

Short course

A vademecum of statistical pattern recognition and machine learning

Classification

Massimo Piccardi
University of Technology, Sydney, Australia

© Massimo Piccardi, UTS 1

Agenda

- Classification
- Classification accuracy
- Binary classification test
- ROC curve
- Classification accuracy for multiple classes
- Maximum-a-posteriori decision rule
- Classification by Bayes' theorem
- Bayes risk
- Minimum risk decision rule
- Generative vs discriminative approaches
- Learning a probabilistic classifier

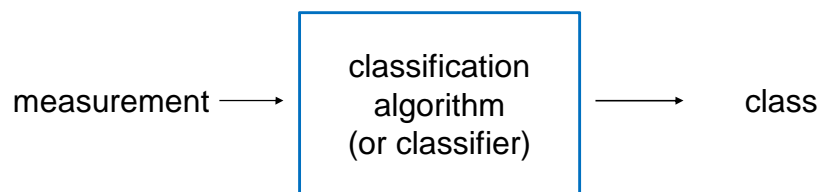
© Massimo Piccardi, UTS 2

Classification

- Let us assume that we are given a set of **classes** (or **categories**)
- Assume that we are also given an object that can be categorised into one of the classes. The assignment of the object to a class is called **classification**
- Classification is performed by obtaining a **measurement** of the object and applying a classification algorithm

© Massimo Piccardi, UTS 3

Classification



- The **class** takes *symbolic* values: “good”, “bad”, “Parisian”, “trustworthy”, “cat”, “dog”...
- The **measurement** can be multiple (e.g., age, income, occupation), and both numerical or symbolic

© Massimo Piccardi, UTS 4

Classification

- Let us assume both the class and the measurement to be random variables, \mathbf{y} and \mathbf{x} , respectively
- We aim to provide a *prediction* (aka *assignment*), $\tilde{\mathbf{y}}$, that is as close as possible to the “ground-truth” class, \mathbf{y}
- $\tilde{\mathbf{y}}$ is another random variable that can be obtained from \mathbf{x} through a deterministic algorithm, $\mathbf{f}(\mathbf{x})$:

$$\tilde{\mathbf{y}} = \mathbf{f}(\mathbf{x}) \quad \text{classifier}$$

© Massimo Piccardi, UTS 5

Statistical classification

- In statistical classification, the fundamental quantity used for classification is:

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

- Deriving $p(\mathbf{y}|\mathbf{x})$ is (properly) called **inference**. Often, it is confused/assimilated with classification itself, but it is only a part of it

© Massimo Piccardi, UTS 6

Parameters

- Where are the parameters? Like any distribution, $p(y|x)$ can be modelled through a chosen family and a set of parameters, as either a point-based or Bayesian estimate:

$$p(y | x, \theta^*)$$
$$p(y | x) \approx \int_{\theta} p(y | x, \theta) p(\theta | Y, X) d\theta$$

- In the following, we omit the explicit dependence on the parameters. We'll address parameters and learning later

© Massimo Piccardi, UTS 7

Classification accuracy

- Before we address how to build a classification algorithm, let us discuss how we will measure its accuracy
- Let us adopt more prominent symbols for y and \tilde{y} :
 - \mathbf{T} , the “true class” of the given object (aka *ground truth*)
 - \mathbf{A} , the class to which the classification algorithm assigns it
- We wish that classification will be as accurate as possible, i.e. \mathbf{A} be the same as \mathbf{T} most of the times (if not all)

© Massimo Piccardi, UTS 8

Binary classification test

- Let us start with *binary classification*, the case where there are only two classes
- In addition, in a binary classification *test*, the two classes are asymmetric, of the type “it has the property/it hasn’t”
- Classes are often called the two hypotheses, **H1** and **H0**
- A binary classification test is often called a *detector*

