

IntroML

ML. 3. (Linear) Classification

DataLab CSIC

Objectives and schedule

General concepts and basic algos in (linear) classification

- Discriminant functions
- Probabilistic discriminative approaches
- Probabilistic generative approaches

Contents

- Key concepts: Likelihood optimisation for classification, Bayesian classification, Stochastic gradient descent, multiple class, metrics, glms (beyond logistic regression)
- Later chapters more on non-linear classification

ISLR 4, ESL 4, Bishop 4

Lab

- LDA, QDA
- Full logistic regression example (incl. Regularisation)
- Text classification (multiple classifiers)

Classification. Three broad approaches

Objective

Divide input space into decision regions so that each input is associated to one (and only one) class

Methods (in this chapter). Decision surfaces are linear in inputs

- Discriminant function
- Probabilistic generative model
- Probabilistic discriminative model

Classification. Three broad approaches (as usual with blurry frontiers)

Objective

Divide input space into decision regions so that each input is associated to one (and only one) class

Methods (in this chapter). Decision surfaces are linear in inputs

- Discriminant function
- Probabilistic generative model
- Probabilistic discriminative model

A decision analytic perspective

K classes. Prior probability for class i

Given class i, feature distribution

$$\pi_i$$
$$f(x|\theta_i) = f_i(x)$$

Given x, posterior probability for class i

$$p_i(x) = \frac{\pi_i f_i(x)}{\sum_{i=1}^K \pi_i f_i(x)} \propto \pi_i f_i(x)$$

Utility for assigning to class j when actually in class i

$$u_{ij}$$

Maximise expected utility

$$\arg \max_j \sum_{i=1}^K u_{ij} \frac{\pi_i f_i(x)}{\sum \pi_i f_i(x)} = \arg \max_j \sum_{i=1}^K u_{ij} \pi_i f_i(x)$$

Discriminant functions

Discriminant functions

A function that takes an input x and assigns it to one of the classes

No direct assessment of class probability given input x

LDA

Feature density given class is normal

Same variance for all classes

$$f_i(x) = \frac{1}{\sqrt{2\pi} \sigma_i} \exp\left(-\frac{1}{2} \frac{(x-\mu_i)^2}{\sigma_i^2}\right)$$
$$\sigma_1^2 = \sigma_2^2 = \dots = \sigma_K^2$$

$$\max_i p_i(x) = \frac{\pi_i \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2} (x-\mu_i)^2\right)}{\sum_{i=1}^K \pi_i \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2} (x-\mu_i)^2\right)} \propto \pi_i \exp\left(-\frac{1}{2\sigma^2} (x-\mu_i)^2\right) \propto \pi_i \exp\left(-\frac{1}{2\sigma^2} (\mu_i^2 - 2\mu_i x)\right)$$

$$\max_i x \frac{\mu_i}{\sigma^2} - \frac{\mu_i^2}{2\sigma^2} + \log(\pi_i)$$

LDA (cont)

K=2, Equal prior probs

Class 1 if

Boundary

Gen, Discriminant function

$$2x(\mu_1 - \mu_2) > \mu_1^2 - \mu_2^2$$

$$x = \frac{\mu_1^2 - \mu_2^2}{2(\mu_1 - \mu_2)} = \frac{\mu_1 + \mu_2}{2}$$

$$\hat{\mathcal{S}}_i(x) = x \frac{\hat{\mu}_i}{\hat{\sigma}^2} - \frac{\hat{\mu}_i^2}{2\hat{\sigma}^2} + \log(\hat{\pi}_i)$$

$$\begin{aligned}\hat{\mu}_i &= \frac{1}{n_i} \sum_{j=1}^{n_i} x_j^i \\ \hat{\sigma}^2 &= \frac{1}{n-K} \sum_{i=1}^K \sum_{j=1}^{n_i} (x_j^i - \hat{\mu}_i)^2 \\ \hat{\pi}_i &= \frac{n_i}{n}\end{aligned}$$

QDA

Feature density given class is normal

Specific variance for each class

$$x|i \sim N(\mu_i, \Sigma_i)$$

Discriminant function

$$\delta_i(x) = -\frac{1}{2} x^T \Sigma_i^{-1} x + x^T \Sigma_i^{-1} \mu_i - \frac{1}{2} \mu_i^T \Sigma_i^{-1} \mu_i - \frac{1}{2} \log |\Sigma_i| + \log \pi_i$$

$$\hat{\Sigma}_i, \hat{\mu}_i, \hat{\pi}_i \rightarrow \operatorname{argmax} \delta_i(x)$$

