# The Gradient and Gradient Descent

Recall the Situation, the Deep/Feed forward
neural net (also called MLP= Multi Layered Perceptron)
has the form

$$F(x/\eta) = F_L \circ \cdots \circ F_1(x) \quad \text{with each}$$

$$F_i(x) = \sigma(A_L x + \vec{b}_L)$$

- The parameters are $\eta = A_1, \eta, A_L, b_1, \cdots, b_L$
- The training data is $x_1, \cdots, x_p$ with correct output
$y_1, y_2, \cdots, y_N$

- The loss or error or objective function is (simplest version)

$$\Phi(\eta) = \frac{1}{N} \sum_{i=1}^{N} \| F(x_L, \eta) - y_L \|^2$$

- Learning consists of adjusting the parameters
$\eta$ so that the error diminishes.

• To accomplish this we need to recall
the tools from multi-variable Calculus,
specifically, the gradient.

## Review of the Gradient

• Now we let $x_1, \ldots, x_n$ revert to the usual Calculus
roles as components of $\vec{x} = (x_1, \ldots, x_n) \in \mathbb{R}^n$

• Let $\Phi$ be a real valued function of $\vec{x}$
so $\Phi(\vec{x}) = \Phi(x_1, \ldots, x_n)$ is a scalar

• We want to understand how $\Phi$ is
changing in the various directions in $\mathbb{R}^n$

◦ As a simple example let $\Phi(x_1, x_2) = 9x_1^2 + x_2^2$

- We think of $\Phi$ as the temperature at each point in the plane for concreteness
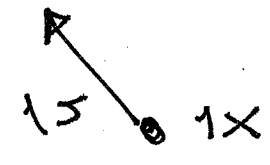
- Starting at the point
$$\vec{x} = (x_1, x_2)$$

We walk in a direction given by the unit vector
$$\vec{u} = (u_1, u_2) \quad [\text{I am writing things as row vectors like is done in Calculus}]$$



- How does $\Phi$ change? We can write a limit

$$\lim_{t \to 0} \frac{\Phi(\vec{x} + t\vec{u}) - \Phi(\vec{x})}{t} = D_{\vec{u}} \Phi(x)$$
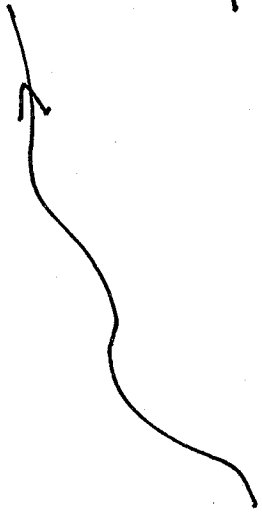
- This is called the directional derivative of $\Phi$ in the direction $\vec{u}$

- How do we compute this?

- Let $\gamma(t)$ by a path with unit speed

  so $\left|\frac{d\gamma}{dt}\right| = 1$ i.e. its velocity

- $\overline{\Phi}$ changes along the path

  We WATCH how $\gamma(t)$

and call this $g(t) = \overline{\Phi}(\gamma(t))$

so if $\gamma(0) = \vec{x}$ we seek $g'(0) =$

where $\vec{u} = \frac{d\gamma(0)}{dt}$

$$D_{\vec{u}}\,\overline{\Phi}(\vec{x})$$

- We compute this from the chain rule

$$g(t) = \Phi(x(t))$$

$$\text{so} \quad \frac{dg(t)}{dt} = \nabla\bar\Phi(x(t)) \cdot \frac{dx(t)}{dt} \quad \text{with} \quad \nabla\Phi = \left[\frac{\partial\Phi}{\partial x_1}, \cdots, \frac{\partial\Phi}{\partial x_n}\right]$$

evaluating at $t=0$, $x(0)=x$, $\frac{dx}{dt} = \vec{u}$

$$\text{so} \quad \frac{dg}{dt} = \nabla\bar\Phi(x) \cdot \vec{u} = D_{\vec{u}}\bar\Phi(x)$$
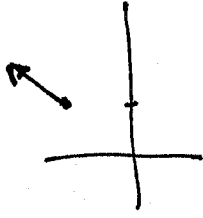
the directional derivative

• Back to $\bar\Phi(x_1, x_2) = 9x_1^2 + x_2^2$, Find $D_{\vec{u}}\bar\Phi(x)$
which $\vec{x} = (1,2)$ and $\vec{u} = \left(\frac{\sqrt{3}}{2}, \frac{1}{2}\right)$. Note $\nabla\bar\Phi = [18x_1, 2x_2]$

$$\cdot \; D_{\vec{u}}f(x) = \nabla\bar\Phi(x) \cdot \vec{u} = (18,4) \cdot \left(\frac{\sqrt{3}}{2}, \frac{1}{2}\right) = 9\sqrt{3} + 2.$$
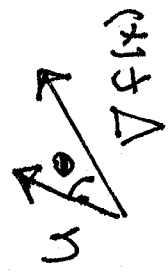
• Recall our goal is to find the direction in which $\overline{\Phi}$ is decreasing most rapidly

