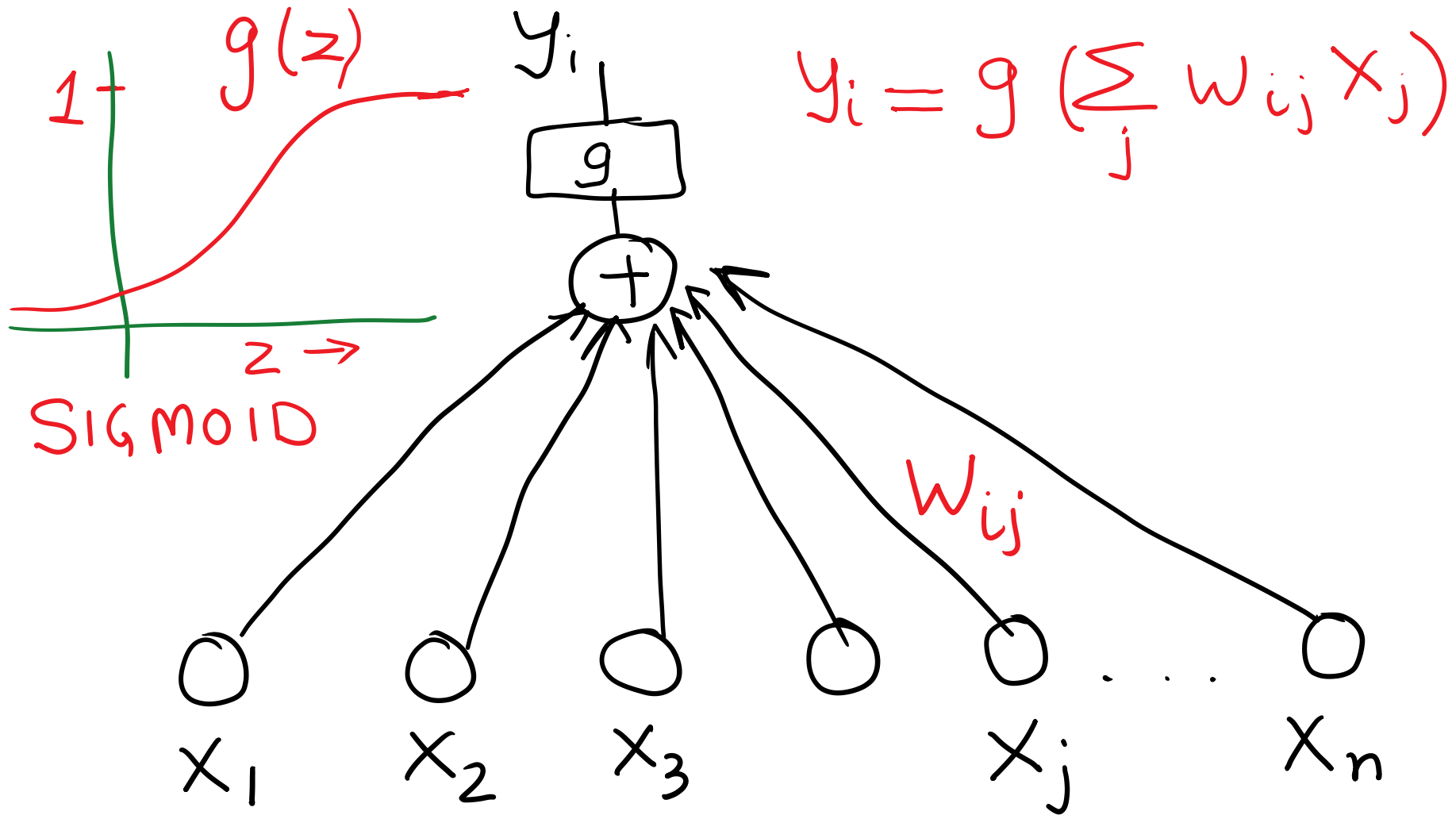# Neural Networks

Jitendra Malik

CS 189

# Mathematical Abstraction



$$Y_i = g\left(\sum_j W_{ij} X_j\right)$$

$g(z)$

$1$

$z \rightarrow$

SIGMOID

$Y_i$

$g$

$W_{ij}$

$X_1 \quad X_2 \quad X_3 \quad \cdots \quad X_j \quad \cdots \quad X_n$

# Single layer neural network

$g(z) = \dfrac{1}{1+e^{-z}}$

$1$

$Z$

$V_1$

$V_2 = g\left(\displaystyle\sum_K W_{jK} x_K\right)$

$W_{21}$   $W_{22}$   $W_{23}$   $W_{24}$   $W_{25}$

$X_1$   $X_2$   $X_3$   $X_4$   $X_5$

# Two layer neural network

$$O_i = g\left( \sum_j W_{ij} \; g\left( \sum_K W_{jK} X_K \right) \right)$$

$O_1$   $O_2$

$O_i$

$W_{ij}$

$W_{11}$   $W_{23}$

$V_1$   $V_2$   $V_3$

$V_j$

$W_{12}$   $W_{35}$   $W_{jK}$

$X_1$   $X_2$   $X_3$   $X_4$   $X_5$   $X_K$

Recall from multivariable calculus that the *gradient* of $f$ is the vector

$$\nabla f(x) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{bmatrix}.$$

The gradient is the direction which leads to a maximal increase of $f$. Similarly, the negative gradient is the direction of *steepest descent*. Gradient descent uses this fact to construct an algorithm: at every step, compute the gradient and follow that direction to minimize $f$.

# Training a neural network

Goal: Find $W$ such that $O_i$ is as close as possible to $Y_i$ (desired output)

Approach: • Define loss function $\mathcal{L}(W)$

• Compute $\nabla_W \mathcal{L}$

• $W_{new} \leftarrow W_{old} - \eta \nabla_W \mathcal{L}$

