

Gradient Descent Solutions to Least-Square Problems

Objectives

1

- explain need for iterative algorithms
- derive gradient descent algorithm
- consider impact of step size on convergence
- introduce notion of convex functions

Iterative solution methods play an important role ²

Features/labels: $\underline{x}_i, d_i, i=1, 2, \dots N$

Classifier or model error: $e^2 = \sum_{i=1}^N (\underline{x}_i^\top \underline{w} - d_i)^2$

$$\underline{A} = \begin{bmatrix} \underline{x}_1^\top \\ \underline{x}_2^\top \\ \vdots \\ \underline{x}_N^\top \end{bmatrix} \quad \underline{d} = \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_N \end{bmatrix} \quad e^2 = \|\underline{A}\underline{w} - \underline{d}\|_2^2$$

Regularized least squares: $\arg \min_{\underline{w}} \|\underline{A}\underline{w} - \underline{d}\|_2^2 + \lambda r(\underline{w})$

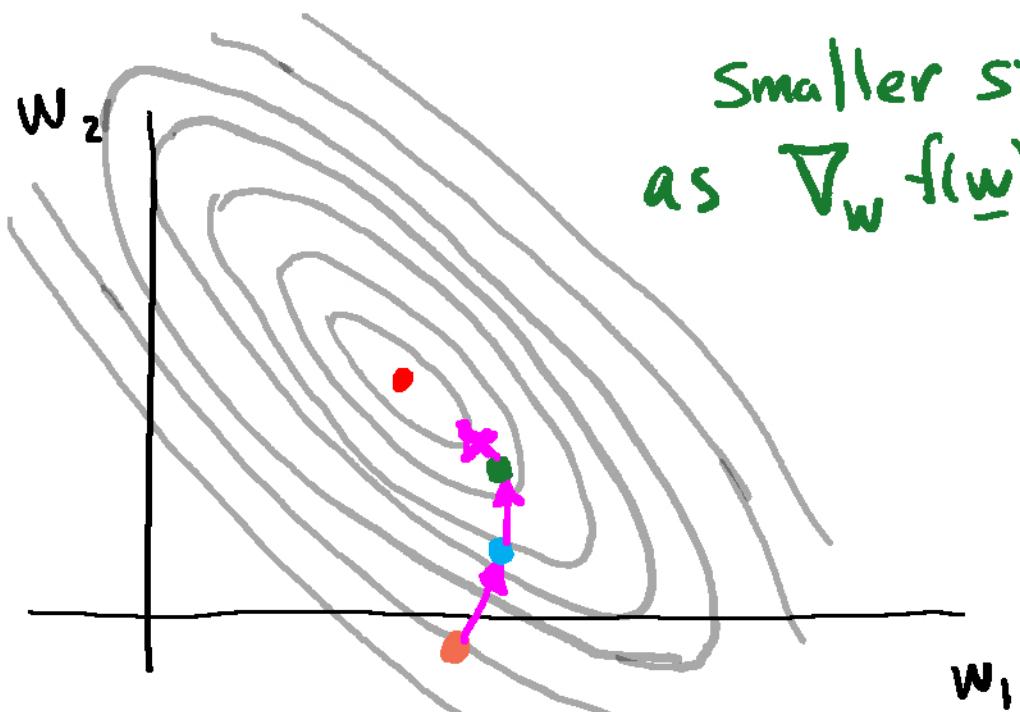
1. Computational cost $(\underline{A}^\top \underline{A})^{-1}$
 2. Closed form solution maybe unavailable
 3. Adapt \underline{w} to new features/labels
- } develop iterative approach

Gradient descent finds the minimum

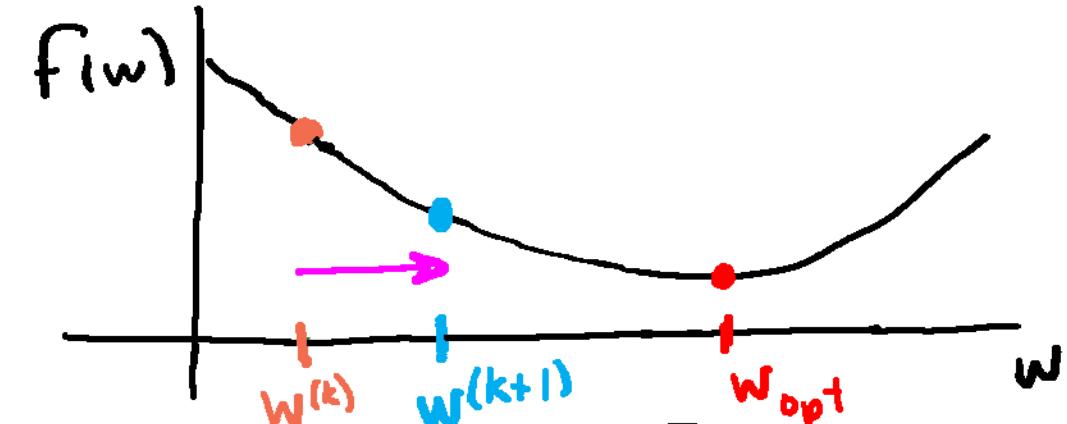
$$f(\underline{w}) = \|\underline{A}\underline{w} - \underline{d}\|_2^2$$

$$\underline{w}^{(k+1)} = \underline{w}^{(k)} - \tau' \nabla_{\underline{w}} f(\underline{w})$$

step size gradient
($\tau' > 0$)



Smaller steps
as $\nabla_{\underline{w}} f(\underline{w}) \downarrow$

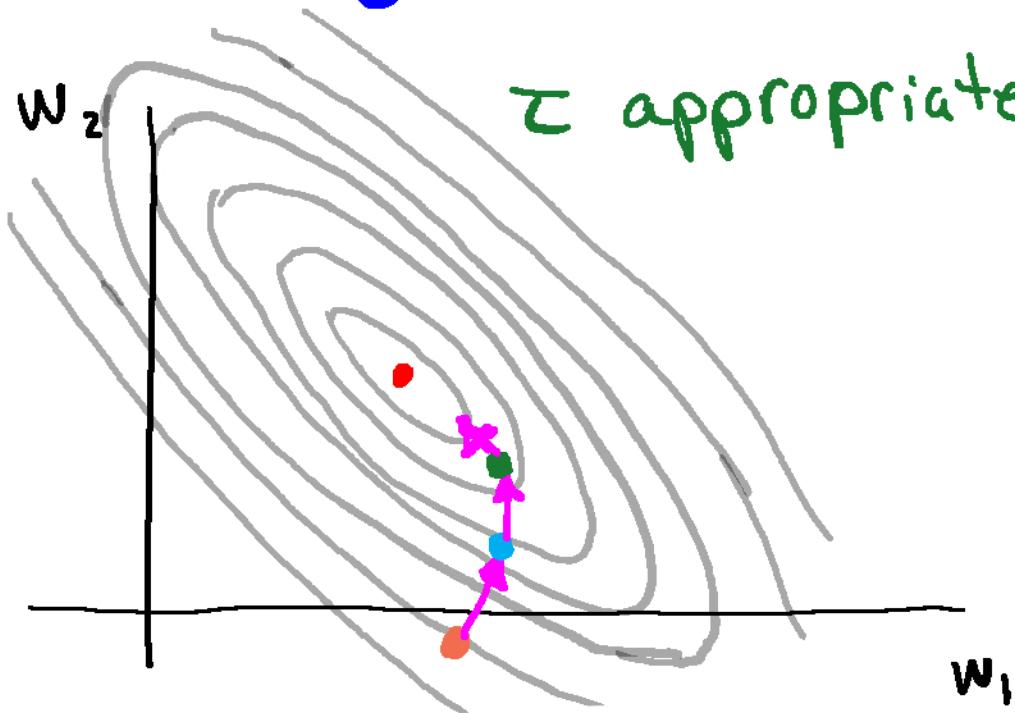


$$\begin{aligned} f(\underline{w}) &= (\underline{A}\underline{w} - \underline{d})^T (\underline{A}\underline{w} - \underline{d}) \\ &= \underline{w}^T \underline{A}^T \underline{A} \underline{w} - 2 \underline{w}^T \underline{A}^T \underline{d} + \underline{d}^T \underline{d} \end{aligned}$$

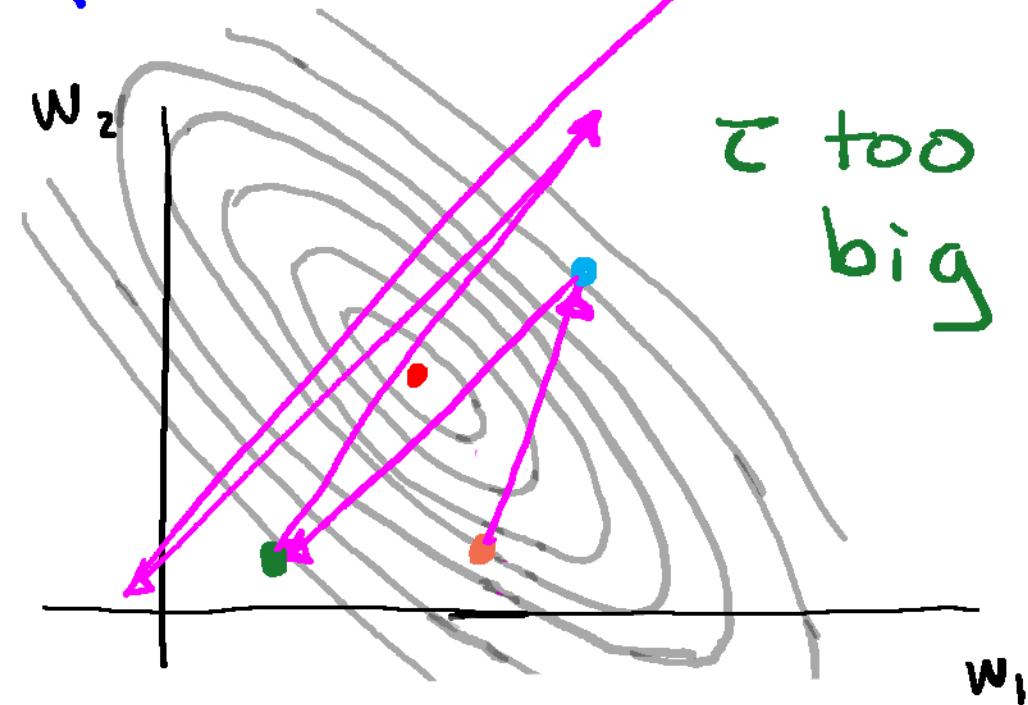
$$\begin{aligned} \nabla_{\underline{w}} f(\underline{w}) &= 2 \underline{A}^T \underline{A} \underline{w} - 2 \underline{A}^T \underline{d} \\ &= 2 \underline{A}^T (\underline{A} \underline{w} - \underline{d}) \end{aligned}$$

$$\begin{aligned} \underline{w}^{(k+1)} &= \underline{w}^{(k)} - \tau \underline{A}^T (\underline{A} \underline{w}^{(k)} - \underline{d}) \\ &\quad \text{(Landweber iteration)} \end{aligned}$$

Convergence behavior depends on τ 4



τ appropriate



τ too big

τ too small: slow convergence
 τ too big: no convergence
unstable!

