



Classification Perceptron

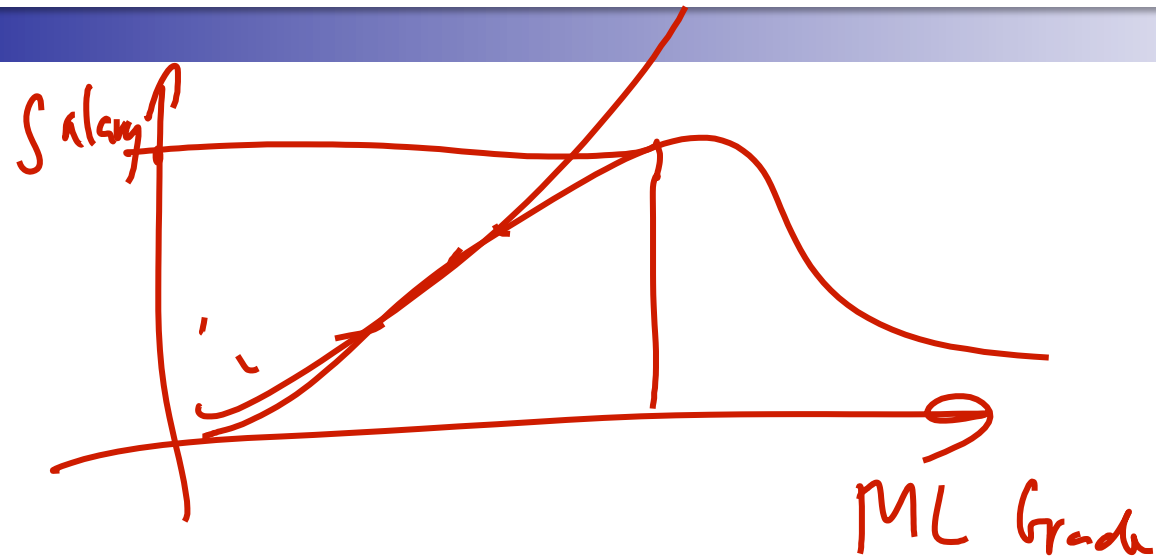
Machine Learning – CSEP546

Carlos Guestrin

University of Washington

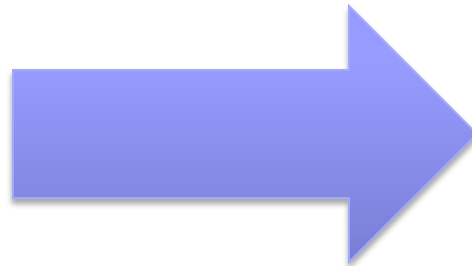
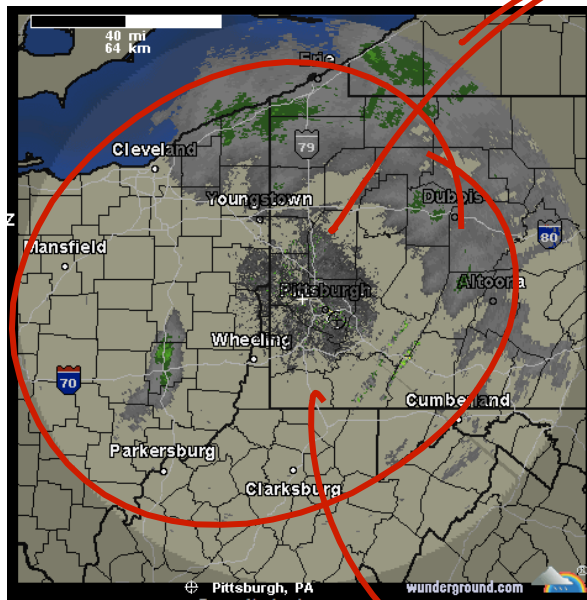
January 21, 2014

©Carlos Guestrin 2005-2014



**THUS FAR, REGRESSION:
PREDICT A CONTINUOUS
VALUE GIVEN SOME INPUTS**

Weather prediction revisited



Y
Set of classes
 $\{1, 2, 3, \dots, k\}$

Temperature

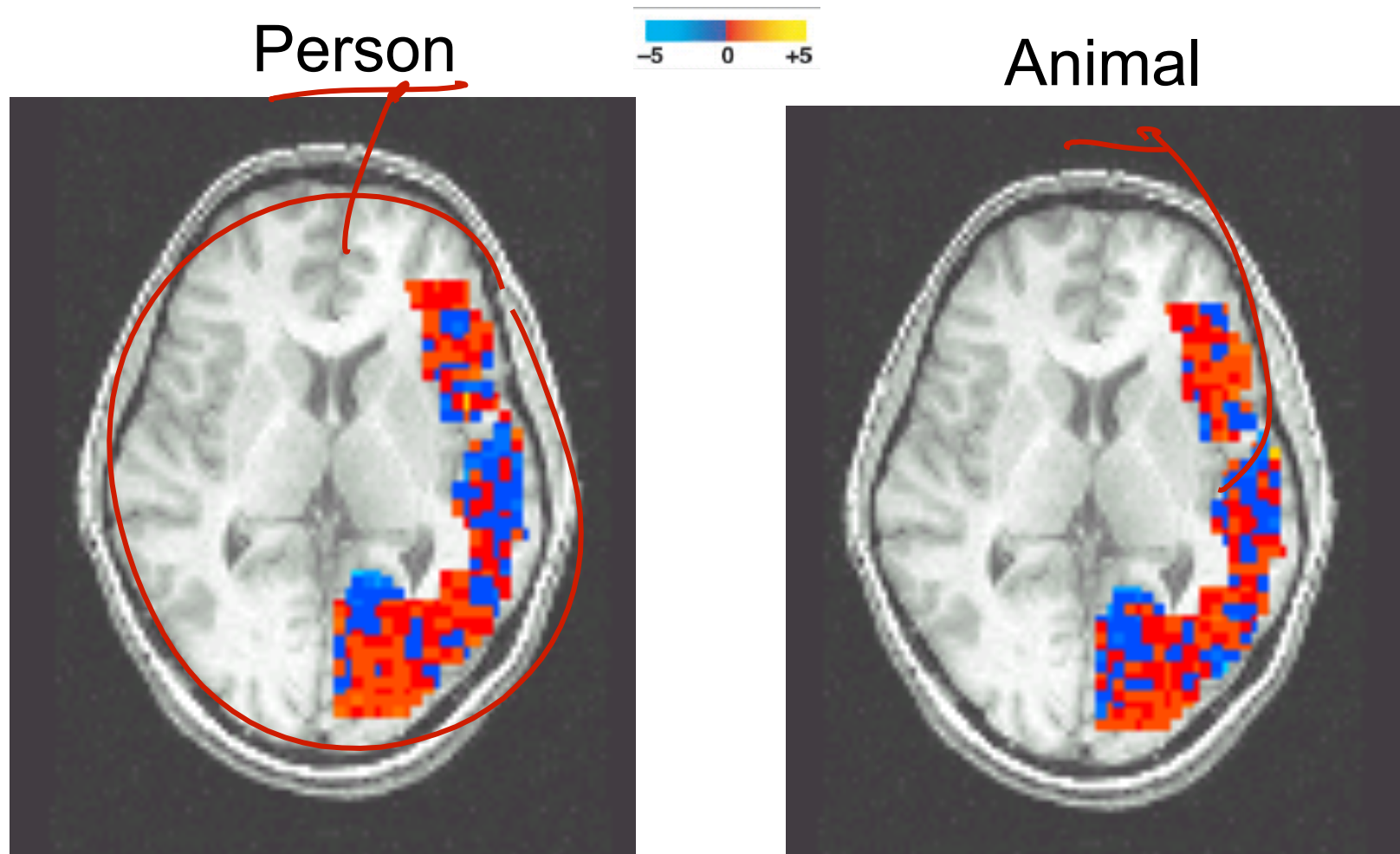
$52^{\circ}F$

$Y \subset \mathbb{R}$

Reading Your Brain, Simple Example

[Mitchell et al.]

