



ISTA 421 + INFO 521

Introduction to Machine Learning

Lecture 3: Moving to Higher dimensions

Clay Morrison

claytonm@email.arizona.edu

Harvill 437A

Phone 621-6609

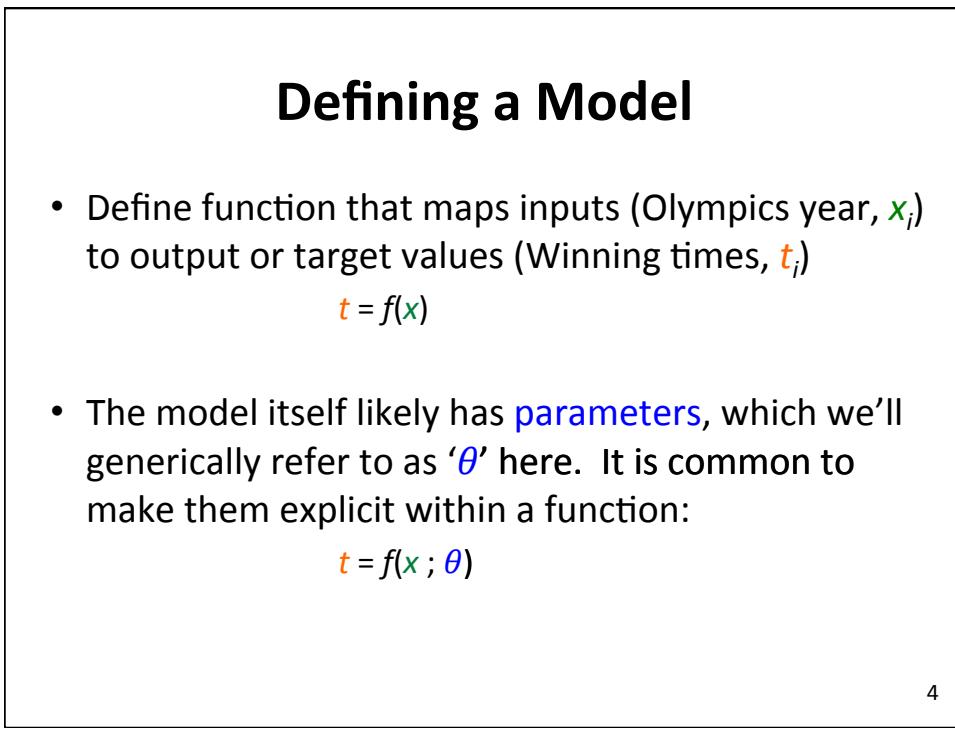
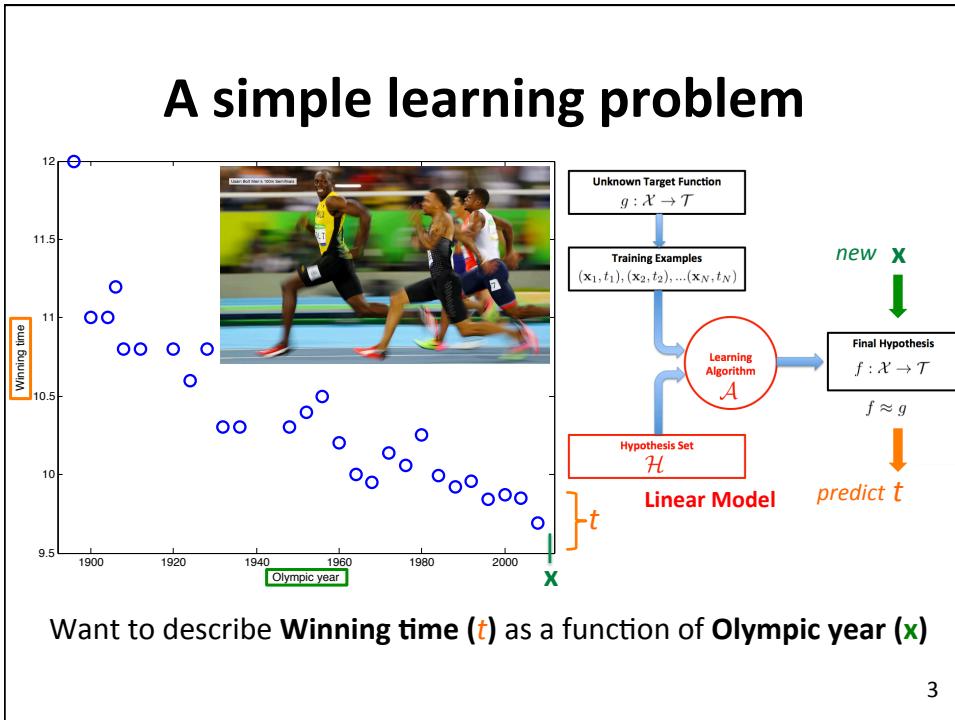
28 August 2017

