# Section 7: Exam Review

CSE 493G1, Spring 2025

Materials prepared by Vishnu Iyengar

# Course Logistics

- EXAM next Tuesday [5/20]
    - **Show up to class!!**
    - Covers content from Lectures 2-12, Assignments 1-3, Sections 1-6
    - Allowed one double sided note sheet on a standard 8.5'x11' paper
- Project Milestone due next Friday [5/23]
- A4 due next Sunday [5/25]

# Exam Review

- **Deep Learning Foundations**
- Optimization and Training Techniques
- Neural Networks and Backpropogation
- Convolutional Neural Networks (CNNs)
- Sequence Models and Interpretability

# Deep Learning Foundations: Image Classification

- 1) Collect a dataset of images and labels
- 2) Use machine learning algorithms to train a classifier
- 3) Evaluate the classifier on new images

```
def train(images, labels):
    # Machine learning!
    return model
```

```
def predict(model, test_images):
    # Use model to predict labels
    return test_labels
```
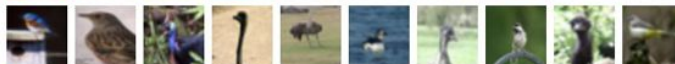
**Example training set**



airplane
automobile
bird
cat
deer

# Deep Learning Foundations: kNN

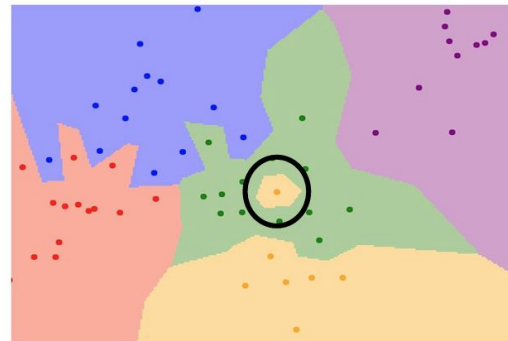```python
import numpy as np

class NearestNeighbor:
  def __init__(self):
    pass

  def train(self, X, y):
    """ X is N x D where each row is an example. Y is 1-dimension of size N """
    # the nearest neighbor classifier simply remembers all the training data
    self.Xtr = X
    self.ytr = y

  def predict(self, X):
    """ X is N x D where each row is an example we wish to predict label for """
    num_test = X.shape[0]
    # lets make sure that the output type matches the input type
    Ypred = np.zeros(num_test, dtype = self.ytr.dtype)

    # loop over all test rows
    for i in xrange(num_test):
      # find the nearest training image to the i'th test image
      # using the L1 distance (sum of absolute value differences)
      distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
      min_index = np.argmin(distances) # get the index with smallest distance
      Ypred[i] = self.ytr[min_index] # predict the label of the nearest example

