# Computer Vision
# Prof. Dr. Songül Varlı

Based on notes of CS231 in Stanford University from Andrej Karpathy, Fei-Fei Li, Justin Johnson

OPTIMIZATION

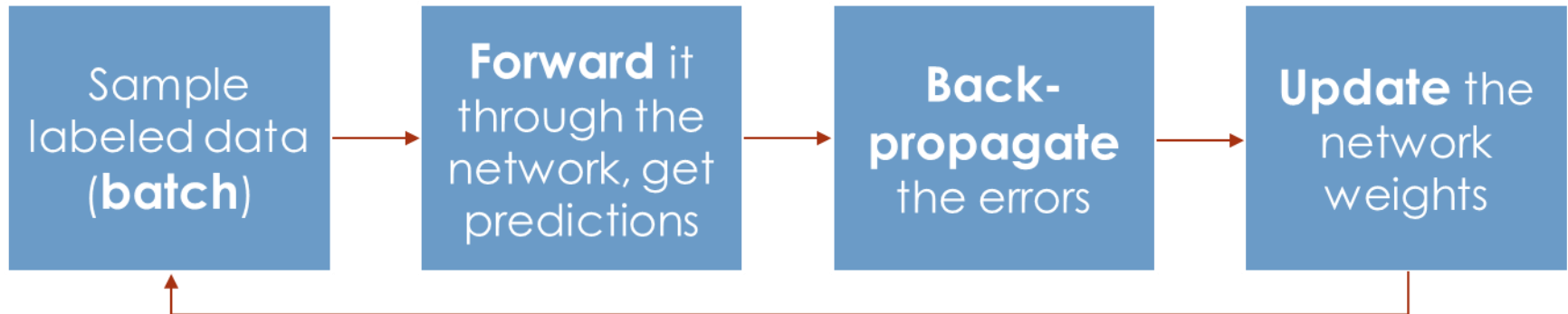REGULARIZATION

DROPOUT
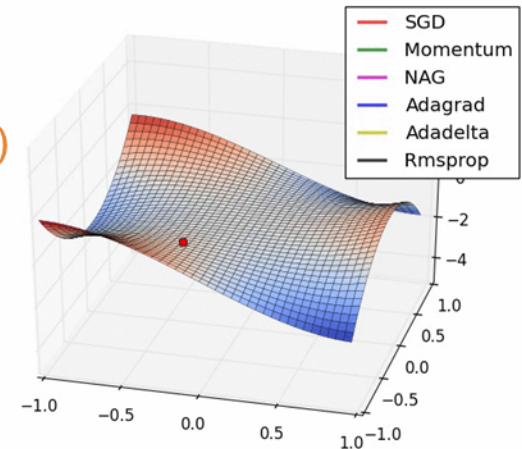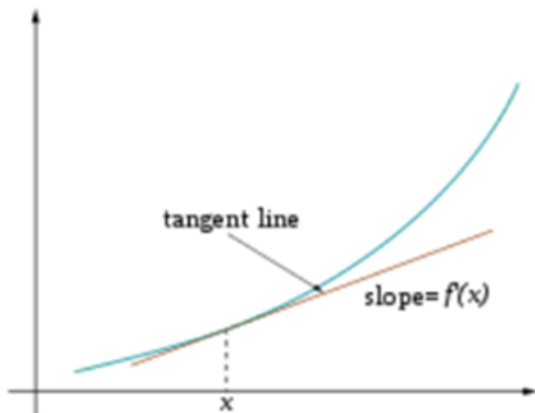
DATA AUGMENTATION

TRANSFER LEARNING

# Training

| Sample labeled data (**batch**) | **Forward** it through the network, get predictions | **Back-propagate** the errors | **Update** the network weights |
|---|---|---|---|

Optimize (min. or max.) **objective/cost function $J(\theta)$**
Generate **error signal** that measures difference between predictions and target values