1

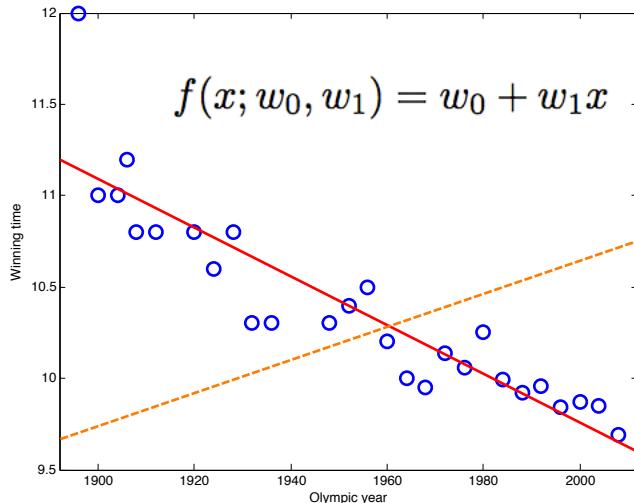
Next Topics

- Linear Model (what is the *model*)
- Loss Function (what is a *good* model)
- Least Squares (finding the *best* model)
- Prediction
- Review LMS fit (the “normal equations”) for the single variable, 2 parameter **linear model**
- Moving to higher dimensions
 - Linear Algebra: matrix operators
 - Some Geometry of Linear Algebra
 - Least Mean Squares in Matrix formulation
 - **The Geometry of LMS solution**
- Nonlinear Response
- Generalization and Overfitting
- Regularized Least Squares

2



Data with line (particular w_0 & w_1)



(The red line happens to be a “best” fit)

(The dashed orange line does not describe the trend in the data very well; not a good fit)

5

Goal: Find the “best” values for the parameters of model according to the loss fn

- If we’re lucky (i.e., the model and the loss fn are “well-behaved”) we can derive an **analytic** solution. Otherwise, we’ll pick some (iterative) optimization method that is appropriate.
- Our first example, using a **mean squared error loss function** with a **linear model** permits a nice analytic solution!
 - Here (and in the book) we’ll first look at the direct, analytic method.
 - Another method: gradient descent
 - Same loss function, but iterative algorithm and can be used in cases where we don’t have an analytic solution for the parameters

6

Least Mean Squares Solution

(for single variable, 2 parameter linear model)

$$f(x; w_0, w_1) = w_0 + w_1 x$$

$$\begin{aligned}
 \mathcal{L} &= \frac{1}{N} \sum_{n=1}^N \mathcal{L}_n(t_n, f(x_n; w_0, w_1)) \\
 &= \frac{1}{N} \sum_{n=1}^N (t_n - f(x_n; w_0, w_1))^2 && \text{The specific loss fn we're working with here} \\
 &= \frac{1}{N} \sum_{n=1}^N (t_n - (w_0 + w_1 x_n))^2 && \text{The specific model we're working with here} \\
 &= \frac{1}{N} \sum_{n=1}^N (w_1^2 x_n^2 + 2w_1 x_n w_0 - 2w_1 x_n t_n + w_0^2 - 2w_0 t_n + t_n^2) && \text{Multiply out and re-arrange to put into an easier-to deal with form.} \\
 &= \frac{1}{N} \sum_{n=1}^N (w_1^2 x_n^2 + 2w_1 x_n (w_0 - t_n) + w_0^2 - 2w_0 t_n + t_n^2)
 \end{aligned}$$

7

Least Mean Squares Solution

(for single variable, 2 parameter linear model)

$$\begin{aligned}
 f(x; w_0, w_1) &= w_0 + w_1 x && \text{Our model family} \\
 \mathcal{L} &= \frac{1}{N} \sum_{n=1}^N (w_1^2 x_n^2 + 2w_1 x_n (w_0 - t_n) + w_0^2 - 2w_0 t_n + t_n^2) && \text{Our loss fn}
 \end{aligned}$$

Our goal: We want values for w_0 and w_1 that will **minimize** this loss function

i.e., we seek values for w_0 and w_1 that will make the loss function be the smallest when we actually sum over all the values of x and t in the dataset.

Because the loss function happens to be quadratic (in the two parameters) we can use a standard method from calculus for finding minima (maxima) directly: **taking the derivative of the function and setting it to zero.**

Our loss function has **two** parameters that we're trying set to minimize the loss fn, so we need to take the **partial derivative** (w.r.t. w_0 and w_1)

What we end up with are two functions, one for w_0 and one for w_1 , and both will work with **any** data and give the best **least mean square** (LMS) fit!

