

CPSC 340: Machine Learning and Data Mining

The Normal Equations

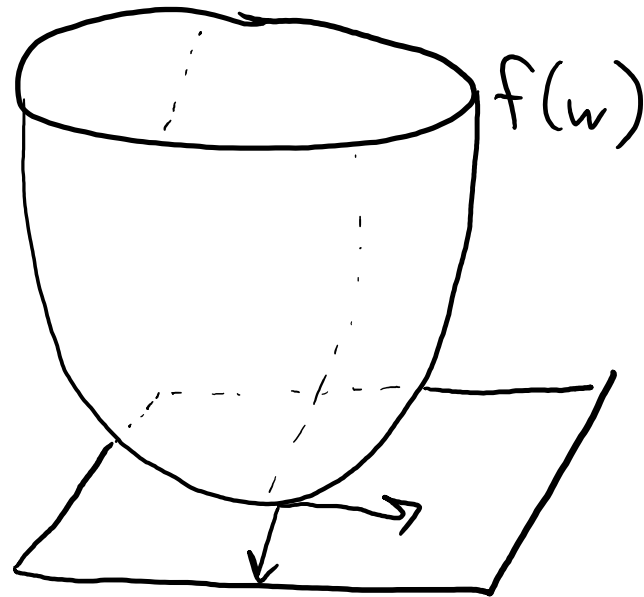
Admin

- a3 posted, due Feb 9
- Midterm Feb 14 in class
- New office hour on Wednesdays, per your feedback
 - In general, check calendar regularly for updates

Gradient and Critical Points in d-Dimensions

- Generalizing “set the derivative to 0 and solve” in d-dimensions:
 - Find ‘w’ where the **gradient** vector **equals the zero vector**.
- **Gradient** is vector with partial derivative ‘j’ in position ‘j’:

$$\nabla f(w) = \begin{bmatrix} \frac{\partial f}{\partial w_1} \\ \frac{\partial f}{\partial w_2} \\ \vdots \\ \frac{\partial f}{\partial w_d} \end{bmatrix}$$



Tangent slope is 0 in every direction at minimizer.

Gradient and Critical Points in d-Dimensions

- Generalizing “set the derivative to 0 and solve” in d-dimensions:
 - Find ‘w’ where the **gradient** vector **equals the zero vector**.
- **Gradient** is vector with partial derivative ‘j’ in position ‘j’:

$$\nabla f(w) = \begin{bmatrix} \frac{\partial f}{\partial w_1} \\ \frac{\partial f}{\partial w_2} \\ \vdots \\ \frac{\partial f}{\partial w_d} \end{bmatrix}$$

For linear least squares:

$$\nabla f(w) = \begin{bmatrix} \sum_{i=1}^n (w^T x_i - y_i) x_{i1} \\ \sum_{i=1}^n (w^T x_i - y_i) x_{i2} \\ \vdots \\ \sum_{i=1}^n (w^T x_i - y_i) x_{id} \end{bmatrix}$$

Claims for linear least square:

1. Finding a ‘w’ where $\nabla f(w) = 0$ can be done by solving a system of linear equations.
2. All ‘w’ where $\nabla f(w) = 0$ are minimizers.

Least Squares in d-Dimensions

- The **linear least squares** model in d-dimensions minimizes:

$$f(w) = \frac{1}{2} \sum_{i=1}^n (w^T x_i - y_i)^2$$

$$\begin{aligned} w^T x_i &= w_1 x_{i1} + w_2 x_{i2} + \dots + w_d x_{id} \\ \frac{d}{dw_1} [w^T x_i] &= x_{i1} + 0 + \dots + 0 \\ &= x_{i1} \end{aligned}$$

- Computing the **partial derivative**:

$$\begin{aligned} \frac{\partial}{\partial w_1} \left[\frac{1}{2} \sum_{i=1}^n (w^T x_i - y_i)^2 \right] &= \frac{1}{2} \sum_{i=1}^n \frac{\partial}{\partial w_1} \left[(w^T x_i - y_i)^2 \right] \\ &= \frac{1}{2} \sum_{i=1}^n 2 (w^T x_i - y_i) \frac{\partial}{\partial w_1} [w^T x_i] \end{aligned}$$

Problem: I can't just set to 0 and solve because it depends on w_2, w_3, \dots, w_d

$$= \sum_{i=1}^n (w^T x_i - y_i) x_{i1}$$

What is the derivative of $w^T x_i$ with respect to w_1 ?