tangent line
slope= f'(x)
x

Use error signal to change the **weights** and get more accurate predictions
Subtracting a fraction of the **gradient** moves you towards the **(local) minimum of the cost function**
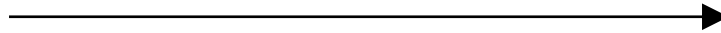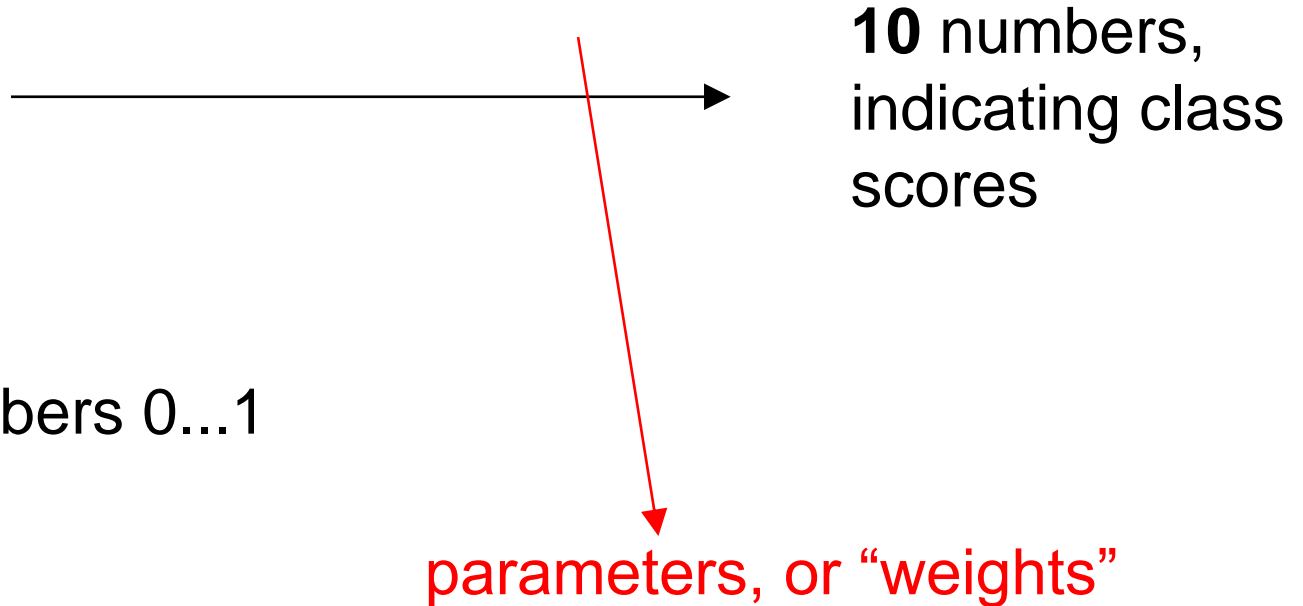
17.05.2021

$$f(x, W) = Wx$$



$\rightarrow$

**10** numbers,
indicating class
scores

**[32x32x3]**
array of numbers 0...1

$$f(x, W) = Wx$$

**3072x1**

**10x1**   **10x3072**



**[32x32x3]**
array of numbers 0...1

**10** numbers,
indicating class
scores

parameters, or "weights"

stretch pixels into single column



input image

| 0.2 | -0.5 | 0.1 | 2.0 |
|---|---|---|---|
| 1.5 | 1.3 | 2.1 | 0.0 |
| 0 | 0.25 | 0.2 | -0.3 |

$W$

| 56 |
|---|
| 231 |
| 24 |
| 2 |

$x_i$

$+$

| 1.1 |
|---|
| 3.2 |
| -1.2 |

$b$

| -96.8 | cat score |
|---|---|
| 437.9 | dog score |
| 61.95 | ship score |

$f(x_i; W, b)$

# Going forward: Loss functions/optimization



| | | | |
|---|---|---|---|
| airplane | -3.45 | -0.51 | 3.42 |
| automobile | -8.87 | **6.04** | 4.64 |
| bird | 0.09 | 5.31 | 2.65 |
| cat | **2.9** | -4.22 | 5.1 |
| deer | 4.48 | -4.19 | 2.64 |
| dog | 8.02 | 3.58 | 5.55 |
| frog | 3.78 | 4.49 | **-4.34** |
| horse | 1.06 | -4.37 | -1.5 |
| ship | -0.36 | -2.09 | -4.79 |
| truck | -0.72 | -2.93 | 6.14 |

TODO:

1. Define a **loss function** that quantifies our unhappiness with the scores across the training data.

1. Come up with a way of efficiently finding the parameters that minimize the loss function. **(optimization)**

# Loss functions/optimization

Suppose: 3 training examples, 3 classes.
For some W the scores $f(x, W) = Wx$ are:



|       |       |       |       |
|-------|-------|-------|-------|
| cat   | **3.2** | 1.3   | 2.2   |
| car   | 5.1   | **4.9** | 2.5   |
| frog  | -1.7  | 2.0   | **-3.1** |

Given an example $(x_i, y_i)$ where $x_i$ is the image and where $y_i$ is the (integer) label, and using the shorthand for the scores vector:

$$s = f(x_i, W)$$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max\left(0, s_j - s_{y_i} + 1\right)$$

# Loss functions/optimization

Suppose: 3 training examples, 3 classes.
For some W the scores $f(x, W) = Wx$ are:



|  | cat | car | frog |
|---|---|---|---|
| cat | **3.2** | 1.3 | 2.2 |
| car | 5.1 | **4.9** | 2.5 |
| frog | -1.7 | 2.0 | **-3.1** |
| Losses: | 2.9 |  |  |

Given an example $(x_i, y_i)$ where $x_i$ is the image and where $y_i$ is the (integer) label, and using the shorthand for the scores vector:

$$s = f(x_i, W)$$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

= max(0, 5.1 - 3.2 + 1)
   +max(0, -1.7 - 3.2 + 1)
= max(0, 2.9) + max(0, -3.9)
= 2.9 + 0
= 2.9

# Loss functions/optimization

Suppose: 3 training examples, 3 classes.
For some W the scores $f(x, W) = Wx$ are:

|  | cat | car | frog |
|------|------|------|------|
| cat | **3.2** | 1.3 | 2.2 |
| car | 5.1 | **4.9** | 2.5 |
| frog | -1.7 | 2.0 | **-3.1** |
| Losses: | 2.9 | 0 | |

Given an example $(x_i, y_i)$ where $x_i$ is the image and where $y_i$ is the (integer) label, and using the shorthand for the scores vector:

$$s = f(x_i, W)$$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

= max(0, 1.3 - 4.9 + 1)
   +max(0, 2.0 - 4.9 + 1)
= max(0, -2.6) + max(0, -1.9)
= 0 + 0
= 0

