

Elements of Machine Learning & Data Science

Winter semester 2025/26

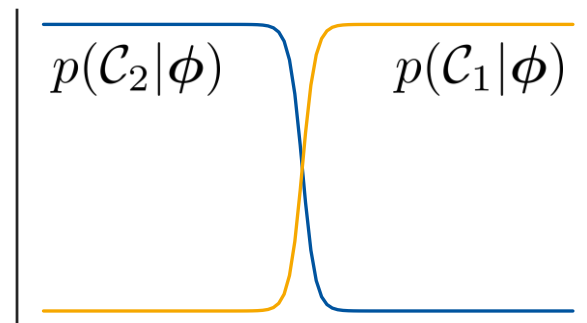
Lecture 13 – Logistic Regression

08.12.2025

Prof. Bastian Leibe

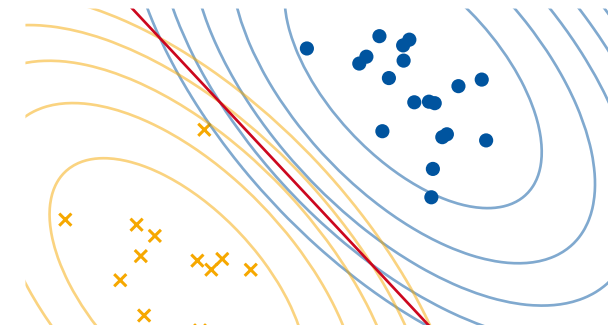
Machine Learning Topics

- 8. Introduction to ML
- 9. Probability Density Estimation
- 10. Linear Discriminants
- 11. Linear Regression
- 12. Logistic Regression**
- 13. Support Vector Machines
- 14. Neural Network Basics



$$y(\phi) = p(\mathcal{C}_1|\phi) = \sigma(\mathbf{w}^\top \phi + w_0)$$

Logistic Regression
Formulation

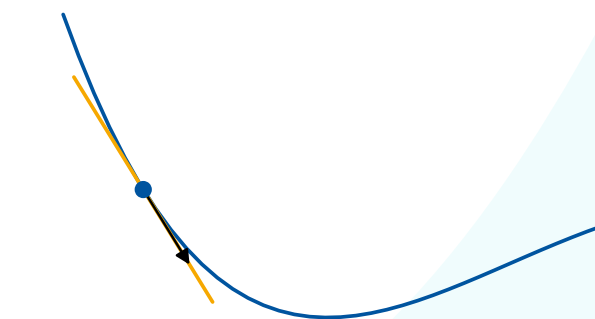


Parameter Efficiency



$$E(\mathbf{w}) = -\ln p(\mathbf{t}|\mathbf{w})$$

Cross-Entropy Error

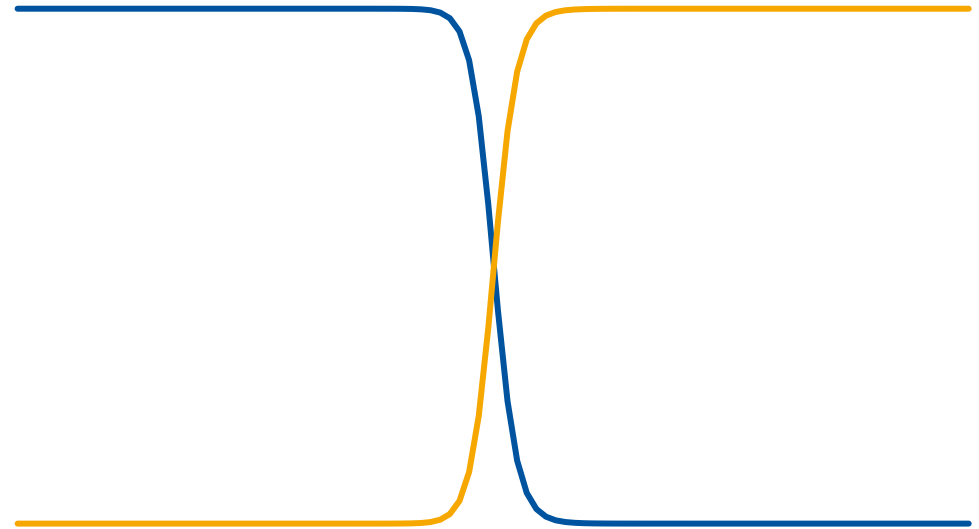


$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E(\mathbf{w})$$

Iterative Optimization

Logistic Regression

1. **Logistic Regression Formulation**
2. Motivation and Background
3. Iterative Optimization
4. First-Order Gradient Descent
5. Second-Order Gradient Descent
6. Error Function Analysis



Motivation

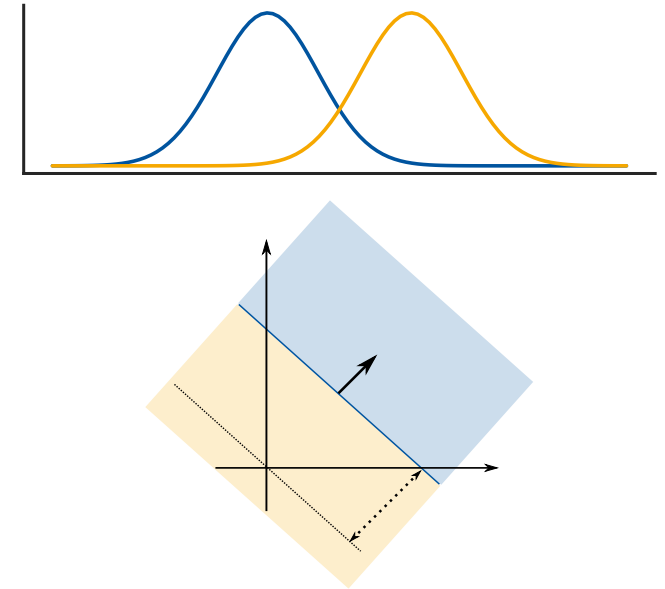
- We have seen how to build probabilistic classifiers using Bayes' Theorem:

$$y_k(\mathbf{x}) = p(\mathcal{C}_k|\mathbf{x}) \propto p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)$$

- We have directly modeled the decision boundary with linear discriminants:

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

- In the following, we will combine those two ideas
 - We will model the posterior $p(\mathcal{C}_k|\mathbf{x})$
 - But we will do that using a linear discriminant function $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$
- The resulting model will be called **logistic regression**.



Reminder: Probabilistic Classification

- Remember what we did in probabilistic classification

- We modeled the likelihood of each class

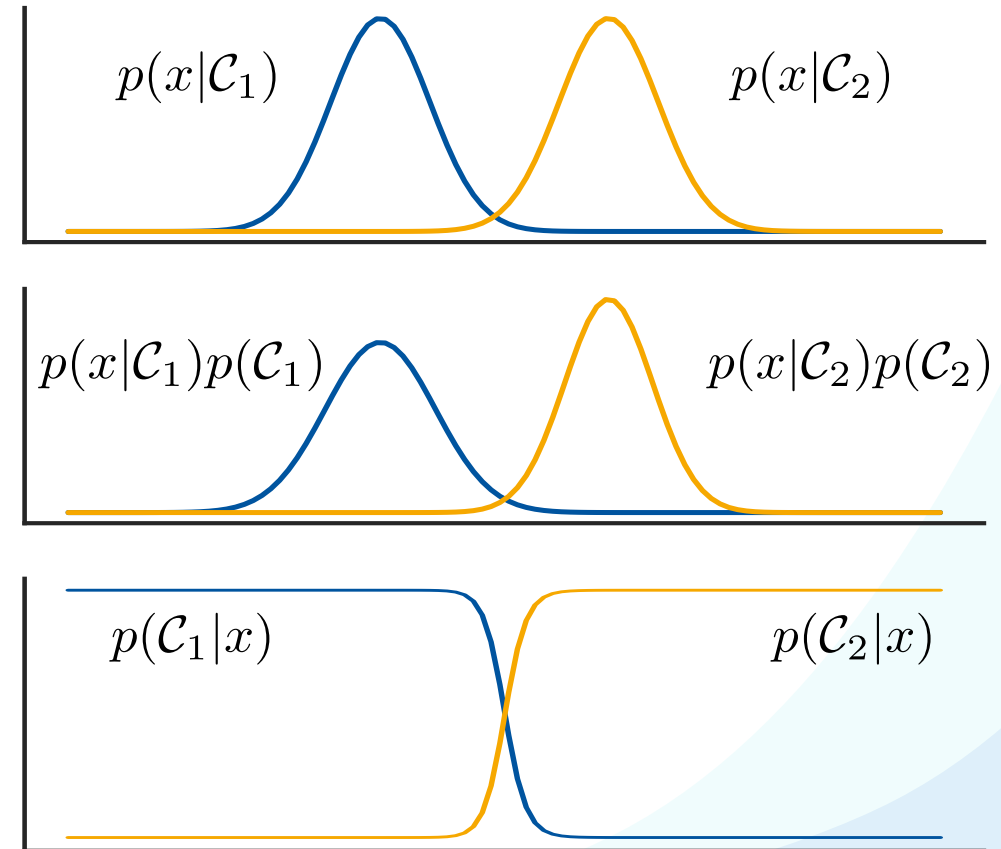
$$p(\mathbf{x}|\mathcal{C}_k)$$

- We scaled the likelihoods with the priors

$$p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)$$

- We normalized to compute the posterior

$$p(\mathcal{C}_1|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1) + p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)}$$



- Let's now start with the posterior and rewrite it

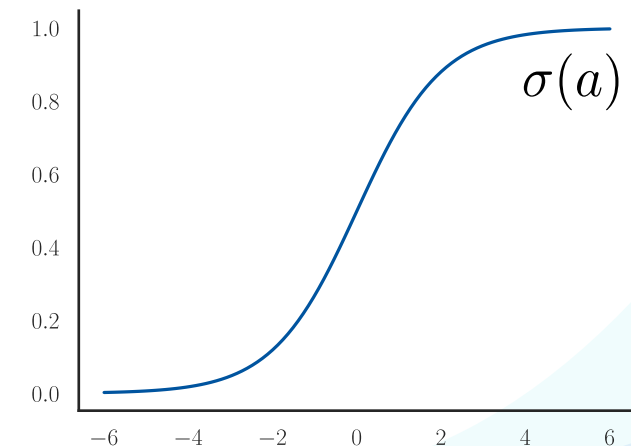
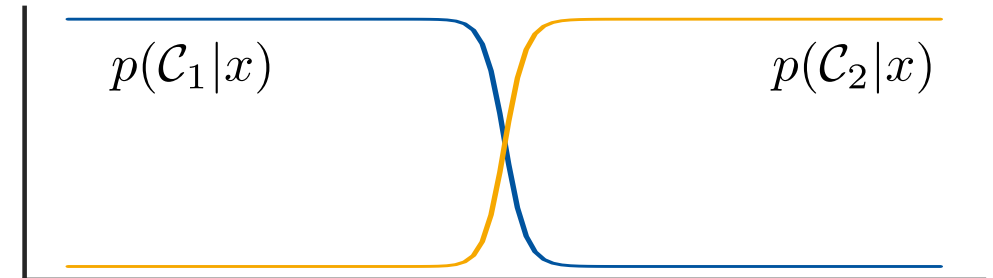
$$\begin{aligned}
 p(\mathcal{C}_1|\mathbf{x}) &= \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1) + p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)} \\
 &= \frac{1}{1 + \frac{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)}{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}} \\
 &= \frac{1}{1 + \exp(-a)} =: \sigma(a)
 \end{aligned}$$

- This is the equation for the **logistic sigmoid** function $\sigma(a)$

⇒ If we set

$$a = \ln \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)}$$

the logistic sigmoid expresses a posterior probability!



Properties of the Logistic Sigmoid

- Definition:

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

- Inverse (also known as **logit** function):

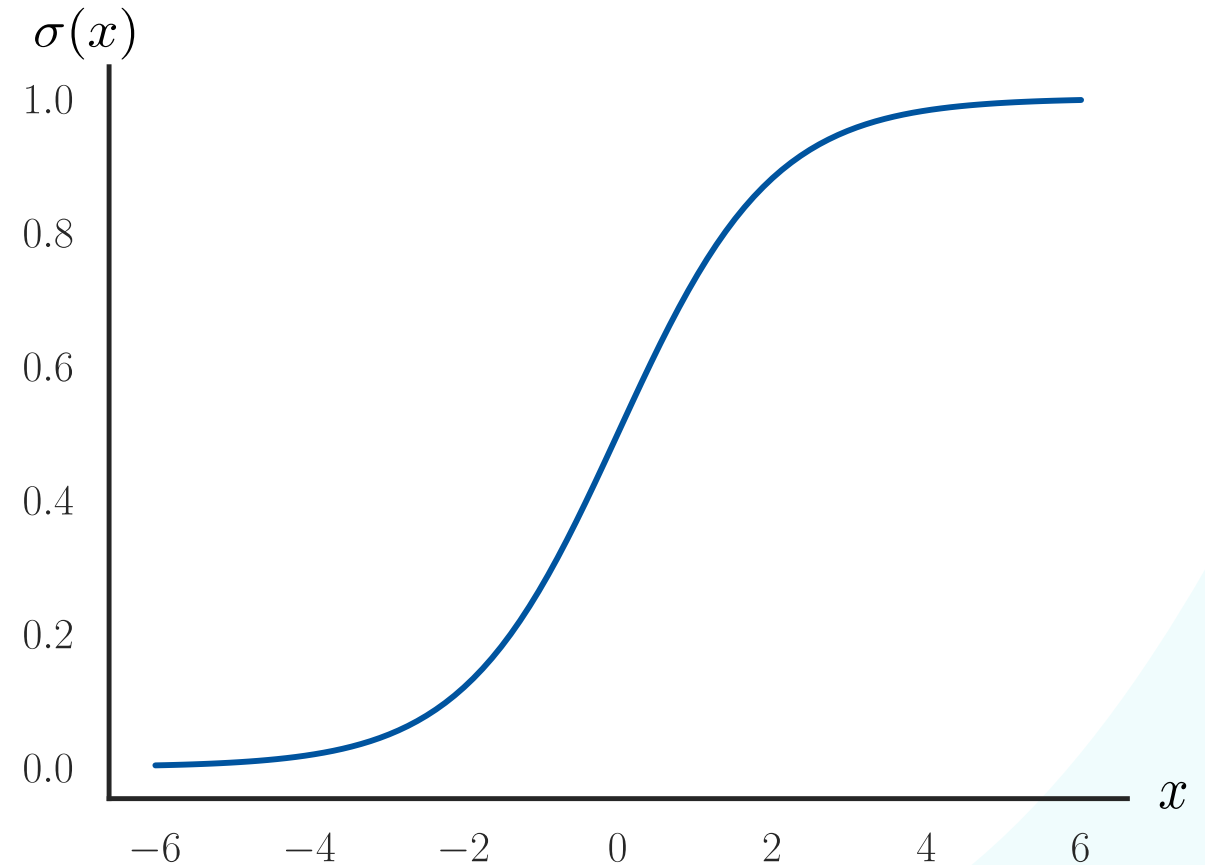
$$a = \ln \left(\frac{\sigma(a)}{1 - \sigma(a)} \right)$$

- Symmetry:

$$\sigma(-a) = 1 - \sigma(a)$$

- Derivative:

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



Logistic Regression

- We now define the **logistic regression** model
 - For the start, let us assume two classes $\mathcal{C}_1, \mathcal{C}_2$.
 - We model the class posteriors $p(\mathcal{C}_k|\mathbf{x})$ as

$$p(\mathcal{C}_1|\mathbf{x}) = y(\mathbf{x}) = \sigma(\mathbf{w}^\top \mathbf{x})$$

$$p(\mathcal{C}_2|\mathbf{x}) = 1 - p(\mathcal{C}_1|\mathbf{x})$$

Logistic sigmoid
activation function:

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



- I.e., we define a linear discriminant model

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

that is meant to represent the class posterior with the help of a logistic sigmoid activation function $\sigma(a)$.

- Our target labels are now $t_n \in \{0, 1\}$.

Error Function

- Consider a data set $\mathcal{D} = \{(\phi_1, t_1), \dots, (\phi_N, t_N)\}$
 - with data points $\Phi = [\phi_1, \dots, \phi_N]^\top$, $\phi_n = \phi(\mathbf{x}_n)$
 - And target labels $\mathbf{t} = [t_1, \dots, t_N]^\top$, $t_n \in \{0, 1\}$

- Maximum likelihood approach

- With $y_n = p(\mathcal{C}_1 | \phi_n) = \sigma(\mathbf{w}^\top \phi_n)$
- We model the probability of the target labels \mathbf{t} given our model parameters \mathbf{w} as

Trick: use t_n as an indicator variable

$$p(\mathbf{t} | \mathbf{w}) = \prod_{n=1}^N \begin{cases} p(\mathcal{C}_1 | \phi_n) & , t_n = 1 \\ p(\mathcal{C}_2 | \phi_n) & , t_n = 0 \end{cases} = \prod_{n=1}^N \begin{cases} y_n & , t_n = 1 \\ (1 - y_n) & , t_n = 0 \end{cases} = \prod_{n=1}^N y_n^{t_n} (1 - y_n)^{1-t_n}$$

- **Maximum likelihood** approach

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

- Define the error function as the **negative log-likelihood**

$$\begin{aligned} E(\mathbf{w}) &= -\ln p(\mathbf{t}|\mathbf{w}) \\ &= -\sum_{n=1}^N (t_n \ln y_n + (1 - t_n) \ln(1 - y_n)) \end{aligned}$$

- This function is known as the **binary cross-entropy error**.

Softmax Regression

- Multi-class extension of logistic regression
 - Generalization to K classes with target labels in 1-of- K notation $\mathbf{t}_n = [0, 1, \dots, 0]^\top$
 - Again, we define a linear discriminant function that models the class posteriors

$$\mathbf{y}(\mathbf{x}; \mathbf{w}) = \begin{bmatrix} p(\mathcal{C}_1 | \mathbf{x}; \mathbf{w}) \\ p(\mathcal{C}_2 | \mathbf{x}; \mathbf{w}) \\ \vdots \\ p(\mathcal{C}_K | \mathbf{x}; \mathbf{w}) \end{bmatrix} = \frac{1}{\sum_{j=1}^K \exp(\mathbf{w}_j^\top \mathbf{x})} \begin{bmatrix} \exp(\mathbf{w}_1^\top \mathbf{x}) \\ \exp(\mathbf{w}_2^\top \mathbf{x}) \\ \vdots \\ \exp(\mathbf{w}_K^\top \mathbf{x}) \end{bmatrix}$$

- This makes use of the **softmax** function as a multi-class extension of the logistic sigmoid

$$\text{softmax}(\mathbf{a}) = \frac{\exp(a_k)}{\sum_{j=1}^K \exp(a_j)}$$

- We can write the **binary cross-entropy error** as

$$\begin{aligned} E(\mathbf{w}) &= - \sum_{n=1}^N (t_n \ln y_n + (1 - t_n) \ln(1 - y_n)) \\ &= - \sum_{n=1}^N \sum_{k=0}^1 (\mathbb{I}(t_n = k) \ln p(\mathcal{C}_k | \mathbf{x}_n; \mathbf{w})) \end{aligned}$$

indicator function

$$\mathbb{I}(\varphi) = \begin{cases} 1 & \text{if } \varphi \text{ is true} \\ 0 & \text{otherwise} \end{cases}$$

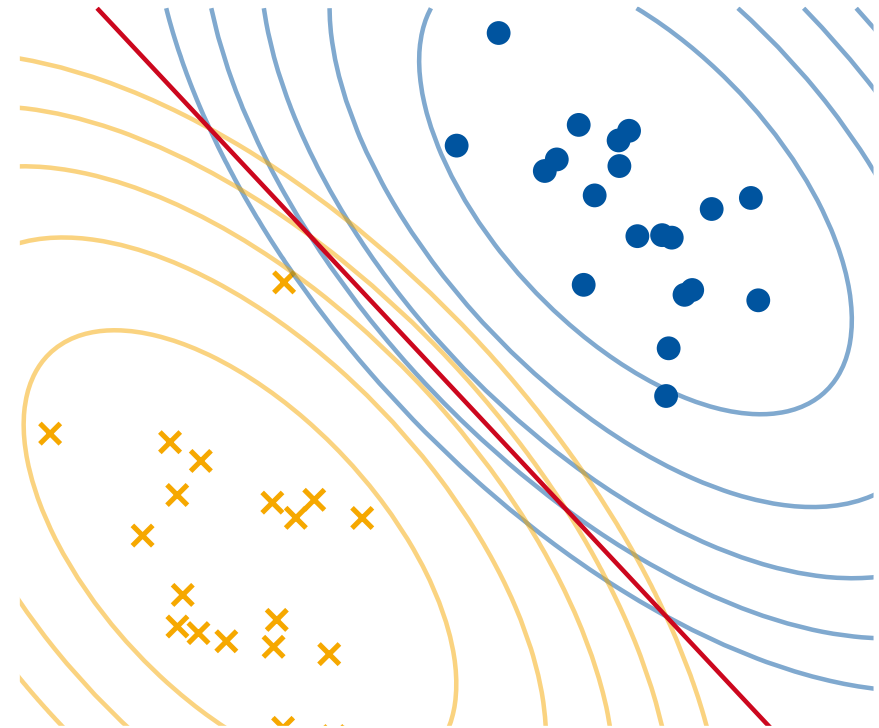
- Using one-hot labels \mathbf{t}_n , the generalization to K classes is:

$$E(\mathbf{w}) = - \sum_{n=1}^N \sum_{k=1}^K \left(\mathbb{I}(t_{kn} = 1) \ln \frac{\exp(\mathbf{w}_k^\top \mathbf{x})}{\sum_{j=1}^K \exp(\mathbf{w}_j^\top \mathbf{x})} \right)$$

- This function is known as the **multi-class cross-entropy error** or **softmax cross-entropy error**.

Logistic Regression

1. Logistic Regression Formulation
2. **Motivation and Background**
3. Iterative Optimization
4. First-Order Gradient Descent
5. Second-Order Gradient Descent
6. Error Function Analysis

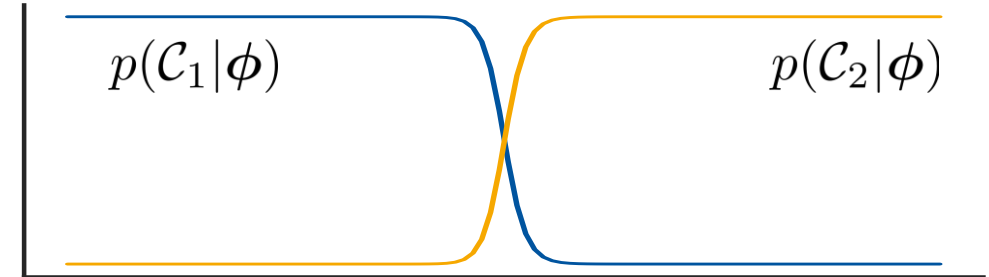


Motivation: Why Logistic Regression?