# Training a single layer neural network

- For binary classification, a good choice of loss function is the cross entropy. For regression, use squared error

$$L = -\sum_{\text{input data}} (y_i \ln O_i + (1 - y_i) \ln(1 - O_i))$$

- We model the activation function *g* as a sigmoid

$$g(z) = \frac{1}{1 + \exp(-z)}$$

- Finding w reduces to logistic regression!

We can use STOCHASTIC GRADIENT DESCENT.

# Training a multi layer neural network

- We compute the gradient with respect to all the weights: from input to hidden layer, and hidden layer to output layer. This requires computing the partial derivative of the error with respect to each weight Wij

- The complexity of computing the gradient in the naïve version is quadratic in the number of weights. The back-propagation algorithm is a trick that enables it to be computed in linear time.

# The Chain Rule from Multivariable Calculus

Given $u(x, y) = x^2 + 2y$ where $x(r, t) = r \sin(t)$ and $y(r,t) = \sin^2(t)$,

determine the value of $\partial u / \partial r$ and $\partial u / \partial t$ using the chain rule.

# The Chain Rule from Multivariable Calculus

Given $u(x, y) = x^2 + 2y$ where $x(r, t) = r \sin(t)$ and $y(r,t) = \sin^2(t)$, determine the value of $\partial u / \partial r$ and $\partial u / \partial t$ using the chain

$$\frac{\partial u}{\partial r} = \frac{\partial u}{\partial x}\frac{\partial x}{\partial r} + \frac{\partial u}{\partial y}\frac{\partial y}{\partial r} = (2x)(\sin(t)) + (2)(0) = 2r\sin^2(t),$$
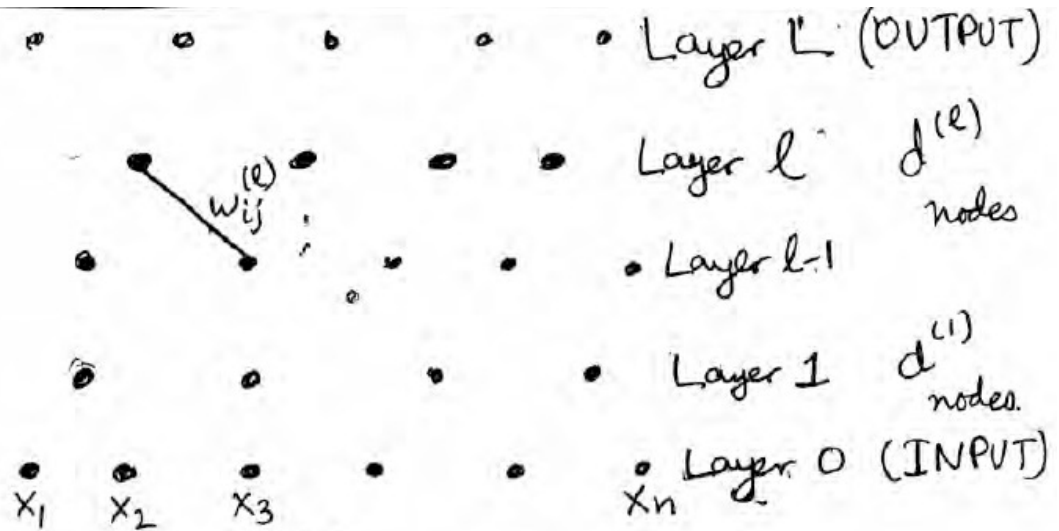
and

