

# ROB 498/599: Deep Learning for Robot Perception (DeepRob)

---

Lecture 10: Training Neural Networks - Part 2

02/11/2026



# Today

P2 Due Feb.15, 2026

Canvas Quiz: Due Feb 18, 2026

- Feedback and Recap (5min)
- Training NNs
  - One-time Setup (Cont'd)
    - Weight Initialization
    - Dropout
  - Training Dynamics
    - Learning Rate scheduling
    - Choosing Hyperparameters
  - If time permits: After Training
    - Model Ensembles, Transfer Learning
- Summary and Takeaways (5min)

# Aha Slides (In-class participation)

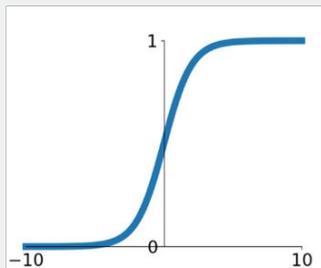
<https://ahaslides.com/RLCCA>



# Recap: Activation Functions

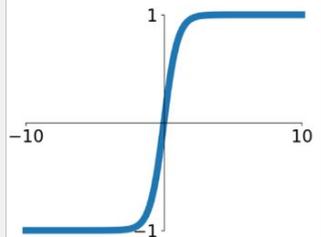
## Sigmoid

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



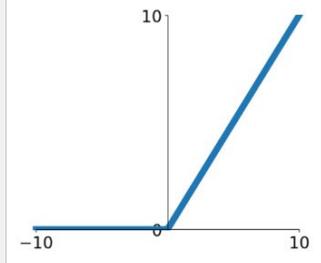
## tanh

$$\tanh(x)$$



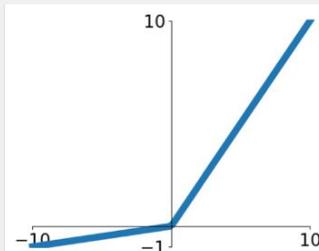
## ReLU

$$\max(0, x)$$



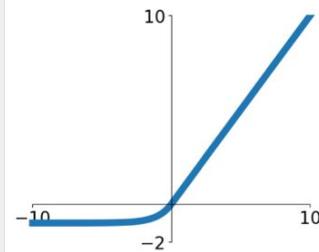
## Leaky ReLU

$$\max(0.1x, x)$$



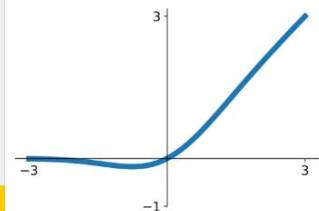
## ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(\exp^x - 1) & x < 0 \end{cases}$$

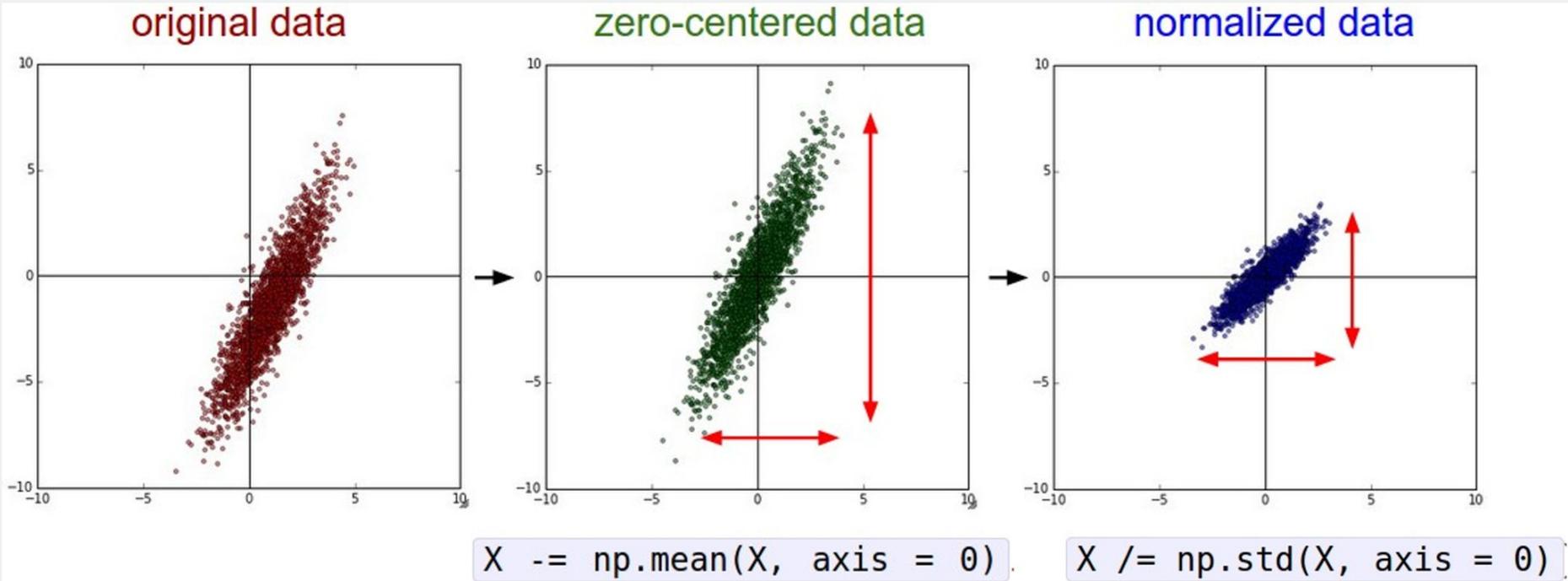


## GELU

$$\approx x\alpha(1.702x)$$



# Recap: Data Preprocessing

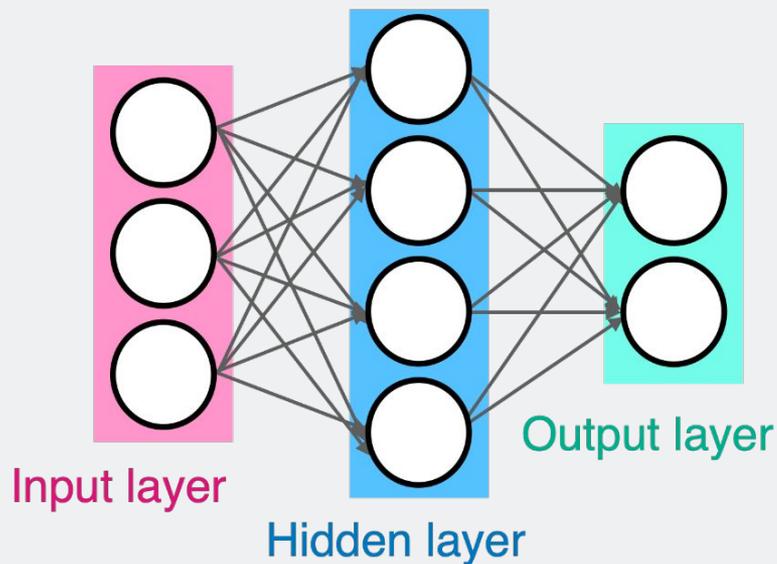


(Assume  $X[N \times D]$  is data matrix, each example in a row)

# Weight Initialization

# Weight Initialization

---



**Q:** What happens if we initialize all  $W=0$ ,  $b=0$ ?

**A:** All outputs are 0, all gradients are the same!  
No “symmetry breaking”

# Weight Initialization

---

Next idea: **small random numbers** (Gaussian with zero mean, std=0.01)

```
W = 0.01 * np.random.randn(Din, Dout)
```

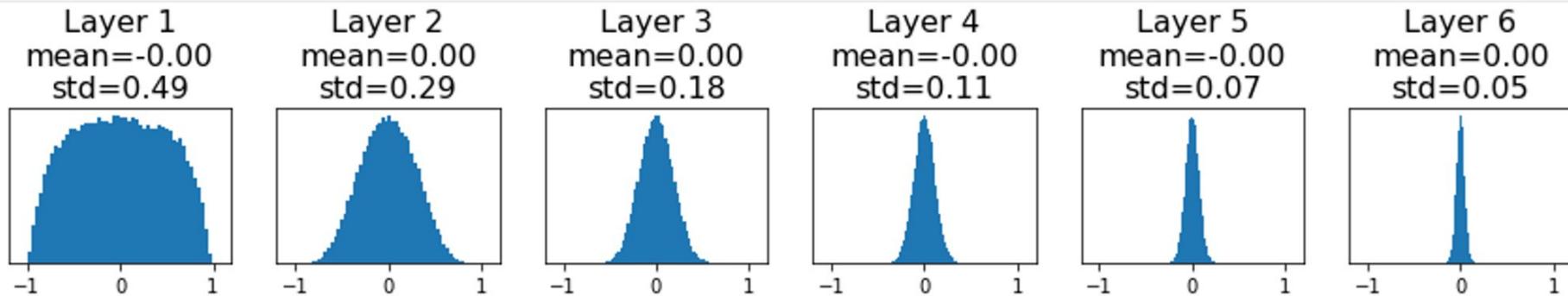
Works ~okay for small networks, but problems with deeper networks.

# Weight Initialization: Activation Statistics

```
dims = [4096] * 7    Forward pass for a 6-layer
hs = []             net with hidden size 4096
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = 0.01 * np.random.randn(Din, Dout)
    x = np.tanh(x.dot(W))
    hs.append(x)
```

All activations tend to zero for deeper network layers

Q: What do the gradients look like?

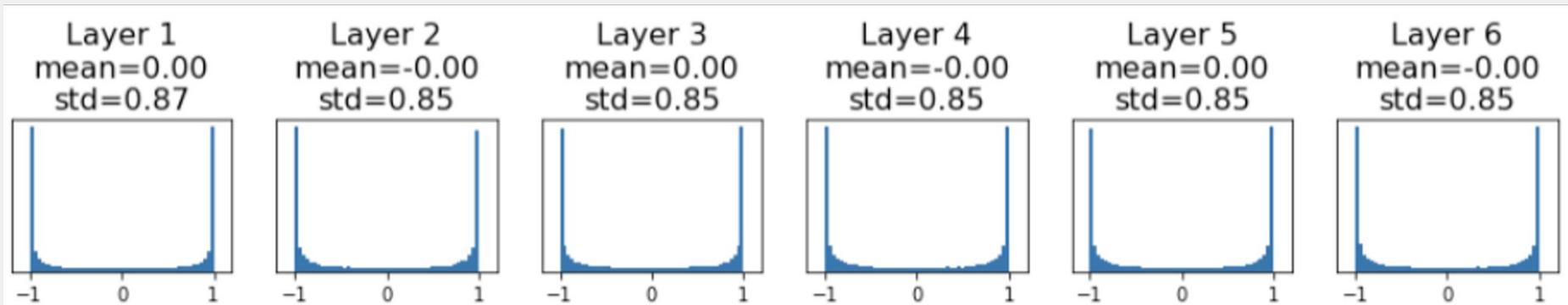


# Weight Initialization: Activation Statistics

```
dims = [4096] * 7    Increase std of initial weights  
hs = []             from 0.01 to 0.05  
x = np.random.randn(16, dims[0])  
for Din, Dout in zip(dims[:-1], dims[1:]):  
    W = 0.05 * np.random.randn(Din, Dout)  
    x = np.tanh(x.dot(W))  
    hs.append(x)
```

All activations saturate

Q: What do the gradients look like?

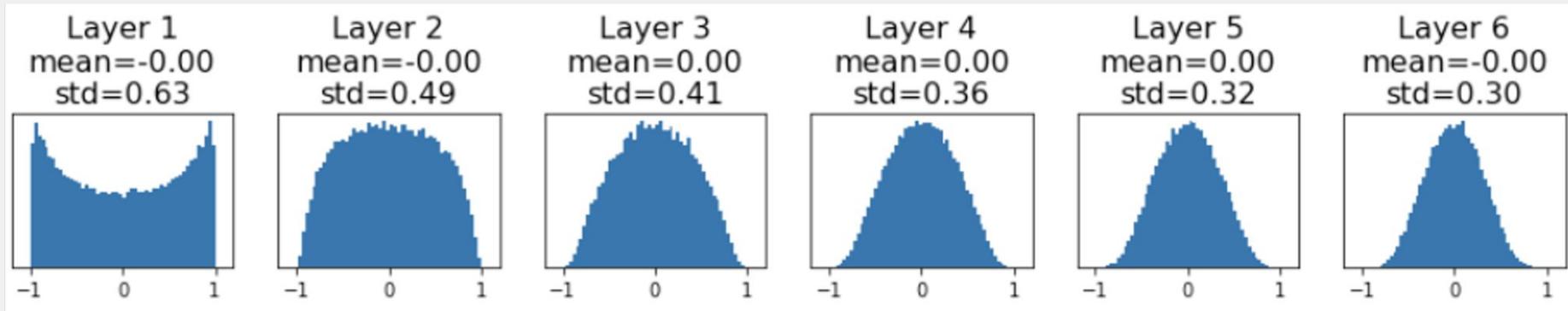


# Weight Initialization: Xavier Initialization

```
dims = [4096] * 7           "Xavier" initialization:
hs = []                     std = 1/sqrt(Din)
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = np.random.randn(Din, Dout) / np.sqrt(Din)
    x = np.tanh(x.dot(W))
    hs.append(x)
```

“Just right”: Activations are nicely scaled for all layers!