Require $0 < \tau < 2/\|\underline{A}\|_{\text{op}}^2$ for convergence 5

Recall $\|\underline{A}\|_{\text{op}} = \|\underline{A}\|_2 = \sigma_{\max}(\underline{A})$

Convergence: $f(\underline{w}^{(k+1)}) < f(\underline{w}^{(k)})$ cost decreases as k increases

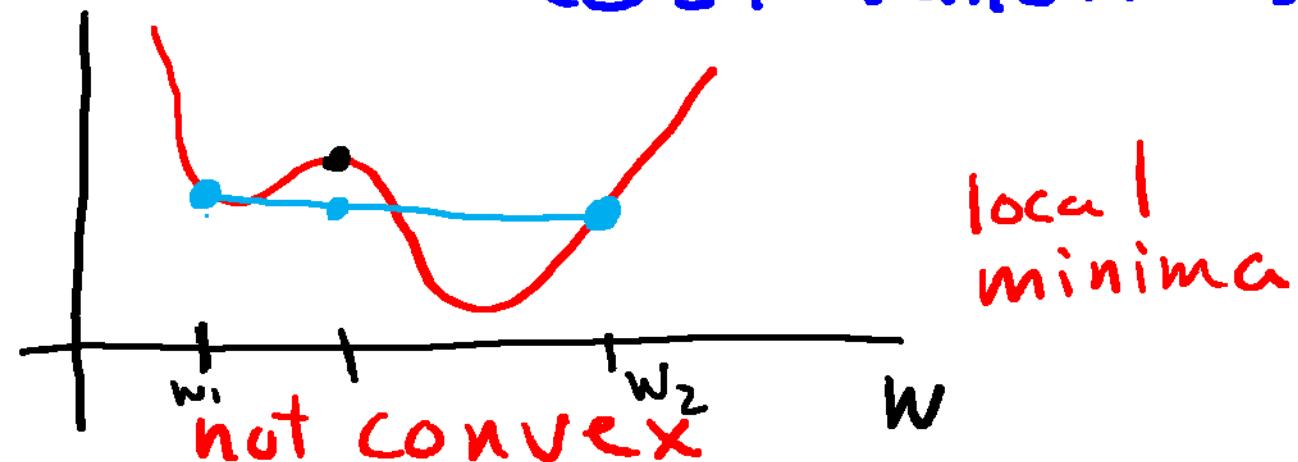
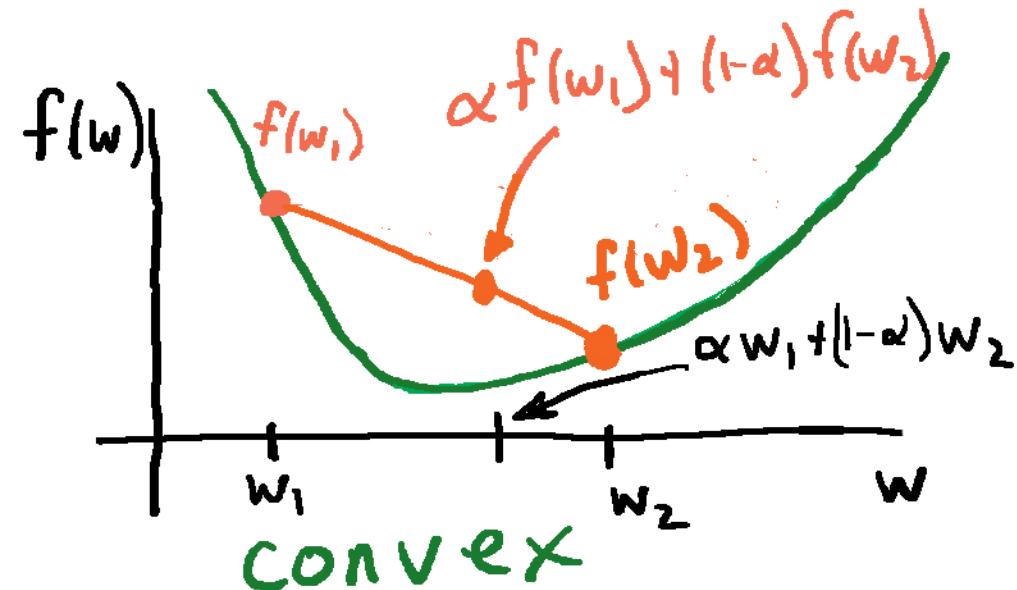
$$\|\underline{A}\underline{w}^{(k+1)} - \underline{d}\|_2^2 < \|\underline{A}\underline{w}^{(k)} - \underline{d}\|_2^2$$

Notes - guaranteed convergence for

$$0 < \tau < 2/\|\underline{A}\|_{\text{op}}^2$$

$$\underline{w}^{(0)} = \underline{0}, \quad \underline{w}^{(k+1)} = \underline{w}^{(k)} - \tau \underline{A}^T (\underline{A}\underline{w}^{(k)} - \underline{d}) \xrightarrow{k} (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{d}$$

Gradient descent is effective for convex cost functions



$$f(\alpha w_1 + (1-\alpha)w_2) \leq \alpha f(w_1) + (1-\alpha) f(w_2); \quad \frac{d^2}{dw^2} f(w) \geq 0$$

$0 < \alpha < 1, \text{ all } w_1, w_2$

Multidimensional case

$$\underline{H}(\underline{w}) \geq 0$$

$$[\underline{H}(\underline{w})]_{ij} = \frac{\partial^2}{\partial w_i \partial w_j} f(\underline{w})$$

**Copyright 2019
Barry Van Veen**

Proximal Gradient Descent Algorithms

Objectives

- derive proximal gradient algorithm for regularized least-squares problems
 - least-squares gradient descent
 - regularize
- apply to ridge regression

Proximal gradient descent solves regularized least-squares problems

$$\min_{\underline{w}} \|\underline{A}\underline{w} - \underline{d}\|_2^2 + \lambda r(\underline{w}) \quad r(\underline{w}): \text{regularizer}$$

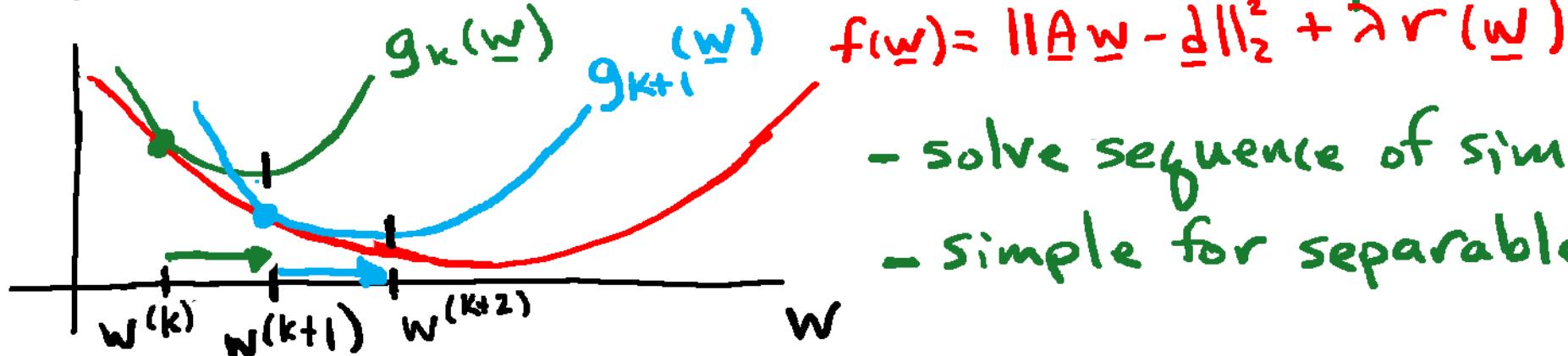
$\lambda > 0$: tuning parameter

Example Convex Regularizers

- Ridge (Tikhonov) $r(\underline{w}) = \|\underline{w}\|_2^2 = \sum_{i=1}^m w_i^2$

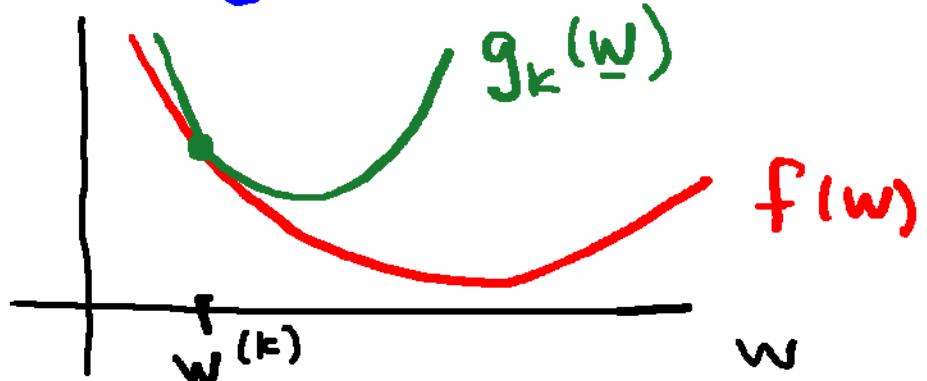
- LASSO (ℓ_1) $r(\underline{w}) = \|\underline{w}\|_1 = \sum_{i=1}^m |w_i|$ not differ.

Proximal Gradient Descent Concept



- solve sequence of simpler problems
- simple for separable $r(\underline{w}) = \sum_i h_i(w_i)$

Find $g_k(\underline{w})$ so $f(\underline{w}) \leq g_k(\underline{w})$, $g_k(\underline{w}^{(k)}) = f(\underline{w}^{(k)})$ 3



minimize $g_k(\underline{w}) \Rightarrow f(\underline{w})$ decreases

$$f(\underline{w}) = \|\underline{d} - \underline{A}\underline{w}\|_2^2 + \lambda r(\underline{w})$$

$$= \|\underline{d} - \underline{A}\underline{w}^{(k)} + (\underline{A}\underline{w}^{(k)} - \underline{A}\underline{w})\|_2^2 + \lambda r(\underline{w})$$

$$\begin{aligned} f(\underline{w}) &= \underbrace{\|\underline{d} - \underline{A}\underline{w}^{(k)}\|_2^2}_{C_k} + \underbrace{\|\underline{A}(\underline{w}^{(k)} - \underline{w})\|_2^2}_{\leq \|\underline{A}\|_{op}^2 \|\underline{w}^{(k)} - \underline{w}\|_2^2} + 2(\underline{d} - \underline{A}\underline{w}^{(k)})^\top \underline{A}(\underline{w}^{(k)} - \underline{w}) + \lambda r(\underline{w}) \\ &\leq \|\underline{A}\|_{op}^2 \|\underline{w}^{(k)} - \underline{w}\|_2^2 + 2\underbrace{(\underline{d} - \underline{A}\underline{w}^{(k)})^\top \underline{A}}_{\underline{V}_k^\top} (\underline{w}^{(k)} - \underline{w}) + \lambda r(\underline{w}) \end{aligned}$$

Define step size $0 < \tau < 1/\|\underline{A}\|_{op}^2 \Rightarrow \frac{1}{\tau} > \|\underline{A}\|_{op}^2$

$$f(\underline{w}) \leq g_k(\underline{w}) = C_k + \frac{1}{\tau} \|\underline{w}^{(k)} - \underline{w}\|_2^2 + 2\underline{V}_k^\top (\underline{w}^{(k)} - \underline{w}) + \lambda r(\underline{w})$$

$g_k(\underline{w})$ is separable
for $r(\underline{w})$ separable: $g_k(\underline{w}) = C_k + \sum_{i=1}^n q_i(w_i)$ no $w_i w_j$ terms

Find $\underline{w}^{(k+1)} = \arg \min_{\underline{w}} g_k(\underline{w})$

$$g_k(\underline{w}) = C_k + \frac{1}{2} \|\underline{w}^{(k)} - \underline{w}\|_2^2 + 2 \underline{v}_k^T (\underline{w}^{(k)} - \underline{w}) + \lambda r(\underline{w})$$

$$\tau g_k(\underline{w}) = \tau C_k + (\underline{w}^{(k)} - \underline{w})^T (\underline{w}^{(k)} - \underline{w}) + 2 \tau \underline{v}_k^T (\underline{w}^{(k)} - \underline{w}) + \lambda \tau r(\underline{w})$$

$$= \tau C_k - \tau^2 \underline{v}_k^T \underline{v}_k + \underbrace{(\tau \underline{v}_k + (\underline{w}^{(k)} - \underline{w}))^T}_{z^{(k)}} \underbrace{(\tau \underline{v}_k + (\underline{w}^{(k)} - \underline{w}))}_{\underline{z}^{(k)}} + \lambda \tau r(\underline{w})$$

$$\underline{w}^{(k+1)} = \arg \min_{\underline{w}} \|\underline{z}^{(k)} - \underline{w}\|_2^2 + \lambda \tau r(\underline{w})$$

$$\underline{z}^{(k)} = \underline{w}^{(k)} + \tau \underline{v}_k$$

$$= \underline{w}^{(k)} + \tau \underline{A}^T (\underline{d} - \underline{A} \underline{w}^{(k)})$$

$$= \underline{w}^{(k)} - \tau \underline{A}^T (\underline{A} \underline{w}^{(k)} - \underline{d})$$

