

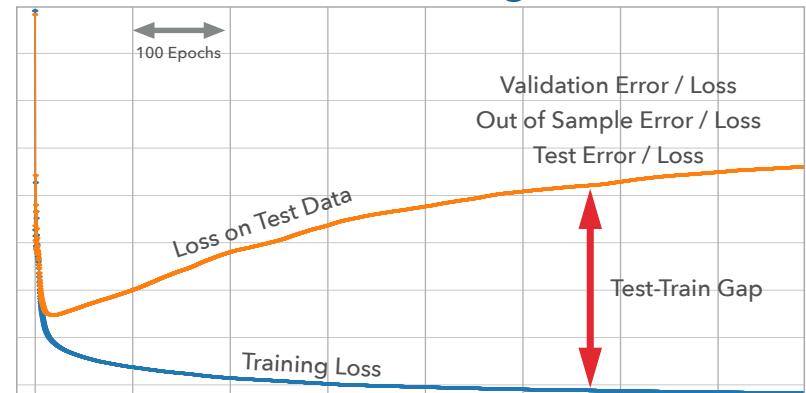
# COS324: INTRODUCTION TO MACHINE LEARNING

Prof. Yoram Singer



Topic: Learning with Regularization

## Overfitting



2

## Generalization Analysis: Ingredients

Examples are **I.I.D** samples from an **unknown** distribution  $D(\mathbf{x}, \mathbf{y})$

Restrict to  $\mathbf{x} \in \mathcal{X} = \{0, 1\}^d$  and  $\mathbf{y} \in \{-1, 1\}$

Marginal distribution:  $D(\mathbf{x}) = D(\mathbf{x}, -1) + D(\mathbf{x}, +1)$

Conditional distribution:  $D(\mathbf{y} | \mathbf{x}) = \frac{D(\mathbf{x}, \mathbf{y})}{D(\mathbf{x})} \in [0, 1]$  also denoted  $\mathbf{P}(\mathbf{y} | \mathbf{x})$

Realizable case:  $\exists h^* : \mathcal{X} \rightarrow \{-1, 1\}$  s.t.  $\mathbf{P}(y = h^*(\mathbf{x}) | \mathbf{x}) = 1$

Unrealizable case:  $0 < \mathbf{P}(\mathbf{y} | \mathbf{x}) < 1$

YORAM SINGER © 2020

3

## I.I.D Samples

- I.I.D: Identically Independently Distributed
- Generalization analysis typically assumes  $\exists D$  :  
unknown distribution  $D(\mathbf{x}, \mathbf{y})$
- Assume again that  $\mathbf{x} \in \{0, 1\}^d$   $\mathbf{y} \in \{-1, 1\}$
- Identically [no dependence on  $i$ ]:  
 $\forall i \in S : D((\mathbf{x}_i, \mathbf{y}_i) = (\mathbf{a}, \mathbf{b}))$  is  $D(\mathbf{a}, \mathbf{b})$
- Independence:  
 $D((\mathbf{x}_i, \mathbf{y}_i) = (\mathbf{a}, \mathbf{b}) \wedge (\mathbf{x}_j, \mathbf{y}_j) = (\mathbf{a}', \mathbf{b}')) = D(\mathbf{a}, \mathbf{b}) D(\mathbf{a}', \mathbf{b}')$

$x_0$	$x_1$	$y$	$D$
0	0	-1	0.1
0	0	1	0.2
0	1	-1	0
0	1	1	0.3
1	0	-1	0.1
1	0	1	0.2
1	1	-1	0.07
1	1	1	0.03

YORAM SINGER © 2020

4

Unrealizable

$x_0$	$x_1$	$y$	$D$
0	0	-1	0.1
0	0	1	0.2
0	1	-1	0
0	1	1	0.15
1	0	-1	0.05
1	0	1	0.2
1	1	-1	0.17
1	1	1	0.13

Realizable

$x_0$	$x_1$	$y$	$D$
0	0	-1	0
0	0	1	0.25
0	1	-1	0
0	1	1	0.2
1	0	-1	0
1	0	1	0.25
1	1	-1	0.3
1	1	1	0

Unrealizable

$x_0$	$x_1$	$y$	$D$
0	0	-1	0.1
0	0	1	0.2
0	1	-1	0
0	1	1	0.15
1	0	-1	0.05
1	0	1	0.2
1	1	-1	0.17
1	1	1	0.13

Realizable

$x_0$	$x_1$	$y$	$D$
0	0	-1	0
0	0	1	0.25
0	1	-1	0
0	1	1	0.2
1	0	-1	0.05
1	0	1	0.2
1	1	-1	0.25
1	1	1	0.3
1	1	1	0

Unrealizable

$x_0$	$x_1$	$y$	$D$
0	0	-1	0.1
0	0	1	0.2
0	1	-1	0
0	1	1	0.15
1	0	-1	0.05
1	0	1	0.2
1	1	-1	0.17
1	1	1	0.13

Realizable

$x_0$	$x_1$	$y$	$D$
0	0	-1	0
0	0	1	0.25
0	1	-1	0
0	1	1	0.2
1	0	-1	0
1	0	1	0.25
1	1	-1	0.3
1	1	1	0

Unrealizable

$x_0$	$x_1$	$y$	$D$
0	0	-1	0.1
0	0	1	0.2
0	1	-1	0
0	1	1	0.15
1	0	-1	0.05
1	0	1	0.2
1	1	-1	0.17
1	1	1	0.13

Realizable

$x_0$	$x_1$	$h^*$	$D$
0	0	-1	0
0	0	1	0.25
0	1	-1	0
0	1	1	0.2
1	0	-1	0
1	0	1	0.25
1	1	-1	0.3
1	1	1	0

$x_0$	$x_1$	$h^*$	$D$
0	0	1	0.25
0	1	1	0.2
1	0	1	0.25
1	1	-1	0.3

$x_0$	$x_1$	$h^*$	D
0	0	1	0.25
0	1	1	0.2
1	0	1	0.25
1	1	-1	0.3



YORAM SINGER © 2020

6

$x_0$	$x_1$	$h^*$	D
0	0	1	0.25
0	1	1	0.2
1	0	1	0.25
1	1	-1	0.3



YORAM SINGER © 2020

$x_0$	$x_1$	y	S
0	0	-1	0
0	0	1	1
0	1	-1	0
0	1	1	1
1	0	-1	0
1	0	1	2
1	1	-1	2
1	1	1	0

6

## Empirical & Population Risk

Empirical risk (or error or loss) of  $f$  w.r.t. dataset of examples  $S$ :

$$\text{err}_S(f) = \frac{1}{|S|} \sum_{(x,y) \in S} \mathbf{1}[f(x) \neq y] \quad // \text{ #errors divided by sample size}$$

Population risk (or error or loss) of  $f$  w.r.t unknown distribution D:

$$\text{err}_D(f) = \sum_x \sum_y D(x, y) \mathbf{1}[f(x) \neq y]$$

❖ Definitions the same for both realizable case & unrealizable case

❖ Realizable case also amounts to:  $\text{err}_D(f) = \sum_x D(x, y) \mathbf{1}[f(x) \neq h^*(x)]$

YORAM SINGER © 2020

YORAM SINGER © 2020

7

$x_0$	$x_1$	y	D
0	0	-1	0.1
0	0	1	0.2
0	1	-1	0
0	1	1	0.15
1	0	-1	0.05
1	0	1	0.1
1	1	-1	0.27
1	1	1	0.13

$f(x)$
1
1
1
1
-1
-1
-1
-1

$x_0$	$x_1$	y	S
0	0	-1	0
0	0	1	1
0	1	-1	0
0	1	1	1
1	0	-1	0
1	0	1	1
1	1	-1	1
1	1	1	0

8

$$\text{err}_D(f) = \sum_{x,y} D(x, y) \mathbf{1}[f(x) \neq y] = 0.33$$

$x_0$	$x_1$	$y$	$D$	$f(x)$
0	0	-1	-0.1	1
0	0	1	0.2	1
0	1	-1	0	1
0	1	1	0.15	1
1	0	-1	0.05	-1
1	0	1	-0.1	-1
1	1	-1	0.27	-1
1	1	1	-0.13	-1

$x_0$	$x_1$	$y$	$S$
0	0	-1	0
0	0	1	1
0	1	-1	0
0	1	1	1
1	0	-1	0
1	0	1	1
1	1	-1	1
1	1	1	0

YORAM SINGER © 2020

8

$$\text{err}_D(f) = \sum_{x,y} D(x, y) \mathbf{1}[f(x) \neq y] = 0.33 \quad \text{err}_S(f) = \frac{1}{|S|} \sum_{(x,y) \in S} \mathbf{1}[f(x) \neq y] = \frac{1}{4}$$

$x_0$	$x_1$	$y$	$D$	$f(x)$
0	0	-1	-0.1	1
0	0	1	0.2	1
0	1	-1	0	1
0	1	1	0.15	1
1	0	-1	0.05	-1
1	0	1	-0.1	-1
1	1	-1	0.27	-1
1	1	1	-0.13	-1

$x_0$	$x_1$	$y$	$S$
0	0	-1	0
0	0	1	1
0	1	-1	0
0	1	1	1
1	0	-1	0
1	0	1	1
1	1	-1	1
1	1	1	0

YORAM SINGER © 2020

8

## Finite Set of Predictors

Only  $k$  predictors — no parameter learning:  $f_1, \dots, f_k$

YORAM SINGER © 2020

## Finite Set of Predictors

Only  $k$  predictors — no parameter learning:  $f_1, \dots, f_k$

Best predictor is  $f_j$  such that,

$$\text{err}_D[f_j(x)] = \epsilon_j \quad \forall i \neq j : \text{err}_D[f_i(x) \neq y] \geq \epsilon_j + \epsilon$$

