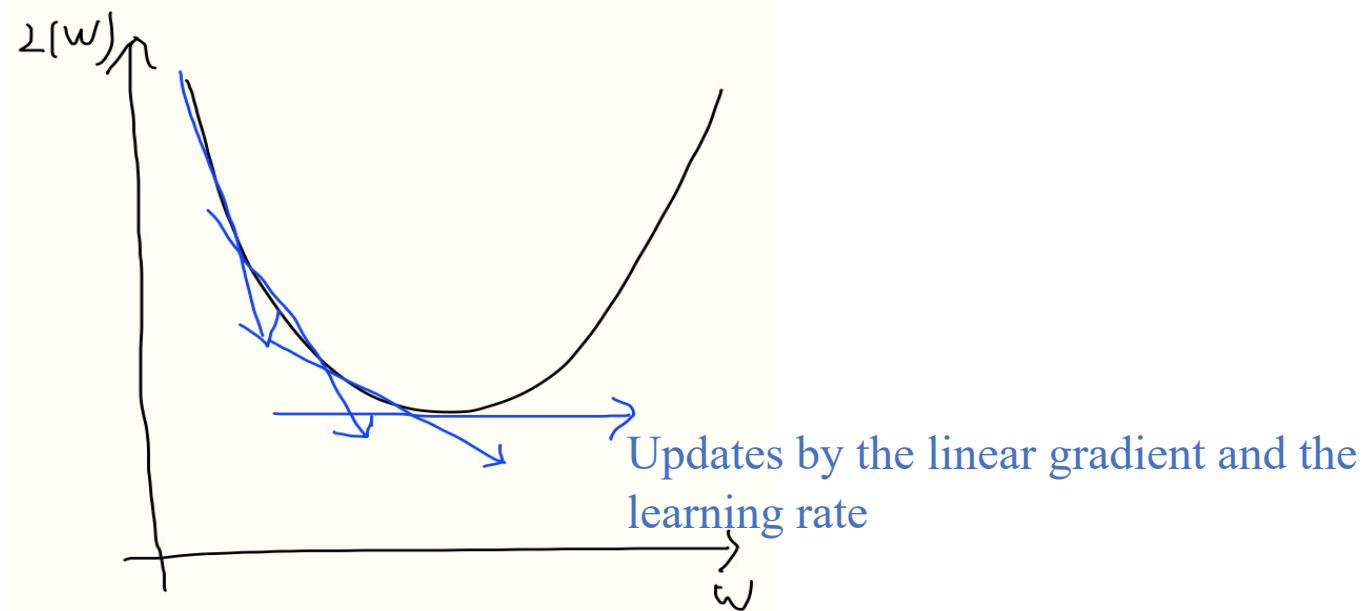
- The goal is to estimate the continuous posterior function.
- Logistic regression is a form of classification.

# Maximum Likelihood

$$-\frac{\partial E(\mathbf{w})}{\partial w} = -\sum_n y_n \left( \frac{\sigma(\mathbf{w}^T \bar{\mathbf{x}}) (1 - \sigma(\mathbf{w}^T \bar{\mathbf{x}})) \bar{\mathbf{x}}_n}{\sigma(\mathbf{w}^T \bar{\mathbf{x}})} \right) - \sum_n \frac{(1 - y_n) (1 - \sigma(\mathbf{w}^T \bar{\mathbf{x}})) \sigma(\mathbf{w}^T \bar{\mathbf{x}}) (-\bar{\mathbf{x}}_n)}{1 - \sigma(\mathbf{w}^T \bar{\mathbf{x}})} \quad (24-1)$$

$$\begin{aligned} 0 &= -\sum_n y_n \bar{\mathbf{x}}_n - \sum_n y_n \sigma(\mathbf{w}^T \bar{\mathbf{x}})(\bar{\mathbf{x}}_n) + \sum_n \sigma(\mathbf{w}^T \bar{\mathbf{x}}) \bar{\mathbf{x}}_n + \sum_n y_n \sigma(\mathbf{w}^T \bar{\mathbf{x}})(\bar{\mathbf{x}}_n) \\ &= \sum_n [\sigma(\mathbf{w}^T \bar{\mathbf{x}}) - y_n] \bar{\mathbf{x}}_n \end{aligned} \quad (24-2)$$

- Can we estimate  $\mathbf{w}$ ? No
- Then how? Use an iterative method



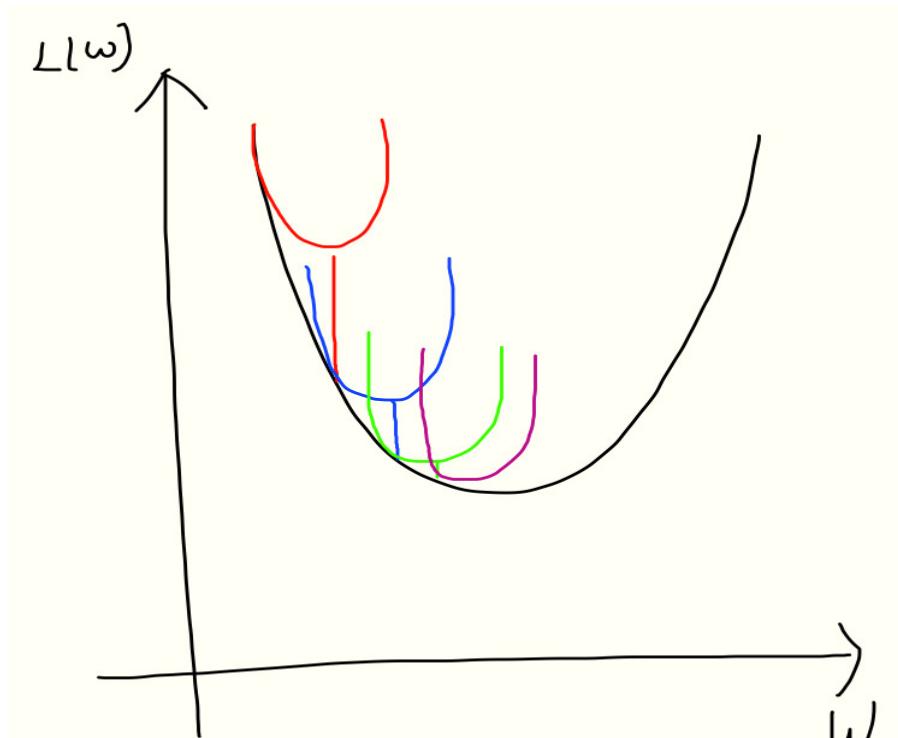
# Newton's Method

Iterative reweighted least square:

$$\mathbf{w} \leftarrow \mathbf{w} - \mathbf{H}^{-1} \nabla E(\mathbf{w}) \quad (25)$$

where  $\nabla E$  is the gradient (column vector) and  $\mathbf{H}$  is the Hessian (matrix)

$$\mathbf{H} = \begin{bmatrix} \frac{\partial^2 E}{\partial^2 w_0} & \cdots & \frac{\partial^2 E}{\partial w_0 \partial w_m} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 E}{\partial w_m \partial w_0} & \cdots & \frac{\partial^2 E}{\partial^2 w_m} \end{bmatrix} \quad (25-1)$$





# Hessian

$$\mathbf{H} = \nabla(\nabla \ln L(w)) \quad (26)$$

$$\begin{aligned} &= \nabla \left( \sum_n [\sigma(\mathbf{w}^T \bar{\mathbf{x}}) - y_n] \bar{\mathbf{x}}_n \right) \\ &= \sum_{n=1}^N \sigma(\mathbf{w}^T \bar{\mathbf{x}}_n) \left( 1 - \sigma(\mathbf{w}^T \bar{\mathbf{x}}_n) \right) \bar{\mathbf{x}}_n \bar{\mathbf{x}}_n^T \\ &= \bar{\mathbf{X}} \mathbf{R} \bar{\mathbf{X}}^T \end{aligned} \quad (27)$$

$$\text{where } \mathbf{R} = \begin{bmatrix} \sigma_1(1 - \sigma_1) & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_N(1 - \sigma_N) \end{bmatrix}$$



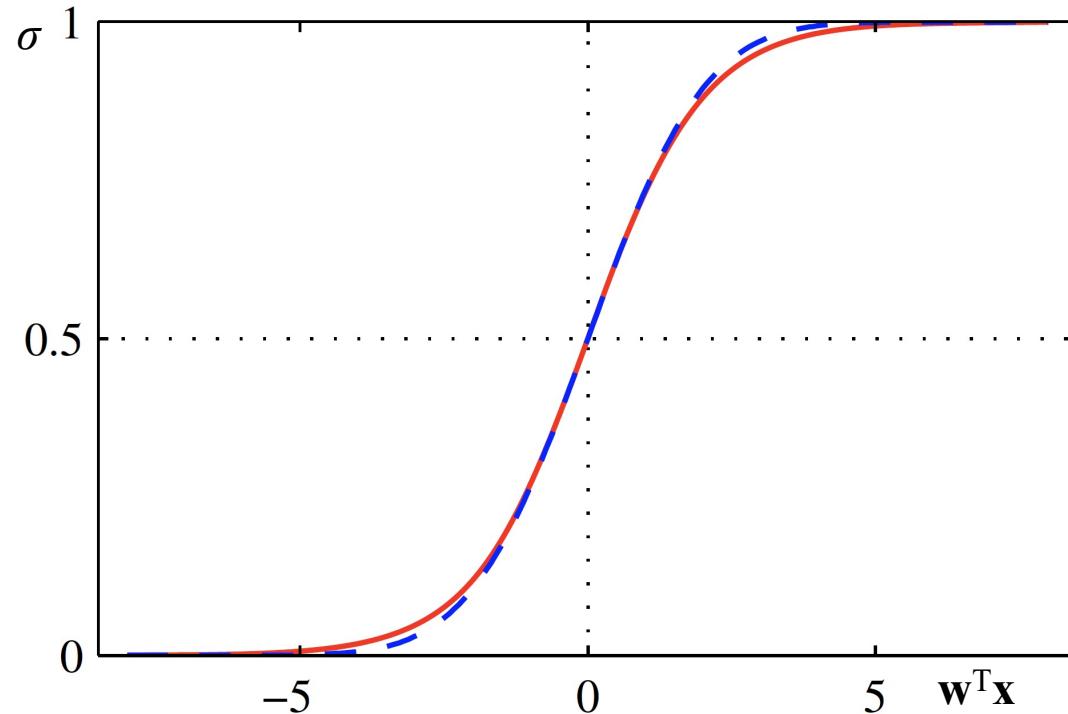
# Issues

Issue: Maximum likelihood can exhibit overfitting: logistic regression can classify each data point arbitrarily well when data is small.

## Reasons:

As the posterior gets toward 1,  $w^T x \rightarrow \infty$  and therefore the magnitude of w must be infinite.

If  $\sigma=1$ , then  $(1-\sigma)=0$  so as  $R=0$  and therefore H becomes singular.





# Regularization

- We can penalize large weights
- How? Add the penalty term  $\lambda$ .

$$\min_w E(w) + \frac{1}{2} \lambda \|w\|^2 \quad (28)$$

$$= \min_w \left\{ - \sum_{n=1}^N \left[ y_n \ln \sigma(\mathbf{w}^T \bar{\mathbf{x}}) + (1 - y_n) \ln (1 - \sigma(\mathbf{w}^T \bar{\mathbf{x}})) \right] \right\} + \frac{1}{2} \lambda \mathbf{w}^T \mathbf{w} \quad (28-1)$$

- Hessian:  $H = \bar{\mathbf{X}} \mathbf{R} \bar{\mathbf{X}}^T + \lambda \mathbf{I}$  (29)
- The term  $\lambda \mathbf{I}$  ensures that  $H$  is not singular ( $\text{eigenvalues} \geq \lambda$ )



# Estimation with Maximum A Posteriori(MAP) from Bayes Rule

Recall,

- Bayesian view: probabilities provide a quantification of uncertainty. Before observing the data, the assumptions about  $w$  are captured in the form of a prior probability distribution  $P(\mathbf{w})$ . The effect of the observed data  $D = \{(x_1, y_1), \dots, (x_N, y_N)\}$  is expressed by  $P(D|\mathbf{w})$ .
- Bayes' theorem:

$$P(\mathbf{w}|D) = \frac{P(D|\mathbf{w})P(\mathbf{w})}{P(D)}$$

- Bayes' theorem in words: posterior  $\propto$  likelihood  $\times$  prior

# Estimation with Maximum A Posteriori(MAP) from Bayes Rule

$$\mathbf{w} = \operatorname{argmax}_{\mathbf{w}} P(\mathbf{w}|(t_1, x_1), \dots, (t_n, x_n)) = \operatorname{argmax}_{\mathbf{w}} \frac{P(\mathbf{w}|(t_1, x_1), \dots, (t_n, x_n))P(\mathbf{w})}{P((t_1, x_1), \dots, (t_n, x_n))}$$

