# Warm up

$$j = 1, \ldots, p$$

$$x_i \rightarrow h_j(x_i) = \cos\left(g_j^T x_i + b_j\right)$$

```
# generate some nonsense data for an example
X = np.random.randn(n,d)        n×d
y = np.random.randn(n)
```

$$X = \begin{bmatrix} - x_1 - \\ \vdots \\ - x_n - \end{bmatrix} \in \mathbb{R}^{n \times d}$$

1 float in NumPy = 8 bytes
$10^6 \approx 2^{20}$ bytes = 1 MB
$10^9 \approx 2^{30}$ bytes = 1 GB

```
# generate the random features
G = np.random.randn(p, d)*np.sqrt(.1)
b = np.random.rand(p)*2*np.pi
```

$$H^T H = \sum_{i=1}^{\hat{}} h_i h_i^T$$

$$H = \begin{bmatrix} - h_1 - \\ \vdots \\ - h_n - \end{bmatrix} \in \mathbb{R}^{n \times p}$$

```
                        cos/
H = np.dot(X, G.T) + b.T        n×p
HTH = np.dot(H.T, H)            p×p
HTy = np.dot(H.T, y)
```
$nd$
$+np$
$+p^2$

```
# construct HTH
HTH = np.zeros((p,p))
HTy = np.zeros(p)
for i in range(n):
    hi = np.dot(X[i,:], G.T)+b
    HTH += np.outer(hi, hi)
    HTy += y[i]*hi
    if i % 1000==0: print(i)
```
$p \times p$

$$nd + p^2$$

```
# construct HTH
HTH = np.zeros((p,p))
HTy = np.zeros(p)
block = p
for i in range(int(np.ceil(n/block))+1):
    Hi = np.dot(X[i*block:min(n,(i+1)*block),:], G.T)+b
    HTH += np.dot(Hi.T, Hi)
    HTy += np.dot(Hi.T, y[i*block:min(n,(i+1)*block)])
```

```
w = np.linalg.solve(HTH + lam*np.eye(p), HTy)
```

For each block compute the memory required in terms of n, p, d.

If d << p << n, what is the most memory efficient program (blue, green, red)?

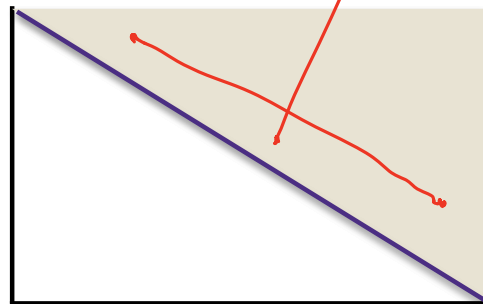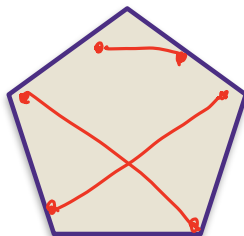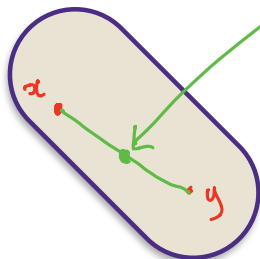If you have unlimited memory, what do you think is the fastest program?

# Convexity

# What is a convex set?
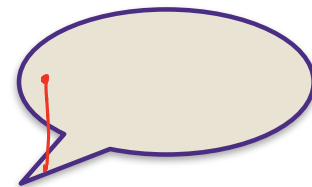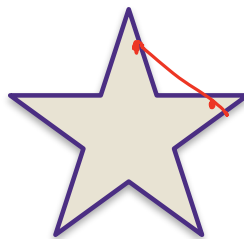
$$K: \{ (x_1, x_2): x_1^2 + x_2^2 \le 10 \}$$
$$(x_1, x_2), (y_1, y_2) \in K$$

A set $K \subset \mathbb{R}^d$ is convex if $(1-\lambda)x + \lambda y \in K$ for all $x, y \in K$ and $\lambda \in [0,1]$

