# DeepRob

Lecture 3
Regularization + Optimization
University of Michigan | Department of Robotics

# Recap: Image Classification
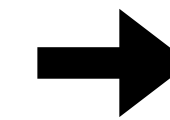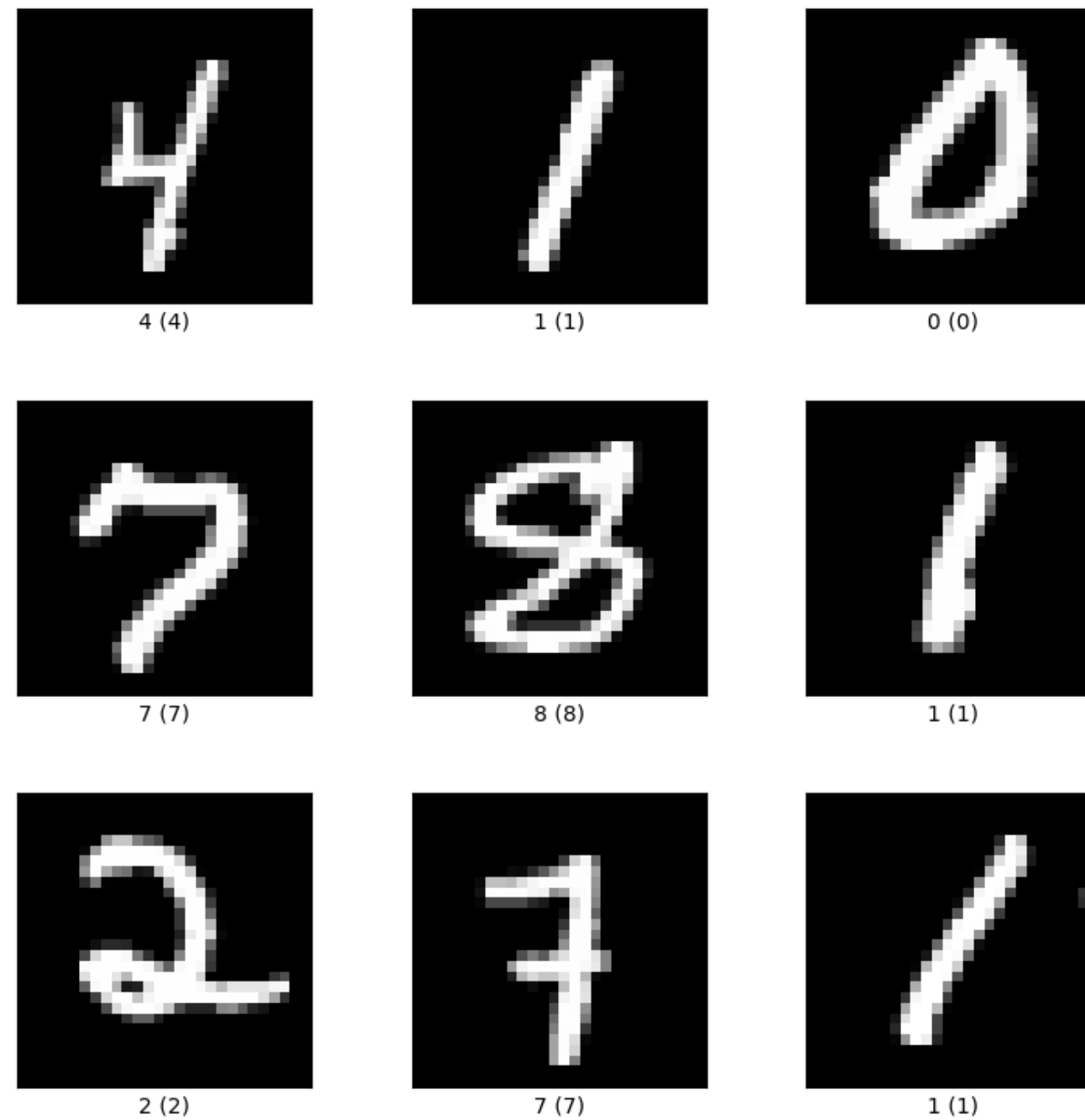
PROPS dataset

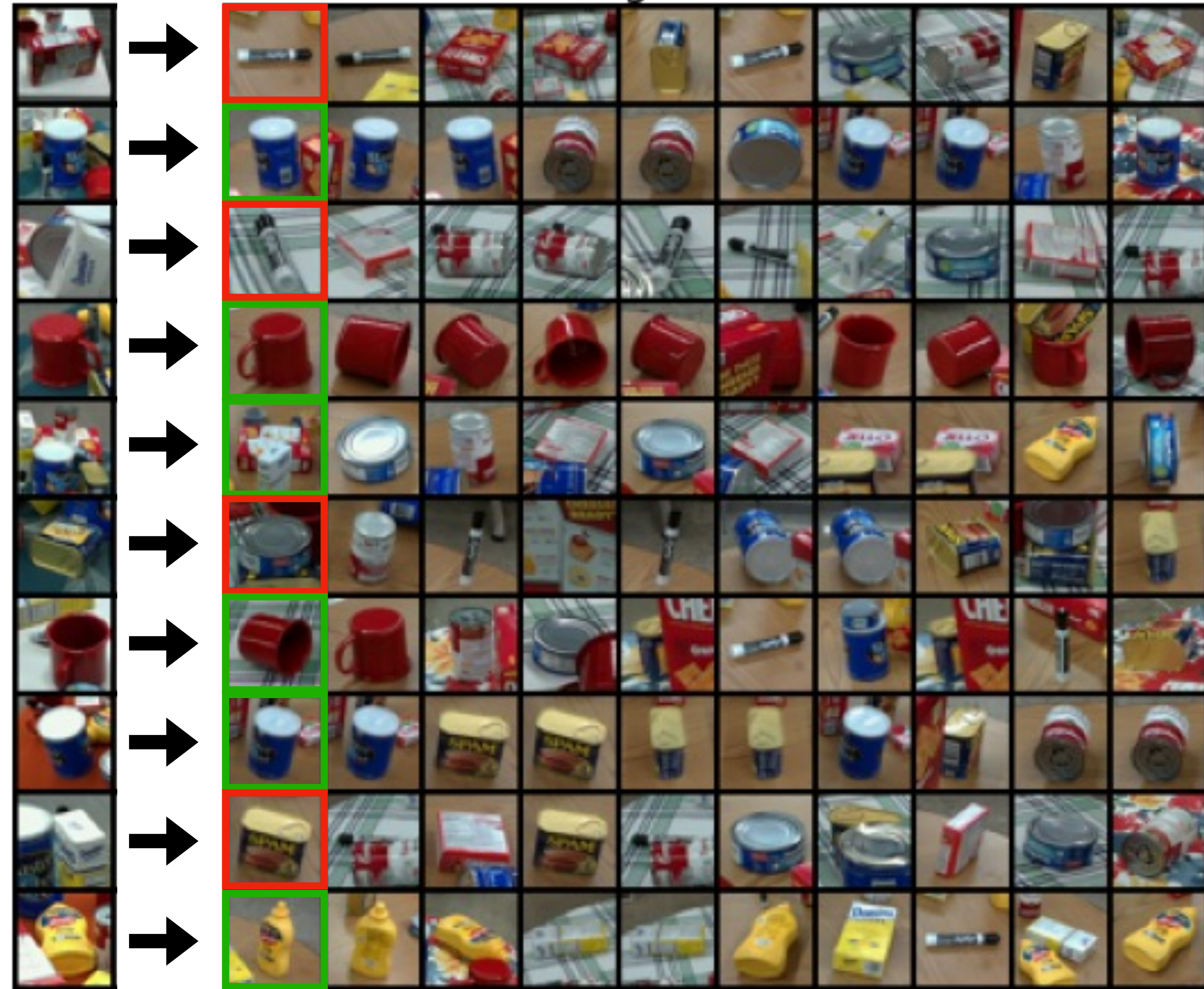# Recap: Image Classification

Labels

MNIST dataset → 

| 4 | 1 | 0 |
|---|---|---|
| 7 | 8 | 1 |
| 2 | 7 | 1 |

# Recap: K Nearest Neighbor

PROPS dataset

# KNN Pseudocode

1. Load training and testing data

2. Choose Hyperparameters (K=?)

3. For each point (image) in test data:

    find the distance to all training data points

    store the distance and sort it

    choose the first K points

    assign a class to the test image based on the majority of the classes

    End

# KNN – Some things to note

1. Hyperparameters: choose from k_choices
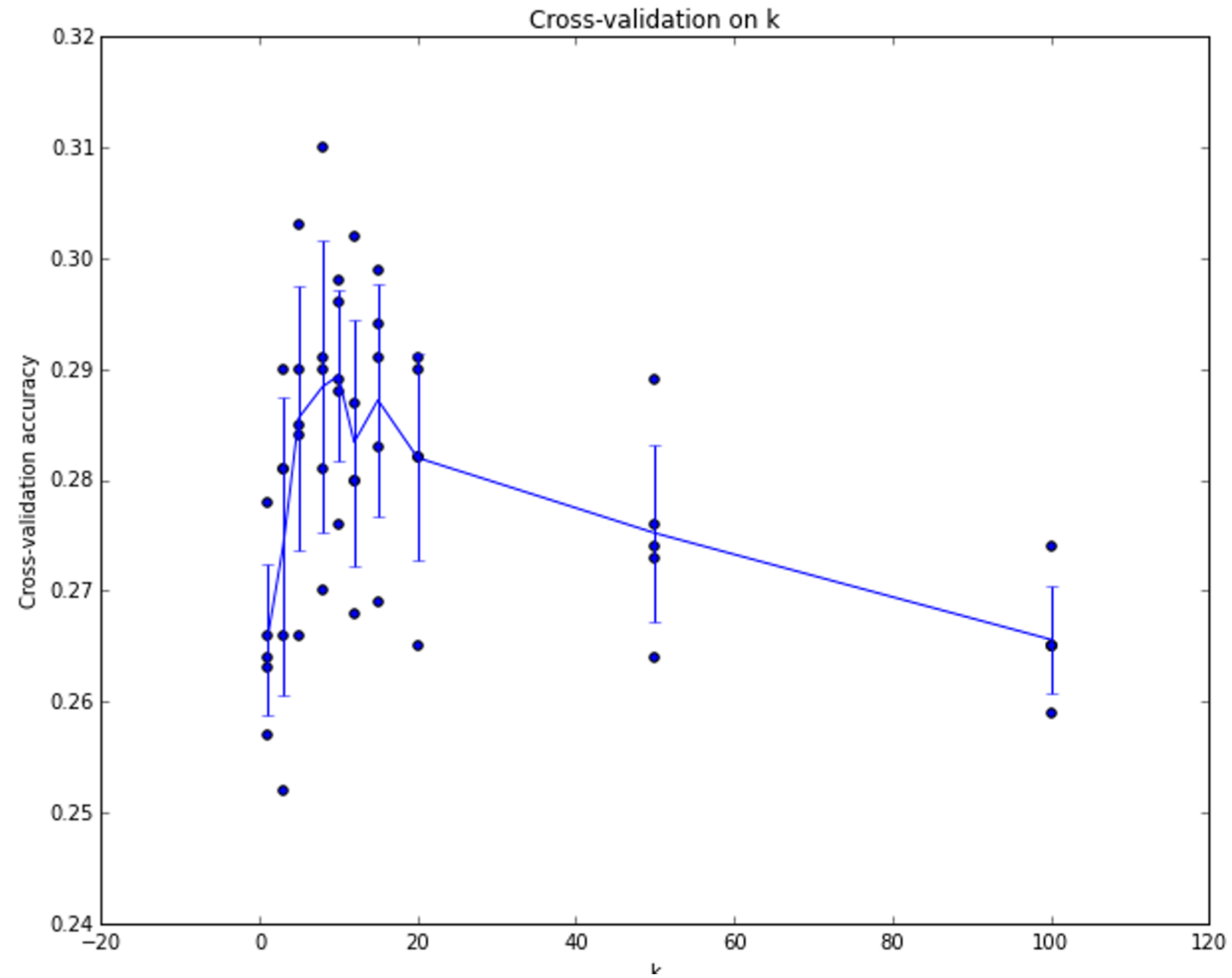
2. Cross-validation  (e.g., 5-fold validation)

   First, split the data into **folds**       `torch.chunk`
   Then, use all but one fold for train and one fold for validation

| fold 1 | fold 2 | fold 3 | fold 4 | fold 5 | test |
|--------|--------|--------|--------|--------|------|
| fold 1 | fold 2 | fold 3 | fold 4 | fold 5 | test |
| fold 1 | fold 2 | fold 3 | fold 4 | fold 5 | test |

# Setting Hyperparameters



Cross-validation on k

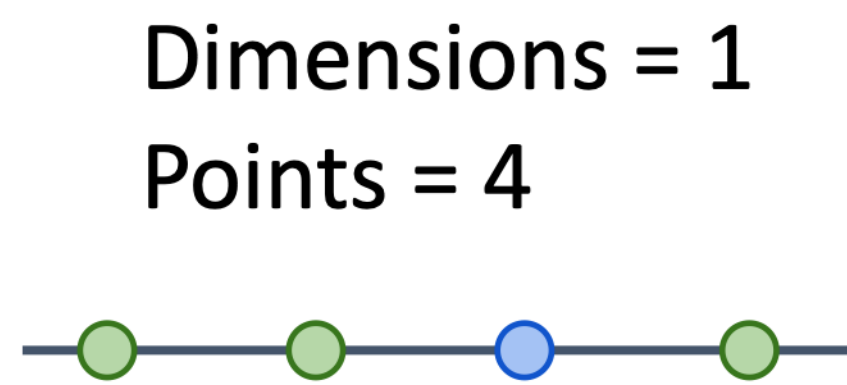Example of 5-fold cross-validation for the value of **k.**

Each point: single outcome.