Least-squares
gradient descent
(Landweber)

Alternate LS gradient descent and regularization 5

$$\underline{w}^{(0)} = \underline{0}, \quad 0 < \tau < \frac{1}{\|\underline{A}\|_{op}^2}$$

initialize

$$\underline{z}^{(k)} = \underline{w}^{(k)} - \tau \underline{A}^T (\underline{A} \underline{w}^{(k)} - \underline{d})$$

LS gradient descent

$$\underline{w}^{(k+1)} = \arg \min_{\underline{w}} \|\underline{z}^{(k)} - \underline{w}\|_2^2 + \lambda \tau r(\underline{w})$$

regularize

$$\text{if } \|\underline{w}^{(k+1)} - \underline{w}^{(k)}\| < \varepsilon \text{ stop}$$

check if converged

Regularization simple for $r(\underline{w})$ separable!

$$\text{if } r(\underline{w}) = \sum_{i=1}^m h_i(w_i)$$

$$\underline{w}^{(k+1)} = \arg \min_{w_i, i=1, \dots, M} \sum_{i=1}^m \left((\underline{z}_i^{(k)} - w_i)^2 + \lambda \tau h_i(w_i) \right)$$

M scalar minimizations

Example: Ridge Regression (Tikhonov) 6

$$f(\underline{w}) = \|\underline{d} - \underline{A}\underline{w}\|_2^2 + \lambda \|\underline{w}\|_2^2$$

LS gradient descent:

$$\underline{z}^{(k)} = \underline{w}^{(k)} - \tau \underline{A}^T (\underline{A}\underline{w}^{(k)} - \underline{d})$$

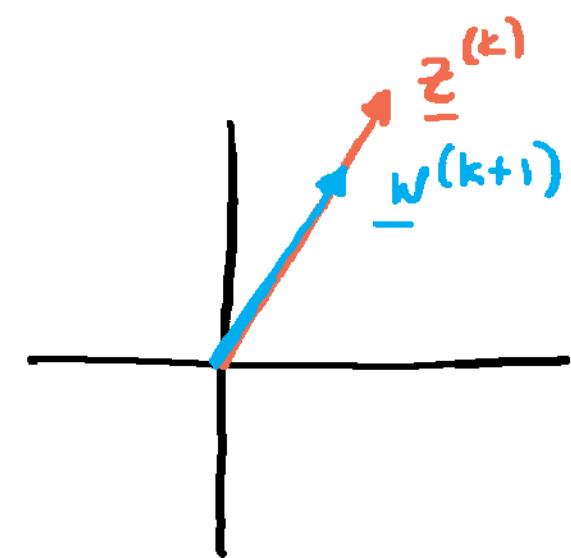
Regularization:

$$\underline{w}^{(k+1)} = \underset{w_i, i=1, \dots, M}{\arg \min} \sum_{i=1}^M (z_i^{(k)} - w_i)^2 + \lambda \tau w_i^2$$

$$\Rightarrow w_i^{(k+1)} = \frac{1}{1 + \lambda \tau} z_i^{(k)}$$

$$\underline{w}^{(k+1)} = \frac{1}{1 + \lambda \tau} \underline{z}^{(k)}$$

"Shrink toward origin"



**Copyright 2019
Barry Van Veen**

Sparse Solutions to Least-Squares Problems Using the LASSO

Objectives

- motivate search for sparse solutions
- introduce ℓ_1 -norm regularization (LASSO)
- overview attributes of ℓ_1 -regularization

Sparse classifiers/models give insight 2

$$(\underline{x}_i, d_i), i=1, \dots, N$$

features, labels

$$\underline{x}_i^\top \underline{w} \approx d_i$$

$$\underline{A}\underline{w} = \begin{bmatrix} \underline{a}_1 & \underline{a}_2 & \cdots & \underline{a}_m \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_m \end{bmatrix} = \sum_{i=1}^m w_i \underline{a}_i$$

\underline{a}_l : lth feature component

Suppose $w_l \approx 0 \Rightarrow \underline{a}_l$ is unimportant

If a small number of w_i are nonzero, only those few features matter! \underline{w} is sparse

$$\|\underline{w}\|_0 = \sum_{i=1}^m \mathbf{1}_{\{w_i \neq 0\}}$$

(number of nonzero elements)

ℓ_0 "norm"

$$\|\underline{a}\underline{w}\|_0 \neq \alpha \|\underline{w}\|_0$$

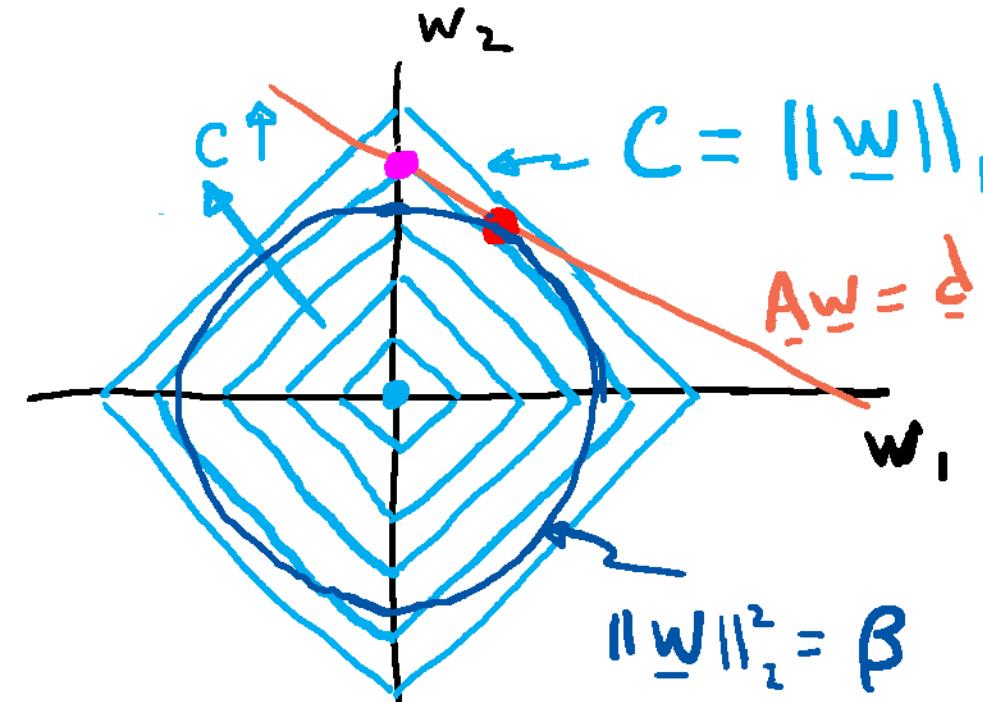
$$\text{Consider } \min_{\underline{w}} \|\underline{w}\|_0 \text{ s.t. } \|\underline{A}\underline{w} - \underline{d}\|_2^2 \leq \epsilon$$

non convex - intractable

Convex relaxation gives tractable problem

$$\min_{\underline{w}} \|\underline{w}\|_1 \text{ s.t. } \|\underline{A}\underline{w} - \underline{d}\|_2^2 \leq \Sigma \quad \text{LASSO: Least Absolute Selection + Shrinkage Operator}$$

- convex



$$c = \|\underline{w}\|_1 = \sum_{i=1}^m |w_i| : |w_1| + |w_2| = c$$

1st quad $w_1 + w_2 = c$

$$\min \|\underline{w}\|_1 \text{ s.t. } \underline{A}\underline{w} = \underline{d}$$

"Corners" on $\|\underline{w}\|_1 \Rightarrow$ sparse sol'n's

$$\min \|\underline{w}\|_2^2 \text{ s.t. } \underline{A}\underline{w} = \underline{d}$$

circular $\|\underline{w}\|_2 \Rightarrow$ non sparse solutions

LASSO is a regularized least-squares problem 4

$\min_{\underline{w}} \|\underline{w}\|_1$, s.t. $\|\underline{A}\underline{w} - \underline{d}\|_2^2 \leq \varepsilon$ is equivalent to

$\min_{\underline{w}} \|\underline{A}\underline{w} - \underline{d}\|_2^2 + \lambda \|\underline{w}\|_1$, for some λ, ε

Note: $\min_{\underline{w}} \|\underline{w}\|_1 + \frac{1}{\lambda} \|\underline{A}\underline{w} - \underline{d}\|_2^2$

LASSO

$$\underline{w}_L = \arg \min_{\underline{w}} \|\underline{A}\underline{w} - \underline{d}\|_2^2 + \lambda \|\underline{w}\|_1$$

Sparse \underline{w}_L

can have small model error

$$\underline{w}_{opt} - \underline{w}_L$$

iterative solution

Ridge Regression

$$\underline{w}_R = \arg \min_{\underline{w}} \|\underline{A}\underline{w} - \underline{d}\|_2^2 + \lambda \|\underline{w}\|_2^2$$

non sparse \underline{w}_R

great prediction error

$$\|\underline{A}\underline{w}_{opt} - \underline{A}\underline{w}_R\|_2^2$$

can solve in closed form

LASSO may be used for model/feature selection 5

$$\underline{w}_L = \arg \min_{\underline{w}} \|\underline{A}\underline{w} - \underline{d}\|_2^2 + \lambda \|\underline{w}\|_1$$

$$S_L = \{i : [\underline{w}_L]_i \neq 0\} \quad \text{selected features}$$

$$\underline{A}\underline{w}_L = \sum_{i=1}^m \underline{a}_i [\underline{w}_L]_i = \sum_{i \in S_L} \underline{a}_i [\underline{w}_L]_i$$

Debiasing $\underline{A}_L = \{\underline{a}_i : i \in S_L\}$

$$\hat{\underline{w}}_L = \arg \min_{\underline{w}} \|\underline{A}_L \underline{w} - \underline{d}\|_2^2 = (\underline{A}_L^\top \underline{A}_L)^{-1} \underline{A}_L^\top \underline{d}$$

avoids shrinkage due to $\|\underline{w}\|_1$

Copyright 2019
Barry Van Veen

Solving ℓ_1 Regularized Least Squares via Proximal Gradient Descent

Objectives

1

- apply proximal gradient approach to solve ℓ_1 -regularized least squares
- derive solution to regularization phase
- explore alternating gradient and soft thresholding steps

The ℓ_1 -regularized least-squares problem 2
can be solved via proximal gradient descent

features/labels (\underline{x}_i, d_i) model $\underline{x}_i^\top \underline{w} \approx d_i$

$$\min_{\underline{w}} \|\underline{A}\underline{w} - \underline{d}\|_2^2 + \lambda \|\underline{w}\|_1 \quad \text{encourages sparse solutions}$$

no closed form solution

Proximal Gradient Descent Algorithm

a) $\underline{z}^{(k)} = \underline{w}^{(k)} - \tau \underline{A}^\top (\underline{A}\underline{w}^{(k)} - \underline{d})$ least squares gradient descent

b) $\underline{w}^{(k+1)} = \arg \min_{\underline{w}} \|\underline{z}^{(k)} - \underline{w}\|_2^2 + \tau \lambda \|\underline{w}\|_1$

Regularization step involves scalar minimization³

$$\min_{\underline{w}} \|\underline{z}^{(k)} - \underline{w}\|_2^2 + \varepsilon \lambda \|\underline{w}\|_1 \Rightarrow \min_{\substack{\underline{w}_i, i=1, \dots, n \\ w_i}} \sum_{i=1}^n (z_i^{(k)} - w_i)^2 + \lambda \varepsilon |w_i|$$

Consider $\min_{w_i} (z_i^{(k)} - w_i)^2 + \lambda \varepsilon |w_i|, \lambda, \varepsilon > 0$

case 1: $w_i \geq 0$

$$\min_{w_i} (z_i - w_i)^2 + \lambda \varepsilon w_i, w_i \geq 0 \quad \text{if } z_i > \frac{\lambda \varepsilon}{2},$$

$$\frac{d}{dw_i} \{(z_i - w_i)^2 + \lambda \varepsilon w_i\} = 0, w_i \geq 0 \quad w_i = z_i - \frac{\lambda \varepsilon}{2}$$

$$-2(z_i - w_i) + \lambda \varepsilon = 0, w_i \geq 0 \quad \text{if } z_i < \frac{\lambda \varepsilon}{2},$$

