# Discriminative classifiers: Logistic Regression, SVMs

CISC 5800

Professor Daniel Leeds

---

## Maximum A Posteriori: a quick review

- Likelihood: $P(D|\theta) = P(D|p) = p^{|H|}(1-p)^{|T|}$
- Prior: $P(\theta) = \frac{\theta^{\alpha-1}(1-\theta)^{\beta-1}}{B(\alpha,\beta)} = P(p) = \frac{p^{\alpha-1}(1-p)^{\beta-1}}{B(\alpha,\beta)}$
- Posterior  Likelihood x prior = $P(D|\theta)P(\theta)$

- MAP estimate:

$\underset{\theta}{\text{argmax}} \log P(D|\theta) + \log P(\theta)$

$\underset{p}{\text{argmax}} \log P(D|p) + \log P(p)$

$$p = \frac{|H| + (\alpha - 1)}{|H| + (\alpha - 1) + |T| + (\beta - 1)}$$

**Choose $\alpha$ and $\beta$ to give the prior belief of Heads bias $p \in [0,1]$**

**Higher $\alpha$: Heads more likely**
**Higher $\beta$: Tails more likely**

2

---

## Estimate each $P(X_i|Y)$ through MAP

Incorporating prior for each class $\beta_j$

$$P(X_i = x_k|Y = y_j) = \frac{\#D(X_i = x_k \wedge Y = y_j) + (\beta_j - 1)}{\#D(Y = y_j) + \sum_m(\beta_m - 1)}$$

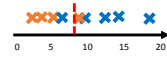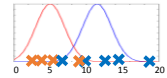$$P(Y = y_j) = \frac{\#D(Y = y_j) + (\beta_j - 1)}{|D| + \sum_m(\beta_m - 1)}$$

**$(\beta_j - 1)$ – "frequency" of class j**
**$\sum_m(\beta_m - 1)$ – "frequencies" of all classes**

**Note:** both X and Y can take on multiple values (binary **and** beyond)

3

---

## Classification strategy: generative vs. discriminative



- Generative, e.g., Bayes/Naïve Bayes:
  - Identify probability distribution for each class
  - Determine class with maximum probability for data example

- Discriminative, e.g., Logistic Regression:
  - Identify boundary between classes
  - Determine which side of boundary new data example exists on



4

---

## Linear algebra: data features

- Vector – list of numbers: each number describes a data **feature**

- Matrix – list of lists of numbers: features for each data point

|  | Document 1 | Document 2 | Document 3 |
|---|---|---|---|
| Wolf | 12 | 8 | 0 |
| Lion | 16 | 10 | 2 |
| Monkey | 14 | 11 | 1 |
| Broker | 0 | # of word | 14 |
| Analyst | 1 | occurrences | 10 |
| Dividend | 1 | 1 | 12 |
| ⋮ | ⋮ | ⋮ | ⋮ |

5

---

## Feature space

- Each data feature defines a dimension in space

|  | Document1 | Document2 | Document3 |
|---|---|---|---|
| Wolf | 12 | 8 | 0 |
| Lion | 16 | 10 | 2 |
| Monkey | 14 | 11 | 1 |
| Broker | 0 | 1 | 14 |
| Analyst | 1 | 0 | 10 |
| Dividend | 1 | 1 | 12 |
| ⋮ | ⋮ | ⋮ | ⋮ |



6

## The dot product

The dot product compares two vectors:

$\cdot\ \boldsymbol{a} = \begin{bmatrix} a_1 \\ \vdots \\ a_n \end{bmatrix}, \boldsymbol{b} = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix}$ $\qquad \boldsymbol{a} \cdot \boldsymbol{b} = \sum_{i=1}^{n} a_i b_i = \boldsymbol{a}^T \boldsymbol{b}$

$\begin{bmatrix} 5 \\ 10 \end{bmatrix} \cdot \begin{bmatrix} 10 \\ 10 \end{bmatrix} = 5 \times 10 + 10 \times 10$
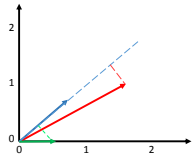
$= 50 + 100 = 150$