• We know that the rate of change of $\overline{\Phi}$ at $x$ in the direction $\vec{u}$ is

$$D_{\vec{u}} \overline{\Phi}(x) = \nabla \overline{\Phi}(x) \cdot \vec{u} = |\nabla \overline{\Phi}(x)| |\vec{u}| \cos \theta$$

where $\theta$ is the angle between

$$\vec{u} \quad \text{and} \quad \nabla \overline{\Phi}(x)$$


$\nabla f(x)$
$u$

Now we now $\cos \theta$ is most positive when

$\theta = 0$ $(\cos(0) = 1)$ and most negative when

$\theta = \pi$ $(\cos(\pi) = -1)$

A function is increasing when $g' > 0$ and decreasing when $g' < 0$ so

$$g(t) = \bar{\Phi}(\mathscr{X}(t))$$ has its maximum

increase when $\theta = 0$ or when $\frac{d\mathscr{X}}{dt} = \vec{u}$ is parallel to

$\nabla\bar{\Phi}(x)$ and has its maximum decrease when

$\frac{d\mathscr{X}}{dt}$ points in the opposite direction

$$D_u \bar{\Phi}(x) = |\nabla\bar{\Phi}(x)| |\vec{u}| \cos\theta = |\nabla\bar{\Phi}(x)| \cos\theta$$

Since $\vec{u}$ is a unit vector the result we want is

since $\vec{u}$ is the direction of maximum decrease of

Theorem The direction of maximum decrease of

$\bar{\Phi}$ at the point $x$ is the direction of $-\nabla\bar{\Phi}(x)$
$\uparrow$
important minus
sign

So back to $\Phi(x_1, x_2) = 9x_1^2 + x_2^2$.

The direction of maximum decrease of $\Phi$ at $(1,2)$ is $-\nabla\Phi(1,2) = -(18,2)$

- Before we get back to the task of diminishing $\Phi$ to decrease the error (learning) we need to learn one more thing about $\nabla\Phi$

- A level set of $\Phi$ is a set of the form $\{x : \Phi(x) = c\}$ for some constant $c$

For $\overline{\Phi}(x_1,x_2) = 9x_1^2 + x_2^2 = C$, dividing by $C > 0$

we get

$$\frac{x_1^2}{\frac{C}{9}} + \frac{x_2^2}{C} = 1 \quad \text{which}$$

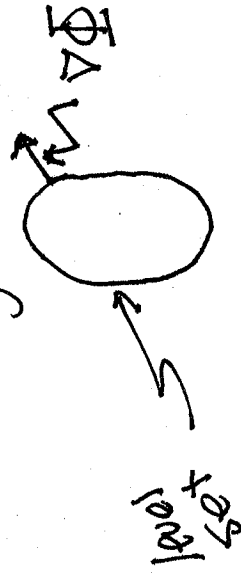is an ellipse with $x_1$ width $\sqrt{\frac{C}{9}}$ and $x_2$ width $\sqrt{C}$



Now recall that $D_u \overline{\Phi}(x) = \nabla \overline{\Phi}(x) \cdot \vec{u} = |\nabla \overline{\Phi}(x)| \cos\theta$

This is zero when $\theta = \frac{\pi}{2}, \frac{3\pi}{2}$.

So the direction of the level set where

$\Phi$ is not changing is perpendicular to $\nabla \Phi$

## Summary :

Given $\underline{\Phi} : \mathbb{R}^n \to \mathbb{R}$ at a point $\vec{x} \in \mathbb{R}^n$

the direction of maximal increase of $\underline{\Phi}$ is given by

$\nabla \underline{\Phi}(x)$ and the direction of maximal decrease is

given by $-\nabla \underline{\Phi}(x)$. Further, for a level

set $L_c = \{ \vec{x} \in \mathbb{R}^n : \underline{\Phi}(\vec{x}) = c \}$ at a point

