

Logistic Regression

Lecturer:
Jiří Matas

Authors:
Ondřej Drbohlav, Jiří Matas

Centre for Machine Perception
Czech Technical University, Prague
<http://cmp.felk.cvut.cz>

Last update: 24/Oct/2019



Logistic Regression



Outline of the talk

- ◆ Motivation
- ◆ Model, relationship between log odds and posteriors
- ◆ Cross entropy objective function
- ◆ Gradient descent for fitting the model
- ◆ Examples
- ◆ Conclusion

Logistic Regression, Motivation

Consider a classification problem with 0/1 loss matrix. Recall that given an observation \mathbf{x} , the optimal Bayesian strategy $q(\mathbf{x})$ decides for a class k which maximizes the posterior:

$$q(\mathbf{x}) = \underset{k}{\operatorname{argmax}} p(k|\mathbf{x}) = \underset{k}{\operatorname{argmax}} p(\mathbf{x}, k) = \underset{k}{\operatorname{argmax}} p(\mathbf{x}|k)p(k). \quad (1)$$

For a binary (2-class) classification,

$$q(\mathbf{x}) = 1 \quad \text{if} \quad p(1|\mathbf{x}) > p(2|\mathbf{x}) \quad \Leftrightarrow \quad \frac{p(1|\mathbf{x})}{p(2|\mathbf{x})} > 1 \quad \Leftrightarrow \quad \ln \frac{p(1|\mathbf{x})}{p(2|\mathbf{x})} > 0, \quad (2)$$

$$q(\mathbf{x}) = 2 \quad \text{if} \quad p(1|\mathbf{x}) < p(2|\mathbf{x}) \quad \Leftrightarrow \quad \frac{p(1|\mathbf{x})}{p(2|\mathbf{x})} < 1 \quad \Leftrightarrow \quad \ln \frac{p(1|\mathbf{x})}{p(2|\mathbf{x})} < 0. \quad (3)$$

The ratio of posteriors $\frac{p(1|\mathbf{x})}{p(2|\mathbf{x})}$ is called **odds ratio**, the logarithm of this ratio, $\ln \frac{p(1|\mathbf{x})}{p(2|\mathbf{x})}$, is called **log odds**.

Logistic Regression, Motivation

Why are we interested in log odds? Because for the following problems, the log odds are linear (therefore simple) function of the observation variable \mathbf{x} :

- ◆ Normal distributions with equal variances
- ◆ Independent features with binary outcomes
- ◆ Multinomial naive Bayes

Normal distributions with equal covariance matrices

$$p(\mathbf{x}|1) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{D}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_1)^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu}_1)}, \quad (4)$$

$$p(\mathbf{x}|2) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_2, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{D}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_2)^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu}_2)}, \quad (5)$$

$$(\boldsymbol{\mu}_1, \boldsymbol{\mu}_2 \in \mathbb{R}^D, \boldsymbol{\Sigma} \in \mathbb{R}^{D \times D}, \boldsymbol{\Sigma} \succ 0, \boldsymbol{\Sigma} = \boldsymbol{\Sigma}^\top) \quad (6)$$

The log odds:

$$\ln \frac{p(1|\mathbf{x})}{p(2|\mathbf{x})} = \ln \frac{p(\mathbf{x}|1)p(1)}{p(\mathbf{x}|2)p(2)} = \ln \frac{p(\mathbf{x}|1)}{p(\mathbf{x}|2)} + \ln \frac{p(1)}{p(2)} = \quad (7)$$

$$= \ln \frac{\cancel{(2\pi)^{-\frac{D}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}}}}{\cancel{(2\pi)^{-\frac{D}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}}}} e^{-\frac{1}{2}\{(\mathbf{x}-\boldsymbol{\mu}_1)^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu}_1) - (\mathbf{x}-\boldsymbol{\mu}_2)^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu}_2)\}} + \ln \frac{p(1)}{p(2)} \quad (8)$$

Logistic Regression, Motivation

$$\ln \frac{p(1|\mathbf{x})}{p(2|\mathbf{x})} = \ln \frac{(2\pi)^{-\frac{D}{2}} |\Sigma|^{-\frac{1}{2}} e^{-\frac{1}{2}\{(\mathbf{x}-\mu_1)^\top \Sigma^{-1}(\mathbf{x}-\mu_1) - (\mathbf{x}-\mu_2)^\top \Sigma^{-1}(\mathbf{x}-\mu_2)\}}}{(2\pi)^{-\frac{D}{2}} |\Sigma|^{-\frac{1}{2}}} + \ln \frac{p(1)}{p(2)} \quad (9)$$

$$= -\frac{1}{2}(\mathbf{x}^\top \Sigma^{-1} \mathbf{x} - \mathbf{x}^\top \Sigma^{-1} \mu_1 - \mu_1^\top \Sigma^{-1} \mathbf{x} + \mu_1^\top \Sigma^{-1} \mu_1) \quad (10)$$

$$+ \frac{1}{2}(\mathbf{x}^\top \Sigma^{-1} \mathbf{x} - \mathbf{x}^\top \Sigma^{-1} \mu_2 - \mu_2^\top \Sigma^{-1} \mathbf{x} + \mu_2^\top \Sigma^{-1} \mu_2) + \ln \frac{p(1)}{p(2)} \quad (11)$$

(Note: $\mathbf{x}^\top \Sigma^{-1} \mathbf{y} = \mathbf{y}^\top \Sigma^{-1} \mathbf{x}$ because $\Sigma = \Sigma^\top$)

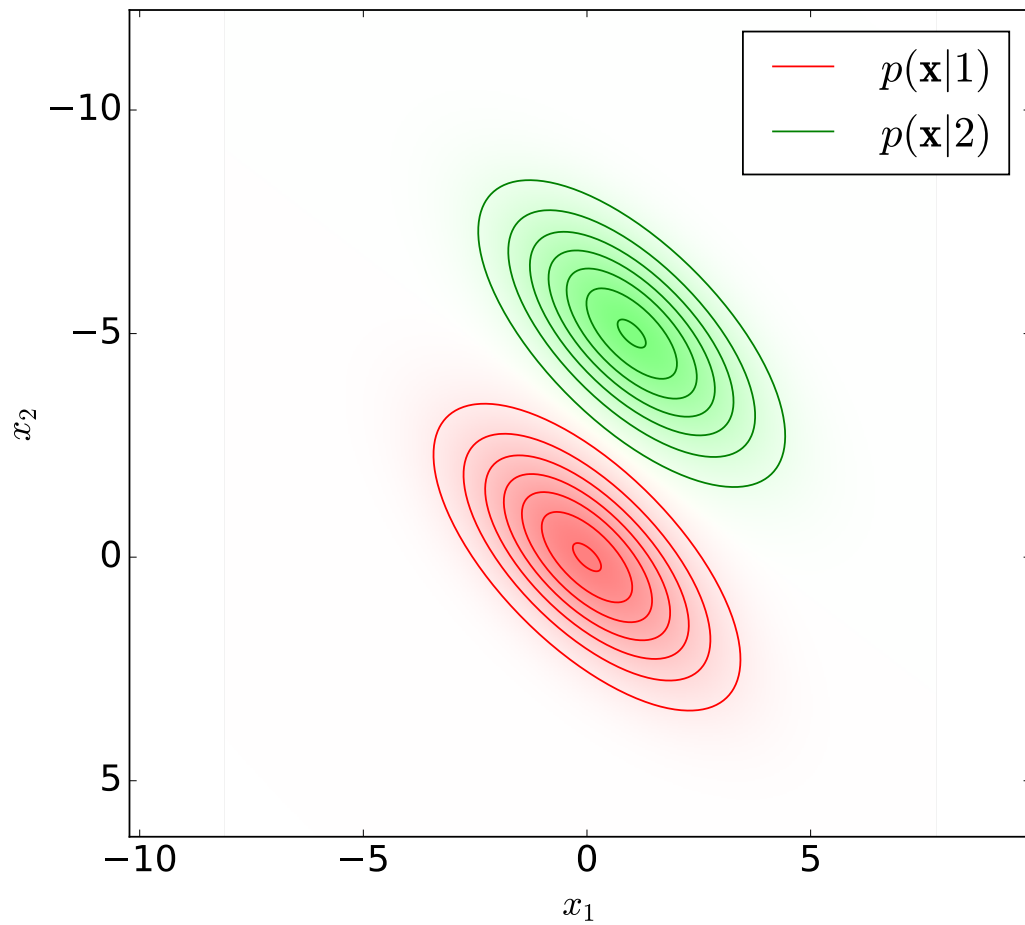
$$= -\frac{1}{2}(-2\mu_1^\top \Sigma^{-1} \mathbf{x} + \mu_1^\top \Sigma^{-1} \mu_1 - (-2\mu_2^\top \Sigma^{-1} \mathbf{x} + \mu_2^\top \Sigma^{-1} \mu_2)) + \ln \frac{p(1)}{p(2)} \quad (12)$$

$$= \underbrace{[(\mu_1 - \mu_2)^\top \Sigma^{-1}] \mathbf{x}}_{\mathbf{w}} + \underbrace{\frac{1}{2}(\mu_1^\top \Sigma^{-1} \mu_1 - \mu_2^\top \Sigma^{-1} \mu_2)}_{w_0} + \ln \frac{p(1)}{p(2)} \quad (13)$$

$$= \mathbf{w} \cdot \mathbf{x} + w_0, \quad (\mathbf{w} \in \mathbb{R}^D, w_0 \in R) \quad (14)$$

Conclusion: When two classes are normally distributed with equal covariance matrices, the log odds are a **linear** function of observation vector \mathbf{x} .

Example



Multinomial normal distribution,
two classes

Centers:

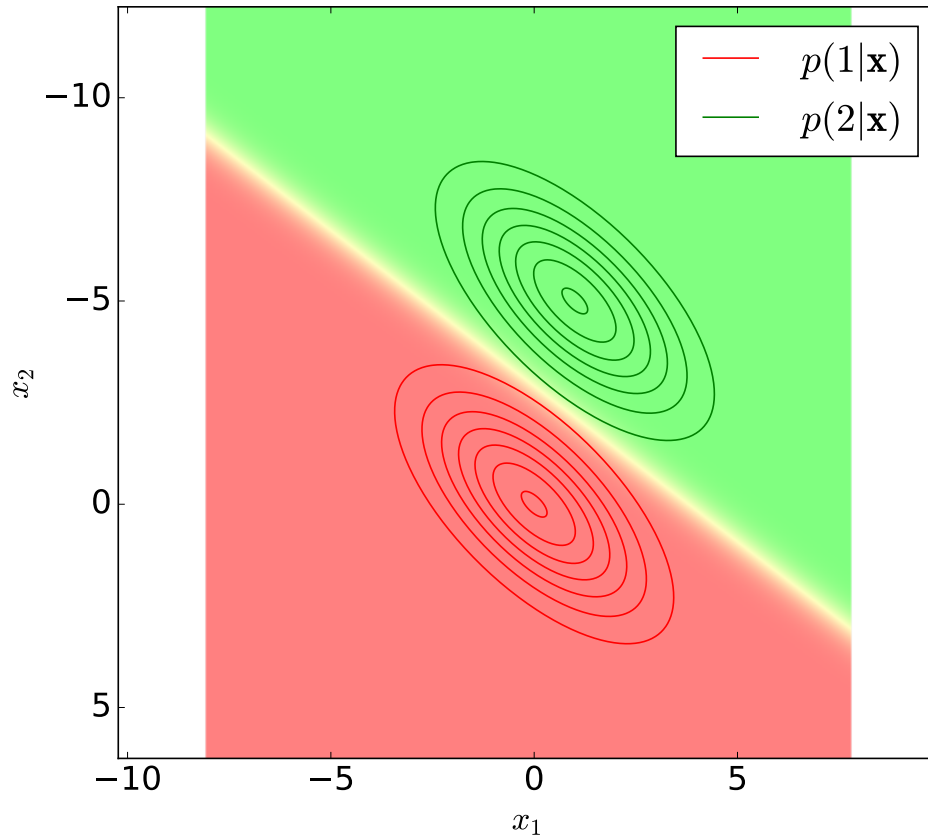
$(0, 0)$, $(1, -5)$.

Cov. matrices (are equal):

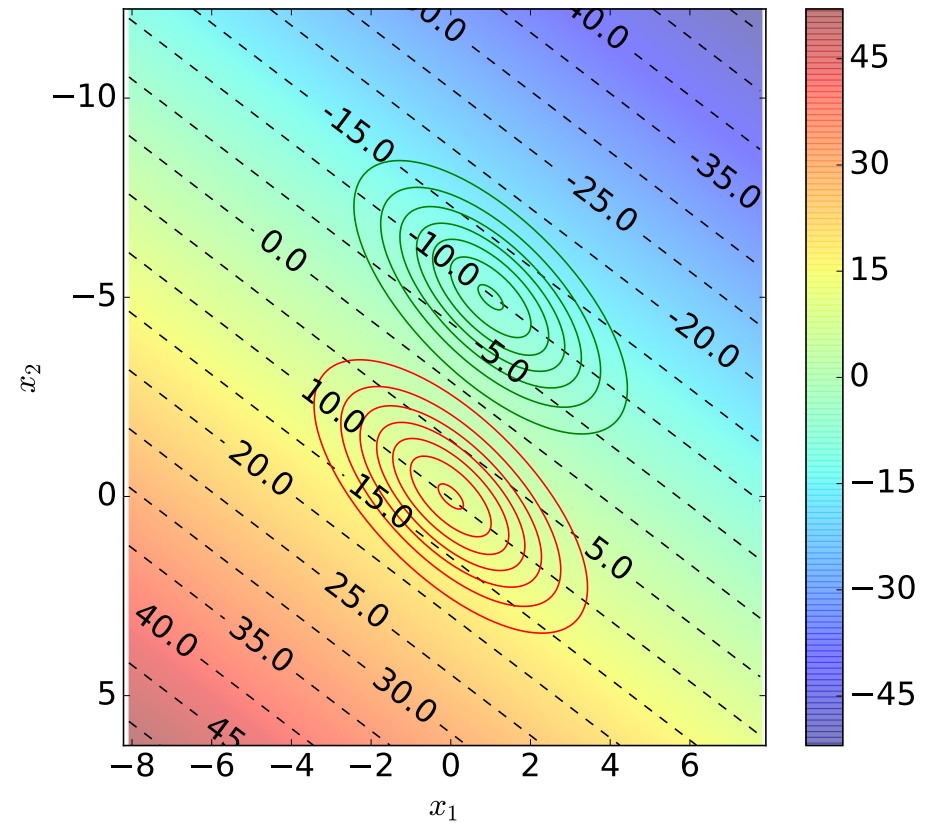
$$\Sigma = \begin{bmatrix} 3 & 2 \\ 2 & 3 \end{bmatrix}$$

Example (contd.)