Pairwise classification accuracy: 85%



Classification

■ Learn: $h: \mathbf{X} \mapsto Y$

- \mathbf{X} – features
- Y – target classes

$\mathbf{X} = (\text{GPA}, \text{grade}, \text{resume}, \dots)$

$Y = \{\text{hired}, \text{not hired}\}$

■ Simplest case: Thresholding

$\mathbf{X} = \text{Load Computer}$

$Y = \text{alarm?}$

$\text{Load}_{X_i} > 99\% \Rightarrow \text{alarm} = \text{true}$

$\text{else} \Rightarrow \text{alarm} = \text{false}$

$X_j > 27^\circ\text{C}$

i

Linear (Hyperplane) Decision Boundaries

$$w_0 + \sum_i w_i x_i \geq 0$$

$\forall x$

$$w_0 + \sum_i w_i x_i \geq 0$$

$$w_0 + \sum_i w_i x_i = 0$$

$$w_0 + \sum_i w_i x_i < 0$$

$x =$ text of email
sender
IP
:
not spam

linear
classification

spam

Learning a Linear Classifier

$$\text{sign}(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ -1 & \text{if } z < 0 \end{cases}$$

- Learn: $h: \mathbf{X} \mapsto Y$
 - \mathbf{X} – features
 - Y – target classes

Binary Classification
 $Y \in \{-1, +1\}$

- Decision rule:

$$w_0 + \sum_{i=1}^n w_i x_i \stackrel{?}{>} 0 : \hat{y} = \text{Sign}(w_0 + \sum_i w_i x_i)$$

Challenge: Data is streaming

- Assumption thus far: **Batch data**

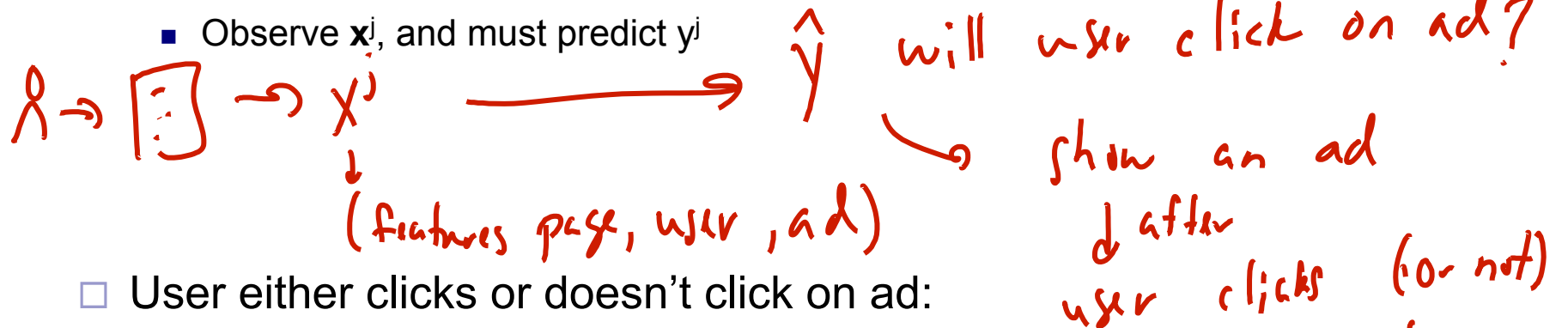
Have all the data now

Side note: streaming algorithms also work in batch case, just read one data point at time

- But, e.g., in click prediction for ads is a streaming data task:

- User enters query, and ad must be selected:

- Observe x^i , and must predict y_i



- User either clicks or doesn't click on ad:

- Label y_i is revealed afterwards

- Google gets a reward if user clicks on ad

Google "loses" if $\mathbb{I}(\hat{y} \neq y_i) = 1 \Rightarrow$ observe y_i

- Weights must be updated for next time:

$w^{(t+1)} \leftarrow w^{(t)} + \Delta$ what's Δ ? Google updates model params

$$w \cdot x = w_0 + \sum_i w_i x_i$$

Online Learning Problem

■ At each time step t :

□ Observe features of data point: $x^{(t)}$

- Note: many assumptions are possible, e.g., data is iid, data is adversarially chosen... details beyond scope of course

□ Make a prediction: $\hat{y}^{(t)} = \text{sign}(w_0^{(t)} + \sum_i w_i^{(t)} x_i^{(t)}) = \text{sign}(w^{(t)} \cdot x^{(t)})$

- Note: many models are possible, we focus on linear models
- For simplicity, use vector notation

$$w_0^{(t)} + \sum_i w_i^{(t)} x_i^{(t)} \geq 0 \Rightarrow \hat{y} = +1$$

else $\hat{y} = -1$

□ Observe true label:

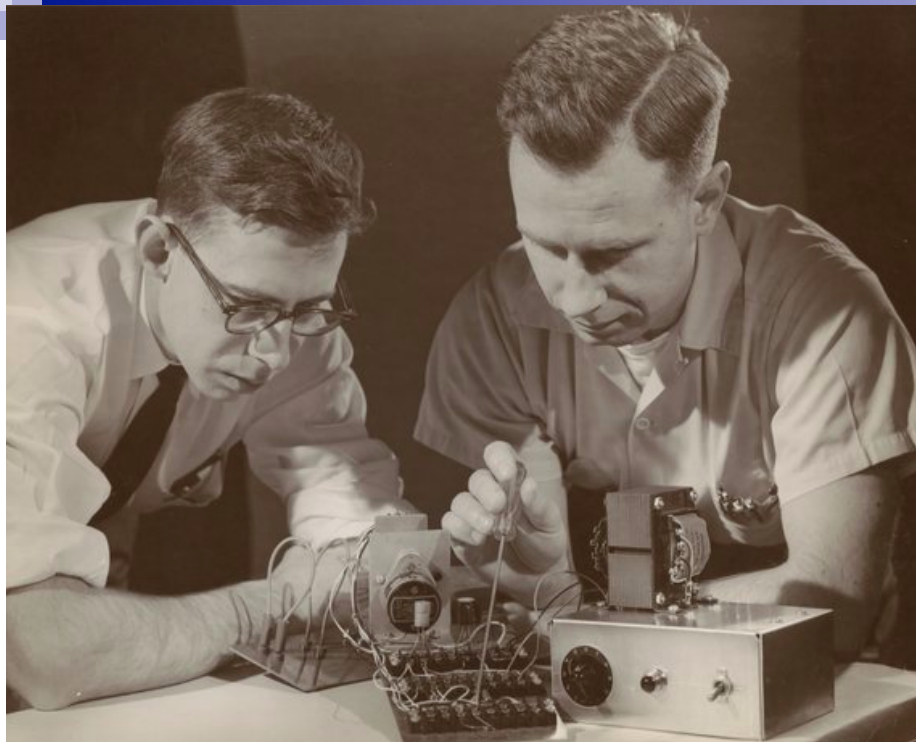
- Note: other observation models are possible, e.g., we don't observe the label directly, but only a noisy version... Details beyond scope of course

observe $y^{(t)} \rightarrow +1$ (clicked) | mistake if $\hat{y}^{(t)} \neq y^{(t)}$

