# L2 正则化以及反对抗能力

# L2正则化是什么
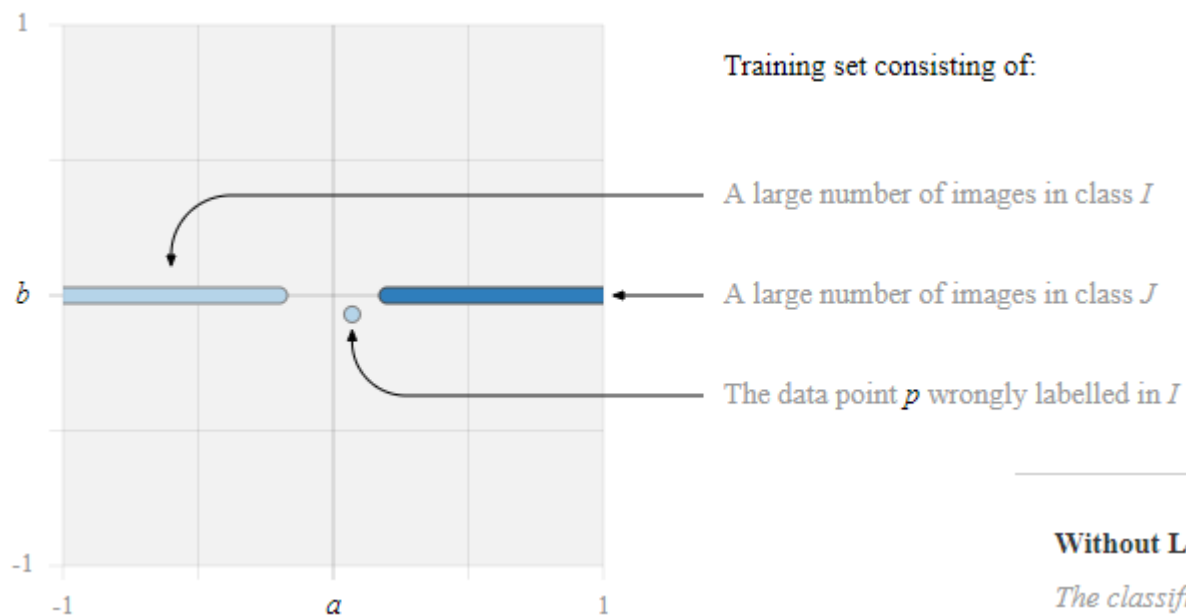
- 在标准损失后面添加penalty
- 作用:
  - 惩罚权重网络，防止过度拟合
  - 还可以做什么?

$$C = C_0 + \frac{\lambda}{2n} \sum_w w^2$$

$$\frac{\partial C}{\partial w} = \frac{\partial C_0}{\partial w} + \frac{\lambda}{n} w$$

$$\frac{\partial C}{\partial b} = \frac{\partial C_0}{\partial b}$$

$$w \rightarrow w - \eta \frac{\partial C_0}{\partial w} - \frac{\eta \lambda}{n} w$$

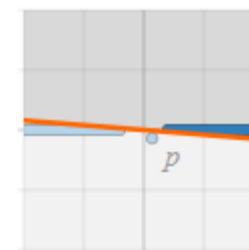$$= \left(1 - \frac{\eta \lambda}{n}\right) w - \eta \frac{\partial C_0}{\partial w}$$

# L2和对抗性，以及推理



Training set consisting of:

A large number of images in class $I$

A large number of images in class $J$

The data point $p$ wrongly labelled in $I$

**Without L2 regularization:**

*The classification boundary is strongly tilted.*

Most of the leeway available to fit the training data resides in the tilting angle of the boundary. Here, the data point $p$ can be classified correctly, but the classifier obtained is then vulnerable to adversarial examples.