Show $((1-\lambda)x_1 + \lambda y_1, (1-\lambda)x_2 + \lambda y_2)$
$\in K$

## Examples of convex sets



## Examples of non-convex functions: anything else

# What is a convex function?

dom($f$): $\{x : f(x)$ defined$\}$

A function $f : \mathbb{R}^d \to \mathbb{R}$ is convex if $f((1-\lambda)x + \lambda y) \leq (1-\lambda)f(x) + \lambda f(y)$ for all $x, y \in \mathbb{R}$ and $\lambda \in [0,1]$

dom($f$)

## Examples of convex functions: "look like bowls"

$x^2 - x$

$f(x)$  $f(y)$

$x$  $y$

$e^x$

$e^{-x}$

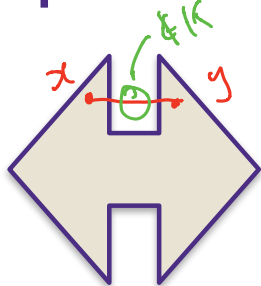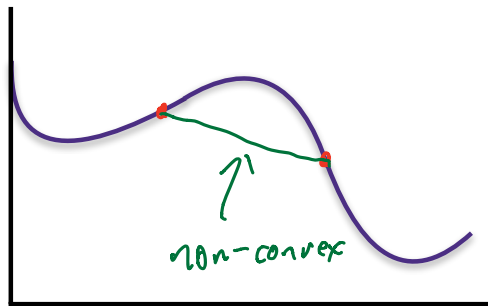## Examples of non-convex functions: anything else

non-convex

# Convex functions and convex sets?

A set $K \subset \mathbb{R}^d$ is convex if $(1 - \lambda)x + \lambda y \in K$ for all $x, y \in K$ and $\lambda \in [0, 1]$

A function $f : \mathbb{R}^d \to \mathbb{R}$ is convex if $f((1 - \lambda)x + \lambda y) \leq (1 - \lambda)f(x) + \lambda f(y)$ for all $x, y \in K$ and $\lambda \in [0, 1]$

A function $f : \mathbb{R}^d \to \mathbb{R}$ is convex if the set $\{(x, t) \in \mathbb{R}^{d+1} : f(x) \leq t\}$ is convex

epigraph($f$)

$\{(x, t) : f(x) \leq t'\}$

$t'$

$f(x)$

$d = 1$

$x$

# More definitions of convexity

A set $K \subset \mathbb{R}^d$ is convex if $(1 - \lambda)x + \lambda y \in K$ for all $x, y \in K$ and $\lambda \in [0, 1]$

A function $f : \mathbb{R}^d \to \mathbb{R}$ is convex if the set $\{(x, t) \in \mathbb{R}^{d+1} : f(x) \leq t\}$ is convex

A function $f : \mathbb{R}^d \to \mathbb{R}$ that is differentiable everywhere is convex if $f(y) \geq f(x) + \nabla f(x)^\top (y - x)$ for all $x, y \in dom(f)$

