

# Introduction to Machine Learning

## TTIC 31020

Prof. Nati Srebro

### Lecture 9:

Part I: Gaussian Process and Bayesian Machine Learning

Part II: Stochastic Gradient Descent

### Generative Mode:

Assumptions on  $P(x|y)$   
 $P(x, y) \in \{p(x, y; \theta) \mid \theta\}$

### Discriminative Probability Model:

Assumption on  $P(y|x)$

$$P(y|x) = \frac{1}{1+e^{h(x)}} \quad h \in \{h_w(x) \mid w\}$$
$$P(y|x) \in \left\{ p(y|x; w) = \text{Ber} \left( \frac{1}{1+e^{h_w(x)}} \right) \mid w \right\}$$

### Predictive Mode:

Assumption on  $h^*$   
 $L_S(h)$  small for some  
 $h \in \{h_w(x) \mid w\}$

More knowledge  
Less data

Less knowledge  
More data

Fit Generative Model, e.g. via

Max Likelihood:  $\arg \max_{\theta} \log P(\{(x_i, y_i)\} \mid \theta)$

or MAP:  $\arg \max_{\theta} \log P(\{(x_i, y_i)\} \mid \theta) + \log P(\theta)$

Task loss, e.g.  $L_S^{01}$ , or conv/tractable surrogate

ERM:  $\arg \min_w L_S(h_w)$

Regularized ERM, e.g.:  $\arg \min_w L_S(h_w) + \lambda \|w\|_2^2$

Scale sensitive, e.g.  $L_S^{marg}$ ,  $L_S^{hinge}$ ,  $L_S^{lgstc}$

Fit Discriminative Model, e.g. via

Max *conditional* Likelihood  $\arg \max_w \log P(\{y_i\} \mid \{x_i\}; w) = \arg \min_w L_S^{lgst}(h_w)$

or conditional MAP:  $\arg \max_w \log P(\{y_i\} \mid \{x_i\}; w) + \log P(w) = \arg \min_w L_S^{lgst}(h_w) + \frac{1}{2m\sigma^2} \|w\|_2^2$

For Gaussian prior,  $w \sim N(0, \sigma^2 I)$

$$P(\{y_i\}|\{x_i\}, w) = \prod \frac{1}{1+e^{-h_w(x_i)}} \quad h_w(x_i) = \langle w, \phi(x_i) \rangle \quad w \sim \mathcal{N}(0, I_D)$$

Can we talk about  $P(\{y_i\}|\{x_i\})$ , or  $h(x)$ , and without talking about  $w$  or  $\phi$ , allowing  $D = \infty$ ?

$h_w(\cdot)$  is a random function. How is it distributed?

How is  $h_w(x_1), h_w(x_2), \dots, h_w(x_m)$  distributed for some  $x_1, \dots, x_m$ ?

$$[h_w(x_1), h_w(x_2), \dots, h_w(x_m)] = [\langle w, \phi(x_1) \rangle, \langle w, \phi(x_2) \rangle, \dots, \langle w, \phi(x_m) \rangle] = w\Phi \sim \mathcal{N}(0, \Phi^T \Phi)$$

Or in other words,  $h_w(\cdot)$  is a **Gaussian Process**

(a random function where every set of values is jointly Gaussian)

A GP is characterized by its Covariance function (and mean, which we take to be  $\mathbb{E}[r(x_i)] = 0$ ):

$$K(x_i, x_j) = \text{Cov}(h(x_i), h(x_j)) = \text{Cov}(\langle w, \phi(x_1) \rangle, \langle w, \phi(x_2) \rangle) = \phi(x_1)^T \text{Cov}(w) \phi(x_2) = \langle \phi(x_1), \phi(x_2) \rangle$$

Gaussian Processes  $\approx$  Kernelized  $\ell_2$ -regularized linear learning

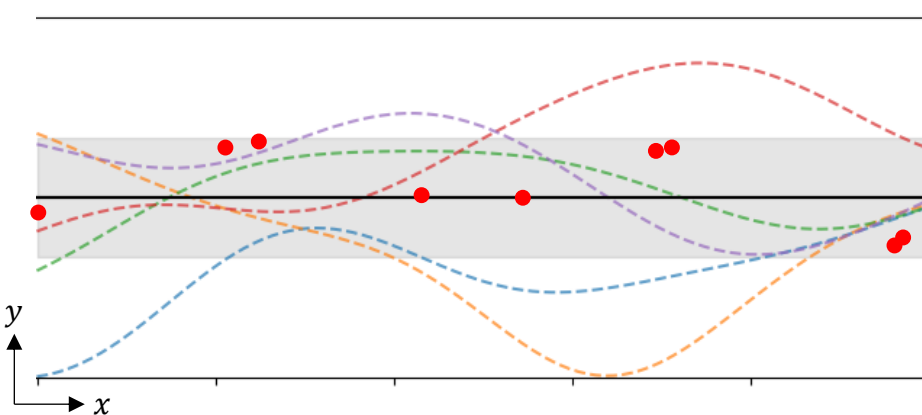
In this case  $\mathbb{E}_{h|S}[h] = \arg \max_h P(h|S)$ , since mode of Gaussian is its mean

Another view: prior over functions given by Hilbert norm in a Reproducing Kernel Hilbert Space (RKHS)

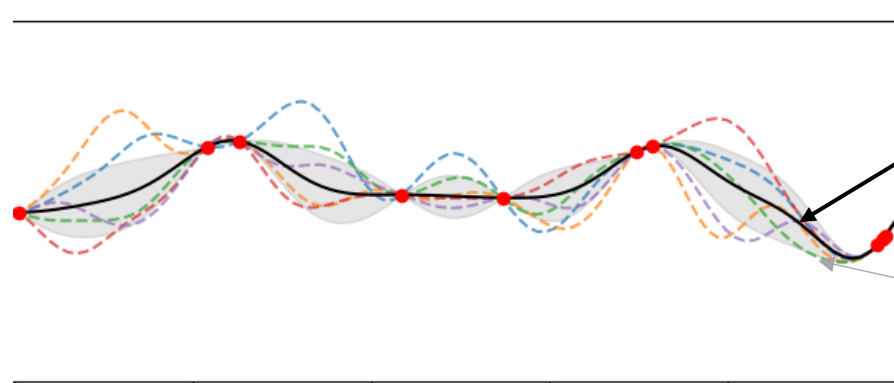
$$p(h) \propto \exp(-\|h\|_K^2)$$

**Beyond scope of course**

$$K(x, x') = e^{-|x-x'|^2}$$



Prior  $P(h)$

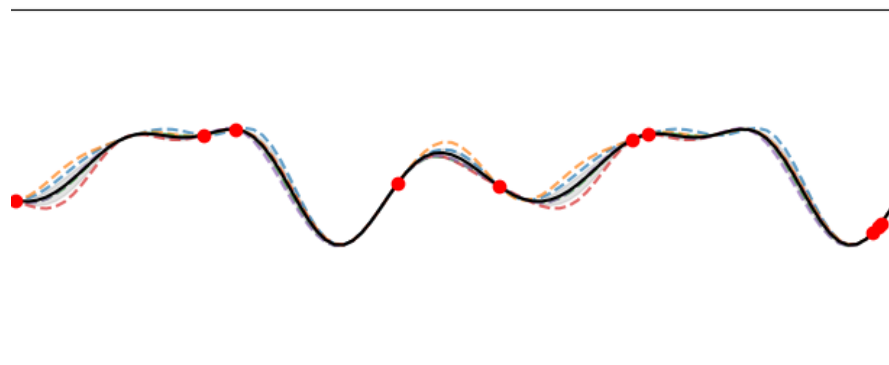
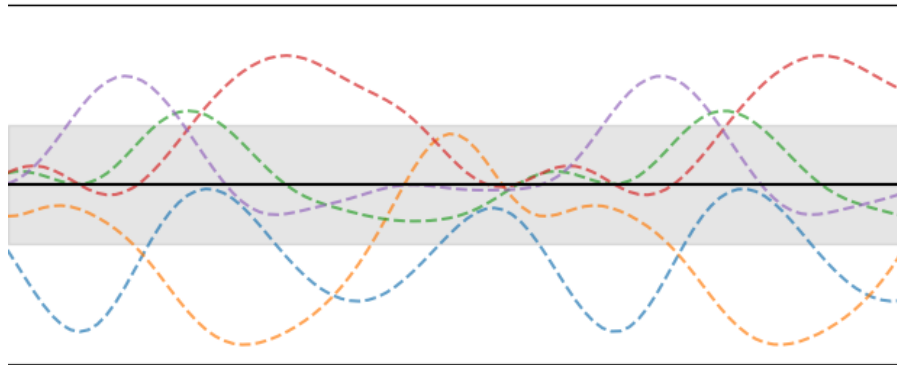


