

Recap SGD (Stochastic Gradient Descent)

- key idea: take the derivative of 1 datapoint at a time and use that to update w
- Moving in the direction to minimize cost function - against the derivative, by the α step size
- Useful when you can't find an analytic solution

for $i = 1, 2, \dots, n$; (shuffle)

$$\vec{w} \leftarrow \vec{w} - \alpha \underbrace{(\vec{w} \cdot \vec{x}_i - y_i)}_{\text{derivative}} \vec{x}_i$$

↑
step size

if $|J(\vec{w}) - J(\vec{w}^{t+1})| < \epsilon$ ← check for convergence, ϵ is very small

Stop!

- Choosing step size alpha - if it's too small, convergence is too slow - if it's too large, we can overshoot the minimum and may fail to converge (or even diverge)

Pros and Cons

Gradient Descent

- requires multiple iterations
- need to choose α
- works well when p is large

Normal Equations

- Non-iterative (computation - sometimes large)
- No need for α
- Slow if p is large

• Can support online learning - just add new data to the dataset w/out

• We have more control, but then have to make decisions

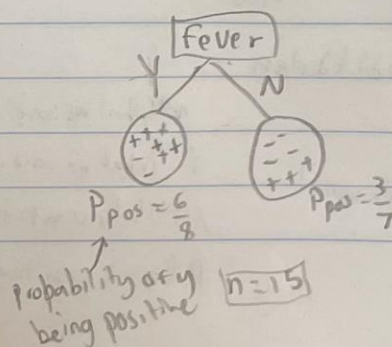
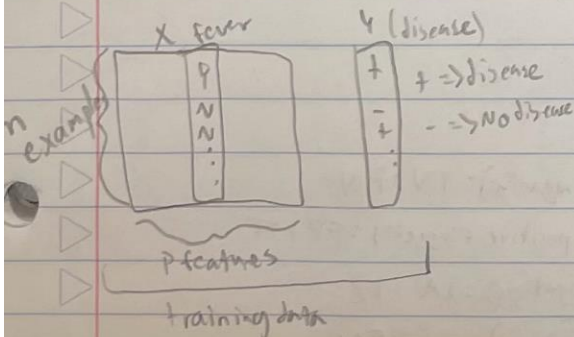
Introduction to Classification

Binary Classification - 2 classes

• ex: transactions that indicate credit card fraud, finding genes under natural selection ← in these, trying to find unusual items - we call these positives

- discrete classes data can be put into

decision tree - w/ a single feature is called a "stump"



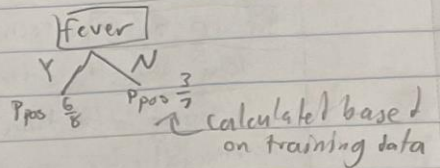
- Use probabilities to classify test examples

threshold $\text{ex } 0.5 \Rightarrow \hat{y}_{\text{test}} = \ominus$ (no disease)

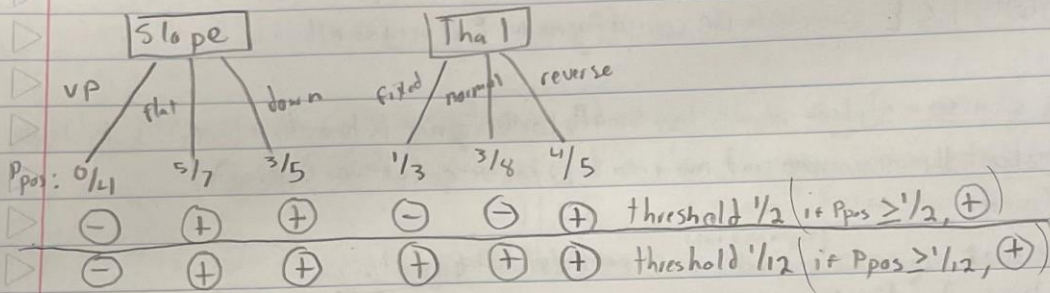
$0.25 \Rightarrow \hat{y}_{\text{test}} = \oplus$ (disease)

if $P_{\text{pos}} \geq \text{threshold}$, then classify example as positive

↳ used to decide how to classify data points



Handout 7



Evaluation Metrics

Goals of Evaluation:

- think about what metrics are important for problem at hand
- compare different methods or models on same problem (ex comparing cost using cost function)
- common set of tools that other users can understand

- Training and Testing - separate data into 'train' and 'test' (-n = num training ex, -m = num testing ex)

• Fit (create) model using training data

• evaluate model using testing data

predicted values

True values	-	+
-	65	15
+	7	13

thresh = 1/2

accuracy = $\frac{65+13}{100} = 78\%$

-	50	30
+	1	19

thresh = 0.25

accuracy = $\frac{50+19}{100} = 69\%$

accuracy = $\frac{\# \text{ correct}}{m}$

Confusion Matrices

Negative Positive

True Class	Negative	True Negative (TN)	False Positive (FP) ← 'false alarm'
	Positive	False Negative (FN) ← 'miss'	True Positive (TP)

N (what we said was negative) = $TN + FN$

P (what we said was positive "flagged") = $FP + TP$

N (total num true negatives) = $TN + FP$

P (total num true positives) = $TP + FN$

$$\text{Error} = \frac{(FN + FP)}{(TN + FP + FN + TP)} = \frac{FN + FP}{N + P}$$

$$\text{False Positive rate} = \frac{FP}{(TN + FP)} = \frac{FP}{N}$$

Accuracy: $1 - \text{Error}$

$$\text{Precision} = \frac{TP}{(FP + TP)} = \frac{TP}{P^*}$$

$$\text{Recall} = \frac{TP}{(FN + TP)} = \frac{TP}{P} \quad (\text{True positive rate})$$

Precision and Recall - Evaluation methods

- of all flagged exs, which ones are actually relevant (i.e. positive) - (purity)
- recall - of all relevant results, which ones did I actually return? (completeness)
- to increase recall - lower threshold to increase results classified as positive, which means higher chance we catch all TPs and reduce FN (but also increases FP which hurts precision)
- to increase precision - more selective in what we flag (this then affects recall as FN \uparrow)
- generally models cannot have both high precision ^{less FP} and high recall