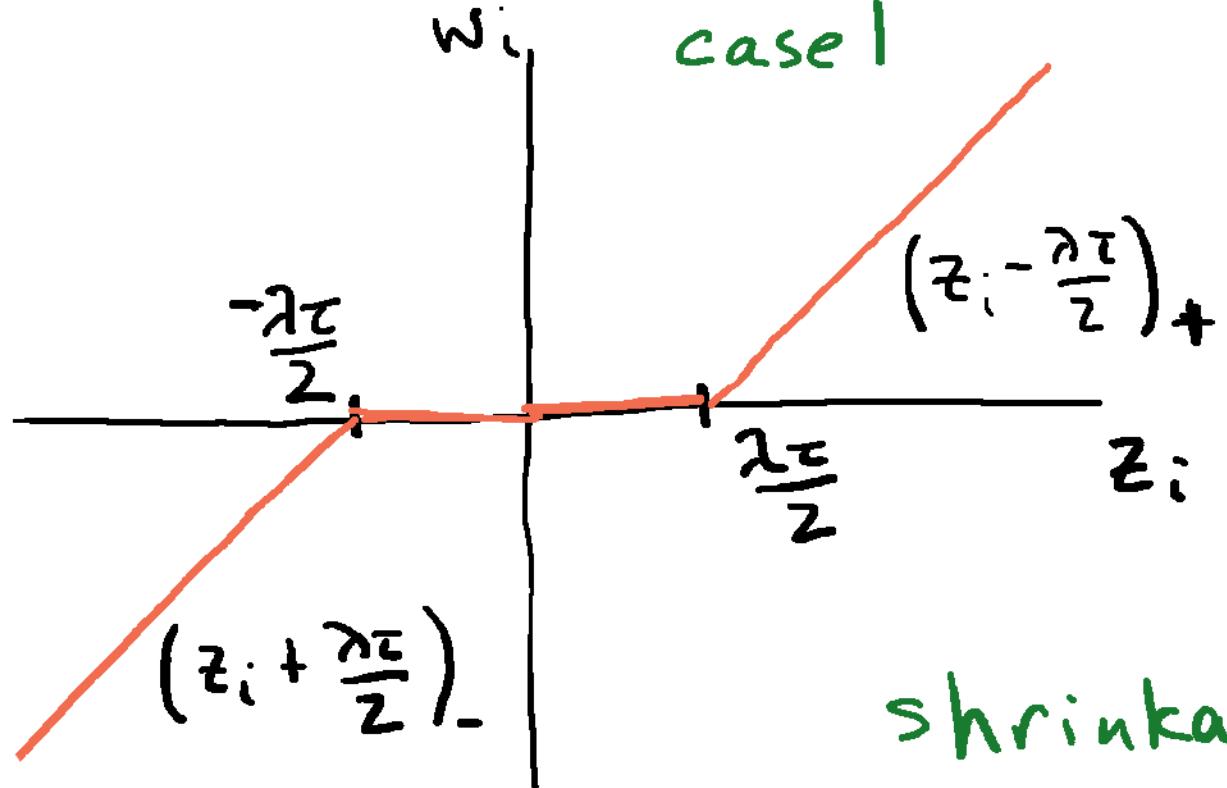
$$w_i = z_i - \frac{\lambda \varepsilon}{2}, w_i \geq 0 \quad w_i = 0$$

$$w_i = (z_i - \frac{\lambda \varepsilon}{2})_+$$

Case 2: $w_i \leq 0$

$$\min_{w_i} (z_i - w_i)^2 - \lambda \tau w_i \Rightarrow \frac{d}{dw_i} \left\{ (z_i - w_i)^2 - \lambda \tau w_i \right\} = 0$$

$$-2(z_i - w_i) - \lambda \tau = 0, w_i \leq 0 \quad w_i = \left(z_i + \frac{\lambda \tau}{2} \right)_-$$



case 2

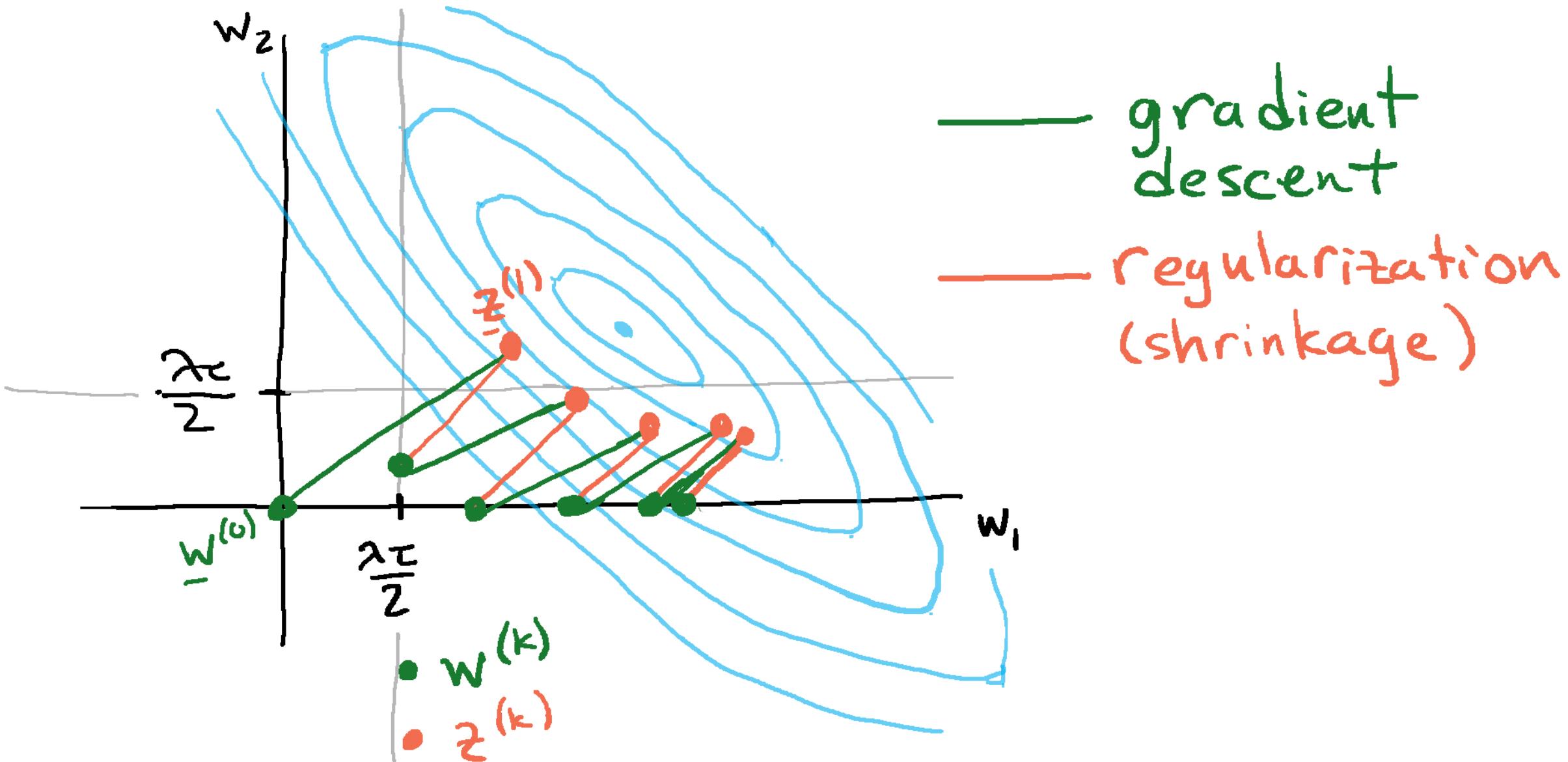
shrinkage

"Soft threshold"

$$w_i = \begin{cases} 0, & -\frac{\lambda \tau}{2} < z_i < \frac{\lambda \tau}{2} \\ z_i - \frac{\lambda \tau}{2}, & z_i > \frac{\lambda \tau}{2} \\ z_i + \frac{\lambda \tau}{2}, & z_i < -\frac{\lambda \tau}{2} \end{cases}$$

$$w_i = \left(|z_i| - \frac{\lambda \tau}{2} \right)_+ \text{sign}(z_i)$$

Algorithm alternates descent and shrinkage 5



Copyright 2019
Barry Van Veen

Hinge Loss for Binary Classifiers

Objectives

- introduce disadvantage of squared error for classification
- introduce hinge loss cost function
- characteristics of hinge loss

Squared error "loss" can be problematic 2

Classifier design

$$\min_{\underline{w}} l(\underline{w}; \underline{A}, \underline{d}) + \lambda r(\underline{w}) \leftarrow \begin{array}{l} \text{regularizer} \\ \text{loss function} \end{array}$$

Squared error loss $ll(\underline{w}; \underline{A}, \underline{d}) = \|\underline{A}\underline{w} - \underline{d}\|_2^2$

Example: dwarf planet vs. planet

| object | Ceres | Eris | Pluto | Mercury | Earth | Jupiter |
|---------------------------------|-------|------|-------|---------|-------|---------|
| x_i radius ($\times 10^6$ m) | 1.0 | 2.3 | 2.4 | 4.9 | 12.8 | 143.0 |
| d_i label | -1 | -1 | -1 | 1 | 1 | 1 |

$$\underline{A} = \begin{bmatrix} 1 & 1 \\ 2.3 & 1 \\ 2.4 & 1 \\ 4.9 & 1 \\ 12.8 & 1 \\ 143 & 1 \end{bmatrix}, \underline{d} = \begin{bmatrix} -1 \\ -1 \\ -1 \\ -1 \\ 1 \\ 1 \end{bmatrix}$$

$$\underline{w}_{LS} = (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{d}$$

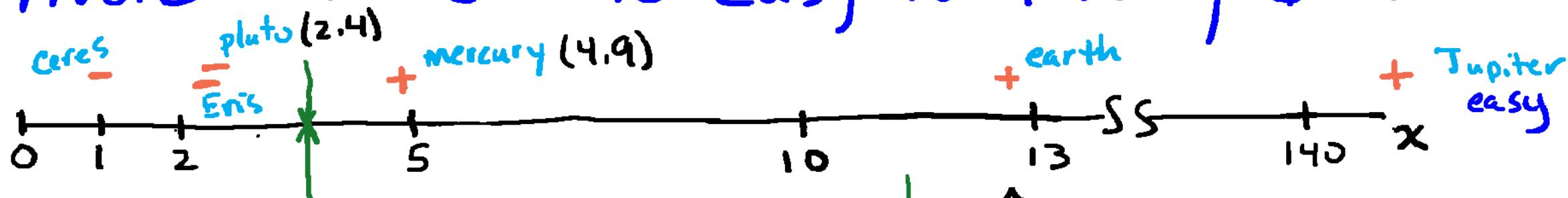
$$\approx 0.01 \begin{bmatrix} 1 \\ -28 \end{bmatrix}$$