Priors: $p(\text{red}) = 0.5$, $p(\text{green}) = 0.5$



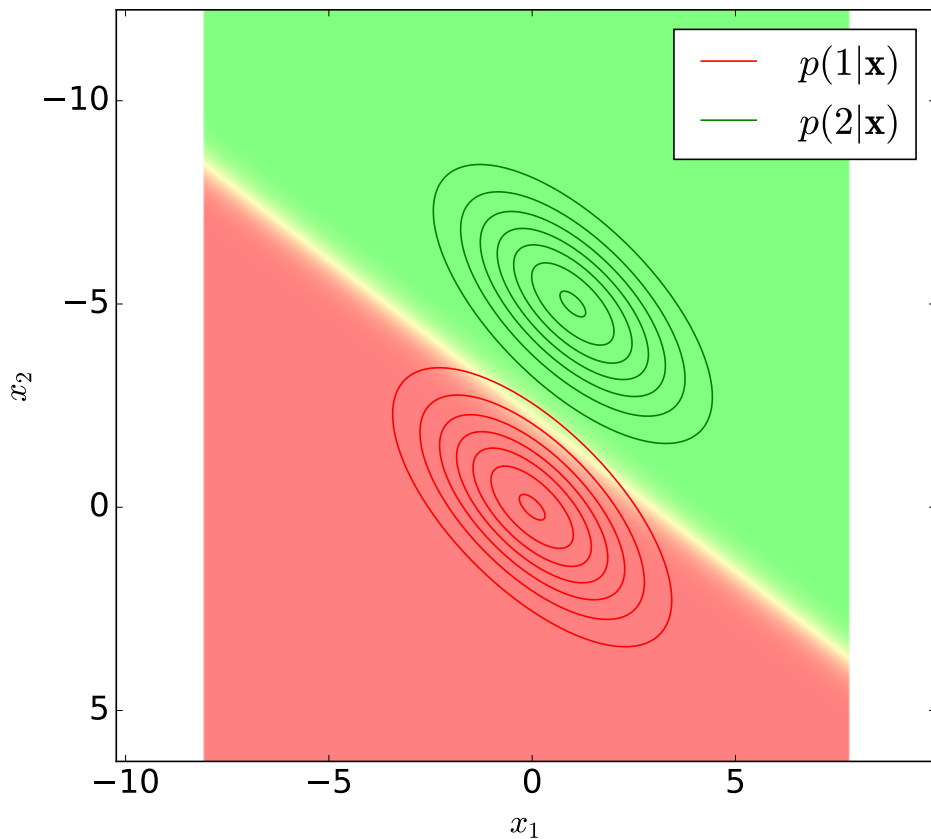
Probabilities $p(1|\mathbf{x})$ and $p(2|\mathbf{x})$.



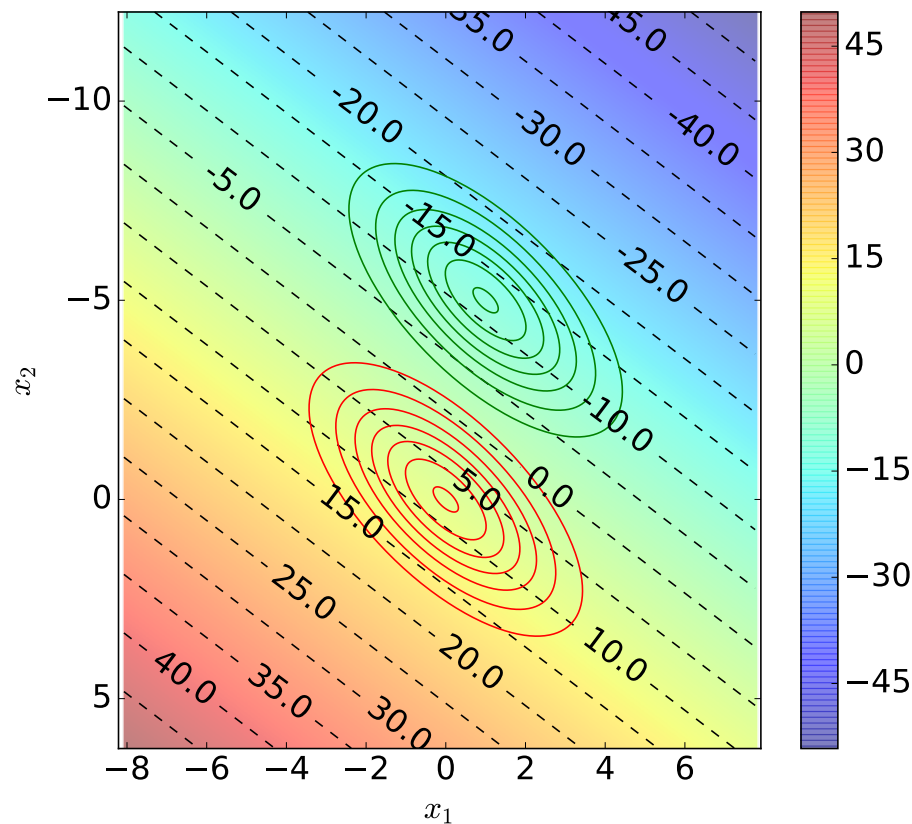
Log odds, shown as a map with contour plot. Note that the log odds are a linear function.

Example (contd.)

Priors: $p(\text{red}) = 0.1$, $p(\text{green}) = 0.9$



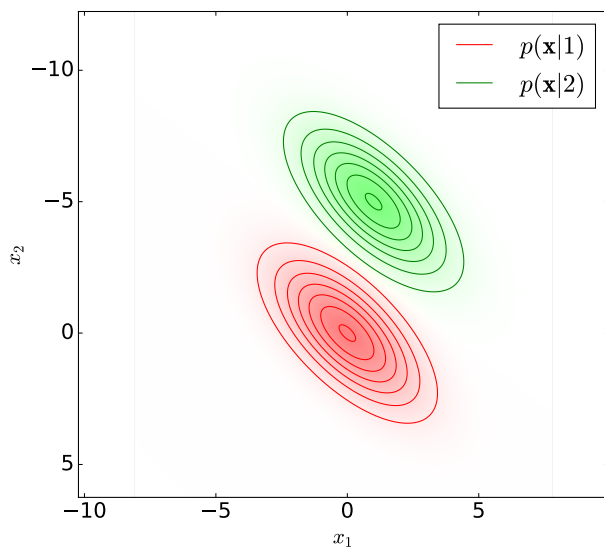
Probabilities $p(1|\mathbf{x})$ and $p(2|\mathbf{x})$.



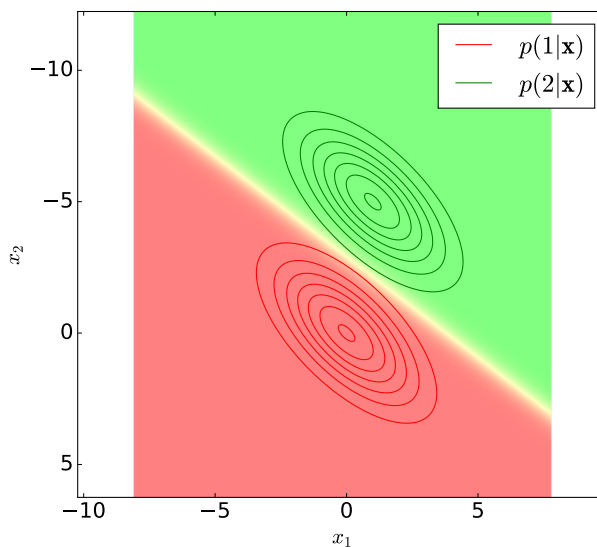
Log odds, shown as a map with contour plot. Note that the log odds are a linear function. Also note the shift of the zero level w.r.t. previous slide.

Example (contd.)

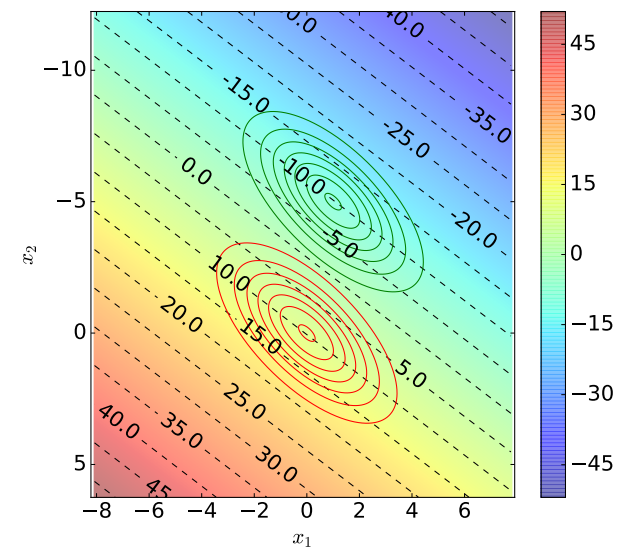
conditionals



posteriors



log odds



It will soon be shown for linear log odds, $a(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + w_0$, the posteriors can be written as

$$p(1|\mathbf{x}) = \frac{1}{1 + e^{-a(\mathbf{x})}} \tag{15}$$

$$p(2|\mathbf{x}) = \frac{1}{1 + e^{a(\mathbf{x})}} \tag{16}$$

The idea of logistic regression is that posteriors are modeled this way directly, without first estimating the priors and conditionals.

Other models with linear log odds (1)

Independent features with binary outcomes.

Feature vector $\mathbf{x} = (x_1, x_2, \dots, x_D) \in \mathbb{R}^D$ $x_i \in \{0, 1\}$ (binary outcomes)

Each feature:

$$p(x_i = 1) = \pi_i, \quad (17)$$

$$p(x_i = 0) = 1 - \pi_i, \quad (18)$$

$$\Rightarrow \text{can be written as: } p(x_i) = \pi_i^{x_i} (1 - \pi_i)^{1-x_i} \quad (19)$$

Conditional probabilities for two classes 1, 2:

$$p(\mathbf{x}|1) = \prod_{i=1}^D \pi_i^{x_i} (1 - \pi_i)^{1-x_i} \quad (20)$$

$$p(\mathbf{x}|2) = \prod_{i=1}^D \kappa_i^{x_i} (1 - \kappa_i)^{1-x_i} \quad (\pi_i, \kappa_i \in (0, 1), x_i \in \{0, 1\}) \quad (21)$$

Other models with linear log odds (1)

$$p(\mathbf{x}|1) = \prod_{i=1}^D \pi_i^{x_i} (1 - \pi_i)^{1-x_i} \quad (22)$$

$$p(\mathbf{x}|2) = \prod_{i=1}^D \kappa_i^{x_i} (1 - \kappa_i)^{1-x_i} \quad (\pi_i, \kappa_i \in (0, 1), x_i \in \{0, 1\}) \quad (23)$$

The log odds are:

$$a(\mathbf{x}) = \ln \frac{p(\mathbf{x}|1)p(1)}{p(\mathbf{x}|2)p(2)} = \sum_{i=1}^D \{x_i \ln \pi_i + (1 - x_i) \ln(1 - \pi_i) - x_i \ln \kappa_i - (1 - x_i) \ln(1 - \kappa_i)\} \quad (24)$$

$$+ \ln \frac{p(1)}{p(2)} = \mathbf{w} \cdot \mathbf{x} + w_0 \quad (\mathbf{w} \in \mathbb{R}^D, w_0 \in \mathbb{R}) \quad (25)$$

Note that the assumption that the features are independent may be quite strong. If this assumption is true (or anyway adopted), we talk about **naive Bayes** approach.

Other models with linear log odds (1)



Example

Problem: male/female classification

x_1 : hair length $>$ 5cm

x_2 : shoe size $>$ 41

Other models with linear log odds (2)

Multinomial naive Bayes

The analysis is similar to the case of binary outcomes. Here, the feature components x_i are not binary but they represent counts, summing to certain constant n ($\sum_{i=1}^D x_i = n$). The probabilities of observing the counts $\{x_i, i = 1, 2, \dots, D\}$ are

$$p(\mathbf{x}|1) = \frac{n!}{\prod_{i=1}^D x_i!} \prod_{i=1}^D \pi_i^{x_i} \quad (26)$$

$$p(\mathbf{x}|2) = \frac{n!}{\prod_{i=1}^D x_i!} \prod_{i=1}^D \kappa_i^{x_i} \quad (x_i \in \mathbb{N}_0, \pi_i > 0, \kappa_i > 0, \sum_i \pi_i = 1, \sum_i \kappa_i = 1, \sum_{i=1}^D x_i = n) \quad (27)$$

It is easy to see that the log odds $a(\mathbf{x})$ are again linear in \mathbf{x} .

Summary

In many real-world problems, the log odds $a(\mathbf{x})$ are a linear function of the observations \mathbf{x} .

Logistic Regression, Model

Idea: Let us look for the log of the ratio of posteriors (log odds) $a(\mathbf{x})$ directly as a linear function of the input vector $\mathbf{x} = (x_1, x_2, \dots, x_D) \in \mathbb{R}^D$ (D is the dimensionality of the feature space):

$$a(\mathbf{x}) = \ln \frac{p(1|\mathbf{x})}{p(2|\mathbf{x})} = \mathbf{w} \cdot \mathbf{x} + w_0, \quad \mathbf{w} = (w_1, w_2, \dots, w_D) \in \mathbb{R}^D, \\ w_0 \in \mathbb{R}, \tag{28}$$

where w_0 is the bias term.

Let us rewrite this as

$$a(\mathbf{x}) = w_0 \cdot \underset{\substack{\uparrow \\ x_0}}{1} + \mathbf{w} \cdot \mathbf{x} = [w_0, w_1, w_2, \dots, w_D] \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_D \end{bmatrix} = \mathbf{w}' \cdot \mathbf{x}'. \quad (x_0 = 1) \tag{29}$$

Note: From now on, we will drop the dash sign and write again only ' \mathbf{x} ' or ' \mathbf{w} ', with the understanding that these **include the zero-index components** $x_0 = 1$ and $w_0 \in \mathbb{R}$ implementing the **bias**.

Logistic Regression, Log Odds and Posteriors (1)

Here is the relationship between the the log odds $a(\mathbf{x})$ and the posterior probabilities $p(1|\mathbf{x})$ and $p(2|\mathbf{x})$.

The log odds $a(\mathbf{x})$ is (remember the bias term is consumed in the \mathbf{x} and \mathbf{w})

$$a(\mathbf{x}) = \ln \frac{p(1|\mathbf{x})}{p(2|\mathbf{x})} = \mathbf{w} \cdot \mathbf{x}. \quad (30)$$

From this, it follows that

$$\frac{p(1|\mathbf{x})}{p(2|\mathbf{x})} = \exp(\mathbf{w} \cdot \mathbf{x}) \quad \frac{p(2|\mathbf{x})}{p(1|\mathbf{x})} = \exp(-\mathbf{w} \cdot \mathbf{x}) \quad (31)$$

$$p(1|\mathbf{x}) = p(2|\mathbf{x}) \exp(\mathbf{w} \cdot \mathbf{x}) \quad p(2|\mathbf{x}) = p(1|\mathbf{x}) \exp(-\mathbf{w} \cdot \mathbf{x}) \quad (32)$$

$$1 = p(1|\mathbf{x}) + p(2|\mathbf{x}) = p(1|\mathbf{x}) (1 + e^{-\mathbf{w} \cdot \mathbf{x}}) = p(2|\mathbf{x}) (1 + e^{\mathbf{w} \cdot \mathbf{x}}) \quad (33)$$

$$p(1|\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}}, \quad p(2|\mathbf{x}) = \frac{1}{1 + e^{\mathbf{w} \cdot \mathbf{x}}}. \quad (34)$$

Logistic Regression, Log Odds and Posteriors (2)

Again,

$$a(\mathbf{x}) = \ln \frac{p(1|\mathbf{x})}{p(2|\mathbf{x})} = \mathbf{w} \cdot \mathbf{x} \quad (35)$$

$$p(1|\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}} = \sigma(\mathbf{w} \cdot \mathbf{x}) \quad (36)$$

$$p(2|\mathbf{x}) = \frac{1}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} = \sigma(-\mathbf{w} \cdot \mathbf{x}) \quad (37)$$

where $\sigma(u) = 1/(1 + \exp(-u))$ is the **logistic sigmoid** function.

It will be advantageous to rename the classes from $(1, 2)$ to $(1, -1)$. Then we can rewrite the equations (36, 37) as

$$p(k|\mathbf{x}) = \frac{1}{1 + e^{-k\mathbf{w} \cdot \mathbf{x}}}, \quad k \in \{-1, 1\} \quad (38)$$

Finding \mathbf{w} : Objective $E(\mathbf{w})$

$$p(k|\mathbf{x}) = \frac{1}{1 + e^{-k\mathbf{w}\cdot\mathbf{x}}}, \quad k \in \{-1, 1\} \quad (39)$$

How do we find \mathbf{w} ?

