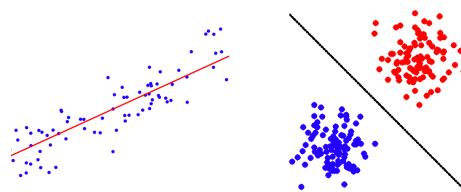


# MACHINE LEARNING

## COMP8220

### 03 – Regression and Classification



## ❖ Linear Regression Introduction

- Linear Regression Model
- Overfitting and Model Selection

## ❖ Classification

- Discriminant Functions
- Perceptron
- Gradient Descent
- K-Nearest Neighbors Classifier

# Lecture Outline

---

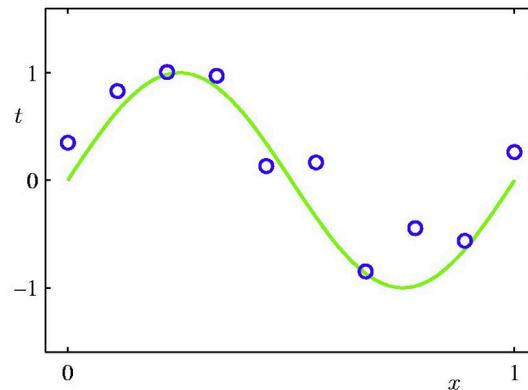
## ❖ Linear Regression Introduction

- Linear Regression Model
- Overfitting and Model Selection

# Regression

- ❖ Suppose we have 10 observations

- $\mathcal{D} = \{\langle x_1, t_1 \rangle, \dots, \langle x_{10}, t_{10} \rangle\}$
- Ground truth:  $t = \sin(2\pi x)$
- Noise added to  $t_i$



- ❖ **Regression** task: estimate  $t = y(x)$  from  $D$



- ❖ A function is linear if

$$f(\alpha x_1 + \beta x_2) = \alpha f(x_1) + \beta f(x_2)$$

- $\alpha, \beta$  are scalars
- Additivity and homogeneity

- ❖ Linear regression model

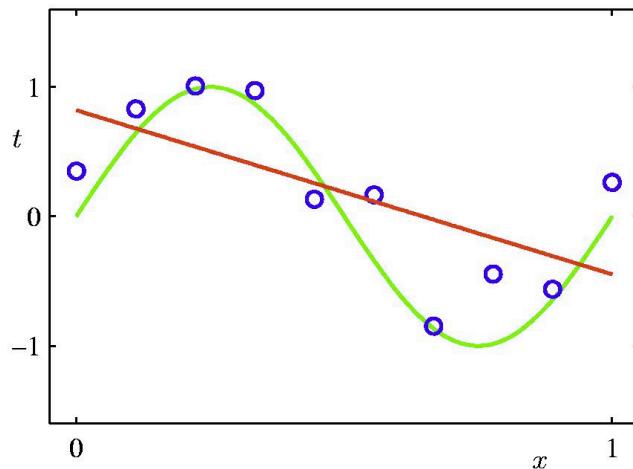
$$y(\mathbf{x}, \omega) = \omega_0 + \omega_1 x_1 + \omega_2 x_2 + \cdots + \omega_D x_D$$

- Note:  $D$  is the # of features, not the # of instances
- The linearity of  $y(\cdot, \omega)$  is with respect to  $\omega$
- Note: we are learning  $\omega$

# The Simplest From

- ❖ Dimensionality  $M = 1$

$$y(x, \omega) = \omega_0 + \omega_1 x$$



- Can only model a “line”
- ❖ How to handle non-linearity?



- ❖ Transform  $x$

$$x \Rightarrow \Phi(x) = (\phi_0(x), \phi_1(x), \phi_2(x), \dots, \phi_M(x))^T$$

- $\phi_j(x), j \geq 0$  form a family of basic functions
- Transformation/feature function is often non-linear

- ❖ Linear basis function models

$$y(x, \omega) = \sum_{j=0}^M \omega_j \phi_j(x) = \omega^T \Phi(x)$$

- $\omega = (\omega_0, \dots, \omega_M)^T, \Phi = (\phi_0, \dots, \phi_M)^T$

# Special Cases

---

- ❖ For  $0 < j \leq D$ ,  $\phi_j(x) = x_j$ ;  $\phi_0(x) = 1$

$$y(\mathbf{x}, \omega) = \sum_{j=0}^D \omega_j x_j$$

- The original form

- ❖  $\phi_j(x) = x^j$ ,  $0 < j \leq M$

$$y(\mathbf{x}, \omega) = \sum_{j=0}^M \omega_j x^j$$

- Power series expansions

# Loss Function for Regression

- ❖ How to learn/choose the parameter vector  $\omega$ ?
  - First, define a loss function  $L(y(x), t)$
  - Then, find out  $y(x)$  to minimize the loss (expectation)

$$y^*(\mathbf{x}) = \arg \min_{y(\mathbf{x})} \mathbb{E}[L] = \int \int L(y(\mathbf{x}), t) p(\mathbf{x}, t) d\mathbf{x} dt$$

- E.g., general squared loss function:

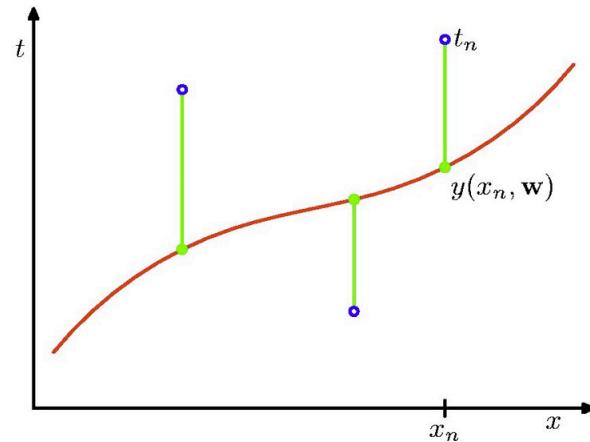
$$L(y(\mathbf{x}), t) = \frac{1}{2} ||y(\mathbf{x}) - t||^2$$

- In practice, we often use the loss function defined on training data

# Least Squares

- ❖ Sum of squared error:

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



- ❖ Then, the specific learning problem

$$\omega^* = \arg \min_{\omega} E(\omega)$$

# Optimization Problem

- ❖ Learning problem → optimization problem

$$\omega^* = \arg \min_{\omega} E(\omega),$$

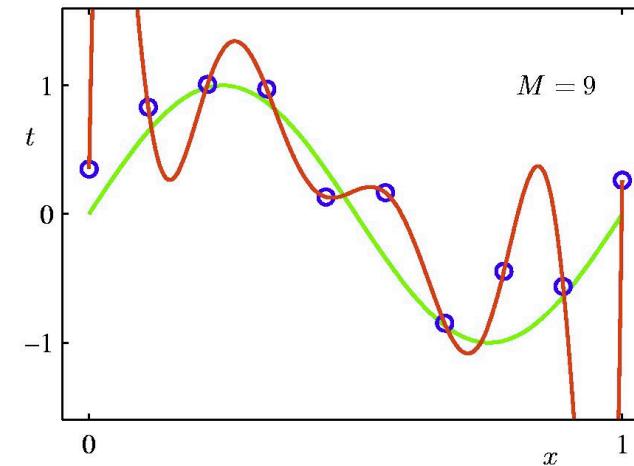
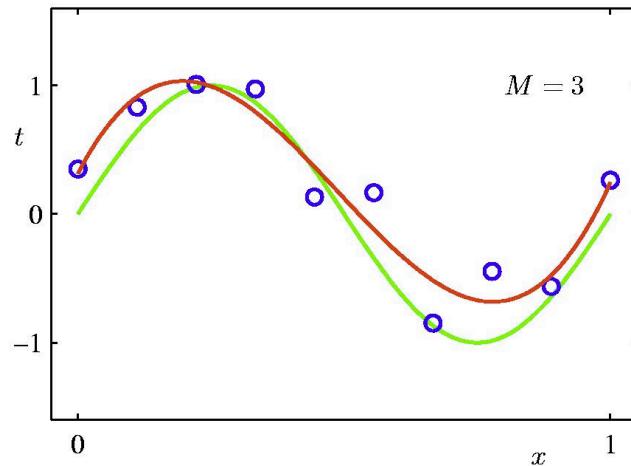
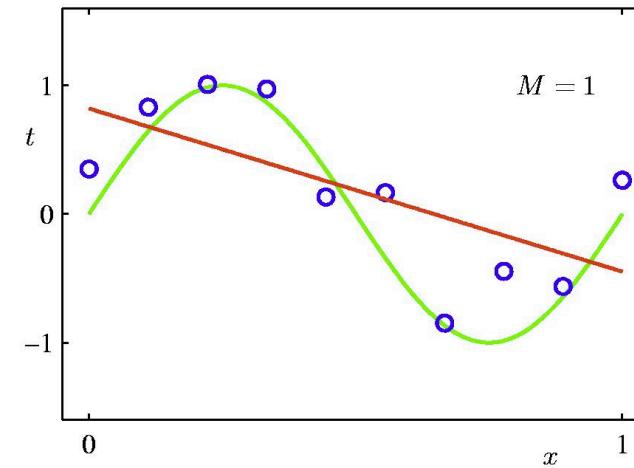
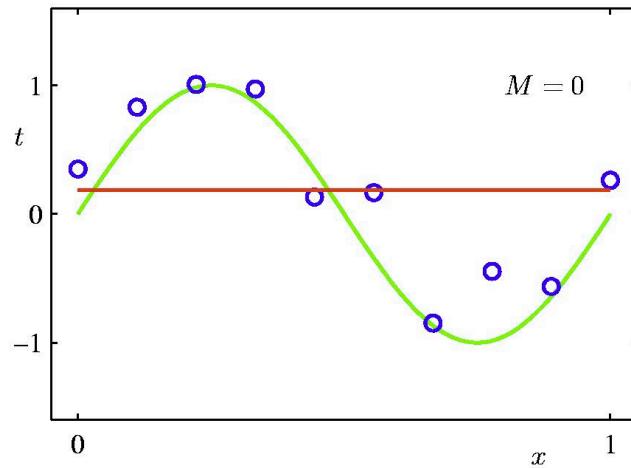
$$E(\omega) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \omega) - t_n\}^2$$

