- SVD  + applications
- least squares
- perceptron alg.

- PAC learning
- Gradient descent & SGD
- Linear programming (maybe)

**Setting:**       Supervised learning setting

classify email msgs $\longrightarrow$ spam
$\searrow$ not spam

Take a sample of msgs, labelled according to spam Y/N.

**Goal:** given labelled sample, come up with a good rule for classifying future msgs
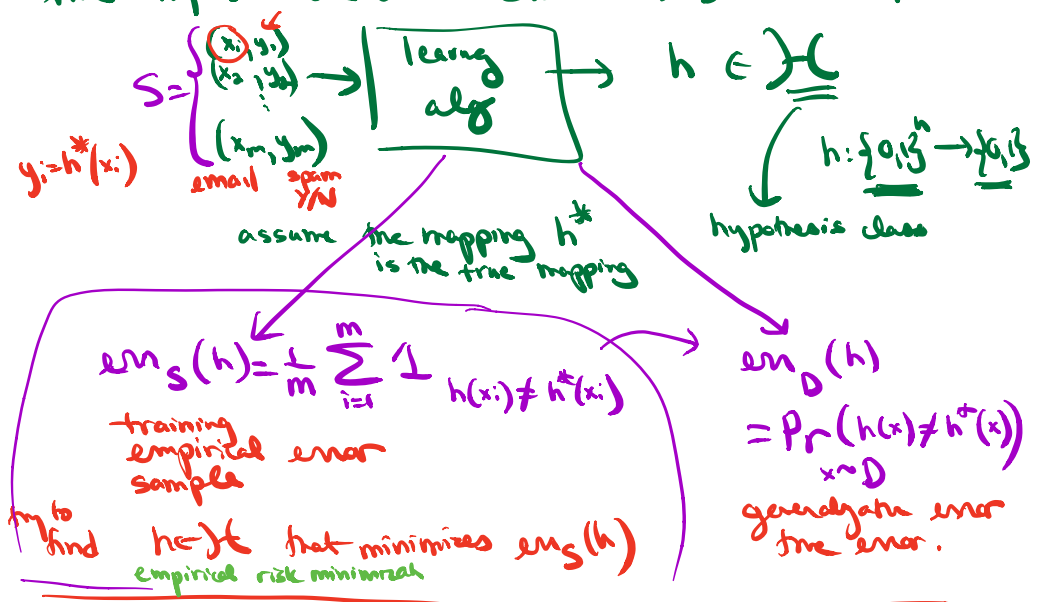
$$h^* : \{0,1\}^n \longrightarrow \{0,1\}$$
true label

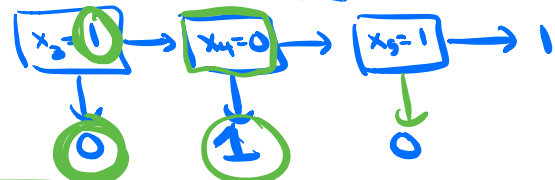| | money | pills | Mr. | bad spelling | known-sender | spam? |
|---|---|---|---|---|---|---|
| 1 | Y | N | Y | Y | Y | Y |
| 2 | Y | N | N | Y | N | N |
| 3 | N | N | N | N | Y | Y |
| 4 | Y | Y | N | N | Y | N |
| 5 | N | N | Y | N | Y | N |
| 6 | Y | N | Y | N | N | Y |
| 7 | N | N | N | N | N | N |
| 8 | N | Y | N | Y | N | Y |

$h^*(x_1)$
$h^*(x_2)$

return SPAM if ¬known and (money or pills)