Least Mean Squares Solution

(for single variable, 2 parameter linear model)

$$\mathcal{L} = \frac{1}{N} \sum_{n=1}^N (w_1^2 x_n^2 + 2w_1 x_n (w_0 - t_n) + w_0^2 - 2w_0 t_n + t_n^2)$$

- Partial derivative for w_0

First, since we're taking the partial w.r.t. w_0 , can drop any terms without w_0 .

$$\frac{1}{N} \sum_{n=1}^N [w_0^2 + 2w_1 x_n w_0 - 2w_0 t_n]$$

$$w_0^2 + 2w_0 w_1 \frac{1}{N} \left(\sum_{n=1}^N x_n \right) - 2w_0 \frac{1}{N} \left(\sum_{n=1}^N t_n \right)$$

Next, move sums inward to put in easier form

$$\frac{\partial \mathcal{L}}{\partial w_0} = 2w_0 + 2w_1 \frac{1}{N} \left(\sum_{n=1}^N x_n \right) - \frac{2}{N} \left(\sum_{n=1}^N t_n \right)$$

Finally, take deriv. w.r.t. w_0

9

Continued...

- Solve for $\frac{\partial \mathcal{L}}{\partial w_0} = 0$

$$2w_0 + 2w_1 \frac{1}{N} \left(\sum_{n=1}^N x_n \right) - \frac{2}{N} \left(\sum_{n=1}^N t_n \right) = 0$$

$$2w_0 = \frac{2}{N} \left(\sum_{n=1}^N t_n \right) - w_1 \frac{2}{N} \left(\sum_{n=1}^N x_n \right)$$

$$w_0 = \frac{1}{N} \left(\sum_{n=1}^N t_n \right) - w_1 \frac{1}{N} \left(\sum_{n=1}^N x_n \right)$$

10

Least Mean Squares Solution

(for single variable, 2 parameter linear model)

$$\mathcal{L} = \frac{1}{N} \sum_{n=1}^N (w_1^2 x_n^2 + 2w_1 x_n (w_0 - t_n) + w_0^2 - 2w_0 t_n + t_n^2)$$

- Partial derivative for w_1 Do the same for w_0 ...

$$\begin{aligned} & \frac{1}{N} \sum_{n=1}^N [w_1^2 x_n^2 + 2w_1 x_n w_0 - 2w_1 x_n t_n] \quad \text{only keep terms with } w_1 \\ & w_1^2 \frac{1}{N} \left(\sum_{n=1}^N x_n^2 \right) + 2w_1 \frac{1}{N} \left(\sum_{n=1}^N x_n (w_0 - t_n) \right) \quad \text{move sums inside} \\ & \frac{\partial \mathcal{L}}{\partial w_1} = 2w_1 \frac{1}{N} \left(\sum_{n=1}^N x_n^2 \right) + \frac{2}{N} \left(\sum_{n=1}^N x_n (w_0 - t_n) \right) \quad \text{now take partial derivative} \\ & = w_1 \frac{2}{N} \left(\sum_{n=1}^N x_n^2 \right) + \frac{2}{N} \left(\sum_{n=1}^N x_n \left(\underbrace{\bar{t} - w_1 \bar{x}}_{\text{solution for } w_0} - t_n \right) \right) \quad \text{plug in solution for } w_0 \\ & = w_1 \frac{2}{N} \left(\sum_{n=1}^N x_n^2 \right) + \bar{t} \frac{2}{N} \left(\sum_{n=1}^N x_n \right) - w_1 \bar{x} \frac{2}{N} \left(\sum_{n=1}^N x_n \right) - \frac{2}{N} \left(\sum_{n=1}^N x_n t_n \right) \quad 11 \end{aligned}$$

- Partial derivative for w_1 continued...

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial w_1} &= w_1 \frac{2}{N} \left(\sum_{n=1}^N x_n^2 \right) + \bar{t} \frac{2}{N} \left(\sum_{n=1}^N x_n \right) - w_1 \bar{x} \frac{2}{N} \left(\sum_{n=1}^N x_n \right) - \frac{2}{N} \left(\sum_{n=1}^N x_n t_n \right) \\ &= 2w_1 \left[\left(\frac{1}{N} \sum_{n=1}^N x_n^2 \right) - \bar{x} \bar{x} \right] + 2\bar{t}\bar{x} - 2 \frac{1}{N} \left(\sum_{n=1}^N x_n t_n \right) \quad \text{replace remaining mean } x \text{ with } \bar{x} \text{ and group } w_1 \text{ terms} \end{aligned}$$

$$\begin{aligned} 2w_1 \left[\left(\frac{1}{N} \sum_{n=1}^N x_n^2 \right) - \bar{x} \bar{x} \right] + 2\bar{t}\bar{x} - 2 \frac{1}{N} \left(\sum_{n=1}^N x_n t_n \right) &= 0 \quad \text{Solve for } w_1 \text{ with } \frac{\partial \mathcal{L}}{\partial w_1} = 0 \dots \\ 2w_1 \left[\left(\frac{1}{N} \sum_{n=1}^N x_n^2 \right) - \bar{x} \bar{x} \right] &= 2 \frac{1}{N} \left(\sum_{n=1}^N x_n t_n \right) - 2\bar{t}\bar{x} \end{aligned}$$