$$\begin{aligned}
\frac{\partial u}{\partial t} &= \frac{\partial u}{\partial x}\frac{\partial x}{\partial t} + \frac{\partial u}{\partial y}\frac{\partial y}{\partial t} \\
&= (2x)(r\cos(t)) + (2)(2\sin(t)\cos(t)) \\
&= (2r\sin(t))(r\cos(t)) + 4\sin(t)\cos(t) \\
&= 2(r^2 + 2)\sin(t)\cos(t) \\
&= (r^2 + 2)\sin(2t).
\end{aligned}$$

# Our goal

- To compute the partial derivative of error with respect to Wij

- We will work this out separately for connections feeding into the output layer and then for any other layer.

- It turns out to be useful to define two variables Sj (the summed input to a node j) $\delta j$ (the partial derivative of error w. r. t Sj)

# The Notation

## Use superscripts to keep track of layers



Layer $L$ (OUTPUT)

Layer $\ell$    $d^{(\ell)}$ nodes

$w_{ij}^{(\ell)}$

Layer $\ell-1$

Layer $1$    $d^{(1)}$ nodes.

Layer $0$ (INPUT)

$x_1 \quad x_2 \quad x_3 \qquad x_n$

Output of neuron $j$ in layer $\ell$

$$x_j^{(\ell)} = g\left( \sum_{i=0}^{d^{(\ell-1)}} w_{ij}^{(\ell)} x_i^{(\ell-1)} \right)$$

$$= g\left( s_j^{(\ell)} \right)$$

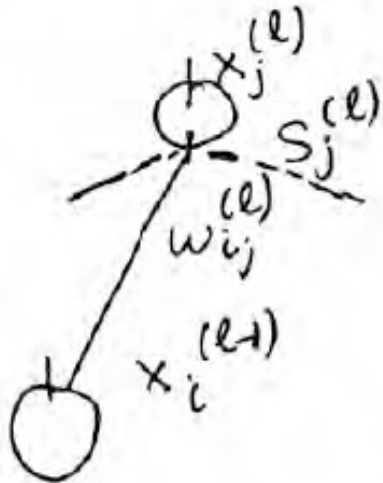$s_j^{(\ell)}$ is the weighted input to node $j$ in layer $\ell$

$w_{ij}^{(\ell)}$ is weight from node $i$ (in layer $\ell-1$) to node $j$ (in layer $\ell$)

# Defining δ

Partial derivative of error with respect to s, the summed input

Error on example $(x_n, y_n)$ is $e(h(x_n), y_n) = e(w)$
We use $h$ to denote the function computed by the neural network.

$$\nabla e(w) = \frac{\partial e(w)}{\partial w_{ij}^{(\ell)}}$$

$$\frac{\partial e(w)}{\partial w_{ij}^{(\ell)}} = \frac{\partial e(w)}{\partial s_j^{(\ell)}} \times \frac{\partial s_j^{(\ell)}}{\partial w_{ij}^{(\ell)}}$$

$$= \delta_j^{(\ell)} \quad x_i^{(\ell-1)}$$

Define $\delta_j^{(\ell)} = \dfrac{\partial e(w)}{\partial s_j^{(\ell)}}$

# Computing δ for neuron in final layer
## (this example is for squared loss, as in regression)

For the final layer $\ell = L$ and $j = 1$. For $e(h(x_n), y_n)$

Squared error, $\dfrac{\partial e(w)}{\partial s_1^{(L)}} = \dfrac{\partial}{\partial s_1^{(L)}} \left\{ \dfrac{1}{2} \left( g(s_1^{(L)}) - y_n \right)^2 \right\}$

**BASE STEP**

$$= \left( g(s_1^{(L)}) - y_n \right) \cdot g'(s_1^{(L)}) \quad --- \textcircled{1}$$

This will change if we use a different loss function.