- ❖ How to solve?  $\frac{\partial E(\omega)}{\partial \omega} = 0$

- Close-form solution can be obtained
- For the basis function model

$$\omega_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

# Different Degrees of Polynomials





# Observations

---

- ❖  $M = 0$  or  $M = 1$ 
  - Only passing a couple of points
  - Models are too simple
- ❖  $M = 9$ 
  - Passing every points but oscillating wildly
  - **Overfitting** problem: Lack of generalisation capability
- ❖  $M = 3$ 
  - Much better solution
  - Good trade-off has been made

# Training/Test Error

---

- ❖ Training/empirical error of the trained model  $\hat{f}$  on a training dataset with size  $N$ , e.g.,

$$E_{emp}(\hat{f}) = \frac{1}{N} \sum_{i=1}^N L(y_i, \hat{f}(x_i))$$

- Learning process usually tries to minimize this error
- ❖ Test error on a test data set with size  $N'$ , e.g.,

$$E_{test}(\hat{f}) = \frac{1}{N'} \sum_{i=1}^{N'} L(y_i, \hat{f}(x_i))$$

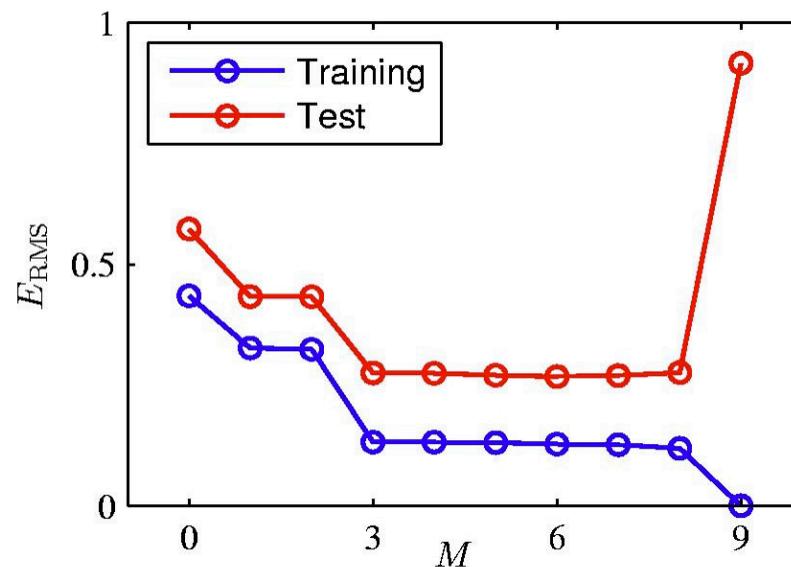
- Test error indicates the **generalisation capability** of a learned model

# Overfitting Issue

- ❖ Test errors with respect to  $M$

- Test data containing 100 points
  - Root-Mean-Square (RMS) error

$$E_{RMS} = \sqrt{2E(\omega^*)/N}$$



- ❖ Observations on  $M$ 
  - Different models generated by varying  $M$
  - $M$  controls the learning capability of the model
  - Unlike model parameter  $\omega$ ,  $M$  can NOT be learned from data (**hyperparameter**)
- ❖ Model selection
  - Models generated by different learning algorithms
  - Models generated by different algorithm parameter (model hyperparameter) configurations
- ❖ Model selection is a non-trial job!

# Regularization

- ❖ Another way to the overfitting issue
  - Penalizing large coefficients in error function
  - E.g., quadratic regularizer

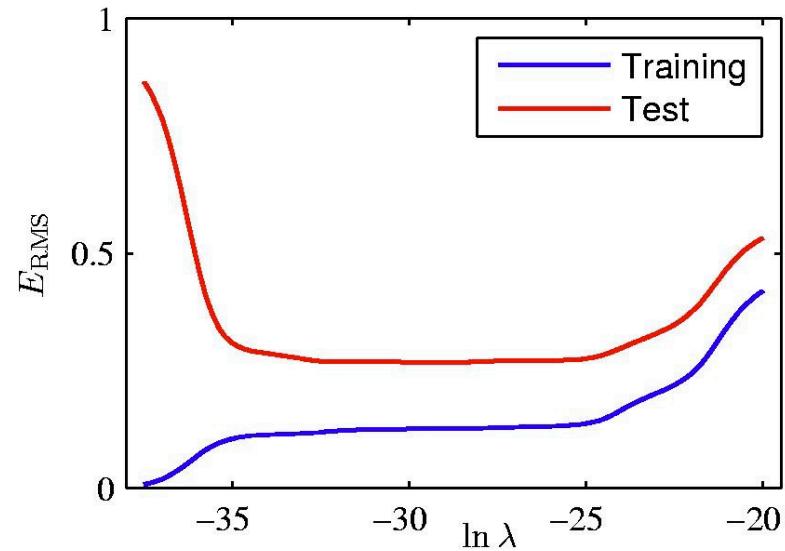
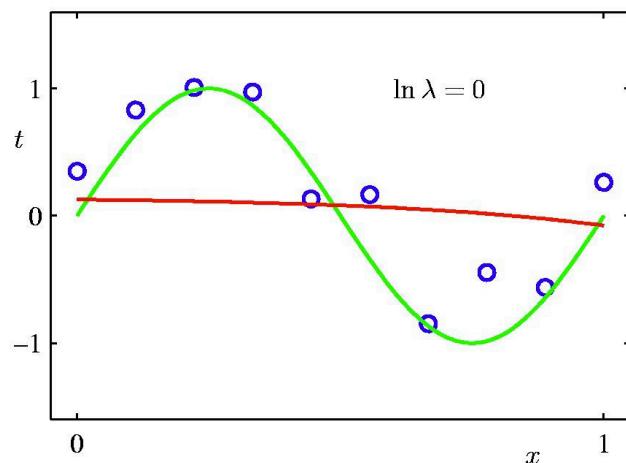
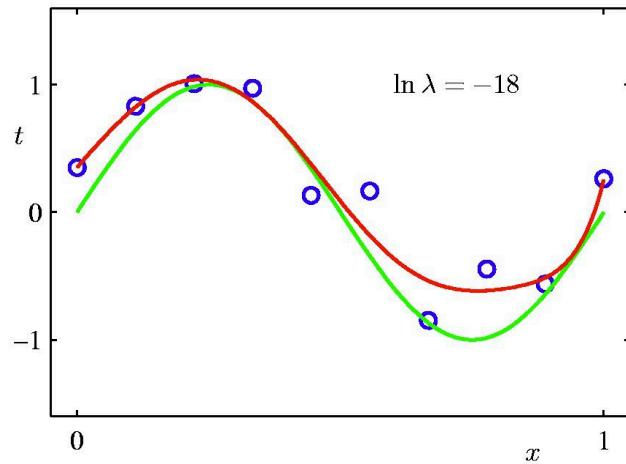
$$\tilde{E}(\omega) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \omega) - t_n\}^2 + \frac{\lambda}{2} \|\omega\|^2,$$

$$\|\omega\|^2 = \omega^T \omega = \omega_0^2 + \omega_1^2 + \dots + \omega_M^2$$

- $\lambda$  governs the trade-offs between the two terms
- A.k.a. ridge regression
- Close-form solution can be obtained
- Lasso regression: the regularizer is  $\|\omega\|_{L_1}$

# Effects of $\lambda$

- ❖ Performance of different  $\lambda$  when  $M = 9$



# Scikit-learn (sklearn)



MACQUARIE  
University



- ❖ Data preparation
- ❖ Training/test data split
  - `x_train, y_train, x_test, y_test`
- ❖ Train the model
  - Build a model: `model = xxxModel.xxxModel()`
  - Train a model: `model.fit(x_train, y_train)`
  - Validate: `cross_validate(xxxmodel, x_train, y_train, k, scoring)`
- ❖ Test the learning model
  - Test the model: `y_pred = xxxModel.predict(x_test)`
  - Can use `y_pred` and `y_test` to calculate the error, e.g.,  
`mean_squared_error(y_test, y_pred)`

## ❖ Linear Regression Introduction

- Linear Regression Model
- Overfitting and Model Selection

## ❖ Classification

- Discriminant Functions
- Perceptron
- Gradient Descent
- K-Nearest Neighbors Classifier

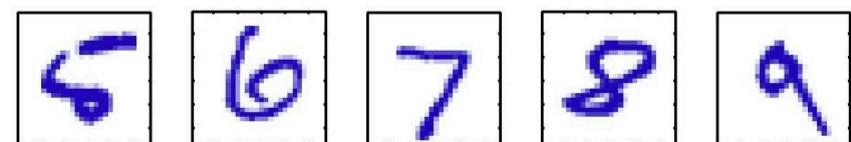
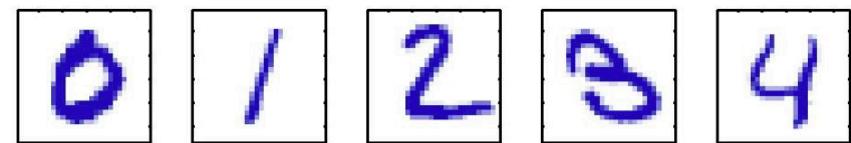
# Classification Problem



- ❖ Given training set  $\mathcal{D} = \{\langle x_1, t_1 \rangle, \langle x_2, t_2 \rangle, \dots, \langle x_N, t_N \rangle\}$
- ❖  $t_i \in \{C_1, \dots, C_K\}$  is discrete: **Classification**
- ❖  $x_i$  can be continuous or discrete
- ❖ Learning problem: to construct  $y: \mathcal{X} \rightarrow \{C_1, \dots, C_K\}$
- ❖ E.g., MNIST data classification