$$w_1 = \frac{\frac{1}{N} \left(\sum_{n=1}^N x_n t_n \right) - \bar{t} \bar{x}}{\left(\frac{1}{N} \sum_{n=1}^N x_n^2 \right) - \bar{x} \bar{x}} = \frac{\left(\frac{1}{N} \sum_{n=1}^N x_n t_n \right) - \left(\frac{1}{N} \sum_{m=1}^N t_n \right) \left(\frac{1}{N} \sum_{m=1}^N x_n \right)}{\left(\frac{1}{N} \sum_{n=1}^N x_n^2 \right) - \left(\frac{1}{N} \sum_{n=1}^N x_n \right)^2} = \frac{\bar{x}\bar{t} - \bar{x}\bar{t}}{\bar{x}^2 - (\bar{x})^2}$$

Solving LMS: Method 1 (analytic)

(for single variable, 2 parameter linear model)

$$\mathcal{L} = \frac{1}{N} \sum_{n=1}^N (w_0^2 x_n^2 + 2w_1 x_n(w_0 - t_n) + w_0^2 - 2w_0 t_n + t_n^2)$$

- Partial derivative for w_0 :

$$\frac{\partial \mathcal{L}}{\partial w_0} = 2w_0 + 2w_1 \frac{1}{N} \left(\sum_{n=1}^N x_n \right) - \frac{2}{N} \left(\sum_{n=1}^N t_n \right)$$

- Set $\frac{\partial \mathcal{L}}{\partial w_0} = 0$ and solve for w_0 :

$$w_0 = \frac{1}{N} \left(\sum_{n=1}^N t_n \right) - w_1 \frac{1}{N} \left(\sum_{n=1}^N x_n \right) = \bar{t} - w_1 \bar{x}$$

The so-called
normal equations!

- Partial derivative for w_1 :

$$\frac{\partial \mathcal{L}}{\partial w_1} = 2w_1 \frac{1}{N} \left(\sum_{n=1}^N x_n^2 \right) + \frac{2}{N} \left(\sum_{n=1}^N x_n(w_0 - t_n) \right)$$

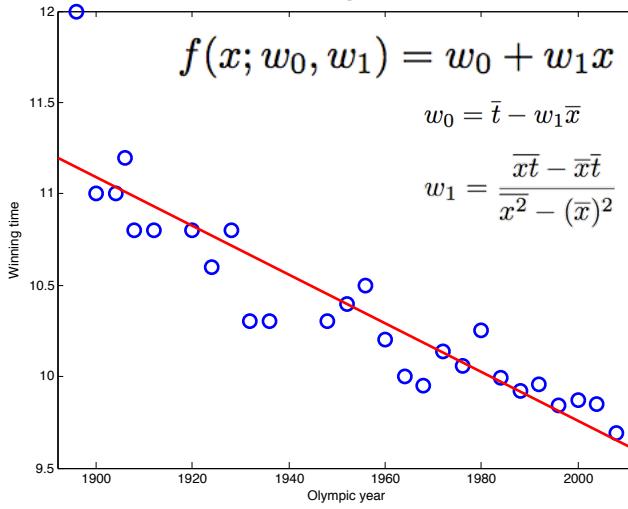
- Plug in w_0 , set $\frac{\partial \mathcal{L}}{\partial w_1} = 0$ and solve for w_1 :

$$w_1 = \frac{\left(\frac{1}{N} \sum_{n=1}^N x_n t_n \right) - \left(\frac{1}{N} \sum_{n=1}^N t_n \right) \left(\frac{1}{N} \sum_{n=1}^N x_n \right)}{\left(\frac{1}{N} \sum_{n=1}^N x_n^2 \right) - \left(\frac{1}{N} \sum_{n=1}^N x_n \right)^2} = \frac{\bar{x}\bar{t} - \bar{t}\bar{x}}{\bar{x}^2 - (\bar{x})^2}$$

