

IAML: Logistic Regression

Nigel Goddard
School of Informatics

Semester 1

- ▶ Logistic function
- ▶ Logistic regression
- ▶ Learning logistic regression
- ▶ Optimization
- ▶ The power of non-linear basis functions
- ▶ Least-squares classification
- ▶ Generative and discriminative models
- ▶ Relationships to Generative Models
- ▶ Multiclass classification
- ▶ Reading: W & F §4.6 (but pairwise classification, perceptron learning rule, Winnow are not required)

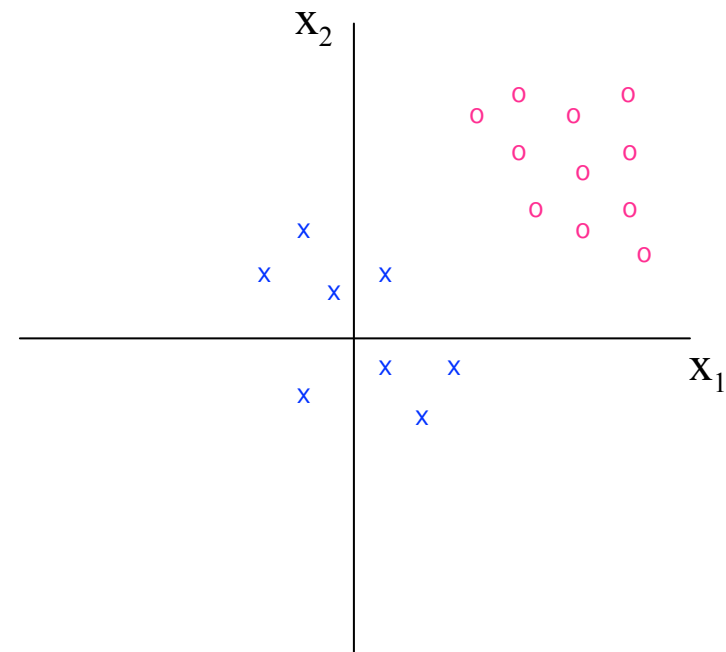
1 / 22

2 / 22

Decision Boundaries

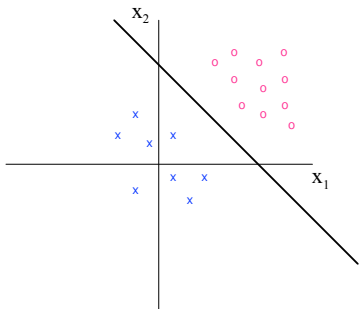
- ▶ In this class we will discuss *linear classifiers*.
- ▶ For each class, there is a *region* of feature space in which the classifier selects one class over the other.
- ▶ The decision boundary is the boundary of this region. (i.e., where the two classes are “tied”)
- ▶ In linear classifiers the decision boundary is a line.

Example Data



3 / 22

4 / 22



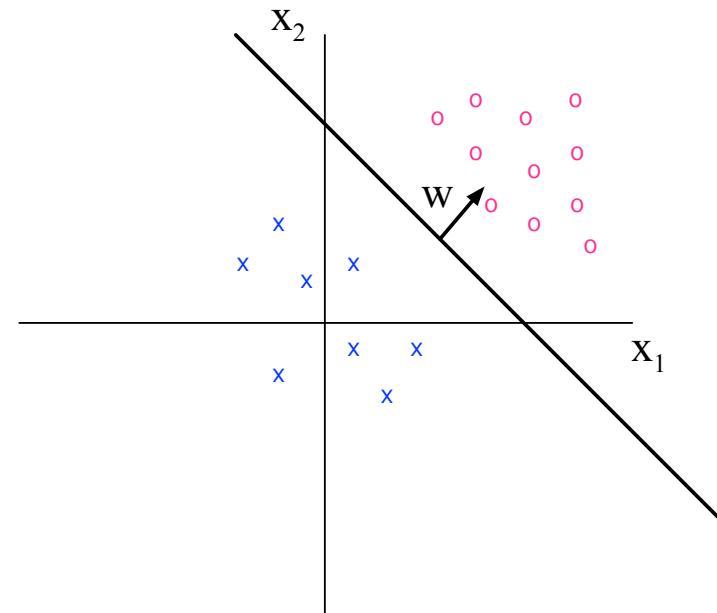
- ▶ In a two-class linear classifier, we learn a function

$$F(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \mathbf{x} + w_0$$

that represents how aligned the instance is with $y = 1$.