We apply the Maximum Likelihood approach for finding \mathbf{w} . Let us have the training set $\mathcal{T} = \{(\mathbf{x}_1, k_1), (\mathbf{x}_2, k_2), \dots, (\mathbf{x}_N, k_N)\}$. The optimal \mathbf{w}^* would be the one maximizing the log-likelihood $l(\mathbf{w})$,

$$l(\mathbf{w}) = \sum_{(\mathbf{x}, k) \in \mathcal{T}} \ln p(\mathbf{x}, k)_{[\mathbf{w}]} = \sum_{(\mathbf{x}, k) \in \mathcal{T}} \ln p(k|\mathbf{x})_{[\mathbf{w}]} + \sum_{(\mathbf{x}, k) \in \mathcal{T}} \ln p(\mathbf{x})_{[\mathbf{w}]}, \quad (40)$$

where for the sake of clarity, dependence on \mathbf{w} is denoted by a subscript $[\mathbf{w}]$.

As there are no assumptions about the form of $p(\mathbf{x})$ as a function of \mathbf{w} , logistic regression employs instead the maximization of *conditional log-likelihood* $l'(\mathbf{w})$

$$l'(\mathbf{w}) = \sum_{(\mathbf{x}, k) \in \mathcal{T}} \ln p(k|\mathbf{x})_{[\mathbf{w}]} = - \sum_{(\mathbf{x}, k) \in \mathcal{T}} \ln(1 + e^{-k\mathbf{w}\cdot\mathbf{x}}) \quad (\text{conditional log likelihood}) \quad (41)$$

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmax}} l'(\mathbf{w}) \quad (\text{optimal } \mathbf{w}^*) \quad (42)$$

Finding \mathbf{w} : Objective $E(\mathbf{w})$

(copied from previous slide)

$$l'(\mathbf{w}) = \sum_{(\mathbf{x},k) \in \mathcal{T}} \ln p(k|\mathbf{x})_{[\mathbf{w}]} = - \sum_{(\mathbf{x},k) \in \mathcal{T}} \ln(1 + e^{-k\mathbf{w} \cdot \mathbf{x}}) \quad (\text{conditional log likelihood}) \quad (43)$$

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmax}} l'(\mathbf{w}) \quad (\text{optimal } \mathbf{w}^*) \quad (44)$$

In order for the optimization to fit into the **minimization** framework, we define the objective function $E(\mathbf{w})$ as the negative conditional log likelihood, $E(\mathbf{w}) = -l'(\mathbf{w})$. This objective function corresponds to **cross entropy**. Let us now analyze the properties of $E(\mathbf{w})$.

Finding \mathbf{w} : Gradient of $E(\mathbf{w})$

$$E(\mathbf{w}) = - \sum_{(\mathbf{x},k) \in \mathcal{T}} \ln p(k|\mathbf{x}) = \sum_{(\mathbf{x},k) \in \mathcal{T}} \ln(1 + e^{-k\mathbf{w} \cdot \mathbf{x}}) \quad (45)$$

(46)

The gradient vector $g(\mathbf{w})$ of E is:

$$g(\mathbf{w}) = \frac{\partial E(\mathbf{w})}{\partial \mathbf{w}} = \sum_{(\mathbf{x},k) \in \mathcal{T}} \frac{e^{-k\mathbf{w} \cdot \mathbf{x}}}{1 + e^{-k\mathbf{w} \cdot \mathbf{x}}} (-k\mathbf{x}) = - \sum_{(\mathbf{x},k) \in \mathcal{T}} \underbrace{\frac{1}{1 + e^{k\mathbf{w} \cdot \mathbf{x}}}}_{p(-k|\mathbf{x})} k\mathbf{x} \quad (47)$$

$$= - \sum_{(\mathbf{x},k) \in \mathcal{T}} (1 - p(k|\mathbf{x})) k\mathbf{x}. \quad (48)$$

We require $g(\mathbf{w}) = 0$ (the necessary condition for optimality). However, it seems that these equations **cannot be solved analytically**. We will need to resort to the numerical optimization methods. Let us continue and check the second order derivatives.

Finding \mathbf{w} : $E(\mathbf{w})$ is convex

The Hessian matrix $H(\mathbf{w})$ of the objective function E is

$$H(\mathbf{w}) = \frac{\partial^2 E(\mathbf{x})}{\partial \mathbf{w}^2} = \frac{\partial g(\mathbf{w})}{\partial \mathbf{w}} = -\frac{\partial}{\partial \mathbf{w}} \sum_{(\mathbf{x}, k) \in \mathcal{T}} \frac{1}{1 + e^{k\mathbf{w} \cdot \mathbf{x}}} k\mathbf{x} \quad (49)$$

$$= \sum_{(\mathbf{x}, k) \in \mathcal{T}} \underbrace{\frac{e^{k\mathbf{w} \cdot \mathbf{x}}}{(1 + e^{k\mathbf{w} \cdot \mathbf{x}})^2} k^2}_{> 0} \mathbf{x}\mathbf{x}^\top = \sum_{(\mathbf{x}, k) \in \mathcal{T}} p(-1|\mathbf{x})p(1|\mathbf{x}) \mathbf{x}\mathbf{x}^\top \quad (50)$$

This is a very important result. It shows that the Hessian matrix $H(\mathbf{w})$ is **positive definite** in every point \mathbf{w} and, therefore, the function $E(\mathbf{w})$ is **convex**. As a consequence, $E(\mathbf{w})$ has a **unique minimum**.

Note. Can you show that $H(\mathbf{w})$ is positive definite?

Finding w : Gradient Descent

Any method of convex optimization can be used to find the optimal w^* . For the examples in this lecture, the following gradient descent method with adaptive step size has been used:

```
# input: x (observations), k (class labels), w_init (initial w)

# init:
w = w_init
step_size = 1.0
E, g = compute_E_and_gradient(x, k, w)

# iterate:
while not TERMINATION_CONDITION:
    E_new, g_new = compute_E_and_gradient(x, k, w - step_size * g)

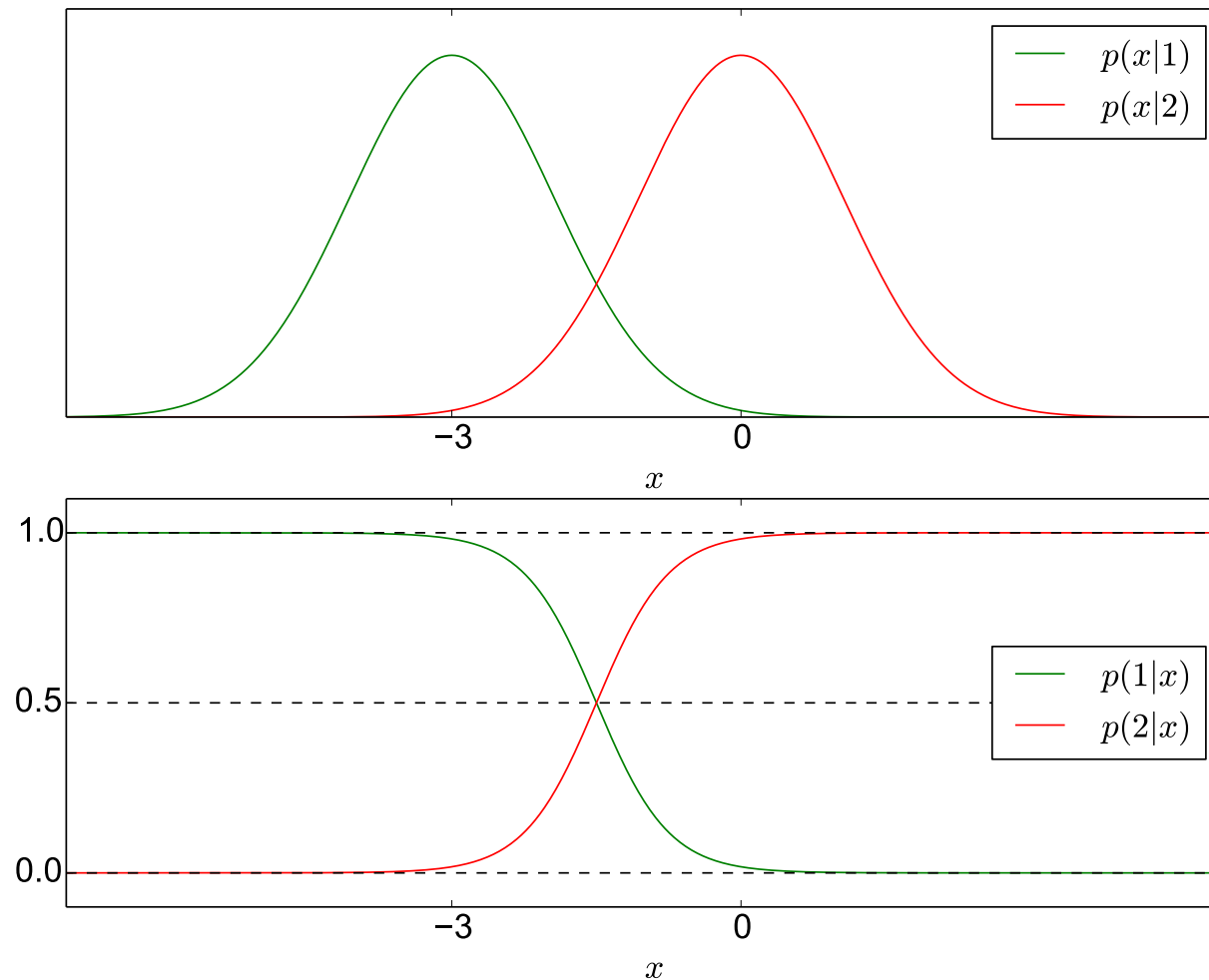
    if E_new < E:
        # success.
        w -= step_size * g
        g = g_new
        E = E_new
        step_size *= 2
    else:
        step_size /= 2

return w
```

Notes:

- i) Iteration is accepted if $E(w)$ decreases. If it hasn't decreased, either the step size is too high (thus it is halved), or optimum has been already found.
- ii) We normalize the gradient by the number of training data N because otherwise its magnitude scales linearly with N , causing the necessity for smaller step sizes with higher N .

Example 1, Two Normal Distributions with Equal Variance (1)



$$p(x|1) = \mathcal{N}(x|\mu_1 = -3, \sigma_1 = 1.5)$$

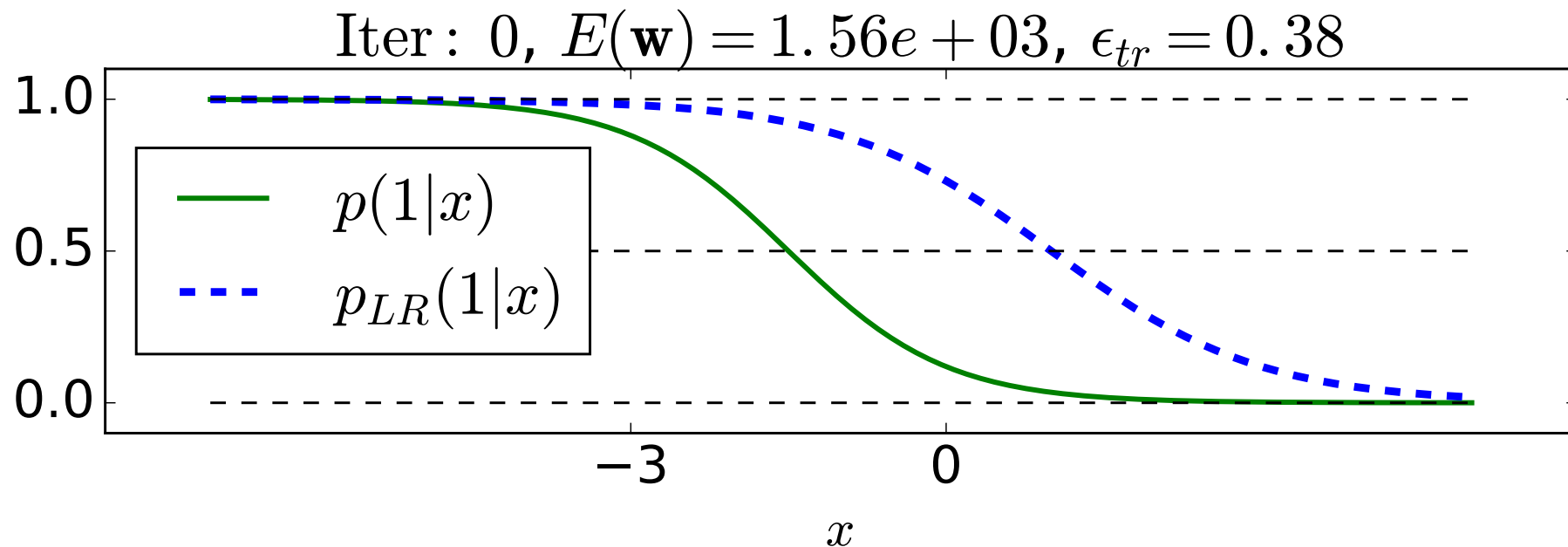
$$p(x|2) = \mathcal{N}(x|\mu_2 = 0, \sigma_2 = 1.5)$$

$$p(1) = p(2) = 0.5. \text{ Bayesian error is } \epsilon_B = 0.16.$$

Example 1, Two Normal Distributions with Equal Variance (2)

Initial state.

Training set: 1000 samples from each of the distributions.



$p(1|x)$: The actual conditional for the 1st class.