- Independence & chain rule of probability.
- $x_i$  is independent of  $\mathbf{w}$  &  $P(\cdot)$  constant and can be dropped.
- $P(t_i, x_i|\mathbf{w}) = P(t_i|x_i, \mathbf{w})P(x_i|\mathbf{w}) = P(t_i|x_i, \mathbf{w})P(x_i) = P(t_i|x_i, \mathbf{w})$

$$= \operatorname{argmax}_{\mathbf{w}} \left[ \prod_{i=1}^n P(t_i|x_i, \mathbf{w}) \right] P(\mathbf{w}) = \operatorname{argmax}_{\mathbf{w}} \sum_{i=1}^n \log[P(t_i|x_i, \mathbf{w})] + \log P(\mathbf{w})$$

- $w \sim N(0, \tau^2)$  assumption

$$= \operatorname{argmin}_{\mathbf{w}} \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i^T \mathbf{w} - t_i)^2 + \frac{1}{\tau^2} \mathbf{w}^T \mathbf{w} = \operatorname{argmin}_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n (x_i^T \mathbf{w} - y_i)^2 + \lambda \|\mathbf{w}\|^2$$

- $\lambda = \sigma^2/n\tau^2$

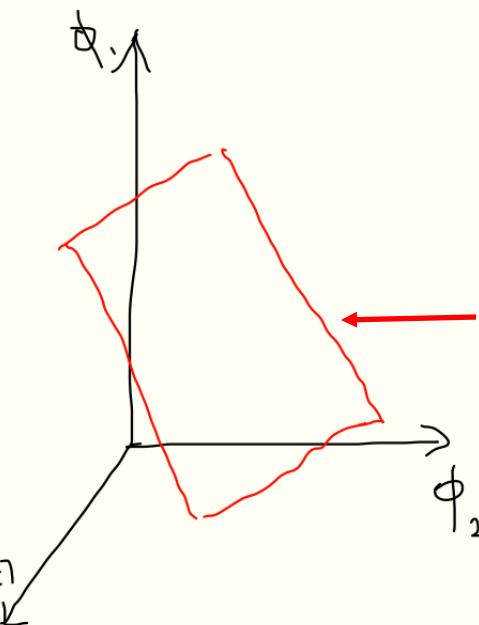
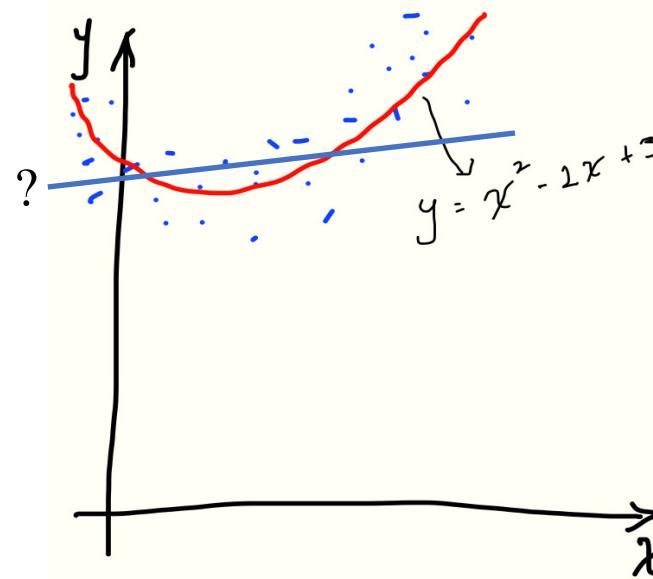
\* This objective is known as Ridge Regression.

# Generalized Linear Models



- Can we do non-linear classification?
- How can we do so? We can map inputs to a different space and so linear classification in that space.

# Basis Functions



Use non-linear basis functions. Examples of basis functions are

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

Gaussian:  $\phi_j(x) = e^{-\frac{(x-\mu_j)^2}{2s_j^2}}$

Sigmoid:  $\phi(x) = \sigma\left(\frac{x - \mu_j}{s_j}\right)$

Other many as long as they seem appropriate.

The new space we have is  $H = \{x \rightarrow \sum_i w_i \phi_i(x) \mid w_i \in \mathbb{R}\}$



# CS 559: Kernel Methods, Supportive Vector Machine, & Empirical Risk

Lecture 6

# Outline

- Kernel Methods
- SVM
- Empirical Risk



# 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.

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

We set the gradient to 0.

$$\nabla E(\mathbf{w}) = \sum_n (\mathbf{w}^T \phi(x_n) - y_n) \phi(x_n) + \lambda \mathbf{w} = 0 \quad (2)$$

$$\mathbf{w} = -\frac{1}{\lambda} \sum_n (\mathbf{w}^T \phi(x_n) - y_n) \phi(x_n) \quad (3)$$

$\mathbf{w}$  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

$$E(a) = \frac{1}{2} \mathbf{a}^T \Phi^T \Phi \Phi^T \Phi \mathbf{a} - \mathbf{a}^T \Phi^T \Phi \mathbf{y} + \frac{\mathbf{y}^T \mathbf{y}}{2} + \frac{\lambda}{2} \mathbf{a}^T \Phi^T \Phi \mathbf{a} \quad (4)$$



# Gram Matrix

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

Substitution K into (6-1):  $\frac{1}{2} \sum_{n=1}^N [w^T \phi(x_n) - y_n]^2 + \frac{1}{2} \lambda w^T w$  where  $w = \Phi a$

$$E(a) = \frac{1}{2} a^T K K a - a^T K y + \frac{y^T y}{2} + \frac{\lambda}{2} a^T K a \quad (5)$$

$$\nabla E(a) = K K a - K y + \lambda K a = 0$$

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

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

Prediction:  $y_* = \phi(x_*)^T w = \phi(x_*)^T \Phi a = K(x_*, X)(K + \lambda I)^{-1} y$  (6)

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



# Kernel Construction Example

Consider a 2-D input space  $\mathbf{x}^T = (x_1, x_2)^T$  and  $\mathbf{z} = (z_1, z_2)$  and suppose we are going to map into a feature space in the degree of 2 using

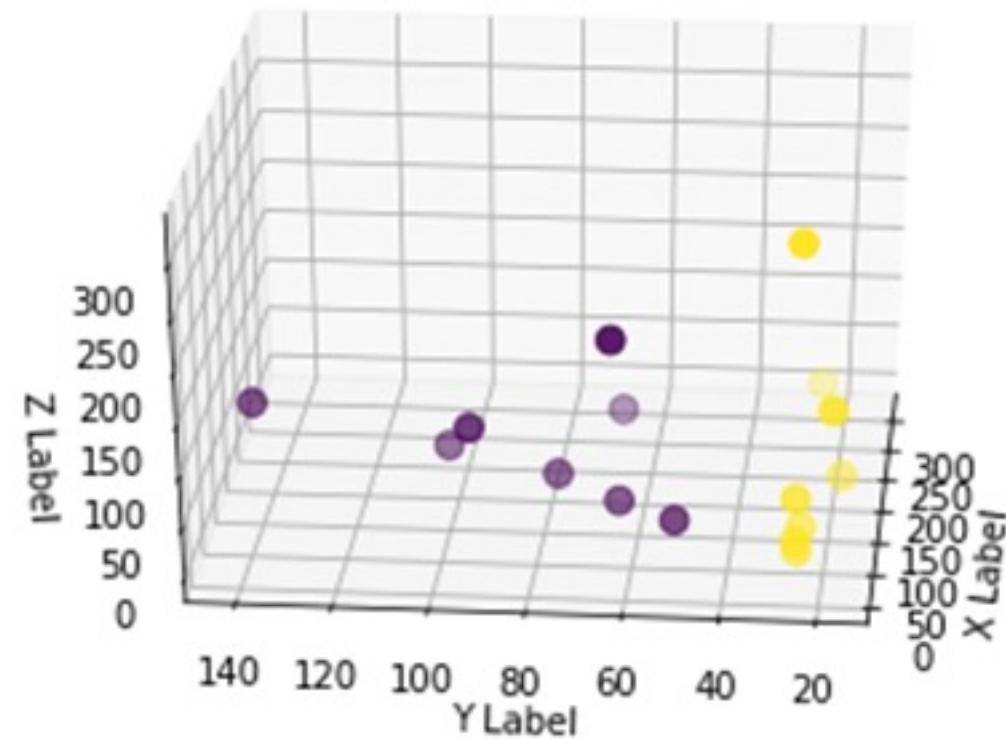
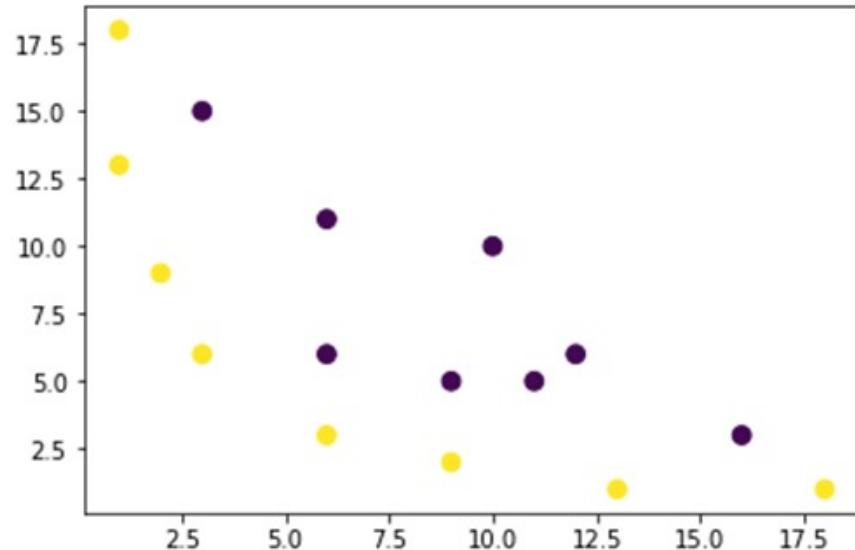
$$k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z})^2$$

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

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.

# Kernel Construction Example





# Kernel Construction Example

Verifying the validity of kernels:  $k(\mathbf{x}, \mathbf{x}') = ck_1(\mathbf{x}, \mathbf{x}')$ .

- Start with the basic function -  $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$
- It follows that  $ck_1(\mathbf{x}, \mathbf{x}') = c\mathbf{u}(\mathbf{x})^T \mathbf{u}(\mathbf{x}')$  where  $\mathbf{u}(\mathbf{x}) = \sqrt{c}\phi(\mathbf{x})$  and  $\mathbf{u}(\mathbf{x}') = \sqrt{c}\phi(\mathbf{x}')$
- So  $ck_1(\mathbf{x}, \mathbf{x}')$  can be expressed as the scalar product of feature vectors and is a valid kernel.