dwarf: $x_i < 28$ (earth!)
planet: $x_i \geq 28$

squared error \rightarrow poor classification

Avoid loss due to "easy-to-classify" data

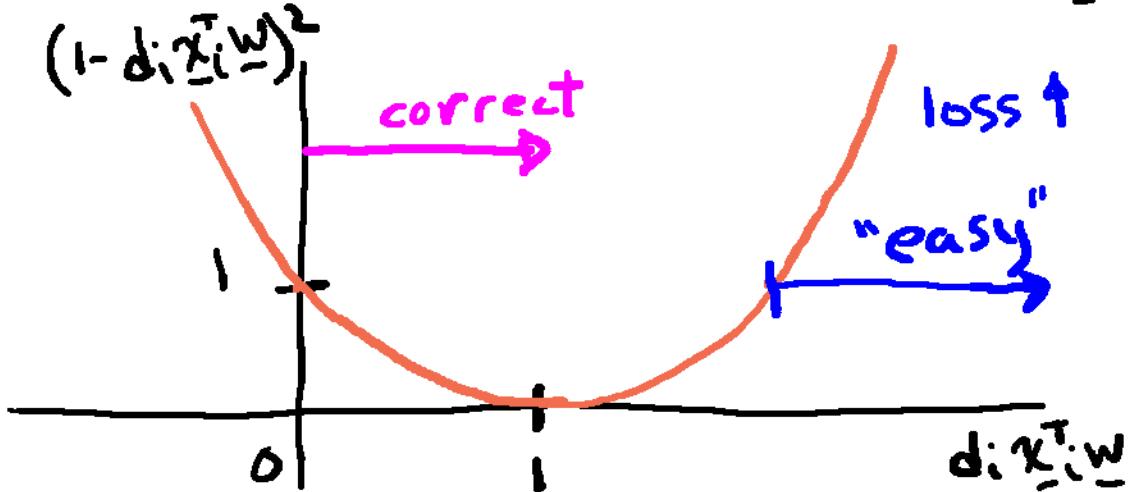
3



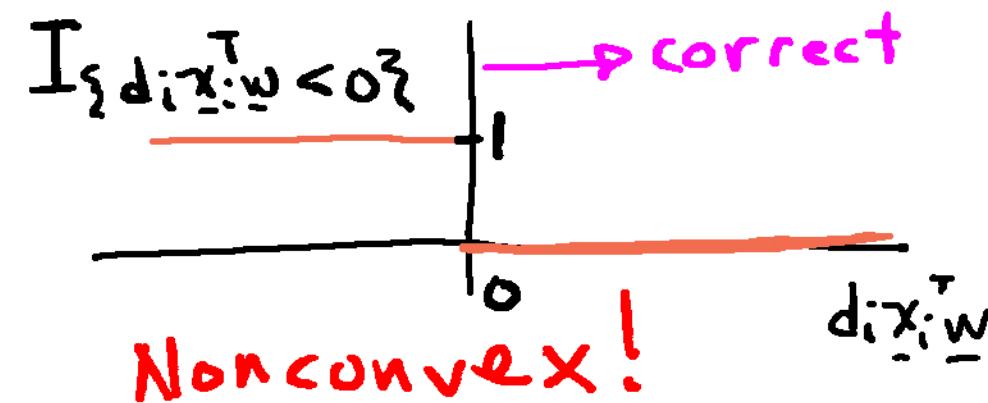
max margin classifier: midpoint $\hat{d} = \text{sign}(x - 3.65)$
 Margin $4.9 - 2.4 = 2.5$ (class separation)

Squared error loss: $\|\underline{A}\underline{w} - \underline{d}\|_2^2 = \sum_{i=1}^N (d_i - \underline{x}_i^\top \underline{w})^2 = \sum_i (1 - d_i \underline{x}_i^\top \underline{w})^2$
 (d_i = ±1)

correct classification: $d_i \underline{x}_i^\top \underline{w} > 0$

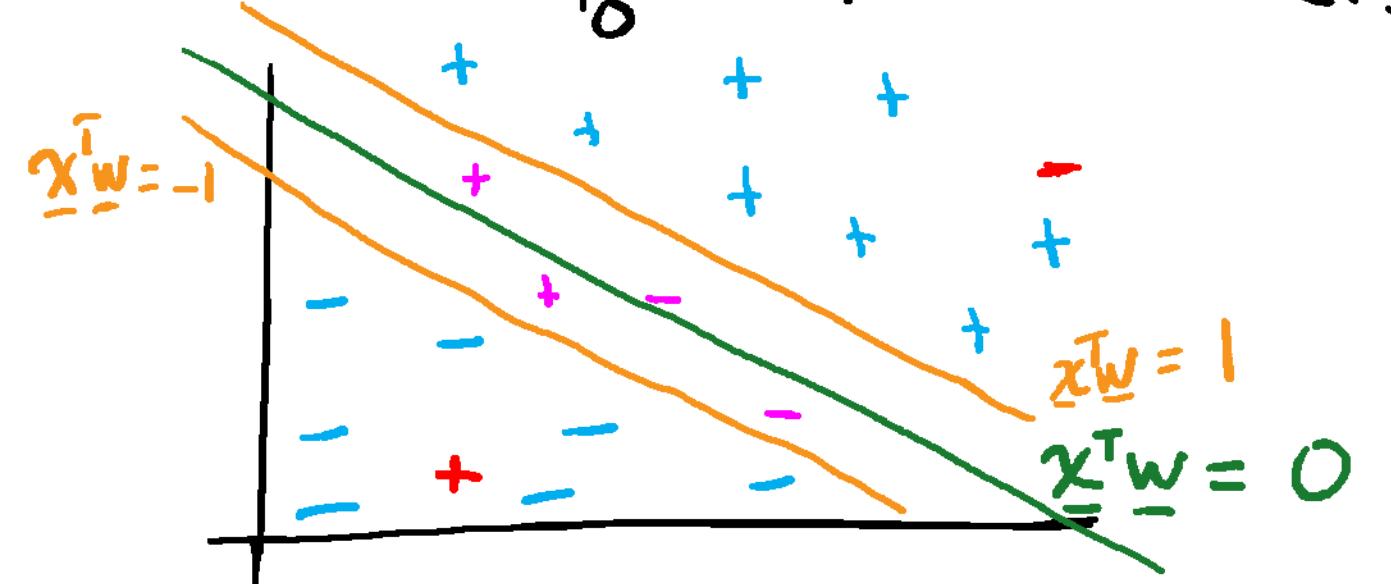
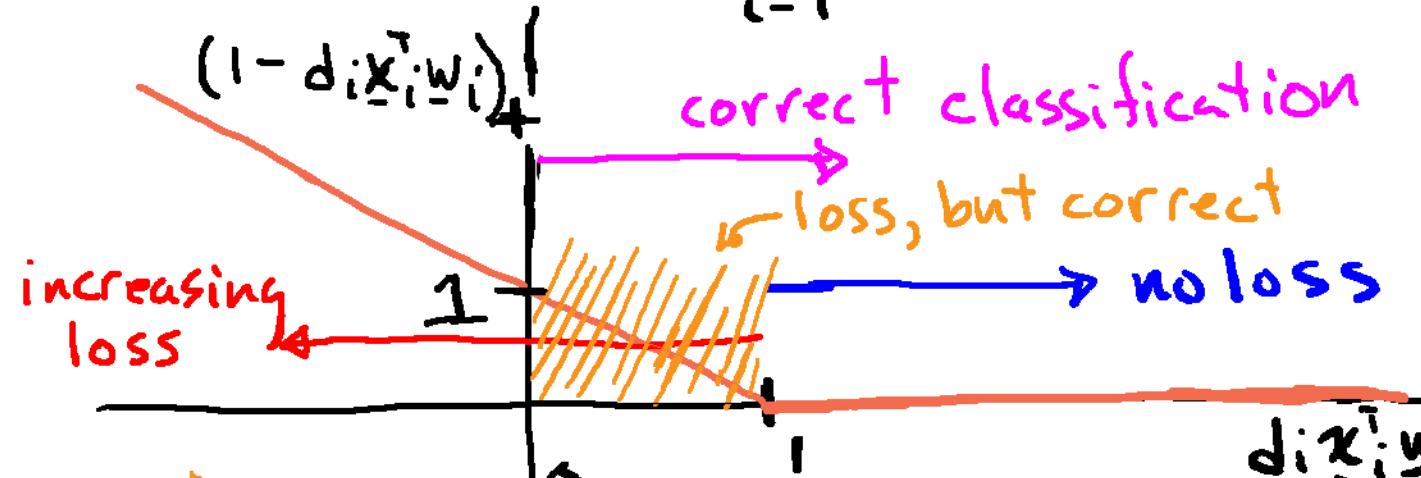


Ideal loss



Hinge loss is convex and has no loss for easy
to classify data 4

$$\ell(\underline{w}; \underline{A}, \underline{d}) = \sum_{i=1}^n (1 - d_i \underline{x}_i^\top \underline{w})_+$$

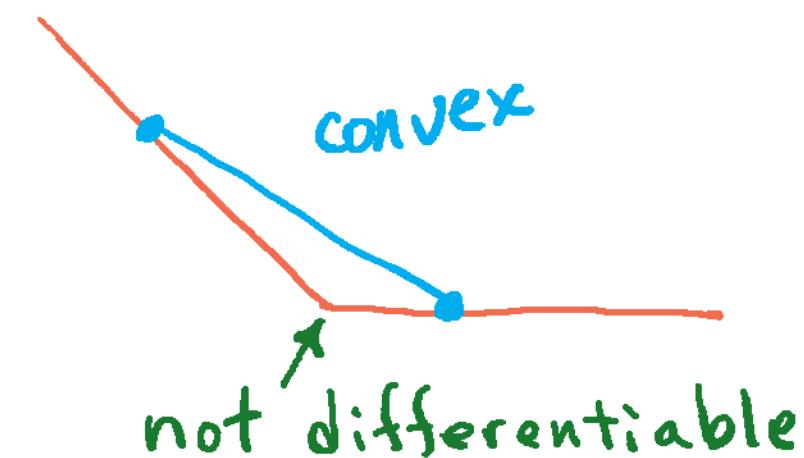


+, - no hinge loss

+, - small hinge loss

+, - large hinge loss

$$(\alpha)_+ = \begin{cases} \alpha & \alpha > 0 \\ 0 & \alpha \leq 0 \end{cases}$$



Hinge loss better approximates ideal:
number of misclassifications

Iterative algorithms required for
finding minimum hinge loss classifier

Copyright 2019
Barry Van Veen

Support Vector Machines for Classification

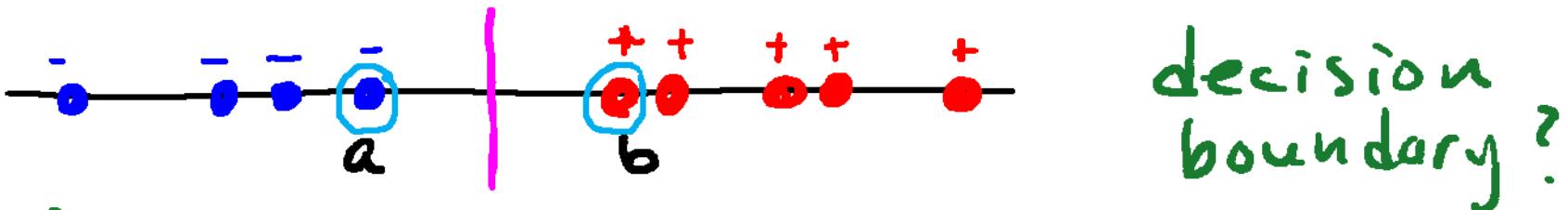
Objectives

1

- Define margin for separable data
- Show Support vector machines
Maximize margin
- Use hinge loss to define support vector machines for non separable data

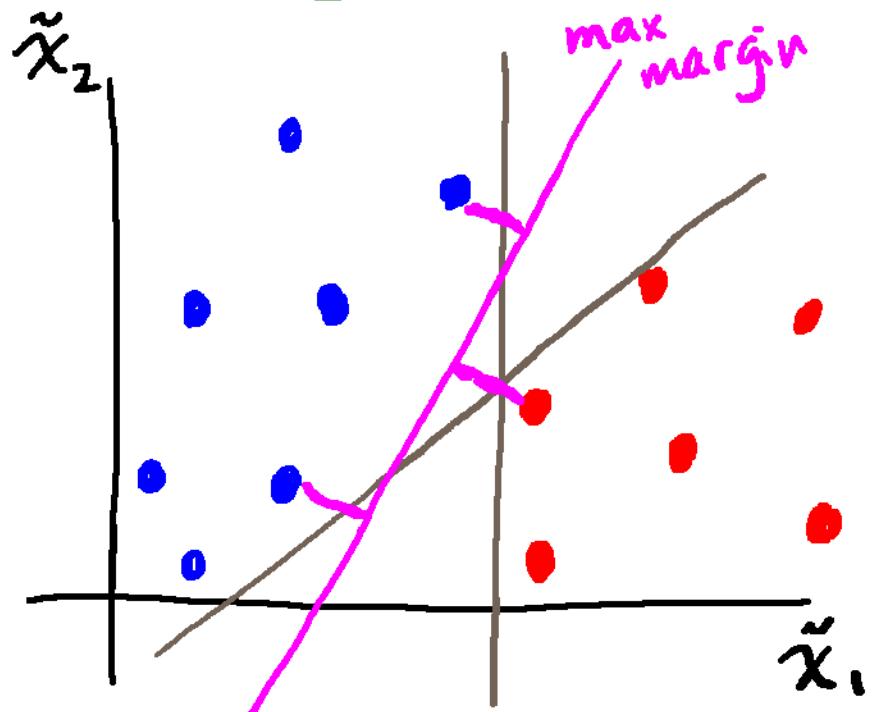
Maximize margin for separable training data 2

Example:



margin: distance from boundary to nearest sample

max margin
boundary: midpoint
only a, b matter



decision boundary?