© Massimo Piccardi, UTS 9

Binary classification test

- Consistently with the idea of a test, based on the **true** labels, the samples can be divided into:
 - Positives (P) and Negatives (N)
- After classification, based on both the **true and assigned** labels, the same samples can be categorised as:
 - True Positives (TP), True Negatives (TN), False Positives (FP), False Negatives (FN)
- Obviously, TP and TN are desirable, FP and FN are not

© Massimo Piccardi, UTS 10

Contingency table

- Using symbols P , N , TP , TN , FP , FN to mean the corresponding numbers, the *contingency table* for the binary classification test is:

		A		
		H0	H1	
T	H0	TN	FP	← N
	H1	FN	TP	← P

- Please note that the following holds:
 - Total number of samples = $P + N = TP + TN + FP + FN$
 - $P = TP + FN$ and $N = TN + FP$

© Massimo Piccardi, UTS 11

Probabilities in a binary test

The contingency table can be used to estimate probabilities in a maximum-likelihood sense. For instance:

- joint probability $p(A = H1, T = H1)$** is given by the number of true positives divided the total number of samples:
 - $p(A = H1, T = H1) = TP/(P + N)$
- marginal probability $p(T = H1)$** can be obtained either directly as $P/(P + N)$ or by marginalising A in the above:
 - $p(T = H1) = p(A = H1, T = H1) + p(A = H0, T = H1) =$
 $= TP/(P + N) + FN/(P + N) = P/(P + N)$
- conditional probability $p(A = H1 | T = H0)$** is given by:
 - $p(A = H1 | T = H0) = FP/N$

© Massimo Piccardi, UTS 12

Binary classification test

- Accuracy of a binary classification test is typically expressed by two numbers:
 - The rate of true H1 classified as H1 (**detection rate**, DR):
 $p(A=H1 \mid T=H1) = TP/P = TP/(FN+TP)$
 - The rate of true H0 classified as H1 (**false alarm rate**, FAR):
 $p(A=H1 \mid T=H0) = FP/N = FP/(TN+FP)$
- Similar rates take different names in different fields:
 - Signal Processing: the above
 - Information Theory: Recall \equiv DR;
Precision $\equiv TP/(TP+FP) = p(T=H1 \mid A=H1)$
 - Medicine: Sensitivity \equiv DR; Specificity $\equiv 1 - FAR$
- $FP \equiv$ *Type 1*, or *alpha*, error; $FN \equiv$ *Type 2*, or *beta*, error

© Massimo Piccardi, UTS 13

Binary classification test: example

Example:

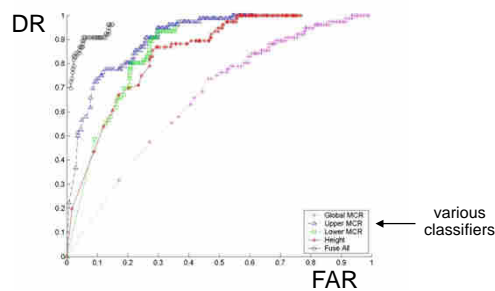
- 100 samples, values as in table
- Marginal probabilities of T:
 - $p_{H0} = p(T=H0) = 88/100 = .88$
 - $p_{H1} = p(T=H1) = 12/100 = .12$
- FAR and DR:
 - $p_{FA} = FAR = 8/88 = 9\%$
 - $p_D = DR = 10/12 = 83\%$
- All the joint probabilities can be built from the above with Bayes' theorem, or derived directly from the table
- The **overall accuracy** is the number of correctly classified samples divided the total number of samples

T	A	H0	H1
H0		80	8
H1		2	10

© Massimo Piccardi, UTS 14

ROC (receiver operating characteristic) curve

- Often the parameters of a classifier can be changed to achieve a better DR at the expense of a worse FAR, or vice versa
- The plot of (FAR, DR) values is a *ROC curve*. The area under it is known as the *area under the curve* (AUC)



© Massimo Piccardi, UTS 15

Classification accuracy for multiple classes

- *Confusion matrix* (example with 4 classes):

T \ A				
	c ₁	c ₂	c ₃	c ₄
c ₁	7	0	3	0
c ₂	2	30	2	6
c ₃	0	0	34	6
c ₄	0	0	0	35

- The confusion matrix provides all information about accuracy, for each class and as a whole
 - what is the overall accuracy above? and the *overall error rate* (1 – accuracy)?
 - what is the detection rate/recall for class C3? and its false alarm rate?