$\rightarrow -1$ (not clicked)

□ Update model:

$$w^{(t+1)} \leftarrow w^{(t)} + \Delta^{(t)} \leftarrow \text{learn next}$$



Rosenblatt 1957

If x_i is categorical \Rightarrow turn into numeric value e.g. binary indicator

The Perceptron Algorithm

[Rosenblatt '58, '62]

- Classification setting: y in $\{-1, +1\}$

- Linear model

- Prediction:

$$\hat{y}^{(t)} = \text{sign}(w^{(t)} \cdot x^{(t)})$$

- Training:

- Initialize weight vector:

$w^{(0)}$ e.g., 0 or random

- At each time step:

- Observe features:

- Make prediction:

- Observe true class:

$$\hat{y}^{(t)} = \text{sign}(w^{(t)} \cdot x^{(t)})$$

$y^{(t)}$

- Update model:

- If prediction is not equal to truth

if $\hat{y}^{(t)} \neq y^{(t)}$

$w^{(t+1)} = w^{(t)}$

else

$w^{(t+1)} = w^{(t)} + y^{(t)} x^{(t)}$

mistake: $\hat{y}^{(t)} \neq y^{(t)}$

\Rightarrow add

$x^{(t)}$ to w
if $y^{(t)} = +1$

or

add

$-x^{(t)}$ to w

if $y^{(t)} = -1$

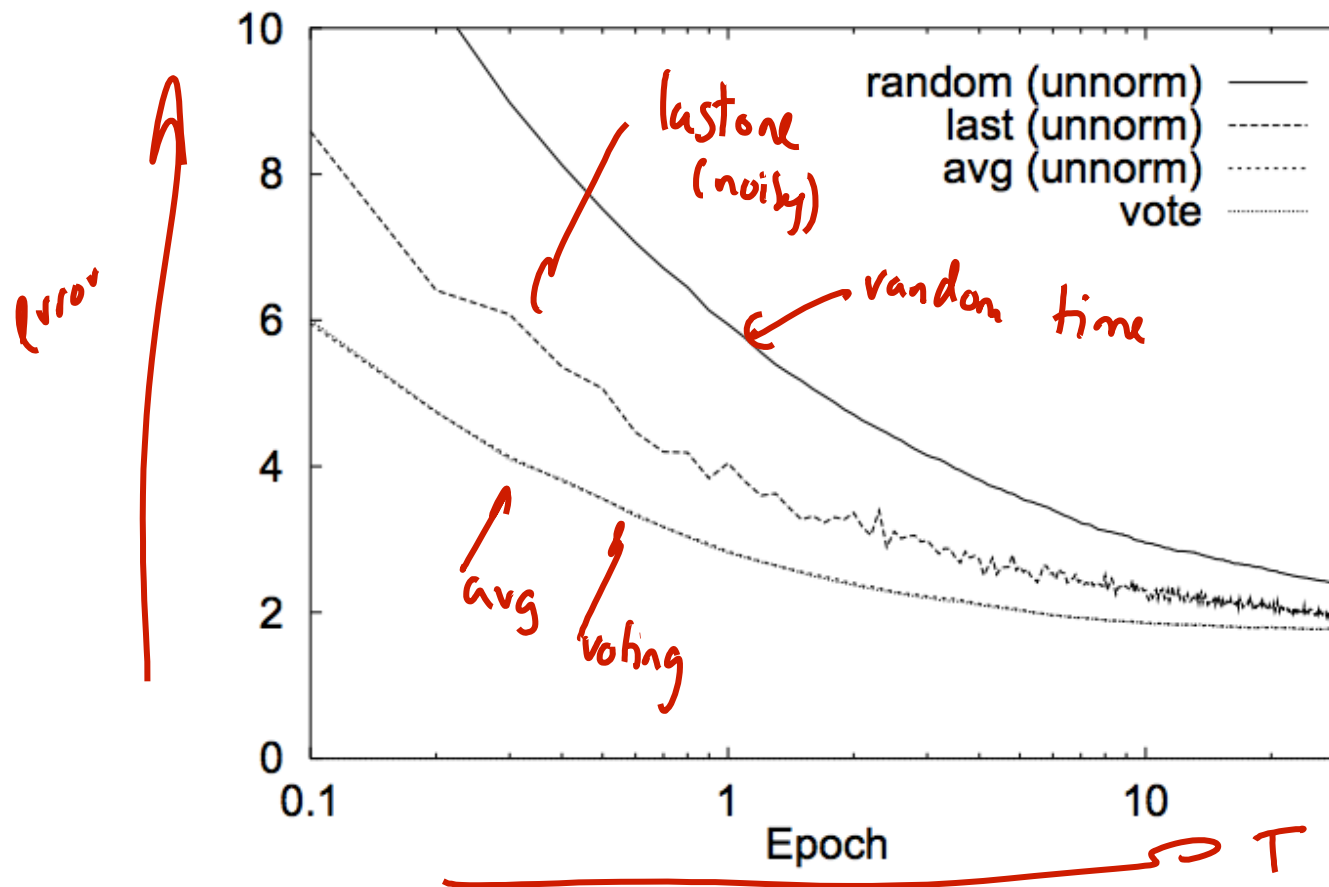
Fundamental Practical Problem for All Online Learning Methods: **Which weight vector to report?**

$w^{(0)}, w^{(1)}, \dots, w^{(T)}$

- Perceptron prediction:
- Suppose you run online learning method and want to sell your learned weight vector... Which one do you sell???

- Last one? $w^{(T)}? \leftarrow \text{very noisy}$
- Random time step? $\leftarrow \text{very noisy}$
- Average: $\hat{w} = \frac{1}{T} \sum_{t=1}^T w^{(t)} \leftarrow \text{use to show ads}$
- Voting see reading

Choice can make a huge difference!!



[Freund & Schapire '99]

Mistake Bounds

- Algorithm “pays” every time it makes a mistake:

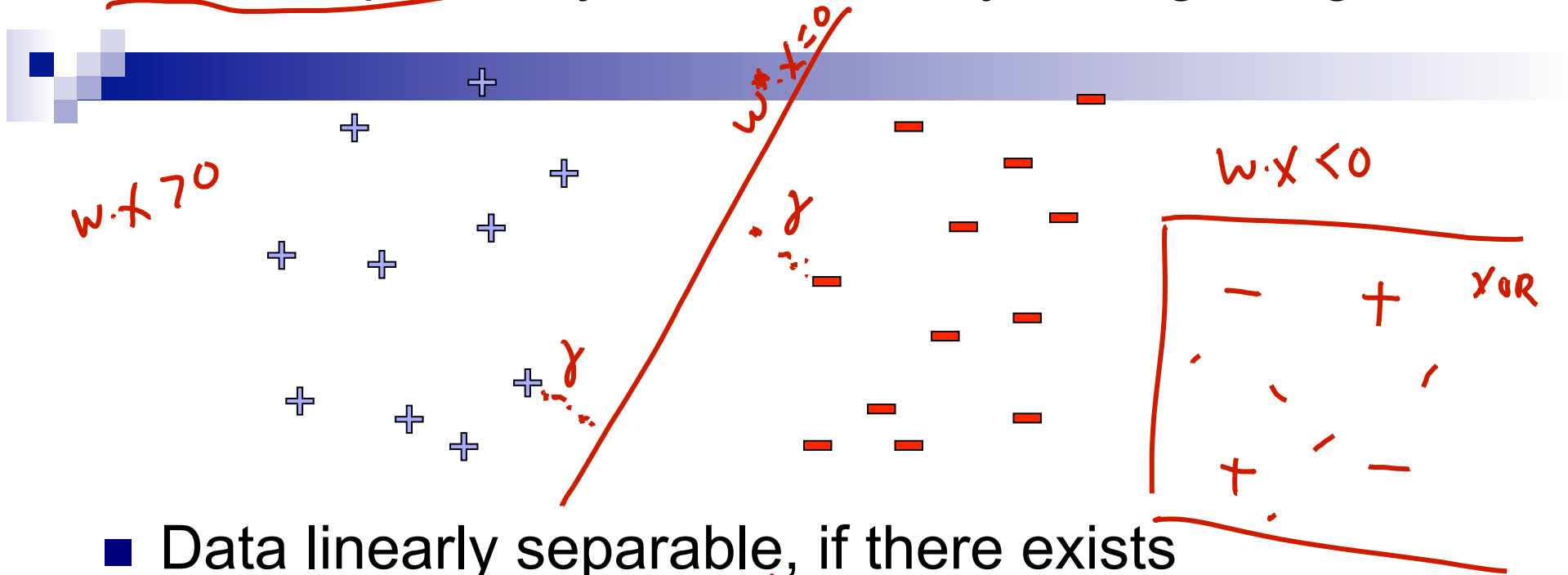
L, loss function is # mistakes

*↑
minimize*

- How many mistakes is it going to make?

mistake bound

Linear Separability: More formally, Using Margin



■ Data linearly separable, if there exists

□ a vector $\exists w^*, \|w^*\|=1$

□ a margin $\gamma \geq 0$

■ Such that

$$\forall x^{(t)} \quad \left. \begin{array}{l} w^* \cdot x^{(t)} \geq \gamma \text{ if } y^{(t)} = +1 \\ w^* \cdot x^{(t)} \leq -\gamma \text{ if } y^{(t)} = -1 \end{array} \right\} \Leftrightarrow y^{(t)} w^* \cdot x^{(t)} \geq \gamma$$

$$\|x\| = \sqrt{\sum_i x_i^2}$$

Perceptron Analysis: Linearly Separable Case

- Theorem [Block, Novikoff]:

- Given a sequence of labeled examples:

$$(x^{(1)}, y^{(1)}) \dots (x^{(2)}, y^{(2)}) \dots$$

- Each feature vector has bounded norm:

$$\forall t \quad \|x^{(t)}\| \leq R$$

- If dataset is linearly separable:

$$\exists w^*, \forall t \quad y^{(t)} w^* \cdot x^{(t)} \geq \gamma \quad \text{for some } \gamma > 0$$

- Then the number of mistakes made by the online perceptron on any such sequence is bounded by

$$\left(\frac{R}{\gamma}\right)^2$$

won!!

constant, doesn't depend on T

Mistake $y=1$ $w \cdot x < 0 \Rightarrow$ mistake $y w \cdot x < 0$ $a \cdot b \leq \|a\| \|b\|$



Perceptron Proof for Linearly Separable case

Assume $w^{(0)} = 0$

- Every time we make a mistake, we get gamma closer to w^* :

- Mistake at time t : $w^{(t+1)} = w^{(t)} + y^{(t)} x^{(t)}$

- Taking dot product with w^* :

- Thus after m mistakes:

by induction: $w^* \cdot w^{(t+1)} \geq m\gamma$ $\Rightarrow w^* \cdot w^{(t+1)} \geq w^* \cdot w^{(t)} + \gamma \geq \gamma$

$$\begin{aligned} w^* \cdot w^{(t+1)} &= w^* \cdot (w^{(t)} + y^{(t)} x^{(t)}) \\ &= w^* \cdot w^{(t)} + y^{(t)} w^* \cdot x^{(t)} \\ &\geq \gamma \end{aligned}$$

- Similarly, norm of $w^{(t+1)}$ doesn't grow too fast:

$$\|w^{(t+1)}\|^2 = \|w^{(t)}\|^2 + 2y^{(t)}(w^{(t)} \cdot x^{(t)}) + \|x^{(t)}\|^2 \leq R^2$$

new norm \quad previous norm \quad < 0 mistake

- Thus, after m mistakes:

$$\|w^{(t+1)}\|^2 \leq m R^2$$

$$\Rightarrow \|w^{(t+1)}\|^2 \leq \|w^{(t)}\|^2 + R^2$$

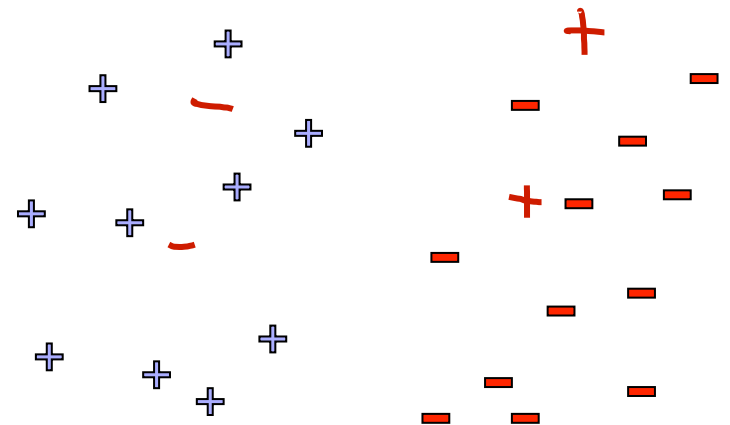
- Putting all together:

$$m\gamma \leq w^* \cdot w^{(t+1)} \leq \|w^*\| \|w^{(t+1)}\| \leq \sqrt{m} R$$

$$\Rightarrow m \leq \left(\frac{R}{\gamma}\right)^2 \text{ would be}$$

Beyond Linearly Separable Case

- Perceptron algorithm is super cool!
 - No assumption about data distribution!
 - Could be generated by an oblivious adversary, no need to be iid
 - Makes a fixed number of mistakes, and it's done for ever!
 - Even if you see infinite data
- However, real world not linearly separable
 - Can't expect never to make mistakes again
 - Analysis extends to non-linearly separable case
 - Very similar bound, see Freund & Schapire
 - Converges, but ultimately may not give good accuracy (make many many many mistakes)



What you need to know



- Notion of online learning
- Perceptron algorithm
- Mistake bounds and proof
- In online learning, report averaged weights at the end



Kernels

Machine Learning – CSEP546

Carlos Guestrin

University of Washington

January 21, 2014

©Carlos Guestrin 2005-2014

Summary Thus Far

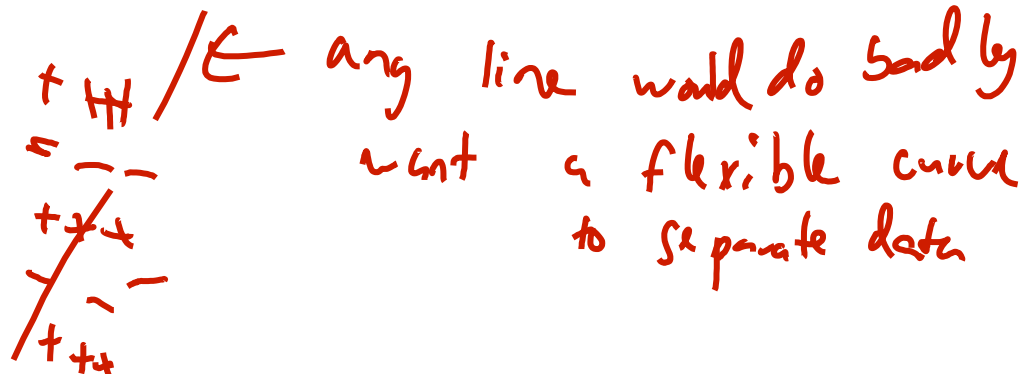
- Perceptron algorithm:

- Extremely simple classifier, works well in practice, *especially when you add regularization*
 - If you generalize it slightly by adding regularization → called a support vector machine (more next time)

- Constant number of mistakes in the linearly separable case

- More general results in the non-linearly separable case

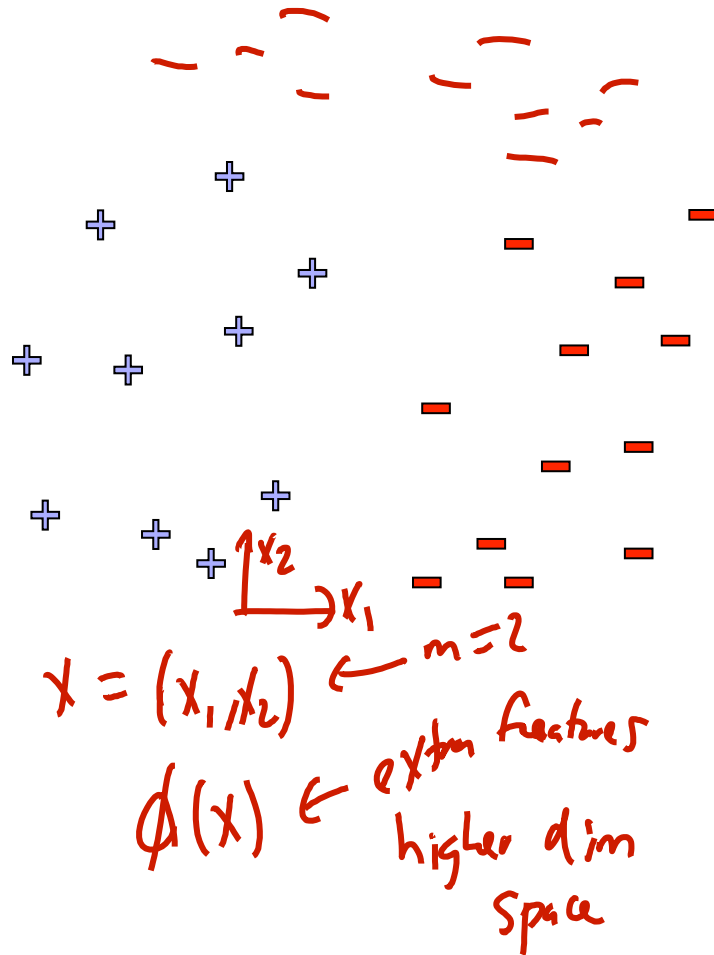
- In general, performance depends on how well we can separate the data



What if the data is not linearly separable?

Use features of features of features....

$$\Phi(\mathbf{x}) : \mathbb{R}^m \mapsto F \subseteq \mathbb{R}^D$$



$$\phi(x) = \begin{pmatrix} x_1 \\ x_2 \\ x_1 x_2 \\ x_1^2 \\ x_2^2 \\ e^{x_1} \\ \log x_2 \\ \vdots \end{pmatrix}$$

$\phi_i(x)$

learn: $\sum_i w_i \phi_i(x)$

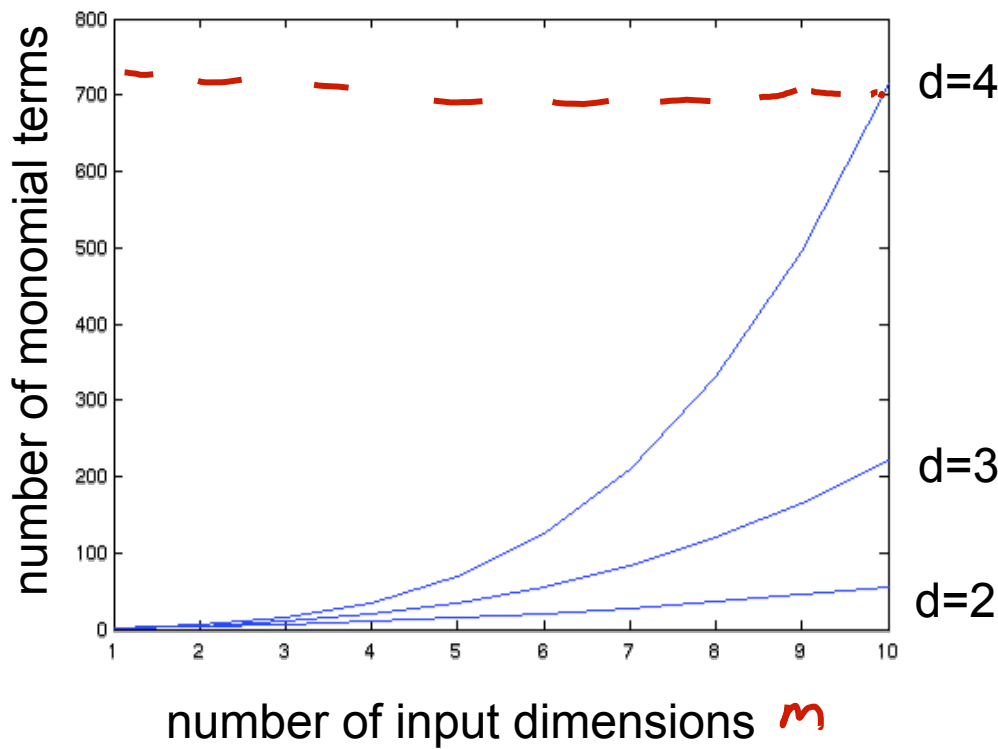
\uparrow D dim
 \uparrow a hyperplane in D dimensions

Feature space can get really large really quickly!

Higher order polynomials

$$\text{num. terms} = \binom{d + m - 1}{d} = \frac{(d + m - 1)!}{d!(m - 1)!}$$

①



m – input features
 d – degree of polynomial

Kernels:

eg. learn high degree poly
 \sim same cost as
 \sim degree 1

grows fast!

$d = 6, m = 100$

about 1.6 billion terms

Perceptron Revisited

- Given weight vector $w^{(t)}$, predict point x by:

$$\hat{y} = \text{sign}(w^{(t)} \cdot x)$$

- Mistake at time t : $w^{(t+1)} \leftarrow w^{(t)} + y^{(t)} x^{(t)}$

- Thus, write weight vector in terms of mistaken data points only: $w^{(0)} = 0$
 - Let $M^{(t)}$ be time steps up to t when mistakes were made:

$$w^{(t)} = \sum_{j \in M^{(t)}} y^{(j)} x^{(j)}$$

- Prediction rule now:

$$\text{sign}(w^{(t)} \cdot x^{(t)}) = \text{sign}\left(x^{(t)} \cdot \sum_{j \in M^{(t)}} y^{(j)} x^{(j)}\right) = \text{sign}\left(\sum_{j \in M^{(t)}} y^{(j)} x^{(t)} \cdot x^{(j)}\right)$$

- When using high dimensional features:

$$\text{sign}(w^{(t)} \cdot \phi(x^{(t)})) = \text{sign}\left(\sum_{j \in M^{(t)}} y^{(j)} \underbrace{\phi(x^{(t)}) \cdot \phi(x^{(j)})}_{\text{Kern: very fast}}\right)$$