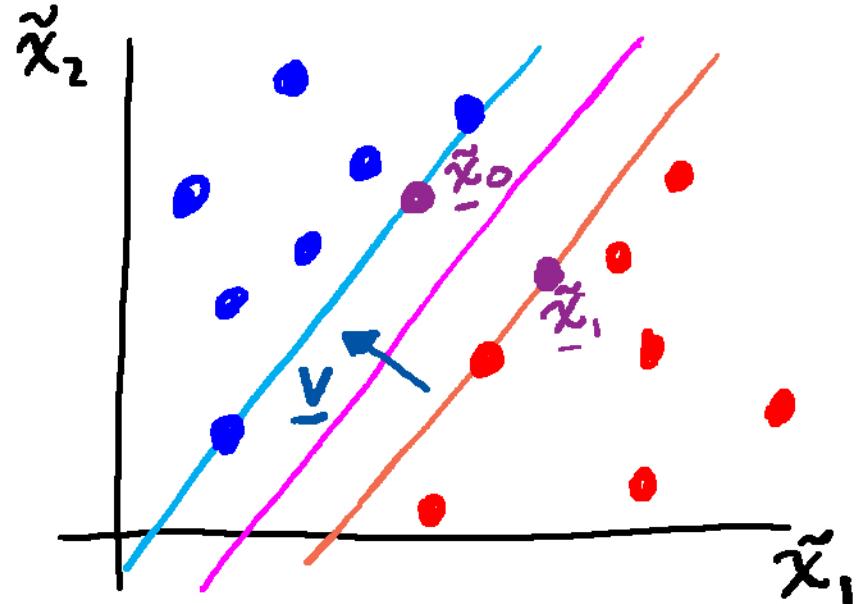
$$\text{feature } \underline{x}^T = [\tilde{x}_1^T \ 1]$$

$$\text{weights } \underline{w}^T = [\tilde{w}^T \ w_0]$$

decision $\hat{d} = \text{sign}(\underline{x}^T \underline{w})$

$$\hat{d} = \begin{cases} 1 & \tilde{x}^T \tilde{w} + w_0 > 0 \\ -1 & \tilde{x}^T \tilde{w} + w_0 < 0 \end{cases}$$

Margin is determined by $\|\tilde{w}\|_2^{-1}$



$$\text{label } "-1": \underline{\tilde{x}}^T \underline{\tilde{w}} + w_0 \leq -1$$

$$\text{label } "+1": \underline{\tilde{x}}^T \underline{\tilde{w}} + w_0 \geq 1$$

$$\text{boundary: } \underline{\tilde{x}}^T \underline{\tilde{w}} + w_0 = 0$$

margin: γ_2 distance between //
measure in direction \underline{v}

Unit normal to boundary plane: $\underline{v} = \frac{\underline{\tilde{w}}}{\|\underline{\tilde{w}}\|_2}$

$$\text{Margin } m = \frac{1}{2} \|\underline{\tilde{x}}_1 - \underline{\tilde{x}}_0\|_2 \quad \underline{\tilde{x}}_1 = \underline{\tilde{x}}_0 + 2m \underline{v}$$

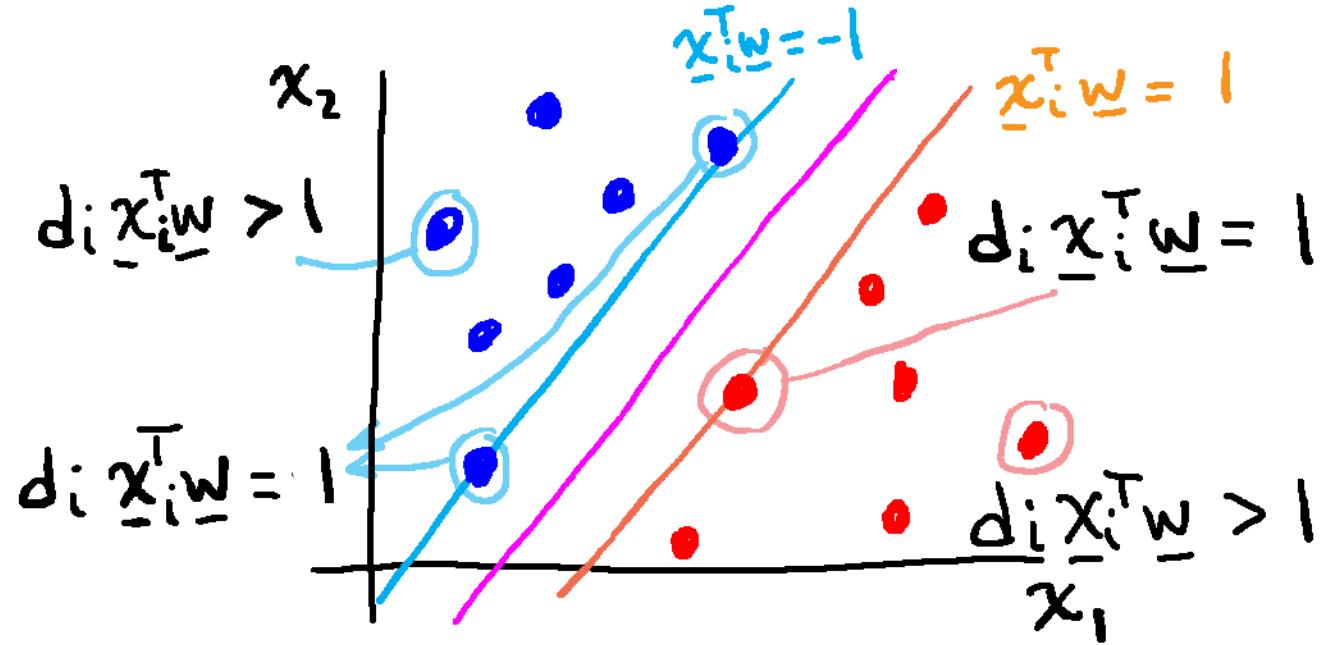
but $\underline{\tilde{x}}_0^T \underline{\tilde{w}} + w_0 = -1$

$$1 = \underline{\tilde{x}}_1^T \underline{\tilde{w}} + w_0 = \underline{\tilde{x}}_0^T \underline{\tilde{w}} + 2m \frac{\underline{\tilde{w}}^T \underline{\tilde{w}}}{\|\underline{\tilde{w}}\|_2} + w_0$$

$$m = \frac{2}{\|\underline{\tilde{w}}\|_2}$$

Support Vector Machine maximizes margin 4

Correct classification: $d_i(\tilde{x}_i^T \tilde{w} + w_0) \geq 1$



SVM:

$$\min_{\underline{\tilde{w}}, w_0} \|\underline{\tilde{w}}\|_2^2 \text{ s.t. } d_i(\tilde{x}_i^T \tilde{w} + w_0) \geq 1 \quad i=1, 2, \dots, N$$

max margin

perfect
classification
unique solution

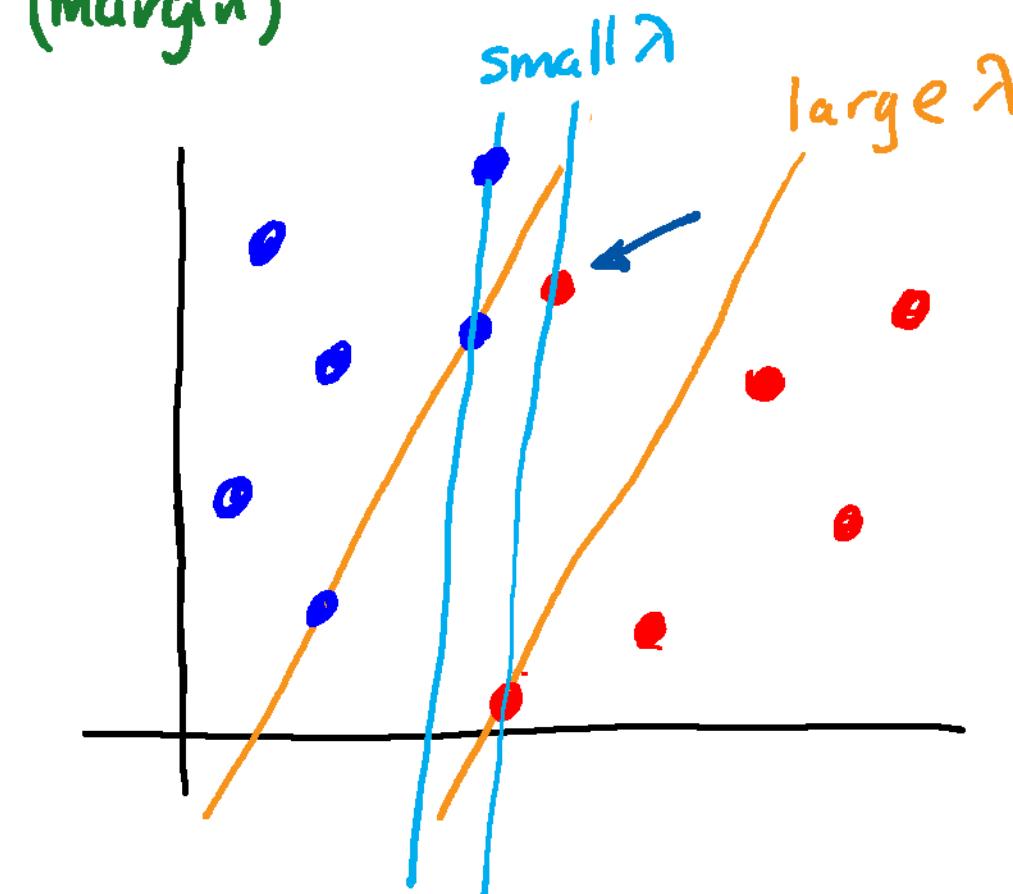
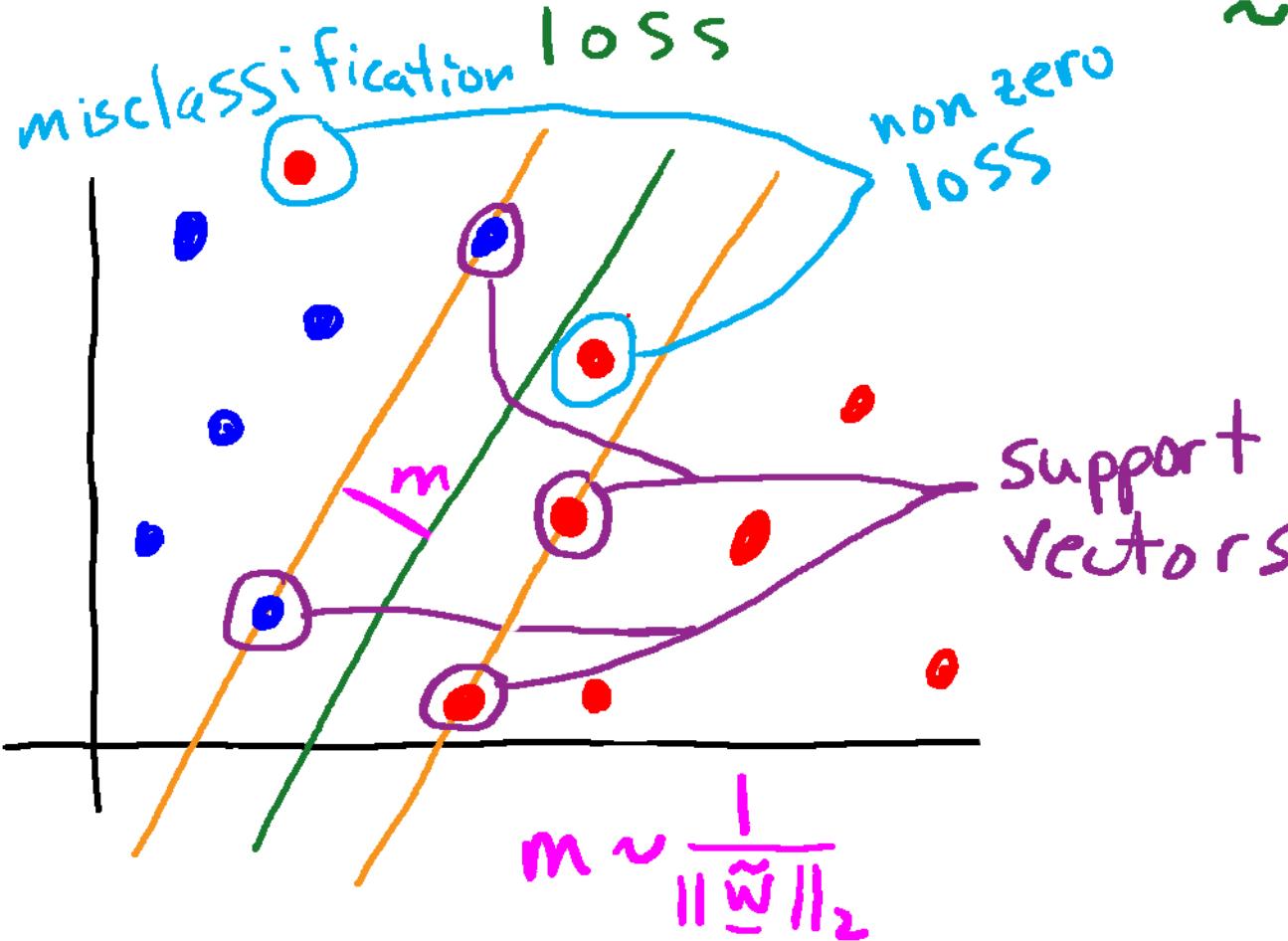
Boundary defined by \underline{x}_i for which $d_i \underline{x}_i^T \underline{w} = 1$