7

## The dot product, continued $\quad \boldsymbol{a} \cdot \boldsymbol{b} = \sum_{i=1}^{n} a_i b_i$

Magnitude of a vector is the sum of the squares of the elements
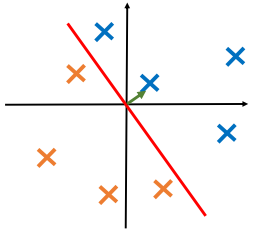
$|\boldsymbol{a}| = \sqrt{\sum_i a_i^2}$

If $\boldsymbol{a}$ has unit magnitude, $\boldsymbol{a} \cdot \boldsymbol{b}$ is the "projection" of $\boldsymbol{b}$ onto $\boldsymbol{a}$



$\begin{bmatrix} 0.71 \\ 0.71 \end{bmatrix} \cdot \begin{bmatrix} 1.5 \\ 1 \end{bmatrix} = .71 \times 1.5 + .71 \times 1$

$\approx 1.07 + .71 = 1.78$

$\begin{bmatrix} 0.71 \\ 0.71 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 0.5 \end{bmatrix} = .71 \times 0 + .71 \times 0.5$
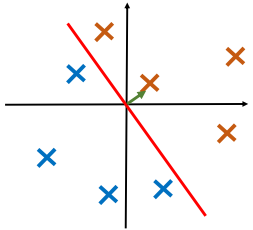
$\approx 0 + .35 = 0.35$  8

## Separating boundary, defined by w



- Separating **hyperplane** splits **class 0** and **class 1**

- Plane is defined by line **w** perpendicular to plan

- Is data point **x** in class 0 or class 1? **w$^T$x** > 0 class 0
  **w$^T$x** < 0 class 1   9

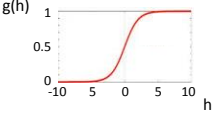## Separating boundary, defined by w

*More typically*



- Separating **hyperplane** splits **class 0** and **class 1**

- Plane is defined by line **w** perpendicular to plan

- Is data point **x** in class 0 or class 1? **w$^T$x** > 0 class **1**
  **w$^T$x** < 0 class **0**   10

## From real-number projection to 0/1 label

- Binary classification: 0 is class A, 1 is class B
- Sigmoid function stands in for p(x|y)

- Sigmoid: $g(h) = \frac{1}{1+e^{-h}}$

- $p(y = 0|x; \theta) = 1 - g(w^T x) = \frac{e^{-w^T x}}{1+e^{-w^T x}}$

- $p(y = 1|x; \theta) = g(w^T x) = \frac{1}{1+e^{-w^T x}}$



$w^T x = \sum_j w_j x_j$   11

## Learning parameters for classification

- Similar to MLE for Bayes classifier
- "Likelihood" for data points $y^1, ..., y^n$ (different from Bayesian likelihood)
  - If $y^i$ in class A, $y^i = 0$, multiply $(1 - g(x^i; w))$
  - If $y^i$ in class B, $y^i = 1$, multiply $(g(x^i; w))$

$$\underset{w}{\operatorname{argmax}} L(y|x; w) = \prod_i \left(1 - g(x^i; w)\right)^{(1-y^i)} g(x^i; w)^{y^i}$$

$$LL(y|x; w) = \sum_i (1 - y^i) \log\left(1 - g(x^i; w)\right) + y^i \log\left(g(x^i; w)\right)$$

$$LL(y|x; w) = \sum_i y^i \log \frac{g(x^i; w)}{1 - g(x^i; w)} + \log\left(1 - g(x^i; w)\right)$$   12

Learning parameters for classification

$$g(h) = \frac{1}{1 + e^{-h}}$$

$$LL(y|x;w) = \sum_i y^i \log \frac{g(x^i;w)}{1 - g(x^i;w)} + \log\left(1 - g(x^i;w)\right)$$

$$LL(y|x;w) = \sum_i y^i \log \frac{\frac{1}{1 + e^{-w^T x^i}}}{1 - \frac{1}{1 + e^{-w^T x^i}}} + \log\left(\frac{e^{-w^T x^i}}{1 + e^{-w^T x^i}}\right)$$

