

Stochastic Gradient Descent

Spring 2020

Image Classification



08	02	22	97	38	15	00	40	00	75	04	05	07	78	52	12	50	77	31	85
49	49	99	40	17	81	18	57	60	87	17	40	98	43	69	45	04	56	62	00
81	49	31	73	55	79	14	29	93	71	40	67	57	83	30	03	49	13	36	65
52	70	95	23	04	60	11	42	69	44	68	56	01	32	56	71	37	02	36	91
22	31	16	71	51	67	82	89	41	92	36	54	22	40	40	28	66	33	13	80
24	47	38	80	99	03	45	02	44	75	33	53	78	36	84	20	35	17	12	50
32	98	81	28	64	23	67	10	26	38	40	67	59	54	70	66	18	38	64	70
67	26	20	68	02	62	12	20	95	63	94	39	63	08	40	91	66	49	94	21
24	55	58	05	66	73	99	26	97	17	78	78	96	83	14	88	34	89	63	72
21	36	23	09	75	00	76	44	20	45	35	14	00	61	33	97	34	31	33	95
78	17	53	28	22	75	31	67	15	94	03	80	04	62	16	14	09	53	56	92
16	39	05	42	96	35	31	47	55	58	88	24	00	17	54	24	36	29	85	57
86	56	00	48	35	71	89	07	05	44	44	37	44	60	21	58	51	54	17	58
19	80	81	68	05	94	47	69	28	73	92	13	86	52	17	77	04	89	55	40
04	52	08	83	97	35	99	16	07	97	57	32	16	26	26	79	33	27	98	66
59	16	68	87	57	62	20	72	03	46	33	67	46	55	12	32	63	93	53	69
04	42	16	73	38	85	39	11	24	94	72	18	08	46	29	32	40	62	76	36
20	69	36	41	72	30	23	88	34	88	99	69	82	67	59	85	74	04	36	16
20	73	35	29	78	31	90	01	74	31	49	71	48	84	81	16	23	57	05	54
01	70	54	71	83	51	54	69	16	92	33	48	61	43	52	01	89	29	47	48

What the computer sees

image classification →
82% cat
15% dog
2% hat
1% mug

Linear model

- Score function
 - Maps raw data to class scores
- Loss function
 - Measures how well predicted classes agree with ground truth labels
 - Multiclass Support Vector Machine loss (SVM loss)
 - Softmax classifier (cross-entropy loss)
- Learning
 - Find parameters of score function that minimize loss function
 - Multiclass Support Vector Machine loss (SVM loss)

Recall: Linear model with SVM loss

- Score function

- Maps raw data to class scores

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

- Loss function

- Measures how well predicted classes agree with ground truth labels

$$L = \frac{1}{N} \sum_i \sum_{j \neq y_i} [\max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + \Delta)] + \lambda \sum_k \sum_l W_{k,l}^2$$

Today

- Learning model parameters with Stochastic Gradient Descent that minimize loss
- Later
 - Different score functions: deep networks
 - Same loss functions and learning algorithm

Outline

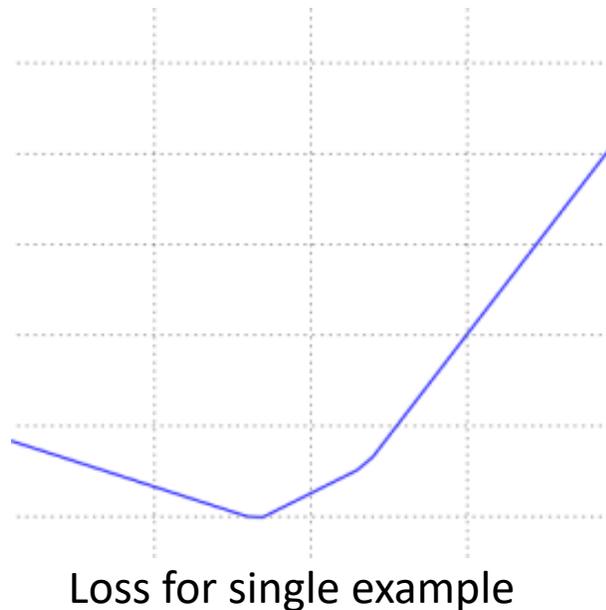
- Visualizing the loss function
- Optimization
 - Random search
 - Random local search
 - Gradient descent
 - Mini-batch gradient descent

Visualizing SVM loss function

- Difficult to visualize fully
 - CIFAR-10 a linear classifier weight matrix is of size $[10 \times 3073]$ for a total of 30,730 parameters
- Can gain intuition by visualizing along rays (1 dimension) or planes (2 dimensions)

Visualizing in 1-D

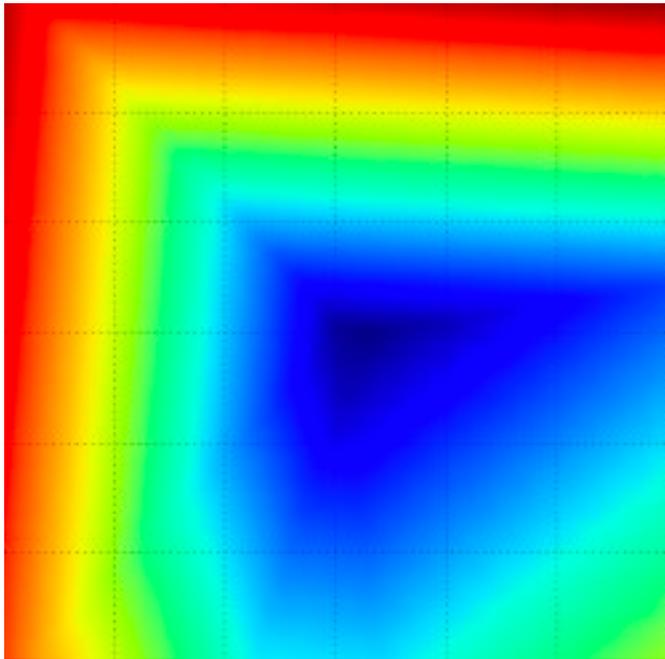
- Generate random weight matrix W
- Generate random direction W_1
- Compute loss along this direction $L(W + aW_1)$



