



ARISTOTLE UNIVERSITY OF THESSALONIKI



FACULTY OF ENGINEERING

# Pattern Recognition & Machine Learning

## *Linear Discriminant Functions and Models*

**Panagiotis C. Petrantonakis**

*Assistant Professor*

Dept. of Electrical and Computer Engineering

[ppetrant@ece.auth.gr](mailto:ppetrant@ece.auth.gr)

Fall Semester

# Until now...

---

- We have seen only **one** out of **three** different approaches of solving decision problems:
  - *Generative probabilistic models:*
    - solve the inference problem of determining the class-conditional density function  $p(\mathbf{x}|\omega_i)$  and the priors  $P(\omega_i)$  and then use the Bayes formula to find posterior class probabilities.

$$P(\omega_i | \mathbf{x}) = \frac{p(\mathbf{x}|\omega_i)P(\omega_i)}{p(\mathbf{x})}$$

- We used training data to infer  $p(\mathbf{x}|\omega_i)$ .

# Until now...

---

- We have seen only **one** out of **three** different approaches of solving decision problems:
  - *Generative probabilistic models:*
    - solve the inference problem of determining the class-conditional density function  $p(\mathbf{x}|\omega_i)$  and the priors  $P(\omega_i)$  and then use the Bayes formula to find posterior class probabilities.

$$P(\omega_i | \mathbf{x}) = \frac{p(\mathbf{x}|\omega_i)P(\omega_i)}{p(\mathbf{x})}$$

- We used training data to infer  $p(\mathbf{x}|\omega_i)$ .
- We have also talked about *discriminant functions*
  - In the above case, we replaced  $P(\omega_i | \mathbf{x})$  with  $y_i(\mathbf{x})$  (in the BDT lecture we use the symbol  $g$  instead of  $y$ ) and under certain assumptions about the class-conditional density  $p(\mathbf{x}|\omega_i)$  we ended up with linear or non-linear *decision boundaries*.
  - We say that the corresponding discriminant functions are *linear or non-linear discriminants*!

# Until now...

---

- We have seen only **one** out of **three** different approaches of solving decision problems:
  - *Generative probabilistic models:*
    - solve the inference problem of determining the class-conditional density function  $p(\mathbf{x}|\omega_i)$  and the priors  $P(\omega_i)$  and then use the Bayes formula to find posterior class probabilities.

$$P(\omega_i | \mathbf{x}) = \frac{p(\mathbf{x}|\omega_i)P(\omega_i)}{p(\mathbf{x})}$$

- We used training data to infer  $p(\mathbf{x}|\omega_i)$ .
- We have also talked about *discriminant functions*
  - In the above case, we replaced  $P(\omega_i | \mathbf{x})$  with  $y_i(\mathbf{x})$  (in the BDT lecture we use the symbol  $g$  instead of  $y$ ) and under certain assumptions about the class-conditional density  $p(\mathbf{x}|\omega_i)$  we ended up with linear or non-linear *decision boundaries*.
  - We say that the corresponding discriminant functions are **linear or non-linear discriminants**!
- In this lecture we will talk solely about linear discriminants of the form
$$y(\mathbf{x}) = \mathbf{w}^t \mathbf{x} + w_0$$

# Today...

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

- 
- We will see the rest **two** approaches of solving decision problems
    - *Linear discriminant functions* that map each input  $\mathbf{x}$  directly onto a class label.
      - For instance, in two-class classification problems  $y(\mathbf{x})$  takes binary values such that  $y(\mathbf{x}) = 0$  represents class  $\omega_1$  and  $y(\mathbf{x}) = 1$  represents class  $\omega_2$  (probabilistic modeling plays no role in this scheme)

# Today...

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

- 
- We will see the rest **two** approaches of solving decision problems
    - *Linear discriminant functions* that map each input  $\mathbf{x}$  directly onto a class label.
      - For instance, in two-class classification problems  $y(\mathbf{x})$  takes binary values such that  $y(\mathbf{x}) = 0$  represents class  $\omega_1$  and  $y(\mathbf{x}) = 1$  represents class  $\omega_2$  (probabilistic modeling plays no role in this scheme)
    - *Discriminative probabilistic models*
      - solve the inference problem of determining the posterior class probabilities  $P(\omega_i|\mathbf{x})$  and then assign each sample  $\mathbf{x}$  to one of the classes

# Today...

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

- 
- We will see the rest **two** approaches of solving decision problems
    - *Linear discriminant functions* that map each input  $\mathbf{x}$  directly onto a class label.
      - For instance, in two-class classification problems  $y(\mathbf{x})$  takes binary values such that  $y(\mathbf{x}) = 0$  represents class  $\omega_1$  and  $y(\mathbf{x}) = 1$  represents class  $\omega_2$  (probabilistic modeling plays no role in this scheme)
    - *Discriminative probabilistic models*
      - solve the inference problem of determining the posterior class probabilities  $P(\omega_i|\mathbf{x})$  and then assign each sample  $\mathbf{x}$  to one of the classes
  - For the former approach we will see the algorithms:
    - Least squares
    - Perceptron
  - For the latter approach we will see the algorithm:
    - Logistic Regression

# Today...

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

- We will see the rest **two** approaches of solving decision problems
  - *Linear discriminant functions* that map each input  $\mathbf{x}$  directly onto a class label.
    - For instance, in two-class classification problems  $y(\mathbf{x})$  takes binary values such that  $y(\mathbf{x}) = 0$  represents class  $\omega_1$  and  $y(\mathbf{x}) = 1$  represents class  $\omega_2$  (probabilistic modeling plays no role in this scheme)
  - *Discriminative probabilistic models*
    - solve the inference problem of determining the posterior class probabilities  $P(\omega_i|\mathbf{x})$  and then assign each sample  $\mathbf{x}$  to one of the classes
- For the former approach we will see the algorithms:
  - Least squares
  - Perceptron
- For the latter approach we will see the algorithm:
  - Logistic Regression
- In **all algorithms** that we will see in this lecture, we use training data to determine vector  $\mathbf{w}$  and  $w_0$ . (**big difference with the probabilistic generative models!**)

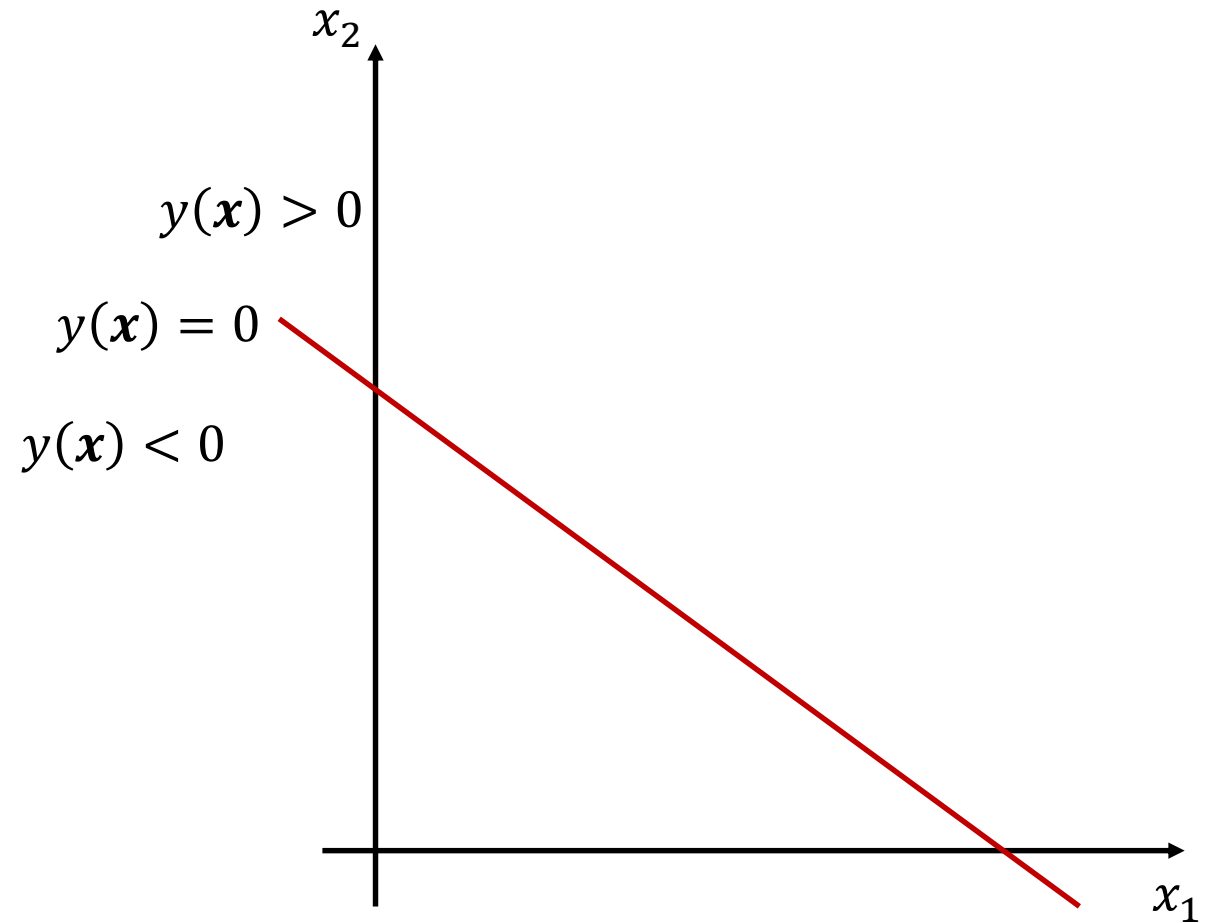


# Linear planes – Two classes

---

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

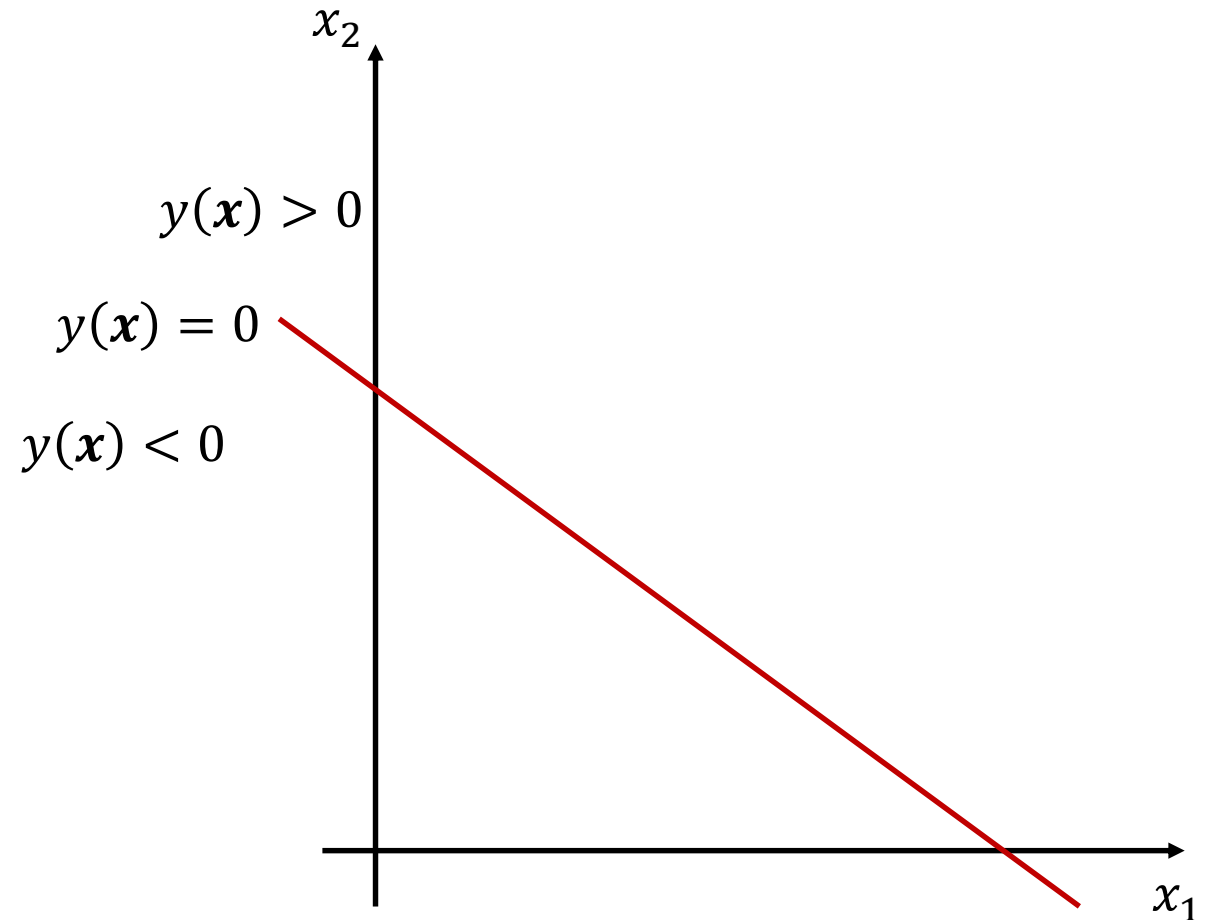
- We assign every input vector to  $\omega_1$  if  $y(\mathbf{x}) \geq 0$  and to  $\omega_2$  otherwise.



# Linear planes – Two classes

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

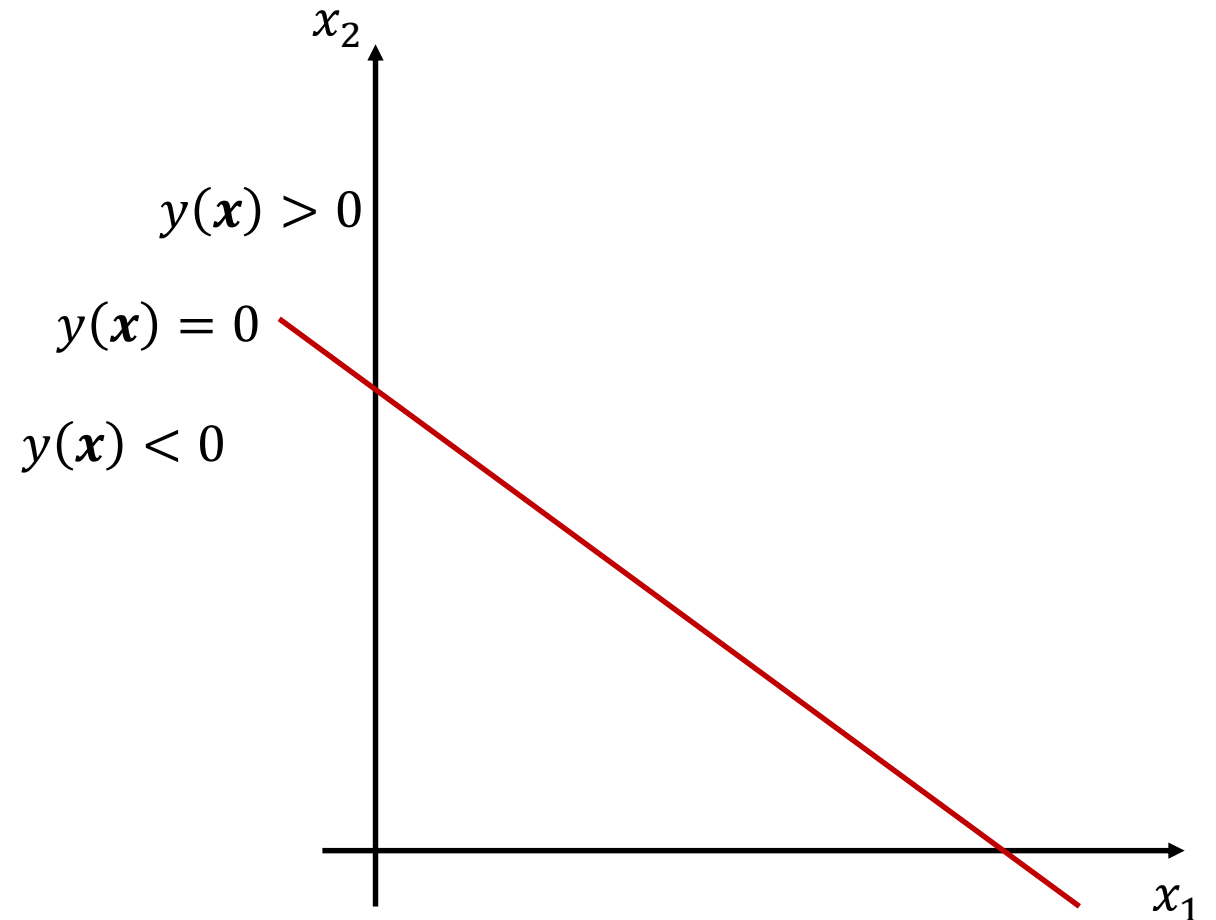
- We assign every input vector to  $\omega_1$  if  $y(\mathbf{x}) \geq 0$  and to  $\omega_2$  otherwise.
- The decision boundary (or surface) is defined by  $y(\mathbf{x}) = 0$ .
  - for a  $d$ -dimensional feature space I get a  $(d - 1)$ -dimensional (hyper) plane.