- ▶ \mathbf{w} are parameters of the classifier that we learn from data.
- ▶ To do classification of an input \mathbf{x} :

$$\mathbf{x} \mapsto (y = 1) \quad \text{if } F(\mathbf{x}, \mathbf{w}) > 0$$



5 / 22

6 / 22

Explanation of Geometric View

Two Class Discrimination

- ▶ The decision boundary in this case is

$$\{\mathbf{x} | \mathbf{w}^T \mathbf{x} + w_0 = 0\}$$

- ▶ \mathbf{w} is a normal vector to this surface
- ▶ (Remember how lines can be written in terms of their normal vector.)
- ▶ Notice that in more than 2 dimensions, this boundary will be a hyperplane.

- ▶ For now consider a two class case: $y \in \{0, 1\}$.
- ▶ From now on we'll write $\mathbf{x} = (1, x_1, x_2, \dots, x_d)$ and $\mathbf{w} = (w_0, w_1, \dots, w_d)$.
- ▶ We will want a linear, probabilistic model. We could try $P(y = 1 | \mathbf{x}) = \mathbf{w}^T \mathbf{x}$. But this is stupid.
- ▶ Instead what we will do is

$$P(y = 1 | \mathbf{x}) = f(\mathbf{w}^T \mathbf{x})$$

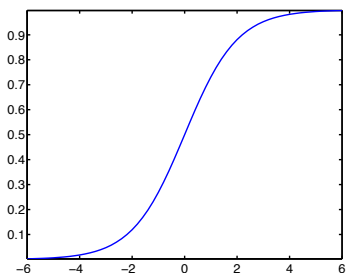
- ▶ f must be between 0 and 1. It will squash the real line into $[0, 1]$
- ▶ Furthermore the fact that probabilities sum to one means

$$P(y = 0 | \mathbf{x}) = 1 - f(\mathbf{w}^T \mathbf{x})$$

7 / 22

8 / 22

- ▶ We need a function that returns probabilities (i.e. stays between 0 and 1).
- ▶ The logistic function provides this
- ▶ $f(z) = \sigma(z) \equiv 1 / (1 + \exp(-z))$.
- ▶ As z goes from $-\infty$ to ∞ , so f goes from 0 to 1, a “squashing function”
- ▶ It has a “sigmoid” shape (i.e. S-like shape)



- ▶ Linear weights + logistic squashing function == logistic regression.
- ▶ We model the class probabilities as

$$p(y = 1 | \mathbf{x}) = \sigma\left(\sum_{j=0}^D w_j x_j\right) = \sigma(\mathbf{w}^T \mathbf{x})$$

- ▶ $\sigma(z) = 0.5$ when $z = 0$. Hence the decision boundary is given by $\mathbf{w}^T \mathbf{x} = 0$.
- ▶ Decision boundary is a $M - 1$ hyperplane for a M dimensional problem.

9 / 22

10 / 22

- ▶ For this slide write $\tilde{\mathbf{w}} = (w_1, w_2, \dots, w_d)$ (i.e., exclude the bias w_0)
- ▶ The bias parameter w_0 shifts the position of the hyperplane, but does not alter the angle
- ▶ The direction of the vector $\tilde{\mathbf{w}}$ affects the angle of the hyperplane. The hyperplane is perpendicular to $\tilde{\mathbf{w}}$
- ▶ The magnitude of the vector $\tilde{\mathbf{w}}$ effects how certain the classifications are
- ▶ For small $\tilde{\mathbf{w}}$ most of the probabilities within the region of the decision boundary will be near to 0.5.
- ▶ For large $\tilde{\mathbf{w}}$ probabilities in the same region will be close to 1 or 0.

