

Deep learning with Neural Networks

Logistic Regression, Stochastic Optimization & Regularization

Pablo Martínez Olmos, pamartin@ing.uc3m.es

A binary classification problem

- Training database: $\mathcal{D} \doteq (\mathbf{x}^{(i)}, y^{(i)})_{i=1}^N \quad \mathbf{x}^{(i)} \in \mathbb{R}^m \quad y^{(i)} \in \{0, 1\}$
- Goal: to propose an **hypothesis function** to estimate the **most likely class** of a new data point

$$y^{(*)} \approx \arg \max_{y \in \{0,1\}} h_{\theta}(\mathbf{x}^{(*)}, y)$$



Parameters to be optimized using
the training set

Logistic Regression: the model

Binary Logistic Regression

- Discriminative classifier:

$$P(y = 1|\boldsymbol{x}) = \frac{1}{1 + e^{-(\boldsymbol{w}^T \boldsymbol{x} + w_0)}} \doteq \sigma(\boldsymbol{w}^T \boldsymbol{x} + w_0)$$

Binary Logistic Regression

- Discriminative classifier:

$$P(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + w_0)}} \doteq \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$



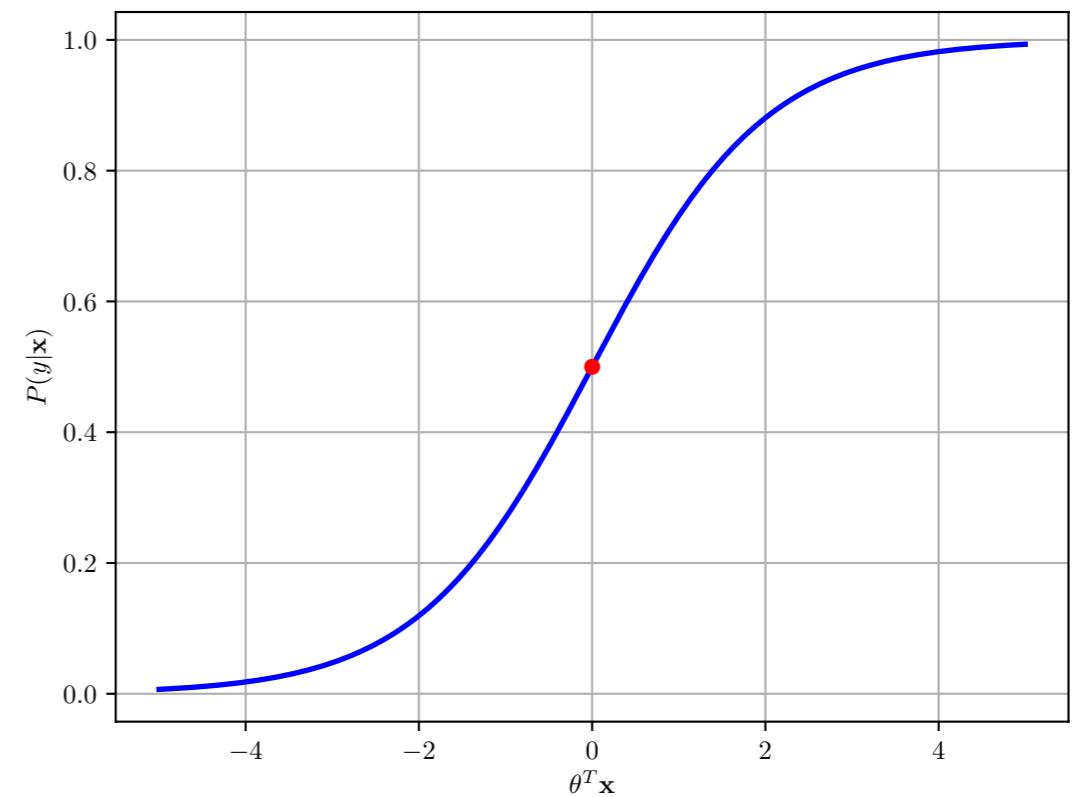
Generalized linear model

Binary Logistic Regression

- Discriminative classifier:

$$P(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + w_0)}} \doteq \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

Generalized linear model



Binary Logistic Regression

- Discriminative classifier:

$$P(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + w_0)}} \doteq \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

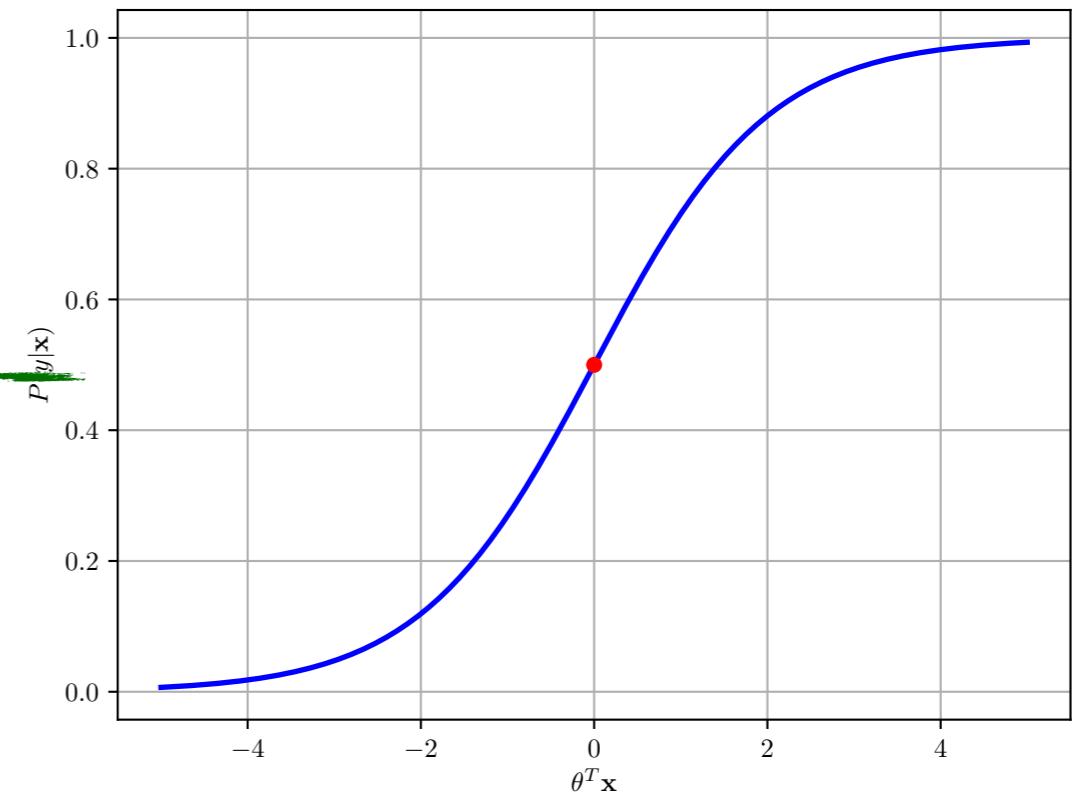
Generalized linear model



- Decision boundary:

$$\{\mathbf{x} \in \mathbb{R}^m : \mathbf{w}^T \mathbf{x} + w_0 = 0\}$$

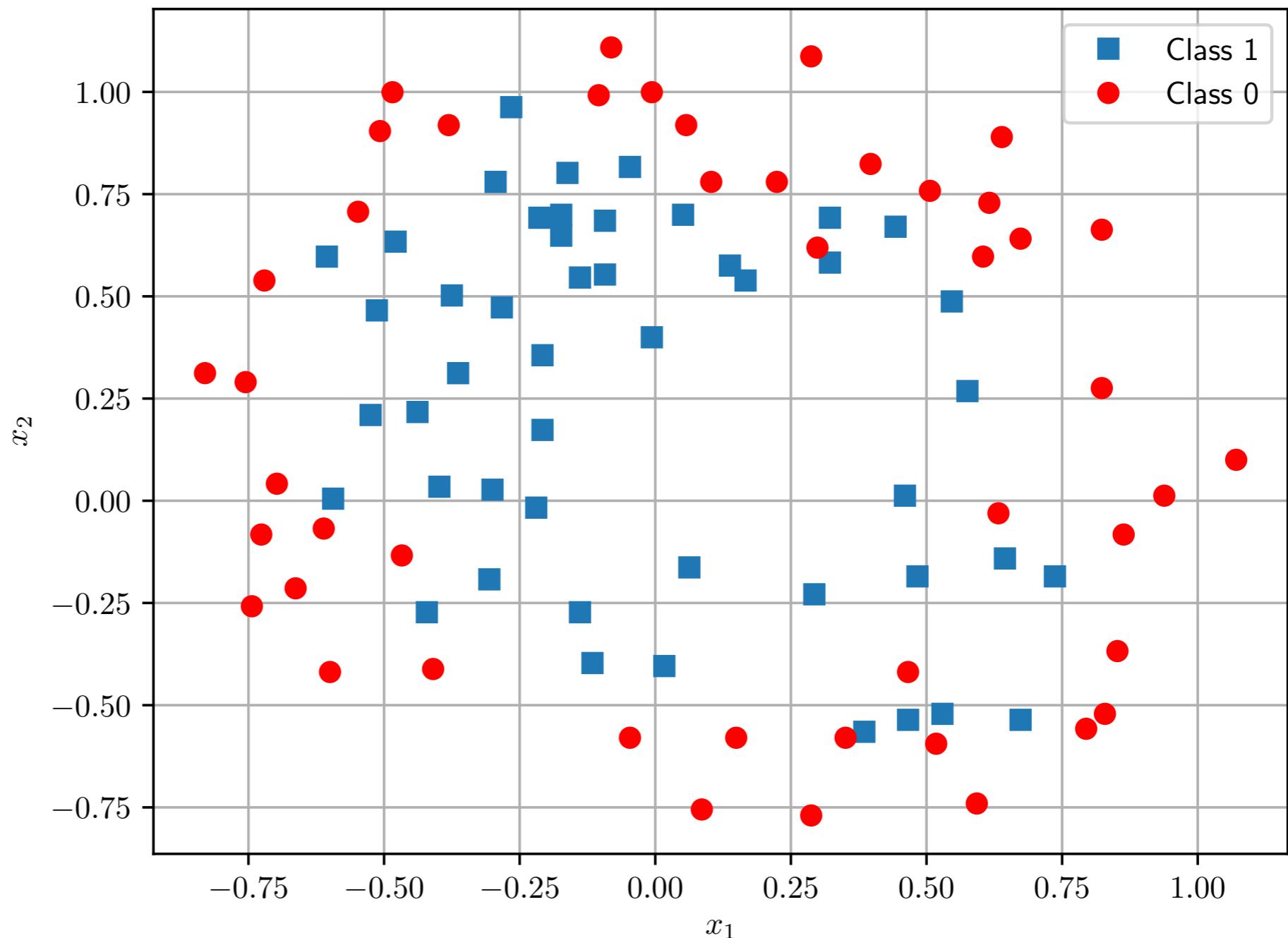
Hyperplane! Linear classifier



Binary Logistic Regression

$$P(y = 1|x) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + w_0)}} \doteq \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

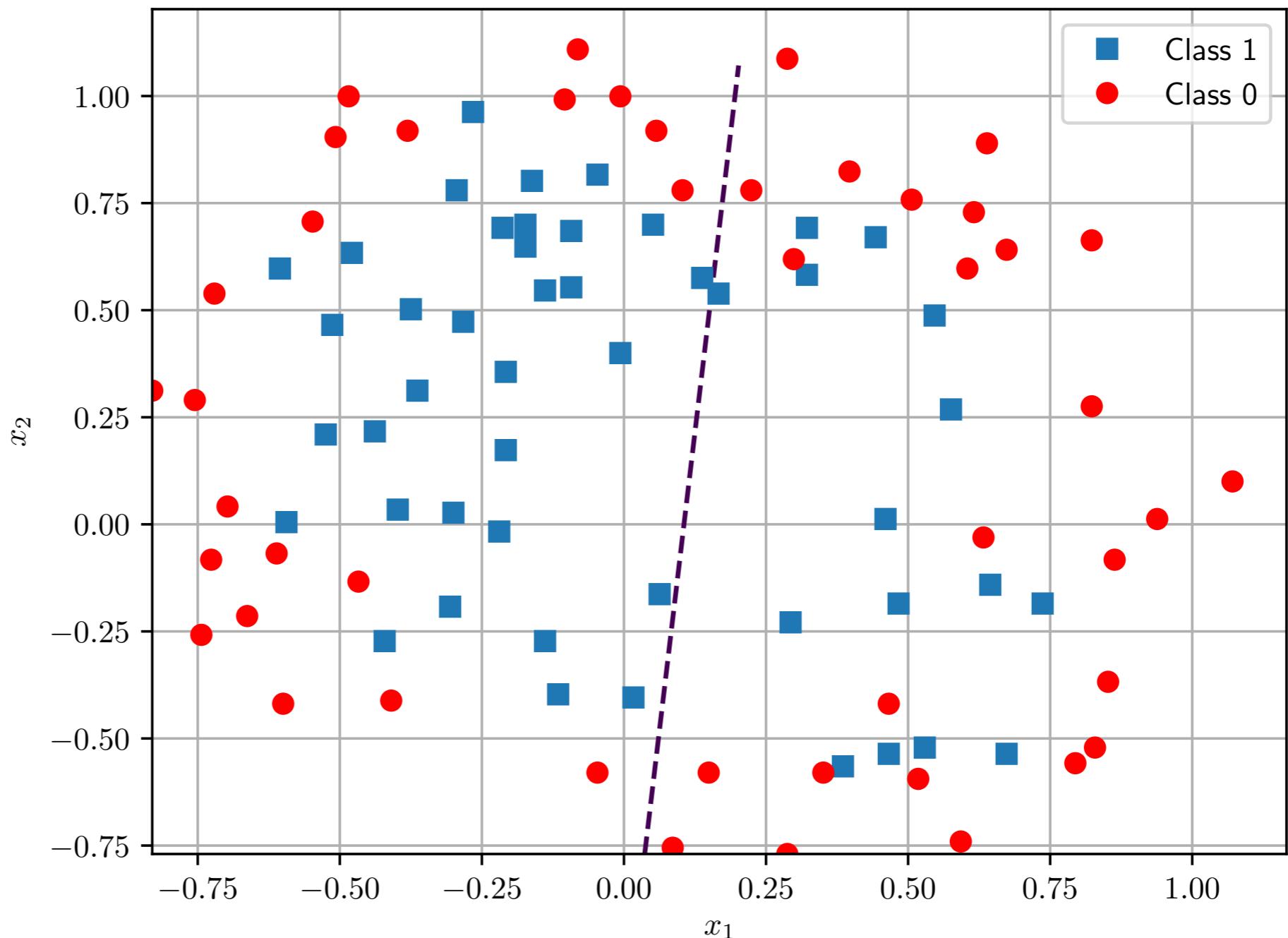
Training Data



Binary Logistic Regression

$$P(y = 1|x) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + w_0)}} \doteq \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

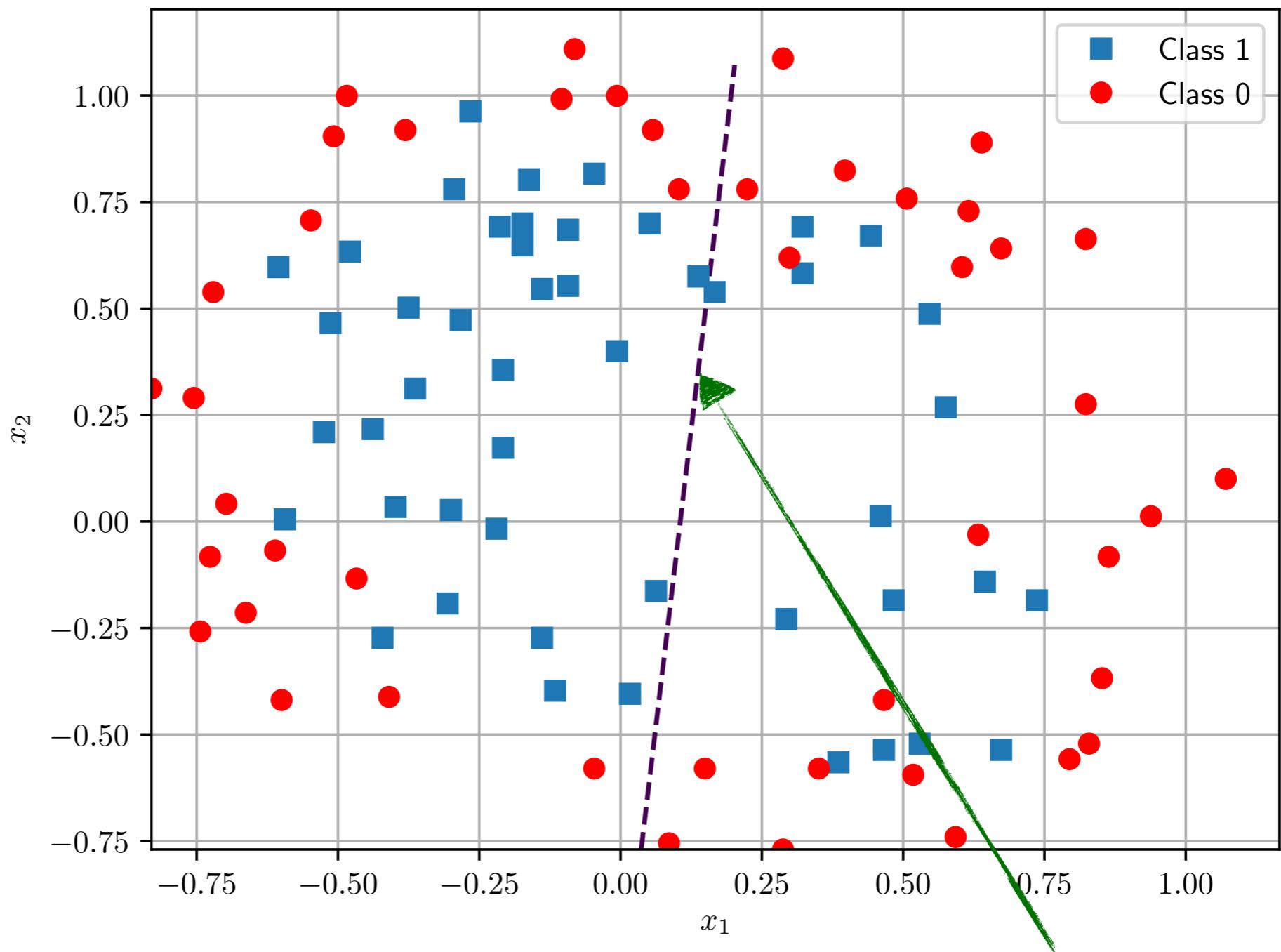
Training Data



Binary Logistic Regression

$$P(y = 1|x) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + w_0)}} \doteq \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

Training Data



LR decision boundary with linear features

Binary Logistic Regression. Adding new features

- The problem is easy to visualize
- An ellipsoid clearly fits better as decision boundary than a straight line



- Expand input space with extra degrees of freedom

- Quadratic features

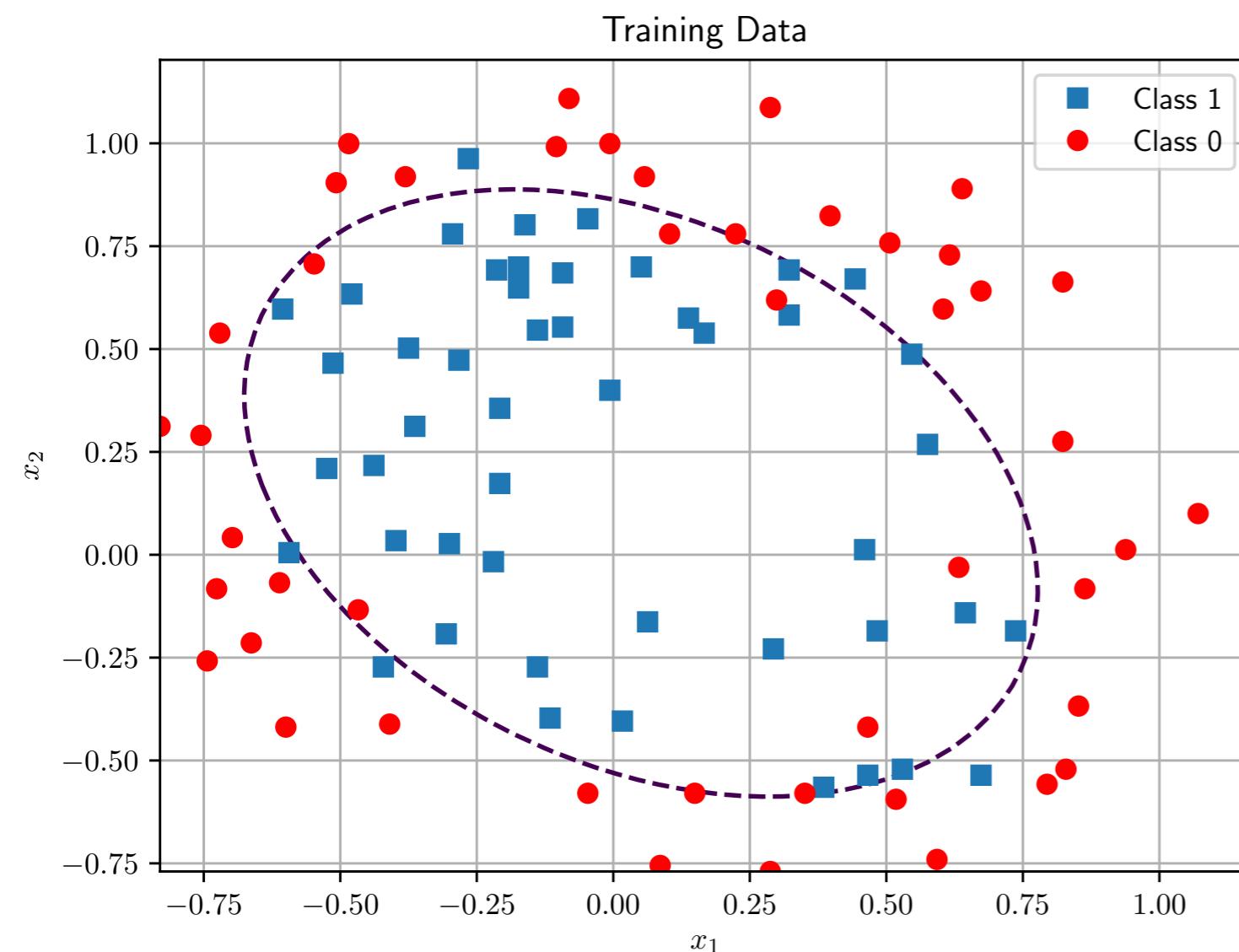
$$\mathbf{x} = [x_1, x_2]$$

$$\phi(\mathbf{x}) = [x_1, x_2, x_1^2, x_2^2, x_1 x_2]$$

$$P(y = 1|\mathbf{x}) = \sigma(\mathbf{w}^T \phi(\mathbf{x}) + w_0)$$

Binary Logistic Regression. Adding new features

- The problem is easy to visualize
- An ellipsoid clearly fits better as decision boundary than a straight line



- Expand input space with extra degrees of freedom

- Quadratic features

$$\mathbf{x} = [x_1, x_2]$$

$$\phi(\mathbf{x}) = [x_1, x_2, x_1^2, x_2^2, x_1 x_2]$$

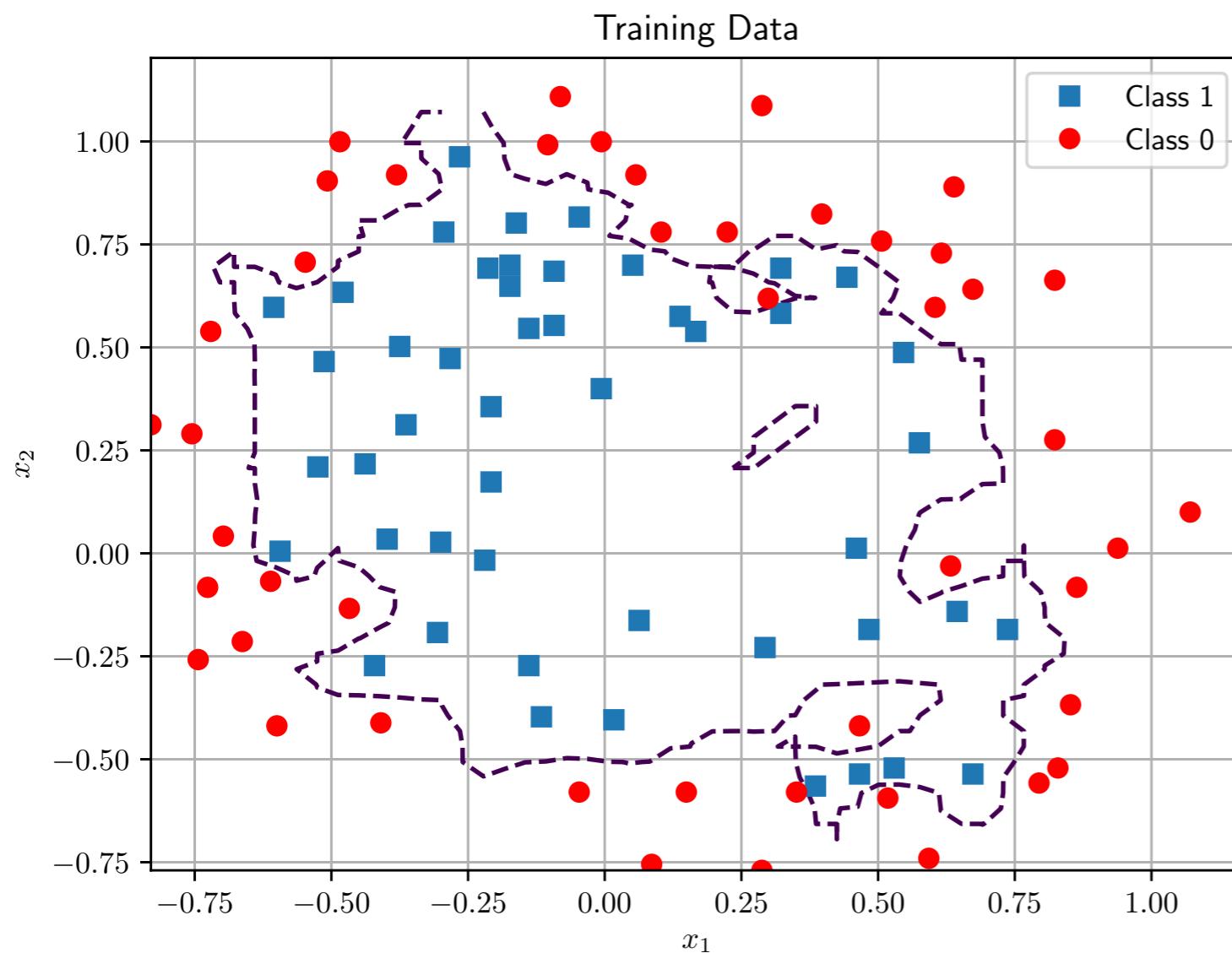
$$P(y = 1|\mathbf{x}) = \sigma(\mathbf{w}^T \phi(\mathbf{x}) + w_0)$$

Binary Logistic Regression. Adding more and more features

- By adding new features we allow the model to fit (even overfit) to the training data
- We need a validation set to evaluate the model's exposure to overfitting

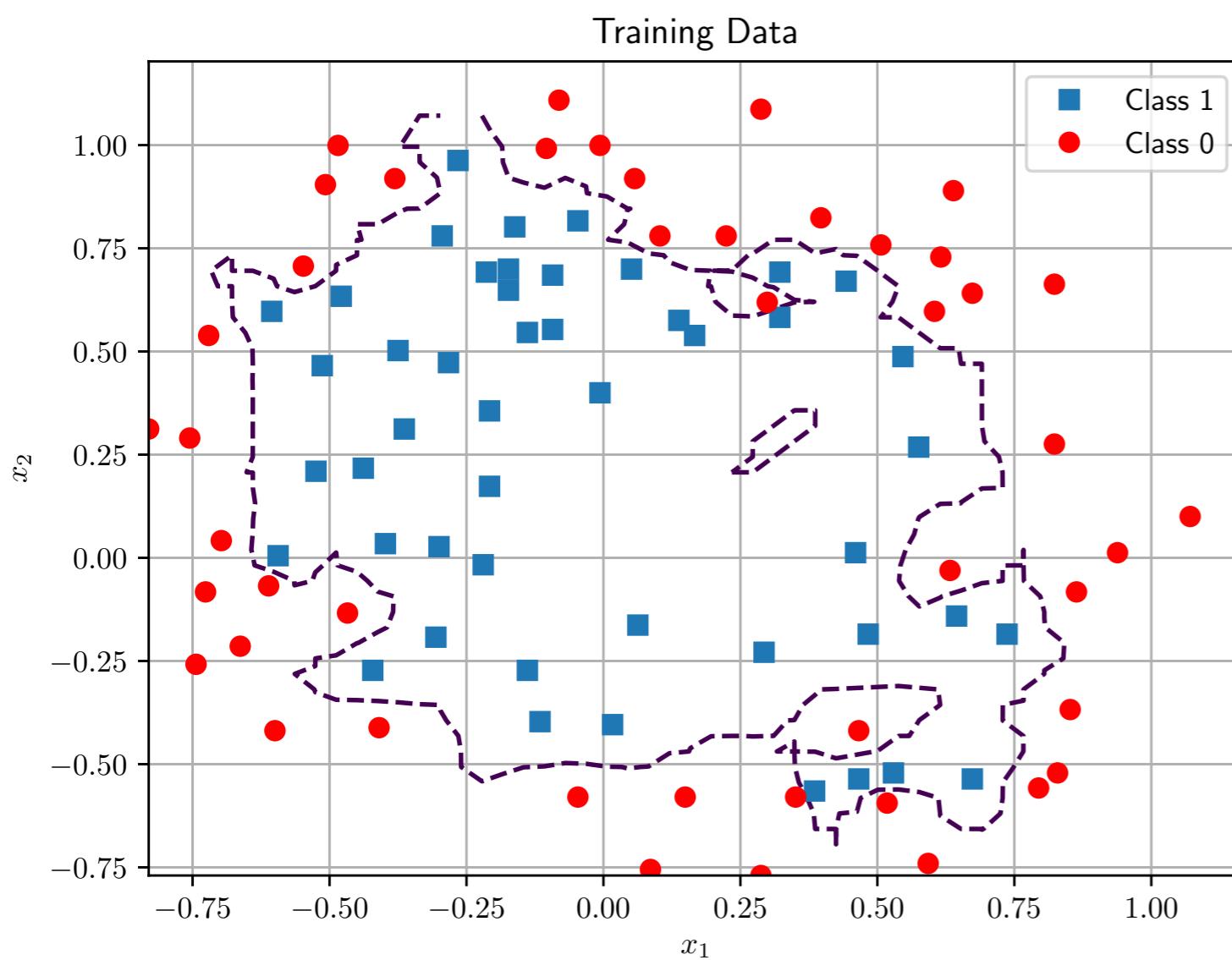
Binary Logistic Regression. Adding more and more features

- By adding new features we allow the model to fit (even overfit) to the training data
- We need a validation set to evaluate the model's exposure to overfitting



Binary Logistic Regression. Adding more and more features

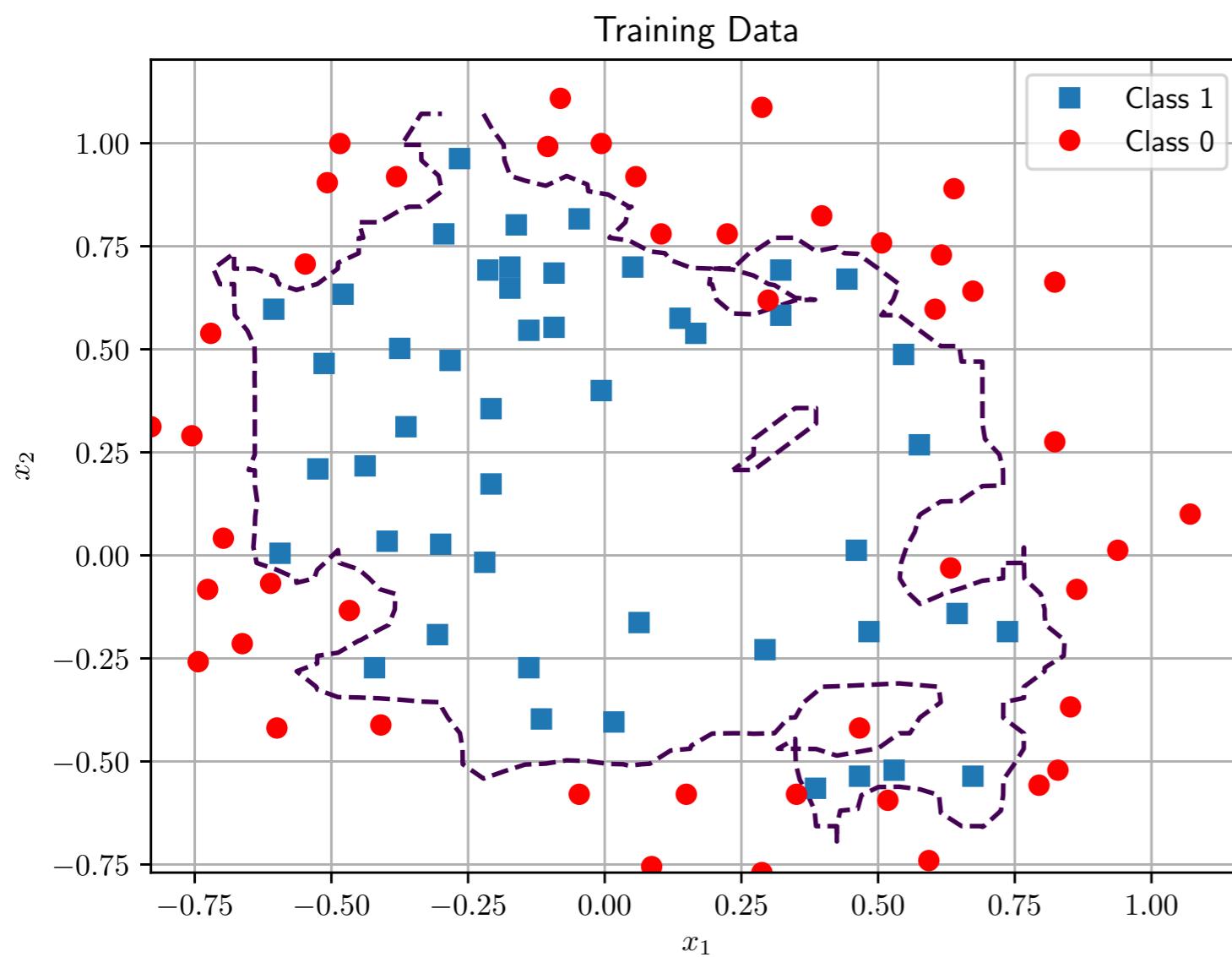
- By adding new features we allow the model to fit (even overfit) to the training data
- We need a validation set to evaluate the model's exposure to overfitting



Provided you find the right set of features!