    return Ypred
```
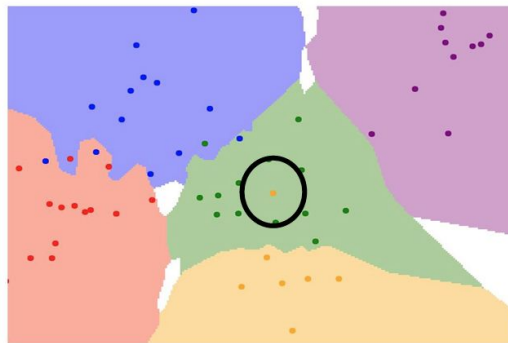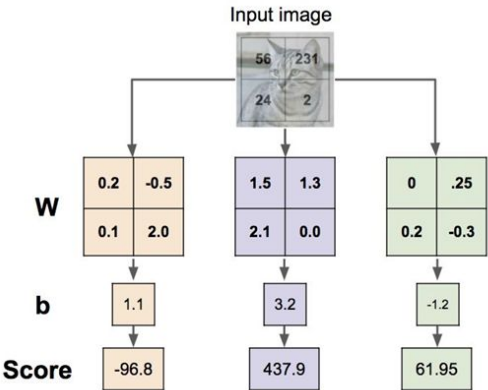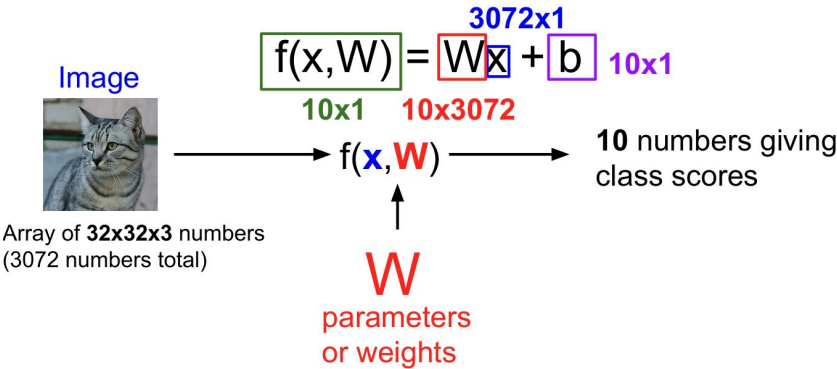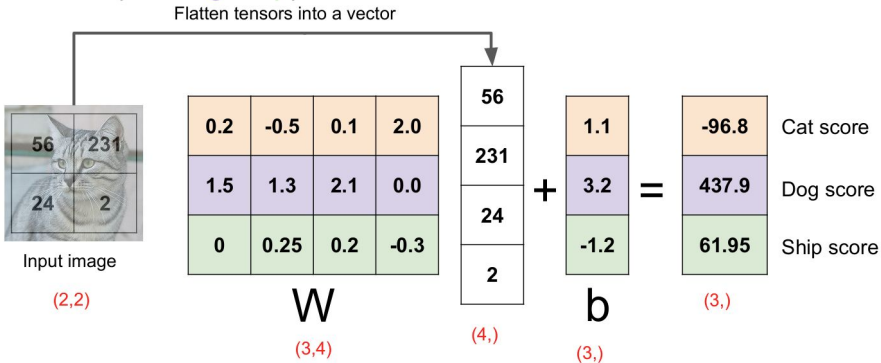


K = 1



K = 3

# Deep Learning Foundations: Linear Classifier
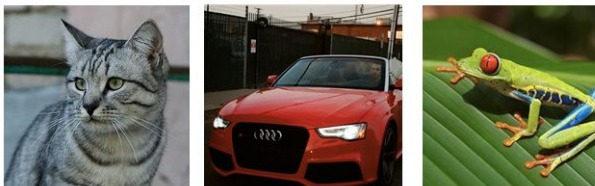
Parametric Approach: Linear Classifier

**Image**



Array of **32x32x3** numbers
(3072 numbers total)

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

10x1    3072x1    10x1

10x1    10x3072

f(**x**,**W**) → **10** numbers giving class scores

**W**
parameters or weights

**Algebraic viewpoint:** Example with an image with 4 pixels, and 3 classes (cat/dog/ship)

Flatten tensors into a vector



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

Input image
(2,2)

W
(3,4)

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

(4,)

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

b
(3,)

| -96.8 | Cat score |
|-------|-----------|
| 437.9 | Dog score |
| 61.95 | Ship score |

(3,)

**Input image**



**W**

| 0.2 | -0.5 | | 1.5 | 1.3 | | 0 | .25 |
|-----|------|--|-----|-----|--|---|-----|
| 0.1 | 2.0 | | 2.1 | 0.0 | | 0.2 | -0.3 |

**b**

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

**Score**

| -96.8 | | 437.9 | | 61.95 |
|-------|--|-------|--|-------|


cat     deer     dog

# Deep Learning Foundations: SVM Loss

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

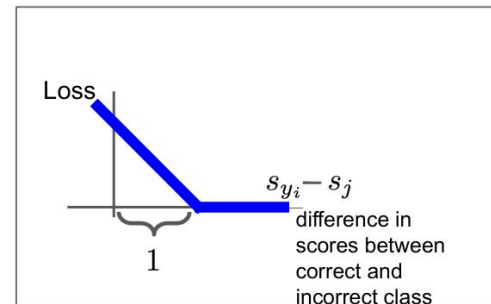|  | | | |
|---|---|---|---|
| cat | **3.2** | 1.3 | 2.2 |
| car | 5.1 | **4.9** | 2.5 |
| frog | -1.7 | 2.0 | **-3.1** |
| Losses: | 2.9 | 0 | |

**Multiclass SVM loss:**

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

**Interpreting Multiclass SVM loss:**

Loss

$s_{y_i} - s_j$

difference in
scores between
correct and
incorrect class

1

# Deep Learning Foundations: Softmax Loss



**Softmax Classifier** (Multinomial Logistic Regression)

Want to interpret raw classifier scores as **probabilities**