# Loss functions/optimization

Suppose: 3 training examples, 3 classes.
For some W the scores $f(x, W) = Wx$ are:



|        | cat   | car   | frog  |
|--------|-------|-------|-------|
| cat    | **3.2** | 1.3   | 2.2   |
| car    | 5.1   | **4.9** | 2.5   |
| frog   | -1.7  | 2.0   | **-3.1** |
| Losses: | 2.9  | 0     | 10.9  |

Given an example $(x_i, y_i)$ where $x_i$ is the image and where $y_i$ is the (integer) label, and using the shorthand for the scores vector:
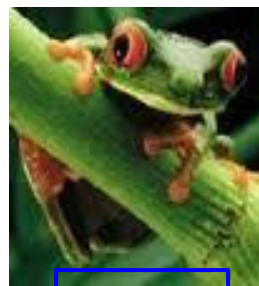
$$s = f(x_i, W)$$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

= max(0, 2.2 - (-3.1) + 1)
   +max(0, 2.5 - (-3.1) + 1)
= max(0, 5.3) + max(0, 5.6)
= 5.3 + 5.6
= 10.9

# Loss functions/optimization

Suppose: 3 training examples, 3 classes. For some W the scores $f(x, W) = Wx$ are:



|        | cat   | car  | frog  |
|--------|-------|------|-------|
| cat    | **3.2** | 1.3  | 2.2   |
| car    | 5.1   | **4.9** | 2.5   |
| frog   | -1.7  | 2.0  | **-3.1** |
| Losses: | 2.9   | 0    | 10.9  |

Given an example $(x_i, y_i)$ where $x_i$ is the image and where $y_i$ is the (integer) label, and using the shorthand for the scores vector:

$$s = f(x_i, W)$$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max\left(0, s_j - s_{y_i} + 1\right)$$

and the full training loss is the mean over all the examples:

$$L = \frac{1}{N} \sum_{i=1}^{N} L_i$$

L = (2.9 + 0 + 10.9)/3  = **4.6**

# Example numpy Code

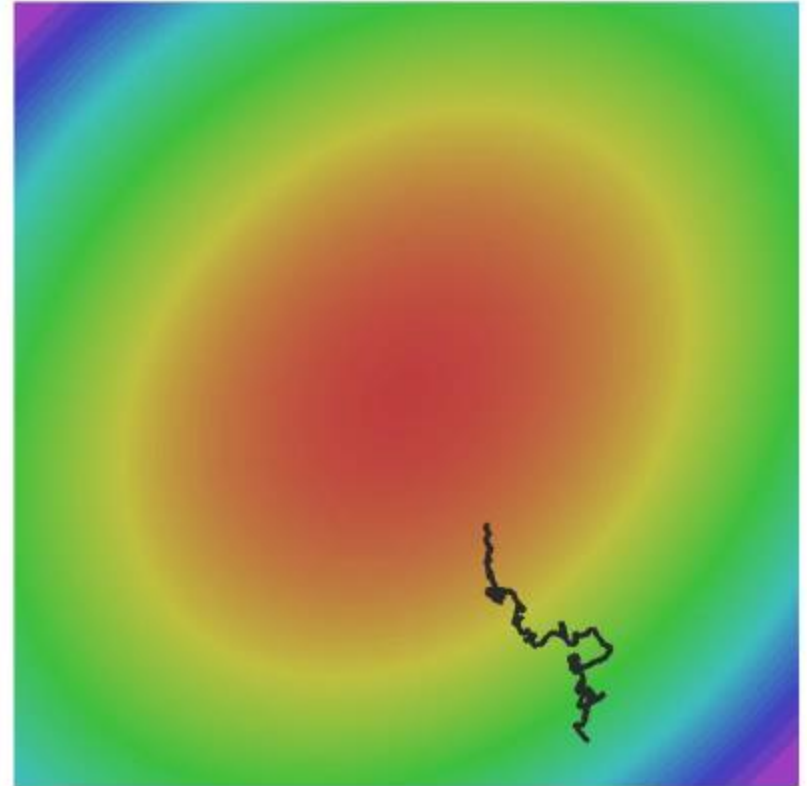$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

```python
def L_i_vectorized(x, y, W):
    scores = W.dot(x)
    margins = np.maximum(0, scores - scores[y] + 1)
    margins[y] = 0
    loss_i = np.sum(margins)
    return loss_i
```

# Optimization: Problems with SGD

Our gradients come from minibatches so they can be noisy!

$$L(W) = \frac{1}{N} \sum_{i=1}^{N} L_i(x_i, y_i, W)$$

$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^{N} \nabla_W L_i(x_i, y_i, W)$$



https://en.wikipedia.org/wiki/Stochastic_gradient_descent

# Example [edit]

Let's suppose we want to fit a straight line $\hat{y} = w_1 + w_2 x$ to a training set with observations $(x_1, x_2, \ldots, x_n)$ and corresponding estimated responses $(\hat{y}_1, \hat{y}_2, \ldots, \hat{y}_n)$ using least squares. The objective function to be minimized is:

$$Q(w) = \sum_{i=1}^{n} Q_i(w) = \sum_{i=1}^{n} \left(\hat{y}_i - y_i\right)^2 = \sum_{i=1}^{n} \left(w_1 + w_2 x_i - y_i\right)^2.$$

The last line in the above pseudocode for this specific problem will become:

$$\begin{bmatrix} w_1 \\ w_2 \end{bmatrix} := \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} - \eta \begin{bmatrix} \frac{\partial}{\partial w_1}(w_1 + w_2 x_i - y_i)^2 \\ \frac{\partial}{\partial w_2}(w_1 + w_2 x_i - y_i)^2 \end{bmatrix} = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} - \eta \begin{bmatrix} 2(w_1 + w_2 x_i - y_i) \\ 2 x_i (w_1 + w_2 x_i - y_i) \end{bmatrix}.$$

Note that in each iteration (also called update), only the gradient evaluated at a single point $x_i$ instead of evaluating at the set of all samples.