For conv layers, Din is  $\text{kernel\_size}^2 \times \text{input\_channels}$



# Weight Initialization: Xavier Initialization

**Derivation:** Variance of output = Variance of input

$$y = Wx$$

$$y_i = \sum_{j=1}^{D_{in}} x_j w_j$$

$$\text{Var}(y_i) = D_{in} * \text{Var}(x_i w_i)$$

[Assume x, w are iid]

$$= D_{in} * (E[x_i^2]E[w_i^2] - E[x_i]^2 E[w_i]^2)$$

[Assume x, w independent]

$$= D_{in} * \text{Var}(x_i) * \text{Var}(w_i)$$

[Assume x, w are zero-mean]

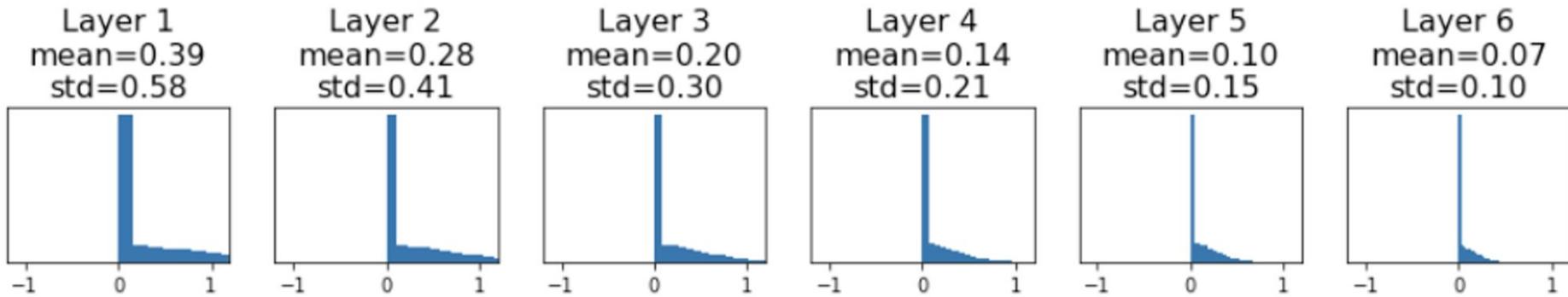
$$\text{If } \text{Var}(w_i) = 1/D_{in} \text{ then } \text{Var}(y_i) = \text{Var}(x_i)$$

# Weight Initialization: Xavier Initialization

```
dims = [4096] * 7      Change from tanh to ReLU
hs = []
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = np.random.randn(Din, Dout) / np.sqrt(Din)
    x = np.maximum(0, x.dot(W))
    hs.append(x)
```

Xavier assumes zero centered activation function

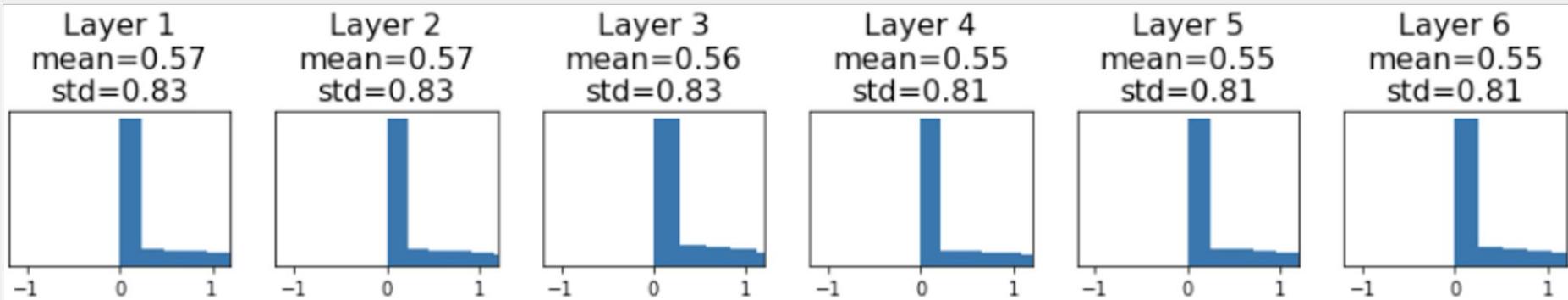
Activations collapse to zero again, no learning :(



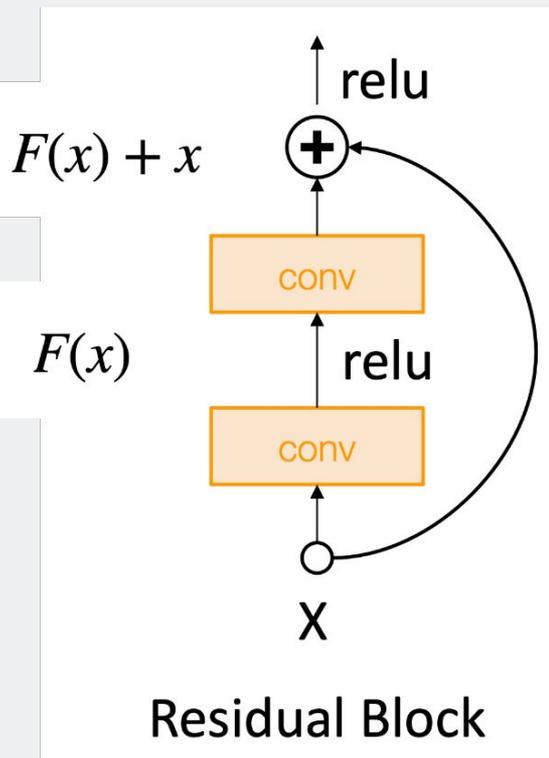
# Weight Initialization: Kaiming/MSRA initialization

```
dims = [4096] * 7 ReLU correction: std = sqrt(2 / Din)
hs = []
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = np.random.randn(Din, Dout) / np.sqrt(Din)
    x = np.maximum(0, x.dot(W))
    hs.append(x)
```

“Just right” - activations nicely scaled for all layers



# Weight Initialization: Residual Networks



If we initialize with MSRA: then  
 $Var(F(x)) = Var(x)$

But then  $Var(F(x) + x) > Var(x)$   
variance grows with each block!

**Solution:** Initialize first conv with MSRA,  
initialize second conv to zero. Then  
 $Var(F(x) + x) = Var(x)$

# Proper Initialization: Active area of research

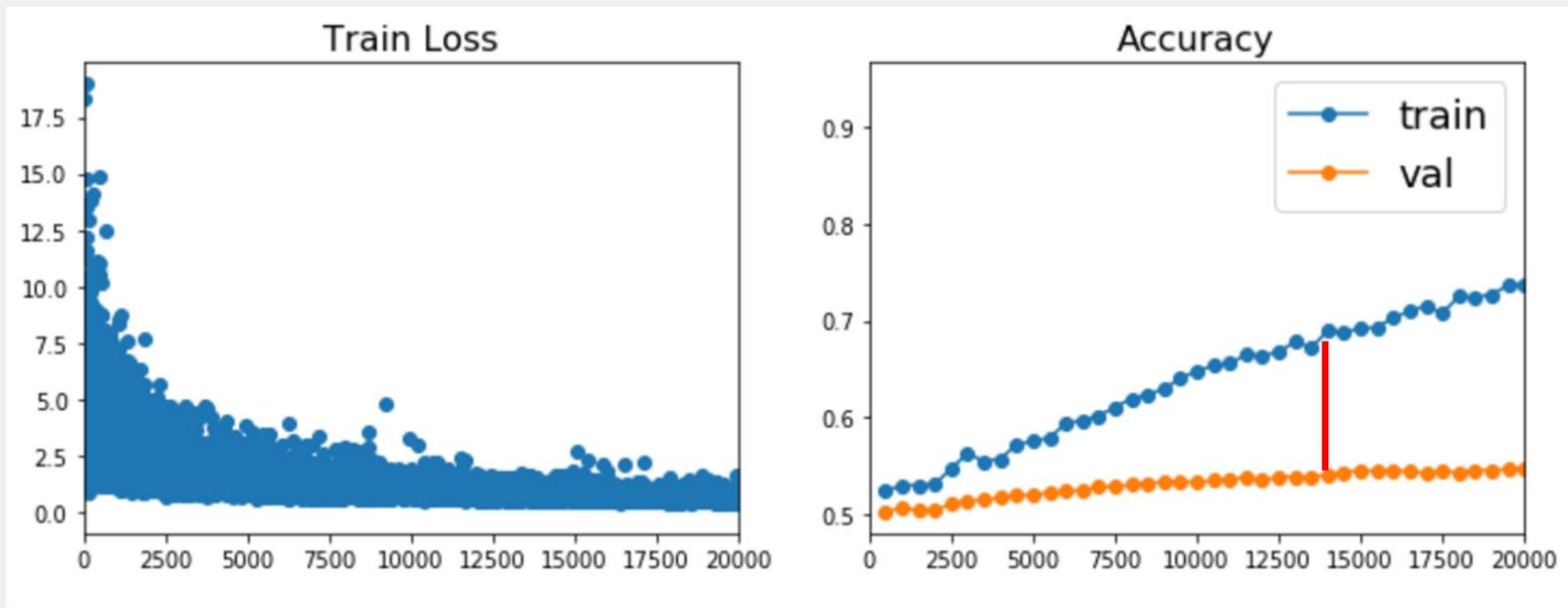
---

- Understanding the difficulty of training deep feedforward neural networks by Glorot and Bengio, 2010
- Exact solutions to the nonlinear dynamics of learning in deep linear neural networks by Saxe et al, 2013
- Random walk initialization for training very deep feedforward networks by Sussillo and Abbott, 2014
- Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification by He et al., 2015
- Data-dependent Initializations of Convolutional Neural Networks by Krähenbühl et al., 2015
- All you need is a good init, Mishkin and Matas, 2015
- Fixup Initialization: Residual Learning Without Normalization, Zhang et al, 2019
- The Lottery Ticket Hypothesis: Finding Sparse, Trainable Neural Networks, Frankle and Carbin, 2019
- .....

# Regularization (Dropout)

# Now your model is training... but it overfits!

---



## Regularization

# Recap: Regularization

---

$$L = \frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1) + \lambda R(W)$$

**In common use:**

**L2 regularization**

$$R(W) = \sum_k \sum_l W_{k,l}^2 \quad (\text{Weight decay})$$

**L1 regularization**

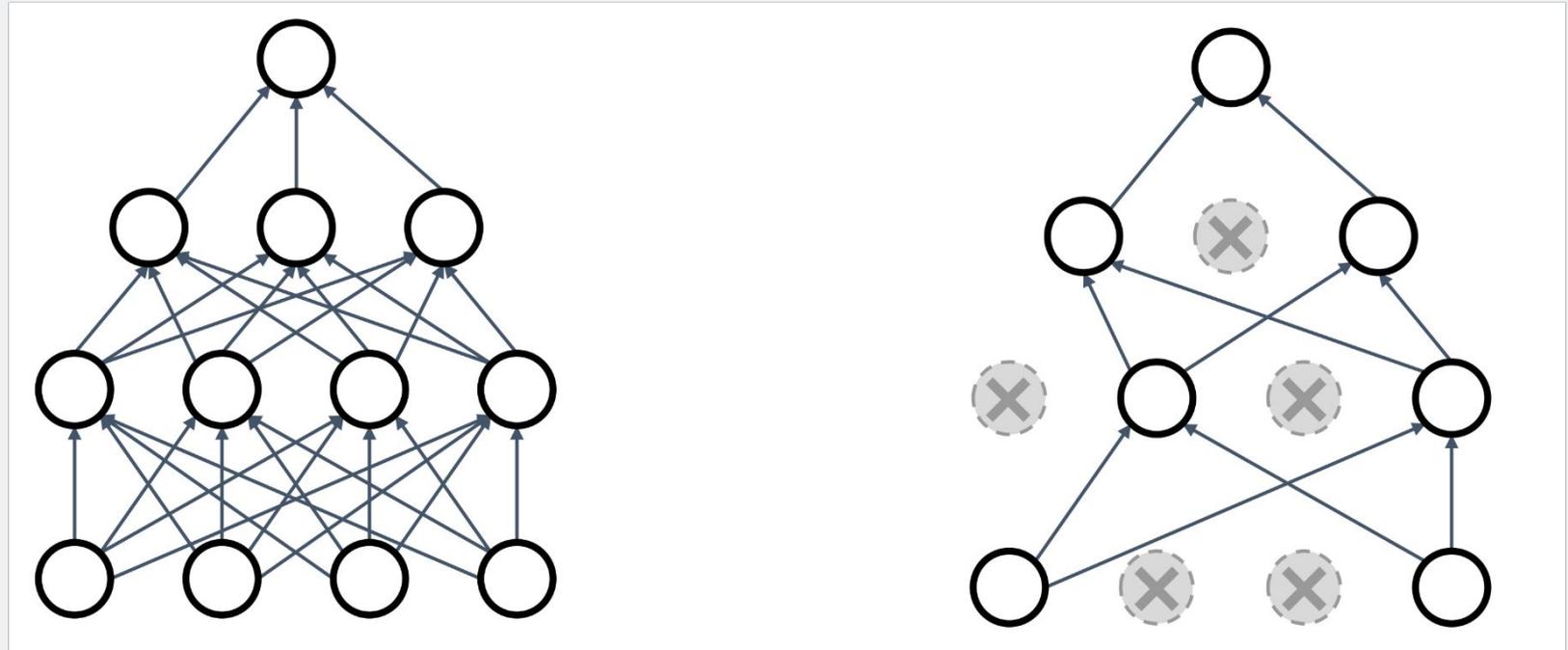
$$R(W) = \sum_k \sum_l |W_{k,l}|$$

**Elastic net (L1 + L2)**

$$R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$$

# Regularization: Dropout

In each forward pass, randomly set some neurons to zero  
Probability of dropping is a hyperparameter; 0.5 is common