- Logistic Regression uses models of the form

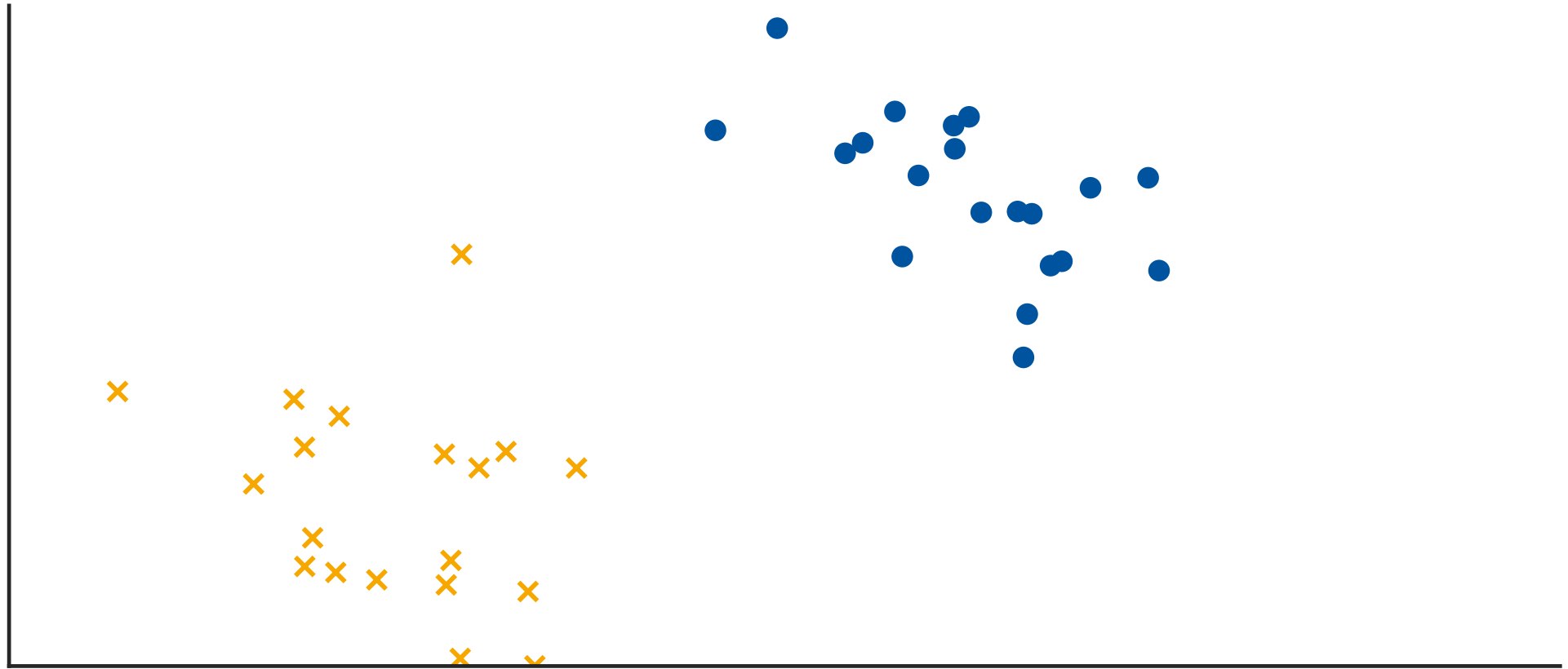
$$p(\mathcal{C}_1|\phi) = y(\phi) = \sigma(\mathbf{w}^\top \phi)$$

$$p(\mathcal{C}_2|\phi) = 1 - p(\mathcal{C}_1|\phi)$$



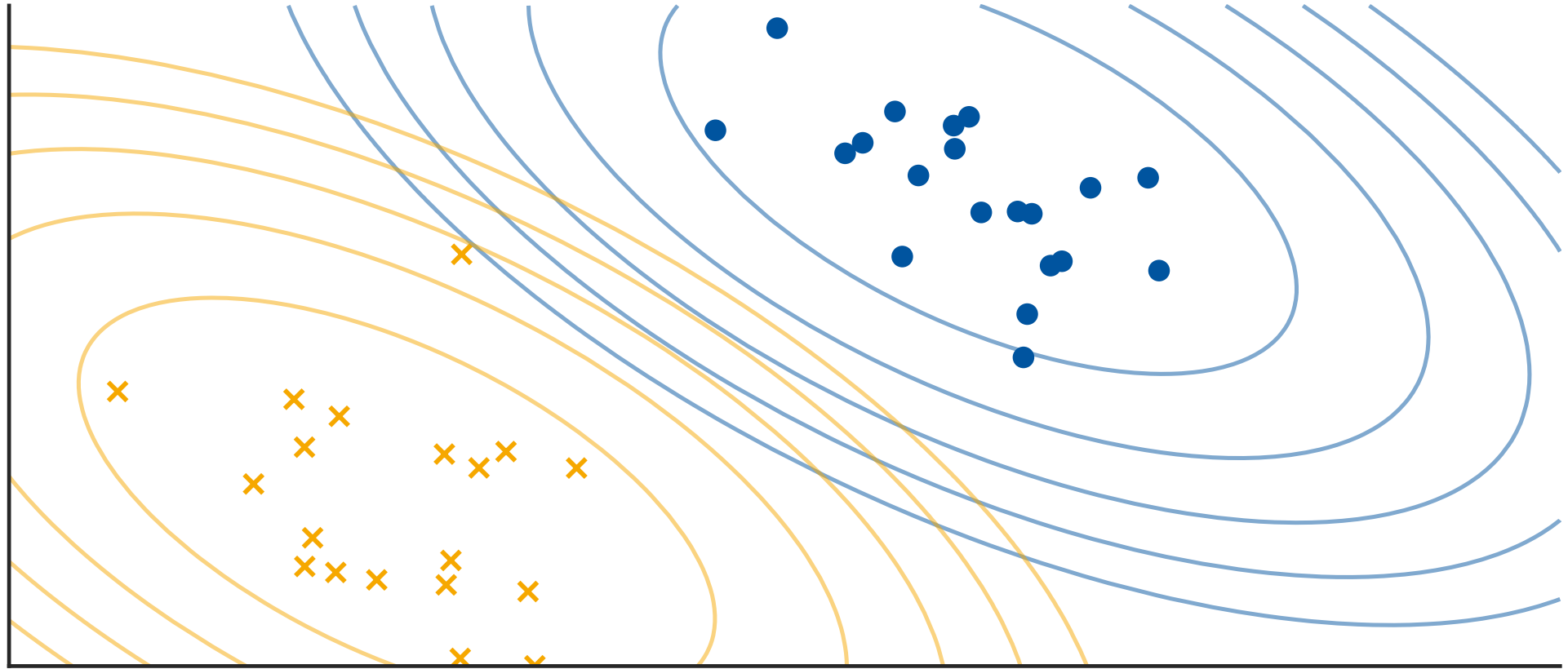
- Interpretation
 - We model the **class posteriors** $p(\mathcal{C}_k|\phi)$, as required to make Bayes optimal decisions.
 - We have seen previously that we can obtain $p(\mathcal{C}_k|\phi) = p(\phi|\mathcal{C}_k)p(\mathcal{C}_k)$.
 - However, here we model $p(\mathcal{C}_k|\phi)$ as a **linear discriminant function** $y(\phi) = \sigma(\mathbf{w}^\top \phi)$ instead.
- Why should we do this?*
 - What advantage does such a model have compared to direct modeling of the probabilities?*

Example



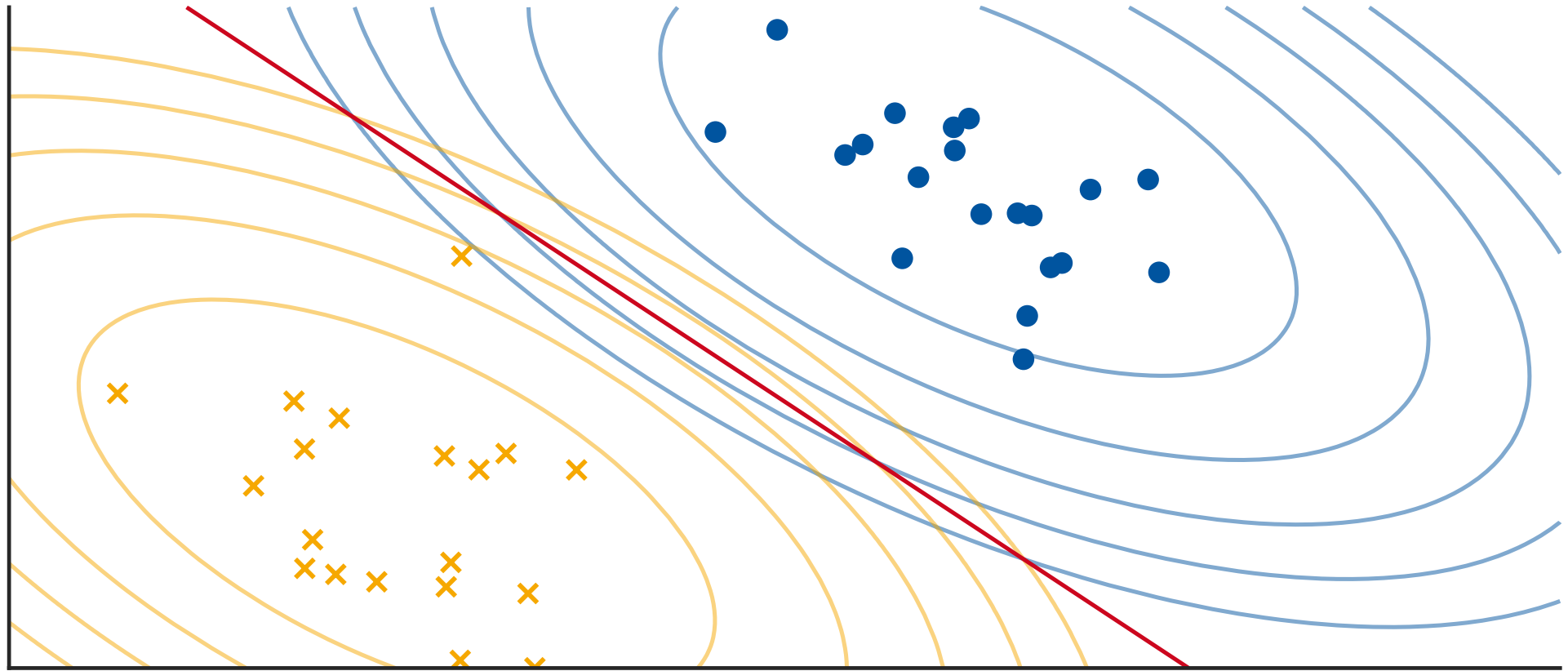
Let's assume the $p(\phi|\mathcal{C}_k)$ are modeled using Gaussians with equal covariances.

Example



Let's assume the $p(\phi|\mathcal{C}_k)$ are modeled using Gaussians with equal covariances.

Example



Let's assume the $p(\phi|\mathcal{C}_k)$ are modeled using Gaussians with equal covariances.

⇒ The decision boundary between them will be linear!

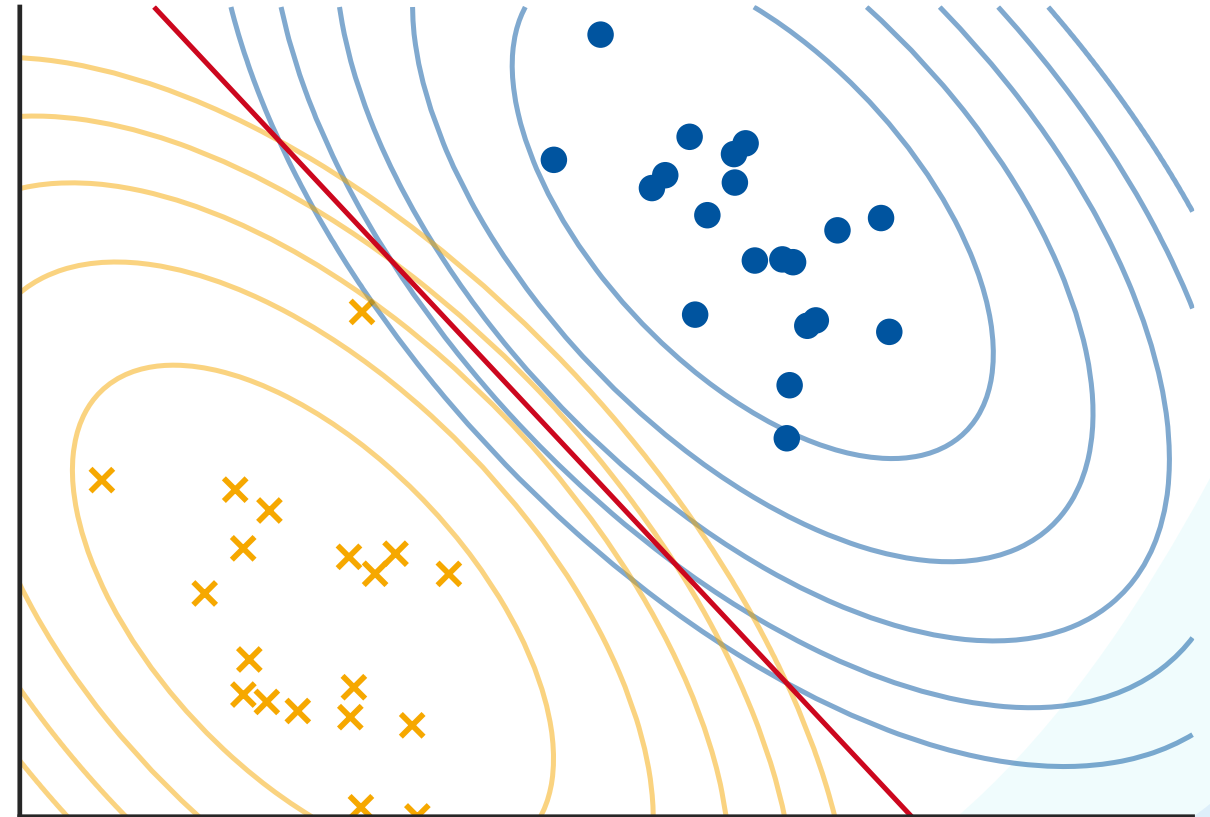
Parameter Efficiency

- #Parameters needed for generative models:
 - Assuming an M -dimensional feature space
 - Prior $p(\mathcal{C}_1)$ 1
 - Means μ_1, μ_2 $2M$
 - Covariances Σ $M(M + 1)/2$

⇒ Total $M(M + 5)/2 + 1$

- #Parameters needed for logistic regression:
 - Weights \mathbf{w} M

⇒ For large M , logistic regression has clear advantages!



Discussion: Logistic Regression

Advantages

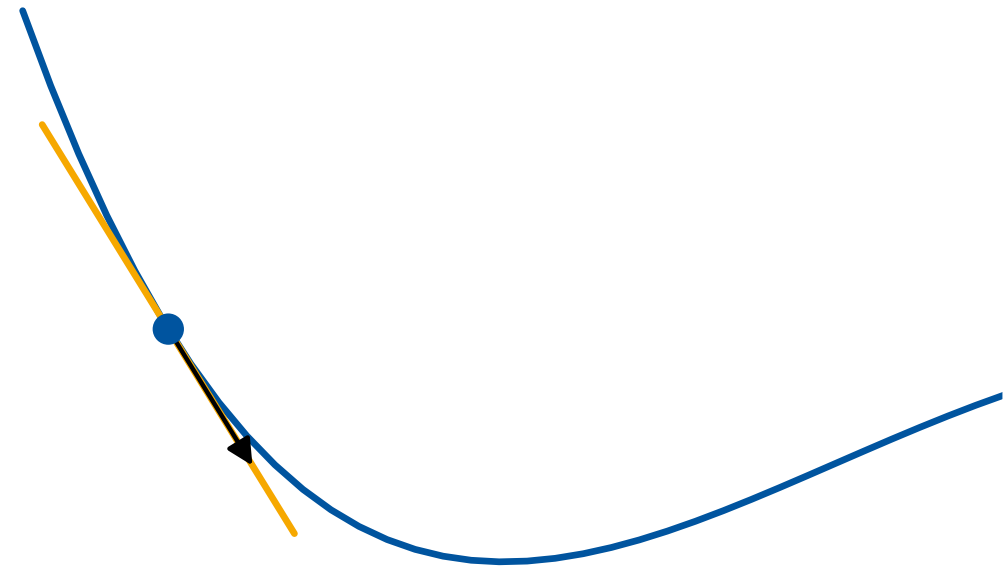
- Nice probabilistic interpretation, directly represents the posterior.
- Requires fewer parameters than modeling the likelihood + prior.
- Cross-Entropy error is convex: unique minimum exists.
- More robust than least-squares.

Limitations

- No closed-form solution, requires iterative optimization approach.

Logistic Regression

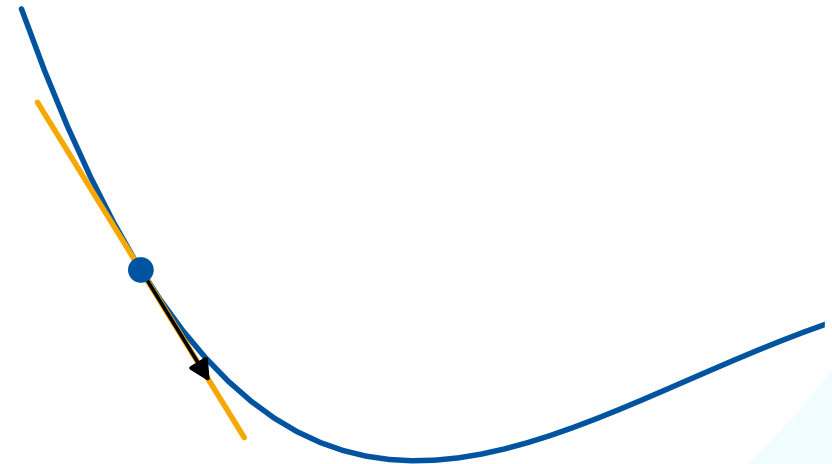
1. Logistic Regression Formulation
2. Motivation and Background
3. **Iterative Optimization**
4. First-Order Gradient Descent
5. Second-Order Gradient Descent
6. Error Function Analysis



Iterative Optimization

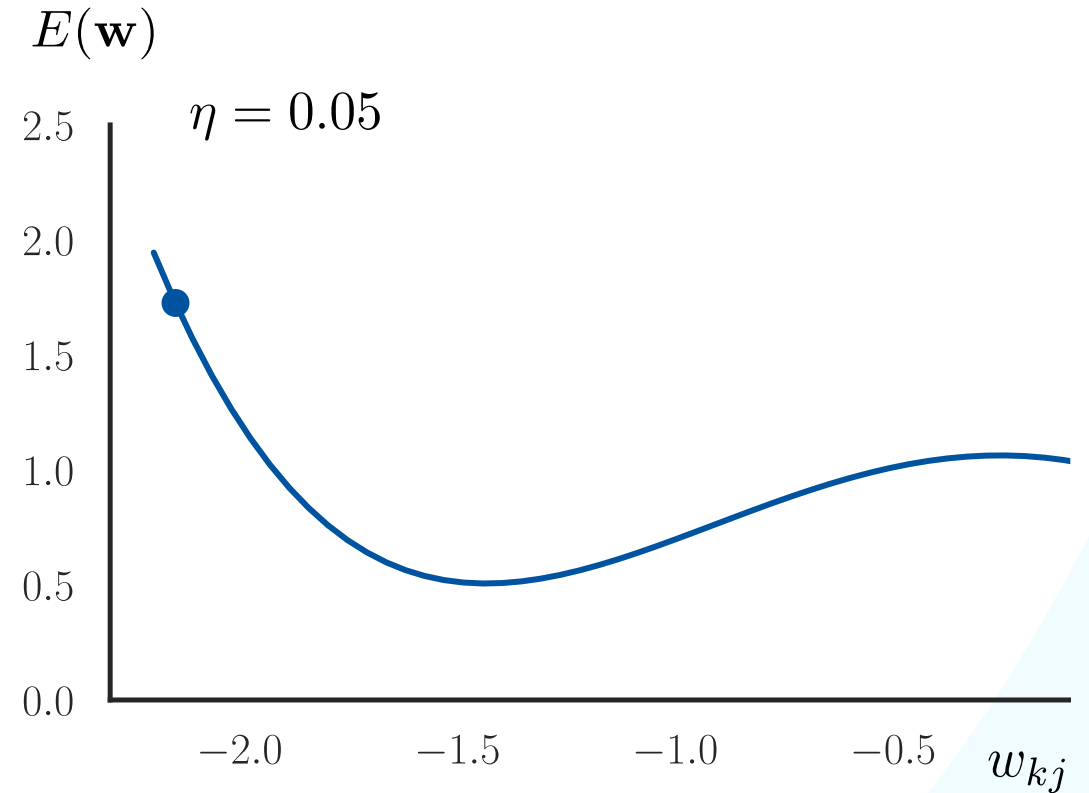
- In general, generalized linear discriminants with nonlinear activation and/or basis functions can no longer be optimized in closed form.
- Instead, we use iterative optimization schemes.
- Here: [Gradient Descent](#).
 - Start with initial guess for parameter values.
 - Move towards a minimum of the error function by following the direction of steepest descent.
 - Iterate until convergence

$$y_k(\mathbf{x}) = g \left(\sum_{j=0}^M w_{kj} \phi_j(\mathbf{x}) \right) = g(\mathbf{w}^T \phi(\mathbf{x}))$$



Idea: Gradient Descent

- Start with an initial guess of parameter values $w_{kj}^{(0)}$.