YORAM SINGER © 2020

9

## Finite Set of Predictors

Only  $k$  predictors — no parameter learning:  $f_1, \dots, f_k$

Best predictor is  $f_j$  such that,

$$\text{err}_D[f_j(x)] = \epsilon_j \quad \forall i \neq j : \text{err}_D[f_i(x) \neq y] \geq \epsilon_j + \epsilon$$

Evaluate errors on  $S$ :  $\text{err}_S(f_r) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}[f_r(x_i) \neq y_i]$

9

## Finite Set of Predictors

Only  $k$  predictors — no parameter learning:  $f_1, \dots, f_k$

Best predictor is  $f_j$  such that,

$$\text{err}_D[f_j(x)] = \epsilon_j \quad \forall i \neq j : \text{err}_D[f_i(x) \neq y] \geq \epsilon_j + \epsilon$$

Evaluate errors on  $S$ :  $\text{err}_S(f_r) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}[f_r(x_i) \neq y_i]$

Choose  $\hat{f}$  with smallest empirical risk:  $\hat{f} = \arg \min_i \text{err}_S(f_i)$

YORAM SINGER © 2020

9

## Generalization

✓ When  $\epsilon_j = 0$ , if  $|S|$  is greater than  $\frac{2 \log(k)}{\epsilon}$ ,

then with very high probability  $\hat{f} = f_j$

✓ When  $\epsilon_j > 0$ , if  $|S|$  is greater than  $\frac{4 \log(k)}{\epsilon^2}$

then with very high probability  $\hat{f} = f_j$

YORAM SINGER © 2020

10

## Overview of Analysis for $\epsilon_j = 0$

Probability that  $\epsilon_i = 0$  is at most  $(1 - \epsilon)^n \leq e^{-\epsilon n}$  [independence of sample]

© 2020 YORAM SINGER 11

## Overview of Analysis for $\epsilon_j = 0$

Probability that  $\epsilon_i = 0$  is at most  $(1 - \epsilon)^n \leq e^{-\epsilon n}$  [independence of sample]

Probability  $\alpha$  that  $\exists i \neq j$  s.t.  $\epsilon_i = 0$  is at most  $\alpha = (k - 1) e^{-\epsilon n}$

© 2020 YORAM SINGER 11

## Overview of Analysis for $\epsilon_j = 0$

Probability that  $\epsilon_i = 0$  is at most  $(1 - \epsilon)^n \leq e^{-\epsilon n}$  [independence of sample]

Probability  $\alpha$  that  $\exists i \neq j$  s.t.  $\epsilon_i = 0$  is at most  $\alpha = (k - 1) e^{-\epsilon n}$

If  $\alpha \leq \frac{1}{k}$  it is unlikely we do not find correct predictor:  $(k - 1) e^{-\epsilon n} \geq \frac{1}{k}$

© 2020 YORAM SINGER 11

## Overview of Analysis for $\epsilon_j = 0$

Probability that  $\epsilon_i = 0$  is at most  $(1 - \epsilon)^n \leq e^{-\epsilon n}$  [independence of sample]

Probability  $\alpha$  that  $\exists i \neq j$  s.t.  $\epsilon_i = 0$  is at most  $\alpha = (k - 1) e^{-\epsilon n}$

If  $\alpha \leq \frac{1}{k}$  it is unlikely we do not find correct predictor:  $(k - 1) e^{-\epsilon n} \geq \frac{1}{k}$

This means that we need  $O\left(\frac{\log(k)}{\epsilon}\right)$  samples

© 2020 YORAM SINGER 11

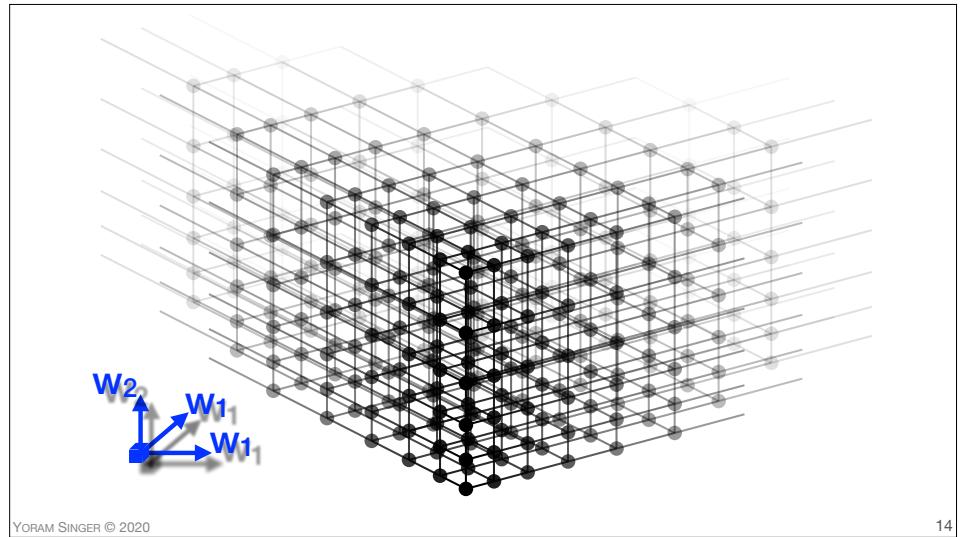
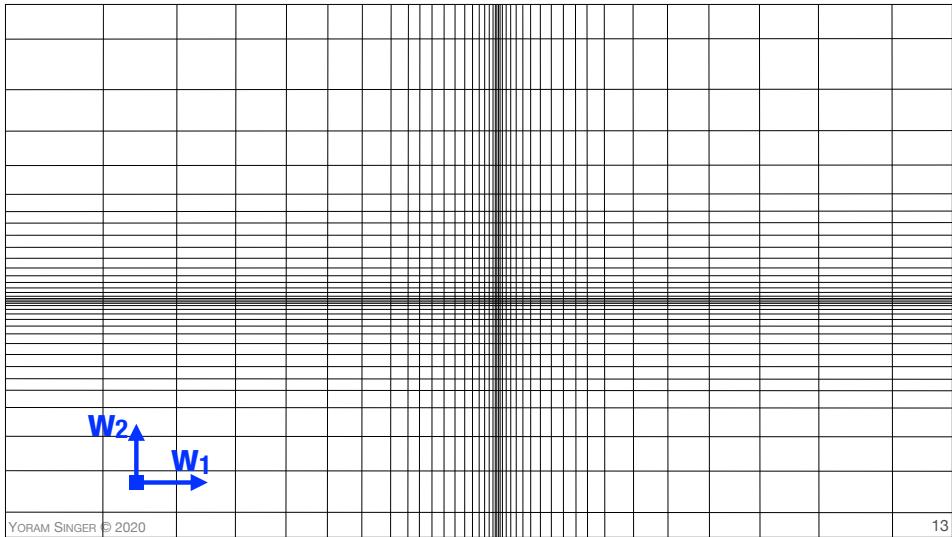
## Finite Precision Reals (Floating Point)



Also : Inf  $\equiv \infty$  and NaN

YORAM SINGER © 2020

12



## Brute Force Enumeration

Given say a neural net structure with  $d$  parameters

## Brute Force Enumeration

Given say a neural net structure with  $d$  parameters

Enumerate all possible floating point values for  $d$  parameters

## Brute Force Enumeration

Given say a neural net structure with  $d$  parameters

Enumerate all possible floating point values for  $d$  parameters

Each parameter has less than  $2^{32}$  different values

## Brute Force Enumeration

Given say a neural net structure with  $d$  parameters

Enumerate all possible floating point values for  $d$  parameters

Each parameter has less than  $2^{32}$  different values

We have “only”  $2^{32d}$  different neural nets  $\Rightarrow f_1, \dots, f_{2^{32d}}$

## Brute Force Enumeration

Given say a neural net structure with  $d$  parameters

Enumerate all possible floating point values for  $d$  parameters

Each parameter has less than  $2^{32}$  different values

We have “only”  $2^{32d}$  different neural nets  $\Rightarrow f_1, \dots, f_{2^{32d}}$

Only  $\tilde{O}(d)$  examples needed to find w.h.p NN w/ the best **generalization**

## Brute Force Enumeration

Given say a neural net structure with  $d$  parameters

Enumerate all possible floating point values for  $d$  parameters

Each parameter has less than  $2^{32}$  different values

We have “only”  $2^{32d}$  different neural nets  $\Rightarrow f_1, \dots, f_{2^{32d}}$

Only  $\tilde{O}(d)$  examples needed to find w.h.p NN w/ the best **generalization**

- We need  $O(\exp(d))$  time

## Brute Force Enumeration

Given say a neural net structure with  $d$  parameters

Enumerate all possible floating point values for  $d$  parameters

Each parameter has less than  $2^{32}$  different values

We have "only"  $2^{32d}$  different neural nets  $\Rightarrow f_1, \dots, f_{2^{32d}}$

Only  $\tilde{O}(d)$  examples needed to find w.h.p NN w/ the best **generalization**

- We need  $O(\exp(d))$  time
- $\tilde{O}(d)$  hides really big constants( $\gg 32$ ) because of  $\epsilon$

