

## Classification

- The goal of classification is to take an input vector  $\mathbf{x}$  and to assign  $y$  to one of  $K$  discrete classes  $C_k$  where  $k = 1, \dots, K$ , i.e.,

$$f : \mathbf{x} \in \mathbb{R}^D \rightarrow y \in \{1, \dots, K\}$$

- The input space is divided into  $K$  regions (one for each class), whose boundaries are called decision boundaries or decision surfaces.

Linear  
classification



## Classification Approaches

- Deterministic approach:** Construct a **discriminant function**  $f(\mathbf{x})$  that directly maps each input vector to a specific class
  - Linear deterministic classifier**
    - SVM
    - Perceptron
- Probabilistic Approach:** Model the **probability distribution**  $\mathbf{x}$  and  $y$  and then use this distribution to make optimal decisions
  - Discriminative approach –  $p(y|\mathbf{x})$
  - Generative approach –  $p(\mathbf{x}, y)$

## Linear Deterministic Classifiers

- Construct a linear discriminate function

$$f(\mathbf{x}, \mathbf{w}) = \mathbf{w}^t \mathbf{x} + w_0$$

- Linear classification

$$y(\mathbf{x}, \mathbf{w}) = \phi(f(\mathbf{x}, \mathbf{w})) \quad y \in \{1, 2, \dots, K\}$$

where  $\phi()$  is the activation function that maps  $\mathbf{x}$  into  $y$  that belongs to one of  $K$  different classes.

## Linear Deterministic Classifiers

- Binary discriminant classifiers map input into two classes
- Multi-class discriminant classifiers map input into  $K$  classes,  $K > 2$ .

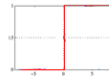
## Binary Discriminant Classifier

- Construct a linear discriminant function

$$f(\mathbf{x}) = \mathbf{w}^t \mathbf{x} + w_0$$

- Choose the *activation function as a step function*, i.e.,

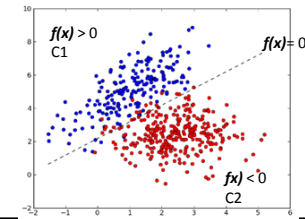
$$y(\mathbf{w}, \mathbf{x}) = \varphi(f(\mathbf{x})) = \begin{cases} C_1 & \text{if } f(\mathbf{x}) > 0 \\ C_2 & \text{else} \end{cases}$$



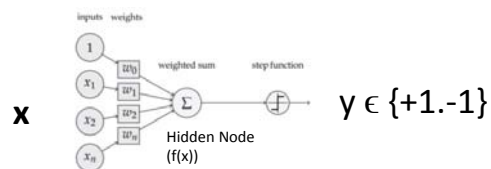
- Input vector  $\mathbf{x}$  is assigned to class  $C_1$  if the discriminant function  $f()$  is larger than 0 and to class  $C_2$  otherwise.

## Binary Discriminant (cont'd)

Discriminant function  $f(\mathbf{x}) = \mathbf{w}^t \mathbf{x} + w_0 = 0$  corresponds to a  $(N - 1)$ -dimensional hyperplane within the  $N$ -dimensional input space.  $\mathbf{w}$  determines the orientation and location of the decision surface. For 2D  $\mathbf{x}$ ,  $\mathbf{w}^t \mathbf{x} + w_0 = 0$  is a line



## The Perceptron Algorithm



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

where nonlinear activation function  $f()$  is given by a step function:

$$y(\mathbf{x}) = \begin{cases} +1 & \text{if } \mathbf{w}^t \mathbf{x} + w_0 > 0 \\ -1 & \text{else} \end{cases}$$

## Binary Support Vector Machine

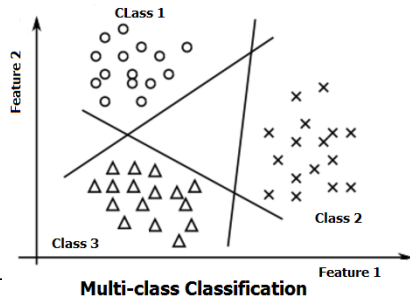
Two-class classification with the linear model is

$$f(\mathbf{x}) = \mathbf{w}^t \mathbf{x} + w_0 = \mathbf{w}^t \mathbf{x} + w_0$$

Given the target  $y = \{-1, +1\}$ , an input  $\mathbf{x}$  is classified into +1 if  $f(\mathbf{x}) > 0$  and to -1 else