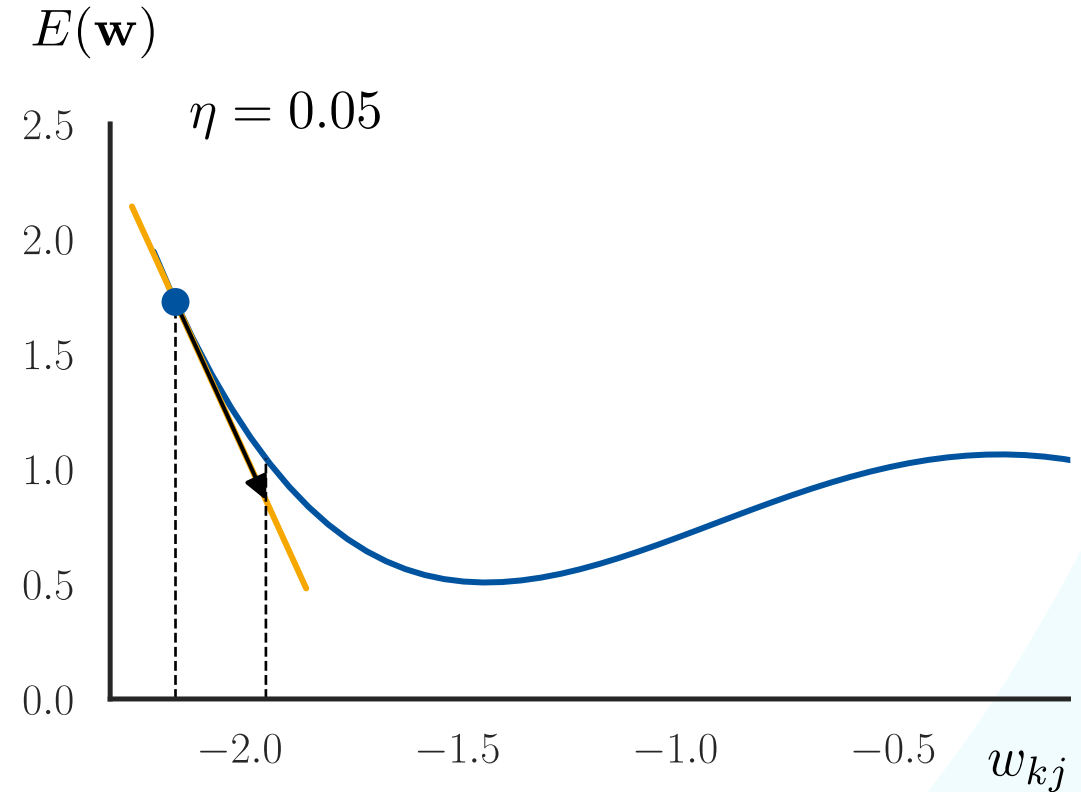


Idea: Gradient Descent

- Start with an initial guess of parameter values $w_{kj}^{(0)}$.
- Follow the gradient to move to a (local) minimum:

$$w_{kj}^{(\tau+1)} = w_{kj}^{(\tau)} - \eta \left. \frac{\partial E(\mathbf{w})}{\partial w_{kj}} \right|_{\mathbf{w}^{(\tau)}}$$

- η is called the **learning rate**.
- This corresponds to a 1st-order Taylor expansion.
 - I.e., we approximate the error function by its tangent plane around the current point $\mathbf{w}^{(\tau)}$.
- Repeat this procedure for a number of steps.

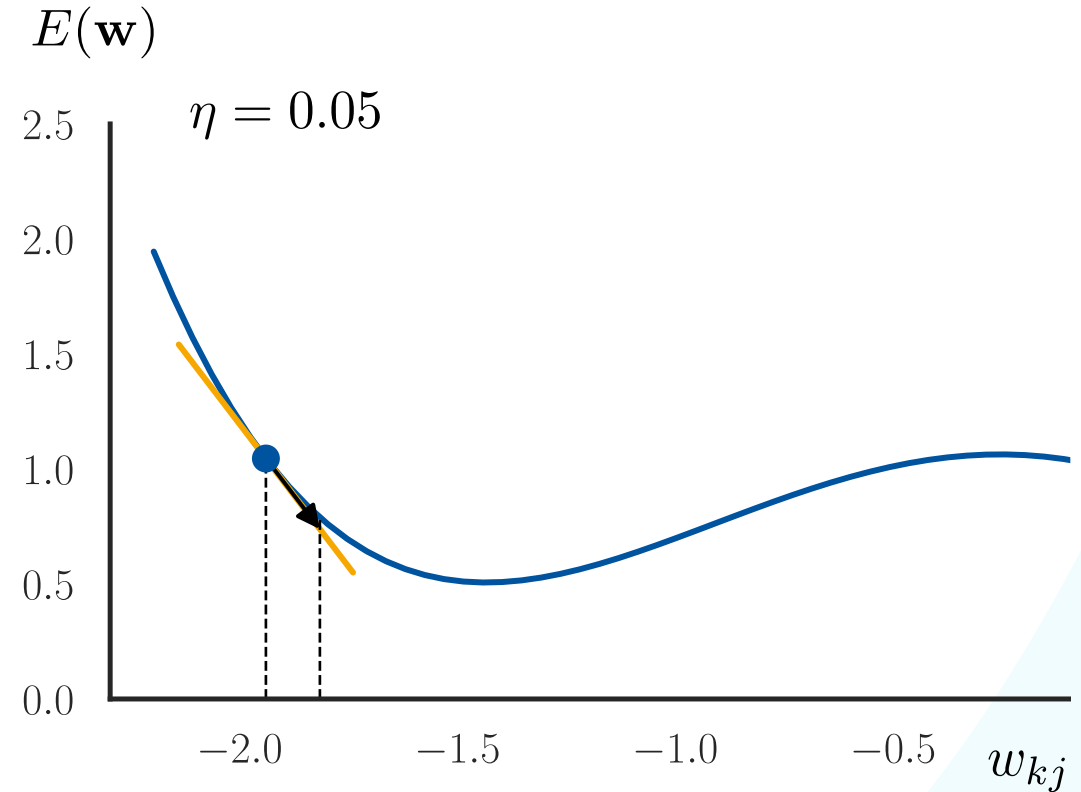


Idea: Gradient Descent

- Start with an initial guess of parameter values $w_{kj}^{(0)}$.
- Follow the gradient to move to a (local) minimum:

$$w_{kj}^{(\tau+1)} = w_{kj}^{(\tau)} - \eta \left. \frac{\partial E(\mathbf{w})}{\partial w_{kj}} \right|_{\mathbf{w}^{(\tau)}}$$

- η is called the **learning rate**.
- This corresponds to a 1st-order Taylor expansion.
 - I.e., we approximate the error function by its tangent plane around the current point $\mathbf{w}^{(\tau)}$.
- Repeat this procedure for a number of steps.

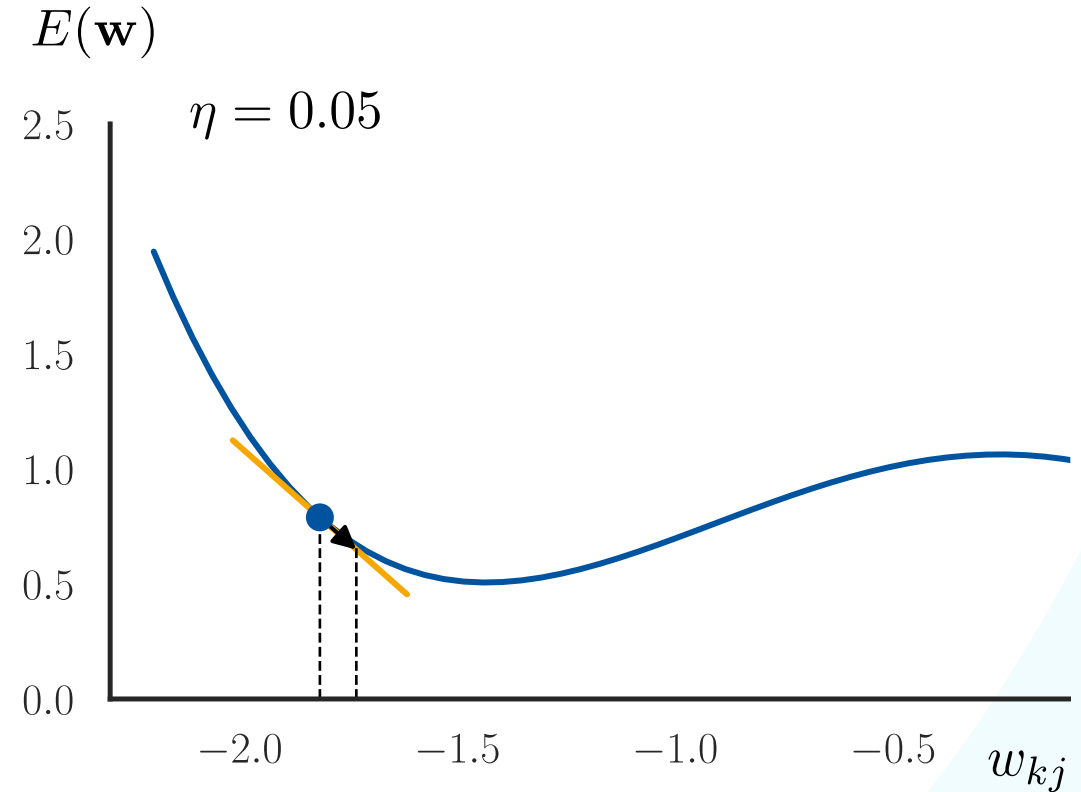


Idea: Gradient Descent

- Start with an initial guess of parameter values $w_{kj}^{(0)}$.
- Follow the gradient to move to a (local) minimum:

$$w_{kj}^{(\tau+1)} = w_{kj}^{(\tau)} - \eta \left. \frac{\partial E(\mathbf{w})}{\partial w_{kj}} \right|_{\mathbf{w}^{(\tau)}}$$

- η is called the **learning rate**.
- This corresponds to a 1st-order Taylor expansion.
 - I.e., we approximate the error function by its tangent plane around the current point $\mathbf{w}^{(\tau)}$.
- Repeat this procedure for a number of steps.

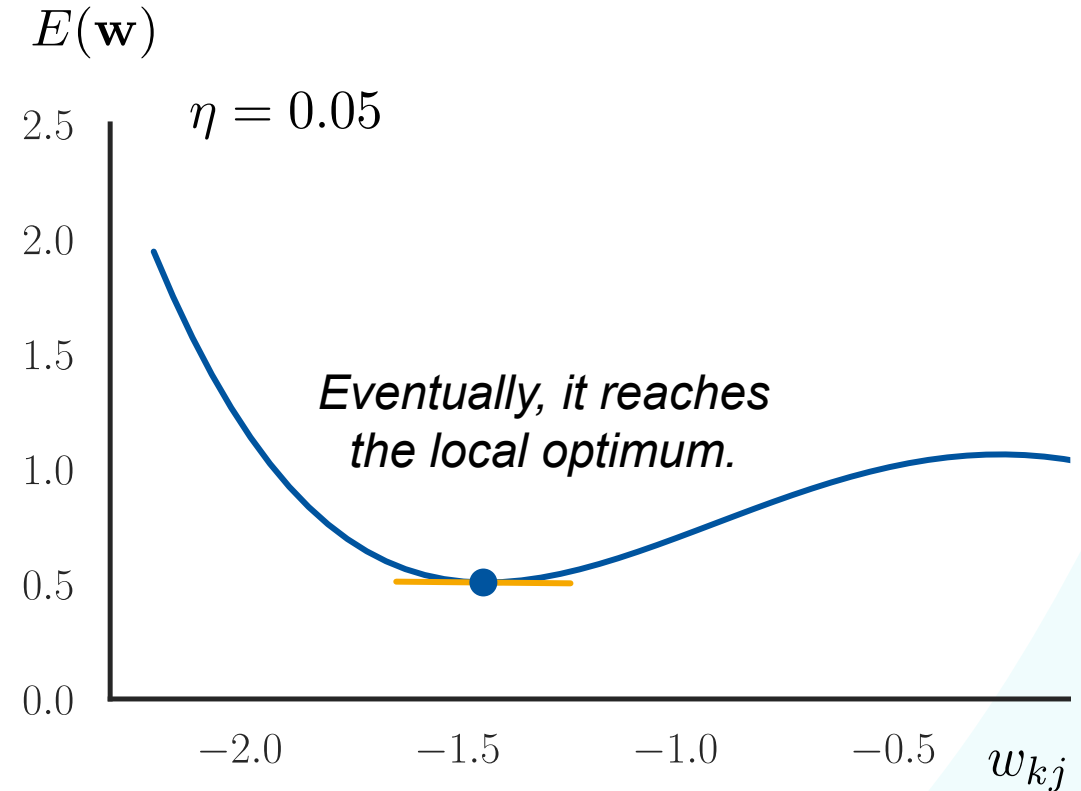


Idea: Gradient Descent

- Start with an initial guess of parameter values $w_{kj}^{(0)}$.
- Follow the gradient to move to a (local) minimum:

$$w_{kj}^{(\tau+1)} = w_{kj}^{(\tau)} - \eta \left. \frac{\partial E(\mathbf{w})}{\partial w_{kj}} \right|_{\mathbf{w}^{(\tau)}}$$

- η is called the **learning rate**.
- This corresponds to a 1st-order Taylor expansion.
 - I.e., we approximate the error function by its tangent plane around the current point $\mathbf{w}^{(\tau)}$.
- Repeat this procedure for a number of steps.



Discussion: Gradient Descent

Advantages

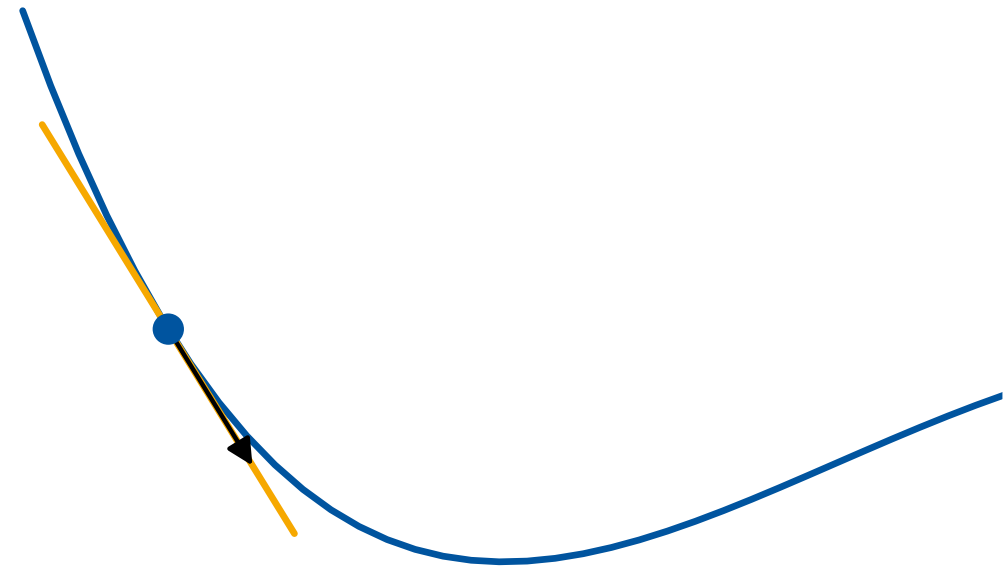
- Simple approach for iterative optimization.
- Approximates the error function by its tangent plane around the current point in order to find the direction of steepest descent.

Limitations

- Local optimization. Unless the error function is convex, will only converge to a local optimum.
- Relatively slow convergence (can be improved by second-order approaches).
- In practice, finding a good step size (**learning rate**) is important for fast convergence.

Logistic Regression

1. Logistic Regression Formulation
2. Motivation and Background
3. Iterative Optimization
4. **First-Order Gradient Descent**
5. Second-Order Gradient Descent
6. Error Function Analysis



First-order Optimization

- Logistic regression uses the **binary cross-entropy error**:

$$E(\mathbf{w}) = - \sum_{n=1}^N (t_n \ln y(\mathbf{x}_n; \mathbf{w}) + (1 - t_n) \ln(1 - y(\mathbf{x}_n; \mathbf{w})))$$

- Properties
 - Convex function, so it has a unique minimum
 - But no closed-form solution
- We need to use iterative methods for optimization
 - Let's try (first-order) gradient descent:

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E(\mathbf{w})$$

Gradient of the Cross-Entropy Error

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

$$\begin{aligned} \nabla E(\mathbf{w}) &= - \sum_{n=1}^N \left(t_n \frac{\frac{\partial}{\partial \mathbf{w}} y_n}{y_n} + (1 - t_n) \frac{\frac{\partial}{\partial \mathbf{w}} (1 - y_n)}{(1 - y_n)} \right) \\ &= - \sum_{n=1}^N \left(t_n \frac{\cancel{y_n} (1 - y_n)}{\cancel{y_n}} \phi_n + (1 - t_n) \frac{\cancel{y_n} (1 - \cancel{y_n})}{\cancel{(1 - y_n)}} \phi_n \right) \\ &= - \sum_{n=1}^N ((t_n - \cancel{t_n y_n} - y_n + \cancel{t_n y_n}) \phi_n) \\ &= \sum_{n=1}^N (y_n - t_n) \phi_n \end{aligned}$$

$$y_n = y(\mathbf{x}_n; \mathbf{w})$$

$$\begin{aligned} \sigma'(a) &= \sigma(a)(1 - \sigma(a)) \\ \frac{\partial y_n}{\partial \mathbf{w}} &= y_n(1 - y_n) \phi_n \end{aligned}$$

$$\phi_n = \phi(\mathbf{x}_n)$$

- The gradient for logistic regression is

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n$$

- We can plug this into gradient descent:

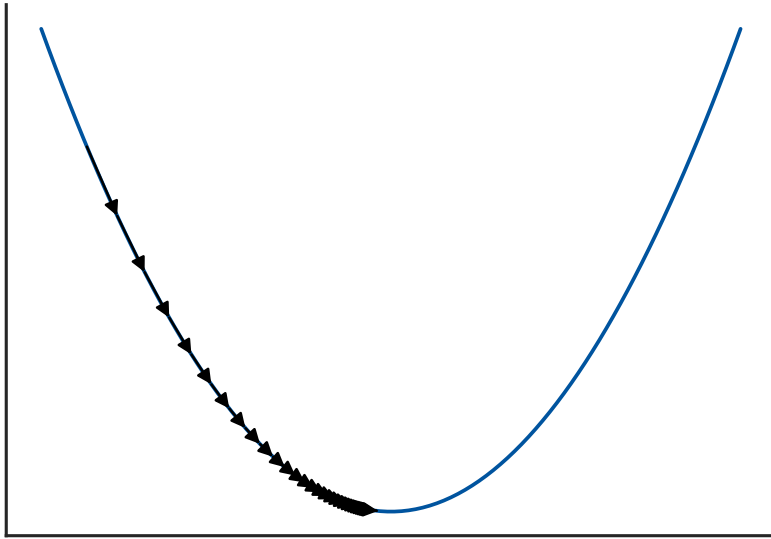
$$\begin{aligned} \mathbf{w}^{(\tau+1)} &= \mathbf{w}^{(\tau)} - \eta \nabla E(\mathbf{w}) \\ &= \mathbf{w}^{(\tau)} - \eta \sum_{n=1}^N (y_n - t_n) \phi_n \end{aligned}$$

How should we choose the learning rate?

- This update rule is known as the **Delta rule** (= **LMS rule**)
 - *Simply feed back the input data points, weighted by the classification error.*

