# loss functions and optimization

① A loss functions tells us how good our current classifier is. Given a dataset of examples

$$\{x_i, y_i\}_{i=1}^{N}$$

image (pixels) ⟶ label (category 1-10)

loss over the dataset is given by

$$L = \frac{1}{N} \sum_i L_i\left(f(x_i, w), y_i\right)$$

predicted value (label) ⟶ (under $f(x_i,w)$)
True value (label) ⟶ ($y_i$)

② multi-class svm loss (support vector machines). usually if we have 2 classes. we could've used binary loss, but since here we have 10 classes, we will use multi-class svm loss.
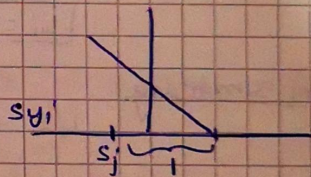
given an example $(x_i, y_i)$ and using the shorthand for the score vector
$$s = f(x_i, w)$$
the svm loss has the form.

correct score (true class)

"Hinge loss"

$$L_i = \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_{y_i} \geq s_j + 1 \\ (s_j - s_{y_i} + 1) & \text{otherwise} \end{cases}$$

⟶ incorrect score

$$= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

$s_{y_i}$ ⟶ $s_j$

Eg.

|       |      |      |      |
|-------|------|------|------|
| Cat   | 3.2  | 1.3  | 2.2  |
| Car   | 5.1  | 4.9  | 2.5  |
| frog  | -1.7 | 2.0  | -3.1 |
| loss  | = 2.9 | 0   | 12.9 |

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$
$$= \max(0, 5.1 - 3.2 + 1) +$$
$$\max(0, -1.7 - 3.2 + 1)$$
$$= \max(0, 2.9) + \max(0, -3.9)$$
$$= 2.9 + 0$$
$$= 2.9$$

⟶ Kind of a quantitative measure of how much our classifier screwed up in this training example.

loss over entire dataset

$$L = \frac{1}{N} \sum_{i=1}^{N} L_i$$
$$= (2.9 + 0 + 12.9)/3$$
$$= 5.27$$

③ choosing w

hyperparameter

$$L(w) = \frac{1}{N} \sum_{i=1}^{N} L_i\left(f(x_i, w), y_i\right) + \lambda R(w)$$

Data loss: model predictions should match the training data

Regularization loss: model should be simple. so that it works on test data.

Generally, we have the following Regularization types

① L1 Regularization $R(w) = \sum_k \sum_l |W_{k,l}|$

② L2 " " $R(w) = \sum_k \sum_l W_{k,l}^2$ ← General.

③ Elastic net (L1 + L2) $R(w) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$

④ max norm regularization.

⑤ Dropout

⑥ Fancier:

④ **Softmax loss** (multinomial logistic Regression)

scores = unnormalised log probabilities of the classes

$$P(Y=k \mid X=x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}} \qquad \boxed{s = f(x_i; w)}$$

softmax
function

we want to maximise the log likelihood, or (for a loss function) to minimize the negative log likelihood of the correct class.

$$\boxed{L_i = -\log \{P(Y=y_i \mid X=x_i)\}}$$

**in summary:**

$$L_i = -\log \left( \frac{e^{s_{y_i}}}{\sum_j e^{s_j}} \right)$$

Eg.
$$\left.\begin{array}{ll} \text{Cat} & 3.2 \\ \text{Car} & 5.1 \\ \text{Frog} & -1.7 \end{array}\right\} \xrightarrow{exp} \left.\begin{array}{l} 24.5 \\ 164.0 \\ 0.18 \end{array}\right\} \xrightarrow{normalise} \left.\begin{array}{l} 0.13 \\ 0.87 \\ 0.00 \end{array}\right\} \begin{array}{l} L_i = -\log(0.13) \\ \\ = 0.89 \end{array}$$

unnormalised          Unnormalised          Probabilities.
log probabilities      probabilities

---

**Recap.**

① we have some dataset of $(x_i, y_i)$
② we have a score function       ③ we have a loss function.

full loss → $$\boxed{L = \frac{1}{N} \sum_{i=1}^{N} L_i + R(w)}$$

$$L_i = -\log \left( \frac{e^{s_{y_i}}}{\sum_j e^{s_j}} \right) \qquad\qquad L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

(softmax)                              (svm)

Q/ How do we find "w" which minimises the loss?
→ optimization.

→ following the slope #.

slope → In 1-d, the derivative of a function
$$\frac{d f(x)}{dx} = \lim_{h \to 0} \frac{f(x+h) - f(x)}{h}$$

→ In multiple dimensions, the gradient is the vector of (Partial-derivatives) along each dimension.

→ The slope in any direction is the <u>dot product</u> of the direction with the gradient. The direction of the steepest descent is the negative gradient.

② Numerical gradient → approximate, slow, easy to write.
   Analytical gradient → Exact, fast, error-prone.

<u>Hint</u> → Always use analytical gradient, but check implementation with numerical gradient. This is called "gradient check".

③ <u>Gradient Descent</u> → algorithm where we use the gradient at every timestep to determine where to step next.

<u>Stochastic gradient Descent</u> → when calculating the loss, or the gradient descent, the full sum is expensive when N becomes very large, like:

$$L(w) = \frac{1}{N} \sum_{i=1}^{N} L_i(x_i, y_i, w) + \lambda R(w)$$

instead, we approximate the sum using "mini batch" of examples 32/64/12

$$\nabla_w L(w) = \frac{1}{N} \sum_{i=1}^{N} \nabla_w L_i(x_i, y_i, w) + \lambda \nabla_w R(w)$$