### Multi-class discriminant classifiers

- The output  $y$  represents  $K$  classes ( $K > 2$ ). The goal is to map input  $\mathbf{x}$  into multiple classes



### Multi-class discriminant classifiers

- Construct  $K$  linear discriminant functions, one for each class :

$$f_k(\mathbf{x}) = \mathbf{w}_k^t \mathbf{x} + w_{k,0}$$

- Apply the max activation function (pick the max), assign  $\mathbf{x}$  to class  $C_k$ , i.e.,

$$y = \arg \max_k (f_k(\mathbf{x}))$$

### Linear Probabilistic Classifiers

Model the **probability distribution  $\mathbf{x}$  and  $y$**  and then use this distribution to make optimal decisions

- Discriminative approach –  $p(y|\mathbf{x})$
- Generative approach –  $p(\mathbf{x}, y)$

### Linear Discriminative Classifiers

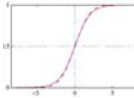
- Construct the probability of  $p(y|\mathbf{x})$  and perform classification based on  $p(y|\mathbf{x})$ , i.e.,

$$y^* = \arg \max_{y \in \{1, 2, \dots, K\}} p(y|\mathbf{x})$$

## Binary Discriminative Classifier

- Construct the discriminant function  
 $f(\mathbf{x}) = \mathbf{w}^t \mathbf{x} + w_0$
- Compute  $p(y|x)$  by choosing the sigmoid activity function, i.e.,  

$$p(y|x, \mathbf{w}) = \sigma(f(\mathbf{x})) = \frac{1}{1 + e^{-y * f(\mathbf{x})}}$$
- The sigmoid function  $\sigma()$  maps the discriminant function  $f(\mathbf{x})$  into a number between 0 and 1 (probability).
- $\mathbf{x}$  is classified into  $y=1$  if  $p(y=1|x) > 0.5$  and  $y=-1$  otherwise.
- This is the so-called logistic regression



## Multiclass Discriminative Models

- Output  $y \in \{1, 2, \dots, K\}$   $K > 2$  and construct  $P(y|x)$  by a *multiclass sigmoid function* ( $\sigma_M$ ) using  $K$  discriminant functions  $f_k(\mathbf{x}, \mathbf{w}_k, w_{k,0})$

$$p(y = k|x) = \sigma_M(f_k(\mathbf{x}, \mathbf{w}_k, w_{k,0})) = \frac{\exp(\mathbf{w}_k^t \mathbf{x} + w_{k,0})}{\sum_{j=1}^K \exp(\mathbf{w}_j^t \mathbf{x} + w_{j,0})}$$

$$k^* = \arg \max_k p(y = k | \mathbf{x})$$

- This is called softmax classifier or multinomial logistic regression

## Classifier Summary

- Deterministic Classifiers**
  - Binary linear classifier discriminant function
    - $f(\mathbf{x}) = \mathbf{w}^t \mathbf{x} + w_0$  - binary discriminant function
    - $y = \phi(f(\mathbf{x})) - \phi()$  is activation function
    - $\phi()$  is step function - binary discriminant classifier
      - Perceptron
      - SVM
  - Multi-class classifier
    - $f_k(\mathbf{x}) = \mathbf{w}_k^t \mathbf{x} + w_{k,0}$  - construct  $K$  discriminant functions
    - $y = \phi(f_k(\mathbf{x})) = \arg \max_k f_k(\mathbf{x})$  - multi-class discriminant classifier

## Classifier Summary

- Discriminative Classifiers  $p(y|x)$** 
  - Binary classifier
    - $p(y|x) = \sigma(yf(\mathbf{x}))$  -  $\phi()$  = sigmoid function - logistic regression
  - Multi-class classifier
    - $p(y=k|x) = \sigma_M(f_k(\mathbf{x})) = \frac{e^{f_k(\mathbf{x})}}{\sum_{k=1}^K e^{f_k(\mathbf{x})}}$
    - $\sigma_M$  is multi-class sigmoid function

## Classifier Learning

- Deterministic classifier learning
  - Binary discriminant classifier learning
  - Multi-class discriminant classifier learning
- Discriminative classifier learning
  - Binary classifier learning
  - Multi-class classifier learning

