

# fastbook ch.1 DEEP LEARNING APPLICATIONS: NLP @ vision & medicine CLASSICAL PROGRAM: → ML APPROACH:

by ⚡ @dk21

BIOLOGY IMAGE GENERATION

RECOMMENDATION SYSTEMS

ROBOTICS LOGISTICS FINANCE

ML: discipline where we define program not by writing it entirely ourselves, but by learning from data

DL: subset of ML, using NN with multiple layers

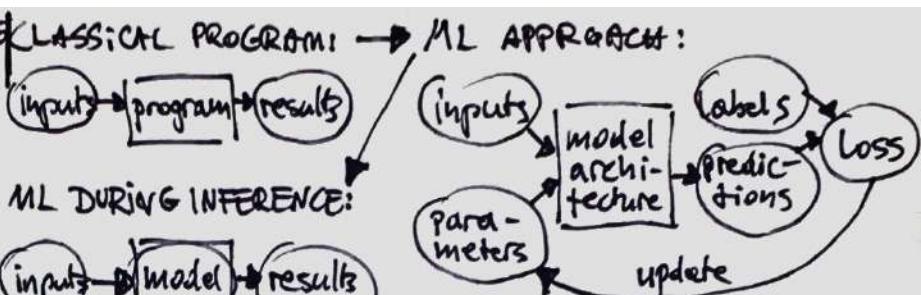
TEACHING APPROACH: TOP-DOWN (VIA D. PERKINS, BASEBALL ANALOGY)

Start with real, practical examples (teach the whole game)

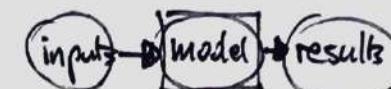
Learn by doing (run/re-implement code)

Deep dive into theory later, as needed - to improve models

Simplify and remove barriers:  
• fastai library  
• Pytorch  
• Jupyter



ML DURING INFERENCE:



LIMITATIONS OF ML:

- Need data with labels
- Can only learn patterns seen in input data
- Predictions vs recommended actions

WATCH OUT:  
FEEDBACK LOOPS!  
e.g. YouTube recommending viral anti-vax videos

Some history:

1943 ~ McCulloch, Pitts: Artificial Neuron

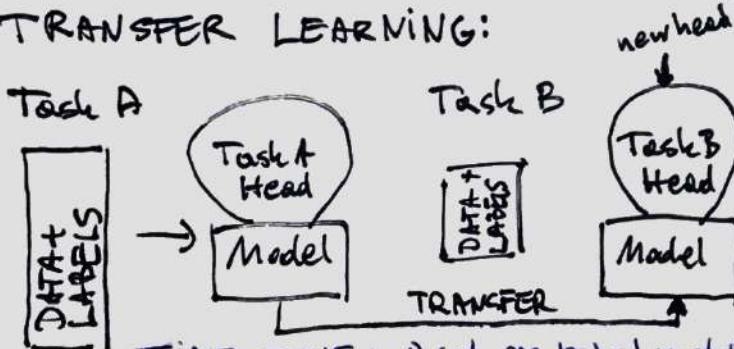
Rosenblatt's device: Mark I Perceptron

Minsky, Papert: Perceptrons - book introducing multiple layer neural networks

1986 - Parallel Distributed Processing - book introducing most of current DL framework

1961/62: Samuel - Artificial Intelligence essay, introducing current ML approach, program beating humans in checkers

TRANSFER LEARNING:



- FINE-TUNE:
- 1) Cut pre-trained model's head
  - 2) Add new head specific to new task
  - 3) Train new head only for 1 epochs
  - 4) Fit entire model for more epochs (more tricks here)

```
text | tabular | collab
from fastai.vision.all import * # fastai library
path = untar_data(URLs.PETS)/'images' # download dataset
def is_cat(x): return x[0].isupper() # labelling function
dls = ImageDataLoaders.from_name_func(
    path, get_image_files(path), valid_pct=0.2, seed=42,
    label_func=is_cat, item_tfms=Resize(224))
# load data, label, split into train/valid, transform
learn = cnn_learner(dls, resnet34, metrics=error_rate)
# load architecture, pretrained model, define metric
learn.fine_tune(1) # finetune ~ fit pretrained model
```

OVERFITTING! Single most important and challenging issue! Model starts memorizing examples, rather than learning to generalize!



As a rule, split data into train/validation set, or maybe test set as well:

TRAINING | VALID | TEST

Train data on this only

Tune model, choose hyperparameters, see performance on

Hide test, use it to estimate performance at the very end



AVOID LEAKAGE: Time series, same subjects in train/test etc...

Different layers in NNs learn to recognize increasingly complex features.

CNN - Convolutional Neural Network example:



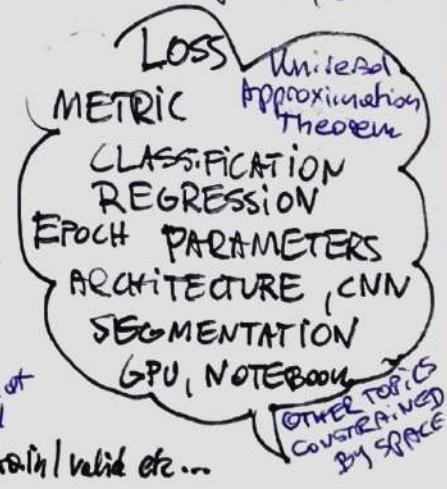
L1: edges

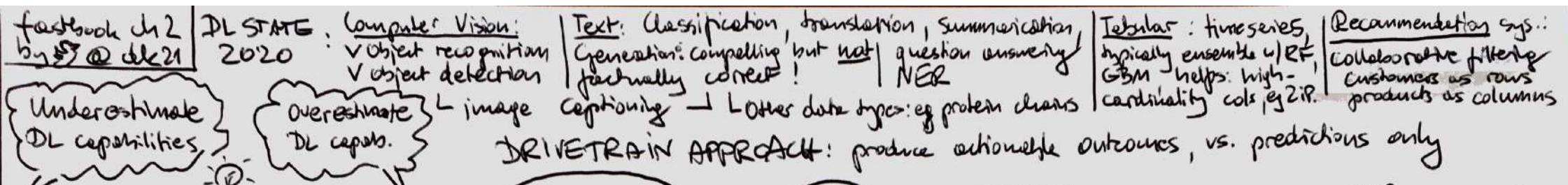


L2: corners, shapes, colors, patterns



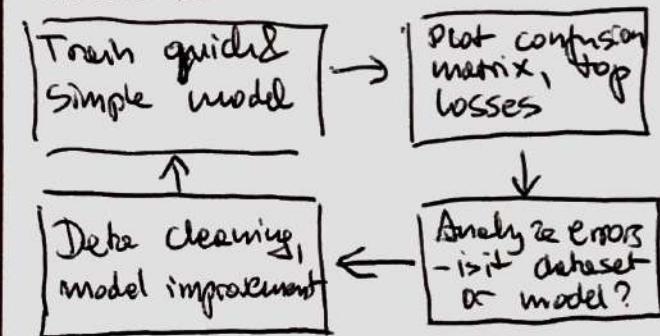
L3+: higher level concepts: faces, objects, etc.





- ① Complete lots of small experiments and work on your own project
- ② Consider data availability
- ③ Iterate E2E - all the way to final product
- ④ Start with tasks that DL is good at.

### WORKFLOW:



### HOW TO AVOID DISASTER?

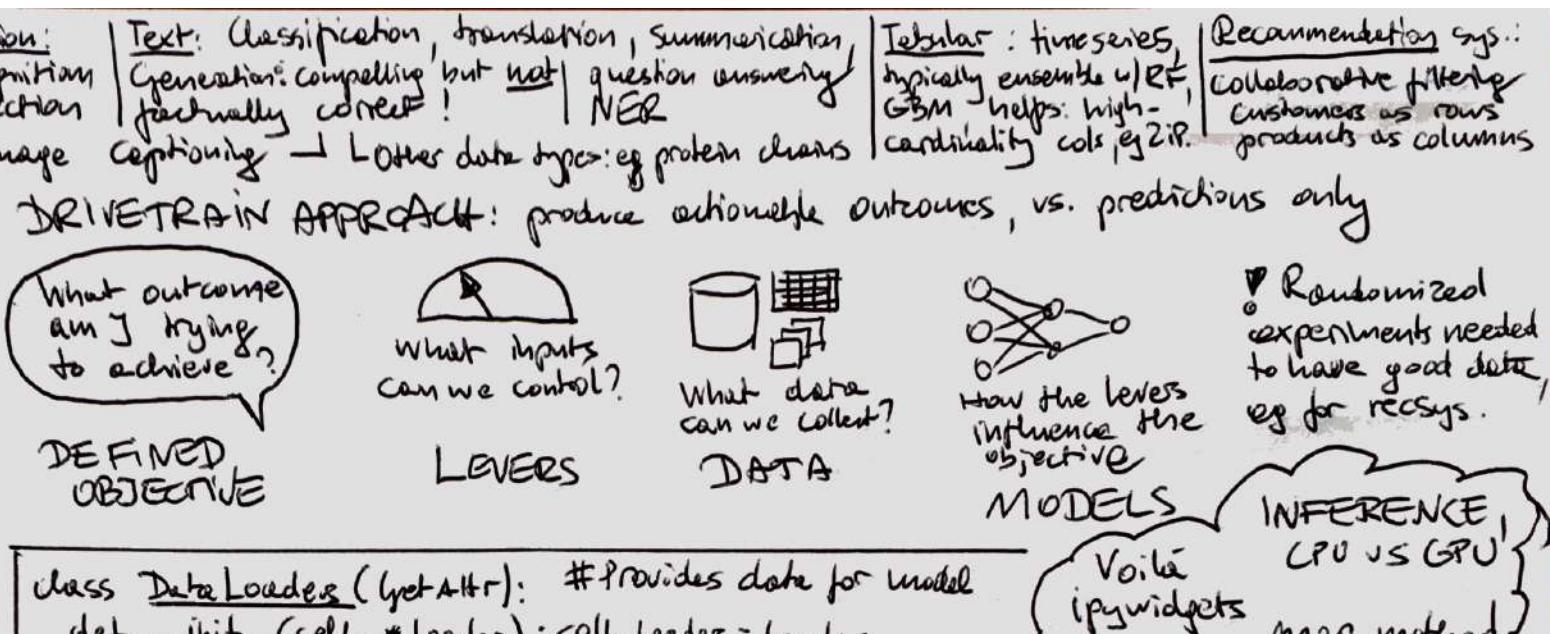
- Out of domain data: training vs production
- Domain shift: data changes over time
- Anticipate unforeseen consequences and feedback loops. What if this went really well?

DATA  
1. Manual process  
- Run model in parallel  
- Humans check all predictions

2. Limited scope deployment  
- Careful human supervision  
- Time or geography limited

3. Gradual expansion  
- Good reporting systems needed  
- Consider what could go wrong!

TIP: Start writing, blog!  
Write for people one step behind you.



```

class DataLoaders(LightningDataModule):
    def __init__(self, *loaders):
        self.loaders = loaders
    def __getitem__(self, i):
        return self.loaders[i]
    train, valid = add_props(lambda i, self: self[i])
  
```

```

bears = DataBlock( # template for creating data loaders (DB)
    blocks=(ImageBlock, CategoryBlock), # type of independent/dependent variable
    get_items=get_image_files, # function takes a path and returns images in that path
    splitter=RandomSplitter(valid_pct=0.3, seed=42), # split train/valid, fix seed
    get_y=parent_label # label images
    item_tfms=Resize(128) # item_tfms run on CPU, eg. image resizing
    dls = bears.dataloaders(path) # provide source of data to (DB) - has path with images
  
```

(space constrained)

RESIZE OPTIONS: crop, squish, pad, RandomResizedCrop & recommended

DATA AUGMENTATION: create random variation in the input data, standard self provided in aug-transforms, can be done on GPU in batch:

```

beats = bears.new(item_tfms=RandomResizedCrop(128, min_scale=0.5), batch_tfms=aug_transform,
      learn.export(): saves both model and parameters. load_learner(path/'export.pkl') = loads model()
  
```

~~Last slide Ch.3~~ ETHICS: the study of right & wrong, how we define & recognize them, understand the connection between action & consequences  
notes by @dk21 DATA ETHICS: complicated, context dependent → learn through examples, like a ~~muscle~~ muscle → develop & practice it!