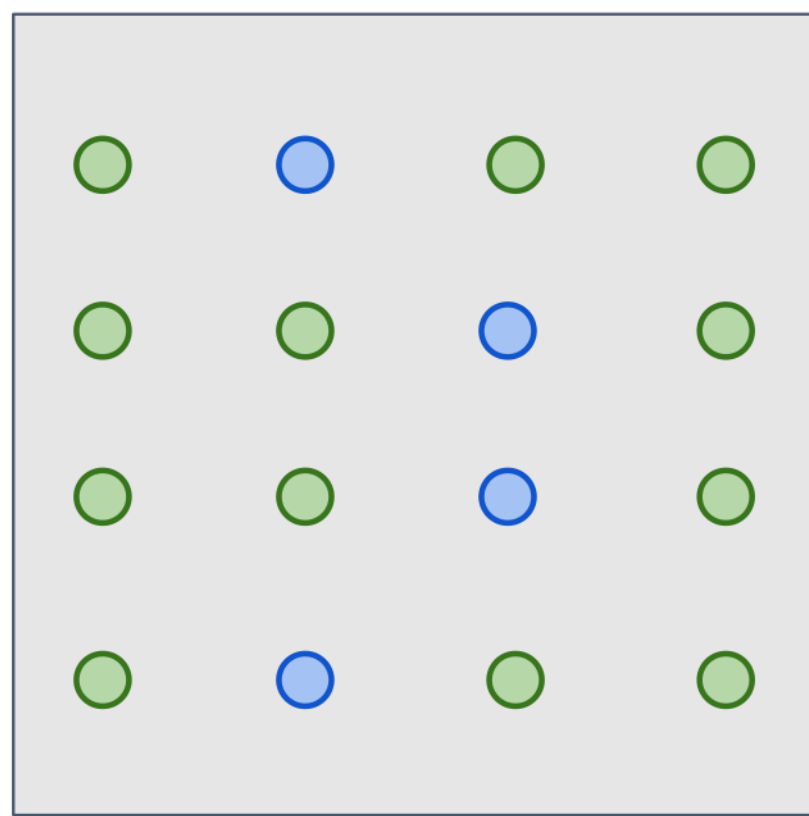The line goes through the mean, bars indicated standard deviation
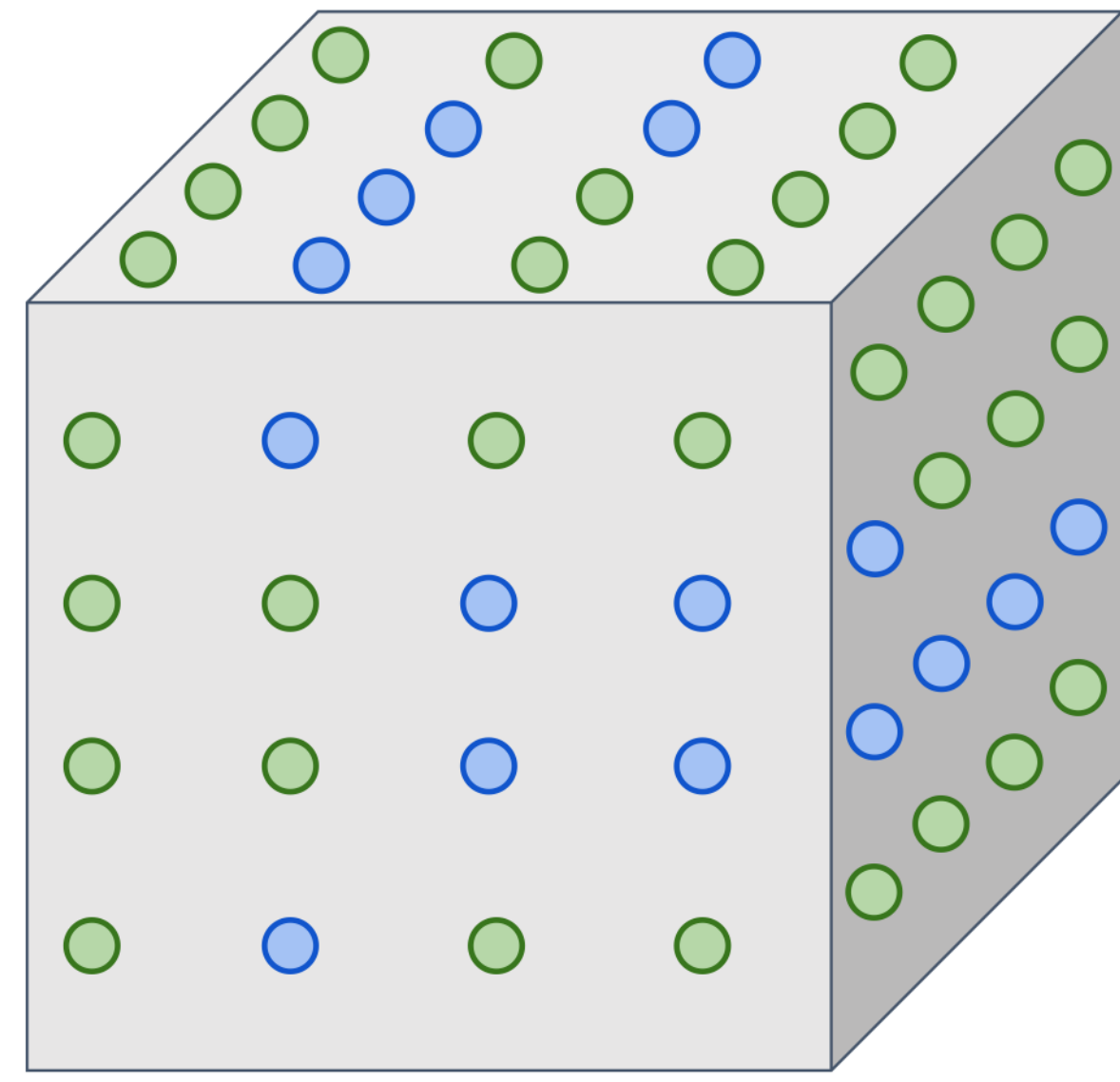
# Problem—Curse of Dimensionality

**Curse of dimensionality:** For uniform coverage of space, number of training points needed grows exponentially with dimension

Dimensions = 1
Points = 4

Dimensions = 2
Points = $4^2$
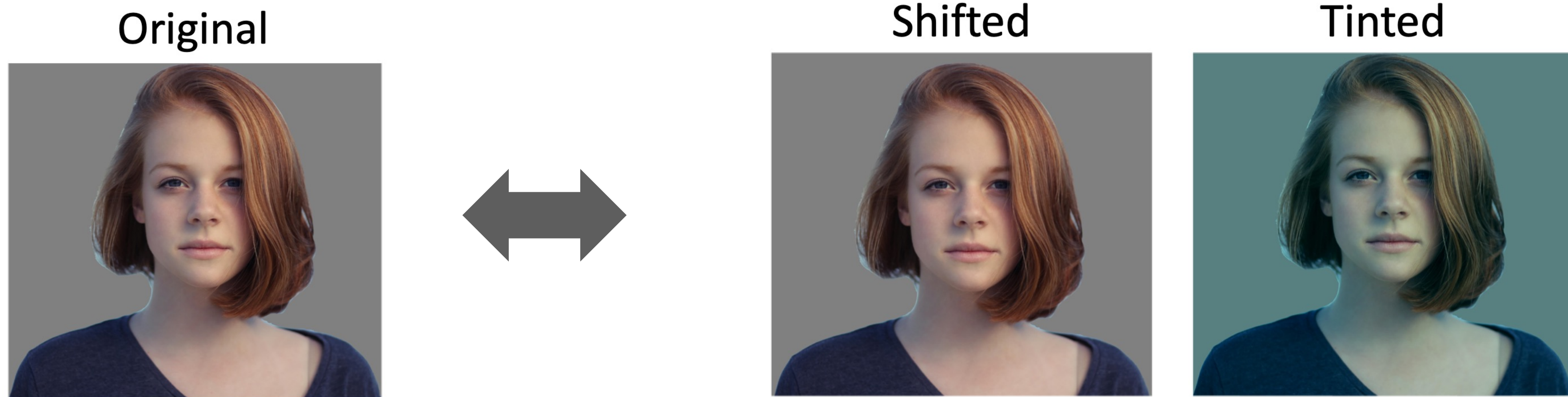
Dimensions = 3
Points = $4^3$

