

# Solving the Least-Squares Problem Using Gradients

# Objectives

- Introduce gradients of linear and quadratic functions of  $\underline{w}$
- Use gradients to solve least-squares problem
- Show solution is a minimizer
- Introduce projection matrices

# The Least-Squares Problem

2

$$\min_{\underline{w}} \|\underline{A}\underline{w} - \underline{d}\|_2^2 \quad \text{Note: } \|\underline{z}\|_2^2 = \underline{z}^T \underline{z}$$

$\underline{w}$  ↑       $\underline{A}$  ↑       $\underline{d}$  ↑  
 N features      P model parameters      N labels

$$(A\underline{w} - \underline{d})^T (A\underline{w} - \underline{d}) \Rightarrow (AB)^T = B^T A^T$$

$$\min_{\underline{w}} \underline{w}^T \underline{A}^T \underline{A} \underline{w} - \underline{w}^T \underline{A}^T \underline{d} - \underline{d}^T \underline{A} \underline{w} + \underline{d}^T \underline{d}$$

$f(\underline{w})$        $f(\underline{w})$        $f(\underline{w})$ : quadratic in  $\underline{w}$

Scalar problem:  $\min_{\underline{w}} \underbrace{a^2 w^2 - 2 \beta w + \gamma^2}_{g(w)}$  set  $\frac{d}{dw} g(w) = 0$

$$2a^2 w - 2\beta = 0 \Rightarrow w = \frac{-\beta}{a^2} \quad a^2 > 0 \Rightarrow \text{concave up, min}$$

\* Gradients: differentiate  $f(\underline{w})$  with respect to vector  $\underline{w}$

$$\frac{\partial}{\partial \underline{w}} f(\underline{w}) = \underline{0}$$

$L_2$  Norm is  $\|z\|_2$

On the previous page it is  $\|z\|_2^2 = z_1^2 + z_2^2 + \dots + z_n^2$

$$z^T z = \begin{bmatrix} z_1 \\ \vdots \\ z_n \end{bmatrix} [z_1 \ \dots \ z_n] = z_1^2 + \dots + z_n^2$$

$$\therefore z^T z = \|z\|_2^2$$

# Gradients

$$\nabla_{\underline{w}} f(\underline{w}) = \left[ \frac{\partial}{\partial w_1} f(\underline{w}) \quad \frac{\partial}{\partial w_2} f(\underline{w}) \cdots \frac{\partial}{\partial w_p} f(\underline{w}) \right]^T \quad \text{3}$$

Grad:

① Linear case

Suppose  $f(\underline{w}) = \underline{w}^T \underline{h} = \underline{h}^T \underline{w} = \sum_{i=1}^p w_i h_i \rightarrow \frac{\partial}{\partial w_j} f(\underline{w}) = h_j$

$$\nabla_{\underline{w}} f(\underline{w}) = [h_1 \ h_2 \ \dots \ h_p]^T = \underline{h} \quad \text{The gradients are just the } h_i \text{ (coefficients of each } w_i\text{)}$$

② Quadratic case

Suppose  $f(\underline{w}) = \underline{w}^T \underline{Q} \underline{w}$  Can show  $\nabla_{\underline{w}} f(\underline{w}) = \underline{Q}^T \underline{w} + \underline{Q} \underline{w}$

Symmetric case:  $\underline{Q} = \underline{Q}^T \Rightarrow \nabla_{\underline{w}} f(\underline{w}) = 2 \underline{Q} \underline{w}$

If  $\underline{Q} = \underline{A}^T \underline{A} \Rightarrow (\underline{A}^T \underline{A})^T = \underline{A}^T \underline{A}$  symmetric

Solution:  $\nabla_{\underline{w}} (\underline{w}^T \underline{A}^T \underline{A} \underline{w} - \underline{w}^T \underline{A}^T \underline{d} - \underline{d}^T \underline{A} \underline{w} + \underline{d}^T \underline{d}) = \underline{0}$

$$2 \underline{A}^T \underline{A} \underline{w} - \underline{A}^T \underline{d} - \underline{A}^T \underline{d} = \underline{0}$$

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

# Solution Attributes

4

$$\underline{w} = (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{d} \quad \underline{A}: N \times P, P \leq N, \text{rank } \underline{A} = P \Rightarrow (\underline{A}^T \underline{A})^{-1} \text{ exists}$$

Minimizer?  $f(\underline{w}) = \underline{w}^T \underline{A}^T \underline{A} \underline{w} - \underline{w}^T \underline{A}^T \underline{d} - \underline{d}^T \underline{A} \underline{w} + \underline{d}^T \underline{d}$

$$f(\underline{w}) = (\underline{w} - (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{d})^T \underline{A}^T \underline{A} (\underline{w} - (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{d}) + \underline{d}^T \underline{d} - \underline{d}^T \underline{A} (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{d}$$

$$= \cancel{\underline{w}^T \underline{A}^T \underline{A} \underline{w}} - \cancel{\underline{w}^T \underline{A}^T \underline{A} (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{d}} - \cancel{\underline{d}^T \underline{A} (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{w}} + \cancel{\underline{d}^T \underline{A} (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{A} (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{d}}$$

$\downarrow w_0$

$$+ \underline{d}^T \underline{d} - \underline{d}^T \underline{A} (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{d}$$

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

So  $f(\underline{w}) = \underline{z}^T(\underline{w}) \underline{A}^T \underline{A} \underline{z}(\underline{w}) + \underline{d}^T \underline{d} - \underline{d}^T \underline{A} (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{d}$

$$\underline{z}(\underline{w}) = \underline{w} - \underline{w}_0; \quad \underline{w}_0 = (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{d}$$

$\underline{A}^T \underline{A} > 0 \Rightarrow \min_{\underline{w}} f(\underline{w}) \text{ when } \underline{w} = \underline{w}_0$

$$\min_{\underline{w}} f(\underline{w}) = \underline{d}^T \underline{d} - \underline{d}^T \underline{A} (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{d}$$

So why are we doing this?

