

# Logistic Regression

# Logistic Regression

- A binary classification model
- Developed in the field of Statistics, not Machine Learning
- Easy to implement
- Very widely used
- Easily extended to multiclass classifications using the OvR technique

Ref.: Wikipedia

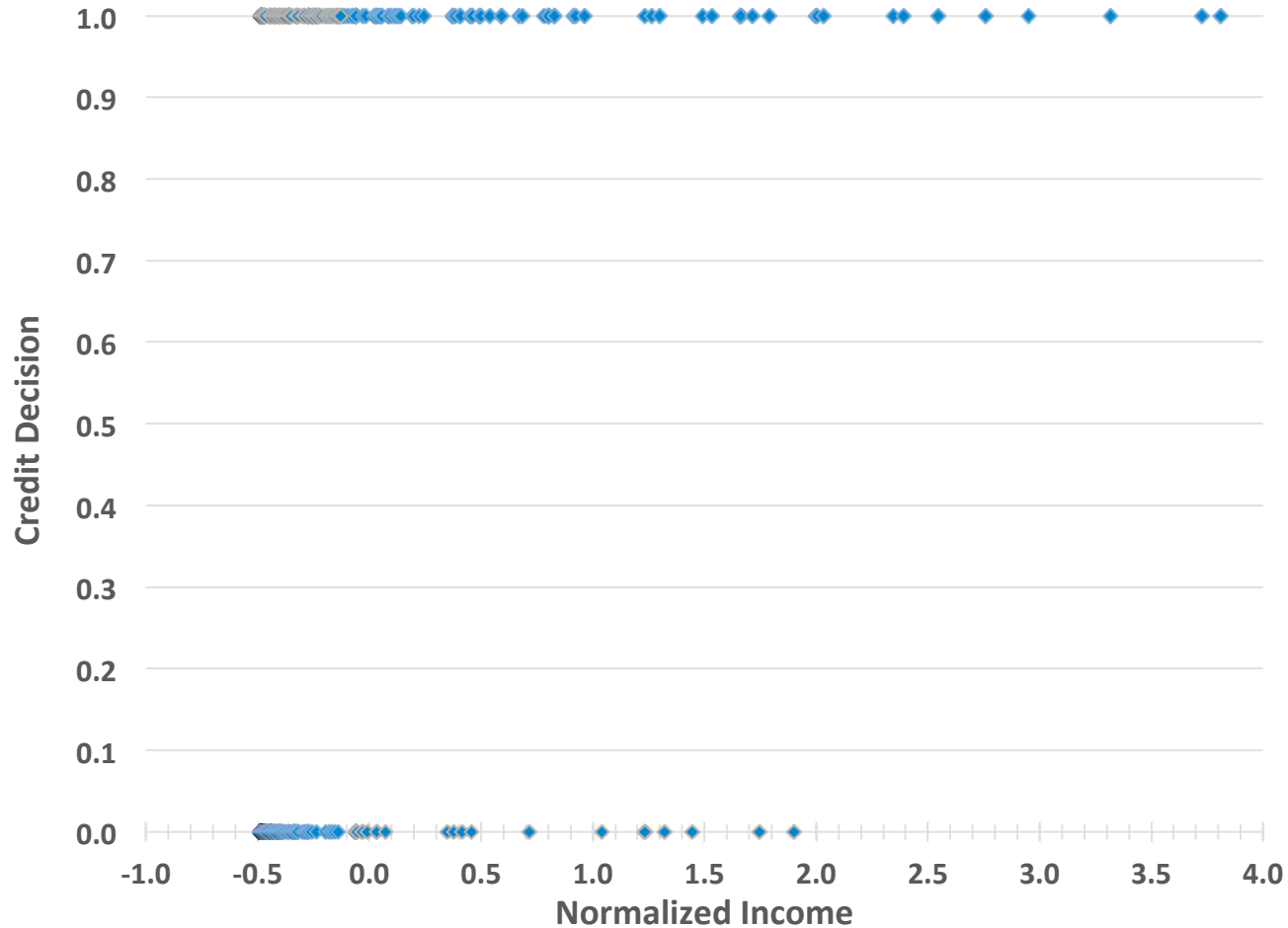
# Learning techniques we have seen so far

- Perceptron
  - Divide classes with a hyperplane
  - Converges only if classes are linearly separable
  - Provides no information about the confidence of the classification of a sample
- Adaline
  - Finds the hyperplane that does the best job of dividing the classes
  - Can converge if classes are not linearly separable
  - Provides no information about the confidence of the classification of a sample
- A Better Technique TBD
  - Finds a mechanism that does the best job of dividing the classes
  - Works if classes are not linearly separable
  - Provides information about the confidence of the classification of a sample

# Credit decision example

- Based on Real Data (not from the United States)
- Available data
  - Normalized income,  $x$
  - Credit decision  $y$ , (0 = denied, 1 = approved)
- Let  $P(y=1 | x)$  be the function whose value is the conditional probability that credit is approved given a normalized income of  $x$ .
- Required characteristics of  $P(y=1 | x)$ 
  - $0 \leq P(y=1 | x) \leq 1$
  - $P(y=0 | x) + P(y=1 | x) = 1$
  - $P(y=1 | x) \leq P(y=1 | x + \varepsilon)$  for all  $x$  and  $\varepsilon > 0$ .
- Desired characteristic of  $P(y=1 | x)$ 
  - Defined by as few parameters as possible

# Credit decision



# Odds ratio

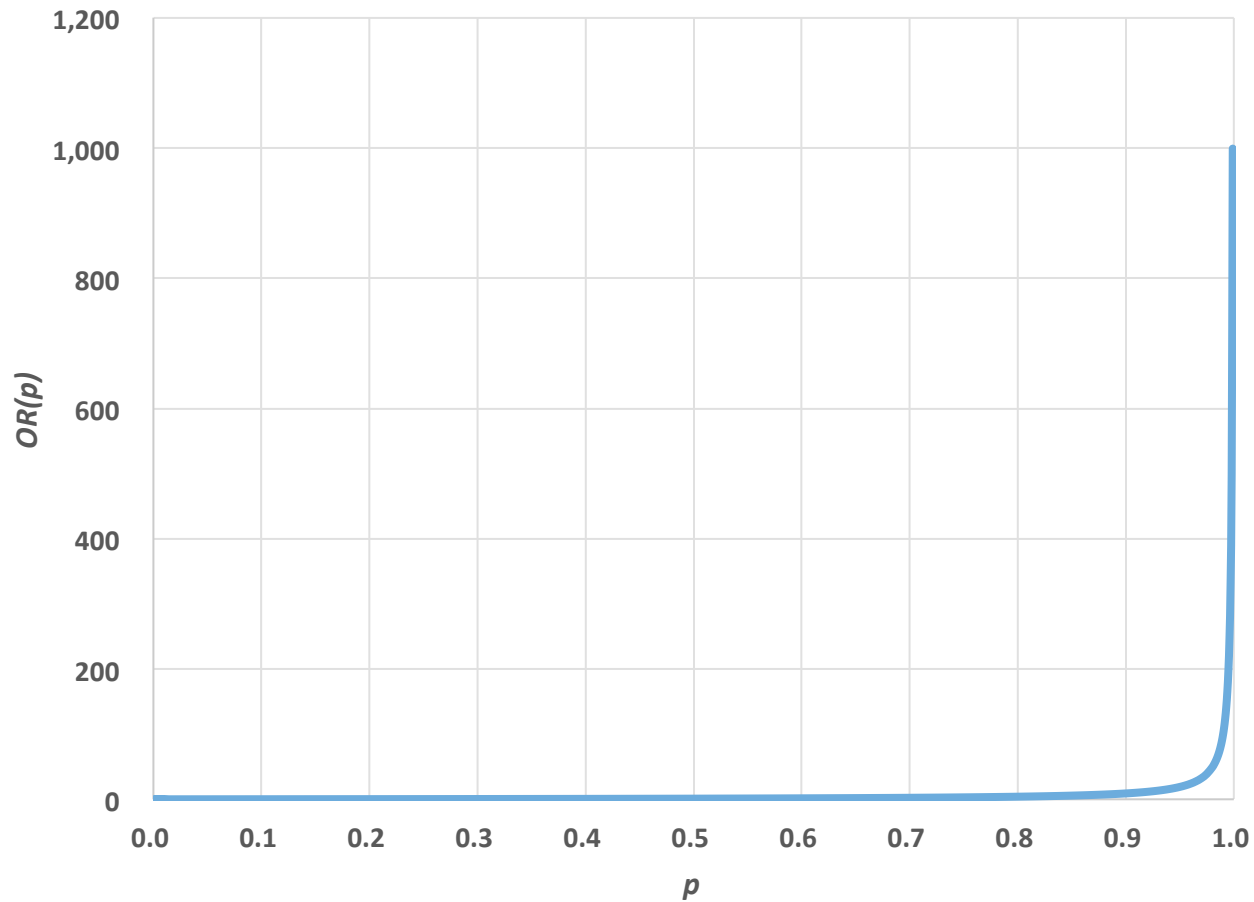
- Let  $p$  be the probability of a “positive” event
- We define the **odds ratio** ( $OR$ ) as the probability of a “positive” events divided by the odds of a “negative” event,