## BUGS, RECOUPSE, ACCOUNTABILITY

- Ex) Buggy algorithm cut healthcare benefits, impacting a vulnerable group
- ▷ Finger-pointing vs taking accountability bureaucracy as a way to evade responsibility
- ▷ Data often contains errors → mechanisms for audit and correction are crucial
- (Ex) Police maintaining database of gang members with no mechanism to correct obvious errors
- Ex) US credit report system - very difficult to correct errors

## FEEDBACK LOOPS

- Ex) Conspiracy theories videos tend to get recommended more on YT | FB  
People watching them tend to watch more online videos  
YT | FB recommendation algorithms suggest more similar videos

## Why you SHOULD CARE?

- Ex. IBM products used in Nazi concentration camps - would you be ok to contribute to killing people?

Ex. VW emission scandal - engineers jailed!

- ML can create feedback loops & amplify bias
- People more likely to assume algorithms are objective and error-free
- Often used at scale, with no appeals process in place
- Considering this will make you a better practitioner!

BIAS Historical bias - people, processes, society are biased - taking real world data includes these biases

MEDICAL - doctor prescriptions differ for white vs black patients  
SALES - different prices by race

- Ex) Searching google for a name that is historically black, you get ads for background checks (suggesting a criminal record)

▷ Systematic imbalance in the make-up of popular datasets, e.g. ImageNet, word embeddings

- Ex) Translating doctor man, nurse ~ women.

Measurement bias - measuring wrong thing, in the wrong way!

Incorporating it into model incorrectly

- Ex) Factors predictive of stroke
  - prior stroke
  - cardiovascular disease
  - accidental injury
  - colonoscopy

these are correlated with people actually going to a doctor, being able to afford it, vs having a stroke

Aggregation bias - eg diabetes treatment based on linear, univariate models, small studies on homogeneous groups, when reality is non-linear, e.g. diff. complications, symptoms across ethnicities

## IDENTIFYING & ADDRESSING ETHICAL ISSUES

- ① Analyze a project you're working on
  - Should we even be doing this?
  - What bias is in the data?
  - Can the code and data be audited?
  - What are error rates for different sub-groups?
  - What is the accuracy of a simple, rule-based alternative?
  - What processes are in place to handle appeals or mistakes?
  - How diverse is the team that built it?

### ② Processes to implement

- Ex. Regular ethical risk sweeps (pen testing)
  - include perspectives of a variety of stakeholders
  - what could bad actors do?
  - who will be directly and indirectly affected?
  - apply ethical lenses! which option...

[RIGHTS] best reflects the rights of all stakeholders

[JUSTICE] treat people equally or proportionately?

[UTILITARIAN] will produce most good & least harm?

[COMMON] best serves community as a whole, good, not just some members

[VIRTUE] leads me to act as the sort of person I want to be

### ③ The power of diversity

- ▷ Similar backgrounds = similar blindspots, → innovation, more risks/solutions considered

④ Role of policy-regulation is important

Ex) FB lack of action during Rohingya genocide, vs quick action to address GDPR

Advocacy is important - support the regulations that you & data scientist believe we need!

# Fundamental Tools and Concepts for Deep Learning

Fastbook Ch. 4  
Notes by @dk21

`new_list = [f(o) for o in a_list if o > phi]`

list comprehension → to do for each element, filter

Tensor shape - length of each axis  
rank - number of axes = len(shape)

Measuring distance in space:

Mean absolute distance (L1 norm)

→ absolute differences

Root mean squared error (RMSE, L2 norm)  
→ mean of square diff's then root

Numpy ARRAY: multidimensional table of data, with all items of the same type, any type. With array of arrays, arrays underneath may have different sizes → "agged array". Operations on regular arrays written in optimized C - much faster than Python.

PyTorch TENSOR - like a numpy array, but has to use simple basic numeric type for all components. Can run on GPU.

BROADCASTING - critical efficient code! PyTorch, when performing operation between tensors of different ranks, will automatically expand tensor with smaller rank to have the same size as the one with larger rank.

PyTorch VIEW change the shape of a tensor without changing contents, -1 parameter: make this axis as big as necessary to fit all data

② - PyTorch matrix multiplication

$\text{batch} @ \text{weights} + \text{bias} \Rightarrow$  fundamental operation MNN

$w * x + b$   
Weights      Biases } parameters

Universal Approximation Theorem: this can represent any function

```
def train_epoch(model, epochs):
    for i in range(epochs):
        train_epoch(model)
        print(validate_epoch(model))
```

```
def train_epoch(model):
    for xb, yb in dl:
        calc_grad(xb, yb, model)
        opt.step()
        opt.zero_grad()
```

Class Basic Optim:

```
def __init__(self, params, lr):
    self.params, self.lr = list(params), lr
    def step(self, *args, **kwargs):
        for p in self.params:
            p.data -= p.grad.data * lr
```

# we use .data so PyTorch won't take gradient  
# of this step

```
def zero_grad(self, *args, **kwargs):
    for p in self.params: p.grad = None
```

```
def calc_grad(xb, yb, model):
    pred = model(xb)
    loss = mnist_loss(preds, yb)
    loss.backward()
```

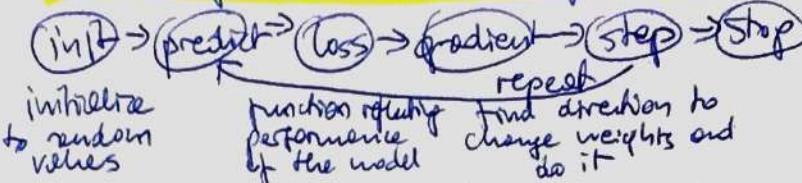
```
def linear(xb): return xb @ weights + bias
    ⇒ linear1 = nn.Linear(28*28, 1)
```

```
def mnist_loss(predictions, targets):
    predictions = predictions.sigmoid()
    return torch.where(targets == 1, 1 - predictions,
                       predictions).mean()
```

```
def init_params(size, std=1):
    return (torch.randn(size)*std).requires_grad_()
```

```
def simple_net(xb):
    res = xb @ w1 + b1
    res = res.max(tensor(0.0))
    res = res @ w2 + b2
    return res
```

SGD Gradient descent process:



initialize to random values

function returning find direction to performance of the model change weights and do it

Derivative of a function tells us how much a change in parameters will change results (rise/run)

GRADIENT: value of loss function's derivative at the point we're predicting.  $f'(x)$ : return  $x+2$

PyTorch can do it for us!

$x_t = \text{tensor}(3)$ , requires\_grad=True  
 $g_t = f(x_t)$   
 $y_t = g_t.backward()$   
 $x_t.grad # tensor(6.)$

STEPPING WITH LEARNING RATE

Multiply the gradient by a small learning rate to decide how much to change parameters:

LOSS FUNCTION: represents how good is the performance of our model. Needs to react to small changes in weights (accuracy isn't good!)

$\text{def sigmoid}(x): \text{return } 1/(1+\text{torch.exp}(-x))$   
→ ensure values between 0 and 1.

Metric - drive human understanding

To step: change the weights/biases - we need to calculate loss on 1 or more data items. 1 is not enough - not much info, not optimized. Whole dataset would be too slow → MINI-BATCH

# items = batch size (! important decision)  
We need to vary examples during training - randomly shuffle dataset on every epoch.

DATA LOADER: takes Python collection, turns it into iterator over batches:

dl = DataLoader(collection, batch\_size=8, shuffle=True)

DATA SET: collection of tuples of independent and dependent variables. In most PyTorch dataset, dataset = list(zip(X-train, y-train))

# Fastbook ch.5 DEEP DIVE INTO MECHANICS OF DL

Learn Use eg in Regex! → RegexLabeller

## PRESIZING

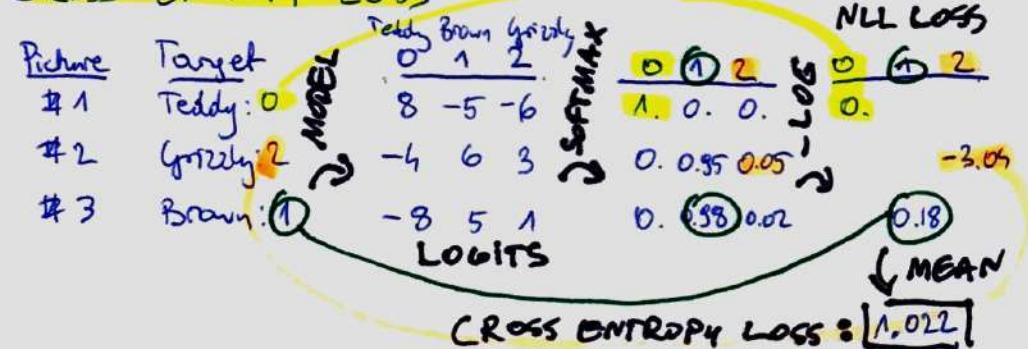
item\_tfms = Resize(460),  
this adds RandomResizedCrop  
batch\_tfms = aug\_transforms(size=224, min\_scale=0.75)

- ① Resize images to relatively large dimension (vs target training dimension)
- ② Compose all common aug operations (incl. resize to target size) into 1, and perform the combined operation on the GPU once at the end of processing  
 → avoids data losses during augmentation  
 → speeds up the process!

## CHECKING AND DEBUGGING DATA BLOCK

Show-batch → inspect data, DataBlock.summary(path)

## CROSS ENTROPY LOSS



⇒ take the softmax, then negative log likelihood of that softmax : ensure final activations are between 0-1, and sum=1  
 $\text{softmax}(x) = \frac{\exp(x)}{\sum \exp(x)}$

- if one activation is slightly higher than others, exp will amplify this - softmax really wants to pick one class - good if each image has definitive label
- may not be good at inference - will boast probs of example relative to other choices, independent of overall confidence - binary output columns, with sigmoid activation may be better?

Log Likelihood : pick loss from column with correct label only, take -log of that to transform 0-1 scale to 0->inf scale  
 $\text{PyTorch: nn.CrossEntropyLoss} \Rightarrow \text{nn.LogSoftmax} + \text{nn.NLLLoss}$

**Learning Rate Finder (Leslie Smith)**

→ start with a very small learning rate  
 → use that for 1 mini-batch, find the loss  
 → increase the lr gradually, e.g. double, per mini-batch, track the loss again  
 → keep doing this until the loss gets worse  
 → good choices: a) divide minimum loss lr by 10, or b) last point where the loss is clearly decreasing (steepest point)

## Unfreezing and transfer learning

→ remove pretrained model's classification head  
 → replace it with classif. head for new task  
 → this will have random weights initially, so we freeze pretrained layers and only train new head  
 → later, we unfreeze, check lr-finder again

**Discriminative learning rates** : pass slice (lr1, lr2) in lr-finder

→ train first layers with smaller lr, last layer with higher lr, range between - multiplicatively equidistant lr's.

## Selecting the number of epochs

- 1) Choose based on time available
- 2) Observe training / validation losses and val. metrics
- 3) Make decisions on metrics, not losses! - initially, validation loss will get worse, because it becomes overconfident, it's still ok if the metrics improve. Only later, model will start to memorize.

Early Stopping - save model after each epoch, select the one with best metric

This is NOT GOOD with 1-cycle training - epochs in the middle have higher lr, so unlikely to find best result. Better - If we overfit - refresh model from scratch with the number of epochs where we had best results.

**Deeper Architectures** : rule of thumb - more layers ⇒ more accurate model, but also risk of overfitting, longer to train, not always better!  
 We can speed this up with mixed precision training: learner.to\_fp16()

$\log(a*b) = \log(a) + \log(b) \rightarrow$  important property  
 • log increases linearly, when underlying signal increases exponentially / multiplicatively  
 • replace multiplication with addition ⇒ numerical stability  
 $\text{Gradient of cross entropy } (\mathbf{a}, \mathbf{b}) = \text{softmax}(\mathbf{a}) - \mathbf{b}$ , diff pred. vs target - linear, allows smoother training.

Factbook dr.6  
notes by @dtk21

multi-one-hot encoding

{ 0 1 0 0 1 0 1 0 0 }

**Mult-label classification:** more than one type of object in an image - more common in practice to have some images with zero or more than 1 category matches.