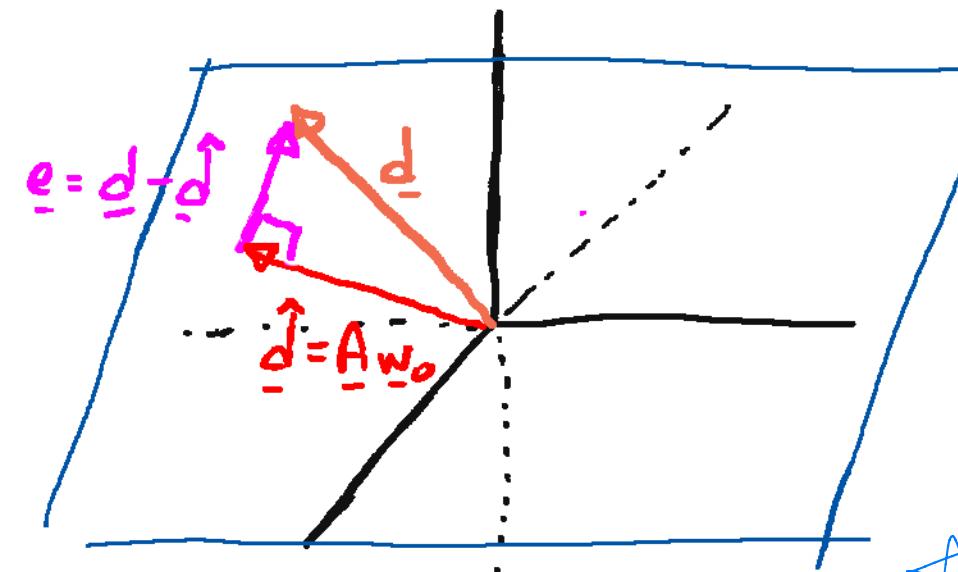
Completing the square :  $f(w) = (w - w_0)^T A^T A (w - w_0) + \text{constant}$

$(w_0 = (A^T A)^{-1} A^T d \text{ is the best solution})$

when  $w = w_0 \rightarrow \text{minimum and that is } 0$

$$\begin{aligned}\therefore \min f(w) &= d^T d - \underbrace{b^T A (A^T A)^{-1} A^T d}_d \\ &= d^T d - d^T d\end{aligned}$$

# Projection and the Pythagorean Theorem



Right triangle

$$\|\underline{d}\|_2^2 = \|\underline{e}\|_2^2 + \|\hat{\underline{d}}\|_2^2$$

$$\Rightarrow \|\underline{e}\|_2^2 = \|\underline{d}\|_2^2 - \|\hat{\underline{d}}\|_2^2$$

$$= \underline{d}^\top \underline{d} - \underline{d}^\top \underline{P}_A \underline{d}$$

$$= \underline{d}^\top (\underline{\mathbb{I}} - \underline{P}_A) \underline{d}$$

$$\hat{\underline{d}} = \underline{A}(\underline{A}^\top \underline{A})^{-1} \underline{A}^\top \underline{d} = \underline{P}_A \underline{d}$$

$\underline{P}_A$  projects  $\underline{d}$  onto  $\text{span}\{\underline{A}^T\}$

$\underline{P}_A = \underline{A}(\underline{A}^\top \underline{A})^{-1} \underline{A}^\top$  "projection matrix"

$$\underline{P}_A^2 = \underline{A}(\underline{A}^\top \underline{A})^{-1} \underline{A}^\top \underline{A} (\underline{A}^\top \underline{A})^{-1} \underline{A}^\top = \underline{A}(\underline{A}^\top \underline{A})^{-1} \underline{A}^\top = \underline{P}_A$$

(No matter what power, always  $\underline{P}_A$  for projection)

$$\star \underline{e} = \underline{d} - \hat{\underline{d}} = (\underline{\mathbb{I}} - \underline{P}_A) \underline{d} = \underline{P}_{A\perp} \underline{d}$$

$\star \underline{P}_{A\perp} = \underline{\mathbb{I}} - \underline{P}_A$  projects onto space  $\perp$  to  $\text{span}\{\underline{A}^T\}$

$$\underline{P}_{A\perp}^2 = (\underline{\mathbb{I}} - \underline{P}_A)(\underline{\mathbb{I}} - \underline{P}_A) = \underline{\mathbb{I}} - 2\underline{P}_A + \underline{P}_A^2$$

$$= \underline{\mathbb{I}} - \underline{P}_A = \underline{P}_{A\perp}$$

$\underline{P}_A$  perpendicular : same as above, no matter how many power raised to it. It'll always be  $\underline{P}_{A\perp}$

$$\|\underline{d}^\top \underline{d} \underline{P}_A^2 - \underline{d}^\top \underline{P}_A \underline{d}\|$$

(This is also a projection matrix)

Copyright 2019  
Barry Van Veen