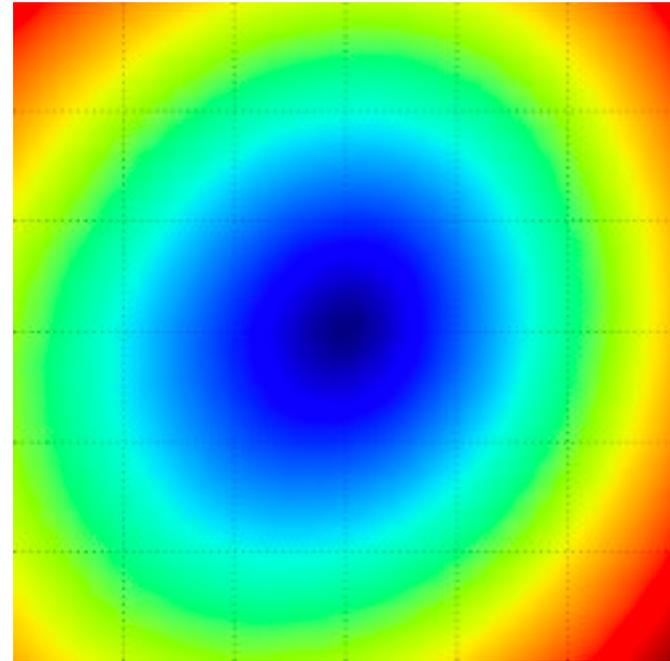
Where is the minima?

Visualizing in 2-D

- Compute loss along plane $L(W + aW_1 + bW_2)$



Loss for single example



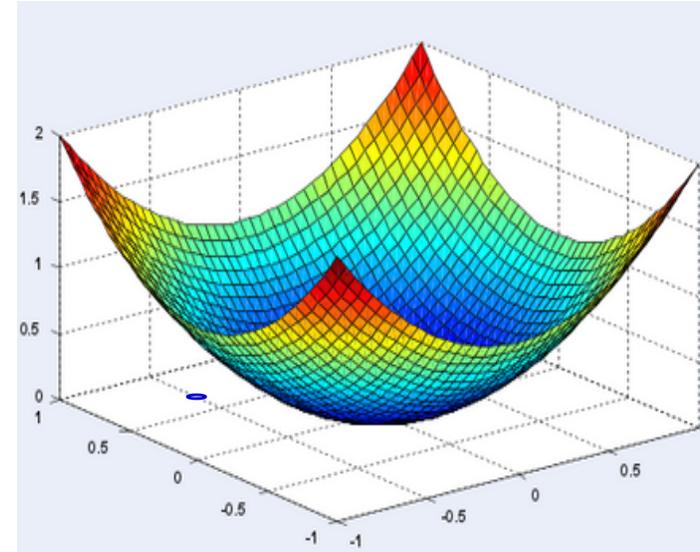
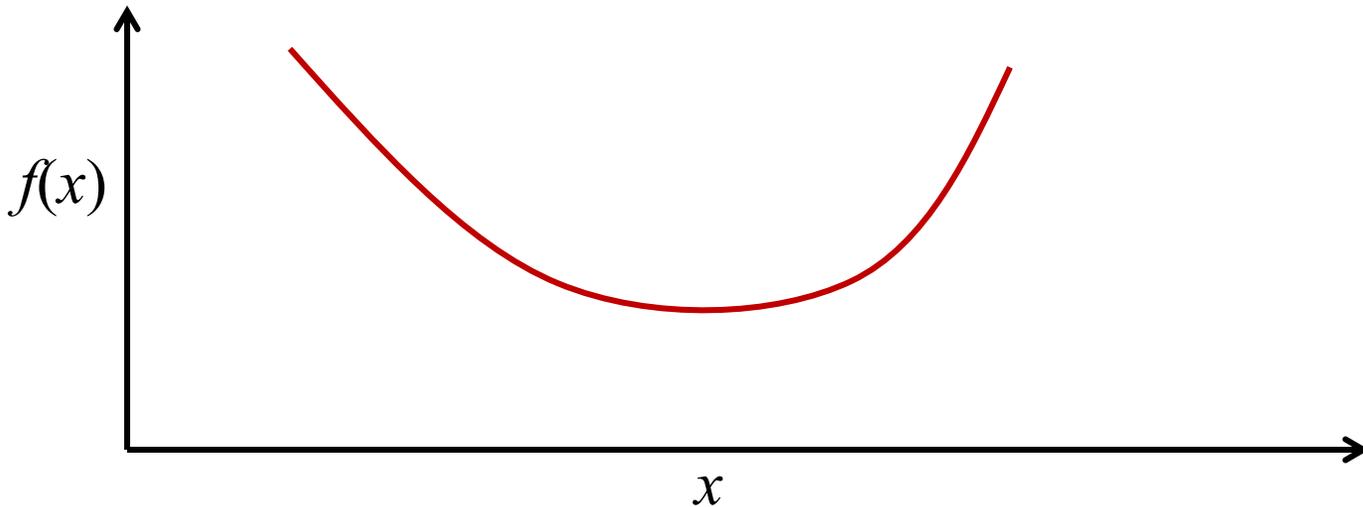
Average loss for 100 examples
(convex function)

How do we find weights that minimize loss?