The key difference compared to standard (Batch) Gradient Descent is that only one piece of data from the dataset is used to calculate the step, and the piece of data is picked randomly at each step.

https://en.wikipedia.org/wiki/Stochastic_gradient_descent

# SGD + Momentum

## SGD

$$x_{t+1} = x_t - \alpha \nabla f(x_t)$$

```
while True:
    dx = compute_gradient(x)
    x += learning_rate * dx
```

## SGD+Momentum

$$v_{t+1} = \rho v_t + \nabla f(x_t)$$

$$x_{t+1} = x_t - \alpha v_{t+1}$$

```
vx = 0
while True:
    dx = compute_gradient(x)
    vx = rho * vx + dx
    x += learning_rate * vx
```

- Build up "velocity" as a running mean of gradients
- Rho gives "friction"; typically rho=0.9 or 0.99

# AdaGrad

```
grad_squared = 0
while True:
    dx = compute_gradient(x)
    grad_squared += dx * dx
    x -= learning_rate * dx / (np.sqrt(grad_squared) + 1e-7)
```

Added element-wise scaling of the gradient based on the historical sum of squares in each dimension

# Adaptive Moment Estimation (ADAM)

## Adam (full form)

```
first_moment = 0
second_moment = 0
for t in range(1, num_iterations):
    dx = compute_gradient(x)
    first_moment = beta1 * first_moment  + (1 - beta1) * dx
    second_moment = beta2 * second_moment + (1 - beta2) * dx * dx
    first_unbias = first_moment / (1 - beta1 ** t)
    second_unbias = second_moment / (1 - beta2 ** t)
    x -= learning_rate * first_unbias / (np.sqrt(second_unbias) + 1e-7))
```
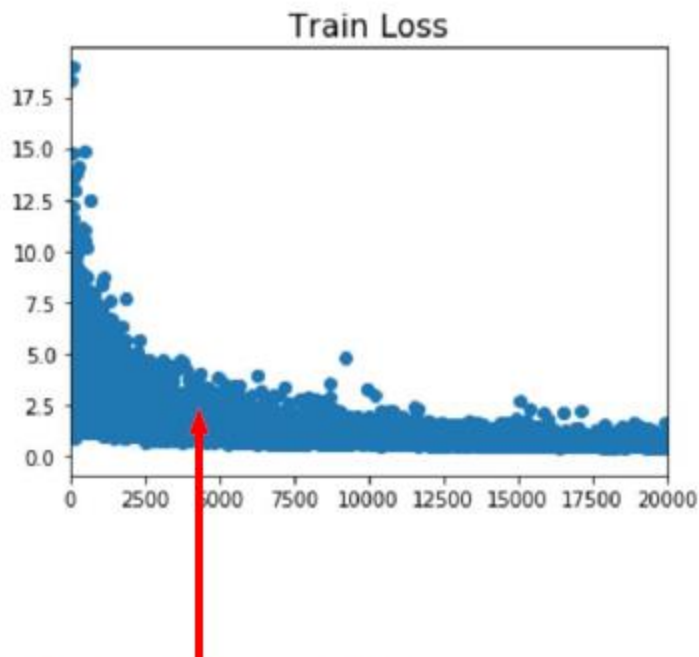
Momentum

Bias correction

AdaGrad / RMSProp

Bias correction for the fact that first and second moment estimates start at zero
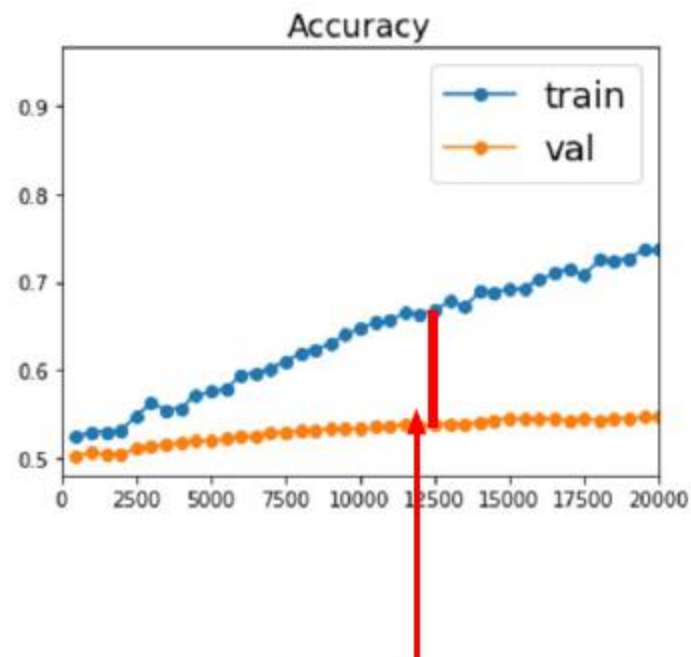
Adam with beta1 = 0.9, beta2 = 0.999, and learning_rate = 1e-3 or 5e-4 is a great starting point for many models!