## Brute Force Enumeration

Given say a neural net structure with  $d$  parameters

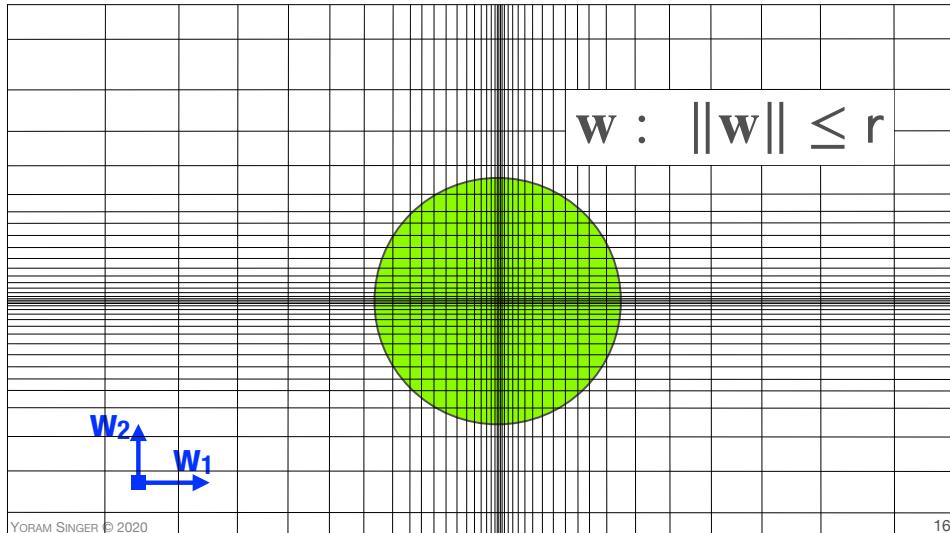
Enumerate all possible floating point values for  $d$  parameters

Each parameter has less than  $2^{32}$  different values

We have "only"  $2^{32d}$  different neural nets  $\Rightarrow f_1, \dots, f_{2^{32d}}$

Only  $\tilde{O}(d)$  examples needed to find w.h.p NN w/ the best **generalization**

- We need  $O(\exp(d))$  time
- $\tilde{O}(d)$  hides really big constants( $\gg 32$ ) because of  $\epsilon$
- For any reasonable amount of examples we are doomed to overfit



## Projection Onto a Ball

Limits space of admissible parameters

Helps in preventing overfitting

Simple & seamless to implement w/ SGD

## Projection Onto a Ball

Limits space of admissible parameters  
Helps in preventing overfitting  
Simple & seamless to implement w/ SGD  
However, semantics of radius is not clear

YORAM SINGER © 2020

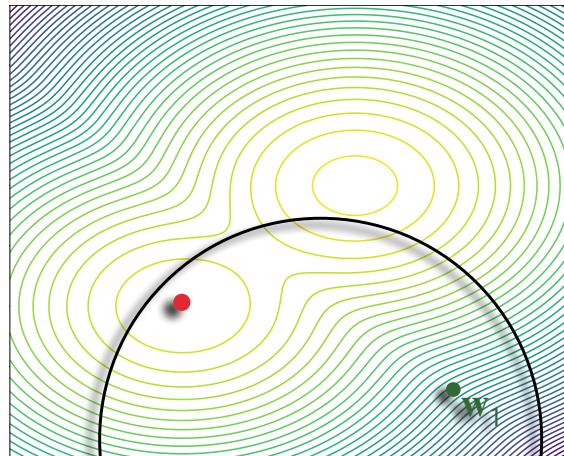
17

## Projection Onto a Ball

Limits space of admissible parameters  
Helps in preventing overfitting  
Simple & seamless to implement w/ SGD  
However, semantics of radius is not clear  
Often ill-defined in non-linear models (e.g NN)

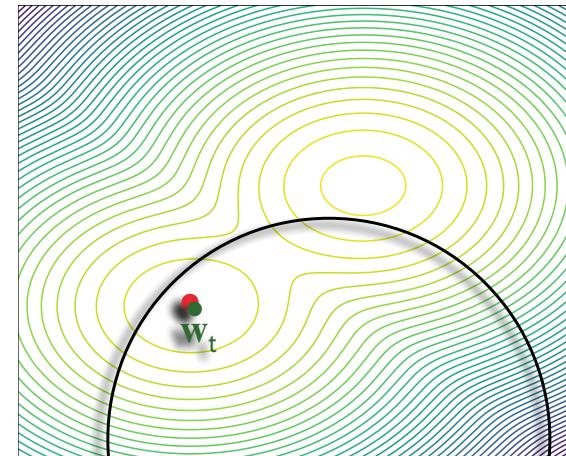
YORAM SINGER © 2020

17



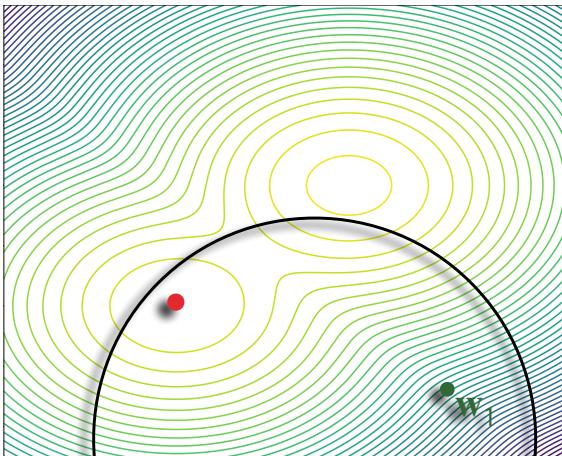
YORAM SINGER © 2020

18



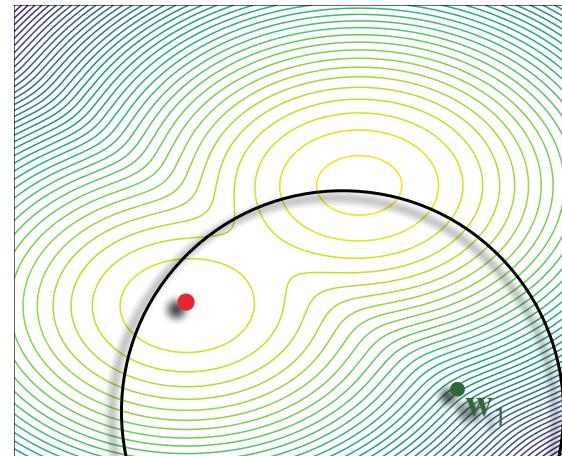
YORAM SINGER © 2020

18



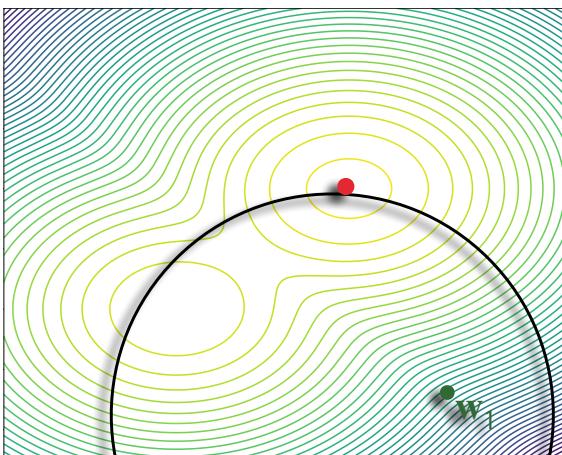
Constrained  
Optimum

19



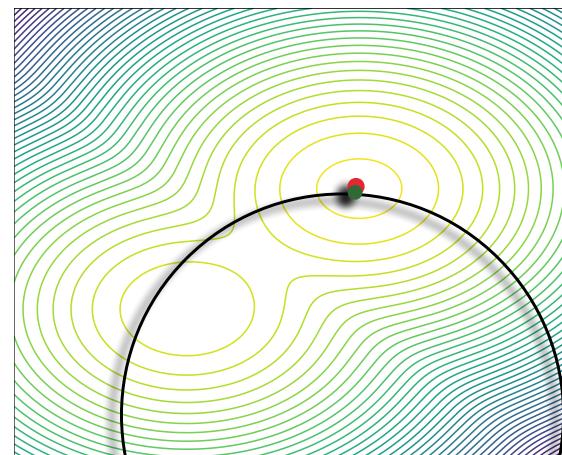
Constrained  
Optimum

19



Constrained  
Optimum

19



Constrained  
Optimum

19

## Regularization

Instead of constrained problem:

$$\min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) \text{ s.t. } \|\mathbf{w}\| \leq r$$

Solve penalized problem:

$$\min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) + \lambda \|\mathbf{w}\|^2 = \mathcal{L}(\mathbf{w}) + \lambda \sum_j w_j^2$$

## Regularization

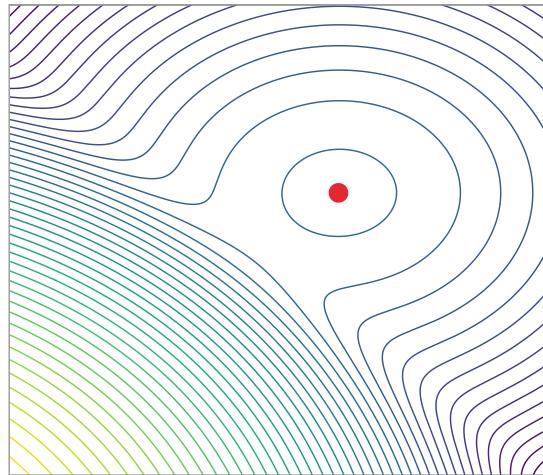
Instead of constrained problem:

$$\min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) \text{ s.t. } \|\mathbf{w}\| \leq r$$

regularization value

Solve penalized problem:

$$\min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) + \lambda \|\mathbf{w}\|^2 = \mathcal{L}(\mathbf{w}) + \lambda \sum_j w_j^2$$



Regularized  
Optimum

## Regularization $\Rightarrow$ $\exists$ Domain Constraints

$$\text{Recall that: } \mathcal{L}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell((\mathbf{x}_i, \mathbf{y}_i), \mathbf{w}) \equiv \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{w})$$

## Regularization $\Rightarrow \exists$ Domain Constraints

$$\text{Recall that: } \mathcal{L}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell((\mathbf{x}_i, y_i), \mathbf{w}) \equiv \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{w})$$

By construction:  $\ell_i(\mathbf{0}) = 1 \Rightarrow \mathcal{L}(\mathbf{0}) = 1$  and  $\ell_i(\mathbf{w}) \geq 0 \Rightarrow \mathcal{L}(\mathbf{w}) \geq 0$

## Regularization $\Rightarrow \exists$ Domain Constraints

$$\text{Recall that: } \mathcal{L}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell((\mathbf{x}_i, y_i), \mathbf{w}) \equiv \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{w})$$