Called Support Vectors

SVM for non separable data uses hinge loss⁶

$$\min_{\underline{w}} \sum_{i=1}^N (1 - d_i \underline{x}_i^T \underline{w})_+ + \lambda \|\underline{w}\|_2^2 \quad (\text{l}_2 \text{ regularization})$$

$\sim (\text{margin})^{-2}$



**Copyright 2019
Barry Van Veen**

Gradient Descent for Support Vector Machines and Subgradients

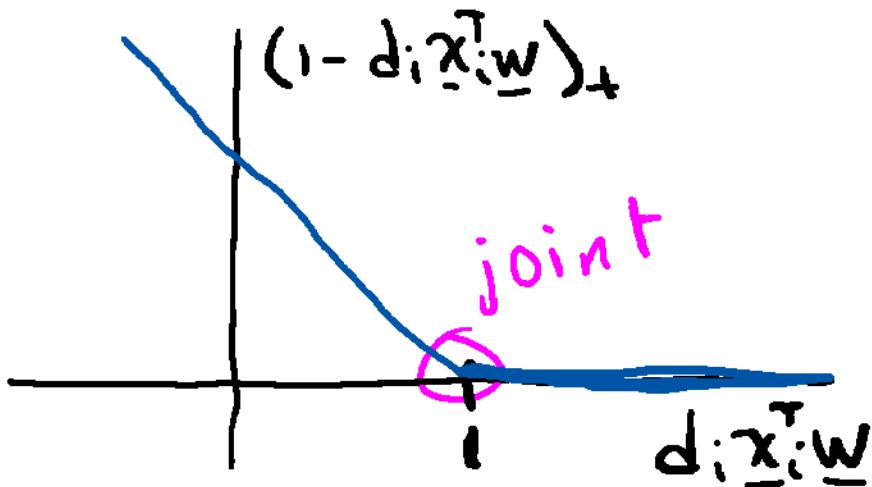
Objectives

- develop a gradient descent algorithm for SVMs
- introduce subgradients for convex but non differentiable cost functions

Support vector machines require iterative algorithms ²

$$\min_{\underline{w}} \sum_{i=1}^N (1 - d_i \underline{x}_i^\top \underline{w})_+ + \lambda \|\underline{w}\|_2^2$$

labels features hinge loss regularization



No closed form solution
Convex

⇒ gradient descent

Problem: hinge loss not differentiable

Subderivatives generalize derivatives

3

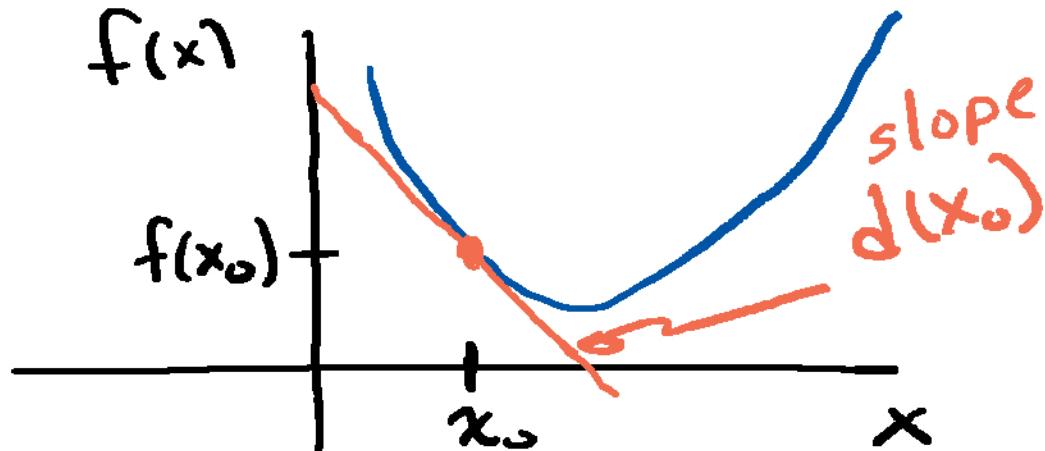
- Convex, but non-differentiable $f(x)$

Derivatives -

$$d(x_0) = \lim_{x \rightarrow x_0} \frac{f(x) - f(x_0)}{x - x_0}$$

Convex:

$$f(x) \geq f(x_0) + d(x_0)(x - x_0)$$



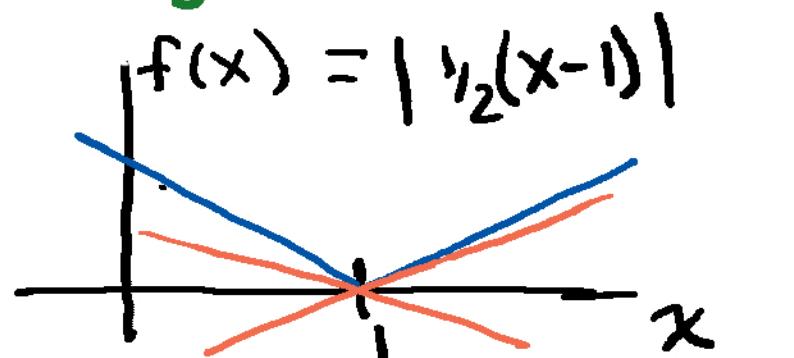
"above tangent line"

Subderivative (convex)

Any $d_s(x_0)$: $f(x) \geq f(x_0) + d_s(x_0)(x - x_0)$

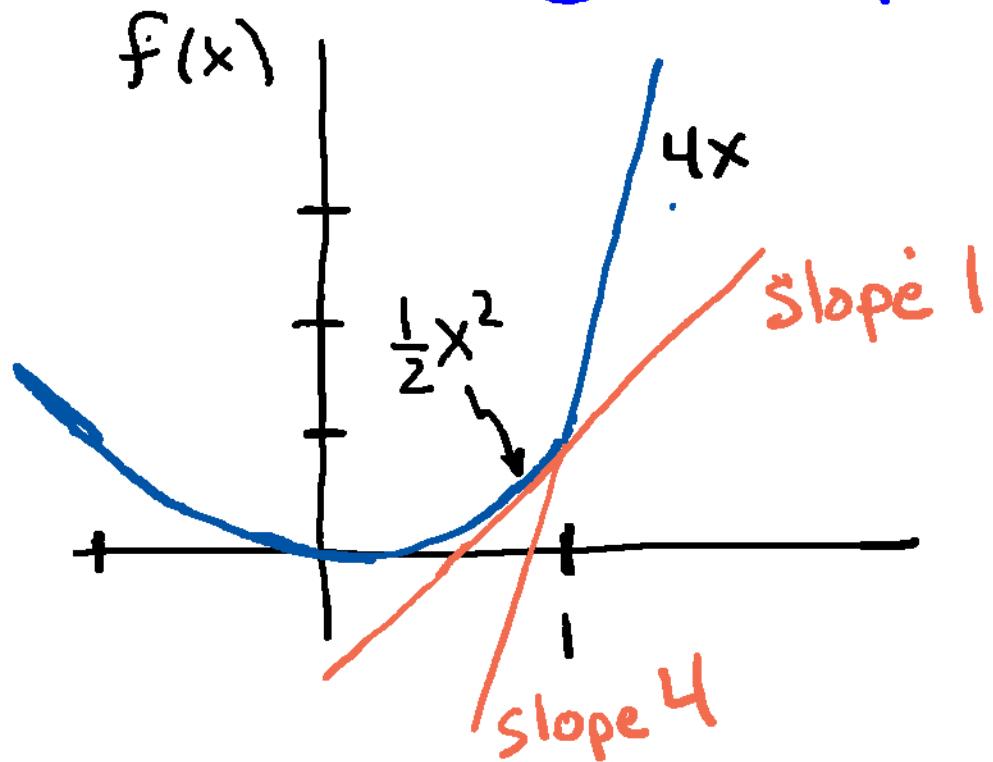
$$x < 1 : d_s(x) = -\frac{1}{2}; x > 1 : d_s(x) = \frac{1}{2}$$

$$-\frac{1}{2} \leq d_s(1) \leq \frac{1}{2}$$



"below f(x)"

Sub derivatives produce "reasonable" downhill directions 4



Example: $f(x) = \begin{cases} \frac{1}{2}x^2 & x < 1 \\ 4x & x \geq 1 \end{cases}$
convex

Subderivative

$$d_s(x) = \begin{cases} x > & x < 1 \\ 4 & x \geq 1 \\ [1, 4] & x = 1 \end{cases}$$

Subgradients generalize gradients
- Convex, nondifferentiable $l(\underline{w})$

Gradients -

$$l(\underline{w}) \geq l(\underline{w}_0) + (\underline{w} - \underline{w}_0)^T \underline{v}(\underline{w}_0) \quad \underline{v}(\underline{w}) = \nabla_{\underline{w}} l(\underline{w})$$

"above tangent plane" ($\sum_{i=1}^m (w_i - w_{0i}) \frac{d}{dw_i} l(\underline{w})$)

Subgradients -

$$\text{Any } \underline{v}(\underline{w}) : l(\underline{w}) \geq l(\underline{w}_0) + (\underline{w} - \underline{w}_0)^T \underline{v}(\underline{w}_0)$$

Gradient descent optimization: replace
gradient with subgradient

Gradient descent for SVMs

$$\ell(\underline{w}) = \sum_{i=1}^N (1 - d_i \underline{x}_i^\top \underline{w})_+ \rightarrow \text{subgradient}$$

$$\ell_i(\underline{w}) = (1 - d_i \underline{x}_i^\top \underline{w})_+ = \begin{cases} 1 - d_i \underline{x}_i^\top \underline{w} & d_i \underline{x}_i^\top \underline{w} < 1 \\ 0 & d_i \underline{x}_i^\top \underline{w} \geq 1 \end{cases}$$