Binary Logistic Regression. Adding more and more features

- By adding new features we allow the model to fit (even overfit) to the training data
- We need a validation set to evaluate the model's exposure to overfitting

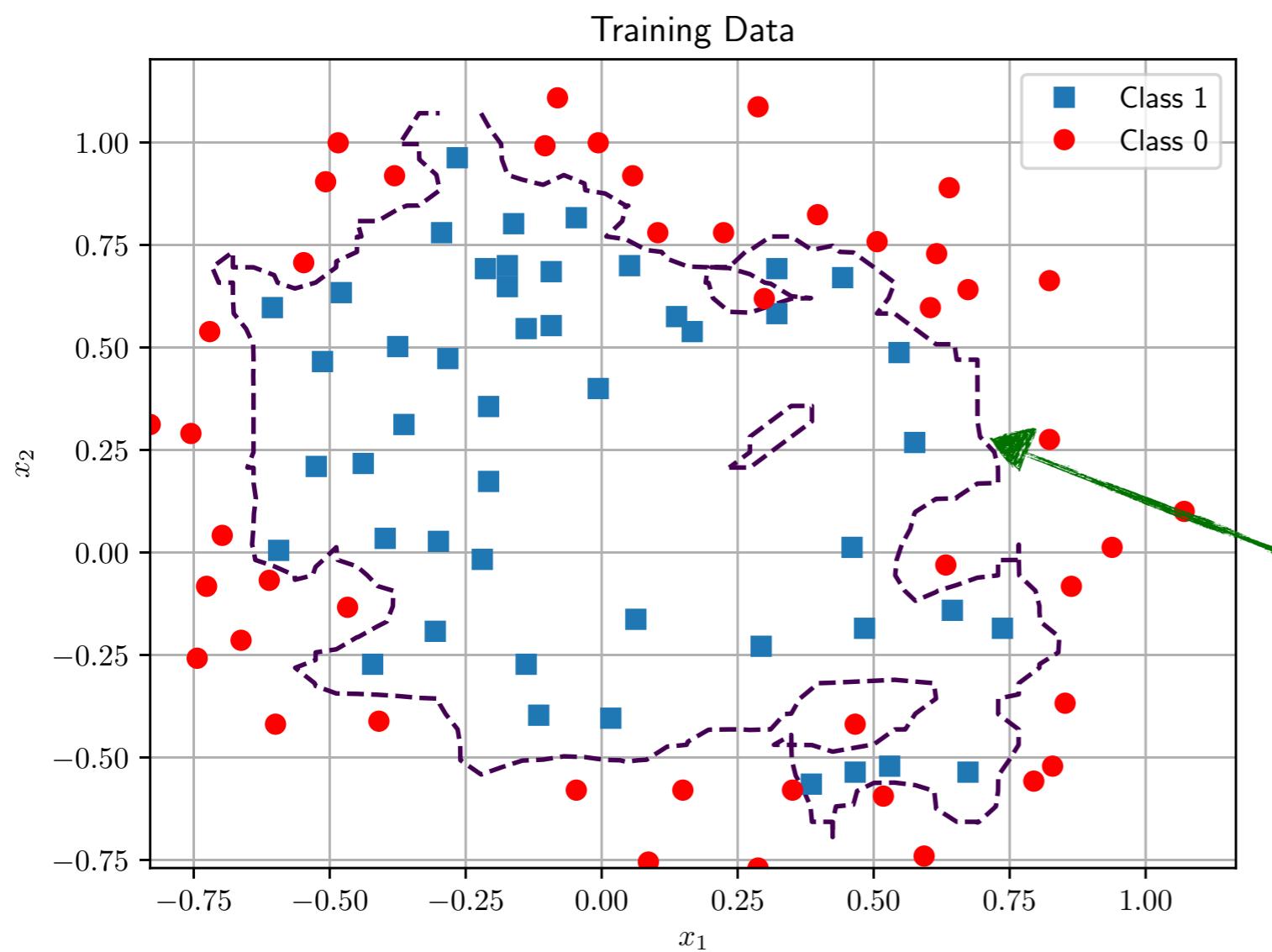


Provided you find the right set of features!

**VERY hard problem in general.
Requires domain-specific knowledge.**

Binary Logistic Regression. Adding more and more features

- By adding new features we allow the model to fit (even overfit) to the training data
- We need a validation set to evaluate the model's exposure to overfitting



Provided you find the right set of features!

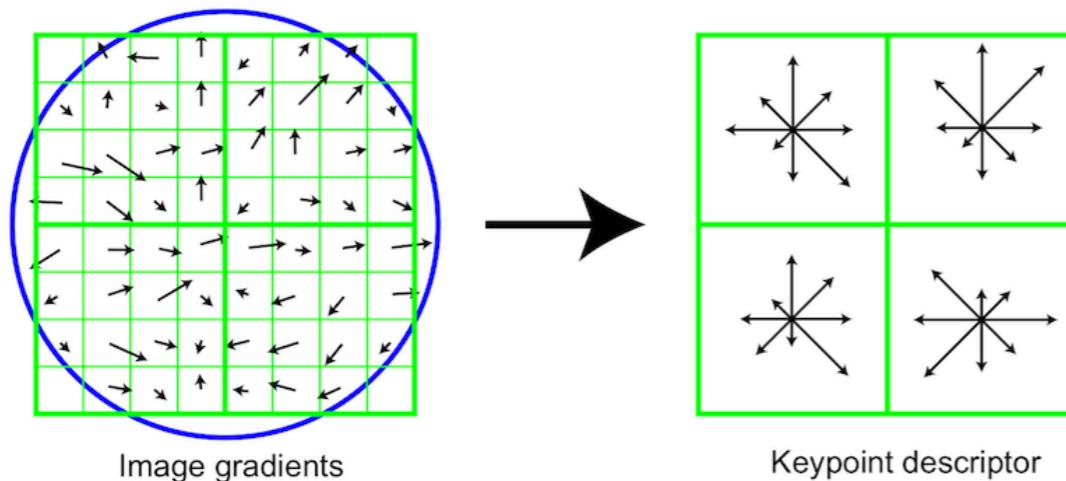
**VERY hard problem in general.
Requires domain-specific knowledge.**

What is the equation of such a weird boundary?

Ways to go ...

Ways to go ...

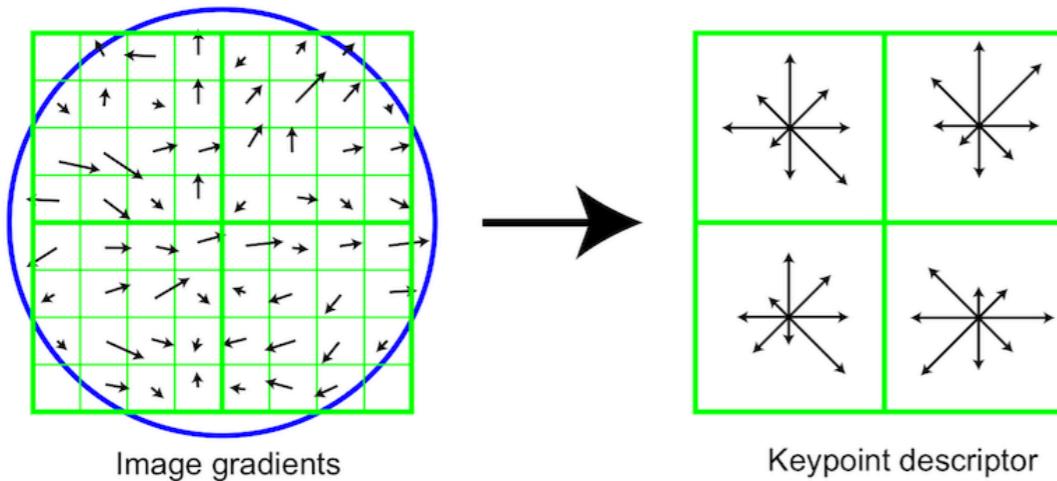
- Manual **feature design**: complex in high-dimensional settings, domain-specific knowledge



E.g. SIFT features in an image
(Image taken from [this excellent blog](#))

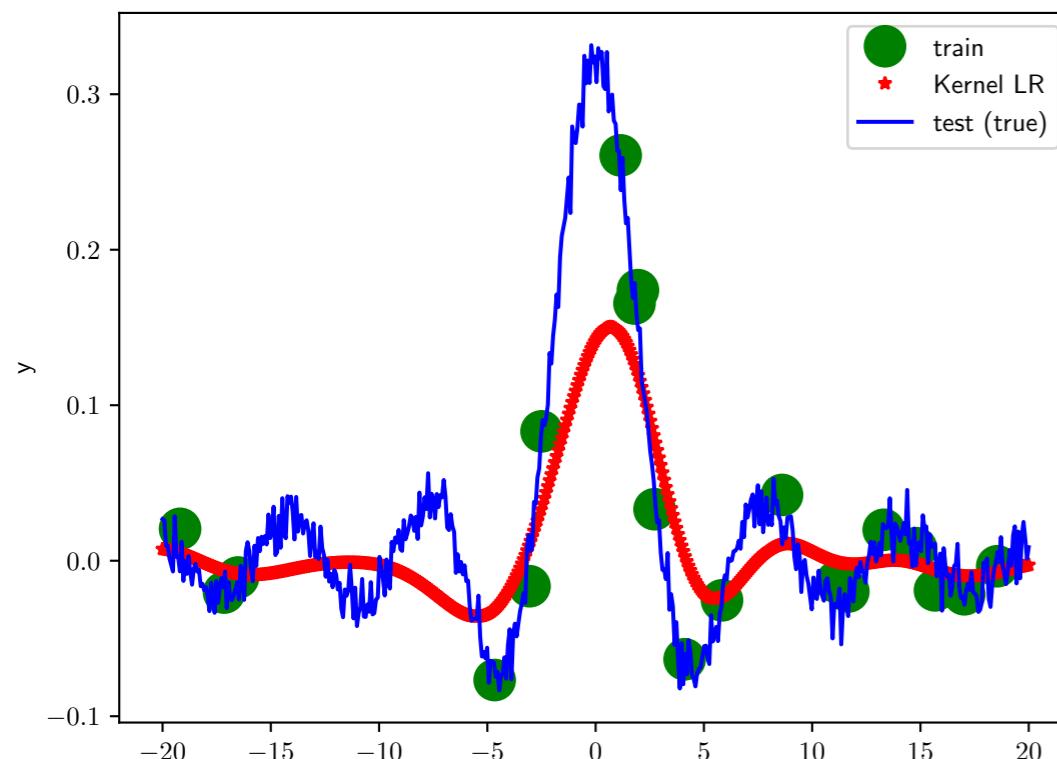
Ways to go ...

- Manual **feature design**: complex in high-dimensional settings, domain-specific knowledge



E.g. SIFT features in an image
(Image taken from [this excellent blog](#))

- **Kernel methods**: you do not think on feature space, but on the data space (~ interpolate with your neighbours). Computationally very complex for large databases.

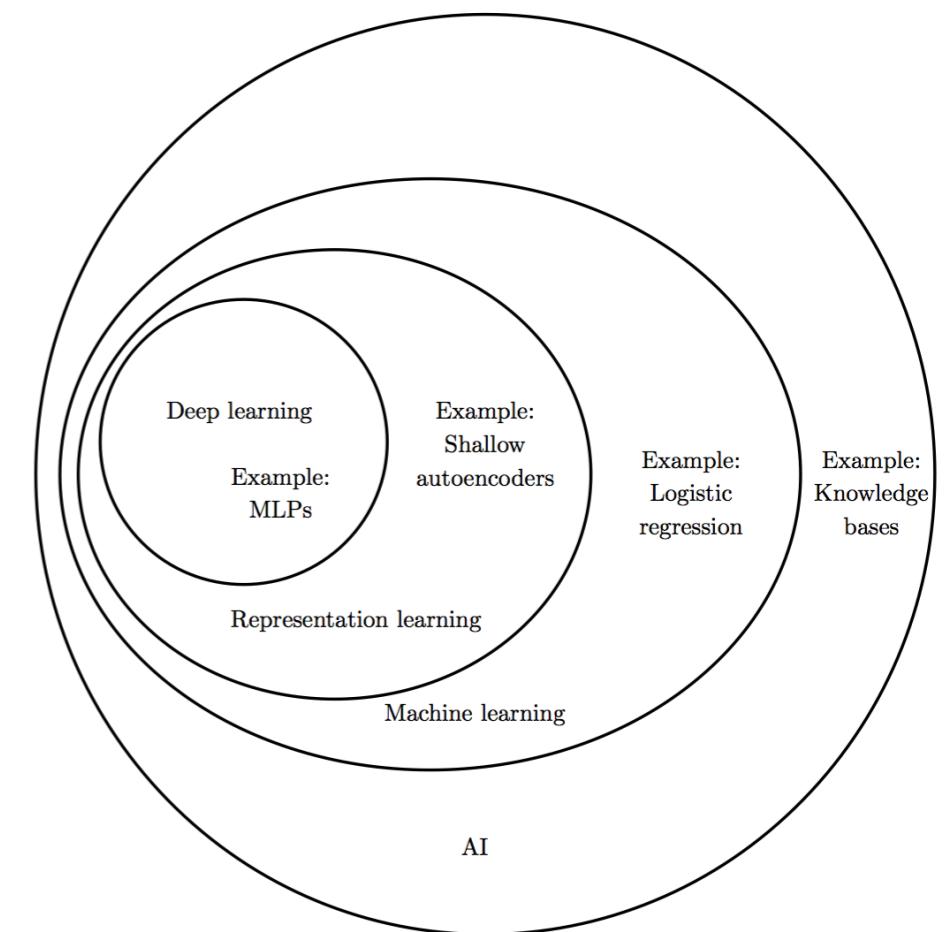
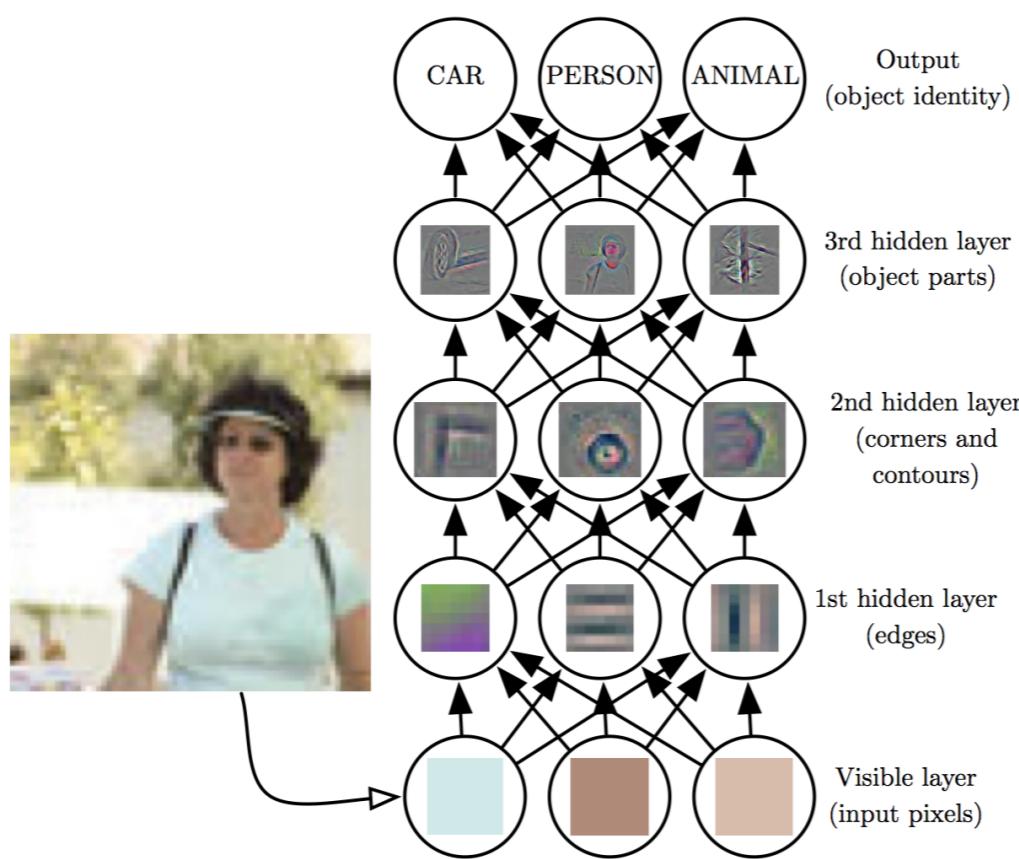


Example of kernel regression with Gaussian kernel

Ways to go ...

- **Representation Learning:** use raw-data and a highly flexible family of transformation functions.

Neural networks belong to this class of methods!!



Source: Deep Learning, Good Fellow et al. 2017

Logistic Regression: training

Binary Logistic Regression. Maximum likelihood training

- Hypothesis function: $P(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + w_0)}} \doteq \sigma(\mathbf{w}^T \mathbf{x} + w_0)$
- Training database:
- (Log)-Likelihood function:

Binary Logistic Regression. Maximum likelihood training

- Hypothesis function: $P(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + w_0)}} \doteq \sigma(\mathbf{w}^T \mathbf{x} + w_0)$
- Training database: $\mathcal{D} \doteq (\mathbf{x}^{(i)}, y^{(i)})_{i=1}^N \quad \mathbf{x}^{(i)} \in \mathbb{R}^m \quad y^{(i)} \in \{0, 1\}$
- (Log)-Likelihood function:

Binary Logistic Regression. Maximum likelihood training

- Hypothesis function: $P(y = 1|\boldsymbol{x}) = \frac{1}{1 + e^{-(\boldsymbol{w}^T \boldsymbol{x} + w_0)}} \doteq \sigma(\boldsymbol{w}^T \boldsymbol{x} + w_0)$
- Training database: $\mathcal{D} \doteq (\boldsymbol{x}^{(i)}, y^{(i)})_{i=1}^N \quad \boldsymbol{x}^{(i)} \in \mathbb{R}^m \quad y^{(i)} \in \{0, 1\}$
- (Log)-Likelihood function:

$$P(y|\boldsymbol{x}) = P(y = 1|\boldsymbol{x})^{\mathbb{1}[y=1|\boldsymbol{x}]} P(y = 0|\boldsymbol{x})^{\mathbb{1}[y=0|\boldsymbol{x}]}$$

Binary Logistic Regression. Maximum likelihood training

- Hypothesis function: $P(y = 1|\boldsymbol{x}) = \frac{1}{1 + e^{-(\boldsymbol{w}^T \boldsymbol{x} + w_0)}} \doteq \sigma(\boldsymbol{w}^T \boldsymbol{x} + w_0)$
- Training database: $\mathcal{D} \doteq (\boldsymbol{x}^{(i)}, y^{(i)})_{i=1}^N \quad \boldsymbol{x}^{(i)} \in \mathbb{R}^m \quad y^{(i)} \in \{0, 1\}$
- (Log)-Likelihood function:

$$P(y|\boldsymbol{x}) = P(y = 1|\boldsymbol{x})^{\mathbb{1}[y=1|\boldsymbol{x}]} P(y = 0|\boldsymbol{x})^{\mathbb{1}[y=0|\boldsymbol{x}]}$$

$$P(\boldsymbol{y}|\boldsymbol{X}) = \prod_{i=1}^N P(y^{(i)} = 1|\boldsymbol{x}^{(i)})$$

Binary Logistic Regression. Maximum likelihood training

- Hypothesis function: $P(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + w_0)}} \doteq \sigma(\mathbf{w}^T \mathbf{x} + w_0)$
- Training database: $\mathcal{D} \doteq (\mathbf{x}^{(i)}, y^{(i)})_{i=1}^N \quad \mathbf{x}^{(i)} \in \mathbb{R}^m \quad y^{(i)} \in \{0, 1\}$
- (Log)-Likelihood function:

$$P(y|\mathbf{x}) = P(y = 1|\mathbf{x})^{\mathbb{1}[y=1|\mathbf{x}]} P(y = 0|\mathbf{x})^{\mathbb{1}[y=0|\mathbf{x}]}$$

$$P(\mathbf{y}|\mathbf{X}) = \prod_{i=1}^N P(y^{(i)} = 1|\mathbf{x}^{(i)})$$

$$\begin{aligned} \log P(\mathbf{y}|\mathbf{X}) &= \sum_{i=1}^N \log P(y^{(i)}|\mathbf{x}^{(i)}) \\ &= \sum_{i=1}^N \mathbb{1}[y^{(i)} = 1] \log P(y^{(i)} = 1|\mathbf{x}^{(i)}) + \mathbb{1}[y^{(i)} = 0] \log P(y^{(i)} = 0|\mathbf{x}^{(i)}) \end{aligned}$$

Binary Logistic Regression. Binary Cross Entropy Function

$$P(y = 1|\boldsymbol{x}) = \frac{1}{1 + e^{-(\boldsymbol{w}^T \boldsymbol{x} + w_0)}} \doteq \sigma(\boldsymbol{w}^T \boldsymbol{x} + w_0)$$

Binary Cross Entropy Function

$$\mathcal{L} = -\log P(\boldsymbol{y}|\boldsymbol{X})$$

$$= - \sum_{i=1}^N \mathbb{1}[y^{(i)} = 1] \log P(y^{(i)} = 1|\boldsymbol{x}^{(i)}) + \mathbb{1}[y^{(i)} = 0] \log P(y^{(i)} = 0|\boldsymbol{x}^{(i)})$$

Binary Logistic Regression. Binary Cross Entropy Function

$$P(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + w_0)}} \doteq \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

Binary Cross Entropy Function

$$\mathcal{L} = -\log P(\mathbf{y}|\mathbf{X})$$

$$= - \sum_{i=1}^N \mathbb{1}[y^{(i)} = 1] \log P(y^{(i)} = 1 | \mathbf{x}^{(i)}) + \mathbb{1}[y^{(i)} = 0] \log P(y^{(i)} = 0 | \mathbf{x}^{(i)})$$

Logistic Regression Training

$$\hat{\mathbf{w}}, \hat{w}_0 = \arg \min_{\mathbf{w}, w_0} \mathcal{L}$$

Binary Logistic Regression. Binary Cross Entropy Function

$$P(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + w_0)}} \doteq \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

Binary Cross Entropy Function

$$\mathcal{L} = -\log P(\mathbf{y}|\mathbf{X})$$

$$= - \sum_{i=1}^N \mathbb{1}[y^{(i)} = 1] \log P(y^{(i)} = 1 | \mathbf{x}^{(i)}) + \mathbb{1}[y^{(i)} = 0] \log P(y^{(i)} = 0 | \mathbf{x}^{(i)})$$

\mathcal{L} is convex!!

Logistic Regression Training

$$\hat{\mathbf{w}}, \hat{w}_0 = \arg \min_{\mathbf{w}, w_0} \mathcal{L}$$

Binary Logistic Regression. Binary Cross Entropy Function

$$P(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + w_0)}} \doteq \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

Binary Cross Entropy Function

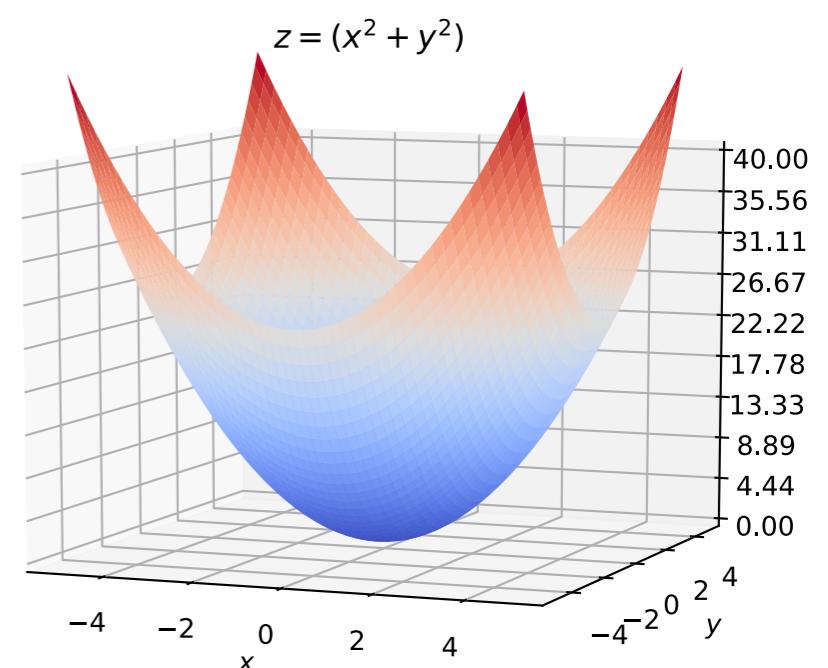
$$\mathcal{L} = -\log P(\mathbf{y}|\mathbf{X})$$

$$= -\sum_{i=1}^N \mathbb{1}[y^{(i)} = 1] \log P(y^{(i)} = 1|\mathbf{x}^{(i)}) + \mathbb{1}[y^{(i)} = 0] \log P(y^{(i)} = 0|\mathbf{x}^{(i)})$$

\mathcal{L} is convex!!

