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

Lecture 4: Regularization and Optimization
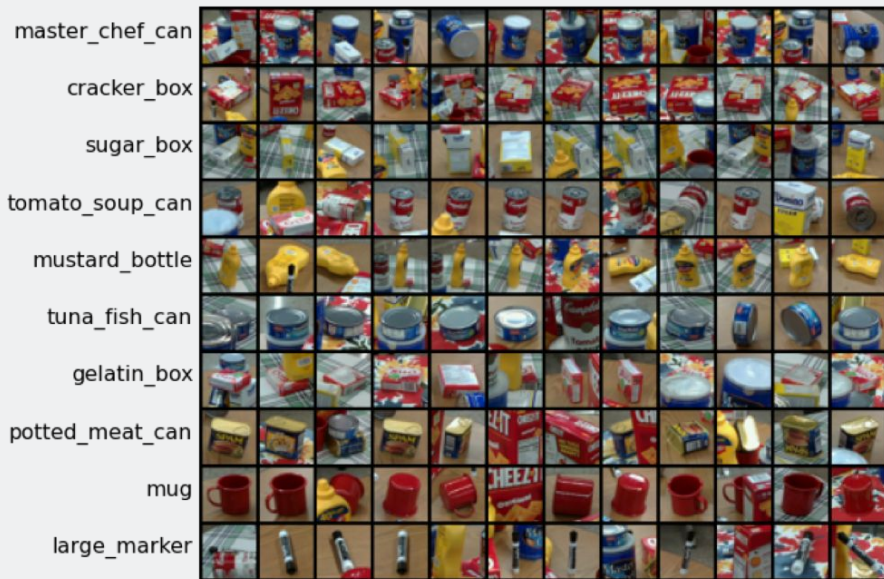
https://deeprob.org/w25/

ROBOTICS

# Today

- Feedback and Recap (5min)
- Regularization and Optimization
  - Regularization (15min)
  - Optimization (20min)
  - Computing Gradients (30min)
- Summary and Takeaways (5min)

# Project 1 - Dataset

## **P**rogress **R**obot **O**bject **P**erception **S**amples **D**ataset
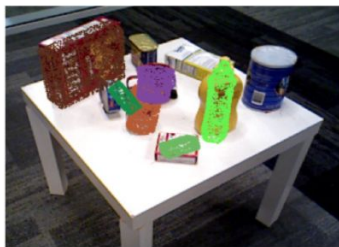


Chen et al., "ProgressLabeller: Visual Data Stream Annotation for Training Object-Centric 3D Perception", IROS, 2022.

**10 classes**
**32x32** RGB images
**50k** training images (5k per class)
**10k** test images (1k per class)

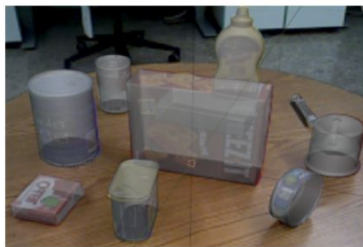# Project 1 - How was this dataset created?



ProgressLabeller: Visual Data Stream Annotation for Training Object-Centric 3D Perception

Xiaotong Chen    Huijie Zhang    Zeren Yu    Stanley Lewis    Odest Chadwicke Jenkins

**Rough Pose Estimates from Pretrained Model** → **6D pose annotation through interactive interface** → **Fine-tuned Pose Estimates** → **Pose-based Robot Grasping**
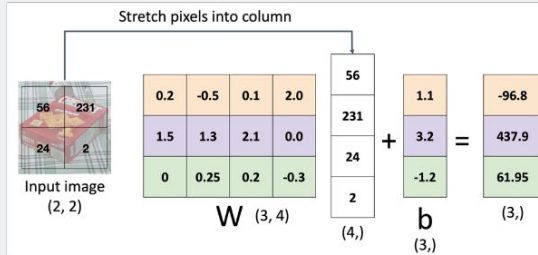
Human Annotator

## Idea:

1. **Record video of scene**

2. **Human labels object pose in selected frames**

3. **Pose labels propagate to (large number of) remaining frames**

# Recap: Linear Classifier - Three Viewpoints
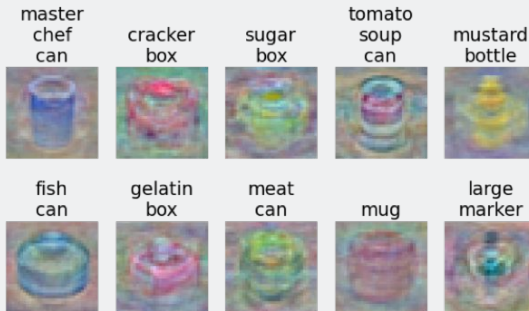
① **Algebraic Viewpoint**

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



② **Visual Viewpoint**

One template per class



③ **Geometric Viewpoint**

Hyperplanes cutting up space

# Recap: Loss Functions

- We have some dataset of (x, y)
- We have a **score function:**
- We have a **loss function**:

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

Linear classifier

Softmax: $L_i = -\log\left(\dfrac{\exp(s_{y_i})}{\sum_j \exp(s_j)}\right)$

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



**M | ROBOTICS**

# Discussion on Last week's Quizzes

(refer to Canvas)
- If you have questions, please come ask!

# How to find the best W and b?

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

Linear classifier

Problem: Loss functions encourage good performance on training data but we care about <u>test</u> data

# Regularization

# Overfitting

A model is **overfit** when it performs too well on the training data, and has <u>poor performance</u> for unseen data

Example: Linear classifier with 1D inputs, 2 classes, and softmax loss

$$s_i = w_i x + b_i$$

$$p_i = \frac{exp(s_i)}{exp(s_1) + exp(s_2)}$$

$$L = -log(p_y)$$

# Overfitting

A model is **overfit** when it performs too well on the training data, and has <u>poor performance</u> for unseen data