LDA vs QDA

With p predictors

Algo	Pars	Flexibility	Variance	Bias
LDA	$p(p+1)/2$	Less	Smaller	Bigger
QDA	$K(p(p+1))/2$	More	Bigger	Smaller

Regularised discriminant analysis

$$\begin{aligned}\Sigma_i(\alpha) &= \alpha \hat{\Sigma}_i + (1-\alpha) \hat{\Sigma} \\ \hat{\Sigma}(\delta) &= \delta \hat{\Sigma} + (1-\delta) \hat{\sigma}^2 I\end{aligned}$$

Computations in discriminant analysis

$\Sigma = U D U^T$ U ORTHONORMAL
 D DIAGONAL POSITIVE EIGENVALUES

$$(x - \hat{\mu})^T \Sigma^{-1} (x - \hat{\mu}) = (U^T (x - \hat{\mu}))^T D^{-1} (U^T (x - \hat{\mu}))$$
$$\log |\Sigma| = \sum \log d_{ii}$$

1. $X^* \leftarrow D^{-1/2} U^T X$
2. CLASSIFY NEW INSTANCE TO THE CLOSEST CENTROID,
MODULO π_i

Probabilistic discriminative models

Probabilistic discriminative models

Compute class probabilities given input, typically through a parametric model

Logistic regression

Data

$$(x_i, y_i), i = 1, \dots, n$$

x_i : EXPLANATORY
 $y_i \in \{0, 1\}$

Model

$$y_i | \theta_i \sim \text{Ber}(\theta_i)$$
$$\text{logit } \theta_i = \log \frac{\theta_i}{1 - \theta_i} = \alpha + \beta x_i$$
$$\theta_i = \text{logit}^{-1}(\alpha + \beta x_i) = \frac{e^{\alpha + \beta x_i}}{1 + e^{\alpha + \beta x_i}} = \frac{1}{1 + e^{-(\alpha + \beta x_i)}}$$

Logistic regression. Likelihood formulation

$$l(\alpha, \beta | (x, y)) = \left[\prod_{i: y_i=1} \theta_i(x_i) \right] \left[\prod_{i: y_i=0} (1 - \theta_i(x_i)) \right] = \prod_{i=1}^n \left[\theta_i(x_i)^{y_i} (1 - \theta_i(x_i))^{1-y_i} \right]$$

$$\begin{aligned} \log(l(\alpha, \beta | (x, y))) &= \sum_{y_i=1} \log \theta_i(x_i) + \sum_{y_i=0} \log (1 - \theta_i(x_i)) \\ &= \sum_{i=1}^n y_i \log(\theta_i(x_i)) + (1 - y_i) \log(1 - \theta_i(x_i)) \end{aligned}$$

Logistic regression. Evaluating the likelihood

Some little tricks

$$\frac{\text{logit}^{-1}(z)}{1 - \text{logit}^{-1}(z)} = \frac{\frac{1}{1 + e^{-z}}}{1 - \frac{1}{1 + e^{-z}}} = \frac{1}{e^{-z}} = e^z \quad \log\left(\frac{\text{logit}^{-1}(z)}{1 - \text{logit}^{-1}(z)}\right) = z$$

$$1 - \text{logit}^{-1}(z) = 1 - \frac{1}{1 + e^{-z}} = \frac{e^{-z}}{1 + e^{-z}} = \frac{1}{1 + e^z}$$

$$\log(1 - \text{logit}^{-1}(z)) = -\log(1 + e^z) = -\log p(e^z)$$

$$\begin{aligned} \log(l(\alpha, \beta | (x_i, y_i))) &= \sum_{y_i=1} (\alpha + \beta x_i) - \sum_{y_i=0} \log p(\exp(\alpha + \beta x_i)) \\ &= \sum y_i (\alpha + \beta x_i) - (1 - y_i) \log p(\exp(\alpha + \beta x_i)) \end{aligned}$$