Logistic Regression Training

$$\hat{\mathbf{w}}, \hat{w}_0 = \arg \min_{\mathbf{w}, w_0} \mathcal{L}$$



Logistic Regression Training

Logistic Regression Training

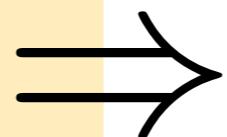
Logistic Regression Training

$$\hat{\mathbf{w}}, \hat{w}_0 = \arg \min_{\mathbf{w}, w_0} \mathcal{L}$$

Logistic Regression Training

Logistic Regression Training

$$\hat{\mathbf{w}}, \hat{w}_0 = \arg \min_{\mathbf{w}, w_0} \mathcal{L}$$

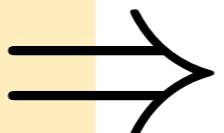


$$\nabla_{\mathbf{w}, w_0} \mathcal{L}(\hat{\mathbf{w}}, \hat{w}_0) = \mathbf{0}$$

Logistic Regression Training

Logistic Regression Training

$$\hat{\mathbf{w}}, \hat{w}_0 = \arg \min_{\mathbf{w}, w_0} \mathcal{L}$$



$$\nabla_{\mathbf{w}, w_0} \mathcal{L}(\hat{\mathbf{w}}, \hat{w}_0) = \mathbf{0}$$

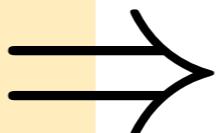
By using that $\frac{d\sigma(a)}{da} = \sigma(a)(1 - \sigma(a))$ it can be shown that

$$\nabla_{\mathbf{w}, w_0} \mathcal{L} = \sum_{i=1}^N \left(\sigma \left(\begin{bmatrix} w_0 & \mathbf{w} \end{bmatrix} \begin{bmatrix} 1.0 \\ \mathbf{x}^{(i)} \end{bmatrix} \right) - y^{(i)} \right) \begin{bmatrix} 1.0 \\ \mathbf{x}^{(i)} \end{bmatrix}$$

Logistic Regression Training

Logistic Regression Training

$$\hat{\mathbf{w}}, \hat{w}_0 = \arg \min_{\mathbf{w}, w_0} \mathcal{L}$$



$$\nabla_{\mathbf{w}, w_0} \mathcal{L}(\hat{\mathbf{w}}, \hat{w}_0) = \mathbf{0}$$

By using that $\frac{d\sigma(a)}{da} = \sigma(a)(1 - \sigma(a))$ it can be shown that

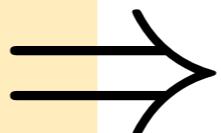
$$\nabla_{\mathbf{w}, w_0} \mathcal{L} = \sum_{i=1}^N \left(\underbrace{\sigma \left(\begin{bmatrix} w_0 & \mathbf{w} \end{bmatrix} \begin{bmatrix} 1.0 \\ \mathbf{x}^{(i)} \end{bmatrix} \right) - y^{(i)}}_{\text{Prediction error}} \right) \begin{bmatrix} 1.0 \\ \mathbf{x}^{(i)} \end{bmatrix}$$

Prediction error
(real scalar)

Logistic Regression Training

Logistic Regression Training

$$\hat{\mathbf{w}}, \hat{w}_0 = \arg \min_{\mathbf{w}, w_0} \mathcal{L}$$



$$\nabla_{\mathbf{w}, w_0} \mathcal{L}(\hat{\mathbf{w}}, \hat{w}_0) = \mathbf{0}$$

By using that $\frac{d\sigma(a)}{da} = \sigma(a)(1 - \sigma(a))$ it can be shown that

$$\nabla_{\mathbf{w}, w_0} \mathcal{L} = \sum_{i=1}^N \left(\underbrace{\sigma \left(\begin{bmatrix} w_0 & \mathbf{w} \end{bmatrix} \begin{bmatrix} 1.0 \\ \mathbf{x}^{(i)} \end{bmatrix} \right) - y^{(i)}}_{\text{Prediction error}} \right) \begin{bmatrix} 1.0 \\ \mathbf{x}^{(i)} \end{bmatrix}$$

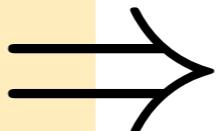
Prediction error
(real scalar)

A curly brace grouping the vector components 1.0 and $\mathbf{x}^{(i)}$ from the equation above.
Vector
(bias + data point)

Logistic Regression Training

Logistic Regression Training

$$\hat{\mathbf{w}}, \hat{w}_0 = \arg \min_{\mathbf{w}, w_0} \mathcal{L}$$



$$\nabla_{\mathbf{w}, w_0} \mathcal{L}(\hat{\mathbf{w}}, \hat{w}_0) = \mathbf{0}$$

By using that $\frac{d\sigma(a)}{da} = \sigma(a)(1 - \sigma(a))$ it can be shown that

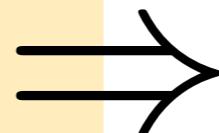
$$\nabla_{\mathbf{w}, w_0} \mathcal{L} = \sum_{i=1}^N \left(\underbrace{\sigma \left(\begin{bmatrix} w_0 & \mathbf{w} \end{bmatrix} \begin{bmatrix} 1.0 \\ \mathbf{x}^{(i)} \end{bmatrix} \right) - y^{(i)}}_{\text{Prediction error (real scalar)}} \right) \underbrace{\begin{bmatrix} 1.0 \\ \mathbf{x}^{(i)} \end{bmatrix}}_{\text{Vector (bias + data point)}}$$

Gradient is a weighted-sum of training data points

Logistic Regression Training

Logistic Regression Training

$$\hat{\mathbf{w}}, \hat{w}_0 = \arg \min_{\mathbf{w}, w_0} \mathcal{L}$$



$$\nabla_{\mathbf{w}, w_0} \mathcal{L}(\hat{\mathbf{w}}, \hat{w}_0) = \mathbf{0}$$

No closed-form solution!

By using that $\frac{d\sigma(a)}{da} = \sigma(a)(1 - \sigma(a))$ it can be shown that

$$\nabla_{\mathbf{w}, w_0} \mathcal{L} = \sum_{i=1}^N \left(\underbrace{\sigma \left(\begin{bmatrix} w_0 & \mathbf{w} \end{bmatrix} \begin{bmatrix} 1.0 \\ \mathbf{x}^{(i)} \end{bmatrix} \right) - y^{(i)}}_{\text{Prediction error (real scalar)}} \right) \underbrace{\begin{bmatrix} 1.0 \\ \mathbf{x}^{(i)} \end{bmatrix}}_{\text{Vector (bias + data point)}}$$

Gradient is a weighted-sum of training data points

Gradient Descent (GD)

$$\mathbf{w}_+ \doteq \begin{bmatrix} w_0 & \mathbf{w} \end{bmatrix} \Rightarrow \nabla_{\mathbf{w}_+} \mathcal{L} = \mathbf{0}$$

$$\mathbf{w}_+^{(\ell)} = \mathbf{w}_+^{(\ell-1)} - (\alpha \nabla \mathcal{L}) \big|_{\mathbf{w}_+^{(\ell-1)}}$$

$$\nabla_{\mathbf{w}_+} \mathcal{L} = \sum_{i=1}^N \left(\sigma \left(\mathbf{w}_+ \begin{bmatrix} 1.0 \\ \mathbf{x}^{(i)} \end{bmatrix} \right) - y^{(i)} \right) \begin{bmatrix} 1.0 \\ \mathbf{x}^{(i)} \end{bmatrix}$$

Gradient Descent (GD)

$$\mathbf{w}_+ \doteq [w_0 \quad \mathbf{w}] \Rightarrow \nabla_{\mathbf{w}_+} \mathcal{L} = \mathbf{0}$$

the step size

$$\mathbf{w}_+^{(\ell)} = \mathbf{w}_+^{(\ell-1)} - (\alpha \nabla \mathcal{L}) \Big|_{\mathbf{w}_+^{(\ell-1)}}$$

$$\nabla_{\mathbf{w}_+} \mathcal{L} = \sum_{i=1}^N \left(\sigma \left(\mathbf{w}_+ \begin{bmatrix} 1.0 \\ \mathbf{x}^{(i)} \end{bmatrix} \right) - y^{(i)} \right) \begin{bmatrix} 1.0 \\ \mathbf{x}^{(i)} \end{bmatrix}$$

Gradient Descent (GD)

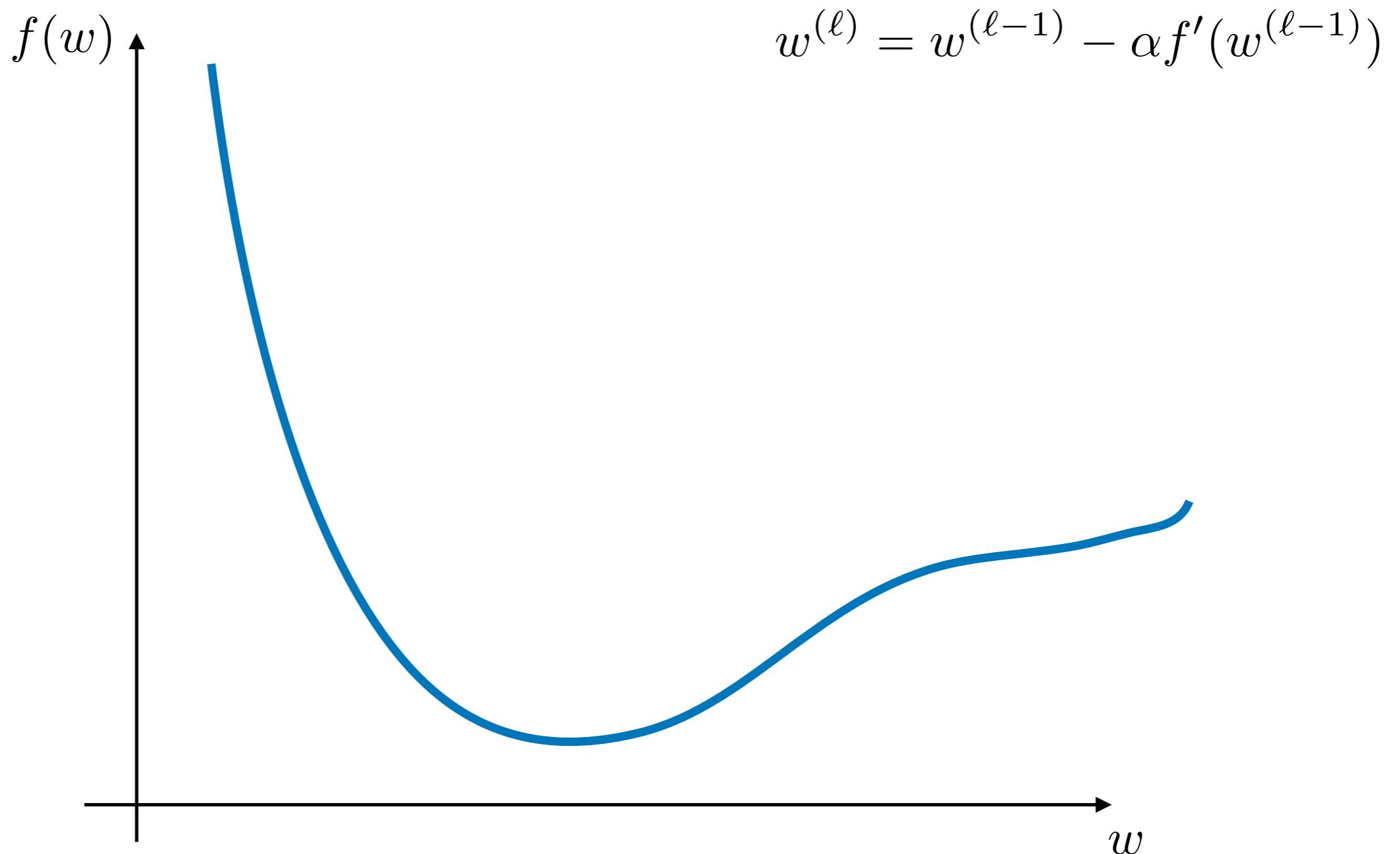
$$\mathbf{w}_+ \doteq [w_0 \quad \mathbf{w}] \Rightarrow \nabla_{\mathbf{w}_+} \mathcal{L} = \mathbf{0}$$

$$\mathbf{w}_+^{(\ell)} = \mathbf{w}_+^{(\ell-1)} - (\alpha \nabla \mathcal{L}) \Big|_{\mathbf{w}_+^{(\ell-1)}}$$

$$\nabla_{\mathbf{w}_+} \mathcal{L} = \sum_{i=1}^N \left(\sigma \left(\mathbf{w}_+ \begin{bmatrix} 1.0 \\ \mathbf{x}^{(i)} \end{bmatrix} \right) - y^{(i)} \right) \begin{bmatrix} 1.0 \\ \mathbf{x}^{(i)} \end{bmatrix}$$

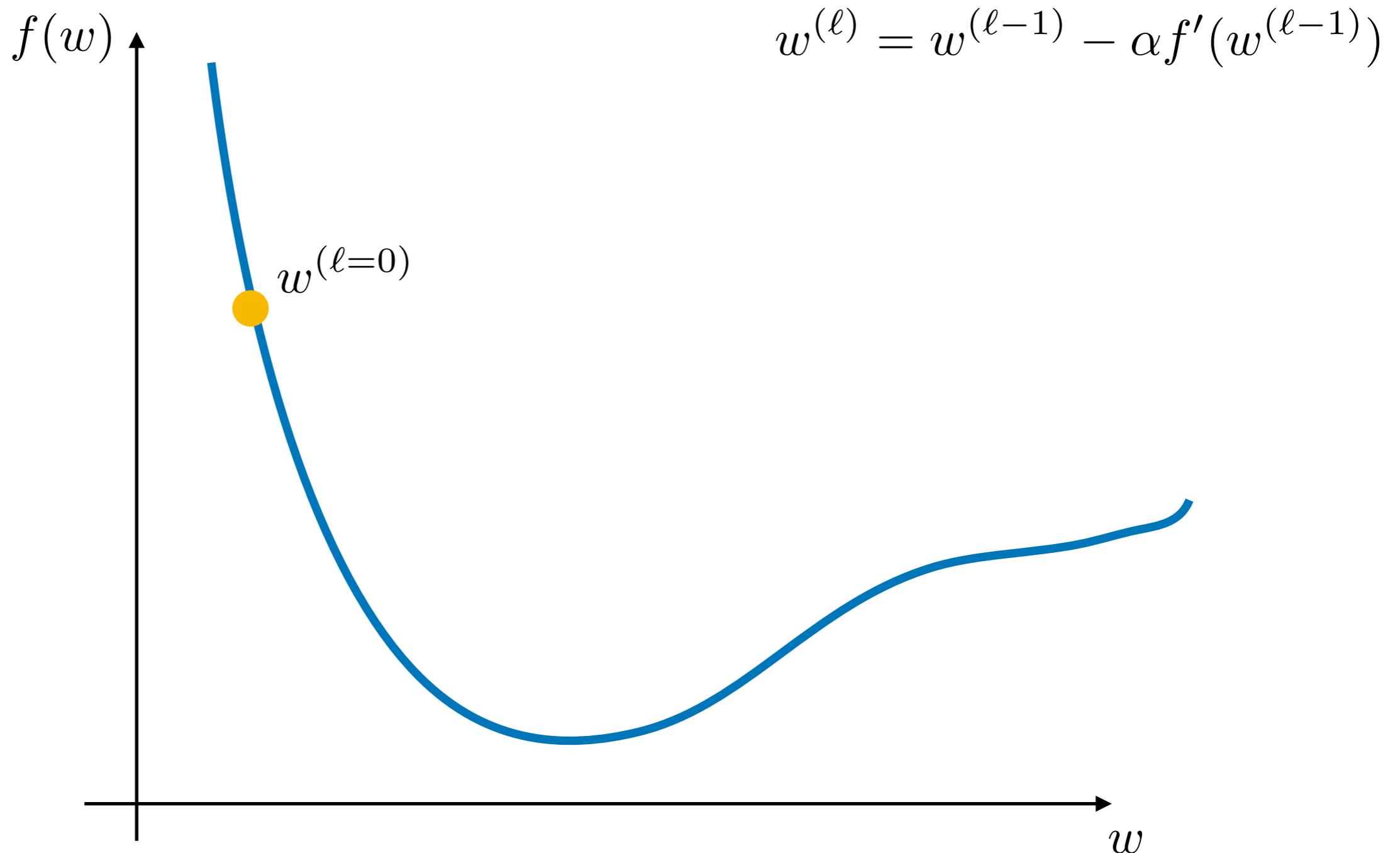
Choosing the step size manually is a **hard problem!**