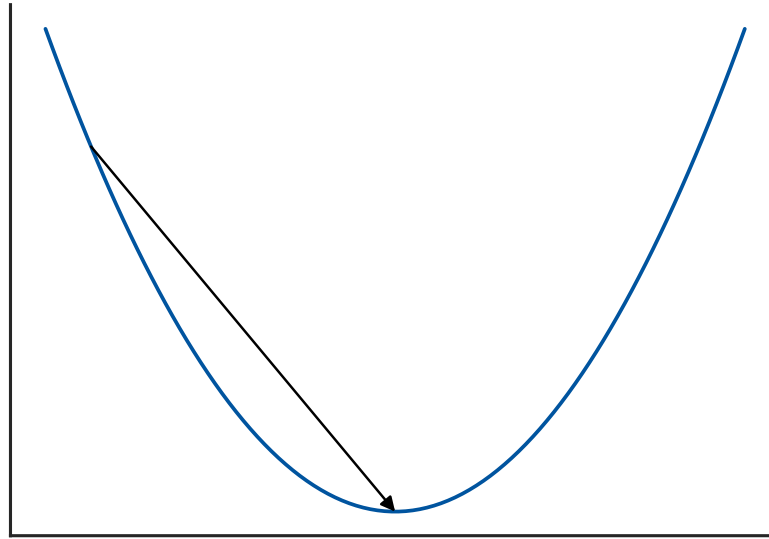
Effects of the learning rate

η too small



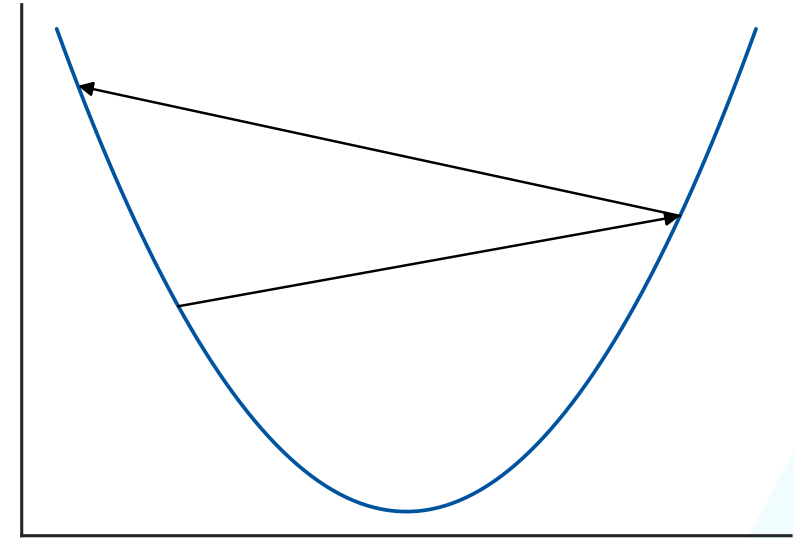
Convergence is slow

η_{opt}



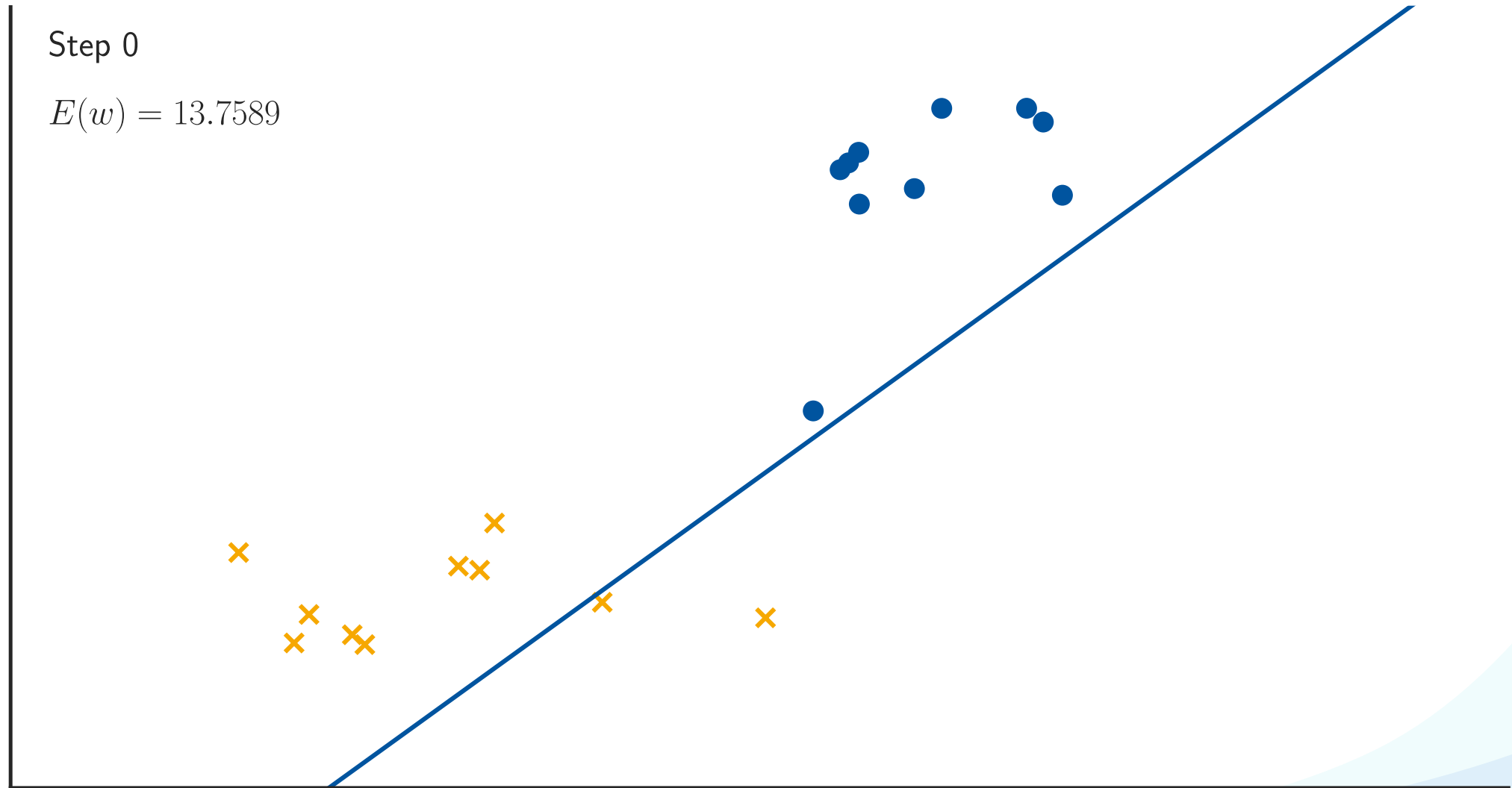
Converges ideally in
a single step

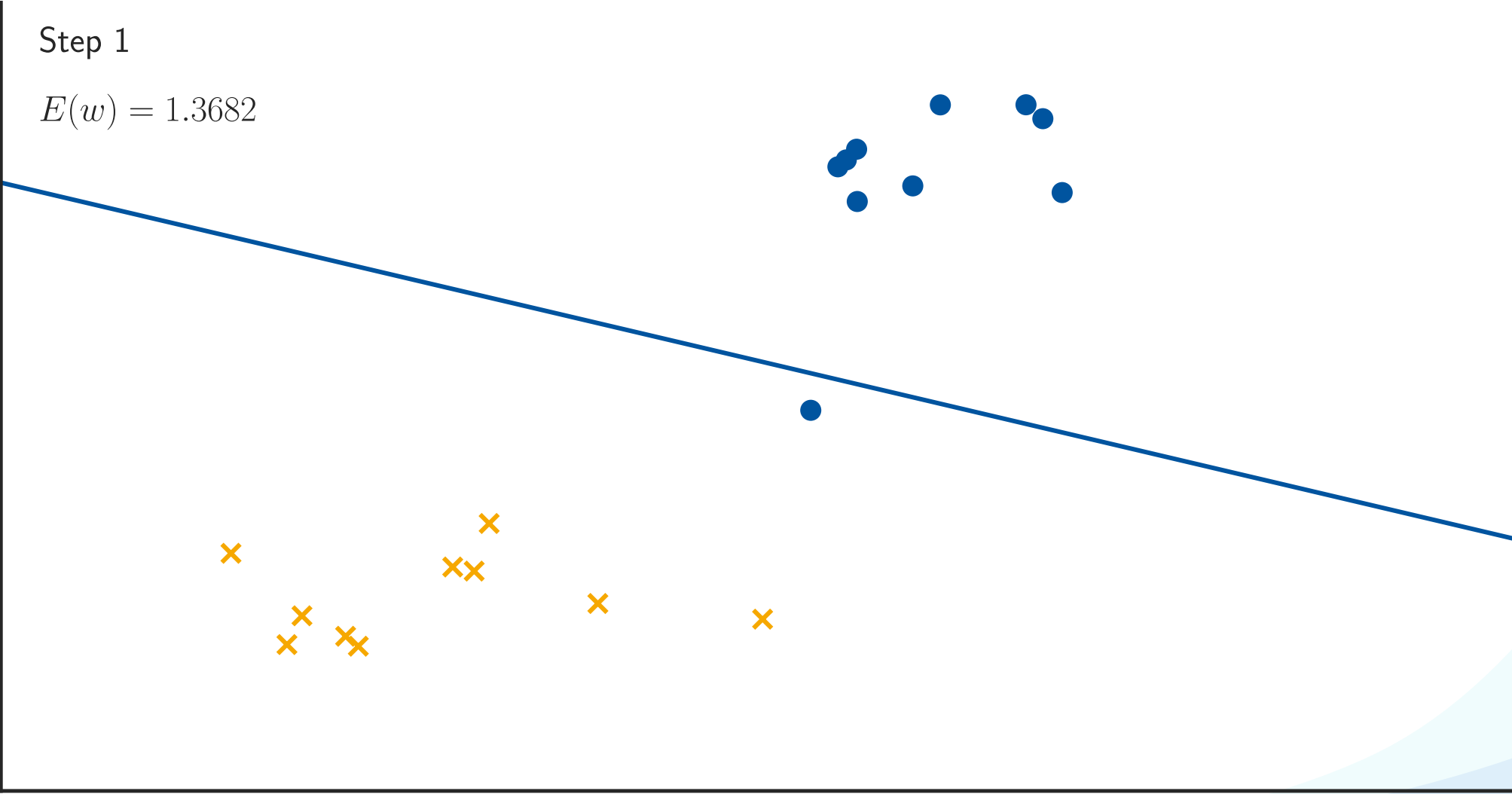
η too large

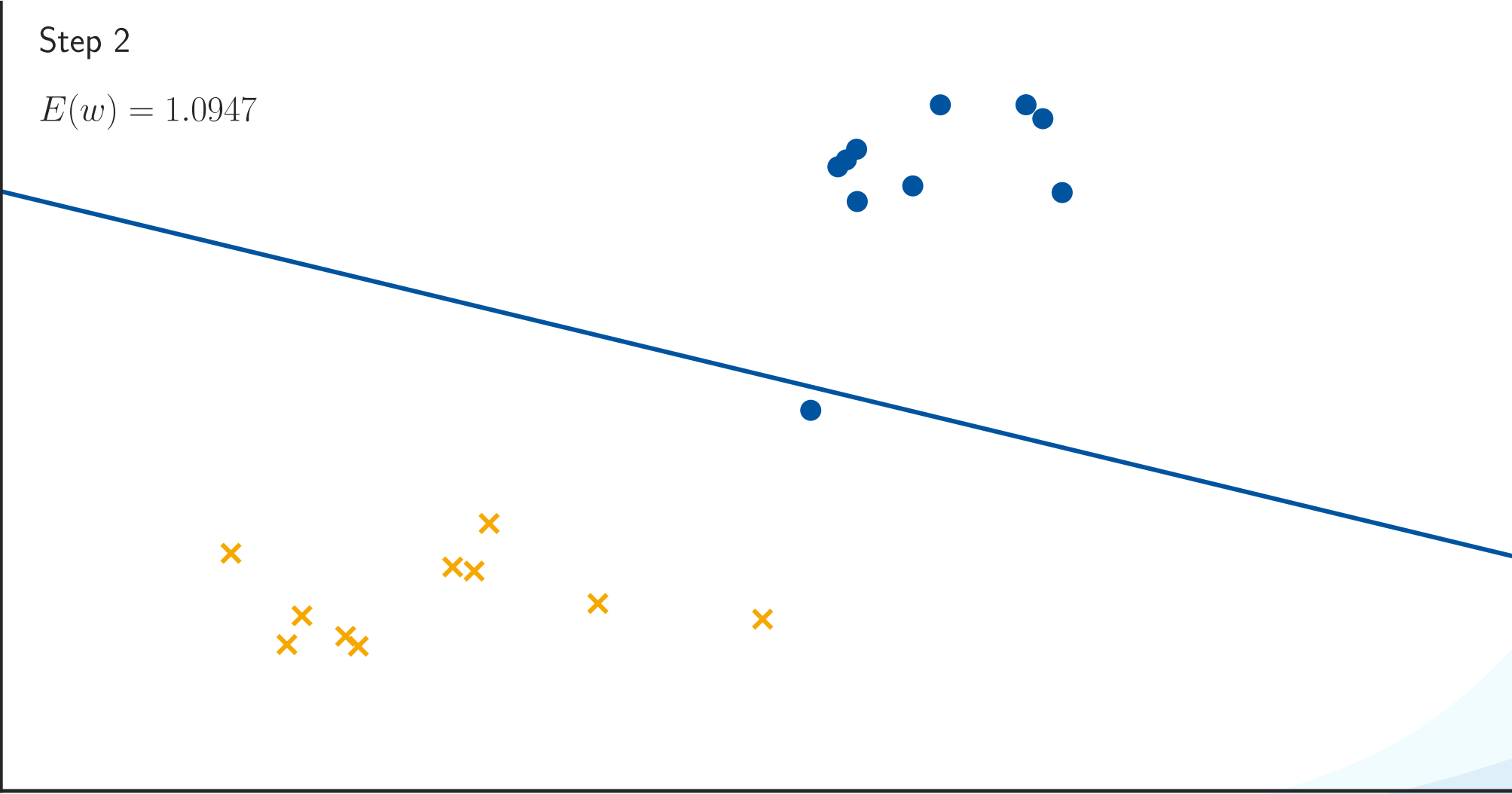


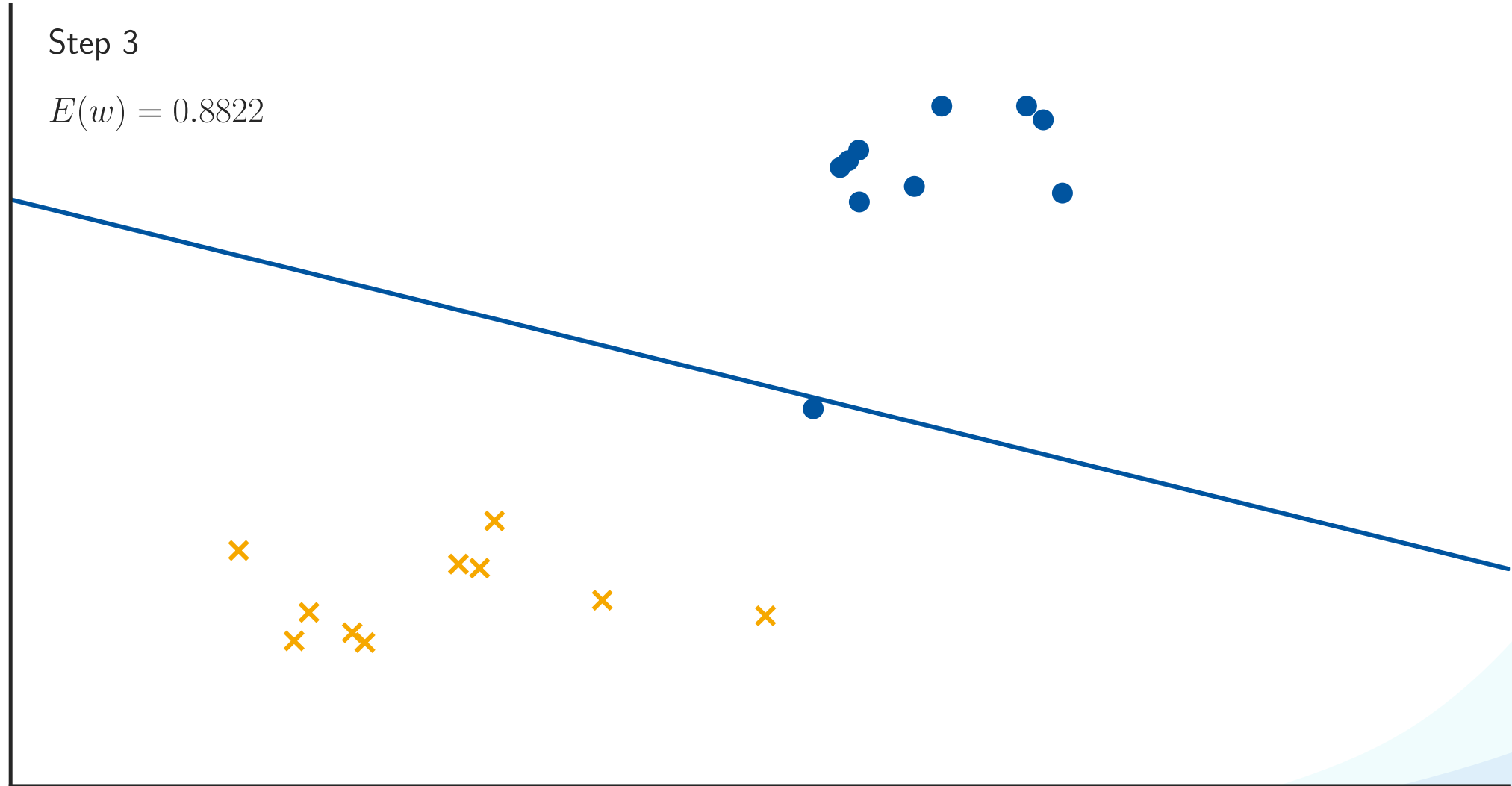
Might not converge

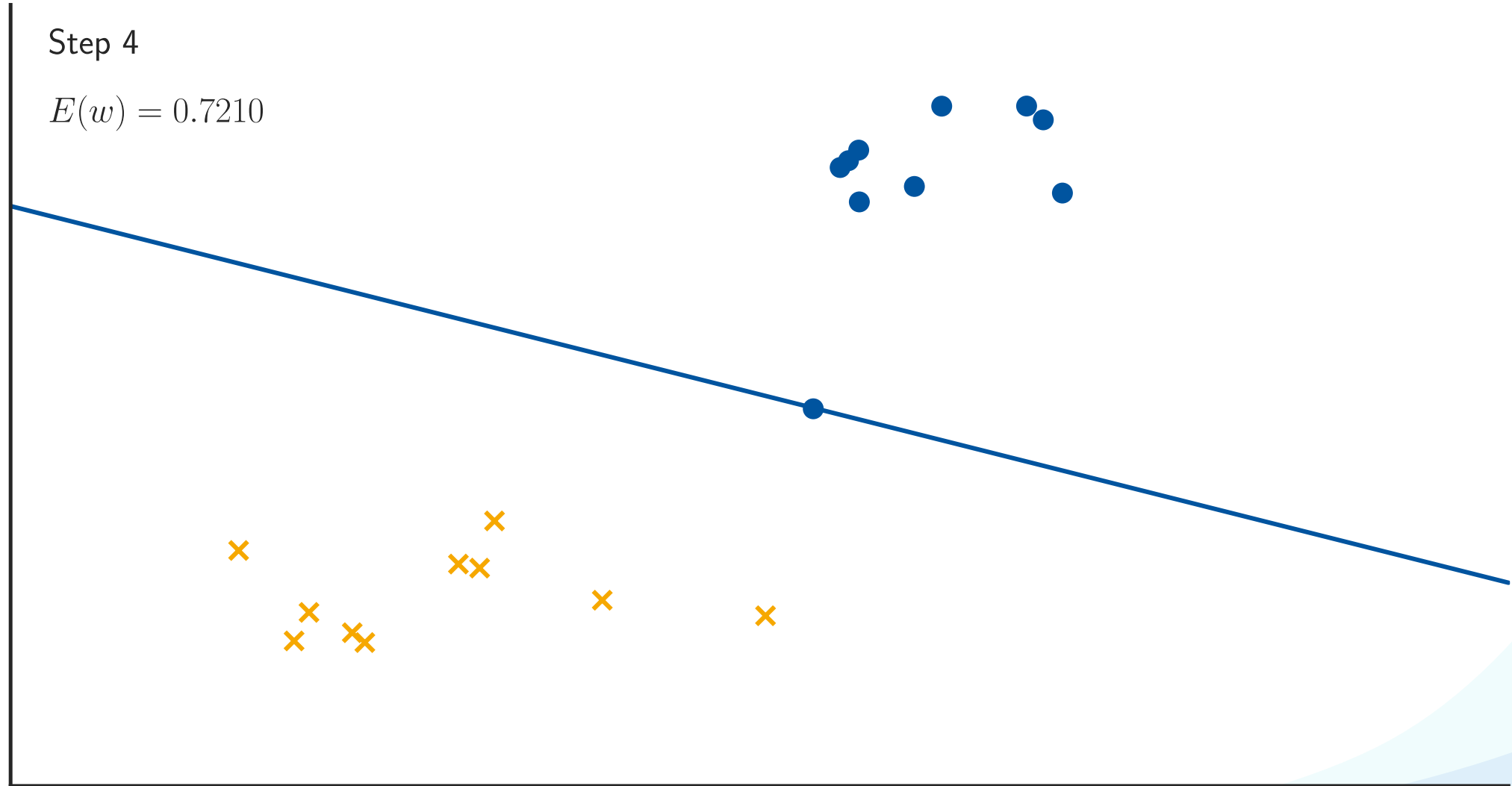
Example: Logistic Regression with Gradient Descent