- $t_i$  is the digital 0 to 9
- $x_i$  is an image

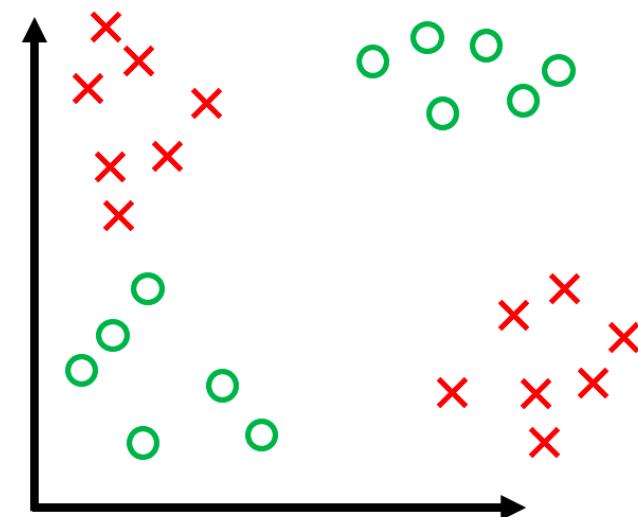
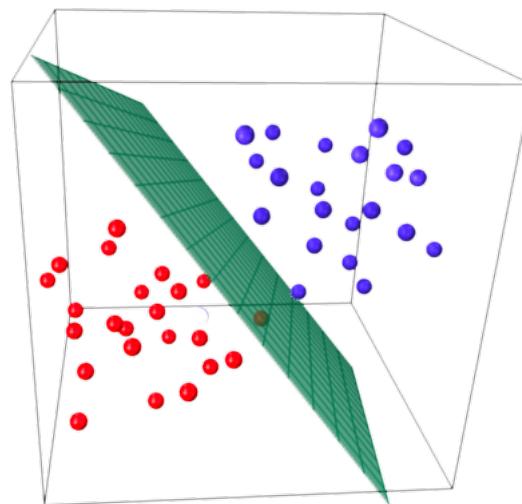
$x_i$  is usually mapped to a vector via a **feature extraction function**  $\phi(x_i)$  explicitly or implicitly



# Linear Separability



- ❖ **Linearly separable:** Classes of labelled data sets can be separated exactly by linear decision surfaces (boundaries)



- ❖ A classic **linearly inseparable** example: XOR data

# 2-class Discriminant Functions

- ❖ The linear decision surface should be a hyperplane

$$\mathbf{w}^T \mathbf{x} + w_0 = 0$$

- M-dimensional feature vector  $\mathbf{x}$
- Classes:  $C_1$  and  $C_2$ , or  $t_i \in \{0, 1\}$
- $\mathbf{w}$ : weight vector (column-based)
- $w_0$  bias or intercept

- ❖ The discriminant function (the predictive model)

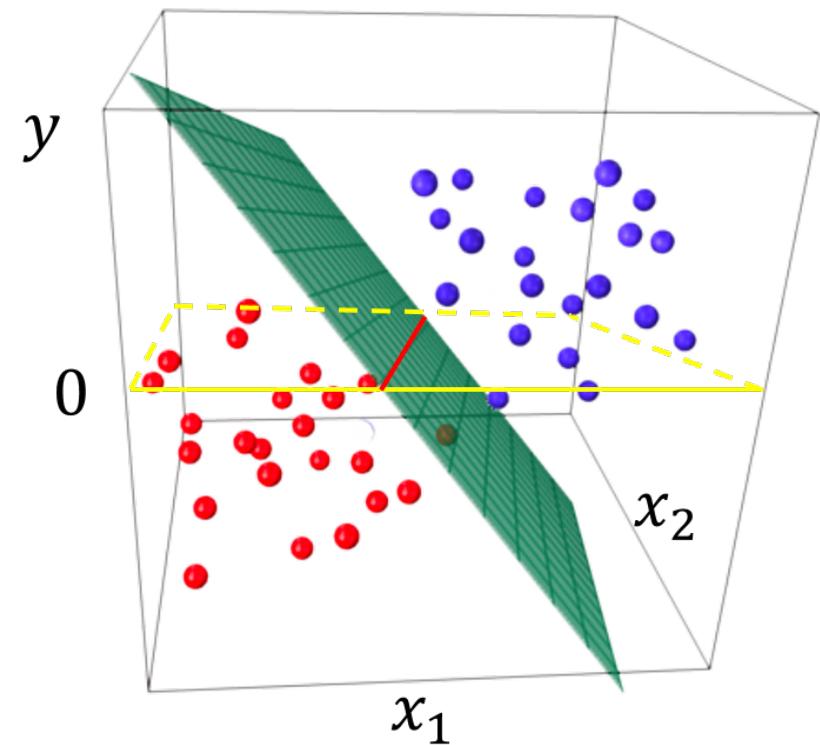
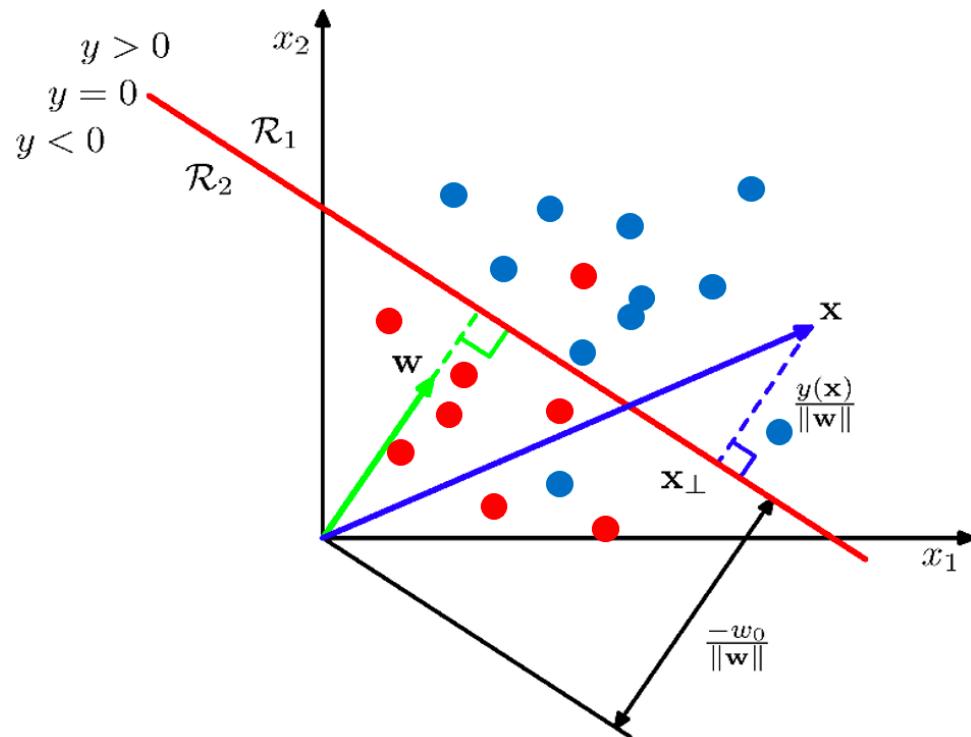
$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

- ❖ Decision rule

$$Label(\mathbf{x}) \leftarrow \begin{cases} C_1, & y(\mathbf{x}) \geq 0 \\ C_2, & y(\mathbf{x}) < 0 \end{cases}$$

# Geometry Perspective

- ❖ E.g., 2-dimensional input space





- ❖  $w$  determines the orientation of the decision surface
  - Assume two points  $x_A$  and  $x_B$  on surface
  - $\because y(x_A) = y(x_B) = 0, w^T(x_A - x_B) = 0$
- ❖ Normal distance from the origin to the decision surface
$$\frac{w^T x_A}{\|w\|} = - \frac{w_0}{\|w\|}$$
  - $w_0$  determines the location of the decision surface
- ❖ Relationship between  $y(x)$  and the signed distance of  $x$  to the decision surface (denoted as  $r$  )

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

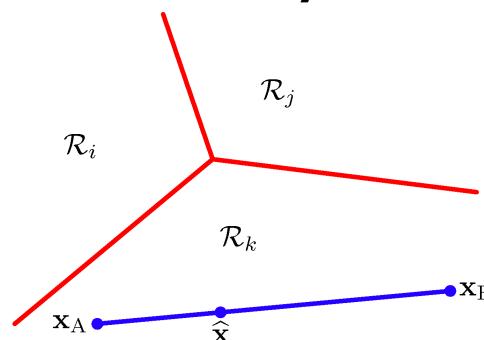
- $\because x = x_\perp + r \frac{w}{\|w\|}$ ; multiply the equation by  $w$  and add  $w_0$

# Multiple-class Cases

- ❖  $t_i \in \{C_1, \dots, C_K\}$
- ❖  $K$ -class discriminant functions

$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}$$

- $1 \leq k \leq K$
- ❖ Decision rule
$$\text{Label}(\mathbf{x}) \leftarrow \arg \max_k \{y_k(\mathbf{x})\}$$
- ❖ Decision regions are always connected and convex



# Least Squares for Classification

- ❖ 1-of- $K$  binary coding scheme (a.k.a. one-hot coding)

$$c_k \Leftrightarrow \left( 0, \dots, 0, \underset{k^{th} position}{1}, 0, \dots, 0 \right)^T$$

- ❖ Matrix denotation for the model

$$y(x) = \tilde{W}^T \tilde{x}$$

- $\tilde{x} \equiv \langle 1, x_1, \dots, x_M \rangle$
- $\tilde{W}$  includes both coefficients and intercepts
- ❖ Given dataset  $\mathcal{D}$ , we can minimize sum-of-squares error to learn the model parameters  $\tilde{W}$

# Least Squares for Classification



$$E_D(\tilde{\mathbf{w}}) = \frac{1}{2} \text{Tr}\{\mathbf{E}^T \mathbf{E}\}$$

$$\mathbf{E} = \tilde{\mathbf{X}} \tilde{\mathbf{w}} - \mathbf{T} = \begin{bmatrix} \mathbf{e}_1^T \\ \vdots \\ \mathbf{e}_N^T \end{bmatrix}, \quad \mathbf{e}_i = \tilde{\mathbf{w}}^T \mathbf{x}_i - t_i$$

$$\mathbf{E}^T \mathbf{E} = \begin{bmatrix} \mathbf{e}_1^T \\ \vdots \\ \mathbf{e}_N^T \end{bmatrix}^T \begin{bmatrix} \mathbf{e}_1^T \\ \vdots \\ \mathbf{e}_N^T \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^N e_{i1}^2 & \sum_{i=1}^N e_{i1} e_{i2} & \cdots \\ \sum_{i=1}^N e_{i2} e_{i1} & \ddots & \vdots \\ \vdots & \cdots & \sum_{i=1}^N e_{ik}^2 \end{bmatrix}$$

# Least Squares for Classification

- ❖ Set the derivative w.r.t.  $\widetilde{W}$  to  $0$  to obtain

$$\widetilde{W} = (\widetilde{X}^T \widetilde{X})^{-1} \widetilde{X}^T T = \widetilde{X}^\dagger T$$