- Random search
 - Try many random weight matrices and pick the best one
 - Performance: poor
- Random local search
 - Start with random weight matrix
 - Try many local perturbations, pick the best one, and iterate
 - Performance: better but still quite poor
- Useful idea: iterative refinement of weight matrix

Optimization basics

The problem of optimization



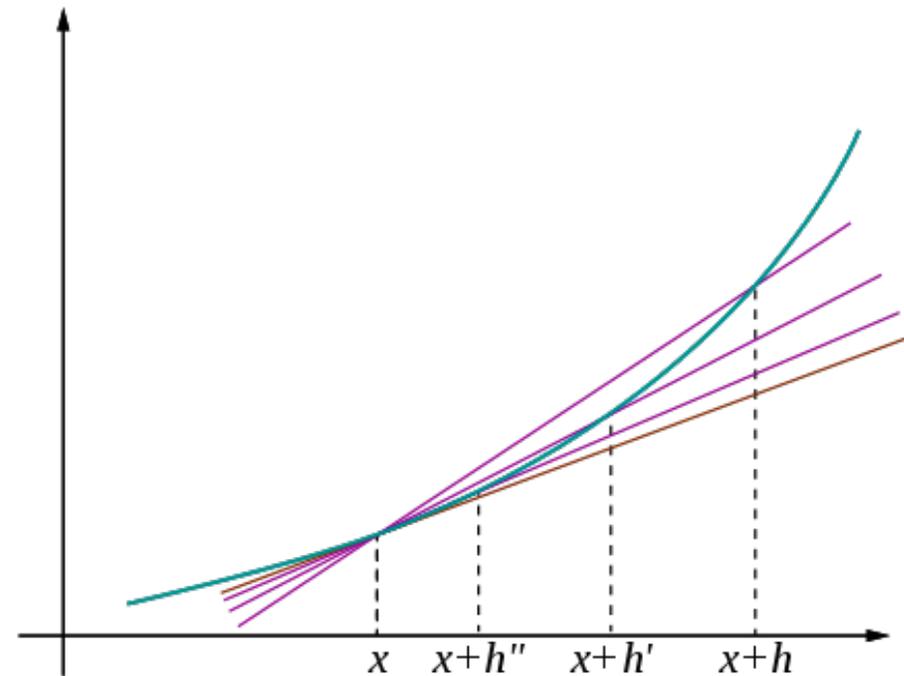
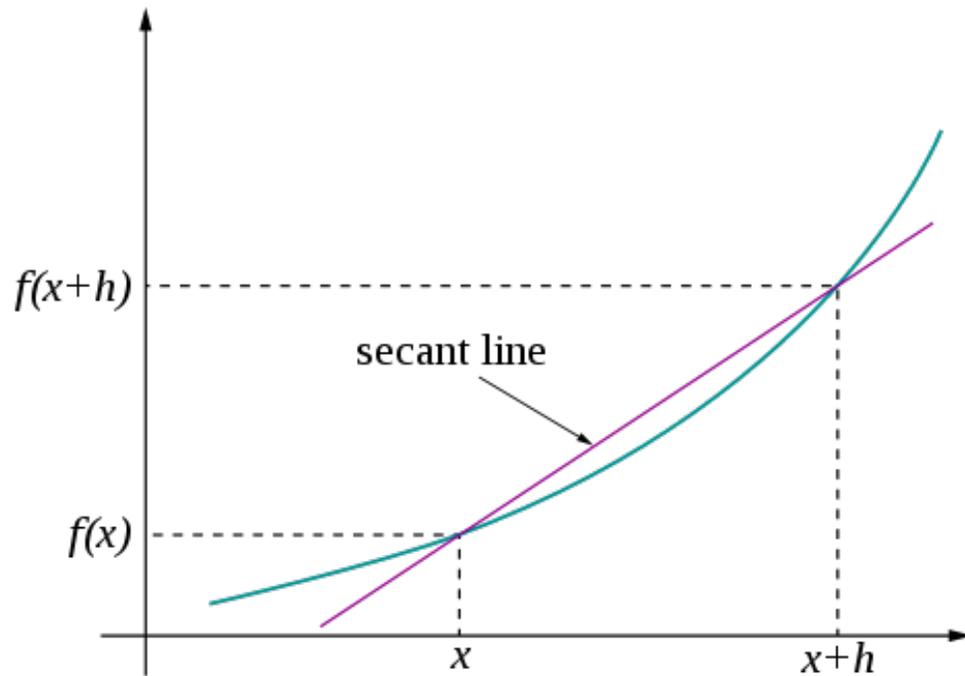
Find the value of x where $f(x)$ is minimum