© Massimo Piccardi, UTS 16

Accuracy and error rate

- The overall **accuracy** and **error rate** ($1 - \text{accuracy}$) can be obtained as a weighted average over all the classes:

$$\text{accuracy} = \sum_{j=1}^M p(A = c_j | T = c_j) p(T = c_j) = \sum_{j=1}^M p(A = c_j, T = c_j)$$

$$\text{error rate} = \sum_{j=1}^M \sum_{i=1, i \neq j}^M p(A = c_i | T = c_j) p(T = c_j) = \sum_{j=1}^M \sum_{i=1, i \neq j}^M p(A = c_i, T = c_j)$$

* See the Notes Page *

© Massimo Piccardi, UTS 17

Empirical vs theoretical rates

- Alternatively the accuracy and error rate can be written as expectations of indicator functions. For instance, for the error rate let us define $\Delta(T, A)$ evaluating to 1 if $A \neq T$ and 0 otherwise:

$$\text{error rate} = \sum_{i=1}^M \sum_{j=1}^M \Delta(T, A) p(A = c_i, T = c_j) = E[\Delta(T, A)]$$

- If we use empirical distributions (i.e., estimates from a sample set) in the expectations, we talk about *empirical rates*. Using exact distributions returns the *theoretical rates*

* Notes Page *

© Massimo Piccardi, UTS 18

Validation

- Most classifiers need a training stage, during which the parameters of the classifier are optimised based on a given set of samples (*training set*)
- If the accuracy is measured on the same set, the risk is that of obtaining an overly optimistic value
- The *test set* must be different from the training set; given a set of samples, this may be split into a training set (e.g., 66%) and a test set (e.g., 34%). No hard rule for the split. The measured accuracy is called **cross-validation accuracy**

© Massimo Piccardi, UTS 19

Validation

- Repeating cross-validation several times and averaging the resulting accuracy may lead to more realistic estimates which depend less on the specific training and test sets
- In *N-fold cross validation*, a given data set is divided into N “slices”: the data in N-1 are used for training and the last for testing; the procedure is repeated N times, and accuracy averaged

© Massimo Piccardi, UTS 20

Minimum error rate as assignment criterion

- Based on a measurement, x , of an object, we want to *optimally assign* it to one of M classes, c_j , $j=1..M$
- We need to choose an assignment criterion: let us choose the **minimum error rate**
- We assume that we know $p(y|x)$ (and $p(x)$, and related marginals and conditionals) or that we can estimate them from a set of labelled samples
- We need a **decision rule**, to complete $p(y|x)$ into a classifier $f(x)$

© Massimo Piccardi, UTS 21

Maximum-a-posteriori decision rule

- It can be proven that the minimum error rate is obtained by applying this decision rule to each sample:

$$\tilde{y} = f(x) = \arg \max_{j=1..M} p(y = c_j | x)$$

- The above criterion is called the **maximum-a-posteriori (MAP)** decision rule
- In the next slides, instead of $p(y = c_j | x)$ we will use short-hand notation $p(c_j | x)$

© Massimo Piccardi, UTS 22

Example

- Let us say that we need to classify a piece of fruit as one of classes {apple, pear, avocado*} based on a measurement of size and weight
- For a given measurement x , we have come to know that the probabilities are:
$$p(y = \text{apple} \mid x) = 0.2$$
$$p(y = \text{pear} \mid x) = 0.5$$
$$p(y = \text{avocado} \mid x) = 0.3$$
- According to the MAP rule, we'll classify this object as a pear