By construction:  $\ell_i(\mathbf{0}) = 1 \Rightarrow \mathcal{L}(\mathbf{0}) = 1$  and  $\ell_i(\mathbf{w}) \geq 0 \Rightarrow \mathcal{L}(\mathbf{w}) \geq 0$

Since we look for a (local) minimum:

$$1 = \mathcal{L}(\mathbf{0}) + \lambda \|\mathbf{0}\|^2 \geq \mathcal{L}(\mathbf{w}^*) + \lambda \|\mathbf{w}^*\|^2 \geq \lambda \|\mathbf{w}^*\|^2$$

## Regularization $\Rightarrow \exists$ Domain Constraints

$$\text{Recall that: } \mathcal{L}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell((\mathbf{x}_i, y_i), \mathbf{w}) \equiv \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{w})$$

By construction:  $\ell_i(\mathbf{0}) = 1 \Rightarrow \mathcal{L}(\mathbf{0}) = 1$  and  $\ell_i(\mathbf{w}) \geq 0 \Rightarrow \mathcal{L}(\mathbf{w}) \geq 0$

Since we look for a (local) minimum:

$$1 = \mathcal{L}(\mathbf{0}) + \lambda \|\mathbf{0}\|^2 \geq \mathcal{L}(\mathbf{w}^*) + \lambda \|\mathbf{w}^*\|^2 \geq \lambda \|\mathbf{w}^*\|^2$$

Therefore  $\mathbf{w}^*$  implicitly satisfies the ball constraint  $\|\mathbf{w}^*\| \leq \frac{1}{\sqrt{\lambda}}$

## Regularization $\Rightarrow \exists$ Domain Constraints

$$\text{Recall that: } \mathcal{L}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell((\mathbf{x}_i, y_i), \mathbf{w}) \equiv \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{w})$$

By construction:  $\ell_i(\mathbf{0}) = 1 \Rightarrow \mathcal{L}(\mathbf{0}) = 1$  and  $\ell_i(\mathbf{w}) \geq 0 \Rightarrow \mathcal{L}(\mathbf{w}) \geq 0$

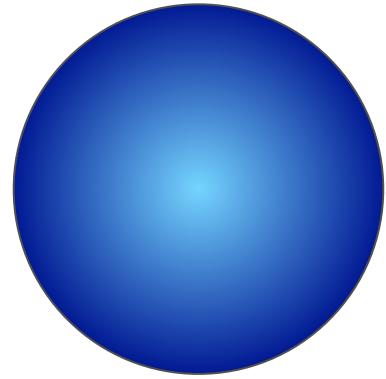
Since we look for a (local) minimum:

$$1 = \mathcal{L}(\mathbf{0}) + \lambda \|\mathbf{0}\|^2 \geq \mathcal{L}(\mathbf{w}^*) + \lambda \|\mathbf{w}^*\|^2 \geq \lambda \|\mathbf{w}^*\|^2$$

Therefore  $\mathbf{w}^*$  implicitly satisfies the ball constraint  $\|\mathbf{w}^*\| \leq \frac{1}{\sqrt{\lambda}}$

However, the converse (domain  $\Rightarrow \exists$  regularization) does not hold

$$\dots + \lambda \|\mathbf{w}\|^2$$



$$\|\mathbf{w}\| \leq \frac{1}{\sqrt{\lambda}}$$

YORAM SINGER © 2020

23

## SGD w/ 2-norm Regularization

$$\frac{\partial}{\partial w_j} \|\mathbf{w}\|^2 = \frac{\partial}{\partial w_j} \sum_i w_i^2 = \frac{\partial w_j^2}{\partial w_j}$$

YORAM SINGER © 2020

24

## SGD w/ 2-norm Regularization

$$\text{Note that } \nabla_{\mathbf{w}} \|\mathbf{w}\|^2 = \left( \frac{\partial w_1^2}{\partial w_1}, \dots, \frac{\partial w_d^2}{\partial w_d} \right) = 2(w_1, \dots, w_d) = 2\mathbf{w}$$

YORAM SINGER © 2020

24

## SGD w/ 2-norm Regularization

$$\text{Note that } \nabla_{\mathbf{w}} \|\mathbf{w}\|^2 = \left( \frac{\partial w_1^2}{\partial w_1}, \dots, \frac{\partial w_d^2}{\partial w_d} \right) = 2(w_1, \dots, w_d) = 2\mathbf{w}$$

$$\text{Add } \lambda \|\mathbf{w}\|^2 \text{ to the objective } \nabla_{\mathbf{w}} [\mathcal{L}(\mathbf{w}) + \lambda \|\mathbf{w}\|^2] = \nabla \mathcal{L}(\mathbf{w}) + 2\lambda \mathbf{w}$$

YORAM SINGER © 2020

24

## SGD w/ 2-norm Regularization

Note that  $\nabla_{\mathbf{w}} \|\mathbf{w}\|^2 = \left( \frac{\partial w_1^2}{\partial w_1}, \dots, \frac{\partial w_d^2}{\partial w_d} \right) = 2(w_1, \dots, w_d) = 2\mathbf{w}$

Add  $\lambda \|\mathbf{w}\|^2$  to the objective  $\nabla_{\mathbf{w}} [\mathcal{L}(\mathbf{w}) + \lambda \|\mathbf{w}\|^2] = \nabla \mathcal{L}(\mathbf{w}) + 2\lambda \mathbf{w}$

In stochastic case we get  $\frac{1}{|S|} \sum_{i \in S} \nabla \ell_i(\mathbf{w}_t) + 2\lambda \mathbf{w}_t$

## SGD w/ 2-norm Regularization

Note that  $\nabla_{\mathbf{w}} \|\mathbf{w}\|^2 = \left( \frac{\partial w_1^2}{\partial w_1}, \dots, \frac{\partial w_d^2}{\partial w_d} \right) = 2(w_1, \dots, w_d) = 2\mathbf{w}$

Add  $\lambda \|\mathbf{w}\|^2$  to the objective  $\nabla_{\mathbf{w}} [\mathcal{L}(\mathbf{w}) + \lambda \|\mathbf{w}\|^2] = \nabla \mathcal{L}(\mathbf{w}) + 2\lambda \mathbf{w}$

In stochastic case we get  $\frac{1}{|S|} \sum_{i \in S} \nabla \ell_i(\mathbf{w}_t) + 2\lambda \mathbf{w}_t$

$\mathbf{g}_t$

## SGD w/ 2-norm Regularization

Note that  $\nabla_{\mathbf{w}} \|\mathbf{w}\|^2 = \left( \frac{\partial w_1^2}{\partial w_1}, \dots, \frac{\partial w_d^2}{\partial w_d} \right) = 2(w_1, \dots, w_d) = 2\mathbf{w}$

Add  $\lambda \|\mathbf{w}\|^2$  to the objective  $\nabla_{\mathbf{w}} [\mathcal{L}(\mathbf{w}) + \lambda \|\mathbf{w}\|^2] = \nabla \mathcal{L}(\mathbf{w}) + 2\lambda \mathbf{w}$

In stochastic case we get  $\frac{1}{|S|} \sum_{i \in S} \nabla \ell_i(\mathbf{w}_t) + 2\lambda \mathbf{w}_t$

SGD update  $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta_t (\mathbf{g}_t + 2\lambda \mathbf{w}_t) = (1 - 2\lambda \eta_t) \mathbf{w}_t - \eta_t \mathbf{g}_t$

## SGD w/ 2-norm Regularization

Note that  $\nabla_{\mathbf{w}} \|\mathbf{w}\|^2 = \left( \frac{\partial w_1^2}{\partial w_1}, \dots, \frac{\partial w_d^2}{\partial w_d} \right) = 2(w_1, \dots, w_d) = 2\mathbf{w}$

Add  $\lambda \|\mathbf{w}\|^2$  to the objective  $\nabla_{\mathbf{w}} [\mathcal{L}(\mathbf{w}) + \lambda \|\mathbf{w}\|^2] = \nabla \mathcal{L}(\mathbf{w}) + 2\lambda \mathbf{w}$

In stochastic case we get  $\frac{1}{|S|} \sum_{i \in S} \nabla \ell_i(\mathbf{w}_t) + 2\lambda \mathbf{w}_t$

SGD update  $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta_t (\mathbf{g}_t + 2\lambda \mathbf{w}_t) = (1 - 2\lambda \eta_t) \mathbf{w}_t - \eta_t \mathbf{g}_t$

weight decay

Suppose  $\ell(y(\mathbf{w} \cdot \mathbf{x})) = [1 - y(\mathbf{w} \cdot \mathbf{x})]_+$

$$\mathcal{L}([1, 0]) = \frac{1}{4}(1 + 0 + 0 + 2) = \frac{3}{4}$$

$$\mathcal{L}([0, 1]) = \frac{1}{4}(2 + 0 + 0 + 1) = \frac{3}{4}$$

$$\mathcal{L}([0.5, 0.5]) = \frac{1}{4}(1.5 + 0 + 0 + 1.5) = \frac{3}{4}$$

$$\mathcal{L}([1, 0]) = \mathcal{L}([0, 1]) = \mathcal{L}([0.5, 0.5]) = \frac{3}{4} \quad \dots \text{no preference} \dots$$

	x[0]	x[1]	y
Ex. 1	0	1	-1
Ex. 2	1	1	1
Ex. 3	-1	-1	-1
Ex. 4	-1	0	1

Regularized version:  $\mathcal{L}(\mathbf{w}) + 10^{-6} \|\mathbf{w}\|^2$

$$\mathcal{L}([0.5, 0.5]) + 10^{-6} [0.5^2 + 0.5^2] = \frac{3}{4} + \frac{1}{2} 10^{-6}$$

$$\mathcal{L}([1, 0]) + 10^{-6} [1^2 + 0^2] = \frac{3}{4} + 10^{-6}$$

$$\mathcal{L}([0, 1]) + 10^{-6} [0^2 + 1^2] = \frac{3}{4} + 10^{-6}$$

|| ||<sub>2</sub>  
regularization  
is Inclusive



## SGD w/ 1-norm Regularization

Recall  $\|\mathbf{w}\|_1 = \sum_j |w_j|$  and derivative at  $w_j = 0$  is not uniquely defined

