# 7.1 Linear Regression and Classification
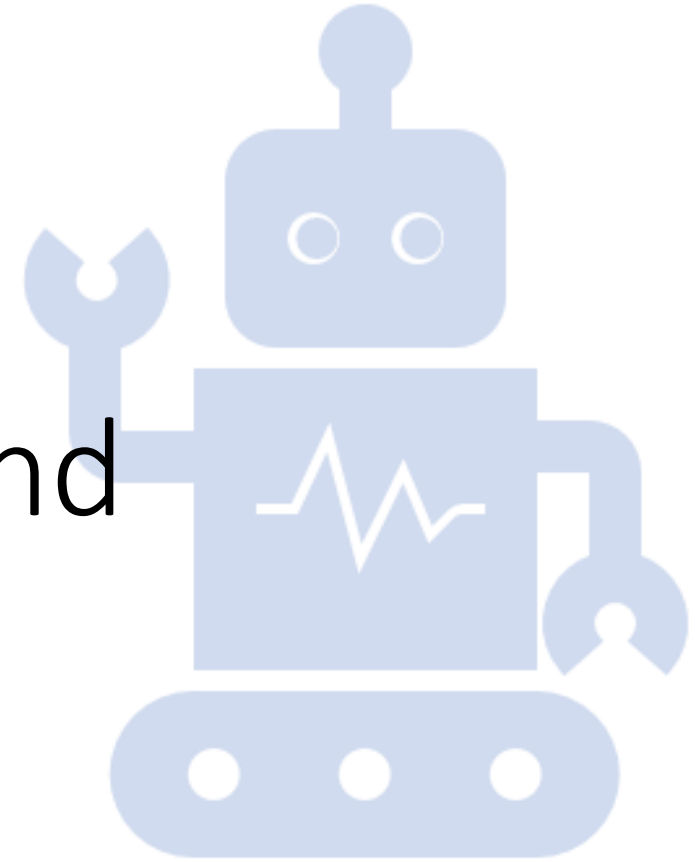
IN3050/IN4050 Introduction to Artificial Intelligence
and Machine Learning

# Today

1. Linear Regression and Classification
2. The Logistic Function and its Derivative
3. The Logistic Regression Classifier
4. Cross-Entropy Loss
5. Training the Logistic Regression Classifier
6. Variants of Gradient Descent
7. Multi-Class Classification

# Supervised learning - Where are we?
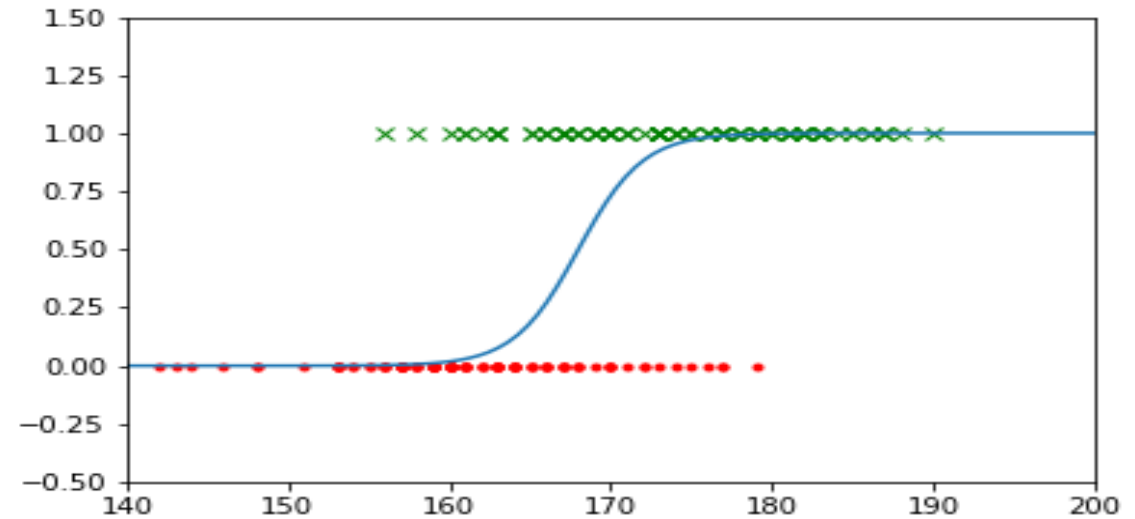
| | Classification | Regression |
|---|---|---|
| Decision tree | Lec.1 (simplified form) | |
| *k* Nearest Neighbors | Lec.5 | |
| Perceptron | Lec.6 | |
| Linear regression | | Lec. 6 |
| Logistic regression | | |
| Neural networks | | |

# Supervised learning - Where are we?

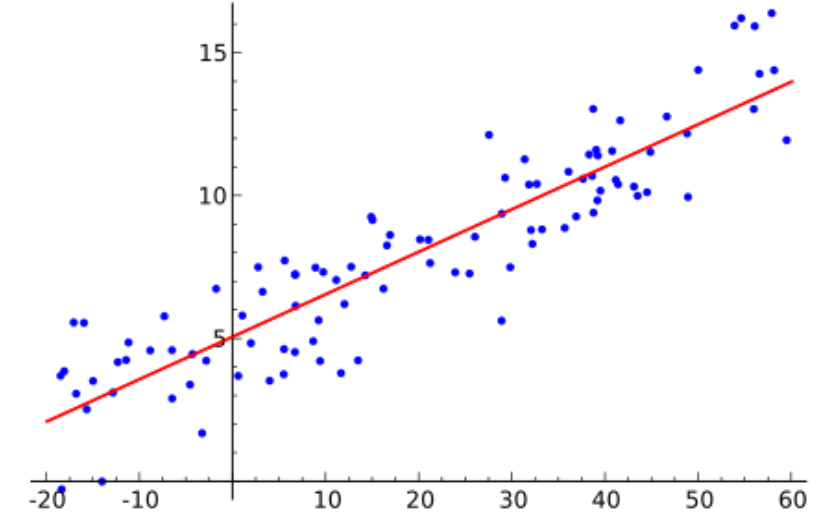| | Classification | Regression |
|---|---|---|
| Decision tree | Lec.1 (simplified form) | |
| *k* Nearest Neighbors | Lec.5 | Possible |
| Perceptron | Lec.6 | |
| Linear regression | today | Lec. 6 |
| Logistic regression | today! | |
| Neural networks | Next week | |

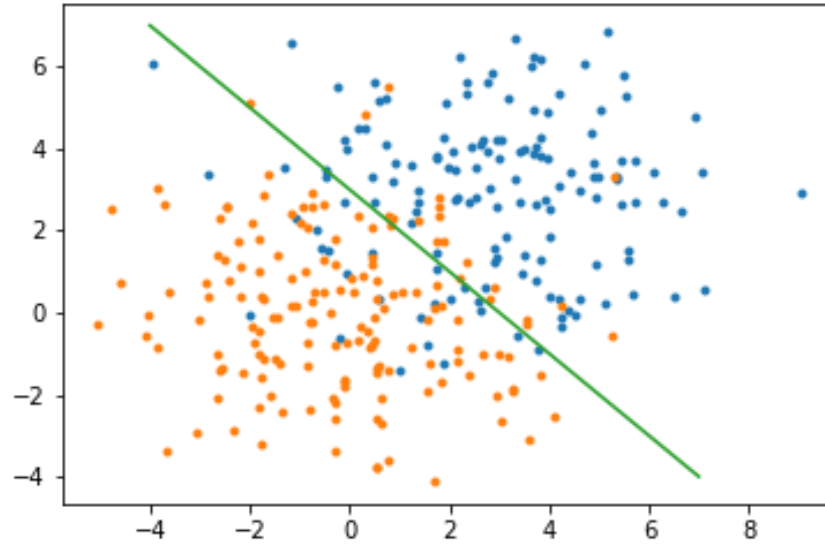# Logistic regression?



**Interesting by itself**

- A classifier
  - (not numerical regression)
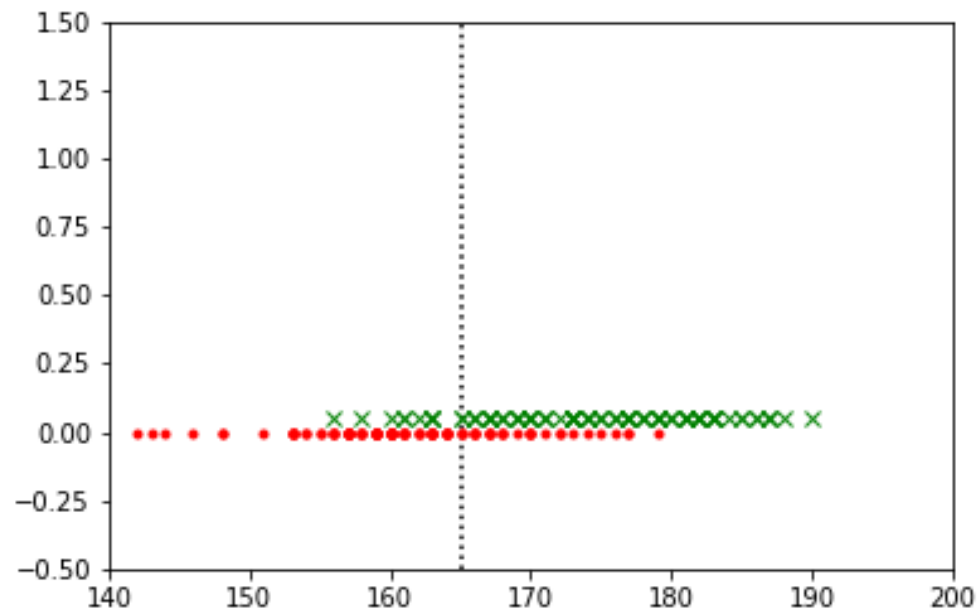- ''Standard'' (''best'') purely linear classifier
- (Not in Marsland)

**Useful tools for neural networks:**

- The logistic function:
  - Its derivative
- Loss function
- Application of the chain rule for derivatives for gradient descent

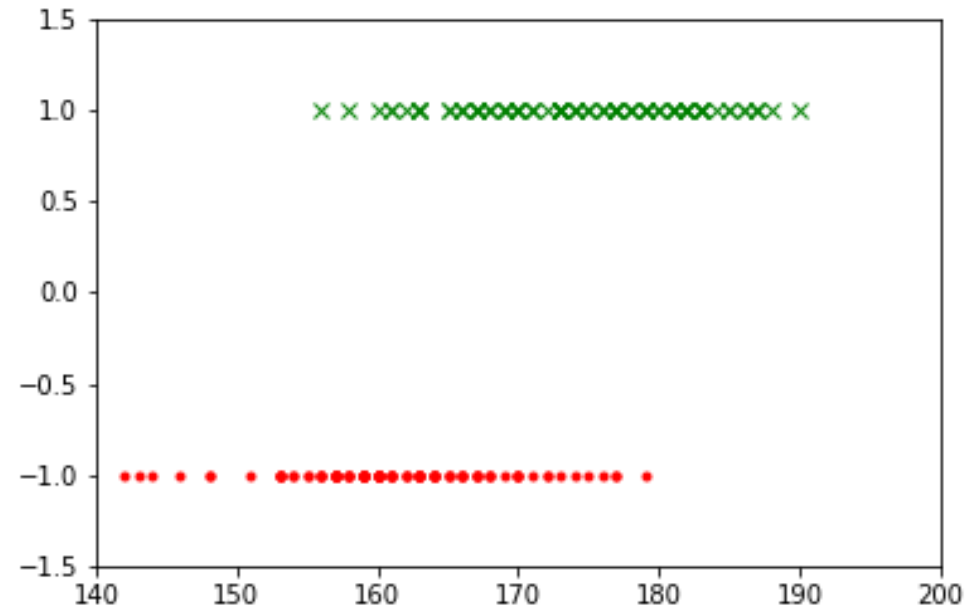| | Linear Classifier | Linear regression |
|---|---|---|
| Number of input variable | Decision boundary | Prediction |
| One | Point | Line |
| Two | Line | Plane |
| Three | Plane | Hyper-plane |
| >3 | Hyper-plane | |
| Update | Perceptron: $$w_i = w_i - \eta(y - t)x_i$$ | $$w_k = w_k - \eta \frac{2}{N} \sum_{j=1}^{N} \left( (t_j - y_j)(-x_{j,k}) \right)$$ |
| Type of y, t | {0,1} | Real numbers |

# Example: predicting gender from height



- The decision boundary should be a number: $c$

- An observation, $n$, is classified
  - *male* if *height_n* > $c$
  - *female* otherwise

- How do we determine $c$?