$$2^{32 X 32} \approx 10^{308}$$

# K-Nearest Neighbors: Seldomly Used on Raw Pixels

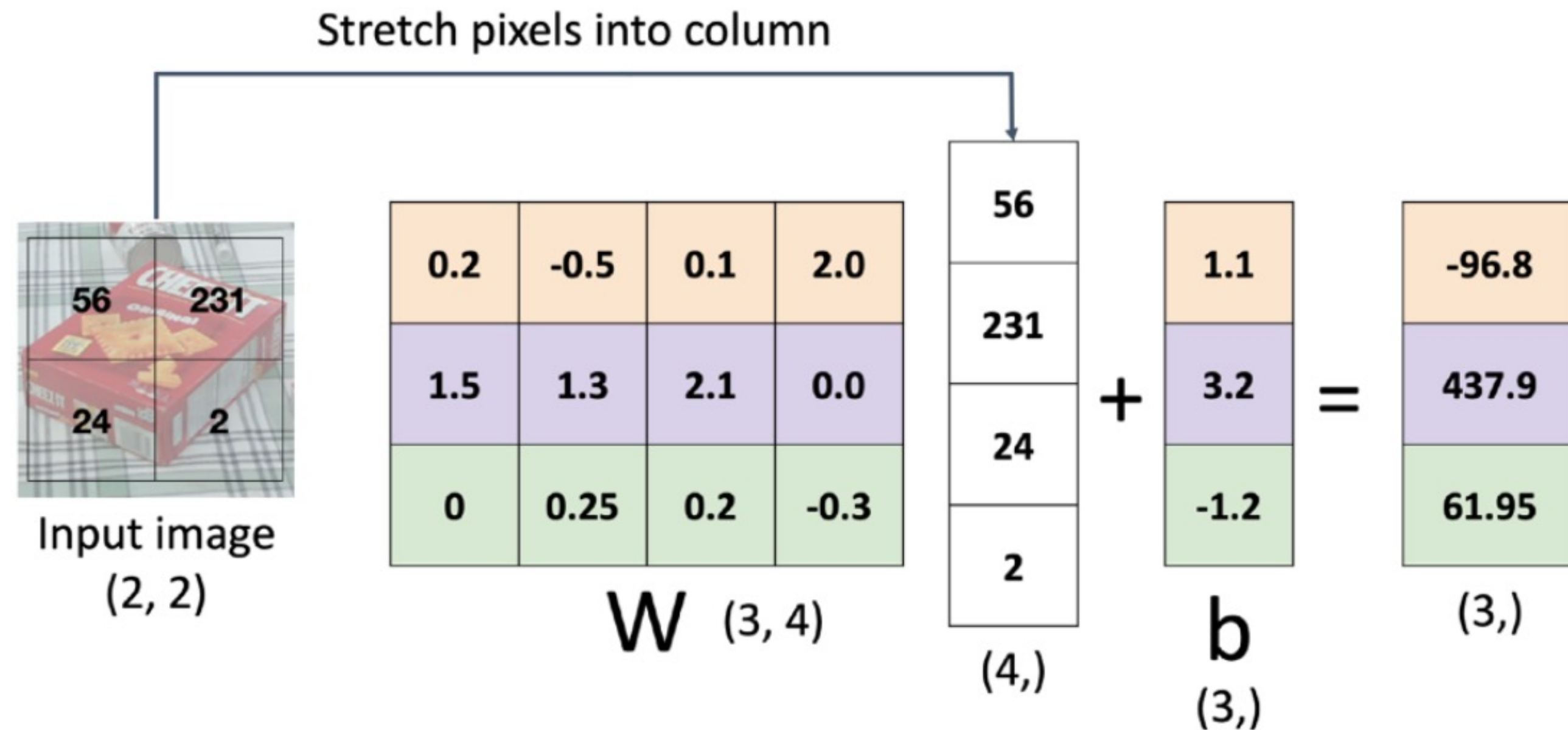Very slow at test time

Distance metrics on pixels are not informative



Original     Shifted     Tinted

Both images have same L2 distance to the original

Algebraic Viewpoint

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



Stretch pixels into column

| 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)

56  231
24  2

W (3, 4)

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

(4,)

$+$

| 1.1 |
| 3.2 |
| -1.2 |

b (3,)

$=$

| -96.8 |
| 437.9 |
| 61.95 |

(3,)

# Recap: Linear Classifier

## Visual Viewpoint

"stretch rows of W into images"

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 |
|-----|---|-----|---|------|

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

master chef can · cracker box · sugar box · tomato soup can · mustard bottle · fish can

gelatin box · meat can · mug · large marker

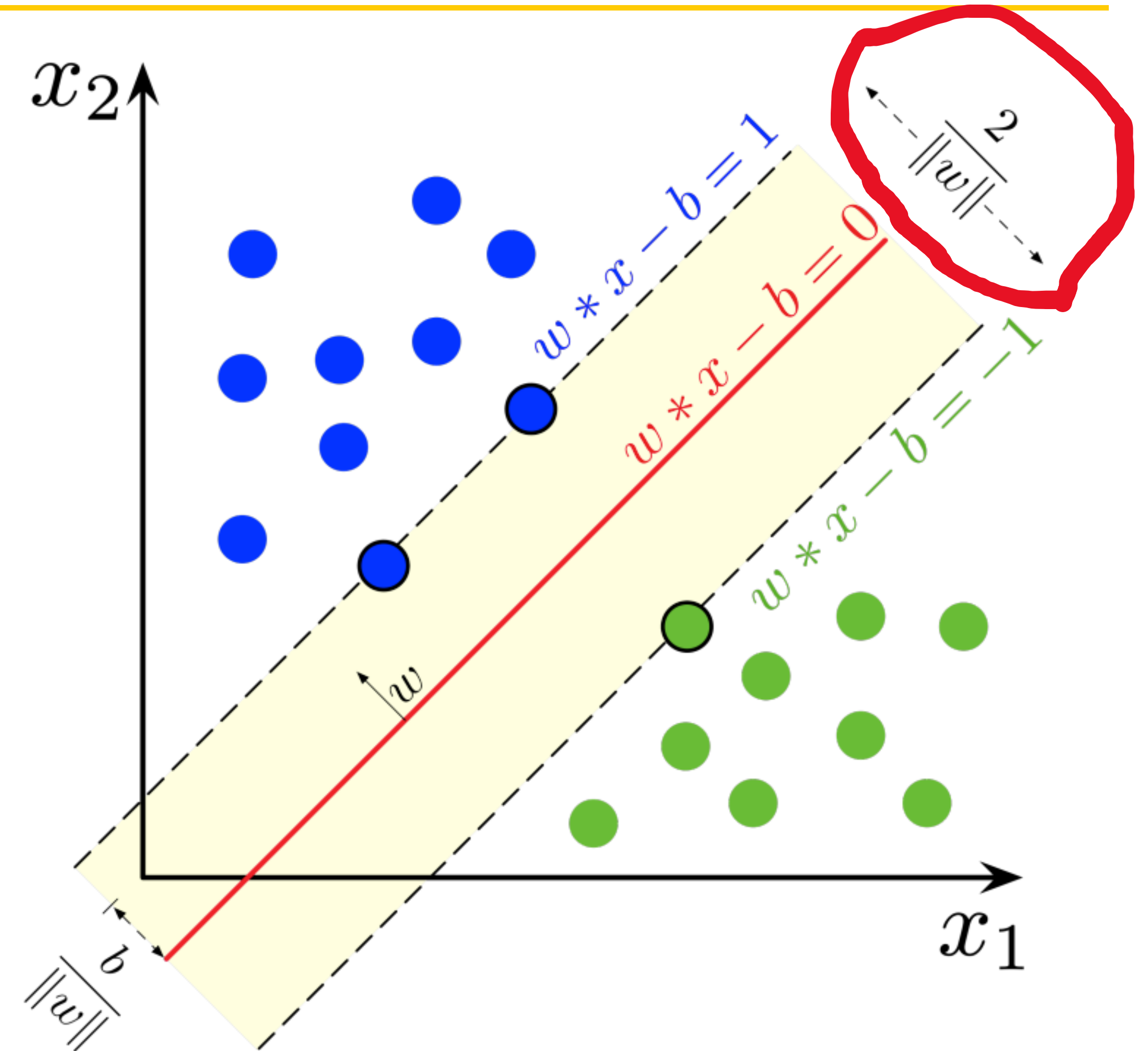# Recap: Linear Classifier

Geometric Viewpoint

Training Data

$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \ldots, (\mathbf{x}_n, y_n)$$

Hyperplane
$$\mathbf{w}^{\mathrm{T}}\mathbf{x} - b = 0$$

# What if there are misclassifications?

Hinge Loss  (soft margin)

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

**Overfitting**

# Back to SVM...

Training Data

$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \ldots, (\mathbf{x}_n, y_n)$$

Hyperplane

$$\mathbf{w}^T \mathbf{x} - b = 0$$

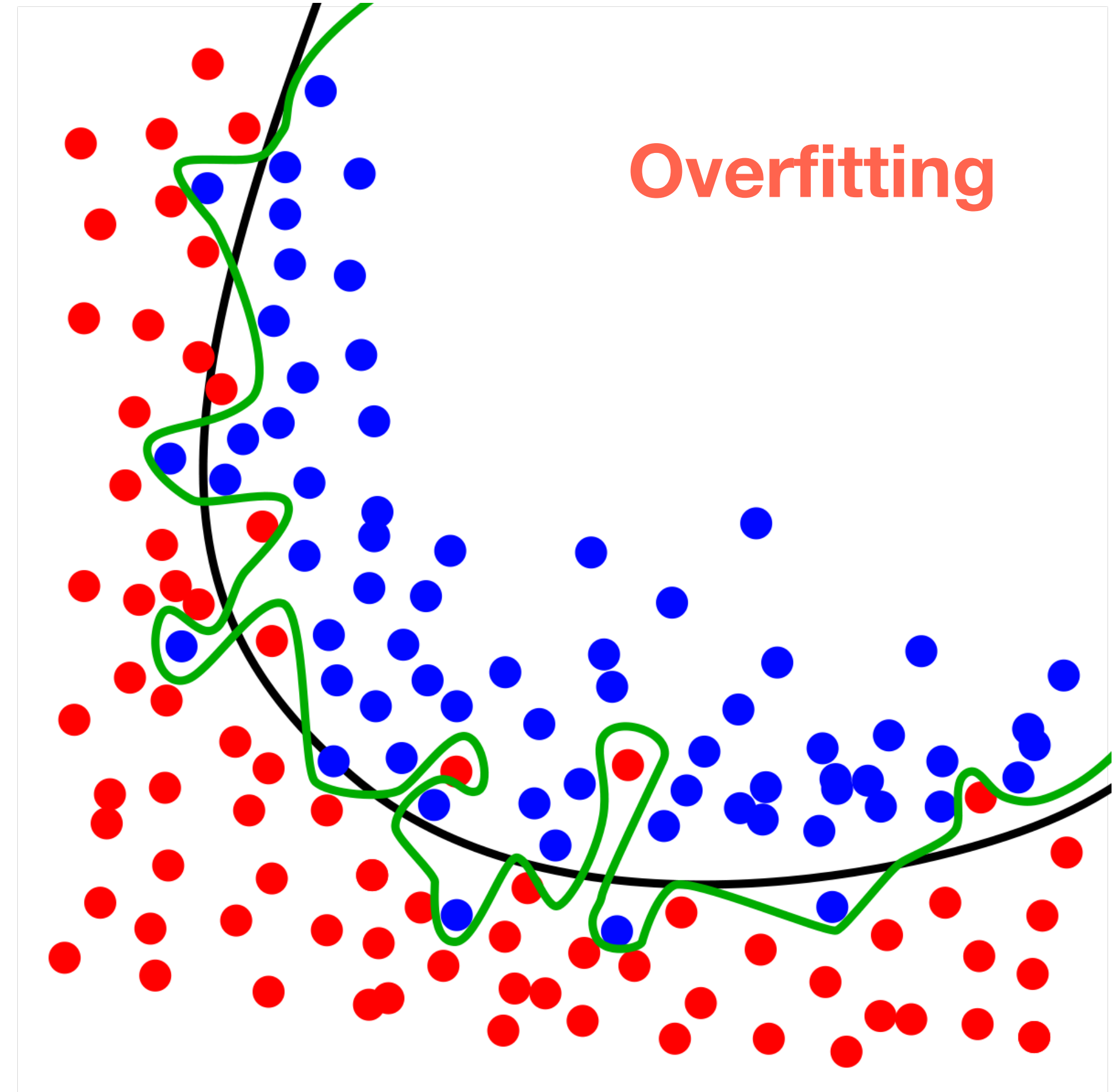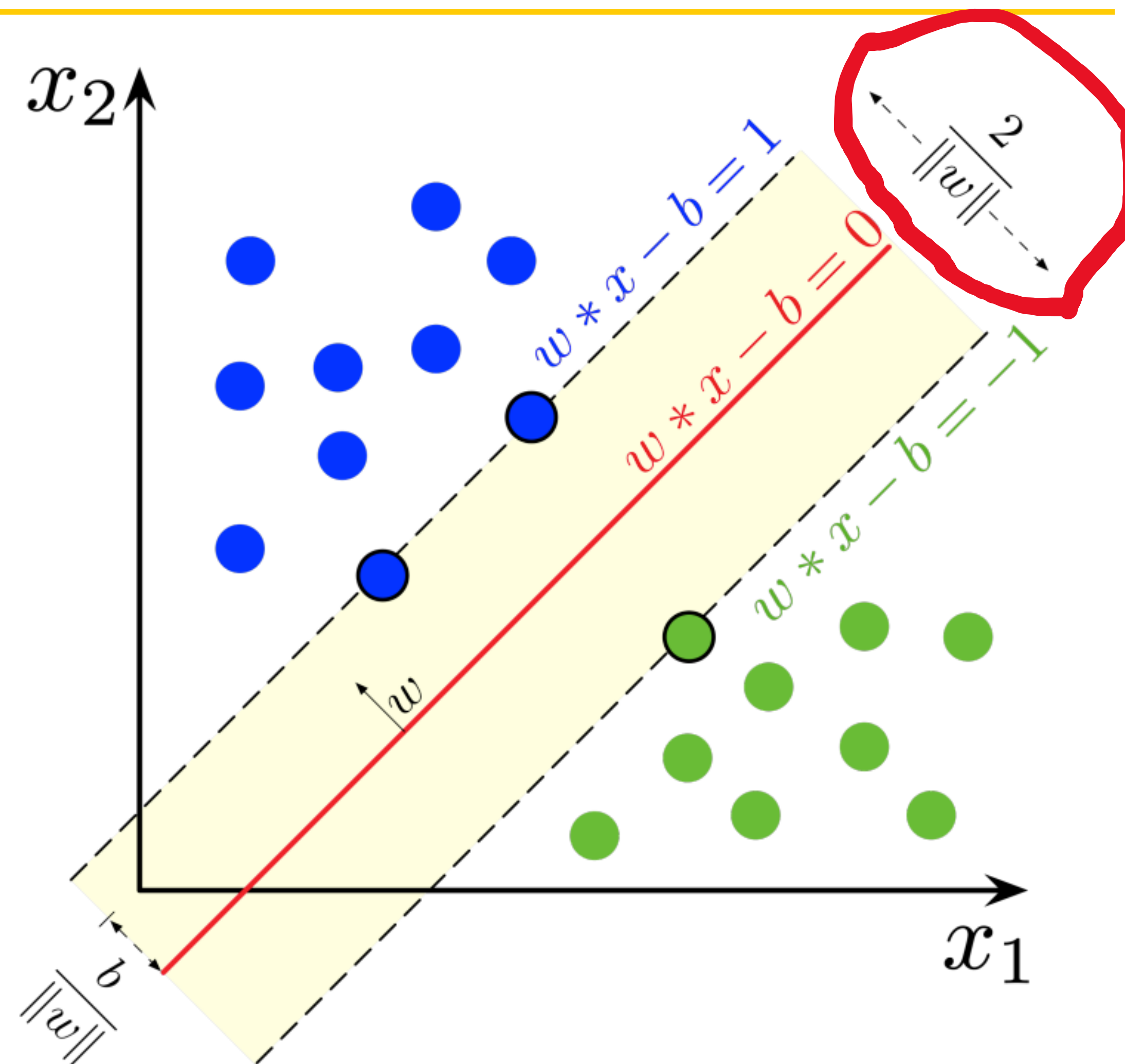Maximize $\dfrac{2}{\|w\|}$ ➡ Minimize $\dfrac{\|w\|}{2}$