Logistic regression. Gradient descent

Basic approach

$$(\alpha, \beta)_{i+1} = (\alpha, \beta)_i - \eta_i \nabla \log(\ell(\alpha, \beta))$$
$$\nabla \log(\ell(\alpha, \beta)) = \begin{pmatrix} \sum_{i=1}^n 1 - \sum_{i=1}^n \frac{\partial}{\partial \alpha} \log \text{tp}(\exp(\alpha + \beta x_i)) \\ \sum_{i=1}^n x_i - \sum_{i=1}^n \frac{\partial}{\partial \beta} \log \text{tp}(\exp(\alpha + \beta x_i)) \end{pmatrix} \quad \alpha(x)$$

Newton-Raphson

$$(\alpha, \beta)_{i+1} = (\alpha, \beta)_i - H^{-1} \nabla \log(\ell(\alpha, \beta))$$

Iterative reweighted least squares

Logistic regression. Stochastic gradient descent

$$J(w) = \frac{1}{n} \sum_{i=1}^n J_i(w)$$
$$w^{t+1} = w^t - \eta_t \nabla_w J(w)$$
$$= w^t - \eta_t \left(\frac{1}{n} \sum_{i=1}^n \nabla_w J_i(w) \right) \quad O(n)$$

SAMPLE minibatch $B \{i_1, \dots, i_B\} \subset \{1, \dots, n\}$

$$w^{t+1} = w^t - \eta_t \frac{1}{B} \sum_{i \in B} J_i(w^t) \quad O(B)$$

When n is super-large

Robbins Monro (1954) stoch. approx.

$$\sum \eta_t = \infty \quad \sum \eta_t^2 < \infty \Rightarrow \|f(w^t) - f^*\| = O(1/t)$$

Fits in memory, smaller complexity, Parallelise, Escape local optima

Logistic regression. Stochastic gradient descent



Robbins Monro (

Fits in memory, smaller complexity, Parallelise, Escape local optima

Logistic regression. Stochastic gradient descent

$$J(w) = \frac{1}{n} \sum_{i=1}^n J_i(w)$$
$$w^{t+1} = w^t - \eta_t \nabla_w J(w)$$
$$= w^t - \eta_t \left(\frac{1}{n} \sum_{i=1}^n \nabla_w J_i(w) \right)$$



Robbins Monro (1954)

Fits in memory, smaller complexity, Parallelise, Escape local optima

Logistic regression. Interpretation

Odds ratio

$$\theta(x) = \frac{e^{\alpha + \beta x}}{1 + e^{\alpha + \beta x}} \quad \text{ODDS} \quad \frac{\theta(x)}{1 - \theta(x)} = e^{\alpha + \beta x} \in (0, \infty)$$
$$\log \frac{\theta(x)}{1 - \theta(x)} = \alpha + \beta x \quad x \rightarrow x+1 \Rightarrow \log \text{ODDS} + \beta$$
$$\Downarrow$$
$$\text{ODDS} * e^{\beta}$$

$$\text{OR} = \frac{\exp(\alpha + \beta(x+1))}{\exp(\alpha + \beta x)} = e^{\beta}$$

Odds Ratio	X vs Y
1	No association
>1	Bigger output probability
<1	Smaller output probability

Logistic regression. Prediction

Just plug-in estimators to assess probabilities

$$\hat{\alpha}, \hat{\beta}$$
$$P(Y=1|\tilde{x}) = \frac{e^{\hat{\alpha} + \hat{\beta}\tilde{x}}}{1 + e^{\hat{\alpha} + \hat{\beta}\tilde{x}}}$$

Cut-off level

$$P(Y=1|\tilde{x}) \geq k \rightarrow 1$$

Logistic regression with regulariser

$$\max \sum_{i=1}^n \left[y_i \log p(x_i | \beta) + (1 - y_i) \log (1 - p(x_i | \beta)) \right] - \lambda \sum_j \beta_j^2$$
$$\max \sum_{i=1}^n \left[y_i \log p(x_i | \beta) + (1 - y_i) \log (1 - p(x_i | \beta)) \right] - \lambda \left| \sum_i \beta_i \right|$$

Bayesian logistic regression

Likelihood

Generic prior

Generic posterior

$$p(\alpha, \beta | y, x) \propto p(\alpha, \beta) \prod_{i=1}^n \left[p(x_i | \beta)^{y_i} (1 - p(x_i | \beta))^{1-y_i} \right]$$