**Both models have perfect accuracy on the training data!**
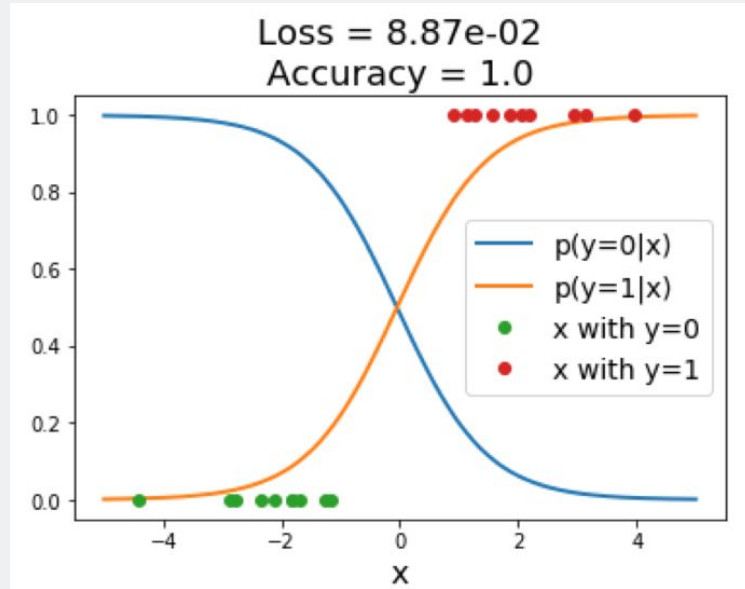
Low loss, but unnatural "cliff" between the training points
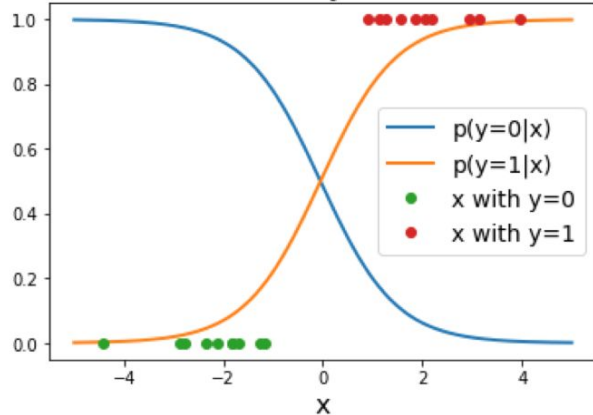
# Overfitting

A model is **overfit** when it performs too well on the training data, and has <u>poor performance</u> for unseen data



Overconfidence in regions with no training data could give **poor generalization**

# Regularization: Beyond Training Error

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

Data loss: Model predictions should match training data

# Regularization: Beyond Training Error

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

Data loss: Model predictions should match training data

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

# Regularization: Beyond Training Error

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

Hyperparameter giving regularization strength

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

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

**Simple examples:**

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

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

More complex:
Dropout
Batch normalization
Cutout, Mixup, Stochastic depth, etc...

**M | ROBOTICS**

# Regularization: Prefer Simpler Models

Example: Linear classifier with 1D inputs, 2 classes, and softmax loss

$$s_i = w_i x + b_i$$

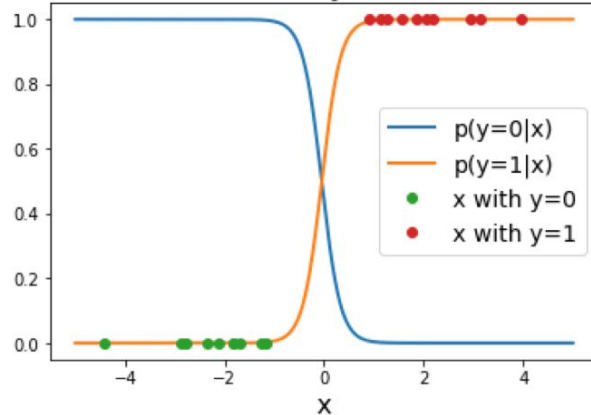$$p_i = \frac{exp(s_i)}{exp(s_1) + exp(s_2)}$$

$$L = -log(p_y) \; +\lambda \sum_i w_i^2$$

Regularization term causes loss to **increase** for model with sharp cliff



ROBOTICS

# Aha Slides
# (In-class participation)

https://ahaslides.com/WJTNO

# Regularization: Expressing Preferences

$x = [1,1,1,1]$

$w_1 = [1,0,0,0]$

$w_2 = [0.25,0.25,0.25,0.25]$

L2 Regularization

$$R(W) = \sum_{k,l} W_{k,l}^2$$

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

Q1: Which weight would the data loss prefer?

Q2: Which weight would the L2 regularization prefer?

Hint: what does it mean by "prefer"?   Higher? Lower?

ROBOTICS

# Optimization

# Finding a good W

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

**Loss function** consists of **data loss** to fit the training data and **regularization** to prevent overfitting

ROBOTICS

# Optimization

$$w^* = \arg\min_w L(w)$$

# Optimization

$$w^* = \arg\min_w L(w)$$



The valley image and the walking man image are in CC0 1.0 public domain

ROBOTICS

# Idea 1: Random Search (bad idea!)

```python
# assume X_train is the data where each column is an example (e.g. 3073 x 50,000)
# assume Y_train are the labels (e.g. 1D array of 50,000)
# assume the function L evaluates the loss function

bestloss = float("inf") # Python assigns the highest possible float value
for num in xrange(1000):
  W = np.random.randn(10, 3073) * 0.0001 # generate random parameters
  loss = L(X_train, Y_train, W) # get the loss over the entire training set
  if loss < bestloss: # keep track of the best solution
    bestloss = loss
    bestW = W
  print 'in attempt %d the loss was %f, best %f' % (num, loss, bestloss)

# prints:
# in attempt 0 the loss was 9.401632, best 9.401632
# in attempt 1 the loss was 8.959668, best 8.959668
# in attempt 2 the loss was 9.044034, best 8.959668
# in attempt 3 the loss was 9.278948, best 8.959668
# in attempt 4 the loss was 8.857370, best 8.857370
# in attempt 5 the loss was 8.943151, best 8.857370
# in attempt 6 the loss was 8.605604, best 8.605604
# ... (truncated: continues for 1000 lines)
```

ROBOTICS

# Idea 1: Random Search (bad idea!)

```python
# Assume X_test is [3073 x 10000], Y_test [10000 x 1]
scores = Wbest.dot(Xte_cols) # 10 x 10000, the class scores for all test examples
# find the index with max score in each column (the predicted class)
Yte_predict = np.argmax(scores, axis = 0)
# and calculate accuracy (fraction of predictions that are correct)
np.mean(Yte_predict == Yte)
# returns 0.1555
```

~15.5 % accuracy on CIFAR-10

# Idea 2: Follow the slope

$$w^* = \arg\min_w L(w)$$



The valley image and the walking man image are in CC0 1.0 public domain

# Idea 2: Follow the slope

$$w^* = \arg\min_w L(w)$$

In 1-dimension, the **derivative** of a function gives the slope:

$$\frac{df}{dx} = \lim_{h \to 0} \frac{f(x+h) - f(x)}{h}$$

In multiple dimensions, the **gradient** is the vector of (partial derivatives) along each dimension

The slope in any direction is the **dot product** of the direction with the gradient. The direction of steepest descent is the **negative gradient**.