## Binary Discriminant Classifier Learning

- Given the training data  $\mathbf{D}=\{\mathbf{x}[m], y[m]\}$ ,  $m=1,2,\dots,M$ , learn the model parameters  $\Theta$  by minimizing the overall loss function  $L(\mathbf{D}:\Theta)$  of the form

$$L(\mathbf{D}:\Theta) = \frac{1}{M} \sum_{m=1}^M l(\mathbf{x}[m], y[m], \Theta) + \lambda R(\Theta)$$

where the first term is the individual loss function and the second term is the regularization

## Binary Classifier Loss Functions

For binary classifier  $y \in \{-1, +1\}$

- 0-1 loss function – the loss is 0 if the predicted labels and groundtruth labels are the same and 1 otherwise

$$l_{0/1} = \begin{cases} 0 & \text{if } y(\mathbf{w}^T \mathbf{x} + w_0) > 0 \\ 1 & \text{else} \end{cases}$$

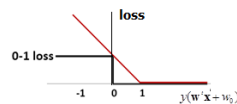
- Hinge loss function

$$l_h = \max(0, 1 - y(\mathbf{w}^T \mathbf{x} + w_0))$$

when  $y(\mathbf{w}^T \mathbf{x} + w_0) > 1$  the hinge loss is zero.

Otherwise the loss increases linearly with  $\mathbf{w}^T \mathbf{x} + w_0$ .

Hinge loss is an upper bound of 0/1 loss



## Classifier Learning with 0/1 Loss Function

- For the 0/1 loss function, the total loss function can be written as

$$L_{0/1}(\mathbf{D} : \mathbf{w}, w_0) = \frac{1}{M} \sum_{m=1}^M l_{0/1}(y[m](\mathbf{w}^T \mathbf{x}[m] + w_0)) + \lambda R(\mathbf{w}, w_0)$$

- The parameters can be obtained by minimizing the loss function, i.e.,

$$\mathbf{w}^*, w_0^* = \arg \min_{\mathbf{w}, w} L_{0/1}(\mathbf{D} : \mathbf{w}, w_0)$$

### Classifier Learning with 0/1 Loss Function

- Gradient-based method cannot solve for the parameters by minimizing the total 0/1 loss function  $L_{0/1}(\mathbf{D} : \mathbf{w}, w_0)$  since the gradient of  $l_{0/1}(y[m](\mathbf{w}'\mathbf{x}[m] + w_0))$  is zero everywhere
- Minimizing the overall 0/1 loss function is intractable
- Solution is to minimize an alternative loss function such as the hinge loss function

### Binary Classifier Learning with Hinge Loss

- For the hinge loss function, the total loss function can be written as

$$L_H(D : \mathbf{w}, w_0) = \frac{1}{M} \sum_{m=1}^M \max(0, 1 - y[m](\mathbf{w}'\mathbf{x}[m] + w_0)) + \lambda R(\mathbf{w}, w_0)$$

- The parameters can be obtained by minimizing the loss function, i.e.,

$$\mathbf{w}^*, w_0^* = \arg \min_{\mathbf{w}, w} L_H(\mathbf{D} : \mathbf{w}, w_0)$$

### Binary Classifier Learning with Hinge Loss

Let

$$\Theta = (\mathbf{w}, w_0)^T \quad \mathbf{X} = (\mathbf{x}, 1)^T \quad z[m] = y[m](\mathbf{w}'\mathbf{x}[m] + w_0)$$

$$L_H(D : \Theta) = \frac{1}{M} \sum_{m=1}^M \max(0, 1 - y[m](\mathbf{w}'\mathbf{x}[m] + w_0)) + \lambda R(\mathbf{w}, w_0)$$

$$= \frac{1}{M} \sum_{m=1}^M \max(0, 1 - y[m]\Theta' \mathbf{X}[m]) + \lambda R(\Theta)$$

$$= \frac{1}{M} \sum_{m=1}^M \max(0, 1 - z[m]) + \lambda R(\Theta)$$

$$\frac{\partial L_H(D : \Theta)}{\partial \Theta} = \frac{1}{M} \sum_{m=1}^M \frac{\partial \max(0, 1 - z[m])}{\partial z} \frac{\partial z}{\partial \Theta} + \lambda \frac{\partial R(\Theta)}{\partial \Theta}$$

$$\frac{\partial \max(0, 1 - z[m])}{\partial z} = \begin{cases} 0 & \text{if } z[m] > 1 \\ -1 & \text{else} \end{cases}$$