Matrix/Norm Notation (MEMORIZE/STUDY THIS)

- To solve the d-dimensional least squares, we use **matrix notation**:
 - We use 'y' as an "n times 1" vector containing target 'y_i' in position 'i'.
 - We use 'x_i' as a "d times 1" vector containing features 'j' of example 'i'.
 - We're now going to be careful to make sure these are **column vectors**.
 - So 'X' is a matrix with the x_i^T in row 'i'.

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$$x_i = \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{id} \end{bmatrix}$$

$$X = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1d} \\ x_{21} & x_{22} & \dots & x_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nd} \end{bmatrix} = \begin{bmatrix} \text{---} & x_1^T & \text{---} \\ \text{---} & x_2^T & \text{---} \\ \vdots & \vdots & \vdots \\ \text{---} & x_n^T & \text{---} \end{bmatrix}$$

Matrix/Norm Notation (MEMORIZE/STUDY THIS)

- To solve the d-dimensional least squares, we use **matrix notation**:
 - Our **prediction** for example 'i' is given by scalar $w^T x_i$.
 - The **matrix-vector product** Xw gives predictions for all 'i' (n times 1 vector).

$$\begin{aligned}w^T x_i &= \sum_{j=1}^d w_j x_{ij} \\ &= w_1 x_{i1} + w_2 x_{i2} + \dots + w_d x_{id}\end{aligned}$$

Also, because $w^T x_i$ is a scalar,
we have $w^T x_i = x_i^T w$.
(e.g., $[5]^T = [5]$)

$$\begin{aligned}Xw &= \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1d} \\ x_{21} & x_{22} & \dots & x_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nd} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} = \begin{bmatrix} x_{11}w_1 + x_{12}w_2 + \dots + x_{1d}w_d \\ x_{21}w_1 + x_{22}w_2 + \dots + x_{2d}w_d \\ \vdots \\ x_{n1}w_1 + x_{n2}w_2 + \dots + x_{nd}w_d \end{bmatrix} \\ &= \begin{bmatrix} \sum_{j=1}^d x_{1j}w_j \\ \sum_{j=1}^d x_{2j}w_j \\ \vdots \\ \sum_{j=1}^d x_{nj}w_j \end{bmatrix} = \begin{bmatrix} x_1^T w \\ x_2^T w \\ \vdots \\ x_n^T w \end{bmatrix} = \begin{bmatrix} w^T x_1 \\ w^T x_2 \\ \vdots \\ w^T x_d \end{bmatrix} \end{aligned}$$

Prediction for example 'i' is in column 'i'

Matrix/Norm Notation (MEMORIZE/STUDY THIS)

- To solve the d-dimensional least squares, we use **matrix notation**:
 - Our **prediction** for example 'i' is given by scalar $w^T x_i$.
 - The **matrix-vector product** Xw gives predictions for all 'i' (n times 1 vector).
 - The **residual vector** r gives $w^T x_i$ minus y_i for all 'i' (n times 1 vector).
 - **Least squares** can be written as the squared L2-norm of the residual.

$$r = \begin{bmatrix} w^T x_1 - y_1 \\ w^T x_2 - y_2 \\ \vdots \\ w^T x_n - y_n \end{bmatrix} = \underbrace{\begin{bmatrix} w^T x_1 \\ w^T x_2 \\ \vdots \\ w^T x_n \end{bmatrix}}_{Xw} - \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}}_y = Xw - y$$
$$\begin{aligned} \sum_{i=1}^n (w^T x_i - y_i)^2 &= \sum_{i=1}^n (r_i)^2 \\ &= \sum_{i=1}^n r_i r_i \\ &= r^T r \\ &= \|r\|^2 = \|Xw - y\|^2 \end{aligned}$$

Matrix Algebra Review (MEMORIZE/STUDY THIS)

- Review of **linear algebra operations** we'll use:
 - If 'a' and 'b' be vectors, and 'A' and 'B' be matrices then:

$$a^T b = b^T a$$

$$\|a\|^2 = a^T a$$

$$(A+B)^T = A^T + B^T$$

$$(AB)^T = B^T A^T$$

$$(A+B)(A+B) = AA + BA + AB + BB$$

$$a^T \underbrace{A}_{\text{vector}} b = b^T \underbrace{A^T}_{\text{vector}} a$$

Sanity check:

ALWAYS CHECK THAT
DIMENSIONS MATCH
(if not, you did something wrong)

Linear Least Squares

Want 'w' that minimizes

$$f(w) = \frac{1}{2} \sum_{i=1}^n (w^T x_i - y_i)^2 = \frac{1}{2} \|Xw - y\|_2^2 = \frac{1}{2} (Xw - y)^T (Xw - y)$$

Let's expand
then compute
gradient.

$$= \frac{1}{2} ((Xw)^T - y^T) (Xw - y)$$

$$= \frac{1}{2} (w^T X^T - y^T) (Xw - y)$$

$$= \frac{1}{2} (w^T X^T (Xw - y) - y^T (Xw - y))$$

$$= \frac{1}{2} (w^T X^T Xw - w^T X^T y - y^T Xw + y^T y)$$

$$= \frac{1}{2} w^T X^T Xw - w^T X^T y + \frac{1}{2} y^T y$$

Sanity check: all of these are scalars.

Linear and Quadratic Gradients

- We've written as a **d-dimensional quadratic**:

$$f(w) = \frac{1}{2} \sum_{i=1}^n (w^T x_i - y_i)^2 = \frac{1}{2} \|Xw - y\|^2 = \frac{1}{2} \underbrace{w^T X^T X w}_{\text{matrix 'A'}} - \underbrace{w^T X^T y}_{\text{vector 'b'}} + \frac{1}{2} \underbrace{y^T y}_{\text{scalar 'c'}} \\ = \frac{1}{2} w^T A w + w^T b + c$$

- How do we compute gradient?

Let's first do it with $d=1$:

$$f(w) = \frac{1}{2} a w^2 + w b + c \\ = \frac{1}{2} a w^2 + w b + c$$

$$f'(w) = a w + b + 0$$

Here are the generalizations to 'd' dimensions:

$$\nabla [c] = 0 \quad (\text{zero vector})$$

$$\nabla [w^T b] = b$$

$$\frac{1}{2} \nabla [w^T A w] = A w \quad (\text{if } A \text{ is } \underline{\text{symmetric}})$$

Full derivations are on webpage in notes on linear and quadratic gradients.

Linear and Quadratic Gradients

- We've written the least squares objective as a quadratic function:

$$\begin{aligned} f(w) &= \frac{1}{2} \sum_{i=1}^n (w^T x_i - y_i)^2 = \frac{1}{2} \|Xw - y\|^2 = \frac{1}{2} \underbrace{w^T X^T X w}_{\text{matrix 'A'}} - \underbrace{w^T X^T y}_{\text{vector 'b'}} + \frac{1}{2} \underbrace{y^T y}_{\text{scalar 'c'}} \\ &= \frac{1}{2} w^T A w + w^T b + c \end{aligned}$$

- Gradient is given by: $\nabla f(w) = Aw + b + 0$

- Using definitions of 'A' and 'b': $= X^T X w - X^T y = 0$

Sanity check: these are both $d \times 1$ vectors.

Normal Equations

- Set gradient equal to zero to find the least squares “critical points”:

$$X^T X w - X^T y = 0$$

- We now move terms not involving ‘w’ to the other side:

$$X^T X w = X^T y$$

- This is a set of ‘d’ linear equations called the normal equations.
 - This a linear system like “Ax = b” from Math 152.
 - You can use Gaussian elimination to solve for ‘w’.
 - In Python, `numpy.linalg.solve` can be used to solve linear systems.

Incorrect Solutions to Least Squares Problem

The least squares objective is $f(w) = \frac{1}{2} \|Xw - y\|^2$

The minimizers of this objective are solutions to the linear system:

$$X^T X w = X^T y$$

The following are not the solutions to the least squares problem:

$$w = (X^T X)^{-1} (X^T y) \quad (\text{only true if } \underline{X^T X \text{ is invertible}})$$

$$w X^T X = X^T y \quad (\text{matrix multiplication is } \underline{\text{not}} \text{ commutative, dimensions don't even match})$$

$$w = \frac{X^T y}{X^T X} \quad (\text{you } \underline{\text{cannot divide by a matrix}})$$

Least Squares Issues

- Issues with least squares model:
 - Solution might **not be unique**.
 - It is **sensitive to outliers**.
 - It always **uses all features**.
 - Data can be so big we **can't store $X^T X$** .
 - It might **predict outside range** of y_i values.
 - It assumes a **linear relationship** between x_i and y_i .