- ▶ Want to set the parameters \mathbf{w} using training data.
- ▶ As before:
 - ▶ Write out the model and hence the likelihood
 - ▶ Find the derivatives of the log likelihood w.r.t the parameters.
 - ▶ Adjust the parameters to maximize the log likelihood.

11 / 22

12 / 22

- ▶ Assume data is independent and identically distributed.
- ▶ Call the data set $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$
- ▶ The likelihood is

$$p(D|\mathbf{w}) = \prod_{i=1}^n p(y = y_i | \mathbf{x}_i, \mathbf{w})$$

$$= \prod_{i=1}^n p(y = 1 | \mathbf{x}_i, \mathbf{w})^{y_i} (1 - p(y = 1 | \mathbf{x}_i, \mathbf{w}))^{1-y_i}$$

- ▶ Hence the log likelihood $L(\mathbf{w}) = \log p(D|\mathbf{w})$ is given by

$$L(\mathbf{w}) = \sum_{i=1}^n y_i \log \sigma(\mathbf{w}^\top \mathbf{x}_i) + (1 - y_i) \log(1 - \sigma(\mathbf{w}^\top \mathbf{x}_i))$$

Fitting this into the general structure for learning algorithms:

- ▶ Define the **task**: classification, discriminative
- ▶ Decide on the **model structure**: logistic regression model
- ▶ Decide on the **score function**: log likelihood
- ▶ Decide on **optimization/search method** to optimize the score function: numerical optimization routine. Note we have several choices here (stochastic gradient descent, conjugate gradient, BFGS).

- ▶ It turns out that the likelihood has a unique optimum (given sufficient training examples). It is *convex*.
- ▶ How to maximize? Take gradient

$$\frac{\partial L}{\partial w_j} = \sum_{i=1}^n (y_i - \sigma(\mathbf{w}^\top \mathbf{x}_i)) x_{ij}$$

- ▶ (Aside: something similar holds for linear regression)

$$\frac{\partial E}{\partial w_j} = \sum_{i=1}^n (\mathbf{w}^\top \phi(\mathbf{x}_i) - y_i) x_{ij}$$

where E is squared error.)

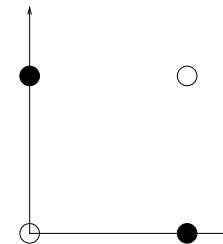
- ▶ Unfortunately, you cannot maximize $L(\mathbf{w})$ explicitly as for linear regression. You need to use a numerical optimisation method, see later.

13/22

14/22

XOR and Linear Separability

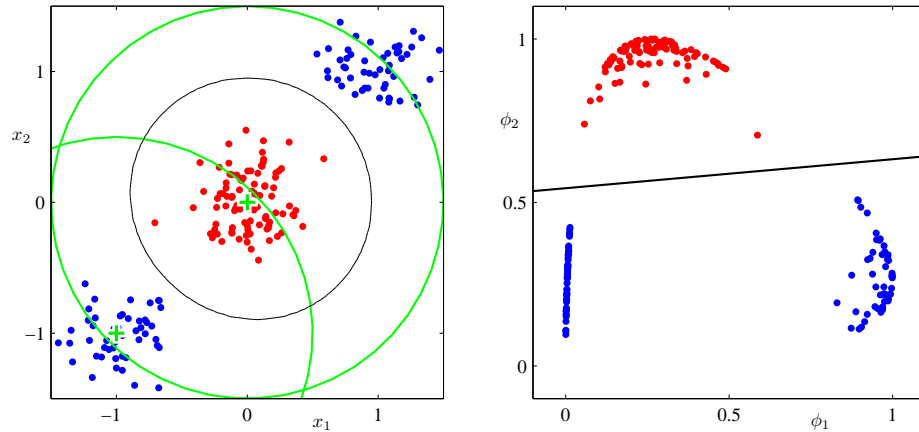
- ▶ A problem is linearly separable if we can find weights so that
 - ▶ $\tilde{\mathbf{w}}^\top \mathbf{x} + w_0 > 0$ for all positive cases (where $y = 1$), and
 - ▶ $\tilde{\mathbf{w}}^\top \mathbf{x} + w_0 \leq 0$ for all negative cases (where $y = 0$)
- ▶ XOR



- ▶ XOR becomes linearly separable if we apply a non-linear transformation $\phi(\mathbf{x})$ of the input — what is one?