# Loss Functions Quantify Preferences

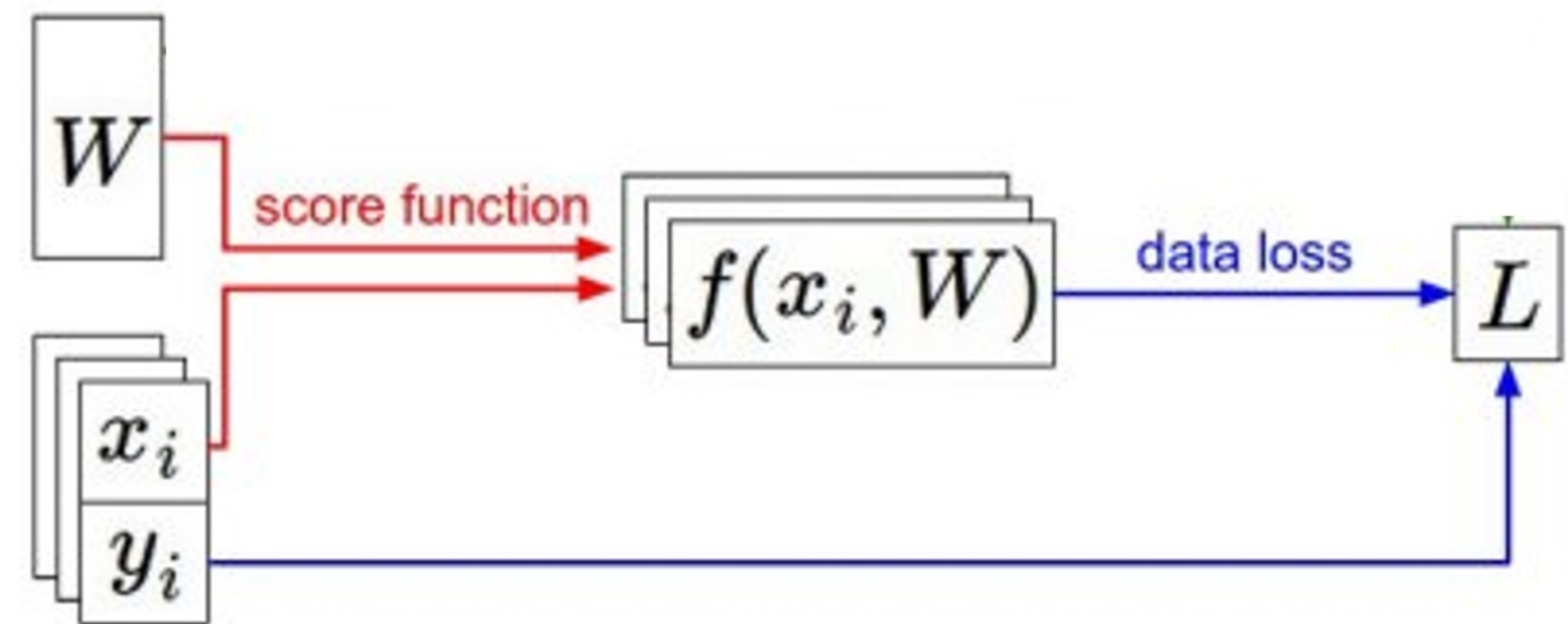**Q: How do we find the best W,b?**

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

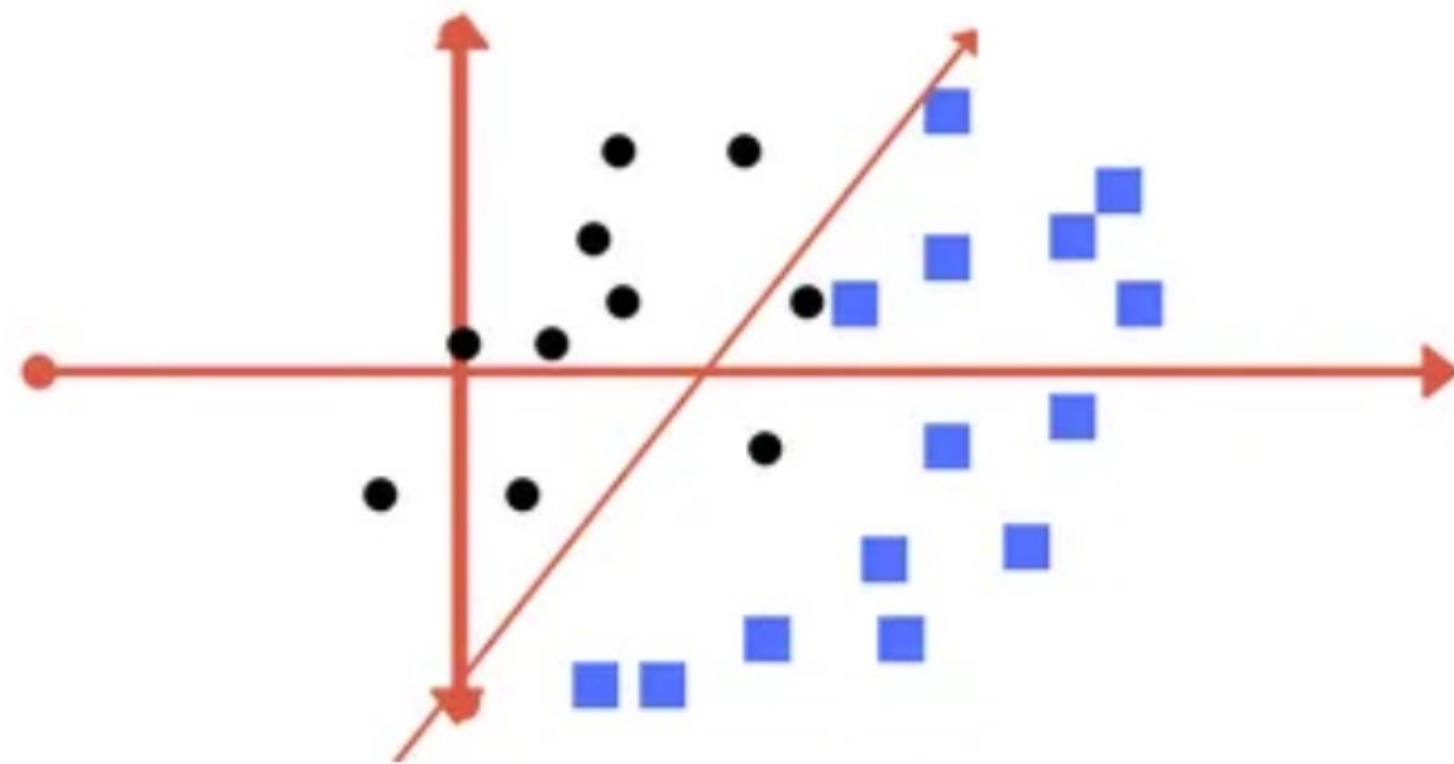SVM: $L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$

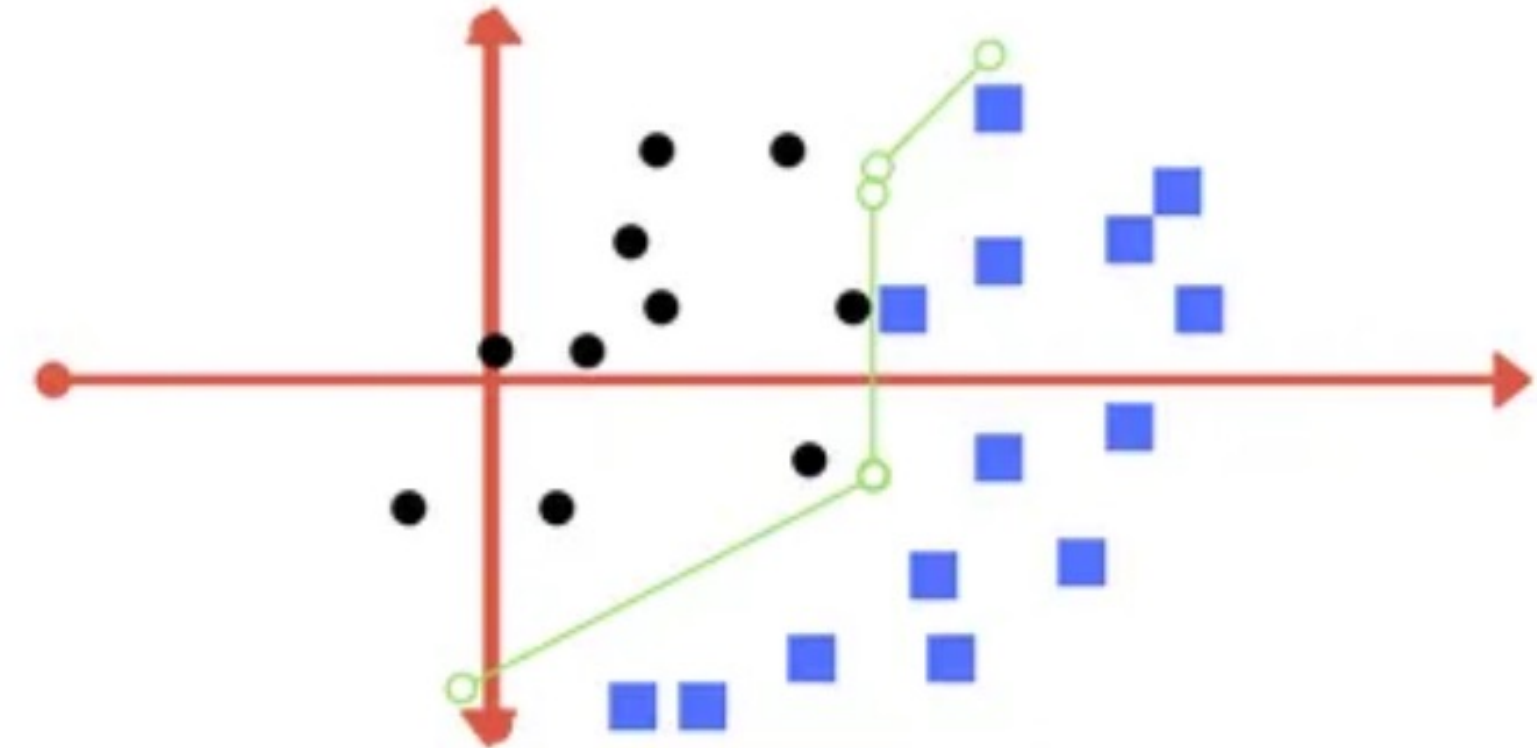$$+\lambda \frac{\|w\|}{2}$$

# Loss Functions Quantify Preferences



Q: Low or High regularization?