# Supportive Vector Machine (SVM)



# 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:  $\{(x_i, t_i) | n = 1, \dots, n\}$
- $t_i \in \{-1, 1\}$

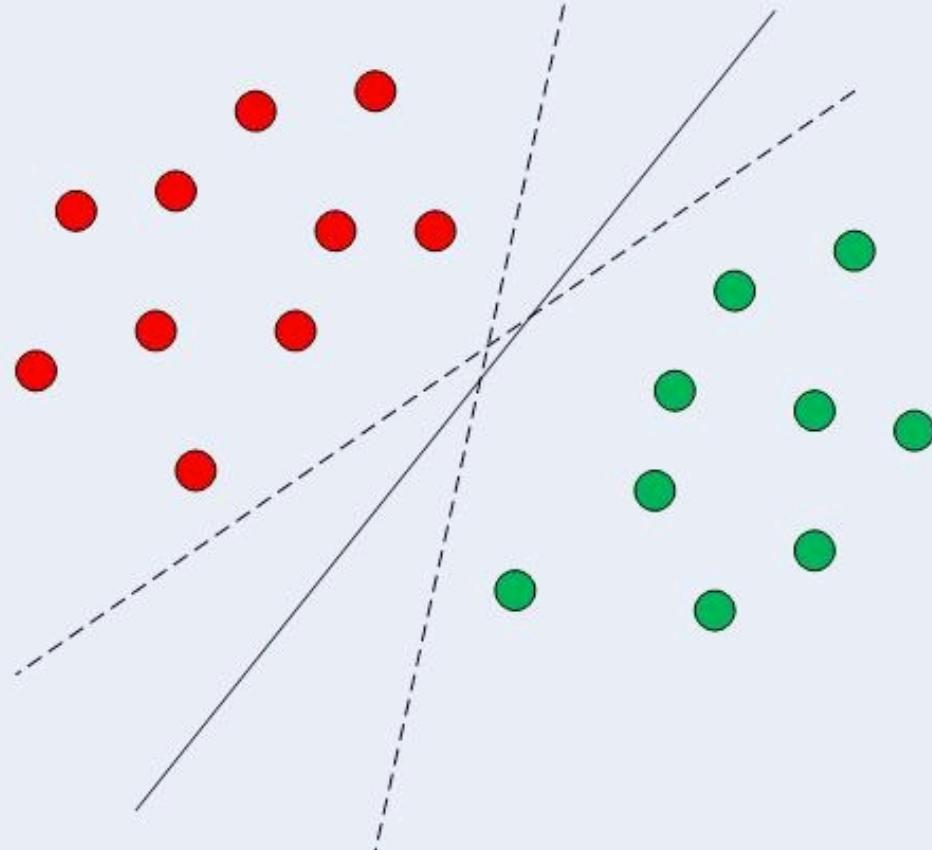
# SVM vs. Perceptron Algorithm

## Perceptron Algorithm

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

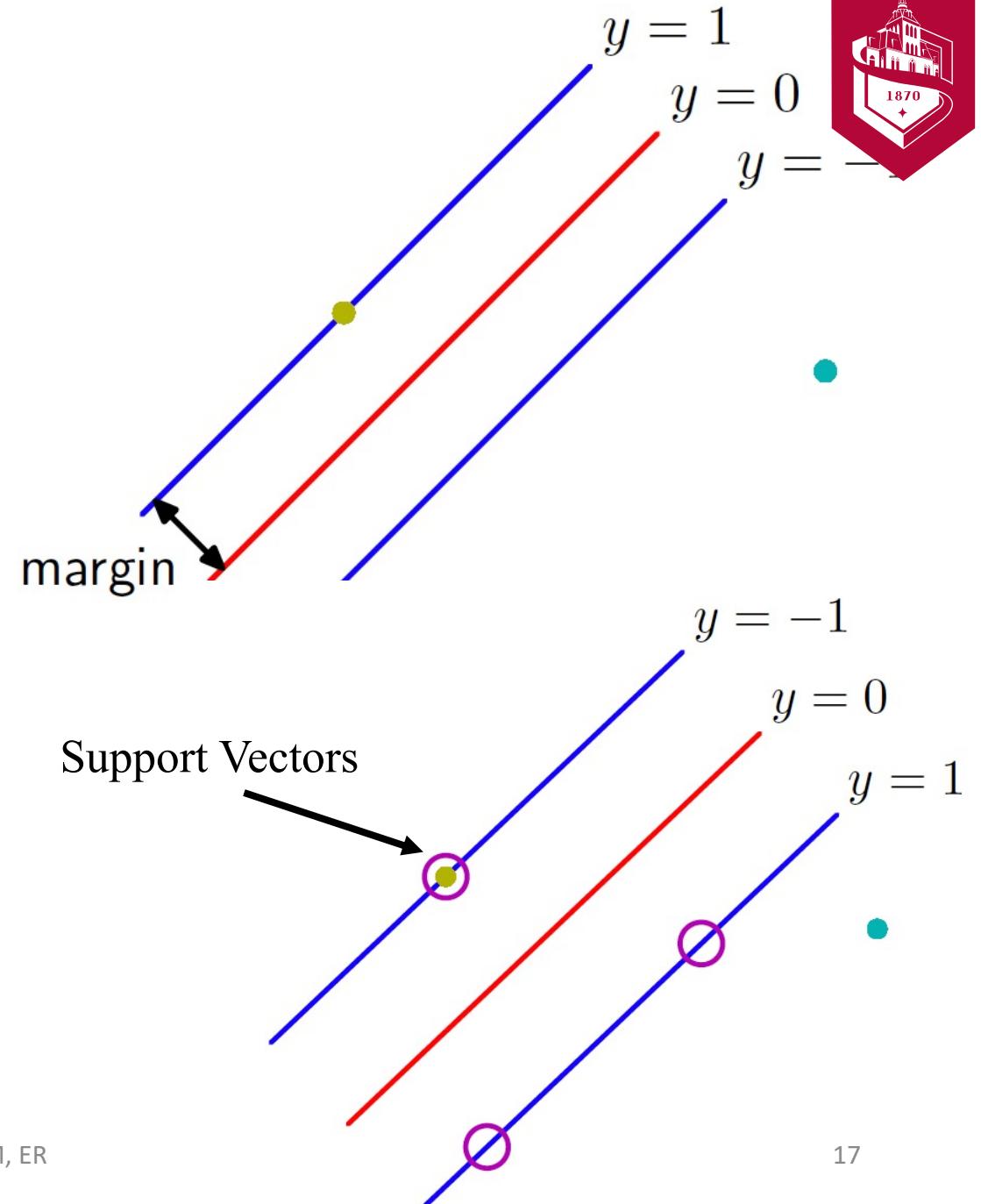
If there are multiple solutions of all of which classify the training data set exactly, then we should try to find the one that give **the smallest generalization error!**

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  $t_n y(x_n) > 0$

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

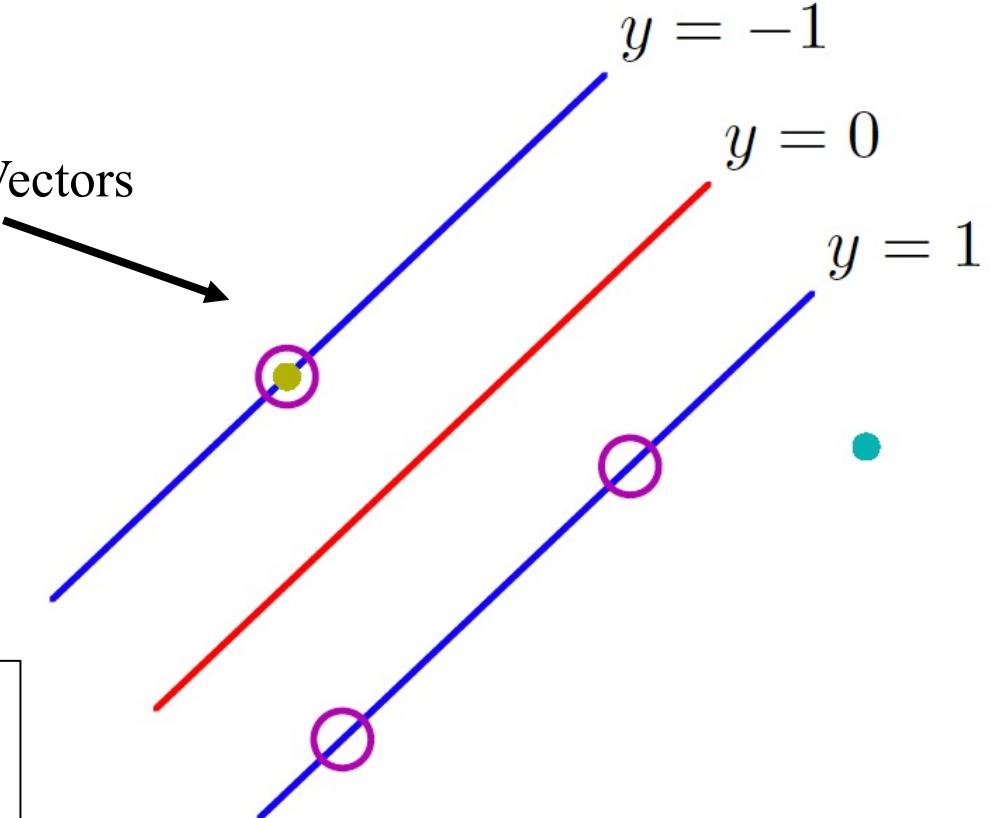
Support Vectors

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  $a_n \geq 0$

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N a_n \{t_n (\mathbf{w}^T \phi(x_n) + b) - 1\} \quad (7)$$

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

Derivative:

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



# Dual Representation

Dual representation of the maximum margin problem

$$L(a) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{m=1}^N \sum_{n=1}^N a_n a_m t_n t_m k(x_n, x_m) \quad (8)$$

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

$$y(\mathbf{x}) = \sum_{n=1}^N a_n t_n k(\mathbf{x}, \mathbf{x}_n) + b \quad (9)$$

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 (9)

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

$S$  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  $b$

$$b = \frac{1}{N_s} \left( \sum_{n \in S} \left( t_n - \sum_{m \in S} a_m t_m k(x_n, x_m) \right) \right) \quad (11)$$

$N_s$  is the total number of SV.

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

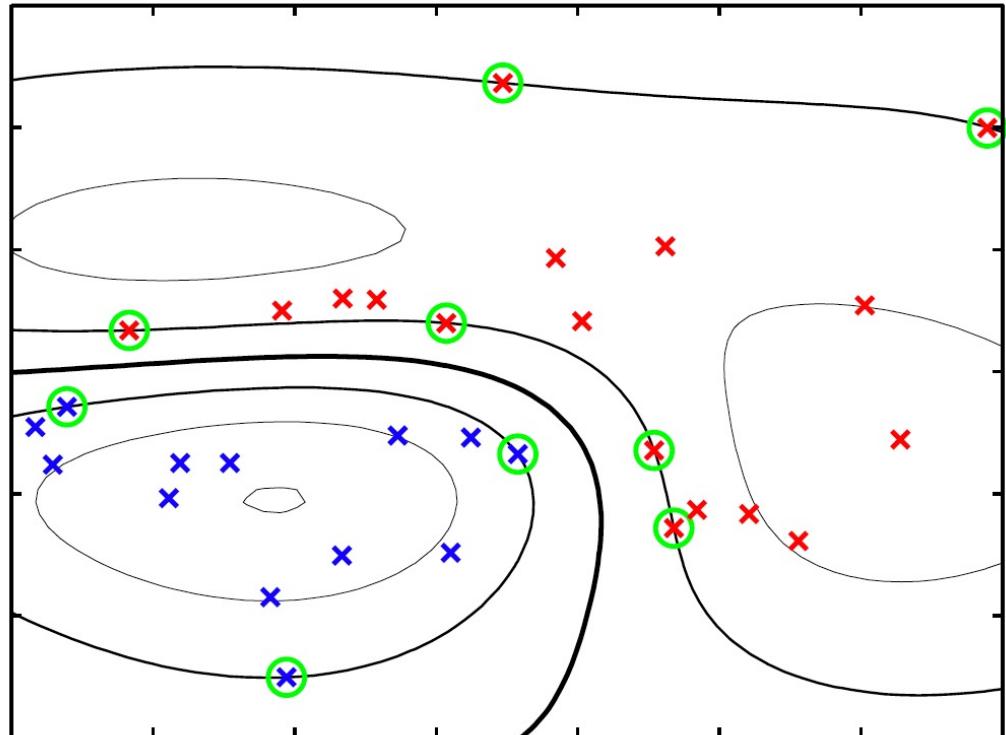
$$\sum_{n=1}^N E_\infty(y(x_n)t_n - 1) + \lambda \left\| w \right\|^2 \quad (12)$$

# Prediction

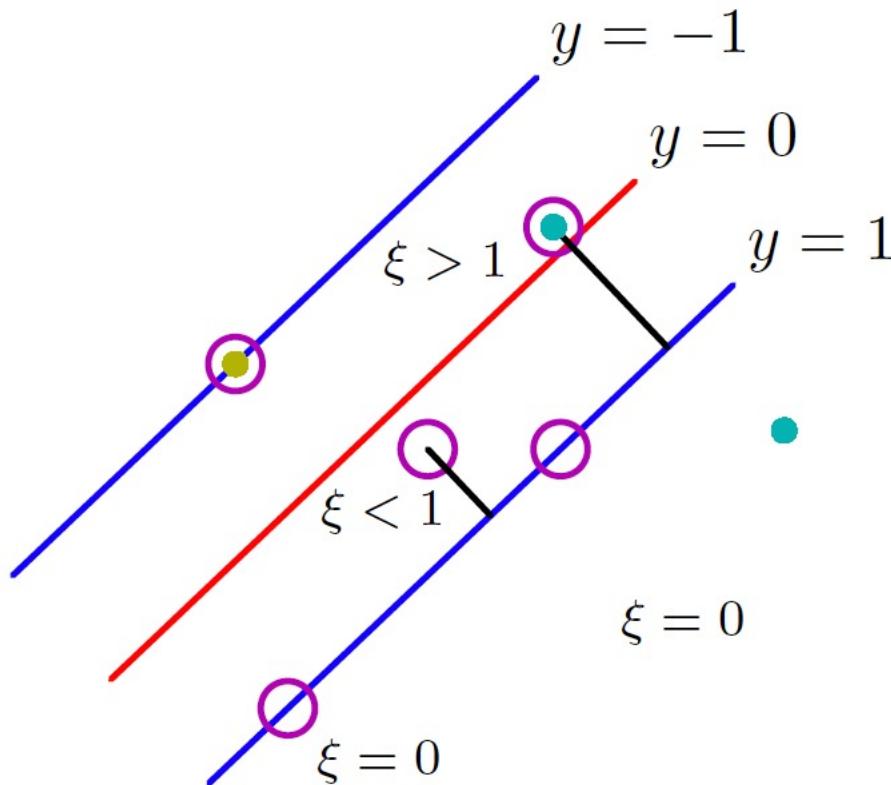
The maximum margin classifier in terms of the minimization of an error function

$$\sum_{n=1}^N E_\infty(y(x_n)t_n - 1) + \lambda \|w\|^2 \quad (13)$$

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



To minimize the margin while softly penalizing points that lie on the wrong side of the margin boundary, we need to minimize

$$C \sum_{n=1}^N \xi_n + \frac{1}{2} \|w\|^2 \quad (14)$$

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

- Misclassified points have  $\xi_n > 1$  that  $\sum_n \xi_n$  is an upper bound on the number of misclassified points!
- $C = 1/\lambda$  and if  $C \rightarrow \infty$ , we can recover SVM for separable data.

