

Binary classification setting:

$$z = \vec{w}^T \vec{x} + b$$

$$y = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{if } z < 0 \end{cases}$$

0-1 loss to measure the error of our linear classifier:

$$L_{0-1}(y, t) = \begin{cases} 0 & \text{if } y = t \\ 1 & \text{if } y \neq t \end{cases}$$

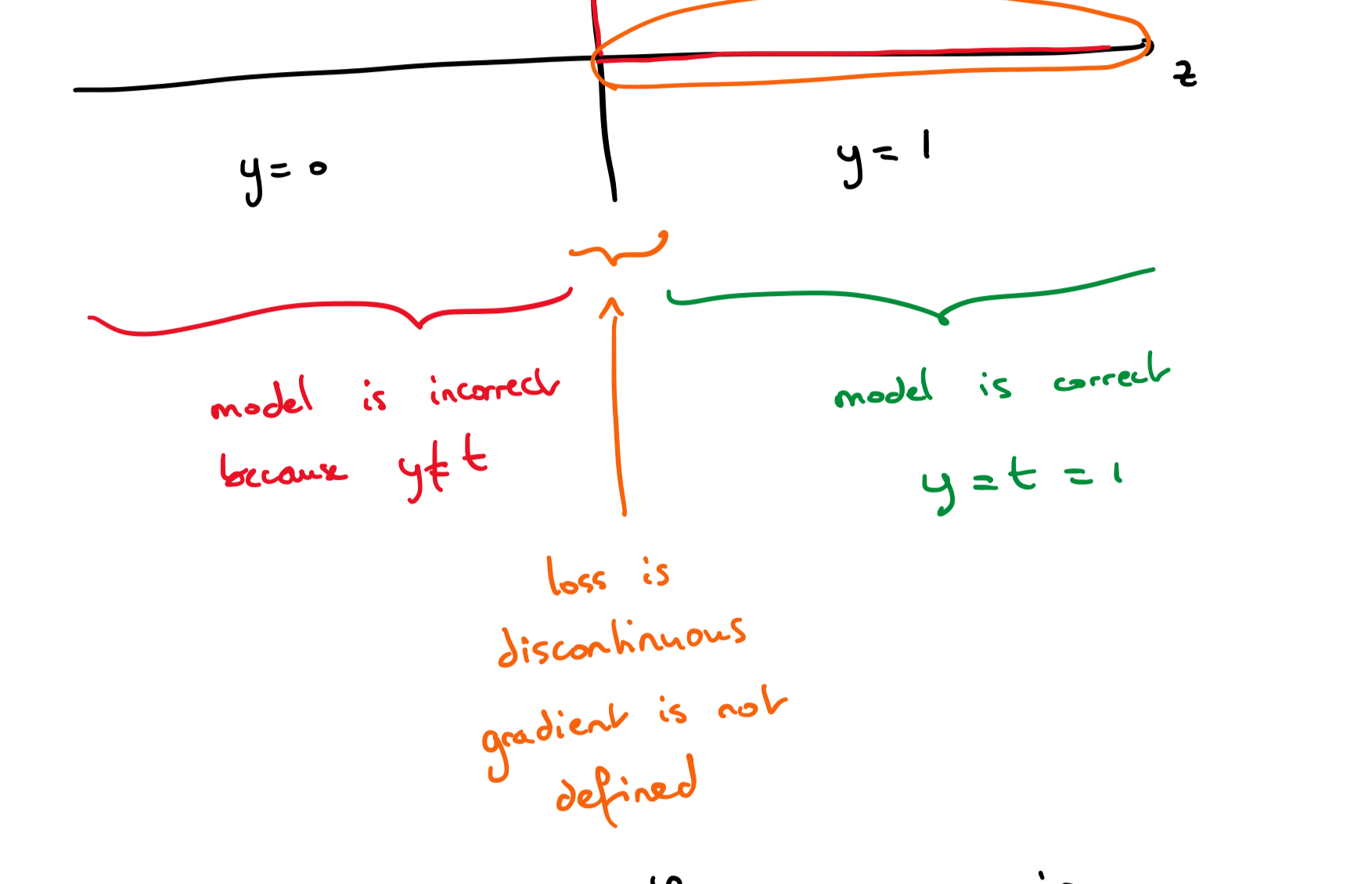
prediction of the model  $\nearrow$  target  $\nwarrow$

$\min L_{0-1} \approx \min$  the number of training examples which are misclassified.