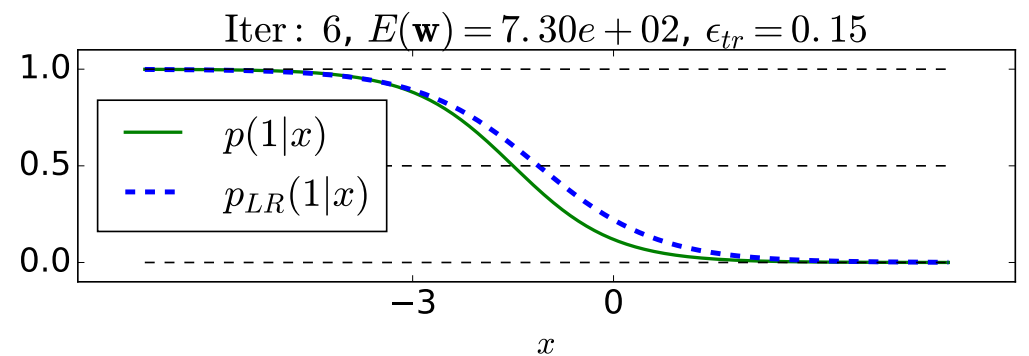
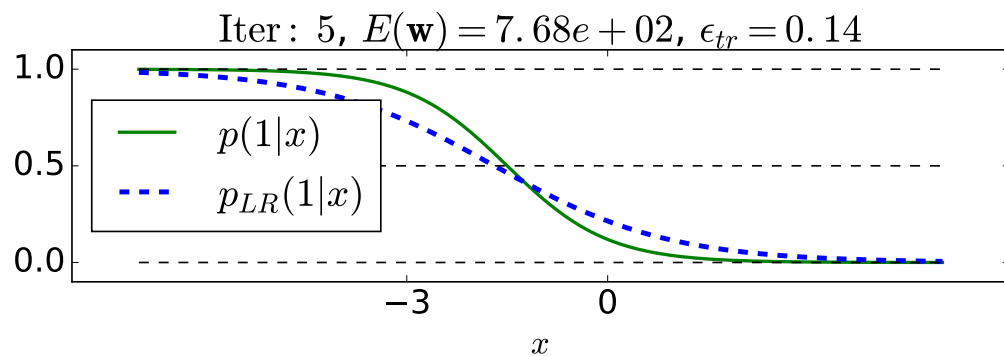
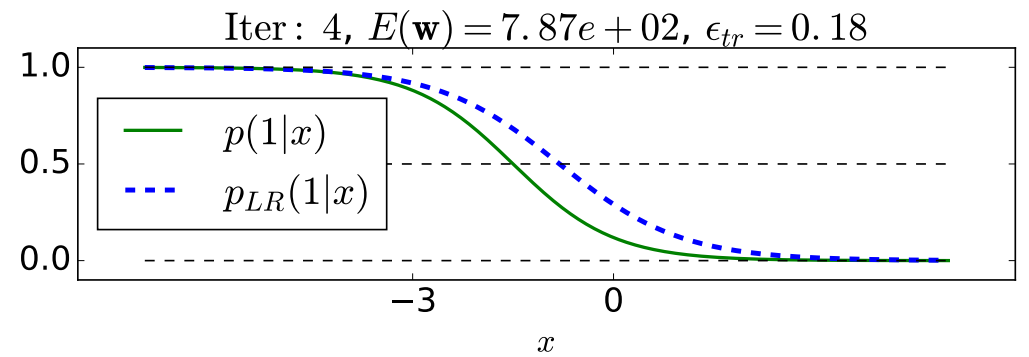
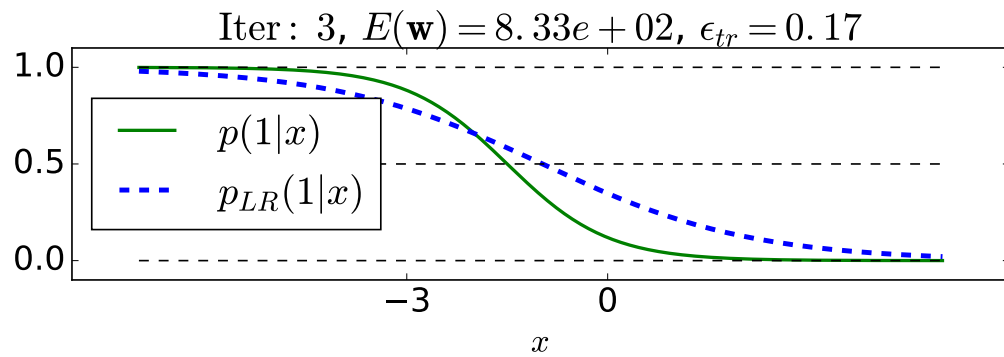
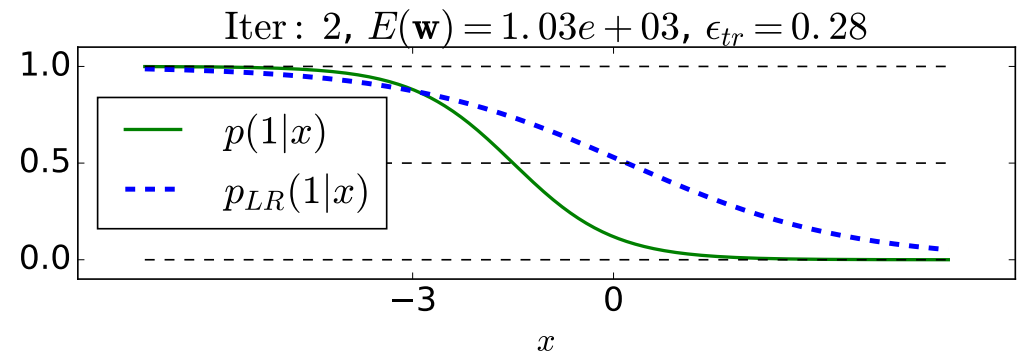
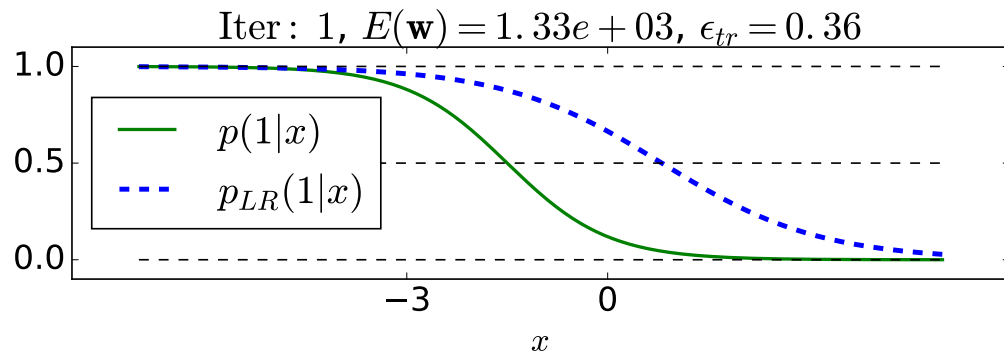
$p_{LR}(1|x)$: The conditional for the 1st class predicted by logistic regression.

$E(\mathbf{w})$: the value of cross entropy.

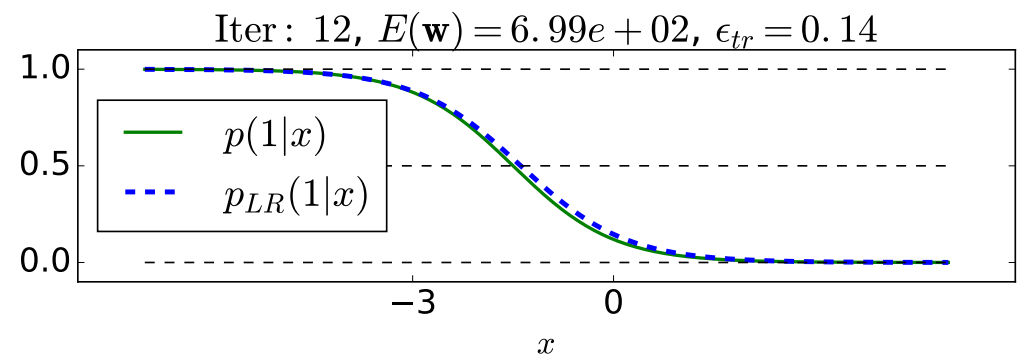
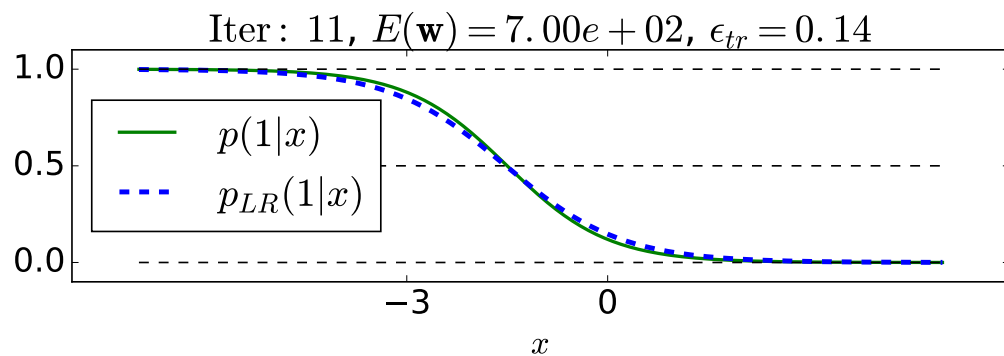
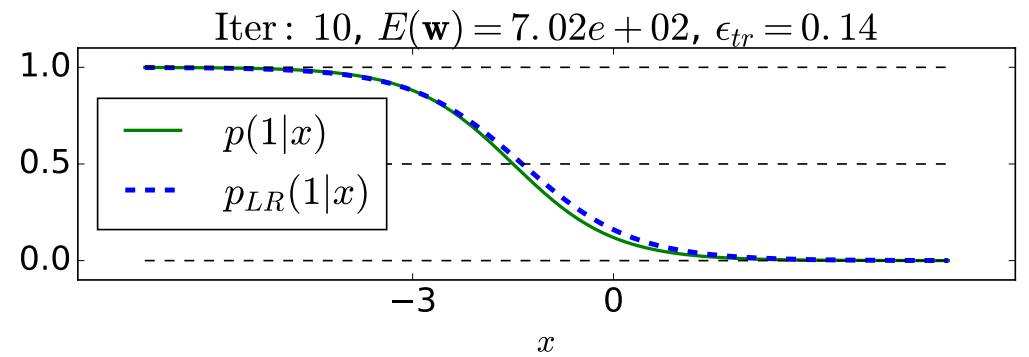
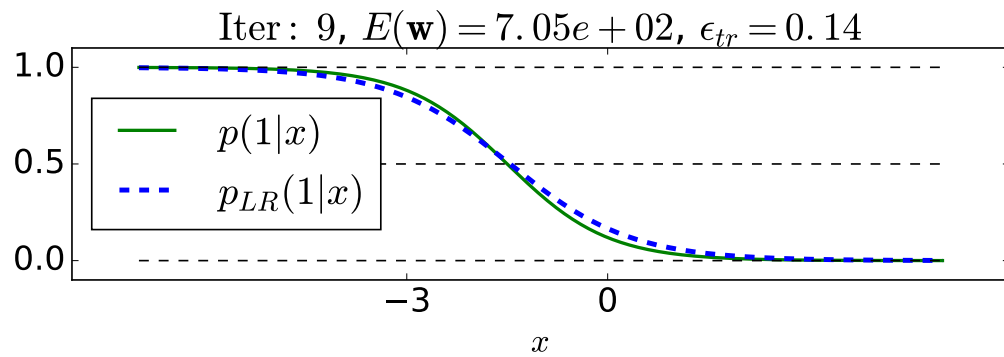
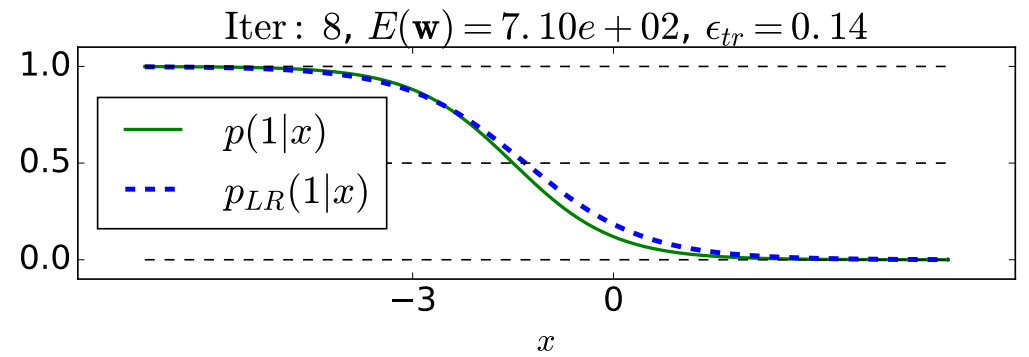
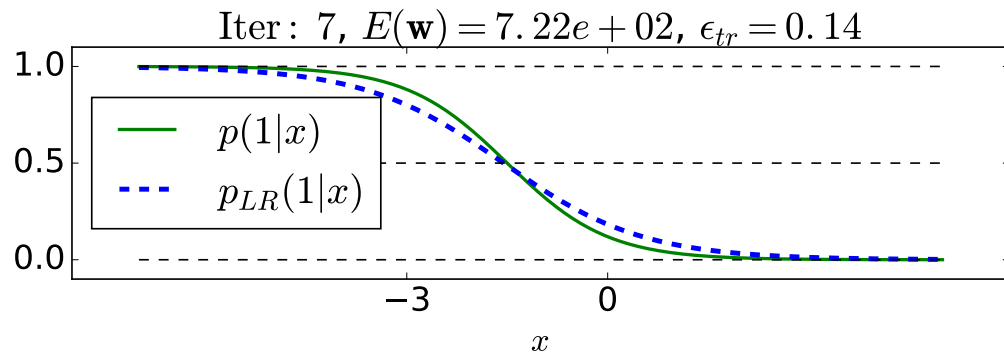
ϵ_{tr} : the training error (error on the training set.)

(initial $w = [1, -1]^T$)

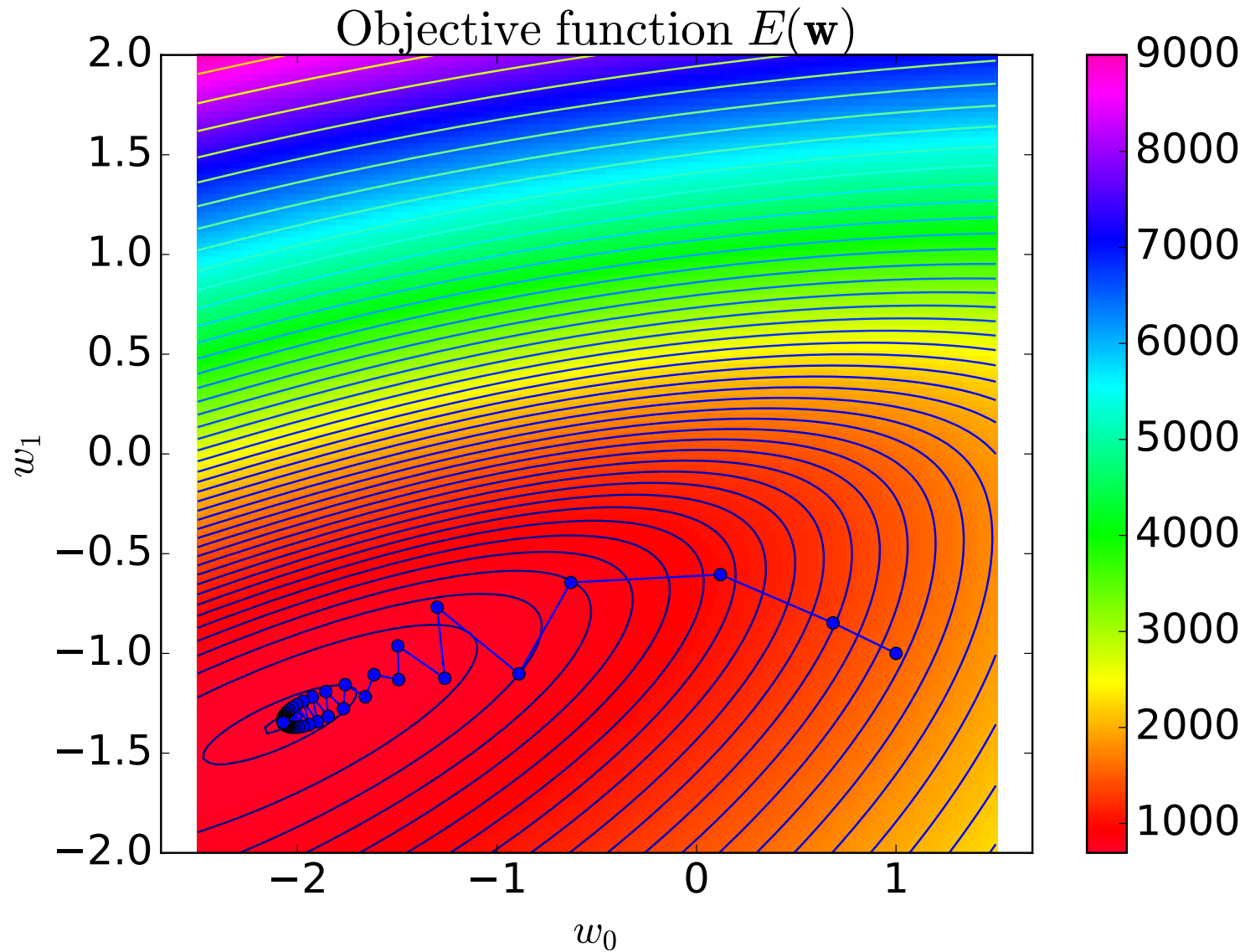
Example 1, Two Normal Distributions with Equal Variance (3)



Example 1, Two Normal Distributions with Equal Variance (4)

