# Training Neural Networks

# Lecture 5 Recap

# Gradient Descent for Neural Networks

Loss function
$$L_i = (\hat{y}_i - y_i)^2$$



input layer

hidden layer

output layer

$$\nabla_{\boldsymbol{W},\boldsymbol{b}} f_{\{\boldsymbol{x},\boldsymbol{y}\}}(\boldsymbol{W}) = \begin{bmatrix} \dfrac{\partial f}{\partial w_{0,0,0}} \\ \cdots \\ \cdots \\ \dfrac{\partial f}{\partial w_{l,m,n}} \\ \cdots \\ \cdots \\ \dfrac{\partial f}{\partial b_{l,m}} \end{bmatrix}$$

optimizing for
weights & biases

$$\hat{y}_i = A\left(b_{1,i} + \sum_j h_j w_{1,i,j}\right)$$

$$h_j = A\left(b_{0,j} + \sum_k x_k w_{0,j,k}\right)$$

Just simple:
$$A(x) = \max(0, x)$$

# Stochastic Gradient Descent (SGD)

$$\boldsymbol{\theta}^{k+1} = \boldsymbol{\theta}^k - \alpha \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}^k, \boldsymbol{x}_{\{1..m\}}, \boldsymbol{y}_{\{1..m\}})$$

$k$ now refers to $k$-th iteration

$$\nabla_{\boldsymbol{\theta}} L = \frac{1}{m} \sum_{i=1}^{m} \nabla_{\boldsymbol{\theta}} L_i$$

$m$ training samples in the current minibatch

Gradient for the $k$-th minibatch

# Gradient Descent with Momentum

$$v^{k+1} = \beta \cdot v^k + \nabla_{\theta} L(\theta^k)$$

accumulation rate ('friction', momentum)

velocity

Gradient of current minibatch

In case grads are going in the same dir, we are accumulating them i.e going faster
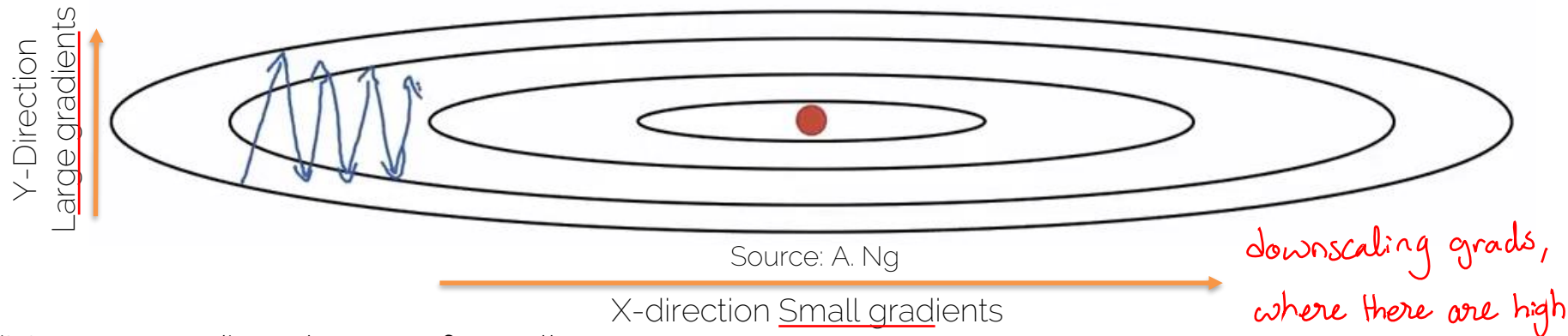
$$\theta^{k+1} = \theta^k - \alpha \cdot v^{k+1}$$

model

learning rate

velocity

Exponentially-weighted average of gradient

Important: velocity $v^k$ is vector-valued!

# RMSProp



Y-Direction Large gradients

Source: A. Ng

X-direction Small gradients

downscaling grads, where there are high variance

(Uncentered) variance of gradients → second momentum

$$\boldsymbol{s}^{k+1} = \beta \cdot \boldsymbol{s}^k + (1-\beta)[\nabla_{\boldsymbol{\theta}} L \circ \nabla_{\boldsymbol{\theta}} L]$$

We're dividing by square gradients:
- Division in Y-Direction will be large
- Division in X-Direction will be small

$$\boldsymbol{\theta}^{k+1} = \boldsymbol{\theta}^k - \alpha \cdot \frac{\nabla_{\boldsymbol{\theta}} L}{\sqrt{\boldsymbol{s}^{k+1}} + \epsilon}$$

Can increase learning rate!

# Adam

- Combines Momentum and RMSProp

$$m^{k+1} = \beta_1 \cdot m^k + (1 - \beta_1)\nabla_{\boldsymbol{\theta}}L(\boldsymbol{\theta}^k) \qquad v^{k+1} = \beta_2 \cdot v^k + (1 - \beta_2)[\nabla_{\boldsymbol{\theta}}L(\boldsymbol{\theta}^k) \circ \nabla_{\boldsymbol{\theta}}L(\boldsymbol{\theta}^k)$$

"mean"     "variance"