If we want to use  $L_{0-1}$  with gradient descent to train our linear classifier,

$$\frac{\partial L_{0-1}}{\partial w_j} = \frac{\partial L_{0-1}}{\partial z} \frac{\partial z}{\partial w_j}$$

is 0 almost everywhere



Gradient descent on the  $L_{0-1}$  is not going to update the weights.

Least squares loss

$$L_{SE}(y, t) = \frac{1}{2} (y - t)^2$$

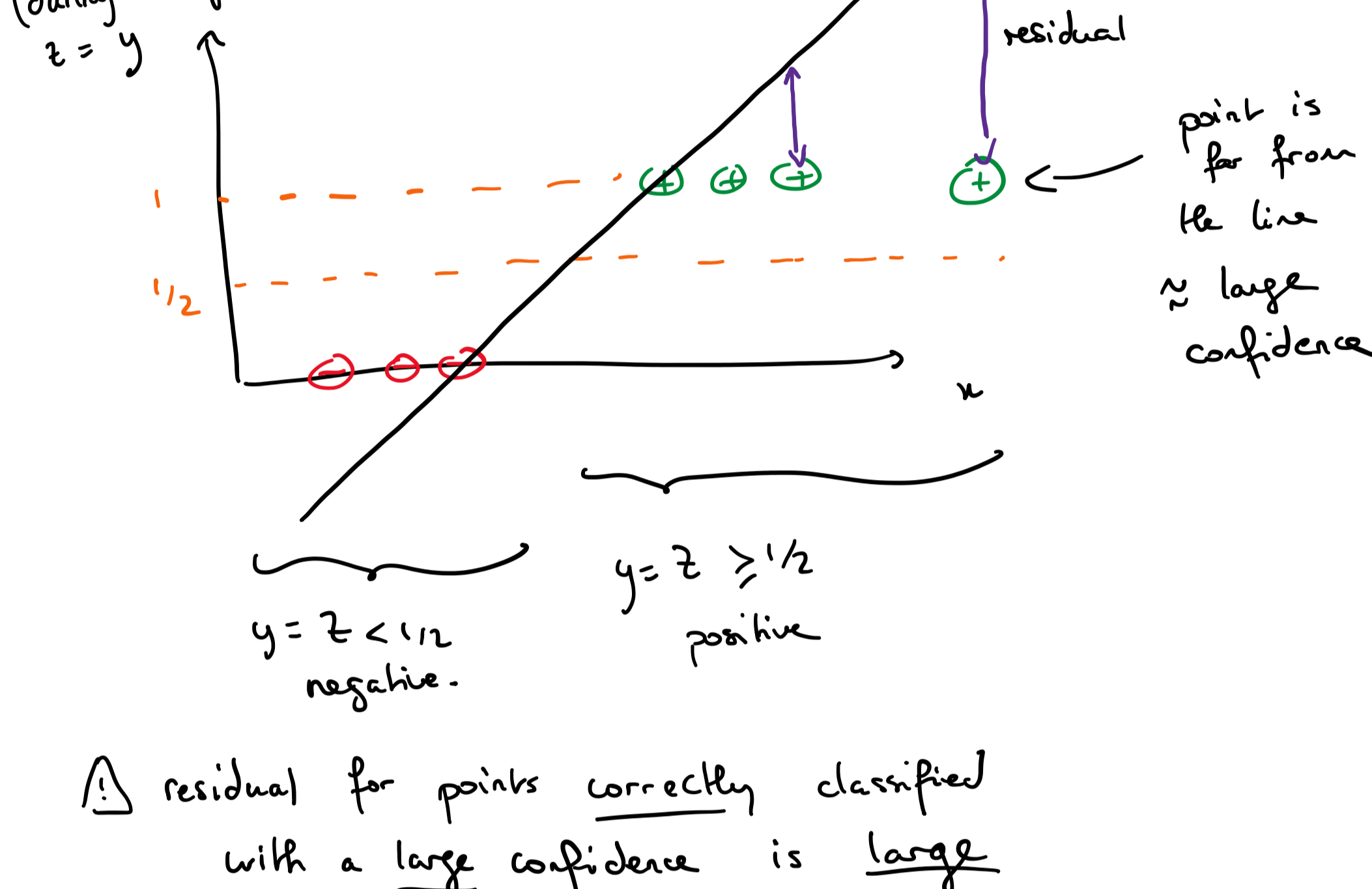
was originally a scalar (in the regression setting)

Can we use this loss in classification as well?