**M** | ROBOTICS

# (example)

Current W:

[0.34,

-1.11,

0.78,

0.12,

0.55,

2.81,

-3.1,

-1.5,

0.33, …]

loss 1.25347

Gradient $\frac{dL}{dW}$

[?,

?,

?,

?,

?,

?,

?,

?,

?, …]

# (example)

$$\frac{df}{dx} = \lim_{h \to 0} \frac{f(x+h) - f(x)}{h}$$

| Current **W**: | **W + h** (first dim): | Gradient $\frac{dL}{dW}$ |
|---|---|---|
| [0.34, | [0.34 + **0.0001**, | [?,  ← **???** |
| -1.11, | -1.11, | ?, |
| 0.78, | 0.78, | ?, |
| 0.12, | 0.12, | ?, |
| 0.55, | 0.55, | ?, |
| 2.81, | 2.81, | ?, |
| -3.1, | -3.1, | ?, |
| -1.5, | -1.5, | ?, |
| 0.33, …] | 0.33, …] | ?, …] |
| loss 1.25347 | loss 1.25322 | |

|ROBOTICS

# Aha Slides
# (In-class participation)

https://ahaslides.com/WJTNO

Q3

# (example)

| Current **W**: | **W + h** (second dim): | Gradient $\frac{dL}{dW}$ |
|---|---|---|
| [0.34, | [0.34, | |
| -1.11, | -1.11 + **0.0001**, | ?, |
| 0.78, | 0.78, | ?, |
| 0.12, | 0.12, | ?, |
| 0.55, | 0.55, | ?, |
| 2.81, | 2.81, | ?, |
| -3.1, | -3.1, | ?, |
| -1.5, | -1.5, | ?, |
| 0.33, …] | 0.33, …] | ?, …] |
| loss 1.25347 | loss 1.25353 | |

**ROBOTICS**

# (example)

| Current **W**: | **W + h** (second dim): | Gradient $\dfrac{dL}{dW}$ |
|---|---|---|
| [0.34,<br>-1.11,<br>0.78,<br>0.12,<br>0.55,<br>2.81,<br>-3.1,<br>-1.5,<br>0.33, …]<br><br><span style="color:red">loss 1.25347</span> | [0.34,<br>-1.11 + **0.0001**,<br>0.78,<br>0.12,<br>0.55,<br>2.81,<br>-3.1,<br>-1.5,<br>0.33, …]<br><br><span style="color:red">loss 1.2535</span><span style="color:blue">3</span> | **0.6**,<br>?,<br>?,<br><br>(1.25353 - 1.25347)/<br>0.0001<br>= 0.6<br><br>$$\frac{df}{dx} = \lim_{h \to 0} \frac{f(x+h) - f(x)}{h}$$ |

# (example)

| Current **W**: | **W + h** (third dim): | Gradient $\frac{dL}{dW}$ |
|---|---|---|
| [0.34, | [0.34, | |
| -1.11, | -1.11, | 0.6, |
| 0.78, | 0.78 + **0.0001**, | 0.0, |
| 0.12, | 0.12, | ? |
| 0.55 | | |
| 2.81, | | |
| -3.1, | | |
| -1.5, | | |
| 0.33, ...] | 0.33, ...] | ?, ...] |
| loss 1.25347 | loss 1.25347 | |

① **Numeric Gradient:**
- Slow: O(# dimensions)
- Approximate

ROBOTICS

# Loss is a function of W

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

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

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

Want $\nabla_w L$

Use calculus to compute an

② **Analytic gradient**

# (example)

Current **W**:

Gradient $\frac{dL}{dW}$

[0.34,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33, …]
loss 1.25347

$\frac{dL}{dW}$ = some function of data and $W$

In practice we will compute $\frac{dL}{dW}$
using back propagation;
see Lecture 6

[-2.5,
0.6,
0.0,
0.2,
0.7,
-0.5,
1.1,
1.3,
-2.1, …]

**M** | ROBOTICS

# Computing Gradients

# Computing Gradients

① **Numeric gradient:** approximate, slow, easy to write

② **Analytic gradient:** exact, fast, error-prone

<u>In practice:</u> Always use analytic gradient, but check implementation with numerical gradient. This is called a **gradient check**.

```python
def grad_check_sparse(f, x, analytic_grad, num_checks=10, h=1e-7):
    """
    sample a few random elements and only return numerical
    in this dimensions.
    """
```

Also check out: https://cs231n.github.io/optimization-1/
https://pytorch.org/docs/stable/notes/gradcheck.html

**M | ROBOTICS**

# Computing Gradients

① **Numeric gradient:** approximate, slow, easy to write

② **Analytic gradient:** exact, fast, error-prone

```
torch.autograd.gradcheck(func, inputs, eps=1e-06, atol=1e-05, rtol=0.001,
raise_exception=True, check_sparse_nnz=False, nondet_tol=0.0)
```
[SOURCE]

Check gradients computed via small finite differences against analytical gradients w.r.t. tensors in `inputs` that are of floating point type and with `requires_grad=True`.

The check between numerical and analytical gradients uses `allclose()`.

# Computing Gradients

① **Numeric gradient:** approximate, slow, easy to write

② **Analytic gradient:** exact, fast, error-prone

---

```
torch.autograd.gradgradcheck(func, inputs, grad_outputs=None, eps=1e-06, atol=1e-
05, rtol=0.001, gen_non_contig_grad_outputs=False, raise_exception=True,          [SOURCE]
nondet_tol=0.0)
```

Check gradients of gradients computed via small finite differences against analytical gradients w.r.t. tensors in `inputs` and `grad_outputs` that are of floating point type and with `requires_grad=True`.

This function checks that backpropagating through the gradients computed to the given `grad_outputs` are correct.

# Gradient Descent

- Iteratively step in the direction of the negative gradient (direction of local steepest descent)

```
# Vanilla gradient descent
w = initialize_weights()
for t in range(num_steps):
    dw = compute_gradient(loss_fn, data, w)
    w -= learning_rate * dw
```

