# **Mastering Model Building**

How to design (and debug) your ML model

Nicole Hartman nicole.hartman@tum.de

#### Train the Trainer workshop 7<sup>th</sup> February 2023





## Who am I?

- ORIGINS Data Science Lab post-doc
- Particle physicist



• Big data -> big opportunities!!



Generative models / Density Ratio Estimation





In the context of science, the well-known adage "a picture is worth a thousand words" might well be "a model is worth a thousand datasets"









### What we'll cover today



Bias / variance trade-off

### Training techniques



Feature choices

Would love feedback + discussions!

Open question

Going deeper





Open question

5 / 67

٦Л

### Working example





## Underfitting

Model is not expressive enough

### Training data





ТШ



>

Overfitting



Small (zero) training error

пп



## **Optimal model complexity**

Fit 2nd order polynomial to quadratic distribution







### Bias / variance tradeoff



### Bias / variance tradeoff: maths 1

- Training dataset  $S = \{x^{(i)}, y^{(i)}\}_{i=1}^{n}$ 
  - Truth labels  $y = h^*(x) + \xi$ 
    - $h^*$  : ground truth function
    - $\boldsymbol{\xi}^{(i)} \sim \mathcal{N}(\boldsymbol{0}, \sigma^2)$
- Train model  $\hat{h}_{S}$  on dataset S
- Consider test point (x,y) and quantify the *expected test error*:

$$\begin{split} MSE(x) &= \mathbb{E}_{S,\xi} \left[ (y - h_S(x))^2 \right] \stackrel{\text{Defn of } y}{=} \mathbb{E} \left[ (h^*(x) + \xi - h_S(x))^2 \right] \\ &= \mathbb{E} \left[ \xi^2 \right] + 2 \mathbb{E}[\xi] \stackrel{0}{\cdot} \mathbb{E}[h^*(x) - h_S(x)] + \mathbb{E} \left[ (h^*(x) - h_S(x))^2 \right] \\ &= \sigma^2 + \mathbb{E} \left[ (h^*(x) - h_S(x))^2 \right] \end{split}$$



### Bias / variance tradeoff: maths 2

- Let  $h_{avg}(x) = \mathbb{E}_{S} |h_{S}(x)|$  the performance of the model trained on infinitely many datasets
- Substitute this back into the MSE expression

 $MSE(x) = \sigma^{2} + \mathbb{E} \left[ (h^{*}(x) - h_{S}(x))^{2} \right]$  $= \sigma^2 + \mathbb{E} \left[ \left( h^*(x) - h_{avg}(x) + h_{avg}(x) - h_S(x) \right)^2 \right]$ No cross-term because  $= \sigma^{2} + (h^{*}(x) - h_{avg}(x))^{2} + \mathbb{E}\left[(h_{avg}(x) - h_{S}(x))^{2}\right]$ Variance Bias<sup>2</sup> Error on this class of models How does this instantiation compare with the other possible ones?

 $\mathbb{E}\left[h_{avg}(x) - h_{S}(x)\right] = 0$ 







Lots of possibilities for the fitted function depending on the random realization of training data.

CS 229 Lecture

### What's the culprit?



CS 229 Lecture

### What's the culprit?





### Starting off...

#### Feature choices

#### Statistical learning theory Bias / variance trade-off

## Training techniques

Open question



Would love feedback + discussions!

ΠП

Going deeper

### How to train?



Want a *large enough* validation set for statistically significant generalization metric.

Want a *large* training dataset to minimize  $\mathscr{L}$ .



Image by brgfx on Freepik







Minimize  $\mathscr{L}$  by SGD on K splits for the dataset  $w = w - \alpha \nabla_w \mathscr{L}$ 





✓ Couple percent gain in accuracy✓ Used in Kaggle competitions!

## K-fold cross validation



### Learning rate

## Minimize $\mathscr{L}$ by SGD $w = w - \alpha \nabla_w \mathscr{L}$

How to choose  $\alpha$ ?



### Label the loss curves!

very high learning rate high learning rate good learning rate low learning rate



...

