### Linear Predictors Part 1

#### COMPSCI 371D — Machine Learning

э

< ロ > < 同 > < 回 > < 回 > < 回 > <

# Outline

- 1 Definitions and Properties
- 2 The Least-Squares Linear Regressor
- 3 The Logistic-Regression Classifier
- Probabilities and the Geometry of Logistic Regression

< 🗇 🕨

A B > A B >

### Definitions

• A linear *regressor* fits an affine function to the data

 $y \approx h(\mathbf{x}) = b + \mathbf{w}^T \mathbf{x}$  for  $\mathbf{x} \in \mathbb{R}^d$  and  $y \in \mathbb{R}$ 

- A linear, binary *classifier* separates the two classes in  $Y = \{c_0, c_1\}$  with a hyperplane in  $\mathbb{R}^d$
- The actual data can be separated only if it is linearly separable (!)
- Multi-class linear classifiers separate any two classes with a hyperplane
- The resulting decision regions are convex and simply connected (polyhedra)

< 同 > < 回 > < 回 > -

# **Properties of Linear Predictors**

- Linear Predictors...
  - ...have a very small  $\mathcal{H}$  with d + 1 parameters (resist overfitting)
  - ... are trained by solving a convex optimization problem (global optimum)
  - ... are fast at inference time (and training is not too slow)
  - ... work well if the data is close to linearly separable

< 同 > < 三 > < 三 >

### The Least-Squares Linear Regressor

- Déjà vu: Polynomial regression with k = 1
  y ≈ h<sub>v</sub>(x) = b + w<sup>T</sup>x for x ∈ ℝ<sup>d</sup>
- Parameter vector  $\mathbf{v} = \begin{bmatrix} b \\ \mathbf{w} \end{bmatrix} \in \mathbb{R}^{d+1}$  $\mathcal{H}$  isomorphic to  $\mathbb{R}^m$  with m = d + 1
- "Least Squares:"  $\ell(y, \hat{y}) = (y \hat{y})^2$
- $\hat{\mathbf{v}} = \operatorname{arg\,min}_{\mathbf{v} \in \mathbb{R}^m} L_T(\mathbf{v})$
- Risk  $L_T(\mathbf{v}) = \frac{1}{N} \sum_{n=1}^N \ell(y_n, h_{\mathbf{v}}(\mathbf{x}_n))$
- We know how to solve this

・ 同 ト ・ ヨ ト ・ ヨ ト ・

### Linear Regression Example



- Left: All of Ames. Residual  $\sqrt{\text{Risk}}$ : \$55,800
- Right: One Neighborhood. Residual  $\sqrt{\text{Risk}}$ : \$23,600
- Left, yellow: Ignore two largest homes

## Binary Classification by Logistic Regression

 $\textbf{\textit{Y}} = \{\textbf{\textit{c}}_0, \textbf{\textit{c}}_1\}$ 

- Multi-class case later
- The logistic-regression classifier is a classifier!
- A linear classifier implemented through regression
- The *logistic* is a particular function

### **Score-Based Classifiers**

 $Y=\{\textit{c}_0,\textit{c}_1\}$ 

- Think of  $c_0$ ,  $c_1$  as numbers:  $Y = \{0, 1\}$
- We saw the idea of level sets: Regress a *score* function *s* such that *s* is large where *y* = 1, small where *y* = 0
- Threshold *s* to obtain a classifier:  $h(\mathbf{x}) = \begin{cases} c_0 & \text{if } s(\mathbf{x}) \leq \text{threshold} \\ c_1 & \text{otherwise.} \end{cases}$
- A linear classifier implemented through regression

・ 同 ト ・ ヨ ト ・ ヨ ト …

Idea 1

•  $s(\mathbf{x}) = b + \mathbf{w}^T \mathbf{x}$ 



- Not so good!
- A line does not approximate a step well
- Why not fit a step function?
- NP-hard unless the data is separable

< 🗇 🕨

# Idea 2

- How about a "soft step?"
- The logistic function



$$f(x) \stackrel{\mathsf{def}}{=} rac{1}{1+e^{-x}}$$

- If a true step moves, the risk does not change until a data point flips label
- If the logistic function moves, the risk changes gradually
- We have a nonzero gradient!
- The optimization problem is no longer combinatorial

# What is a Logistic Function in *d* Dimensions?

- We want a *linear* classifier
- The level crossing must be a hyperplane
- Level crossing: Solution to  $s(\mathbf{x}) = 1/2$
- Shape of the crossing depends on s
- Compose an affine  $a(\mathbf{x}) = c + \mathbf{u}^T \mathbf{x}$   $(a : \mathbb{R}^d \to \mathbb{R})$ ...with a monotonic f(a) that crosses 1/2  $(f : \mathbb{R} \to \mathbb{R})$  $s(\mathbf{x}) = f(a(\mathbf{x})) = f(c + \mathbf{u}^T \mathbf{x})$
- Then, if f(α) = 1/2, the equation s(x) = 1/2 is the same as c + u<sup>T</sup>x = α
- A hyperplane!
- Let f be the logistic function

э

COMPSCI 371D — Machine Learning

Ξ.