**Hyperparameters:**

- Weight initialization method

- Number of steps

- Learning rate

Q4: guarantee?   https://ahaslides.com/WJTNO

Negative gradient direction
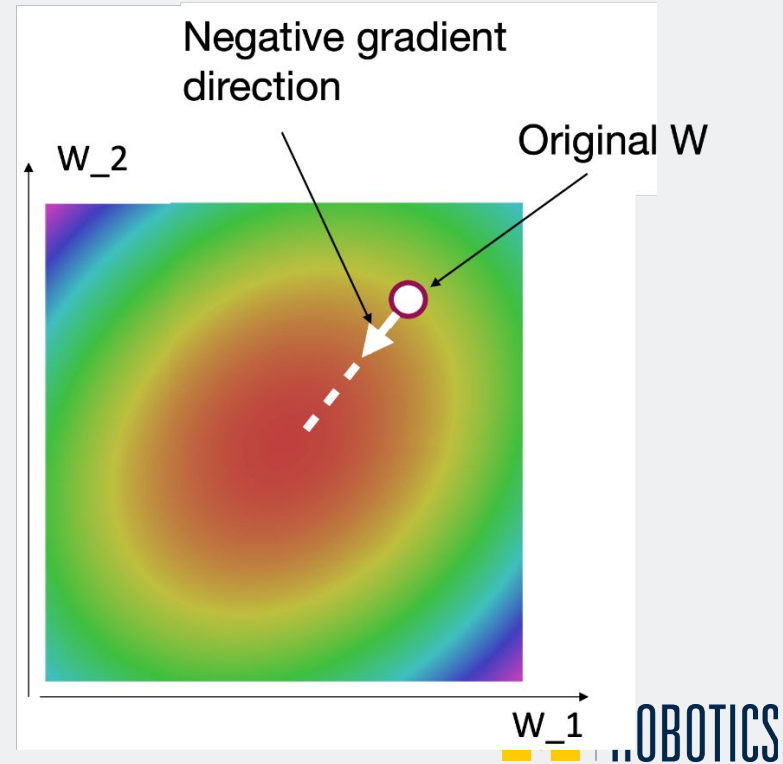
Original W

W_2

W_1

.OBOTICS

# Gradient Descent

- Iteratively step in the direction of the negative gradient (direction of local steepest descent)

```
# Vanilla gradient descent
w = initialize_weights()
for t in range(num_steps):
  dw = compute_gradient(loss_fn, data, w)
  w -= learning_rate * dw
```

**Hyperparameters:**

- Weight initialization method
- Number of steps
- Learning rate

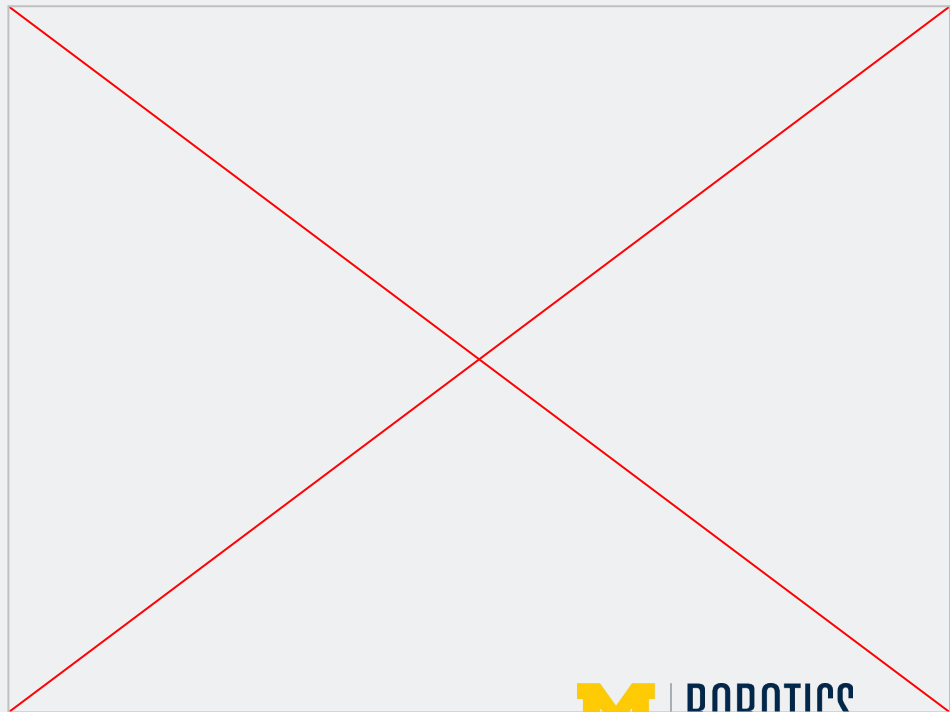# Batch Gradient Descent

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

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

Full sum expensive when N is large!

# Stochastic Gradient Descent (SGD)

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

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

```
# Stochastic gradient descent
w = initialize_weights()
for t in range(num_steps):
  minibatch = sample_data(data, batch_size)
  dw = compute_gradient(loss_fn, minibatch, w)
  w -= learning_rate * dw
```

Full sum expensive when N is large!

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

**Hyperparameters:**
- Weight initialization
- Number of steps
- Learning rate
- Batch size
- Data sampling

M|ROBOTICS

# Stochastic Gradient Descent (SGD)

$$L(W) = \mathbb{E}_{(x,y) \sim p_{data}}[L(x, y, W)] + \lambda R(W)$$

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

Think of loss as an expectation over the full **data distribution** $p_{data}$

Approximate expectation via sampling

$$\nabla_W L(W) = \nabla_W \mathbb{E}_{(x,y) \sim p_{data}}[L(x, y, W)] + \lambda R(W)$$

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

For reference: an interactive web demo:
http://vision.stanford.edu/teaching/cs231n-demos/linear-classify/

**ROBOTICS**

# Aha Slides
# (In-class participation)

https://ahaslides.com/WJTNO

Q5: drawbacks/problem w/ SGD

# Problem with SGD ①

What if loss changes quickly in one direction and slowly in another?

What does gradient decent do?



Loss function has high condition number: ratio of largest to smallest singular value of the Hessian matrix is large

# Problem with SGD ①

What if loss changes quickly in one direction and slowly in another?

What does gradient decent do?

Very slow progress along shallow dimension, jitter along steep direction



Loss function has high condition number: ratio of largest to smallest singular value of the Hessian matrix is large

# Problem with SGD ②



Local Minimum

Saddle point

What if the loss function has a local minimum or saddle point?

ROBOTICS

# Problem with SGD ②



Local Minimum

Saddle point

What if the loss function has a local minimum or saddle point?

Zero gradient, gradient descent gets stuck

# Problem with SGD ③

Our gradients come from mini batches so they can be noisy!

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

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

# Problem with SGD



Local Minimum

Saddle point

What if the loss function has a local minimum or saddle point?

Batched gradient descent always computes same gradients

SGD computes noisy gradients, may help to escape saddle points

# More than SGD…

# SGD + Momentum

## SGD

$$w_{t+1} = w_t - \alpha \nabla L(w_t)$$

```
for t in range(num_steps):
    dw = compute_gradient(w)
    w -= learning_rate * dw
```



## SGD + Momentum

$$v_{t+1} = \rho v_t + \nabla L(w_t)$$
$$w_{t+1} = w_t - \alpha v_{t+1}$$

```
v = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    v = rho * v + dw
    w -= learning_rate * v
```

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

ROBOTICS

# SGD + Momentum

## Momentum update:



Velocity

Actual step

Gradient

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

## SGD + Momentum

$$v_{t+1} = \rho v_t + \nabla L(w_t)$$
$$w_{t+1} = w_t - \alpha v_{t+1}$$

```
v = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    v = rho * v + dw
    w -= learning_rate * v
```

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

# SGD + Momentum

SGD + Momentum

$$v_{t+1} = \rho v_t - \alpha \nabla L(w_t)$$

$$w_{t+1} = w_t + v_{t+1}$$

```
v = 0
for t in range(num_steps):
  dw = compute_gradient(w)
  v = rho * v - learning_rate * dw
  w += v
```