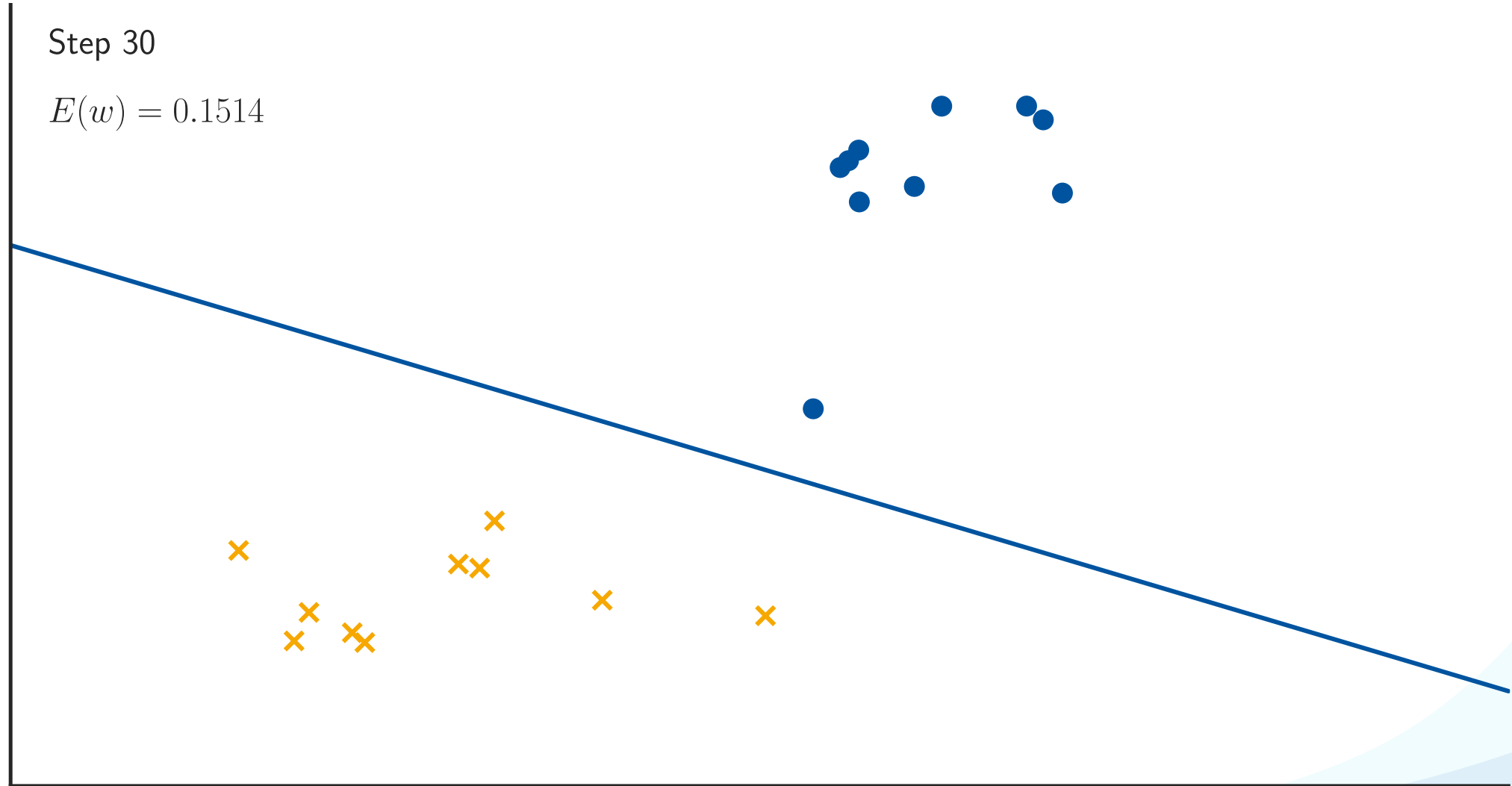


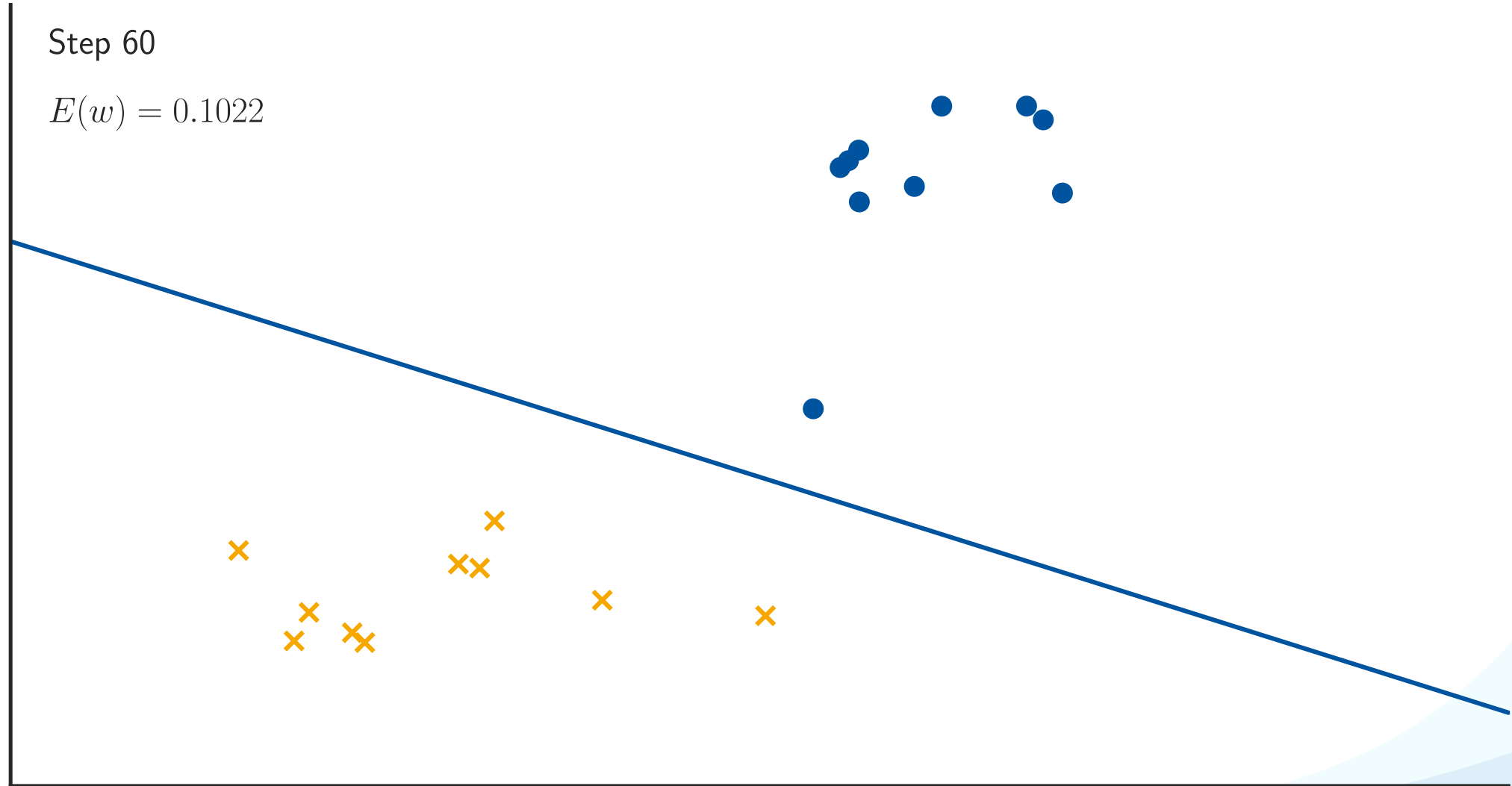


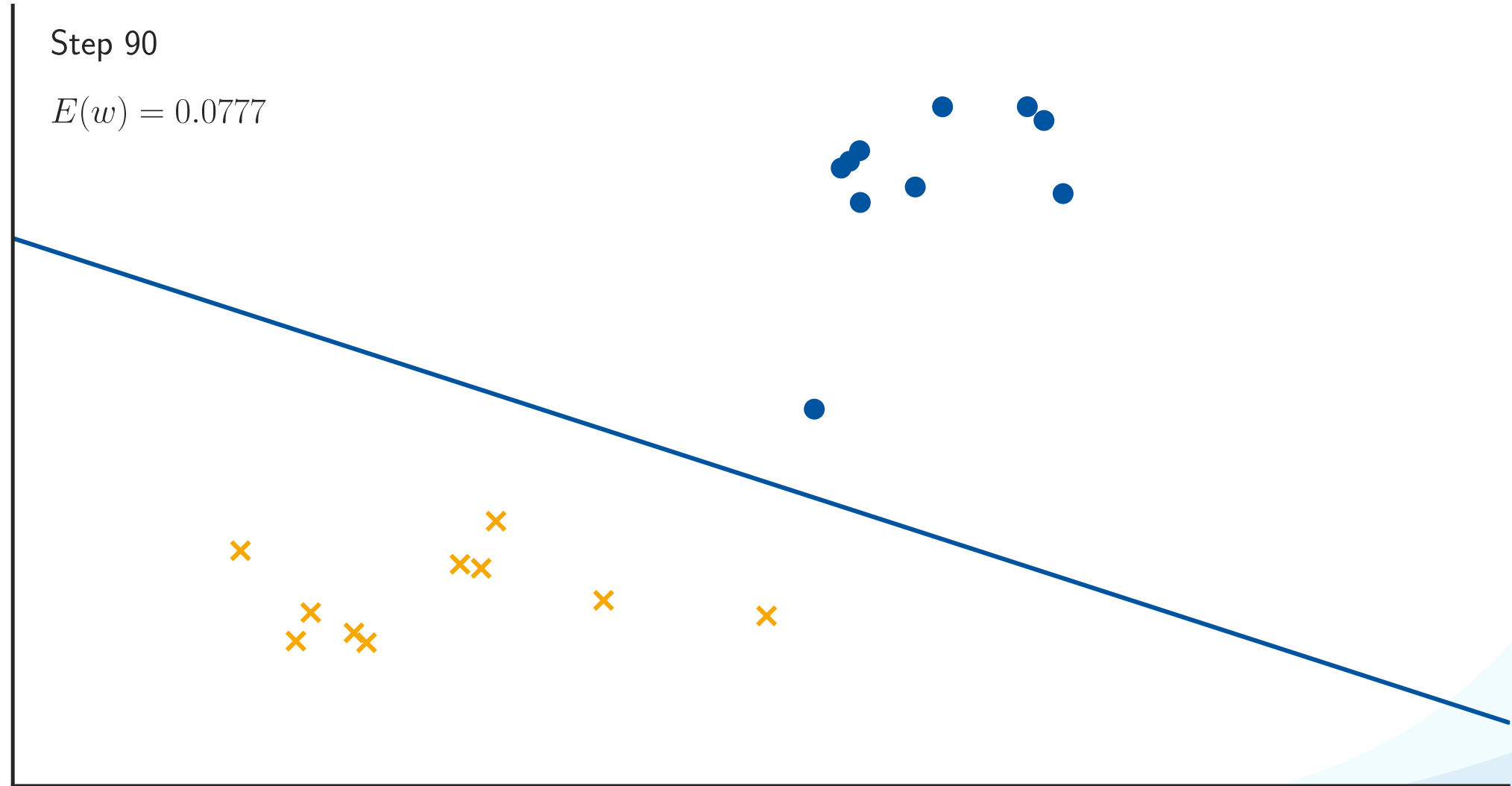












Discussion: Logistic Regression with Gradient Descent

Advantages

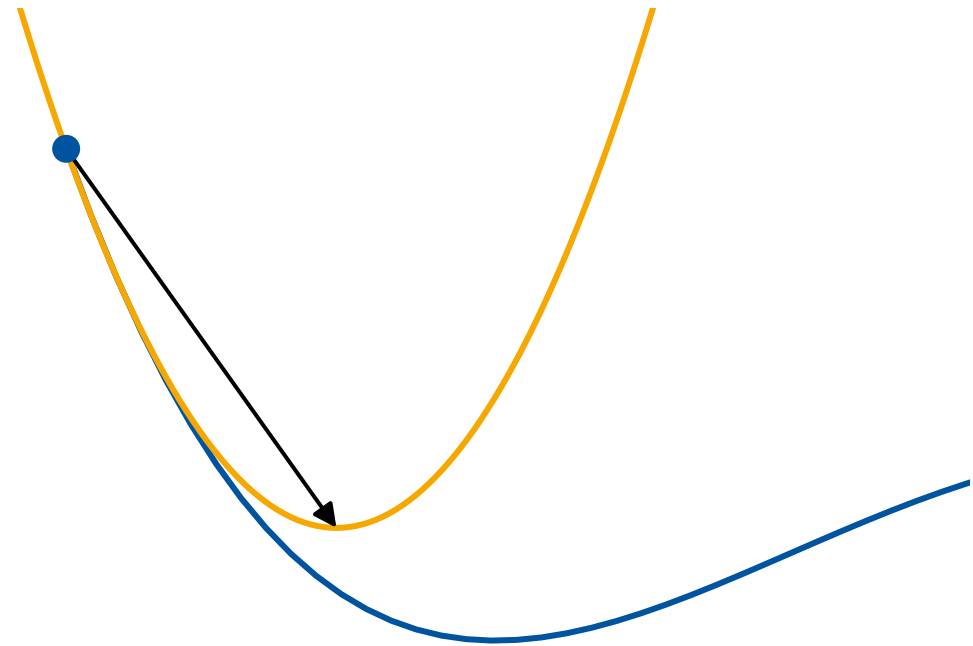
- Simple iterative optimization scheme with a familiar update rule ([Delta rule](#)).

Limitations

- Slow convergence
- Need to choose a suitable learning rate.

Logistic Regression

1. Logistic Regression Formulation
2. Motivation and Background
3. Iterative Optimization
4. First-Order Gradient Descent
5. **Second-Order Gradient Descent**
6. Error Function Analysis

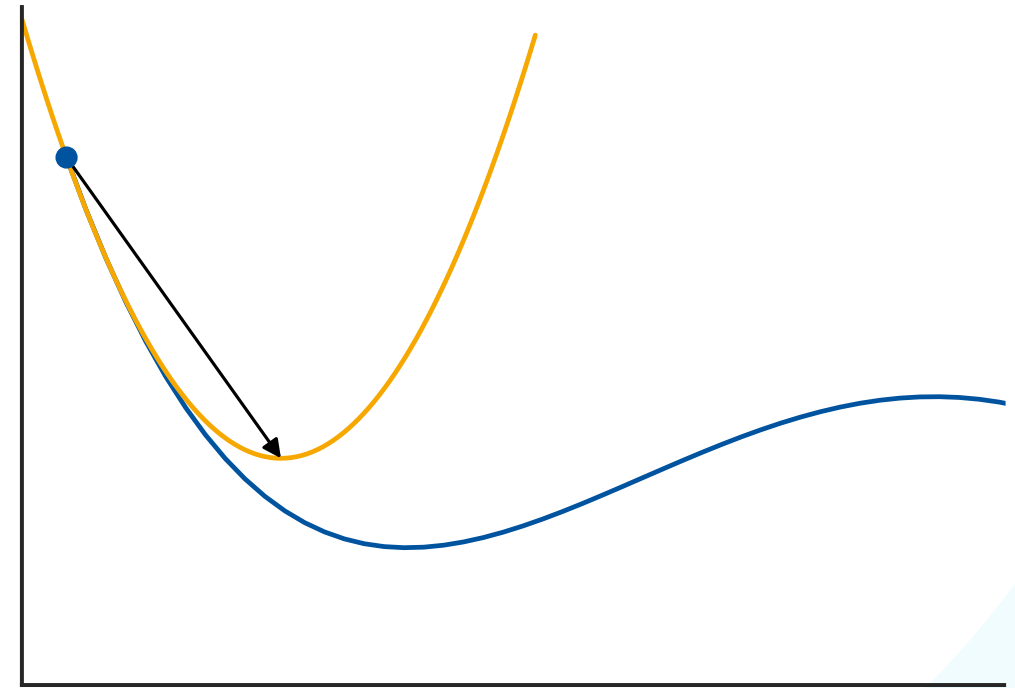


Second-Order Optimization

- So far, we have optimized the cross-entropy error with gradient descent:

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E(\mathbf{w})$$

- This is a first-order approximation, and it heavily depends on the learning rate η .
- Instead, we can apply a second-order optimization scheme that converges faster and is independent of the learning rate.



Newton-Raphson Gradient Descent

- Second-order **Newton-Raphson** update scheme:

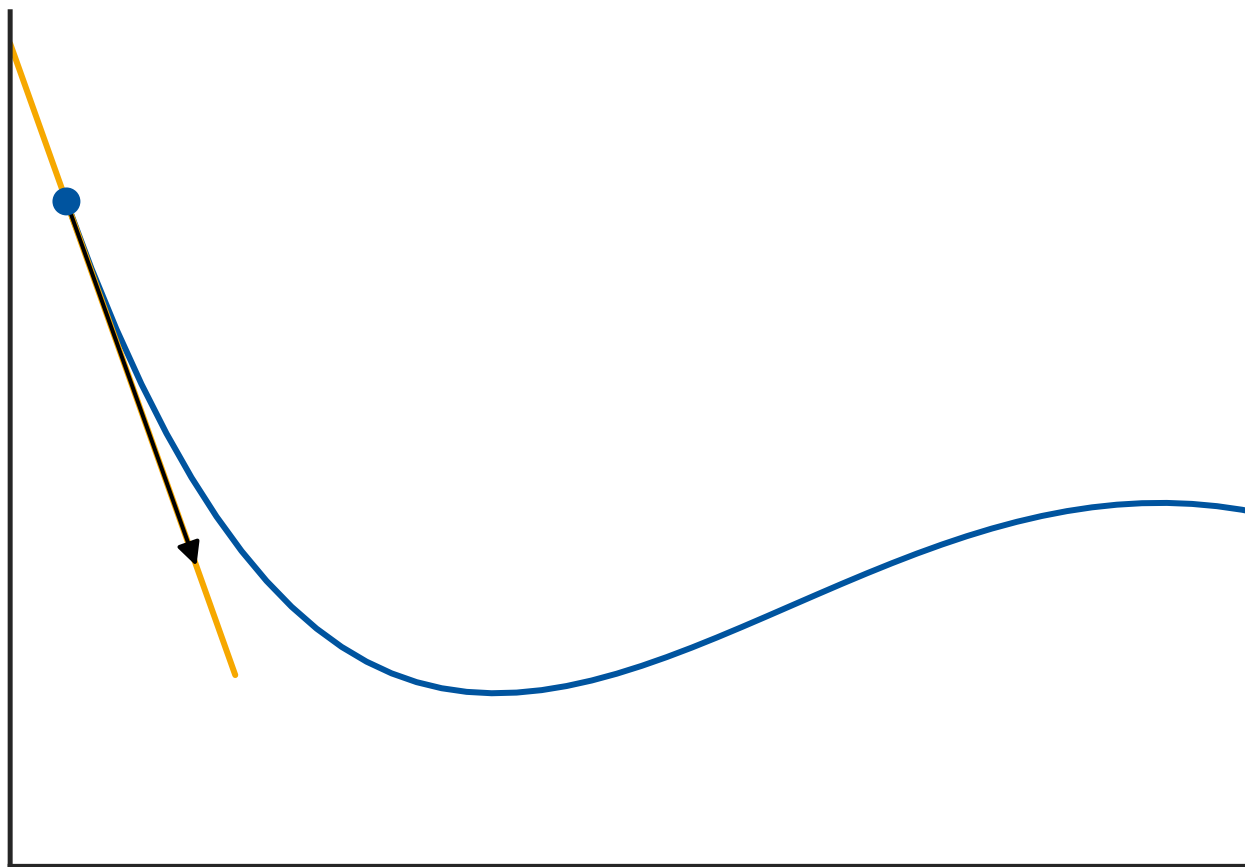
$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \mathbf{H}^{-1} \nabla E(\mathbf{w})$$

- Here, $\mathbf{H} = \nabla \nabla E(\mathbf{w})$ is the Hessian matrix, i.e., the matrix of second derivatives:

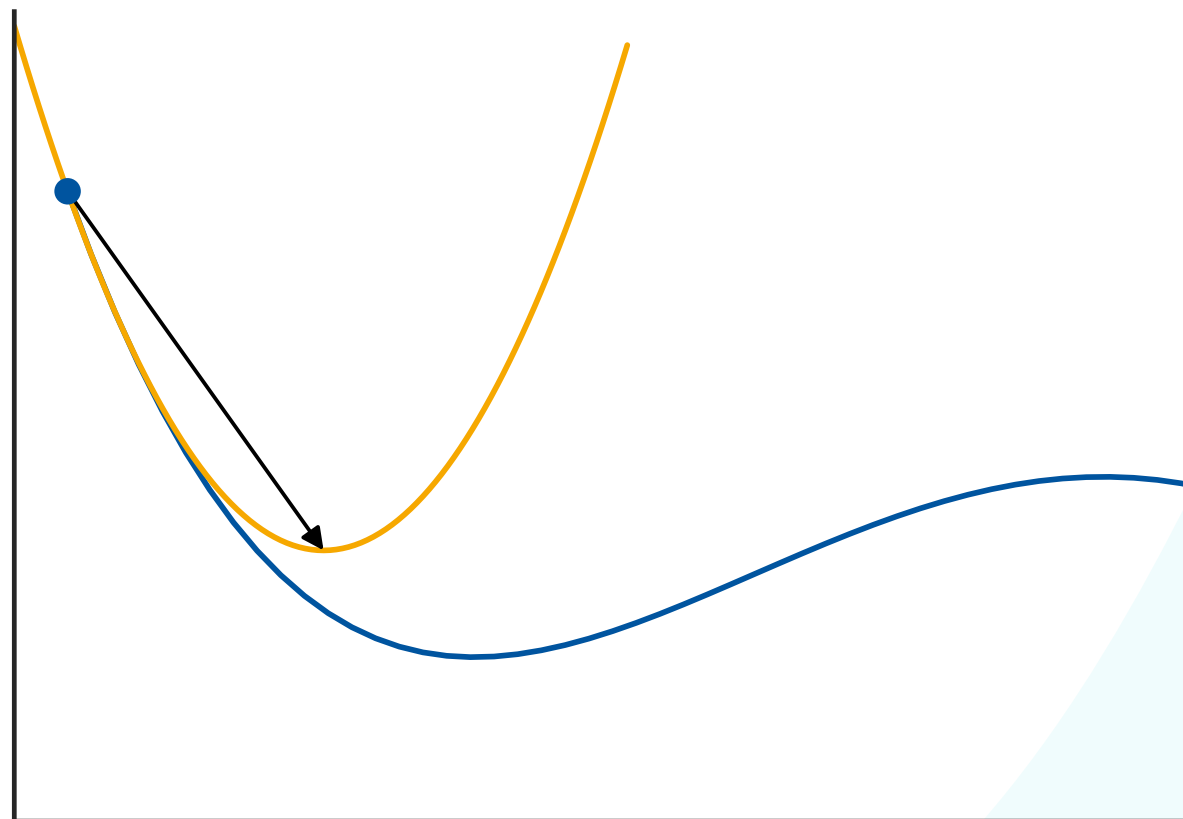
$$\mathbf{H}_{ij} = \frac{\partial^2 E(\mathbf{w})}{\partial w_i \partial w_j}$$

- Properties
 - Local quadratic approximation
 - Much faster convergence by taking into account the curvature of the error function.