### **Batch size**



- Bigger batches reduces the error on the MC estimate
- As large as possible to still fit on the GPU.
- Powers of 2 for memory efficiency E.g, 256, 512, 1024



Andrej Karpathy 🤣 post

The most dramatic optimization to nanoGPT so far (~25% speedup) is to simply increase vocab size from 50257 to 50304 (nearest multiple of 64). This calculates added useless dimensions but goes down a different kernel path with much higher occupancy. Careful with your Powers of 2.

### Intimately tied to learning rate!

If you increase the batch size by a factor of 2, scale  $\alpha$  by  $\frac{1}{2}$  for a fair comparison

CS231n Lecture 7

## Early stopping



### Hyperparameter search

Already many options...



Random Search for Hyper-Parameter Optimization Bergstra and Bengio, 2012 Image from CS231n lecture

### Hyper-parameter search



### Hyperparameter strategies



### Fast prototyping

- 1. Start with a subset of the training dataset
- 2. Find parameters for a model that overfits
- 3. Start the random search around this point

Scan in log space







Coarse scan:

3 activations {sigmoid, ReLU, ELU}

- x **3** layers {5, 10, 20}
- x 4 nodes {10, 50, 250, 200}
- x **4** learning rates {1e-2, 3e-3, 1e-3, 3e-4}
- x **2** {with and w/o scheduler}
- x 10 K-fold cross validation (K=10)

= 3840 trainings !!!



### Loss landscape

### Linear functions



(convex loss)

Neural networks







# Models in the Deep Learning Era





### The Deep Learning Revolution

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



### Natural Language Processing

31 / 67



What about our model complexity · / discussion?

Train the Trainer Deep Learning philosophy... Bigger = better

ResNets: 1512.03385 Slide CS231n lecture

### **ResNet building block**



### **ResNet: performance**



34



# **Regularization Techniques**

How to restrict the optimization problem to help the NN generalize better.

•

35 / 67

### Regularization

 $\mathscr{L}_{\beta}(w) = \mathscr{L}(w) + \beta \mathscr{R}(w)$ 

Fit the training data

Penalize complicated models

Hyperparameter governing tradeoff of the two objectives.

## Occam's razor for ML

When multiple models describe the training data... choose the simplest one!


L2 Regularization (most common for NNs)

$$\mathscr{L}_{\beta}(w) = \mathscr{L}(w) + \beta |w|^2$$

Encourages weights to be small

"Weight decay"

$$w = w - \alpha \nabla_{w} \mathscr{L}_{\lambda}$$
$$= w - \alpha \nabla_{w} (\mathscr{L} + \beta w^{2})$$
$$= (1 - 2\beta)w - \alpha \nabla_{w} \mathscr{L}$$
Decay

Include in random search

(Log scale, e.g,  $\beta$  = 0, 1e-6, 1e-4)



ЛШ

### L1 regularization (inducing sparsity)

 $\mathscr{L}_{\beta}(w) = \mathscr{L}(w) + \beta |w|$ 

#### Encourages weights to be small







Srivastava et al, "Dropout: A simple way to prevent neural networks from overfitting", JMLR 2014 CS231n Lecture 7

### Dropout: intro

Issue: Don't want the NN to rely heavily on individual features





Srivastava et al, "Dropout: A simple way to prevent neural networks from overfitting", JMLR 2014 CS231n Lecture 7

### Dropout: idea

Issue: Don't want the NN to rely heavily on individual features



- At training time, zero out some neurons with dropout fraction p Hyperparameter we need to optimize!
- Encourages learning robust features





At **test time**, use all the neurons for prediction!!

Like training an ensemble of the NNs without being as  $\in \in \in$ 

### Dropout: performance



### **Batch Normalization**

#### What about the features in the hidden layers?

"You want zero mean and unit variance operations? Just make them so!"



At **training time**, normalize over the activations of the minibatch

Additional learnable parameters for the scale and shift:  $\gamma$  and  $\beta$ .

Original paper: inserted after Fully Connected layers, and before nonlinearity.

Some debate, e.g, 1912.04259 suggests other placements better in some cases.