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

<https://www.cs.toronto.edu/~hinton/absps/JMLRdropout.pdf>

# Regularization: Dropout

```
p = 0.5 # probability of keeping a unit active. higher = less dropout
```

```
def train_step(X):
```

```
    """ X contains the data """
```

```
    # forward pass for example 3-layer neural network
```

```
    H1 = np.maximum(0, np.dot(W1, X) + b1)
```

```
    U1 = np.random.rand(*H1.shape) < p # first dropout mask
```

```
    H1 *= U1 # drop!
```

```
    H2 = np.maximum(0, np.dot(W2, H1) + b2)
```

```
    U2 = np.random.rand(*H2.shape) < p # second dropout mask
```

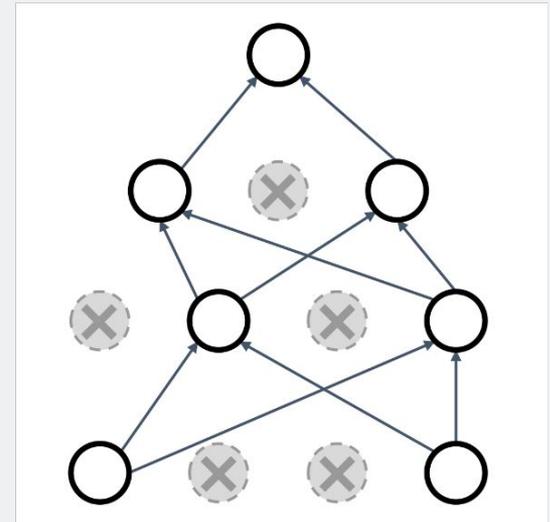
```
    H2 *= U2 # drop!
```

```
    out = np.dot(W3, H2) + b3
```

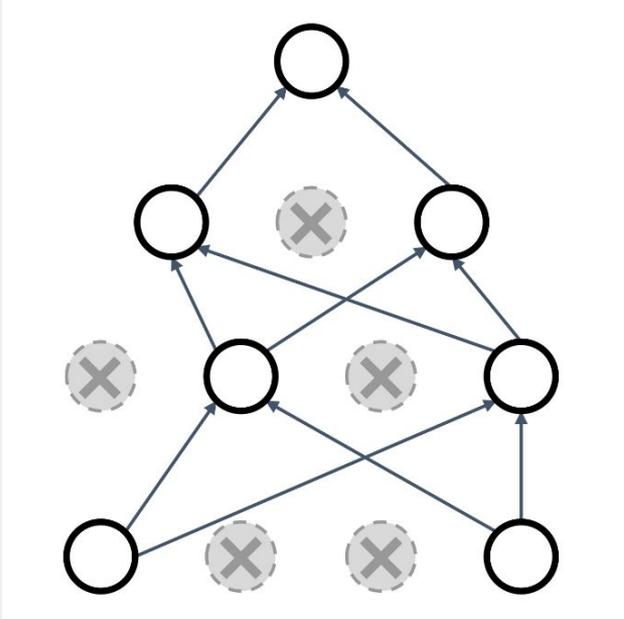
```
    # backward pass: compute gradients... (not shown)
```

```
    # perform parameter update... (not shown)
```

Example forward pass with a 3-layer network using dropout



# Interpreting Dropout

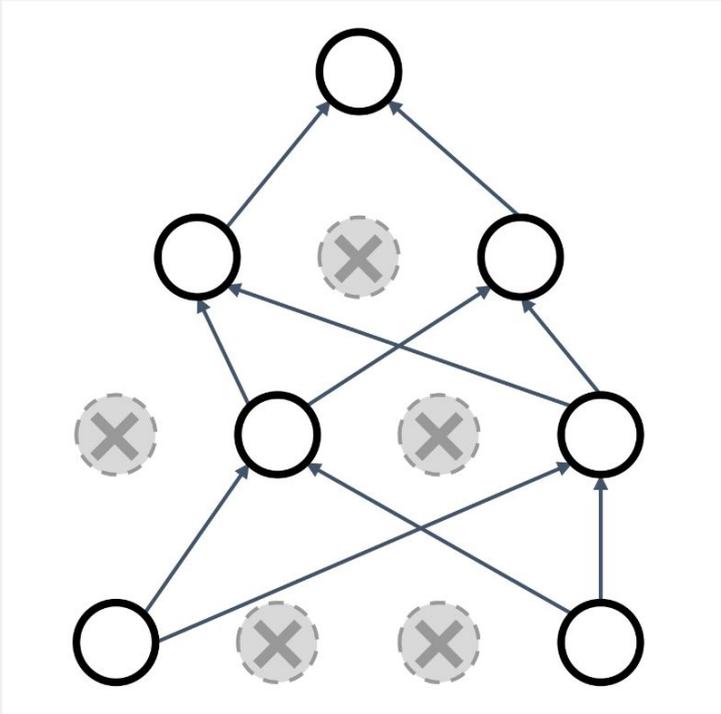


Forces the network to have a redundant representation; prevents **co-adaptation** of features



# Interpreting Dropout

---



Another interpretation:

Dropout is training a large *ensemble* of models (that share parameters).

Each binary mask is one model

An FC layer with 4096 units has  $2^{4096} \sim 10^{1233}$  possible masks!

Only  $\sim 10^{82}$  atoms in the universe...

# Dropout: Test Time

---

Dropout makes our output random!

$$\mathbf{y} = f_w(\mathbf{x}, \mathbf{z})$$

Output label      Input image      Random mask

Want to “average out” the randomness at test-time

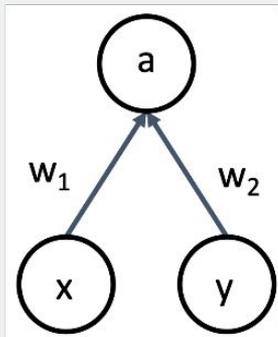
$$\mathbf{y} = f(\mathbf{x}, \mathbf{z}) = \mathbb{E}_{\mathbf{z}}[f(\mathbf{x}, \mathbf{z})] = \int p(\mathbf{z})f(\mathbf{x}, \mathbf{z})d\mathbf{z}$$

But this integral seems hard...

# Dropout: Test Time

Want to approximate  
the integral

$$y = f(x, z) = \mathbb{E}_z[f(x, z)] = \int p(z)f(x, z)dz$$



Consider a single neuron:

At test time we have:  $\mathbb{E}[a] = w_1x + w_2y$

During training time we have:  $\mathbb{E}[a] = \frac{1}{4}(w_1x + w_2y) + \frac{1}{4}(w_1x + 0y)$

$$+ \frac{1}{4}(0x + 0y) + \frac{1}{4}(0x + w_2y)$$

$$= \frac{1}{2}(w_1x + w_2y)$$

At test time, drop nothing and **multiply**  
by dropout probability

# Dropout: Test Time

---

```
def predict(X):  
    # ensembled forward pass  
    H1 = np.maximum(0, np.dot(W1, X) + b1) * p # NOTE: scale the activations  
    H2 = np.maximum(0, np.dot(W2, H1) + b2) * p # NOTE: scale the activations  
    out = np.dot(W3, H2) + b3
```

At test time all neurons are active always

=> We must scale the activations so that for each neuron:

**Output at test time = Expected output at training time**

# Dropout: Summary

---

```
""" Vanilla Dropout: Not recommended implementation (see notes below) """  
  
p = 0.5 # probability of keeping a unit active. higher = less dropout  
  
def train_step(X):  
    """ X contains the data """  
  
    # forward pass for example 3-layer neural network  
    H1 = np.maximum(0, np.dot(W1, X) + b1)  
    U1 = np.random.rand(*H1.shape) < p # first dropout mask  
    H1 *= U1 # drop!  
    H2 = np.maximum(0, np.dot(W2, H1) + b2)  
    U2 = np.random.rand(*H2.shape) < p # second dropout mask  
    H2 *= U2 # drop!  
    out = np.dot(W3, H2) + b3  
  
    # backward pass: compute gradients... (not shown)  
    # perform parameter update... (not shown)  
  
def predict(X):  
    # ensembled forward pass  
    H1 = np.maximum(0, np.dot(W1, X) + b1) * p # NOTE: scale the activations  
    H2 = np.maximum(0, np.dot(W2, H1) + b2) * p # NOTE: scale the activations  
    out = np.dot(W3, H2) + b3
```

Drop in forward pass

Scale at test time

# More common: “Inverted dropout”

```
p = 0.5 # probability of keeping a unit active. higher = less dropout
```

```
def train_step(X):
```

```
    # forward pass for example 3-layer neural network
```

```
    H1 = np.maximum(0, np.dot(W1, X) + b1)
```

```
    U1 = (np.random.rand(*H1.shape) < p) / p # first dropout mask. Notice /p!
```

```
    H1 *= U1 # drop!
```

```
    H2 = np.maximum(0, np.dot(W2, H1) + b2)
```

```
    U2 = (np.random.rand(*H2.shape) < p) / p # second dropout mask. Notice /p!
```

```
    H2 *= U2 # drop!
```

```
    out = np.dot(W3, H2) + b3
```

```
    # backward pass: compute gradients... (not shown)
```

```
    # perform parameter update... (not shown)
```

```
def predict(X):
```

```
    # ensembled forward pass
```

```
    H1 = np.maximum(0, np.dot(W1, X) + b1) # no scaling necessary
```

```
    H2 = np.maximum(0, np.dot(W2, H1) + b2)
```

```
    out = np.dot(W3, H2) + b3
```

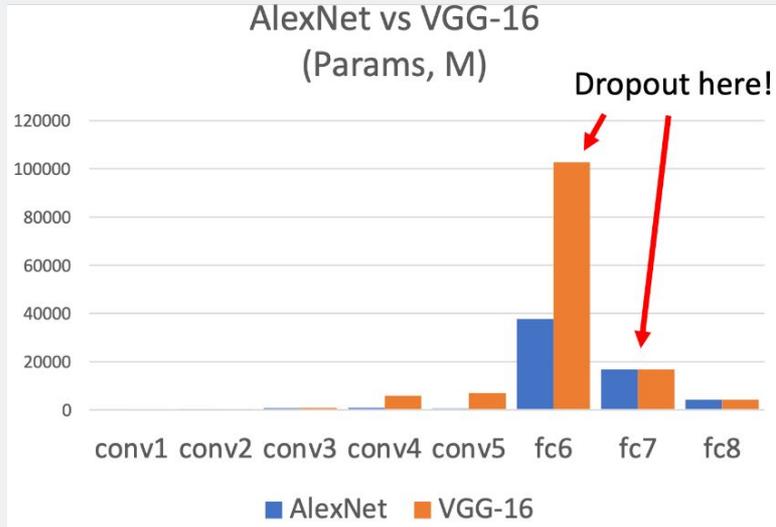
Drop and scale  
during training

test time is unchanged!

# Dropout architectures

---

Recall AlexNet, VGG have most of their parameters in **fully-connected layers**; usually Dropout is applied there



Later architectures (GoogLeNet, ResNet, etc) use global average pooling instead of fully-connected layers: they don't use dropout at all!

# Regularization: A common pattern

---

**Training:** Add some kind of randomness

$$y = f_w(x, z)$$

For ResNet and later, often L2 and Batch Normalization are the only regularizers!

