

## **3F8: Inference**

### **Classification**

**José Miguel Hernández–Lobato and Richard E. Turner**

Department of Engineering  
University of Cambridge

Lent Term

# What is classification?

It is the same as regression, but with **discrete outputs**:  $y_n \in \{1, \dots, C\}$ , where  $C$  is the number of **classes**. Often  $C = 2$ .

Same goals as in regression. The patterns to identify consists of

- A partition of the input space into  $C$  **decision regions**, one for each class.
- Each new input is assigned the class of its corresponding decision region.
- We would also like a measure of **confidence** in the decisions.

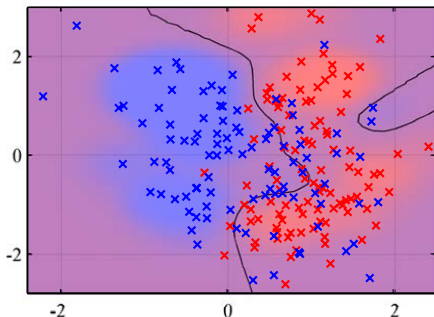
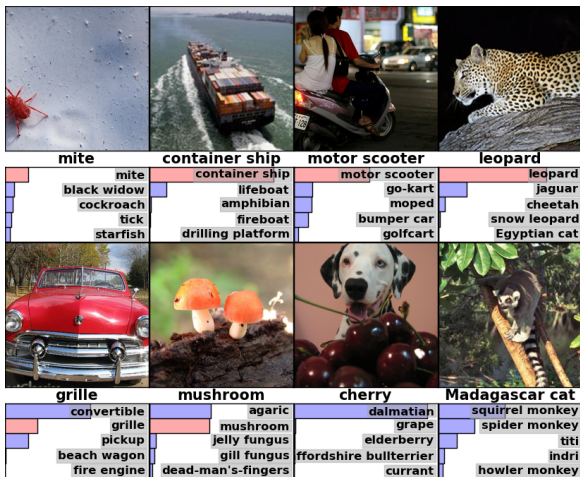


Figure: C. Bishop. *Pattern Recognition and Machine Learning*, 2006.

# Real-world example

ImageNet: about 22,000 classes and 15 million high-resolution images.



A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In NIPS, 2012.

# Why not use methods for regression?

Let the output  $\mathbf{y}_n$  be a  $C$ -dimensional vector with a **one-hot-encoding** of the class for  $\tilde{\mathbf{x}}_n$  ( $y_{n,c} = 1$  if the class is  $c$  and  $y_{n,c} = 0$ , otherwise).

We can then solve  $C$  linear regression problems, one for each class:

$$\mathbf{W} = \left( \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}^T \mathbf{Y},$$

with  $\mathbf{Y} = (\mathbf{y}_1; \dots; \mathbf{y}_N)^T$  and then predict the class for the highest entry of  $\tilde{\mathbf{x}}_\star^T \mathbf{W}$ .

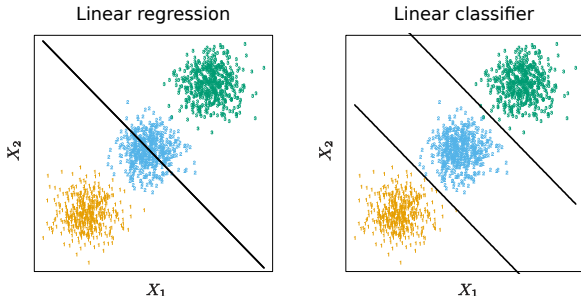


Figure: G. James, D. Witten, T. Hastie and R. Tibshirani. *An Introduction to statistical learning*, 2013.

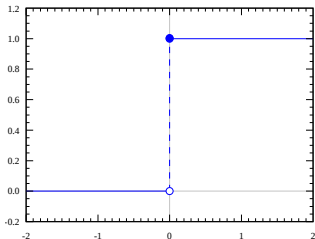
# Deterministic linear classification

Works by mapping the output of the linear model into **discrete class labels**.

Assume  $y_n \in \{0, 1\}$  (binary classification). Then, we can define

$$y_n = H(\mathbf{w}^T \tilde{\mathbf{x}}),$$

where  $H(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ 0 & \text{otherwise} \end{cases}$  is the Heaviside step function:



What is the **geometric interpretation** of  $\mathbf{w}$ ?

**Problem:** deterministic predictions.

- Missclassification errors not allowed.
- Inference is hard: what is the MLE?

# Probabilistic linear classification

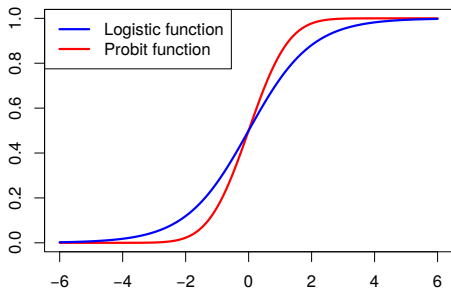
Works by mapping the output of the linear model into **class probabilities**.

$$p(y_n = 1 | \tilde{\mathbf{x}}, \mathbf{w}) = \sigma(\mathbf{w}^T \tilde{\mathbf{x}}),$$

where  $\sigma(\cdot)$  is a monotonically increasing function that maps  $\mathbb{R}$  into  $[0, 1]$ .