Try at home
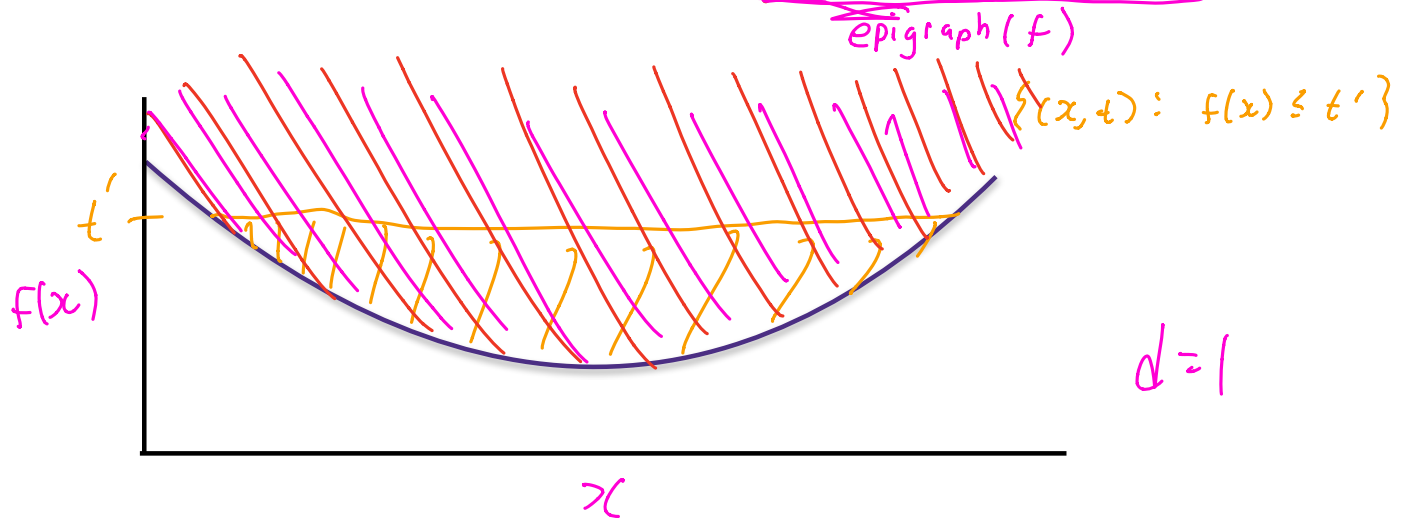$f(x) = e^x$ show convex using $\uparrow$
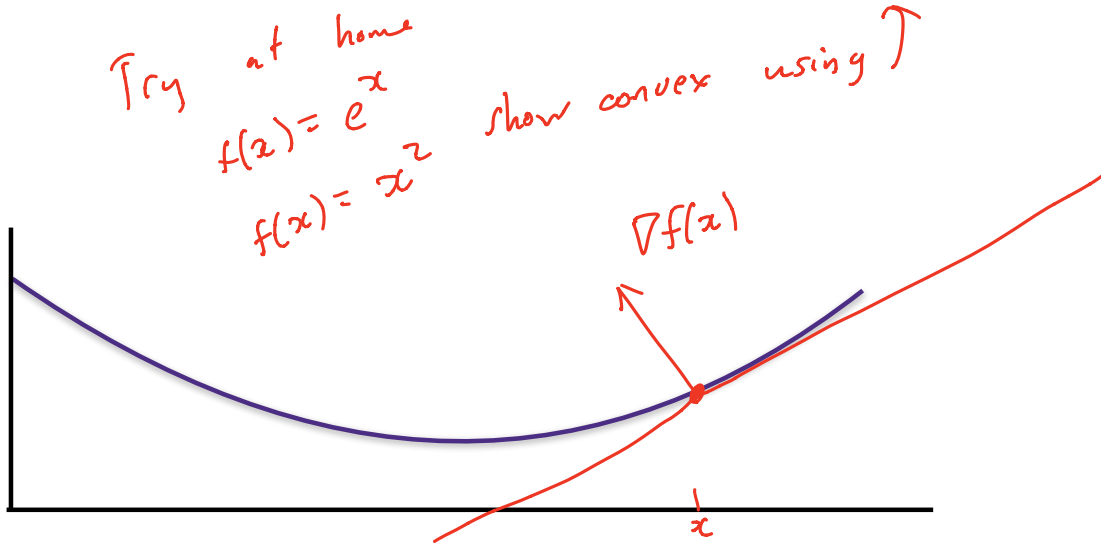$f(x) = x^2$

$\nabla f(x)$

$x$

# More definitions of convexity

A set $K \subset \mathbb{R}^d$ is convex if $(1 - \lambda)x + \lambda y \in K$ for all $x, y \in K$ and $\lambda \in [0, 1]$

A function $f : \mathbb{R}^d \to \mathbb{R}$ is convex if the set $\{(x, t) \in \mathbb{R}^{d+1} : f(x) \leq t\}$ is convex

A function $f : \mathbb{R}^d \to \mathbb{R}$ that is differentiable everywhere is convex if $f(y) \geq f(x) + \nabla f(x)^\top (y - x)$ for all $x, y \in dom(f)$

A function $f : \mathbb{R}^d \to \mathbb{R}$ that is twice-differentiable everywhere is convex if $\nabla^2 f(x) \succeq 0$ for all $x \in dom(f)$

$f(x) = \frac{1}{2}x^2$
$f'(x) = x$
$f''(x) = 1 > 0 \implies f$ is convex

$[\nabla^2 f(x)]_{i,j}$
$= \dfrac{\partial^2 f(x)}{\partial x_i \, \partial x_j}$