- $m^{k+1}$ and $v^{k+1}$ are initialized with zero
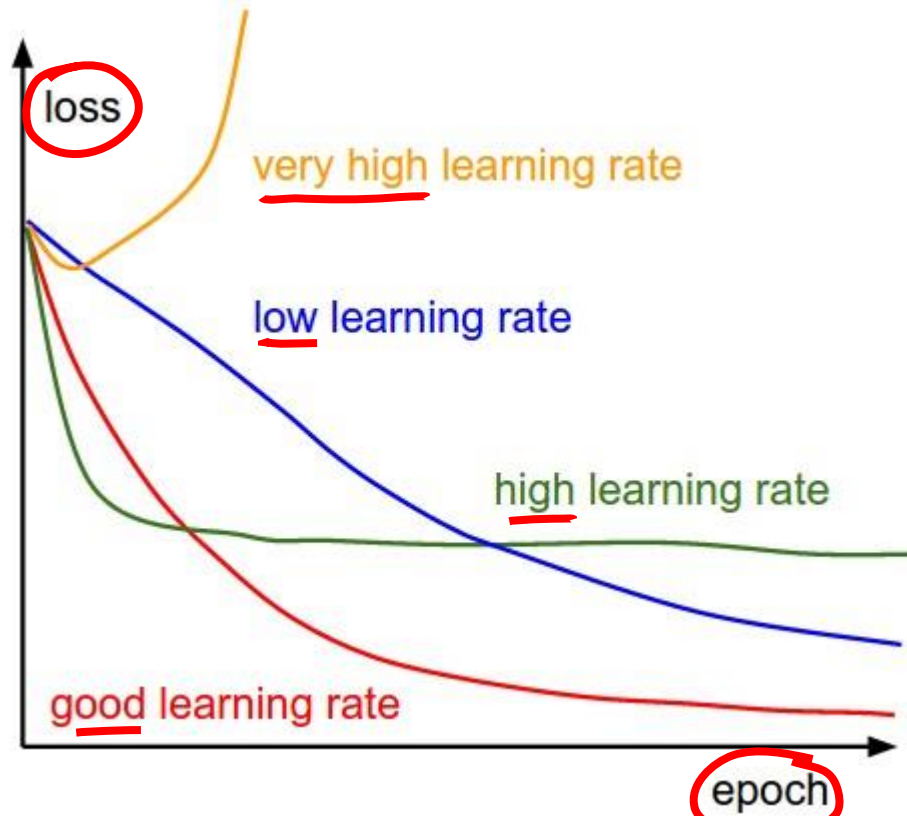  → bias towards zero
  → Typically, bias-corrected moment updates

$$\widehat{m}^{k+1} = \frac{m^{k+1}}{1 - {\beta_1}^{k+1}} \qquad \widehat{v}^{k+1} = \frac{v^{k+1}}{1 - {\beta_2}^{k+1}} \qquad \longrightarrow \qquad \boldsymbol{\theta}^{k+1} = \boldsymbol{\theta}^k - \alpha \cdot \frac{\widehat{m}^{k+1}}{\sqrt{\widehat{v}^{k+1}} + \epsilon}$$

# Training Neural Nets

# Learning Rate: Implications

- What if too high?

- What if too low?



loss

very high learning rate

low learning rate

high learning rate

good learning rate

epoch

overshooting optimum; fluctuating alot
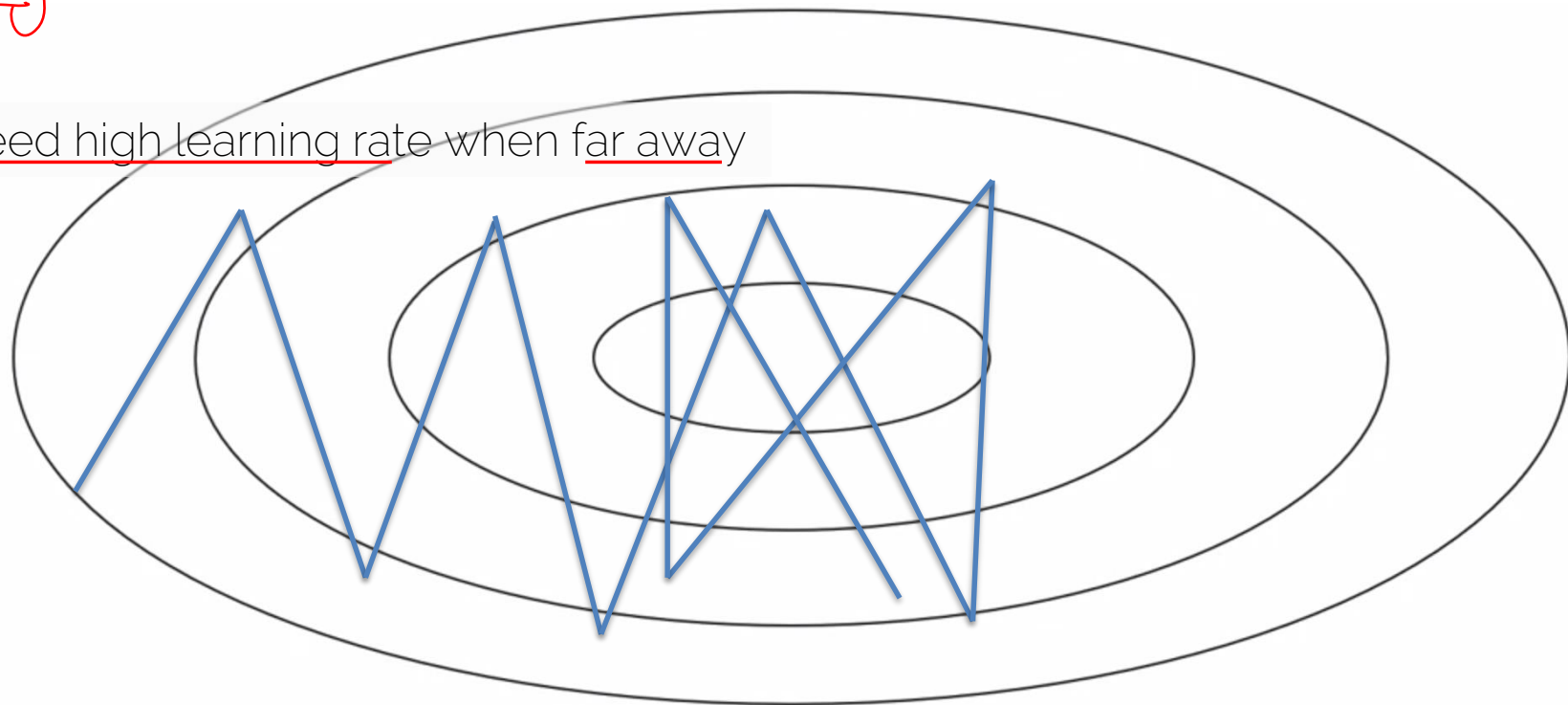
Source: http://cs231n.github.io/neural-networks-3/