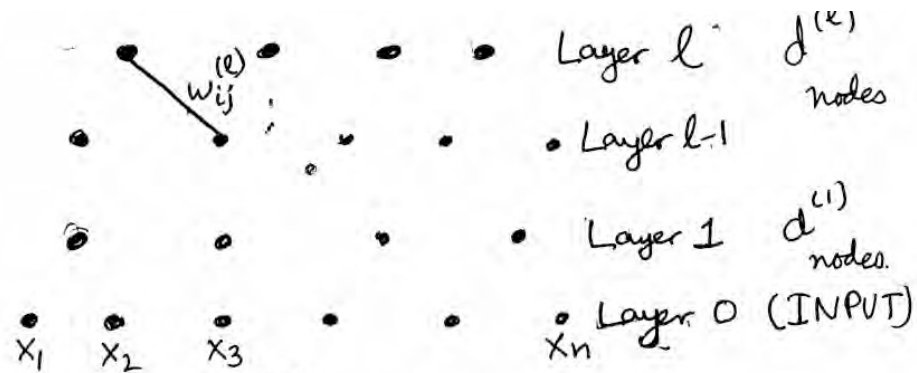# Computing δ for neuron in intermediate layer
## This is just the chain rule from calculus

INDUCTION STEP

$$\delta_i^{(\ell-1)} = \frac{\partial e(w)}{\partial s_i^{(l-1)}}$$

$$= \sum_j \frac{\partial e(w)}{\partial s_j^{(\ell)}} \times \frac{\partial s_j}{\partial x_i^{(\ell-1)}} \times \frac{\partial x_i^{(\ell-1)}}{\partial s_i^{(l-1)}}$$

# Computing $\delta$ for neuron in intermediate layer

Layer $\ell$    $d^{(\ell)}$
nodes

$w_{ij}^{(\ell)}$

Layer $\ell$-1

Layer 1    $d^{(1)}$ nodes.

Layer 0 (INPUT)

$x_1$   $x_2$   $x_3$        $x_n$

$$x_j^{(\ell)} = g\left( \sum_{i=0}^{d^{(\ell-1)}} w_{ij}^{(\ell)} x_i^{(\ell-1)} \right)$$

$$= g\left( s_j^{(\ell)} \right)$$

$s_j^{(\ell)}$ is the weighted input to node $j$ in layer $\ell$

INDUCTION STEP

$$\delta_i^{(\ell-1)} = \frac{\partial e(w)}{\partial s_i^{(\ell-1)}}$$

$$= \sum_j \frac{\partial e(w)}{\partial s_j^{(\ell)}} \times \frac{\partial s_j}{\partial x_i^{(\ell-1)}} \times \frac{\partial x_i^{(\ell-1)}}{\partial s_i^{(\ell-1)}}$$

# Computing δ for neuron in intermediate layer
## Deriving the basic step of "backpropagation"

$$\boxed{\text{INDUCTION STEP}}$$

$$\delta_i^{(\ell-1)} = \frac{\partial e(w)}{\partial s_i^{(l-1)}}$$

$$= \sum_j \frac{\partial e(w)}{\partial s_j^{(\ell)}} \times \frac{\partial s_j}{\partial x_i^{(\ell-1)}} \times \frac{\partial x_i^{(\ell-1)}}{\partial s_i^{(l-1)}}$$

$$\delta_i^{(\ell-1)} = \sum_j \delta_j^{(\ell)} \; w_{ij}^{(\ell)} \; g'\left(s_i^{(\ell-1)}\right)$$

$$= g'\left(s_i^{(\ell-1)}\right) \sum_j w_{ij}^{(\ell)} \delta_j^{(\ell)} \qquad - \cdot \cdot - \textcircled{2}$$

# The full algorithm

We have thus defined an iterative algorithm for computing $\delta$ for every node

, compute it for nodes on output layer $L$ using ①
• compute it for a layer $(\ell-1)$ using $\delta$ values at a layer $(\ell)$ using ②

Once this process has concluded, we can use

$$\frac{\partial e(w)}{\partial w_{ij}^{(\ell)}} = \delta_j^{(\ell)} \, x_i^{(\ell-1)} \qquad \text{to compute}$$

the derivative with respect to any weight $w_{ij}^{(\ell)}$

# Summary

- Note that we have computed all of the partial derivatives in one backward pass. The time complexity is proportional to number of weights

- We could have done this, one partial derivative at a time, using finite differences. This is a good way to debug your code for backprop. But it costs more, because you need to repeat the computation as many times as you have weights.