$$\frac{\partial z}{\partial \Theta} = \frac{\partial [y[m](\mathbf{w}'\mathbf{x}[m] + w_0)]}{\partial \Theta} = \frac{\partial [y[m]\Theta' \mathbf{X}[m]]}{\partial \Theta} = y[m]\mathbf{X}[m]$$

$$\frac{\partial L_H(D : \Theta)}{\partial \Theta} = \frac{-1}{M} \sum_{m=1}^M y[m]\mathbf{X}[m] I[z[m] > 1] + \lambda \frac{\partial R(\Theta)}{\partial \Theta}, \text{ where } I[z[m]] = \begin{cases} 0 & \text{if } z[m] > 1 \\ 1 & \text{else} \end{cases}$$

Similarly, we can compute  $\frac{\partial R(\Theta)}{\partial \Theta}$ . Given  $\frac{\partial L_H(D : \Theta)}{\partial \Theta}$ , as there is no analytic solution, gradient descent can apply to update  $\Theta$  iteratively.

### Multiclass Discriminative Classifier Learning

- Suppose  $y \in \{1, 2, \dots, K\}$ , we can construct a discriminate function for each class, i.e.,  $f_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k,0}$
- A generalized hinge loss for K class can be constructed as follows
 
$$l_h(\mathbf{x}, y = k, \mathbf{W}) = \max(0, 1 + \max_{k' \neq k} (\mathbf{w}_{k'}^T \mathbf{x} + w_{k',0}) - \mathbf{w}_k^T \mathbf{x} - w_{k,0})$$
 where  $\mathbf{W} = (\mathbf{w}_1, w_{1,0}, \mathbf{w}_2, w_{2,0}, \dots, \mathbf{w}_K, w_{K,0})^T$
- $\mathbf{W}$  can be learnt by minimizing the generalized hinge loss function

## Discriminative Classifier Learning

- Binary logistic regression learning
- Multi-class logistic regression learning

## Binary Logistic Regression Learning

- Given the training data  $\mathbf{D}=\{\mathbf{x}[m], y[m]\}$ ,  $m=1,2,\dots,M$  and  $y[m] \in \{-1,+1\}$ , learn the model parameters  $\Theta$  by minimizing the negative log conditional likelihood, i.e.,

$$\begin{aligned}
 -LCL(\mathbf{D} : \Theta) &= -\sum_{m=1}^M \log p(y[m] | \mathbf{x}[m]) = -\sum_{m=1}^M \log \sigma(yf(\mathbf{x}[m])) \\
 &= -\sum_{m=1}^M \log \left\{ \frac{1}{1 + e^{-y[m](\mathbf{w}'\mathbf{x}[m] + w_0)}} \right\} \\
 &= \sum_{m=1}^M \log \{ 1 + e^{-y[m]\Theta'\mathbf{X}[m]} \}
 \end{aligned}$$

where  $\Theta = \begin{bmatrix} \mathbf{w} \\ w_0 \end{bmatrix}$      $\mathbf{X}[m] = \begin{bmatrix} \mathbf{x}[m] \\ 1 \end{bmatrix}$      $\sigma$  is the sigmoid function

## Binary Logistic Regression Learning (cont'd)

- $\Theta$  is found by minimizing the negative log likelihood function, i.e.,

$$\Theta^* = \arg \min_{\Theta} -LCL(\mathbf{D} : \Theta)$$

- There is no closed-form solution. Gradient descent can be used to estimate  $\Theta$  iteratively, i.e.,

$$\Theta^t = \Theta^{t-1} - \eta \frac{\partial -LCL(\mathbf{D} : \Theta^t)}{\partial \Theta} \Big|_{\Theta^t}$$

where

$$\begin{aligned}
 \frac{\partial -LCL(\mathbf{D} : \Theta^t)}{\partial \Theta} &= \frac{\partial \sum_{m=1}^M \log \{ 1 + e^{-y[m]\Theta^t'\mathbf{X}[m]} \}}{\partial \Theta} \\
 &= -\sum_{m=1}^M \frac{e^{-y[m]\Theta^t'\mathbf{X}[m]}}{1 + e^{-y[m]\Theta^t'\mathbf{X}[m]}} y[m]\mathbf{X}[m]
 \end{aligned}$$

