

→ Log loss → When ~~\hat{y}_i~~ CE is small → high P

Error function: $-\frac{1}{m} \cdot \sum_{i=1}^m (1-y_i) (\ln(1-\hat{y}_i)) + y_i \cdot \ln(\hat{y}_i)$

→ Sigmoid is continuous

- One hot encoding

Multi Class CE: $-\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^n y_{ij} \cdot \ln(\hat{y}_{ij})$

Gradient Descent

$$E(w/b) := -\frac{1}{m} \sum (1-y_i) \ln(1-\sigma(w^T x^{(i)} + b)) + y_i \ln(\sigma(w^T x^{(i)} + b))$$

α has to be reasonable to not end up in local min.

$$\hat{x} = \sigma(w^T x + b) \leftarrow \text{bad}$$

$$\hat{y} = \sigma(w_1 x_1 + \dots + w_n x_n + b)$$

$$\nabla E = (\delta E / \delta w_1, \dots, \delta E / \delta w_n, \delta E / \delta b)$$

$$\alpha = 0.1$$

$$w_i' = w_i - \alpha \frac{\delta E}{\delta w_i}$$

$$b' = b - \alpha \frac{\delta E}{\delta b}$$

$$\hat{y}' = \sigma(w'^T x + b') \leftarrow \text{better}$$

$$\sigma'(x) = \frac{e^{-x}}{(1+e^{-x})^2} = \frac{1}{1+e^{-x}} \cdot \frac{e^{-x}}{1+e^{-x}} = \sigma(x) \cdot (1 - \sigma(x))$$

$$\frac{\delta}{\delta b} E = -(y - \hat{y}) \Rightarrow \text{closer the label to the pred, smaller the gradient}$$

$$\sigma(x) = \frac{1}{1+e^{-x}}$$

$$E = (y - \hat{y})^2 \quad \text{Error for one pred.}$$

(3)

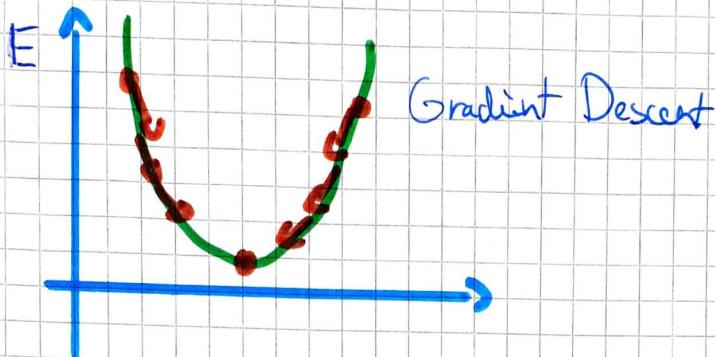
Total Error:

$$E = \frac{1}{2} \sum (y^m - \hat{y}^m)^2$$

\downarrow

$$= f(\sum w_i x_i^m)^2$$

\rightarrow SSE Sum Squared Error



$$w_i = w_i + \Delta w_i$$

$$\Delta w_i \propto -\frac{\delta E}{\delta w_i} \rightarrow \text{The gradient}$$

$$\frac{\delta E}{\delta w_i} = \frac{\delta}{\delta w_i} \frac{1}{2} (y - \hat{y})^2$$

$$= (y - \hat{y}) \cdot \frac{\delta}{\delta w_i} (y - \hat{y}) = - (y - \hat{y}) f'(h) x_i$$

$$\hat{y} = f(h) \quad h = \sum w_i x_i$$

$$\text{error} = (y - \hat{y})$$

$$\delta = (y - \hat{y}) f'(h) \quad \text{error term} = \text{error} \cdot \text{output} \cdot (1 - \text{output})$$

$$\text{new } w_i = w_i + \eta \delta x_i \quad \Delta w = \alpha \cdot \text{error} \cdot x_i$$

Multiple Outputs: Δw : \uparrow delta

$$\delta_j = (y_j - \hat{y}_j) f'(h_j)$$

$$\Delta w_{ij} = \eta \delta_j x_i$$

Derivative of activation
fn.