For set / sequence based data (Nikolai's talk) <u>LayerNorm</u> also useful

 $\sqrt{\text{Var}\left[x^{(k)}\right]}$  $y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$ 

 $\hat{x}^{(k)}$ 

### Data augmentation : motivation

### Recall: More training data reduces variance.



**Q:** How can you modify your training data to artificially increase your dataset size?

CS231n Lecture 7

### Data augmentation : examples

### 1 Horizontal flips





#### 2 Random crops and scales



ТΠ









### Data augmentation — physics



### Less used in practice...

With a simulator and we can often get as many training examples as we want



### Data augmentation — alternatives

Issue: Larger models with more data take longer to train

➡ Remove the variation to <u>train faster</u>.





## Architecture design to preserve invariants



Lorentz Group Equivariant Block (LGEB)

LorentzInvariance: Lorentz Net and ParT 2202.03772 Azimuthal Symmetry: 2107.02908 + many others



#### Fine-tuning / transfer learning **Original Task** Modified task Particle ID Particle ID 95% accurate 75% accurate x 20 million x 200 thousand (1) Start with weights optimized **Fine-tuning** with the larger dataset. (2) Adjust weights with the smaller dataset. Particle ID (3) Enjoy increase 95% accurate in performance! x 200 thousand 47 / 67

### Fine-tuning / Transfer learning: chatGPT

0





I want you to act as a Julia REPL. I will type commands and you will reply with what the REPL should show. I want you to only reply with the REPL output inside one unique code block, and nothing else. Do not write explanations. Do not type commands unless I instruct you to do so. When I need to tell you something in English I will do so by putting text inside curly brackets {like this}. My first command is VERSION.

\$		ငြို Copy code
	julia> VERSION v"1.7.0"	
0	cholesky([2.0 0.5; 0.5 1.2])	
		۹ Comu eo de
	julia> cholesky([2.0 0.5; 0.5 1.2])	
	LowerIrlangular{Float64,Array{Float64,2}}	
	[0.447 0.785]	

Apparently... chatGPT has learned math and **julia** syntax



### Fine-tuning / Transfer learning





### Weight sharing

# Alternative architectures designed to *reuse weights* as suited for the input data.

"Don't relearn what you don't need"



#### See next talk by Judith!





Positional Encoding

Input

Embedding

Inputs

Positional Encoding

Output

Embedding

Outputs (shifted right)

1912.02757

### **Ensembles**



### **Ensembles:** Application



### **Ensembles: Application**



### Starting off...

#### Statistical learning theory Bias / variance trade-off

Open question

#### Feature choices



### Training techniques



Would love feedback + discussions!

ПΠ



Going deeper



### Log transform: motivation



How does this help? 20% speed up in training time!



### Sample dependence

**Issue:** Want a classifier that is performant over a range of energies



#### CS 229 Lecture

### Ablation studies: What has the model learned?

### Example — spam classification





### Ablation studies



Component	Accuracy	
Overall system	99.9%	
Spelling correction	99.0	
Sender host features	98.9%	
Email header features	98.9%	
Email text parser features	95%	
Javascript parser	94.5%	
Features from images	94.0%	



Email text parser: most important feature!

[baseline]

TH I

### Saliency maps

- NN: nonlinear function
- Approximate as a linear classifier by using a *Taylor expansion*.

$$S_{c}(I) \approx \theta^{T}I + b$$
$$\theta = \frac{\partial S_{c}}{\partial I} \bigg|_{I_{0}}$$

Saliency map: Plot the  $|\theta|$  for each of these inputs





60 / 67



### Saliency maps

### **Physics**

Understand what the model has learned about this particle ID task.





### **Maths**

Use saliency maps to postulate **new** conjectures which could then become new math theorems!

#### Article

### Advancing mathematics by guiding human intuition with AI

https://doi.org/10.1038/s41586-021-04086-x Received: 10 July 2021 Accepted: 30 September 2021 Alex Davies<sup>153</sup>, Petar Veličković<sup>1</sup>, Lars Buesing<sup>1</sup>, Sam Blackwell<sup>1</sup>, Daniel Zheng<sup>1</sup>, Nenad Tomašev<sup>1</sup>, Richard Tanburn<sup>1</sup>, Peter Battaglia<sup>1</sup>, Charles Blundell<sup>1</sup>, András Juhász<sup>2</sup>, Marc Lackenby<sup>2</sup>, Geordie Williamson<sup>3</sup>, Demis Hassabis<sup>1</sup> & Pushmeet Kohli<sup>154</sup>