## SGD w/ 1-norm Regularization

Recall  $\|\mathbf{w}\|_1 = \sum_j |w_j|$  and derivative at  $w_j = 0$  is not uniquely defined

Start with a simple single variable problem:  $\min_w \mathcal{L}(w) + \lambda w$  s.t.  $w \geq 0$

## SGD w/ 1-norm Regularization

Recall  $\|\mathbf{w}\|_1 = \sum_j |w_j|$  and derivative at  $w_j = 0$  is not uniquely defined

Start with a simple single variable problem:  $\min_w \mathcal{L}(w) + \lambda w$  s.t  $w \geq 0$

$$\text{Let } g_t = \frac{d\mathcal{L}(w)}{dw} \Big|_{w=w_t} \quad \text{For } w > 0 \text{ derivative well defined } \frac{d}{dw} \lambda w = \lambda$$

## SGD w/ 1-norm Regularization

Recall  $\|\mathbf{w}\|_1 = \sum_j |w_j|$  and derivative at  $w_j = 0$  is not uniquely defined

Start with a simple single variable problem:  $\min_w \mathcal{L}(w) + \lambda w$  s.t  $w \geq 0$

$$\text{Let } g_t = \frac{d\mathcal{L}(w)}{dw} \Big|_{w=w_t} \quad \text{For } w > 0 \text{ derivative well defined } \frac{d}{dw} \lambda w = \lambda$$

If  $w_t$  and  $w_{t+1}$  are positive update is :  $w_{t+1} \leftarrow w_t - \eta_t(g_t + \lambda)$

## SGD w/ 1-norm Regularization

Recall  $\|\mathbf{w}\|_1 = \sum_j |w_j|$  and derivative at  $w_j = 0$  is not uniquely defined

Start with a simple single variable problem:  $\min_w \mathcal{L}(w) + \lambda w$  s.t  $w \geq 0$

$$\text{Let } g_t = \frac{d\mathcal{L}(w)}{dw} \Big|_{w=w_t} \quad \text{For } w > 0 \text{ derivative well defined } \frac{d}{dw} \lambda w = \lambda$$

If  $w_t$  and  $w_{t+1}$  are positive update is :  $w_{t+1} \leftarrow w_t - \eta_t(g_t + \lambda)$

However if  $w_t - \eta_t(g_t + \lambda) \leq 0 \Leftrightarrow w_t - \eta_t g_t \leq \eta_t \lambda$  we must stop at 0

## SGD w/ 1-norm Regularization

Recall  $\|\mathbf{w}\|_1 = \sum_j |w_j|$  and derivative at  $w_j = 0$  is not uniquely defined

Start with a simple single variable problem:  $\min_w \mathcal{L}(w) + \lambda w$  s.t  $w \geq 0$

$$\text{Let } g_t = \frac{d\mathcal{L}(w)}{dw} \Big|_{w=w_t} \quad \text{For } w > 0 \text{ derivative well defined } \frac{d}{dw} \lambda w = \lambda$$

If  $w_t$  and  $w_{t+1}$  are positive update is :  $w_{t+1} \leftarrow w_t - \eta_t(g_t + \lambda)$

However if  $w_t - \eta_t(g_t + \lambda) \leq 0 \Leftrightarrow w_t - \eta_t g_t \leq \eta_t \lambda$  we must stop at 0

$$w_{t+1} \leftarrow \max \{w_t - \eta_t g_t - \eta_t \lambda, 0\}$$

## SGD w/ 1-norm Regularization

Similarly constraining  $\mathbf{w} \leq 0$  while penalizing for  $\lambda |\mathbf{w}|$  yields the problem:

$$\min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) - \lambda \mathbf{w} \text{ s.t. } \mathbf{w} \leq 0$$

## SGD w/ 1-norm Regularization

Similarly constraining  $\mathbf{w} \leq 0$  while penalizing for  $\lambda |\mathbf{w}|$  yields the problem:

$$\min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) - \lambda \mathbf{w} \text{ s.t. } \mathbf{w} \leq 0$$

Assuming  $\mathbf{w}_t$  and  $\mathbf{w}_{t+1}$  are *negative* update is :  $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta_t(\mathbf{g}_t - \lambda)$

## SGD w/ 1-norm Regularization

Similarly constraining  $\mathbf{w} \leq 0$  while penalizing for  $\lambda |\mathbf{w}|$  yields the problem:

$$\min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) - \lambda \mathbf{w} \text{ s.t. } \mathbf{w} \leq 0$$

Assuming  $\mathbf{w}_t$  and  $\mathbf{w}_{t+1}$  are *negative* update is :  $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta_t(\mathbf{g}_t - \lambda)$

However if  $\mathbf{w}_t - \eta_t(\mathbf{g}_t - \lambda) \geq 0 \Leftrightarrow \mathbf{w}_t - \eta_t \mathbf{g}_t \geq -\eta_t \lambda$  we must stop at 0

## SGD w/ 1-norm Regularization

Similarly constraining  $\mathbf{w} \leq 0$  while penalizing for  $\lambda |\mathbf{w}|$  yields the problem:

$$\min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) - \lambda \mathbf{w} \text{ s.t. } \mathbf{w} \leq 0$$

Assuming  $\mathbf{w}_t$  and  $\mathbf{w}_{t+1}$  are *negative* update is :  $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta_t(\mathbf{g}_t - \lambda)$

However if  $\mathbf{w}_t - \eta_t(\mathbf{g}_t - \lambda) \geq 0 \Leftrightarrow \mathbf{w}_t - \eta_t \mathbf{g}_t \geq -\eta_t \lambda$  we must stop at 0

## SGD w/ 1-norm Regularization

Similarly constraining  $w \leq 0$  while penalizing for  $\lambda |w|$  yields the problem:

$$\min_w \mathcal{L}(w) - \lambda w \text{ s.t. } w \leq 0$$

Assuming  $w_t$  and  $w_{t+1}$  are negative update is :  $w_{t+1} \leftarrow w_t - \eta_t(g_t - \lambda)$

However if  $w_t - \eta_t(g_t - \lambda) \geq 0 \Leftrightarrow w_t - \eta_t g_t \geq -\eta_t \lambda$  we must stop at 0

$$w_{t+1} \leftarrow \min \{w_t - \eta_t g_t + \eta_t \lambda, 0\}$$

## SGD w/ 1-norm Regularization

Combining the two problems:

$$w^+ = \arg \min_w \mathcal{L}(w) + \lambda w \text{ s.t. } w \geq 0 \quad w^- = \arg \min_w \mathcal{L}(w) - \lambda w \text{ s.t. } w \leq 0$$

## SGD w/ 1-norm Regularization

Combining the two problems:

$$w^+ = \arg \min_w \mathcal{L}(w) + \lambda w \text{ s.t. } w \geq 0 \quad w^- = \arg \min_w \mathcal{L}(w) - \lambda w \text{ s.t. } w \leq 0$$

- ▶  $w^+ > 0 \text{ & } w^- = 0 \Rightarrow w^* > 0$

## SGD w/ 1-norm Regularization

Combining the two problems:

$$w^+ = \arg \min_w \mathcal{L}(w) + \lambda w \text{ s.t. } w \geq 0 \quad w^- = \arg \min_w \mathcal{L}(w) - \lambda w \text{ s.t. } w \leq 0$$

- ▶  $w^+ > 0 \text{ & } w^- = 0 \Rightarrow w^* > 0$
- ▶  $w^+ = 0 \text{ & } w^- < 0 \Rightarrow w^* < 0$

## SGD w/ 1-norm Regularization

Combining the two problems:

$$\mathbf{w}^+ = \arg \min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) + \lambda \mathbf{w} \text{ s.t } \mathbf{w} \geq 0 \quad \mathbf{w}^- = \arg \min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) - \lambda \mathbf{w} \text{ s.t } \mathbf{w} \leq 0$$

- $\mathbf{w}^+ > 0 \text{ & } \mathbf{w}^- = 0 \Rightarrow \mathbf{w}^* > 0$
- $\mathbf{w}^+ = 0 \text{ & } \mathbf{w}^- < 0 \Rightarrow \mathbf{w}^* < 0$
- $\mathbf{w}^+ = 0 \text{ & } \mathbf{w}^- = 0 \Rightarrow \mathbf{w}^* = 0$

## SGD w/ 1-norm Regularization

Construct  $\min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) + \lambda |\mathbf{w}|$  by combining the two cases

## SGD w/ 1-norm Regularization

Construct  $\min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) + \lambda |\mathbf{w}|$  by combining the two cases

If  $-\eta_t \lambda \leq \mathbf{w}_t - \eta_t \mathbf{g}_t \leq \eta_t \lambda \Leftrightarrow |\mathbf{w}_t - \eta_t \mathbf{g}_t| \leq \eta_t \lambda$  then

## SGD w/ 1-norm Regularization

Construct  $\min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) + \lambda |\mathbf{w}|$  by combining the two cases