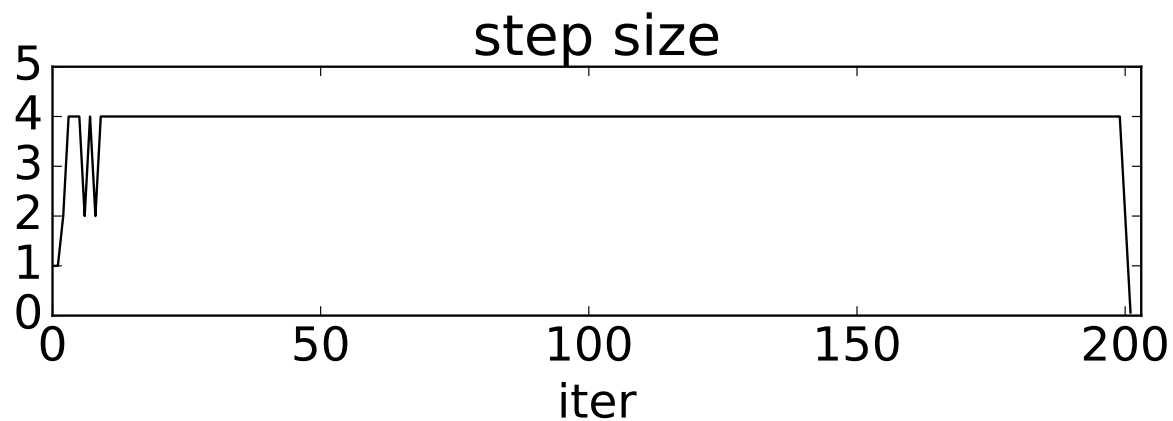
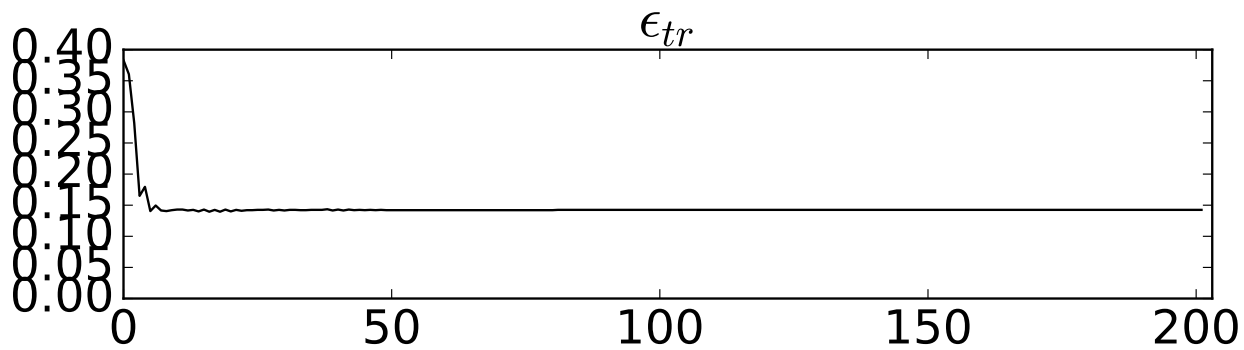
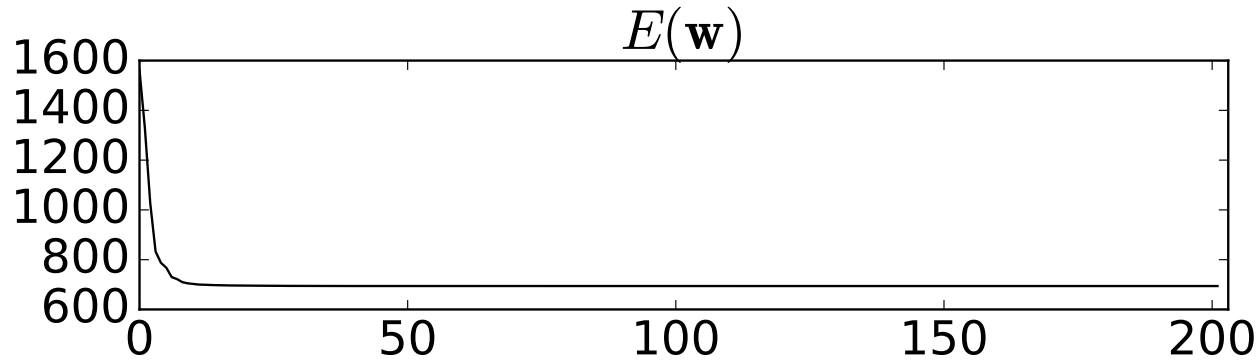


Example 1, Two Normal Distributions with Equal Variance (5)



The cross-entropy $E(\mathbf{w})$ and the progress of \mathbf{w} with iterations.

Example 1, Two Normal Distributions with Equal Variance (6)



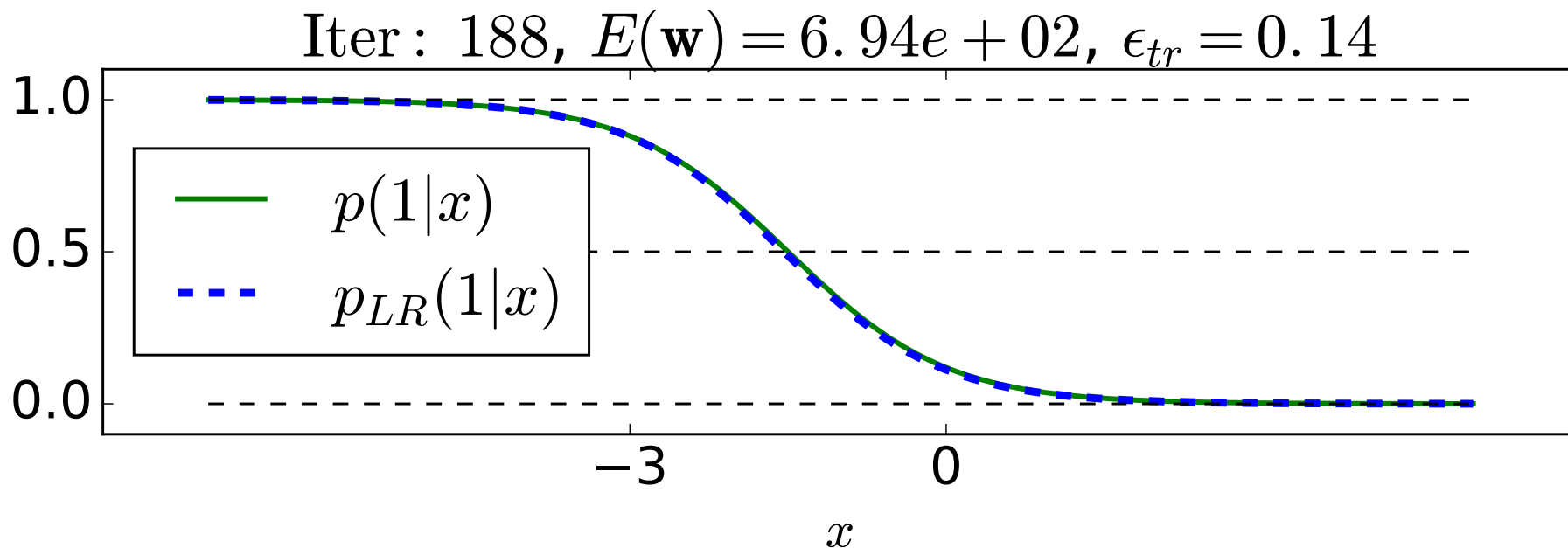
Example 1, Two Normal Distributions with Equal Variance (7)

Things to note:

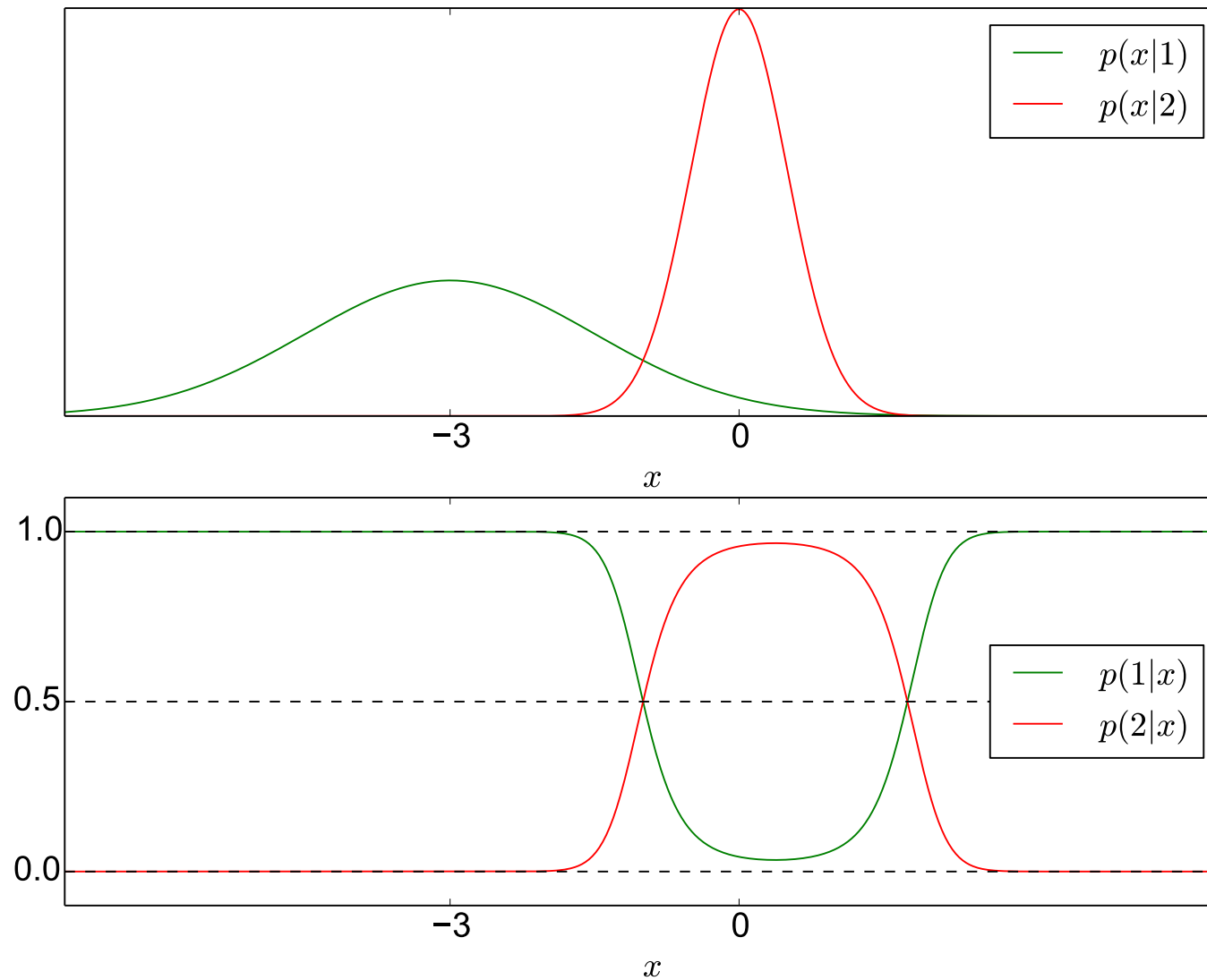
- ◆ ϵ_{tr} does not monotonically decrease with iterations. $E(\mathbf{w})$ does.
- ◆ Some intermediate ϵ_{tr} 's as well as the final one are lower than the Bayesian error $\epsilon_B = 0.16$. This is not a contradiction of the theory.

Converged state:

$$w = [-2.07, -1.35]^T.$$



Example 2, Non-Equal Variance but Different Mean (1)



$$p(x|1) = \mathcal{N}(x|\mu_1 = -3, \sigma_1 = 1.5)$$

$$p(x|2) = \mathcal{N}(x|\mu_2 = 0, \sigma_2 = 0.5)$$

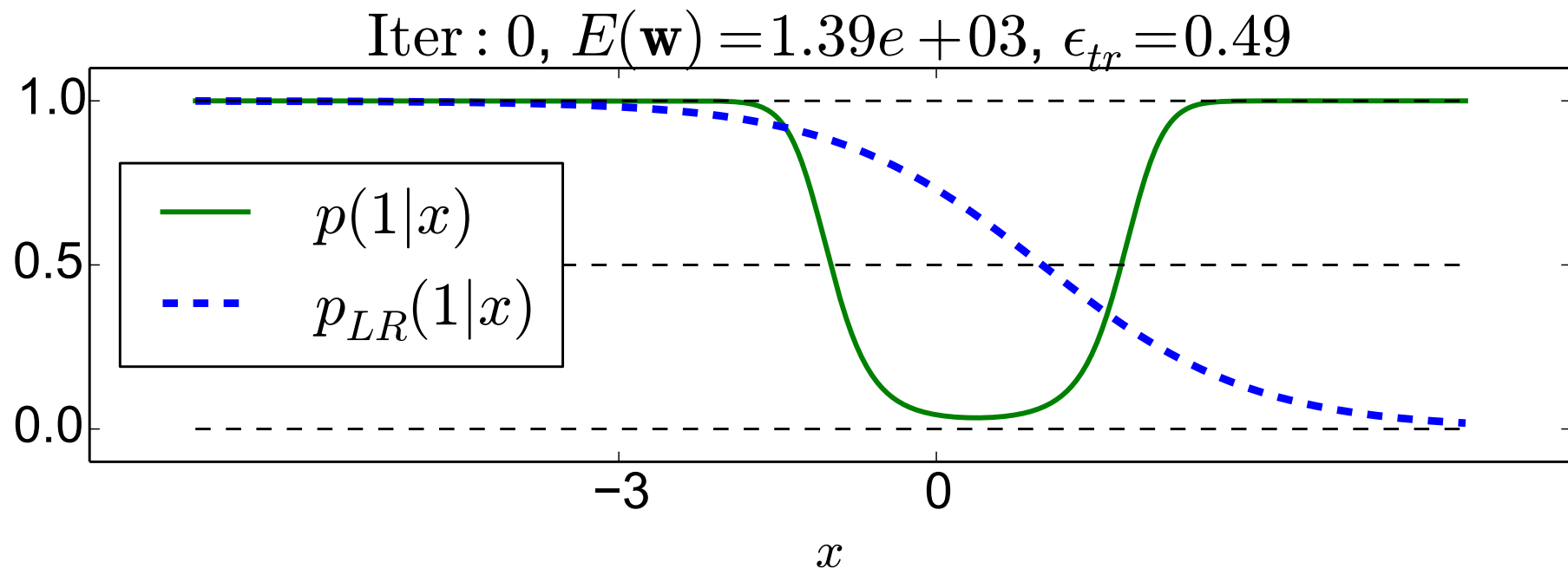
$p(1) = p(2) = 0.5$. **Bayesian error is** $\epsilon_B = 0.057$.