① distn over inputs $x \in X$

each sample $x_i$ is drawn indep from ①
see $m$ samples

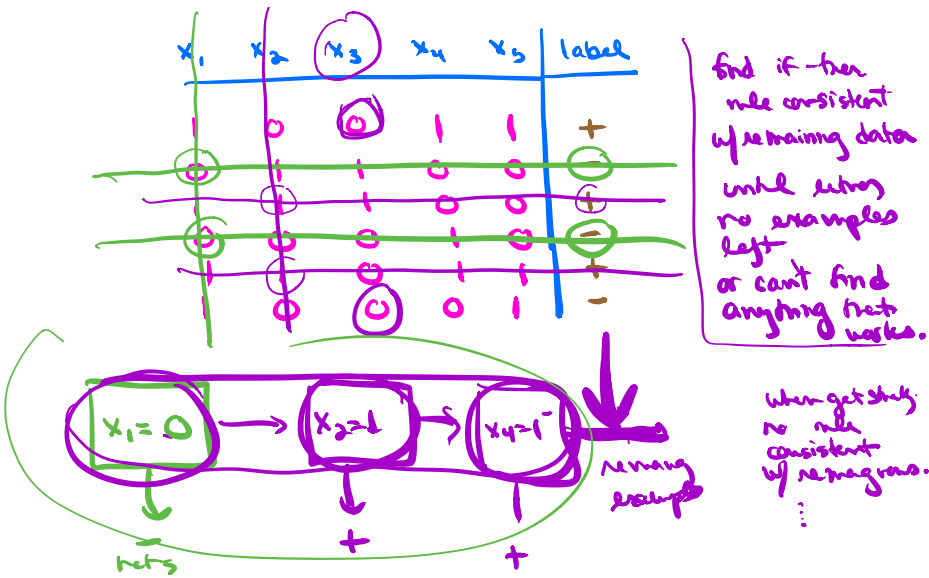future inputs are also drawn from same distn

$$S = \left\{ \begin{array}{c} (x_1, y_1) \\ (x_2, y_2) \\ \vdots \\ (x_m, y_m) \end{array} \right\} \rightarrow \boxed{\begin{array}{c} \text{learning} \\ \text{alg} \end{array}} \rightarrow h \in \mathcal{H}$$

$y_i := h^*(x_i)$

email spam Y/N

$h : \{0,1\}^n \rightarrow \{0,1\}$

hypothesis class

assume the mapping $h^*$ is the true mapping

$$err_S(h) = \frac{1}{m} \sum_{i=1}^{m} \mathbb{1}_{h(x_i) \neq h^*(x_i)}$$

training empirical error sample

try to find $h \in \mathcal{H}$ that minimizes $err_S(h)$
empirical risk minimizer

$$err_D(h) = \Pr_{x \sim D}(h(x) \neq h^*(x))$$

generalization error true error.

$\mathcal{H}$ : decision lists   $\vec{x} \in \{0,1\}^n = (x_1, x_2, \ldots, x_n)$

$$\boxed{x_2 = 1} \rightarrow \boxed{x_4 = 0} \rightarrow \boxed{x_g = 1} \rightarrow 1$$

$\downarrow$ 0   $\downarrow$ 1   $\downarrow$ 0

$$\boxed{|\mathcal{H}| = n! \cdot (2 \cdot 2)^n = n! \, 4^n}$$

| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | label |
|-------|-------|-------|-------|-------|-------|
|       | 0     | 0     | 1     | 1     | +     |
| 0     | 1     | 1     | 0     | 0     | +     |
| 0     | 1     | 0     | 1     | 0     | +     |
| 0     | 1     | 0     | 1     | 1     | +     |
| 0     | 1     | 0     | 0     | 1     | −     |

find if-then rule consistent w/ remaining data until either no examples left or can't find anything that works.



$x_1 = 0 \rightarrow x_2 = 1 \rightarrow x_4 = 1 \rightarrow$ remaining examples

$-$ nets    $+$    $+$

When get stuck, no rule consistent w/ remaining ex. ...

have a nice alg for finding consistent DL if such exists.

Confidence, generalization — Claim: if $|S|$ was $\geq \boxed{\phantom{xx}}$ then w.h.p. $err_D(h)$ small.

Consider some DL $h \in \mathcal{H}$ that $\boxed{err_D(h) \geq \varepsilon}$

$h$ misclass has $\geq \varepsilon$    $h$ classifies correctly

Prob( h was consistent w/ our sample)

$$\leq (1-\varepsilon)^{|S|}$$

$\mathcal{H}$

$$Pr\left(\exists\ DL\ h\ \text{with}\ err_D(h) \geq \varepsilon\ \text{but}\ err_S(h) = 0\right)$$

$$\leq \sum_{\substack{h \in \mathcal{H} \\ s.t.\ err_D(h) \geq \varepsilon}} Pr\left(err_S(h) = 0\right) \leq |\mathcal{H}|(1-\varepsilon)^{|S|}$$