**With L2 regularization:**

*The classification boundary is not tilted.*

L2 regularization reduces overfitting by allowing some training samples to be misclassified. When enough regularization is used, the data point $p$ is ignored and the classifier obtained is robust to adversarial examples.

# 几个定义

$ln(z+1)$

$$s(\boldsymbol{x}) := \boldsymbol{w} \cdot \boldsymbol{x} + b$$

$\boldsymbol{x}$ is classified in $\left| \begin{array}{l} I \text{ if } s(\boldsymbol{x}) \leq 0 \\ J \text{ if } s(\boldsymbol{x}) \geq 0 \end{array} \right.$

$$R(\boldsymbol{w}, b) := \frac{1}{n} \sum_{(\boldsymbol{x},y) \in T} f(y\, s(\boldsymbol{x}))$$

$$d(\boldsymbol{x}) := \hat{\boldsymbol{w}} \cdot \boldsymbol{x} + b' \qquad \text{where} \qquad \hat{\boldsymbol{w}} := \frac{\boldsymbol{w}}{\|\boldsymbol{w}\|} \quad b' := \frac{b}{\|\boldsymbol{w}\|}$$

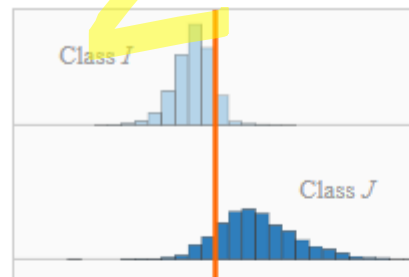$$\text{and} \quad s(\boldsymbol{x}) = \|\boldsymbol{w}\|\, d(\boldsymbol{x})$$

Hence, the norm $\|\boldsymbol{w}\|$ can be interpreted as a scaling parameter for the loss function in the expression of the empirical risk:

$$R(\boldsymbol{w}, b) = \frac{1}{n} \sum_{(\boldsymbol{x},y) \in T} f\big( \underbrace{\|\boldsymbol{w}\| \times y\, d(\boldsymbol{x})}_{\text{scaling parameter for } f} \big)$$

Let us define the *scaled loss function* $f_{\|\boldsymbol{w}\|} : z \to f(\|\boldsymbol{w}\| \times z)$.

$s(\boldsymbol{x})$

Class $I$

Class $J$
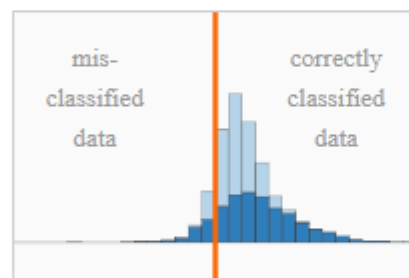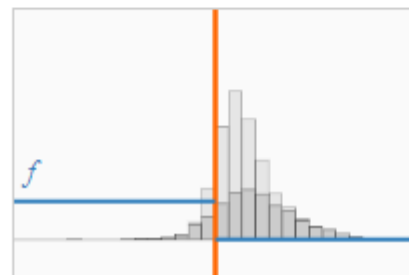
If we plot the histograms of the raw scores over the training set, we typically get two clusters of opposite signs

$y\, s(\boldsymbol{x})$

mis-classified data

correctly classified data

Multiplying by the label allows us to distinguish the correctly classified data from the misclassified data

$f(y\, s(\boldsymbol{x}))$

$f$

We can then attribute a penalty to each training point $\boldsymbol{x}$ by applying a *loss function* to $y\, s(\boldsymbol{x})$

$$d(\boldsymbol{x}) := \hat{\boldsymbol{w}} \cdot \boldsymbol{x} + b' \qquad \text{where} \qquad \hat{\boldsymbol{w}} := \frac{\boldsymbol{w}}{\|\boldsymbol{w}\|} \qquad b' := \frac{b}{\|\boldsymbol{w}\|}$$

$$\text{and} \quad s(\boldsymbol{x}) = \|\boldsymbol{w}\| \, d(\boldsymbol{x})$$

Hence, the norm $\|\boldsymbol{w}\|$ can be interpreted as a scaling parameter for the loss function in the expression of the empirical risk:

$$R(\boldsymbol{w}, b) = \frac{1}{n} \sum_{(\boldsymbol{x},y) \in T} f\big( \underbrace{\|\boldsymbol{w}\| \times y \, d(\boldsymbol{x})}_{\text{scaling parameter for } f} \big)$$