**Testing:** Average out randomness (sometimes approximate)

$$y = f(x, z) = \mathbb{E}_z[f(x, z)] = \int p(z)f(x, z)dz$$

**Example:** Batch Normalization

**Training:** Normalize using stats from random mini batches

**Testing:** Use fixed stats to normalize

# Regularization: A common pattern

---

Training: Add some randomness

Testing: Marginalize over randomness

## Examples:

- Dropout
- Batch Normalization
- Data Augmentation

# Regularization: DropConnect

---

**Training:** Drop random connections between neurons (set weight=0)

**Testing:** Use all the connections

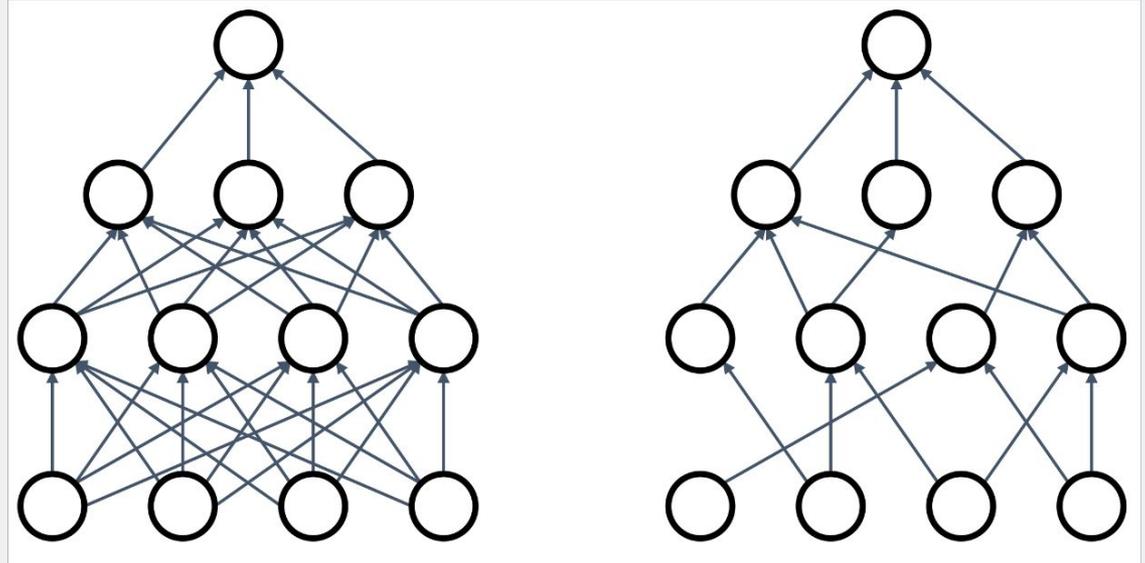
## Examples:

Dropout

Batch Normalization

Data Augmentation

DropConnect



# Regularization: Fractional Pooling

---

**Training:** Use randomized pooling regions

**Testing:** Average predictions over different samples

## Examples:

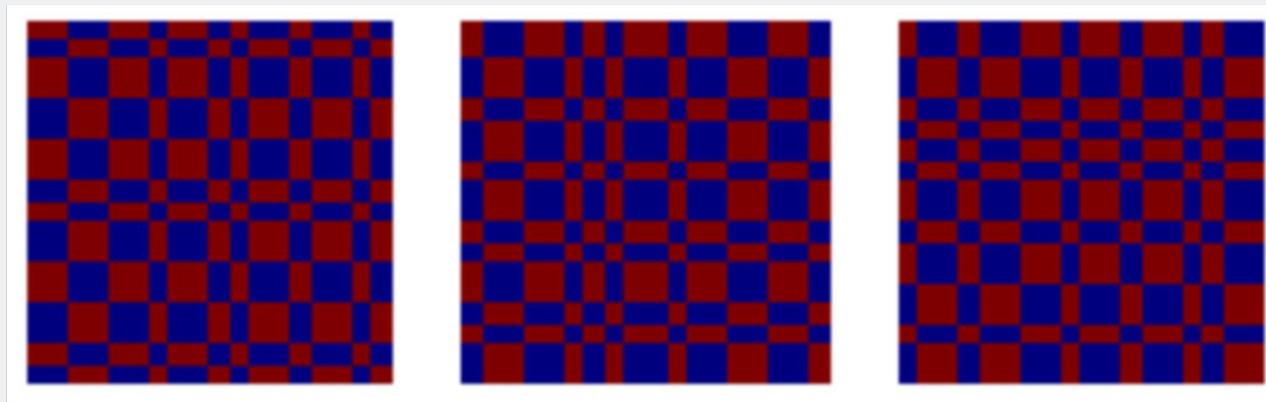
Dropout

Batch Normalization

Data Augmentation

DropConnect

Fractional Max Pooling



Graham, "Fractional Max Pooling", arXiv 2014

<https://arxiv.org/abs/1412.6071>

# Regularization: Stochastic Depth

**Training:** Skip some residual blocks in ResNet

**Testing:** Use the whole network

## Examples:

Dropout

Batch Normalization

Data Augmentation

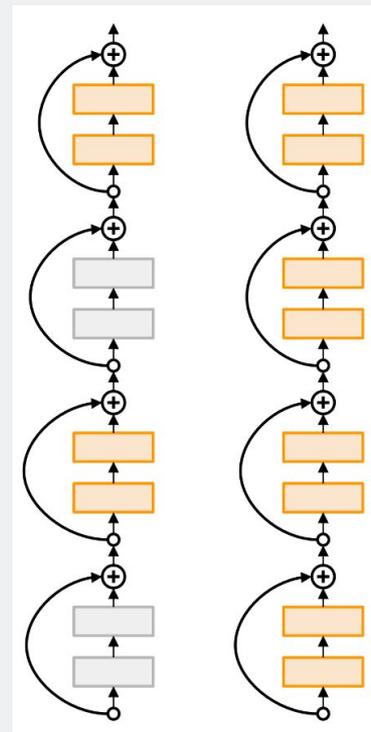
DropConnect

Fractional Max Pooling

Stochastic Depth

## Starting to become common in recent architectures:

- Pham et al, “Very Deep Self-Attention Networks for End-to-End Speech Recognition”, INTERSPEECH 2019
- Tan and Le, “EfficientNetV2: Smaller Models and Faster Training”, ICML 2021
- Fan et al, “Multiscale Vision Transformers”, ICCV 2021
- Bello et al, “Revisiting ResNets: Improved Training and Scaling Strategies”, NeurIPS 2021
- Steiner et al, “How to train your ViT? Data, Augmentation, and Regularization in Vision Transformers”, arXiv 2021



# Regularization: CutOut

---

**Training:** Set random image regions to 0

**Testing:** Use the whole image

## Examples:

Dropout

Batch Normalization

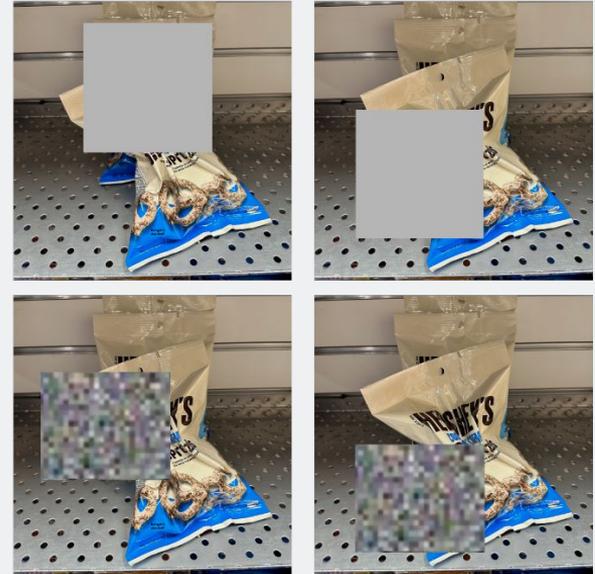
Data Augmentation

DropConnect

Fractional Max Pooling

Stochastic Depth

Cutout / Random Erasing



Replace random regions with  
mean value or random values

# Regularization: Mixup

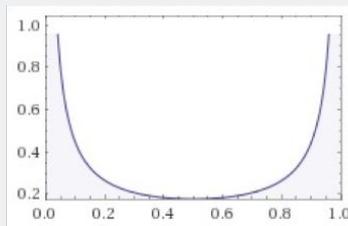
**Training:** Train on random blends of images

**Testing:** Use original images

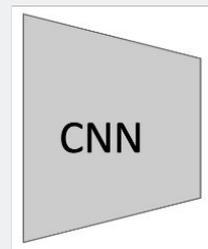
## Examples:

- Dropout
- Batch Normalization
- Data Augmentation
- DropConnect
- Fractional Max Pooling
- Stochastic Depth
- Cutout / Random Erasing

Mixup



Sample blend probability from a beta distribution  $\text{Beta}(a, b)$  with  $a=b=0$  so blend weights are close to 0/1



Target label:  
Pretzels: 0.6  
Robot: 0.4



Randomly blend the pixels of pairs of training images, e.g. 60% pretzels, 40% robot

# Regularization: CutMix

**Training:** Train on random blends of images

**Testing:** Use original images

## Examples:

Dropout

Batch Normalization

Data Augmentation

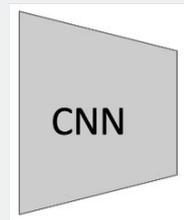
DropConnect

Fractional Max Pooling

Stochastic Depth

Cutout / Random Erasing

Mixup / CutMix



Target label:  
Pretzels: 0.6  
Robot: 0.4

Replace random crops of one image with another, e.g. 60% of pixels from pretzels, 40% from robot

# Regularization: Label Smoothing

**Training:** Train on smooth labels

**Testing:** Use original images

## Examples:

Dropout

Batch Normalization

Data Augmentation

DropConnect

Fractional Max Pooling

Stochastic Depth

Cutout / Random Erasing

Mixup / CutMix

Label Smoothing



## Standard Training

Pretzels: 100%

Robot: 0%

Sugar: 0%

## Label Smoothing

Pretzels: 90%

Robot: 5%

Sugar: 5%

Set target distribution to be  $1 - \frac{K-1}{K}\epsilon$  on the correct category and  $\epsilon/K$  on all other categories, with  $K$  categories and  $\epsilon \in (0,1)$ .

Loss is cross-entropy between predicted and target distribution.

# Data Augmentation

(example - cvpr2024)

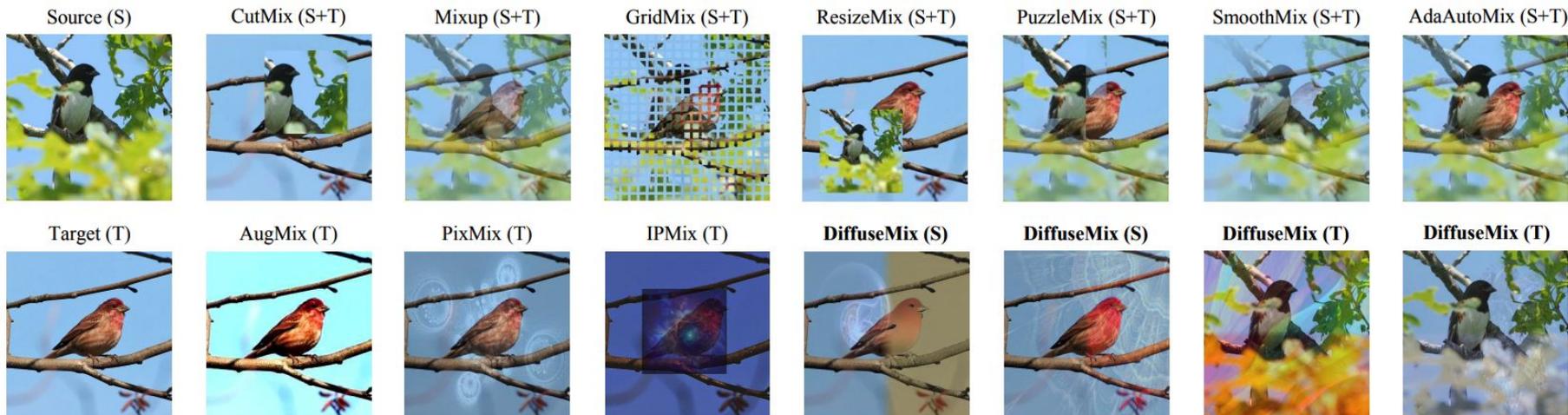


Figure 1. **Top row:** existing mixup methods *interpolate* two different training images [22, 48]. **Bottom row:** label-preserving methods. For each input image, DIFFUSEMIX employs *conditional prompts* to obtain generated images. The input image is then concatenated with a generated image to obtain a hybrid image. Each hybrid image is blended with a random fractal to obtain the final training image.