Intuition

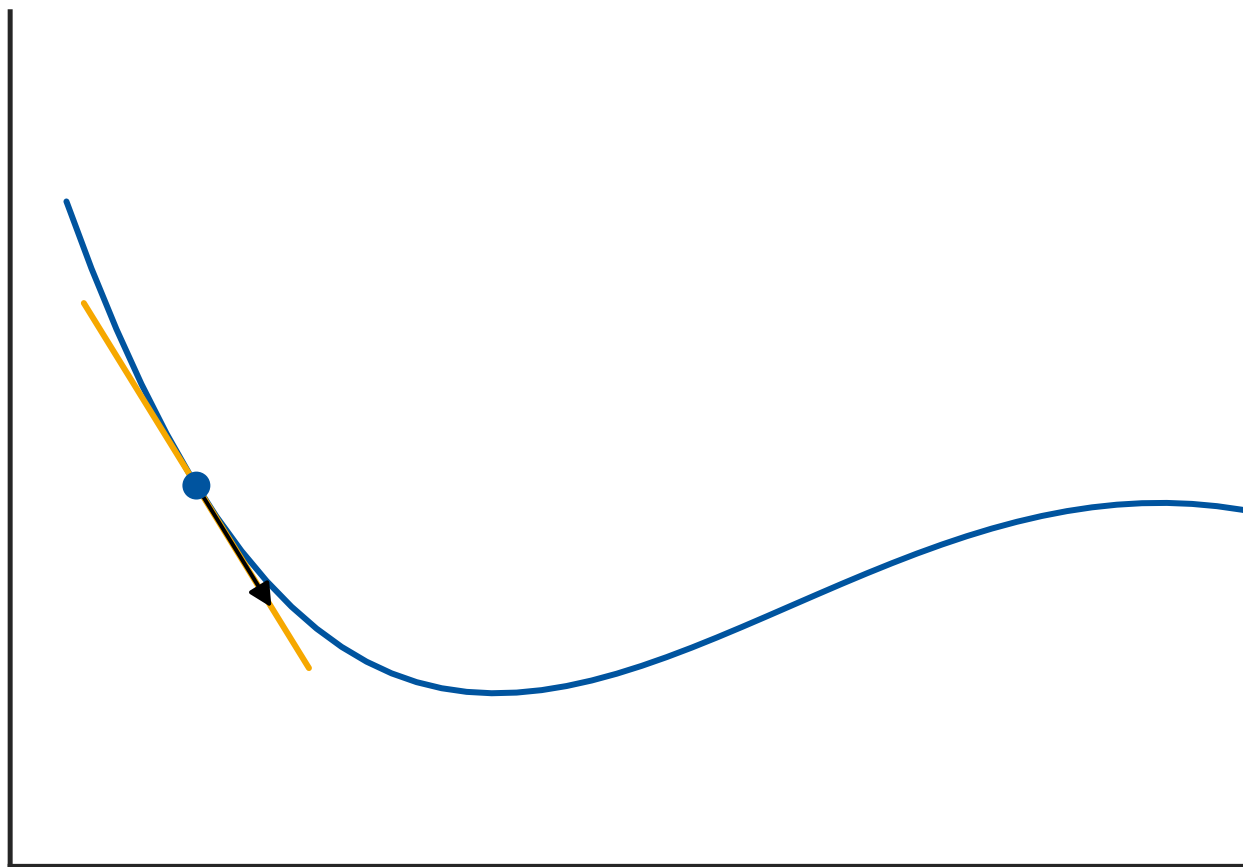


First-Order
 $\eta = 0.005$

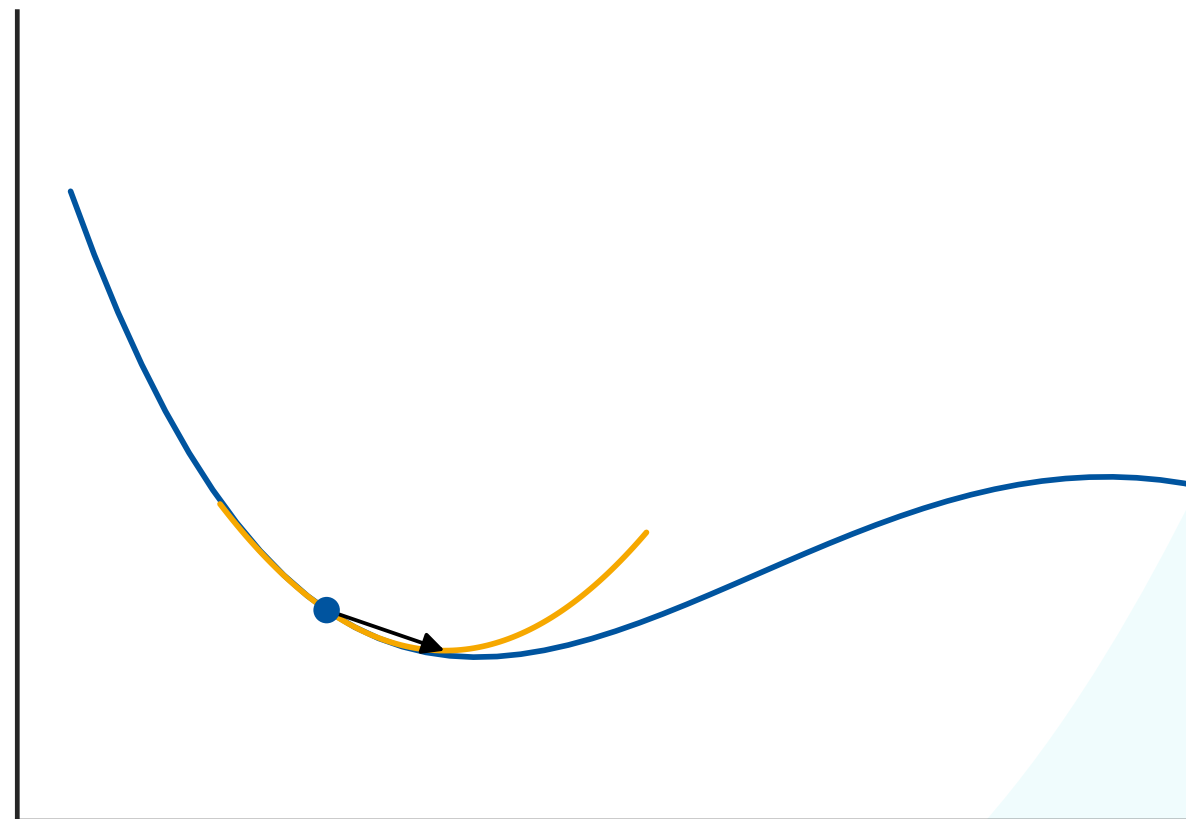


Second-Order

Intuition

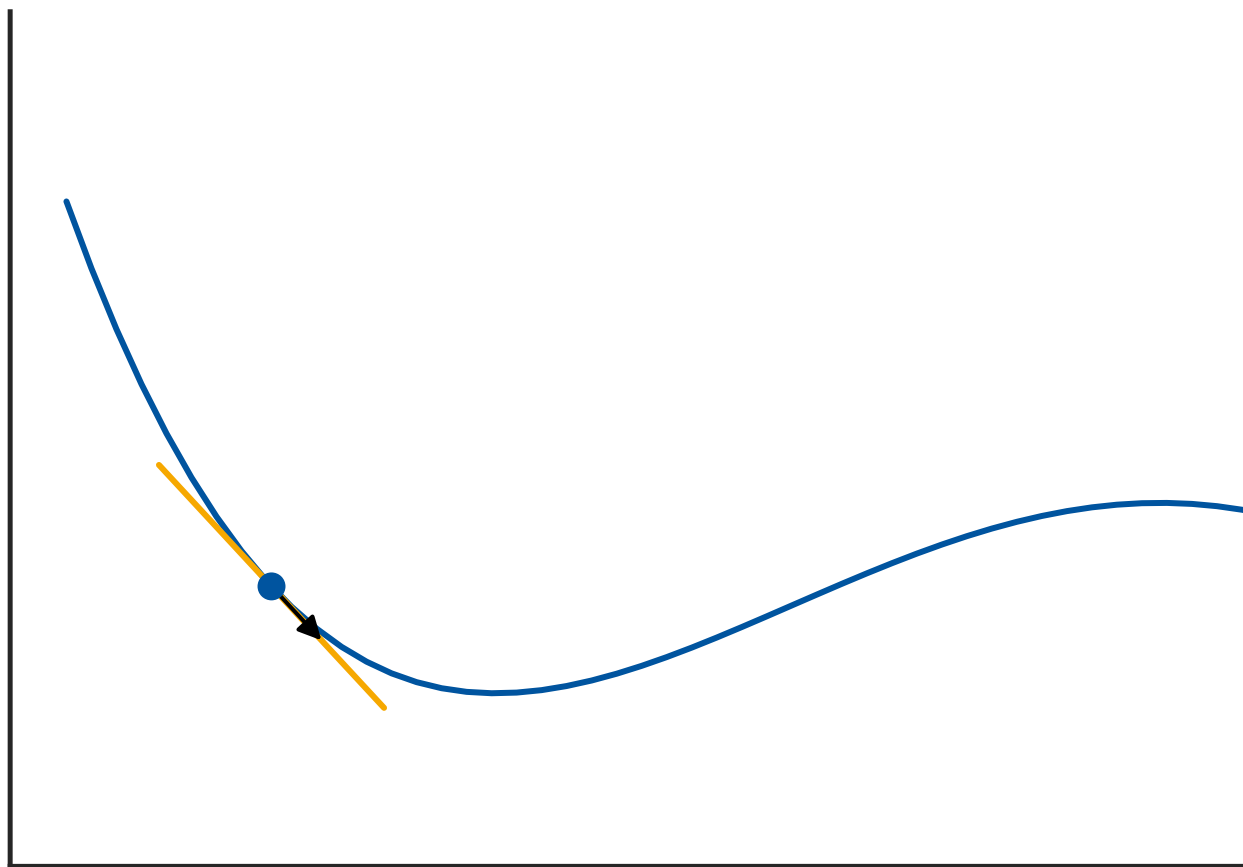


First-Order
 $\eta = 0.005$

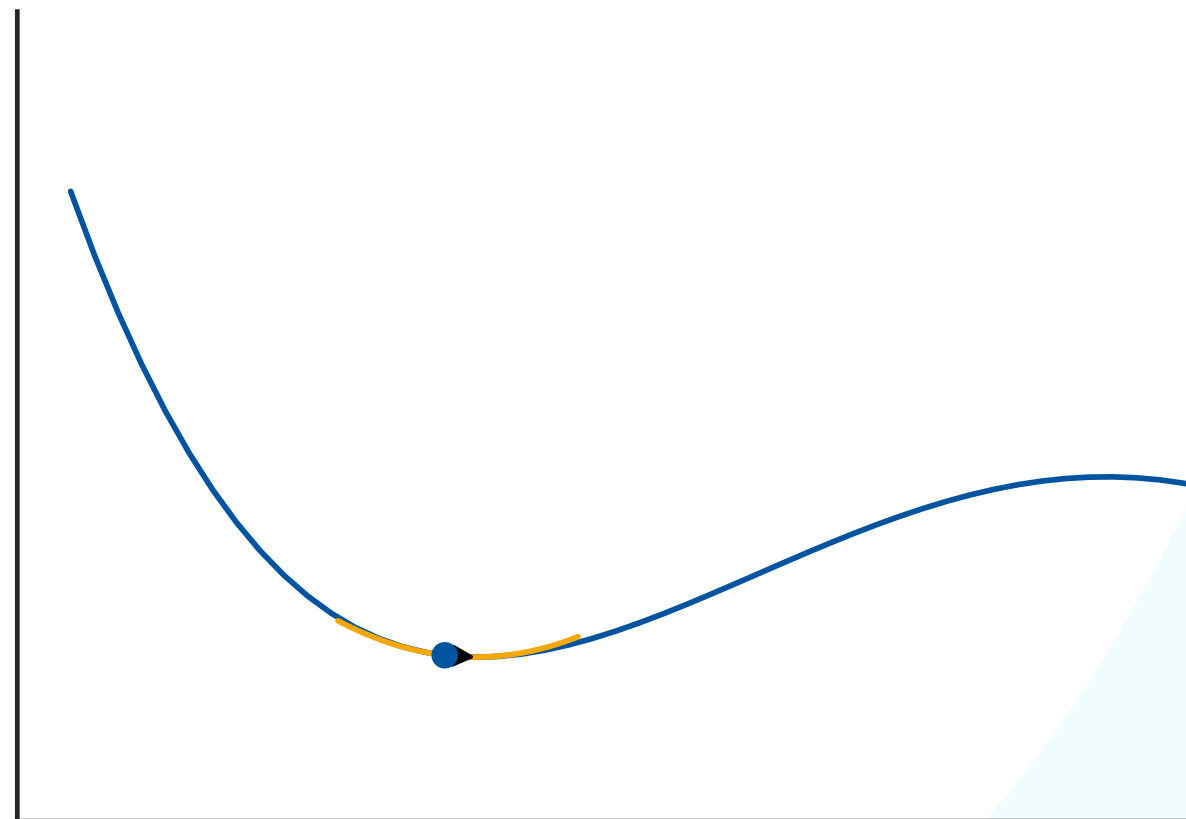


Second-Order

Intuition

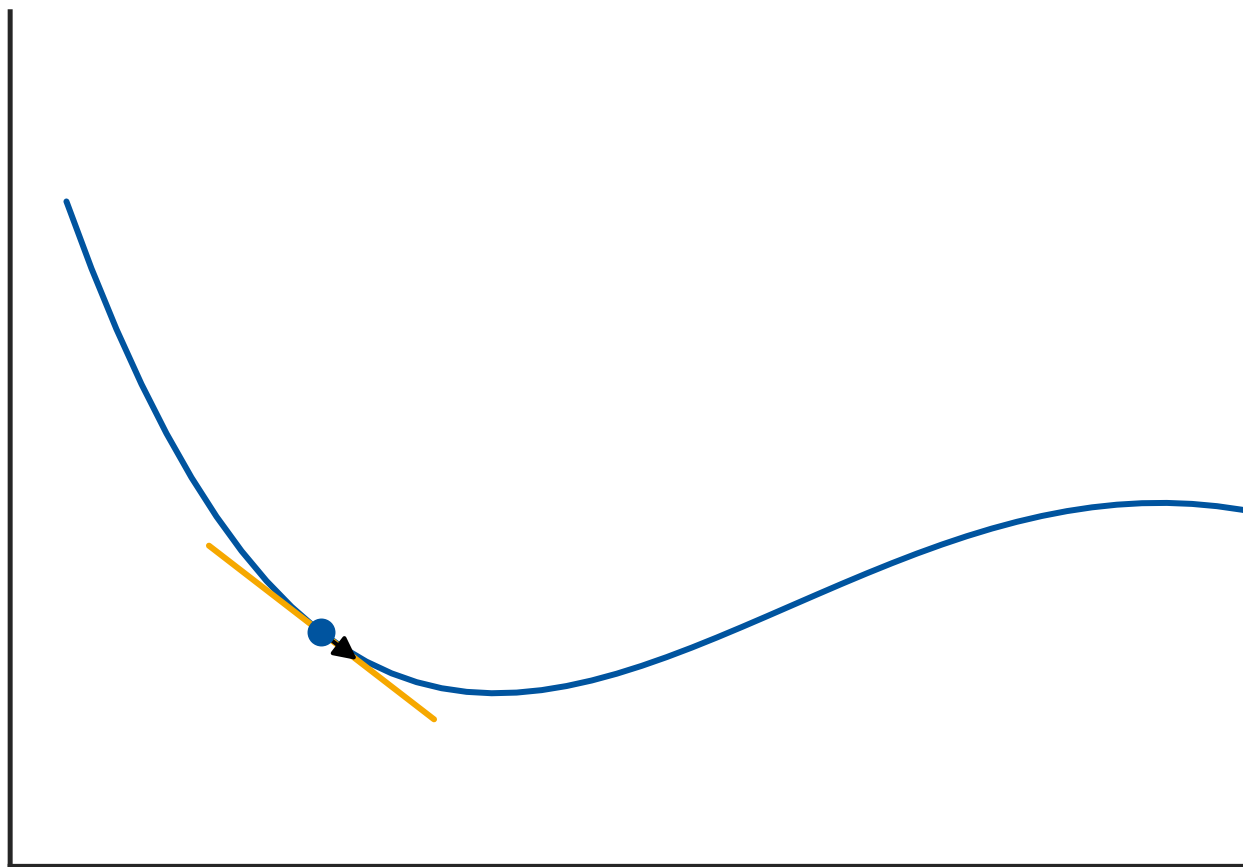


First-Order
 $\eta = 0.005$

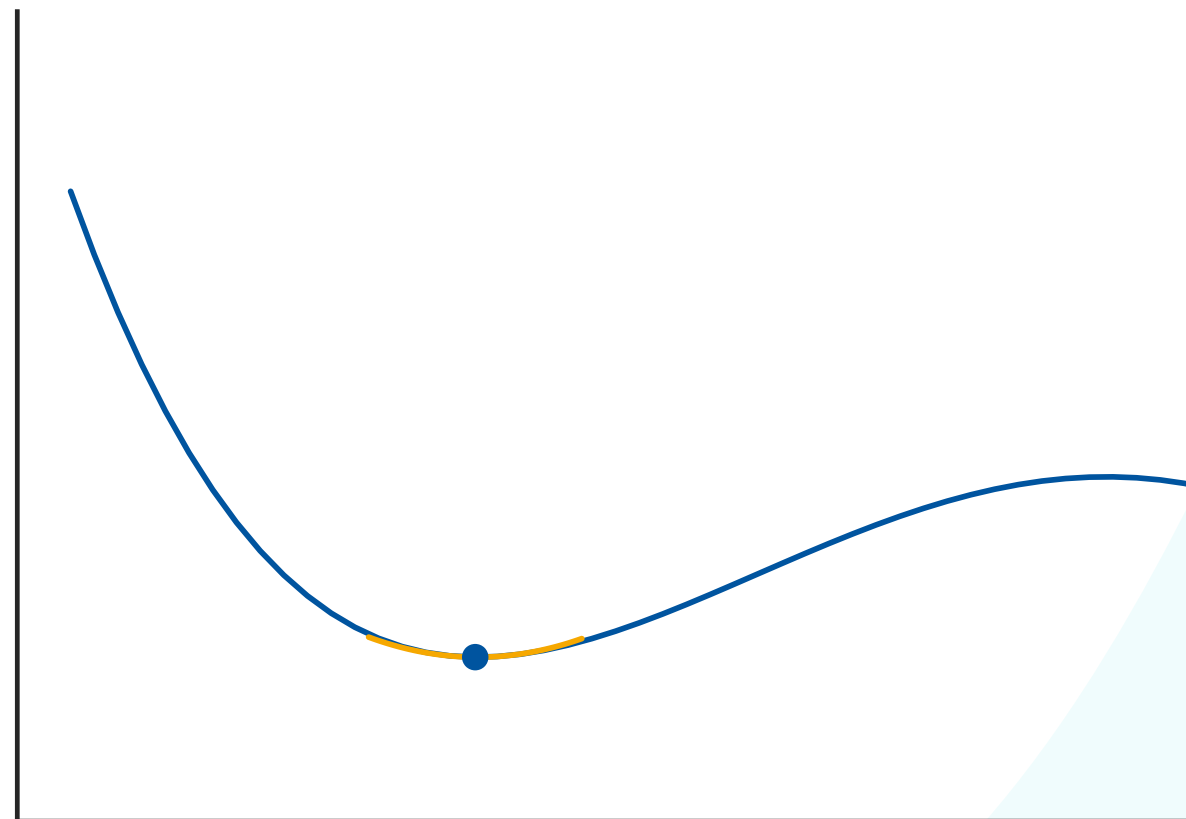


Second-Order

Intuition

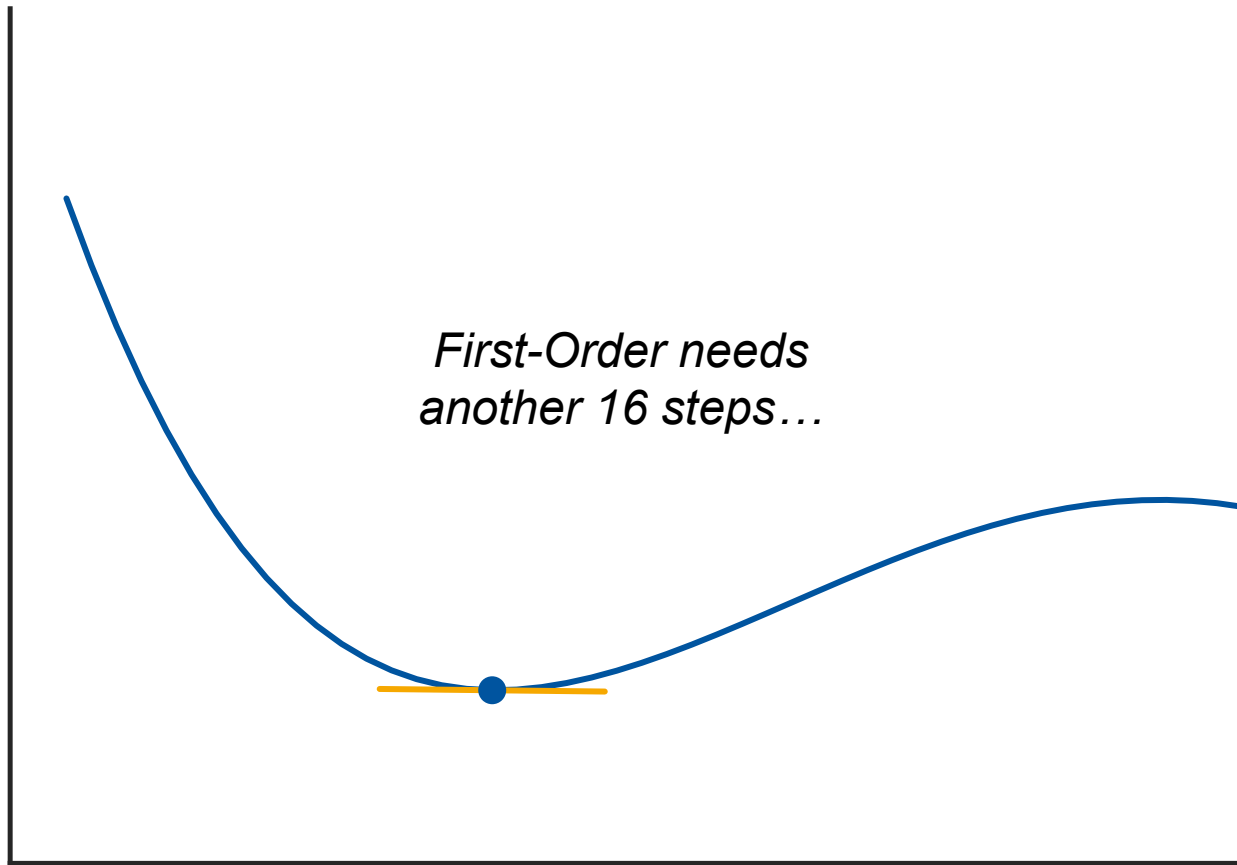


First-Order
 $\eta = 0.005$

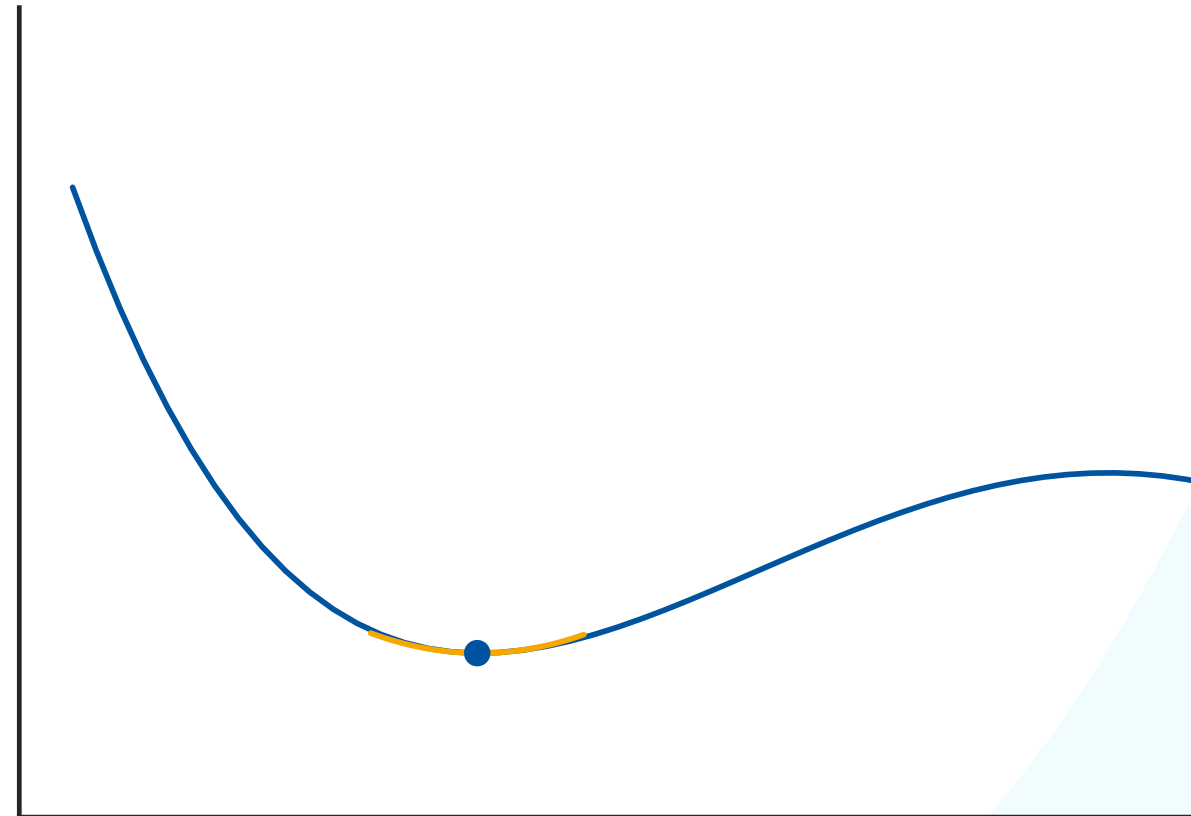


Second-Order

Intuition

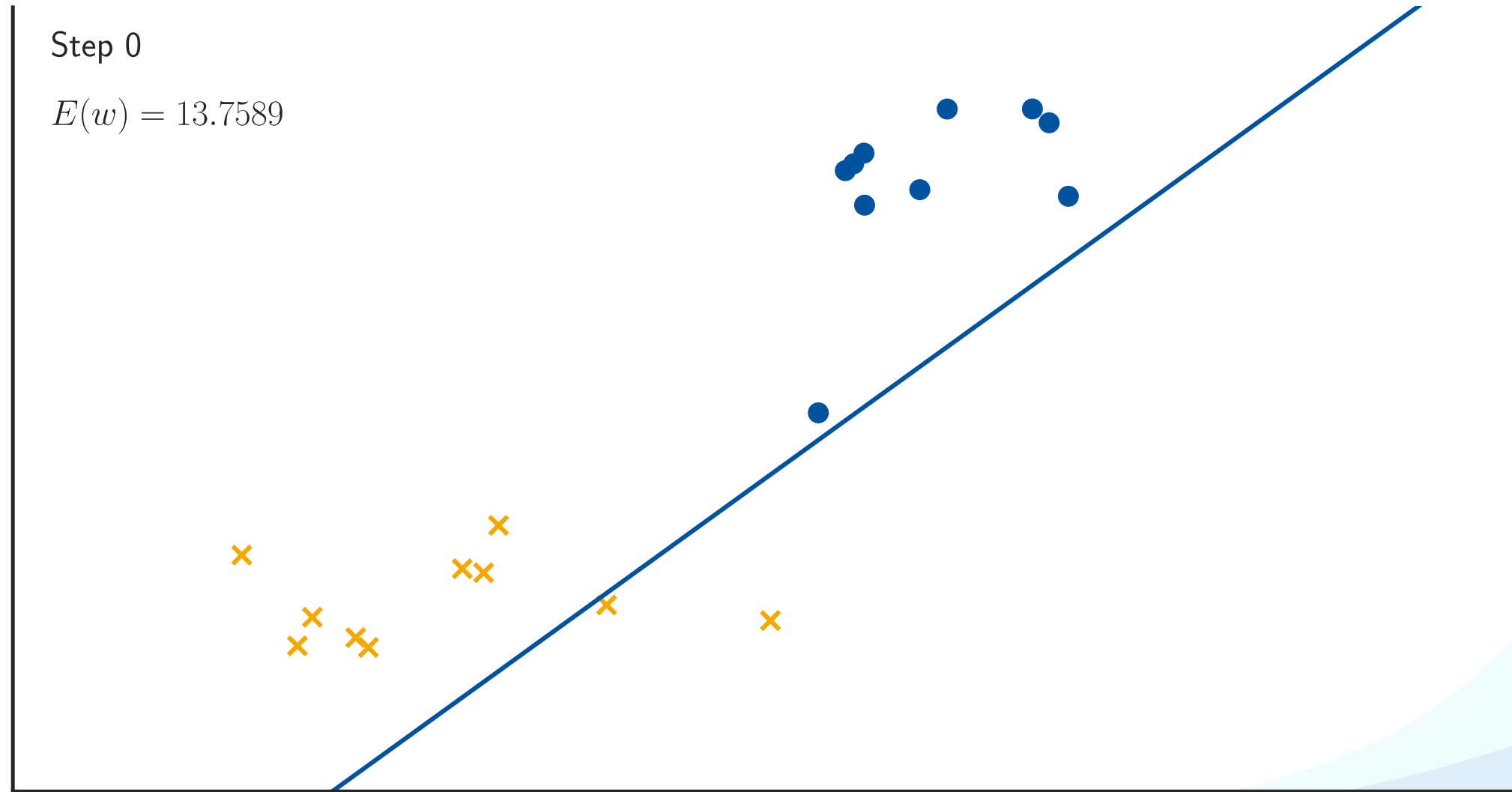