For example:

- The logistic function:  $\sigma(x) = 1/(1 + \exp(-x))$ .
- The probit function or Gaussian CDF:  $\sigma(x) = \int_{-\infty}^x \mathcal{N}(z|0, 1) dz$ .



# Logistic regression (classification)

Assume  $\sigma(x)$  is the **logistic function** and that  $y_n \in \{-1, 1\}$ . Then

$$p(y_n|\mathbf{x}_n, \mathbf{w}) = \sigma(y_n \mathbf{w}^T \tilde{\mathbf{x}}_n),$$

since  $1 - \sigma(x) = \sigma(-x)$ . For  $\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N$ , the **log-likelihood** is

$$\mathcal{L}(\mathbf{w}) = \log p(y_1, \dots, y_N | \mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{w})$$

$$\mathcal{L}(\mathbf{w}) = \sum_{n=1}^N \log p(y_n | \mathbf{x}_n, \mathbf{w}) = \sum_{n=1}^N \log \sigma(y_n \mathbf{w}^T \tilde{\mathbf{x}}_n).$$

We can then use  $d\sigma(x)/dx = \sigma(x)(1 - \sigma(x))$  to obtain the gradient:

$$\frac{d\mathcal{L}(\mathbf{w})}{d\mathbf{w}} = \sum_{n=1}^N y_n \underbrace{(1 - \sigma(y_n \mathbf{w}^T \tilde{\mathbf{x}}_n))}_{\text{Error}} \tilde{\mathbf{x}}_n.$$

No closed-form solution for MLE, but gradient has **geometric interpretation**.

# Gradient ascent and stochastic gradient ascent

The batch gradient ascent rule to maximize  $\mathcal{L}(\mathbf{w})$  is

$$\mathbf{w}^{\text{new}} = \mathbf{w}^{\text{old}} + \alpha \sum_{n=1}^N y_n (1 - \sigma(y_n \mathbf{w}^{\text{T}} \tilde{\mathbf{x}}_n)) \tilde{\mathbf{x}}_n .$$

where  $\alpha > 0$  is the **learning rate**.

When  $N$  is very large evaluating the exact gradient at each step is wasteful.

**Solution:** use stochastic gradient ascent (SGD):

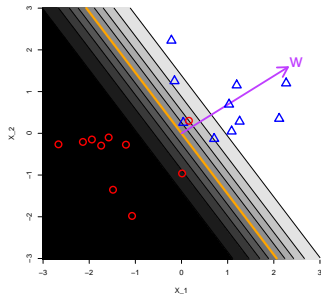
$$\mathbf{w}^{\text{new}} = \mathbf{w}^{\text{old}} + \alpha \frac{N}{|\mathcal{S}|} \sum_{n \in \mathcal{S}} y_n (1 - \sigma(y_n \mathbf{w}^{\text{T}} \tilde{\mathbf{x}}_n)) \tilde{\mathbf{x}}_n .$$

where  $\mathcal{S}$  is a set with the indexes of data points in a **minibatch** sampled uniformly at random from the training data before each update step.

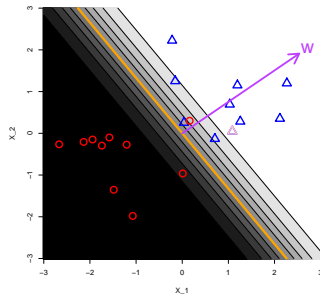


# Example logistic regression

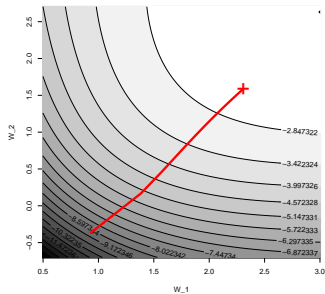
Batch Gradient Ascent



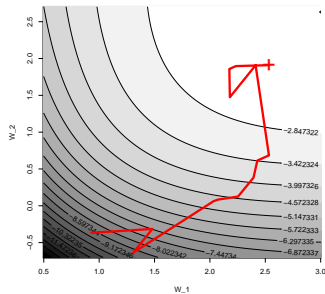
Stochastic Gradient Ascent  $|S| = 1$



Log-likelihood surface



Log-likelihood surface



# Linear classification with more than 2 classes

We can map multiple outputs to discrete class labels using the **max function**:

$$y_n = \arg \max_{k \in \{1, \dots, C\}} \mathbf{w}_k^T \tilde{\mathbf{x}},$$

but this has similar problems as in the deterministic binary classification case.

Instead, use the **soft-max function** to map the outputs into class probabilities:

$$p(y_n = k | \mathbf{w}_1, \dots, \mathbf{w}_K, \tilde{\mathbf{x}}_n) = \frac{\exp(\mathbf{w}_k^T \tilde{\mathbf{x}}_n)}{\sum_{k'=1}^K \exp(\mathbf{w}_{k'}^T \tilde{\mathbf{x}}_n)}.$$

Equivalent also to logistic regression when  $C = 2$ .

# Non-linear logistic regression

Replace  $\mathbf{x}$  with non-linear functions of the inputs  $\phi(\mathbf{x}) = (\phi_1(\mathbf{x}), \dots, \phi_M(\mathbf{x}))^T$ .

**Inference does not change**, just replace each  $\mathbf{x}_n$  with the new  $\phi(\mathbf{x}_n)$ .

For example,  $(x_1, x_2) \rightarrow (x_1, x_2, x_1x_2, x_1^2, x_2^2)$ .