$(1-x) \leq e^{-x}$

$$|\mathcal{H}| e^{-\varepsilon|S|}$$

How big does S need to be so that $> \delta$

$\varepsilon, \delta$

$$|\mathcal{H}|(1-\varepsilon)^{|S|} \leq \delta$$

$$n!\,4^n\, e^{-\varepsilon|S|} \leq \delta$$

$$\boxed{\frac{n^n 4^n}{} \leq e^{\varepsilon|S|}}$$

$$n \ln n + n \ln H + \ln\left(\tfrac{1}{\delta}\right) \leq \varepsilon |S|$$

$$\frac{2}{\varepsilon}\left(n \ln n + \ln\left(\tfrac{1}{\delta}\right)\right) \leq |S|$$

$\leftarrow \ln |\mathcal{H}|$

0.01        0.01

if $|S| = \Omega\left(\tfrac{1}{\varepsilon} n \ln n + \tfrac{1}{\delta}\right)$ then

$\Pr\left( \exists h \in \mathcal{H} \text{ s.t. } \text{err}_D(h) \geq \varepsilon \text{ and } \text{err}_S(h) = 0 \right) \leq \delta$

If we can find $h \in \mathcal{H}$ that is consistent w/ sample

then rule $h$ we find

is        probably        approximately        correct

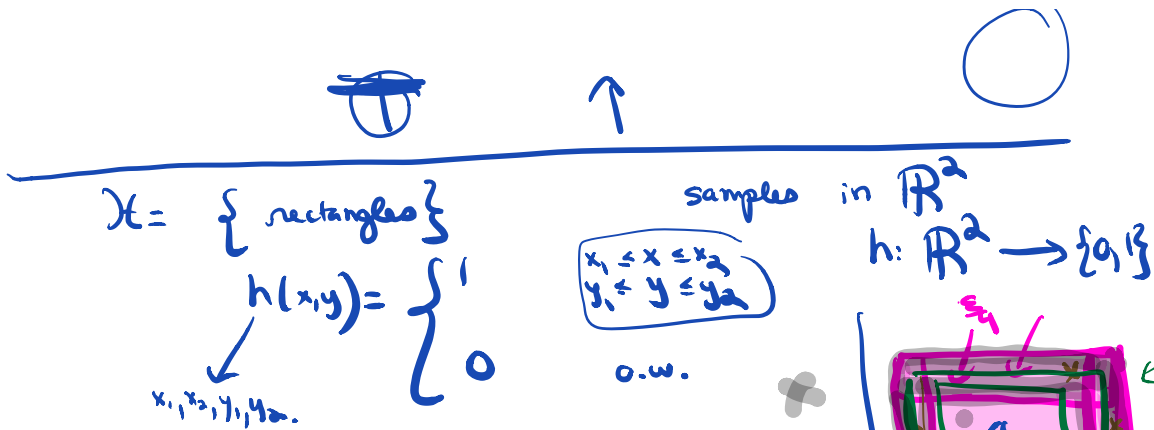$\geq 1 - \delta$        error $\leq \varepsilon$        Turing award.

$\boxed{PAC - \text{learning}}$        Leslie Valiant.

$\mathcal{H}$

If $\boxed{|S| \geq \tfrac{1}{\varepsilon}\left(\ln|\mathcal{H}| + \ln\left(\tfrac{1}{\delta}\right)\right)}$        Sample complexity

then w prob $\geq 1 - \delta$, any $h \in \mathcal{H}$

that has $\text{err}_D(h) \geq \varepsilon$        will have $\text{err}_S(h) > 0$

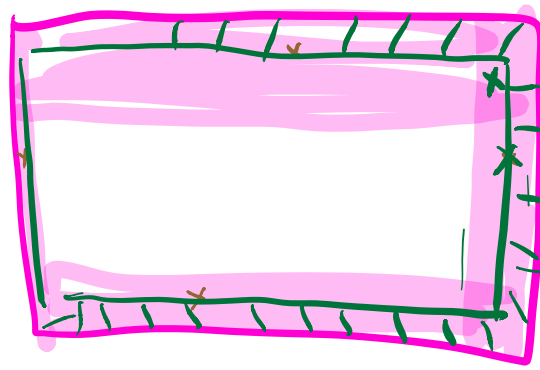assumed $\exists h \in \mathcal{H}$ s.t. $\text{err}_S(h) = 0$

Thm:

If $\boxed{|S| \geq \tfrac{1}{2\varepsilon^2}\left(\ln|\mathcal{H}| + \ln\left(\tfrac{1}{\delta}\right)\right)}$        Sample complexity

then w prob $\geq 1 - \delta$,        for every $h \in \mathcal{H}$

$$|\text{err}_S(h) - \text{err}_D(h)| \leq \varepsilon$$

$h^*$

$\mathcal{H} = \{$ rectangles $\}$

$h(x,y) = \begin{cases} 1 \\ \\ 0 \end{cases}$

$x_1 \leq x \leq x_2$
$y_1 \leq y \leq y_2$

o.w.