If few parameters, numerical integration
else MCMC or Laplace approximations

Multiclass logistic regression

OR with K classes

Probabilities

Parameters

Log-likelihood

$$\begin{aligned}\log \frac{P(1|x)}{P(K|x)} &= \beta_{10} + \beta_1^T x \quad \log \frac{P(2|x)}{P(K|x)} = \beta_{20} + \beta_2^T x \quad \dots \quad \log \frac{P(K-1|x)}{P(K|x)} = \beta_{(K-1)0} + \beta_{(K-1)}^T x \\ \Pr(K|x) &= \frac{\exp(\beta_{K0} + \beta_K^T x)}{1 + \sum_{k=1}^{K-1} \exp(\beta_{k0} + \beta_k^T x)} \quad K=1, \dots, K-1 \quad \Pr(K|x) = \frac{1}{1 + \sum_{k=1}^{K-1} \exp(\beta_{k0} + \beta_k^T x)} \\ \theta &= (\beta_{10}, \beta_1^T, \dots, \beta_{(K-1)0}, \beta_{(K-1)}^T) \\ \ell(\theta) &= \sum_{i=1}^N \log p_{y_i}(x_i | \theta) \quad p_K(x_i | \theta) = \Pr(K | x_i, \theta)\end{aligned}$$

Probit regression

Close to logistic regression

$$Pr(1|x) = \Phi(\beta_0 + \beta_1 x) \quad \Phi \leftrightarrow N(0,1)$$
$$\Phi(z) = \frac{1}{2} \left\{ 1 + \operatorname{erf} \left(\frac{z}{\sqrt{2}} \right) \right\}$$

glms

Exponential family

$$p(y|\eta) = b(y) \exp(\eta^T \phi(y)) h(\eta)$$

e.g., Bernoulli

$$\begin{aligned} \theta^y (1-\theta)^{1-y} &= 1 \cdot \left(\frac{\theta}{1-\theta}\right)^y (1-\theta) \\ &= 1 \cdot \exp\left(\left(\log \frac{\theta}{1-\theta}\right) y\right) \cdot (1-\theta) \end{aligned}$$

glm

$$\begin{aligned} &\text{GLM} \\ &y|x; \theta \sim \text{ExpFamily}(\eta) \quad \eta_0(x) = E(y|x, \theta) \quad \eta = \theta^T x \end{aligned}$$

Probabilistic generative models

Probabilistic generative models

Estimate class probabilities and input distribution given class

Compute class probability given input via Bayes formula

Naive Bayes

Variational autoencoders

Naive Bayes

$$Pr(k|x) = \frac{n_k f_k(x)}{\sum_i n_i f_i(x)}$$
$$f_k(x) = \prod_{j=1}^p f_{kj}(x_j)$$

NB ass.

Naive Bayes. Parameter estimation

Continuous

$$X_j | K \sim N(\mu_{kj}, \sigma_{kj}^2)$$

Counts

$$Pr(X_j = x | K) = \frac{n_{kj} + 1}{n_k + h}$$

Categorical

$$Pr(X_i = x_i | K) = \frac{\sum_c I(x_i = x_i) I(Y = y_k)}{\sum_c I(Y = y_k)}$$

Classification metrics

Ok... but how do we assess classifier performance???

(More in labs!!!)

Confusion matrix

		Predicted class	
		+	-
Actual class	+	TP True Positives	FN False Negatives Type II error
	-	FP False Positives Type I error	TN True Negatives

Confusion matrix

		Predicted class	
		+	-
Actual class	+	TP True Positives	FN False Negatives Type II error
	-	FP False Positives Type I error	TN True Negatives

Metric	Formula	Interpretation
Accuracy	$\frac{TP + TN}{TP + TN + FP + FN}$	Overall performance of model
Precision	$\frac{TP}{TP + FP}$	How accurate the positive predictions are
Recall Sensitivity	$\frac{TP}{TP + FN}$	Coverage of actual positive sample
Specificity	$\frac{TN}{TN + FP}$	Coverage of actual negative sample

Imbalanced problems

Accuracy not sufficient in imbalanced problems

- Example. CTR (click-through rate, click 1, no-click 0)
 - About 10^8 ads (observations), only 80000 clicks
 - A model classifying always as 0, accuracy of 99.92%

Unequal class distribution inherent in

Fraud detection

Anomaly detection

Credit default prediction

Conversion prediction

Intrusion detection....

