# 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 $\eta = A_1, \ldots, A_L, b_1, \ldots, b_L$

- The training data is $x_1, \ldots, x_p$ with correct output
  $y_1, y_2, \ldots, y_N$

- The loss or error or objective function is (simplest version)
  $$\Phi(\eta) = \frac{1}{N} \sum_{L=1}^{N} \| F(x_L; \eta) - y_L \|^2$$

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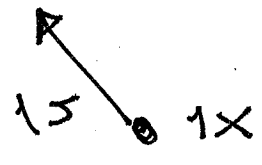  $\eta$ so that the error diminishes.

- 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 various directions in $\mathbb{R}^n$

  o 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

. Startiny at the point ~~$\vec{X_o} = \cancel{(X_1,X_2)}$~~
$\vec{X} = (X_1, X_2)$

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



. How does $\phi$ ~~change~~? We can write a limit

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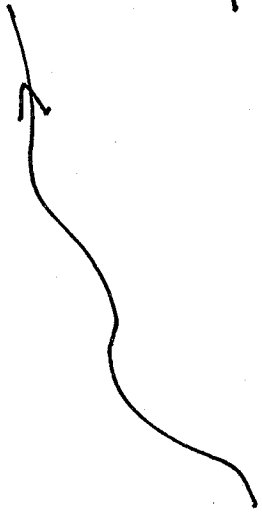
• This is called the directional derivative of $\phi$
in the direction $\vec{u}$

- How do we compute This?

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

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

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

    We WATCH how $\overline{\Phi}$ changes along the path

    $\gamma(t)$

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

    we seek $g'(0) = D_t \overline{\Phi}(t) \cdot D_\gamma \overline{\Phi}(t)$

    So if $\gamma(0) = \vec{x}$

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

- We compute This from the chain rule

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

so $\dfrac{dg(t)}{dt} = \nabla\Phi(x(t))\cdot\dfrac{dx(t)}{dt}$ with $\nabla\Phi = \left[\dfrac{\partial\Phi}{\partial x_1}, \ldots, \dfrac{\partial\Phi}{\partial x_n}\right]$
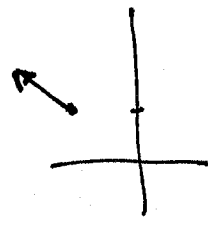
evaluating at $t=0$, $x(0)=x$, $\dfrac{dx}{dt}=\vec{u}$

so $\dfrac{dg}{dt} = \nabla\Phi(x)\cdot\vec{u} = D_{\vec{u}}\Phi(x)$

the directional derivative

• Back to $\Phi(x_1,x_2) = 9x_1^2 + x_2^2$, First ~~~~~ $D_{\vec{u}}\Phi(x)$

which $\vec{x} = (1,2)$ and $\vec{u} = \left(\dfrac{\sqrt{3}}{2}, 1/2\right)$, Note $\nabla\Phi = [18x_1, 2x_2]$

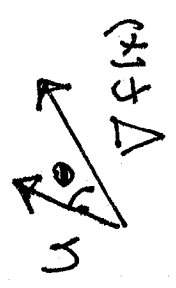$D_{\vec{u}}f(x) = \nabla\Phi(x)\cdot\vec{u} = (18,4)\cdot\left(\dfrac{\sqrt{3}}{2}, 1/2\right) = 9\sqrt{3}+2$.

• Recall our goal is to find the direction in which $\Phi$ is decreasing most rapidly

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

$$D_u \Phi(x) = \nabla \Phi(x) \cdot \vec{u} = |\nabla \Phi(x)||\vec{u}| \cos \Theta$$

where $\Theta$ is the angle between $\nabla \Phi(x)$ and $\vec{u}$



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

$$g(t) = \bar{\phi}(\gamma(t))$$ has its MAXIMUM

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

$\nabla\bar\phi(x)$ and has its MAXIMUM decrease when

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

Since $D_u \bar\phi(x) = |\nabla\bar\phi(x)||\vec{u}|\cos\theta = |\nabla\phi(x)||\vec{u}|\cos\theta$ the result we want is

$|\nabla\phi(x)|$

Since $\vec{u}$ is a unit vector

Theorem The direction of MAXIMUM decrease of

$\bar\phi$ at the point $x$ is the direction of $-\nabla\bar\phi(x)$

Important MINUS sign

9

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

The direction of maximum decrease of $\Phi$ at

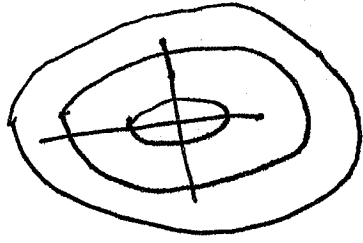$(1,2)$ is $-\nabla\Phi(1,2) = -(18,2)$

• 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 $\frac{\sqrt{C}}{3}$ and $x_2$ width $\sqrt{C}$



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

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

So the direction of the level set where

$\overline{\Phi}$ is not changing is perpendicular to $\nabla\overline{\Phi}$

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

the direction of maximal increase of $\phi$ is given by

$\nabla \phi (x)$ and the direction of maximal decrease is

given by $-\nabla \phi (x)$. Further, for a level

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

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

the level set [provided level set is nice

at $x$, no corners, etc]