Let us define the *scaled loss function* $f_{\|\boldsymbol{w}\|} : z \to f(\|\boldsymbol{w}\| \times z)$.

# WHEN $\|$ W $\|$ IS LARGE

$$f_{\|\boldsymbol{w}\|} \big( y \, d(\boldsymbol{x}) \big) \underset{\|\boldsymbol{w}\| \to +\infty}{\approx} \|\boldsymbol{w}\| \, \max \left( -y \, d(\boldsymbol{x}), 0 \right)$$

enience, we name the set of misclassified data:

$$M := \{ (\boldsymbol{x}, y) \in T \mid y \, d(\boldsymbol{x}) \le 0 \}$$

in then write the empirical risk as:

$$R(\boldsymbol{w}, b) \underset{\|\boldsymbol{w}\| \to +\infty}{\approx} \|\boldsymbol{w}\| \left( \frac{1}{n} \sum_{(\boldsymbol{x},y) \in M} \big( -y \, d(\boldsymbol{x}) \big) \right)$$

ession contains a term which we call the *error distance*:

$$d_{\mathrm{err}} := \frac{1}{n} \sum_{(\boldsymbol{x},y) \in M} \big( -y \, d(\boldsymbol{x}) \big)$$

*ive* and can be interpreted as the average distance by which ning sample is misclassified by $\mathcal{C}$ (with a null contribution for ctly classified data). It is related—although not exactly t—to the training error.[3]

e have:

$$\text{minimize: } R(\boldsymbol{w}, b) \underset{\|\boldsymbol{w}\| \to +\infty}{\Longleftrightarrow} \text{minimize: } d_{\mathrm{err}}$$

## softplus loss

$$d(\mathbf{x}) := \hat{\mathbf{w}} \cdot \mathbf{x} + b' \qquad \text{where} \qquad \hat{\mathbf{w}} := \frac{\mathbf{w}}{\|\mathbf{w}\|} \qquad b' := \frac{b}{\|\mathbf{w}\|}$$
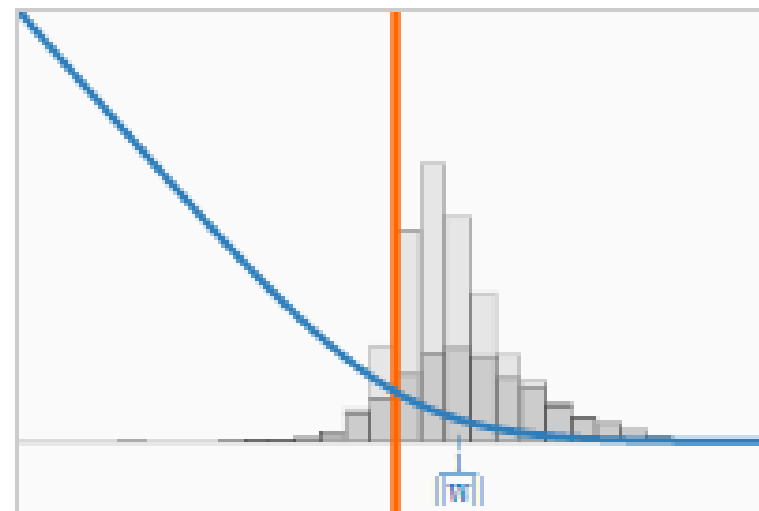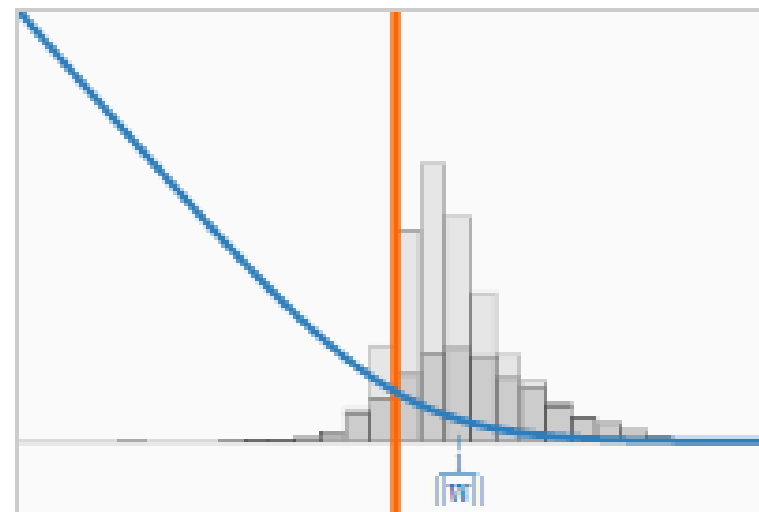
$$\text{and} \quad s(\mathbf{x}) = \|\mathbf{w}\| \, d(\mathbf{x})$$