# Linear regression as a classifier

1. Consider the prediction of classes as prediction of the two numbers 1, -1, resp.

2. Fit a linear regressor to these data (minimizing) MSE
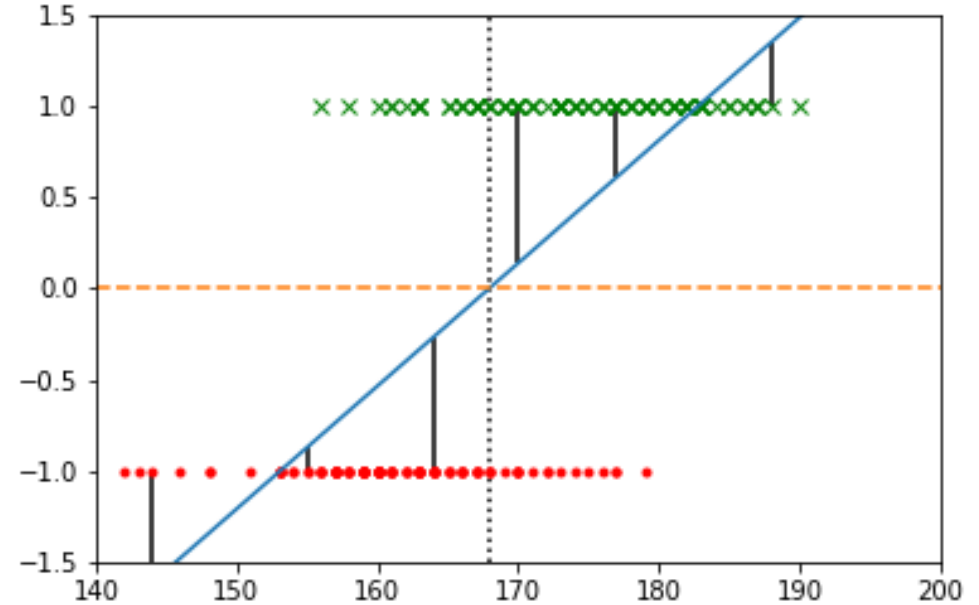
# Linear regression as a classifier

1.  Consider the prediction of classes as prediction of the two numbers 1, -1, resp.

2.  Fit a linear regressor to these data (minimizing) MSE
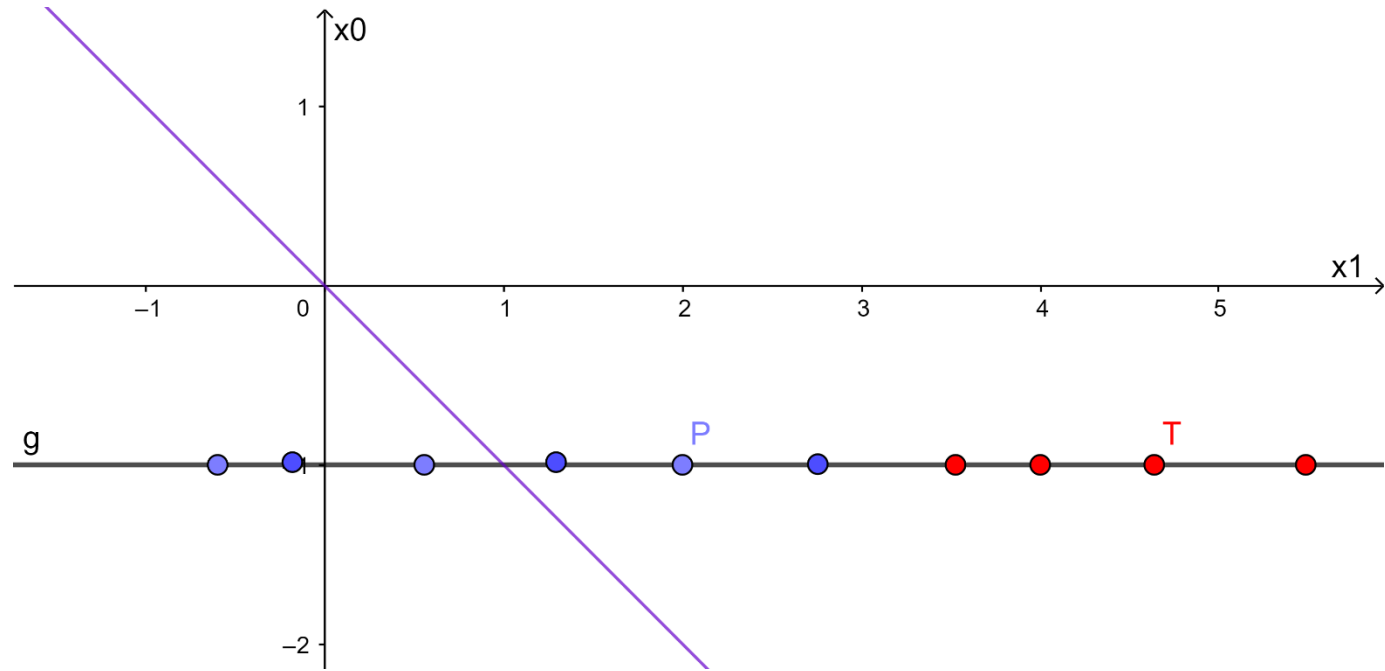
# Linear regression as a classifier

1. Consider the prediction of classes as prediction of the two numbers 1, -1, resp.

2. Fit a linear regressor to these data (minimizing) MSE

3. Predict
   - Positive class if $y > 0$ and
   - Negative class, otherwise

- Hence, decision boundary is dotted yellow line

# Example



- Consider example from last week.
- Compare Lin.reg.-classifier to perceptron
- Assume stochastic gradient descent:
  We update for one datapoint at a time

# Example



- We are in the middle of training
- Learning rate: $\eta = 0.1$
- We have so far, the following weights for the decisions:
- Positive class provided
  $$h = -w_0 + w_1 x_1 = 1 - x_1 > 0$$
  - i.e., $w_0 = -1$ and $w_1 = -1$

- Consider the point Q=(-1, -2):
  - Correctly classified
  - Perceptron: Do nothing

- Lin.reg.classifier:
- $h(Q) = 1 - 1(-2) = 3$
  - $w_0 = w_0 - \eta(y - t)x_0 = -1 - 0.1(3 - 1)(-1) = -0.8$
  - $w_1 = w_1 - \eta(y - t)x_1 = -1 - 0.1(3 - 1)(-2) = -0.6$

# Limitations

- For example
  - moving 7 (out of) 100 pos 100 steps to the right
  - the decision boundary is moved
    - from 168
    - to 171.5
  - the accuracy (on the 200 training set) goes
    - from 0.81
    - to 0.78
- Should these outliers have such an effect?



By the way:
We have here used 0 and 1 for the two classes.
This works equally fine.
Prediction: Positive class for $y > 0.5$

# Linear regression as a classifier

- The MSE seems to punish correctly classified items too severely.

# The ''correct'' decision boundary



- The (Heaviside) step function
- But:
  - How do we find the best one?
  - Not a differentiable function

# 7.2 The Logistic Function and its Derivative

IN3050/IN4050 Introduction to Artificial Intelligence and Machine Learning

# The ''correct'' decision boundary



- The (Heaviside) step function
- But:
  - How do we find the best one?
  - Not a differentiable function

# The sigmoid curve



- An approximation to the ideal decision boundary
- Differentiable
  - Gradient descent
- Mistakes further from the decision boundary are punished harder

An observation, *n*, is classified
- *male* if f(*height_n) > 0.5*
- *female* otherwise

# Exponential function - Logistic function

$$y = e^z$$

$$y = \frac{1}{1 + e^{-z}} = \frac{e^z}{e^z + 1}$$

# The logistic function

- $y = \dfrac{1}{1+e^{-z}} = \dfrac{e^z}{e^z+1}$
- A sigmoid curve
  - Other functions also make sigmoid curves e.g., $y = \tanh(z)$
- Maps $(-\infty, \infty)$ to $(0,1)$
- Monotone
- Can be used for transforming numeric values into probabilities

# The derivative of the logistic function

- $y = f(x) = \dfrac{1}{1+e^{-x}}$

- This has the form $y = g\big(h(x)\big)$
  where $g(z) = \dfrac{1}{z}$ and $z = h(x) = 1 + e^{-x}$

- Hence $f'(x) = g'(z)h'(x) = \dfrac{-1}{(1+e^{-x})^2}(-e^{-x}) =$

- $\dfrac{e^{-x}+1-1}{(1+e^{-x})^2} = \dfrac{e^{-x}+1}{(1+e^{-x})^2} - \dfrac{1}{(1+e^{-x})^2} = y - y^2 = y(1-y)$

- We will use this also in the multi-layer neural networks

# 7.3 The Logistic Regression Classifier

IN3050/IN4050 Introduction to Artificial Intelligence

and Machine Learning

# Logistic Regression:

- First sum of weighted inputs :
- $z = \sum_{i=0}^{m} w_i x_i = \boldsymbol{w} \cdot \boldsymbol{x}$
- Apply the logistic function σ to this sum

$$y = \sigma(z) = \frac{1}{1 + e^{-\vec{w} \cdot \vec{x}}}$$

- For $\boldsymbol{x} = \vec{x}$ predict
  - as the positive class if $y > 0.5$,
  - otherwise, the negative class

# Comparison: activation function

- Perceptron: step function
- Linear regression: identity
- Logistic regression:
  the logistic function

# With two features

- Two features: $x_1, x_2$
- Apply weights: $w_0, w_1, w_2$
- Let $y = -w_0 + w_1 x_1 + w_2 x_2$
- Apply the logistic function, $\sigma$, and check whether
  - $\sigma(y) = \frac{1}{1+e^{-y}} > 0.5$

Geometrically: Folding a plane along a sigmoid
The decision boundary is the intersection of this surface and the plane $p = \sigma(y) = 0.5$:
This turns out to be a straight line

# Example with two features



- Example:
  - Heights and weights
  - Acc.: = 0.95
- Blue line = decision boundary
  - Points above it gets a value > 0.5

# Understanding logistic regression 1

The following 3 slides attempt to give you an understanding of logistic regression models.

- The model is probability-based

- There are two classes t=1, t=0

- For an observation $x = \vec{x}$, we wonder:

- *How probable is it that this $\vec{x}$ belongs to class 1, and how probable is it that it belongs to class 0?*

- i.e., what are $P(t = 1|\vec{x})$ and $P(t = 0|\vec{x})$? Which is largest?

# Understanding logistic regression 2

- What are $P(t = 1|\vec{x})$ and $P(t = 0|\vec{x})$? Which is largest?

- Consider the odds: $\dfrac{P(t=1|\vec{x})}{P(t=0|\vec{x})} = \dfrac{P(t=1|\vec{x})}{1-P(t=1|\vec{x})}$
  - If this is >1, $\vec{x}$ most probably belongs to t=1, otherwise t=0
  - The odds varies between 0 and infinity

- Take the logarithm of this, $\log \dfrac{P(t=1|\vec{x})}{1-P(t=1|\vec{x})}$
  - If this is >0, $\vec{x}$ most probably belongs to t=1
  - This varies between minus infinity and plus infinity

# Understanding logistic regression 3

- $\log \frac{P(t=1|\vec{x})}{1-P(t=1|\vec{x})} > 0$ ?

- Try to find a linear expression for this, $\log\left(\frac{P(t=1|\vec{x})}{1-P(t=1|\vec{x})}\right) = \vec{w} \cdot \vec{x} > 0$

- Given such a linear expression, then

  - $\frac{P(t=1|\vec{x})}{1-P(t=1|\vec{x})} = e^{\vec{w}\cdot\vec{x}}$

- Solving this with respect to $P(t=1|\vec{x})$ yields

  - $P(t=1|\vec{x}) = \frac{e^{\vec{w}\cdot\vec{x}}}{1+e^{\vec{w}\cdot\vec{x}}} = \frac{1}{1+e^{-\vec{w}\cdot\vec{x}}}$

# A probabilistic classifier

- The logistic regression will ascribe a probability to all instances for the class t=1 (and for t=0)

- We turn it into a classifier by ascribing class t=1 if and only if
$$P(t = 1|\vec{x}) > 0.5$$

- We could also choose other cut-offs, e.g., if the classes are not equally important.



Probability of passing exam versus hours of studying

source: Wikipedia

# 7.4 Cross-Entropy Loss

IN3050/IN4050 Introduction to Artificial Intelligence and Machine Learning

# How to find the best curve?



- What are the best choices of $a$ and $b$ in $\frac{1}{1+e^{-(ax+b)}}$ ?
- Geometrically $a$ and $b$ determine the
  - Midpoint (b)
  - Steepness (a)
- of the curve
- What are the best choices of $\vec{w}$

$$y = P(t = 1|\vec{x}) = \frac{1}{1+e^{-\vec{w}\cdot\vec{x}}}$$

# Learning in the logistic regression model



- A training instance consists of
  - a feature vector $\vec{x}$
  - a label (class), $t$, which is 1 or 0.
- With a set of weights, $\vec{w}$, the classifier will assign
  - $y = P(t = 1|\vec{x}) = \frac{1}{1+e^{-\vec{w}\cdot\vec{x}}}$ to this training instance $\vec{x}$
  - where $P(t = 0|\vec{x}) = 1 - y$
- Goal: find $\vec{w}$ that maximize $P(t|\vec{x})$ of all training inst.s $(\vec{x}, t)$

# Loss function

- In machine learning we decide on an objective for the training.