$$OR(p) = \frac{p}{(1 - p)}$$

Note that the  $OR \in [0, \infty)$

Ref.: Wikipedia

# Odds Ratio



# logit function

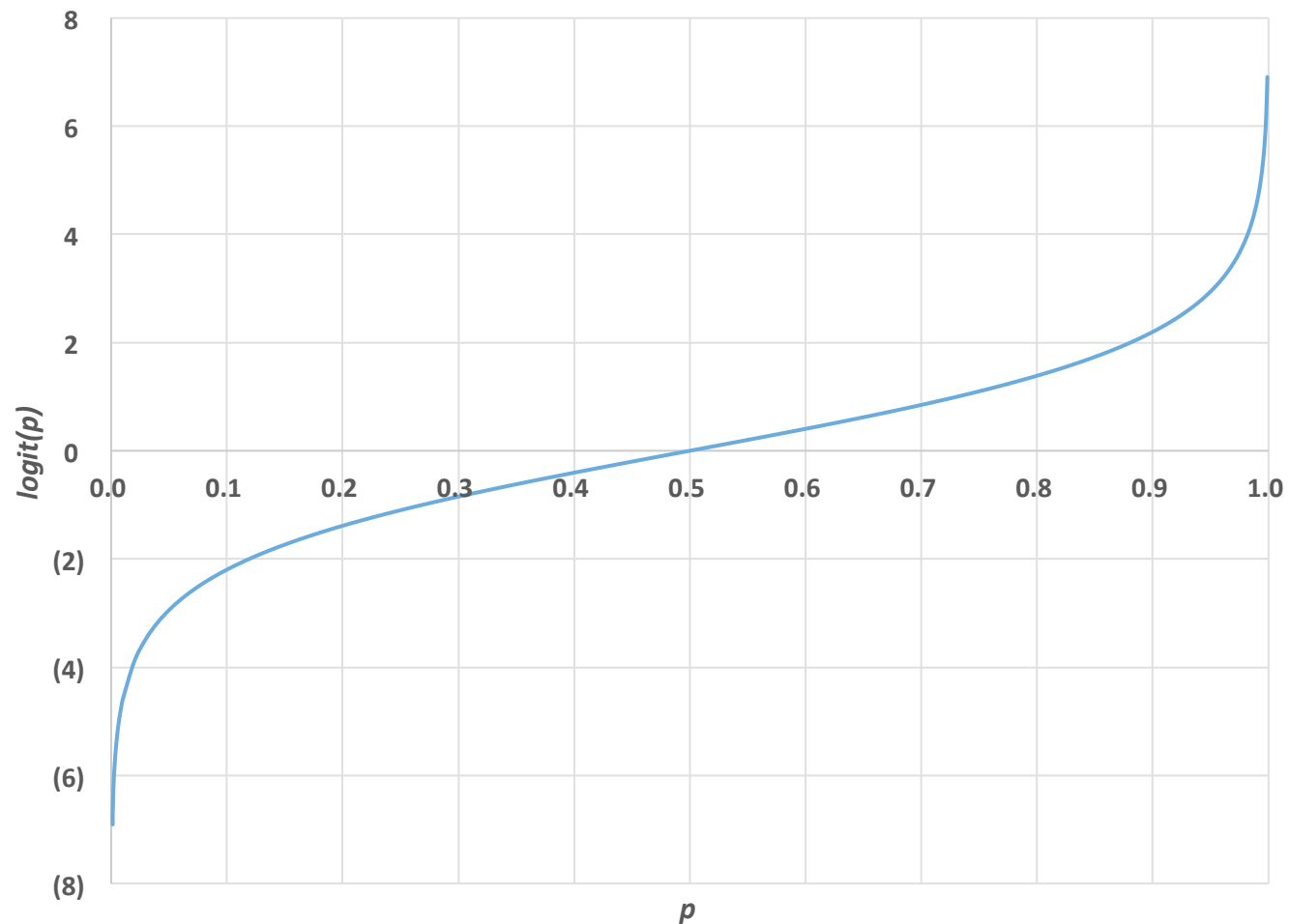
- Let  $p$  be the probability of a “positive” event
- We define the **logit function** as the log of the odds ratio,

$$\text{logit}(p) = \ln(OR(p)) = \ln\left(\frac{p}{(1-p)}\right)$$

Note that  $\text{logit}(p) \in (-\infty, +\infty)$



# logit function



# Solve for the inverse of the logit function

$$y = \text{logit}(p) = \ln\left(\frac{p}{1-p}\right)$$

$$e^y = e^{\ln\left(\frac{p}{1-p}\right)} = \left(\frac{p}{1-p}\right)$$

$$(1-p)e^y = p$$

$$e^y - pe^y = p$$

$$\begin{aligned} e^y &= p + pe^y \\ &= p(1 + e^y) \end{aligned}$$

# Solve for the inverse of the logit function

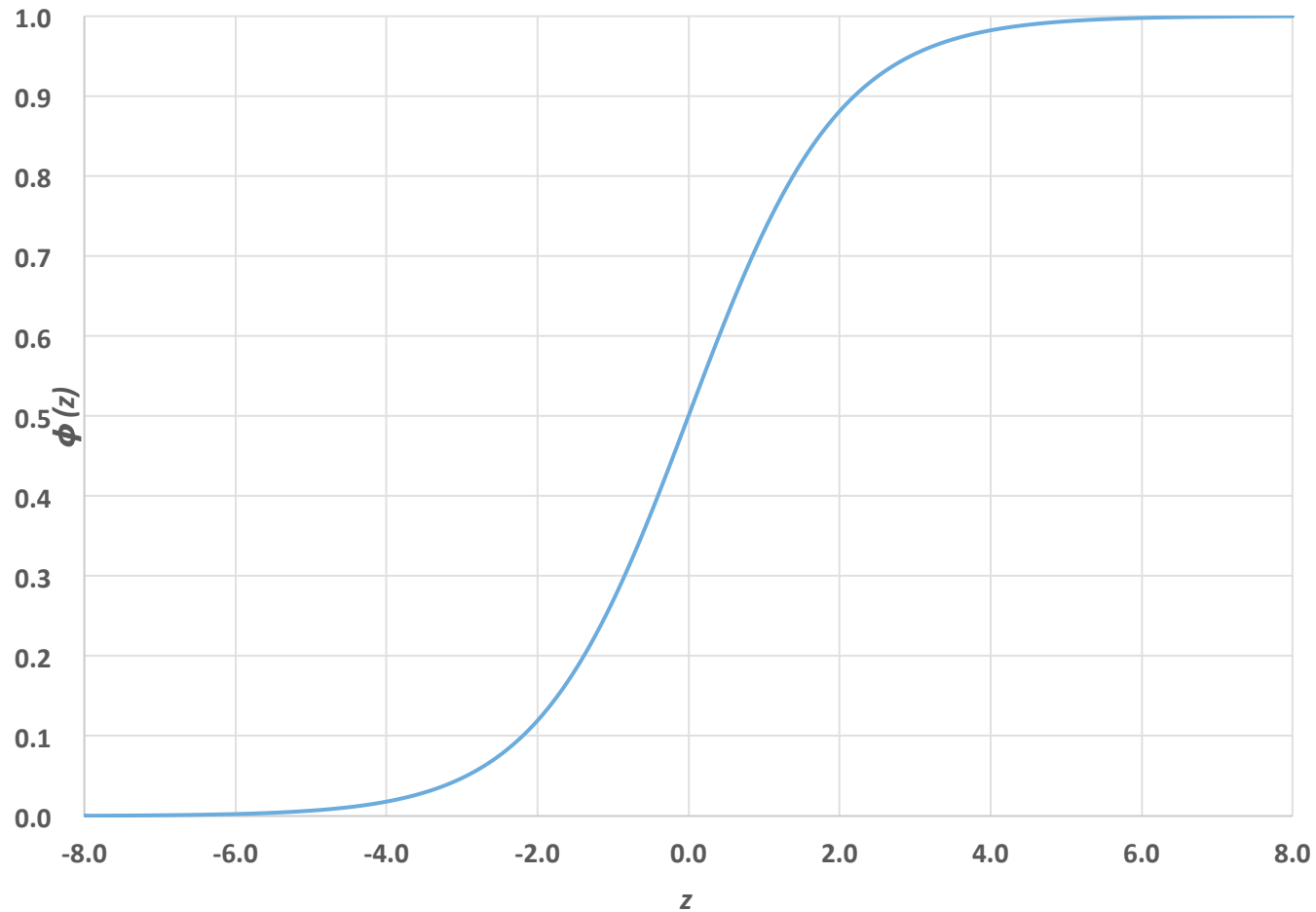
$$e^y = p(1 + e^y)$$

$$\begin{aligned} p &= \left( \frac{e^y}{1 + e^y} \right) \\ &= \left( \frac{e^y}{1 + e^y} \right) \left( \frac{e^{-y}}{e^{-y}} \right) \\ &= \left( \frac{1}{1 + e^{-y}} \right) \end{aligned}$$

