



# COMP9444: Neural Networks and Deep Learning

## Week 3a. Backprop Variations

Alan Blair

School of Computer Science and Engineering

September 24, 2024

# Outline

- Cross Entropy
- Maximum Likelihood
- Softmax
- Weight Decay
- Bayesian Inference and MAP Estimation
- Second Order Methods

# Cross Entropy

For *function approximation*, we normally use the *sum squared error* (SSE) loss:

$$E = \frac{1}{2} \sum_i (z_i - t_i)^2$$

where  $z_i$  is the output of the network, and  $t_i$  is the target output.

However, for *classification* tasks, where the target output  $t_i$  is either 0 or 1, it is more logical to use the *cross entropy* loss:

$$E = \sum_i (-t_i \log(z_i) - (1 - t_i) \log(1 - z_i))$$

The motivation for these loss functions can be explained using the mathematical concept of *Maximum Likelihood*.

# Maximum Likelihood

Let  $H$  be a class of *hypotheses* for predicting observed *data*  $D$ .

$\text{Prob}(D | h)$  = probability of data  $D$  being generated under hypothesis  $h \in H$ .

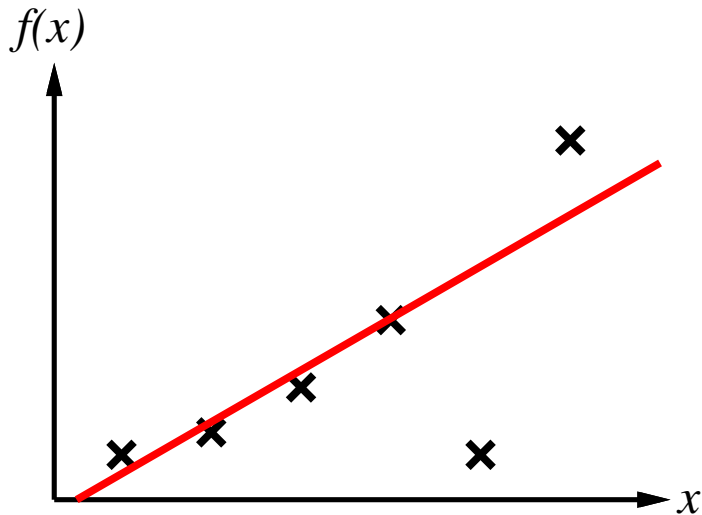
$\log \text{Prob}(D | h)$  is called the *likelihood* of  $D$ , given  $h$ .

ML Principle: Choose  $h \in H$  which maximizes this likelihood,

i.e. maximize  $\text{Prob}(D | h)$  [or, maximize  $\log \text{Prob}(D | h)$ ]

Here, the data  $D$  are the target values  $\{t_i\}$  corresponding to input features  $\{x_i\}$ , and each hypothesis  $h$  is a function  $f()$  determined by a neural network with specified weights or, to give a simpler example,  $f()$  could be a straight line with a specified slope and  $y$ -intercept.

# Least Squares Fit



# Derivation of Least Squares

Due to the Central Limit Theorem, an accumulation of small errors will tend to produce “noise” in the form of a Gaussian distribution.

Suppose the data are generated by a linear function  $f()$  plus Gaussian noise with mean zero and standard deviation  $\sigma$ . Then

$$\begin{aligned}\text{Prob}(D|h) = \text{Prob}(\{t_i\} | f) &= \prod_i \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(t_i - f(x_i))^2} \\ \log \text{Prob}(\{t_i\} | f) &= \sum_i \left( -\frac{1}{2\sigma^2} (t_i - f(x_i))^2 - \log(\sigma) - \frac{1}{2} \log(2\pi) \right) \\ f_{\text{ML}} &= \operatorname{argmax}_{f \in H} \log \text{Prob}(\{t_i\} | f) \\ &= \operatorname{argmin}_{f \in H} \sum_i (t_i - f(x_i))^2\end{aligned}$$

(Note: we do not need to know  $\sigma$ )

# Derivation of Cross Entropy

For binary classification tasks, the target value  $t_i$  is either 0 or 1.