- $\widetilde{X}^\dagger$  is the pseudo-inverse of the matrix of  $\widetilde{X}$

- ❖ The discriminant function

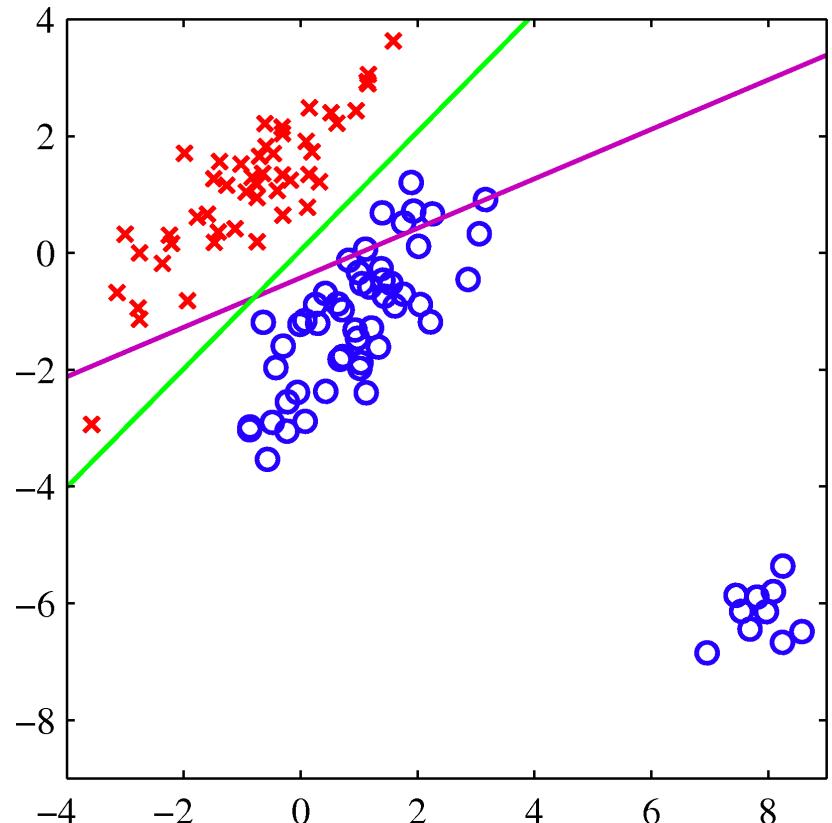
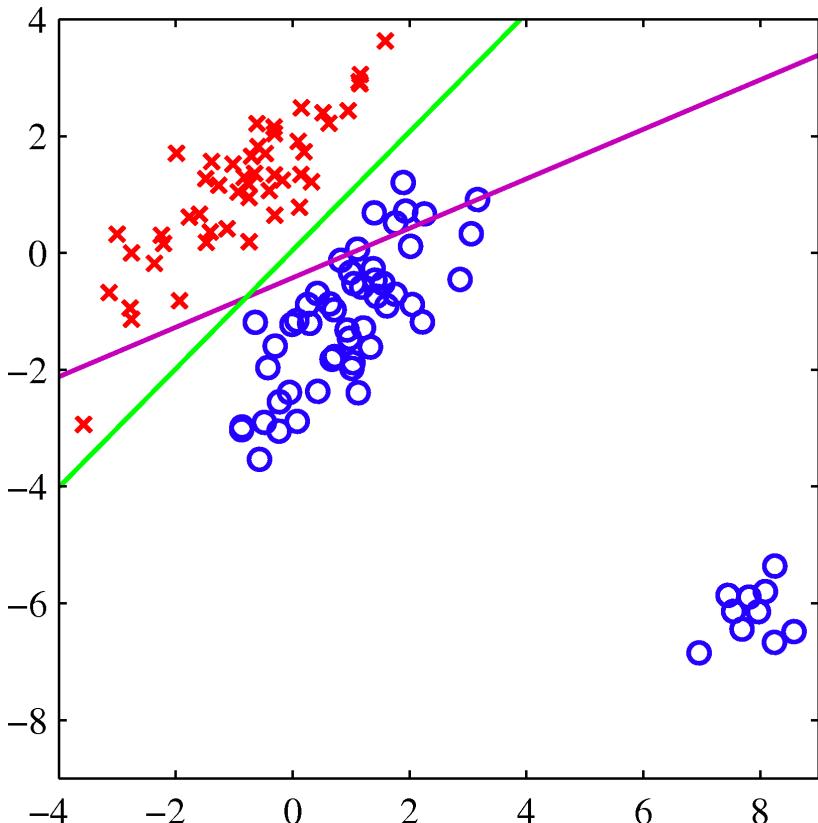
$$y(x) = \widetilde{W}^T \widetilde{x} = T^T (\widetilde{X}^\dagger)^T \widetilde{x}$$

- ❖ Issues: poor robustness w.r.t. outliers

- Least squares method corresponds to maximum likelihood under Gaussian conditional distribution
- Binary target vectors clearly fail to follow Gaussian

# Least Squares for Classification

## ❖ Robustness w.r.t. Outliers



- ❖ The Perceptron two-class model

$$y(\mathbf{x}) = f(\mathbf{w}^T \phi(\mathbf{x}))$$

- ❖  $f(\cdot)$  is a step function

$$f(a) = \begin{cases} +1, & a \geq 0 \\ -1, & a < 0 \end{cases}$$

- ❖ Target value encoding +1 for  $C_1$ , and -1 for  $C_2$

- ❖ Observations

- If  $\mathbf{x}_i$  in  $C_1$ (i.e.,  $t_n = +1$ ),  $\mathbf{w}^T \phi(\mathbf{x}_i) > 0$
- If  $\mathbf{x}_i$  in  $C_2$ (i.e.,  $t_n = -1$ ),  $\mathbf{w}^T \phi(\mathbf{x}_i) < 0$

- ❖  $\Rightarrow \mathbf{w}^T \phi(\mathbf{x}_i) t_n > 0$

- ❖ **Perceptron criterion:** the error function used to determine  $\mathbf{w}$ 
  - Correct classification: 0
  - Misclassification error:  $-\mathbf{w}^T \phi(\mathbf{x}_i) t_i$
- ❖ We hope to drive  $-\mathbf{w}^T \phi(\mathbf{x}_i) t_i$  to 0
- ❖ The Perceptron criterion function

$$E_P(\mathbf{w}) = - \sum_{i \in \mathcal{M}} \mathbf{w}^T \phi(\mathbf{x}_i) t_i$$

- $\mathcal{M}$  denotes the set of misclassified patterns
- ❖ Key: How to determine  $\mathcal{M}$ ?



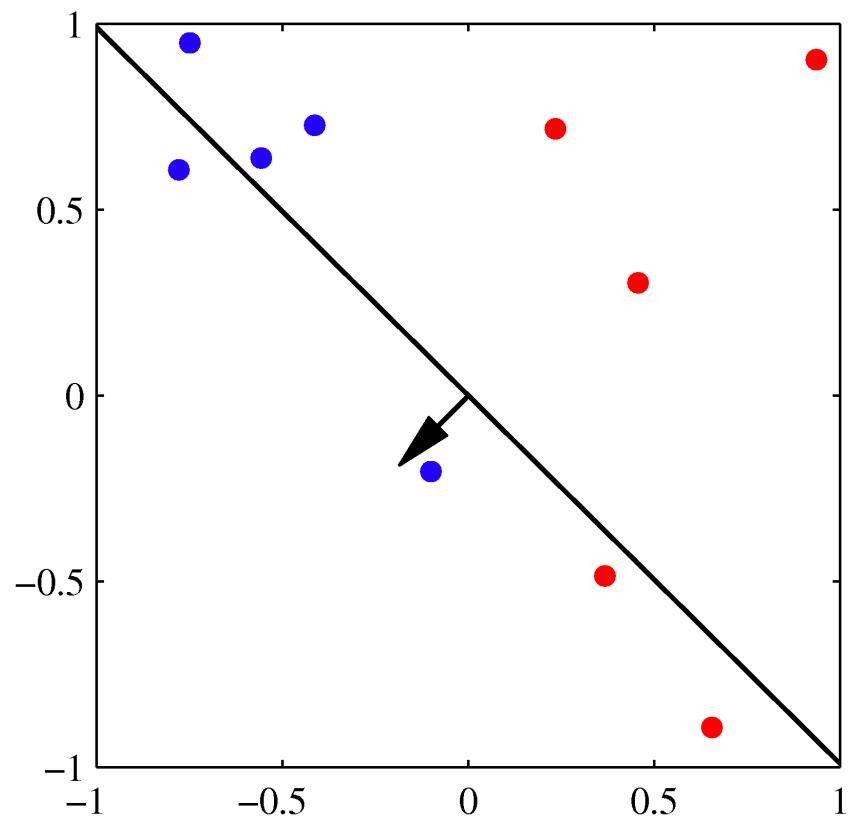
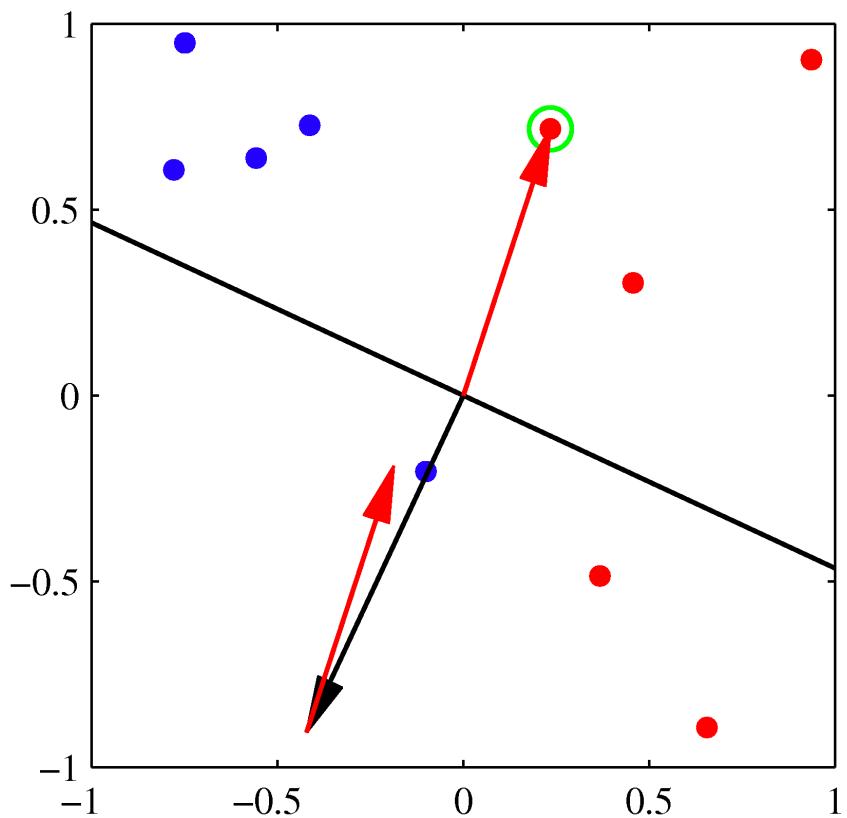
- ❖ Rewrite  $E_P(\mathbf{w})$  as