Posterior  $P(h|S)$

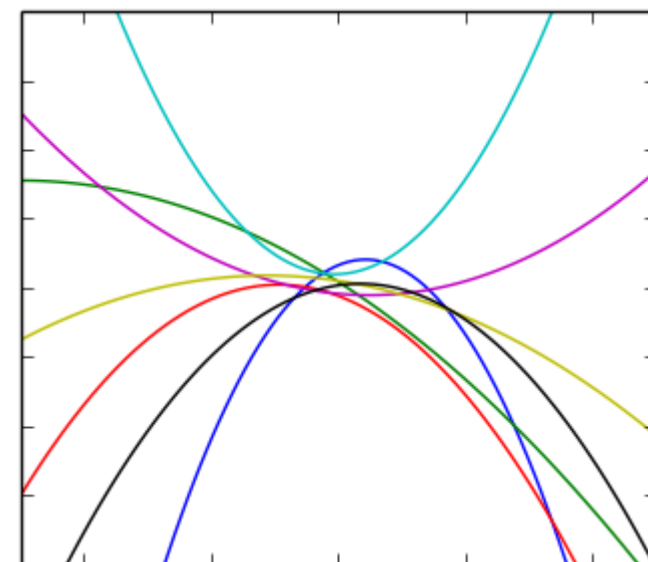
$$\begin{aligned} \text{MAP: } \max_h P(h|L_S(h) = 0) \\ = \arg \max \|w\| \text{ s.t. } L_S(h) = 0 \\ = \mathbb{E}[h(x)|L_S(h) = 0] \end{aligned}$$

$$\sqrt{\text{Var}[h(x)|S]}$$

$$K(x, x') = e^{-\sin^2(x-x')}$$



$$K(x, x') = (xx' + 1)^2$$



Beyond scope of course

# Advantages of Probabilistic Approach

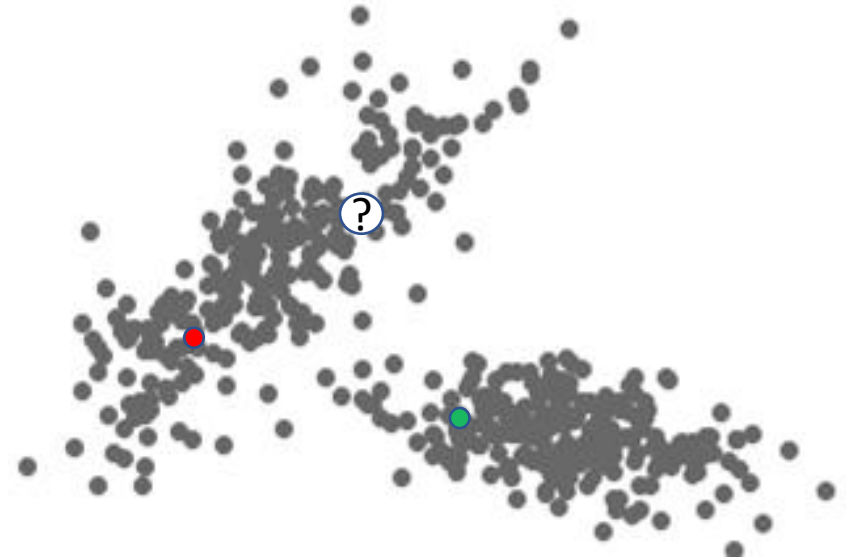
- **Probabilistic language** (conditional independence, hierarchical probabilistic models) often convenient and natural for expressing assumptions about reality
- Allows leveraging other assumptions, and integrating multiple data sources, in a principled way
  - E.g., knowledge about  $P(x|w)$  can be useful if we have lots of *unlabeled* data ( $x_i$  without  $y_i$ )  
(there are other ways to leverage unlabeled data, but prob approach directly implies recipe)

?



# Advantages of Probabilistic Approach

- **Probabilistic language** (conditional independence, hierarchical probabilistic models) often convenient and natural for expressing assumptions about reality
- Allows leveraging other assumptions, and integrating multiple data sources, in a principled way
  - E.g., knowledge about  $P(x|w)$  can be useful if we have lots of *unlabeled* data ( $x_i$  without  $y_i$ )  
(there are other ways to leverage unlabeled data, but prob approach directly implies recipe)



# Advantages of Probabilistic Approach

- **Probabilistic language** (conditional independence, hierarchical probabilistic models) often convenient and natural for expressing assumptions about reality
- Allows leveraging other assumptions, and integrating multiple data sources, in a principled way
  - E.g., knowledge about  $P(x|w)$  can be useful if we have lots of *unlabeled* data ( $x_i$  without  $y_i$ )  
(there are other ways to leverage unlabeled data, but prob approach directly implies recipe)
  - Language models:  $P(\text{text}; \theta) = \prod_i p(\text{text}[i] \mid \text{text}[:i-1]; \theta)$   
training:  $\arg \max_{\theta} P(\{\text{doc}_1, \text{doc}_2, \dots, \text{doc}_m\}; \theta) \equiv \arg \min_{\theta} \sum_t \sum_i -\log p(\text{doc}_t[i] \mid \text{doc}_t[:i-1]; \theta)$
  - Leveraging other tasks, other observations of reality, ...
- Yields meaningful “certainties”  $P(y|x)$ , not only  $h(x)$ , taking into account  $P(y|x, w)$ , but more importantly also uncertainty in estimating  $w$  (variance in the posterior  $P(w|S)$ )
- Actual right thing to do: integrate over posterior  $p(w|S)$

- MAP approach:
  - $\hat{w} = \arg \max_w p(w|\{y_i\}; \{x_i\})$  or  $\hat{w} = \arg \max_w p(w|\{y_i\}, \{x_i\})$  e.g.  $\arg \min_w L^{lgt}(r_w) + \lambda \|w\|^2$
  - Build predictor  $h(x_{query}) = h_{\hat{w}}(x_{query})$  based on  $P(y_{query}|x_{query}; \hat{w})$  e.g.  $h_{\hat{w}}(w) = \text{sign}(r_{\hat{w}})$
- This depends not only on model class, but also choice of parametrization
  - Consider parametrizing using  $u = e^w$  (ie  $w = \log u$ ), with the same prior *distribution* over models, i.e.  $P(u \in [a, b]) = P(w \in [\log a, \log b])$ 
    - ➔ The *density*:  $p_u(u) = \frac{dw}{du} p_w(w) = \frac{1}{u} p_w(\log u)$
  - The posterior *distribution* over models is the same, i.e.  $P(u \in [a, b]|S) = P(w \in [\log a, \log b]|S)$ 
    - ➔ The posterior *density*:  $p_{u|S}(u|S) = \frac{1}{u} p_{w|S}(w|S)$
  - Maximizing  $p_{u|S}$  and  $p_{w|S}$  is different!
- True Probabilistic (Bayesian) approach:
  - Use  $P(y_{query}|\mathcal{S}, x_{query}) = \int_w P(y_{query}|x_{query}, w) P(w|\mathcal{S}) = \mathbb{E}_{w \sim P(w|\mathcal{S})} [P(y_{query}|x_{query}, w)]$
  - E.g.  $h(x_{query}) = \underset{y}{\operatorname{argmin}} \mathbb{E}_{y \sim P(y_{query}|\mathcal{S})} [\text{loss}(y, y_{query})] = \mathbb{I}[P(y_{query}|\mathcal{S}) > 0.5]$
- Much of the challenge of Bayesian Inference/Bayesian ML: calculating the integral

Beyond scope of course



# Probabilistic Bayesian Approach

MAP:  $\hat{w} = \arg \max_w p(w|\{y_i\}; \{x_i\})$  or  $\hat{w} = \arg \max_w p(w|\{y_i\}, \{x_i\})$

then predict using  $r_{\hat{w}}(x)$ , pretending, e.g.,  $P(y|x) = P(y|x, \hat{w}) = \frac{1}{1+\exp(-r_{\hat{w}}(x))}$

Actual “right thing to do”™: integrate over posterior  $p(w|S)$

Use  $P(y_{tst}|x_{tst}, \{y_i\}, \{x_i\}) = \int P(y_{tst}|x_{tst}, w)p(w|\{y_i\}, \{x_i\})dw$

- Sometimes (significantly) improve prediction
- More “correct” since unlike MAP, doesn’t depend on parametrization/base measure  
→ can fix “broken” situations with MAP
- More importantly: more correct uncertainty estimate
- But: integral generally extremely intractable (though worth approximating?)
- And..... it’s the “right thing to do” only if you truly believe  $P(y, x|w), P(w)$

Beyond scope of course

Conditional likelihood  
 $-\log P(\{y_i\}|\{x_i\}; w)$

predictor class  
E.g.  $h_w = \langle w, \phi(x) \rangle$

"prior"  
 $R(w) = -\log P(w)$

$$\min_w L_S(h_w) + \lambda R(w)$$

$$L_S(h_w) = \frac{1}{m} \sum \ell(h_w(x_i); y_i)$$
$$\ell(h_w(x_i); y_i) = -\log P(x|y)$$

Captures, or is surrogate, for task-loss (e.g.  $\ell^{01}$ )

If  $R(w) \propto \text{scale}$  : Scale sensitive (e.g. Lipschitz)

For tractability: convex

$\ell(z; y)$  convex in  $z$  AND  $h_w = \langle w, \phi(x) \rangle$  linear in  $w$