→ X is $n \times d$
so X^T is $d \times n$
and $X^T X$ is $d \times d$.

Least Squares cost

- Forming matrix $X^T X$ costs $O(nd^2)$
 - because $X^T X$ has d^2 elements and each is a sum of n numbers.
- Solving system $X^T X w = X^T y$ costs $O(d^3)$
 - because we are solving a d -by- d linear system.
- Overall cost is $O(nd^2 + d^3)$
 - Which term dominates depends on how ‘ n ’ compares to ‘ d ’
 - $n > d$ is the standard case
 - $d > n$ is a bit trickier, solution not unique (“underdetermined” system)
 - Put another way, we have ‘ n ’ equations and ‘ d ’ unknowns/variables
 - Imagine our 2d plots with $n < 2$ points... that would be just one point
 - Remember it’s not correct to write $O(nd^2) + O(d^3)$

Non-Uniqueness: Colinearity

- Imagine have two features that are identical for all examples.
- Then these features are called **collinear**.
- I can increase weight on one feature, and decrease it on the other, **without changing predictions**.
- Thus the solution is not unique.
- But, any 'w' where $\nabla f(w) = 0$ is a global optimum, due to **convexity**.
- We will revisit the uniqueness issue soon when we cover **regularization** in a couple lectures.

Convexity of Linear Regression

- Consider linear regression objective with squared error:

$$f(w) = \|Xw - y\|^2$$

- This is a **convex function composed with linear**:

Let $g(r) = \|r\|^2$, which is convex because it's a squared norm.

Then $f(w) = g(Xw - y)$, which is convex because it's a convex function composed with the linear function $h(w) = Xw - y$.

Summary

- Normal equations: solution of least squares as a linear system.
 - Solve $(X^T X)w = (X^T y)$.
- Solution might not be unique because of **collinearity**.
- But any solution is optimal because of **convexity**.
- **Convex functions**:
 - Set of functions with property that $\nabla f(w) = 0$ implies 'w' is a global min.
 - Can (usually) be identified using a few simple rules.

Convexity, min, and argmin

- If a function is convex, then all stationary points are global optima.
- However, **convex functions don't necessarily have stationary points:**
 - For example, $f(x) = a \cdot x$, $f(x) = \exp(x)$, etc.
- Also, **more than one 'x' can achieve the global optimum:**
 - For example, $f(x) = c$ is minimized by any 'x'.

Bonus Slide: Householder(-ish) Notation

- **Householder notation:** set of (fairly-logical) conventions for math.

Use greek letters for scalars: $\alpha = 1$, $\beta = 3.5$, $\gamma = \pi$

Use first/last lowercase letters for vectors: $w = \begin{bmatrix} 0.1 \\ 0.2 \end{bmatrix}$, $x = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$, $y = \begin{bmatrix} 2 \\ -1 \end{bmatrix}$, $a = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$, $b = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}$

↳ Assumed to be column-vectors.

Use first/last uppercase letters for matrices: X, Y, W, A, B

Indices use i, j, k .

Sizes use m, n, d, p , and k

← hopefully meaning of 'k' is obvious from context

Sets use S, T, U, V

Functions use f, g , and h .

When I write x_i I mean "grab row 'i' of X and make a column-vector with its values." 21

Bonus Slide: Householder(-ish) Notation

- **Householder notation:** set of (fairly-logical) conventions for math:

Our ultimate least squares notation:

$$f(w) = \frac{1}{2} \|Xw - y\|^2$$

But if we agree on notation we can quickly understand:

$$g(x) = \frac{1}{2} \|Ax - b\|^2$$

If we use random notation we get things like:

$$H(\beta) = \frac{1}{2} \|R\beta - p_n\|^2$$

Is this the same model?

When does least squares have a unique solution?

- We said that least squares solution is not unique if we have repeated columns.
- But there are other ways it could be non-unique:
 - One column is a scaled version of another column.
 - One column could be the sum of 2 other columns.
 - One column could be three times one column minus four times another.
- Least squares solution is unique if and only if all columns of X are “linearly independent”.
 - No column can be written as a “linear combination” of the others.
 - Many equivalent conditions (see Strang’s linear algebra book):
 - X has “full column rank”, $X^T X$ is invertible, $X^T X$ has non-zero eigenvalues, $\det(X^T X) > 0$.
 - Note that we **cannot have independent columns if $d > n$** .