Our setting: x represents weights, $f(x)$ represents loss function

In two stages

- Function of single variable
- Function of multiple variables

Derivative of a function of single variable



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

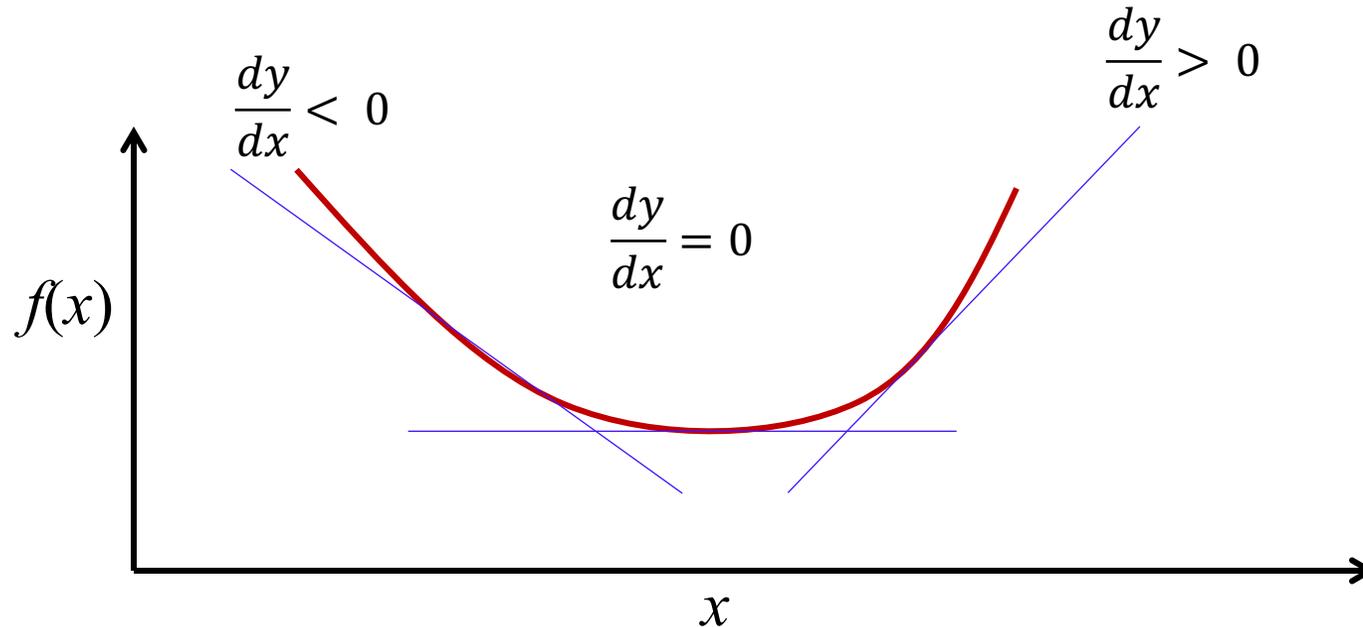
Derivatives

$$\frac{d}{dx}(x^2) = 2x$$

$$\frac{d}{dx}(e^x) = e^x$$

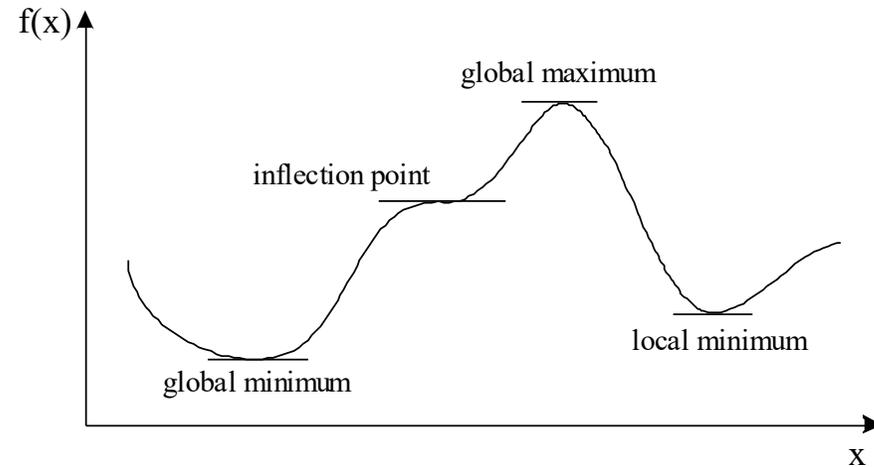
$$\frac{d}{dx}(\ln x) = \frac{1}{x} \text{ if } x > 0$$

Finding minima



Increase x if derivative negative, decrease if positive
i.e., take step in direction opposite to sign of derivative
(key idea of gradient descent)

Doesn't always work



- Theoretical and empirical evidence that gradient descent works quite well for deep networks

In two stages

- Function of single variable
- Function of multiple variables

Partial derivatives

The **partial derivative** of an n-ary function $f(x_1, \dots, x_n)$ in the direction x_i at the point (a_1, \dots, a_n) is defined to be:

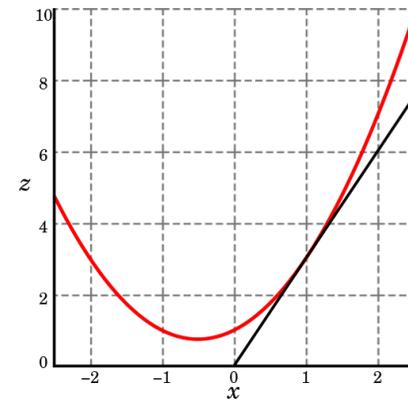
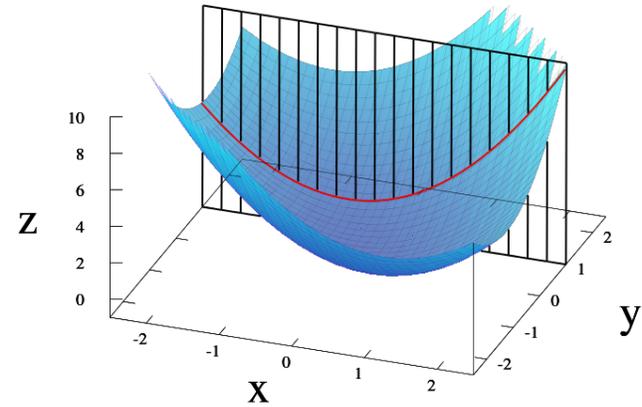
$$\frac{\partial f}{\partial x_i}(a_1, \dots, a_n) = \lim_{h \rightarrow 0} \frac{f(a_1, \dots, a_i + h, \dots, a_n) - f(a_1, \dots, a_i, \dots, a_n)}{h}.$$

Partial derivative example

$$z = f(x, y) = x^2 + xy + y^2.$$

$$\frac{\partial z}{\partial x} = 2x + y.$$

At (1, 1), the slope is 3

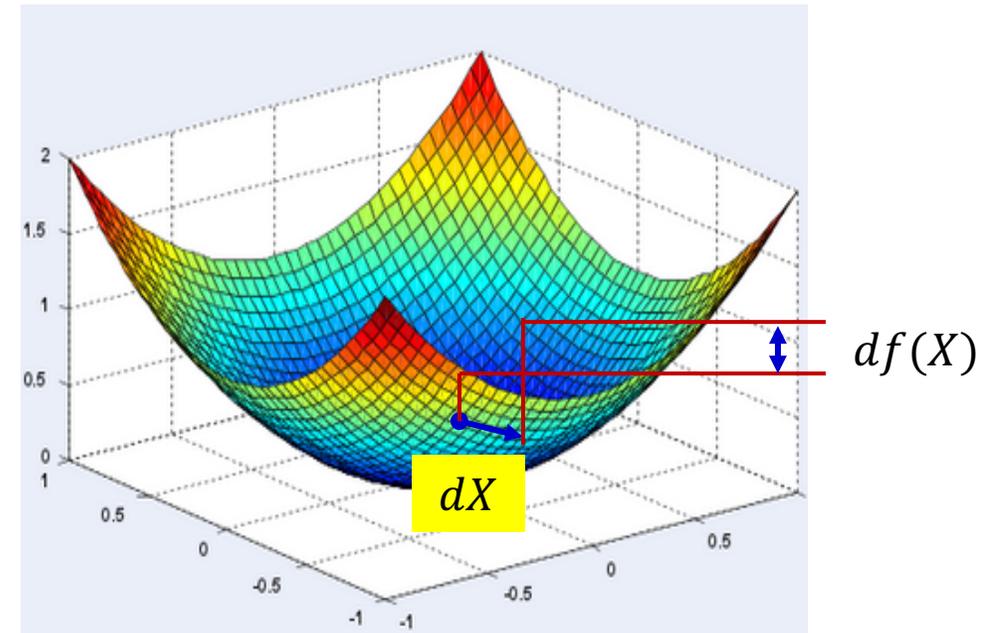


The gradient of a scalar function

- The *gradient* $\nabla f(X)$ of a scalar function $f(X)$ of a multi-variate input X is a multiplicative factor that gives us the change in $f(X)$ for tiny variations in X