$$J(\mathbf{w}) = E_P(\mathbf{w}) = \sum_{i=1}^n J_i(\mathbf{w})$$

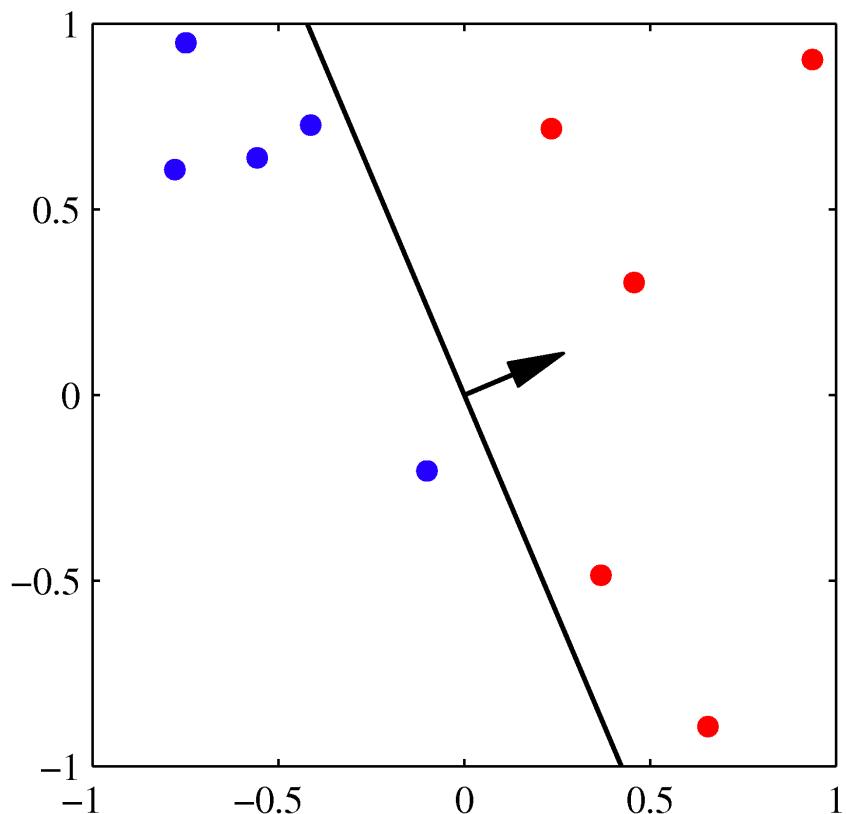
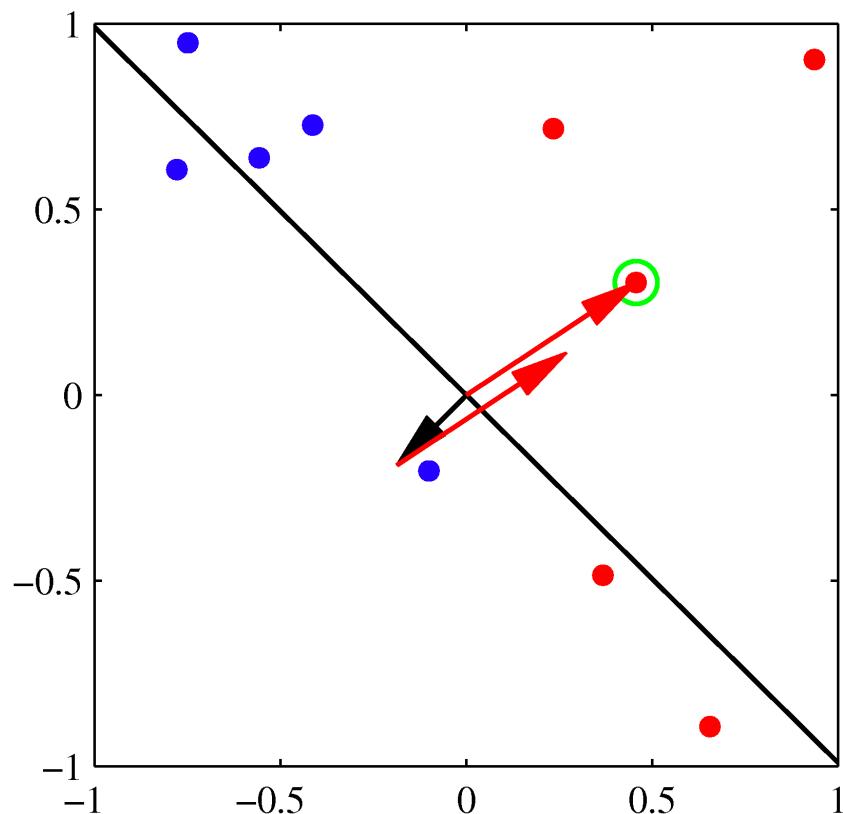
$$\blacksquare J_i(\mathbf{w}) = \begin{cases} 0, & \mathbf{w}^T \phi(\mathbf{x}_i) t_i \geq 0 \\ -\mathbf{w}^T \phi(\mathbf{x}_i) t_i, & \mathbf{w}^T \phi(\mathbf{x}_i) t_i < 0 \end{cases}$$

- ❖ Then, we can use gradient descent method to estimate  $\mathbf{w}$
- ❖ Perceptron convergence theorem guarantees solutions if existing (data set is linearly separable)

# Perceptron



# Perceptron

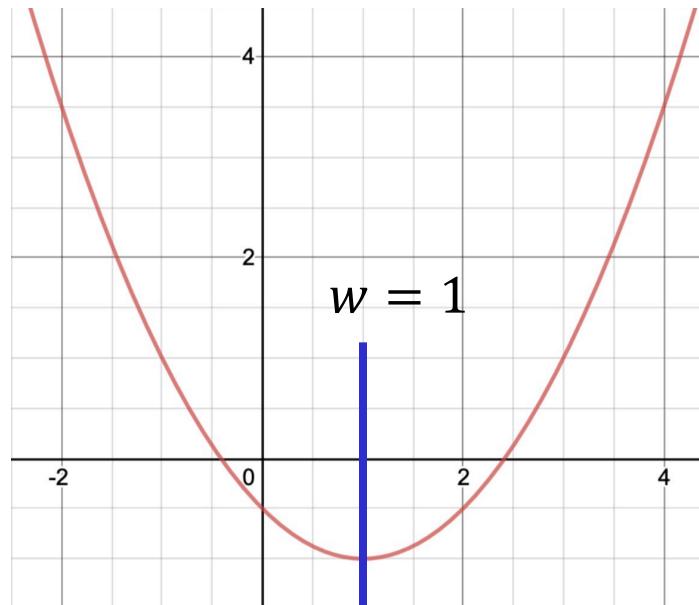


# Gradient Descent

- ❖ Question: which  $w$  minimize the following function?

$$f(w) = \frac{1}{2}(w - 1)^2 - 1$$

- ❖ Option I: observe its curve (middle/high school)



# Gradient Descent (Cont'd)



- ❖ Option 2: analytical method with derivative (calculus)

$$f'(w) = \frac{1}{2} \cdot 2 \cdot (w - 1) \cdot 1 - 0 = w - 1$$

- To find stationary points, let

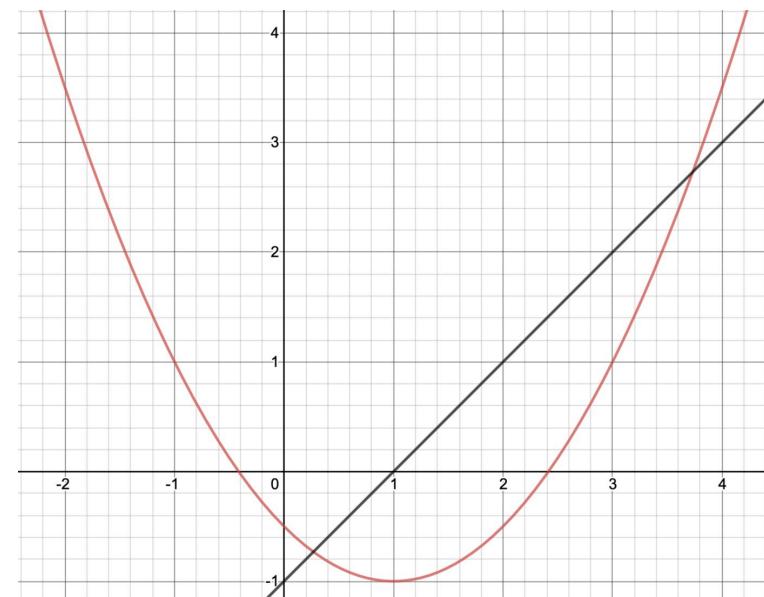
$$f'(w) = 0$$

$$f'(w) = w - 1 = 0$$

$$\Rightarrow w^* = 1$$

$$f''(w) = (w - 1)' = 1 > 0$$

- $w^* = 1$  is the value minimizing the function
- **Linear regression** adopt this method to compute weights



# Gradient Descent (Cont'd)

## ❖ Option 3: Numerical method

- Randomly initialize the value ( $w^{(0)}$ ), and optimize it gradually
- Needs an update rule (still make use of derivative)

$$w^{(\tau+1)} \leftarrow w^{(\tau)} - \eta f'(w) \Big|_{w^{(\tau)}} = w^{(\tau)} - \eta(w^{(\tau)} - 1)$$

- $\eta$  is step size (or learning rate)

- Stops after the result is good enough
  - E.g., difference between two iterations is small enough