MSE

$$E = \frac{1}{2m} \sum_{\mu} (y^{\mu} - \hat{y}^{\mu})^2$$

Algo:

$$\text{set } \Delta w_i = 0$$

for each record in training data:

$$\hat{y} = f(\sum_i w_i x_i)$$

$$\delta = (y - \hat{y}) \cdot f'(\sum_i w_i x_i)$$

$$\Delta w_i = \Delta w_i + \delta x_i$$

- update weights $w_i = w_i + \eta \Delta w_i / m$

• repeat for e epochs

→ can be updated each step instead of each epoch

→ recommended to use normal distributed:

$$\text{weights} = \text{np.random.normal scale} = 1/\sqrt{m}, \text{size} = n \times m$$

MLP

$$\text{hidden_layer_in} = \text{np.dot}(X, \text{weights_input_to_hidden})$$

$$\text{hidden_layer_out} = \text{sigmoid}(\text{hidden_layer_in})$$

$$\text{output_layer_in} = \text{np.dot}(\text{hidden_layer_out}, \text{weights_hidden_to_output})$$

$$\text{output_layer_out} = \text{sigmoid}(\text{output_layer_in})$$

Bach prop

5

1. MLP Code

2. error = target - output

$$\text{out_put_error_term} = \text{error} \cdot \text{output} \cdot (1 - \text{output})$$

$$\text{hidden_error_term} = \text{wps_dot} (\text{out_put_error_term}, \\ \text{weights_hidden_output}) \cdot \text{hidden_layer_output} \\ (1 - \text{hidden_layer_output})$$

$$\text{delta_w_h_o} = \alpha \cdot \text{out_error_term} \cdot \text{hidden_layer_output}$$

$$\text{delta_w_i_h} = \alpha \cdot \text{hidden_error_term} \cdot x[:, \text{None}]$$

Complete Bach prop:

loop start

$$\text{hidden_input} = x \cdot \text{del_w_input_hidden}^*$$

$$\text{hidden_output} = o(\text{hidden_input})$$

$$\text{out_put} = o(\text{hidden_output} \cdot w_hidden_output))$$

$$\text{error} = y - \text{out_put}$$

$$\text{out_err_t} = \text{error} \cdot \text{out_put} \cdot (1 - \text{out_put})$$

$$\text{hidden_error} = \text{out_put_error_term} \cdot w_hidden_output$$

$$\text{hidden_error_term} = \text{hidden_error} \cdot \text{hidden_output} \cdot (1 - \text{hidden_output})$$

$$\text{del_w_hidden_output_t} = \text{out_err_t} \cdot \text{hidden_output}$$

$$\text{del_w_hidden_input_t} = \text{hidden_error_term} \cdot x[:, \text{None}]$$

$$\text{weights_input_hidden} = \alpha \cdot \text{del_w_input_hidden} / n_records$$

$$\text{weights_hidden_output} = \alpha \cdot \text{del_w_hidden_output} / n_records$$

(6)

Training

- Train / Test set
- simple model is better!

Overfitting

- Too complex, too specific
- error due to variance
- instead of studying we memorize