$$df(X) = \nabla f(X) dX$$

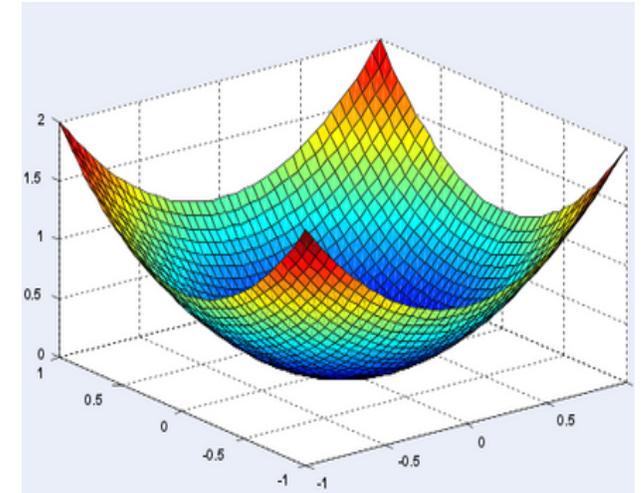
$$\nabla f(X) = df(X) / dX$$



Gradients of scalar functions with multivariate inputs

- Consider $f(X) = f(x_1, x_2, \dots, x_n)$

- $\nabla f(X) = \left[\frac{\partial f(X)}{\partial x_1} \quad \frac{\partial f(X)}{\partial x_2} \quad \dots \quad \frac{\partial f(X)}{\partial x_n} \right]$



Computing gradients analytically

$$f(x, y) = x + y \quad \rightarrow \quad \frac{\partial f}{\partial x} = 1 \quad \frac{\partial f}{\partial y} = 1$$

Computing gradients analytically

$$f(x, y) = xy \quad \rightarrow \quad \frac{\partial f}{\partial x} = y \quad \frac{\partial f}{\partial y} = x$$

$$\nabla f = \left[\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right] = [y, x]$$

Derivatives measure sensitivity

$$f(x, y) = xy$$

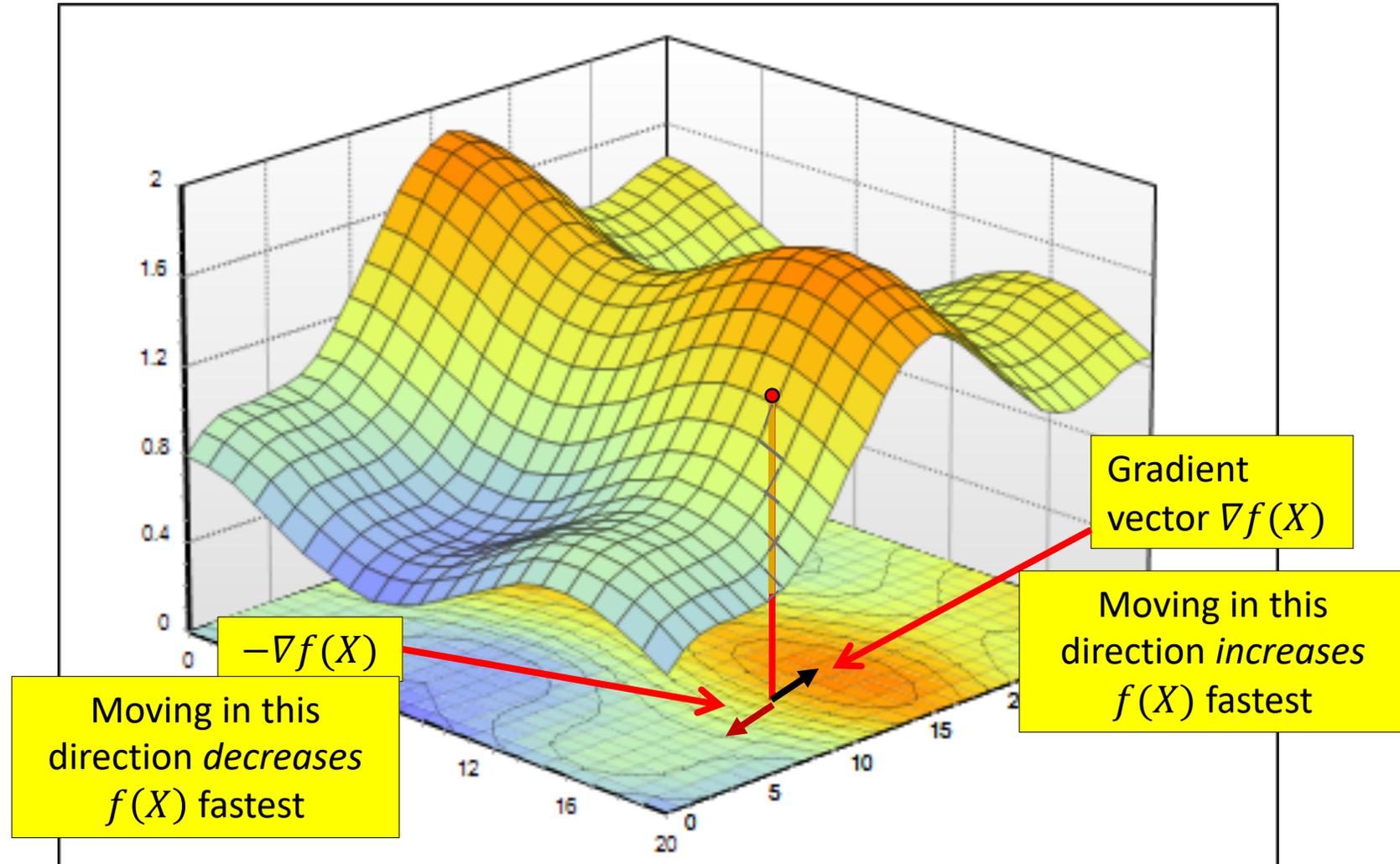
$$x = 4, y = -3 \quad f(x, y) = -12$$

$$\frac{\partial f}{\partial x} = -3$$

If we were to increase x by a tiny amount, the effect on the whole expression would be to decrease it (due to the negative sign), and by three times that amount.