13

Linear Least Mean Square

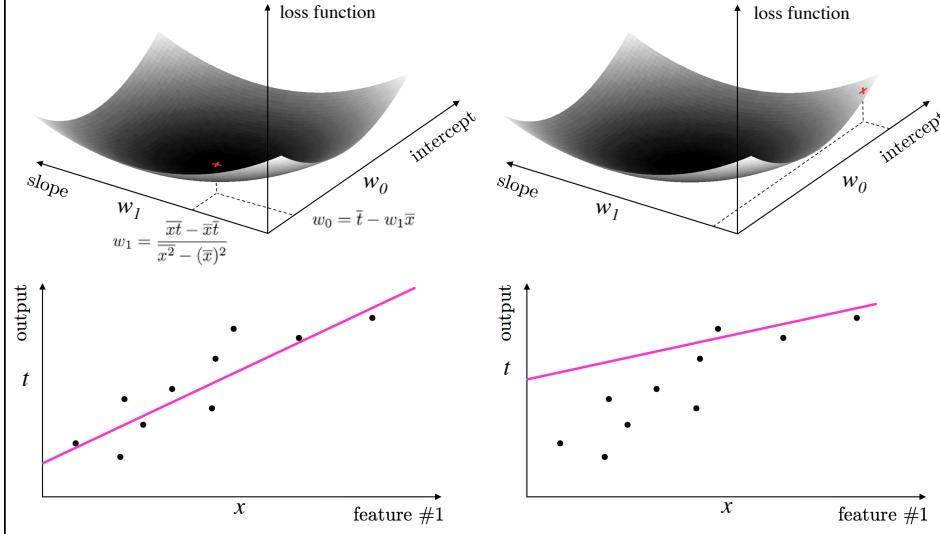
Normal Equations



14

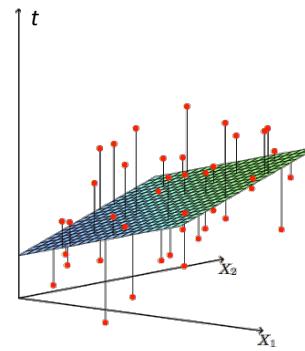
Visualizing the 2-dim loss fn surface

$$f(x; w_0, w_1) = w_0 + w_1 x \quad \mathcal{L} = \frac{1}{N} \sum_{n=1}^N (w_1^2 x_n^2 + 2w_1 x_n (w_0 - t_n) + w_0^2 - 2w_0 t_n + t_n^2)$$



How about more than 1 input?

- Most problems will involve more than just the relationship between 1 input attribute and a target.
- Extending our linear models to higher dimensions is desirable. For 2 inputs it is easy to visualize the geometry: now the “line” is a plane in 3D
- In general, a (regression) linear model with n input variables and $n+1$ parameters (the w 's, with their values determined) is an n -dimensional “hyperplane” embedded in $n+1$ dimensions.



16

Things quickly get messy as we increase the dimensions...

- Suppose we want a richer predictive model for the Olympic data: not only include the best overall time for the gold, but also the personal-best times of each sprinter that raced (s_1, \dots, s_8)
 - This is a 9 dimensional hyperplane with 10 parameters:
- $$\begin{aligned} t = f(x, s_1, \dots, s_8; w_0, \dots, w_9) &= w_0 + w_1 x + w_2 s_1 + w_3 s_2 \\ &\quad + w_4 s_3 + w_5 s_4 + w_6 s_5 \\ &\quad + w_7 s_6 + w_8 s_7 + w_9 s_8 \end{aligned}$$
- The math is fundamentally the same, but to derive the normal equations, we need to take 10 partial derivatives, then have 10 equations to re-arrange and substitute back in...

17

Dot product $\sum_{n=1}^N x_n^2$

Matrix Product

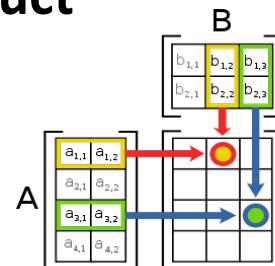
Let $\mathbf{C} = \mathbf{AB}$, where

\mathbf{A} is a $N \times P$ matrix

\mathbf{B} is a $P \times M$ matrix

\mathbf{C} is a $N \times M$ matrix

and each entry C of \mathbf{C} is: $C_{ij} = \sum_k A_{ik} B_{kj}$



$$\begin{bmatrix} b_{11} & b_{12} & b_{13} \\ b_{21} & b_{22} & b_{23} \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} a_{11}b_{11} + a_{12}b_{21} & a_{11}b_{12} + a_{12}b_{22} & a_{11}b_{13} + a_{12}b_{23} \\ a_{21}b_{11} + a_{22}b_{21} & a_{21}b_{12} + a_{22}b_{22} & a_{21}b_{13} + a_{22}b_{23} \end{bmatrix}$$

(1st edition of FCML, has small error on p.18: copies first line, but index for a's become a_{21}, a_{22})

18

Inner versus Outer Product

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad \mathbf{x}^\top = [x_1, x_2, \dots, x_n]$$

Inner (dot) product

$$\mathbf{x}^\top \mathbf{x} = \sum_{n=1}^N x_n^2$$

Outer product

$$\mathbf{x} \mathbf{x}^\top = \begin{bmatrix} x_1^2 & x_1 x_2 & \cdots & x_1 x_n \\ x_2 x_1 & x_2^2 & \cdots & x_2 x_n \\ \vdots & \vdots & \ddots & \vdots \\ x_n x_1 & x_n x_2 & \cdots & x_n^2 \end{bmatrix}$$