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

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$ Softmax Function

Probabilities must be >= 0

Probabilities must sum to 1

$$L_i = -\log P(Y = y_i | X = x_i)$$

| | logits | | unnormalized probabilities | | probabilities | |
|---|---|---|---|---|---|---|
| cat | **3.2** | | **24.5** | | **0.13** | |
| car | 5.1 | exp | 164.0 | normalize | 0.87 | |
| frog | -1.7 | | 0.18 | | 0.00 | |

$L_i$ = -log(0.13) = **2.04**

Unnormalized log-probabilities / logits

unnormalized probabilities

probabilities

**Maximum Likelihood Estimation**
Choose weights to maximize the likelihood of the observed data

# Deep Learning Foundations: Regularization

## Regularization

$\lambda$ = regularization strength (hyperparameter)

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

**Data loss**: Model predictions should match training data

**Regularization**: Prevent the model from doing *too* well on training data

**Simple examples**

L2 regularization: $R(W) = \sum_k \sum_l W_{k,l}^2$

L1 regularization: $R(W) = \sum_k \sum_l |W_{k,l}|$

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

# Deep Learning Foundations: Summary

## Recap

How do we find the best W?

- We have some dataset of (x,y)
- We have a **score function:** $s = f(x; W) \overset{\text{e.g.}}{=} Wx$
- We have a **loss function:**

Softmax

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

SVM

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

$$L = \frac{1}{N} \sum_{i=1}^{N} L_i + R(W)$$ Full loss

# Exam Review

- Deep Learning Foundations
- **Optimization and Training Techniques**
- Neural Networks and Backpropogation
- Convolutional Neural Networks (CNNs)
- Sequence Models and Interpretability

# Optimization & Training Techniques: Gradient Desc



$$L = \frac{1}{N} \sum_{i=1}^{N} L_i + \sum_k W_k^2$$

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

$$s = f(x; W) = Wx$$

want $\nabla_W L$