SGD + Momentum

$$v_{t+1} = \rho v_t + \nabla L(w_t)$$

$$w_{t+1} = w_t - \alpha v_{t+1}$$

```
v = 0
for t in range(num_steps):
  dw = compute_gradient(w)
  v = rho * v + dw
  w -= learning_rate * v
```

You may see SGD+Momentum formulated different ways, but they are equivalent - give same sequence of w

Sutskever et al, "On the importance of initialization and momentum in deep learning," ICML 2013

ROBOTICS

# SGD + Momentum

Local Minima     Saddle Points



Poor Conditioning





SGD     SGD+Momentum

# SGD + Momentum

Momentum update:



Velocity

Actual step

Gradient

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

Nesterov, "A method of solving a convex programming problem with convergence rate O(1/k^2),", 1983"
Nesterov, "Introductory lectures on convex optimization: a basic course," 2004
Sutskever et al, "On the importance of initialization and momentum in deep learning," ICML 2013

# Nesterov Momentum

Momentum update:



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

## Nesterov Momentum



"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

Nesterov, "A method of solving a convex programming problem with convergence rate O(1/k^2),", 1983"
Nesterov, "Introductory lectures on convex optimization: a basic course," 2004
Sutskever et al, "On the importance of initialization and momentum in deep learning," ICML 2013

**M** | ROBOTICS

# Nesterov Momentum

Annoying, usually we
want to update in terms of $w_t, \nabla L(w_t)$



Velocity

Gradient

Actual step

$$v_{t+1} = \rho v_t - \alpha \nabla L(w_t + \rho v_t)$$

$$w_{t+1} = w_t + v_{t+1}$$

"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

ROBOTICS

# Nesterov Momentum

Annoying, usually we
want to update in terms of $w_t, \nabla L(w_t)$

Velocity

Gradient

Actual step

```
v = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    old_v = v
    v = rho * v - learning_rate * dw
    w -= rho * old_v - (1 + rho) * v
```

$$v_{t+1} = \rho v_t - \alpha \nabla L(\boxed{w_t + \rho v_t})$$

$$w_{t+1} = w_t + v_{t+1}$$

Change of variables
and rearrange:    $\tilde{w}_t = w_t + \rho v_t$

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

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

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

ROBOTICS

# Nesterov Momentum



SGD

SGD+Momentum

Nesterov

ROBOTICS

# AdaGrad

```
grad_squared = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    grad_squared += dw * dw
    w -= learning_rate * dw / (grad_squared.sqrt() + 1e-7)
```

- Added element-wise scaling of the gradient based on the historical sum of squares in each dimension
- "Per-parameter learning rates" or "adaptive learning rates"

# AdaGrad

```
grad_squared = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    grad_squared += dw * dw
    w -= learning_rate * dw / (grad_squared.sqrt() + 1e-7)
```

Problem: AdaGrad will slow over many iterations



**Q: What happens with AdaGrad?**

Progress along "steep" directions is damped; progress along "flat" directions is accelerated

Duchi et al, "Adaptive sub gradient methods for online learning and stochastic optimization," JMLR 2011

# RMSProp: "Leaky AdaGrad"

```
grad_squared = 0
for t in range(num_steps):
  dw = compute_gradient(w)
  grad_squared += dw * dw
  w -= learning_rate * dw / (grad_squared.sqrt() + 1e-7)
```

AdaGrad

```
grad_squared = 0
for t in range(num_steps):
  dw = compute_gradient(w)
  grad_squared = decay_rate * grad_squared + (1 - decay_rate) * dw * dw
  w -= learning_rate * dw / (grad_squared.sqrt() + 1e-7)
```

RMSProp

M | ROBOTICS

# RMSProp: "Leaky AdaGrad"

# RMSProp + Momentum ("Almost" Adam)

```python
moment1 = 0
moment2 = 0
for t in range(1, num_steps + 1):  # Start at t = 1
  dw = compute_gradient(w)
  moment1 = beta1 * moment1 + (1 - beta1) * dw
  moment2 = beta2 * moment2 + (1 - beta2) * dw * dw
  w -= learning_rate * moment1 / (moment2.sqrt() + 1e-7)
```

# RMSProp + Momentum ("Almost" Adam)

```
moment1 = 0
moment2 = 0
for t in range(1, num_steps + 1):  # Start at t = 1
  dw = compute_gradient(w)
  moment1 = beta1 * moment1 + (1 - beta1) * dw
  moment2 = beta2 * moment2 + (1 - beta2) * dw * dw
  w -= learning_rate * moment1 / (moment2.sqrt() + 1e-7)
```

Adam

Momentum

```
v = 0
for t in range(num_steps):
  dw = compute_gradient(w)
  v = rho * v + dw
  w -= learning_rate * v
```

SGD+Momentum

# RMSProp + Momentum ("Almost" Adam)

```python
moment1 = 0
moment2 = 0
for t in range(1, num_steps + 1):  # Start at t = 1
    dw = compute_gradient(w)
    moment1 = beta1 * moment1 + (1 - beta1) * dw
    moment2 = beta2 * moment2 + (1 - beta2) * dw * dw
    w -= learning_rate * moment1 / (moment2.sqrt() + 1e-7)
```

Adam

Momentum

AdaGrad / RMSProp

```python
grad_squared = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    grad_squared = decay_rate * grad_squared + (1 - decay_rate) * dw * dw
    w -= learning_rate * dw / (grad_squared.sqrt() + 1e-7)
```

RMSProp

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

**M | ROBOTICS**

# RMSProp + Momentum ("Almost" Adam)

```python
moment1 = 0
moment2 = 0
for t in range(1, num_steps + 1):  # Start at t = 1
    dw = compute_gradient(w)
    moment1 = beta1 * moment1 + (1 - beta1) * dw
    moment2 = beta2 * moment2 + (1 - beta2) * dw * dw
    w -= learning_rate * moment1 / (moment2.sqrt() + 1e-7)
```

Adam

Momentum

AdaGrad / RMSProp

Q: What happens at t=1?
(Assume beta2 = 0.999)

ROBOTICS

# RMSProp + Momentum ("Almost" Adam)

```python
moment1 = 0
moment2 = 0
for t in range(1, num_steps + 1):  # Start at t = 1
    dw = compute_gradient(w)
    moment1 = beta1 * moment1 + (1 - beta1) * dw
    moment2 = beta2 * moment2 + (1 - beta2) * dw * dw
    moment1_unbias = moment1 / (1 - beta1 ** t)
    moment2_unbias = moment2 / (1 - beta2 ** t)
    w -= learning_rate * moment1_unbias / (moment2_unbias.sqrt() + 1e-7)
```

Momentum

AdaGrad / RMSProp

Bias correction

**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, 5e-4, 1e-4 is a great starting point for many models!

Kingma, D. P. (2014). Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980.

M | ROBOTICS

# Adam: Very common in practice!

for input to the CNN; each colored pixel in the image yields a 7D one-hot vector. **Following common practice, the network is trained end-to-end using stochastic gradient descent with the Adam optimizer [22].** We anneal the learning rate to 0 using a half cosine schedule without restarts [28].