19

Simple Linear Model in Matrix Notation

- First, express our original 1-variable, 2-param model in matrix notation:

$$\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} \quad \mathbf{x}_n = \begin{bmatrix} 1 \\ x_n \end{bmatrix}$$

$$f(x_n; w_0, w_1) = \mathbf{w}^\top \mathbf{x}_n = w_0 + w_1 x_n$$

$$\mathcal{L} = \frac{1}{N} \sum_{n=1}^N (t_n - \mathbf{w}^\top \mathbf{x}_n)^2$$

20

Simple Linear Model in

$$\mathcal{L} = \frac{1}{N} \sum_{n=1}^N (t_n - \mathbf{w}^\top \mathbf{x}_n)^2$$

Matrix Notation

- Next, express the operations involving all of the data (the inputs \mathbf{x}_n and the targets t_n):

$$\begin{aligned}
 \mathbf{w} &= \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}, \quad \mathbf{x}_n = \begin{bmatrix} 1 \\ x_n \end{bmatrix} & \mathbf{t} - \mathbf{X}\mathbf{w} &= \begin{bmatrix} t_1 - w_0 - w_1 x_1 \\ t_2 - w_0 - w_1 x_2 \\ \vdots \\ t_N - w_0 - w_1 x_N \end{bmatrix} \\
 \mathbf{X} = \text{Design Matrix} &= \begin{bmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_N^\top \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_N \end{bmatrix}, \quad \mathbf{t} = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{bmatrix} & (\mathbf{t} - \mathbf{X}\mathbf{w})^\top (\mathbf{t} - \mathbf{X}\mathbf{w}) &= (t_1 - (w_0 + w_1 x_1))^2 + (t_2 - (w_0 + w_1 x_2))^2 + \dots \\
 &&&+ (t_N - (w_0 + w_1 x_N))^2 \\
 \mathbf{X}\mathbf{w} &= \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_N \end{bmatrix} \times \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} = \begin{bmatrix} w_0 + w_1 x_1 \\ w_0 + w_1 x_2 \\ \vdots \\ w_0 + w_1 x_N \end{bmatrix} & &= \sum_{n=1}^N (t_n - (w_0 + w_1 x_n))^2 \\
 &&&= \sum_{n=1}^N (t_n - f(x_n; w_0, w_1))^2 \\
 \mathcal{L} &= \frac{1}{N} (\mathbf{t} - \mathbf{X}\mathbf{w})^\top (\mathbf{t} - \mathbf{X}\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N (t_n - \mathbf{w}^\top \mathbf{x}_n)^2 = \frac{1}{N} \sum_{n=1}^N (t_n - (w_0 + w_1 x_n))^2
 \end{aligned}$$

Much nicer! The $\mathbf{x}^\top \mathbf{y}$ operation allows us to drop the sums!

21

Simple Linear Model in

Matrix Notation

- Now that we have the matrix version of the loss function, “just” take derivative...

$$\begin{aligned}
 \mathcal{L} &= \frac{1}{N} (\mathbf{X}\mathbf{w} - \mathbf{t})^\top (\mathbf{X}\mathbf{w} - \mathbf{t}) \\
 &= \frac{1}{N} ((\mathbf{X}\mathbf{w})^\top - \mathbf{t}^\top)(\mathbf{X}\mathbf{w} - \mathbf{t}) \\
 \text{note: book accidentally drops the } 1/N \text{ here} &\quad \rightarrow = \frac{1}{N} (\mathbf{X}\mathbf{w})^\top \mathbf{X}\mathbf{w} - \frac{1}{N} \mathbf{t}^\top \mathbf{X}\mathbf{w} - \frac{1}{N} (\mathbf{X}\mathbf{w})^\top \mathbf{t} + \frac{1}{N} \mathbf{t}^\top \mathbf{t} \\
 &= \frac{1}{N} \mathbf{w}^\top \mathbf{X}^\top \mathbf{X}\mathbf{w} - \frac{2}{N} \mathbf{w}^\top \mathbf{X}^\top \mathbf{t} + \frac{1}{N} \mathbf{t}^\top \mathbf{t}
 \end{aligned}$$

$\mathbf{t}^\top \mathbf{X}\mathbf{w}$ and $\mathbf{w}^\top \mathbf{X}^\top \mathbf{t}$ are the transpose of one another
... and both products come out to be scalars (i.e., “1x1 matrices”, where
transpose of one is the same as the other), so the products
are the same and can be combined.