$\vec{x}$ in $L_c$, $\nabla \underline{\Phi}(\vec{x})$ is perpendicular to

the level set [ provided level set is nice

at $x$, no corners, etc ]
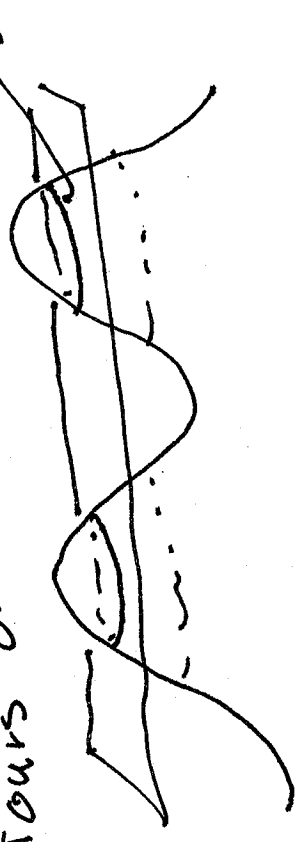
- How do we use this information to decrease

the error $\Phi$

- Let's stick with the role of $\vec{x}$ as the independent
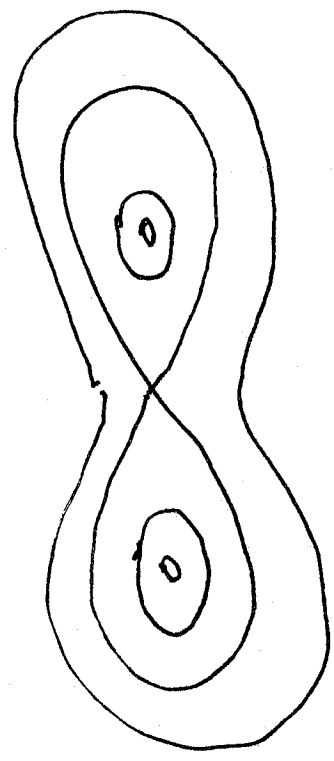  variable

- As an example, Let $\Phi(x_1, x_2)$ be the height
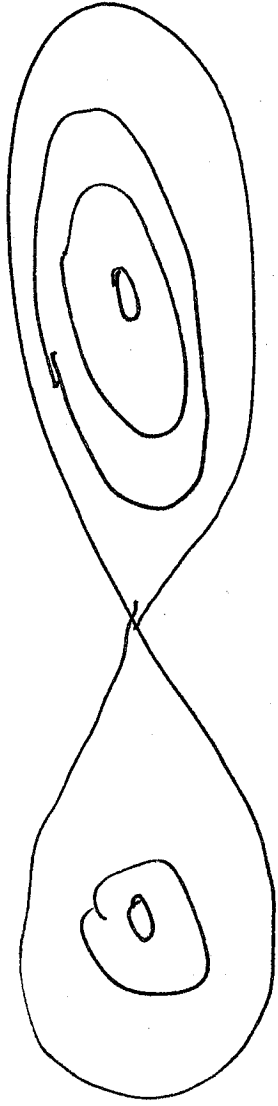  of a terrain. So the level sets are called
  contours or a contour map
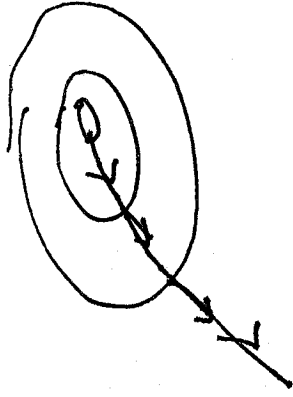
  level sets

  graph of $\Phi$, two
  mountains

  you want to get
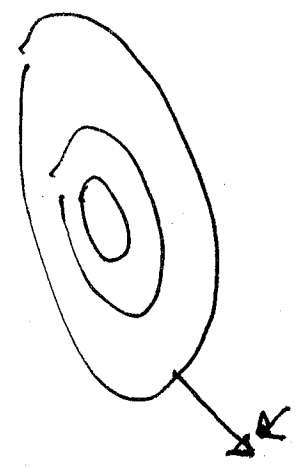  off the mountain
  as quickly as possible

So at each step you go in the direction of steepest descent ( biggest decrease in $\Phi$ ) which is $-\nabla\Phi$ and this will be perpendicular to the contour lines



If your path is $R(t)$ its satisfies the differential equation $\frac{d}{dt}R(t) = -\nabla\Phi(R(t))$

rather than solve this, we discretize and take

since we have to

take discret steps.

So if your initial position is $\vec{x}_0$ and your

steps have length $h$ after one step you are

at $\vec{x}_0 + h(-\nabla\vec{\Phi}(x_0)) = \vec{x}_1$, your new position



$x_0$

$x_1$    $h$ is magnified

[We are assuming that your step on your next step

by $|\Delta\vec{\Phi}(x_0)|$.

$\vec{x}_2 = \vec{x}_1 - h \nabla\vec{\Phi}(x_1)$, etc.

So gradient descent is given by

- Requires initial point $\bar{x}_0$ and subroutine to compute $\nabla \Phi$ and step size $h$

  for $i = 1$ to $n$

  $$x_i = x_{i-1} - h \nabla \underline{\Phi}(x_{i-1})$$

  end.

This is the basic outline, but it raises many questions

(1) How do you choose $h$=stepsize or learning rate?

(2) How do you choose, or other halting condition?