Bakhtin, van der Maaten, Johnson, Gustafson, and Girshick, NeurIPS 2019

We train all models using Adam [23] with learning rate $10^{-4}$ and batch size 32 for 1 million iterations; training takes about 3 days on a single Tesla P100. For each mini-batch we first update $f$, then update $D_{img}$ and $D_{obj}$.

Johnson, Gupta, and Fei-Fei, CVPR 2018

ganized into three residual blocks. **We train for 25 epochs using Adam [27] with learning rate $10^{-4}$ and 32 images per batch on 8 Tesla V100 GPUs.** We set the `cubify` thresh-

Gkioxari, Malik, and Johnson, ICCV 2019

sampled with each bit drawn uniformly at random. **For gradient descent, we use Adam [29] with a learning rate of $10^{-3}$ and default hyperparameters. All models are trained with batch size 12.** Models are trained for 200 epochs, or 400 epochs if being trained on multiple noise layers.

Zhu, Kaplan, Johnson, and Fei-Fei, ECCV 2018

16 dimensional vectors. We iteratively train the Generator and Discriminator **with a batch size of 64 for 200 epochs using Adam [22] with an initial learning rate of 0.001.**

Gupta, Johnson, et al, CVPR 2018

➡ Adam with beta1 = 0.9, beta2 = 0.999, and learning_rate = 1e-3, 5e-4, 1e-4 is a great starting point for many models!

M | ROBOTICS

# Adam: Very common in practice!

Additional References:

https://towardsdatascience.com/a-visual-explanation-of-gradient-descent-methods-momentum-adagrad-rmsprop-adam-f898b102325c

https://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf

ROBOTICS

# Optimization Algorithms Comparison

| Algorithm | Tracks first moments (Momentum) | Tracks second moments (Adaptive learning rates) | Leaky second moments | Bias correction for moment estimates |
|---|---|---|---|---|
| SGD | ✗ | ✗ | ✗ | ✗ |
| SGD+Momentum | ✓ | ✗ | ✗ | ✗ |
| Nesterov | ✓ | ✗ | ✗ | ✗ |
| AdaGrad | ✗ | ✓ | ✗ | ✗ |
| RMSProp | ✗ | ✓ | ✓ | ✗ |
| Adam | ✓ | ✓ | ✓ | ✓ |

# L2 Regularization vs. Weight Decay

**Optimization Algorithm**

$$L(w) = L_{data}(w) + L_{reg}(w)$$

$$g_t = \nabla L(w_t)$$

$$s_t = optimizer(g_t)$$

$$w_{t+1} = w_t - \alpha s_t$$

L2 Regularization and Weight Decay are equivalent for SGD, SGD+Momentum so people often use the terms interchangeably!

But they are not the same for adaptive methods (AdaGrad, RMSProp, Adam, etc)

Loshchilov and Hunter, "Decoupled Weight Decay Regularization," ICLR 2019

**L2 Regularization**

$$L(w) = L_{data}(w) + \lambda |w|^2$$

$$g_t = \nabla L(w_t) = \nabla L_{data}(w_t) + 2\lambda w_t$$

$$s_t = optimizer(g_t)$$

$$w_{t+1} = w_t - \alpha s_t$$

**Weight decay**

$$L(w) = L_{data}(w)$$

$$g_t = \nabla L_{data}(w_t)$$

$$s_t = optimizer(g_t) + 2\lambda w_t$$

$$w_{t+1} = w_t - \alpha s_t$$

ROBOTICS

# AdamW: Decouple Weight Decay

**Algorithm 2** Adam with L$_2$ regularization and Adam with decoupled weight decay (AdamW)

1: **given** $\alpha = 0.001$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 10^{-8}$, $\lambda \in \mathbb{R}$

2: **initialize** time step $t \leftarrow 0$, parameter vector $\boldsymbol{\theta}_{t=0} \in \mathbb{R}^n$, first moment vector $\boldsymbol{m}_{t=0} \leftarrow \boldsymbol{0}$, second moment vector $\boldsymbol{v}_{t=0} \leftarrow \boldsymbol{0}$, schedule multiplier $\eta_{t=0} \in \mathbb{R}$

3: **repeat**

4:     $t \leftarrow t + 1$

5:     $\nabla f_t(\boldsymbol{\theta}_{t-1}) \leftarrow \text{SelectBatch}(\boldsymbol{\theta}_{t-1})$          ▷ select batch and return the corresponding gradient

6:     $\boldsymbol{g}_t \leftarrow \nabla f_t(\boldsymbol{\theta}_{t-1}) \; +\lambda\boldsymbol{\theta}_{t-1}$

7:     $\boldsymbol{m}_t \leftarrow \beta_1\boldsymbol{m}_{t-1} + (1 - \beta_1)\boldsymbol{g}_t$          ▷ here and below all operations are element-wise

8:     $\boldsymbol{v}_t \leftarrow \beta_2\boldsymbol{v}_{t-1} + (1 - \beta_2)\boldsymbol{g}_t^2$

9:     $\hat{\boldsymbol{m}}_t \leftarrow \boldsymbol{m}_t/(1 - \beta_1^t)$          ▷ $\beta_1$ is taken to the power of $t$

10:    $\hat{\boldsymbol{v}}_t \leftarrow \boldsymbol{v}_t/(1 - \beta_2^t)$          ▷ $\beta_2$ is taken to the power of $t$

11:    $\eta_t \leftarrow \text{SetScheduleMultiplier}(t)$          ▷ can be fixed, decay, or also be used for warm restarts

12:    $\boldsymbol{\theta}_t \leftarrow \boldsymbol{\theta}_{t-1} - \eta_t \left( \alpha\hat{\boldsymbol{m}}_t/(\sqrt{\hat{\boldsymbol{v}}_t} + \epsilon) \; +\lambda\boldsymbol{\theta}_{t-1} \right)$