Dot-product of polynomials

$$u = \begin{pmatrix} u_1 \\ u_2 \end{pmatrix}$$
$$v = \begin{pmatrix} v_1 \\ v_2 \end{pmatrix}$$

$\Phi(u) \cdot \Phi(v)$ = polynomials of degree exactly d

$$d=1 \quad \phi(u) \cdot \phi(v) = \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} \cdot \begin{pmatrix} v_1 \\ v_2 \end{pmatrix} = u_1 v_1 + u_2 v_2 = u \cdot v$$

$$d=2 \quad \phi(u) \cdot \phi(v) = \begin{pmatrix} u_1^2 \\ u_1 u_2 \\ u_2 u_1 \\ u_2^2 \end{pmatrix} \cdot \begin{pmatrix} v_1^2 \\ v_1 v_2 \\ v_2 v_1 \\ v_2^2 \end{pmatrix} = u_1^2 v_1^2 + 2 u_1 u_2 v_1 v_2 + u_2^2 v_2^2$$
$$= (u_1 v_1 + u_2 v_2)^2 = (u \cdot v)^2$$

"proof by one step of induction"

Polynomial of degree exactly d

$$\phi(u) \phi(v) = (u \cdot v)^d = K(u, v)$$

✓ "kernel trick" $\begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$

Finally the Kernel Trick!!!

(Kernelized Perceptron)

- Every time you make a mistake, remember $(\mathbf{x}^{(t)}, y^{(t)})$

↳ keep index $M^{(t)}$ into mistakes
keep mistakes as a list

- Kernelized Perceptron prediction for \mathbf{x} :

$$\underbrace{\text{sign}(\mathbf{w}^{(t)} \cdot \phi(\mathbf{x}))}_{\hat{y}} = \text{sign} \left(\sum_{j \in M^{(t)}} y^{(j)} \underbrace{\phi(\mathbf{x}^{(j)}) \cdot \phi(\mathbf{x})}_{\downarrow} \right)$$

remember all mistakes
& use kernels for dot product.

$$= \text{sign} \left(\sum_{j \in M^{(t)}} y^{(j)} k(\mathbf{x}^{(j)}, \mathbf{x}) \right)$$

Polynomial kernels

$$\phi(u) = \begin{pmatrix} 1 \\ u_1 \\ u_2 \\ u_1^2 \\ u_1 u_2 \\ \vdots \end{pmatrix}$$

- All monomials of degree d in $O(d)$ operations:

$$\Phi(u) \cdot \Phi(v) = (u \cdot v)^d = \text{polynomials of degree exactly } d$$

- How about all monomials of degree up to d ?

□ Solution 0: $\phi(u) \cdot \phi(v) = \sum_{i=0}^d \binom{d}{i} (u \cdot v)^i$

□ Better solution:

$$d=2: (u \cdot v)^0 + (u \cdot v)^1 + (v \cdot u)^1 + (u \cdot v)^2 = (u \cdot v + 1)^2$$

proof by "induction"

$$\phi(u) \cdot \phi(v) = (u \cdot v + 1)^d = \kappa(u, v) \quad !!$$

Common kernels

- Polynomials of degree exactly d

$$K(\mathbf{u}, \mathbf{v}) = (\mathbf{u} \cdot \mathbf{v})^d$$

- Polynomials of degree up to d

$$K(\mathbf{u}, \mathbf{v}) = (\mathbf{u} \cdot \mathbf{v} + 1)^d$$

- Gaussian (squared exponential) kernel

$$K(\mathbf{u}, \mathbf{v}) = \exp\left(-\frac{\|\mathbf{u} - \mathbf{v}\|^2}{2\sigma^2}\right)$$

- Sigmoid

$$K(\mathbf{u}, \mathbf{v}) = \tanh(\eta \mathbf{u} \cdot \mathbf{v} + \nu)$$

radial basis function

↓
infinite dimensional
 $\phi(x)$

param choose by cross validation

many ;

What you need to know



- Notion of online learning
- Perceptron algorithm
- Mistake bounds and proofs
- The kernel trick
- Kernelized Perceptron
- Derive polynomial kernel
- Common kernels
- In online learning, report averaged weights at the end



Naïve Bayes

Machine Learning – CSEP546

Carlos Guestrin

University of Washington

January 21, 2014

©Carlos Guestrin 2005-2014

Classification

- **Learn:** $h: \mathbf{X} \mapsto Y$

- \mathbf{X} – features
- Y – target classes

- Thus far: just a decision boundary

$$\hat{y} = \text{sign}(w \cdot x) \leftarrow \text{yes/no decision}$$

- What if you want probability of each class? $P(Y|X)$

$$\underline{P(Y = \text{spam} \mid x \in \text{text of email})}$$

$$\hat{y} = \underset{y}{\text{argmax}} P(Y = y \mid x \in \text{text of email})$$

Ad Placement Strategies

- Companies bid on ad prices

$$c_1 \rightarrow \$10$$

$$c_2 \rightarrow \$20$$

$$c_3 \rightarrow \$100$$

- Which ad wins? (many simplifications here)

□ Naively: $c_3 \rightarrow \$100$

□ But: paid on click only

□ Instead:

e.g.

$$p(\text{click} | c_3, \text{'big data'}) = 0.01 \Rightarrow E[\$] = 0.01 * \$100 = \$1$$

$$p(\text{click} | c_1, \text{'big data'}) = 0.5 \Rightarrow E[\$] = 0.5 * \$10 = \$5$$

The screenshot shows a Google search for "big data". The search bar at the top contains "big data" and a search button. Below the search bar, there are tabs for "Web", "Images", "Maps", "Shopping", "News", "More", and "Search tools". The search results are displayed in a grid. On the left, there are several search results from various sources like SAS, Dell, Oracle, and Wikipedia. On the right, there are several ads related to "big data", including "Big Data Cloud Analytics", "Big Data Monitoring", "New: Big Data in 2013", "Future Data Management", "NetApp@ Big Data", "PROS@ Big Data Research", "Extend Big Data With UJA", and "Big Data Solutions".

Key Task: Estimating Click Probabilities

- What is the probability that user i will click on ad j
- Not important just for ads:
 - ☐ Optimize search results
 - ☐ Suggest news articles
 - ☐ Recommend products
- Methods much more general, useful for:
 - ☐ Classification
 - ☐ Regression
 - ☐ Density estimation

Learning Problem for Click Prediction

- Prediction task: $X \rightarrow [0,1]$ $P(\text{click} = 1 | X)$
- Features:
 $X = (\text{facts ad}, \text{facts user}, \text{facts page})$
- Data: $(X, Y) \rightarrow (\text{webpage1}, \text{user17}, \text{ad27}, \text{time}=12) \rightarrow \text{click}=\text{true}$
 - Batch: fixed data set
 - Online: data stream
- Many approaches (e.g., logistic regression, SVMs, naïve Bayes, decision trees, boosting,...)
 - Focus on ~~naïve Bayes~~ and logistic regression; captures main concepts, ideas generalize to other approaches

Bayes Rule

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$

click → $P(Y|X)$
 page features → $P(X|Y)$
 likelihood fm → $P(X|Y)$
 prior → $P(Y)$

Which is shorthand for:

$$(\forall i, j) P(Y = y_i | X = x_j) = \frac{P(X = x_j | Y = y_i) P(Y = y_i)}{P(X = x_j)}$$

UWCSF → $P(X = x_j | Y = y_i)$
 click=1 → $P(Y = y_i)$
 normalization → $P(X = x_j)$
 normalizer ⇒ probs add up to 1

How hard is it to learn the optimal classifier?

	x_1 Gender	x_2 Age	x_3 Location	x_4 Income	x_5 Referrer	x_6 New or Returning	y Clicked?
	F	Young	US	High	Google	New	1
	M	Middle	US	Low	Direct	New	0
	F	Old	BR	Low	Google	Returning	1
	M	Young	BR	Low	Bing	Returning	0

- Data =

$$P(Y=1 | X_1, X_2, X_3, X_4, X_5, X_6)$$

$$= \frac{P(X_1 | Y) P(Y) / P(X)}{P(X)}$$

- How do we represent these? How many parameters?