Kingma and Ba, "Adam: A method for stochastic optimization", ICLR 2015

# Beyond Training Error



Train Loss

Accuracy

Better optimization algorithms
help reduce training loss

But we really care about error on new
data - how to reduce the gap?

# Regularization: Add term to loss

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

In common use:

**L2 regularization**  $R(W) = \sum_k \sum_l W_{k,l}^2$  (Weight decay)
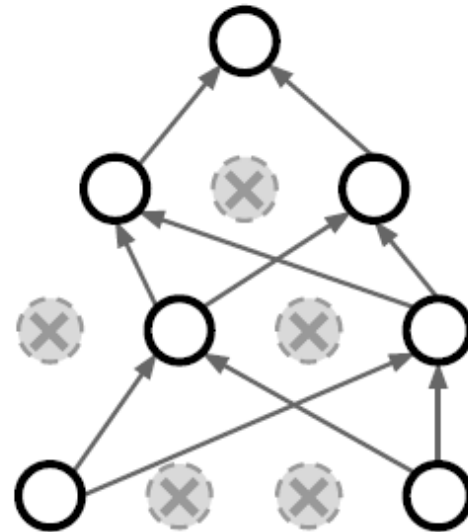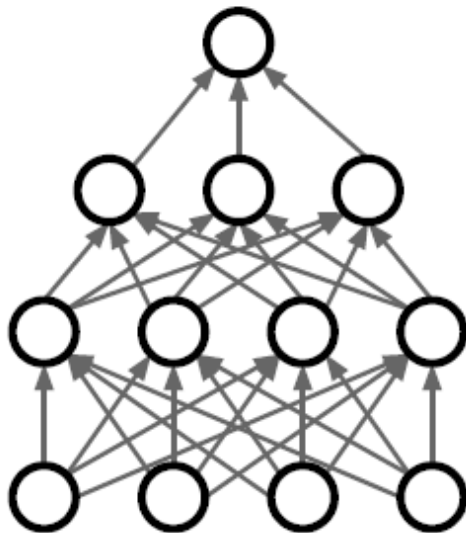
L1 regularization  $R(W) = \sum_k \sum_l |W_{k,l}|$

Elastic net (L1 + L2)  $R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$

# Dropout

## 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
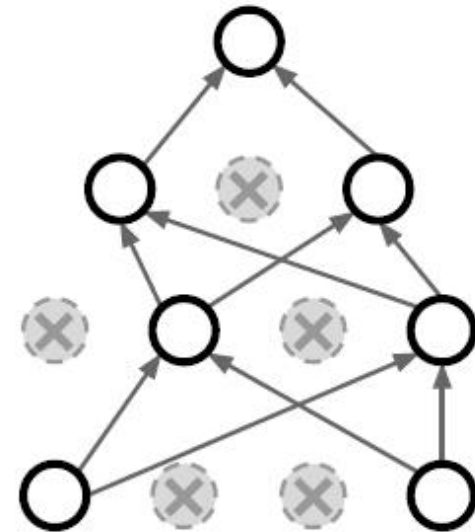
# Regularization: Dropout

Example forward pass with a 3-layer network using dropout

```python
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)
```
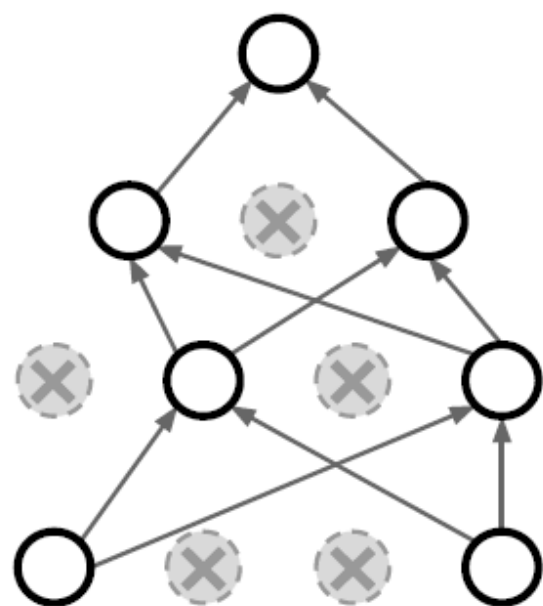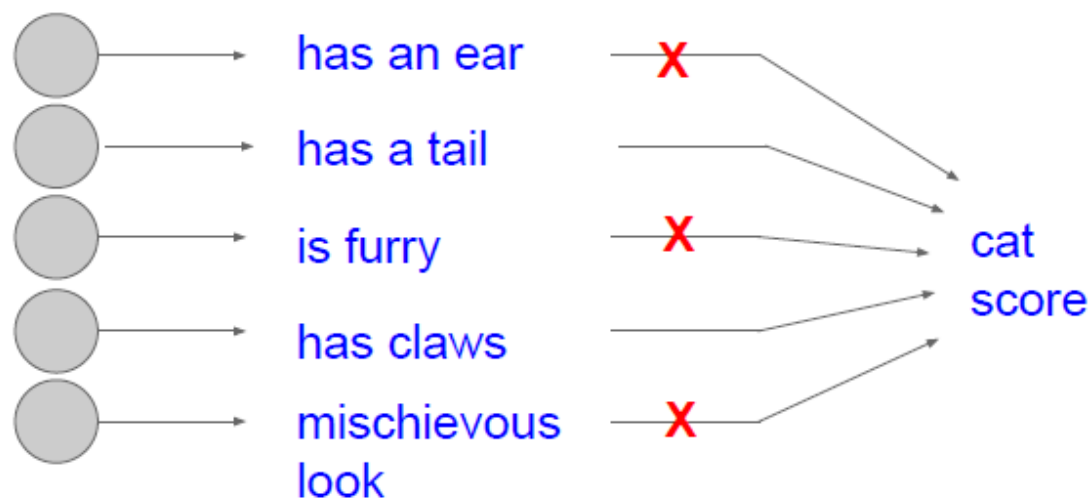
# Regularization: Dropout
## How can this possibly be a good idea?