# Learning Rate

*Ideally*



Need high learning rate when far away

Need low learning rate when close

# Learning Rate Decay

- $\alpha = \dfrac{1}{1+decay\_rate * epoch} \cdot \alpha_0$

$\alpha$ gets smaller, as epoch increase

  – E.g., $\alpha_0 = 0.1$, $decay\_rate = 1.0$

→ Epoch 0:   0.1

→ Epoch 1:   0.05

→ Epoch 2:   0.033

→ Epoch 3:   0.025

...

**Learning Rate over Epochs**

*Same idea : start falling quickly " high α ", & then relatively keep it constant*

# Many options:

- **Step decay** $\alpha = \alpha - t \cdot \alpha$ (only every n steps)
  - T is decay rate (often 0.5)

- **Exponential decay** $\alpha = t^{epoch} \cdot \alpha_0$
  - t is decay rate (t < 1.0)

- $\alpha = \dfrac{t}{\sqrt{epoch}} \cdot a_0$
  - t is decay rate

- Etc.

# Training Schedule

Manually specify learning rate for entire training process

- Manually set learning rate every n-epochs
- How?
  - Trial and error (the hard way)
  - Some experience (only generalizes to some degree)

Consider: #epochs, training set size, network size, etc.

# Basic Recipe for Training

- Given ground dataset with ground labels
  - $\{x_i, y_i\}$
    - $x_i$ is the $i^{th}$ training image, with label $y_i$
    - Often $\dim(x) \gg \dim(y)$ (e.g., for classification)
    - $i$ is often in the 100-thousands or millions
  - Take network $f$ and its parameters $w, b$
  - Use SGD (or variation) to find optimal parameters $w, b$
    - Gradients from backpropagation

# Gradient Descent on Train Set

- Given large train set with ($n$) training samples $\{x_i, y_i\}$
  - Let's say 1 million labeled images
  - Let's say our network has 500k parameters


- Gradient has 500k dimensions
- $n = 1\ million$
- ✱ Extremely expensive to compute

✱ Even with minibatches, it is still expensive!
↳ majority of computation is gradient.

# Learning

**✱** <u>Learning means generalization</u> to unknown dataset

- (So far no 'real' learning)
- I.e., train on known dataset → test with optimized parameters on unknown dataset

- Basically, we hope that based on the train set, the optimized parameters will give similar results on different data (i.e., test data)

**✱ GOAL** → getting good descriptors, that are agnostic to specific dataset

↘ making training harder i.e avoiding data memoriZation

- <mark>Training set</mark> ('*train*'):
  - Use for training your neural network
- <mark>Validation set</mark> ('*val*'):
  - Hyperparameter optimization
  - Check generalization progress
- <mark>Test set</mark> ('*test*'):
  - Only for the very end
  - ➤ NEVER TOUCH DURING DEVELOPMENT OR TRAINING

- Typical splits
  - Train (60%), Val (20%), Test (20%)
  - Train (80%), Val (10%), Test (10%)

*Important*

*• During training:*
  - Train error comes from average minibatch error
  - → Typically take subset of validation every n iterations

* Doing so is not expensive! Computing grad through backprop is by an order of magnitude more expensive.

# Basic Recipe for Machine Learning

- Split your data

- Split your data

60%          20%     20%

| train | validation | test |
| --- | --- | --- |

**Example scenario**

Ground truth error ...... 1%

Training set error ....... 5%

Val/test set error ...... 8%

*Bias* (underfitting)

*Variance* (overfitting)

Generalization Gap

**Very Important !**



Training error high? → Yes → Bigger model / Train longer / New model architecture — (Bias)

↓ No

Train-Dev error high? → Yes → *Overfitting* — More data / Regularization / New model architecture — (Variance)

↓ No

Dev error high? → Yes → Make training data more similar to test data. / Data synthesis (Domain adaptation.) / New model architecture — (Train-test data mismatch)

↓ No

Test error high? → Yes → More dev set data — (Overfit dev set)

↓ No

Done

Credits: A. Ng

# Over- and Underfitting



Underfitted

Appropriate

Overfitted

Source: Deep Learning by Adam Gibson, Josh Patterson, O'Reily Media Inc., 2017

Underfitting zone | Overfitting zone

generalization error

generalization gap

training error

Source: https://srdas.github.io/DLBook/ImprovingModelGeneralization.html

# Learning Curves

* **Training graphs**
  - <u>Accuracy</u>
  - <u>Loss</u>



* Small batches are the reason of these noisy fluctuations "high variance"