- The basic steps in backpropagation can be expressed as matrix multiplications, and GPUs are great at that. This is what enabled the scaling up of deep learning in the 2010s.

# Basics of optimization

- Distinction between local and global minima
- Convex functions
- Gradient descent
- Gradient descent enables us to find local minima, but for a convex function, a local minimum is also a global minimum

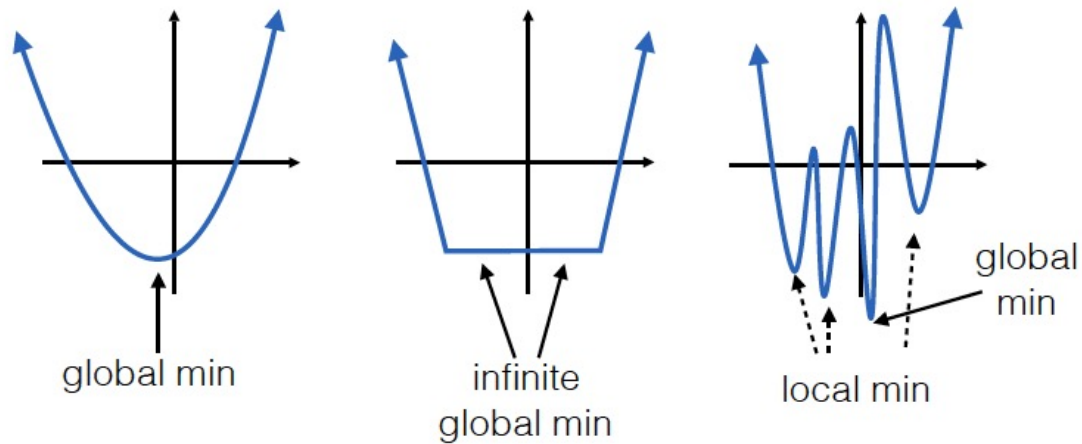$$\text{minimize}_w \quad f(w)$$



Figure 1: Examples of minima of functions. In the first illustration, there is a unique minimizer. In the second, there are an infinite number of minimizers, but all local minimizers are global minimizers. In the third example, there are many local minimizers that are not global minimizers.
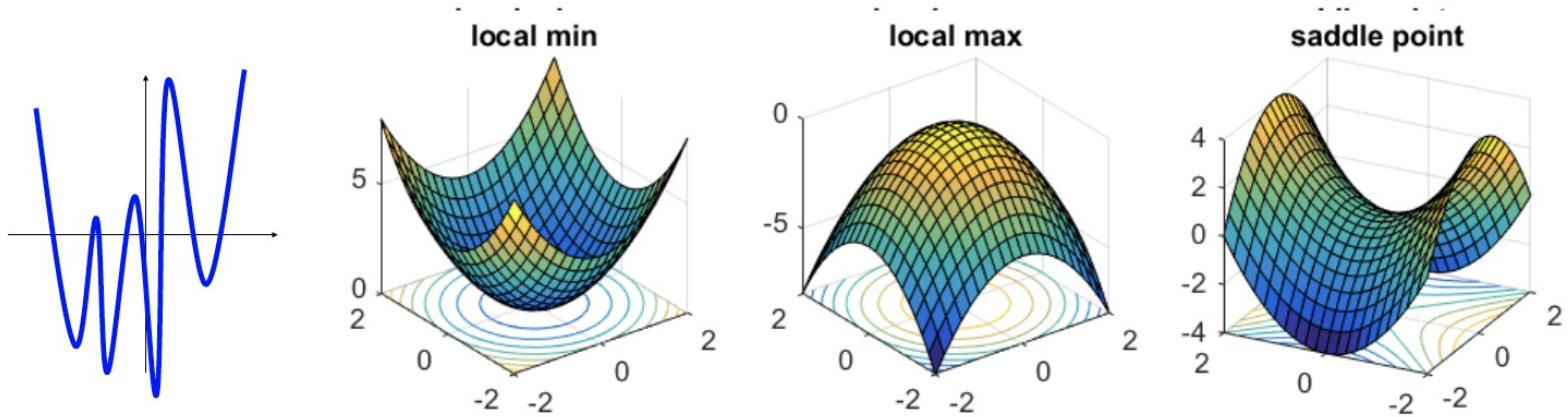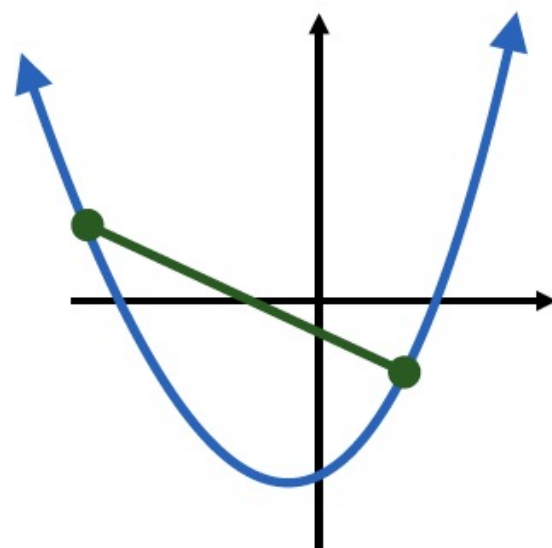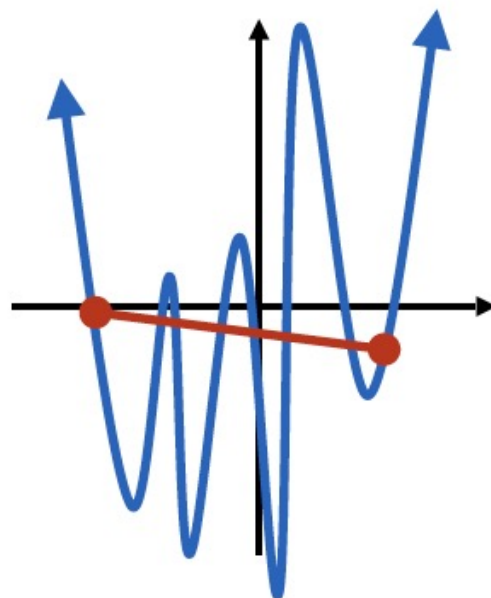
# Stationary points for functions of two variables