$\Rightarrow \ell(h_w(x); y)$  and hence  $L_S(h_w)$  convex in  $w$

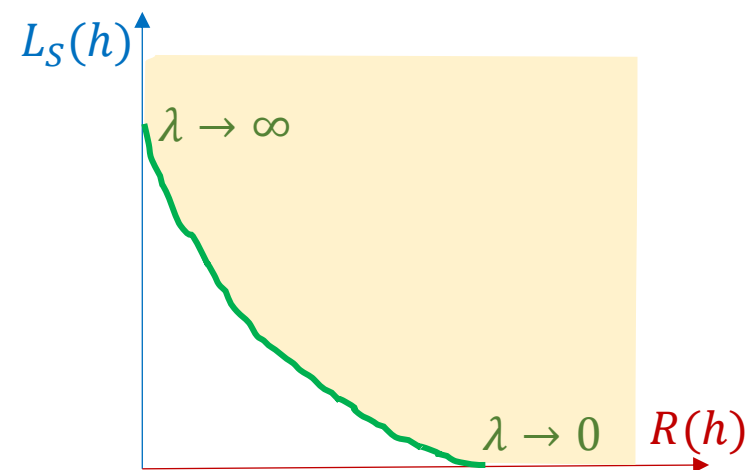
Inductive bias, simplicity among  
predictors  $h_w \in \mathcal{H}$  ( $\approx$ all functions?)

e.g.  $R(w) = \|w\|_2$

Or other norms, e.g.  $\|w\|_1$

Or non-norm, e.g.  $\|w\|_0 = |\{i | w[i] \neq 0\}|$

For tractability:  $R(w)$  convex in  $w$



$$h_w = w \cdot x + w^2$$

$$\min_w L_S(h_w) + \lambda R(w)$$

$$L_S(h_w) = \frac{1}{m} \sum \ell(h_w(x_i); y_i)$$

$$\ell(h_w(x_i); y_i) = -\log P(x|y)$$

Captures, or is surrogate, for task-loss (e.g.  $\ell^{01}$ )

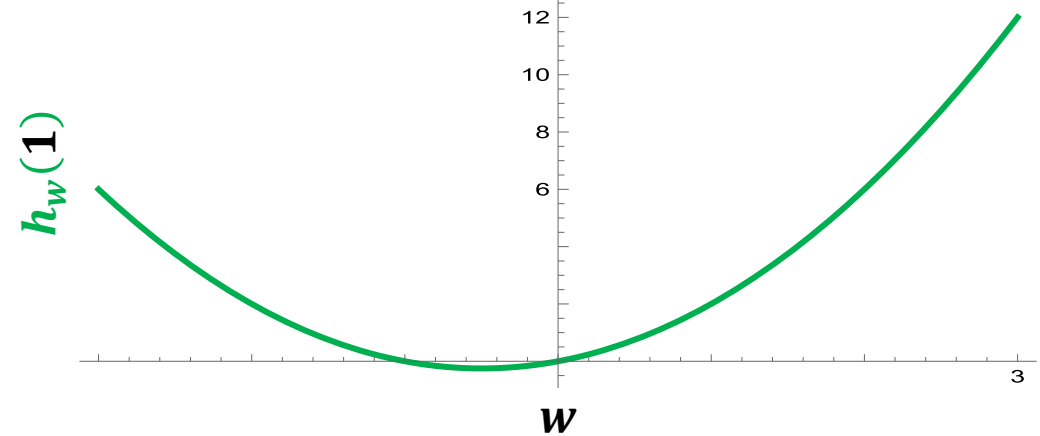
If  $R(w) \propto \text{scale}$  : Scale sensitive (e.g. Lipschitz)

For tractability: convex

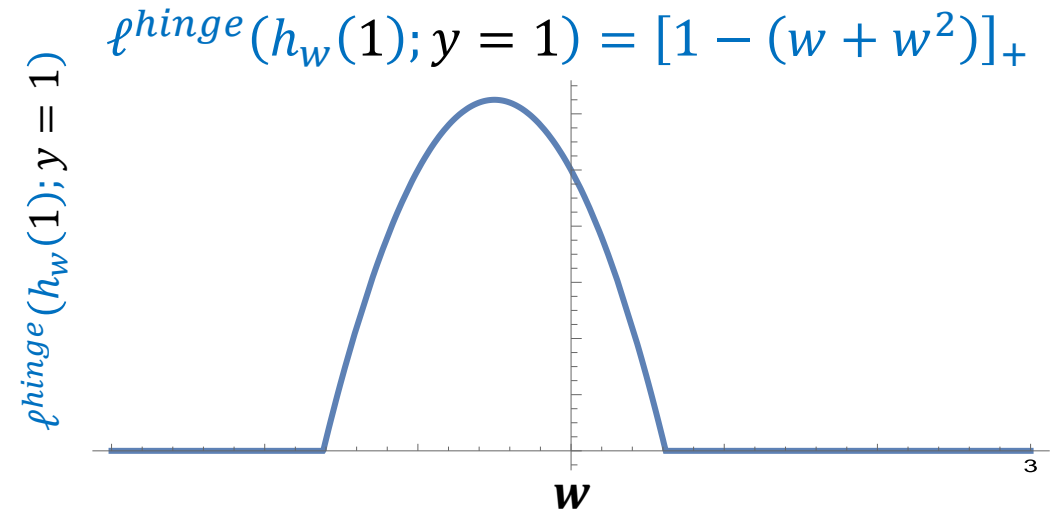
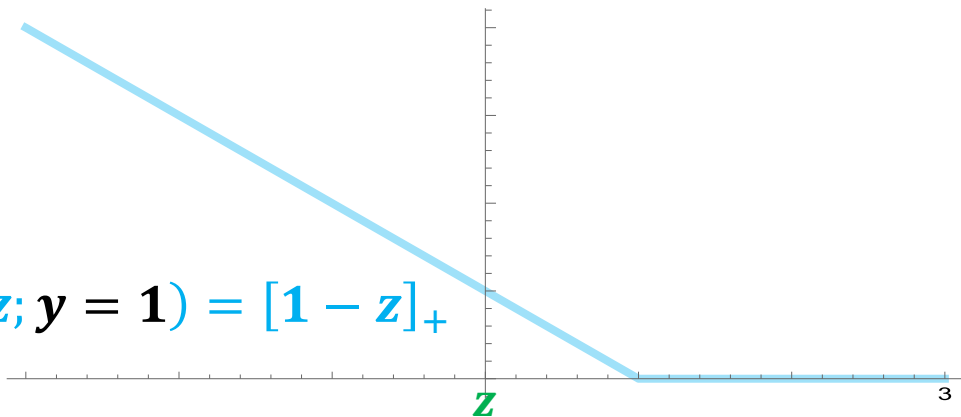
$\ell(z; y)$  convex in  $z$

→  $\ell(h_w(x); y)$  and hence  $L_S(h_w)$  convex in  $w$

$$h_w(x) = w \cdot x[1] + w^2$$



$$\ell^{hinge}(z; y = 1) = [1 - z]_+$$



$$h_w = \langle w, \phi(x) \rangle$$

convex(affine) is convex  
 $\sum$ convex is convex

$$\min_w L_S(h_w) + \lambda R(w)$$

$$L_S(h_w) = \frac{1}{m} \sum \ell(h_w(x_i); y_i)$$

$$\ell(h_w(x_i); y_i) = -\log P(x|y)$$

Captures, or is surrogate, for task-loss (e.g.  $\ell^{01}$ )

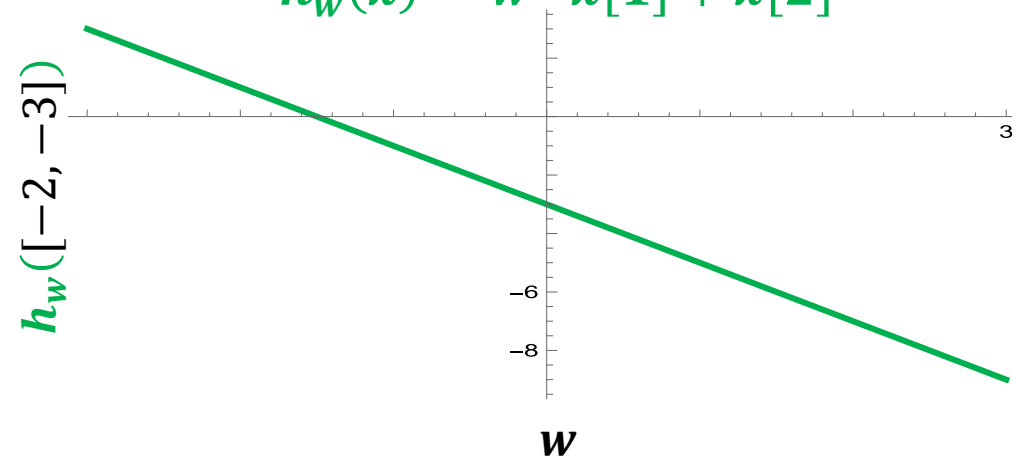
If  $R(w) \propto \text{scale}$  : Scale sensitive (e.g. Lipschitz)

For tractability: convex

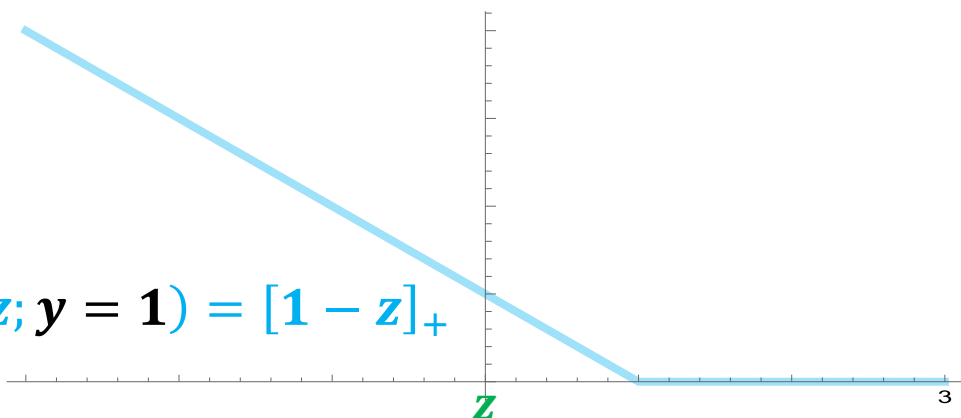
$\ell(z; y)$  convex in  $z$  AND  $h_w = \langle w, \phi(x) \rangle$  linear in  $w$