$f(x) = e^x$
$f'(x) = e^x$
$f''(x) e^x > 0 \implies e^x$ is convex

$A \succeq 0$ if
$x^\top A x \geq 0 \; \forall x$

# Why do we care about convexity?

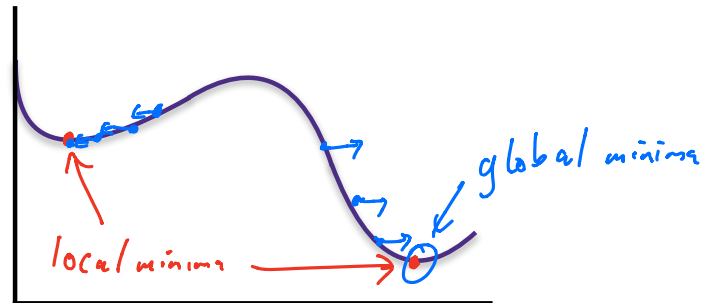Convex functions

- All local minima are global minima
- Efficient to optimize (e.g., gradient descent)

Convex Function

Non-convex Function



local minima        global minima

# Gradient Descent

Initialize: $w_0 = 0$ (or randomly)

for $t = 1, 2, \ldots$

$$w_{t+1} = w_t - \eta \nabla f(w_t)$$

step size

Convex Function

Non-convex Function

# Sub-Gradient Descent

Initialize: $w_0 = 0$

for $t = 1, 2, \ldots$

    Find any $g_t$ such that $f(y) \geq f(w_t) + g_t^\top (y - w_t)$

    $w_{t+1} = w_t - \eta g_t$

$f(x) = |x|$

$f$ is not diff. at $x = 0$

but $g \in [-1, 1]$ is a sub-gradient

for $|x|$ at $0$    since

$|y| \geq |x| + g(y - x)$

sub-gradients are not unique

$g$ is a subgradient at $x$ if $f(y) \geq f(x) + g^T (y - x)$

## Convex Function

Convex fn's have sub-gradients defined everywhere.

f is not differentible

## Non-convex Function

If $f$ is diff. at $x$

then $g = \nabla f(x)$

sub-gradients not necessarily exist for non-convex fns.

# Coordinate descent

Initialize: $w_0 = 0$

for $t = 1, 2, \ldots$

    Let $i_t = (t \% d) + 1$

$$w_{t+1}^{(i_t)} = w_t^{(i_t)} - \eta_t \frac{\partial f(w)}{\partial w^{(i_t)}}\bigg|_{w=w_t}$$



Special case:

$$\text{Choose } \eta_t: \qquad \alpha_t = \underset{\alpha}{\text{argmin}}\; f(w_t + e_{i_t}\alpha)$$

$$w_{t+1}^{(i_t)} = w_t + e_{i_t}\alpha$$

# Machine Learning Problems

- **Given data:**

$$\{(x_i, y_i)\}_{i=1}^{n} \qquad x_i \in \mathbb{R}^d \qquad y_i \in \mathbb{R}$$

each

If $\ell_i(w)$ is convex the sum is convex

- **Learning a model's parameters:**

$$\sum_{i=1}^{n} \ell_i(w)$$

Logistic Loss: $\ell_i(w) = \log(1 + \exp(-y_i\, x_i^T w))$

Squared error Loss: $\ell_i(w) = (y_i - x_i^T w)^2$

Gradient Descent:

iteration
t

$$w_{t+1} = w_t - \eta \nabla_w \left( \frac{1}{n} \sum_{i=1}^{n} \ell_i(w) \right) \Big|_{w=w_t}$$

$$= w_t - \eta\, \frac{1}{n} \sum_{i=1}^{n} \nabla_w \ell_i(w) \Big|_{w=w_t}$$

# Optimization summary

$f(x) = x^3$

$f'(x) = 0$ af $x = 0$

- **You can always run <u>gradient descent</u> whether f is convex or not. But you only have guarantees if f is convex**
- **Many bells and whistles can be added onto gradient descent such as momentum and dimension-specific step-sizes (Nesterov, Adagrad, ADAM, etc.)**