Q: Low or High regularization?

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

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

$+\mathrm{C}\dfrac{\|w\|^{\wedge}2}{2}$

# General Case: Adding Regularization Term

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

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

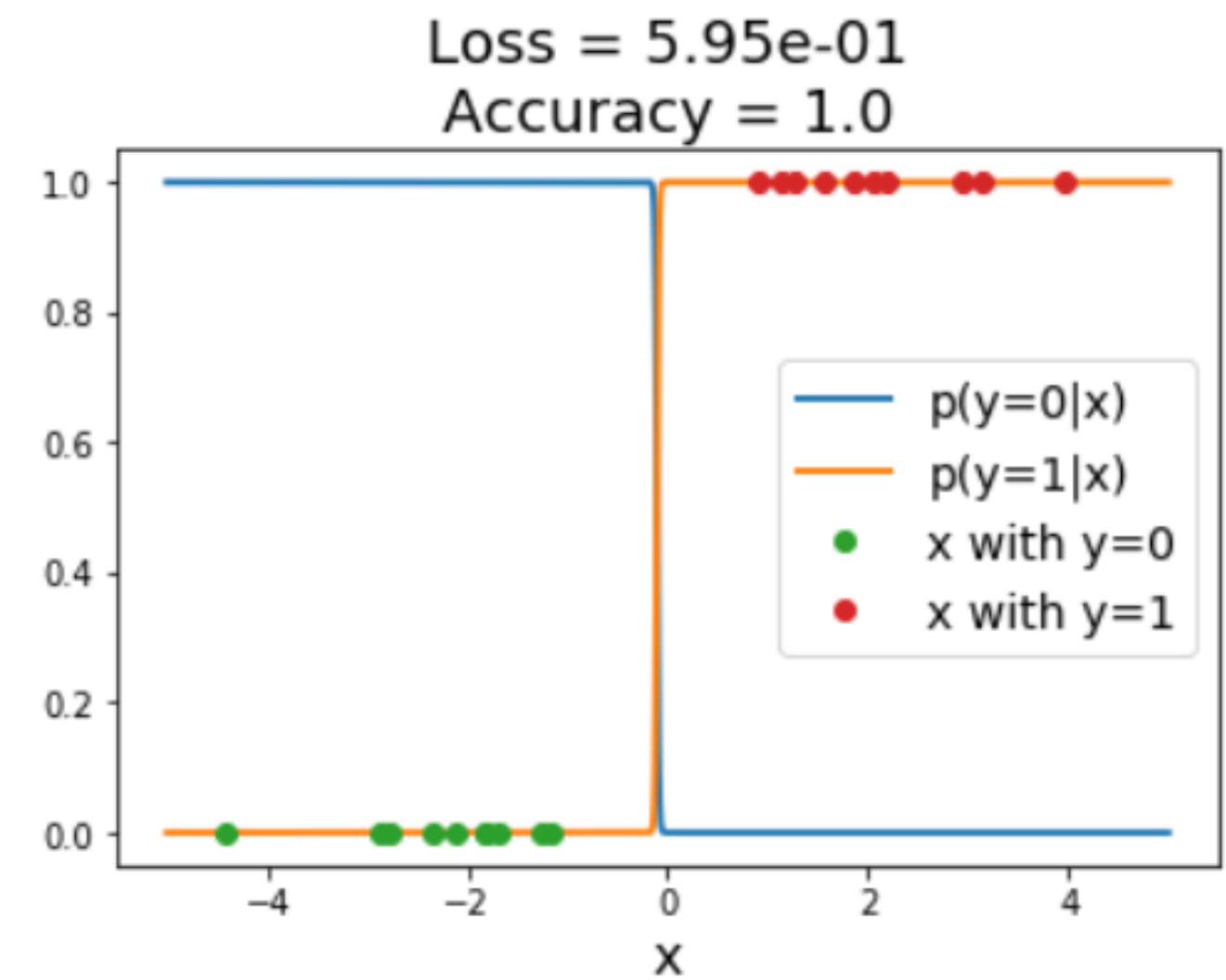L1 regularization: $R(W) = \sum_{k,l} |W_{k,l}|$
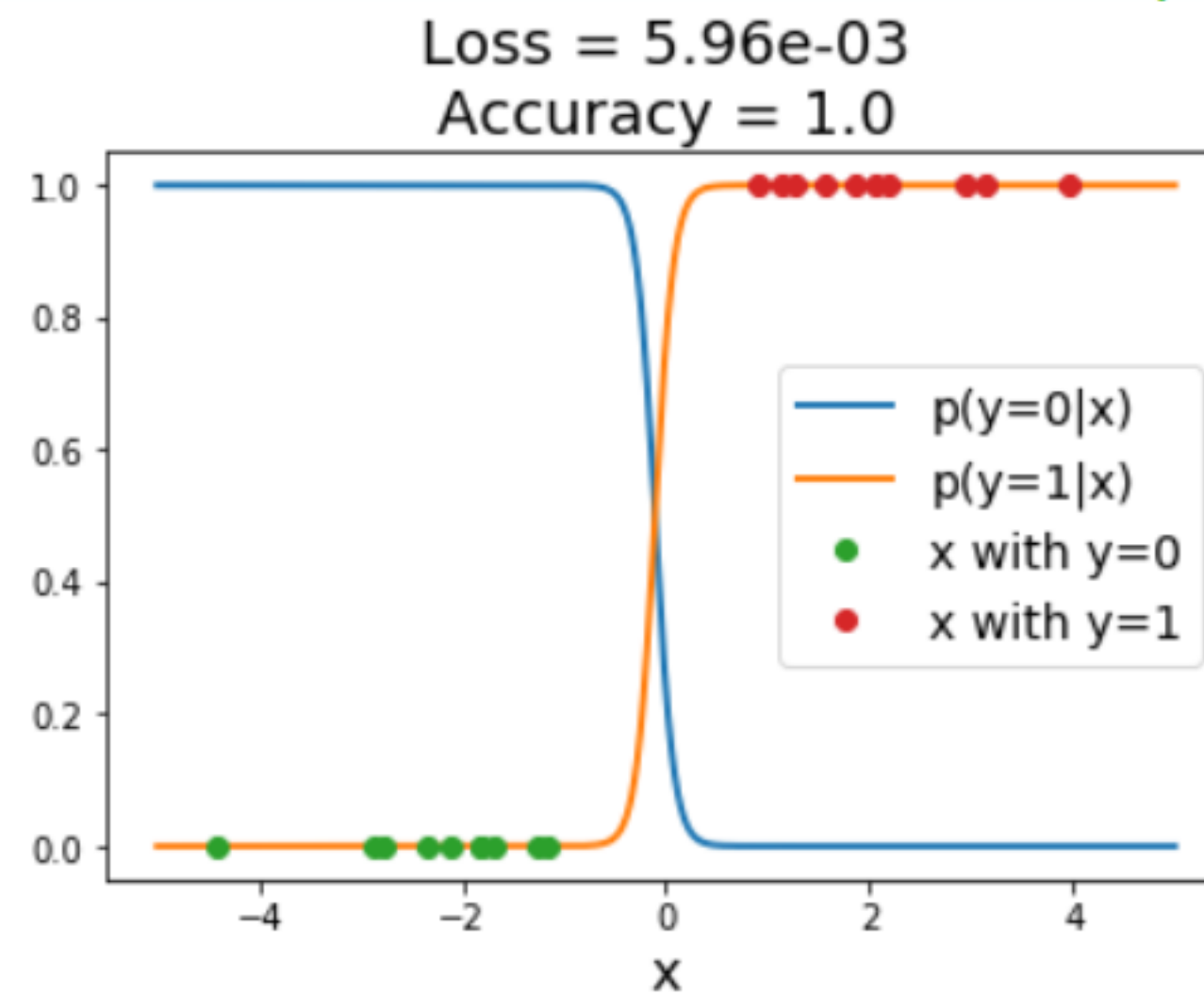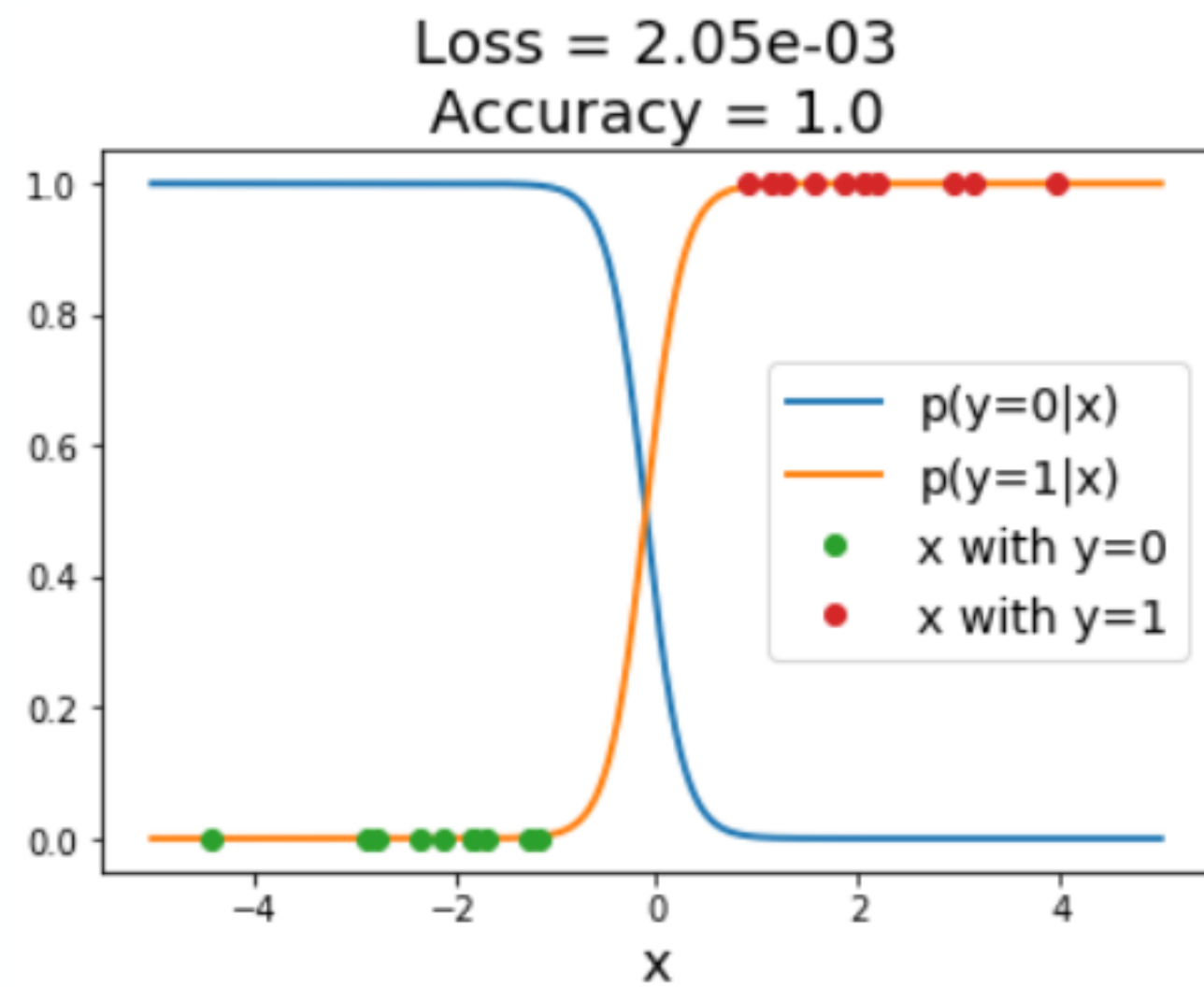
# Regularization: Example