**This is called the logistic function and sometimes also the sigmoid function**

$$\phi(z) = \left( \frac{1}{1 + e^{-z}} \right)$$

# logistic function



# Limits of the logistic function

$$\phi(z) = \left( \frac{1}{1 + e^{-z}} \right)$$

$$\begin{aligned} \lim_{z \rightarrow -\infty} \phi(z) &= \left( \frac{1}{1 + e^{-\infty}} \right) \\ &= \left( \frac{1}{1 + e^{\infty}} \right) \\ &= \left( \frac{1}{1 + \infty} \right) \\ &= \frac{1}{\infty} \\ &= 0 \end{aligned}$$

# Limits of the logistic function

$$\phi(z) = \left( \frac{1}{1 + e^{-z}} \right)$$

$$\begin{aligned} \lim_{z \rightarrow +\infty} \phi(z) &= \left( \frac{1}{1 + e^{-\infty}} \right) \\ &= \left( \frac{1}{1 + 0} \right) \\ &= \left( \frac{1}{1} \right) \\ &= 1 \end{aligned}$$

# Attributes of the logistic function

$$\phi(z) = \left( \frac{1}{1 + e^{-z}} \right)$$

$$\begin{aligned}\phi(0) &= \left( \frac{1}{1 + e^{-0}} \right) \\ &= \left( \frac{1}{1 + 1} \right) \\ &= \frac{1}{2}\end{aligned}$$



# Symmetry of the logistic function

$$\begin{aligned}\phi(-z) + \phi(z) &= \left( \frac{1}{1 + e^{-z}} \right) + \left( \frac{1}{1 + e^z} \right) \\&= \left( \frac{1}{1 + e^z} \right) + \left( \frac{1}{1 + e^{-z}} \right) \\&= \frac{(1 + e^{-z}) + (1 + e^z)}{(1 + e^z)(1 + e^{-z})} \\&= \frac{(2 + e^{-z} + e^z)}{(1 + e^z + e^{-z} + e^z e^{-z})} \\&= \frac{(2 + e^{-z} + e^z)}{(1 + e^z + e^{-z} + e^0)}\end{aligned}$$

# Symmetry of the logistic function

$$\begin{aligned}\phi(-z) + \phi(z) &= \frac{(2 + e^{-z} + e^z)}{(1 + e^z + e^{-z} + e^0)} \\ &= \frac{(2 + e^{-z} + e^z)}{(1 + e^z + e^{-z} + 1)} \\ &= \frac{(2 + e^{-z} + e^z)}{(2 + e^z + e^{-z})} \\ &= 1\end{aligned}$$

# Symmetry of the logistic function

$$\phi(-z) + \phi(z) = 1$$

$$\phi(+z) = 1 - \phi(-z)$$

$$\phi(-z) = 1 - \phi(+z)$$

$$\left( \phi(-z) - \frac{1}{2} \right) = \left( \frac{1}{2} - \phi(+z) \right)$$

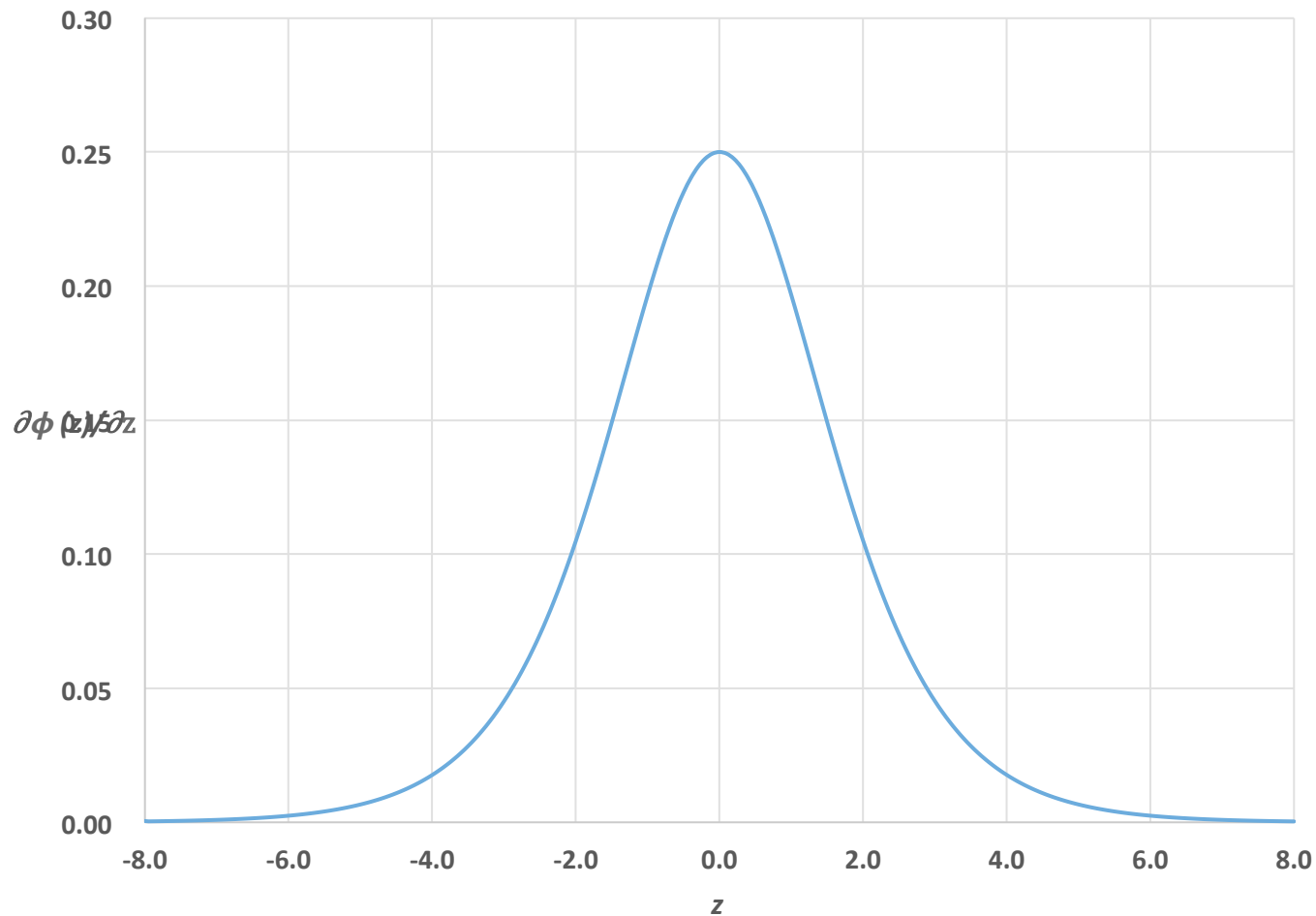
# Derivative of the logistic function

$$\begin{aligned}\frac{\partial}{\partial z} \phi(z) &= \frac{\partial}{\partial z} \left( \frac{1}{1 + e^{-z}} \right) \\&= \frac{1}{(1 + e^{-z})^2} \frac{\partial}{\partial z} (1 + e^{-z}) \\&= \frac{1}{(1 + e^{-z})^2} e^{-z} \\&= \frac{1}{(1 + e^{-z})^2} \frac{(1 + e^{-z})}{(1 + e^{-z})} e^{-z} \frac{e^z}{e^z} \\&= \frac{1}{(1 + e^{-z})} \frac{1}{(1 + e^{-z})} \frac{1}{e^z}\end{aligned}$$

# Derivative of the logistic function

$$\begin{aligned}\frac{\partial}{\partial z} \phi(z) &= \frac{1}{(1 + e^{-z})} \frac{1}{(1 + e^{-z})} \frac{1}{e^z} \\&= \frac{1}{(1 + e^{-z})} \frac{1}{(e^z + e^{-z} e^z)} \\&= \frac{1}{(1 + e^{-z})} \frac{1}{(e^z + 1)} \\&= \frac{1}{(1 + e^{-z})} \frac{1}{(1 + e^z)} \\&= \phi(z) \phi(-z) \\&= \phi(z) (1 - \phi(z))\end{aligned}$$