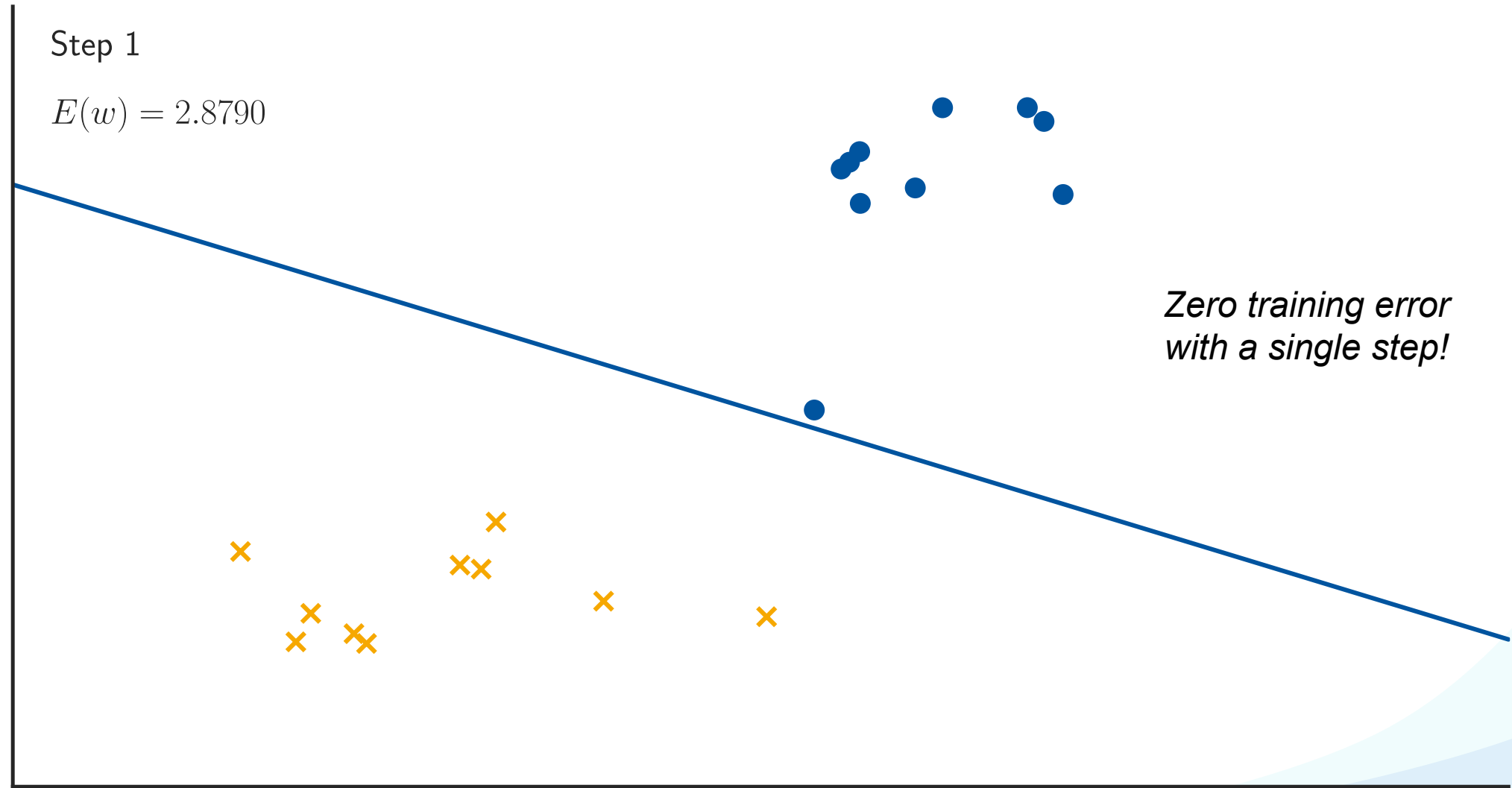
First-Order
 $\eta = 0.005$

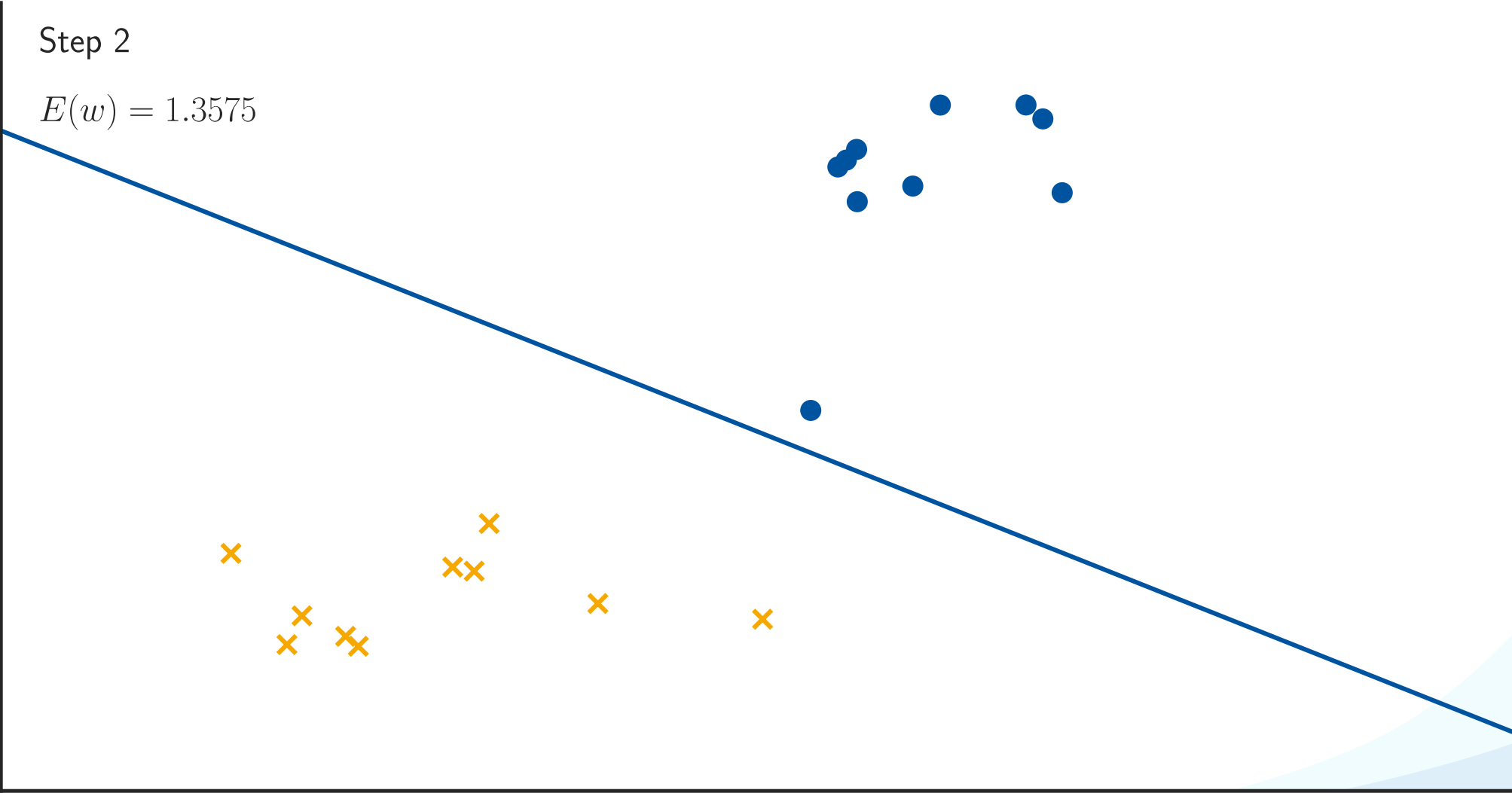


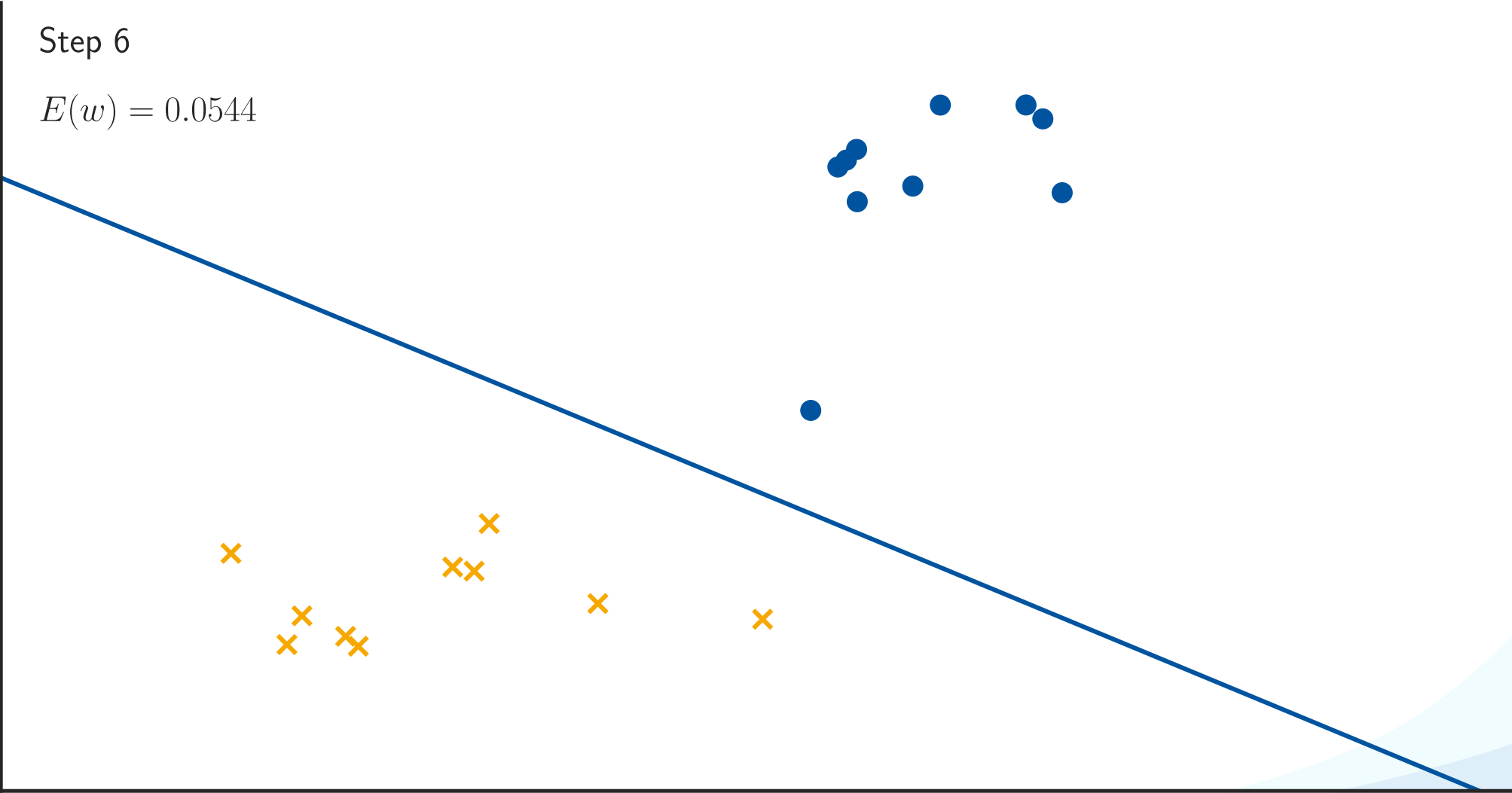
Second-Order

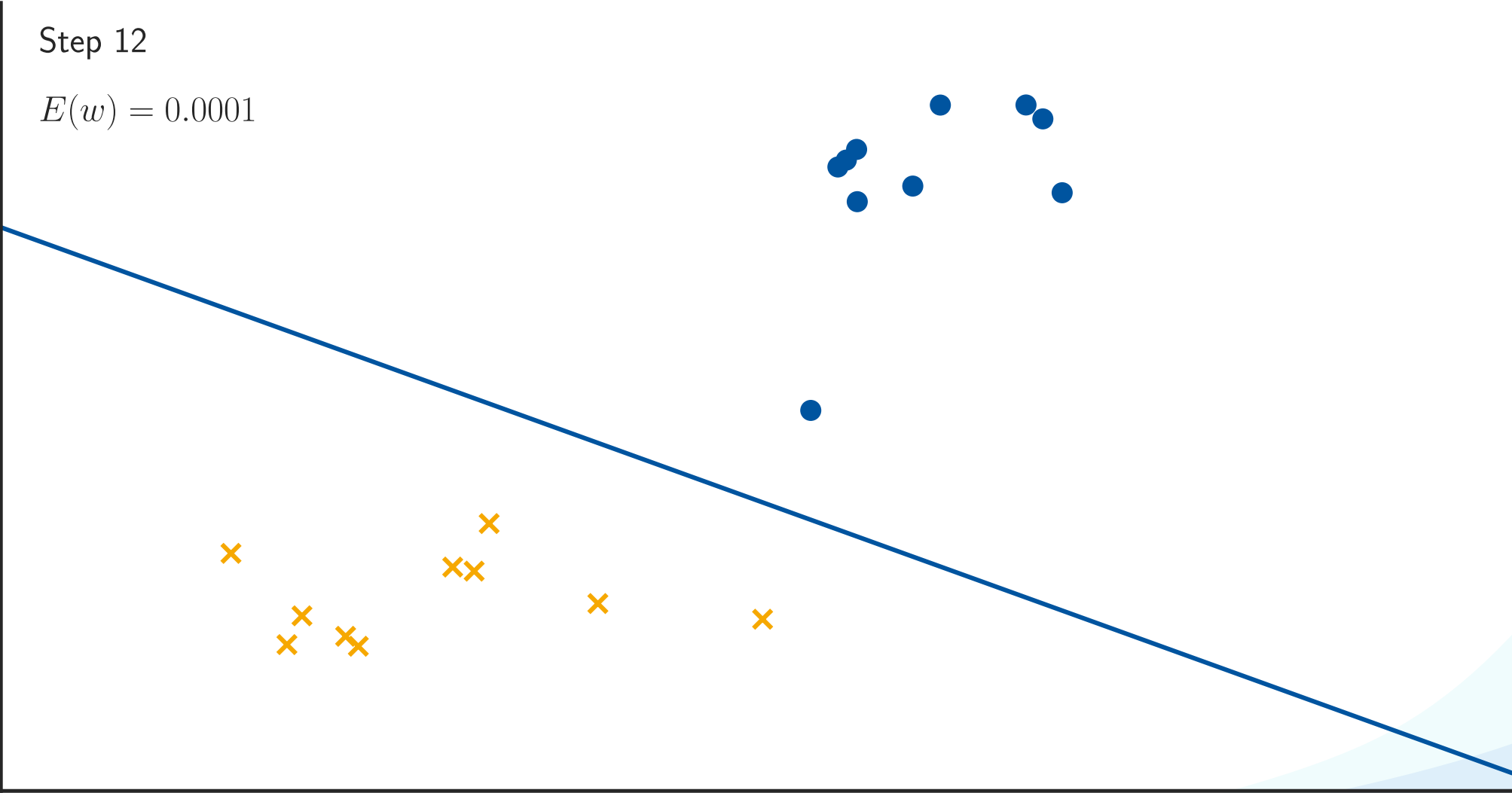
Example: Logistic Regression with 2nd-order Optimization (IRLS)











Discussion: Second-Order Optimization

Advantages

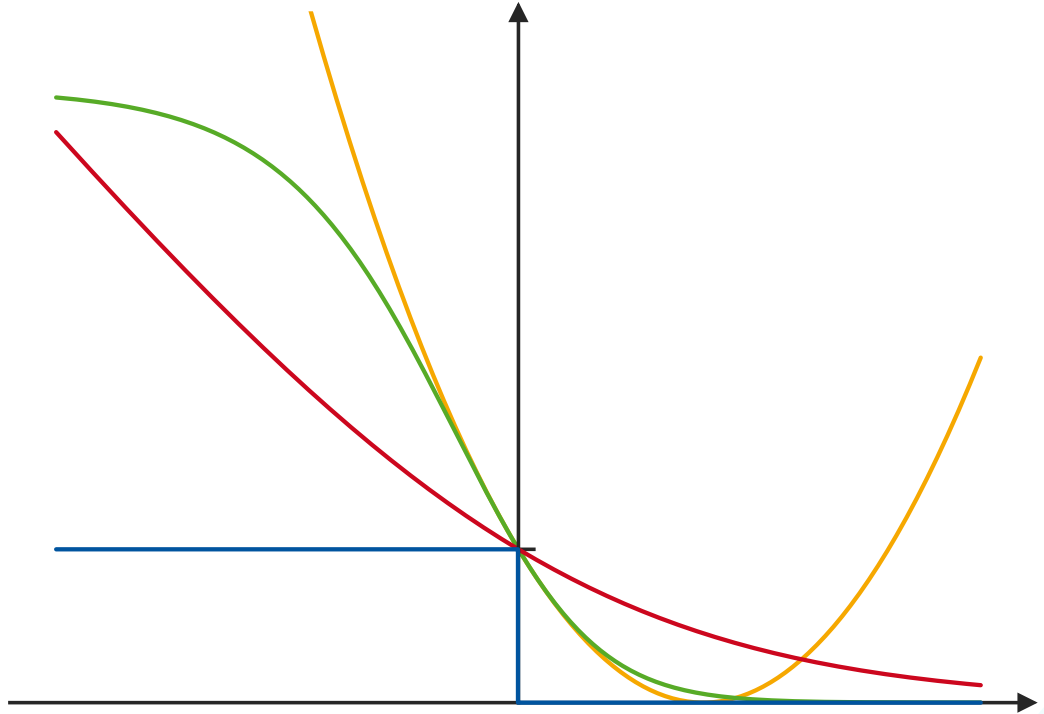
- Faster convergence than first-order methods

Limitations

- Second-order approach, relies on computing second derivatives.
- Computing (and inverting) the Hessian matrix is expensive for problems with many parameters.