It makes sense to interpret the output  $f(x_i)$  of the neural network as the *probability* of the true value being 1, i.e.

$$\begin{aligned}P(1 | f(x_i)) &= f(x_i) \\P(0 | f(x_i)) &= (1 - f(x_i)) \\ \text{i.e. } P(t_i | f(x_i)) &= f(x_i)^{t_i} (1 - f(x_i))^{(1-t_i)}\end{aligned}$$

$$-\log P(\{t_i\} | f) = \sum_i (-t_i \log f(x_i) - (1 - t_i) \log(1 - f(x_i)))$$

$$f_{\text{ML}} = \operatorname{argmin}_{f \in H} \sum_i (-t_i \log f(x_i) - (1 - t_i) \log(1 - f(x_i)))$$

(Can also be generalized to multiple classes.)

# Cross Entropy and Backprop

Cross Entropy loss is often used in combination with sigmoid activation at the output node, which guarantees an output strictly between 0 and 1, and also makes the backprop computations a bit simpler, as follows:

$$E = \sum_i (-t_i \log(z_i) - (1 - t_i) \log(1 - z_i))$$

$$\frac{\partial E}{\partial z} = -\frac{t_i}{z_i} + \frac{1 - t_i}{1 - z_i} = \frac{z_i - t_i}{z_i(1 - z_i)}$$

$$\text{If } z = \frac{1}{1 + e^{-s}}, \quad \frac{\partial E}{\partial s_i} = \frac{\partial E}{\partial z_i} \frac{\partial z_i}{\partial s_i} = z_i - t_i$$



# Cross Entropy and KL-Divergence

If we consider  $p_i = \langle 1 - t_i, t_i \rangle$ ,  $q_i = \langle 1 - f(x_i), f(x_i) \rangle$  as discrete probability distributions, the Cross Entropy loss can be written as:

$$\begin{aligned} -\log P(\{t_i\} | f) &= \sum_i (-t_i \log f(x_i) - (1 - t_i) \log (1 - f(x_i))) \\ &= \sum_i [(t_i (\log(t_i) - \log f(x_i)) + (1 - t_i)(\log(1 - t_i) - \log(1 - f(x_i))) \\ &\quad + (-t_i \log(t_i) - (1 - t_i) \log(1 - t_i))] \\ &= \sum_i [D_{\text{KL}}(p_i \| q_i) + H(p_i)] \end{aligned}$$

Since  $H(p_i)$  is fixed, minimizing the Cross Entropy loss is the same as minimizing  $\sum_i D_{\text{KL}}(p_i \| q_i)$ .

# Cross Entropy and Outliers

SSE and Cross Entropy behave a bit differently when it comes to outliers.

SSE is more likely to misclassify outliers, because the loss function for each item is bounded between 0 and 1.

Cross Entropy is more likely to keep outliers correctly classified, because the loss function grows logarithmically (unbounded) as the difference between the target and network output approaches 1.

For this reason, Cross Entropy works particularly well for classification tasks that are *unbalanced* in terms of negative items vastly outnumbering positive ones (or vice versa).

# Softmax

- classification task with  $N$  classes
- neural network with  $N$  outputs  $z_1, \dots, z_N$
- assume the network's estimate for the probability of the correct class being  $j$  is proportional to  $\exp(z_j)$
- because the probabilities must add up to 1, we need to *normalize* by dividing by their sum:

$$\begin{aligned}\text{Prob}(i) &= \frac{\exp(z_i)}{\sum_{j=1}^N \exp(z_j)} \\ \log \text{Prob}(i) &= z_i - \log \sum_{j=1}^N \exp(z_j)\end{aligned}$$

# Log Softmax and Backprop

If the correct class is  $k$ , we can treat  $-\log \text{Prob}(k)$  as our cost function, and the gradient is

$$\frac{d}{dz_i} \log \text{Prob}(k) = \delta_{ik} - \frac{\exp(z_i)}{\sum_{j=1}^N \exp(z_j)} = \delta_{ik} - \text{Prob}(i),$$

where  $\delta_{ik}$  is the *Kronecker delta*.