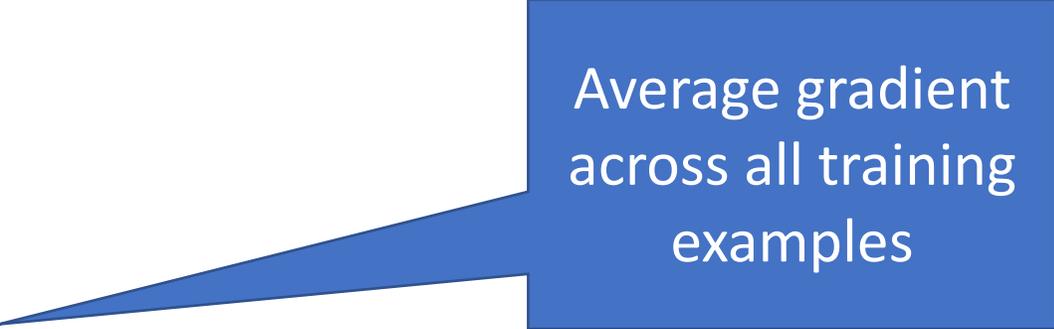
Finding minima

Take step in direction opposite to sign of gradient



Gradient descent algorithm

- Initialize:
 - x^0
 - $k = 0$
- Do
 - $x^{k+1} = x^k - \eta^k \nabla f(x^k)$
 - $k = k + 1$
- Until $|f(x^k) - f(x^{k-1})| \leq \varepsilon$

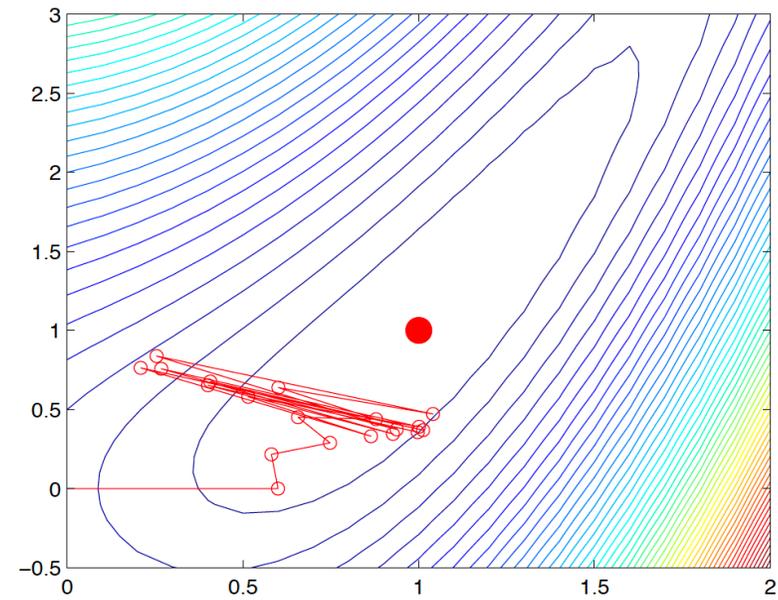
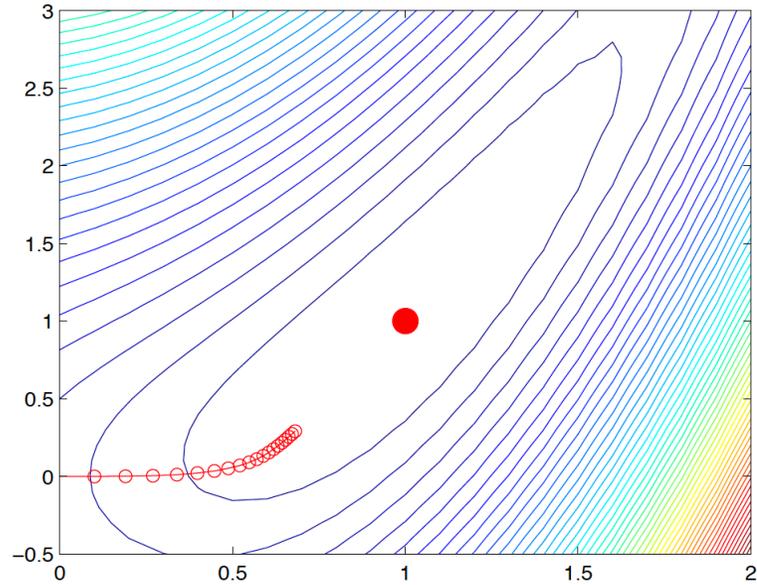


Average gradient
across all training
examples

$$f(x^k) = \frac{1}{N} \sum_{i=1}^N f_i(x^k)$$

$$\nabla f(x^k) = \frac{1}{N} \sum_{i=1}^N \nabla f_i(x^k)$$

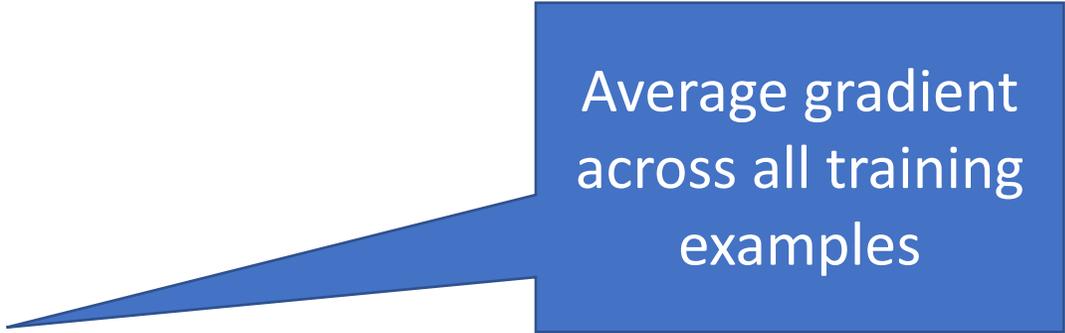
Step size affects convergence of gradient descent