## Binary Logistic Regression Learning (cont'd)

Alternatively, the gradient can be computed as

$$\begin{aligned}
 \frac{\partial -LCL(\mathbf{D} : \Theta^t)}{\partial \Theta} &= -\sum_{m=1}^M \frac{\partial \log \sigma(y[m]f(\mathbf{x}[m]))}{\partial \Theta} \\
 &= -\sum_{m=1}^M \frac{1}{\sigma(y[m]f(\mathbf{x}[m]))} \frac{\partial \sigma(y[m]f(\mathbf{x}[m]))}{\partial (y[m]f(\mathbf{x}[m]))} \frac{\partial (y[m]f(\mathbf{x}[m]))}{\partial \Theta} \\
 &= -\sum_{m=1}^M \frac{1}{\sigma(y[m]f(\mathbf{x}[m]))} \sigma(y[m]f(\mathbf{x}[m]))(1 - \sigma(y[m]f(\mathbf{x}[m]))) y[m]\mathbf{X}[m] \\
 &= -\sum_{m=1}^M (1 - \sigma(y[m]f(\mathbf{x}[m]))) y[m]\mathbf{X}[m]
 \end{aligned}$$

## Multiclass Logistic Regression Learning

- Suppose  $y \in \{1, 2, \dots, K\}$ , we can construct a discriminant function for each class, i.e.,  $f_k(\mathbf{x}) = \mathbf{w}_k^t \mathbf{x} + w_{k,0} = \Theta_k^t \mathbf{X}$ , where  $\Theta_k = [\mathbf{w}_k^t, w_{k,0}]^t$   $\mathbf{X} = [\mathbf{x} \ 1]^t$
- For output  $y$ , we use 1-of-K encoding, where we use a  $K \times 1$  binary vector  $\mathbf{y}$ , whose contains a single 1 for element  $k$  (the correct class) and 0 elsewhere.
- For 5 classes and the input belongs to class 2,  $\mathbf{y}$  is

$$\mathbf{y} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

## Multiclass Logistic Regression Learning

Given training data  $D = \{\mathbf{X}[m], y[m]\}$ ,  $m=1, 2, \dots, M$ , learn the parameters  $\Theta$  by minimizing the total negative log conditional likelihood

$$\begin{aligned} L(\mathbf{D} : \Theta) &= -\sum_{m=1}^M \log p(y[m] | \mathbf{X}[m], \Theta) = -\sum_{m=1}^M \log \prod_{k=1}^K [p(y[m] = k | \mathbf{X}[m], \Theta)]^{y[m][k]} \\ &= -\sum_{m=1}^M \sum_{k'=1}^K y[m][k'] \log p(y[m] = k' | \mathbf{X}[m], \Theta) \\ &= -\sum_{m=1}^M \sum_{k'=1}^K y[m][k'] \log \frac{\exp(\Theta_{k'}^t \mathbf{X}[m])}{\sum_{j=1}^K \exp(\Theta_j^t \mathbf{X}[m])} \end{aligned}$$

## Multiclass Logistic Regression Learning

$$\begin{aligned} L(\mathbf{D} : \Theta) &= -\sum_{m=1}^M \log p(y[m] | \mathbf{X}[m], \Theta) = -\sum_{m=1}^M \log \prod_{k=1}^K [p(y[m] = k | \mathbf{X}[m], \Theta)]^{y[m][k]} \\ &= -\sum_{m=1}^M \log \prod_{k=1}^K [\sigma_M(f(\mathbf{X}[m], \Theta_k))]^{y[m][k]} = -\sum_{m=1}^M \log \prod_{k=1}^K \left( \frac{\exp(\Theta_k^t \mathbf{X}[m])}{\sum_{j=1}^K \exp(\Theta_j^t \mathbf{X}[m])} \right)^{y[m][k]} \\ &= -\sum_{m=1}^M \sum_{k=1}^K y[m][k] \log \frac{\exp(\Theta_k^t \mathbf{X}[m])}{\sum_{j=1}^K \exp(\Theta_j^t \mathbf{X}[m])} \\ &= -\sum_{m=1}^M \sum_{k=1}^K y[m][k] \left( \Theta_k^t \mathbf{X}[m] - \log \sum_{j=1}^K \exp(\Theta_j^t \mathbf{X}[m]) \right) \end{aligned}$$