Use DataBlock API to construct DataLoaders object from Pandas df:

- start by creating and testing Datasets
- create DataLoaders after that's working

dblock = DataBlock(get\_x= lambda r: r['frame'],  
get\_y= lambda r: r['labels'])

dsets = dblock.datasets(df)

dsets.train[0], len(dsets.train), len(dsets.valid)

! Lambda functions (defined inline) are good for iterating, but not compatible with serialization!  
Need verbose functions to export Learner after training.

def get\_x(r): return path/'train'/r['frame']

def get\_y(r): return r['labels'].split(' ')

def Splitter(df):

train = df.index[df['is-valid']].tolist()  
valid = df.index[df['is-valid']].tolist()

return train, valid

dblock = DataBlock(blocks=(ImageBlock, MultiCategoryBlock),

splitter=splitter,

get\_x=get\_x,

get\_y=get\_y,

item\_tfms=RandomResizedCrop(128, min-scale=0.35))

dls = dblock.dataloaders(df)

dls.show\_batch(rows=1, ncols=3)

## BINARY CROSS ENTROPY (BCE)

def binary\_cross\_entropy(inputs, targets):

inputs = inputs.sigmoid()

return -torch.where(targets == 1, inputs, 1 - inputs) \* log().mean()

Picture	Target	Logits	Sigmoid	Loss
# 1	Teddy Bear (grizzly)	8 -5 -6	1 0.0 0	0. 0.01 0.
# 2	0 1 0	-4 6 3	0.0 1 0.5	0.02 0. 3.05
# 3	0 1 1	-8 5 1	0 0.99 0.03	0.001 0.31

positive targets: -log(sigmoid)

neg target: -log(1-sigmoid)

BCE Loss = 0.3777

PyTorch: nn.BCELoss (without sigmoid) or nn.BCEWithLogitsLoss

loss\_func = nn.BCEWithLogitsLoss()

loss = loss\_func(activs, targets)

Accuracy - single label | Accuracy - multi label

def accuracy(inp, targ, axis=-1):  
pred = inp.argmax(dim=axis)  
return (pred == targ).float().mean()

def acc\_multi(inp, targ, thresh=0.5,  
sigmoid=True):

if sigmoid: inp = inp.sigmoid()  
return ((inp > thresh) == targ.bool()).float().mean()

PARTIAL example (Python):

acc\_0.2 = partial(acc\_multi, thresh=0.2)

## REGRESSION

Example - Image regression, key point detection:

biwi = DataBlock(blocks=(ImageBlock, PointsBlock),  
splitter=FuncSplitter(lambda o: o.parent.name == 'B'),  
get\_items=get\_image\_files,

get\_y=get\_y,

item\_tfms=[ $\text{AugTransforms}(\text{size}=(240, 320))$ ,

NormalizeFromStats(\*biwi.stats)]

Flexible API + Transfer learning  
= POWER!

Rather than focusing on domains, focus on:

Independent Var | Dependent Var

Image	Text (caption)
Text	Image (from opt.)
Image + Text + Behavior	Product - Purchase-prob.

+ Loss function!

Finding the best threshold:

xs = torch.linspace(0.05, 0.95, 25)

accs = [acc\_multi(preds, targets,  
thresh=i, sigmoid=False)  
for i in xs]

plt.plot(xs, accs);

acc

↑ thresh  
ok to choose hyperparam based on valid set if the function looks smooth

learn = CnnLearner(dls, resnet15,  
y\_range=(-1, 1))

dls.loss\_func

→ MSELoss()

→ right function for regression!

! Pass y-range to learner to force outputs into range: def sigmoid\_range(x, lo, hi): return torch.sigmoid(x)\*(hi-lo)+lo

## Fastbook Ch. 7 | Training a SOTA Model

- ① If your dataset is big, experiment on a subset of it
  - iterate at faster speed
  - the more experiments you can do, the better
  - subset should be representative → generalize
- (Ex) Imagenette: 10 classes from Imagenet

Normalization  $\Rightarrow$  mean = 0, std dev = 1

- helps the model train
- especially important when using pretrained models - distributed with stats used for normalization, use them for inference or transfer learning

check:  $x, y = dls.one\_batch()$

•  $x.mean(dim=[0, 2, 3]), x.std(dim=[0, 2, 3])$

(average over all axes except channel = 1)

fastai: add Normalize transform in batch\_tfms

def get\_dls(bs, size):

dblock = DataBlock(blocks=(ImageBlock, CategoryBlock),

get\_items = get\_image\_files,

get\_y = parent\_label,

item\_tfms = Resize(460),

batch\_tfms = [aug\_transforms(size=size,

min\_scale=0.75),

Normalize.from\_stats(\*imagenet\_stats)])

return dblock.dataloaders(path, bs=bs)

## Progressive resizing

Gradually using larger & larger images as we train similar to transfer learning: finetune after resizing  
Works also as data augmentation

dls = get\_dls(128, 128)

# create learner, fit one-cycle

learn, dls = get\_dls(64, 224)

# learn.fine\_tune(epocs, lr)

May hurt performance in transfer learning if pretrained model / dataset are similar to our dataset.

## Test Time Augmentation (TTA)

During inference or validation, creating multiple versions of each image using data augmentation, and then taking the average or maximum of the predictions.

- (the default in fastai is center-cropping, for validation = largest square centered)
- problematic if relevant objects near edges
- squish/stretch may be difficult to handle
- TTA solves these problems

preds, targs = learn.tta() # default is unaugmented center crop + 4 randomly augmented images, applied on valid dset.

## Mixup

For each image img

- 1) select another dset image at random
- 2) pick a weight at random
- 3) take a weighted average of the img image and the selected image
- 4) take a weighted average of the img labels with the selected image's labels (targets need to be one-hot encoded)

In fastai: callbacks are used to inject custom behavior in training loop

model = xresnet50()

learn = Learner(dls, model, loss\_func = CrossEntropyLossFlat, metrics=accuracy, cb = Mixup) # callback

→ more epochs are needed to train for good accuracy (e.g. Imagenette LB - mixup in models  $\geq 80$  epochs)

→ can be used on activations inside models (NLP use cases)

## Label Smoothing (LS)

Problem with OHE: overconfidence, labels are always 0 or 1 even if there is nuance or uncertainty. Leads to overfitting, probabilities at inference not meaningful.

LS: replace all 1s with a number bit less than 1  
 $\text{--} 1 - 0.5 - \text{--} \text{more than } 0$

→ then train. Leads to:

- training more robust, even with noisy data
  - models that generalize better
- ① Start with OHE
  - ② Replace all 0s with  $\frac{\epsilon}{N}$  - no. of classes
  - ③ Replace 1 with  $1 - \epsilon$  and  $\frac{\epsilon}{N}$

In practice we don't change, or one-hot encode labels, but apply this in the loss function.

model = xresnet50()

learn = Learner(dls, model, loss\_func = LabelSmoothingCrossEntropy(), metrics=accuracy)

learn.fit\_one\_cycle(5, 3e-5)

→ more epochs are needed to train for good accuracy.

## Summary

① Establish your environment for quick iteration (experimentation

- subset of dataset
- validation generalizes to test/prod.

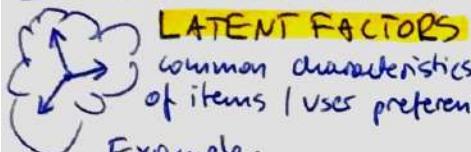
② Start with a simple, strong baseline

③ Run many experiments:

- Augmentation
- Lots functions
- Inspect detail results for insights

Check/read research papers

**COLLABORATIVE FILTERING DEEP GIVE** = look at what products the current user has used or liked, find other users that have used or liked similar products, then recommend other products that those users have used or liked. Generalize: items vs products, e.g. diagnoses, links etc.



Example:

$$\text{Movie is } [ \begin{matrix} \text{action} & \text{sci-fi} & \text{old} \\ 0.9 & 0.98 & -0.9 \end{matrix} ] \text{ A}$$

$$\text{User likes } [ \begin{matrix} 0.8 & 0.9 & -0.6 \end{matrix} ] \text{ B}$$

Dot Product  $A \cdot B =$  multiply vectors together, then sum up the result  
We don't know latent factors → need to learn them!

Embedding from scratch in PyTorch

Multiplying by a one-hot encoded matrix, using the computational shortcut that it can be implemented by simply indexing directly.

	$f_1$	$f_2$	$f_3$
$U_1$	1	2	3
$U_2$	3	1	2
$U_3$	1	1	2
$U_4$	2	2	3

$$U = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} \leftarrow \text{vector user3}$$

$$\leftarrow \text{matrix } F$$

$$F[3, :] \Rightarrow F^T * w$$

We index into the embedding matrix using an integer, but calculate the derivative as if we were multiplying the matrix with OLT vector

**Bootstrap NN problem** (cold start)  
→ pick some user to represent avg taste  
→ use a tabular model based on user metadata to construct initial embedding!  
! Representation bias & feedback loops are a risk - monitor keep humans in the loop, general / robust...

**dls = Colab's DotProduct**.from\_dls(ratings, item\_name=  
= 'title', bs=64)

$$n\_users = \text{len(dls.classes['user'])}$$

$$n\_movies = \text{len(dls.classes['title'])}$$

$$n\_factors = 50$$

Parameters in PyTorch - wrapper automatically calls requires\_grad\_, and initializes return nn.Parameter(torch.zeros(size).normal\_(0, 0.01))

**Class DotProductBasis (Module):** inherit from Module

**def \_\_init\_\_(self, n\_users, n\_movies, n\_factors, y\_range=(0, 5.5)):**

**self.user\_factors = create\_params([n\_users, n\_factors])**

**self.movie\_factors = Embedding(n\_movies, n\_factors)**

**self.user\_bias = create\_params([n\_users])**

**self.movie\_bias = Embedding(n\_movies, 1)**

**self.movie\_factors = create\_params([n\_movies, n\_factors])**

**self.movie\_bias = create\_params([n\_movies])**

**self.y\_range = y\_range** in PyTorch forward method

**def forward(self, x):** is called whenever module is called, passing along any parameters included in the cell

**users = self.user\_factors[x[:, 0]]**

**movies = self.movie\_factors[x[:, 1]]**

**res = (users \* movies).sum(dim=1)**

**res += self.user\_bias[x[:, 0]] + self.movie\_bias[x[:, 1]]**

**return SigmoidRange(res, \*self.y\_range)**

**PCA - principal component analysis**

pull out underlying directions in latent factor matrix

**movie\_w = learn.model.movie\_factors[idxs].cpu().detach()**

**movie\_pca = movie\_w.pca(3)**

**fac0, fac1, fac2 = movie\_pca.t()** (OOP in Python)

**x = fac0[movie\_ids]** \*\* leverages

**y = fac2[movie\_ids]** @ delegates decorator

→ plt scatter plot to visualize.

**Weight decay | L2 regularization**

Adding to our loss function the sum of all the weights squared, to encourage weights to stay small to prevent overfitting to wd | weight decay - parameter, how much we add loss: loss\_with\_wd = loss + wd \* (parameters\*\*2).sum() parameters.grad += wd \* 2 \* parameters - same

In fastai, we pass it in a cell to fit, learn.fit\_one\_cycle(5, 5e-3, wd=0.1)

**Embedding distance**

Movie similarity can be defined by the similarity of users that like those movies, distance between embedding vectors can define that similarity

**movie\_factors = learn.model.i\_weight.weight**

**idx = dls.classes['title'].map('MovieTitle1')**

**distances = nn.LosineSimilarity(dim=1)(movie\_factors, movie\_factors[idx][None])**

**idx = distances.argsort(descending=True)[1]**

**dls.classes['title'][idx]**

**Collab NN** - dot product was PMF (probabilistic matrix factorization) approach, there is an option to do it with deep learning!

**class CollabNN(Module):**

**def \_\_init\_\_(...)**  
...

**self.layers = nn.Sequential(**

**nn.Linear(user\_sz[1] + item\_sz[1], gn\_out),**

**nn.ReLU(),**

**nn.Linear(n\_out, 1))**

...