13: **until** *stopping criterion is met*

14: **return** optimized parameters $\boldsymbol{\theta}_t$

Loshchilov and Hunter, "Decoupled Weight Decay Regularization," ICLR 2019

ROBOTICS

# AdamW: Decouple Weight Decay

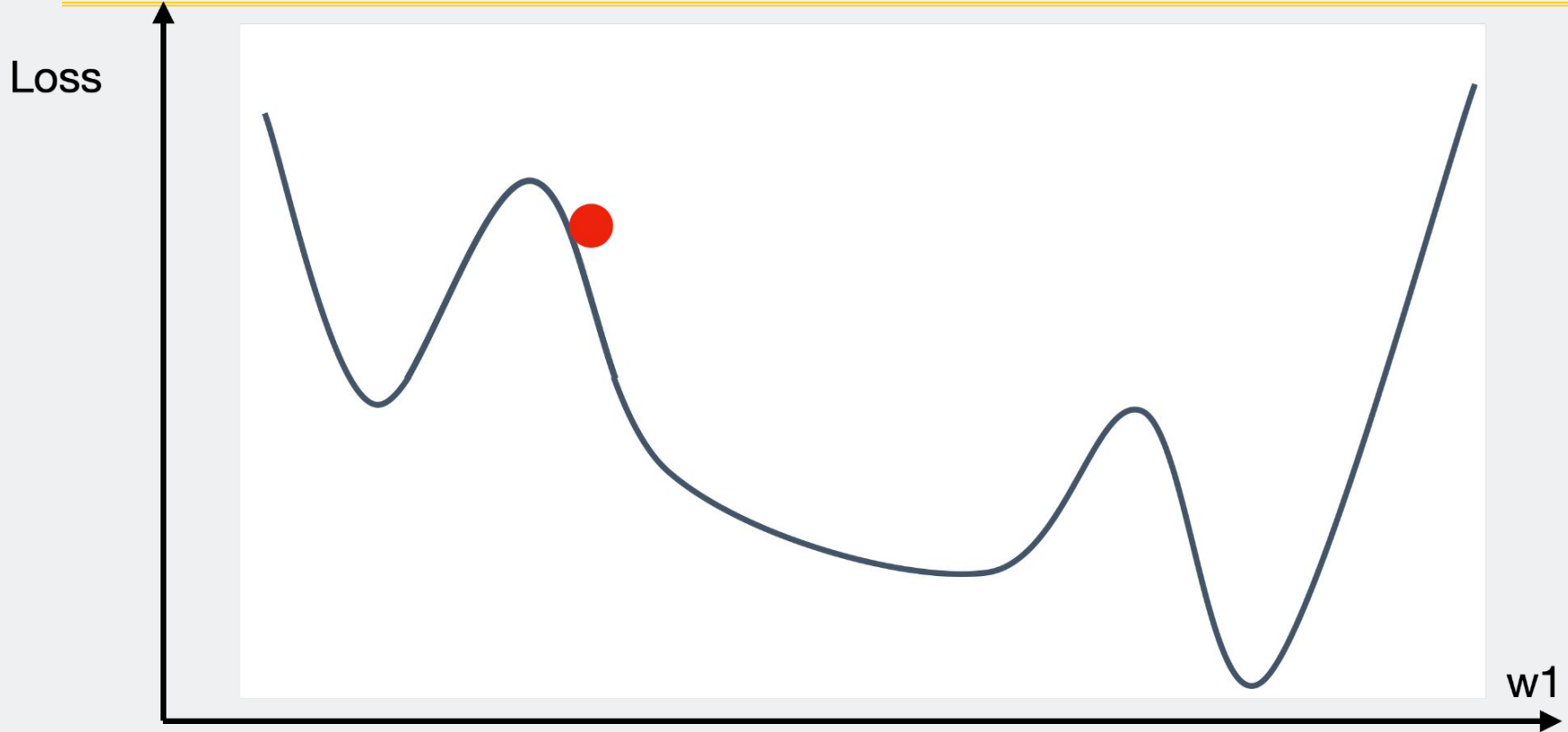**Algorithm 2** Adam with L$_2$ regularization and Adam with decoupled weight decay (AdamW)

1: **given** $\alpha = 0.001, \beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}, \lambda \in \mathbb{R}$
2: **initialize** time step $t \leftarrow 0$, parameter vector $\boldsymbol{\theta}_{t=0} \in \mathbb{R}^n$, first moment vector $\boldsymbol{m}_{t=0} \leftarrow \boldsymbol{0}$, second moment vector $\boldsymbol{v}_{t=0} \leftarrow \boldsymbol{0}$, schedule multiplier $\eta_{t=0} \in \mathbb{R}$
3:
4:
5:
6:
7:
8:
9:
10:
11:      $\eta_t \leftarrow \text{SetScheduleMultiplier}(t)$     $\triangleright$ can be fixed, decay, or also be used for warm restarts
12:      $\boldsymbol{\theta}_t \leftarrow \boldsymbol{\theta}_{t-1} - \eta_t \left( \alpha \hat{\boldsymbol{m}}_t / (\sqrt{\hat{\boldsymbol{v}}_t} + \epsilon) + \lambda \boldsymbol{\theta}_{t-1} \right)$
13: **until** *stopping criterion is met*
14: **return** optimized parameters $\boldsymbol{\theta}_t$

AdamW should/could probably be your "default" optimizer for new problems

Loshchilov and Hunter, "Decoupled Weight Decay Regularization," ICLR 2019

M | ROBOTICS

# Second-Order Optimization

# So Far: First-Order Optimization



Loss

w1

# So Far: First-Order Optimization

1. Use gradient to make linear approximation
2. Step to minimize the approximation



Loss

w1

M | ROBOTICS

# Second-Order Optimization



1. Use gradient and Hessian to make quadratic approximation
2. Step to minimize the approximation

Loss

w1

ROBOTICS

# Second-Order Optimization



1. Use gradient and Hessian to make quadratic approximation
2. Step to minimize the approximation

Take bigger steps in areas of low curvature

Loss

w1

# Second-Order Optimization

Second-order Taylor Expansion:

$$L(w) \approx L(w_0) + (w - w_0)^T \nabla_w L(w_0) + \frac{1}{2}(w - w_0)^T H_w L(w_0)(w - w_0)$$

Solving for the critical point we obtain the Newton parameter update:

$$w^* = w_0 - \mathbf{H}_w L(w_0)^{-1} \nabla_w L(w_0)$$

**Q: Why is this impractical?**

Hessian has O(N^2) elements

Inverting takes O(N^3)

N = (Tens or Hundreds of) Millions

# Second-Order Optimization

$$w^* = w_0 - \mathbf{H}_w L(w_0)^{-1} \nabla_w L(w_0)$$

- Quasi-Newton methods (BGFS most popular): *instead of inverting the Hessian ((O(n^3)), approximate inverse Hessian with rank 1 updates over time (O(n^2) each).*

- **L-BFGS** (Limited memory BFGS): *Does not form/store the full inverse Hessian*

# Second-Order Optimization: L-BFGS

- **Usually works very well in full batch, deterministic mode** i.e. if you have a single, deterministic f(x) then L-BFGS will probably work very nicely.

- **Does not transfer very well to mini-batch setting.** Gives bad results. Adapting second-order methods to large-scale, stochastic setting is an active area of research.