# Linear planes – Two classes

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

- We assign every input vector to  $\omega_1$  if  $y(\mathbf{x}) \geq 0$  and to  $\omega_2$  otherwise.
- The decision boundary (or surface) is defined by  $y(\mathbf{x}) = 0$ .
  - for a  $d$ -dimensional feature space I get a  $(d - 1)$ -dimensional (hyper) plane.
- $\mathbf{w}$  is the weight vector and  $w_0$  is the bias (or threshold)
- Can you tell how  $\mathbf{w}$  is related to every vector on the decision boundary  $y(\mathbf{x}) = 0$ ?



# Linear planes – Two classes

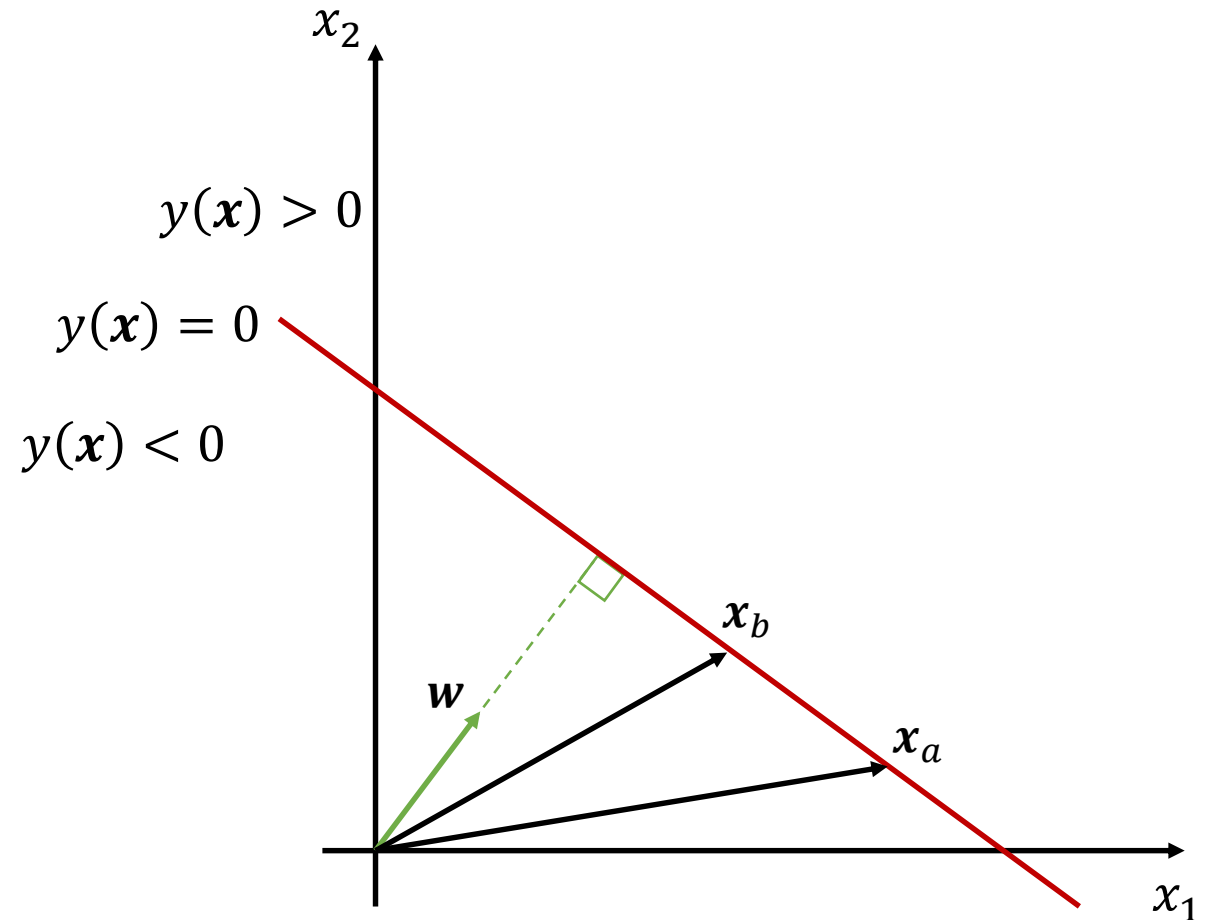
$$y(\mathbf{x}) = \mathbf{w}^t \mathbf{x} + w_0$$

- If we consider two vectors  $\mathbf{x}_a$  and  $\mathbf{x}_b$  that both lie on the decision boundary then:

$$\mathbf{w}^t \mathbf{x}_a + w_0 = \mathbf{w}^t \mathbf{x}_b + w_0 \Leftrightarrow$$

$$\mathbf{w}^t (\mathbf{x}_a - \mathbf{x}_b) = 0$$

- Thus,  $\mathbf{w}$  is orthogonal to the decision surface and to every vector on it.
  - Determines the orientation of the decision surface!
- What is the interpretation of the bias term  $w_0$ ?



# Linear planes – Two classes

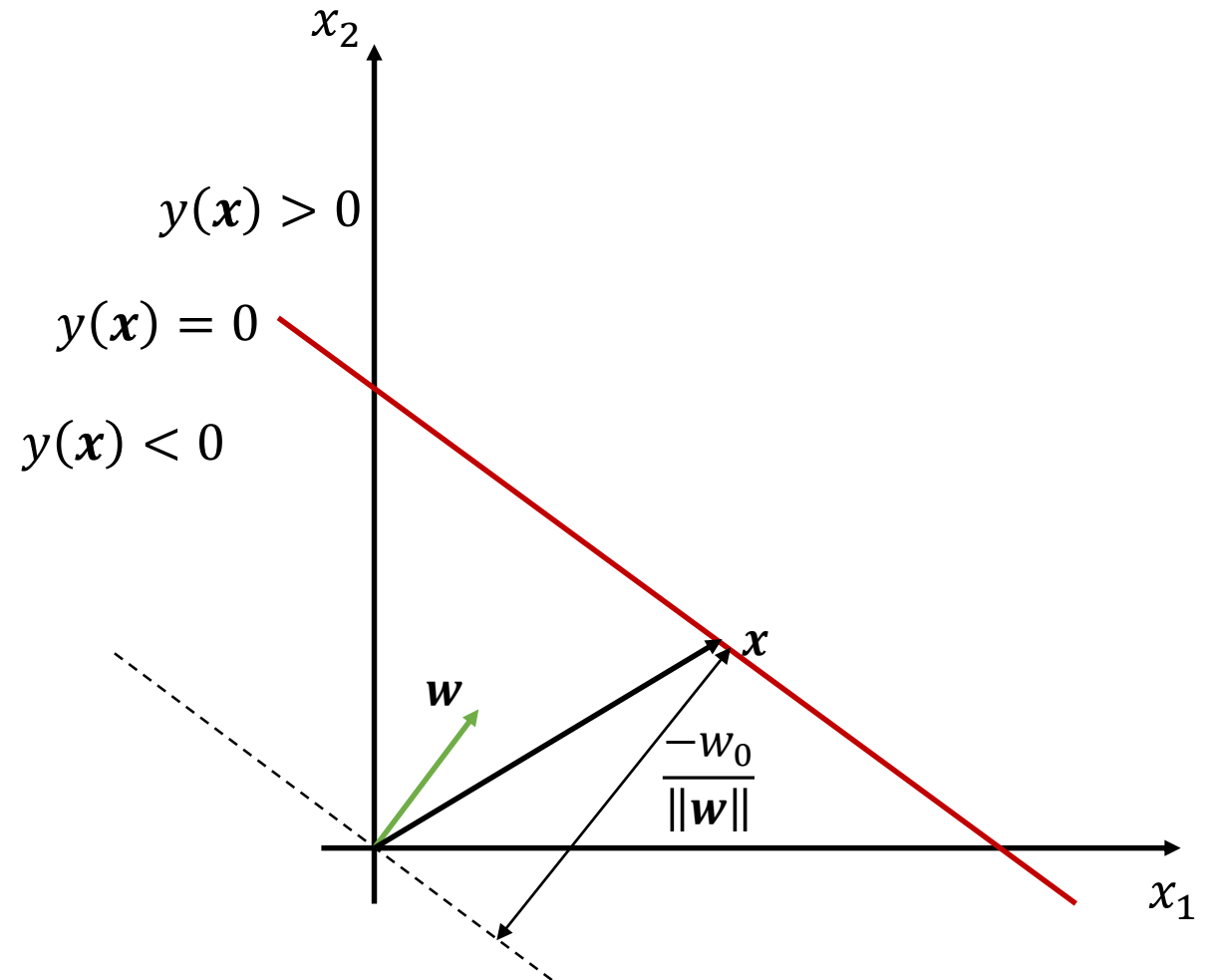
$$y(\mathbf{x}) = \mathbf{w}^t \mathbf{x} + w_0$$

- If we consider a vector  $\mathbf{x}$  that lie on the decision boundary then,  $y(\mathbf{x}) = 0$  and:

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

$$\frac{\mathbf{w}^t \mathbf{x}}{\|\mathbf{w}\|} = \frac{-w_0}{\|\mathbf{w}\|}$$

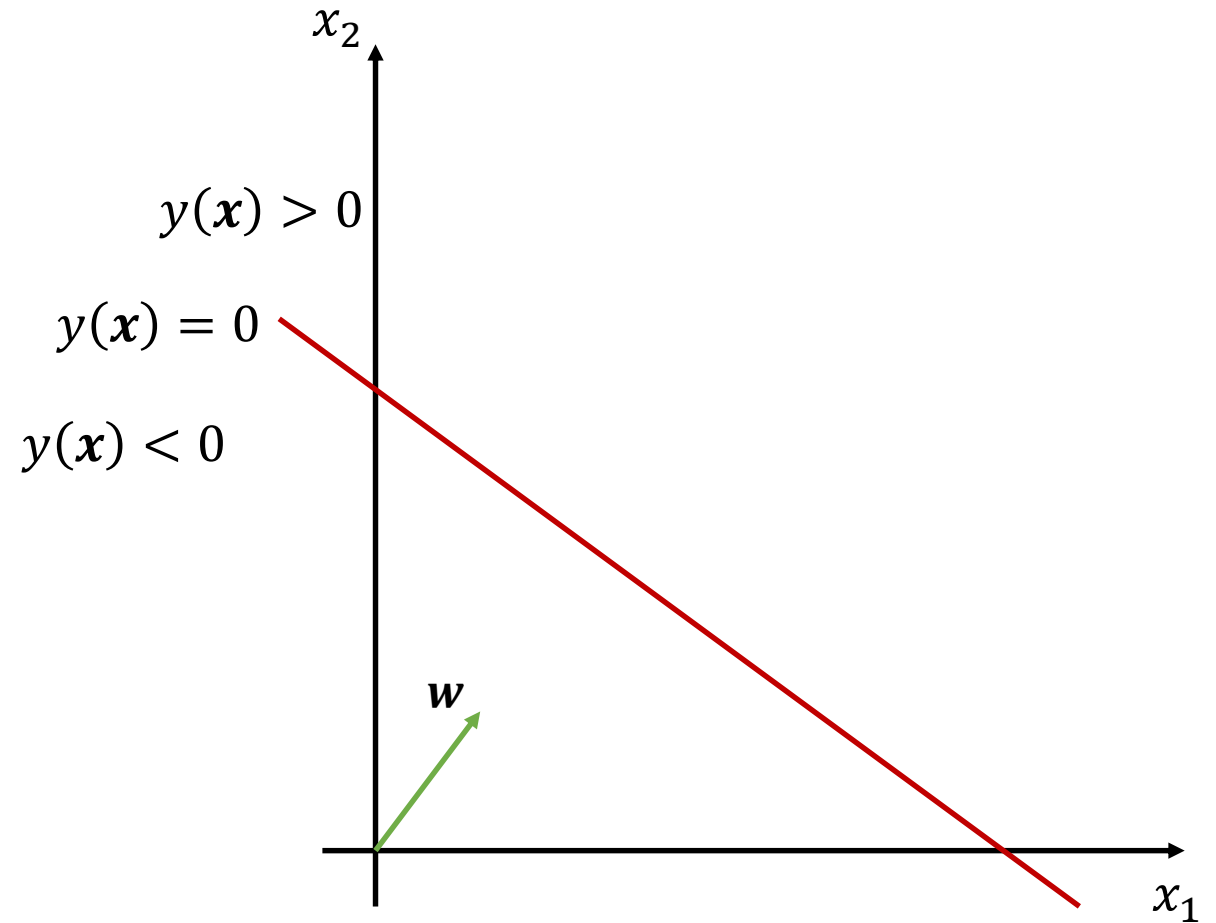
- Thus the bias term is actually the (signed) distance of the decision surface from the origin.
- It shifts the boundary away from the origin.



# Linear planes – Two classes

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

- Finally what is the interpretation of the  $y(\mathbf{x})$  value?



# Linear planes – Two classes

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

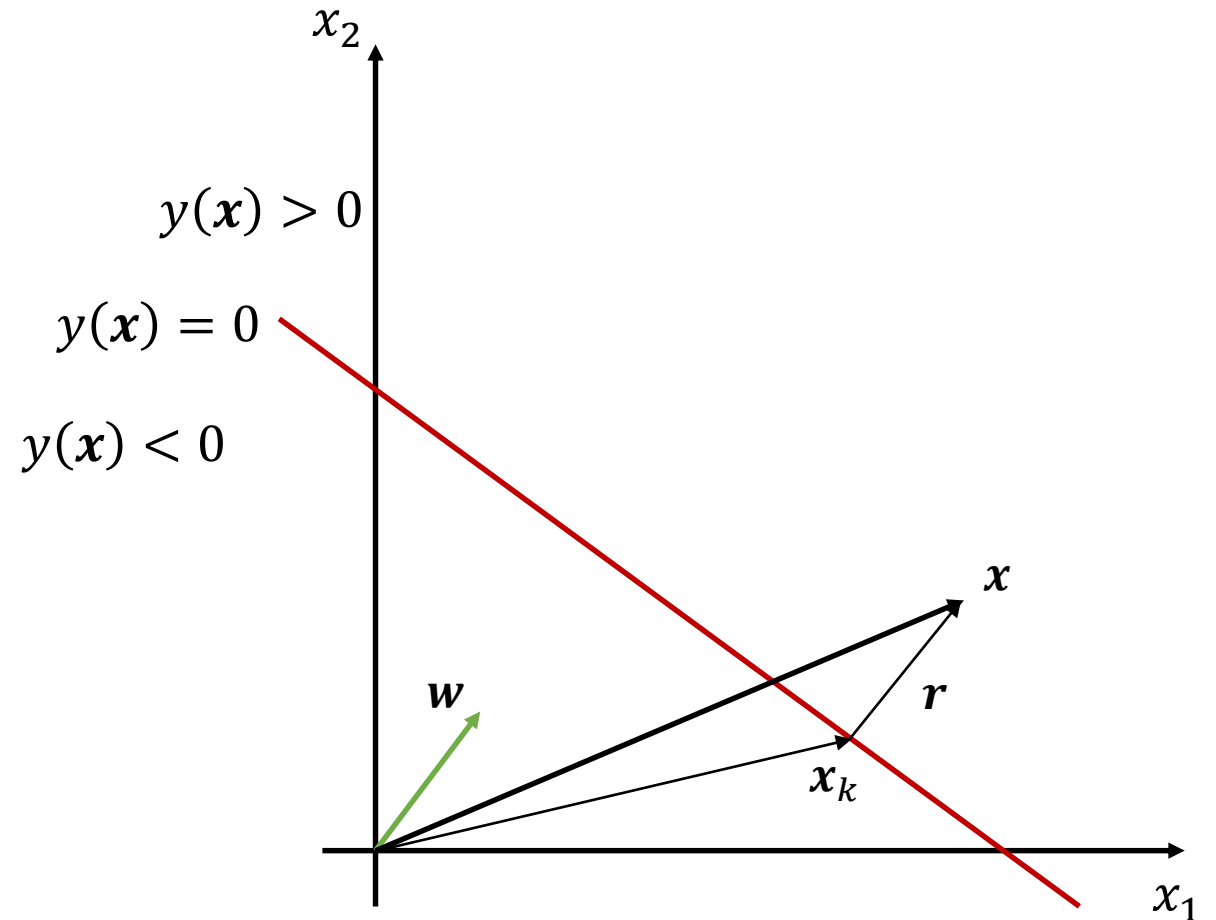
- Finally what is the interpretation of the  $y(\mathbf{x})$  value?
- If  $\mathbf{x}$  is an arbitrary sample vector and  $\mathbf{x}_k$  is its projection on the decision surface then:

$$\mathbf{x} = \mathbf{x}_k + r \frac{\mathbf{w}}{\|\mathbf{w}\|} \Leftrightarrow$$

$$\mathbf{w}^T \mathbf{x} + w_0 = \mathbf{w}^T \mathbf{x}_k + r \frac{\mathbf{w}^T \mathbf{w}}{\|\mathbf{w}\|} + w_0 \Leftrightarrow$$

$$y(\mathbf{x}) = r \frac{\|\mathbf{w}\|^2}{\|\mathbf{w}\|} \Leftrightarrow r = \frac{y(\mathbf{x})}{\|\mathbf{w}\|}$$

- Thus,  $y(\mathbf{x})$  determines the distance of the sample from the decision surface.