Example 2, Non-Equal Variance but Different Mean (2)

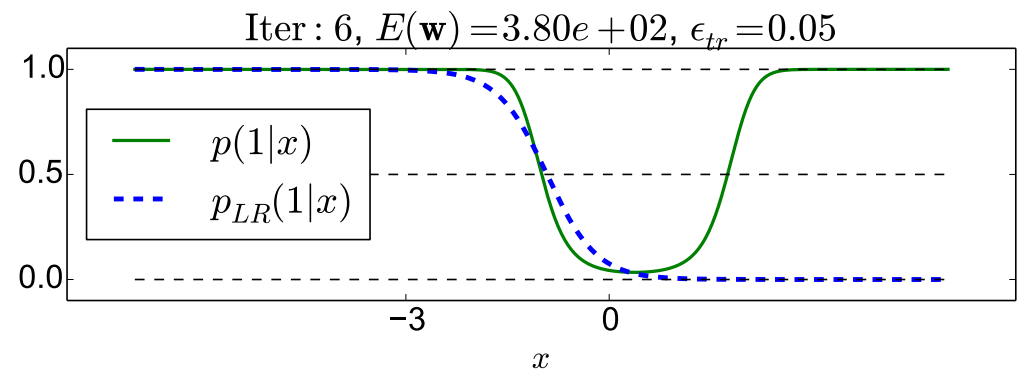
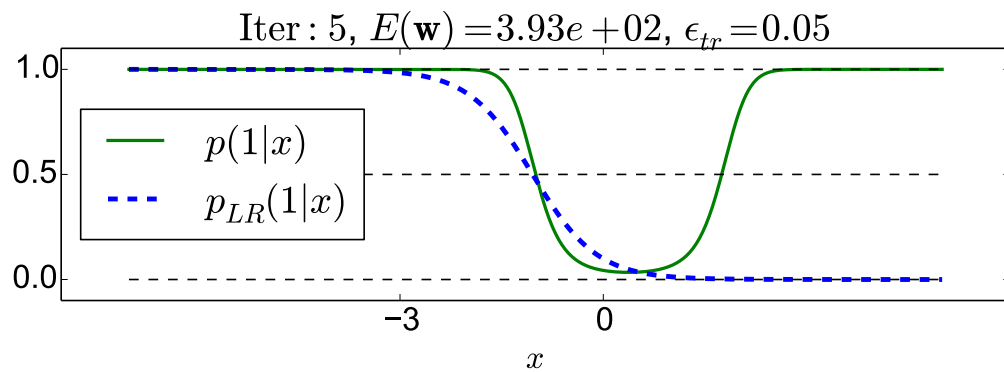
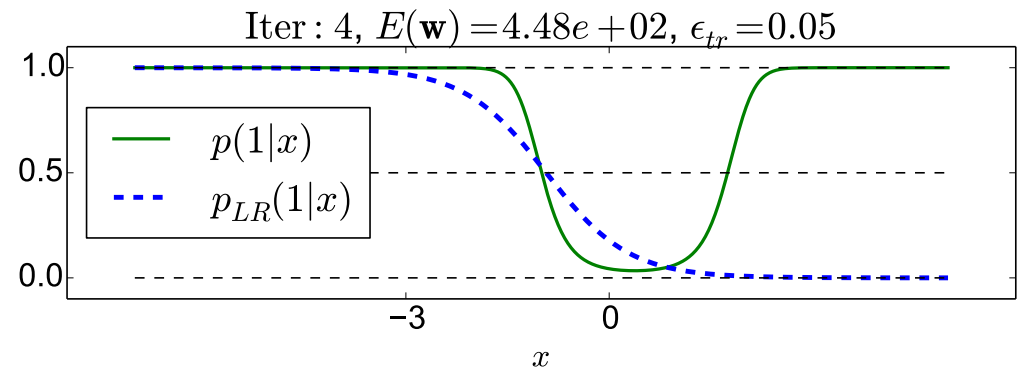
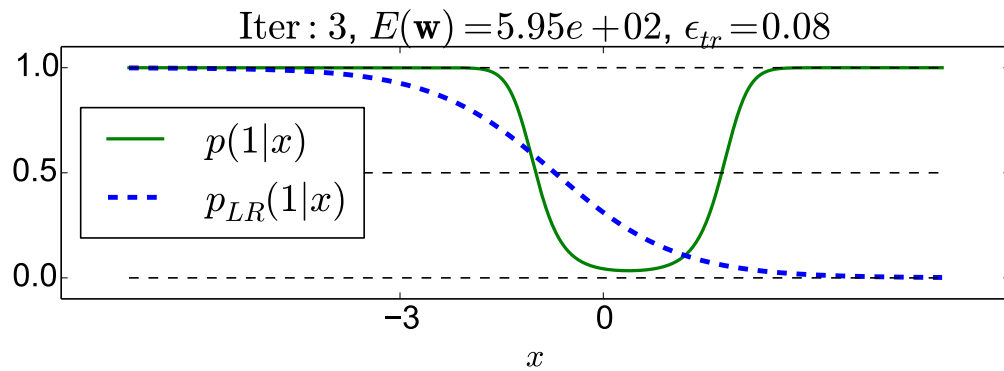
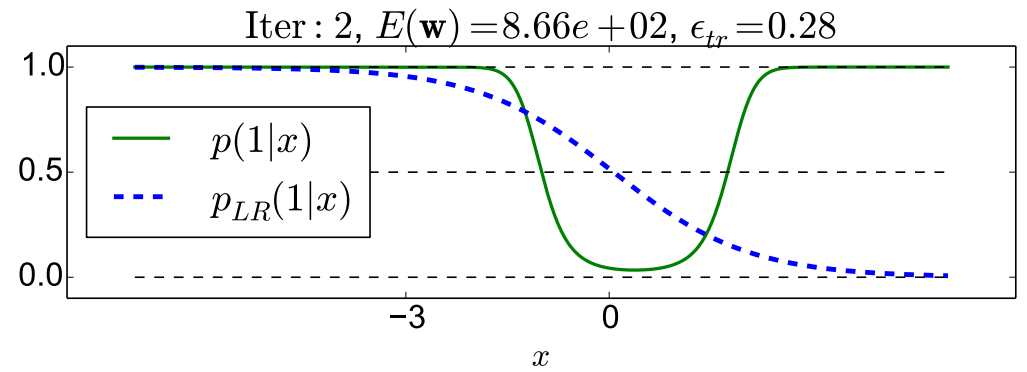
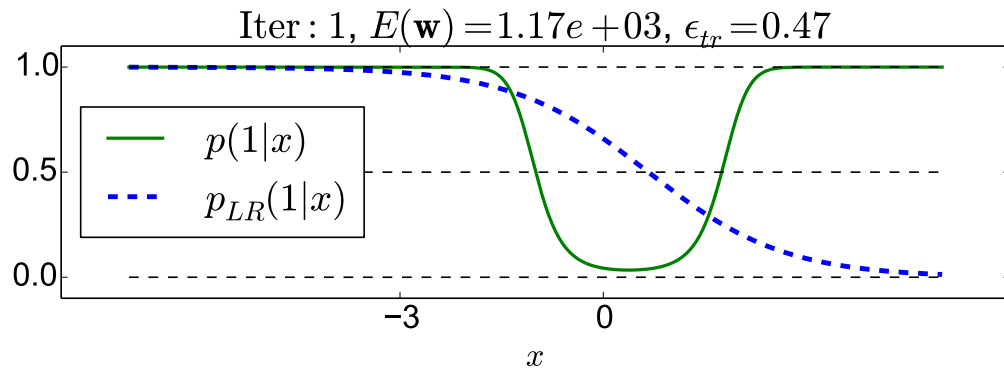
Initial state.

$$w = [1, -1]^T.$$

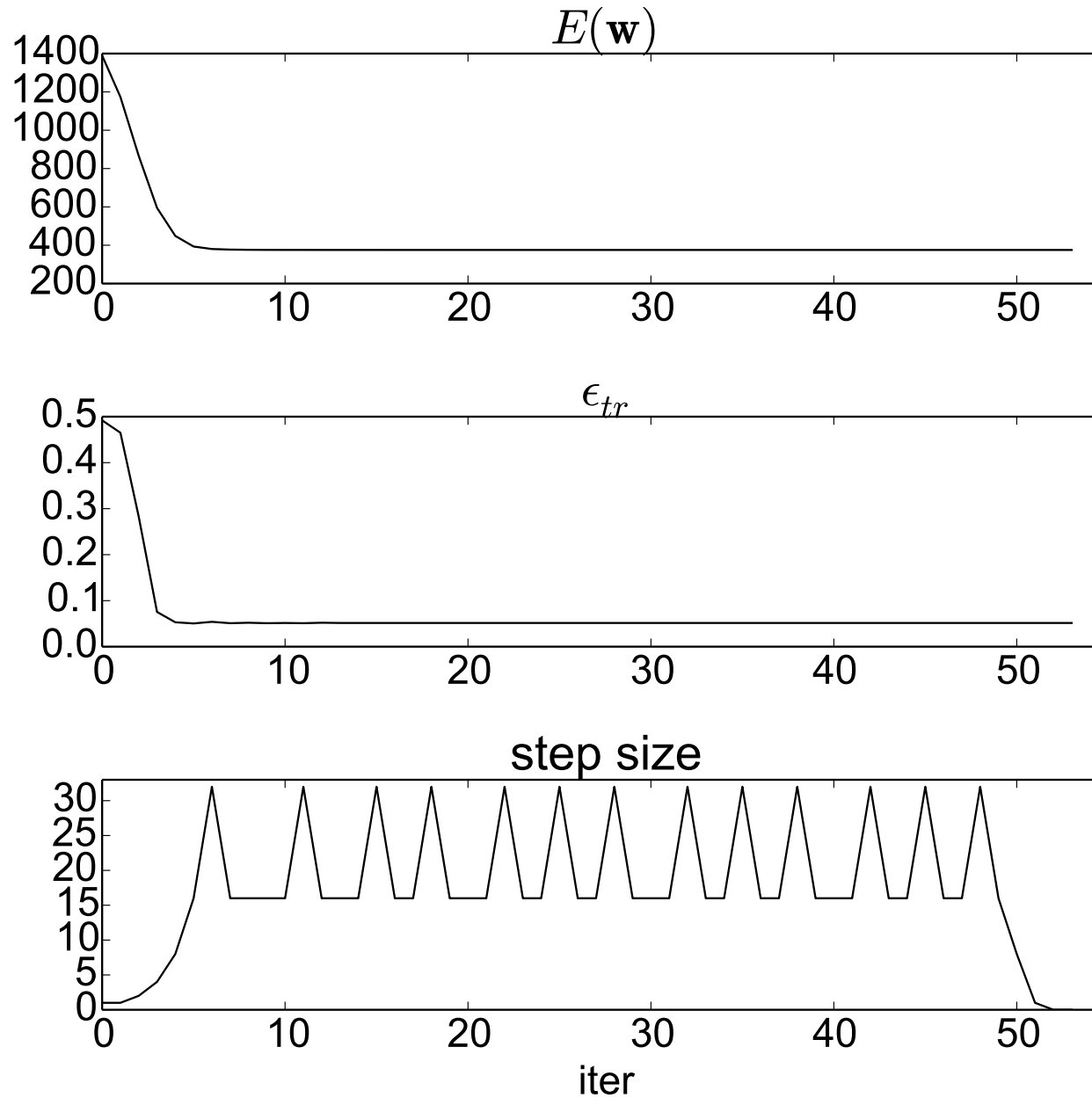
Training set: 1000 samples from each of the distributions.



Example 2, Non-Equal Variance but Different Mean (3)



Example 2, Non-Equal Variance but Different Mean (4)



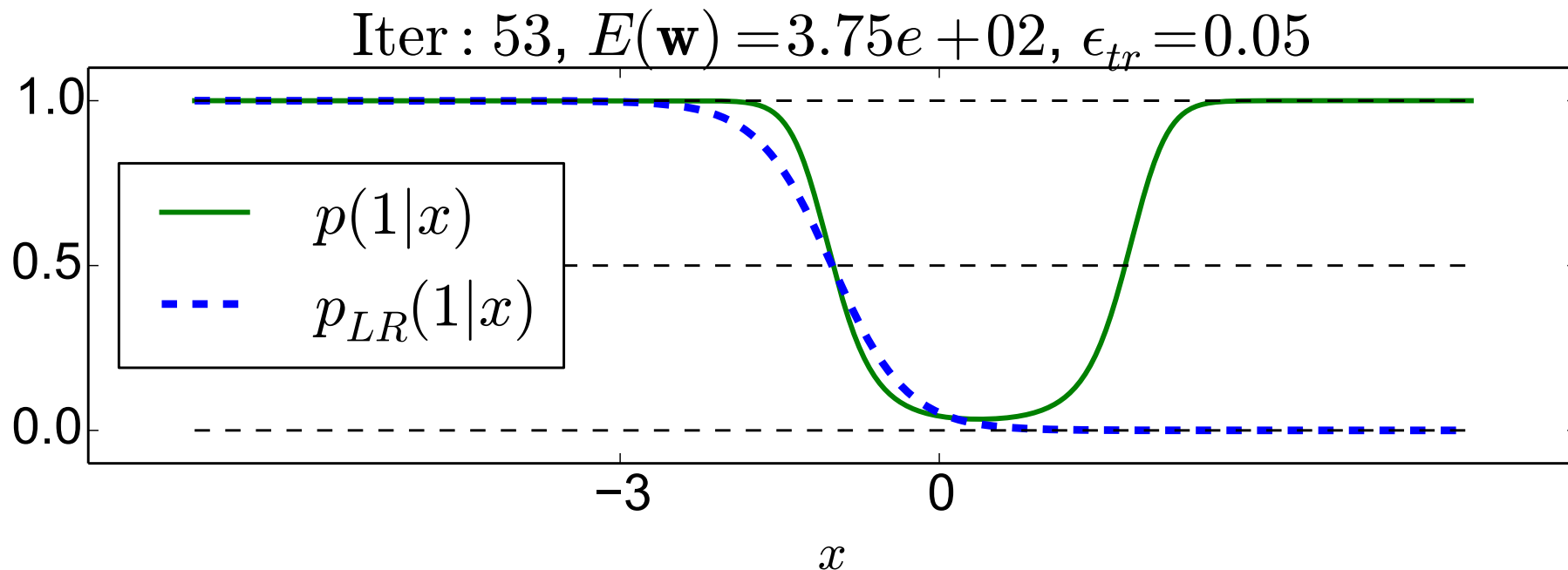
Example 2, Non-Equal Variance but Different Mean (5)

Things to note:

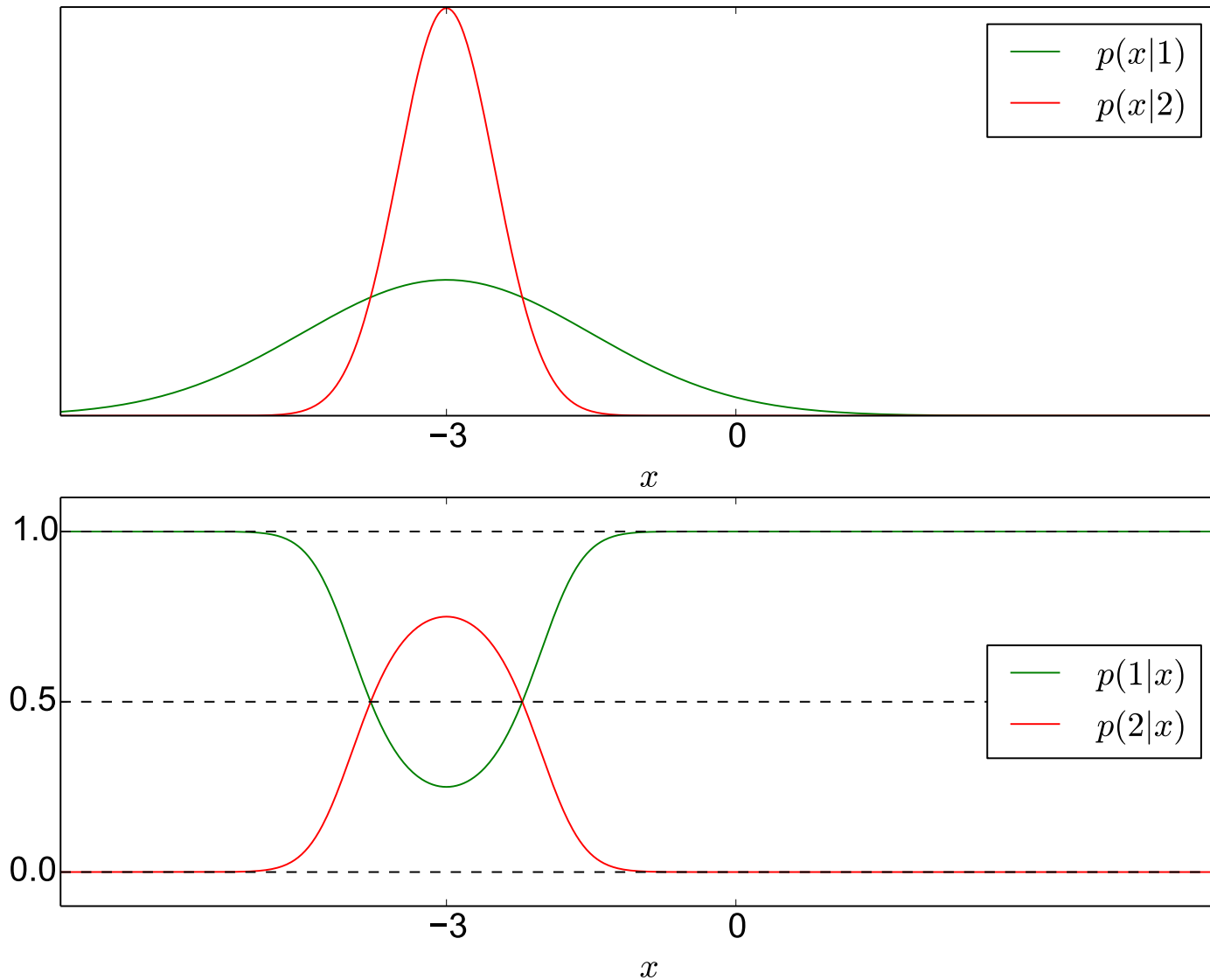
- ◆ The logistic regression cannot provide the two thresholds the optimal decision strategy requires. But it provides the one threshold which matters more in reducing the classification error (here the left one.)