$$z = \vec{w}^T \vec{x} + b$$

$$y = \begin{cases} 1 & \text{if } z \geq 1/2 \\ 0 & \text{if } z < 1/2 \end{cases}$$

now we don't really need to implement the threshold during training



residual for points correctly classified with a large confidence is large

loss will be large for these points whereas (if the loss reflects how wrong the model is) it should be small

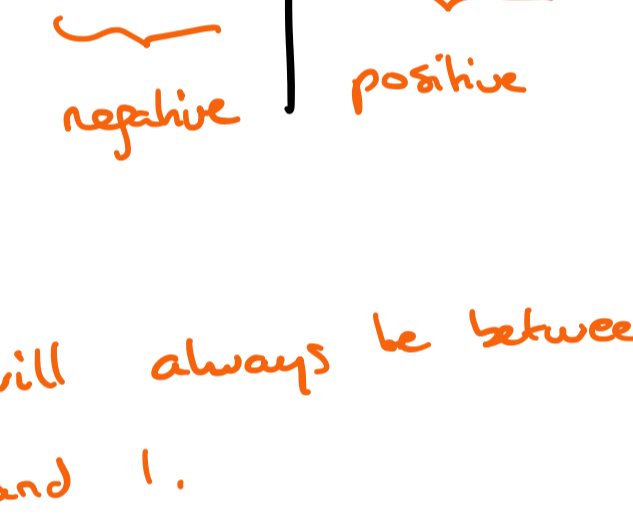
We can't use the  $L_{SE}$  to train our binary classifier with gradient descent.

Logistic function

Turn our linear model into a log-linear model.

Logistic function (non-linearity)

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



Log-linear model:

$$z = \vec{w}^T \vec{x} + b$$

$$y = \sigma(z)$$

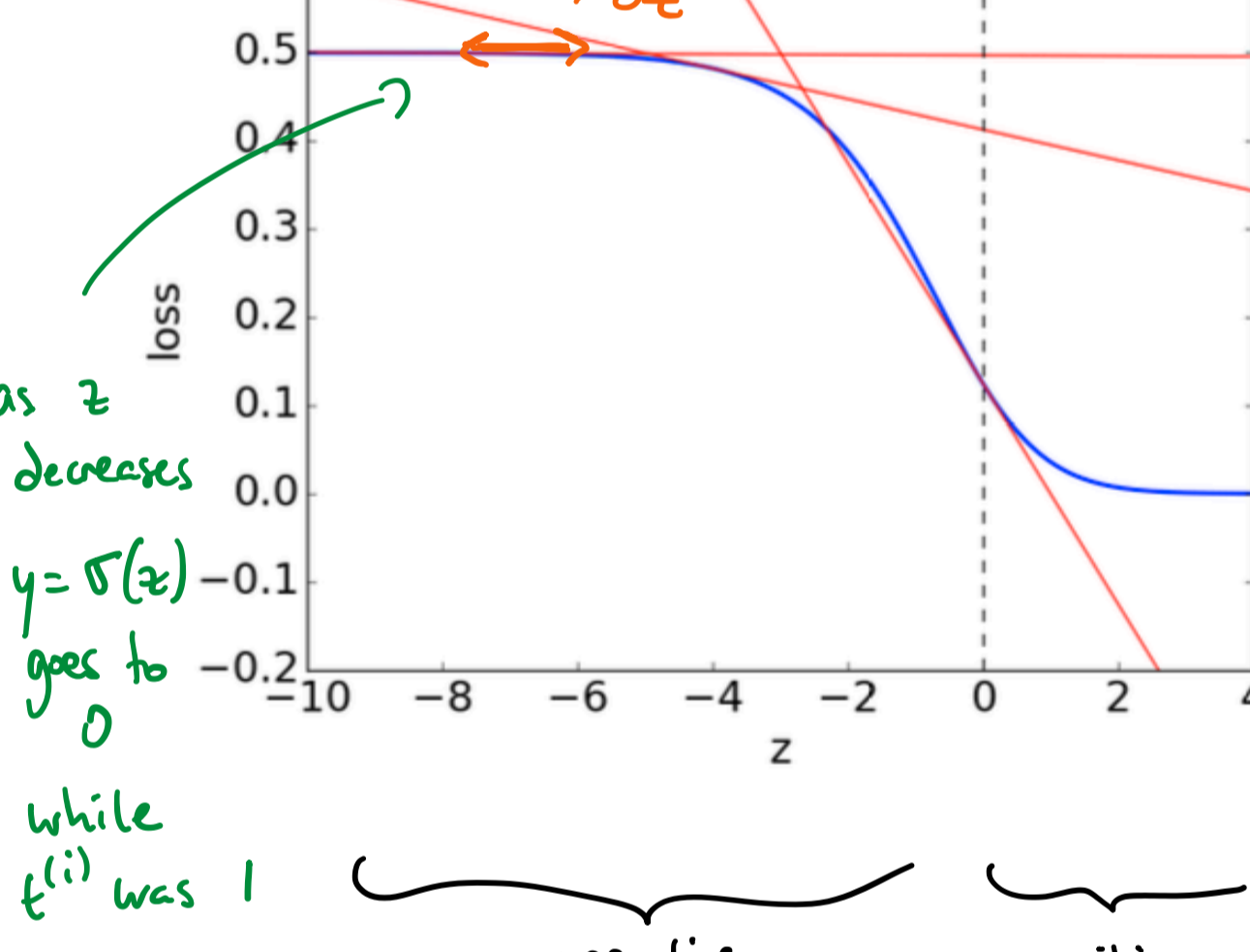
logits

Combined with our  $L_{SE}$ .

y will always be between 0 and 1. As before, we can interpret  $y < 1/2$  as negative

$y \geq 1/2$  as positive

Here for an example with label  $t^{(i)} = 1$



while  $t^{(i)}$  was 1 so  $L_{SE}$  increases towards  $1/2$

$$L_{SE} = \frac{1}{2} (y - t)^2 \rightarrow \frac{1}{2}$$

$y \rightarrow 0$        $t = 1$

Recall  $w_j \leftarrow w_j - \alpha \frac{\partial L}{\partial w_j}$

Using the chain rule, we have:

$$\frac{\partial L}{\partial w_j} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial w_j}$$

For a positive example ( $t^{(i)} = 1$ ), when model is confident in a negative prediction ( $z \rightarrow -\infty$ )

then  $\frac{\partial L}{\partial z} \rightarrow 0$

$\frac{\partial L}{\partial w_j} \rightarrow 0$  / instable.

Gradient descent won't be able to optimize our weights properly if the model's parameters  $w$  leads to a very incorrect prediction on one of the training examples.