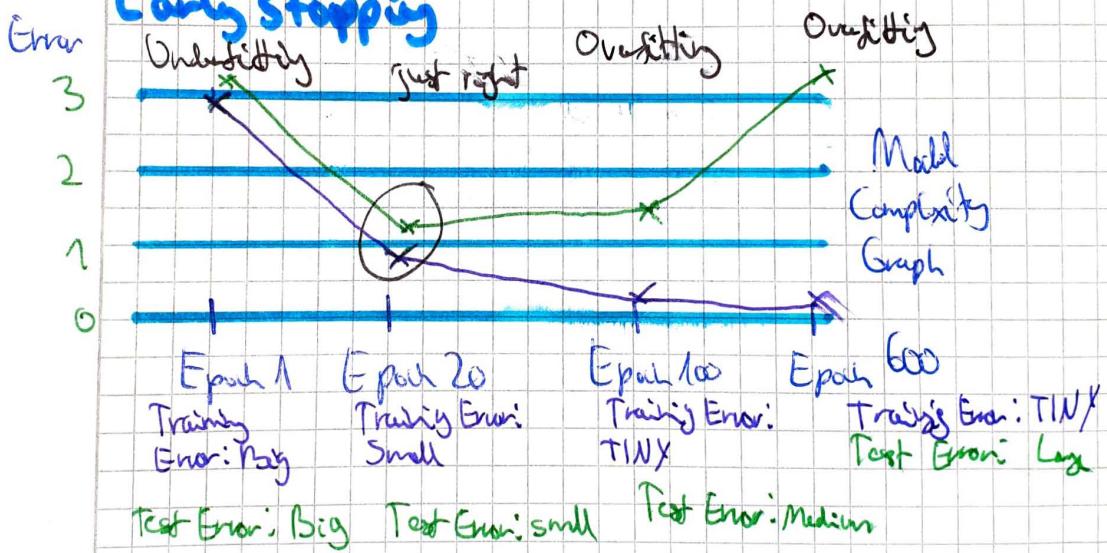
High Variance

Underfitting

- error due bias
- oversimplified

High Bias

Early stopping



"The whole problem of bad model can consist of themselves and good ones are so full of doubts."

Large coefficients → Large Overfitting

Solution: Regularization (L_1 Reg. → good for feature selection)

Add $\lambda (w_1 + \dots + w_n)$) or

$\lambda (w_1^2 + \dots + w_n^2)$ if ~~the~~ training models

weights are large we penalize

Dropout

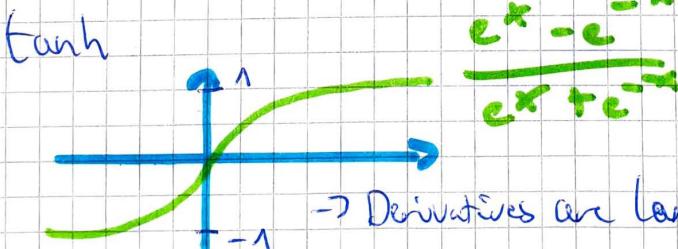
- turn parts off (Nodes) so other Nodes are trained better
- Different Nodes are Dropped
- 0.5 is recommended by researchers, Some nodes may never turn off though which is fine
- best after layers w/ high # of params (not pooling!!)

Gradient Descent

- no local minimum, don't want to get stuck
- Random restart to find global minimum

Vanishing Gradient Problem:

- Very tiny steps especially w/ Sigmoid \rightarrow Backprop increases that
- New activation Function:



→ Derivatives are larger



- Final Output mostly Sigmoid

Batch GD

(8)

- better than stochastic
- Guaran
- Small subsets of data
- Data is split into batches:
Run points from one batch and correct weights
based on error. → Small steps less accurate than
large (stochastic), but almost the

Learning Rate Decay

LR too big:

may miss minimum, fast

LR too small:

Very slow, get to minimum

→ if error too big, decrease LR

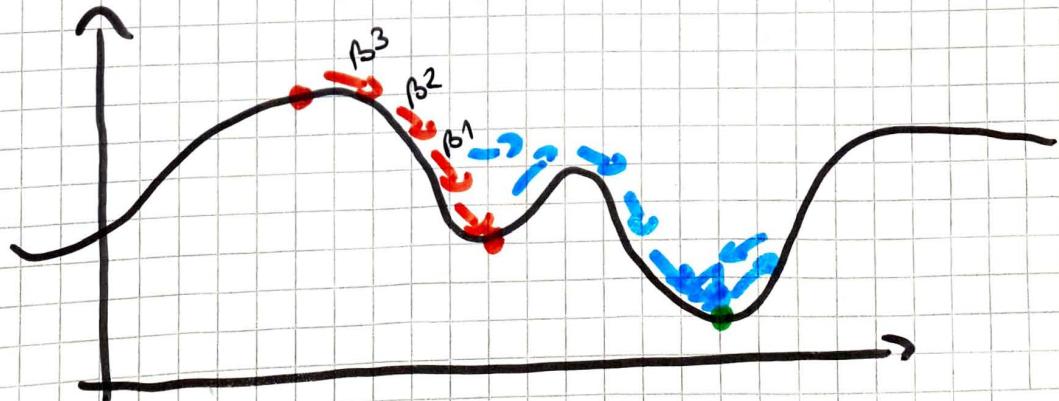
Momentum

→ constant β between 0 and 1

Step → average of previous steps → weighting
last steps more

Step(n) → Step(n) + β Step(n-1)