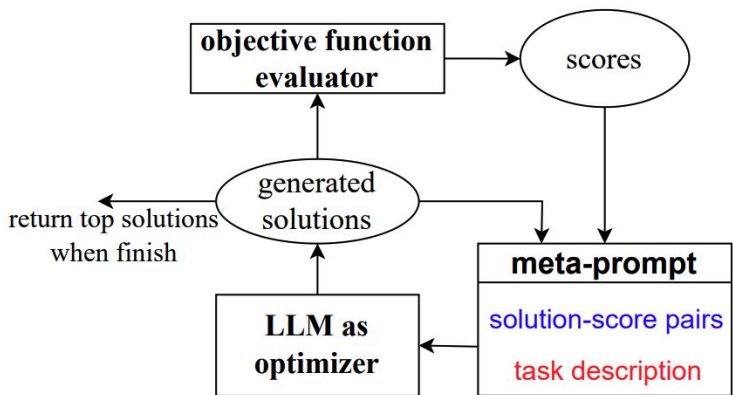Le et al, "On optimization methods for deep learning," ICML 2011
Ba et al, "Distributed second-order optimization using Kronecker-factored approximations," ICLR 2017

# In-Practice (Take-aways)

- Adam is a good default choice in many cases. SGD+Momentum can outperform Adam but may require more tuning.
- If you can afford to do full batch updates then try out L-BFGS (and don't forget to disable all sources of noise)

- ## Large Language Models as Optimizers



Table 1: Top instructions with the highest GSM8K zero-shot test accuracies from prompt optimization with different optimizer LLMs. All results use the pre-trained `PaLM 2-L` as the scorer.

| Source | Instruction | Acc |
|---|---|---|
| *Baselines* | | |
| (Kojima et al., 2022) | Let's think step by step. | 71.8 |
| (Zhou et al., 2022b) | Let's work this out in a step by step way to be sure we have the right answer. | 58.8 |
| | (empty string) | 34.0 |
| *Ours* | | |
| `PaLM 2-L-IT` | Take a deep breath and work on this problem step-by-step. | **80.2** |
| `PaLM 2-L` | Break this down. | 79.9 |
| `gpt-3.5-turbo` | A little bit of arithmetic and a logical approach will help us quickly arrive at the solution to this problem. | 78.5 |
| `gpt-4` | Let's combine our numerical command and clear thinking to quickly and accurately decipher the answer. | 74.5 |

"meta-prompt"

Large Language Models as Optimizers. Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V Le, Denny Zhou, Xinyun Chen. https://iclr.cc/virtual/2024/poster/19209

- Neural Topic Modeling as Multi-Objective Contrastive Optimization

$$\min_{\alpha} \left\{ \left\| \alpha \nabla_\theta \mathcal{L}_{\text{InfoNCE}}(\theta) + (1 - \alpha) \nabla_\theta \mathcal{L}_{\text{ELBO}}(\theta, \phi) \right\|_2^2 \Big| \alpha \geq 0 \right\}$$

| | | NTM+CL | Our Model |
|---|---|---|---|
| shuttle lands on planet | job career ask development | 0.0093 | 0.0026 |
| | star astronaut planet light moon | 0.8895 | 0.9178 |
| shuttle lands on planet *zeppelin/scardino* | job career ask development *zeppelin/s-cardino* | 0.9741/0.9413 | 0.0064/0.0080 |
| | star astronaut planet light moon | 0.1584/0.2547 | 0.8268/0.7188 |

https://openreview.net/pdf?id=HdAoLSBYXj

ROBOTICS

# State-of-the-Art

- Neural Topic Modeling as Multi-Objective Contrastive Optimization

$$\min_{\alpha} \left\{ \left\| \alpha \nabla_\theta \mathcal{L}_{\text{InfoNCE}}(\theta) + (1-\alpha) \nabla_\theta \mathcal{L}_{\text{ELBO}}(\theta, \phi) \right\|_2^2 \,\middle|\, \alpha \geq 0 \right\}$$

$$f(\mathbf{x}, \mathbf{y}) = \frac{g_\varphi(\mathbf{x})^T g_\varphi(\mathbf{y})}{\| g_\varphi(\mathbf{x}) \| \| g_\varphi(\mathbf{y}) \|} / \tau$$
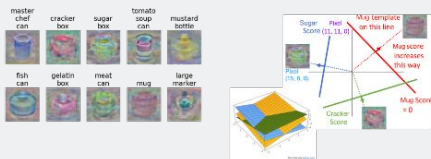
$$\min_{\theta, \phi} \mathcal{L}_{\text{ELBO}} = -\mathbb{E}_{q_\theta(\mathbf{z}|\mathbf{x})} \left[ \log p_\phi(\mathbf{x}|\mathbf{z}) \right] + \mathbb{KL} \left[ q_\theta(\mathbf{z}|\mathbf{x}) \,\|\, p(\mathbf{z}) \right]$$

https://openreview.net/pdf?id=HdAoLSBYXj

ROBOTICS

# Summary

# Summary

- Use **Linear Models** for image classification problems.

- Use **Loss Functions** to express preferences over different choices of weights.

- Use **Regularization** to prevent overfitting to training data.

- Use **Stochastic Gradient Descent** to minimize our loss functions and train the model.
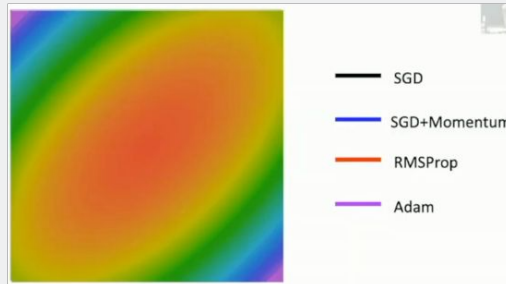
$$s = f(x; W) = Wx$$



$$L_i = -\log\left(\frac{\exp^{s_{y_i}}}{\sum_i \exp^{s_j}}\right) \quad \textbf{Softmax}$$

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

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

```
v = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    v = rho * v + dw
    w -= learning_rate * v
```



| | SGD |
| | SGD+Momentum |
| | RMSProp |
| | Adam |

Next up: Neural Networks

ROBOTICS

# Due dates

**Canvas Assignment:** 20250122 Optimization Quiz

**Scored - individual** (as part of in-class activity points)

Due Sunday Jan. 26, 2025

**P1 (KNN and Linear Classifier)**

**5 submissions per day - Start today!!!**

**Due Feb. 2, 2025**

**M | ROBOTICS**

# Enrollment/Waitlist

**Please send us your UniqName  (or reply via email)**
**by Thursday Jan. 23 5pm EST**
**if you intend to enroll in the class**

1.  **ROB 498 or 599**
2.  **Your UniqName**

**https://piazza.com/class/m4pgejar4ua2qf/post/33**

ROBOTICS

# Office Hour Calendar Now Available

https://calendar.google.com/calendar/u/0?cid=Y18zZDZhOGMyMTg0Y2I3ZDA4ZmIwZDg4OGM1OWNiNTU0OGViNzczMTZiOTg3ZTE3YmFlYjFkZDkwOWRhZWQyZTc2QGdyb3VwLmNhbGVuZGFyLmdvb2dsZS5jb20

You can add this calendar to your UM google calendar.

ROBOTICS