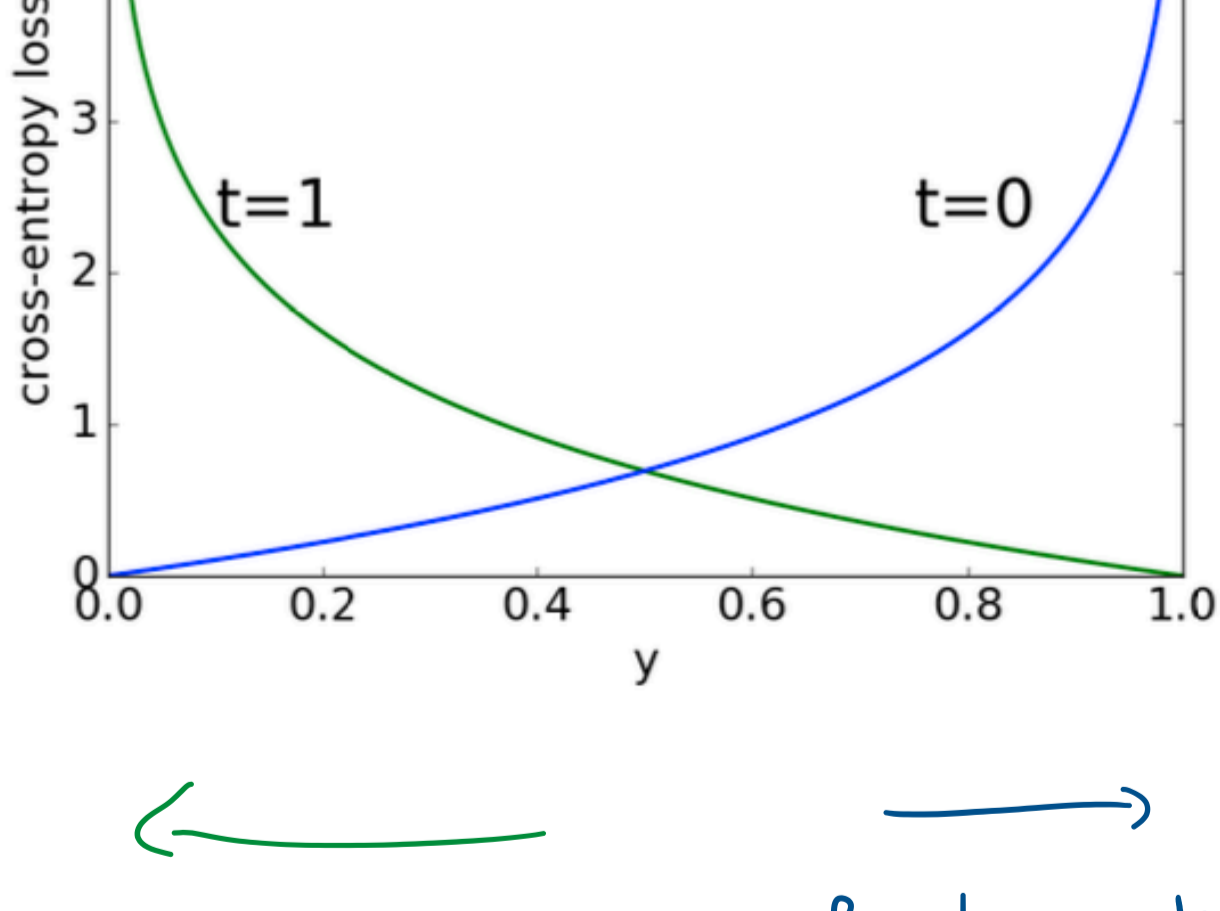
Cross-entropy loss

$$L_{CE}(y, t) = \begin{cases} -\log y & \text{if } t = 1 \\ -\log(1-y) & \text{if } t = 0 \end{cases}$$

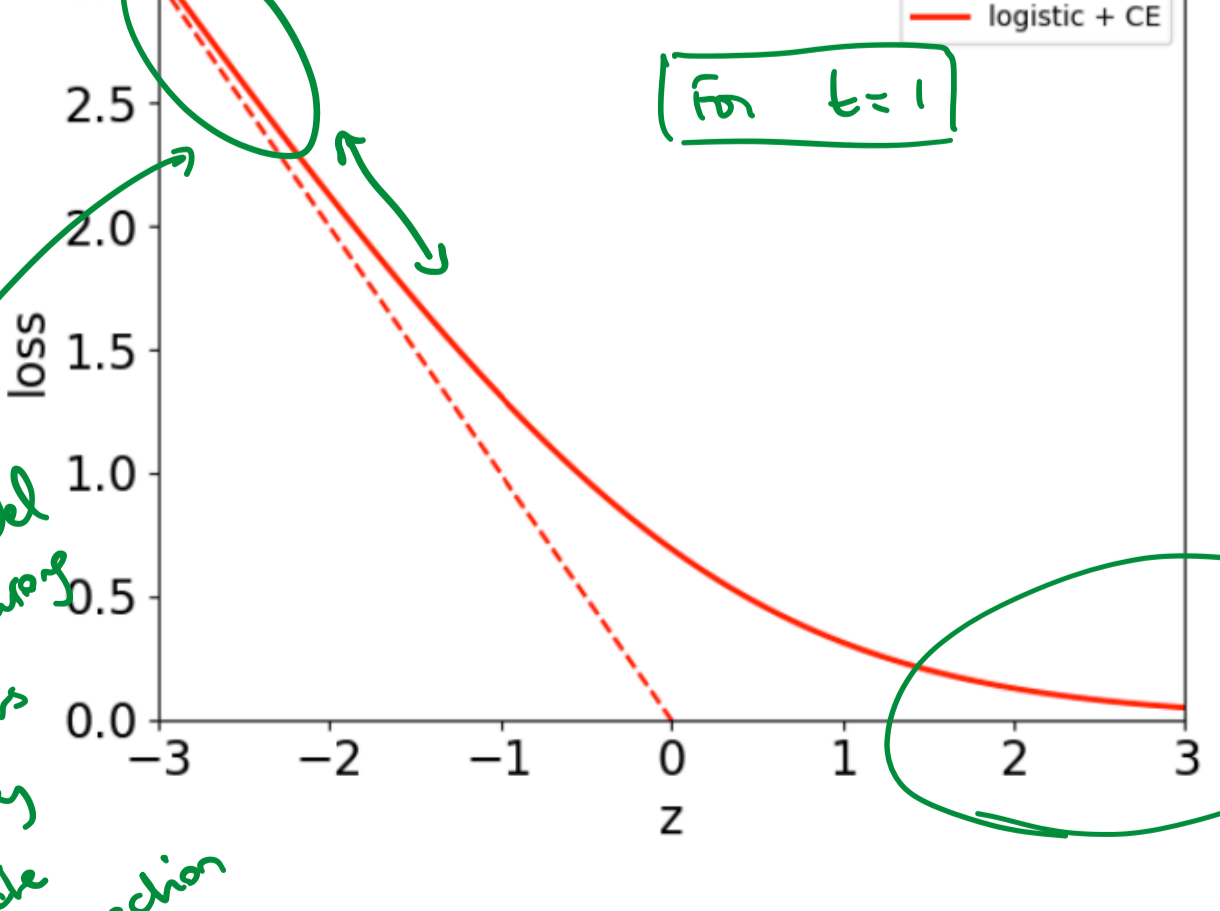
interpreting  $y \in [0, 1]$  (because we applied the logistic function) as a probability.