Gradient Descent



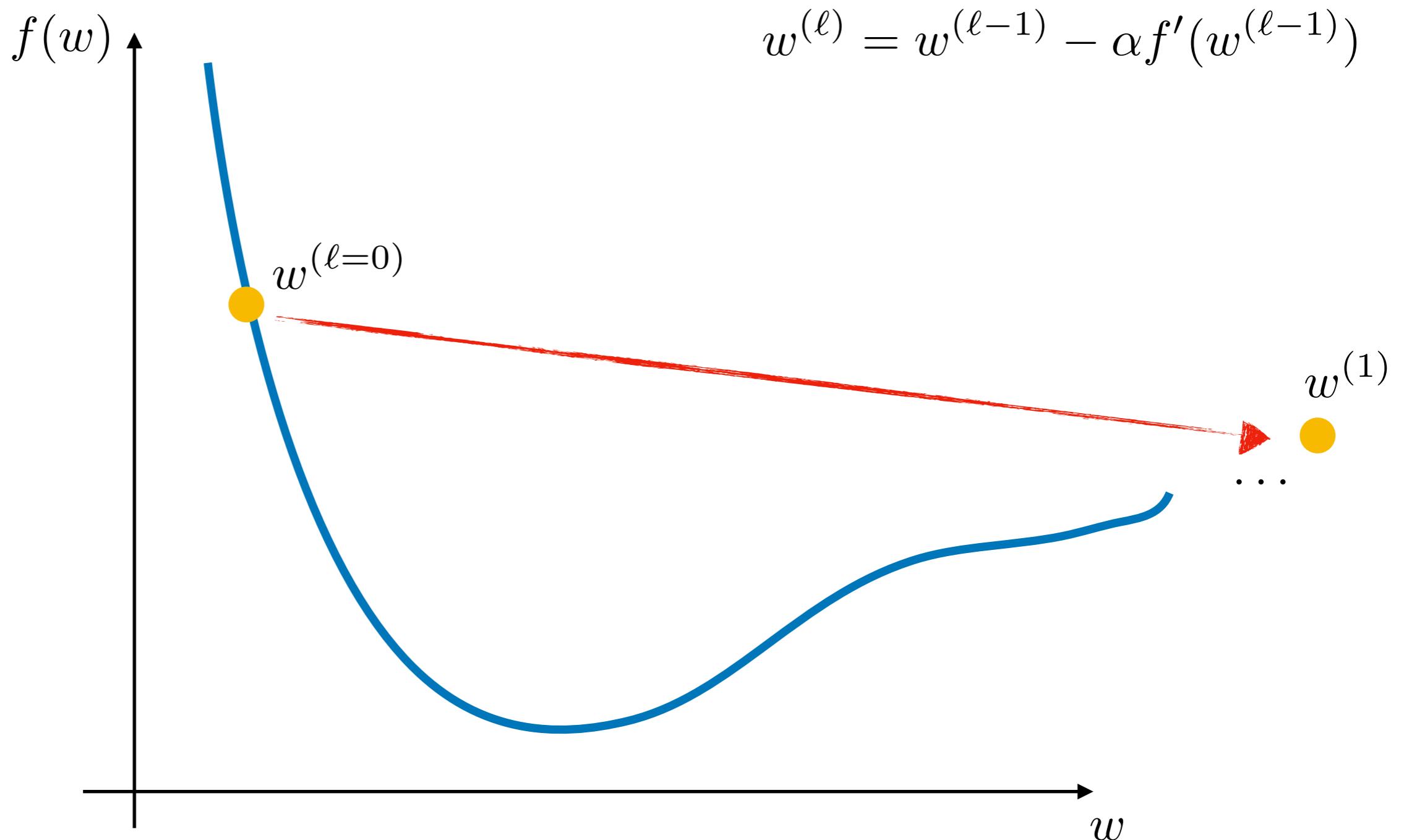
If α is too large

Gradient Descent



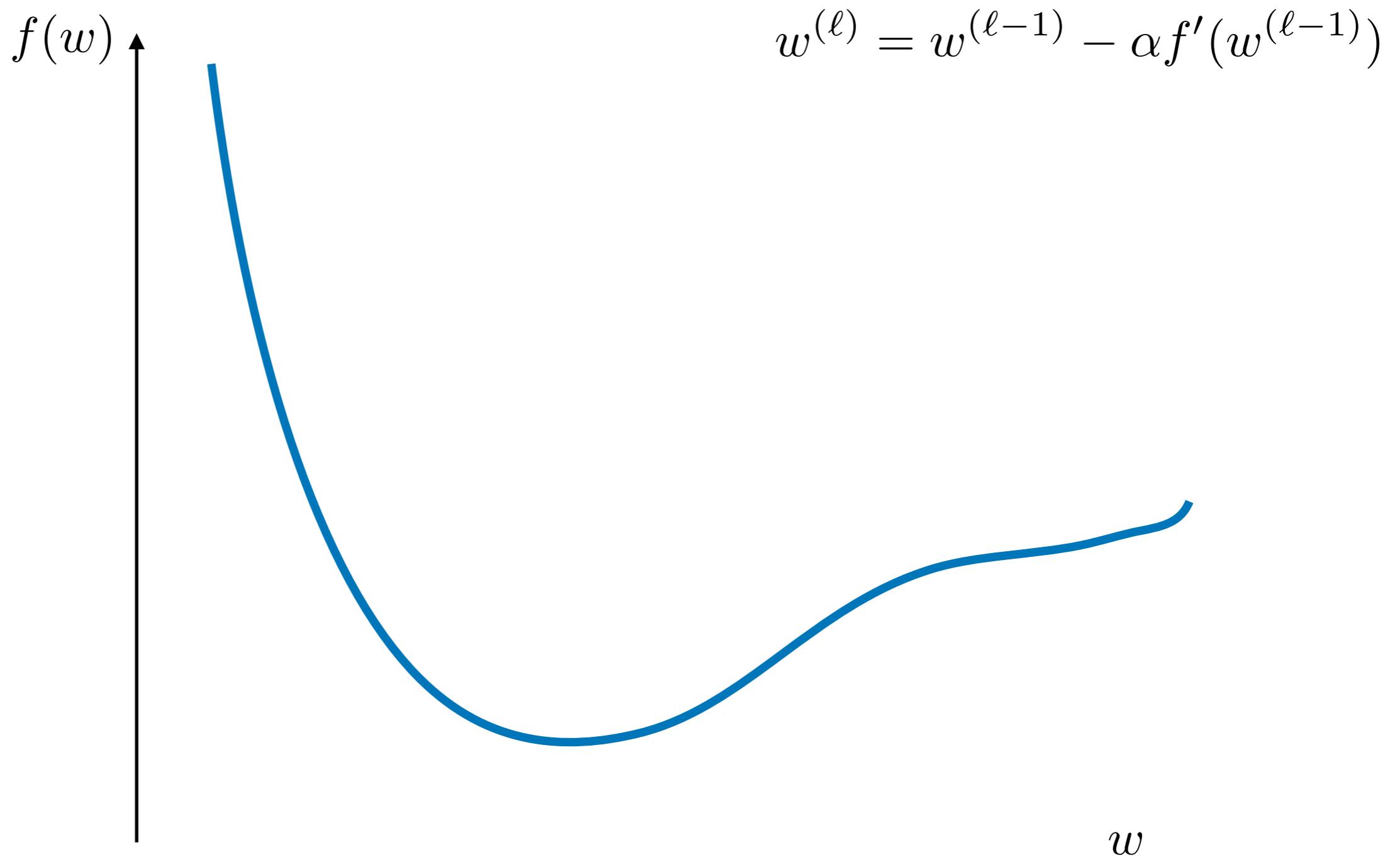
If α is too large

Gradient Descent



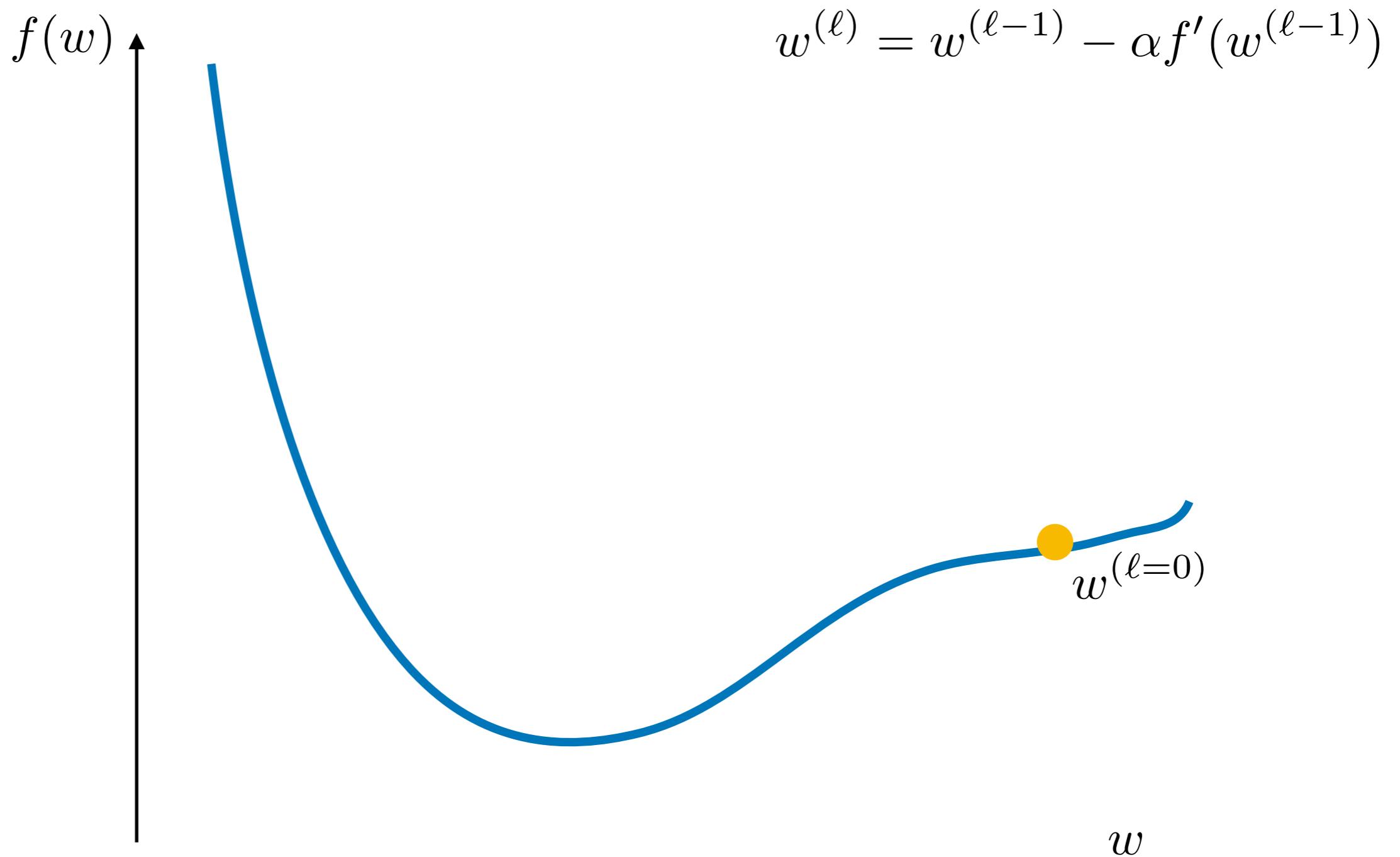
If α is too large

Gradient Descent



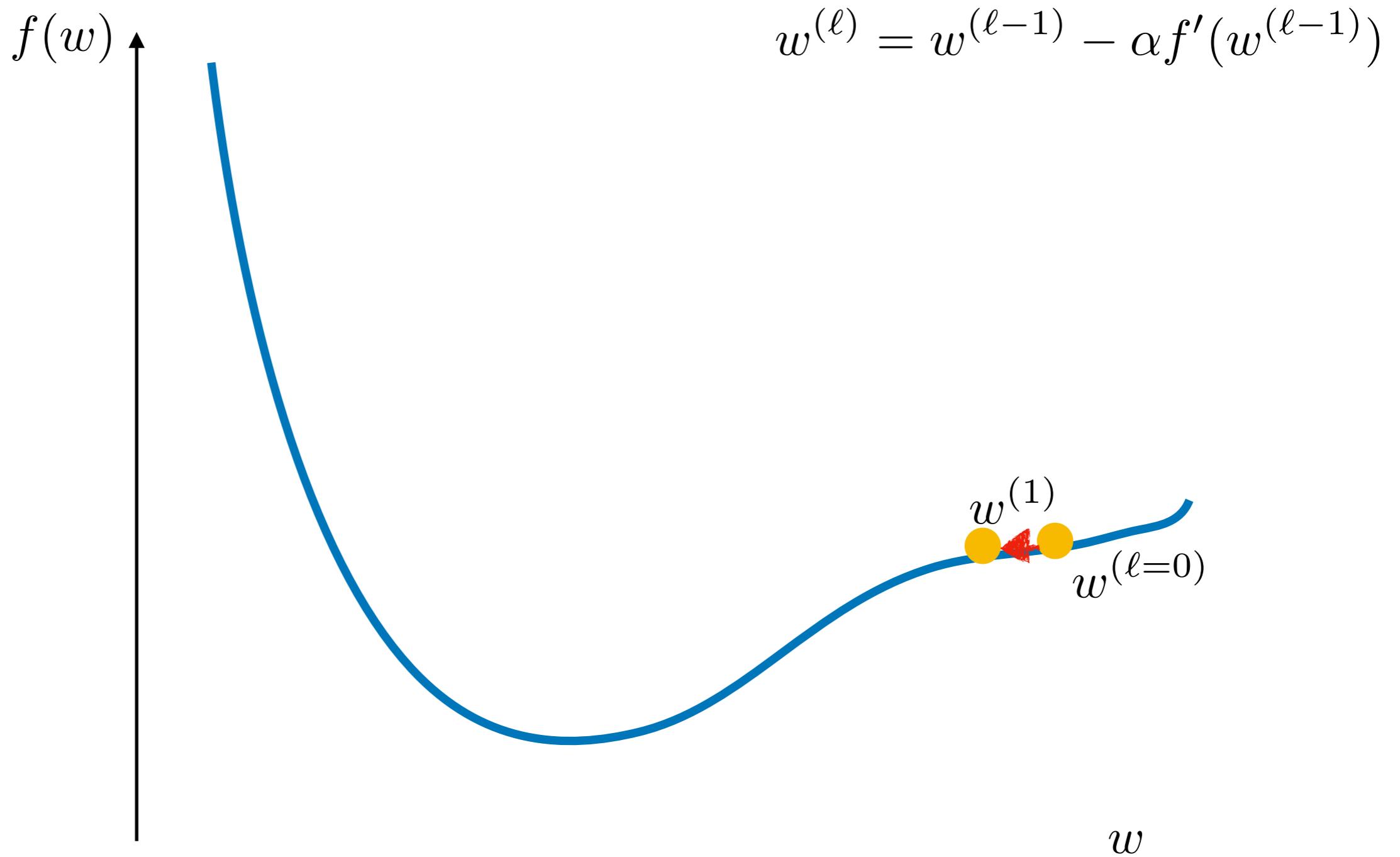
If α is too small

Gradient Descent



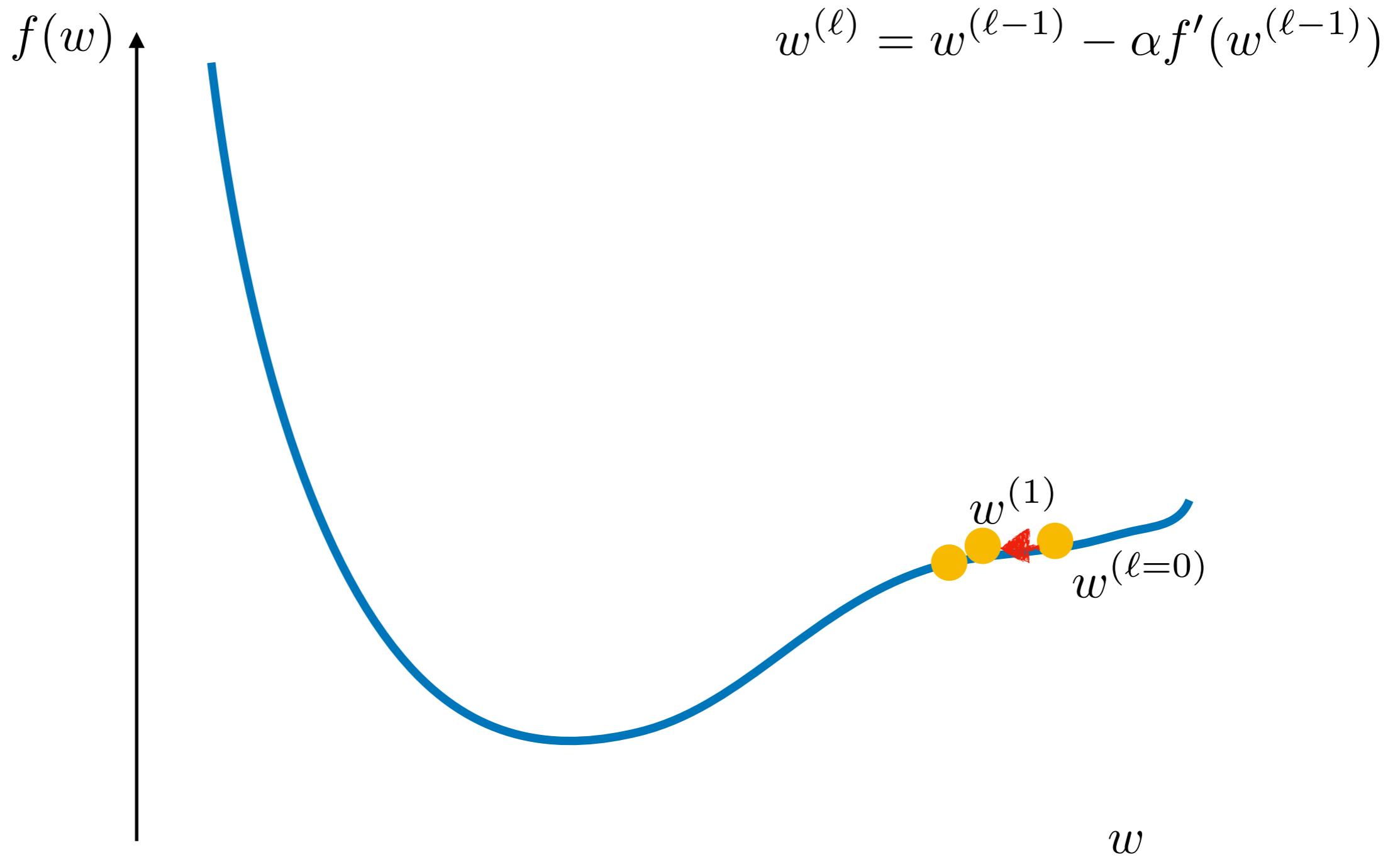
If α is too small

Gradient Descent



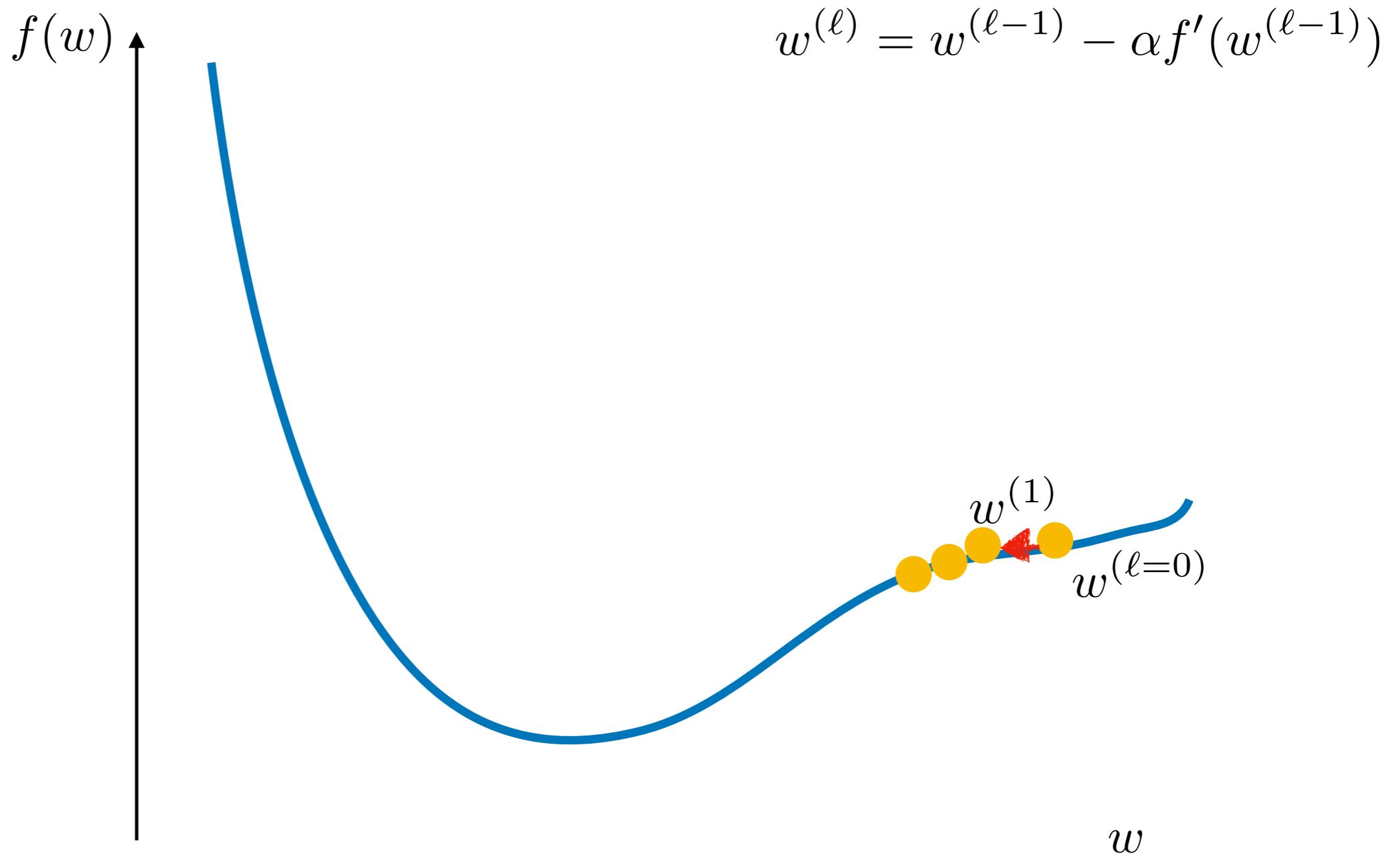
If α is too small

Gradient Descent



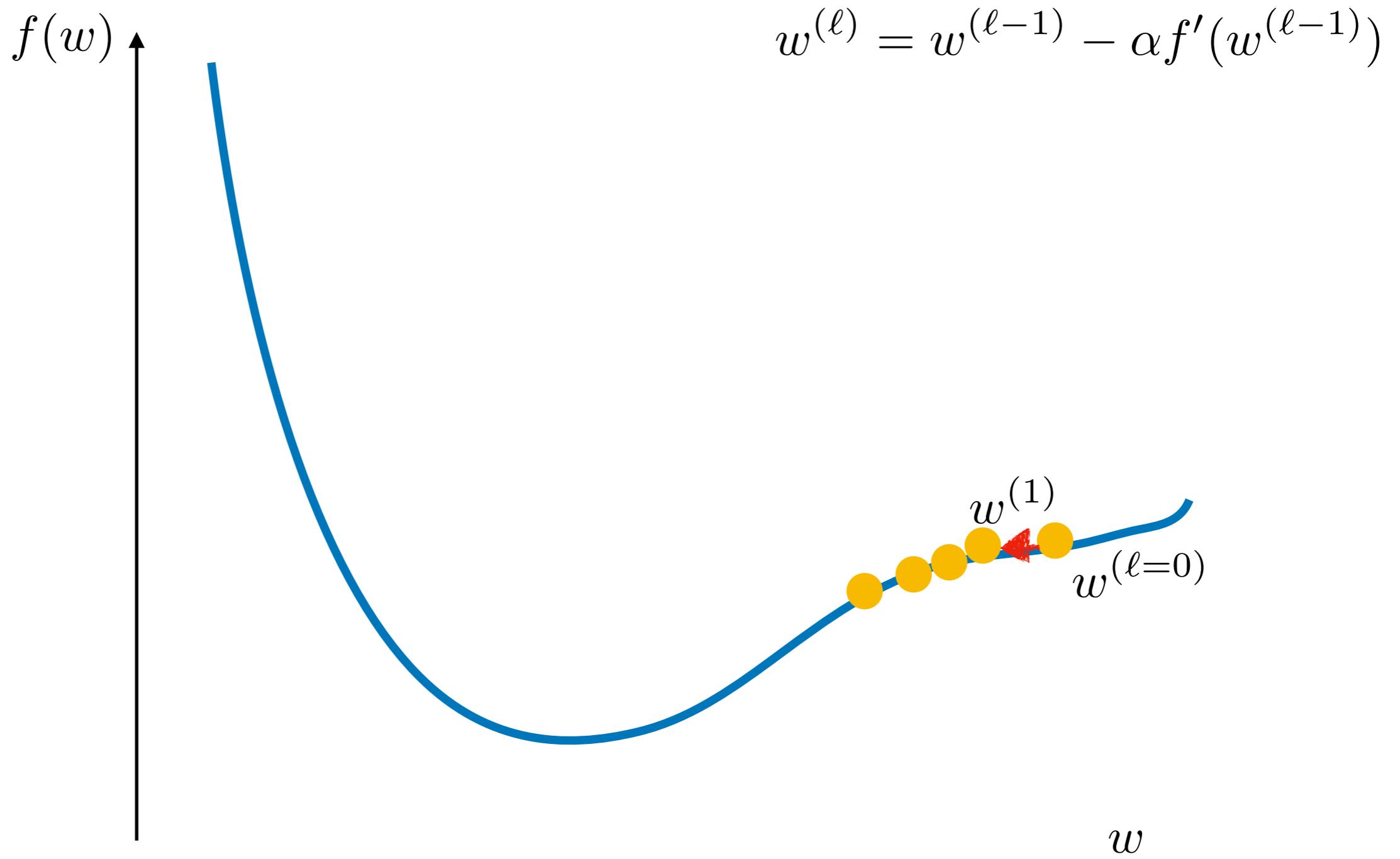
If α is too small

Gradient Descent



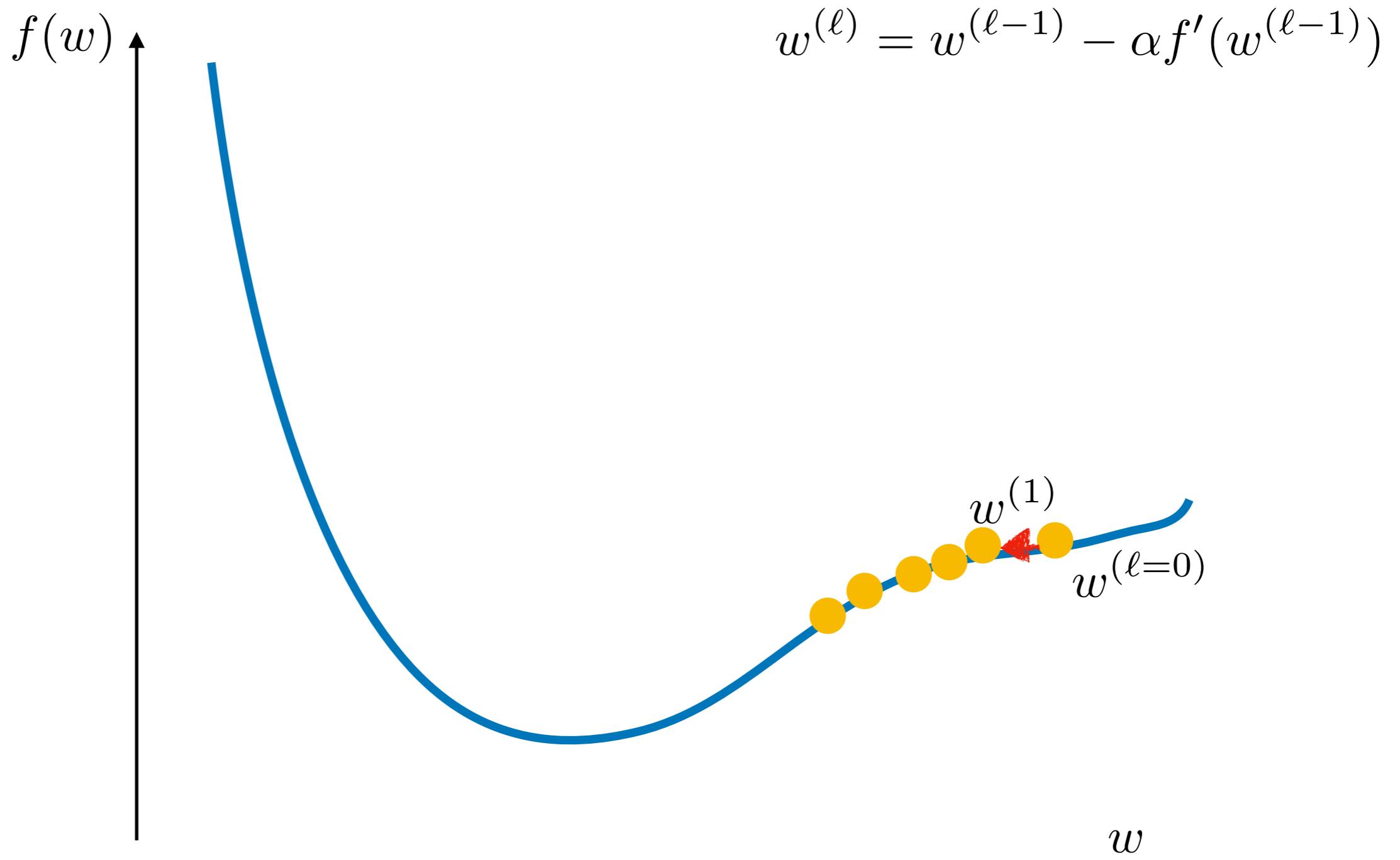
If α is too small

Gradient Descent



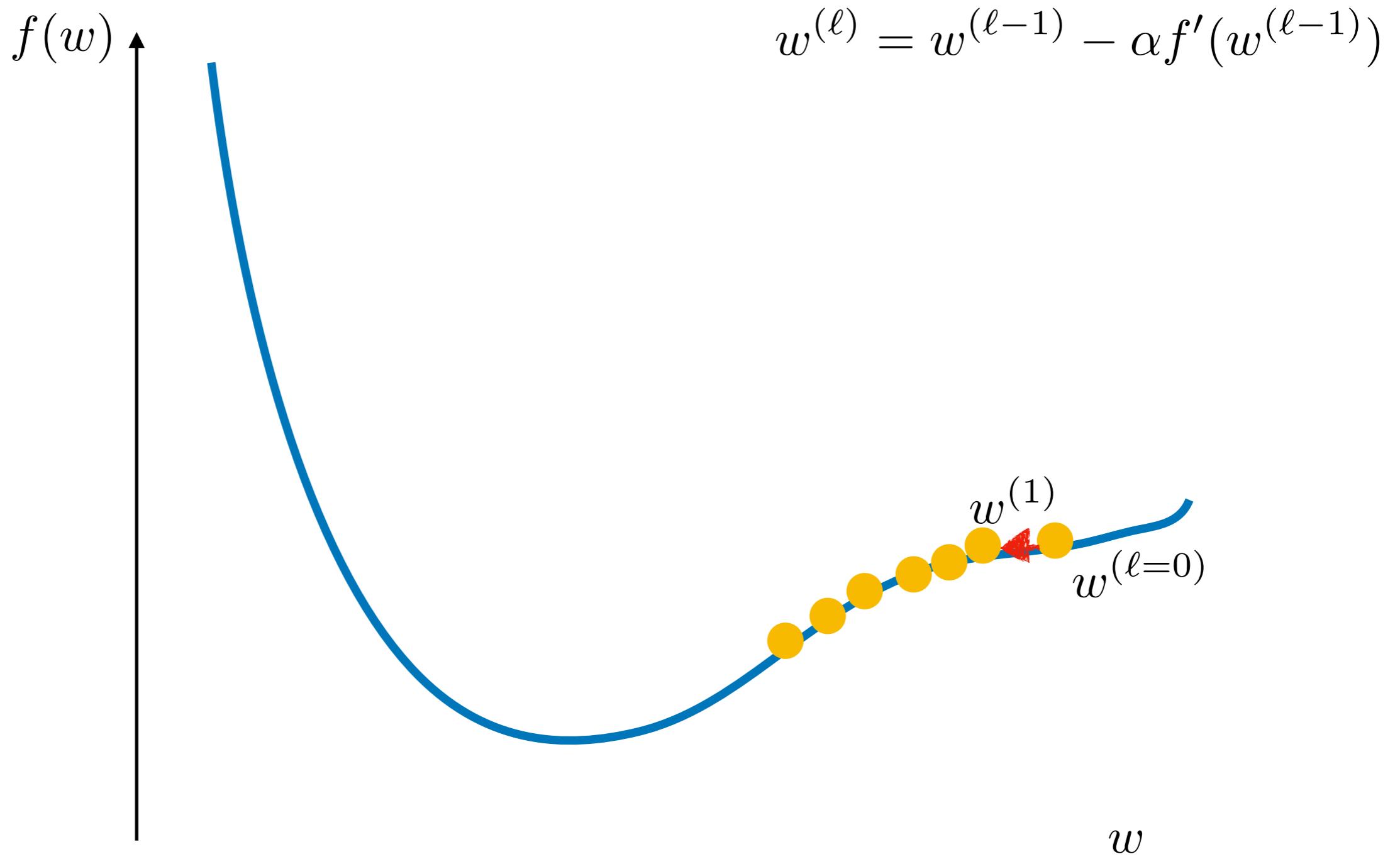
If α is too small

Gradient Descent



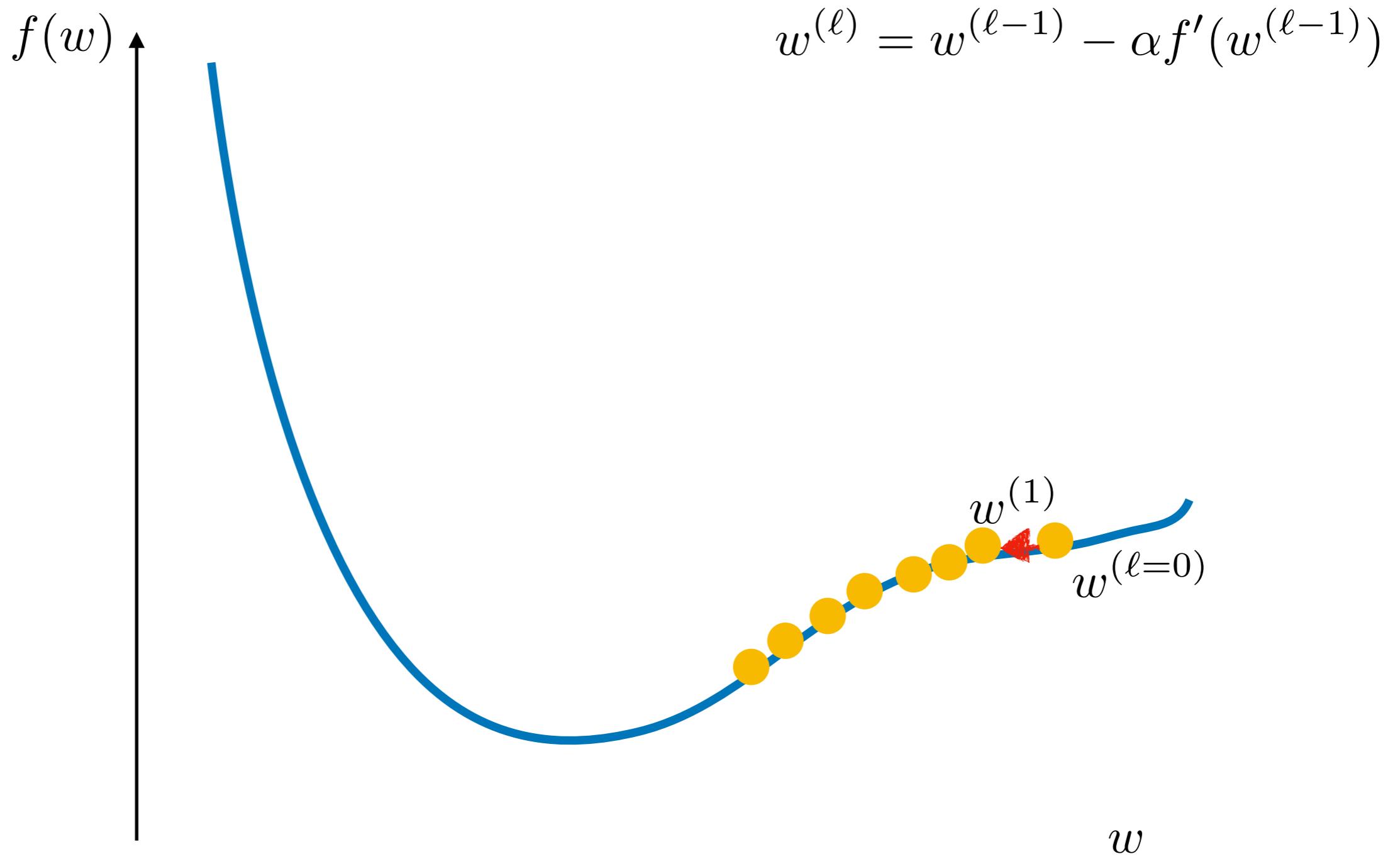
If α is too small

Gradient Descent



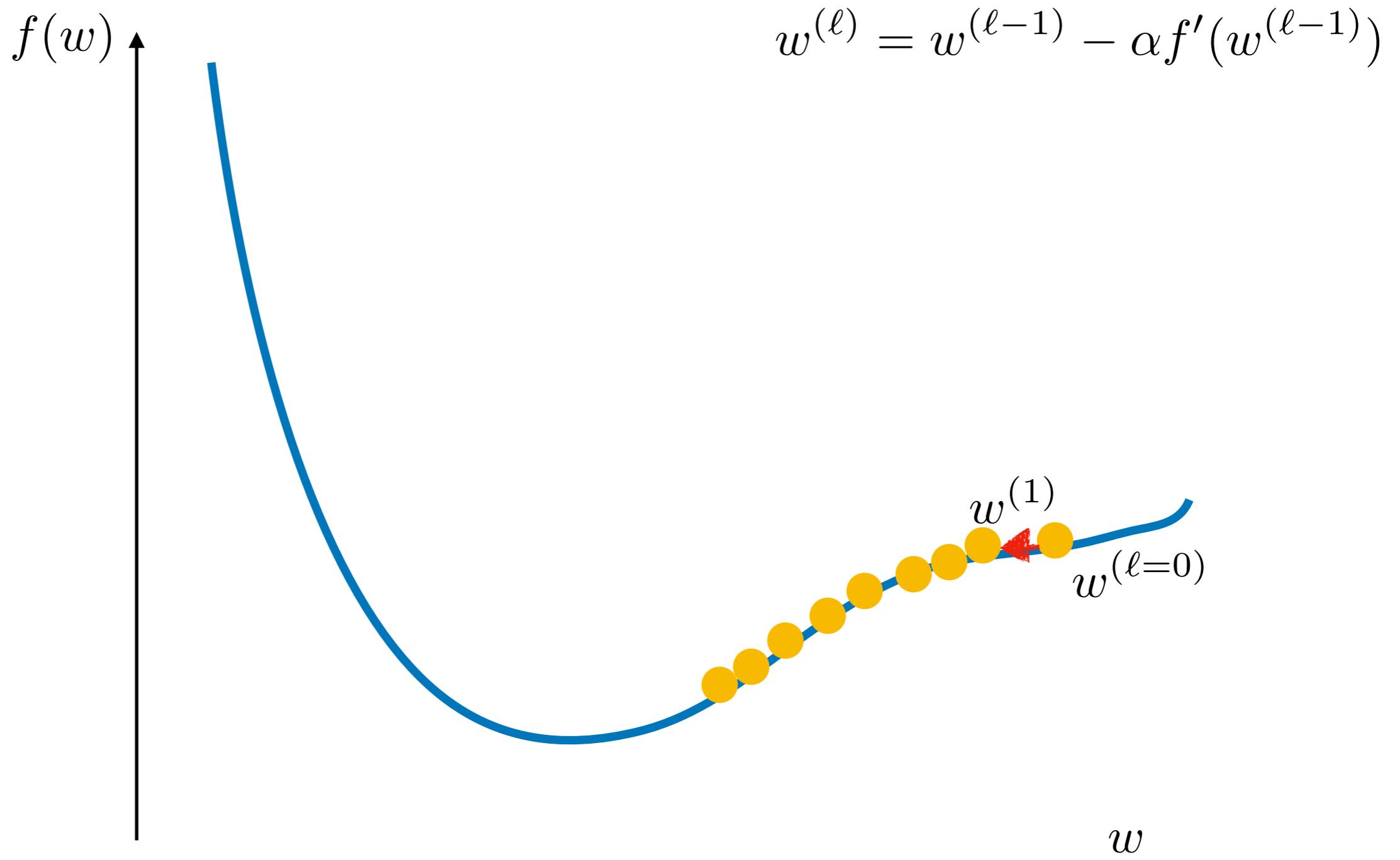
If α is too small

Gradient Descent



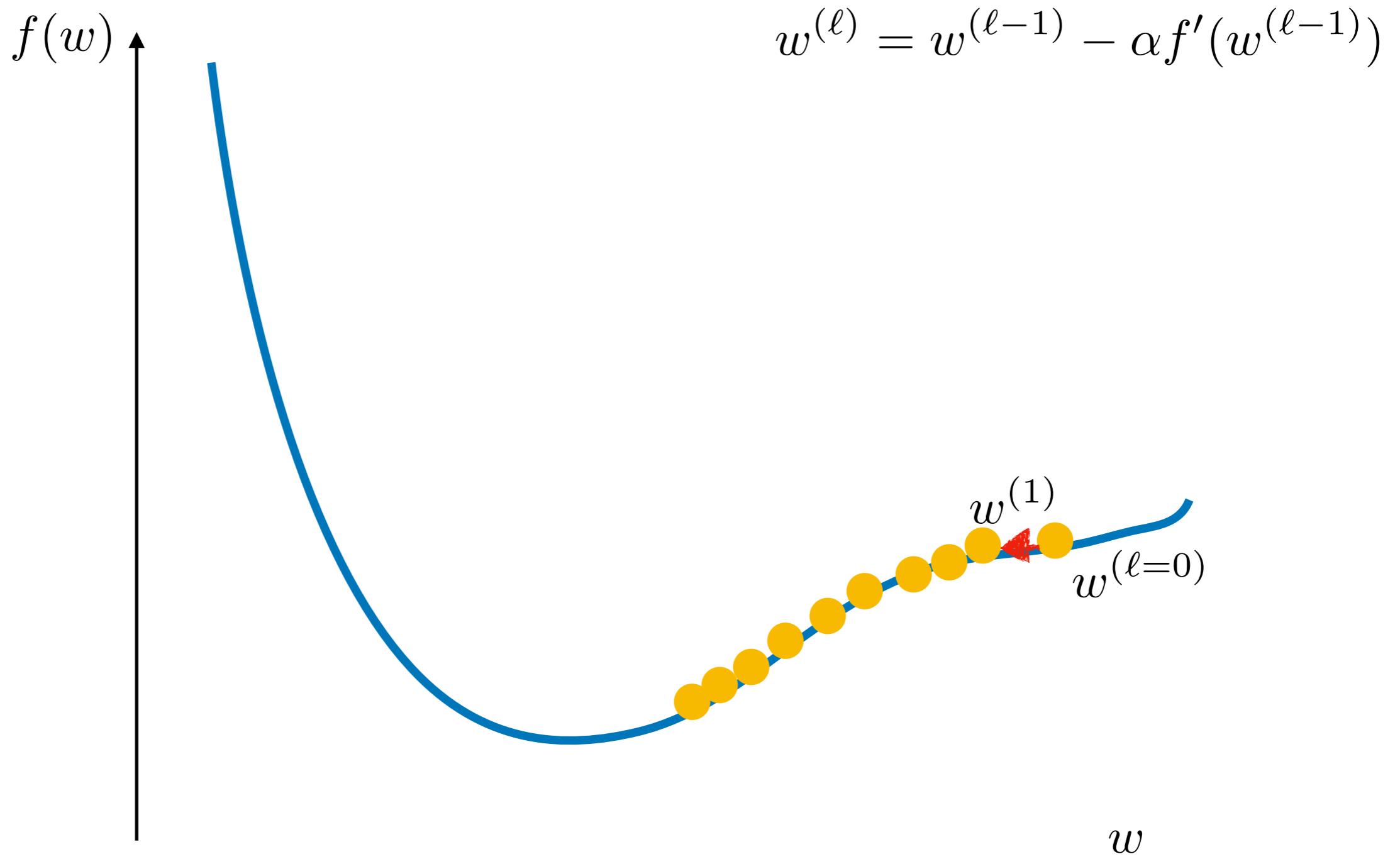
If α is too small

Gradient Descent



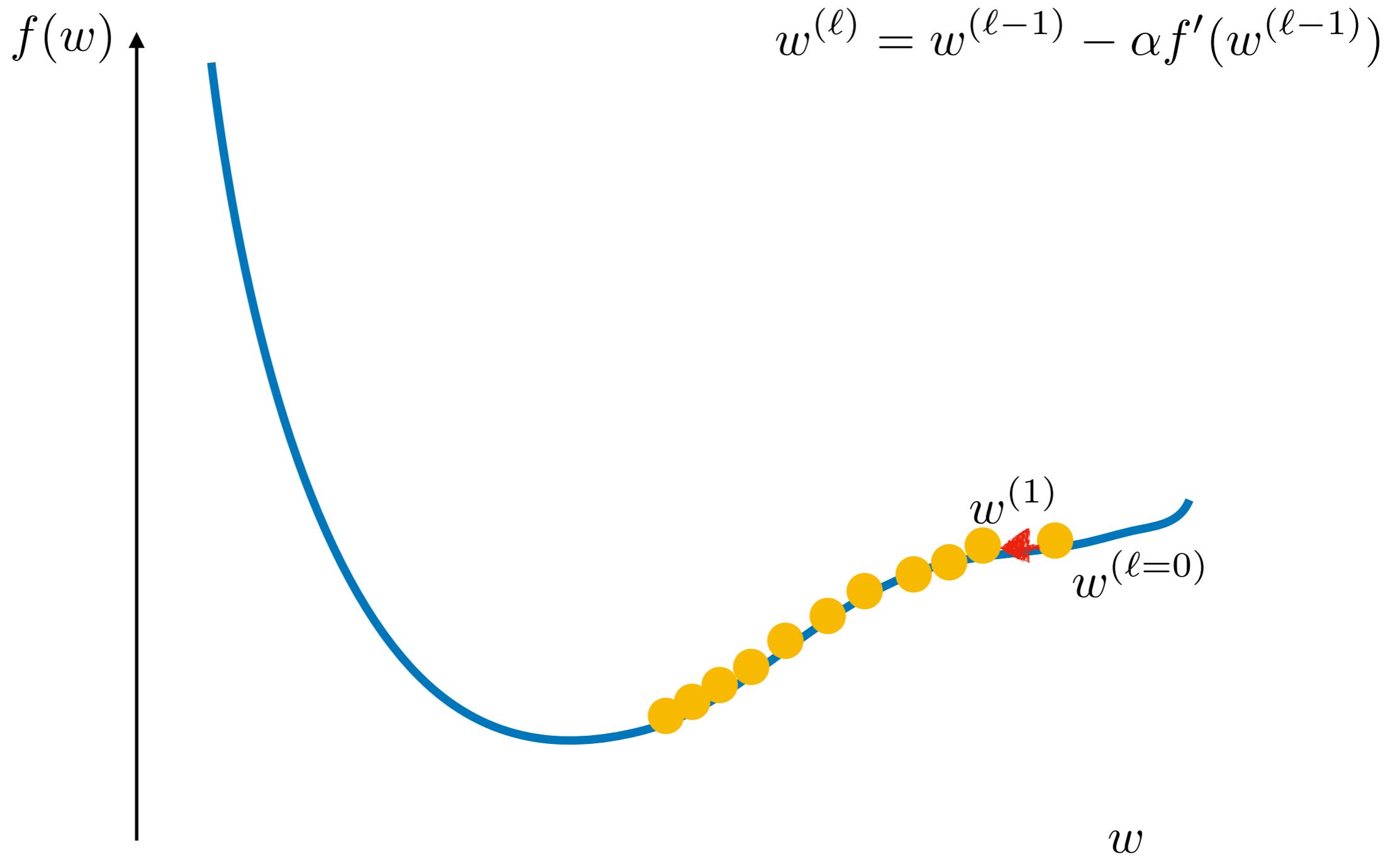
If α is too small

Gradient Descent



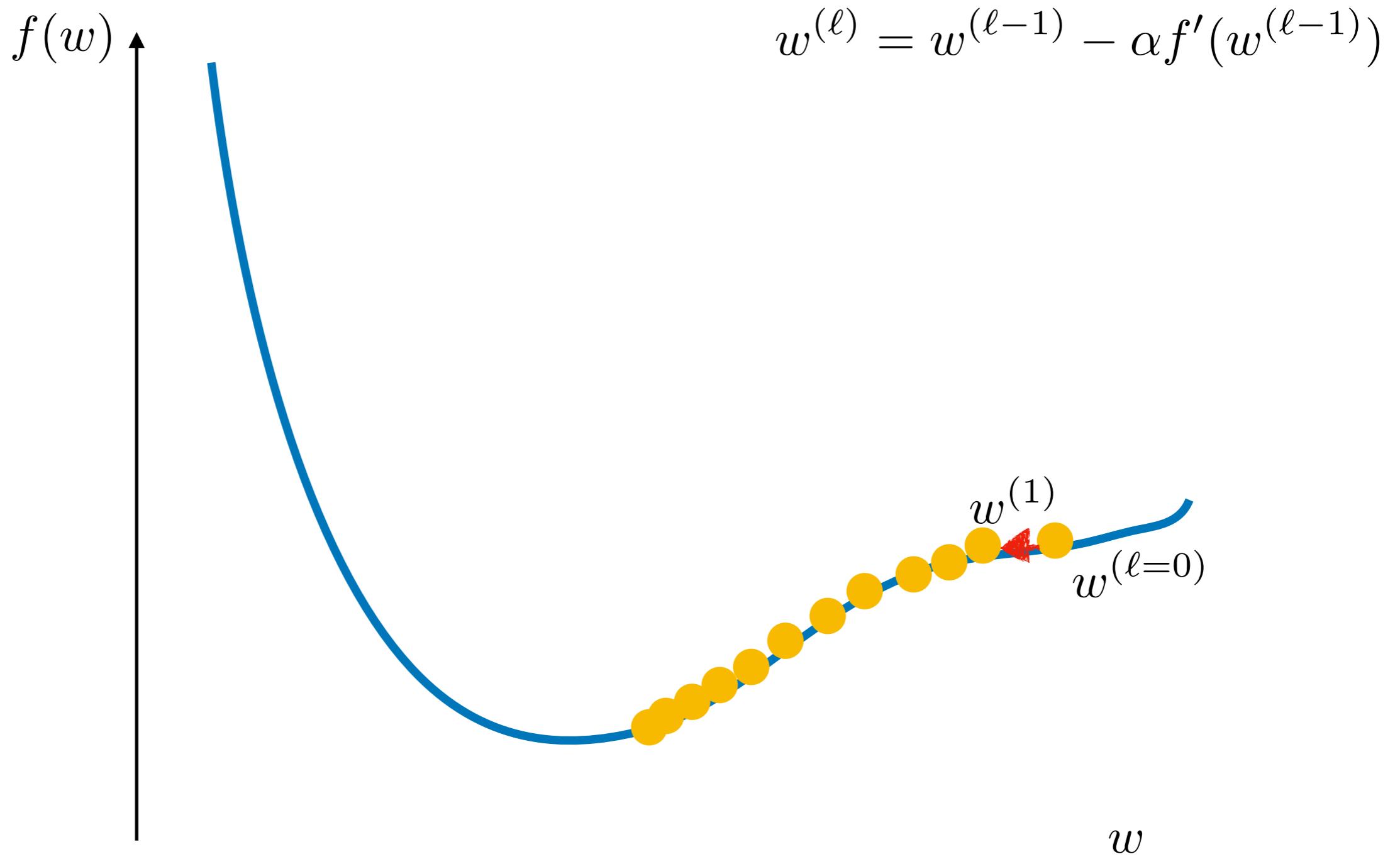
If α is too small

Gradient Descent



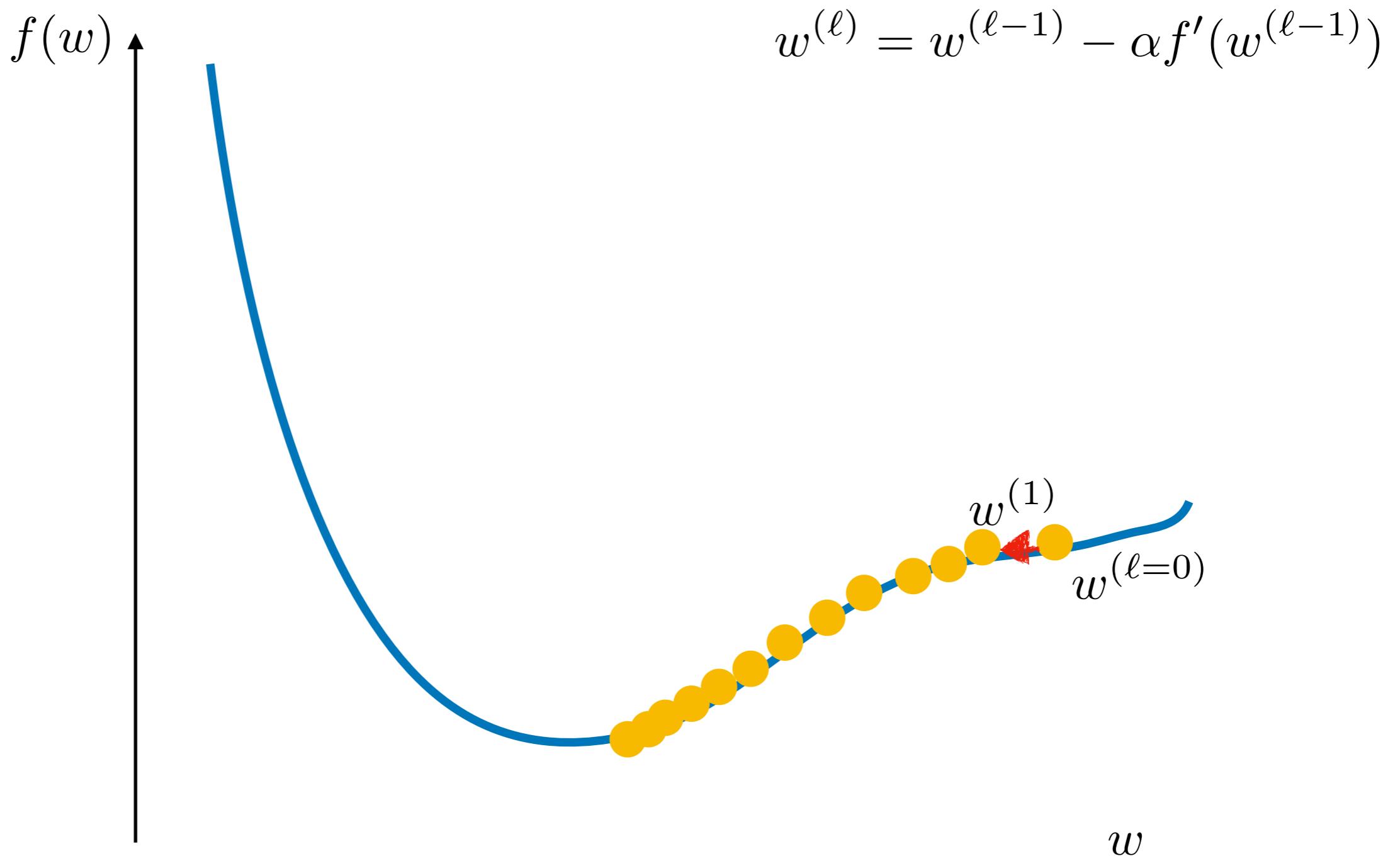
If α is too small

Gradient Descent



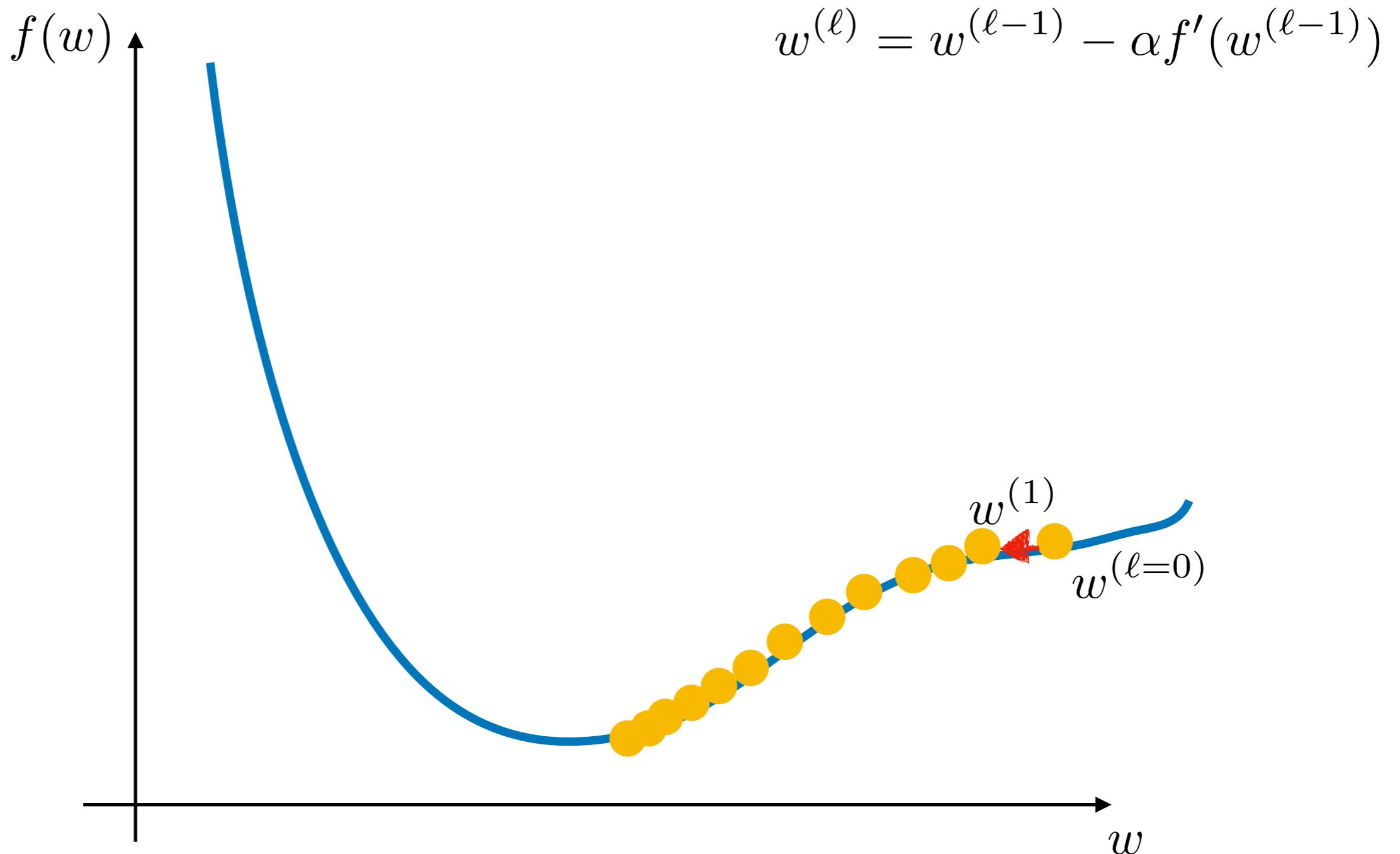
If α is too small

Gradient Descent



If α is too small

Gradient Descent



If α is too small

Line Search methods: automatically running the step size

$$\mathbf{w}_+^{(\ell)} = \mathbf{w}_+^{(\ell-1)} - \left(\alpha^{(\ell)} \nabla \mathcal{L} \right) |_{\mathbf{w}_+^{(\ell-1)}}$$

- Estimate the largest possible step size that guarantees that the function decreases
- Every state-of-the-art Deep Learning library contains implementations of various algorithms to optimize gradient descent
- Momentum, Adam, Adagrad, ...

Check out in Aula Global two excellent posts on SGD methods for Deep Learning

([Link1](#), [Link2](#))

Logistic Regression: Stochastic gradient descent

Stochastic Optimization (Mini-batch Optimization)

$$\mathbf{w}_+ \doteq \begin{bmatrix} w_0 & \mathbf{w} \end{bmatrix} \Rightarrow \nabla_{\mathbf{w}_+} \mathcal{L} = \mathbf{0}$$

$$\mathbf{w}_+^{(\ell)} = \mathbf{w}_+^{(\ell-1)} - (\alpha \nabla \mathcal{L}) \big|_{\mathbf{w}_+^{(\ell-1)}}$$

$$\nabla_{\mathbf{w}_+} \mathcal{L} = \sum_{i=1}^N \left(\sigma \left(\mathbf{w}_+ \begin{bmatrix} 1.0 \\ \mathbf{x}^{(i)} \end{bmatrix} \right) - y^{(i)} \right) \begin{bmatrix} 1.0 \\ \mathbf{x}^{(i)} \end{bmatrix}$$