15/22

16/22



Using two Gaussian basis functions $\phi_1(\mathbf{x})$ and $\phi_2(\mathbf{x})$

Figure credit: Chris Bishop, PRML

As for linear regression, we can transform the input space if we want $\mathbf{x} \rightarrow \phi(\mathbf{x})$

17/22

- ▶ Notice that we have done something very different here than with naive Bayes.
- ▶ Naive Bayes: Modelled how a class “generated” the feature vector $p(\mathbf{x}|y)$. Then could classify using

$$p(y|\mathbf{x}) \propto p(\mathbf{x}|y)p(y)$$

- . This called is a *generative* approach.
- ▶ Logistic regression: Model $p(y|\mathbf{x})$ directly. This is a *discriminative* approach.
- ▶ Discriminative advantage: Why spend effort modelling $p(\mathbf{x})$? Seems a waste, we’re always given it as input.
- ▶ Generative advantage: Can be good with missing data (remember how naive Bayes handles missing data). Also good for detecting outliers. Or, sometimes you really do want to generate the input.

18/22

Generative Classifiers can be Linear Too

Multiclass classification

Two scenarios where naive Bayes gives you a linear classifier.

1. *Gaussian data with equal covariance.* If $p(\mathbf{x}|y = 1) \sim N(\mu_1, \Sigma)$ and $p(\mathbf{x}|y = 0) \sim N(\mu_2, \Sigma)$ then

$$p(y = 1|\mathbf{x}) = \sigma(\tilde{\mathbf{w}}^T \mathbf{x} + w_0)$$

for some $(w_0, \tilde{\mathbf{w}})$ that depends on μ_1, μ_2, Σ and the class priors

2. *Binary data.* Let each component x_j be a Bernoulli variable i.e. $x_j \in \{0, 1\}$. Then a Naïve Bayes classifier has the form

$$p(y = 1|\mathbf{x}) = \sigma(\tilde{\mathbf{w}}^T \mathbf{x} + w_0)$$

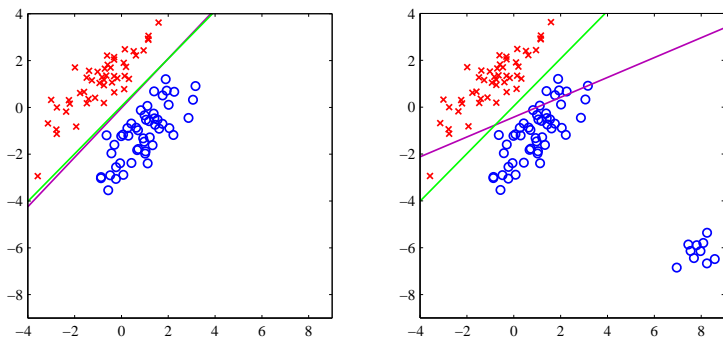
3. Exercise for keepers: prove these two results

- ▶ Create a different weight vector \mathbf{w}_k for each class, to classify into k and not- k .
- ▶ Then use the “softmax” function

$$p(y = k|\mathbf{x}) = \frac{\exp(\mathbf{w}_k^T \mathbf{x})}{\sum_{j=1}^C \exp(\mathbf{w}_j^T \mathbf{x})}$$

- ▶ Note that $0 \leq p(y = k|\mathbf{x}) \leq 1$ and $\sum_{j=1}^C p(y = j|\mathbf{x}) = 1$
- ▶ This is the natural generalization of logistic regression to more than 2 classes.

- ▶ Logistic regression is more complicated algorithmically than linear regression
- ▶ Why not just use linear regression with 0/1 targets?



Green: logistic regression; magenta, least-squares regression

Figure credit: Chris Bishop, PRML

- ▶ The logistic function, logistic regression
- ▶ Hyperplane decision boundary
- ▶ Linear separability
- ▶ We still need to know how to *compute* the maximum of the log likelihood. Coming soon!