$$LL(y|x;w) = \sum_i y^i \log \frac{1}{1 + e^{-w^T x^i} - 1} + \log\left(\frac{e^{-w^T x^i}}{1 + e^{-w^T x^i}}\right)$$

$$LL(y|x;w) = \sum_i y^i w^T x^i - w^T x^i - \log\left(1 + e^{-w^T x^i}\right)$$

13

---

Learning parameters for classification

$$w^T x = \sum_j w_j x_j$$

$$g'(h) = \frac{e^{-h}}{(1 + e^{-h})^2}$$

$$LL(y|x;w) = \sum_i y^i w^T x^i - w^T x^i + \log\left(g(w^T x^i)\right)$$

$$\frac{\partial}{\partial w_j} LL(y|x;w) = \sum_i y^i x_j^i - x_j^i + \frac{x_j^i e^{-w^T x^i}}{1 + e^{-w^T x^i}}$$

$$\frac{\partial}{\partial w_j} LL(y|x;w) = \sum_i x_j^i\left(y^i - \left(1 - \left(1 - g(w^T x^i)\right)\right)\right)$$

$$\frac{\partial}{\partial w_j} LL(y|x;w) = \sum_i x_j^i\left(y^i - g(w^T x^i)\right)$$

14

---

Iterative gradient asscent

$y^i$ – true data label
$g(w^T x^i)$ – computed data label

- Begin with initial guessed weights w
- For each data point $(y^i, x^i)$, update each weight $w_j$

$$w_j \leftarrow w_j + \varepsilon x_j^i\left(y^i - g(w^T x^i)\right)$$

- Choose $\varepsilon$ so change is not too big or too small – "**step size**"

**Intuition**

- $x_j^i\left(y^i - g(w^T x^i)\right)$
  - If $y^i$=1 and $g(w^T x^i)$=0, and $x_j^i$>0, make $w_j$ larger and push $w^T x^i$ to be larger
  - If $y^i$=0 and $g(w^T x^i)$=1, and $x_j^i$>0, make $w_j$ smaller and push $w^T x^i$ to be smaller

15

---

Separating boundary, defined by w



- Separating **hyperplane** splits **class 0** and **class 1**

- Plane is defined by line **w** perpendicular to plan

- Is data point **x** in class 0 or class 1? $w^T x > 0$ class **1**
  $w^T x < 0$ class **0**

16

---

But, where do we place the boundary?


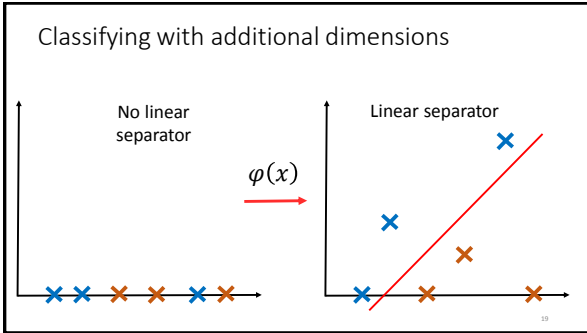
Logistic regression:

$$LL(y|x;w):$$
$$\sum_i (y^i - 1)w^T x^i - \log\left(1 + e^{-w^T x^i}\right)$$

- Each data point $x^i$ considered for boundary **w**

- Outlier data pulls boundary towards it

17

---

Max margin classifiers



- Focus on boundary points

- Find largest margin between boundary points on both sides

- Works well in practice

- We can call the boundary points "**support vectors**"

18

3

## Classifying with additional dimensions

No linear separator

Linear separator

$\varphi(x)$

19

## Mapping function(s)

- Map from low-dimensional space $x = (x_1, x_2)$ to higher dimensional space $\varphi(x) = (x_1, x_2, x_1^2, x_2^2, x_1x_2)$