* Can not do large batches, for memory reasons

* do some smoothing fn e.g exp avg

**NB** One good practice is to decrease $\alpha$, once loss starts to decrease slowly.

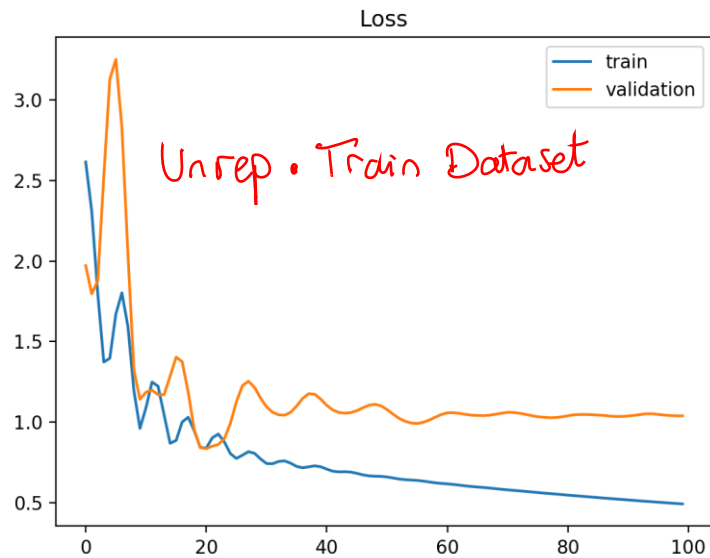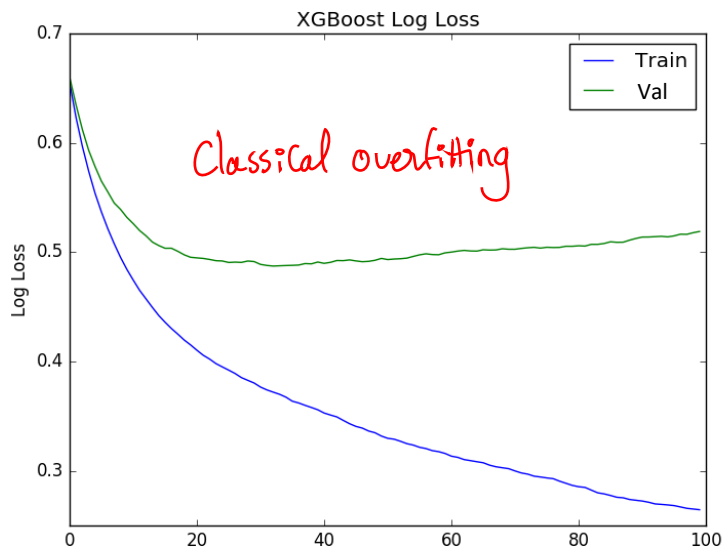Sometimes, however, iterating over more epochs would be ideal.

Typically, α shall be high at the beginning, then decrease over time.



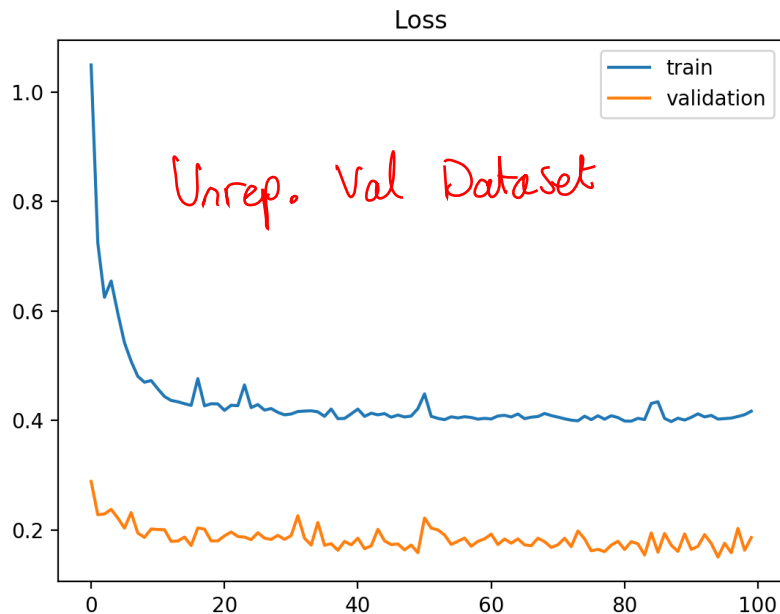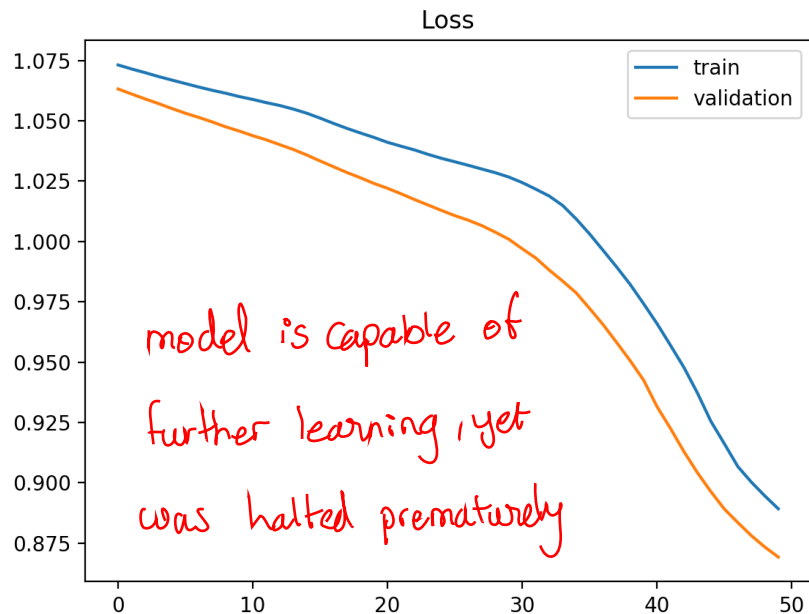model loss

high α at the beginning,
should've decreased over time;
to avoid such plateau near the end

# Overfitting Curves



XGBoost Log Loss

Classical overfitting

Loss

Unrep. Train Dataset

# Other Curves

### Loss (left chart)

model is capable of
further learning, yet
was halted prematurely

### Loss (right chart)

Unrep. Val Dataset

Underfitting (loss still decreasing)

Validation Set is easier than Training set

↳ flat loss regardless of training
↳ training continues to decrease until the end

# To Summarize

*Important*

- Underfitting
  - Training and validation losses decrease even at the end of training
- Overfitting
  - Training loss decreases and validation loss increases
- Ideal Training
  - Small gap between training and validation loss, and both go down at same rate (stable without fluctuations).