$\Theta^* = \arg \min_{\Theta} L(\mathbf{D} : \Theta)$   $\sigma_m$  is the multi-class sigmoid function

## Multiclass Logistic Regression Learning

$$\frac{\partial L(\mathbf{D} : \Theta)}{\partial \Theta} = \begin{bmatrix} \frac{\partial L(\mathbf{D} : \Theta)}{\partial \Theta_1} \\ \frac{\partial L(\mathbf{D} : \Theta)}{\partial \Theta_2} \\ \vdots \\ \frac{\partial L(\mathbf{D} : \Theta)}{\partial \Theta_K} \end{bmatrix}$$



## Multiclass Logistic Regression Learning

$$\begin{aligned}
 \frac{\partial L(\mathbf{D}; \Theta)}{\partial \Theta_k} &= \frac{\partial \sum_{m=1}^M \sum_{k'=1}^K y[m][k'] \log \frac{\exp(\Theta_k^T \mathbf{X}[m])}{\sum_{j=1}^K \exp(\Theta_j^T \mathbf{X}[m])}}{\partial \Theta_k} \\
 &= - \frac{\partial \sum_{m=1}^M \sum_{k'=1}^K \{y[m][k'] \Theta_k^T \mathbf{X}[m] - y[m][k'] \log \sum_{j=1}^K \exp(\Theta_j^T \mathbf{X}[m])\}}{\partial \Theta_k} \\
 &= - \sum_{m=1}^M \{y[m][k] \mathbf{X}[m] - \frac{\partial \sum_{k'=1}^K y[m][k'] \log \sum_{j=1}^K \exp(\Theta_j^T \mathbf{X}[m])}{\partial \Theta_k}\} \\
 &= - \sum_{m=1}^M \{y[m][k] \mathbf{X}[m] - \frac{\exp(\Theta_k^T \mathbf{X}[m]) \mathbf{X}[m]}{\sum_{j=1}^K \exp(\Theta_j^T \mathbf{X}[m])}\} = - \sum_{m=1}^M \mathbf{X}[m] \{y[m][k] - \sigma_M(\Theta_k^T \mathbf{X}[m])\}
 \end{aligned}$$

## Multiclass Logistic Regression Learning

$\Theta$  can be solved iteratively gradient descent

$$\begin{aligned}
 \nabla_{\Theta_k} L(\mathbf{D}; \Theta) &= \frac{\partial L(\mathbf{D}; \Theta)}{\partial \Theta_k} = - \frac{\partial \sum_{m=1}^M \sum_{k'=1}^K \left[ y[m][k'] \left( \Theta_k^T \mathbf{X}[m] - \log \sum_{j=1}^K \exp(\Theta_j^T \mathbf{X}[m]) \right) \right]}{\partial \Theta_k} \\
 &= - \sum_{m=1}^M \sum_{k'=1}^K \left[ y[m][k'] \left( \Theta_k^T \mathbf{X}[m] - \log \sum_{j=1}^K \exp(\Theta_j^T \mathbf{X}[m]) \right) \right] = - \sum_{m=1}^M \left\{ y[m][k] \mathbf{X}[m] - \frac{\partial \sum_{k'=1}^K y[m][k'] \log \sum_{j=1}^K \exp(\Theta_j^T \mathbf{X}[m])}{\partial \Theta_k} \right\} \\
 &= - \sum_{m=1}^M \left\{ y[m][k] \mathbf{X}[m] - \sum_{k'=1}^K y[m][k'] \frac{\exp(\Theta_k^T \mathbf{X}[m])}{\sum_{j=1}^K \exp(\Theta_j^T \mathbf{X}[m])} \mathbf{X}[m] \right\} = - \sum_{m=1}^M \left\{ y[m][k] \mathbf{X}[m] - \frac{\exp(\Theta_k^T \mathbf{X}[m])}{\sum_{j=1}^K \exp(\Theta_j^T \mathbf{X}[m])} \mathbf{X}[m] \sum_{k'=1}^K y[m][k'] \right\} \\
 &= - \sum_{m=1}^M \mathbf{X}[m] \left\{ y[m][k] - \frac{\exp(\Theta_k^T \mathbf{X}[m])}{\sum_{j=1}^K \exp(\Theta_j^T \mathbf{X}[m])} \right\} \\
 \Theta_k^{t+1} &= \Theta_k^t - \eta \nabla_{\Theta_k} L(\mathbf{D}; \Theta) \\
 &\text{Repeat this for } k=1, 2, \dots, K
 \end{aligned}$$