Stochastic Optimization (Mini-batch Optimization)

$$\mathbf{w}_+ \doteq [w_0 \quad \mathbf{w}] \Rightarrow \nabla_{\mathbf{w}_+} \mathcal{L} = \mathbf{0}$$

$$\mathbf{w}_+^{(\ell)} = \mathbf{w}_+^{(\ell-1)} - (\alpha \nabla \mathcal{L}) \Big|_{\mathbf{w}_+^{(\ell-1)}}$$

Expensive for very large databases!!

$$\nabla_{\mathbf{w}_+} \mathcal{L} = \sum_{i=1}^N \left(\sigma \left(\mathbf{w}_+ \begin{bmatrix} 1.0 \\ \mathbf{x}^{(i)} \end{bmatrix} \right) - y^{(i)} \right) \begin{bmatrix} 1.0 \\ \mathbf{x}^{(i)} \end{bmatrix}$$

Stochastic Optimization (Mini-batch Optimization)

$$\mathbf{w}_+ \doteq [w_0 \quad \mathbf{w}] \Rightarrow \nabla_{\mathbf{w}_+} \mathcal{L} = \mathbf{0}$$

$$\mathbf{w}_+^{(\ell)} = \mathbf{w}_+^{(\ell-1)} - (\alpha \nabla \mathcal{L}) \Big|_{\mathbf{w}_+^{(\ell-1)}}$$

Expensive for very large databases!!

$$\nabla_{\mathbf{w}_+} \mathcal{L} = \sum_{i=1}^N \left(\sigma \left(\mathbf{w}_+ \begin{bmatrix} 1.0 \\ \mathbf{x}^{(i)} \end{bmatrix} \right) - y^{(i)} \right) \begin{bmatrix} 1.0 \\ \mathbf{x}^{(i)} \end{bmatrix}$$

Mini-batch optimization

Select at random a mini-batch \mathcal{B} of data at every SGD iteration

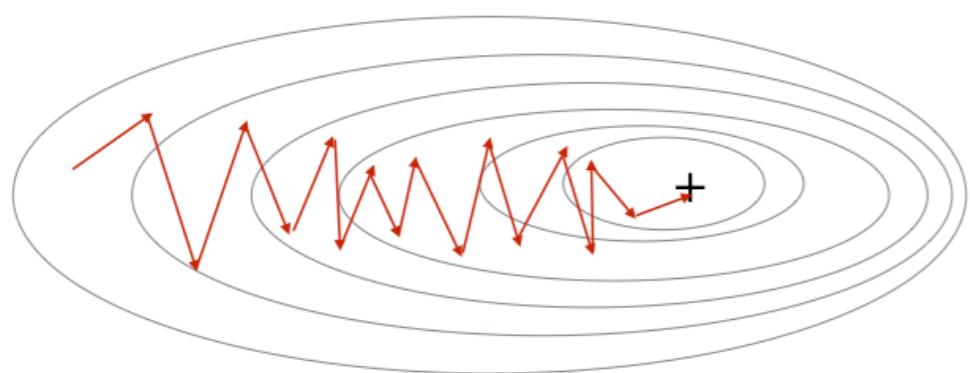
Stochastic Optimization (Mini-batch Optimization)

$$\mathbf{w}_+ \doteq [w_0 \quad \mathbf{w}] \Rightarrow \nabla_{\mathbf{w}_+} \mathcal{L} = \mathbf{0}$$

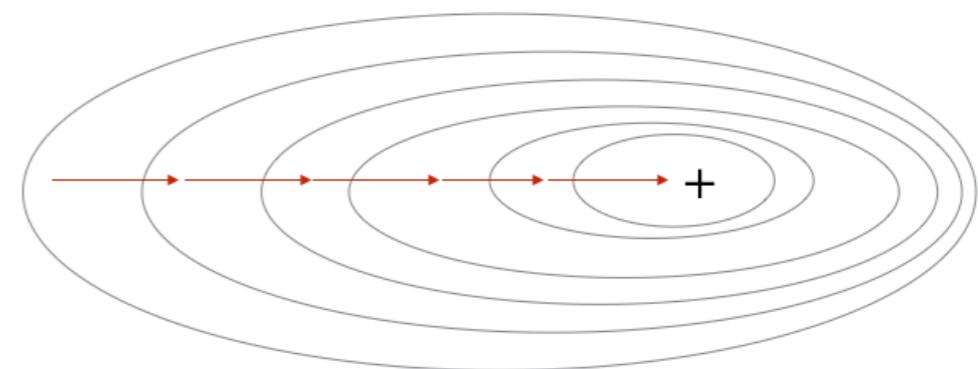
$$\mathbf{w}_+^{(\ell)} = \mathbf{w}_+^{(\ell-1)} - (\alpha \nabla \mathcal{L}) \big|_{\mathbf{w}_+^{(\ell-1)}}$$

$$\nabla_{\mathbf{w}_+} \mathcal{L} \approx \sum_{i \in \mathcal{B}} \left(\sigma \left(\mathbf{w}_+ \begin{bmatrix} 1.0 \\ \mathbf{x}^{(i)} \end{bmatrix} \right) - y^{(i)} \right) \begin{bmatrix} 1.0 \\ \mathbf{x}^{(i)} \end{bmatrix}$$

Stochastic Gradient Descent



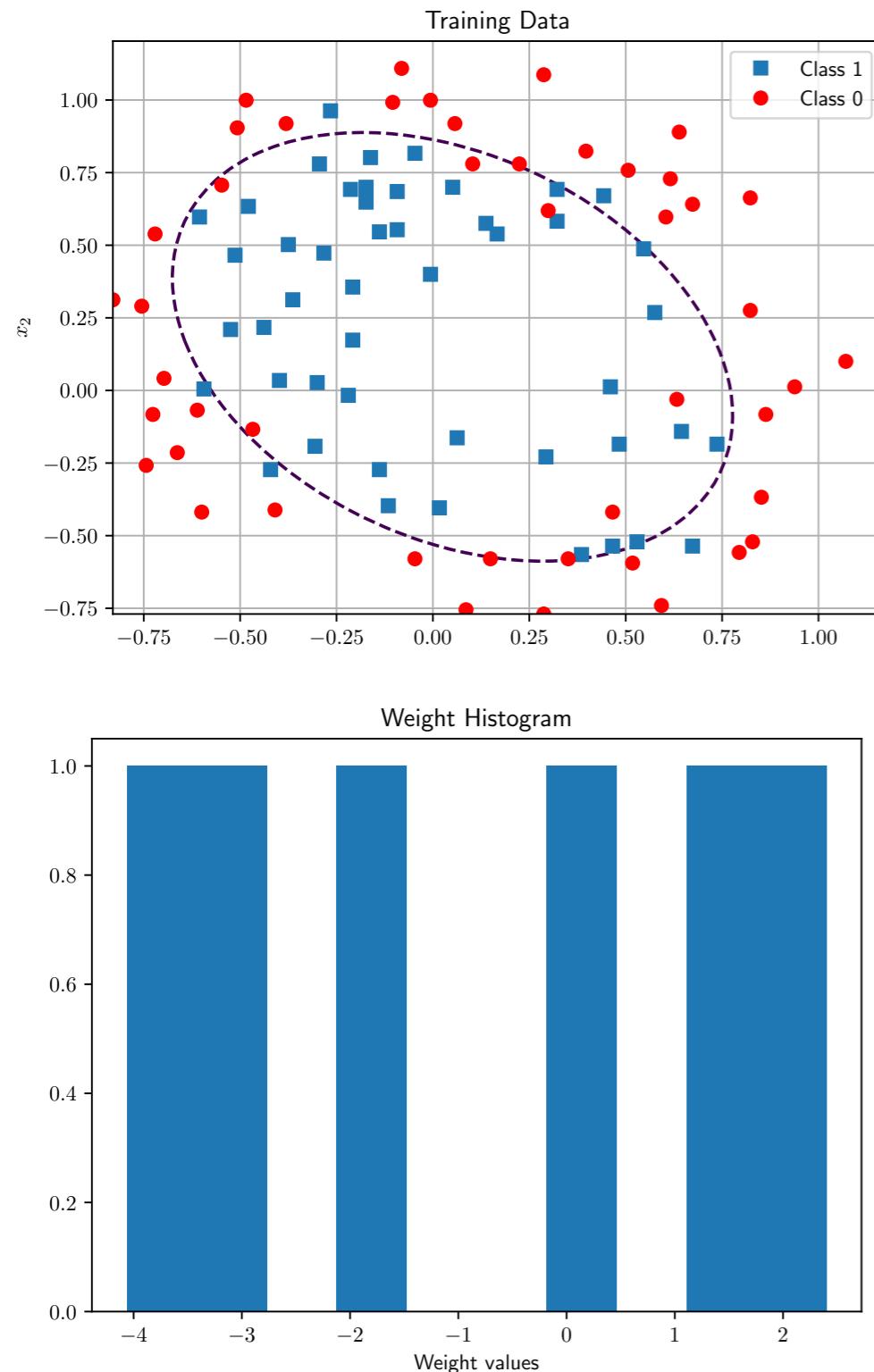
Gradient Descent



Source: [this post](#)

Logistic Regression: regularizing the weight norm

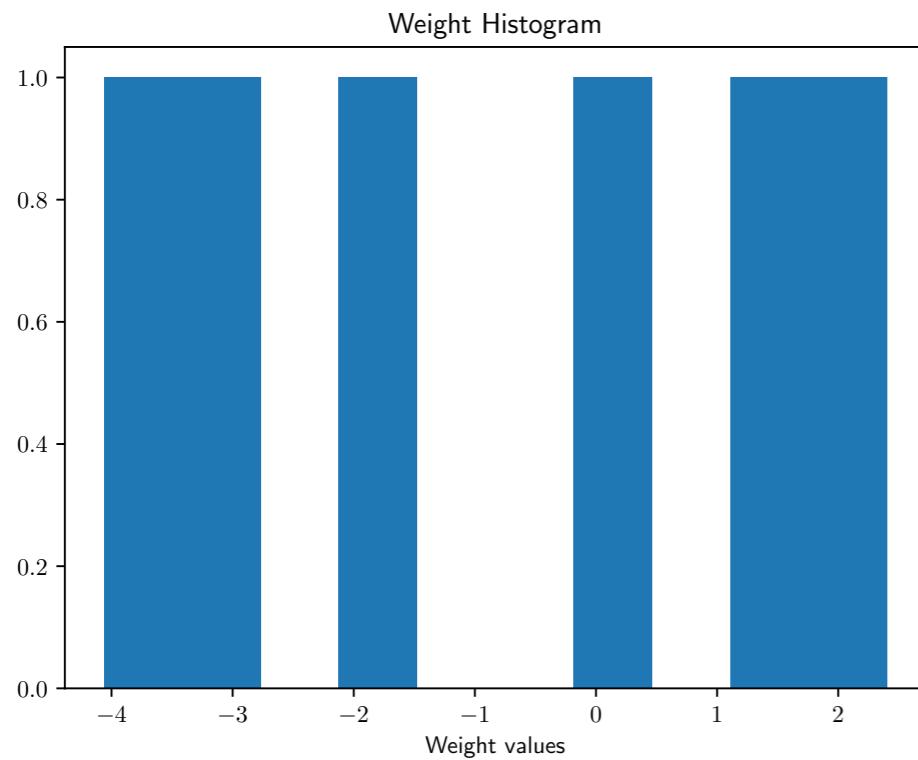
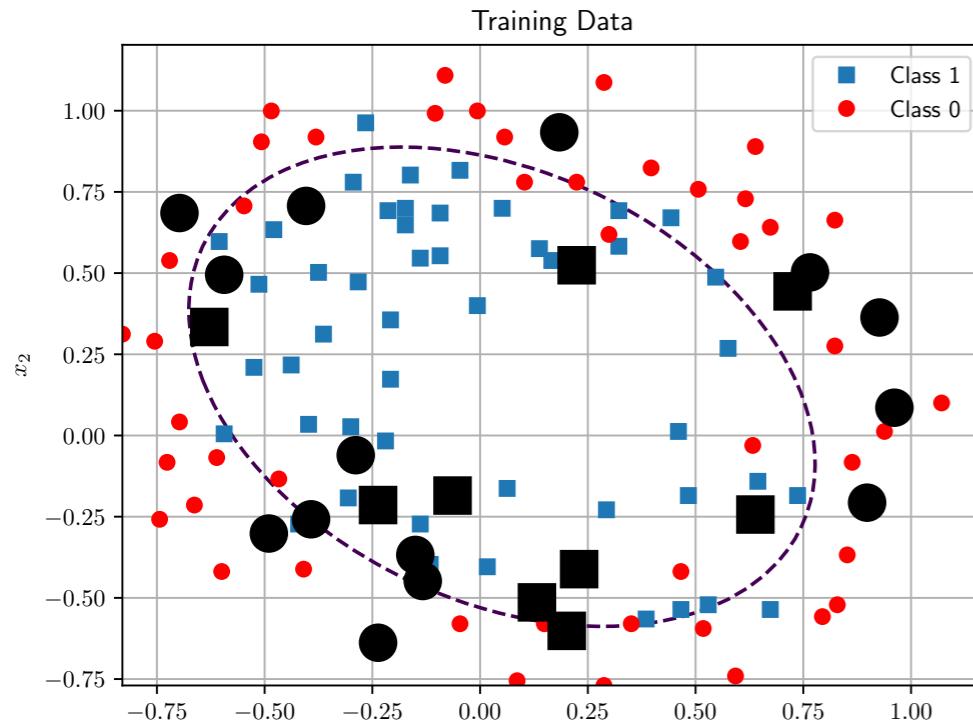
Binary Logistic Regression. Adding new features



$$\phi(\mathbf{x}) = [x_1, x_2, x_1^2, x_2^2, x_1 x_2]$$

$\phi(\mathbf{x})$ All polynomial terms up to degree 6

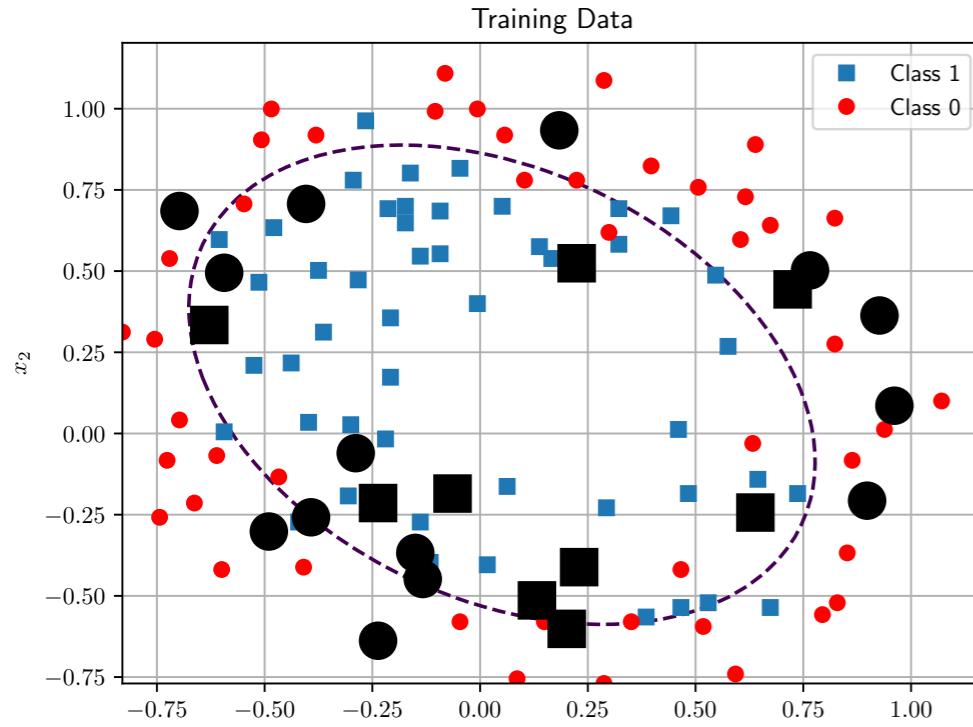
Binary Logistic Regression. Adding new features