$\Rightarrow \ell(h_w(x); y)$  and hence  $L_S(h_w)$  convex in  $w$

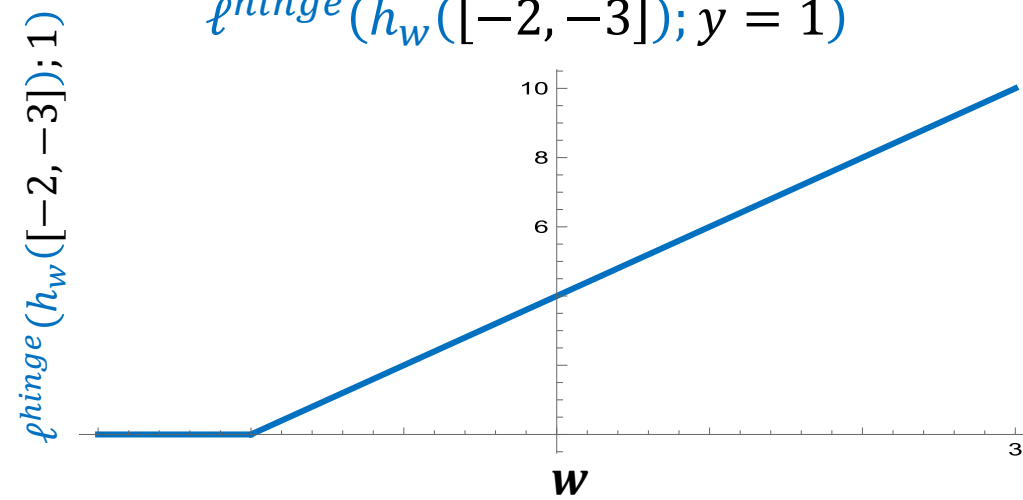
$$h_w(x) = w \cdot x[1] + x[2]$$



$$\ell^{hinge}(z; y = 1) = [1 - z]_+$$



$$\ell^{hinge}(h_w([-2, -3]); y = 1)$$



# Sub-Gradient Descent

$$\min F(w)$$

Initialize  $w^{(0)} = 0$

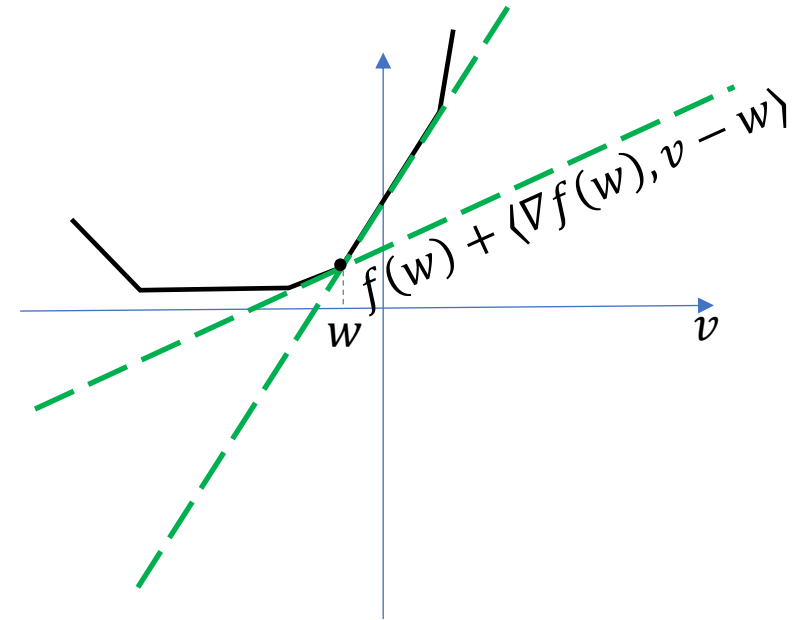
At iteration t:

- $w^{(t+1)} \leftarrow w^{(t)} - \eta_t \nabla F(w^{(t)})$

# Sub-Gradients

- Definition:  $\nabla f(w) \in \mathbb{R}^d$  is a **sub-gradient** of  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  at  $w$  if for all  $v \in \mathbb{R}^d$ :  
$$f(v) \geq f(w) + \langle \nabla f(w), v - w \rangle$$

linear lower bound passing through  $f(w)$
- Claim:  $f$  is convex if and only if  $f$  has a sub-gradient at each point.
- Claim: If  $f$  is convex and differentiable, its only subgradient at each point is the gradient.
  - At non-diff points, there may be many sub-gradients. The set of all sub-gradients is the “differential set”  $\partial f(w)$
- Sub-Gradient Descent:  $w^{(t+1)} \leftarrow w^{(t)} - \eta_t \cdot \nabla f(w^{(t)})$



# (Sub)Gradient Descent

$$\min F(w) \quad \text{s.t. } w \in \mathcal{W}$$

Initialize  $w^{(0)} = 0$

At iteration  $t$ :

- Obtain  $g^{(t)} = \nabla F(w^{(t)}) \in \partial F(w^{(t)})$
- $w^{(t+1)} \leftarrow \Pi_{\mathcal{W}}(w^{(t)} - \eta_t g^{(t)})$

Return  $w^{(T)}$

$$\Pi_{\mathcal{W}}(w) = \arg \min_{v \in \mathcal{W}} \|v - w\|_2$$

e.g.,  $\Pi_{\|w\|_2 \leq B}(w) = B \frac{w}{\|w\|_2}$  if  $\|w\|_2 \geq B$

## Gradients for (R)ERM

$$F(w) = \frac{1}{m} \sum_i \text{loss}(\langle w, \phi(x_i) \rangle, y_i) + \lambda R(w)$$

$$\nabla F(w) = \frac{1}{m} \sum_i \text{loss}'(\langle w, \phi(x_i) \rangle, y_i) \phi(x_i) + \lambda \nabla R(w)$$

# (Sub)Gradient Descent

$$\min F(w) \quad \text{s.t. } w \in \mathcal{W}$$

Initialize  $w^{(0)} = 0$

At iteration  $t$ :

- Obtain  $g^{(t)} = \nabla F(w^{(t)}) \in \partial F(w^{(t)})$
- $w^{(t+1)} \leftarrow \Pi_{\mathcal{W}}(w^{(t)} - \eta_t g^{(t)})$

Return  $w^{(T)}$

$\Pi_{\mathcal{W}}(w) = \arg \min_{v \in \mathcal{W}} \|v - w\|_2$   
e.g.,  $\Pi_{\|w\|_2 \leq B}(w) = B \frac{w}{\|w\|_2}$  if  $\|w\|_2 \geq B$

**Analysis:** If  $F$  is convex and  $\rho$ -Lipschitz (i.e.  $\|\nabla F(w)\|_2 \leq \rho$ ) then

using appropriate step-size (line search or  $\eta_t = \sqrt{\frac{B^2}{\rho^2 T}}$  or  $\eta_t = \sqrt{\frac{B^2}{\rho^2 t}}$ ):

$$F(w^{(T)}) \leq \inf_{w \in \mathcal{W}, \|w\|_2 \leq B} F(w) + \sqrt{\frac{B^2 \rho^2}{T}}$$

i.e., to ensure  $F(w) \leq F(w^*) + \epsilon$ , need:

$$T \geq \frac{B^2 \rho^2}{\epsilon^2}$$



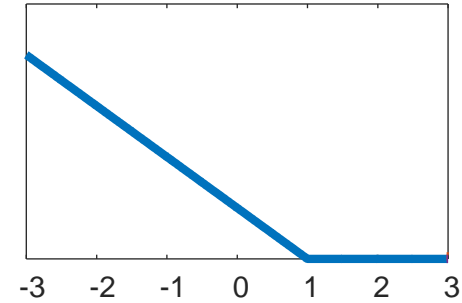
# Analysis of Projected Sub-Gradient Descent for SVM

$$\min \frac{1}{m} \sum_i \text{hinge}(y_i \langle w, \phi(x_i) \rangle) \quad \text{s.t.} \quad \|w\|_2 \leq B$$

$$\|\nabla F(w)\|_2 =$$

$$\leq \frac{1}{m} \sum_i \underbrace{|\text{hinge}'(\cdot)|}_{\leq 1} \cdot \underbrace{|y_i|}_{\leq 1} \cdot \|\phi(x_i)\|_2$$

Recall we need  
 $\|\phi(x)\| \leq R$



- Conclusion: to ensure  $L_S(w^{(T)}) - L_S(\hat{w}) \leq \epsilon_{opt}$ :

$$T = \frac{R^2 B^2}{\epsilon_{opt}^2} \text{ iterations}$$

- Runtime for each iteration:

$O(m)$  vector operations

# (Sub)Gradient Descent

$$\min F(w) \quad \text{s.t. } w \in \mathcal{W}$$

Initialize  $w^{(0)} = 0$

At iteration  $t$ :

- Obtain  $g^{(t)} = \nabla F(w^{(t)}) \in \partial F(w^{(t)})$
- $w^{(t+1)} \leftarrow \Pi_{\mathcal{W}}(w^{(t)} - \eta_t g^{(t)})$

Return  $w^{(T)}$

$$\Pi_{\mathcal{W}}(w) = \arg \min_{v \in \mathcal{W}} \|v - w\|_2$$

e.g.,  $\Pi_{\|w\|_2 \leq B}(w) = B \frac{w}{\|w\|_2}$  if  $\|w\|_2 \geq B$

**Analysis:** If  $F$  is convex and  $\rho$ -Lipschitz (i.e.  $\|\nabla F(w)\|_2 \leq \rho$ ) then

using appropriate step-size (line search or  $\eta_t = \sqrt{\frac{B^2}{\rho^2 t}}$  or  $\eta_t = \sqrt{\frac{B^2}{\rho^2 t}}$ ):

$$F(w^{(T)}) \leq \inf_{w \in \mathcal{W}, \|w\|_2 \leq B} F(w) + \sqrt{\frac{B^2 \rho^2}{T}}$$

i.e., to ensure  $F(w) \leq F(w^*) + \epsilon$ , need:

$$T = \frac{B^2 \rho^2}{\epsilon^2} = \frac{B^2 R^2}{\epsilon^2} \text{ iterations}$$

$$F(w) = \frac{1}{m} \sum_i \text{hinge}(y_i \langle w, \phi(x_i) \rangle)$$

$\|\phi(c_i)\| \leq R$

To calculate  $g^{(t)} = \nabla F(w^{(t)}) \in \partial F(w^{(t)})$ :

$O(m)$  vector operations =  $O(md)$  runtime per iteration

# Stochastic (Sub)Gradient Descent (SGD)

$$\min F(w) \quad \text{s.t. } w \in \mathcal{W}$$

Initialize  $w^{(0)} = 0$

At iteration  $t$ :

- Obtain  $g^{(t)}$  s.t.  $\mathbb{E}[g^{(t)}] \in \partial F(w^{(t)})$
- $w^{(t+1)} \leftarrow \Pi_{\mathcal{W}}(w^{(t)} - \eta_t g^{(t)})$

Return  $\bar{w}^{(T)} = \frac{1}{T} \sum_{t=1}^T w^{(t)}$

$g^{(1)}, g^{(2)}, \dots$  independent

Formally:

$$\mathbb{E}[g^{(t)} | g^{(1)} \dots g^{(t-1)}] \in \partial F(w^{(t)})$$

$$\Pi_{\mathcal{W}}(w) = \arg \min_{v \in \mathcal{W}} \|v - w\|_2$$

e.g.,  $\Pi_{\|w\|_2 \leq B}(w) = B \frac{w}{\|w\|_2}$  if  $\|w\|_2 \geq B$

**Analysis:** If  $F$  is convex and  $\rho$ -Lipschitz (i.e.  $\|\nabla F(w)\|_2 \leq \rho$ ) then

using appropriate step-size (line search or  $\eta_t = \sqrt{\frac{B^2}{\rho^2 T}}$  or  $\eta_t = \sqrt{\frac{B^2}{\rho^2 t}}$ ):

$$F(w^{(T)}) \leq \inf_{w \in \mathcal{W}, \|w\|_2 \leq B} F(w) + \sqrt{\frac{B^2 \rho^2}{T}}$$

i.e., to ensure  $F(w) \leq F(w^*) + \epsilon$ , need:

$$T = \frac{B^2 \rho^2}{\epsilon^2} = \frac{B^2 R^2}{\epsilon^2} \text{ iterations}$$

Only need independent unbiased estimates  $g^{(t)}$  of (sub)gradient, i.e. s.t.  $\mathbb{E}[g^{(t)}] = \nabla F(w^{(t)}) \in \partial F(w^{(t)})$

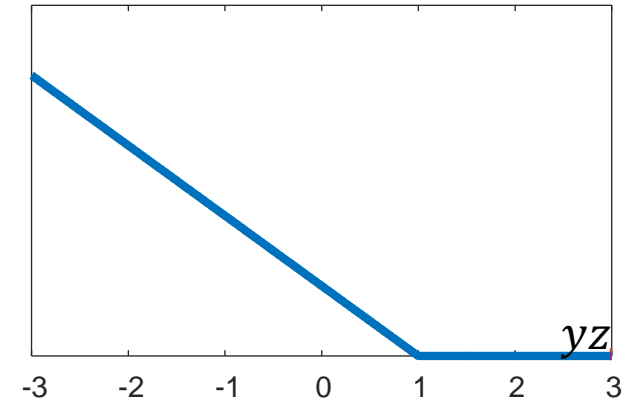
# Obtaining Independent Sub-Gradients Estimates

$$F(w) = \frac{1}{m} \sum_i \text{loss}(\langle w, \phi(x_i) \rangle, y_i)$$
$$\nabla F(w) = \frac{1}{m} \sum_i \text{loss}'(\langle w, \phi(x_i) \rangle, y_i) \phi(x_i)$$

Pick  $i \sim \text{Unif}(1..m)$  at random and use:  $g = \text{loss}'(\langle w, \phi(x_i) \rangle, y_i) \phi(x_i)$

$$\rightarrow \mathbb{E}[g] = \nabla F(w)$$

For  $\ell^{\text{hinge}}(z, y) = [1 - yz]_+$ ,  $\text{loss}'(\langle w, \phi(x_i) \rangle, y_i) = \begin{cases} -y_i, & y_i \langle w, \phi(x_i) \rangle < 1 \\ 0, & y_i \langle w, \phi(x_i) \rangle > 1 \end{cases}$



$$\rightarrow g = \text{loss}'(\langle w, \phi(x_i) \rangle, y_i) \phi(x_i) = \begin{cases} -y_i \phi(x_i) & , y_i \langle w, \phi(x_i) \rangle < 1 \\ 0 & , y_i \langle w, \phi(x_i) \rangle > 1 \end{cases}$$

# Obtaining Independent Sub-Gradients Estimates

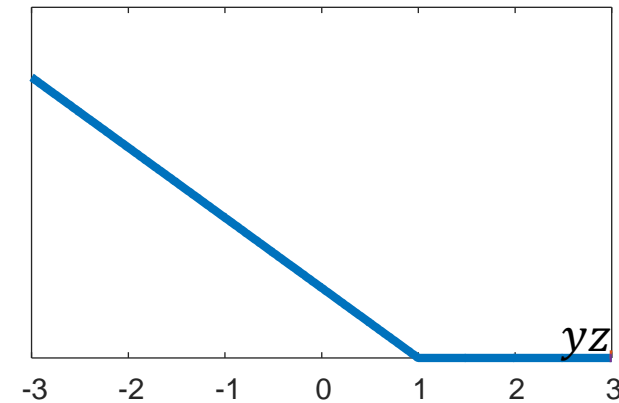
$$F(w) = \frac{1}{m} \sum_i \text{loss}(\langle w, \phi(x_i) \rangle, y_i) + \lambda R(w)$$
$$\nabla F(w) = \frac{1}{m} \sum_i \text{loss}'(\langle w, \phi(x_i) \rangle, y_i) \phi(x_i) + \lambda \nabla R(w)$$

Pick  $i \sim \text{Unif}(1..m)$  at random and use:  $g = \text{loss}'(\langle w, \phi(x_i) \rangle, y_i) \phi(x_i) + \lambda \nabla R(w)$

$$\rightarrow \mathbb{E}[g] = \nabla F(w)$$

For  $\ell^{\text{hinge}}(z, y) = [1 - yz]_+$ ,  $\text{loss}'(\langle w, \phi(x_i) \rangle, y_i) = \begin{cases} -y_i, & y_i \langle w, \phi(x_i) \rangle < 1 \\ 0, & y_i \langle w, \phi(x_i) \rangle > 1 \end{cases}$

For  $R(w) = \frac{1}{2} \|w\|_2^2$ ,  $\nabla R(w) = w$



$$\rightarrow g = \text{loss}'(\langle w, \phi(x_i) \rangle, y_i) \phi(x_i) + \lambda w = \begin{cases} -y_i \phi(x_i) + \lambda w, & y_i \langle w, \phi(x_i) \rangle < 1 \\ \lambda w, & y_i \langle w, \phi(x_i) \rangle > 1 \end{cases}$$

$\{g^{(t)}\}$  independent  
And  $\mathbb{E}[g^{(t)}] = \nabla L_S(w^{(t)})$

# SGD for SVM

$$\min L_S(w) \text{ s.t. } \|w\|_2 \leq B$$

$$\text{Use } g^{(t)} = \nabla_w \text{hinge}(y_i \langle w, \phi_i(x) \rangle) = \begin{cases} -y_i \phi(x_i), & y_i \phi(x_i) < 1 \\ 0, & y_i \phi(x_i) > 1 \end{cases} \text{ for random } i$$

Initialize  $w^{(0)} = 0$

At iteration t:

- Obtain  $g^{(t)}$  s.t.  $\mathbb{E}[g^{(t)}] \in \partial F(w^{(t)})$
- $w^{(t+1)} \leftarrow \Pi_{\mathcal{W}}(w^{(t)} - \eta_t g^{(t)})$

$$\text{Return } \bar{w}^{(T)} = \frac{1}{T} \sum_{t=1}^T w^{(t)}$$

Initialize  $w^{(0)} = 0$

At iteration t:

- Pick  $i \in 1 \dots m$  at random
- If  $y_i \langle w^{(t)}, \phi(x_i) \rangle < 1$ ,  
 $w^{(t+1)} \leftarrow w^{(t)} + \eta_t y_i \phi(x_i)$   
 else:  $w^{(t+1)} \leftarrow w^{(t)}$

- If  $\|w^{(t+1)}\|_2 > B$ , then  $w^{(t+1)} \leftarrow B \frac{w^{(t+1)}}{\|w^{(t+1)}\|_2}$

$$\text{Return } \bar{w}^{(T)} = \frac{1}{T} \sum_{t=1}^T w^{(t)}$$

$$g^{(t)} = \nabla_w [1 - y_i w^{(t)} \langle w, \phi_i(x) \rangle]_+$$

$$w^{(t+1)} \leftarrow \Pi_{\|w\| \leq B}(w^{(t)} - \eta_t g^{(t)})$$

$$\|g^{(t)}\|_2 \leq R \Rightarrow L_S(\bar{w}^{(T)}) \leq L_S(\hat{w}) + \sqrt{\frac{B^2 R^2}{T}}$$

(in expectation over randomness in algorithm)

$\{g^{(t)}\}$  independent  
 And  $\mathbb{E}[g^{(t)}] = \nabla(L_S(w^{(t)}) + \frac{\lambda}{2}\|w\|_2^2)$

# SGD for SVM

$$\min L_S(w) + \frac{\lambda}{2}\|w\|_2^2$$

$$\text{Use } g^{(t)} = \nabla_w \text{hinge}(y_i \langle w, \phi_i(x) \rangle) = \begin{cases} -y_i \phi(x_i) + \lambda w, & y_i \phi(x_i) < 1 \\ \lambda w, & y_i \phi(x_i) > 1 \end{cases} \text{ for random } i$$

Initialize  $w^{(0)} = 0$

At iteration t:

- Obtain  $g^{(t)}$  s.t.  $\mathbb{E}[g^{(t)}] \in \partial F(w^{(t)})$
- $w^{(t+1)} \leftarrow \Pi_{\mathcal{W}}(w^{(t)} - \eta_t g^{(t)})$

Return  $\bar{w}^{(T)} = \frac{1}{T} \sum_{t=1}^T w^{(t)}$

Initialize  $w^{(0)} = 0$

At iteration t:

- Pick  $i \in 1 \dots m$  at random
- If  $y_i \langle w^{(t)}, \phi(x_i) \rangle < 1$ ,  
 $w^{(t+1)} \leftarrow w^{(t)} + \eta_t y_i \phi(x_i)$   
 else:  $w^{(t+1)} \leftarrow w^{(t)}$
- $w^{(t+1)} \leftarrow w^{(t+1)} - \lambda \eta_t w^{(t)}$

Return  $\bar{w}^{(T)} = \frac{1}{T} \sum_{t=1}^T w^{(t)}$

$$g^{(t)} = \nabla_w [1 - y_i w^{(t)} \phi(x_i)]_+ + \lambda w$$

$$w^{(t+1)} \leftarrow w^{(t)} - \eta_t g^{(t)}$$

# Stochastic vs Batch

- Batch Gradient descent:

$O\left(\frac{B^2 R^2}{\epsilon^2}\right)$  iterations,  **$O(m)$  vector operations per iteration**,  $O(d)$  per vector operation

➔ Runtime:  $O\left(\frac{B^2 R^2}{\epsilon^2} d \cdot m\right)$

- Stochastic Gradient Descent:

$O\left(\frac{B^2 R^2}{\epsilon^2}\right)$  iterations,  **$O(1)$  vector operations per iteration**,  $O(d)$  per vector operation

➔ Runtime:  $O\left(\frac{B^2 R^2}{\epsilon^2} d\right)$

Intuitive argument: if only taking simple gradient steps, might as well be stochastic



# Stochastic vs Batch

$$\min L_S(w) = \frac{1}{m} \sum_i \text{loss}(w \text{ on } (x_i, y_i))$$

$w \leftarrow w - \eta \sum g_i$	{	$w \leftarrow w - \eta g_1$	$g_1 = \nabla \text{loss}(w \text{ on } (x_1, y_1))$	$\mathbf{x}_1, y_1$	$g_1 = \nabla \text{loss}(w \text{ on } (x_1, y_1))$	}
		$w \leftarrow w - \eta g_2$	$g_2 = \nabla \text{loss}(w \text{ on } (x_2, y_2))$	$\mathbf{x}_2, y_2$	$g_2 = \nabla \text{loss}(w \text{ on } (x_2, y_2))$	
		$w \leftarrow w - \eta g_3$	$g_3 = \nabla \text{loss}(w \text{ on } (x_3, y_3))$	$\mathbf{x}_3, y_3$	$g_3 = \nabla \text{loss}(w \text{ on } (x_3, y_3))$	
		$w \leftarrow w - \eta g_4$	$g_4 = \nabla \text{loss}(w \text{ on } (x_4, y_4))$	$\mathbf{x}_4, y_4$	$g_4 = \nabla \text{loss}(w \text{ on } (x_4, y_4))$	
		$w \leftarrow w - \eta g_5$	$g_5 = \nabla \text{loss}(w \text{ on } (x_5, y_5))$	$\mathbf{x}_5, y_5$	$g_5 = \nabla \text{loss}(w \text{ on } (x_5, y_5))$	
		$\vdots$		$\vdots$		
		$w \leftarrow w - \eta g_{m-1}$	$g_m = \nabla \text{loss}(w \text{ on } (x_m, y_m))$	$\mathbf{x}_m, y_m$	$g_m = \nabla \text{loss}(w \text{ on } (x_m, y_m))$	
		$w \leftarrow w - \eta g_m$				

$\nabla L_S(w) = \frac{1}{m} \sum g_i$   
 $w \leftarrow w - \eta \frac{1}{m} \sum g_i$

# Stochastic vs Batch

$$\min L_S(w) = \frac{1}{m} \sum_i \text{loss}(w \text{ on } (x_1, y_1))$$

$w \leftarrow w - \eta \sum g_i$	$\left\{ \begin{array}{l} w \leftarrow w - \eta g_1 \\ w \leftarrow w - \eta g_2 \\ w \leftarrow w - \eta g_3 \end{array} \right.$	$g_1 = \nabla \text{loss}(w \text{ on } (x_1, y_1))$	$\mathbf{x}_1, y_1$	$g_1 = \nabla \text{loss}(w \text{ on } (x_1, y_1))$	$\left. \begin{array}{l} g_2 = \nabla \text{loss}(w \text{ on } (x_2, y_2)) \\ g_3 = \nabla \text{loss}(w \text{ on } (x_3, y_3)) \\ g_4 = \nabla \text{loss}(w \text{ on } (x_4, y_4)) \\ g_5 = \nabla \text{loss}(w \text{ on } (x_5, y_5)) \end{array} \right\} \nabla L_S(w) = \frac{1}{m} \sum g_i$
		$g_2 = \nabla \text{loss}(w \text{ on } (x_2, y_2))$	$\mathbf{x}_2, y_2$	$g_2 = \nabla \text{loss}(w \text{ on } (x_2, y_2))$	
		$g_3 = \nabla \text{loss}(w \text{ on } (x_3, y_3))$	$\mathbf{x}_3, y_3$	$g_3 = \nabla \text{loss}(w \text{ on } (x_3, y_3))$	
		$g_4 = \nabla \text{loss}(w \text{ on } (x_4, y_4))$	$\mathbf{x}_4, y_4$	$g_4 = \nabla \text{loss}(w \text{ on } (x_4, y_4))$	
		$g_5 = \nabla \text{loss}(w \text{ on } (x_5, y_5))$	$\mathbf{x}_5, y_5$	$g_5 = \nabla \text{loss}(w \text{ on } (x_5, y_5))$	
		$\vdots$	$\vdots$		
$w \leftarrow w - \eta \sum g_i$	$\left\{ \begin{array}{l} w \leftarrow w - \eta g_4 \\ w \leftarrow w - \eta g_5 \end{array} \right.$				$\left. \begin{array}{l} g_m = \nabla \text{loss}(w \text{ on } (x_m, y_m)) \end{array} \right\}$
		$\vdots$			
$w \leftarrow w - \eta \sum g_i$	$\left\{ \begin{array}{l} w \leftarrow w - \eta g_{m-1} \\ w \leftarrow w - \eta g_m \end{array} \right.$	$g_m = \nabla \text{loss}(w \text{ on } (x_m, y_m))$	$\mathbf{x}_m, y_m$	$g_m = \nabla \text{loss}(w \text{ on } (x_m, y_m))$	$w \leftarrow w - \eta \frac{1}{m} \sum g_i$

Mini-Batch SGD: use  $g_t = \frac{1}{b} \sum_{j \in \mathcal{B}_t} \nabla \ell(h_{w^{(t)}}(x_j), y_j)$