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



has an ear ✗

has a tail

is furry ✗

has claws

mischievous look ✗

cat score

# Regularization: Dropout

How can this possibly be a good idea?



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

```
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:
<u>output at test time</u> = <u>expected output at training time</u>

# Dropout Summary

```python
""" 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**

# Regularization: Data Augmentation



Load image
and label

"cat"

Compute
loss

CNN

This image by Nikita is
licensed under CC-BY 2.0

# Regularization: Data Augmentation



Load image and label

"cat"

Transform image

CNN

Compute loss

# Data Augmentation
## Horizontal Flips

# Data Augmentation
## Random crops and scales

**Training**: sample random crops / scales

ResNet:

1. Pick random L in range [256, 480]
2. Resize training image, short side = L
3. Sample random 224 x 224 patch

# Data Augmentation
## Color Jitter

Simple: Randomize
contrast and brightness

# Data Augmentation
## Get creative for your problem!

Random mix/combinations of :
- translation
- rotation
- stretching
- shearing,
- lens distortions, …  (go crazy)

# Transfer Learning

"You need a lot of a data if you want to train/use CNNs"

# Transfer Learning with CNNs

Donahue et al, "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014
Razavian et al, "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014

## 1. Train on Imagenet

| |
|:---:|
| FC-1000 |
| FC-4096 |
| FC-4096 |
| MaxPool |
| Conv-512 |
| Conv-512 |
| MaxPool |
| Conv-512 |
| Conv-512 |
| MaxPool |
| Conv-256 |
| Conv-256 |
| MaxPool |
| Conv-128 |
| Conv-128 |
| MaxPool |
| Conv-64 |
| Conv-64 |
| Image |

# Transfer Learning with CNNs

Donahue et al, "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014
Razavian et al, "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014



17.05.2021
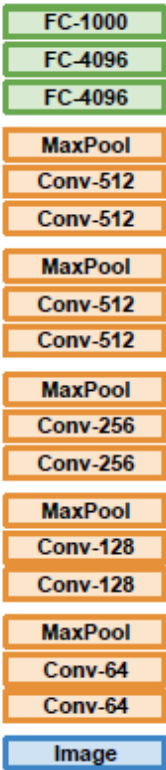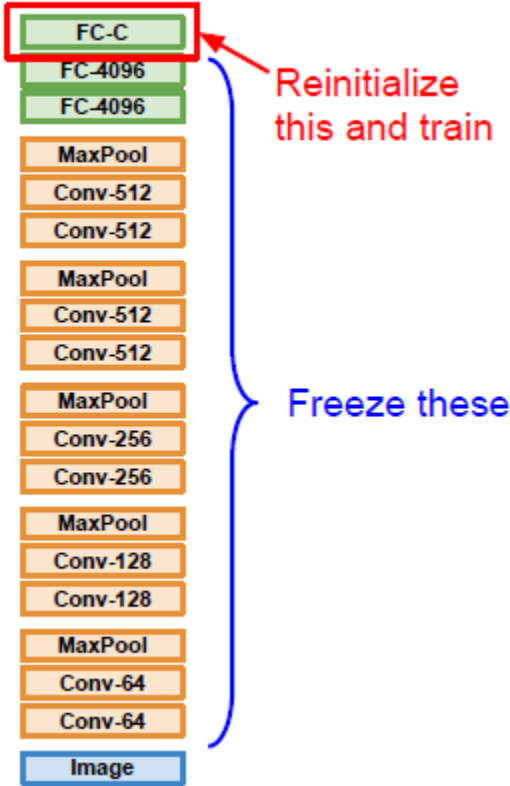
# Transfer Learning with CNNs