$$L(w, b, \xi, a, \mu) = \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N \xi_n + \sum_{n=1}^N a_n \{t_n y_n(x_n) - 1 + \xi_n\} - \sum_{n=1}^N \mu_n \xi_n \quad (15)$$

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

- |                                      |   |
|--------------------------------------|---|
| 1) $a_n = 0$                         | 3) $a_n \{t_n y_n(x_n) - 1 + \xi_n\} = 0$ |
| 2) $t_n y_n(x_n) - 1 + \xi_n \geq 0$ | 4) $\mu_n \geq 0$                         |
| 5) $\xi_n \geq 0$                    | 6) $\mu_n \xi_n = 0$                      |

# 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

Using these results to eliminate from the Lagrangian, the dual Lagrangian forms as below:

$$\tilde{L}(a) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{m=1}^N \sum_{n=1}^N a_n a_m t_n t_m k(x_n, x_m) \quad (20)$$

where  $a_n, \mu_n \geq 0$  implies that  $a_n \leq C$ . Then we need to maximize respect to  $\{a_n\}$  subject to

$$0 \leq a_n \leq C$$

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



# Overlapping Class Distributions – soft SVM

Solution

$$t_n \left( \sum_{m \in S} a_m t_m k(x_n, x_m) + b \right) = 1 \quad (21)$$

$$b = \frac{1}{N_M} \left( \sum_{n \in M} \left( t_n - \sum_{m \in S} a_m t_m k(x_n, x_m) \right) \right) \quad (22)$$

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

What condition needs on  $\xi_n$ ?

# 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)

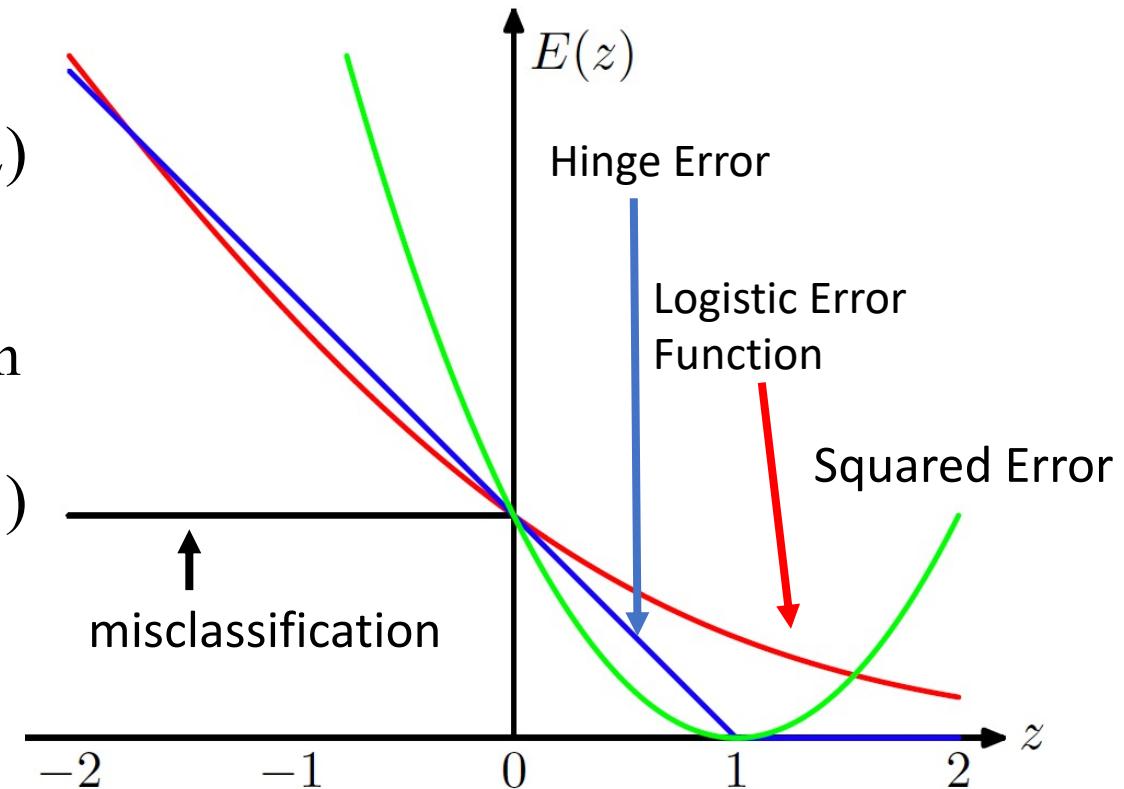
$$\sum_{n=1}^N E_{SV}(y_n t_n) + \lambda \|w\|^2$$

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

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

(6-42)

(6-43)





# 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

$$L = \frac{1}{2} \sum_{n=1}^N \{y_n - t_n\}^2 + \frac{\lambda}{2} \|w\|^2 \quad (23)$$

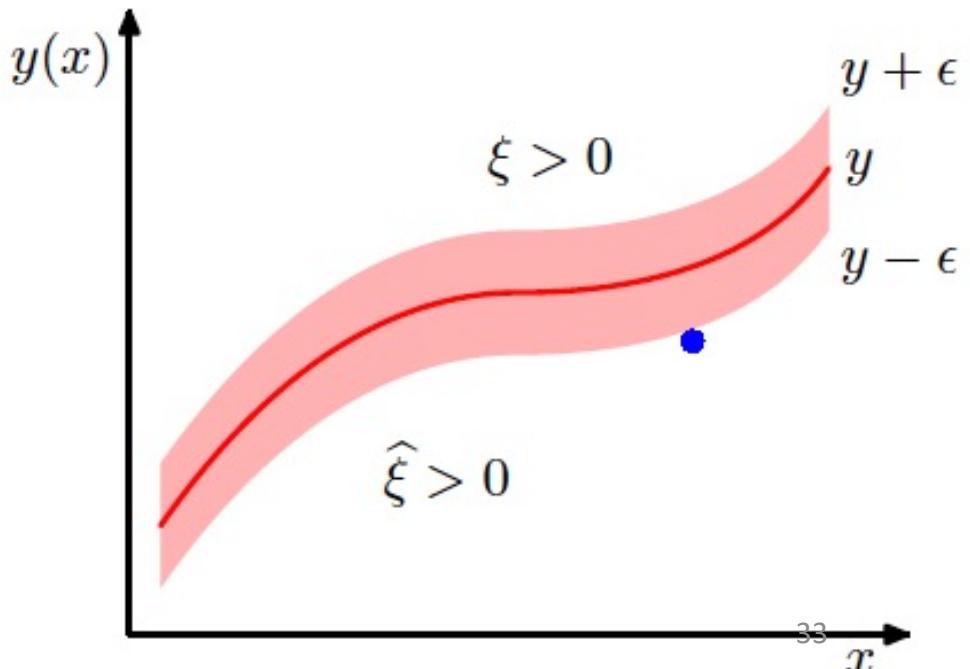
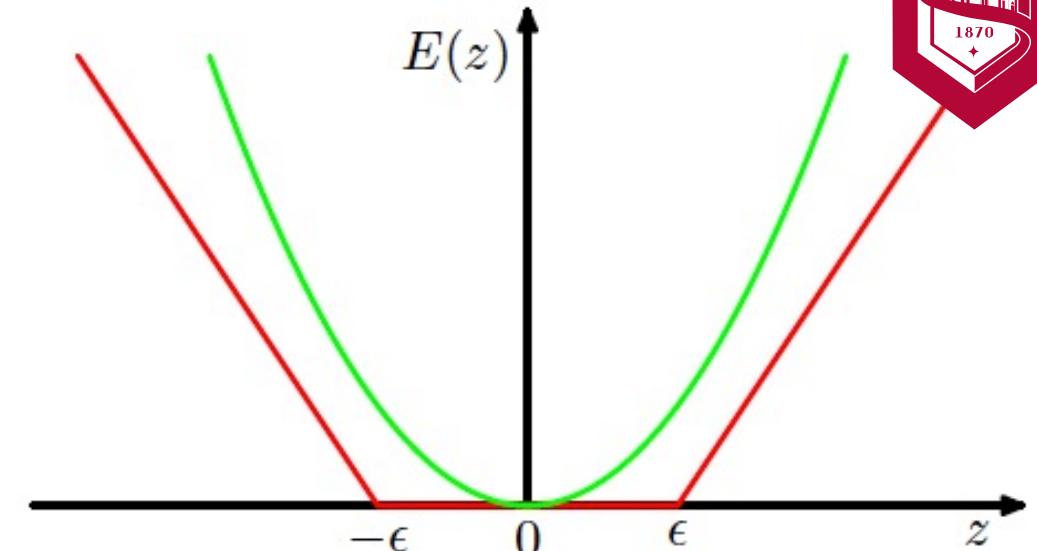
$\epsilon$ -insensitive error function (Vapnik, 1995)

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

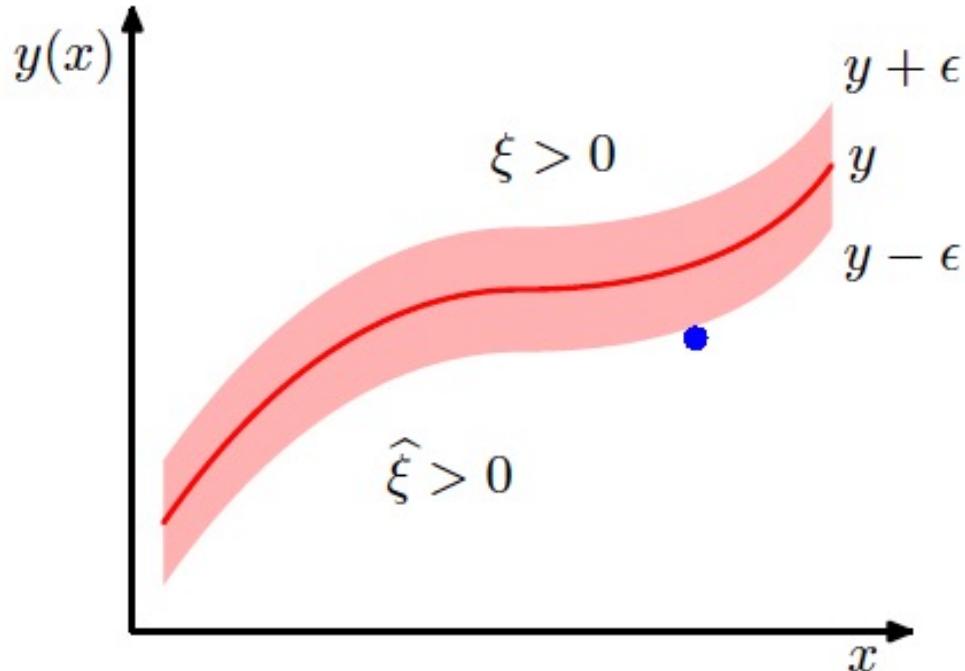
$$C \sum_{n=1}^N E_\epsilon(y_n(x_n) - t_n) + \frac{\lambda}{2} \|w\|^2 \quad (24)$$

$\xi_n > 0$  corresponds to a point  $t_n > y(x_n) + \epsilon$  and  
 $\hat{\xi}_n \geq 0$  for  $t_n < y(x_n) - \epsilon$

$$C \sum_{n=1}^N (\xi_n + \hat{\xi}_n) + \frac{1}{2} \|w\|^2 \quad (25)$$



# SVM for regression



$\xi_n > 0$  corresponds to a point  $t_n > y(x_n) + \epsilon$  and  
 $\hat{\xi}_n \geq 0$  for  $t_n < y(x_n) - \epsilon$

$$L = C \sum_{n=1}^N (\xi_n + \hat{\xi}_n) + \frac{1}{2} \|w\|^2 \quad (25)$$

Eqn 25 must be minimized subject to  $\xi_n, \hat{\xi}_n \geq 0$ .  
 How about  $t_n$ ?

$$\begin{aligned} t_n &\leq y_n + \epsilon + \xi_n \\ t_n &\geq y_n - \epsilon - \hat{\xi}_n \end{aligned}$$

How can we optimize  $L$ ?