# Least Squares

---

- We will see now how can we use the Least Squares approach to estimate  $\mathbf{w}$  and  $w_0$  (training!).



# Least Squares

---

- We will see now how can we use the Least Squares approach to estimate  $\mathbf{w}$  and  $w_0$  (training!).
- Assume we have  $c$  classes and  $N$  training samples. Each one of the samples is associated with a target (label) which is structured using one-hot encoding.
  - For instance, if we have 5 classes and a sample  $x_n$  has  $label = 3$  then the target vector using one-hot encoding is  $\mathbf{t}_n = [0 \ 0 \ 1 \ 0 \ 0]$ . We can gather all such target encoding in one matrix  $\mathbf{T}$ , where every row is the one-hot encoding of the target of the respective sample, i.e.  $\mathbf{T} \in \mathbb{R}^{N \times c}$ .

# Least Squares

---

- We will see now how can we use the Least Squares approach to estimate  $\mathbf{w}$  and  $w_0$  (training!).
- Assume we have  $c$  classes and  $N$  training samples. Each one of the samples is associated with a target (label) which is structured using one-hot encoding.
  - For instance, if we have 5 classes and a sample  $\mathbf{x}_n$  has *label* = 3 then the target vector using one-hot encoding is  $\mathbf{t}_n = [0 \ 0 \ 1 \ 0 \ 0]$ . We can gather all such target encoding in one matrix  $\mathbf{T}$ , where every row is the one-hot encoding of the target of the respective sample, i.e.  $\mathbf{T} \in \mathbb{R}^{N \times c}$ .
- Each of the  $c$  classes has its own linear discriminant:
- For simplicity we will adopt the shorter notation:

$$y_i(\mathbf{x}) = \tilde{\mathbf{w}}_i^t \tilde{\mathbf{x}}, \text{ where } \tilde{\mathbf{w}}_i = \begin{bmatrix} w_{i0} \\ w_{i1} \\ \vdots \\ w_{id} \end{bmatrix} \text{ and } \tilde{\mathbf{x}} = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}$$

# Least Squares

---

$$y_i(\mathbf{x}) = \mathbf{w}_i^t \mathbf{x}$$

**Note:** for the rest of the lecture we will use  $\mathbf{w}, \mathbf{x}$  but we will mean  $\tilde{\mathbf{w}}$  and  $\tilde{\mathbf{x}}$ !

# Least Squares

---

$$y_i(\mathbf{x}) = \mathbf{w}_i^t \mathbf{x}$$

**Note:** for the rest of the lecture we will use  $\mathbf{w}, \mathbf{x}$  but we will mean  $\tilde{\mathbf{w}}$  and  $\tilde{\mathbf{x}}$ !

- We can use matrix notation to include all the discriminants:

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

where the columns of  $\mathbf{W}$  contain the vectors  $\mathbf{w}_i$  for each class, i.e.  $\mathbf{W} \in \mathbb{R}^{(d+1) \times c}$ , and

$$\mathbf{y}(\mathbf{x}) = \begin{bmatrix} y_1(\mathbf{x}) \\ y_2(\mathbf{x}) \\ \vdots \\ y_c(\mathbf{x}) \end{bmatrix}$$

- Thus the classification problem is formed as:

assign  $\mathbf{x}$  to class  $\omega_i$  if  $i = \underset{j}{\operatorname{argmax}} y_j(\mathbf{x})$

# Least Squares

---

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

- Now that we have formulated our problem we only need to determine matrix  $\mathbf{W}$ .

# Least Squares

---

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

- Now that we have formulated our problem we only need to determine matrix  $\mathbf{W}$ .
- Our objective is to make, for each sample  $\mathbf{x}_n$ , as minimum as possible the difference:

$$\mathbf{W}^T \mathbf{x}_n - \mathbf{t}_n^T$$

# Least Squares

---

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

- Now that we have formulated our problem we only need to determine matrix  $\mathbf{W}$ .
- Our objective is to make, for each sample  $\mathbf{x}_n$ , as minimum as possible the difference:

$$\mathbf{W}^T \mathbf{x}_n - \mathbf{t}_n^T$$

if we use the matrices  $\mathbf{T}$  and  $\mathbf{X}$ , where every row of  $\mathbf{X}$  is a sample of the training set, i.e.  $\mathbf{X} \in \mathbb{R}^{N \times (d+1)}$ , then our optimization problem can be described as the minimization of the sum-of-squares error function:

$$E(\mathbf{W}) = \frac{1}{2} \text{tr}[(\mathbf{XW} - \mathbf{T})^T (\mathbf{XW} - \mathbf{T})]$$

# Least Squares

---

- In order to minimize the sum-of-squares error function (least squares approach):

$$E(\mathbf{W}) = \frac{1}{2} \text{tr}[(\mathbf{XW} - \mathbf{T})^T (\mathbf{XW} - \mathbf{T})]$$

We will find the respective derivative and set equal to zero:

$$\frac{\partial E(\mathbf{W})}{\partial \mathbf{W}} = 0$$

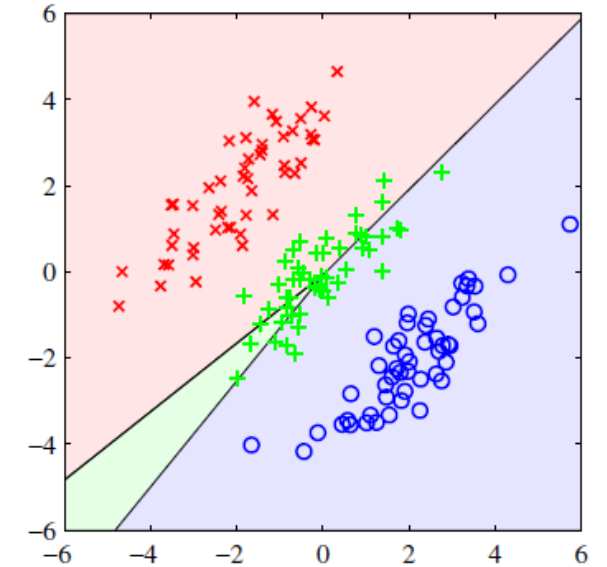
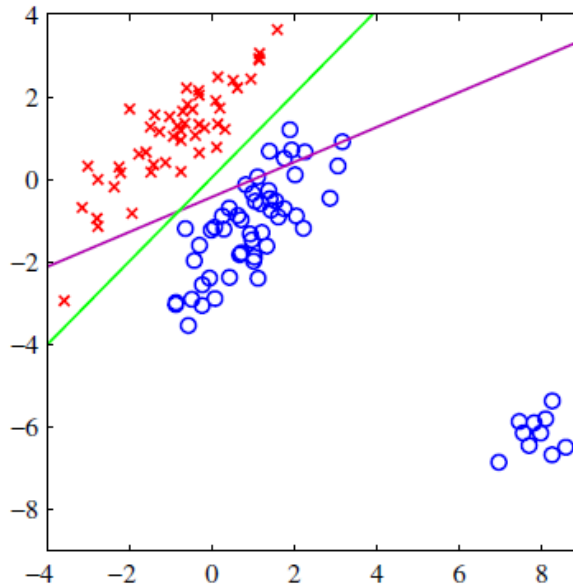
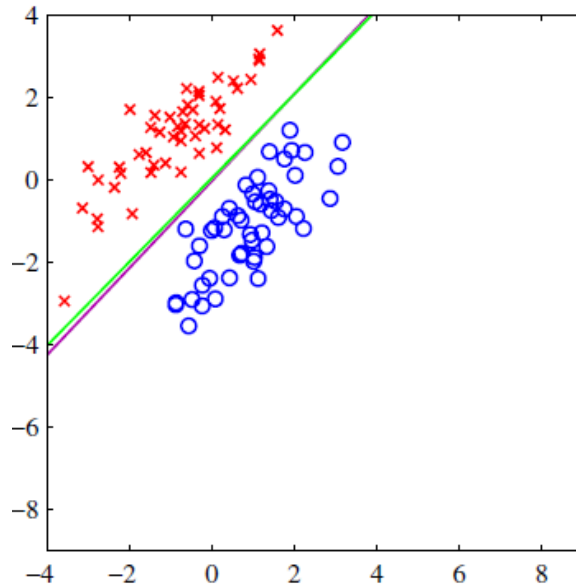
we will come up with the solution:

$$\mathbf{W}_{LS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{T} = \mathbf{X}^\dagger \mathbf{T}$$

- Classification using the discriminant function:  $\mathbf{y}_{LS}(\mathbf{x}) = \mathbf{W}_{LS} \mathbf{x}$



# Least Squares – Problems!



Bishop

- The decision boundaries are very sensitive to outliers.
- for  $c > 2$  there are masking phenomena (decision regions become very small or completely ignored).
- $y_{LS}(x)$  are not real probabilities and sometimes can have values outside of the range  $[0,1]$ .
- Very costly to compute!

# Perceptron

---

- Perceptron algorithm is **only for two-class** problems.

# Perceptron

---

- Perceptron algorithm is only for two-class problems.
- Assuming input  $\mathbf{x} \in \mathbb{R}^{d+1}$  (again we consider the augmented samples  $[1, \mathbf{x}]$ )

# Perceptron

---

- Perceptron algorithm is only for two-class problems.
- Assuming input  $\mathbf{x} \in \mathbb{R}^{d+1}$  (again we consider the augmented samples  $[1, \mathbf{x}]$ )
- We also assume that we have targets  $t \in \{\omega_1, \omega_2\} \rightarrow t \in \{-1, 1\}$

# Perceptron

---

- Perceptron algorithm is only for two-class problems.
- Assuming input  $\mathbf{x} \in \mathbb{R}^{d+1}$  (again we consider the augmented samples  $[1, \mathbf{x}]$ )
- We also assume that we have targets  $t \in \{\omega_1, \omega_2\} \rightarrow t \in \{1, -1\}$
- Thus, we can make a prediction based on outcome:  $y(\mathbf{x}) = f(\mathbf{w}^T \mathbf{x})$  where
$$f(a) = \begin{cases} 1, & a > 0 \\ -1, & a \leq 0 \end{cases}$$
- With this adoption I decide to assign to class  $\omega_1$  if  $\mathbf{w}^T \mathbf{x} > 0$  and to class  $\omega_2$  otherwise.

$$t \in \{\omega_1, \omega_2\} \rightarrow t \in \{-1, 1\}$$

$$\omega_1 \text{ if } \mathbf{w}^T \mathbf{x} \geq 0$$

# Perceptron - Training

---

- Assume we have  $N$  training samples,  $x_n$ ,  $n = 1, \dots, N$  with the respective targets (labels)  $t_n$ .

- Objective: find  $\mathbf{w}$  such that for all  $(x_n, t_n)$  pairs I get:

$$\mathbf{w}^T x_n t_n > 0$$

- What is the objective function that we should minimize?

$$t \in \{\omega_1, \omega_2\} \rightarrow t \in \{-1, 1\}$$

$$\omega_1 \text{ if } \mathbf{w}^T \mathbf{x} \geq 0$$

# Perceptron - Training

---

- Assume we have  $N$  training samples,  $\mathbf{x}_n, n = 1, \dots, N$  with the respective targets (labels)  $t_n$ .
- Objective: find  $\mathbf{w}$  such that for all  $(\mathbf{x}_n, t_n)$  pairs I get:

$$\mathbf{w}^T \mathbf{x}_n t_n > 0$$

- What is the objective function that we should minimize?

$$E_P(\mathbf{w}) = - \sum_{n \in \mathcal{M}} \mathbf{w}^T \mathbf{x}_n t_n$$

where  $\mathcal{M} = \{n: \mathbf{w}^T \mathbf{x}_n t_n \leq 0\}$ .

# Perceptron - Training

---

$$E_P(\mathbf{w}) = -\sum_{n \in \mathcal{M}} \mathbf{w}^t \mathbf{x}_n t_n \text{ where } \mathcal{M} = \{n: \mathbf{w}^t \mathbf{x}_n t_n \leq 0\}$$

- $E_P(\mathbf{w})$  is piecewise linear. Why?



# Perceptron - Training

---

$$E_P(\mathbf{w}) = -\sum_{n \in \mathcal{M}} \mathbf{w}^t \mathbf{x}_n t_n \text{ where } \mathcal{M} = \{n: \mathbf{w}^t \mathbf{x}_n t_n \leq 0\}$$

- $E_P(\mathbf{w})$  is piecewise linear. Why?
  - If we change  $\mathbf{w}$  linearly then  $E_P(\mathbf{w})$  changes also linearly until the point where there is a change in the number of misclassified samples!
  - At these points the partial derivative of  $E_P(\mathbf{w})$  cannot be defined and, thus, the derivative of  $E_P(\mathbf{w})$  is a discontinuous function!
- Thus, how shall we proceed with the training of the perceptron algorithm?

# Perceptron - Training

---

$$E_P(\mathbf{w}) = -\sum_{n \in \mathcal{M}} \mathbf{w}^t \mathbf{x}_n t_n \text{ where } \mathcal{M} = \{n: \mathbf{w}^t \mathbf{x}_n t_n \leq 0\}$$

- $E_P(\mathbf{w})$  is piecewise linear. Why?
  - If we change  $\mathbf{w}$  linearly then  $E_P(\mathbf{w})$  changes also linearly until the point where there is a change in the number of misclassified samples!
  - At these points the partial derivative of  $E_P(\mathbf{w})$  cannot be defined and, thus, the derivative of  $E_P(\mathbf{w})$  is a discontinuous function!
- Thus, how shall we proceed with the training of the perceptron algorithm?
- Gradient Descent!

# Gradient Descent

---

- The gradient of a function encodes all directional derivatives.
  - It points towards the direction of the steepest ascent.
  - the magnitude of the gradient is the rate of increase in that direction.
  - It is always perpendicular the contours of a function.

# Gradient Descent

---

- The gradient of a function encodes all directional derivatives.
  - It points towards the direction of the steepest ascent.
  - the magnitude of the gradient is the rate of increase in that direction.
  - It is always perpendicular the contours of a function.
- Whenever the loss is a sum of error terms for each datapoint we can use Gradient Descent, e.g.:

$$E(\mathbf{w}) = - \sum_{n \in \mathcal{M}} \mathbf{w}^t \mathbf{x}_n t_n = - \sum_{n \in \mathcal{M}} E_n(\mathbf{w})$$

- In general we have:  $E(\mathbf{w}) = \sum_{n=1}^N E_n(\mathbf{w})$

# Stochastic gradient descent

---

- **Gradient descent:**  $\mathbf{w}^{i+1} = \mathbf{w}^i - \eta \nabla E(\mathbf{w}^i)$ 
  - where  $E(\mathbf{w}^i) = \sum_{n=1}^N E_n(\mathbf{w}^i)$
- We can extend the above concept by adding stochasticity in our Gradient Descent algorithm (Stochastic Gradient Descent):

$$\mathbf{w}^{i+1} = \mathbf{w}^i - \eta \nabla E_n(\mathbf{w}^i)$$

or

$$\mathbf{w}^{i+1} = \mathbf{w}^i - \eta \nabla \sum_{l=1}^L E_l(\mathbf{w}^i) \text{ (minibatch)}$$

- Either with single points or minibatches what I compute is an approximation of the full gradient.