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



Loss = 2.05e-03
Accuracy = 1.0

Loss = 5.96e-03
Accuracy = 1.0

Loss = 5.95e-01
Accuracy = 1.0

# Regularization: Expressing Preference

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

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

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

$$w_1^T x = w_2^T x = 1$$

L2 Regularization

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

L2 Regularization prefers weights to be "spread out"

Same predictions, so data loss will always be the same

# How to find 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

Optimization

$$w* = \arg \min_{w} L(w)$$

# 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
# ... (trunctated: continues for 1000 lines)
```

# 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!
not bad but not great…  (SOTA is ~95%)

# Idea #2: Follow the slope

# Idea #2: Follow the slope

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**.

Current W:

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

Gradient $\dfrac{dL}{dW}$

[?,
?,
?,
?,
?,
?,
?,
?,
?, …]

Current **W**:

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

**W + h** (first dim):

[0.34 + **0.0001**,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33, …]
loss 1.25322

Gradient $\dfrac{dL}{dW}$

[-2.5,
?,
?,

(1.25322 - 1.25347)/
0.0001
= -2.5

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

?, …]

Current **W**:

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

**W + h** (second dim):

[0.34,
-1.11 + **0.0001**,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33, …]
loss 1.25353

Gradient $\dfrac{dL}{dW}$

[-2.5,
**0.6**,
?,
?,

(1.25353 - 1.25347)/
0.0001
= 0.6

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

29

# Computing Gradients

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

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

In practice: 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.
    """
```

# Computing Gradients

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

- **Numeric gradient:** approximate, slow, easy to write
- **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.
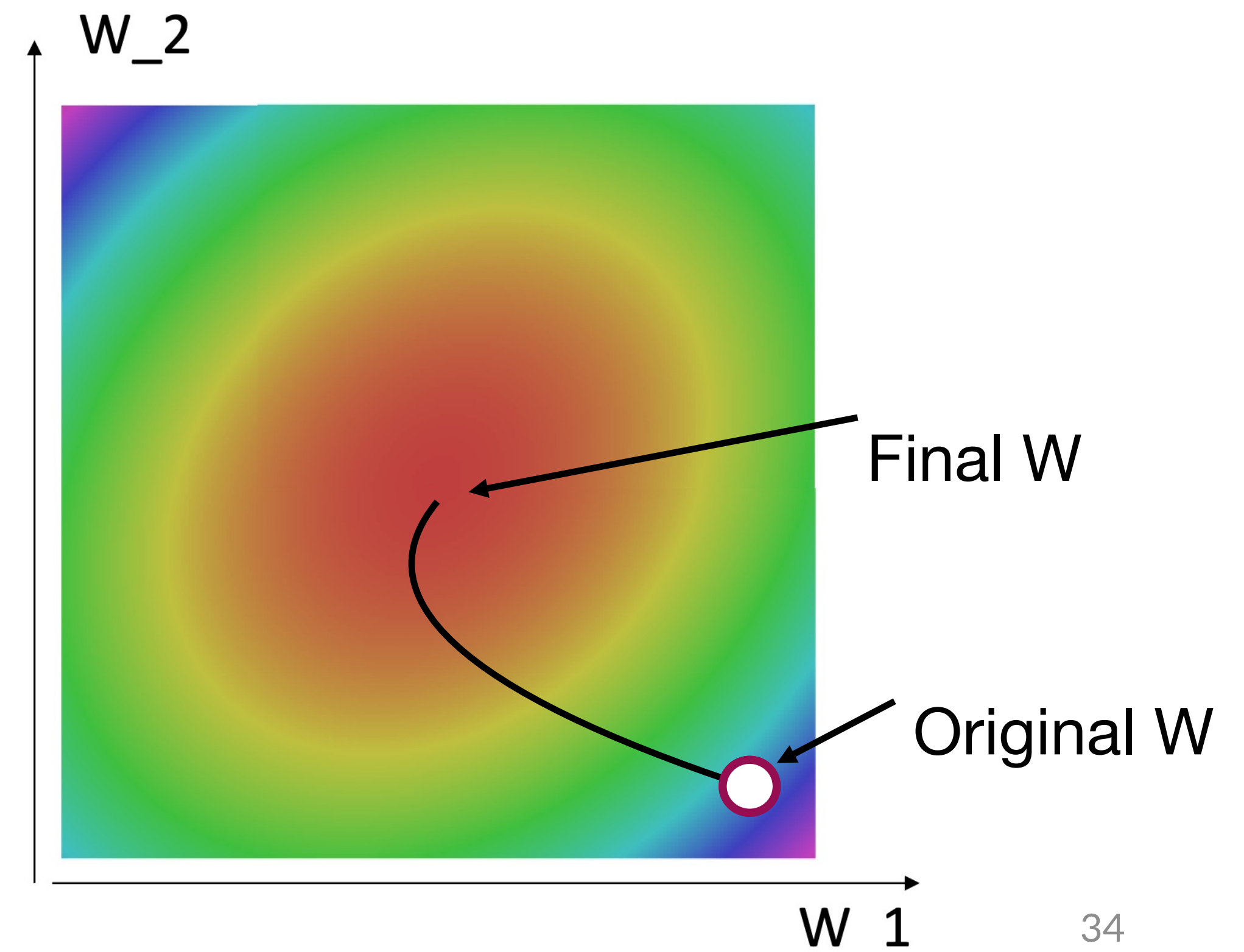
# Gradient Descent

- 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



Negative gradient direction

Original W

W_2

W_1

# 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



W_2

Final W

Original W

W_1

# Batch Gradient Descent

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

Full sum expensive when N is large!

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

# Stochastic Gradient Descent (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 **minibatch** of examples 32/64/128 common

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

```
# 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
```

# 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

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

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

# Interactive Web Demo

# Problems with SGD

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
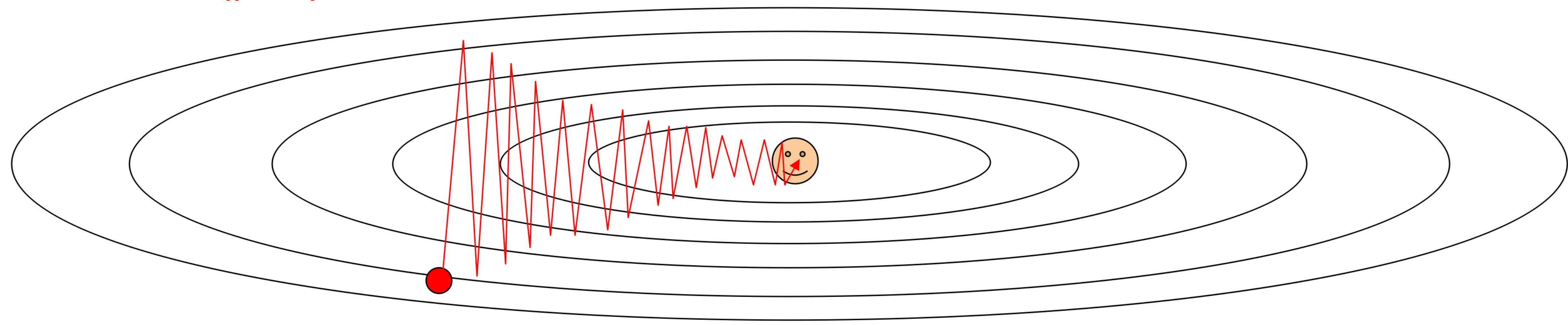
# Problems with SGD

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

What does gradient decent do?

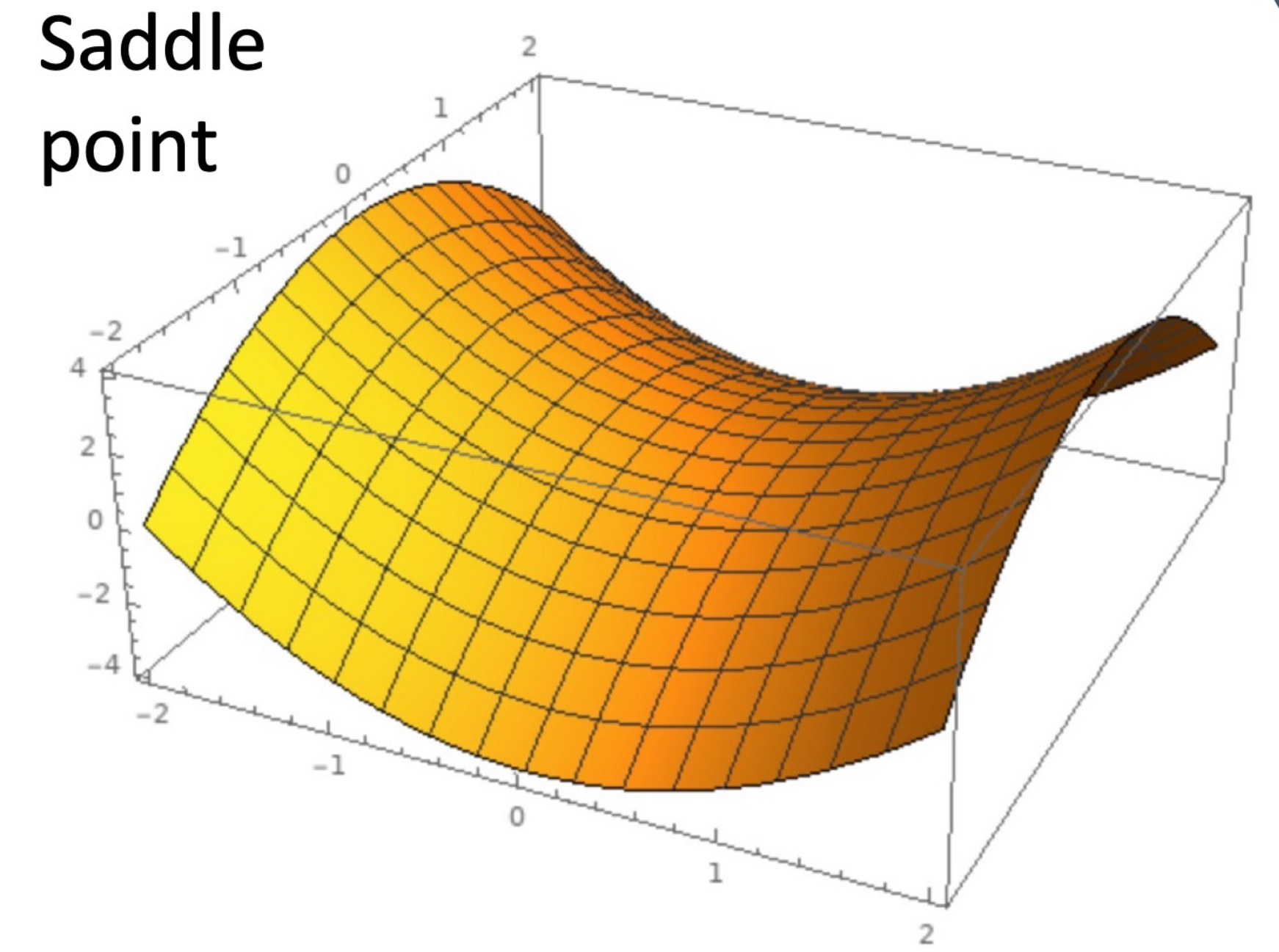Very slow progress along shallow dimension, jitter along steep