# Bad Signs

- Training error not going down
- Validation error not going down
- Performance on validation better than on training set
- Tests on train set different than during training

# Bad Practice

- Training set contains test data
- Debug algorithm on test data

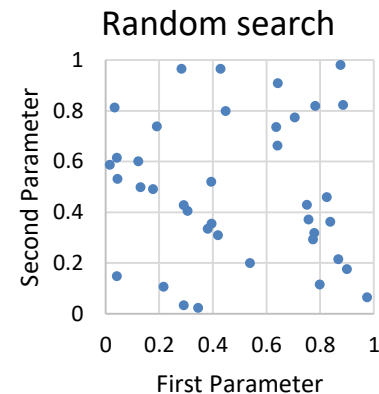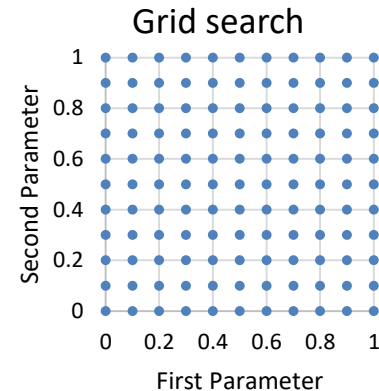Never touch during development or training

# Hyperparameters

*anything but **model** parameters*

- Network architecture (e.g., num layers, #weights)

- Number of iterations

- Learning rate(s)  (i.e., solver parameters, decay, etc.)

- Regularization (more later next lecture)

- Batch size

- ...

* Overall:
learning setup + optimization = hyperparameters

# Hyperparameter Tuning

- Methods:
  - **Manual** search:
    - most common ☺
  - **Grid** search (structured, for 'real' applications)
    - Define ranges for all parameters spaces and select points
    - Usually pseudo-uniformly distributed
    - → Iterate over all possible configurations
  - **Random** search:
    Like grid search but one picks points at random in the predefined ranges

**Grid search**



First Parameter

**Random search**



First Parameter

# How to Start

ALWAYS

- Start with single training sample
  - Check if output correct *i.e debug loss fn, not the input*
  - Overfit → train accuracy should be 100% because input just memorized
- Increase to handful of samples (e.g., **4**)
  - Check if input is handled correctly

  *no need to extract features, just learn the input*
- Move from overfitting to more samples
  - 5, 10, 100, 1000, …
  - At some point, you should see generalization

ONE

FEW

MANY
…

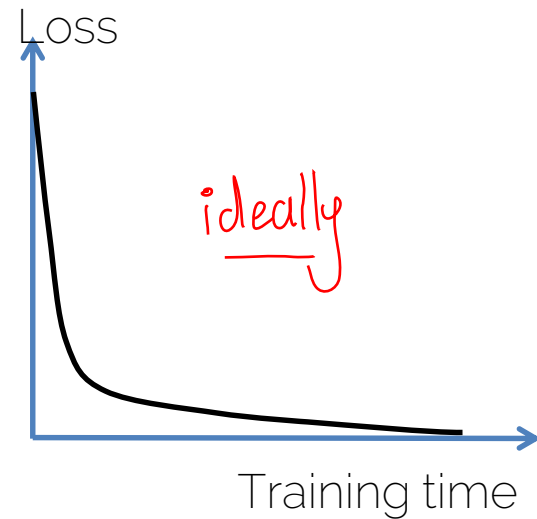# Find a Good Learning Rate
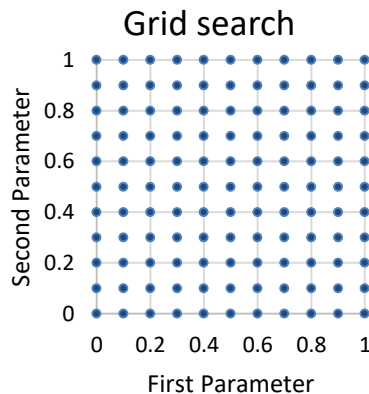
- Use all training data with small weight decay
- Perform initial loss sanity check e.g., $\log(C)$ for softmax with $C$ classes
- Find a learning rate that makes the loss drop significantly (exponentially) within 100 iterations
- Good learning rates to try: 1e-1, 1e-2, 1e-3, 1e-4

Loss

ideally

Training time

# Coarse Grid Search

usually out of compute budget

- Choose a few values of learning rate and weight decay around what worked from

- Train a few models for a few epochs.

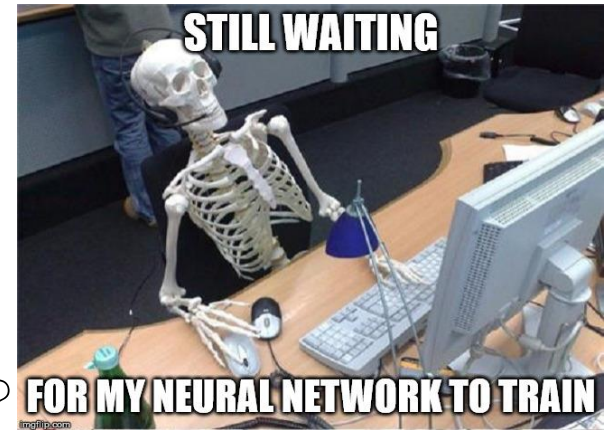- Good weight decay to try: 1e-4, 1e-5, 0

### Grid search

# Refine Grid

- Pick best models found with coarse grid.