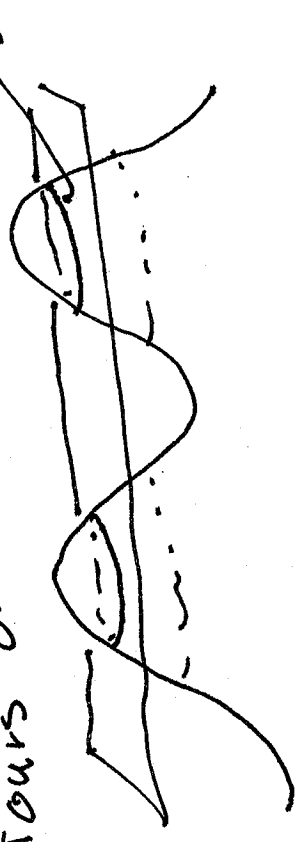- How do we use this information to decrease the error $\Phi$

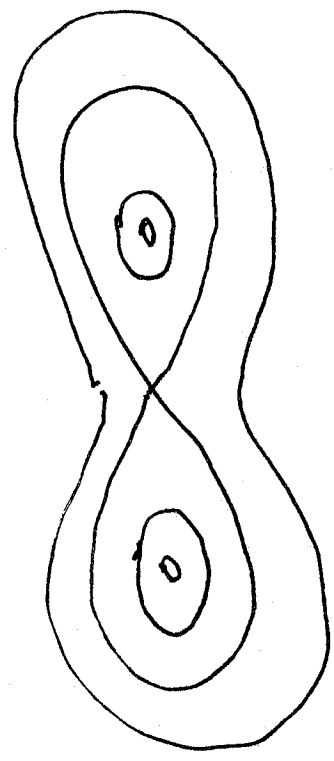- Let's stick with the role of $\vec{x}$ as the independent variable

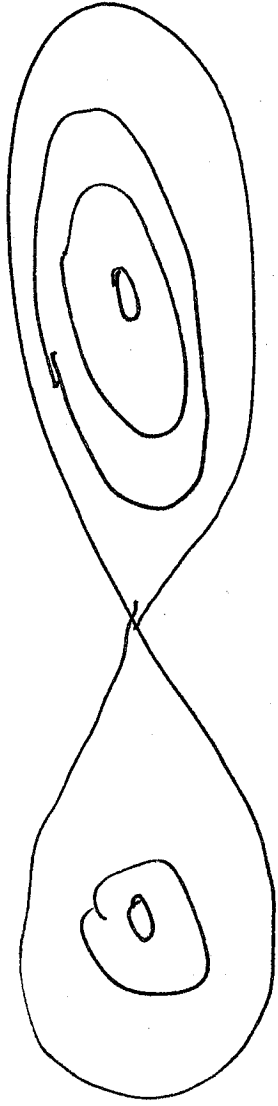- As an example, Let $\Phi(x_1, x_2)$ be the height of a terrain. So the level sets are called contours or a contour map
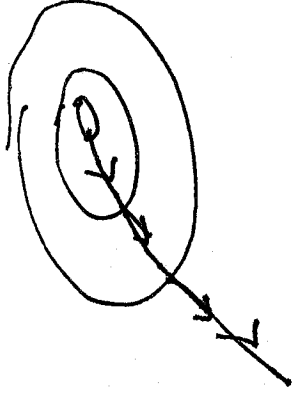
level sets

graph of $\Phi$, two mountains
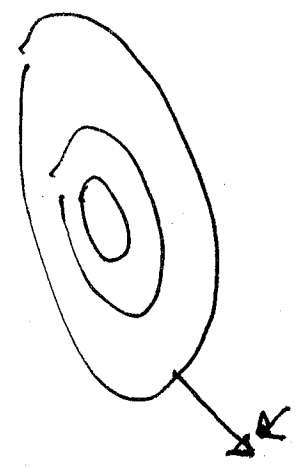
You WANT TO get off the mountain as quickly as possible

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
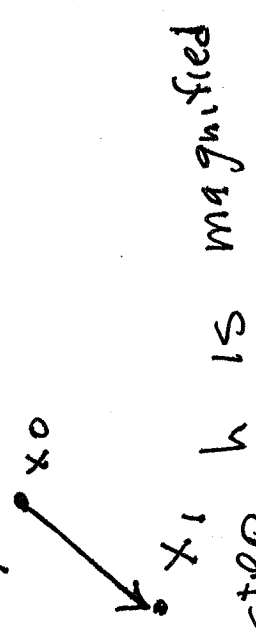


If your path is $R(t)$ its satisfies the differential equation

$$\frac{dR(t)}{dt} = -\nabla\Phi(R(t))$$

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

$$\vec{x_0} + h(-\nabla\bar{\phi}(x_0)) = \vec{x_1}, \text{ your new position}$$



[ We are assuming that your step $h$ is magnified by $|\nabla\bar{\phi}(x_0)|$. On your next step

$$\vec{x_2} = \vec{x_1} - h\nabla\bar{\phi}(x_1), \text{ etc.}$$

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?