If  $-\eta_t \lambda \leq \mathbf{w}_t - \eta_t \mathbf{g}_t \leq \eta_t \lambda \Leftrightarrow |\mathbf{w}_t - \eta_t \mathbf{g}_t| \leq \eta_t \lambda$  then

$$\mathbf{w}_{t+1} = 0$$

## SGD w/ 1-norm Regularization

Construct  $\min_w \mathcal{L}(w) + \lambda |w|$  by combining the two cases

If  $-\eta_t \lambda \leq w_t - \eta_t g_t \leq \eta_t \lambda \Leftrightarrow |w_t - \eta_t g_t| \leq \eta_t \lambda$  then

$$w_{t+1} = 0$$

If  $w_t - \eta_t g_t > \eta_t \lambda$  then

30

## SGD w/ 1-norm Regularization

Construct  $\min_w \mathcal{L}(w) + \lambda |w|$  by combining the two cases

If  $-\eta_t \lambda \leq w_t - \eta_t g_t \leq \eta_t \lambda \Leftrightarrow |w_t - \eta_t g_t| \leq \eta_t \lambda$  then

$$w_{t+1} = 0$$

If  $w_t - \eta_t g_t > \eta_t \lambda$  then

$$w_{t+1} \leftarrow w_t - \eta_t(g_t + \lambda)$$

YORAM SINGER © 2020

30

## SGD w/ 1-norm Regularization

Construct  $\min_w \mathcal{L}(w) + \lambda |w|$  by combining the two cases

If  $-\eta_t \lambda \leq w_t - \eta_t g_t \leq \eta_t \lambda \Leftrightarrow |w_t - \eta_t g_t| \leq \eta_t \lambda$  then

$$w_{t+1} = 0$$

If  $w_t - \eta_t g_t > \eta_t \lambda$  then

$$w_{t+1} \leftarrow w_t - \eta_t(g_t + \lambda)$$

If  $w_t - \eta_t g_t < -\eta_t \lambda$  then

30

## SGD w/ 1-norm Regularization

Construct  $\min_w \mathcal{L}(w) + \lambda |w|$  by combining the two cases

If  $-\eta_t \lambda \leq w_t - \eta_t g_t \leq \eta_t \lambda \Leftrightarrow |w_t - \eta_t g_t| \leq \eta_t \lambda$  then

$$w_{t+1} = 0$$

If  $w_t - \eta_t g_t > \eta_t \lambda$  then

$$w_{t+1} \leftarrow w_t - \eta_t(g_t + \lambda)$$

If  $w_t - \eta_t g_t < -\eta_t \lambda$  then

$$w_{t+1} \leftarrow w_t - \eta_t(g_t - \lambda)$$

YORAM SINGER © 2020

30

## SGD w/ 1-norm Regularization

To update  $\min_w \mathcal{L}(w) + \lambda \|w\|_1$  we perform

$$\text{i. } w_j^{t+1/2} \leftarrow w_j^t - \eta_t g_j^t$$

$$\text{ii. } w_j^{t+1} \leftarrow (\text{sign}(w_j^{t+1/2}) [ |w_j^{t+1/2}| - \eta_t \lambda ])_+$$

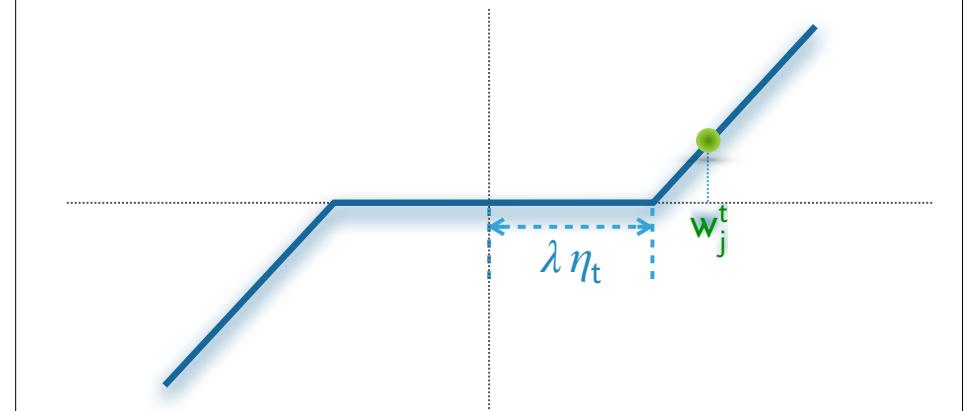
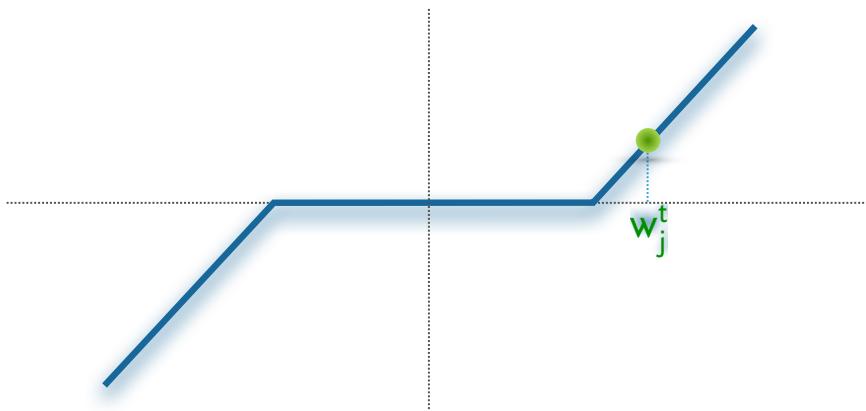
## SGD w/ 1-norm Regularization

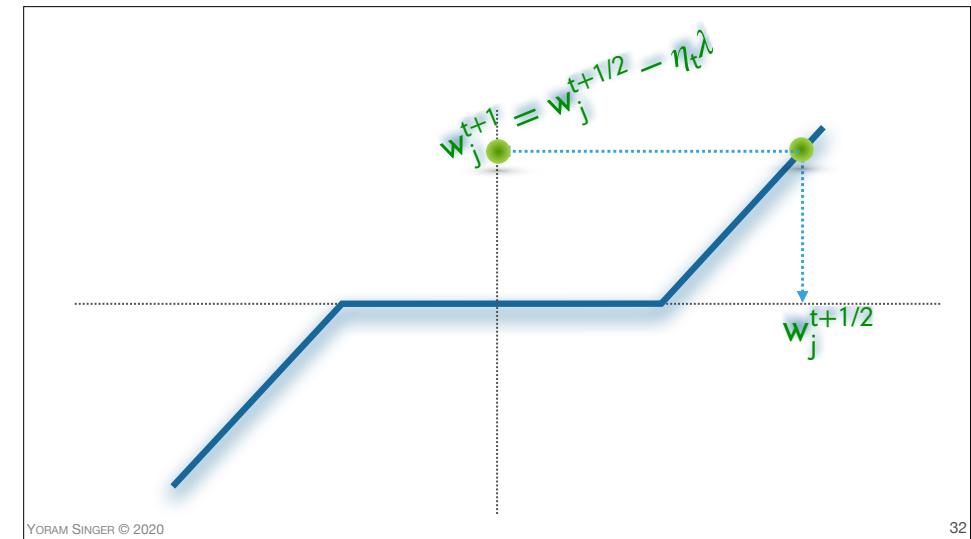
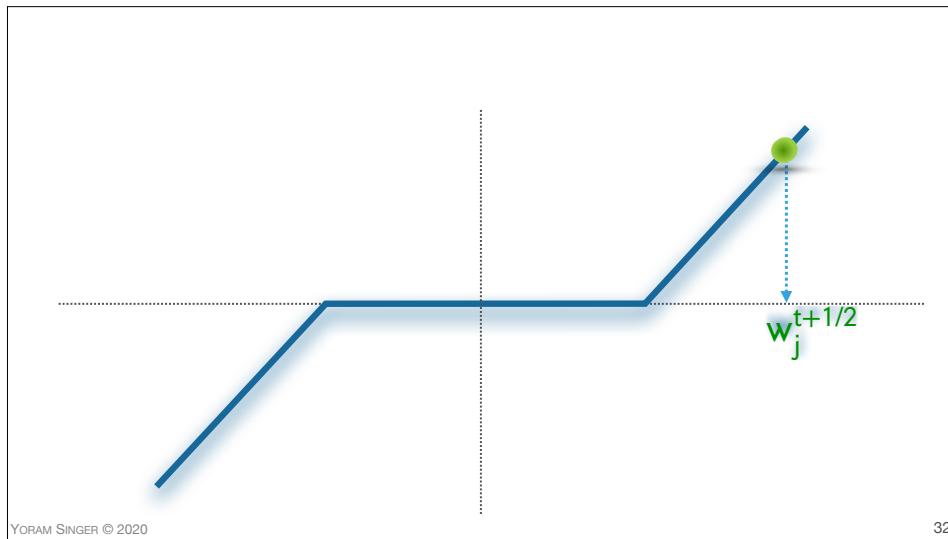
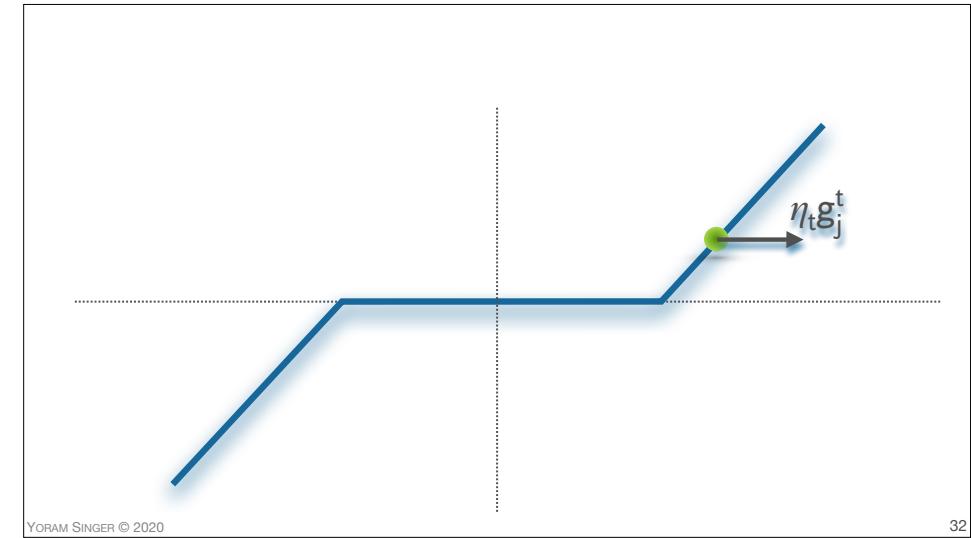
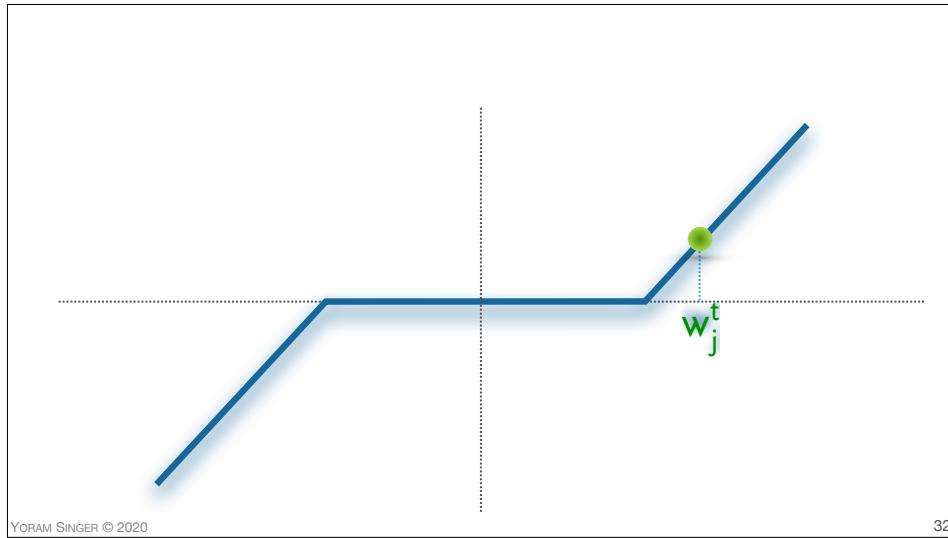
To update  $\min_w \mathcal{L}(w) + \lambda \|w\|_1$  we perform