```
# Vanilla Gradient Descent

while True:
  weights_grad = evaluate_gradient(loss_fun, data, weights)
  weights += - step_size * weights_grad # perform parameter update
```

# Optimization & Training Techniques: SGD

## Stochastic Gradient Descent (SGD)

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

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

Full sum expensive when N is large!

Approximate sum using a **minibatch** of examples
32 / 64 / 128 common

```
# Vanilla Minibatch Gradient Descent

while True:
    data_batch = sample_training_data(data, 256) # sample 256 examples
    weights_grad = evaluate_gradient(loss_fun, data_batch, weights)
    weights += - step_size * weights_grad # perform parameter update
```

# Optimizers & Training Techniques: Optimizers

## 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
```

Velocity

actual step

Gradient

Combine gradient at current point with
velocity to get step used to update weights

# Nesterov Momentum

That's it!

Step 1: Calculate the velocity at *t+1*
Step 2: Update the parameters using the velocities at *t+1* and *t*

$$v_{t+1} = \rho v_t - \alpha \nabla f(\tilde{x}_t)$$

$$\tilde{x}_{t+1} = \tilde{x}_t - \rho v_t + (1 + \rho)v_{t+1}$$

$$= \tilde{x}_t + v_{t+1} + \rho(v_{t+1} - v_t)$$

Velocity

Gradient

actual step

"Look ahead" to the point where updating using velocity would take us; compute gradient there and mix it with velocity to get actual update direction

# **Optimizers & Training Techniques: Optimizers**

**AdaGrad:**

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

**RMSProp:**

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

## 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

# Optimizers & Training Techniques: LR Schedules

Phases of learning...

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

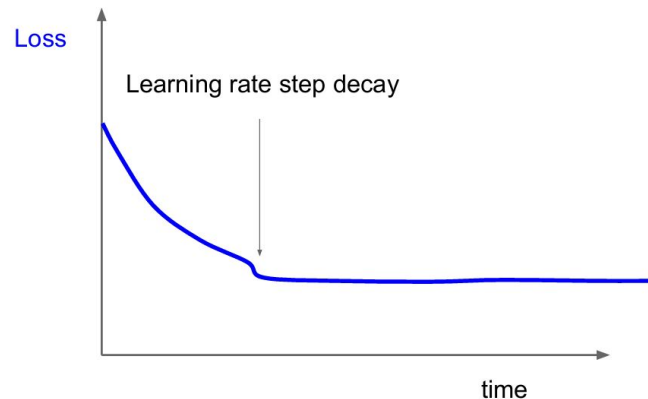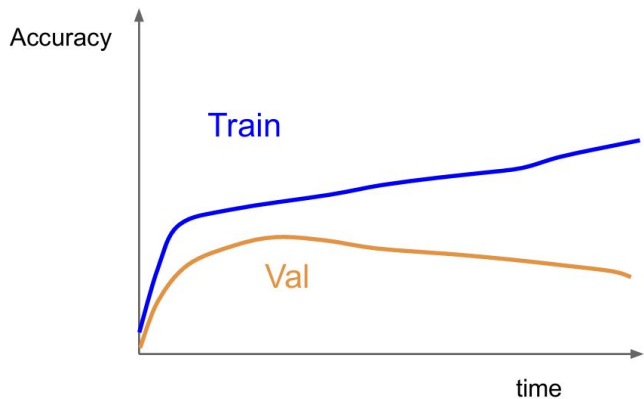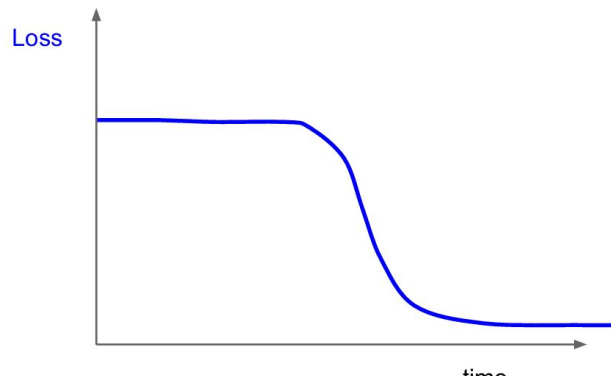**Cosine:** $\alpha_t = \frac{1}{2}\alpha_0\left(1 + \cos(t\pi/T)\right)$
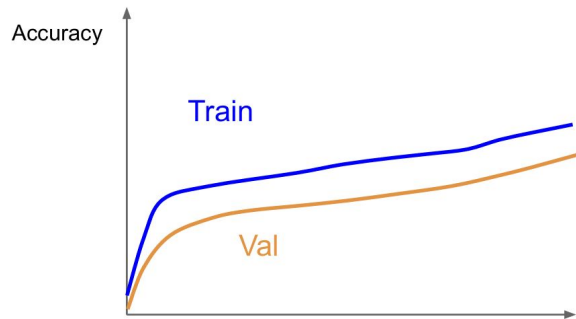
**Linear:** $\alpha_t = \alpha_0(1 - t/T)$