* avocados are classified as fruit despite their primary use in savoury food; in that, they share the fate of tomatoes

© Massimo Piccardi, UTS 23

Classification by Bayes' theorem

- Instead of using $p(c_j \mid x)$ directly, applying Bayes' theorem allows us to write:

$$p(x, c_j) = p(x \mid c_j)p(c_j) = p(c_j \mid x)p(x)$$
$$\rightarrow p(c_j \mid x) = \frac{p(x \mid c_j)p(c_j)}{p(x)}$$

- Why is this theorem so fundamental? Because it can prove simpler to learn the probabilities on the right separately

© Massimo Piccardi, UTS 24

Classification by Bayes' theorem

- This formula is so canonical that its quantities have names:

$$p(c_j | x) = \frac{p(x | c_j)p(c_j)}{p(x)}$$

Diagram labels with arrows pointing to the formula components:

- $p(c_j | x)$: posterior probability
- $p(x | c_j)$: (class-conditional) likelihood
- $p(c_j)$: prior probability
- $p(x)$: evidence

- Note that the same class label is obtained if we maximise either $\mathbf{p(c_j|x)}$ (*posterior*) or $\mathbf{p(c_j,x)}$ (*joint*); $p(x)$ (*evidence*) doesn't count in the assignment

© Massimo Piccardi, UTS 25

Notes

- Why names *prior* and *posterior*?
 - $\mathbf{p(c_j)}$ is the probability of the class *before* (i.e. prior to) the measurement;
 - $\mathbf{p(c_j | x)}$ is the probability of the class *after* (i.e. a posteriori of) the measurement
- As expected, the $p(c_j|x)$ add up to 1 as shown by:

$$p(c_j | x) = \frac{p(x | c_j)p(c_j)}{p(x)} = \frac{p(x | c_j)p(c_j)}{\sum_j p(x | c_j)p(c_j)}$$

Diagram label with an arrow pointing to the denominator of the second fraction: *evidence of x*

© Massimo Piccardi, UTS 26

Example

- You hear an engine at night: is that a car (C) or a motorbike (M)?
- Let us assume that we measured the prior probabilities:
 $p(C) = 0.90$, $p(M) = 0.10$
- Let us also assume that we can model likelihoods $p(\text{noise} | C)$, $p(\text{noise} | M)$ from aural samples of cars and motorbikes
- First audio clip 🗣️: $p(\text{noise} | M) = 0.60$, $p(\text{noise} | C) = 0.20$
 - $p(M | \text{noise}) \propto p(\text{noise} | M) * p(M) = 0.60 * 0.10 = 0.06$
 - $p(C | \text{noise}) \propto p(\text{noise} | C) * p(C) = 0.20 * 0.90 = 0.18 \rightarrow \text{stick with car}$
- Second audio clip 🗣️: $p(\text{noise} | M) = 1.2$; $p(\text{noise} | C) = 0.10$
 - $p(M | \text{noise}) \propto 1.2 * 0.1 = 0.12$
 - $p(C | \text{noise}) \propto 0.10 * 0.90 = 0.09 \rightarrow \text{opt for bike}$

© Massimo Piccardi, UTS 27

Bayes risk

- In many cases, the *cost*, or *loss*, associated with a misclassification varies with the combination of the true and assigned classes
- In such cases, errors are weighted by a loss function, $\Delta(\mathbf{y} = \mathbf{c}_i, \tilde{\mathbf{y}} = \mathbf{c}_j)$ (or Δ_{ij} for short), that is the cost of misclassifying class c_i (true class) as c_j (assigned class)