- N data points guaranteed to be separable in space of N-1 dimensions or more

$$w = \sum_i \alpha_i \varphi(x_i) y_i$$

Classifying $x_j$:

$$\sum_i \alpha_i y_i \varphi^T(x_i)\varphi(x_j) + b$$

20

## Discriminative classifiers

Find a separator to minimize classification error

- Logistic Regression
- Support Vector Machines

21

## Logistic Regression review

Logistic function
$$g(h) = \frac{1}{1 + e^{-h}}$$

- $p(y = 0|x; \theta) = 1 - g(w^T x) = \frac{e^{-w^T x}}{1 + e^{-w^T x}}$
- $p(y = 1|x; \theta) = g(w^T x) = \frac{1}{1 + e^{-w^T x}}$
- Maximize likelihood:
  - $\underset{w}{\text{argmax}}\, L(y|x; w) = \prod_i \left(1 - g(x^i; w)\right)^{(1-y^i)} g(x^i; w)^{y^i}$
  - Likelihood is $P(D|\theta): D = \{(x^i, y^i)\}, \theta = w$
  - Update $w: \frac{\partial}{\partial w_j} LL(y|x; w) = \sum_i x_j^i (y^i - g(w^T x^i))$

22

## MAP for discriminative classifier

- MLE: P(y=1|x;w) ~ g(wᵀx), P(y=0|x;w) ~ 1-g(wᵀx)

- MAP: P(y=1,w|x) ∝ P(y=1|x;w) P(w) ~ g(wᵀx) ???
        (different from Bayesian posterior)

- P(w) priors
  - L2 regularization – minimize all weights
  - L1 regularization – minimize number of non-zero weights

23
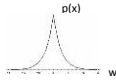
## MAP – L2 regularization

- P(y=1,w|x) ∝ P(y=1|x;w) P(w):

$$L(y, w|x) = \prod_i \left(1 - g(x^i; w)\right)^{(1-y^i)} g(x^i; w)^{y^i} \prod_j e^{-\frac{w_j^2}{2\lambda}}$$

$$LL(y, w|x) = \sum_i y^i w^T x^i - w^T x^i + \log\left(g(w^T x^i)\right) - \sum_j \frac{w_j^2}{2\lambda}$$

$$\frac{\partial}{\partial w_j} LL(y, w|x) = \sum_i x_j^i (y^i - g(w^T x^i)) - \frac{w_j}{\lambda}$$

**Prevent wᵀx from getting too large**

24

4

## MAP – L1 regularization


p(x)

w

- $P(y=1,w|x) \propto P(y=1|x;w)\,P(w)$:

$$L(y,w|x) = \prod_i \left(1 - g(x^i;w)\right)^{(1-y^i)} g(x^i;w)^{y^i} \prod_j e^{-\frac{|w_j|}{\lambda}}$$

$$LL(y,w|x) = \sum_i y^i w^T x^i - w^T x^i + \log\left(g(w^T x^i)\right) - \sum_j \frac{|w_j|}{\lambda}$$

$$\frac{\partial}{\partial w_j} LL(y,w|x) = \sum_i x_j^i\left(y^i - g(w^T x^i)\right) - \frac{\text{sign}(w_j)}{\lambda}$$

**Force most dimensions to 0**

25

## Parameters for learning

$$w_j \leftarrow w_j + \varepsilon\left[x_j^i\left(y^i - g(w^T x^i)\right) - \frac{w_j}{\lambda N}\right]$$

- Regularization: selecting $\lambda$ influences the strength of your bias
- Gradient ascent: selecting $\varepsilon$ influences the effect of individual data points in learning
- Bayesian: selecting $\beta_j$ indicates the strength of the class prior beliefs

- $\lambda, \varepsilon, \beta_j$ are parameters controlling our learning

26

## Multi-class logistic regression: class probability

Recall binary class:

- $p(y = 0|x;\theta) = 1 - g(w^T x) = \frac{e^{-w^T x}}{1+e^{-w^T x}}$
- $p(y = 1|x;\theta) = g(w^T x) = \frac{1}{1+e^{-w^T x}}$

Multi-class – m classes:

- $p(y = j|x;\theta) = \frac{1}{e^{-w_j^T x} + \sum_{k=1}^{m-1} \frac{e^{-w_j^T x}}{e^{-w_k^T x}}}$
- $p(y = m|x;\theta) = 1 - \sum_{j=1}^{m-1} p(y = j|x;\theta)$

27

## Multi-class logistic regression: likelihood

Recall binary class:

- $L(y|x;\theta) = \prod_i p\left(y^i = 0|x^i;\theta\right)^{(1-i)} p\left(y^i = 1|x^i;\theta\right)^i$
- $\frac{\partial}{\partial w_j} LL(y|x;w) = \sum_i x_j^i\left(y^i - g(w^T x^i)\right)$

$$\delta(a) = \begin{cases} 1 & \text{if } a = 0 \\ 0 & \text{otherwise} \end{cases}$$

Multi-class:

- $L(y|x;\theta) = \prod_{i,k} p\left(y^i = k|x^i;\theta\right)^{\delta(y^i-k)}$
- $\frac{\partial}{\partial w_j} LL(y|x;w) = \sum_{i,k} x_j^i\left(\delta(y^i - k) - p\left(y^i = k|x^i;\theta\right)\right)$

28

## Multi-class logistic regression: update rule

Recall binary class:
- $w_j \leftarrow w_j + \varepsilon x_j^i\left(y^i - g(w^T x^i)\right)$

Multi-class:

$$\delta(a) = \begin{cases} 1 & \text{if } a = 0 \\ 0 & \text{otherwise} \end{cases}$$

- $w_{j,k} \leftarrow w_{j,k} + \varepsilon x_j^i\left(\delta(y^i - k) - p\left(y^i = k|x^i;\theta\right)\right)$

29

## Logistic regression: how many parameters?

- N features
- M classes

- Learn $w_k$ for each of M-1 classes: N (M-1) parameters

- Actually: $w^T x = \sum_j w_j x_j$
- Would be better to allow offset from origin: $w^T x + b$ : N+1 parameters per class
- Total (N+1) (M-1) parameters

30

## Max margin classifiers



- Focus on boundary points
- Find largest margin between boundary points on both sides
- Works well in practice
- We can call the boundary points **"support vectors"**

31

## Maximum margin definitions



Classify as +1 if $w^T x + b \geq 1$
Classify as -1 if $w^T x + b \leq -1$
Undefined if $-1 \leq w^T x + b \leq 1$

- M is the margin width
- $x^+$ is a +1 point closest to boundary, $x^-$ is a -1 point closest to boundary
- $x^+ = \lambda w + x^-$
- $|x^+ - x^-| = M$

32

## $\lambda$ derivation



- $w^T x^- + b = -1$
- $w^T x^+ + b = +1$
- $x^+ = \lambda w + x^-$

- $w^T x^+ + b = +1$
- $w^T (\lambda w + x^-) + b = +1$
- $\lambda w^T w + w^T x^- + b = +1$
- $\lambda w^T w - 1 - b + b = +1$
- $\lambda = \frac{2}{w^T w}$

33

## M derivation



- $w^T x^- + b = -1$
- $w^T x^+ + b = +1$
- $x^+ = \lambda w + x^-$
- $|x^+ - x^-| = M$

- $M = |\lambda w + x^- - x^-| = |\lambda w| = \lambda |w|$
- $M = \lambda \sqrt{w^T w}$
- $M = \frac{2}{w^T w} \sqrt{w^T w} = \frac{2}{\sqrt{w^T w}}$

maximize $M$          minimize $w^T w$

34

## Support vector machine (SVM) optimization

- $\text{argmax}_{w,b} \, M = \frac{2}{\sqrt{w^T w}}$

$\text{argmin}_{w,b} \, w^T w$
    subject to
        $w^T x + b \geq 1$     **for x in class +1**
        $w^T x + b \leq -1$     **for x in class -1**

Optimization with constraints: $\frac{\partial}{\partial w_j} f(w_j) = 0$ with Lagrange multipliers.

- Gradient descent
- Matrix calculus

35

## Alternate SVM formulation



$$w = \sum_i \alpha^i \, x^i y^i$$

Support vectors $x^i$ have $\alpha^i > 0$

$y_i$ are the data labels +1 or -1

To classify sample $x^j$, compute:
$$w^T x^j + b = \sum_i \alpha^i \, y^i (x^i)^T x^j + b$$

$\alpha^i \geq 0 \; \forall i$     $\sum_i \alpha^i \, y^i = 0$

36

6

# Benefits of generative methods

- $P(D|\theta)$ and $P(\theta|D)$ can generate non-linear boundary

- E.g.: Gaussians with multiple variances



43