# Regularization: Summary

---

**Training:** Add some randomness

**Testing:** Marginalize over randomness

## Examples:

Dropout

Batch Normalization

Data Augmentation

DropConnect

Fractional Max Pooling

Stochastic Depth

Cutout / Random Erasing

Mixup / CutMix

Label Smoothing

- Use DropOut for large fully-connected layers
- Data augmentation is always a good idea
- Use BatchNorm for CNNs (but not ViTs)
- Try Cutout, Mixup, CutMix, Stochastic Depth, Label Smoothing to squeeze out a bit of extra performance

# Recap

P2: Due Feb.15, 2026

Canvas Quiz: Due Feb 18, 2026

## 1. One time setup:

Last time

- Activation functions, data preprocessing, weight initialization, regularization

## 2. Training dynamics:

Today

- Learning rate schedules; large-batch training; hyperparameter optimization

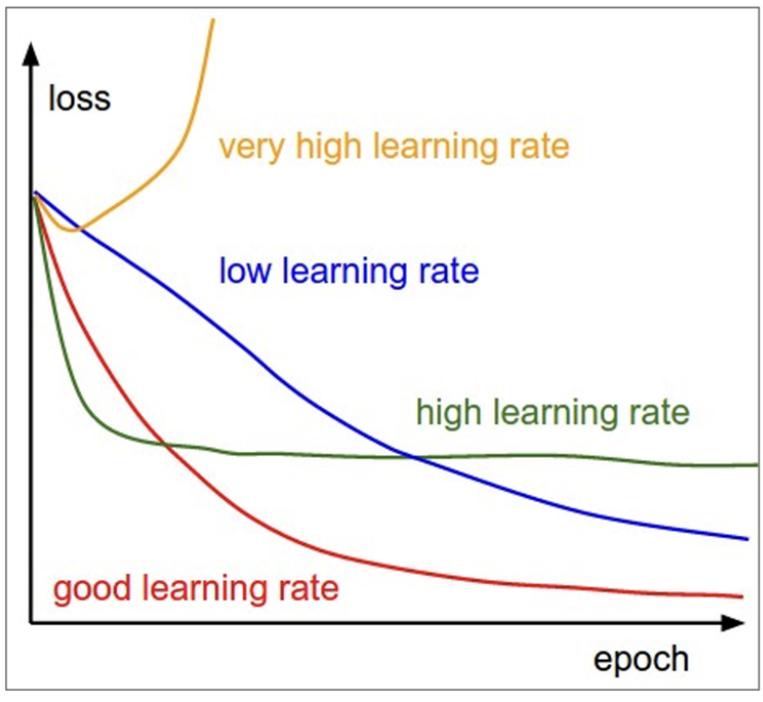
## 3. After training:

- Model ensembles, transfer learning

# Learning Rate Scheduling

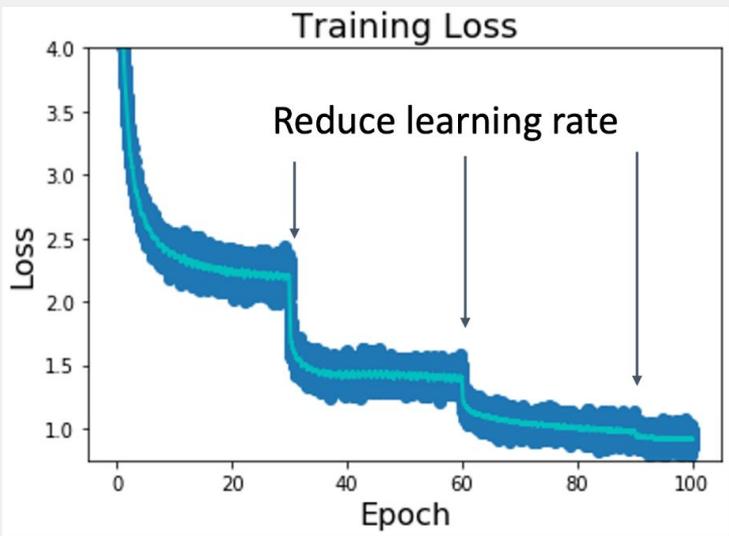
SGD, SGD+Momentum, Adagrad, RMSProp, Adam  
all have **learning rate** as hyper parameter

---

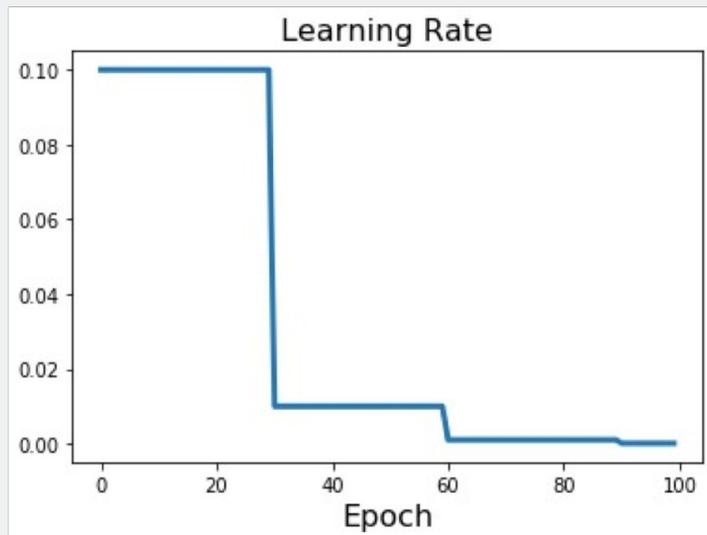


Q: Which one of these learning rates is **best** to use?

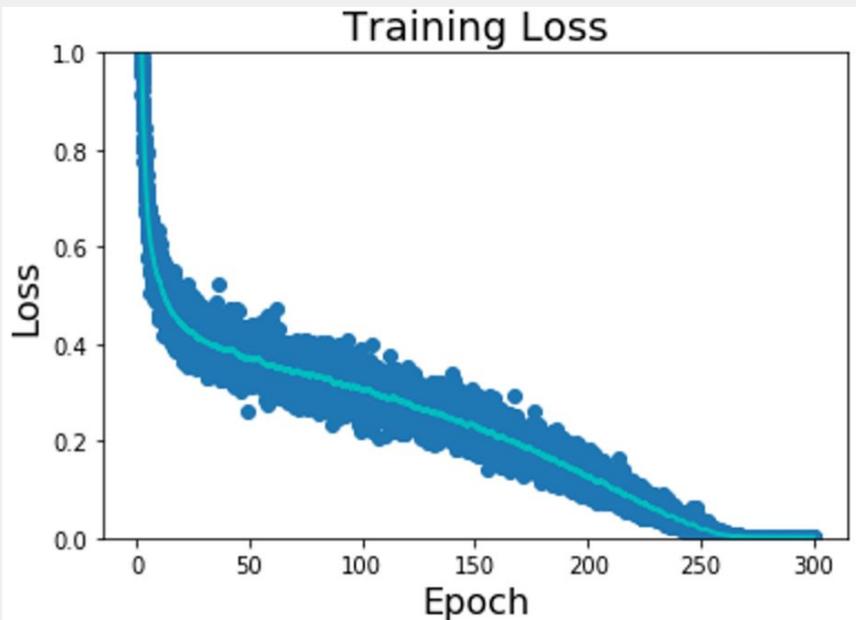
# Learning Rate Decay: Step



**Step:** Reduce learning rate at a few fixed points. E.g. for ResNets, multiply LR by 0.1 after epochs 30, 60, and 90.

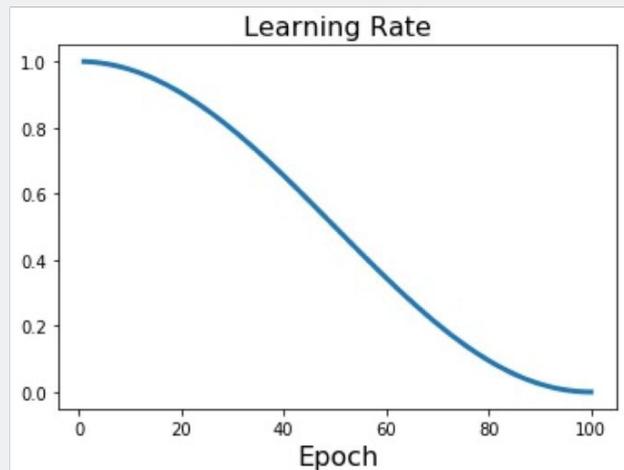


# Learning Rate Decay: Cosine



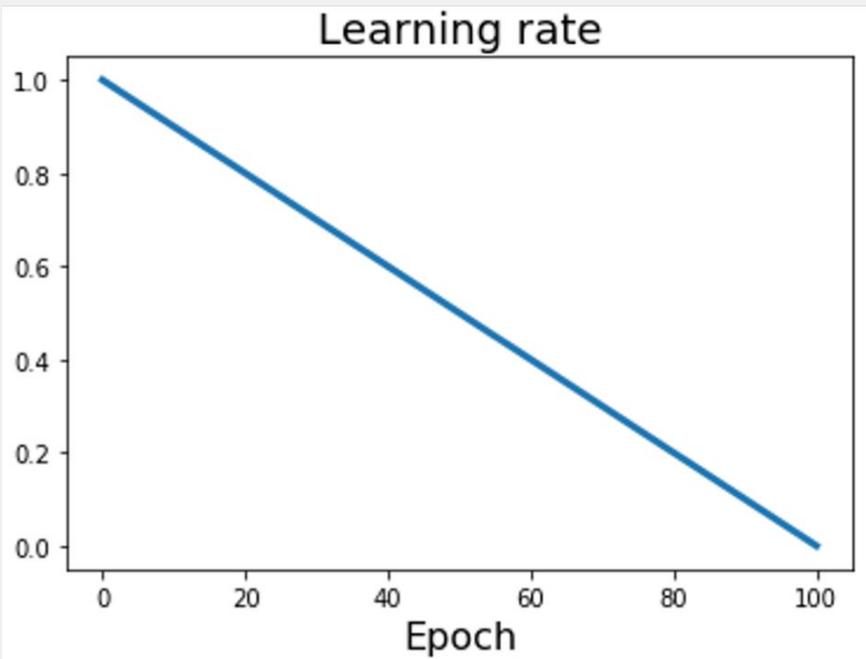
**Step:** Reduce learning rate at a few fixed points. E.g. for ResNets, multiply LR by 0.1 after epochs 30, 60, and 90.

**Cosine:** 
$$\alpha_t = \frac{1}{2}\alpha_0\left(1 + \cos\left(\frac{t\pi}{T}\right)\right)$$



Loshchilov and Hutter, "SGDR: Stochastic Gradient Descent with Warm Restarts", ICLR 2017  
Radford et al, "Improving Language Understanding by Generative Pre-Training", 2018  
Feichtenhofer et al, "SlowFast Networks for Video Recognition", ICCV 2019  
Radosavovic et al, "On Network Design Spaces for Visual Recognition", ICCV 2019  
Child et al, "Generating Long Sequences with Sparse Transformers", arXiv 2019

# Learning Rate Decay: Linear



**Step:** Reduce learning rate at a few fixed points. E.g. for ResNets, multiply LR by 0.1 after epochs 30, 60, and 90.

**Cosine:**  $\alpha_t = \frac{1}{2}\alpha_0(1 + \cos(\frac{t\pi}{T}))$

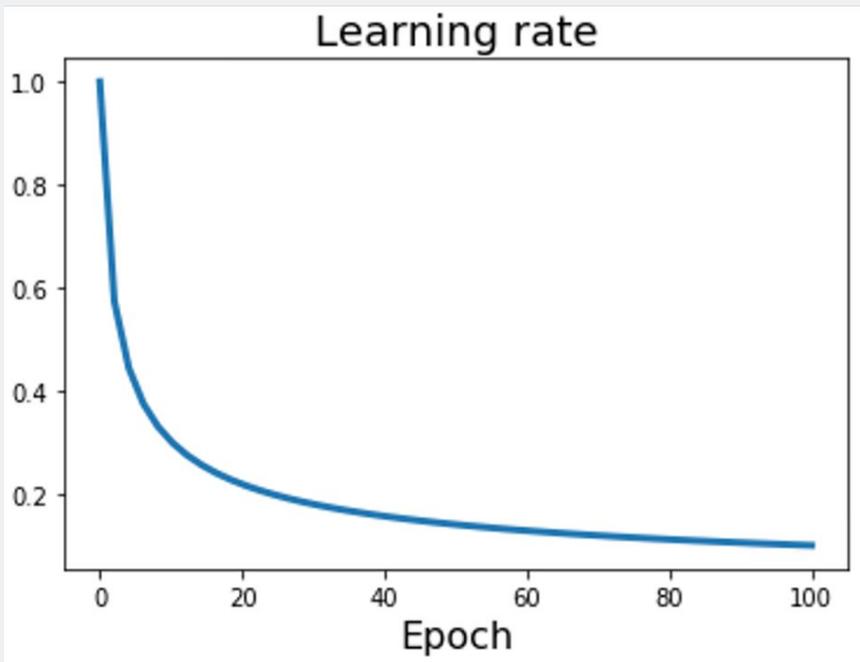
**Linear:**  $\alpha_t = \alpha_0(1 - \frac{t}{T})$

Devlin et al, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding", NAACL 2018

Liu et al, "RoBERTa: A Robustly Optimized BERT Pretraining Approach", 2019 Yang et al, "XLNet: Generalized Autoregressive

Pretraining for Language Understanding", NeurIPS 2019

# Learning Rate Decay: Inverse Sqrt



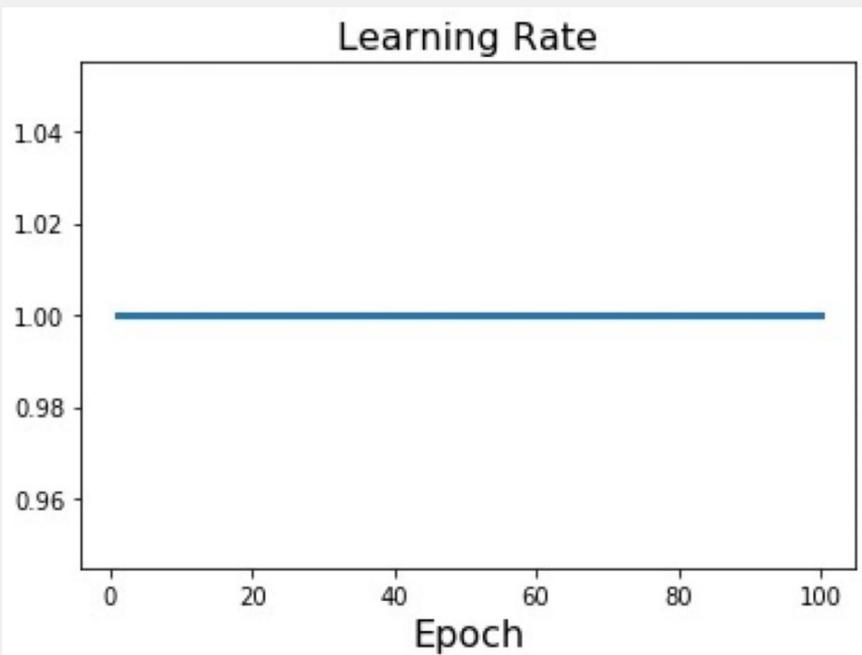
**Step:** Reduce learning rate at a few fixed points. E.g. for ResNets, multiply LR by 0.1 after epochs 30, 60, and 90.

