

Machine Learning Notes

Zach Virgilio

October 31, 2025

To Learn

1. Statistical Basics

- Mean vs (Sample) S.D.
- Expectation vs variance vs Pop. S.D.
- Normal/Multivariate Normal Distributions
- Bernoulli and Binomial Distributions
- Covariance and Correlation (what they measure and what they miss)

2. Estimation and Inference

- Estimators: biased vs unbiased, consistent vs inconsistent
- Variance of estimators (beginning of bias-variance tradeoff)
- Uniform minimum variance unbiased estimators
- maximum likelihood estimators
- maximum a posteriori
- Point vs interval estimation
- Hypothesis testing: p-values, Type I/II errors, power
- Bayesian vs Frequentist thinking

3. Information Theory

- Entropy - uncertainty, information content, optimal encoding
- Cross-Entropy - measuring the distance between distributions
- Binary vs Categorical Cross entropy
- KL divergence (why isn't it symmetric)
- Connection between cross entropy and MLE
- Mutual information (used in feature selection mechanisms)
- Perplexity (how language models measure uncertainty)

4. Models Able to Derive on a Board

- Linear regression, closed form, assumptions, when it breaks
- Logistic regression, why sigmoid, why cross-entropy, connection to odds ratios
- Softmax regression (multinomial logistics) - foundation of modern classifiers
- K-Means (EM algorithm in disguise)
- Expectation maximization (principle behind unsupervised algorithms)
- Naive Bayes (why naive, why it still works)
- Language models: n-grams, smoothing, why transformers are still doing probability

5. Generalizations

- Bias-Variance Tradeoff (explain w/o formula, probably also talk about double-descent)
- Underfitting vs Overfitting - recognize it from loss curves
- Regularization:
 - a) L1 (Lasso) - sparsity and feature selection
 - b) L2 (Ridge) - shrinkage and stability
 - c) Elastic net - when you need both
 - d) Dropout - why random works
 - e) Early stopping - simplest regularizer
 - f) Data augmentation - regularization through diversity
- Cross-validation - k-fold, stratified, time-series splits
- Train/Val/Test splits - why you can't cheat
- Model selection vs model assessment

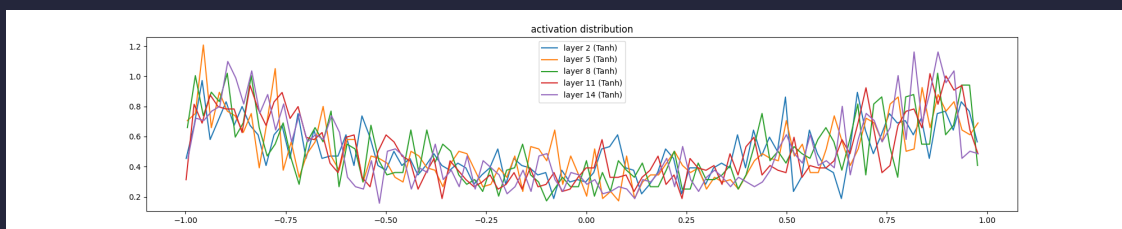
6. Connections that Matter

- Why is MSE loss equivalent to MLE under Gaussian noise?
- Why is cross-entropy loss equivalent to MLE for classification?
- How does L2 regularization relate to MAP with a Gaussian prior?
- How does L1 regularization relate to MAP with a Laplace prior?
- Why does logistic regression use both sigmoid and cross-entropy?
- What's the relationship between PCA and eigenvalues?
- How do attention weights relate to probability distributions?
- Why use log probabilities in practice? (Numerical stability + MLE connection)

Makemore notes

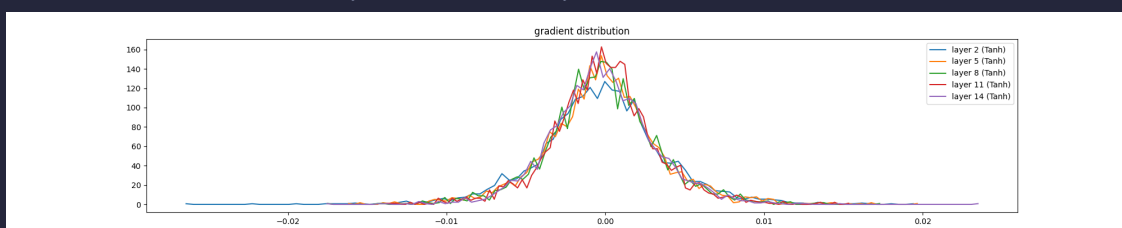
Diagnostic Tools and Visualization

1. Activation saturation of non-linear layers, i.e. tanh or ReLU etc. For tanh layers, want to make sure the distribution is fairly even, and not tailed towards ± 1 like in:



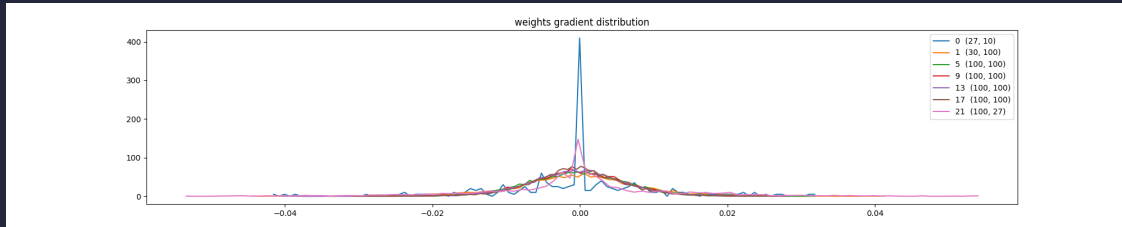
Produced with 1. This is important, because neurons with activations of ± 1 will have gradients of 0, as $\frac{d}{dx} \tanh x = 1 - (\tanh x)^2$, leading to dead neurons which don't learn. The same phenomena occurs for sigmoid non-linearity. These are forward pass statistics

2. Gradients of the non-linear layers. For tanh they can look like:



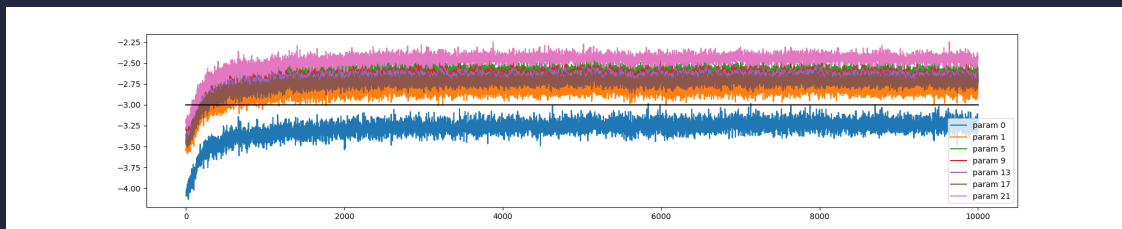
Produced with 2. These are backward pass statistics since they involve the gradient

3. The ratio between the gradient and the parameter. Should be similar for each layer after the first and before the last. Restricting to 2-dim parameters, i.e. weights of the linear layers, not biases or γ and β in batch norm.



Produced with 3.

4. The size of the updates relative to the parameters. As an estimate, want the \log_{10} to be around 10^{-3} . This can be more informative than the gradient to data ratio, since what matters in learning is the size of the update.



Produced with 4.

Question

Why do we look at the ratio of grad standard deviations to parameter standard deviation (similarly for updates), instead of absolute values.

Batch Normalization

Batch normalization is a technique used to stabilize training. It computes the mean and sd of the batch and uses that to normalize the batch to a standard normal distribution. This helps ensure neuron activation is more stable. For example, if the non-linear layer is tanh, it will take values of ± 1 for inputs that are too far away from 0. Oversaturation could keep some neurons always on or off. It also means that adding a bias in a linear layer that comes before a batch normalization step is redundant as that bias is added to the norm and immediately subtracted, rendering it useless.

Backpropagation

Be aware when a column is replicated before an elementwise operation is performed, you need to sum the gradients along the duplicated rows. If a value x feeds into to different paths, need to sum those gradients together. Keep track to make sure torch broadcasting is properly backpropagated, this can be done by checking shape.

Containers

Containers from PyTorch are ways of organizing layers into lists or dicts.

Experimental Harness

When training takes longer and longer, need a way to test hyper-parameters such as learning rates and tune those without running the full training over and over.

Self attention trick

```
1 import torch
2
3 torch.manual_seed(1337)
4 B, T, C = 4, 8, 2 # batch, time (block_size), channels (how many outputs)
5 x = torch.randn(B, T, C)
6 x.shape
```

want the 8 tokens in a block to 'talk' with each other. The fifth token in a block doesn't need info about sixth, seventh, eighth. Information only flows forwards in time. We want the average of all the vectors corresponding to the previous tokens. This is called *bag of words* (bow).

```
1 xbow = torch.zeros((B, T, C))
2 for b in range(B):
3     for t in range(T):
4         xprev = x[b,:t+1] # size (t, C)
5         xbow[b,t] = torch.mean(xprev, 0)
```

But this is very inefficient. We can speed it up with matrix multiplication.