- We can do that in terms of a loss function.

- The goal of the training is to minimize the loss function.

- Example: linear regression
  - Loss: Mean Square Error

- We can choose between various loss functions.

- The choice is partly determined by the learner.

- For logistic regression we choose (simplified) cross-entropy loss

# Footnote: Notation

- I observe that I haven't been consequent in notation
- I fluctuate between boldface $\boldsymbol{x}$ and non-bold with an arrow $\vec{x}$. There are no (intended) differences between the two, $\boldsymbol{x} = \vec{x}$
- I have also fluctuated between $\boldsymbol{x}_j$ and $\vec{x}^{(j)}$ for vector number $j$ in the input set. Again, the two ways of writing amount to the same.

# The money game

- I will give you 10 multiple-choice questions. You must answer all.

- I give you a million NOK before the game.

- In each round, you must bet your remaining money on the alternatives. Say there are 3 answers in the first round. You could bet any of the following, e.g.

| | Your bet | | | You keep | | |
|---|---|---|---|---|---|---|
| | Answer A | Answer B | Answer C | If A correct | If B correct | If C correct |
| Strategy 1 | 1,000,000 | 0 | 0 | 1,000,000 | 0 | 0 |
| Strategy 2 | 400,000 | 300,000 | 300,000 | 400,000 | 300,000 | 300,000 |
| Strategy 3 | 800,000 | 150,000 | 50,000 | 800,000 | 150,000 | 50,000 |

- You proceed to the next round with the money you keep.
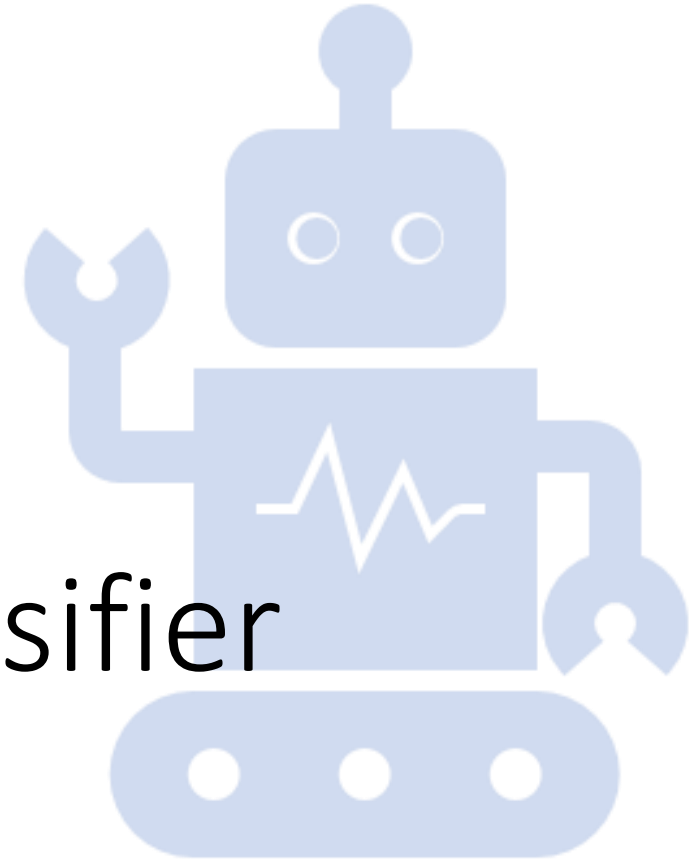
- What would be the best strategy?

# Cross-entropy loss

- The underlying idea is that we want to <span style="color:red">maximize the joint probability of all the predictions we make</span>
  - $\prod_{i=1}^{N} P\left(t^{(i)} \mid \vec{x}^{(i)}\right)$, over all the training data $i = 1, 2, \dots, N$
    - (since the training data are independent)
- This is the same as <span style="color:red">maximizing</span>
  - $\log \prod_{i=1}^{N} P\left(t^{(i)} \middle| \vec{x}^{(i)}\right) = \sum_{i=1}^{N} \log P\left(t^{(i)} \middle| \vec{x}^{(i)}\right)$
- This is the same as <span style="color:red">minimizing</span>
  - $L_{CE}(\vec{w}) = -\log \prod_{i=1}^{N} P(t^{(i)} | \vec{x}^{(i)}) = \sum_{i=1}^{N} -\log P(t^{(i)} | \vec{x}^{(i)})$
  - Which is an instance of what is called the cross-entropy loss

# More on cross-entropy loss

- When t = 1, $P(t \mid \vec{x}) = y = \frac{1}{1 + e^{-\vec{w} \cdot \vec{x}}}$

- When t= 0, $P(t \mid \vec{x}) = 1 - y$

- Since
  - $y^t = y$ when t = 1
  - $y^t = 1$ when t = 0
  - $(1-y)^{(1-t)} = 1$ when t = 1
  - $(1-y)^{(1-t)} = (1-y)$ when t = 0

- $P(t \mid \vec{x}) = y^t (1-y)^{(1-t)}$, whether t = 1 or t = 0

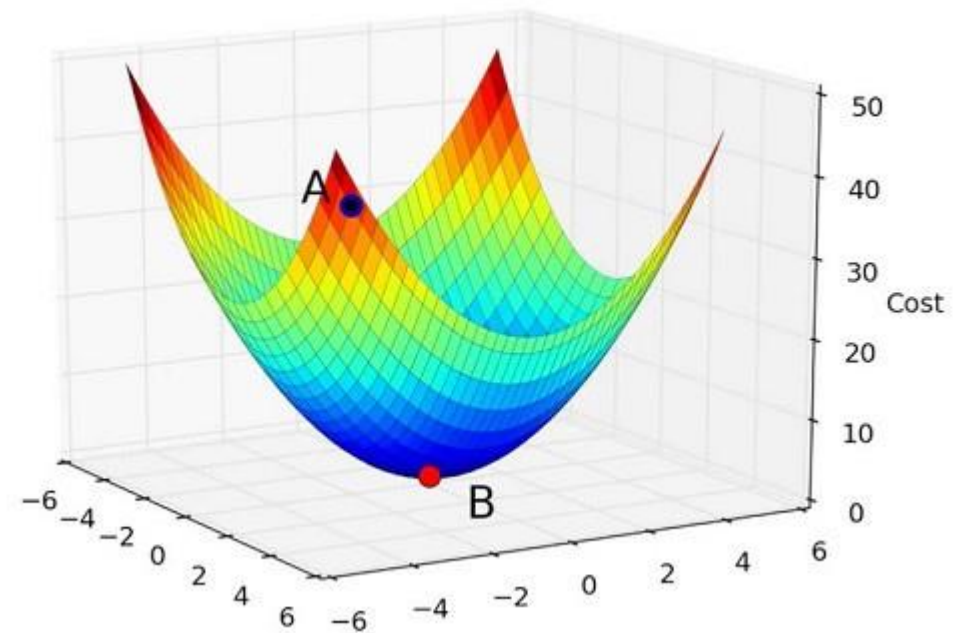# 7.5 Training the Logistic Regression Classifier

IN3050/IN4050 Introduction to Artificial Intelligence
and Machine Learning

# Gradient descent

- The loss function tells us which model is best.

- How do we find it?

- No closed-form solution, i.e., formula as there are for linear regression,

- Good news:
  - The log-loss function is convex: you are not stuck in local minima
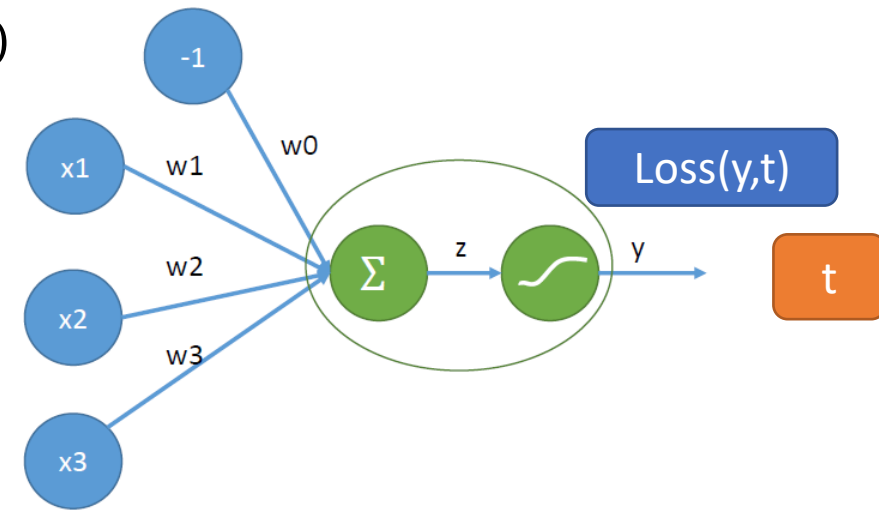  - We know which way to go

# The gradient

- We have
  - $L_{CE}(\vec{w}) = -\log \prod_{i=1}^{N} P(t^{(i)} | \vec{x}^{(i)}) = \sum_{i=1}^{N} -\log P(t^{(i)} | \vec{x}^{(i)})$
  - $P(t|\vec{x}) = y^t (1-y)^{(1-t)}$
  - $y = \sigma(\vec{w} \cdot \vec{x}) = \frac{1}{1 + e^{-\vec{w} \cdot \vec{x}}}$
- We shall find
  - $\frac{\partial}{\partial w_j} L_{CE}(\vec{w})$ for each $w_j$
  - since $\frac{\partial}{\partial w_j} L_{CE}(\vec{w}) = \sum_{i=1}^{N} -\frac{\partial}{\partial w_j} \log P(t^{(i)} | \vec{x}^{(i)})$
  - we can consider what this looks like for one pair $\left(\vec{x}^{(i)}, t^{(i)}\right)$ at a time
  - $-\frac{\partial}{\partial w_i} \log P(t|\vec{x}) = -\frac{\partial}{\partial w_i}\left(\log\left(y^t(1-y)^{(1-t)}\right)\right) =$
    $$-\frac{\partial}{\partial w_i}\left(t \log(y) + (1-t)\log(1-y)\right)$$
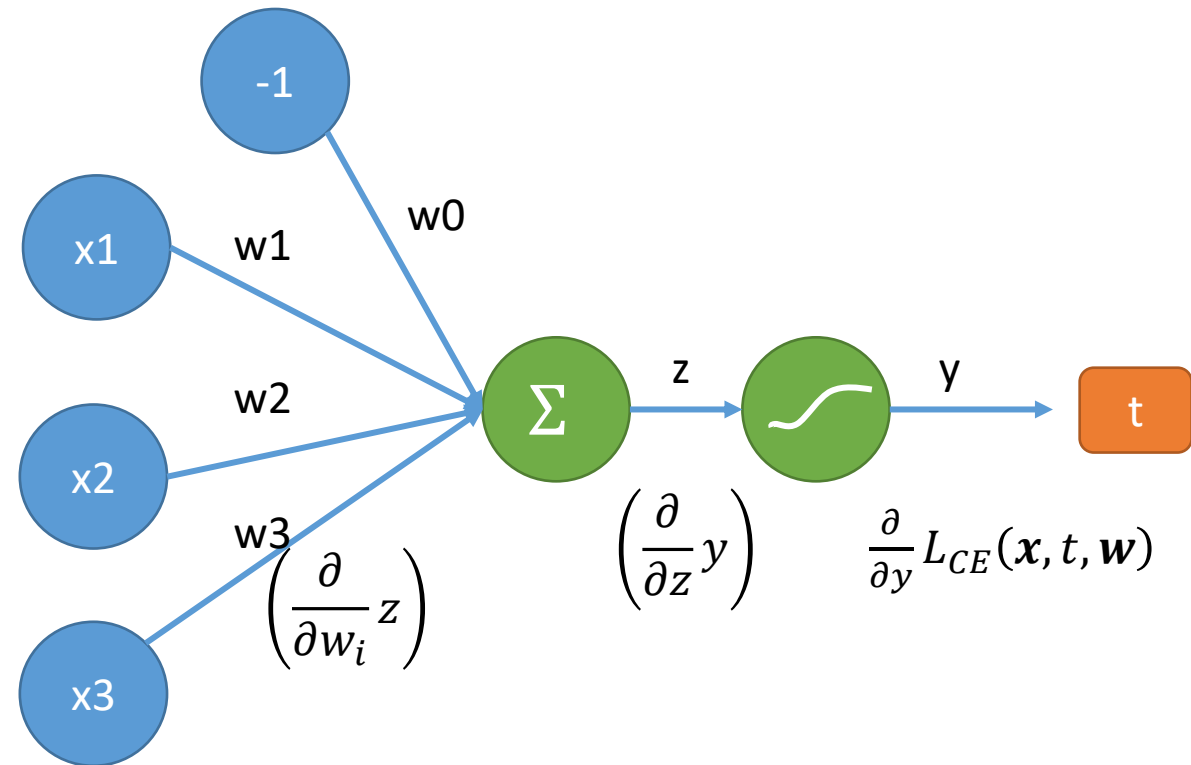
# Derivative: the chain rule