- Refine grid search around these models.

- Train them for longer (10-20 epochs) without learning rate decay

- Study loss curves <- most important debugging tool!

# Timings

- How long does each <u>iteration</u> take?
  - Get **precise** timings!
  - If an iteration exceeds **500**ms, things get dicey
- <u>Look for bottlenecks</u>
  - Dataloading: smaller resolution, compression, train from SSD
  - Backprop
- <u>Estimate total time</u>
  - How long until you see some pattern?
  - How long till convergence?


STILL WAITING FOR MY NEURAL NETWORK TO TRAIN

# Network Architecture

- Frequent mistake: "*Let's use this super big network, train for two weeks and we see where we stand.*"

* Instead: start with simplest network possible
  - Rule of thumb divide **#**layers you started with by 5
- Get debug cycles down
  - Ideally, minutes

# Debugging

- Use train/validation/test curves
  - Evaluation needs to be consistent
  - Numbers need to be comparable

* Only make **one change at a time**
  - "I've added 5 more layers and double the training size, and now I also trained 5 days longer. Now it's better, but why?"

# Common Mistakes in Practice

- <u>Did not overfit to single batch first</u>

- Forgot to <u>toggle train/eval mode</u> for network
  - Check later when we talk about dropout…

- Forgot to <u>call</u> **.zero_grad()** (*in PyTorch*) before calling **.backward()**