## Multiclass Logistic Regression Learning

Alternatively, the gradient can be compute as

$$\begin{aligned}
 \frac{\partial L(\mathbf{D}; \Theta)}{\partial \Theta_k} &= - \frac{\partial \sum_{m=1}^M \log \left[ \prod_{k'=1}^K [\sigma_M(f(\mathbf{X}[m], \Theta_k))]^{y[m][k']} \right]}{\partial \Theta_k} = - \frac{\partial \sum_{m=1}^M \sum_{k'=1}^K y[m][k'] \log \sigma_M(f(\mathbf{X}[m], \Theta_k))}{\partial \Theta_k} \\
 &= - \sum_{m=1}^M \sum_{k'=1}^K y[m][k'] \frac{\partial \log \sigma_M(f(\mathbf{X}[m], \Theta_k))}{\partial \Theta_k} = - \sum_{m=1}^M \left[ y[m][k] \frac{\partial \log \sigma_M(f(\mathbf{X}[m], \Theta_k))}{\partial \Theta_k} + \sum_{k'=1, k' \neq k}^K y[m][k'] \frac{\partial \log \sigma_M(f(\mathbf{X}[m], \Theta_k))}{\partial \Theta_k} \right] \\
 &= - \sum_{m=1}^M \left[ y[m][k] (1 - \sigma_M(f(\mathbf{X}[m], \Theta_k)) \mathbf{X}[m]) - \sum_{k'=1, k' \neq k}^K y[m][k'] \frac{\sigma_M(f(\mathbf{X}[m], \Theta_k)) \sigma_M(f_k(\mathbf{X}[m], \Theta_k))}{\sigma_M(f(\mathbf{X}[m], \Theta_k))} \mathbf{X}[m] \right] \\
 &= - \sum_{m=1}^M \left[ y[m][k] (1 - \sigma_M(f(\mathbf{X}[m], \Theta_k)) \mathbf{X}[m]) - \sum_{k'=1, k' \neq k}^K y[m][k'] \sigma_M(f(\mathbf{X}[m], \Theta_k)) \mathbf{X}[m] \right] \\
 &= - \sum_{m=1}^M \left[ y[m][k] (1 - \sigma_M(f(\mathbf{X}[m], \Theta_k)) \mathbf{X}[m]) - \sum_{k'=1, k' \neq k}^K y[m][k'] \sigma_M(f(\mathbf{X}[m], \Theta_k)) \mathbf{X}[m] \right] \\
 &= - \sum_{m=1}^M \mathbf{X}[m] [y[m][k] - \sigma_M(f(\mathbf{X}[m], \Theta_k))]
 \end{aligned}$$

Note

$$\frac{\partial \sigma_M(z)}{\partial z} = \sigma_M(z) (1 - \sigma_M(z))$$

## Multiclass Logistic Regression Learning

Add regularization

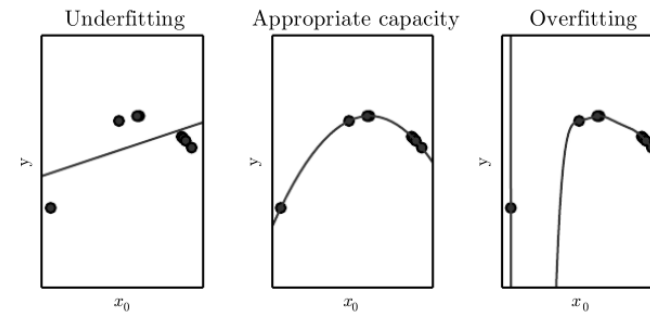
$$\Theta_k^{t+1} = \Theta_k^t - \eta [\nabla_{\Theta_k} L(\mathbf{D}; \Theta) + \lambda \nabla_{\Theta_k} R(\Theta_k)]$$

where  $R(\Theta)$  can be L1-norm or squared L2-norm

## Capacity, Overfitting, and Underfitting

- Capacity-a model's ability to fit a variety of functions. It is determined by the number of parameters the model has.
- Overfitting-a model performs well on training data but poorly on unseen testing data
- Underfitting-a model performs poorly both on training and testing data

## Capacity, Overfitting, and Underfitting



## Capacity, Overfitting, and Underfitting