## ❖ Unlike Option 2, we don't need to solve $f'(w) = 0$

- When a function is more complicated, solving the equation system analytically is really challenging
- So, this method is more general with an approximate result

# Gradient Descent (Cont'd)



## ❖ Option 3: Numerical method

- Show an example when  $w^{(0)} = 0$ , and  $\eta = \frac{1}{2}$

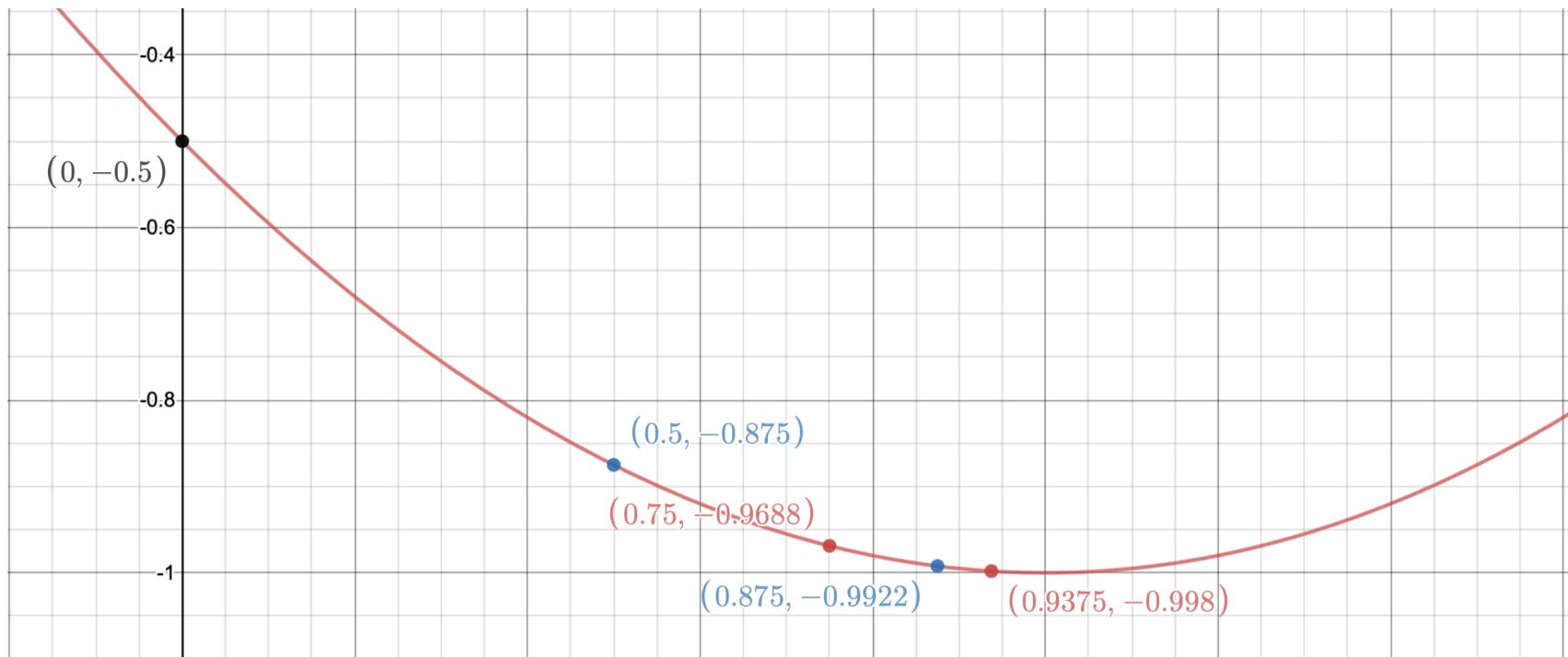
$\tau$	$w$	$f'(w)$	$f(w)$	$ f^{(\tau+1)}(w) - f^{(\tau)}(w) $
0	0	-1	$-\frac{1}{2}$	
1	$0 - \frac{1}{2}(-1) = \frac{1}{2}$	$-\frac{1}{2}$	$-\frac{7}{8}$	0.375
2	$\frac{1}{2} - \frac{1}{2}\left(-\frac{1}{2}\right) = \frac{3}{4}$	$-\frac{1}{4}$	$-\frac{31}{32}$	0.093
3	$\frac{3}{4} - \frac{1}{2}\left(-\frac{1}{4}\right) = \frac{7}{8}$	$-\frac{1}{8}$	$-\frac{127}{128}$	0.0234
4	$\frac{7}{8} - \frac{1}{2}\left(-\frac{1}{8}\right) = \frac{15}{16}$	$-\frac{1}{16}$	$-\frac{511}{512}$	0.0058
...	...	...	...	...

# Gradient Descent (Cont'd)



## ❖ Option 3: Numerical method

- $w^* \leftarrow w^{(4)} = \frac{15}{16}$  is a good approximate to the true value 1

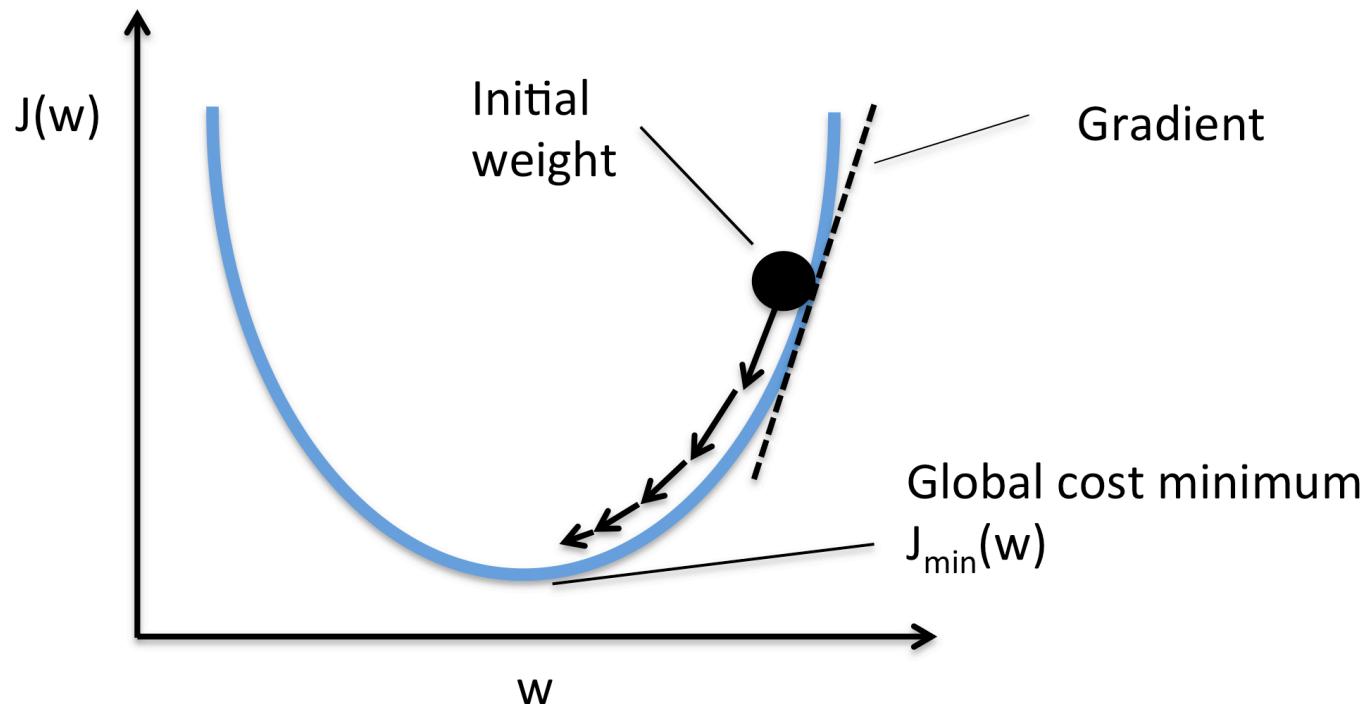


# Gradient Descent (Cont'd)



- ❖ General rule for a cost function  $J(w)$

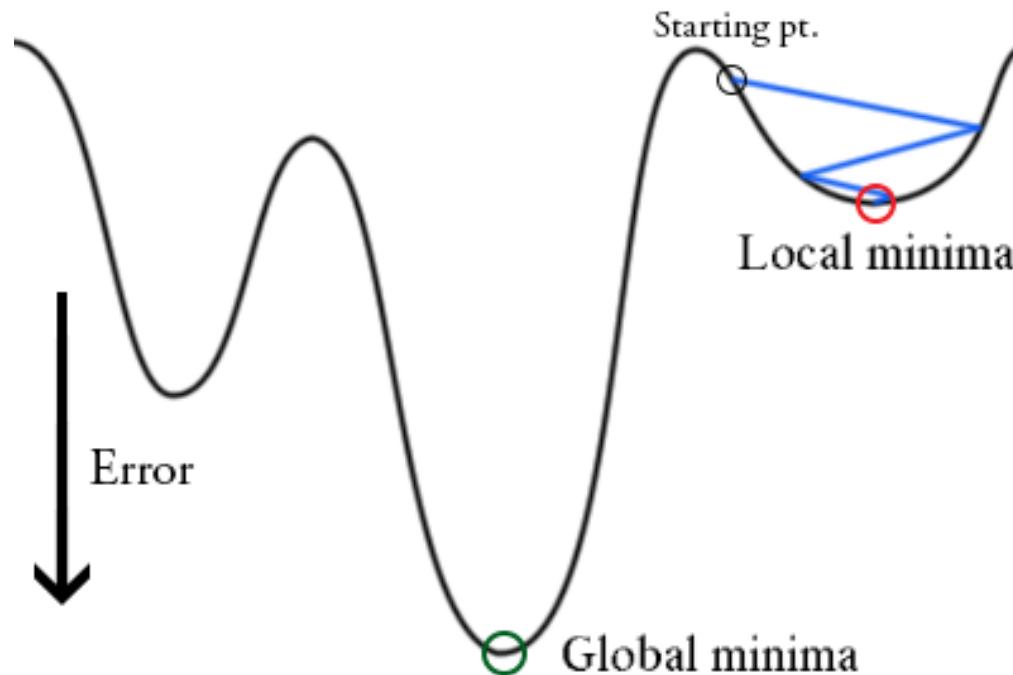
$$w^{(\tau+1)} \leftarrow w^{(\tau)} - \eta \cdot J'(w) \Big|_{w^{(\tau)}}$$