© Massimo Piccardi, UTS 28

Examples of loss functions

- In the simplest case, every type of error is given equal weight (e.g., 1). We can write the loss function (known as *zero-one loss*) as:

$$\Delta(T,A) \text{ (1 if } A \neq T, 0 \text{ otherwise)}$$

- Classifying a defective manufactured part as good may carry a higher cost than classifying a good part as defective (asymmetry: $\Delta_{ij} \neq \Delta_{ji}$)
- Classifying a colour *red* as *green* may have a higher cost than classifying *red* as *orange* because green is more different than orange from red
- In certain cases, classes may have a natural numerical value (e.g., the intensity levels of a pixel). In this case, Δ_{ij} could be set to the square error (SE): $\Delta_{ij} = (i - j)^2$

© Massimo Piccardi, UTS 29

Bayes risk

- The *overall loss*, also known as *Bayes' risk*, *integrated risk*, *expected loss* etc can be written as:

$$Bayes' risk = \sum_y \sum_{\tilde{y}} \Delta(y, \tilde{y}) p(y, \tilde{y}) = E[\Delta(y, \tilde{y})]$$

- In other words, it is the expected value of cost function $\Delta(y, \tilde{y})$
- The overall loss is more commonly defined in terms of variables y and x , but this definition suits our review
- For the zero-one loss, it is the same as the error rate

* Notes Page *

© Massimo Piccardi, UTS 30

Minimum risk decision rule

- It can be proven that the minimum expected loss is achieved by the following decision rule:

$$\tilde{y} = f(x) = \arg \min_u \left(\sum_y \Delta(y, u) p(y | x) \right)$$

- The rule says: for every class u (auxiliary variable), compute the *expected loss from this assignment* and choose the class causing the minimum expected loss
- This decision rule is called the **minimum risk decision rule**

© Massimo Piccardi, UTS 31

Example

- A hypothetical loss function for our example:

<i>true / assigned</i>	<i>apple</i>	<i>pear</i>	<i>avocado</i>
$\Delta(y, u) =$ <i>apple</i>	0	2	2
<i>pear</i>	0.2	0	0.3
<i>avocado</i>	0.4	1	0

© Massimo Piccardi, UTS 32

Example

- Expected loss for choosing $u = \text{apple}$:
$$\Delta(y = \text{apple}, u = \text{apple}) p(y = \text{apple} | x) = 0 * 0.2 +$$
$$+ \Delta(y = \text{pear}, u = \text{apple}) p(y = \text{pear} | x) = 0.2 * 0.5 +$$
$$+ \Delta(y = \text{avocado}, u = \text{apple}) p(y = \text{avocado} | x) = 0.4 * 0.3 =$$
$$= 0.22$$
- Similarly, we compute the loss for $u = \text{pear}$ as equal to 0.7; the loss for $u = \text{avocado}$ is 0.55
- Based on the minimum risk decision rule, *we assign the piece of fruit to apple with a large margin*

© Massimo Piccardi, UTS 33

Generative vs discriminative classifiers

- Classification by Bayes' theorem estimates $p(x|y)$ and $p(y)$ explicitly, and obtains $p(y|x)$ via Bayes' theorem: this is known as a *generative* approach
- The reason for the *generative* name is that one can obtain (i.e. generate) samples of x by sampling either $p(x|y)$ or $p(x)$
- As we have said many times, a model like $p(y|x)$ alone does not contain any probabilistic information on x (and therefore cannot be sampled)

© Massimo Piccardi, UTS 34

Discriminative approaches

- In *discriminative* probabilistic approaches, $p(y|x)$ is estimated directly. A typical model is the *logistic regression classifier*
- One can also renounce probability altogether and build a *discriminant function*, $f(x)$, without reference to $p(y|x)$
- Discriminative approaches in general can focus on class boundaries rather than entire densities: they may prove more accurate. As a most notable case, in a later lecture we'll present the *support vector machine*

