



## Heuristics for faster training

disadvantage of NNs: take much longer to train than other methods  
\* for any of these solutions, cross validation/empirical analysis required

Solutions:

1) fix vanishing grad. problem

2) SGD: faster than BGD on large, redundant datasets

- BGD: Walks downhill on 1 cost fn

- SGD: Very Short Step downhill on one point's lost fn. and then so on...

(cost fn: Sum of loss functions over all sample pts.

↳ One batch step  $\approx$  n Stochastic Steps, roughly same amount of time

disadvantage of SGD: quick to learn redundant info, e.g. a ton of 9's in MNIST

i.e. better than  
BGD

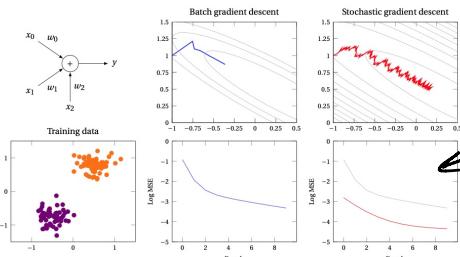


Figure 19.1: Comparison between stochastic gradient descent and batched gradient descent for a simple one-layer neural network

- one epoch processes every training pt. once

- training usually takes multiple epochs

- if large Sample Size (w/ high redundancy), then training can take less than one epoch

Sol. 3) **normalizing data:** Center each feature so  $M=0$ , Scale so  $\sigma=1$

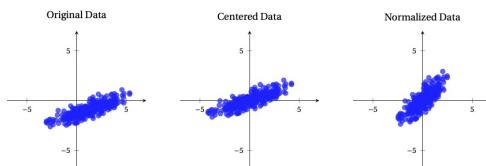


Figure 19.2: Example of centering and normalizing data

empirically, centering data  $\rightarrow$  easier for hidden units to get in good operating region of Sigmoid & relu

Scaling features  $\rightarrow$  better conditioned objective fn.  $\rightarrow$  gradient descent converges faster

- Skewed data  $\rightarrow$  Objective fn w/ ill-conditioned/highly eccentric hessian  $\Rightarrow$  Slow GD convergence

Sol: Scaling or momentum (soon)

4) "centering" the hidden units

$$\hookrightarrow \text{replace } S(\cdot) \text{ w/ } \tanh r = \frac{e^r - e^{-r}}{e^r + e^{-r}} = 2S(2r) - 1$$

$\hookrightarrow$  now, range of  $(-1, 1)$  instead of  $(\infty, \infty)$

- new derivative:  $S' = 1 - \tanh^2 r$

- good output target values:  $(-0.7, 0.7)$

5) different learning rates for each layer of weights

- earlier layers have smaller gradients  $\rightarrow$  larger learning rates

6) Emphasizing Schemes: have network learn/emphasize uncommon examples before the redundant & common examples

- present examples from rare classes more often or w/ bigger LR  
- can do same for misclassified examples

Warning: backfires if there are bad outliers

7) 2nd-order optimization

- Newton's method is impractical since Hessian too large & expensive to compute

**Nonlinear Conjugate Gradient Method:** Works well for small networks w/ small datasets using regression

- Should only be used w/ BGD since this is slow w/ redundant data,

**Stochastic Levenberg Marquardt:** approximates a diagonal Hessian

- authors claim 3x faster convergence vs SGD

8) acceleration schemes e.g. RMSprop, Adam, AMSGrad

## Heuristics for avoiding bad local minima

- 1) fix vanishing grad. problem
- 2) SGD: random motion lets you avoid bad/shallow local minima
  - even if at local minima of cost fn., loss fn. of each pt. will try to push you in a different direction
- 3) **Momentum:** GD can be modified to change a velocity  $\vec{v}_w$  slowly

```
1  $\Delta \vec{w} \leftarrow -\epsilon \nabla J(\vec{w})$ 
2 repeat
3    $\vec{w} \leftarrow \vec{w} + \Delta \vec{w}$ 
4    $\Delta \vec{w} \leftarrow -\epsilon \nabla J(\vec{w}) + \beta \Delta \vec{w}$ 
```

- performs well for BGD & SGD
- choosing  $\beta$ : conflicting advice
  - some set  $\beta = 0.9$ , others  $\beta \approx 0$
  - Hinton: start at 0.5 & increase to 0.4 as grad becomes small
  - large  $\beta \rightarrow$  small LR to compensate
    - can still use large LR at first for reasonable initial velocity
- problem w/ momentum: once it finds a minimum, it oscillates around it
  - ↳ that said, more likely to find good minimum in 1st place

$\beta \neq 1$ , how much momentum persists between iterations

## Heuristics to avoid overfitting:

- 1) ensemble of NNs w/ random initial weights & bagging
  - Stop to train
- 2) L2 regularization / weight decay: add  $\frac{1}{2} \|\mathbf{W}\|^2$  to cost fn  
 $b: \sum \frac{\partial J}{\partial w_i}$  has extra  $-2\epsilon_1 w_i$  term
- 3) Dropout "emulates an ensemble in one network"
  - Temporarily disable random subset of units during training, as well as all edges in and out of those units
  - Seems good if  $P(\text{hidden unit}) = 1/2, p(\text{input}) < p(\text{hidden unit})$
  - w/ SGD, frequently change which random subset of units are disabled
  - authors claim better than L2 regularization
  - Some benefits of ensemble, but faster

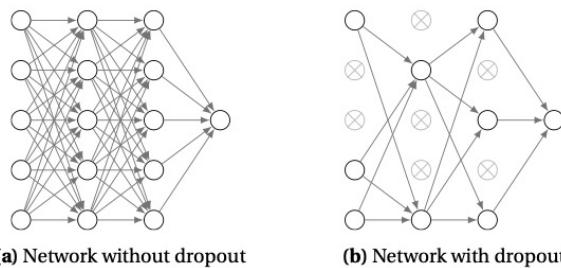


Figure 19.3: Example of a network with and without dropout

## 4) fewer hidden units

- # hidden units is a hyperparameter used to adjust B/V tradeoff
  - too few  $\rightarrow$  high bias
  - too many  $\rightarrow$  low B, high V

## Convolutional Neural Networks

NNs are overparameterized: too many weights, not enough data

- If you have a ton of weights  $\rightarrow$  need more data, slow compute

### 2 New properties of NN design:

- local connectivity: hidden unit (in earlier layers) only connect to small patch of units from the prev. layer

- improves overparametrization problem

- speeds up training & classification

- Shared Weights: groups of hidden units share same mask/filter/kernel

- each mask operates on every patch of an image

- # units in 1st hidden layer = # masks  $\cdot$  # patches

- if NN learns to detect edges on one patch, then every patch has an edge detector

Since the mask that detects edges is applied to every patch

Convolution: same linear transformation applied to different parts of input by shifting

- Shared Weights: suppose gradient descent starts to develop an edge detector

$\hookrightarrow$  this edge detector is being trained on every part of every image, not just one spot

$\hookrightarrow$  good since edges appear on different patches of diff. images

$\hookrightarrow$  loc. no longer matters

NNs: Hidden units are features that we learn

CNNs: learn features from multiple patches simultaneously and then applying those features everywhere