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial w_0} \\ \frac{\partial \mathcal{L}}{\partial w_1} \end{bmatrix}$$

$f(\mathbf{w})$	$\frac{\partial f}{\partial \mathbf{w}}$
$\mathbf{w}^\top \mathbf{x}$	\mathbf{x}
$\mathbf{x}^\top \mathbf{w}$	\mathbf{x}
$\mathbf{w}^\top \mathbf{w}$	$2\mathbf{w}$
$\mathbf{w}^\top \mathbf{C}\mathbf{w}$	$2\mathbf{C}\mathbf{w}$

Some useful identities

$$\begin{aligned}
 \frac{\partial \mathcal{L}}{\partial \mathbf{w}} &= \frac{2}{N} \mathbf{X}^\top \mathbf{X}\mathbf{w} - \frac{2}{N} \mathbf{X}^\top \mathbf{t} = 0 \\
 \mathbf{X}^\top \mathbf{X}\mathbf{w} &= \mathbf{X}^\top \mathbf{t}
 \end{aligned}$$

$$\mathbf{I}\mathbf{w} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t} \quad \begin{aligned}
 &\text{The matrix normal equation!} \\
 &\text{(guaranteed unique solution when} \\
 &\text{the } n \text{ column vectors of } \mathbf{X} \\
 &\text{are linearly independent)}
 \end{aligned}$$

But what does this mean??

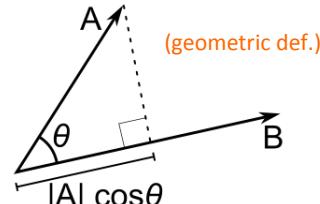
22

Relation of $\mathbf{a}^T \mathbf{b}$ to Geometry

- $\mathbf{a}^T \mathbf{b}$ is special (also $\mathbf{a} \cdot \mathbf{b}$) – the *dot product* (aka *scalar product*; the *inner product* for the Euclidean space)

$$\mathbf{a} \cdot \mathbf{b} = \sum_{i=1}^n a_i b_i = a_1 b_1 + a_2 b_2 + \cdots + a_n b_n \quad (\text{algebraic def.})$$

- Plays a role in defining
 - Euclidean distance (norm)
 - Angles



$$\mathbf{a} \cdot \mathbf{a} = \|\mathbf{a}\|^2$$

$$\mathbf{a} \cdot \mathbf{b} = \|\mathbf{a}\| \|\mathbf{b}\| \cos \theta$$

$$\theta = \arccos \left(\frac{\mathbf{a} \cdot \mathbf{b}}{\|\mathbf{a}\| \|\mathbf{b}\|} \right).$$

The dot product of vectors that are 90° (or more generally, orthogonal) is = 0

23

Geometry of Linear Systems and their Solution

A linear equation expresses a constraint between variables

A system of linear equations – more constraints!

$$2w_0 - w_1 = 0$$

$$-w_0 + 2w_1 = 3$$

$$u \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} = \begin{bmatrix} 0 \\ 3 \end{bmatrix}$$

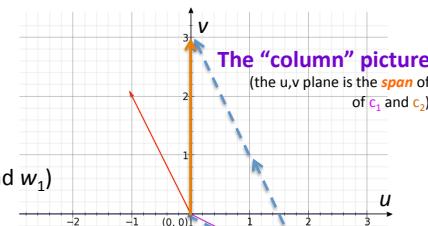
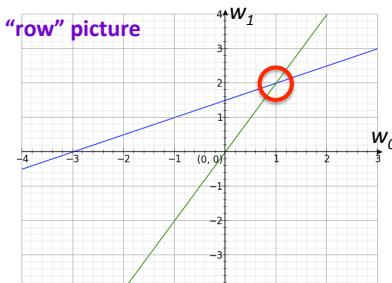
$$\mathbf{X} \mathbf{w} = \mathbf{t} \quad \mathbf{w} = \mathbf{X}^{-1} \mathbf{t}$$

“Solving” the linear system involves finding the linear combinations (i.e., the amounts w_0 and w_1) of \mathbf{c}_1 and \mathbf{c}_2 that equal the column vector $(0, 3)^T$.

Solve using your favorite method (e.g., Gaussian Elimination)

Here, the solution happens to be $w_0=1$ (\mathbf{c}_1), $w_1=2$ ($+ 2\mathbf{c}_2$)

The w_0 's and w_1 's corresponds to the point where the two lines cross!



24

Geometry of Linear Systems and their Solution

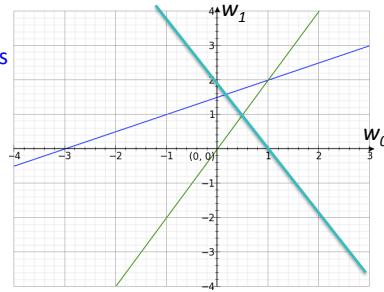
But what happens if we're **over-constrained**?