# Derivative of the logistic function



# Derivative of the logistic function

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

$$\frac{\partial}{\partial z} \phi(0) = \phi(0)(1 - \phi(0))$$

$$= \left(\frac{1}{2}\right)\left(\frac{1}{2}\right)$$

$$= \frac{1}{4}$$

# Limits of the derivative of the logistic function

$$\frac{\partial}{\partial z} \phi(z) = \phi(z)\phi(-z)$$

$$\begin{aligned}\lim_{z \rightarrow +\infty} \frac{\partial}{\partial z} \phi(z) &= \left( \lim_{z \rightarrow +\infty} \frac{\partial}{\partial z} \phi(z) \right) \left( \lim_{z \rightarrow +\infty} \frac{\partial}{\partial z} \phi(-z) \right) \\ &= \left( \lim_{z \rightarrow +\infty} \frac{\partial}{\partial z} \phi(z) \right) \left( \lim_{z \rightarrow -\infty} \frac{\partial}{\partial z} \phi(z) \right) \\ &= (1)(0) \\ &= 0\end{aligned}$$

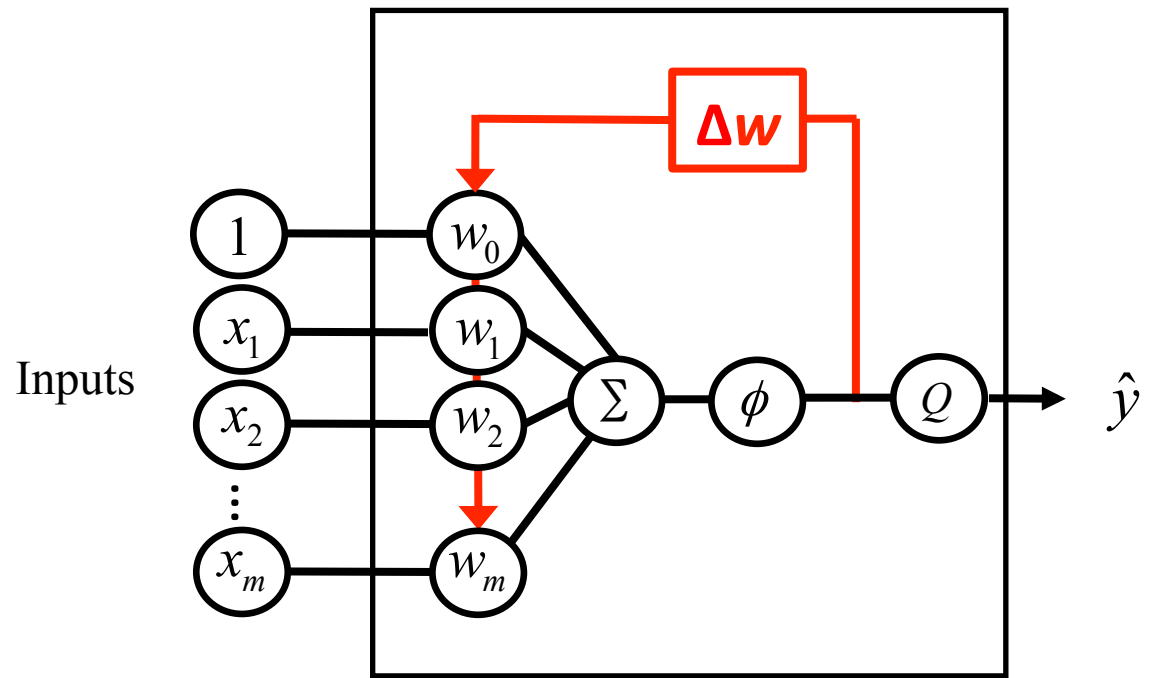
$$\begin{aligned}\lim_{z \rightarrow -\infty} \frac{\partial}{\partial z} \phi(z) &= \left( \lim_{z \rightarrow -\infty} \frac{\partial}{\partial z} \phi(z) \right) \left( \lim_{z \rightarrow -\infty} \frac{\partial}{\partial z} \phi(-z) \right) \\ &= \left( \lim_{z \rightarrow -\infty} \frac{\partial}{\partial z} \phi(z) \right) \left( \lim_{z \rightarrow +\infty} \frac{\partial}{\partial z} \phi(z) \right) \\ &= (0)(1) \\ &= 0\end{aligned}$$



# Back to the credit decision example

- Required characteristics of  $P(y=1 | \mathbf{x})$ 
  - $0 \leq P(y=1 | \mathbf{x}) \leq 1$
  - $P(y=0 | \mathbf{x}) + P(y=1 | \mathbf{x}) = 1$
  - $P(y=1 | \mathbf{x}) \leq P(y=1 | \mathbf{x} + \varepsilon)$  for all  $\mathbf{x}$  and  $\varepsilon > 0$ .
- Desired characteristic of  $P(y=1 | \mathbf{x})$ 
  - Defined by as few parameters as possible
- Let  $\text{logit}(P(y = 1 | \mathbf{x})) = w_0x_0 + w_1x_1 + \dots + w_mx_m$ 
$$= \sum_{i=0}^m w_i x_i = \mathbf{w}^T \mathbf{x}$$
$$= z$$
- Then  $P(y = 1 | \mathbf{x}) = \phi(z)$

# Training an Adaline / Logistic Regression



- $\mathbf{x}$  input
- $\mathbf{w}$  weight
- $\Sigma$   $z = \mathbf{w}^T \mathbf{x}$
- $\phi$  activation function
- $Q$  quantizer
- $\hat{y}$  computed output value

**$\Delta w$**  = adjustments to  $w$

# Activation functions

Adaline  $\phi(z) = z$

Logistic Regression  $\phi(z) = \frac{1}{1 + e^{-z}}$

# Cost functions

Adaline  $J(\mathbf{w}) = \frac{1}{2} \sum_i \left( y^{(i)} - \phi(z^{(i)}) \right)^2$

Logistic Regression  $J(\mathbf{w}) = ?$

# Likelihood function

$$L(\mathbf{w}) = P(\mathbf{y} \mid \mathbf{x} : \mathbf{w})$$

$$= \prod_{i=1}^n P(\mathbf{y}^{(i)} \mid \mathbf{x}^{(i)}; \mathbf{w})$$

$$= \prod_{i=1}^n \left( \phi(z^{(i)}) \right)^{y^{(i)}} \left( 1 - \phi(z^{(i)}) \right)^{1-y^{(i)}}$$

# Log-likelihood and cost functions

$$\begin{aligned} l(\mathbf{w}) &= \ln(L(\mathbf{w})) \\ &= \sum_{i=1}^n \left[ y^{(i)} \ln(\phi(z^{(i)})) + (1 - y^{(i)}) \ln(1 - \phi(z^{(i)})) \right] \end{aligned}$$

$$\begin{aligned} J(\mathbf{w}) &= - l(\mathbf{w}) \\ &= - \sum_{i=1}^n \left[ y^{(i)} \ln(\phi(z^{(i)})) + (1 - y^{(i)}) \ln(1 - \phi(z^{(i)})) \right] \end{aligned}$$

# Cost function for a single-sample instance

$$\begin{aligned} J(\mathbf{w}) &= -l(\mathbf{w}) \\ &= -\sum_{i=1}^n \left[ y^{(i)} \ln(\phi(z^{(i)})) + (1 - y^{(i)}) \ln(1 - \phi(z^{(i)})) \right] \end{aligned}$$

$$J(\phi(z), y; \mathbf{w}) = -y \ln(\phi(z)) - (1 - y) \ln(1 - \phi(z))$$

# Cost function for a single-sample instance

$$\begin{aligned} J(\phi(z), y; \mathbf{w}) &= -y \ln(\phi(z)) - (1-y) \ln(1-\phi(z)) \\ &= \begin{cases} -\ln(\phi(z)) & \text{for } y=1 \\ -\ln(1-\phi(z)) = -\ln(\phi(-z)) & \text{for } y=0 \end{cases} \end{aligned}$$