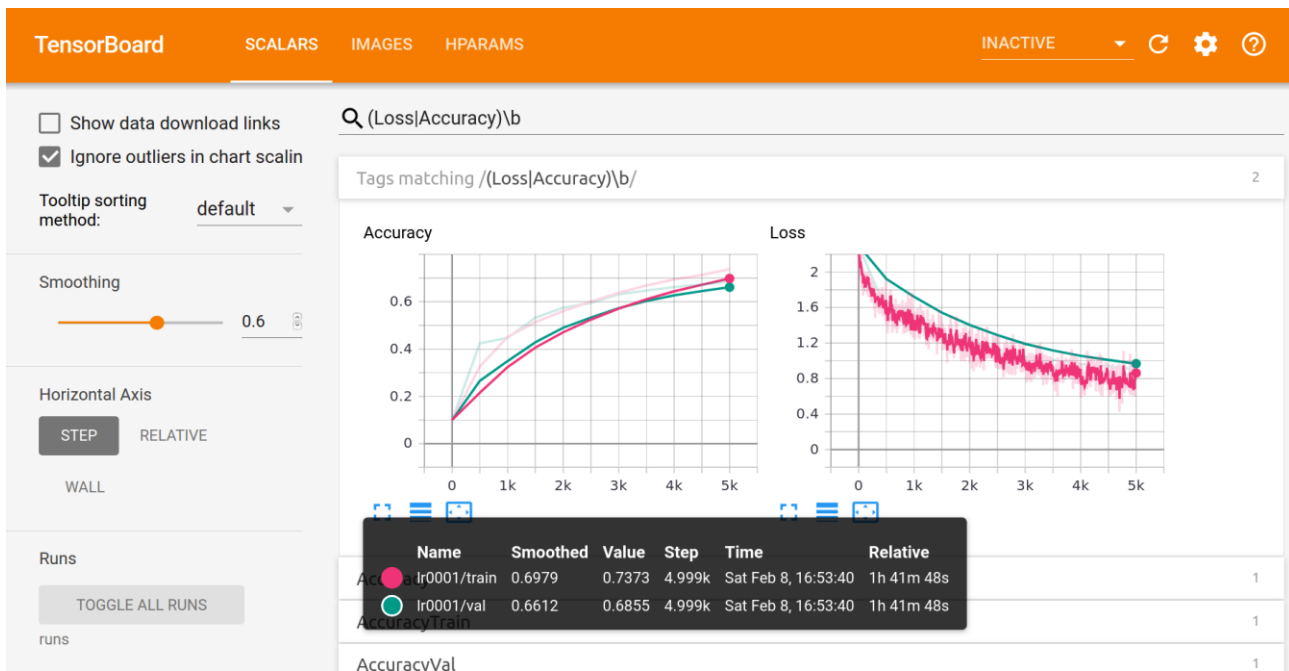- <u>Passed softmaxed outputs</u> to a <u>loss func</u>tion that <u>expects raw logits</u>
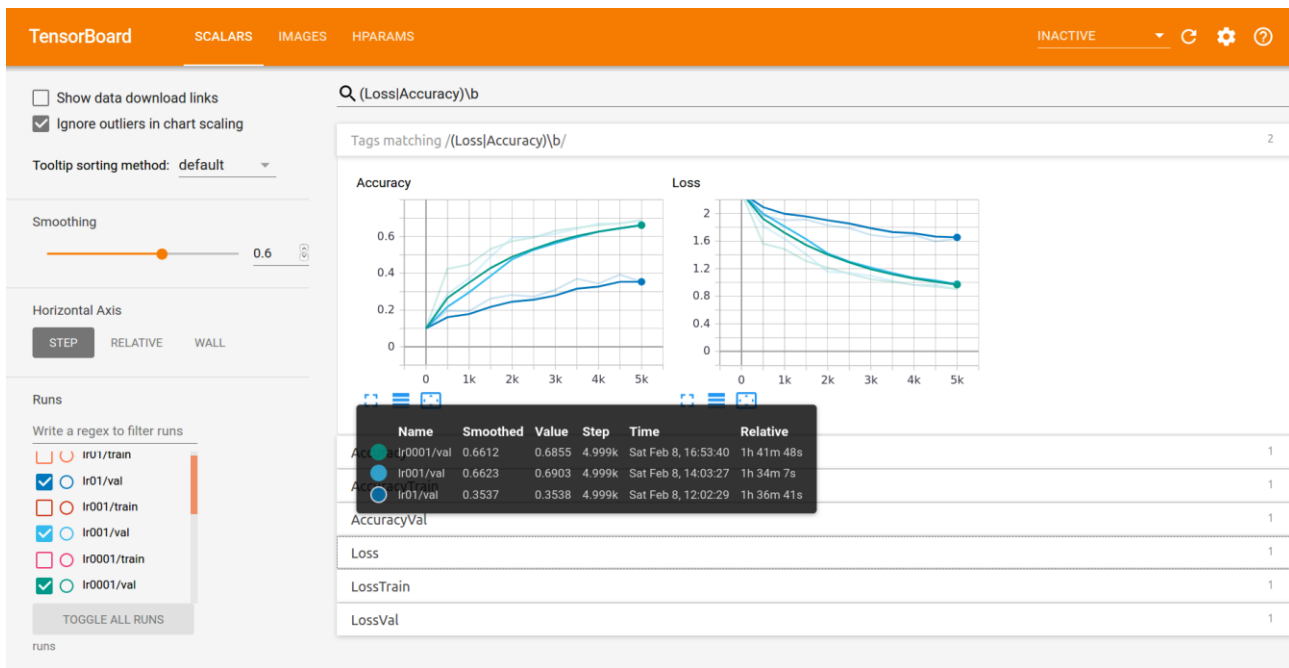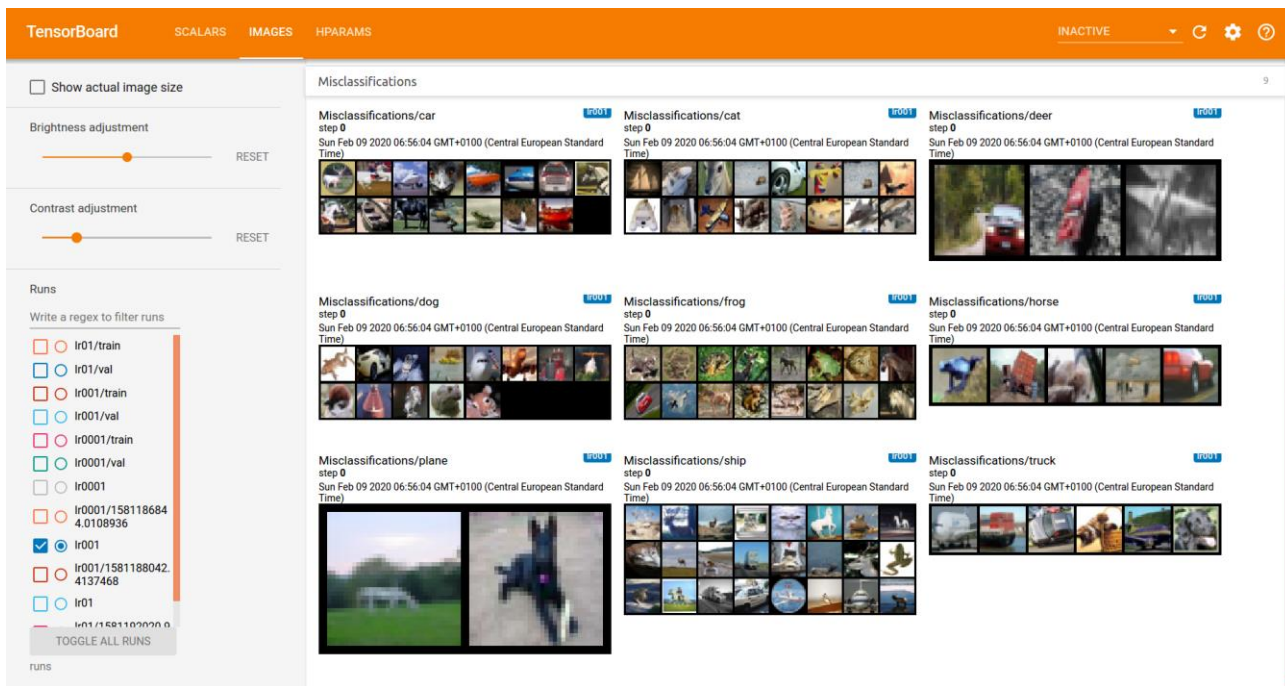
# Tensorboard: Visualization in Practice