Figure 3: Examples of stationary points.

convex          nonconvex

- A function $f$ is *convex* if for all $w_1$, $w_2$ in $\mathbb{R}^d$ and $t \in [0, 1]$,

$$f(tw_1 + (1-t)w_2) \leq tf(w_1) + (1-t)f(w_2) .$$

Recall from multivariable calculus that the *gradient* of $f$ is the vector

$$\nabla f(x) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{bmatrix}.$$

The gradient is the direction which leads to a maximal increase of $f$. Similarly, the negative gradient is the direction of *steepest descent*. Gradient descent uses this fact to construct an algorithm: at every step, compute the gradient and follow that direction to minimize $f$.

# Gradient Descent

The gradient descent method follows the simple algorithmic procedure:

1. Choose $x_0 \in \mathbb{R}^d$ and set $k = 0$

2. Choose $t_k > 0$ and set $x_{k+1} = x_k - t_k \nabla f(x_k)$ and $k = k + 1$,

3. Repeat 2 until converged.

# Two ways to characterize a "descent direction"

- $v$ is a *descent direction* for $f$ at $x_0$ if $f(x_0 + tv) < f(x_0)$ for some $t > 0$.

For continuously differentiable functions, it's easy to tell if $v$ is a descent direction: if $v^T \nabla f(x_0) < 0$ then $v$ is a descent direction.

To see this note that by Taylor's theorem, $f(x_0+tv) = f(x_0)+t\nabla f(x_0+\tilde{t}v)^T v$ for some $\tilde{t} \in [0,t]$. By continuity, if $t$ is small, we'll have $\nabla f(x_0 + \tilde{t}v)^T v < 0$. Therefore $f(x_0 + tv) < f(x_0)$ and $v$ is a descent direction.

Note that among all directions with unit norm, the *steepest descent* possible is given when we move in the direction of the negative gradient.

# Taylor's theorem from multivariable calculus

One of the most important theorems in calculus is *Taylor's Theorem*, which allows us to approximate smooth functions by simple polynomials. The following simplified version of Taylor's Theorem is used throughout optimization. This form of Taylor's theorem is sometimes called the multivariable mean-value theorem. It states that

**Theorem 1 (Taylor's Theorem)** *If $f$ is continuously differentiable, then*

$$f(x) = f(x_0) + \nabla f(tx + (1-t)x_0)^T (x - x_0) \quad \text{for some } t \in [0, 1]$$

$$f(x_0 + tv) = f(x_0) + t \nabla f(x_0 + \tilde{t}v)^T v$$