# Stochastic Gradient Descent - Algorithm

---

- SGD steps:
  - Initialize  $w^0$ , choose  $\eta$
  - Iterate over datapoints:  $\mathbf{w}^{i+1} = \mathbf{w}^i - \eta \nabla \sum_{l=1}^L E_l(\mathbf{w}^i)$

# Stochastic Gradient Descent - Algorithm

---

- SGD steps:
  - Initialize  $w^0$ , choose  $\eta$
  - Iterate over datapoints:  $\mathbf{w}^{i+1} = \mathbf{w}^i - \eta \nabla \sum_{l=1}^L E_l(\mathbf{w}^i)$
- If  $E$  is convex in  $\mathbf{w}$  and  $\eta$  is small enough SGD will converge
  - very small  $\eta$  slow convergence
  - large  $\eta$ : probably it will not converge
- Scheduling for  $\eta$  is important (maybe decrease  $\eta$  with iterations)

# Stochastic Gradient Descent - Algorithm

---

- SGD steps:
  - Initialize  $w^0$ , choose  $\eta$
  - Iterate over datapoints:  $\mathbf{w}^{i+1} = \mathbf{w}^i - \eta \nabla \sum_{l=1}^L E_l(\mathbf{w}^i)$
- If  $E$  is convex in  $\mathbf{w}$  and  $\eta$  is small enough SGD will converge
  - very small  $\eta$  slow convergence
  - large  $\eta$ : probably it will not converge
- Scheduling for  $\eta$  is important (maybe decrease  $\eta$  with iterations)
- Single point gradients are not necessarily equal to zero as the full gradient may be. Thus this also helps to escape local minima.



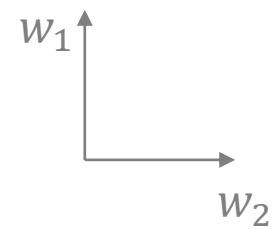
# Stochastic Gradient Descent - Algorithm

---

- SGD steps:
  - Initialize  $w^0$ , choose  $\eta$
  - Iterate over datapoints:  $\mathbf{w}^{i+1} = \mathbf{w}^i - \eta \nabla \sum_{l=1}^L E_l(\mathbf{w}^i)$
- If  $E$  is convex in  $\mathbf{w}$  and  $\eta$  is small enough SGD will converge
  - very small  $\eta$  slow convergence
  - large  $\eta$ : probably it will not converge
- Scheduling for  $\eta$  is important (maybe decrease  $\eta$  with iterations)
- Single point gradients are not necessarily equal to zero as the full gradient may be. Thus this also helps to escape local minima.
- Initialization plays crucial role!
  - maybe perform multiple initializations