Recall:

$$\begin{bmatrix} 1 & 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} a & b & c \end{bmatrix} = a + b + c$$

We can use a lower triangular matrix:

$$\begin{bmatrix} 1 & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{bmatrix} \cdot \begin{bmatrix} a & d \\ b & e \\ c & f \end{bmatrix} = \begin{bmatrix} \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \cdot \begin{bmatrix} a & b & c \end{bmatrix} & \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \cdot \begin{bmatrix} d & e & f \end{bmatrix} \\ \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & 0 \end{bmatrix} \cdot \begin{bmatrix} a & b & c \end{bmatrix} & \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & 0 \end{bmatrix} \cdot \begin{bmatrix} d & e & f \end{bmatrix} \\ \begin{bmatrix} \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{bmatrix} \cdot \begin{bmatrix} a & b & c \end{bmatrix} & \begin{bmatrix} \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{bmatrix} \cdot \begin{bmatrix} d & e & f \end{bmatrix} \end{bmatrix} = \begin{bmatrix} a & d \\ \frac{a+b}{2} & \frac{d+e}{2} \\ \frac{a+b+c}{3} & \frac{d+e+f}{3} \end{bmatrix}$$

We can instantly compute the averages. This lets us vectorize the previous process:

```
1 wei = torch.tril(torch.ones(T, T)) # tril = triangular, lower
2 wei = wei / wei.sum(1, keepdim=True) # divide by each rows total
3 xbow2 = wei @ x
```

The dimensions don't line up, we have: $(T, T) \cdot (B, T, C)$. So torch will replicate the (T, T) B many times giving dimensions: $(B, T, T) \cdot (B, T, C)$ and then multiply each batch element in parallel, so we have a $(T, T) \cdot (T, C) = (T, C)$ computation for each element in the batch B .

There is a third way to do this, using *softmax*.

```
1 tril = torch.tril(torch.ones(T, T))
2 wei = torch.zeros(T, T)
3 wei = wei.masked_fill(tril == 0, float('-inf')) #all elements that are zero in tril (lower
4 #triangular) become -infinity
5 wei = F.softmax(wei, dim=-1) # e to each value, then normalize the rows to sum to 1
6 xbow3 = wei @ x
```

Can generalize this idea to more general weighted averages based on how important some past tokens are to the current token.

Here is the implementation for a single self attention head. Every single token will emit two vectors: a query vector and a key vector. The query vector is 'what am I looking for', the key vectors is 'what do I contain'. The query vector of a token dot products with the key vector of all previous tokens in the block and this should be the information stored in *wei*, the weights. If the key and query are aligned (in space, which should be learned), their dot product will be high, ensuring that the previous token with that key has more influence on the current vector querying the context.

We also produce a *value* parameter. We don't aggregate x , instead we aggregate value applied to x .

Attention is a communication mechanism between nodes in a directed graph. A given node aggregates information with a weight sum from all the nodes that point towards it. There is no notion of space (such as in convolution). Attention simply acts over a set of vectors, hence the need to positionally encode tokens. This is why there is a positional encoding added to the token embeddings. You don't need to have the lower triangular masking in more general settings, such as in an encoder block in sentiment analysis.

Self-attention because keys, queries and values all come from the same x . In general, you could have a case where queries come from x , but keys and values come from an external source of context, where would like to pull that into our nodes.

Scaled attention divides QK^T by $\sqrt{\text{head size}}$. This preserves the variance of wei at 1. This is important, especially at initialization, to ensure that wei has diffuse values. When being fed into *softmax*, any values that are too positive or negative will cause wei to converge to one-hot vectors.

```

1 import torch
2 import torch.nn as nn
3 from torch.nn import functional as F
4
5 torch.manual_seed(1337)
6
7 B, T, C = 4, 8, 32 # batch, time (block_size), channels (how many outputs)
8 x = torch.randn(B, T, C)
9
10 # single head self-attention
11 head_size = 16 # this is a hyper parameter
12 key = nn.Linear(C, head_size, bias=False) #linear layer with no bias is a matrix multiply with
    fixed weights
13 query = nn.Linear(C, head_size, bias=False)
14 value = nn.Linear(C, head_size, bias=False)
15 k = key(x)      # (B, T, 16)
16 q = query(x)    # (B, T, 16)
17
18 # this is the step where the communication occurs
19 wei = q @ k.transpose(-2, -1) # transpose last two dimensions, not the first (B, T, 16) @ (B, 16,
    T) --> (B, T, T)
20 wei = wei * head_size**-0.5   # scaled attention
21
22 tril = torch.tril(torch.ones(T, T))
23 wei = wei.masked_fill(tril == 0, float('-inf')) #all elements that are zero in tril (lower
    triangular) become -infinity
24 wei = F.softmax(wei, dim=-1) # e to each value, then normalize the rows to sum to 1
25
26 v = value(x)      # (B, T, 16)
27 out = wei @ v      # (B, T, T) @ (B, T, 16) --> (B, T, 16)