- Prior, $P(Y)$:

- Suppose Y is composed of k classes

$$P(Y) = \begin{array}{c|ccc} Y & 1 & 2 & 3 \\ \hline P(Y) & .2 & .3 & .5 \end{array}$$

$k-1$ params

because probs. add up to 1

- Likelihood, $P(X|Y)$:

- Suppose X is composed of d binary features

$$P(X=x | Y=y)$$

For each Y : $2^d - 1$ params

total: $k(2^d - 1)$

$Y=1$

$d=2$

	$x_1=1, x_2=1$	TF	FT	FF
$P(X=x Y=1)$	0.1	0.3	0.2	0.4

2^d columns

- Complex model ! High variance with limited data!!!

Conditional Independence

independence.
 $X \perp Y$
 $P(X, Y) = P(X)P(Y)$

- X is **conditionally independent** of Y given Z, if the probability distribution governing X is independent of the value of Y, given the value of Z
 $(\forall i, j, k) P(X = i | Y = j, Z = k) = P(X = i | Z = k)$

→ ■ e.g., $P(\text{Thunder} | \text{Rain}, \text{Lightning}) = P(\text{Thunder} | \text{Lightning})$

$T \perp R | L$: Thunder independent of rain given lightning
 $(\Rightarrow) P(T, R | L) = P(T | L) P(R | L)$

- Equivalent to:

$$P(X, Y | Z) = P(X | Z) P(Y | Z)$$

What if features are independent?

- Predict Thunder $\leftarrow Y \in \{T, F\}$ $P(Y|x_1, x_2) = \frac{P(x_1, x_2 | Y) P(Y)}{P(x_1, x_2)}$
- From two **conditionally Independent** features

□ Lightning $\leftarrow x_1$

□ Rain $\leftarrow x_2$

No assumptions
estimate

$$P(x_1, x_2 | Y) = P(R, L | T)$$

$$\leftarrow 2(2^2 - 1) = 6 \text{ params}$$

But $R \perp L | T \Rightarrow P(R, L | T) = P(R | T) P(L | T)$

params: $2(2^1 - 1) = 2$

$2(2^1 - 1) = 2$

total 4 params

The Naïve Bayes assumption

- Naïve Bayes assumption:

- Features are independent given class:

$$\begin{aligned} \underline{P(X_1, X_2|Y)} &= P(X_1|X_2, Y)P(X_2|Y) \\ &= P(X_1|Y)P(X_2|Y) \end{aligned}$$

- More generally:

$$\boxed{P(X_1 \dots X_d | Y) = \prod_{i=1}^d P(X_i | Y)}$$

- How many parameters now?

- Suppose \mathbf{X} is composed of d binary features

with no assumptions:

$K(2^d - 1)$ params

now:
NB

$d \cdot K$ params

assumption \Rightarrow
more bias, less variance

exponential reduction in # of params!!

The Naïve Bayes Classifier

■ Given:

□ Prior $P(Y)$

$P(Y = \text{click})$

$P(Y = \text{not click})$

□ d conditionally independent features \mathbf{X} given the class Y

□ For each X_i , we have likelihood $P(X_i|Y) \leftarrow P(X_i = \text{'carbs'} | Y = \text{not click})$

■ Decision rule:

$$\hat{y}^* = h_{NB}(\mathbf{x}) = \arg \max_y P(y) P(x_1, \dots, x_d | y)$$

$$\arg \max_y P(Y=y|\mathbf{x}) = \arg \max_y \frac{P(y) P(\mathbf{x}|y)}{P(\mathbf{x})}$$

doesn't change decision

$$= \arg \max_y P(y) \prod_i P(x_i|y)$$

try for each value of y

predict y with highest product

■ If assumption holds, NB is optimal classifier!

$$P(X|y) = \frac{P(X,y)}{P(y)}$$

MLE for the parameters of NB

- Given dataset *N data point*
 - Count(A=a,B=b) == number of examples where A=a and B=b

- MLE for NB, simply:

- Prior: $P(Y=y) = \frac{\text{count}(Y=y)}{N}$

$$P(Y = \text{click})$$

$$P(Y = \text{hired})$$

- Likelihood: $P(X_i=x_i|Y=y) = \frac{\text{count}(X_i=x_i, Y=y)}{\text{count}(Y=y)}$

$$P(X_i = \text{'carbs'} | Y = \text{not clicked})$$

$$P(\text{GPA} = \text{'high'} | Y = \text{not hired})$$

$$\hookrightarrow = \frac{\text{count}(\text{GPA} = \text{high}, Y = \text{not hired})}{\text{count}(Y = \text{not hired})} = \frac{\text{count}(\text{GPA} = \text{high}, Y = \text{not hired})/N}{\text{count}(Y = \text{not hired})/N}$$

Subtleties of NB classifier 1 – Violating the NB assumption

- Usually, features are not conditionally independent:

$$P(X_1 \dots X_d | Y) \neq \prod_i P(X_i | Y)$$

- Actual probabilities $P(Y|\mathbf{X})$ often biased towards 0 or 1
- Nonetheless, NB is the single most used classifier out there
 - NB often performs well, even when assumption is violated
 - [Domingos & Pazzani '96] discuss some conditions for good performance

Subtleties of NB classifier 2 – Insufficient training data

- What if you never see a training instance where $X_1=a$ when $Y=b$?

$\text{Count}(a,b)$ \rightarrow \square e.g., $Y=\{\text{SpamEmail}\}$, $X_1=\{\text{'CSEP546'}\}$
 $\text{Count}(b)$ \rightarrow \square $P(X_1=a \mid Y=b) = 0$

$$P(Y|x) = P(Y) \prod_i P(x_i | Y)$$

if the is zero
always product
of stuff

- Thus, no matter what the values X_2, \dots, X_d take:
 - $\square P(Y=b \mid X_1=a, X_2, \dots, X_d) = 0$

$$\text{Smooth Counts}(x_i=x_i, Y=y) = \text{Count}(x_i, y) + \alpha \cdot \text{Uniform}(\#Y)$$

- “Solution”: smoothing

- \square Add “fake” counts, usually uniformly distributed
- \square Equivalent to “Bayesian Learning”

smoothing
parameter

pick by
cross validation

$$\frac{1}{k \mid Y_i \mid}$$

values x_i
takes

Text classification

- Classify e-mails

- $Y = \{\text{Spam}, \text{NotSpam}\}$

- Classify news articles

- $Y = \{\text{what is the topic of the article?}\}$

Sports
← *politics*
:

- Classify webpages

- $Y = \{\text{student, professor, project, ...}\}$

- What about the features **X**?

- The text!

what's X?