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

# Problems with SGD

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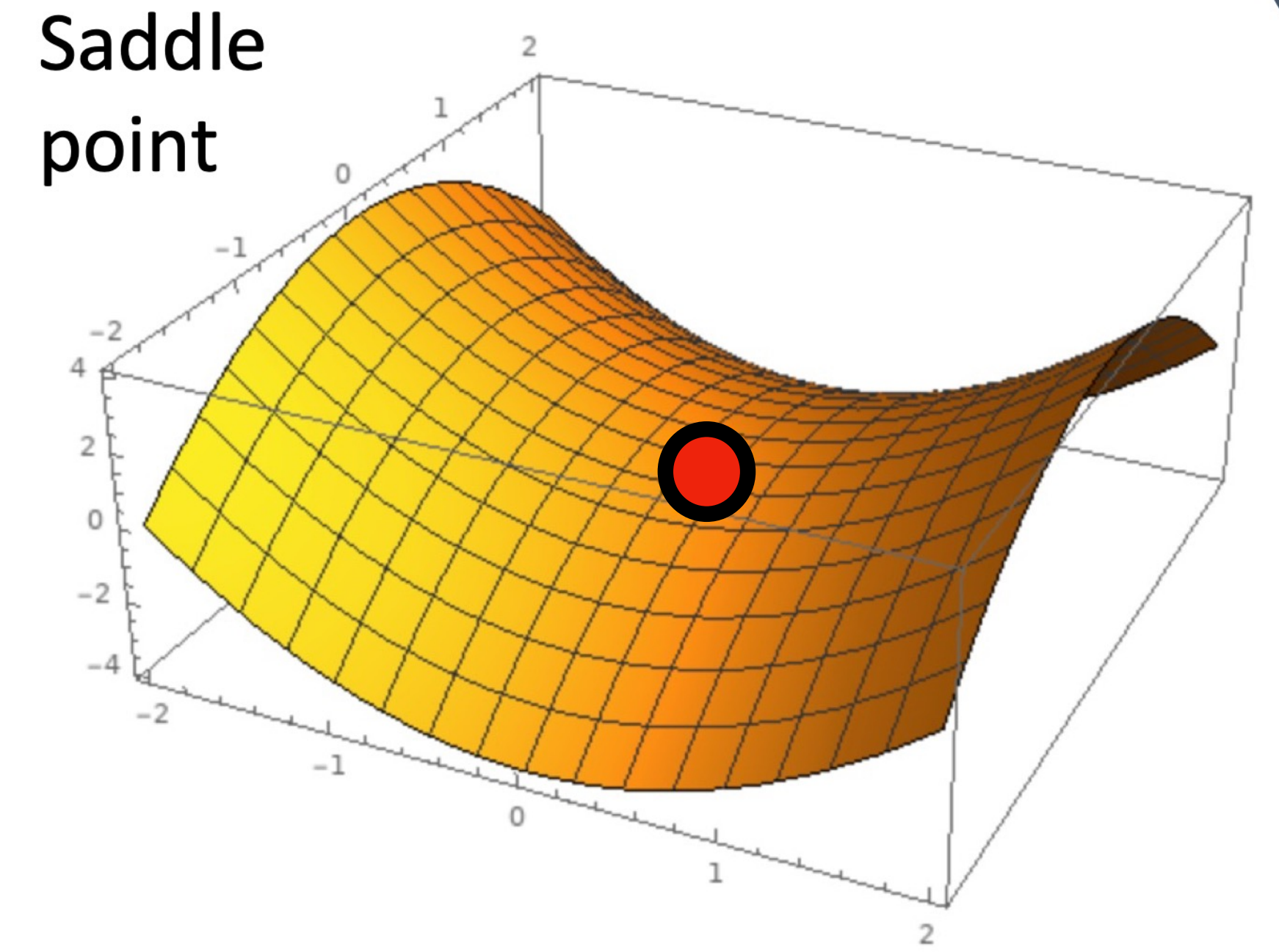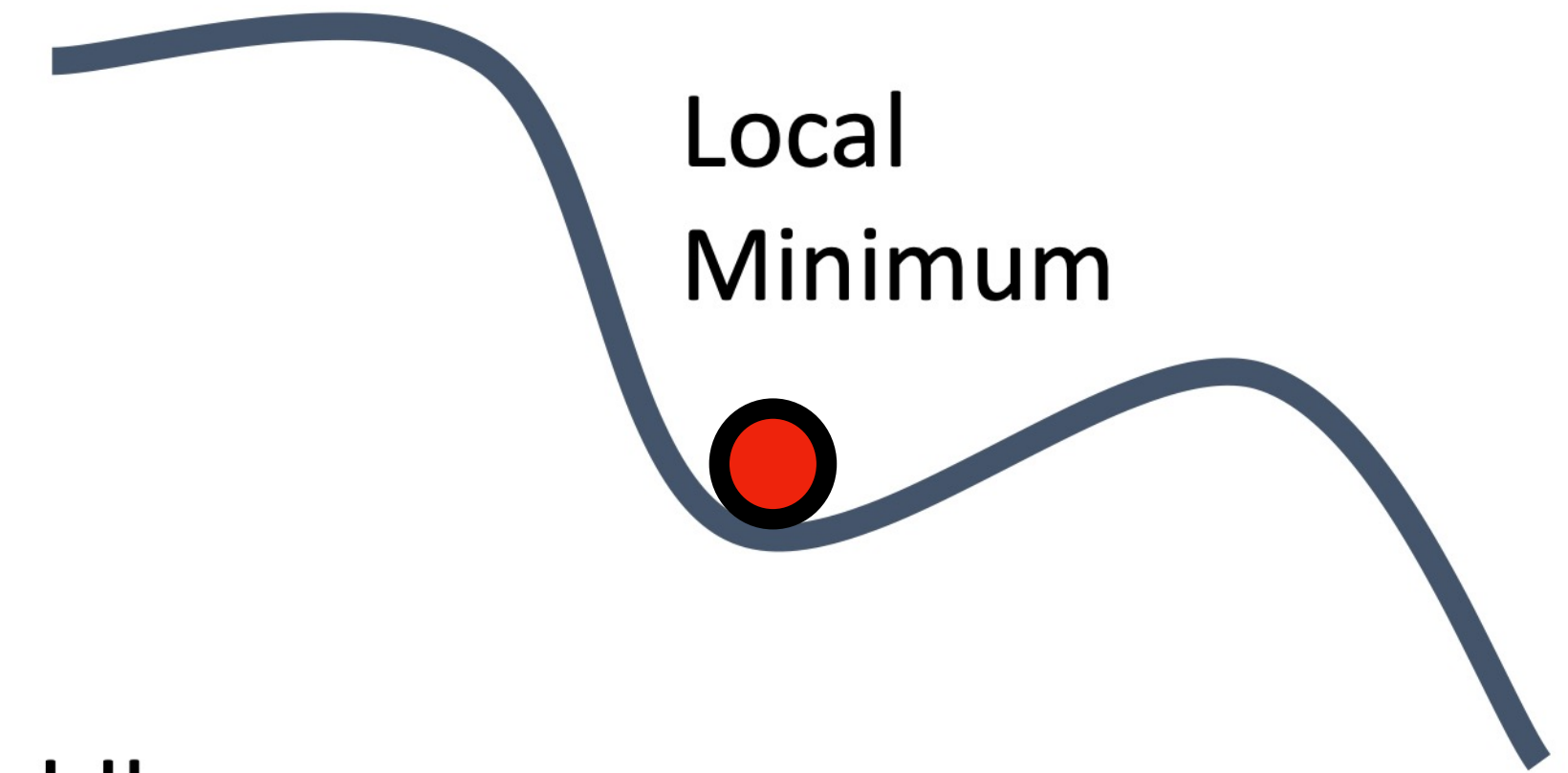**local minimum** or **saddle point**?

Local
Minimum

Saddle
point

# Problems with SGD

What if the loss function has a
**local minimum** or **saddle point**?
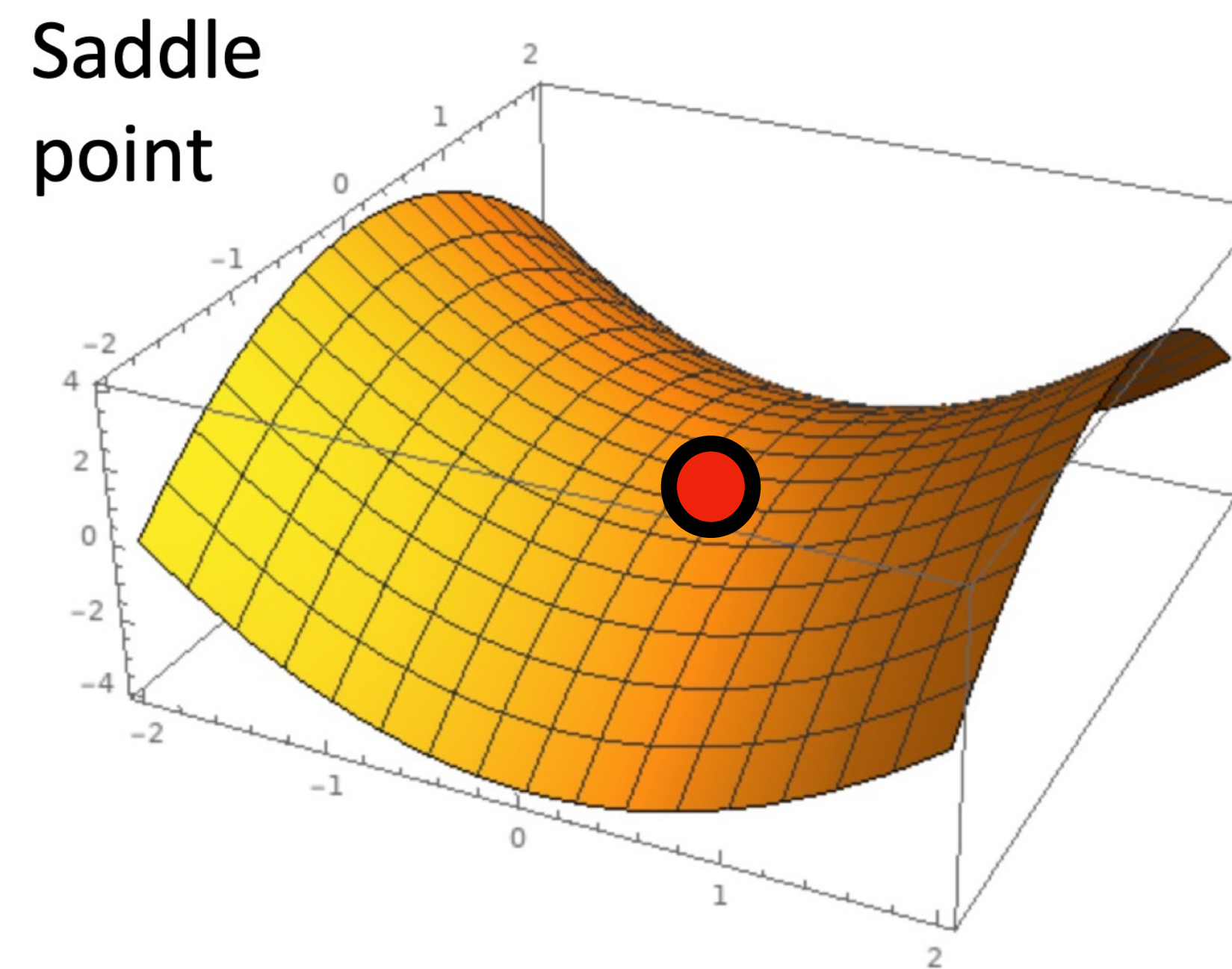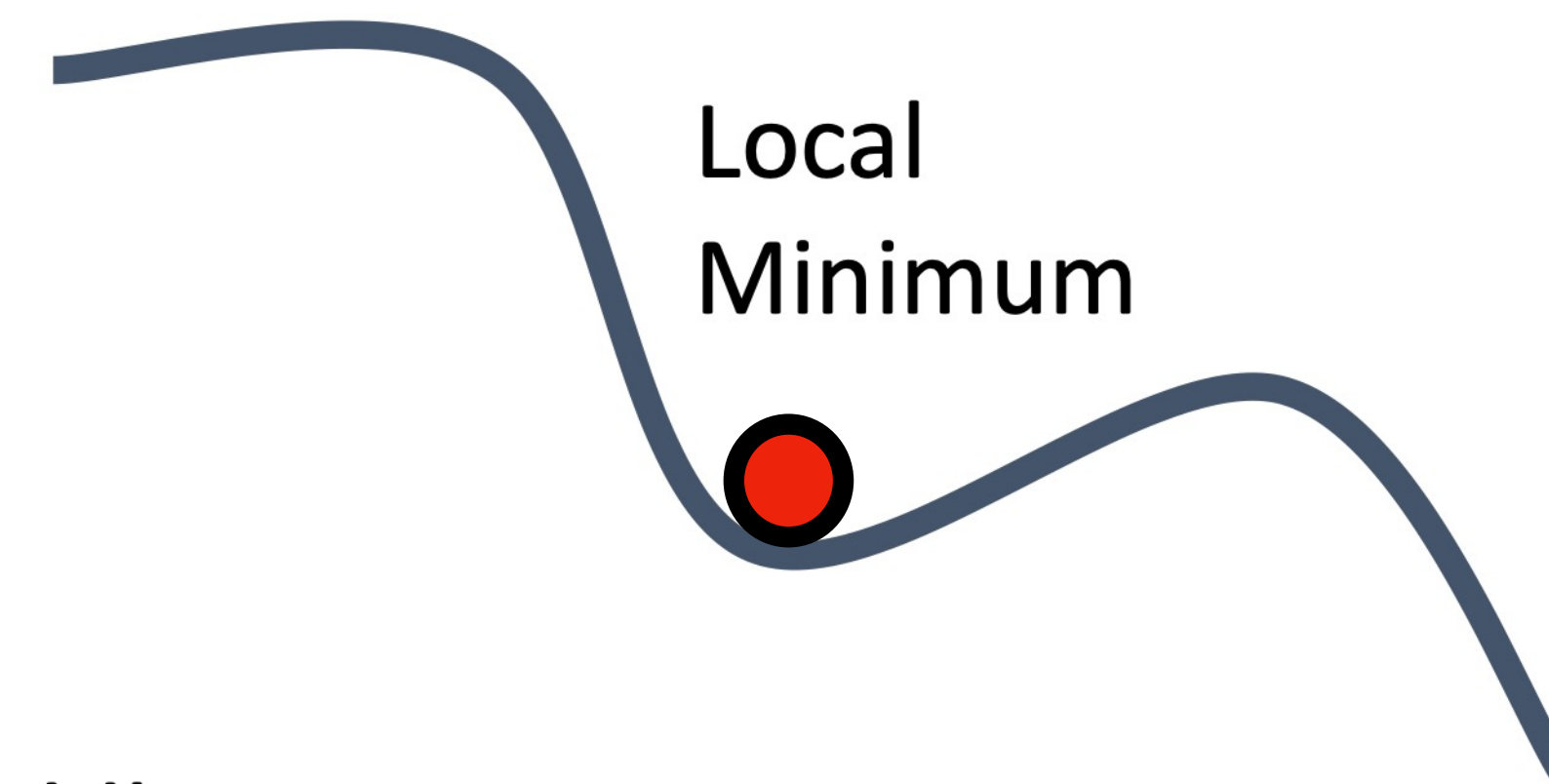
Zero gradient, gradient descent gets stuck