Donahue et al, "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014
Razavian et al, "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014
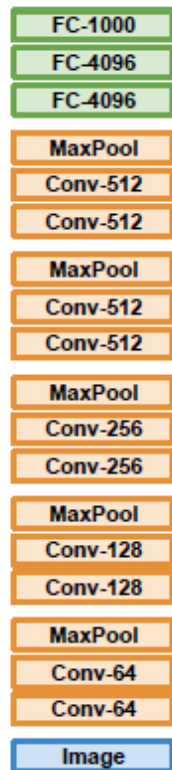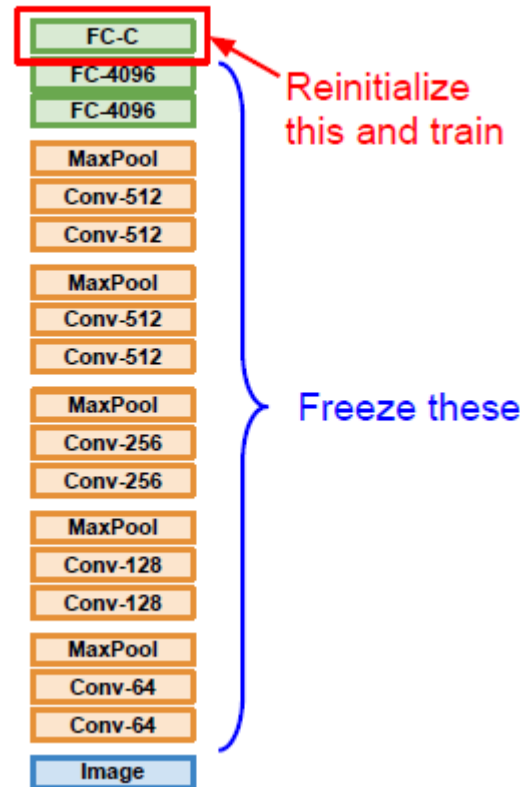


**1. Train on Imagenet**

FC-1000
FC-4096
FC-4096

MaxPool
Conv-512
Conv-512

MaxPool
Conv-512
Conv-512

MaxPool
Conv-256
Conv-256

MaxPool
Conv-128
Conv-128

MaxPool
Conv-64
Conv-64

Image

**2. Small Dataset (C classes)**

FC-C
FC-4096
FC-4096

MaxPool
Conv-512
Conv-512

MaxPool
Conv-512
Conv-512

MaxPool
Conv-256
Conv-256

MaxPool
Conv-128
Conv-128

MaxPool
Conv-64
Conv-64

Image

Reinitialize this and train

Freeze these

**3. Bigger dataset**

FC-C
FC-4096
FC-4096

MaxPool
Conv-512
Conv-512

MaxPool
Conv-512
Conv-512

MaxPool
Conv-256
Conv-256

MaxPool
Conv-128
Conv-128

MaxPool
Conv-64
Conv-64

Image

Train these

With bigger dataset, train more layers

Freeze these

Lower learning rate when finetuning; 1/10 of original LR is good starting point

# ImageNet Classification with Deep Convolutional Neural Networks
*[Krizhevsky, Sutskever, Hinton, 2012]*
"AlexNet"

## Architecture:

- CONV1
- MAX POOL1
- NORM1
- CONV2
- MAX POOL2
- NORM2
- CONV3
- CONV4
- CONV5
- MAX POOL3
- FC6
- FC7
- FC8

Input: 227x227x3 images (224x224 before padding)
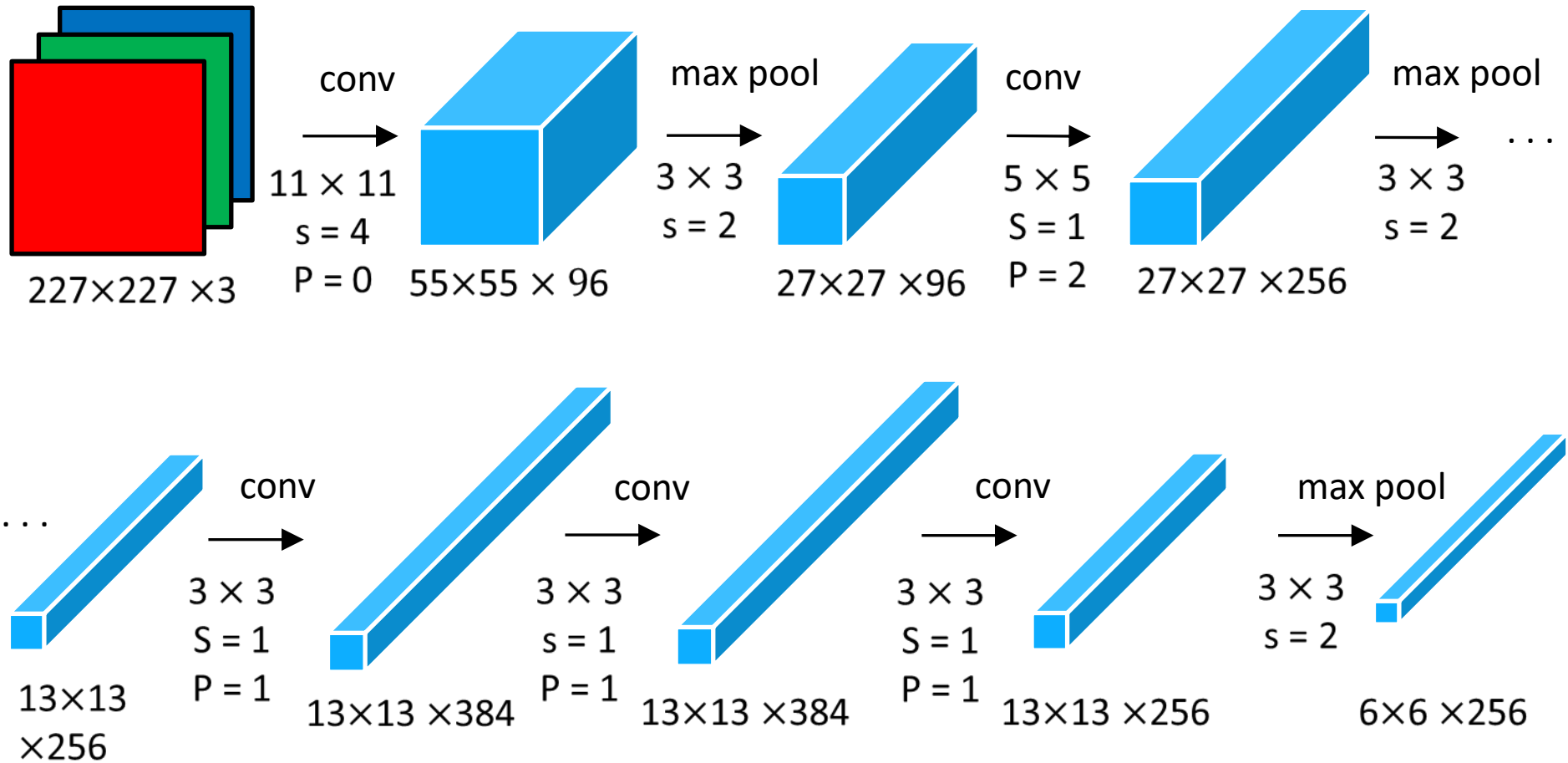
First layer: 96 11x11 filters applied at stride 4