$x_1, x_2, y_1, y_2$.

samples in $\mathbb{R}^2$

$h: \mathbb{R}^2 \rightarrow \{0,1\}$

$\frac{\varepsilon}{4}$ $\frac{\varepsilon}{4}$

$1$

$\geq \frac{\varepsilon}{4}$

Suppose that $\exists$ perfect classifier in $\mathcal{H}$

green inside pink

$0$

output smallest possible bounding rectangle.

If $\exists$ sample in each 4 little pink rectangles, then $\Pr($ error on random draw $\leq \varepsilon)$

$\Pr($ green rectangle has error $\geq \varepsilon)$

if get a sample in each of 4 pink subrectangles then $\Pr($ making a mistake$)$

$\leq 4 \cdot \frac{\varepsilon}{4} \leq \varepsilon$

$\Pr_u \frac{\varepsilon}{4}$

$\Pr\left(\text{err}_D \left(\text{smallest bounding rectangle}\right) > \varepsilon\right)$

$\leq \Pr(\exists$ pink subrectangle w/ no sample in it$)$

$$\leq 4 \cdot \Pr(\text{no sample in particular sub rectangle})$$

$$= 4\left(1 - \frac{\epsilon}{4}\right)^{|S|} \leq \delta$$

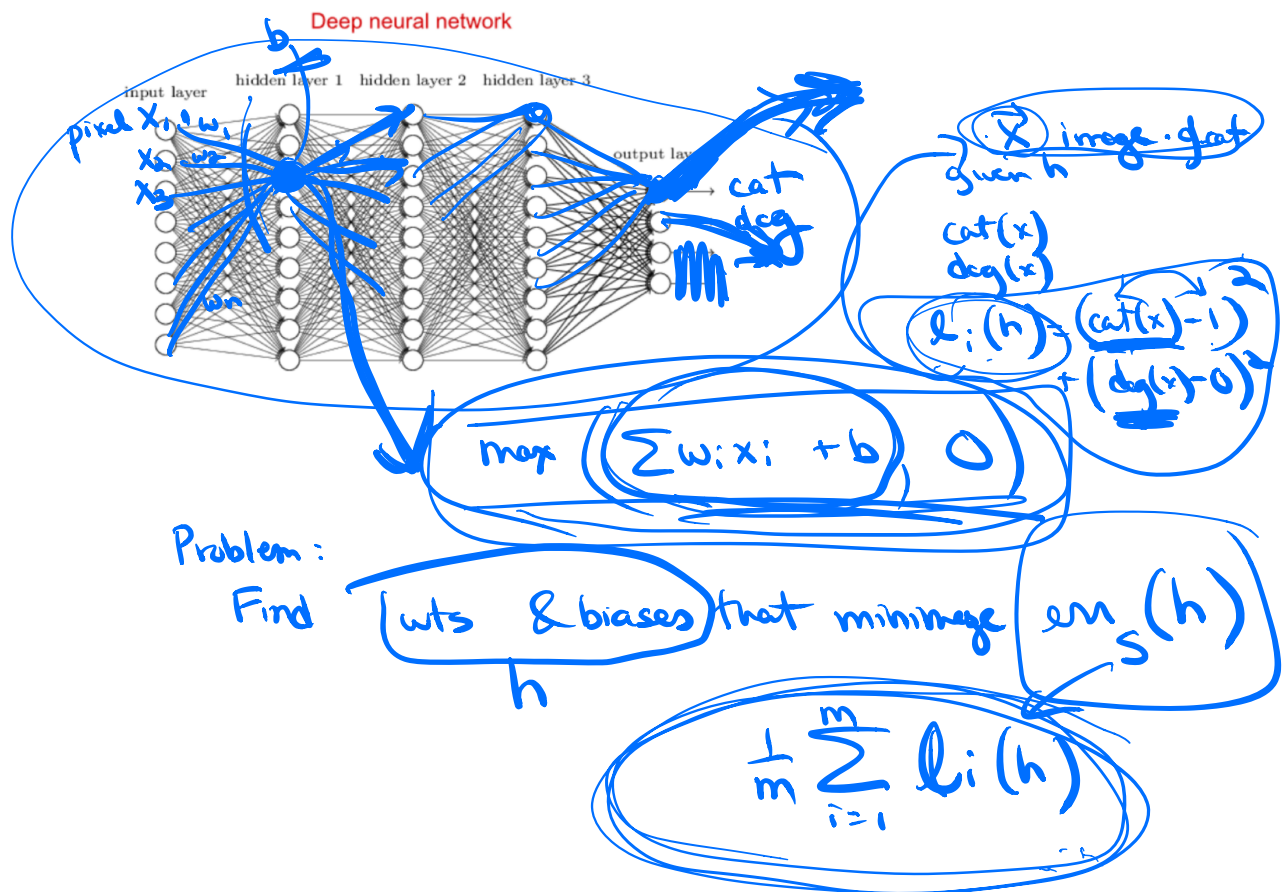$$\boxed{\frac{4}{\epsilon} \ln\left(\frac{4}{\delta}\right) \leq |S|}$$