# SVM for regression - Optimize by slack variable



$$L = C \sum_{n=1}^N (\xi_n + \hat{\xi}_n) + \frac{1}{2} \|w\|^2 - \sum_{n=1}^N (\mu_n \xi_n + \hat{\mu}_n \hat{\xi}_n) \\ - \sum_{n=1}^N a_n (\epsilon + \xi_n + y_n - t_n) - \sum_{n=1}^N \hat{a}_n (\epsilon + \hat{\xi}_n - y_n + t_n) \quad (26)$$

where  $a_n, \hat{a}_n, \mu_n, \hat{\mu}_n \geq 0$ .

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

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

$$2) \frac{\partial L}{\partial b} = 0 \Rightarrow \sum_{n=1}^N (a_n - \hat{a}_n) = 0$$

$$3) \frac{\partial L}{\partial \xi_n} = 0 \Rightarrow a_n = C - \mu_n$$

$$4) \frac{\partial L}{\partial \hat{\xi}_n} = 0 \Rightarrow \hat{a}_n = C - \hat{\mu}_n$$

# SVM for regression - Optimize by slack variable



$$\tilde{L}(a, \hat{a})$$

$$= -\frac{1}{2} \sum_{m=1}^N \sum_{n=1}^N (a_n - \hat{a}_n)(a_m - \hat{a}_m) k(\mathbf{x}_n, \mathbf{x}_m) - \epsilon \sum_{n=1}^N (a_n + \hat{a}_n) + \sum_{n=1}^N (a_n - \hat{a}_n) t_n \quad (27)$$

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

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

$$b = t_n - \epsilon - \mathbf{w}^T \phi(\mathbf{x}_n) = t_n - \epsilon - \sum_{m=1}^N (a_m - \hat{a}_m) k(\mathbf{x}_n, \mathbf{x}_m) \quad (28)$$

$$L(a, \hat{a}) = -\frac{1}{2} \sum_{m=1}^N \sum_{n=1}^N (a_n - \hat{a}_n)(a_m - \hat{a}_m) k(\mathbf{x}_n, \mathbf{x}_m) + \sum_{n=1}^N (a_n - \hat{a}_n) t_n \quad (29)$$



# 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:

$$p(t|x, w, \beta) = \mathcal{N}(t|y(x), \beta^{-1}) \quad (6-50)$$

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.

$$y(x) = \sum_{i=1}^N w_n k(x, x_n) + b \quad (6-52)$$

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  $\mathbf{X}$ . The likelihood function is

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

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

$$p(\mathbf{t}|\mathbf{X}, \boldsymbol{\alpha}, \beta) = \int p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta)p(\mathbf{w}|\boldsymbol{\alpha})d\mathbf{w} \quad (6-53)$$



# 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

$$p(\mathbf{t}|\mathbf{X}, \boldsymbol{\alpha}, \beta) = \int p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta)p(\mathbf{w}|\boldsymbol{\alpha})d\mathbf{w} \quad (6-54)$$

The log likelihood also can be formatted as

$$\ln p(\mathbf{t}|\mathbf{X}, \boldsymbol{\alpha}, \beta) = \ln \mathcal{N}(\mathbf{t}|\mathbf{0}, \mathbf{C}) - \frac{1}{2}\{N \ln(2\pi) + \ln|\mathbf{C}| + \mathbf{t}^T \mathbf{C}^{-1} \mathbf{t}\} \quad (6-55)$$

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  $\beta$ :

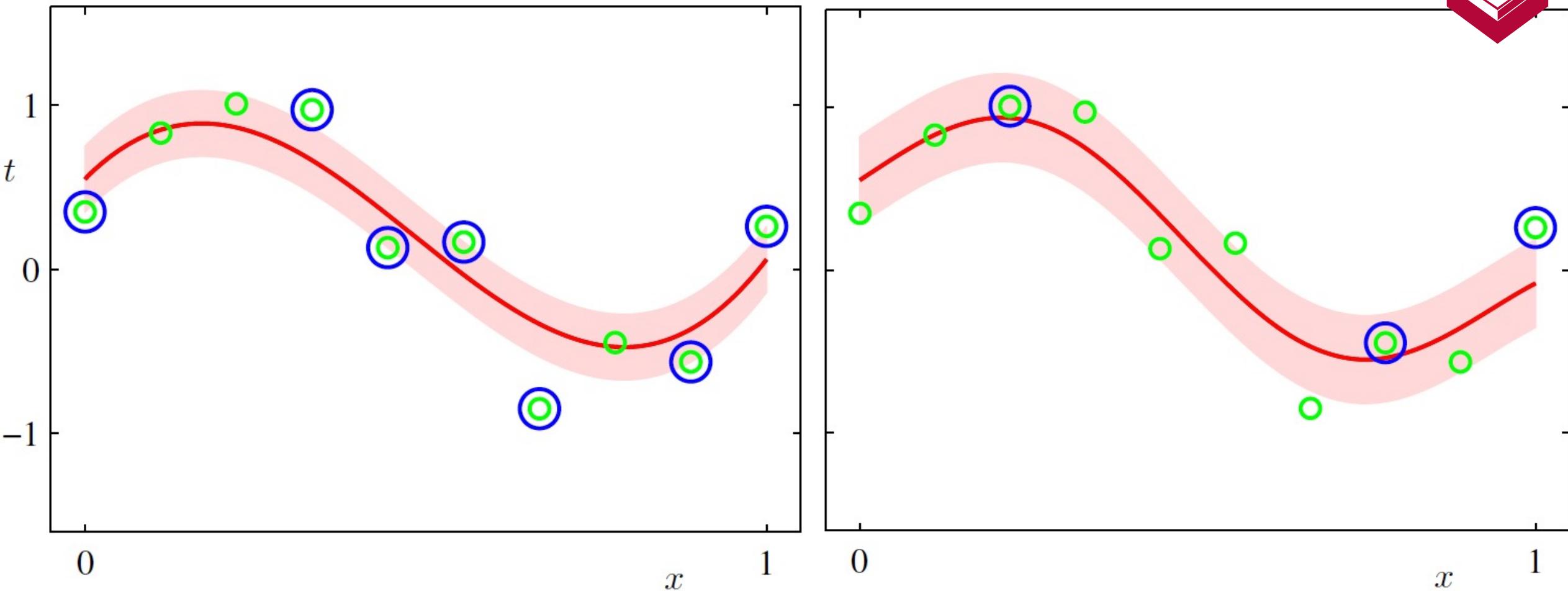
$$\alpha_i^{new} = \frac{\gamma_i}{m_i^2}$$

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

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

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

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:  $\ln p(\mathbf{w}|\mathbf{t}) = -\frac{1}{2}(\mathbf{w} - \mathbf{m}_0)^T S_0^{-1}(\mathbf{w} - \mathbf{m}_0) + \sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\} + \text{Const.}$  when the Gaussian prior of  $\mathbf{w}$  is  $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_0, S_0)$

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

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



# Empirical Risk – Loss Functions



# Empirical Risk Minimization

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

Loss                      Regularization

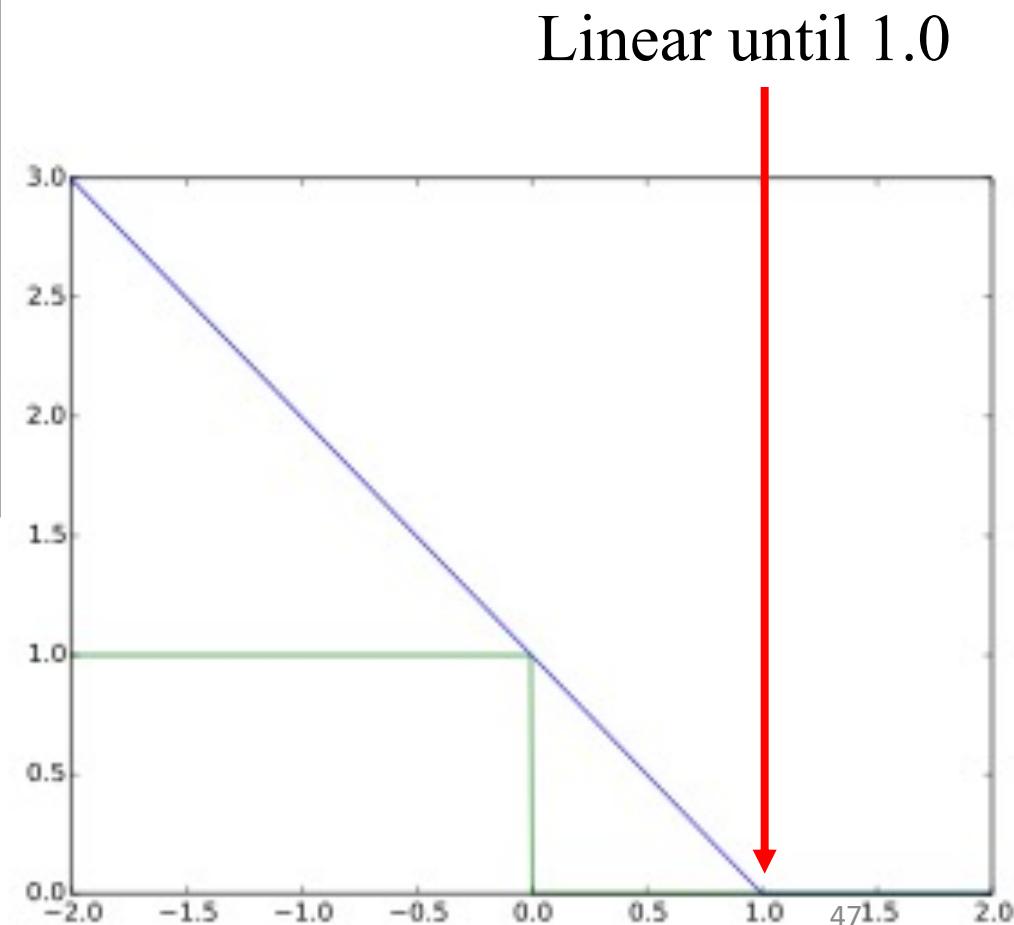
- The loss function – a continuous function penalizing training error
- The regularization – a continuous function penalizing classifier complexity



# Classification: Hinge-Loss

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

- Usage: Supportive Vector Machine
- 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  $p = 2$ .