Published online: 1 December 2021 Open access

Check for updates

The practice of mathematics involves discovering patterns and using these to formulate and prove conjectures, resulting in theorems. Since the 1960s, mathematicians have used computers to assist in the discovery of patterns and formulation of conjectures<sup>1</sup>, most famously in the Birch and Swinnerton-Dyer conjecture<sup>2</sup>, a Millennium Prize Problem<sup>3</sup>, Here we provide examples of new fundamental results in pure mathematics that have been discovered with the assistance of machine learning-demonstrating a method by which machine learning can aid mathematicians in discovering new conjectures and theorems. We propose a process of using machine learning to discover potential patterns and relations between mathematical objects, understanding them with attribution techniques and using these observations to guide intuition and propose conjectures. We outline this machine-learning-guided framework and demonstrate its successful application to current research questions in distinct areas of pure mathematics, in each case showing how it led to meaningful mathematical contributions on important open problems: a new connection between the algebraic and geometric structure of knots, and a candidate algorithm predicted by the combinatorial invariance conjecture for symmetric groups<sup>4</sup>. Our work may serve as a model for collaboration between the fields of mathematics and artificial intelligence (AI) that can achieve surprising results by leveraging the respective strengths of mathematicians and machine learning.

Nature 600, 70-74 (2021)





### Loss curves — what are the problems? <u>Hypotheses</u>

Slow start: initialization learning rate too small

Applied the negative of the gradients

Not converged yet: need longer training

Not learning: gradients not applied to the weights

Overfit: model too large / dataset too small



17

#### Loss curves



- val

2500

iterations

2000

















1912.02292 Daniela Witten @daniela\_witten post

### **Double Descent**





# Backup





CS 229 notes



### High bias: diagnostics

$$MSE(x) = \sigma^{2^{*}} + (h^{*}(x) - h_{avg}(x))^{2} + \mathbb{E}\left[(h_{avg}(x) - h_{S}(x))^{2}\right]$$

The training error on the linear model still large, even when there is *no noise* on the training data.





### Learning rate schedule

# Instead of training independent models, use multiple snapshots of a single model during training!



Loshchilov and Hutter, "SGDR: Stochastic gradient descent with restarts", arXiv 2016 Huang et al, "Snapshot ensembles: train 1, get M for free", ICLR 2017 Figures copyright Yixuan Li and Geoff Pleiss, 2017. Reproduced with permission.



Cyclic learning rate schedules can make this work even better!

 $\alpha < ?$ 



From CS 229 notes Section 9.2

### Batch size — caveat



### Alternative philosophy...

Smaller mini-batches encourage stochasticity and converge to flatter minima which generalize better.


ΠП

#### **ResNets: Motivation**



Complex model doing *worse* than the simple one on the training dataset!

Doesn't make sense from the bias / variance picture we built up earlier! Implies an *optimization problem*.

#### Input representation



Collection of tracks:  $X_i$ :  $i = \{ 1, ..., n \}$ 

Each track has features:  $X_i \in \mathbb{R}^m$ 

#### Jet has labels Y

- Quark / gluon tagging:  $Y \in \{q, g\}$
- Higgs tagging:  $Y \in \{H, top, QCD\}$
- Top tagging:  $Y \in \{top, QCD\}$
- $\bullet \quad b\text{-tagging:} \ Y \in \{b, \, c, \, l\}$

n\* m ~ O(100)



SLAC

### Ensembles

Good way to interpret models!! Let's us *probe* whether experiments are independent of each other



#### **Debugging learning algorithms**

Motivating example:

- Anti-spam. You carefully choose a small set of 100 words to use as features. (Instead of using all 50000+ words in English.)
- Bayesian logistic regression, implemented with gradient descent, gets 20% test error, which is unacceptably high.

$$\max_{\theta} \sum_{i=1}^{m} \log p(y^{(i)}|x^{(i)},\theta) - \lambda ||\theta||^2$$

• What to do next?