<ロ> <同> <ヨ> <ヨ>

# Example



- Gold line: Regression problem  $\mathbb{R} \to \mathbb{R}$
- Black line: Classification problem ℝ<sup>2</sup> → ℝ (result of running a logistic-regression classifier)
- Labels: Good (red squares, y = 1) or poor quality (blue circles, y = 0) homes
- All that matters is how far a point is from the black line

## A Probabilistic Interpretation



- All that matters is how far a point is from the black line
- $s(\mathbf{x}) = f(\Delta(\mathbf{x}))$  where  $\Delta$  is a *signed* distance
- We could interpret the score  $s(\mathbf{x})$  as "the probability that y = 1:"  $f(\Delta(\mathbf{x})) = \mathbb{P}[y = 1]$

## Ingredients for the Regression Part

- Determine the distance Δ of a point **x** ∈ X from a hyperplane χ, and the side of χ on which the point is on (Geometry: *affine functions* as unscaled, signed distances)
- Specify a monotonically increasing function that turns Δ into a probability (Choice based on convenience: the *logistic function*)
- Define a loss function  $\ell(y, \hat{y})$  such that the minimum risk yields the optimal classifier (Ditto, matches function in previous bullet to obtain a *convex* risk: the *cross-entropy loss*)

# Normal to a Hyperplane

• Hyperplane  $\chi$ :  $b + \mathbf{w}^T \mathbf{x} = 0$  (w.l.o.g.  $b \le 0$ )

 $\mathbf{a}_1, \mathbf{a}_2 \in \chi \Rightarrow \mathbf{c} = \mathbf{a}_1 - \mathbf{a}_2$  parallel to  $\chi$ 

- Subtract  $b + \mathbf{w}^T \mathbf{a}_1 = 0$  from  $b + \mathbf{w}^T \mathbf{a}_2 = 0$
- Obtain  $\mathbf{w}^T \mathbf{c} = \mathbf{0}$  for any  $\mathbf{a}_1, \mathbf{a}_2 \in \chi$
- w is perpendicular to  $\chi$

-

・ロト ・ 同ト ・ ヨト ・ ヨト

# Distance of a Hyperplane from the Origin



- Unit-norm version of **w**:  $\mathbf{n} = \frac{\mathbf{w}}{\|\mathbf{w}\|}$
- Rewrite  $\chi$ :  $b + \mathbf{w}^T \mathbf{x} = 0$  (w.l.o.g.  $b \le 0$ ) as  $\mathbf{n}^T \mathbf{x} = \beta$  where  $\beta = -\frac{b}{\|\mathbf{w}\|} \ge 0$
- Line along **n**:  $\mathbf{x} = \alpha \mathbf{n}$  for  $\alpha \in \mathbb{R}$  (parametric form)  $\alpha$  is the signed distance from the origin
- Replace into eq. for  $\chi$ :  $\alpha \mathbf{n}^T \mathbf{n} = \beta$  that is,  $\alpha = \beta \ge \mathbf{0}$
- In particular,  $\mathbf{x}_0 = \beta \mathbf{n}$
- $\beta$  is the distance of  $\chi$  from the origin

### Signed Distance of a Point from a Hyperplane



$$\mathbf{n}^T \mathbf{x} = \beta$$
 where  $\beta = -\frac{b}{\|\mathbf{w}\|} \ge 0$  and  $\mathbf{n} = \frac{\mathbf{w}}{\|\mathbf{w}\|} \mathbf{x}_0 = \beta \mathbf{n}$ 

- In one half-space,  $\mathbf{n}^T \mathbf{x} \ge \beta$
- Distance of **x** from  $\chi$  is  $\mathbf{n}^T \mathbf{x} \beta \ge \mathbf{0}$
- In other half-space,  $\mathbf{n}^T \mathbf{x}' \leq \beta$
- Distance of  $\mathbf{x}'$  from  $\chi$  is  $\beta \mathbf{n}^T \mathbf{x}' \ge \mathbf{0}$
- On decision boundary,  $\mathbf{n}^T \mathbf{x} = \beta$
- Δ(x) <sup>def</sup> = n<sup>T</sup>x − β is the signed distance of x from the hyperplane

COMPSCI 371D — Machine Learning

### Summary

If **w** is nonzero (which it has to be), the distance from the origin of the hyperplane  $\chi$  with equation  $b + \mathbf{w}^T \mathbf{x} = 0$  is

$$\beta \stackrel{\mathsf{def}}{=} \frac{|\boldsymbol{b}|}{\|\mathbf{w}\|}$$

(a nonnegative number) and the quantity

$$\Delta(\mathbf{x}) \stackrel{\mathsf{def}}{=} rac{b + \mathbf{w}^{\mathsf{T}} \mathbf{x}}{\|\mathbf{w}\|}$$

is the *signed distance* of point  $\mathbf{x} \in X$  from hyperplane  $\chi$ . Specifically, the distance of  $\mathbf{x}$  from  $\chi$  is  $|\Delta(\mathbf{x})|$ , and  $\Delta(\mathbf{x})$  is nonnegative if and only if  $\mathbf{x}$  is on the side of  $\chi$  pointed to by  $\mathbf{w}$ . Let us call that side the *positive half-space* of  $\chi$ .