- We shall find

- $-\dfrac{\partial}{\partial w_i} \log P(t|\vec{x}) = -\dfrac{\partial}{\partial w_i}\big(t \log(y) + (1-t)\log(1-y)\big)$

- $= -\dfrac{\partial}{\partial y}\big(t \log(y) + (1-t)\log(1-y)\big)\Big(\dfrac{\partial}{\partial w_i} y\Big)$ by the chain rule for derivatives

- $\dfrac{\partial}{\partial y}\big(t \log(y) + (1-t)\log(1-y)\big) = \dfrac{t}{y} - \dfrac{(1-t)}{(1-y)} = \dfrac{t(1-y)-y(1-t)}{y(1-y)} = \dfrac{(t-y)}{y(1-y)}$

# The derivative of the logistic function

- $y = \sigma(\vec{w} \cdot \vec{x}) = \frac{1}{1+e^{-\vec{w}\cdot\vec{x}}} = \frac{1}{1+e^{-z}}$ , where $z = \vec{w} \cdot \vec{x}$

- $\frac{\partial}{\partial w_i} y = \left(\frac{\partial}{\partial z} y\right)\left(\frac{\partial}{\partial w_i} z\right)$

- $\frac{\partial}{\partial z} y = y(1-y)$ (the logistic function)

- $\frac{\partial}{\partial w_i} z = x_i$

- $\frac{\partial}{\partial w_i} y = y(1-y)x_i$

# Putting it together graphically

- $\frac{\partial}{\partial w_i} L_{CE}(\boldsymbol{x}, t, \boldsymbol{w}) =$

- $\frac{\partial}{\partial y} L_{CE}(\boldsymbol{x}, t, \boldsymbol{w}) \left( \frac{\partial}{\partial z} y \right) \left( \frac{\partial}{\partial w_i} z \right)$

# Putting it all together

- $\frac{\partial}{\partial w_i} L_{CE}(\boldsymbol{x}, t, \boldsymbol{w}) = -\frac{\partial}{\partial w_i} \log P(t|\vec{x}) = -\frac{\partial}{\partial w_i} \left( t \log(y) + (1-t)\log(1-y) \right)$

- $= -\frac{\partial}{\partial y} \left( t \log(y) + (1-t)\log(1-y) \right) \left( \frac{\partial}{\partial w_i} y \right)$

- $= -\frac{(t-y)}{y(1-y)} y(1-y) x_i = -(t-y) x_i$

- A long journey – but the result is simple

- Adding together (matrix multiplication) for all the training data yields the gradient

- $(\nabla f)_i = \frac{\partial}{\partial w_i} L_{CE}(X, T, \boldsymbol{w}) = \sum_{j=1}^{N} -(t_j - y_j) x_{j,i}$
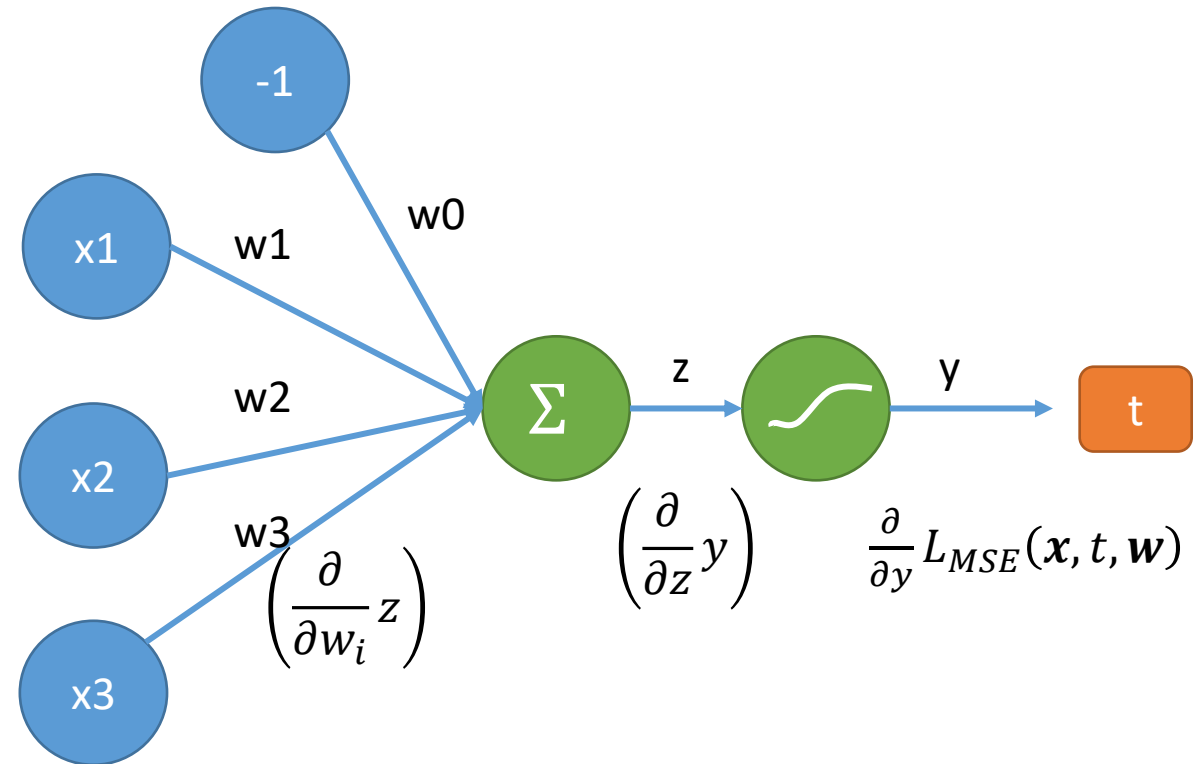
# Afterthoughts: LogReg+MSE-Loss?

- Could we have used Mean Squared Error as the loss function for the Logistic Regression classifier instead?
  - YES

- Would it work equally well?
  - NO

- Why?
  - I will show you

# What would be different?

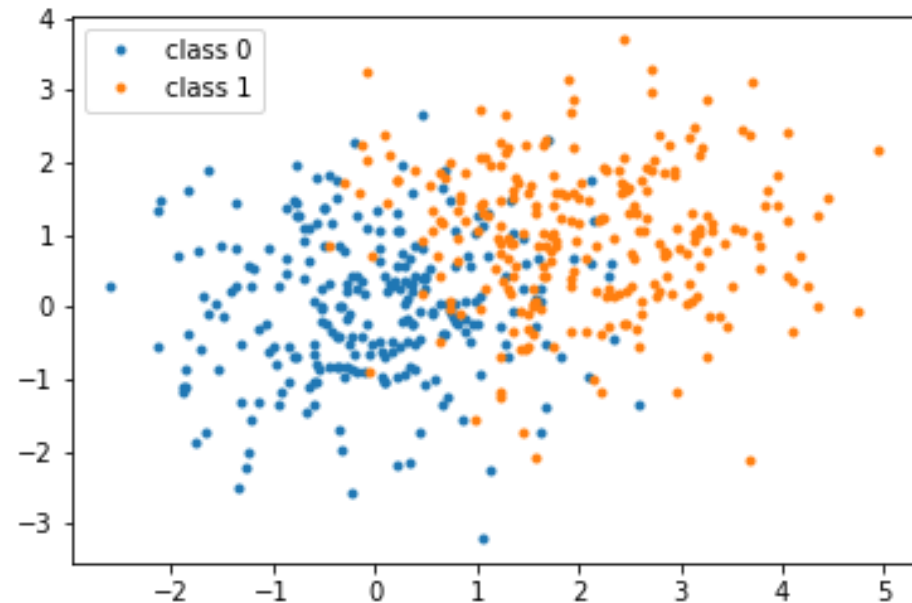- $\dfrac{\partial}{\partial w_i} L_{MSE}(\boldsymbol{x}, t, \boldsymbol{w}) =$

- $\dfrac{\partial}{\partial y} L_{MSE}(\boldsymbol{x}, t, \boldsymbol{w}) \left( \dfrac{\partial}{\partial z} y \right) \left( \dfrac{\partial}{\partial w_i} z \right) =$
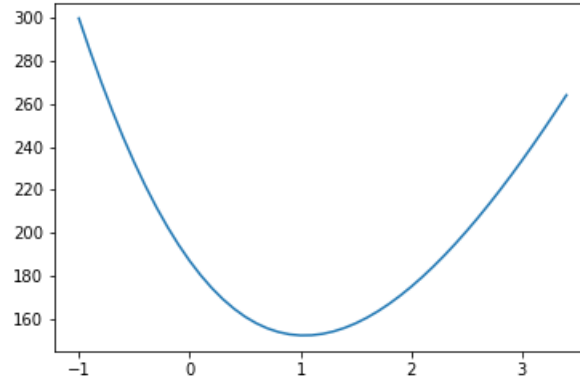
- $2(y - t)y(1 - y)x_i$

# Properties of the two loss functions

- We will consider the two loss functions on the same data set:
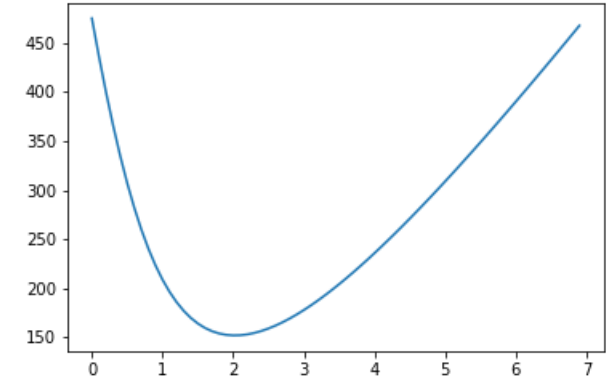  - 2 features + bias
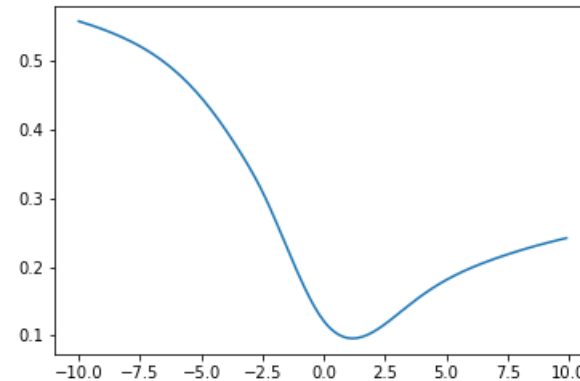  - used weekly exercises 6

# Comparison

- The CE-loss is convex
  - The only minimum is global

- The MSE-loss is not convex when applied to logistic regression



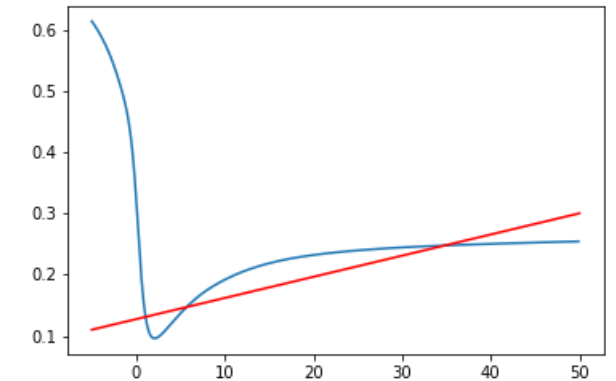$\lambda w_1 Loss_{CE}(\boldsymbol{x}, t, (-2.51, w_1, 2.02))$
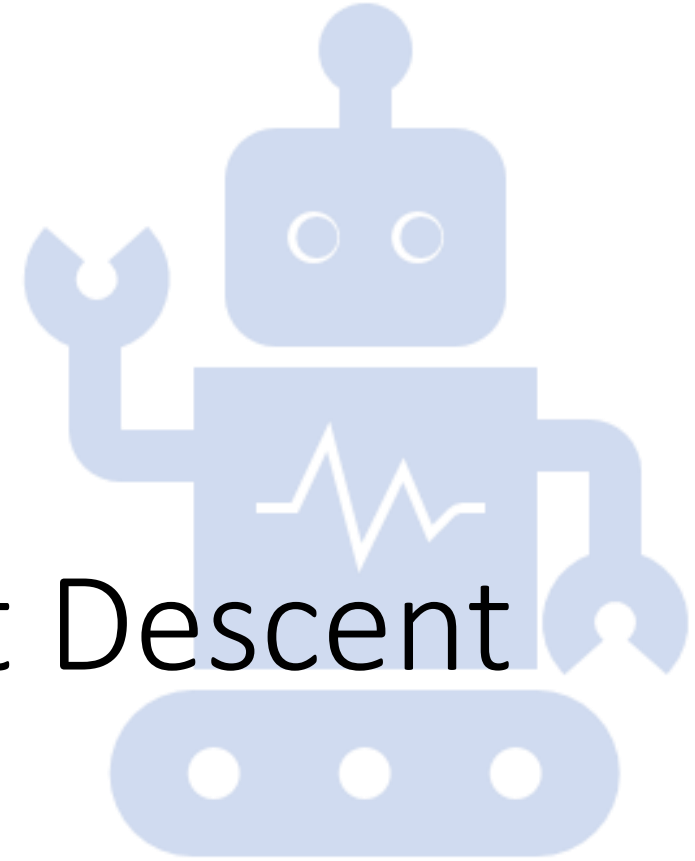
$\lambda w_2 Loss_{CE}(\boldsymbol{x}, t, (-2.51, 1.04, w_2))$