Error $+ \beta^2 \text{Step}(n-2) + \dots$



g Sentiment Classification

- transform text to numeric problem so we can formulate correlations

1. Creating dataset \rightarrow predictive theory for correlation

2. Validate \rightarrow transform to output

3. iterate several times (increase correlation)

X: what we know
label

Y: What we want to know
Comment

\rightarrow Predictive Theory

\rightarrow Which words are associated w/ sentiment

2) word2index = $\{ \}$

for i , word in enumerate(words):

word2index[word] = i

input layer word count = num of words in review (set)

for word in review.split(" "):

layer_0[0][word2index[word]] += 1

Good Initialization Method

$y = 1/\sqrt{h}$ $h = \text{number of inputs in the layer}$

\hookrightarrow close to 0, but not too small

Neural Noise

⑩

→ Weighting interplays

→ many words are filler words → noise as
these are most frequent

→ Remove noise from input data

→ Binary representation

→ Set value of it to 1 if word exists
↑
 wordIndex[word]

→ Which representations are relevant, which
noise affects it most?

Solve inefficiencies

→ There is overhead

→ Slow to train

→ What's stuff we can shave out?

→ e.g. Layer 0 has about 2000 weights

(\hookrightarrow) Large multiplication w/ 0s doesn't change
anything \Rightarrow shave it out

→ use part of the matrix that are not 0 to
multiply

→ do not multiply by 1, it doesn't
change anything

1) \rightarrow pre process training reviews so we can directly
handle indices w/ non-zero inputs

2) Only use logic that we non zero and update weight
accordingly: $\text{for } i \in \text{reviews}$
 $\text{set weights_v_1[i] = ...}$

Pytorch tools

(11)

Autograd \Rightarrow calculates all the gradients
and keeps track of them

- backward() for Backward pass
through operations that created \mathcal{L}

Optimizer

→ use SGD

→

General Process

~~to import~~ import NN, Optim
import torch nn, functions as F

- make forward pass
 - use output to calculate loss \mathcal{L}
 - perform backward pass with $\mathcal{L}.backward()$
 - take a step w/ the optimizer to update the weights
- Need to optimizer.zero_grad()
when we do multiple passes to zero the gradient after one pass

→ Softmax gives a log probability

→ need an exponential for actual probabilities

13

Inference & Validation

7

- Reduce Overfit by using Dropout
- accuracy measure used $\frac{TP + FP}{\text{all examples}}$
- . view * softmax classes shape to shape of Sh. Shape test
- ps. topk (1, dim = 1) returns most likely class
 (↳ tuple top-K values & indices (class index))
- Save model constantly and use early stopping
- Dropout needs to be off during testing

Save Model

import fc-model → just model

Save state-dict → contains weights, bias, ...

→ torch.save(model.state_dict(), 'checkpoint.pth')

→ torch.load to local state dict

model.load_state_dict() to load dict in to model

However: model needs to be built just like it was when we trained \Rightarrow

checkpoint dict with all info

\rightarrow load checkpoint

E

Load (my) Data

→ datasets. ImageFolder ('path', transform = transforms)

root / class1 / img1

root / class2 / img2

(13)

→ transforms.Compose (

→ DataLoader to generate data in batches

→ iter (DataLoader) :

Augmentation → from torchvision import transforms

→ Rotate, Scale, Change orientations,
brightness

→ transforms.Compose [...]

Transfer learning

→ Pre-trained Networks (my Net in this case)

Import models with torchvision (also datasets)

→ train w/ Cuda (GPU)

Code:

```
device = torch.device ("cuda:0" if  
torch.cuda.is_available() else "cpu")
```

→ freeze params so we don't backprop

→ only train last layers