$\mathcal{B}_t$  = random subset of  $b$  training indices

$$F(w) = \frac{1}{m} \sum_i \text{loss}(\langle w, \phi(x_i) \rangle, y_i) \quad \nabla F(w) = \frac{1}{m} \sum_i \text{loss}'(\langle w, \phi(x_i) \rangle, y_i) \phi(x_i)$$

At each iteration, independent (sub)gradient estimate  $g_t$  s.t.  $\mathbb{E}[g_t] = \nabla F(w^{(t)})$

$$g_t = \frac{1}{b} \sum_{i \in \mathcal{B}_t} \text{loss}'(\langle w^{(t)}, \phi(x_i) \rangle, y_i) \phi(x_i) \quad \mathcal{B}_t = \text{random subset of } \mathbf{b} \text{ training indices}$$

Initialize  $w^{(0)} = 0$

At iteration  $t$ :

- Obtain  $g^{(t)}$  s.t.  $\mathbb{E}[g^{(t)}] \in \partial F(w^{(t)})$
- $w^{(t+1)} \leftarrow \Pi_{\mathcal{W}}(w^{(t)} - \eta_t g^{(t)})$

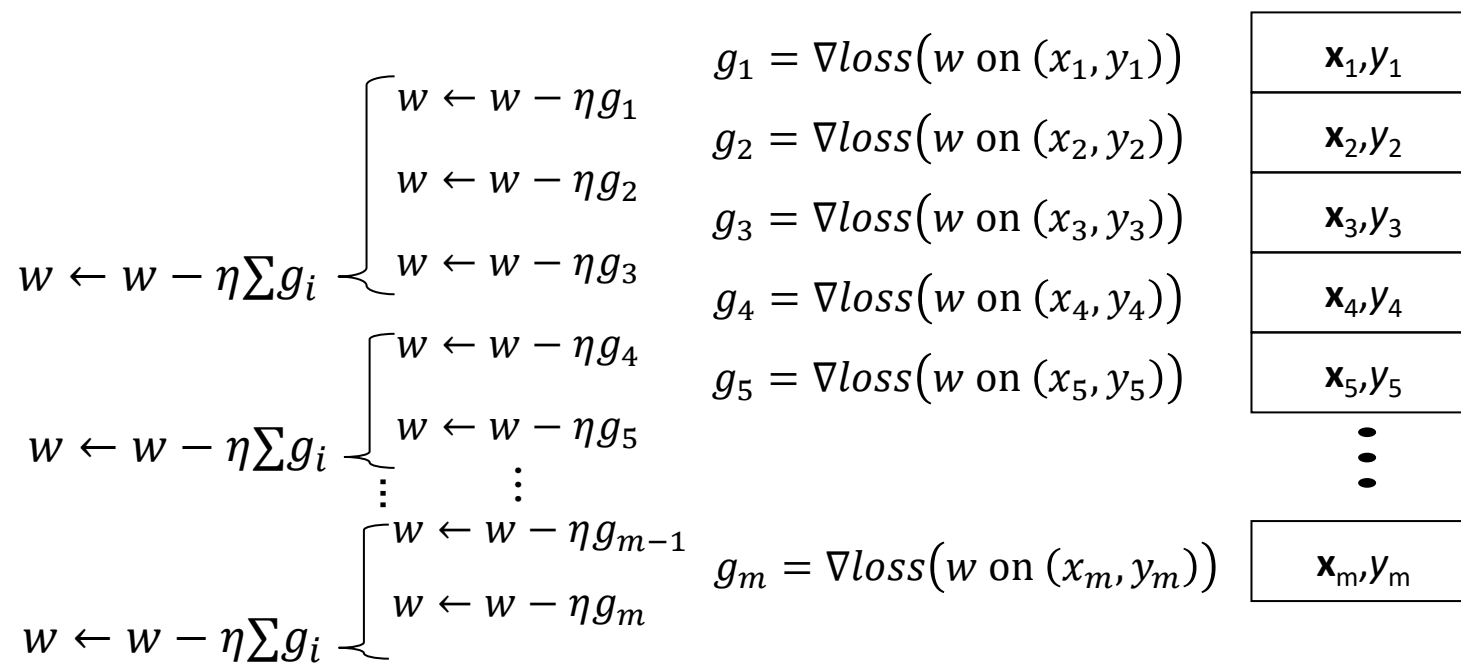
Initialize  $w^{(0)} = 0$

At iteration  $t$ :

- Pick  $\mathcal{B}_t \subseteq [1..m]$ ,  $|\mathcal{B}_t| = \mathbf{b}$  at random
- $w^{(t+1)} = w^{(t)}$
- For  $i \in \mathcal{B}_t$ 

$$w^{(t+1)} \leftarrow w^{(t)} + \frac{\eta_t}{b} \text{loss}'(\langle w^{(t)}, \phi(x_i) \rangle, y_i) \phi(x_i)$$
- If  $\|w^{(t+1)}\|_2 > B$ , then  $w^{(t+1)} \leftarrow B \frac{w^{(t+1)}}{\|w^{(t+1)}\|_2}$

$O\left(\frac{B^2 R^2}{\epsilon^2}\right)$  iterations,  $O(\mathbf{b})$  vector operations per iteration  $\Rightarrow O\left(\frac{B^2 R^2}{\epsilon^2} \mathbf{b} d\right)$



Initialize  $w^{(0)} = 0$

At iteration  $t$ :

- Obtain  $g^{(t)}$  s.t.  $\mathbb{E}[g^{(t)}] \in \partial F(w^{(t)})$
- $w^{(t+1)} \leftarrow \Pi_{\mathcal{W}}(w^{(t)} - \eta_t g^{(t)})$