**Output volume size?**

$$(N-F)/s+1 = (227-11)/4+1 = 55 \rightarrow [55 \times 55 \times 96]$$
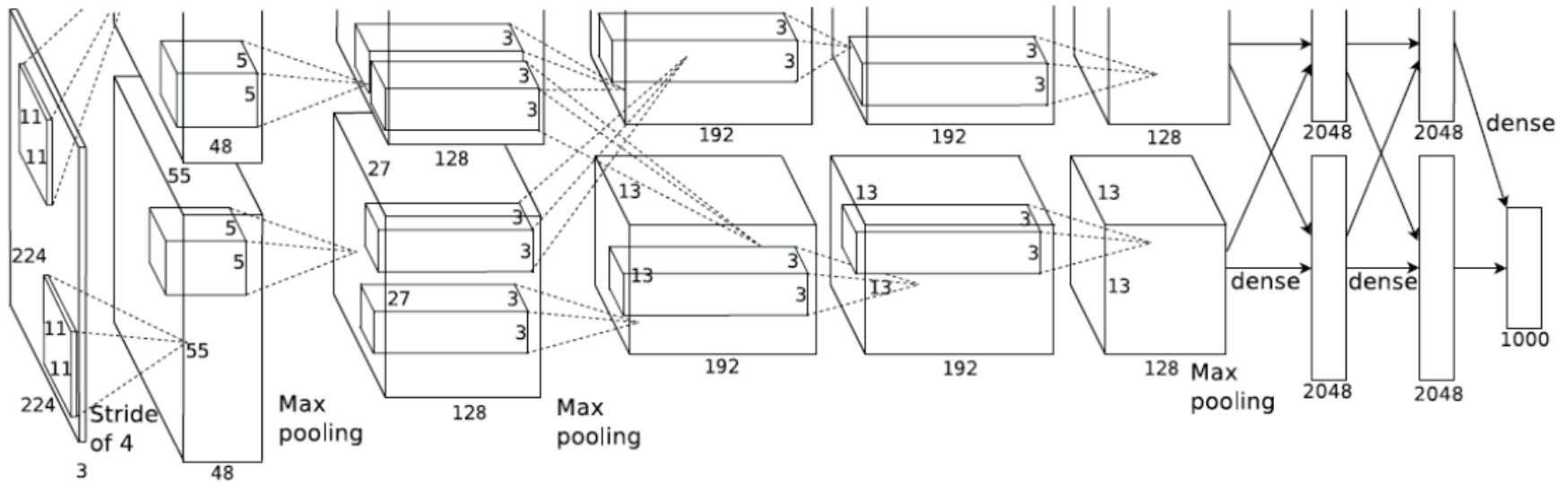
**Number of parameters in this layer?**

$$(11*11*3)*96 = 35K$$

# AlexNet



227×227 ×3    conv   $11 \times 11$   s = 4   P = 0   55×55 × 96   max pool   $3 \times 3$   s = 2   27×27 ×96   conv   $5 \times 5$   S = 1   P = 2   27×27 ×256   max pool   $3 \times 3$   s = 2   . . .

. . .   13×13 ×256   conv   $3 \times 3$   S = 1   P = 1   13×13 ×384   conv   $3 \times 3$   s = 1   P = 1   13×13 ×384   conv   $3 \times 3$   S = 1   P = 1   13×13 ×256   max pool   $3 \times 3$   s = 2   6×6 ×256

# AlexNet

➢ Deep CNN architecture proposed by **Krizhevsky** [*Krizhevsky NIPS 2012*].

  – 5 convolutional layers (with pooling and ReLU)

  – 3 fully-connected layers

  – won ImageNet Large Scale Visual recognition Challenge 2012

  – top-1 validation error rate of 40.7%

# AlexNet