© Massimo Piccardi, UTS 35

Classification by regression

- In probabilistic approaches to classification, we model $p(y|x)$
- A simpler approach is to just put the classes in correspondence with regions of real numbers (or \mathbb{R}^n) by building a **discriminant function**, $f(x)$, that maps the measurement to the regions
- A simple example is a binary classifier for classes $\{A, B\}$ that associates class A with $f(x) \geq 0$ and class B with $f(x) < 0$
- During training, the classes are given numerical values such as (+1, -1) or (+1, 0) or else; it is, in a way, classification by regression
- Examples of classifiers following this approach are linear classifiers (including the linear support vector machine), quadratic classifiers, AdaBoost and many others

© Massimo Piccardi, UTS 36

Classification approaches: summary

generative $\tilde{y} = \arg \max_y (p(x|y)p(y)) \quad (1)$

discriminative $\tilde{y} = \arg \max_y p(y|x) \quad (2)$

discriminant $\tilde{y} = f(x) \quad (3)$

- If a loss function different from the zero-one loss needs to be used, the decision rule for (1) and (2) has to be as with the Bayes risk. For (3), the loss function has to be suitably accounted for inside $f(x)$

© Massimo Piccardi, UTS 37

Learning

- How to learn a probabilistic classifier? Let us start by choosing to use Bayes' theorem to provide classification:

$$\tilde{y} = \arg \max_y (p(x|y)p(y))$$

- Let us now make the parameters of the involved distributions explicit:

$$\tilde{y} = \arg \max_y (p(x|y, \theta_{xy})p(y|\theta_y))$$

- θ_{xy} represents the parameters of the likelihood of class y and θ_y represents the parameters of the prior

© Massimo Piccardi, UTS 38

Learning

- Distribution $p(x|y, \theta_{xy})$ can be, for instance, a multivariate Gaussian. In that case, θ_{xy} represents its mean and covariance. Or it can be any other distribution in x
- There will be one likelihood *per class* (with M classes, M likelihoods, each with its own sets of parameters, θ_{xy})
- Distribution $p(y|\theta_y)$ is certainly a discrete distribution since the class variable, y , is discrete. Its parameters, θ_y , are the M probability values over the simplex
- For clarity, there is only one prior distribution overall
- θ notes the union set of all θ_{xy} and θ_y

© Massimo Piccardi, UTS 39

Objective function: likelihood

- Given a training set of N pairs of measurements and class labels, $\{X, Y\} = \{x_i, y_i\}$, $i = 1 \dots N$, we choose to train the classifier by maximising the following objective function:

$$\theta^* = \arg \max_{\theta} (p(X | Y, \theta_{xy}) p(Y | \theta_y))$$

- This objective function is called the **(joint) likelihood of $\{x, y\}$** ($p(X, Y|\theta)$ for short) and is the same function used for the inference. In this way, we choose the parameters that maximise our inference function over the given training set

© Massimo Piccardi, UTS 40

Maximum likelihood

- Further, we make it explicit that the pairs are i.i.d. given the parameters:

$$\theta^* = \arg \max_{\theta} \prod_{i=1}^N (p(x_i | y_i, \theta_{xy}) p(y_i | \theta_y))$$

- Let us then move to logarithmic scale:

$$\theta^* = \arg \max_{\theta} \sum_{i=1}^N (\ln p(x_i | y_i, \theta_{xy}) + \ln p(y_i | \theta_y))$$

© Massimo Piccardi, UTS 41

Maximum likelihood

- To maximise the objective function we differentiate it in each parameter and equate to 0
- Now, the founding point: when we differentiate it in θ_{xy} for class, say, l , all other terms cancel:

$$\begin{aligned} \frac{\partial}{\partial \theta_{xy=l}} \sum_{i=1}^N (\ln p(x_i | y_i, \theta_{xy}) + \ln p(y_i | \theta_y)) &= \\ &= \frac{\partial}{\partial \theta_{xy=l}} \sum_{i=1}^{N_l} (\ln p(x_i | y_i = l, \theta_{xy=l})) \end{aligned}$$