$\lambda w_1 Loss_{MSE}(\boldsymbol{x}, t, (-2.51, w_1, 2.02))$

$\lambda w_2 Loss_{MSE}(\boldsymbol{x}, t, (-2.51, 1.04, w_2))$

# 7.6 Variants of Gradient Descent

IN3050/IN4050 Introduction to Artificial Intelligence
and Machine Learning

# Variants of gradient descent

**Batch training:**
- Calculate the loss for the whole training set, and the gradient for this
- Make one move in the correct direction
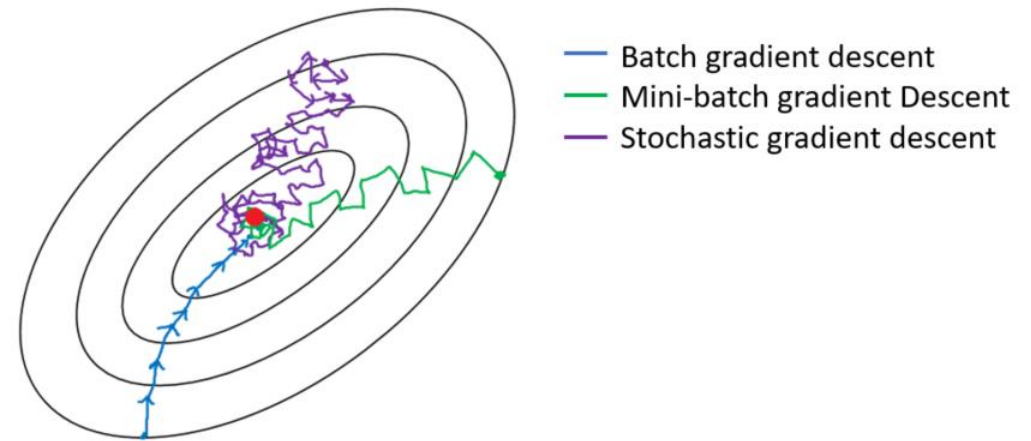- Repeat (an epoch)

- Can be slow

**Stochastic gradient descent:**
- Pick one item
- Calculate the loss for this item
- Calculate the gradient for this item and move in the opposite direction

- Each move does not have to be in towards the direction of the gradient for the whole set.

- But the overall effect may be good

- Can be faster

# Variants of gradient descent

**Mini-batch training:**

- Pick a subset of the training set of a certain size
- Calculate the loss for this subset
- Make one move in the direction opposite of this gradient
- Repeat (an epoch)

- A good compromise between the two extremes
- (The other two are subcases of this)



— Batch gradient descent
— Mini-batch gradient Descent
— Stochastic gradient descent

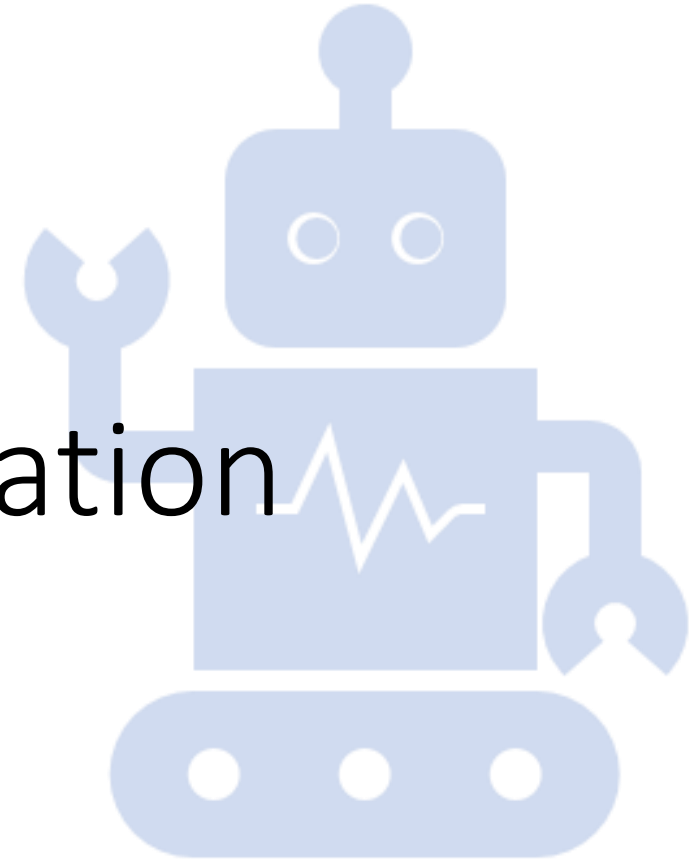https://suniljangirblog.wordpress.com/2018/12/13/variants-of-gradient-descent/
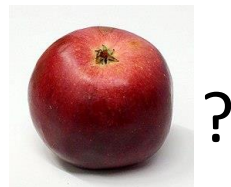
# 7.7 Multi-Class Classification One vs. Rest

IN3050/IN4050 Introduction to Artificial Intelligence

and Machine Learning

# Multi-class classification

## Classification

- Assign a label (class) from a finite set of labels to an observation



?

- So far, many algorithms and examples have been binary: *yes-no, 1-0*

- But many classification tasks are multi-class:
  - To each observation $x$ choose one label from a finite set **C**

- What is different?

# Multi-class classification

- A finite set, **C**, of $n$ different labels ($n > 2$)
- To each observation $x$ choose one label from the set **C**

We will consider two approaches:

- *One vs. rest classifier*
  - (also called one vs. all)
- *Multinomial logistic regression,* or *softmax regression*

# Algorithms

| Binary | Multi-class |
|--------|-------------|
| • Perceptron | • Decision tree |
| • Linear Regression | • *k*NN |
| • Logistic Regression | • Naive Bayes |
| | ============ |
| | • Perceptron |
| | • Multinomial Logistic Regression |

# 1-of-N or "one hot encoding"

- The labels might be categorical:
  - 'apple', tomato', 'dog', 'horse'
- The algorithms demand numerical attributes.
- First attempt
  - 'apple' = 1
  - 'tomato' = 2
  - 'dog' = 3
  - etc.
- Why isn't this a good idea?

- Better:
  - 'apple' = (1, 0, 0, 0, 0, 0)
  - 'tomato' = (0, 1, 0, 0, 0, 0)
  - 'dog' = (0, 0, 1, 0, 0, 0)
  - etc.
- Both the target and the predicted value are vectors.

# From multi-label to multi-class
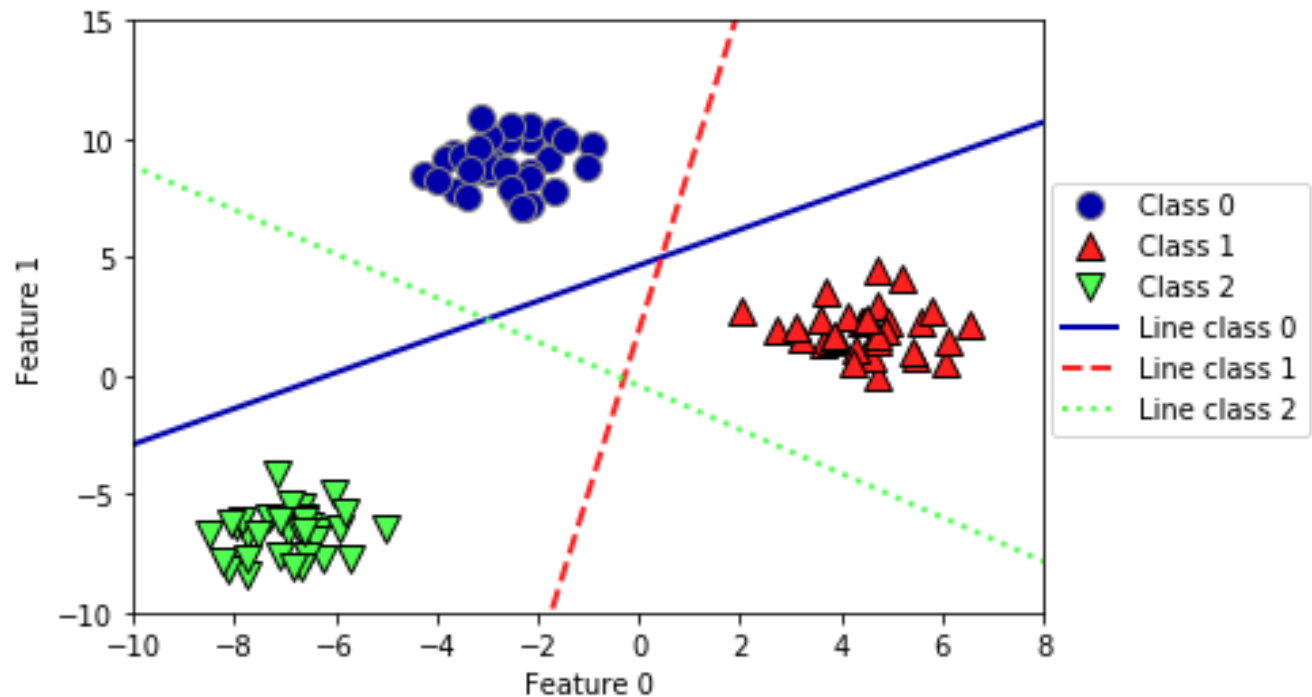
| Multi-label classifier | Multi-class classifier |
|---|---|
| • Make $n$ different classifiers, one for each class<br><br>• For classifier $j$:<br>   • consider class $j$ the positive class<br>   • all other items in the negative class<br>   • train a classifier $f_j$<br><br>• Application<br>   • Assign a label $c_j$ to an item if and only if it is classified as positive by $f_j$. | • "To each observation $x$ choose one label from the set $\mathbf{C}$"<br><br>• How can a multi-label classifier be turned into a multi-class classifier? |

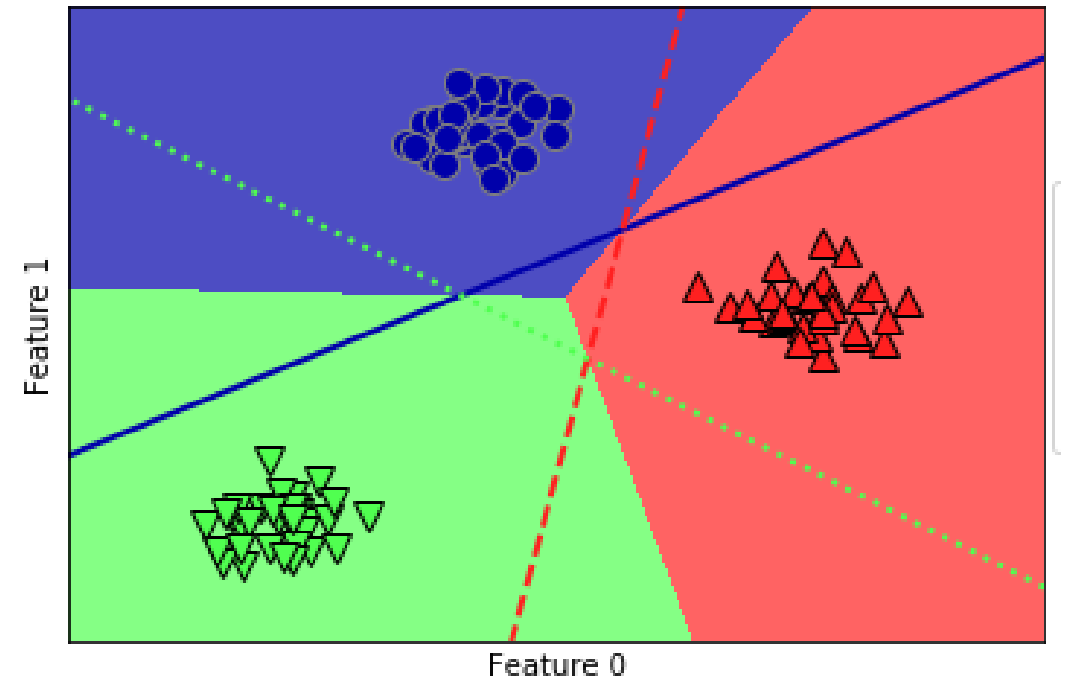# One vs. rest (also called one vs. all)

- Start like the multi-label classifier: make one classifier for each class.