**def forward(self, x):**

**embs = self.user\_factors(x[:, 0]) self.item\_factors(x[:, 1])**

**x = self.layers(torch.cat(embs, dim=1))**

**return SigmoidRange(x, \*self.y\_range)**

Fastbook Ch 3 TABULAR MODELING Data as a table; predict value  
notes by @the21

## Variables

- continuous: numerical data, feed directly to model
- categorical: discrete levels (eg movie IDs), need to convert to numbers first

Ordinal - categories with natural ordering

df['ord\_cat'].cat.set\_categories([order, ordered=True, inplace=True])

CAT: represent via one-hot-encoding or entity embedding:  
→ reduces memory usage and speeds up NN vs DTE  
→ reveals intrinsic properties of variable - similar values  
close to each other in embedding space

Decision Trees TIP: avoid OHE categories for DTs / RFs

- 1) Loop through each column in dataset (greedy approach)
- 2) For each col, loop through each possible level of that col.
- 3) Try splitting data into 2 groups at that level
- 4) For regression: find any value of dependent var. for each of 2 groups, see how close it is to the actual value of dep. variable for each of the items in that group
- 5) After looping thru all cols / levels, pick split point with the best predictions
- 6) For each group based on this split, repeat the process

Random Forests (Breiman 1984, 2001)

- 1) Randomly choose a subset of rows and subset of columns
- 2) Train a model using this subset (decision tree)
- 3) Serve that model, return to step 1 several times
- 4) To make a prediction, predict using all served models  
then take average of those predictions ← BAGGING

Important - errors of individual models are not correlated,  
so the average of those errors is ZERO

Out of Bag Error: measure prediction error on the training set  
by only including in the calculation of a row's error the rows where that row was not included in training

BOOSTING. another approach to ensembling (vs BAGGING)

- 1) Train a small model that underfits your dataset
- 2) Calculate predictions in the training set for this model
- 3) Subtract predictions from targets = RESIDUALS
- 4) go back to step 1, now use the residuals as targets
- 5) continue until a stopping criterion: max no trees, valid error getting worse etc.

GBMs, GDTs, → risk of overfitting  
XGBoost → very sensitive to hyperparameter choices

TABULAR MODELING WORKFLOW

## Model

### INTERPRETATION

- 1) How confident are we in our prediction for a particular row?
- 2) What were the most important factors influencing prediction?

Start with RF - easiest to both res- train to hyperparameters little preprocessing

Use RF model for feature selection, PDP analysis

Then, try NNs or GBMs

try adding embeddings of cat. variables to the data

Which columns are effectively redundant with each other?

How do predictions vary as we vary each column?

### TREE VARIANCE

Check std deviation of predictions across trees:

preds = np.stack([t.predict(valid\_xs) for t in m.estimators\_])

preds\_std = preds.std(0)

high std ⇒ low confidence

### REMOVING LOW IMPORTANCE VARIABLES

Generally the 1st step to improve the model is simplifying it, so that it's easier to study, rollout maintainance →

Remove columns / variables of low importance

→ Retrain the model → check impact on accuracy

### REMOVING REDUNDANT FEATURES

cluster - columns (xs) - shows similarity of columns

We determine similarity by calculating the rank correlation - all values are replaced by their rank, then correlation is calculated

→ try removing each of potentially redundant features, one at a time then multiple from overlapping groups, observe OOB score / accuracy

### TREE INTERPRETER + WATERFALL CHARTS

What were the most imp factors for predicting with a particular row of data? How did they influence that prediction? Calculated similarly to feature importances.

→ display feature contributions with waterfall chart

→ use it to provide useful information to users of your data product - reasons behind predictions

VS Extrapolation - Decision Trees / RF can never

predict values outside of the range of training data (e.g. trend). NN can help. Also finding

out-of-distribution data: 1) combine train / valid sets

2) Use RF to predict if a row is in train or valid set

3) Get feature importances - for the columns that

differ significantly between train / valid by removing them and see how it impacts accuracy. It can

improve it, and make the model more robust.

→ Try to avoid using old data - it may no longer be predictive.

FEATURE IMPORTANCE 3 in. feature importances -  
Loop through each tree, then recursively explore  
each branch - check what feature was used for  
that split, and how much the model improves  
as the result. The improvement weighted by  
number of rows in that group is added to  
importance score for that feature. Sum across  
all branches of all trees, normalize scores so  
that they add to 1.

### PARTIAL DEPENDENCE PLOTS = PDP

PDP Visualize how variables effect our predictions: if a row varied as nothing other than the feature in question, how would it impact the dependent variable?

- replace every value in Year\_Made col. with 1950
- calculate predicted sale price for every auction, take average over all auctions
- Repeat for 1951, 1952, ..., 2020
- This isolates the effect of Year\_Made from sklearn.inspection import The Book plot - partial\_dependence

### DATA LEAKAGE

= giving model information about the target which normally should not be available at the time of prediction. How to detect it?

- check if the accuracy is too good to be true
- look for imp predictors that don't make sense
- look for PDP plots that don't make sense in practice

### Using a Neural Network

- 1) Decide which cols should be treated as cat vs cont
- 2) Create embeddings for categorical variables
- 3) Add normalization (in proc)
- 4) Consider adding y-range to regression models
- 5) Adjust hidden layer sizes to size of dataset

fastai: TabularPandas (handles df + convenience)

Tabular Proc Tabular\_learner Other:

Categorical TabularModel dtreeviz

Fill Missing Study source code add\_dataloader



**Language Model:** model trained to guess the next word in a text, after reading the words before

**Self-supervised learning:** training a model using labels that are embedded in the independent variable, rather than requiring external labels. Usually used for pretraining in transfer learning.

## TEXT PREPROCESSING IN 3 STEPS

### ① TOKENIZATION

= converting a text into a list of tokens

a) **word-based:** split sentence on spaces and language-specific rules

b) **subword-based:** analyze a corpus of documents to find the most commonly occurring groups of letters. These substrings become the vocab.

c) **character-based**

fastai adds some functionality with the Tokenizer class, e.g. special tokens:

xxbos : beginning of text

xxusej : next word begins with a capital

xxunk : next word is unknown

• see rules: defaults.text\_proc\_rules

Setup creates the vocab that is used in tokenization

Spacy = WordTokenizer() [txt] - collection of text documents

toks = first(spacy([txt])) displays first n print(coll\_repr(toks, 30)) elems of collection

tkn = Tokenizer(spacy) adds special print(coll\_repr(tkn(txt), 31)) tokens etc

sp = SubwordTokenizer(vocab\_sz=52)

sp.setup(texts)

### ② NUMERICALIZATION

= mapping tokens to integers

- 1) make a list of all possible levels of a variable (the vocab)
- 2) replace each level with its index in the vocab

num = Numericalize()  
num.Setup(tokens)  
coll\_repr(num, vocab, 20)

class methods

### ③ PUTTING TEXT INTO BATCHES for a language model

- we want LM to read text in order
- we use a model that maintains a state - remembers what it read previously when predicting what comes next

- 1) transform indiv. texts into a stream by concatenating them, shuffle docs order before each epoch
- 2) cut this stream into a certain number of batches (batch size)  
mini-streams preserve order of tokens
- 3) Each time step read seq-len from mini-streams

- dependent variable is offset from the independent variable by 1 token // LM Data Loader

⚠ Potential to generate disinformation campaigns, flood social media with fake content etc → see ch. 3 on ethics

collecting items in a batch

- use padding to make texts all the same size
- sort(ish) docs by length, prior to each epoch - batch together docs with similar lengths, pad to length of longest doc in a batch

get\_imdb = partial(get\_text\_files, folders=['train', 'test', 'unsup'])

ds\_lm = DataBlock(  
when fastbook passed fastai handles tokenization and numericalization)

blocks = TextBlock().from\_folder(path, is\_lm=True),

get\_items = get\_imdb, splitter = RandomSplitter(0.1),  
.dataloaders(path, path=path, bs=128, seq\_len=80)

dls\_lm.show\_batch(max\_n=2) Words that are not in the  
vocab of pretrained lm will be added with random embeddings

### # fine-tuning language model

learn = LanguageModelLearner(dls\_lm, AWD\_LSTM,  
drop\_mult=0.3, metrics=[accuracy, Perplexity()]).to\_fp16()

loss function: cross entropy

perplexity metric: exponential of  $\text{torch.exp}(\text{cross_entropy})$

learn.fit\_one\_cycle(1, 2e-2) automatically frozen, dls\_lm will only train embeddings

learn.save('1epoch') use fit-one-cycle to save/load intermediate model results

learn.unfreeze() learn.fit\_one\_cycle(10, 2e-3) ENCODER: model without task-specific final layers, like body in CNNs

learn.save\_encoder('finetuned') TEXT GENERATION

TEXT = 'I liked this movie because'  
N\_WORDS = 40 pred = learn.predict(TEXT, N\_WORDS, temperature=0.25)

if creating the classifier dataloader ! is\_lm=False

dls\_clas = DataBlock(  
pass vocab to use fine-tuned encoder)

blocks = (TextBlock().from\_folder(path, vocab=dls\_lm.vocab), CategoryBlock),  
get\_y = parent\_label

get\_items = partial(get\_text\_files, folders=['train', 'test']),  
splitter = GrandparentSplitter(valid\_name='test')

, dataloaders(path, path=path, bs=128, seq\_len=72)

learn = TextClassifierLearner(dls\_clas, AWD\_LSTM, drop\_mult=0.5,  
metrics=accuracy).to\_fp16()

learn = learn.load\_encoder('finetuned') + gradual unfreezing

learn.fit\_one\_cycle(1, 2e-2) + discriminative learning rates

learn.freeze\_to(-2)

learn.fit\_one\_cycle(1, slice(1e-2 / (2.6 \*\* 4), 1e-2))

... (freeze\_to(-3)) ...

learn.unfreeze()

learn.fit\_one\_cycle(2, slice(1e-3 / (2.6 \*\* 4), 1e-3))

Fastbook Ch. 11 Data Munging with  
notes by @dh21 [ Data Munging with  
fastai's Mid-Level API ]

Fastai is built on a layered API:

Top layer = applications - train a model  
in 5 lines of code, eg:

TextDataLoaders.from\_folder()

- Mid level API:
- create new DataLoaders
  - apply just part of transforms
  - has the callback system, which allows to customize training loop any way we like
  - has general optimizers

Transform: an object that behaves like a function, has optional setup method to initialize hidden state (eg vocab) and an optional decode that reverses the function. Special behavior - always gets applied over tuples (input, target)

Examples: Tokenizer, Numericalize

- 1) Create object
  - 2) well setup method
  - 3) apply to input by calling object as a function
  - 4) decode result back to understandable representation
- useage, pushes encap-sulates these steps in the Transform class

Write your own Transform

1) Write a function / + decorator

def f(x: int): return x+1

function only gets applied to ints

tfn = Transform(f)

2) Decorator - Python syntax for passing a function to another function (or callable) so that behaves like a function

@Transform

def f(x: int): return x+1

3) If we need setup or decode, then need to subclass Transform and implement encode

FROM: → → → →

path = untar\_data(URLs.IMDB)

dls = DataBlock(

blocks=(TextBlock.from\_folder(path),  
CategoryBlock))

get\_y = parent\_label,

get\_items = partial(get\_text\_files,  
folders=['train', 'test'])

splits = GrandParentSplitter(valid\_name='test')(files)

dsets = Datasets(files, tfms, splits=splits)

dls = dsets.dataloaders(dl\_type=SortedDL,

before\_batch=pad\_input)

TO: (equivalent, but can be customized):  
tfms = [[Tokenizer.from\_folder(path), Numericalize],  
[parent\_label, Categorize]]

files = get\_text\_files(path, folders=['train', 'test'])

Splits = GrandParentSplitter(valid\_name='test')(files)

dsets = Datasets(files, tfms, splits=splits)

dls = dsets.dataloaders(dl\_type=SortedDL,  
before\_batch=pad\_input)

Datasets: apply 2 (or more) Pipelines in parallel to the same raw object and build a tuple with the result. → automatically do setup for us → when indexed, return a tuple with result of each pipeline!

x\_tfms = [Tokenizer.from\_folder(path), Numericalize]

y\_tfms = [parent\_label, Categorize]

dsets = Datasets(files, [x\_tfms, y\_tfms], splits=splits)

x, y = dsets.valid[0]

→ convert it to dataloader (here: need to pad input):

dls = dsets.dataloaders(bs=64, before\_batch=pad\_input)

DataLoader: collects the items from our datasets into batches. Input ways to customize:

1) after\_item: applied on each item after grabbing it inside the dataset (n\_item\_tfms in DataBlock)

2) before\_batch: applied on list of items before they are collated. Ideal place to pad items to the same size.