```

Residual Connections

When training very large/deep neural networks, optimization begins to degrade with size. This is not explained by overfitting along, as even training errors grows. A way to counter this is with residual connections, where identity blocks are added in. The blocks take the inputs, unaltered, and add them to outputs at later stages in the training. This ensures that performance does not become worse. It also places little strain on optimization as the gradient of addition is simple. You can fork off the residual pathway at will, and project back in by addition. The blocks that fork off are initialized to have very small gradients, so early on in training there is a quick gradient path back from the outputs to the inputs. Only later in training do these blocks come on line. You want a clean residual stream, no normalizations inside *why?* *This is a quote from Karpathy.*

Dropout layers prevent overfitting. Every forward/backward pass, randomly turns off some neurons, ends up training an ensemble of networks that all end up working together.

Makemore Code Snippets

```

1 # visualize histograms of layer activation
2 plt.figure(figsize=(20, 4))
3 legends = []
4
5 for i, layer in enumerate(layers[:-1]): #exclude output layer

```

```

6     if isinstance(layer, Tanh):
7         t = layer.out
8         print('layer %d (%10s): mean %+.2f, std %.2f, saturated: %.2f%%' % (i,
9             layer.__class__.__name__,
10             t.mean(), t.std(),
11             (t.abs()>0.97).float().mean()*100))
12         hy, hx = torch.histogram(t, density=True)
13         plt.plot(hx[:-1].detach(), hy.detach())
14         legends.append(f'layer {i} ({layer.__class__.__name__})')
15 plt.legend(legends)
16 plt.title('activation distribution')
17 # plt.show()

```

Listing 1: Activation graph with matplotlib

```

1     # histogram of gradients
2 plt.figure(figsize=(20, 4))
3 legends = []
4 for i, layer in enumerate(layers[:-1]): #exclude output layer
5     if isinstance(layer, Tanh):
6         t = layer.out.grad
7         print('layer %d (%10s): mean %+.2f, std %e' % (i,
8             layer.__class__.__name__,
9             t.mean(), t.std()))
10        hy, hx = torch.histogram(t, density=True)
11        plt.plot(hx[:-1].detach(), hy.detach())
12        legends.append(f'layer {i} ({layer.__class__.__name__})')
13 plt.legend(legends)
14 plt.title('gradient distribution')

```

Listing 2: Gradients of tanh layers with matplotlib

```

1 # parameter values visualization
2 # scale of gradient compared to actual values
3 plt.figure(figsize=(20, 4))
4 legends = []
5 for i, p in enumerate(parameters): #exclude output layer
6     t = p.grad
7     if p.ndim == 2:
8         print('weight %10s | mean %+.f | std %e | grad:data ration %e' % (tuple(p.shape),
9             t.mean(), t.std(), t.std() / p.std() ))
10        hy, hx = torch.histogram(t, density=True)
11        plt.plot(hx[:-1].detach(), hy.detach())
12        legends.append(f'{i} {tuple(p.shape)}')
13 plt.legend(legends)
14 plt.title('weights gradient distribution')

```

Listing 3: Scale between gradients and values of parameters

```

1 #before training code
2 ud=[]
3 ##### TRAINING #####
4 # in training loop
5 # lr is learning rate
6 # this tracks the log base 10
7 with torch.no_grad():
8     ud.append([(lr*p.grad.std() / p.data.std()).log10().item() for p in parameters])
9 #####
10 # at each update, keep track of log of update size compared to data size
11 # these should be low or else we are over updating
12 # last layer can be large since that layer was artificially compressed
13 plt.figure(figsize=(20,4))
14 legends = []
15 for i,p in enumerate(parameters):

```

```

16     if p.ndim == 2:
17         plt.plot([ud[j][i] for j in range(len(ud))])
18         legends.append('param %d' % i)
19 plt.plot([0, len(ud)], [-3, -3], 'k') # ratios should be 1e-3
20 plt.legend(legends)
21 plt.show()

```

Listing 4: Size of update scaled by size of the parameter