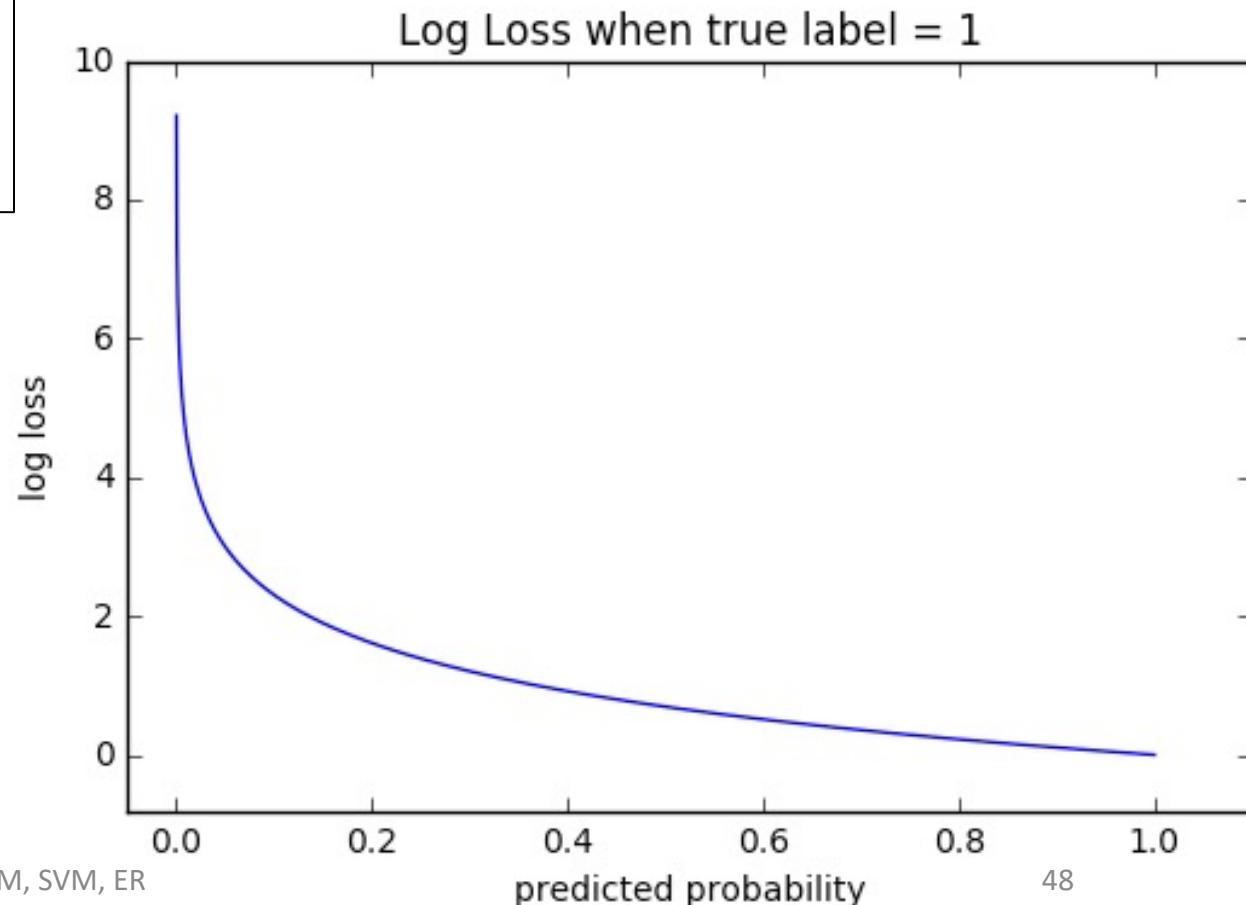




# Classification: Log-Loss

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

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

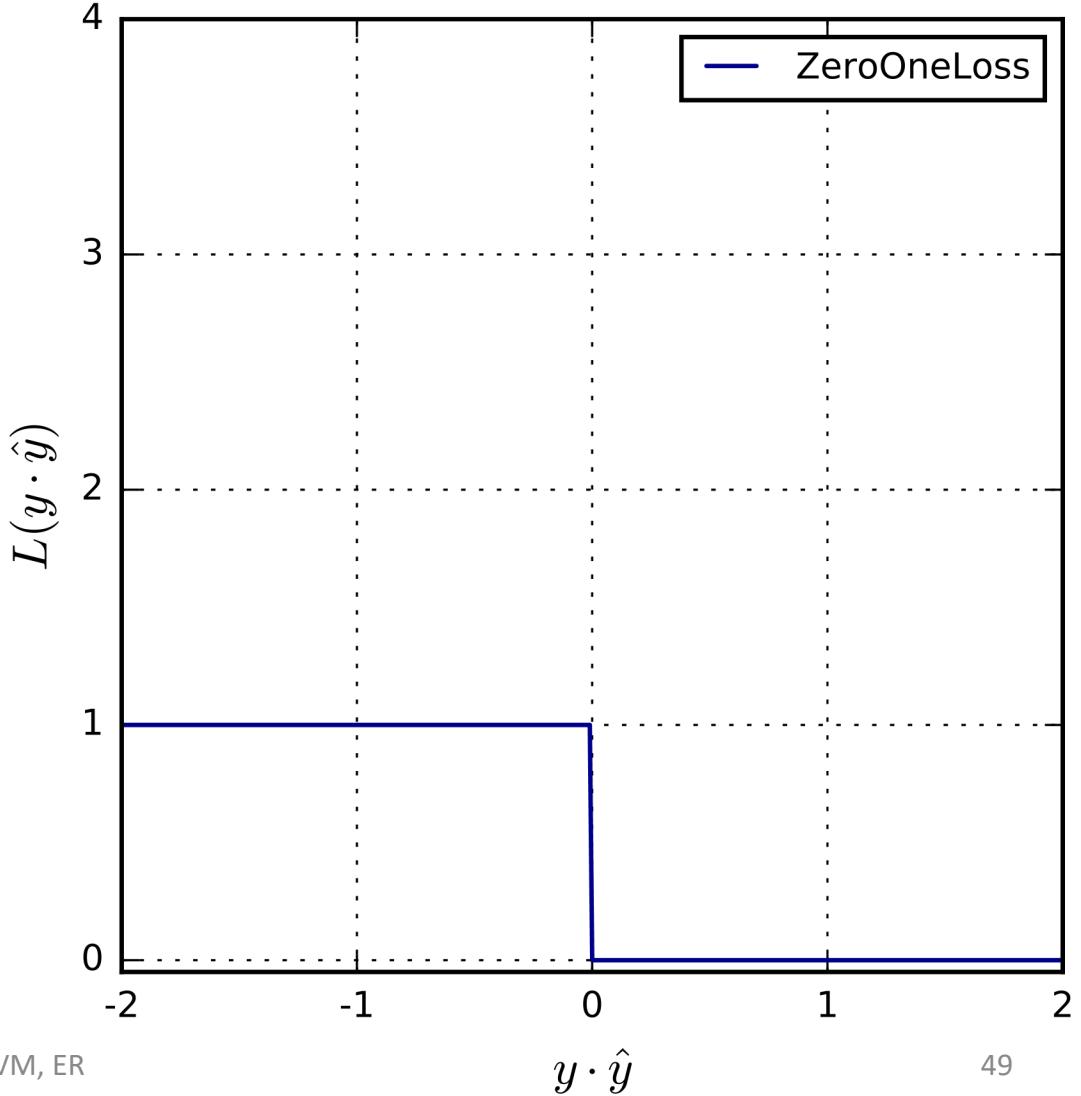




# Classification: Zero-One Loss

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

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





# 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.

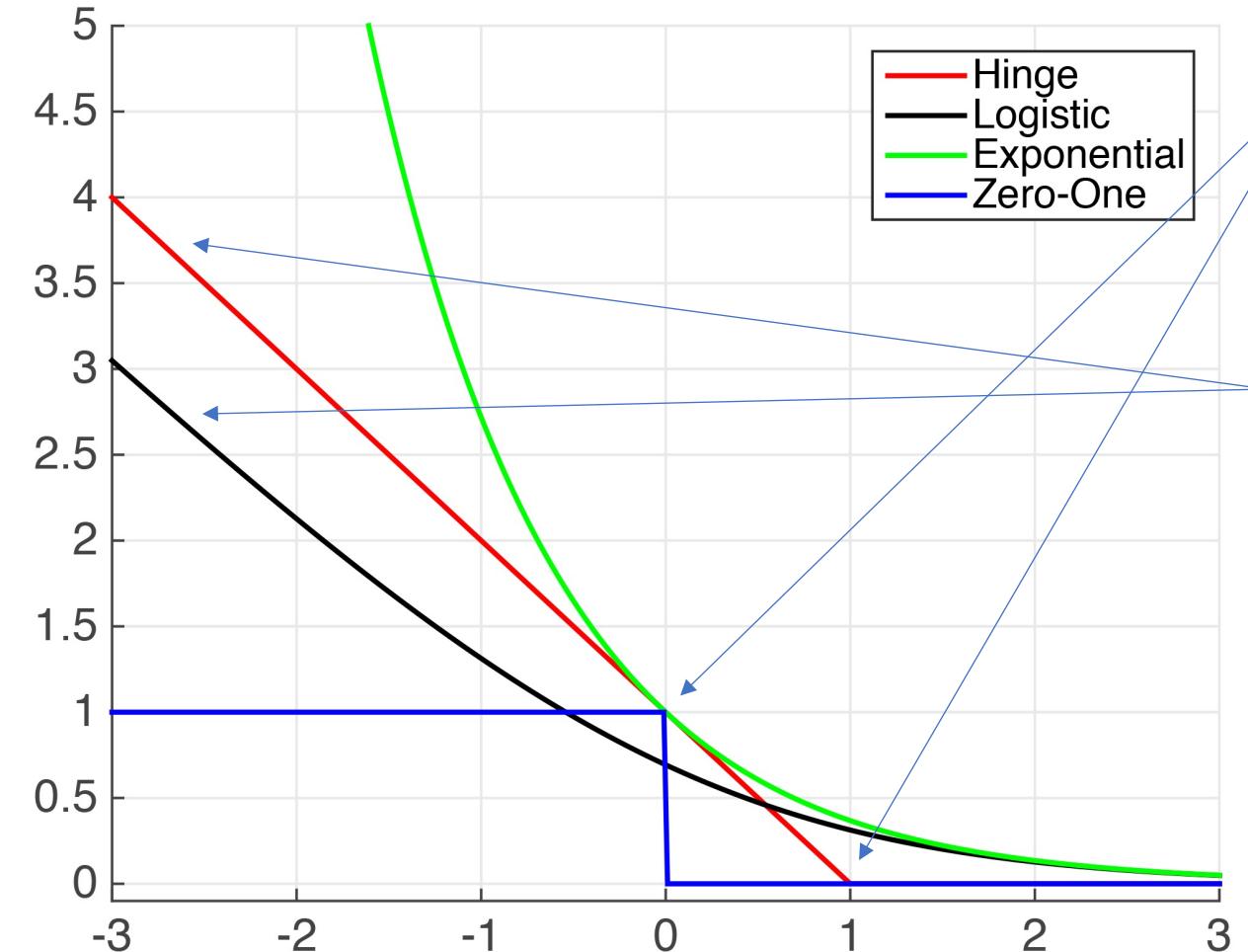
Exponential Loss:

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

- 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?

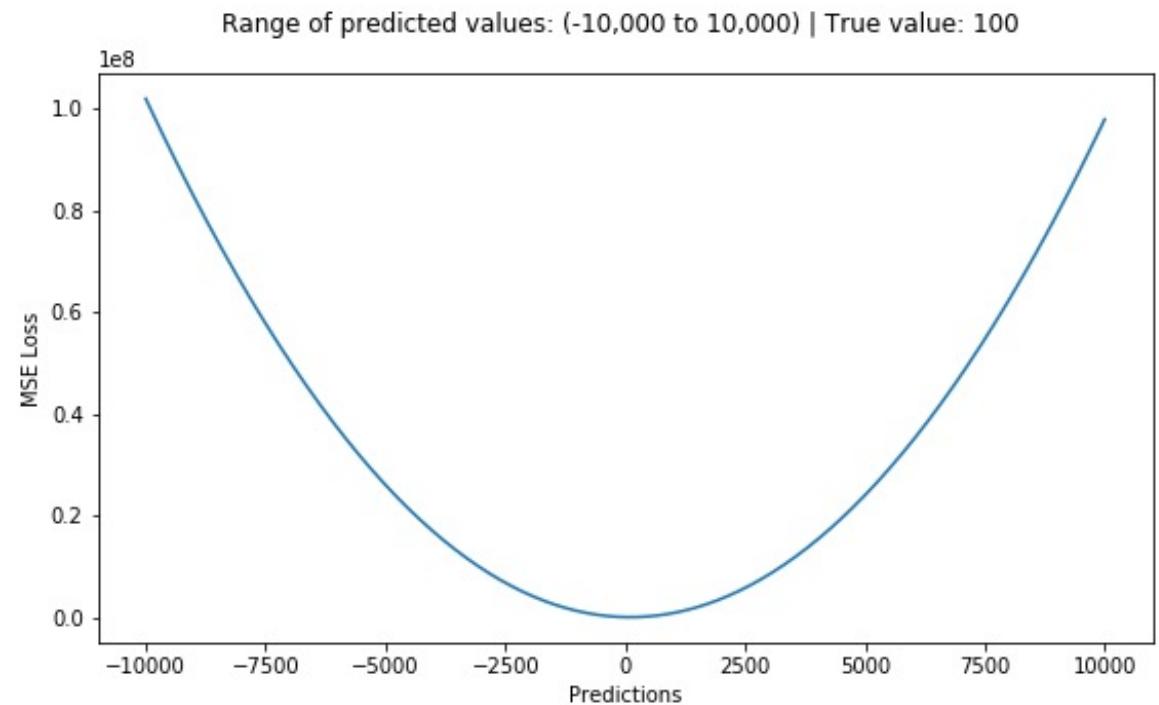
What can you say about the hinge-loss and the log-loss as  $z \rightarrow -\infty$ ?