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

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

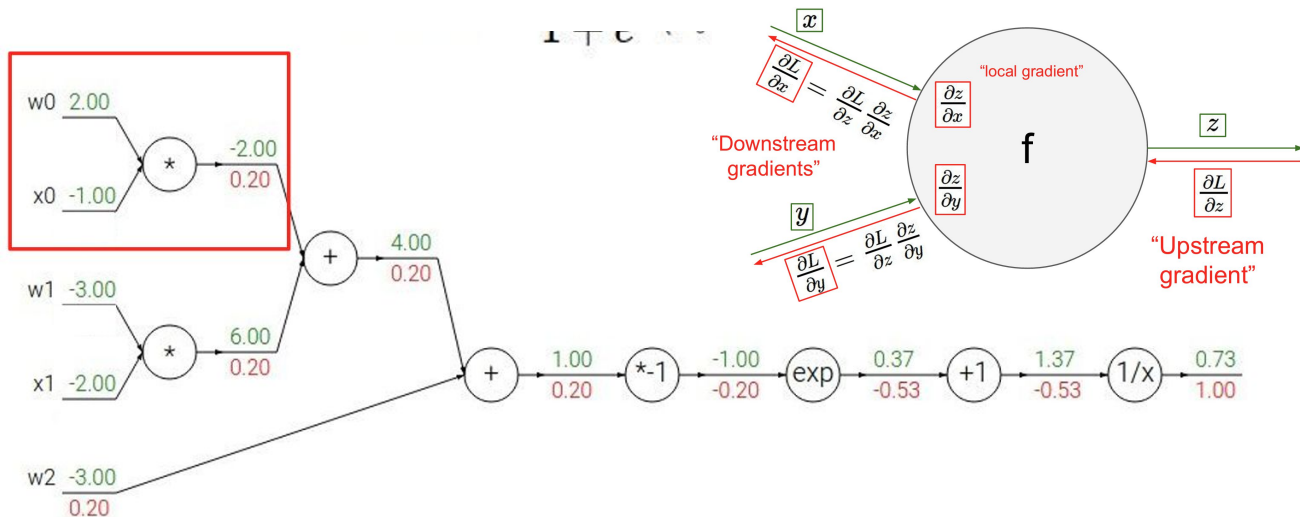# Optimizers & Training techniques: Choosing hyperparamaters



**What should you do?**

# Exam Review

- Deep Learning Foundations
- Optimization and Training Techniques
- **Neural Networks and Backpropogation**
- Convolutional Neural Networks (CNNs)
- Sequence Models and Interpretability

# Neural Networks and Backpropogation



$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x \qquad \bigg| \qquad f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a \qquad \bigg| \qquad f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

# Neural Networks & Backpropogation: Activ. Fns.

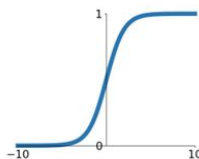Are neurons saturated?

Are outputs zero-centered?

Is it computationally efficient?

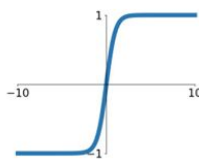Does it "kill" gradients?

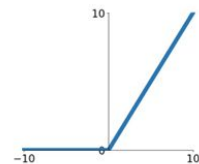## Activation functions

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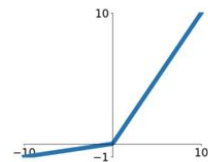$\sigma(x) = \frac{1}{1+e^{-x}}$
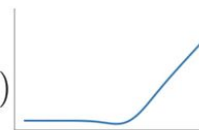
**tanh**
$\tanh(x)$