3) after\_batch: applied on the batch as a whole after its construction (n\_batch\_tfms)

Application Example: Siamese Pair  
Siamese model takes two images and has to determine if they are same class or not.

Pipeline: compose several Transforms together

tfms = Pipeline([t1, num]) define it by passing a list of transforms

t = tfms(txt) automatically applies transforms

tfms.decode(t) decode result of encoding

TfmDLists and Datasets, Transformed Collections

cut = int(len(files)\*0.8)

splits = [list(range(cut)), list(range(cut, len(files)))]

tls = TfmDLists(files, [Tokenizer.from\_folder(path),  
Numericalize], splits=splits)

t = tls[0] TfmDLists = class that groups

tls.decode(t) data as a set of raw items

tls.show(t) (filenames, dt rows...) and Pipeline of Transforms. At initiation, well setup in each Transform in order. We make into it to get results

of pipeline on raw elements. Can handle train/valid: tls.valid[0]. Convert it to dataloader w/ dataloader method

Fastbook ch 12 | A language model from scratch | TIP: When starting to work on a new problem, find the simplest defenses notes by @dk21 | TIP: to try out methods quickly & easily, interpret results. Eg.: Human Numbers

## #First language model:

class LMModel1 (Module):

def \_\_init\_\_(self, vocab\_size, n\_hidden):

self.i\_h = nn.Embedding(vocab\_size, n\_hidden)

self.h\_h = nn.Linear(n\_hidden, n\_hidden)

self.h\_o = nn.Linear(n\_hidden, vocab\_size)

def forward(self, x):

$h = F.relu(self.i_h(x[:, 0]))$

$h = h + self.i_h(x[:, 1])$

$h = F.relu(self.h_h(h))$

$h = h + self.i_h(x[:, 2])$

$h = F.relu(self.h_h(h))$

return self.h\_o(h)

learn = Learner(dls, LMModel1(len(vocab), 64),

loss\_func=F.cross\_entropy, metric=accuracy)

learn.fit\_one\_cycle(5, 1e-3)

# compare vs naive baseline: always predict

the most common token

RNN

① Refactor with a for-loop  $\Rightarrow$  Recurrent NN

= we can now apply to token sequences of different lengths.

Hidden state = activations that are updated at each step of a recurrent NN

② Maintaining the State of RNN

- don't initialize hidden state to zero for each sample (move initialization of hs to \_\_init\_\_)

- BUT: that makes our NN as deep as the number of tokens (consider unrolled representation)

- calculating derivatives will be slow / memory intensive

- SOLUTION: tell Pytorch we don't want to calculate

all gradients just n steps. Use detach method to remove gradient history.

Backpropagation Through Time

Treating a NN with effectively one layer per time step, refactored using a loop, as one big model, and calculating gradients on it in the usual way. To avoid running out of memory and time, we use truncated BPTT, which detaches the history of computation steps in the hidden state every few time steps.

Make sure the model sees text samples in correct order!

class LSTMCell (Module):

def \_\_init\_\_(self, ni, nh):

self.forget\_gate = nn.Linear(ni+nh, nh)

self.input\_gate = nn.Linear(ni+nh, nh)

self.cell\_gate = nn.Linear(ni+nh, nh)

self.output\_gate = nn.Linear(ni+nh, nh)

def forward(self, input, state):

①  $h, c = \text{state}$

②  $h = \text{torch.cat}([h, \text{input}], \text{dim}=1)$

③ forget = torch.sigmoid(self.forget\_gate(h))

④  $c = c * \text{forget}$

⑤  $ihp = \text{torch.sigmoid}(\text{self.input_gate}(h))$

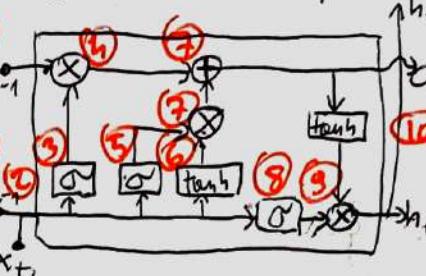
⑥ Cell = torch.tanh(self.cell\_gate(h))

⑦  $c = c + ihp * \text{Cell}$

⑧ out = torch.sigmoid(self.output\_gate(h))

⑨  $h = \text{out} * \text{torch.tanh}(c)$

⑩ return h, (h, c)



# Final language model

class LMModel2 (Module):

def \_\_init\_\_(self, vocab\_size, n\_hidden, n\_layers):

self.i\_h = nn.Embedding(vocab\_size, n\_hidden)

self.rnn = nn.LSTM(n\_hidden, n\_hidden, n\_layers)

self.drop = nn.Dropout(p) batch\_first=True

self.h\_o = nn.Linear(n\_hidden, vocab\_size)

self.h\_o.weight = self.i\_h.weight

self.h = [torch.zeros(n\_layers, bs, n\_hidden)

for i in range(2)]

def forward(self, x):

raw\_h = self.rnn(self.i\_h(x), self.h)

out = self.drop(raw\_h)

self.h = [h.detach() for h in self.h]

return self.h\_o(out), raw\_h, out

def reset(self):

for h in self.h: h.zero\_()

learn = Learner(dls, LMModel2(len(vocab), 64, 2, 0.5))

loss\_func = crossEntropyLossFlat() metric=accuracy

Cbs = ModelResetter.RNNRegularizer(alpha=2, beta=1)

learn.fit\_one\_cycle(15, 1e-2, wd=0.1)

⑥ Regularizing an LSTM

DROPOUT: randomly change some activations

to zero at training time, so that all neurons activate more towards the output.

This changes scale of activations, so we need to rescale them (divide by 1-p)

class Dropout(Module):

def \_\_init\_\_(self, p): self.p = p # set by model.train

def forward(self, x):

if not self.training: return x # zeros and ones

mask = x.new(\*x.shape).bernoulli\_(1-p)

return x \* mask.div\_(1-p)

Activation Regularization / Temporal Act Reg (TAR)

Like weight decay, we want to make final activations from LSTM as small as possible.

loss += alpha \* activations.pow(2).mean()

Temporal AR: keep diff between consecutive token

activations small

loss += beta \* (acts[:, :-1] - acts[:, :-2]).pow(2).mean()

this + other tricks (weight tying) in AWD-LSTM

③ Creating more signal

= predict the next word after every input word

dependent variable = input shifted by 1

model will return outputs of shape batchsize\*seqlen

targets are in shape batchsize\*seqlen \* vocab\_size

flatten with tensor.view, use F.cross\_entropy.

④ Multi-layer RNNs

repeat for 1 to n-1

⑤ Exploding/vanishing activations

many layers  $\Rightarrow$

many matrix multiplications

extremely large or small numbers

floating point numbers are an approximation

gradients end up as zero or infinity

model stops training

SOL's = batch norm / residual connections / careful init

for RNN! GRU / LSTM layers!

LSTM adds another hidden state (cell state)

to remember what happened earlier in a sequence

(long - short term memory)

Feature Engineering: creating new transformations of the input data in order to make it easier to model.  
eg: add\_offset, images: edge detectors

Convolution: applying a kernel across an image (multiplication + addition)  
→ activation map

```
top-edge = tensor([[-1, -1, -1],
                   [0, 0, 0],
                   [1, 1, 1]]).float()
```

```
def apply_kernel(row, col, kernel):
    return (img[row-1:row+2, col-1:col+2]*kernel).sum()
```

Pytorch  $\Rightarrow$  F. conv2d

Input tensor shape:

(minibatch, n-channels, n-rows (height), n-cols (width))

Weight tensor = kernels/filters of shape:

(out-channels, in-channels, n-rows/kernel, n-cols/kernel)

Padding - we use it to ensure that the output activation map is the same size as original image

Stride - we use it to decrease the size of our outputs, how much we move kernel at each step

CNN → we learn the values of the kernels (eg SGD)

→ we use convolutions in the network architecture

Channels / Features - size of 2nd axis of weight matrix, ie number of activations per grid cell after a convolution. Channels can also refer to input data (n colors)

Receptive Field - area of an image that is involved in the calculation of a layer

- the more stride 2 convolutions we have before a layer, the larger the receptive field for an activation in that layer, also more semantically rich features, so we need more weights / parameters (channels)

Color Images: 3 channels: R | G | B : in one sliding

Input kernel output window, on each channel, we multiply elements of window by elements of corresponding filter, then sum the results and sum over all filters

$$\text{kernel tensor shape} = \text{ch-in} \times \text{ksize} \times \text{ksize}$$

```
def conv(ni, nf, ks=3, act=True):
    layers = [nn.Conv2d(ni, nf, stride=2, kernel_size=ks, padding=(ks-1)//2)]
    if act:
        layers.append(nn.BatchNorm2d(nf))
    layers.append(nn.ReLU())
    return nn.Sequential(*layers)
```

```
Simple_CNN = Sequential(
    Conv(1, 4), # 1x1x1
    Conv(4, 8), # 7x7
    Conv(8, 16), # 4x4
    Conv(16, 32), # 2x2
    Conv(32, 2, act=False), # 1x1
    Flatten())
```

Same as tensor.squeeze, but as a module

nn.Conv2d - module equivalent of F.conv2d, more convenient when creating an architecture because it creates the weight matrix automatically whenever we instantiate it.

learn = Learner(dls, simple\_cnn, loss\_func=F.cross\_entropy, metrics=[accuracy])

learn.summary() / Gibbons:

learn.activation\_stats, color\_dim(-2)

→ another way to see training stability

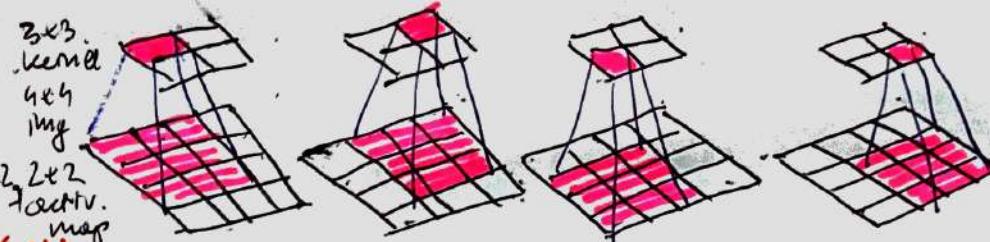
→ observe for change in non-zero activations

increased collapse = bad

smooth increase = good / PYTHON

Nested list comprehension  
[[i,j) for j in range(1,5)] for i in range(1,5)]

use unsqueeze to insert a unit axis into a tensor at specific position



## Improving Training Stability

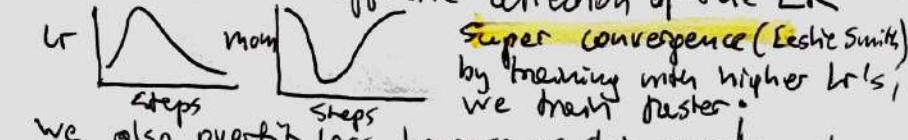
Let's use ActivationStats callback to look inside our models when they're training from fresh. callbacks import ... fit ... obs = ActivationStats(with\_hist=True)

learn.activation\_stats.plot\_layer\_stats(0) → model should have consistent / smooth mean and std of layer activations during training → activations near zero are especially problematic - computation is doing nothing, zeros propagate to further layers

① Increase Batch Size: - gradients are more accurate (from more BVT: fewer batches per epoch, so fewer chances to update weights)

② 1 cycle training: initial weights are not suited well to our task, so we don't want a high LR - training would diverge quickly. We also don't want high LR at the end, so we don't skip over the minimum → we separate learning rate schedule in 2 phases:

- 1) LR grows from min to max value (warm up)
- 2) LR decreases back to min value (annealing)
- 3) the momentum (impact of previous steps / direction in SGD) varies in the opposite direction of the LR



We also overfit less, because we skip over the sharp local minima to end up in smoother part of the loss.