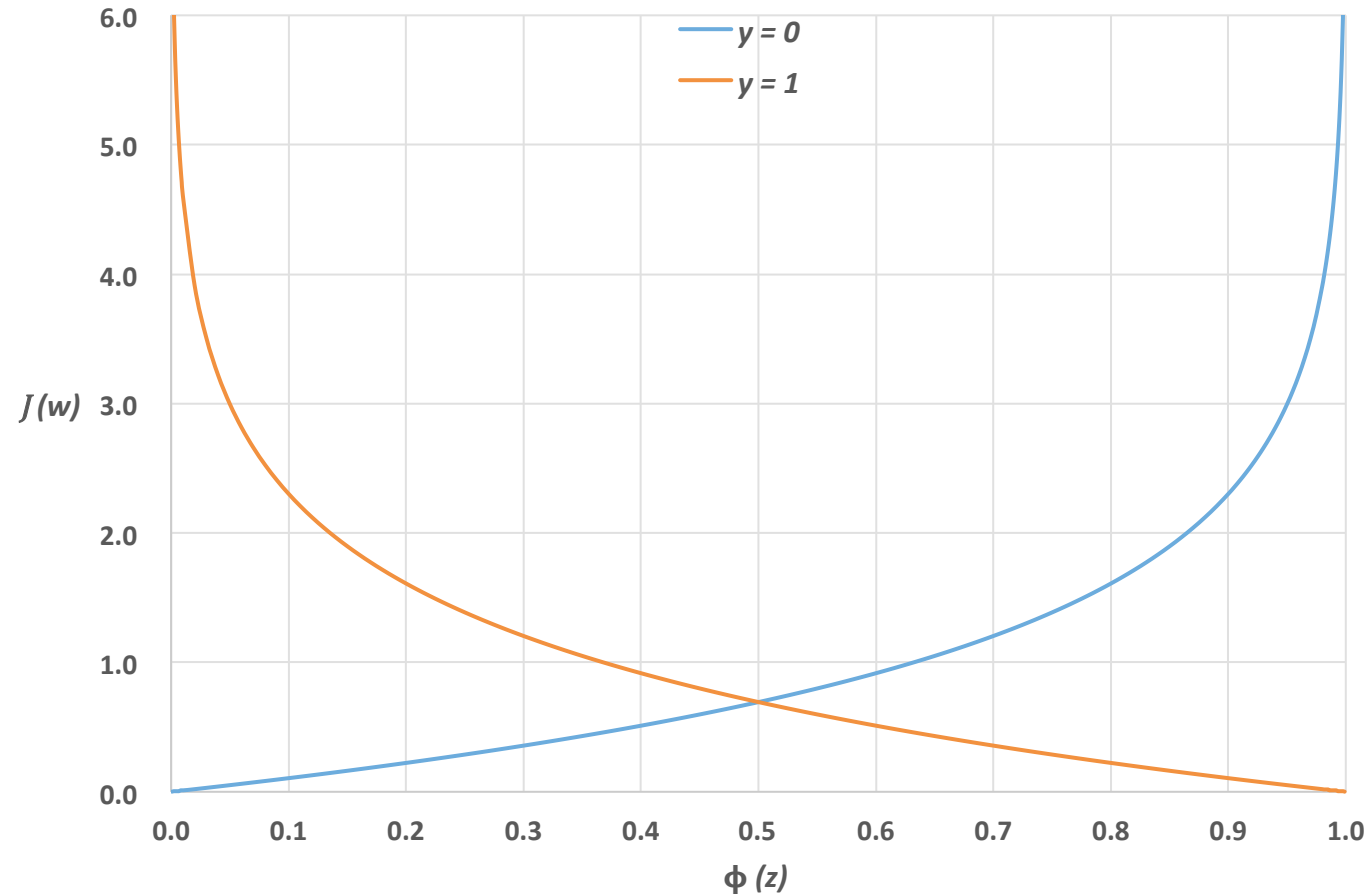
# Derivative of the cost function for a single-sample instance

$$\begin{aligned} J(\phi(z), y; \mathbf{w}) &= -y \ln(\phi(z)) - (1-y) \ln(1-\phi(z)) \\ \frac{\partial J(\phi(z), y; \mathbf{w})}{\partial \phi(z)} &= -y \frac{\partial \ln(\phi(z))}{\partial \phi(z)} - (1-y) \frac{\partial \ln(1-\phi(z))}{\partial \phi(z)} \\ &= -\frac{y}{\phi(z)} - (1-y) \frac{\partial \ln(\phi(-z))}{\partial \phi(z)} \\ &= -\frac{y}{\phi(z)} + \frac{(1-y)}{\phi(-z)} \end{aligned}$$

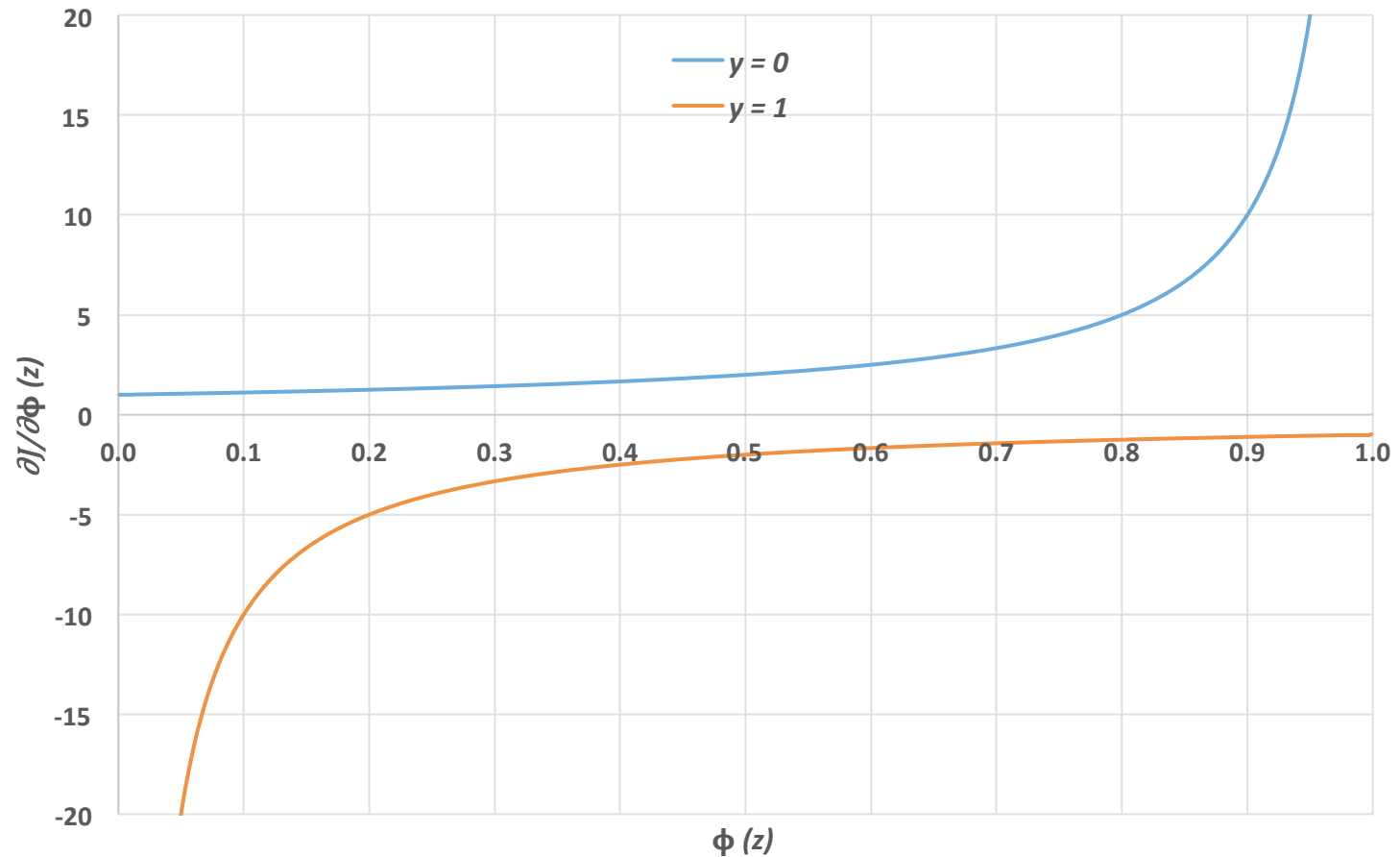
# Derivative of the cost function for a single-sample instance

$$\begin{aligned}\frac{\partial J(\phi(z), y; \mathbf{w})}{\partial \phi(z)} &= -\frac{y}{\phi(z)} + \frac{(1-y)}{\phi(-z)} \\ &= \begin{cases} -\frac{1}{\phi(z)} & \text{for } y=1 \\ +\frac{1}{\phi(-z)} & \text{for } y=0 \end{cases}\end{aligned}$$

# Cost function for a single-sample instance



# Derivative of cost function for a single-sample instance



# Gradient in logistic regression

$$J(\mathbf{w}) = -\left(y \ln(\phi(z)) + (1-y) \ln(1-\phi(z))\right)$$

$$\begin{aligned} \frac{\partial}{\partial w_j} J(\mathbf{w}) &= -\left(\left(y \frac{1}{\phi(z)}\right) \frac{\partial}{\partial w_j} (\phi(z)) - \left((1-y) \frac{1}{(1-\phi(z))}\right) \frac{\partial}{\partial w_j} (\phi(z))\right) \\ &= -\left(\left(y \frac{1}{\phi(z)}\right) - \left((1-y) \frac{1}{(1-\phi(z))}\right)\right) \frac{\partial}{\partial w_j} (\phi(z)) \end{aligned}$$