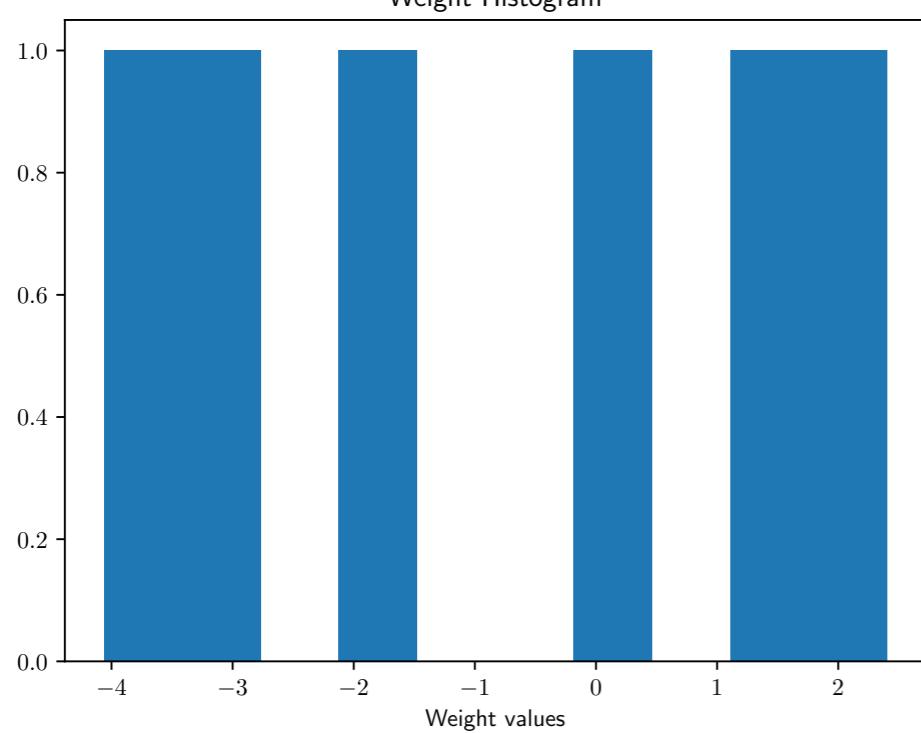
$$\phi(\mathbf{x}) = [x_1, x_2, x_1^2, x_2^2, x_1 x_2]$$

$\phi(\mathbf{x})$ All polynomial terms up to degree 6

Binary Logistic Regression. Adding new features



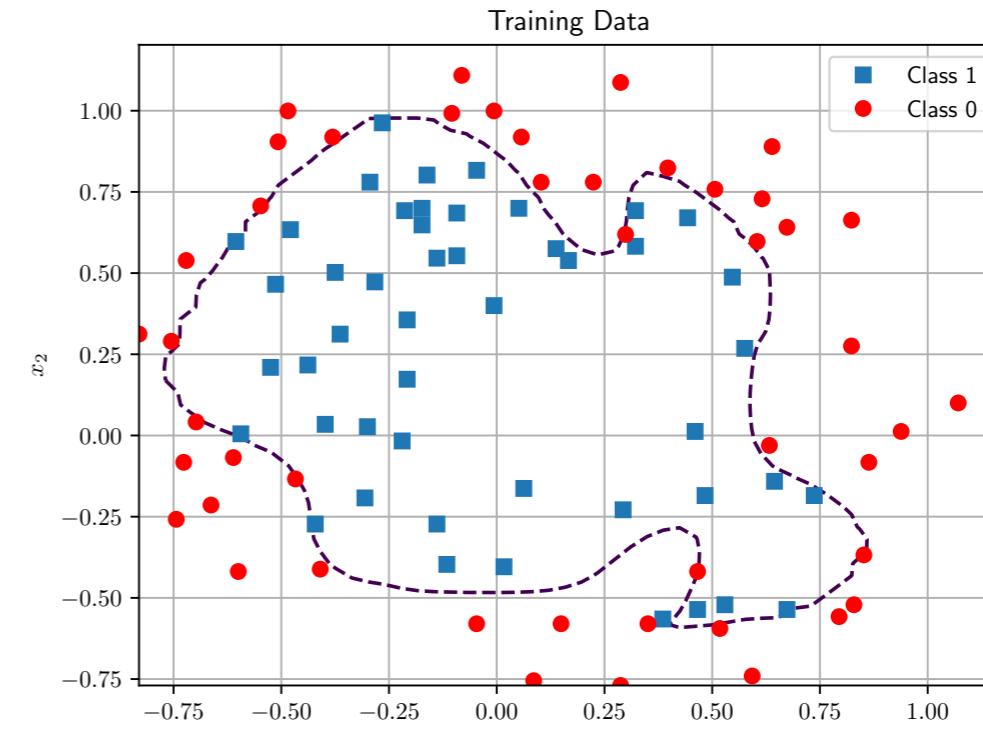
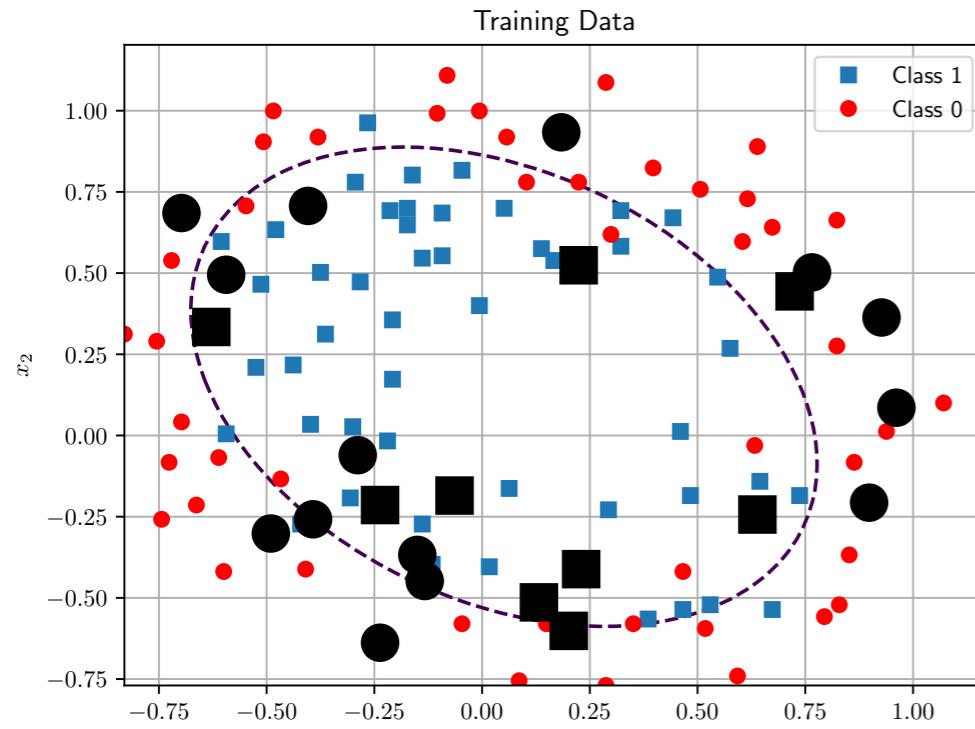
Test accuracy = 65.2 %



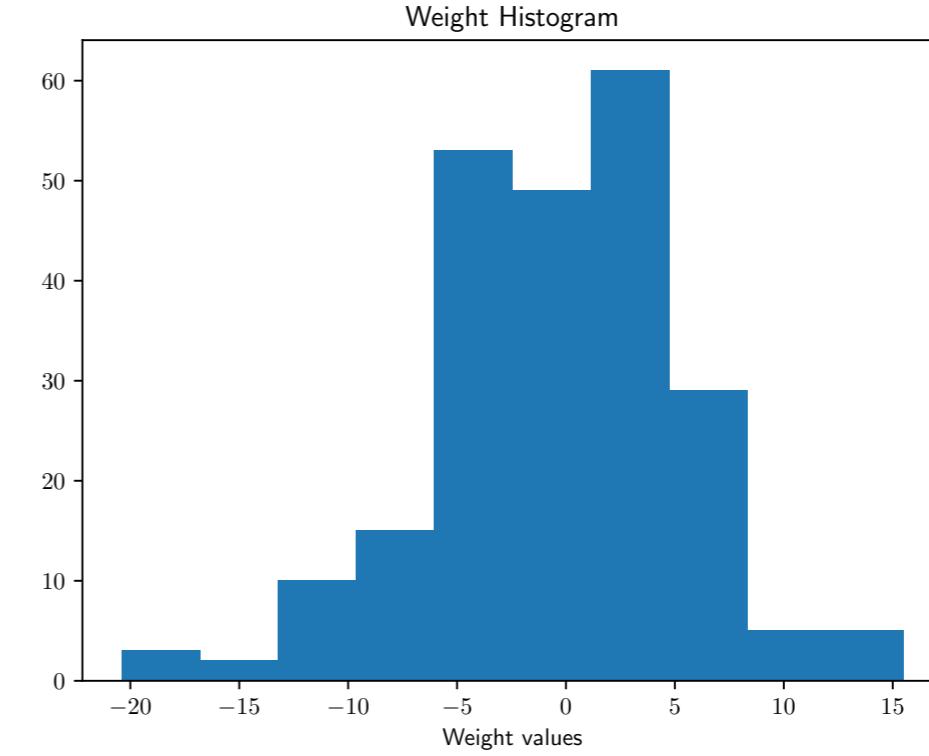
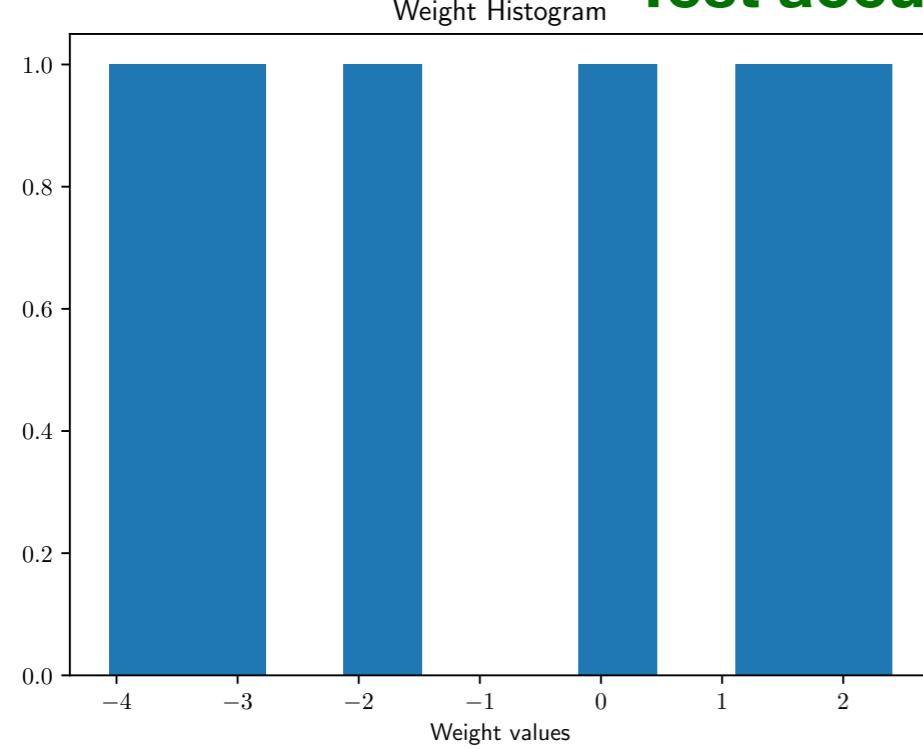
$$\phi(\mathbf{x}) = [x_1, x_2, x_1^2, x_2^2, x_1 x_2]$$

$\phi(\mathbf{x})$ All polynomial terms up to degree 6

Binary Logistic Regression. Adding new features



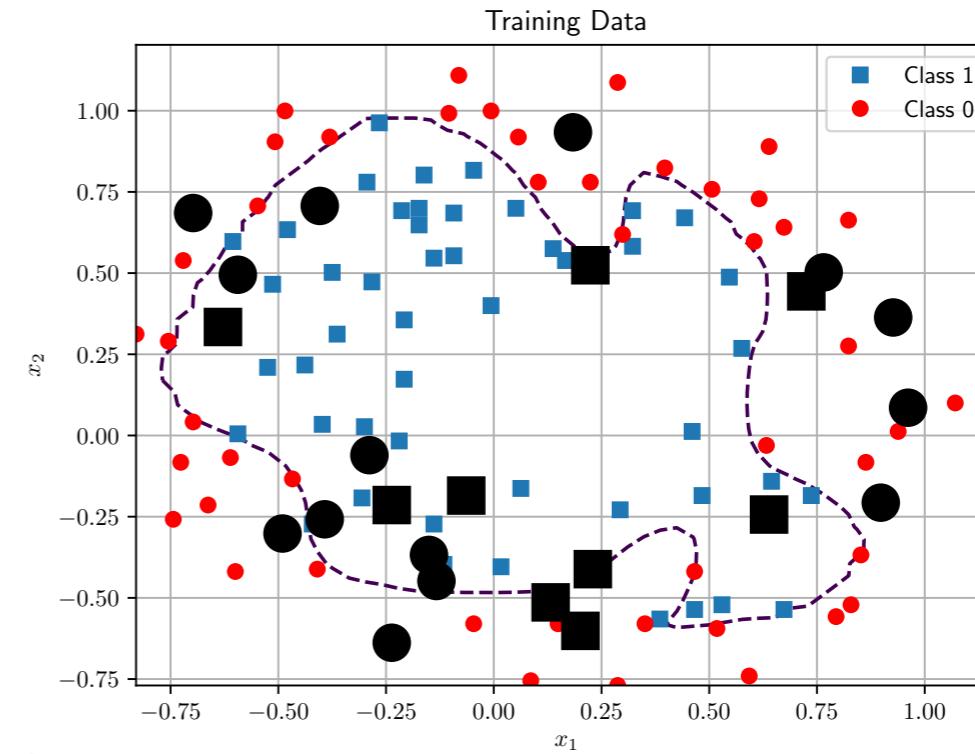
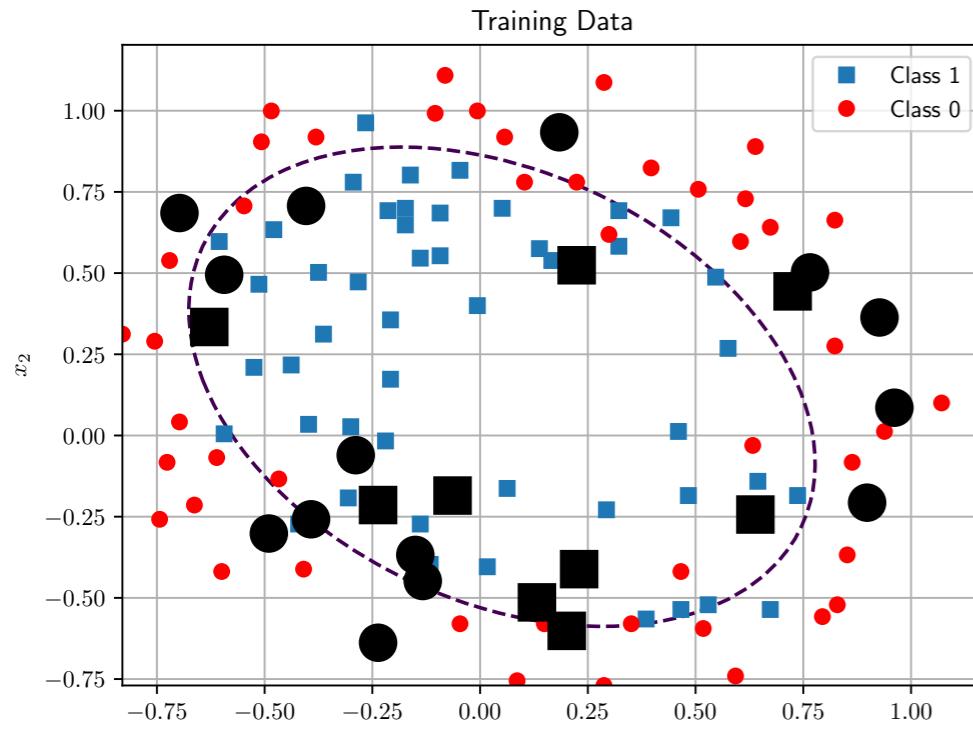
Test accuracy = 65.2 %



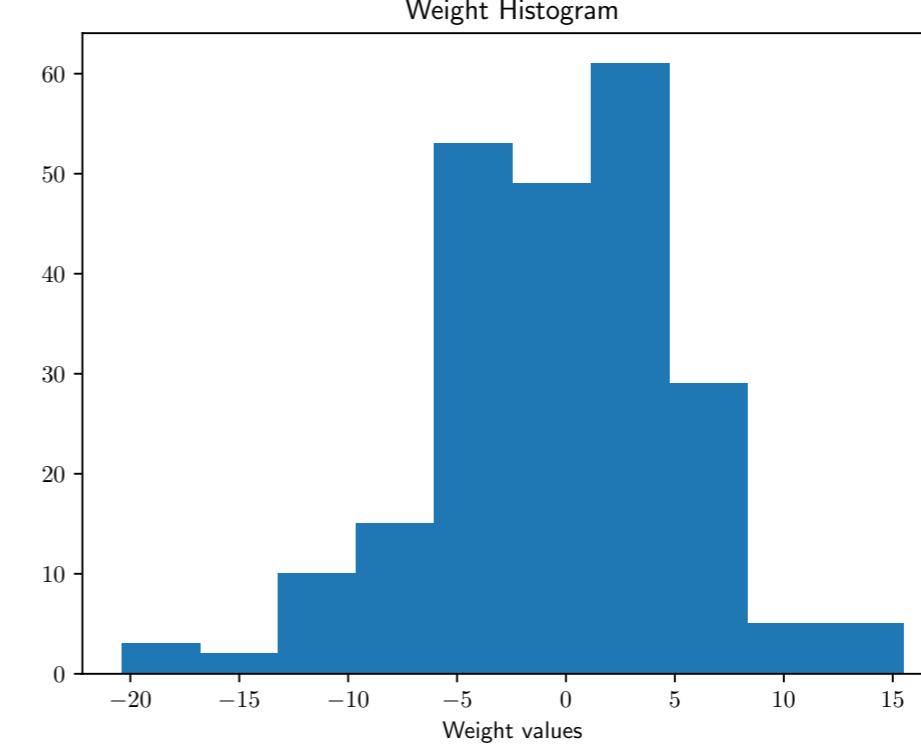
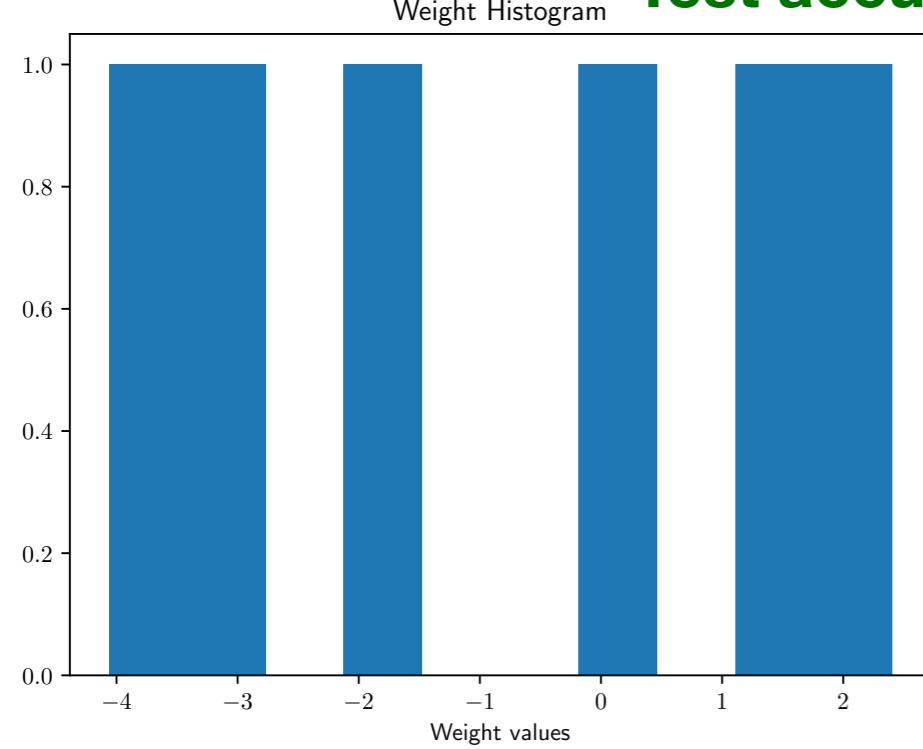
$$\phi(\mathbf{x}) = [x_1, x_2, x_1^2, x_2^2, x_1 x_2]$$

$\phi(\mathbf{x})$ All polynomial terms up to degree 6

Binary Logistic Regression. Adding new features



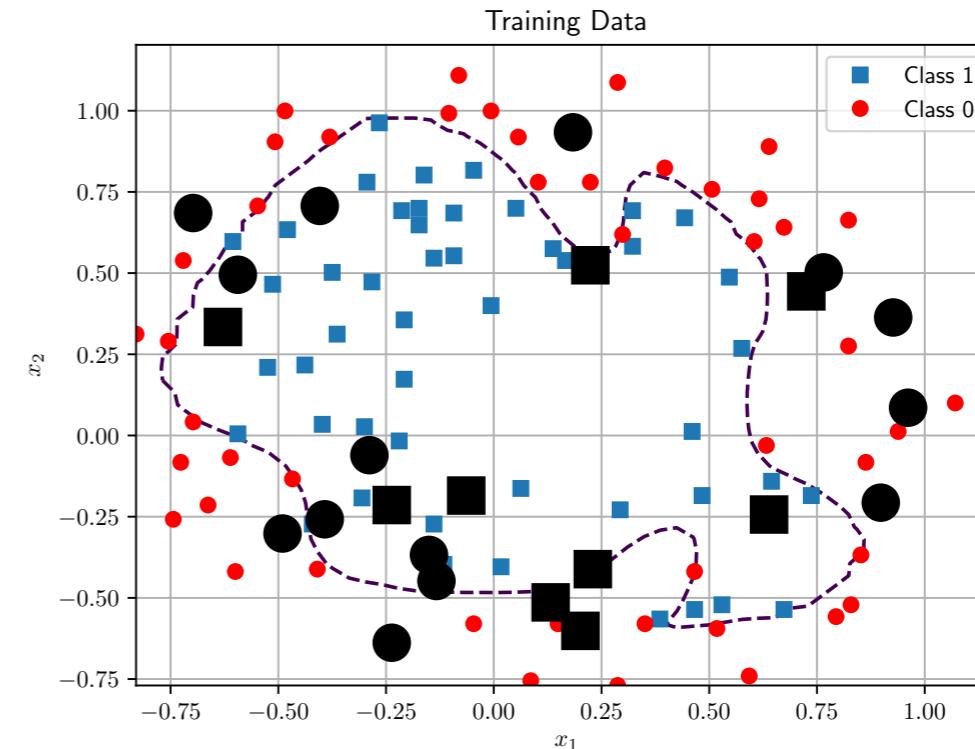
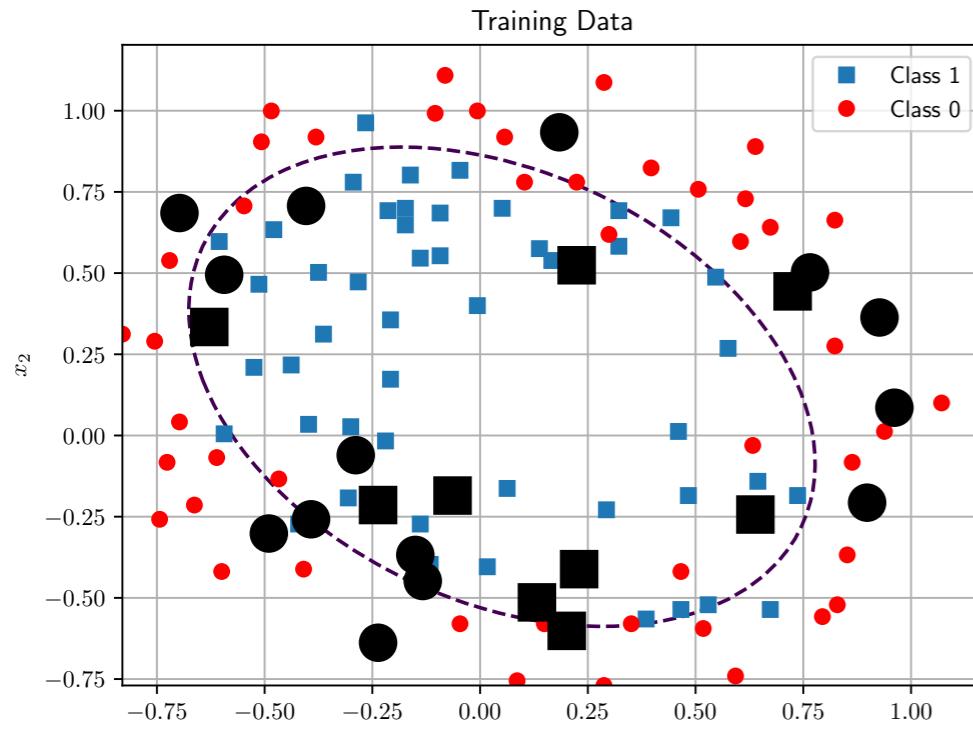
Test accuracy = 65.2 %



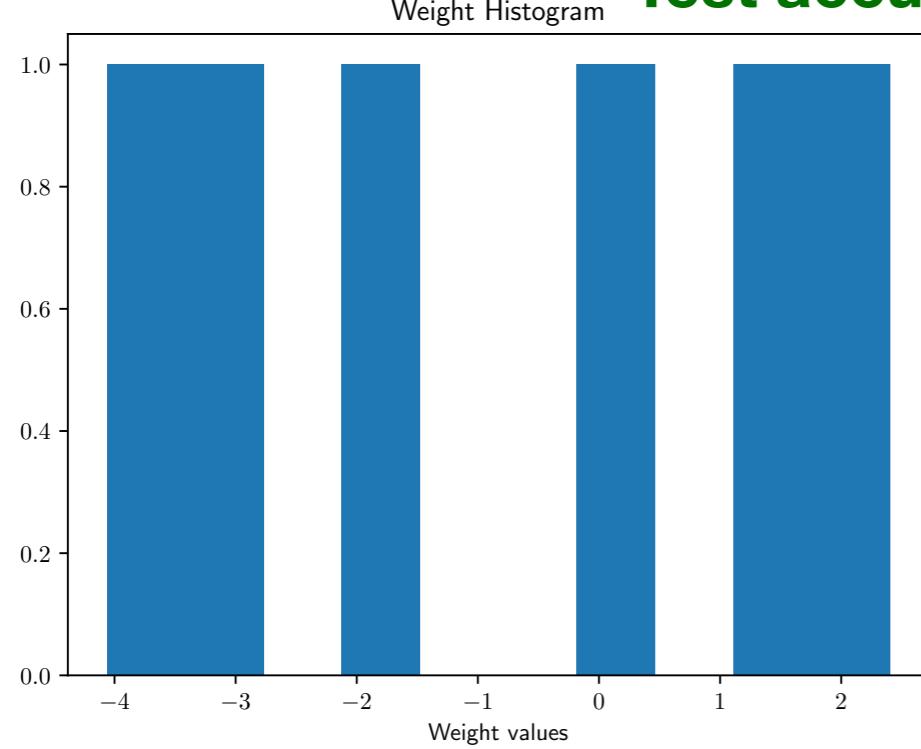
$$\phi(\mathbf{x}) = [x_1, x_2, x_1^2, x_2^2, x_1 x_2]$$

$\phi(\mathbf{x})$ All polynomial terms up to degree 6

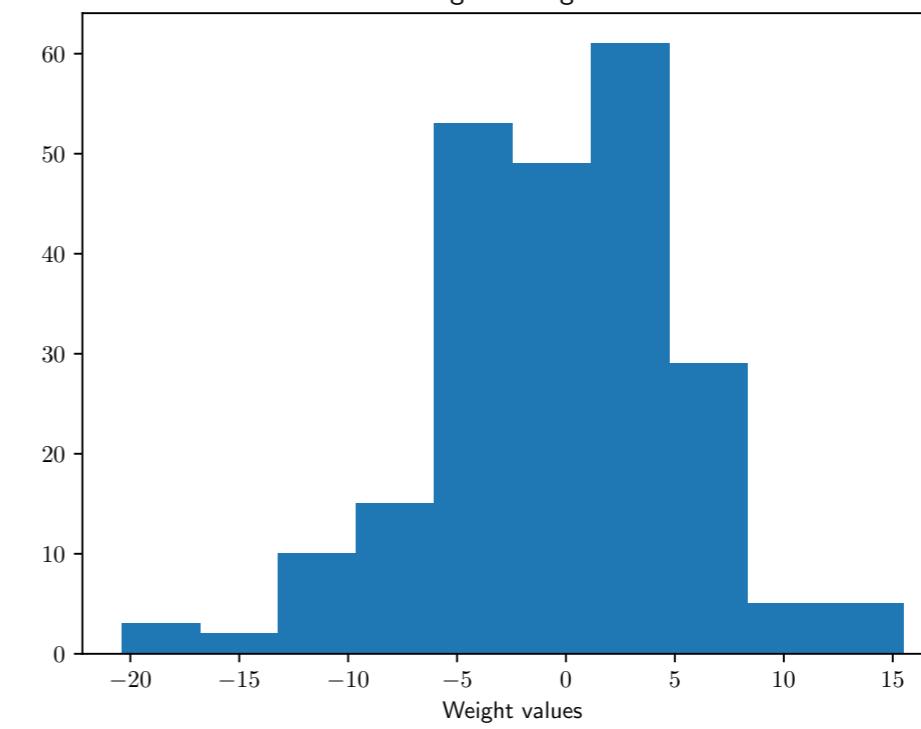
Binary Logistic Regression. Adding new features