CVX — software package for convex opt

# Stochastic Gradient Descent

# Machine Learning Problems

- **Given data:**

$$\{(x_i, y_i)\}_{i=1}^{n} \qquad x_i \in \mathbb{R}^d \qquad y_i \in \mathbb{R}$$

- **Learning a model's parameters:** $\sum_{i=1}^{n} \ell_i(w)$

Gradient Descent:

$$w_{t+1} = w_t - \eta \nabla_w \left( \frac{1}{n} \sum_{i=1}^{n} \ell_i(w) \right) \bigg|_{w=w_t}$$

What if $n = 10^8$

$$= w_t - \eta \frac{1}{n} \sum_{i=1}^{n} \nabla_w \ell_i(w) \big|_{w=w_t}$$

# Machine Learning Problems

- **Given data:**

$$\{(x_i, y_i)\}_{i=1}^n \qquad x_i \in \mathbb{R}^d \qquad y_i \in \mathbb{R}$$
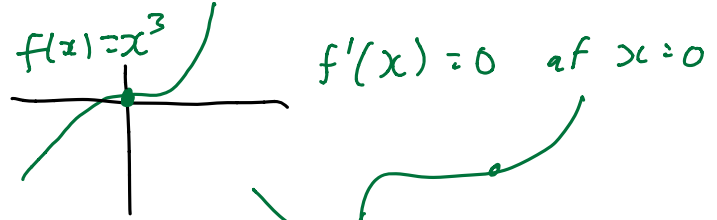
- **Learning a model's parameters:** $\sum_{i=1}^n \ell_i(w)$

Gradient Descent:

$$w_{t+1} = w_t - \eta \nabla_w \left( \frac{1}{n} \sum_{i=1}^n \ell_i(w) \right)\Big|_{w=w_t}$$

Stochastic Gradient Descent:

$$w_{t+1} = w_t - \eta \nabla_w \ell_{I_t}(w)\Big|_{w=w_t}$$

$I_t$ drawn uniform at random from $\{1, \ldots, n\}$

$$\mathbb{E}[\nabla \ell_{I_t}(w)] = \sum_{i=1}^n \underbrace{\mathbb{P}(I_t = i)}_{=\frac{1}{n}} \nabla \ell_i(\omega) = \frac{1}{n} \sum_{i=1}^n \nabla \ell_i(\omega) =$$

# Machine Learning Problems
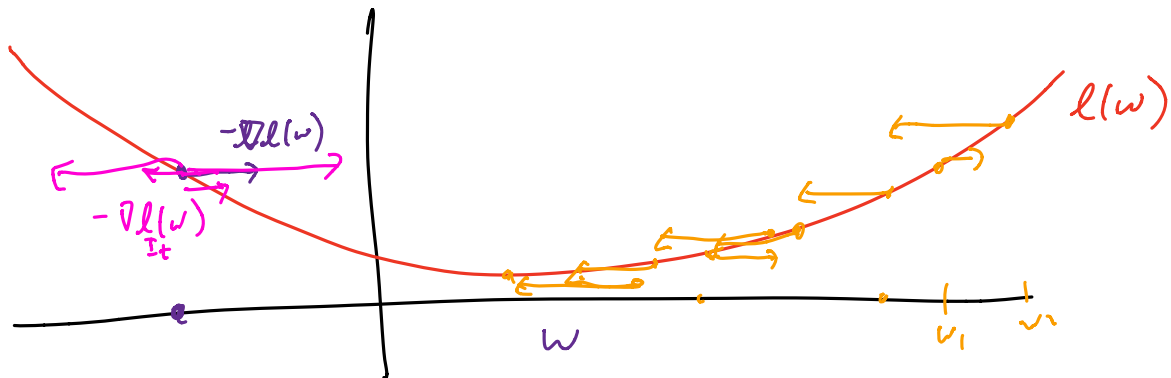
- **Learning a model's parameters:** $\sum_{i=1}^{n} \ell_i(w)$