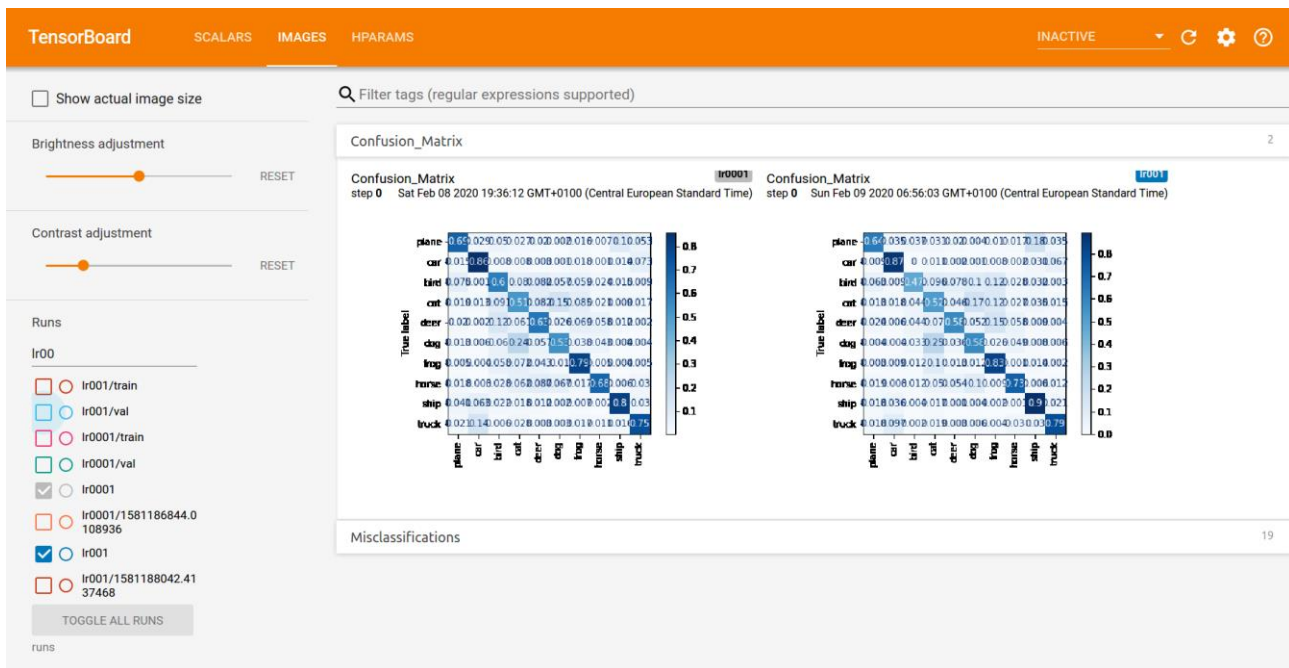# Tensorboard: Compare Train/Val Curves
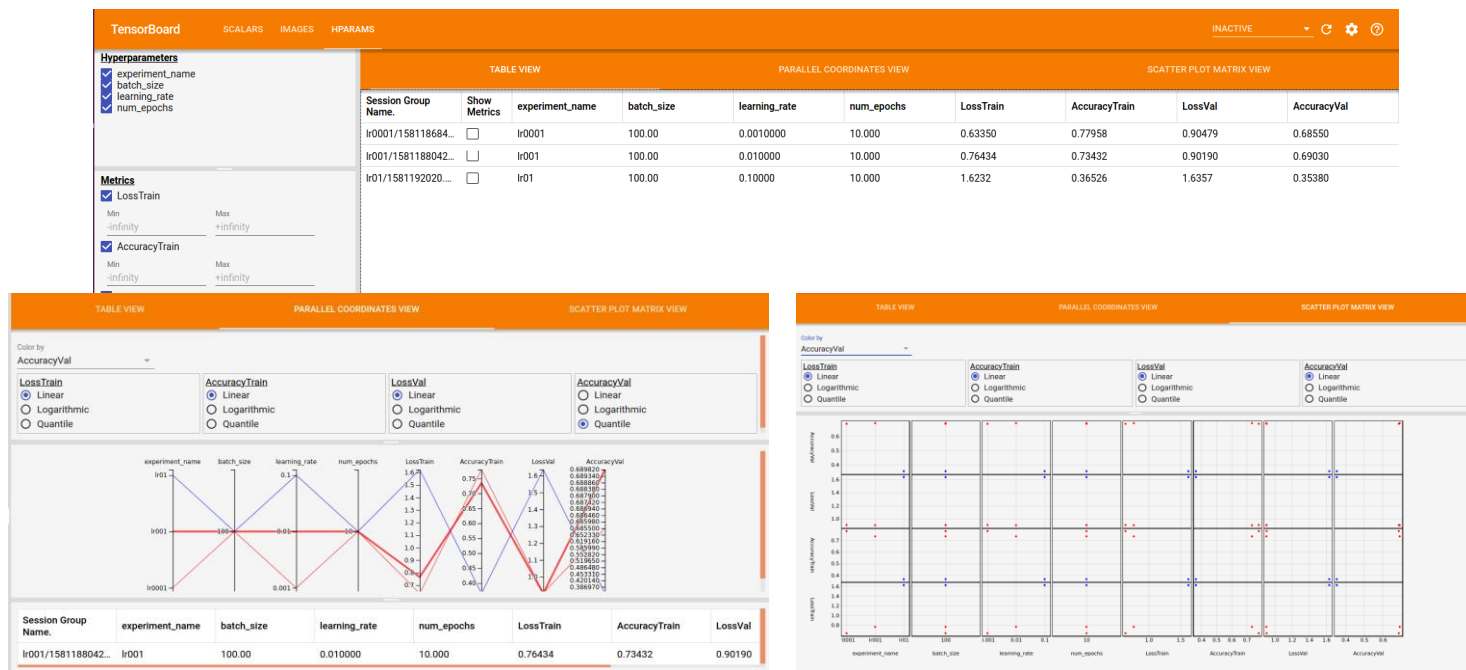
# Tensorboard: Compare Different Runs

# Tensorboard: Visualize Model Predictions

# Tensorboard: Visualize Model Predictions

# Tensorboard: Compare Hyperparameters

# Next Lecture

- Next lecture
  - More about training neural networks: output functions, loss functions, activation functions

- Check the exercises ☺

See you next week ☺

# References

- Goodfellow et al. "Deep Learning" (2016),
  - Chapter 6: Deep Feedforward Networks

- Bishop "Pattern Recognition and Machine Learning" (2006),
  - Chapter 5.5: Regularization in Network Nets

- http://cs231n.github.io/neural-networks-1/

- http://cs231n.github.io/neural-networks-2/

- http://cs231n.github.io/neural-networks-3/