- It is easy to decide for items which fall into exactly one class

- But what if they fall into
  - More than one class?
  - No classes?



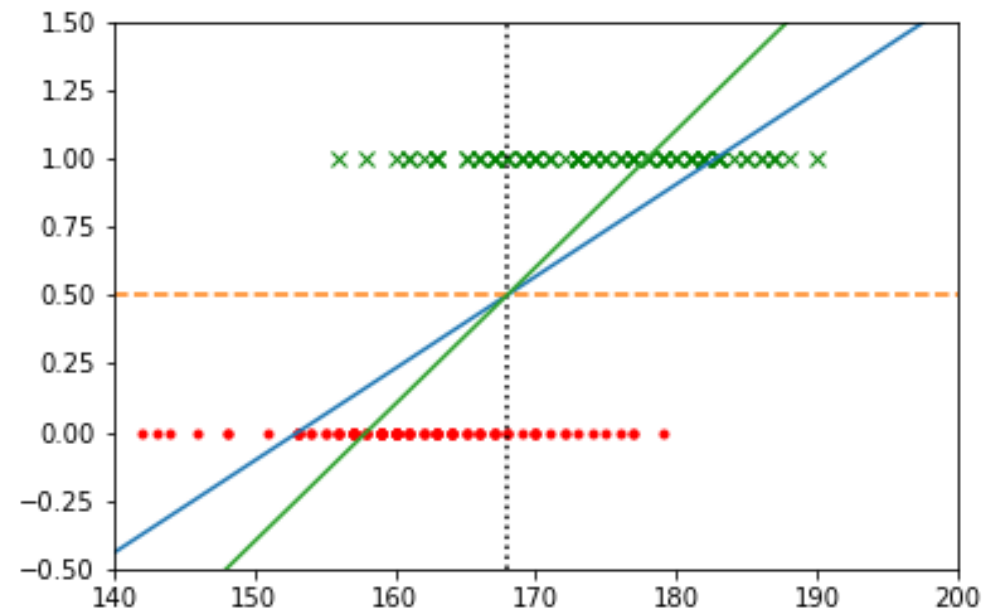https://github.com/amueller/introduction_to_ml_with_python

# One vs. rest

- If each classifier predicts a score, compare the scores for the classes

- Choose  the class with the highest score.

- E.g., log. reg.:
  - Probability of being red: 0.8
  - Probability of being blue: 0.7
  - Choose red



https://github.com/amueller/introduction_to_ml_with_python
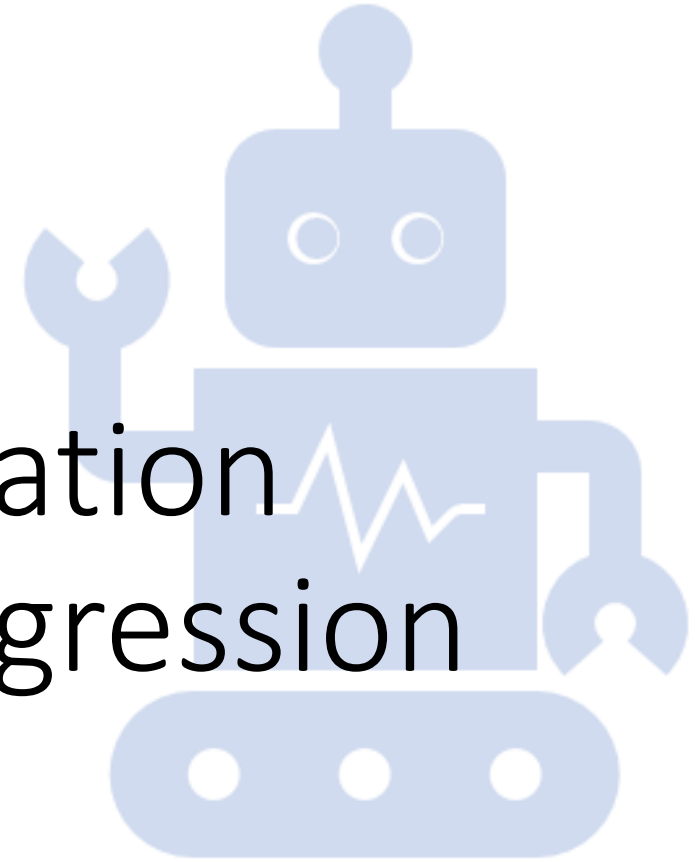
# Footnote: Linear Regression Classifier

- We could in principle also make a multi-class linear regression classifier.

- Beware:
  - Different numerical weights for the same classifier
  - Makes it difficult to compare numbers across classifiers

- Solution:
  - Normalize the weights

- Not part of the syllabus

- Stick to LogReg for one vs. rest.

# 7.8 Multi-Class Classification Multinomial Logistic Regression

IN3050/IN4050 Introduction to Artificial Intelligence

and Machine Learning
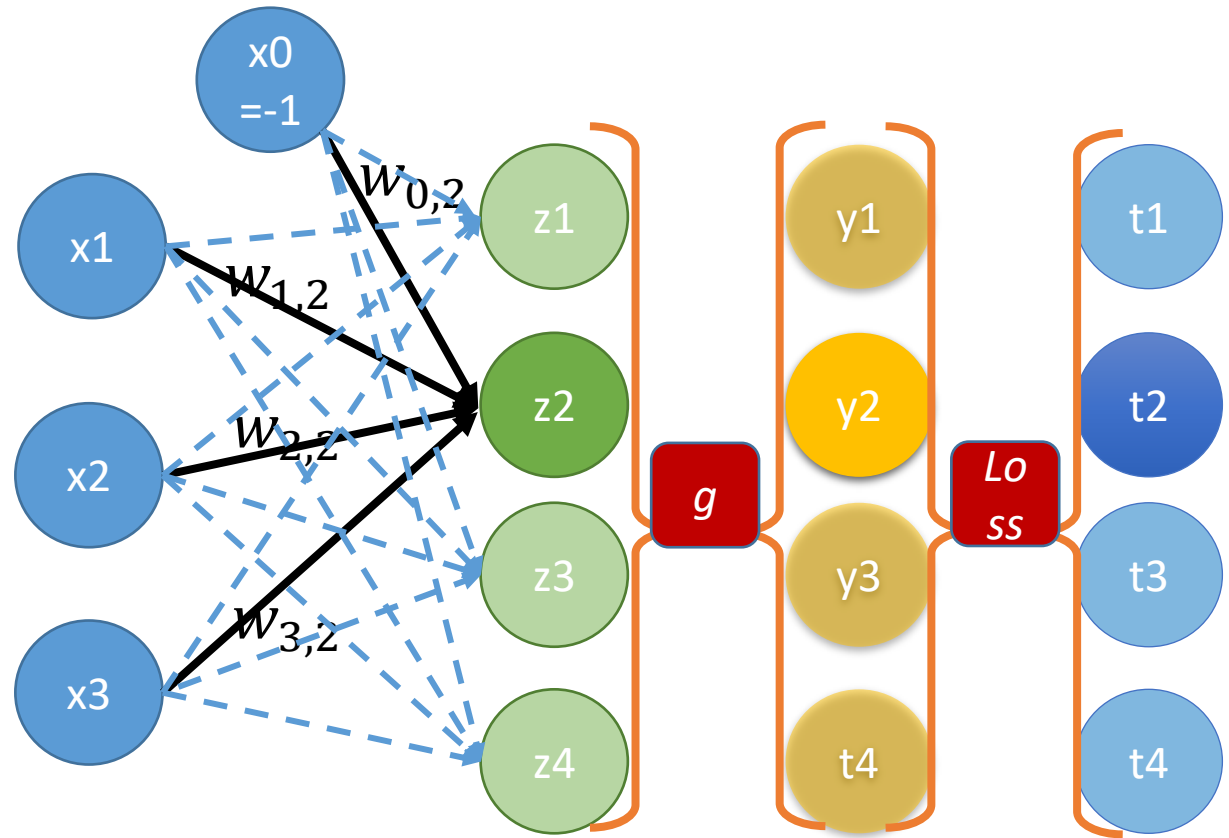
# Towards Multinomial Logistic Regression

- On the way, we will apply a bird's eye perspective
- Comparing perceptron, linear regression, logistic regression
- Compare to Marsland
- Shortly describe a multi-class perceptron
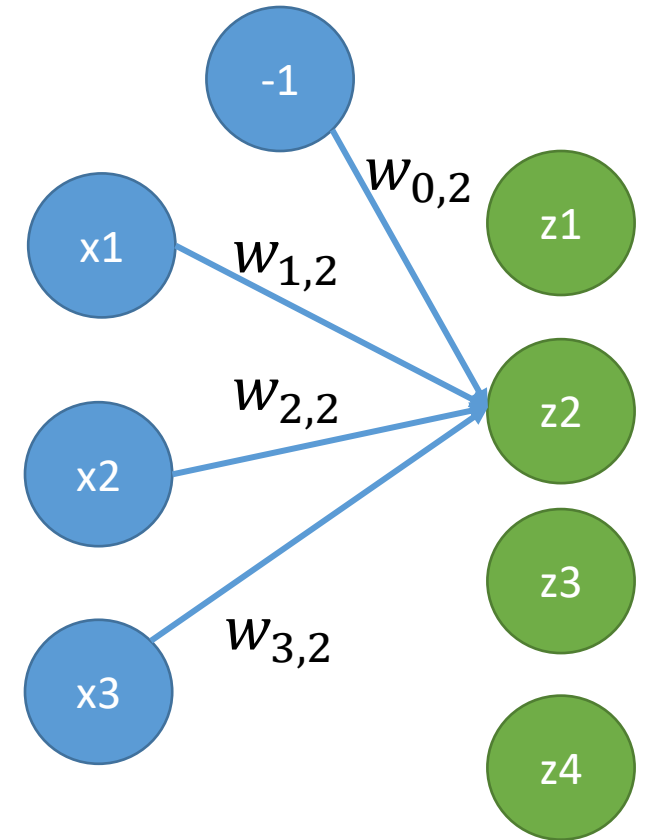
# A general view

- Common
  - $z_j = \sum_{i=0}^{n} w_{i,j} x_i$
  - $w_{i,j}$ is the weight into node j from node i
    - Some index in opposite order

| Binary classifiers | | |
|---|---|---|
| **Classifier** | **g** | **Loss** |
| Perceptron | Step | 0-1 loss |
| Lin. Regr. | Identity | MSE |
| Log.Regr. | Logistic | Cross-entropy |

# Connections going into a node

$$\begin{bmatrix} x_0 & x_1 & x_2 & \cdots & x_m \end{bmatrix} \begin{bmatrix} w_{0,1} & w_{0,2} & \cdots & w_{0,n} \\ w_{1,1} & w_{1,2} & \cdots & w_{1,n} \\ w_{2,1} & w_{2,2} & \cdots & w_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{m,1} & w_{m,2} & \cdots & w_{m,n} \end{bmatrix} = \begin{bmatrix} z_1 & z_2 & \cdots & z_n \end{bmatrix}$$
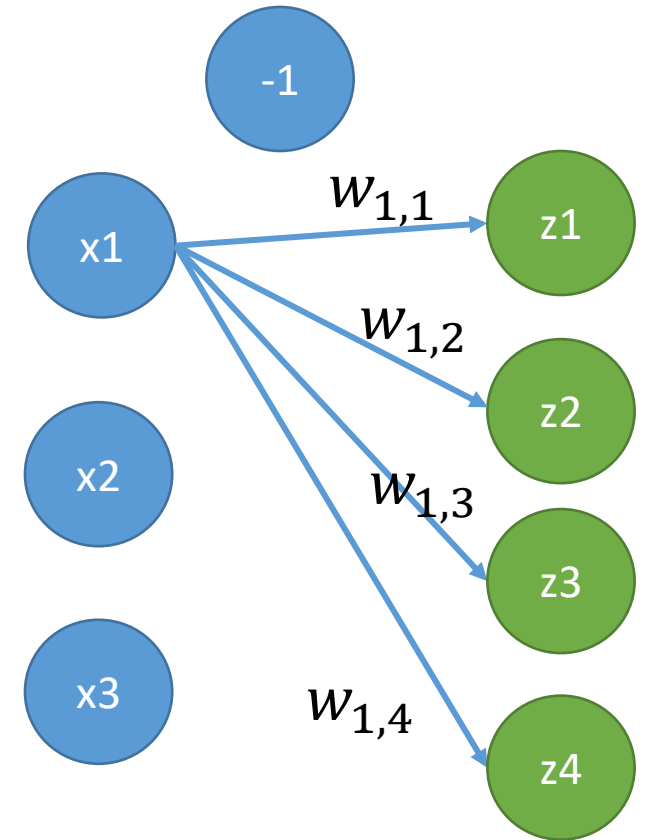
# Connections going out of a node