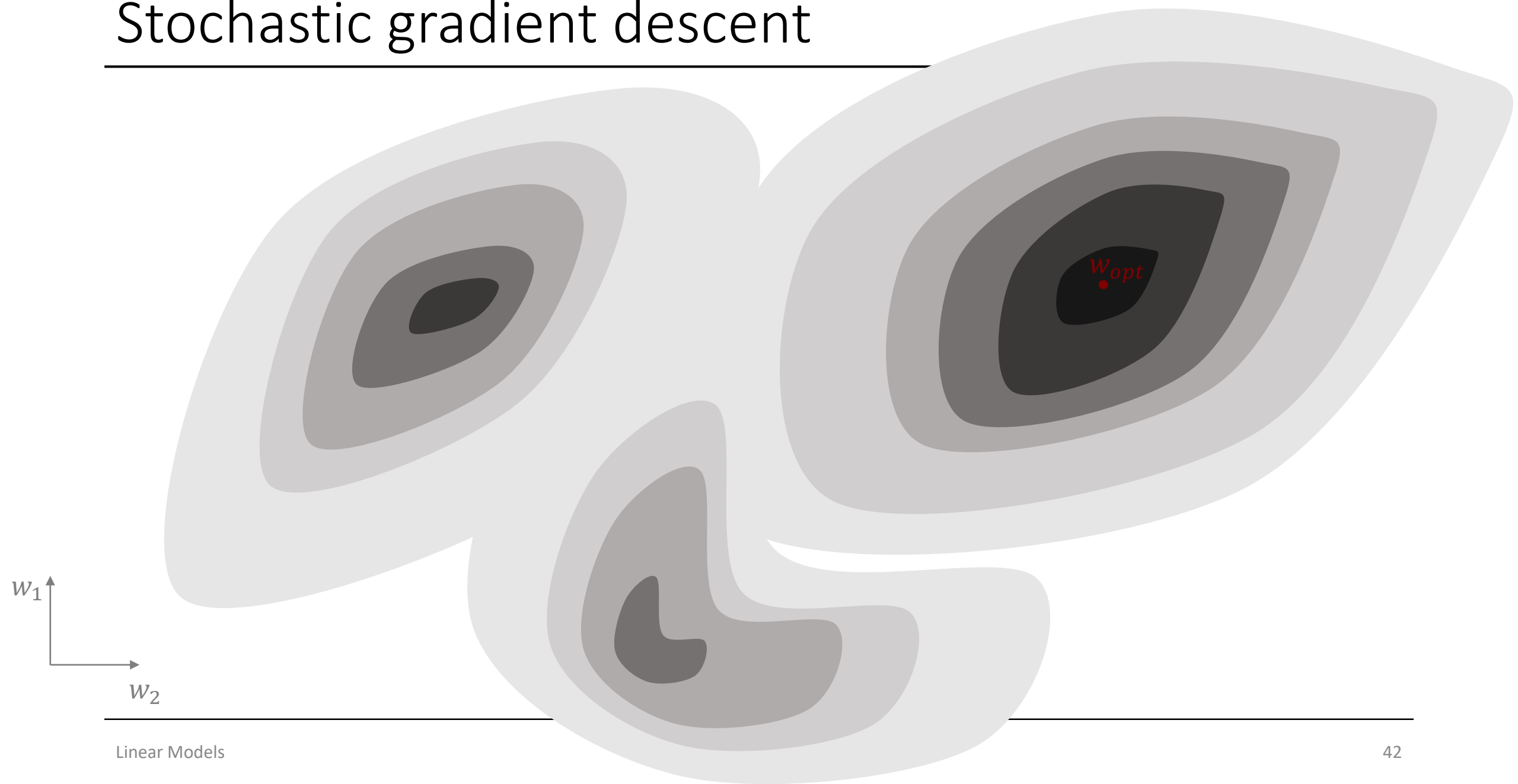
# Stochastic gradient descent

---



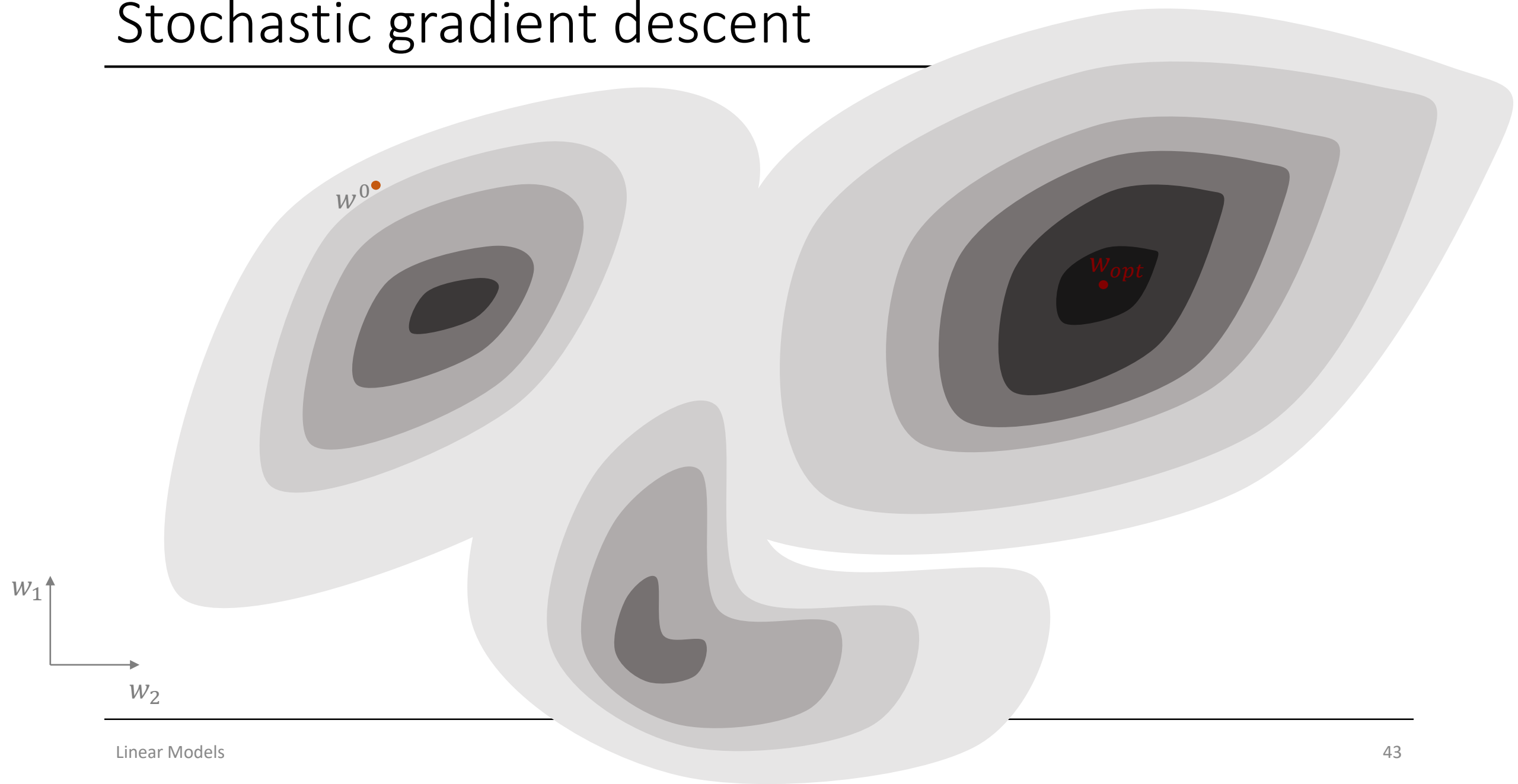
# Stochastic gradient descent

---



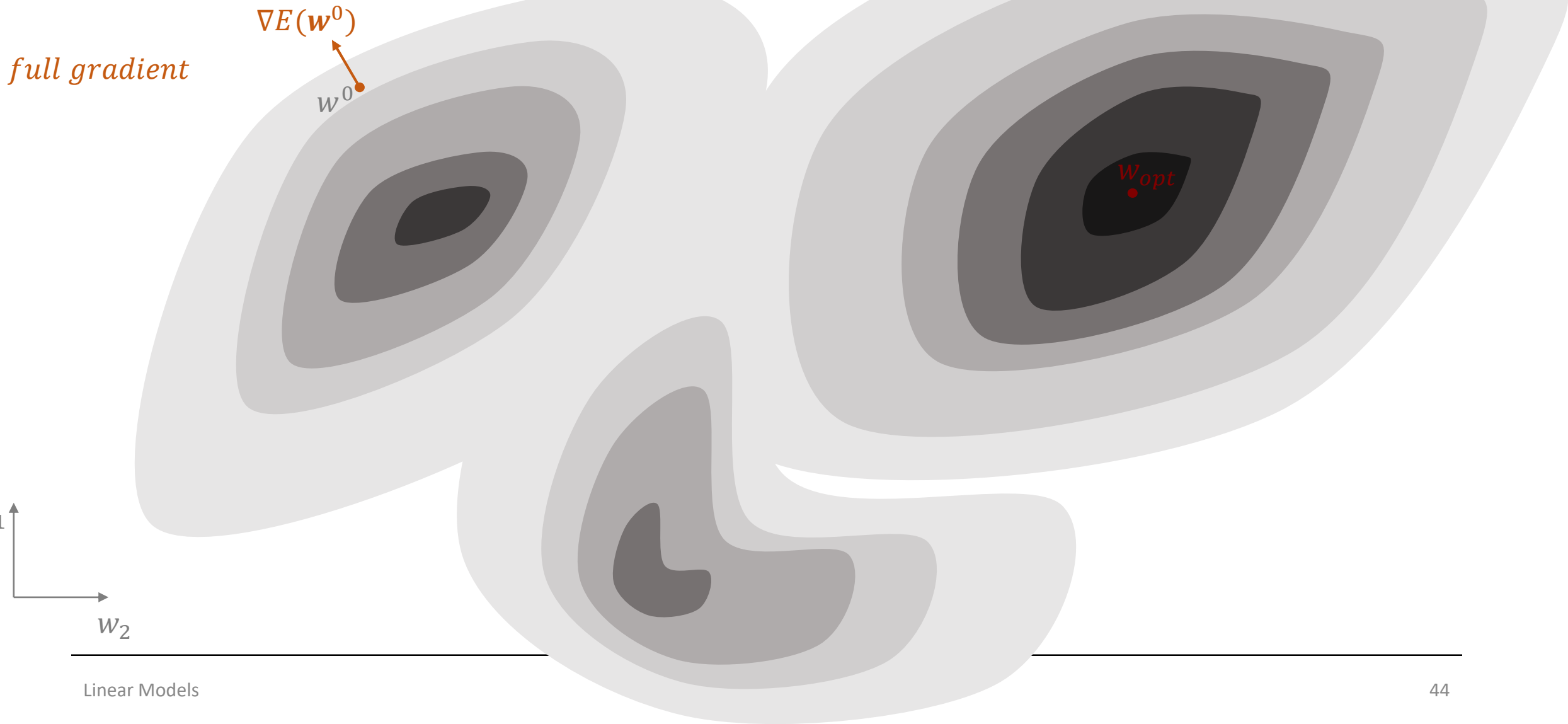
# Stochastic gradient descent

---



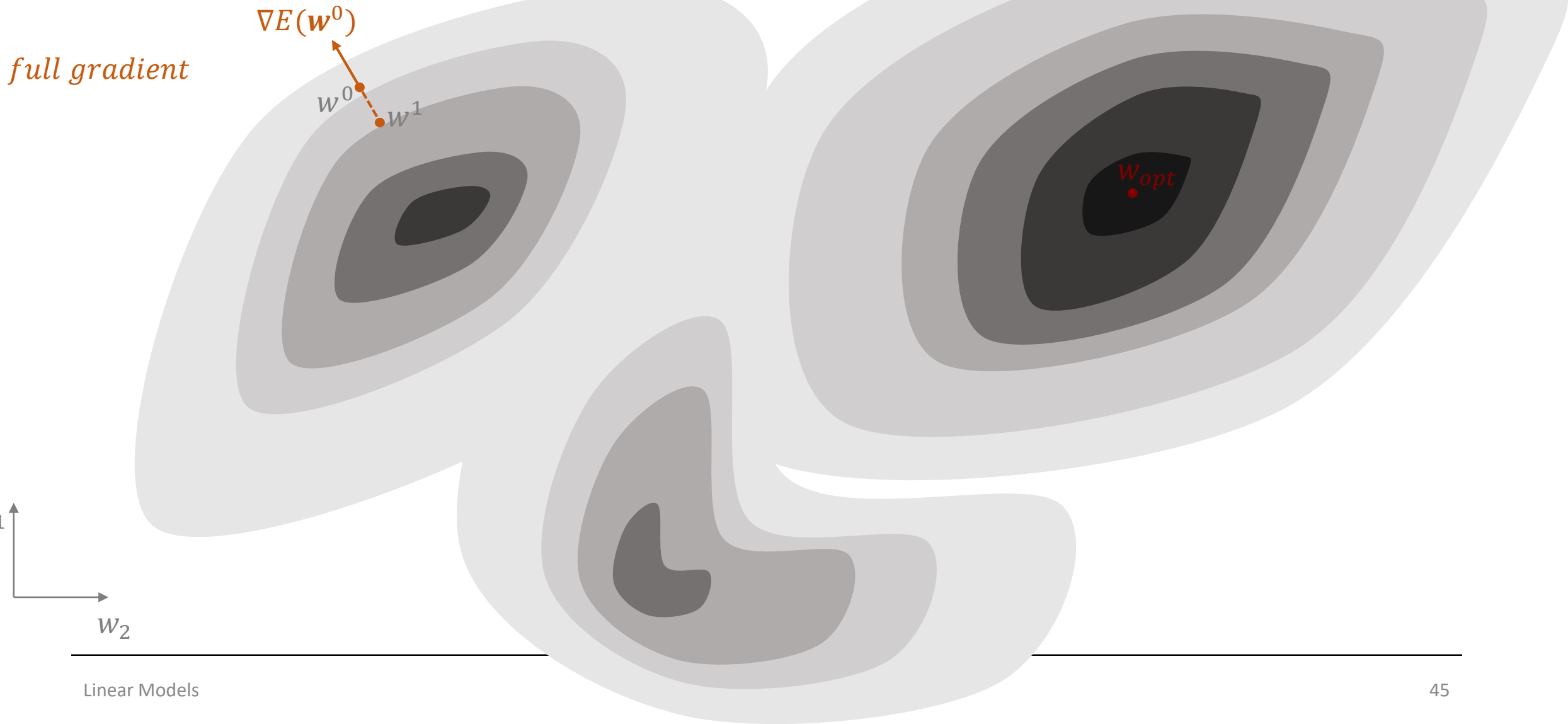
# Stochastic gradient descent

---



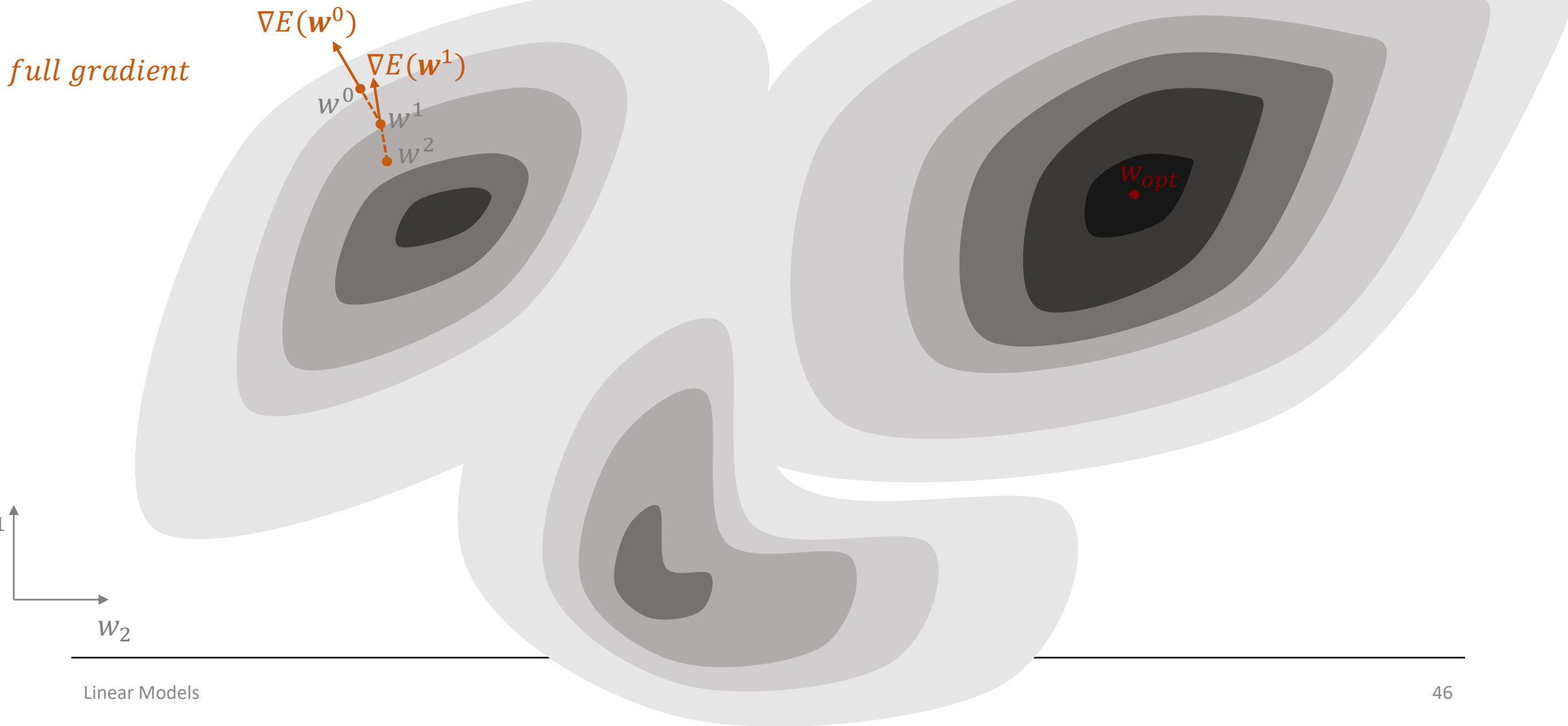
# Stochastic gradient descent

---

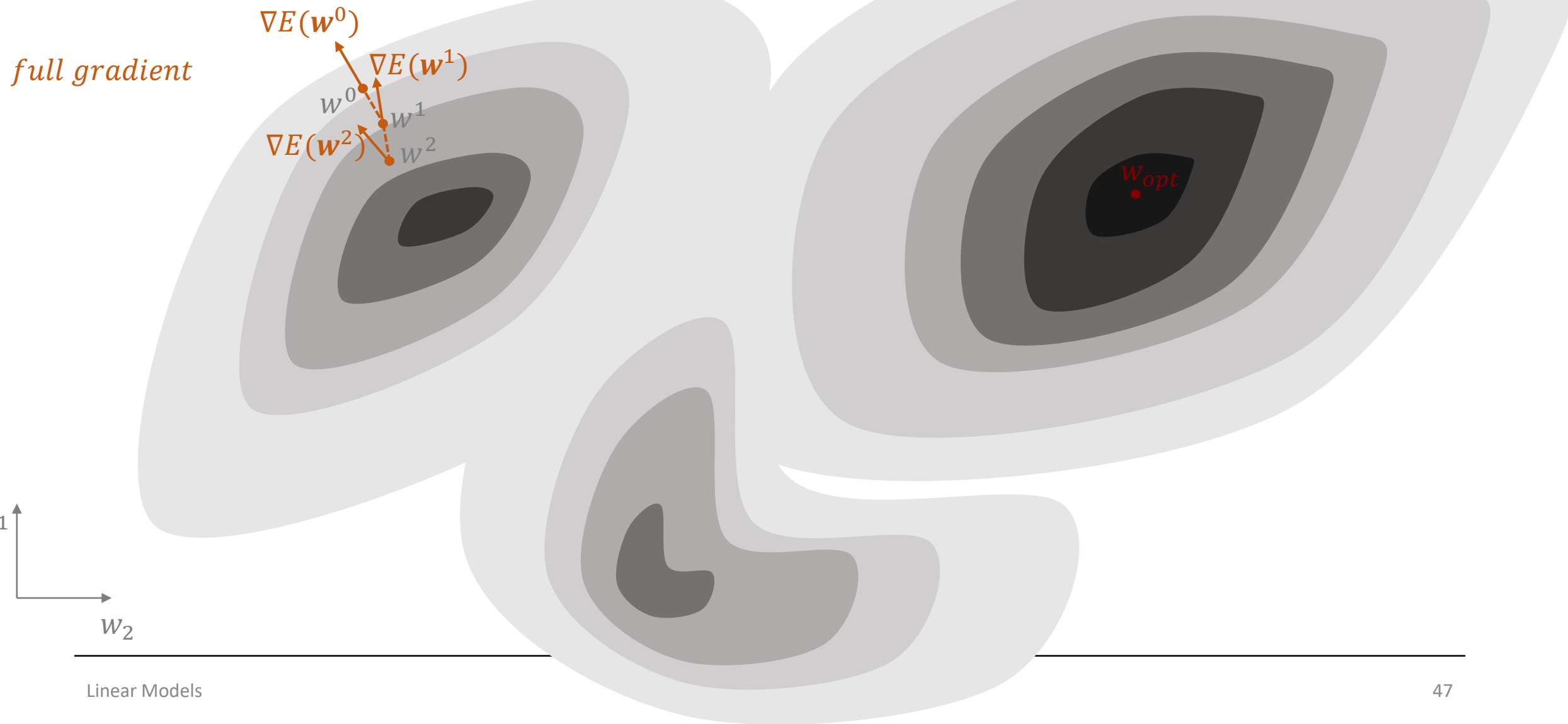


# Stochastic gradient descent

---

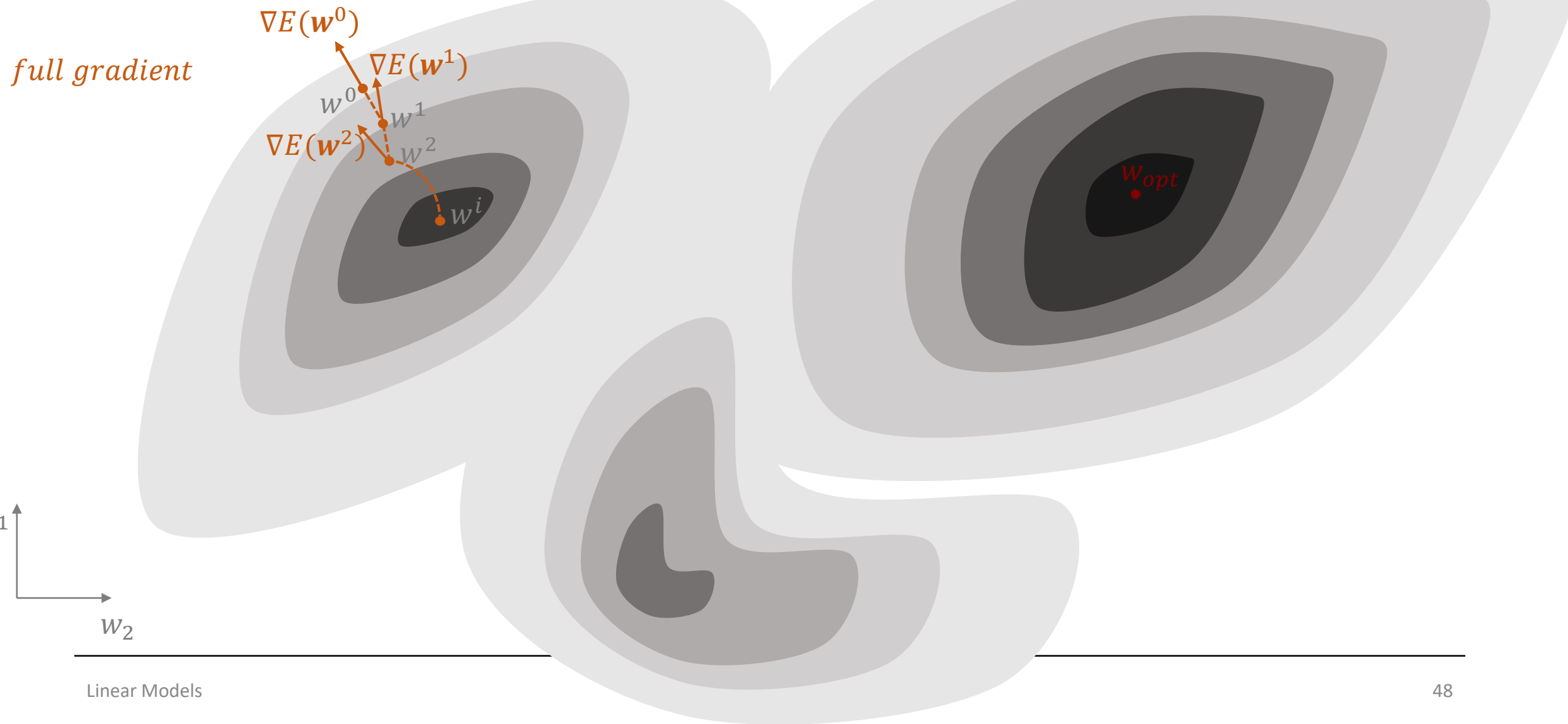


# Stochastic gradient descent





# Stochastic gradient descent

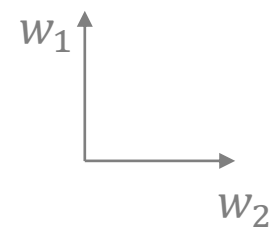


# Stochastic gradient descent

---

*full gradient*

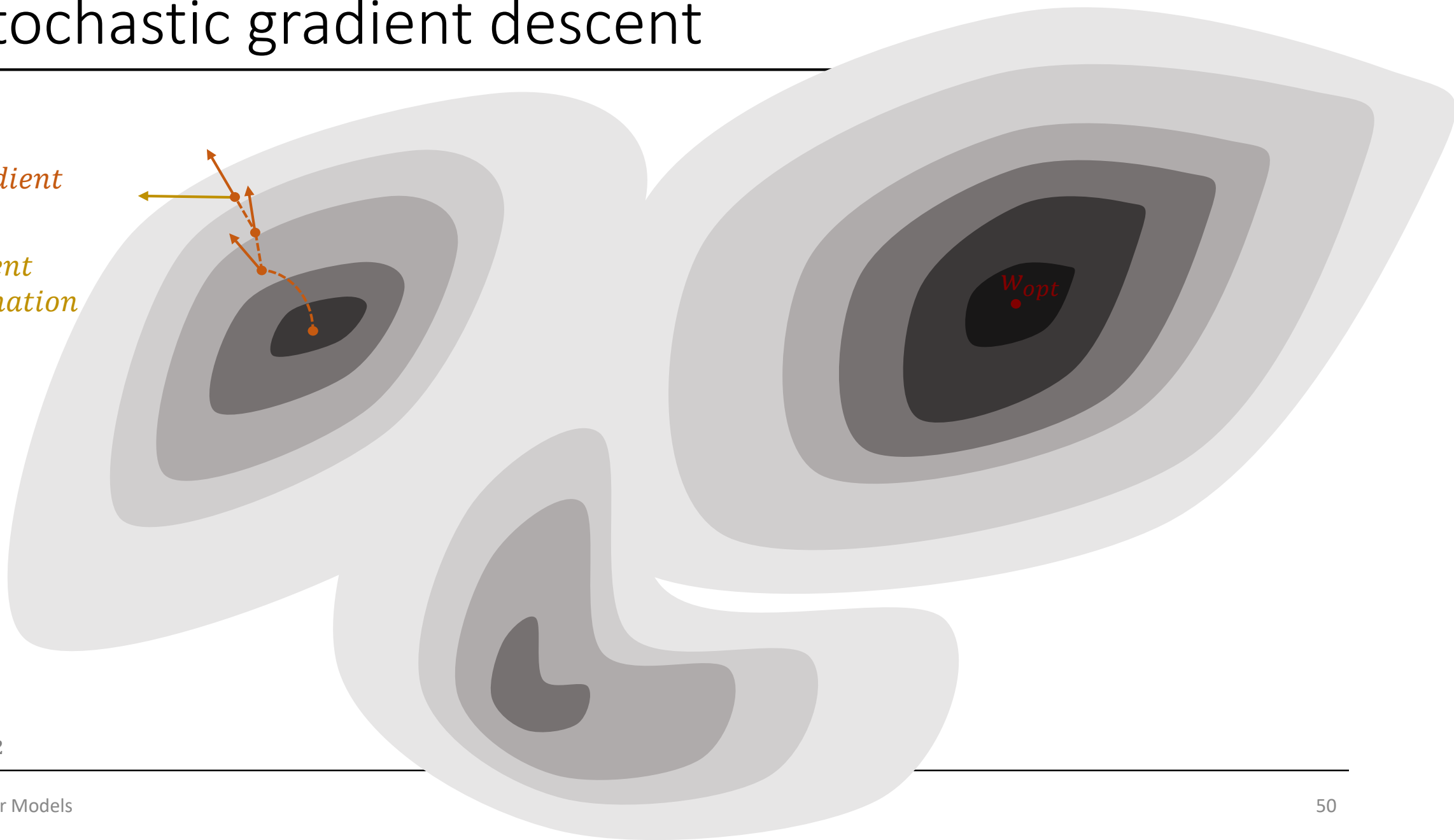
*gradient approximation*



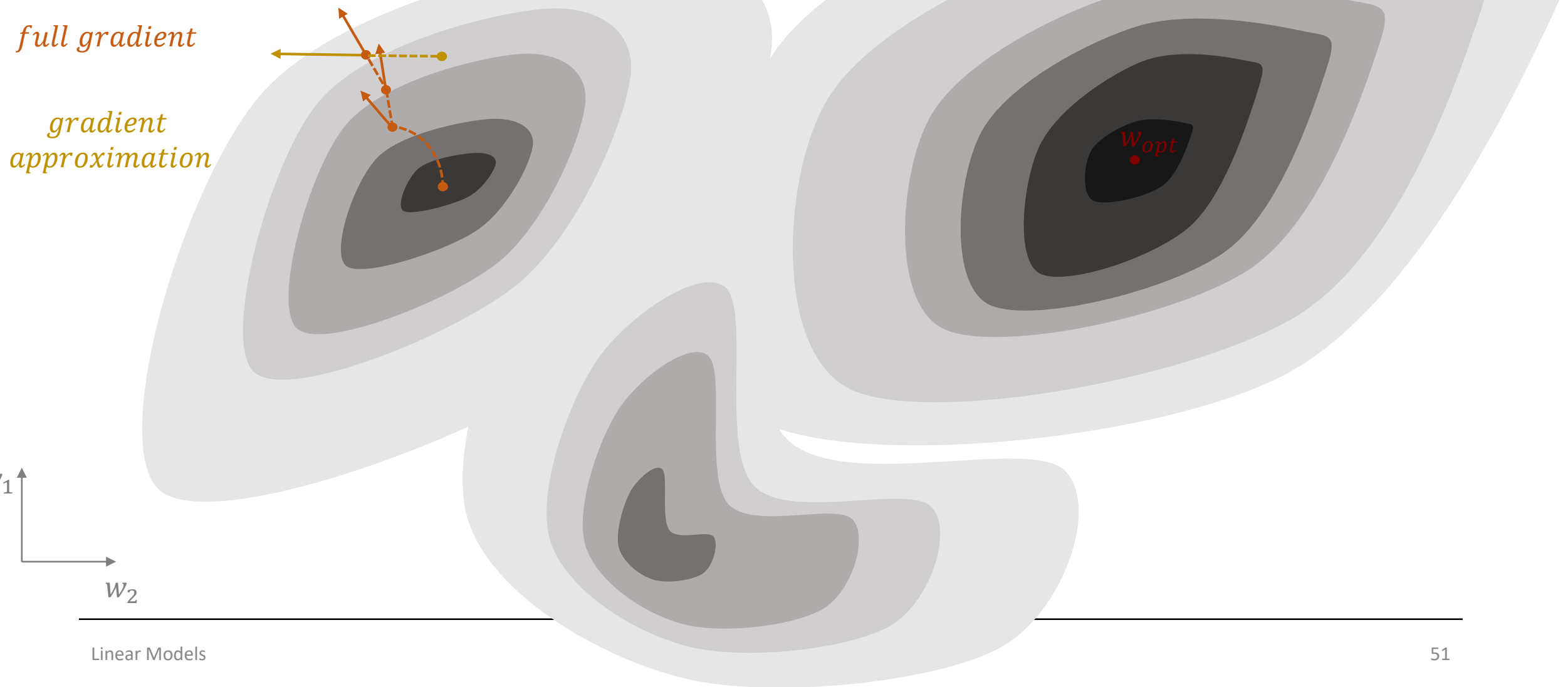
# Stochastic gradient descent

*full gradient*  
*gradient approximation*

$w_1$   
 $w_2$



# Stochastic gradient descent



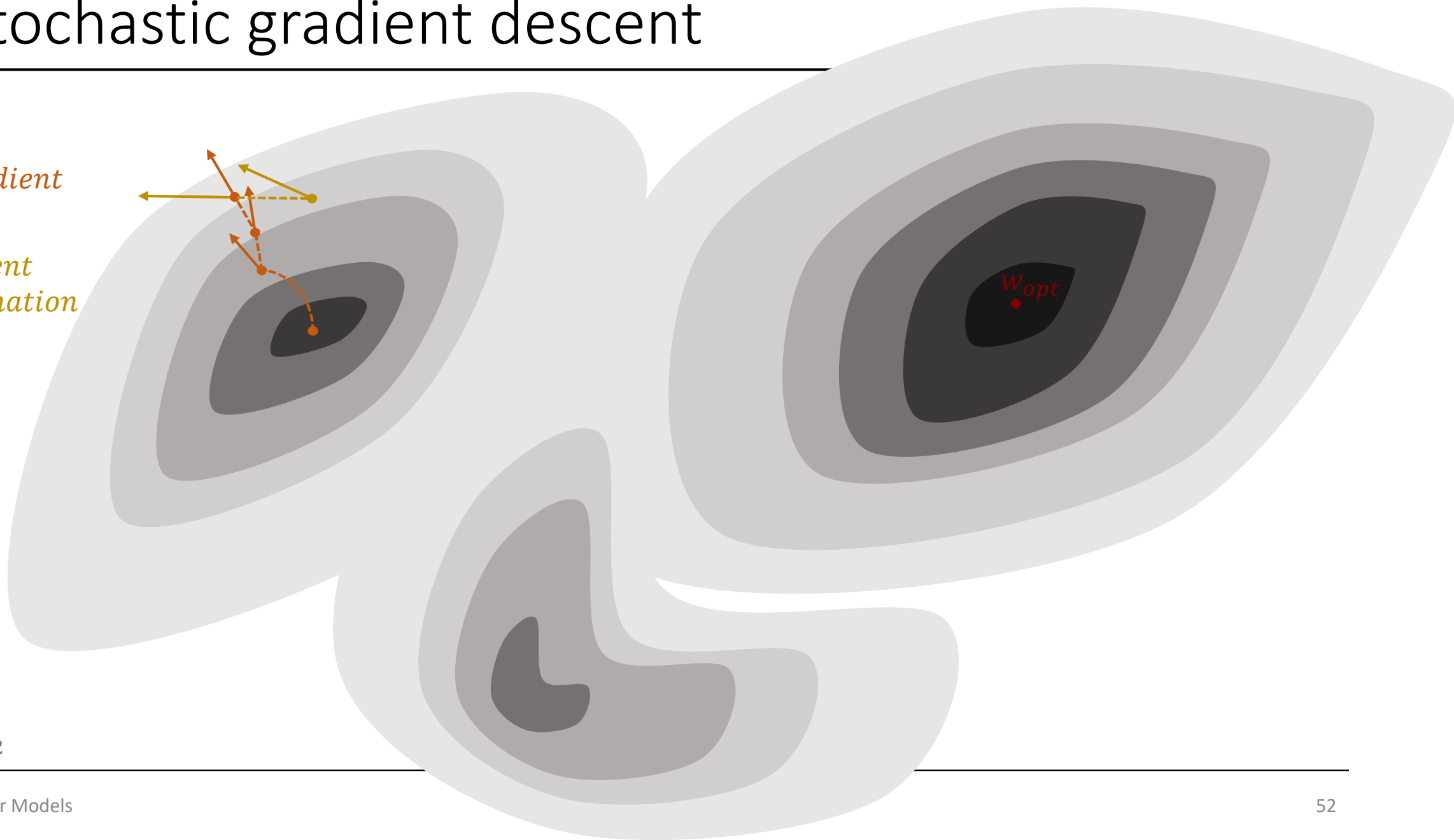
# Stochastic gradient descent

*full gradient*

*gradient approximation*

$w_1$

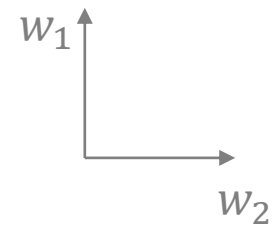
$w_2$



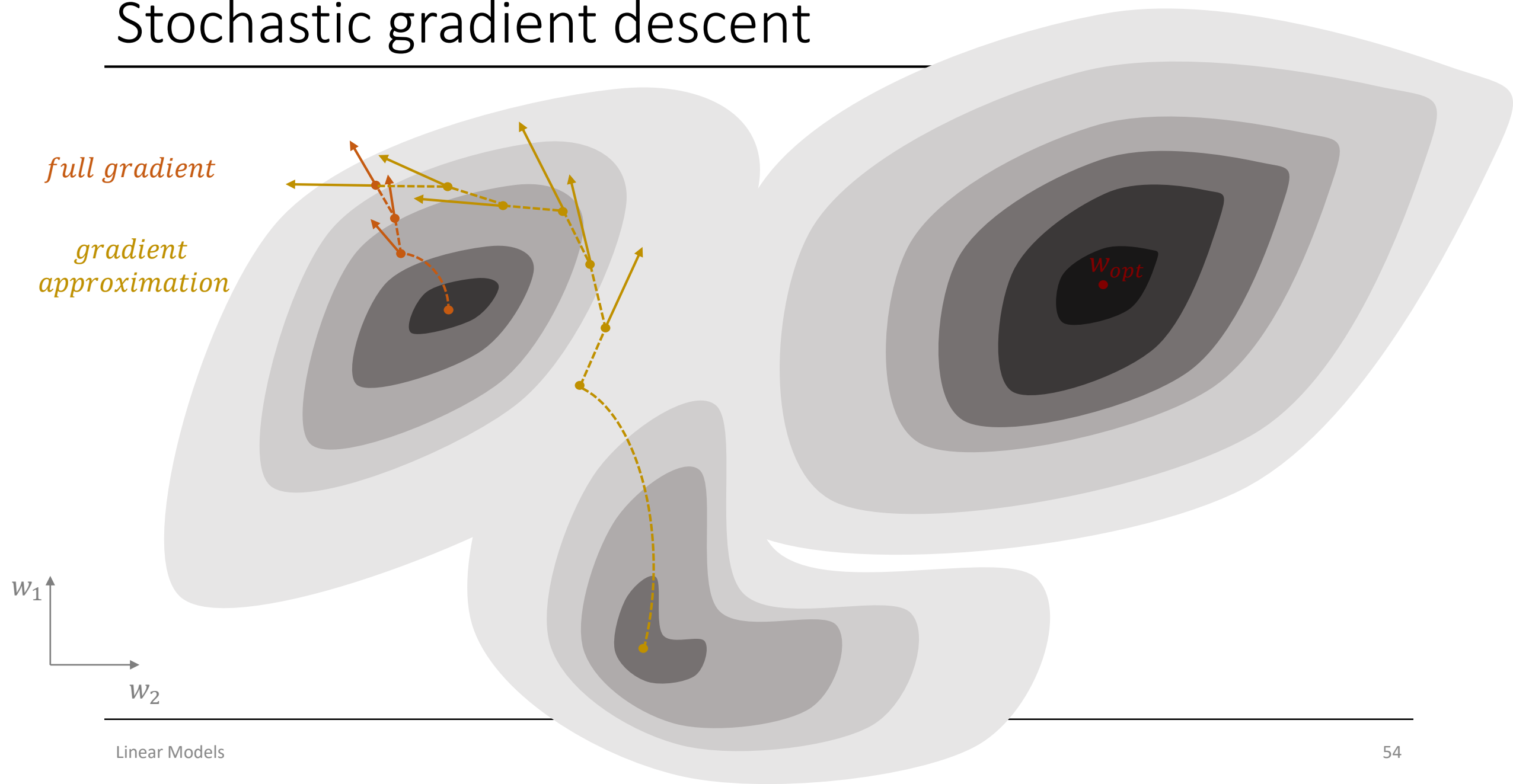
# Stochastic gradient descent

*full gradient*

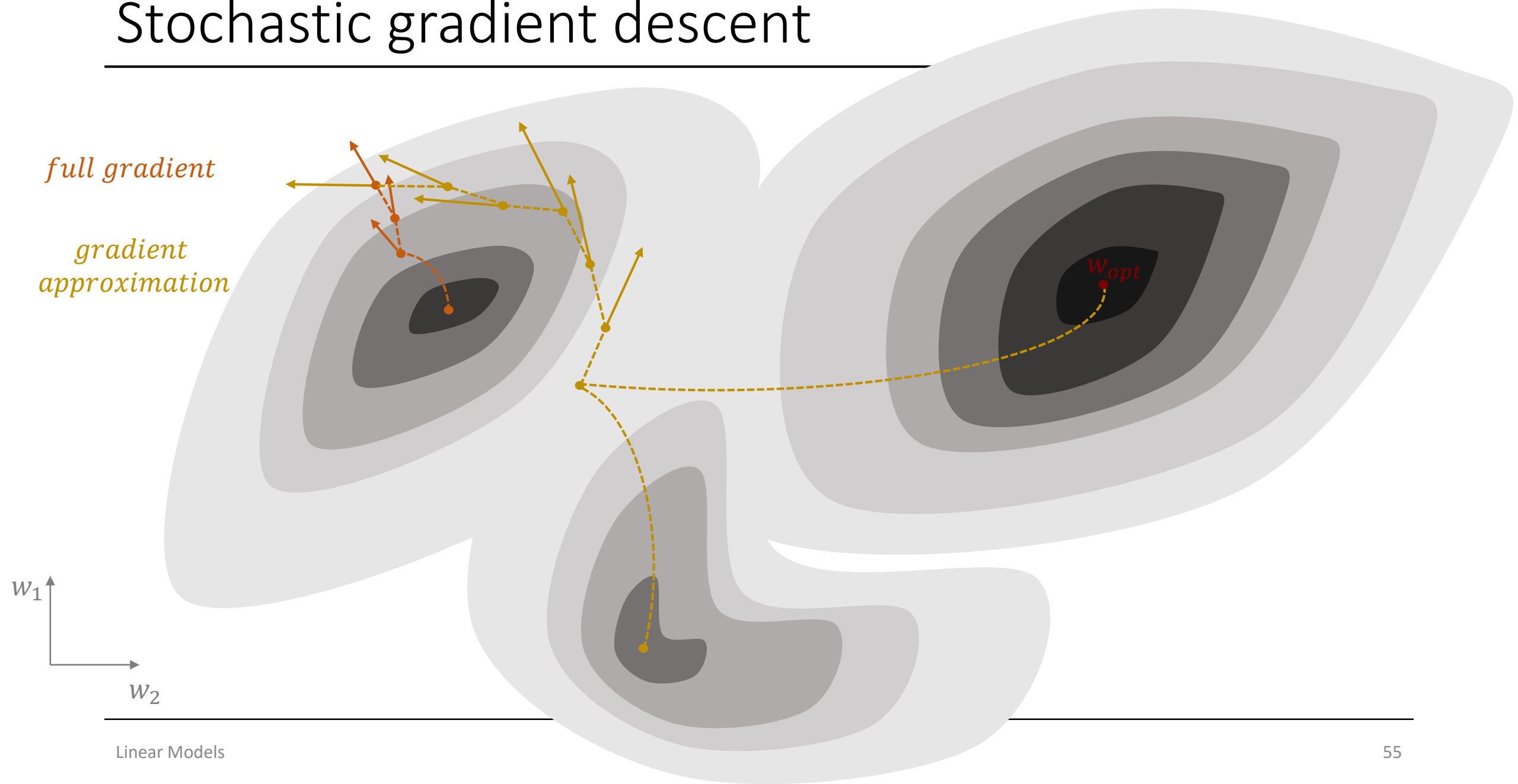
*gradient approximation*