**ReLU**
$\max(0, x)$

**Leaky ReLU**
$\max(0.1x, x)$

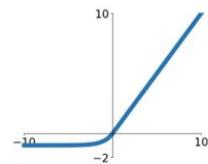**GeLU**
$0.5x\left(1 + \tanh\left[\sqrt{2/\pi}(x + 0.044715x^3)\right]\right)$
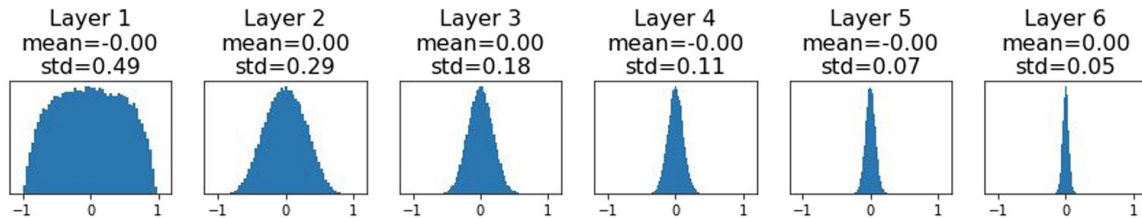
**ELU**
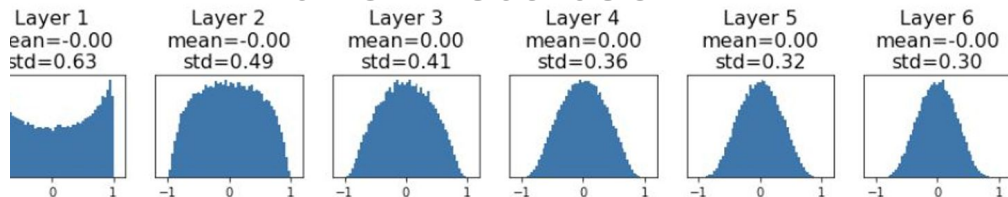$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$

# Neural Networks & Backprop: Weight Initialization



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

### Xavier Initialization



W = np.random.randn(Din, Dout) / np.sqrt(Din)

**Let:** $y = x_1 w_1 + x_2 w_2 + \ldots + x_{Din} w_{Din}$

**Assume:** $Var(x_1) = Var(x_2) = \ldots = Var(x_{Din})$

**We want:** $Var(y) = Var(x_i)$

$Var(y) = Var(x_1 w_1 + x_2 w_2 + \ldots + x_{Din} w_{Din})$
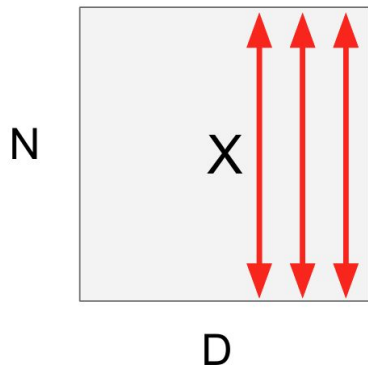$= Din\ Var(x_i w_i)$
$= Din\ Var(x_i)\ Var(w_i)$

[Assume all $x_i$, $w_i$ are iid]

So, $Var(y) = Var(x_i)$ only when $Var(w_i) = 1/Din$

# Neural Networks & Backprop: Normalizations

## Batch Normalization

[Ioffe and Szegedy, 2015]

**Input:** $x : N \times D$



N

X

D

$$\mu_j = \frac{1}{N} \sum_{i=1}^{N} x_{i,j}$$

Per-channel mean, shape is D
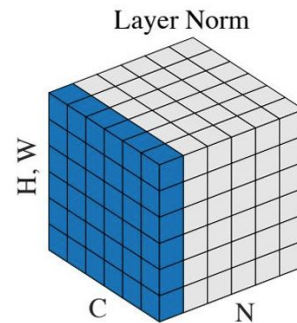
$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^{N} (x_{i,j} - \mu_j)^2$$

Per-channel var, shape is D

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

Normalized x, Shape is N x D



Batch Norm

H, W

C

N

Layer Norm

H, W

C

N

# Neural Networks & Backprop: Dropout



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

# Exam Review

- Deep Learning Foundations
- Optimization and Training Techniques
- Neural Networks and Backpropogation
- **Convolutional Neural Networks (CNNs)**
- Sequence Models and Interpretability