③ Batch Normalization nn.BatchNorm2d (int)

Take average of the mean and std devs of the activations of a layer and using those to normalize the activations

Potential problem: some acts need to be high to predict well, so we add 2 learnable parameters, gamma and beta /

y - activation vector after normalization

BN returns gamma \* y + beta

During training - use mean and std dev of the batch to normalize data. During Validation - use running mean of statistics calculated during training. → training is more smooth, also works as regularizer, model generalizes better

# Freshbooks Ch 14 | ResNets

notes by @dhu21

Before 2013 (incl. VGG): Sequence of conv layers, then flatten activation matrix into a vector

## FCN - Fully Convolutional Network:

Sequence of conv layers, some of them stride 2, followed by an adaptive average pooling layer, a flatten layer to remove the last axes, and finally a linear layer.

`def avg_pool(x): return x.mean((2,3))`  
Pytorch: nn. Adaptive Avg Pool 2d  
→ converts a grid of activations into a single activation per image

`def block(ni, nf): return ConvLayer(ni, nf, stride=2)`  
`def get_model():`  
    return nn.Sequential(  
        block(3, 16),  
        block(16, 32),  
        block(32, 64),  
        block(64, 128),  
        block(128, 256),  
        nn.AdaptiveAvgPool2d(1),  
        Flatten(),  
        nn.Linear(256, dls.c))

? Is this good for OCR, eg MNIST? No!  
We slice image into small pieces, jumble them up, and decide whether on average each piece looks like our target class...

FCN - good for objects that don't have a single correct orientation or size, eg. photos

2015: Kaiming He et al. Deep Residual Learning...

Insight: Deeper NN's led to worse performance

vs shallower nets, both train and test.

## SCIENTIFIC DISCOVERY

- 1) Start with experimental observation
- 2) Watch out for unexpected results
- 3) Try to figure it out, with tenacity

LOSS LANDSCAPE

ResNet

Skip Connection: skip over every second convolution, so:

$$x + \text{conv2}(\text{conv1}(x))$$

ResNet block returns  $y = x + \text{block}(x)$

We are asking the block to predict the difference between  $x$  &  $y$  = residual.  
Later trick: use zero for the initial value of genuine in the last batch norm layer of each block - we begin training with a true identity path through the ResNet blocks.

How to handle different dims (grid size, channels)?

1x1 Convolution: Conv with a kernel size of 1, only doing a dot product over the channels of each input pixel, not combining pixels

`def conv_block(ni, nf, stride=1):`

return nn.Sequential(  
    ConvLayer(ni, nf, stride=stride),  
    ConvLayer(nf, nf, act\_cls=None, norm\_type=  
                    NormType.BatchNormZero))

Class ResBlock (Module):

`def __init__(self, ni, nf, stride=1):`

    Self.convs = conv\_block(ni, nf, stride)  
    Self.idconv = noop if ni == nf else  
        ConvLayer(ni, nf, 1, act\_cls=None)  
    Self.pool = noop if stride == 1 else  
        nn.AvgPool2d(2, ceil\_mode=True)

`def forward(self, x):`

    return F.relu(self.convs(x) + self.idconv(self.pool(x)))

# noop - no operation, return unchanged input

→ First layers: majority of the computation, few parameters. Eg 128 px image; stride 1: apply kernel to each of 128x128 pixels

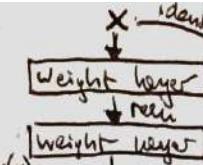
$3 \times \text{input Hs} \times 32 \text{ output Hs} \times 3 \times 3 \text{ kernel} = 864 \text{ params}$

Last layers: majority of parameters, few computations

$256 \text{ input Hs} \times 512 \text{ output Hs} \times 3 \times 3 \text{ kernel} = 1,175,648 \text{ params}$

Apply kernel to grid size 4x4 or 2x2

→ keep first layers as simple as possible (STEM)



def \_resnet\_stem(stems):

return [

    ConvLayer(sizes[i], sizes[i+1], 3, stride=2 if i == 0 else 1)  
        for i in range(len(sizes) - 1)

] + [nn.MaxPool2d(kernel\_size=3, stride=2, padding=1)]

class ResNet(nn.Sequential):

`def __init__(self, n_out, layers, expansion=1):`

    stem = \_resnet\_stem([3, 32, 32, 64])

    Self.block\_S2s = [64, 64, 128, 256, 512]

    for i in range(1, 5): Self.block\_S2s[i] \*= expansion  
    blocks = [self.\_make\_layer(i) for i in enumerate(layers)]

super().\_\_init\_\_(stem, \*blocks)

    nn.AdaptiveAvgPool2d(1), Flatten(),

    nn.Linear(self.block\_S2s[-1], n\_out))

`def _make_layer(self, idx, n_layers):`  
    stride = 1 if idx == 0 else 2

    ch\_in, ch\_out = self.block\_S2s[idx : idx + 2]

    return nn.Sequential(\*[  
        ResBlock(ch\_in if i == 0 else ch\_out, ch\_out,  
                    stride if i == 0 else 1),  
                    for i in range(n\_layers)])

`def conv_block(ni, nf, stride):`

return nn.Sequential(  
    ConvLayer(ni, nf // 4, 1),  
    ConvLayer(nf // 4, nf // 4, stride=stride),  
    ConvLayer(nf // 4, nf, 1, act\_cls=None, norm\_type=NormType.BatchNormZero))

# ResNet 18: m = ResNet(dls.c, [2, 2, 2, 2])

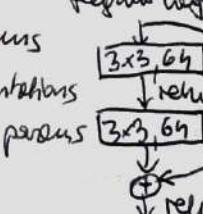
# ResNet 50: pass 4 as expansion parameter,  
start with 4x fewer channels, end with 4x more  
m = ResNet(dls.c, [3, 4, 6, 3], 4)

dls = get\_data(URLs.IMAGENETTE\_320, pretrue=320, resize=224)

learn = get\_learner(m)

learn.fit(2, one\_cycle(20, 3e-3)) # 100 epochs, ...

Regular layers:



Bottleneck layers: 1x1 convs are much faster! This block executes faster so we can add more filters in the same amount of time. 1x1 conv first diminishes, then restores the no. of channels.

CNN-learner

- pass an architecture to use as the body of the network (usually ResNet)
- download and load pretrained weights
- for transfer learning, cut the head ResNet: everything from Adaptive Avg Pooling layer
- use dictionary to determine head vs body: model - head [resnet50] = {'cut': -2, 'split': <function ...>, 'stats': [[...], [...]]}
- create new head: create\_head(20, 2) Sequential(...)
- Adaptive ConcatPool2d()
  - (op): AdaptiveAvgPool2d(output\_size=1)
  - (inp): AdaptiveMaxPool2d(output\_size=1)
- used by default in fastai CNN head
- fastai adds 2 linear layers by default in CNN head (helps transfer learning across domains)
- bn\_final parameter: use batchnorm as the final layer to help the model scale output acts.

UNet-learner

- used for segmentation, super-resolution, colorization, style-transfer
- output is an image/pixel grid/predicted label for each pixel (generative vision models)
- same approach as CNN-cut off head, replace it with a new one for the generation branch
- How to increase grid size?
  - a) nearest neighbor interpolation (Pytorch layer) # So-so
  - b) transposed convolution / stride half convolution (insert zero padding between all the pixels in the input)
  - c) add skip connections - from the activations in the body of the ResNet (U-net)
- fastai Dynamic UNet - autogenerates architecture with the right size for our image size

Practical DL workflow:

- 1) Start with a small model and data sample, get to the point where you overfit
- 2) Add more data
- 3) Data augmentation
- 4) Generalizable architectures
- 5) Regularization
- 6) Reduce architecture complexity

Custom Model: a Siamese network

```
class SiameseModel(Module):
    def __init__(self, encoder, head):
        self.encoder, self.head = encoder, head
    def forward(self, x1, x2):
        ftrs = torch.cat([self.encoder(x1),
                         self.encoder(x2)], dim=1)
        return self.head(ftrs)
```

```
encoder = create_body(resnet34, cut=-2)
```

```
head = create_head(512*4, 2, ps=0.5)
```

```
model = SiameseModel(encoder, head)
```

```
# Convert boolean targets to integers in loss fn:
def loss_func(out, targ):
    return nn.CrossEntropyLoss()(out, targ.long())
```

```
# Custom Splitter: how to split the model
```

into parameter groups, so we can train only the head of the model in transfer:

```
def siamese_splitter(model):
```

```
    return [params(model.encoder), params(model.head)]
```

```
learn = Learner(dls, model, loss_func=loss_func,
```

```
Splitter=SiameseSplitter, metrics=accuracy)
```

```
# need to call freeze/unfreeze manually
```

NLP (and LSTM)

- Select the stacked RNN for the encoder in the language model (= activation for every word in the input)
- To create a classifier, we use "BPTT for text clif":  
 → divide a doc into fixed-length batches of size b  
 → initialize each batch with final state of prev batch  
 → keep track of hidden states for mean and max padding  
 → gradients are back-propagated to the batches whose hidden states contributed to the final prediction
- classifier contains a for-loop over each batch of a sequence
- we maintain state and store activations of each batch
- at the end, we use avg and max concatenated padding over RNN sequences
- each text needs to be treated separately (because of various lengths, we add padding)  
 (After randomness, so texts of similar size are batched together)

Encoder

forward method:  
 if self.n\_emb != 0:  
 $x = [e(x\_cat[:, i]) \text{ for } i, e \in \text{enumerate}(\text{self.embds})]$   
 $x = \text{torch.cat}(x, 1) \text{ # concat embeddings into a single tensor}$   
 $x = \text{self.emb\_drop}(x) \text{ # dropout}$   
 if self.n\_cont != 0:  
 $x\_cont = \text{self.bn\_cont}(x\_cont) \text{ # norm}$   
 $x = \text{torch.cat}([x, x\_cont], 1) \text{ if self.n\_emb != 0 else } x\_cont$   
 return self.layers(x)

## THE TRAINING PROCESS

### OPTIMIZERS & CALLBACKS

Baseline: `def get_learner(**kwargs):`

```
return nn.Learner(dls, model=34, pretrained=False,
    metrics=accuracy, **kwargs).to_fp16() # turn off momentum
```

SGD: `learner = get_learner(opt_func=SGD)`  
`learner.fit_one_cycle(3, 0.03, mom=(0, 0, 0))`

**Generic Optimizer**: 2 key methods:

```
def zero_grad(self):
    for p, _ in self.all_params():
        p.grad.detach_() # remove history of gradient computation
        p.grad = zero_()
```

`def step(self):`

```
for p, pg, state, hyper in self.all_params():
    for cb in self.callbacks: # loop through callbacks
        state = _update(state, cb(p, **{**state, **hyper}))
```

`self.state[p] = state` calls the callbacks to update params

# SGD callback:

```
def sgd_cb(p, lr, **kwargs): p.data.add_(-lr, p.grad.data)
```

# Pytorch add\_ with 2 parameters → multiply them before adding

`opt_func = partial(Optimizer, cb=[sgd_cb])` # SGD

### Momentum

Iron ball rolling down a mountain will pass over bumps and holes, a ping pong ball is more likely to always follow direction of the gradient and get stuck (heavier = more momentum)

→ use moving average instead of only the current gradient

`weight_avg = beta * weight_avg + (1 - beta) * weight`

`new_weight = weight - lr * weight_avg`

`beta`: how much momentum we use (default: 0.9, SGD: 0)

`weight_avg` → we need to store moving avg for each param!

`fit_one_cycle momentum`:

$0.95 \rightarrow 0.85 \rightarrow 0.95$



With large beta, we might miss that the gradients changed direction and roll over small local minima.

This is actually a better place!

New inputs will hit the beginning, unbiased avg looks more like gradients after more steps similar to moving average.

$w_{avg} = beta * w_{avg} + (1 - beta) * w_{grad}$  is like  $w_{unbiased\_avg} = w_{avg} / (1 - (beta * (i + 1)))$  # ratio depends:  $beta_1: 0.9$   $beta_2: 0.999$   $eps: 1e-8$ .  $beta_1: 0.9$   $beta_2: 0.999$   $1e-5$   $eps$  also limits max value of adjusted lr!