**Cosine:**  $\alpha_t = \frac{1}{2}\alpha_0(1 + \cos(\frac{t\pi}{T}))$

**Linear:**  $\alpha_t = \alpha_0(1 - \frac{t}{T})$

**Inverse sqrt:**  $\alpha_t = \alpha_0/\sqrt{t}$

# Learning Rate Decay: Constant



**Step:** Reduce learning rate at a few fixed points. E.g. for ResNets, multiply LR by 0.1 after epochs 30, 60, and 90.

**Cosine:**  $\alpha_t = \frac{1}{2}\alpha_0(1 + \cos(\frac{t\pi}{T}))$

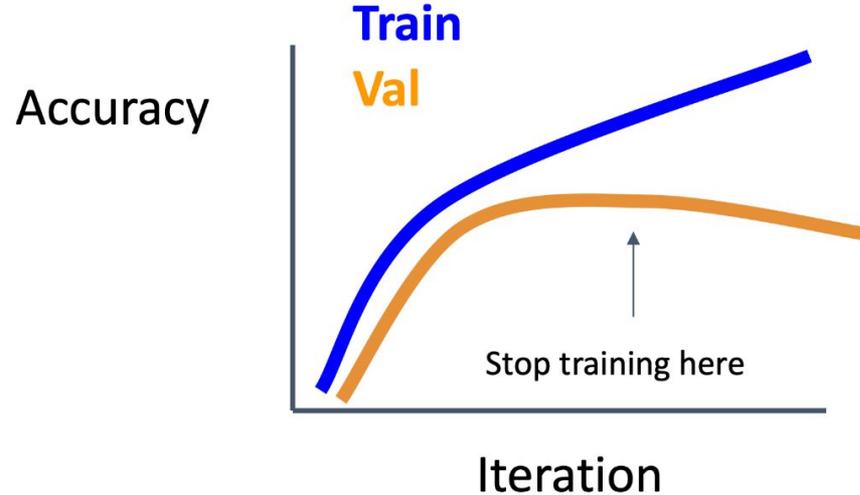
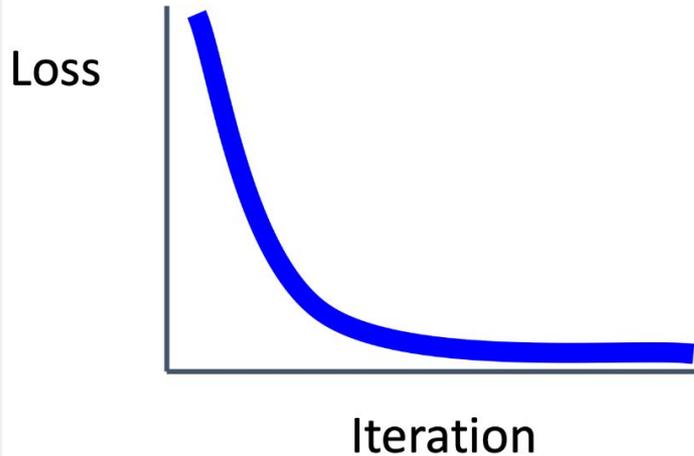
**Linear:**  $\alpha_t = \alpha_0(1 - \frac{t}{T})$

**Inverse sqrt:**  $\alpha_t = \alpha_0/\sqrt{t}$

**Constant:**  $\alpha_t = \alpha_0$

# How long to train? Early Stopping

---



Stop training the model when accuracy on the validation set decreases Or train for a long time, but always keep track of the model snapshot that worked best on val. **Always a good idea to do this!**

# Choosing Hyperparameters

# Choosing Hyperparameters: Grid Search

---

Choose several values for each hyper parameter  
(Often space choices log-linearly)

## **Example:**

Weight decay:  $[1 \times 10^{-4}, 1 \times 10^{-3}, 1 \times 10^{-2}, 1 \times 10^{-1}]$

Learning rate:  $[1 \times 10^{-4}, 1 \times 10^{-3}, 1 \times 10^{-2}, 1 \times 10^{-1}]$

Evaluate all possible choices on this **hyperparameter grid**

# Choosing Hyperparameters: Random Search

---

Choose several values for each hyper parameter  
(Often space choices log-linearly)

## **Example:**

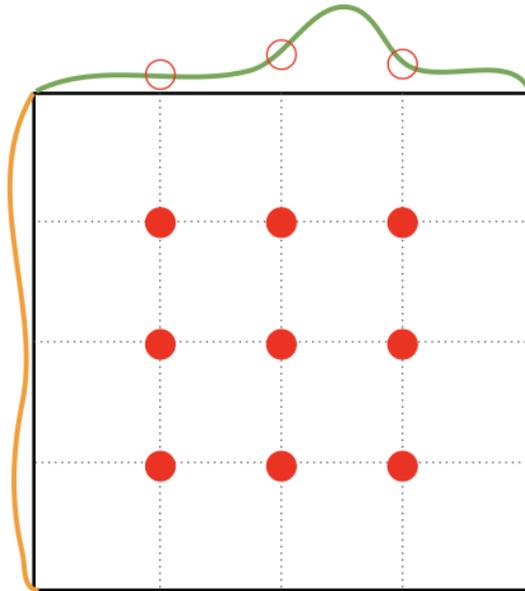
Weight decay: log-uniform on  $[1 \times 10^{-4}, 1 \times 10^{-1}]$

Learning rate: log-uniform on  $[1 \times 10^{-4}, 1 \times 10^{-1}]$

Run many different trials

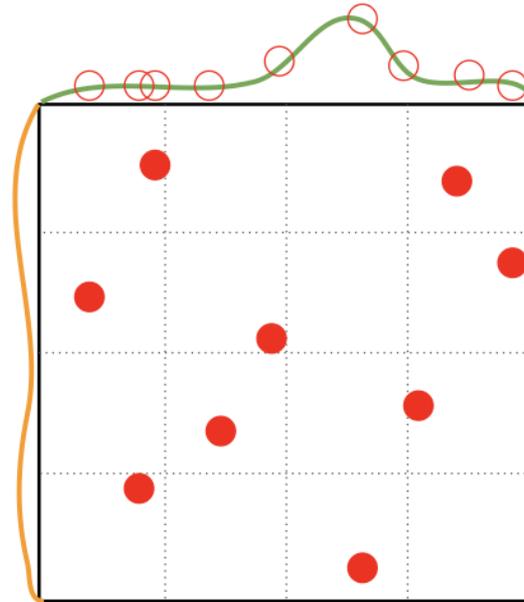
# Hyperparameters: Random vs Grid Search

**Grid Layout**



Important  
Parameter

**Random Layout**

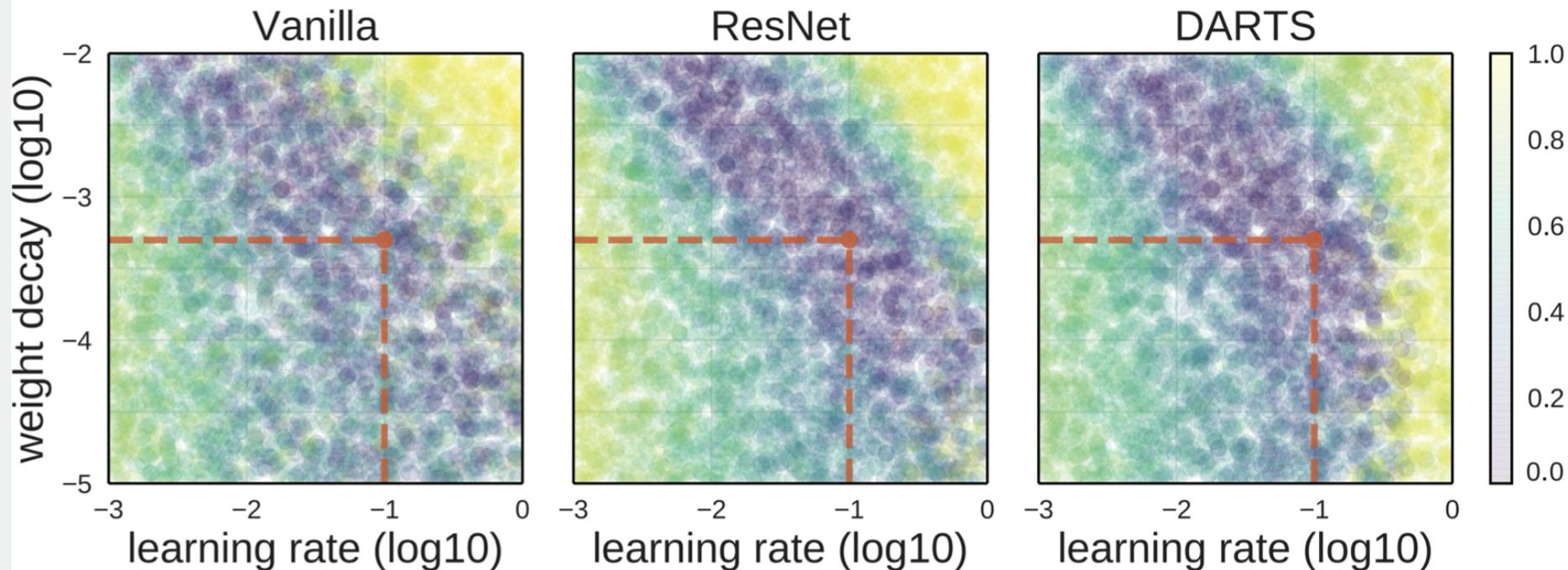


Important  
Parameter

Unimportant  
Parameter

Bergstra and Bengio,  
"Random Search for  
Hyper-Parameter  
Optimization", JMLR 2012

# Choosing Hyperparameters: Random Search



Radosavovic et al, "On Network Design Spaces for Visual Recognition", ICCV 2019

# Choosing Hyperparameters

(without tons of GPUs)

# Choosing Hyperparameters

---

## **Step 1:** Check initial loss

Turn off weight decay, sanity check loss at initialization  
e.g.  $\log(C)$  for softmax with  $C$  classes

# Choosing Hyperparameters

---

**Step 1:** Check initial loss

**Step 2:** Overfit a small sample

Try to train to 100% training accuracy on a small sample of training data (~5-10 mini batches); fiddle with architecture, learning rate, weight initialization. Turn off regularization.

Loss not going down? LR too low, bad initialization

Loss explodes to Inf or NaN? LR too high, bad initialization

# Choosing Hyperparameters

---

**Step 1:** Check initial loss

**Step 2:** Overfit a small sample

**Step 3:** Find LR that makes loss go down

Use the architecture from the previous step, use all training data, turn on small weight decay, find a learning rate that makes the loss drop significantly within ~100 iterations

Good learning rates to try:  $1e-1$ ,  $1e-2$ ,  $1e-3$ ,  $1e-4$

# Choosing Hyperparameters

---

**Step 1:** Check initial loss

**Step 2:** Overfit a small sample

**Step 3:** Find LR that makes loss go down

**Step 4:** Coarse grid, train for ~1-5 epochs

Choose a few values of learning rate and weight decay around what worked from Step 3, train a few models for ~1-5 epochs

Good learning rates to try:  $1e-4$ ,  $1e-5$ , 0

# Choosing Hyperparameters

---

**Step 1:** Check initial loss

**Step 2:** Overfit a small sample

**Step 3:** Find LR that makes loss go down

**Step 4:** Coarse grid, train for ~1-5 epochs

**Step 5:** Refine grid, train longer

Pick best models from Step 4, train them for longer (~10-20 epochs) without learning rate decay