#### **Fixing the learning algorithm**

Bayesian logistic regression:

 $\max_{\theta} \sum_{i=1}^{m} \log p(y^{(i)}|x^{(i)}, \theta) - \lambda ||\theta||_{1}^{2} e, \text{ perhaps a jet classification ex?}$ 

- · Common approach: Try improving the algorithm in different ways.
  - Try getting more training examples.
  - Try a smaller set of features.
  - Try a larger set of features.
  - Try changing the features: Email header vs. email body features.
  - Run gradient descent for more iterations.
  - Try Newton's method.
  - Use a different value for λ.
  - Try using an SVM.
- This approach might work, but it's very time-consuming, and largely a matter of luck whether you end up fixing what the problem really is.



#### **Diagnostic for bias vs. variance**

Bia

Better approach:

- Run diagnostics to figure out what the problem is.
- Fix whatever the problem is.

Bayesian logistic regression's test error is 20% (unacceptably high).

Suppose you suspect the problem is either:

- Overfitting (high variance).
- Too few features to classify spam (high bias).

Diagnostic:

- Variance: Training error will be much lower than test error.
- Bias: Training error will also be high.

### **Overfitting pictures**



### Classification



Regression

### Deep Sets: for b-tagging



#### **CNNs**

#### Translational invariance



Strong inductive bias: or domain knowledge we inject to efficiently converge to a solution.

Y

### **Recurrent neural networks**

### Model the jet as a sequence

- Efficient representation: same tranformation at every time step
- Sequential nature: variable # of tracks → <u>fixed</u> dim vector for the jet



### Transformers

(shifted right)

"Any time you have a deep set you can substitute in a transformer as it's always more expressive."



### Input processing: $\Delta R$

#### Ex: b-tagging input features



-0.4 - 0.3 - 0.2 - 0.1

Training with  $\Delta R(trk,jet)$  encodes a symmetry in the input representation - and makes learning features easier.

SLAC

0.4

0.3

0.1

0.0

Δη

0.2

# Saliency Maps

### Model diagnostic



1312.6034

### Saliency maps

#### What has DIPS learned about b-jets?



**SLAC** 

ATL-PHYS-PUB-2020-014

#### Saliency definition

#### What has DIPS learned about b-jets?

• Consider b-jets failing the 77% WP



- Average over jets with 8 tracks
- Sort the tracks by s<sub>d0</sub> for the average



87

## **Saliency for DIPS**

#### What has DIPS learned about b-jets?

- ✓ Want at least 5 high impact parameter tracks
- ✓ Wants harder leading track as expected from b-decay
- ✓ Larger opening angle corresponding to geometrical constraint from more displaced tracks
- ✓ Want good quality for displaced tracks



b-jets with 8 associated tracks failing a threshold corresponding to a 77% b-tagging efficiency 0.4 nSCTHits **ATLAS** Simulation Preliminary  $V_{inputs} D_b$  $\sqrt{s}$  = 13 TeV,  $t\bar{t}$ nPixHits 0.3 shared SCT hits 0.2 split pixel hits shared pixel hits -0.1 split IBL hits -0.0 shared IBL hits IBL hits -0.1 -PIX1 hits  $\log \Delta R$ -0.2  $\log p_T^{frac}$ -0.3 *S*<sub>20</sub>  $s_{d0}$ -0.4 2 5 6 8 3 4 7 Tracks sorted by s<sub>d0</sub> Smallest sd0 Largest sd0

### The Deep Learning Revolution



An Analysis of Deep Neural Network Models for Practical Applications, 2017.

Figures copyright Alfredo Canziani, Adam Paszke, Eugenio Culurciello, 2017. Reproduced with permission.

### Choosing the right features

### Feature extraction + learning



- Feature extraction: set  $\mathcal{F}$  based on domain knowledge
- Learning: set  $f_{\mathbf{w}} \in \mathcal{F}$  based on data



### Effect of hypothesis class size



As the hypothesis class size increases...

Approximation error decreases because:

taking min over larger set

Estimation error increases because:

harder to estimate something more complex

How do we control the hypothesis class size?

#### Loss curves — what are the problems?



#### Problem: Not shuffling data, periodical patterns in loss curve



Problem: val set too small, statistics not meaningful



Get nans in the loss after a number of iterations: caused by numerical instability in models



### **Batch Normalization**





Fig. 1. Different arrangements of layers used in this study. (a) arrangement 1, (b) arrangement 2 and (c) arrangement 3. (The Pooling layer is dashed because it might not always be present after the third layer in these arrangements.)



ТЛП

#### **Double Descent**