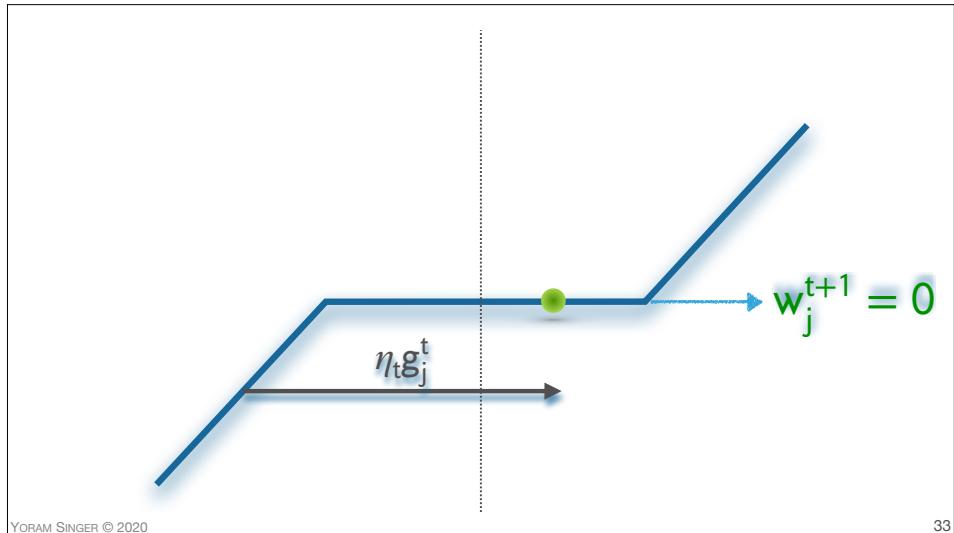
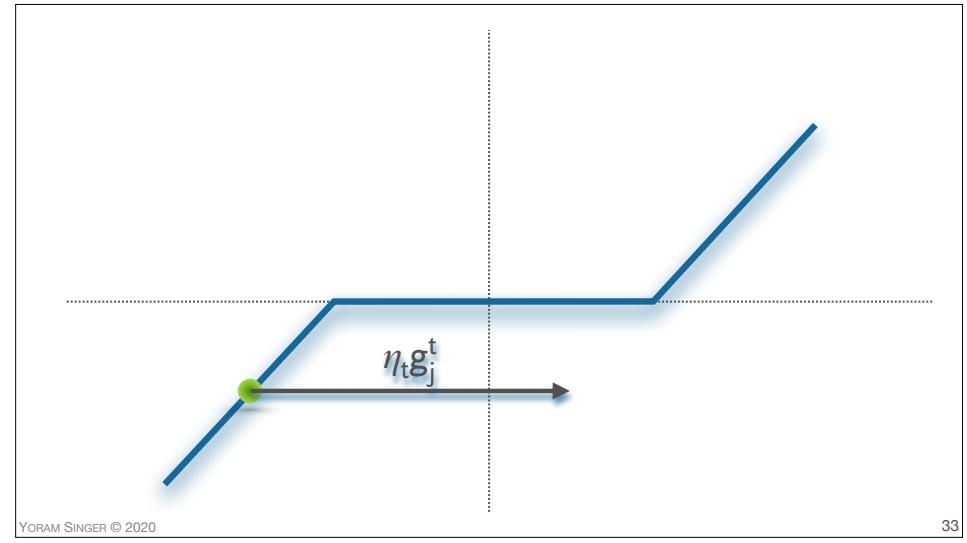
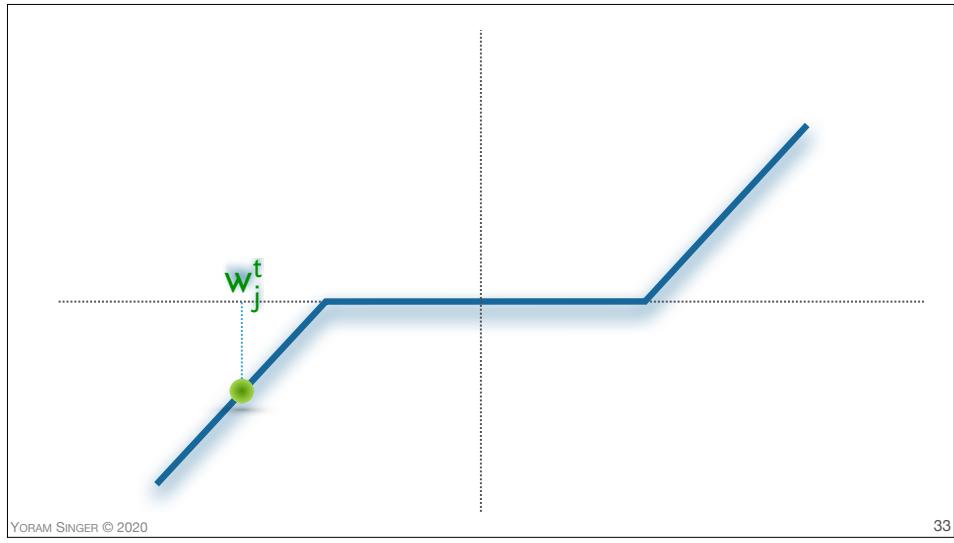
$$\text{i. } w_j^{t+1/2} \leftarrow w_j^t - \eta_t g_j^t$$

$$\text{ii. } w_j^{t+1} \leftarrow (\text{sign}(w_j^{t+1/2}) [ |w_j^{t+1/2}| - \eta_t \lambda ])_+$$

$$\text{i. } w_j^{t+1/2} \leftarrow w_j^t - \eta_t g_j^t \quad \text{ii. } w_j^{t+1} \leftarrow (\text{sign}(w_j^{t+1/2}) [ |w_j^{t+1/2}| - \eta_t \lambda ])_+$$







Suppose  $\ell(y(\mathbf{w} \cdot \mathbf{x})) = [1 - y(\mathbf{w} \cdot \mathbf{x})]_+$

$$\mathcal{L}(1, 0) = \frac{1}{4}(1 + 0 + 0 + 2) = \frac{3}{4}$$

$$\mathcal{L}(0, 1) = \frac{1}{4}(2 \cdot 0.08 + 0 + 0 + 1) = \frac{3}{4} + 0.02$$

$$\mathcal{L}(0.5, 0.5) = \frac{1}{4}(1 \cdot 1.54 + 0 + 0 + 1 \cdot 1.5) = \frac{3}{4} + 0.01$$

$$\mathcal{L}(0, 1) > \mathcal{L}(0.5, 0.5) > \mathcal{L}(1, 0)$$

$x[0]$	$x[1]$	$y$
0	1.08	-1
1	1	1
-1	-1	-1
-1	0	1

YORAM SINGER © 2020

34

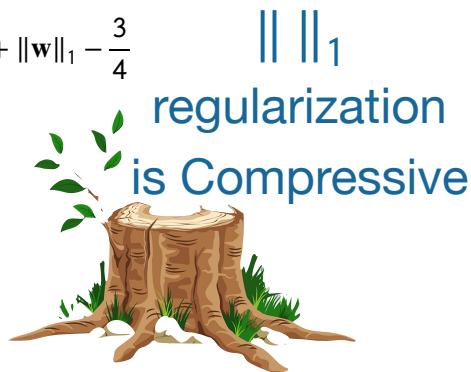
Regularized objective:  $\mathcal{Q}(\mathbf{w}) = \mathcal{L}(\mathbf{w}) + \|\mathbf{w}\|_1 - \frac{3}{4}$

$$\mathcal{Q}(1, 0) = \frac{3}{4} - \frac{3}{4} + 1 = 1$$

$$\mathcal{Q}(0, 1) = 0.02 + 1 = 1.02$$

$$\mathcal{Q}(0.5, 0.5) = 0.01 + 1 = 1.01$$

$$\mathcal{Q}(0, 1) > \mathcal{Q}(0.5, 0.5) > \mathcal{Q}(1, 0)$$



YORAM SINGER © 2020

35

Regularized objective:  $\mathcal{Q}(\mathbf{w}) = \mathcal{L}(\mathbf{w}) + \|\mathbf{w}\|_1 - \frac{3}{4}$