This gradient pushes up the correct class  $i = k$  in proportion to the difference between its assigned probability and 1, and it pushes down the incorrect classes  $i \neq k$  in proportion to the probabilities assigned to them by the network.

# Softmax, Boltzmann and Sigmoid

If you have studied mathematics or physics, you may be interested to know that Softmax is related to the Boltzmann Distribution, with the negative of output  $z_i$  playing the role of the “energy” for “state”  $i$ .

The Sigmoid function can also be seen as a special case of Softmax, with two classes and one output, as follows:

Consider a simplified case where there is a choice between two classes, Class 0 and Class 1. We consider the output  $z$  of the network to be associated with Class 1 and we imagine a fixed “output” for Class 0 which is always equal to zero. In this case, the Softmax becomes:

$$\text{Prob}(1) = \frac{e^z}{e^z + e^0} = \frac{1}{1 + e^{-z}}$$

# Weight Decay

Sometimes we add a penalty term to the loss function which encourages the neural network weights  $w_j$  to remain small:

$$E = \frac{1}{2} \sum_i (z_i - t_i)^2 + \frac{\lambda}{2} \sum_j w_j^2$$

This can prevent the weights from “saturating” to very high values.

It is sometimes referred to as “elastic weights” because the weights experience a force as if there were a spring pulling them back towards the origin according to Hooke’s Law.

The scaling factor  $\lambda$  needs to be determined from experience, or empirically.

# Bayesian Inference

$H$  is a class of hypotheses.

$\text{Prob}(D | h)$  = probability of data  $D$  being generated under hypothesis  $h \in H$ .

$\text{Prob}(h | D)$  = probability that  $h$  is correct, given that data  $D$  were observed.

Bayes' Theorem:

$$\begin{aligned}\text{Prob}(h | D)\text{Prob}(D) &= \text{Prob}(D | h) \text{Prob}(h) \\ \text{Prob}(h | D) &= \frac{\text{Prob}(D | h) \text{Prob}(h)}{\text{Prob}(D)}\end{aligned}$$

$\text{Prob}(h)$  is called the *prior* because it is our estimate of the probability of  $h$  *before* the data have been observed.

$\text{Prob}(h | D)$  is called the *posterior* because it is our estimate of the probability of  $h$  *after* the data have been observed.

# Weight Decay as MAP Estimation

We assume a Gaussian prior distribution for the weights, i.e.

$$P(w) = \prod_j \frac{1}{\sqrt{2\pi}\sigma_0} e^{-w_j^2/2\sigma_0^2}$$

Then

$$P(w | t) = \frac{P(t | w)P(w)}{P(t)} = \frac{1}{P(t)} \prod_i \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(z_i - t_i)^2} \prod_j \frac{1}{\sqrt{2\pi}\sigma_0} e^{-w_j^2/2\sigma_0^2}$$

$$\log P(w | t) = -\frac{1}{2\sigma^2} \sum_i (z_i - t_i)^2 - \frac{1}{2\sigma_0^2} \sum_j w_j^2 + \text{constant}$$

$$\begin{aligned} w_{\text{MAP}} &= \operatorname{argmax}_{w \in H} \log P(w | t) \\ &= \operatorname{argmin}_{w \in H} \left( \frac{1}{2} \sum_i (z_i - t_i)^2 + \frac{\lambda}{2} \sum_j w_j^2 \right), \text{ where } \lambda = \sigma^2/\sigma_0^2 \end{aligned}$$

This is known as Maximum A Posteriori (MAP) estimation.



## Second Order Methods

Some optimization methods involve computing *second order* partial derivatives of the loss function with respect to each *pair* of weights:

$$\frac{\partial^2 E}{\partial w_i \partial w_j}$$

- Conjugate Gradients
  - approximate the landscape with a quadratic function (paraboloid) and jump to the minimum of this quadratic function
- Natural Gradients (Amari, 1995)
  - use methods from information geometry to find a “natural” re-scaling of the partial derivatives

These methods are not normally used for deep learning, because the number of weights is too high. In practice, the Adam optimizer tends to provide similar benefits with low computational cost.