Converged state.

$$w = [-2.88, -2.85]^T.$$



Example 3, Non-Equal Variance and the Same Mean (1)



$$p(x|1) = \mathcal{N}(x|\mu_1 = -3, \sigma_1 = 1.5)$$

$$p(x|2) = \mathcal{N}(x|\mu_2 = -3, \sigma_2 = 0.5)$$

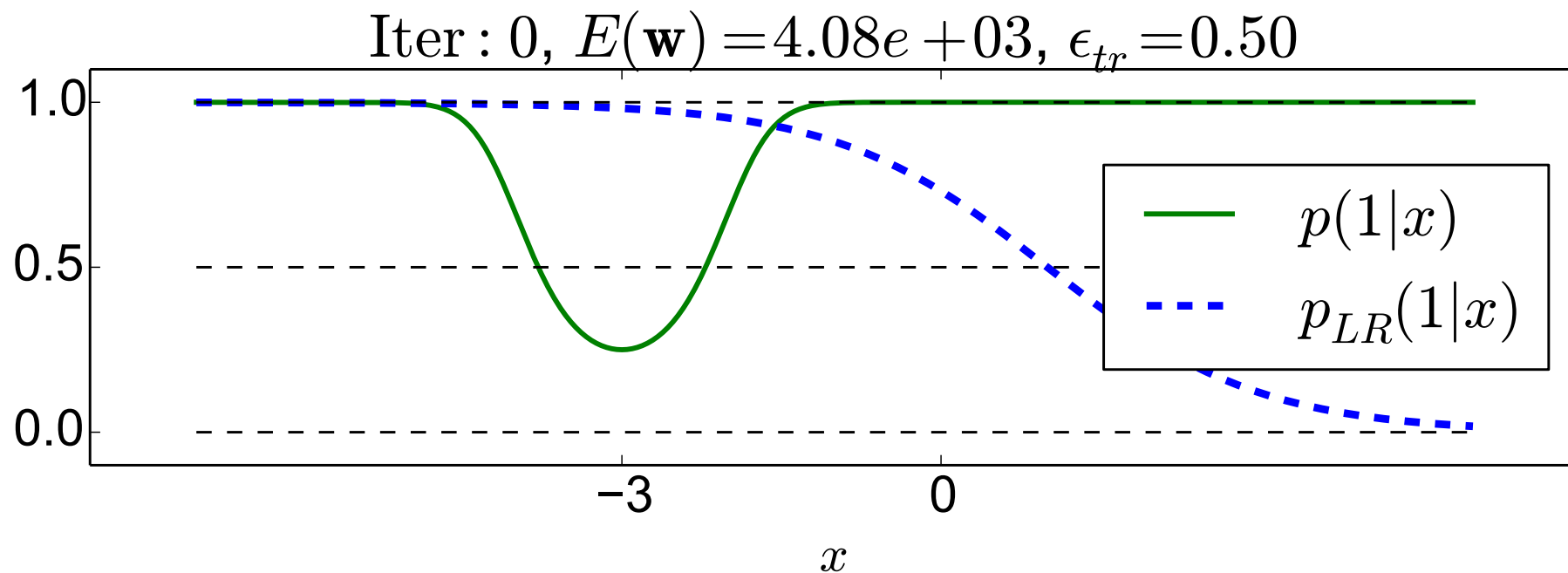
$$p(1) = p(2) = 0.5. \text{ Bayesian error is } \epsilon_B = 0.26.$$

Example 3, Non-Equal Variance and the Same Mean (2)

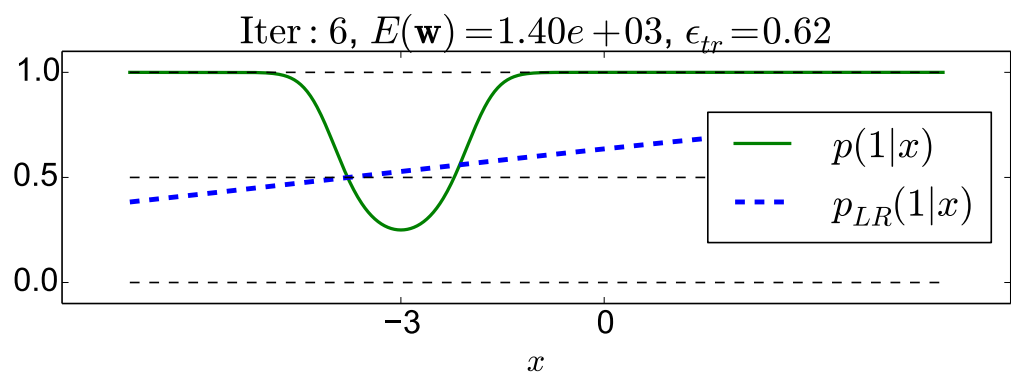
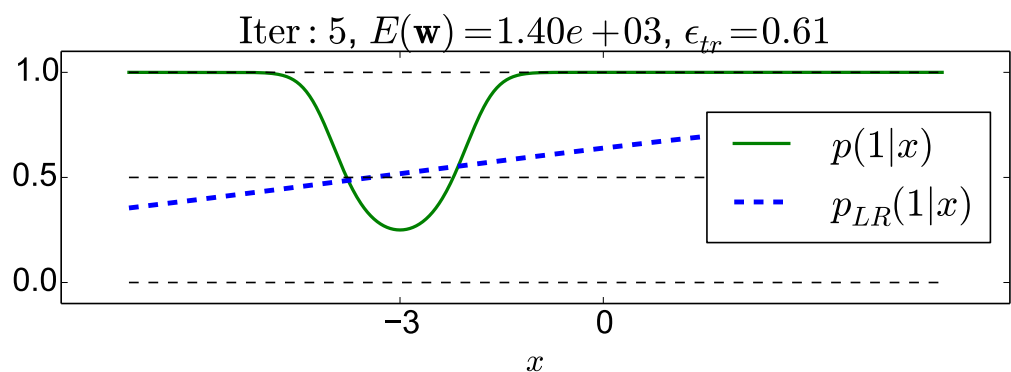
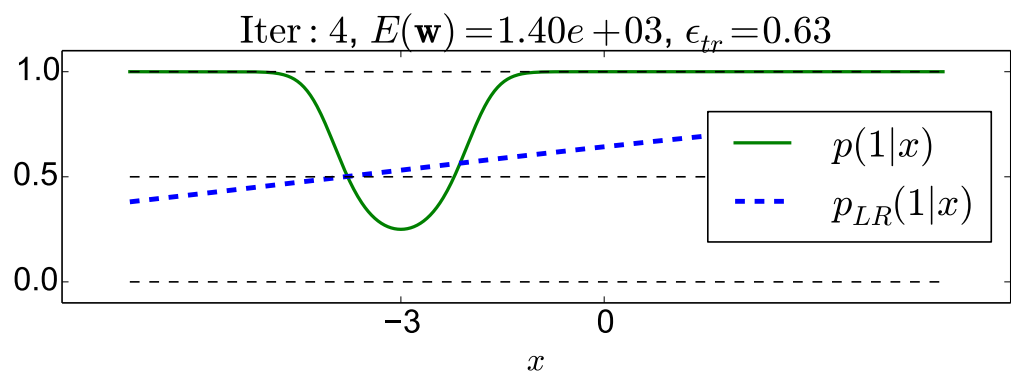
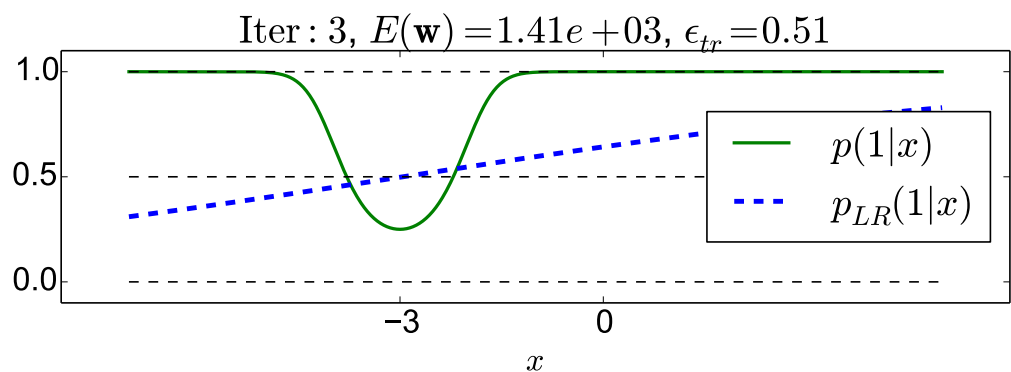
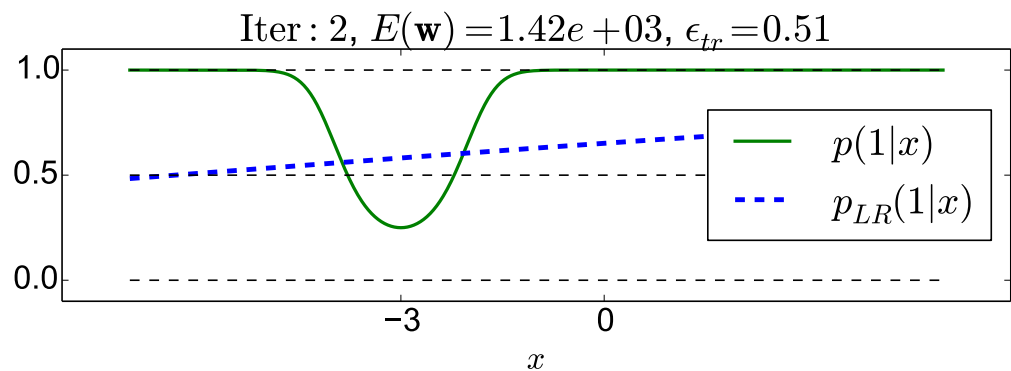
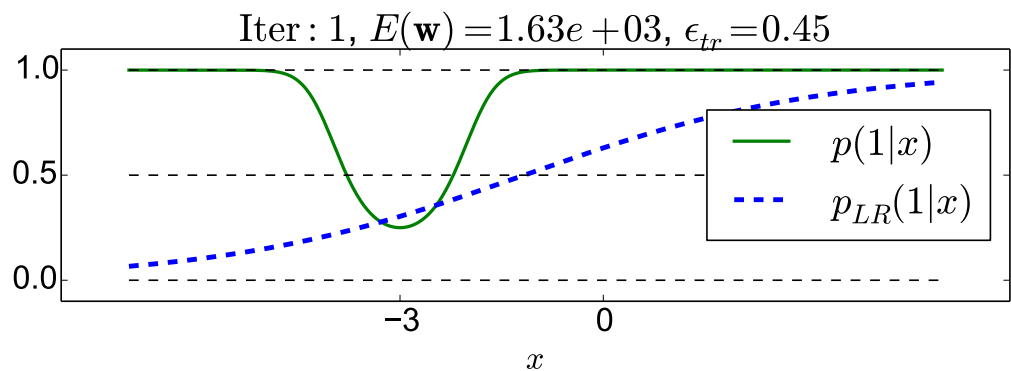
Initial state.

$$w = [1, -1]^T.$$

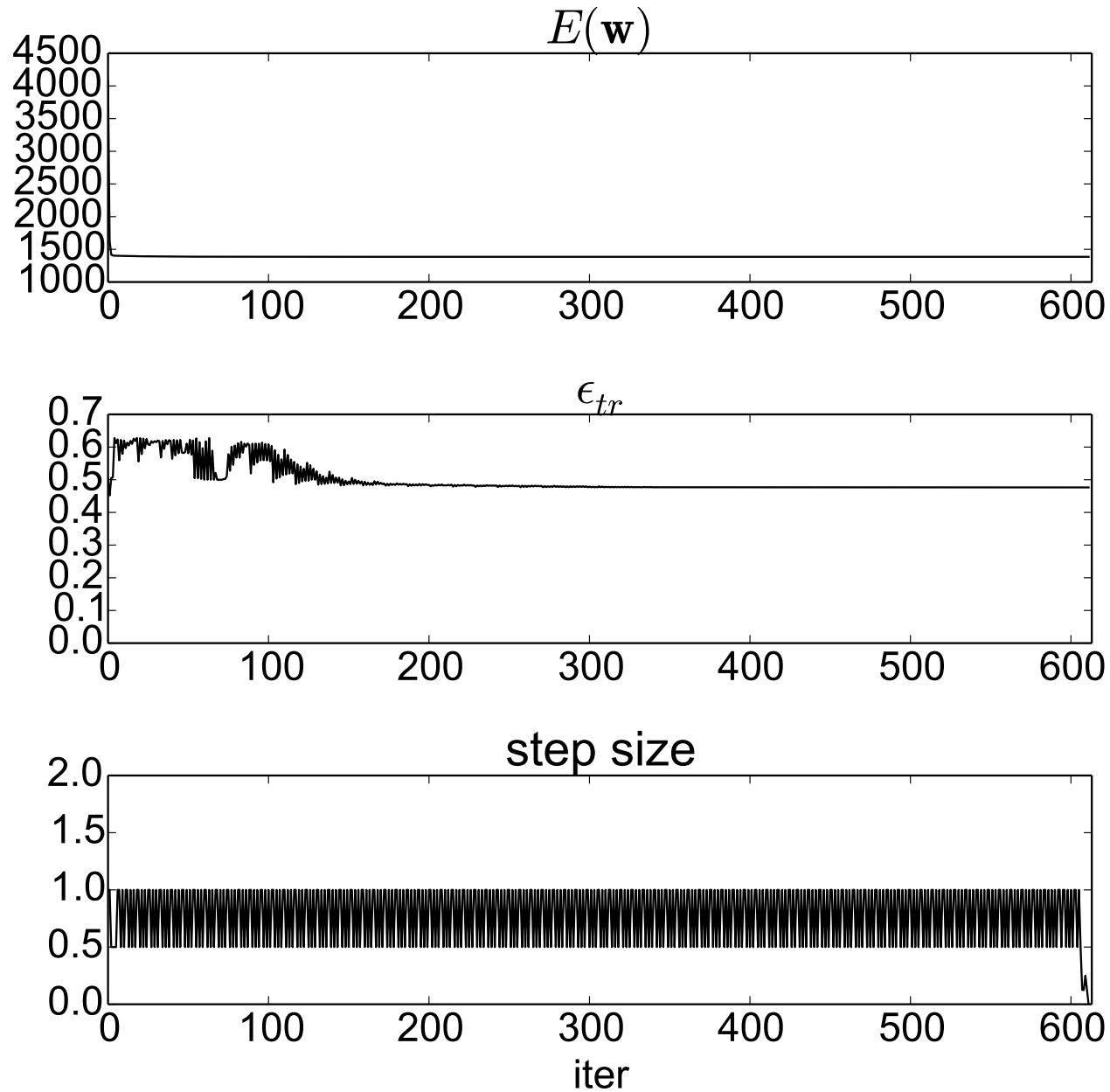
Training set: 1000 samples from each of the distributions.