Test accuracy = 65.2 %



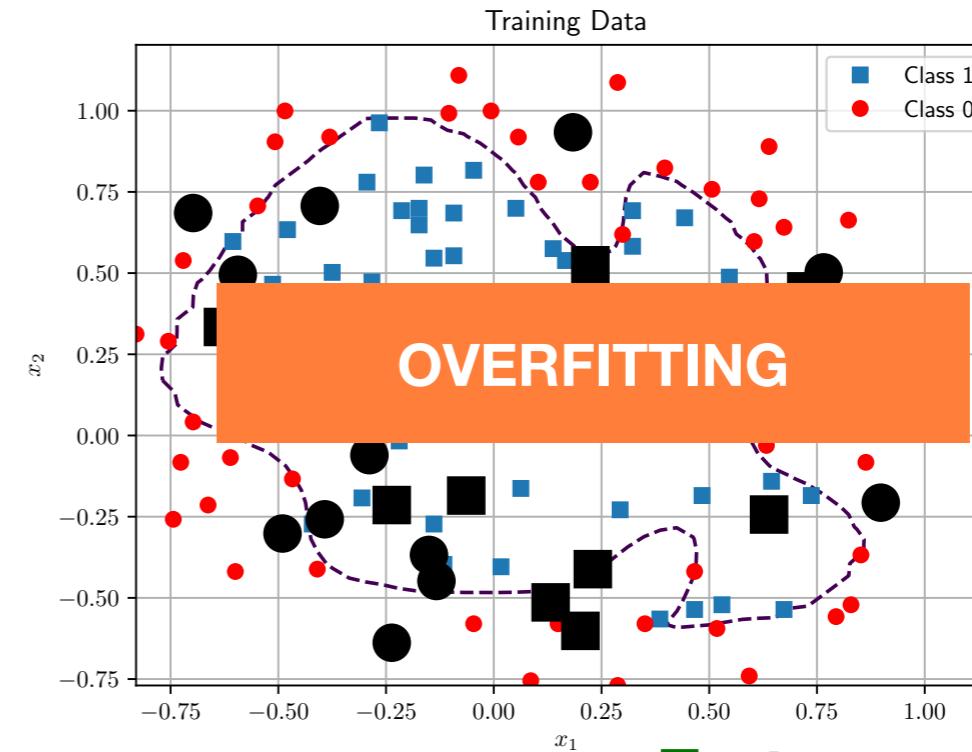
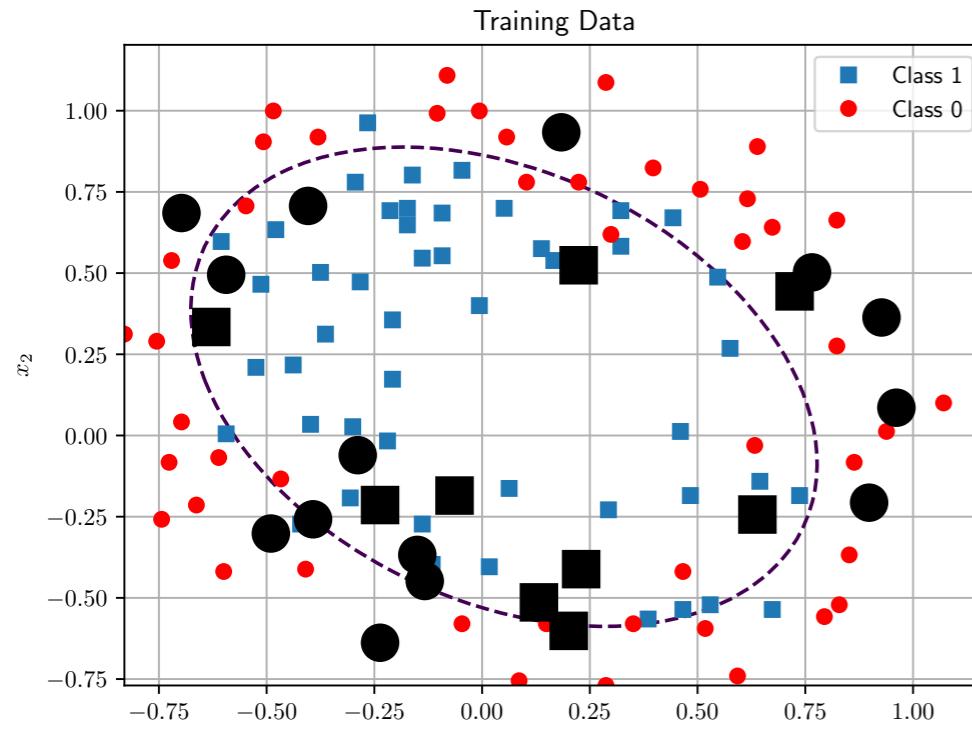
Test accuracy = 60.8 %



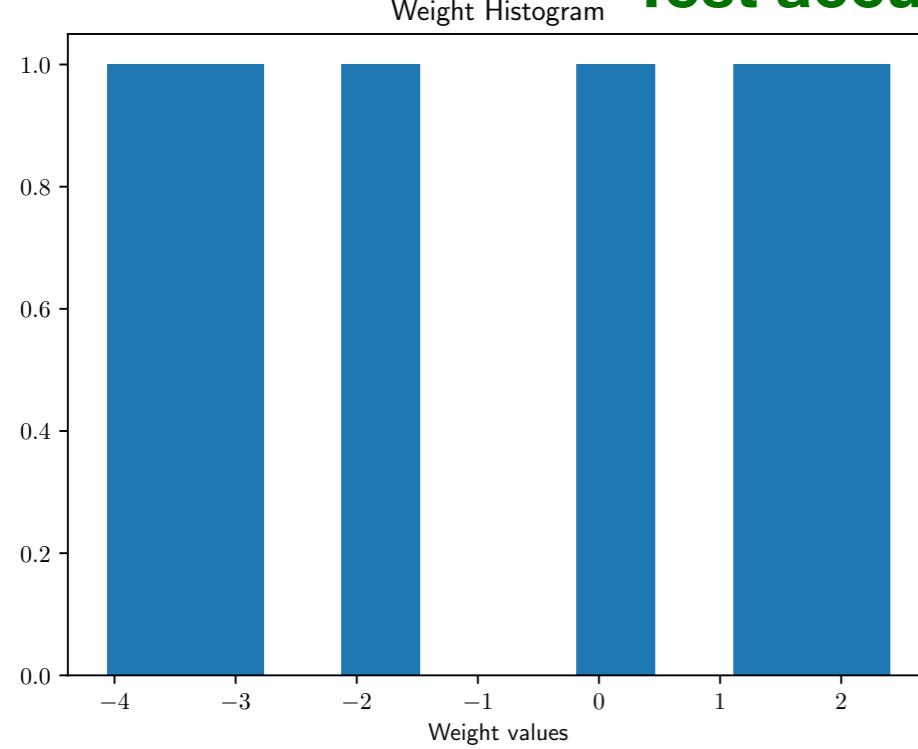
$$\phi(\mathbf{x}) = [x_1, x_2, x_1^2, x_2^2, x_1 x_2]$$

$\phi(\mathbf{x})$ All polynomial terms up to degree 6

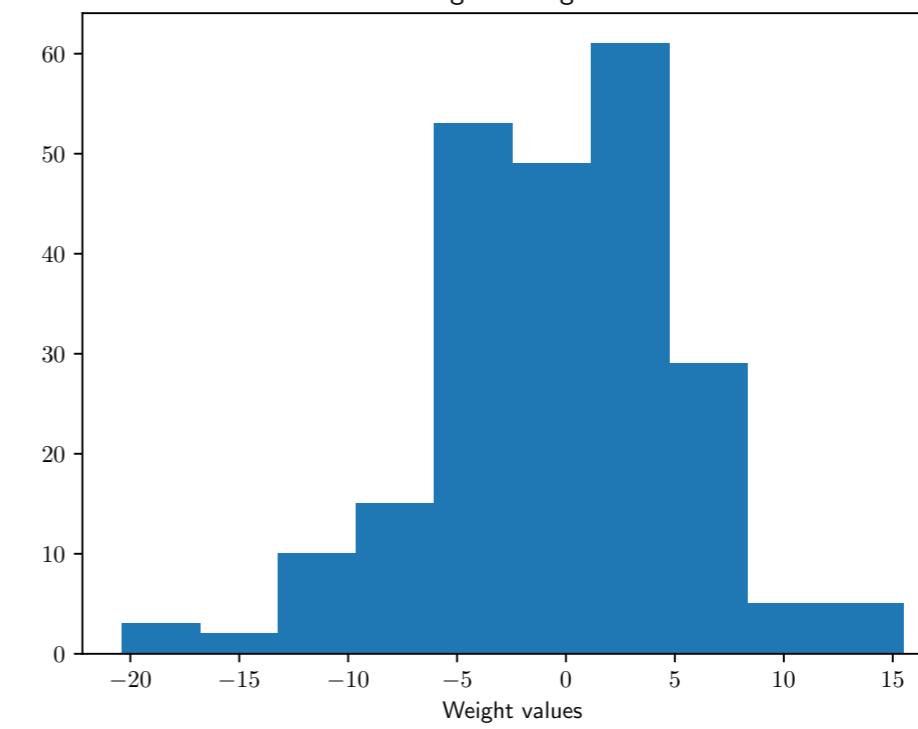
Binary Logistic Regression. Adding new features



Test accuracy = 65.2 %



Test accuracy = 60.8 %



$$\phi(\mathbf{x}) = [x_1, x_2, x_1^2, x_2^2, x_1 x_2]$$

$\phi(\mathbf{x})$ All polynomial terms up to degree 6

Weight norm regularization

- Large weights are associated to overfitting issues
- Instead of manually reducing the number of features, it is easier to penalize solution for which

$$\|w\|_2^2 = \sum w_j^2 \quad \text{L2 regularization}$$

or

$$\|w\|_1 = \sum |w_j| \quad \text{L1 regularization}$$

are very large.

Weight norm regularization

- Large weights are associated to overfitting issues
- Instead of manually reducing the number of features, it is easier to penalize solution for which

$$\|\mathbf{w}\|_2^2 = \sum w_j^2 \quad \text{L2 regularization}$$

or

$$\|\mathbf{w}\|_1 = \sum |w_j| \quad \text{L1 regularization}$$

are very large.

- Regularized Logistic Regression:

$$\mathcal{L} = -\log P(\mathbf{y}|\mathbf{X}) + \lambda \|\mathbf{w}\|_1$$

$$\mathcal{L} = -\log P(\mathbf{y}|\mathbf{X}) + \lambda \|\mathbf{w}\|_2^2$$

Weight norm regularization

- Large weights are associated to overfitting issues
- Instead of manually reducing the number of features, it is easier to penalize solution for which

$$\|\mathbf{w}\|_2^2 = \sum w_j^2 \quad \text{L2 regularization}$$

or

$$\|\mathbf{w}\|_1 = \sum |w_j| \quad \text{L1 regularization}$$

are very large.

- Regularized Logistic Regression:

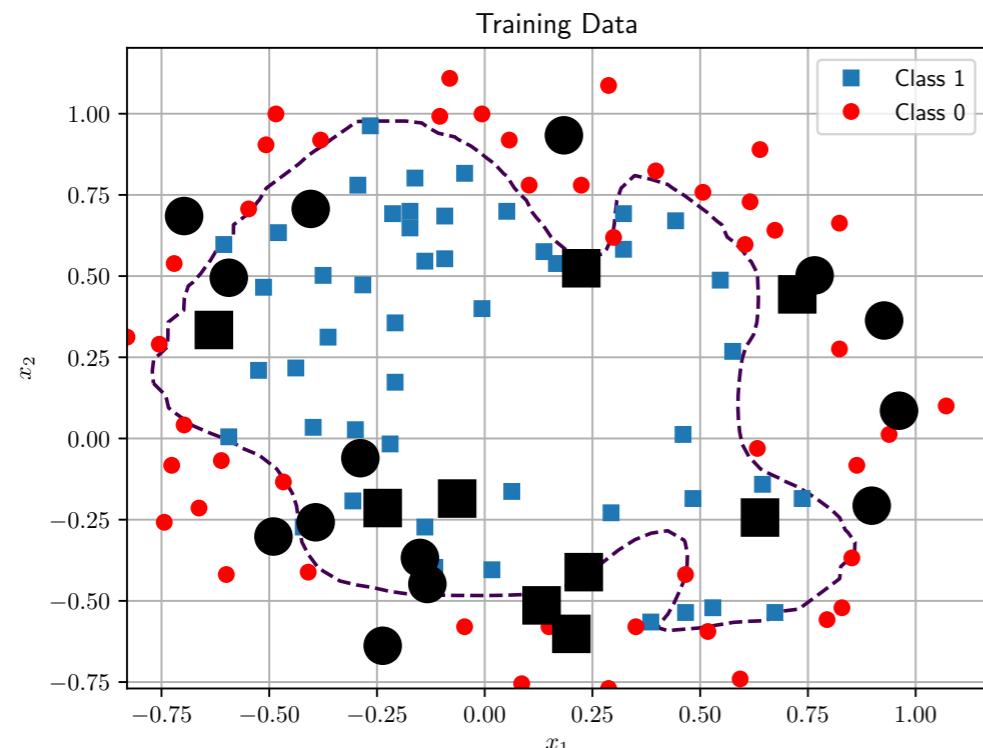
$$\mathcal{L} = -\log P(\mathbf{y}|\mathbf{X}) + \lambda \|\mathbf{w}\|_1$$

Regularization Parameter

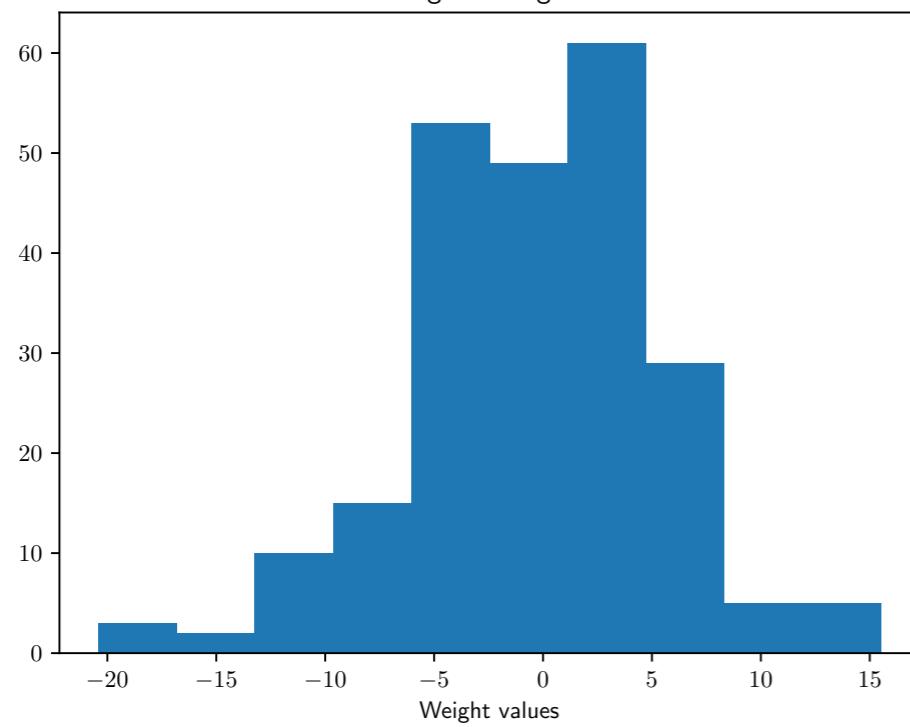
Cross validation!!

$$\mathcal{L} = -\log P(\mathbf{y}|\mathbf{X}) + \lambda \|\mathbf{w}\|_2^2$$

Weight norm regularization

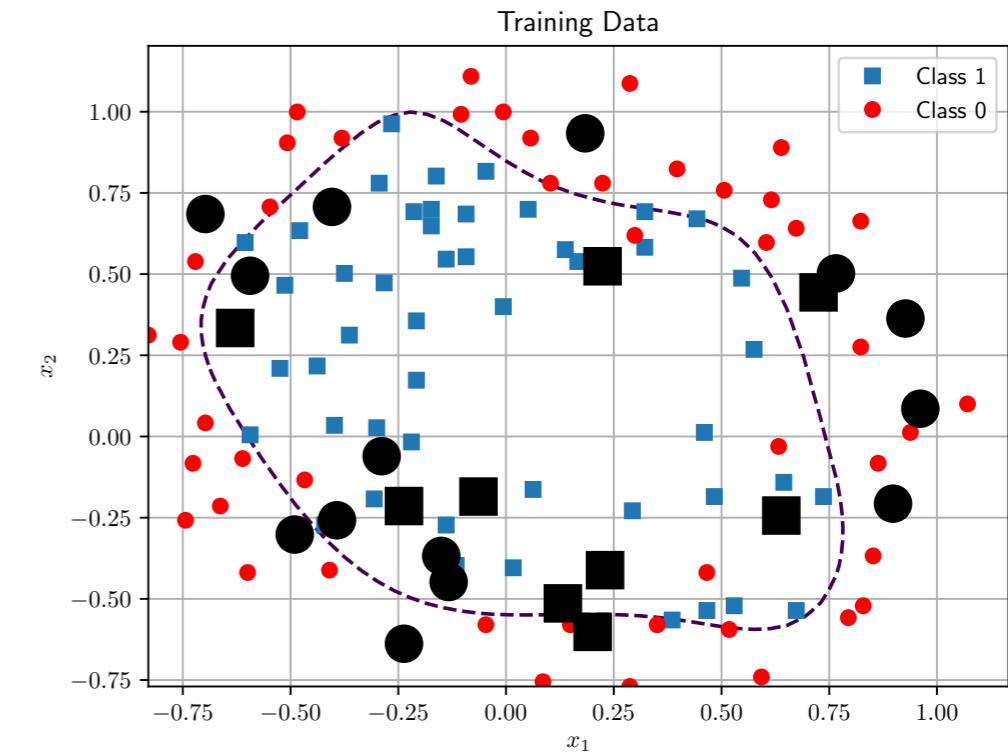
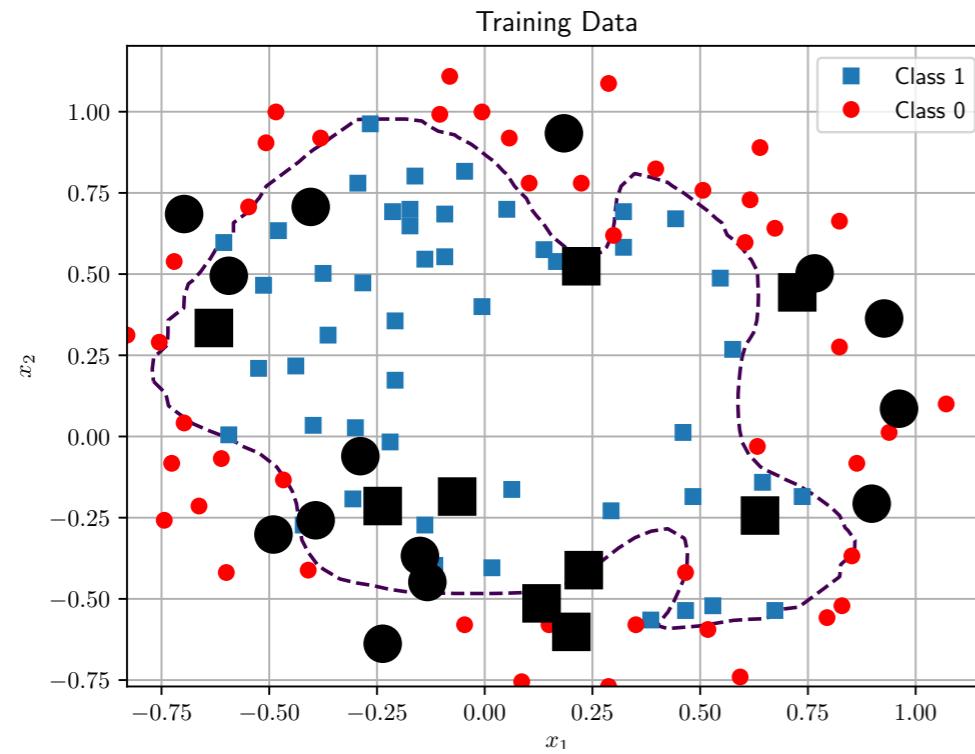


Test accuracy = 60.8 %



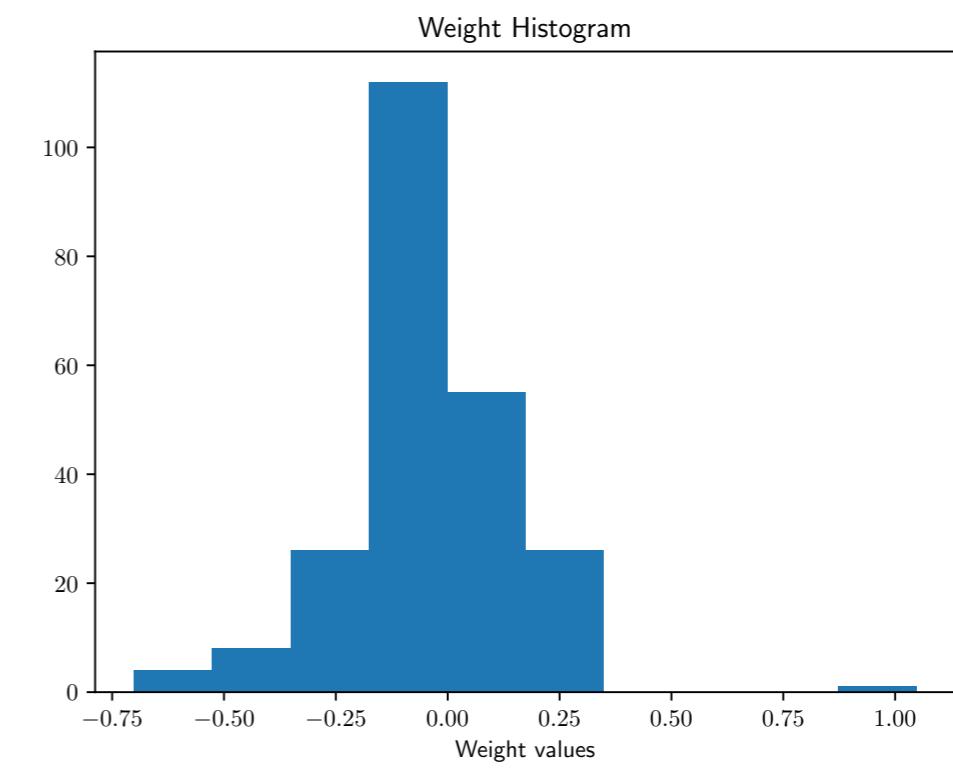
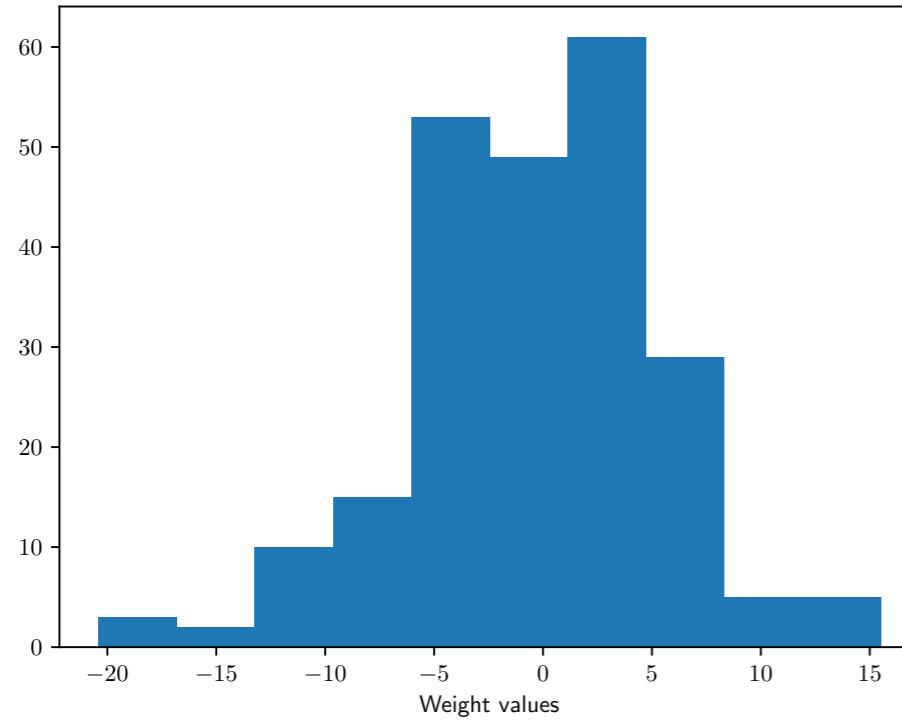
$\phi(\mathbf{x})$ All polynomial terms up to degree 6

Weight norm regularization



Weight Histogram

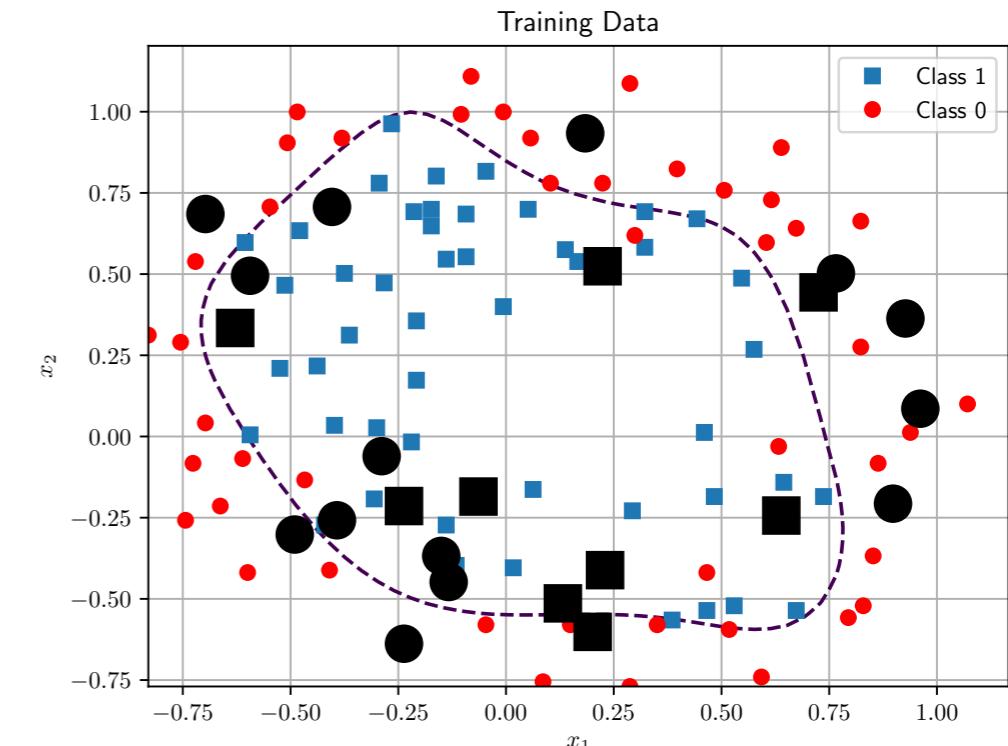
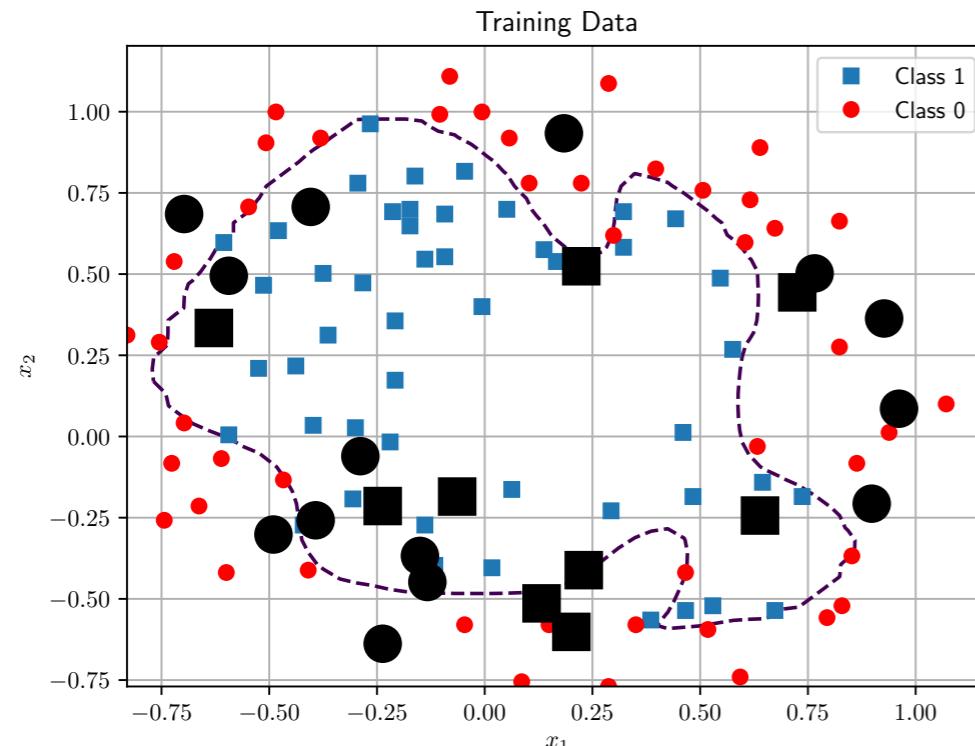
Test accuracy = 60.8 %