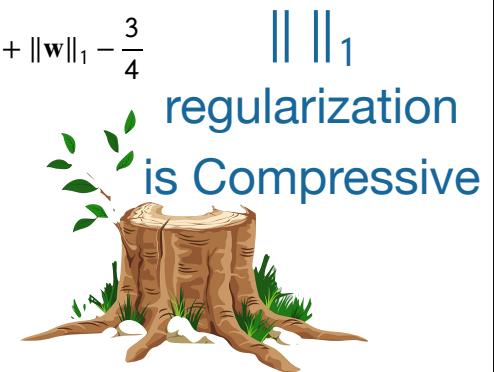
$$\mathcal{Q}(1, 0) = \frac{3}{4} - \frac{3}{4} + 1 = 1$$

$$\mathcal{Q}(0, 1) = 0.02 + 1 = 1.02$$

$$\mathcal{Q}(0.5, 0.5) = 0.01 + 1 = 1.01$$

$$\mathcal{Q}(0, 1) > \mathcal{Q}(0.5, 0.5) > \mathcal{Q}(1, 0)$$

$$\mathcal{Q}_2(\mathbf{w}) = \mathcal{L}(\mathbf{w}) + \|\mathbf{w}\|^2 - \frac{3}{4} \quad \mathcal{Q}_2(1, 0) = 1 \quad \mathcal{Q}_2(0, 1) = 1.02$$



YORAM SINGER © 2020

35

Regularized objective:  $\mathcal{Q}(\mathbf{w}) = \mathcal{L}(\mathbf{w}) + \|\mathbf{w}\|_1 - \frac{3}{4}$

$$\mathcal{Q}(1, 0) = \frac{3}{4} - \frac{3}{4} + 1 = 1$$

$$\mathcal{Q}(0, 1) = 0.02 + 1 = 1.02$$

$$\mathcal{Q}(0.5, 0.5) = 0.01 + 1 = 1.01$$

$$\mathcal{Q}(0, 1) > \mathcal{Q}(0.5, 0.5) > \mathcal{Q}(1, 0)$$

$$\mathcal{Q}_2(\mathbf{w}) = \mathcal{L}(\mathbf{w}) + \|\mathbf{w}\|^2 - \frac{3}{4} \quad \mathcal{Q}_2(1, 0) = 1 \quad \mathcal{Q}_2(0, 1) = 1.02$$

$$\mathcal{Q}_2(0.5, 0.5) = \textcolor{red}{0.51}$$

YORAM SINGER © 2020

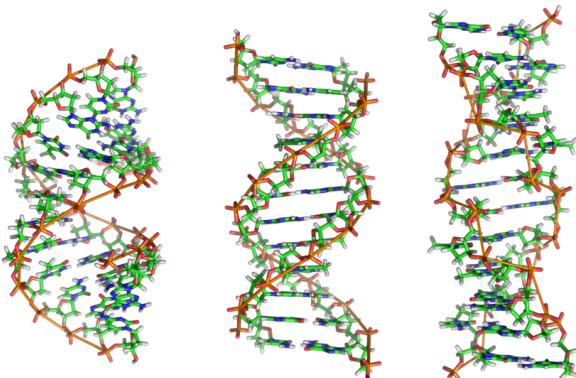
35

	$\ \mathbf{w}\ _2^2$	$\ \mathbf{w}\ _1$
Computation	Scaling	Shrinkage
Promotes	Diversity	Sparsity
$n \gg d$	✓✓✓	✗
$n \sim d$	✓✓	✓✓
$d \gg n$	✓	✓✓✓

YORAM SINGER © 2020

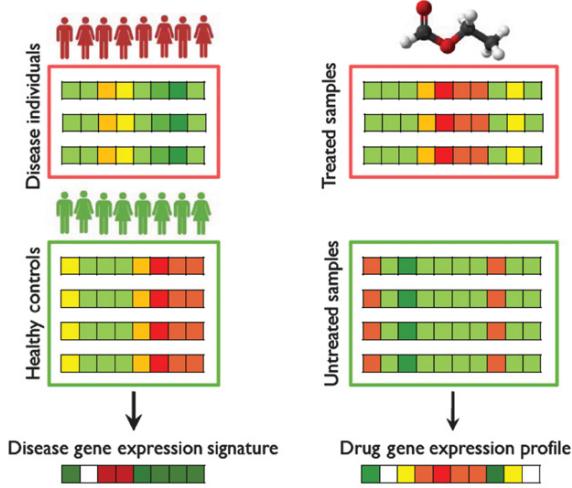
36

## Gene Expression



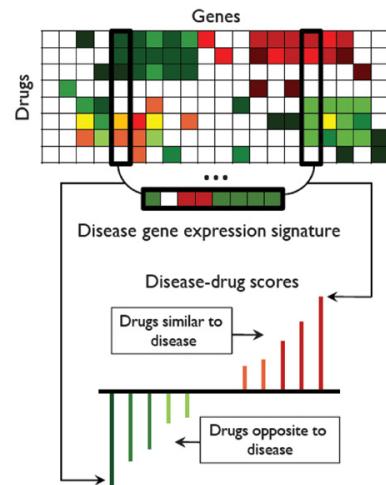
YORAM SINGER © 2020

37



YORAM SINGER © 2020

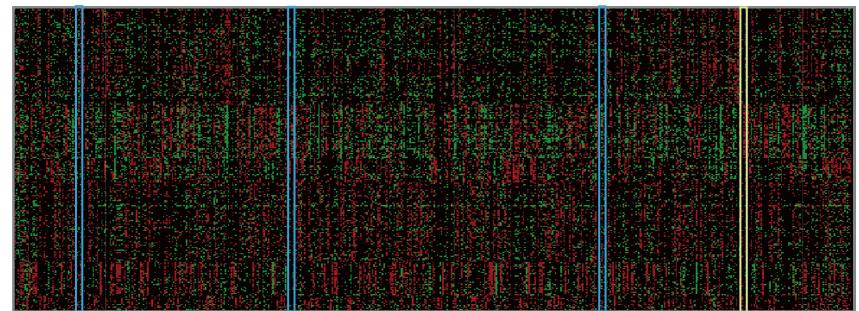
38



YORAM SINGER © 2020

39

## Experimental Expression Matrix

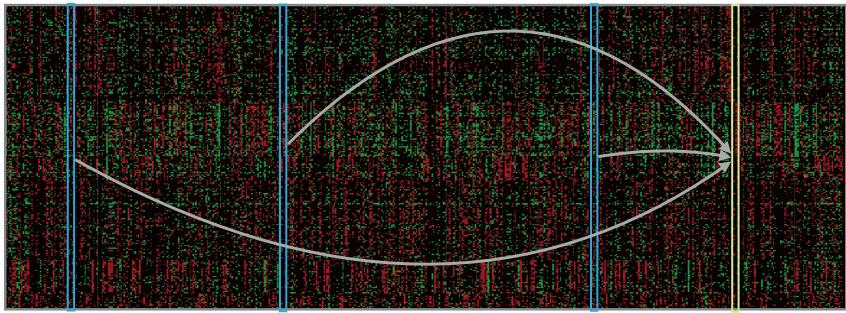


Find few (#genes ≪ #drugs)  
highly predictive genes

YORAM SINGER © 2020

40

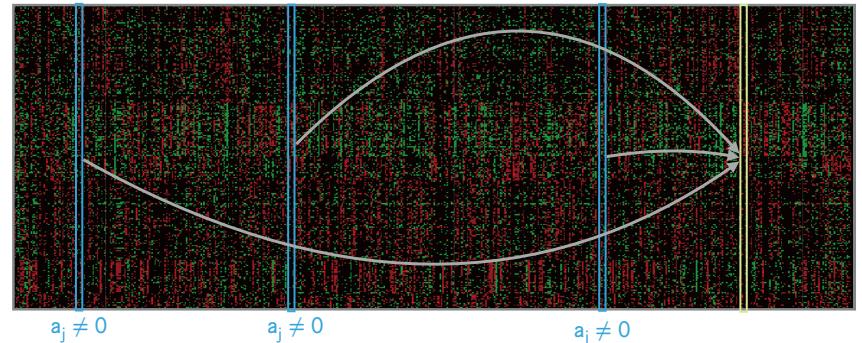
## Experimental Expression Matrix



Find few (#genes ≪ #drugs)  
highly predictive genes

YORAM SINGER © 2020

## Experimental Expression Matrix



$a_j \neq 0$   
 $a_j \neq 0$   
 $a_j \neq 0$

Find few (#genes ≪ #drugs)  
highly predictive genes

$$Q(\mathbf{a}) = \mathcal{L}(\mathbf{a}) + \lambda \|\mathbf{a}\|_1 = \left\| \sum_j a_j \mathbf{c}_j - \mathbf{c}_r \right\|^2 + \lambda \|\mathbf{a}\|_1$$

40

```
# GD with L1 & L2 regularization
def proxgd(X, eta, l1, l2):
    w = np.zeros((X.shape[1], 1))
    ls = []
    for e in range(len(eta)):
        et = eta[e]
        z = X @ w
        ls.append(loss(z, w, l1, l2))
        q = 1 / (1 + np.exp(-z))
        w = w - et * np.mean(X * q, axis=0, keepdims=True).T
        if l1 > 0:
            w = np.sign(w) * np.maximum(abs(w) - et * l1, 0)
        if l2 > 0:
            w = w / (1 + et * l2)
    return w, ls
```

YORAM SINGER © 2020

41