# Stochastic gradient descent



# Stochastic gradient descent





# Perceptron – Training : SGD

- Perceptron criterion:

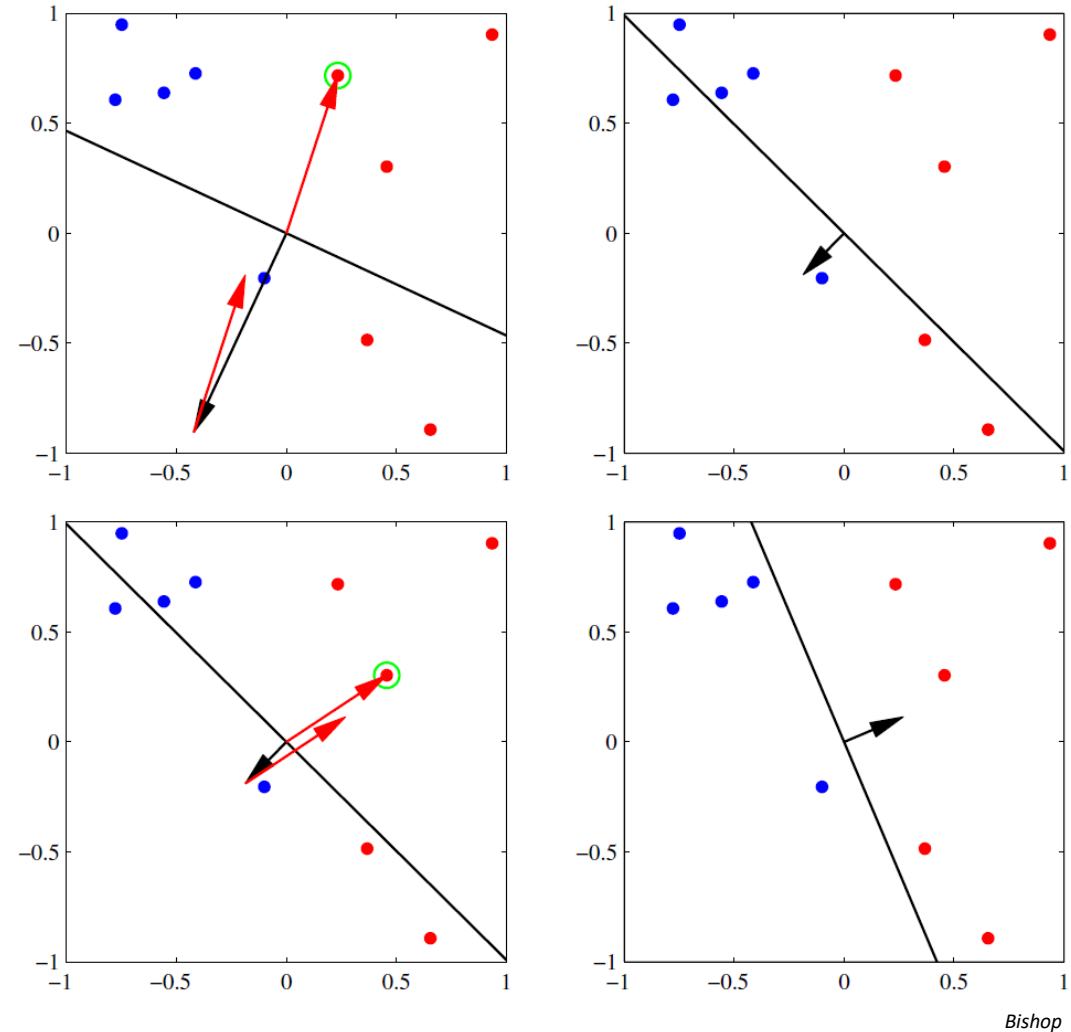
$$E(\mathbf{w}) = - \sum_{n \in \mathcal{M}} \mathbf{w}^t \mathbf{x}_n t_n = - \sum_{n \in \mathcal{M}} E_n(\mathbf{w})$$

- Sequential learning:

- Single sample SGD:

$$\mathbf{w}^{i+1} = \mathbf{w}^i - \eta \nabla E_n(\mathbf{w}^i)$$

$$\mathbf{w}^{i+1} = \mathbf{w}^i + \eta \mathbf{x}_n t_n$$



Bishop

# Least Mean Squares (LMS) algorithm

---

- If we replace the perceptron criterion with a least squares one we get:

$$E_{LS}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}^t \mathbf{x})^2$$

- If we estimate the corresponding derivative:

$$\frac{\partial E_{LS}(\mathbf{w})}{\partial \mathbf{w}} = \sum_{n=1}^N \frac{\partial \frac{1}{2} (t_n - \mathbf{w}^t \mathbf{x})^2}{\partial \mathbf{w}} = \sum_{n=1}^N -(t_n - \mathbf{w}^t \mathbf{x}) \mathbf{x}^t$$

- We get the respective Sequential Learning algorithm (SGD):

$$\mathbf{w}^{i+1} = \mathbf{w}^i + \eta \left( t_n - (\mathbf{w}^i)^t \mathbf{x} \right) \mathbf{x}^t$$

# LMS algorithm - Notes

---

$$\mathbf{w}^{i+1} = \mathbf{w}^i + \eta \left( t_n - (\mathbf{w}^i)^t \mathbf{x} \right) \mathbf{x}^t$$

- LMS is a two-class algorithm as perceptron is.