$\phi(x)$ All polynomial terms up to degree 6

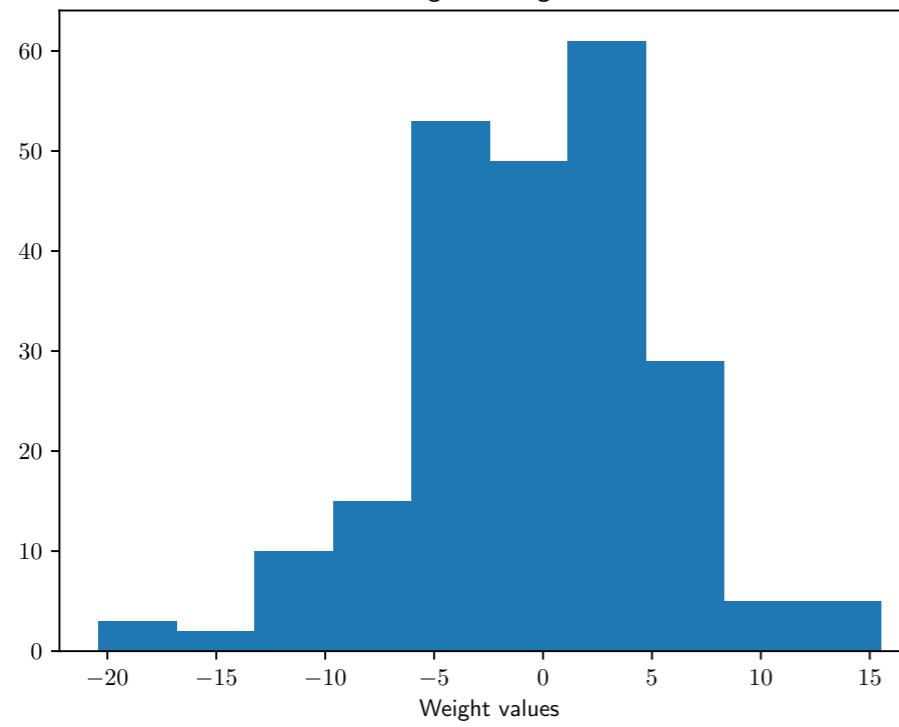
$\lambda = 1.0, L_2$

Weight norm regularization



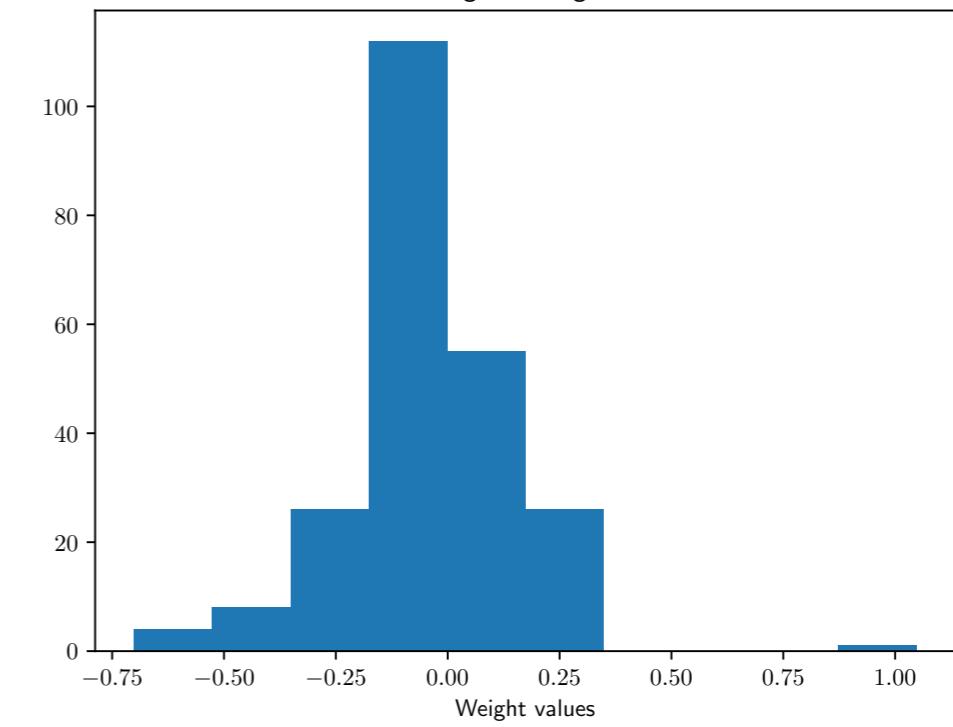
Weight Histogram

Test accuracy = 60.8 %



Weight Histogram

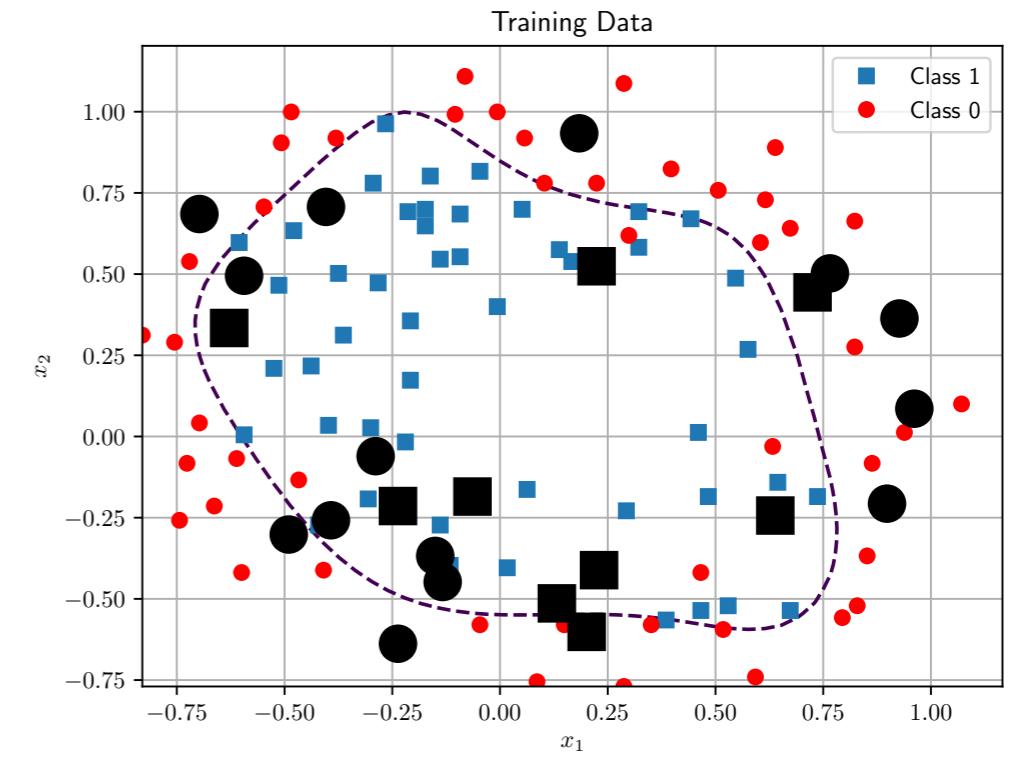
Test accuracy = 65.2 %



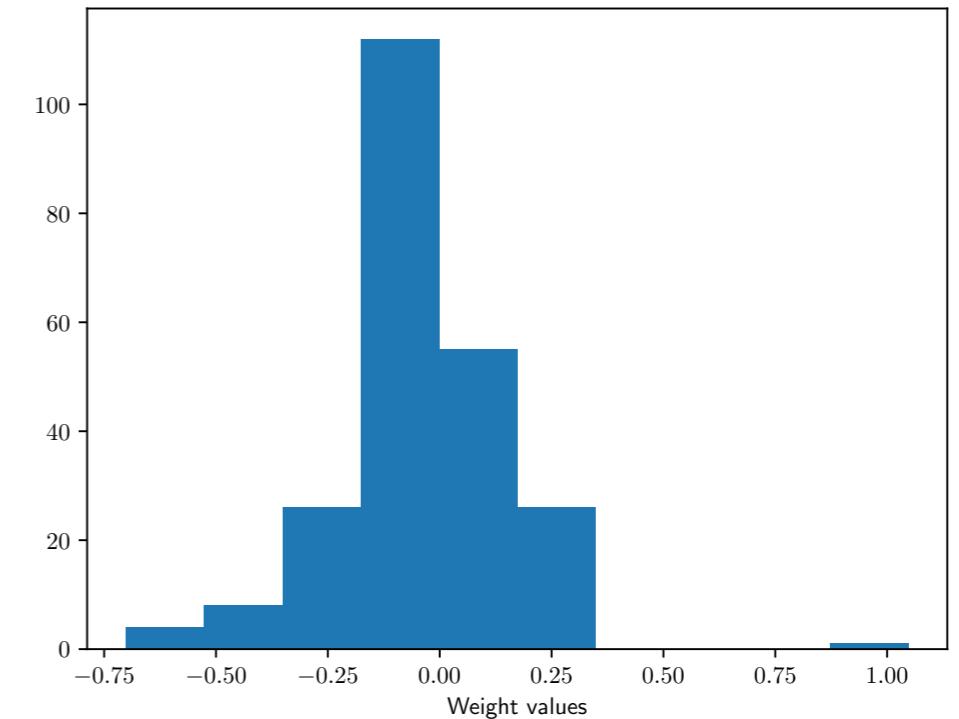
$\phi(\mathbf{x})$ All polynomial terms up to degree 6

$\lambda = 1.0, L_2$

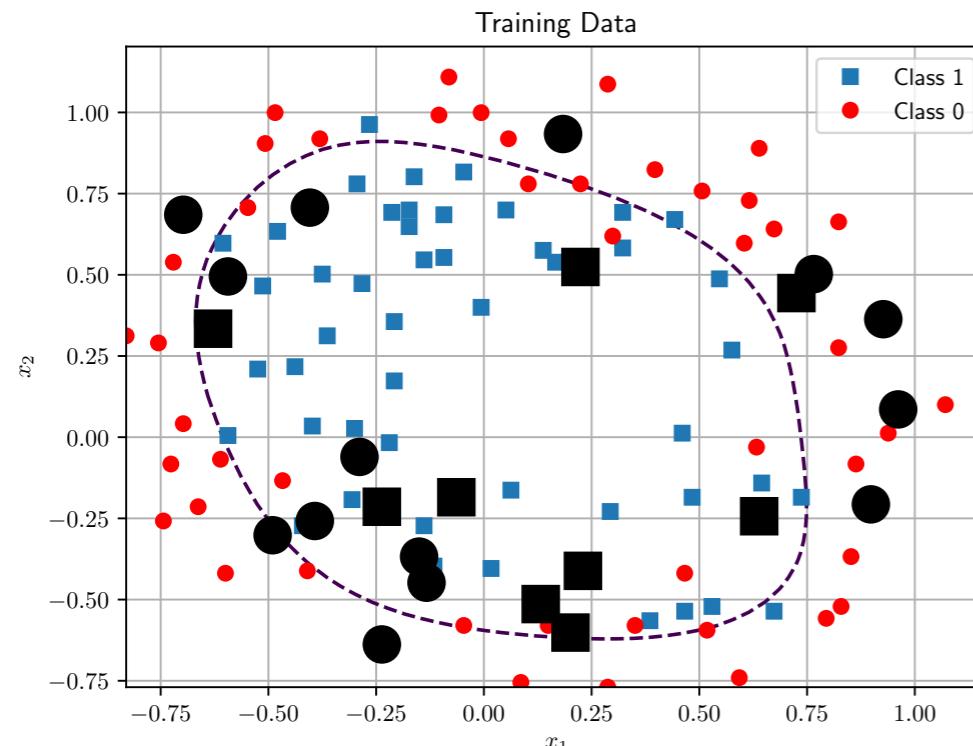
Weight norm regularization



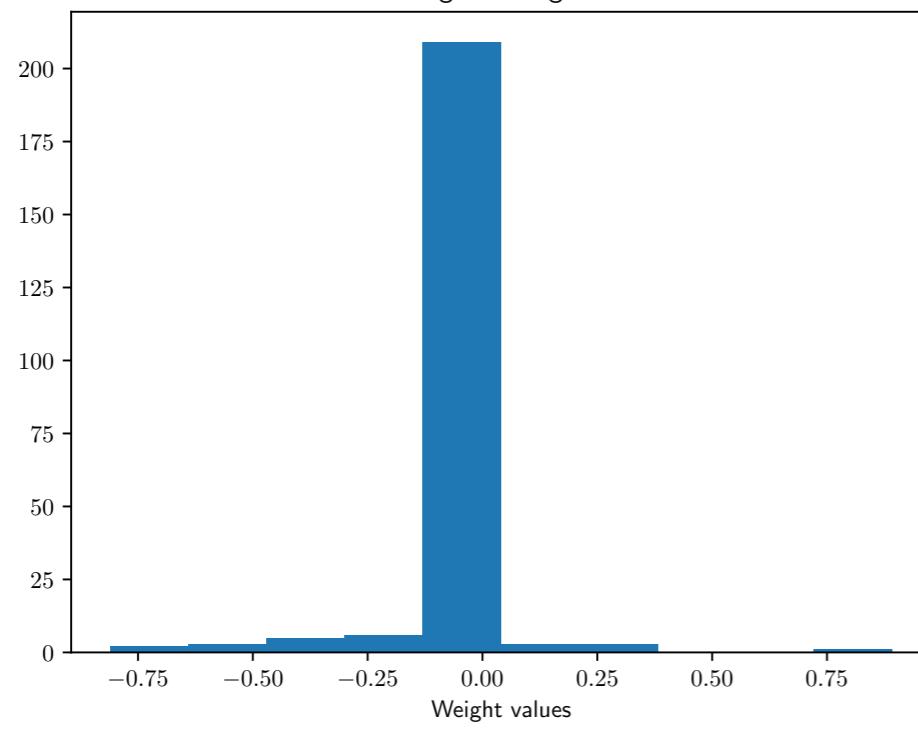
Test accuracy = 65.2 %



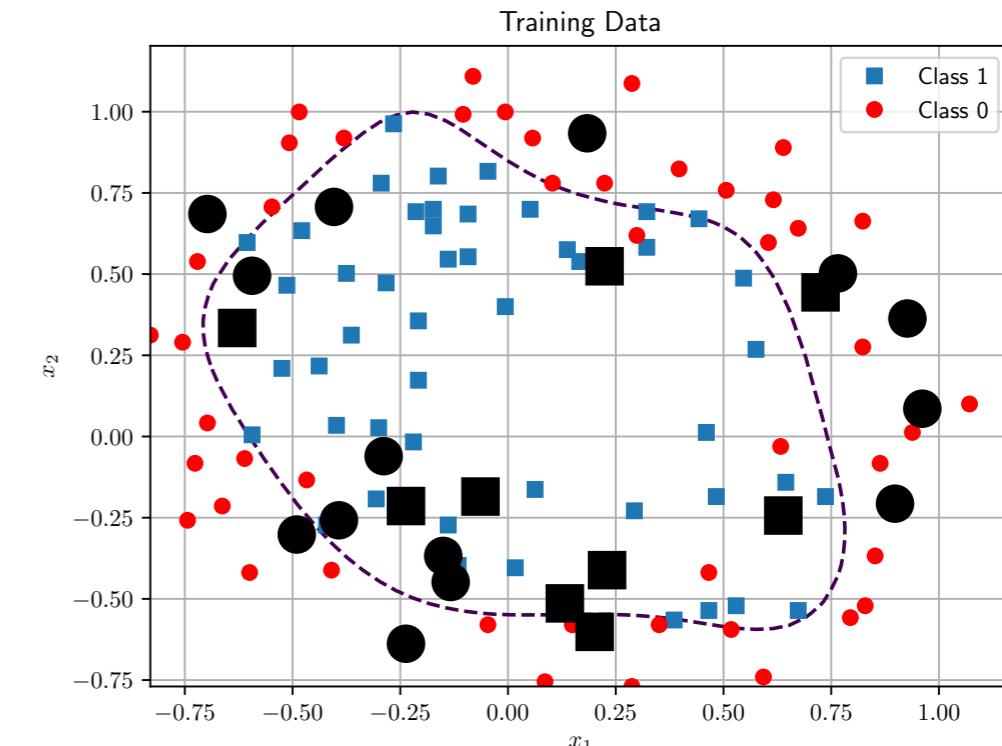
Weight norm regularization



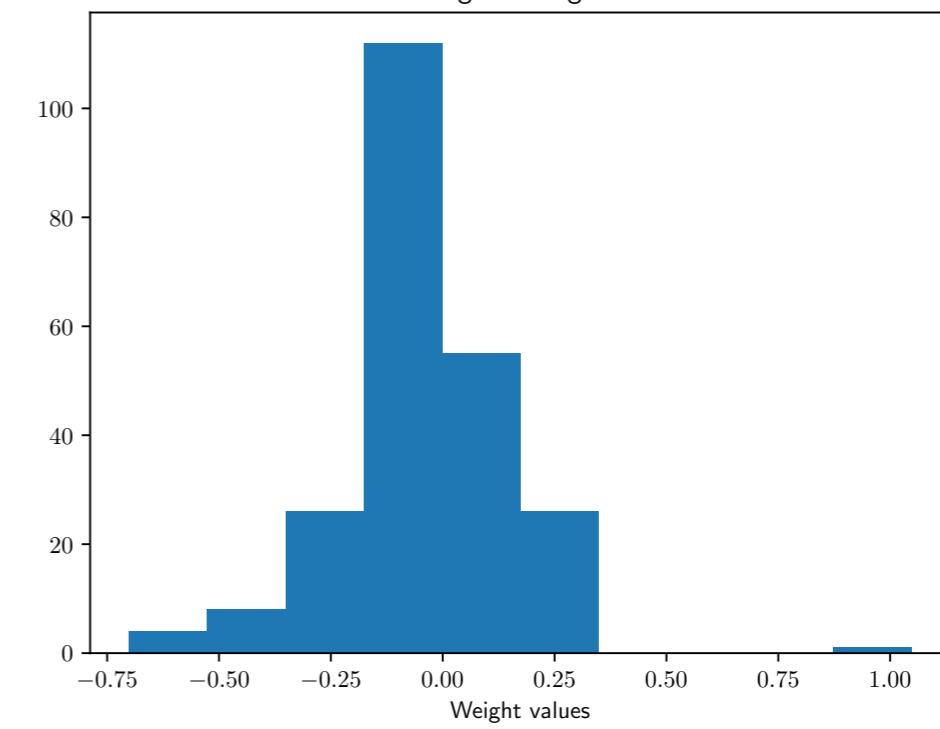
Weight Histogram



$$\lambda = 1.0, \quad L_1$$



Weight Histogram



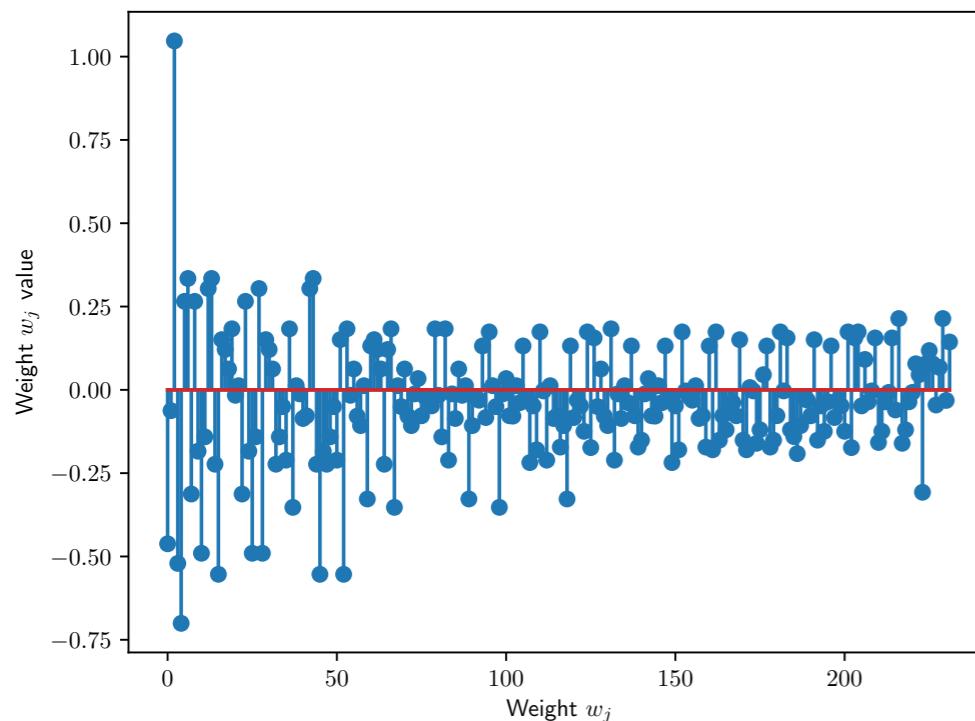
$$\lambda = 1.0, \quad L_2$$

L1 norm regularization induces sparsity

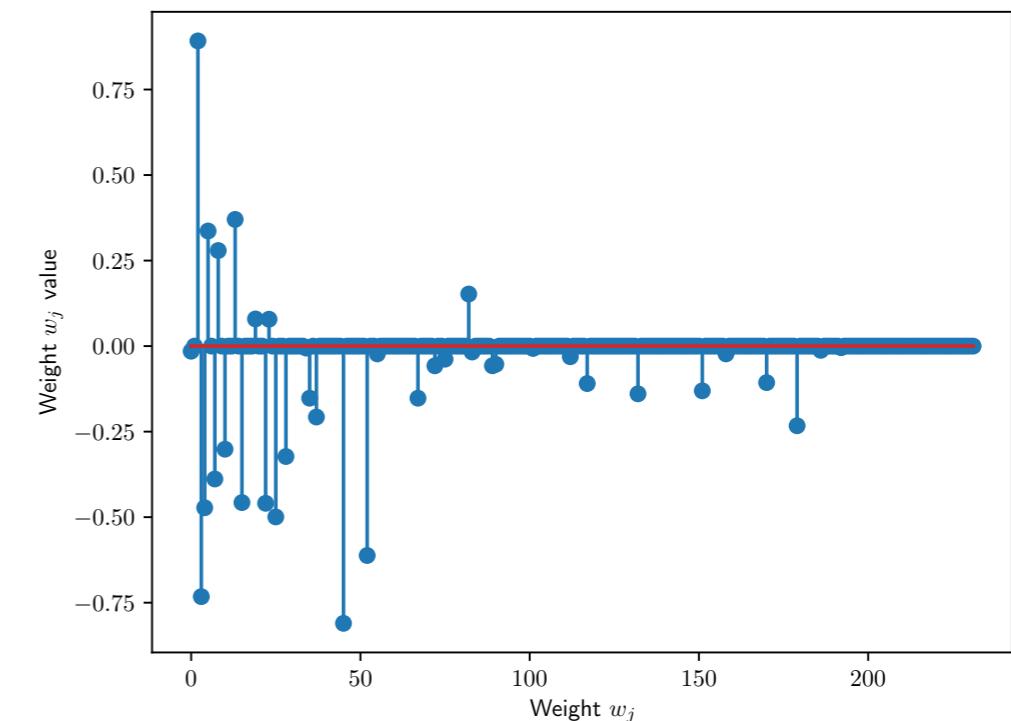
- With L1 regularization, many weights go to zero. Good for feature selection!
- Numerical optimization with L1-type of regularizers is more challenging

Check out in Aula Global this explanation about sparsity with L1 (by Emily Fox):

[Link](#)



$$\lambda = 1.0, \quad L_2$$



$$\lambda = 1.0, \quad L_1$$

Logistic Regression: multiple classes

Binary Logistic Regression

- Discriminative classifier:

$$P(y = 1|\boldsymbol{x}) = \frac{1}{1 + e^{-(\boldsymbol{w}^T \boldsymbol{x} + w_0)}} \doteq \sigma(\boldsymbol{w}^T \boldsymbol{x} + w_0)$$

Binary Logistic Regression

- Discriminative classifier:

$$P(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + w_0)}} \doteq \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$



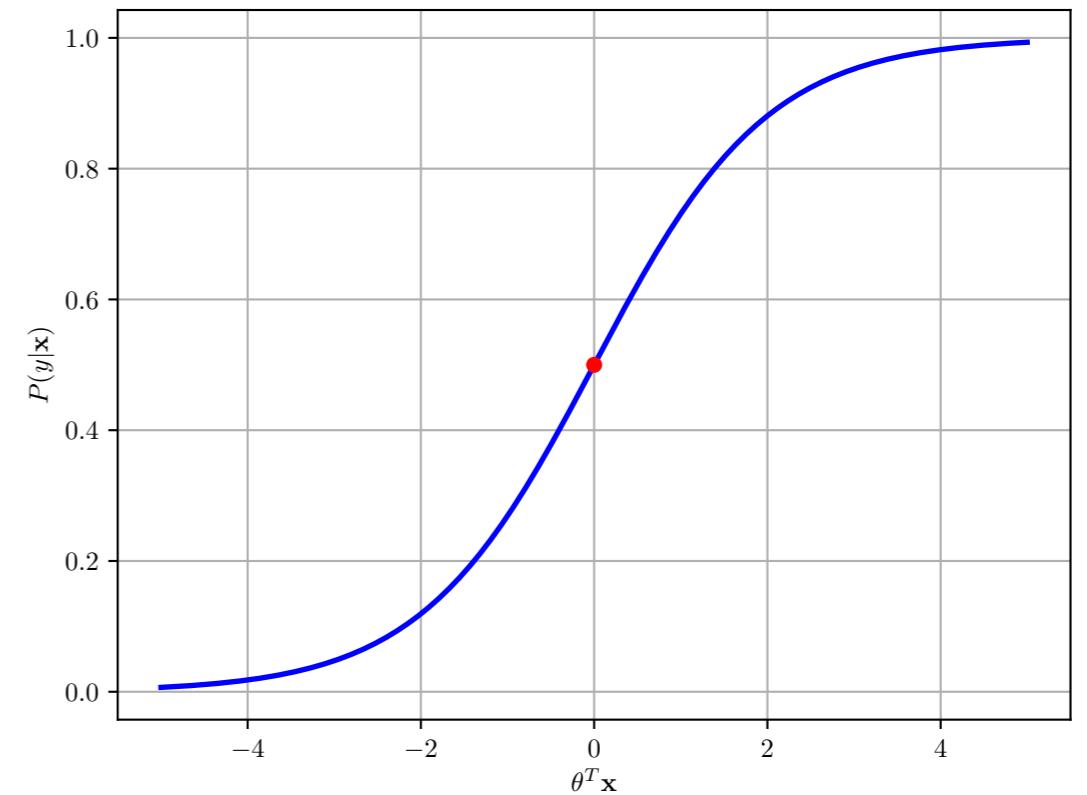
Generalized linear model

Binary Logistic Regression

- Discriminative classifier:

$$P(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + w_0)}} \doteq \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

Generalized linear model



Binary Logistic Regression

- Discriminative classifier:

$$P(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + w_0)}} \doteq \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

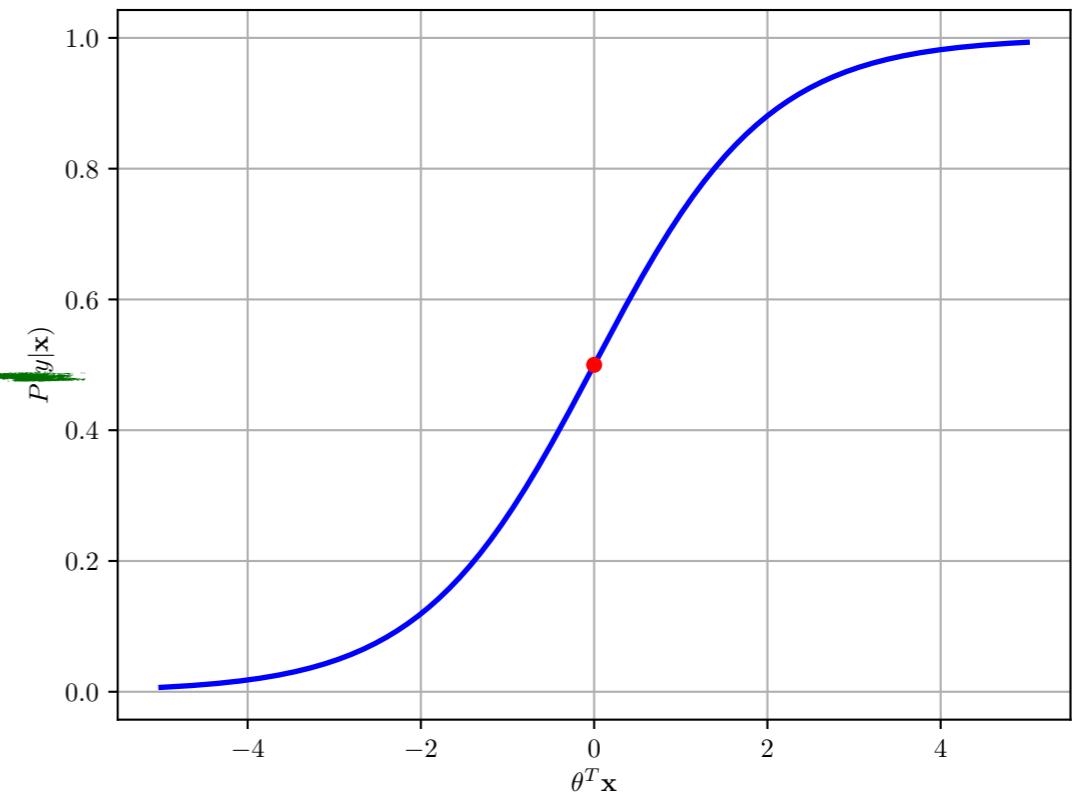
Generalized linear model



- Decision boundary:

$$\{\mathbf{x} \in \mathbb{R}^m : \mathbf{w}^T \mathbf{x} + w_0 = 0\}$$

Hyperplane! Linear classifier



Multiclass Logistic Regression

- Training database: $\mathcal{D} \doteq (\boldsymbol{x}^{(i)}, y^{(i)})_{i=1}^N \quad \boldsymbol{x}^{(i)} \in \mathbb{R}^m \quad y^{(i)} \in \{1, \dots, K\}$
- Discriminative classifier based on the **Softmax function**:

$$P(y = k | \boldsymbol{x}) = \frac{e^{-(z_k)}}{\sum_{j=1}^K e^{-(z_j)}}, \quad \boldsymbol{z} = \boldsymbol{W} \begin{bmatrix} 1.0 \\ \boldsymbol{x}^{(i)} \end{bmatrix}$$

- Similar gradient expression

Multiclass Logistic Regression

- Training database: $\mathcal{D} \doteq (\mathbf{x}^{(i)}, y^{(i)})_{i=1}^N \quad \mathbf{x}^{(i)} \in \mathbb{R}^m \quad y^{(i)} \in \{1, \dots, K\}$
- Discriminative classifier based on the **Softmax function**:

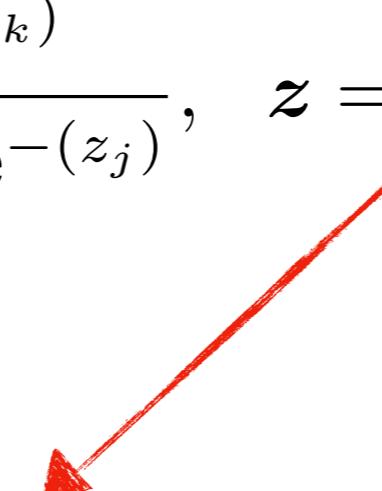
$$P(y = k | \mathbf{x}) = \frac{e^{-(z_k)}}{\sum_{j=1}^K e^{-(z_j)}}, \quad z = \mathbf{W} \begin{bmatrix} 1.0 \\ \mathbf{x}^{(i)} \end{bmatrix}$$

$$K \times (m + 1)$$

- Similar gradient expression

Multiclass Logistic Regression

- Training database: $\mathcal{D} \doteq (\mathbf{x}^{(i)}, y^{(i)})_{i=1}^N \quad \mathbf{x}^{(i)} \in \mathbb{R}^m \quad y^{(i)} \in \{1, \dots, K\}$
- Discriminative classifier based on the **Softmax function**:

$$P(y = k|\mathbf{x}) = \frac{e^{-(z_k)}}{\sum_{j=1}^K e^{-(z_j)}}, \quad \mathbf{z} = \mathbf{W} \begin{bmatrix} 1.0 \\ \mathbf{x}^{(i)} \end{bmatrix}$$


- Loss function (Cross entropy):

$$\mathcal{L} = -\log P(\mathbf{y}|\mathbf{X}) = -\sum_{n=1}^N \sum_{k=1}^K \mathbb{1}[y^{(n)} == k] \log P(y = k|\mathbf{x})$$

- Similar gradient expression