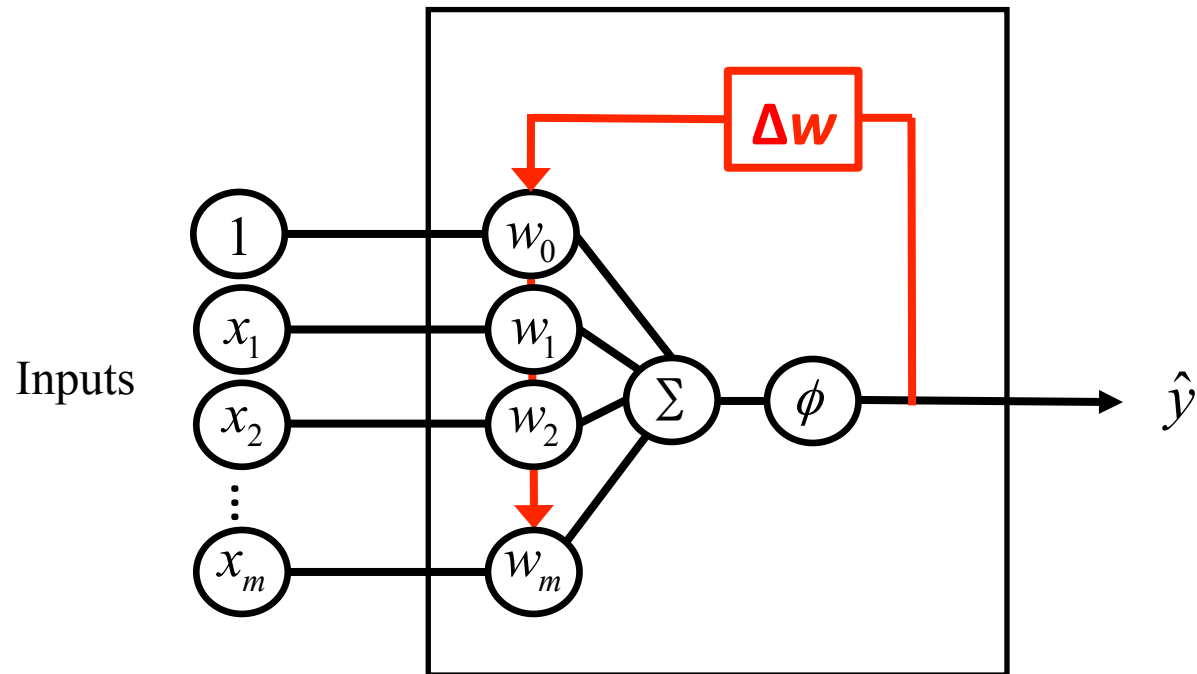
$$\text{since } \frac{\partial}{\partial x} \ln(f(x)) = \frac{1}{f(x)} \frac{\partial}{\partial x} f(x)$$

# Gradient in logistic regression

Using our previous result that  $\frac{\partial}{\partial z} \phi(z) = \phi(z)(1 - \phi(z))$ ,

$$\begin{aligned}\frac{\partial}{\partial w_j} J(\mathbf{w}) &= - \left( y \frac{1}{\phi(z)} + (1-y) \frac{1}{(1-\phi(z))} \right) \frac{\partial}{\partial w_j} (\phi(z)) \\&= - \left( \left( y \frac{1}{\phi(z)} \right) + \left( (1-y) \frac{-1}{(1-\phi(z))} \right) \right) \frac{\partial}{\partial w_j} (\phi(z)) \\&= - \left( \left( y \frac{1}{\phi(z)} \right) + \left( (y-1) \frac{1}{(1-\phi(z))} \right) \right) \phi(z)(1-\phi(z)) \frac{\partial}{\partial w_j} z \\&= - \left( (y(1-\phi(z))) + ((y-1)\phi(z)) \right) \frac{\partial}{\partial w_j} z \\&= - (y - y\phi(z) + y\phi(z) - \phi(z)) x_j \\&= - (y - \phi(z)) x_j\end{aligned}$$

# Training a logistic regression model (or a Adaline model or a Perceptron)



- $\mathbf{x}$  input
- $\mathbf{w}$  weight
- $\Sigma$   $z = \mathbf{w}^T \mathbf{x}$
- $\phi$  activation function
- $\hat{y}$  computed output value

**$\Delta w$**  = adjustments to  $w$

# Training an Logistic regression model (or an Adaline model)

To train a Logistic Regression model we update the weights to minimize the cost function,  $J(\mathbf{w})$ , by moving in direction of the negative of the gradient of  $J(\mathbf{w})$  by an amount proportional to the magnitude of the gradient.

$$\mathbf{w} := \mathbf{w} + \Delta\mathbf{w} \text{ where}$$

$$\Delta\mathbf{w} = -\eta \nabla J(\mathbf{w})$$

and where  $\eta$  is the learning rate such that  $0 < \eta < 1$ .



# Training a logistic regression model

$\mathbf{w} := \mathbf{w} + \Delta \mathbf{w}$  where

$$\Delta \mathbf{w} = -\eta \nabla J(\mathbf{w})$$

$$\begin{aligned} \Delta w_j &= -\eta \frac{\partial J(\mathbf{w})}{\partial w_j} \\ &= \eta \sum_i \left( y^{(i)} - \phi(z^{(i)}) \right) x_j^{(i)} \end{aligned}$$

and where  $z^{(i)} = \mathbf{w}^T \mathbf{x}^{(i)}$

# Logistic regression training algorithm

1. Initialize the weights,  $\mathbf{w}$ , to 0 or small random numbers
2. Compute the gradient of the cost function,  $\nabla J(\mathbf{w})$ , by summing over all (or a fixed number of randomly selected) training samples,  $(\mathbf{x}^{(i)}, y^{(i)})$ ,

$$\frac{\partial J(\mathbf{w})}{\partial w_j} = - \sum_i \left( y^{(i)} - \phi(z^{(i)}) \right) x_j^{(i)}$$

# Logistic regression training algorithm

3. Update the weights,

$w_j := w_j + \Delta w_j$ , where

$$\begin{aligned}\Delta w_j &= -\eta \frac{\partial J(\mathbf{w})}{\partial w_j} \\ &= \eta \sum_i \left( y^{(i)} - \phi(z^{(i)}) \right) x_j^{(i)}\end{aligned}$$

and where  $\eta$  is the learning rate such that  $0 < \eta < 1$ .

# Logistic regression training algorithm

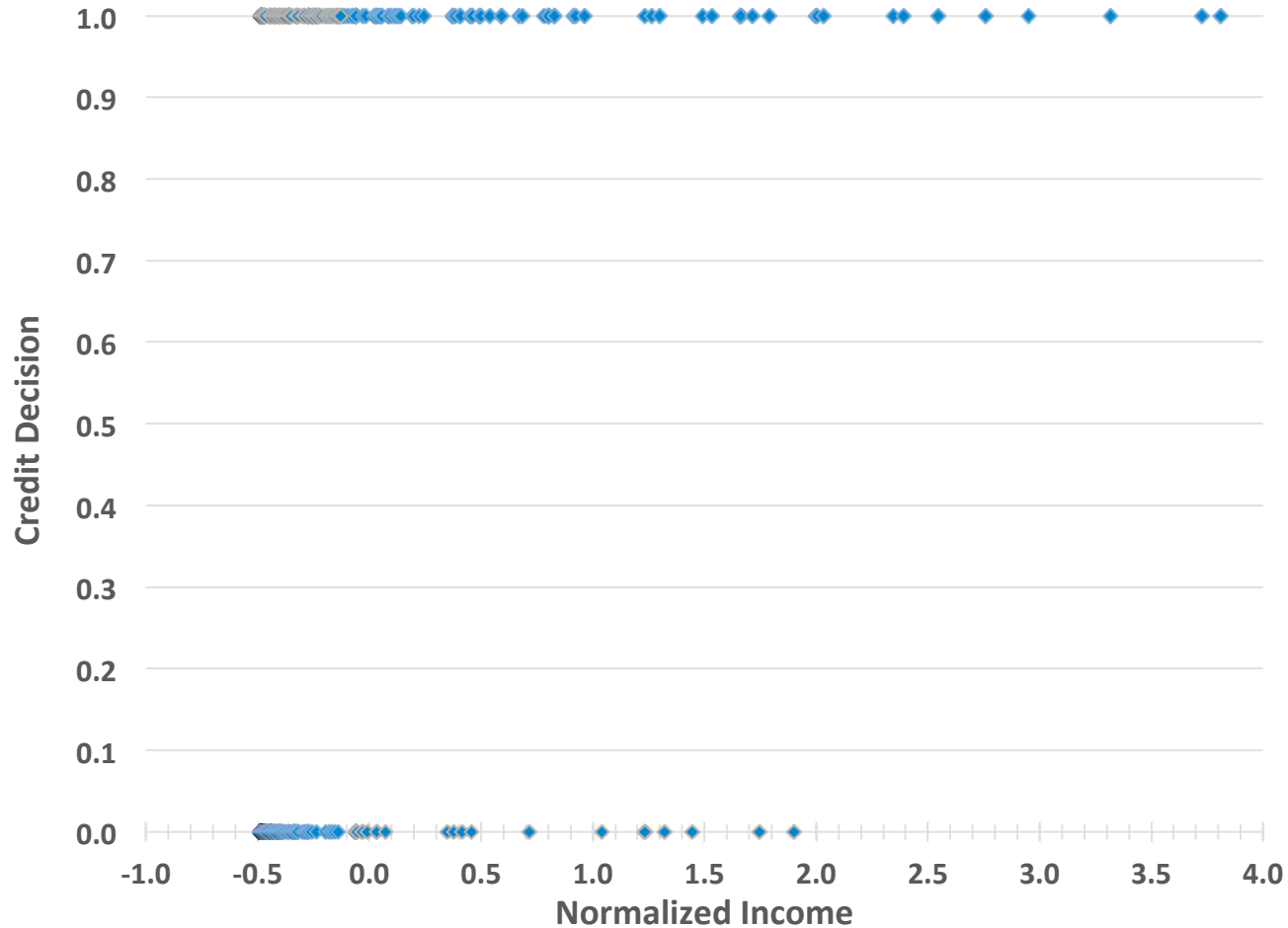
4. Repeat steps 2. and 3. until the weights converge,  
that is, until

$$\|\Delta \mathbf{w}\| < \varepsilon \text{ (or } \|\Delta \mathbf{w}\|_1 < \varepsilon), \text{ where}$$