# LMS algorithm - Notes

---

$$\mathbf{w}^{i+1} = \mathbf{w}^i + \eta \left( t_n - (\mathbf{w}^i)^t \mathbf{x} \right) \mathbf{x}^t$$

- LMS is a two-class algorithm as perceptron is.
- Remember that for the least squares solution that we previously investigated (multiclass problem) we got:

$$\mathbf{W}_{LS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{T} = \mathbf{X}^\dagger \mathbf{T}$$

# LMS algorithm - Notes

---

$$\mathbf{w}^{i+1} = \mathbf{w}^i + \eta \left( t_n - (\mathbf{w}^i)^t \mathbf{x} \right) \mathbf{x}^t$$

- LMS is a two-class algorithm as perceptron is.
- Remember that for the least squares solution that we previously investigated (multiclass problem) we got:

$$\mathbf{W}_{LS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{T} = \mathbf{X}^\dagger \mathbf{T}$$

- Matrix inverse is a computationally difficult to get. ( $\mathcal{O}(N^3)$ ).
- With LMS we can overcome this problem and process data in minibatches or on a single sample base (not all data are necessary to be available at once!)

# LMS algorithm - Notes

---

$$\mathbf{w}^{i+1} = \mathbf{w}^i + \eta \left( t_n - (\mathbf{w}^i)^t \mathbf{x} \right) \mathbf{x}^t$$

- LMS is a two-class algorithm as perceptron is.
- Remember that for the least squares solution that we previously investigated (multiclass problem) we got:

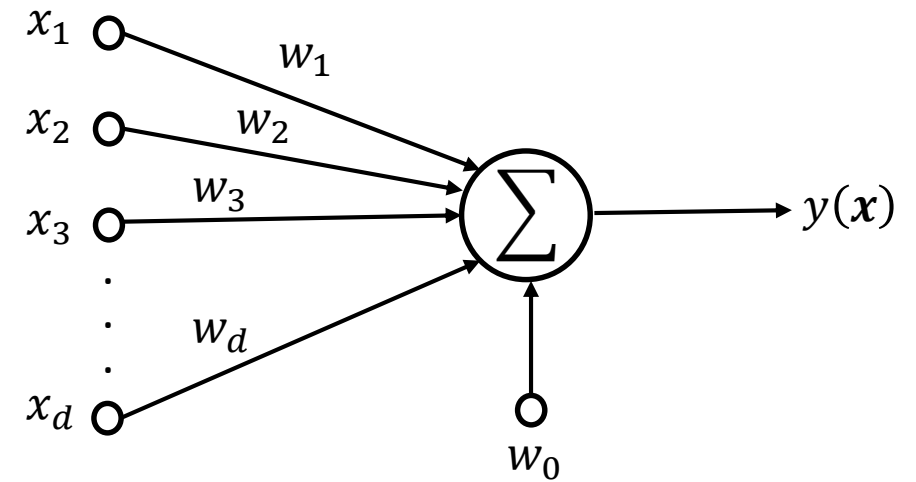
$$\mathbf{W}_{LS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{T} = \mathbf{X}^\dagger \mathbf{T}$$

- Matrix inverse is a computationally difficult to get. ( $\mathcal{O}(N^3)$ ).
- With LMS we can overcome this problem and process data in minibatches or on a single sample base (not all data are necessary to be available at once!)
- LMS algorithm is also called Widrow-Hoff algorithm.

# Neuron

---

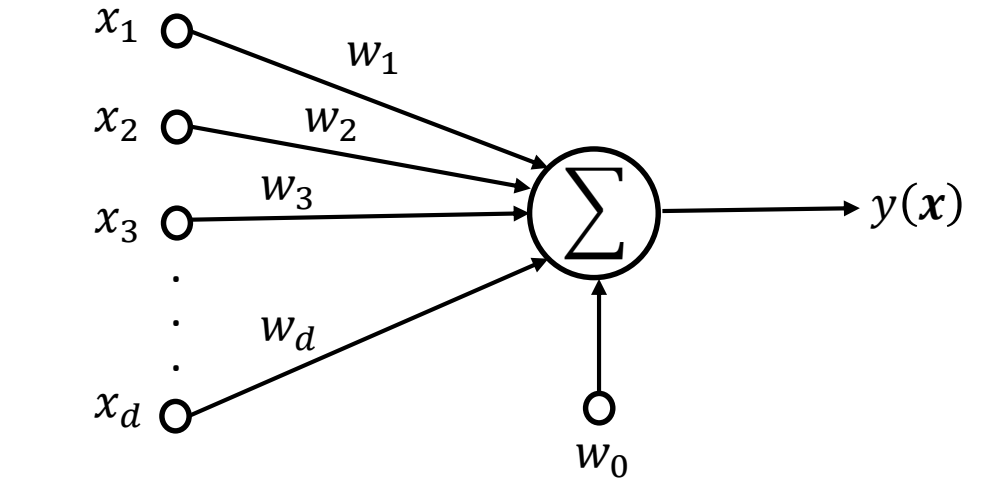
$$y(\mathbf{x}) = \mathbf{w}^t \mathbf{x} + w_0$$



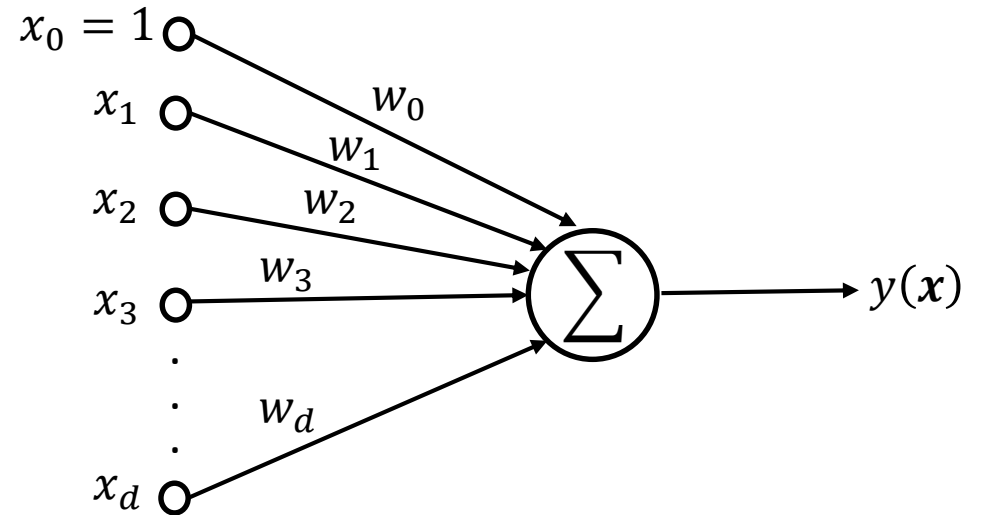
# Neuron

---

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



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



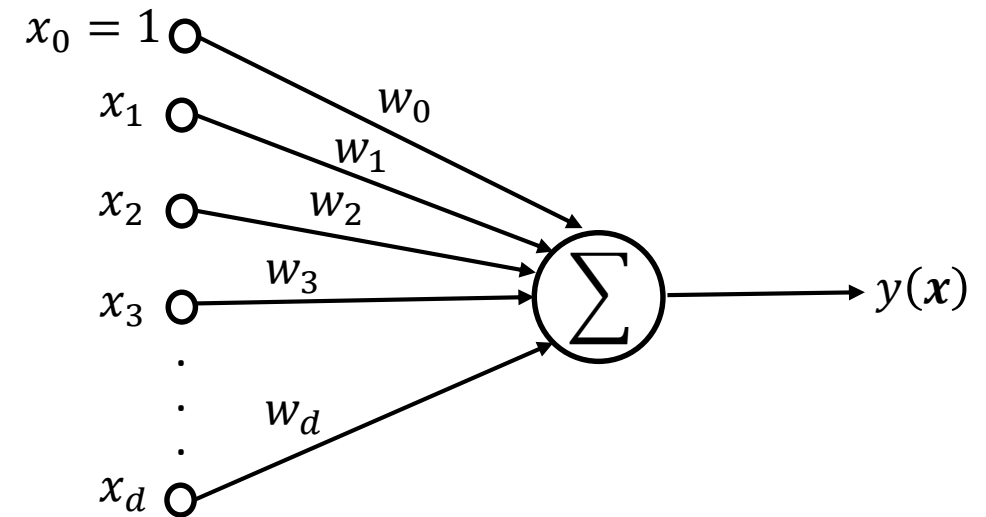
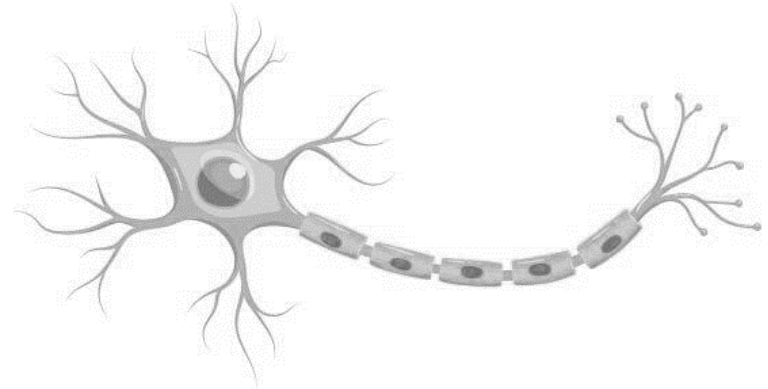


# Neuron

---

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

- In essence, perceptron and LMS algorithms are two different approaches for the training of the model of a neuron.
- Perceptron: Rosenblatt (trained with perceptron algorithm).
- Adaline (adaptive linear element): Widrow and Hoff (trained with LMS).

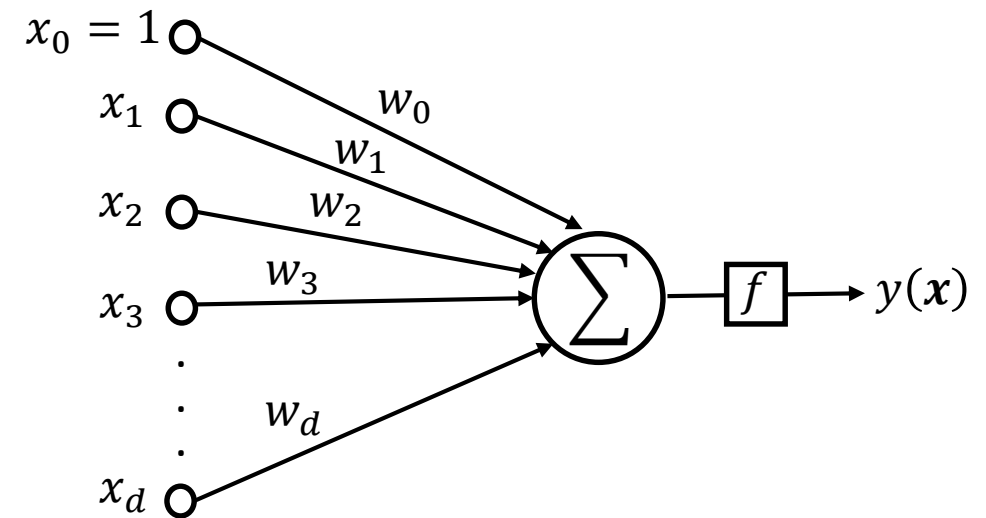
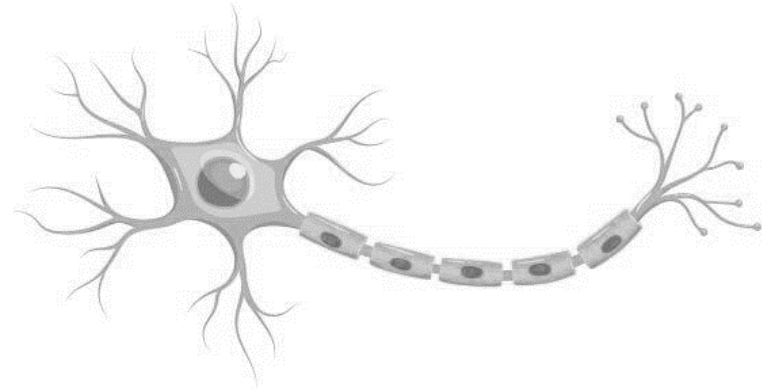


# Neuron

---

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

- In essence, perceptron and LMS algorithms are two different approaches for the training of the model of a neuron.
- Perceptron: Rosenblatt (trained with perceptron algorithm).
- Adaline (adaptive linear element): Widrow and Hoff (trained with LMS).
- This model is the fundamental module of Neural Networks using different activation functions.



# Classification with basis functions

---

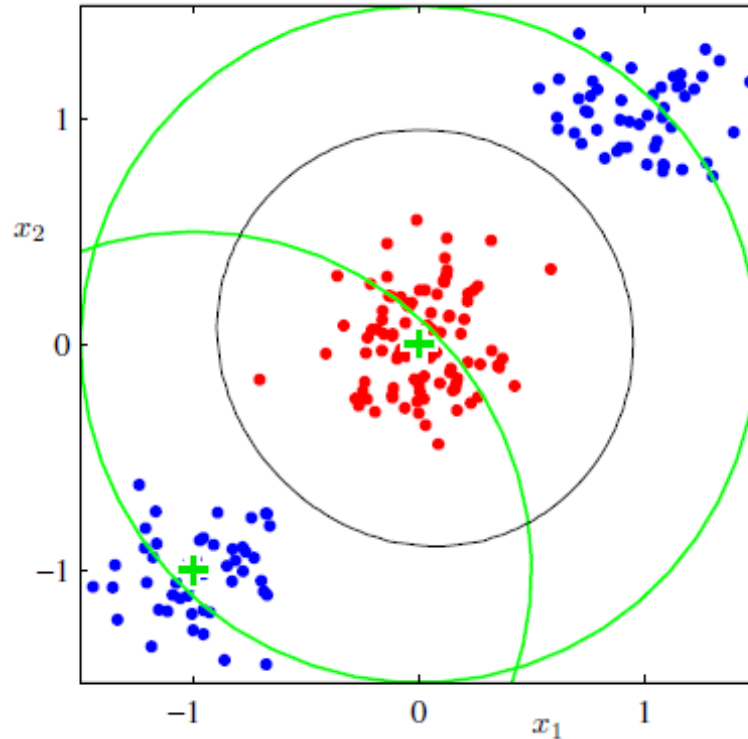
- So far we have considered models that work directly with the **original input** vector  $\mathbf{x}$ , e.g.,  $y(\mathbf{x}) = \mathbf{w}^t \mathbf{x}$ .
- Nevertheless, all of the algorithms are also applicable when we first non-linearly transform the input using a vector of **basis functions**  $\boldsymbol{\varphi}(\mathbf{x})$ , e.g.:

$$y(\mathbf{x}) = \mathbf{w}^t \boldsymbol{\varphi}(\mathbf{x})$$

- With this trick classes that are linearly separable in the feature space  $\boldsymbol{\varphi}(\mathbf{x})$  do not need to be linearly separable in the observation space  $\mathbf{x}$ .
- **Note:** Basis functions are fixed, **not learned**! They need to be selected carefully. They correspond also to sources of inductive bias.