# Choosing Hyperparameters - Summary

---

**Step 1:** Check initial loss

**Step 2:** Overfit a small sample

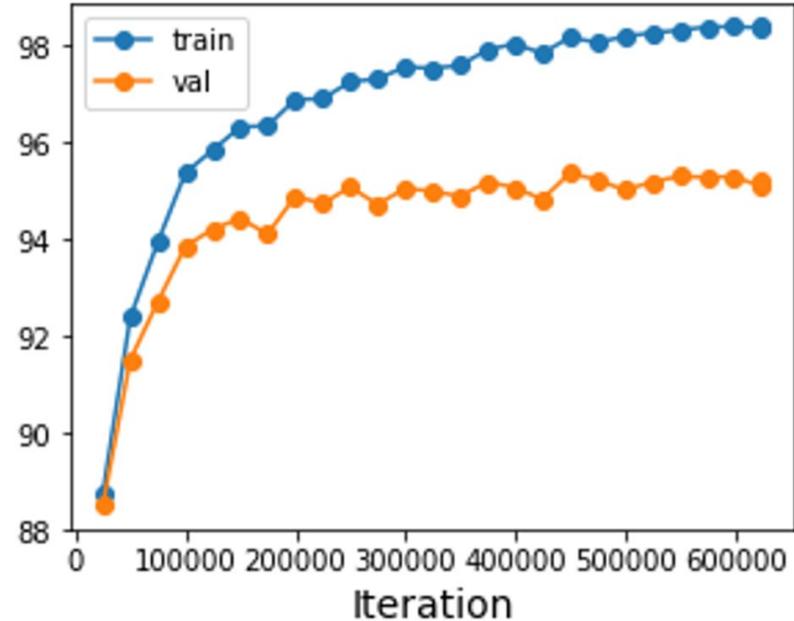
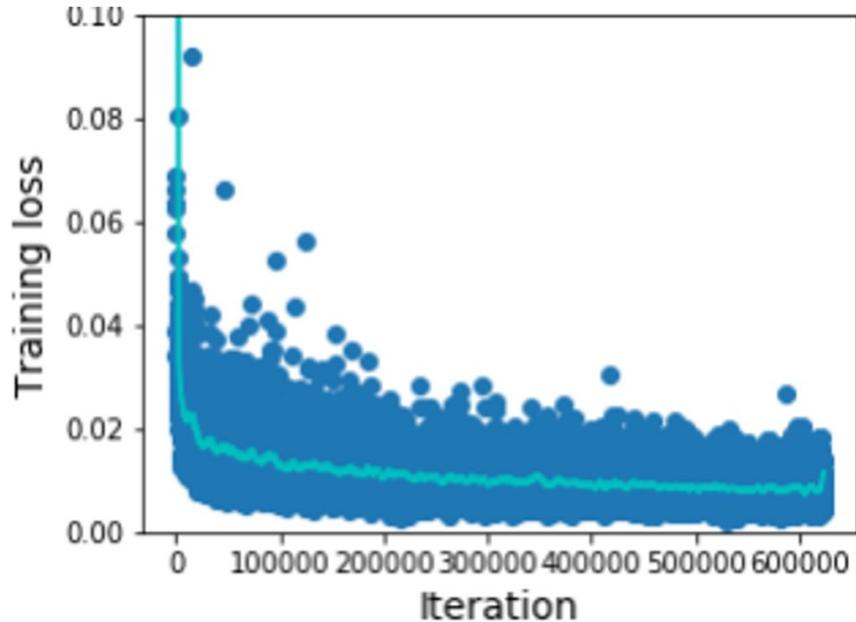
**Step 3:** Find LR that makes loss go down

**Step 4:** Coarse grid, train for ~1-5 epochs

**Step 5:** Refine grid, train longer

**Step 6:** Look at learning curves

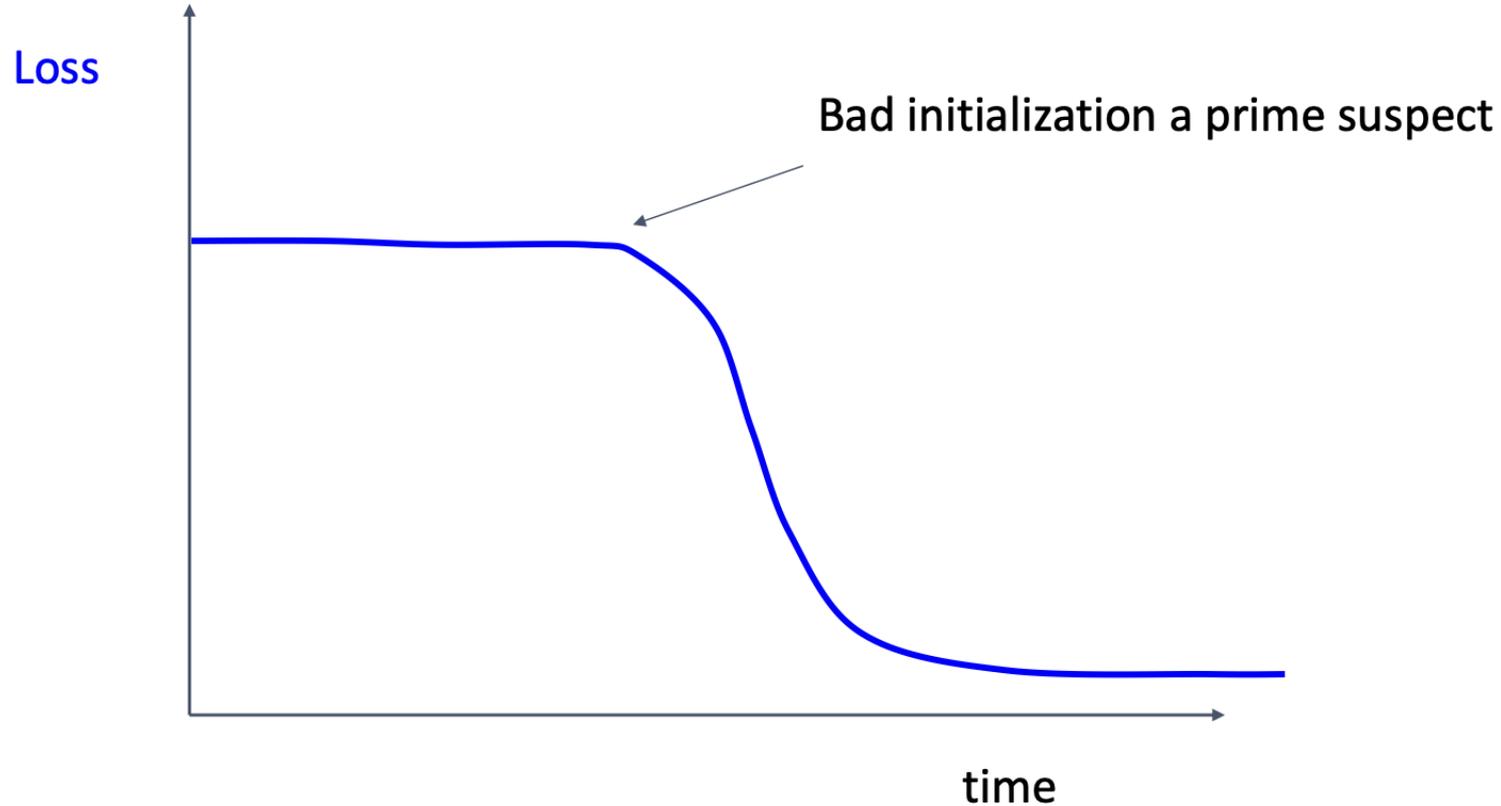
# Looking at Learning Curves



Losses may be noisy, use a scatter plot and also plot moving average to see trends better