$$
\begin{bmatrix} x_0 & x_1 & x_2 & \cdots & x_m \end{bmatrix}
\begin{bmatrix}
w_{0,1} & w_{0,2} & \cdots & w_{0,n} \\
w_{1,1} & w_{1,2} & \cdots & w_{1,n} \\
w_{2,1} & w_{2,2} & \cdots & w_{2,n} \\
\vdots & \vdots & \ddots & \vdots \\
w_{m,1} & w_{m,2} & \cdots & w_{m,n}
\end{bmatrix}
=
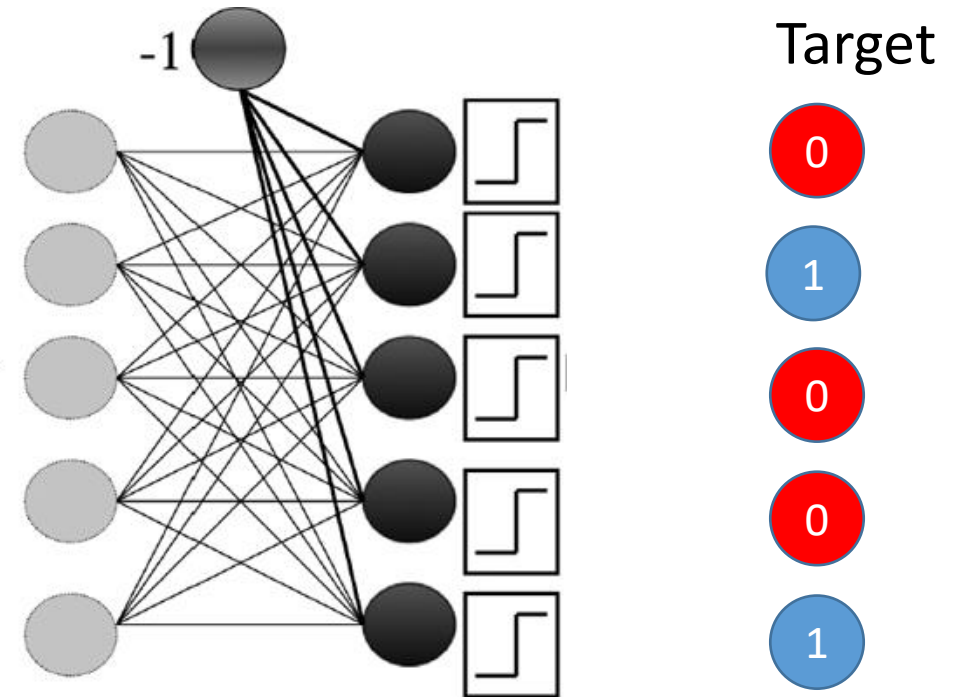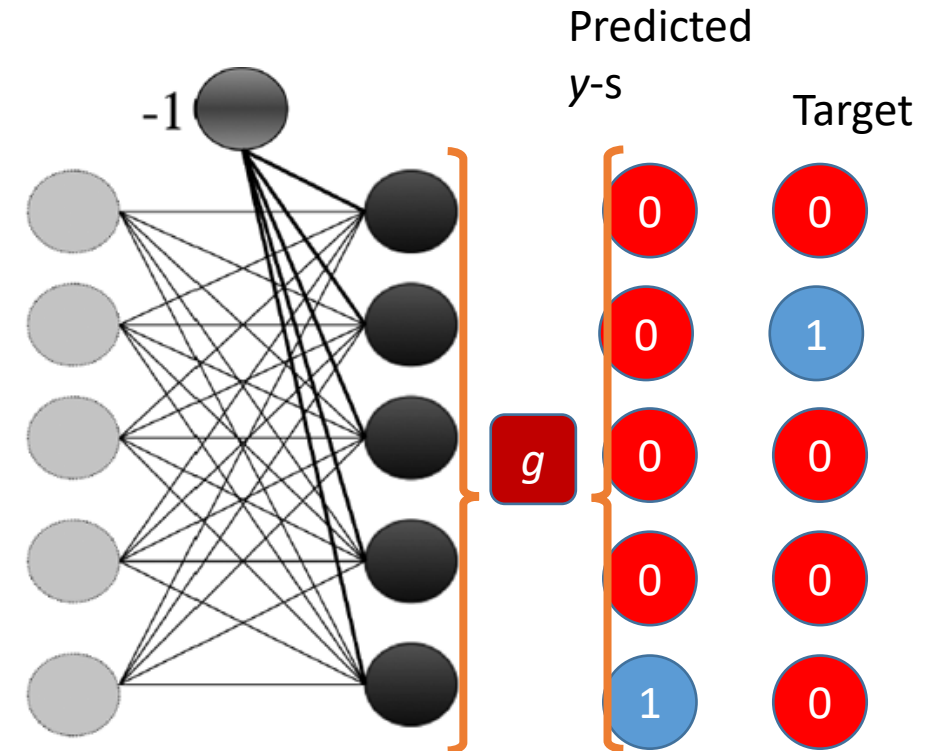\begin{bmatrix} z_1 & z_2 & \cdots & z_n \end{bmatrix}
$$

# Multi-label perceptron

- Marsland's description of perceptron:
  - Possible targets: (0,1,0,0,1)
  - $y_j = g(z_j) = \begin{array}{l} 1 \ if \ z_j > 0 \\ 0 \ if \ z_j \leq 0 \end{array}$
  - Loss is 0-1 loss for each $j$
- Describes a multi-label classifier
  - Each $y_j$ depends only on the $w_{h,k}$ where $k = j$
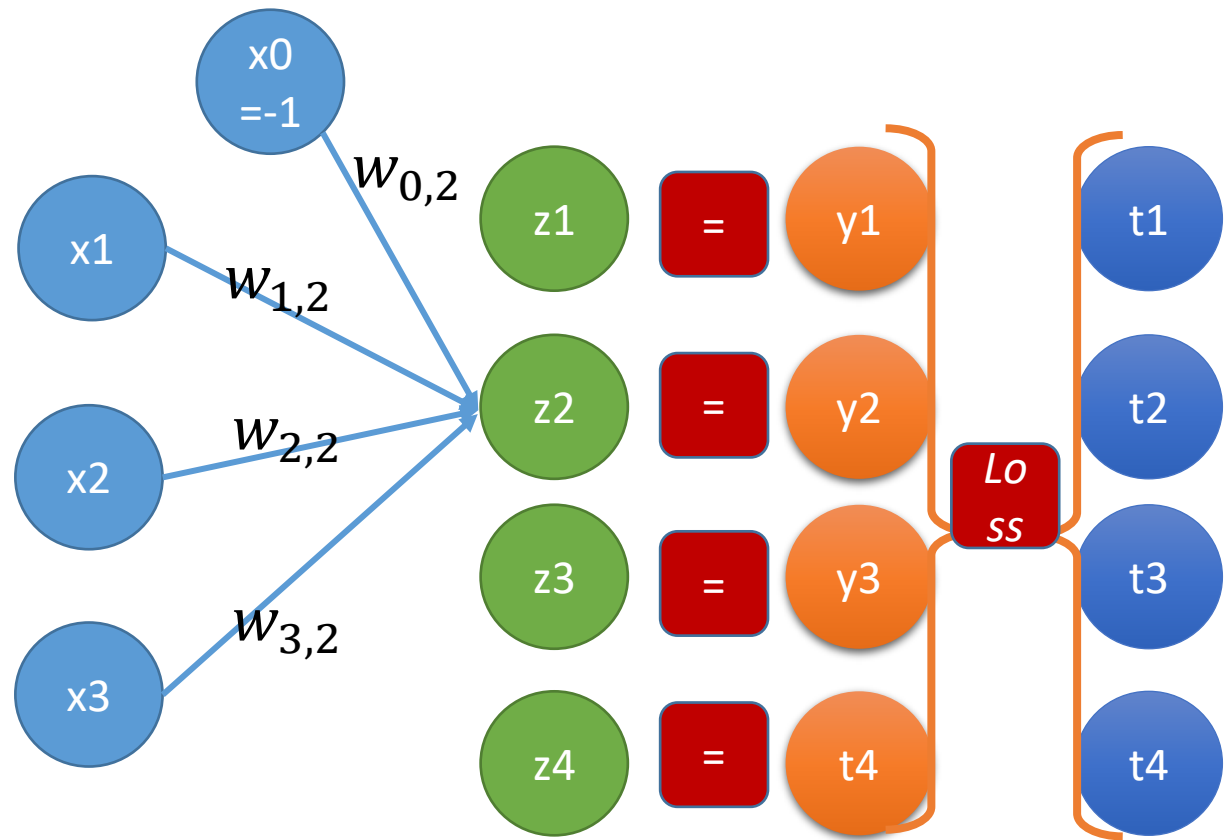  - This could have been described as $n$ independent binary perceptrons

# Multi-class perceptron

- The target contains one 1, the rest are 0s

- $g(z_1, z_2, \ldots, z_n) =$

- $argmax\ (z_1, z_2, \ldots, z_n)$
  - The index with the max value

- The update rule (0-1 loss)
  - $w_{i,j} = w_{i,j} - \eta(y_j - t_j)x_i$
  - will correct for $j = 2$ and $j = 5$
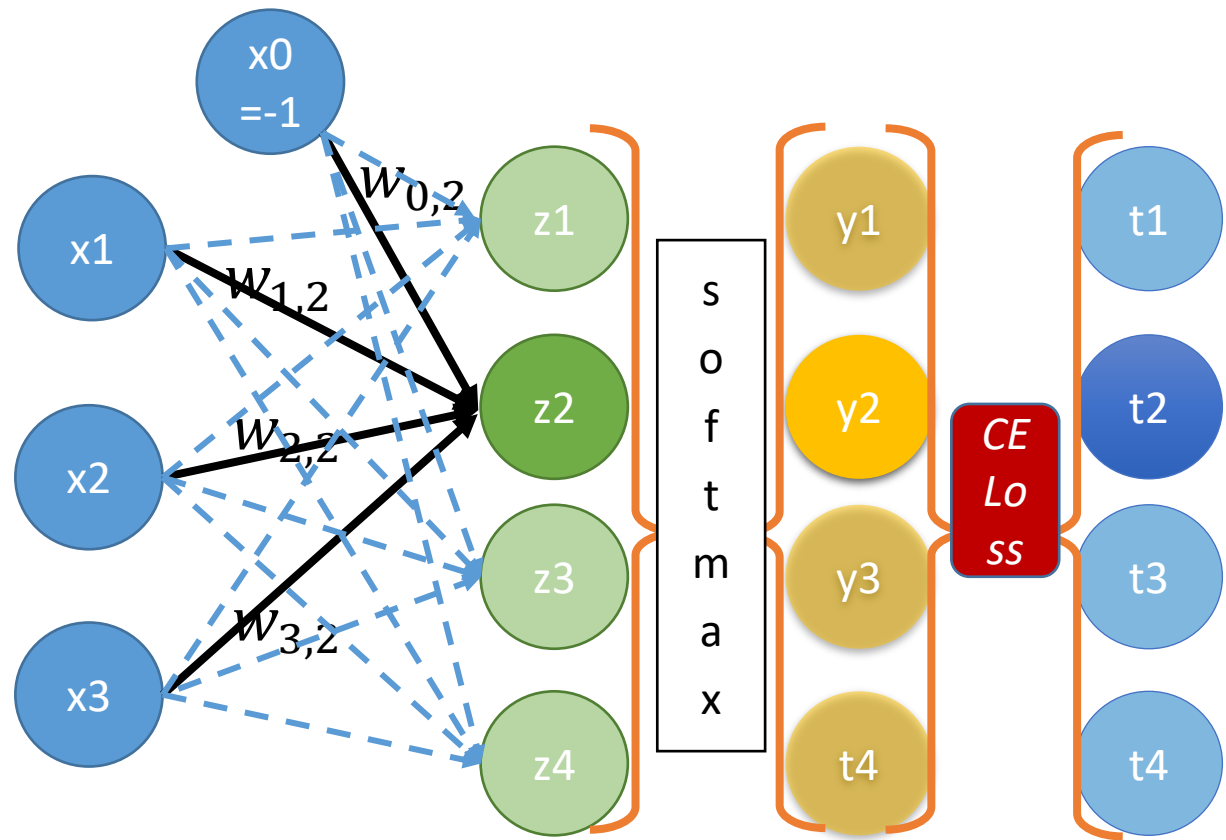  - leaves the other weights unaltered