Hence, the norm $\|\mathbf{w}\|$ can be interpreted as a scaling parameter for the loss function in the expression of the empirical risk:

$$R(\mathbf{w}, b) = \frac{1}{n} \sum_{(\mathbf{x}, y) \in T} f\big( \underbrace{\|\mathbf{w}\| \times y \, d(\mathbf{x})}_{\text{scaling parameter for } f} \big)$$

Let us define the *scaled loss function* $f_{\|\mathbf{w}\|} : z \to f(\|\mathbf{w}\| \times z)$.

More precisely, both losses satisfy:[4]

$$f_{\|\mathbf{w}\|}\left( y \, d(\mathbf{x}) \right) \underset{\|\mathbf{w}\| \to 0}{\approx} \alpha - \beta \, \|\mathbf{w}\| \, y \, d(\mathbf{x})$$

for some positive values $\alpha$ and $\beta$.

We can then write the empirical risk as:

$$R(\mathbf{w}, b) \underset{\|\mathbf{w}\| \to 0}{\approx} \alpha - \beta \, \|\mathbf{w}\| \left( \frac{1}{n} \sum_{(\mathbf{x}, y) \in T} y \, d(\mathbf{x}) \right)$$

This expression contains a term which we call the *adversarial distance*:

$$d_{\mathrm{adv}} := \frac{1}{n} \sum_{(\mathbf{x}, y) \in T} y \, d(\mathbf{x})$$

It is the mean distance between the images in $T$ and the classification boundary $\mathcal{C}$ (with a negative contribution for the misclassified images). It can be viewed as a measure of robustness to adversarial perturbations: when $d_{\mathrm{adv}}$ is high, the number of misclassified images is limited and the correctly classified images are far from $\mathcal{C}$.

## softplus loss



Finally we have:

$$\text{minimize: } R(\mathbf{w}, b) \underset{\|\mathbf{w}\| \to 0}{\Longleftrightarrow} \text{maximize: } d_{\mathrm{adv}}$$

- λ值越大，||W||越小，分类效果越差
  - 同时惩罚了正确分类，但是保证了正确分类离分类边界更远
- λ值越小，||W||越大，分类效果越好
  - 仅仅只惩罚错误分类

$$L(\boldsymbol{w}, b) := \underbrace{R(\boldsymbol{w}, b)}_{\text{empirical risk}} + \underbrace{\lambda \|\boldsymbol{w}\|^2}_{\text{L2 regularization}}$$

# APPENDIX

- https://thomas-tanay.github.io/post--L2-regularization/