Stochastic Gradient Descent:

$$w_{t+1} = w_t - \eta \nabla_w \ell_{I_t}(w)\Big|_{w=w_t}$$

$I_t$ drawn uniform at random from $\{1, \dots, n\}$

$$\mathbb{E}[\nabla \ell_{I_t}(w)] = \nabla \ell(w)$$

# Stochastic Gradient Descent

$f(w_T) - f(w_*) \leq \varepsilon$  flops to reach $\varepsilon$-good soln

GD "typically" converges

like $e^{-T}$

fewer iters  $n \log(1/\varepsilon)$

indep. of $n \longrightarrow 1/\varepsilon$

way smaller cost per iteration  SGD

"typically" converges like $\frac{1}{T}$ or $\frac{1}{\sqrt{T}}$

## Theorem

Let $\boxed{w_{t+1} = w_t - \eta \nabla_w \ell_{I_t}(w)\Big|_{w=w_t}}$   $I_t$ drawn uniform at random from $\{1, \ldots, n\}$   so that

$$\mathbb{E}\big[\nabla \ell_{I_t}(w)\big] = \frac{1}{n}\sum_{i=1}^{n}\nabla \ell_i(w) =: \nabla \ell(w)$$

If $\|w_* - w_0\|_2^2 \leq R$   and   $\sup_w \max_i \|\nabla \ell_i(w)\|_2 \leq G$   then

$$\mathbb{E}[\ell(\bar{w}) - \ell(w_*)] \leq \frac{R}{2T\eta} + \frac{\eta G}{2} \leq \sqrt{\frac{RG}{T}} \qquad \eta = \sqrt{\frac{R}{GT}}$$

$$\bar{w} = \frac{1}{T}\sum_{t=1}^{T} w_t$$

(In practice use last iterate)

# Stochastic Gradient Descent

Proof

$$\mathbb{E}[||w_{t+1} - w_*||_2^2] = \mathbb{E}[||w_t - \eta \nabla \ell_{I_t}(w_t) - w_*||_2^2]$$

# Stochastic Gradient Descent

Proof

$$\mathbb{E}[||w_{t+1} - w_*||_2^2] = \mathbb{E}[||w_t - \eta\nabla\ell_{I_t}(w_t) - w_*||_2^2]$$

$$= \mathbb{E}[||w_t - w_*||_2^2] - 2\eta\mathbb{E}[\nabla\ell_{I_t}(w_t)^T(w_t - w_*)] + \eta^2\mathbb{E}[||\nabla\ell_{I_t}(w_t)||_2^2]$$

$$\leq \mathbb{E}[||w_t - w_*||_2^2] - 2\eta\mathbb{E}[\ell(w_t) - \ell(w_*)] + \eta^2 G$$

$$\mathbb{E}[\nabla\ell_{I_t}(w_t)^T(w_t - w_*)] = \mathbb{E}\big[\mathbb{E}[\nabla\ell_{I_t}(w_t)^T(w_t - w_*)|I_1, w_1, \ldots, I_{t-1}, w_{t-1}]\big]$$

$$= \mathbb{E}\big[\nabla\ell(w_t)^T(w_t - w_*)\big]$$

$$\geq \mathbb{E}\big[\ell(w_t) - \ell(w_*)\big] \qquad \swarrow \quad \text{convexity}$$

$$\sum_{t=1}^{T}\mathbb{E}[\ell(w_t) - \ell(w_*)] \leq \frac{1}{2\eta}\big(\mathbb{E}[||w_1 - w_*||_2^2] - \mathbb{E}[||w_{T+1} - w_*||_2^2] + T\eta^2 G\big)$$

$$\leq \frac{R}{2\eta} + \frac{T\eta G}{2}$$

# Stochastic Gradient Descent

**Jensen's inequality**:
For any random $Z \in \mathbb{R}^d$ and convex function $\phi : \mathbb{R}^d \to \mathbb{R}$, $\phi(\mathbb{E}[Z]) \leq \mathbb{E}[\phi(Z)]$

$$\mathbb{E}[\ell(\bar{w}) - \ell(w_*)] \leq \frac{1}{T} \sum_{t=1}^{T} \mathbb{E}[\ell(w_t) - \ell(w_*)] \qquad \bar{w} = \frac{1}{T} \sum_{t=1}^{T} w_t$$