# Multi-output linear regression

- $y_j = g(z_j) = z_j$

- MSE-loss:

$$\sum_{k=1}^{N} (\sum_{j=1}^{n} (y_{k,j} - t_{k,j})^2)$$

  - $n$ output nodes

  - $N$ input items

- $y_j$ independent of $w_{h,k}$ $h \neq k$,

  - hence corresponds to $n$ independent models

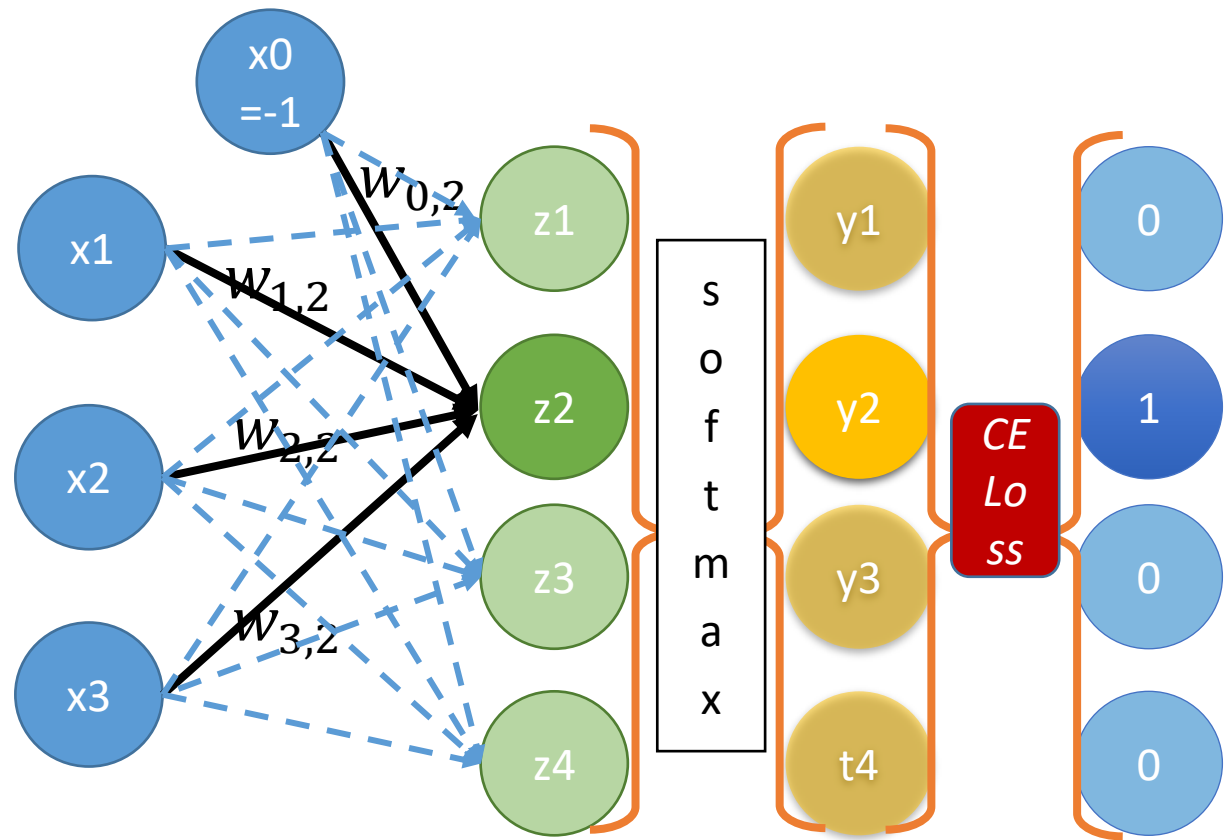  - (Gets more interesting for multi-layer networks)

# Multinomial Logistic Regression

- $z_j = \sum_{i=0}^{n} w_{i,j} x_i$

- Apply the softmax-function, $S$, where

  - $y_j = (S(z_1, \ldots, z_n))_j = \dfrac{e^{z_j}}{\sum_{k=1}^{n} e^{z_k}}$

- Observe:

  - $y_j$ depends on all the $z_k$
  - If $z_h > z_k$ then $y_h > y_k$
  - $0 < y_j < 1$
  - $\sum_{j=1}^{n} y_j = 1$
  - A probability distribution

  - $P(C_j | \vec{x}) = \dfrac{e^{\overrightarrow{w_j} \cdot \vec{x}}}{\sum_{k=1}^{n} e^{\overrightarrow{w_k} \cdot \vec{x}}}$
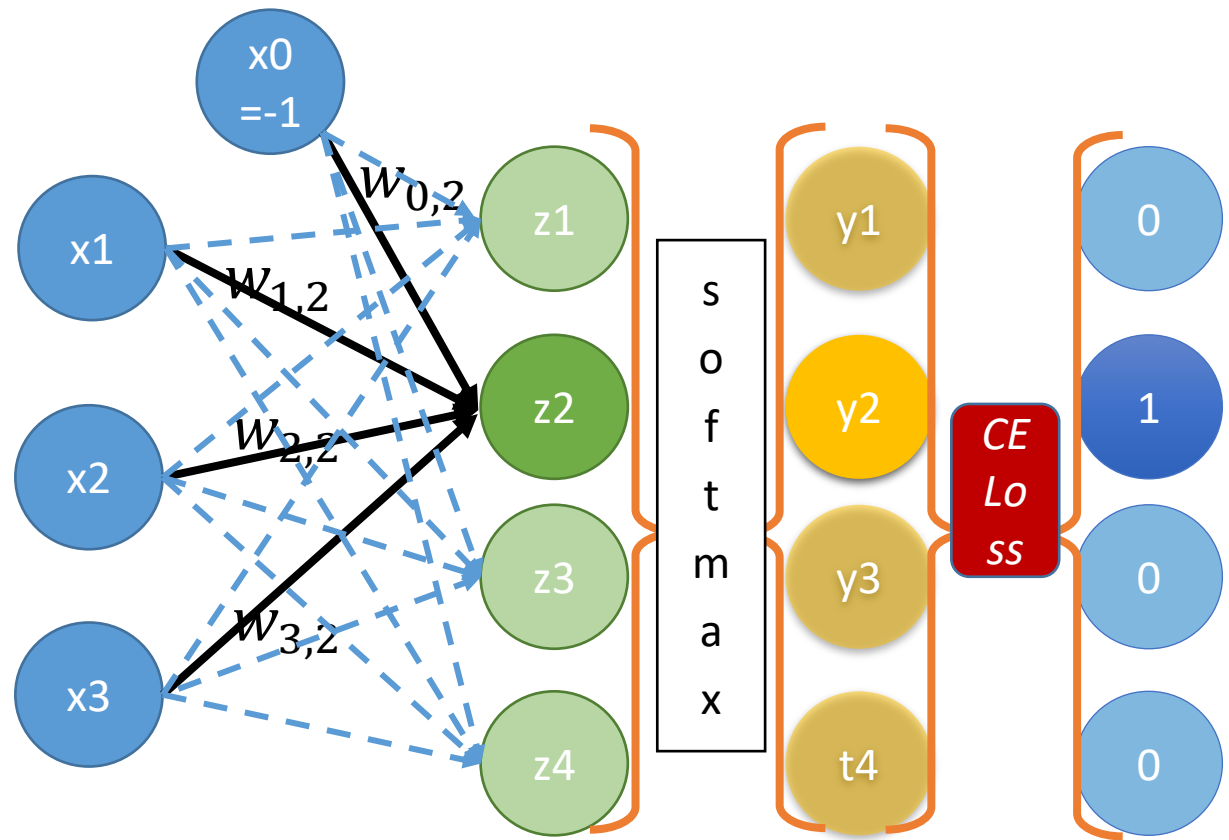
# Training Multinomial Logistic Regression 1

- The target has the form $(0,0, \dots, 0,1,0, \dots, 0)$, say
  - $t_s = 1$ and $t_j = 0$ for $j \neq s$

- We compare
  - $\mathbf{y} = (y_1, y_2, \dots y_n)$

- to the target labels
  - $\mathbf{t} = (t_1, t_2, \dots t_n)$

- using cross-entropy loss
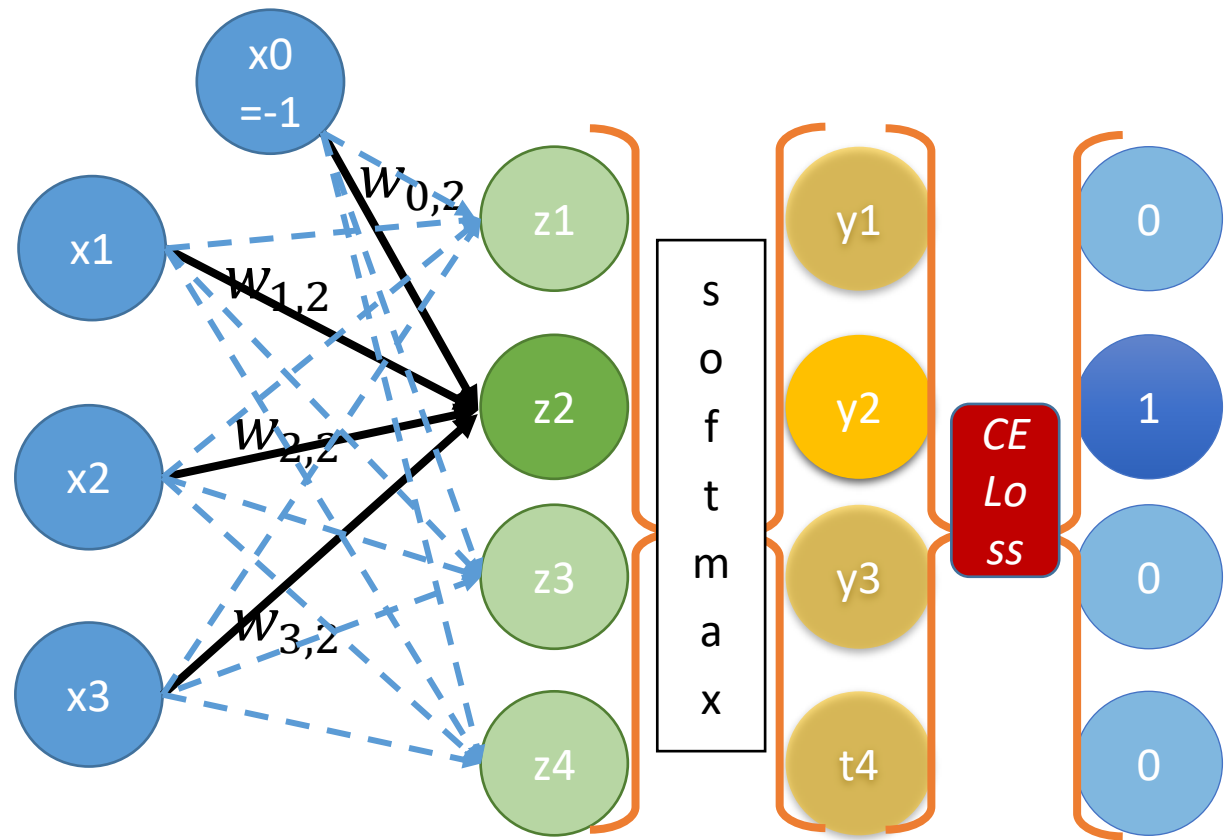  - $L_{CE}(\boldsymbol{y}, \boldsymbol{t}) = -\sum_{j=1}^{n} t_j \log y_j = -\log y_s$

# Training Multinomial Logistic Regression 2

- $y_j = \frac{e^{z_j}}{\sum_{k=1}^{n} e^{z_k}}$

- $L_{CE}(\boldsymbol{y}, \boldsymbol{t}) = -\sum_{j=1}^{n} t_j \log y_j = -\log y_s$

- Goal: to find $\frac{\partial}{\partial w_{i,j}} L_{CE}(\boldsymbol{x}, t, \boldsymbol{w})$ for all $w_{i,j}$

- Use the chain-rule for derivatives

- A little more complicated than for LogReg

- The result is simple, though

- $\frac{\partial}{\partial w_{i,j}} L_{CE}(\boldsymbol{x}, t, \boldsymbol{w}) = (y_j - t_j) x_i$

# Training Multinomial Logistic Regression 3

- $w_{i,j} = w_{i,j} + \eta(t_j - y_j)x_i$

- if $t_s = 1$:
  - $w_{i,s} = w_{i,s} + \eta(1 - y_s)x_i$
  - $w_{i,j} = w_{i,j} - \eta(y_j)x_i$
    - for $j \neq s$

- Here
  - $z_j = \sum_{i=0}^{m} w_{i,j}x_i$
  - $y_j = \dfrac{e^{z_j}}{\sum_{k=1}^{n} e^{z_k}}$

- Observe: All weights are updated for each observation

# Applying Multinomial Logistic Regression

- We can use this as a probabilistic classifier
  - $P(C_j | \vec{x}) = \dfrac{e^{\overrightarrow{w_j} \cdot \vec{x}}}{\sum_{k=1}^{n} e^{\overrightarrow{w_k} \cdot \vec{x}}}$

- To make hard decisions use

$$argmax_{j=1,\ldots,n} \frac{e^{\overrightarrow{w_j} \cdot \vec{x}}}{\sum_{k=1}^{n} e^{\overrightarrow{w_k} \cdot \vec{x}}}$$