Local
Minimum

Saddle
point

# Problems with SGD

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

Local
Minimum

Saddle
point

<span style="color:red">Batched gradient descent always
computes same gradients</span>

<span style="color:red">SGD computes noisy gradients,
may help to escape saddle points</span>

# 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

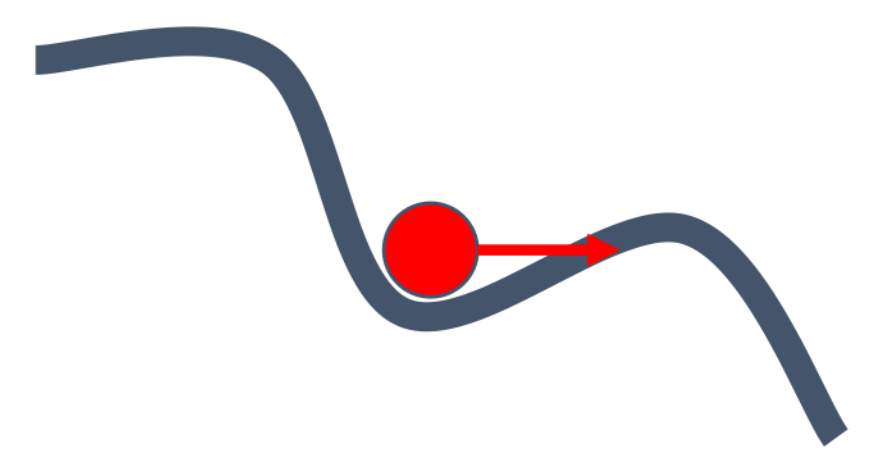$$v_{t+1} = \rho v_t + \nabla L(w_t)$$
$$w_{t+1} = w_t - \alpha v_{t+1}$$

```
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
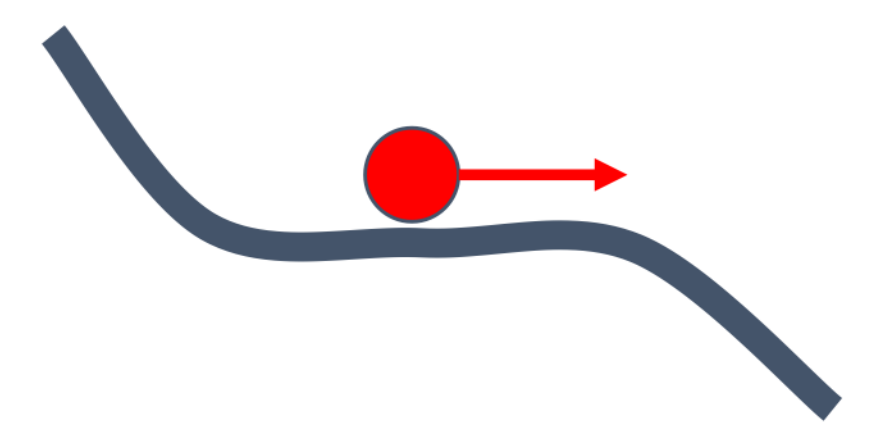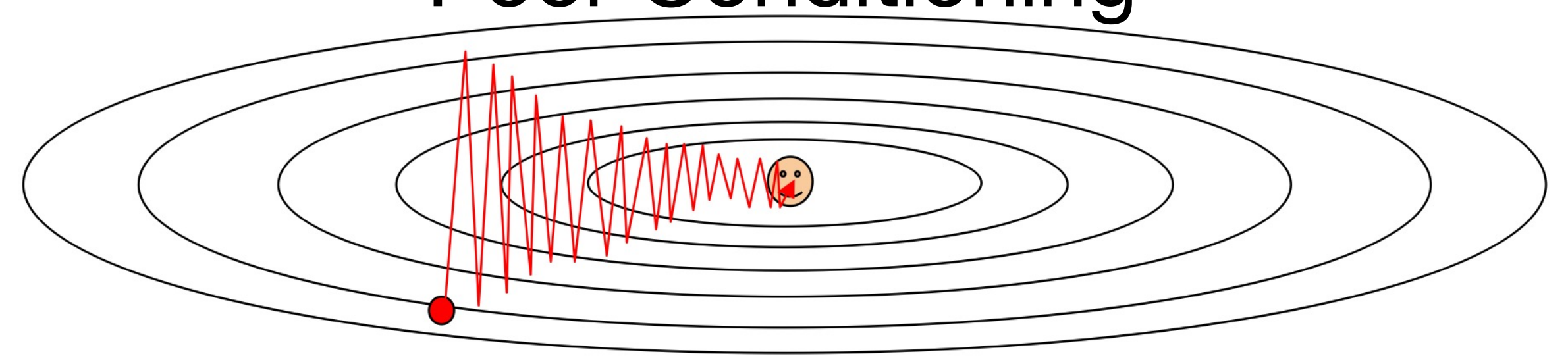- Rho gives "friction"; typically rho = 0.9 or 0.99
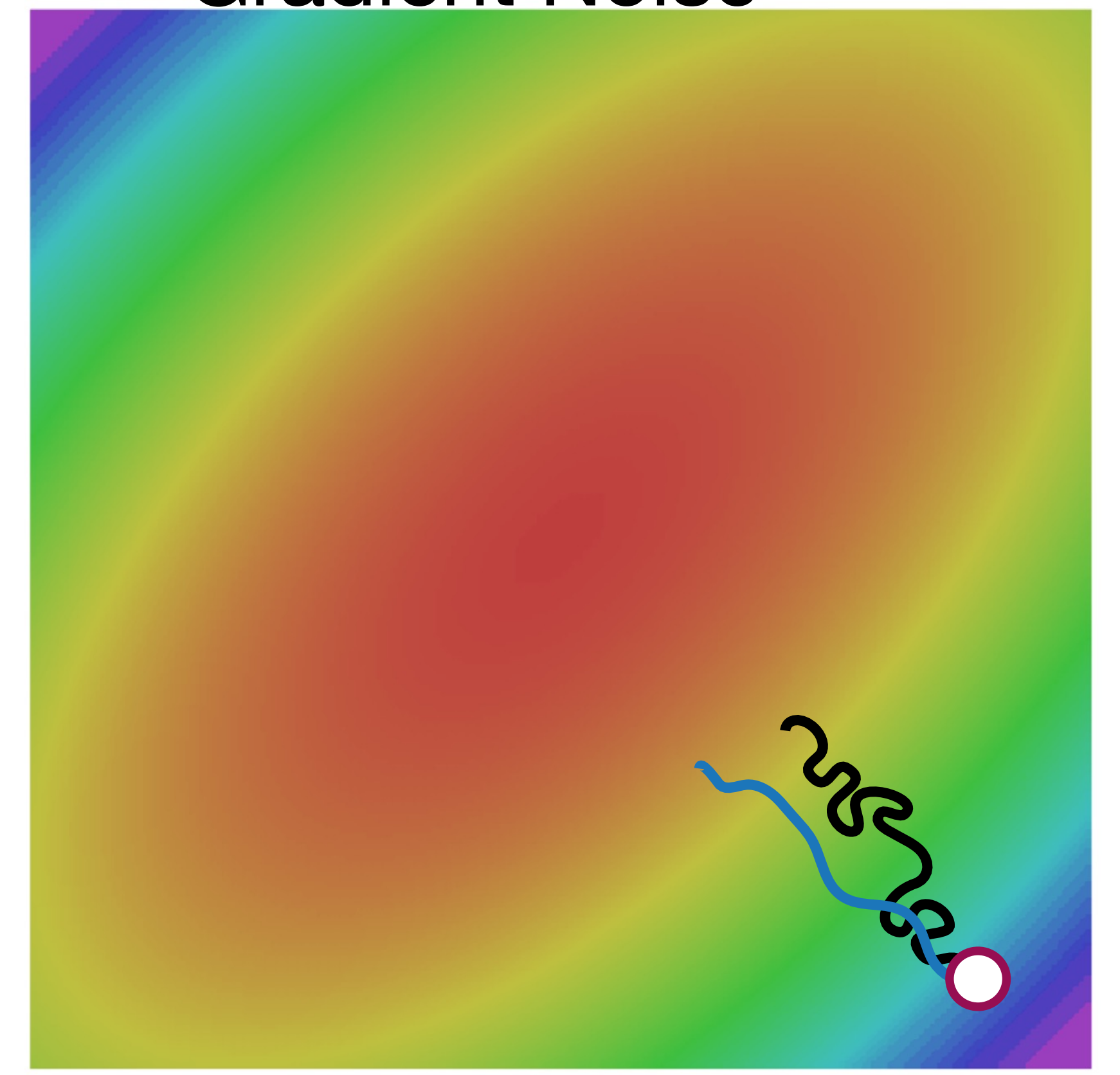
# SGD + Momentum

## Local Minima

## Saddle Points

## Poor Conditioning

## Gradient Noise



———— SGD    ———— SGD+Momentum

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

# SGD + Momentum

## Momentum update:



**Velocity**

**Actual step**

**Gradient**

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

## Nesterov Momentum



**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

# DeepRob

Lecture 3
Regularization + Optimization
University of Michigan | Department of Robotics