Murphy, Machine Learning, Fig 8.2

Gradient descent algorithm

- Initialize:
 - x^0
 - $k = 0$
- Do
 - $x^{k+1} = x^k - \eta^k \nabla f(x^k)$
 - $k = k + 1$
- Until $|f(x^k) - f(x^{k-1})| \leq \varepsilon$



Average gradient
across all training
examples

Challenge: Not scalable for very large data sets

Challenge to discuss later: How to choose step size?

Mini-batch gradient descent

- Initialize:

- x^0
- $k = 0$

- Do

- $x^{k+1} = x^k - \eta^k \nabla f(x^k)$
- $k = k + 1$

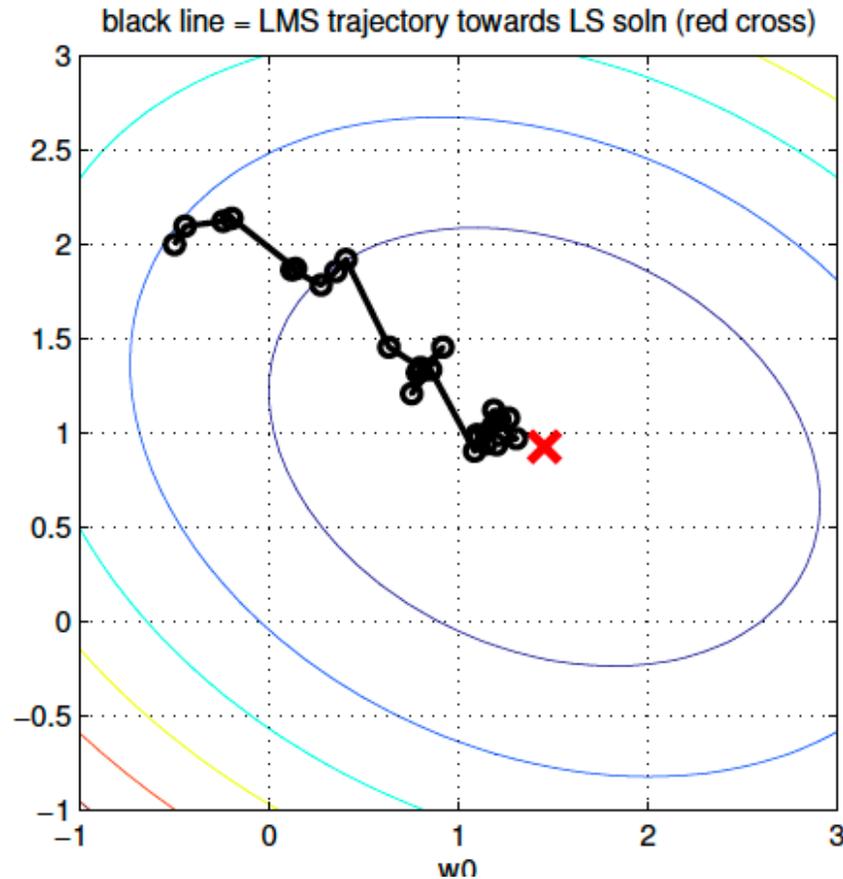
- Until $|f(x^k) - f(x^{k-1})| \leq \varepsilon$

Special case: Stochastic or online gradient descent →
use single training example in each update step

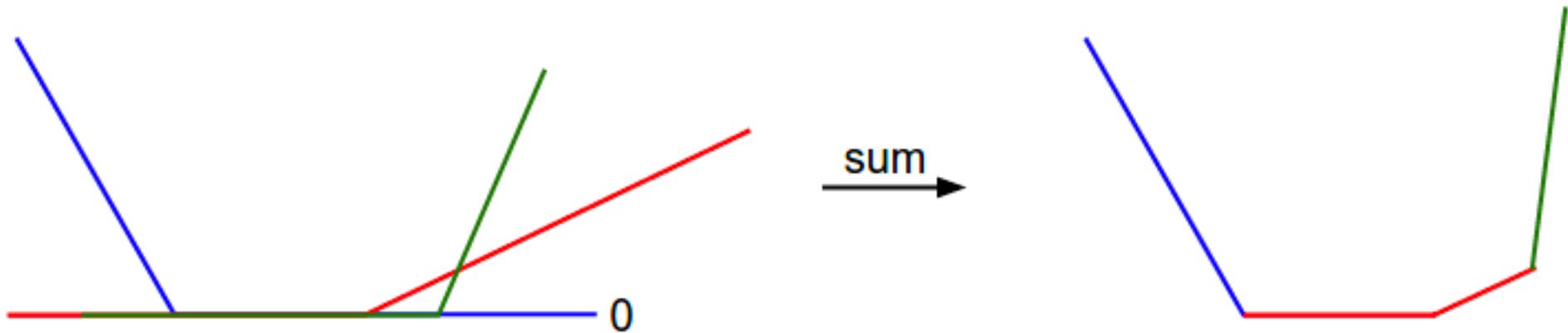
Average gradient
over small batches
of training
examples (e.g.,
sample of 256
examples)