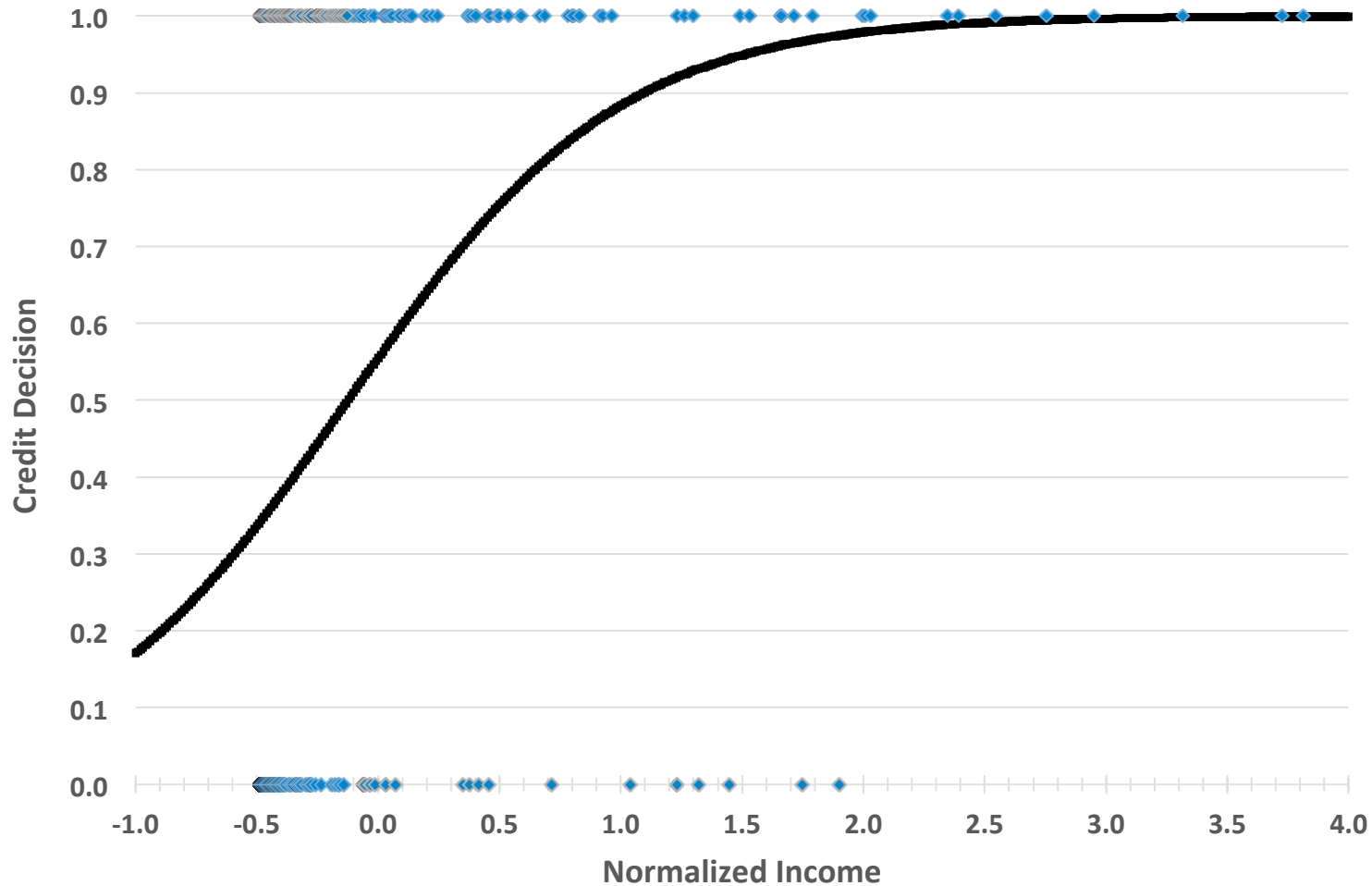
$\varepsilon$  is the convergence threshold,  $\varepsilon > 0$ .

or for a set number of iterations.