- This is equivalent to maximising the likelihood for this class alone, using only the samples from that class

© Massimo Piccardi, UTS 42

Maximum likelihood

- Likewise, when we differentiate it in θ_y , all non-relevant terms cancel:

$$\begin{aligned}\frac{\partial}{\partial \theta_y} \sum_{i=1}^N (\ln p(x_i | y_i, \theta_{xy}) + \ln p(y_i | \theta_y)) &= \\ &= \frac{\partial}{\partial \theta_y} \sum_{i=1}^N (\ln p(y_i | \theta_y))\end{aligned}$$

- This is equivalent to maximising the likelihood of a common discrete distribution!

© Massimo Piccardi, UTS 43

Discriminative training

- Conversely, let us now assume that we *do not want* to use Bayes' theorem to provide classification. We use directly:

$$\tilde{y} = \arg \max_y p(y | x, \theta)$$

where we have made the parameters explicit

- Given the same training set as before, it is now desirable to train θ by maximising this inference function over the training set:

$$\theta^* = \arg \max_{\theta} p(Y | X, \theta)$$

© Massimo Piccardi, UTS 44

Objective function: conditional likelihood

- This objective function is known as the **conditional likelihood** and the corresponding maximisation as *maximum conditional likelihood estimation* (MCLE)
- By applying the i.i.d. assumption:

$$\theta^* = \arg \max_{\theta} \prod_{i=1}^N p(y_i | x_i, \theta)$$

© Massimo Piccardi, UTS 45

Maximum conditional likelihood

- Just to illustrate our point, let us assume that the distribution factorises into the same terms as before:

$$\begin{aligned} \theta^* &= \arg \max_{\theta} \prod_{i=1}^N \left(\frac{p(x_i | y_i, \theta) p(y_i | \theta)}{p(x_i | \theta)} \right) = \\ &= \arg \max_{\theta} \prod_{i=1}^N \left(\frac{p(x_i | y_i, \theta_{xy}) p(y_i | \theta_y)}{\sum_y p(x_i | y, \theta_{xy}) p(y | \theta_y)} \right) \end{aligned}$$

the second equality coming from Bayes' theorem and marginalisation. It can be immediately seen that this distribution is significantly more complex *due to the denominator*

© Massimo Piccardi, UTS 46

Maximum conditional likelihood

- Let us move to logarithmic scale:

$$\theta^* = \arg \max_{\theta} \sum_{i=1}^N \left(\ln p(x_i | y_i, \theta_{xy}) + \ln p(y_i | \theta_y) \right) - \ln \sum_y p(x_i | y, \theta_{xy}) p(y | \theta_y)$$

- When we differentiate in, say, $\theta_{xy=l}$, these parameters appear in all N addenda!, not just those with the samples of class l
- In addition, the parameters of the other classes, $\theta_{xy \neq l}$, appear in the same equations, forcing us to a *simultaneous optimisation of all parameters*

© Massimo Piccardi, UTS 47

Learning: summary

- Maximising the *joint likelihood* is as simple as maximising the likelihood of each class separately
- Maximising the *conditional likelihood* requires maximising in all the parameters *simultaneously*, over samples from all classes *at once*, and the expression is more complicated than the separate likelihoods
- However, it is paramount to say that the advantage of training $p(y|x)$ directly is that we are not required to fully train likelihoods $p(x)$ or $p(x|y)$. Maximum conditional likelihood becomes a suitable objective for the *exponential family of distributions*. This is the topic of a following lecture

© Massimo Piccardi, UTS 48

Possible further readings

- Bishop's PRML book, chapter 4.2 Probabilistic Generative Models and 4.3 Probabilistic Discriminative Models
- Andrew Ng, Michael Jordan, "On discriminative vs. generative classifiers: A comparison of logistic regression and naive Bayes", NIPS 2001, pp. 841-848
- Tony Jebara, "Machine Learning: Discriminative and Generative", Springer, 2003