Faster
convergence

Stochastic gradient descent convergence



SVM loss visualization



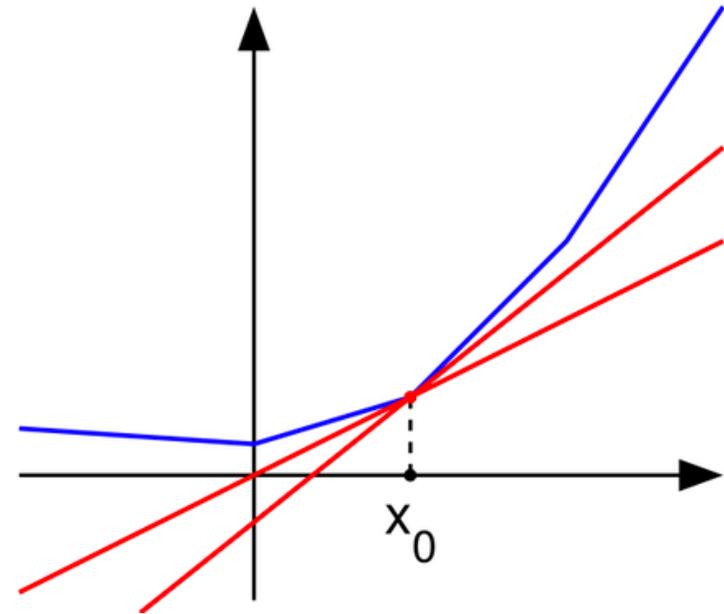
Challenge: Gradient does not exist

Computing subgradients analytically

The set of subderivatives at x_0 for a convex function is a nonempty closed interval $[a, b]$, where a and b are the one-sided limits:

$$a = \lim_{x \rightarrow x_0^-} \frac{f(x) - f(x_0)}{x - x_0}$$

$$b = \lim_{x \rightarrow x_0^+} \frac{f(x) - f(x_0)}{x - x_0}$$



Computing subgradients analytically

$$f(x, y) = \max(x, y) \quad \rightarrow \quad \frac{\partial f}{\partial x} = \mathcal{I}(x \geq y) \quad \frac{\partial f}{\partial y} = \mathcal{I}(y \geq x)$$

The (sub)gradient is 1 on the input that is larger and 0 on the other input

Subgradient of SVM loss

$$L_i = \sum_{j \neq y_i} [\max(0, w_j^T x_i - w_{y_i}^T x_i + \Delta)]$$

$$\nabla_{w_{y_i}} L_i = - \left(\sum_{j \neq y_i} \mathcal{I}(w_j^T x_i - w_{y_i}^T x_i + \Delta > 0) \right) x_i$$

Number of classes that didn't meet the desired margin

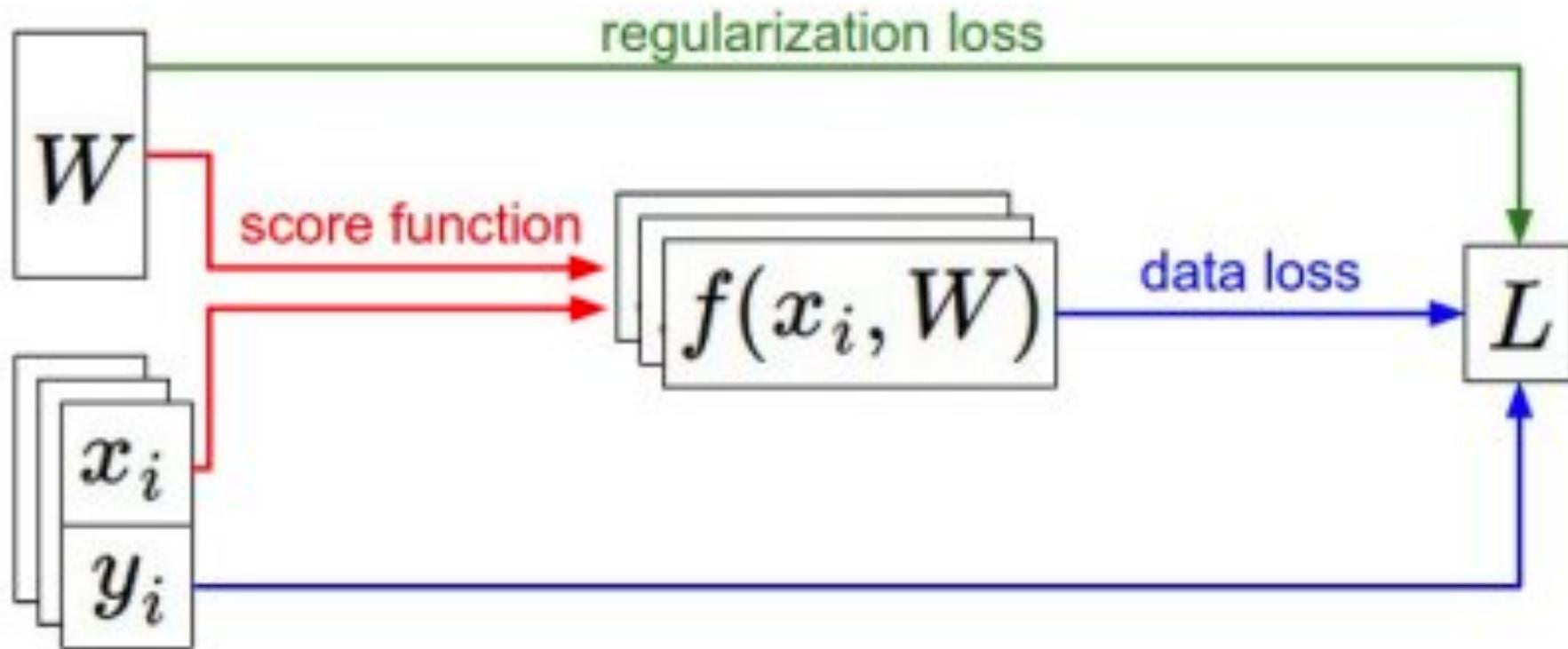
$$\nabla_{w_j} L_i = \mathcal{I}(w_j^T x_i - w_{y_i}^T x_i + \Delta > 0) x_i$$

j-th class didn't meet the desired margin

Review derivatives

- Please review rules for computing derivatives and partial derivatives of functions, including the chain rule
 - <https://www.khanacademy.org/math/multivariable-calculus/multivariable-derivatives>
- You will need to use them in HW1!

Summary



Acknowledgment

- Based on material from
 - Stanford CS231n <http://cs231n.github.io/>
 - CMU 11-785 Course
 - Spring 2019 Course