---

find $h \in \mathcal{H}$ to minimize $\text{err}_S(h)$
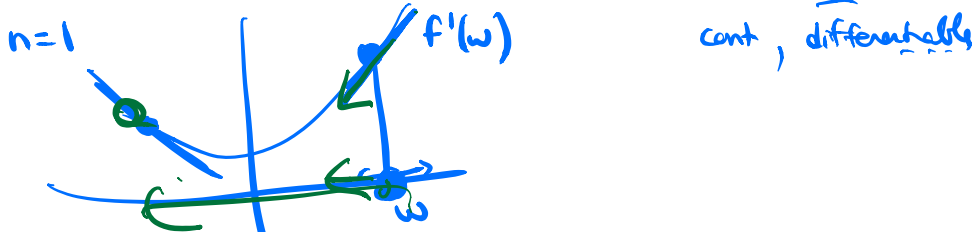
optimization problem. $\Rightarrow$ $= \left(\frac{1}{m}\right) \sum_{i=1}^{m} \text{loss}(h, i\text{th sample})$

$\ell_i(h)$



Deep neural network

input layer   hidden layer 1   hidden layer 2   hidden layer 3

pixel $x_1$, $w_1$
$x_2$, $w_2$
$x_3$
$w_n$

output layer

cat
dog

$\hat{x}$ image of cat
given $h$

$\text{cat}(x)$
$\text{dog}(x)$

$\ell_i(h) = (\text{cat}(x) - 1)^2 + (\text{dog}(x) - 0)^2$

$\max\left(\sum w_i x_i + b, \; 0\right)$

Problem:
Find $\boxed{\text{wts \& biases}}$ that minimize $\text{err}_S(h)$
$h$

$\frac{1}{m} \sum_{i=1}^{m} \ell_i(h)$

# Gradient descent

method for "trying" to minimize a fn. $f: \underline{\underline{\mathbb{R}^n}} \to \mathbb{R}$

cont, differentiably

$n=1$

$f'(w)$

$f(w+s) \approx f(w) + s \cdot f'(w)$

$$\frac{f(w+s) - f(w)}{s} \approx f'(s)$$
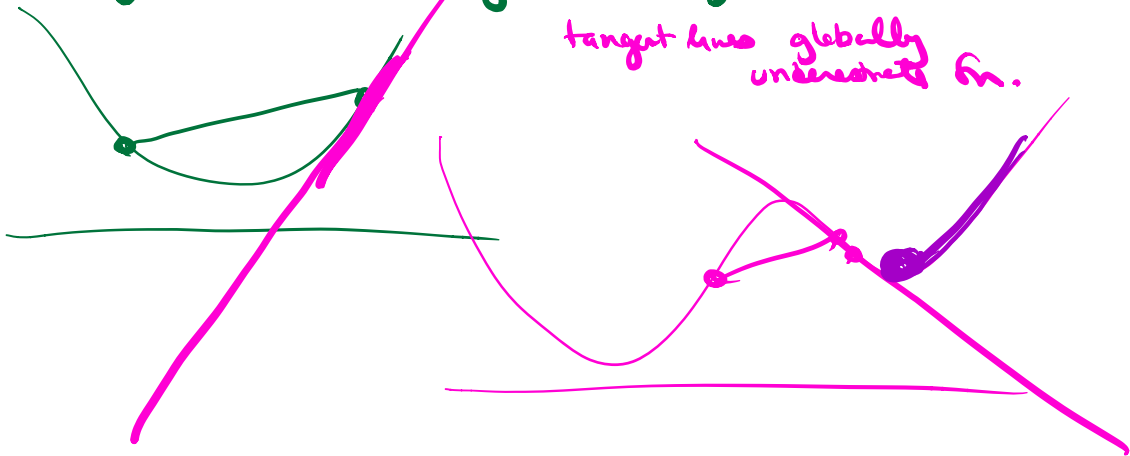
$> 0 \implies s < 0$