**Decoupled Weight Decay**: 2 formulations:

```
def average_grad(p, mom, grad_avg=None, **kwargs):
    if grad_avg is None: grad_avg = torch.zeros_like(p.grad.data)
    return {grad_avg: grad_avg + mom + p.grad.data}
```

```
def momentum_step(p, lr, grad_avg, **kwargs):
    p.data.add_(-lr, grad_avg)
```

`opt_func = partial(Optimizer, cb=[average_grad, momentum_step], mom=0.9)`

### RMS Prop (Istvan)

Adaptive learning rate: Each param gets its own lr, controlled by global lr:

- if gradients close to zero for a while → the loss is flat → need higher lr
- if gradients all over the place → risk of divergence → pick lower lr
- we determine this by using moving average of gradients squared (low avg → higher lr)

$w_{square\_avg} = alpha * w_{square\_avg} + (1 - alpha) * (w_{grad}^2)$   
 $new_w = w - lr * w_{grad} / \text{math.sqrt}(w_{square\_avg} + eps)$   
 current gradient divided by square root of moving avg

```
def average_sqr_grad(p, sqr_mom, sqr_avg=None, **kwargs):
    if sqr_avg is None: sqr_avg = torch.zeros_like(p.grad.data)
    return {sqr_avg: sqr_avg * sqr_mom + (p.grad.data**2) * (1 - sqr_mom)}
```

```
def rms_prop_step(p, lr, sqr_avg, eps, grad_avg=None, **kwargs):
    denom = sqr_avg * sqrt(1 - (1 - sqr_mom))
    p.data.addcdiv_(-lr, p.grad, denom)
```

`opt_func = partial(Optimizer, cb=[average_sqr_grad, rms_prop_step], sqr_mom=0.99, eps=1e-7)`

**Adam**: SGD with momentum + RMSProp together!

→ moving avg of gradient as direction difference! using unbiased average

→ adaptive lr for each parameter

### CALLBACKS

Make changes to the training process without changing the training loop!

Callback: a piece

of code we can inject into another piece of code at a preferred point

Training loop: for  $xb, yb$  in  $dl$ :

```
try:
    self._split(b);
    self._pred = self.model(*self.xb); # after-pred
    self._loss = self.loss_func(self._pred, *self.yb); # after-loss
    if not self.training: return
    self._backward(); # after-backward
    self._opt.step(); # after-step
    self._opt.zero_grad(); # after-batch
except CancelBatchException: # after-cancel-batch
    finally: # after-batch
```

- callbacks can access attributes of Learner and modify them
- callbacks can raise an Exception, e.g. cancel training if the loss becomes Inf or NaN
- callbacks can be called in particular order (specify: run\_before, run\_after)

Firstbook ch 17 | A Neural Net from The Foundations  
 notes by @dru1 | PyTorch Basics

**The forward and backward passes**  
 $x = \text{torch.randn}(200, 100); y = \text{torch.randn}(200)$   
 We have to scale our weights so that activations have mean = 0 and std = 1. std < 0 or > 0 will lead to 0 or inf in our activations after many multiplications.  $\rightarrow$  Initialization

**PyTorch basics**  
 A neuron:  $\text{out} = \sum_{i=1}^n x_i w_i + b$  adds an bias to the sum of weighted inputs

output =  $\text{sum}([\mathbf{x} * \mathbf{w} \text{ for } \mathbf{x}, \mathbf{w} \text{ in } \text{zip}(\text{inputs}, \text{weights})]) + \text{bias}$   
 We need to feed this into a non-linear activation function, e.g.  
 $\text{def relu}(x): \text{return } x \text{ if } x \geq 0 \text{ else } 0$

**Fully connected (dense, linear) layer**: all inputs are linked to all neurons via a dot product.  
 $\mathbf{x}$ : inputs matrix (batch-size x n-inputs)  
 $\mathbf{w}$ : weights matrix (n-neurons x n-inputs)  
 $\mathbf{b}$ : biases vector (n-neurons)

$y = \mathbf{x} @ \mathbf{w} + \mathbf{b}$  (batch-size x n-neurons)  
 Output of a fully connected layer

**Matrix multiplication from scratch**  
 $\text{def matmul}(\mathbf{a}, \mathbf{b}):$   
 $\quad \mathbf{ac} = \mathbf{a}.\text{shape} \# \mathbf{n\_rows} \times \mathbf{n\_cols}$   
 $\quad \mathbf{bc} = \mathbf{b}.\text{shape}$   
 $\quad \text{assert ac == bc}$   
 $\quad \mathbf{c} = \text{torch.zeros}(\mathbf{ac}, \mathbf{bc})$   
 $\quad \text{for i in range(ac):}$   
 $\quad \quad \text{for j in range(bc):}$   
 $\quad \quad \quad \mathbf{c}[i,j] += \mathbf{a}[i,k] * \mathbf{b}[k,j]$

} see below

return  $\mathbf{c}$

→ with elementwise arithmetic:  
 $\mathbf{c}[i,j] = (\mathbf{a}[i,:] * \mathbf{b[:,j]}).\text{sum}()$

→ with broadcasting:  
 $\mathbf{c}[i] = (\mathbf{a}[i].\text{unsqueeze}(-1) * \mathbf{b}).\text{sum(dim>0)}$

**Elementwise arithmetic**: Basic operators ( $+, -, *, /, <, >, ==$ ) can be applied elementwise on tensors of the same size.

- If we want to know if all elements of tensor  $\mathbf{a}$  are  $<$  tensor  $\mathbf{b}$ , or if tensors are equal, we need to combine elementwise with  $\text{torch.all}(\mathbf{a} < \mathbf{b}), \text{all}(\mathbf{a} == \mathbf{b}), \text{all}(\mathbf{a})$

- reduction operations:  $\text{all}, \text{sum}, \text{mean}$  return rank 0 tensor (1 element) - convert to Python Boolean/number with  $\text{item}()$ :  $(\mathbf{a} > \mathbf{b}).\text{mean}().\text{item}()$

**Einstein summation**  
 $\text{def matmul}(\mathbf{a}, \mathbf{b}): \text{return torch.einsum('ik, kj \rightarrow ij', \mathbf{a}, \mathbf{b})$

- Used for combining products and sums  
 - operands dimensions on the left, separated by commas  
 - results dims on the right  
 - repeated indices are implicitly summed over.

**Chain rule**:  $(g \circ f)'(x) = g'(f(x)) f'(x)$

$y = g(u) \Rightarrow \frac{dy}{dx} = \frac{dy}{du} * \frac{du}{dx}$  SymPy:  
 $u = f(x) \Rightarrow \frac{du}{dx} = \frac{\partial u}{\partial x}$   $u_3$  for

From SymPy import symbols, diff  
 $sx, sy = \text{symbols}('sx sy')$  symbolic  
 $\text{diff}(sx * sx^2, sx) \# 2 * sx$  Computation

**The forward and backward passes**  
 $x = \text{torch.randn}(200, 100); y = \text{torch.randn}(200)$   
 We have to scale our weights so that activations have mean = 0 and std = 1. std < 0 or > 0 will lead to 0 or inf in our activations after many multiplications.  $\rightarrow$  Initialization

$w1 = \text{torch.randn}(100, 50) / \text{sqrt}(100) \rightarrow 1/\text{sqrt}(2/100)$   
 $w2 = \text{torch.randn}(50, 1) / \text{sqrt}(50) \rightarrow 1/\text{sqrt}(2/50)$   
 div by  $1/\sqrt{n_{in}}$ ,  $n_{in}$  = number inputs (float-tanh act.)  
 div by  $1/\sqrt{2/n_{in}}$   $\rightarrow$  learning rate, for relu activation

**Gradients and the backward pass (chain rule)**  
 To compute all gradients for the update, we begin from the model output (end work backward through the layers ( $\Rightarrow$  backpropagation)). We can automate it by having each function we implemented provide its backward step and populate gradients in an attribute of each tensor (like PyTorch  $.grad$ ), e.g.:

**Class Lin()**  
 $\text{def __init__}(\text{self}, \mathbf{w}, \mathbf{b}): \text{self.w, self.b} = \mathbf{w}, \mathbf{b}$   
 $\text{def __call__}(\text{self}, \mathbf{x}): \text{--- call --- makes our self. mp = mp class callable in Python}$   
 $\text{self.out} = \mathbf{x} @ \text{self.w} + \text{self.b}$   
 $\text{return self.out}$

**def backward(self):**  
 $\text{self.mp.g} = \text{self.out.g} @ \text{self.w.t()}$   
 $\text{self.w.g} = \text{self.mp.t()} @ \text{self.out.g}$   
 $\text{self.b.g} = \text{self.out.g.sum(0)}$

**Class Model():**  
 $\text{def __init__}(\text{self}, \mathbf{w1}, \mathbf{b1}, \mathbf{w2}, \mathbf{b2}):$   
 $\quad \text{self.layers} = [\text{Lin}(\mathbf{w1}, \mathbf{b1}), \text{ReLU}(), \text{Lin}(\mathbf{w2}, \mathbf{b2})]$   
 $\quad \text{self.loss} = \text{Mse}()$

**def \_\_call\_\_(self, x, target):**  
 $\quad \text{for l in self.layers: } x = l(x)$   
 $\quad \text{return self.loss}(x, target)$

**def backward(self):**  
 $\quad \text{self.loss.backward()}$   
 $\quad \text{for l in reversed(self.layers): } l.\text{backward}()$

**model = Model(w1, b1, w2, b2)** # instantiate  
 $\text{loss} = \text{model}(\mathbf{x}, \mathbf{y})$  # forward pass  
 $\text{model.backward()}$  # backward pass

Class Activation Map (CAM) Zhou et al.  
uses the output of the last convolutional layer (just before average pooling) together with the predictions to give us a heatmap visualization of why the model made its decision. - at each point/position of our final conv layer, we have as many filters as in the last linear layer

- we can compute the dot product of these activations with the final weights to get - for each location on our feature map - the score of the feature that was used to make a decision

Pytorch Hook: like python callbacks

use it to get access to activations inside the model while it's training (inject code into the forward and backward calculations)

Example: cats vs dogs model

```
learn = cnn_learner(dls, resnet34, metrics=error_rate)
learn.fine_tune(1)
class Hook():
    def __init__(self, m):
        self.hook = m.register_forward_hook(self.hook_func)
    def hook_func(self, m, i, o):
        self.stored = o.detach().clone()
hook_output = Hook()
hook = learn.model[0].register_forward_hook(hook_output.hook_func)
img = PIL.Image.open('images/cat.jpg')
x, _ = first(dls.test_dl([img]))
```

with torch.no\_grad(): grab a batch and feed it  
output = learn.model.eval()(x) through our model/  
act = hook\_output.stored[0] access stored activations/  
F.softmax(output, dim=-1) verify prediction!

x.shape → [1, 3, 224, 224] Einsum does dot  
product of weight  
cam\_map = torch.einsum('ck,kij→cij', learn.model[1][1].weight, act) matrix (2x  
learn.model[1][1].weight, act) number of activations  
cam\_map.shape → [2, 7, 7] with the activations  
7x7 feature map for each img (batch size x activations  
class - high/low activations x rows x columns)

Visualize CAM map: We need to decode input x as it was normalized by DataLoader, also cast to TensorImage

```
x_dec = TensorImage(
    dls.train.decode((x))[0][0])
ax = plt.subplots()
x_dec.show(ctx=ax)
```

ex. In show(cam\_map[1].detach(), gpu(),  
alpha=0.6, extent=(0, 224, 224, 0),  
interpolation='bilinear', cmap='magma');

hook.remove() remove when done to avoid leaking memory, OR:

Context manager: Python construct that calls `_enter_` when the object is created in a `with` clause, and `_exit_` at the end of the `with` clause

→ have the Hook class be a context manager, registering the hook when we enter it, and removing it when we exit

```
class Hook():
    def __init__(self, m):
        self.hook = m.register_forward_hook(self.hook_func)
    def hook_func(self, m, i, o):
        self.stored = o.detach().clone()
    def __enter__(self, *args):
        return self
    def __exit__(self, *args):
        self.hook.remove()
```

With Hook(learn.model[0]) as hook:  
with torch.no\_grad():  
 output = learn.model.eval()(x.cuda())  
 act = hook.stored



## Gradient-CAM (Selvaraju et al.)

CAM works only for the last layer. Grad-CAM uses the gradients of the final activation for the desired class (gradients of the output of the last layer with respect to input of that layer = layer weights, since it's a linear layer). In deeper layers, gradients are no longer equal to weights, calculated by PT in the backward pass, but not stored. We can register a hook on the backward pass to store them.

class HookBwd():
 def \_\_init\_\_(self, m):
 self.hook = m.register\_backward\_hook(self.hook\_func)

```
def hook_func(self, m, ip, op):
    self.stored = op[0].detach().clone()
```

def \_\_enter\_\_(self, \*args):
 return self

def \_\_exit\_\_(self, \*args):
 self.hook.remove()

cls = 1 # cat = True we can use it on any layer - here [-2]

with HookBwd(learn.model[0][-2]) as hook:  
with Hook(learn.model[0][-2]) as hook:

output = learn.model.eval()(x.cuda())
 act = hook.stored
 output[0, cls].backward()
 grad = hook.stored

w = grad[0].mean(dim=[1, 2], keepdim=True)
cam\_map = (w \* act[0]).sum(0)

\* We can't use output.backward() → gradients make sense only with respect to a scalar (normally our loss) and output here is a rank-2 tensor. But if we pick a single image and a single class, we can calculate the gradients of weights or activations w.r.t. that single value. We'll use these gradients as weights. Weights for Grad-CAM = average of our gradients across the feature map.