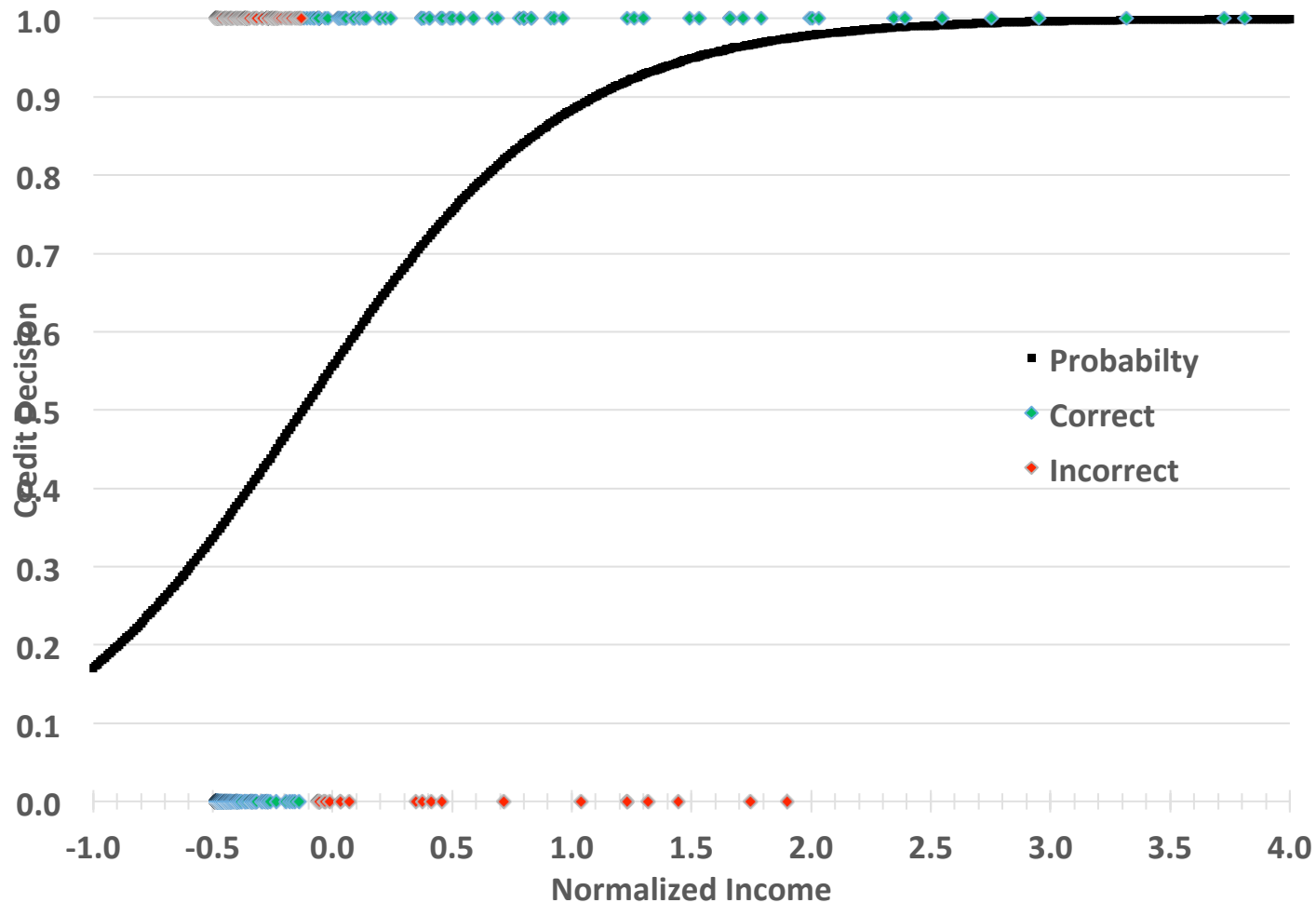
# Credit decision



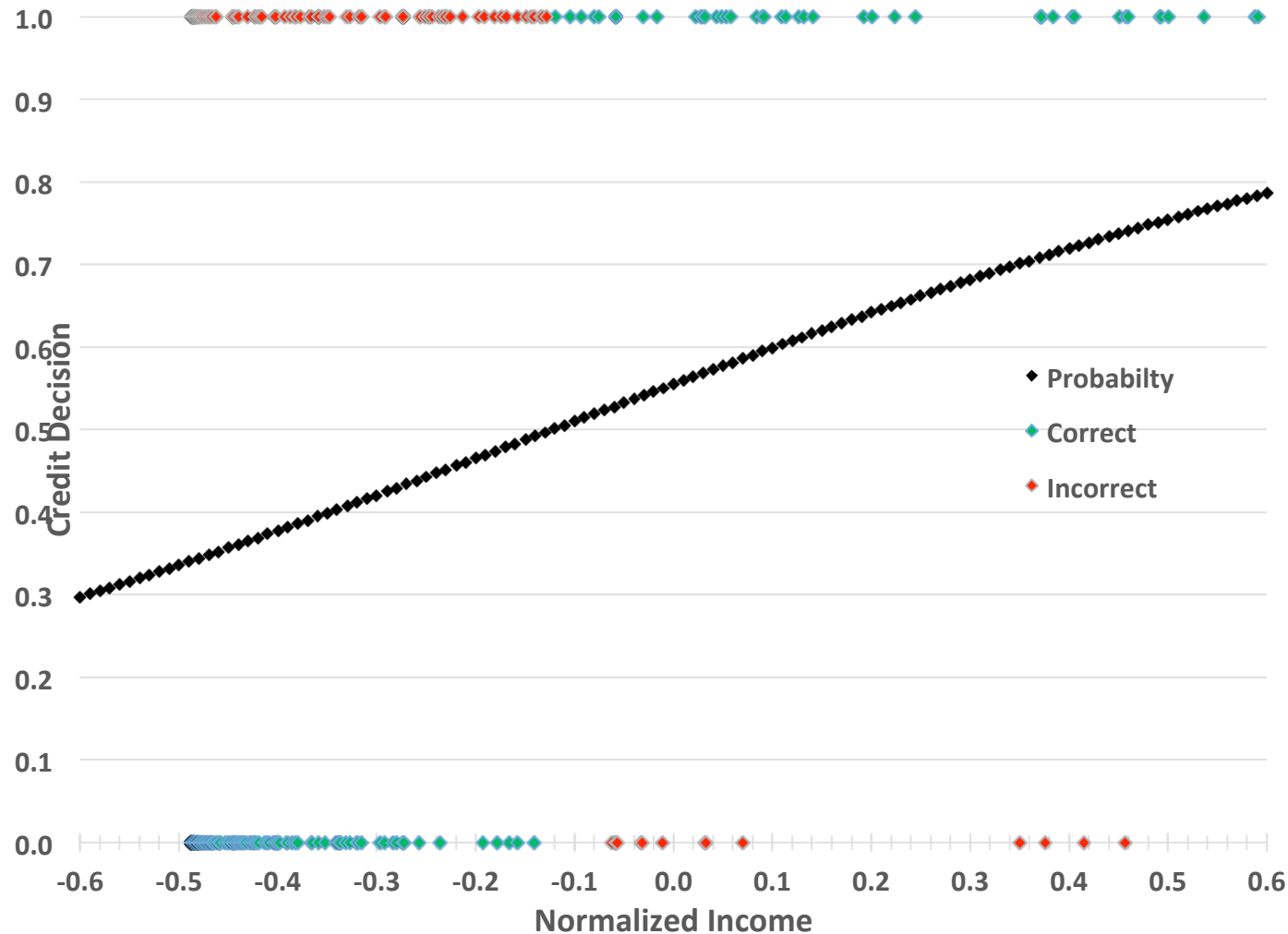
# Credit decision using logistic regression



# Credit decision using logistic regression



# Credit decision using logistic regression





# Regularization

**Regularization** is the process of introducing additional information in order to solve an ill-posed problem or to prevent **overfitting**.

Regularization can enhance the prediction accuracy and/or interpretability of the model it produces.

Regularization often encourages (or requires) models to be simpler rather than more complex.

# Logistic regression cost functions

## Without Regularization

$$J(\mathbf{w}) = - \sum_{i=1}^n \left[ y^{(i)} \ln \left( \phi \left( z^{(i)} \right) \right) + \left( 1 - y^{(i)} \right) \ln \left( 1 - \phi \left( z^{(i)} \right) \right) \right]$$

## With L2 Regularization

$$J_{\lambda}(\mathbf{w}) = - \sum_{i=1}^n \left[ y^{(i)} \ln \left( \phi \left( z^{(i)} \right) \right) + \left( 1 - y^{(i)} \right) \ln \left( 1 - \phi \left( z^{(i)} \right) \right) \right] + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

$$J_{L2}(\mathbf{w}) = -C \sum_{i=1}^n \left[ y^{(i)} \ln \left( \phi \left( z^{(i)} \right) \right) + \left( 1 - y^{(i)} \right) \ln \left( 1 - \phi \left( z^{(i)} \right) \right) \right] + \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

where  $C = \frac{1}{\lambda}$

## With L1 Regularization

$$J_{L1}(\mathbf{w}) = -C \sum_{i=1}^n \left[ y^{(i)} \ln \left( \phi \left( z^{(i)} \right) \right) + \left( 1 - y^{(i)} \right) \ln \left( 1 - \phi \left( z^{(i)} \right) \right) \right] + \|\mathbf{w}\|_1$$

# Gradient of cost functions

## Without Regularization

$$\frac{\partial}{\partial w_j} J(\mathbf{w}) = - (y - \phi(z)) x_j$$

## With L2 Regularization

$$\frac{\partial}{\partial w_j} J_{L2}(\mathbf{w}) = - C (y - \phi(z)) x_j + |w_j|$$

## With L1 Regularization

$$\frac{\partial}{\partial w_j} J_{L1}(\mathbf{w}) = - C (y - \phi(z)) x_j + 1$$

# Alternate cost functions

## Textbook

$$\begin{aligned} J(\mathbf{w}) &= - \sum_{i=1}^n \left[ y^{(i)} \ln \left( \phi \left( z^{(i)} \right) \right) + \left( 1 - y^{(i)} \right) \ln \left( 1 - \phi \left( z^{(i)} \right) \right) \right] \\ &= - \sum_{i=1}^n \left[ y^{(i)} \ln \left( \frac{1}{1 + \exp \left( -z^{(i)} \right)} \right) + \left( 1 - y^{(i)} \right) \ln \left( \frac{1}{1 + \exp \left( +z^{(i)} \right)} \right) \right] \end{aligned}$$

## Scikit - learn Documentation

$$J(\mathbf{w}) = \sum_{i=1}^n \left[ \ln \left( \exp \left( -y^{(i)} z^{(i)} \right) + 1 \right) \right]$$

**These alternatives are equivalent but not identical.**

# Scikit-learn Logistic Regression Class

**class sklearn.linear\_model.LogisticRegression(...)**

Parameter	Default		Parameter	Defaults
penalty	'l2'		random_state	None
dual	False		solver	'liblinear'
tol	0.0001		max_iter	100
C	1.0		multi_class	'ovr'
fit_intercept	True		verbose	0
intercept_scaling	1		warm_start	False
class_weight	None		n_jobs	1

# Scikit-learn Logistic Regression Attributes

Attribute	Description
coef_	Coefficient of the features in the decision function.
intercept_	Intercept (a.k.a. bias) added to the decision function
n_iter_	Actual number of iterations for all classes.

# Scikit-learn Logistic Regression Methods

Method	Description
<a href="#"><code>decision_function(X)</code></a>	Predict confidence scores for samples.
<a href="#"><code>densify()</code></a>	Convert coefficient matrix to dense array format.
<a href="#"><code>fit(X, y[, sample_weight])</code></a>	Fit the model according to the given training data.
<a href="#"><code>get_params([deep])</code></a>	Get parameters for this estimator.
<a href="#"><code>predict(X)</code></a>	Predict class labels for samples in X.
<a href="#"><code>predict_log_proba(X)</code></a>	Log of probability estimates.
<a href="#"><code>predict_proba(X)</code></a>	Probability estimates.
<a href="#"><code>score(X, y[, sample_weight])</code></a>	Returns the mean accuracy on the given test data and labels.
<a href="#"><code>set_params(**params)</code></a>	Set the parameters of this estimator.
<a href="#"><code>sparsify()</code></a>	Convert coefficient matrix to sparse format.