$$= -t \log y - (1-t) \log(1-y)$$

if  $t=0$  this term cancels out  
if  $t=1$   $1-t=0$  so this term cancels out.



for  $t=1$ , smaller values of  $y$  (up to 0) are wrong  
for  $t=0$ , larger values of  $y$  (up to 1) are wrong, we should instead  $y \rightarrow 0$

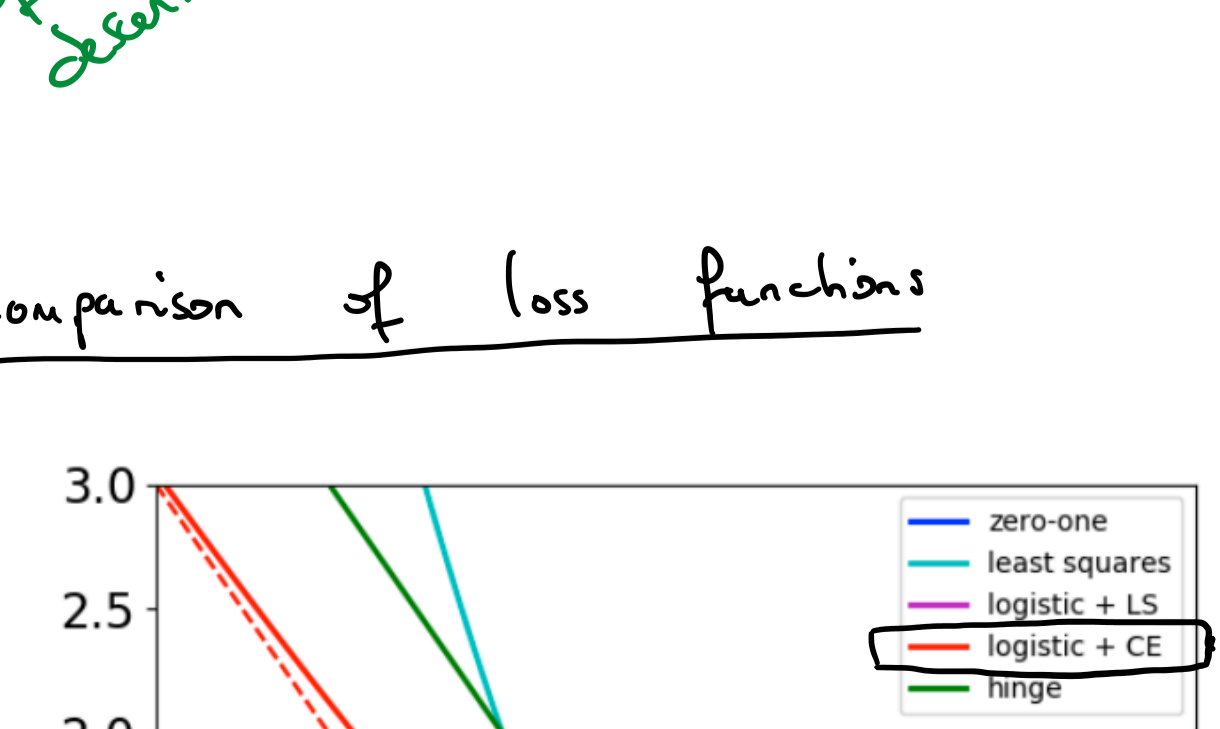


$$z = \vec{w}^T \vec{x} + b$$

$$y = \sigma(z)$$

model is correct loss goes to 0.

Comparison of loss functions



target  $t = 1$   
positive values of  $z$  are correct predictions

Hinge loss:

$$L_H(y, t) = \max(0, 1 - ty)$$