Example 3, Non-Equal Variance and the Same Mean (3)



Example 3, Non-Equal Variance and the Same Mean (4)



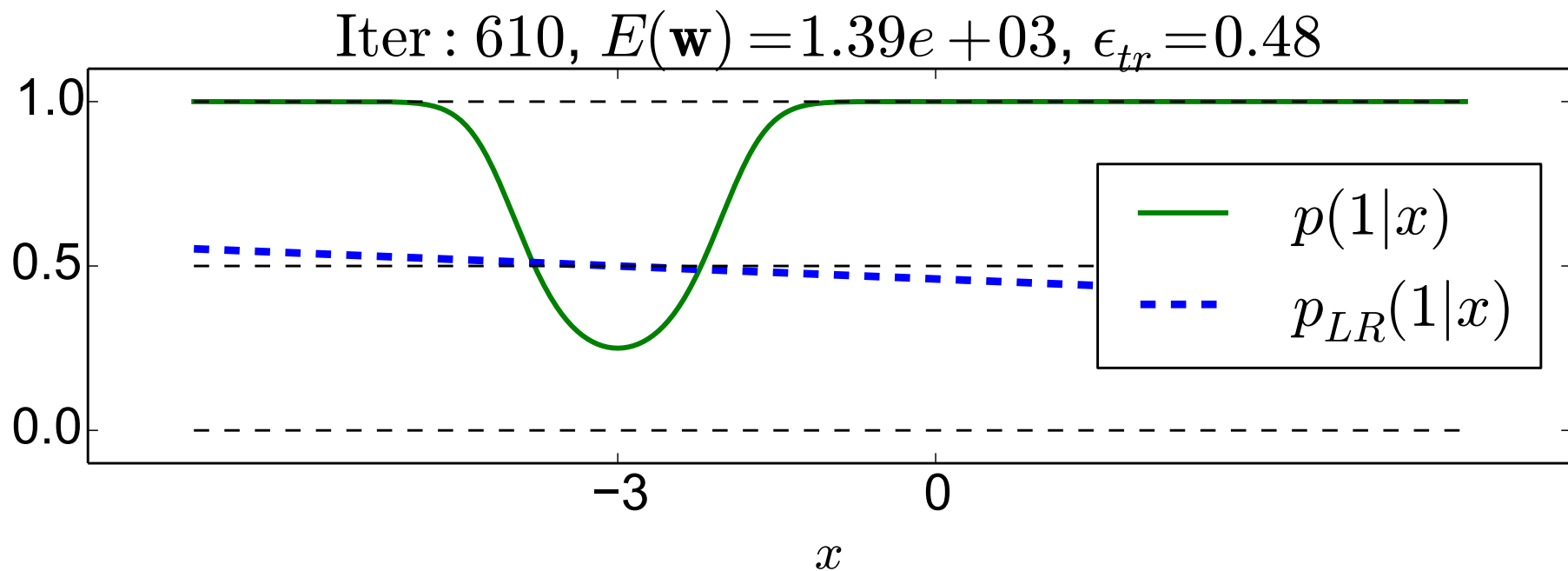
Example 3, Non-Equal Variance and the Same Mean (5)

Things to note:

- ◆ Failure case. The logistic regression cannot provide a good fit to the log odds in this case.

Final state:

$$w = [-0.161, -0.053]^T.$$



Logistic Regression with Multiple Classes (1)

The logistic regression can be generalized to multiple classes as follows.

Each class $k \in \{1, 2, \dots, K\}$ has an associated weight vector \mathbf{w}_k .

The conditional probability for the k -th function is computed using the **softmax** function:

$$p(k|\mathbf{x}) = \frac{e^{\mathbf{w}_k \cdot \mathbf{x}}}{e^{\mathbf{w}_1 \cdot \mathbf{x}} + e^{\mathbf{w}_2 \cdot \mathbf{x}} + \dots + e^{\mathbf{w}_K \cdot \mathbf{x}}}. \quad (51)$$

Things to note:

- ◆ The above term sums to 1 (summing over k .)
- ◆ For two classes only, we get the same terms as previously, with $\mathbf{w} = \mathbf{w}_1 - \mathbf{w}_2$:

$$p(1|\mathbf{x}) = \frac{e^{\mathbf{w}_1 \cdot \mathbf{x}}}{e^{\mathbf{w}_1 \cdot \mathbf{x}} + e^{\mathbf{w}_2 \cdot \mathbf{x}}} = \frac{1}{1 + e^{-(\mathbf{w}_1 - \mathbf{w}_2) \cdot \mathbf{x}}} = \sigma((\mathbf{w}_1 - \mathbf{w}_2) \cdot \mathbf{x}) \quad (52)$$

$$p(2|\mathbf{x}) = \frac{e^{\mathbf{w}_2 \cdot \mathbf{x}}}{e^{\mathbf{w}_1 \cdot \mathbf{x}} + e^{\mathbf{w}_2 \cdot \mathbf{x}}} = \frac{1}{1 + e^{(\mathbf{w}_1 - \mathbf{w}_2) \cdot \mathbf{x}}} = \sigma(-(\mathbf{w}_1 - \mathbf{w}_2) \cdot \mathbf{x}) \quad (53)$$

Logistic Regression with Multiple Classes (2)

Conditional probabilities:

$$p(k|\mathbf{x}) = \frac{e^{\mathbf{w}_k \cdot \mathbf{x}}}{e^{\mathbf{w}_1 \cdot \mathbf{x}} + e^{\mathbf{w}_2 \cdot \mathbf{x}} + \dots + e^{\mathbf{w}_K \cdot \mathbf{x}}}. \quad (54)$$

The training set $\mathcal{T} = \{(\mathbf{x}_1, k_1), (\mathbf{x}_2, k_2), \dots, (\mathbf{x}_N, k_N)\}$ (as before);

The set of parameters to find: $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K]$;

The conditional log likelihood $l'(\mathbf{W})$:

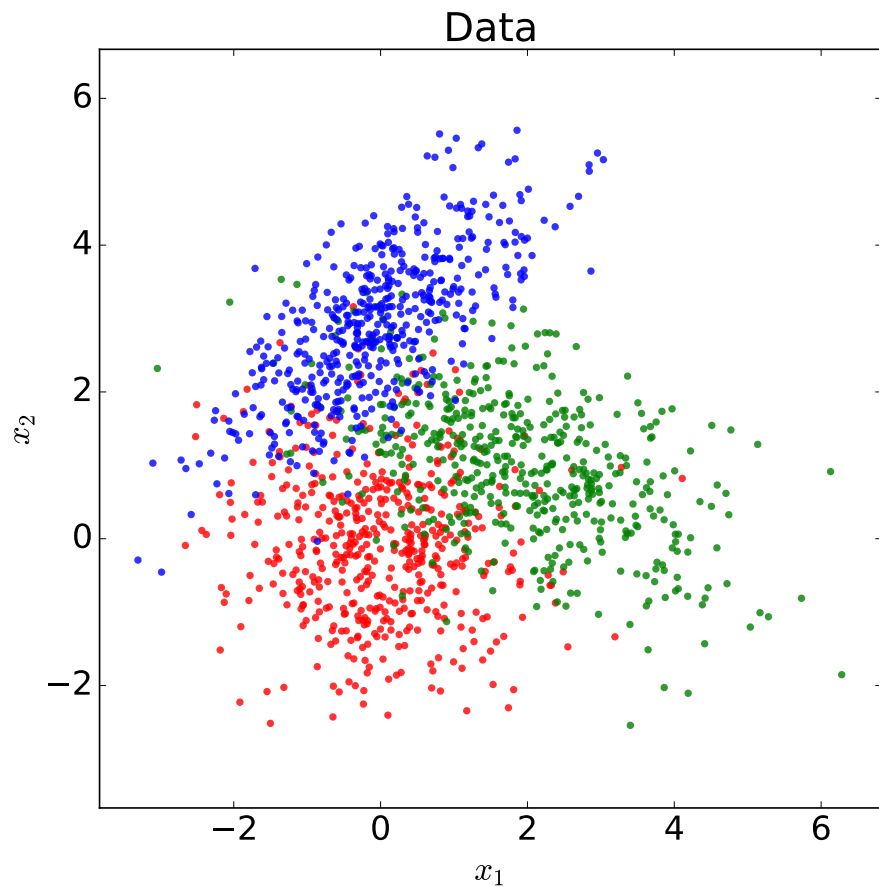
$$l'(\mathbf{W}) = \sum_{(\mathbf{x}, k) \in \mathcal{T}} \ln p(k|\mathbf{x}) = \sum_{(\mathbf{x}, k) \in \mathcal{T}} \mathbf{w}_k \cdot \mathbf{x} - \sum_{(\mathbf{x}, k) \in \mathcal{T}} \ln(e^{\mathbf{w}_1 \cdot \mathbf{x}} + e^{\mathbf{w}_2 \cdot \mathbf{x}} + \dots + e^{\mathbf{w}_K \cdot \mathbf{x}}) \quad (55)$$

Optimal parameters \mathbf{W}^* :

$$\mathbf{W}^* = \operatorname{argmax}_{\mathbf{W}} l'(\mathbf{W}) = \operatorname{argmin}_{\mathbf{W}} E(\mathbf{W}) \quad (56)$$

($E(\mathbf{W})$ is cross entropy, as before)

Logistic Regression with Multiple Classes, Example (1)



Each class: 500 samples from multivariate normal distribution.

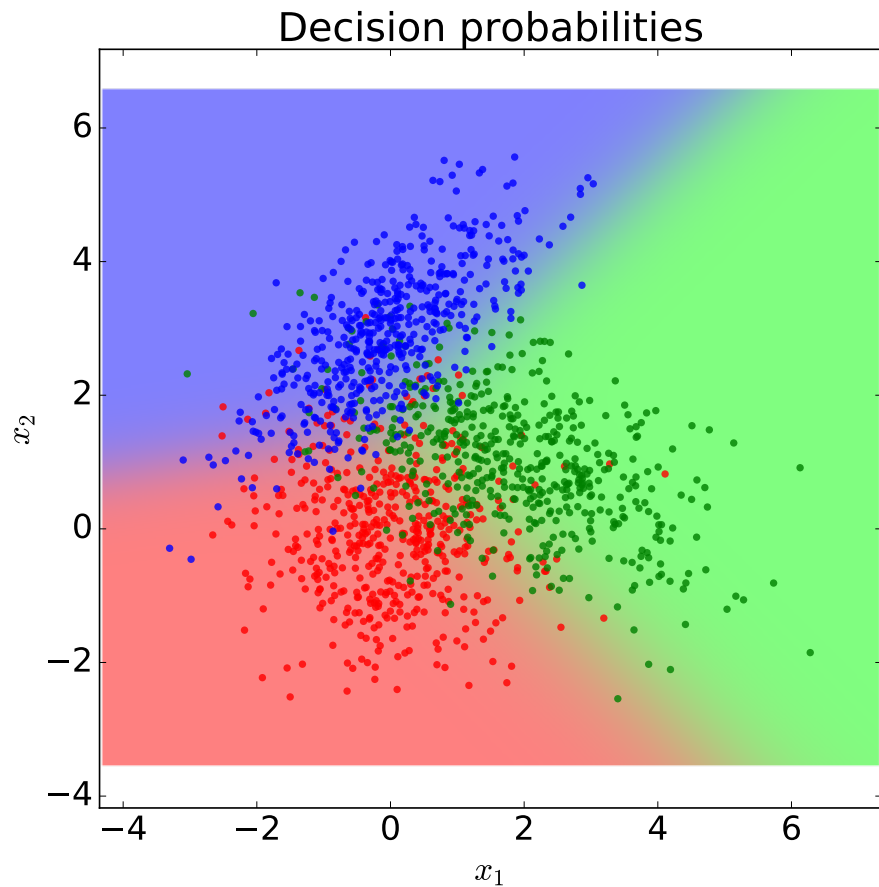
Centers:

$(0, 0)$, $(2, 1)$, $(0, 3)$.

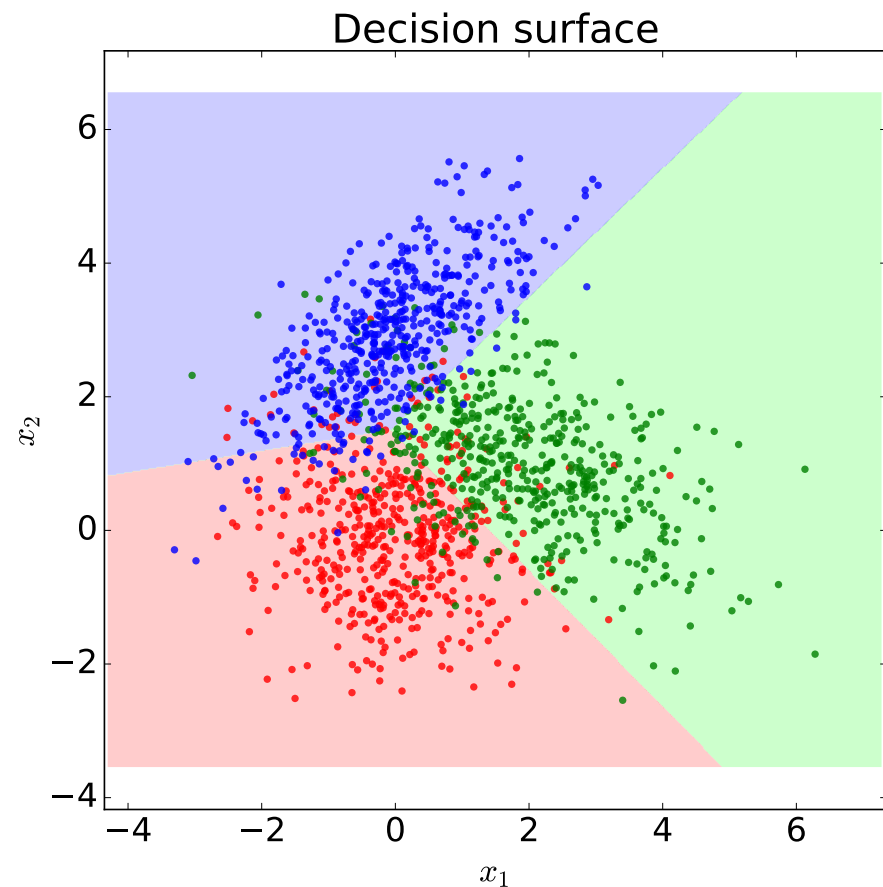
Cov. matrices:

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \begin{bmatrix} 2 & -0.7 \\ -0.7 & 1 \end{bmatrix}, \begin{bmatrix} 1 & 0.7 \\ 0.7 & 1 \end{bmatrix}.$$

Logistic Regression with Multiple Classes, Example (2)



Conditional probabilities $p(k|\mathbf{x})$ (coded by color intensity)



Decisions, coded by class' colors. Note the linearity of decision boundaries.