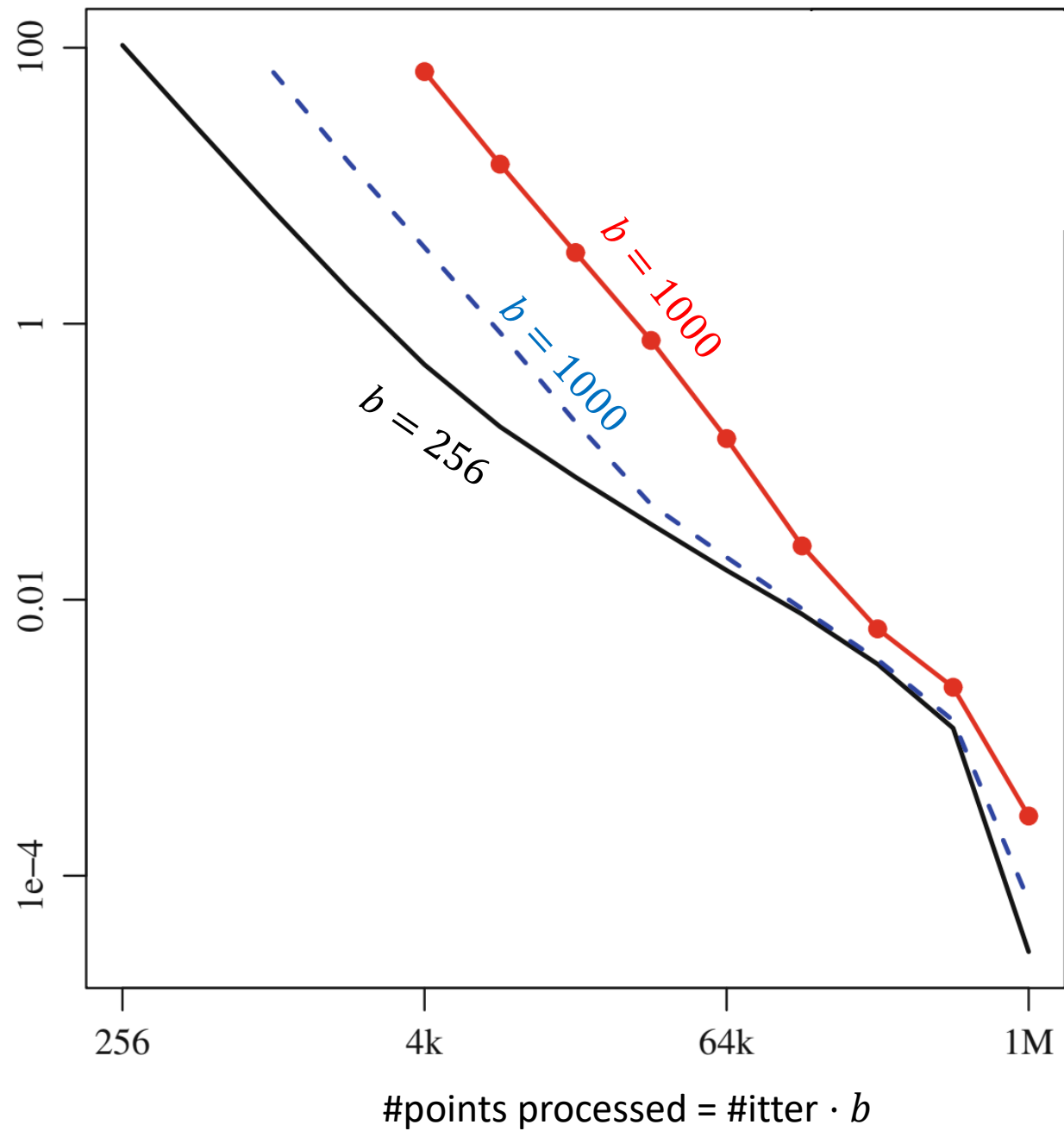
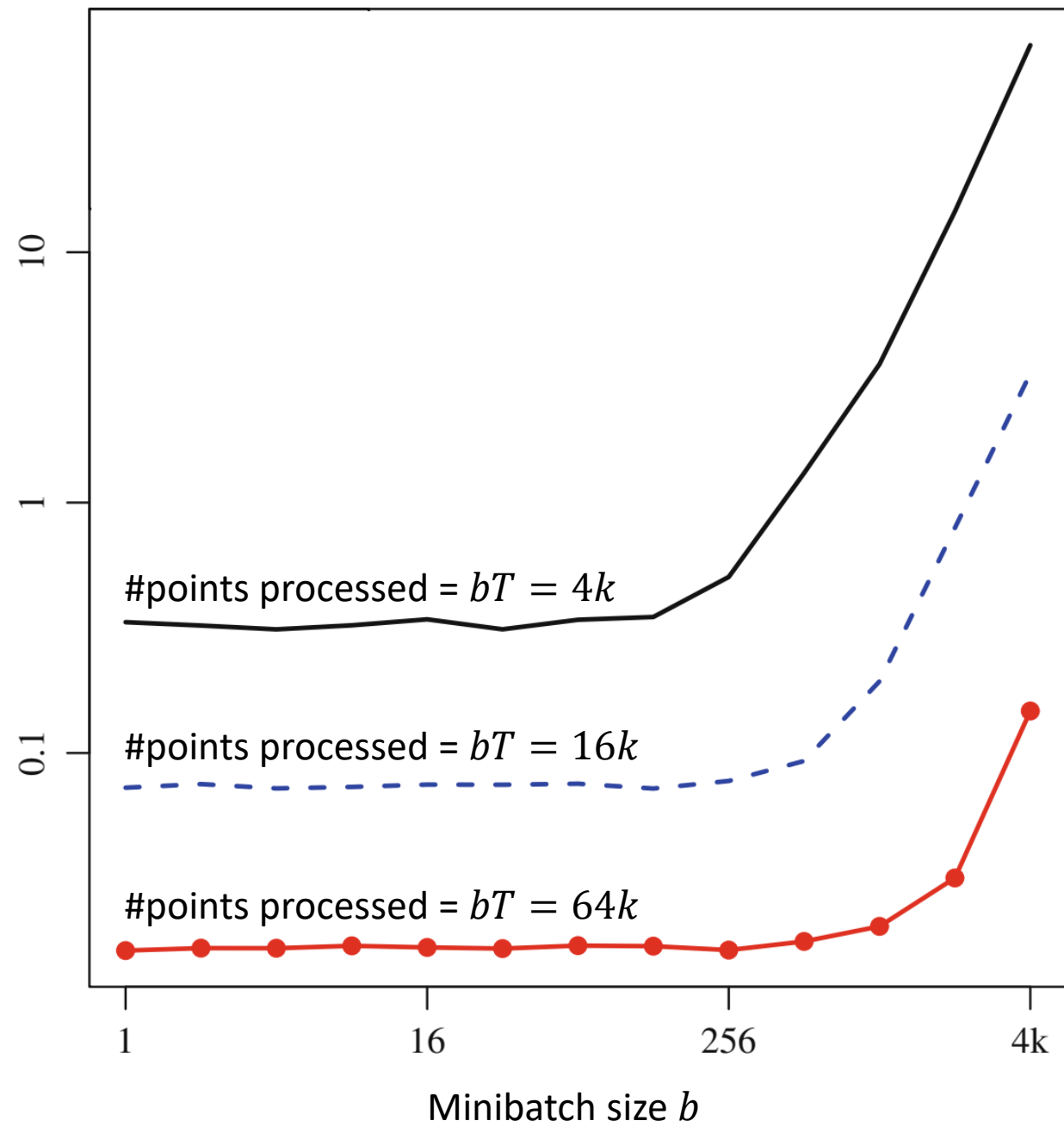
Initialize  $w^{(0)} = 0$

At iteration  $t$ :

- Pick  $\mathcal{B}_t \subseteq [1..m]$ ,  $|\mathcal{B}_t| = \mathbf{b}$  at random
- $w^{(t+1)} = w^{(t)}$
- For  $i \in \mathcal{B}_t$ 

$$w^{(t+1)} \leftarrow w^{(t)} + \frac{\eta_t}{b} \text{loss}'(\langle w^{(t)}, \phi(x_i) \rangle, y_i) \phi(x_i)$$
- If  $\|w^{(t+1)}\|_2 > B$ , then  $w^{(t+1)} \leftarrow B \frac{w^{(t+1)}}{\|w^{(t+1)}\|_2}$

$O\left(\frac{B^2 R^2}{\epsilon^2}\right)$  iterations,  $O(\mathbf{b})$  vector operations per iteration  $\Rightarrow O\left(\frac{B^2 R^2}{\epsilon^2} \mathbf{b} d\right)$



- Stochastic Gradient Descent Runtime:  $O\left(\frac{B^2 R^2}{\epsilon_{opt}^2} d\right)$
- Batch Gradient Descent Runtime:  $O\left(\frac{B^2 R^2}{\epsilon_{opt}^2} d \cdot m\right)$
- Mini-Batch Gradient Descent Runtime:  $O\left(\frac{B^2 R^2}{\epsilon_{opt}^2} b \cdot m\right)$

Optimal  $b$  (in theory and in practice if counting #vector ops):  $b = 1$

In practice, moderate  $b$  just as good as  $b = 1$ , reduces overhead, allows parallelization

- GD actually better for smooth objectives, and under other assumptions
- Also alternate methods, eg Newton, approx. Newton (including BFGS) for smooth objective, Interior Point methods for handling non-smoothness/constraints
- Goal of much of optimization:  $\log 1/\epsilon$  dependence (and many of the above achieve this)
- How small should  $\epsilon_{opt}$  be?
- What about  $L_{\mathcal{D}}(w)$ , which is what we really care about?

# Overall Analysis of $L_{\mathcal{D}}(w)$

- Recall for ERM:  $L_{\mathcal{D}}(\hat{w}) \leq L_{\mathcal{D}}(w^*) + 2 \sup_w |L_{\mathcal{D}}(w) - L_S(w)|$

$$\hat{w} = \arg \min_{\|w\| \leq B} L_S(w)$$

$$w^* = \arg \min_{\|w\| \leq B} L_{\mathcal{D}}(w)$$

- For  $\epsilon_{opt}$  suboptimal ERM  $\bar{w}$ :

$$L_{\mathcal{D}}(\bar{w}) \leq \underbrace{L_{\mathcal{D}}(w^*)}_{\epsilon_{approx}} + \underbrace{2 \sup_w |L_{\mathcal{D}}(w) - L_S(w)|}_{\epsilon_{est}} + \underbrace{(L_S(\bar{w}) - L_S(\hat{w}))}_{\epsilon_{opt}}$$
$$\epsilon_{est} \leq 2 \sqrt{\frac{B^2 R^2}{m}} \quad \epsilon_{opt} \leq \sqrt{\frac{B^2 R^2}{T}}$$

- Take  $\epsilon_{opt} \approx \epsilon_{est}$ , i.e.  $\#iter\ T \approx sample\ size\ m$
- To ensure  $L_{\mathcal{D}}(w) \leq L_{\mathcal{D}}(w^*) + \epsilon$ :

$$T, m = O\left(\frac{B^2 R^2}{\epsilon^2}\right)$$

# SGD as an Ultimate Optimization Algorithm

- Runtime of SGD is  $O(T \cdot d) = O\left(\frac{B^2 R^2}{\epsilon^2} d\right) = O(m \cdot d)$
- Linear in the number of required examples  $m$ 
  - vs full GD which requires quadratic time  $O(Tmd) = O(m^2 d)$
  - or interior point or matrix inversion methods, which require cubic or worse runtime
- SGD runtime is linear in time it takes to read data set
  - ➔ Can't be improved beyond small constant factor **without additional assumptions**



# SGD as a Learning Algorithm: SGD on $L_{\mathcal{D}}(w)$

$$\min_w L_{\mathcal{D}}(w)$$

use  $g^{(t)} = \nabla_w \text{loss}(h_{w^{(t)}}(x); y)$  for random  $y, x \sim \mathcal{D} \rightarrow \mathbb{E}[g^{(t)}] = \nabla L_{\mathcal{D}}(w)$

Initialize  $w^{(0)} = 0$

At iteration  $t$ :

- Draw  $x_t, y_t \sim \mathcal{D}$

- $w^{(t+1)} \leftarrow w^{(t)} - \eta_t \nabla_w \text{loss}(h_{w^{(t)}}(x); y) = w^{(t)} - \eta_t \text{loss}'(\langle w^{(t)}, \phi(x) \rangle; y) \phi(x)$

Return  $\bar{w}^{(T)} = \frac{1}{T} \sum_{t=1}^T w^{(t)}$

$$h_w(x) = \langle w, \phi(x) \rangle$$

SGD suboptimality guarantee:  $L_{\mathcal{D}}(\bar{w}^{(T)}) \leq \inf_{\|w\|_2 \leq B} L_{\mathcal{D}}(w) + \sqrt{\frac{B^2 R^2}{T}}$

$$\rightarrow m = T = O\left(\frac{B^2 R^2}{\epsilon^2}\right)$$