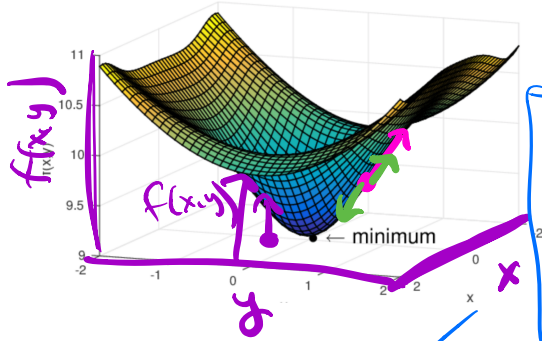$< 0 \implies s > 0$

$w_0 :=$ arbitrary

for $t = 1, \ldots$ "done"

$w_{t+1} := w_t - \eta_t f'(w)$

para

takes opposite direction of derivative

If fn is convex then w/ appropriately selected values of $\eta_t$ guaranteed to converge to global min

tangent lines globally underestimate fn.

$f : \mathbb{R}^n \to \mathbb{R}$

$f(x,y) = x^2 + 2xy + 4y^2$

$\frac{\partial f(x,y)}{\partial x} = 2x + 2y$

$\frac{\partial f(x,y)}{\partial y} = 2x + 8y$

$(1,1)$

$f(\vec{w} + \vec{s}) \approx f(\vec{w}) + \vec{s} \cdot \nabla f(\vec{w})$

$(w_1, \ldots, w_n) \quad (s_1, s_2, \ldots, s_n)$

$(w_1 + s_1, w_2 + s_2 \ldots)$

gradient of f $\begin{pmatrix} \frac{\partial f(\vec{w})}{\partial x_1} \\ \frac{\partial f}{\partial s_2}(\vec{w}) \\ \vdots \end{pmatrix}$

$\nabla f(1,1) = \begin{pmatrix} 4 \\ 10 \end{pmatrix}$

$f(1 + s_1, 1 + s_2) \approx f(1,1) + 4s_1 + 10s_2$

$(s_1, s_2) \quad \begin{pmatrix} 4 \\ 10 \end{pmatrix}$

$\vec{s} \cdot \nabla f$

$= \|\vec{s}\| \|\nabla f\| \cos\theta$

$-1$

$\nabla f$

want to move
in direction opposite

$$\vec{w}_{t+1} = \vec{w}_t + \eta_t \nabla f(w_t)$$

direction of negative gradient.

gradient descent.

$err_S(w_1, \ldots, w_n)$

$f: \mathbb{R}^n \to \mathbb{R}$

Cost of single update

fn of $n$ ; $m$  #-sample pts

fn we're trying to minimize

$$\frac{1}{m} \sum_{i=1}^{m} \ell_i(\vec{w})$$

error on $i^{th}$ sample pt.

$$\nabla f = \frac{1}{m} \sum_{i=1}^{m} \nabla \ell_i(w)$$

$n$  #wts & biases
$m$  #images

at each $t$,
pick one random image
uniformly at random

$I_t \in \{1, \ldots, m\}$

$$E\left( \nabla \ell_{I_t}(w) \right) = \sum_{i=1}^{m} Pr(\text{select } i) \nabla \ell_i(w)$$

$$= \frac{1}{m} \sum_{i=1}^{m} \nabla \ell_i(w)$$

$$= \nabla \text{ loss fn}$$

Define $w^*$ to be min $g$ fn $\left(err_S(\vec{w})\right)$

$$\mathbb{E}\left[err_S(\vec{w}) - err_S(w^*)\right] \leq \sqrt{\frac{RG}{T}} \text{ consts.}$$

run SGD for $T$ steps

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

$\leq \varepsilon$

$\|w_0 - w^*\|^2 \leq R$

$\max \|\nabla \ell_i(w)\|^2 \leq G$