

10) Metaparameters

Metaparam vs. param

- **Parameter**: feature learned by the model algorithm (ex. weight, bias)
- **meta-parameter**: feature u set for the model (ex. architecture, loss function)

Data Normalization

- **Purpose**: Ensure weights remain numerically stable: no gradient issues
(data samples are in the same range w/o changing their relationship)

① **Z-scoring**: $z_i = \frac{x_i - \bar{x}}{\sigma_x}$ such that $\bar{z}=0, \sigma_z=1$ used for gaussian data

② **Min-max scaling**: $\tilde{x} = \frac{x - \min x}{\max x - \min x} \Rightarrow x^* = a + \tilde{x}(b-a)$ used for bounded data
scale to $[0,1]$ scale to $[a,b]$ (optional)

Batch Normalization

- Not normalize inputs x_i but inputs to each layer this time

• **Process**: $y = \sigma(\tilde{x}^T w)$, where $\tilde{x} = \underbrace{r x + b}_{\text{model param that scale stdev \& shift mean}}$
 \uparrow input to next layer \uparrow normalized input

- Should be turned off during testing since batch_size differs

Activation Functions

- Add nonlinearities into layers: they won't collapse into 1 layer

• **Common ones**: ① Sigmoid $\sigma(x) = \frac{1}{1+e^{-x}}$ ② tanh $tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$ ③ ReLU $ReLU(x) = \max(0, x)$ Good with normalized inputs

• **ReLU variants**: ① Leaky ReLU $LeakyReLU(x) = \begin{cases} x & x > 0 \\ \alpha x & x \leq 0 \end{cases}$ usually $\alpha=0.01 \rightarrow \alpha x$ ② ReLU-N $ReLU_N(x) = \begin{cases} x & x > N \\ \min(\max(0, x), N) & otherwise \end{cases}$ usually $N=6$

Loss Functions

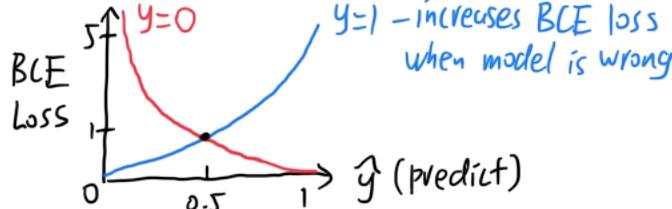
- Metric to backpropagate (adjust weights)

more general (BCE is when $N=1, C=2$)

• **Common ones**: ① **BCE** $L = -(y \log(\hat{y}) + (1-y) \log(1-\hat{y}))$ ② **CCE** $L = -\sum_{i=1}^N \sum_{k=1}^C y_i^{(k)} \log(\hat{y}_i^{(k)})$

data sample $\rightarrow N$ # category

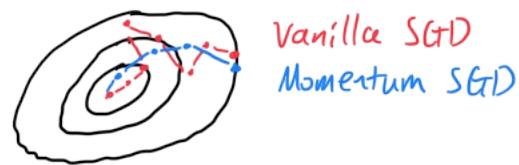
③ **MLE** $L = \frac{1}{N} \sum (\hat{y} - y)^2$



- Softmax variant: log-softmax $S(x) = \log(e^x / \sum e^{x_i})$ is good for CCE
penalizes errors more sensitively

Optimizers

- algorithms that adjust weights during backprop
- ① Vanilla SGD: change weights after each sample \rightarrow outliers point leads to (used when all samples are similar) \rightarrow **Volatile changes**
- ② Mini-batch SGD: change weights after N samples, avg losses across them
- ③ Momentum SGD: change weights by a weighted avg of previous & current costs
- Formula: $V_t = (1-\beta)dJ + \beta V_{t-1}$, $W := W - \alpha V_t$
 $\beta \in [0, 1]$ where $\beta=0$ is vanilla SGD
- ④ RMSprop: vanilla SGD that updates lr by **total energy/rms(gradient)**
- Formula: $\text{rms} = \sqrt{\frac{1}{m} \sum_{i=1}^m X_i^2}$, $S_t = (1-\beta)(dJ)^2 + \beta S_{t-1}$, $W := W - \frac{\alpha}{\sqrt{S_t + \epsilon}} dJ$ $\downarrow E^{-8}$ (ensure denoising)
- Interpretation: Large gradients = volatile changes \rightarrow slows down training
Small gradients = vanishing changes \rightarrow takes larger step
- ⑤ Adam (adaptive momentum): RMSprop + Momentum
- Formula: $V = (1-\beta_1)dJ + \beta_1 V_{t-1}$, $S = (1-\beta_2)(dJ)^2 + \beta_2 S_{t-1}$, $W := W - \frac{\alpha}{\sqrt{S + \epsilon}} V$
usually: $\beta_1 = 0.9$ $\beta_2 = 0.999$ $V = \frac{s}{1-\beta_1^t}$ $S = \frac{1}{1-\beta_2^t}$



Schedulers

- Algorithms that decay lr over time

Ex) torch.optim.lr_scheduler.StepLR (optimizer, step_size, gamma)
interval to decay applied factor when decay