# CNNs



32x32x3 image
5x5x3 filter
32
32
3

activation map

28
28
1

convolve (slide) over all spatial locations

## Convolution layer: summary

Let's assume input is $W_1$ x $H_1$ x C
Conv layer needs 4 hyperparameters:
- Number of filters **K**
- The filter size **F**
- The stride **S**
- The zero padding **P**

This will produce an output of $W_2$ x $H_2$ x K where:
- $W_2 = (W_1 - F + 2P)/S + 1$
- $H_2 = (H_1 - F + 2P)/S + 1$

Number of parameters: $F^2KC$ and K biases

# Pooling

## Pooling layer

- makes the representations smaller and more manageable
- operates over each activation map independently:

Single depth slice



224x224x64 → pool → 112x112x64



224 — downsampling → 112

## Pooling layer: summary

Let's assume input is $W_1 \times H_1 \times C$
Conv layer needs 2 hyperparameters:

- The spatial extent **F**
- The stride **S**

This will produce an output of $W_2 \times H_2 \times C$ where:

- $W_2 = (W_1 - F)/S + 1$
- $H_2 = (H_1 - F)/S + 1$

Number of parameters: 0

# CNNs/Pooling

# Exam Review

- Deep Learning Foundations
- Optimization and Training Techniques
- Neural Networks and Backpropogation
- Convolutional Neural Networks (CNNs)
- **Sequence Models and Interpretability**

# RNNs



$$h_t = f_W(h_{t-1}, x_t)$$

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$y_t = W_{hy}h_t$$

i.e.
Character-Level
Language
Model

# RNNs (Cont.)



## Vanilla RNN Gradient Flow

Gradients over multiple time steps:

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013

$$\frac{\partial L}{\partial W} = \sum_{t=1}^{T} \frac{\partial L_t}{\partial W} \qquad \frac{\partial h_t}{\partial h_{t-1}} = tanh'(W_{hh} h_{t-1} + W_{xh} x_t) W_{hh}$$
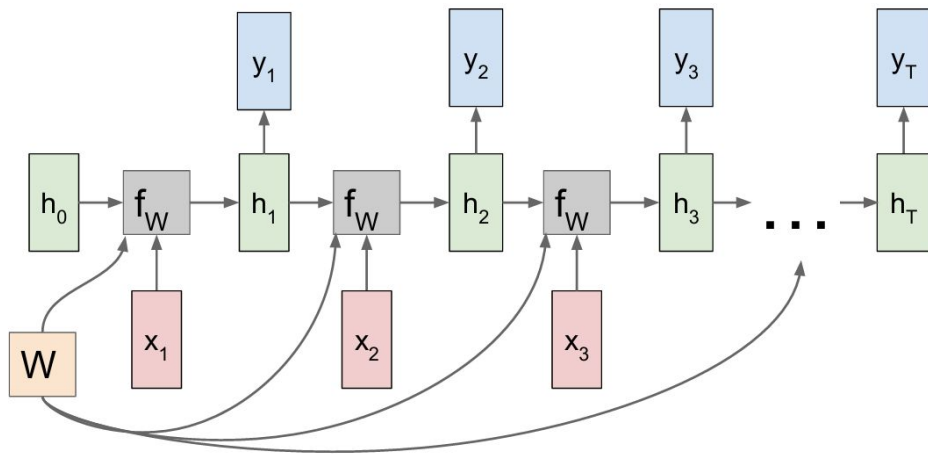
$$\frac{\partial L_T}{\partial W} = \frac{\partial L_T}{\partial h_T} \frac{\partial h_t}{\partial h_{t-1}} \cdots \frac{\partial h_1}{\partial W} = \frac{\partial L_T}{\partial h_T} \left( \prod_{t=2}^{T} \frac{\partial h_t}{\partial h_{t-1}} \right) \frac{\partial h_1}{\partial W}$$

# LSTMs

**i**: <u>Input gate</u>, whether to write to cell
**f**: Forget gate, Whether to erase cell
**o**: <u>Output gate</u>, How much to reveal cell
**g**: Info gate, How much to write to cell



$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

Solves the vanishing gradient
problem for the cell memory blocks!