# Regression: Mean Square Loss

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

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

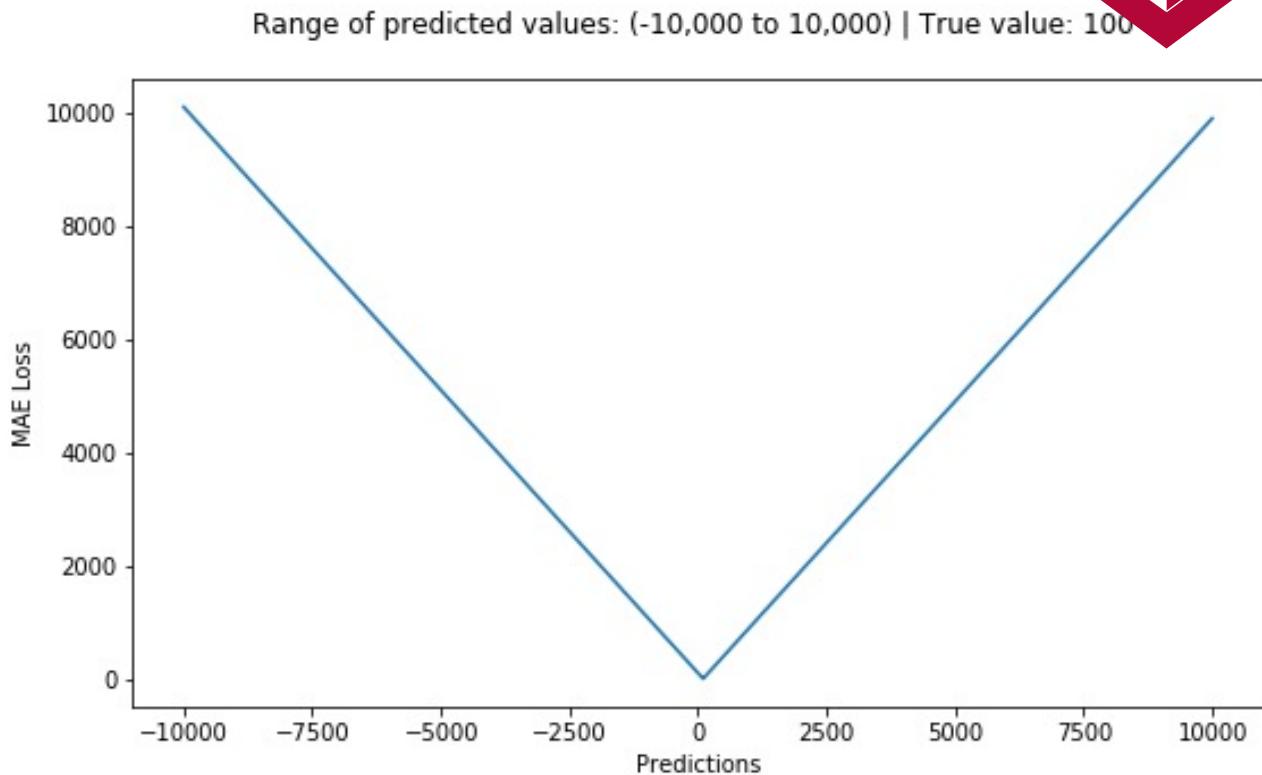




# Regression: Absolute Loss

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

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

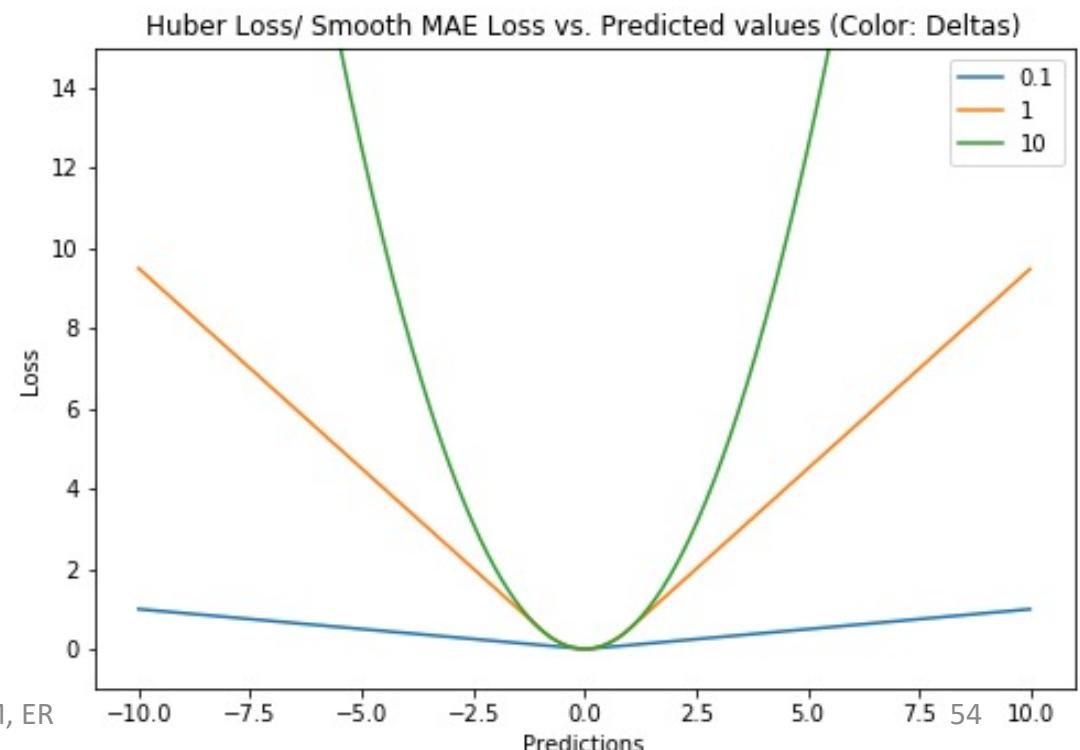




# Regression: Huber Loss

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

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



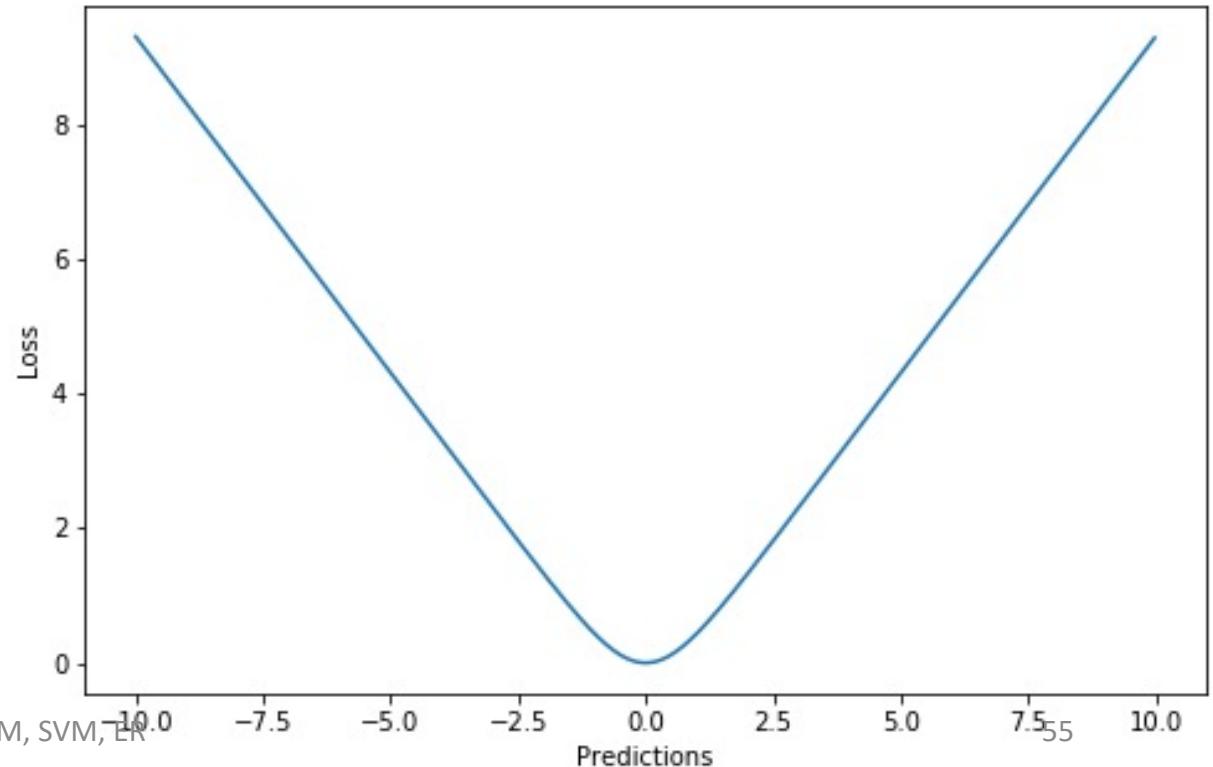


# Regression: Squared Loss

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

- Similar to Huber Loss.
- Differentiable everywhere.

Log-Cosh Loss vs. Predictions





# 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.

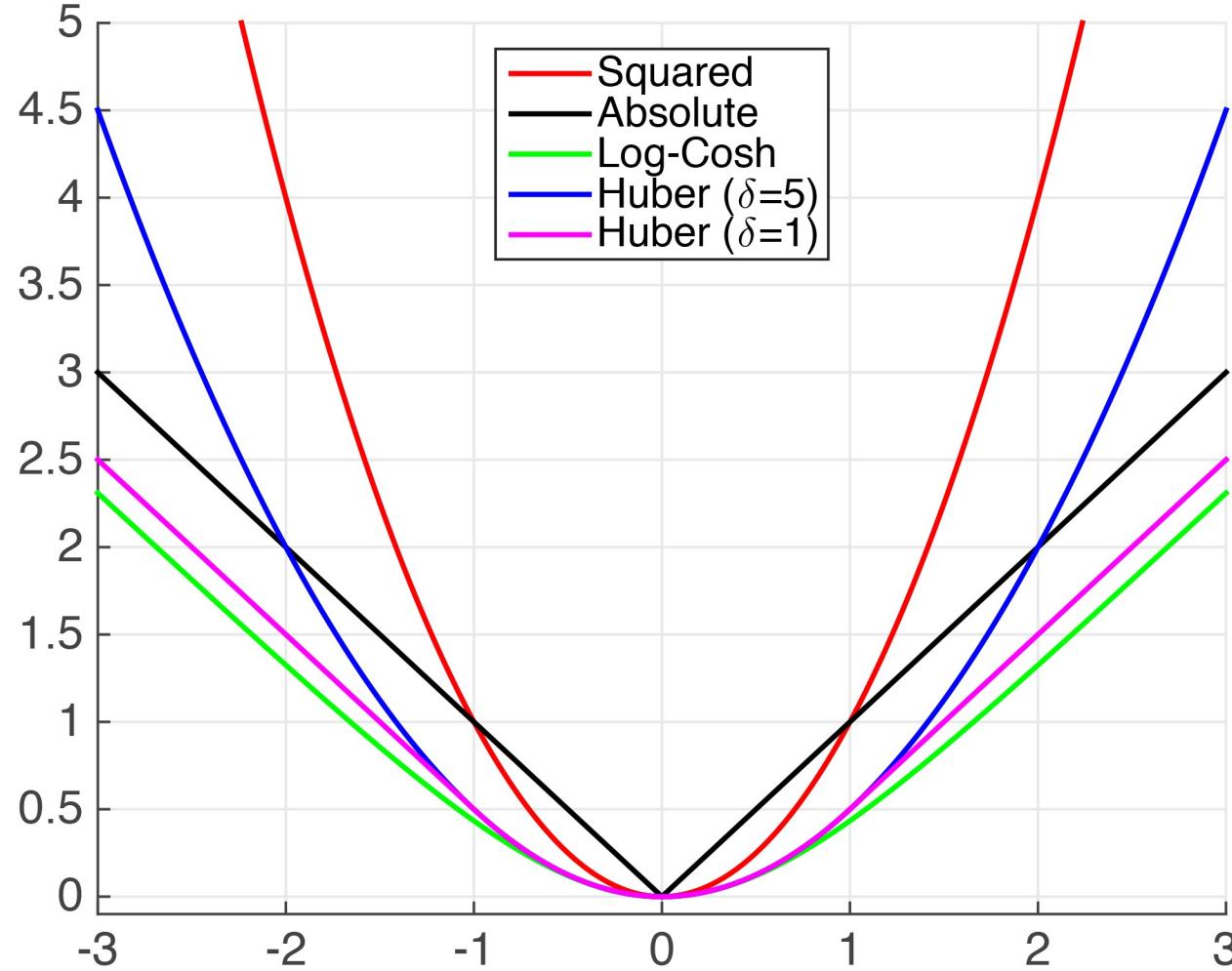
Log-Cosh Loss:

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

- Similar to Huber Loss.
- Differentiable everywhere.



# Regression: Loss Functions Summary



See how sensitive functions are!