GOAL: Find a solution that is **closest** to (minimizes the distance between) the crossing points!

$$2w_0 - w_1 = 0$$

$$-w_0 + 2w_1 = 3$$

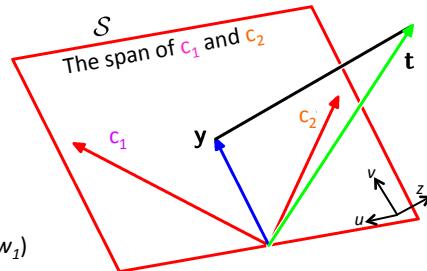
$$2w_0 + w_1 = 2$$



$$\begin{matrix} u & \begin{bmatrix} 2 & -1 \\ -1 & 2 \\ 2 & 1 \end{bmatrix} & \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} & = & \begin{bmatrix} 0 \\ 3 \\ 2 \end{bmatrix} \\ v & & & & \\ z & & & & \end{matrix}$$

$X \mathbf{w} = \mathbf{t}$

$$X \hat{\mathbf{w}} = \mathbf{y}$$



"Solving" the linear system involves finding the linear combinations (i.e., the amounts w_0 and w_1) of \mathbf{c}_1 and \mathbf{c}_2 that equal the column vector $(0, 3, 2)^T$.

Finding the projection

$X \mathbf{w} = \mathbf{t}$...what we started with, but no solution

$X \hat{\mathbf{w}} = \mathbf{y}$...something we can solve

In particular, we want \mathbf{y} to be the **closest** to \mathbf{t} but in the column space (the span) of \mathbf{c}_1 and \mathbf{c}_2 (which are the columns of X)

But we know the shortest distance between the end of \mathbf{t} and the span is a vector \mathbf{e} that is **orthogonal** (right angles!) to \mathbf{c}_1 and \mathbf{c}_2 (i.e., orthogonal to X)

That only happens when $\mathbf{e}^T \mathbf{X} = \mathbf{X}^T \mathbf{e} = 0$

Let's solve for when $\mathbf{X}^T \mathbf{e} = 0$

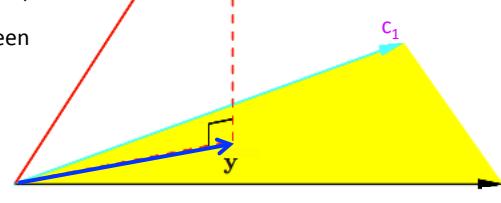
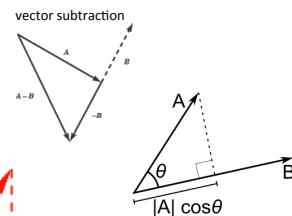
$$\mathbf{X}^T(\mathbf{t} - \mathbf{y}) = 0 \quad \text{plug in for } \mathbf{e}$$

$$\mathbf{X}^T(\mathbf{t} - \mathbf{X}\hat{\mathbf{w}}) = 0 \quad \text{substitute original } X\hat{w}=y \text{ (since we don't know } y)$$

$$\mathbf{X}^T\mathbf{t} - \mathbf{X}^T\mathbf{X}\hat{\mathbf{w}} = 0 \quad \text{multiply through...}$$

$$\mathbf{X}^T\mathbf{X}\hat{\mathbf{w}} = \mathbf{X}^T\mathbf{t}$$

$$\hat{\mathbf{w}} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{t} \quad \dots \text{ah hah!} \quad \mathbf{I}\mathbf{w} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{t}$$



26

Finding the projection

$\mathbf{X} \mathbf{w} = \mathbf{t}$...what we started with, but no solution
 $\mathbf{X} \hat{\mathbf{w}} = \mathbf{y}$...something we can solve

This is just what minimizing the squared loss did!

$$\mathcal{L} = \frac{1}{N} (\mathbf{t} - \mathbf{X}\mathbf{w})^T (\mathbf{t} - \mathbf{X}\mathbf{w})$$

This “error” vector \mathbf{e} is the shortest distance
 (Note: we can’t differentiate the “absolute value”
 so the squaring of the difference – of the *length*
 of \mathbf{e} – was a better choice for doing it the calc way)

That only happens when $\mathbf{e}^T \mathbf{X} = \mathbf{X}^T \mathbf{e} = 0$

Let’s solve for when $\mathbf{X}^T \mathbf{e} = 0$

$\mathbf{X}^T (\mathbf{t} - \mathbf{y}) = 0$	plug in for \mathbf{e}
$\mathbf{X}^T (\mathbf{t} - \mathbf{X}\hat{\mathbf{w}}) = 0$	substitute original $\mathbf{X}\hat{\mathbf{w}} = \mathbf{y}$ (since we don’t know \mathbf{y})
$\mathbf{X}^T \mathbf{t} - \mathbf{X}^T \mathbf{X}\hat{\mathbf{w}} = 0$	multiply through...
$\mathbf{X}^T \mathbf{X}\hat{\mathbf{w}} = \mathbf{X}^T \mathbf{t}$	
$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}$... ah hah! $\mathbf{I}\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}$

Parameters of hyperplane

$2w_0 - w_1 = 0$	$-w_0 + 2w_1 = 3$	$2w_0 + w_1 = 2$
inputs		targets

27