Subgradient

$$\nabla_i(\underline{w}) = \begin{cases} -d_i \underline{x}_i & d_i \underline{x}_i^\top \underline{w} < 1 \\ 0 & d_i \underline{x}_i^\top \underline{w} \geq 1 \end{cases} = -d_i \underline{x}_i I_{\{d_i \underline{x}_i^\top \underline{w} < 1\}}$$

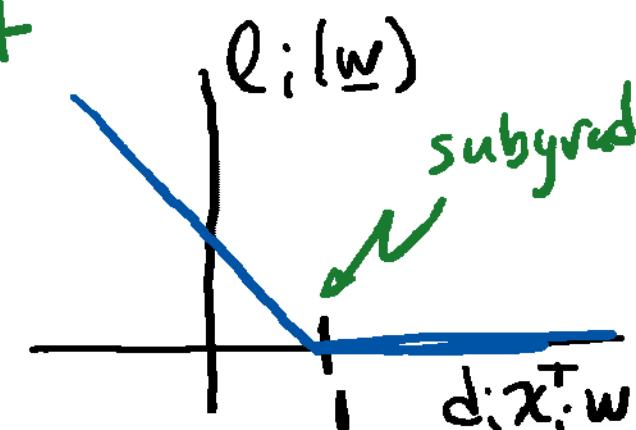
indicator function

Cost $f(\underline{w}) = \ell(\underline{w}) + \lambda \|\underline{w}\|_2^2$

$$\Rightarrow \nabla f(\underline{w})|_{\underline{w}^{(k)}} = \sum_{i=1}^N (-d_i \underline{x}_i I_{\{d_i \underline{x}_i^\top \underline{w}^{(k)} < 1\}}) + 2\lambda \underline{w}^{(k)}$$

Gradient descent

$$\underline{w}^{(k+1)} = \underline{w}^{(k)} - \tau \nabla f(\underline{w})|_{\underline{w}^{(k)}}$$



Copyright 2019
Barry Van Veen

Stochastic Gradient Descent

Objectives

- Simplify gradient descent update
- Common methods for cycling through data
- Benefits
- Examples

Stochastic gradient descent updates weights 2 using part of the data

$$f(\underline{w}) = l(\underline{w}) + \lambda r(\underline{w}) \quad \underline{w}^{(k+1)} = \underline{w}^{(k)} - \frac{\epsilon}{2} \nabla_{\underline{w}} f(\underline{w})$$

"loss" "regularize" gradient

squared error hinge loss $(d_i, x_i), i=1, \dots, N$

$$l(\underline{w}) = \sum_{i=1}^N (d_i - \underline{x}_i^T \underline{w})^2 \quad l(\underline{w}) = \sum_{i=1}^N (1 - d_i \underline{x}_i^T \underline{w})_+$$

$\nabla_{\underline{w}} l(\underline{w}) = -2 \sum_{i=1}^N (d_i - \underline{x}_i^T \underline{w}) \underline{x}_i \quad \nabla_{\underline{w}} l(\underline{w}) = - \sum_{i=1}^N I_{\{d_i \underline{x}_i^T \underline{w} < 1\}} \underline{x}_i$

depends on all the data

SGD: $f(\underline{w}) = \sum_{i=1}^N f_i(\underline{w})$ Define $i_k, k=1, 2, \dots$

$$\underline{w}^{(k+1)} = \underline{w}^{(k)} - \frac{\epsilon}{2} \nabla_{\underline{w}} f_{i_k}(\underline{w}^{(k)})$$

depends on one sample $(d_{i_k}, \underline{x}_{i_k})$

SGD cycles through training data

3

1) Cyclical (incremental gradient descent)

$$i_k = k \bmod N \quad \text{e.g. } i_k = 1, 2, 3, 4, 1, 2, 3, 4, 1, 2, 3 \dots$$

2) Random permutation (reshuffle every N rounds)

$$i_k = 2, 4, 1, 3, \boxed{2, 1, 4, 3}, 4, 3, 1, 2 \dots$$

3) Stochastic gradient descent (uniformly at random)

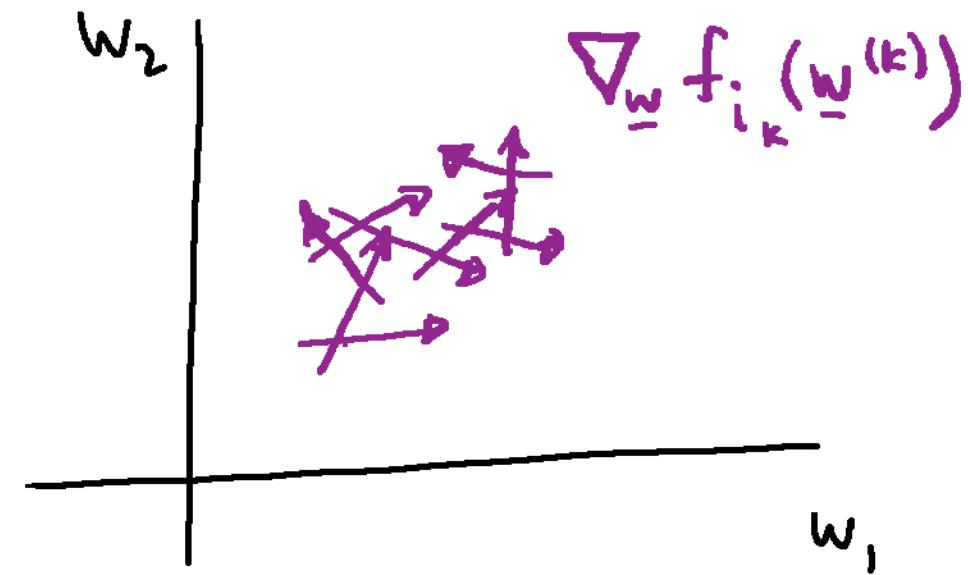
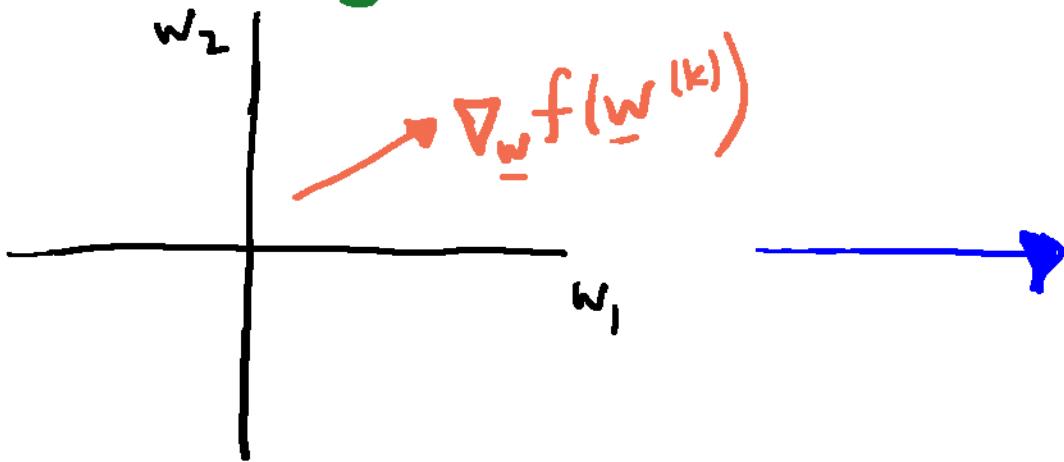
$$i_k = \text{uniform}\{1, 2, \dots, N\} \quad i_k = 2, 1, 3, 1, 4, 4, 2, 3, 1, 3 \dots$$

Update by $-\frac{\eta}{N} \nabla_{\underline{w}} f_{i_k}(\underline{w})$ at each iteration

On average gives gradient $E\{\nabla_{\underline{w}} f_{i_k}(\underline{w})\} \approx \frac{\nabla_{\underline{w}} f(\underline{w})}{N}$

SGD has computational benefits

- 1) Computing $\nabla_{\underline{w}} f_{i_k}(\underline{w}^{(k)})$ is easier/faster than $\nabla_{\underline{w}} f(\underline{w}^{(k)})$
- 2) May not be able to store $\underline{x}_i, i=1, \dots, N$ in memory
- 3) Noisy gradient $\nabla_{\underline{w}} f_{i_k}(\underline{w}^{(k)})$ introduces s added regularization



Example: Ridge Regression

5

$$f(\underline{w}) = \sum_{i=1}^N (d_i - \underline{x}_i^\top \underline{w})^2 + \lambda \|\underline{w}\|_2^2 = \sum_{i=1}^N \left\{ (d_i - \underline{x}_i^\top \underline{w})^2 + \frac{\lambda}{N} \|\underline{w}\|_2^2 \right\}$$

$f_i(\underline{w})$

$$\begin{aligned} \nabla_{\underline{w}} f_i(\underline{w}) &= \nabla_{\underline{w}} \left[(d_i - \underline{x}_i^\top \underline{w})^2 + \frac{\lambda}{N} \underline{w}^\top \underline{w} \right] \\ &= -2(d_i - \underline{x}_i^\top \underline{w}) \underline{x}_i + 2\frac{\lambda}{N} \underline{w} \end{aligned}$$

$$\begin{aligned} \underline{w}^{(k+1)} &= \underline{w}^{(k)} - \frac{\tau}{2} \nabla_{\underline{w}^{(k)}} f_{i_k}(\underline{w}^{(k)}) \\ &= \underline{w}^{(k)} + \tau (d_{i_k} - \underline{x}_{i_k}^\top \underline{w}^{(k)}) \underline{x}_{i_k} - \frac{\tau \lambda}{N} \underline{w}^{(k)} \end{aligned}$$

VS.

$$\underline{w}^{(k+1)} = \underline{w}^{(k)} + \tau \underline{A}^\top (\underline{A} \underline{w}^{(k)} - \underline{d}) - \lambda \tau \underline{w}^{(k)}$$

$\underline{A}: N \times M$

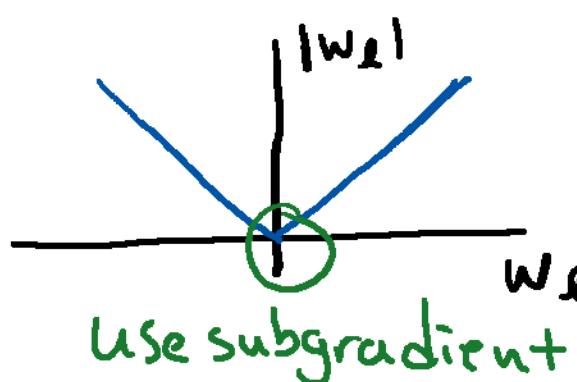
Example: Gradient descent for LASSO

6

$$f(\underline{w}) = \sum_{i=1}^N (d_i - \underline{x}_i^\top \underline{w})^2 + \lambda \|\underline{w}\|_1 = \sum_{i=1}^N \left\{ (d_i - \underline{x}_i^\top \underline{w})^2 + \frac{\lambda}{N} \|\underline{w}\|_1 \right\}$$

Consider $\nabla_{\underline{w}} \sum_{\ell=1}^M |w_\ell|$

Write $\nabla_{\underline{w}} \|\underline{w}\|_1 = \text{sign}(\underline{w})$



$$\frac{d}{dw_\ell} |w_\ell| = \begin{cases} \text{sign}(w_\ell) & w_\ell \neq 0 \\ [-1, 1] & w_\ell = 0 \end{cases}$$

"0" popular

$$\nabla_{\underline{w}} f_i(\underline{w}) = -2(d_i - \underline{x}_i^\top \underline{w}) \underline{x}_i + \frac{\lambda}{N} \text{sign}(\underline{w})$$

$$\underline{w}^{(k+1)} = \underline{w}^{(k)} + \tau (d_{i_k} - \underline{x}_{i_k}^\top \underline{w}^{(k)}) \underline{x}_{i_k} - \frac{\lambda \tau}{N} \text{sign}(\underline{w}^{(k)})$$

Copyright 2019
Barry Van Veen