# Classification with basis functions - Example

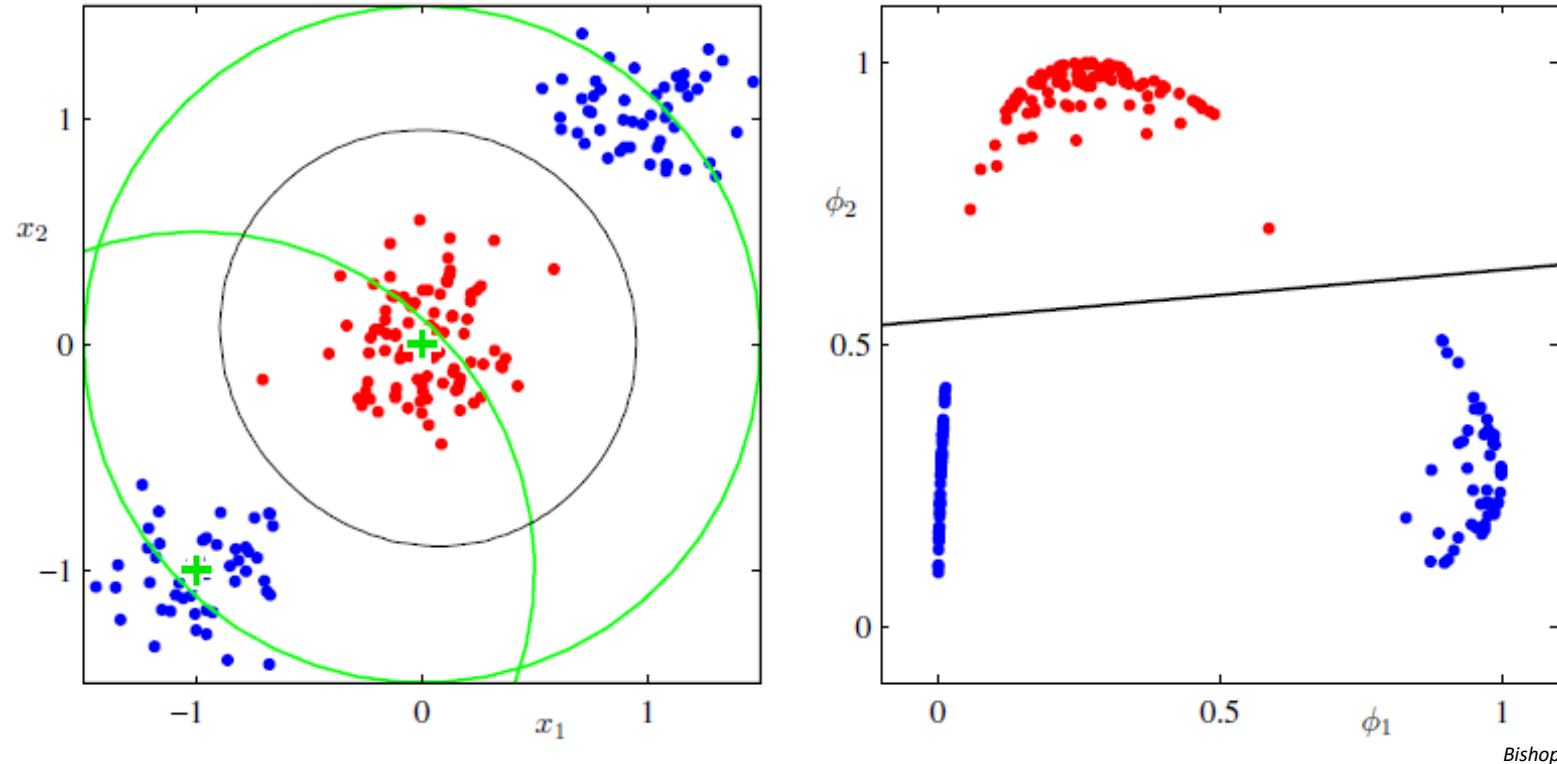
---



**Left:** original input space  $(x_1, x_2)$

$$\varphi_1(\mathbf{x}) = \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^T(\mathbf{x} - \boldsymbol{\mu}_1)\right) \text{ and } \varphi_2(\mathbf{x}) = \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_2)^T(\mathbf{x} - \boldsymbol{\mu}_2)\right)$$

# Classification with basis functions - Example



**Left:** original input space  $(x_1, x_2)$  , **Right:** space of two gaussian basis functions with centers shown by the green crosses:

$$\varphi_1(\mathbf{x}) = \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^T(\mathbf{x} - \boldsymbol{\mu}_1)\right) \text{ and } \varphi_2(\mathbf{x}) = \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_2)^T(\mathbf{x} - \boldsymbol{\mu}_2)\right)$$

# Classification strategies (revisited)

---

- Probabilistic Generative Models
  - class-conditional densities  $p(\mathbf{x}|\omega_i)$
  - prior class probabilities  $P(\omega_i)$
  - Via Bayes we get posterior class probabilities  $P(\omega_i|\mathbf{x})$
- Linear Discriminant functions
  - Direct mapping of input to target  $t = y(\mathbf{x}, \mathbf{w})$
- Probabilistic Discriminative Models
  - Determining directly the posterior class probabilities  $P(\omega_i|\mathbf{x})$
- Logistic Regression

# Probabilistic Generative Model (revisited)

---

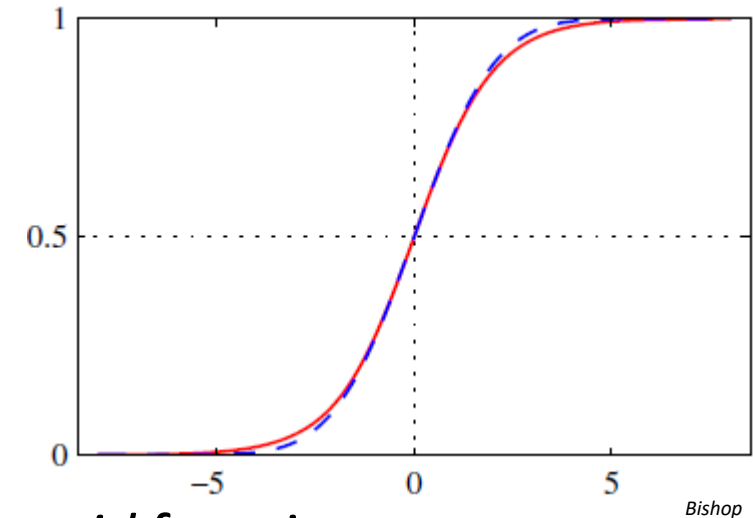
- For two classes:

$$P(\omega_1|x) = \frac{p(x|\omega_1)P(\omega_1)}{p(x)} = \frac{p(x|\omega_1)P(\omega_1)}{p(x|\omega_1)P(\omega_1) + p(x|\omega_2)P(\omega_2)}$$

we can rewrite this equation as:

$$P(\omega_1|x) = \frac{1}{1+\exp(-a)} = \sigma(a)$$

where  $a = \ln \frac{p(x|\omega_1)P(\omega_1)}{p(x|\omega_2)P(\omega_2)}$  and  $\sigma(a)$  is the *logistic sigmoid* function.



# Probabilistic Generative Model (revisited)

---

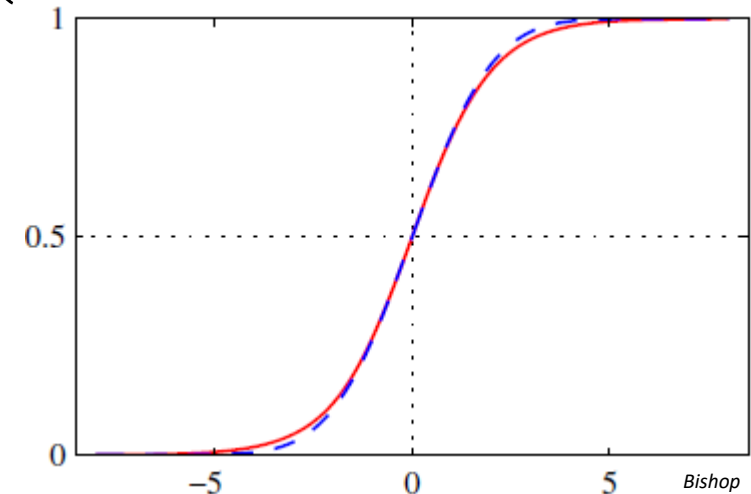
- For number of classes  $c > 2$ :

$$P(\omega_i|x) = \frac{p(x|\omega_i)P(\omega_i)}{p(x)} = \frac{p(x|\omega_i)P(\omega_i)}{\sum_j p(x|\omega_j)P(\omega_j)}$$

we can rewrite this equation as:

$$P(\omega_i|x) = \frac{\exp(a_i)}{\sum_j \exp(a_j)}$$

where  $a_i = \ln(p(x|\omega_i)P(\omega_i))$ . This is known as the *normalized exponential* or the *softmax function*.





# Probabilistic Generative Model (revisited)

---

- Recall that for, e.g., a two-class problem:

$$P(\omega_1|x) = \frac{p(x|\omega_1)P(\omega_1)}{p(x|\omega_1)P(\omega_1) + p(x|\omega_2)P(\omega_2)} = \sigma \left( \ln \frac{p(x|\omega_1)P(\omega_1)}{p(x|\omega_2)P(\omega_2)} \right)$$

# Probabilistic Generative Model (revisited)

---

- Recall that for, e.g., a two-class problem:

$$P(\omega_1|x) = \frac{p(x|\omega_1)P(\omega_1)}{p(x|\omega_1)P(\omega_1) + p(x|\omega_2)P(\omega_2)} = \sigma \left( \ln \frac{p(x|\omega_1)P(\omega_1)}{p(x|\omega_2)P(\omega_2)} \right)$$

$$g(\mathbf{x}) = g_1(\mathbf{x}) - g_2(\mathbf{x})$$

$$g(\mathbf{x}) = P(\omega_1 | \mathbf{x}) - P(\omega_2 | \mathbf{x})$$

$$g(\mathbf{x}) = \ln \frac{p(\mathbf{x}|\omega_1)}{p(\mathbf{x}|\omega_2)} + \ln \frac{P(\omega_1)}{P(\omega_2)}$$

*slide 38 of Bayes Theory lecture*

# Probabilistic Generative Model (revisited)

---

- Recall that for, e.g., a two-class problem:

$$P(\omega_1|x) = \frac{p(x|\omega_1)P(\omega_1)}{p(x|\omega_1)P(\omega_1) + p(x|\omega_2)P(\omega_2)} = \sigma \left( \ln \frac{p(x|\omega_1)P(\omega_1)}{p(x|\omega_2)P(\omega_2)} \right)$$

and if we adopt Gaussian class conditional densities with same covariance matrix across classes we end up with:

$$P(\omega_1|x) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

where

$$\mathbf{w} = \mathbf{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \text{ and}$$

$$w_0 = -\frac{1}{2}\boldsymbol{\mu}_1^T \mathbf{\Sigma}^{-1} \boldsymbol{\mu}_1 + \frac{1}{2}\boldsymbol{\mu}_2^T \mathbf{\Sigma}^{-1} \boldsymbol{\mu}_2 + \ln \frac{P(\omega_1)}{P(\omega_2)}$$

- We can use Maximum Likelihood to estimate  $\boldsymbol{\mu}_1$ ,  $\boldsymbol{\mu}_2$  and  $\mathbf{\Sigma}$ .

$$g(\mathbf{x}) = g_1(\mathbf{x}) - g_2(\mathbf{x})$$

$$g(\mathbf{x}) = P(\omega_1 | \mathbf{x}) - P(\omega_2 | \mathbf{x})$$

$$g(\mathbf{x}) = \ln \frac{p(\mathbf{x}|\omega_1)}{p(\mathbf{x}|\omega_2)} + \ln \frac{P(\omega_1)}{P(\omega_2)}$$

slide 38 of Bayes Theory lecture

# Logistic Regression

---

- Thus we have shown that the posterior class probability **under certain conditions** can be written as a logistic sigmoid acting on a linear function of the feature vector:

$$P(\omega_1|\mathbf{x}) = y(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$$

with  $P(\omega_2|x) = 1 - P(\omega_1|x)$ .

- Using basis functions we can rewrite this as:

$$P(\omega_1|\boldsymbol{\varphi}(\mathbf{x})) = y(\boldsymbol{\varphi}(\mathbf{x})) = \sigma(\mathbf{w}^T \boldsymbol{\varphi}(\mathbf{x}))$$

with  $P(\omega_2|\boldsymbol{\varphi}(\mathbf{x})) = 1 - P(\omega_1|\boldsymbol{\varphi}(\mathbf{x}))$ .

- In statistics this model is known as *logistic regression*. However, this is a model for **classification not regression**.

# Logistic Regression

---

$$P(\omega_1|\boldsymbol{\varphi}(\mathbf{x})) = y(\boldsymbol{\varphi}(\mathbf{x})) = \sigma(\mathbf{w}^T \boldsymbol{\varphi}(\mathbf{x}))$$

with  $P(\omega_2|\boldsymbol{\varphi}(\mathbf{x})) = 1 - P(\omega_1|\boldsymbol{\varphi}(\mathbf{x}))$ .

- For a  $d$  –dimensional feature space  $\boldsymbol{\varphi}$ , this model has  $d + 1$  adjustable parameters (weight of the linear function).
- If we have tried to fit gaussian class conditional densities we would need to estimate  $2d$  parameters for the means and  $\mathcal{O}(d^2)$  parameters for the shared covariance matrix!
- Thus, even if the logistic regression parameters scale linearly with  $d$  the probabilistic generative model approach scales quadratically with  $d$ !
- There is clear advantage to work with the logistic regression model directly!

# Logistic Regression–Training(Maximum Likelihood)

---

$$P(\omega_1|\boldsymbol{\varphi}(\mathbf{x})) = y(\boldsymbol{\varphi}(\mathbf{x})) = \sigma(\mathbf{w}^T \boldsymbol{\varphi}(\mathbf{x}))$$

- In order to adjust the weights of the logistic regression model we can try Maximum Likelihood.
- We will make use of the fact that:

$$\frac{d\sigma(a)}{da} = \sigma(a)(1 - \sigma(a))$$

- Thus with a training set of  $N$  samples of the form  $\{\boldsymbol{\varphi}(\mathbf{x}_n), t_n\}$  with  $t_n \in \{0,1\}$  the likelihood function can be written as:

$$p(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N y_n^{t_n} (1 - y_n)^{1-t_n},$$

where  $y_n = P(\omega_1|\boldsymbol{\varphi}(\mathbf{x}_n)) = \sigma(\mathbf{w}^T \boldsymbol{\varphi}(\mathbf{x}_n))$  and  $\mathbf{t} = (t_1, \dots, t_N)$

---

# Logistic Regression–Training(Maximum Likelihood)

---

- We can **define an error function** and instead of maximizing the likelihood:

$$p(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N y_n^{t_n} (1 - y_n)^{1-t_n},$$

we can **minimize the error function** which is the negative logarithm of the likelihood:

$$E(\mathbf{w}) = -\ln p(\mathbf{t}|\mathbf{w}) = -\sum_{n=1}^N [t_n \ln y_n + (1 - t_n) \ln(1 - y_n)]$$

- This error function is called *cross entropy loss* and its gradient is:

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \boldsymbol{\varphi}(\mathbf{x}_n)$$

# Logistic Regression–Training(SGD)

---

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \boldsymbol{\varphi}(\mathbf{x}_n) = \sum_{n=1}^N \nabla E_n(\mathbf{w})$$

- Unfortunately, we cannot derive a closed form solution for  $\mathbf{w}$  as  $y(\boldsymbol{\varphi}(\mathbf{x})) = \sigma(\mathbf{w}^T \boldsymbol{\varphi}(\mathbf{x}))$  is nonlinear on  $\mathbf{w}$ .
- Nevertheless,  $E(\mathbf{w})$  is a convex function of  $\mathbf{w}$  and we could use a sequential learning algorithm for the logistic regression model training.
- In particular with SGD we get:

$$\mathbf{w}^{i+1} = \mathbf{w}^i - \eta \nabla E_n(\mathbf{w}^i) \quad \rightarrow \quad \mathbf{w}^{i+1} = \mathbf{w}^i - \eta (y_n - t_n) \boldsymbol{\varphi}(\mathbf{x}_n)$$

- **Classification:** a test sample is assigned to  $\omega_1$  if  $P(\omega_1 | \boldsymbol{\varphi}(\mathbf{x})) = \sigma(\mathbf{w}^T \boldsymbol{\varphi}(\mathbf{x})) \geq \frac{1}{2}$





ARISTOTLE UNIVERSITY OF THESSALONIKI



FACULTY OF ENGINEERING

# Questions?

*Pattern Recognition & Machine Learning*

*Linear Discriminant Functions and Models*