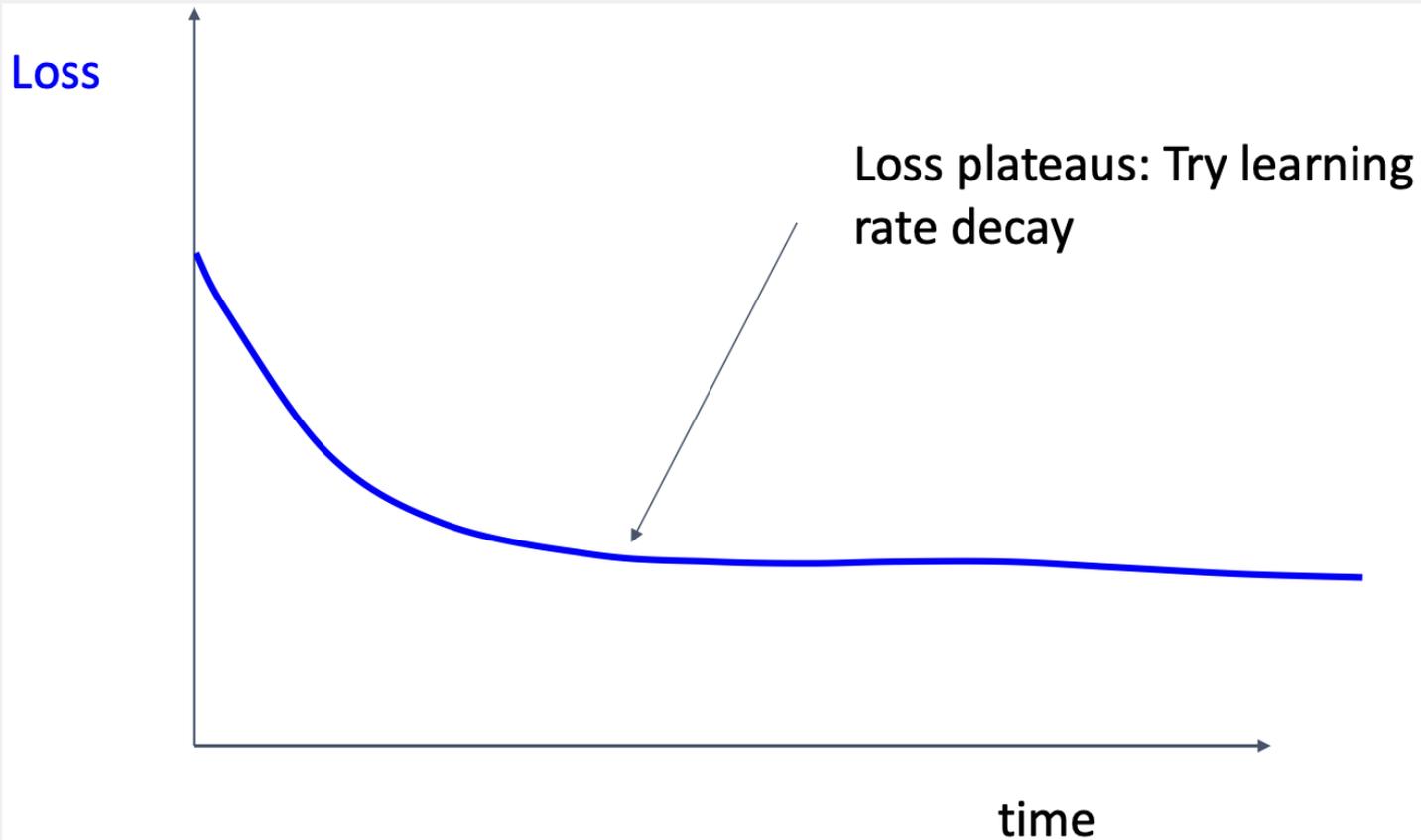
# Choosing Hyperparameters

---



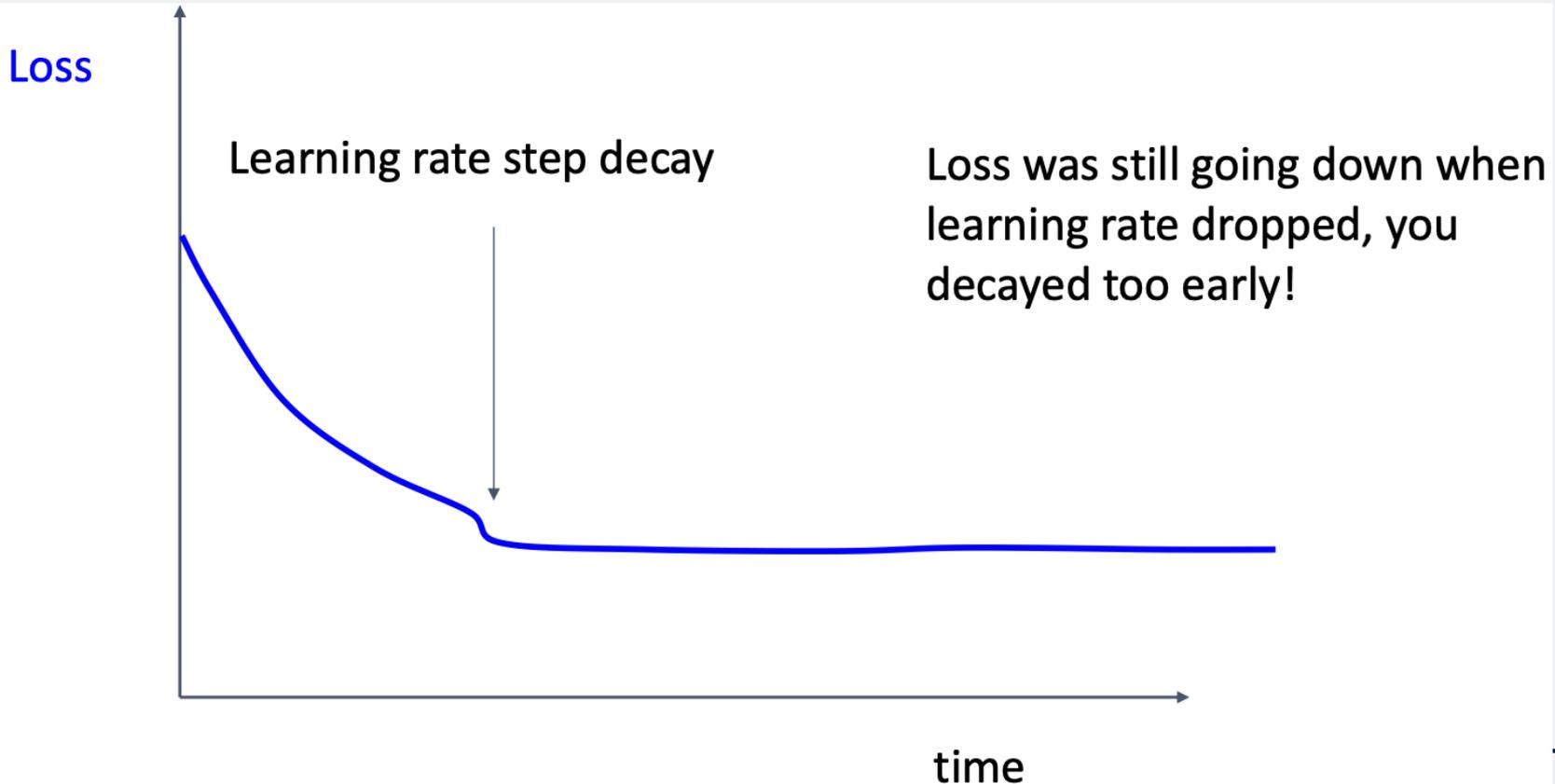
# Choosing Hyperparameters

---



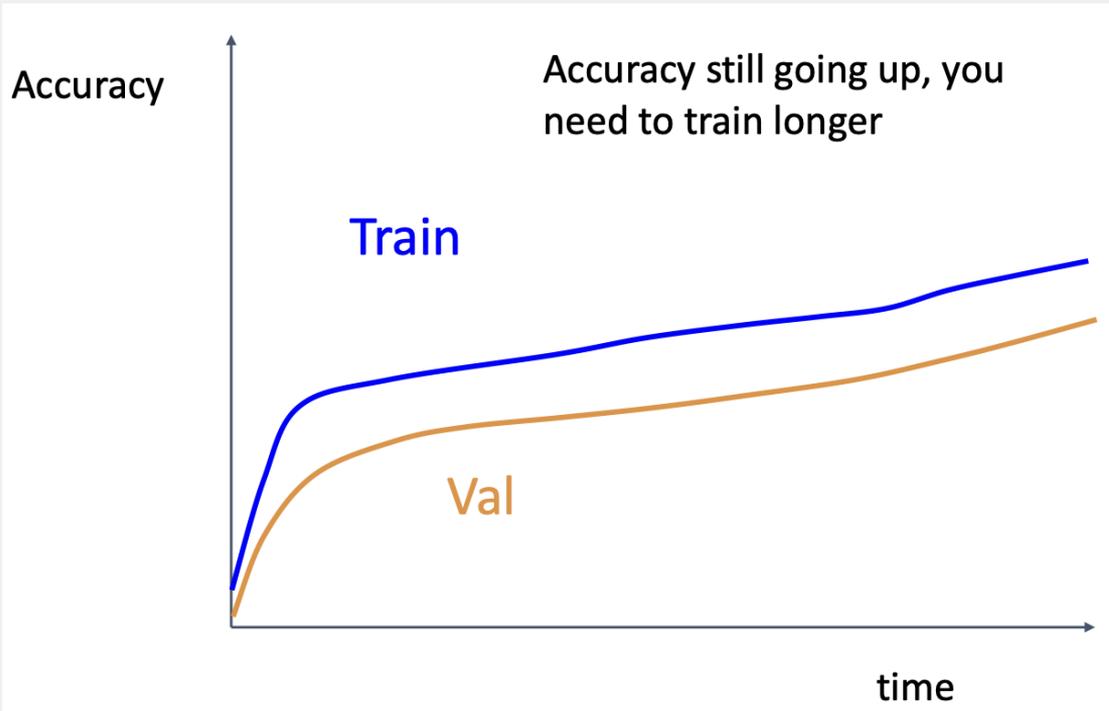
# Choosing Hyperparameters

---



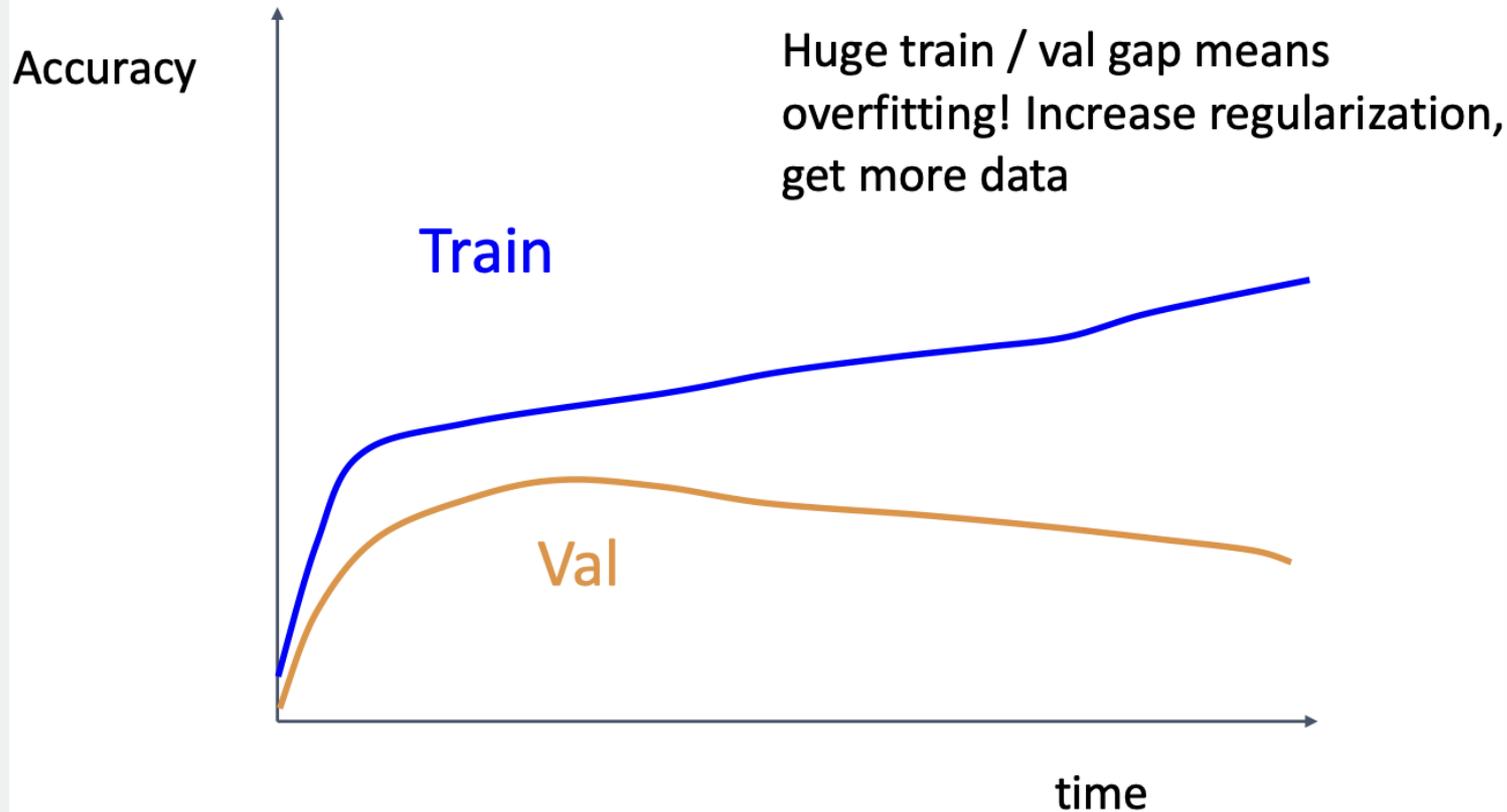
# Choosing Hyperparameters

---



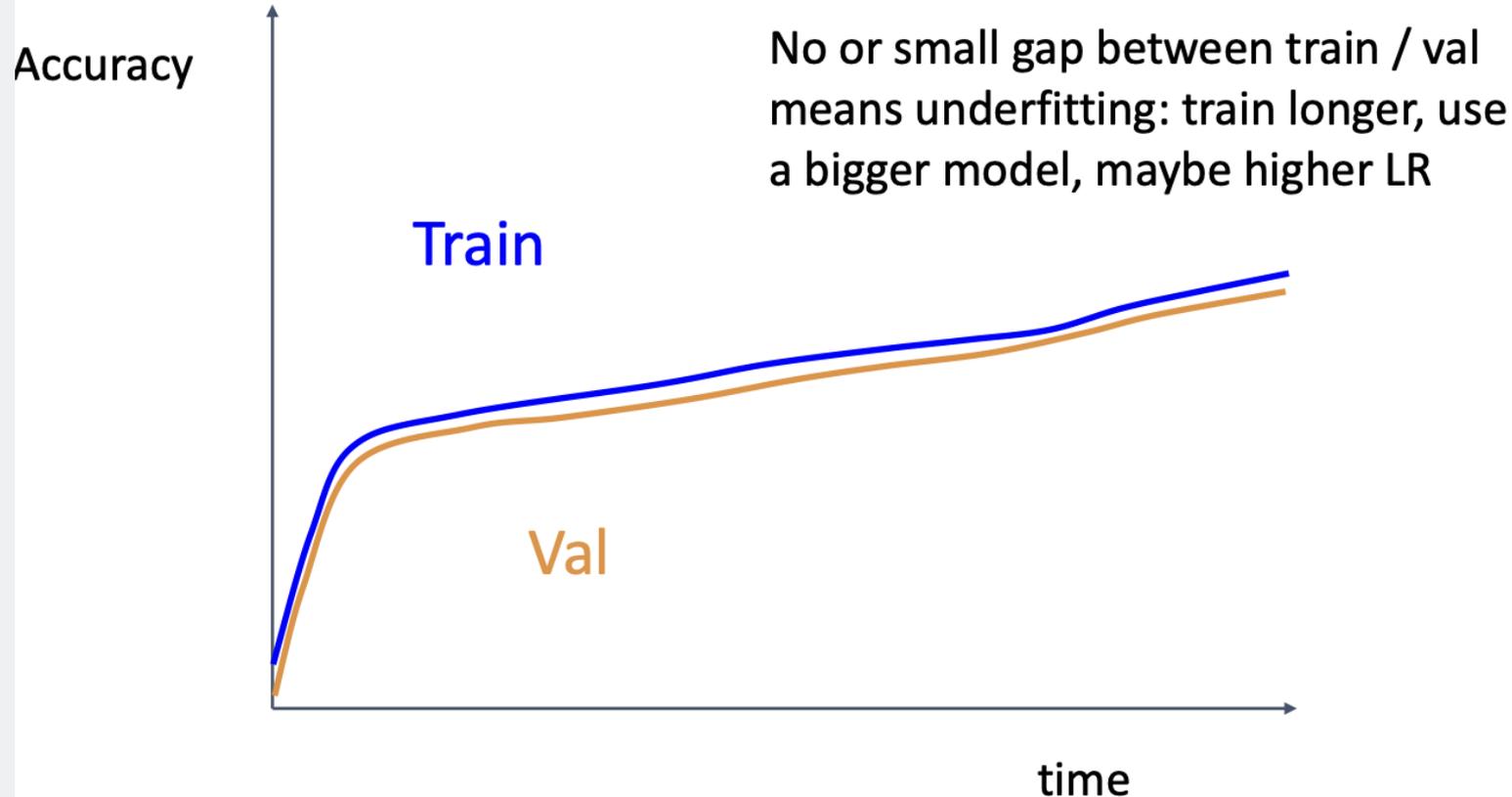
# Choosing Hyperparameters

---



# Choosing Hyperparameters

---



# Test question

---

1. You start training your Neural Network but the total loss (cross entropy loss + regularization loss) is almost completely flat from the start. What could be the cause?
  - (a) The learning rate could be too low
  - (b) The regularization strength could be too high
  - (c) The class distribution could be very uneven in the dataset
  - (d) The weight initialization scale could be incorrectly set

<https://ahaslides.com/RLCCA>



# Choosing Hyperparameters

---

Step 1: Check initial loss

Step 2: Overfit a small sample

Step 3: Find LR that makes loss go down

Step 4: Coarse grid, train for ~1-5 epochs

Step 5: Refine grid, train longer

Step 6: Look at ~~learning curves~~ loss curves

Step 7: GOTO step 5

# Track ratio of weight update / weight magnitude

---

```
# assume parameter vector W and its gradient vector dW
param_scale = np.linalg.norm(W.ravel())
update = -learning_rate*dW # simple SGD update
update_scale = np.linalg.norm(update.ravel())
W += update # the actual update
print update_scale / param_scale # want ~1e-3
```

Ratio between the updates and values:  $\sim 0.0002 / 0.02 = 0.01$  (about okay)  
**want this to be somewhere around 0.001 or so**

# Hyperparameters to play with

---

- Network architecture
- Learning rate, its decay schedule, update type
- Regularization (L2/ Dropout strength)

Neural networks practitioner  
Music = loss function

---



[This image](#) by Paolo Guereta is licensed under [CC-BY 2.0](#)

# Cross-validation “command center”

---



- Tensorboard
- wandb.ai

<https://docs.wandb.ai/tutorials/pytorch/>

# Summary

---

P2: Due Feb.15, 2026

Canvas Quiz: Due Feb 18, 2026

## 1. One time setup:

- Activation functions, data preprocessing, weight initialization, regularization

## 2. Training dynamics:

- Learning rate schedules; hyperparameter optimization

## 3. After training:

- Model ensembles, transfer learning, large-batch training

**Next Time**