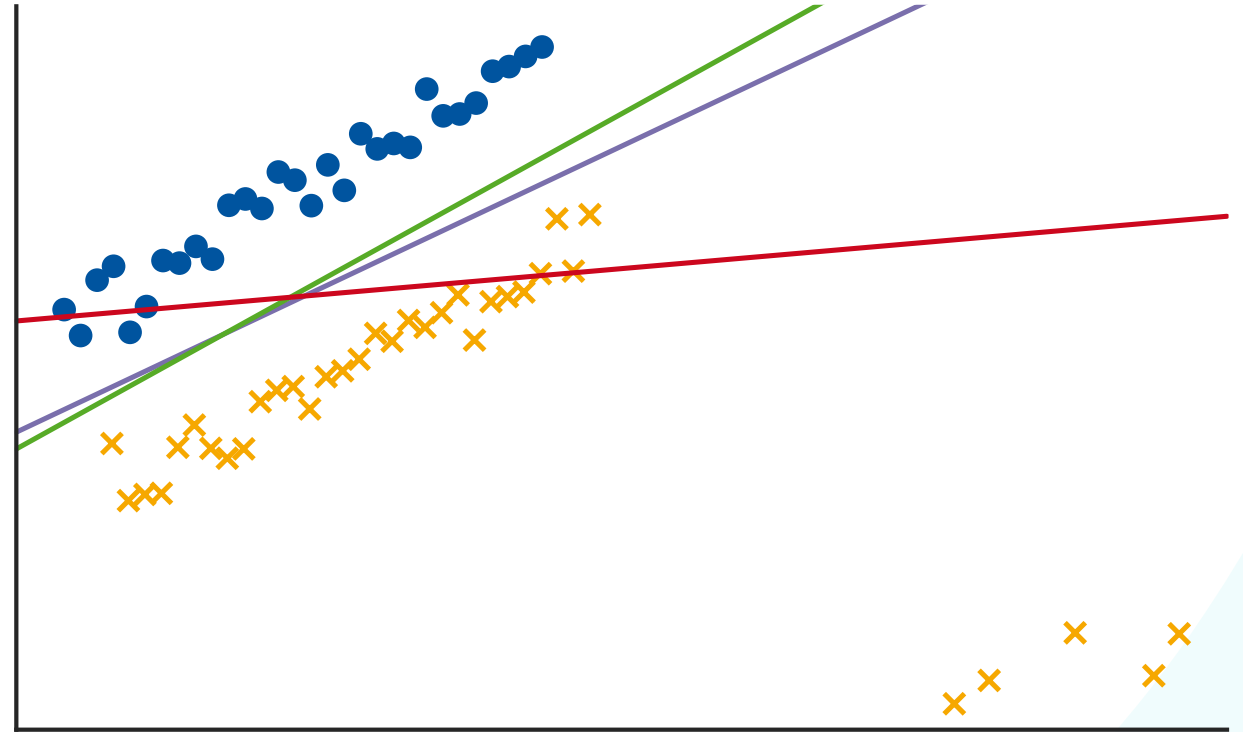
Logistic Regression

1. Logistic Regression Formulation
2. Motivation and Background
3. Iterative Estimation
4. First-Order Gradient Descent
5. Second-Order Gradient Descent
6. **Error Function Analysis**

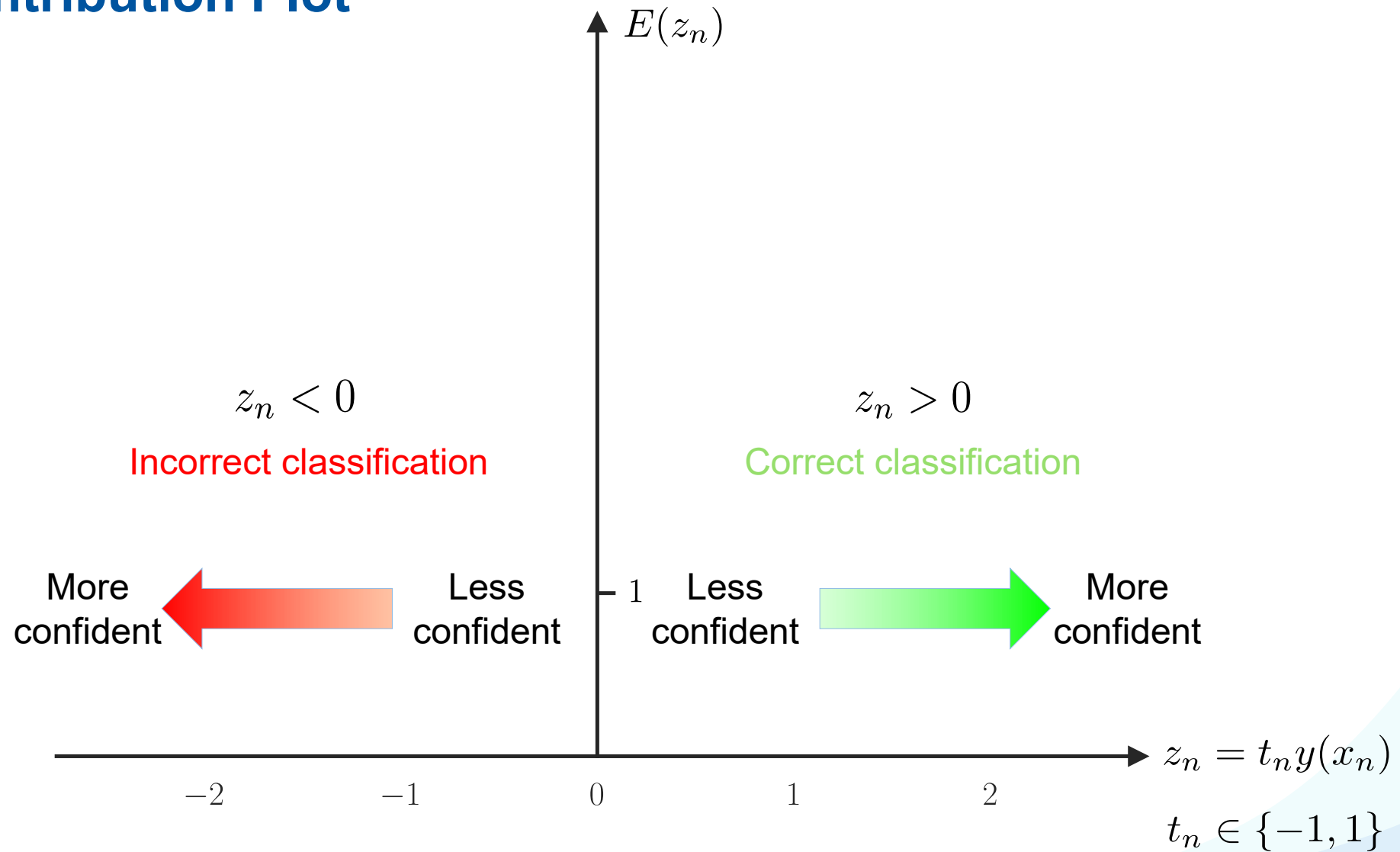


Error Function Analysis

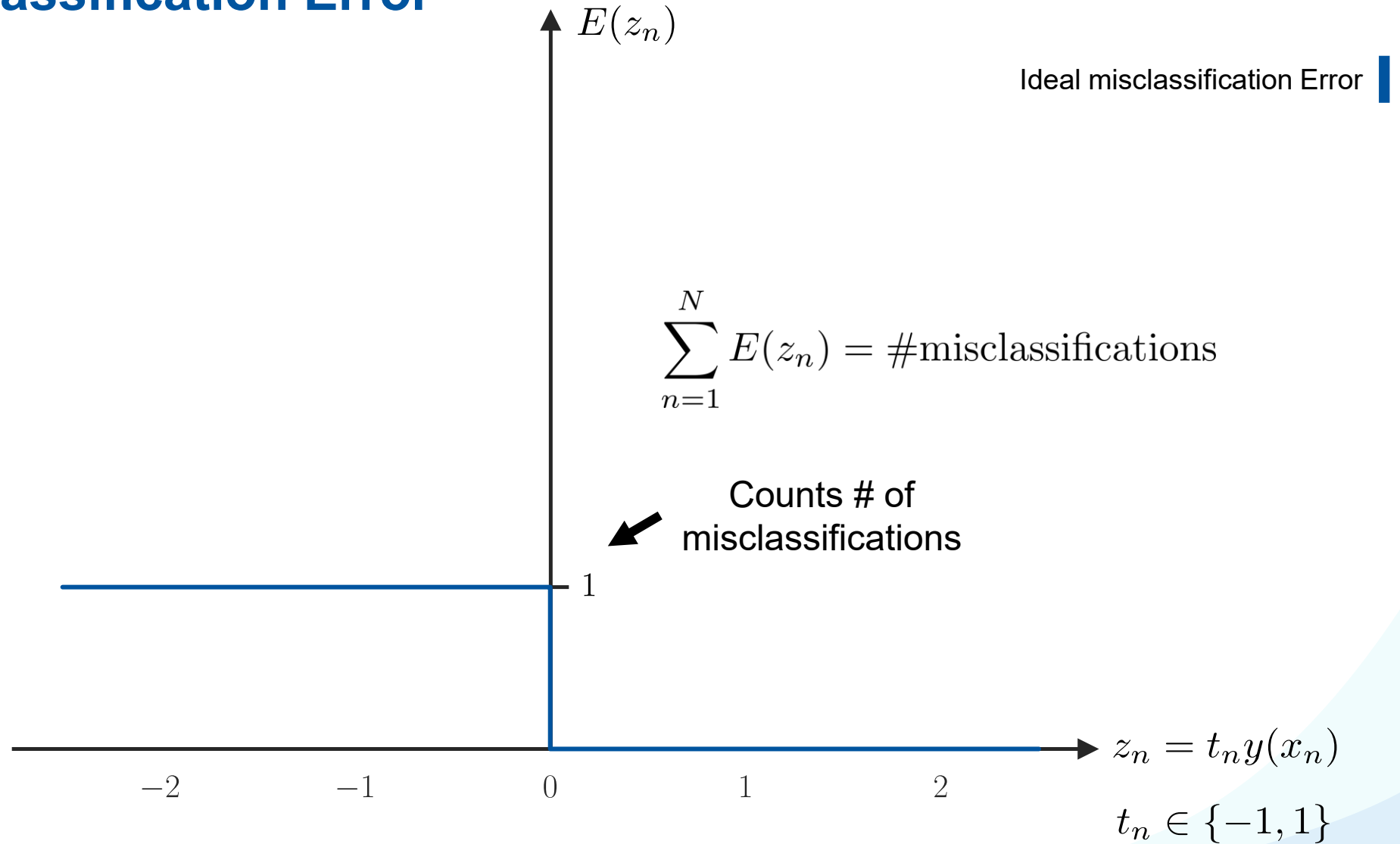
- We have seen how to learn **generalized linear discriminant** models by optimizing an error function.
 - We observed problems with **least-squares classification** based on the squared error function.
 - We have seen that **logistic regression** behaves more robustly.
- *Let's analyze the cross-entropy error in more detail...*



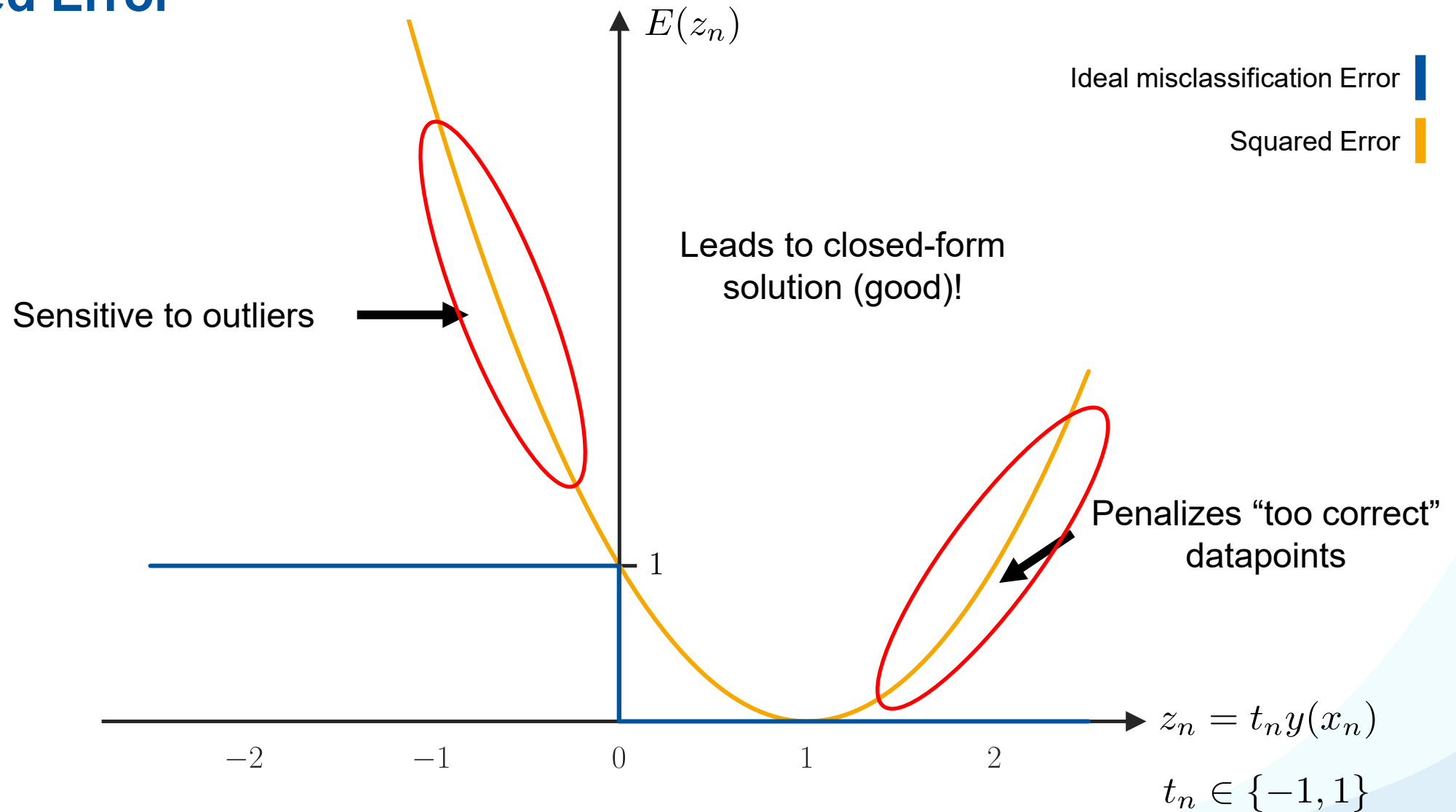
Error Contribution Plot



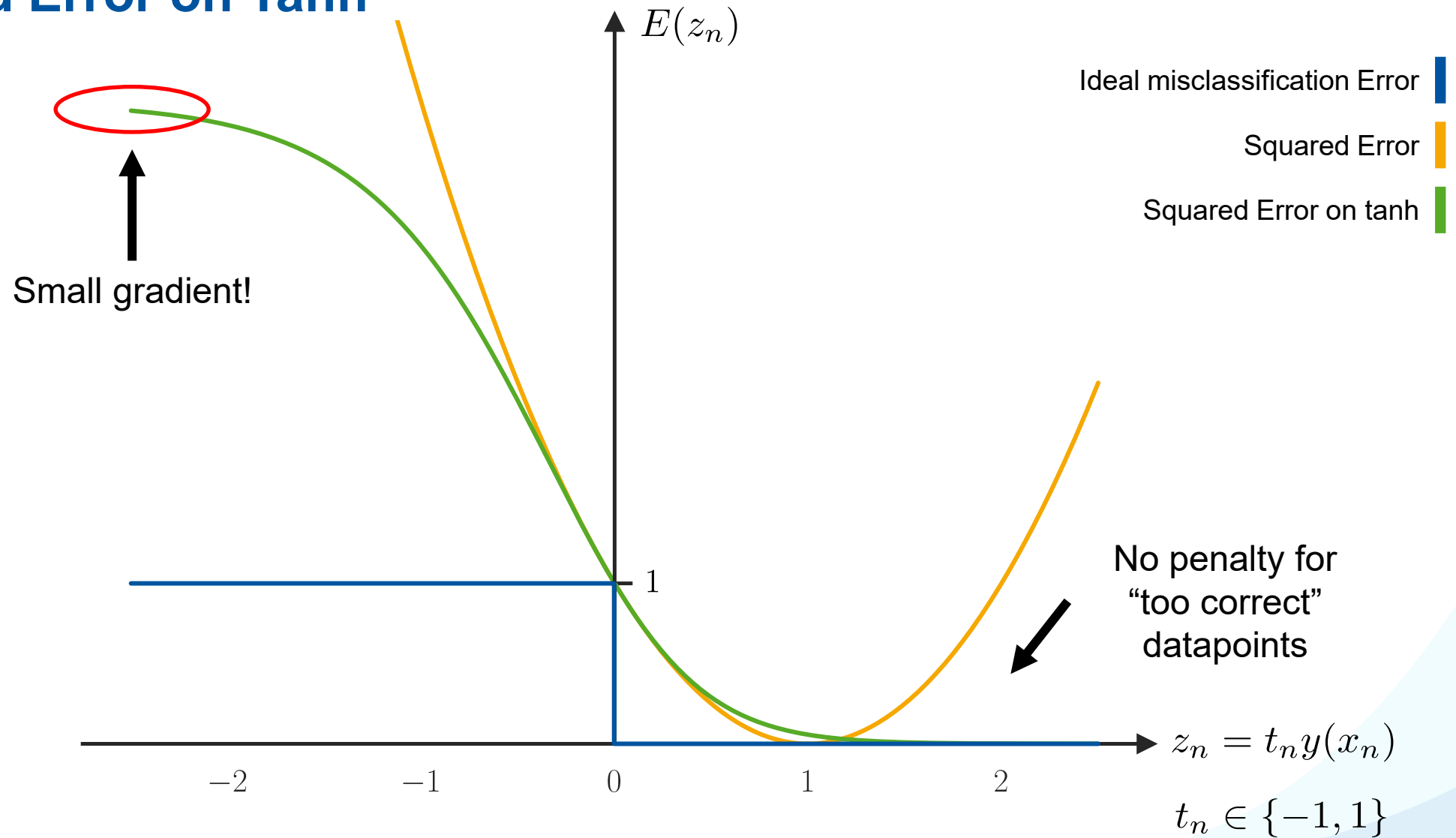
Ideal Misclassification Error



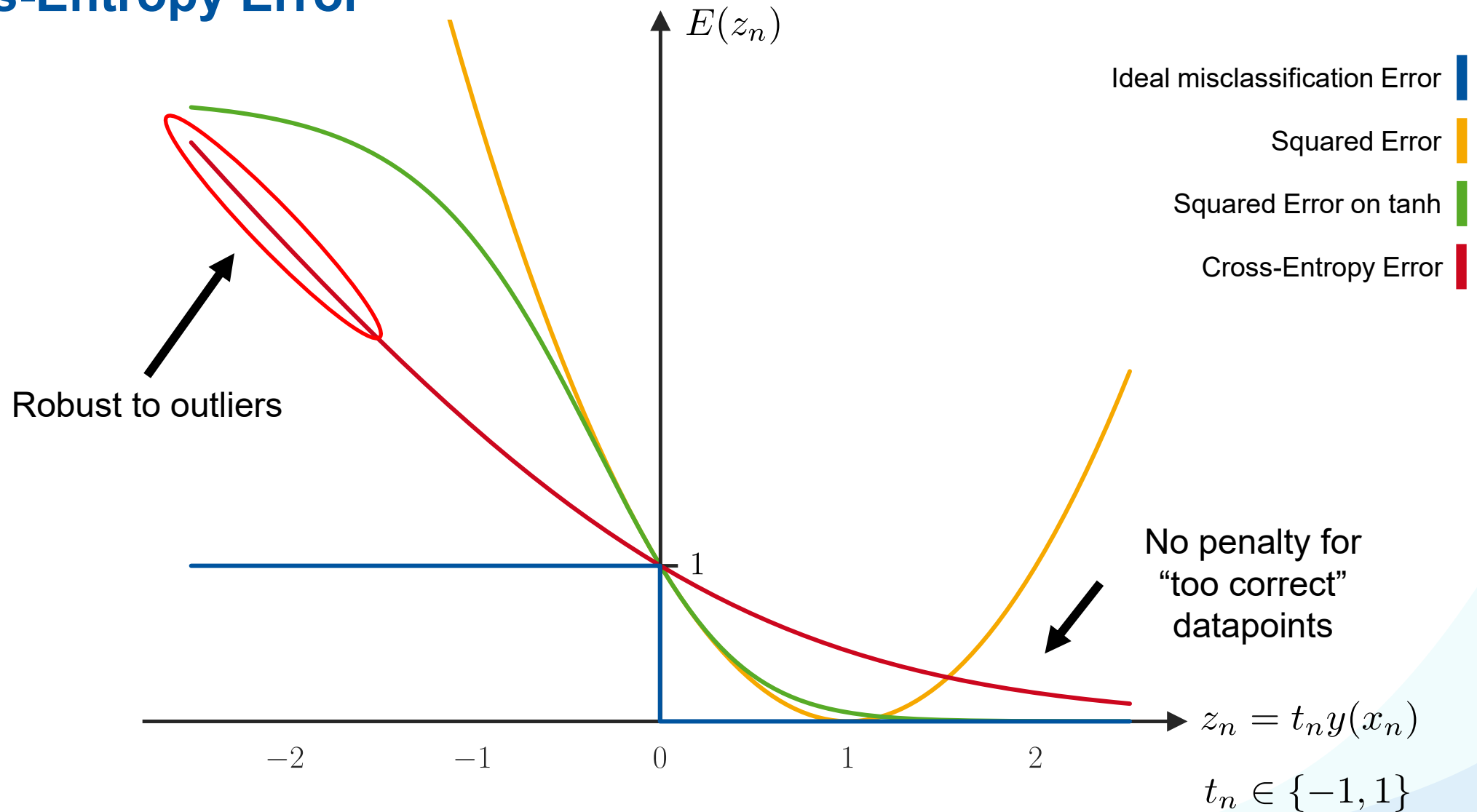
Squared Error



Squared Error on Tanh



Cross-Entropy Error



Discussion: Cross-Entropy Error

Advantages

- Minimizer of this error corresponds to class posteriors
- Convex function, unique minimum exists
- Robust to outliers

Limitations

- No closed-form solution, requires iterative estimation

References and Further Reading

- More information about [Logistic Regression](#) is available in Chapter 4.3 of Bishop's book.

Christopher M. Bishop
Pattern Recognition and Machine Learning
Springer, 2006