# Regularization

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  $B$

# Regularization – Lasso and Ridge

$l_2$  Regularization:

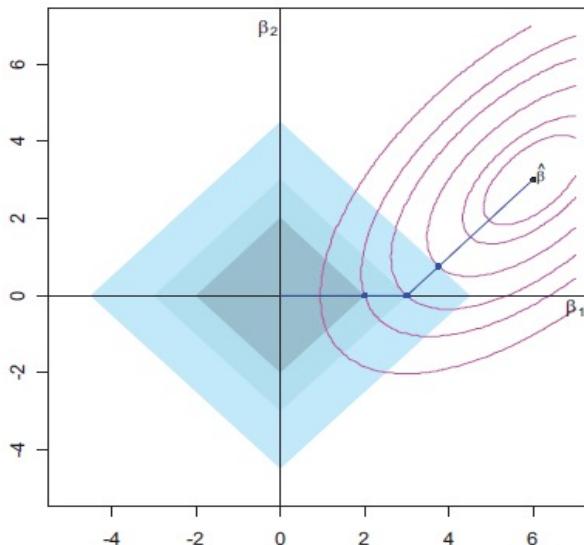
$$r(w) = w^T w = \|w\|^2$$

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

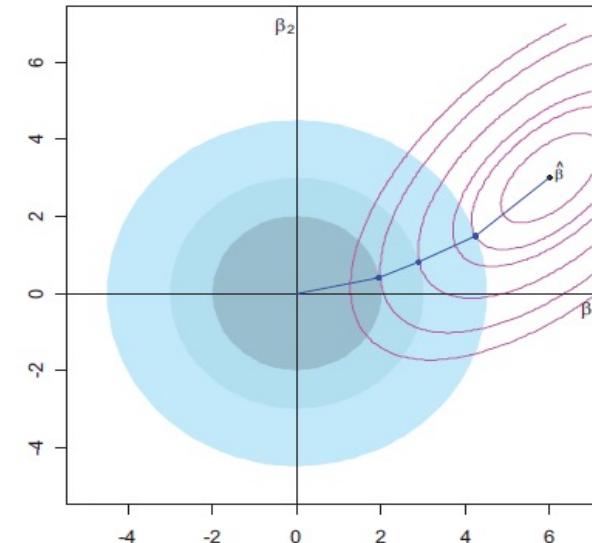
$l_1$  Regularization:

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

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



$l_1$



$l_2$    Lecture 6 - KM, SVM, ER

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

- Solve with gradient descent
- Probability estimation

Linear SVM

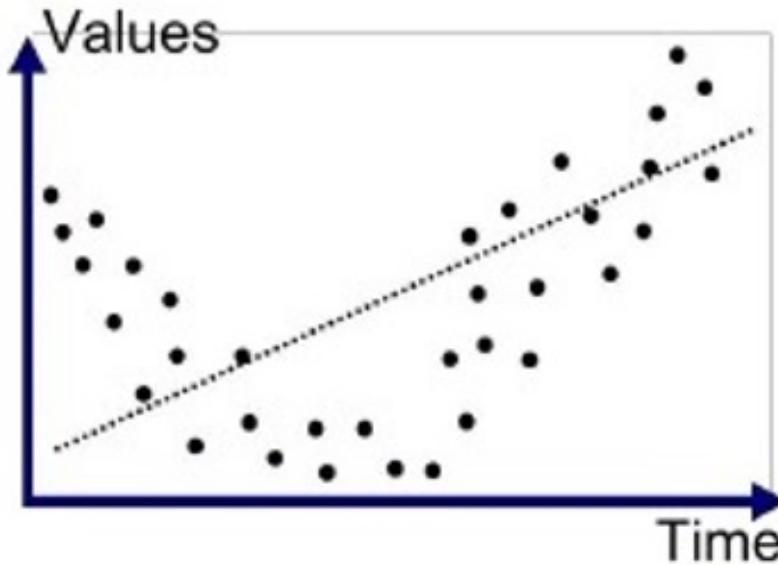
$$\min_w C \sum_{i=1}^N \max[1 - y_i(w^T x_i + b), 0] + \lambda \|w\|^2$$

Lecture 6 - KM, SVM, ER

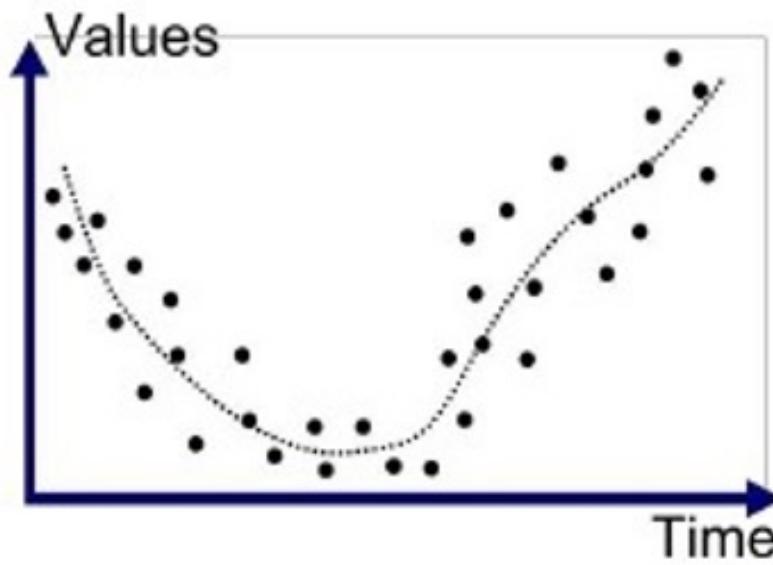
- Typically  $l_2$  regularized sometimes  $l_1$
- Quadratic program
- Kernelize – Solve very efficiently



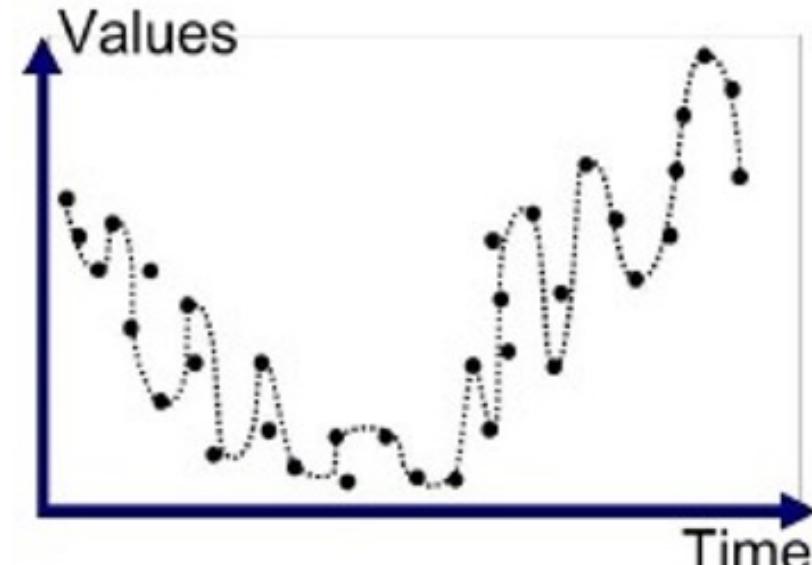
# Overfitting vs. Underfitting



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?

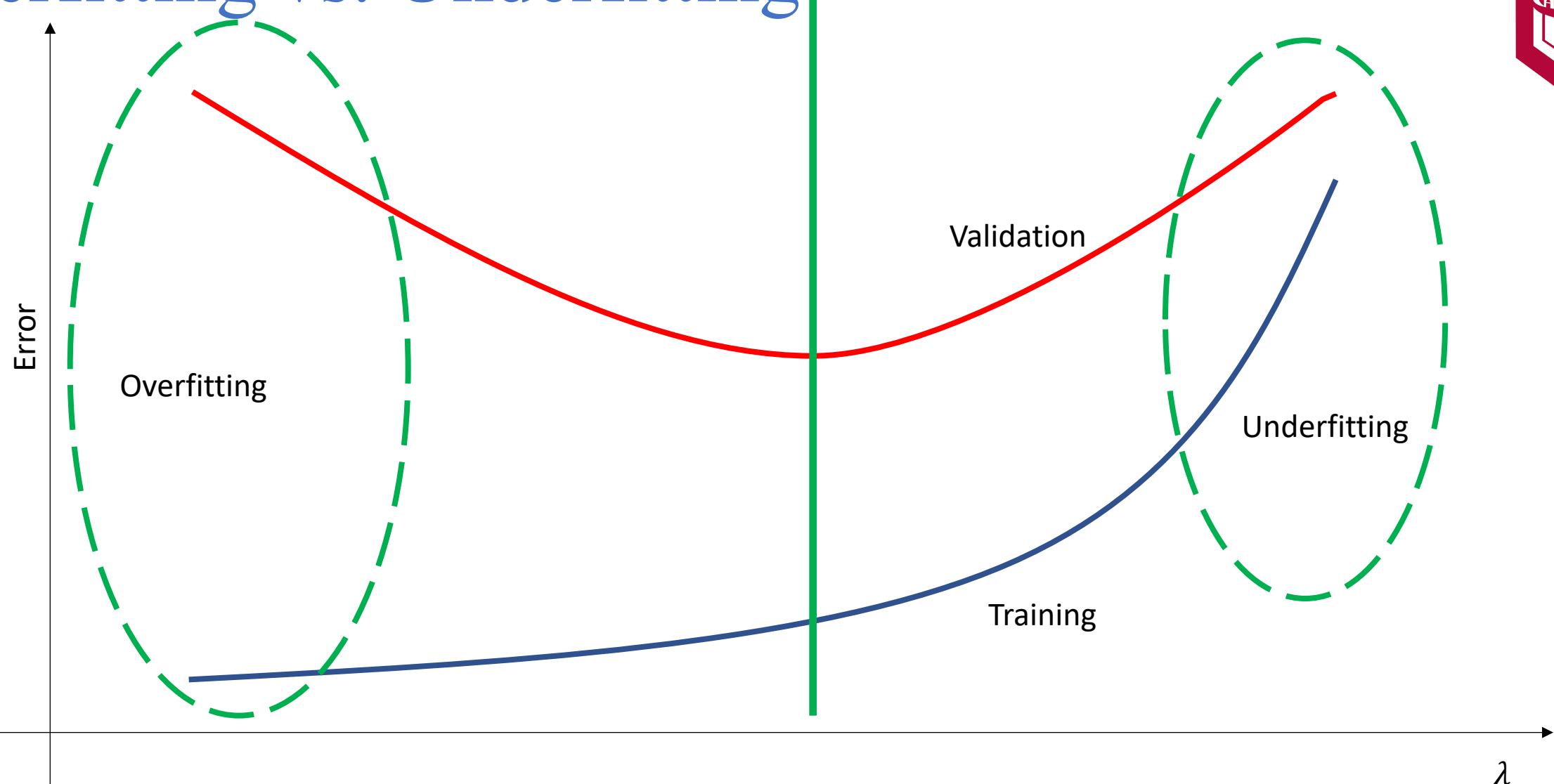
Leads to two possible scenarios: **Underfitting** or **Overfitting**

**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

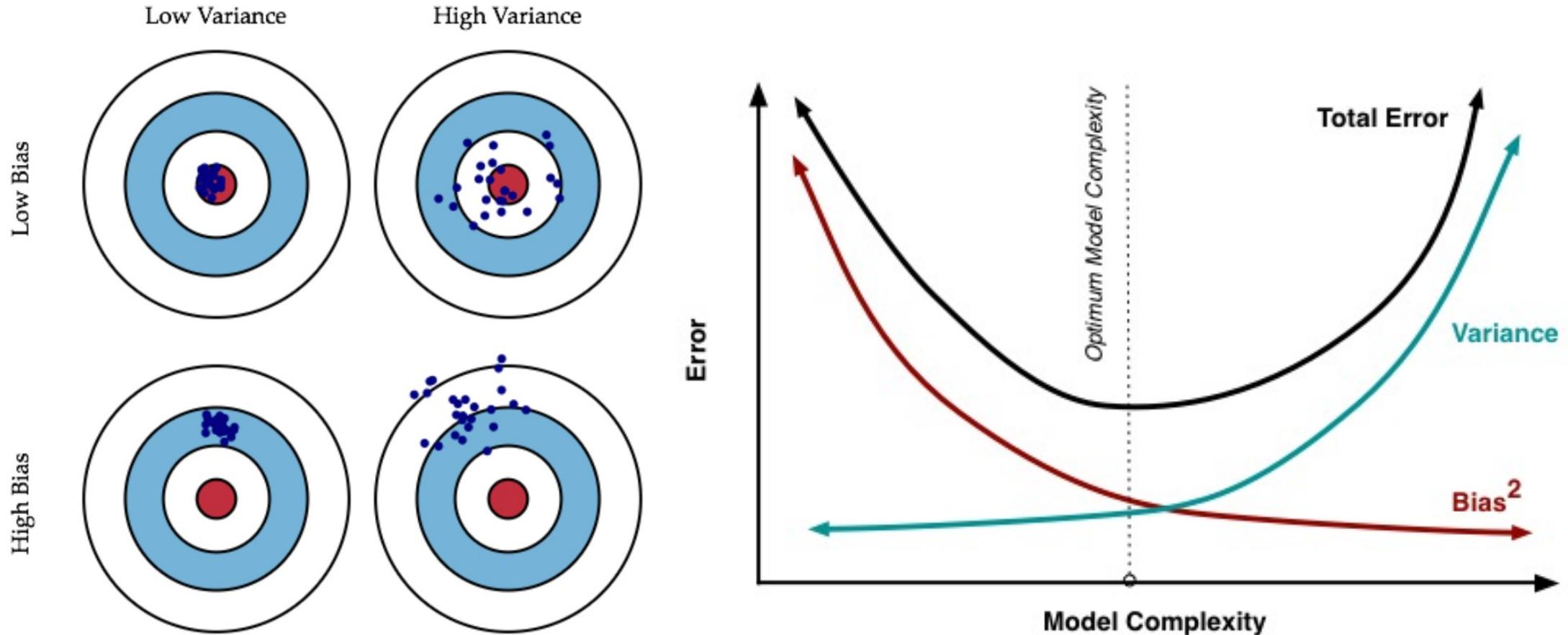


Figure: Graphical illustration of bias vs. variance