# Stochastic Gradient Descent

Proof

**Jensen's inequality**:
For any random $Z \in \mathbb{R}^d$ and convex function $\phi : \mathbb{R}^d \to \mathbb{R}$, $\phi(\mathbb{E}[Z]) \leq \mathbb{E}[\phi(Z)]$

$$\mathbb{E}[\ell(\bar{w}) - \ell(w_*)] \leq \frac{1}{T} \sum_{t=1}^{T} \mathbb{E}[\ell(w_t) - \ell(w_*)]$$

$$\bar{w} = \frac{1}{T} \sum_{t=1}^{T} w_t$$

$$\mathbb{E}[\ell(\bar{w}) - \ell(w_*)] \leq \frac{R}{2T\eta} + \frac{\eta G}{2} \leq \sqrt{\frac{RG}{T}}$$

$$\eta = \sqrt{\frac{R}{GT}}$$

# Mini-batch SGD

Instead of one iterate, average B stochastic gradient together

Advantages:
- de-noises gradient   reducing variance of gradient
- Matrix computations ← Computing B gradients at a time
- Parallelization

$\uparrow$

Compute gradients by sending a partition
of the batch to each computer

# Stochastic Gradient Descent: A Learning perspective

# Learning Problems as Expectations

> Minimizing loss in training data:
- Given dataset:
  > Sampled iid from some distribution p(**x,y**) on features:
- Loss function, e.g., hinge loss, logistic loss,…
- We often minimize loss in training data:

$$\ell_{\mathcal{D}}(\mathbf{w}) = \frac{1}{N} \sum_{j=1}^{N} \ell(\mathbf{w}, \mathbf{x}^j)$$

> However, we should really minimize expected loss on all data:

$$\ell(\mathbf{w}) = E_{\mathbf{x}}\left[\ell(\mathbf{w}, \mathbf{x})\right] = \int p(\mathbf{x})\ell(\mathbf{w}, \mathbf{x})d\mathbf{x}$$

> So, we are approximating the integral by the average on the training data

# Gradient descent in Terms of Expectations

> **"True" objective function:**

$$\ell(\mathbf{w}) = E_{\mathbf{x}}\left[\ell(\mathbf{w}, \mathbf{x})\right] = \int p(\mathbf{x})\ell(\mathbf{w}, \mathbf{x})d\mathbf{x}$$

> **Taking the gradient:**

> **"True" gradient descent rule:**

> **How do we estimate expected gradient?**

# Warm up - Revisited

$$w_{t+1} = w_t - \frac{2}{2} \nabla \left( y_{I_t} - h_{g_t}(x_{I_t}) \right)^2$$

$$nd + p + pd$$

$$I_t \sim [p]$$

$$I_t \sim [n]$$

1 float in NumPy = 8 bytes
$10^6 \approx 2^{20}$ bytes = 1 MB
$10^9 \approx 2^{30}$ bytes = 1 GB

```python
# generate some nonsense data for an example
X = np.random.randn(n,d)
y = np.random.randn(n)
```

```python
# generate the random features
G = np.random.randn(p, d)*np.sqrt(.1)
b = np.random.rand(p)*2*np.pi
```

```python
H = np.dot(X, G.T) + b.T
HTH = np.dot(H.T, H)
HTy = np.dot(H.T, y)
```

```python
# construct HTH
HTH = np.zeros((p,p))
HTy = np.zeros(p)
for i in range(n):
    hi = np.dot(X[i,:], G.T)+b
    HTH += np.outer(hi, hi)
    HTy += y[i]*hi
    if i % 1000==0: print(i)
```

```python
# construct HTH
HTH = np.zeros((p,p))
HTy = np.zeros(p)
block = p
for i in range(int(np.ceil(n/block))+1):
    Hi = np.dot(X[i*block:min(n,(i+1)*block),:], G.T)+b
    HTH += np.dot(Hi.T, Hi)
    HTy += np.dot(Hi.T, y[i*block:min(n,(i+1)*block)])
```

```python
w = np.linalg.solve(HTH + lam*np.eye(p), HTy)
```

For each block compute the memory required in terms of n, p, d.

If d << p << n, what is the most memory efficient program (blue, green, red)?

If you have unlimited memory, what do you think is the fastest program?