Fastbook ch 19(1) A fastai Learner notes by @dk21 from scratch, p.1

## DATA

untar\_data

get\_image\_files: glob vs os.walk

Path.parent (pathlib)

L.val[2].idx, L.map

Python: ~ → filter = [True, False]  
a = alist[filter], b = alist[~filter]

## DATASET

class Dataset:

```
def __init__(self, fns): self.fns = fns
def __len__(self): return len(self.fns)
def __getitem__(self, i):
    im = Image.open(self.fns[i]).resize((64, 64))
    .convert('RGB')
    y = v2i[self.fns[i].parent.name]
    return tensor(im).float()/255, tensor(y)
```

→ fns returns x, y as a tuple that we need to collate into a minibatch

def collate(idxs, ds):

```
xb, yb = zip(*[ds[i] for i in idxs])
return torch.stack(xb), torch.stack(yb)
```

class DataLoader:

```
def __init__(self, ds, bs=128, shuffle=False, n_workers=1):
    self.ds, self.bs, self.shuffle, self.n_workers = ds, bs, shuffle, n_workers
```

```
def __len__(self): return (len(self.ds)-1)//self.bs+1
```

def \_\_iter\_\_(self):

idxs = L.range(self.ds)

if self.shuffle: idxs = idxs.shuffle()

```
    chunks = [idxs[n:n+self.bs] for n in range(
        0, len(self.ds), self.bs)]
```

with ProcessPoolExecutor(self.n\_workers) as ex:  
 yield from ex.map(collate, chunks, ds=self.ds)

→ We need to do our processing in parallel - opening and decoding a JPEG image is a slow process, one CPU core is not enough to decode images fast enough to keep a GPU busy.

n\_workers = min(16, defaults.cpus)

We'll need image stats for normalization:

stats = [xb.mean((0, 1, 2)), xb.std((0, 1, 2))]

class Normalize:

```
def __init__(self, stats): self.stats = stats
def __call__(self, x):
    if x.device != self.stats[0].device:
        self.stats = to_device(self.stats, x.device)
    return (x - self.stats[0]) / self.stats[1]
```

norm = Normalize(stats)

```
def tfm_x(x): return norm(x).permute((0, 3, 1, 2))
```

→ permute the axis order from PIL (HWC) to Pytorch (NCHW)

## MODULE AND PARAMETER

class Module:

```
def __init__(self):
    self.hooks, self.params, self.children, self.training =
        None, [], [], False
    def register_parameters(self, *ps): self.params += ps
    def register_modules(self, *ms): self.children += ms
```

@property

def training(self): return self.training

@training.setter

```
def training(self, v):
    self.training = v
    for m in self.children: m.training = v
```

def parameters(self):

```
    return self.params + sum([m.parameters()
        for m in self.children], []) # recursion
```

def setattr(self, k, v):

```
    super().__setattr__(k, v)
```

if isinstance(v, Parameter): self.register\_parameters(v)

if isinstance(v, Module): self.register\_modules(v)

def \_\_call\_\_(self, \*args, \*\*kwargs):

res = self.forward(\*args, \*\*kwargs)

if self.hook is not None: self.hook(res, args)

return res

def cuda(self):

```
    for p in self.parameters: p.data = p.data.cuda()
```

class Parameter(Tensor): \*

```
def new_(self, x): return Tensor._make_subclass(
    Parameter, x, True)
```

```
def __init__(self, *args, **kwargs): self.requires_grad_()
```

\* Parameter doesn't add any functionality, only calls requires\_grad for us. It's used as a marker to show what to include in parameters!

SIMPLE CNN =>

class ConvLayer(Module):

```
def __init__(self, ni, nf, stride=1, bias=True, act=True):
    super().__init__()
    self.w = Parameter(torch.zeros(nf, ni, 3, 3))
    self.b = Parameter(torch.zeros(nf)) if bias else None
    self.act, self.stride = act, stride
    init = nn.init.kaiming_normal_ if act else nn.init.xavier_normal_
    init(self.w)
```

def forward(self, x):

```
x = F.conv2d(x, self.w, self.b, stride=self.stride, padding=1)
if self.act: x = F.relu(x)
return x
```

class Linear(Module):

```
def __init__(self, ni, nf):
```

```
super().__init__()
self.w = Parameter(torch.zeros(nf, ni))
self.b = Parameter(torch.zeros(nf))
```

```
nn.init.xavier_normal_(self.w)
def forward(self, x): return x @ self.w.t() + self.b
```

class Sequential(Module):

```
def __init__(self, *layers):
```

```
super().__init__()
self.layers = layers
```

```
def register_modules(*layers):
```

```
def forward(self, x):
```

```
for l in self.layers: x = l(x)
return x
```

class AdaptivePool(Module):

```
def forward(self, x): return x.mean((2, 3))
```

def simple\_cnn():

```
return Sequential(
    ConvLayer(3, 16, stride=2), #32
    ConvLayer(16, 32, stride=2), #16
    ConvLayer(32, 64, stride=2), #8
    ConvLayer(64, 128, stride=2), #4
    AdaptivePool(),
    Linear(128, 10))
```

Fastbook ch.19 (2) A fastai Learner notes by @dhr21 from Scratch, p.2

## LOSS

```
def nll(input, target): return -input[range(target.shape[0]), target].mean()

1) def logsoftmax(x): return (x.exp() / (x.exp().sum(-1, keepdim=True))).log()

2) def log_softmax(x): return x - x.exp().sum(-1, keepdim=True)).log() # log(1/x) = log(1) - log(x)

3) def logsumexp(x):  
    m = x.max(-1)[0] # LogSumExp trick --  
    return m + (x - m[:, None]).exp().sum(-1).log()

def logsoftmax(x): return x - x.logsumexp(-1, keepdim=True) .. for numeric stability

def cross_entropy(preds, yb): return nll(logsoftmax(preds), yb).mean()
```

## CALLBACKS

class Callback (GetAttr): default='learner'

class SetupLearnerCB (Callback):

```
def before_batch(self):  
    xb, yb = to_device(self.batch)  
    self.learner.batch = tfm_x(xb), yb  
def before_fit(self): self.model.cuda()
```

class TrackResults (Callback):

```
def before_epoch(self):  
    self.dccs, self.losses, self.ns = [], [], []
```

```
def after_epoch(self):  
    n = sum(self.ns)  
    print(f'Epoch {self.epoch} {self.model.training}'  
        f'Sum({self.losses})/{n}, sum({self.accs})/{n}')
```

```
def after_batch(self):  
    xb, yb = self.batch  
    acc = (self.preds.argmax(dim=1) == yb).float().sum()
```

self.accs.append(acc)

n = len(xb)

self.losses.append(self.loss \* n)

self.ns.append(n)

cls = [SetupLearnerCB(), TrackResults()]

```
learn = Learner(simple_cnn(), dls, cross_entropy,  
                 lr=0.1, cbfs=cls)
```

## LEARNER

```
class SGD:  
    def __init__(self, params, lr, wd=0): store_attr(self, 'params, lr, wd')  
    def step(self):  
        for p in self.params:  
            p.data -= (p.grad.data + p.data * self.wd) * self.lr  
            p.grad.data.zero_()
```

class DataLoaders:  
 def \_\_init\_\_(self, \*dls): self.train, self.valid = dls

class Learner:  
 def \_\_init\_\_(self, model, dls, loss\_func, lr, cbs, opt\_func=SGD):  
 store\_attr(self, 'model, dls, loss\_func, lr, cbs, opt\_func')

```
for cb in cbs: cb.learner = self  
def one_batch(self): → every callback knows  
    → self('before_batch') what learner it is used in  
    xb, yb = self.batch → get info from learner  
    self.preds = self.model(xb) → change things in  
    self.loss = self.loss_func(self.preds, yb) learner
```

```
if self.model.training: when Learner calls self,  
    self.loss.backward() it calls __call__, which  
    self.opt.step() uses getattr(cb, name)
```

```
→ self('after-batch') on each callback in  
    def one_epoch(self, train): self.cbs where the  
    self.model.training = train attribute (here: method)  
    → self('before-epoch') is defined.  
    dl = self.dls.train if train else self.dls.valid  
    for self.num, self.batch in enumerate(progress_bar(  
        dl, leave=False)):
```

```
        self.one_batch() eg. self('before-fit') will call  
    → self('after-epoch') cb.before_fit() if defined  
    def fit(self, n_epochs):  
        → self('before-fit')  
        self.opt = self.opt_func(self.model.parameters(), self.lr)  
        self.n_epochs = n_epochs
```

try:  
 for self.epoch in range(n\_epochs):

```
        self.one_epoch(True)  
        self.one_epoch(False)  
    except CancelFitException: pass
```

```
    → self('after-fit')  
    def __call__(self, name):  
        for cb in self.cbs: setattr(cb, name, noop)()
```

## SCHEDULING THE LEARNING RATE

class LR Finder (Callback):

```
def before_fit(self):  
    self.losses, self.lrs = [], []
```

self.learner.lr = 1e-6

def before\_batch(self):

if not self.model.training: return  
self.opt.lr \*= 1.2

def after\_batch(self):

if not self.model.training: return  
if self.opt.lr > 10 or torch.isnan(self.loss):

raise CancelFitException

self.losses.append(self.loss.item())

self.lrs.append(self.opt.lr)

class OneCycle (Callback):

```
def __init__(self, base_lr): self.base_lr = base_lr
```

def before\_fit(self): self.lrs = []

def before\_batch(self):

if not self.model.training: return  
n = len(self.dls.train)

bn = self.epoch + n + self.num

mn = self.n\_epochs \* n

pct = bn / mn

pct\_start, div\_start = 0.25, 10

if pct < pct\_start:

pct = pct\_start

lr = (1 - pct) \* self.base\_lr / div\_start + pct \* self.base\_lr

else:  
 pct = (pct - pct\_start) / (1 - pct\_start)

lr = (1 - pct) \* self.base\_lr

self.opt.lr = lr

self.lrs.append(lr)

\* GetAttr is a helper class that implements  
Python's \_\_getattribute\_\_ and \_\_dir\_\_: if we try to access  
attribute that doesn't exist, it passes the request  
to \_\_default\_\_, e.g:

self.model.cuda → self.learner.model.cuda

(only works for getting attributes, not setting)

active on  
Fastai Forum

2020 : Deep Dive  
into Machine Learning,  
Deep Learning, Data  
Science

active on  
Kaggle Forum

Huge thanks  
to Jeremy Howard,  
Sylvain Gugger and  
Rachel Thomas  
& fastai community  
for the knowledge  
and inspiration!

I started a blog:  
skok.ai

I published a paper:  
'Polbert: Attaching Polish  
NLP Tasks with Transformers'

NEW:  
to be  
announced ↗

Fastbook  
notes

NEW: Kaggle competition  
notes

HELP

WRITE

Concluding  
Thoughts

GATHER

BUILD

Co-  
Organizing  
Polish NLP  
Meetup Group

NEW: fastai community  
Meetup group

ML Community of Practice  
@ work

2021: carry on,  
build meaningful & ethical  
data products and  
contribute to  
the community.

Poleval: winning Polish  
NLP competition  
task: Word Sense  
Disambiguation

Contribute to  
Hugging Face library  
(datasets)

Work projects: Image Similarity,  
Object detection,  
Classification,  
Random Forest - tabular