Features X are entire document – X_i for i^{th} word in article

Article from rec.sport.hockey

Path: cantaloupe.srv.cs.cmu.edu!das-news.harvard.e
From: xxx@yyy.zzz.edu (John Doe)
Subject: Re: This year's biggest and worst (opinion)
Date: 5 Apr 93 09:53:39 GMT

I can only comment on the Kings, but the most obvious candidate for pleasant surprise is Alex Zhitnik. He came highly touted as a defensive defenseman, but he's clearly much more than that. Great skater and hard shot (though wish he were more accurate). In fact, he pretty much allowed the Kings to trade away that huge defensive liability Paul Coffey. Kelly Hrudefy is only the biggest disappointment if you thought he was any good to begin with. But, at best, he's only a mediocre goaltender. A better choice would be Tomas Sandstrom, though not through any fault of his own, but because some thugs in Toronto decided

NB for Text classification

- $P(\mathbf{X}|Y)$ is huge!!!

- Article at least 1000 words, $\mathbf{X}=\{X_1, \dots, X_{1000}\}$
- X_i represents i^{th} word in document, i.e., the domain of X_i is entire vocabulary, e.g., Webster Dictionary (or more), 10,000 words, etc.

- NB assumption helps a lot!!!

- $P(X_i=x_i|Y=y)$ is just the probability of observing word x_i in a document on topic y

$$\underline{h_{NB}(\mathbf{x})} = \arg \max_y P(y) \prod_{i=1}^{LengthDoc} P(x_i|y)$$

Handwritten notes:

- Red arrow from "prior prob of Spam" points to $P(y)$
- Red arrow from "prob observing a particular word in Spam emails" points to $P(x_i|y)$

Bag of words model

- Typical additional assumption – **Position in document doesn't matter**: $P(X_i=x_i|Y=y) = P(X_k=x_i|Y=y)$
 - “Bag of words” model – order of words on the page ignored
 - Sounds really silly, but often works very well!

$$P(y) \prod_{i=1}^{LengthDoc} P(x_i|y)$$

When the lecture is over, remember to wake up the person sitting next to you in the lecture room.

Bag of words model

- Typical additional assumption – **Position in document doesn't matter**: $P(X_i=x_i|Y=y) = P(X_k=x_i|Y=y)$
 - “Bag of words” model – order of words on the page ignored
 - Sounds really silly, but often works very well!

$$P(y) \prod_{i=1}^{LengthDoc} P(x_i|y)$$

in is lecture lecture next over person remember room
sitting the the the to to up wake when you

Bag of Words Approach



Count
Q

N-grams
consecutive
sets of words
'University of
Washington'

aardvark	0
about	2
all	2
Africa	1
apple	0
anxious	0
...	
gas	1
...	
oil	1
...	
Zaire	0

NB with Bag of Words for text classification

■ Learning phase:

□ Prior $P(Y)$

- Count how many documents you have from each topic (+ prior)

□ $P(X_i|Y)$

- For each topic, count how many times you saw word in documents of this topic (+ prior)

■ Test phase:

□ For each document

- Use naïve Bayes decision rule

$$h_{NB}(\mathbf{x}) = \arg \max_y P(y) \prod_{i=1}^{LengthDoc} P(x_i|y)$$

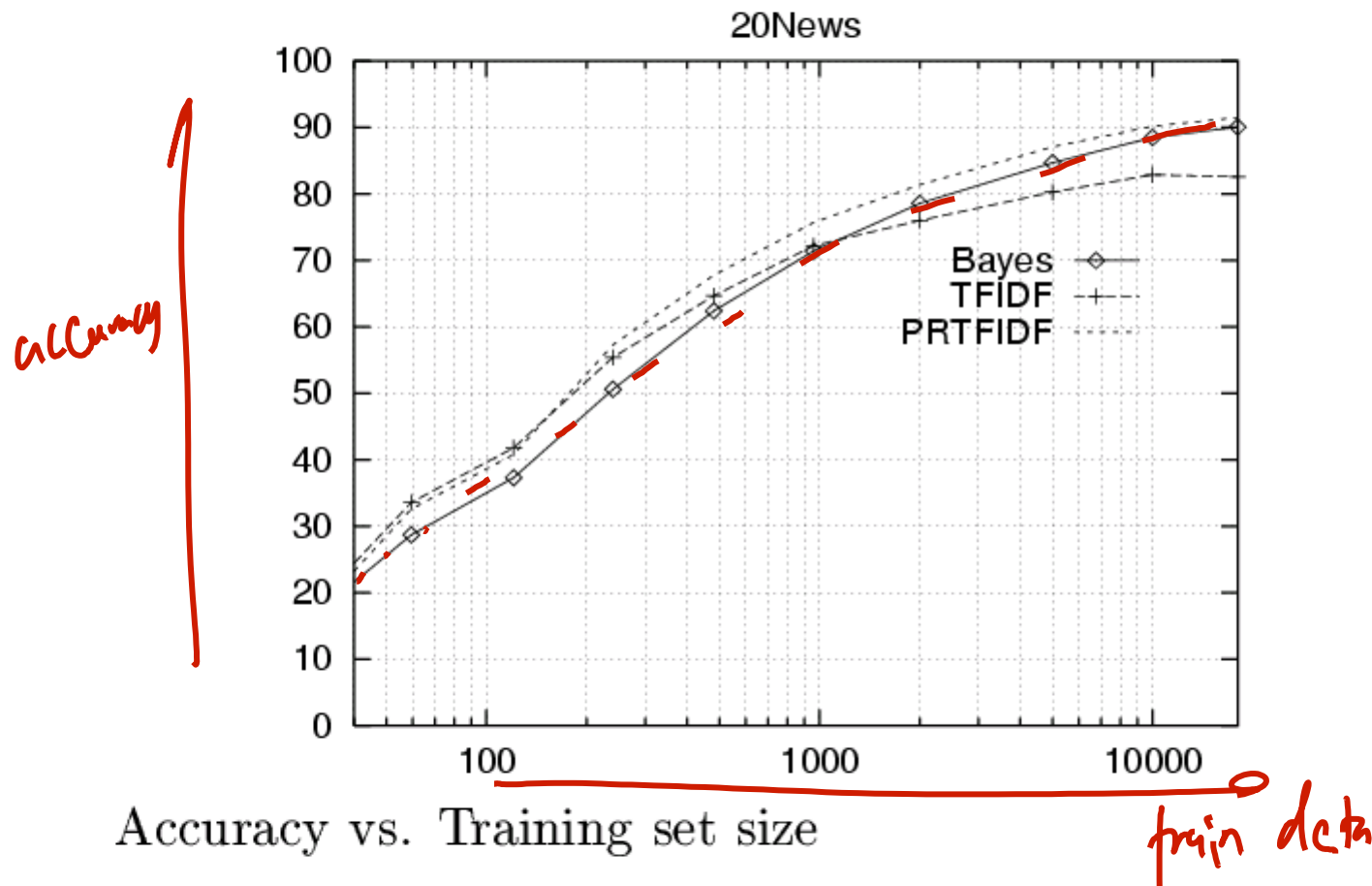
Twenty News Groups results

Given 1000 training documents from each group
Learn to classify new documents into
which newsgroup it came from

comp.graphics	misc.forsale
comp.os.ms-windows.misc	rec.autos
comp.sys.ibm.pc.hardware	rec.motorcycles
comp.sys.mac.hardware	rec.sport.baseball
comp.windows.x	rec.sport.hockey
alt.atheism	sci.space
soc.religion.christian	sci.crypt
talk.religion.misc	sci.electronics
talk.politics.mideast	sci.med
talk.politics.misc	
talk.politics.guns	

Naive Bayes: 89% classification accuracy

Learning curve for Twenty News Groups



What you need to know



- Click prediction problem
- Probabilities rather than classification
- Naïve Bayes model
 - Assumption
 - Formulation
- Application to text data
 - Bag of words model