# Gradient Descent (Cont'd)



- ❖ What if the function is **non-convex**?
  - Can get stuck in local minima



# Gradient Descent (Cont'd)

- ❖ Extend to multiple variable function
  - This is a much more common case in machine learning
- ❖ Beyond ‘derivative’, we need ‘gradient’
- ❖ For a function  $f(\mathbf{w}) = f(w_1, \dots, w_i, \dots, w_n)$
- ❖ Gradient of  $f(\mathbf{w})$

$$\nabla f(\mathbf{w}) = \left( \frac{\partial f(\mathbf{w})}{\partial w_1}, \dots, \frac{\partial f(\mathbf{w})}{\partial w_i}, \dots, \frac{\partial f(\mathbf{w})}{\partial w_n} \right)$$

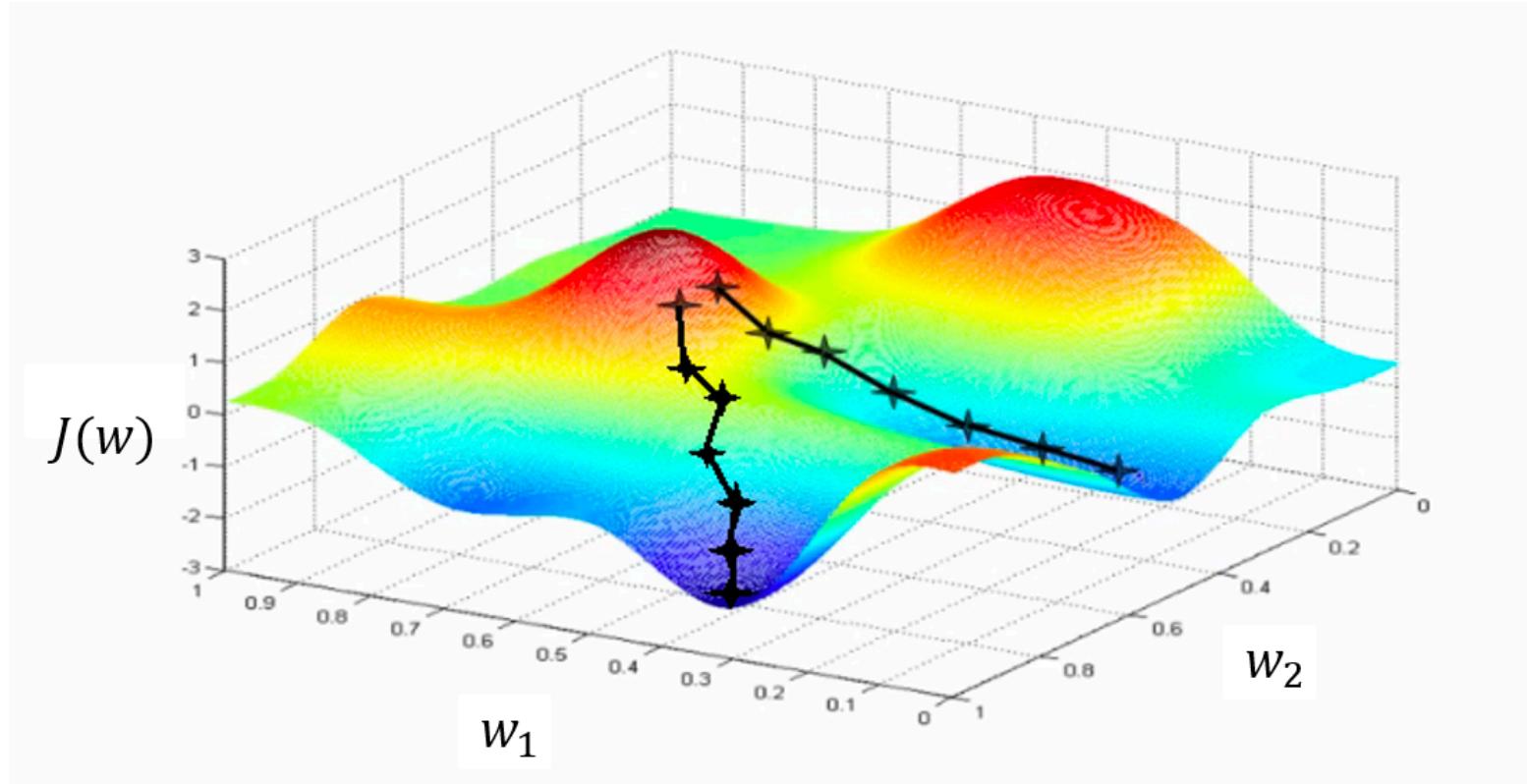
- $\frac{\partial f(\mathbf{w})}{\partial w_i}$  is partial derivative w.r.t.  $w_i$
- **Direction:** direction of steepest increase
- **Magnitude:** rate of change in that direction

# Gradient Descent (Cont'd)



- ❖ Update rule for a cost function  $J(\mathbf{w})$

$$\mathbf{w}^{(\tau+1)} \leftarrow \mathbf{w}^{(\tau)} - \eta \cdot \nabla J(\mathbf{w}) \Big|_{\mathbf{w}^{(\tau)}}$$





## ❖ Batch Gradient Descent

- Uses the whole training set to compute gradients
- Slow when the training set is large

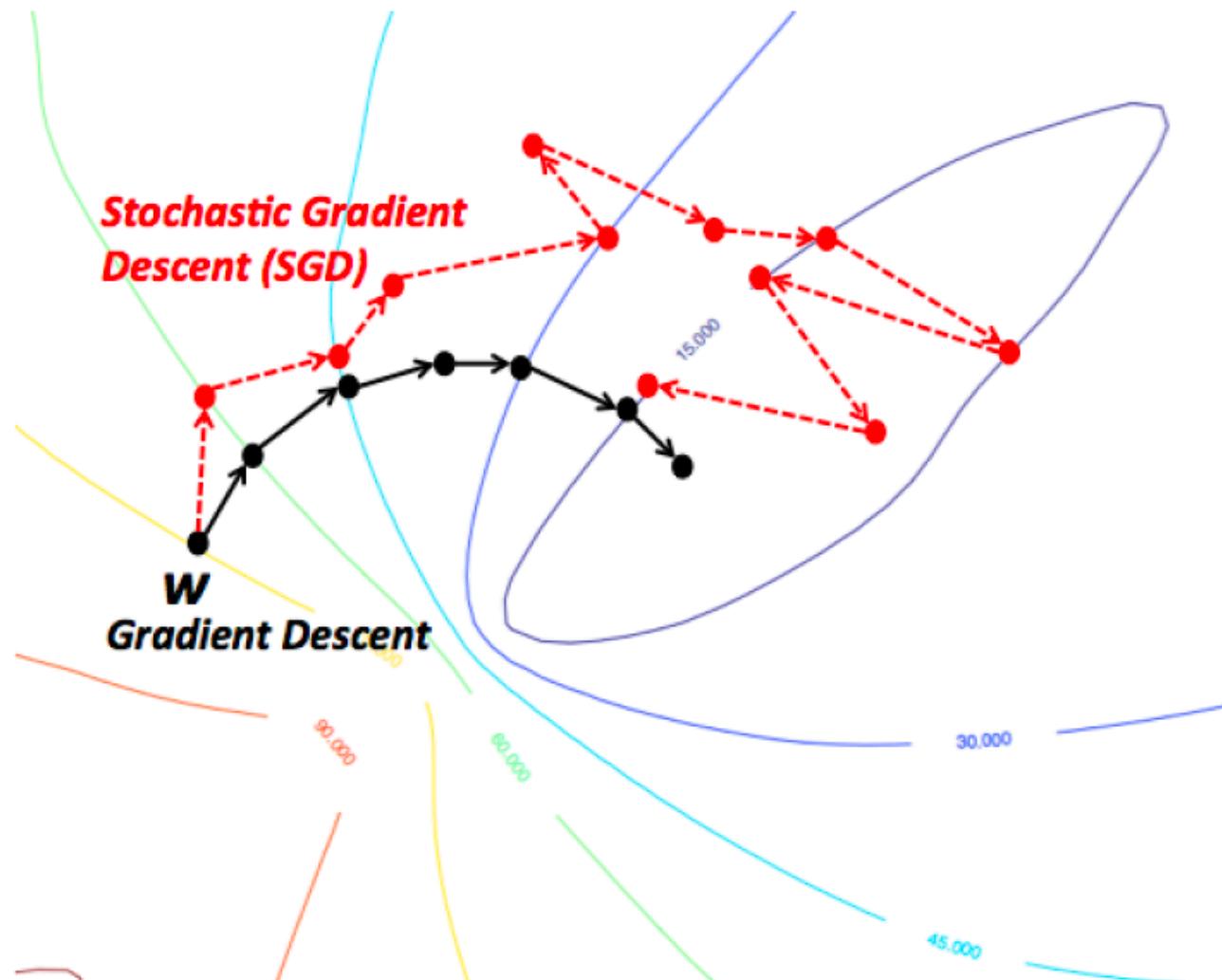
## ❖ Stochastic Gradient Descent

- Picks a random instance at each step to compute gradients
- Fast but not optimal
- Learning rate is determined by a learning schedule

## ❖ Mini-batch Gradient Descent

- Computes the gradient on a small random set of instances
- Gets performance boost over SGD when using GPUs

# Stochastic Gradient Descent

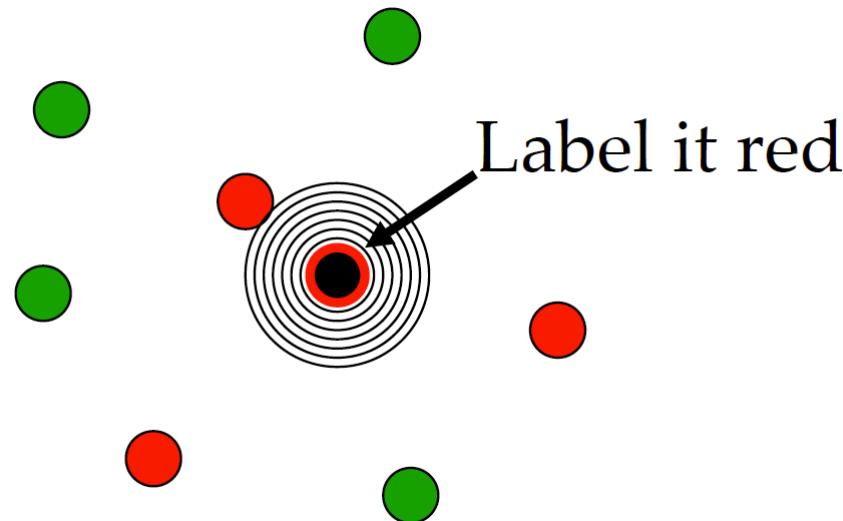


# Nearest Neighbors Classification

- ❖ Idea: instance-based learning
  - Similar examples have similar labels
  - Classify new examples like similar training examples
- ❖ Algorithm
  - Given some new example  $x$  for which we need to predict its class  $y$
  - Find the most similar training examples
  - Classify  $x$  “like” these most similar examples
- ❖ Questions:
  - How to determine similarity?
  - How many similar training examples to consider?