Imbalanced problems

F1 score	$\frac{2TP}{2TP + FP + FN}$	Hybrid metric useful for unbalanced classes
----------	-----------------------------	---

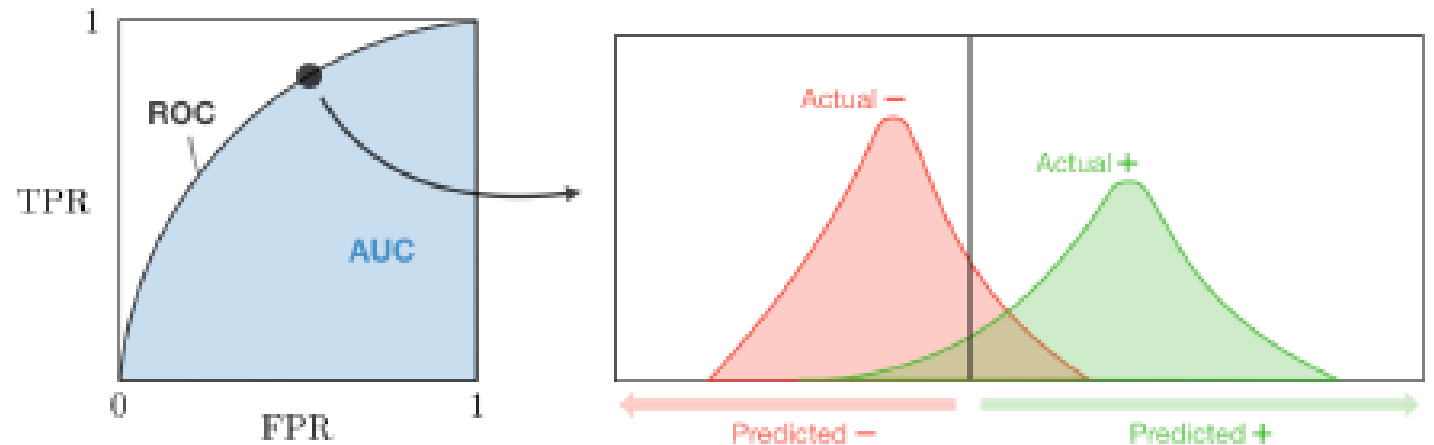
Oversampling SMOTE Synthetic Minority Oversampling Tech

Undersampling

Bayesian methods

ROC (receiver operating characteristics) AUC (area under curve)

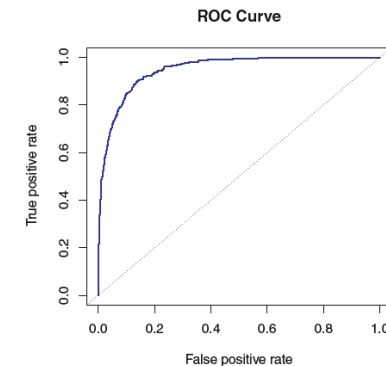
Metric	Formula	Equivalent
True Positive Rate TPR	$\frac{TP}{TP + FN}$	Recall, sensitivity
False Positive Rate FPR	$\frac{FP}{TN + FP}$	1-specificity



More generally, need to take into account costs

Fraud, no fraud

Cannabinoids. Decisions entail researching or not researching drug



Calibration

Discrimination vs Calibration

The Achilles heel of predictive analytics

<https://bmcmmedicine.biomedcentral.com/articles/10.1186/s12916-019-1466-7>

When you say probability of A is 70% you are right 70% of the times

Back to the decision analytic perspective...

K classes. Prior probability for class i

Given class i, feature distribution

Given x, posterior probability for class i

Utility for assigning to class j when actually in class i

Maximise expected utility

$$\pi_i$$
$$f(x|\theta_i) = f_i(x)$$

$$p_i(x) = \frac{\pi_i f_i(x)}{\sum_{i=1}^K \pi_i f_i(x)} \propto \pi_i f_i(x)$$

$$u_{ij}$$

$$\arg \max_j \sum_{i=1}^K u_{ij} \frac{\pi_i f_i(x)}{\sum \pi_i f_i(x)} = \arg \max_j \sum_{i=1}^K u_{ij} \pi_i f_i(x)$$

Classification. To be seen

C+R. Decision trees, random forests, boosting

C (+R). Support vector machines

R +C. Perceptrons. Neural networks. Deep neural nets

See you next week

introml@icmat.es

Stuff at

https://datalab-icmat.github.io/courses_stats.html