# Attention and Transformers

## Image Captioning with RNNs & <span style="color:red">Attention</span>

Compute alignments scores (scalars):

$$e_{t,i,j} = f_{att}(h_{t-1}, z_{i,j})$$

$f_{att}(.)$ is an MLP

Alignment scores:
H x W

| $e_{1,0,0}$ | $e_{1,0,1}$ | $e_{1,0,2}$ |
| $e_{1,1,0}$ | $e_{1,1,1}$ | $e_{1,1,2}$ |
| $e_{1,2,0}$ | $e_{1,2,1}$ | $e_{1,2,2}$ |

Attention:
H x W

| $a_{1,0,0}$ | $a_{1,0,1}$ | $a_{1,0,2}$ |
| $a_{1,1,0}$ | $a_{1,1,1}$ | $a_{1,1,2}$ |
| $a_{1,2,0}$ | $a_{1,2,1}$ | $a_{1,2,2}$ |

Normalize to get attention weights:

$$a_{t,:,:} = softmax(e_{t,:,:})$$

$0 < a_{t,i,j} < 1$,
attention values sum to 1

Compute context vector:

$$c_t = \sum_{i,j} a_{t,i,j}\, z_{t,i,j}$$

CNN

| $z_{0,0}$ | $z_{0,1}$ | $z_{0,2}$ |
| $z_{1,0}$ | $z_{1,1}$ | $z_{1,2}$ |
| $z_{2,0}$ | $z_{2,1}$ | $z_{2,2}$ |

$h_0$

$c_1$

Extract spatial features from a pretrained CNN

Features:
H x W x D

X

Xu et al. "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention". ICML 2015

# Attention and Transformers



## Image Captioning with RNNs & Attention

Alignment scores: H x W

Attention: H x W

$e_{t,i,j} = f_{att}(h_{t-1}, z_{i,j})$

$a_{t,:,:} = softmax(e_{t,:,:})$

$c_t = \sum_{i,j} a_{t,i,j} z_{t,i,j}$

**Decoder**: $y_t = g_v(y_{t-1}, h_{t-1}, c_t)$
New context vector at every time step

Extract spatial features from a pretrained CNN

CNN

Features: H x W x D

Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

# Attention and Transformers
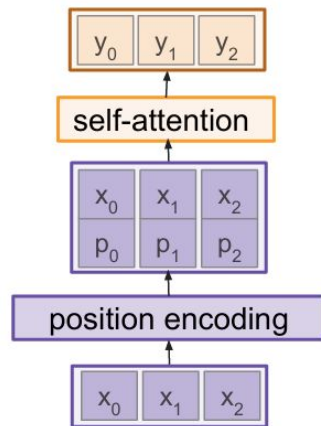
## Self attention layer



**Outputs:**
context vectors: $\mathbf{y}$ (shape: $D_v$)

**Operations:**
Key vectors: $\mathbf{k} = \mathbf{x}\mathbf{W}_k$
Value vectors: $\mathbf{v} = \mathbf{x}\mathbf{W}_v$
Query vectors: $\mathbf{q} = \mathbf{x}\mathbf{W}_q$
Alignment: $e_{i,j} = q_j \cdot k_i / \sqrt{D}$
Attention: $\mathbf{a} = \text{softmax}(\mathbf{e})$
Output: $y_j = \sum_i a_{i,j} v_i$

**Inputs:**
Input vectors: $\mathbf{x}$ (shape: N x D)

## Positional encoding



Concatenate special positional encoding $p_j$ to each input vector $x_j$

We use a function $pos: \text{N} \rightarrow \text{R}^d$
to process the position $j$ of the vector into a d-dimensional vector

So, $p_i = pos(j)$

**Good Luck!**

**Practice Exam will be posted later today**