# 1-Nearest Neighbor

- ❖ One of the simplest classifiers
- ❖ Basic idea: label a new point the same as the closest known data instance
- ❖ E.g.,



# 1-Nearest Neighbor (Cont'd)

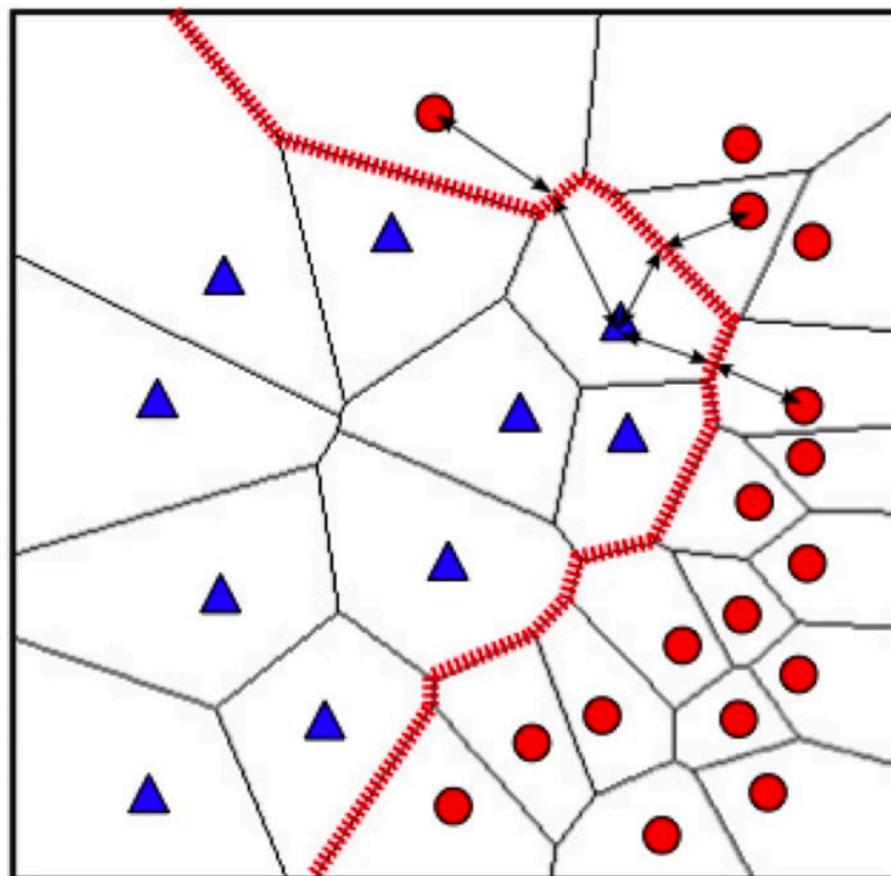
---



- ❖ A distance metric (to measure similarity)
  - Euclidean distance is commonly-used
  - When different units are used for each dimension
    - Standardization can apply
  - For categorical data, we can use hamming distance
    - $d(x_1, x_2) = \text{number of features on which } x_1 \text{ and } x_2 \text{ differ}$
  - Others (e.g., cosine, Manhattan)
- ❖ How many nearby neighbors to look at?
  - Only one
- ❖ How to fit with training data instances?
  - Just predict the same output as the nearest neighbor

# 1-Nearest Neighbor (Cont'd)

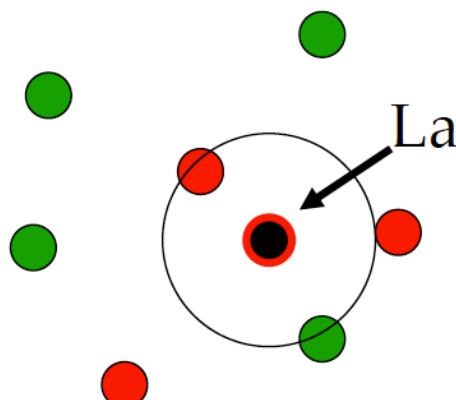
- ❖ The resultant space partition is a **Voronoi diagram**



# K-Nearest Neighbor

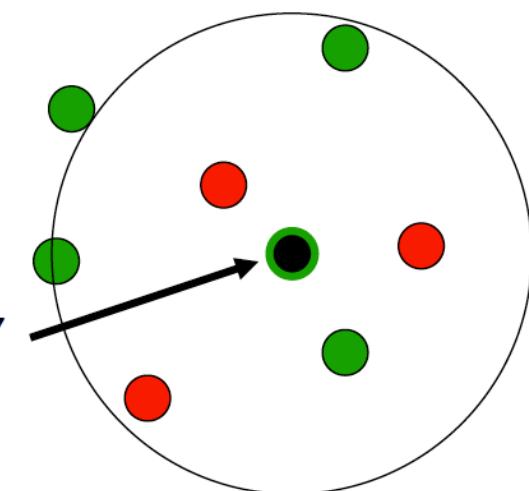


- ❖ Generalizes 1 NN to smooth away **noise** in the labels
- ❖ A new data instance is now assigned the **most frequent** label of its  $k$  nearest neighbors
- ❖ E.g.,  $k = 3$ , and  $k = 7$



Label it red, when  $k = 3$

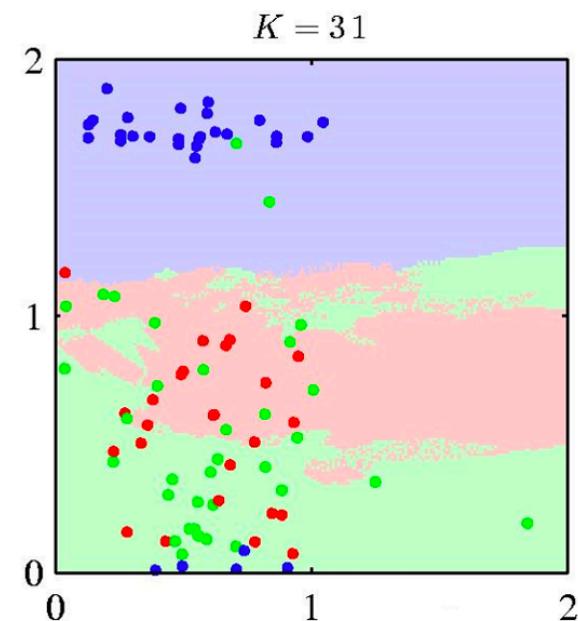
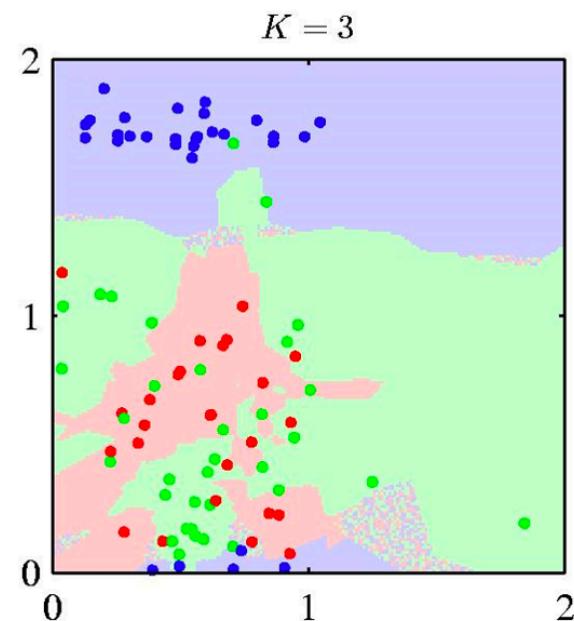
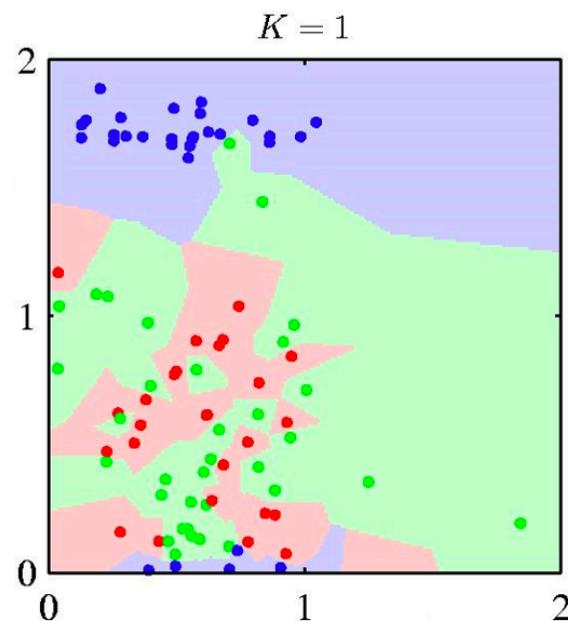
Label it green, when  $k = 7$



# *K*-Nearest Neighbor (Cont'd)



- ❖  $K$  acts as a **smoother** and controls **model complexity**



# Selecting $k$

---



- ❖ Increase  $k$ 
  - Make KNN less sensitive to noise
  - Avoid overfitting
- ❖ Decrease  $k$ 
  - Allow capturing finer structure of space
  - Avoid underfitting
- ❖ Pick  $k$  not too large, but not too small
  - Data-specific
  - This is model selection by hyperparameter tuning
  - Cross validation can be used to find a suitable  $k$
  - Theoretically, this is related to **bias-variance tradeoff**